

Machine Learning Exercise Sheet 04

Linear Regression

Homework

Least squares regression

Problem 1: Let's assume we have a dataset where each datapoint, (\mathbf{x}_i, y_i) is weighted by a scalar factor which we will call t_i . We will assume that $t_i > 0$ for all i . This makes the sum of squares error function look like the following:

$$E_{\text{weighted}}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N t_i [\mathbf{w}^T \phi(\mathbf{x}_i) - y_i]^2$$

Find the equation for the value of \mathbf{w} that minimizes this error function.

Furthermore, explain how this weighting factor, t_i , can be interpreted in terms of

- 1) the variance of the noise on the data and
- 2) data points for which there are exact copies in the dataset.

If we define $\mathbf{T} = \text{diag}(t_1, \dots, t_N)$ to be a diagonal matrix containing the weighting coefficients, then we can write the weighted sum-of-squares cost function in the form

$$E_{\text{weighted}}(\mathbf{w}) = \frac{1}{2} (\Phi \mathbf{w} - \mathbf{y})^T \mathbf{T} (\Phi \mathbf{w} - \mathbf{y})$$

Setting the derivative with respect to \mathbf{w} to zero, and re-arranging, then gives

$$\mathbf{w}_{\text{weighted}}^* = (\Phi^T \mathbf{T} \Phi)^{-1} \Phi^T \mathbf{T} \mathbf{y}$$

which reduces to the standard solution for the case $\mathbf{T} = \mathbf{I}$. I.e.

$$\mathbf{w}_{\text{ML}} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y}$$

If you remember back to when we modeled the likelihood using a Gaussian, our likelihood had the following form:

$$p(\mathbf{y} \mid \Phi, \mathbf{w}, \beta) = \prod_{i=1}^N \mathcal{N}(y_i \mid \mathbf{w}^T \phi(\mathbf{x}_i), \beta^{-1})$$

After applying the logarithm and using the standard form for the univariate Gaussian our equation looked like this:

$$\begin{aligned}\ln p(\mathbf{y} \mid \Phi, \mathbf{w}, \beta) &= \sum_{i=1}^N \ln \mathcal{N}(y_i \mid \mathbf{w}^T \phi(\mathbf{x}_i), \beta^{-1}) \\ &= \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta E_{\text{LS}}(\mathbf{w})\end{aligned}$$

Where $E_{\text{LS}}(\mathbf{w})$ is the standard sum of squares error function (not to be confused with the $E_{\text{weighted}}(\mathbf{w})$ we defined earlier). Remember that $E_{\text{LS}}(\mathbf{w})$ is defined as follows:

$$E_{\text{LS}}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (\mathbf{w}^T \phi(\mathbf{x}_i) - y_i)^2$$

When we compare E_{LS} with E_{weighted} and the effect of swapping the two in the previous likelihood equation we can see that T_i can be regarded as a precision (inverse variance) parameter, particular to the data point (\mathbf{x}_i, y_i) , that either replaces or scales β .

Alternatively, t_i can be regarded as an *effective* number of replicated observations of data point (\mathbf{x}_i, y_i) ; this becomes particularly clear if we consider $E_{\text{weighted}}(\mathbf{w})$ with t_i taking positive integer values, although it is valid for any $t_i > 0$.

Ridge regression

Problem 2: Show that the following holds: The ridge regression estimates can be obtained by ordinary least squares regression on an augmented dataset: Augment the design matrix $\Phi \in \mathbb{R}^{N \times M}$ with M additional rows $\sqrt{\lambda} \mathbf{I}_{M \times M}$ and augment \mathbf{y} with M zeros.

Ordinary least squares minimizes $(\mathbf{y} - \Phi \mathbf{w})^T (\mathbf{y} - \Phi \mathbf{w})$. For ridge regression we need to minimize $(\mathbf{y} - \Phi \mathbf{w})^T (\mathbf{y} - \Phi \mathbf{w}) + \lambda \mathbf{w}^T \mathbf{w}$. If we define $\hat{\Phi} = \begin{pmatrix} \Phi \\ \sqrt{\lambda} \mathbf{I} \end{pmatrix}$ and $\hat{\mathbf{y}} = \begin{pmatrix} \mathbf{y} \\ \mathbf{0}_M \end{pmatrix}$, we can formulate the ridge regression objective as minimizing $(\hat{\mathbf{y}} - \hat{\Phi} \mathbf{w})^T (\hat{\mathbf{y}} - \hat{\Phi} \mathbf{w})$.

