

CSCI 4155/6505: Machine Learning with Robotics

Thomas P. Trappenberg
Dalhousie University

Acknowledgements

These lecture notes have been inspired by several great sources, which I commend as further readings. In particular, Andrew Ng from Stanford University has several lecture notes on Machine Learning (CS229) and Artificial Intelligence: Principles and Techniques (CS221). His lecture notes and video links to his lectures are available on his web site (<http://robotics.stanford.edu/~ang>). Excellent book on the theory of machine learning are *Introduction to Machine Learning* by Ethem Alpaydin, 2nd edition, MIT Press 2010, and *Pattern Recognition and Machine Learning* by Christopher Bishop, Springer 2006. A wonderful book on Robotics with a focus of Bayesian models is *Probabilistic Robotics* by S. Thrun, W. Burgard, and D. Fox, MIT Press 2005, and the standard book on RL is *Reinforcement Learning: An Introduction* by Richard Sutton and Andrew Barto, MIT press, 1998. The standard book for AI, *Artificial Intelligence: A Modern Approach* by Stuart Russell and Peter Norvig, 2nd edition, Prentice Hall, 2003, does also include some chapters on Machine Learning.

Several people have contributed considerably to this lecture notes. In particular, Leah Brown and Ian Graven have created most of the original Lego Mindstorm implementations and tutorials, Paul Hollensen and Patrick Connor made several contributions, and Chris Maxwell helped with some implementation issues.

Contents

I INTRODUCTION	
1 Introduction	3
1.1 Some history of AI and machine learning	3
1.2 Machine Learning and the probabilistic framework	5
1.3 Robotics and control theory	6
1.4 Sensing and acting	7
II BACKGROUND	
2 Mathematical tools	13
2.1 Probability theory	13
2.2 Vector and matrix notations	25
2.3 Basic calculus	28
3 Programming with Matlab	30
3.1 The MATLAB programming environment	30
3.2 Main programming constructs	31
3.3 Creating MATLAB programs	37
3.4 Graphics	38
3.5 A first project: modelling the world	40
3.6 Alternative programming environments: Octave and Scilab	42
4 Basic robotics with LEGO NXT	45
4.1 Building a basic robot with LEGO NXT	45
4.2 Basic MATLAB NXT toolbox commands	49
4.3 First examples	54
4.4 Classical control theory	56
4.5 Adaptive Controle	59
4.6 Configuration space	61
4.7 Planning as search	64
III SUPERVISED LEARNING	
5 Regression and maximum likelihood	69
5.1 Regression of noisy data	69

5.2	Probabilistic models and maximum likelihood	70
5.3	Maximum a posteriori estimates	73
5.4	Minimization procedures	74
5.5	Gradient Descent	75
5.6	Non-linear regression and the bias-variance tradeoff	80
6	Classification	84
6.1	Logistic regression	84
6.2	Generative models	86
6.3	Discriminant analysis	87
6.4	Naive Bayes	90
7	Graphical models	92
7.1	Causal models	92
7.2	Bayes Net toolbox	94
7.3	Temporal Bayesian networks: Markov Chains and Bayes filters	96
7.4	Application and generalization: Localisation example	99
8	General learning machines	100
8.1	The Perceptron	100
8.2	Multilayer perceptron (MLP)	102
8.3	LIP regression	105
8.4	Support Vector Machines (SVM)	106
8.5	SV-Regression and implementation	114
8.6	Supervised Line-following	116
IV REINFORCEMENT LEARNING		
9	Markov Decision Process	121
9.1	Learning from reward and the credit assignment problem	121
9.2	The Markov Decision Process	122
9.3	The Bellman equation	125
10	Temporal Difference learning and POMDP	131
10.1	Temporal Difference learning	131
10.2	Temporal difference methods for optimal control	132
10.3	Robot exercise with reinforcement learning	133
10.4	POMDP	135
10.5	Model-based RL	135

V UNSUPERVISED LEARNING

11 Unsupervised learning	139
11.1 K-means clustering	139
11.2 Mixture of Gaussian and the EM algorithm	141
11.3 The Boltzmann machine	145

Part I

Introduction

1 Introduction

Both, machine learning and robotics have been important topics in the area of artificial intelligence. This introductory chapter outlines some history of these areas before outlining the direction that have recently advanced these areas considerably. We will also outline some basic robotics terminology and technology as background to the main machine learning focus of this course. The following x chapters reviews some additional material that we will use later, including probability theory, programming with Matlab, and how to use the Lego Mindstorm NXT.

1.1 Some history of AI and machine learning

Artificial Intelligence(AI) has many sub-discipline such as knowledge representation, expert systems, search, etc. This course will focus on how machines can learn to improve their performance or learn how to solve problems. **Machine learning** is now widely respected scientific area with a growing number of applications. Indeed, major breakthroughs in autonomous robotics, new types of consumer electronics, and new approaches to information processing have recently been enabled by new machine learning advancements.

The history of AI is tightly interwoven with the history of machine learning. For example, Arthur Samuel's checkers program from the 1950s, which has been celebrated as an early AI hallmark, was able to learn from experience and thereby was able to outperform its creator. Basic Neural Networks, such as Bernhard Widrow's ADALINE (1959), Karl Steinbuch's Lernmatrix (around 1960), and Frank Rosenblatt's Perceptron (late 1960s), have sparked the imaginations of researcher about brain-like information processing systems. And Richard Bellman's Dynamic Programming (1953) has created a lot of excitement during the 1950s and 60s, and is now considered one of the foundations reinforcement learning.

Biological systems of often been and inspiration for understanding learning systems and visa versa. Donald Hebb's book *The Organization of Behavior* (1943) has been very influential not only in the cognitive science community, but has been marvelled in the early computer science community. While Hebb speculated how self-organizing mechanisms can enable human behavior, it was Eduardo Caianiello's influential paper *Outline of a theory of thought-processes and thinking machines* (1961) who quantified these ideas into two important equations, the neuron equation, describing how the functional elements of the networks behave, and the memmonic equation that describe how these systems learn. This opened the doors to more theoretical investigations. But such investigations came to a sudden halt after Marvin Minsky and Seymore Papert published their book *Perceptrons* in 1969 with a proof that simple perceptrons can not learn all problems. At this time, likely somewhat triggered by the vacuum in AI

research by the departure of learning networks, mainstream AI shifted towards expert systems.



Fig. 1.1 Some AI pioneers AI. From top left to bottom right: Alan Turing, Frank Rosenblatt, Geoffrey Hinton, Arthur Lee Samuel, Richard Bellman, Judea Pearl.

Neural Networks became again popular in the mid 1980 after the backpropagation algorithms was popularized by David Rumelhart, Geoffrey Hinton and Ronald Williams (1986). There was a brief period of extreme hype with claims that neural networks can now forecast the stock-market and how these learning systems will quickly become little brains. The hype backslashed somewhat. Neural Networks predictive power in scientific explanations became questions as they seem to always fit any data, and early claims of the future progress has not substantiated. But the understanding of both these problems, which are related to **generalizability** and **scalability** have matured this field considerably since.

The recent progress was made possible through several factors, but likely most of all through a more rigorous mathematical formulations and the embedding of such methods with stochastic theories. In particular, important learning theories become widely known and developed further after Vladimir Vapnik published his book *The Nature of Statistical Learning Theory* in 1995. And more generally, Bayesian methods have clarified and transformed the field, with many important advancements outlined in this book (Judea Pearl, 1985).

1.2 Machine Learning and the probabilistic framework

Traditional AI provides many useful approaches for solving specific problems. For example, search algorithms can be used to navigate mazes or to find scheduling solutions. And expert systems can manage a large data-base of expert knowledge that can be used to argue (infer) about specific situations. While such strategies might be well suited for specific applications, learning systems are usually more general and can be applied to situations for which closed solutions are not known. Also, a major challenge in many AI applications has been that systems change over time and that systems encounter situations for which they were not designed. Thus, some form of adaptation to changing situations and generalizations to unseen environments are important for many systems for which AI is being considered. While many AI systems have adaptive components, we are specifically concentrating on the theory of learning machines in this course.

Our aim of learning machines is to learn from the environment, either through instructions, by reward or punishment, or just by exploring the environment and learning about typical objects or spatio-temporal relations. For the systematic discussion of learning systems it is useful to distinguish three types of learning circumstances, namely supervised learning, unsupervised learning, and reinforcement learning. **Supervised learning** is characterized by using explicit examples with labels that the system should use to learn to predict labels of previously unseen data. The training labels can be seen as being supplied by a teacher who tells the system exactly the desired answer. **Reinforcement learning** is somewhat similar in that the system receives some feedback from the environment, but this feedback is only reward or punishment for previously taken actions and thus typically delay in time. The goal of reinforcement learning is to discover the action, or a sequence of actions, which maximize reward over time. Finally, the aim of **unsupervised learning** is to find useful representations of data based on regularities of data without labels. Smart representation of data is often the key of smart solutions, and this type of learning is often more applicable since no labels are required. While it is useful to distinguish such classical learning systems, they can also augment each other. Our ultimate goal is to model the world and to use such models to make ‘smart’ decisions in order to maximize reward.

Learning systems address can help to solve problems for which more direct solutions are not known, another common problem in AI applications is that systems are generally **unreliable** and that environments are **uncertain**. For example, a program might read data in a specific format, but some user might supply corrupted files. Indeed, production software is often lengthy only to consider all kind of situations that could occur, and we now realize that considering all possibilities is often impossible. Unreliable inputs is very apparent in robotics where sensors are often noisy or have severe limitations. Also, the state of a system might be unclear due to limited data or the inability to process data sufficiently in time due to limited resources.

Major progress in many AI areas, in particular in robotics and machine learning, has been made by using concepts of (Bayesian) probability theory. This language of probability theory is appropriate since it acknowledges the fundamental limitations we have in real world applications (such as limited computational resources or inaccurate sensors). The language of probability theory has certainly helped to unify much of related areas and improved communication between researchers. Furthermore, we will

see that the representation of uncertain states as a probabilistic map will be very useful. Probability theory will also provide us with the solution for a basic computational need, that of combining prior knowledge with new evidence.

1.3 Robotics and control theory

In this course, we will demonstrate many of the discussions, algorithms and associated problems with the help of computer simulations and robots. Computer simulations are a great way to experiment with many of the ideas, but physical implementations have often the advantage to show more clearly the challenges in real applications. So our emphasis in this book is using fairly general machine learning techniques to solve Robotics tasks even though more direct engineering solutions might be possible. While this is not always the way robots are controlled today, it is also true that machine learning methods become increasingly important in robotics.

So, what is robotics?

"Robotics is the science of perceiving and manipulating the physical world through computer-controlled devices"¹.

We use the word **robot** or **agent** to describe a system which interacts actively with the environment through **sensors** and **actuators**. The 'brain' of the robot is often called the **controller**. An agent can be implemented in software or hardware, and we often use the term **robot** more specifically for a hardware implementation of an agent, though this does generally include software components. **Active interactions** means that the robot's position in the environment or the environment itself changes through the action of the robot. **Mobile robots** are a good example of this and are mainly used in the examples of this course. In contrast, a vision system that uses a digital camera to recognize objects is not a robot in our definition.

The word 'robot' is sometimes credited to Isaac Asimov (1941) or to the Czech writer Karel Čapek (1921). The dream of having machines that can act more autonomously for human benefits is certainly older, and the Unimate (1961) is often referred to as the first industrial robot. Robots are now invaluable in manufacturing, and research is ongoing to make robots more autonomous and to make those machines more robust and able to work in hostile environments. Robotics has many components, including mechanical and electrical engineering, computer or machine vision, and even behavioral studies have become prominent in this field.

Robots are intended to make useful actions to achieve some goals, so robotics is all about finding and safely executing those appropriate actions. This area is generally called **control theory**. A functioning robotics system needs to control on many different levels, controlling the low level functions such as the proper rotation of motors up to ensuring that high level tasks are accomplished. We will formalize some control theory in a later chapter but it is useful to get some of the larger picture build the basis of current robotics systems.

There are two extreme approaches to robotics control, the **deliberative approach** and the **reactive approach**. In the deliberative approach we gather all available information plan action carefully based on all available actions. Such a **planning process**

¹Probabilistic Robotics, Sebastian Thrun, Wolfram Burgard, and Dieter Fox, MIT press 2006

usually takes time and is based on searching through all available alternatives. The advantage of such systems is that they usually provide superior actions, but the search for such actions can be time consuming and might thus not be applicable in all situations. Also, deliberative systems require a large amount of knowledge about the environment that might not be available in certain applications.

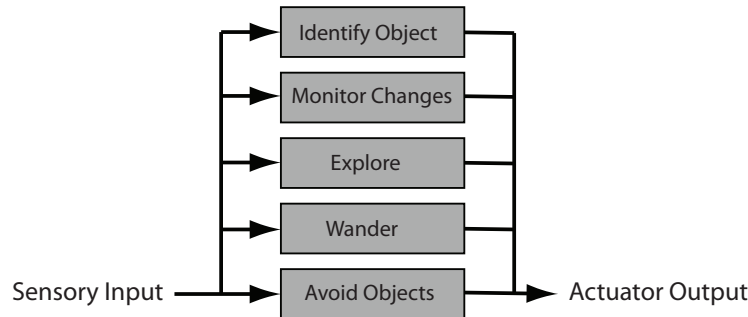


Fig. 1.2 An example of an subsumption architecture. From Maja J. Matari c, *The Robotocs Primer*, MIT Press 2007

Reactive systems have a more direct approach of translating sensory information in actions. For example, a typical rule in such systems might be to turn the robot when a proximity sensor senses some objects in its path. To generate more complex behavior, such reactive systems typically build response systems by combining lower level control systems with higher level functions. A very successful approach has there been the **subsumption architecture**, which is illustrated in Figure 1.2. Such systems are typically build bottom-up in that more complex functions use lower-level functions to achieve more complex tasks. The higher level modules can chose to use the lower level function or can inhibit them if necessary.

1.4 Sensing and acting

As stressed above, robots act in the physical world through actuators based on sensory information. So you might be interested to briefly review some sensing and acting technology.

Actuators are mainly motors to move the robot around (locomotion) or to move limbs to crasp objects (manipulation). Motors for continuous rotations are typically DC (direct current) motors, but in robotics we often need to move a limb to a specific location. Motors that can turn to a specific position are called **servo motors**. Such motors are based on greas and position sensors together with some electronics for control or the desired rotation angle. Figure 1.3.

Sensors come in many varieties. Table 1.1 gives some examples of what kind of sensing technology is often used to sense (measure) certain physical properties. The basic Lego sensors used in this course are also shown in Figure 1.3. In addition we will use a web camera, and for special projects we could also use the Kinect ensors

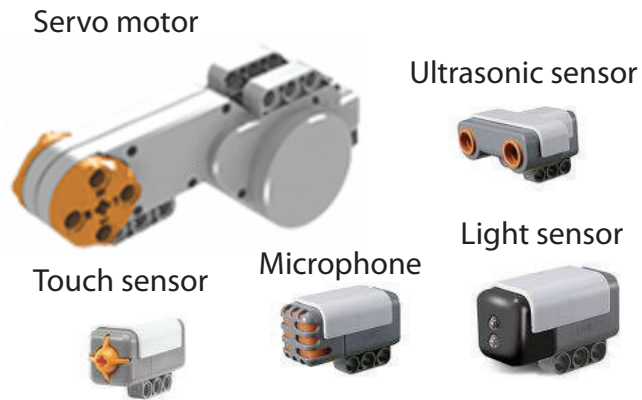


Fig. 1.3 The actuator and sensors of the basic Lego Mindstorm NXT robotics set.

Table 1.1 Some sensors and the information they measure. From From Maja J. Matarić, *The Robotics Primer*, MIT Press 2007

Physical Property	Sensing Technology
Contact	bump, switch
Distance	ultrasound, radar, infra red
Light level	photocells, cameras
Sound levels	microphones
Strain	strain gauges
Rotation	encoders, potentiometers
Acceleration	accelerometers, gyroscopes
Magnetism	compass
Smell	chemical sensors
Temperature	thermal, infra red
Inclination	inclinometers, gyroscopes
Pressure	pressure gauge
Altitude	altimeter

from Microsoft and special Lego sensors including a GPS and an accelerometer (see Figure 1.4).



Fig. 1.4 Some additional sensors that we could use such as the Microsoft Kinect and a GPS for the Lego NXT.

Part II

Background

2 Mathematical tools

2.1 Probability theory

As outlined in Chapter 1, a major milestone for the modern approach to machine learning is to acknowledge our limited knowledge about the world and the unreliability of sensors and actuators. It is then only natural to consider quantities in our approaches as **random variables**. While a regular variable, once set, has only one specific value, a random variable will have different values every time we ‘look’ at it (draw an example from the distributions). Just think about a light sensor. We might think that an ideal light sensor will give us only one reading while holding it to a specific surface, but since the peripheral light conditions change, the characteristics of the internal electronic might change due to changing temperatures, or since we move the sensor unintentionally away from the surface, it is more than likely that we get different readings over time. Therefore, even internal variables that have to be estimated from sensors, such as the state of the system, is fundamentally a random variable.

A common misconception about randomness is that one can not predict values of random values. Some values might be more likely than others, and, while we might not be able to predict a specific value when drawing a random number, it is possible to say something like how often a certain number will appear when drawing many examples. We might even be able to state how confident we are with this number, or, in other words, how variable these predictions are. The complete knowledge of a random variable, that is, how likely each value is for a random variable x , is captured by the **probability density function** $pdf(x)$. We discuss some specific examples of pdfs below. In these examples we assume that we know the pdf, but in many practical applications we must estimate this function. Indeed, estimation of pdfs is at the heart if not the central tasks of machine learning. If we would know the ‘world pdf’, the probability function of all possible events in the world, then we could predict as much as possible in this world.

Most of the systems discussed in this course are **stochastic models** to capture the uncertainties in the world. Stochastic models are models with random variables, and it is therefore useful to remind ourselves about the properties of such variables. This chapter is a refresher on concepts in probability theory. Note that we are mainly interested in the language of probability theory rather than statistics, which is more concerned with hypothesis testing and related procedures.

2.1.1 Random numbers and their probability (density) function

Probability theory is the theory of **random numbers**. We denoted such numbers by capital letters to distinguish them from regular numbers written in lower case. A random variable, X , is a quantity that can have different values each time the variable is inspected, such as in measurements in experiments. This is fundamentally different

to a regular variable, x , which does not change its value once it is assigned. A random number is thus a new mathematical concept, not included in the regular mathematics of numbers. A specific value of a random number is still meaningful as it might influence specific processes in a deterministic way. However, since a random number can change every time it is inspected, it is also useful to describe more general properties when drawing examples many times. The frequency with which numbers can occur is then the most useful quantity to take into account. This frequency is captured by the mathematical construct of a **probability**. Note that there is often a debate if random numbers should be defined solely on the basis of a frequency measurement, or if there they should be treated as a special kind of objects. This philosophical debate between ‘Frequentists’ and ‘Bayesians’ is of minor importance for our applications.

We can formalize the idea of expressing probabilities of drawing specific values for random variable with some compact notations. We speak of a **discrete random variable** in the case of discrete numbers for the possible values of a random number. A **continuous random variable** is a random variable that has possible values in a continuous set of numbers. There is, in principle, not much difference between these two kinds of random variables, except that the mathematical formulation has to be slightly different to be mathematically correct. For example, the **probability function**,

$$P_X(x) = P(X = x) \quad (2.1)$$

describes the frequency with which each possible value x of a discrete variable X occurs. Note that x is a regular variable, not a random variable. The value of $P_X(x)$ gives the fraction of the times we get a value x for the random variable X if we draw many examples of the random variable.² From this definition it follows that the frequency of having any of the possible values is equal to one, which is the normalization condition

$$\sum_x P_X(x) = 1. \quad (2.2)$$

In the case of continuous random numbers we have an infinite number of possible values x so that the fraction for each number becomes infinitesimally small. It is then appropriate to write the probability distribution function as $P_X(x) = p_X(x)dx$, where $p_X(x)$ is the **probability density function** (pdf). The sum in eqn 2.2 then becomes an integral, and normalization condition for a continuous random variable is

$$\int_x p_X(x)dx = 1. \quad (2.3)$$

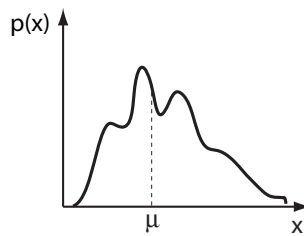
We will formulate the rest of this section in terms of continuous random variables. The corresponding formulas for discrete random variables can easily be deduced by replacing the integrals over the pdf with sums over the probability function. It is also possible to use the δ -function, outlined in Appendix 2.2, to write discrete random processes in a continuous form.

²Probabilities are sometimes written as a percentage, but we will stick to the fractional notation.

2.1.2 Moments: mean, variance, etc.

In the following we only consider independent random values that are drawn from identical pdfs, often labeled as iid (independent and identically distributed) data. That is, we do not consider cases where there is a different probabilities of getting certain numbers when having a specific number in a previous trial. The static probability density function describes, then, all we can know about the corresponding random variable.

Let us consider the arbitrary pdf, $p_X(x)$, with the following graph:



Such a distribution is called **multimodal** because it has several peaks. Since this is a pdf, the area under this curve must be equal to one, as stated in eqn 2.3. It would be useful to have this function parameterized in an analytical format. Most pdfs have to be approximated from experiments, and a common method is then to fit a function to the data. We can also view this approximation as a learning problem, that is, how can we learn the pdf from data? We will return to this question later.

Finding a precise form of a pdf is difficult, and we became thus used to describing random variables with a small set of numbers that are meant to capture some properties. For example, we might ask what the most frequent value is when drawing many examples. This number is given by the largest peak value of the distribution. It is often more useful to know something about the average value itself when drawing many examples. A common quantity to know is thus the expected arithmetic average of those numbers, which is called the **mean, expected value, or expectation value** of the distribution, defined by

$$\mu = \int_{-\infty}^{\infty} xp(x)dx. \quad (2.4)$$

This formula formalizes the calculation of adding all the different numbers together with their frequency.

A careful reader might have noticed a little oddity in our discussion. On the one hand we are saying that we want to characterize random variables through some simple measurements because we do not know the pdf, yet the last formula uses the pdf $p(x)$ that we usually don't know. To solve this apparent oddity we need to be more careful and talk about the **true underlying functions** and the **estimation** of such functions. If we would know the pdf that governs the random variable X , then equation 2.4 is the definition of the mean. However, in most applications we do not know the pdf, but we can define an approximation of the mean from measurements. For example, if we measure the frequency p_i of values in certain intervals around values x_i , then we can estimate the true mean μ by

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i p_i. \quad (2.5)$$

It is a common practice to denote an estimate of a quantity by adding a hat symbol to the quantity name. Also, note that we have used here a discretization procedure to approximate a random variable that can be continuous in the most general case. Also note that we could enter here again the philosophical debate. Indeed, we have treated the pdf as fundamental and described the arithmetic average like an estimation of the mean. This might be viewed as *Bayesian*. However, we could also be pragmatic and say that we only have a collection of measurements so that the numbers are the ‘real’ thing, and that pdfs are only a mathematical construct. We will continue with a Bayesian description but note that this makes no difference at the end when using it in specific applications.

The mean of a distribution is not the only interesting quantity that characterizes a distribution. For example, we might want to ask what the **median** value is for which it is equally likely to find a value lower or larger than this value. Furthermore, the spread of the pdf around the mean is also very revealing as it gives us a sense of how spread the values are. This spread is often characterized by the standard deviation (std), or its square, which is called **variance**, σ^2 , and is defined as

$$\sigma^2 = \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx. \quad (2.6)$$

This quantity is generally not enough to characterize the probability function uniquely; this is only possible if we know all moments of a distribution, where the n th **moment about the mean** is defined as

$$m^n = \int_{-\infty}^{\infty} (x - \mu)^n f(x) dx. \quad (2.7)$$

The **variance** is the second moment about the mean,

$$\sigma^2 = \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx. \quad (2.8)$$

Higher moments specify further characteristics of distributions such as terms with third-order exponents (lie a quantity called skewness) or fourth-order (such as a quantity called kurtosis). Moments higher than this have not been given explicit names. Knowing all moments of a distribution is equivalent to knowing the distribution precisely, and knowing a pdf is equivalent to knowing everything we could know about a random variable.

In case the distribution function is not given, moments have to be estimated from data. For example, the mean can be estimated from a sample of measurements by the **sample mean**,

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad (2.9)$$

and the variance from the **sample variance**,

$$s_1^2 = \frac{1}{n} \sum_{i=1}^n (\bar{x} - x_i)^2. \quad (2.10)$$

We will discuss later that these are the appropriate maximum likelihood estimates of these parameters. Note that the sample mean is an **unbiased estimate** while the sample variance is **biased**. A statistic is said to be biased if the mean of the sampling distribution is not equal to the parameter that is intended to be estimated. It can be shown that $E(s_1^2) = \frac{1}{n}\sigma^2$, and we can therefore adjust for the bias with a different normalization. It is hence common to use the **unbiased sample variance**

$$s_2^2 = \frac{1}{n-1} \sum_{i=1}^n (\bar{x} - x_i)^2, \quad (2.11)$$

as estimator of the variance.

Finally, it is good to realize that knowing all moments uniquely specifies a pdf. But the reverse is also true, that is, an incomplete list of moments does not uniquely define a pdf. Note that the all higher moments are zero for the Gaussian distributions, which means that the mean and variance uniquely define the distribution. This is however not the case for other distributions, and the usefulness of reporting these statistics can then be questioned.

2.1.3 Examples of probability (density) functions

There is an infinite number of possible pdfs. However, some specific forms have been very useful for describing some specific processes and have thus been given names. The following are just some examples of discrete and several continuous distributions. Most examples are discussed as one-dimensional distributions except the last example, which is a higher dimensional distribution.

2.1.3.1 Bernoulli distribution

A Bernoulli random variable is a variable from an experiment that has two possible outcomes: success with probability p ; or failure, with probability $(1 - p)$.

Probability function:

$$P(\text{success}) = p; P(\text{failure}) = 1 - p$$

mean: p

variance: $p(1 - p)$

2.1.3.2 Multinomial distribution

This is the distribution of outcomes in n trials that have k possible outcomes. The probability of each outcome is thereby p_i .

Probability function:

$$P(x_i) = n! \prod_{i=1}^k (p_i^{x_i} / x_i!)$$

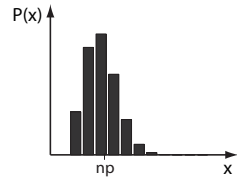
mean: np_i

variance: $np_i(1 - p_i)$

An important example is the Binomial distribution ($k = 2$), which describes the the number of successes in n Bernoulli trials with probability of success p . Note that the binomial coefficient is defined as

$$\binom{n}{x} = \frac{n!}{x!(n-x)!} \quad (2.12)$$

and is given by the MATLAB function `nchoosek`.



Probability function:

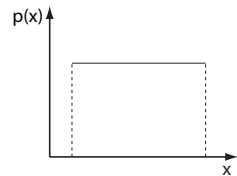
$$P(x) = \binom{n}{x} p^x (1-p)^{n-x}$$

mean: np

variance: $np(1-p)$

2.1.3.3 Uniform distribution

Equally distributed random numbers in the interval $a \leq x \leq b$. Pseudo-random variables with this distribution are often generated by routines in many programming languages.



Probability density function:

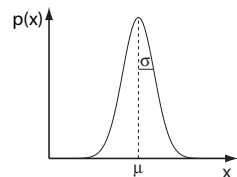
$$p(x) = \frac{1}{b-a}$$

mean: $(a+b)/2$

variance: $(b-a)^2/12$

2.1.3.4 Normal (Gaussian) distribution

Limit of the binomial distribution for a large number of trials. Depends on two parameters, the mean μ and the standard deviation σ . The importance of the normal distribution stems from the central limit theorem outlined below.



Probability density function:

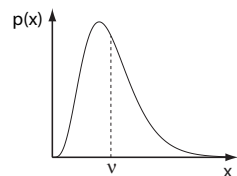
$$p(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

mean: μ

variance: σ^2

2.1.3.5 Chi-square distribution

The sum of the squares of normally distributed random numbers is chi-square distributed and depends on a parameter ν that is equal to the mean. Γ is the gamma function included in MATLAB as `gamma`.



Probability density function:

$$p(x) = \frac{x^{(\nu-2)/2} e^{-x/2}}{2^{\nu/2} \Gamma(\nu/2)}$$

mean: ν

variance: 2ν

2.1.3.6 Multivariate Gaussian distribution

We will later consider density functions of a several random variables, x_1, \dots, x_n . Such density functions are functions in higher dimensions. An important example is the multivariate Gaussian (or Normal) distribution given by

$$p(x_1, \dots, x_n) = p(\mathbf{x}) = \frac{1}{(\sqrt{2\pi})^n \sqrt{|\det(\boldsymbol{\Sigma})|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right). \quad (2.13)$$

This is a straight forward generalization of the one-dimensional Gaussian distribution mentioned before where the mean is now a vector, $\boldsymbol{\mu}$ and the variance generalizes to a covariance matrix, $\boldsymbol{\Sigma} = [\text{Cov}[X_i, X_j]]_{i=1,2,\dots,k;j=1,2,\dots,k}$ which must be symmetric and positive semi-definit. An example with mean $\boldsymbol{\mu} = (1 \ 2)^T$ and covariance $\boldsymbol{\Sigma} = \begin{pmatrix} 1 & 0.5 \\ 0.5 & 1 \end{pmatrix}$ is shown in Fig.2.1.

2.1.4 Cumulative probability (density) function and the Gaussian error function

We have mainly discussed probabilities of single values as specified by the probability (density) functions. However, in many cases we want to know the probabilities of having values in a certain range. Indeed, the probability of a specific value of a continuous random variable is actually infinitesimally small (nearly zero), and only the probability of a range of values is finite and has a useful meaning of a probability. This integrated version of a probability density function is the probability of having a value x for the random variable X in the range of $x_1 \leq x \leq x_2$ and is given by

$$P(x_1 \leq X \leq x_2) = \int_{x_1}^{x_2} p(x) dx. \quad (2.14)$$

Note that we have shortened the notation by replacing the notation $P_X(x_1 \leq X \leq x_2)$ by $P(x_1 \leq X \leq x_2)$ to simplify the following expressions. In the main text we often need to calculate the probability that a normally (or Gaussian) distributed variable has values between $x_1 = 0$ and $x_2 = y$. The probability of eqn 2.14 then becomes a function of y . This defines the **Gaussian error function**

$$\frac{1}{\sqrt{2\pi}\sigma} \int_0^y e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx = \frac{1}{2} \text{erf}\left(\frac{y-\mu}{\sqrt{2}\sigma}\right). \quad (2.15)$$

The name of this function comes from the fact that this integral often occurs when calculating confidence intervals with Gaussian noise and is often abbreviated as erf. This Gaussian error function for normally distributed variables (Gaussian distribution with mean $\mu = 0$ and variance $\sigma = 1$) is commonly tabulated in books on statistics. Programming libraries also frequently include routines that return the values for specific arguments. In MATLAB this is implemented by the routine `erf`, and values for the inverse of the error function are returned by the routine `erfinv`.

