

# MAD Assignment 3

Ask Jensen

13. december 2021

## Indhold

<b>1 Problem 1</b>	<b>1</b>
1.1 (a) . . . . .	1
1.2 (b) . . . . .	3
<b>2 Problem 2</b>	<b>4</b>
<b>3 Problem 3</b>	<b>4</b>
3.1 (a) . . . . .	4
<b>4 Problem 4</b>	<b>5</b>
4.1 (a) . . . . .	5
4.2 (b) . . . . .	5
4.3 (c) . . . . .	6

# 1 Problem 1

## 1.1 (a)

```
1 import numpy.matlib
2
3 def pca(data):
4     # Creating "clone" of matrix
5     data_cent = np.full_like(data,0)
6
7     # Iterate the matrix subtracting the mean diatom
8     # from each row
9     for i in range(780):
10         data_cent[:,i] = diatoms[:,i] - mean_diatom
11
12     # Create the covariance matrix
13     cov_matrix = np.cov(data_cent)
14
15     # Calculate the eigenvectors and eigenvalues
16     PCevals, PCevecs = np.linalg.eigh(cov_matrix)
17
18     # linalg.eigh returns the vectors and values
19     # in the wrong order.
20     # Np.flip will reverse the order so it is correct and
21     # corresponding to the exercise requirements
22     PCevals = np.flip(PCevals)
23     PCevecs = np.flip(PCevecs, axis=1)
24     return PCevals, PCevecs, data_cent
25
26 PCevals, PCevecs, data_cent = pca(diatoms)
```

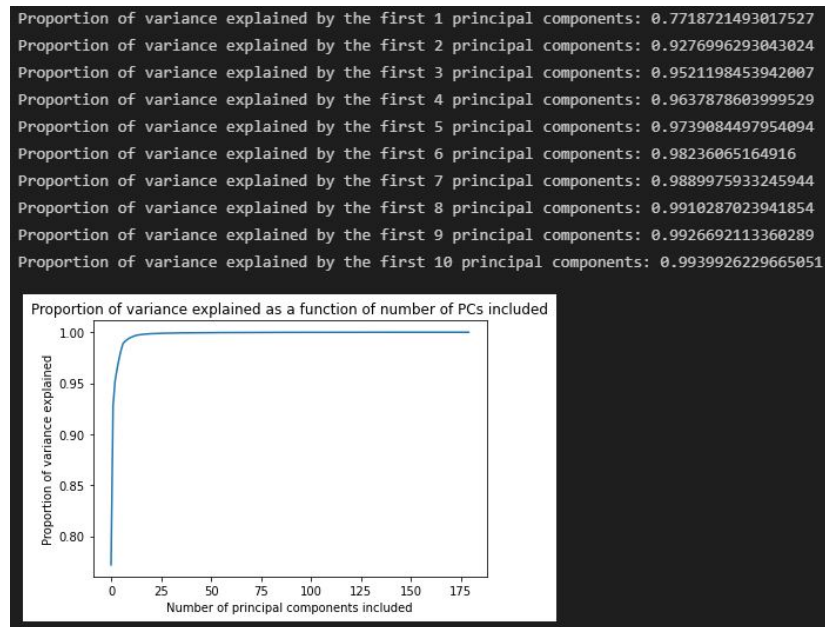


Figure 1: Values for the first 10 proportions of variance, and the corresponding graph

The proportion and the figure shows the context between the principal components and the amount of variance each iteration captures.

## 1.2 (b)

```

1      # gets the fourth eigenvector
2      e4 = PCevecs[:, 3]
3      # gets the fourth eigenvalue
4      lambda4 = PCevals[3]
5      # In case the naming std is confusing --
6      # the eigenvalues have a statistical interpretation
7      # print(std4)
8      std4 = np.sqrt(lambda4)
9
10     # Makes matrix filled with zeros
11     diatoms_along_pc = np.zeros((7, 180))
12
13     # Iterates the length of the matrix
14     # For each row, add the mean diatom with added
15     # values
16     for i in range(7):
17         diatoms_along_pc[i] = mean_diatom + ( e4 * std4 * (i-3))
18
19     # Plotting each diatom
20     for i in range(7):
21         plot_diatom(diatoms_along_pc[i])
22
23     plt.title('Diatom shape along PC1')

```

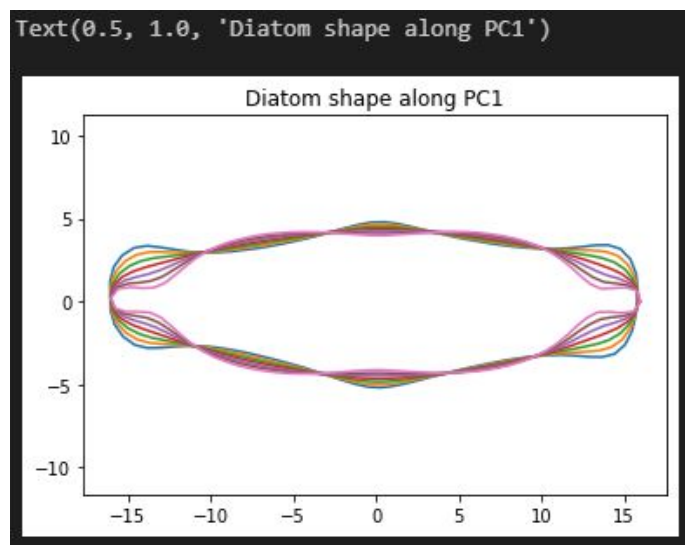


Figure 2: Plotted diatom

## 2 Problem 2

assesses the given claim  $E[(X - \mu)^4] \geq \sigma^4$

$X$  has the mean  $\mu$  and the variance  $\sigma^4$  which can be rewritten as  $\sigma^4 = (Var(X))^2$   
 $E[(X - \mu)^4]$  can be rewritten as  $E(g(x))$  where  $g(x) = (x - \mu)^4$

