

# Statistical Machine Learning - Exercise 3



TECHNISCHE  
UNIVERSITÄT  
DARMSTADT

Dominik Marino - 2468378, Pertamina Kunz - 2380210  
August 8, 2020

---

## Contents

---

<b>1</b>	<b>Linear Regression</b>	<b>2</b>
1.1	1a) . . . . .	2
1.2	1b) . . . . .	4
1.3	1c) . . . . .	5
1.4	1d) . . . . .	7
1.5	1e) . . . . .	10
1.6	1f) . . . . .	12
<b>2</b>	<b>Linear Classification</b>	<b>19</b>
2.1	2a) . . . . .	19
2.2	2b) . . . . .	19
<b>3</b>	<b>Principal Component Analysis</b>	<b>24</b>
3.1	3a) . . . . .	24
3.2	3b) . . . . .	25
3.3	3c) . . . . .	26
3.4	3d) . . . . .	26
3.5	3e) . . . . .	27
3.6	3f) . . . . .	28

---

## 1 Linear Regression

---

### 1.1 1a)

---

1. Explain: What is the ridge coefficient and why do we use it?

- The ridge coefficient minimize a penalized residual sum of squares and it is used in Ridge Regression (Variation of Linear Regression). We use a ridge coefficient when there are many correlated variables in a linear regression model. The coefficients can become determined and exhibit high variance.

- **Tutor solution:**

The ridge coefficient is a (typically small) positive number that determines how much Tikhonov regularization shall be used. It is used to improve numerical stability and to prevent overfitting

2. Derive the optimal model parameters by minimizing the squared error loss function

Mean squared Error:

$$L(y, \hat{y}) = \frac{1}{N} \sum_{i=1}^N (y - \hat{y}_i)^2$$

Model:

$$\begin{aligned} y_i &= x^T w + \epsilon \\ 0 &= x^T w + \epsilon - y_i \\ &= x^T w - y_i \end{aligned}$$

with loss function Mean squared Error:

$$\begin{aligned} \hat{w} &= \operatorname{argmin}_w \frac{1}{2} \|x^T w - y_i\|^2 + \frac{\lambda}{2} \|w\|^2 \\ \nabla_w \frac{1}{2} \|x^T w - y_i\|^2 + \frac{\lambda}{2} \|w\|^2 \end{aligned}$$

$$\begin{aligned} L(w, \lambda) &= \frac{1}{2} (Xw - y)^T (Xw - y) + \frac{\lambda}{2} w^T w \\ &= \frac{1}{2} (w^T X^T - y^T) (Xw - y) + \frac{1}{2} w^T w \\ &= \frac{1}{2} (w^T X^T Xw - 2w^T X^T y + y^T y) + \frac{1}{2} w^T w \end{aligned}$$

We suppose  $w^T w \rightarrow w^2$

$$\begin{aligned} \frac{\partial L}{\partial w} &= X^T Xw - X^T y + \lambda w \\ &= (X^T X + \lambda I)w - X^T y \\ 0 &= (X^T X + \lambda I)w - X^T y \\ x^T y &= (X^T X + \lambda I)w \\ X^T y &= (X^T X + \lambda I)w \\ w &= (X^T X + \lambda I)^{-1} X^T y \end{aligned}$$

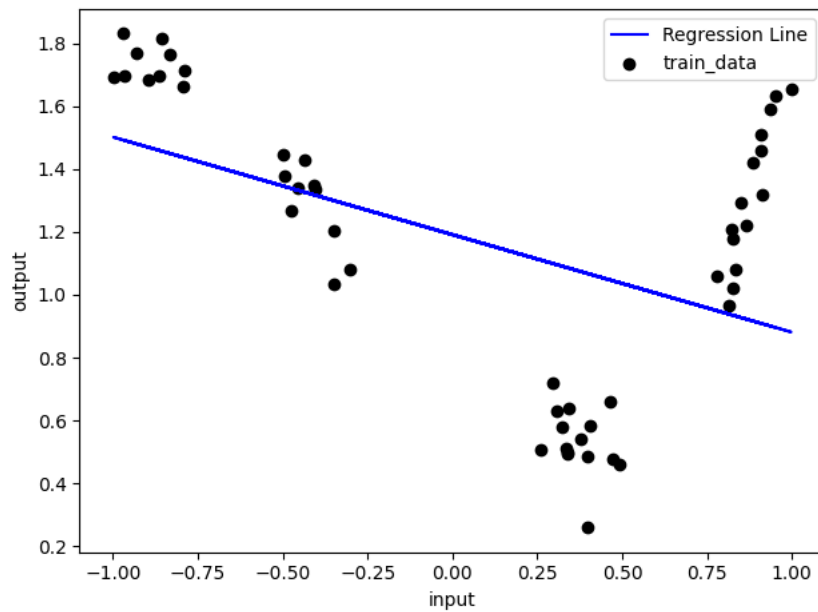
3. Report the root mean squared error of the train and test data under your linear model with linear features

Root Mean Squared Error (RMSE):

$$L(y, \hat{y}) = \sqrt{\frac{1}{N} \sum_{i=1}^N (y - \hat{y}_i)^2}$$

	Train Set	TestSet
Error	0.41217801567361084	0.38428816992597886

4. Include a single plot that shows the training data as black dots and the predicted function as a blue line



```

1 def RMSD(predictions, targets):
2     return np.sqrt(((predictions - targets) ** 2).mean())

```

Listing 1: Root Mean Squared Error

```

1 def linear_regression(X, y, alpha, d = 1):
2     X = polynomial_matrix(X, d)
3     n, m = X.shape
4     I = np.identity(m)
5     return np.linalg.inv(X.T @ X + alpha * I) @ X.T @ y

```

Listing 2: Optimization Model

```

1 def polynomial_matrix(X, d):
2     x = []
3     for i in range(0, d + 1):
4         x.append(np.power(X, i))
5     return np.asarray(x).T

```

Listing 3: Compute Polynomial Matrix

```

1 def compute_prediction(X, w_hat, d=1):
2     return polynomial_matrix(X, d) @ w_hat

```

Listing 4: Compute prediction y

```

1 w_hat = linear_regression(X_train, y_train, ridge_coefficient)
2 y_pred_train = compute_prediction(X_train, w_hat)
3 e_train = RMSD(y_pred_train, y_train)
4
5 y_pred_test = compute_prediction(X_test, w_hat)
6 e_test = RMSD(y_pred_test, y_test)
7
8 w_hat_test = linear_regression(X_test, y_test, ridge_coefficient)
9 y_pred_test = compute_prediction(X_test, w_hat_test)
10 e_test_ = RMSD(y_pred_test, y_test)
11
12 print("RMSD Train {}, \t RMSD Test {}, {}".format(e_train, e_test, e_test_))
13 plt.figure()
14 plt.plot(X_train, y_pred_train, c='b', label='Regression Line')
15 show_data(X_train, y_train, 'train_data')
16 plt.show()

```

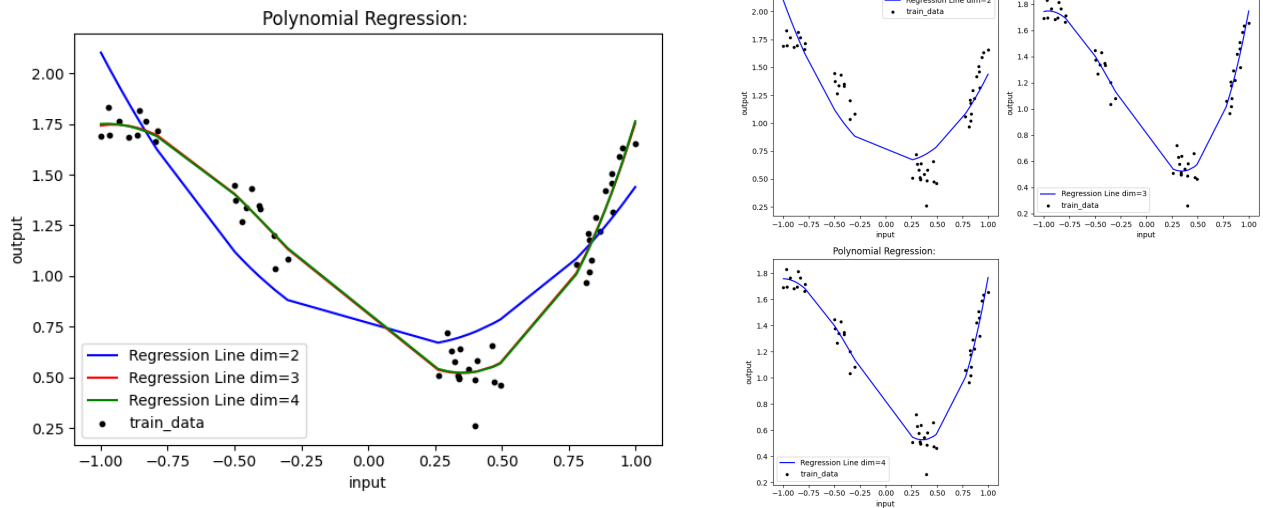
Listing 5: main

## 1.2 1b)

- Report the root mean squared error of the training data and of the testing data under your model with polynomial features

degree	Trainingsdata	Testdata
2	0.2120144726596861	0.21687242714148738
3	0.08706821295481752	0.10835803719738046
4	0.08701261306638178	0.1066623982096429

- Include a single plot that shows the training data as black dots and the predicted function as a blue line.



- Why do we call this method linear regression despite using polynomials?

- This model depends on a linear model. We are using polynomial but it has no influence on our regression coefficient. The polynomials do not change the model but can be described as a multiple combination of linear regression coefficients.
- Tutor solution:**

We still call this linear regression because the data model  $y_i = \mathbf{p}(x_i)^T \mathbf{w} + \epsilon_i$  remains linear with respect to the parameter  $\mathbf{w}$ . This is particularly important for Bayesian linear regression.

```

1 #Root Mean Squared Error
2 def RMSD(predictions, targets):
3     return np.sqrt(((predictions - targets) ** 2).mean())
4
5 def show_data(x, y, name):
6     plt.scatter(x, y, color='k', label=name, s=10)
7     plt.legend()
8     plt.xlabel('input')
9     plt.ylabel('output')
10
11 def linear_regression(X, y, alpha, d = 1):
12     X = polynomial_matrix(X, d)
13     n, m = X.shape
14     I = np.identity(m)
15     return np.linalg.inv(X.T @ X + alpha * I) @ X.T @ y
16
17 def compute_prediction(X, w_hat, d=1):
18     return polynomial_matrix(X, d) @ w_hat
19
20 def polynomial_matrix(X, d):
21     x = []
22     for i in range(0, d + 1):
23         x.append(np.power(X, i))
24     return np.asarray(x).T
25

```

Listing 6: Same like above

```

1 degrees = [2, 3, 4]
2 color_map = ['b', 'r', 'g']
3

```

```

4 plt.subplots_adjust(wspace=0.5, hspace=0.5)
5 plt.figure(figsize=(10, 10))
6
7 i = 221
8 show_data(X_train, y_train, 'train_data')
9 for d, color in zip(degrees, color_map):
10     w_hat = linear_regression(X_train, y_train, ridge_coefficient, d)
11     y_pred = compute_prediction(X_train, w_hat, d)
12     error = RMSD(y_pred, y_train)
13     print('1b) Error: ', error)
14
15     sorted_zip = sorted(zip(X_train, y_pred))
16     x, y = zip(*sorted_zip)
17     plt.subplot(i)
18     i = i + 1
19     show_data(X_train, y_train, 'train_data')
20     plt.plot(x, y, c=color, label='Regression Line dim=' + str(d))
21     plt.legend()
22 plt.title(label='Polynomial Regression: ')
23 plt.show()
24

```

Listing 7: main

### 1.3 1c)

- For each polynomial degree, report the average train, validation and test RMSEs among all folds

Dimension	Training	Validation	Test
2	0.20943990135161356	0.22488505597839778	0.21835094011192718
3	0.08620813857069498	0.09271100111785265	0.10927671570049995
4	0.0854725162035436	0.09841566883354019	0.10867173876433567

- Explain: Do the resulting numbers meet your expectations? Why (not)?

