

# Exercises week 35

## Deriving and Implementing Ordinary Least Squares

This week you will be deriving the analytical expressions for linear regression, building up the model from scratch. This will include taking several derivatives of products of vectors and matrices. Such derivatives are central to the optimization of many machine learning models. Although we will often use automatic differentiation in actual calculations, to be able to have analytical expressions is extremely helpful in case we have simpler derivatives as well as when we analyze various properties (like second derivatives) of the chosen cost functions.

Vectors are always written as boldfaced lower case letters and matrices as upper case boldfaced letters. You will find useful the notes from week 35 on derivatives of vectors and matrices. See also the textbook of Faisal et al, chapter 5 and in particular sections 5.3-5.5 at

<https://github.com/CompPhysics/MachineLearning/blob/master/doc/Textbooks/MathMLb>

### Learning goals

After completing these exercises, you will know how to

- Take the derivatives of simple products between vectors and matrices
- Implement OLS using the analytical expressions
- Create a feature matrix from a set of data
- Create a feature matrix for a polynomial model
- Evaluate the MSE score of various model on training and test data, and comparing their performance

### Deliverables

Complete the following exercises while working in a jupyter notebook. Then, in canvas, include

- The jupyter notebook with the exercises completed
- An exported PDF of the notebook  
([https://code.visualstudio.com/docs/datascience/jupyter-notebooks#\\_export-your-jupyter-notebook](https://code.visualstudio.com/docs/datascience/jupyter-notebooks#_export-your-jupyter-notebook))

## How to take derivatives of Matrix-Vector expressions

In these exercises it is always useful to write out with summation indices the various quantities. Take also a look at the weekly slides from week 35 and the various

examples included there.

As an example, consider the function

$$f(\mathbf{x}) = \mathbf{A}\mathbf{x},$$

which reads for a specific component  $f_i$  (we define the matrix  $\mathbf{A}$  to have dimension  $n \times n$  and the vector  $\mathbf{x}$  to have length  $n$ )

$$f_i = \sum_{j=0}^{n-1} a_{ij}x_j,$$

which leads to

$$\frac{\partial f_i}{\partial x_j} = a_{ij},$$

and written out in terms of the vector  $\mathbf{x}$  we have

$$\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} = \mathbf{A}.$$

## Exercise 1 - Finding the derivative of Matrix-Vector expressions

a) Consider the expression

$$\frac{\partial(\mathbf{a}^T \mathbf{x})}{\partial \mathbf{x}},$$

Where  $\mathbf{a}$  and  $\mathbf{x}$  are column-vectors with length  $n$ .

What is the *shape* of the expression we are taking the derivative of?

What is the *shape* of the thing we are taking the derivative with respect to?

What is the *shape* of the result of the expression?

### Answer a)

We are taking the derivative of  $\mathbf{a}^T \mathbf{x}$ , where  $\mathbf{a}^T$  is a row-vector with length  $n$ , and  $\mathbf{x}$  is a column vector with length  $n$ . Multiplied together they compute to a single number (shape  $(1 \times 1)$ ).

We are taking the derivative with respect to  $\mathbf{x}$ . It is a column-vector with length  $n$ , so it has the shape  $(n \times 1)$ .

The expression can be written as

$$\left[ \frac{\partial(\mathbf{a}^T \mathbf{x})}{\partial x_0} \quad \frac{\partial(\mathbf{a}^T \mathbf{x})}{\partial x_1} \quad \cdots \quad \frac{\partial(\mathbf{a}^T \mathbf{x})}{\partial x_{n-1}} \right],$$

and because the  $\mathbf{a}^T \mathbf{x}$  is a single number, the shape of the result of the expression is this row-vector with shape  $(1 \times n)$

---

**b) Show that**

$$\frac{\partial(\mathbf{a}^T \mathbf{x})}{\partial \mathbf{x}} = \mathbf{a}^T,$$

**Answer b)**

We have the expression

$$\mathbf{a}^T \mathbf{x} = \sum_{i=0}^{n-1} a_i x_i$$

Then we have

$$\begin{aligned} \frac{\partial(\mathbf{a}^T \mathbf{x})}{\partial \mathbf{x}} &= \left[ \frac{\partial(\mathbf{a}^T \mathbf{x})}{\partial x_0} \quad \frac{\partial(\mathbf{a}^T \mathbf{x})}{\partial x_1} \quad \dots \quad \frac{\partial(\mathbf{a}^T \mathbf{x})}{\partial x_{n-1}} \right] \\ &= \left[ \frac{\partial}{\partial x_0} \left( \sum_{i=0}^{n-1} a_i x_i \right) \quad \frac{\partial}{\partial x_1} \left( \sum_{i=0}^{n-1} a_i x_i \right) \quad \dots \quad \frac{\partial}{\partial x_{n-1}} \left( \sum_{i=0}^{n-1} a_i x_i \right) \right] \\ &= [a_0 \quad a_1 \quad \dots \quad a_{n-1}] \\ &= \underline{\underline{\mathbf{a}^T}} \end{aligned}$$


