

Machine Learning, 2021 Spring

Homework 4

Due on 23:59 MAY 5, 2021

Please submit your homework in “pdf” format. Submit the supplementary materials (e.g., files for code) in an **extra** “zip” file.

Problem 1

For a random variable z , let \bar{z} denote its mean, i.e., $\bar{z} = \mathbb{E}[z]$.

Suppose we are given a dataset $\mathcal{D} = \{x^{(i)}, y^{(i)}\}_{i=1}^m$ drawn i.i.d. from some unknown distribution $P(X, Y)$. Given x , the expected label is defined as

$$\bar{y}(x) = \mathbb{E}_{y|x}[Y] \quad (1)$$

which denotes the label we would expect to obtain.

Next, we run some learning algorithm, such as SVM, linear regression, from which we learned our hypothesis function $h_{\mathcal{D}}$.

Now for a new data point (x, y) sampled from $P(X, Y)$ and out of \mathcal{D} , we want to investigate the expected error between the predicted value $h_{\mathcal{D}}(x)$ and the observation y , i.e.,

$$\mathbb{E}_{\mathcal{D}, x, y}[(y - h_{\mathcal{D}}(x))^2]. \quad (2)$$

This error can be decomposed into three parts namely: variance, bias, and noise, where the expectation is taken over all possible training set \mathcal{D} and all (x, y) . Here

$$\begin{aligned} \text{bias}^2 &= \mathbb{E}_x[(\bar{y}(x) - \bar{h}(x))^2] \\ \text{variance} &= \mathbb{E}_{x, \mathcal{D}}[(\bar{h}(x) - h_{\mathcal{D}}(x))^2] \\ \text{noise} &= \mathbb{E}_{x, y}[(y - \bar{y}(x))^2] \end{aligned} \quad (3)$$

where $\bar{h}(x) = \mathbb{E}_{\mathcal{D}}[h_{\mathcal{D}}(x)]$ is the “average approximator” by averaging classifiers on all possible training dataset \mathcal{D} .

The error bias is the amount by which the expected model prediction differs from the true value or target; while variance measures how inconsistent are the predictions from one another, over different training sets, not whether they are accurate or not. Models that exhibit small variance and high bias underfit the truth target. Models that exhibit high variance and low bias overfit the truth target.

The data scientist’s goal is to simultaneously reduce bias and variance as much as possible in order to obtain as accurate model as is feasible. However, there is a trade-off to be made when selecting models of different flexibility or complexity and in selecting appropriate training sets to minimize these sources of error.

Question: [2 pts]

Show that

$$\mathbb{E}_{\mathcal{D}, x, y}[(y - h(x))^2] = \text{variance} + \text{bias}^2 + \text{noise}. \quad (4)$$