- The resulting numbers meet my expectation, because if we look up to the graphs in 1b, we trained our model over all the data. With a higher dimension we get a smoother line to fit the data better and it leads to a smaller error. We see also with a dimension of 4 that the Trainingserror gets smaller and the error of the validation set and trainingsset get higher. So it may lead to overfit.

**Tutor solution:**

As the polynomial degree increases, we increase the expensiveness of the model. Therefore, we expect the train RMSE to decrease and the validation RMSE to increase eventually due to overfitting. The results match our expectations because the train RMSE consistently decreases as we increase the degree of the polynomial, but the validation RMSE starts to increase again as we increase the degree from 3 to 4

- Which polynomial degree should be chosen for the given data? Why?

- Wie should choose a degree of 3. The average errors in Validation and Test is the error with the dimension 3 the smallest value. With a dimension of 4 the Trainingserror gets at its lowest, but we dont matter the error. We are more interested in the Validation- and Testerror. We also see, that the error of the training-set gets lower with a higher dimension meanwhile the validaiaon- and test-set-error gets an higher error. Therefore, we can support our assumption that with a higher dimension it leads to overfit.

**Tutor solution:**

Based on the validation results, we should choose 3 as the polynomial degree because it comes with the lowest validation RMSE. Although train RMSE and test RMSE are lower for 4, the train RMSE is prone to overfitting and the test RMSE must not be used for any model selection. Using the test data for model selection violates the 'golden rule' of machine learning.

```

1 def cross_validation(X_data, y_data, i, size=0.2):
2     train_size = len(y_data) * size
3     X_val, y_test, X_train, y_val = [], [], [], []
4
5     min = round(train_size * i)
6     max = round(train_size * i + train_size)
7     for idx, (X, y) in enumerate(zip(X_data, y_data)):
8         if idx >= min and idx < max:
9             X_val.append(X)
10            y_test.append(y)
11        else:
12            X_train.append(X)

```

```

13     y_val.append(y)
14
15     return np.asarray(X_train), np.asarray(X_val), np.asarray(y_val), np.asarray(y_test)
16

```

Listing 8: Cross-Validation: Separate the data into a training and validation set

```

1  #Root Mean Squared Error
2  def RMSD(predictions, targets):
3      return np.sqrt(((predictions - targets) ** 2).mean())
4
5  def show_data(x, y, name):
6      plt.scatter(x, y, color='k', label=name, s=10)
7      plt.legend()
8      plt.xlabel('input')
9      plt.ylabel('output')
10
11 def linear_regression(X, y, alpha, d = 1):
12     X = polynomial_matrix(X, d)
13     n, m = X.shape
14     I = np.identity(m)
15     return np.linalg.inv(X.T @ X + alpha * I) @ X.T @ y
16
17 def compute_prediction(X, w_hat, d=1):
18     return polynomial_matrix(X, d) @ w_hat
19
20
21 error_2, error_3, error_4 = [], [], []
22 for i in range(5):      #5 subsets
23     X_train, X_val, y_train, y_val = cross_validation(train_data[:, 0], train_data[:, 1], i)
24     for d in degrees:    #degrees = [2, 3, 4]
25         w_hat = linear_regression(X_train, y_train, ridge_coefficient, d)
26         y_pred = compute_prediction(X_train, w_hat, d)
27         error_train = RMSD(y_pred, y_train)
28
29         y_pred = compute_prediction(X_val, w_hat, d)
30         error_val = RMSD(y_pred, y_val)
31
32         y_pred = compute_prediction(X_test, w_hat, d)
33         error_test = RMSD(y_pred, y_test)
34
35         if(d == 2):
36             error_2.append([error_train, error_val, error_test])
37         elif(d == 3):
38             error_3.append([error_train, error_val, error_test])
39         else:
40             error_4.append([error_train, error_val, error_test])
41
42 #Calculate Average of
43 error_2 = np.asarray(error_2)
44 error_3 = np.asarray(error_3)
45 error_4 = np.asarray(error_4)
46 sets = ['Train', 'Val', 'Test']
47 for i, s in zip(range(error_2.shape[1]), sets):
48     print("Average value of Error: {} --> dim2: {}, dim3: {}, dim4: {}".format(s, error_2[:, i].mean(), error_3
49    [:, i].mean(), error_4[:, i].mean()))
50

```

Listing 9: Compute Error of train-, validation- and test-set

---

## 1.4 1d)

---

1. State the posterior distribution of the model parameters  $p(w \mid X, y)$  (no derivation required)

$$p(w \mid X, y) = p(y \mid X, w, \beta)p(w \mid \alpha)$$
$$p(w \mid X, y) = \mathcal{N}(y \mid X, \beta^{-1})\mathcal{N}(0, \alpha^{-1})$$

- Tutor solution:

$$p(w \mid X, y) = \mathcal{N}(\mu_n, \Lambda_n^{-1})$$
$$\mu_n = \sigma^{-2} \Lambda_n^{-1} X^T y$$
$$\Lambda_n^{-1} = \sigma^{-2} X^T X + \lambda I$$

2. State the predictive distribution  $p(y_* \mid X_*, X, y)$  (no derivation required)

$$p(y_* \mid X_*, X, y) = \int p(y_* \mid X_*, \theta) p(\theta, X, y) d\theta$$

Predictive distribution:

$$p(y_* \mid X_*, X, y) = N(y_* \mid \mu(X_*), \sigma^2(X_*))$$
$$\mu(X_*) = \Phi^T(X_*) \left( \frac{\alpha}{\beta} I + \Phi \Phi^T \right)^{-1} \Phi^T y$$
$$\sigma^2(X_*) = \frac{1}{\beta} + \Phi^T(X_*) (\alpha I + \beta \Phi \Phi^T)^{-1} \Phi(X_*)$$

- Tutor solution:

$$p(\mathbf{y}_* \mid \mathbf{X}_*, \mathbf{X}, y) = \int p(\mathbf{y}_* \mid \mathbf{X}_*, y) p(\mathbf{w} \mid \mathbf{X}, y) d\mathbf{w}$$
$$= \mathcal{N}(\mathbf{X}_* \mu_n, \sigma^2 + \mathbf{X}_* \Lambda_n^{-1} \mathbf{X}_*^T)$$

3. Report the RMSE of the train and test data under your Bayesian model (use the predictive mean)

Train	Test
0.4121779259165973	0.38434085452132943

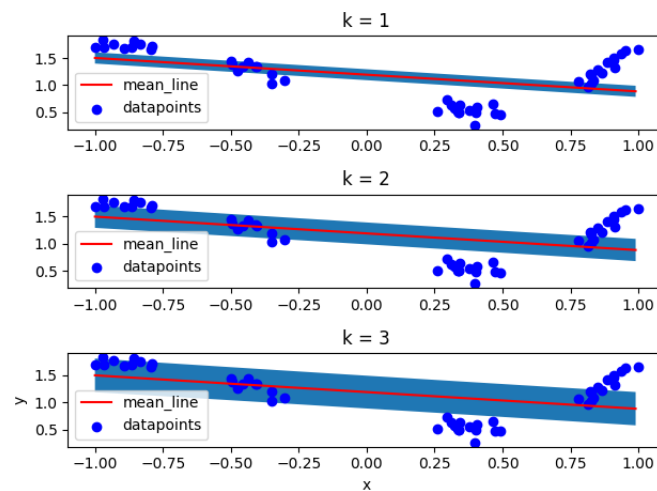
4. Report the average log-likelihood of the train and test data under your Bayesian model

Train	Test
-87.70894140971299	-90.76197537573506

- Tutor Solution:

Train	Test
-6.83469956991346	-5.774748828572732

- Include a single plot that shows the training data as black dots, the mean of the predictive distribution as blue line and 1, 2 and 3 standard deviations of the predictive distribution in shades of blue (you can use matplotlib's fill between function for that)



- Explain the differences between linear regression and Bayesian linear regression

- In Bayesian regression, we look for a predictive probability or the sampling to look for posterior probability distribution of the model. So if  $x$  is a new point, we look up for the probability of the  $y$ -value corresponding to the given  $x$ . In linear regression we are using a linear function model to estimate the unknown parameters from its data.

**Tutor solution:**

In linear regression, we compute the optimal value for the parameters  $\hat{\mathbf{w}}$  by setting the gradient of the squared error loss function to zero. Problematically, this is equivalent to the maximum likelihood point estimate under Gaussian assumptions. However, in Bayesian linear regression, instead of the maximum likelihood estimate, we compute the full posterior distribution of the parameters  $\mathbf{w}$  using Bayes rule. To the end, we define the data likelihood  $p(\mathcal{D}) = \int p(\mathcal{D} | \mathbf{w})p(\mathbf{w})d\mathbf{w}$ . Main take-away: linear regression computes a single vector for  $\mathbf{w}$ , Bayesian linear regression computes a full probability distribution for  $\mathbf{w}$

```

1 def log_likelihood(p):
2     return np.sum(np.log(p)) / len(p)
3
4 def likelihood(mu, sigma, x):
5     factor = 1 / np.sqrt(2*np.pi*sigma**2)
6     return factor * np.exp(-0.5 * ((x - mu) / sigma)**2)
7
8
9
10 def bayesian_linear_regression(phi, y, alpha, beta):
11     print(alpha, beta, " alpha and beta")
12     n, m = phi.shape
13     I = np.identity(m)
14
15     mean_matrix = np.linalg.inv((alpha / beta) * I + phi.T @ phi) @ phi.T @ y
16     variance_matrix = np.linalg.inv(alpha*I + beta * phi.T @ phi)
17
18     def predictive_mean(x):
19         return x.T @ mean_matrix
20
21     def predictive_variance(x):
22         return 1/beta + x.T @ variance_matrix @ x
23
24     means = []
25     variance = []
26     likeli = []
27     for x in phi:
28         m = predictive_mean(x)
29         means.append(m)
30         v = predictive_variance(x)
31         variance.append(v)
32         likeli.append(likelihood(m, np.sqrt(v), x))
33     std = np.asarray(np.sqrt(variance))
34     means = np.asarray(means)

