

Data Analysis and Machine Learning: Linear Regression and more Advanced Regression Analysis

Morten Hjorth-Jensen^{1,2}

¹Department of Physics, University of Oslo

²Department of Physics and Astronomy and National Superconducting Cyclotron Laboratory, Michigan State University

May 10, 2019

Why Linear Regression (aka Ordinary Least Squares)

Fitting a continuous function with linear parameterization in terms of the parameters β .

- Method of choice for fitting a continuous function!
- Gives an excellent introduction to central Machine Learning features with **understandable pedagogical** links to other methods like **Neural Networks, Support Vector Machines** etc
- Analytical expression for the parameters β
- Analytical expressions for statistical properties like mean values, variances, confidence intervals and more
- Analytical relation with probabilistic interpretations
- Easy to introduce basic concepts like bias-variance tradeoff, cross-validation, resampling and regularization techniques and many other ML topics
- Easy to code! And links well with classification problems and logistic regression and neural networks
- Allows for **easy** hands-on understanding of gradient descent methods
- and many more features

Regression analysis, overarching aims

Regression modeling deals with the description of the sampling distribution of a given random variable y and how it varies as function of another variable or a set of such variables $\mathbf{x} = [x_0, x_1, \dots, x_{n-1}]^T$. The first variable is called the **dependent**, the **outcome** or the **response** variable while the set of variables \mathbf{x} is called the independent variable, or the predictor variable or the explanatory variable.

A regression model aims at finding a likelihood function $p(\mathbf{y}|\mathbf{x})$, that is the conditional distribution for \mathbf{y} with a given \mathbf{x} . The estimation of $p(\mathbf{y}|\mathbf{x})$ is made using a data set with

- n cases $i = 0, 1, 2, \dots, n-1$
- Response (target, dependent or outcome) variable y_i with $i = 0, 1, 2, \dots, n-1$
- p so-called explanatory (independent or predictor) variables $\mathbf{x}_i = [x_{i0}, x_{i1}, \dots, x_{ip-1}]$ with $i = 0, 1, 2, \dots, n-1$ and explanatory variables running from 0 to $p-1$. See below for more explicit examples.

The goal of the regression analysis is to extract/exploit relationship between \mathbf{y} and \mathbf{X} in or to infer causal dependencies, approximations to the likelihood functions, functional relationships and to make predictions, making fits and many other things.

Regression analysis, overarching aims II

Consider an experiment in which p characteristics of n samples are measured. The data from this experiment, for various explanatory variables p are normally represented by a matrix \mathbf{X} .

The matrix \mathbf{X} is called the *design matrix*. Additional information of the samples is available in the form of \mathbf{y} (also as above). The variable \mathbf{y} is generally referred to as the *response variable*. The aim of regression analysis is to explain \mathbf{y} in terms of \mathbf{X} through a functional relationship like $y_i = f(\mathbf{X}_{i,*})$. When no prior knowledge on the form of $f(\cdot)$ is available, it is common to assume a linear relationship between \mathbf{X} and \mathbf{y} . This assumption gives rise to the *linear regression model* where $\boldsymbol{\beta} = [\beta_0, \dots, \beta_{p-1}]^T$ are the *regression parameters*.

Linear regression gives us a set of analytical equations for the parameters β_j .

Examples

In order to understand the relation among the predictors p , the set of data n and the target (outcome, output etc) \mathbf{y} , consider the model we discussed for describing nuclear binding energies.

There we assumed that we could parametrize the data using a polynomial approximation based on the liquid drop model. Assuming

$$BE(A) = a_0 + a_1A + a_2A^{2/3} + a_3A^{-1/3} + a_4A^{-1},$$

we have five predictors, that is the intercept, the A dependent term, the $A^{2/3}$ term and the $A^{-1/3}$ and A^{-1} terms. This gives $p = 0, 1, 2, 3, 4$. Furthermore we have n entries for each predictor. It means that our design matrix is a $p \times n$ matrix \mathbf{X} .

Here the predictors are based on a model we have made. A popular data set which is widely encountered in ML applications is the so-called [credit card default data from Taiwan](#). The data set contains data on $n = 30000$ credit card holders with predictors like gender, marital status, age, profession, education, etc. In total there are 24 such predictors or attributes leading to a design matrix of dimensionality 24×30000

General linear models

Before we proceed let us study a case from linear algebra where we aim at fitting a set of data $\mathbf{y} = [y_0, y_1, \dots, y_{n-1}]$. We could think of these data as a result of an experiment or a complicated numerical experiment. These data are functions of a series of variables $\mathbf{x} = [x_0, x_1, \dots, x_{n-1}]$, that is $y_i = y(x_i)$ with $i = 0, 1, 2, \dots, n-1$. The variables x_i could represent physical quantities like time, temperature, position etc. We assume that $y(x)$ is a smooth function.

Since obtaining these data points may not be trivial, we want to use these data to fit a function which can allow us to make predictions for values of y which are not in the present set. The perhaps simplest approach is to assume we can parametrize our function in terms of a polynomial of degree $n-1$ with n points, that is

$$y = y(x) \rightarrow y(x_i) = \tilde{y}_i + \epsilon_i = \sum_{j=0}^{n-1} \beta_j x_i^j + \epsilon_i,$$

where ϵ_i is the error in our approximation.

Rewriting the fitting procedure as a linear algebra problem

For every set of values y_i, x_i we have thus the corresponding set of equations

$$\begin{aligned} y_0 &= \beta_0 + \beta_1 x_0^1 + \beta_2 x_0^2 + \dots + \beta_{n-1} x_0^{n-1} + \epsilon_0 \\ y_1 &= \beta_0 + \beta_1 x_1^1 + \beta_2 x_1^2 + \dots + \beta_{n-1} x_1^{n-1} + \epsilon_1 \\ y_2 &= \beta_0 + \beta_1 x_2^1 + \beta_2 x_2^2 + \dots + \beta_{n-1} x_2^{n-1} + \epsilon_2 \\ &\dots\dots\dots \\ y_{n-1} &= \beta_0 + \beta_1 x_{n-1}^1 + \beta_2 x_{n-1}^2 + \dots + \beta_{n-1} x_{n-1}^{n-1} + \epsilon_{n-1}. \end{aligned}$$

Rewriting the fitting procedure as a linear algebra problem, more details

Defining the vectors

$$\mathbf{y} = [y_0, y_1, y_2, \dots, y_{n-1}]^T,$$

and

$$\boldsymbol{\beta} = [\beta_0, \beta_1, \beta_2, \dots, \beta_{n-1}]^T,$$

and

$$\boldsymbol{\epsilon} = [\epsilon_0, \epsilon_1, \epsilon_2, \dots, \epsilon_{n-1}]^T,$$

and the design matrix

$$\mathbf{X} = \begin{bmatrix} 1 & x_0^1 & x_0^2 & \dots & \dots & x_0^{n-1} \\ 1 & x_1^1 & x_1^2 & \dots & \dots & x_1^{n-1} \\ 1 & x_2^1 & x_2^2 & \dots & \dots & x_2^{n-1} \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & x_{n-1}^1 & x_{n-1}^2 & \dots & \dots & x_{n-1}^{n-1} \end{bmatrix}$$

we can rewrite our equations as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}.$$

Generalizing the fitting procedure as a linear algebra problem

We are obviously not limited to the above polynomial. We could replace the various powers of x with elements of Fourier series or instead of x_i^j we could have $\cos(jx_i)$ or $\sin(jx_i)$, or time series or other orthogonal functions. For every set of values y_i, x_i we can then generalize the equations to

$$\begin{aligned} y_0 &= \beta_0 x_{00} + \beta_1 x_{01} + \beta_2 x_{02} + \dots + \beta_{n-1} x_{0n-1} + \epsilon_0 \\ y_1 &= \beta_0 x_{10} + \beta_1 x_{11} + \beta_2 x_{12} + \dots + \beta_{n-1} x_{1n-1} + \epsilon_1 \\ y_2 &= \beta_0 x_{20} + \beta_1 x_{21} + \beta_2 x_{22} + \dots + \beta_{n-1} x_{2n-1} + \epsilon_2 \\ &\dots\dots\dots \\ y_i &= \beta_0 x_{i0} + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_{n-1} x_{in-1} + \epsilon_i \\ &\dots\dots\dots \\ y_{n-1} &= \beta_0 x_{n-1,0} + \beta_1 x_{n-1,1} + \beta_2 x_{n-1,2} + \dots + \beta_{n-1} x_{n-1,n-1} + \epsilon_{n-1}. \end{aligned}$$

Note that we have $p = n$ here. The matrix is symmetric. This is generally not the case!

Generalizing the fitting procedure as a linear algebra problem

We redefine in turn the matrix \mathbf{X} as

$$\mathbf{X} = \begin{bmatrix} x_{00} & x_{01} & x_{02} & \dots & \dots & x_{0,n-1} \\ x_{10} & x_{11} & x_{12} & \dots & \dots & x_{1,n-1} \\ x_{20} & x_{21} & x_{22} & \dots & \dots & x_{2,n-1} \\ \dots & \dots & \dots & \dots & \dots & \dots \\ x_{n-1,0} & x_{n-1,1} & x_{n-1,2} & \dots & \dots & x_{n-1,n-1} \end{bmatrix}$$

and without loss of generality we rewrite again our equations as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}.$$

The left-hand side of this equation is known. Our error vector $\boldsymbol{\epsilon}$ and the parameter vector $\boldsymbol{\beta}$ are our unknown quantities. How can we obtain the optimal set of β_i values?

Optimizing our parameters

We have defined the matrix \mathbf{X} via the equations

$$\begin{aligned} y_0 &= \beta_0 x_{00} + \beta_1 x_{01} + \beta_2 x_{02} + \dots + \beta_{n-1} x_{0,n-1} + \epsilon_0 \\ y_1 &= \beta_0 x_{10} + \beta_1 x_{11} + \beta_2 x_{12} + \dots + \beta_{n-1} x_{1,n-1} + \epsilon_1 \\ y_2 &= \beta_0 x_{20} + \beta_1 x_{21} + \beta_2 x_{22} + \dots + \beta_{n-1} x_{2,n-1} + \epsilon_1 \\ &\dots\dots\dots \\ y_i &= \beta_0 x_{i0} + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_{n-1} x_{i,n-1} + \epsilon_i \\ &\dots\dots\dots \\ y_{n-1} &= \beta_0 x_{n-1,0} + \beta_1 x_{n-1,1} + \beta_2 x_{n-1,2} + \dots + \beta_{n-1} x_{n-1,n-1} + \epsilon_{n-1}. \end{aligned}$$

As we noted above, we stayed with a system with the design matrix $\mathbf{X} \in \mathbb{R}^{n \times n}$, that is we have $p = n$. For reasons to come later (algorithmic arguments) we will hereafter define our matrix as $\mathbf{X} \in \mathbb{R}^{p \times n}$, with the predictors referring to the row numbers and the entries n being the column elements.

