

# CSCI 4155/6505 (2016): Machine Learning

**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. 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. Finally, the book by Kevin Murphy, *Machine Learning: A Probabilistic Perspective* is highly recommended for more in-depth studies.

Several people have contributed considerably to this lecture notes. In particular I would thank my PhD students Paul Hollensen and Patrick Connor.



# Contents

---

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	The basic idea behind supervised ML	1
1.2	Training, validating and testing	2
1.3	Unsupervised and reinforcement learning	4
1.4	Causal learning, overfitting and regression	6
<b>2</b>	<b>ML programming with Python</b>	<b>10</b>
2.1	General scientific programming in Python	10
2.2	Cross validation example from Intro	13
2.3	Classification with support vector machine using scikit-learn	15
2.4	Other classification methods including MLP with Tensorflow	16
<b>3</b>	<b>Basic probability theory</b>	<b>19</b>
3.1	Random numbers and their probability (density) function	19
3.2	Moments: mean, variance, etc.	21
3.3	Examples of probability (density) functions	23
3.4	Bernoulli distribution	23
3.5	Cumulative probability (density) function and the Gaussian error function	25
3.6	Functions of random variables and the central limit theorem	27
3.7	Measuring the difference between distributions	28
3.8	Density functions of multiple random variables	29
<b>4</b>	<b>Regression and maximum likelihood</b>	<b>32</b>
4.1	Trends in stochastic data	32
4.2	Probabilistic models and maximum likelihood	35
4.3	Maximum a posteriori estimates	37
4.4	Multivariate causal modelling	38
4.5	Discrete probabilistic modeling in Python using LEA	41
<b>5</b>	<b>MLP</b>	<b>43</b>
5.1	The historical threshold perceptron	43

5.2	The sigmoid perceptron	45
5.3	Multilayer perceptron (MLP)	49
5.4	Stochastic MLPs and Cross-Entropy Loss for sigmoidal classification	53

# 1 Introduction

---

## 1.1 The basic idea behind supervised ML

In the sense of its words, Machine Learning (ML) is the area that tries to build intelligent machines by defining a machine with specific operational abilities and training it with examples to perform specific tasks. This might sound like a niche area in science and you might wonder why there is now so much interest in this discipline, both academically and in industry. The reason is that ML is really about modeling data that provide the basis of advanced object recognition and data mining and is hence the analytical engine in areas such as data science, big data, data analytics and science in general.

ML has a long history with traces far back in time. One of the first recognized exciting realizations of the promise of learning machines came in the late 1950s and early 1960s with work like Arthur Samuel's self-learning checkers program and Frank Rosenblatt's perceptron. A second wave came in the 1980s and 90s with multilayer perceptrons and recurrent networks. And since 2006 we have a third wave of neural networks, that of deep learning which we will discuss more in this course.

ML is not restricted to neural networks. Indeed, the development of statistical machine learning and Bayesian networks has influenced the field strongly in the last 20 years and has been essential in much of its progress as well as in the deeper understanding of machine learning. This course will hence also introduce these more general ideas.

It is common and somewhat useful to distinguish three areas of machine learning, that of supervised learning, unsupervised learning, and reinforcement learning. Much of what is currently most associated with the success of ML is supervised learning, sometimes also called predictive learning. The basic task of supervised learning is that of taking as input a vector  $\mathbf{x}$  of measurements, such as some medical measurements or robotic sensor data, and predicting an output value  $y$  such as the state of a patient's health or the location of obstacles. Note that we follow here a common notation of denoting a vector, matrix or tensor with bold faced letters, whereas we use regular fonts for scalars. We usually call the input vector a feature vector as the components of this are typically a set feature values of an object. The output could be also a multi-dimensional object such as a vector or tensor itself.

Mathematically we can denote the relations between the input and the output as a function

$$y = f(\mathbf{x}). \tag{1.1}$$

We consider the function above as a description of the **true underlying world**, and our main goal is to find out this precise relation. In the above formula we considered a single output value and several input values for illustration purposes, although we

see later that we can extend this readily to multiple output values. This would then correspond to a vector function.

The challenge for machine learning to find this function or at least to approximate it sufficiently. Machine learning has several approaches to deal with this. One approach that we will predominantly follow for much of the course is to define a general parameterized function

$$\hat{y} = \hat{f}(\mathbf{x}; \theta). \quad (1.2)$$

This formula describes that we make a parameterized hypothesis in which we specified a function  $\hat{f}$  that depend on parameters to approximate the desired input-output relation. This function is called a **model**. The model is generally an approximation of a system to study specific aspects of its behaviour. This often means that not all of the underlying world has to be captured in depth. For example, a building engineer might make a model of a bridge to tests its static without including the ascetic aspects that an architect might emphasize in a model. In our context the word model is synonymous with approximation.

We have indicated that this model is an approximation of the desired relation by using a hat symbol. However, we frequently drop this symbol when the relation is clear from the context, for example when the function contains parameters. The parameters are specified in this function by including the parameter as vector  $\theta$  behind a semicolon in the function arguments. More appropriately, the formula defines a set of functions in the parameter space. A good solution represented by a point in this parameter space is an approximation that can be used for predicting output values (labels)  $\hat{y}$  for specific input values.

We will later go one step further by considering the more general case when we might not be able to predict an exact value but at least the probability that a certain value will occur. Indeed, it is quite common that the process under investigation includes stochastic (random) or unknown factors. True underlying world model is thus better described by a probability density function

$$P(Y = y|\mathbf{x}). \quad (1.3)$$

Formulating ML learning in a probabilistic context has been most useful and provides us with the formalization that created the most insight into this field. In the stochastic framework we are then modelling a density function

$$p(\hat{Y} = \hat{y}|\mathbf{x}; \theta). \quad (1.4)$$

For now we will follow the function approximation formalization, but we return to the probabilistic framework later.

