

Statistical Machine Learning 2018

Assignment 4

Deadline: 11th of January 2019

Christoph Schmidl
s4226887

c.schmidl@student.ru.nl

Mark Beijer
s4354834

mbeijer@science.ru.nl

January 5, 2019

Exercise 1 - Logistic regression (weight 5)

Part 1 - The IRLS algorithm

Many machine learning problems require minimizing / maximizing some function $f(x)$. For this, an alternative to the familiar gradient descent technique, is the so called Newton-Raphson iterative method:

$$\mathbf{x}^{(n+1)} = \mathbf{x}^{(n)} - \mathbf{H}^{-1} \nabla f(\mathbf{x}^{(n)}) \quad (1)$$

where \mathbf{H} represents the Hessian matrix of second derivatives of $f(\mathbf{x})$, see Bishop, §4.3.3.

1.1.1

Derive an expression for the minimization / maximization of the function $f(x) = \sin(x)$, using the Newton-Raphson iterative optimization scheme (1), and verify (using Matlab, just up to, e.g., five iterations) how quickly it converges when starting from $x^{(0)} = 1$. What happens when you start from $x^{(0)} = -1$?

Hint: The Hessian of a 1-dimensional function $f(x)$ is just the second derivative f'' . So, the Newton-Raphson iterative method reduces in 1-d to

$$x^{(n+1)} = x^{(n)} - \frac{f'(x^{(n)})}{f''(x^{(n)})} \quad (2)$$

Answer:

Expression for the minimization / maximization of the function $f(x) = \sin(x)$:

$$\begin{aligned} x^{(n+1)} &= x^{(n)} - \frac{\sin'(x^{(n)})}{\sin''(x^{(n)})} \\ &= x^{(n)} - \frac{\cos(x^{(n)})}{-\sin(x^{(n)})} \\ &= x^{(n)} + \frac{\cos(x^{(n)})}{\sin(x^{(n)})} \end{aligned}$$

```
1 import numpy as np
2 import sympy as sp
3
4 # Exercise 1.1.1
5
6 ## See also: https://docs.sympy.org/latest/tutorial/calculus.html
7
```

```

8 def f(x):
9     return np.sin(x)
10
11 def iterative_optimization(x):
12     """ Using Newton-Raphson iterative optimization scheme """
13     return x + (np.cos(x)/np.sin(x))
14
15 x_t = 1
16 n_iterations = 5
17
18 print("Starting the optimization process with x_0: {}".format(x_t))
19
20 for i in range(n_iterations):
21     x_t = iterative_optimization(x_t)
22     print("Step {}: {}".format(i, x_t))
23
24 # Starting the optimization process with x_0: 1
25 # Step 0: 1.6420926159343308
26 # Step 1: 1.5706752771612507
27 # Step 2: 1.5707963267954879
28 # Step 3: 1.5707963267948966
29 # Step 4: 1.5707963267948966
30
31 # Starting the optimization process with x_0: -1
32 # Step 0: -1.6420926159343308
33 # Step 1: -1.5706752771612507
34 # Step 2: -1.5707963267954879
35 # Step 3: -1.5707963267948966
36 # Step 4: -1.5707963267948966

```

If we take $x_0 = 1$, then the algorithm converges after 4 steps towards 1.5707963267948966 which is an approximation of $\frac{\pi}{2}$.

If we take $x_0 = -1$, then the algorithm converges after 4 steps towards -1.5707963267948966 which is an approximation of $-\frac{\pi}{2}$.

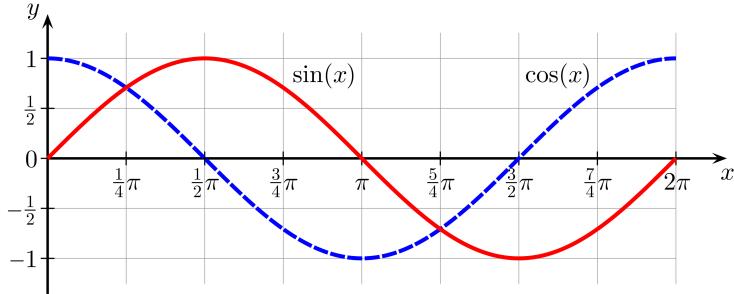


Figure 1: Sine and Cosine

When we take a look at figure 1 then we can see that the maximum and minimum of the sine function is indeed at $(x = \frac{\pi}{2}, f(x) = 1)$ and $(x = -\frac{\pi}{2}, f(x) = -1)$ with possible shiftings to the right and left based on the periodic nature of the sine function.

Note: Also see https://en.wikipedia.org/wiki/Newton%27s_method#Applications

1.1.2

We want to apply this method to the logistic regression model for classification (see Bishop, §4.3.2):

$$p(\mathcal{C}_1|\phi) = y(\phi) = \sigma(w^T \phi) \quad (3)$$

For a data set $\{\phi_n, t_n\}_{n=1}^N$, with $t_n \in \{0, 1\}$, using $y_n = p(\mathcal{C}_1|\phi_n)$ the corresponding cross entropy error function to minimize is

$$E(\mathbf{w}) = - \sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\} \quad (4)$$

With one basis function ϕ and the dummy basis function 1, the feature vector in (3) becomes $\phi = [1, \phi]^T$. The weight vector including the bias term is then also two dimensional, $\mathbf{w} = [w_0, w_1]^T$. Expressions for

the gradient $\nabla E(\mathbf{w})$ and Hessian \mathbf{H} in terms of the data set are given in Bishop, eq.4.96-98. As both are implicitly dependent on the weights \mathbf{w} , they have to be recalculated after each step: hence this is known as the 'Iterative Reweighted Least Squares' algorithm.

Consider the following data set: $\{\phi_1, t_1\} = \{0.3, 1\}$, $\{\phi_2, t_2\} = \{0.44, 0\}$, $\{\phi_3, t_3\} = \{0.46, 1\}$ and $\{\phi_4, t_4\} = \{0.6, 0\}$, and initial weight vector $\mathbf{w}^{(0)} = [1.0, 1.0]^T$.

Show using e.g. a Matlab implementation that for this situation the IRLS algorithm converges in a few iterations to the optimal solution $\hat{\mathbf{w}}^T \approx [9.8, -21.7]$, and show that this solution corresponding to a decision boundary $\phi = 0.45$ in the logistic regression model. (The IRLS algorithm should take about five lines of Matlab code inside a loop + initialization).

Answer:

$$p(\mathcal{C}_1|\phi) = y(\phi) = \sigma(w^T \phi) = \frac{1}{1 + \exp(-(w^T \phi))}$$

As stated in Bishop, page 208, eq 4.99 and 4.100: The Newton-Raphson update formula for the logistic regression model becomes

$$\begin{aligned}\mathbf{w}^{new} &= \mathbf{w}^{old} - (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T (\mathbf{y} - \mathbf{t}) \\ &= (\Phi^T \mathbf{R} \Phi)^{-1} \{ \Phi^T \mathbf{R} \Phi \mathbf{w}^{(old)} - \Phi^T (\mathbf{y} - \mathbf{t}) \} \\ &= (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T \mathbf{R} \mathbf{z}\end{aligned}$$

where \mathbf{z} is an N-dimensional vector with elements

$$\mathbf{z} = \Phi \mathbf{w}^{old} - \mathbf{R}^{-1} (\mathbf{y} - \mathbf{t})$$

The following Python code converges after 6 iterations towards $\hat{\mathbf{w}}^T = [9.78227684, -21.73839298]$ which is pretty close to the optimal solution.

```

1 # Exercise 1.1.2
2
3 def sigmoid(x):
4     """ The standard logistic function. np.exp also accepts arrays"""
5     return 1.0 / (1 + np.exp(-x))
6
7 def gradient_of_error(phi, y, t):
8     """ Gradient (first-order derivatives) of the error function, see Bishop page 207, eq.
9     4.96 """
10    return np.dot(phi.T, y - t)
11
12 def hessian_of_error(phi, y):
13     """ Hessian (second-order derivatives) of the error function, see Bishop page 207, eq.
14     4.97 """
15    R = np.diag(np.ravel(y * (1 - y)))
16    return np.dot(phi.T, np.dot(R, phi))
17
18 def cross_entropy_error(y, t):
19     # Implementation of cross entropy error = E(w)
20     result = 0
21     for n in range(0, t.shape[0]):
22         result += t[n] * np.log(y[n]) + (1 - t[n]) * np.log(1 - y[n])
23
24 w = np.array([[1.0], [1.0]]) # two-dimensional weight vector
25 x = np.array([0.3, 0.44, 0.46, 0.6])
26 t = np.array([[1], [0], [1], [0]]) # targets
27 phi = np.array([[1, x_element] for x_element in x]) # feature vector
28
29 for i in range(10):
30     y = sigmoid(np.dot(phi, w)) # class estimates
31     current_gradient = gradient_of_error(phi, y, t)
32     current_hessian = hessian_of_error(phi, y)
```

```

33
34     w = w - np.dot(np.linalg.inv(current_hessian), current_gradient)
35
36     print("Iteration {}: \nphi={} \ny={}, \ncurrent_gradient={}, \ncurrent_hessian={}, \nw
37     ={}\\n".format(
38         i, phi, y, current_gradient, current_hessian, w))
39 # Converges after 6 iterations
40 # w=[[ 9.78227684][-21.73839298]]

