

Statistical Machine Learning: Exercise 3



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Group A: Kexin Wang (2540047), Paul Philipp Seitz (2337506)

Summer Term 2022

Task 1: Linear Regression(44Points)

1a) Linear Features (12 Points)

1).
The ridge coefficient is a reduced factor of the simple linear regression coefficients that never attains zero values but very small values. It is used to prevent overfitting of the data and improve the numerical stability.

2).^[1]
Due to the quadratic term in the squared error function, a larger difference between the estimated value and the observed value will contribute a larger loss. First, when training the model, we want to find the parameter $\hat{\mathbf{w}}$, which minimizes the total loss over all training samples. We need to minimize the squared Euclidean norm of \mathbf{w} , scaled by the ridge coefficient λ .

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_2^2$$

Second, setting the derivative of the loss with respect to \mathbf{w} to 0 yields the analytical solution.

$$\begin{aligned} \frac{\partial}{\partial \mathbf{w}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_2^2 &= 0 \\ \Rightarrow 2\mathbf{X}^\top \mathbf{X} \mathbf{w} - 2\mathbf{X}^\top \mathbf{y} + 2\lambda \mathbf{w} &= 0 \end{aligned}$$

Finally, solving the equation for \mathbf{w} yields.

$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}.$$

3).

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 test = np.loadtxt('dataSets/lin_reg_test.txt')
5 train = np.loadtxt('dataSets/lin_reg_train.txt')
6
7 x_train = train[:, :-1]
8 y_train = train[:, -1]
9 x_test = test[:, :-1]
10 y_test = test[:, -1]
11
12 def rmse(y_pred, y):
13     return np.sqrt(np.mean((y_pred - y)**2))
14
15 def linear_ridge_regression(plot_path=None, ridge = 0.01):
16     # See 1a.2
17     n = x_train.shape[1]
```

```

18 w = np.linalg.inv(x_train.T @ x_train + ridge * np.eye(n)) @ x_train.T @ y_train
19
20 y_pred_train = x_train @ w
21 y_pred_test = x_test @ w
22
23 train_rmse = rmse(y_pred_train, y_train)
24 test_rmse = rmse(y_pred_test, y_test)
25
26 print('Train RMSE: ', train_rmse)
27 print('Test RMSE: ', test_rmse)
28
29 x = np.arange(-1, 1, 1e-2)
30 y_pred = x * w[0] + w[1]
31
32 plt.plot(train[:, :-1], y_train, 'o', color='k', label='Training data')
33 plt.plot(x, y_pred, color='tab:blue', label='Linear prediction')
34 plt.legend()
35 plt.xlabel('x')
36 plt.ylabel('y')
37 if plot_path:
38     plt.savefig(plot_path, bbox_inches='tight')
39 else:
40     plt.show()
41
42
43 # Append one dimension for the bias term
44 x_train = np.append(x_train, np.ones((x_train.shape[0], 1)), axis=1)
45 x_test = np.append(x_test, np.ones((x_test.shape[0], 1)), axis=1)
46 linear_ridge_regression('1a.png')

```

4).

Train RMSE: 0.4121780156736108
Test RMSE: 0.38428816992597875

5).

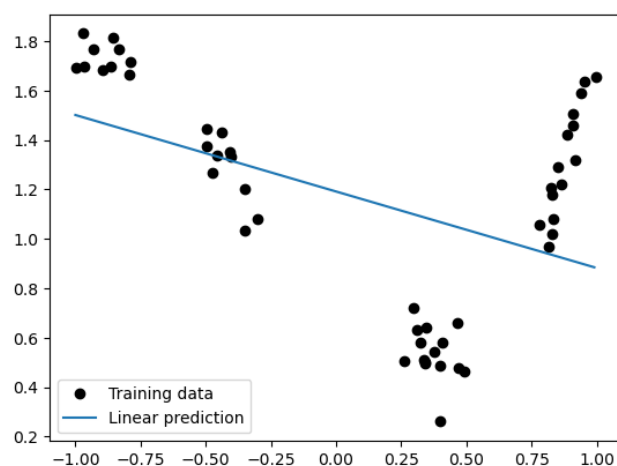


Figure 1: Training data and linear prediction

1b) Polynomial Features (8 Points)

1).