## 1.2 Training, validating and testing

Coming up with the right parameterized approximation function is the hard problem in machine learning, and we will later discuss several choices. There are also methods to systematically develop the approximation function from the data, generally called non-parametric methods. However, we assume for now that we have a parameterized



approximation function. The question is then how we determine good parameters. This is where the learning process comes in.

In supervised learning we must be given some examples of input-output relation from which we learn. We can think about these examples as given by a teacher. The teacher data called the **training set** are used to directly determine the parameters of the model. We can denote this training set as

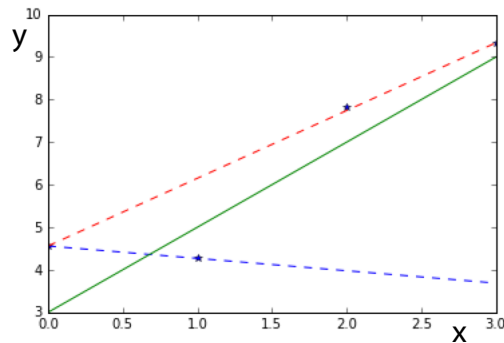
$$\{\mathbf{x}(i), y(i)\}, \quad (1.5)$$

where the superscript  $i$  labels the specific training example. These indices are enclosed in brackets to not confuse them with exponents.

There are different ways – usually called a training algorithm – to determine the parameters of a model. Let us illustrate this with an example where we assume we have a single input feature,  $x$ , and we hypothesize that the output  $y$  is linearly related to  $x$ . Mathematically we write this linear model as

$$\hat{y} = ax + b, \quad (1.6)$$

where  $a$  is the slope of the linear function and  $b$  is the  $y$ -axis intercept or bias of this function. Using the training data to determine good parameters is also called regression in this context.



**Fig. 1.1** A form of linear regression of data and cross-validation.

Let us illustrate an example regression procedure (there are different possible regression procedures) with the help of an example. For this we choose the true underlying model

$$y = 2x + 3 + \eta, \quad (1.7)$$

where  $\eta$  is a normal distributed random variable. We added this random addition to the perfect linear model to include right away a typical challenge in ML, that of having imprecise measurements and hence noisy training data. The other way to interpret the model is actually to accept the ‘world’ as stochastic and hence we are considering a stochastic model. In any case, from this model we chose some data points by sampling,

$$(0, 4.56), (1, 4.27), (2, 7.81), (3, 9.33) \quad (1.8)$$

We now make the assumption of a parameterized linear function

$$\hat{y} = ax + b, \quad (1.9)$$

with two parameters  $a$  and  $b$ , and use a simple method to determine the two parameters. As we have a linear equation with only two unknowns, we only need two data points to determine their values. Using the two first data points as training set we get

$$a = \frac{y^{(2)} - y^{(1)}}{x^{(2)} - x^{(1)}} = -0.29 \quad (1.10)$$

$$b = y^{(1)} - ax^{(1)} = 4.56 \quad (1.11)$$

This is not a very good agreement with the ideal values of  $a = 2$  and  $b = 3$ . How about using the first and last data point as training set. The estimate of the  $y$ -intercept stays the same, but the slope estimation becomes  $a = 1.59$  which is already much better.

But how would we know what the best solution is when we do not already know the solution? The answer is a procedure called **cross validation** where we use the data not used for training to validate the goodness of prediction on other data. Thus, the data not used directly to estimate the model parameters are called the **validation set**. We chose here to evaluate the goodness of the model with the mean square error (MSE) on the data not used for training

$$MSE = \frac{1}{N_{\text{val}}} \sum_{i=1}^{N_{\text{val}}} (\hat{y} - y)^2, \quad (1.12)$$

where  $N_{\text{val}}$  is the number of validation data. If we calculate the MSE for all possible data combinations we can choose as final estimation the estimation that leads to the smallest MSE of the unseen validation data points.

Cross validation is generally used to tune the **hyper-parameters**. Hyper-parameters are parameters of the learning algorithms. In our specific examples these are the choice of which data points to choose. Later we will discuss other learning algorithms that have hyper-parameter like a learning rate or the number of learning steps.

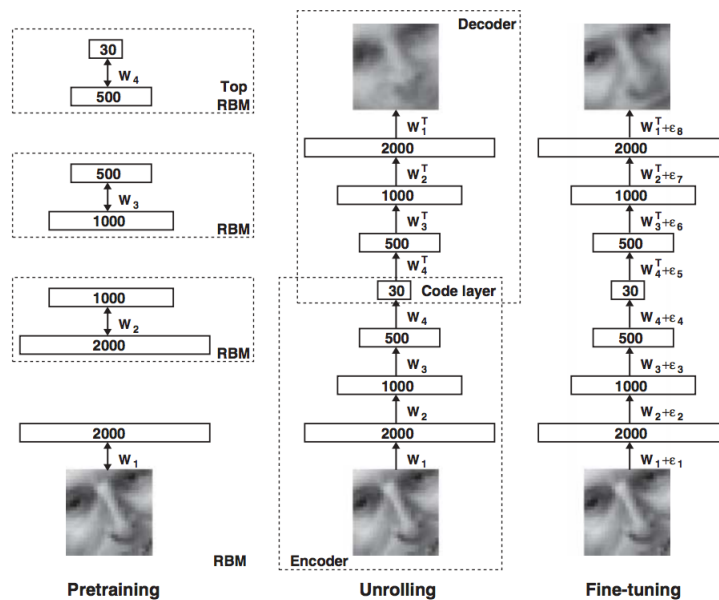
Finally, the ultimate test is using data that were not given to us during a learning phase to test the prediction of the model with the parameters that we estimated during the learning and cross-validation procedure. These data that are usually not given to the researchers during a training phase is recalled the **test set**. It is essential to be very clear about the differences between the training set, the validation set and the test set.

### 1.3 Unsupervised and reinforcement learning

Up to now we have discussed supervised learning where a teacher provides detailed examples of what the output of a machine should be. In **reinforcement learning** there is still some feedback from a teacher or the environment in the form of indicators that show how desirable the outputs of the learner is. This is often described as reward or punishment. However, the learning process requires some exploration as the learner must find the required action to attain a rewarding position. Such a learning

requirement represents a common task in robotics and human learning, although a supervised learning system could be in itself a component of such a learning agenda, for example in the required vision system.