It is possible to Jensen's inequality if  $g''(x) \geq 0$

$$g(x) = (x - \mu)^4$$

$$g'(x) = 4(x - \mu)^3$$

$$g''(x) = 12(x - \mu)^2$$

$g(x)$  is convex, since the second derivative of the function is quadratic. Hence it will always be greater than zero. Which means that it is possible to make use of Jensen's inequality.

$$E((X - \mu)^4) \geq (E(X - \mu))^4$$

$$(E(X - \mu))^4$$

$$((E(X - \mu))^2)^2 = (Var(X))^2$$

thus the claim is true, and shown by Jensen's inequality

$$E[(X - \mu)^4] \geq \sigma^4$$

## 3 Problem 3

### 3.1 (a)

```

1  for i in range(nexp):
2      # simulates n realizations from a Gaussian
3      # with mean mu and var sigma^2
4      x = np.random.normal(mu, sigma, n)
5      # TODO: adapt for b)
6      sig = np.sqrt(np.var(x, ddof=1))
7      # computes the 0.5% quantile of a Gaussian, roughly -2.576
8      fac1 = scipy.stats.norm.ppf((1-gamma)/2, 0, 1)
9      # computes the 99.5% quantile of a Gaussian, roughly 2.576
10     fac2 = scipy.stats.norm.ppf((1-gamma)/2 + gamma, 0, 1)
11
12     # computes the 0.5 quantile using the t-test
13     fac3 = scipy.stats.t.ppf((1+gamma)/2, n-1)
14     # computes the 99.5 quantile using the t-test
15     fac4 = scipy.stats.t.ppf((1+gamma)/2-gamma, n-1)
16     xmean = np.mean(x) # Sample mean
17     a = xmean - fac2*sig/np.sqrt(n)
18     b = xmean - fac1*sig/np.sqrt(n)
19     ac = xmean - fac3*sig/np.sqrt(n) # TODO: adapt for c)
20     bc = xmean - fac4*sig/np.sqrt(n) # TODO: adapt for c)

```

```

99.0%-confidence interval:
b) Not matching in 360 (out of 10000) experiments, 3.6%
c) Not matching in 97 (out of 10000) experiments, 0.97%

```

Figur 3: Result from 10000 experiments

## 4 Problem 4

### 4.1 (a)

In this exercise I'm asked to choose the null hypothesis.

My Null hypothesis is;  $H_0 : \mu_0 = 0$  My alternative hypothesis is,  $H_A : \mu \neq \mu_0$

Which means, that I assume that there is no difference in flowering time, since the value  $X_3 - Y_3 = -0.5$  shows that the scientists claim does not hold for all of the samples. With the specified alternative hypothesis, I would have to perform a two-sided t-test

$\mu_0 = 0$  since the assumption is, that there is no difference between the two types of flowers

### 4.2 (b)

Performing the corresponding t-test (Assuming that I have to perform the corresponding test using my claim from (a)).

Following the "six steps" from the lecture. I will be starting from step 3 since both step one and two are defined in question (a).

The dataset is  $X_i - Y_i$ , which gives me

1	0.5	-0.5	1.5	0.5
---	-----	------	-----	-----

Calculating the observed mean:

$$\frac{1+0.5-0.5+1.5+0.5}{5} = 0.6$$

Calculating the standard deviation for my sample

$$S = \sqrt{\text{Var}(Z)}$$

$$\text{Var}(Z) = \sum_{i=1}^5 \frac{(x_i - \bar{x})^2}{n-1}$$

$$\text{Var}(Z) = \frac{(1-0.6)^2 + (0.5-0.6)^2 + (-0.5-0.6)^2 + (1.5-0.6)^2 + (0.5-0.6)^2}{4} = \frac{2.2}{4} = 0.55$$

$$S = \sqrt{0.55} = 0.7416$$

$$t = \frac{\bar{x} - \mu_0}{S/\sqrt{n}} = \frac{0.6 - 0}{0.7416/\sqrt{5}} = \frac{0.6}{0.3317} \approx 1.81$$

I've calculated  $c_1$  and  $c_2$  using the following code I made.

```

1  from scipy.stats import t
2
3  # Mean of my sample
4  mean = 0.6
5
6  # Samples (5 flowers)
7  n_samples = 5
8
9  # I need to divide alpha since I'm making use of
10 # the two-side test
11 alpha_val = 0.05 / 2
12 std_deviation = 0.7416
13
14 # Performing the t.ppf, in order to find c1 and c2
15 c1 = t.ppf(alpha_val, n_samples-1,
16            loc= mean, scale = std_deviation)
17
18 c2 = t.ppf((1 - alpha_val), n_samples-1,
19            loc= mean, scale = std_deviation)
20
21 print("c1: lower_cutoff", c1)
22 print("c2: upper_cutoff", c2)

```

which gives me:

$$c_1 \approx -1.4590 \quad \text{and} \quad c_2 \approx 2.6590$$

Since  $c_1 < t < c_2$  meaning that  $t$  is in the "acceptance" region, I can accept that the flowering time with the two types of flowers appears to be indifferent.

### 4.3 (c)

No, it is not that simple. If the scientist were to multiply the whole dataset with a number  $k$  it would change the calculations of  $t$ . The standard deviation will cause problems since the degrees of freedom will not be the same

$$Var^* = \sum_{i=1}^{k \cdot n} \frac{(x_i - \mu)^2}{kn - 1} \neq k \cdot \sum_{i=1}^n \frac{(x_i - \mu)^2}{n - 1}$$