Problem 3: Derive the closed form solution for ridge regression error function

$$E_{\text{ridge}}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (\mathbf{w}^T \phi(\mathbf{x}_i) - y_i)^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

Additionally, discuss the scenario when the number of training samples N is smaller than the number of basis functions M . What computational issues arise in this case? How does regularization address them?

$$\begin{aligned} E_{\text{ridge}}(\mathbf{w}) &= \frac{1}{2} \sum_{i=1}^N (\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_i) - y_i)^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} \\ &= \frac{1}{2} (\boldsymbol{\Phi} \mathbf{w} - \mathbf{y})^T (\boldsymbol{\Phi} \mathbf{w} - \mathbf{y}) + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} \end{aligned}$$

Taking the gradient

$$\begin{aligned} \nabla_{\mathbf{w}} E_{\text{ridge}}(\mathbf{w}) &= \boldsymbol{\Phi}^T \boldsymbol{\Phi} \mathbf{w} - \boldsymbol{\Phi}^T \mathbf{y} + \lambda \mathbf{w} \\ &= (\boldsymbol{\Phi}^T \boldsymbol{\Phi} + \lambda \mathbf{I}) \mathbf{w} - \boldsymbol{\Phi}^T \mathbf{y} \end{aligned}$$

Set it to zero

$$\begin{aligned} (\boldsymbol{\Phi}^T \boldsymbol{\Phi} + \lambda \mathbf{I}) \mathbf{w} &= \boldsymbol{\Phi}^T \mathbf{y} \\ \mathbf{w} &= (\boldsymbol{\Phi}^T \boldsymbol{\Phi} + \lambda \mathbf{I})^{-1} \boldsymbol{\Phi}^T \mathbf{y} \end{aligned}$$

If $N < M$ the covariance matrix $\boldsymbol{\Phi}^T \boldsymbol{\Phi} \in \mathbb{R}^{M \times M}$ will be singular, therefore not invertible. (this may happen even if $N \geq M$, e.g. when some features are correlated).

When regularization is used, $\lambda \mathbf{I}$ is added to the covariance matrix, thus fixing the potential degeneracy issue and making the problem tractable.

Multi-output linear regression

Problem 4: In class, we only considered functions of the form $f : \mathbb{R}^n \rightarrow \mathbb{R}$. What about the general case of $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$? For linear regression with multiple outputs, write down the loglikelihood formulation and derive the MLE of the parameters.

The observation \mathbf{y}_i is a vector with $\mathbf{y}_i \sim \mathcal{N}(\mathbf{W}\mathbf{x}_i, \Sigma)$, $\mathbf{W} \in \mathbb{R}^{n \times m}$, $\Sigma \in \mathbb{R}^{m \times m}$, covariance Σ is known and fixed for all possible observations. For n i.i.d observed pairs $(\mathbf{x}_i, \mathbf{y}_i)$, the likelihood is

$$p(\mathbf{y}_i | \mathbf{W}, \Sigma) = \prod_i \mathcal{N}(\mathbf{y}_i | \mathbf{W}^T \mathbf{x}_i, \Sigma),$$

and thus the negative log-likelihood is

$$\ln p(\mathbf{y}_i | \mathbf{W}, \Sigma) = \text{const} + \frac{1}{2} \sum_i (\mathbf{y}_i - \mathbf{W}^T \mathbf{x}_i)^T \Sigma^{-1} (\mathbf{y}_i - \mathbf{W}^T \mathbf{x}_i)$$