The last basic category of **unsupervised learning** is another important type of learning. This the of learning is more directly related to supervised learning except that the training data sets has no labels, That is, we are for example given a large data set of images but without detailed descriptions what the photographs represent. So one might ask how this can be of any use. However, there is a lot of information in the pictures itself as they give examples of how natural images look like or that there are usually edges and a certain distribution of colors. In other words, we can learn a lot about the statistics of the feature values from a large collection of unlabelled examples.



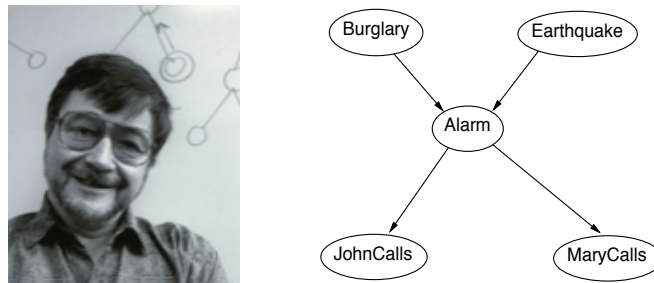
**Fig. 1.2** A form of linear regression of data and cross-validation.

Some knowledge of the statistics of data is important to guide good representations of the data. We will see later that a good representation of data is a key element of solving mapping problems such as object recognition. Indeed, we can view the success of deep learning as the result of learning good layers of representations between the raw sensory input and the desired semantic space. It is therefore **representational learning** based on unsupervised and supervised learning that had build the foundation of much progress. For example, Figure 1.2 shows a famous work of Geoff Hinton and Ruslan Salakhutdinov published in Science 2006 of a neural network with several layers that transforms an input image to itself. At first it uses some unsupervised learning techniques that learns to reconstruct each input layer with a hidden layer. Such a network is called a **restricted Boltzmann machine**. Next we stack these layers on top of each other and train it like a supervised learning problem where the desired

output (label) is the image itself. This is called an **autoencoder**. While the ability of translating an input image to itself might be questionable, the interesting part of such an autoencoder is that it has a number of representational layers in-between and that it has a bottleneck in the middle. This layer provides us with a **compressed representation** that is also interesting as it has been shown that it provides an interesting similarity measure between different inputs that can represent some semantic components. It is this kind of learning that provides a stream to more 'intelligent' processing.

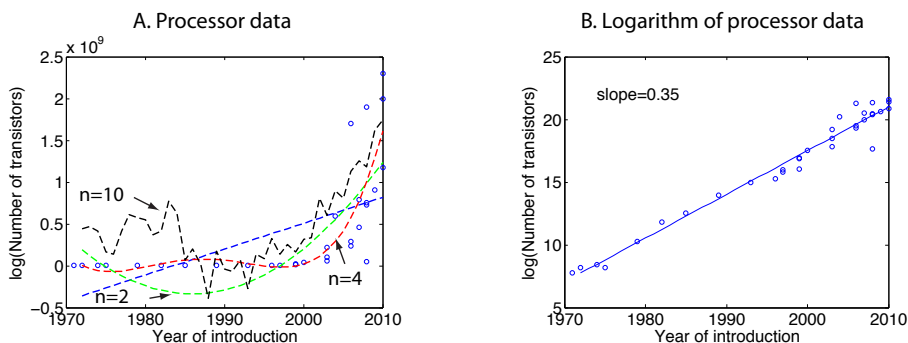
## 1.4 Causal learning, overfitting and regression

I would like to close the introductory overview with one more concept that is central in machine learning. For this we go back to the supervised learning of regression. As mentioned previously finding the right parameterized function is somewhat difficult. There is an important area of **causal learning** that tries to provide specific probabilistic models of the components that provide the necessary foundations of the inference engine. Inference here means a system that can 'argue' about a solution in a probabilistic sense. Such systems fall generally in the domain of Bayesian learning, and we will include some introduction to these important systems in this course. Figure 1.3 shows an famous example from Judea Pearl, one of the inventors of this important modelling framework.



**Fig. 1.3** Judea Pearl and causal modeling

Above we have discussed a case where we assumed a linear function, but regression with more general non-linear functions brings another level of challenges. An example of data that do not follow a linear trend is shown in Fig.1.4A. There, the number of transistors of microprocessors is plotted against the year each processor was introduced. This plot includes a line showing a linear regression, which is of course not very good. It is however interesting to note that this linear approximation shows some systematic deviation or **bias** suggest that we have to take more complex functions into account. 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 understanding of scientific methods are required.



**Fig. 1.4** Data showing the number of transistors in microprocessors plotted against the year they were introduced. (A) Data and some linear and polynomial fits of the data. (B) Logarithm of the data and linear fit of these data.

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.1.4B. 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

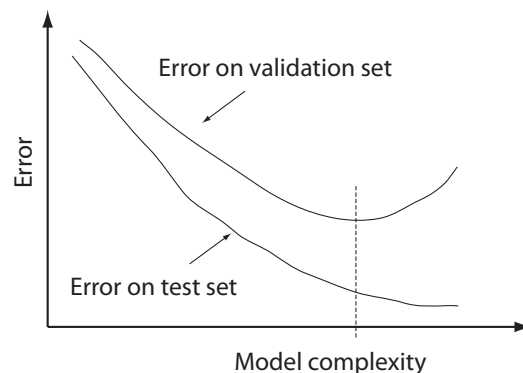
$$y = \theta_0 x^0 + \theta_1 x^1 + \theta_2 x^2 + \dots + \theta_n x^n \quad (1.13)$$

We can again use a regression method to determine the parameters from the data by minimizing the LMS error function between the hypothesis and the data. The LMS regression of the transistor data to a polynomial of orders  $n = 2, 4, 10$  are shown in Fig.??A 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. A schematic figure showing the bias-variance tradeoff is shown in Fig. 1.5. 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. 1.5** Illustration of bias-variance tradeoff.