Optimizing our parameters, more details

With the above we use the design matrix to define the approximation $\tilde{\mathbf{y}}$ via the unknown quantity $\boldsymbol{\beta}$ as

$$\tilde{\mathbf{y}} = \mathbf{X}^T \boldsymbol{\beta},$$

and in order to find the optimal parameters β_i instead of solving the above linear algebra problem, we define a function which gives a measure of the spread

between the values y_i (which represent hopefully the exact values) and the parametrized values \tilde{y}_i , namely

$$Q(\beta) = \sum_{i=0}^{n-1} (y_i - \tilde{y}_i)^2 = (\mathbf{y} - \tilde{\mathbf{y}})^T (\mathbf{y} - \tilde{\mathbf{y}}),$$

or using the matrix \mathbf{X} and in a more compact matrix-vector notation as

$$Q(\beta) = (\mathbf{y} - \mathbf{X}^T \beta)^T (\mathbf{y} - \mathbf{X}^T \beta).$$

Note: It is normal to have a term $1/n$ in front of this equation since it can then be linked to statistical properties like the variance of a given variable. We come back to this in our discussion of the so-called bias-variance trade-off. It is also common to define the function Q as

$$Q(\beta) = \frac{1}{2n} \sum_{i=0}^{n-1} (y_i - \tilde{y}_i)^2,$$

since when taking the first derivative with respect to the unknown parameters β , the factor of 2 cancels out.

Interpretations and optimizing our parameters

The function

$$Q(\beta) = (\mathbf{y} - \mathbf{X}\beta)^T (\mathbf{y} - \mathbf{X}\beta),$$

can be linked to the variance of the quantity y_i if we interpret the latter as the mean value. When linking below with the maximum likelihood approach below, we will indeed interpret y_i as a mean value (see exercises)

$$y_i = \langle y_i \rangle = \beta_0 x_{i,0} + \beta_1 x_{i,1} + \beta_2 x_{i,2} + \cdots + \beta_{n-1} x_{i,n-1} + \epsilon_i,$$

where $\langle y_i \rangle$ is the mean value. Keep in mind also that till now we have treated y_i as the exact value. Normally, the response (dependent or outcome) variable y_i the outcome of a numerical experiment or another type of experiment and is thus only an approximation to the true value. It is then always accompanied by an error estimate, often limited to a statistical error estimate given by the standard deviation discussed earlier. In the discussion here we will treat y_i as our exact value for the response variable.

In order to find the parameters β_i we will then minimize the spread of $Q(\beta)$ by requiring

$$\frac{\partial Q(\beta)}{\partial \beta_j} = \frac{\partial}{\partial \beta_j} \left[\sum_{i=0}^{n-1} (y_i - \beta_0 x_{i,0} - \beta_1 x_{i,1} - \beta_2 x_{i,2} - \cdots - \beta_{n-1} x_{i,n-1})^2 \right] = 0,$$

which results in

$$\frac{\partial Q(\beta)}{\partial \beta_j} = -2 \left[\sum_{i=0}^{n-1} x_{ij} (y_i - \beta_0 x_{i,0} - \beta_1 x_{i,1} - \beta_2 x_{i,2} - \cdots - \beta_{n-1} x_{i,n-1}) \right] = 0,$$

or in a matrix-vector form as

$$\frac{\partial Q(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = 0 = \mathbf{X}^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}).$$

Interpretations and optimizing our parameters

We can rewrite

$$\frac{\partial Q(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = 0 = \mathbf{X}^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}),$$

as

$$\mathbf{X}^T \mathbf{y} = \mathbf{X}^T \mathbf{X} \boldsymbol{\beta},$$

and if the matrix $\mathbf{X}^T \mathbf{X}$ is invertible we have the solution

$$\boldsymbol{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}.$$

Interpretations and optimizing our parameters

The residuals $\boldsymbol{\epsilon}$ are in turn given by

$$\boldsymbol{\epsilon} = \mathbf{y} - \tilde{\mathbf{y}} = \mathbf{y} - \mathbf{X}\boldsymbol{\beta},$$

and with

$$\mathbf{X}^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = 0,$$

we have

$$\mathbf{X}^T \boldsymbol{\epsilon} = \mathbf{X}^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = 0,$$

meaning that the solution for $\boldsymbol{\beta}$ is the one which minimizes the residuals. Later we will link this with the maximum likelihood approach.

Back to Nuclear Binding Energies

The χ^2 function

Normally, the response (dependent or outcome) variable y_i the outcome of a numerical experiment or another type of experiment and is thus only an approximation to the true value. It is then always accompanied by an error estimate, often limited to a statistical error estimate given by the standard deviation discussed earlier. In the discussion here we will treat y_i as our exact value for the response variable.

Introducing the standard deviation σ_i for each measurement y_i , we define now the χ^2 function (omitting the $1/n$ term) as

$$\chi^2(\boldsymbol{\beta}) = \sum_{i=0}^{n-1} \frac{(y_i - \tilde{y}_i)^2}{\sigma_i^2} = (\mathbf{y} - \tilde{\mathbf{y}})^T \frac{1}{\boldsymbol{\Sigma}^2} (\mathbf{y} - \tilde{\mathbf{y}}),$$

where the matrix $\boldsymbol{\Sigma}$ is a diagonal matrix with σ_i as matrix elements.

The χ^2 function

In order to find the parameters β_i we will then minimize the spread of $\chi^2(\boldsymbol{\beta})$ by requiring

$$\frac{\partial \chi^2(\boldsymbol{\beta})}{\partial \beta_j} = \frac{\partial}{\partial \beta_j} \left[\sum_{i=0}^{n-1} \left(\frac{y_i - \beta_0 x_{i,0} - \beta_1 x_{i,1} - \beta_2 x_{i,2} - \cdots - \beta_{n-1} x_{i,n-1}}{\sigma_i} \right)^2 \right] = 0,$$

which results in

$$\frac{\partial \chi^2(\boldsymbol{\beta})}{\partial \beta_j} = -2 \left[\sum_{i=0}^{n-1} \frac{x_{ij}}{\sigma_i} \left(\frac{y_i - \beta_0 x_{i,0} - \beta_1 x_{i,1} - \beta_2 x_{i,2} - \cdots - \beta_{n-1} x_{i,n-1}}{\sigma_i} \right) \right] = 0,$$

or in a matrix-vector form as

$$\frac{\partial \chi^2(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = 0 = \mathbf{A}^T (\mathbf{b} - \mathbf{A}\boldsymbol{\beta}).$$

where we have defined the matrix $\mathbf{A} = \mathbf{X}/\boldsymbol{\Sigma}$ with matrix elements $a_{ij} = x_{ij}/\sigma_i$ and the vector \mathbf{b} with elements $b_i = y_i/\sigma_i$.

The χ^2 function

We can rewrite

$$\frac{\partial \chi^2(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} = 0 = \mathbf{A}^T (\mathbf{b} - \mathbf{A}\boldsymbol{\beta}),$$

as

$$\mathbf{A}^T \mathbf{b} = \mathbf{A}^T \mathbf{A} \boldsymbol{\beta},$$

and if the matrix $\mathbf{A}^T \mathbf{A}$ is invertible we have the solution

$$\boldsymbol{\beta} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{b}.$$

The χ^2 function

If we then introduce the matrix

$$\mathbf{H} = (\mathbf{A}^T \mathbf{A})^{-1},$$

we have then the following expression for the parameters β_j (the matrix elements of \mathbf{H} are h_{ij})

$$\beta_j = \sum_{k=0}^{p-1} h_{jk} \sum_{i=0}^{n-1} \frac{y_i}{\sigma_i} \frac{x_{ik}}{\sigma_i} = \sum_{k=0}^{p-1} h_{jk} \sum_{i=0}^{n-1} b_i a_{ik}$$

We state without proof the expression for the uncertainty in the parameters β_j as (we leave this as an exercise)

$$\sigma^2(\beta_j) = \sum_{i=0}^{n-1} \sigma_i^2 \left(\frac{\partial \beta_j}{\partial y_i} \right)^2,$$

resulting in

$$\sigma^2(\beta_j) = \left(\sum_{k=0}^{p-1} h_{jk} \sum_{i=0}^{n-1} a_{ik} \right) \left(\sum_{l=0}^{p-1} h_{jl} \sum_{m=0}^{n-1} a_{ml} \right) = h_{jj}!$$

The χ^2 function

The first step here is to approximate the function y with a first-order polynomial, that is we write

$$y = y(x) \rightarrow y(x_i) \approx \beta_0 + \beta_1 x_i.$$

By computing the derivatives of χ^2 with respect to β_0 and β_1 show that these are given by

$$\frac{\partial \chi^2(\beta)}{\partial \beta_0} = -2 \left[\sum_{i=0}^{n-1} \left(\frac{y_i - \beta_0 - \beta_1 x_i}{\sigma_i^2} \right) \right] = 0,$$

and

$$\frac{\partial \chi^2(\beta)}{\partial \beta_1} = -2 \left[\sum_{i=0}^{n-1} x_i \left(\frac{y_i - \beta_0 - \beta_1 x_i}{\sigma_i^2} \right) \right] = 0.$$

The χ^2 function

For a linear fit we don't need to invert a matrix!! Defining

$$\gamma = \sum_{i=0}^{n-1} \frac{1}{\sigma_i^2},$$

$$\gamma_x = \sum_{i=0}^{n-1} \frac{x_i}{\sigma_i^2},$$

$$\gamma_y = \sum_{i=0}^{n-1} \left(\frac{y_i}{\sigma_i^2} \right),$$

$$\gamma_{xx} = \sum_{i=0}^{n-1} \frac{x_i x_i}{\sigma_i^2},$$

$$\gamma_{xy} = \sum_{i=0}^{n-1} \frac{y_i x_i}{\sigma_i^2},$$

we obtain

$$\beta_0 = \frac{\gamma_{xx} \gamma_y - \gamma_x \gamma_{xy}}{\gamma \gamma_{xx} - \gamma_x^2},$$

$$\beta_1 = \frac{\gamma_{xy} \gamma - \gamma_x \gamma_y}{\gamma \gamma_{xx} - \gamma_x^2}.$$

This approach (different linear and non-linear regression) suffers often from both being underdetermined and overdetermined in the unknown coefficients β_i . A better approach is to use the Singular Value Decomposition (SVD) method discussed below. Or using Lasso and Ridge regression. See below.

Simple regression model

We are now ready to write our first program which aims at solving the above linear regression equations. We start with data we have produced ourselves, in this case normally distributed random numbers along the x -axis. These numbers define then the value of a function $y(x) = 4 + 3x + N(0, 1)$. Thereafter we order the x values and employ our linear regression algorithm to set up the best fit. Here we find it useful to use the numpy function `c_` arrays where arrays are stacked along their last axis after being upgraded to at least two dimensions with ones post-pended to the shape. The following examples help in understanding what happens

We see that, as expected, a linear fit gives a seemingly (from the graph) good representation of the data.

