

Homework 1 Fall Semester 2020

Data Analysis and Machine Learning FYS-STK3155/FYS4155

Department of Physics, University of Oslo, Norway

Aug 19, 2020

Exercise, Setting up various Python environments

The first exercise here is of a mere technical art. We want you to have

- git as a version control software and to establish a user account on a provider like GitHub. Other providers like GitLab etc are equally fine. You can also use the University of Oslo [GitHub facilities](#).
- Install various Python packages

We will make extensive use of Python as programming language and its myriad of available libraries. You will find IPython/Jupyter notebooks invaluable in your work. You can run **R** codes in the Jupyter/IPython notebooks, with the immediate benefit of visualizing your data. You can also use compiled languages like C++, Rust, Fortran etc if you prefer. The focus in these lectures will be on Python.

If you have Python installed (we recommend Python3) and you feel pretty familiar with installing different packages, we recommend that you install the following Python packages via **pip** as

1. `pip install numpy scipy matplotlib ipython scikit-learn sympy pandas pillow`

For **Tensorflow**, we recommend following the instructions in the text of [Aurelien Geron, Hands-On Machine Learning with Scikit-Learn and TensorFlow](#), O'Reilly

We will come back to **tensorflow** later.

For Python3, replace **pip** with **pip3**.

For OSX users we recommend, after having installed Xcode, to install **brew**. Brew allows for a seamless installation of additional software via for example

1. `brew install python3`

For Linux users, with its variety of distributions like for example the widely popular Ubuntu distribution, you can use **pip** as well and simply install Python as

1. `sudo apt-get install python3` (or `python` for Python2.7)

If you don't want to perform these operations separately and venture into the hassle of exploring how to set up dependencies and paths, we recommend two widely used distributions which set up all relevant dependencies for Python, namely

- [Anaconda](#),

which is an open source distribution of the Python and R programming languages for large-scale data processing, predictive analytics, and scientific computing, that aims to simplify package management and deployment. Package versions are managed by the package management system **conda**.

- [Enthought canopy](#)

is a Python distribution for scientific and analytic computing distribution and analysis environment, available for free and under a commercial license.

We recommend using **Anaconda**.

*

Exercise 1: Our first Python encounter

This exercise has as its aim to write a small program which reads in data from a **csv** file on the equation of state for dense nuclear matter. The file is localized at <https://github.com/mhjensen/MachineLearningMSU-FRIB2020/blob/master/doc/pub/Regression/ipynb/datafiles/EoS.csv>. Thereafter you will have to set up the design matrix \mathbf{X} for the n datapoints and a polynomial of degree 3. The steps are:

- Write a Python code which reads the in the above mentioned file.
- Use for example **pandas** to order your data and find out how many data points there are.
- Set thereafter up the design matrix with dimensionality $n \times p$ where $p = 4$ and where you have defined a polynomial of degree $p - 1 = 3$. Print the matrix and check that the numbers are correct.

We recommend looking at the examples in the [regression slides](#).

Solution. `import os import numpy as np import pandas as pd import matplotlib.pyplot as plt from sklearn.model_selection import train_test_split 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')
def R2(y_data, y_model) : return 1 - np.sum((y_data - y_model)**2) / np.sum((y_data -
np.mean(y_data))**2)
def MSE(y_data, y_model) : n = np.size(y_model) return np.sum((y_data -
y_model)**2) / n
infile = open(data_path("EoS.csv"), 'r')
Read the EoS data as csv file and organized into two arrays with density and
energies EoS = pd.read_csv(infile, names = ('Density', 'Energy')) EoS['Energy'] =
pd.to_numeric(EoS['Energy'], errors='coerce') EoS = EoS.dropna() Energies =
EoS['Energy'] Density = EoS['Density'] The design matrix now as function of various polytrops X =
np.zeros((len(Density), 5)) X[:, 0] = 1 X[:, 1] = Density * (2.0/3.0) X[:, 2] =
Density * X[:, 3] = Density**2 X[:, 4] = Density**3 Wesplit the data into test and training data X_train, X_test, y_train, y_test =
train_test_split(X, Energies, test_size = 0.2) matrix inversion to find beta beta =
np.linalg.pinv(X_train.T @ X_train) @ X_train.T @ y_train and then make the prediction y_tilde =
X_test @ beta print("Training R2") print(R2(y_train, y_tilde)) print("Training MSE") print(MSE(y_train, y_tilde))
X_test @ beta print("Test R2") print(R2(y_test, y_predict)) print("Test MSE") print(MSE(y_test, y_predict))

```

*

Exercise 2: making your own data and exploring scikit-learn

We will generate our own dataset for a function $y(x)$ where $x \in [0, 1]$ and defined by random numbers computed with the uniform distribution. The function y is a quadratic polynomial in x with added stochastic noise according to the normal distribution $\mathcal{N}(\mu, \sigma)$. The following simple Python instructions define our x and y values (with 100 data points). $x = \text{np.random.rand}(100, 1)$ $y = 2.0 + 5 * x * x + 0.1 * \text{np.random.randn}(100, 1)$

1. Write your own code (following the examples under the [regression slides](#)) for computing the parametrization of the data set fitting a second-order polynomial.
2. Use thereafter **scikit-learn** (see again the examples in the regression slides) and compare with your own code.
3. Using scikit-learn, compute also the mean square error, a risk metric corresponding to the expected value of the squared (quadratic) error defined as

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

and the R^2 score function. 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

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

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

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

You can use the functionality included in scikit-learn. If you feel for it, you can use your own program and define functions which compute the above two functions. Discuss the meaning of these results. Try also to vary the coefficient in front of the added stochastic noise term and discuss the quality of the fits.

Solution. The code here is an example of where we define our own design matrix and fit parameters β .