Deep neural networks are a form of high dimensional non-linear fitting function, and preventing overfitting is therefore a very important component in deep learning. Deep networks have many free parameters, and large data sets (big data) has therefore been important for the recent progress in this area in combination with other techniques to prevent overfitting such as a technique called dropout that we will discuss later. In general one can think about techniques to prevent overfitting as restricting the possible range of the parameters. Indeed, learning from data already represents providing information about the values of the parameters, and restricting such ranges further is a key element in machine learning. This area is generally discussed under the heading of **regularization**.

Machine learning methods are often easy to apply through application programs that implement these techniques. This is good news. However, a deeper understanding of the methods is necessary to make these applications and their conclusions appropriate. The machine learning algorithms will come up with some predictions, but if these predictions are sensible is important to comprehend and evaluate. Machine learning

education needs therefore to go beyond learning how to run an application program, and this course thrives to find a balance between practical applications and their theoretical foundation.

## 2 ML programming with Python

---

This chapter is a brief introduction to scientific programming with the Python programming environment and more specific examples of using ML libraries. The basic idea behind this chapter is to jump right away into some examples. So we will intentionally only cover some essential basics to keep us going. We will continue to refine programming issues throughout the course and will talk about the science behind the algorithms later.

### 2.1 General scientific programming in Python

#### 2.1.1 Resources and installation

Python is a high level programming language that gains increasing popularity in the machine learning community (Matlab has been dominating before). We assume some familiarity with programming concepts and concentrate on the specific environment and supporting libraries for this class. A comprehensive documentation and tutorials are available at <https://www.python.org>. Also, more specific resources for scientific computing with Python are:

```
http://docs.scipy.org/doc/numpy-dev/user/quickstart.html
http://docs.scipy.org/doc/numpy-dev/user/numpy-for-matlab-users.html
http://www.scipy-lectures.org/index.html
http://matplotlib.org/users/beginner.html
```

We will be using the Ubuntu operating system with Python 3 and supporting programs. An image with these components is provided so that you can install this on your own computer or use computers in our faculty. Please see the helps desk for any problems with the installation.

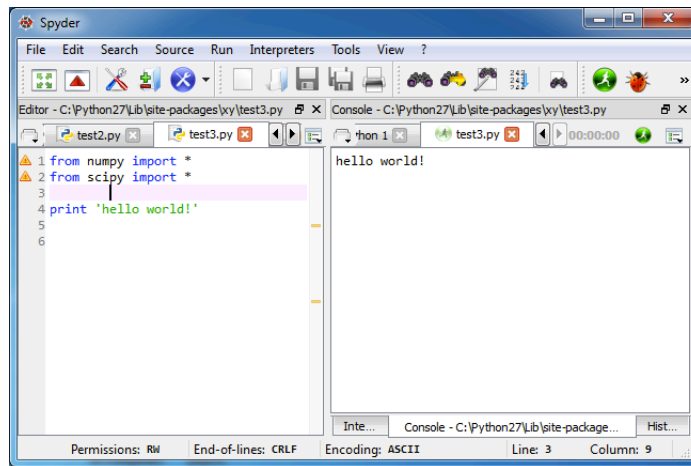
#### 2.1.2 The Spyder programming environment

We will be using a programming environment called Spyder that provides a graphical user interface to basic tools such as an editor and an python interpreter. Start Spyder and you should see the programing environment similar to the one shown in Figure 2.1. On the left is a editor window in which we can write the program. On the right is the console that executes and interpreted the code.

#### 2.1.3 Main programming constructs

The following lines of course are intended to show the syntax of the basic programming constructs that we need in this course. We will be using Python as a scientific





**Fig. 2.1** The Spyder programming environment for Python.

programming language, and we will always import the `pylab` library that includes a lot of useful functions.

```
from pylab import *
```

Next we consider the basic data types that we are using. We are mainly concerned here with numerical data of which a scalar is the simplest example,

```
#basic data types
scalar=4
print( scalar )
```

Note that we can include comment lines with the hashtag symbol. We also included a `print` function that will report the value of the variable `scalar` that we defined here. Note that the type of the variables are dynamically assigned in python.

Most of the time we need to work on a large collection of data so that we need a construct to access the data collection. For we use some form of lists. There are slightly different concepts of such constructs in python. The basic one dimensional list is given by enclosing a semicolon-separated list in square brackets such as

```
list = [ 1 , 2 , 3 ]
```

However, basically need to perform well defined mathematical operations with these list, which makes these one dimensional list formally a vector. The basic data structure for a collection of data is usually called an array in computer science. Thus, the mathematical concept of a **vector** is a one-dimensional array with some operations defined to it. To confuse this matter a bit more, we will use the numpy construct of an array to implement a vector. `numpy` is a collection of numerical functions that is included in `pylab`. The `numpy` function `array()` turns a Python list into a vector,

```
vector=array ([ 1 , 2 , 3 ])
print( vector )
print( vector [ 1 ])
```

As shown in the last line, we can access an element of the array with indices in square brackets, and the first element in an array has the index 0.

Of course, we can generalize such data collections to higher dimension arrangements. For example, a two dimensional array with the appropriate definition of mathematical operations is called a **matrix** and can be defined and accessed in python like

```
matrix=array ([[1 ,2 ,3] ,[4 ,5 ,6]])
print( matrix )
print( matrix [1][2])
```

Corresponding mathematical constructs in higher dimensions are called a **tensor** that we will talk about later. A basic matrix multiplication, also called a dot product, is implemented as function `dot(a,b)` and in Python 3 also as operator `@`,

```
matrix2=array ([[5 ,5 ,6] ,[7 ,8 ,9]])
result=matrix @ matrix2.T
print( result )
```

So far we have discussed the basic data types that we need. Besides these numerical data types there are of course others such as logical or characters. Please consult Python documentation for these data types when needed. We now mention three more basic programming constructs, that of loops, logical statements, and functions.

