

4. Linear Regression

Mittwoch, 13. November 2019 11:31

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

$$x_i \sim \text{Uniform}(-10, 10)$$
$$y_i = ax_i^3 + bx_i^2 + cx_i + d + \epsilon_i, \quad \text{where } \epsilon_i \sim \mathcal{N}(0, 1) \quad \text{and } a, b, c, d \in \mathbb{R}.$$

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

Bias: How well can model fit data?

Variance: How much does prediction change when data changes (a little)?

Bias-variance trade-off

a) Bias high, low variance

Straight line is underfitting the data

b) Bias low, variance low

Same model as generating process

c) Bias low, variance high

Overfitting the data, model too powerful

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).

$$g(\vec{x}) = \vec{v}^T \phi(\vec{x}) = \vec{v}^T (\vec{A}^T \vec{x}) = \underbrace{(\vec{A} \vec{v})^T}_{\vec{w}} \vec{x} = \vec{w}^T \vec{x}$$
$$\vec{v} = \vec{A}^{-1} \vec{w}$$

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

A not invertible, so any 2 vector could be transformed to the same point, so Bob cannot assign different values to these \Rightarrow Bob's model at most as powerful as Alice's. $\Rightarrow \text{NO}$

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

Homework

1 Least squares regression

Problem 1: Let's assume we have a dataset where each datapoint, (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:

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

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.

$$E = \frac{1}{2} (\phi \vec{w} - \vec{y})^T \mathbb{T} (\phi \vec{w} - \vec{y})$$

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

$$\vec{w}^* = (\phi^T \mathbb{T} \phi)^{-1} \phi^T \mathbb{T} \vec{y}$$

$$p(\vec{y} | \phi, \vec{w}, \beta) = \prod_i \mathcal{N}(y_i | \vec{w}^T \phi(\vec{x}_i), \beta^{-1})$$

\nwarrow precision

$$\ln p(\vec{y} | \phi, \vec{w}, \beta) = \sum_i \ln \mathcal{N}(y_i | \vec{w}^T \phi(\vec{x}_i), \beta^{-1}) = \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) \cancel{\beta E_{LS}(\vec{w})}$$

$$\beta E_{LS}(\vec{w}) = \frac{1}{2} \sum_i \beta (\vec{w}^T \phi(\vec{x}_i) - y_i)^2$$

\downarrow
 $t_i \rightarrow$ precision per sample!

t_i is the effective number of replicated observations

2 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.

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?

3 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.

$$\begin{aligned} \vec{y}_i &\sim \mathcal{N}(\mathbf{w} \vec{x}_i, \Sigma) \\ p(Y | W, \Sigma) &= \prod_i \mathcal{N}(\mathbf{w} \vec{x}_i, \Sigma) \\ \ln p(Y | W, \Sigma) &= \text{const} + \sum_i \frac{1}{2} (\vec{y}_i - \mathbf{w} \vec{x}_i)^T \Sigma^{-1} (\vec{y}_i - \mathbf{w} \vec{x}_i) = \\ \text{Cholesky decomposition: } \Sigma &= LL^T, \quad \Sigma^{-1} = L^{-T} L^{-1} \quad (\text{since } \Sigma \text{ is symmetric and PSD}) \\ &= C + \sum_i \frac{1}{2} (L^{-1}(\vec{y}_i - \mathbf{w} \vec{x}_i))^T (L^{-1}(\vec{y}_i - \mathbf{w} \vec{x}_i)) = \\ &= \left[L^{-1}(\vec{y}_i - \mathbf{w} \vec{x}_i) = L^{-1}\vec{y}_i - L^{-1}\mathbf{w} \vec{x}_i = \vec{y}_i - \hat{W} \vec{x}_i \right] \\ &= C + \sum_i \frac{1}{2} (\vec{y}_i - \hat{W} \vec{x}_i)^T (\vec{y}_i - \hat{W} \vec{x}_i) = \\ &= C + \frac{1}{2} \text{Tr}(\hat{Y} - \hat{X} \hat{W})^T (\hat{Y} - \hat{X} \hat{W}) \\ \hat{W}_{\text{MLE}} &= (\hat{X}^T \hat{X})^{-1} \hat{X}^T \hat{Y}, \quad W_{\text{MLE}} = L \hat{W}_{\text{MLE}} \end{aligned}$$