```

In order to show that $\phi = 0.45$ is indeed the decision boundary, we show that the probability of a data point with $\phi = 0.45$ is the same for both classes. That means that both classes have a probability of 0.5 since $p(C_1|\phi) + p(C_2|\phi) = 1$.

$$p(C_1|\phi) = y(\phi) = \sigma(\mathbf{w}^T \phi) = \sigma\left(\begin{pmatrix} 9.78227684 \\ -21.73839298 \end{pmatrix}^T \begin{pmatrix} 1 \\ 0.45 \end{pmatrix}\right) \approx 0.50$$

We can also prove that in a slightly different way using the optimal solution:

$$\begin{aligned} p(C_1|\phi) &= \frac{1}{1 + \exp(-w^T \phi)} \\ 0.5 &= \frac{1}{1 + \exp(-[9.8, -21.7]^T [1, \phi])} \\ 0.5 &= \frac{1}{1 + \exp(-9.8 + 21.7\phi)} \\ 2 &= 1 + \exp\{-9.8 + 21.7\phi\} \\ 1 &= \exp\{-9.8 + 21.7\phi\} \\ \ln(1) &= -9.8 + 21.7\phi \\ 0 &= -9.8 + 21.7\phi \\ 21.7\phi &= 9.8 \\ \phi &= \frac{9.8}{21.7} \approx 0.45 \end{aligned}$$

Useful internet resources which helped solving this exercise:

- Second Order Optimization - The Math of Intelligence 2: https://www.youtube.com/watch?v=UIFMLK2nj_w
- Logistic Regression - The Math of Intelligence (Week 2): <https://www.youtube.com/watch?v=D8alok2P468>
- <https://thelaziestprogrammer.com/sharrington/math-of-machine-learning/solving-logreg-newtons-method>
- <http://cs229.stanford.edu/notes/cs229-notes1.pdf>
- <https://www.stat.cmu.edu/~cshalizi/350/lectures/26/lecture-26.pdf>
- <https://www.stat.cmu.edu/~cshalizi/402/lectures/14-logistic-regression/lecture-14.pdf>
- <https://www.khanacademy.org/math/multivariable-calculus/applications-of-multivariable-derivatives/quadratic-approximations/a/the-hessian>
- <https://web.stanford.edu/group/sisl/k12/optimization/#index.md>
- <http://openclassroom.stanford.edu/MainFolder/DocumentPage.php?course=MachineLearning&doc=exercises/ex4/ex4.html>

Part 2 - Two-class classification using logistic regression

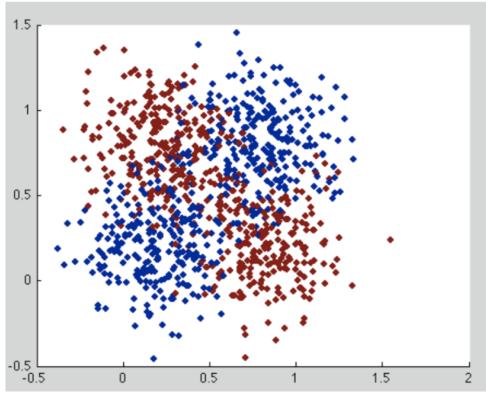


Figure 2: Two class data for logistic regression

Two-class classification using logistic regression in the IRLS algorith,. Data consists of 1000 pairs $\{x_1, x_2\}$ with corresponding class labels $C_1 = 0$ or $C_2 = 1$. Load it into Matlab using

```
data = load('a010_irlsdata.txt', '-ASCII');
X = data(:,1:2); Y = data(:,3);
```

1.2.1

Make a scatter plot of the data, similar to Figure 2. (Have a look at Matlab file `a010plotideas.m` in Brightspace for some ideas to make such a scatter plot and the plots later on.) Do you think logistic regression can be a good approach to classification for this type of data? Explain why.

Answer:

The following Python code generates figure 3:

```
1 # Exercise 1.2.1
2
3 data = np.loadtxt("data/a010_irlsdata.txt")
4
5 print("Shape of loaded data: {}".format(data.shape))
6 X = data[:,0:2] # feature vector
7 print("Shape of feature vector: {}".format(X.shape))
8
9 Y = data[:,2:3] # label vector
10 print("Shape of label vector: {}".format(Y.shape))
11
12 # Plot
13 plt.scatter(X[:,0], X[:,1], c=Y.ravel())
14 plt.xlabel('X')
15 plt.ylabel('Y')
16 plt.savefig('figure_1_2_1.png')
17 plt.show()
```

We do not think that Logistic Regression is a good approach to classify this type of data because the data does not seem to be linearly separable.

The decision boundary drawn by Logistic regression is a linear combination of features and weights denoted by $w_0 + \sum_i w_i X_i$ and therefore is only able to draw a hyperplane through the space to predict one class of each side respectively. The presented data cannot be divided by such a plane except when it would be possible to apply kernel methods like in Support Vector Machines which maps the problem into a new space.

See also: <https://homes.cs.washington.edu/~marcotcr/blog/linear-classifiers/>

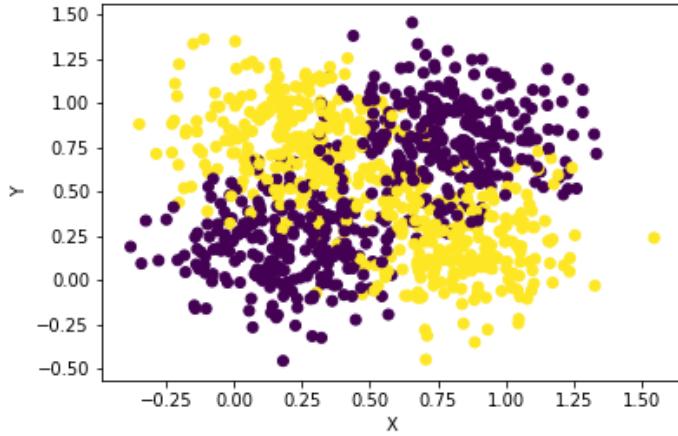


Figure 3: Scatter plot of two class data of exercise 1 - Part 2

1.2.2

Modify the Iterative Reweighted Least Squares algorithm from part 1 to calculate the optimal weights for this data. Use again a dummy basis function. Initialize with the weight vector $\mathbf{w}^T = [0, 0, 0]$. With these initial weights, what are the class probabilities according to the logistic regression model (i.e., before optimization)?

Answer:

$$\begin{aligned}
 p(C_1|\phi) &= \frac{1}{1 + \exp(-\mathbf{w}^T \phi)} \\
 &= \frac{1}{1 + \exp(-[0.0, 0.0, 0.0]^T [1, \phi, \phi])} \\
 &= \frac{1}{1 + \exp([0.0, 0.0, 0.0])} \\
 &= \frac{1}{1 + \exp(0)} \\
 &= \frac{1}{1 + 1} \\
 &= \frac{1}{2}
 \end{aligned}$$

According to the logistic regression model with the weight vector $\mathbf{w}^T = [0, 0, 0]$, both classes have a probability of 0.5.

1.2.3

Run the algorithm. Make a scatter plot of the data, similar to figure 2, but now with color that represent the data point probabilities $P(C = 1, X_n)$ according to the model after optimization. Compare the cross entropy error with the initial value. Did it improve? Much? Explain your findings.

Answer:

The following Python code generates figure 4:

```

1 # Exercise 1.2.3
2
3 w = np.array([[0.0], [0.0], [0.0]]) # three-dimensional weight vector
4 x = X # data
5 t = Y # targets/labels, has to look like np.array([ [1], [0], [1], [0] ])
6
7 print("Shape of x: {}".format(x.shape[0]))
8 print("Shape of t: {}".format(t.shape[0]))
9

```

```

10 phi = np.array([ [1, x_element[0], x_element[1]] for x_element in x ]) # feature vector
11
12 y = sigmoid(np.dot(phi, w))
13 print("Cross entropy error before optimization: {}".format(cross_entropy_error(y, t)))
14
15 for i in range(0, 1000):
16     y = sigmoid(np.dot(phi, w)) # class estimates
17     current_gradient = gradient_of_error(phi, y, t)
18     current_hessian = hessian_of_error(phi, y)
19
20     w = w - np.dot(np.linalg.inv(current_hessian), current_gradient)
21
22 print("Cross entropy error after optimization: {}".format(cross_entropy_error(y, t)))
23 print("Weights: {}".format(w))
24
25 # Normalised [0,1]
26 y = (y - np.min(y))/np.ptp(y)
27
28 cmap = plt.get_cmap('Blues')
29 plt.scatter(x[:, 0], x[:, 1], c = np.squeeze(cmap(y)))
30 plt.xlabel('X')
31 plt.ylabel('Y')
32 plt.savefig('figure_1_2_3.png')
33 plt.show()
34
35 # Cross entropy error before optimization: [693.14718056]
36 # Cross entropy error after optimization: [692.96935948]
37 # Weights: [[ 0.00440664]
38 # [-0.02139153]
39 # [-0.04930069]]

```

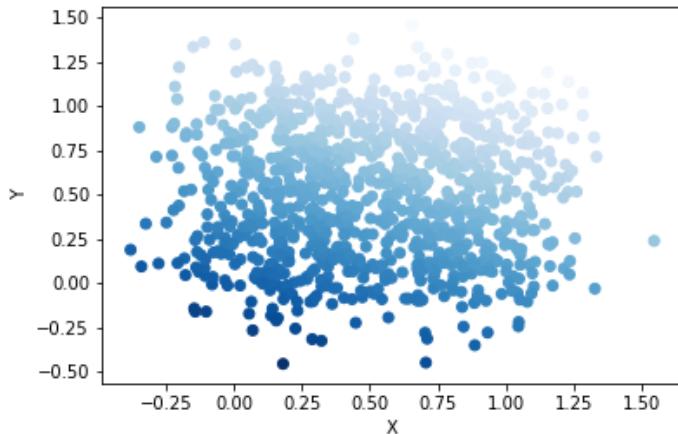


Figure 4: Scatter plot of two class data with class probabilities

The cross entropy error only improved slightly from 693.14718056 to 692.96935948.