To loop through some code one can use the following construct,

```
for i in range(4):
    print(i)
```

which starts at `i=0` and goes in steps of one until `i=3`. Note that Python is sensitive to the code position; the indented code represents the block of statements executed inside the loop.

A conditional statement takes the form

```
if scalar <1:
    print( " true " )
else :
    print( " false " )
```

Again note the indentation to specify the block of code for each condition.

To structure code better, specifically to define program constructs that can be reused, we have the option to define functions like

```
def func( arg1 , arg2=10):
    arg=arg1+arg2
    return arg;
```

Variables are passed by reference.

One final example of basic programming we need is that of plotting graphs. Plotting graphs is a useful scientific tool, and an example of a basic line plot can be given in the following code.

```
#plotting
```

```
x=arange(100) #same as array(range(10))
y=sin(0.1*x)
plot(x,y)
```

When you submit plots in an assignment or paper, you always need axis labels to know what is plotted. This can be done with

```
xlabel("x")
ylabel("y")
```

## Exercises

1. Write a Python function that takes a character string and prints out the character string in reverse order.
2. Write a Python program that uses 3-dimensional plotting routines to plot a two dimensional Gaussian function.
3. more

## 2.2 Cross validation example from Intro

To practice Python programming and to deepen our understanding of cross validation, we will now review the program that was used to produce the linear model with cross validation of the example.

In the code below we start by generating the training set consisting of 4 data points that are derived from a line  $y = 2x + 3$  with added Gaussian noise,

```
from pylab import *
# training set
n=4
x=array(range(4)); y=2*x+3+randn(n)
plot(x,y,'*')
```

For the learning tasks we chose a linear model  $\hat{y} = ax + b$ , see it as a wise choice, with two parameters, the slope  $a$  the the intercept  $b$ . Our task is now to determine the values for these parameters from the data. Since we have only two unknown we only need two data point to determine, so let us choose the first two,

```
# one example
a=(y[1]-y[0])/(x[1]-x[0])
b=y[0]-a*x[0]
yhat=a*x+b
plot(x,yhat,'b—')
ytrue=2*x+3
plot(x,ytrue,'g-')
```

We plotted here this specific solution in black and well as the best possible solution in green which we know as we know what the parameters were that have been used

to generate the data and also since the Gaussian noise is unbiased (symmetric around zero)

Of course, this solution is only one possible solution since we could have used any other pair to determine the parameters. Indeed, we should try out all and use all the remaining points to see how good one specific solution will predict the reminder. This is exactly the essence of cross validation.

To determine all the possible combination we use a preferred function from the `itertools` collection,

```
# cross validation
import itertools
c = list(itertools.combinations(x, 2))
```

The list `c` contains now all possible pairs. We then loop over all the pairs and determine the parameters for each choice, and also calculate the error for predicting the other data points not used in the determination of the parameters,

```
#try out all possible pairs
error=[]
for i in range(len(c)):
    #train fold
    k=c[i][0]
    l=c[i][1]
    a=(y[l]-y[k])/(x[l]-x[k])
    b=y[k]-a*x[k]

    er=0
    for j in range(n):
        if j!=k and j!=l:
            er=er+(y[j]-a*x[j]-b)**2
    error.extend([er])
```

This ends the loop. We then take the pair with the minimal cross validation error as our final answer,

```
#search for best pair with smallest cross validation error
i=error.index(min(error))
k=c[i][0]
l=c[i][1]
#and use this as answer
a=(y[l]-y[k])/(x[l]-x[k])
b=y[k]-a*x[k]
yhat=a*x+b
plot(x, yhat, 'r—')
```

## 2.3 Classification with support vector machine using scikit-learn

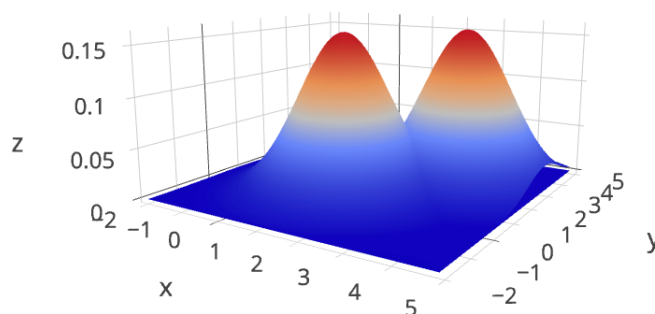
We will now show an explicit example of classification using a support vector machine from the scikit-learn collection of machine learning methods at <http://scikit-learn.org/>. This library started as a Google Summer of Code project by David Cournapeau and developed into an open source library. We will later have a look of what kind of algorithms are implemented, but for now we are just using one of the methods for classification called support vector machine. The SVM in scikit-learn is actually a wrapper to the very popular SVMLIB implementation by Chih-Chung Chang and Chih-Jen Lin. We will go through the code here with some explanations.

We begin as usual by importing libraries we need and to create the training set.

```
from pylab import *
from sklearn import svm

# training
n=100
x1=array([randn(n)+1,randn(n)+1]); y1=zeros(n)
x2=array([randn(n)+3,randn(n)+3]); y2=zeros(n)+1
x = hstack((x1, x2)).T
y = hstack((y1, y2))
```

In real world application the data set is of course usually supplied by a third party often through a data file. Here we simulate an example that consists of two 2-dimensional Gaussian classes each with a unit covariance matrix and different means. The mean of the first class is  $\mu_1 = (1, 1)$  and the second class has  $\mu_2 = (3, 3)$ . The distributions of these two classes are shown in Fig. 2.2.



