# A2 - MAD

## Julian Wulff

### November 2019

## Exercise 1

### a)

In order to derive the optimal least squares parameter value $\hat{\mathbf{w}}$ of the function:

$$\mathcal{L} = \frac{1}{N} \sum_{n=1}^{N} \alpha_n (f(\mathbf{x}_n; \mathbf{w}) - t_n)^2 = \frac{1}{N} \sum_{n=1}^{N} \alpha_n (\mathbf{w}^T \mathbf{x}_n - t_n)^2$$

We will start by rewriting the function in terms of various vectors and matrices, which will be easier to manipulate. We already now how to express

$$\mathcal{L} = \frac{1}{N} \sum_{n=1}^{N} (\mathbf{w}^T \mathbf{x}_n - t_n)^2 = \frac{1}{N} (\mathbf{X}\mathbf{w} - \mathbf{t})^\top (\mathbf{X}\mathbf{w} - \mathbf{t})$$

In order to include $\alpha_n$ in this expression we start by defining a matrix $\mathbf{A}$ as diagonal matrix containing the weights $\alpha_1, \ldots, \alpha_n$, in order to multiply these onto our expression. In order to let each individual $\alpha_n$ appear one time each in the final scalar we will multiply $\mathbf{A}$ on to the expression as showm below:

$$\mathcal{L} \frac{1}{N} ((\mathbf{X}\mathbf{w} - \mathbf{t})^\top \mathbf{A} (\mathbf{X}\mathbf{w} - \mathbf{t}))$$

Which gives us the desired expression. We will now multiply out the parenthesis in order continue:

$$\mathcal{L} = \frac{1}{N} (\mathbf{X}\mathbf{w} - \mathbf{t})^\top \mathbf{A} (\mathbf{X}\mathbf{w} - \mathbf{t}) \tag{1}$$

$$= \frac{1}{N} ((\mathbf{X}\mathbf{w}^\top) - \mathbf{t}^\top) \mathbf{A} (\mathbf{X}\mathbf{w} - \mathbf{t}) \tag{2}$$

$$= \frac{1}{N} ((\mathbf{X}\mathbf{w}^\top) \mathbf{A} - \mathbf{t}^\top \mathbf{A}) (\mathbf{X}\mathbf{w} - \mathbf{t}) \tag{3}$$

$$= \frac{1}{N} (\mathbf{X}\mathbf{w}^\top) \mathbf{A} \mathbf{X}\mathbf{w} - \frac{1}{N} \mathbf{t}^\top \mathbf{A} \mathbf{X}\mathbf{w} - \frac{1}{N} (\mathbf{X}\mathbf{w})^\top \mathbf{A} \mathbf{t} + \frac{1}{N} \mathbf{t}^\top \mathbf{A} \mathbf{t} \tag{4}$$

$$= \frac{1}{N} \mathbf{w}^\top \mathbf{X}^\top \mathbf{A} \mathbf{X}\mathbf{w} - \frac{1}{N} \mathbf{t}^\top \mathbf{A} \mathbf{X}\mathbf{w} - \frac{1}{N} \mathbf{w}^\top \mathbf{X}^\top \mathbf{A} \mathbf{t} + \frac{1}{N} \mathbf{t}^\top \mathbf{A} \mathbf{t} \tag{5}$$

Since the two terms $\mathbf{t}^\top \mathbf{A}\mathbf{X}\mathbf{w}$ and $\mathbf{w}^\top \mathbf{X}^\top \mathbf{A}\mathbf{t}$ are the transpose of one another, and also scalars, we conclude that they must be the same and therefore we combine them.

$$\mathcal{L} = \frac{1}{N}\mathbf{w}^\top \mathbf{X}^\top \mathbf{A}\mathbf{X}\mathbf{w} - \frac{2}{N}\mathbf{w}^\top \mathbf{X}^\top \mathbf{A}\mathbf{t} + \frac{1}{N}\mathbf{t}^\top \mathbf{A}\mathbf{t} \tag{6}$$

$$= \frac{1}{N}(\mathbf{w}^\top \mathbf{X}^\top \mathbf{A}\mathbf{X}\mathbf{w} - 2\mathbf{w}^\top \mathbf{X}^\top \mathbf{A}\mathbf{t} + \mathbf{t}^\top \mathbf{A}\mathbf{t}) \tag{7}$$

