Image Analysis, Assignment 3

## 1&2 Your own classifier & pre-coded machine learning techniques

```
function classification_data = class_train(X, Y)
classification_data = [X;Y];
end


function Y = classify(X,CD)
train_set = CD(1:size(CD,1)-1,:);
diff = zeros(1,size(train_set,2));
for i = 1:size(train_set,2)
    diff(i) = norm(X - train_set(:,i));
end
[value,position] = min(diff);
Y = CD(end,position);
end
```

mean_err_rate_test ✕

1x4 double

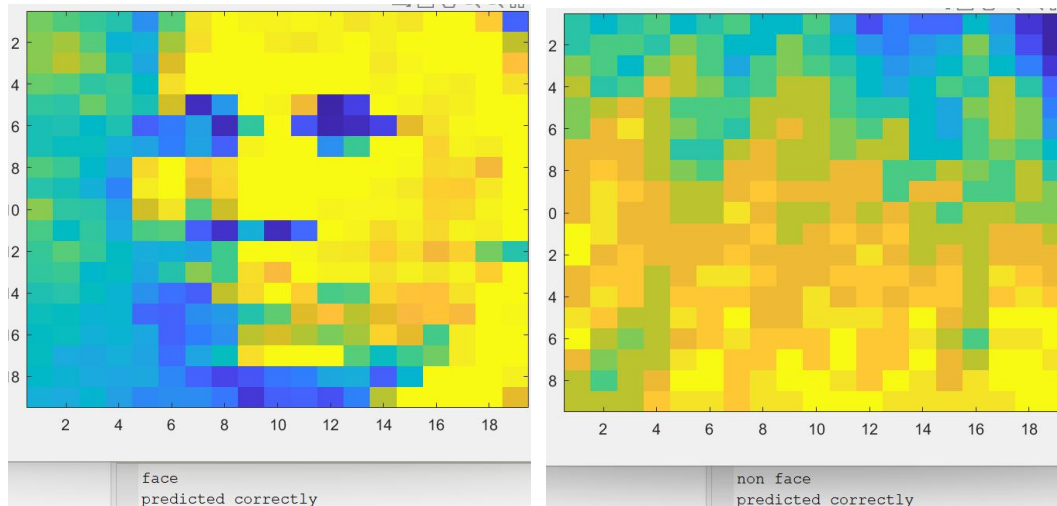| 1 | 2 | 3 | 4 |
|---|---|---|---|
| 0.1575 | 0.1525 | 0.0433 | 0.1575 |

mean_err_rate_train ✕

1x4 double

| 1 | 2 | 3 | 4 |
|---|---|---|---|
| 0 | 0.0115 | 0 | 0 |

1:my own method. 2:regression tree. 3:SVM 4:Nearest neighbor

The error rates of training data is 0, because I use nearest neighbor method, the classify result is just choosing the nearest data in the training data which can find exactly the same picture when we test it on the training data, so the label can not be wrong.

The error rates of test data is higher than training data. And among them, the SVM classifier(3rd) has the lowest error rate, and the Regression tree Classifier and Nearest Neighbor Classifier have almost the same error rates.

My own NN classifier has the same error rate with the NN method by built in function.

face
predicted correctly

non face
predicted correctly

## 3 Testing a simple CNN model



This model is much more complex than the previous models, this is a CNN model contains 3 convolutional layers, 2 max pooling layers and a fully connected layer. And it takes over 5 minutes to run this model, it's a long time compared to the previous methods.

The error rate in test data is lower than regression tree and NN and my own classifier, and a bit higher than the SVM classifier.



