

ISTA 421 / INFO 521 – Homework 2

Due: Friday, Sept 15, 5pm

15 points total Undergraduate / 20 points total Graduate

Pooja lakshmi Narayan

Graduate

Instructions

In this assignment you are required to write 2 scripts in python. They will be submitted as 2 separate files, although you are free to copy any parts of code from one script to the next as desired. The names for the script files are specified in problems 1 and 4 below.

In your submission, include a text file called `code.txt` that includes details about your code environment, including python version, and any instructions for executing your code from the command line.

Included in the homework 2 release are two sample scripts

`fitpoly_incomplete.py` and `cv_demo_incomplete.py`

and two data files

`womens100.csv`, and `synthdata2017.csv`

The sample scripts are provided for your convenience – you may use any part of the code in those scripts for your submissions. Note that neither will run just as provided: you must fill in the calculation of \mathbf{w} . The data files are provided in `.csv` (comma separated values) format. The script `fitpoly_incomplete.py` contains the function `read_data` which shows how to load the data files in python as a numpy array.

All problems after problem 1 require that you provide some “written” answer (in some cases also figures; you must include an informative caption with each figure), so you will also submit a `.pdf` of your written answers. (You can use L^AT_EX or any other system (including handwritten; plots, of course, must be program-generated) as long as the final version is in PDF.)

The final submission will include (minimally) the two scripts and a PDF version of your written part of the assignment. You are required to create either a `.zip` or tarball (`.tar.gz` / `.tgz`) archive of all of the files for your submission and submit the archive to the d2l dropbox by the date/time deadline above.

NOTE: Problems 5 and 6 are required for Graduate students only; Undergraduates may complete them for extra credit equal to the point value.

FCMA refers to the course text: Rogers and Girolami (2016), *A First Course in Machine Learning, Second Edition*.

For general notes on using L^AT_EX to typeset math, see: <http://en.wikibooks.org/wiki/LaTeX/Mathematics>

1. [3 points] Adapted from **Exercise 1.2** of FCMA p.35:

Write a Python script that can find the parameters \mathbf{w} for an arbitrary dataset of x_n, t_n pairs. You will use this script to answer problem 2, which only requires fitting a simple line to the data (i.e., you only need to fit parameters w_0 and w_1); however, in problem 5 you will need to fit higher-order polynomial models (e.g., $t = w_0 + w_1x + w_2x^2 + \dots$), so you must make your script generalized to handle higher-order polynomials. The script `fitpoly_incomplete.py` is provided to help get you started. `fitpoly_incomplete.py` provides helper functions to read data, plot data, and plot the model (once you've determined the weight vector \mathbf{w}), but the function for computing linear least-squares fit, `fitpoly` (starting on line 166), is *incomplete*. `fitpoly` takes as input the (one-dimensional) data vector \mathbf{x} , the target values vector \mathbf{t} , and a non-negative integer `model_order`, which represents the highest polynomial order term of the model; `fitpoly` is intended to return the \mathbf{w} vector (as a numpy array).

Recommended: (If needed) review the Appendix to HW 1, the brief tutorial to numpy arrays!

Just to state the obvious: the objective of this exercise is for you to implement the linear least squares fit solution (i.e., the normal equation) in their general linear algebra form. **DO NOT** use existing least squares solvers, such as `numpy.linalg.lstsq`, or scikit learn's `sklearn.linear_model.LogisticRegression` as your implemented solution; however, it is certainly fine to use these to help *verify* your implementation's output.

You will submit your script as a stand-alone file called `fitpoly.py`.

2. [2 point] Adapted from **Exercise 1.6** of FCMA p.35:

Table 1.3 (p.13) of FCMA lists the women's 100m gold medal Olympic results – this data is provided in the file `womens100.csv` in the data folder. Using your script from problem 1, find the 1st-order polynomial model (i.e., a line with parameters w_0 and w_1) that minimizes the squared loss of this data. Report the model here as a linear equation. Plot the data with your best-fit model and include the plot in your answer, with an informative caption.