Another special case of eqn 2.14 is when x_1 in the equation is equal to the lowest possible value of the random variable (usually $-\infty$). The integral in eqn 2.14 then

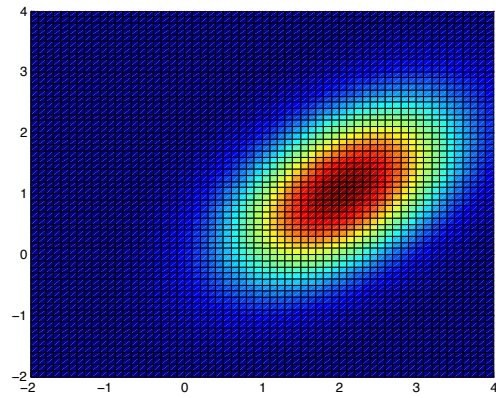


Fig. 2.1 Multivariate Gaussian with mean $\mu = (1 \ 2)^T$ and covariance $\Sigma = (1 \ 0.5; 0.5 \ 1)$.

corresponds to the probability that a random variable has a value smaller than a certain value, say y . This function of y is called the **cumulative density function (cdf)**,³

$$P^{\text{cum}}(x < y) = \int_{-\infty}^y p(x) dx, \quad (2.16)$$

which we will utilize further below.

³Note that this is a probability function, not a density function.

2.1.5 Functions of random variables and the central limit theorem

A function of a random variable X ,

$$Y = f(X), \quad (2.17)$$

is also a random variable, Y , and we often need to know what the pdf of this new random variable is. Calculating with functions of random variables is a bit different to regular functions and some care has to be given in such situations. Let us illustrate how to do this with an example. Say we have an equally distributed random variable X as commonly approximated with pseudo-random number generators on a computer. The probability density function of this variable is given by

$$p_X(x) = \begin{cases} 1 & \text{for } 0 \leq x \leq 1, \\ 0 & \text{otherwise.} \end{cases} \quad (2.18)$$

We are seeking the probability density function $p_Y(y)$ of the random variable

$$Y = e^{-X^2}. \quad (2.19)$$

The random number Y is **not** Gaussian distributed as we might think naively. To calculate the probability density function we can employ the cumulative density function eqn 2.16 by noting that

$$P(Y \leq y) = P(e^{-X^2} \leq y) = P(X \geq \sqrt{-\ln y}). \quad (2.20)$$

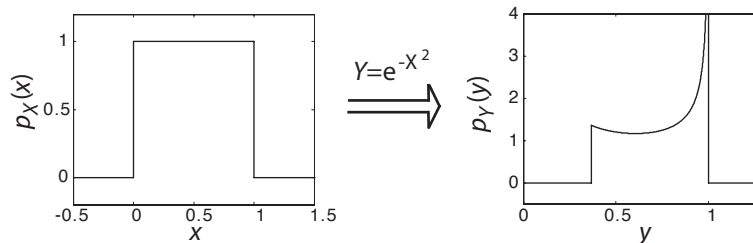
Thus, the cumulative probability function of Y can be calculated from the cumulative probability function of X ,

$$P(X \geq \sqrt{-\ln y}) = \begin{cases} \int_{\sqrt{-\ln y}}^1 p_X(x) dy = 1 - \sqrt{-\ln y} & \text{for } e^{-1} \leq y \leq 1, \\ 0 & \text{otherwise.} \end{cases} \quad (2.21)$$

The probability density function of Y is the derivative of this function,

$$p_Y(y) = \begin{cases} 1 - \sqrt{-\ln y} & \text{for } e^{-1} \leq y \leq 1, \\ 0 & \text{otherwise.} \end{cases} \quad (2.22)$$

The probability density functions of X and Y are shown below.



A special function of random variables, which is of particular interest it can approximate many processes in nature, is the sum of many random variables. For example, such a sum occurs if we calculate averages from measured quantities, that is,

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i, \quad (2.23)$$

and we are interested in the probability density function of such random variables. This function depends, of course, on the specific density function of the random variables X_i . However, there is an important observation summarized in the **central limit theorem**. This theorem states that the average (normalized sum) of n random variables that are drawn from any distribution with mean μ and variance σ is approximately normally distributed with mean μ and variance σ/n for a sufficiently large sample size n . The approximation is, in practice, often very good also for small sample sizes. For example, the normalized sum of only seven uniformly distributed pseudo-random numbers is often used as a pseudo-random number for a normal distribution.

2.1.6 Measuring the difference between distributions

An important practical consideration is how to measure the similarity of difference between two density functions, say the density function p and the density function q . Note that such a measure is a matter of definition, similar to distance measures of real numbers or functions. However, a proper distance measure, d , should be zero if the items to be compared, a and b , are the same, it's value should be positive otherwise, and a distance measure should be symmetrical, meaning that $d(a, b) = d(b, a)$. The following popular measure of similarity between two density functions is not symmetric and is hence not called a distance. It is called **Kulbach–Leibler divergence** and is given by

$$d^{\text{KL}}(p, q) = \int p(x) \log\left(\frac{p(x)}{q(x)}\right) dx \quad (2.24)$$

$$= \int p(x) \log(p(x)) dx - \int p(x) \log(q(x)) dx \quad (2.25)$$

This measure is zero if $p = q$. This measure is related to the information gain or relative entropy in information theory.

2.1.7 Density functions of multiple random variables

So far, we have discussed mainly probability (density) functions of single random variables. As mentioned before, we use random variables to describe data such as sensor readings in robots. Of course, we often have then more than one sensor and also other quantities that we describe by random variables at the same time. Thus, in many applications we consider multiple random variables. The quantities described by the random variables might be independent, but in many cases they are also related. Indeed, we will later talk about how to describe various types of relations. Thus, in order to talk about situations with multiple random variables, or multivariate statistics,

it is useful to know basic rules. We start by illustrating these basic multivariate rules with two random variables since the generalization from there is usually quite obvious. But we will also talk about the generalization to more than two variables at the end of this section.

2.1.7.1 Basic definitions

We have seen that probability theory is quite handy to model data, and probability theory also considers multiple random variables. The total knowledge about the co-occurrence of specific values for two random variables X and Y is captured by the

$$\text{joined distribution: } p(x, y) = p(X = x, Y = y). \quad (2.26)$$

This is a two dimensional functions. The two dimensions refers here to the number of variables, although a plot of this function would be a three dimensional plot. An example is shown in Fig.2.2. All the information we can have about a stochastic system is encapsulated in the joined pdf. The slice of this function, given the value of one variable, say y , is the

$$\text{conditional distribution: } p(x|y) = p(X = x|Y = y). \quad (2.27)$$

A conditional pdf is also illustrated in Fig.2.2 If we sum over all realizations of y we get the

$$\text{marginal distribution: } p(x) = \int p(x, y)dy. \quad (2.28)$$

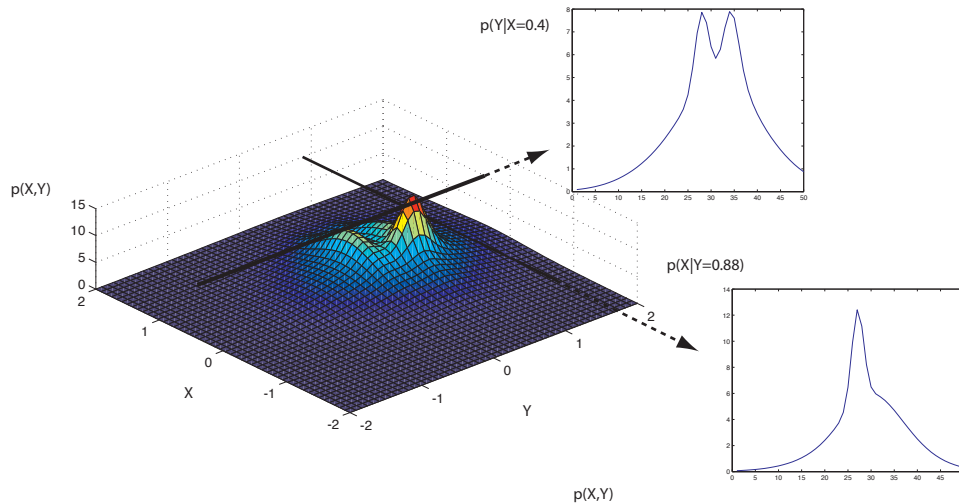


Fig. 2.2 Example of a two-dimensional probability density function (pdf) and some examples of conditional pdfs.

If we know some functional form of the density function or have a parameterized hypothesis of this function, than we can use common statistical methods, such as

maximum likelihood estimation, to estimate the parameters as in the one dimensional cases. If we do not have a parameterized hypothesis we need to use other methods, such as treating the problem as discrete and building histograms, to describe the density function of the system. Note that parameter-free estimation is more challenging with increasing dimensions. Considering a simple histogram method to estimate the joint density function where we discretize the space along every dimension into n bins. This leads to n^2 bins for a two-dimensional histogram, and n^d for a d -dimensional problem. This exponential scaling is a major challenge in practice since we need also considerable data in each bin to sufficiently estimate the probability of each bin.

2.1.7.2 The chain rule

As mentioned before, if we know the joint distribution of some random variables we can make the most predictions of these variables. However, in practice we have often to estimate these functions, and we can often only estimate conditional density functions. A very useful rule to know is therefore how a joint distribution can be decompose into the product of a conditional and a marginal distribution,

$$p(x, y) = p(x|y)p(y) = p(y|x)p(x), \quad (2.29)$$

which is sometimes called the **chain rule**. Note the two different ways in which we can decompose the joint distribution. This is easily generalized to n random variables by

$$p(x_1, x_2, \dots, x_n) = p(x_n|x_1, \dots, x_{n-1})p(x_1, \dots, x_{n-1}) \quad (2.30)$$

$$= p(x_n|x_1, \dots, x_{n-1}) * \dots * p(x_2|x_1) * p(x_1) \quad (2.31)$$

$$= \prod_{i=1}^n p(x_i|x_{i-1}, \dots, x_1) \quad (2.32)$$

but note that there are also different decompositions possible. We will learn more about this and useful graphical representations in Chapter ??.

Estimations of processes are greatly simplified when random variables are independent. A random variable X is independent of Y if

$$p(x|y) = p(x). \quad (2.33)$$

Using the chain rule eq.2.29, we can write this also as

$$p(x, y) = p(x)p(y), \quad (2.34)$$

that is, the joint distribution of two independent random variables is the product of their marginal distributions. Similar, we can also define conditional independence. For example, two random variables X and Y are conditionally independent of random variable Z if

$$p(x, y|z) = p(x|z)p(y|z). \quad (2.35)$$

Note that total independence does not imply conditional independence and visa versa, although this might hold true for some specific examples.

2.1.7.3 How to combine prior knowledge with new evidence: Bayes rule

One of the most common tasks we will encounter in the following is the integration of prior knowledge with new evidence. For example, we could have an estimate of

the location of an agent and get new (noisy) sensory data that adds some suggestions for different locations. A similar task is the fusion of data from different sensors. The general question we have to solve is how to weight the different evidence in light of the reliability of this information. Solving this problem is easy in a probabilistic framework and is one of the main reasons that so much progress has been made in probabilistic robotics.

How prior knowledge should be combined with prior knowledge is an important question. Luckily, basically already know how to do it best in a probabilistic sense. Namely, if we divide this chain rule eq. 2.29 by $p(x)$, which is possible as long as $p(x) > 0$, we get the identity

$$p(x|y) = \frac{p(y|x)p(x)}{p(y)}, \quad (2.36)$$

which is called **Bayes theorem**. This theorem is important because it tells us how to combine a **prior** knowledge, such as the expected distribution over a random variable such as the state of a system, $p(x)$, with some evidence called the likelihood function $p(y|x)$, for example by measuring some sensors reading y when controlling the state, to get the **posterior** distribution, $p(y|x)$ from which the new estimation of state can be derived. The marginal distribution $p(y)$, which does not depend on the state X , is the proper normalization so that the left-hand side is again a probability.

Exercises

1. Use your favourite plotting program to plot a Gaussian, a uniform, and the Chi-square distribution (probability density function). Include units on the axis.
2. Explain if the random variables X and Y are independent if their marginal distribution is $p(x) = 3x^2 + \log(x)$ and $p(y) = 3y^2 + \log(y)$ and the joined distributions is $p(x, y) = 3x^2y^2 + \log(xy)$.
3. (From Thrun, Burgard and Fox, Probabilistic Robotics) A robot uses a sensor that can measure ranges from $0m$ to $3m$. For simplicity, assume that the actual ranges are distributed uniformly in this interval. Unfortunately, the sensors can be faulty. When the sensor is faulty it constantly outputs a range below $1m$, regardless of the actual range in the sensor's measurement cone. We know that the prior probability for a sensor to be faulty is $p = 0.01$. Suppose the robot queries its sensors N times, and every single time the measurement value is below $1m$. What is the posterior probability of a sensor fault, for $N = 1, 2, \dots, 10$. Formulate the corresponding probabilistic model.

2.2 Vector and matrix notations

We frequently use vector and matrix notation in this book as it is extremely convenient for specifying neural network models. It is a shorthand notation for otherwise lengthy looking formulas, and formulas written in this notation can easily be entered into MATLAB. We consider three basic data types:

1. Scalar:

$$\mathbf{a} \text{ for example } 41 \quad (2.37)$$

2. Vector:

$$\mathbf{a} \text{ or component-wise } \begin{pmatrix} a_1 \\ a_2 \\ a_3 \end{pmatrix} \text{ for example } \begin{pmatrix} 41 \\ 7 \\ 13 \end{pmatrix} \quad (2.38)$$

3. Matrix:

$$\mathbf{a} \text{ or component-wise } \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \\ a_{31} & a_{32} \end{pmatrix} \text{ for example } \begin{pmatrix} 41 & 12 \\ 7 & 45 \\ 13 & 9 \end{pmatrix} \quad (2.39)$$

We used bold face to indicate both a vector and a matrix because the difference is usually apparent from the circumstances. A matrix is just a collection of scalars or vectors. We talk about an $n \times m$ matrix where n is the number of rows and m is the number of columns. A scalar is thus a 1×1 matrix, and a vector of length n can be considered an $n \times 1$ matrix. A similar collection of data is called **array** in computer science. However, a matrix is different because we also define operations on these data collections. The rules of calculating with matrices can be applied to scalars and vectors.

We define how to add and multiply two matrices so that we can use them in algebraic equations. The **sum of two matrices** is defined as the sum of the individual components

$$(\mathbf{a} + \mathbf{b})_{ij} = \mathbf{a}_{ij} + \mathbf{b}_{ij}. \quad (2.40)$$

For example, \mathbf{a} and \mathbf{b} are 3×2 matrices, then

$$\mathbf{a} + \mathbf{b} = \begin{pmatrix} a_{11} + b_{11} & a_{12} + b_{12} \\ a_{21} + b_{21} & a_{22} + b_{22} \\ a_{31} + b_{31} & a_{32} + b_{32} \end{pmatrix} \quad (2.41)$$

Matrix multiplication is defined as

$$(\mathbf{a} * \mathbf{b})_{ij} = \sum_k \mathbf{a}_{ik} \mathbf{b}_{kj}. \quad (2.42)$$

The matrix multiplication is hence only defined as multiplication matrices \mathbf{a} and \mathbf{b} where the number of columns of the matrix \mathbf{a} is equal to the number of rows of matrix \mathbf{b} . For example, for two square matrices with two rows and two columns, their product is given by

$$\mathbf{a} * \mathbf{b} = \begin{pmatrix} a_{11}b_{11} + a_{12}b_{21} & a_{11}b_{12} + a_{12}b_{22} \\ a_{21}b_{11} + a_{22}b_{21} & a_{21}b_{12} + a_{22}b_{22} \end{pmatrix} \quad (2.43)$$

A handy rule for matrix multiplications is illustrated in Fig. 2.3. Each component in the resulting matrix is calculated from the sum of two multiplicative terms. The rule for multiplying two matrices is tedious but straightforward and can easily be implemented in a computer. It is the default when multiplying variables of the matrix type in MATLAB. If we want to multiply each component of a matrix by the corresponding component in a second matrix, we just have to include the operator ‘.’ between

$$\begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} \begin{pmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{pmatrix} = \begin{pmatrix} a_{11}b_{11} + a_{12}b_{21} & a_{11}b_{12} + a_{12}b_{22} \\ a_{21}b_{11} + a_{22}b_{21} & a_{21}b_{12} + a_{22}b_{22} \end{pmatrix}$$

Fig. 2.3 Illustration of a matrix multiplication. Each element in the resulting matrix consists of terms that are taken from the corresponding row of the first matrix and column of the second matrix. Thus in the example we calculate the highlighted element from the components of the first row of the first matrix and the second column of the second matrix. From these rows and columns we add all the terms that consist of the element-wise multiplication of the terms.

the matrices in MATLAB where the dot in front of the multiplication sign indicates ‘component-wise’.

Another useful definition is the **transpose** of a matrix. This operation, indicated usually by a superscript t or a prime ($'$). The later is used in MATLAB. Taking the transpose of a matrix means that the matrix is rotated 90 degrees; the first row becomes the first column, the second row becomes the second column, etc. For example, the transpose of the example in 2.39 is

$$\mathbf{a}' = \begin{pmatrix} 41 & 7 & 13 \\ 12 & 45 & 9 \end{pmatrix} \quad (2.44)$$

The transpose of a vector transforms a column vector into a row vector and vice versa.

As already mentioned, matrices were invented to simplify the notations for systems of coupled algebraic equations. Consider, for example, the system of three equations

$$41x_1 + 12x_2 = 17 \quad (2.45)$$

$$7x_1 + 45x_2 = -83 \quad (2.46)$$

$$13x_1 + 9x_2 = -5. \quad (2.47)$$

This can be written as

$$\mathbf{a}\mathbf{x} = \mathbf{b} \quad (2.48)$$

with the matrix \mathbf{a} as in the example of 2.39, the vector $\mathbf{x} = (x_1 \ x_2)'$, and the vector $\mathbf{b} = (17 \ -83 \ -5)'$.

The solution of this linear equation system is equivalent to finding the inverse of matrix \mathbf{a} which we write as \mathbf{a}^{-1} . The inverse of the matrix is defined by

$$\mathbf{a}^{-1}\mathbf{a} = \mathbf{1}, \quad (2.49)$$

where the matrix $\mathbf{1}$ is the unit matrix that has element of one on the diagonal and zeros otherwise. Multiplying equation 2.48 from left with \mathbf{a}^{-1} is hence

$$\mathbf{x} = \mathbf{a}^{-1}\mathbf{b} \quad (2.50)$$

The inverse of a matrix, if this exists, can be found by the MATLAB function `inv()`.

2.3 Basic calculus

2.3.1 Differences and sums

We are often interested how a variable change with time. Let us consider the quantity $x(t)$ where we indicated that this quantity depends on time. The change of this variable from time t to time $t' = t + \Delta t$ is then

$$\Delta x = x(t + \Delta t) - x(t). \quad (2.51)$$

The quantity Δt is the finite difference in time. For a continuously changing quantity we could also think about the instantaneous change value, dx , by considering an infinitesimally small time step. Formally,

$$dx = \lim_{\Delta t \rightarrow 0} \Delta x = \lim_{\Delta t \rightarrow 0} (x(t + \Delta t) - x(t)). \quad (2.52)$$

The infinitesimally small time step is often written as dt . Calculating with such infinitesimal quantities is covered in the mathematical discipline of calculus, but on the computer we have always finite differences and we need to consider very small time steps to approximate continuous formulation. With discrete time steps, differential become differences and integrals become sums

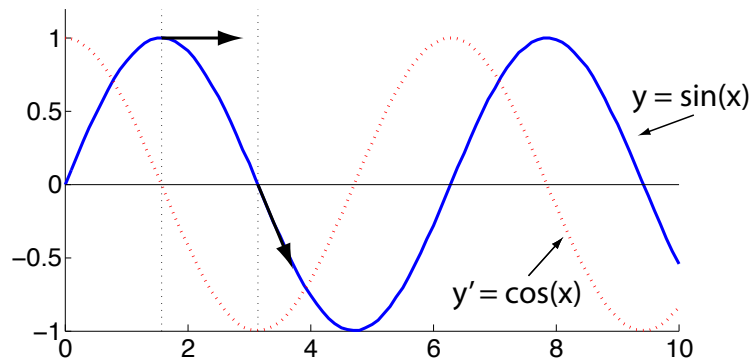
$$dx \rightarrow \Delta x \quad (2.53)$$

$$\int dx \rightarrow \Delta x \sum \quad (2.54)$$

Note the factor of Δx in front of the summation in the last equation. It is easy to forget this factor when replacing integrals with sums.

2.3.2 Derivatives

The derivative of a quantity y that depends on x is the slope of the function $y(x)$. This derivative can be defined as the limiting process equation 2.52 and is commonly written as $\frac{dy}{dx}$ or as y' .



It is useful to know some derivatives of basic functions.

$$y = e^x \rightarrow y' = e^x \quad (2.55)$$

$$y = \sin(x) \rightarrow y' = \cos(x) \quad (2.56)$$

$$y = x^n \rightarrow y' = nx^{n-1} \quad (2.57)$$

$$y = \log(x) \rightarrow \frac{1}{x} \quad (2.58)$$

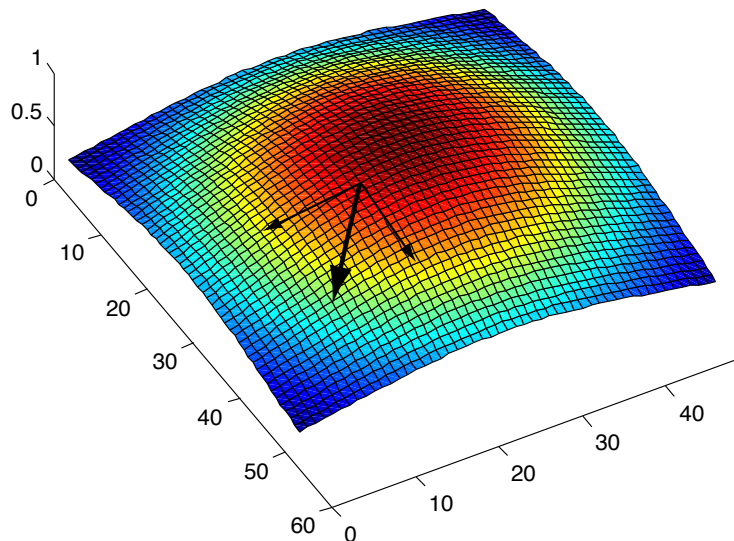
as well as the chain rule

$$y = f(x) \rightarrow y' = \frac{dy}{dx} = \frac{dy}{df} \frac{df}{dx}. \quad (2.59)$$

2.3.3 Partial derivative and gradients

A function that depends on more than one variable is a higher dimensional function. An example is the two-dimensional function $z(x, y)$. The slope of the function in the direction x (keeping y constant) is defined as $\frac{\partial z}{\partial x}$ and in the direction of y (keeping x constant) as $\frac{\partial z}{\partial y}$. The **gradient** is the vector that point in the direction of the maximal slope and has a length proportional to the slope,

$$\text{grad} z = \begin{pmatrix} \frac{\partial z}{\partial x} \\ \frac{\partial z}{\partial y} \end{pmatrix}. \quad (2.60)$$



3 Programming with Matlab

This chapter is a brief introduction to programming with the Matlab programming environment. We assume thereby little programming experience, although programmers experienced in other programming languages might want to scan through this chapter. MATLAB is an interactive programming environment for scientific computing. This environment is very convenient for us for several reasons, including its interpreted execution mode, which allows fast and interactive program development, advanced graphic routines, which allow easy and versatile visualization of results, a large collection of predefined routines and algorithms, which makes it unnecessary to implement known algorithms from scratch, and the use of matrix notations, which allows a compact and efficient implementation of mathematical equations and machine learning algorithms. MATLAB stands for matrix laboratory, which emphasizes the fact that most operations are array or matrix oriented. Similar programming environments are provided by the open source systems called **Scilab** and **Octave**. The Octave system seems to emphasize syntactic compatibility with MATLAB, while Scilab is a fully fledged alternative to MATLAB with similar interactive tools. While the syntax and names of some routines in Scilab are sometimes slightly different, the distribution includes a converter for MATLAB programs. Also, the Matlab web page provides great videos to learn how to use Matlab at <http://www.mathworks.com/demos/matlab/...getting-started-with-matlab-video-tutorial.html>.

3.1 The MATLAB programming environment

MATLAB⁴ is a programming environment and collection of tools to write programs, execute them, and visualize results. MATLAB has to be installed on your computer to run the programs mentioned in the manuscript. It is commercially available for many computer systems, including Windows, Mac, and UNIX systems. The MATLAB web page includes a set of brief tutorial videos, also accessible from the **demos** link from the MATLAB desktop, which are highly recommended for learning MATLAB.

As already mentioned, there are several reasons why MATLAB is easy to use and appropriate for our programming need. MATLAB is an interpreted language, which means that commands can be executed directly by an interpreter program. This makes the time-consuming compilation steps of other programming languages redundant and allows a more interactive working mode. A disadvantage of this operational mode is that the programs could be less efficient compared to compiled programs. However, there are two possible solutions to this problem in case efficiency becomes a concern. The first is that the implementations of many MATLAB functions is very efficient

⁴MATLAB and Simulink are registered trademarks, and MATLAB Compiler is a trademark of The MathWorks, Inc.

and are themselves pre-compiled. MATLAB functions, specifically when used on whole matrices, can therefore outperform less well-designed compiled code. It is thus recommended to use matrix notations instead of explicit component-wise operations whenever possible. A second possible solution to increase the performance is to use the MATLAB compiler to either produce compiled MATLAB code in `.mex` files or to translate MATLAB programs into compilable language such as C.

A further advantage of MATLAB is that the programming syntax supports matrix notations. This makes the code very compact and comparable to the mathematical notations used in the manuscript. MATLAB code is even useful as compact notation to describe algorithms, and it is hence useful to go through the MATLAB code in the manuscript even when not running the programs in the MATLAB environment. Furthermore, MATLAB has very powerful visualization routines, and the new versions of MATLAB include tools for documentation and publishing of codes and results. Finally, MATLAB includes implementations of many mathematical and scientific methods on which we can base our programs. For example, MATLAB includes many functions and algorithms for linear algebra and to solve systems of differential equations. Specialized collections of functions and algorithms, called a ‘toolbox’ in MATLAB, can be purchased in addition to the basic MATLAB package or imported from third parties, including many freely available programs and tools published by researchers. For example, the MATLAB Neural Network Toolbox incorporates functions for building and analysing standard neural networks. This toolbox covers many algorithms particularly suitable for connectionist modelling and neural network applications. A similar toolbox, called NETLAB, is freely available and contains many advanced machine learning methods. We will use some toolboxes later in this course, including the LIBSVM toolbox and the MATLAB NXT toolbox to program the Lego robots.

Starting MATLAB opens the MATLAB desktop as shown in Fig. 3.1 for MATLAB version 7. The MATLAB desktop is comprised of several windows which can be customized or undocked (moving them into an own window). A list of these tools are available under the **desktop menu**, and includes tools such as the **command window**, **editor**, **workspace**, etc. We will use some of these tools later, but for now we only need the **MATLAB command window**. We can thus close the other windows if they are open (such as the **launch pad** or the **current directory window**); we can always get them back from the **desktop** menu. Alternatively, we can undock the command window by clicking the arrow button on the command window bar. An undocked command window is illustrated on the left in Fig. 3.2. Older versions of MATLAB start directly with a command window or simply with a MATLAB command prompt `>>` in a standard system window. The command window is our control centre for accessing the essential MATLAB functionalities.

3.2 Main programming constructs

3.2.1 Basic variables in MATLAB

The MATLAB programming environment is interactive in that all commands can be typed into the command window after the command prompt (see Fig. 3.2). The

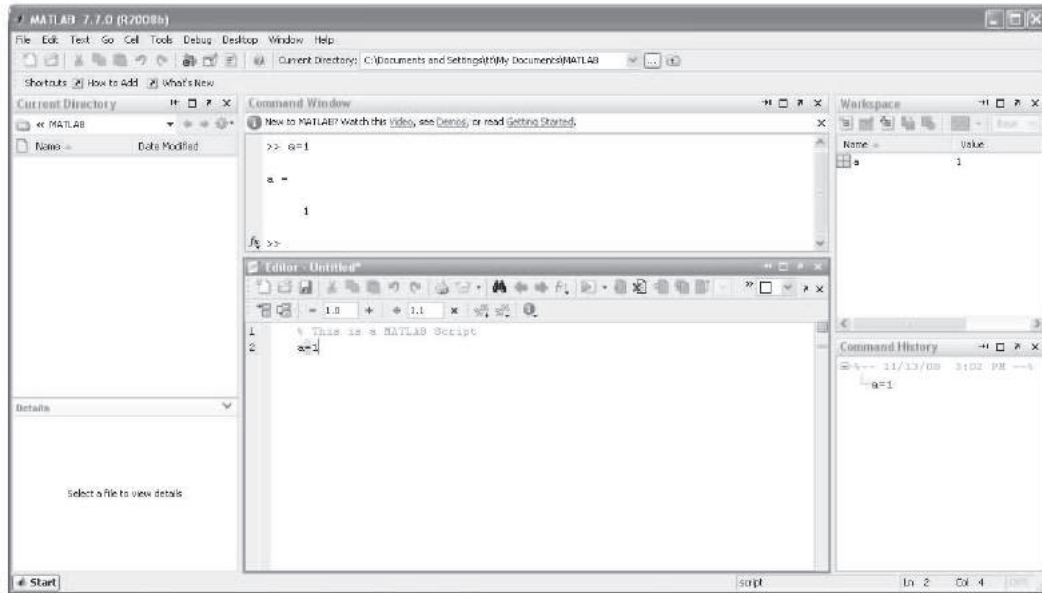


Fig. 3.1 The MATLAB *desktop window* of MATLAB Version 7.