From this expression we find the derivative using the identities in Table 1.4.[1] Loking at the first term of $\mathcal{L}$ we have $\mathbf{w}^\top \mathbf{X}^\top \mathbf{A}\mathbf{X}\mathbf{w}$. We know that $\mathbf{X}^\top \mathbf{X}$ gives a symmetric matrix. Since $A$ is a diagonal matrix we now that $\mathbf{X}^\top \mathbf{A}\mathbf{X}$ must still be symmetric. For this reason we can use the 4th identity from table 1.4 which gives us $2\mathbf{X}^\top \mathbf{A}\mathbf{X}\mathbf{w}$. Looking at the second term of $\mathcal{L}$ we have $2\mathbf{w}^\top \mathbf{X}^\top \mathbf{A}\mathbf{t}$. We observe that $\mathbf{X}^\top \mathbf{A}\mathbf{t}$ equals a vector that does not depend on $\mathbf{w}$. Using the 1st identity from table 1.4 we get $\mathbf{X}^\top \mathbf{A}\mathbf{t}$. Looking at the last term of $\mathcal{L}$ we see that $\mathbf{t}^\top \mathbf{A}\mathbf{t}$ does not depend on $\mathbf{w}$, which is why it disappears when deriving with respect to $\mathbf{w}$. From this we get the following derivative and solve for 0:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \frac{2}{N}\mathbf{X}^\top \mathbf{A}\mathbf{X}\mathbf{w} - \frac{2}{N}\mathbf{X}^\top \mathbf{A}\mathbf{t} = 0 \tag{8}$$

$$\mathbf{X}^\top \mathbf{A}\mathbf{X}\mathbf{w} = \mathbf{X}^\top \mathbf{A}\mathbf{t} \tag{9}$$

In order to isolate $\mathbf{w}$ we will cancel $\mathbf{X}^\top \mathbf{A}\mathbf{X}$ out by premultiplying its inverse of both sides giving us:

$$\mathbf{I}\mathbf{w} = (\mathbf{X}^\top \mathbf{A}\mathbf{X})^{-1}\mathbf{X}^\top \mathbf{A}\mathbf{t} \tag{10}$$

As $\mathbf{I}\mathbf{w} = \mathbf{w}$ we are left with a matrix equation for $\hat{\mathbf{w}}$, the value of $\mathbf{w}$ that minimises the loss:

$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{A}\mathbf{X})^{-1}\mathbf{X}^\top \mathbf{A}\mathbf{t} \tag{11}$$

## b)

Since we use $\alpha_n = t_n^2$ as the additional weights i expect that the big $t_n$ values will have a bigger impact on the $w_n$ values such that the spread will be less for high $t_n$ values and and bigger for small $t_n$ values. This is indeed what we observe when we look at figure 1.

See section 1.2 to see code for the solution.

---

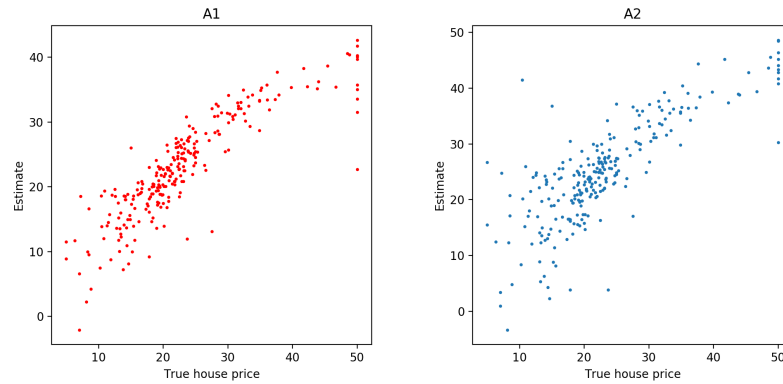[1] S. Rogers  Mark Girolami, A First Course In Machine Learning, 2nd edtion, p.23

Figure 1: Figures containing scatter plot for true house prices vs estimated house prices. One the left is the plot from assignment 1, to the left is the plot using the weighted average loss found using the deriation found in a).

# Exercise 2

## a)

The best value of $lamda = 0.0000003430$ and its $loss = 0.0624312040$
For this model $w0 = 36.3776282284$ and $w1 = -0.0133110046$
For the model with $lamda = 0$, $w0 = 36.4164559025$ and $w1 = -0.0133308857$