Solution. Model linear equation: $40.924155 - 0.015072 x$

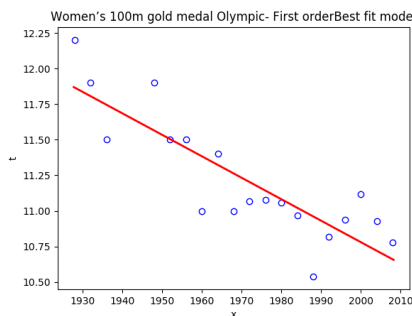


Figure 1: 1st-order polynomial model that minimized squared loss

3. [2 point] Adapted from **Exercise 1.9** of FCMA p.36:

Use your python script from problem 1 to load the data stored in the file `synthdata2016.csv` (in the data folder). Fit a 3rd order polynomial function – $f(x; \mathbf{w}) = w_0 + w_1x + w_2x^2 + w_3x^3$ – to this data (if you extended the `fitpoly_incomplete.py` script, then `model_order= 3`). Report the best-fit model parameters as a linear equation. Plot the data and your model and include the plot in your answer (be sure to include an informative caption to your plot).

Solution. Model linear equation: $-1002.322285 + 8.759129 x - 24.636503 x^2 - 4.889470 x^3$

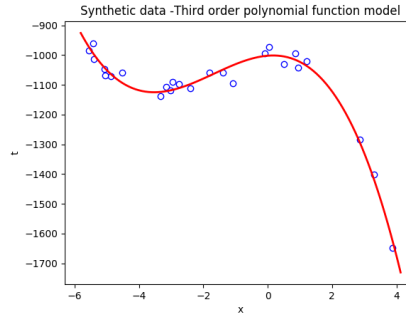


Figure 2: 3rd-order polynomial model that minimized squared loss

4. [8 points] Write a script that implements K-fold cross-validation to choose the polynomial order (between orders 0 and 7) with the best predictive error for the `synthdata2017.csv`. The provided script `cv_demo_incomplete.py` implements the synthetic data experiment described in Ch 1 (pp.31-32) of the book; you are welcome to use and adapt any part of this code you like; keep in mind that this script won't successfully execute until you add the general (matrix form) normal equation calculation on line 416. Note also that in the synthetic data experiment in `cv_demo_incomplete.py`, 1000 test data points are generated in addition to the 100 data points used for 10-fold cross-validation in the demo; for this problem (problem 4), you won't have this independent test set, only the data from `synthdata2017.csv` on which to perform K-fold cross-validation. Also, for this problem, you will perform 5-fold cross-validation, in addition to Leave-One-Out-CV (LOOCV).

Run your script with 5-fold cross-validation and LOOCV multiple times, each while randomizing the order of the data (see the `randomize_data` option of the `run_cv` function). Which model order do the two cross-validation methods predict as the best order for predictive accuracy? Do the two different cross-validation runs always agree?

Report the best-fit model parameters for the best model order according to LOOCV, and plot this model with the data. Include a plot of the training loss and CV-loss for the 8 different (0..7) polynomial model orders for **one example each** of 5-fold cross-validation and LOOCV (i.e., you will include four plots: (1) 5-fold CV and (2) related training loss, (3) LOOCV, and (4) related training loss). You can use the provided `plot_cv_results` to plot the training loss and CV loss (whether 5-fold or LOOCV) as a pair of plots.

You will submit your script as a stand-alone file called `cv.py`

Solution.

When 5-fold cross validation is run multiple times it predicts model order 3 or model order 4. But the LOOCV is consistent and predicts the model order 4 as the best fitting model! The two different cross validation runs do not agree always.

Figure 3 For the best fitting model order, i.e.4 in LOOCV the model parameters are: $w = [-1005.6642645 \ 7.66419981 \ -24.34186004 \ -4.81583084]$.