We think the problem lies within the nature of the data that it is not linearly separable which we can also see in the transparency/probabilities in figure 4. The decision boundary seems to be drawn through the middle (linearly) although the data is not separable that way. Based on the fact that only linear basis functions are used for this non linearly separable data explains that outcome. A non-linear basis function would probably yield better results.

1.2.4

Introduce two Gaussian basis functions as features ϕ_1, ϕ_2 , similar to Bishop, fig.4.12. Use identical isotropic covariance matrices $\Sigma = \sigma^2 I$ with $\sigma^2 = 0.2$, and center the basis functions around $\mu_1 = (0, 0)$ and $\mu_2 = (1, 1)$. Make a scatter plot of the data in the feature domain. Do you think logistic regression can be a good approach to classification with these features? Explain why.

Answer:

The following Python code generates figure 5:

```

1 # Exercise 1.2.4
2
3 # Use identical, isotropic covariance matrices Sigma = sigma^2 * I with sigma^2 = 0.2
4 # and center the basis functions around mu_1 = (0,0) and mu_2 = (1,1)
5
6 def gaussian_basis_function(data, origin, variance):
7     sigma = variance * np.identity(np.array(data).shape[0])
8     a = (data - origin).T
9     normalizer = 1.0 / (2 * np.pi)**(a.shape[0]/2) * 1.0 / np.linalg.det(sigma)**(1/2)
10    return normalizer * np.exp(-(1.0/2)*np.dot(a.T, np.dot(inv(sigma), a)))
11
12 phi = np.array([[1, gaussian_basis_function(xx, (0, 0), 0.2), gaussian_basis_function(xx,
13                                         (1, 1), 0.2)] for xx in x])
14 plt.scatter(phi[:, 1], phi[:, 2], c = t.squeeze())
15 plt.xlabel('X')
16 plt.ylabel('Y')
17 plt.savefig('figure_1_2_4.png')
18 plt.show()

```

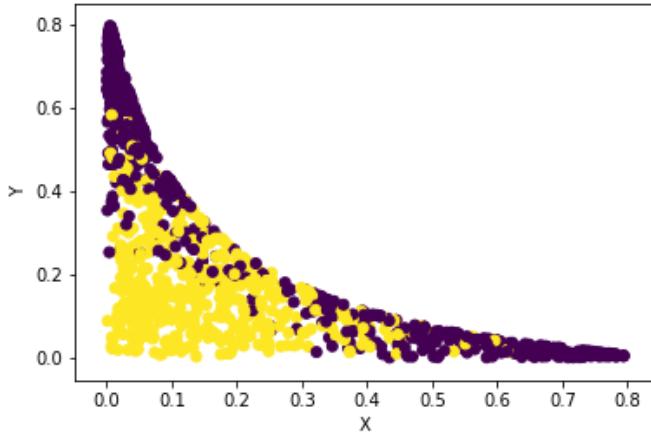


Figure 5: Scatter plot of two class data in new feature space using gaussian basis functions

The new feature space is a lot more suitable to a linear classifier because a linear decision boundary can be drawn which separates most of the two classes from each other, although not every point seems to be clearly separable. Nevertheless, this makes logistic regression an appropriate choice.

1.2.5

Modify the IRLS algorithm to use the features $\{\phi_1, \phi_2\}$ and the dummy basis function. Initialize with the weight vector $\mathbf{w}^T = [0, 0, 0]$.

Run the algorithm. Make a scatter plot of the data, similar to Figure 2, but now with colors that represent the data point probabilities $P(C = 1|X_n)$ according to this second model (after optimization). Compare the cross entropy error with the initial value. Did it improve? Much? Explain your findings.

Answer:

The following Python code generates figure 6:

```

1 # Exercise 1.2.5
2
3 y = sigmoid(np.dot(phi, w))
4 print("Cross entropy error before optimization: {}".format(cross_entropy_error(y, t)))
5
6 for i in range(0, 1000):
7     y = sigmoid(np.dot(phi, w)) # class estimates
8     current_gradient = gradient_of_error(phi, y, t)
9     current_hessian = hessian_of_error(phi, y)
10
11    w = w - np.dot(np.linalg.inv(current_hessian), current_gradient)

```

```

12 print("Cross entropy error after optimization: {}".format(cross_entropy_error(y, t)))
13 print("Weights: {}".format(w))
14
15 # Normalised [0,1]
16 y = (y - np.min(y))/np.ptp(y)
17
18 cmap = plt.get_cmap('Blues')
19 plt.scatter(x[:, 0], x[:, 1], c = np.squeeze(cmap(y)))
20 plt.xlabel('X')
21 plt.ylabel('Y')
22 plt.savefig('figure_1_2_5.png')
23 plt.show()
24
25 # Cross entropy error before optimization: [690.64901362]
26 # Cross entropy error after optimization: [346.50408046]
27 # Weights: [[ 7.10834887]
28 # [-15.42143347]
29 # [-15.53830921]
30 # [-15.53830921]]

```

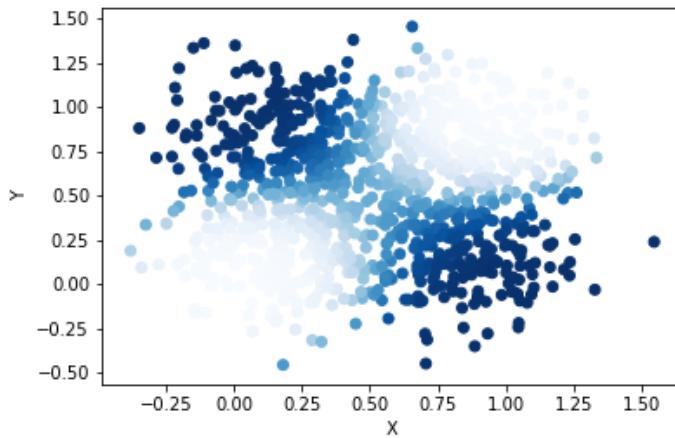


Figure 6: Scatter plot of two class data and their probabilities in new feature space using gaussian basis functions

The cross entropy error improved significantly to 346.50408046. Figure 6 shows that the gaussian basis functions are able to separate a big amount of the data at its origin of (0,0) and (1,1) and only data points in the middle have uncertain probabilities around 0.5. The cross entropy error would probably rise again if the origin of the gaussian basis functions shifted towards other origins.

Exercise 2 - Neural network regression (weight 5)

We train a neural network using backpropagation, to learn to mimic a 2D multimodal probability density. First, we implement the network and test its regression capabilities on a standard Gaussian; then we train it on the real data set. Visualization of the network output plays an important role in monitoring the progress.

2.1

Create a plot of an isotropic 2D Gaussian $y = 3 \cdot \mathcal{N}(\mathbf{0} | \frac{2}{5} \mathbf{I}_2)$ centered at the origin using the `meshgrid()`, `mvnpdf()` and `surf()` functions. Sample the density at 0.1 intervals over the range $[-2, 2] \times [-2, 2]$ and store the data in column vector variables \mathbf{X} (2D) and \mathbf{Y} (1D).

