

"Machine Learning and Computational Statistics"

3rd Homework

Exercise 1:

(a) Use the Lagrangian function of the ridge regression problem

$$\min_{\theta} L(\theta) = \sum_{n=1}^N (y_n - \theta^T x_n)^2 + \lambda \|\theta\|^2$$

and show that the solution satisfies the equation

$$\left(\sum_{n=1}^N x_n x_n^T + \lambda I \right) \hat{\theta} = \sum_{n=1}^N y_n x_n \quad (\text{A})$$

Hints: Take the gradient of $L(\theta)$ with respect to θ equate to $\mathbf{0}$ and solve.

$$\text{It is } (\theta^T x) x = (x x^T) \theta$$

It is $Az - \lambda z = Az - \lambda Iz = (A - \lambda I)z$, where A is a matrix, z is a vector and λ is a scalar and I is the identity matrix.

$$\text{It is } \frac{\partial \|\theta\|^2}{\partial \theta} = \frac{\partial (\theta^T \theta)}{\partial \theta} = 2\theta$$

(b) Prove that the above solution can be expressed in matrix form as

$$\hat{\theta} = (X^T X + \lambda I)^{-1} X^T y$$

$$\text{where } X = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_N^T \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & \cdots & x_{1l} \\ 1 & x_{21} & \cdots & x_{2l} \\ 1 & \vdots & \ddots & \vdots \\ 1 & x_{N1} & \cdots & x_{Nl} \end{bmatrix} \quad \text{and} \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

Hint: Prove that $X^T X = \sum_{n=1}^N x_n x_n^T + \lambda I$ and $X^T y = \sum_{n=1}^N y_n x_n$

Exercise 2:

Consider a 1-dimensional parameter estimation problem, where the true parameter value is θ_o . Let $\hat{\theta}_{MVU}$ be a minimum variance unbiased estimator of θ_o . Consider the parametric set F of all estimators of the form

$$\hat{\theta}_b = (1 + \alpha) \hat{\theta}_{MVU}, \quad (1)$$

with $\alpha \in \mathbb{R}$. Recall that

$$\text{MSE}(\hat{\theta}) = E[(\hat{\theta} - E(\hat{\theta}))^2] + (E(\hat{\theta}) - \theta_o)^2 \quad (2)$$

- What can you infer from the fact that $\hat{\theta}_{MVU}$ is an unbiased estimator of θ_o ?
- Prove that all $\hat{\theta}_b$'s of F , for $\alpha \neq 0$, are biased estimators of θ_o .
- Find the $MSE(\hat{\theta}_{MVU})$ using eq. (2). Explain why this value cannot be zero for finite N .
- Express the $MSE(\hat{\theta}_b)$ in terms of $MSE(\hat{\theta}_{MVU})$ (Substitute eq. (1) to eq. (2) and after some algebra, utilize the results of (c) above).
- Determine the range of values of the parameter α that result to estimators with MSE lower than $MSE(\hat{\theta}_{MVU})$ (Consider the inequality $MSE(\hat{\theta}_b) < MSE(\hat{\theta}_{MVU})$, substitute $MSE(\hat{\theta}_b)$ from (d) and you will end up with a second order polynomial wrt α , which should be negative. Determine its roots and define the range of values of α where polynomial is negative).
- Prove that for any value of α in the range defined in (e), it is $|\hat{\theta}_b| < |\hat{\theta}_{MVU}|$ (Show first that $|1+\alpha| < 1$ and then utilize eq. (1)).
- Determine the value α^* of the parameter α that corresponds to the estimator giving the lowest MSE (Consider the expression of $MSE(\hat{\theta}_b)$ derived in (d), take the derivative wrt α , equate to 0 and solve).
- Explain why in practice α^* cannot be determined.