Since Σ is symmetric and positive semi-definite we can use the Cholesky decomposition $\Sigma = \mathbf{L}\mathbf{L}^T$ to obtain $\Sigma^{-1} = \mathbf{L}^{-T}\mathbf{L}^{-1}$. With this we can simplify the two factors via

$$\mathbf{L}^{-1}(\mathbf{y}_i - \mathbf{W}^T \mathbf{x}_i) = (\mathbf{L}^{-1}\mathbf{y}_i - \mathbf{L}^{-1}\mathbf{W}^T \mathbf{x}_i) = (\hat{\mathbf{y}}_i - \hat{\mathbf{W}}^T \mathbf{x}_i).$$

Using this transformation we get

$$\begin{aligned} \ln p(\mathbf{y}_i | \mathbf{W}, \Sigma) &= \text{const} + \frac{1}{2} \sum_i (\hat{\mathbf{y}}_i - \hat{\mathbf{W}}^T \mathbf{x}_i)^T (\hat{\mathbf{y}}_i - \hat{\mathbf{W}}^T \mathbf{x}_i) \\ &= \text{const} + \frac{1}{2} \sum_i \sum_j (\hat{\mathbf{y}}_i - \hat{\mathbf{W}}^T \mathbf{x}_i)_j (\hat{\mathbf{y}}_i - \hat{\mathbf{W}}^T \mathbf{x}_i)_j \\ &= \text{const} + \frac{1}{2} \sum_i \sum_j (\hat{\mathbf{Y}}_{ij} - (\mathbf{X}\hat{\mathbf{W}})_{ij}) (\hat{\mathbf{Y}}_{ij} - (\mathbf{X}\hat{\mathbf{W}})_{ij}) \\ &= \text{const} + \frac{1}{2} \sum_i \sum_j \left((\hat{\mathbf{Y}} - \mathbf{X}\hat{\mathbf{W}})^T \right)_{ji} (\hat{\mathbf{Y}} - \mathbf{X}\hat{\mathbf{W}})_{ij} \\ &= \text{const} + \frac{1}{2} \sum_j \left((\hat{\mathbf{Y}} - \mathbf{X}\hat{\mathbf{W}})^T (\hat{\mathbf{Y}} - \mathbf{X}\hat{\mathbf{W}}) \right)_{jj} \\ &= \text{const} + \frac{1}{2} \text{Tr} \left[(\hat{\mathbf{Y}} - \mathbf{X}\hat{\mathbf{W}})^T (\hat{\mathbf{Y}} - \mathbf{X}\hat{\mathbf{W}}) \right], \end{aligned}$$

where $\hat{\mathbf{Y}}$ and \mathbf{X} are matrices that have the vectors $\hat{\mathbf{y}}_i$ and \mathbf{x}_i as their rows. Note that the trace essentially just denotes the sum over all dimensions of the output domain. Taking the derivative of the negative log-likelihood with respect to $\hat{\mathbf{W}}$ and setting it to 0 gives $\hat{\mathbf{W}}_{\text{MLE}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \hat{\mathbf{Y}}$. So these are m single least square problems, one for every *column* of $\hat{\mathbf{Y}}$. Finally, transforming $\hat{\mathbf{W}}_{\text{MLE}}$ back gives $\mathbf{W}_{\text{MLE}} = \hat{\mathbf{W}}_{\text{MLE}} \mathbf{L}^T$.

Comparison of Linear Regression Models

Problem 5: We want to perform regression on a dataset consisting of N samples $\mathbf{x}_i \in \mathbb{R}^D$ with corresponding targets $y_i \in \mathbb{R}$ (represented compactly as $\mathbf{X} \in \mathbb{R}^{N \times D}$ and $\mathbf{y} \in \mathbb{R}^N$).

Assume that we have fitted an L_2 -regularized linear regression model and obtained the optimal weight vector $\mathbf{w}^* \in \mathbb{R}^D$ as

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \frac{1}{2} \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}_i - y_i)^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

Note that there is no bias term.

Now, assume that we obtained a new data matrix \mathbf{X}_{new} by scaling all samples by the same positive factor $a \in (0, \infty)$. That is, $\mathbf{X}_{new} = a\mathbf{X}$ (and respectively $\mathbf{x}_i^{new} = a\mathbf{x}_i$).

