

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

Morten Hjorth-Jensen<sup>1,2</sup>

<sup>1</sup>Department of Physics, University of Oslo

<sup>2</sup>Department of Physics and Astronomy and National Superconducting Cyclotron Laboratory, Michigan State University

Aug 19, 2020

## To do list

- Clean up Boston Housing data with correlations, add perhaps another more science related example. Add more info about correlation matrix, how to set it up and what it means
- show how the data is set up and show the design matrix for this data set
- Bring up covid data, say eg recent UK data on race etc
- Same for the Ising model
- Clean up part on statistics and resampling methods
- add material on ridge and lasso and link to SVD, show proofs, use doughnut example
- clean up math formalism, hats to boldface
- add material about kernel regression and splines
- rewrite Cross validation part and add discussion of train, validate and test, think of placing in statistics part or merge statistics into the reg material. add more formal stat on bootstrapping
- add discussion of training, validation and testing using ridge and lasso regression. demonstrate usage of gridsearch in validation case.

## Why Linear Regression (aka Ordinary Least Squares and family)

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

- 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 fitting parameters  $\beta$
- 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

For more discussions of Ridge and Lasso regression, [Wessel van Wieringen's article](#) is highly recommended. Similarly, [Mehta et al's article](#) is also recommended.

## 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  $\beta_j$ .

## 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

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

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  $24 \times 30000$ . This is however a classification problem and we will come back to it when we discuss Logistic Regression.

## 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  $\epsilon_i$  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

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

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

## Generalizing the fitting procedure as a linear algebra problem

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

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

**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  $\beta_i$  values?

## Optimizing our parameters

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

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

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

## Our model for the nuclear binding energies

In our [introductory notes](#) we looked at the so-called [liquid drop model](#). Let us remind ourselves about what we did by looking at the code.

We restate the parts of the code we are most interested in.

```
# Common imports
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from IPython.display import display
import os

# Where to save the figures and data files
PROJECT_ROOT_DIR = "Results"
FIGURE_ID = "Results/FigureFiles"
DATA_ID = "DataFiles/"

if not os.path.exists(PROJECT_ROOT_DIR):
    os.mkdir(PROJECT_ROOT_DIR)

if not os.path.exists(FIGURE_ID):
    os.makedirs(FIGURE_ID)

if not os.path.exists(DATA_ID):
    os.makedirs(DATA_ID)

def image_path(fig_id):
    return os.path.join(FIGURE_ID, fig_id)

def data_path(dat_id):
    return os.path.join(DATA_ID, dat_id)

def save_fig(fig_id):
    plt.savefig(image_path(fig_id) + ".png", format='png')
```

```

infile = open(data_path("MassEval2016.dat"), 'r')

# Read the experimental data with Pandas
Masses = pd.read_fwf(infile, usecols=(2,3,4,6,11),
                      names=('N', 'Z', 'A', 'Element', 'Ebinding'),
                      widths=(1,3,5,5,5,1,3,4,1,13,11,11,9,1,2,11,9,1,3,1,12,11,1),
                      header=39,
                      index_col=False)

# Extrapolated values are indicated by '#' in place of the decimal place, so
# the Ebinding column won't be numeric. Coerce to float and drop these entries.
Masses['Ebinding'] = pd.to_numeric(Masses['Ebinding'], errors='coerce')
Masses = Masses.dropna()
# Convert from keV to MeV.
Masses['Ebinding'] /= 1000

# Group the DataFrame by nucleon number, A.
Masses = Masses.groupby('A')
# Find the rows of the grouped DataFrame with the maximum binding energy.
Masses = Masses.apply(lambda t: t[t.Ebinding==t.Ebinding.max()])
A = Masses['A']
Z = Masses['Z']
N = Masses['N']
Element = Masses['Element']
Energies = Masses['Ebinding']

# Now we set up the design matrix X
X = np.zeros((len(A),5))
X[:,0] = 1
X[:,1] = A
X[:,2] = A**(2.0/3.0)
X[:,3] = A**(-1.0/3.0)
X[:,4] = A**(-1.0)
# Then nice printout using pandas
DesignMatrix = pd.DataFrame(X)
DesignMatrix.index = A
DesignMatrix.columns = ['1', 'A', 'A^(2/3)', 'A^(-1/3)', '1/A']
display(DesignMatrix)

```

With  $\beta \in \mathbb{R}^{p \times 1}$ , it means that we will hereafter write our equations for the approximation as

$$\tilde{y} = X\beta,$$

throughout these lectures.

