

## Assignment 1

Homework assignments will be done individually: each student must hand in their own answers. Use of partial or entire solutions obtained from others or online is strictly prohibited. Electronic submission on Canvas is mandatory.

1. **Maximum Likelihood estimator** (10 points) Assuming data points are independent and identically distributed (i.i.d.), the probability of the data set given parameters:  $\mu$  and  $\sigma^2$  (the likelihood function):

$$P(\mathbf{x}|\mu, \sigma^2) = \prod_{n=1}^N \mathcal{N}(x_n|\mu, \sigma^2)$$

Please calculate the solution for  $\mu$  and  $\sigma^2$  using Maximum Likelihood (ML) estimator.

Sol:  $\therefore f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$

$$\begin{aligned} \therefore P(\mathbf{x}|\mu, \sigma^2) &= \prod_{n=1}^N f(x_n) \\ &= \prod_{n=1}^N \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_n-\mu)^2}{2\sigma^2}} \\ &= (2\pi)^{-\frac{N}{2}} \cdot (\sigma)^{-N} \cdot e^{-\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n-\mu)^2} \end{aligned}$$

$$\ln P(\mathbf{x}|\mu, \sigma^2) = -\frac{N}{2} \ln 2\pi - \frac{N}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} \sum_{n=1}^N (x_n-\mu)^2$$

$$\therefore \text{we have } \begin{cases} \frac{\partial \ln P(\mathbf{x}|\mu, \sigma^2)}{\partial \mu} = \frac{1}{\sigma^2} \left( \sum_{n=1}^N x_n - N\mu \right) = 0 \\ \frac{\partial \ln P(\mathbf{x}|\mu, \sigma^2)}{\partial \sigma^2} = -\frac{N}{2\sigma^2} + \frac{\sum_{n=1}^N (x_n-\mu)^2}{2\sigma^4} = 0 \end{cases}$$

$$\therefore \text{we get } \begin{cases} \mu = \frac{1}{N} \sum_{n=1}^N x_n = \bar{x} \\ \sigma^2 = \frac{\sum_{n=1}^N (x_n - \bar{x})^2}{N} \end{cases}$$

- 
2. **Maximum Likelihood** (10 points) We assume there is a true function  $f(x)$  and the target value is given by  $y = f(x) + \epsilon$  where  $\epsilon$  is a Gaussian distribution with mean 0 and variance  $\sigma^2$ . Thus,

$$p(y|x, w, \beta) = \mathcal{N}(y|f(x), \beta^{-1})$$

where  $\beta^{-1} = \sigma^2$ .

Assuming the data points are drawn independently from the distribution, we obtain the likelihood function:

$$p(y|x, w, \beta) = \prod_{n=1}^N \mathcal{N}(y_n|f(x), \beta^{-1})$$

Please show that maximizing the likelihood function is equivalent to minimizing the sum-of-squares error function.

$$\text{Sol: } p(y|x, w, \beta) = \prod_{n=1}^N \mathcal{N}(y_n|f(x), \beta^{-1})$$

$$\therefore \ln p(y|x, w, \beta) = -\frac{N}{2} \ln 2\pi - \frac{N}{2} \ln \beta - \beta \cdot \frac{1}{2} \sum_{n=1}^N (y_n - f(x))^2$$

$$\therefore E_0(w) = \frac{1}{2} \sum_{n=1}^N (y_n - w^T x_n)^2 = \frac{1}{2} \sum_{n=1}^N (y_n - f(x_n))^2$$

$\therefore$  ML function equal to minimizing the sum of squares error function.