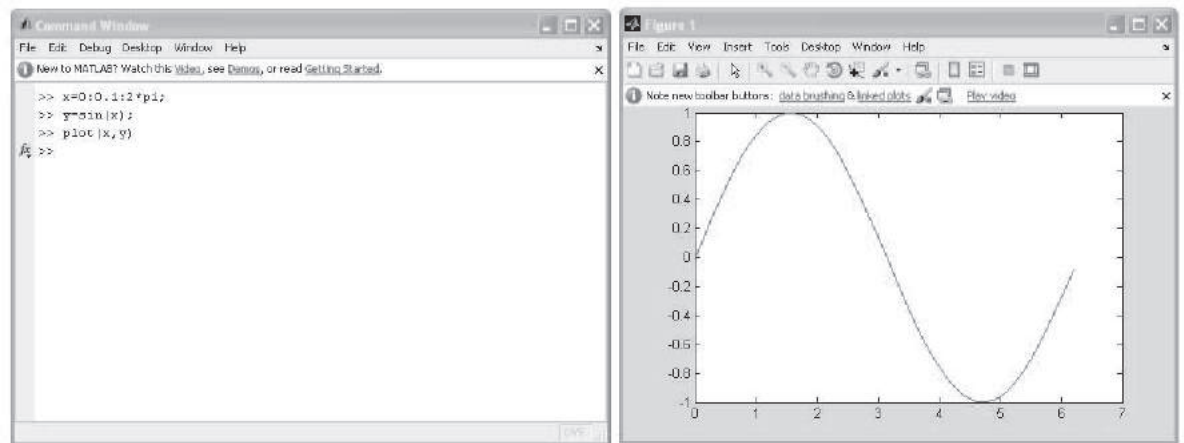


Fig. 3.2 A MATLAB *command window* (left) and a MATLAB *figure window* (right) displaying the results of the function `plot_sin` developed in the text.

commands are interpreted directly, and the result is returned to (and displayed in) the command window. For example, a variable is created and assigned a value with the = operator, such as

```
>> a=3
```

```
a =
```

```
3
```

Ending a command with semicolon (;) suppresses the printing of the result on screen. It is therefore generally included in our programs unless we want to view some intermediate results. All text after a percentage sign (%) is not interpreted and thus treated as comment,

```
>> b='Hello World!'; % delete the semicolon to echo Hello World!
```

This example also demonstrates that the type of a variable, such as being an integer, a real number, or a string, is determined by the values assigned to the elements. This is called **dynamic typing**. Thus, variables do not have to be declared as in some other programming languages. While dynamic typing is convenient, a disadvantage is that a mistyped variable name can not be detected by the compiler. Inspecting the list of created variables is thus a useful step for debugging.

All the variables that are created by a program are kept in a buffer called **workspace**. These variable can be viewed with the command `whos` or displayed in the **workspace** window of the MATLAB desktop. For example, after declaring the variables above, the `whos` command results in the responds

```
>> whos
  Name      Size      Bytes  Class  Attributes
  a         1x1         8  double
  b         1x12        24  char
```

It displays the name, the size, and the type (class) of the variable. The size is often useful to check the orientation of matrices as we will see below. The variables in the workspace can be used as long as MATLAB is running and as long as it is not cleared with the command `clear`. The workspace can be saved with the command `save filename`, which creates a file `filename.mat` with internal MATLAB format. The saved workspace can be reloaded into MATLAB with the command `load filename`. The load command can also be used to import data in ASCII format. The workspace is very convenient as we can run a program within a MATLAB session and can then work interactively with the results, for example, to plot some of the generated data.

Variables in MATLAB are generally matrices (or data arrays), which is very convenient for most of our purposes. Matrices include scalars (1×1 matrix) and vectors ($1 \times N$ matrix) as special cases. Values can be assigned to matrix elements in several ways. The most basic one is using square brackets and separating rows by a semicolon within the square brackets, for example (see Fig. 3.2),

```
>> a=[1 2 3; 4 5 6; 7 8 9]
```

```
a =
```

```
1     2     3
4     5     6
7     8     9
```


A vector of elements with consecutive values can be assigned by column operators like

```
>> v=0:2:4
```

```
v =
```

```
    0    2    4
```

Furthermore, the MATLAB desktop includes an **array editor**, and data in ASCII files can be assigned to matrices when loaded into MATLAB. Also, MATLAB functions often return matrices which can be used to assign values to matrix elements. For example, a uniformly distributed random 3×3 matrix can be generated with the command

```
>> b=rand(3)
```

```
b =
```

```
    0.9501    0.4860    0.4565
    0.2311    0.8913    0.0185
    0.6068    0.7621    0.8214
```

The multiplication of two matrices, following the matrix multiplication rules, can be done in MATLAB by typing

```
>> c=a*b
```

```
c =
```

```
    3.2329    4.5549    2.9577
    8.5973   10.9730    6.8468
   13.9616   17.3911   10.7360
```

This is equivalent to

```
c=zeros(3);
for i=1:3
    for j=1:3
        for k=1:3
            c(i,j)=c(i,j)+a(i,k)*b(k,j);
        end
    end
end
```

which is the common way of writing matrix multiplications in other programming languages. Formulating operations on whole matrices, rather than on the individual components separately, is not only more convenient and clear, but can enhance the programs performance considerably. Whenever possible, operations on whole matrices should be used. This is likely to be the major change in your programming style when converting from another programming language to MATLAB. The performance disadvantage of an interpreted language is often negligible when using operations on

whole matrices.

The transpose operation of a matrix changes columns to rows. Thus, a row vector such as v can be changed to a column vector with the MATLAB transpose operator (`'`),

```
>> v'
```

```
ans =
```

```
0
2
4
```

which can then be used in a matrix-vector multiplication like

```
>> a*v'
```

```
ans =
```

```
16
34
52
```

The inconsistent operation $a*v$ does produce an error,

```
>> a*v
```

```
??? Error using ==> mtimes
```

```
Inner matrix dimensions must agree.
```

Component-wise operations in matrix multiplications (`*`), divisions (`/`) and potentiation `^` are indicated with a dot modifier such as

```
>> v.^2
```

```
ans =
```

```
0    4   16
```

The most common operators and basic programming constructs in MATLAB are similar to those in other programming languages and are listed in Table 3.1.

3.2.2 Control flow and conditional operations

Besides the assignments of values to variables, and the availability of basic data structures such as arrays, a programming language needs a few basic operations for building loops and for controlling the flow of a program with conditional statements (see Table 3.1). For example, the **for loop** can be used to create the elements of the vector v above, such as

```
>> for i=1:3; v(i)=2*(i-1); end
```

```
>> v
```

```
v =
```

Table 3.1 Basic programming constructs in MATLAB.

Programming construct	Command	Syntax
Assignment	=	a=b
Arithmetic operations	add	a+b
	multiplication	a*b (matrix), a.*b (element-wise)
	division	a./b (matrix), a./b (element-wise)
	power	a.^b (matrix), a.^b (element-wise)
Relational operators	equal	a==b
	not equal	a~=b
	less than	a<b
Logical operators	AND	a & b
	OR	a b
Loop	for	for index=start:increment:end statement end
	while	while expression statement end
Conditional command	if statement	if logical expressions statement elseif logical expressions statement else statement end
Function		function [x,y,...]= name (a,b,...)

```
0    2    4
```

Table 3.1 lists, in addition, the syntax of a while loop. An example of a conditional statement within a loop is

```
>> for i=1:10; if i>4 & i<=7; v2(i)=1; end; end
>> v2
```

```
v2 =
```

```
0    0    0    0    1    1    1
```

In this loop, the statement `v2(i)=1` is only executed when the loop index is larger than 4 and less or equal to 7. Thus, when `i=5`, the array `v2` with 5 elements is created, and since only the elements `v2(5)` is set to 1, the previous elements are set to 0 by default. The loop adds then the two element `v2(6)` and `v2(7)`. Such a vector can also be created by assigning the values 1 to a specified range of indices,

```
>> v3(4:7)=1
```

```
v3 =
```

```
0 0 0 1 1 1 1
```

A 1×7 array is thereby created with elements set to 0, and only the specified elements are overwritten with the new value 1. Another method to write compact loops in MATLAB is to use vectors as index specifiers. For example, another way to create a vector with values such as `v2` or `v3` is

```
>> i=1:10
```

```
i =
```

```
1 2 3 4 5 6 7 8 9 10
```

```
>> v4(i>4 & i<=7)=1
```

```
v4 =
```

```
0 0 0 0 1 1 1
```

3.3 Creating MATLAB programs

If we want to repeat a series of commands, it is convenient to write this list of commands into an ASCII file with extension `.m`. Any ASCII editor (for example; WordPad, Emacs, etc.) can be used. The MATLAB package contains an editor that has the advantage of colouring the content of MATLAB programs for better readability and also provides direct links to other MATLAB tools. The list of commands in the ASCII file (e.g. `prog1.m`) is called a **script** in MATLAB and makes up a MATLAB program. This program can be executed with a run button in the MATLAB editor or by calling the name of the file within the command window (for example, by typing `prog1`). We assumed here that the program file is in the current directory of the MATLAB session or in one of the search paths that can be specified in MATLAB. The MATLAB desktop includes a ‘current directory’ window (see desktop menu). Some older MATLAB versions have instead a ‘path browser’. Alternatively, one can specify absolute path when calling a program, or change the current directories with UNIX-style commands such as `cd` in the command window (see Fig. 3.3).

Functions are another way to encapsulate code. They have the additional advantage that they can be pre-compiled with the MATLAB Compiler™ available from MathWorks, Inc. Functions are kept in files with extension `.m` which start with the command line like

```
function y=f(a,b)
```

where the variables `a` and `b` are passed to the function and `y` contains the values returned by the function. The return values can be assigned to a variable in the calling MATLAB



Fig. 3.3 Two editor windows and a command window.

script and added to the workspace. All other internal variables of the functions are local to the function and will not appear in the workspace of the calling script.

MATLAB has a rich collection of predefined functions implementing many algorithms, mathematical constructs, and advanced graphic handling, as well as general information and help functions. You can always search for some keywords using the useful command `lookfor` followed by the keyword (search term). This command lists all the names of the functions that include the keywords in a short description in the function file within the first **comment lines** after the function declaration in the function file. The command `help`, followed by the function name, displays the first block of comment lines in a function file. This description of functions is usually sufficient to know how to use a function. A list of some frequently used functions is listed in Table 3.3.

3.4 Graphics

MATLAB is a great tool for producing scientific graphics. We want to illustrate this by writing our first program in MATLAB: calculating and plotting the sine function. The program is

```

x=0:0.1:2*pi;
y=sin(x);
plot(x,y)

```

Name	Brief description	Name	Brief description
abs	absolute functions	mod	modulus function
axis	sets axis limits	num2str	converts number to string
bar	produces bar plot	ode45	ordinary differential equation solver
ceil	round to larger interger	ones	produces matrix with unit elements
colormap	colour matrix for surface plots	plot	plot lines graphs
cos	cosine function	plot3	plot 3-dimensional graphs
diag	diagonal elements of a matrix	prod	product of elements
disp	display in command window	rand	uniformly distributed random variable
errorbar	plot with error bars	randn	normally distributed random variable
exp	exponential function	randperm	random permutations
fft	fast Fourier transform	reshape	reshaping a matrix
find	index of non-zero elements	set	sets values of parameters in structure
floor	round to smaller integer	sign	sign function
hist	produces histogram	sin	sine function
int2str	converts integer to string	sqrt	square root function
isempty	true if array is empty	std	calculates standard deviation
length	length of a vector	subplot	figure with multiple subfigures
log	logarithmic function	sum	sum of elements
lsqcurvefit	least mean square curve fitting (statistics toolbox)	surf	surface plot
max	maximum value and index	title	writes title on plot
mix	minimum value and index	view	set viewing angle of 3D plot
mean	calculates mean	xlabel	label on x-axis of a plot
meshgrid	creates matrix to plot grid	ylabel	label on y-axis of a plot
		zeros	creates matrix of zero elements

Table 3.2 MATLAB functions used in this course. The MATLAB command `help cmd`, where `cmd` is any of the functions listed here, provides more detailed explanations.

The first line assigns elements to a vector x starting with $x(1) = 0$ and incrementing the value of each further component by 0.1 until the value 2π is reached (the variable `pi` has the appropriate value in MATLAB). The last element is $x(63) = 6.2$. The second line calls the MATLAB function `sin` with the vector x and assigns the results to a vector y . The third line calls a MATLAB plotting routine. You can type these lines into an ASCII file that you can name `plot_sin.m`. The code can be executed by typing `plot_sin` as illustrated in the **command window** in Fig. 3.2, provided that the MATLAB session points to the folder in which you placed the code. The execution of this program starts a **figure window** with the plot of the sine function as illustrated on the right in Fig. 3.2.

The appearance of a plot can easily be changed by changing the attributes of the plot. There are several functions that help in performing this task, for example, the function `axis` that can be used to set the limits of the axis. New versions of MATLAB provide window-based interfaces to the attributes of the plot. However, there are also two basic commands, `get` and `set`, that we find useful. The command `get(gca)` returns a list with the axis properties currently in effect. This command is useful for finding out what properties exist. The variable `gca` (get current axis) is the **axis handle**, which is a variable that points to a memory location where all the attribute variables are

kept. The attributes of the axis can be changed with the `set` command. For example, if we want to change the size of the labels we can type `set(gca, 'fontsize', 18)`. There is also a handle for the current figure `gcf` that can be used to get and set other attributes of the figure. MATLAB provides many routines to produce various special forms of plots including plots with error-bars, histograms, three-dimensional graphics, and multi-plot figures.

3.5 A first project: modelling the world

Suppose there is a simple world with a creature that can be in three distinct states, sleep (state value 1), eat (state value 2), and study (state value 3). An agent, which is a device that can sense environmental states and can generate actions, is observing this creature with poor sensors, which add white (Gaussian) noise to the true state. Our aim is to build a model of the behaviour of the creature which can be used by the agent to observe the states of the creature with some accuracy despite the limited sensors. For this exercise, the function `creature_state()` is available on the course page on the web. This function returns the current state of the creature. Try to create an agent program that predicts the current state of the creature. In the following we discuss some simple approaches.

A simulation program that implements a specific agent `a` with simple world model (a model of the creature), which also evaluates the accuracy of the model, is given in Table 3.3. This program, also available on the web, is provided in file `main.m`. This program can be downloaded into the working directory of MATLAB and executed by typing `main` into the command window, or by opening the file in the MATLAB editor and starting it from there by pressing the icon with the green triangle. The program reports the percentage of correct perceptions of the creature's state.

Line 1 of the program uses a comment indicator (%) to outline the purpose of the program. Line 2 clears the workspace to erase all eventual existing variables, and sets a counter for the number of correct perceptions to zero. Line 4 starts a loop over 1000 trials. In each trial, a creature state is pulled by calling the function `creature_state()` and recording this state value in variable `x`. The sensory state `s` is then calculated by adding a random number to this value. The value of the random number is generated from a normal distribution, a Gaussian distribution with mean zero and unit variance, with the MATLAB function `randn()`.

We are now ready to build a model for the agent to interpret the sensory state. In the example shown, this model is given in Lines 8–12. This model assumes that a sensory value below 1.5 corresponds to the state of a sleeping creature (Line 9), a sensory value between 1.5 and 2.5 corresponds to the creature eating (Line 10), and a higher value corresponds to the creature studying (Line 11). Note that we made several assumptions by defining this model, which might be unreasonable in real-world applications. For example, we used our knowledge that there are three states with ideal values of 1, 2, and 3 to build the model for the agent. Furthermore, we used the knowledge that the sensors are adding independent noise to these states in order to come up with the decision boundaries. The major challenge for real agents is to build models without this explicit knowledge. When running the program we find that a little bit over 50%

Table 3.3 Program main.m

```

1  % Project 1: simulation of agent which models simple creature
2  clear; correct=0;
3
4  for trial=1:1000
5      x=creature_state();
6      s=x+randn();
7
8      %% perception model
9      if (s<1.5) x_predict=1;
10     elseif (s<2.5) x_predict=2;
11     else x_predict=3;
12     end
13
14     %% calculate accuracy
15     if (x==x_predict) correct=correct+1; end
16 end
17
18 disp(['percentage correct: ',num2str(correct/1000)]);

```

of the cases are correctly perceived by the agent. While this is a good start, one could do better. Try some of your own ideas . . .

. . . Did you succeed in getting better results? It is certainly not easy to guess some better model, and it is time to inspect the data more carefully. For example, we can plot the number of times each state occurs. For this we can write a loop to record the states in a vector,

```
>> for i=1:1000; a(i)=creature_state(); end
```

and then plot a histogram with the MATLAB function `hist()`,

```
>> hist(a)
```

The result is shown in Fig. 3.4. This histogram shows that not all states are equally likely as we implicitly assumed in the above agent model. The third state is indeed much less likely. We could use this knowledge in a modified model in which we predict that the agent is sleeping for sensory states less than 1.5 and is eating otherwise. This modified model, which completely ignores study states, predicts around 65% of the states correctly. Many machine learning methods suffer from such ‘explaining away’ solutions for imbalanced data, as further discussed in Chapter ??.

It is important to recognize that 100% accuracy is not achievable with the inherent limitations of the sensors. However, higher recognition rates could be achieved with better world (creature + sensor) models. The main question is how to find such a model. We certainly should use observed data in a better way. For example, we could use several observations to estimate how many states are produced by function `creature_state()` and their relative frequency. Such parameter estimation is a basic form of learning from data. Many models in science take such an approach by proposing a parametric model and estimating parameters from the data by model fitting. The main

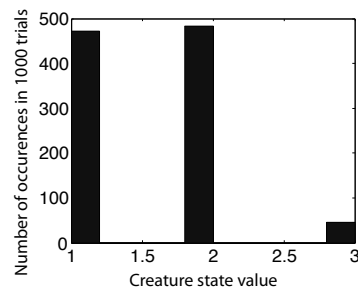


Fig. 3.4 The MATLAB *desktop window* histogram of states produced by function `creature_state()` from 1000 trials.

challenge with this approach is how complex we should make the model. It is much easier to fit a more complex model with many parameters to example data, but the increased flexibility decreases the prediction ability of such models. Much progress has been made in machine learning by considering such questions, but those approaches only work well in limited worlds, certainly much more restricted than the world we live in. More powerful methods can be expected by learning how the brain solves such problems.

3.6 Alternative programming environments: Octave and Scilab

We briefly mention here two programming environments that are very similar to Matlab and that can, with certain restrictions, execute Matlab scripts. Both of these programming systems are open source environments and have general public licenses for non-commercial use.

The programming environment called **Octave** is freely available under the GNU general public license. Octave is available through links at <http://www.gnu.org/software/octave/>. The installation requires the additional installation of a graphics package, such as `gnuplot` or `Java graphics`. Some distributions contain the `SciTE` editor which can be used in this environment. An example of the environment is shown in Fig. 3.5

Scilab is another scientific programming environment similar to MATLAB. This software package is freely available under the CeCILL software license, a license compatible to the GNU general public license. It is developed by the Scilab consortium, initiated by the French research centre INRIA. The Scilab package includes a MATLAB import facility that can be used to translate MATLAB programs to Scilab. A screen shot of the Scilab environment is shown in Fig. 3.6. A Scilab script can be run from the execute menu in the editor, or by calling `exec("filename.sce")`.

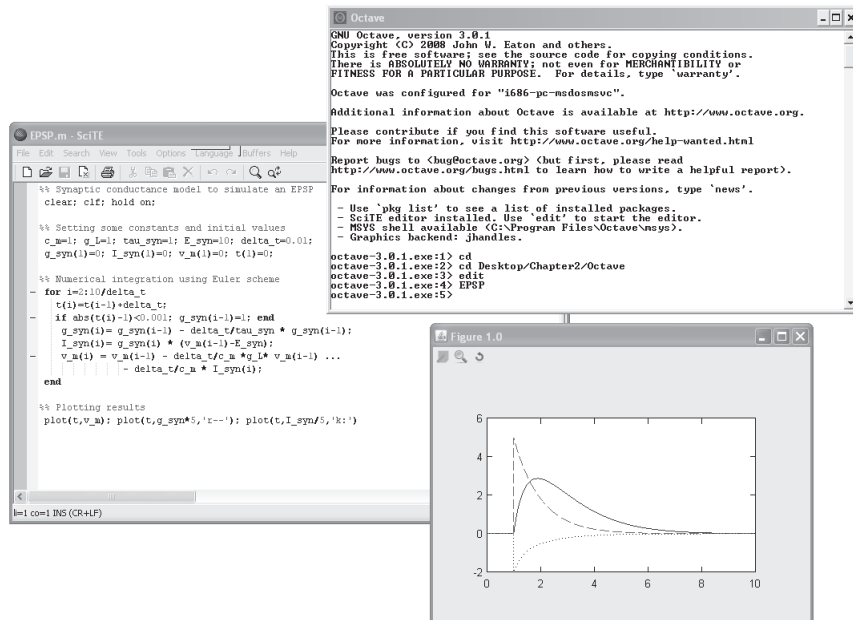


Fig. 3.5 The Octave programming environment with the main console, and editor called *SciTE*, and a graphics window.

Exercises

1. Write a Matlab function that takes a character string and prints out the character string in reverse order.
2. Write a Matlab program that uses 3-dimensional plotting routines to plot a two dimensional Gaussian function.
3. Following up from the 3rd question of the last assignment (the faulty sensor question), use Matlab to plot the believe that the sensor is faulty against the number of measurements.

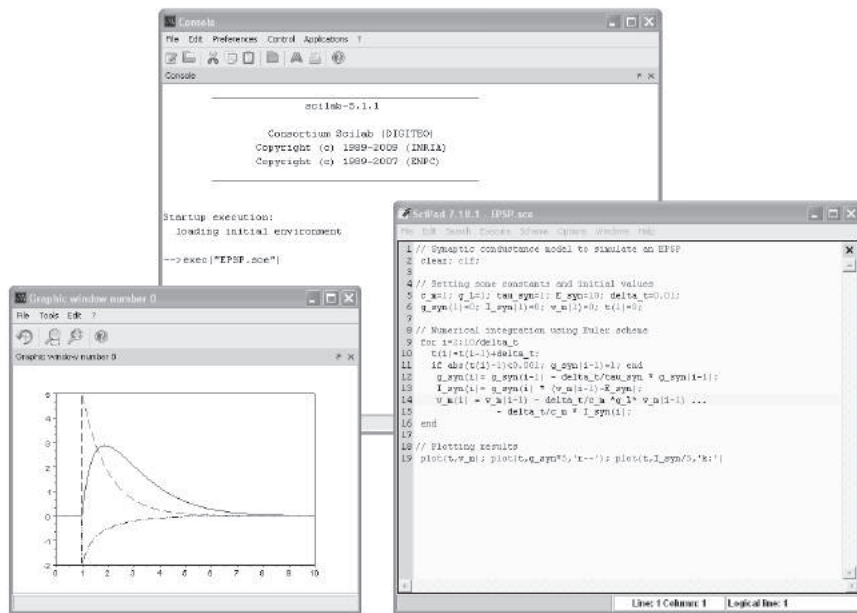


Fig. 3.6 The Scilab programming environment with *console*, and editor called *SciPad*, and a graphics window.

4 Basic robotics with LEGO NXT

4.1 Building a basic robot with LEGO NXT

We will use the LEGO mindstorm system in this course. Besides having common LEGO building blocks to construct different designs, this system consist of a micro-processor, called the **brick**, which is programmable and which controls the sensors and actuators. The actuators are stepping motors that can be told to run for a specific duration, a specific number of rotations, or with a specific speed. Our tool kit also includes several sensors, a light sensor that can be used to measure the wavelength of reflecting light and also small distances, an ultrasonic sensor to measure larger distances, a touch sensor, and a microphone. The motors can also be used to sense some externally applied movements.

We will use a basic tribot design as shown in Fig.4.1 for most of the explorations in this course. We will build the basic tribot in the following tutorial and will outline an example of a program to avoids walls and then work on a program so that the tribot can follow a line. After this we will use a robot arm configuration to discuss the important concept of a configuration space and demonstrate path-planing.



Fig. 4.1 Basic Lego Robot with microprocessor, two motors, and a light sensor.

4.1.1 Building the tribot

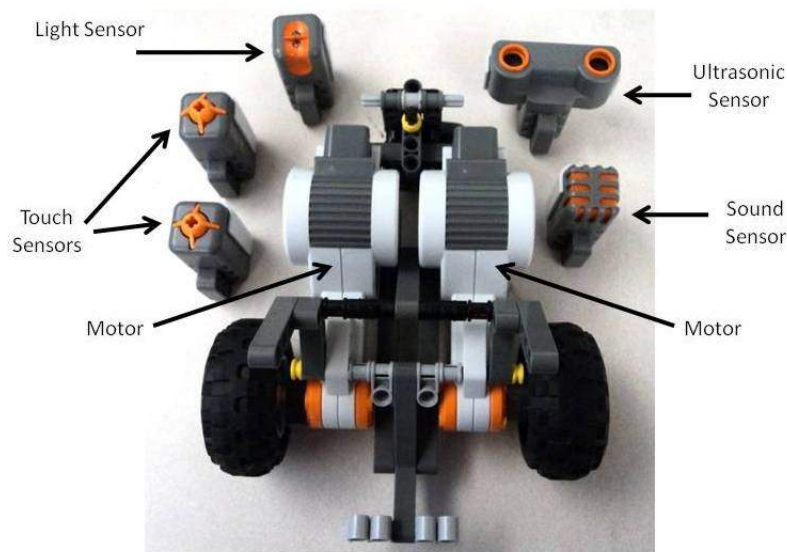
1. **Charge the battery**

Before constructing the robot, plug the battery pack into the NXT brick controller and plug it into an outlet to recharge the battery.

2. **Construct the NXT**

Follow the instruction booklet included with the Lego kit to construct the robot. There are mini programs illustrated in the manual that can be implemented directly on the NXT brick. Although this is not the method that will be used to control the NXT, you may implement them if you have time in order to get an idea of how the sensors and motors function. The instruction manual is organized as follows:

Section	Pages
NXT brick base construction	8-23
Sound sensor mounting	24-27
Ultrasonic sensor mounting	28-31
Light sensor mounting	32-39
Touch sensor mounting	40-45



Sensors & constructed NXT brick base

4.1.2 Mindstorms NXT toolbox installation

We will use some software to control the Lego actuator and gather information from their sensors within the Matlab programming environment. To enable this we need to install software developed at the German university called 'RWTH Aachen', which in turn uses some other drivers that we need to install. Most of the software should be installed in our Lab, but we will outline briefly some of the installation issues in case you want to install them under your own system or if some problems exist with the current installation. The following software installation instructions are adapted from RWTH Aachen University's NXT Toolbox website:

<http://www.mindstorms.rwth-aachen.de/trac/wiki/Download4.04>

Currently, the toolbox is not available for USB connections are not available on Mac because of 64-bit compatibility issues; therefore, you will only be able to connect to the NXT via Bluetooth if you are using a Mac.

1. Check NXT Firmware version

Check what version of NXT Firmware is running on the NXT brick by going to "Settings" > "NXT Version". Firmware version ("FW") should be 1.28 or 1.29. If it does not, it needs to be updated (Note: The NXT toolbox website claims version 1.26 will work, however it will not, and 1.31 has not worked in previous tests)

To update the firmware:

The Lego Mindstorms Education NXT Programming software is required to update the firmware. In the NXT Programming software, look under "tools" > "Update NXT Firmware" > "browse", select the firmware's directory, click "download".

2. Setting up a USB Connection [Windows Users only]

The toolbox requires different installation instructions depending on whether you are running 32 or 64-bit Matlab. To check which version you are running, in Matlab go to "Help" > "About".

On 32-bit Matlab [Windows Users only]

- **Install USB (Fantom) Driver**

If the Lego Mindstorms Education NXT Programming software is already on your computer, this should already be installed. Otherwise, download it from: <http://mindstorms.lego.com/support/updates/>

- * If you run into problems with the Fantom Library on windows go to this site for help:

<http://bricxcc.sourceforge.net/NXTFantomDriverHelp.pdf>

- * If you have Windows 7 Starter edition the standard setup file will not run properly. To install the Fantom Driver go into Products and then LEGO_NXT_Driver_32 and run LegoMindstormsNXTdriver32.

- **Download the Mindstorms NXT Toolbox 4.04:**

Download: <http://www.mindstorms.rwth-aachen.de/trac/wiki/Download4.04>

- * Save and extract the files anywhere, but do not change the directory structure.

- * The folder will appear as "RWTHMindstormsNXT"

On 64-bit Matlab [Windows Users only]

- **Download libusb-win32**

Download: <http://sourceforge.net/projects/libusb-win32/files/>

- * Save and extract the files anywhere.

- * Open the extracted folder and go to "bin". Run the **inf-wizard** and follow the instructions. This will install a device driver for the NXT connection.

- **Download the Mindstorms NXT Toolbox 4.04:**

Download: [http://projects.cs.dal.ca/hallab/wiki/index.php/RWTH_Mindstorms_NXT_Toolbox_\(64-bit\)](http://projects.cs.dal.ca/hallab/wiki/index.php/RWTH_Mindstorms_NXT_Toolbox_(64-bit))

- * Save and extract the files anywhere, but do not change the directory structure.

- * The folder will appear as "RWTHMindstormsNXT64"

3. Install NXT Toolbox into Matlab

In Matlab: "File" > "SetPath" > "Add Folder", and browse and select "RWTH-MindstormsNXT" or "RWTHMindstormsNXT64" - the which is a **subfolder** of file you saved in the previous step.

- Also add the "tools" folder, which is a subdirectory of the RWTHMindstormsNXT folder.
- Click "save" when finished.

4. Download MotorControl to NXT brick

Use the USB cable for this step. Lego's Fantom Driver is also required for this step. This will already be installed if you have the Lego Mindstorms Education NXT Programming software or if you downloaded it when setting up a USB connection with 32-bit Matlab in Windows. Otherwise, download it from: <http://mindstorms.lego.com/en-us/support/files/default.aspx#Driver>

- **Windows:** Download the NBC compiler (<http://bricxcc.sourceforge.net/nbc/>). Unzip the folder and move the file "nbc.exe" to RWTHMindstormsNXT/tools/MotorControl. Under RWTHMindstormsNXT/tools/MotorControl, double click TransferMotorControlBinaryToNXT, and follow the onscreen instructions. If this fails, try using NeXTTool instead of nbc; download from <http://bricxcc.sourceforge.net/utilities.html> (be careful to download NeXTTool not NextTools!). Again, unzip the folder and move the file "NeXTTool.exe" to the MotorControl folder.
- **Mac:** Download the NeXT Tools for Mac OS X from <http://bricxcc.sourceforge.net/utilities.html>. Run the toolbar and open the NXT Explorer (the globe in the toolbar). With the arrow key at the top, transfer the file MotorControl22.rxe (found in RWTHMindstormsNXT/tools/MotorControl) to the brick.