$$\begin{aligned}
& 1. E_{D, X, y}[(y - \hat{h}(x))^2] \\
&= E_{D, X, y}[(y - \bar{y}(x) + \bar{y}(x) - \hat{h}(x))^2] \\
&= E_{D, X, y}[(y - \bar{y}(x))^2 + 2(y - \bar{y}(x))(\bar{y}(x) - \hat{h}(x)) + (\bar{y}(x) - \hat{h}(x))^2] \\
&= E_{D, X, y}[(y - \bar{y}(x))^2] + E_{D, X, y}[(\bar{y}(x) - \hat{h}(x))^2] + 2E_{D, X, y}[(y - \bar{y}(x))(\bar{y}(x) - \hat{h}(x))] \\
&= E_{X, y}[(y - \bar{y}(x))^2] \quad [\text{Since } y \text{ and } \bar{y}(x) \text{ are functions of } x \text{ and } y] + E_{D, X}[(\bar{y}(x) - \hat{h}(x))^2] + 2 \int_{\mathcal{D}} \int_{\mathcal{X}} (y - \bar{y}(x))(\bar{y}(x) - \hat{h}(x)) f(x) f(y|x) f(D) dy dx dD \\
&= \text{noise} + E_{D, X}[(\bar{y}(x) - \hat{h}(x) + \bar{h}(x) - \hat{h}(x))^2] + 2 \int_{\mathcal{D}} \int_{\mathcal{X}} (\bar{y}(x) - \hat{h}(x)) f(x) f(D) \int_{\mathcal{Y}} (y - \bar{y}(x)) f(y|x) dy dx dD \\
&= \text{noise} + E_{D, X}[(\bar{y}(x) - \hat{h}(x))^2 + 2(\bar{y}(x) - \hat{h}(x))(\bar{h}(x) - \hat{h}(x)) + (\bar{h}(x) - \hat{h}(x))^2] + 0 \quad [\because \int_{\mathcal{Y}} (y - \bar{y}(x)) f(y|x) dy = \int_{\mathcal{Y}} y f(y|x) dy - \bar{y}(x) \int_{\mathcal{Y}} f(y|x) dy = \bar{y}(x) - \bar{y}(x) = 0] \\
&= \text{noise} + E_{D, X}[(\bar{y}(x) - \hat{h}(x))^2] + E_{D, X}[(\bar{h}(x) - \hat{h}(x))^2] + 2E_{D, X}[(\bar{y}(x) - \hat{h}(x))(\bar{h}(x) - \hat{h}(x))] \quad [\text{Since } \bar{y}(x) \text{ and } \bar{h}(x) \text{ are functions of } x] + \text{variance} + 2 \int_{\mathcal{D}} \int_{\mathcal{X}} (\bar{y}(x) - \hat{h}(x))(\bar{h}(x) - \hat{h}(x)) f(x) f(D) dx dD \\
&= \text{noise} + \text{bias}^2 + \text{variance} + 2 \int_{\mathcal{D}} \int_{\mathcal{X}} (\bar{y}(x) - \hat{h}(x)) f(x) \int_{\mathcal{Y}} (\bar{h}(x) - \hat{h}(x)) f(y|x) dy dx dD = 0 \quad [\because \int_{\mathcal{D}} (\bar{h}(x) - \hat{h}(x)) f(x) dx = \bar{h}(x) \int_{\mathcal{D}} f(x) dx - \int_{\mathcal{D}} \hat{h}(x) f(x) dx = \bar{h}(x) - \bar{h}(x) = 0] \\
&= \text{variance} + \text{bias}^2 + \text{noise} \quad \text{Q.E.D}
\end{aligned}$$

Problem 2

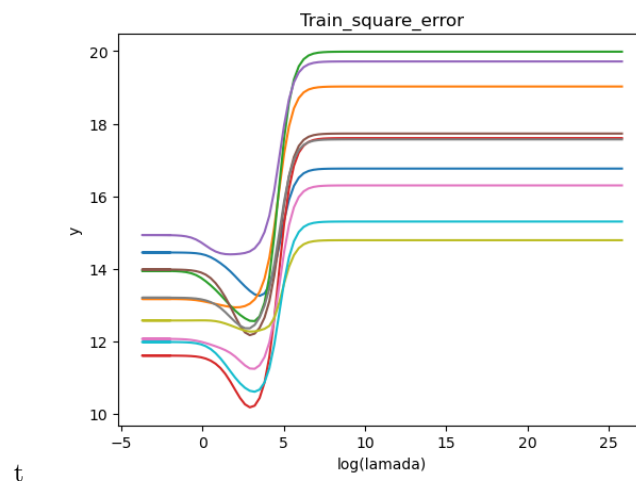
Given the training dataset “data/crime-train.txt” and the test dataset “data/crime-test.txt” (For more information about the datasets, you may refer to “README.md”).

We’d like to use the training dataset to fit a model which can predict the crime rate in new communities, and evaluate model performance on the test set. As there are a considerable number of input variables, overfitting is a serious issue. In order to avoid this, we will use the L2 regularization.

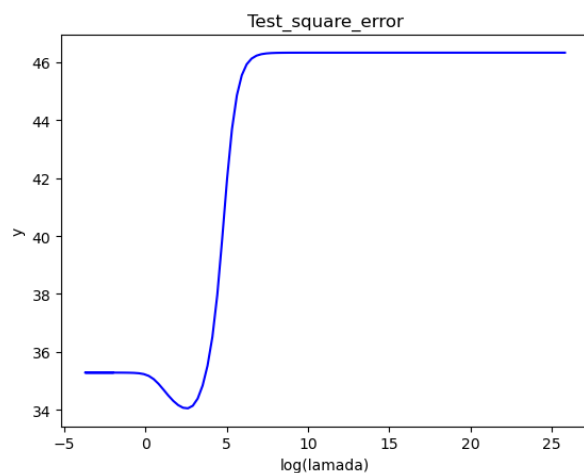
The main goal of this homework is to give you some experience using L2 regularization as a method for variable selection and using 10-folder cross-validation as a technique to get an insight on how the model will generalize to an independent dataset. Your function should accept a scalar value of λ , a vector-valued response variable (\mathbf{y}), a matrix of input variables (\mathbf{X}), and an initial vector of weights (\mathbf{w}_0). It should output a vector of coefficient values ($\hat{\mathbf{w}}$).