Simple regression model, now using scikit-learn

We can repeat the above algorithm using **scikit-learn** as follows

Simple linear regression model using scikit-learn

We start with perhaps our simplest possible example, using **scikit-learn** to perform linear regression analysis on a data set produced by us. What follows is a simple Python code where we have defined function y in terms of the variable x . Both are defined as vectors of dimension 1×100 . The entries to the vector \mathbf{x} are given by random numbers generated with a uniform distribution with entries $x_i \in [0, 1]$ (more about probability distribution functions later). These values are then used to define a function $y(x)$ (tabulated again as a vector) with a linear dependence on x plus a random noise added via the normal distribution.

The Numpy functions are imported used the **import numpy as np** statement and the random number generator for the uniform distribution is called using the function **np.random.rand()**, where we specify that we want 100 random variables. Using Numpy we define automatically an array with the specified number of elements, 100 in our case. With the Numpy function **randn()** we can compute random numbers with the normal distribution (mean value μ equal to zero and variance σ^2 set to one) and produce the values of y assuming a linear dependence as function of x

$$y = 2x + N(0, 1),$$

where $N(0, 1)$ represents random numbers generated by the normal distribution. From **scikit-learn** we import then the **LinearRegression** functionality and make a prediction $\tilde{y} = \alpha + \beta x$ using the function **fit(x,y)**. We call the set

of data (\mathbf{x}, \mathbf{y}) for our training data. The Python package **scikit-learn** has also a functionality which extracts the above fitting parameters α and β (see below). Later we will distinguish between training data and test data.

For plotting we use the Python package **matplotlib** which produces publication quality figures. Feel free to explore the extensive [gallery](#) of examples. In this example we plot our original values of x and y as well as the prediction **ypredict** (\tilde{y}), which attempts at fitting our data with a straight line.

The Python code follows here.

Simple linear regression model

This example serves several aims. It allows us to demonstrate several aspects of data analysis and later machine learning algorithms. The immediate visualization shows that our linear fit is not impressive. It goes through the data points, but there are many outliers which are not reproduced by our linear regression. We could now play around with this small program and change for example the factor in front of x and the normal distribution. Try to change the function y to

$$y = 10x + 0.01 \times N(0, 1),$$

where x is defined as before.

Less noise

Does the fit look better? Indeed, by reducing the role of the normal distribution we see immediately that our linear prediction seemingly reproduces better the training set. However, this testing 'by the eye' is obviously not satisfactory in the long run. Here we have only defined the training data and our model, and have not discussed a more rigorous approach to the **cost** function.

How to study our fits

We need more rigorous criteria in defining whether we have succeeded or not in modeling our training data. You will be surprised to see that many scientists seldomly venture beyond this 'by the eye' approach. A standard approach for the *cost* function is the so-called χ^2 function

$$\chi^2 = \frac{1}{n} \sum_{i=0}^{n-1} \frac{(y_i - \tilde{y}_i)^2}{\sigma_i^2},$$

where σ_i^2 is the variance (to be defined later) of the entry y_i . We may not know the explicit value of σ_i^2 , it serves however the aim of scaling the equations and make the cost function dimensionless.

Minimizing the cost function

Minimizing the cost function is a central aspect of our discussions to come. Finding its minima as function of the model parameters (α and β in our case) will be a recurring theme in these series of lectures. Essentially all machine learning algorithms we will discuss center around the minimization of the chosen cost function. This depends in turn on our specific model for describing the data, a typical situation in supervised learning. Automatizing the search for the minima of the cost function is a central ingredient in all algorithms. Typical methods which are employed are various variants of **gradient** methods. These will be discussed in more detail later. Again, you'll be surprised to hear that many practitioners minimize the above function "by the eye", popularly dubbed as 'chi by the eye'. That is, change a parameter and see (visually and numerically) that the χ^2 function becomes smaller.

Relative error

There are many ways to define the cost function. A simpler approach is to look at the relative difference between the training data and the predicted data, that is we define the relative error as

$$\epsilon_{\text{relative}} = \frac{|\mathbf{y} - \tilde{\mathbf{y}}|}{|\mathbf{y}|}.$$

We can modify easily the above Python code and plot the relative error instead

Depending on the parameter in front of the normal distribution, we may have a small or larger relative error. Try to play around with different training data sets and study (graphically) the value of the relative error.

The richness of scikit-learn

As mentioned above, **scikit-learn** has an impressive functionality. We can for example extract the values of α and β and their error estimates, or the variance and standard deviation and many other properties from the statistical data analysis.

Here we show an example of the functionality of scikit-learn.

Functions in scikit-learn

The function **coef** gives us the parameter β of our fit while **intercept** yields α . Depending on the constant in front of the normal distribution, we get values near or far from $\alpha = 2$ and $\beta = 5$. Try to play around with different parameters in front of the normal distribution. The function **meansquarederror** gives us the mean square error, a risk metric corresponding to the expected value of the squared (quadratic) error or loss defined as

$$MSE(\mathbf{y}, \tilde{\mathbf{y}}) = \frac{1}{n} \sum_{i=0}^{n-1} (y_i - \tilde{y}_i)^2,$$

The smaller the value, the better the fit. Ideally we would like to have an MSE equal zero. The attentive reader has probably recognized this function as being similar to the χ^2 function defined above.

Other functions in scikit-learn

The **r2score** function computes R^2 , the coefficient of determination. It provides a measure of how well future samples are likely to be predicted by the model. Best possible score is 1.0 and it can be negative (because the model can be arbitrarily worse). A constant model that always predicts the expected value of \mathbf{y} , disregarding the input features, would get a R^2 score of 0.0.

If \tilde{y}_i is the predicted value of the i -th sample and y_i is the corresponding true value, then the score R^2 is defined as

$$R^2(\mathbf{y}, \tilde{\mathbf{y}}) = 1 - \frac{\sum_{i=0}^{n-1} (y_i - \tilde{y}_i)^2}{\sum_{i=0}^{n-1} (y_i - \bar{y})^2},$$

where we have defined the mean value of \mathbf{y} as

$$\bar{y} = \frac{1}{n} \sum_{i=0}^{n-1} y_i.$$

The mean absolute error and other functions in scikit-learn

Another quantity will meet again in our discussions of regression analysis is mean absolute error (MAE), a risk metric corresponding to the expected value of the absolute error loss or what we call the l_1 -norm loss. In our discussion above we presented the relative error. The MAE is defined as follows

$$\text{MAE}(\mathbf{y}, \tilde{\mathbf{y}}) = \frac{1}{n} \sum_{i=0}^{n-1} |y_i - \tilde{y}_i|.$$

Finally we present the squared logarithmic (quadratic) error

$$\text{MSLE}(\mathbf{y}, \tilde{\mathbf{y}}) = \frac{1}{n} \sum_{i=0}^{n-1} (\log_e(1 + y_i) - \log_e(1 + \tilde{y}_i))^2,$$

where $\log_e(x)$ stands for the natural logarithm of x . This error estimate is best to use when targets having exponential growth, such as population counts, average sales of a commodity over a span of years etc.

Cubic polynomial in scikit-learn

We will discuss in more detail these and other functions in the various lectures. We conclude this part with another example. Instead of a linear x -dependence we study now a cubic polynomial and use the polynomial regression analysis tools of scikit-learn. Add description of the various python commands.

Polynomial Regression

Linking the regression analysis with a statistical interpretation

Before we proceed, and to link with our discussions of Bayesian statistics to come, it is useful to derive the standard regression analysis equations using a statistical interpretation. This allows us also to derive quantities like the variance and other expectation values in a rather straightforward way.

It is assumed that $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$ and the ε_i are independent, i.e.:

$$\text{Cov}(\varepsilon_{i_1}, \varepsilon_{i_2}) = \begin{cases} \sigma^2 & \text{if } i_1 = i_2, \\ 0 & \text{if } i_1 \neq i_2. \end{cases}$$

The randomness of ε_i implies that \mathbf{Y}_i is also a random variable. In particular, \mathbf{Y}_i is normally distributed, because $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$ and $\mathbf{X}_{i,*} \beta$ is a non-random scalar. To specify the parameters of the distribution of \mathbf{Y}_i we need to calculate its first two moments.

Expectation value and variance

Its expectation equals:

$$\mathbb{E}(Y_i) = \mathbb{E}(\mathbf{X}_{i,*} \beta) + \mathbb{E}(\varepsilon_i) = \mathbf{X}_{i,*} \beta,$$

while its variance is

$$\begin{aligned} \text{Var}(Y_i) &= \mathbb{E}\{[Y_i - \mathbb{E}(Y_i)]^2\} = \mathbb{E}(Y_i^2) - [\mathbb{E}(Y_i)]^2 \\ &= \mathbb{E}[(\mathbf{X}_{i,*} \beta + \varepsilon_i)^2] - (\mathbf{X}_{i,*} \beta)^2 \\ &= \mathbb{E}[(\mathbf{X}_{i,*} \beta)^2 + 2\varepsilon_i \mathbf{X}_{i,*} \beta + \varepsilon_i^2] - (\mathbf{X}_{i,*} \beta)^2 \\ &= (\mathbf{X}_{i,*} \beta)^2 + 2\mathbb{E}(\varepsilon_i) \mathbf{X}_{i,*} \beta + \mathbb{E}(\varepsilon_i^2) - (\mathbf{X}_{i,*} \beta)^2 \\ &= \mathbb{E}(\varepsilon_i^2) = \text{Var}(\varepsilon_i) = \sigma^2. \end{aligned}$$

Hence, $Y_i \sim \mathcal{N}(\mathbf{X}_{i,*} \beta, \sigma^2)$.

The singular value decomposition

A general $m \times n$ matrix \mathbf{A} can be written in terms of a diagonal matrix \mathbf{D} of dimensionality $n \times n$ and two orthogonal matrices \mathbf{U} and \mathbf{V} , where the first has dimensionality $m \times m$ and the last dimensionality $n \times n$. We have then

$$\mathbf{A} = \mathbf{U} \mathbf{D} \mathbf{V}^T$$

From standard regression to Ridge regressions

One of the typical problems we encounter with linear regression, in particular when the matrix \mathbf{X} (our so-called design matrix) is high-dimensional, are problems with near singular or singular matrices. The column vectors of \mathbf{X} may be

linearly dependent, normally referred to as super-collinearity. This means that the matrix may be rank deficient and it is basically impossible to model the data using linear regression. As an example, consider the matrix

$$\mathbf{X} = \begin{bmatrix} 1 & -1 & 2 \\ 1 & 0 & 1 \\ 1 & 2 & -1 \\ 1 & 1 & 0 \end{bmatrix}$$

The columns of \mathbf{X} are linearly dependent. We see this easily since the first column is the row-wise sum of the other two columns. The rank (more correct, the column rank) of a matrix is the dimension of the space spanned by the column vectors. Hence, the rank of \mathbf{X} is equal to the number of linearly independent columns. In this particular case the matrix has rank 2.