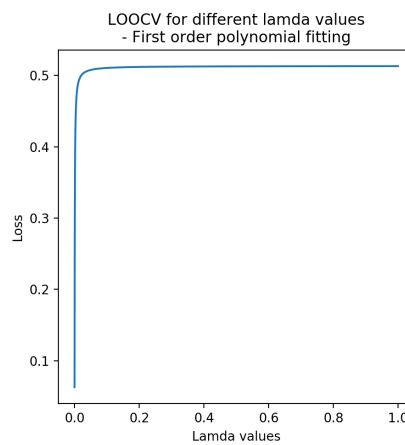Figure 2

**b)**

The best value of $lamda = 0.0000205651$ and its $loss = 0.0521430656$
For this model $w0 = 0.0188300350$, $w1 = 9.1127782080$, $w2 = -0.0138600662$,
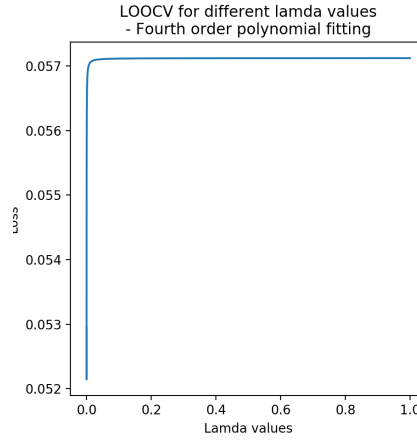$w3 = 0.0000070338$, $w4 = -0.0000000012$



Figure 3

# Exercise 3

**a)**

We are given the CDF:

$$F(x) = \begin{cases} 0 & \leq 0 \\ 1 - e^{-\beta x^{\alpha}} & > 0 \end{cases}$$

Since the CDF $F(X)$ is a continuous functions, we obtain the PDF by taking
its derivative. Therefore we have

$$F'(x) = \begin{cases} 0 & \leq 0 \\ \beta \alpha e^{-\beta x^{\alpha}} x^{\alpha-1} & > 0 \end{cases}$$

**b)**

With $\alpha = 2$ and $\beta = \frac{1}{4}$. To find the probability that the chip works long than
four years we use the CDF to compute

$$1 - F(4) \approx 0.01831563888873422$$

To find the probability that the chip stops working in the time interval $[5; 10]$
years we use the CDF to compute

4

$$F(10) - F(4) \approx 0.0019304541223398308$$

# Exercise 4

## a)

Given a person with no history of conviction ($NC$). The expected mean sentence duration he will have to spend in prison if

- $NC$ talks to the police: $0.002 * (1 - 0.5) * 5 * 0.75 = 0.00375$ years

- $NC$ stays silent: $0.001 * (1 - (0.5 * \frac{1}{4})) * 5 = 0.004375$ years

Given a person with a history of conviction ($C$). The expected mean sentence duration he will have to spend in prison is:

- $C$ talks to the police: $0.005 * (1 - 0.1) * 0.5 * 0.75 = 0.003125$

- $C$ stays silent: $0.001 * (1 - (0.1 * \frac{1}{4})) * 5 = 0.004875$

# 1 Appendix

## 1.1 Exercise 1 files

```python
import numpy as np

class LinearWeightedRegression():
    def __init__(self):

        pass

    def fit(self, X, t):
        # Make sure that we have N-dimensional Numpy arrays (ndarray)
        X = np.array(X).reshape((len(X), -1))
        t = np.array(t).reshape((len(t), 1))

        # Prepend a column of ones til the X matrix
        ones = np.ones((X.shape[0], 1))
        X = np.concatenate((ones, X), axis=1)

        # Transform the w vector into a diagonal vector
        A = np.diag(t[:,0] ** 2)

        self.w = np.linalg.solve(X.T @ A @ X, X.T @ A @ t)

    def predict(self, X):
        X = np.array(X).reshape((len(X), -1))

        ones = np.ones((X.shape[0], 1))
        X = np.concatenate((ones, X), axis=1)
```

```
1026
             predictions = np.dot(X, self.w)
1028
             return predictions
```

lineweighreg.py

```
1000 import numpy as np
     import lineweighreg
1002 import matplotlib.pyplot as plt

1004 # Load data
     train_data = np.loadtxt("boston_train.csv", delimiter=",")
1006 test_data = np.loadtxt("boston_test.csv", delimiter=",")

1008 # Choose all rows and all cols minus the last, choose all rows and
         the last col
     X_train, t_train = train_data[:,:-1], train_data[:,-1]
1010 X_test, t_test = test_data[:,:-1], test_data[:,-1]

1012 # Make sure that we have N-dimensional Numpy arrays (ndarray)
     t_train = t_train.reshape((len(t_train), 1))
1014 t_test = t_test.reshape((len(t_test), 1))

1016 # Fit the model
     model = linweighreg.LinearWeightedRegression()
1018 model.fit(X_train, t_train)

1020 # Make predictions
     predictions = model.predict(X_test)
1022
     # Plot the data
1024 plt.figure(figsize=(5,5))
     plt.scatter(t_test, predictions, s=3)
1026 plt.title("A2")
     plt.xlabel("True house price")
1028 plt.ylabel("Estimate")
     plt.show()
```

exercise1.py

## 1.2   Exercise 2 files

```
1000 import numpy as np

1002 class LinearRegression():
         def __init__(self, lam=0.0):
1004
             self.lam = lam
1006
         def fit(self, X, t):
1008
             X = np.array(X).reshape((len(X), -1))
1010         t = np.array(t).reshape((len(t), 1))
```

