

Exercise Sheet 04

CORBUCCI WCA
03726968

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.

Minimize the error function wrt. \mathbf{w} :

$$\begin{aligned} \mathbf{w}^* &= \operatorname{argmin}_{\mathbf{w}} E_{\text{weighted}}(\mathbf{w}) \\ &= \operatorname{argmin}_{\mathbf{w}} \frac{1}{2} \sum_{i=1}^N t_i [\mathbf{w}^T \phi(\mathbf{x}_i) - y_i]^2 \end{aligned}$$

I define the vector $\mathbf{T} \in \mathbb{R}^N$ which contains $t_i \in \mathbb{R} \quad \forall i \quad 1 < i < N$.

$$= \operatorname{argmin}_{\mathbf{w}} \frac{1}{2} (\mathbf{w} - \mathbf{y})^T \mathbf{T} (\mathbf{w} - \mathbf{y})$$

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

I compute the gradient:

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

$$= \nabla_{\omega} \frac{1}{2} (\omega^T \Phi^T - y^T) T (\Phi \omega - y)$$

$$= \frac{1}{2} (\omega^T \Phi^T T - y^T T) (\Phi \omega - y)$$

$$\nabla_{\omega} \frac{1}{2} \left(\omega^T \Phi^T T \Phi \omega - \omega^T \Phi^T T y - y^T \Phi T \omega + y^T T y \right)$$

Differentiate with respect to ω :

$$\frac{1}{2} \omega^T (\Phi^T T \Phi + \Phi^T T \Phi) - \frac{1}{2} (2 \Phi^T T y) =$$

$$= \omega^T (\Phi^T T \Phi) - (\Phi^T T y) = 0$$

Set the gradient to 0 to find ω^*

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

How can t_i can be interpreted in term of:

(a) the variance of the noise on the data

Using the weighting factor t_i and knowing the variance σ^2 of the noise data we can balance the noise term setting $t_i = \frac{1}{\sigma^2}$. This leads to a more precise regression.

(b) Data points for which there are exact copies in the dataset

I think that if we have a data point which has some copies, the model will consider these data points as more important with respect to the others data points.

Maybe we can use the weight to balance the importance of these points so that the model will not consider them more important than the others.

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.

$$E_{LS} = \frac{1}{2} \sum_{i=1}^N \left[\omega^\top \phi(x_i) - y_i \right]^2 + \frac{1}{2} \sum_{i=N+1}^{N+M} \left[\omega^\top \phi(x_i) - y_i \right]^2$$

if $N+1 \leq i \leq N+M$

$$\phi(x_i) = \begin{pmatrix} \sqrt{\lambda} \mathbf{I}_{M \times M} \\ \vdots \\ \sqrt{\lambda} \mathbf{I}_{M \times M} \end{pmatrix} = V$$

Consider the second part of the E_{LS} :

$$\begin{aligned} & \frac{1}{2} \sum_{i=N+1}^{N+M} (\omega^\top V - y_i)^\top (\omega^\top V - y_i) \\ & \quad \text{equal to } \phi \\ &= \frac{1}{2} \left(\omega_0 \sqrt{\lambda} \mathbf{I}_{M \times M} + \dots + \omega_M \sqrt{\lambda} \mathbf{I}_{M \times M} \right)^2 \\ &= \frac{1}{2} \left(\sqrt{\lambda} \mathbf{I}_{M \times M} (\omega_0 + \dots + \omega_M) \right)^2 \\ &= \frac{1}{2} \lambda (\omega_0 + \dots + \omega_M)^2 \\ &= \frac{\lambda}{2} \|\omega\|_2^2 \end{aligned}$$

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?

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

Compute the gradient

$$\nabla_{\mathbf{w}} \frac{1}{2} (\mathbf{w} - \mathbf{y})^T (\mathbf{w} - \mathbf{y}) + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

$$\nabla_{\mathbf{w}} \frac{1}{2} (\mathbf{w}^T \mathbf{w} - 2\mathbf{w}^T \mathbf{y} + \mathbf{y}^T \mathbf{y}) + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

$$\nabla_{\mathbf{w}} \frac{1}{2} (2\mathbf{w}^T \mathbf{w} - 2\mathbf{y}^T \mathbf{y}) + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

$$\mathbf{w}^T \mathbf{w} - \mathbf{y}^T \mathbf{y} + \frac{\lambda}{2} 2\mathbf{w}^T = 0$$

Put the gradient equal to 0 to find \mathbf{w}^*

$$\Phi^T \Phi w + \lambda w^T = \Phi^T y$$

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

→ Identity matrix

Consider the case $N < M$

I think that we could have overfitting because our model will fit every data point that are in our dataset.

To fit every data points the coefficient w takes large positive values and large negative values.

the regularization (if we choose the right value for λ) solves this issue because it avoid that the coefficients can reach large positive or negative values

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.