Figure 4 Plots of the training loss and 5-fold CV-loss for the 8 different (0..7) polynomial model orders

Figure 5 Plots of the training loss and LOOCV CV-loss for the 8 different (0..7) polynomial model orders

5. [2 points – **Required only for Graduates**] **Exercise 1.10** from FCMA p.36

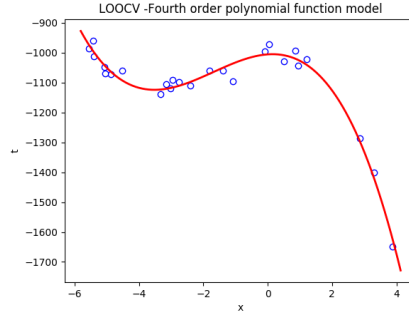


Figure 3: best-fit model parameters for the best model order according to LOOCV

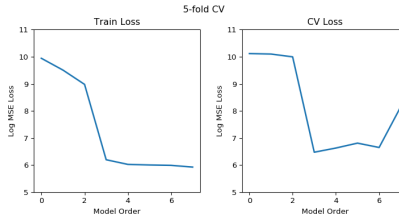


Figure 4: 5-fold CV and related training loss

Derive the optimal least squares parameter value, $\hat{\mathbf{w}}$, for the total training loss:

$$\mathcal{L} = \sum_{n=1}^N (t_n - \mathbf{w}^\top \mathbf{x}_n)^2$$

How does the expression compare with that derived from the average (mean) loss? (Hint: Express this loss in the **full** matrix form and derive the normal equation.)

Solution. Using the full matrix notation: $\mathbf{x}_n = [1, x_n]^\top$, $\mathbf{t} = [t_1, \dots, t_N]^\top$, $\mathbf{w} = [w_0, w_1]^\top$, and

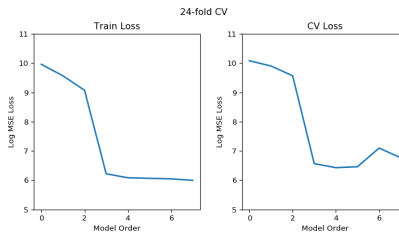


Figure 5: LOOCV CV and related training loss

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^\top \\ \mathbf{x}_2^\top \\ \vdots \\ \mathbf{x}_N^\top \end{bmatrix}$$

$$\mathbf{X}\mathbf{w} = \begin{bmatrix} w_0 + w_1x_1 \\ w_0 + w_1x_2 \\ \vdots \\ w_0 + w_1x_N \end{bmatrix}$$

$$\mathbf{t} - \mathbf{X}\mathbf{w} = \begin{bmatrix} t_1 - w_0 - w_1x_1 \\ t_2 - w_0 - w_1x_2 \\ \vdots \\ t_N - w_0 - w_1x_N \end{bmatrix}$$

$$(\mathbf{t} - \mathbf{X}\mathbf{w})^\top (\mathbf{t} - \mathbf{X}\mathbf{w}) = (t_1 - (w_0 + w_1x_1))^2 + (t_2 - (w_0 + w_1x_2))^2 + \dots + (t_N - (w_0 + w_1x_N))^2$$

Thus the total loss:

$$\mathcal{L} = (\mathbf{t} - \mathbf{X}\mathbf{w})^\top (\mathbf{t} - \mathbf{X}\mathbf{w}) = \sum_{n=1}^N (t_n - \mathbf{w}^\top \mathbf{x}_n)^2$$

To find the optimal least squares parameter value, $\hat{\mathbf{w}}$, take the derivative of the loss wrt the weight and equate to zero.

