

Machine Learning Exercise Sheet 04

Linear Regression

Exercise sheets consist of two parts: In-class exercises and homework. The in-class exercises will be solved and discussed during the tutorial. The homework is for you to solve at home and further engage with the lecture content. There is no grade bonus and you do not have to upload any solutions. Note that the order of some exercises might have changed compared to last year's recordings.

In-class Exercises

Problem 1: 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
- b) Polynomial regression with degree 3
- c) Polynomial regression with degree 10

Problem 2: Given is a training set consisting of samples $\mathbf{X} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)^T$ with respective regression targets $\mathbf{y} = (y_1, y_2, \dots, y_N)^T$ 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).
- 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.

Problem 3: See Jupyter notebook `inclass_04_notebook.ipynb`.

P1

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

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

underfitting

- ① Linear regression: High bias / Low variance \Leftrightarrow Linear model is too rigid to represent the Uniform and high-dim data
- ② Polynomial 3: Low bias / Low variance \Leftrightarrow because y_i is also a polynomial with 3 dim
- ③ Polynomial 10: Low bias / High variance \Leftrightarrow Polynomial with 10 degree model is too flexible to fit the y_i , large fit noise

P2 $X = (x_1, \dots, x_n)^T$ $y = (y_1, \dots, y_n)^T$ $x_i \in \mathbb{R}^{p \times 1}$ $y_i \in \mathbb{R}$. $X \in \mathbb{R}^{n \times p}$

A $f(x_i) = w^T x_i$

B $g(x_i) = v^T A^T x_i$

$$E_{LS} = \frac{1}{2} \sum_{i=1}^N (w^T x_i - y_i)^2$$

$$= \frac{1}{2} (XW - Y)^T (XW - Y)$$

$$E_{LS} = \frac{1}{2} \sum_{i=1}^N (v^T A^T x_i - y_i)^2$$

$$= \frac{1}{2} (XAV - Y)^T (XAV - Y)$$

$$g(x_i) = w^{*T} A^{-T} A^T x_i$$

$$= w^{*T} x_i$$

$$f(x) = w^{*T} x_i$$

$$\therefore f(x) = g(x)$$

$$w^* = (X^T X)^{-1} X^T y$$

$$\frac{1}{2} (A^T V^T X^T X V A - 2 A^T V^T X^T y + Y^T Y)$$

$$A^T X^T X A V - A^T X^T y \stackrel{!}{=} 0$$

$$V^* = (A^T X^T X A)^{-1} A^T X^T y$$

$$= A^{-1} (X^T X)^{-1} A^{-T} A^T X^T y$$

$$= A^{-1} (X^T X)^{-1} X^T y$$

$$= A^T w^*$$

Because $V^* = A^{-1} w^*$

so V^* is a subset of w^*

can't achieve a lower error.

P4

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

$$\frac{1}{2} [w^T \Phi^T \Phi w - 2 w^T \Phi^T y + y^T y]$$

$$w^* = \underset{w}{\text{argmin}} E_{\text{weighted}}(w) = \underset{w}{\text{argmin}} \frac{1}{2} [\Phi w - y]^T [\Phi w - y]$$

$$= \Phi^T \Phi w - \Phi^T y \stackrel{!}{=} 0$$

$$w^* = (\Phi^T \Phi)^{-1} \Phi^T y$$

assume $T = \text{diag}(t_1, \dots, t_n)$

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

$$= \frac{1}{2} [w^T \Phi^T T \Phi w - y^T T y]$$

$$= \frac{1}{2} [w^T \Phi^T T \Phi w - 2 w^T \Phi^T T y + y^T T y]$$

$$\frac{dE_{\text{weighted}}}{dw} = \Phi^T T \Phi w - \Phi^T T y \stackrel{!}{=} 0$$

$$w^* = (\Phi^T T \Phi)^{-1} \Phi^T T y$$

1) $y_i \sim \mathcal{N}(w^T \phi(x_i), \beta^{-1})$

$$P(Y | x, w, \beta) = \prod_{i=1}^N p(y_i | x, w, \beta)$$

$$\underset{w}{\text{argmax}} P(Y | x, w, \beta) = \underset{w}{\text{argmin}} - \ln \prod_{i=1}^N p(y_i | x, w, \beta) = E_{\text{weighted}}(w)$$

$$= - \sum_{i=1}^N \ln \sqrt{\frac{\beta}{2\pi}} + \frac{1}{2} \sum_{i=1}^N \beta (w^T \phi(x_i) - y_i)^2$$