```



```

35
36 print("Average Error: ", RMSD(means, y))
37 print("average Likelihood: ", log_likelihood(np.asarray(likeli)[: , 1]))
38
39 c = [1, 2, 3]
40 X = np.arange(-1, 1, 0.01)
41 X = np.c_[np.ones(len(X)), X]
42
43 means = []
44 variances = []
45 for n in X:
46     means.append(predictive_mean(n))
47     variances.append(predictive_variance(n))
48 means = np.asarray(means)
49 std = np.sqrt(np.asarray(variances))
50
51 fig, axs = plt.subplots(len(c))
52 for idx, k in enumerate(c):
53     axs[idx].plot(X[:, 1], means, label='mean_line', color='r')
54     axs[idx].fill_between(X[:, 1], (means - k * std), (means + k * std))
55     axs[idx].scatter(phi[:, 1], y, label='datapoints', c='b')
56     axs[idx].set_title('k = {}'.format(k))
57     axs[idx].legend()
58     plt.xlabel('x')
59     plt.ylabel('y')
60     plt.legend()
61 plt.show()
62
63

```

Listing 10: Bayesian Linear Regression

```

1 def calc_error_of_train(phi_train, phi_test, y_train, y_test, alpha, beta):
2     n, m = phi_train.shape
3     I = np.identity(m)
4
5     mean_matrix = np.linalg.inv((alpha / beta) * I + phi_train.T @ phi_train) @ phi_train.T @ y_train
6     variance_matrix = np.linalg.inv(alpha * I + beta * phi_train.T @ phi_train)
7
8     def predictive_mean(x):
9         return x.T @ mean_matrix
10
11     def predictive_variance(x):
12         return 1 / beta + x.T @ variance_matrix @ x
13
14     means = []
15     variance = []
16     likeli = []
17     for x in phi_test:
18         m = predictive_mean(x)
19         means.append(m)
20         v = predictive_variance(x)
21         variance.append(v)
22         likeli.append(likelihood(m, np.sqrt(v), x))
23     std = np.asarray(np.sqrt(variance))
24     means = np.asarray(means)
25     print("1e) Average Error: ", RMSD(means, y_test))
26     print("1e) average Likelihood: ", log_likelihood(np.asarray(likeli)[: , 1]))
27
28

```

Listing 11: Evaluation under Trained Model

```

1 ridge_coefficient = 0.01
2 X_train = train_data[:, 0]
3 y_train = train_data[:, 1]
4 X_test = test_data[:, 0]
5 y_test = test_data[:, 1]
6 bayesian_linear_regression(np.c_[np.ones(len(X_train)), X_train], y_train, ridge_coefficient, 1/0.01)
7 bayesian_linear_regression(np.c_[np.ones(len(X_test)), X_test], y_test, ridge_coefficient, 1/0.01)
8 calc_error(np.c_[np.ones(len(X_train)), X_train], np.c_[np.ones(len(X_test)), X_test], y_train, y_test,
9             ridge_coefficient, 1 / 0.01)

```

Listing 12: main

## 1.5 1e)

1. Report the RMSE of the train and test data under your Bayesian model with SE features. (1)

Train	Test
0.08188803394360888	0.16297037244258858 / sol: 0.138874367..

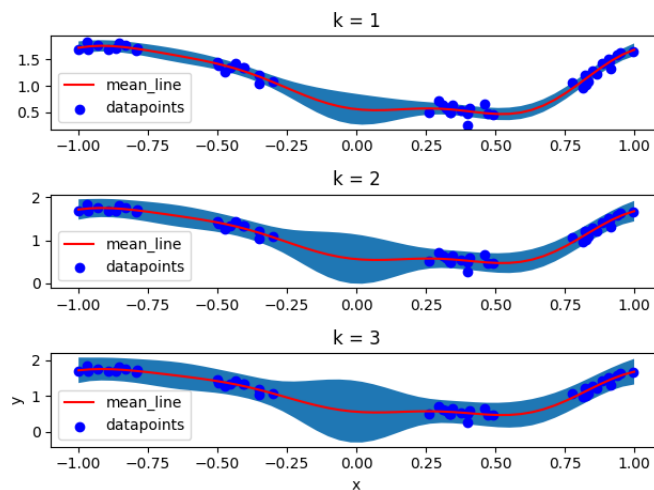
2. Report the average log-likelihood of the train and test data under your Bayesian model with SE features.

Train	Test
-38.05692344875588	-23.936399413849045

- Tutor solution:

Train	Test
1.013735727793537	0.57764976548710

3. Include a single plot that shows the training data as black dots, the mean of the predictive distribution as blue line and 1, 2 and 3 standard deviations of the predictive distribution in shades of blue (you can use matplotlib's fill between function for that)



4. How can SE features be interpreted from a statisticians point of view? What are  $\alpha$  and  $\beta$  in that context?

- it shows us the relationship between the independent variable  $x$  and the dependent variable  $y$ . In task 1b) the polynomial regression is a nonlinear model to the data, like a statistical estimation problem. The polynomial regression is a special case of multiple linear regression. With squared exponential features we can create a much smoother multiple linear regression.

- \*  $\alpha$ : distrust of a strong  $\beta$

- \*  $\beta$ : precision of the noise

– **Tutor solution:**

Squared exponential featrues are equivalent to Gaussian basis function where  $\alpha$  represents the mean of these functions and  $\beta$  is the precision, i.e. inverse variance.

```
1 def squad_exponential_matrix(k, X, beta = 10):
2     phi = []
3     for i in range(k):
4         alpha = np.ones(len(X)) * i * 0.1 - 1
5         vec = np.exp(-0.5 * beta * (X - alpha) ** 2)
6         phi.append(vec)
7     return np.asarray(phi).T
8
```

Listing 13: create SE matrix

```
1 def squad_exponential_matrix(k, X, beta = 10):
2     phi = []
3     for i in range(k):
4         alpha = np.ones(len(X)) * i * 0.1 - 1
5         vec = np.exp(-0.5 * beta * (X - alpha) ** 2)
6         phi.append(vec)
7     return np.asarray(phi).T
8
9 def log_likelihood(p):
10     return np.sum(np.log(p)) / len(p)
11
12 def likelihood(mu, sigma, x):
13     factor = 1 / np.sqrt(2*np.pi*sigma**2)
14     return factor * np.exp(-0.5 * ((x - mu)/ sigma)**2)
15
16 def bayesian_regression(phi, X_data, y, alpha, beta):
17     n, m = phi.shape
18     I = np.identity(m)
19     mean_matrix = np.linalg.inv((alpha / beta) * I + phi.T @ phi) @ phi.T @ y
20     variance_matrix = np.linalg.inv(alpha * I + beta * phi.T @ phi)
21
22     def predictive_mean(x):
23         return x.T @ mean_matrix
24
25     def predictive_variance(x):
26         return 1 / beta + x.T @ variance_matrix @ x
27
28     means = []
29     variances = []
30     likeli = []
31     for row in phi:
32         m = predictive_mean(row)
33         means.append(m)
34         v = predictive_variance(row)
35         variances.append(v)
36         likeli.append(likelihood(m, np.sqrt(v), row))
37
38     means = np.asarray(means)
39     print("1e) Average Error: ", RMSD(means, y))
40     print("1e) average Likelihood: ", log_likelihood(np.asarray(likeli)[: , 1]))
41
42     c = [1, 2, 3]
43     X_val = np.linspace(-1, 1, len(y))
44     X_val = squad_exponential_matrix(20, X_val)
45
46     means = []
47     variances = []
48     for x in X_val:
49         means.append(predictive_mean(x))
50         variances.append(predictive_variance(x))
51
52     means = np.asarray(means)
53     std = np.sqrt(np.asarray(variances))
```

```

54 X = np.linspace(-1, 1, len(y))
55
56
57 fig, axs = plt.subplots(len(c))
58 for idx, k in enumerate(c):
59     axs[idx].plot(X, means, label='mean_line', color='r')
60     axs[idx].fill_between(X, (means - k * std), (means + k * std))
61     axs[idx].scatter(X_data, y, label='datapoints', c='b')
62     axs[idx].set_title('k = {}'.format(k))
63     axs[idx].legend()
64     plt.xlabel('x')
65     plt.ylabel('y')
66     plt.legend()
67 plt.show()
68

```

Listing 14: bayesian regression

```

1 phi_train = squad_exponential_matrix(20, X_train)
2 phi_test = squad_exponential_matrix(20, X_test)
3 calc_error_of_train(phi_train, phi_test, y_train, y_test, ridge_coefficient, 1/0.01)
4 bayesian_regression(phi_train, X_train, y_train, ridge_coefficient, 1/0.01)
5 bayesian_regression(phi_test, X_test, y_test, ridge_coefficient, 1/0.01)
6

```

Listing 15: main

## 1.6 1f)

- What is the difference between the marginal likelihood  $p(y | X)$  and the likelihood  $p(y | X, w)$ ? (1)
  - The marginal likelihood provides a principled and automatic way of model comparison. So it is used to select between models. For each model we take all the posterior function to reject sampling. Likelihood describes the density of the estimation of the data.
  - tutor solution:**  
The marginal likelihood  $p(y | X)$  does not depend on the model parameters  $w$  because they have been marginalized out via integration. Intuitively, the marginal likelihood represents the probability of the data averaged over possible model parameters  $w$ , as specified by the prior  $p(w)$ . The likelihood  $p(y | X, w)$  depends on a specific  $w$ .
- For each  $\beta$ , report RMSE and average log-likelihood of the train and test data and the log-marginal likelihood

beta	RMSE Train	RMSE Test
1	0.08981902717812963	0.12423823423981933 / sol: 0.1083688207135218
10	0.08241791522223116	0.16897517789411245 / sol: 0.1425781435565231
100	0.08188803394360888	0.16297037244258858 / sol: 1.2300696160551463

$\beta$	log-likelihood (Train)	log-likelihood (Test)	$\beta$	log-marginal likelihood
1	0.9616249230979097	0.8030792476697927	1	27.170162750726966
10	1.013119934600136	0.5605016434925285	10	11.915076107400346
100	1.0401611897246368	-0.4814402200628525	100	-14.575065157319287

- For each  $\beta$ , include a single plot that shows the training data as black dots, the mean of the predictive distribution as blue line and 1, 2 and 3 standard deviations of the predictive distribution in shades of blue (you can use matplotlib's fill between function for that)

