

# Learning from data: Linear Regression

Christian Forssén<sup>1</sup>

Morten Hjorth-Jensen<sup>2,3</sup>

<sup>1</sup>Department of Physics, Chalmers University of Technology, Sweden

<sup>2</sup>Department of Physics, University of Oslo

<sup>3</sup>Department of Physics and Astronomy and National Superconducting Cyclotron Laboratory, Michigan State University

Aug 30, 2020

## 1 Linear regression

### 1.1 Why Linear Regression (aka Ordinary Least Squares)

Fitting a continuous function with linear parameterization in terms of the parameters  $\theta$ .

- Often used 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 fitting parameters  $\theta$
- 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

**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  $M$  aims at finding a likelihood function  $p(\mathbf{y}|\mathbf{x}, M, \mathcal{D}_n)$ , that is the conditional distribution for  $\mathbf{y}$  given the independent variable  $\mathbf{x}$  and a model  $M$  that has been trained on a data set  $\mathcal{D}_n$ . The data set consists of:

- $n$  cases  $i = 0, 1, 2, \dots, n - 1$
- Response variable  $y_i$  with  $i = 0, 1, 2, \dots, n - 1$ . These are sometimes referred to as target, dependent or outcome.
- For each case there will be  $p$  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 a relationship between  $\mathbf{y}$  and  $\mathbf{x}$  or to infer causal dependencies, and to make fits, and predictions, and many other things.

The  $p$  explanatory variables for the  $n$  cases in the data set are normally represented by a matrix  $\mathbf{X}$ .

The matrix  $\mathbf{X}$  is called the *design matrix*. In addition, each case is also represented by its *response variable*  $\mathbf{y}$ . The aim of regression analysis is to explain  $\mathbf{y}$  in terms of  $\mathbf{X}$  through a functional relationship like  $y_i = f(\mathbf{X}_i, *)$ .

It is common to assume a linear relationship between  $\mathbf{X}$  and  $\mathbf{y}$ . This assumption gives rise to the *linear regression model* where  $\boldsymbol{\theta} = [\theta_0, \dots, \theta_{p-1}]^T$  are the *regression parameters*. Linear regression gives us a set of analytical equations for the parameters  $\theta_j$ .

#### Example: Liquid-drop model for nuclear binding energies.

In order to understand the relation among the predictors  $p$ , the set of data  $\mathcal{D}_n$  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

$$BE(A, N, Z) = a_0 + a_1 A + a_2 A^{2/3} + a_3 Z^2 A^{-1/3} + a_4 (N - Z)^2 A^{-1},$$

we have five predictors, that is the intercept (constant term, aka bias), the  $A$  dependent term, the  $A^{2/3}$  term and the  $Z^2 A^{-1/3}$  and  $(N - Z)^2 A^{-1}$  terms. Although the predictors are somewhat complicated functions of  $A, N, Z$ , we note that the  $p = 5$  regression parameters  $\theta = (a_0, a_1, a_2, a_3, a_4)$  enter linearly. Furthermore we have  $n$  cases. It means that our design matrix is a  $p \times n$  matrix  $\mathbf{X}$ .

## 1.2 Polynomial basis functions

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 variable  $x$  so that for the data set we have  $\mathbf{x} = [x_0, x_1, \dots, x_{n-1}]$  and  $y_i = y(x_i)$  with  $i = 0, 1, 2, \dots, n - 1$ . The variable  $x_i$  could represent a physical quantity like time, temperature, position etc. We assume that  $y(x)$  is a smooth function.

Now, we don't know  $y(x)$  but we want to use the data that we have 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  $f(x)$  of degree  $p - 1$ . Since we realize that our polynomial model might not represent  $y(x)$  perfectly we also add an error term

$$y = y(x) \rightarrow y(x_i) = f(x_i) + \epsilon_i = \sum_{j=0}^{p-1} \theta_j x_i^j + \epsilon_i,$$

where  $\epsilon_i$  is the error in our approximation.

