

ASSIGNMENT 4

THEODOR STOICAN

03725732

Problem 1

$$E(w) = \frac{1}{2} \sum_{i=1}^N t_i \cdot (w^T \phi(x_i) - y_i)^2$$

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

$$E(w) = \frac{1}{2} \cdot (\phi(X) w - y)^T T (\phi(X) w - y),$$

where $T = \begin{bmatrix} t_1 & 0 & \dots & 0 \\ \vdots & t_2 & & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ 0 & \ddots & & t_N \end{bmatrix}$

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

$$= \nabla_w \frac{1}{2} \left(w^T \phi(X)^T T \phi(X) w - \underbrace{w^T \phi(X)^T T y}_{\text{NxN}} \right) -$$

$$- y^T T \phi(X) w + y^T T y$$

$$\begin{aligned}
 &= \nabla_w \frac{1}{2} \left(w^T \phi(x)^T T \phi(x) w - 2 w^T \phi(x)^T T y + y^T T y \right) \\
 &= \phi(x)^T T \phi(x) w - \phi(x)^T T y
 \end{aligned}$$

To find the minimum: $\nabla_w E(w) = 0$

$$\Rightarrow \phi(x)^T T \phi(x) w - \phi(x)^T T y = 0$$

$$\Rightarrow \phi(x)^T T \phi(x) w \equiv \phi(x)^T T y$$

$$\Rightarrow w^* = (\phi(x)^T T \phi(x))^{-1} \phi(x)^T T y$$

1) By using t_i to scale each point, we get to the same effect as we were modifying the variance of the noise. For example, by choosing particularly large t_i factors (especially for large differences), it's like increasing the variance of the noise for those particular points. Similarly, by choosing small t_i factors, it's like decreasing the variance of the noise for such points.

2) By using t_i to scale the error for a certain data point, we get the same effect as having t_i identical data points (if $t_i \in \mathbb{N}$, of course). So, by doing this for an arbitrary amount of points in the dataset, we essentially augment the dataset with duplicate points.

Problem 2

After augmenting ϕ :

$$\phi = \begin{bmatrix} \phi_0(x_n) & \phi_1(x_n) & \cdots & \phi_m(x_n) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(x_N) & \phi_1(x_N) & \cdots & \phi_m(x_N) \\ \sqrt{\lambda} & 0 & \cdots & 0 \\ 0 & 0 & \cdots & \sqrt{\lambda} \end{bmatrix}$$

The loss function will now become:

$$E(w) = \frac{1}{2} \sum_{i=1}^{N+M} (w^T \phi(x_i) - y_i)^2,$$

taking into account that the last M rows in y are zeros.

$$\Rightarrow E(w) = \frac{1}{2} \sum_{i=1}^N (w^\top \phi(x_i) - y_i)^2 + \frac{1}{2} \sum_{i=N+1}^{N+M} (w^\top \phi(x_i) - 0)^2$$

$$\Rightarrow E(w) = \frac{1}{2} \sum_{i=1}^N (w^\top \phi(x_i) - y_i)^2 + \frac{1}{2} \sum_{i=N+1}^{N+M} (\sqrt{\lambda} \cdot w_i)^2$$

$$= \text{--- } y \text{ ---} + \frac{1}{2} \cdot \lambda \cdot \sum_{i=1}^M (w_i)^2$$

$$= \text{--- } y \text{ ---} + \frac{\lambda}{2} \|w\|_2^2 \rightarrow \text{ridge regression}$$

Problem 3

$$E(w) = \frac{1}{2} (w^\top \phi(x) - y)^\top (w^\top \phi(x) - y) + \frac{\lambda}{2} w^\top w$$

From Problem 1, we know that :

$$\nabla_w E(w) = \nabla_w \left(\frac{1}{2} (w^\top \phi(x))^\top \phi(x) w - 2 w^\top \phi(x)^\top y + y^\top y + \lambda w^\top w \right)$$

$$\nabla_w E(w) = \phi(x)^\top \phi(x) w - \phi(x)^\top y + \lambda w$$

$$\nabla_w E(w) = 0 \Rightarrow \phi(x)^T \phi(x) w + \lambda w = \phi(x)^T y$$

$$\Rightarrow (\phi(x)^T \phi(x) + \lambda I) w = \phi(x)^T y$$

$$\Rightarrow w^* = [\phi(x)^T \phi(x) + \lambda I]^{-1} \phi(x)^T y$$

Problem 4

For single output linear regression, we had:

$$P(y_i | f_w(x_i), \beta) = N(y_i | f_w(x_i), \beta^{-1}) - \text{the likelihood of each single target } y_i.$$