3. **MAP estimator** (15 points) Given input values  $\mathbf{x} = (x_1, \dots, x_N)^T$  and their corresponding target values  $\mathbf{y} = (y_1, \dots, y_N)^T$ , we estimate the target by using function  $f(x, \mathbf{w})$  which is a polynomial curve. Assuming the target variables are drawn from Gaussian distribution:

$$p(y|x, \mathbf{w}, \beta) = \mathcal{N}(y|f(x, \mathbf{w}), \beta^{-1})$$

and a prior Gaussian distribution for  $\mathbf{w}$ :

$$p(\mathbf{w}|\alpha) = \left(\frac{\alpha}{2\pi}\right)^{(M+1)/2} \exp\left(-\frac{\alpha}{2} \mathbf{w}^T \mathbf{w}\right)$$

Please prove that maximum posterior (MAP) is equivalent to minimizing the regularized sum-of-squares error function. Note that the posterior distribution of  $\mathbf{w}$  is  $p(\mathbf{w}|\mathbf{x}, \mathbf{y}, \alpha, \beta)$ . Hint: use Bayes' theorem.

Sol: Bayes' function:  $P(\mathbf{w}|\mathbf{A}) = \frac{P(\mathbf{A}|\mathbf{w})P(\mathbf{w})}{\sum_{j=1}^N P(\mathbf{A}|\mathbf{w}_j)P(\mathbf{w}_j)}$

$$\begin{aligned} \therefore \text{MAP} &= \arg\max P(\mathbf{w}|\mathbf{y}) \\ &= \arg\max \left( \frac{P(\mathbf{y}|\mathbf{x})P(\mathbf{w})}{\sum_{j=1}^N P(\mathbf{y}|\mathbf{x}_j)P(\mathbf{w}_j)} \right) \\ &= \arg\max \prod_{i=1}^N P(y_i|x_i)P(\mathbf{w}) \end{aligned}$$

$$\therefore P(\mathbf{y}|\mathbf{x}, \mathbf{w}, \beta) = \frac{1}{\sqrt{2\pi}\beta} e^{-\frac{\beta}{2} (y_i - f(x_i; \mathbf{w}))^2}$$

$$P(\mathbf{w}|\alpha) = \left(\frac{\alpha}{2\pi}\right)^{(M+1)/2} e^{-\frac{\alpha}{2} \mathbf{w}^T \mathbf{w}}$$

$$\text{MAP} = \left(\frac{\alpha}{2\pi}\right)^{-\frac{N}{2}} \cdot \beta^{\frac{N}{2}} \cdot e^{-\frac{1}{2}\beta \sum_{i=1}^N (y_i - f(x_i; \mathbf{w}))^2} \cdot \left(\frac{\alpha}{2\pi}\right)^{\frac{M+1}{2}} e^{-\frac{\alpha}{2} \mathbf{w}^T \mathbf{w}}$$

$$\ln \text{MAP} = -\frac{N}{2} \ln 2\pi + \frac{N}{2} \ln \beta - \frac{1}{2} \beta \sum_{i=1}^N (y_i - f(x_i; \mathbf{w}))^2 + \left(\frac{M+1}{2}\right) \ln \left(\frac{\alpha}{2\pi}\right) - \frac{\alpha}{2} \mathbf{w}^T \mathbf{w}$$

$$\therefore -\frac{N}{2} \ln 2\pi, \frac{N}{2} \ln \beta, \frac{M+1}{2} \ln \frac{\alpha}{2\pi}, -\frac{\alpha}{2} \mathbf{w}^T \mathbf{w} \text{ are constant}$$

$$\therefore \text{we need to minimize } \frac{1}{2} \beta \sum_{i=1}^N (y_i - f(x_i; \mathbf{w}))^2 \text{ to make MAP max.}$$

4. Linear model (20 points) Consider a linear model of the form:

$$f(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{i=1}^D w_i x_i$$

together with a sum-of-squares error/loss function of the form:

$$L_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{f(\mathbf{x}_n, \mathbf{w}) - y_n\}^2$$

Now suppose that Gaussian noise  $\epsilon_i$  with zero mean and variance  $\sigma^2$  is added independently to each of the input variables  $x_i$ . By making use of  $E[\epsilon_i] = 0$  and  $E[\epsilon_i \epsilon_j] = \delta_{ij} \sigma^2$ , show that minimizing  $L_D$  averaged over the noise distribution is equivalent to minimizing the sum-of-squares error for noise-free input variables with the addition of a weight-decay regularization term, in which the bias parameter  $w_0$  is omitted from the regularizer.