Answer:

This can be done by using the following code:

```

1 def Gaus(X):
2     FirstPart = 1/np.power(2*np.pi, len(X)/2)*1/(np.power(np.linalg.det(Sigma), 0.5))
3     Exp = -(1/2)*np.matmul(np.transpose(X-Mu), np.matmul(np.linalg.inv(Sigma), X-Mu))
4     return FirstPart*np.exp(Exp)

```

```

5 def Plot(X,Y,Name):
6     x,y = np.hsplit(X,2)
7     x = np.reshape(x,(41,41))
8     y = np.reshape(y,(41,41))
9     z = np.reshape(Y,(41,41))
10    fig = plt.figure(Name)
11    ax = fig.gca(projection='3d')
12    surf = ax.plot_surface(x, y, z, cmap=cm.coolwarm, linewidth=1, antialiased=True)
13    fig.colorbar(surf, shrink=0.5, aspect=5, pad=-0.07)
14    ax.view_init(azim=-120,elev = 70)
15    ax.set_zlabel('\nY', fontsize=50)
16    ax.set_ylabel('\n$X_1$', fontsize=50)
17    ax.set_xlabel('\n\n$X_0$', fontsize=50)
18
19
20 Sigma = (2/5)*np.identity(2)
21 Mu = np.zeros(2)
22 x = y = np.arange(-2,2+0.1,0.1)
23 x,y = np.meshgrid(x,y)
24 X = []
25 for i in range(len(x)):
26     for j in range(len(x[0])):
27         X.append(np.array([x[i][j],y[i][j]]))
28 X = np.array(X)
29 Y = np.array([Gaus(x) for x in X])

```

The function "Gaus"(line 1) simply follows the definition of the multi-variable Gaussian as given at page 78 of Bishop:

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{(2\pi)^{\frac{D}{2}}} \frac{1}{|\Sigma|^{\frac{1}{2}}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)} \quad (5)$$

For plotting, we need to reshape it because this is the way that matplotlib plots, it takes 2 2D arrays for X and Y and 1 2D array for z. The other lines create the colorbar and axis labels. The construction with the for-loop ensures that X contains all points required uniquely. The plot of the Gaussian can be seen at figure 7:

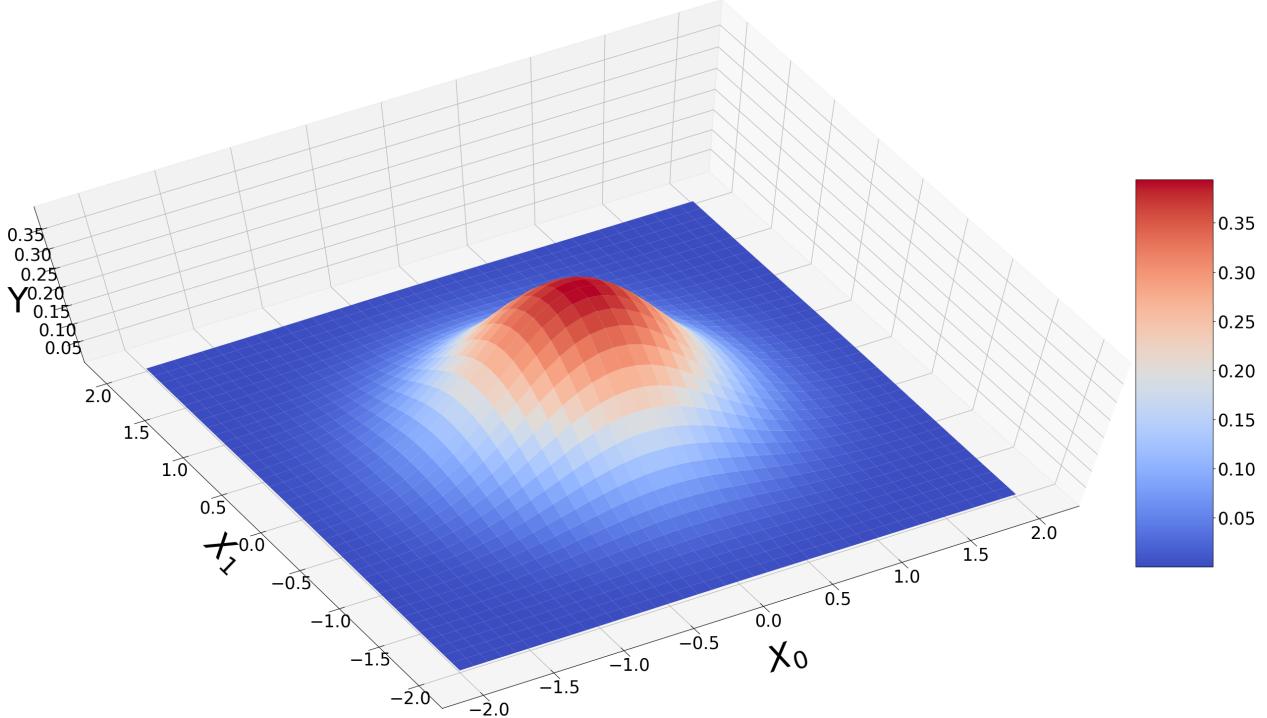


Figure 7: Gaussian function $Y = 3 \cdot \mathcal{N}(\mathbf{0} | \frac{2}{5} \mathbf{I}_2)$ graphed.

2.2

Implement a 2-layer neural network with $D = 2$ input nodes, $K = 1$ output nodes and M hidden nodes in the intermediate layer that can be trained using a sequential error backpropagation procedure, as described

in Bishop §5.3. Use $tanh(\cdot)$ activation functions for the hidden nodes and a linear activation function (regression) for the output node. Introduce appropriate weights and biases, and set the learning rate parameter $\eta = 0.1$. Initialize the weights to random values in the interval $[-0.5, 0.5]$. Plot a 2D graph of the initial output of the network over the same $[-2, 2] \times [-2, 2]$ grid as the Gaussian (again using `surf()`).

Answer:

For the network I created a class in Python:

```

1  class NeuralNW():
2      def __init__(self, X, Y, eta):
3          self.w1 = np.random.random_sample((D+1,M))-0.5
4          self.w2 = np.random.random_sample(M+1)-0.5
5          self.X = X
6          self.Y = Y
7          self.TrainX = []
8          for x in self.X:
9              self.TrainX.append(np.array([1, x[0], x[1]]))
10         self.TrainX = np.array(self.TrainX)
11
12     def setNetworkOnce(self, Num):
13         self.HiddenLayer = np.tanh(np.dot(self.w1[0], self.TrainX[Num][0]) + np.dot(self.w1
14 [1], self.TrainX[Num][1]) + np.dot(self.w1[2], self.TrainX[Num][2]))
15         self.HiddenLayer = np.insert(self.HiddenLayer, 0, 1)
16         self.Output = np.sum(self.HiddenLayer * self.w2)
17
18     def TrainNetwork(self, Num):
19         self.setNetworkOnce(Num)
20         Delta2 = self.Output - Y[Num]
21         Delta1 = (1 - np.square(self.HiddenLayer)) * self.w2 * Delta2
22         Der2 = Delta2 * self.HiddenLayer
23         Delta1 = Delta1[1:]
24         Der1 = np.array([Delta1 * self.TrainX[Num][0], Delta1 * self.TrainX[Num][1], Delta1 * self
25 .TrainX[Num][2]])
26         self.w2 -= eta * Der2
27         self.w1 -= eta * Der1
28
29     def TrainNetworkAll(self):
30         for i in range(len(self.X)):
31             self.TrainNetwork(i)
32
33     def PlotOutput(self, Name):
34         Out = []
35         for i in range(len(X)):
36             self.setNetworkOnce(i)
37             Out.append(self.Output)
38         Out = np.array(Out)
39         Plot(self.X, Out, Name)

```

Initialization

When it is initialized it creates two sets of weights. The first (`w1`) are the weights from the inputs to the hidden layer. There are D inputs + 1 bias, so there are $(D+1)$ sets of weights leading to the M hidden layers. This I represent in a two dimensional $(D+1, M)$ matrix, where `w1[0][5]` stands for the weight between x_0 (bias) and z_5 , of the hidden layer. The hidden layer is connected to the one output by $M+1$ weights, since we have another bias z_0 . The `random.sample` function gives a random number between $[0,1]$ so I lower it by 0.5 to make it between $[-0.5, 0.5]$. Then it stores the `X` and `Y` values and for convenience creates an array "self.X" which already has x_0 as bias, which is always 1.

Set network

To run something through the model you can set all the layers of the network, including the output, by calling the "setNetworkOnce" function. This runs the nth (given by "Num") entry of the `X` array. The bias is set at 1 using "`np.insert(self.HiddenLayer, 0, 1)`".

TrainNetwork & TrainNetworkAll

This will run the nth point (given by `Num`) of the vector `X` and set the weights according to the 'simple example' given at page 245 of Bishop. `TrainNetworkAll` will call the function `TrainNetwork` for all points in `X`.

PlotOutput

This will put every point in \mathbf{X} through the network and plot it's output.