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

$$\Rightarrow \beta = t_i \Rightarrow y_i \sim \mathcal{N}(y_i | w^T \phi(x_i), t_i^{-1})$$

t_i seen as the effective number of replicated observation



$$p(x | y = c) = \mathcal{N}(x | \mu_c, \Sigma)$$

(3)

$$= \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu_c)^T \Sigma^{-1} (x - \mu_c) \right\} \quad (4)$$

$$p_j \Phi \in \mathbb{R}^{N \times M} \quad \lambda$$

$$\hat{\Phi} = \begin{pmatrix} \Phi \\ \sqrt{\lambda} I_M \end{pmatrix} \quad \text{and} \quad \hat{y} = \begin{pmatrix} y \\ 0_M \end{pmatrix}.$$

Ridge regression

$$y_i = w^T \phi(x_i)$$

$$\begin{aligned} E_{LS}(w) &= \frac{1}{2} \sum_{i=1}^N (w^T \phi(x_i) - y_i)^2 + \frac{\lambda}{2} \|w\|_2^2 \\ &= \frac{1}{2} [\Phi w - y]^T [\Phi w - y] + \frac{\lambda}{2} \|w\|_2^2 \\ &= \frac{1}{2} [w^T \Phi^T \Phi w - 2w^T \Phi^T y + y^T y] + \frac{\lambda}{2} \|w\|_2^2 \end{aligned}$$

$$\begin{aligned} \frac{dE_{LS}}{dw} &= \Phi^T \Phi w - \Phi^T y + \lambda w \stackrel{!}{=} 0 \\ (\Phi^T \Phi + \lambda I) w^* &= \Phi^T y \\ w^* &= (\Phi^T \Phi + \lambda I)^{-1} \Phi^T y \end{aligned}$$

Linear Regression

$$y_i = w^T x_i$$

$$\hat{y} = w^T \hat{\Phi}(x_i)$$

$$\begin{aligned} E_{LS} &= \frac{1}{2} \sum_{i=1}^N (w^T \phi(x_i) - y_i)^2 \\ &= \frac{1}{2} \left[\begin{pmatrix} \Phi \\ \sqrt{\lambda} I_M \end{pmatrix} w - \begin{pmatrix} y \\ 0_M \end{pmatrix} \right]^T \left[\begin{pmatrix} \Phi \\ \sqrt{\lambda} I_M \end{pmatrix} w - \begin{pmatrix} y \\ 0_M \end{pmatrix} \right] \\ &= \frac{1}{2} \left[w^T \begin{pmatrix} \Phi^T \\ \sqrt{\lambda} I_M \end{pmatrix} \begin{pmatrix} \Phi \\ \sqrt{\lambda} I_M \end{pmatrix} w - 2w^T \begin{pmatrix} \Phi^T \\ \sqrt{\lambda} I_M \end{pmatrix} \begin{pmatrix} y \\ 0_M \end{pmatrix} + c \right] \\ \frac{dE_{LS}}{dw} &= \begin{pmatrix} \Phi^T \\ \sqrt{\lambda} I_M \end{pmatrix} \begin{pmatrix} \Phi \\ \sqrt{\lambda} I_M \end{pmatrix} w - \begin{pmatrix} \Phi^T \\ \sqrt{\lambda} I_M \end{pmatrix} \begin{pmatrix} y \\ 0_M \end{pmatrix} \stackrel{!}{=} 0 \\ (\Phi^T \Phi + \lambda I_M) \cdot w &= \Phi^T y \\ w^* &= (\Phi^T \Phi + \lambda I_M)^{-1} \Phi^T y \end{aligned}$$

p6 Because the training error will decrease if the model (with high degree polynomial) perfectly fit all the data

Bart is 49

Split data into training and val data
then compute the validation error

$$p_7 \quad P(w, \beta | D) \propto P(Y | \Phi, w, \beta) P(w | \beta)$$

$$\begin{aligned} \log P(w, \beta | D) &= \sum_{i=1}^N \ln \left[\sqrt{\frac{\beta}{2\pi}} \exp \left(-\frac{\beta}{2} (w^T \phi(x_i) - y_i)^2 \right) \right] + \sum_{i=1}^M \ln \left[\sqrt{\frac{\beta}{2\pi}} \sqrt{\beta^{-1} s_0} = \sigma \right. \\ &\quad \left. + \ln \frac{b_0^{a_0}}{\Gamma(a_0)} + (a_0 - 1) \ln \beta - b_0 \beta \right] \end{aligned}$$

$$N = \frac{1}{\sqrt{2\pi}\sigma} \exp \left(-\frac{1}{2\sigma^2} (x - \mu)^2 \right)$$

$$\frac{1}{\sqrt{2\pi}} \cdot \frac{1}{\sqrt{\beta^{-1} s_0}}$$

$$\log \sqrt{\frac{\beta}{2\pi \cdot s_0}}$$

$$\begin{aligned} &= \frac{1}{2} \log \left(\frac{\beta}{2\pi \cdot s_0} \right) \\ &= \frac{1}{2} \left[\log \beta - \frac{\log 2\pi}{c} - \log s_0 \right] \end{aligned}$$

$$P(w, \beta | D) = N(w | m_N, \beta^{-1} s_N) \text{Gamma}(\beta | a_N, b_N)$$

Homework

Least squares regression

Problem 4: 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.