$$\begin{aligned} \text{Sol: } L_D(\mathbf{w}) &= \frac{1}{2} \sum_{n=1}^N \left\{ w_0 + \sum_{i=1}^D w_i x_i - y_n \right\}^2 \\ &= \frac{1}{2} \sum_{n=1}^N \left\{ w_0 + \sum_{i=1}^D w_i (x_i + \epsilon_i) - y_n \right\}^2 \\ &= \frac{1}{2} \sum_{n=1}^N \left\{ w_0 + \sum_{i=1}^D w_i x_i - y_n + \sum_{i=1}^D w_i \epsilon_i \right\}^2 \\ &= \frac{1}{2} \sum_{n=1}^N \left\{ y(\mathbf{x}_n, \mathbf{w}) - y_n + \sum_{i=1}^D w_i \epsilon_i \right\}^2 \\ &= \frac{1}{2} \sum_{n=1}^N \left\{ [y(\mathbf{x}_n, \mathbf{w}) - y_n]^2 + 2 \left( \sum_{i=1}^D w_i \epsilon_i \right) (y(\mathbf{x}_n, \mathbf{w}) - y_n) + \left( \sum_{i=1}^D w_i \epsilon_i \right)^2 \right\} \end{aligned}$$

$$\begin{aligned} E \left[ \left( \sum_{i=1}^D w_i \epsilon_i \right)^2 \right] &= E \left[ \sum_{i=1}^D \sum_{j=1}^D w_i \epsilon_i w_j \epsilon_j \right] = \sum_{i=1}^D \sum_{j=1}^D w_i w_j E[\epsilon_i \epsilon_j] \\ &= \sum_{i=1}^D \sum_{j=1}^D w_i w_j \delta_{ij} \sigma^2 \\ &= \sigma^2 \sum_{i=1}^D w_i^2 \end{aligned}$$

$$\begin{aligned} E \left[ 2 \left( \sum_{i=1}^D w_i \epsilon_i \right) (y(\mathbf{x}_n, \mathbf{w}) - y_n) \right] &= E[\epsilon_i] \cdot 2 \sum_{i=1}^D w_i (y(\mathbf{x}_n, \mathbf{w}) - y_n) = 0 \\ \therefore \text{that minimizing } L_D \text{ over} \end{aligned}$$

$$E \left[ \left( \sum_{i=1}^D w_i \epsilon_i \right)^2 \right] = \sigma^2 \sum_{i=1}^D w_i^2 \text{ is a constant, } E \left[ 2 \left( \sum_{i=1}^D w_i \epsilon_i \right) (y(\mathbf{x}_n, \mathbf{w}) - y_n) \right] = 0$$

$\therefore$  minimizing  $L_D$  averaged over the noise distribution is equivalent to minimizing the sum of squares error for

5. Linear regression (45 points) Please choose one of the below problems. You will need to submit your code.

a) UCI Machine Learning: Facebook Comment Volume Data Set

Please apply both Lasso regression and Ridge regression algorithms on this dataset for predicting the number of comments in next H hrs (H is given in the feature). You do not need to use all the features. Use K-fold cross validation and report the mean squared error (MSE) on the testing data. You need to write down every step in your experiment.

b) UCI Machine Learning: Bike Sharing Data Set

Please apply both Lasso regression and Ridge regression algorithms on this dataset for predicting the count of total rental bikes including both casual and registered. You do not need to use all the features. Use K-fold cross validation and report the mean squared error (MSE) on the testing data. You need to write down every step in your experiment.

b) Step 1: Download UCI Machine Learning: Bike Sharing Data Set  
and create a new python file

Step 2: ~~Get feature from~~  
Get feature from .csv file as X and set a label  
as y.

Step 3: Using train-test-split get X\_train, X\_test,  
Y\_train and Y\_test.

Step 4: Create ridge and lasso linear regression (Using sklearn import these  
two model)

Step 5: Create "cross\_val\_score()" this method to achieve k-fold cross  
validation  
generate y-predict.

Step 6: Create "mean\_squared\_error" this method to achieve MSE  
on testing data.