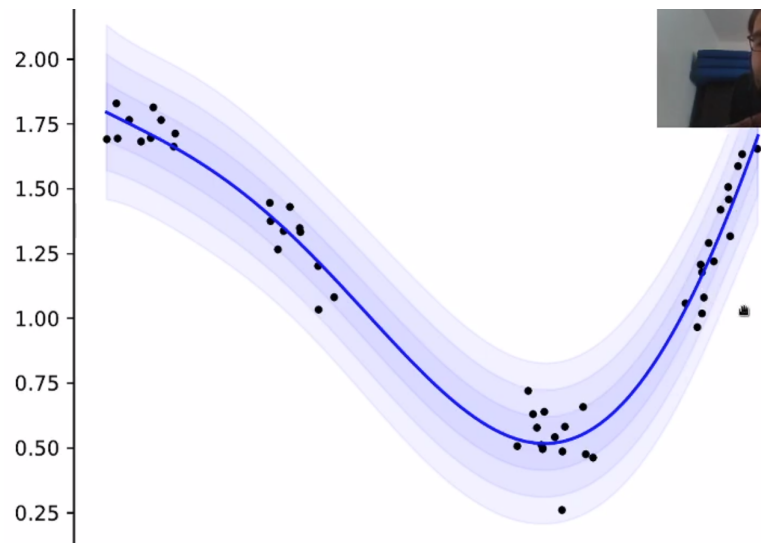


Figure 1:  $\beta = 1$

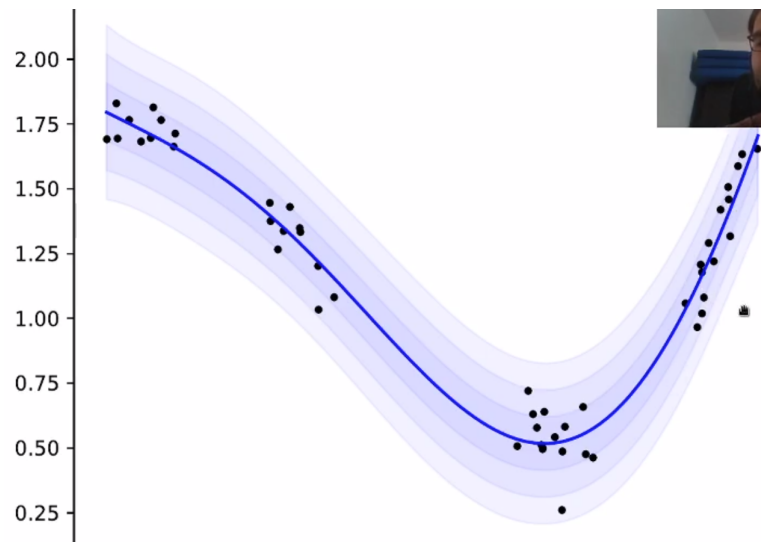


Figure 2:  $\beta = 10$

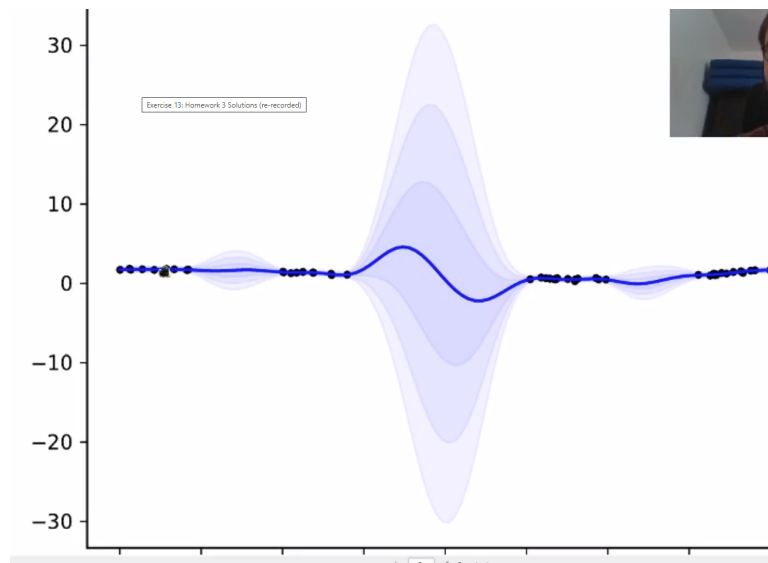


Figure 3:  $\beta = 100$

4. According to the grid search, which value for  $\beta$  is the best? Why?
  - According our result, i would take  $\beta = 100$ , because it covers better the train and test data. The Error is acceptable.
  - **Tutor solution:**  
According to the grid search,  $\beta = 1$  is the best value because it yielded the highest log-marginal likelihood
5. Compare the log-marginal likelihood values to the average train and test log-likelihood values. What do you observe? Is the log-marginal likelihood a 'good' score function compared to the train log-likelihood?
  - **Tutor solution:**  
The average train log-likelihood consistently increases as  $\beta$  increases, whereas the log-marginal likelihood decreases as  $\beta$  increases. In particular, consistent with the log-marginal likelihood, the average test log-likelihood also decreases as  $\beta$  increases. Therefore, to prevent overfitting, the log-marginal likelihood can be used for model selection without looking at the test data. This property makes the log-marginal likelihood a good score function.

```

1  beta = [1, 10, 100]
2  for b in beta:
3      baysian_regression(phi_train, X_train, y_train, 0.01, b)
4      baysian_regression(phi_test, X_test, y_test, 0.01, b)
5      calc_error_of_train(phi_train, phi_test, y_train, y_test, ridge_coefficient, b)

```

Listing 16: main

```

1  import numpy as np
2  import matplotlib.pyplot as plt
3
4  #Root Mean Squared Error
5  def RMSD(predictions, targets):
6      return np.sqrt(((predictions - targets) ** 2).mean())
7
8  def show_data(x, y, name):
9      plt.scatter(x, y, color='k', label=name, s=10)
10     plt.legend()
11     plt.xlabel('input')
12     plt.ylabel('output')
13
14  def linear_regression(X, y, alpha, d = 1):
15     X = polynomial_matrix(X, d)
16     n, m = X.shape
17     I = np.identity(m)
18     return np.linalg.inv(X.T @ X + alpha * I) @ X.T @ y
19
20  def compute_prediction(X, w_hat, d=1):
21     return polynomial_matrix(X, d) @ w_hat

```

```

22
23 def polynomial_matrix(X, d):
24     x = []
25     for i in range(0, d + 1):
26         x.append(np.power(X, i))
27     return np.asarray(x).T
28
29 def cross_validation(X_data, y_data, i, size=0.2):
30     train_size = len(y_data) * size
31     X_val, y_test, X_train, y_val = [], [], [], []
32
33     min = round(train_size * i)
34     max = round(train_size * i + train_size)
35     for idx, (X, y) in enumerate(zip(X_data, y_data)):
36         if (idx >= min and idx < max):
37             X_val.append(X)
38             y_test.append(y)
39         else:
40             X_train.append(X)
41             y_val.append(y)
42
43     return np.asarray(X_train), np.asarray(X_val), np.asarray(y_val), np.asarray(y_test)
44
45 def log_likelihood(p):
46     return np.sum(np.log(p)) / len(p)
47
48 def likelihood(mu, sigma, x):
49     factor = 1 / np.sqrt(2*np.pi*sigma**2)
50     return factor * np.exp(-0.5 * ((x - mu) / sigma)**2)
51
52
53 def bayesian_linear_regression(phi, y, alpha, beta):
54     print(alpha, beta, " alpha und beta")
55     n, m = phi.shape
56     I = np.identity(m)
57
58     mean_matrix = np.linalg.inv((alpha / beta) * I + phi.T @ phi) @ phi.T @ y
59     variance_matrix = np.linalg.inv(alpha*I + beta * phi.T @ phi)
60
61     def predictive_mean(x):
62         return x.T @ mean_matrix
63
64     def predictive_variance(x):
65         return 1/beta + x.T @ variance_matrix @ x
66
67     means = []
68     variance = []
69     likeli = []
70     for x in phi:
71         m = predictive_mean(x)
72         means.append(m)
73         v = predictive_variance(x)
74         variance.append(v)
75         likeli.append(likelihood(m, np.sqrt(v), x))
76     std = np.asarray(np.sqrt(variance))
77     means = np.asarray(means)
78
79     print("Average Error: ", RMSD(means, y))
80     print("average Likelihood: ", log_likelihood(np.asarray(likeli)[: , 1]))
81
82     c = [1, 2, 3]
83     X = np.arange(-1, 1, 0.01)
84     X = np.c_[np.ones(len(X)), X]
85
86     means = []
87     variances = []
88     for n in X:
89         means.append(predictive_mean(n))
90         variances.append(predictive_variance(n))
91     means = np.asarray(means)
92     std = np.sqrt(np.asarray(variances))
93
94     fig, axs = plt.subplots(len(c))
95     for idx, k in enumerate(c):
96         axs[idx].plot(X[:, 1], means, label='mean_line', color='r')
97         axs[idx].fill_between(X[:, 1], (means - k * std), (means + k * std))
98         axs[idx].scatter(phi[:, 1], y, label='datapoints', c='b')
99         axs[idx].set_title('k = {}'.format(k))
100         axs[idx].legend()
101         plt.xlabel('x')

```

```

102     plt.ylabel('y')
103     plt.legend()
104     plt.show()
105
106 def bayesian_regression(phi, X_data, y, alpha, beta):
107     n, m = phi.shape
108     I = np.identity(m)
109     mean_matrix = np.linalg.inv((alpha / beta) * I + phi.T @ phi) @ phi.T @ y
110     variance_matrix = np.linalg.inv(alpha * I + beta * phi.T @ phi)
111
112     def predictive_mean(x):
113         return x.T @ mean_matrix
114
115     def predictive_variance(x):
116         return 1 / beta + x.T @ variance_matrix @ x
117
118     means = []
119     variances = []
120     likeli = []
121     for row in phi:
122         m = predictive_mean(row)
123         means.append(m)
124         v = predictive_variance(row)
125         variances.append(v)
126         likeli.append(likelihood(m, np.sqrt(v), row))
127
128     means = np.asarray(means)
129     print("1e Average Error: ", RMSD(means, y))
130     print("1e average Likelihood: ", log_likelihood(np.asarray(likeli)[: , 1]))
131
132     c = [1, 2, 3]
133     X_val = np.linspace(-1, 1, len(y))
134     X_val = squad_exponential_matrix(20, X_val)
135
136     means = []
137     variances = []
138     for x in X_val:
139         means.append(predictive_mean(x))
140         variances.append(predictive_variance(x))
141
142     means = np.asarray(means)
143     std = np.sqrt(np.asarray(variances))
144
145     X = np.linspace(-1, 1, len(y))
146
147     fig, axs = plt.subplots(len(c))
148     for idx, k in enumerate(c):
149         axs[idx].plot(X, means, label='mean_line', color='r')
150         axs[idx].fill_between(X, (means - k * std), (means + k * std))
151         axs[idx].scatter(X_data, y, label='datapoints', c='b')
152         axs[idx].set_title('k = {}'.format(k))
153         axs[idx].legend()
154         plt.xlabel('x')
155         plt.ylabel('y')
156         plt.legend()
157     plt.show()
158
159 def squad_exponential_matrix(k, X, beta = 10):
160     phi = []
161     for i in range(k):
162         alpha = np.ones(len(X)) * i * 0.1 - 1
163         vec = np.exp(-0.5 * beta * (X - alpha) ** 2)
164         phi.append(vec)
165     return np.asarray(phi).T
166
167 def doBayesianRegression(phi, X_data, y, alpha, beta, sigma = 0.1):
168     n, m = phi.shape
169     I = np.identity(m)
170
171     mu = sigma**-2 * np.linalg.inv((alpha/beta)*I) @ phi.T @ y
172     sig = sigma**-2 * phi.T @ phi + (alpha / beta) * I
173
174     def logMarginalLikelihood(x):
175         return ((phi.shape[1] + 1) / 2) * np.log(alpha/beta) \
176             - (phi.shape[0] / 2) * np.log(sig ** 2) \
177             - 0.5 * np.linalg.norm(y - x @ mu)**2 / sig**2 \
178             + 0.5 * (alpha/beta) * mu.T @ mu \
179             - 0.5 * np.log(sig) \
180             - 0.5 * phi.shape[0]*np.log(2*np.pi)
181