5. Setting up a Bluetooth connection [optional if you're connecting via USB]

- To connect to the NXT via bluetooth you must first turn on the bluetooth in the NXT and make sure that the visibility is set to on. Then use the bluetooth device on your computer to search for your specific NXT. The name of your NXT can be found at the top center of your NXT's screen.
- Create a connection between the computer and the NXT. When you create the connection between the NXT and the bluetooth device the NXT will ask for a passkey (usually either 0000 or 1234 on the NXT screen and press the orange button. The computer will then ask for the same passkey. To test the connection, type the command `COM_OpenNXT('bluetooth.ini');` in the Matlab command window. The command should run without any red error messages.
- Make sure that the bluetooth.ini file is present. There are sample files for Windows and Linux (Mac) in the main RWTH toolbox folder. If it is not present, create one by running the command `COM_MakeBTConfigFile` in Matlab.
- Check if the **serial ports** are correct in the "bluetooth.ini" file are correct. In **Windows**, go to the Control Panel, click "Hardware and Sound">"Devices and Printers". Find the NXT device, right click it and go to "Properties". The COMPort will be listed under the "Services" tab (i.e. "COM3"). On **Mac**, The port name can be found by typing `ls -ltr /dev` in a terminal

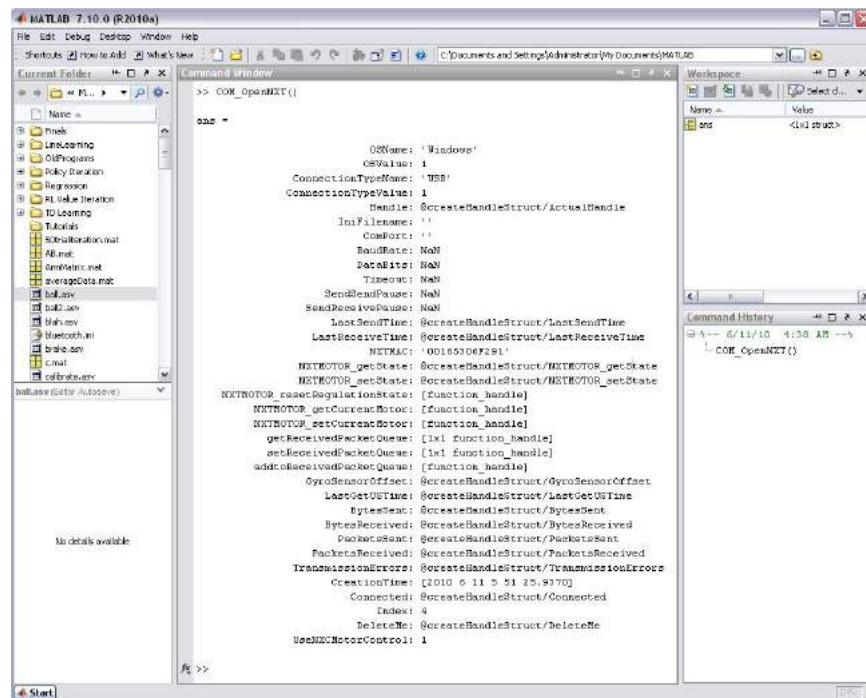
window, or by going to "System Preferences">"Bluetooth Devices". Select your NXT and click on the gear symbol (left hand bottom corner). The COMPort should be listed there.

- Turning the NXT off and back on again can help. After every failed `COM_OpenNXT('bluetooth.ini')`; command type `COM_CloseNXT('all')`; to close the failed connection for a clean new attempt.
- To switch on a debug mode enter the command `DebugMode on` before entering the command `COM_OpenNXT('bluetooth.ini');` .

6. Does it work?

In Matlab, enter the commands below into the command window. The command should execute without error and the NXT should play a sound.

```
h=COM_OpenNXT(); COM_SetDefaultNXT(h);
NXT_PlayTone(400,300);
```



`h=COM_OpenNXT()`; To test a USB connection you need to include an argument in the `COM_OpenNXT()` function like `COM_OpenNXT('bluetooth.ini')` `COM_SetDefaultNXT(h)`; Note that there are some examples included with the RWTH Mindstorm's NXT Toolbox, under `RWTHMindstormsNXT/demos`.

4.2 Basic MATLAB NXT toolbox commands

The instructions in this section have been adapted from RWTH's website. Installation instructions from:

<http://www.mindstorms.rwth-aachen.de/trac/wiki/Download4.03>

Coding instructions from:

<http://www.mindstorms.rwth-aachen.de/trac/wiki/Documentation>

You can find more instruction on installation and usage of the RWTH Mindstorms NXT Toolbox from both of these sites.

4.2.1 Startup NXT

The first thing to do is make sure the workspace is clear. Enter:

```
COM_CloseNXT('all');
close all;
clear all;
```

To start, enter:

```
hNXT=COM_OpenNXT;      %hNXT is an arbitrary name
COM_SetDefaultNXT(hNXT); %sets opened NXT as the
                        %default handle
```

4.2.2 NXT Motors

Motors are treated as objects. To create one, enter:

```
motorA = NXTMotor('a'); %motorA is an arbitrary name, 'a' is
                        %the port the motor connected to
```

This will give:

```
NXTMotor object properties:
      Port(s): 0 (A)
      Power: 0
      SpeedRegulation: 1 (on)
      SmoothStart: 0 (off)
      TachoLimit: 0 (no limit)
      ActionAtTachoLimit: 'Brake' (brake, turn off when stopped)
```

4.2.3 Basic Motor Commands & Properties

Below is a list of these properties and how to change them:

Power

Determines speed of the motor

```
motorA.Power=50; % value must be between -100 and 100 (negative
                % will cause the motor to rotate in reverse)
```

SpeedRegulation

If the motor encounters some sort of load, the motor will (if possible) increase its power to keep a constant speed

```
motorA.SpeedRegulation=true; % either true or false, or
                             alternatively, 1 for true, 0 for
                             false
```

SmoothStart

Causes the motor to slowly accelerate and build up to full speed.

Works only if ActionAtTachoLimit is not set to 'coast' and if TachoLimit>0

```
motorA.SmoothStart= true; % either true or false, or
                          % 1 for true, 0 for false
```

ActionAtTachoLimit

Determines how the motor will come to rest after the TachoLimit has been reached.

There are three options:

1. 'brake': the motor brakes
2. 'Holdbrake': the motor brakes, and then holds the brakes
3. 'coast' the motor stops moving, but there is no braking

```
motorA.ActionAtTachoLimit='coast';
```

TachoLimit

Determines how far the motor will turn

```
motorA.TachoLimit= 360; % input is in terms of degrees
```

Alternative Motor Initiation

Motors can also be created this way:

```
motorA=NXTMotor('a', 'Power', 50, 'TachoLimit', 360);
```

4.2.4 Other Motor Commands

SendToNXT

This is required to send the settings of the motor to the robot so the motors will actually run.

```
motorA.SendToNXT();
```

Stop

Stops the motor. There are two ways to do this:

1. 'off' will turn off the motor, letting it come to rest by coasting.
2. 'brake' will turn cause the motor to be stopped by braking, however the motors will need to be turned off after the braking.

```
motorA.Stop('off');
```

ReadFromNXT();

Returns a list of information pertaining to a motor

```
motorA.ReadFromNXT();
```

Entering `motorA.ReadFromNXT.Position()`; will return the position of the motor in degrees.

ResetPosition

Resets the position of the motor back to 0

```
motorA.ResetPosition();
```

WaitFor

Program will wait for motor to finish current command. For example:

```
motorA=('a', 'Power', 30, 'TachoLimit', 360);
motorA.SendToNXT();
motorA.SendToNXT();
```

This command will cause problems as the motor can only process one command at a time. Instead, the following should be entered:

```
motorA=('a', 'Power', 30, 'TachoLimit', 360)
motorA.SendToNXT();
motorA.WaitFor();
motorA.SendToNXT();
```

The exception to this is if TachoLimit of the motor is set to 0.

4.2.5 Using Two Motors At Once

Some operations, for example driving forward and backwards, require the simultaneous

```
use of two motors. Entering: B=NXTMotor('b', 'Power', 50, 'TachoLimit', 360);
C=NXTMotor('c', 'Power', 50, 'TachoLimit', 360);
B.SendToNXT();
C.SendToNXT();
```

will start the bot moving, but the signals for both motors to start at will not be sent at exactly the same time, so the robot will curve a little and fail to drive in a straight line.

Instead, you should enter:

```
BC=NXTMotor('bc', 'Power', 50, 'TachoLimit', 360);
```

OR

```
BC = NXTMotor('bc');
BC.Power=50;
BC.TachoLimit=360;
```

Turning left or right can be achieved by only running one motor at a time, or by moving both motors, but one slower than the other.

4.2.6 Sensors

The following commands are used to open a sensor, plugged into port 1:

```
OpenSwitch(SENSOR_1);      % initiates touch sensor
OpenSound(SENSOR_1, 'DB'); % initiates sound sensor, using
                          % either 'DB' or 'DBA'
OpenLight(SENSOR_1, 'ACTIVE'); % initiates light sensor as
                              % either 'ACTIVE' or 'INACTIVE', The following com-
                              % plugged into Port 1
OpenUltrasonic(SENSOR_1); % initiates ultrasonic sensor
                          % plugged into Port 1
```

mands are used to get values from the sensor plugged into port 2:

```
GetSwitch(SENSOR_2);      % returns 1 if pressed, 0 if depressed
GetSound(SENSOR_2);      % returns a value ranging from 0-1023
GetLight(SENSOR_2);      % returns a value ranging from 0 to
                          % a few thousand
GetUltrasonic(SENSOR_2); % returns a value in cm
```

To close a sensor, ex. Sensor 1:

```
CloseSensor(SENSOR_1); %properly closes the sensor
```

4.2.7 Direct NXT Commands

PlayTone

Plays a tone at a specified frequency for a specified amount of time

```
NXT_PlayTone(400,300); % Plays a tone at 400Hz for 300 ms
```

KeepAlive

Send this command every once in a while to prevent the robot from going into sleep mode:

```
NXT_SendKeepAlive('dontreply');
```

Send this command to see how long the robot will stay awake, in milliseconds:

```
[status SleepTimeLimit] = NXT_SendKeepAlive('reply');
```

GetBatteryLevel

Returns the voltage left in the battery in millivolts

```
NXT_GetBatteryLevel;
```

StartProgram/StopProgram

To run programs written on LEGO Mindstorms NXT software, enter:

```
NXT_StartProgram('MyDemo.rxe') % the file extension '.rxe' can be  
                                % omitted, it will then be automatically  
                                % added
```

Entering `NXT_StopProgram` stops the program mid-run.

4.3 First examples

The following exercises are intended to explore how to use RWTH's Mindstorms NXT Toolbox.

4.3.1 Example 1: Wall avoidance

The following is a simple example of how to drive a robot and use the ultrasonic sensor. The robot will drive forward until it is around 20 cm away from a barrier (i.e. a wall), stop, beep, turn right, and continue moving forward. The robot will repeat this 5 times. Attach the Ultrasonic sensor and connect it to port 1. The study and run the following program.

```

COM.CloseNXT('all');           %cleans up workspace
close all;
clear all;
hNXT=COM.OpenNXT();           % initiates NXT, hNXT is an arbitrary name
COM.SetDefaultNXT(hNXT);      %sets default handle

OpenUltrasonic(SENSOR_1);

forward=NXTMotor('BC'); &% setting motors B           C to drive forward
forward.Power=50;
forward.TachoLimit=0;
turnRight=NXTMotor('B'); &% setting motor B to turn right
turnRight.Power=50;
turnRight.TachoLimit=360;
for i= 1:5
    while GetUltrasonic(SENSOR_1)>20
        forward.SendToNXT();           %sends command for robot to move forward
                                         %TachoLimit=0; no need for a WaitFor() statement

    end %while
    forward.Stop('brake');             %robot brakes from going forward
    NXT_PlayTone(400,300);            %plays a note
    turnRight.SendToNXT;              %sends the command to turn right
    turnRight.WaitFor;                %TachoLimit is not 0; WaitFor() statement required
end %for
turnRight.Stop('off');              %properly closes motors
forward.Stop('off');
CloseSensor(SENSOR_1);              %properly closes the sensor
COM.CloseNXT(hNXT);                 % properly closes the NXT
close all;
clear all;

```

4.3.2 Example 2: Line following

The next exercise is writing your own program that uses readings from its light sensor to drive the NXT and follow a line.



Setup:

1. Mount light sensor, facing downwards on front of NXT and plugged into Port 3.
2. Mount a switch sensor on the NXT, plugged into Port 2.
3. Use a piece of dark tape (i.e. electrical tape) to mark a track on a flat, light coloured surface. Make sure the tape and the surface are coloured differently enough that the light sensor returns reasonably different values between the two surfaces.
4. Write a program so that the tribot follows the line.

4.4 Classical control theory

4.4.1 Inverse plant dynamics and feedback control

In this section we discuss control systems that are at the brains of robots. Control systems are necessary to guide actions in an uncertain environment. The principle idea of a control system is to use sensory and motivational information to produce goal directed behavior.

A control system is characterized by a **control plant** and **controller**. A control plant is the dynamical object that should be controlled, such as a robot arm. The state of this object is described by a state vector

$$\mathbf{z}^T(t) = (1, \mathbf{x}(t), \mathbf{x}^{(1)}(t), \mathbf{x}^{(2)}(t), \dots), \quad (4.1)$$

where $\mathbf{x}(t)$ are basic coordinates such as the positions on a plane for a land-based robot and its heading direction at a particular time. We also included a constant part in the description of the plant, as well as higher derivatives by writing the i -th derivative as $\mathbf{x}^{(i)}(t)$.

The state of the plant is influenced by a **control command** $\mathbf{u}(t)$. A control command can be, for example, sending a specific current to motors, or the initiation of some sequences of inputs to the robot. The effect of a control command $\mathbf{u}(t)$ when the plant is in state $\mathbf{z}(t)$ is given by the **plant equation**

$$\mathbf{z}(t + 1) = \mathbf{f}(\mathbf{z}(t), \mathbf{u}(t)), \quad (4.2)$$

Learning the plant function \mathbf{f} will be an important point in adaptive control discussed later, and we will discuss some examples below.

The **control problem** is to find the appropriate commands to reach **desired states** $\mathbf{z}^*(t)$. We assume for now that this desired state is a specific point in the state space, also called **setpoint**. Such a control problem with a single desired state is called **point-to-point control**. The control problem is called **tracking** if the desired state is changing. In an ideal situation we might be able to calculate the appropriate control commands. For example, if the plant is linear in the control commands, that is, if \mathbf{f} has the form

$$\mathbf{z}(t + 1) = \mathbf{f}(\mathbf{z}(t), \mathbf{u}(t)), \quad (4.3)$$

and \mathbf{g} has an inverse, then it is possible to calculate the command to reach the desired state as

$$\mathbf{u}^* = \mathbf{f}^{-1}(\mathbf{z})\mathbf{z}^*. \quad (4.4)$$

A block diagram of this strategy which shows the controller and plant is shown in Figure 4.2.



Fig. 4.2 The elements of a basic control system.

As an example, let us consider the tribot moving a specified distance from a starting point. Let $u = t$ be the motor command for the Tribot to move forward for t seconds with a certain motor power. If the tribot moves a distance of d_1 in one second, then we expect the tribot to move a distance of $d * t$ in t seconds. The plant equation for this tribot movement is hence

$$x = x_0 + d * t, \quad (4.5)$$

To find out the parameter d we can mark the initial location of the tribot and let it run for 1 second. The parameter can then be determined from the end location x by $d = (x - x_0)/t$. Note that we learned a parameter from measurement in the environment. This is at the heart of machine learning and the formula is already learning algorithm.

In the next experiment want the robot to move from a start position to a desired position of 60cm away. Ideally this could be achieved by letting the motor run by $t = 60/d$ seconds. Try it out. While the tribot might come close to the desired state, a perfect match is not likely. There are several reasons for such failures. One is that our measurement of the dynamics might be not accurate enough. We also made the assumption that the rotations of the wheels are linear in time, but it might be that the wheels slow down after a while due to power loss or heating of the gears which alter physical properties, etc. And most of all, disturbances in the environment can through the robot off track (such as a mean instructor). All these influences on the controlled object are indicated in the figure by a disturbance signal to the plant.

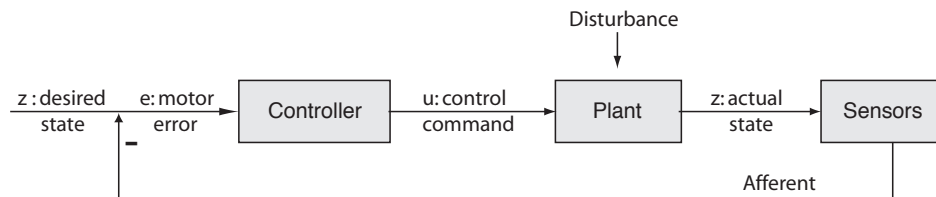


Fig. 4.3 The elements of a basic control system with negative feedback.

Of course, if we do not reach the desired state we can initiate a new movement to compensate for the discrepancy between the desired and actual state of the controlled object. For this we need a measurement of the new position (which actually might also

contain some error) from which we can calculate the new desired distance to travel. We call the distance between the desired location x^* and the actual location x the **displacement error**

$$e(t) = x^*(t) - x(t) \quad (4.6)$$

We can also iterate the procedure until the distance to the desired state is sufficiently small. Such a controller is called a **feedback controller** and is shown in Fig. 4.3. The desired state is the input to the system, and the controller uses the desired state to determine the appropriate control command. The controller can be viewed as an **inverse plant model** as it takes a state signal and produce the right command so that the controlled object ends up in the desired state. The motor command thus causes a new state of the object that we have labelled ‘actual state’ in Fig. 4.3. The actual state is then measured by sensors and subtracted from the desired state. If the difference is zero, the system has reached the desired state. If it is different than zero than this difference is the new input to the control system to generate the correction move.

Exercise

Make the robot to move to a specific distance from a wall by using a feedback controller with the ultrasonic sensor.

4.4.2 PID control

The negative feedback controller doe often work in minimizing the position error of a robotic systems. However, the movement is often jerky and overshooting a setpoint. There are simple additions of the basic feedback controller that will making the reaching of a setpoint better. Here we discuss briefly a very common feedback controller call **PID coontroller** for reasons that will become clear shortly.

The basic idea of a PID controller is to not only use the current error between the desired state and estimated actual state, but to take some history into account. For example, when the correction in each state takes a long time, that is, if the sum of the errors is large, then we should make larger changes so that we reach the setpoint faster. In a continuous system, the sum over the last errors becomes and integral. Such a component in the controller should help to accelerate the reaching of the setpoint. However, such a component can also increase overshooting an leads to cycles which can even make the system unstable. To remedy this it is also common to take the rate of change in the error into account. Such a term is corresponds to the derivative in a continuous system. The name for a PID controller is actually the acronym for **P**roportional, **I**ntegral and **D**erivative and is illustrated in figure ??

The motor command generated by a PID controller is given by

$$u(t) = k_P e(t) + k_I \int e(t) dt + k_D \frac{de(t)}{dt}, \quad (4.7)$$

where we have weighted each of the components with different constants. Appropriate choices of these constants are often chosen by trial and error.

We have only scratched the surface of classical control theory at this point. In order to make controllers robust and applicable for practical applications, many other

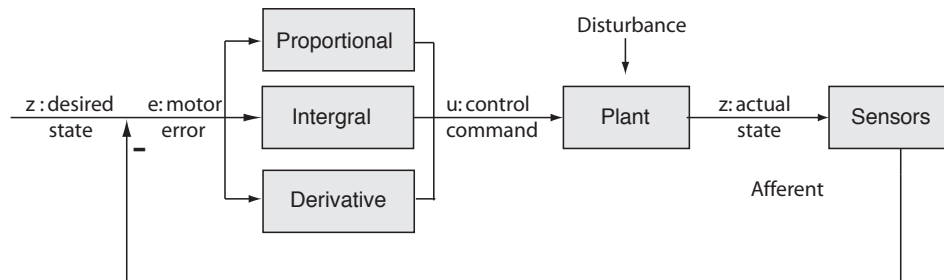


Fig. 4.4 The elements of a PID control system.

considerations have to be made. For example, we are often not only interested in minimizing the motor error but actually minimizing a cost function such as minimizing the time to reach a setpoint or to minimize the energy used to reach the setpoint. Corresponding methods are the subject of **optimal control** theory. While we will not follow classical optimal control theory here, we will come back to topic in the context of reinforcement learning later in the course.

Another major problem for many applications is that the plant dynamic can change over time and has to be estimated from data. This is subject of the area of **adaptive control** to which we will return after introducing the main machine learning methods.

The control theory outlines so far has several other drawbacks in the computer environment. In particular, it treats the feedback signal as the supreme knowledge not taking in consideration that it can be wrong. For example, in the exercise above with measured the distance to a wall to drive the tribot to a particular distance. If we put our hand in between the wall and the tribot, the controller would treat our hand as the wall and would try to adjust the position accordingly. A smarter controller might ask if this is consistent with previous measurements and its movements it made lately. A smarter controller could therefore benefit from internal models and an acknowledgement, and corresponding probabilistic treatment, that the sensor information is not always reliable. Even more, a smart controller should be able to learn from the environment how to judge certain information. This will be the our main strategy to follow in this course.

4.5 Adaptive Control

In many cases we do not know the plant dynamics in a closed analytic form so that we have to estimate the dynamics and approximate the inverse. Also, it is common that the dynamics of the plant changes over time due to fatigue of the elements or disturbances in the environment. A solution to these problems is to let the controller learn from the data during operation.

The general idea is outlined in Figure ???. The controller knows already what input it receives, but the new stage is that the controller receives feedback from the environment from the sensor. This information is used to adapt the controller as indicated by the

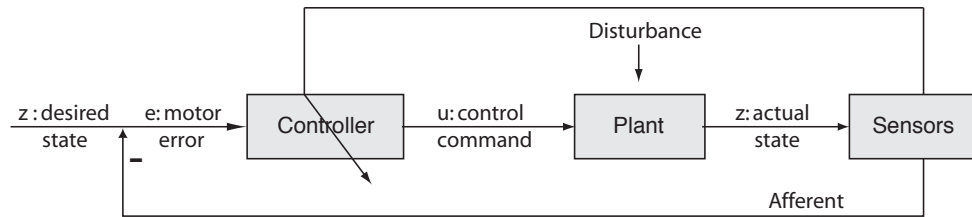


Fig. 4.5 Two advanced control systems which incorporate models for corrective adjustments in the system. These models are trained with sensory feedback, which is indicated by the arrows going through these components.

crossing arrow. General architectures of learning functions from examples corresponds to supervised learning that will be covered in the next chapter.

Note that there are now two feedback loops, for previously discussed negative feedback to initiate corrective movements, and the newly introduced feedback for the adaptive control. The negative feedback is instantaneous and a low level reactive system for motor control. The adaptive feedback often works on a larger time scale to adapt the controller to slowly changing plants.

Two refined schemes for adaptive motor control with slow sensory feedback are illustrated in Figs. 4.6. The first one employs some subsystems that mimic the dynamic of the controlled object and the behaviour of the sensory system. These subsystems are called **forward models**. If these models are good approximations of the real systems they are modelling, then we can use the output of these systems, instead of the slow sensory response, as feedback signal. The models have, of course, to be gauged against the real system during ongoing learning. Thus, the sensory feedback is used to change the behaviour of the forward model. The forward model, which is divided into a **dynamic** and **output** component in the figure, influences the sensory feedback to improve performance in the main control loop. This scheme works as long as the changes in the systems that they mimic are much slower than the time scale of the movement that is controlled.

The second scheme, shown in Fig. 4.6B, employs an inverse model instead of the forward model in the previous scheme and is therefore called an **inverse model controller**. The inverse model, which is incorporated as a side-loop to the standard feedback controller, learns to correct the computation of the motor command generator. The reason that this scheme works is similar to that for the previous example. The sensory feedback can be used to make the inverse model controller more accurate while it provides the necessary correction signals in time to be incorporated into the motor command.

Exercise

Write a PID controller that let the tribot follow a wall. Experiment with different control parameters and report your findings.

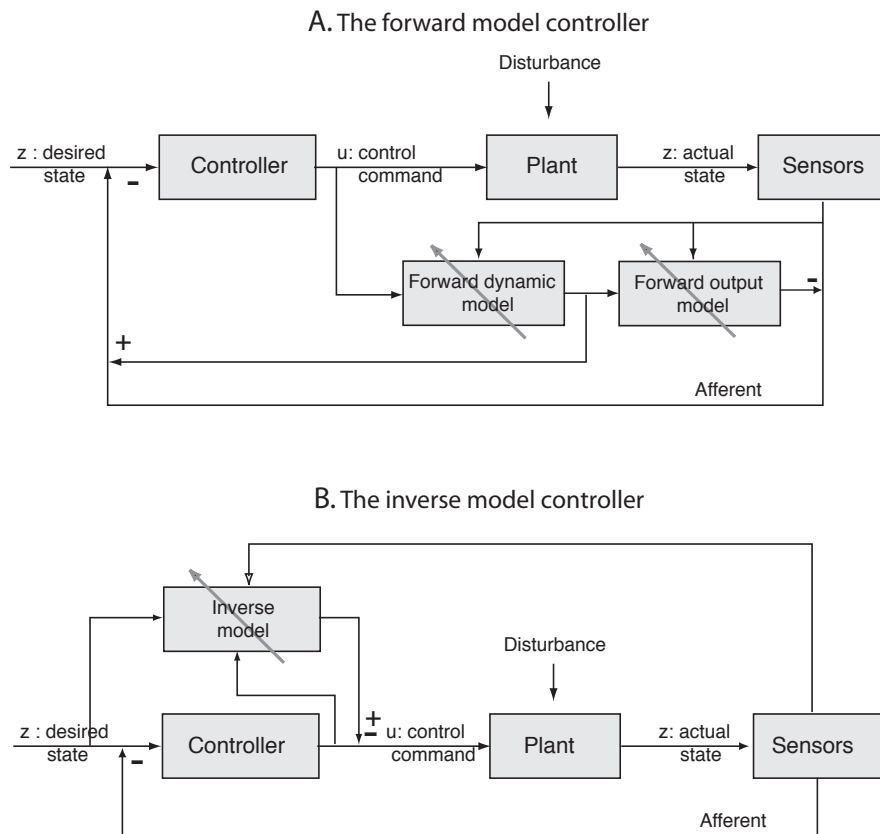


Fig. 4.6 Two advanced control systems which incorporate models for corrective adjustments in the system. These models are trained with sensory feedback, which is indicated by the arrows going through these components.

4.6 Configuration space

Very important for most algorithms to control a robot is the description of its state or the possible states of the system. The physical state of an agent can be quite complicated. For example, the Lego components of the tribot can have different positions and can shift in operation; the battery will have different levels of charge during the day, lightening conditions change, etc. Thus, description of the physical state would need a lot of variables, and such a description space would be very high dimensional, which is one important source of the computational challenges we face. To manage this problem we need to consider abstractions.

Abstraction is an important concept in computer science and science in general. To abstract means to simplify the system in a way that it can be described in as simple terms as possible to answer the specific questions under consideration. This philosophy is sometimes called the **principle of parsimony**, also known as **Occam's razor**. Basically, we want a model as simple as possible, while still capturing the main

aspects of the system that the model should capture.

For example, let us consider a robot that should navigate around some obstacles from point A to point B as shown in Fig.4.7A. The simplifications we make in the following is that we take consider the movement in only in a two-dimensional (2D) plane. The robot will have a position described by an x and y coordinate, and a heading direction described by an angle α . Note that we ignore for now the physical extension of the robot, that is, we consider it here only as a point object. If this is a problem, we can add an extension parameter later. We also ignore many of the other physical parameters mention above. The current state of robot is hence described by three real numbers, $[x, y, \alpha]$, and the space of all possible values is called the **configuration space** or **state space**. An obstacle, as shown in Fig.4.7A with the grey area, can be represented as state values that are not allowed by the system.

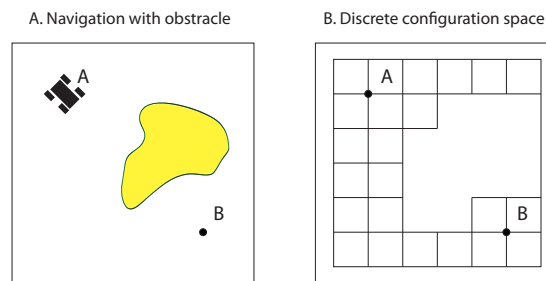


Fig. 4.7 (A) A physical space with obstacles in which an agent should navigate from point A to point B . (B) Discretized configuration space.

We have reduced the description of an robot navigation system from a nearly infinite physical space to a three-dimensional (3D) configuration space by the process of abstraction. However, there are still an infinite number of possible states in the state space if the state variables are real numbers. We will shortly see that we have often to search the state space, and an infinite space then often becomes problematic. A common solution to this problem is to make further abstractions and to allow the state variables to have only a finite number of states. Thus, we commonly **discretize** the state space. An example is shown in Fig.4.7B, where we used a grid to represent the possible values of the robot positions. Such a discretization is very useful in practice, and we can also make the discretization error increasingly small by decreasing the grid parameter, Δx , which describes the minimal distance between two states. We will use this grid discretization for planing a path through the configuration space later in this chapter.

4.6.1 Robot Arm State Space Visualizer

The goal of the following exercise is to graph the state space in which a double jointed robot arm is able to move. To do this, use two motors of the NXT Lego kit as joints of this arm, similarly to the one pictured in Fig.4.8. Also attach a touch sensor that we will use as a button in the following. Secure it on the table so that you can move its pointing finger around by turning the externally moving the motors. We will use the

ability of measuring this movement in the motors by the Lego system. Also, use a box, a coffee mug or some other items to create some obstacles for the arm to move around.



Fig. 4.8 (A) Robot arm with two joints. (B)

Use the Matlab program in Table 4.1 to visualize the state space. Start the Matlab program with the robot arm touching one of the obstacles. Then move the arm all around one of the obstacles. This will allow the NXT to create set of data that will display which combinations of angles cannot be travelled through by the robot arm. To map more than one obstacle without mapping the angles in between the two obstacles, push the button. The button pauses the program for 4 seconds so that you can move the arm from one obstacle to the next without mapping the angles in between the two.

4.6.2 Exercise

1. Write a program that produces a map of the state space in form of a matrix which has zeros of entries that can be reached and ones as entries for states that have obstacles in it.