In your analysis, include:

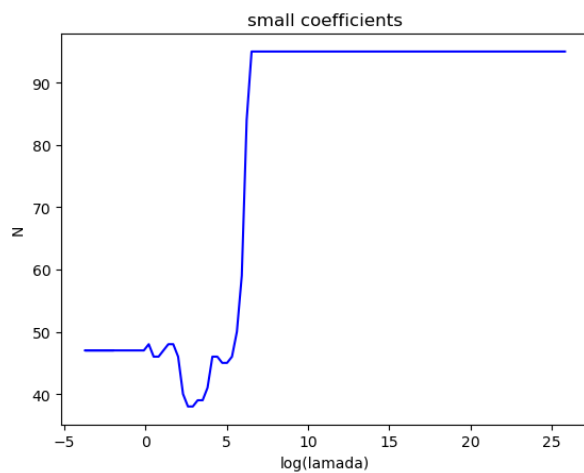
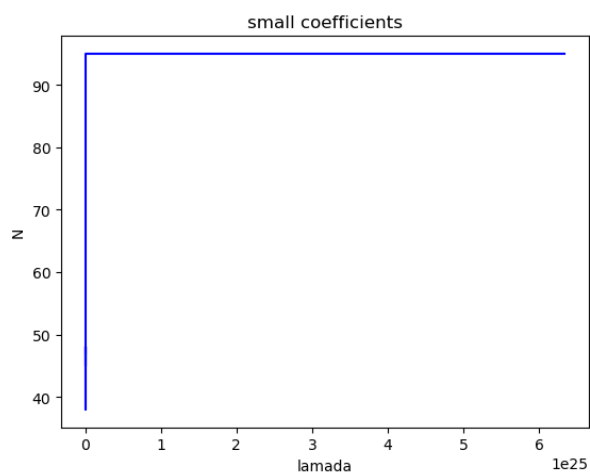
1. A plot of $\log(\lambda)$ against the squared error in the 10-folder splitted training data. [1 pts]



2. A plot of $\log(\lambda)$ against the squared error in the test data. [0.5 pts]



3. A plot of λ against the number of small coefficients (you can set a threshold), and a brief commentary on the task of selecting λ . [1 pts]



(1).I think the graph with λ as x-axis is not so clear about the trend so I also add the graph of $\log(\lambda)$ as x-axis.

(2).The threshold:0.0001

Brief analysis:

(1).Since λ is the penalty coefficient of norm of coefficient,when it is small, it means not give to much restrict to large coefficient,so the variance maybe big because of the extreme situation which had learned out some large coefficients.The generalization is not too bad in this range.

(2).When λ is not too small or too big,it gives some restricts to avoid the extreme situation,which cause the variance down.And it also not give to much restricts so most of the situation of coefficient can be learnt.The generalization is best in this range.[In this case is about $1e3-1e5$].

(3).When λ is too big,strict penalty on the coefficient makes them tend to 0.Square error is big and the generalization is bad.

4. For the λ that gave the best test set performance, which variable had the largest (most positive) coefficient? What about the most negative? Discuss briefly. [0.5 pts]

(1).best $\lambda=838.8608$

(2).largest coefficient: PctIlleg=0.03244458909756975

smallest coefficient: PctKids2Par=-0.020147945079381974

(3).The same as the analysis in 3,best λ will appear when it is not too small or too big.So it can give some restricts to extreme situation to decrease the variance and also not give too much restricts so that most of the coefficient can be learnt.The generalization is best when λ in this range.

Problem 3

The goal in the prediction problem is to be able to make prediction for the target variable t given some new value of the input variable x on the basis of a set of training data comprising N input values $\mathbf{x} = (x_1, \dots, x_N)^T$ and their corresponding target variable $\mathbf{t} = (t_1, \dots, t_N)^T$.

We assume that, given the value of x , the corresponding value of t has a Gaussian distribution with a mean equal to the value $y(x, w)$ and the variance σ , where $y(x, w)$ is the prediction function. For example, for the linear regression, the $y(x, \mathbf{w}) = w_0 + w_1 x$.