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

$$\begin{aligned} y_0 &= \theta_0 + \theta_1 x_0^1 + \theta_2 x_0^2 + \cdots + \theta_{p-1} x_0^{p-1} + \epsilon_0 \\ y_1 &= \theta_0 + \theta_1 x_1^1 + \theta_2 x_1^2 + \cdots + \theta_{p-1} x_1^{p-1} + \epsilon_1 \\ y_2 &= \theta_0 + \theta_1 x_2^1 + \theta_2 x_2^2 + \cdots + \theta_{p-1} x_2^{p-1} + \epsilon_2 \\ &\dots\dots\dots \\ y_{n-1} &= \theta_0 + \theta_1 x_{n-1}^1 + \theta_2 x_{n-1}^2 + \cdots + \theta_{p-1} x_{n-1}^{p-1} + \epsilon_{n-1}. \end{aligned}$$

Defining the vectors

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

and

$$\boldsymbol{\theta} = [\theta_0, \theta_1, \theta_2, \dots, \theta_{p-1}]^T,$$

and

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

and the design matrix

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

we can rewrite our equations as

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

The above design matrix is called a [Vandermonde matrix](#).

### General basis functions.

We are obviously not limited to the above polynomial expansions. 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

$$\begin{aligned}
y_0 &= \theta_0 x_{00} + \theta_1 x_{01} + \theta_2 x_{02} + \cdots + \theta_{p-1} x_{0p-1} + \epsilon_0 \\
y_1 &= \theta_0 x_{10} + \theta_1 x_{11} + \theta_2 x_{12} + \cdots + \theta_{p-1} x_{1p-1} + \epsilon_1 \\
y_2 &= \theta_0 x_{20} + \theta_1 x_{21} + \theta_2 x_{22} + \cdots + \theta_{p-1} x_{2p-1} + \epsilon_2 \\
&\dots\dots\dots \\
y_i &= \theta_0 x_{i0} + \theta_1 x_{i1} + \theta_2 x_{i2} + \cdots + \theta_{p-1} x_{ip-1} + \epsilon_i \\
&\dots\dots\dots \\
y_{n-1} &= \theta_0 x_{n-1,0} + \theta_1 x_{n-1,1} + \theta_2 x_{n-1,2} + \cdots + \theta_{p-1} x_{n-1,p-1} + \epsilon_{n-1}.
\end{aligned}$$

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

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

and without loss of generality we rewrite again our equations as

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

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

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

$$\begin{aligned}
y_0 &= \theta_0 x_{00} + \theta_1 x_{01} + \theta_2 x_{02} + \cdots + \theta_{p-1} x_{0p-1} + \epsilon_0 \\
y_1 &= \theta_0 x_{10} + \theta_1 x_{11} + \theta_2 x_{12} + \cdots + \theta_{p-1} x_{1p-1} + \epsilon_1 \\
y_2 &= \theta_0 x_{20} + \theta_1 x_{21} + \theta_2 x_{22} + \cdots + \theta_{p-1} x_{2p-1} + \epsilon_1 \\
&\dots\dots\dots \\
y_i &= \theta_0 x_{i0} + \theta_1 x_{i1} + \theta_2 x_{i2} + \cdots + \theta_{p-1} x_{ip-1} + \epsilon_i \\
&\dots\dots\dots \\
y_{n-1} &= \theta_0 x_{n-1,0} + \theta_1 x_{n-1,1} + \theta_2 x_{n-1,2} + \cdots + \theta_{p-1} x_{n-1,p-1} + \epsilon_{n-1}.
\end{aligned}$$

Note that the design matrix  $\mathbf{X} \in \mathbb{R}^{n \times p}$ , with the predictors referring to the column numbers and the entries  $n$  being the row elements.

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

$$\tilde{\mathbf{y}} = \mathbf{X}\boldsymbol{\theta},$$

and in order to find the optimal parameters  $\theta_i$  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 parameterized values  $\tilde{y}_i$ , namely

$$C(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=0}^{n-1} (y_i - \tilde{y}_i)^2 = \frac{1}{n} \left\{ (\mathbf{y} - \tilde{\mathbf{y}})^T (\mathbf{y} - \tilde{\mathbf{y}}) \right\},$$

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

$$C(\boldsymbol{\theta}) = \frac{1}{n} \left\{ (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) \right\}.$$

This function is one possible way to define the so-called **cost function**.

It is also common to define the cost function as

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

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

The function

$$C(\boldsymbol{\theta}) = \frac{1}{n} \left\{ (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) \right\},$$

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

$$y_i = \langle y_i \rangle = \theta_0 x_{i,0} + \theta_1 x_{i,1} + \theta_2 x_{i,2} + \cdots + \theta_{n-1} x_{i,n-1} + \epsilon_i,$$

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  $\theta_i$  we will then minimize the spread of  $C(\boldsymbol{\theta})$ , that is we are going to solve the problem

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^p} \frac{1}{n} \left\{ (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) \right\}.$$