Super-collinearity of an $(n \times p)$ -dimensional design matrix \mathbf{X} implies that the inverse of the matrix $\mathbf{X}^T \mathbf{x}$ (the matrix we need to invert to solve the linear regression equations) is non-invertible. If we have a square matrix that does not have an inverse, we say this matrix singular. The example here demonstrates this

$$\mathbf{X} = \begin{bmatrix} 1 & -1 \\ 1 & -1 \end{bmatrix}.$$

We see easily that $\det(\mathbf{X}) = x_{11}x_{22} - x_{12}x_{21} = 1 \times (-1) - 1 \times (-1) = 0$. Hence, \mathbf{X} is singular and its inverse is undefined. This is equivalent to saying that the matrix \mathbf{X} has at least an eigenvalue which is zero.

Fixing the singularity

If our design matrix \mathbf{X} which enters the linear regression problem

$$\boldsymbol{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}, \quad (1)$$

has linearly dependent column vectors, we will not be able to compute the inverse of $\mathbf{X}^T \mathbf{X}$ and we cannot find the parameters (estimators) β_i . The estimators are only well-defined if $(\mathbf{X}^T \mathbf{X})^{-1}$ exists. This is more likely to happen when the matrix \mathbf{X} is high-dimensional. In this case it is likely to encounter a situation where the regression parameters β_i cannot be estimated.

The *ad hoc* approach which was introduced in the 70s was simply to add a diagonal component to the matrix to invert, that is we change

$$\mathbf{X}^T \mathbf{X} \rightarrow \mathbf{X}^T \mathbf{X} + \lambda \mathbf{I},$$

where \mathbf{I} is the identity matrix.

Fitting vs. predicting when data is in the model class

We start by considering the case $f(x) = 2x$.

Then the data is clearly generated by a model that is contained within all three model classes we are using to make predictions (linear models, third order polynomials, and tenth order polynomials).

Run the code for the following cases:

1. For $f(x) = 2x$, $Ntrain = 10$ and $\sigma = 0$ (noiseless case), train the three classes of models (linear, third-order polynomial, and tenth order polynomial) for a training set when $x \in [0, 1]$. Make graphs comparing fits for different order of polynomials. Which model fits the data the best?
2. Do you think that the data that has the least error on the training set will also make the best predictions? Why or why not? Can you try to discuss and formalize your intuition? What can go right and what can go wrong?
3. Check your answer by seeing how well your fits predict newly generated test data (including on data outside the range you fit on, for example $x \in [0, 1.2]$) using the code below. How well do you do on points in the range of x where you trained the model? How about points outside the original training data set?
4. Repeat the above for $f(x) = 2x$, $Ntrain = 10$, and $\sigma = 1$. What changes?

Repeat the exercises above for $f(x) = 2x$, $Ntrain = 100$, and $\sigma = 1$. What changes? Summarize what you have learned about the relationship between model complexity (number of parameters), goodness of fit on training data, and the ability to predict well.

Fitting versus predicting when data is not in the model class

Thus far, we have considered the case where the data is generated using a model contained in the model class. Now consider $f(x) = 2x - 10x^5 + 15x^{10}$. Notice that the for linear and third-order polynomial the true model $f(x)$ is not contained in model class.

1. Do better fits lead to better predictions?
2. What is the relationship between the true model for generating the data and the model class that has the most predictive power? How is this related to the model complexity? How does this depend on the number of data points $Ntrain$ and σ ?

Summarize what you think you learned about the relationship of knowing the true model class and predictive power.

An example code without the model assessment part

Generating test data

How can we effectively evaluate the various models?

In Ridge regression and the subsequent discussion of its properties the bias or penalty parameter is considered known or ‘given’. In practice, it is unknown and the user needs to make an informed decision on its value. How do we do that? Much of the same considerations apply to the Lasso method.

Code examples for Ridge and Lasso Regression

A second-order polynomial with Ridge and Lasso

Resampling methods

Resampling methods are an indispensable tool in modern statistics. They involve repeatedly drawing samples from a training set and refitting a model of interest on each sample in order to obtain additional information about the fitted model. For example, in order to estimate the variability of a linear regression fit, we can repeatedly draw different samples from the training data, fit a linear regression to each new sample, and then examine the extent to which the resulting fits differ. Such an approach may allow us to obtain information that would not be available from fitting the model only once using the original training sample.

Resampling approaches can be computationally expensive

Resampling approaches can be computationally expensive, because they involve fitting the same statistical method multiple times using different subsets of the training data. However, due to recent advances in computing power, the computational requirements of resampling methods generally are not prohibitive. In this chapter, we discuss two of the most commonly used resampling methods, cross-validation and the bootstrap. Both methods are important tools in the practical application of many statistical learning procedures. For example, cross-validation can be used to estimate the test error associated with a given statistical learning method in order to evaluate its performance, or to select the appropriate level of flexibility. The process of evaluating a model’s performance is known as model assessment, whereas the process of selecting the proper level of flexibility for a model is known as model selection. The bootstrap is widely used.

Why resampling methods ?

Statistical analysis.

- Our simulations can be treated as *computer experiments*. This is particularly the case for Monte Carlo methods

- The results can be analysed with the same statistical tools as we would use analysing experimental data.
- As in all experiments, we are looking for expectation values and an estimate of how accurate they are, i.e., possible sources for errors.

Statistical analysis

- As in other experiments, many numerical experiments have two classes of errors:
 - Statistical errors
 - Systematical errors
- Statistical errors can be estimated using standard tools from statistics
- Systematical errors are method specific and must be treated differently from case to case.

Statistics

The *probability distribution function (PDF)* is a function $p(x)$ on the domain which, in the discrete case, gives us the probability or relative frequency with which these values of X occur:

$$p(x) = \text{prob}(X = x)$$

In the continuous case, the PDF does not directly depict the actual probability. Instead we define the probability for the stochastic variable to assume any value on an infinitesimal interval around x to be $p(x)dx$. The continuous function $p(x)$ then gives us the *density* of the probability rather than the probability itself. The probability for a stochastic variable to assume any value on a non-infinitesimal interval $[a, b]$ is then just the integral:

$$\text{prob}(a \leq X \leq b) = \int_a^b p(x)dx$$

Qualitatively speaking, a stochastic variable represents the values of numbers chosen as if by chance from some specified PDF so that the selection of a large set of these numbers reproduces this PDF.

Statistics, moments

A particularly useful class of special expectation values are the *moments*. The n -th moment of the PDF p is defined as follows:

$$\langle x^n \rangle \equiv \int x^n p(x) dx$$

The zero-th moment $\langle 1 \rangle$ is just the normalization condition of p . The first moment, $\langle x \rangle$, is called the *mean* of p and often denoted by the letter μ :

$$\langle x \rangle = \mu \equiv \int x p(x) dx$$

Statistics, central moments

A special version of the moments is the set of *central moments*, the n -th central moment defined as:

$$\langle (x - \langle x \rangle)^n \rangle \equiv \int (x - \langle x \rangle)^n p(x) dx$$

The zero-th and first central moments are both trivial, equal 1 and 0, respectively. But the second central moment, known as the *variance* of p , is of particular interest. For the stochastic variable X , the variance is denoted as σ_X^2 or $\text{var}(X)$:

$$\sigma_X^2 = \text{var}(X) = \langle (x - \langle x \rangle)^2 \rangle = \int (x - \langle x \rangle)^2 p(x) dx \quad (2)$$

$$= \int (x^2 - 2x\langle x \rangle + \langle x \rangle^2) p(x) dx \quad (3)$$

$$= \langle x^2 \rangle - 2\langle x \rangle \langle x \rangle + \langle x \rangle^2 \quad (4)$$

$$= \langle x^2 \rangle - \langle x \rangle^2 \quad (5)$$

The square root of the variance, $\sigma = \sqrt{\langle (x - \langle x \rangle)^2 \rangle}$ is called the *standard deviation* of p . It is clearly just the RMS (root-mean-square) value of the deviation of the PDF from its mean value, interpreted qualitatively as the *spread* of p around its mean.

Statistics, covariance

Another important quantity is the so called covariance, a variant of the above defined variance. Consider again the set $\{X_i\}$ of n stochastic variables (not necessarily uncorrelated) with the multivariate PDF $P(x_1, \dots, x_n)$. The *covariance* of two of the stochastic variables, X_i and X_j , is defined as follows:

$$\begin{aligned} \text{cov}(X_i, X_j) &\equiv \langle (x_i - \langle x_i \rangle)(x_j - \langle x_j \rangle) \rangle \\ &= \int \dots \int (x_i - \langle x_i \rangle)(x_j - \langle x_j \rangle) P(x_1, \dots, x_n) dx_1 \dots dx_n \end{aligned} \quad (6)$$

with

$$\langle x_i \rangle = \int \dots \int x_i P(x_1, \dots, x_n) dx_1 \dots dx_n$$

Statistics, more covariance

If we consider the above covariance as a matrix $C_{ij} = \text{cov}(X_i, X_j)$, then the diagonal elements are just the familiar variances, $C_{ii} = \text{cov}(X_i, X_i) = \text{var}(X_i)$. It turns out that all the off-diagonal elements are zero if the stochastic variables are uncorrelated. This is easy to show, keeping in mind the linearity of the expectation value. Consider the stochastic variables X_i and X_j , ($i \neq j$):

$$\text{cov}(X_i, X_j) = \langle (x_i - \langle x_i \rangle)(x_j - \langle x_j \rangle) \rangle \quad (7)$$

$$= \langle x_i x_j - x_i \langle x_j \rangle - \langle x_i \rangle x_j + \langle x_i \rangle \langle x_j \rangle \rangle \quad (8)$$

$$= \langle x_i x_j \rangle - \langle x_i \rangle \langle x_j \rangle - \langle \langle x_i \rangle x_j \rangle + \langle \langle x_i \rangle \langle x_j \rangle \rangle \quad (9)$$

$$= \langle x_i x_j \rangle - \langle x_i \rangle \langle x_j \rangle - \langle x_i \rangle \langle x_j \rangle + \langle x_i \rangle \langle x_j \rangle \quad (10)$$

$$= \langle x_i x_j \rangle - \langle x_i \rangle \langle x_j \rangle \quad (11)$$

Statistics, independent variables

If X_i and X_j are independent, we get $\langle x_i x_j \rangle = \langle x_i \rangle \langle x_j \rangle$, resulting in $\text{cov}(X_i, X_j) = 0$ ($i \neq j$).

Also useful for us is the covariance of linear combinations of stochastic variables. Let $\{X_i\}$ and $\{Y_i\}$ be two sets of stochastic variables. Let also $\{a_i\}$ and $\{b_i\}$ be two sets of scalars. Consider the linear combination:

$$U = \sum_i a_i X_i \quad V = \sum_j b_j Y_j$$

By the linearity of the expectation value

$$\text{cov}(U, V) = \sum_{i,j} a_i b_j \text{cov}(X_i, Y_j)$$

Statistics, more variance

Now, since the variance is just $\text{var}(X_i) = \text{cov}(X_i, X_i)$, we get the variance of the linear combination $U = \sum_i a_i X_i$:

$$\text{var}(U) = \sum_{i,j} a_i a_j \text{cov}(X_i, X_j) \quad (12)$$

And in the special case when the stochastic variables are uncorrelated, the off-diagonal elements of the covariance are as we know zero, resulting in:

$$\text{var}(U) = \sum_i a_i^2 \text{cov}(X_i, X_i) = \sum_i a_i^2 \text{var}(X_i)$$

$$\text{var}\left(\sum_i a_i X_i\right) = \sum_i a_i^2 \text{var}(X_i)$$

which will become very useful in our study of the error in the mean value of a set of measurements.