**Fig. 2.2** Two 2d-Gaussian curves with unit covariance and different means.

We now define a classifier model. We are using a Support Vector Classifier, as specific support vector machine for classification, with two parameters that we will discuss only later, namely we are using a linear kernel and regularization parameter  $C = 1$ ,

```
SVC = svm.SVC(kernel='linear', C=1)
```

```
SVC. fit (x, y)
```

The second line implements the learning, that is it takes the examples in arrays `x` and `y` and fit the model to it. At this point we have a trained model `SVC` that we can use to predict data. We will test its performance with some new sample data,

```
# testing
x1=array ([ randn (n)+1 ,randn (n)+1]); y1=zeros (n)
x2=array ([ randn (n)+3 ,randn (n)+3]); y2=zeros (n)+1
x = hstack ((x1 ,x2)).T
y = hstack ((y1 ,y2))
```

that we also plot with different symbols and color. We use the model for predicting the labels for the class with the command

```
a=SVC. predict (x)
```

and calculate the percentage of correct labels with

```
print (" Percentage _Correct : " , (n-sum ( abs (y-a)))/n)
```

Finally, we also like to plot the results

```
plot (x1 [0 ,:] ,x1 [1 ,:] , 'xr' )
plot (x2 [0 ,:] ,x2 [1 ,:] , 'ob' )
show ()
```

## 2.4 Other classification methods including MLP with Tensorflow

The final example here is using two more classifiers in addition to the SVM on the same two-Gaussian example, namely a random forrest (RF) classifier and a multilayer perceptron MLP). The RF is also implemented in sklearn and is hence verst similar. We are only changing the name of the model. For the MLP we use Googles Tensorflow implementation which is also quite similar to the sklearn notation. The only difference is that the model has of course different parameters and hyper-parameters.

```
from pylab import *
from sklearn import svm
from sklearn.ensemble import RandomForestClassifier
import tensorflow
from tensorflow.contrib import skflow
```

```
# training
n=100
x1=array ([ randn (n)+1 ,randn (n)+1]); y1=zeros (n, int)
x2=array ([ randn (n)+3 ,randn (n)+3]); y2=zeros (n, int)+1

x = hstack ((x1 ,x2)).T
y = hstack ((y1 ,y2))
```

```

SVC = svm.SVC(kernel='linear', C=1)
SVC.fit(x, y)

RF = RandomForestClassifier(n_estimators=10)
RF.fit(x, y)

MLP = skflow.TensorFlowDNNClassifier(
    hidden_units=[10, 20, 10],
    n_classes=2,
    batch_size=128,
    steps=500,
    learning_rate=0.05)
MLP.fit(x, y)

# testing
x1=array([randn(n)+1,randn(n)+1]); y1=zeros(n)
x2=array([randn(n)+3,randn(n)+3]); y2=zeros(n)+1

plot(x1[0,:],x1[1,:], 'xr')
plot(x2[0,:],x2[1,:], 'ob')
show()

x = hstack((x1,x2)).T
y = hstack((y1,y2))

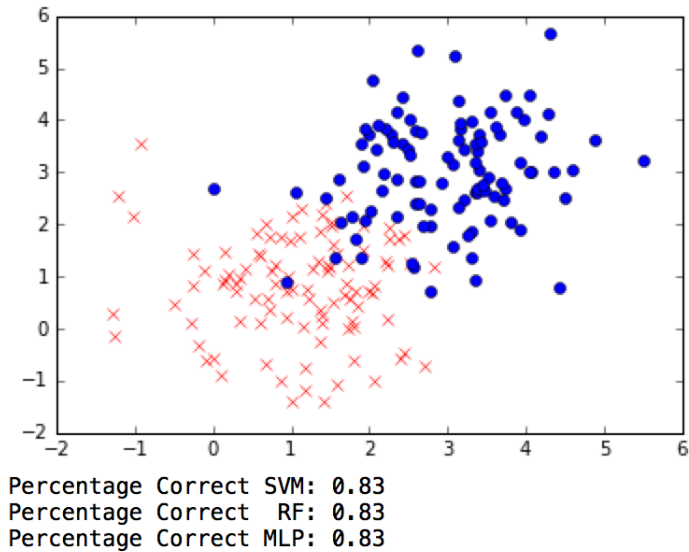
a=SVC.predict(x)
print("Percentage_Correct_SVM: ", (n-sum(abs(y-a)))/n)
b=RF.predict(x)
print("Percentage_Correct_RF: ", (n-sum(abs(y-b)))/n)
c=MLP.predict(x)
print("Percentage_Correct_MLP: ", (n-sum(abs(y-c)))/n)

```

The result of running the program is shown in Fig. 2.3. All three classifiers give the same result in this run which is close to the optimal result in this example. When running this program repeatedly there will be slight differences in the answers. We will later discuss the stochastic nature of machine learning.

## Exercise

1. What is the accuracy of the optimal solution of the two-Gaussian problem with the parameters used in the above example? Calculate this bound analytically.
2. This exercise is a little project where you should apply a binary machine learning classification to a problem of your choice. You might have an own data set



**Fig. 2.3** The plot shows the data points in the two-Gaussian example that consists of two Gaussian classes with the same unit co-variance but different mean values. The problem is that these classes overlap. Below the figure is the percentage correct of three machine learning classifiers, that of a Support Vector Machine (SVM), a Random Forrest (RF) classifier, and a multilayer perceptron (MLP).

or you could choose a data set from the UCI machine learning repository at <http://archive.ics.uci.edu/ml/>. Note that many of the problems are not binary, but you could often change these data sets to a binary problem. For example, a high value could be changed into two classes of small and large.



## 3 Basic probability theory

---

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.

### 3.1 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**,

$$P_X(x) = P(X = x) \quad (3.1)$$

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.<sup>1</sup> From this definition it follows that the frequency of having any of the possible values is equal to one, which is the normalization condition

$$\sum_x P_X(x) = 1. \quad (3.2)$$

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 3.2 then becomes an integral, and normalization condition for a continuous random variable is

$$\int_x p_X(x)dx = 1. \quad (3.3)$$

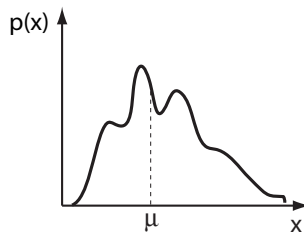
We will formulate the rest of this section in terms of continuous random variables. The corresponding formulas for discrete random variables can easily be deduced by replacing the integrals over the pdf with sums over the probability function. It is also possible to use the  $\delta$ -function, outlined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

<sup>1</sup>Probabilities are sometimes written as a percentage, but we will stick to the fractional notation.