## Optimizing our parameters, more details

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

$$\tilde{y} = X\beta,$$

and in order to find the optimal parameters  $\beta_i$  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

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

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

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

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

It is also common to define the function  $Q$  as

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

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

## Interpretations and optimizing our parameters

The function

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

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 below, we will indeed interpret  $y_i$  as a mean value

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

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

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

In practical terms it means we will require

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



which results in

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

or in a matrix-vector form as

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

## Interpretations and optimizing our parameters

We can rewrite

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

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

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

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 above case we have that  $p \ll n$ , in our case  $p = 5$  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. The methods discussed here and for many other supervised learning algorithms like classification with logistic regression or support vector machines, exhibit dimensionalities 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:** Do you think the example we have at hand here (the nuclear binding energies) can lead to problems in inverting the matrix  $\mathbf{X}^T \mathbf{X}$ ? What kind of problems can we expect?

## Some useful matrix and vector expressions

The following matrix and vector relation will be useful here and for the rest of the course. Vectors are always written as boldfaced lower case letters and matrices as upper case boldfaced letters.

$$\begin{aligned} \frac{\partial (\mathbf{b}^T \mathbf{a})}{\partial \mathbf{a}} &= \mathbf{b}, \\ \frac{\partial (\mathbf{a}^T \mathbf{A} \mathbf{a})}{\partial \mathbf{a}} &= (\mathbf{A} + \mathbf{A}^T) \mathbf{a}, \\ \frac{\partial \text{tr}(\mathbf{B} \mathbf{A})}{\partial \mathbf{A}} &= \mathbf{B}^T, \\ \frac{\partial \log |\mathbf{A}|}{\partial \mathbf{A}} &= (\mathbf{A}^{-1})^T. \end{aligned}$$

## Interpretations and optimizing our parameters

The residuals  $\epsilon$  are in turn given by

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

and with

$$X^T (y - X\beta) = 0,$$

we have

$$X^T \epsilon = X^T (y - X\beta) = 0,$$

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

Let us now return to our nuclear binding energies and simply code the above equations.

## Own code for Ordinary Least Squares

It is rather straightforward to implement the matrix inversion and obtain the parameters  $\beta$ . After having defined the matrix  $X$  we simply need to write

```
# matrix inversion to find beta
beta = np.linalg.inv(X.T.dot(X)).dot(X.T).dot(Energies)
# and then make the prediction
ytilde = X @ beta
```

Alternatively, you can use the least squares functionality in **Numpy** as

```
fit = np.linalg.lstsq(X, Energies, rcond=None)[0]
ytilde_np = np.dot(fit, X.T)
```

And finally we plot our fit with and compare with data

```
Masses['Eapprox'] = ytilde
# Generate a plot comparing the experimental with the fitted values
fig, ax = plt.subplots()
ax.set_xlabel(r'$A = N + Z$')
ax.set_ylabel(r'$E_{\mathrm{bind}} \backslash, / \mathrm{MeV}$')
ax.plot(Masses['A'], Masses['Ebinding'], alpha=0.7, lw=2,
        label='Ame2016')
ax.plot(Masses['A'], Masses['Eapprox'], alpha=0.7, lw=2, c='m',
        label='Fit')
ax.legend()
save_fig("Masses2016OLS")
plt.show()
```

## Adding error analysis and training set up

We can easily test our fit by computing the  $R^2$  score that we discussed in connection with the functionality of *scikit\_Learn* in the introductory slides. Since we are not using *scikit-Learn* here we can define our own  $R^2$  function as

```
\PYGdefault{k}{def} \PYGdefault{n+nf}{R2}\PYGdefault{p}{(}\PYGdefault{n}{y\PYGdefault{Zus}{d}
\PYGdefault{k}{return} \PYGdefault{l+m+mi}{1} \PYGdefault{o}{\PYGdefault{Zhy}{}} \PYGdefault
```