- a) Find the weight vector \mathbf{w}_{new} that will produce the same predictions on \mathbf{X}_{new} as \mathbf{w}^* produces on \mathbf{X} .

Predictions of a linear regression model are generated as $\hat{y} = \mathbf{w}^T \mathbf{x}$.

This means that we need to ensure that $\mathbf{w}^{*T} \mathbf{x}_i = \mathbf{w}_{new}^T \mathbf{x}_i^{new}$ or equivalently $\mathbf{w}^{*T} \mathbf{x}_i = \mathbf{w}_{new}^T a \mathbf{x}_i$. Solving for \mathbf{w}_{new} we get $\mathbf{w}_{new} = \frac{\mathbf{w}^*}{a}$.

- b) Find the regularization factor $\lambda_{new} \in \mathbb{R}$, such that the solution \mathbf{w}_{new}^* of the new L_2 -regularized linear regression problem

$$\mathbf{w}_{new}^* = \arg \min_{\mathbf{w}} \frac{1}{2} \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}_i^{new} - y_i)^2 + \frac{\lambda_{new}}{2} \mathbf{w}^T \mathbf{w}$$

will produce the same predictions on \mathbf{X}_{new} as \mathbf{w}^* produces on \mathbf{X} .

Provide a mathematical justification for your answer.

The closed form solution for \mathbf{w}^* on the original data \mathbf{X} is

$$\mathbf{w}^* = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

The closed form solution for \mathbf{w}_{new}^* on the new data \mathbf{X}_{new} is

$$\begin{aligned} \mathbf{w}_{new}^* &= (\mathbf{X}_{new}^T \mathbf{X}_{new} + \lambda_{new} \mathbf{I})^{-1} \mathbf{X}_{new}^T \mathbf{y} \\ &= a(a^2 \mathbf{X}^T \mathbf{X} + \lambda_{new} \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \end{aligned}$$

by setting $\lambda_{new} = a^2 \lambda$, we get

$$\begin{aligned} &= a(a^2 \mathbf{X}^T \mathbf{X} + a^2 \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \\ &= \frac{1}{a} (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \\ &= \frac{1}{a} \mathbf{w}^* \end{aligned}$$

Which (according to our answer in part (a) of this problem) will produce the same predictions on \mathbf{X}_{new} as \mathbf{w}^* does on \mathbf{X} , as desired.

Equivalent solution

$$\begin{aligned}
 \mathbf{w}_{new}^* &\stackrel{!}{=} \frac{\mathbf{w}^*}{a} = \frac{1}{a} \arg \min_{\mathbf{w}} \frac{1}{2} \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}_i - y_i)^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} \\
 &= \frac{1}{a} \arg \min_{\mathbf{w}} \frac{1}{2} \sum_{i=1}^N \left(\frac{\mathbf{w}^T}{a} a \mathbf{x}_i - y_i \right)^2 + \frac{a^2 \lambda}{2} \frac{\mathbf{w}^T}{a} \frac{\mathbf{w}}{a} \\
 &= \frac{a}{a} \arg \min_{\mathbf{w}_{new} = \frac{\mathbf{w}}{a}} \frac{1}{2} \sum_{i=1}^N (\mathbf{w}_{new}^T \mathbf{x}_i^{new} - y_i)^2 + \frac{a^2 \lambda}{2} \mathbf{w}_{new}^T \mathbf{w}_{new} \\
 &\stackrel{!}{=} \mathbf{w}_{new}^* = \arg \min_{\mathbf{w}_{new}} \frac{1}{2} \sum_{i=1}^N (\mathbf{w}_{new}^T \mathbf{x}_i^{new} - y_i)^2 + \frac{\lambda_{new}}{2} \mathbf{w}_{new}^T \mathbf{w}_{new}
 \end{aligned}$$

For this equality to hold we need to match the regularization term by setting $\lambda_{new} = a^2 \lambda$.

Programming Task: Least squares regression

Problem 6: Load the notebook `04_homework_linear_regression.ipynb` from Piazza. Fill in the missing code and run the notebook. Convert the evaluated notebook to pdf and add it to the printout of your homework.