### 3.2 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 3.3. 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

$$\mu = \int_{-\infty}^{\infty} xp(x)dx. \quad (3.4)$$

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 3.4 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  $\mu$  by

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i p_i. \quad (3.5)$$

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**,  $\sigma^2$ , and is defined as

$$\sigma^2 = \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx. \quad (3.6)$$

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

$$m^n = \int_{-\infty}^{\infty} (x - \mu)^n f(x) dx. \quad (3.7)$$

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

$$\sigma^2 = \int_{-\infty}^{\infty} (x - \mu)^2 f(x) dx. \quad (3.8)$$

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**,

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad (3.9)$$

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

$$s_1^2 = \frac{1}{n} \sum_{i=1}^n (\bar{x} - x_i)^2. \quad (3.10)$$

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**

$$s_2^2 = \frac{1}{n-1} \sum_{i=1}^n (\bar{x} - x_i)^2, \quad (3.11)$$

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, and incomplete list of moments does not uniquely define a pdf. Note that the Gaussian distributions is fully characterized by the first two moments that are given by the mean  $\mu$  and the variance  $\sigma^2$ . This is however not the case for other distributions, and the usefulness of reporting only the statistics of the first two moments can then be questioned.

### 3.3 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 is a list of examples with some discrete and several continuous distributions. Most examples are discussed as one-dimensional distributions except the last example, which is a higher dimensional distribution.

### 3.4 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:  $p(1 - p)$

#### 3.4.1 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:  $np_i$

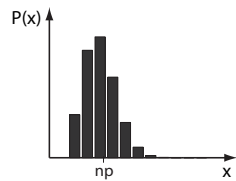
variance:  $np_i(1 - p_i)$

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

$$\binom{n}{x} = \frac{n!}{x!(n-x)!} \quad (3.12)$$

and is given by the MATLAB function `nchoosek`.



Probability function:

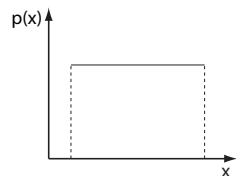
$$P(x) = \binom{n}{x} p^x (1-p)^{n-x}$$

mean:  $np$

variance:  $np(1-p)$

### 3.4.2 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:

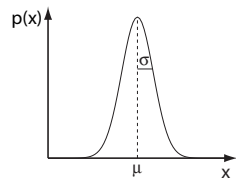
$$p(x) = \frac{1}{b-a}$$

mean:  $(a+b)/2$

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

### 3.4.3 Normal (Gaussian) distribution

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



Probability density function:

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

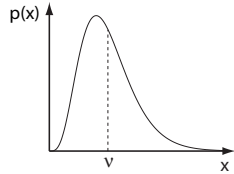
mean:  $\mu$

variance:  $\sigma^2$

### 3.4.4 Chi-square distribution

The sum of the squares of normally distributed random numbers is chi-square distributed and depends on a parameter  $\nu$  that is equal to the mean.  $\Gamma$  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:  $\nu$

variance:  $2\nu$

### 3.4.5 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

$$p(x_1, \dots, x_n) = p(\mathbf{x}) = \frac{1}{(\sqrt{2\pi})^n \sqrt{|\det(\Sigma)|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1}(\mathbf{x} - \mu)\right). \tag{3.13}$$

This is a straight forward generalization of the one-dimensional Gaussian distribution mentioned before where the mean is now a vector,  $\mu$  and the variance generalizes to a covariance matrix,  $\Sigma = [\text{Cov}[X_i, X_j]]_{i=1,2,\dots,k; j=1,2,\dots,k}$  which must be symmetric and positive semi-definit. An example with mean  $\mu = (1 \ 2)^T$  and covariance  $\Sigma = \begin{pmatrix} 1 & 0.5 \\ 0.5 & 1 \end{pmatrix}$  is shown in Fig,3.1.

## 3.5 Cumulative probability (density) function and the Gaussian error function

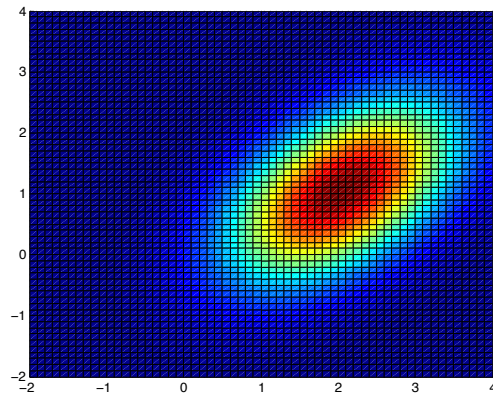
We have mainly discussed probabilities of single values as specified by the probability (density) functions. However, in many cases we want to know the probabilities of having values in a certain range. Indeed, the probability of a specific value of a continuous random variable is actually infinitesimally 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

$$P(x_1 \leq X \leq x_2) = \int_{x_1}^{x_2} p(x) dx. \tag{3.14}$$

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 3.14 then becomes a function of  $y$ . This defines the **Gaussian error function**

$$\frac{1}{\sqrt{2\pi}\sigma} \int_0^y e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx = \frac{1}{2} \text{erf}\left(\frac{y-\mu}{\sqrt{2}\sigma}\right). \tag{3.15}$$

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`.



**Fig. 3.1** Multivariate Gaussian with mean  $\mu = (1 \ 2)^T$  and covariance  $\Sigma = (1 \ 0.5; 0.5 \ 1)$ .

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 3.14 is when  $x_1$  in the equation is equal to the lowest possible value of the random variable (usually  $-\infty$ ). The integral in eqn 3.14 then 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),<sup>2</sup>

$$P^{\text{cum}}(x < y) = \int_{-\infty}^y p(x)dx, \quad (3.16)$$

which we will utilize further below.