Exercise 3:

Consider a set N pairs (y_n, \mathbf{x}_n) , $n=1, \dots, N$, satisfying the equation

$$y_n = \theta_o^T \mathbf{x}_n + \eta_n, \quad (3)$$

where η_n is normally distributed **zero mean** i.i.d. noise. As it is known, the LS estimator satisfies the equation

$$\left(\sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T \right) \theta = \sum_{n=1}^N y_n \mathbf{x}_n \quad (4)$$

Consider now the special case where the θ is a scalar and $\mathbf{x}_n=1$ for all n . In this case, eq. (3) becomes

$$y_n = \theta_o + \eta_n. \quad (3)^1$$

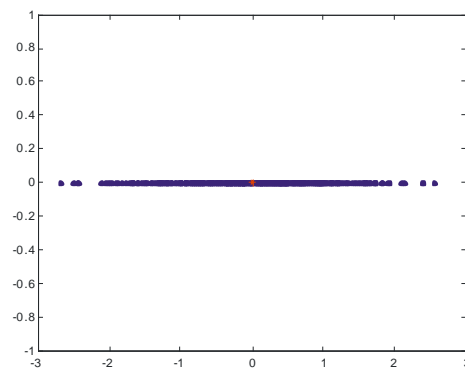
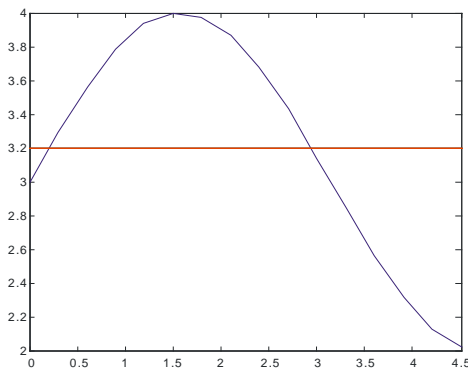
¹ In regression this implies a line perpendicular to the y-axis (left figure), while in classification implies a set of points spread around θ_o (right figure)

Let y_n denote the rv that models y_n . Let also $\bar{y} = \frac{1}{N} \sum_{n=1}^N y_n$ be the mean of N statistically independent rv's that follow the same pdf. Each one of the y_n 's models a y_n .

- Derive the LS estimator of θ_o from eq. (4) for this case (Note that in eq. (4) all x_n 's are now scalars equal to 1).
- Prove that y_n is an unbiased estimator of θ_o (Show that $E[y_n] = \theta_o$).
- Prove that \bar{y} is an unbiased estimator of θ_o (Show that $E[\bar{y}] = \theta_o$).
- Since \bar{y} is the LS estimator for the 1-dim. case, it is known that it is minimum variance unbiased estimator. From now on it will be denoted by $\hat{\theta}_{MVU}$.
- Prove that the ridge regression estimator for the present 1-dim. case is expressed as $\hat{\theta} = \frac{\sum_{n=1}^N y_n}{N + \lambda}$ (use the formula (A) proved in exercise 1a).
- Denoting as $\hat{\theta}$ the ridge regression estimator, express it in terms of $\hat{\theta}_{MVU}$.
- Prove that $\hat{\theta}$ is a biased estimator (Show that $E[\hat{\theta}] \neq \theta_o$).
- Verify that $|\hat{\theta}| < |\hat{\theta}_{MVU}|$.
- (*) Recalling from exercise 3 that $\hat{\theta}_b = (1 + \alpha)\hat{\theta}_{MVU}$ express the quantity α in terms of λ . Then, utilizing exercise 3(e), derive the range of values of λ for which the MSE of the ridge regression estimator is less than that of the least squares estimator.

(*) This is not obligatory.

Exercise 4 (regularization - python code):



Consider the data set given in the attachment file (also the code for reading from python is also given). Specifically, it consists of 10 data pairs of the form (y_i, x_i) , $i=1, \dots, 10$. All y_i 's are accumulated in the vector **y** while all x_i 's are accumulated in the vector **x**.

The aim is to unravel the relation between x_i 's and y_i 's.

- (a) Plot the data.
- (b) Fit a 8th degree polynomial on the data using the LS estimator and plot the results (data points and the curve resulting from the fit). Output also the estimates of the parameters of the polynomial.
- (c) Fit a 8th degree polynomial on the data using the **regularized** LS estimator and plot the results (data points and the curve resulting from the fit). Output also the estimates of the parameters of the polynomial.
- (d) Discuss briefly on the results.

Hint for (b), (c): The X matrix that needs to be constructed will contain 10 rows, one for each x_i . Each row will have the form $[1, x_i, x_i^2, x_i^3, x_i^4, x_i^5, x_i^6, x_i^7, x_i^8]$.