Multi output: $f_{\omega}(x) : \mathbb{R}^n \rightarrow \mathbb{R}^M$

In this case we have

$$Y = \omega X + \epsilon$$

↓ ↓ ↓ ↗
 matrix

Likelihood of the entire dataset:

$$P(y|X, \omega, \beta) = \prod_{i=1}^N P(y_i|f_{\omega}(x_i), \beta)$$

$$(\omega_{ML}, \beta_{ML}) = \underset{\omega, \beta}{\operatorname{argmax}} P(y|X, \omega, \beta)$$

$$= \underset{\omega, \beta}{\operatorname{argmin}} -\ln P(y|X, \omega, \beta)$$

$$F_{ML} = -\ln P(y|X, \omega, \beta)$$

$$= -\ln \left[\prod_{i=1}^N N(y_i | f_{\omega}(x_i), \beta^{-1}) \right]$$

$$= -\ln \left[\prod_{i=1}^N \sqrt{\frac{\beta}{2\pi}} \exp \left(-\frac{\beta}{2} (\omega x_i - y_i)^2 \right) \right]$$

$$= -\sum_{i=1}^N \ln \left(\sqrt{\frac{\beta}{2\pi}} \right) - \frac{\beta}{2} (\omega x_i - y_i)^2$$

$$= \frac{\beta}{2} \sum (\omega x_i - y_i)^2 - \frac{N}{2} \ln \beta + \frac{N}{2} \ln 2\pi$$

Derive w.r.t β and ω

$$\omega_{ML} = \frac{\beta}{2} \sum (\omega x_i - y_i)^2$$

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

$$\beta_{ML} = \operatorname{argmin} \left[\frac{\beta}{2} \sum (\omega_{ML} x_i - y_i)^2 - \frac{N}{2} \ln \beta + \frac{N}{2} \ln 2\pi \right]$$

set the derivative to 0

$$\frac{1}{2} \sum (\omega_{ML} x_i - y_i)^2 - \frac{N}{2\beta} = 0$$

$$\frac{1}{\beta_{ML}} = \frac{1}{2N} \sum (\omega_{ML} x_i - y_i)^2$$

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} .
- 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 want a \mathbf{w}_{new} that produces the same predictions on \mathbf{X}_{new} as \mathbf{w}^* :

$$f_{\mathbf{w}^*}(\mathbf{x}) = f_{\mathbf{w}_{new}}(\mathbf{x})$$

$$\mathbf{w}_0^* + \mathbf{w}_1^* x_1 + \dots + \mathbf{w}_D^* x_D = \mathbf{w}_{0,new} + \mathbf{w}_{1,new} a x_1 + \dots + \mathbf{w}_{D,new} a x_{D,new}$$

$$\mathbf{w}_0 + \mathbf{w}^T \mathbf{x} = \mathbf{w}_{new}^T \cdot \mathbf{a} \cdot \mathbf{x}$$

$$\mathbf{w}_{new} = \frac{1}{a} \cdot \mathbf{w}^*$$

(b)

$$\hat{\omega}_{\text{new}} = \underset{\omega}{\operatorname{argmin}} \frac{1}{2} \sum_{i=1}^n (\omega^T x_i^{\text{new}} - y_i)^2 + \frac{\lambda_{\text{new}}}{2} \omega^T \omega$$

From the Problem 3 we know that the closed form solution for ridge regression is:

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

$$\text{We know that } \hat{\omega}_{\text{new}} = \frac{1}{q} \hat{\omega}^*$$

$$\frac{1}{q} \hat{\omega}^* = (X^T X + \lambda_{\text{new}} I)^{-1} X^T y$$

$$\frac{1}{q} (X^T X + \lambda_{\text{new}} I) \hat{\omega}^* = X^T y$$

$$\frac{1}{q} X^T \hat{\omega}^* X + \frac{1}{q} \lambda_{\text{new}} \hat{\omega}^* = X^T y$$

$$\frac{1}{q} \lambda_{\text{new}} \hat{\omega}^* = X^T y - \frac{1}{q} X^T \hat{\omega}^* X$$

$$\lambda_{\text{new}} \hat{\omega}^* = q X^T y - X^T \hat{\omega}^* X$$

exercise_04_notebook

November 10, 2019

1 Programming assignment 2: Linear regression

```
In [10]: import numpy as np  
  
from sklearn.datasets import load_boston  
from sklearn.model_selection import train_test_split
```

1.1 Your task

In this notebook code skeleton for performing linear regression is given. Your task is to complete the functions where required. You are only allowed to use built-in Python functions, as well as any numpy functions. No other libraries / imports are allowed.