$$\begin{aligned} \mathcal{L} &= (\mathbf{t} - \mathbf{X}\mathbf{w})^\top (\mathbf{t} - \mathbf{X}\mathbf{w}) \\ &= \mathbf{t}^\top \mathbf{t} - \mathbf{t}^\top \mathbf{X}\mathbf{w} - (\mathbf{X}\mathbf{w})^\top \mathbf{t} + (\mathbf{X}\mathbf{w})^\top \mathbf{X}\mathbf{w} \\ &= \mathbf{t}^\top \mathbf{t} - 2\mathbf{w}^\top \mathbf{X}^\top \mathbf{t} + \mathbf{w}^\top \mathbf{X}^\top \mathbf{X}\mathbf{w} \end{aligned}$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = -2\mathbf{X}^\top \mathbf{t} + 2\mathbf{X}^\top \mathbf{X}\mathbf{w} = 0$$

$$2\mathbf{X}^\top \mathbf{X}\mathbf{w} = 2\mathbf{X}^\top \mathbf{t}$$

$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{t}$$

This expression is the same as the one that is derived from the average (mean) loss.

6. [3 points – **Required only for Graduates**] **Exercise 1.11** from FCMA p.36

The following expression is known as the *weighted* average loss:

$$\mathcal{L} = \frac{1}{N} \sum_{n=1}^N \alpha_n (t_n - \mathbf{w}^\top \mathbf{x}_n)^2$$

where the influence of each data point is controlled by its associated parameter. Assuming that each α_n is fixed, derive the optimal least squares parameter value $\hat{\mathbf{w}}$. (Hint: When expressing in the full matrix form, the *alpha*'s become a matrix...)

Solution. Given α expressed as a full matrix (N cross N) :

$$\boldsymbol{\alpha} = \begin{bmatrix} \alpha_1 & 0 \dots 0 \\ 0 & \alpha_2 \dots 0 \\ \vdots & \vdots \\ 0 & 0 \dots \alpha_N \end{bmatrix}$$

Similar to the above problem 5,

$$\mathbf{t} - \mathbf{X}\mathbf{w} = \begin{bmatrix} t_1 - w_0 - w_1 x_1 \\ t_2 - w_0 - w_1 x_2 \\ \vdots \\ t_N - w_0 - w_1 x_N \end{bmatrix}$$

$$\begin{aligned} \mathcal{L} &= \frac{1}{N} (\mathbf{t} - \mathbf{X}\mathbf{w})^\top \boldsymbol{\alpha} (\mathbf{t} - \mathbf{X}\mathbf{w}) \\ &= \frac{1}{N} (\mathbf{t}^\top \boldsymbol{\alpha} \mathbf{t} - \mathbf{t}^\top \boldsymbol{\alpha} \mathbf{X}\mathbf{w} - (\mathbf{X}\mathbf{w})^\top \boldsymbol{\alpha} \mathbf{t} + (\mathbf{X}\mathbf{w})^\top \boldsymbol{\alpha} \mathbf{X}\mathbf{w}) \\ &= \frac{1}{N} (\mathbf{t}^\top \boldsymbol{\alpha} \mathbf{t} - 2\mathbf{w}^\top \mathbf{X}^\top \boldsymbol{\alpha} \mathbf{t} + \mathbf{w}^\top \mathbf{X}^\top \boldsymbol{\alpha} \mathbf{X}\mathbf{w}) \end{aligned}$$

To find the optimal least squares parameter value, $\hat{\mathbf{w}}$, take the derivative of the loss wrt the weight and equate to zero.

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \frac{-2}{N} \mathbf{X}^\top \boldsymbol{\alpha} \mathbf{t} + \frac{2}{N} \mathbf{X}^\top \boldsymbol{\alpha} \mathbf{X}\mathbf{w} = 0$$

$$\mathbf{X}^\top \boldsymbol{\alpha} \mathbf{t} = \mathbf{X}^\top \boldsymbol{\alpha} \mathbf{X}\mathbf{w}$$

$$\hat{\mathbf{w}} = (\mathbf{X}^\top \boldsymbol{\alpha} \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\alpha} \mathbf{t}$$