Table 4.1 State space visualization program

```

for counter = 1:10000
    %The angles for both wheels are read
    bottomDeg=bottom.ReadFromNXT.Position;
    topDeg=top.ReadFromNXT.Position;
    %the angles are then stored in angleWheel
    angleWheel(1,counter)=bottomDeg;
    angleWheel(2,counter)=topDeg;
    %This plot will show all the angles read from the
    %wheels, because hold is on
    plot(topDeg,bottomDeg,'o');
    %If the switch is pushed then pause the angle reader
    if GetSwitch(SENSOR_1)==1
        disp('paused');
        pause(4);
        disp('unpaused');
    end
end
end

```

4.7 Planning as search

Planning a movement is basic requirement in robotics. In particular, we want to find the **shortest path** from a starting position to a goal position. The previous exercise provides us with a **map** that we can use to plan the movement. One of the elements of the array is dedicated as the starting position, and one as the goal position. We now need to write a search algorithm that finds the shortest path.

There are many search algorithms that we could utilize for this task. For example, a basic search algorithm is **depth-first-search** that chooses a neighboring node, then moves to the next, etc, until it reaches a goal. If the algorithm ends up in a dead-end, then the algorithm tracks back to the last node with another options and goes down this node. The **breadth-first-search** algorithms searches instead first if each of the neighboring nodes of a current node is the goal state before advancing to the next level. Of course, both need to keep track of which routs have been tried already to ensure that we do not continuously try failed paths. These algorithms would sooner or later find a path between the starting position and the goal if it exists since these algorithms would try out all possible paths. However, these algorithms do not guarantee to find the shortest path, and these algorithms are also usually not so efficient.

To find better solutions we could take more information into account. For example, in the grid world, we can use some distance measure between the current state and the goal, such as the **Euclidean distance**

$$g(\text{current}) = \sqrt{(x_{\text{goal}} - x_{\text{current}})^2 + (y_{\text{goal}} - y_{\text{current}})^2}, \quad (4.8)$$

to guide the search⁵. That is, if we chose the next node to expand, we would chose the node that has the smallest distance to the goal node. Such a **heuristic search algorithm** is sometimes called a **greedy best-first search**. Valid heuristics, which are estimations of the expected cost that must not underestimate the real cost, can greatly accelerate search performance. We can take such strategy a step further by also taken into account in the cost of reaching the current state from the starting state,

$$h(\text{current}) = \sqrt{(x_{\text{start}} - x_{\text{current}})^2 + (y_{\text{start}} - y_{\text{current}})^2}, \quad (4.9)$$

so that the total heuristic cost of node ‘current’ is

$$f(\text{current}) = h(\text{current}) + g(\text{current}). \quad (4.10)$$

The algorithm that uses this heuristics while taking track of paths that are still ‘open’ or that are ‘closed’ since they lead to dead-ends, is called an A* star search algorithm. An Matlab example and visualization of the A* search algorithm is provided on the course web page.

4.7.1 Exercise: Robot Arm State Space Traversal

The goal of this exercise is to make the double jointed NXT arm progress through the state space. You need the double-jointed robot arm and the obstacles from the state space visualization exercise above.

- Make a program which can move the arm from the starting coordinate to the ending coordinate while avoiding the obstacles in the state space.

⁵This distance measure is also called the L_2 norm. Other measures such as the Manhattan distance, the Minkowski distance, or the L_1 norm also frequently used.

Part III

Supervised learning

5 Regression and maximum likelihood

In the previous chapter we used some measurements to estimate a plant function for the tribot when driving it forward with a certain power for different amounts of time. This chapter starts a more thorough discussion of such learning problems, that of supervised learning. We start by considering regression with noisy data and formalize our learning objectives. Along the line we will also learn some useful techniques to find extrema of a function and how to apply such methods to classification problems.

5.1 Regression of noisy data

An important type of learning is **supervised learning**. In supervised learning, examples of input-output relations are given and our goal is to discover the relations between these data so that we can make **predictions** of previously unseen data.

More formally, let us consider **training data** that are given to a learner by a teacher. The training data are often denoted by the inputs \mathbf{x} , such as motor commands, and the outputs y , such as the reaction of a plant. In the following we consider m training examples, the pairs of values

$$\{(\mathbf{x}^{(i)}, y^{(i)}); i = 1 \dots m\}. \quad (5.1)$$

We use here an index counter i to label the different training examples. We use brackets around it to distinguish it from an exponent. The above notation describes mapping examples from a n dimensional input to a 1-dimensional output, hence the vector notation for \mathbf{x} and the scalar notation for y . The 1-dimensional output is mainly for convenience as the generalization to a higher-dimensional output function only clutters the notations slightly.

Let us use again consider the example of a tribot running forward for a certain amount of time. In Figure 5.1 we show several measurements of the distance traveled when both motors are driven forward with a certain amount of power for different times in milliseconds. The shown data set is provided in file `tribotPlantData.mat`. Looking at the plot reveals that there seems to be a systematic relation between the time of running the motor and the distance traveled, and that this trend seems linear. This hypothesis can be quantified as a parameterized function,

$$\hat{h}(x; \theta) = \theta_1 + \theta_2 x. \quad (5.2)$$

This notation means that the hypothesis h is a function of the quantity x , and the hypothesis includes all possible straight lines, where each line can have a different offset θ_1 (intercept with the y -axis), and slope θ_2 .

We typically collect parameters in a **parameter vector** θ . We only considered on input variable x above, but we can easily generalize this to higher dimensional

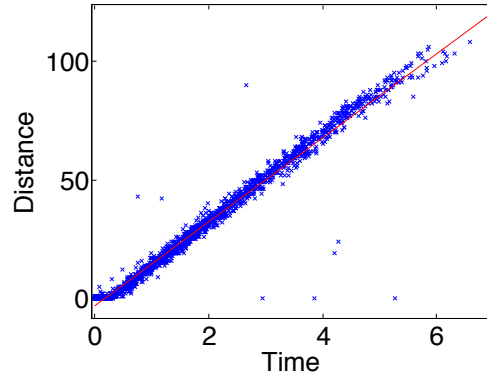


Fig. 5.1 Measurements of distance travelled by the tribot when running the motor for different number of milliseconds.

problems where more input **attributes** are given. For example, we could not only vary the time the motor is running, but also vary the power with which we run the motor. The corresponding 3-dimensional plot is shown in Figure 5.1. This variation of the power also seem to have a linear effect on the distance traveled. To compress our notations we will also use the convention that we consider a constant input in the first component of the input vector,

$$\mathbf{x} = \begin{pmatrix} x_1 = 1 \\ x_2 \\ x_3 \end{pmatrix}. \quad (5.3)$$

We can write the hypothesis as

$$\hat{h}(\mathbf{x}; \theta) = \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3, \quad (5.4)$$

or a linear hypothesis with n attributes as

$$\hat{h}(\mathbf{x}; \theta) = \theta_1 x_1 + \dots + \theta_{n+1} x_{n+1} = \sum_i \theta_i x_i = \theta^T \mathbf{x}. \quad (5.5)$$

The vector θ^T is the transpose of the vector θ .

Now we should remember that our goal was to learn the function that describes these data. In the approach we have taken here we first made a **hypothesis** in form of a parameterized function, $h(\mathbf{x}; \theta)$, and the learning part boils down to determining appropriate values for the parameters. After learning this function we can use it to predict the reaction of the plant even for motor commands for which no training examples were given. The remaining question is how we find appropriate values for the parameters. However, before we do this we will also be more faithful to the data and acknowledge fluctuation around our initial hypothesis.

5.2 Probabilistic models and maximum likelihood

So far, we have only modelled the mean trend of the data, and we should investigate more the fluctuations around this mean trend. Fig.5.2 is a plot of the histogram of

the differences between the actual data and the hypothesis regression line. This look

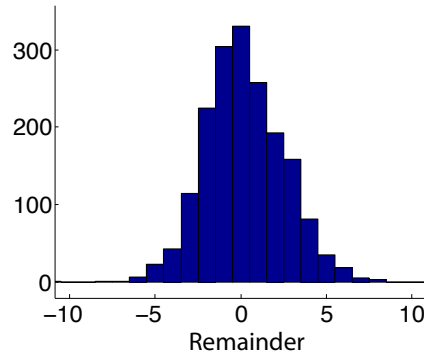


Fig. 5.2 Histogram of the difference between the data points and the fitted hypothesis, $(y_i - \theta_1 - \theta_2 x_i)$.

a bit Gaussian, which is likely to be a common finding though not necessarily the only possible. With this additional conjecture, we should revise our hypothesis. More precisely, we acknowledge that the data are noisy and that we can only give a probability of finding certain values. Specifically, we assume here that the data follow a certain trend, our previous deterministic hypothesis $\hat{h}(\mathbf{x}; \theta)$ with **additive noise**, η ,

$$h(x; \theta) = \hat{h}(\mathbf{x}; \theta) + \eta, \quad (5.6)$$

where the random variable η is Gaussian distributed in the example above,

$$p(\eta) = N(\mu, \sigma) \quad (5.7)$$

We can then also write the probabilistic hypothesis in the example as a Gaussian model of the data distributed with a mean that depends on the variable x ,

$$h(y|x; \theta, \sigma) = N(\mu = \hat{h}(x), \sigma) \quad (5.8)$$

$$= \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y - \theta^T \mathbf{x})^2}{2\sigma^2}\right) \quad (5.9)$$

This functions specifies the probability of values for y , given an input x and the parameters θ and σ .

Specifying a model with a density function is an important step in modern modelling and machine learning. In this type of thinking, we treat data from the outset as fundamentally stochastic, that is, data can be different even in situations that we deem identical. This randomness may come from an **irreducible indeterminacy**, that is, true randomness in the world that can not be penetrated by further knowledge, or this noise might represent **epistemological limitations** such as the lack of knowledge of hidden processes or limitations in observing states directly. The only important fact for us is that we have to live with these limitations. This acknowledgement together with the corresponding language of probability theory has helped to make large progress in the machine learning area.

We will now turn to the important principle that will guide our learning process which corresponds here to estimating the parameters of the model. While the parameterized hypothesis so far describes the form of the data, we need to estimate values for the parameters to make real predictions. We therefore consider again the examples for the input-output pairs, our training set $\{(x^{(i)}, y^{(i)}); i = 1 \dots m\}$. The important principle that we will now follow is to choose the parameters so that the examples we have are most likely. This is called **maximum likelihood estimation**. To formalize this principle, we need to think about how to combine probabilities for several observations. If the observations are independent, then the joined probability of several observations is the product of the individual probabilities,

$$p(y_1, y_2, \dots, y_m | x_1, x_2, \dots, x_m; \theta) = \prod_i^m p(y_i | x_i; \theta). \quad (5.10)$$

Note that y_i are still random variables in the above formula. We now use our training examples as specific observations for each of these random variables, and introduce the **Likelihood function**

$$L(\theta) = \prod_i^m h(\theta; y^{(i)}, x^{(i)}). \quad (5.11)$$

The p on the right hand side is now not a density function, but it is a regular function (with the same form as our parameterized hypothesis) of the parameters θ for the given values $y^{(i)}$ and $x^{(i)}$. Instead of evaluating this large product, it is common to use the logarithm of the likelihood function, so that we can use the sum over the training examples,

$$l(\theta) = \log L(\theta) = \sum_i^m \log(h(\theta; y^{(i)}, x^{(i)})). \quad (5.12)$$

Since the log function strictly monotonically rising, the maximum of L is also the maximum of l . The maximum (log-)likelihood can thus be calculated from the examples as

$$\theta^{\text{MLE}} = \arg \max_{\theta} l(\theta). \quad (5.13)$$

We might be able to calculate this analytically or use one of the search algorithms to find a maximum from this function.

Let us apply this to the linear regression discussed above. The log-likelihood function for this example is

$$l(\theta) = \log \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T \mathbf{x}^{(i)})^2}{2\sigma^2}\right) \quad (5.14)$$

$$= \sum_{i=1}^m \left(\log \frac{1}{\sqrt{2\pi}\sigma} - \frac{(y^{(i)} - \theta^T \mathbf{x}^{(i)})^2}{2\sigma^2} \right) \quad (5.15)$$

$$= -\frac{m}{2} \log 2\pi\sigma - \sum_{i=1}^m \frac{(y^{(i)} - \theta^T \mathbf{x}^{(i)})^2}{2\sigma^2}. \quad (5.16)$$

Thus, the log was chosen so that we can use the sum in the estimate instead of dealing with big numbers based on the product of the examples.

Let us now consider the special case in which we assume that the constant σ , the variance of the data, is the same for all x and thus has a fixed value given to us. We can thus concentrate on the estimation of the other parameters θ . Since the first term in the expression 5.16, $-\frac{m}{2} \log 2\pi\sigma$, is independent of θ , maximizing the log-likelihood function is equivalent to minimizing a quadratic error term

$$E = \frac{1}{2}(y - h((x; (\theta))))^2 \iff p(y|\mathbf{x}; \theta) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{(y - h((x; (\theta))))^2}{2}\right) \quad (5.17)$$

This **error function** or **cost function** was a frequently used criteria called **Least Mean Square (LSM)** regression for parameters estimation when considering deterministic hypothesis. In terms of our probabilistic view, the LSM regression is equivalent to MLE for gaussian data with constant variance. When the variance is a free parameter, then we need to minimize equation 5.16. instead.

We have discussed Gaussian distributed data in most of this section, but one can similarly find correspondences error functions for other distributions. For example, a **polynomial error function** correspond more generally to a density model of the form

$$E = \frac{1}{p} \|y - h((x; (\theta)))\|^p \iff p(y|\mathbf{x}; \theta) = \frac{1}{2\Gamma(1/p)} \exp(-\|y - h((x; (\theta)))\|^p). \quad (5.18)$$

Later we will specifically discuss and use the **ϵ -insensitive error function**, where errors less than a constant ϵ do not contribute to the error measure, only errors above this value,

$$E = \|y - h((x; (\theta)))\|_\epsilon \iff p(y|\mathbf{x}; \theta) = \frac{p}{2(1-\epsilon)} \exp(-\|y - h((x; (\theta)))\|_\epsilon). \quad (5.19)$$

Since we already acknowledged that we do expect that data are noisy, it is somewhat logical to not count some deviations from the expectation as errors. It also turns out that this error function is much more robust than other error functions.

5.3 Maximum a posteriori estimates

In the maximum likelihood estimation we assumed that we have no prior knowledge of the parameters θ . However, we sometimes might know which values of the parameters are impossible or less likely. This prior knowledge can be summarized in the prior distribution $p(\theta)$, and the next question is how to combine this prior knowledge in the maximum likelihood scheme. Combining prior knowledge with some evidence is described by Bayes' theorem. Thus, let us consider again that we have some observations (x, y) from specific realizations of the parameters, which is given by $(p(x, y|\theta))$, and the prior about the possible values of the parameters, given by $p(\theta)$. The prior is in this situation sometimes called the **regularizer**, restricting possible values in a specific domain. We want to know the distribution of parameters given the observation, $p(\theta|x, y)$, which can be calculated from Bayes's theorem,

$$p(\theta|x, y) = \frac{p(x, y|\theta)p(\theta)}{\int_{\theta \in \Theta} p(x, y|\theta)p(\theta)d\theta}, \quad (5.20)$$

where Θ is the domain of the possible parameter values. We can now use this expression to estimate the most likely values for the parameters. For this we should notice that the

denominator, which is called the **partition function**, does not depend on the parameters θ . The most likely values for the parameters can thus be calculated without this term and is given by the **maximum a posteriori (MAP)** estimate,

$$\theta^{\text{MAP}} = \arg \max_{\theta} p(x, y | \theta) p(\theta). \quad (5.21)$$

This is, in a Bayesian sense, the most likely value for the parameters, where, of course, we now treat the probability function as a function of the parameters (e.g., a likelihood function).

A final caution: ML and MAP estimates give us a **point estimate**, a single answer of the most likely values of the parameters. This is often useful as a first guess and is commonly used to **make decisions** about which actions to take. However, it is possible that other sets of parameters values might have only a little smaller likelihood value, and should therefore also be considered. Thus, one limit of the estimation methods discussed here is that they do not take distribution of answers into account, which is more common in more advanced Bayesian methods.

5.4 Minimization procedures

We have discussed the important principle of maximum likelihood estimation to learn parameters from data so that the data are most likely under our hypothesis. This principle tells us how to use the training examples to come-up with some reasonable values for the parameters. To execute these principles we have to find the maximize of a (log)likelihood function or minimize a corresponding cost functions such as LMS to determine those parameters.

In this technique we chose values for θ that will minimize the sum of the **mean square error (MSE)** from the line to each data point. In general we consider a **cost function** to be minimized, which is in this case the square function

$$E(\theta) = \frac{1}{2} (y - h(x; \theta))^2 \quad (5.22)$$

$$\approx \frac{1}{2m} \sum_i (y^{(i)} - h(x^{(i)}; \theta))^2. \quad (5.23)$$

Note that the objective function is a function of the parameters. Also, there is a small subtlety in the above equation since we wrote the general form of the objective function in line 5.22 and considered its approximation with the data, considered to be independent, in line 5.23. With this objective function, we reduced the learning problem to a search problem of finding the parameter values that minimize this objective function,

$$\theta = \arg \min_{\theta} E(\theta) \quad (5.24)$$

We will demonstrate practical solutions to this search problem with three important methods. The first method is an analytical one. You might remember from basic mathematics classes that finding a minimum of a function $f(x)$ is characterized by

$\frac{df}{dx} = 0$, and $\frac{d^2f}{dx^2} > 0$. Here we have a vector function since the cost function depends on several parameters. The derivative then becomes a gradient

$$\nabla_{\theta} E(\theta) = \begin{pmatrix} \frac{\partial E}{\partial \theta_0} \\ \cdot \\ \cdot \\ \frac{\partial E}{\partial \theta_n} \end{pmatrix}. \quad (5.25)$$

It is useful to collect the training data in a large matrix for the x values, and a vector for the y values,

$$X = (\mathbf{x}^{(1)} \dots \mathbf{x}^{(m)}) \quad Y = (y^{(1)} \dots y^{(m)}). \quad (5.26)$$

We can then write the cost function as

$$E(\theta) = \frac{1}{2m} (Y - X\theta)(Y - X\theta)^T, \quad (5.27)$$

The parameters for which the gradient is zero is then given by the **normal equation**

$$\theta = (XX^T)^{-1}XY^T. \quad (5.28)$$

We have still to make sure these parameters are the minimum and a maximum value, but this can easily be done and is also obvious when plotting the result.

The above analytic method is optimal in the requirement of computational time. However, it requires that we can analytically solve an equation system. This was easy in the linear case but fails for more complex functions. The next methods are more widely applicable.

The second method we mainly mention for illustrative reasons is random search. This is a very simple algorithm, but worthwhile considering to compare them to the other solutions. In this algorithm, a new random values for the parameters are tried, and these new parameters replace the old ones if the new values result in a smaller error value (see Matlab code in Tab.5.3). The benefit of this method is that it can be applied to any function, but in practice it takes much too long to find a solution.

5.5 Gradient Descent

The final method discussed here for finding a minimum of a function $E(\theta)$ is **Gradient Descent**. This method will often be used in the following and it will thus be reviewed here in more detail.

Gradient Descent starts at some initial value for the parameters, θ , and improves the values iteratively by making changes to the parameters along the negative gradient of the cost function,

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} E(\theta). \quad (5.29)$$

The constant α is called a **learning rate**. The principle idea behind this method is illustrated for a general cost function with one parameter in Fig.5.4.

Table 5.1 Program randomSearch.m

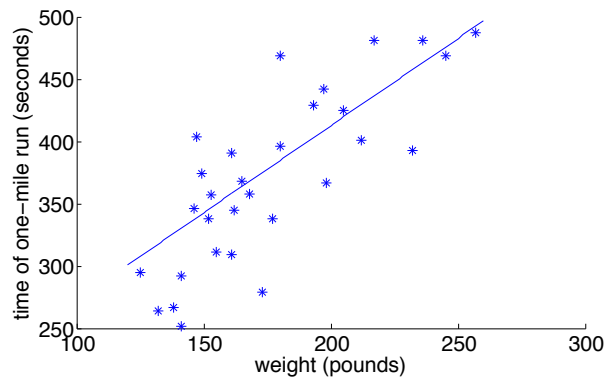
```

%% Linear regression with random search
clear; clf; hold on;

load healthData;
E=1000000;
for trial=1:100
    thetaNew=[100*rand()+50, 3*rand()];
    Enew=0.5*sum((y-(thetaNew(1)+thetaNew(2)*x)).^2);
    if Enew<E; E=Enew; theta=thetaNew; end
end

plot(x,y,'*')
plot(120:260,theta(1)+theta(2)*(120:260))
xlabel('weight (pounds)')
ylabel('time of one-mile run (seconds)')

```

**Fig. 5.3** Health data with linear least-mean-square (LMS) regression from random search.

The gradient is simple the slope (local derivative) for a function with one variable, but with functions in higher dimensions (more variables), the gradient is the local slope along the direction of the steepest ascent, and since we are interested here in minimizing the cost function we make changes along the negative gradient. For large gradients, this method takes large steps, whereas the effective step-width becomes smaller near a minimum. Gradient descent works often well for local optimization, but it can get stuck in local minima. The corresponding update rule in case of the LMS error function,

$$E(\theta) = \frac{1}{2m} \sum_{i=1}^m \left(y^{(i)} - h(x^{(i)}; \theta) \right)^2, \quad (5.30)$$

with a general hypothesis function $h(x^{(i)}; \theta)$ is given by

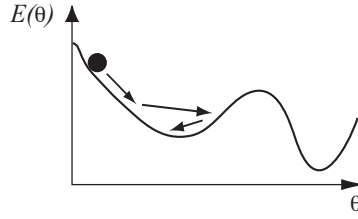


Fig. 5.4 Illustration of error minimization with a gradient descent method on a one-dimensional error surface $E(\theta)$.

$$\theta_k \leftarrow \theta_k - \frac{\alpha}{m} \sum_{i=1}^m (y^{(i)} - h(x^{(i)}; \theta)) \frac{\partial h(\theta)}{\partial \theta_k}. \quad (5.31)$$

For a linear regression function, the update rule for the two parameters θ_0 and θ_1 are therefore:

$$h(x^{(i)}; \theta) = \theta_0 + \theta_1 x^{(i)} \quad (5.32)$$

$$\theta_0 \leftarrow \theta_0 - \alpha \frac{\partial E(\theta)}{\partial \theta_0} \quad (5.33)$$

$$\theta_1 \leftarrow \theta_1 - \alpha \frac{\partial E(\theta)}{\partial \theta_1} \quad (5.34)$$

$$\frac{\partial E(\theta)}{\partial \theta_0} = \frac{1}{m} \sum_{i=1}^m (y^{(i)} - \theta_0 - \theta_1 x^{(i)}) (-1) \quad (5.35)$$

$$\frac{\partial E(\theta)}{\partial \theta_1} = \frac{1}{m} \sum_{i=1}^m (y^{(i)} - \theta_0 - \theta_1 x^{(i)}) (-x^{(i)}), \quad (5.36)$$

which lead to the final rule:

$$\theta_0 \leftarrow \theta_0 + \frac{\alpha}{m} \sum_{i=1}^m (y_i - \theta_0 - \theta_1 x_i) \quad (5.37)$$

$$\theta_1 \leftarrow \theta_1 + \frac{\alpha}{m} \sum_{i=1}^m (y_i - \theta_0 - \theta_1 x_i) x_i. \quad (5.38)$$

Note that the learning rate α has to be chosen small enough for the algorithm to converge. An example is shown in Fig. 5.5, where the dashed line shows the initial hypothesis, and the solid line the solution after 100 updates.

In the algorithm above we calculate the average gradient over all examples before updating the parameters. This is called a **batch algorithm** or **synchronous update** since the whole batch of training data is used for each updating step and the update is only made after seeing all training data. This might be problematic in some applications as the training examples have to be stored somewhere and have to be recalled continuously. A much more applicable method, also thought to be more biologically realistic, is to use each training example when it comes in and disregards it right after. In this way we do not have to store all the data. Such algorithms are called **online algorithms** or **asynchronous update**. Specifically, in the example above, we calculate the

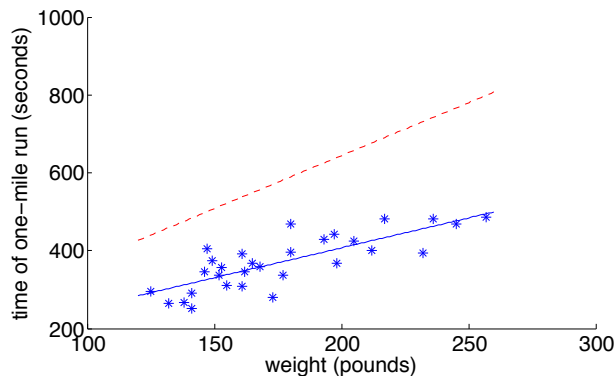


Fig. 5.5 Health data with linear least-mean-square (LMS) regression with gradient descent. The dashed line shows the initial hypothesis, and the solid line the solution after 100 updates.

change for each training examples and update the parameters for this training example before moving to the next example. While we might have to run through the short list of training examples in this specific example, it can still be considered online since we need only one training example at each training step and a list could be supplied to us externally. There are also variations of this algorithms depending on the order we use the training examples (e.g. random or sequential), although this should not be crucial in for the examples discussed here.

5.5.1 LinearRegressionExampleCode

```

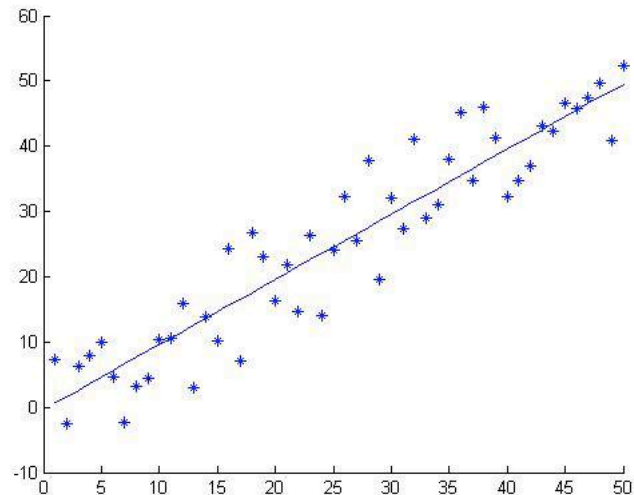
% Linear regression with gradient descent
clear all; clc; hold on;

load SampleRegressionData; m=50; alpha=0.001;
theta1=rand*100*((rand<0.5)*2-1);
theta2=rand*100*((rand<0.5)*2-1);

for trial=1:50000
    sum1 = sum(y - theta1 - theta2 *x);
    sum2 = sum((y - theta1 - theta2 *x).* x);
    theta1 = theta1 + (2*alpha/m) * sum1;
    theta2 = theta2 + (2*alpha/m) * sum2;
    sum1 = 0; sum2 = 0;
end

plot(x, y, '*');
plot(x, x*theta2+theta1);
hold off;

```



A final remark: We have used in this section the popular cost function MSE to calculate the regression in the above example, but it is interesting to think about this choice a bit more. If we change this function, then we would value outliers in a different way. The MSE is a quadratic function and does therefore weight heavily large deviations from the regression curve. An alternative approach would be to say that these might be outliers to which we should not give so much weight. In this case we might want that the increase of the error value decrease with large distances, for example by using the logarithm function $\log(|y - h|)$ as cost (error) function. So, the MSE is certainly not the only choice. We will see below that there is a natural choice for cost functions if we take the stochastic nature of the data into account.

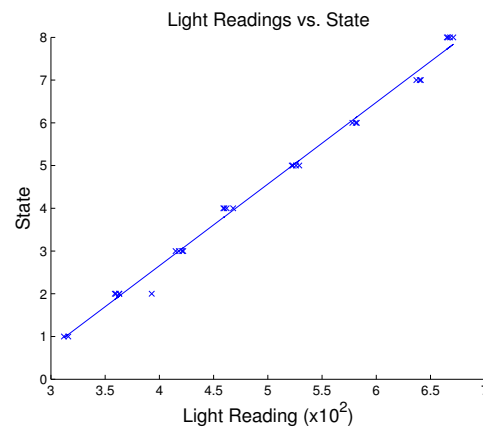
5.5.2 Calibrating state estimator with linear regression

One of the basic requirements for a robot or agent is to estimate the state it is in such as its position in space. In this tutorial we will use the color sensor and a floor with different shades of gray to help the robot to determine its position on a surface. Specifically, we use a sheet of gray bars of various gradients that is given in file `StateSheet.pdf`. The NXT can read this model as each gray bar will give a different light reading when the its light sensor is held over it. However, we need to pair each light reading with a number which represents the state value. Of course, the reading of the light sensors will fluctuate somewhat, we would like to repeat the measurements and find a function that translates each light reading to a state value.

To set up the tribot for this experiment, do the following:

1. Mount a switch sensors on the NXT.
2. Mount light sensor on front of the NXT, facing downward.
3. Because the state space used in this tutorial only involves driving forward and backward, it is a good idea to use additional lego pieces to lock the back wheel so it stays straight. If it is left to freely rotate, the NXT will start to turn and travel off the state sheet which will interfere with results.

4. Tape the state sheet to a table or other flat surface.
5. Setup an NXT motor object that will drive the robot forward in a straight line. Power should be between 20-30, TachoLimit around 20; motors should brake once it reaches the TachoLimit.
6. Create two matrices, one for light readings (from the NXT's light sensor), and the other for states. They must be the same size.
7. Have the NXT move forward, take a light reading from its sensor, and store it in the light reading matrix. The user should then be prompted to input which state the NXT is in. Store this value in the state matrix. Repeat this until you have a few readings from every state on the state sheet.
8. Using one of the methods discussed above, use linear regression to find a linear relationship between the light reading matrix and its associated state.



9. Test the data. Have the NXT take in a light sensor reading, and using the regressed line, interpolate the state and a round function to round the value to the nearest whole number). Does the regressed value give the correct state?

The state sheet will be used in chapter ?? when experimenting with reinforcement learning.

5.6 Non-linear regression and the bias-variance tradeoff

So far we only discussed the case where our hypothesis is linear in the feature values. We will now discuss regression with more general non-linear functions. As an example of data that do not seem to follow a linear trend is shown in Fig.5.6A. There, the number of transistors of microprocessors is plotted against the year each processor was introduced. We included a linear fit to show that this does not seem to describe the data well. Indeed, there seem to be systematic bias that suggest that we have to take more complex functions into account.

xxx Example of polynomial regression!