```
1 def linear_ridge_regression(phi, d, plot_path=None, ridge = 0.01):
2     # See 1a.2
3     phi_train = phi(x_train)
4     phi_test = phi(x_test)
5
6     n = phi_train.shape[1]
7     w = np.linalg.inv(phi_train.T @ phi_train + ridge * np.eye(n)) @ phi_train.T @ y_train
8
9     y_pred_train = phi_train @ w
10    y_pred_test = phi_test @ w
11
12    train_rmse = rmse(y_pred_train, y_train)
13    test_rmse = rmse(y_pred_test, y_test)
14
15    print('Train RMSE (d={}): '.format(d), train_rmse)
16    print('Test RMSE (d={}): '.format(d), test_rmse)
17
18    x = np.arange(-1, 1, 1e-2)
19    y_pred = phi(x) @ w
20
21    plt.plot(x_train, y_train, 'o', color='k', label='Training data (d={}').format(d))
22    plt.plot(x, y_pred, color='tab:blue', label='Linear prediction (d={}').format(d))
23    plt.legend()
24    plt.xlabel('x')
25    plt.ylabel('y')
26
27    if plot_path:
28        plt.savefig(plot_path, bbox_inches='tight')
29        plt.clf()
30    else:
31        plt.show()
32
33 def polynomial_feature_projection(x, degree):
34     # Polynomial projection: x -> [1, x, x^2, ..., x^degree]
35     poly = np.array([np.power(x, i) for i in range(degree+1)]).T
36     return poly.reshape([x.shape[0], degree+1])
37
38
39 for d in [2,3,4]:
40     poly_phi = lambda x: polynomial_feature_projection(x, d)
41     linear_ridge_regression(poly_phi, d, '1b_d{}.png'.format(d))
```

2).

```
Train RMSE (d=2):    0.21201447265968612
Test RMSE (d=2):    0.21687242714148727
Train RMSE (d=3):    0.08706821295481752
Test RMSE (d=3):    0.10835803719738049
Train RMSE (d=4):    0.08701261306638186
Test RMSE (d=4):    0.10666239820964747
```

3).

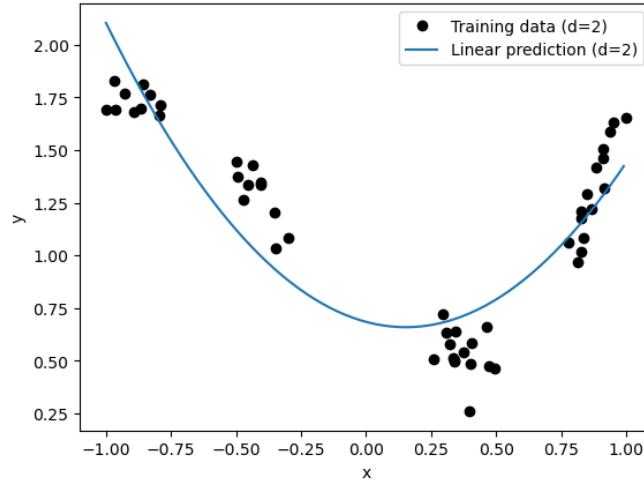


Figure 2: Training data and polynomial prediction (d=2)

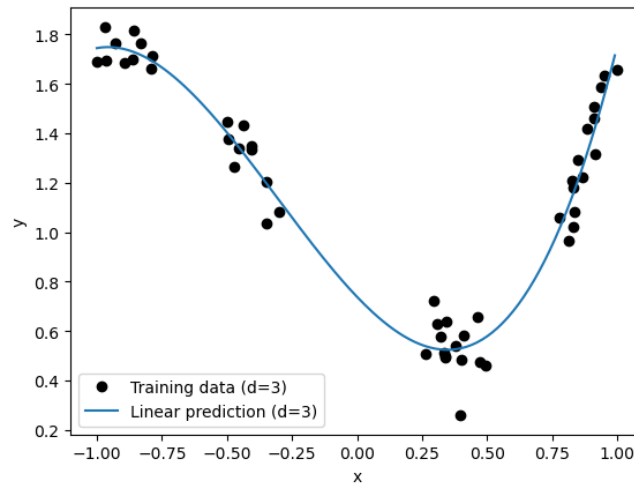


Figure 3: Training data and polynomial prediction (d=3)

- 4).
Even though the features are projected to a polynomial space as a pre-processing step, the regression itself is still a linear function as it is composed of linear operations (multiplication with weight matrix and addition of bias).

1c) Bayesian Linear Regression (11 Points)

1).

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}) = \mathcal{N}(\mu_n, \Lambda_n^{-1}).$$

$$\sigma = 0.1, \lambda = 0.01$$

$$\mu_n = \sigma^{-2} * \Lambda_n^{-1} \mathbf{X}^T \mathbf{y}$$

$$\Lambda_n^{-1} = \sigma^{-2} \mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}$$

2).

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

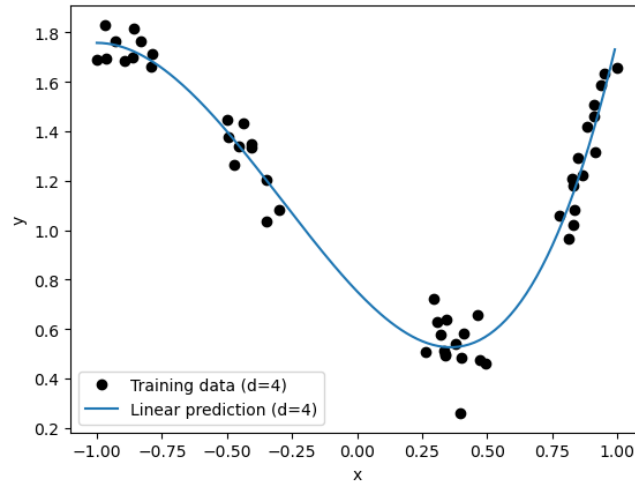


Figure 4: Training data and polynomial prediction (d=4)

3).