## 4 Line fit

```matlab
% Fit a line to these data points with least squares
train_xm = [xm,ones(N,1)];
train_ym = ym;
p_ls = train_xm\train_ym;
plot(x_fine, p_ls(1) * x_fine + p_ls(2))
% Fit a line to these data points using RANSAC.
bestins = 0;
theshold = 1;
for i=1:100
    idx = randperm(N,2);
    k = (ym(idx(1),1)-ym(idx(2),1))/(xm(idx(1),1)-xm(idx(2),1));
    m = ym(idx(1),1)- k*xm(idx(1),1);
    line = [k -1 m];
    distance = abs(line*[xm';ym';ones(1,N)]);
    ins = sum(distance<theshold);
    if ins > bestins
        bestins = ins;
        trueline = line;
    end
end
pa(1) = trueline(1);
pa(2) = trueline(3);
plot(x_fine,pa(1)*x_fine+pa(2),'k--')
```

Need 2 points to sample to estimate the line model in RANSAC.
Least squares method can be affected by noisy points (outliers), and can not fit the line well, and RANSAC approach can refit a line to the inliers, and this line is robust to the outliers that RANSAC identified and ignored.

```matlab
% WRITE CODE BELOW TO COMPUTE THE 4 ERRORS
LS_squared_vertical_errors = 0;
LS_squared_orthogonal_errors = 0;
RANSAC_squared_vertical_errors = 0;
RANSAC_squared_orthogonal_errors = 0;
for i =1:N

    LS_squared_vertical_errors = (p_ls(1) * xm(i) + p_ls(2) - ym(i))^2 +
LS_squared_vertical_errors;
    LS_squared_orthogonal_errors =
(abs(p_ls(1)*xm(i)+p_ls(2)-ym(i))/sqrt(p_ls(1)^2+1))^2 +
LS_squared_orthogonal_errors;
    RANSAC_squared_vertical_errors = (modelInliers(1) * xm(i) + modelInliers(2)
- ym(i))^2 + RANSAC_squared_vertical_errors;
    RANSAC_squared_orthogonal_errors = (abs(modelInliers(1)*xm(i) +
modelInliers(2)-ym(i))/sqrt(modelInliers(1)^2+1))^2 +
RANSAC_squared_orthogonal_errors;
end
LS_squared_vertical_errors
LS_squared_orthogonal_errors
RANSAC_squared_vertical_errors
RANSAC_squared_orthogonal_errors
```

LS_squared_vertical_errors =
  1.4839e+03
LS_squared_orthogonal_errors =
  367.7324
RANSAC_squared_vertical_errors =
  1.9028e+03
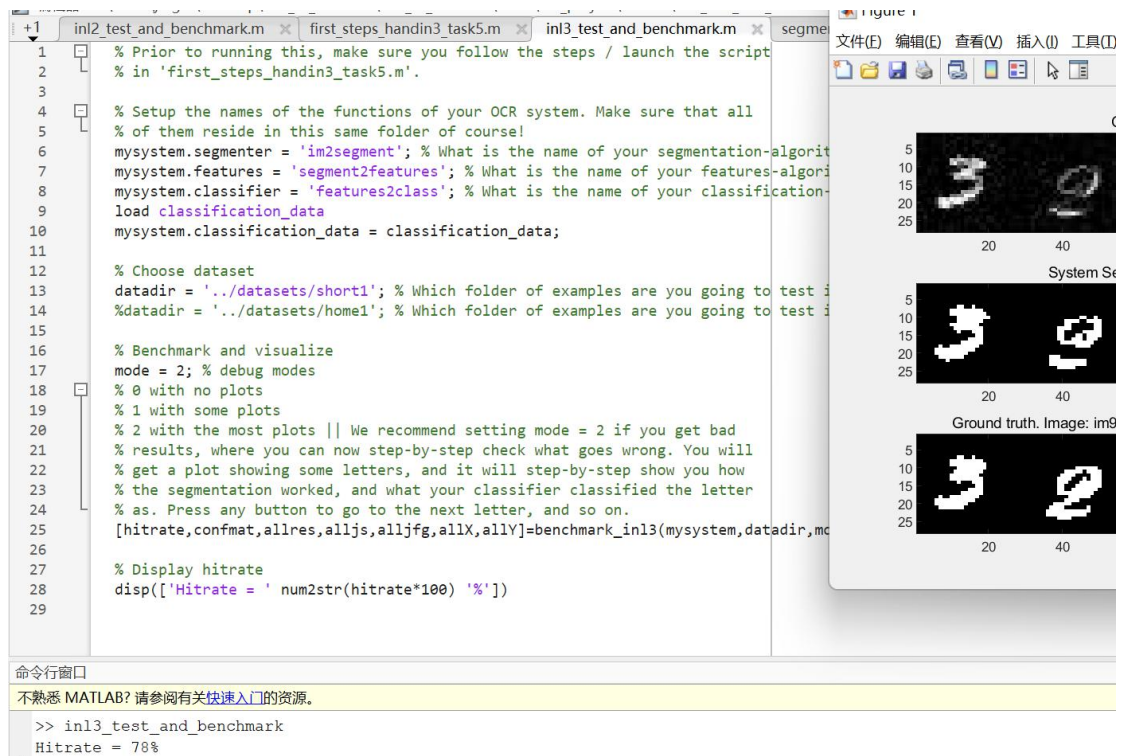RANSAC_squared_orthogonal_errors =
  192.2856

Among the four error, squared_orthogonal_error (TLS) for RANSAC is the lowest.

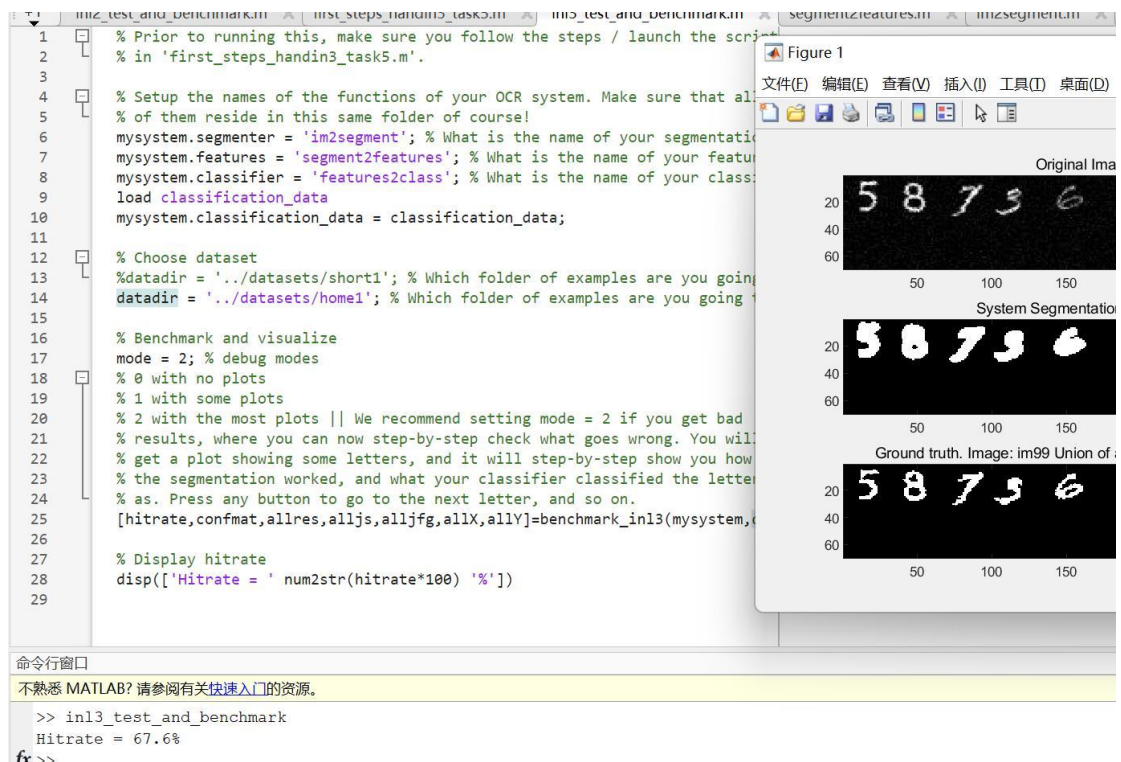The two LS error are much higher than the two TLS error.

So RANSAC is the better approach than LS, because its TLS is lower.

## 5 OCR system construction and system testing

For the short1, Hit rate = 78%

```
1    % Prior to running this, make sure you follow the steps / launch the script
2    % in 'first_steps_handin3_task5.m'.
3
4    % Setup the names of the functions of your OCR system. Make sure that all
5    % of them reside in this same folder of course!
6    mysystem.segmenter = 'im2segment'; % What is the name of your segmentation-algori
7    mysystem.features = 'segment2features'; % What is the name of your features-algori
8    mysystem.classifier = 'features2class'; % What is the name of your classification-
9    load classification_data
10   mysystem.classification_data = classification_data;
11
12   % Choose dataset
13   datadir = '../datasets/short1'; % Which folder of examples are you going to test i
14   %datadir = '../datasets/home1'; % Which folder of examples are you going to test i
15
16   % Benchmark and visualize
17   mode = 2; % debug modes
18   % 0 with no plots
19   % 1 with some plots
20   % 2 with the most plots || We recommend setting mode = 2 if you get bad
21   % results, where you can now step-by-step check what goes wrong. You will
22   % get a plot showing some letters, and it will step-by-step show you how
23   % the segmentation worked, and what your classifier classified the letter
24   % as. Press any button to go to the next letter, and so on.
25   [hitrate,confmat,allres,alljs,alljfg,allX,allY]=benchmark_inl3(mysystem,datadir,mo
26
27   % Display hitrate
28   disp(['Hitrate = ' num2str(hitrate*100) '%'])
29
```

命令行窗口
不熟悉 MATLAB? 请参阅有关 快速入门 的资源。

```
>> inl3_test_and_benchmark
Hitrate = 78%
```

For the home1, Hit rate =67.6 %

```
1    % Prior to running this, make sure you follow the steps / launch the script
2    % in 'first_steps_handin3_task5.m'.
3
4    % Setup the names of the functions of your OCR system. Make sure that all
5    % of them reside in this same folder of course!
6    mysystem.segmenter = 'im2segment'; % What is the name of your segmentatio
7    mysystem.features = 'segment2features'; % What is the name of your featur
8    mysystem.classifier = 'features2class'; % What is the name of your class
9    load classification_data
10   mysystem.classification_data = classification_data;
11
12   % Choose dataset
13   %datadir = '../datasets/short1'; % Which folder of examples are you going
14   datadir = '../datasets/home1'; % Which folder of examples are you going
15
16   % Benchmark and visualize
17   mode = 2; % debug modes
18   % 0 with no plots
19   % 1 with some plots
20   % 2 with the most plots || We recommend setting mode = 2 if you get bad
21   % results, where you can now step-by-step check what goes wrong. You wil
22   % get a plot showing some letters, and it will step-by-step show you how
23   % the segmentation worked, and what your classifier classified the lette
24   % as. Press any button to go to the next letter, and so on.
25   [hitrate,confmat,allres,alljs,alljfg,allX,allY]=benchmark_inl3(mysystem,
26
27   % Display hitrate
28   disp(['Hitrate = ' num2str(hitrate*100) '%'])
29
```

命令行窗口
不熟悉 MATLAB? 请参阅有关 快速入门 的资源。

```
>> inl3_test_and_benchmark
Hitrate = 67.6%
fx >>
```

When I test the short1 data, my segment function works well and give the 78% hit rate.
But when I test the home1 data, there are more noise in the image, so the segment function always recognizes some noise point as a digit, so I have to adjust the gaussian filter's standard deviation from 0.5 to 1.25, and then it can work.

We do not demand a complete report for these assignments, but more text and structure would have gotten you much further.

1) Comment your own classifier specifically, for example why we get exactly 0 for training error.

3) Comment the method, ex model complexity, running times etc.

4) "How is the line model estimated for the minimal dataset in the RANSAC approach?" I.e. how would you calculate the line if you do not use polyfit?

"Atleast2pointstosampletoestimatethelinemodelinRANSAC" No, exactly 2 points.

Please address my comments and resubmit.