```
import os
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split

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

def R2(y_data, y_model):
    return 1 - np.sum((y_data - y_model)**2) / np.sum((y_data - np.mean(y_data))**2)

def MSE(y_data, y_model):
    n = np.size(y_model)
    return np.sum((y_data - y_model)**2) / n

x = np.random.rand(100)
y = 2.0 + 5 * x * x + 0.1 * np.random.randn(100)

# The design matrix now as function of a given polynomial
X = np.zeros((len(x), 3))
X[:, 0] = 1.0
X[:, 1] = x
X[:, 2] = x**2

# We split the data in test and training data
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2)

# matrix inversion to find beta
beta = np.linalg.inv(X_train.T @ X_train) @ X_train.T @ y_train
print(beta)

# and then make the prediction
ytilde = X_train @ beta
print("Training R2")
print(R2(y_train, ytilde))
print("Training MSE")
print(MSE(y_train, ytilde))

ytilde_test = X_test @ beta
print("Test R2")
print(R2(y_test, ytilde_test))
print("Test MSE")
print(MSE(y_test, ytilde_test))
```

*

Exercise 3: mean values and variances in linear regression

This exercise deals with various mean values and variances in linear regression method (here it may be useful to look up chapter 3, equation (3.8) of [Trevor Hastie, Robert Tibshirani, Jerome H. Friedman, The Elements of Statistical Learning, Springer](#)).

The assumption we have made is that there exists a function $f(\mathbf{x})$ and a normal distributed error $\varepsilon \sim \mathcal{N}(0, \sigma^2)$ which describes our data

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

We then approximate this function with our model from the solution of the linear regression equations (ordinary least squares OLS), that is our function f is approximated by $\tilde{\mathbf{y}}$ where we minimized $(\mathbf{y} - \tilde{\mathbf{y}})^2$, with

$$\tilde{\mathbf{y}} = \mathbf{X}\beta.$$

The matrix \mathbf{X} is the so-called design matrix.

paragraphparagraph>paragraph>-0.5em

a) Show that the expectation value of \mathbf{y} for a given element i

$$\mathbb{E}(y_i) = \mathbf{X}_{i,*} \boldsymbol{\beta},$$

and that its variance is

$$\text{Var}(y_i) = \sigma^2.$$

Hence, $y_i \sim \mathcal{N}(\mathbf{X}_{i,*} \boldsymbol{\beta}, \sigma^2)$, that is \mathbf{y} follows a normal distribution with mean value $\mathbf{X}\boldsymbol{\beta}$ and variance σ^2 .

Solution. We can calculate the expectation value of \mathbf{y} for a given element i

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

while its variance is

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

Hence, $y_i \sim \mathcal{N}(\mathbf{X}_{i,*} \boldsymbol{\beta}, \sigma^2)$, that is \mathbf{y} follows a normal distribution with mean value $\mathbf{X}\boldsymbol{\beta}$ and variance σ^2 (not be confused with the singular values of the SVD).

paragraphparagraph>paragraph>-0.5em

b) With the OLS expressions for the parameters $\boldsymbol{\beta}$ show that

$$\mathbb{E}(\boldsymbol{\beta}) = \boldsymbol{\beta}.$$

Solution.

$$\mathbb{E}(\boldsymbol{\beta}) = \mathbb{E}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}] = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{E}[\mathbf{Y}] = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} = \boldsymbol{\beta}.$$

This means that the estimator of the regression parameters is unbiased.

paragraphparagraph>paragraph>-0.5em

c) Show finally that the variance of $\boldsymbol{\beta}$ is

$$\text{Var}(\boldsymbol{\beta}) = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}.$$

Solution. The variance of β is

$$\begin{aligned}
\text{Var}(\beta) &= \mathbb{E}\{[\beta - \mathbb{E}(\beta)][\beta - \mathbb{E}(\beta)]^T\} \\
&= \mathbb{E}\{[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} - \beta][(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} - \beta]^T\} \\
&= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbb{E}\{\mathbf{Y} \mathbf{Y}^T\} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} - \beta \beta^T \\
&= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \{\mathbf{X} \beta \beta^T \mathbf{X}^T + \sigma^2 \mathbf{I}_n\} \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} - \beta \beta^T \\
&= \beta \beta^T + \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} - \beta \beta^T = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1},
\end{aligned}$$

where we have used that $\mathbb{E}(\mathbf{Y} \mathbf{Y}^T) = \mathbf{X} \beta \beta^T \mathbf{X}^T + \sigma^2 \mathbf{I}_n$. From $\text{Var}(\beta) = \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1}$, one obtains an estimate of the variance of the estimate of the j -th regression coefficient: $\sigma^2(\hat{\beta}_j) = \sigma^2 \sqrt{[(\mathbf{X}^T \mathbf{X})^{-1}]_{jj}}$. This may be used to construct a confidence interval for the estimates.

In a similar way, we can obtain analytical expressions for say the expectation values of the parameters β and their variance when we employ Ridge regression, allowing us again to define a confidence interval.