```

1 def extend_x(data):
2     # Append one dimension for the bias term
3     n = data.shape[0]
4     X = np.vstack([data[:, 0], np.ones((1, n))])
5     y = data[:, 1]
6
7     return X, y
8
9 def compute_w(X, y, ridge):
10    # Compute w from posterior distribution (see 1c.1)
11    return np.linalg.inv(X @ X.T + ridge) @ X @ y
12
13 def predict_y(x, w):
14    # Compute y_pred from predictive distribution (see 1c.2)
15    return np.array([x @ w for x in x.T])
16
17 def rmse(y_pred, y):
18    return np.sqrt(np.mean((y_pred - y)**2))
19
20 def compute_var(x_train, x_test, a, beta):
21    inverse = np.linalg.inv((a * np.eye(x_train.shape[0]) + beta * (x_train @ x_train.T)))
22    sigma = [(1/beta) + x @ inverse @ x.T for x in x_test.T]
23    return sigma
24
25 def gaussian_distance(y, y_pred, sigma):
26    # Compute gaussian distance for each (y, y_pred, sig)
27    p = [1 / np.sqrt(2 * np.pi * sig) * np.exp((-1) * pow((y - m), 2)) / (2 * sig)) for y, m,
28        sig in zip(y, y_pred, sigma)]
29    return p
30
31 train = np.loadtxt('dataSets/lin_reg_train.txt')
32 test = np.loadtxt('dataSets/lin_reg_test.txt')
33 x_train, y_train = extend_x(train)
34 x_test, y_test = extend_x(test)
35
36 # Compute ridge
37 d = x_train.shape[0]

```