Statistics and stochastic processes

A *stochastic process* is a process that produces sequentially a chain of values:

$$\{x_1, x_2, \dots, x_k, \dots\}.$$

We will call these values our *measurements* and the entire set as our measured *sample*. The action of measuring all the elements of a sample we will call a stochastic *experiment* since, operationally, they are often associated with results of empirical observation of some physical or mathematical phenomena; precisely an experiment. We assume that these values are distributed according to some PDF $p_X(x)$, where X is just the formal symbol for the stochastic variable whose PDF is $p_X(x)$. Instead of trying to determine the full distribution p we are often only interested in finding the few lowest moments, like the mean μ_X and the variance σ_X .

Statistics and sample variables

In practical situations a sample is always of finite size. Let that size be n . The expectation value of a sample, the *sample mean*, is then defined as follows:

$$\bar{x}_n \equiv \frac{1}{n} \sum_{k=1}^n x_k$$

The *sample variance* is:

$$\text{var}(x) \equiv \frac{1}{n} \sum_{k=1}^n (x_k - \bar{x}_n)^2$$

its square root being the *standard deviation of the sample*. The *sample covariance* is:

$$\text{cov}(x) \equiv \frac{1}{n} \sum_{kl} (x_k - \bar{x}_n)(x_l - \bar{x}_n)$$

Statistics, sample variance and covariance

Note that the sample variance is the sample covariance without the cross terms. In a similar manner as the covariance in Eq. (6) is a measure of the correlation between two stochastic variables, the above defined sample covariance is a measure of the sequential correlation between succeeding measurements of a sample.

These quantities, being known experimental values, differ significantly from and must not be confused with the similarly named quantities for stochastic variables, mean μ_X , variance $\text{var}(X)$ and covariance $\text{cov}(X, Y)$.

Statistics, law of large numbers

The law of large numbers states that as the size of our sample grows to infinity, the sample mean approaches the true mean μ_X of the chosen PDF:

$$\lim_{n \rightarrow \infty} \bar{x}_n = \mu_X$$

The sample mean \bar{x}_n works therefore as an estimate of the true mean μ_X .

What we need to find out is how good an approximation \bar{x}_n is to μ_X . In any stochastic measurement, an estimated mean is of no use to us without a measure of its error. A quantity that tells us how well we can reproduce it in another experiment. We are therefore interested in the PDF of the sample mean itself. Its standard deviation will be a measure of the spread of sample means, and we will simply call it the *error* of the sample mean, or just sample error, and denote it by err_X . In practice, we will only be able to produce an *estimate* of the sample error since the exact value would require the knowledge of the true PDFs behind, which we usually do not have.

Statistics, more on sample error

Let us first take a look at what happens to the sample error as the size of the sample grows. In a sample, each of the measurements x_i can be associated with its own stochastic variable X_i . The stochastic variable \bar{X}_n for the sample mean \bar{x}_n is then just a linear combination, already familiar to us:

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$$

All the coefficients are just equal $1/n$. The PDF of \bar{X}_n , denoted by $p_{\bar{X}_n}(x)$ is the desired PDF of the sample means.

Statistics

The probability density of obtaining a sample mean \bar{x}_n is the product of probabilities of obtaining arbitrary values x_1, x_2, \dots, x_n with the constraint that the mean of the set $\{x_i\}$ is \bar{x}_n :

$$p_{\bar{X}_n}(x) = \int p_X(x_1) \cdots \int p_X(x_n) \delta\left(x - \frac{x_1 + x_2 + \cdots + x_n}{n}\right) dx_n \cdots dx_1$$

And in particular we are interested in its variance $\text{var}(\bar{X}_n)$.

Statistics, central limit theorem

It is generally not possible to express $p_{\bar{X}_n}(x)$ in a closed form given an arbitrary PDF p_X and a number n . But for the limit $n \rightarrow \infty$ it is possible to make an approximation. The very important result is called *the central limit*

theorem. It tells us that as n goes to infinity, $p_{\bar{X}_n}(x)$ approaches a Gaussian distribution whose mean and variance equal the true mean and variance, μ_X and σ_X^2 , respectively:

$$\lim_{n \rightarrow \infty} p_{\bar{X}_n}(x) = \left(\frac{n}{2\pi \text{var}(X)} \right)^{1/2} e^{-\frac{n(x - \bar{x}_n)^2}{2\text{var}(X)}} \quad (13)$$

Statistics, more technicalities

The desired variance $\text{var}(\bar{X}_n)$, i.e. the sample error squared err_X^2 , is given by:

$$\text{err}_X^2 = \text{var}(\bar{X}_n) = \frac{1}{n^2} \sum_{ij} \text{cov}(X_i, X_j) \quad (14)$$

We see now that in order to calculate the exact error of the sample with the above expression, we would need the true means μ_{X_i} of the stochastic variables X_i . To calculate these requires that we know the true multivariate PDF of all the X_i . But this PDF is unknown to us, we have only got the measurements of one sample. The best we can do is to let the sample itself be an estimate of the PDF of each of the X_i , estimating all properties of X_i through the measurements of the sample.

Statistics

Our estimate of μ_{X_i} is then the sample mean \bar{x} itself, in accordance with the central limit theorem:

$$\mu_{X_i} = \langle x_i \rangle \approx \frac{1}{n} \sum_{k=1}^n x_k = \bar{x}$$

Using \bar{x} in place of μ_{X_i} we can give an *estimate* of the covariance in Eq. (14)

$$\text{cov}(X_i, X_j) = \langle (x_i - \langle x_i \rangle)(x_j - \langle x_j \rangle) \rangle \approx \langle (x_i - \bar{x})(x_j - \bar{x}) \rangle,$$

resulting in

$$\frac{1}{n} \sum_l \left(\frac{1}{n} \sum_k (x_k - \bar{x}_n)(x_l - \bar{x}_n) \right) = \frac{1}{n} \frac{1}{n} \sum_{kl} (x_k - \bar{x}_n)(x_l - \bar{x}_n) = \frac{1}{n} \text{cov}(x)$$

Statistics and sample variance

By the same procedure we can use the sample variance as an estimate of the variance of any of the stochastic variables X_i

$$\text{var}(X_i) = \langle x_i - \langle x_i \rangle \rangle \approx \langle x_i - \bar{x}_n \rangle,$$

which is approximated as

$$\text{var}(X_i) \approx \frac{1}{n} \sum_{k=1}^n (x_k - \bar{x}_n) = \text{var}(x) \quad (15)$$

Now we can calculate an estimate of the error err_X of the sample mean \bar{x}_n :

$$\begin{aligned} \text{err}_X^2 &= \frac{1}{n^2} \sum_{ij} \text{cov}(X_i, X_j) \\ &\approx \frac{1}{n^2} \sum_{ij} \frac{1}{n} \text{cov}(x) = \frac{1}{n^2} n^2 \frac{1}{n} \text{cov}(x) \\ &= \frac{1}{n} \text{cov}(x) \end{aligned} \quad (16)$$

which is nothing but the sample covariance divided by the number of measurements in the sample.

Statistics, uncorrelated results

In the special case that the measurements of the sample are uncorrelated (equivalently the stochastic variables X_i are uncorrelated) we have that the off-diagonal elements of the covariance are zero. This gives the following estimate of the sample error:

$$\text{err}_X^2 = \frac{1}{n^2} \sum_{ij} \text{cov}(X_i, X_j) = \frac{1}{n^2} \sum_i \text{var}(X_i),$$

resulting in

$$\text{err}_X^2 \approx \frac{1}{n^2} \sum_i \text{var}(x) = \frac{1}{n} \text{var}(x) \quad (17)$$

where in the second step we have used Eq. (15). The error of the sample is then just its standard deviation divided by the square root of the number of measurements the sample contains. This is a very useful formula which is easy to compute. It acts as a first approximation to the error, but in numerical experiments, we cannot overlook the always present correlations.

Statistics, computations

For computational purposes one usually splits up the estimate of err_X^2 , given by Eq. (16), into two parts

$$\text{err}_X^2 = \frac{1}{n} \text{var}(x) + \frac{1}{n} (\text{cov}(x) - \text{var}(x)),$$

which equals

$$\frac{1}{n^2} \sum_{k=1}^n (x_k - \bar{x}_n)^2 + \frac{2}{n^2} \sum_{k < l} (x_k - \bar{x}_n)(x_l - \bar{x}_n) \quad (18)$$

The first term is the same as the error in the uncorrelated case, Eq. (17). This means that the second term accounts for the error correction due to correlation between the measurements. For uncorrelated measurements this second term is zero.

Statistics, more on computations of errors

Computationally the uncorrelated first term is much easier to treat efficiently than the second.

$$\text{var}(x) = \frac{1}{n} \sum_{k=1}^n (x_k - \bar{x}_n)^2 = \left(\frac{1}{n} \sum_{k=1}^n x_k^2 \right) - \bar{x}_n^2$$

We just accumulate separately the values x^2 and x for every measurement x we receive. The correlation term, though, has to be calculated at the end of the experiment since we need all the measurements to calculate the cross terms. Therefore, all measurements have to be stored throughout the experiment.

Statistics, wrapping up 1

Let us analyze the problem by splitting up the correlation term into partial sums of the form:

$$f_d = \frac{1}{n-d} \sum_{k=1}^{n-d} (x_k - \bar{x}_n)(x_{k+d} - \bar{x}_n)$$

The correlation term of the error can now be rewritten in terms of f_d

$$\frac{2}{n} \sum_{k < l} (x_k - \bar{x}_n)(x_l - \bar{x}_n) = 2 \sum_{d=1}^{n-1} f_d$$

The value of f_d reflects the correlation between measurements separated by the distance d in the sample samples. Notice that for $d = 0$, f is just the sample variance, $\text{var}(x)$. If we divide f_d by $\text{var}(x)$, we arrive at the so called *autocorrelation function*

$$\kappa_d = \frac{f_d}{\text{var}(x)}$$

which gives us a useful measure of pairwise correlations starting always at 1 for $d = 0$.