Ridge regression

Problem 5: Show that ridge regression on a design matrix $\Phi \in \mathbb{R}^{N \times M}$ with regularization strength λ is equivalent to ordinary least squares regression with an augmented design matrix and target vector

$$\hat{\Phi} = \begin{pmatrix} \Phi \\ \sqrt{\lambda} I_M \end{pmatrix} \quad \text{and} \quad \hat{\mathbf{y}} = \begin{pmatrix} \mathbf{y} \\ \mathbf{0}_M \end{pmatrix}.$$

Implementation

Problem 6: John Doe is a data scientist, and he wants to fit a polynomial regression model to his data. For this, he needs to choose the degree of the polynomial that works best for his problem.

Unfortunately, John hasn't attended IN2064, so he writes the following code for choosing the optimal degree of the polynomial:

```
X, y = load_data()
best_error = -1
best_degree = None

for degree in range(1, 50):
    w = fit_polynomial_regression(X, y, degree)
    y_predicted = predict_polynomial_regression(X, w, degree)
    error = compute_mean_squared_error(y, y_predicted)
    if (error <= best_error) or (best_error == -1):
        best_error = error
        best_degree = degree

print("Best degree is " + str(best_degree))
```

Assume that the functions are implemented correctly and do what their name suggests.

- a) Explain briefly why this code doesn't do what it's supposed to do.

- b) Describe a possible way to fix the problem with this code. (You don't need to write any code, just describe the approach.)

Bayesian linear regression

Bishop 3.12

In the lecture we made the assumption that we already knew the precision (inverse variance) for our Gaussian distributions. What about when we don't know the precision and we need to put a prior on that as well as our Gaussian prior that we already have on the weights of the model?

Problem 7: It turns out that the conjugate prior for the situation when we have an unknown mean and unknown precision is a normal-gamma distribution (See section 2.3.6 in Bishop). This means that if our likelihood is as follows:

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

Then the conjugate prior for both \mathbf{w} and β is

$$p(\mathbf{w}, \beta) = \mathcal{N}(\mathbf{w} \mid \mathbf{m}_0, \beta^{-1} \mathbf{S}_0) \text{Gamma}(\beta \mid a_0, b_0)$$

Show that the posterior distribution takes the same form as the prior, i.e.

$$p(\mathbf{w}, \beta \mid \mathcal{D}) = \mathcal{N}(\mathbf{w} \mid \mathbf{m}_N, \beta^{-1} \mathbf{S}_N) \text{Gamma}(\beta \mid a_N, b_N)$$

Also be sure to give the expressions for \mathbf{m}_N , \mathbf{S}_N , a_N , and b_N .

Hint: Expand $\log p(\mathbf{w}, \beta \mid \mathcal{D})$ once with the prior and likelihood and once with the presumed posterior form. The resulting expressions have to be equal, so you should be able to match all terms in the two expansions against each other and then read off the parameters of the posterior distribution.

Problem 8: 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?

Comparison of Linear Regression Models

Problem 9: 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

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

pg 8 $E_{\text{ridge}} = \frac{1}{2} \sum_{i=1}^N (w^T \phi(x_i) - y_i)^2 + \frac{\lambda}{2} w^T w$

argmin $\frac{dE}{dw} = (\Phi^T \Phi + \lambda I) w = \Phi^T y$
 $w^* = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T y$

\uparrow
 $\in \mathbb{R}^{N \times M}$

$\Phi^T \Phi \in N \times M$

If $N < M$ not invertible
singularity

add λI to eliminate the

pg ① $\hat{y} = w^* x_i = w_{\text{non}}^* x_{i, \text{non}}$

$X w^* = X_{\text{non}} w_{\text{non}}^* = a X_{\text{non}} w_{\text{non}}^*$

$w_{\text{non}}^* = \frac{1}{a} w^*$

② $w_{\text{non}}^* = a (a^2 X^T X + \lambda_{\text{non}} I)^{-1} X^T y$

$= a (a^2 x^T x + a^2 \lambda I)^{-1} x^T y$

$= \frac{1}{a} (x^T x + \lambda I)^{-1} x^T y$

$= \frac{1}{a} w^* = w_{\text{non}}^*$

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} .
- 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.

Programming Task

Problem 10: Download the notebook `exercise_04_linear_regression.ipynb` from Moodle. Fill in the missing code and follow the instructions in the notebook to append the solution to your PDF submission.

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.