```

38 a = 0.01
39 std = 0.1
40 beta = (std**2)**(-1)
41 ridge = a / beta * np.eye(d)
42
43 w = compute_w(x_train, y_train, ridge)
44 train_y_pred = predict_y(x_train, w)
45 test_y_pred = predict_y(x_test, w)
46
47 # 1c.4
48 print("Train RSME: ", str(rmse(train_y_pred, y_train)))
49 print("Test RSME: ", str(rmse(test_y_pred, y_test)))
50
51 # 1c.5
52 var = compute_var(x_train, x_test, a, beta)
53 train_log_likelihood = np.mean(np.log(gaussian_distance(y_train, train_y_pred, var)))
54 test_log_likelihood = np.mean(np.log(gaussian_distance(y_test, test_y_pred, var)))
55
56 print("Train Log Likelihood: ", str(train_log_likelihood))
57 print("Test Log Likelihood: ", str(test_log_likelihood))
58
59 # 1c.6
60 x = np.linspace(-1, 1, 1000)
61 x = np.vstack([x, np.ones(1000)])
62
63 y_pred = predict_y(x, w)
64
65 plt.scatter(x_train[0], y_train, c='k', label = 'Training data')
66 plt.plot(x[0], y_pred, c='tab:blue', label = 'Function prediction')
67
68 # Plot 1, 2 and 3 standard deviations in different shades
69 plt.fill_between(x[0], y_pred - std, y_pred + std, color='tab:blue', alpha=1/3)
70 plt.fill_between(x[0], y_pred - 2 * std, y_pred + 2 * std, color='tab:blue', alpha=1/3)
71 plt.fill_between(x[0], y_pred - 3 * std, y_pred + 3 * std, color='tab:blue', alpha=1/3)
72
73 plt.xlabel('x')
74 plt.ylabel('y')
75 plt.legend()
76 plt.savefig('1c.png', bbox_inches='tight')
77 plt.clf()

```

4).

Train RSME: 0.4121779259165973
 Test RSME: 0.38434085452132943

5).

Train Log Likelihood: -6.867540282525448
 Test Log Likelihood: -5.774752732594602

6).

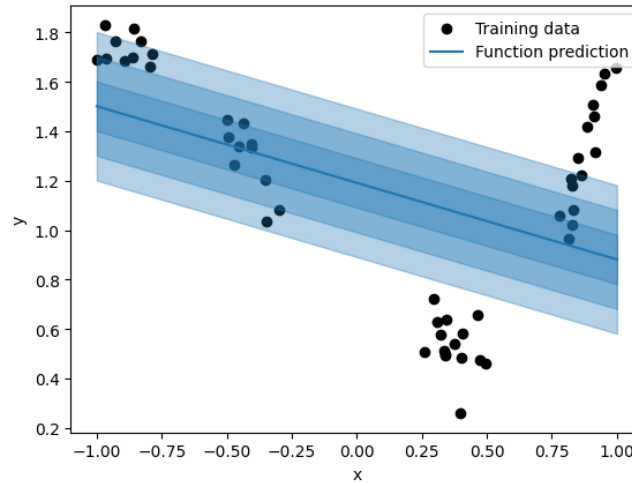


Figure 5: Training data and bayesian prediction

7).

Linear regression:

The optimal value of the parameter $\hat{\mathbf{w}}$ is calculated by setting the gradient of the squared error loss function to zero, which means that the linear regression only calculates a single vector for \mathbf{w} .

Bayesian linear regression:

We define the data likelihood $p(\mathcal{D}|\mathbf{w})$, the prior $p(\mathbf{w})$ and the marginal likelihood $p(\mathcal{D}) = \int p(\mathcal{D}|\mathbf{w})p(\mathbf{w})d\mathbf{w}$. By computing these values, we can compute the full posterior distribution of parameters \mathbf{w} . Thus the Bayesian linear regression computes a full probability distribution for \mathbf{w} .

1d) Squared Exponential Features (9 Points)

1).

```

1 def squared_exponential(data, k, beta):
2     n = data.shape[0]
3     X = data[:, 0]
4     y = data[:, 1]
5
6     alphas = [j * 0.1 - 1 for j in range(k)]
7     X_se = np.array([np.exp(-beta/2 * (Xi - alphaj)**2) for Xi in X for alphaj in alphas])
8
9     # Append one dimension for bias
10    X_se = np.vstack([X_se.reshape((n,k)).T, np.ones(n)])
11    return X_se, y
12
13
14 train = np.loadtxt('dataSets/lin_reg_train.txt')
15 test = np.loadtxt('dataSets/lin_reg_test.txt')
16
17 a = 0.01
18 beta = 10
19 std = 0.1
20 k = 20

```