```



```

182
183 def calc_error_of_train(phi_train, phi_test, y_train, y_test, alpha, beta):
184     n, m = phi_train.shape
185     I = np.identity(m)
186
187     mean_matrix = np.linalg.inv((alpha / beta) * I + phi_train.T @ phi_train) @ phi_train.T @ y_train
188     variance_matrix = np.linalg.inv(alpha * I + beta * phi_train.T @ phi_train)
189
190     def predictive_mean(x):
191         return x.T @ mean_matrix
192
193     def predictive_variance(x):
194         return 1 / beta + x.T @ variance_matrix @ x
195
196     means = []
197     variance = []
198     likeli = []
199     for x in phi_test:
200         m = predictive_mean(x)
201         means.append(m)
202         v = predictive_variance(x)
203         variance.append(v)
204         likeli.append(likelihood(m, np.sqrt(v), x))
205     std = np.asarray(np.sqrt(variance))
206     means = np.asarray(means)
207     print("1e) Average Error: ", RMSD(means, y_test))
208     print("1e) average Likelihood: ", log_likelihood(np.asarray(likeli)[: , 1]))
209
210
211 def main():
212     #inizialization
213     ridge_coefficient = 0.01
214     train_data = np.genfromtxt('dataSets/lin_reg_train.txt', delimiter=' ')
215     test_data = np.loadtxt('dataSets/lin_reg_test.txt', delimiter=' ')
216
217     X_train = train_data[:, 0] #np.c_[np.ones(len(train_data[:, 0])), train_data[:, 0]]
218     y_train = train_data[:, 1]
219
220     X_test = test_data[:, 0] #np.c_[np.ones(len(test_data[:, 0])), test_data[:, 0]]
221     y_test = test_data[:, 1]
222
223     """ 1a) """
224     w_hat = linear_regression(X_train, y_train, ridge_coefficient)
225     y_pred_train = compute_prediction(X_train, w_hat)
226     e_train = RMSD(y_pred_train, y_train)
227
228     y_pred_test = compute_prediction(X_test, w_hat)
229     e_test = RMSD(y_pred_test, y_test)
230
231     w_hat_test = linear_regression(X_test, y_test, ridge_coefficient)
232     y_pred_test = compute_prediction(X_test, w_hat_test)
233     e_test_ = RMSD(y_pred_test, y_test)
234
235     print("RMSD Train {}, \t RMSD Test {}, {}".format(e_train, e_test, e_test_))
236
237     plt.figure()
238     plt.plot(X_train, y_pred_train, c='b', label='Regression Line')
239     show_data(X_train, y_train, 'train_data')
240     plt.show()
241
242     """ 1b) """
243     degrees = [2, 3, 4]
244     color_map = ['b', 'b', 'b']
245
246     plt.subplots_adjust(wspace=0.5, hspace=0.5)
247     plt.figure(figsize=(10, 10))
248
249     i = 221
250     show_data(X_train, y_train, 'train_data')
251     for d, color in zip(degrees, color_map):
252         w_hat = linear_regression(X_train, y_train, ridge_coefficient, d)
253         y_pred = compute_prediction(X_train, w_hat, d)
254         error = RMSD(y_pred, y_train)
255         print('1b) Error: ', error)
256
257         sorted_zip = sorted(zip(X_train, y_pred))
258         x, y = zip(*sorted_zip)
259         plt.subplot(i)
260         i = i + 1
261         show_data(X_train, y_train, 'train_data')

```

```

262     plt.plot(x, y, c=color, label='Regression Line dim=' + str(d))
263     plt.legend()
264     plt.title(label='Polynomial Regression: ')
265     plt.show()
266
267     error_2, error_3, error_4 = [], [], []
268     """ 1c) """
269     for i in range(5):
270         X_train, X_val, y_train, y_val = cross_validation(train_data[:, 0], train_data[:, 1], i)
271         for d in degrees:
272             w_hat = linear_regression(X_train, y_train, ridge_coefficient, d)
273             y_pred = compute_prediction(X_train, w_hat, d)
274             error_train = RMSD(y_pred, y_train)
275
276             y_pred = compute_prediction(X_val, w_hat, d)
277             error_val = RMSD(y_pred, y_val)
278
279             y_pred = compute_prediction(X_test, w_hat, d)
280             error_test = RMSD(y_pred, y_test)
281
282             if(d == 2):
283                 error_2.append([error_train, error_val, error_test])
284             elif(d == 3):
285                 error_3.append([error_train, error_val, error_test])
286             else:
287                 error_4.append([error_train, error_val, error_test])
288
289             #print("Dim: {}, i: {} --> Error Train: {}, Error Val: {}, Error Test: {}".format(d, i, error_train,
290 error_val, error_test))
291 #
292 #Calculate Average of
293 error_2 = np.asarray(error_2)
294 error_3 = np.asarray(error_3)
295 error_4 = np.asarray(error_4)
296 set = ['Train', 'Val', 'Test']
297 for i, set in zip(range(error_2.shape[1]), set):
298     print("Average value of Error: {} --> dim2: {}, dim3: {}, dim4: {}".format(set, error_2[:, i].mean(),
299 error_3[:, i].mean(), error_4[:, i].
300 mean()))
301
302 """ 1d) """
303 ridge_coefficient = 0.01
304 X_train = train_data[:, 0]
305 y_train = train_data[:, 1]
306 X_test = test_data[:, 0]
307 y_test = test_data[:, 1]
308 bayesian_linear_regression(np.c_[np.ones(len(X_train)), X_train], y_train, ridge_coefficient, 1/0.01)
309 bayesian_linear_regression(np.c_[np.ones(len(X_test)), X_test], y_test, ridge_coefficient, 1/0.01)
310 calc_error_of_train(np.c_[np.ones(len(X_train)), X_train], np.c_[np.ones(len(X_test)), X_test], y_train, y_test,
311 ridge_coefficient, 1 / 0.01)
312
313 """ 1e) """
314 phi_train = squad_exponential_matrix(20, X_train)
315 phi_test = squad_exponential_matrix(20, X_test)
316 calc_error_of_train(phi_train, phi_test, y_train, y_test, ridge_coefficient, 1/0.01)
317 bayesian_regression(phi_train, X_train, y_train, ridge_coefficient, 1/0.01)
318 bayesian_regression(phi_test, X_test, y_test, ridge_coefficient, 1/0.01)
319
320 """ 1f) """
321 beta = [1, 10, 100]
322 for b in beta:
323     bayesian_regression(phi_train, X_train, y_train, 0.01, b)
324     bayesian_regression(phi_test, X_test, y_test, 0.01, b)
325     calc_error_of_train(phi_train, phi_test, y_train, y_test, ridge_coefficient, b)
326
327 if __name__ == '__main__':
328     main()

```

Listing 17: full code

---

## 2 Linear Classification

---

### 2.1 2a)

---

Explain the difference between discriminative and generative models and give an example for each case. Which model category is generally easier to learn and why?

Both predict the conditional probability but both models learn different probabilities.

A Discriminative model models the decision boundary between the classes by assuming some form for  $P(\text{Class}|\mathbf{X})$  and estimating parameters of  $P(\text{Class}|\mathbf{X})$  directly from training data. E.g. logistic regression, SVM, neural networks, nearest neighbour, Conditional Random Fields (CRF)s.

A Generative Model models the actual distribution of each class by assuming some forms for  $P(\text{Class})$  and  $P(\mathbf{X}|\text{Class})$ , estimating the parameters of  $P(\mathbf{X}|\text{Class})$  and  $P(\text{Class})$  directly from training data and using Bayes rule to calculate  $P(\text{Class}|\mathbf{X})$ . E.g. Naive Bayes, Bayesian networks, Markov random fields, Hidden Markov Models (HMM).

Generally discriminative model is easier to learn because it doesn't involve estimating the prior of the classes and we don't care whether we fit the class-conditional well.

**Tutor solution:** Discriminative models categorize a sample to a class, but unlike the generative models, can not generate samples from that class. Discriminative models learn how to categorize the data while generative models can learn the distribution of the data,  $p(\vec{X} | C_i)$ . A special case of discriminant models can model the conditional distribution  $p(C_i | \vec{X})$ . Examples of generative models are Bayesian classifiers; examples of discriminant models are linear discriminant functions. Generally, discriminative models are easier to learn as they only have to learn how to assign classes to inputs and do not learn the probability distribution of the data

---

### 2.2 2b)

---

Discriminative models map a vector  $\mathbf{x}$  to one of the  $K$  available classes. For example a Bayes classifier that maps  $\mathbf{x}$  according to a set of decision boundaries.

Given  $N$   $M$ -dimensional samples  $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ , where in this case  $N = 137$  and  $M = 2$ . We know we have  $K = 3$  classes, each class has  $N_i$  samples:  $N_1 = 50$ ,  $N_2 = 43$ , and  $N_3 = 44$ . Stacking the samples into a wide matrix  $\mathbf{X} \in \mathbb{R}^{M \times N}$  such that each column represents one sample, we want to obtain a transformation of  $\mathbf{X}$  to  $\mathbf{Y} = \mathbf{W}^T \mathbf{X}$  through projecting the samples in  $\mathbf{X}$  onto a plane (hyperplane of dimension  $K - 1 = 2$ ), where  $\mathbf{W}$  is the projection matrix used to project  $\mathbf{X}$  to  $\mathbf{Y}$ ,

$$\mathbf{W} = [\mathbf{w}_1 | \dots | \mathbf{w}_{K-1}] = [\mathbf{w}_1 | \mathbf{w}_2]$$

The optimal projection matrix  $\mathbf{W}^*$  is the one whose columns are the eigenvectors corresponding to the largest eigenvalues of the generalized eigenvalue problem

$$\mathbf{S}_W^{-1} \mathbf{S}_B \mathbf{W}^* = \lambda \mathbf{W}^*,$$

where **within-class scatter** for  $K$ -classes:

$$\mathbf{S}_W = \sum_{i=1}^K \mathbf{S}_i,$$

where

$$\mathbf{S}_i = \sum_{\mathbf{x} \in K_i} (\mathbf{x} - \boldsymbol{\mu}_i) (\mathbf{x} - \boldsymbol{\mu}_i)^T$$

the class means

$$\boldsymbol{\mu}_i = \frac{1}{N_i} \sum_{\mathbf{x} \in K_i} \mathbf{x}$$

and the **between-class scatter**

$$\mathbf{S}_B = \sum_{i=1}^K N_i (\boldsymbol{\mu}_i - \boldsymbol{\mu}) (\boldsymbol{\mu}_i - \boldsymbol{\mu})^T,$$

where the overall mean

$$\boldsymbol{\mu} = \frac{1}{N} \sum_{\forall \mathbf{x}} \mathbf{x}.$$

The provided dataset, its classification, and class means are depicted in Figure 4.