---

**c) Show that**

$$\frac{\partial(\mathbf{a}^T \mathbf{A} \mathbf{a})}{\partial \mathbf{a}} = \mathbf{a}^T (\mathbf{A} + \mathbf{A}^T),$$

**Answer c)**

First we find the solution for an element  $k$

$$\begin{aligned}
\frac{\partial(\mathbf{a}^T \mathbf{A} \mathbf{a})}{\partial a_k} &= \frac{\partial}{\partial a_k} \left( \sum_{i=0}^{n-1} a_i \sum_{j=0}^{n-1} A_{ij} a_j \right) \\
&= \sum_{i=0}^{n-1} \left( \frac{\partial a_i}{\partial a_k} \sum_{j=0}^{n-1} A_{ij} a_j + a_i \sum_{j=0}^{n-1} \left( \frac{\partial a_j}{\partial a_k} A_{ij} \right) \right) \\
&= \sum_{i=0}^{n-1} \left( \frac{\partial a_i}{\partial a_k} \sum_{j=0}^{n-1} A_{ij} a_j \right) + \sum_{i=0}^{n-1} \left( a_i \sum_{j=0}^{n-1} \left( \frac{\partial a_j}{\partial a_k} A_{ij} \right) \right) \\
&= \sum_{j=0}^{n-1} A_{kj} a_j + \sum_{i=0}^{n-1} A_{ik} a_i \\
&= \sum_{i=0}^{n-1} a_i (A_{ki} + A_{ik}) \\
&= \sum_{i=0}^{n-1} a_i (A_{ik} + A_{ik}^T) \\
&= \sum_{i=0}^{n-1} a_i (A + A^T)_{ik}
\end{aligned}$$

Then, setting in for all  $k$  values in the resulting matrix we get

$$\begin{aligned}
\frac{\partial(\mathbf{a}^T \mathbf{A} \mathbf{a})}{\partial \mathbf{a}} &= \left[ \sum_{i=0}^{n-1} a_i (A + A^T)_{i0} \quad \sum_{i=0}^{n-1} a_i (A + A^T)_{i1} \quad \dots \quad \sum_{i=0}^{n-1} a_i (A + A^T)_{in-1} \right] \\
&= \underline{\underline{\mathbf{a}^T (\mathbf{A} + \mathbf{A}^T)}}
\end{aligned}$$


---

## Exercise 2 - Deriving the expression for OLS

The ordinary least squares method finds the parameters  $\boldsymbol{\theta}$  which minimizes the squared error between our model  $\mathbf{X}\boldsymbol{\theta}$  and the true values  $\mathbf{y}$ .

To find the parameters  $\boldsymbol{\theta}$  which minimizes this error, we take the derivative of the squared error expression with respect to  $\boldsymbol{\theta}$ , and set it equal to 0.

**a)** Very briefly explain why the approach above finds the parameters  $\boldsymbol{\theta}$  which minimizes this error.

**Answer a)** The square error is a convex function, which means that wherever the derivative is equal to 0 is the minimum of the function.

---

We typically write the squared error as

$$\|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|^2$$

which we can rewrite in matrix-vector form as

$$(\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})$$

**b)** If  $\mathbf{X}$  is invertible, what is the expression for the optimal parameters  $\boldsymbol{\theta}$ ? (**Hint:** Don't compute any derivatives, but solve  $\mathbf{X}\boldsymbol{\theta} = \mathbf{y}$  for  $\boldsymbol{\theta}$ )

**Answer b)** If  $\mathbf{X}$  is invertible, then  $\boldsymbol{\theta} = \mathbf{X}^{-1}\mathbf{y}$ .

---

**c)** Show that

$$\frac{\partial(\mathbf{x} - \mathbf{A}\mathbf{s})^T (\mathbf{x} - \mathbf{A}\mathbf{s})}{\partial \mathbf{s}} = -2(\mathbf{x} - \mathbf{A}\mathbf{s})^T \mathbf{A},$$

**Answer c)** First we define

$$\begin{aligned} \mathbf{v} &= \mathbf{x} - \mathbf{A}\mathbf{s} \\ w &= \mathbf{v}^T \mathbf{v} \end{aligned}$$