```

21 ridge = a / beta * np.eye(k + 1)
22
23 x_train, y_train = squared_exponential(train, k, beta)
24 x_test, y_test = squared_exponential(test, k, beta)
25
26 # 1d.3
27 w = compute_w(x_train, y_train, ridge)
28 train_y_pred = predict_y(x_train, w)
29 test_y_pred = predict_y(x_test, w)
30
31 # 1d.4
32 print("Train RSME: ", str(rmse(train_y_pred, y_train)))
33 print("Test RSME: ", str(rmse(test_y_pred, y_test)))
34
35 # 1d.5
36 var = compute_var(x_train, x_test, a, beta)
37 train_log_likelihood = np.mean(np.log(gaussian_distance(y_train, train_y_pred, var)))
38 test_log_likelihood = np.mean(np.log(gaussian_distance(y_test, test_y_pred, var)))
39
40 print("Train Log Likelihood: ", str(train_log_likelihood))
41 print("Test Log Likelihood: ", str(test_log_likelihood))
42
43
44 # Restore non-SE data for visualization
45 x_train, y_train = extend_x(train)
46
47 x = np.linspace(-1, 1, 1000)
48 x = np.vstack([x, np.ones(1000)])
49 x_se = squared_exponential(x.T, k, beta)[0]
50 y_pred = predict_y(x_se, w)
51
52 sigma = np.sqrt(compute_var(x_se, x_se, a, beta))
53
54 plt.scatter(x_train[0], y_train, c='k', label = 'Training data')
55 plt.plot(x[0], y_pred, c='tab:blue', label = 'Function prediction')
56
57 plt.fill_between(x[0], y_pred - std, y_pred + std, color='tab:blue', alpha=1/3)
58 plt.fill_between(x[0], y_pred - 2 * std, y_pred + 2 * std, color='tab:blue', alpha=1/3)
59 plt.fill_between(x[0], y_pred - 3 * std, y_pred + 3 * std, color='tab:blue', alpha=1/3)
60
61 plt.xlabel('x')
62 plt.ylabel('y')
63 plt.legend()
64 plt.savefig('1d.png', bbox_inches='tight')
65 plt.clf()

```

2).

Train RSME: 0.08204602273485714
 Test RSME: 0.1266499667226416

3).

Train Log Likelihood: -0.04591945825929642
 Test Log Likelihood: -0.07191737612058605

4).

The SE transformation leads to the best results

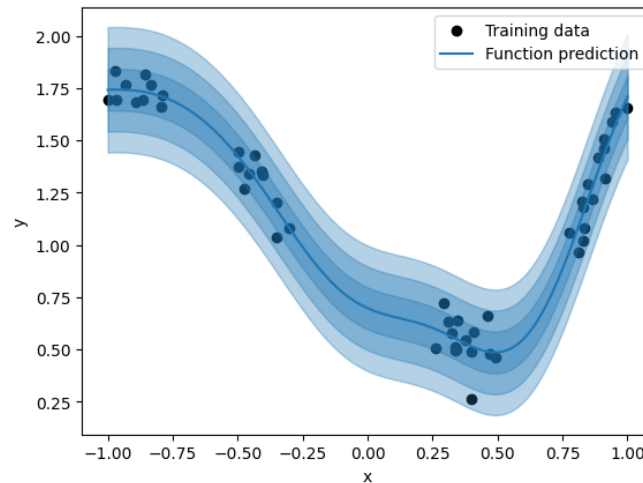


Figure 6: Training data and SE prediction

5).

The SE feature transformation resembles a Gaussian distribution.

α is like the mean μ as it gets subtracted from \mathbf{X}_i

β can be interpreted as $\frac{1}{\sigma^2}$, the multiplicative inverse of the variance σ^2

1e) Cross validation (4 Points)

Cross-validation is used to assess how well a model generalizes to data that is independent from the training set. The most common CV is K-Fold where instead of splitting the data into two, the data is split K-fold and then in each iteration some subsets (typically K-1) are used for training and the remaining subsets are used for validation.

In order to reduce variability of the partitioning K-Fold can be repeated multiple times with different partitioning and the evaluation results are combined (e.g. average). Another variant of CV is Leave-P-out which is basically N-Fold CV where N is the size of the training data and P subsets are used for validation. CV can further be enhanced by shuffling the data upfront and after each round.

Pros:

- Validates on multiple folds of data thus providing a measure for the model's generalization ability
- Can be used to balance out prediction labels if the dataset's class distribution is unbalanced
- Makes efficient use of the available data (does not require any additional validation data)
- Allows for hyperparameter optimization without needing to reserve some data exclusively for hyperparameter tuning

Cons:

- Requires that training and validation set are drawn from the same population and does therefor not work well for sequential data such as time series
- Computationally expensive because the model needs to be trained and evaluated multiple times

1f) Grid Search + Empirical Bayes (15 Points)

Task 2: Linear Classification (16 Points)

2a) Discriminative and Generative Models (4 Points)

The discriminant model learns $p(y|x)$ directly, or learns the category labels directly from the feature space. The common discriminant models are logistic regression.

The joint probability distribution $P(X, Y)$ is learned from the data, and then the probability distribution $P(Y|X)$ is derived from $P(Y|X) = \frac{P(X, Y)}{P(X)}$ as the model for prediction. This method represents the generative relationship between a given input X and the resulting output Y . Commonly used generative models include plain Bayesian, etc.

Generally, discriminant models are easier to learn as they only need to learn how to assign classes to inputs and do not learn the probability distribution of the data.

2b) Linear Discriminant Analysis (12 Points)

We decided to use a Maximum Likelihood discriminator^[2] for the LDA prediction