Thus, we have

$$p(t|x, \mathbf{w}, \sigma) = \mathcal{N}(t|y(x, \mathbf{w}), \sigma) \quad (5)$$

Here we only consider the case of a single real-valued variable x . Now you need to use the training data $\{\mathbf{x}, \mathbf{t}\}$ to determine the parameter \mathbf{w} and σ by maximum likelihood.

1. Show that maximizing the log likelihood is equal to minimizing the sum-of-squares error function. [1 pts]
2. More, if we assume that the polynomial coefficients \mathbf{w} is distributed as the Gaussian distribution of the form

$$p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha \mathbf{I}) \quad (6)$$

where α is the parameter of the distribution. Then what is the formulation of the prediction problem? And give us the regularization parameter. Please show us the induction of the procedure.

(Hint. Using Bayes' theorem) [1.5 pts]

$$\begin{aligned}
3. a) L(t; x, w, \sigma) &= \prod_{i=1}^N p(t_i | x_i, w, \sigma) \\
&= \prod_{i=1}^N \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(y(x_i, w) - t_i)^2}{2\sigma^2}} \\
\ln L(t; x, w, \sigma) &= \sum_{i=1}^N \ln \left(\frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(y(x_i, w) - t_i)^2}{2\sigma^2}} \right) \\
&= -\frac{N}{2} \ln \sigma^2 - \frac{N}{2} \ln(2\pi) - \frac{1}{2\sigma^2} \|xw - t\|^2 \quad x = (x_1, \dots, x_N)^T, x_i = \begin{pmatrix} 1 \\ x_i \end{pmatrix}_{i=1 \dots N}, w = \begin{pmatrix} w_0 \\ w_1 \end{pmatrix} \\
\max_w \ln L(t; x, w, \sigma) &= -\frac{N}{2} \ln \sigma^2 - \frac{N}{2} \ln(2\pi) - \frac{1}{2\sigma^2} \|xw - t\|^2 \Leftrightarrow \max_w -\frac{1}{2\sigma^2} \|xw - t\|^2 \Leftrightarrow \min_w \frac{1}{2} \|xw - t\|^2 \text{ [sum-of-square error]} \\
&\text{Q.E.D.} \quad \text{function}
\end{aligned}$$

$$\begin{aligned}
(2). \because p(w|t, x, \sigma) &= \frac{p(t|x, w, \sigma) p(w)}{p(t)} \propto p(t|x, w, \sigma) p(w) \\
L(w|t, x, \sigma) &= \prod_{i=1}^N p(t_i | w, x_i, \sigma) p(w) = \prod_{i=1}^N p(t_i | w, x_i, \sigma) \prod_{j=0}^1 p(w_j) \\
w^* &= \arg \max_w L(w|t, x, \sigma) \\
&= \arg \max_w e^{\sum_{i=1}^N \left(-\frac{(y(x_i, w) - t_i)^2}{2\sigma^2} \right)} \prod_{j=0}^1 e^{-\frac{w_j^2}{2\lambda^2}} \\
&= \arg \max_w \left[\sum_{i=1}^N \left(-\frac{(y(x_i, w) - t_i)^2}{2\sigma^2} \right) + \sum_{j=0}^1 \left(-\frac{w_j^2}{2\lambda^2} \right) \right] \\
&= \arg \min_w \|xw - t\|^2 + \frac{1}{2\lambda^2} \|w\|_2^2 \quad x = (x_1, \dots, x_N)^T, x_i = \begin{pmatrix} 1 \\ x_i \end{pmatrix}_{i=1 \dots N}, w = \begin{pmatrix} w_0 \\ w_1 \end{pmatrix} \\
&\text{which is a } L_2\text{-regularization problem}
\end{aligned}$$

Problem 4

In the following problems, we explore the the relationship between the size of the validation set and the expected error.

Let us first look at how the validation set is created. The first step is to partition the data set \mathcal{D} of fixed size N into a training set \mathcal{D}_{train} of size $(N - K)$ and a validation set \mathcal{D}_{val} of size K . We select $N - K$ points at random for training and the remaining for validation. Figure 1 depicts the relationship among \mathcal{D} , \mathcal{D}_{train} and \mathcal{D}_{val} . Suppose we have M models $\mathcal{H}_1, \dots, \mathcal{H}_M$. Given $m \in \{1, \dots, M\}$, denote a $g_m \in \mathcal{H}_m$ as the hypothesis chosen based on data set \mathcal{D} , $g_m^- \in \mathcal{H}_m$ as the hypothesis chosen based on data set \mathcal{D}_{train} as shown in Figure 2.