Finding the right function is one of the most difficult tasks, and there is not a simple algorithm that can give us the answer. This task is therefore an important area where experience, a good understanding of the problem domain, and a good

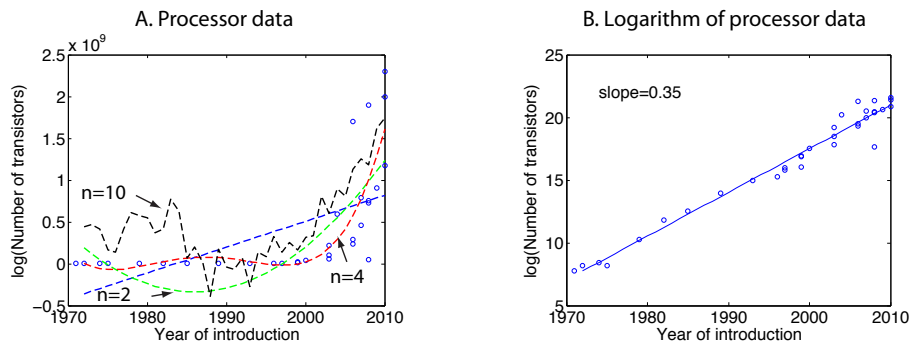


Fig. 5.6 Data from showing the number of transistors in microprocessors plotted the year that the processor was introduced. (A) Data and some linear and polynomial fits of the data. (B) Logarithm of the data and linear fit of these data.

understanding of scientific methods are required. This section does therefore try to give some recommendations when generalizing regression to the non-linear domain. These comments are important to understand for applying machine learning techniques in general.

It is often a good idea to visualize data in various ways since the human mind is often quite advanced in ‘seeing’ trends and patterns. Domain-knowledge thereby very valuable as specialists in the area from which the data are collected can give important advice or they might have specific hypothesis that can be investigated. It is also good to know common mechanisms that might influence processes. For example, the rate of change in basic growth processes is often proportional to the size of the system itself. Such a situation leads to exponential growth. (Think about why this is the case). Such situations can be revealed by plotting the functions on a logarithmic scale or the logarithm of the function as shown in Fig. 5.6B. A linear fit of the logarithmic values is also shown, confirming that the average growth of the number of transistors in microprocessors is exponential.

But how about more general functions. For example, we can consider a polynomial of order n , that can be written as

$$y = \theta_0 x^0 + \theta_1 x^1 + \theta_2 x^2 + \dots + \theta_n x^n \quad (5.39)$$

We can again use LMS regression to determine the parameters from the data by minimizing the LMS error function between the hypothesis and the data. This is implemented in Matlab as function `polyfit(x,y,n0)`. The LMS regression of the transistor data to a polynomial for orders $n = 2, 4, 10$ are shown in Fig. 5.6A as dashed lines.

A major question when fitting data with fairly general non-linear functions is the order that we should consider. The polynomial of order $n = 4$ seems to somewhat fit the data. However, notice there are systematic deviations between the curve and the data points. For example, all the data between years 1980 and 1995 are below the fitted curve, while earlier data are all above the fitted curve. Such a systematic **bias** is typical when the order of the model is too low. However, when we increase the order, then we usually get large fluctuations, or **variance**, in the curves. This fact is also

called **overfitting** the data since we have typically too many parameters compared to the number of data points so that our model starts describing individual data points with their fluctuations that we earlier assumed to be due to some noise in the system. This difficulty to find the right balance between these two effects is also called the **bias-variance tradeoff**.

The bias-variance tradeoff is quite important in practical applications of machine learning because the complexity of the underlying problem is often not known. It then becomes quite important to study the performance of the learned solutions in some detail. For this it is useful to split the data set into **training set**, which is used to estimate the parameters of the model, and a **validation set** that can be used to study the performance on data that have not been used in the training procedure, that is, how the machine performs in **generalizes**. A schematic figure showing the bias-variance tradeoff is shown in Fig.5.7. The plot shows the error rate as evaluated by the training data (dashed line) and validation curve (solid line) when considering models with different complexities. When the model complexity is lower than the true complexity of the problem, then it is common to have a large error both in the training set and in the evaluation due to some systematic bias. In the case when the complexity of the model is larger than the generative model of the data, then it is common to have a small error on the training data but a large error on the generalization data since the predictions are becoming too much focused on the individual examples. Thus varying the complexity of the data, and performing experiments such as training the system on different number of data sets or for different training parameters or iterations can reveal some of the problems of the models.

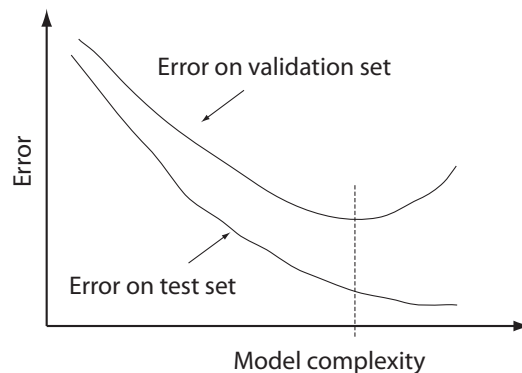


Fig. 5.7 Illustration of bias-variance tradeoff.

Using some of the data to validate the learning model is essential for many machine learning methods. An important question is then how many data we should keep to validate versus train the model. If we use too many data for validation, then we might have too less data for accurate learning in the first place. On the other end, if we have too few data for validation than this might not be very representative. In practice we are often using some **cross-validation** techniques to minimize the trade-off. That is, we use the majority of the data for training, but we repeat the selection of the validation data several times to make sure that the validation was not just a result of outliers.

The repeated division of the data into a training set and validation set can be done in different ways. For example, in **random subsampling** we just use random subsample for each set and repeat the procedure with other random samples. More common is **k -fold cross-validation**. In this technique we divide the data set into k -subsamples samples and use $k - 1$ subsamples for training and one subsample for validation. In the next round we use another subsample for validating the training. A common choice for the number of subsamples is $k = 10$. By combining the results for the different runs we can often reduce the variance of our prediction while utilizing most data for learning.

We can sometimes even help the learning process further. In many learning examples it turns out that some data are easy to learn while others are much harder. In some techniques called **boosting**, data which are hard to learn are over-sampled in the learning set so that the learning machine has more opportunities to learn these examples. A popular implementation of such an algorithm is **AdaBoost** (adaptive Boosting).

Exercise

1. Compare the analytic solution to the numerical solutions of random search and gradient descent search in the linear regression of the health data in section 5.1.
2. Is a linear fit of the logarithm of a function equal to an exponential fit of the original function? Discuss.
3. Explain briefly how a binary classification method can be used to support multiclass classification of more than two classes.

6 Classification

This chapter follows closely: <http://cs229.stanford.edu/notes/cs229-notes2.pdf>

6.1 Logistic regression

An important special case of learning problems is **classification** in which features are mapped to a finite number of possible categories. We discuss in this section **binary classification**, which is the case of two target classes where the target function (y -values) has only two possible values such as 0 and 1.

Let us first consider a random number which takes the value of 1 with probability ϕ and the value 0 with probability $1 - \phi$ (the probability of being either of the two choices has to be 1.) Such a random variable is called Bernoulli distributed. Tossing a coin is a good example of a process that generates a Bernoulli random variable and we can use maximum likelihood estimation to estimate the parameter ϕ from such trials. That is, let us consider m tosses in which h heads have been found. The log-likelihood of having h heads ($y = 1$) and $1 - h$ tails ($y = 0$) is

$$l(\phi) = \log(\phi^h (1 - \phi)^{m-h}) \quad (6.1)$$

$$= h \log(\phi) + (m - h) \log(1 - \phi). \quad (6.2)$$

To find the maximum with respect to ϕ we set the derivative of l to zero,

$$\frac{dl}{d\phi} = \frac{h}{\phi} - \frac{m-h}{1-\phi} \quad (6.3)$$

$$= \frac{h}{\phi} - \frac{m-h}{1-\phi} \quad (6.4)$$

$$= 0 \quad (6.5)$$

$$\rightarrow \phi = \frac{h}{m} \quad (6.6)$$

As you might have expected, the maximum likelihood estimate of the parameter ϕ is the fraction of heads in m trials.

Now let us discuss the case when the probability of observing a head or tail, the parameter ϕ , depends on an attribute x , as usual in a stochastic (noisy) way. More specifically,

$$h(y = 1|\mathbf{x}; \theta) = \hat{h}(\mathbf{x}; \theta) \quad (6.7)$$

$$h(y = 0|\mathbf{x}; \theta) = 1 - \hat{h}(\mathbf{x}; \theta) \quad (6.8)$$

which we can combine as

$$h(y|\mathbf{x}; \theta) = (\hat{h}(\mathbf{x}; \theta))^y (1 - \hat{h}(\mathbf{x}; \theta))^{1-y} \quad (6.9)$$

The corresponding log-likelihood function is

$$l(\theta) = \sum_{i=1}^m y^{(i)} \log(\hat{h}(\mathbf{x}; \theta)) + (1 - y^{(i)}) \log(1 - \hat{h}(\mathbf{x}; \theta)). \quad (6.10)$$

To find the corresponding maximum we can use the gradient ascent algorithm, which is equivalent to eq 5.5 with a changed sign,

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} l(\theta). \quad (6.11)$$

To calculate the gradient we can calculate the partial derivative of the log-likelihood function with respect to each parameters,

$$\frac{\partial l(\theta)}{\partial \theta_j} = \left(y \frac{1}{\hat{h}} - (1 - y) \frac{1}{1 - \hat{h}} \right) \frac{\partial \hat{h}(\theta)}{\partial \theta_j} \quad (6.12)$$

where we dropped indices for better readability.

An example is illustrated in Fig.6.1 with 100 examples plotted with star symbols. The data suggest that it is far more likely that the class is $y = 0$ for small values of x and that the class is $y = 1$ for large values of x , and the probabilities are more similar in between. We put forward the hypothesis that the transition between the low and high probability region is smooth and qualify this hypothesis as parameterized density function known as a **logistic** or **sigmoid** function

$$\hat{h} = g(\theta^T \mathbf{x}) = \frac{1}{1 + \exp(-\theta^T \mathbf{x})}. \quad (6.13)$$

As before, we can then treat this density function as function of the parameters θ for the given data values (likelihood function), and use maximum likelihood estimation to estimate values for the parameters so that the data are most likely. The density function with sigmoidal offset $\theta_0 = 2$ and slope $\theta_1 = 4$ is plotted as solid line in Fig.6.1.

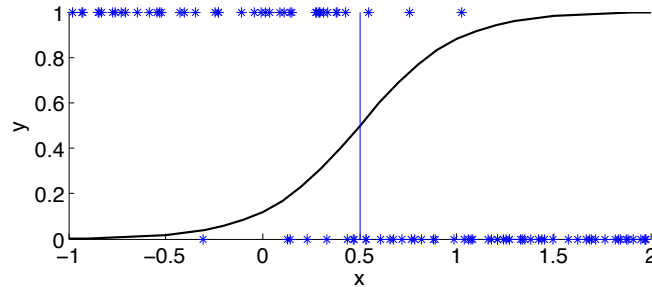


Fig. 6.1 Binary random numbers (stars) drawn from the density $p(y = 1) = \frac{1}{1 + \exp(-\theta_1 x - \theta_0)}$ (solid line).

We can now calculate the derivative of the hypothesis \hat{h} with respect to the parameters for the specific choice of the logistic functions. This is given by

$$\frac{\partial \hat{h}}{\partial \theta} = \frac{\partial}{\partial \theta} \frac{1}{1 + e^{-\theta x}} \quad (6.14)$$

$$= \frac{1}{(1 + e^{-\theta x})^2} e^{-\theta x} \quad (6.15)$$

$$= \frac{1}{(1 + e^{-\theta x})} \left(1 - \frac{1}{(1 + e^{-\theta x})}\right) \quad (6.16)$$

$$= \hat{h}(1 - \hat{h}) \quad (6.17)$$

Using this in equation 6.12 and inserting it into equation 6.11 give the learning rule

$$\theta_j \leftarrow \theta_j + \alpha (y^{(i)} - \hat{h}(\mathbf{x}^{(i)}, \theta)) x_j^{(i)} \quad (6.18)$$

How can we use the knowledge (estimate) of the density function to do classification? The obvious choice is to predict the class with the higher probability, given the input attribute. This **bayesian decision point**, x^t , or **dividing hyperplane** in higher dimensions, is give by

$$p(y = 1|x^t) = p(y = 0|x^t) = 0.5 \rightarrow x^t \theta^T \mathbf{x}^t = 0. \quad (6.19)$$

We have here considered binary classification with linear decision boundaries as logistic regression, and we can also generalize this method to problems with non-linear decision boundaries by considering hypothesis with different functional forms of the decision boundary. However, coming up with specific functions for boundaries is often difficult in practice, and we will discuss much more practical methods for binary classification later in this course.

6.2 Generative models

In the previous sections we have introduced the idea that understanding the world should be based on a model of the world in a probabilistic sense. That is, building knowledge about the world really means estimating a large density function about the world. So far we have used such stochastic model mainly for a recognition model that take feature values \mathbf{x} and make a prediction of an output y . Given the stochastic nature of the systems we want to model, the models where formulated as parameterized functions that represent the conditional probability $p(y|\mathbf{x}; \theta)$. Of course, learning such models is a big task. Indeed, we had to assume that we know already the principle form of the distribution, and we used only simple model with low-dimensional feature vectors. The learning tasks of humans to be able to function in the real world seems much more daunting, and even training a robot in more restricted environment seems still beyond our current ability. While the previous models illustrate the principle problem in supervised learning, much of the rest of this course discusses more practical methods.

In the last section we discussed a classification task where the aim of the model was to discriminate between classes based on the feature values. Such models are called

discriminative models because they try to discriminate between possible outcomes based on the input values. Building a discriminative model directly from example data can be a daunting task as we have to learn how each item is distinguished from every other possible item. A different strategy, which seems much more resembling human learning, is to learn first about the nature of specific classes and then use this knowledge when faced with a classification task. For example, we might first learn about chairs, and independently about tables, and when we are shown pictures with different furnitures we can draw on our knowledge to classify them. Thus, in this chapter we start discussing **generative models** of individual classes, given by $p(\mathbf{x}|y; \theta)$.

Generative models can be useful in its own right, and are also important to guide learning as discussed later, but for now we are mainly interested in using these models for classification. Thus, we need to ask how we can combine the knowledge about the different classes to do classification. Of course, the answer is provided by Bayes' theorem. In order to make a discriminative model from the generative models, we need to the **class priors know**, e.g. what the relative frequencies of the classes is, and can then calculate the probability that an item with features \mathbf{x} belong to a class y as

$$p(y|\mathbf{x}; \theta) = \frac{p(\mathbf{x}|y; \theta)p(y)}{p(\mathbf{x})}. \quad (6.20)$$

We can use this directly in the case of classification. The Bayesian decision criterion of predicting the class with the largest posterior probability is then:

$$\arg \max_y p(y|\mathbf{x}; \theta) = \arg \max_y \frac{p(\mathbf{x}|y; \theta)p(y)}{p(\mathbf{x})} \quad (6.21)$$

$$= \arg \max_y p(\mathbf{x}|y; \theta)p(y), \quad (6.22)$$

where we have used the fact that the denominator does not depend on y and can hence be ignores. In the case of binary classification, this reads:

$$\arg \max_y p(y|\mathbf{x}; \theta) = \arg \max_y (p(\mathbf{x}|y = 0; \theta)p(y = 0) + p(\mathbf{x}|y = 1; \theta)p(y = 1)). \quad (6.23)$$

While using generative models for classification seem to be much more elaborate, we will see later that there are several arguments which make generative models attractive for machine learning, and we will argue that generative models are do more closely resemble human brain processing principles.

6.3 Discriminant analysis

We will now discuss some common examples of using generative models in classification. The methods in this section go back to a paper by R. Fisher in 1936. In the following examples we consider that there are k classes, and we first assume that each class has members which are Gaussian distribution over the n feature value. An example for $n = 2$ is shown in Fig.??A.

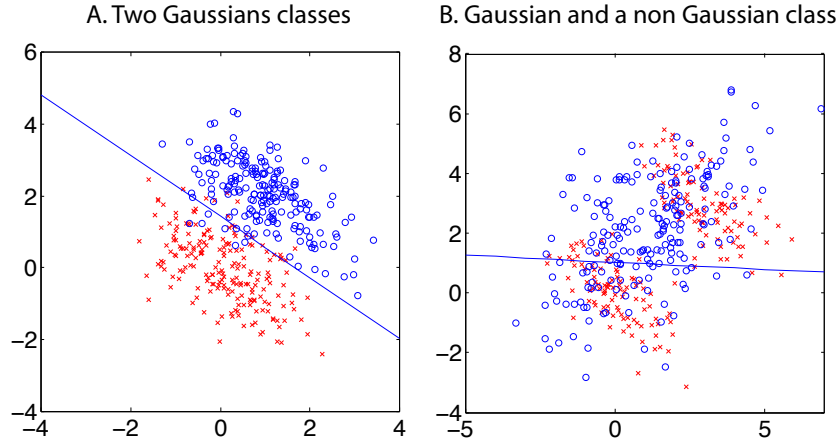


Fig. 6.2 Linear Discriminant analysis on a two class problem with different class distributions.

Each of the classes have a certain class prior

$$p(y = k) = \phi_k \quad (6.24)$$

, and each class itself is multivariate Gaussian distributed, generally with different means, μ_k and variances, Σ_k ,

$$p(\mathbf{x}|y = k) = \frac{1}{\sqrt{2\pi}^n \sqrt{|\Sigma_k|}} e^{-\frac{1}{2}(\mathbf{x}-\mu_k)^T \Sigma_k^{-1}(\mathbf{x}-\mu_k)} \quad (6.25)$$

$$(6.26)$$

Since we have supervised data with examples for each class, we can use maximum likelihood estimation to estimate the most likely values for the parameters $\theta = (\phi_k, \mu_k, \Sigma_k)$. For the class priors, this is simply the relative frequency of the training data,

$$\phi_k = \frac{1}{m} \sum_{i=1}^m \mathbb{1}(y^{(i)} = k) \quad (6.27)$$

where the function $\mathbb{1}(y^{(i)} = k) = 1$ if the i th example belongs to class k , and $\mathbb{1}(y^{(i)} = k) = 0$ otherwise. The estimates of the means and variances within each class are given by

$$\mu_k = \frac{\sum_{i=1}^m \mathbb{1}(y^{(i)} = k) \mathbf{x}^{(i)}}{\sum_{i=1}^m \mathbb{1}(y^{(i)} = k)} \quad (6.28)$$

$$\Sigma_k = \frac{\sum_{i=1}^m \mathbb{1}(y^{(i)} = k) (\mathbf{x}^{(i)} - \mu_{y^{(i)}})(\mathbf{x}^{(i)} - \mu_{y^{(i)}})^T}{\sum_{i=1}^m \mathbb{1}(y^{(i)} = k)}. \quad (6.29)$$

With these estimates, we can calculate the optimal (in a Bayesian sense) decision rule, $G(x; \theta)$, as a function of \mathbf{x} with parameters θ , namely

$$G(x) = \arg \max_k p(y = k|\mathbf{x}) \quad (6.30)$$

$$= \arg \max_k [p(\mathbf{x}|y = k; \theta)p(y = k)] \quad (6.31)$$

$$= \arg \max_k [\log(p(\mathbf{x}|y = k; \theta)p(y = k))] \quad (6.32)$$

$$= \arg \max_k [-\log(\sqrt{2\pi}^n \sqrt{|\Sigma_0|}) - \frac{1}{2}(\mathbf{x} - \mu_k)^T \Sigma_k^{-1}(\mathbf{x} - \mu_k) + \log(\phi_k)] \quad (6.33)$$

$$= \arg \max_k [-\frac{1}{2}\mathbf{x}^T \Sigma_k^{-1} \mathbf{x} - \frac{1}{2}\mu_k^T \Sigma_k^{-1} \mu_k + \mathbf{x}^T \Sigma_k^{-1} \mu_k + \log(\phi_k)], \quad (6.34)$$

since the first term in equation 6.33 does not depend on k and we can multiply out the other terms. With the maximum likelihood estimates of the parameters, we have all we need to make this decision.

In order to calculate the decision boundary between classes l and k , we make the common additional assumption that the covariance matrices of the classes are the same,

$$\Sigma_k =: \Sigma. \quad (6.35)$$

The decision boundary is then

$$\log\left(\frac{\phi_k}{\phi_l}\right) - \frac{1}{2}(\mu_k - \mu_l)^T \Sigma^{-1}(\mu_k - \mu_l) - \mathbf{x} \Sigma^{-1}(\mu_k - \mu_l) = 0. \quad (6.36)$$

The first two terms do not depend on x and can be summarized as constant \mathbf{a} . We can also introduce the vector

$$\mathbf{w} = -\Sigma^{-1}(\mu_k - \mu_l). \quad (6.37)$$

With these simplifying notations it is easy to see that this decision boundary is a linear,

$$\mathbf{a} + \mathbf{w}\mathbf{x} = 0, \quad (6.38)$$

and this method with the Gaussian class distributions with equal variances is called **Linear Discriminant Analysis (LDA)**. The vector \mathbf{w} is perpendicular to the decision surface. Examples are shown in Figure ???. If we do not make the assumption of equal variances of the classes, then we have a quadratic equation for the decision boundary, and the method is then called **Quadratic Discriminant Analysis (QDA)**. With the assumptions of LDA, we can calculate the contrastive model directly using Bayes rule.

$$p(y = k|\mathbf{x}; \theta) = \frac{\phi_k \frac{1}{\sqrt{2\pi}^n \sqrt{|\Sigma|}} e^{-\frac{1}{2}(\mathbf{x} - \mu_k)^T \Sigma_k^{-1}(\mathbf{x} - \mu_k)}}{\phi_k \frac{1}{\sqrt{2\pi}^n \sqrt{|\Sigma|}} e^{-\frac{1}{2}(\mathbf{x} - \mu_k)^T \Sigma_k^{-1}(\mathbf{x} - \mu_k)} + \phi_l \frac{1}{\sqrt{2\pi}^n \sqrt{|\Sigma|}} e^{-\frac{1}{2}(\mathbf{x} - \mu_l)^T \Sigma_l^{-1}(\mathbf{x} - \mu_l)}} \quad (6.39)$$

$$= \frac{1}{1 + \frac{\phi_l}{\phi_k} \exp^{-\theta^T \mathbf{x}}}, \quad (6.40)$$

where θ is an appropriate function of the parameters μ_k , μ_l , and Σ . Thus, the contrastive model is equivalent to logistic regression discussed in the previous chapter, although we use parametrisations and the two methods will therefore usual give different results on specific data sets. So which method should be used? In LDA we made the assumption that each class is Gaussian distributed. If this is the case, then LDA is the best method we can use. Discriminant analysis is also popular since it often works well even when the classes are not strictly gaussian. However, as can be seen in Figure ??B,

it can also produce quite bad results if the data are multimodal distributed. Logistic regression is somewhat more general since it does not make the assumption that the class distributions are Gaussian. However, so far we have mainly looked at linear models and logistic regression would have also problems with the data shown in Figure ??B.

Finally, we should note that Fisher's original method was slightly more general than the examples discussed here since he did not assume Gaussian distributions. Instead considered within-class variances compared to between-class variances, something which resembles a signal-to-noise ratio. In **Fisher discriminant analysis (FDA)**, the separating hyperplane is defined as

$$\mathbf{w} = -(\Sigma_k + \Sigma_l)^{-1}(\mu_k - \mu_l). \quad (6.41)$$

which is the same as in LDA in the case of equal covariance matrices.

6.4 Naive Bayes

In the previous example we used mainly a two dimensional feature vector to discuss a classification problems as generative model. We will give an example that has a much higher dimensionality.

In the following example we want to make a spam filter that classifies email messages as either spam ($y = 1$) or non-spam ($y = 0$) emails. To do this we need first a method to represent the problem in a suitable way. We chose here to represent a text (email in this situation) as **vocabulary** vector. A vocabulary is simply the list of all possible words that we consider, and the text is represented by this vector with entries 1 if the word can be found in the list or an entry 0 if not, e.g.

$$\mathbf{x} = \begin{pmatrix} 1 \\ 0 \\ 0 \\ \cdot \\ \cdot \\ \cdot \\ 1 \\ \cdot \\ \cdot \\ 0 \end{pmatrix} \begin{matrix} \text{a} \\ \text{aardvark} \\ \text{aardwolf} \\ \cdot \\ \cdot \\ \cdot \\ \text{buy} \\ \cdot \\ \cdot \\ \text{zygmurgy} \end{matrix} \quad (6.42)$$

We are here only considering values 0 and 1 instead, for example, counting how often the corresponding word appears. The difference is that each entry is a binomial random variable and would be a multinomial in the other example, though the methods generalize directly to the other case. Note that this feature vector is typically very high dimensional. Let us consider here that our vocabulary has 50.000 word, which is a typical size of common languages.

We now want to build a discriminative model from some training examples. That is, we want to model

$$p(\mathbf{x}|y) = p(x_1, x_2, \dots, x_{50000}|y). \quad (6.43)$$

This is a very high dimensional density function which has $2^{50.000} - 1$ parameters (the -1 comes from the normalization condition). We can factorize this conditional density function with the chain rule

$$p(x_1, x_2, \dots, x_{50000}|y) = p(x_1|y)p(x_2|y, y_1)\dots p(x_{50000}|y, x_1, \dots, x_{49999}). \quad (6.44)$$

While the right hand side has only 50.000 factors, there are still $2^{50.000} - 1$ parameters we have to learn. However, we no make a strong assumption namely that all the words are conditionally independent in each text, that is,

$$p(x_1|y)p(x_2|y, y_1)\dots p(x_{50000}|y, x_1, \dots, x_{49999}) = p(x_1|y)p(x_2|y)\dots p(x_{50000}|y). \quad (6.45)$$

This is called the **Naive Bayes (NB) assumption**. Hence, we can write the conditional probability as a factor of terms with 50.000 parameters

$$p(\mathbf{x}|y) = \prod_{i=1}^{50000} p(x_i|y). \quad (6.46)$$

To estimate these parameters we can apply again maximum likelihood estimation, which gives

$$\phi_{j,y=1} = \frac{\sum_{i=1}^m 1\{x_j^{(i)} = 1 \wedge y^{(i)} = 1\}}{\sum_{i=1}^m 1\{y^{(i)} = 1\}} \quad (6.47)$$

$$\phi_{j,y=0} = \frac{\sum_{i=1}^m 1\{x_j^{(i)} = 1 \wedge y^{(i)} = 0\}}{\sum_{i=1}^m 1\{y^{(i)} = 0\}} \quad (6.48)$$

$$\phi_{y=1} = \frac{\sum_{i=1}^m 1\{y^{(i)} = 1\}}{m}. \quad (6.49)$$

With these parameters we can now calculate the probability that email \mathbf{x} is spam as

$$p(y = 1|\mathbf{x}) = \frac{\prod_{i=1}^m \phi_{i,y=1} p h i_{y=1}}{\prod_{i=1}^m \phi_{i,y=1} p h i_{y=1} + \prod_{i=1}^m \phi_{i,y=0} p h i_{y=0}}. \quad (6.50)$$

In practice this often works well when the Naive Bayes assumption is appropriate, that is, if there are no strong correlation between the features. Finally, note that there is a slight problem if some of the words, say x_{100} , are not part of the training set. In this case we get an estimate that the probability of this word ever occurring is zero, $\phi_{100,y=1} = 0$ and $\phi_{100,y=0} = 0$, and hence $p(y = 1|x) = \frac{0}{0}$. A common trick, called **Laplace smoothing** is to add one occurrence of this word in every case, which will insert a small probability proportional to your training examples to the estimates,

$$\phi_{j,y=1} = \frac{\sum_{i=1}^m 1\{x_j^{(i)} = 1 \wedge y^{(i)} = 1\} + 1}{\sum_{i=1}^m 1\{y^{(i)} = 1\} + 2} \quad (6.51)$$

$$\phi_{j,y=0} = \frac{\sum_{i=1}^m 1\{x_j^{(i)} = 1 \wedge y^{(i)} = 0\} + 1}{\sum_{i=1}^m 1\{y^{(i)} = 0\} + 2}. \quad (6.52)$$

We will later compare the Naive Bayes classification with other classification methods.

7 Graphical models

7.1 Causal models

In the regression example that looked at health data relating the weight of subjects to running times, we considered the weight to be a random variable to capture the uncertainties in the data. There we treated the running time as the dependent variable. Since this is also a measured quantity, it should also be treated as random variable. In general we should anyhow consider more complex models with many more factors described as random variables, and we are most interested in describe the relations of these variables in a model.

We have seen that probability theory is quite handy to model data, and probability theory also considers multiple random variables. The total knowledge about the co-occurrence of specific values for two random variables X and Y is captured by the **joined probability function** $p(X, Y)$. This is a symmetric function,

$$p(X, Y) = p(Y, X), \quad (7.1)$$

There is an interesting limitation of joined density functions of multiple variables which only describe the co-occurrence of specific values of the random variables. However, we want to use our knowledge about the world, captured by a model, to reason about possible events. For this we want to add knowledge or hypotheses about **causal relations**. For example, a fire alarm should be triggered by a fire, although there is some small chance that the alarm will not sound when the unit is defect. However, it is (hopefully) unlikely that the sound of a fire alarm will trigger a fire. It is useful to illustrate such casual relations with graphs such as



In such **graphical models**, the nodes represent random variables, and the links between them represent causal relations with conditional probabilities, $p(A|F)$. Since we use arrows on the links we are discussing here **directed graphs**, and we are also restricting our discussions here to graphs that have no loops, so called **acyclic graphs**. **Directed acyclic graphs** are also called **DAGs**.

Graphical causal models have been advanced largely by Judea Pearl, and the following example is taken from his book⁶. The model is shown in Figure 7.1. Each of the five nodes stands for a random binary variable (Burglary $B=\{\text{yes,no}\}$, Earthquake $E=\{\text{yes,no}\}$, Alarm $A=\{\text{yes,no}\}$, JohnCalls $J=\{\text{yes,no}\}$, MaryCalls $M=\{\text{yes,no}\}$) The figure also include **conditional probability tables (CPTs)** that specify the conditional probabilities represented by the links between the nodes.

⁶Judea Pearl, 'Causality: Models, Reasoning and Inference', Cambridge University Press 2000, 2009'.

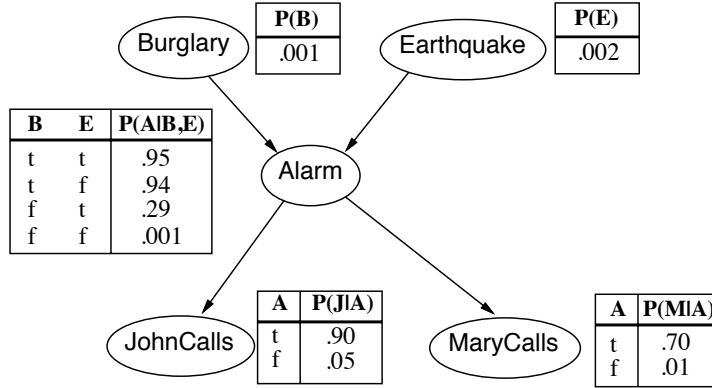


Fig. 7.1 Example of causal model a two-dimensional probability density function (pdf) and some examples of marginal pdfs.