```
1 X = np.loadtxt('dataSets/ldaData.txt')
2 y = np.empty(X.shape[0], dtype=int)
3 y[0:50] = 0
4 y[50:93] = 1
5 y[93:] = 2
6
7 class_colors = ['tab:red', 'tab:green', 'tab:blue']
8 X_classes = [X[0:50, :], X[50:93, :], X[93:, :]]
9
10 # Original dataset
11 for i in range(3):
12     plt.scatter(X_classes[i][:, 0], X_classes[i][:, 1], color=class_colors[i],
13                 label='C'+str(i))
14
15 plt.title('Original Dataset')
16 plt.xlabel("x1")
17 plt.ylabel("x2")
18 plt.legend()
19 plt.savefig('2b_original.png', bbox_inches='tight')
20 plt.clf()
21
22 # LDA assumes normal distribution with equal covariance in all classes of X
23 def gaussian(data, mu, var):
24     return np.array([1/np.sqrt(2*np.pi*var) * np.exp(-(x - mu)**2 / (2*var)) for x in data])
25
26 n_classes = 3
27 n, d = X.shape
28
29 # Scattering within and between classes
30 mean = np.mean(X, axis=0).reshape(1, d)
31
32 X_means = [np.mean(C, axis=0) for C in X_classes]
33 S_W = sum(sum(np.outer(x - mu, (x - mu).T) for x in C) for C, mu in zip(X_classes, X_means))
34
35 S_B = np.zeros((2,2))
36 for class_mean in X_means:
```

```

37     S_B += (class_mean - mean) @ (class_mean - mean).T
38
39     # Computes eigenvalues and eigenvectors sorted in descending order
40     eigenvalues, eigenvectors = np.linalg.eig(np.linalg.inv(S_W) @ S_B)
41     w = eigenvectors[:, :1].T # take first dimension (linear projection)
42
43     # Linear projection of classes and their means and variances
44     linear_X_classes = [w @ C.T for C in X_classes]
45     linear_X_means = [np.mean(C) for C in linear_X_classes]
46     linear_X_variances = [np.var(C) for C in linear_X_classes]
47
48     # Assign X to the class with maximum likelihood discriminator
49     linear_X = np.concatenate([linear_X_classes[i][0] for i in range(n_classes)])
50     likelihood = [gaussian(linear_X, mean, variance) for mean, variance in zip(linear_X_means,
51                                     linear_X_variances)]
52     y_pred = np.argmax(likelihood, axis=0)
53
54     c = [class_colors[y] for y in y_pred]
55
56     plt.scatter(X[:, 0], X[:, 1], c=c)
57
58     plt.title("LDA Classification")
59     plt.xlabel("x1")
60     plt.ylabel("x2")
61     plt.savefig('2b_lda.png', bbox_inches='tight')
62     plt.clf()
63
64     print('Misclassifications:', np.sum(y_pred != y))

```

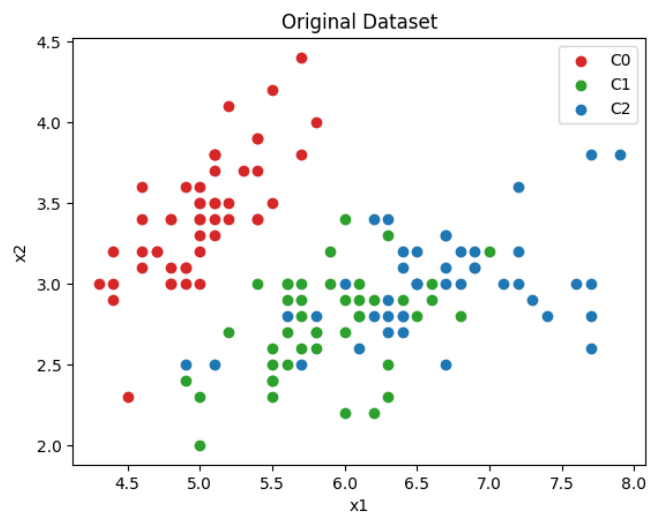


Figure 7: Original Dataset

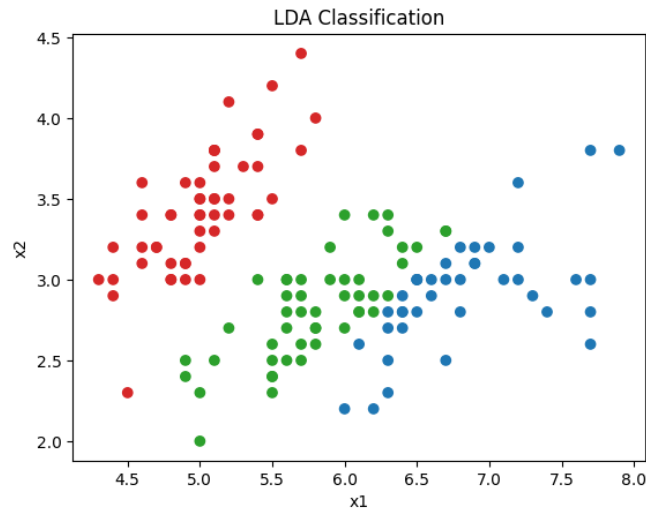


Figure 8: LDA Classification with ML Discriminator

The linear decision boundaries are easy to spot in the second plot. In total 25 samples are misclassified.

Task 3: Principal Component Analysis(28Points)

3a) Data Normalization (3 Points)