Statistics, final expression

The sample error (see eq. (18)) can now be written in terms of the autocorrelation function:

$$\begin{aligned}
 \text{err}_X^2 &= \frac{1}{n} \text{var}(x) + \frac{2}{n} \cdot \text{var}(x) \sum_{d=1}^{n-1} \frac{f_d}{\text{var}(x)} \\
 &= \left(1 + 2 \sum_{d=1}^{n-1} \kappa_d \right) \frac{1}{n} \text{var}(x) \\
 &= \frac{\tau}{n} \cdot \text{var}(x)
 \end{aligned} \tag{19}$$

and we see that err_X can be expressed in terms the uncorrelated sample variance times a correction factor τ which accounts for the correlation between measurements. We call this correction factor the *autocorrelation time*:

$$\tau = 1 + 2 \sum_{d=1}^{n-1} \kappa_d \tag{20}$$

Statistics, effective number of correlations

For a correlation free experiment, τ equals 1. From the point of view of eq. (19) we can interpret a sequential correlation as an effective reduction of the number of measurements by a factor τ . The effective number of measurements becomes:

$$n_{\text{eff}} = \frac{n}{\tau}$$

To neglect the autocorrelation time τ will always cause our simple uncorrelated estimate of $\text{err}_X^2 \approx \text{var}(x)/n$ to be less than the true sample error. The estimate of the error will be too *good*. On the other hand, the calculation of the full autocorrelation time poses an efficiency problem if the set of measurements is very large.

Log-likelihood

A popular strategy is to choose a penalty parameter that yields a good but parsimonious model. Information criteria measure the balance between model fit and model complexity. One possibility is Aikake's information criterion (AIC). The AIC measures model fit by the log-likelihood and model complexity is measured by the number of parameters used by the model. The number of model parameters in regular regression simply corresponds to the number of covariates in the model. Or, by the degrees of freedom consumed by the model, which is equivalent to the trace of the hat matrix. For ridge regression it thus seems natural to define model complexity analogously by the trace of the ridge hat matrix. This yields the AIC for the linear regression model with ridge estimates:

$$\begin{aligned}
\text{AIC}(\lambda) &= 2p - 2\log(\mathbf{L}) \\
&= 2\text{tr}[\mathbf{H}(\lambda)] - 2\log\{L[\boldsymbol{\beta}(\lambda), \boldsymbol{\sigma}^2(\lambda)]\} \\
&= 2 \sum_{j=1}^p \frac{d_{jj}^2}{d_{jj}^2 + \lambda} + 2n \log[\sqrt{2\pi} \boldsymbol{\sigma}(\lambda)] + \frac{1}{\boldsymbol{\sigma}^2(\lambda)} \sum_{i=1}^n [y_i - \mathbf{X}_{i,*} \boldsymbol{\beta}(\lambda)]^2.
\end{aligned}$$

The value of λ which minimizes $\text{AIC}(\lambda)$ corresponds to the ‘optimal’ balance of model complexity and overfitting.

Cross-validation

Instead of choosing the penalty parameter to balance model fit with model complexity, cross-validation requires it (i.e. the penalty parameter) to yield a model with good prediction performance. Commonly, this performance is evaluated on novel data. Novel data need not be easy to come by and one has to make do with the data at hand. The setting of ‘original’ and novel data is then mimicked by sample splitting: the data set is divided into two (groups of samples). One of these two data sets, called the *training set*, plays the role of ‘original’ data on which the model is built. The second of these data sets, called the *test set*, plays the role of the ‘novel’ data and is used to evaluate the prediction performance (often operationalized as the log-likelihood or the prediction error or its square or the R2 score) of the model built on the training data set. This procedure (model building and prediction evaluation on training and test set, respectively) is done for a collection of possible penalty parameter choices. The penalty parameter that yields the model with the best prediction performance is to be preferred. The thus obtained performance evaluation depends on the actual split of the data set. To remove this dependence the data set is split many times into a training and test set. For each split the model parameters are estimated for all choices of λ using the training data and estimated parameters are evaluated on the corresponding test set. The penalty parameter that on average over the test sets performs best (in some sense) is then selected.

Computationally expensive

The validation set approach is conceptually simple and is easy to implement. But it has two potential drawbacks:

- The validation estimate of the test error rate can be highly variable, depending on precisely which observations are included in the training set and which observations are included in the validation set.
- In the validation approach, only a subset of the observations, those that are included in the training set rather than in the validation set are used to fit the model. Since statistical methods tend to perform worse when trained on fewer observations, this suggests that the validation set error

rate may tend to overestimate the test error rate for the model fit on the entire data set.

Various steps in cross-validation

When the repetitive splitting of the data set is done randomly, samples may accidentally end up in a fast majority of the splits in either training or test set. Such samples may have an unbalanced influence on either model building or prediction evaluation. To avoid this k -fold cross-validation structures the data splitting. The samples are divided into k more or less equally sized exhaustive and mutually exclusive subsets. In turn (at each split) one of these subsets plays the role of the test set while the union of the remaining subsets constitutes the training set. Such a splitting warrants a balanced representation of each sample in both training and test set over the splits. Still the division into the k subsets involves a degree of randomness. This may be fully excluded when choosing $k = n$. This particular case is referred to as leave-one-out cross-validation (LOOCV).

How to set up the cross-validation for Ridge and/or Lasso

- Define a range of interest for the penalty parameter.
- Divide the data set into training and test set comprising samples $\{1, \dots, n\} \setminus i$ and $\{i\}$, respectively.
- Fit the linear regression model by means of ridge estimation for each λ in the grid using the training set, and the corresponding estimate of the error variance $\sigma_{-i}^2(\lambda)$, as

$$\beta_{-i}(\lambda) = (\mathbf{X}_{-i,*}^\top \mathbf{X}_{-i,*} + \lambda \mathbf{I}_{pp})^{-1} \mathbf{X}_{-i,*}^\top \mathbf{y}_{-i}$$

- Evaluate the prediction performance of these models on the test set by $\log\{L[y_i, \mathbf{X}_{i,*}; \beta_{-i}(\lambda), \sigma_{-i}^2(\lambda)]\}$. Or, by the prediction error $|y_i - \mathbf{X}_{i,*} \beta_{-i}(\lambda)|$, the relative error, the error squared or the R2 score function.
- Repeat the first three steps such that each sample plays the role of the test set once.
- Average the prediction performances of the test sets at each grid point of the penalty bias/parameter by computing the *cross-validated log-likelihood*. It is an estimate of the prediction performance of the model corresponding to this value of the penalty parameter on novel data. It is defined as

$$\frac{1}{n} \sum_{i=1}^n \log\{L[y_i, \mathbf{X}_{i,*}; \beta_{-i}(\lambda), \sigma_{-i}^2(\lambda)]\}.$$

- The value of the penalty parameter that maximizes the cross-validated log-likelihood is the value of choice. Or we can use the MSE or the R2 score functions.

Predicted Residual Error Sum of Squares

Another approach in the LOOCV scheme is to use the so-called Predicted Residual Error Sum of Squares (PRESS).

We can define the optimal penalty parameter to minimize

$$\lambda_{\text{opt}} = \arg \min_{\lambda} \frac{1}{n} \sum_{i=1}^n [y_i - \mathbf{X}_{i,*} \boldsymbol{\beta}_{-i}(\lambda)]^2.$$

The LOOCV prediction performance can be expressed analytically in terms of the known quantities derived from the design matrix and the parameters $\boldsymbol{\beta}$.

Resampling methods: Jackknife and Bootstrap

Two famous resampling methods are the **independent bootstrap** and the **jackknife**.

The jackknife is a special case of the independent bootstrap. Still, the jackknife was made popular prior to the independent bootstrap. And as the popularity of the independent bootstrap soared, new variants, such as the **dependent bootstrap**.

The Jackknife and independent bootstrap work for independent, identically distributed random variables. If these conditions are not satisfied, the methods will fail. Yet, it should be said that if the data are independent, identically distributed, and we only want to estimate the variance of \bar{X} (which often is the case), then there is no need for bootstrapping.

Resampling methods: Jackknife

The Jackknife works by making many replicas of the estimator $\hat{\theta}$. The jackknife is a resampling method, we explained that this happens by scrambling the data in some way. When using the jackknife, this is done by systematically leaving out one observation from the vector of observed values $\mathbf{x} = (x_1, x_2, \dots, x_n)$. Let \mathbf{x}_i denote the vector

$$\mathbf{x}_i = (x_1, x_2, \dots, x_{i-1}, x_{i+1}, \dots, x_n),$$

which equals the vector \mathbf{x} with the exception that observation number i is left out. Using this notation, define $\hat{\theta}_i$ to be the estimator $\hat{\theta}$ computed using \mathbf{x}_i .

Resampling methods: Jackknife estimator

To get an estimate for the bias and standard error of $\hat{\theta}$, use the following estimators for each component of $\hat{\theta}$

$$\widehat{\text{Bias}}(\hat{\theta}, \theta) = (n-1) \left(-\hat{\theta} + \frac{1}{n} \sum_{i=1}^n \hat{\theta}_i \right) \quad \text{and} \quad \hat{\sigma}_{\hat{\theta}}^2 = \frac{n-1}{n} \sum_{i=1}^n \left(\hat{\theta}_i - \frac{1}{n} \sum_{j=1}^n \hat{\theta}_j \right)^2.$$

Jackknife code example

Resampling methods: Bootstrap

Bootstrapping is a nonparametric approach to statistical inference that substitutes computation for more traditional distributional assumptions and asymptotic results. Bootstrapping offers a number of advantages:

1. The bootstrap is quite general, although there are some cases in which it fails.
2. Because it does not require distributional assumptions (such as normally distributed errors), the bootstrap can provide more accurate inferences when the data are not well behaved or when the sample size is small.
3. It is possible to apply the bootstrap to statistics with sampling distributions that are difficult to derive, even asymptotically.
4. It is relatively simple to apply the bootstrap to complex data-collection plans (such as stratified and clustered samples).

Resampling methods: Bootstrap background

Since $\hat{\theta} = \hat{\theta}(\mathbf{X})$ is a function of random variables, $\hat{\theta}$ itself must be a random variable. Thus it has a pdf, call this function $p(\mathbf{t})$. The aim of the bootstrap is to estimate $p(\mathbf{t})$ by the relative frequency of $\hat{\theta}$. You can think of this as using a histogram in the place of $p(\mathbf{t})$. If the relative frequency closely resembles $p(\mathbf{t})$, then using numerics, it is straight forward to estimate all the interesting parameters of $p(\mathbf{t})$ using point estimators.

Resampling methods: More Bootstrap background

In the case that $\hat{\theta}$ has more than one component, and the components are independent, we use the same estimator on each component separately. If the probability density function of X_i , $p(x)$, had been known, then it would have been straight forward to do this by:

1. Drawing lots of numbers from $p(x)$, suppose we call one such set of numbers $(X_1^*, X_2^*, \dots, X_n^*)$.

2. Then using these numbers, we could compute a replica of $\hat{\theta}$ called $\hat{\theta}^*$.

By repeated use of (1) and (2), many estimates of $\hat{\theta}$ could have been obtained. The idea is to use the relative frequency of $\hat{\theta}^*$ (think of a histogram) as an estimate of $p(\mathbf{t})$.

Resampling methods: Bootstrap approach

But unless there is enough information available about the process that generated X_1, X_2, \dots, X_n , $p(x)$ is in general unknown. Therefore, [Efron in 1979](#) asked the question: What if we replace $p(x)$ by the relative frequency of the observation X_i ; if we draw observations in accordance with the relative frequency of the observations, will we obtain the same result in some asymptotic sense? The answer is yes.

