

Machine Learning and Pattern Recognition Coursework

Zhen Yi, s1563190

November 24, 2015

NOTE — BEFORE YOU START

The coursework is formatted by \LaTeX , if you want a **doc** format copy, please contact me without any hesitation.

Environment for this coursework

For this coursework, I used **Matlab** as the development tool, and I have some references such as the official document of **Matlab**, the papers and the book which are helpful in theory part.

Solution — 2 parts

Part 1

1.a) For this question, I used command `histogram(std.tr, 64)`. Please refer to Figure 1. The reason why 64 bins is that although we scaled all the data to 64 bins, we don't change the consistency of data. Thus, we still can split the data to 64 bins.

This plot tells us the approximate distribution of the standard deviation of the

patches. It looks like that the most of standard deviation coverge to zero, thus we can assume that most of the patches are flat patches.

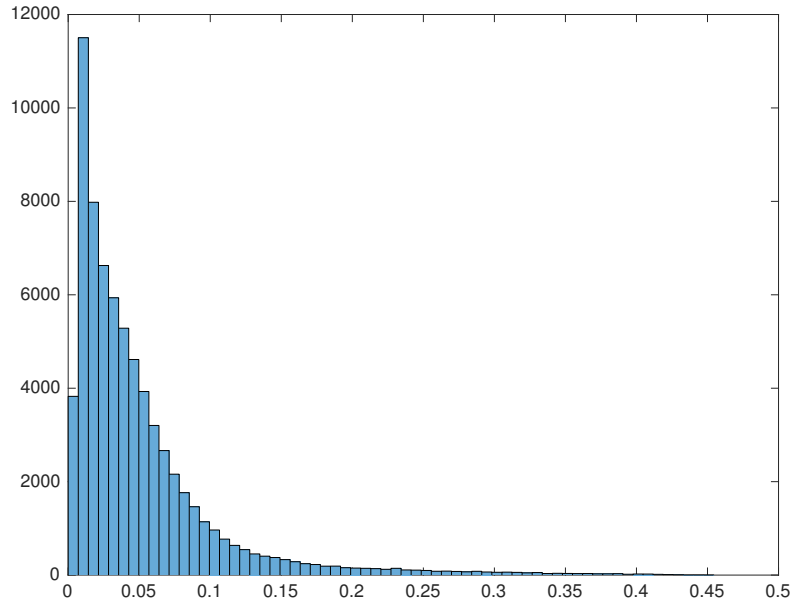


Figure 1: Histogram(using 64 bins)

1.b) This question ask us to predict the target pixel, because we have already had the result of standard deviation. We can simply use such result. The meaning of standard deviaion is that if the closer it is to zero, the more flat it is. Thus, we can use the function **mean**, which means that most the patch can be seen as flat. Thus, the value of the target pixel closes to the mean of the whole patch.

```
1 targetPixel = mean(xtr_1a, 2); %(1.b)
```

1.c) Figure 2 shows one flat and one non-flat image patch, respectively.

```
1 flatPatch = find(std_tr ≤ err);
2 nonflatPatch = find(std_tr > err);
3 nononeFlatPatch = [xtr_1a(nonflatPatch(10000), :) zeros(1, 18)]; ...
   % find a non-flat image patch locates in 2000
```

```

4 nononeFlatPatch = reshape(nononeFlatPatch, [35, 30]);
5
6 oneFlatPatch = [xtr_1a(flatPatch(50), :) zeros(1, 18)]; % find a ...
    non-flat image patch locates in 100
7 oneFlatPatch = reshape(oneFlatPatch, [35, 30]);
8
9 colormap gray;
10
11 subplot(1,2,1), imagesc(oneFlatPatch, [0,1]);
12 subplot(1,2,2), imagesc(nononeFlatPatch, [0,1]);

```

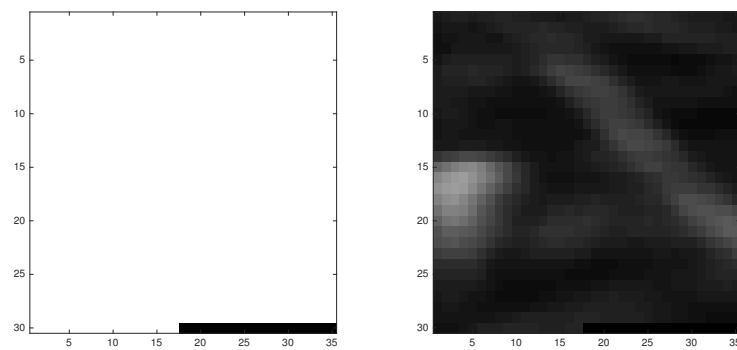


Figure 2: Flat and Non-flat pictures

Note: Here is the full Matlab code for the question 1.

```

1 %load('imgregdata.mat');
2
3 xtr_1a = xtr / 63.0; % scaling the training data
4 ytr_1a = ytr / 63.0;
5
6 xte_1a = xte / 63.0; % scaling the test data
7 yte_1a = yte / 63.0;
8
9 std_tr = std(xtr_1a,0,2); % standard derivation
10
11 histogram(std_tr, 64); %(1.a)
12
13 err = 0 / 64.0;
14 targetPixel = mean(xtr_1a, 2); %(1.b)
15
16 flatPatch = find(std_tr <= err);
17 nonflatPatch = find(std_tr > err);
18 nononeFlatPatch = [xtr_1a(nonflatPatch(10000), :) zeros(1, 18)]; ...
    % find a non-flat image patch locates in 2000
19 nononeFlatPatch = reshape(nononeFlatPatch, [35, 30]);
20

```

```

21 oneFlatPatch = [xtr_1a(flatPatch(50), :) zeros(1, 18)]; % find a ...
    non-flat image patch locates in 100
22 oneFlatPatch = reshape(oneFlatPatch, [35, 30])';
23
24 colormap gray;
25
26 subplot(1,2,1), imagesc(oneFlatPatch, [0,1]);
27 subplot(1,2,2), imagesc(nononeFlatPatch, [0,1]);

```

2.a) Figure 3 shows that the points cluster near a point (0.25, 0.25, 0.25) which closes to zero. It looks like a comet. Most of the points cluster around the such centre point (0.25, 0.25, 0.25), the lower bound and upper bound of this body is (0, 0.5]. It also has a tail from (0.5, 1.0).

Substantially, it can be seen as positive linear correlation (from (0,0,0) to (1,1,1)). And the part of the points which below (0.5, 0.5, 0.5) is more related to the linear correlation than the points which above (0.5, 0.5, 0.5).

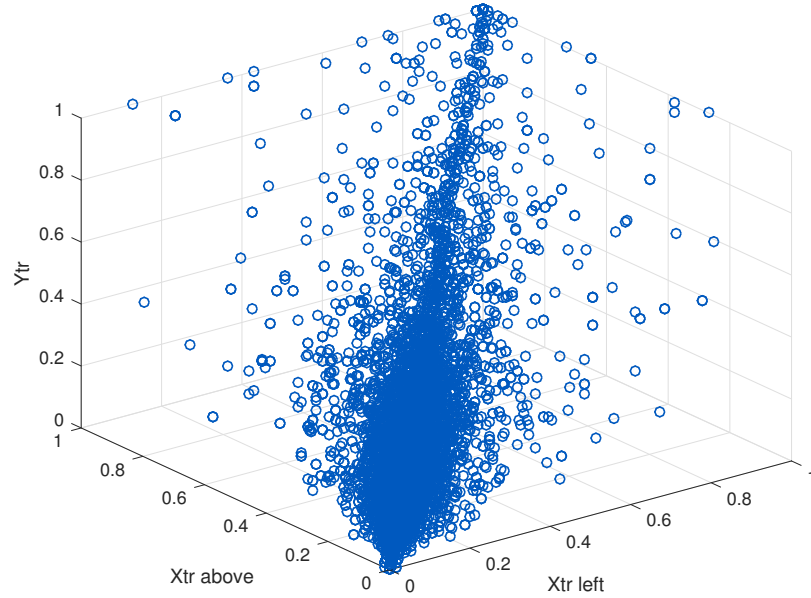


Figure 3: Subsample

2.b) The kernel code is $w = \text{inv}(X' * X) * X' * y_{tr_nf}$. We can get the solution of w using X and the vector y :

$$\begin{aligned}\log P(y|x, w) &= -\frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2}(y - w^T X)^2 \\ L(w) &= \sum_{i=1}^n (\log P(y_i|x_i, w)) \\ \frac{\partial}{\partial w} L(w) &= \frac{\partial}{\partial w} \left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - w^T X)^2 \right) \\ \hat{w} &= (X^T X)^{-1} X^T y\end{aligned}$$

```

1 X = [xtr_nf(:,end) xtr_nf(:,end - 34) (ones(1, length(xtr_nf)))'];
2
3 w = inv(X' * X) * X' * ytr_nf;
4 % w = regress(ytr_nf,X); I used this function to check my answer

```

2.c) The linear regression predictor is $y = X * w$. RMSE of training data is 0.0506, RMSE of test data is 0.0503.

My linear regressor use the closed-form maximum likelihood for weight vector \vec{w} , it can work in the simple data sample, however it may not as powerful as we predicting. Because, it doesn't include any iteration to reduce the error.

Now, we need to analyze why RMSE of test data is less than the training data's. I think, there are more outliers in the training data. Thus, such circumstance happens.

$$w = \begin{bmatrix} 0.4606 \\ 0.5241 \\ 0.0026 \end{bmatrix}$$

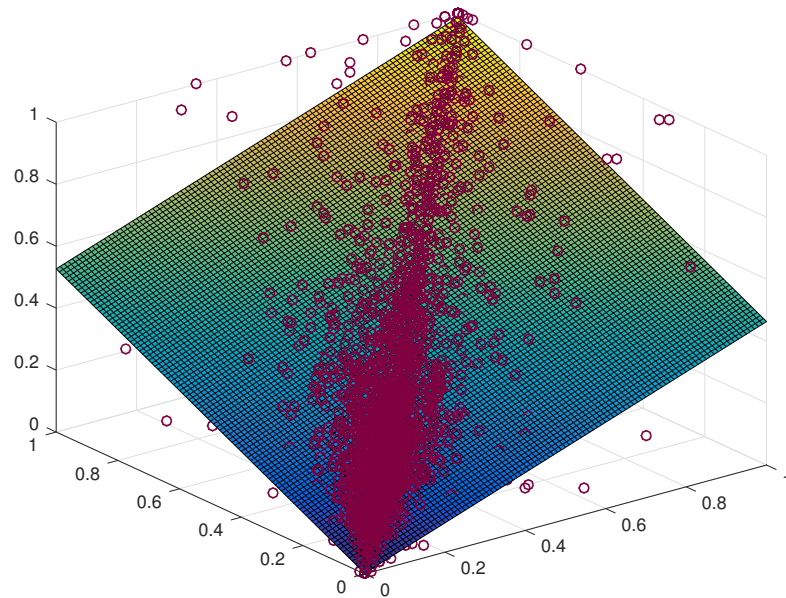


Figure 4: Surface and 3-D plot of test data

Note: Here is the full Matlab code for the question 2.

```

1  load('imgregdata.mat');
2
3  position = datasample(1:length(xtr_nf), 10000); %subsample the ...
   position
4
5  xtr_left = [];
6  xtr_above = [];
7  ytr_new = [];
8
9  for i = 1 : 10000
10     xtr_left = [xtr_left; xtr_nf(position(i), end)];
11     xtr_above = [xtr_above; xtr_nf(position(i), end - 34)];
12     ytr_new = [ytr_new; ytr_nf(position(i))];
13 end;
14
15 scatter3(xtr_left, xtr_above, ytr_new, 'MarkerEdgeColor',[0 .35 ...
   .75]);
16
17 xlabel('Xtr left');
18 ylabel('Xtr above');
19 zlabel('Ytr');

```

```

20
21 % 2.b
22 X = [xtr_nf(:,end) xtr_nf(:,end - 34) (ones(1, length(xtr_nf)))'];
23
24 w = inv(X' * X) * X' * ytr_nf;
25 % w = regress(ytr_nf,X); I used this function to check my answer
26
27 % 2.c
28 ytr_regress = X * w;
29 RMSE_ytr = sqrt(mean((ytr_regress - ytr_nf).^2));
30 yte_regress = [xte_nf(:,end) xte_nf(:,end - 34) ...
    ones(length(xte_nf), 1)] * w;
31 RMSE_yte = sqrt(mean((yte_regress - yte_nf).^2));
32
33 figure,
34 [dim1, dim2] = meshgrid(0:0.01:1, 0:0.01:1);
35 ysurf = [[dim1(:), dim2(:)], ones(numel(dim1), 1)] * w;
36 surf(dim1, dim2, reshape(ysurf, size(dim1)))
37
38 hold on
39
40 scatter3(xte_nf(:,end), xte_nf(:,end - 34), yte_nf(:), ...
    'MarkerEdgeColor',[0.5 .0 .25]);

```

3.a) In this question, I used 10-fold cross validation and only map the argument of training data to function. Figure 5 shows the RMSE of 6 different nbf functions.

I have run the **matlab** program 10 times, and plot them in one diagram. You can see that it is fluctuant. Therefore, we should choose the most stable one from the all 6 possible choices. By observation, I will choose 10 for next problem.

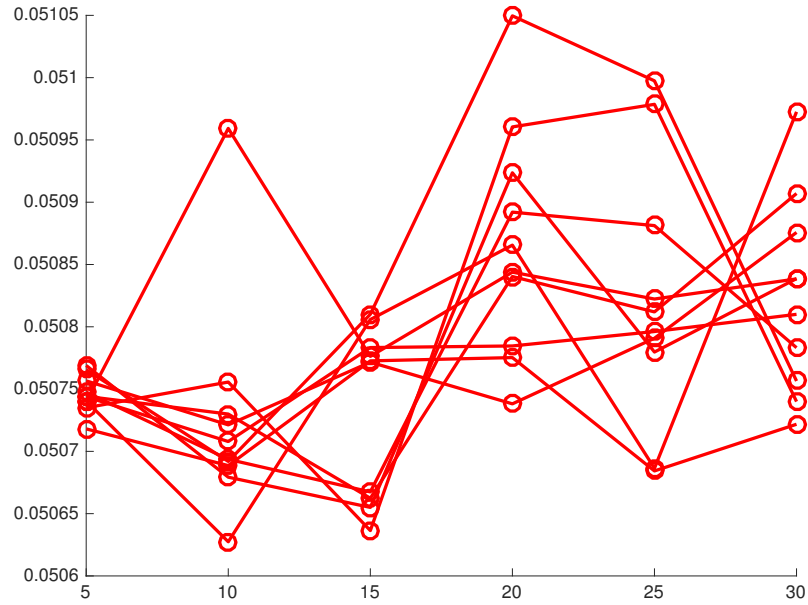


Figure 5: cross-validation RMSE

Note: Here is the full Matlab code for the question 3.a.

```

1  load('imgregdata.mat');
2  %code for creating RBF and making predictions
3  %nbf is the number of basis functions
4  %dim is the dimensionality of input space
5
6  nbf = [5,10,15,20,25,30];
7  format long; % set display precision
8  options = foptions;
9  options(1) = 1;
10 options(14) = 5;
11 TrainX = xtr_nf(:,[end end - 34]);
12 %TestX = xte_nf(:,[end end - 34]);
13
14 for j= 1:10
15     for i = 1:6
16         dim = 2;
17         net = rbf(dim, nbf(i), 1, 'gaussian');
18         %10-fold cross validation
19         rbff=@(XTRAIN,ytrain,XTEST)(rbffwd(rbftrain(net, options, ...
20             XTRAIN, ytrain), XTEST));
20         RMSE(i) = sqrt(crossval('mse',TrainX, ytr_nf, 'Predfun',rbff));

```



```

21 end
22
23 hold on;
24 plot(nbf, RMSE, 'r', 'LineWidth', 2, 'Marker', 'o', 'MarkerSize', 10);
25 end

```

3.b) For this question, the result is better the linear regression, but they are very close. And the method of RBF is totally different from linear regression, because of non-linear feature. And the result of test is better than the result of training.

$$RMSE_{training} = 0.0505$$

$$RMSE_{test} = 0.0501$$

```

1 load('imgregdata.mat');
2 %code for creating RBF and making predictions
3 %nbf is the number of basis functions
4 %dim is the dimensionality of input space
5 net = rbf(2, 10, 1, 'gaussian');
6
7 options = foptions;
8 options(1) = 1;
9 options(14) = 5;
10
11 rbff= rbfftrain(net, options, xtr_nf(:,[end end - 34]), ytr_nf);
12 ypred_tr = rbffwd(rbff, xtr_nf(:,[end end - 34]));
13 ypred_te = rbffwd(rbff, xte_nf(:,[end end - 34]));
14
15 RMSE_ytr = sqrt(mean((ypred_tr - ytr_nf).^2));
16 RMSE_yte = sqrt(mean((ypred_te - yte_nf).^2));

```

4) In this question, we implement a linear regressor with full pixels. The result are : RMSE of training data is 0.0371, RMSE of test data is 0.0456.

The RMSEs are less than previous result which produced by that linear regressor. The more dimension of the train data is, the less RESE we can get. Therefore, we got the better solution.

Of course, we can get more features from the larger amount of information. The reason why we use non-flat training data is that we should reduce the outliers. For this example, the over-fitting occurs, thus the gap between RMSE of training data and RMSE of test data is big.

Note: Here is the full **Matlab** code for the question 4.

```

1 X = [xtr_nf (ones(1, length(xtr_nf)))'];

```

```

2
3 w = inv(X' * X) * X' * ytr_nf;
4 %w = regress(ytr_nf,X); %I used this function to check my answer
5
6 ytr_regress = X * w;
7 RMSE_ytr = sqrt(mean((ytr_regress - ytr_nf).^2));
8 yte_regress = [xte_nf ones(length(xte_nf), 1)] * w;
9 RMSE_yte = sqrt(mean((yte_regress - yte_nf).^2));

```

5.a) In this question, I used the well trained MLP to compute both of training data and test data. The RMSE of training data is 0.0333, the RMSE of test data is 0.0473. The RMSE of training data is less than such data of linear regression. However, the RMSE of test data is greater than such data of linear regression. Only for this example, we can assert that sometimes the performance of NN is not as good as linear regression. I guess the reason is that, they used different training data or the NN is more appropriate training data, but not test data. This may cause the over-fitting. Therefore, the result of test data is worse than that of linear regression.

```

1 %load('imgregdata.mat');
2
3 ypred_tr = mlpfwd(net, xtr_nf);
4 rmse_NNsuball_tr = sqrt(mean((ytr_nf - ypred_tr).^2));
5
6 ypred = mlpfwd(net, xte_nf);
7 rmse_NNsuball_te = sqrt(mean((yte_nf - ypred).^2));

```

5.b) Referring to the table. I report the result of different seed. The gap between the RMSE of training data and RMSE of test data gets larger. Because we have tried 5 seeds. We have the different initialization of the corresponding seeds. And the results of the performance are similar. So we get the conclusion is that limited data cause the over-fitting problem.

Seed	2015	2016	2017	2018	2019
RMSE train	0.0387	0.0363	0.0346	0.0334	0.0347
RMSE test	0.0515	0.0504	0.0515	0.0516	0.0527

```

1 rng(2019, 'twister');
2 nhid = 10;
3 net = mlp(size(xtr_nf, 2), nhid, 1, 'linear');
4
5 options = zeros(1, 18);
6 options(1) = 1;
7 options(9) = 1;
8 options(14) = 200;

```

```

9
10 tic
11 [net, options] = netopt(net, options, xtr_nf(1: 5000, :), ...
    ytr_nf(1:5000,:), 'scg');
12 toc
13
14 ypred_tr = mlpfwd(net, xtr_nf(1:5000,:));
15 rmse_NNsuball_tr = sqrt(mean((ytr_nf(1:5000) - ypred_tr).^2));
16
17 ypred = mlpfwd(net, xte_nf);
18 rmse_NNsuball_te = sqrt(mean((yte_nf - ypred).^2));

```

6) First of all, I think we should compare the methods in different aspects. We cannot compare them directly, because there exists some different factors which influenced the results.

- Firstly, let us compare the linear regression and RBF network using 2 neighbouring pixels. The gap between two RMSEs is very small, which means we can think that they have the similar performance of the result. However, the efficiency of linear regression is much better than the RBF's. Because we do need to choose the number of functions in RBF network. Thus, at least in this training data set, the linear regression is better.
- Secondly, in question 4 we used whole image patch to train the linear regressor. And the performance of such regressor is better than the one which used 2 neighbouring pixels (because of less RMSE). The reason why that happens is, in my opinion, the amount of the information is larger. Thus, the cost function can get more features to do the gradient descent (although we used the pseudo-method here). But the RMSE of test data is still higher than the training one. I think this may be the problem of over-fitting. It is a serious problem. We should decide the initial parameters carefully.
- Thirdly, let us compare neural network with linear regression. In problem 4 and 5, we used all the dimension of the training data to train the neural network and regressor. Obviously, the method of neural network is much complicated than the linear regression. And the serious problem of over-fitting occurs in the neural network. Thus, the performance of NN is not as good as linear regression.
- Finally, there are some different schemes that can strengthen the performance of above methods. 1) We can use different features of the training data in linear regression. For example, we can use the whole row on the left of the predicting pixel and whole column on the above of the predicting pixel. The reason is that, we can do the experiment to find which feature

is the best one of the problem. 2) We can add the prior probability. For example, bayes regression and bayes network.

Part 2

1.a) For this question, the code is been showed as follow. I used the basic method to implement the target which augment a 'bias' weight. Thus, we can avoid adding extra inputs and outputs to the likelihood function.

```
1 load('text_data.mat');
2 %a: augment the data
3 x_te = [x_test ones(length(x_test), 1)];
4 x_tr = [x_train ones(length(x_train), 1)];
```

1.b) The follow routine **t21b** is negative-log-likelihood function:

```
1 function [Lp, dLp_dw] = t21b(ww, xx, yy)
2 yy = (yy==1)*2 - 1;
3
4 sigmas = 1./(1 + exp(-yy.*(xx*ww))); % Nx1
5 Lp = -sum(log(sigmas));
6
7 if nargout > 1
8     dLp_dw = -(((1-sigmas).*yy)' * xx)';
9 end
```

We can use the above function to compute the data we need. Now, let us compare the two different sets. Firstly, the accuracy of training data is less than that of test data. It is strange. However, it might be correct because of the features of data or other factors.

For the mean log probability(MLP for short), the baseline is $\log 0.5 = -0.693147$. Thus, both the MLP are bigger than the baseline. And the performance of the test data is still better than the training one.

$$\begin{aligned} \text{Accuracy of test} &= 0.9083 \\ \text{Accuracy of training} &= 0.8344 \\ \text{Mean log probability of test} &= -0.2881 \\ \text{Mean log probability of training} &= -0.4400 \end{aligned}$$

Here is the full **matlab** code of 1.b.

```

1 %b
2 b_tr = rand(101, 1) * 10;
3 [X, fX, i] = minimize(b_tr, 't21b', 100, x_tr, y_train);
4
5 y_prob = 1./(1 + exp(-(x_te*X)));
6 tot = 0;
7 for i = 1 : length(y_prob)
8     if (y_prob(i) > 0.5)
9         temp = 1;
10    else
11        temp = -1;
12    end
13    if (temp == y_test(i))
14        tot = tot + 1;
15    end
16 end
17 y_prob_ml = 1./(1 + exp(-y_test.*(x_te*X)));
18 accuracy = tot / length(y_test);
19 y_mlog = mean(log(y_prob_ml));
20
21 y_prob_tr = 1./(1 + exp(-(x_tr*X)));
22 tot_t = 0;
23 for i = 1 : length(y_prob_tr)
24     if (y_prob_tr(i) > 0.5)
25         temp = 1;
26     else
27         temp = -1;
28     end
29     if (temp == y_train(i))
30         tot_t = tot_t + 1;
31     end
32 end
33 y_prob_tr_ml = 1./(1 + exp(-y_train.*(x_tr*X)));
34 accuracy_t = tot_t / length(y_train);
35 y_mlog_t = mean(log(y_prob_tr_ml));
36 %standard error
37 std_err_a = sqrt(accuracy.*(1-accuracy)/size(y_prob, 1));
38 std_err_atr = sqrt(accuracy_t.*(1-accuracy_t)/size(y_prob_tr, 1));
39 std_err = sqrt(var(log(y_prob_ml))/size(y_prob, 1));
40 std_err_tr = sqrt(var(log(y_prob_tr_ml))/size(y_prob_tr, 1));

```

1.c) Here we limited the data to first 100 training cases. The accuracy and the mean log probability descend dramatically. I think the reason is that the training cases is not enough to train the log regressor well. For this example, the first 100 cases cannot represent the general feature of the whole data sets.

$$\begin{aligned}
 \text{Accuracy of test} &= 0.7520 \\
 \text{Mean log probability of test} &= -inf
 \end{aligned}$$

Here is the full **matlab** code of 1.c.

```

1 load('text_data.mat');
2
3 x_te = [x_test ones(length(x_test), 1)];
4 x_tr = [x_train ones(length(x_train), 1)];
5
6 b_tr = ones(101, 1);
7 [X, fX, i] = minimize(b_tr, 't21b', 10000, x_tr(1:100,:), ...
    y_train(1:100,:));
8
9 y_prob_tr = 1./(1 + exp(-(x_te*X)));
10 tot_t = 0;
11 for i = 1 : size(y_prob_tr, 1)
12     if (y_prob_tr(i) > 0.5)
13         temp = 1;
14     else
15         temp = -1;
16     end
17     if (temp == y_test(i))
18         tot_t = tot_t + 1;
19     end
20 end
21 accuracy = tot_t / size(y_test, 1);
22 y_mlog_t = mean(log(y_prob_tr));

```

2.a) For this question, we need to compute the derivatives of both \mathbf{w} and ϵ . Such result can be used as the function which compute the log regression with the label noise.

The original function showed as below:

$$P(y|\mathbf{x}, \mathbf{w}, \epsilon) = (1 - \epsilon)\sigma(y\mathbf{w}^T \mathbf{x}) + \epsilon/2, \quad y \in \{-1, +1\}$$

Firstly, we should compute the derivative of \mathbf{w} :

$$\begin{aligned}
\mathcal{L}(\mathbf{w}, \epsilon) &= \sum_{n=1}^N \log P(y^{(n)}|x^{(n)}, \mathbf{w}, \epsilon) \\
\nabla_w \mathcal{L}(\mathbf{w}, \epsilon) &= \frac{\partial}{\partial w} \sum_{n=1}^N \log P(y^{(n)}|x^{(n)}, \mathbf{w}, \epsilon) \\
\nabla_w \mathcal{L}(\mathbf{w}, \epsilon) &= \frac{\partial y^{(n)} \mathbf{w}^T x^{(n)}}{\partial w} \frac{\partial}{\partial y^{(n)} \mathbf{w}^T x^{(n)}} \sum_{n=1}^N \log((1 - \epsilon)\sigma(y^{(n)} \mathbf{w}^T \mathbf{x}^{(n)}) + \epsilon/2) \\
\nabla_w \mathcal{L}(\mathbf{w}, \epsilon) &= y^{(n)} \mathbf{x}^{(n)} \sum_{n=1}^N \left[\frac{(1 - \epsilon)\sigma(y^{(n)} \mathbf{w}^T x^{(n)}) (1 - \sigma(y^{(n)} \mathbf{w}^T x^{(n)}))}{(1 - \epsilon)\sigma(y^{(n)} \mathbf{w}^T \mathbf{x}^{(n)}) + \epsilon/2} \right]
\end{aligned}$$

Secondly, we should compute the derivative of ϵ :

$$\begin{aligned}\mathcal{L}(\mathbf{w}, \epsilon) &= \sum_{n=1}^N \log P(y^{(n)}|x^{(n)}, \mathbf{w}, \epsilon) \\ \nabla_{\epsilon} \mathcal{L}(\mathbf{w}, \epsilon) &= \frac{\partial}{\partial \epsilon} \sum_{n=1}^N \log P(y^{(n)}|x^{(n)}, \mathbf{w}, \epsilon) \\ \nabla_{\epsilon} \mathcal{L}(\mathbf{w}, \epsilon) &= \frac{\partial}{\partial \epsilon} \sum_{n=1}^N \log((1 - \epsilon)\sigma(y^{(n)}\mathbf{w}^T \mathbf{x}^{(n)}) + \epsilon/2) \\ \nabla_{\epsilon} \mathcal{L}(\mathbf{w}, \epsilon) &= \sum_{n=1}^N \left[\frac{\frac{1}{2} - \sigma(y^{(n)}\mathbf{w}^T \mathbf{x}^{(n)})}{(1 - \epsilon)\sigma(y^{(n)}\mathbf{w}^T \mathbf{x}^{(n)}) + \epsilon/2} \right]\end{aligned}$$

Here are the **matlab** code of the above two derivatives.

```
1 function [Lp, dLp_dw] = t22f1(ww, eps, xx, yy)
2 yy = (yy==1)*2 - 1;
3 % derivative of w
4 sigmas = 1./(1 + exp(-yy.*(xx*ww))); % Nx1
5 Lp = sum(log((1-eps)*sigmas+(eps/2)));
6
7 if nargin > 1
8     dLp_dw = ...
9         ((1-eps)*sigmas.*(1-sigmas)./( (1-eps)*sigmas+eps/2 ).*yy)' ...
10         * xx)';
11 end
```

```
1 function [Lp, dLp_dw] = t22f2(eps, ww, xx, yy)
2 yy = (yy==1)*2 - 1;
3 %derivative of epsilon
4 sigmas = 1./(1 + exp(-yy.*(xx*ww))); % Nx1
5 Lp = sum(log((1-eps)*sigmas+(eps/2)));
6
7 if nargin > 1
8     dLp_dw = sum((0.5-sigmas)./( (1-eps)*sigmas+(eps/2) ));
9 end
```

Now we can use the routine **checkgrad.m** to check the derivatives of above two functions. This routine provides the difference between my gradients and their numerical approximation. However, because of the precision, I found that there might not exist the largest difference.

```
1 load('text_data.mat');
2 b_tr = rand(101, 1) * 10;
3 x_te = [x_test ones(length(x_test), 1)];
4 x_tr = [x_train ones(length(x_train), 1)];
5 checkgrad('t22f1', b_tr, 1e-3, rand(1, 1), x_tr, y_train);
6 checkgrad('t22f2', rand(1, 1), 1e-3, b_tr, x_tr, y_train);
```

2.b) In this problem, we will write the parameter of ϵ as the result of taking the logistic sigmoid of an unconstrained parameter a :

$$\epsilon = \sigma(a) = \frac{1}{1+\exp(-a)}$$

Now, we compute the derivative of a :

$$\begin{aligned}\mathcal{L}(\mathbf{w}, a) &= \sum_{n=1}^N \log P(y^{(n)}|x^{(n)}, \mathbf{w}, \sigma(a)) \\ \nabla_a \mathcal{L}(\mathbf{w}, a) &= \frac{\partial}{\partial a} \sum_{n=1}^N \log P(y^{(n)}|x^{(n)}, \mathbf{w}, \sigma(a)) \\ \nabla_a \mathcal{L}(\mathbf{w}, a) &= \frac{\partial}{\partial a} \sum_{n=1}^N \log((1 - \sigma(a))\sigma(y^{(n)}\mathbf{w}^T \mathbf{x}^{(n)}) + \sigma(a)/2) \\ \nabla_a \mathcal{L}(\mathbf{w}, a) &= \sum_{n=1}^N \left[\frac{\sigma(a)(1-\sigma(a))(\frac{1}{2}-\sigma(y^{(n)}\mathbf{w}^T \mathbf{x}^{(n)}))}{(1-\sigma(a))\sigma(y^{(n)}\mathbf{w}^T \mathbf{x}^{(n)})+\sigma(a)/2} \right]\end{aligned}$$

Here is the full **matlab** code of the function of 2.b.

```

1  load('text_data.mat');
2  %a: augment the data
3  x_te = [x_test ones(length(x_test), 1)];
4  x_tr = [x_train ones(length(x_train), 1)];
5
6  [X, fX, i] = minimize([ones(101,1)',0]', @t23f, 10000, x_tr, ...
    y_train);
7
8  w = X(1: end -1);
9  eps = 1./(1+exp(-w(end)));
10 y_prob = 1./(1 + exp(-(x_te*w)));
11 tot = 0;
12 for i = 1 : length(y_prob)
13     if (y_prob(i) > 0.5)
14         temp = 1;
15     else
16         temp = -1;
17     end
18     if (temp == y_test(i))
19         tot = tot + 1;
20     end
21 end
22
23 y_prob_ml = 1./(1 + exp(-y_test.*(x_te*w)));
24 accuracy = tot / length(y_test);
25 y_mlog = mean(log(y_prob_ml));

```


Now, we use the program to compute accuracy and mean log probability(MLP for short). We can get the new result as below:

$$\begin{aligned} accuracy &= 0.9108 \\ mean \quad log \quad probability &= -1.3486 \end{aligned}$$

For this result, we found that the MLP is less than the last problem. And the accuracy is similar to that problem. We can observe the predicting value of y , and the value are much more extreme. The reason why that happens is that, I think, we use the constrained parameter to build the sigmoid function of ϵ . On the other hand, because the training data of the problem may not always reliable. Thus, in this problem, we lead into the label noise model. If $s_n = 1$ then the label is chosen with a uniform random choice and ignore the original \mathbf{x} features. Therefore, we may guess the correct choice which is not reliable by balance the original function with constrained parameter.

Here is the full **matlab** code of the function of derivative.

```
1 function [Lp, dLp_dw] = t22f1n(w, x, y)
2 y = (y==1)*2 - 1;
3 a = w(end);
4 ww = w(1: end - 1);
5 eps = 1./(1+exp(-a));
6 sigmas = 1./(1 + exp(-y.*(x*ww))); % Nx1
7 Lp = -sum(log((1-eps)*sigmas+(eps/2)));
8
9 a = -(((1-eps)*sigmas.*(1-sigmas)./((1-eps)*sigmas+eps/2).*y)' * ...
10 x)';
11 b = -sum((0.5-sigmas)*eps*(1-eps)./((1-eps)*sigmas+(eps/2)));
12 dLp_dw = [a', b]';
```

3.a) The largest log-likelihood that any model can given a training set of N binary outcomes equal to the log-likelihood if a model were somehow able to predict every training label correctly. Thus, it is just 1.

$$\begin{aligned} &\text{Let } w = 0 \text{ in} \\ \log P(\epsilon, \mathbf{w}, \log \lambda | data) &= \mathcal{L}(\mathbf{w}, \epsilon) - \lambda \mathbf{w}^T \mathbf{w} + \frac{D}{2} \log \lambda + const \end{aligned}$$

We can get following result:

$$\log P(\epsilon, 0, \log \lambda | data) = \frac{1}{2} + \frac{D}{2} \log \lambda + const$$

By questions 1 and 2 in part 2, we have already known the derivative of the $\mathcal{L}(\mathbf{w}, \epsilon)$ showed as below:

$$\nabla_{\epsilon} \mathcal{L}(\mathbf{w}, \epsilon) = \sum_{n=1}^N \left[\frac{\frac{1}{2} - \sigma(y^{(n)} \mathbf{w}^T x^{(n)})}{(1-\epsilon)\sigma(y^{(n)} \mathbf{w}^T \mathbf{x}^{(n)}) + \epsilon/2} \right]$$

Now, we need to solve the extrema of above function:

$$\begin{aligned} 0 &= \sum_{n=1}^N \left[\frac{\frac{1}{2} - \sigma(y^{(n)} \mathbf{w}^T x^{(n)})}{(1-\epsilon)\sigma(y^{(n)} \mathbf{w}^T \mathbf{x}^{(n)}) + \epsilon/2} \right] \\ &\quad \vdots \\ &\quad \sigma(y \mathbf{w}^T x) = 0 \\ &\quad w = 0 \end{aligned}$$

To sum up, the function $\log \lambda$ is the strictly increasing function. Thus, the log-posterior above is globally maximized by setting the weights to zero and making λ infinite.

3.b) For this question, we should implement the function and make the choice of the parameters. The parameters of the `slice_sample.m` should be decided carefully. For this question, I set all the value of w is 1, and the value of ϵ and $\log \lambda$ are 0.5 and $\log 1$. The width which I set is below 0.2(0.1 for plot the scatter plot). The reason is that, we can get more precise result.

Figure 6 shows the scatter plot of $\log \lambda$ against ϵ . We can assert that, in this question, our posterior beliefs about them are independent. Because we do the simple addition of the two different parts $\mathcal{L}(w, \epsilon)$ and $\frac{D}{2} \log \lambda$. Therefore they do not affect each other.

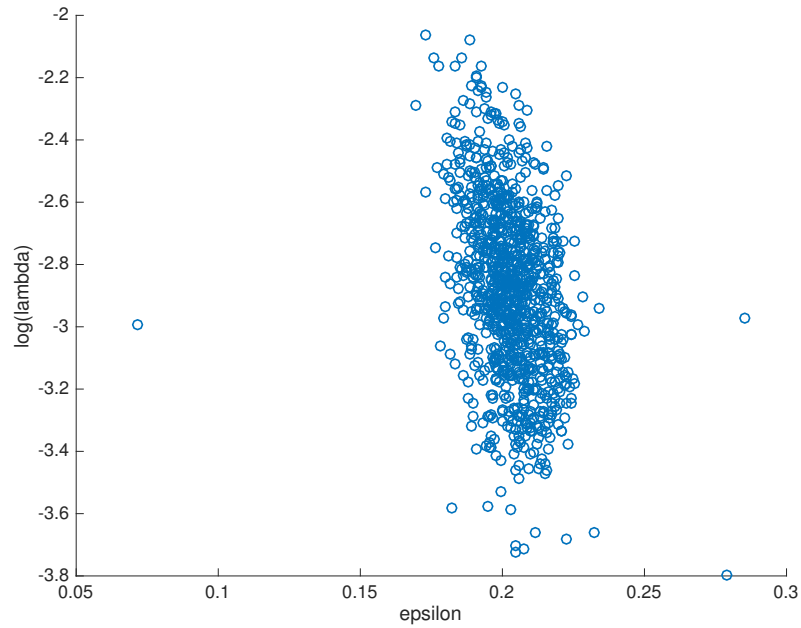


Figure 6: The $\log(\lambda)$ and ϵ

```

1 load('text_data.mat');
2 %a: augment the data
3 x_te = [x_test ones(length(x_test), 1)];
4 x_tr = [x_train ones(length(x_train), 1)];
5
6 samples = slice_sample(1000, 0, @t23f, [ones(101,1)', 0.5, ...
    log(1)'], 0.2, true, x_tr, y_train);
7
8 xlabel('epsilon');
9 ylabel('log(lambda)');
10 hold on;
11 scatter(samples(end-1,:), samples(end,:));

```

```

1 function logP = t23f(w, x, y)
2 y = (y==1)*2 - 1;
3 loglambda = w(end);
4 eps = w(end-1);
5 ww = w(1: end - 2);
6 sigmas = 1./(1 + exp(-y.*(x*ww))); % Nx1
7 Lp = sum(log((1-eps)*sigmas+(eps/2)));
8 lambda = exp(loglambda);
9 logP = Lp - lambda*(ww')*ww + length(w)/2*loglambda ;

```

3.c) Because of the following estimate:

$$\int f(x)P(x) dx \approx \frac{1}{S} \sum_{s=1}^S f(x^{(s)}).$$

Above equation tells us that we should compute all of the prediction by each \vec{w} and do the sum of the results (for this question, we should compute 1000 results of each sample). Finally, we compute the mean of it.

We can compute the data, and get the accuracy 0.9151. So it is better than the previous model. Because we compute the mean of the 1000 samples. And some of the prediction of the samples maybe very good. That's why we got the better result.

Here is the full **matlab** code of 3.

```
1 load('text_data.mat');
2 %a: augment the data
3 x_te = [x_test ones(length(x_test), 1)];
4 x_tr = [x_train ones(length(x_train), 1)];
5
6 samples = slice_sample(1000, 0, @t23f, [ones(101,1)', 0.5, ...
    log(1)]', 0.2, true, x_tr, y_train);
7
8 xlabel('epsilon');
9 ylabel('log(lambda)');
10 hold on;
11 scatter(samples(end-1,:), samples(end,:));
12
13 for i = 1 : 1000
14     y_pred = y_pred + x_te * samples(1:end-2,i);
15 end
16
17 y_prob = mean(y_pred, 2);
18
19 tot = 0;
20 for i = 1 : length(y_prob)
21     if (y_prob(i) > 0.5)
22         temp = 1;
23     else
24         temp = -1;
25     end
26     if (temp == y_test(i))
27         tot = tot + 1;
28     end
29 end
```

```
30  
31 accuracy = tot / length(y_prob);
```

Conclusion

What I have learned Basically, I learned how to use **matlab** to write some simple routines. Futhur more, the topic of comparing the different models and how to set the initial guess data are all very important. The ability to compare the two or more models is what I have grasped.

Where should I make the improvement I have to say that this is the most tough coursework which I have seen. I work hard and spend so much time on it. However, the methods which I have used are not good enough. And some of the principles cannot be undertood deeply by me. Thus, I should do more on the theoritical part. Ann then, I think I can report better. The best practice is that I should study more resource both on papers and examples.