The joined distribution of the five variables can be factories in various ways following the chain rule mentioned before (equations 2.30), for example as

$$p(B, E, A, J, M) = P(B|E, A, J, M)P(E|A, J, M)P(A|J, M)P(J|M)P(M) \quad (7.2)$$

However, the the causal model represents a specific factorization of the joined probability functions, namely

$$p(B, E, A, J, M) = P(B)P(E)P(A|B, E)P(J|A)P(M|A), \quad (7.3)$$

which is much easier to handle. For example, if we do not know the conditional probability functions, we need to run many more experiments to estimate the various conditions ($2^4 + 2^3 + 2^2 + 2^1 + 2^0 = 31$) instead of the reduced conditions in the causal model ($1 + 1 + 2^2 + 2 + 2 = 10$). It is also easy to use the casual model to do inference (drawing conclusions), for specific questions. For example, say we want to know the probability that there was no earthquake or burglary when the alarm rings and both John and Mary call. This is given by

$$\begin{aligned} P(B = f, E = f, A = t, J = t, M = t) &= \\ &= P(B = f)P(E = f,)P(A = t|B = f, E = f)P(J = t|A = t)P(M = t|A = t) \\ &= 0.998 * 0.999 * 0.001 * 0.7 * 0.9 \\ &= 0.00062 \end{aligned}$$

Although we have a casual model where parents variables influence the outcome of child variables, we can also use a child evidence to infer some possible values of parent evidence. For example, let us calculate the probability that the alarm rings given that John calls, $P(A = t|J = t)$. For this we should first calculate the probability that the alarm rings as we need this later. This is given by

$$\begin{aligned} P(A = t) &= P(A = t|B = t, E = t)P(B = t)P(E = t) + \dots \\ &P(A = t|B = t, E = f)P(B = t)P(E = f) + \dots \end{aligned}$$

$$\begin{aligned}
& P(A = t|B = f, E = t)P(B = f)P(E = t) + \dots \\
& P(A = t|B = f, E = f)P(B = f)P(E = f) \\
& = 0.95 * 0.001 * 0.002 + 0.94 * 0.001 * 0.998 + \dots \\
& 0.29 * 0.999 * 0.002 + 0.001 * 0.999 * 0.998 \\
& = 0.0025
\end{aligned}$$

We can then use Bayes' rule to calculate the required probability,

$$\begin{aligned}
P(A = t|J = t) &= \frac{P(J = t|A = t)P(A = t)}{P(J = t|A = t)P(A = t) + P(J = t|A = f)P(A = f)} \\
&= \frac{0.90.0025}{0.90.0025 + 0.050.9975} \\
&= 0.0434
\end{aligned}$$

We can similarly apply the rules of probability theory to calculate other quantities, but these calculations can get cumbersome with larger graphs. It is therefore useful to use numerical tools to perform such inference. A Matlab toolbox for Bayesian networks is introduced in the next section.

While inference is an important application of causal models, inferring causality from data is another area where causal models revolutionize scientific investigations. Many traditional methods evaluate co-occurrences of events to determine dependencies, such as a correlation analysis. However, such a correlation analysis is usually not a good indication of causality. Consider the example above. When the alarm rings it is likely that John and Mary call, but the event that John calls is mutually independent of the event that Mary calls. Yet, when John calls it is also statistically more likely to observe the event that Mary calls. Sometimes we might just be interested in knowing about the likelihood of co-occurrence, for which a correlation analysis can be a good start, but if we are interested in describing the causes of the observations, then we need another approach. Some algorithms have been proposed for **structural learning**, such as an algorithm called **inferred causation (IC)**, which deduces the most likely causal structure behind given data is.

7.2 Bayes Net toolbox

An Matlab implementation of various algorithms for inference, parameters estimation and inferred causation is provided by Kevin Murphy in the Bayes Net toolbox⁷. We will demonstrate some of its features on the burglary/earthquake example above.

The first step is to create a graph structure for the DAG We have five nodes. The nodes are given numbers, but we also use variables with capital letter names to refer to them. The DAG is then a matrix with entries 1 where directed links exist.

```

N=5;% number of nodes
B=1; E=2; A=3; J=4; M=5;
dag = zeros(N,N);

```

⁷The toolbox can be downloaded at <http://code.google.com/p/bnt>.

```

dag(B,A)=1;
dag(E,A)=1;
dag(A,[J M])=1;

```

The nodes represent discrete random variables with two possible state. We only discuss here discrete random variables, although the toolbox contains methods for continuous random variables. For the discrete case we have to specify the number of possible states of each variable, and we can then create the corresponding Bayesian network,

```

% Make bayesian network
node_sizes=[2 2 2 2 2]; %binary nodes
bnet=mk_bnet(dag,node_sizes); %make bayesian net

```

The next step is to provide the numbers for the conditional probability distributions, which are the conditional probability tables for discrete variables. For this we provide the numbers in a vector according to the following convention. Say we specify the probabilities for node 3, which is conditionally dependent on nodes 1 and 2. We then provide the probabilities in the following order:

Node 1	Node 2	P(Node 3=X)
F	F	F
T	F	F
F	T	F
T	T	F
F	F	T
T	F	T
F	T	T
T	T	T

For our specific examples, the CPT are thus specified as

```

bnet.CPD{B} = tabular_CPD(bnet,B,[0.999 0.001]);
bnet.CPD{E} = tabular_CPD(bnet,E,[0.998 0.002]);
bnet.CPD{A} = tabular_CPD(bnet,A,[0.999 0.06 0.71 0.05 0.001 0.94 0.29 0.95]);
bnet.CPD{J} = tabular_CPD(bnet,J,[0.95 0.10 0.05 0.90]);
bnet.CPD{M} = tabular_CPD(bnet,M,[0.99 0.30 0.01 0.70]);

```

We are now ready to calculate some inference. For this we need to specify a specific inference engine. There are several algorithms implemented, a variety of exact algorithm as well as approximate procedures in case the complexity of the problem is too large. Here we use the basic exact inference engine, the **junction tree algorithm**, which is based on a message passing system.

```

engine=jtree_inf_engine(bnet);

```

While this is an exact inference engine, there are other engines, such as approximate engines, that might be employed for large graphs when other methods fail.

As an example of an inference we recalculate the example above, that of calculate the probability that the alarm rings given that John calls, $P(A = t | J = t)$. For this we have to enter some evidence, namely that $J = t$, into a cell array and add this to the inference engine,

```

evidence=cell(1,N);
evidence{J}=2;

```

```
[engine,loglik]=enter_evidence(engine,evidence)
```

We can then calculate the marginal distribution for a variable, given the evidence as,

```
marg=marginal_nodes(engine,A)
```

```
p=marg.T(2)
```

It is now very easy to calculate other probabilities.

The Bayesnet toolbox also includes routines to handle some continuous models such as models with Gaussian nodes. In addition, there are routines to do parameter estimation, including point estimates, such maximum likelihood estimation, and also full bayesian priors. Finally, he toolbox includes routines to inferred causation through structural learning.

Exercise:

In the above example, shown in Figure 7.1, calculate the probability that there was a burglary, given that John and Marry called.

7.3 Temporal Bayesian networks: Markov Chains and Bayes filters

In many situations we are interested in how things change over time. For this we need to build time-dependent causal models, also called **dynamic bayesian models (DBNs)**. We will discuss here a specific family of such models in which the states at one time, t , only depend on the random variables at the previous time, $t - 1$. This condition is called a **Markov condition**, and the corresponding models are called **markov chains**. A very common situation is captured in the example shown in Figure 7.2. In this example we have three time-dependent random variables called $u(t)$, $x(t)$, and $z(t)$. We used different grey shades for the nodes to indicate if they are observed or not. Only $u(t)$ and $z(t)$ are observed, whereas $x(t)$ is an un-observed, **hidden** or **latent** variable. Such Markov chains are called **Hidden Markov Models (HMMs)**.

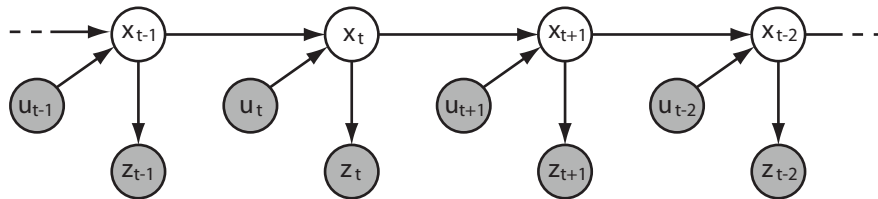


Fig. 7.2 A hidden markov model (HMM) with time dependent state variable x , motor control u and sensor readings z .

The HMM shown in figure 7.2 capture the basic operation of a mobile robot in which the variable $u(t)$ represents the motor command at time t and the variable

$z(t)$ represents sensor readings at time t . The motor command is the cause that the robot goes into a new state $x(t)$. The problem illustrated in the figure is that the new state of the robot, such as its location or posture, can not be observed directly. In the following we will specifically discuss **robot localization**, where we must infer the possible state from the motor command and sensor reading. We have previously used a state sheet to ‘measure’ the state (location) of the robot with the light sensor (see section 5.5.2). The model above allows us a much more sophisticated guess of the location by taking not only the current reading into account, but by combining this information with our previous guess of the location and the knowledge of the motor command. More specifically, we want to calculate the **believe** of the state, which is a probability (density) function over the state space

$$bel(x_t) = p(x_t | bel(x_0), u_{1:t}, z_{1:t}), \quad (7.4)$$

give an initial believe and the history of motor command and sensor reading. This believe can be calculated recursively from the previous believes and the new information, which is a form of **believe propagation**. It is often convenient to break this process down into calculating a new **prediction** from the previous believe and the motor command, and then to add the knowledge provided by the new sensor reading. To predict of probability of a new state from the previous believes and the current motor command, we need to marginalize over the previous states, that is, we need to multiply the previous believe with the probability of the new state, given the previous state and motor command, and to sum (integrate) over it,

$$pred(x_t) = \frac{1}{N_x} \sum_{x_{t-1}} p(x_t | u_t, x_{t-1}) bel(x_{t-1}), \quad (7.5)$$

where N_x is the number of states. This prediction can be combined with the sensor reading to give a new believe,

$$bel(x_t) = \frac{p(z_t | x_t) pred(x_t)}{\sum_{x_t} p(z_t | x_t) pred(x_t)}, \quad (7.6)$$

where we included the normalization over all possible states to get a number representing probabilities. This is called **Bayes filtering**. With a Bayes filter we can calculate the new believe from the previous believe, the current motor command, and the latest sensor reading. This is the best we can do with the available knowledge to estimate the positions of a system. The application of Bayes filtering to a robot localization is also called **Markov localization** and is illustrated in Figure 7.3. The limitation in practice is often the that we have to sum (integrate) over all the possible states, which might be a very large sum if it has to be done explicitly. In some cases we can do this analytically, as shown in the next sections.

7.3.1 The Kalman filter

We have, so far, only assumed that our process complies with the Markov condition in that the new state only depends on the previous state and current motor command.

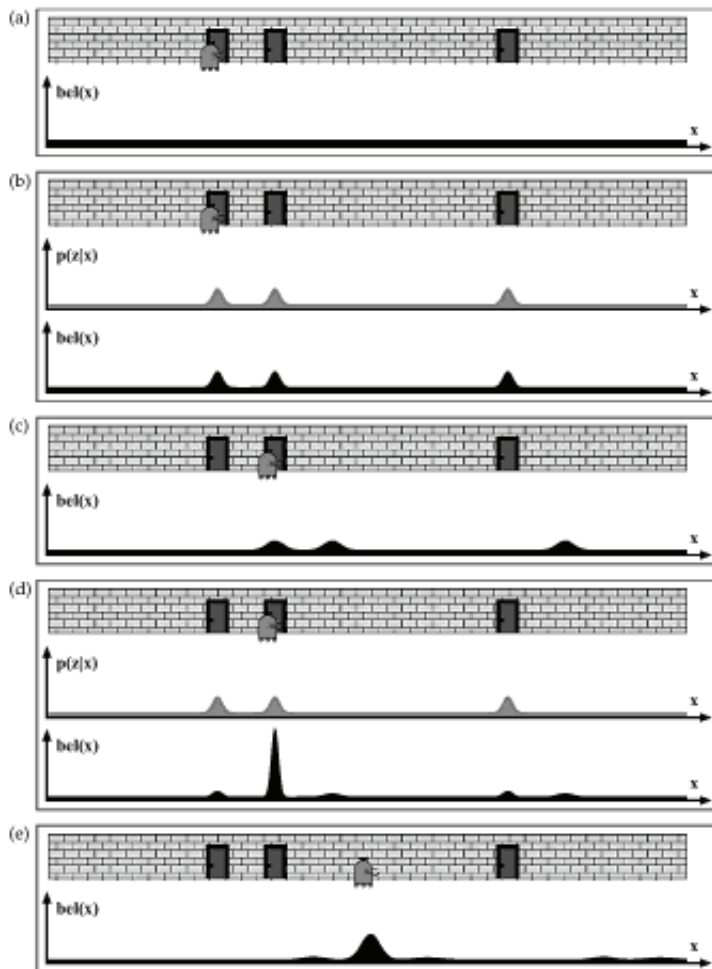


Fig. 7.3 A illustration of markov localization [from *Probabilistic model*, Thrun, Burgard, Fox, MIT press 2006].

In many cases we might have a good estimate of the state at some point with some variance. We now assume that our believes are Gaussian distributed,

$$p(\mathbf{x}_{t-1}) = \frac{1}{(\sqrt{2\pi})^n \sqrt{\det(\boldsymbol{\Sigma}_{t-1})}} \exp\left(-\frac{1}{2}(\mathbf{x}_{t-1} - \boldsymbol{\mu}_{t-1})^T \boldsymbol{\Sigma}_{t-1}^{-1} (\mathbf{x}_{t-1} - \boldsymbol{\mu}_{t-1})\right). \quad (7.7)$$

We also consider first the case where the transition probability $p(\mathbf{x}_t | \mathbf{u}_t, \mathbf{x}_{t-1})$ is linear in the previous state and the motor command, up to Gaussian noise ϵ_t ,

$$\bar{\mathbf{x}}_t = \mathbf{A}_t \mathbf{x}_{t-1} + \mathbf{B}_t \mathbf{u}_t + \epsilon_t; \quad \epsilon_t \sim N(0, \mathbf{Q}_t), \quad (7.8)$$

where \mathbf{A}_t and \mathbf{B}_t are time dependent matrices. The gaussian noise ϵ_t has thereby mean zero and covariance \mathbf{Q}_t . This setting is very convenient since the posterior

after moving a Gaussian uncertainty in such a linear fashion is again a Gaussian. The believe after incorporating the movement is hence a Gaussian probability over states with parameters

$$\bar{\mu}_t = \mathbf{A}_t \mu_{t-1} + \mathbf{B}_t \mathbf{u}_t \quad (7.9)$$

$$\bar{\Sigma}_t = \mathbf{A}_t \Sigma_{t-1} \mathbf{A}_t^T + \mathbf{Q}_t. \quad (7.10)$$

We now need to take the measurement probability into account, $P(\mathbf{z}_t | \bar{\mathbf{x}}_t)$, which we also assume to be linear in its argument up to a Gaussian noise δ_t ,

$$\mathbf{z}_t = \mathbf{C}_t \bar{\mathbf{x}}_t + \delta_t; \quad \delta_t \sim N(0, \mathbf{R}_t). \quad (7.11)$$

With this measurement update, we can calculate our new state estimate, parameterized by μ_t and Σ_t , as

$$\bar{\mathbf{K}}_t = \bar{\Sigma}_t \mathbf{C}_t^T (\mathbf{C}_t \bar{\Sigma}_t \mathbf{C}_t^T + \mathbf{R}_t)^{-1} \quad (7.12)$$

$$\mu_t = \bar{\mu}_t + \bar{\mathbf{K}}_t (\mathbf{z}_t - \mathbf{C}_t \bar{\mu}_t) \quad (7.13)$$

$$\Sigma_t = (\mathbf{I} - \bar{\mathbf{K}}_t \mathbf{C}_t) \bar{\Sigma}_t, \quad (7.14)$$

where \mathbf{I} is the identity matrix.

The basic Kalman filter makes several simplifying assumptions such as Gaussian believes, Gaussian noise, and linear relations. There are many generalizations of this method. For example, it is possible to extend the method to non-linear transformations while still demanding that the posteriors are Gaussian. While this introduces some errors and is this only an approximate method, this **extended Kalman filter (EKF)** is often very useful in practical applications. Also, there are several non-parametric methods such as **particle filters**.

7.4 Application and generalization: Localisation example

Include robot example with Karman filter, mention EKF, and show particle filters.

Exercises:

1. A mobile robot is instructed to travel 3 meters along a line in each time step. The true distance traveled is gaussian distributed around the intended position with standard deviation of 2 meters; the positions sensors are also unreliable with Gaussian that has a standard deviation of 4 meters; What is the average absolute difference between the true position and the estimated position when using only the sensor information, only the information inherent in the motor command, and both sources of information?
2. The lego robot is traveling along a line that has dark stripes at various positions. These distance between the stripes stripes are doubling for each consecutive stipe ($d(i) = 2 * d(i - 1)$). Estimate the position of the Lego tribot when traveling straight ahead if the robot is positioned at a random initial position.

8 General learning machines

The previous models, through the formulation of specific hypothesis functions, have been designed specifically for each applications. However, much more practical would be to have more general machines that can learn without making very specific functional assumptions. But how can we do this? The general idea is to provide a very general functions with many parameters that will be adjusted through learning. Of course, the real problem is then to not over-fit the model by using appropriate restrictions and also to make the learning efficient so that it can be used to large problem size. This chapter starts with a brief historical introduction to general learning machines and neural networks. We then discuss support vector machines and more rigorous learning theories.

8.1 The Perceptron

There was always a strong interest of AI researchers in **real intelligence**, that is, to understand the human mind. For example, both Alan Turing and John von Neumann worked more directly on biological systems in their last years before their early passing, and human behaviour and the brain have always been of interest to AI researchers. Of course, we want to understand how the brain is working in its own right, but it is also a great example of a quite general and successful learning machine.

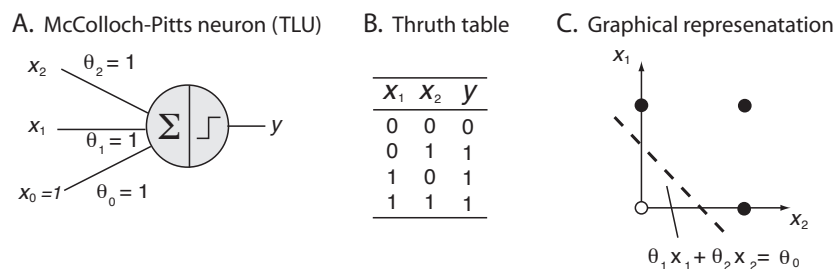


Fig. 8.1 Representation of the boolean OR function with a McCulloch-Pitts neuron (TLU).

A seminal paper, which has greatly influenced the development of early learning machines, is the 1943 paper by Warren McCulloch and Walter Pitts. In this paper, they proposed a simple model of a neuron, called the **threshold logical unit**, now often called the **McCulloch–Pitts neuron**. Such a unit is shown in Fig. 8.1A with three input channels, although it could have an arbitrary number of input channels. Input values are labeled by x with a subscript for each channel. Each channel has also a **weight**

parameter, θ_i . The McCulloch–Pitts neuron operates in the following way. Each input value is multiplied with the corresponding weight value, and these weighted values are then summed. If the weighted summed input is larger than a certain threshold value, $-\theta_0$, then the output is set to one, and zero otherwise, that is,

$$h(\mathbf{x}; \theta) = \begin{cases} 1 & \text{if } \sum_{i=0}^n \theta_i x_i = \theta^T \mathbf{x} > 0 \\ 0 & \text{otherwise} \end{cases}. \quad (8.1)$$

Such an operation resembles, to some extent, a neuron in that a neuron is also summing synaptic inputs and fires (has a spike in its membrane potential) when the membrane potential reaches a certain level that opens special voltage-gated ion channels. McCulloch and Pitts introduced this unit as a simple neuron model, and they argued that such a unit can perform computational tasks resembling boolean logic. This is demonstrated in Fig. 8.1 for a threshold unit that implements the Boolean OR function. The symbol h is used in these lecture notes since the output of this neuron represents the **hypothesis** that this neuron implements given the parameters θ . Also note that the non-linear step-function used in this neuron model corresponds to hypothesis for classification.

The next major developments in this area were done by Frank Rosenblatt and his engineering colleague Charles Wightman (Fig. 8.2), using such elements to build a machine that Rosenblatt called the **perceptron**. As can be seen in the figures, they worked on a machine that can perform letter recognition, and that the machine consisted of a lot of cables, forming a network of simple, neuron-like elements.

The most important challenge for the team was to find a way how to adjust the parameters of the model, the connection weights θ_i , so that the perceptron would perform a task correctly. The procedure was to provide to the system a number of examples, let's say m input data, $\mathbf{x}^{(i)}$ and the corresponding desired outputs, $y^{(i)}$. The procedure they used thus resembles supervised learning. The learning rule they used is called the **perceptron learning rule**,

$$\theta_j := \theta_j + \alpha \left(y^{(i)} - h_{\theta}(\mathbf{x}_i) \right) x_j^{(i)}, \quad (8.2)$$

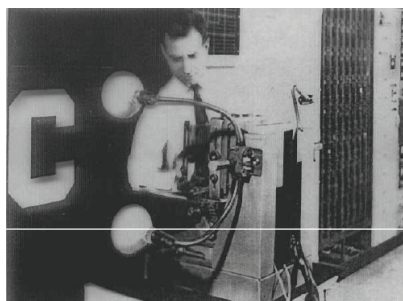
which is also related to the Widrow-Hoff learning rule, the Adaline rule, and the delta rule. These learning rules are nearly identical, but are sometimes used in slightly different contexts. It is often called the delta rule because the difference between the desired and actual output (difference between actual (training) data and hypothesis) to guide the learning. When multiplying out the difference with the inputs we have end up with the product of the activity between the inputs and output values for each synaptic channel minus the same term for the desired output,

$$\Delta w_{ij} \propto (y_i x_j - \hat{y}_i x_j), \quad (8.3)$$

which reinforces the correlation between the input and the desired output and reduces the correlation between input and the actual output. Such a learning rule is also called Hebbian after the famous NovaScotian Donald Hebb. Of course, we now view this learning rule as a special case of a gradient descent rule for a linear hypothesis function. Although this rule is not ideal for a perceptron with non-linear functions, it turned out that it still works in some cases since it corresponds to taking a step towards minimizing MSE, albeit with a wrong gradient.



Frank Rosenblatt



Charles Wightman

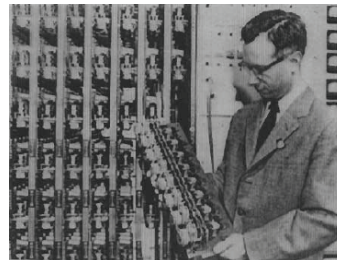


Fig. 8.2 Neural Network computers in the late 1950s.

There was a lot of excitement during the 1960s in the AI and psychology community about such learning system that resemble some brain functions. However, Minsky and Peppert showed in 1968 that such perceptrons can not represent all possible boolean functions (sometimes called the XOR problem). While it was known at this time that this problem could be overcome by a layered structure of such perceptrons (called **multilayer perceptrons**), a learning algorithms was not widely known at this time. This nearly abolished the field of learning machines, and the AI community concentrated on rule-based systems in the following years.

8.2 Multilayer perceptron (MLP)

The generalization of a delta rule, known as error-backpropagation, was finally introduced and popularized by Rumelhart, Hinton and Williams in 1986, although many others including Paul Werbos and Sunichi Amari have used it before. This popularization resulted in the explosion of the field of **Artificial Neural Networks**.

A multilayer perceptron with a layer of n^{in} input nodes, a layer of n^{h} hidden nodes, and a layer of n^{out} output nodes, is shown in Figure 8.3. The input layer is merely just relaying the inputs, while the hidden and output layer do active calculations. Such a network is thus called a 2-layer network. The term hidden nodes comes from the fact that these nodes do not have connections to the external world such as the input and

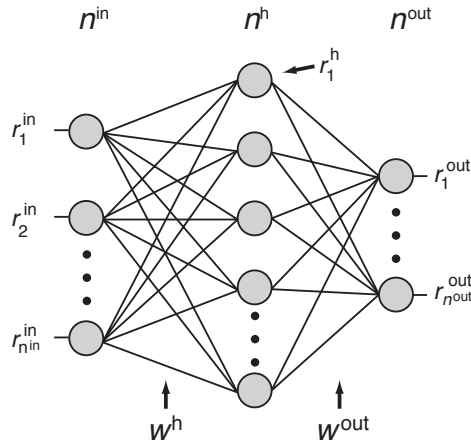


Fig. 8.3 The standard architecture of a feedforward multilayer network with one hidden layer, in which input values are distributed to all hidden nodes with weighting factors summarized in the weight matrix \mathbf{w}^h . The output values of the nodes of the hidden layer are passed to the output layer, again scaled by the values of the connection strength as specified by the elements in the weight matrix \mathbf{w}^{out} . The parameters shown at the top, n^{in} , n^h , and n^{out} , specify the number of nodes in each layer, respectively.

output nodes. Instead of the step function used in the McCulloch-Pitts model above, most such networks use a sigmoidal non-linearity,

$$g(x) = \tanh(\beta x) = 2 \frac{1}{1 + e^{-\beta x}} - 1, \quad (8.4)$$

to allow for continuous values of the nodes. The network is thus a graphical representation of a nonlinear function of the form

$$\hat{y} = g^{\text{out}}(\mathbf{w}^{\text{out}} g^h(\mathbf{w}^h \mathbf{x})). \quad (8.5)$$

It is easy to include more hidden layers in this formula. For example, the operation rule for a four-layer network with three hidden layers and one output layer can be written as

$$\hat{y} = g^{\text{out}}(\mathbf{w}^{\text{out}} g^{\text{h3}}(\mathbf{w}^{\text{h3}} g^{\text{h2}}(\mathbf{w}^{\text{h2}} g^{\text{h1}}(\mathbf{w}^{\text{h1}} \mathbf{x}))). \quad (8.6)$$

Let us discuss a special case of a multilayer mapping network where all the nodes in all hidden layers have linear activation functions ($g(x) = x$). Eqn 8.6 then simplifies to

$$\begin{aligned} \mathbf{r}^{\text{out}} &= g^{\text{out}}(\mathbf{w}^{\text{out}} \mathbf{w}^{\text{h3}} \mathbf{w}^{\text{h2}} \mathbf{w}^{\text{h1}} \mathbf{r}^{\text{in}}) \\ &= g^{\text{out}}(\tilde{\mathbf{w}} \mathbf{r}^{\text{in}}). \end{aligned} \quad (8.7)$$

In the last step we have used the fact that the multiplication of a series of matrices simply yields another matrix, which we labelled $\tilde{\mathbf{w}}$. Eqn 8.7 represents a single-layer network as discussed before. It is therefore essential to include non-linear activation functions, at least in the hidden layers, to make possible the advantages of hidden layers that we are about to discuss. We could also include connections between different hidden

layers, not just between consecutive layers as shown in Fig. 8.3, but the basic layered structure is sufficient for the following discussions.

Which functions can be approximated by multilayer perceptrons? The answer is, in principle, any. A multilayer feedforward network is a **universal function approximator**. This means that, given enough hidden nodes, any mapping functions can be approximated with arbitrary precision by these networks. The remaining problems are to know how many hidden nodes we need, and to find the right weight values. Also, the general approximator characteristics does not tell us if it is better to use more hidden layers or just to increase the number of nodes in one hidden layer. These are important concerns for practical engineering applications of those networks. These questions are related to the bias-variance tradeoff in non-linear regression as already discussed in section ??.

To train the these networks we consider again minimizing MSE which would be appropriate for Gaussian noisy data around the mean described by the model. The learning rule is then given by a gradient descent on this error function. Specifically, the gradient of the MSE error function with respect to the output weights is given by

$$\begin{aligned}
 \frac{\partial E}{\partial w_{ij}^{\text{out}}} &= \frac{1}{2} \frac{\partial}{\partial w_{ij}^{\text{out}}} \sum_i (r_i^{\text{out}} - y_i)^2 \\
 &= \frac{1}{2} \frac{\partial}{\partial w_{ij}^{\text{out}}} \sum_i (g^{\text{out}}(\sum_j w_{ij}^{\text{out}} r_j^{\text{h}}) - y_i)^2 \\
 &= g^{\text{out}'}(h_i^{\text{h}}) (\sum_j w_{ij}^{\text{out}} r_j^{\text{h}} - y_i) r_j^{\text{h}} \\
 &= \delta_i^{\text{out}} r_j^{\text{h}},
 \end{aligned} \tag{8.8}$$

where we have defined the delta factor

$$\begin{aligned}
 \delta_i^{\text{out}} &= g^{\text{out}'}(h_i^{\text{h}}) (\sum_j w_{ij}^{\text{out}} r_j^{\text{h}} - y_i) \\
 &= g^{\text{out}'}(h_i^{\text{h}}) (r_i^{\text{out}} - y_i).
 \end{aligned} \tag{8.9}$$

Eqn 8.8 is just the delta rule as before because we have only considered the output layer. The calculation of the gradients with respect to the weights to the hidden layer again requires the chain rule as they are more embedded in the error function. Thus we have to calculate the derivative

$$\begin{aligned}
 \frac{\partial E}{\partial w_{ij}^{\text{h}}} &= \frac{1}{2} \frac{\partial}{\partial w_{ij}^{\text{h}}} \sum_i (r_i^{\text{out}} - y_i)^2 \\
 &= \frac{1}{2} \frac{\partial}{\partial w_{ij}^{\text{h}}} \sum_i (g^{\text{out}}(\sum_j w_{ij}^{\text{out}} g^{\text{h}}(\sum_k w_{jk}^{\text{h}} r_k^{\text{in}})) - y_i)^2.
 \end{aligned} \tag{8.10}$$

After some battle with indices (which can easily be avoided with analytical calculation programs such as MAPLE or MATHEMATICA), we can write the derivative in a form similar to that of the derivative of the output layer, namely

$$\frac{\partial E}{\partial w_{ij}^{\text{h}}} = \delta_i^{\text{h}} r_j^{\text{in}}, \tag{8.11}$$

Table 8.1 Summary of error-back-propagation algorithm

Initialize weights arbitrarily
Repeat until error is sufficiently small
Apply a sample pattern to the input nodes: $r_i^0 := r_i^{\text{in}} = \xi_i^{\text{in}}$
Propagate input through the network by calculating the rates of nodes in successive layers l : $r_i^l = g(h_i^l) = g(\sum_j w_{ij}^l r_j^{l-1})$
Compute the delta term for the output layer:
$\delta_i^{\text{out}} = g'(h_i^{\text{out}})(\xi_i^{\text{out}} - r_i^{\text{out}})$
Back-propagate delta terms through the network:
$\delta_i^{l-1} = g'(h_i^{l-1}) \sum_j w_{ji}^l \delta_j^l$
Update weight matrix by adding the term: $\Delta w_{ij}^l = \epsilon \delta_i^l r_j^{l-1}$

when we define the delta term of the hidden term as

$$\delta_i^{\text{h}} = g^{\text{h}'}(h_i^{\text{in}}) \sum_k w_{ik}^{\text{out}} \delta_k^{\text{out}}. \quad (8.12)$$

The error term δ_i^{h} is calculated from the error term of the output layer with a formula that looks similar to the general update formula of the network, except that a signal is propagating from the output layer to the previous layer. This is the reason that the algorithm is called the **error-back-propagation algorithm**. The algorithm is summarized in Table 8.1.