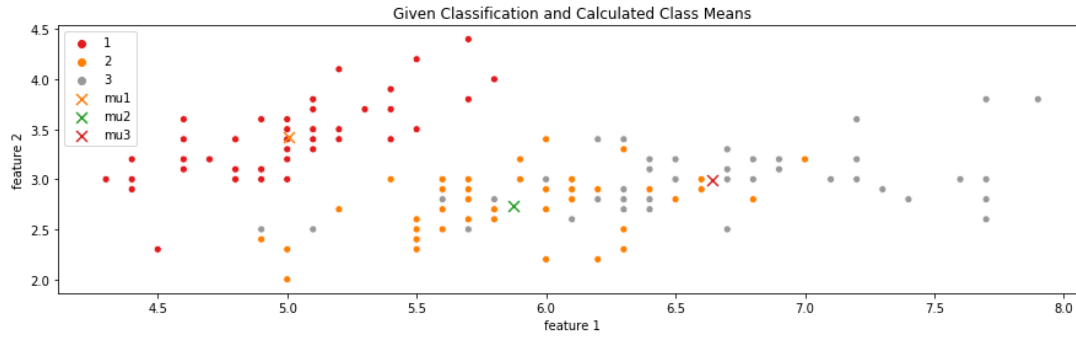


Figure 4: Given dataset and classification

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4
5 def gauss(x,m,var):
6     return np.exp(-(x-m)**2/2/var) / np.sqrt(2*np.pi*var)
7
8 def multivar_gauss(x,m,S):
9     p = m.shape[0]
10    return (2*np.pi)**(-p/2) / np.sqrt(np.linalg.det(S)) * np.exp(-0.5*np.matmul(np.transpose(x-m), np.matmul(np.linalg.
11    inv(S),(x-m)) ))
12
13 x = np.genfromtxt("ldaData.txt")
14 N = x.shape[0]
15 # No. of samples per class
16 N1 = 50
17 N2 = 43
18 N3 = 44
19
20 # data per class
21 x1 = x[0:N1,:]
22 x2 = x[N1:N1+N2,:]
23 x3 = x[N1+N2:,:]
24
25 # The class means
26 mu1 = np.mean(x1, axis=0).reshape((2,1))
27 mu2 = np.mean(x2, axis=0).reshape((2,1))
28 mu3 = np.mean(x3, axis=0).reshape((2,1))
29
30 # class vector
31 c = np.concatenate((np.repeat(1,50), np.repeat(2,43), np.repeat(3,44)))
32
33 plt.figure(figsize=(15,4))
34 sns.scatterplot(x[:,0], x[:,1], hue=c, palette='Set1', legend='full')
35 plt.scatter(mu1[0,0], mu1[1,0], marker='x', s=80, label='mu1')
36 plt.scatter(mu2[0,0], mu2[1,0], marker='x', s=80, label='mu2')
37 plt.scatter(mu3[0,0], mu3[1,0], marker='x', s=80, label='mu3')
38 plt.title('Given Classification and Calculated Class Means')
39 plt.xlabel('feature 1')
40 plt.ylabel('feature 2')
41 plt.legend()
42 plt.show()

```

Listing 18: Reading and Plotting of the Given Classification

### Using 1-D projection for new classification

The density of the projected samples onto the two eigenvectors are shown in Figure 5.

Using the **first projection vector** (corresponding to the eigenvector with the highest eigenvalue) for maximum separation, find the class of each sample according to Bayes classification,

$$\begin{aligned}
 \hat{f}(y) &= \arg \max_i p(K_i|y) \\
 &= \arg \max_i p(y|K_i) \hat{\pi}_i,
 \end{aligned}$$

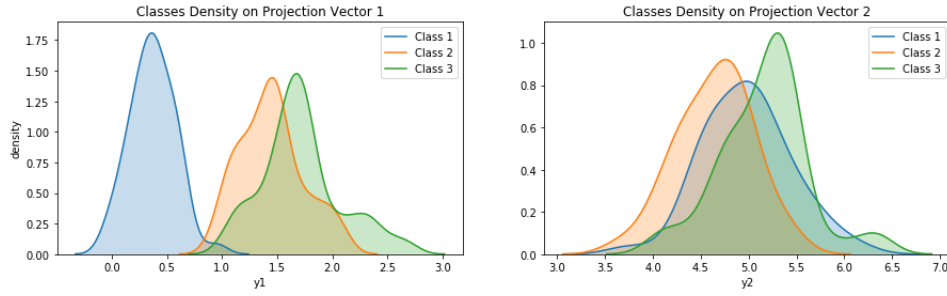


Figure 5: Projection onto the eigenvectors.

where assuming Gaussian,  $p(y|K_i) = \frac{1}{\sqrt{2\pi\sigma_i^2}} \exp\left(-\frac{(y-\mu_i)^2}{2\sigma_i^2}\right)$ , where  $\mu_i$  and  $\sigma_i^2$  are the mean and variance of the projected samples that belong to the  $i$ th class, respectively, and the prior probability for class  $i$ ,  $\pi_i$ ,  $i = 1, 2, 3$ , is estimated by the fraction of training samples of class  $i$ :

$$\begin{aligned}\hat{\pi}_1 &= \frac{N_1}{N} = \frac{50}{137} \\ \hat{\pi}_2 &= \frac{N_2}{N} = \frac{43}{137} \\ \hat{\pi}_3 &= \frac{N_3}{N} = \frac{44}{137}.\end{aligned}$$

The estimated classes are shown in Figure 6, **number of misclassified points: 30**.

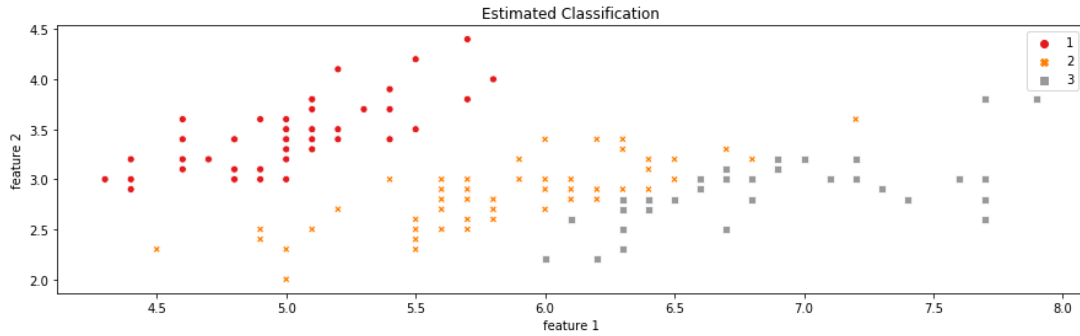


Figure 6: Classification by projection to the first eigenvector.

We observe in the figure that class 1 has clearer separation from the other classes, therefore there is only one misclassified point for this class (1 to another), while classes 2 and 3 are more mixed up, therefore 29 of the 30 errors are misclassifications between these two classes (2 to 3, 3 to 2).

```
43 # Overall mean
44 mu = np.mean(x, axis=0).reshape((2,1))
45
46 # Class covariance matrices
47 S1 = np.cov(np.transpose(x1))
48 S2 = np.cov(np.transpose(x2))
49 S3 = np.cov(np.transpose(x3))
50
51 # Within-class scatter matrix
52 S_W = S1+S2+S3
53
54 # Between-class scatter matrix
55 S_B1 = N1 * np.matmul( (mu1-mu), np.transpose(mu1-mu) )
56 S_B2 = N2 * np.matmul( (mu2-mu), np.transpose(mu2-mu) )
57 S_B3 = N3 * np.matmul( (mu3-mu), np.transpose(mu3-mu) )
58 S_B = S_B1 + S_B2 + S_B3
59
60 ISW_SB = np.dot(np.linalg.inv(S_W), S_B)
61 # The projection vectors are the eigenvectors
62 # eigvals are already sorted
```

```

63 eigvals, eigvecs = np.linalg.eig(ISW_SB)
64
65 ##### USING 1-D PROJECTION #####
66 # 3 classes => 2 projection vectors
67 w1 = eigvecs[:,0].reshape(2,1)
68 w2 = eigvecs[:,1].reshape(2,1)
69
70 # Projection onto the first projection vector
71 y_w1 = np.dot(x,w1).reshape(N)
72 # Projection onto the second projection vector
73 y_w2 = np.dot(x,w2).reshape(N)
74
75 plt.figure(figsize=(15,4))
76 plt.subplot(1,2,1)
77 sns.kdeplot(y_w1[0:N1], shade=True, label='Class 1')
78 sns.kdeplot(y_w1[N1:N1+N2], shade=True, label='Class 2')
79 sns.kdeplot(y_w1[N1+N2:], shade=True, label='Class 3')
80 plt.legend()
81 plt.title("Classes Density on Projection Vector 1")
82 plt.xlabel('y1')
83 plt.ylabel('density')
84
85 plt.subplot(1,2,2)
86 sns.kdeplot(y_w2[0:N1], shade=True, label='Class 1')
87 sns.kdeplot(y_w2[N1:N1+N2], shade=True, label='Class 2')
88 sns.kdeplot(y_w2[N1+N2:], shade=True, label='Class 3')
89 plt.legend()
90 plt.title("Classes Density on Projection Vector 2")
91 plt.xlabel('y2')
92 plt.show()
93
94 # Priors
95 pi1 = 1/N1
96 pi2 = 1/N2
97 pi3 = 1/N3
98
99 y1 = y_w1[0:N1]
100 y2 = y_w1[N1:N1+N2]
101 y3 = y_w1[N1+N2:]
102 # The new class means
103 y_mu1 = np.mean(y1)
104 y_mu2 = np.mean(y2)
105 y_mu3 = np.mean(y3)
106 # The new class variance
107 y_var1 = np.var(y1)
108 y_var2 = np.var(y2)
109 y_var3 = np.var(y3)
110
111 # Classifier
112 cnew = np.zeros(N, dtype='int8')
113 for n in range(N):
114     Post1 = gauss(y_w1[n], y_mu1, y_var1) * pi1
115     Post2 = gauss(y_w1[n], y_mu2, y_var2) * pi2
116     Post3 = gauss(y_w1[n], y_mu3, y_var3) * pi3
117     cnew[n] = np.argmax((Post1, Post2, Post3)) + 1
118
119 # Plot new classification result
120 plt.figure(figsize=(15,4))
121 sns.scatterplot(x[:,0], x[:,1], hue=cnew, style=cnew, palette='Set1', legend='full')
122 plt.xlabel('feature 1')
123 plt.ylabel('feature 2')
124 plt.title('Estimated Classification')
125 plt.legend()
126 plt.show()
127
128 print('Misclassified points: ' + str(np.sum(c!=cnew)))
129
130 wrongidx = np.argwhere(c!=cnew)
131 for i in range(len(wrongidx)):
132     print('pt (' + str(x[wrongidx[i],0]) + ', ' + str(x[wrongidx[i],1]) +
133           '), given: ' + str(c[wrongidx[i]]) +
134           ', classified as: ' + str(cnew[wrongidx[i]]))

```

Listing 19: Finding optimal projection vectors and new classification using 1D projection onto the first eigenvector.

### Using 2-D projection for new classification

Using both eigenvectors  $\mathbf{W} = [\mathbf{w}_1 | \mathbf{w}_2]$ , the projected samples are  $\mathbf{Y} = \mathbf{W}^T \mathbf{X}$ . The new classes are now determined by

$$\hat{f}(\mathbf{y}) = \arg \max_i p(\mathbf{y} | K_i) \hat{\pi}_i,$$

where  $\mathbf{y}$  is each sample in  $\mathbf{Y}$  and assuming multivariate Gaussian,  $p(\mathbf{y} | K_i) \sim \mathcal{N}_p(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$  with  $\mathbf{y}, \boldsymbol{\mu}_i \in \mathbb{R}^M$ ,  $\boldsymbol{\Sigma}_i \in \mathbb{R}^{M \times M}$ :

$$p(\mathbf{y} | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) = (2\pi)^{-\frac{M}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp \left[ -\frac{1}{2} (\mathbf{y} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1} (\mathbf{y} - \boldsymbol{\mu}_i) \right],$$

with  $\boldsymbol{\mu}_i$  and  $\boldsymbol{\Sigma}_i$  are the class mean and covariance, estimated from the projected samples.

The estimated classes are shown in Figure 7, **number of misclassified points: 19**.

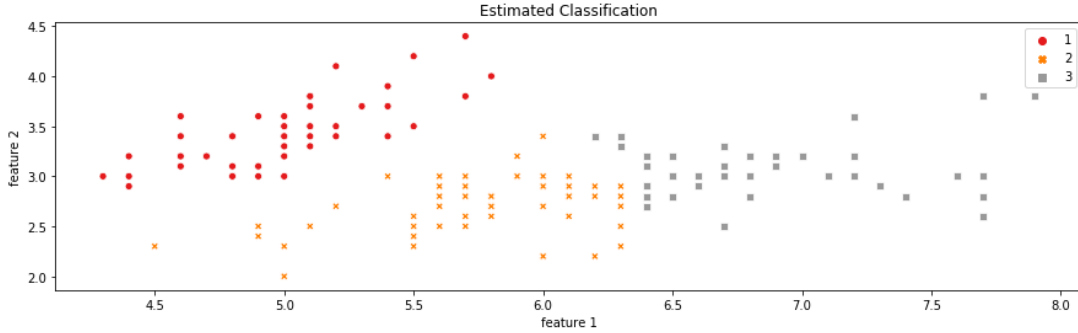


Figure 7: Classification by projection to both eigenvectors.