The initial output is for $M = 8$ can be seen in the following figure 8:

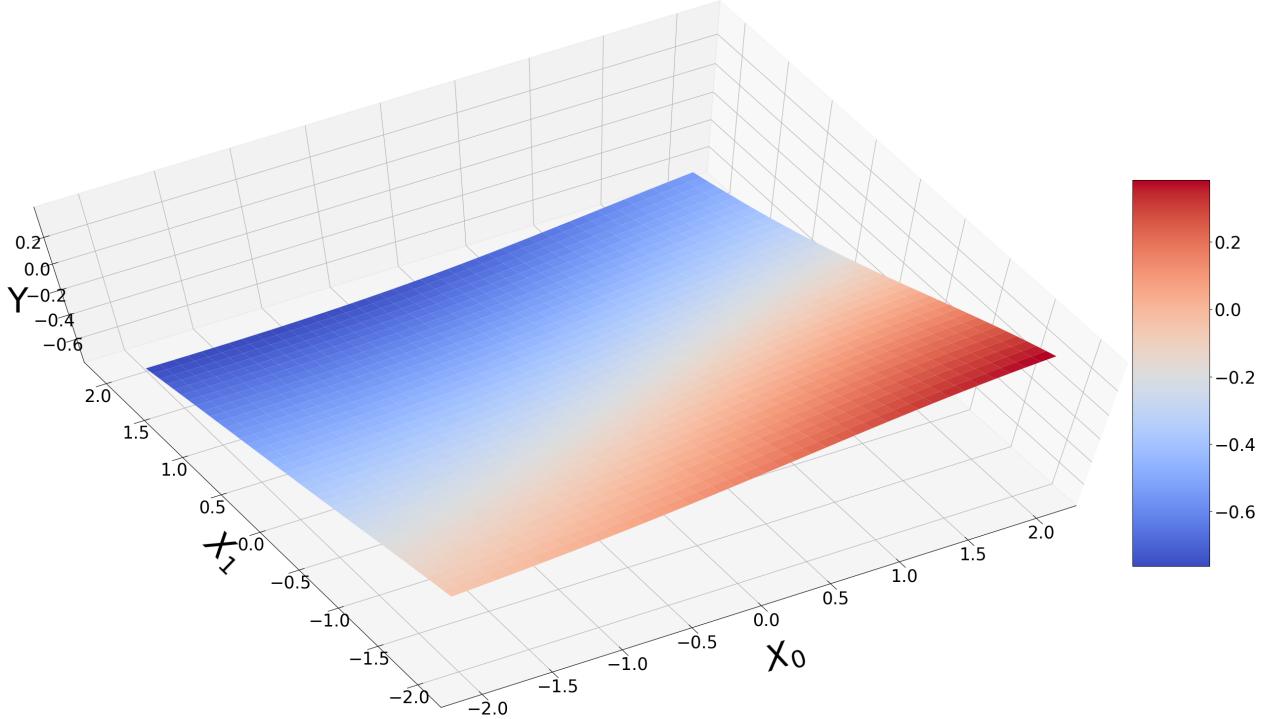


Figure 8: Initial output of the network for $m = 8$.

2.3

Train the network for $M = 8$ hidden nodes on the previously stored \mathbf{X} and \mathbf{Y} values (the $\{x_1, x_2\}$ input coordinates and corresponding output probability density y), by repeatedly looping over all datapoints and updating the weights in the network after each point. Repeat for at least 500 complete training cycles and monitor the progress of the training by plotting the output of the network over the \mathbf{X} grid after each full cycle. Verify the output starts to resemble the Gaussian density after some 200 cycles (all be it with lots of 'wobbles').

Answer:

The code I used for this is:

```
1 eta = 0.1
2 D = 2
3 M = 8
4 Network = NeuralNW(X,Y,eta)
5 Now = time()
6
7 Network.PlotOutput("First")
8
9 for i in range(2000):
10     if (i == 199):
11         Network.PlotOutput("Second")
12         print(time()-Now)
13     Network.TrainNetworkAll()
14
15 Network.PlotOutput("Third")
16 print(time()-Now)
```

First we create an instance of the Network class. Then we plot the initial output and run it for some iterations, this case, 200. Then we plot the final output and the time it took to run the model.

200 iterations took me around 19.7 seconds and created the follow figure 9:

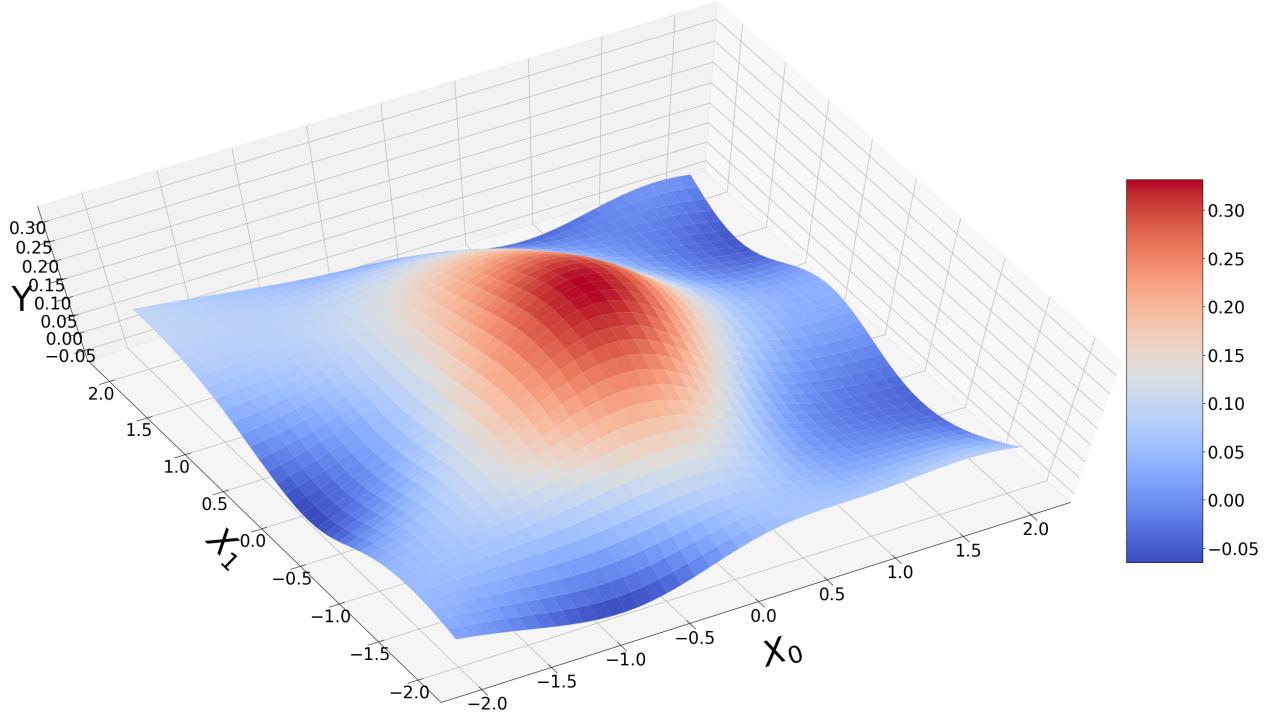


Figure 9: Output after 200 iterations of training.

Then we ran it for 2000 iterations and the results you can see at figure 10

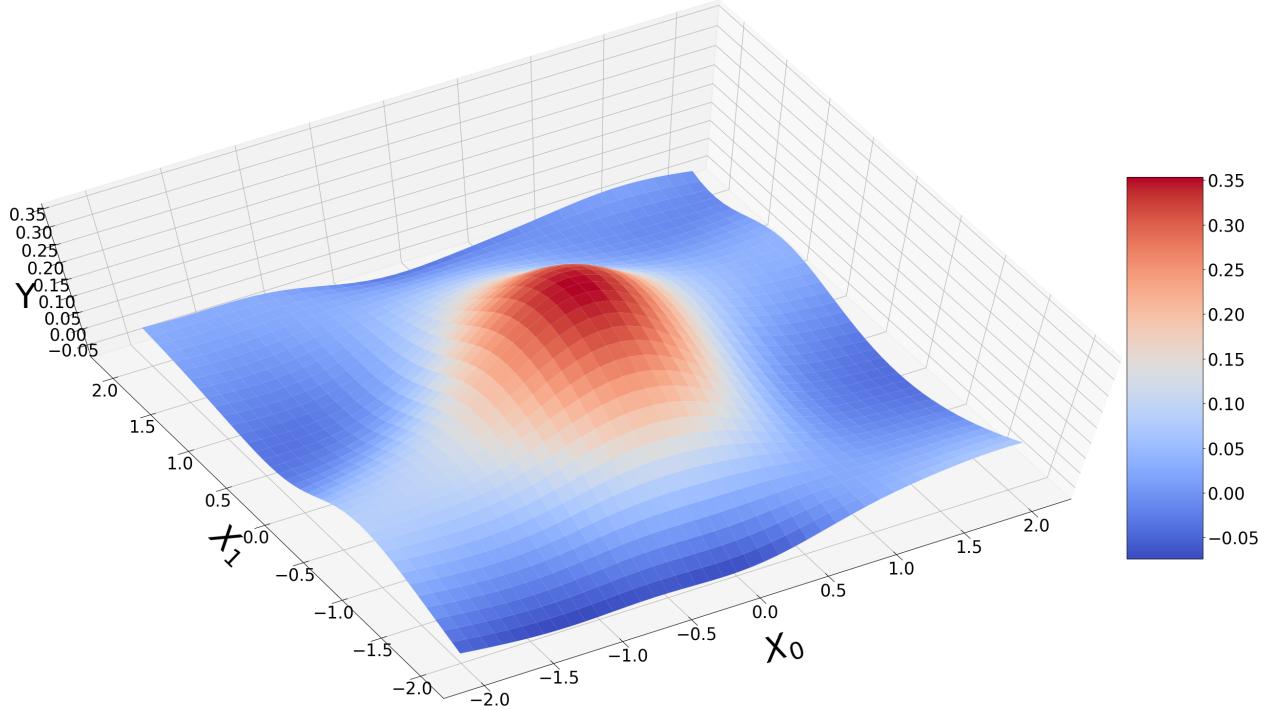


Figure 10: Output after 2000 iterations of training.