6

```
1012            ones = np.ones((X.shape[0], 1))
              X = np.concatenate((ones, X), axis=1)
1014
              diag = self.lam * len(X) * np.identity(X.shape[1])
1016            a = np.dot(X.T, X) + diag
              b = np.dot(X.T, t)
1018            self.w = np.linalg.solve(a,b)

1020        def predict(self, X):
              X = np.array(X).reshape((len(X), -1))
1022
              ones = np.ones((X.shape[0], 1))
1024            X = np.concatenate((ones, X), axis=1)

1026            predictions = np.dot(X, self.w)

1028            return predictions
```

linereg.py

```
1000 import numpy as np
     import linereg
1002 import matplotlib.pyplot as plt

1004 # DATA
     data = np.loadtxt("men-olympics-100.txt")
1006
     X, t = data[:, 0], data[:, 1]
1008 X = X.reshape((len(X), 1))
     t = t.reshape((len(t), 1))
1010
     lamdaValues = np.logspace(-8, 0, 100, base=10)
1012
     def LOOCV(X, t, l):
1014     errors = []
         for lam in lamdaValues:
1016         model = linereg.LinearRegression(lam=lam)
             modelError = 0
1018         for i in range(X.shape[0]):
                 X_train = np.delete(X, i, 0)
1020             t_train = np.delete(t, i, 0)

1022             X_pred = X[i]
                 X_pred = X_pred.reshape((-1, len(X_pred)))
1024             t_pred = t[i]
                 t_pred = t_pred.reshape((len(t_pred), 1))
1026
                 model.fit(X_train, t_train)
1028
                 prediction = model.predict(X_pred)
1030
                 error = (prediction - t_pred)**2
1032             modelError += error[0][0]
             modelError = modelError / X.shape[0]
1034         errors.append(modelError)
         return errors
1036
```

7

```python
     # First order polynomial
1038 firstOrderLOOCV = LOOCV(X, t, lamdaValues)
     bestValueindex = np.argmin(firstOrderLOOCV)
1040 bestLamda = lamdaValues[bestValueindex]

1042 model = linereg.LinearRegression(lam=bestLamda)
     model.fit(X, t)
1044 print("======= First order polynomial fitting =======")
     print("The best value of lamda=%.10f and its loss=%.10f" %
1046     (bestLamda, firstOrderLOOCV[bestValueindex]))
     print("For this model w0=%.10f and w1=%.10f \n" %
1048     (model.w[0][0], model.w[1][0]))

1050 modelLam0 = linereg.LinearRegression()
     modelLam0.fit(X, t)
1052 print("For the model with lamda=0 w0=%.10f and w1=%.10f\n" %
         (modelLam0.w[0][0], modelLam0.w[1][0]))

1054
     # Plot data
1056 plt.figure(figsize=(5,5))
     plt.plot(lamdaValues, firstOrderLOOCV)
1058 plt.title("LOOCV for different lamda values\n- First order
         polynomial fitting")
     plt.xlabel("Lamda values")
1060 plt.ylabel("Loss")
     plt.show()

1062
     # Fourth order polynomial
1064 def augment(X, max_order):
         X_augmented = X

1066
         for i in range(2, max_order+1):
1068         X_augmented = np.concatenate([X_augmented, X**i], axis=1)
         return X_augmented

1070

1072 Xnew = augment(X, 4)
     fourthOrderLOOCV = LOOCV(Xnew, t, lamdaValues)
1074 bestValueindex = np.argmin(fourthOrderLOOCV)
     bestLamda = lamdaValues[bestValueindex]

1076
     modelFourthOrder = linereg.LinearRegression(lam=bestLamda)
1078 modelFourthOrder.fit(Xnew, t)
     print("======= Fourth order polynomial fitting =======")
1080 print("The best value of lamda=%.10f and its loss=%.10f" %
         (bestLamda, fourthOrderLOOCV[bestValueindex]))
1082 print("For this model w0=%.10f, w1=%.10f, w2=%.10f, w3=%.10f, w4
     =%.10f \n" %
         (modelFourthOrder.w[0][0], modelFourthOrder.w[1][0],
1084      modelFourthOrder.w[2][0], modelFourthOrder.w[3][0],
     modelFourthOrder.w[4][0]))

1086
     # Plot data
1088 plt.figure(figsize=(5,5))
     plt.plot(lamdaValues, fourthOrderLOOCV)
1090 plt.title("LOOCV for different lamda values\n- Fourth order
         polynomial fitting")
```

```python
plt.xlabel("Lamda values")
plt.ylabel("Loss")
plt.show()
```

exercise2.py