We use the partial derivatives

$$\frac{\partial w}{\partial \mathbf{v}} = 2\mathbf{v}^T,$$

$$\frac{\partial \mathbf{v}}{\partial \mathbf{s}} = -\mathbf{A},$$

and the chain rule to get

$$\begin{aligned} \frac{\partial(\mathbf{x} - \mathbf{A}\mathbf{s})^T (\mathbf{x} - \mathbf{A}\mathbf{s})}{\partial \mathbf{s}} &= \frac{\partial w}{\partial \mathbf{s}} \\ &= \frac{\partial w}{\partial \mathbf{v}} \frac{\partial \mathbf{v}}{\partial \mathbf{s}} \\ &= 2\mathbf{v}^T (-\mathbf{A}) \\ &= \underline{\underline{-2(\mathbf{x} - \mathbf{A}\mathbf{s})^T \mathbf{A}}} \end{aligned}$$


---

**d)** Using the expression from **c)**, but substituting back in  $\boldsymbol{\theta}$ ,  $\mathbf{y}$  and  $\mathbf{X}$ , find the expression for the optimal parameters  $\boldsymbol{\theta}$  in the case that  $\mathbf{X}$  is not invertible, but  $\mathbf{X}^T \mathbf{X}$  is, which is most often the case.

$$\hat{\boldsymbol{\theta}}_{OLS} = \dots$$

**Answer d)**

By taking the expression in **c)** and replacing  $\mathbf{x}$  with  $\mathbf{y}$ ,  $\mathbf{A}$  with  $\mathbf{X}$ , and  $\mathbf{s}$  with  $\boldsymbol{\theta}$  we get an expression for the derivative of the squared error with respect to  $\boldsymbol{\theta}$ :

$$\frac{\partial (\|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|^2)}{\partial \boldsymbol{\theta}} = \frac{\partial(\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = -2(\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T \mathbf{X}$$

We get the optimal parameters  $\hat{\boldsymbol{\theta}}_{OLS}$  by looking at where this expression is equal to zero:

$$-2\left(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\theta}}_{OLS}\right)^T \mathbf{X} = 0.$$

By rearranging to

$$\left(\hat{\boldsymbol{\theta}}_{OLS}^T \mathbf{X}^T \mathbf{X}\right)^T = (\mathbf{y}^T \mathbf{X})^T$$

$$\mathbf{X}^T \mathbf{X} \hat{\boldsymbol{\theta}}_{OLS} = \mathbf{X}^T \mathbf{y}$$

then finally inverting  $\mathbf{X}^T \mathbf{X}$  we obtain an expression for the optimal parameters

$$\underline{\underline{\hat{\boldsymbol{\theta}}_{OLS} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}}}$$

## Exercise 3 - Creating feature matrix and implementing OLS using the analytical expression

With the expression for  $\hat{\boldsymbol{\theta}}_{OLS}$ , you now have what you need to implement OLS regression with your input data and target data  $\mathbf{y}$ . But before you can do that, you need to set up your input data as a feature matrix  $\mathbf{X}$ .

In a feature matrix, each row is a datapoint and each column is a feature of that data. If you want to predict someone's spending based on their income and number of children, for instance, you would create a row for each person in your dataset, with the monthly income and the number of children as columns.

We typically also include an intercept in our models. The intercept is a value that is added to our prediction regardless of the value of the other features. The intercept tries to account for constant effects in our data that are not dependant on anything else. In our current example, the intercept could account for living expenses which are typical regardless of income or childcare expenses.

We calculate the optimal intercept by including a feature with the constant value of 1 in our model, which is then multiplied by some parameter  $\theta_0$  from the OLS method into the optimal intercept value (which will be  $\theta_0$ ). In practice, we include the intercept in our model by adding a column of ones to the start of our feature matrix.

In [236... `import numpy as np`

In [237... `n = 20`  
`income = np.array([116., 161., 167., 118., 172., 163., 179., 173., 162.,`  
`children = np.array([5, 3, 0, 4, 5, 3, 0, 4, 4, 3, 3, 5, 1, 0, 2, 3, 2, 1`  
`spending = np.array([152., 141., 102., 136., 161., 129., 99., 159., 160.,`

**a)** Create a feature matrix  $\mathbf{X}$  for the features income and children, including an intercept column of ones at the start.