2.4

Permute the \mathbf{X} and \mathbf{Y} arrays to a random order using the `randperm()` function, keeping corresponding x and y together. Repeat the network training session using this randomized data set. Verify that convergence is now much quicker. Can you understand why? Try out the effect of different numbers of hidden nodes, different initial weights and different learning rates on speed and quality of the network training. Explain your results. **Answer:**

We permute by doing the following:

```

1 def TrainNetworkRand(self):
2     for i in np.random.permutation(len(self.X)):
3         self.TrainNetwork(i)

```

This is a function of the NeuralNetwork class. Instead of training every data point from start to finish this randomizes the order, which is equal to randomizing X and Y. (It's a little bit better since it randomizes every time it's training)

We get an indication of the convergence by using the following error function:

```

1 def Distance(self):
2     Out = []
3     for i in range(len(X)):
4         self.setNetworkOnce(i)
5         Out.append(self.Output)
6     Out = np.array(Out)
7     return np.sum(np.abs(self.Y - Out))

```

This takes the output for every point and calculates how far it is from the target, and sums it. The following code we used to show the convergence for a random permuation is indeed faster:

```

1 NetworkRand = NeuralNW(X,Y,eta)
2 for i in range(20):
3     NetworkRand.TrainNetworkRand()
4     print(NetworkRand.Distance())
5
6 NetworkNorm = NeuralNW(X,Y,eta)
7 for i in range(20):
8     NetworkNorm.TrainNetworkAll()
9     print(NetworkNorm.Distance())

```

A typical result would be something like 50 for the randomized training and 100-250 for the sequential training. By choosing random points the difference between the output and the desired output is higher, therefore more training can be done. With the desired output being close to the target not much training will be done since the gradient is not as high.

Different parameters

We've tested the convergence by choosing different parameters and looking at the Error after 100 iterations of random training. First we've tested for the amount of nodes M using the following code:

```

1 ErrorsM = []
2 Mlist = np.linspace(1,50,50,dtype=np.int)
3
4 for M in Mlist:
5     print(M)
6     NetworkTestM = NeuralNW(X,Y,eta)
7     for i in range(100):
8         NetworkTestM.TrainNetworkRand()
9     ErrorsM.append(NetworkTestM.Distance())
10
11 plt.figure("Errors M")
12 plt.plot(Mlist,ErrorsM)
13 plt.xlabel('M')
14 plt.ylabel('Error')
15
16 plt.figure("Time M")
17 plt.plot(Mlist,TimeList)
18 plt.xlabel('M')
19 plt.ylabel('t(s)')

```

This creates an list of 1,2,...,50, to test Neural networks with these amount of hidden layers. Then it trains the network and append the error. It also keeps track of how long the training took. Finally it plots the error and the execution times. The result can be viewed in figure 11.

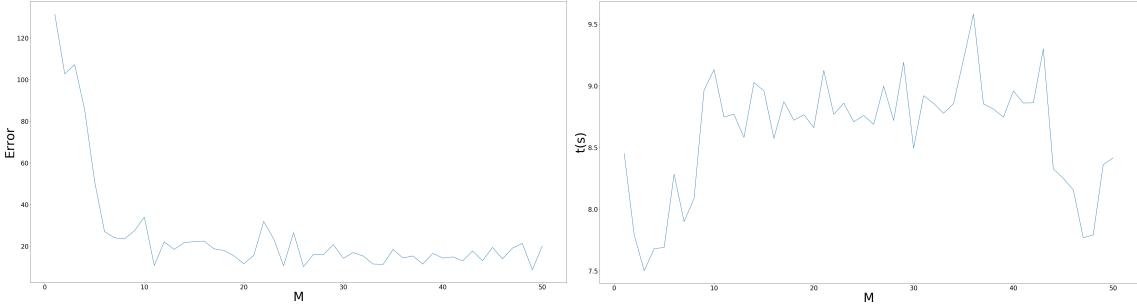


Figure 11: The error and time for various amount of hidden layers M . The time taken is fairly constant and the error seems to be constant after 10-20 hidden layers. So it's best to take around 20 hidden layers.

We also tested the effect of different η , for which we took a logarithmically, meaning that the log of the different η 's reveals a linear scale.

We used the following code:

```

1 Errorseta = []
2 M = 8
3 EtaLog = np.linspace(-4,2,50)
4
5 Etas = np.power(10,EtaLog)
6
7 for eta in Etas:
8     print(eta)
9     NetworkTesteta = NeuralNW(X,Y,eta)
10    for i in range(100):
11        NetworkTesteta.TrainNetworkRand()
12    Errorseta.append(NetworkTesteta.Distance())
13
14 plt.figure("Errors eta")
15 plt.plot(Etas,Errorseta)
16 plt.xlabel('$\eta$')
17 plt.ylabel('Error')
```

So we took 50 different values for η between 10^{-4} and 10^2 . The plot we got from this calculation is figure 12:

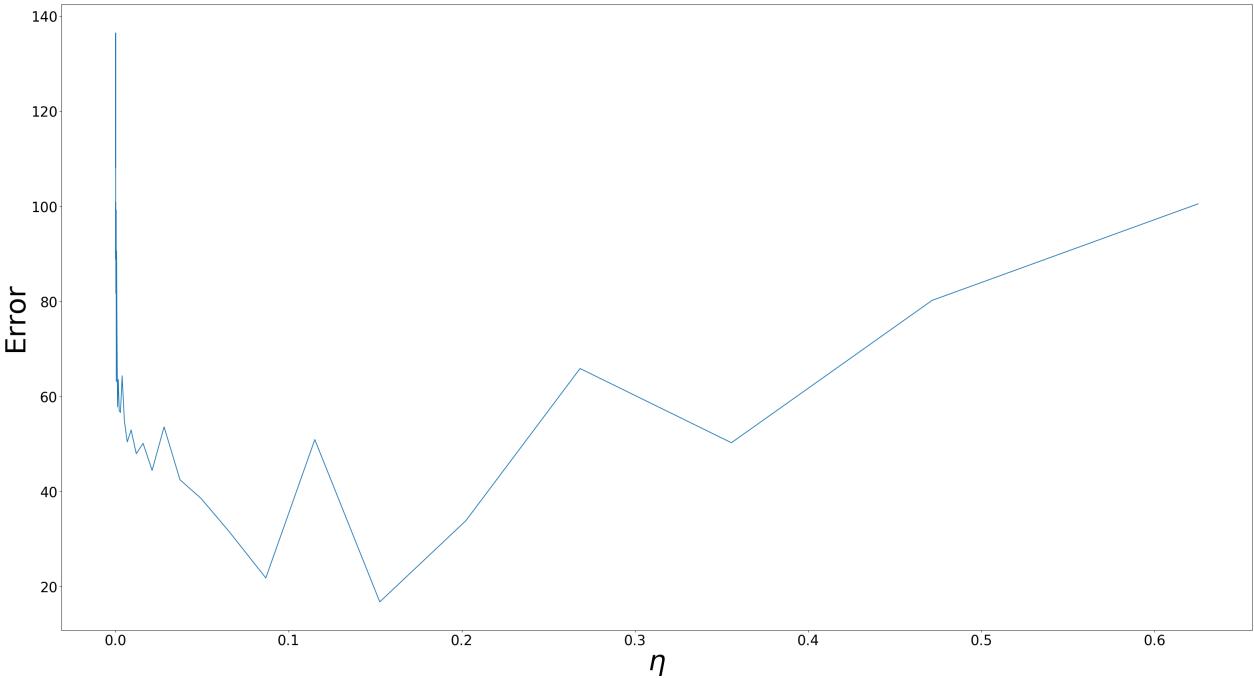


Figure 12: Errors for different values of η . A value around 0.1-1.5 seems to be the best.

Finally to test the influence of the starting values of the weights we took a random value between $[-0.5 + i, 0.5 + i]$ with i varying from -5 to 5. So we can test if a lower random weight or a higher random weight would be beneficial. We did it using the following code:

```

1 ErrorW = []
2 Weights = []
3 eta = 0.1
4 for i in range(-5,5):
5     print(i)
6     Weights.append([])
7     for j in range(10):
8         w1 = np.random.sample((D+1,M))-0.5+i
9         w2 = np.random.sample(M+1)-0.5+i
10        Weights[-1].append((w1,w2))
11
12
13 for partWeight in Weights:
14     ErrorW.append(0)
15     for w1,w2 in partWeight:
16         NetworkTestW = NeuralNW(X,Y,eta)
17         NetworkTestW.setWeights(w1,w2)
18         for i in range(100):
19             NetworkTestW.TrainNetworkRand()
20         ErrorW[-1] += NetworkTestW.Distance()
21
22 ErrorW = np.array(ErrorW)/10
23
24 WPlot = range(-5,5)
25 plt.figure("Errors in W")
26 plt.plot(WPlot,ErrorW)
27 plt.xlabel('Initial weights bias')
28 plt.ylabel('Error')
29
30 def setWeights(self,w1,w2):
31     self.w1 = w1
32     self.w2 = w2

```

We tested very different w 10 times to reduce the bias for taking random values for w. The number i we called the 'bias'. The code from 30-32 is part of the Neuralnetwork class and is used to set the weights. The plot we got from this can be seen at figure 13.