### 3.6 Functions of random variables and the central limit theorem

A function of a random variable  $X$ ,

$$Y = f(X), \quad (3.17)$$

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

$$p_X(x) = \begin{cases} 1 & \text{for } 0 \leq x \leq 1, \\ 0 & \text{otherwise.} \end{cases} \quad (3.18)$$

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

$$Y = e^{-X^2}. \quad (3.19)$$

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 3.16 by noting that

$$P(Y \leq y) = P(e^{-X^2} \leq y) = P(X \geq \sqrt{-\ln y}). \quad (3.20)$$

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

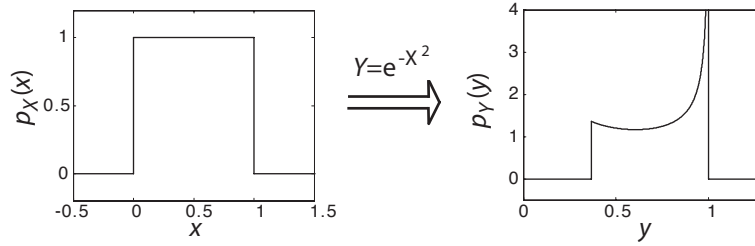
$$P(X \geq \sqrt{-\ln y}) = \begin{cases} \int_{\sqrt{-\ln y}}^1 p_X(x)dy = 1 - \sqrt{-\ln y} & \text{for } e^{-1} \leq y \leq 1, \\ 0 & \text{otherwise.} \end{cases} \quad (3.21)$$

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

$$p_Y(y) = \begin{cases} 1 - \sqrt{-\ln y} & \text{for } e^{-1} \leq y \leq 1, \\ 0 & \text{otherwise.} \end{cases} \quad (3.22)$$

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

<sup>2</sup>Note that this is a probability function, not a density function.



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,

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i, \quad (3.23)$$

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  $\mu$  and variance  $\sigma$  is approximately normally distributed with mean  $\mu$  and variance  $\sigma/n$  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.

### 3.7 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

$$d^{\text{KL}}(p, q) = \int p(x) \log\left(\frac{p(x)}{q(x)}\right) dx \quad (3.24)$$

$$= \int p(x) \log(p(x)) dx - \int p(x) \log(q(x)) dx \quad (3.25)$$

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

### 3.8 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.

#### 3.8.1 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

$$\text{joined distribution: } p(x, y) = p(X = x, Y = y). \quad (3.26)$$

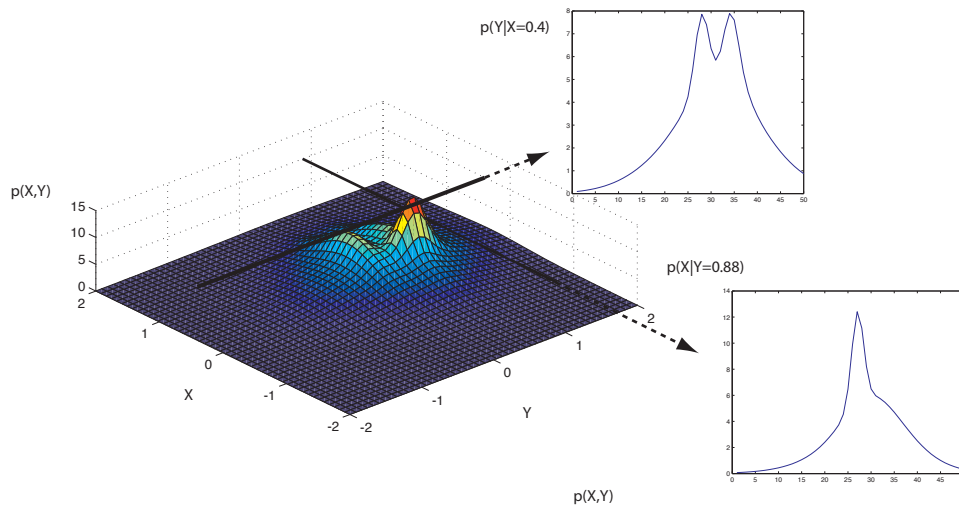
This is a two dimensional functions. The two dimensions refers here to the number of variables, although a plot of this function would be a three dimensional plot. An example is shown in Fig.3.2. All the information we can have about a stochastic system is encapsulated in the joined pdf. The slice of this function, given the value of one variable, say  $y$ , is the

$$\text{conditional distribution: } p(x|y) = p(X = x|Y = y). \quad (3.27)$$

A conditional pdf is also illustrated in Fig.3.2 If we sum over all realizations of  $y$  we get the

$$\text{marginal distribution: } p(x) = \int p(x, y)dy. \quad (3.28)$$

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 joined 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.



**Fig. 3.2** Example of a two-dimensional probability density function (pdf) and some examples of conditional pdfs.

### 3.8.2 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 decomposed into the product of a conditional and a marginal distribution,

$$p(x, y) = p(x|y)p(y) = p(y|x)p(x), \tag{3.29}$$

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

$$p(x_1, x_2, \dots, x_n) = p(x_n|x_1, \dots, x_{n-1})p(x_1, \dots, x_{n-1}) \tag{3.30}$$

$$= p(x_n|x_1, \dots, x_{n-1}) * \dots * p(x_2|x_1) * p(x_1) \tag{3.31}$$

$$= \prod_{i=1}^n p(x_i|x_{i-1}, \dots, x_1) \tag{3.32}$$

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

$$p(x|y) = p(x). \tag{3.33}$$

Using the chain rule eq.3.29, we can write this also as

$$p(x, y) = p(x)p(y), \tag{3.34}$$

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

$$p(x, y|z) = p(x|z)p(y|z). \quad (3.35)$$

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

### 3.8.3 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. 3.29 by  $p(x)$ , which is possible as long as  $p(x) > 0$ , we get the identity

$$p(x|y) = \frac{p(y|x)p(x)}{p(y)}, \quad (3.36)$$

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.