In [238... `X = np.zeros((n, 3))`  
`X[:, 0] = 1`

```
X[:, 1] = income
X[:, 2] = children
```

**b)** Use the expression from **3d)** to find the optimal parameters  $\hat{\beta}_{OLS}$  for predicting spending based on these features. Create a function for this operation, as you are going to need to use it a lot.

```
In [239... def OLS_parameters(X, y):
    X_transpose = np.transpose(X)
    return np.linalg.inv(X_transpose @ X) @ X_transpose @ y
```

## Exercise 4 - Fitting a polynomial

In this course, we typically do linear regression using polynomials, though in real world applications it is also very common to make linear models based on measured features like you did in the previous exercise.

When fitting a polynomial with linear regression, we make each polynomial degree ( $x, x^2, x^3, \dots, x^p$ ) its own feature.

```
In [240... n = 100
x = np.linspace(-3, 3, n)
y = np.exp(-x**2) + 1.5 * np.exp(-(x-2)**2) + np.random.normal(0, 0.1)
```

**a)** Create a feature matrix  $\mathbf{X}$  for the features  $x, x^2, x^3, x^4, x^5$ , including an intercept column of ones at the start. Make this into a function, as you will do this a lot over the next weeks.

```
In [241... def polynomial_features(x, p):
    n = len(x)
    X = np.zeros((n, p + 1))
    X[:, 0] = 1
    for degree in range(1, p+1):
        X[:, degree] = x ** degree
    return X

X = polynomial_features(x, 5)
```

**b)** Use the expression from **3d)** to find the optimal parameters  $\hat{\beta}_{OLS}$  for predicting  $y$  based on these features. If you have done everything right so far, this code will not need changing.

```
In [242... beta = OLS_parameters(X, y)
```

**c)** Like in exercise 4 last week, split your feature matrix and target data into a training split and test split.

```
In [243... from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(X, y)
```

**d)** Train your model on the training data (find the parameters which best fit) and compute the MSE on both the training and test data.

```
In [266... from sklearn.metrics import mean_squared_error

beta = OLS_parameters(X_train, y_train)

prediction_train = X_train @ beta
mse_train = mean_squared_error(y_train, prediction_train)

prediction_test = X_test @ beta
mse_test = mean_squared_error(y_test, prediction_test)

print(f"MSE of training data: {mse_train}")
print(f"MSE of testing data: {mse_test}")
```

MSE of training data: 6.463755225724075e-05

MSE of testing data: 0.0002625423637803146

**e)** Do the same for each polynomial degree from 2 to 10, and plot the MSE on both the training and test data as a function of polynomial degree. The aim is to reproduce Figure 2.11 of [Hastie et al.](#) Feel free to read the discussions leading to figure 2.11 of [Hastie et al.](#)

```
In [267... import matplotlib.pyplot as plt
from_degree = 2
to_degree = 10

degrees = list(range(from_degree, to_degree+1))
mse_train_values = np.zeros(to_degree-from_degree+1)
mse_test_values = np.zeros(to_degree-from_degree+1)
for degree in degrees:
    X = polynomial_features(x, degree)
    X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=0)
    beta = OLS_parameters(X_train, y_train)

    prediction_train = X_train @ beta
    mse_train = mean_squared_error(y_train, prediction_train)

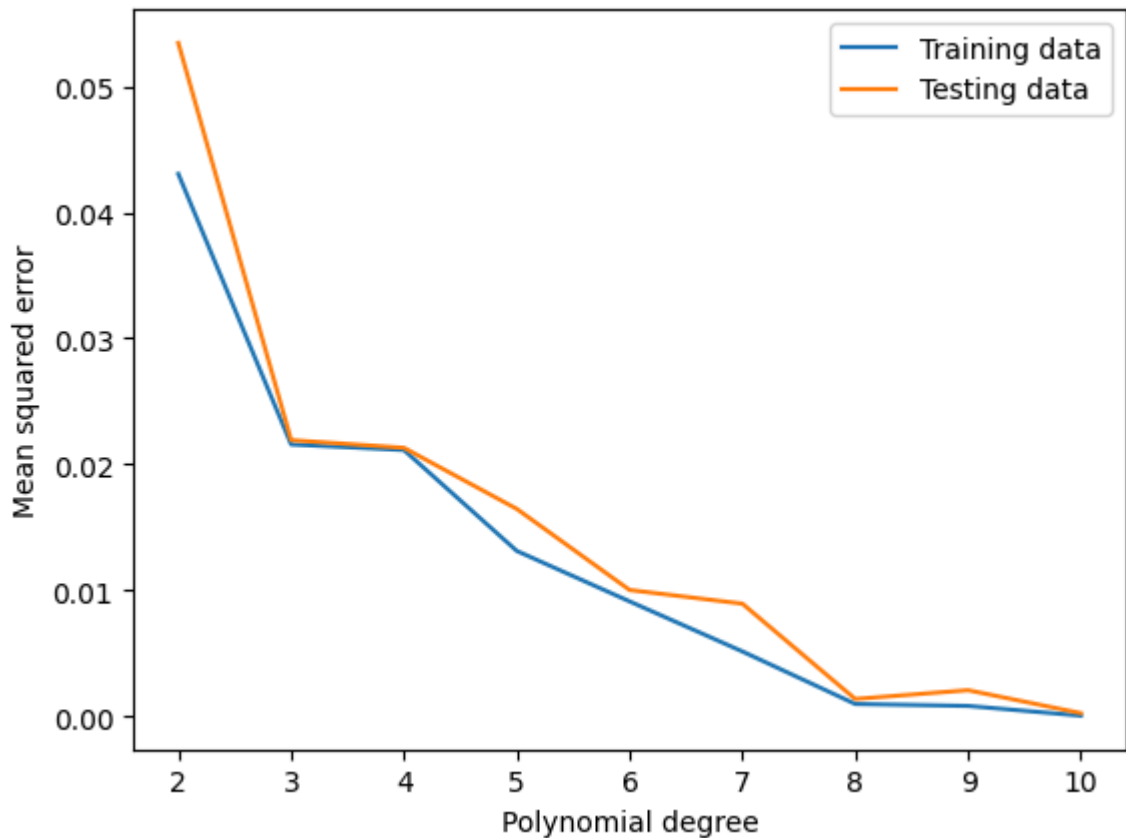
    prediction_test = X_test @ beta
    mse_test = mean_squared_error(y_test, prediction_test)

    mse_train_values[degree-from_degree] = mse_train
    mse_test_values[degree-from_degree] = mse_test

    # plt.title(f"Degree {degree} polynomial")
    # plt.scatter(x, y)
    # plt.plot(x, X @ beta, color="red")
    # plt.show()

plt.plot(degrees, mse_train_values, label="Training data")
plt.plot(degrees, mse_test_values, label="Testing data")
plt.xlabel("Polynomial degree")
plt.ylabel("Mean squared error")
plt.legend()
plt.show()
```





**f)** Interpret the graph. Why do the lines move as they do? What does it tell us about model performance and generalizability?

**Answer f)**

The graph shows how both the MSE of the training data and test data decreases as the degree of polynomial we use to approximate the real function increases. We see that the test data MSE is approximately the same as training data, and both get very low as the polynomial degree increases. This tells us that the model works for both data it has seen before (training data) and new data (test data).

The model has therefore good generalizability (at least for values in the same range as the train and test data (-3 to 3)), and performs well with very low MSE when the polynomial degree reaches 10.

---

## Exercise 5 - Comparing your code with sklearn

When implementing different algorithms for the first time, it can be helpful to double check your results with established implementations before you go on to add more complexity.

**a)** Make sure your `polynomial_features` function creates the same feature matrix as sklearn's `PolynomialFeatures`.

(<https://scikit-learn.org/stable/modules/generated/sklearn.preprocessing.PolynomialFeatures.html>)

```
In [302... from sklearn.preprocessing import PolynomialFeatures
from sklearn.metrics import mean_absolute_error
X_sklearn = PolynomialFeatures(5).fit_transform(x.reshape(-1, 1))
X = polynomial_features(x, 5)

print(f"The largest difference between the two matrices is {abs(X-X_sklea
```

The largest difference between the two matrices is 5.68e-14

### Answer a)

The largest difference between the two matrices is an extremely small number. This is likely due to imprecision in floating point operations.

That means that the sklearn implementation creates the same feature matrix as my own `polynomial_features` function.

---

**b)** Make sure your `OLS_parameters` function computes the same parameters as sklearn's LinearRegression with `fit_intercept` set to False, since the intercept is included in the feature matrix. Use `your_model_object.coef_` to extract the computed parameters.

([https://scikit-learn.org/stable/modules/generated/sklearn.linear\\_model.LinearRegression.html](https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LinearRegression.html))

```
In [315... from sklearn.linear_model import LinearRegression

theta = OLS_parameters(X, y)
theta_sklearn = LinearRegression(fit_intercept=False).fit(X, y).coef_
print(f"Own implementation: \t\t{theta}")
print(f"sklearn implementation: \t{theta_sklearn}")
```

```
Own implementation:      [ 0.83946674  0.27464654 -0.02326439  0.05
342623 -0.0034652 -0.0087781 ]
sklearn implementation:  [ 0.83946674  0.27464654 -0.02326439  0.05
342623 -0.0034652 -0.0087781 ]
```

### Answer b)

My function gives the exact same result as the sklearn implementation here.

---