$$\begin{aligned} \text{Tr}[(Y - XW)^T (Y - XW)] &= \\ = \sum_j \sum_i (Y_{ij} - X_{ij}W_{ij})^2 &= \\ = \sum_j \sum_i (Y_{ij} - X_{ij}W_{ij})_{ij} &= \\ = \sum_j \sum_i (Y_{ij} - [XW]_{ij})^2 &= \\ \text{samples} \quad \text{dimensionality of } \vec{y} &= \\ = \sum_i (\vec{y}_i - \vec{x}_i^T W)^T (\vec{y}_i - \vec{x}_i^T W) & \end{aligned}$$

4 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, $\mathbf{X}_{\text{new}} = a\mathbf{X}$ (and respectively $\mathbf{x}_i^{\text{new}} = a\mathbf{x}_i$).

- 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} .

$$\begin{aligned} \vec{y} &= \vec{w}^T \vec{x}_i = \vec{w}_{\text{new}}^T \vec{x}_i^{\text{new}} = \vec{w}_{\text{new}}^T \alpha \vec{x}_i \\ \Leftrightarrow \vec{w}_{\text{new}} &= \frac{\vec{w}}{\alpha} \end{aligned}$$

- 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.

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

$$\begin{aligned} \mathbf{w}_{new}^* &= (\mathbf{X}_{new}^T \mathbf{X}_{new} + \lambda_{new} \mathbf{I})^{-1} \mathbf{X}_{new}^T \vec{y} = \\ &= a (\mathbf{X}^T \mathbf{X} + \lambda_{new} \mathbf{I})^{-1} \mathbf{X}^T \vec{y} = \\ &= \underbrace{\frac{1}{a} (\mathbf{X}^T \mathbf{X} + \boxed{\lambda_{new} \mathbf{I}})^{-1} \mathbf{X}^T \vec{y}}_{\mathbf{w}^*, \text{ because } \mathbf{w}_{new}^* = \frac{\mathbf{w}^*}{a}} \\ &\quad \lambda = \frac{\lambda_{new}}{a^2} \Leftrightarrow \lambda_{new} = a^2 \lambda \end{aligned}$$

II) $\mathbf{w}_{new}^* \stackrel{!}{=} \frac{\mathbf{w}^*}{a} = \frac{1}{a} \underset{\mathbf{w}}{\operatorname{argmin}} \frac{1}{2} \sum_i (\vec{w}^T \vec{x}_i - y_i)^2 + \frac{\lambda}{2} \vec{w}^T \vec{w} =$

$$\begin{aligned} &= \boxed{\frac{1}{a}} \underset{\vec{w}}{\operatorname{argmin}} \frac{1}{2} \sum_i \left(\frac{\vec{w}^T \vec{x}_i}{a} - y_i \right)^2 + \frac{a^2 \lambda}{2} \left(\frac{\vec{w}^T \vec{w}}{a} \right) = \\ &= \boxed{a} \underset{\substack{\vec{w}_{new} = \frac{\vec{w}}{a} \\ \text{result: } \vec{w}}}{\operatorname{argmin}} \frac{1}{2} \sum_i (\vec{w}_{new}^T \vec{x}_i - y_i)^2 + \boxed{\frac{a^2 \lambda}{2}} \vec{w}_{new}^T \vec{w}_{new} \\ &\quad \text{result: } \boxed{\frac{\vec{w}}{a}} \qquad \qquad \qquad \lambda_{new} = a^2 \lambda \end{aligned}$$

5 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](#).