```

138 ##### USING 2-D PROJECTION #####
139 W = eigvecs # projection matrix
140 Y = np.dot(x,W)
141 Y1 = Y[0:N1,:]
142 Y2 = Y[N1:N1+N2,:]
143 Y3 = Y[N1+N2:,:]
144 # The new class means
145 Y_mu1 = np.mean(Y1, axis=0)
146 Y_mu2 = np.mean(Y2, axis=0)
147 Y_mu3 = np.mean(Y3, axis=0)
148 # The new class covariance
149 YS1 = np.cov(Y1.T)
150 YS2 = np.cov(Y2.T)
151 YS3 = np.cov(Y3.T)
152
153 cnew = np.zeros(N, dtype='int8')
154 for n in range(N):
155     Post1 = multivar_gauss(Y[n:], Y_mu1, YS1) * pi1
156     Post2 = multivar_gauss(Y[n:], Y_mu2, YS2) * pi2
157     Post3 = multivar_gauss(Y[n:], Y_mu3, YS3) * pi3
158     cnew[n] = np.argmax((Post1, Post2, Post3)) + 1
159
160 print('Misclassified points: ' + str(np.sum(c!=cnew)))
161
162 # Plot new classification result
163 plt.figure(figsize=(15,4))
164 sns.scatterplot(x[:,0], x[:,1], hue=cnew, style=cnew, palette='Set1', legend='full')
165 plt.xlabel('feature 1')
166 plt.ylabel('feature 2')
167 plt.title('Estimated Classification')
168 plt.legend()
169 plt.show()

```

Listing 20: New classification using 2-D projection onto both eigenvectors.

\*Note that while in the text the samples are collected in a wide matrix, in the code tall matrices are used.

---

### 3 Principal Component Analysis

---

Given a collection of data points  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}, \mathbf{x}^{(n)} \in \mathbb{R}^M$ , perform a low-dimensional representation

$$\mathbf{x}^{(n)} = \mathbf{B}\mathbf{a}^{(n)} + \mathbf{c} + \mathbf{v}^{(n)}, \quad n = 1, \dots, N,$$

where

- $\mathbf{B} \in \mathbb{R}^{M \times D}$  is a basis matrix,  $D < \min\{M, N\}$
- $\mathbf{a}^{(n)} \in \mathbb{R}^D$  is the coefficient for  $\mathbf{x}^{(n)}$ ,
- $\mathbf{c} \in \mathbb{R}^M$  is the base or mean,
- $\mathbf{v}^{(n)} \in \mathbb{R}^M$  is the noise/ modeling error.

PCA:

- Choose  $\mathbf{c} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}^{(n)}$
- Let  $\mathbf{x}'^{(n)} = \mathbf{x}^{(n)} - \mathbf{c}$ , collate the data into matrix  $\mathbf{X}' \in \mathbb{R}^{N \times M}$ , collate  $\mathbf{a}^{(n)}$  into  $\mathbf{A} \in \mathbb{R}^{N \times D}$ , and solve

$$\min_{\mathbf{A}, \mathbf{B}} \left\| \mathbf{X}' - \mathbf{A}\mathbf{B}^\top \right\|_F^2,$$

$$\text{where } \|\mathbf{Y}\|_F \text{ is the Frobenius norm } \|\mathbf{Y}\|_F = \sqrt{\sum_{i,j} |y_{ij}|^2} = \sqrt{\text{Tr}(\mathbf{Y}^\top \mathbf{Y})}.$$

In this question,  $M = 4$  for 4 features (sepal length, sepal width, petal length, and petal width) and  $N = 150$  observations.

We approached this minimization problem by using the projection of the (normalized) data onto the eigenvectors of the data's covariance matrix. The covariance matrix is used because it captures both the spread and the orientation. It therefore can be represented by a vector (eigenvector) that points into the direction of the spread of the data and by a magnitude (eigenvalue) that represents the spread in this direction.

Since we want to look for the vector that points into the direction of the largest variance, we first choose the largest eigenvalue and moving towards the smaller ones until we arrive to the amount of variance explained that we desire.

---

#### 3.1 3a)

---

Normalizing is important, because in many cases (like this one), we want to minimize the error based on the distance measure  $\|\mathbf{X}' - \mathbf{A}\mathbf{B}^\top\|_F^2$  that put the same weight (importance) to all of its dimensions. Therefore the features need to be of the same scale otherwise the features that are relatively small will have little to no contributions in shaping the optimal parameters, as the calculations will favor the features with the biggest scale.

To normalize each dimension to have zero mean and unity variance, we standardize it with:

$$x'_j = \frac{x_j - \bar{x}_j}{\text{std}(x_j)}, \quad j = 1, \dots, M$$

where  $\bar{x}_j$  is the sample mean and  $\text{std}(x_j)$  is the sample standard deviation of dimension  $j$ .

```
1 import numpy as np
2
3 def get_data(name :str) -> np.ndarray:
4     file = open(name, 'r')
5     X = []
6     y=[]
7     for i in file.readlines():
8         line = i.split(',')
9         X.append([float(x) for x in line[0:4]])
10        y.append(float(line[4]))
11    file.close()
12    return np.array(X), np.array(y)
13
14 ## 3a
15 X,y = get_data('iris.txt') #read data
16 X = (X - np.mean(X, axis=0)) / np.std(X, axis=0) #Normalize X
17 # Sanity check
18 print(np.mean(X, axis=0))
19 print(np.std(X, axis=0))
```

Listing 21: Reading and normalizing data.



### 3.2 3b)

Let  $C$  be the covariance matrix of the samples,  $\lambda_1, \dots, \lambda_M$  be the eigenvalues of  $C$  and  $\mathbf{u}_1, \dots, \mathbf{u}_M$  be the corresponding eigenvectors, i.e.

$$C\mathbf{u}_j = \lambda_j \mathbf{u}_j.$$

Order  $\lambda_j$ 's such that  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_M$ , the largest eigenvalue  $\lambda_1$  gives us the maximal variance, and the corresponding eigenvector  $\mathbf{u}_1$  gives the direction with maximal variance.

Since we want to look for the vector that points into the direction of the largest variance, we first choose the eigenvector that corresponds to the largest eigenvalue and moving towards the smaller ones until we arrive to the amount of variance explained that we desire.

Let the chosen eigenvectors be  $\mathbf{u}_1, \dots, \mathbf{u}_D$  and  $B = [\mathbf{u}_1, \dots, \mathbf{u}_D] \in \mathbb{R}^{M \times D}$ , the projection of the normalized samples onto the selected eigenvectors:

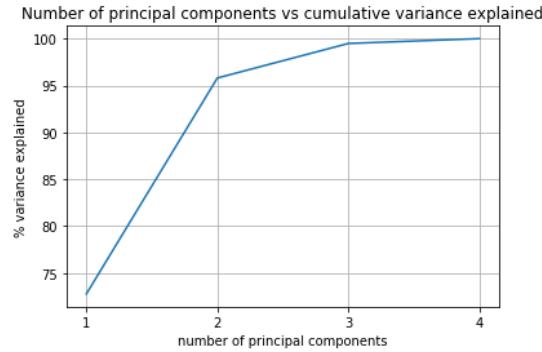
$$\mathbf{a}^{(n)} = \mathbf{x}^{(n)} B^\top.$$

The fraction of variance explained by principal components is the ratio between the variance of these principal components and the total variance. Each variance is calculated column-wise on the projected data onto the eigenvectors.

The cumulative variance explained is

$$\text{cve}(D) = \frac{\sum_{j=1}^D \text{Var}(\mathbf{a}_j)}{\sum_{j=1}^M \text{Var}(\mathbf{a}_j)} = \frac{\sum_{j=1}^D \lambda_j}{\sum_{j=1}^M \lambda_j},$$

with  $\mathbf{a}_j$  the  $j$ th column of  $A = X' B_M^\top$  with  $B_M = [\mathbf{u}_1, \dots, \mathbf{u}_M]$  and  $X' \in \mathbb{R}^{N \times M}$  the collated samples.



As we can see from the image, slightly above 95% variations in the data are explained with as few as **2 principal components**.

```

21 ##### 3b
22 import matplotlib.pyplot as plt
23
24 M = X.shape[1] # No. of features
25 Sigma = np.cov(np.transpose(X)) # Covariance matrix
26 eigvals, eigvecs = np.linalg.eig(Sigma) # eigvals are already sorted
27
28 # total variance
29 var_total = np.sum( np.var(np.dot(X,eigvecs), axis=0) )
30 # cumulative variance explained
31 cve = np.zeros(M)
32 for D in np.arange(1,M+1):
33     B = eigvecs[:,0:D] # matrix of chosen eigenvectors
34     A = np.dot(X,B) # projection of the data X to it
35     var_D = np.sum( np.var(A, axis=0) ) # explained variance
36     cve[D-1] = var_D / var_total * 100 # cumul.percent variance explained
37
38 # Plotting
39 plt.plot(np.arange(1,M+1),cve)
40 plt.title('Number of principal components vs cumulative variance explained')
41 plt.xlabel('number of principal components')
42 plt.xticks(np.arange(1,M+1))
43 plt.grid()
44 plt.ylabel('% variance explained')
45 plt.show()

```

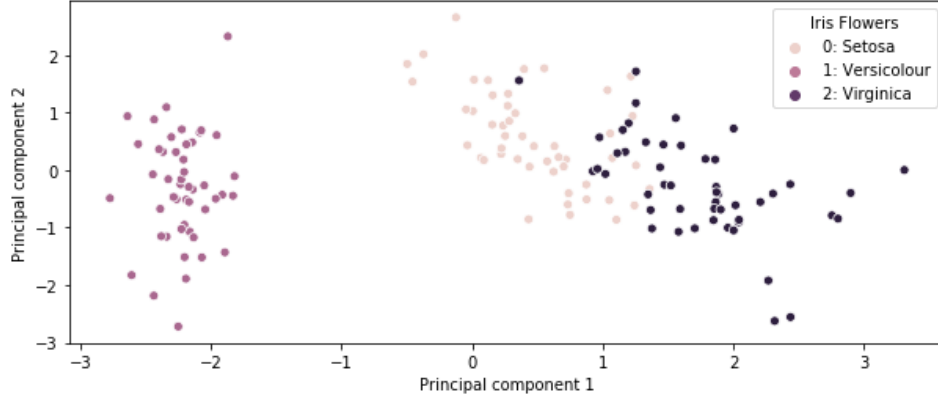
Listing 22: Finding cumulative variance explained for all possible number of principal components.

### 3.3 3c)

Using  $D = 2$  as found from the previous question, the projected samples are

$$\mathbf{A} = \mathbf{X}'([\mathbf{u}_1, \dots, \mathbf{u}_2]),$$

with  $\mathbf{A} \in \mathbb{R}^{N \times 2}$ .



The figure shows that with half the original dimension, we still have pretty good separability since the two principal components capture more than 95% of variations.