Now, for multi-output linear regression, the likelihood will be a multi-variate Gaussian distribution (for a specific target y_i):

$$P(y_i | f_w(x_i), S) = N(y_i | f_w(x_i), S^{-1}), \text{ where } S = \sum$$

$$y_i \in \mathbb{R}^m, x_i \in \mathbb{R}^n$$

Therefore, the likelihood of the whole dataset \mathcal{Y} is:

$$P(Y | X, w, S) = \prod_{i=1}^N P(y_i | f_w(x_i), S), \text{ where } Y \in \mathbb{R}^{m \times n}$$

$$\Rightarrow w_{ML}, S_{ML} = \underset{w, S}{\operatorname{argmax}} P(Y | X, w, S)$$

$$= \underset{w, S}{\operatorname{argmin}} -\ln P(Y | X, w, S)$$

$$E_{ML}(w, S) = -\ln P(Y/X, w, S) \text{ - ML error function}$$

$$= -\ln \left[\prod_{i=1}^N P(Y_i | f_w(x_i), S) \right]$$

$$= -\ln \left[\prod_{i=1}^N \frac{\exp(-\frac{1}{2} (Y_i - w^T \phi(x_i))^T S (Y_i - w^T \phi(x_i)))}{\sqrt{(2\pi)^m |S|}} \right]$$

$$= -\sum_{i=1}^N \ln \left[\frac{1}{\sqrt{(2\pi)^m |S|}} \cdot \exp(-\frac{1}{2} (Y_i - w^T \phi(x_i))^T S (Y_i - w^T \phi(x_i))) \right]$$

$$= \frac{S}{2} \sum_{i=1}^N (Y_i - w^T \phi(x_i))^T (Y_i - w^T \phi(x_i)) - N \ln \pi + \frac{N}{2} \ln (2\pi)^m / |S|$$

$$= \frac{S}{2} \sum_{i=1}^N (Y_i - w^T \phi(x_i))^T (Y_i - w^T \phi(x_i)) + \frac{N}{2} \ln (\pi)^m / |S|$$

$$w_{ML} = \underset{w}{\operatorname{argmin}} E_{ML}(w, S)$$

$$= \underset{w}{\operatorname{argmin}} \left\{ \frac{S}{2} \sum_{i=1}^N (Y_i - w^T \phi(x_i))^T (Y_i - w^T \phi(x_i)) + \underbrace{\frac{N}{2} \ln (\pi)^m / |S|}_{\text{Const}} \right\}$$

$$= \underset{w}{\operatorname{argmin}} \left[\frac{1}{2} \sum_{i=1}^N \left(Y_i - w^T \phi(x_i) \right)^T \left(Y_i - w^T \phi(x_i) \right) \right] -$$

We dropped S since by minimizing this function we'll minimize the previous one as well

Here, w_{ML} for multi-output linear regression is identical with w_{ML} for single-output l. regression:

$$w_{ML} = (\phi^T \phi)^{-1} \phi^T Y,$$

except for the fact
that now $w_{ML} \in \mathbb{R}^{(p+1) \times m}$,
 $Y \in \mathbb{R}^{N \times m}$, $\phi \in \mathbb{R}^{N \times (p+1)}$

Now, we determine S_{ML} :

$$S_{ML} = \underset{\Sigma}{\operatorname{argmin}} E_{ML}(w_{ML}) S$$

$$= \underset{S}{\operatorname{argmin}} \left[\frac{S}{2} \sum_{i=1}^N \left(Y_i - w^T \phi(x_i) \right)^T \left(Y_i - w^T \phi(x_i) \right) + \frac{N}{2} \ln(\hat{\sigma}^2) \right]$$

At the same time, we know that:

$$|\Sigma^{-1}| = \frac{1}{|\Sigma|} \Rightarrow |\Sigma| = \frac{1}{|\Sigma^{-1}|}$$

And :

$$\frac{d|\Sigma^{-1}|}{d(\Sigma^{-1})} = |\Sigma^{-1}| \text{tr}(\Sigma)$$

$$\begin{aligned}\frac{dE_{ML}}{dS} &= \frac{1}{2} \sum_{i=1}^N (Y_i - w^T \phi(x_i))^T (Y_i - w^T \phi(x_i)) - \left(\frac{N}{2} \ln |S| \right) \\ &= \frac{1}{2} \sum_{i=1}^N (Y_i - w^T \phi(x_i))^T (Y_i - w^T \phi(x_i)) - \frac{N}{2|\Sigma|} |\Sigma|^T |\Sigma| \text{tr}(\Sigma^{-1}) \\ &= 0\end{aligned}$$

$$\Rightarrow \text{tr}(\Sigma) = \frac{1}{N} \sum_{i=1}^N (Y_i - w^T \phi(x_i))^T (Y_i - w^T \phi(x_i))$$

If we assume that the multivariate Gaussian is isotropic, then Σ will have non-zero elements only on the main diagonal (each one equal to $\frac{\text{tr}(\Sigma)}{N}$).

exercise_04_notebook (1)

November 10, 2019

1 Programming assignment 2: Linear regression

```
[14]: 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>

```
[15]: 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

```
[16]: def fit_least_squares(X, y):
    """Fit ordinary least squares model to the data.

    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
    return np.linalg.pinv(X) @ y
```

1.5 Task 2: Fit ridge regression

```
[17]: 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
return np.linalg.inv(np.transpose(X) @ X + reg_strength * np.identity(X.
->shape[1])) @ np.transpose(X) @ y

```

1.6 Task 3: Generate predictions for new data

```
[18]: 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
    return X @ w

```

1.7 Task 4: Mean squared error

```
[19]: 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
return np.mean((y_true - y_pred) ** 2)
```

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.

```
[20]: # 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 = 23.964571384956837

MSE for Ridge regression = 21.03493121591825

[]: