

Assignment 4

Juraj Šušnjara (s4846559)
Sebastian Tiesmeyer (s4373162)

November 2016

1 Exercise 1 - Gaussian processes for regression

We consider a noisy model of the form:

$$t_n = y_n + \epsilon_n \quad (1)$$

Noise process has a Gaussian distribution:

$$p(t_n|y_n) = \mathcal{N}(t_n|y_n, \beta^{-1}) \quad (2)$$

We will use following kernel function for Gaussian process regression:

$$k(\mathbf{x}, \mathbf{x}') = \theta_0 \exp\left(-\frac{\theta_1}{2} \|\mathbf{x} - \mathbf{x}'\|^2\right) + \theta_2 + \theta_3 \mathbf{x}^T \mathbf{x}' \quad (3)$$

1.1

Implementation of equation 3 in MATLAB is:

```
1 function k = kernel( x, xt, theta )
2
3     k = theta(1).*exp((-theta(2)./2).*norm(x-xt).^2)+
4             theta(3)+theta(4).*(x'*xt);
5
6 end
```

1.2

Following script is used for computation of Gram matrix $\mathbf{K}(\mathbf{X}, \mathbf{X})$:

```

1 N = 101;
2 X = linspace(-1,1,N);
3 K = zeros(N, N);
4 theta = [1 1 1 1];
5
6 for i=1:N
7     for j=1:N
8         K(i,j) = kernel(X(i), X(j), theta);
9     end
10 end

```

1.3

The size of \mathbf{K} is $101 * 101$, $d = R^2$.

A given $[m*m]$ matrix M is semipositive definite if:

$$z^T M z \geq 0 \quad (4)$$

for every given $[m*1]$ vector z.

The Gram matrix is the cross product of all vectors of the vector set X:

$$G = X^T X \quad (5)$$

Substituting G in equation (4) results in a squared absolute of the cross product of X and z, which is always zero or positive:

$$z^T X^T X z = \|X z\|_2^2 \geq 0 \quad (6)$$

1.4

Plots of 5 random samples produced from Gaussian process prior $y(\mathbf{X}) \sim \mathcal{N}(0, \mathbf{K}(\mathbf{X}, \mathbf{X}))$ are presented on figure 1.

1.5

On figure 2 we can see different plots produced from a Gaussian process prior with different θ values. Plots are drawn by following MATLAB code snippet:

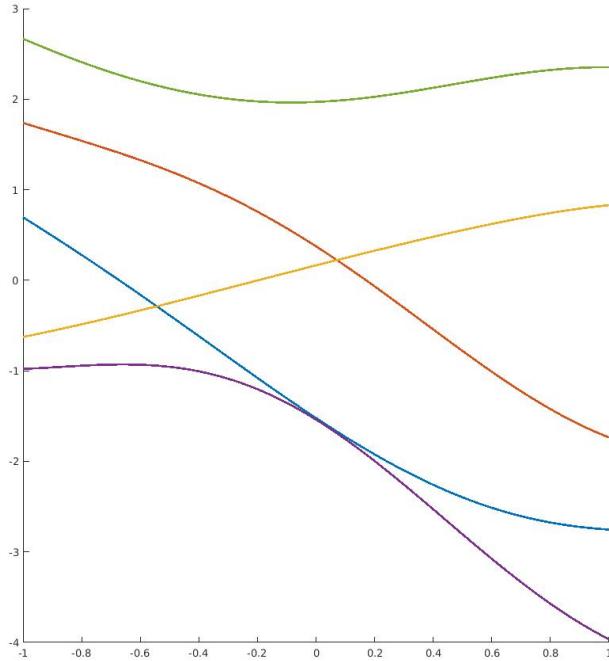


Figure 1: Plots of 5 different samples from Gaussian process prior where $\theta = (1, 1, 1, 1)$

```

1 mu = zeros(101,1);
2 Y = mvnrnd(mu, K);
3 plot(X,Y);

```

Bonus: The different θ s emphasize or suppress different parts of the kernel equation which leads to different covariance matrices. θ_2 is just a constant that raises the value of every element, so the covariance of all (x_1, x_2) -combinations increases regardless of their euclidean distance. The sample curves will show small local and global covariances, be relatively leveled, flat and smooth. θ_0 and θ_1 both influence the first term. θ_0 increases covariance proportionally, so local and global variances in/decrease the same and the range of $y(x)$ is basically scaled up. θ_1 decreases the local variance by magnifying the absolute distance indicator $\|x - x'\|$ ($- >$ when $[-0.1:0.1]$ were scaled up in Bishop 6.5 [3], we would retrieve the original function). θ_3 promotes global devariance of (x_1, x_2) -combinations with opposite signs.

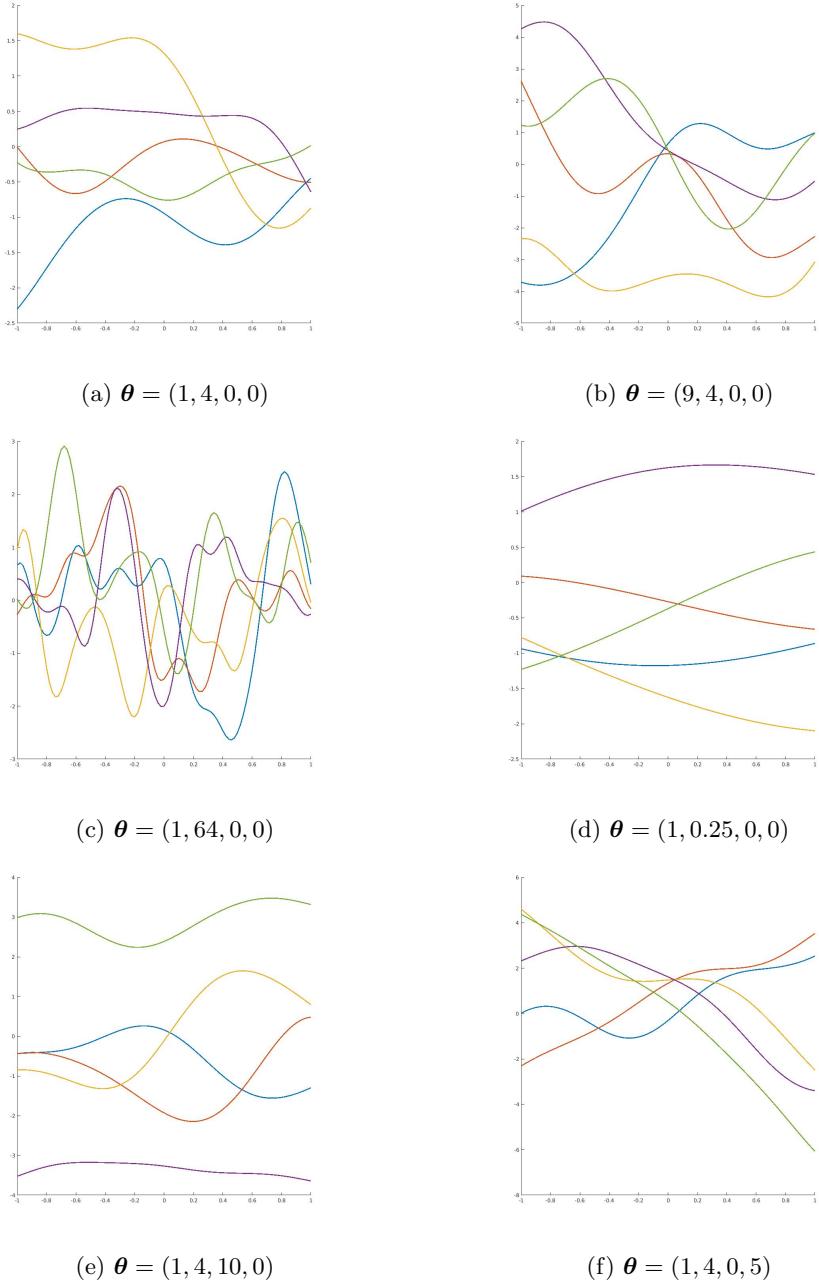


Figure 2: Plots of samples from Gaussian process prior with different θ values

1.6

Covariance matrix \mathbf{C} corresponding to the marginal distribution of the training target values $p(\mathbf{t}) = \mathcal{N}(\mathbf{t}|\mathbf{0}, \mathbf{C})$ is calculated using the expression 7 and shown in expression 8

$$C(x_n, x_m) = k(x_n, x_m) + \beta^{-1} \delta_{nm} \quad (7)$$

$$C = \begin{pmatrix} 3.25 & 1.68 & 1.58 & 1.97 \\ 1.68 & 3.04 & 2.06 & 1.94 \\ 1.58 & 2.06 & 3.09 & 1.89 \\ 1.97 & 1.94 & 1.89 & 3.01 \end{pmatrix} \quad (8)$$

1.7

To find the conditional distribution $p(t_{N+1}|\mathbf{t})$ we start with writing joint distribution $p(\mathbf{t}_{N+1})$ which is shown in expression 9.

$$p(\mathbf{t}_{N+1}) = \mathcal{N}(\mathbf{t}_{N+1}|\mathbf{0}, \mathbf{C}_{N+1}) \quad (9)$$

where

$$\begin{aligned} \mathbf{C}_{N+1} &= \begin{pmatrix} \mathbf{C}_N & \mathbf{k} \\ \mathbf{k}^T & c \end{pmatrix} \\ c &= k(x_{N+1}, x_{N+1} + \beta^{-1}) \\ \mathbf{k} &= k(x_n, x_{N+1}) \end{aligned} \quad (10)$$

Finally, mean and covariance are calculated by following expressions:

$$m(x_{N+1}) = \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{t} \quad (11)$$

$$\sigma^2(x_{N+1}) = c - \mathbf{k}^T \mathbf{C}_N^{-1} \mathbf{k} \quad (12)$$

Based on given data set, mean and covariance of a new target value t corresponding to input $x = 0$ are given in expressions 13 and 14.

$$m(0) = -0.0252 \quad (13)$$

$$\sigma(0) = 1.2357 \quad (14)$$

1.8

This depends on the mean of the input data - t needs to have mean zero as well. To achieve this (and we want to, since our assumption a priori is that the target process = mean(0)), one can adjust $f(x)$ by setting θ_2 to be -mean(t).

2 Exercise 2 - Neural network regression

2.1 Plotting a 3D Gaussian

MATLAB code to plot a 3D Gaussian with $\mu = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$ and $\Sigma = \begin{pmatrix} 0.4 & 0 \\ 0 & 0.4 \end{pmatrix}$

```
1 figure(1)
2 clf;
3
4 %Initiate X Y values and covmat:
5 [X,Y] = meshgrid(-2:.1:2);
6 Sigma = (2/5)*eye(2);
7
8 %function handle to calculate a target value:
9 tgt = @(X,Y) reshape(mvnpdf([X(:) Y(:)], [0 0], Sigma), size(X));
10
11 %Calculate Z-values & plotification:
12 Z= reshape(tgt(X(:),Y(:)),size(X))
13 surf(X,Y,Z)
```

Output:

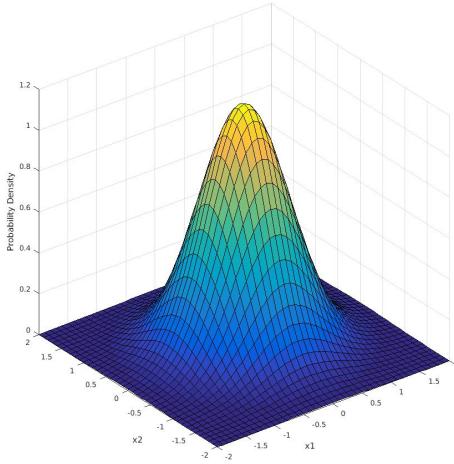


Figure 3: Plot of an isotropic 2D Gaussian $y = 3 \cdot \mathcal{N}(\mathbf{x}|0, \frac{2}{5} I_2)$

2.2 Initiating a network

```
1 %% --- Implement a 2-layer ANN: ---
2 M=8;           %hidden layer size
3 eta = .1;      %learning rate
4
5 %function handles for forward propagation:
6 fwdp1 =@(x,w) tanh(([x,1])*w);
7 fwdp2 =@(a,w) [a,1]*w';
8
9 %function handles for backprop:
10 bckp2 =@(y,t) y-t;
11 bckp1 =@(w,z,d)(1-([z 1].^2)).*(w*d);
12
13 %training rounds
14 reps = 500;
15
16 %initiate weights (one added for bias...)
17 w1 = rand(3,M)-0.5;
18 w2 = rand(1,M+1)-0.5;
```

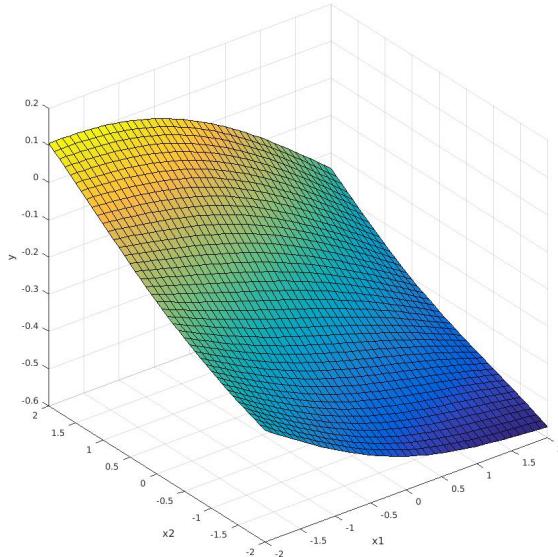


Figure 4: Output of neural network with random initial weights

Output of neural network is presented in figure 4. Neural network has 2 layers, 2 input nodes, 8 hidden nodes and 1 output node. Activation function for hidden nodes is $\tanh(\cdot)$ whereas output node uses linear activation function. All weights are random values in $[-0.5, 0.5]$ interval.

2.3 Training a network

On figure 5 we can see outputs of neural network after certain number of training cycles. On figure 4 we can see output given the random initial weights. After that all weights are updated using sequential error backpropagation procedure. Each cycle represents one iteration through whole data set. Weights are updated after each data point. On figures we can see that after 100 cycles plot starts to look like Gaussian density.

MATLAB code for iterative ANN training:

```

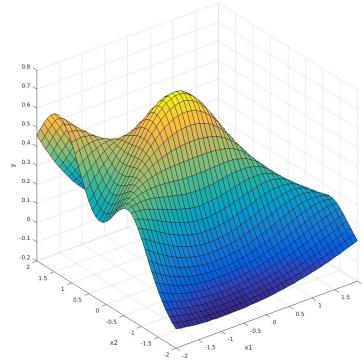
1 %training iterator:
2 for rep = 1:(reps*length(X(:)))
3     %x,t points for this iteration (basically parsing through XY-grid
4     %over and over)
5     x=[X(mod(rep-1,length(X1))+1) Y(mod(rep-1,length(X1))+1)];
6     t = tgt(x(1),x(2));
7
8     %fwdprop:
9     z = fwdp1(x,w1);
10    y = fwdp2(z,w2);
11
12    %backprop:
13    d2 = bckp2(y,t) ;
14    d1 = bckp1(w2,z,d2);
15    d2s = [d2s,d2];
16    %update weights:
17    w1 = w1-(eta*repmat(d1(1:M),[3 1]).*repmat([x 1]',[1 M]));
18    w2 = w2-(eta*repmat(d2,[1 M+1]).*repmat([z 1],[1 1]));
19
20    %Update plot every Nth iteration:
21    N=50;
22    if ~mod(rep,length(X(:))*N)
23        for i=1:length(X1)
24            Z1(i)=fwdp2(fwdp1([X1(i) Y1(i)], w1),w2);
25            %Z(i) = tgt(X1(i),Y1(i));
26        end
27
28        subplot(2,5,(rep/length(X(:)))/N)

```

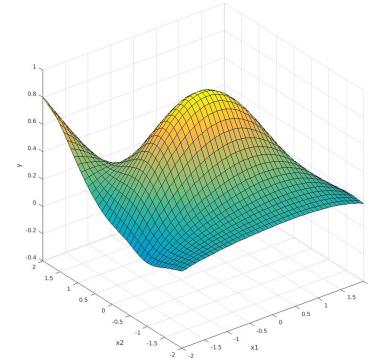
```

29      surf(X,Y,Z1);
30      title((rep/length(X(:)))); 
31      drawnow
32  end
33 end

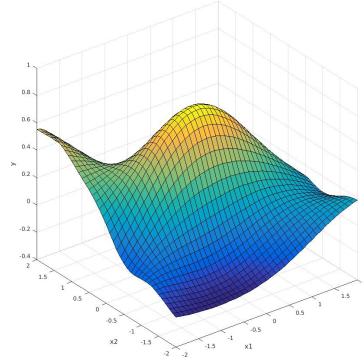
```



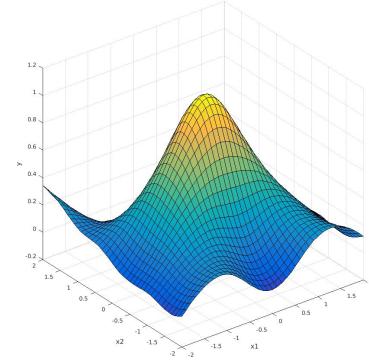
(a) Output after 20 cycles



(b) Output after 100 cycles



(c) Output after 200 cycles



(d) Output after 500 cycles

Figure 5: Output of neural network after 20, 100, 200 and 500 cycles with learning rate 0.1 and 8 units in hidden layer

2.4 A less patterned input

If we randomly permute our training input and output sets we can see that the convergence is much quicker. Even after 10 cycles plot resembles Gaussian density. That can be seen on figure 6. Randomized data set obviously leads to faster convergence. That is because if we have data set ordered in some way, weights will update slowly according to presented order. So, for example if we have ordered clusters of data where every cluster represents different class, after iterating through first cluster weights will be updated according to that cluster, than when we start iterating through second cluster weights will start to update according to that cluster etc. It is better to randomly pick data set points so our learning algorithm won't suffer from any effects of order in data.

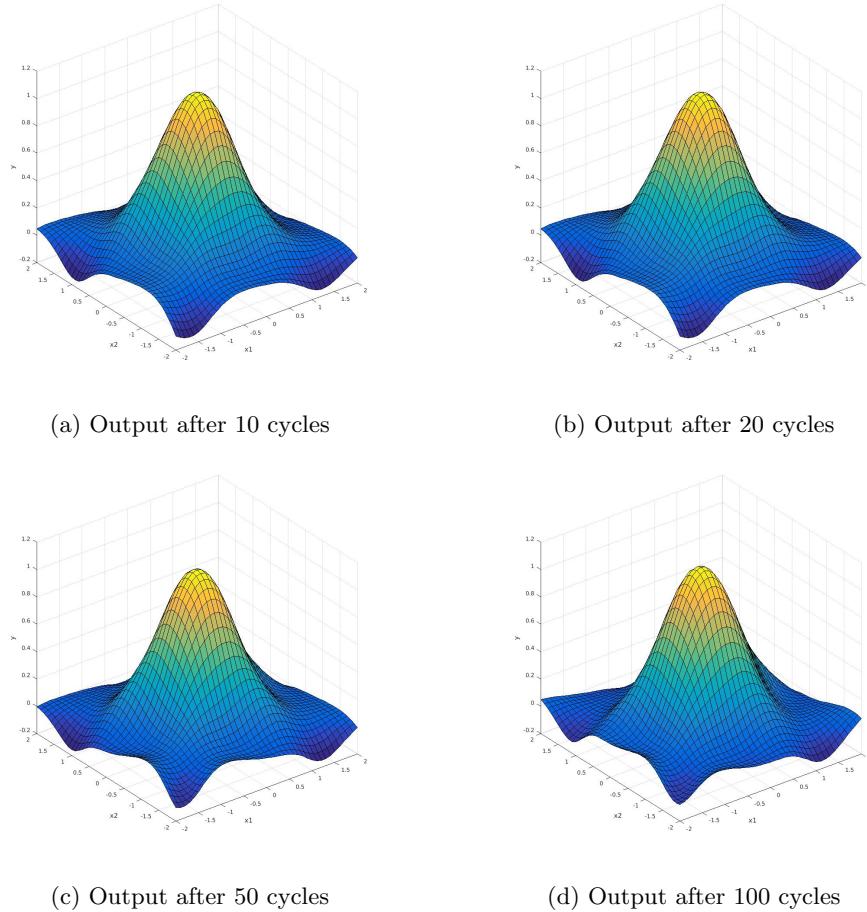


Figure 6: Output of neural network after 10, 20, 50 and 100 cycles with learning rate 0.1 and 8 units in hidden layer

If we increase learning rate to 0.4 we can notice that algorithm doesn't converge at all as we can see on figure 7. On the other hand, if we decrease learning rate to 0.001 we can notice that algorithm converge but much slower as presented in figure 8.

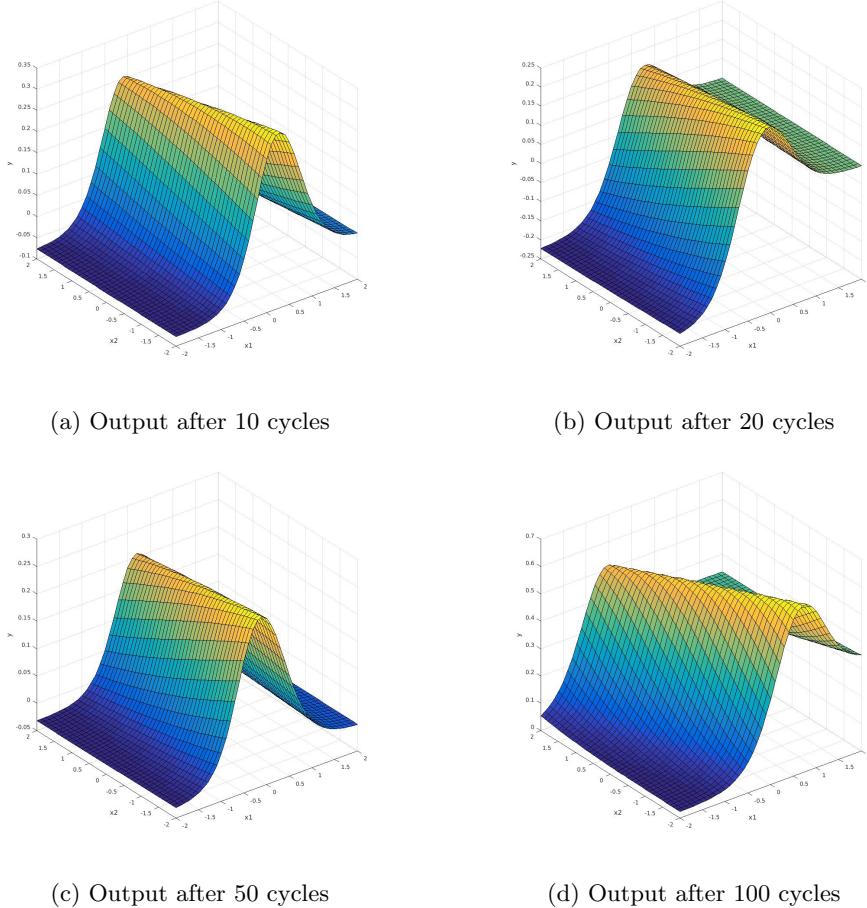
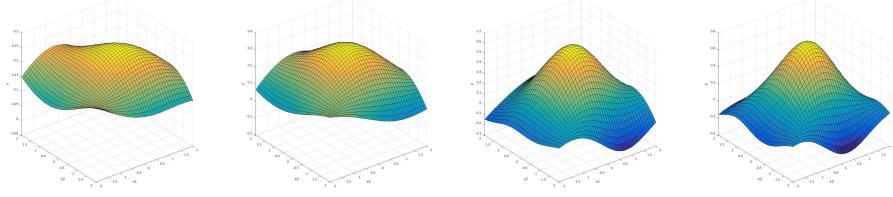


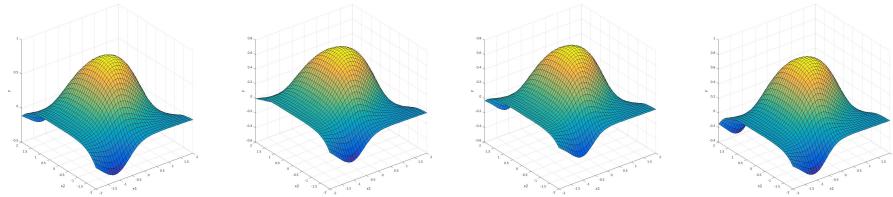
Figure 7: Output of neural network after 10, 20, 50 and 100 cycles with learning rate 0.4 and 8 units in hidden layer

We have also tried changing number of units in hidden layer in order to see its effect on learning. From figures 8 and 9 we can see plots for networks with 3 and 25 hidden units, respectively. It is obvious that network with more hidden units performs better in this case. But adding many units in hidden layer doesn't help that much, results are similar with 8, 15, 25 or more units. But reducing number of hidden units reduces network's learning capabilities and its expressiveness:



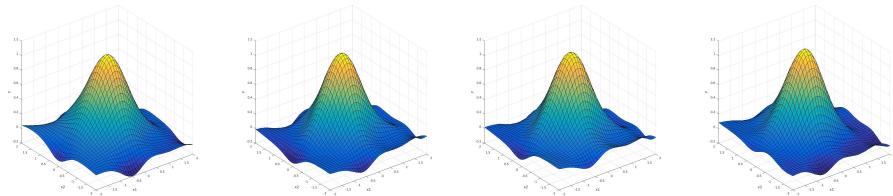
(a) Output after 10 (b) Output after 20 (c) Output after 50 (d) Output after 100
cycles cycles cycles cycles

Figure 8: Output of neural network after 10, 20, 50 and 100 cycles with learning rate 0.001 and 8 units in hidden layer



(e) Output after 10 (f) Output after 20 (g) Output after 50 (h) Output after 100
cycles cycles cycles cycles

Figure 8: Output of neural network after 10, 20, 50 and 100 cycles with learning rate 0.1 and 3 units in hidden layer



(a) Output after 10 (b) Output after 20 (c) Output after 50 (d) Output after 100
cycles cycles cycles cycles

Figure 9: Output of neural network after 10, 20, 50 and 100 cycles with learning rate 0.1 and 25 units in hidden layer

2.5

Required plot is shown in figure 10.

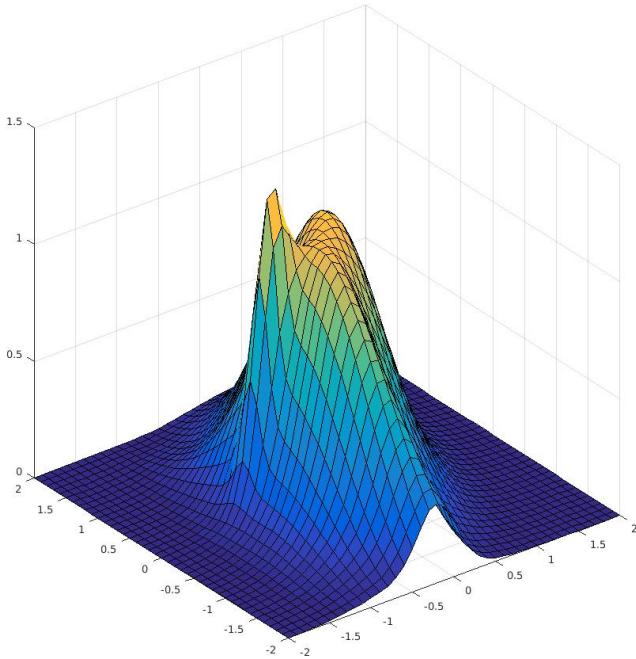


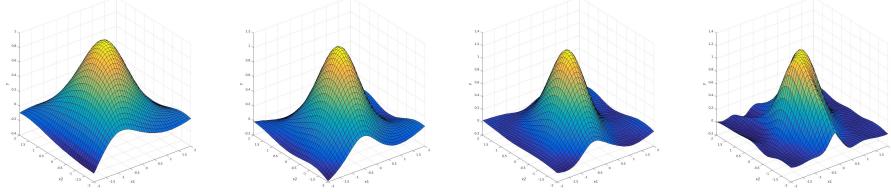
Figure 10: 2D plot of the target probability density function

2.6

We used network with 50 units in hidden layer and we set learning rate to 0.005 to train our neural network for data set given in previous task. On figure 11 we can see plots of our results after certain number of training cycles. If we compare our results with plot presented in 10 (which is actually a plot from given data set) we can notice one big difference: our plots are missing that "peak anomaly" that can be seen on figure 10.

2.7

Using *Netlab* toolbox we created network with 50 units in hidden layer, we trained our data set using scaled conjugate gradients. Final output is presented in figure 12. As we can see, our results (figure 11) are similar to this result, but this one looks slightly better at approximating function from given data set.



(a) Output after 20 cycles (b) Output after 40 cycles (c) Output after 100 cycles (d) Output after 2000 cycles

Figure 11: Output of neural network after 20, 40, 100 and 2000 cycles with learning rate 0.005 and 50 units in hidden layer

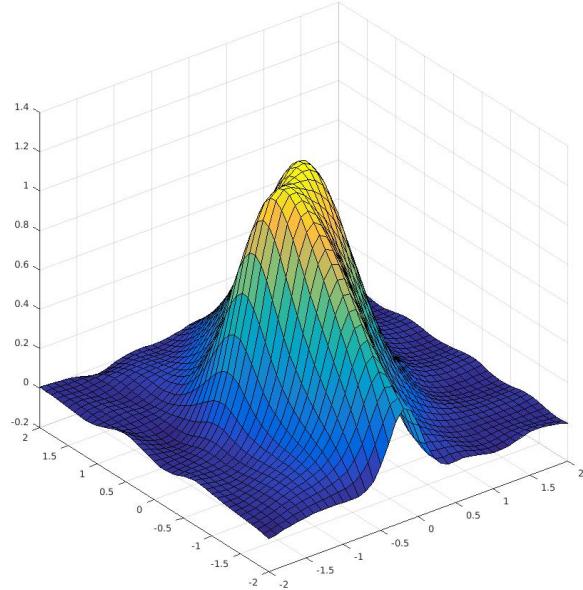


Figure 12: 2D plot of neural network output that was trained and implemented using *Netlab* toolbox

3 Exercise 3 - EM and doping

3.1

In figure 13 we can see histogram representing distribution of each variable (x_1, x_2, x_3, x_4). Figure 14 depicts a 3D plot of all 4 variables where dot color

stands for variable x_4 . Based on that graph we made our prediction of data clusters (black ellipses on the graph). We think that there might be 4 different classes and that the "middle" ellipse represents athletes that have substance X in their blood because that cluster shows strong correlation between variables x_1 and x_2 .

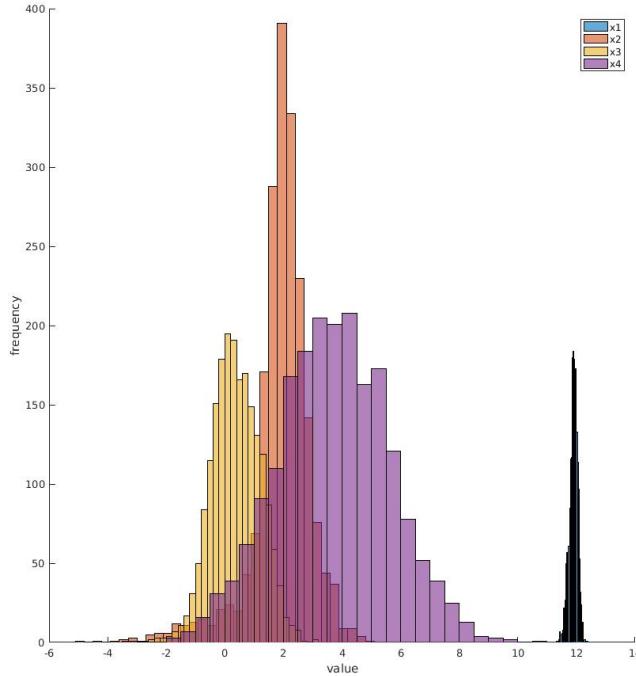


Figure 13: Histogram of given data (x_1, x_2, x_3, x_4)

3.2

EM algorithm:

- Initialize the means μ_k , covariances Σ_k and mixing coefficients π_k , and evaluate the initial value of the log likelihood.
- **E step.** Evaluate the responsibilities using the current parameter values

$$\gamma(z_{nk}) = \frac{\pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \quad (15)$$

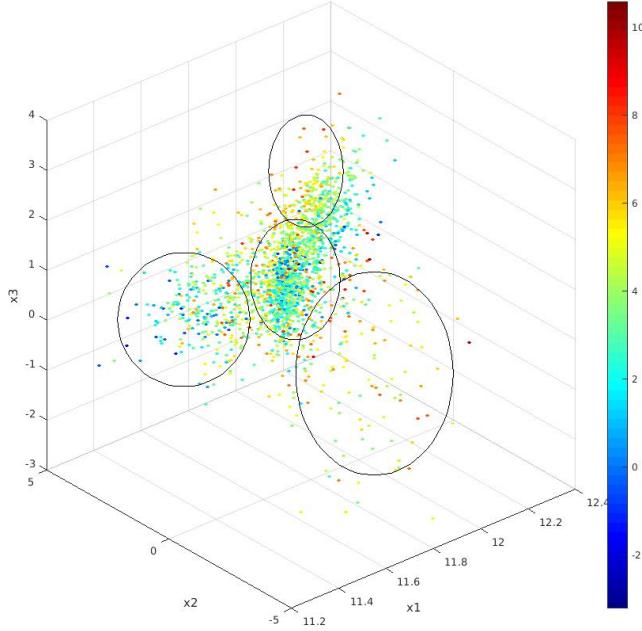


Figure 14: 3D scatter plot representing all 4 variables (x_1, x_2, x_3, x_4) of given data where color represents variable x_4

- **M step.** Re-estimate the parameters using the current responsibilities

$$\boldsymbol{\mu}_k^{new} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n \quad (16)$$

$$\boldsymbol{\Sigma}_k^{new} = \frac{1}{N_k} \sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \boldsymbol{\mu}_k^{new})(\mathbf{x}_n - \boldsymbol{\mu}_k^{new})^T \quad (17)$$

$$\pi_k^{new} = \frac{N_k}{N} \quad (18)$$

where

$$N_k = \sum_{n=1}^N \gamma(z_{nk}) \quad (19)$$

- Evaluate the log likelihood

$$\ln(p(\mathbf{X}|\boldsymbol{\mu}, \boldsymbol{\Sigma}, \boldsymbol{\pi})) = \sum_{n=1}^N \ln\left(\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)\right) \quad (20)$$

and check for convergence of either the parameters or the log likelihood.
If the convergence criterion is not satisfied return to step 2.

Next we present our implementation of EM algorithm in MATLAB where initialization is done in the way described in this exercise.

```

1 X = load('a011_mixdata.txt', '-ASCII');
2 [N, D] = size(X);
3
4 %INITIALIZATION
5
6 %number of classes is set to 2
7 K = 2;
8 pi = [0.5 0.5]';
9 meanx1 = mean(X(:,1));
10 meanx2 = mean(X(:,2));
11 meanx3 = mean(X(:,3));
12 meanx4 = mean(X(:,4));
13 means = [meanx1 meanx2 meanx3 meanx4];
14
15 %columns represents variable from x1 to x4
16 %rows represents K different classes
17 mu = zeros(K,D);
18 %sigma(:,:,k) represents covariance matrix for k-th class
19 sigma = zeros(D, D, K);
20 for i=1:K
21     mu(i,:) = means + rand(1,4).*2-1;
22     sigma(:,:,i) = diag(4.*rand(1,D)+2);
23 end
24
25 initLH = evalLogLH(X, mu, sigma, pi);
26 fprintf('Init log likelihood = %f\n', initLH);
27
28 maxIter = 100;
29 for i=1:maxIter
30     %Expectation step
31     resp = responsibility(X,mu,sigma,pi);
32
33     %Maximization step
34     Nk = sum(resp);
35     for k=1:K
36         mu(k,:) = (resp(:,k)'*X)./Nk(k);
37     end
38

```

```

39   for k=1:K
40     temp = zeros(D,D);
41     for n=1:N
42       temp = temp + resp(n,k) .* (X(n,:)-mu(k,:))' * (X(n,:)-mu(k,:));
43     end
44     sigma(:,:,k) = temp ./ Nk(k);
45   end
46
47   pi = Nk ./ N;
48
49   LH = evalLogLH(X, mu, sigma, pi);
50   fprintf('Step %d: log likelihood = %f\n', i, LH);
51 end

```

Plot routine that plots the x_1, x_2 coordinates of data points and color each data point according to the most probable component according to the mixture model is given by following MATLAB code:

```

1  %probs is a NxK matrix where N is number of data set points and K is
2  %number of classes. Every row of this matrix represents probabilities
3  %of assigning given point to class represented by column.
4  probs = zeros(N,K);
5  for k=1:K
6    probs(:,k) = mvnpdf(X, mu(k,:), sigma(:,:,k));
7  end
8  [M, I] = max(probs');
9
10 %scale I
11 a = 0; b = 1;
12 mn = min(I); mx = max(I);
13
14 mycolormap = colormap('winter');
15 inp = 0:63;
16 d64 = inp./63;
17 I = ((b-a).*(I-mn))./(mx-mn) + a;
18 c = interp1(d64, mycolormap, I);
19
20 dotsize = 10;
21 sca = scatter(X(:,1), X(:,2), dotsize,c,'fill');
22 colorbar;
23 xlabel('x1'); ylabel('x2');

```

3.3

For K=2 log likelihood converges to -7058.126254 in 40-50 steps. On figure 15 we can see scatter plot of the x_1 and x_2 coordinates colored according to the most probable component. Even with different initializations we got very similar results for K=2.

Correlation coefficient is -0.053 for blue class on graph and -0.323 for green class on graph. This means there is no correlation in the blue class, and a medium negative correlation in the green class. Neither class shows the characteristic strong positive correlation between x_1 and x_2 . Therefore it is safe to say that the data is not made up of only two Gaussians.

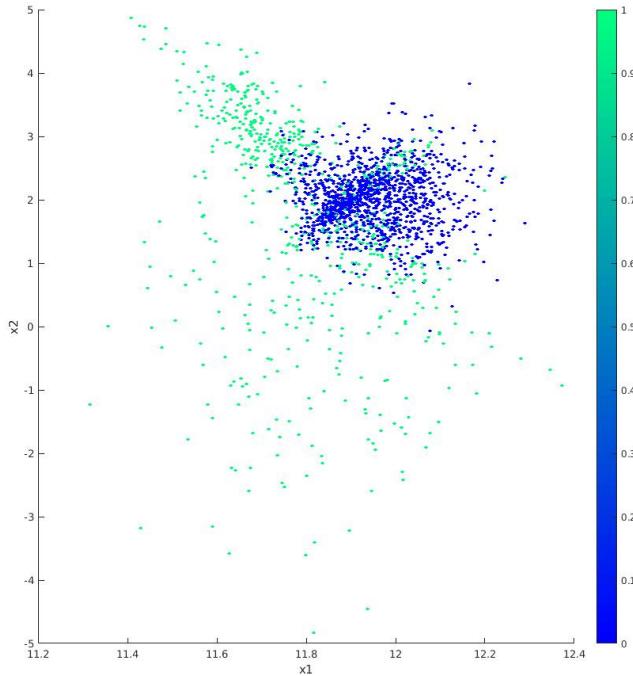


Figure 15: Scatter plot of data points colored according to most probable component for K=2

3.4

For K=3 log likelihood converges to -6551.130294. On figure 16 we can see scatter plot of the x_1 and x_2 coordinates colored according to the most probable component. Correlations for the green, blue and red classes are respectively: 0.081, -0.72 and -0.017. Still, there is no strong positive correlation in the components.

For K=4 log likelihood converges to -5989.752406. On figure 17 we can see scatter plot of the x_1 and x_2 coordinates colored according to the most probable component. Correlations for red, dark blue, light blue and yellow classes are respectively: 0.92, 0.036, -0.89 and -0.053. These show the strong positive correlation between x_1 and x_2 in the red class and the strong negative correlation that was hypothesized as a possible group for clean athletes in the light blue class. From the plots for various values of K we can assume that K = 4 is the most likely case. The portion of all data points in the red class is 0.1958. Therefore it is safe to say that the estimate of 1-in-5 of the users using doping is credible.

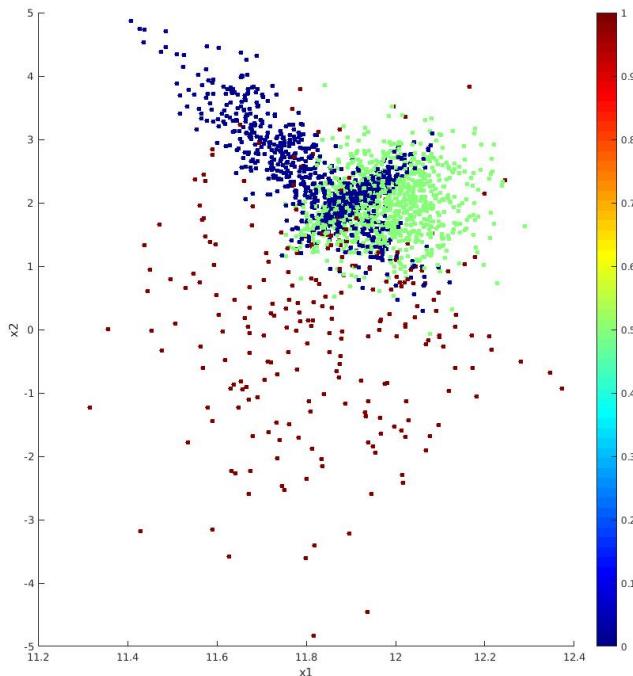


Figure 16: Scatter plot of data points colored according to most probable component for K=3

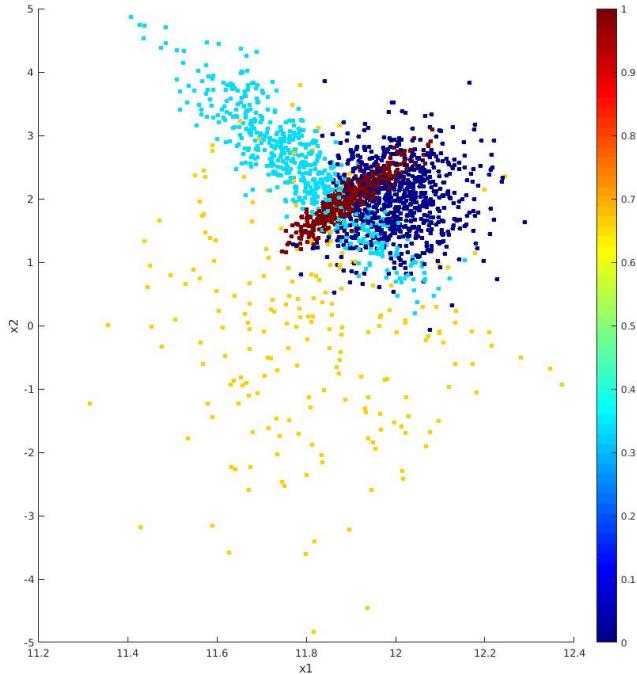


Figure 17: Scatter plot of data points colored according to most probable component for $K=4$

3.5

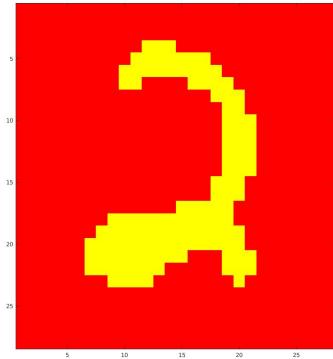
In matrix presented in expression 21 we can see probability of samples A, B, C and D belonging to certain classes. Rows represents samples A, B, C and D while columns represents probabilities of belonging to class 1 (dark blue), 2 (light blue), 3 (yellow) or 4 (red), respectively. From that matrix we can conclude that subject C took the drug X because sample C is most likely to belong to class 4 (red). Subject B tried to tamper with the test because probability of belonging to any class for that sample is 0 or very close to 0.

$$\begin{pmatrix} 0.1218 & 0.0857 & 0.0012 & 0.0001 \\ 0.0000 & 0.0000 & 0.0000 & 0.0000 \\ 0.0000 & 0.0000 & 0.0000 & 0.2484 \\ 0.1089 & 0.0000 & 0.0007 & 0.0000 \end{pmatrix} \quad (21)$$

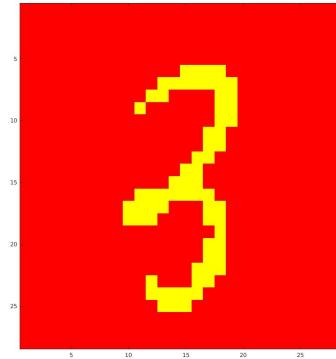
4 Exercise 4 - Handwritten digit recognition

4.1

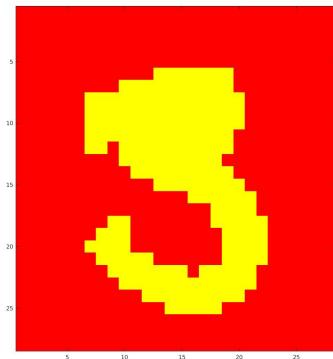
Some of the digits from data set are given in figure 18. We think it will be hard to classify them correctly since every number can be handwritten in many different ways. For example, in figure 18 we can see that number 3 is written in 2 completely different ways.



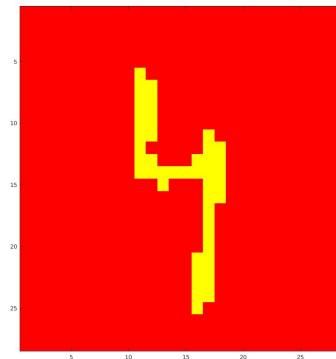
(a) Two



(b) Three



(c) Three



(d) Four

Figure 18: Examples of handwritten digits from data set

4.2

The algorithm is very similar to the algorithm that was used in the first exercise. The only difference being that a Bernoulli pdf is used instead of a Gaussian, and

the computation of Σ is omitted. Images are drawn after each iteration using following MATLAB snippet:

```
1 mudraw = mu.*255;
2 d = sqrt(D);
3 for k=1:K
4     image(reshape(mudraw(k,:),d,d));
5     colormap(gray(255));
6     pause;
7 end
```

4.3

In figure 19 greyscale images of pixel means are presented. First row stands for means after first iteration and second row for means after 40 iterations. Algorithm converges quickly and from the images we can clearly see numbers 2, 3 and 4. Even with different initializations we got similar results and algorithm converged quickly.

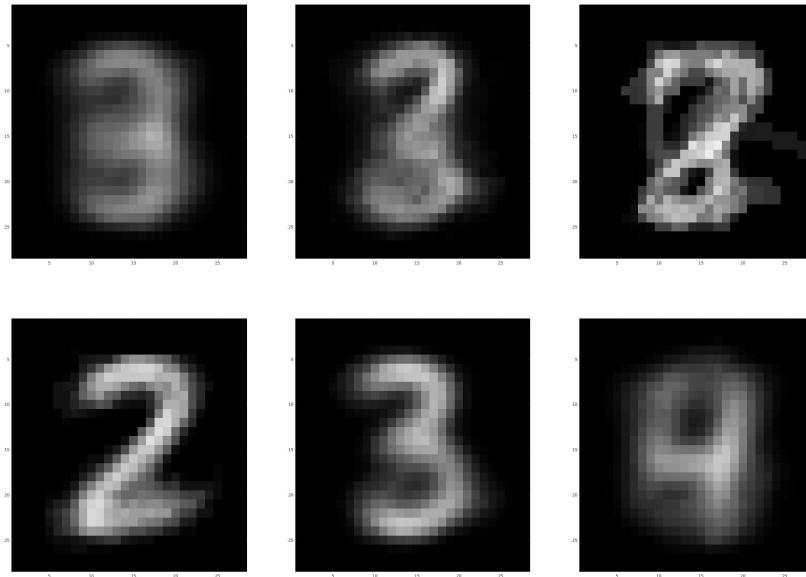


Figure 19: Greyscale images of pixels means where $K = 3$

4.4

In figure 20 we can see images of means with 2 classes and in figure 21 images of means with 4 classes. If we use $K > 3$ we get interesting results. Most of the time, we get cluster prototypes for different variations of the same digit, e.g. a curly 2 vs. a straight 2 or a closed 4 vs. an open 4 etc. For 2 classes we usually got results of some kind of hybrids between some digits.

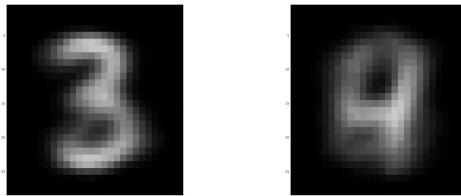


Figure 20: Greyscale images of pixels means where $K = 2$

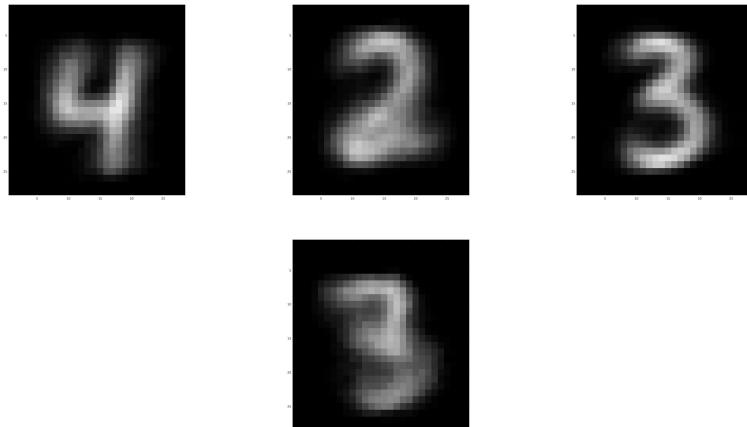


Figure 21: Greyscale images of pixels means where $K = 4$

Our algorithm correctly classified 91.75% of images. In figure 22 we can see one of the missclassified images. Our algorithm classified it as number 4 while it is actually number 3. Every human would say that it is number 3 but it isn't a "proper" shape of the number so it has been missclassified by algorithm. Based on that we can conclude that it is hard for algorithm to classify different shapes of particular numbers correctly.

If we initialize our algorithm with true mean values than it converges more quickly and there are less missclassified images.

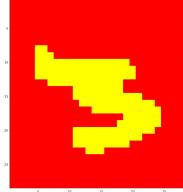


Figure 22: Missclassified number 3

4.5

We created a small dataset of 15 integers in the MNIST style and calculated its log likelihoods on our bernoulli class prototypes, which can be seen in figure 23. 7/15 (0.4667) of the images could be given the right label. Noticeably, the scale of the digits in the image is significant. Some of the digits that were drawn bigger and close to the border were missclassified with infinite $p(x)$.

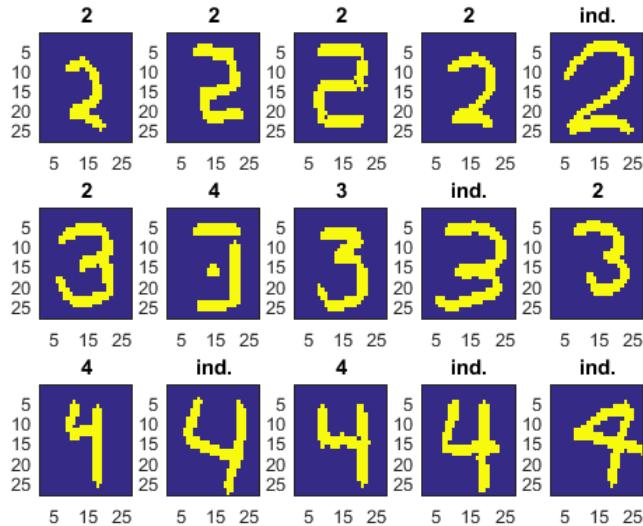


Figure 23: MAP- classifications of a self-created MNIST-style data set.