In practical terms it means we will require

$$\frac{\partial C(\boldsymbol{\theta})}{\partial \theta_j} = \frac{\partial}{\partial \theta_j} \left[ \frac{1}{n} \sum_{i=0}^{n-1} (y_i - \theta_0 x_{i,0} - \theta_1 x_{i,1} - \theta_2 x_{i,2} - \cdots - \theta_{n-1} x_{i,n-1})^2 \right] = 0,$$

which results in

$$\frac{\partial C(\boldsymbol{\theta})}{\partial \theta_j} = -\frac{2}{n} \left[ \sum_{i=0}^{n-1} x_{ij} (y_i - \theta_0 x_{i,0} - \theta_1 x_{i,1} - \theta_2 x_{i,2} - \cdots - \theta_{n-1} x_{i,n-1}) \right] = 0,$$

or in a matrix-vector form as

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

We can rewrite

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

as

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

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

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

We note also that since our design matrix is defined as  $\mathbf{X} \in \mathbb{R}^{n \times p}$ , the product  $\mathbf{X}^T \mathbf{X} \in \mathbb{R}^{p \times p}$ . In the liquid drop model example from the Intro lecture, we had  $p = 5$  ( $p \ll n$ ) meaning that we end up with inverting a small  $5 \times 5$  matrix. This is a rather common situation, in many cases we end up with low-dimensional matrices to invert, which allow for the usage of direct linear algebra methods such as **LU** decomposition or **Singular Value Decomposition** (SVD) for finding the inverse of the matrix  $\mathbf{X}^T \mathbf{X}$ .

**Small question:** What kind of problems can we expect when inverting the matrix  $\mathbf{X}^T \mathbf{X}$ ?

### 1.3 Training scores

We can easily test our fit by computing various **training scores**. Several such measures are used in machine learning applications. First we have the **Mean-Squared Error** (MSE)

$$\text{MSE}(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n (y_{\text{data},i} - y_{\text{model},i}(\boldsymbol{\theta}))^2,$$

where we have  $n$  training data and our model is a function of the parameter vector  $\boldsymbol{\theta}$ .

Furthermore, we have the **mean absolute error** (MAE) defined as.

$$\text{MAE}(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n |y_{\text{data},i} - y_{\text{model},i}(\boldsymbol{\theta})|,$$

And the  $R^2$  score, also known as *coefficient of determination* is

$$R^2(\boldsymbol{\theta}) = 1 - \frac{\sum_{i=1}^n (y_{\text{data},i} - y_{\text{model},i}(\boldsymbol{\theta}))^2}{\sum_{i=1}^n (y_{\text{data},i} - \bar{y}_{\text{model}}(\boldsymbol{\theta}))^2},$$

where  $\bar{y}_{\text{model}}(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n y_{\text{model},i}(\boldsymbol{\theta})$  is the mean of the model predictions.

**The  $\chi^2$  function.**

Normally, the response (dependent or outcome) variable  $y_i$  is 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.

Introducing the standard deviation  $\sigma_i$  for each measurement  $y_i$  (assuming uncorrelated errors), we define the  $\chi^2$  function as

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

where the matrix  $\boldsymbol{\Sigma}$  is a diagonal  $n \times n$  matrix with  $\sigma_i^2$  as matrix elements.



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

$$\frac{\partial \chi^2(\boldsymbol{\theta})}{\partial \theta_j} = \frac{\partial}{\partial \theta_j} \left[ \frac{1}{n} \sum_{i=0}^{n-1} \left( \frac{y_i - \theta_0 x_{i,0} - \theta_1 x_{i,1} - \theta_2 x_{i,2} - \cdots - \theta_{n-1} x_{i,n-1}}{\sigma_i} \right)^2 \right] = 0,$$

which results in

$$\frac{\partial \chi^2(\boldsymbol{\theta})}{\partial \theta_j} = -\frac{2}{n} \left[ \sum_{i=0}^{n-1} \frac{x_{ij}}{\sigma_i} \left( \frac{y_i - \theta_0 x_{i,0} - \theta_1 x_{i,1} - \theta_2 x_{i,2} - \cdots - \theta_{n-1} x_{i,n-1}}{\sigma_i} \right) \right] = 0,$$

or in a matrix-vector form as

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

where we have defined the matrix  $\mathbf{A} = \mathbf{X}\boldsymbol{\Sigma}^{-1/2}$  with matrix elements  $a_{ij} = x_{ij}/\sigma_i$  and the vector  $\mathbf{b}$  with elements  $b_i = y_i/\sigma_i$ .

We can rewrite

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

as

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

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

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

If we then introduce the matrix

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

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

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

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

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

resulting in

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