Figure 1: Using a validation set to estimate E_{out}

Figure 2: Using a validation set for model selection

Now we can evaluate each model on the validation set to obtain the validation errors E_1, \dots, E_M , where

$$E_m = E_{val}(g_m^-), \quad m = 1, \dots, M \quad (7)$$

The validation error can be used as an estimation of the out-of-sample error $E_{out}(g_m^-)$ for each \mathcal{H}_m . It is now a simple matter to select the model with lowest validation error. Let m^* be the index of the model which achieves the minimum validation error. So for \mathcal{H}_{m^*} , $E_{m^*} \leq E_m$ for $m = 1, \dots, M$.

Please answer the following questions:

Figure 3: The out-of-sample error $E_{out}(g_{m^*}^-)$ and the validation error $E_{val}(g_{m^*}^-)$ versus the size K of \mathcal{D}_{val} .

Problem 4a

Referring to the Figure 3, why are both curves increasing with K ? Why do they converge to each other with increasing K ? [1 pts]

Problem 4b

Referring to the Figure 4, answer the following 3 problems:

1. $\mathbb{E}[E_{out}(g_{m^*}^-)]$ is initially decreasing. How can this be, if $\mathbb{E}[E_{out}(g_m^-)]$ is increasing in K for each m ? [0.5 pts]
2. We see that $\mathbb{E}[E_{out}(g_{m^*}^-)]$ is initially decreasing, and then it starts to increase. What are possible reason for this? [0.5 pts]
3. When $K = 1$, $\mathbb{E}[E_{out}(g_{m^*}^-)] < \mathbb{E}[E_{out}(g_{m^*})]$. How can this be, if the learning curves for both models are decreasing? [0.5 pts]

4.(a). 1° Because the size of the train set : $N-k$ is decreasing when K is increasing

2° Because the size of the validation set : K is increasing, by law of large numbers $E_{val}(g_{m^*}^-) \rightarrow E_{out}(g_{m^*}^-)$, $K \rightarrow \infty$

(b). 1° Since \bar{g}_{m^*} is the smallest E_{in} among M hypotheses, $E_{in}(\bar{g}_{m^*}) \approx E_{out}(\bar{g}_{m^*})$ for small N and larger K , fixed N , $E_{in}(\bar{g}_{m^*})$ is decreasing when K is increasing, so $E[E_{out}(\bar{g}_{m^*})]$ is initially decreasing.

2° The reason for initially decreasing is discussed in (a). The reason for $E[E_{out}(\bar{g}_{m^*})]$ starts to increase when K is large is because the size of the train set : $N-k$ is decreasing when K is increasing, when K is too large, set for training is too small.

3° The reason for this is that when $K=1$, \bar{g}_{m^*} and \bar{g}_{m^*} has almost the same size of the train set to produce almost the same hypothesis, that is $\bar{g}_{m^*} \approx \bar{g}_{m^*}$, but the choice for \bar{g}_{m^*} is based on $\arg \min_{g_{m^*, N=1 \dots N}} E_{val}(g_{m^*})$ which guarantee the smallest E_{in} among M hypothesis, that means given a guarantee to not too big. $E[E_{out}(\bar{g}_{m^*})]$ [decrease the effect of noise] but $E[E_{out}(\bar{g}_{m^*})]$ doesn't have this guarantee, so when the size of train set is similar ($K \ll N$), $E[E_{out}(\bar{g}_{m^*})] < E[E_{out}(\bar{g}_{m^*})]$ can be possible.

Figure 4: Generalized error versus validation set size K . validation $g_{m^*}^-$: the model trained on \mathcal{D}_{train} with the validation set for model selection. in sample g_{m^*} : the model trained on \mathcal{D} without validation set. validation g_{m^*} : the model that is trained on \mathcal{D} at first for selecting m^* and then is retrained on \mathcal{D} with fixed m^* . The dotted line optimal is the optimal model selection, if we could select the model based on the true out of sample error.

References

- [1] Abu-Mostafa, Yaser S., 1957-. Learning From Data : a Short Course. [United States] :AMLBook.com, 2012.