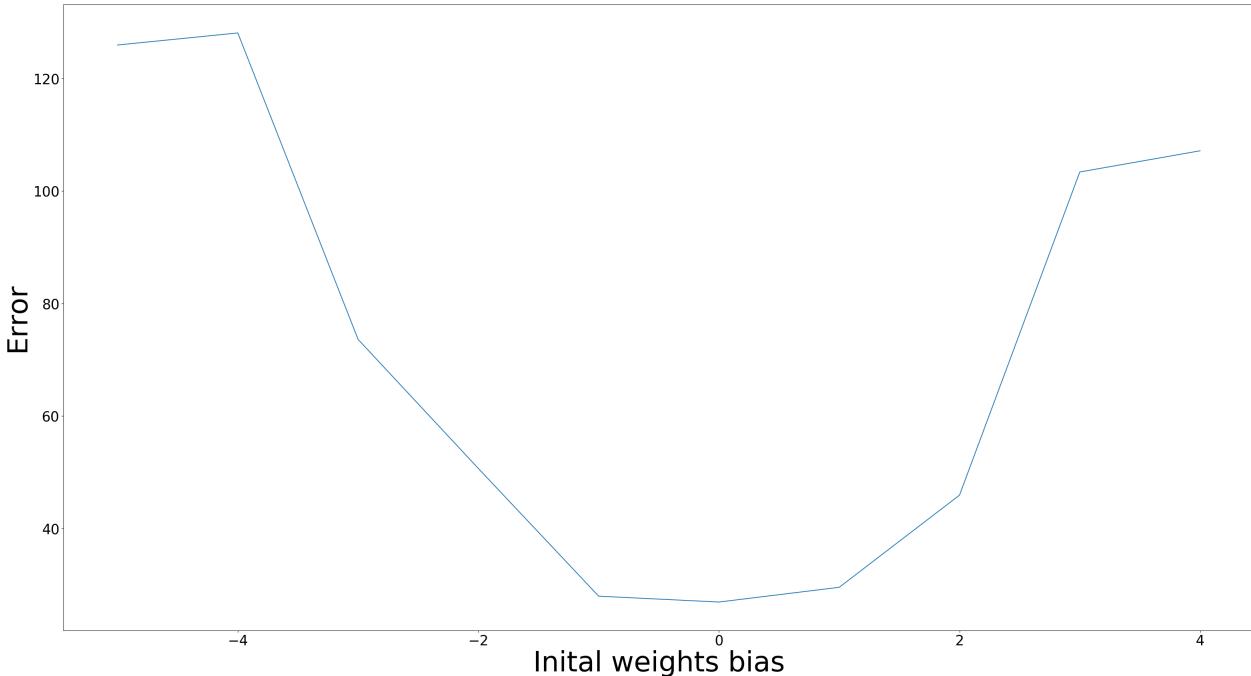


Figure 13: The bias on the random value of the initial weights. A bias of 0 seems to be the best.

To summarize: Only the amount of hidden networks could be a little bit higher in order to get a better convergence, but the values for w and η given in the assignment seemed to be quite close to the optimal, as we can conclude from our tests.

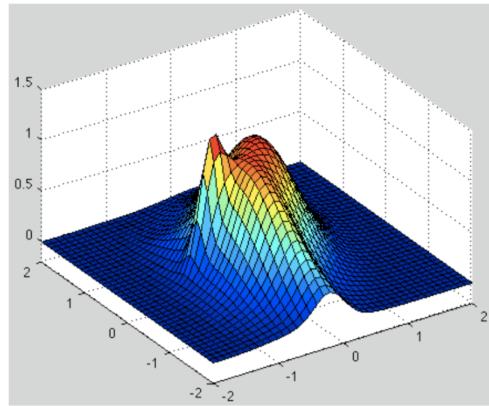


Figure 14: Multi-modal probability density

After these preliminaries we are now going to train the network on the real data set.
Load the data using

```
data = load('a017_NNpdfGaussMix.txt', '-ASCII');
X = data(:,1:2); Y = data(:,3);
```

2.5

Create a 2D-plot of the target probability density function. Notice that the data is in the correct sequence to use in `surf()`.

Answer:

This has been done with the following code:

```
1 Data = np.genfromtxt('data/a017_NNpdfGaussMix.txt')
2 X = Data[:,0:2]
3 Y = Data[:,2]
4 Plot(X,Y,"2.5")
```

Figure 15 shows the plot:

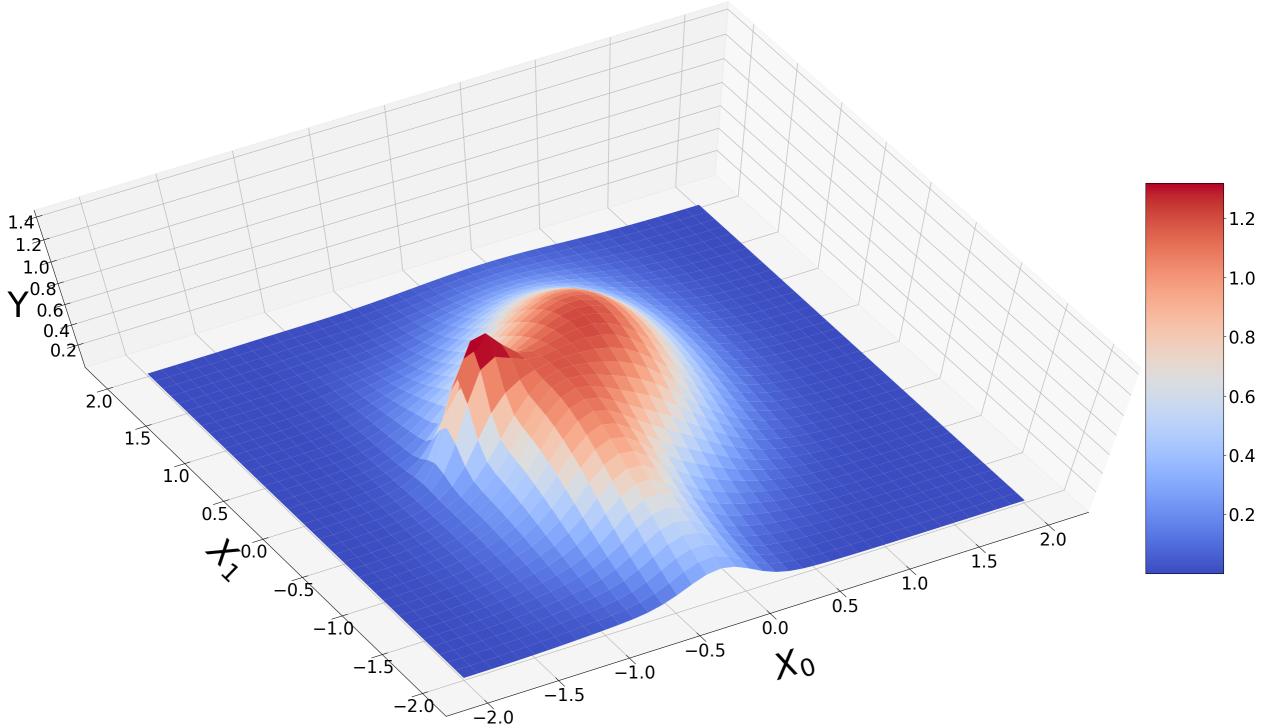


Figure 15: Plot of the data.

2.6

Train the network on this data set. Use at least 40 hidden nodes and a learning rate parameter no higher than $\eta = 0.01$. Make sure the input data is properly randomized. Run the training phase for at least 2000 complete cycles and follow the progress by plotting the updated network output after every 20 full cycles. How does the final output of the network compare to the target distribution in the data? Explain. How could you improve the neural network in terms of speed of convergence and/or quality of the approximation?

Exercise 3 - Gaussian processes (weight 5)

Part 1 - Sampling from Gaussian stochastic processes

One widely used kernel function for Gaussian process regression is given by the exponential of quadratic form, with the addition of constant and linear terms (eq. 6.63 Bishop):

$$k(\mathbf{x}, \mathbf{x}') = \theta_0 \exp\left(-\frac{\theta_1}{2} \|\mathbf{x} - \mathbf{x}'\|^2\right) + \theta_2 + \theta_3 \mathbf{x}^T \mathbf{x}' \quad (6)$$

We denote by $\boldsymbol{\theta} = (\theta_0, \theta_1, \theta_2, \theta_3)$ the hyperparameter vector governing the kernel function k .

3.1.1

Implement the kernel given by Equation (6) in Matlab as a function of \mathbf{x}, \mathbf{x}' and $\boldsymbol{\theta}$. Note that \mathbf{x} can have any dimension.