Instead of generating the histogram for the relative frequency of the observation X_i , just draw the values $(X_1^*, X_2^*, \dots, X_n^*)$ with replacement from the vector \mathbf{X} .

Resampling methods: Bootstrap steps

The independent bootstrap works like this:

1. Draw with replacement n numbers for the observed variables $\mathbf{x} = (x_1, x_2, \dots, x_n)$.
2. Define a vector \mathbf{x}^* containing the values which were drawn from \mathbf{x} .
3. Using the vector \mathbf{x}^* compute $\hat{\theta}^*$ by evaluating $\hat{\theta}$ under the observations \mathbf{x}^* .
4. Repeat this process k times.

When you are done, you can draw a histogram of the relative frequency of $\hat{\theta}^*$. This is your estimate of the probability distribution $p(t)$. Using this probability distribution you can estimate any statistics thereof. In principle you never draw the histogram of the relative frequency of $\hat{\theta}^*$. Instead you use the estimators corresponding to the statistic of interest. For example, if you are interested in estimating the variance of $\hat{\theta}$, apply the estimator $\hat{\sigma}^2$ to the values $\hat{\theta}^*$.

Code example for the Bootstrap method

The following code starts with a Gaussian distribution with mean value $\mu = 100$ and variance $\sigma = 15$. We use this to generate the data used in the bootstrap analysis. The bootstrap analysis returns a data set after a given number of bootstrap operations (as many as we have data points). This data set consists of estimated mean values for each bootstrap operation. The histogram generated by the bootstrap method shows that the distribution for these mean values is also a Gaussian, centered around the mean value $\mu = 100$ but with standard deviation σ/\sqrt{n} , where n is the number of bootstrap samples (in this case the

same as the number of original data points). The value of the standard deviation is what we expect from the central limit theorem.

Resampling methods: Blocking

The blocking method was made popular by [Flyvbjerg and Pedersen \(1989\)](#) and has become one of the standard ways to estimate $V(\hat{\theta})$ for exactly one $\hat{\theta}$, namely $\hat{\theta} = \bar{X}$.

Assume $n = 2^d$ for some integer $d > 1$ and X_1, X_2, \dots, X_n is a stationary time series to begin with. Moreover, assume that the time series is asymptotically uncorrelated. We switch to vector notation by arranging X_1, X_2, \dots, X_n in an n -tuple. Define:

$$\mathbf{X} = (X_1, X_2, \dots, X_n).$$

The strength of the blocking method is when the number of observations, n is large. For large n , the complexity of dependent bootstrapping scales poorly, but the blocking method does not, moreover, it becomes more accurate the larger n is.

Blocking Transformations

We now define blocking transformations. The idea is to take the mean of subsequent pair of elements from \vec{X} and form a new vector \vec{X}_1 . Continuing in the same way by taking the mean of subsequent pairs of elements of \vec{X}_1 we obtain \vec{X}_2 , and so on. Define \vec{X}_i recursively by:

$$\begin{aligned} (\vec{X}_0)_k &\equiv (\vec{X})_k \\ (\vec{X}_{i+1})_k &\equiv \frac{1}{2} \left((\vec{X}_i)_{2k-1} + (\vec{X}_i)_{2k} \right) \quad \text{for all } 1 \leq i \leq d-1 \end{aligned} \quad (21)$$

The quantity \vec{X}_k is subject to k **blocking transformations**. We now have d vectors $\vec{X}_0, \vec{X}_1, \dots, \vec{X}_{d-1}$ containing the subsequent averages of observations. It turns out that if the components of \vec{X} is a stationary time series, then the components of \vec{X}_i is a stationary time series for all $0 \leq i \leq d-1$.

We can then compute the autocovariance, the variance, sample mean, and number of observations for each i . Let $\gamma_i, \sigma_i^2, \bar{X}_i$ denote the autocovariance, variance and average of the elements of \vec{X}_i and let n_i be the number of elements of \vec{X}_i . It follows by induction that $n_i = n/2^i$.

Blocking Transformations

Using the definition of the blocking transformation and the distributive property of the covariance, it is clear that since $h = |i - j|$ we can define

$$\begin{aligned}\gamma_{k+1}(h) &= \text{cov}((X_{k+1})_i, (X_{k+1})_j) \\ &= \frac{1}{4} \text{cov}((X_k)_{2i-1} + (X_k)_{2i}, (X_k)_{2j-1} + (X_k)_{2j}) \\ &= \frac{1}{2} \gamma_k(2h) + \frac{1}{2} \gamma_k(2h+1) \quad h = 0 \\ &= \frac{1}{4} \gamma_k(2h-1) + \frac{1}{2} \gamma_k(2h) + \frac{1}{4} \gamma_k(2h+1) \quad \text{else}\end{aligned} \quad (22)$$

The quantity \mathbf{X} is asymptotic uncorrelated by assumption, \mathbf{X}_k is also asymptotic uncorrelated. Let's turn our attention to the variance of the sample mean $V(\bar{X})$.

Blocking Transformations, getting there

We have

$$V(\bar{X}_k) = \frac{\sigma_k^2}{n_k} + \underbrace{\frac{2}{n_k} \sum_{h=1}^{n_k-1} \left(1 - \frac{h}{n_k}\right) \gamma_k(h)}_{\equiv e_k} = \frac{\sigma_k^2}{n_k} + e_k \quad \text{if } \gamma_k(0) = \sigma_k^2. \quad (24)$$

The term e_k is called the **truncation error**:

$$e_k = \frac{2}{n_k} \sum_{h=1}^{n_k-1} \left(1 - \frac{h}{n_k}\right) \gamma_k(h). \quad (25)$$

We can show that $V(\bar{X}_i) = V(\bar{X}_j)$ for all $0 \leq i \leq d-1$ and $0 \leq j \leq d-1$.

Blocking Transformations, final expressions

We can then wrap up

$$\begin{aligned}n_{j+1} \bar{X}_{j+1} &= \sum_{i=1}^{n_{j+1}} (\mathbf{X}_{j+1})_i = \frac{1}{2} \sum_{i=1}^{n_j/2} (\mathbf{X}_j)_{2i-1} + (\mathbf{X}_j)_{2i} \\ &= \frac{1}{2} [(\mathbf{X}_j)_1 + (\mathbf{X}_j)_2 + \cdots + (\mathbf{X}_j)_{n_j}] = \underbrace{\frac{n_j}{2}}_{=n_{j+1}} \bar{X}_j = n_{j+1} \bar{X}_j.\end{aligned} \quad (26)$$

By repeated use of this equation we get $V(\bar{X}_i) = V(\bar{X}_0) = V(\bar{X})$ for all $0 \leq i \leq d-1$. This has the consequence that

$$V(\bar{X}) = \frac{\sigma_k^2}{n_k} + e_k \quad \text{for all } 0 \leq k \leq d-1. \quad (27)$$

Fyvbjerg and Petersen demonstrated that the sequence $\{e_k\}_{k=0}^{d-1}$ is decreasing, and conjecture that the term e_k can be made as small as we would like by making k (and hence d) sufficiently large. The sequence is decreasing (Master of Science thesis by Marius Jonsson, UiO 2018). It means we can apply blocking transformations until e_k is sufficiently small, and then estimate $V(\bar{X})$ by $\hat{\sigma}_k^2/n_k$.

Code examples for Blocking, Jackknife and bootstrap

The bias-variance tradeoff

We begin with an unknown function $y = f(x)$ and fix a *hypothesis set* \mathcal{H} consisting of all functions we are willing to consider, defined also on the domain of f . This set may be uncountably infinite (e.g. if there are real-valued parameters to fit). The choice of which functions to include in \mathcal{H} usually depends on our intuition about the problem of interest. The function $f(x)$ produces a set of pairs (x_i, y_i) , $i = 1 \dots N$, which serve as the observable data. Our goal is to select a function from the hypothesis set $h \in \mathcal{H}$ which approximates $f(x)$ as best as possible, namely, we would like to find $h \in \mathcal{H}$ such that $h \approx f$ in some strict mathematical sense which we specify below. If this is possible, we say that we *learned* $f(x)$. But if the function $f(x)$ can, in principle, take any value on *unobserved* inputs, how is it possible to learn in any meaningful sense?

Training and testing data

We will discuss the bias-variance tradeoff in the context of continuous predictions such as regression. However, many of the intuitions and ideas discussed here also carry over to classification tasks. Consider a dataset \mathcal{L} consisting of the data $\mathbf{X}_{\mathcal{L}} = \{(y_j, \mathbf{x}_j), j = 1 \dots N\}$. Let us assume that the true data is generated from a noisy model

$$y = f(\mathbf{x}) + \epsilon$$

where ϵ is normally distributed with mean zero and standard deviation σ_ϵ .

Procedure to find a predictor

We have a statistical procedure (e.g. least-squares regression) for forming a predictor $\mathbf{g}_{\mathcal{L}}(\mathbf{x})$ that gives the prediction of our model for a new data point \mathbf{x} . This estimator is chosen by minimizing a cost function which we take to be the squared error

$$\mathcal{C}(\mathbf{X}, \mathbf{g}(\mathbf{x})) = \sum_i (y_i - \mathbf{g}_{\mathcal{L}}(\mathbf{x}_i))^2.$$

What we want

We are interested in the generalization error on all data drawn from the true model, not just the error on the particular training dataset \mathcal{L} that we have in

hand. This is just the expectation of the cost function over many different data sets $\{\mathcal{L}_j\}$. Denote this expectation value by $E_{\mathcal{L}}$. In other words, we can view $\mathbf{g}_{\mathcal{L}}$ as a stochastic functional that depends on the dataset \mathcal{L} and we can think of $E_{\mathcal{L}}$ as the expected value of the functional if we drew an infinite number of datasets $\{\mathcal{L}_1, \mathcal{L}_2, \dots\}$.

The expected generalization error

We would also like to average over different instances of the “noise” ϵ and we denote the expectation value over the noise by E_{ϵ} . Thus, we can decompose the expected generalization error as

$$\begin{aligned}
E_{\mathcal{L}, \epsilon}[\mathcal{C}(\mathbf{X}, \mathbf{g}(\mathbf{x}))] &= E_{\mathcal{L}, \epsilon} \left[\sum_i (y_i - \mathbf{g}_{\mathcal{L}}(\mathbf{x}_i))^2 \right] \\
&= E_{\mathcal{L}, \epsilon} \left[\sum_i (y_i - f(\mathbf{x}_i) + f(\mathbf{x}_i) - \mathbf{g}_{\mathcal{L}}(\mathbf{x}_i))^2 \right] \\
&= \sum_i E_{\epsilon}[(y_i - f(\mathbf{x}_i))^2] + E_{\mathcal{L}, \epsilon}[(f(\mathbf{x}_i) - \mathbf{g}_{\mathcal{L}}(\mathbf{x}_i))^2] + 2E_{\epsilon}[y_i - f(\mathbf{x}_i)]E_{\mathcal{L}}[f(\mathbf{x}_i) - \mathbf{g}_{\mathcal{L}}(\mathbf{x}_i)] \\
&= \sum_i \sigma_{\epsilon}^2 + E_{\mathcal{L}}[(f(\mathbf{x}_i) - \mathbf{g}_{\mathcal{L}}(\mathbf{x}_i))^2], \tag{28}
\end{aligned}$$

where in the last line we used the fact that our noise has zero mean and variance σ_{ϵ}^2 and the sum over i applies to all terms.