Such multilayer perceptrons were able to learn nonlinear relations in data and had some success in application. However, the general problem of overfitting and the question of optimality, and well as the applicability to large data problems with more complex functions, has hampered the field since the early successes. This area of artificial neural networks become now absorbed into the more general area of machine learning with new methods that have clarified field and have led to more applicable methods that we will discuss below. We therefore follow in the next sections a more contemporary path.

Before leaving this area it is useful to point out some more general observations. Artificial neural networks have certainly been one of the first successful methods for nonlinear regression, implementing nonlinear hypothesis of the form $h(\mathbf{x}; \theta) = g(\theta^T x)$. The corresponding mean square loss function,

$$L \propto (y - g(\theta^T x))^2 \quad (8.13)$$

is then also a general nonlinear function of the parameters. Minimizing such a function is generally difficult. However, we could consider instead hypothesis that are linear in the parameters, $h(\mathbf{x}; \theta) = \theta^T \phi(\mathbf{x})$, so that the MSE loss function is quadratic in the parameters,

$$L \propto (y - \theta^T \phi(\mathbf{x}))^2. \quad (8.14)$$

The corresponding quadratic optimization problem can be solved much more efficiently. This line of ideas are further developed in support vector machines discussed next. An interesting and central further issue is how to choose the non-linear function ϕ . This will be an important ingredient for nonlinear support vector machines and unsupervised learning discussed below.

8.3 Support Vector Machines (SVM)

8.3.1 Linear classifiers with large margins

In this section we briefly outline the basic idea behind Support Vector Machines (SVM) that are currently thought to be the best general purpose supervised classifier algorithm. SVMs, and the underlying statistical learning theory, has been worked out by Vladimir Vapnik since the early 1960, but some breakthroughs were also made in the late 1990 with some collaborators like Corinna Cortes, Chris Burges, Alex Smola, and Bernhard Schölkopf to name but a few, and SVM have since become very popular and hard to beat. While we outline some of the underlying formulas, we do not derive all the steps but will try to give some intuition. A more thorough treatment can be found in the references on the web page. The here we just want to provide the big picture, but need to show some formulas to highlight some of the discussion.

The basic SVMs are concerned with binary classification. Figure 8.4 shows an example of two classes, depicted by different symbols, in a two dimensional attribute space. We distinguish here attributes from features as follows. Attributes are the raw measurements, where as features can be made up by combining attributes. For example, the attributes x_1 and x_2 could be combines in a feature vector $(x_1, x_1x_2, x_2, x_1^2, x_2^2)^T$. This will become important later, but it is important to introduce the notation here. Our training set consists again of m data with attribute values $\mathbf{x}^{(i)}$ and labels $y^{(i)}$. We chose here the labels of the two classes as $y \in \{-1, 1\}$, as this will nicely simplify some equations.

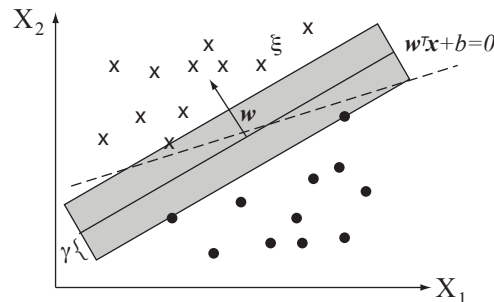


Fig. 8.4 Illustration of linear support vector classification.

The two classes in the figure 8.4 can be separated by a line, which can be parameterized by

$$w_1x_1 + w_2x_2 - b = \mathbf{w}^T \mathbf{x} - b = 0. \quad (8.15)$$

While the first equation shows the lines equation with its components in two dimensions, the next expression is the same in any dimension. Of course, in three dimension we would talk about a plane. In general, we will talk about a **hyperplane** in any dimensions. The particular hyperplane is the dividing or separating hyperplane between the two classes. We also introduce what the **margin** γ , which is the perpendicular distance between the dividing hyperplane and the closest point.

The main point to realize now is that the dividing hyperplane that maximizes the margin, the so called **maximum margin classifier**, is the best solution we can find. Why is that? We should assume that the training data, shown in the figure, are some unbiased examples of the true underlying density function describing the distribution of points within each class. It is then likely that new data points, which we want to classify, are close to the already existing data points. Thus, if we make the separating hyperplane as far as possible from each point, than it is most likely to not make wrong classification. Or, with other words, a separating hyperplane like the one shown as dashed line in the figure, is likely to generalize much worse than the maximum margin hyperplane. So the maximum margin hyperplane is the best generalizer for binary classification for the training data.

What is if we can not divide the data with a hyperplane and we have to consider non-linear separators. Don't we then run into the same problems as outlined before, specifically the bias-variance tradeoff? Yes, indeed, this will still be the challenge, and our aim is really to work on this problem. But before going there it is useful to formalize the linear separable case in some detail as the representation of the optimization problem will be a key in applying some tricks later.

Learning a linear maximum margin classifier on labeled data means finding the parameters (w) and b that maximizes the margin. For this we could compute the distances of each point from the hyperplane, which is simply a geometric exercise,

$$\gamma^{(i)} = y^{(i)} \left(\left(\frac{\mathbf{w}}{\|\mathbf{w}\|} \right)^T \mathbf{x}^{(i)} + \frac{b}{\|\mathbf{w}\|} \right). \quad (8.16)$$

The vector $\mathbf{w}/\|\mathbf{w}\|$ is the normal vector of the hyperplane, a vector of unit length perpendicular to the hyperplane. We overall margin we want to maximize is the distance to the closest point,

$$\gamma = \min_i \gamma^{(i)}. \quad (8.17)$$

By looking at equation 8.21 we see that maximizing γ is equivalent to minimizing $\|\mathbf{w}\|$, or, equivalently, of minimizing

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2. \quad (8.18)$$

More precisely, we want to maximize this margin under the constraint that no training data lies within the margin,

$$\mathbf{w}^T \mathbf{x} + b \geq 1 \quad \text{for } y^{(i)} = 1 \quad (8.19)$$

$$\mathbf{w}^T \mathbf{x} + b \leq -1 \quad \text{for } y^{(i)} = -1, \quad (8.20)$$

which can nicely be combines with our choice of class representation as

$$y^{(i)} (\mathbf{w}^T \mathbf{x} + b) \leq 1. \quad (8.21)$$

Thus we have a quadratic minimization problem with linear inequalities as constraint. Taking a constrain into account can be done with a **Lagrange formalism**. For this we

simply add the constraints to the main objective function with parameters α_i called Lagrange multipliers,

$$\mathcal{L}^P(\mathbf{w}, b, \alpha_i) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^m \alpha_i [y^{(i)}(\mathbf{w}^T \mathbf{x} + b) - 1]. \quad (8.22)$$

The Lagrange multipliers determine how well the constraints are observed. In the case of $\alpha_i = 0$, the constraints do not matter. In order to conserve the constraints, we should thus make these values as big as we can. Finding the maximum margin classifier is given by

$$p^* = \min_{\mathbf{w}, b} \max_{\alpha_i} \mathcal{L}^P(\mathbf{w}, b, \alpha_i) \leq p^* = \max_{\alpha_i} \min_{\mathbf{w}, b} \mathcal{L}^D(\mathbf{w}, b, \alpha_i) = d^*. \quad (8.23)$$

In this formula we also added the formula when interchanging the min and max operations, and the reason for this is the following. It is straight forward to solve the optimization problem on the left hand side, but we can also solve the related problem on the right hand side which turns out to be essential when generalizing the method to nonlinear cases below. Moreover, the equality holds when the optimization function and the constraints are convex⁸. So, if we minimize \mathcal{L} by looking for solutions of the derivatives $\frac{\partial \mathcal{L}}{\partial \mathbf{w}}$ and $\frac{\partial \mathcal{L}}{\partial b}$, we get

$$\mathbf{w} = \sum_{i=1}^m \alpha_i y^{(i)} \mathbf{x}^{(i)} \quad (8.24)$$

$$0 = \sum_{i=1}^m \alpha_i y^{(i)} \quad (8.25)$$

Substituting this into the optimization problem we get

$$\max_{\alpha_i} \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} y^{(i)} y^{(j)} \alpha_i \alpha_j \mathbf{x}^{(i)T} \mathbf{x}^{(j)}, \quad (8.26)$$

subject to the constraints

$$\alpha_i \geq 0 \quad (8.27)$$

$$\sum_{i=1}^m \alpha_i y^{(i)} = 0. \quad (8.28)$$

From this optimization problem it turns out that the α_i 's of only a few examples, those ones that are lying on the margin, are the only ones with have $\alpha_i \neq 0$. The corresponding training examples are called **support vectors**. The actual optimization can be done with several algorithms. In particular, John Platt developed the sequential minimal optimization (SMO) algorithm that is very efficient for this optimization problem. Please note that the optimization problem is convex and can thus be solved very efficiently without the danger of getting stuck in local minima.

⁸Under these assumptions there are other conditions that hold, called the **Karush-Kuhn-Tucker** conditions, that are useful in providing proof in the convergence of these methods outlined here.

Once we found the support vectors with corresponding α_i 's, we can calculate (w) from equation 8.29 and b from a similar equation. Then, if we are given a new input vector to be classified, this can then be calculated with the hyperplane equation 8.20 as

$$y = \begin{cases} 1 & \text{if } \sum_{i=1}^m \alpha_i y^{(i)} \mathbf{x}^{(i)T} \mathbf{x} > 0 \\ -1 & \text{otherwise} \end{cases} \quad (8.29)$$

Since this is only a sum over the support vectors, classification becomes very efficient after training.

8.3.2 Soft margin classifier

So far we only discussed the linear separable case. But how about the case when there are overlapping classes? It is possible to extend the optimization problem by allowing some data points to be in the margin while penalizing these points somewhat. We include therefore some **slag variables** ξ_i that reduce the effective margin for each data point, but we add to the optimization a penalty term that penalizes if the sum of these slag variables are large,

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i \xi_i, \quad (8.30)$$

subject to the constrains

$$y^{(i)} (\mathbf{w}^T \mathbf{x} + b) \geq 1 - \xi_i \quad (8.31)$$

$$\xi_i \geq 0 \quad (8.32)$$

The constant C is a free parameter in this algorithm. Making this constant large means allowing less points to be in the margin. This parameter must be tuned and it is advisable to at least try to vary this parameter to verify that the results do not dramatically depend on an initial choice.

8.3.3 Nonlinear Support Vector Machines

We have treated the case of overlapping classes while assuming that the best we can do is still a linear separation. But what if the underlying problem is separable, f only with a more complex function. We will now look into the non-linear generalization of the SVM.

When discussing regression we started with the linear case and then discussed non-linear extensions such as regressing with polynomial functions. For example, a linear function in two dimensions (two attribute values) is given by

$$y = w_0 + w_1 x_1 + w_2 x_2, \quad (8.33)$$

and an example of a non-linear function, that of an polynomial of 3rd order, is given by

$$y = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_1 x_2 + w_4 x_1^2 + w_5 x_2^2. \quad (8.34)$$

The first case is a linear regression of a feature vector

$$\mathbf{x} = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}. \quad (8.35)$$

We can also view the second equation as that of linear regression on a feature vector

$$\mathbf{x} \rightarrow \phi(\mathbf{x}) = \begin{pmatrix} x_1 \\ x_2 \\ x_1 x_2 \\ x_1^2 \\ x_2^2 \end{pmatrix}, \quad (8.36)$$

which can be seen as a mapping $\phi(\mathbf{x})$ of the original attribute vector. We call this mapping a **feature map**. Thus, we can use the above maximum margin classification method in non-linear cases if we replace all occurrences of the attribute vector \mathbf{x} with the mapped feature vector $\phi(\mathbf{x})$. There are only two problems remaining. One is that we have again the problem of overfitting as we might use too many feature dimensions and corresponding free parameters w_i . The second is also that with an increased number of dimensions, the evaluation of the equations becomes more computational intensive. However, there is a great trick to alleviate the later problem in the case when the methods only rely on dot products, like in the case of the formulation in the dual problem. In this, the function to be minimized, equation 8.31 with the feature maps, only depends on the dot products $\phi(\mathbf{x}^{(i)})^T \phi(\mathbf{x}^{(j)})$. Also, when predicting the class for a new input vector \mathbf{x} from equation 8.29 when adding the feature maps, we only need the resulting values for the dot products $\phi(\mathbf{x}^{(i)})^T \phi(\mathbf{x})$ which can sometimes be represented as function called **Kernel function**,

$$K(\mathbf{x}, \mathbf{z}) = \phi(\mathbf{x})^T \phi(\mathbf{z}). \quad (8.37)$$

Instead of actually specifying a feature map, which is often a guess to start, we could actually specify a Kernel function. For example, let us consider a quadratic feature map

$$K(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z} + c)^2. \quad (8.38)$$

We can then try to write this in the form of equation 8.42 to find the corresponding feature map. That is

$$K(\mathbf{x}, \mathbf{z}) = (\mathbf{x}^T \mathbf{z})^2 + 2c\mathbf{x}^T \mathbf{z} + c^2 \quad (8.39)$$

$$= \left(\sum_i x_i z_i \right)^2 + 2c \sum_i x_i z_i + c^2 \quad (8.40)$$

$$= \sum_j \sum_i (x_i x_j) (z_i z_j) + \sum_i (\sqrt{(2c)} x_i) (\sqrt{(2c)} z_i) + cc \quad (8.41)$$

$$= \phi(\mathbf{x})^T \phi(\mathbf{z}), \quad (8.42)$$

with

$$\phi(\mathbf{x}) = \begin{pmatrix} x_1x_1 \\ x_1x_2 \\ \dots \\ x_nx_1 \\ x_nx_2 \\ \dots \\ \sqrt{2c}x_1 \\ \sqrt{2c}x_2 \\ \dots \\ c \end{pmatrix}, \quad (8.43)$$

The dimension of this feature vector is $O(n^2)$ for n original attributes. Thus, evaluating the dot product in the mapped feature space is much more time consuming than calculating the Kernel function which is just the square of the dot product of the original attribute vector. The dimensionality of Kernels with higher polynomials is quickly rising, making the benefit of the Kernel method even more impressive.

While we have derived the corresponding feature map for a specific Kernel function, this task is not always easy and not all functions are valid Kernel functions. We have also to be careful that the Kernel functions still lead to convex optimization problems. In practice, only a small number of Kernel functions is used. Besides the polynomial Kernel mentioned before, one of the most popular is the Gaussian Kernel,

$$K(\mathbf{x}, \mathbf{z}) = \exp - \frac{\|\mathbf{x} - \mathbf{z}\|^2}{2\gamma^2}, \quad (8.44)$$

which corresponds to an infinitely large feature map.

As mentioned above, a large feature space corresponds to a complex model that is likely to be prone to overfitting. We must therefore finally look into this problem. The key insight here is that we are already minimizing the sum of the components of the parameters, or more precisely the square of the norm $\|\mathbf{w}\|^2$. This term can be viewed as **regularization** which favours a smooth decision hyperplane. Moreover, we have discussed two extremes in classifying complicated data, one was to use Kernel functions to create high-dimensional non-linear mappings and hence have a high-dimensional separating hyperplane, the other method was to consider a low-dimensional separating hyperplane and interpret the data as overlapping. The last method includes a parameter C that can be used to tune the number of data points that we allow to be within the margin. Thus, we can combine these two approaches to classify non-linear data with overlaps where the soft margins will in addition allow us to favour more smooth dividing hyperplanes.

8.3.4 Regularization and parameter tuning

In practice we have to consider several free parameters when applying support vector machines. First, we have to decide which Kernel function to use. Most packages have a number of choices implemented. We will use for the following discussion the Gaussian Kernel function with width parameter γ . Setting a small value for γ and allowing for a large number of support vectors (small C), corresponds to a complex model. In contrast, larger width values and regularization constant C will increase the stiffness of the model and lower the complexity. In practice we have to tune these

parameters to get good results. To do this we need to use some form of validation set, as discussed in section 5.6, and k-times cross validation is often implemented in the software packages. An example of the SVM performance (accuracy) on some examples (Iris Data set from the UCI repository; From Broadman and Trappenberg, 2006) is shown in figure 8.5 for several values of γ and C . It is often typical that there is a large area where the SVM works well and has only little variations in terms of performance. This robustness has helped to make SVMs practical methods that often outperform other methods. However, there is often also an abrupt onset of the region where the SVM fails, and some parameter tuning is hence required. While just trying a few settings might be sufficient, some more systematic methods such as grid search or simulated annealing also work well.

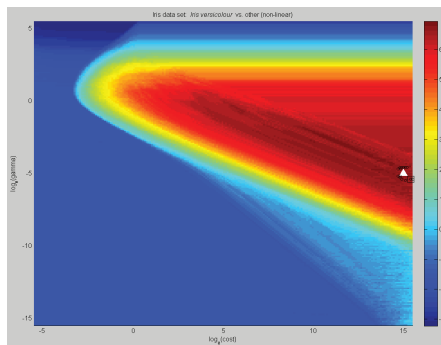


Fig. 8.5 Illustration of SVM accuracy for different values of parameters C and γ .

8.3.5 Statistical learning theory and VC dimension

SVMs are good and practical classification algorithms for several reasons, including the advantage of being convex optimization problem that than can be solved with quadratic programming, have the advantage of being able to utilize the Kernel trick, have a compact representation of the decision hyperplane with support vectors, and turn out to be fairly robust with respect to the hyper parameters. However, in order to be good learners, they need to moderate the variance-bias tradeoff discussed in section 5.6. A great theoretical contributions of Vapnik and colleagues was the embedding of supervised learning into statistical learning theory and to derive some bounds that make statements on the average ability to learn from data. We outline here briefly the ideas and state some of the results. We discuss this issue here in the context of binary classification, although similar observations can be made in the case of multiclass classification and regression.

We start again by stating our objective, which is to find a hypothesis which minimized the generalization error. To state this a bit more differentiated and to use the nomenclature common in these discussions, we call the error function here the **risk function** R . In particular, the **expected risk** for a binary classification problem is the probability of misclassification,

$$R(h) = P(h(x) \neq y) \quad (8.45)$$

Of course, we generally do not know this density function, though we need to approximate this with our validation data. We assume thereby again that the samples are iid (independent and identical distributed) data, and can then estimate what is called the **empirical risk**,

$$\hat{R}(h) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}(h(\mathbf{x}^{(i)}; \theta) = y^{(i)}). \quad (8.46)$$

We use here again m as the number of examples, but note that this is here the number of examples in the validations set, which is the number of all training data minus the ones used for training. Also, we will discuss this empirical risk further, but note that it is better to use the regularized version that incorporates a smoothness constrain such as

$$\hat{R}_{rmreg}(h) = \frac{1}{m} \sum_i \mathbb{1}(h(\mathbf{x}^{(i)}; \theta) = y^{(i)}) - \lambda \|\mathbf{w}\|^2 \quad (8.47)$$

in the case of SVM, where λ is a regularization constant. Thus, wherever $\hat{R}(h)$ is used in the following, we can replace this with $\hat{R}_{rmreg}(h)$. **Empirical risk minimization** is the process of finding the hypothesis \hat{h} that minimizes the empirical risk,

$$\hat{h} = \arg \min_h \hat{R}(h). \quad (8.48)$$

The empirical risk is the MLE of the mean of a Bernoulli-distributed random variable with true mean $R(h)$. Thus, the empirical risk is itself a random variable for each possible hypothesis h . Let us first assume that we have k possible hypothesis h_i . We now draw on a theorem by Hoeffding called the **Hoeffding inequality** that provides and upper bound for the sum of random numbers to its mean,

$$P(|R(h_i) - \hat{R}(h_i)| > \gamma) \leq 2 \exp(-2\gamma^2 m). \quad (8.49)$$

This formula states that there is a certain probability that we make an error larger than γ for each hypothesis of the empirical risk compared to the expected risk, although the good news is that this probability is bounded and that the bound itself becomes exponentially smaller with the number of validation examples. This is already an interesting results, but we now want to know the probability that some, out of all possible hypothesis, are less than γ . Using the fact that the probability of the union of several events is always less or equal to the sum of the probabilities, one can show that with probability $1 - \delta$ the error of a hypothesis is bounded by

$$|R(h) - \hat{R}(h)| \leq \sqrt{\frac{1}{2m} \log \frac{2k}{\delta}}. \quad (8.50)$$

This is a great results since it shows how the error of using an estimate the risk, the empirical risk that we can evaluate from the validation data, is getting smaller with training examples and with the number of possible hypothesis.

While the error scales only with the log of the number of possible hypothesis, the values goes still to infinite when the number of possible hypothesis goes to infinite, which much more resembles the situation when we have parameterized hypothesis. However, Vapnik was able to show the following generalization in the infinite case,

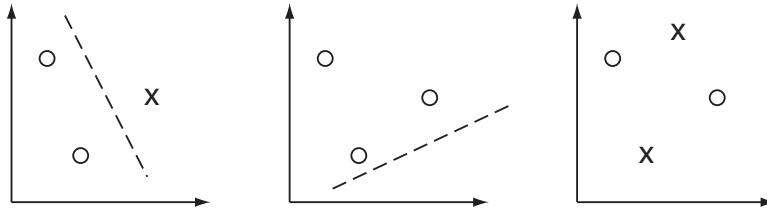


Fig. 8.6 Illustration of VC dimensions for the class of linear functions in two dimensions.

which is that given a hypothesis space with **Vapnik-Chervonencis** dimension $VC(\{h\})$, then, with probability $1 - \delta$, the error of the empirical risk compared to the expected risk (true generalization error) is

$$|R(h) - \hat{R}(h)| \leq O \left(\sqrt{\frac{VC}{m} \log \frac{m}{VC} + \frac{1}{m} \log \frac{1}{\delta}} \right). \quad (8.51)$$

The VC dimension is thereby a measure of how many points can be divided by a member of the hypothesis set for all possible label combinations of the point. For example, consider three arbitrary points in two dimensions as shown in figure 8.6, and let us consider the hypothesis class of all possible lines in two dimensions. I can always divide the three points under any class membership condition, of which two examples are also shown in the figure. In contrast, it is possible to easily find examples with four points that can not be divided by a line in two dimensions. The VC dimension of lines in two dimensions is hence $VC = 3$.⁹

8.4 SV-Regression and implementation

8.4.1 Support Vector Regression

While we have mainly discussed classification in the last few sections, it is time to consider the more general case of regression and to connect these methods to the general principle of maximum likelihood estimation outlined in the previous chapter. It is again easy to illustrate the method for the linear case before generalizing it to the non-linear case similar to the strategy followed for SVMs.

We have already mentioned in section 5.2 the ϵ -insensitive error function which does not count deviations of data from the hypothesis that are less than ϵ from the hypothesis. This is illustrated in figure 8.7. The corresponding optimization problem is

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_i (\xi_i + \xi_i^*), \quad (8.52)$$

subject to the constraints

$$y^{(i)} - \mathbf{w}^T \mathbf{x} - b \leq \xi_i \quad (8.53)$$

⁹Three points of different classes can not be separated by a single line, but these are singular points that are not effective in the definition of VC dimension.

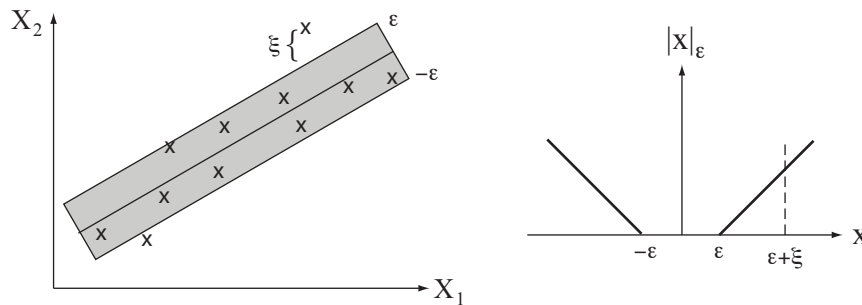


Fig. 8.7 Illustration of support vector regression and the ϵ -insensitive cost function.

$$y^{(i)} - \mathbf{w}^T \mathbf{x} - b \geq \xi_i^* \quad (8.54)$$

$$\xi_i, \xi_i^* \geq 0 \quad (8.55)$$

The dual formulations does again only depend on scalar products between the training examples, and the regression line can be also be expressed by a scalar product between the support vectors and the prediction vector,

$$h(\mathbf{x}; \alpha_i, \alpha_i^*) = \sum_{i=1}^m (\alpha_i - \alpha_i^*) \mathbf{x}_i^T \mathbf{x}. \quad (8.56)$$

This, we can again use Kernels to generalize the method to non-linear cases.

8.4.2 Implementation

There are several SVM implementations available, and SVMs are finally becoming a standard component of data mining tools. We will use the implementation called LIBSVM which was written by Chih-Chung Chang and Chih-Jen Lin and has interfaces to many computer languages including Matlab. There are basically two functions that you need to use, namely `model=svmtrain(y,x,options)` and `svmpredict(y,x,model,options)`. The vectors \mathbf{x} and \mathbf{y} are the training data or the data to be tested. `svmtrain` uses k -fold cross validation to train and evaluate the SVM and returns the trained machine in the structure `model`. The function `svmpredict` uses this model to evaluate the new data points. Below is a list of options that shows the implemented SVM variants. We have mainly discussed C-SVC for the basic soft support vector classification, and `epsilonSVR` for support vector regression.

```
-s svm_type : set type of SVM (default 0)
0 -- C-SVC
1 -- nu-SVC
2 -- one-class SVM
3 -- epsilon-SVR
4 -- nu-SVR
-t kernel_type : set type of kernel function (default 2)
0 -- linear: u'*v
1 -- polynomial: (gamma*u'*v + coef0)^degree
```

```

2 -- radial basis function: exp(-gamma*|u-v|^2)
3 -- sigmoid: tanh(gamma*u'*v + coef0)
-d degree : set degree in kernel function (default 3)
-g gamma : set gamma in kernel function (default 1/num_features)
-r coef0 : set coef0 in kernel function (default 0)
-c cost : set the parameter C of C-SVC, epsilon-SVR, and nu-SVR (default 1)
-n nu : set the parameter nu of nu-SVC, one-class SVM, and nu-SVR (default 0.5)
-p epsilon : set the epsilon in loss function of epsilon-SVR (default 0.1)
-m cachesize : set cache memory size in MB (default 100)
-e epsilon : set tolerance of termination criterion (default 0.001)
-h shrinking: whether to use the shrinking heuristics, 0 or 1 (default 1)
-b probability_estimates: whether to train a SVC or SVR model for probability estimates, 0 or 1 (default 0)
-wi weight: set the parameter C of class i to weight*C, for C-SVC (default 1)

```

The k in the `-g` option means the number of attributes in the input data.

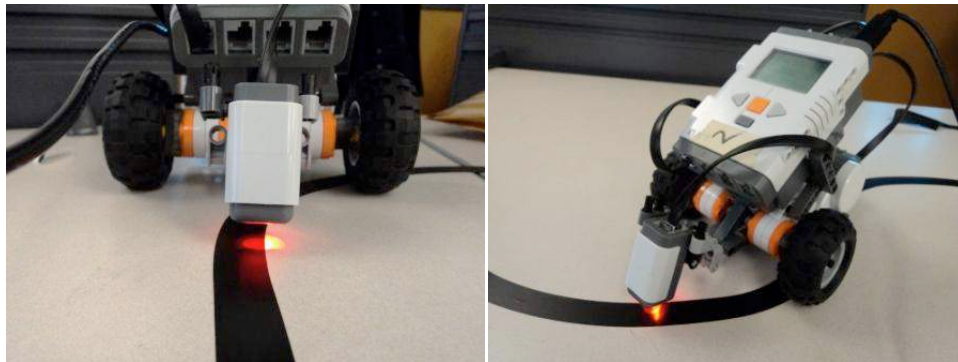
8.5 Supervised Line-following

8.5.1 Objective

The objective of this experiment is to investigate Supervised Learning through teaching a robot how to follow a line using a Support Vector Machine (SVM).

8.5.2 Setup

- Mount light sensor on front of NXT, plugged into Port 1
- Use a piece of dark tape (i.e. electrical tape) to mark a track on a flat surface. Make sure the tape and the surface are coloured differently enough that the light sensor returns reasonably different values between the two surfaces.
- This program requires a MATLAB extension that can use Support Vector Machines. Download:



8.5.3 Program

Data collection requires the user to manually move the wheels of the NXT. When the training begins, start the NXT so the light sensor's beam is either on the tape or the surface. Zig zag the NXT so the beam travels on and off the tape by moving either the right or the left wheel, one at a time. Record the positions of the left and right wheels, as well as the light sensor's reading during frequent intervals. It is important to make sure the wheel not in motion stays as stationary as possible to obtain the optimal training set of data.

After data collection, find the difference between the right and the left wheel positions for each time sample taken, and use the SVM to create a model between these differences and the light sensor readings. For instance:

```
model = svmtrain(delta,lightReading,'-s 8 -g 0 -b 1 -e 0.1 -q');
```

To implement the model, place the NXT on the line and use SVMpredict to input a light sensor reading and drive the robot left or right depending on the returned value of the SVMpredict.

```
lightVal=GetLight(SENSOR_1);
if svmpredict(0,lightVal,model,'0')>0
    left.Stop();
    right.SendToNXT();
else
    right.Stop();
    left.SendToNXT();
end
```