Note: We suggest that you use Anaconda for installing Python and Jupyter, as well as for managing packages. We recommend that you use Python 3.

For more information on Jupyter notebooks and how to convert them to other formats, consult the Jupyter documentation and nbconvert documentation.

In-class Exercises

Problem 7: Assume that we are given a dataset, where each sample x_i and regression target y_i is generated according to the following process

$$x_i \sim \text{Uniform}(-10, 10)$$

$$y_i = ax_i^3 + bx_i^2 + cx_i + d + \epsilon_i, \quad \text{where} \quad \epsilon_i \sim \mathcal{N}(0, 1) \quad \text{and} \quad a, b, c, d \in \mathbb{R}.$$

The 3 regression algorithms below are applied to the given data. Your task is to say what the bias and variance of these models are (low or high). Provide a 1-2 sentence explanation to each of your answers.

a) Linear regression

Bias: high. Variance: low.

A straight line cannot capture a degree 3 polynomial (thus underfitting the data).

b) Polynomial regression with degree 3

Bias: low. Variance: low.

The model is same as the data generating process. We can achieve a good fit.

c) Polynomial regression with degree 10

Bias: low. Variance: high.

Since we are using a polynomial regression with a degree much higher compared to the data generating process, the model will overfit the data.

Problem 8: Given is a training set consisting of samples $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$ with respective regression targets $\mathbf{y} = \{y_1, y_2, \dots, y_N\}$ where $\mathbf{x}_i \in \mathbb{R}^D$ and $y_i \in \mathbb{R}$.

Alice fits a linear regression model $f(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i$ to the dataset using the closed form solution for linear regression (normal equations).

Bob has heard that by transforming the inputs \mathbf{x}_i with a vector-valued function ϕ , he can fit an alternative function, $g(\mathbf{x}_i) = \mathbf{v}^T \phi(\mathbf{x}_i)$, using the same procedure (solving the normal equations). He decides to use a linear transformation $\phi(\mathbf{x}_i) = \mathbf{A}^T \mathbf{x}_i$, where $\mathbf{A} \in \mathbb{R}^{D \times D}$ has full rank.

a) Show that Bob's procedure will fit the same function as Alice's original procedure, that is $f(\mathbf{x}) = g(\mathbf{x})$ for all $\mathbf{x} \in \mathbb{R}^D$ (given that \mathbf{w} and \mathbf{v} minimize the training set error).

By fitting $g(\mathbf{x}) = \mathbf{v}^T \phi(\mathbf{x})$, Bob can only set the function equal to linear combinations of the inputs $(\mathbf{v}^T \mathbf{A}^T) \mathbf{x} = \mathbf{w}^T \mathbf{x}$, where $\mathbf{w} = \mathbf{A} \mathbf{v}$.

Moreover, just like Alice, all linear combinations are available: any function Alice fits can be matched by setting $\mathbf{v} = \mathbf{A}^{-1} \mathbf{w}$.

As both Alice and Bob are selecting the function that best matches the outputs from the same set of functions, and with the same cost function, they will select the same function.

- b) Can Bob's procedure lead to a lower training set error than Alice's if the matrix \mathbf{A} is not invertible? Explain your answer.

Since \mathbf{A} is square and not invertible, then multiple input vectors are transformed to the same output vector. It's not possible for Bob to assign different function values to two such inputs, whereas Alice can. Bob can no longer fit all the same functions as Alice. Furthermore, Alice can still match any of Bob's solutions by setting $\mathbf{w} = \mathbf{A}\mathbf{v}$. Since the error function is quadratic and PSD it has a unique (set of) optima. So either Alice's (set of) optima will include Bob's solution or it will be lower. In summary, Bob's training error might be worse than Alice's, but it can't be better.

Note that we are only talking about training error in this example, not test error. Bob might manage to find a model that generalizes better than Alice's, but Alice will always be able to fit the training data at least as well as Bob.

Problem 9: See Jupyter notebook `inclass_04_notebook.ipynb`.