```

1 data = np.loadtxt('dataSets/iris.txt', delimiter=',')
2
3 X_original = data[:, :-1]
4 X_original_mean = np.mean(X_original, axis=0)
5 X_original_std = np.std(X_original, axis=0)
6
7 y = data[:, -1]
8
9 print('Mean (original): ', X_original_mean)
10 print('Variance (original): ', X_original_std)
11
12 # Normalization
13 X = (X_original - X_original_mean) / X_original_std
14
15 print('Mean (normalized): ', np.mean(X, axis=0))
16 print('Variance (normalized): ', np.var(X, axis=0))

```

```

Mean (original): [5.84333333 3.054      3.75866667 1.19866667]
Variance (original): [0.82530129 0.43214658 1.75852918 0.76061262]
Mean (normalized): [-1.69031455e-15 -1.63702385e-15 -1.48251781e-15 -1.62314606e-15]
Variance (normalized): [1. 1. 1. 1.]

```

Normalizing the data is important when features have different distributions to make sure that all features are on the same scale. Otherwise features with higher values by nature would have greater impact on the model than lower-value features for example but typically we don't want our model to have initial assumptions about the relevance of individual features as wrong biases can hinder the learning success.

3b) Principal Component Analysis (8 Points)

```
1 def pca(X, D=None):
2     # Principal Component Analysis with D components
3
4     # Compute eigenvalues and eigenvectors of covariance matrix
5     C = np.cov(X, rowvar=False)
6     eigenvalues, eigenvectors = np.linalg.eig(C)
7
8     # Fraction of the variance that eigenbasis with D components captures
9     explained_vars = np.cumsum(eigenvalues) / np.sum(eigenvalues)
10
11     if D is None:
12         D = np.argmax(explained_vars >= 0.95) + 1
13
14     # Take D eigenvectors with the highest eigenvalues as eigenbasis
15     B = eigenvectors[:, :D]
16     a = B.T @ X.T
17
18     return B, a, D, explained_vars
19
20
21 # Compute PCA with 95% variance explained
22 B, a, D, explained_vars = pca(X)
23 print('At least {} components are required to explain 95% of the variance.'.format(D))
24
25 # Plot the explained variance
26 plt.plot(range(1, 5), explained_vars)
27 plt.xticks(range(1, 5))
28
29 # Plot the 95% boundary
30 xx, yy = np.meshgrid(np.linspace(1, 5, 1000), np.linspace(np.min(explained_vars), 1, 1000))
31 plt.contourf(xx, yy, yy >= 0.95, 2, cmap=cm.gray, alpha=0.5)
32
33 plt.title("Cumulative variance explained")
34 plt.xlabel("N. of components")
35 plt.ylabel("% variance explained")
36 plt.savefig('3b.png', bbox_inches='tight')
37 plt.clf()
```

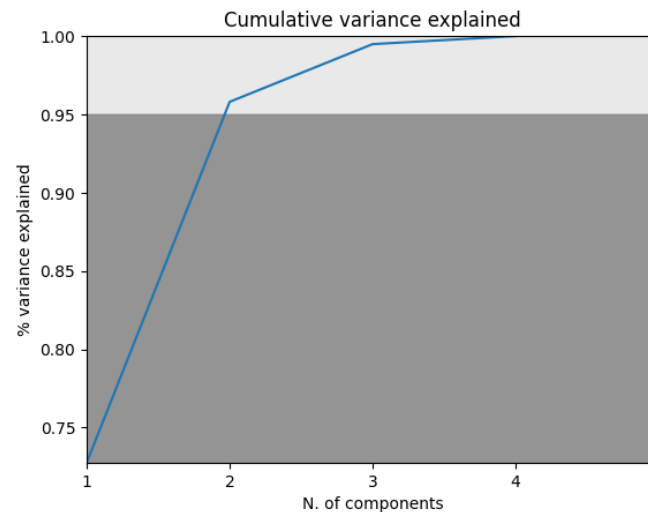


Figure 9: Cumulative variance explained by PCA

The PCA requires at least 2 components in order to explain 95% of the variance.

3c) Low Dimensional Space (6 Points)