Elaborating a little bit more

It is also helpful to further decompose the second term as follows:

$$\begin{aligned}
E_{\mathcal{L}}[(f(\mathbf{x}_i) - \mathbf{g}_{\mathcal{L}}(\mathbf{x}_i))^2] &= E_{\mathcal{L}}[(f(\mathbf{x}_i) - E_{\mathcal{L}}[\mathbf{g}_{\mathcal{L}}(\mathbf{x}_i)] + E_{\mathcal{L}}[\mathbf{g}_{\mathcal{L}}(\mathbf{x}_i)] - \mathbf{g}_{\mathcal{L}}(\mathbf{x}_i))^2] \\
&= E_{\mathcal{L}}[(f(\mathbf{x}_i) - E_{\mathcal{L}}[\mathbf{g}_{\mathcal{L}}(\mathbf{x}_i)])^2] + E_{\mathcal{L}}[(\mathbf{g}_{\mathcal{L}}(\mathbf{x}_i) - E_{\mathcal{L}}[\mathbf{g}_{\mathcal{L}}(\mathbf{x}_i)])^2] \\
&\quad + 2E_{\mathcal{L}}[(f(\mathbf{x}_i) - E_{\mathcal{L}}[\mathbf{g}_{\mathcal{L}}(\mathbf{x}_i)])(\mathbf{g}_{\mathcal{L}}(\mathbf{x}_i) - E_{\mathcal{L}}[\mathbf{g}_{\mathcal{L}}(\mathbf{x}_i)])] \\
&= (f(\mathbf{x}_i) - E_{\mathcal{L}}[\mathbf{g}_{\mathcal{L}}(\mathbf{x}_i)])^2 + E_{\mathcal{L}}[(\mathbf{g}_{\mathcal{L}}(\mathbf{x}_i) - E_{\mathcal{L}}[\mathbf{g}_{\mathcal{L}}(\mathbf{x}_i)])^2]. \tag{29}
\end{aligned}$$

The bias

The first term is called the bias

$$Bias^2 = \sum_i (f(\mathbf{x}_i) - E_{\mathcal{L}}[\mathbf{g}_{\mathcal{L}}(\mathbf{x}_i)])^2$$

and measures the deviation of the expectation value of our estimator (i.e. the asymptotic value of our estimator in the infinite data limit) from the true value.

The variance

The second term is called the variance

$$Var = \sum_i E_{\mathcal{L}}[(g_{\mathcal{L}}(\mathbf{x}_i) - E_{\mathcal{L}}[g_{\mathcal{L}}(\mathbf{x}_i)])^2],$$

and measures how much our estimator fluctuates due to finite-sample effects. Combining these expressions, we see that the expected out-of-sample error of our model can be decomposed as

$$E_{\text{out}} = E_{\mathcal{L}, \epsilon}[\mathcal{C}(\mathbf{X}, \mathbf{g}(\mathbf{x}))] = Bias^2 + Var + Noise.$$

The bias-variance tradeoff summarizes the fundamental tension in machine learning, particularly supervised learning, between the complexity of a model and the amount of training data needed to train it. Since data is often limited, in practice it is often useful to use a less-complex model with higher bias – a model whose asymptotic performance is worse than another model – because it is easier to train and less sensitive to sampling noise arising from having a finite-sized training dataset (smaller variance).

Example code for Bias-Variance tradeoff

Understanding what happens

Summing up

The above equations tell us that in order to minimize the expected test error, we need to select a statistical learning method that simultaneously achieves low variance and low bias. Note that variance is inherently a nonnegative quantity, and squared bias is also nonnegative. Hence, we see that the expected test MSE can never lie below $Var(\epsilon)$, the irreducible error.

What do we mean by the variance and bias of a statistical learning method? The variance refers to the amount by which our model would change if we estimated it using a different training data set. Since the training data are used to fit the statistical learning method, different training data sets will result in a different estimate. But ideally the estimate for our model should not vary too much between training sets. However, if a method has high variance then small changes in the training data can result in large changes in the model. In general, more flexible statistical methods have higher variance.

The one-dimensional Ising model

The one-dimensional Ising model with nearest neighbor interaction, no external field and a constant coupling constant J is given by

$$H = -J \sum_k^L s_k s_{k+1}, \quad (30)$$

where $s_i \in \{-1, 1\}$ and $s_{N+1} = s_1$. The number of spins in the system is determined by L . For the one-dimensional system there is no phase transition.

We will look at a system of $L = 40$ spins with a coupling constant of $J = 1$. To get enough training data we will generate 10000 states with their respective energies.

Here we use linear (ordinary least squares), ridge and LASSO regression to predict the energy in the nearest neighbor one-dimensional Ising model on a ring, i.e., the endpoints wrap around. We will use the linear regression models to fit a value for the coupling constant to achieve this.

Reformulating the problem to suit regression

A more general form for the one-dimensional Ising model is

$$H = - \sum_j^L \sum_k^L s_j s_k J_{jk}. \quad (31)$$

Here we allow for interactions beyond the nearest neighbors and a more adaptive coupling matrix. This latter expression can be formulated as a matrix-product on the form

$$H = XJ, \quad (32)$$

where $X_{jk} = s_j s_k$ and J is the matrix consisting of the elements $-J_{jk}$. This form of writing the energy fits perfectly with the form utilized in linear regression, viz.

$$y = X\omega + \epsilon, \quad (33)$$

Linear regression

The problem at hand is to try to fit the equation

$$y = f(x) + \epsilon, \quad (34)$$

where $f(x)$ is some unknown function of the data x and ϵ is normally distributed with mean zero noise with standard deviation σ_ϵ . Our job is to try to find a predictor which estimates the function $f(x)$. In linear regression we assume that we can formulate the problem as

$$y = X\omega + \epsilon, \quad (35)$$

where X and ω are now matrices. Our job at hand is now to find a **cost function** C , which we wish to minimize in order to find the best estimate of ω .

Ordinary least squares

In the ordinary least squares method we choose the cost function

$$C(X, \omega) = \|X\omega - y\|^2 = (X\omega - y)^T(X\omega - y) \quad (36)$$

We then find the extremal point of C by taking the derivative with respect to ω and setting it to zero, i.e.,

$$\frac{dC}{d\omega} = 0. \quad (37)$$

This yields the expression for ω to be

$$\omega = \frac{X^T y}{X^T X}, \quad (38)$$

which immediately imposes some requirements on X as there must exist an inverse of $X^T X$. If the expression we are modelling contains an intercept, i.e., a constant expression we must make sure that the first column of X consists of 1.

Singular Value decomposition

Doing the inversion directly turns out to be a bad idea as the matrix $X^T X$ is singular. An alternative approach is to use the **singular value decomposition**. Using the definition of the Moore-Penrose pseudoinverse we can write the equation for ω as

$$\omega = X^+ y, \quad (39)$$

where the pseudoinverse of X is given by

$$X^+ = \frac{X^T}{X^T X}. \quad (40)$$

Using singular value decomposition we have that $X = U\Sigma V^T$, where $X^+ = V\Sigma^+ U^T$. This reduces the equation for ω to

$$\omega = V\Sigma^+ U^T y. \quad (41)$$

Note that solving this equation by actually doing the pseudoinverse (which is what we will do) is not a good idea as this operation scales as $\mathcal{O}(n^3)$, where n is the number of elements in a general matrix. Instead, doing QR -factorization and solving the linear system as an equation would reduce this down to $\mathcal{O}(n^2)$ operations.

Before passing in the data to the function we append a column with ones to the training data.

Fitting with scikit-learn

Next we fit a **LinearRegression**-model from Scikit-learn for comparison.

Extracting the J -matrix from both our own method and the Scikit-learn model where we make sure to remove the intercept.

A way of looking at the coefficients in J is to plot the matrices as images.

We can see that our model for the least squares method performs close to the benchmark from Scikit-learn. It is interesting to note that OLS considers both $J_{j,j+1} = -0.5$ and $J_{j,j-1} = -0.5$ as valid matrix elements for J .

Ridge regression

Having explored the ordinary least squares we move on to ridge regression. In ridge regression we include a **regularizer**. This involves a new cost function which leads to a new estimate for the weights ω . This results in a penalized regression problem. The cost function is given by

$$C(X, \omega; \lambda) = \|X\omega - y\|^2 + \lambda \|\omega\|^2 = (X\omega - y)^T (X\omega - y) + \lambda \omega^T \omega. \quad (42)$$

Finding the extremum of this function yields the weights

$$\omega(\lambda) = \frac{X^T y}{X^T X + \lambda} \rightarrow \frac{\omega_{\text{LS}}}{1 + \lambda}, \quad (43)$$

where ω_{LS} is the weights from ordinary least squares. The last assumption assumes that X is orthogonal, which it is not. We will therefore resort to solving the equation as it stands on the left hand side.

LASSO regression

In the **Least Absolute Shrinkage and Selection Operator** (LASSO)-method we get a third cost function.

$$C(X, \omega; \lambda) = \|X\omega - y\|^2 + \lambda \|\omega\| = (X\omega - y)^T (X\omega - y) + \lambda \sqrt{\omega^T \omega}. \quad (44)$$

Finding the extremal point of this cost function is not so straight-forward as in least squares and ridge. We will therefore rely solely on the function “Lasso” from Scikit-learn.

It is quite striking how LASSO breaks the symmetry of the coupling constant as opposed to ridge and OLS. We get a sparse solution with $J_{j,j+1} = -1$.

Performance of the different models

In order to judge which model performs best at varying values of λ (for ridge and LASSO) we compute R^2 which is given by

$$R^2 = 1 - \frac{(y - \mathbf{y})^2}{(y - \bar{y})^2}, \quad (45)$$

where y is a vector with the true values of the energy, \mathbf{y} is the predicted values of y from the models and \bar{y} is the mean of \mathbf{y} .

This is the same metric used by Scikit-learn for their regression models when scoring.

Performance as function of the regularization parameter

We see how the different models perform for a different set of values for λ .

We can see that LASSO quite fast reaches a good solution for low values of λ , but will "wither" when we increase λ too much. Ridge is more stable over a larger range of values for λ , but eventually also fades away.

Finding the optimal value of λ

To determine which value of λ is best we plot the accuracy of the models when predicting the training and the testing set. We expect the accuracy of the training set to be quite good, but if the accuracy of the testing set is much lower this tells us that we might be subject to an overfit model. The ideal scenario is an accuracy on the testing set that is close to the accuracy of the training set.

From the above figure we can see that LASSO with $\lambda = 10^{-2}$ achieve a very good accuracy on the test set. This by far surpasses the other models for all values of λ .