1.2 Exporting the results to PDF

Once you complete the assignments, export the entire notebook as PDF and attach it to your homework solutions. The best way of doing that is 1. Run all the cells of the notebook. 2. Export/download the notebook as PDF (File -> Download as -> PDF via LaTeX (.pdf)). 3. Concatenate your solutions for other tasks with the output of Step 2. On a Linux machine you can simply use `pdfunite`, there are similar tools for other platforms too. You can only upload a single PDF file to Moodle.

Make sure you are using `nbconvert` Version 5.5 or later by running `jupyter nbconvert --version`. Older versions clip lines that exceed page width, which makes your code harder to grade.

1.3 Load and preprocess the data

In this assignment we will work with the Boston Housing Dataset. The data consists of 506 samples. Each sample represents a district in the city of Boston and has 13 features, such as crime rate or taxation level. The regression target is the median house price in the given district (in \$1000's).

More details can be found here: <http://lib.stat.cmu.edu/datasets/boston>

```
In [11]: X , y = load_boston(return_X_y=True)  
  
# Add a vector of ones to the data matrix to absorb the bias term  
# (Recall slide #7 from the lecture)  
X = np.hstack([np.ones([X.shape[0], 1]), X])  
# From now on, D refers to the number of features in the AUGMENTED dataset
```

```

# (i.e. including the dummy '1' feature for the absorbed bias term)

# Split into train and test
test_size = 0.2
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=test_size)

```

1.4 Task 1: Fit standard linear regression

```

In [12]: def fit_least_squares(X, y):
    """Fit ordinary least squares model to the data.
    D is the number of feature in the augmented dataset
    Parameters
    -----
    X : array, shape [N, D]
        (Augmented) feature matrix.
    y : array, shape [N]
        Regression targets.

    Returns
    -----
    w : array, shape [D]
        Optimal regression coefficients (w[0] is the bias term).

    """
    # TODO
    XT = X.T
    w = (np.linalg.inv(XT.dot(X))).dot(XT.dot(y))
    return w

```

1.5 Task 2: Fit ridge regression

```

In [13]: def fit_ridge(X, y, reg_strength):
    """Fit ridge regression model to the data.

    Parameters
    -----
    X : array, shape [N, D]
        (Augmented) feature matrix.
    y : array, shape [N]
        Regression targets.
    reg_strength : float
        L2 regularization strength (denoted by lambda in the lecture)

    Returns
    -----
    w : array, shape [D]
        Optimal regression coefficients (w[0] is the bias term).

```

```

"""
# TODO
XT = X.T
w = (np.linalg.inv(np.dot(XT, X) + (reg_strength)*np.identity(len(XT))))
     .dot(XT.dot(y))
return w

```

1.6 Task 3: Generate predictions for new data

```
In [14]: def predict_linear_model(X, w):
    """Generate predictions for the given samples.

    Parameters
    -----
    X : array, shape [N, D]
        (Augmented) feature matrix.
    w : array, shape [D]
        Regression coefficients.

    Returns
    -----
    y_pred : array, shape [N]
        Predicted regression targets for the input data.

"""
# TODO

y_pred = [w.dot(x) for x in X]
return y_pred
```

1.7 Task 4: Mean squared error

```
In [15]: def mean_squared_error(y_true, y_pred):
    """Compute mean squared error between true and predicted regression targets.

    Reference: `https://en.wikipedia.org/wiki/Mean_squared_error`

    Parameters
    -----
    y_true : array
        True regression targets.
    y_pred : array
        Predicted regression targets.

    Returns
    -----
    mse : float
        Mean squared error.

```

```

"""
# TODO
mse = np.square(np.subtract(y_true, y_pred)).mean()
return mse

```

1.8 Compare the two models

The reference implementation produces * MSE for Least squares ≈ 23.98 * MSE for Ridge regression ≈ 21.05

Your results might be slightly (i.e. $\pm 1\%$) different from the reference solution due to numerical reasons.

```
In [16]: # Load the data
np.random.seed(1234)
X, y = load_boston(return_X_y=True)
X = np.hstack([np.ones([X.shape[0], 1]), X])
test_size = 0.2
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=test_size)

# Ordinary least squares regression
w_ls = fit_least_squares(X_train, y_train)
y_pred_ls = predict_linear_model(X_test, w_ls)
mse_ls = mean_squared_error(y_test, y_pred_ls)
print('MSE for Least squares = {}'.format(mse_ls))

# Ridge regression
reg_strength = 1
w_ridge = fit_ridge(X_train, y_train, reg_strength)
y_pred_ridge = predict_linear_model(X_test, w_ridge)
mse_ridge = mean_squared_error(y_test, y_pred_ridge)
print('MSE for Ridge regression = {}'.format(mse_ridge))

MSE for Least squares = 23.964571384952194
MSE for Ridge regression = 21.034931215919112
```