```
1 class_colors = ['tab:red', 'tab:green', 'tab:blue']
2 c = [class_colors[int(i)] for i in y]
3
4 # Plot X w.r.t. the min. PCA basis that can explain 95% of the variance (2D)
5 # See 3b for the computation of a
6 plt.scatter(a[0], a[1], c=c)
7
8 plt.title("PCA with {} components".format(enough))
9 plt.xlabel("x1")
10 plt.ylabel("x2")
11 plt.savefig('3c.png', bbox_inches='tight')
12 plt.clf()
```

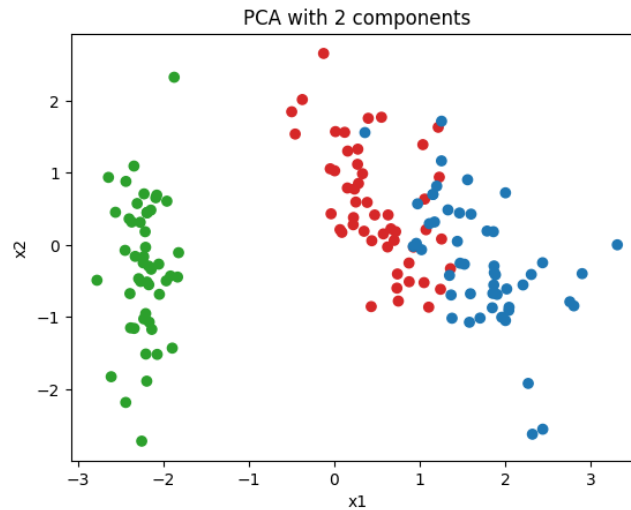


Figure 10: PCA with 2 components

3d) Projection to the Original Space (6 Points)

We considered the different normalization techniques^[3] and decided to use range normalization

```

1 def nrmse(X_pred, X):
2     rmse = np.sqrt(np.mean((X - X_pred)**2, axis=0))
3
4     # Normalize RMSE with the range
5     y_max = np.max(X, axis=0)
6     y_min = np.min(X, axis=0)
7     return rmse / (y_max - y_min)
8
9 for D in range(1, 5):
10     B, a, _, _ = pca(X, D)
11
12     # X is reconstructed from a by applying reverse transformation
13     X_reconstructed = (B @ a).T
14     nrmse_D = nrmse(X_reconstructed, X)
15
16     print('NRMSE (D={}): {}'.format(D, nrmse_D))

```

```

NRMSE (D=1): [0.10397936 0.16086196 0.03835783 0.08311765]
NRMSE (D=2): [0.06403369 0.0170341 0.03788015 0.08070068]
NRMSE (D=3): [0.00862221 0.00320869 0.03427903 0.02381893]
NRMSE (D=4): [1.42394418e-16 5.34157918e-17 1.40099782e-16 1.50415871e-16]

```

N. of components	x_1	x_2	x_3	x_4
1	0.10397936	0.16086196	0.03835783	0.08311765
2	0.06403369	0.0170341	0.03788015	0.08070068
3	0.00862221	0.00320869	0.03427903	0.02381893
4	1.42394418e-16	5.34157918e-17	1.40099782e-16	1.50415871e-16

Lower NRMSE values are an indicator for lower residual variance between the original and the reconstructed dataset. The results show more components lead to a lower reconstruction error which makes sense because more information from the original data is preserved by the PCA.

3e) Kernel PCA (5 Points)

KPCA is, in order to obtain better linear separability, the original data matrix is considered to be mapped to a new high-dimensional space, and then PCA is performed in the new high-dimensional space to reduce the dimensionality, but this processing requires a large number of computations, and the Kernel function can significantly reduce the amount of computations without affecting the effect of high-dimensional PCA.

The general process is:

- Select kernel functions: linear kernel functions, polynomial kernel functions, Gaussian kernel functions, etc.
- Construct the standard kernel matrix.

$$\tilde{K} = K - \mathbf{1}_{1/n}K - K\mathbf{1}_{1/n} + \mathbf{1}_{1/n}K\mathbf{1}_{1/n}$$

- Solve the feature vectors.

$$\tilde{K}\alpha_i = \lambda_i\alpha_i$$

- Mapping arbitrary original space data to the values of the space composed of principal components.

$$y_i = \sum_{j=1}^n \alpha_{ji} K(x, x_j), j = 1, \dots, d$$

The KPCA will turn to have more computations if we have more data points than input variables.

References

- [1] *Linear regression*. URL: https://en.wikipedia.org/wiki/Linear_regression.
- [2] *LDA Discriminator*. URL: https://en.wikipedia.org/wiki/Linear_discriminant_analysis#Discrimination_rules.
- [3] *RSME Normalization*. URL: https://en.wikipedia.org/wiki/Root-mean-square_deviation#Normalization.