```
##### 3c
import seaborn as sns
# p=2, eigvecs calculated in 3b
B = eigvecs[:,0:2]
X_transformed = np.dot(X,B) # projection of the data X
# Plotting
plt.figure(figsize=(10,4))
ax = sns.scatterplot(X_transformed[:,0], X_transformed[:,1], hue=y, legend="full")
leg_handles = ax.get_legend_handles_labels()[0]
ax.legend(leg_handles, ['0: Setosa', '1: Versicolour', '2: Virginica'], title='Iris Flowers')
plt.xlabel('Principal component 1')
plt.ylabel('Principal component 2')
plt.show()
```

Listing 23: Transformation to lower dimension with 2 principal components.

### 3.4 3d)

Let the (unnormalized) back-transformed samples be  $\tilde{\mathbf{X}}$ ,

$$\tilde{\mathbf{X}} = (\mathbf{A}\mathbf{B}^\top) \odot \sigma_X \oplus \mu_X,$$

where  $\mathbf{B} = [\mathbf{u}_1, \dots, \mathbf{u}_D]$ , and  $\mu_X, \sigma_X \in \mathbb{R}^{N \times D}$  contain the feature-wise mean and standard deviation of the original dataset, respectively, repeated row-wise. The operators  $\odot$  and  $\oplus$  indicate element-wise multiplication and addition.

The NRMSE for the  $j$ th dimension is then <sup>1</sup>

$$\text{NRMSE}_j = \frac{\text{RMSE}_j}{x_{j,\max} - x_{j,\min}} = \sqrt{\frac{1}{N} \sum_{n=1}^N \left( x_j^{(n)} - \tilde{x}_j^{(n)} \right)^2} / (x_{j,\max} - x_{j,\min}),$$

where  $\tilde{x}_j^{(n)}$  the  $n$ th sample (row) of  $j$ th dimension (column) of  $\tilde{\mathbf{X}}$  for each  $j = 1, 2, 3, 4$ , and  $x_{j,\max} - x_{j,\min}$  is the range of the (unnormalized) samples in  $j$ th dimension.

No. of components	$x_1$	$x_2$	$x_3$	$x_4$
1	0.1040	0.1609	0.0384	0.0831
2	0.0640	0.0170	0.0379	0.0808
3	0.0086	0.0032	0.0343	0.0238
4	0	0	0	0

<sup>1</sup>[https://en.wikipedia.org/wiki/Root-mean-square\\_deviation](https://en.wikipedia.org/wiki/Root-mean-square_deviation)

```

62 ##### 3d
63 def nrmse(x, y): #normalized root mean square error between vectors X and Y
64     return np.sqrt(np.mean((x-y)**2)) / np.ptp(x)
65
66 X_ori,y_ori = get_data('iris.txt') #reread unnormalized data
67 means = np.mean(X_ori, axis=0)
68 stdevs = np.std(X_ori, axis=0)
69 # eigenvcs calculated in 3b
70 for D in np.arange(1,M+1):
71     B = eigvecs[:,0:D]
72     X_transformed = np.dot(X,B) # X is the normalized data
73     X_backtransformed = np.dot(X_transformed, B.T) * stdevs + means #unnormalize
74     print("Number of PC: " + str(D))
75     for j in np.arange(M):
76         print( 'x' + str(j+1) + ': ' + str(nrmse(X_ori[:,j], X_backtransformed[:,j])) )

```

Listing 24: NMRSE calculation

### 3.5 3e)

Let  $\mathbf{W}$  be the whitening matrix that satisfies  $\mathbf{W}^\top \mathbf{W} = \mathbf{\Sigma}^{-1}$  where  $\mathbf{\Sigma}$  is the covariance of the design matrix  $\mathbf{X}$  that contains the samples with mean 0. This results in  $\mathbf{Y} = \mathbf{W}\mathbf{X}$  whose covariance matrix is an identity matrix.

1. Explain the difference between PCA and ZCA whitening.

Mahalanobis or ZCA whitening:  $\mathbf{W} = \mathbf{\Sigma}^{-1/2}$ . It uses the sample covariance directly <sup>2</sup>.

PCA whitening: Divide each dimension of the projected samples with the square root of the respective eigenvalue, i.e.

$$[\mathbf{A}_{\text{PCAwhite}}]_j = \frac{[\mathbf{A}]_j}{\sqrt{\lambda_j}}, j = 1, \dots, D \text{ with } D \leq M, \text{ where } \mathbf{A} = \mathbf{X}'\mathbf{B} = \mathbf{X}'([\mathbf{u}_1, \dots, \mathbf{u}_D]).$$

$\mathbf{u}_j, \lambda_j$  are the eigenvectors and eigenvalues of the covariance matrix of  $\mathbf{X}'$ , and  $\mathbf{X}'$  the centered dataset with mean 0 in each dimension. PCA whitening uses the eigenvalues and eigenvectors of the (centered) sample covariance and it involves reducing the dimension at the same time.

Through singular value decomposition, PCA and ZCA whitenings differ in rotation.

2. State the equation(s) to compute the ZCA whitening parameters, given the data.

Given data  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}$  that are collated into  $\mathbf{X} \in \mathbb{R}^{N \times M}$ .

$$\hat{\mathbf{W}} = \left[ \frac{1}{N-1} \sum_{n=1}^N (\mathbf{x}_n - \hat{\boldsymbol{\mu}})(\mathbf{x}_n - \hat{\boldsymbol{\mu}})^\top \right]^{-1/2}, \text{ where } \hat{\boldsymbol{\mu}} \text{ is the sample mean.}$$

3. State the equation(s) to whiten a (new) data example  $\mathbf{x}$ , given the ZCA parameters.

ZCA-whitened data:  $\mathbf{Y} = \hat{\mathbf{W}}\mathbf{X}$

4. Compute and report the ZCA whitening parameters for the unnormalized IRIS data (including numerical values!).

The whitening parameters  $\hat{\mathbf{W}}$ :

$$\hat{\mathbf{W}} = \begin{bmatrix} 2.79802739 & -0.94581698 & -1.22308867 & 0.37350572 \\ -0.94581698 & 3.03704291 & 0.86978754 & -0.53844006 \\ -1.22308867 & 0.86978754 & 1.93387305 & -2.02303252 \\ 0.37350572 & -0.53844006 & -2.02303252 & 4.81572824 \end{bmatrix}$$

```

78 ##### 3e
79 from scipy.linalg import sqrtm
80 X,y = get_data('iris.txt') #reread data, unnormalized
81 Sigma = np.cov(np.transpose(X))
82 e = 1e-5
83 W = np.linalg.inv(sqrtm(Sigma + e))
84
85 Y = np.dot(X,W)
86 print(np.round(np.cov(Y.T),3))

```

Listing 25: ZAC whitening.

Covariance of  $\mathbf{Y}$  after whitening:

$$\begin{bmatrix} 1. & -0. & 0. & -0. \\ -0. & 1. & 0. & -0. \\ 0. & 0. & 1. & 0. \\ -0. & -0. & 0. & 1. \end{bmatrix}$$

<sup>2</sup><https://martin-thoma.com/zca-whitening/>

### Tutor solution:

- In PCE whitening, the data transformation to identity covariance is not unique. Any orthogonal matrix  $R$  can be multiplied from the left without changing the identity covariance. To make the transformation unique, in ZCA whitening,  $R$  is defined to be the matrix of Eigenvectors.

Let

$$\begin{aligned}\mu_x &= \frac{1}{n} \sum_{i=1}^n x_i \\ \Sigma_x &= \frac{1}{n} (X - \mu_x)^T (X - \mu_x) = USV \text{ (using SVD)} \\ \sqrt{S_{ii}} &:= \sqrt{S_{ii} + \epsilon}\end{aligned}$$

- then

$$\begin{aligned}\mu_{ZCA} &= \mu_x \\ \Sigma_{ZCA} &= U(\sqrt{S})^{-1}U^T\end{aligned}$$

where  $\mu_{ZCA}$  and  $\Sigma_{ZCA}$  are the ZCA whitening parameters.

To whiten an input  $x$ , we compute:

$$\hat{x} = (x - \mu_{ZCA})\Sigma_{ZCA}^T$$

---

## 3.6 3f)

Kernel PCA (KPCA) is extension of PCA using kernel method. It is useful for linearly separating data in higher dimensions, when the original data cannot be linearly separated.

When a linear kernel is used, then KPCA is (regular) PCA. **Limitations:** If the data is linearly separable, KPCA can do no better than PCA (maybe even worse due to overfitting). PCA generally is computationally cheaper than KPCA.

Consider mapping data points  $x$  to  $f(x)$  in a higher dimensional feature space. Implementation steps (assuming the samples have zero mean/centered)<sup>3</sup>:

- Choose a kernel  $k(x_m, x_n)$ , e.g.:
  - Gaussian,  $k(x_m, x_n) = \exp\left(-\frac{\|x_m - x_n\|^2}{2\sigma^2}\right)$
  - Polynomial  $k(x_m, x_n) = (x_m \cdot x_n)^k$
- Calculate

$$K = \begin{bmatrix} \dots & \dots & \dots \\ \vdots & k(x_i, x_j) & \vdots \\ \dots & \dots & \dots \end{bmatrix}_{N \times N}$$

- Find the eigenvalues  $\lambda_i$  and eigenvectors  $u_i$  of  $K$ ,  $i = 1, \dots, N$ .
- For each  $x$ , obtain its principal components in the feature space  $\sum_{n=1}^N u_{i,n} k(x, x_n)$
- **Tutor solution:**

KPCA reformulates the problem of PCA in a higher -dimensional space using a kernel function. Instead of finding the eigenvectors of the covariance of the dataset, it finds the eigenvector of the covariance of a kernel prjection  $\hat{K}$  of the dataset. The projection  $\hat{K}$  is a matrix of entries  $k_{ij} = \mathcal{K}(\vec{x}_i, \vec{y}_j)$ , where  $\mathcal{K}$  is a kernel. Applying a kernel transformation, KPCA is able to capture non-linear correlation between input variables. However, it can turn to be more computationally expensive if we have more data points than input variables (which is usually the case), as the matrix  $\hat{K}$  has  $n^2$  entries, where  $n$  is the number of data points

---

<sup>3</sup><http://fourier.eng.hmc.edu/e161/lectures/kernelPCA/node4.html>