Answer:

3.1.2

We first consider the univariate case. For the parameter values $\boldsymbol{\theta} = (1, 1, 1, 1)$ and $N = 101$ equally spaced points \mathbf{X} in the interval $[-1, 1]$, compute the Gram matrix $\mathbf{K}(\mathbf{X}, \mathbf{X})$ (eq. 6.54 Bishop). What is the dimension of \mathbf{K} ? How can we show that \mathbf{K} is positive semidefinite?

Note: Even when \mathbf{K} is positive definitive, some of its eigenvalues may be too small to accurately compute (same for the determinant). This may pose a problem when generating a multivariate Gaussian distribution using \mathbf{K} as its covariance matrix. You can alleviate this issue by adding a small diagonal term to \mathbf{K} .

Answer:

3.1.3

We will now use the previously computed matrix $\mathbf{K}(\mathbf{X}, \mathbf{X})$ to produce samples from the Gaussian process prior $\mathbf{y}(\mathbf{X}) \sim \mathcal{N}(\mathbf{0}, \mathbf{K}(\mathbf{X}, \mathbf{X}))$, with \mathbf{X} being the previously determined N equally spaced points. Generate five functions $\mathbf{y}(\mathbf{X})$ with Matlab and plot them against the N input values \mathbf{X} . Repeat the process (remember to compute a new \mathbf{K} each time) for the hyperparameter configurations from Bishop, Figure 6.5:

$$\boldsymbol{\theta} \in \{(1, 4, 0, 0), (9, 4, 0, 0), (1, 64, 0, 0), (1, 0.25, 0, 0), (1, 4, 10, 0), (1, 4, 0, 5)\}.$$

Describe the differences between the plots. Explain in which way each of the kernel parameters affects the generated samples.

Answer:

3.1.4

We now move to the bivariate case. Instead of an interval, we now consider a 2-D grid of equally spaced points of size $N = 21 \times 21$ in $[-1, 1] \times [-1, 1]$. We collect all these grid points in a data matrix \mathbf{X} , where each one of the 441 observations has two dimensions. What is the dimension of \mathbf{K} now? What does this tell you about the scalability of sampling multivariate functions from Gaussian processes in higher dimensions?

Answer:

3.1.5

Using the same kernel from (6), compute the Gram matrix $\mathbf{K}(\mathbf{X}, \mathbf{X})$ on the grid for each hyperparameter configuration $\boldsymbol{\theta} \in \{(1, 1, 1, 1), (1, 10, 1, 1), (1, 1, 1, 10)\}$. For each \mathbf{K} , generate and plot four random surfaces from the Gaussian process prior $\mathbf{y}(\mathbf{X}) \sim \mathcal{N}(\mathbf{0}, \mathbf{K}(\mathbf{X}, \mathbf{X}))$. Compare the observed differences to the univariate case.

Answer:

Part 2 - Gaussian processes for regression

We would like to apply Gaussian process models to the problem of regression (Bishop 6.4.2). We consider a noisy model of the form

$$t_n = y_n + \epsilon_n,$$

where $y_n = y(x_n)$ and ϵ_n are i.i.d. samples from a random noise variable on the observed target values. Furthermore, we assume that the noise process has a Gaussian distribution given by:

$$p(t_n | y_n) = \mathcal{N}(t_n | y_n, \beta^{-1}) \quad (7)$$

Going back to a one-dimensional input space, we consider the following training data consisting of four data points

$$\mathcal{D}\{(x_1 = -0.5, t_1 = 0.5), (x_2 = 0.2, t_2 = -1), (x_3 = 0.3, t_3 = 3), (x_4 = -0.1, t_4 = -2.5)\}$$

3.2.1

Just as before, compute the Gram matrix of the training data for $\theta = (1, 1, 1, 1)$. Then, taking $\beta = 1$ in Equation 7, compute the covariance matrix \mathbf{C} corresponding to the marginal distribution of the training target values: $p(\mathbf{t}) = \mathcal{N}(\mathbf{t}|\mathbf{0}, \mathbf{C})$.

Answer:

3.2.2

Using the previous results, compute the mean and the covariance of the conditional distribution $p(t|\mathbf{t})$ of a new target value t corresponding to the input $x = 0$. Which equations from Bishop do you need?

Answer:

3.2.3

Does the mean of the conditional distribution $p(t|\mathbf{t})$ go to zero in the limit $x \rightarrow \pm\infty$? If so, explain why this happens. If not, how would you set the parameters θ of the kernel function to make it happen?

Answer:

Exercise 4 - EM and doping (weight 5)

In a certain hypothetical sport, banned substance 'X' has become popular as a performance enhancing drug, as its presence is hard to establish in blood samples directly. Recently, it has been discovered that users of the drug tend to show a strong positive correlation between concentrations of two other quantities, x_1 and x_2 , present in the blood. In contrast, 'clean' athletes tend to fall in one of two or three groups, that either show no or a negative correlation between x_1 and x_2 . Unfortunately, as each sample contains only a single, instantaneous, measurement for each variable, it is not possible to establish this correlation from the sample. However, in many cases it is possible to distinguish to which *class* a certain sample belongs by also looking at the values of two other measured variables, x_3 and x_4 : certain combinations of measured values are often typical for one class but highly unusual for others.

After a high profile event, a large scale test has resulted in 2000 samples. Rumours suggest the number of positives could be as high as 20%. However, the exact relationship between different classes and typical x values is still not clear. This is where the EM-algorithm comes in...

The blood sample measurements are modelled as a mixture of K Gaussians, one for each class

$$p(x|\mu, \Sigma, \pi) = \sum_{k=1}^K \pi_k \mathcal{N}(x|\mu_k, \sigma_k) \quad (8)$$

where $\mathbf{x} = [x_1, x_2, x_3, x_4]$ represents the values for the measured quantities in the blood sample, $\mu = \{\mu_1, \dots, \mu_K\}$ and $\Sigma = \{\Sigma_1, \dots, \Sigma_K\}$ are the means and covariance matrices of the Gaussians for each class, and $\pi = \{\pi_1, \dots, \pi_K\}$ are the mixing coefficients in the overall data set.

Load the data using

```
data = load('a011_mixdata.txt', '-ASCII');
```

and set N to the number of datapoints and D to the number of variables in the dataset X.

4.1

Try to give an estimate of the number, size and shape of the classes in the data, by plotting the distribution of the variables, e.g. using `hist()`, `scatter()` or `scatter3()`.

Answer:

4.2

Implement an EM-algorithm using the description and formulas given in Bishop, §9.2.2. Use variable K for the number of classes and choose a priori equal mixing coefficients π_k . Initialize the means \mathbf{m}_k , to random values around the sample mean of each variable, e.g. set $\mu_{k,1}$ to $\bar{x}_1 + [-1 \leq \epsilon \leq +1]$. Initialize the Σ_k to diagonal matrices with reasonably high variances, e.g. $4 * \text{rand}() + 2$, to avoid very small responsibilities in the first step. Make sure the EM-loop runs over at least 100 iterations. Display relevant quantities, at least the log likelihood (9.28), after each step so you can monitor progress and convergence. Write a plot routine that plots the x_1, x_2 coordinates of the data points, and color each data point according to the most probable component according to the mixture model.

Answer:

4.3

Set $K = 2$, initialize your random generator and run the EM-algorithm on the data. Describe what happens. Try different random initializations and compare results.

(Should converge within 50 steps to two clusters, accounting for $\pm 1/3$ resp. $2/3$ of the data). Plot the x_1, x_2 coordinates colored according to the most probable component. Compute the correlation coefficients

$$p_{12} = \frac{\text{cov}[x_1, x_2]}{\sqrt{\text{var}[x_1]\text{var}[x_2]}} \quad (9)$$

of each of the components (i.e., use their covariance matrices to compute variances and covariances in (9), see also (Bishop, eq. (2.93))). Does either class show the characteristic strong¹ positive correlation for $\{x_1, x_2\}$?

Answer:

4.4

Increase the number of classes to $K = 3$ and rerun your algorithm on the data, again trying different random initializations. Plot the x_1, x_2 coordinates colored according to the most probable component and compute the correlation coefficients of each of the components. Check both your plot and your coefficients if one of the clusters now displays the strong positive $\{x_1, x_2\}$ correlation we are looking for.

Increase to $K = 4$, do the same, and see if this improves your result (in terms of detection of the doping-clusters). Based on your findings, is the rumoured 1-in-5 estimate for users of X credible?

Answer:

4.5

Having found the offending cluster in the data using the EM-algorithm, we are now presented with four samples $\{A, B, C, D\}$, with values for $[x_1, x_2, x_3, x_4]$ given as

$$\begin{aligned} A &= [11.85, 2.2, 0.5, 4.0] \\ B &= [11.95, 3.1, 0.0, 1.0] \\ C &= [12.00, 2.5, 0.0, 2.0] \\ D &= [12.00, 3.0, 1.0, 6.3] \end{aligned}$$

One of these is from a subject who took drug X, and one is from a subject who tried to tamper with the test by artificially altering one or more of the x_i levels in his/her blood sample.

¹According to Wikipedia, the correlation is none if $|p| < 0.1$, small if $0.1 < |p| < 0.3$, medium if $0.3 < |p| < 0.5$ and strong if $|p| > 0.5$

Identify which sample belongs to the suspected user and which one belongs to the 'fraud'.

Answer: