

Ans: i

$$p(y_i=1/f(x_i)) = \frac{1}{1+e^{-\beta(x_i)}}$$

$$= \sigma(f(x_i))$$

; $\sigma(x)$ is sigmoid function.

$$p(y_i=-1/f(x_i)) = 1 - \frac{1}{1+e^{-\beta(x_i)}}$$

$$= \frac{e^{-\beta(x_i)}}{1+e^{-\beta(x_i)}}$$

$$= \frac{1}{1+e^{\beta(x_i)}}$$

$$= \sigma(-f(x_i))$$

$$p(y_i/f(x_i)) = \sigma(y_i f(x_i))$$

Assuming data points to be independent and identically distributed (i.i.d):

$$\text{log likelihood, } L = \prod_{i=1}^N p(y_i/f(x_i))$$

$$\text{Negative log likelihood, } -\ln L = - \sum_{i=1}^N \ln[p(y_i/f(x_i))]$$

$$= - \sum_{i=1}^N \ln[\sigma(y_i f(x_i))]$$

$$= - \sum_{i=1}^N \ln \left[\frac{1}{1+e^{-y_i f(x_i)}} \right]$$

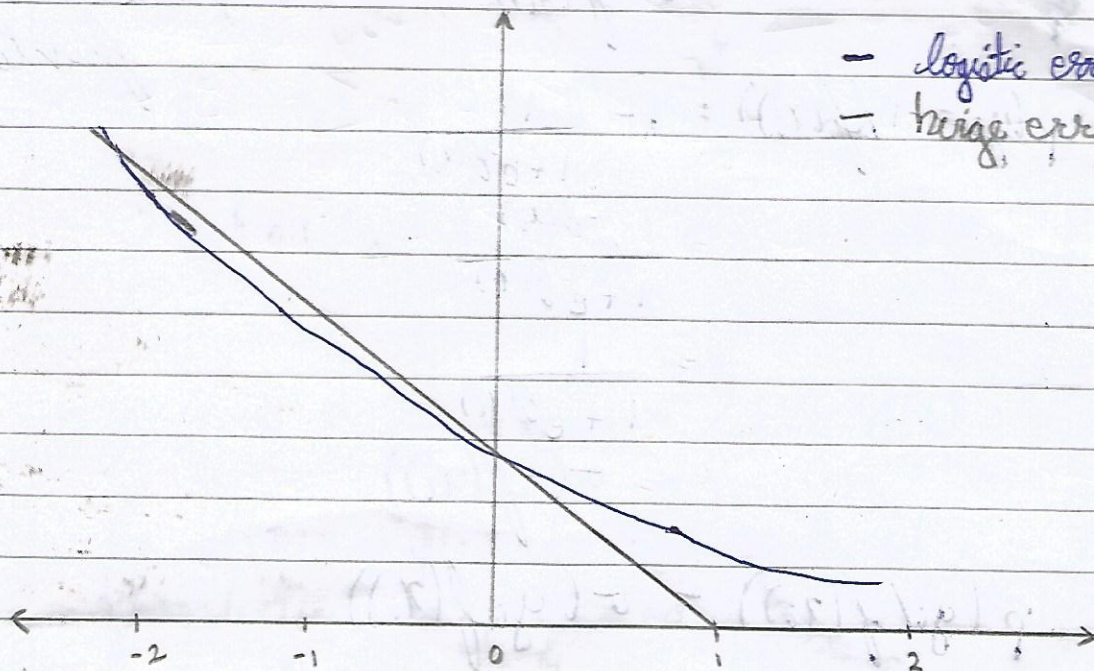
$$= \sum_{i=1}^N \ln[1+e^{-y_i f(x_i)}]$$

$$= \sum_{i=1}^N E_{LR}(y_i, f(x_i))$$

$$\text{where } E_{LR}(y_i, f(x_i)) = \ln(1 + \exp(-y_i f(x_i)))$$

With addition of a quadratic regularization term;

$$\rightarrow \sum_{i=1}^N \text{Eur}(y_i, f(x_i)) + \lambda \|w\|^2$$



- logistic error function
- hinge error function

$$\text{reg Hinge error} = \sum_{n=1}^N \text{Esv}(y_n, t_n) + \lambda \|w\|^2$$

$$\text{where } \text{Esv}(y_n, t_n) = [1 - y_n t_n]_+$$

Both logistic and hinge loss are continuous approximations of misclassification error. Hinge loss is more sparse and covers less area than logistic error function.

~~Ans:~~ Theory Extra Credit

Ans:

$$L(w, b, \epsilon_1, \dots, \epsilon_n, \hat{\epsilon}_1, \dots, \hat{\epsilon}_n) =$$

$$C \sum_{i=1}^N (\epsilon_i + \hat{\epsilon}_i) + \frac{1}{2} \|w\|^2 - \sum_{i=1}^N (\mu_i \epsilon_i + \hat{\mu}_i \hat{\epsilon}_i)$$

$$- \sum_{i=1}^N a_i (\epsilon + \epsilon_i + f(x_i) - y_i) - \sum_{i=1}^N \hat{a}_i (\epsilon + \hat{\epsilon}_i - f(x_i) + y_i) \quad (1)$$

$$\# \Rightarrow \frac{dL}{dw} = w - \sum_{i=1}^N a_i \phi(x) + \sum_{i=1}^N \hat{a}_i \phi(x) = 0$$

Here $f(x) = w^T \phi(x) + b$ was replaced and then differentiation was performed. ; $w = \sum_{i=1}^N \phi(x_i) [a_i - \hat{a}_i]$ (2)

$$\# \frac{dL}{dx} = - \sum_{i=1}^N a_i + \sum_{i=1}^N \hat{a}_i = 0$$

$$\sum_{i=1}^N (\hat{a}_i - a_i) = 0 \quad (3)$$

$$\# \frac{dL}{d\epsilon_i} = C - \sum_{i=1}^N \mu_i - \sum_{i=1}^N a_i = 0$$

$$C = \sum_{i=1}^N (\mu_i + a_i) \quad (4)$$

$$\# \frac{dL}{d\hat{\epsilon}_i} = C - \sum_{i=1}^N \hat{\mu}_i - \sum_{i=1}^N \hat{a}_i = 0$$

$$C = \sum_{i=1}^N (\hat{\mu}_i + \hat{a}_i) \quad (5)$$

Now, combining like terms in 1;

$$\begin{aligned}
 & \rightarrow \sum_{i=1}^N \epsilon_i (C - \mu_i - a_i) + \sum_{i=1}^N \hat{\epsilon}_i [C - \hat{\mu}_i - \hat{a}_i] \\
 & + \epsilon \left[-\sum_{i=1}^N a_i - \sum_{i=1}^N \hat{a}_i \right] + \left[\sum_{i=1}^N a_i y_i - \sum_{i=1}^N \hat{a}_i \hat{y}_i \right] \\
 & + \frac{1}{2} \|\omega\|^2 + \sum_{i=1}^N [\hat{a}_i \omega^T \phi(x_i) - a_i \omega^T \phi(x_i)]
 \end{aligned}$$

Now, from 4, $C = \mu_i + a_i$
 from 5, $C = \hat{\mu}_i + \hat{a}_i$

→ ⑥ and ⑦ becomes 0, where ⑥ $\sum_{i=1}^N \epsilon_i (C - \mu_i - a_i)$
 ⑦ $\sum_{i=1}^N \hat{\epsilon}_i (C - \hat{\mu}_i - \hat{a}_i)$

Thus;

$$\begin{aligned}
 & \rightarrow -\epsilon \sum_{i=1}^N [a_i + \hat{a}_i] + \left[\sum_{i=1}^N (a_i - \hat{a}_i) y_i \right] \\
 & + \frac{1}{2} \|\omega\|^2 + \sum_{i=1}^N [\hat{a}_i \omega^T \phi(x_i) - a_i \omega^T \phi(x_i)]
 \end{aligned}$$

from 2;

$$\|\omega\|^2 = \omega^T \omega$$

~~$$\begin{aligned}
 & \frac{1}{2} \sum_{i=1}^N \phi(x_i) [a_i - \hat{a}_i] + \sum_{j=1}^N \phi(x_j) [\hat{a}_j - a_j] \\
 & + \sum_{i=1}^N
 \end{aligned}$$~~

$$\begin{aligned}
 & - \varepsilon \sum_{i=1}^N (a_i + \hat{a}_i) + \sum_{i=1}^N (a_i - \hat{a}_i) y_i \\
 & + \frac{1}{2} \left[\sum_{i=1}^N \phi(x_i)^T (a_i - \hat{a}_i) \sum_{j=1}^N \phi(x_j) (a_j - \hat{a}_j) \right] \\
 & + \sum_{i=1}^N (\hat{a}_i^T \phi(x_i) - a_i^T \phi(x_i)) \\
 & - \varepsilon \sum_{i=1}^N (a_i + \hat{a}_i) + \sum_{i=1}^N (a_i - \hat{a}_i) y_i \\
 & + \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N k(x_i, x_j) (a_i - \hat{a}_i) (a_j - \hat{a}_j) \\
 & + \sum_{i=1}^N \left[\hat{a}_i^T \sum_{j=1}^N \phi(x_j) (a_j - \hat{a}_j) - a_i^T \sum_{j=1}^N \phi(x_j) (a_j - \hat{a}_j) \right] ; k(x_i, x_j) = \phi(x_i)^T \phi(x_j)
 \end{aligned}$$

$$\begin{aligned}
 & - \varepsilon \sum_{i=1}^N (a_i + \hat{a}_i) + \sum_{i=1}^N (a_i - \hat{a}_i) y_i \\
 & + \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N k(x_i, x_j) (a_i - \hat{a}_i) (a_j - \hat{a}_j) \\
 & - \sum_{i=1}^N \sum_{j=1}^N \phi(x_j)^T (\hat{a}_j - a_j) (\hat{a}_i - a_i)
 \end{aligned}$$

$$\begin{aligned}
 & - \varepsilon \sum_{i=1}^N (a_i + \hat{a}_i) + \sum_{i=1}^N (a_i - \hat{a}_i) y_i \\
 & + \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N k(x_i, x_j) (a_i - \hat{a}_i) (a_j - \hat{a}_j) \\
 & - \sum_{i=1}^N \sum_{j=1}^N k(x_i, x_j) (a_i - \hat{a}_i) (a_j - \hat{a}_j) \\
 & - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N (a_i - \hat{a}_i) (a_j - \hat{a}_j) k(x_i, x_j) \\
 & - \varepsilon \left(\sum_{i=1}^N (a_i + \hat{a}_i) + \sum_{i=1}^N (a_i - \hat{a}_i) y_i \right) ; \text{ Hence proved.}
 \end{aligned}$$

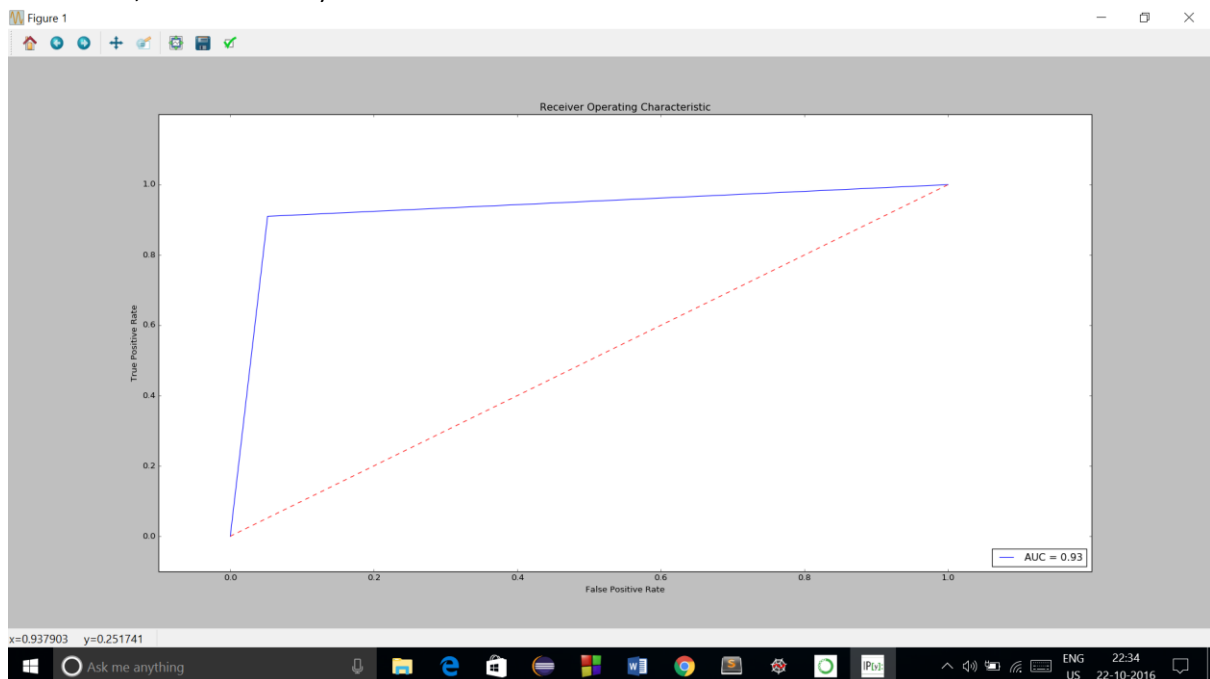
Ans 2).

2.1).

First I have made two list of lists containing image of 28*28 pixels and label belonging to that image. Later I selected those image and labels where label is equal to 3 or 8. 2000 samples per class from the training set belonging to 3 or 8 and a test set of 500 samples per class belonging to 3 or 8 from the testing set were selected. Then soft-margin linear SVM formulation was used and five-fold cross-validation with grid search was performed for finding an appropriate regularization parameter, C. Soft-margin SVM was trained from that C.

- `grid.best_estimator_`

```
SVC(C=0.1, cache_size=200, class_weight=None, coef0=0.0,  
decision_function_shape=None, degree=3, gamma='auto', kernel='linear',  
max_iter=-1, probability=False, random_state=None, shrinking=True,  
tol=0.001, verbose=False)
```



C = 0.1, Accuracy = 0.51

C = 1, Accuracy = 0.5

C = 10, Accuracy = 0.5

C = 100, Accuracy = 0.5

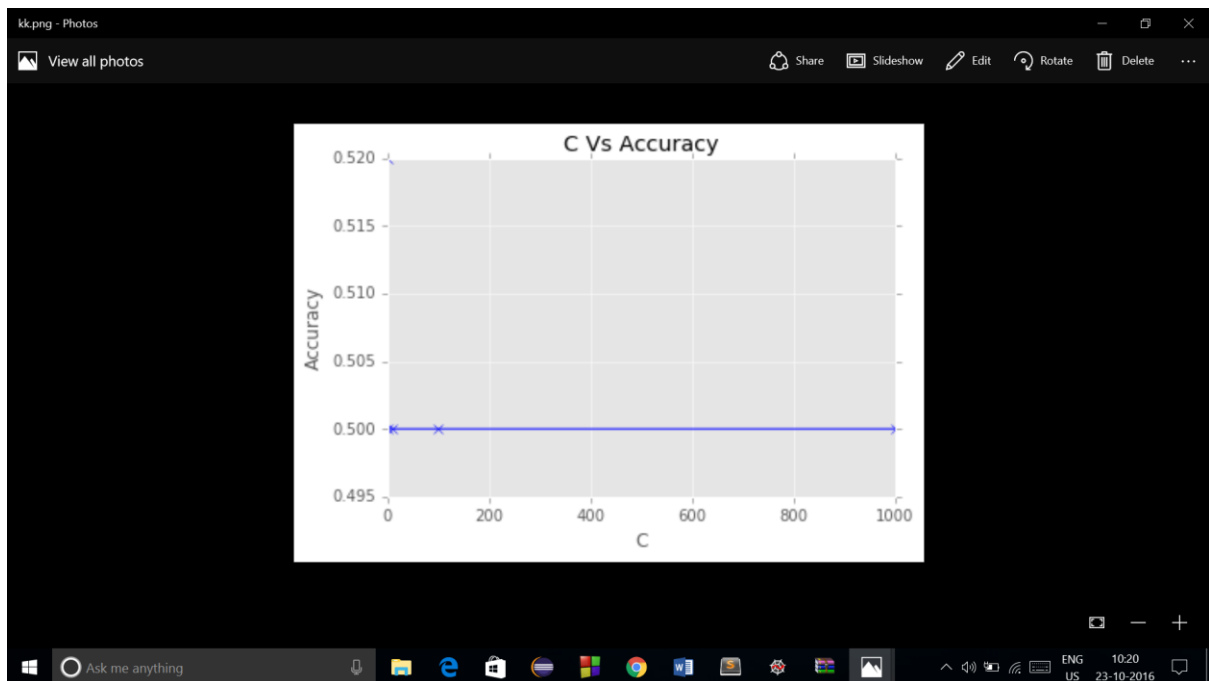
C = 1000, Accuracy = 0.5

```
param_grid = { 'kernel': ['linear'], 'C': [0.1, 1, 10, 100, 1000] }
```

```
svr = SVC()

grid = GridSearchCV(svr, param_grid, cv = 5)

grid.fit(trainimage, trainlabel)
```



2.2):

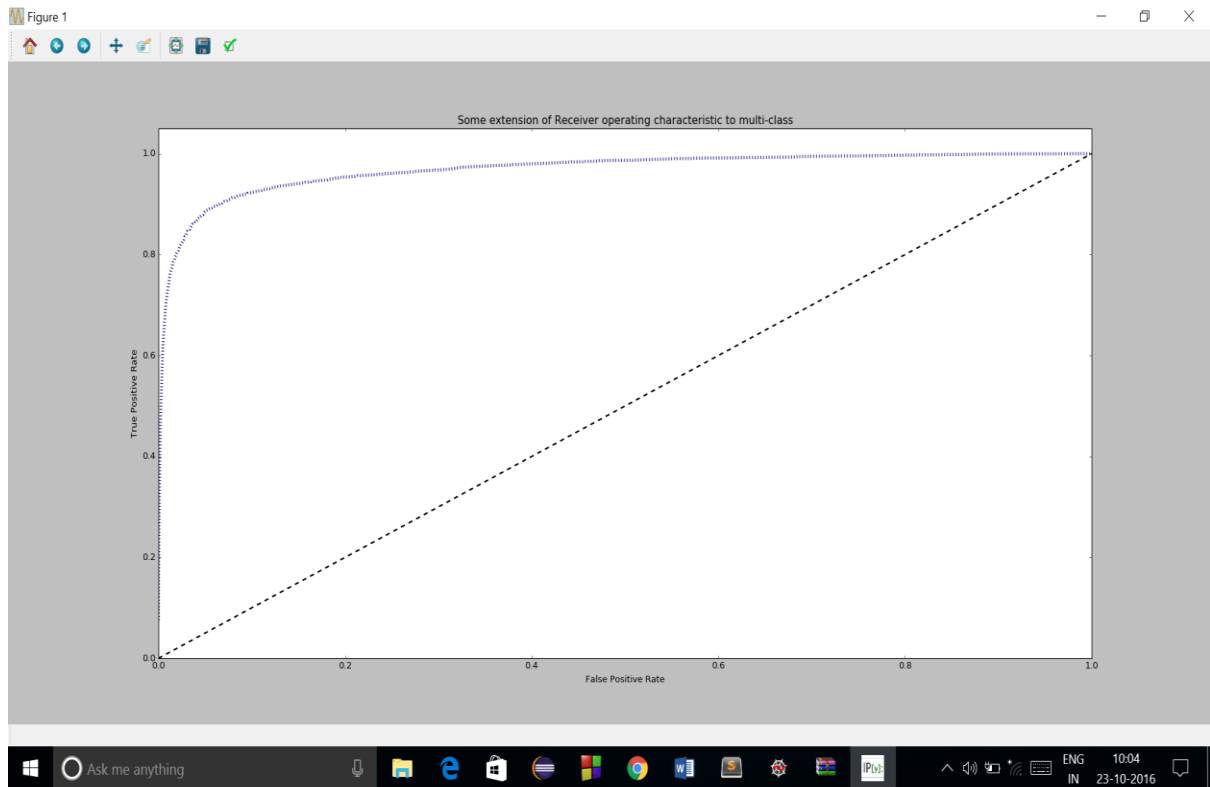
First I have made two list of lists containing image of 28*28 pixels and label belonging to that image. 2000 samples per class from the training set and a test set of 500 samples per class from the testing set were selected. Then soft-margin linear SVM formulation was used and five-fold cross-validation with grid search was performed for finding an appropriate regularization parameter, C. Soft-margin SVM was trained from that C. There was a gridsearch performed to find C.

```
param_grid = {
    'estimator__C': [0.1, 1, 10, 100, 1000]
}

svr = OneVsRestClassifier( LinearSVC() )

grid = GridSearchCV(svr, param_grid, cv = 5)

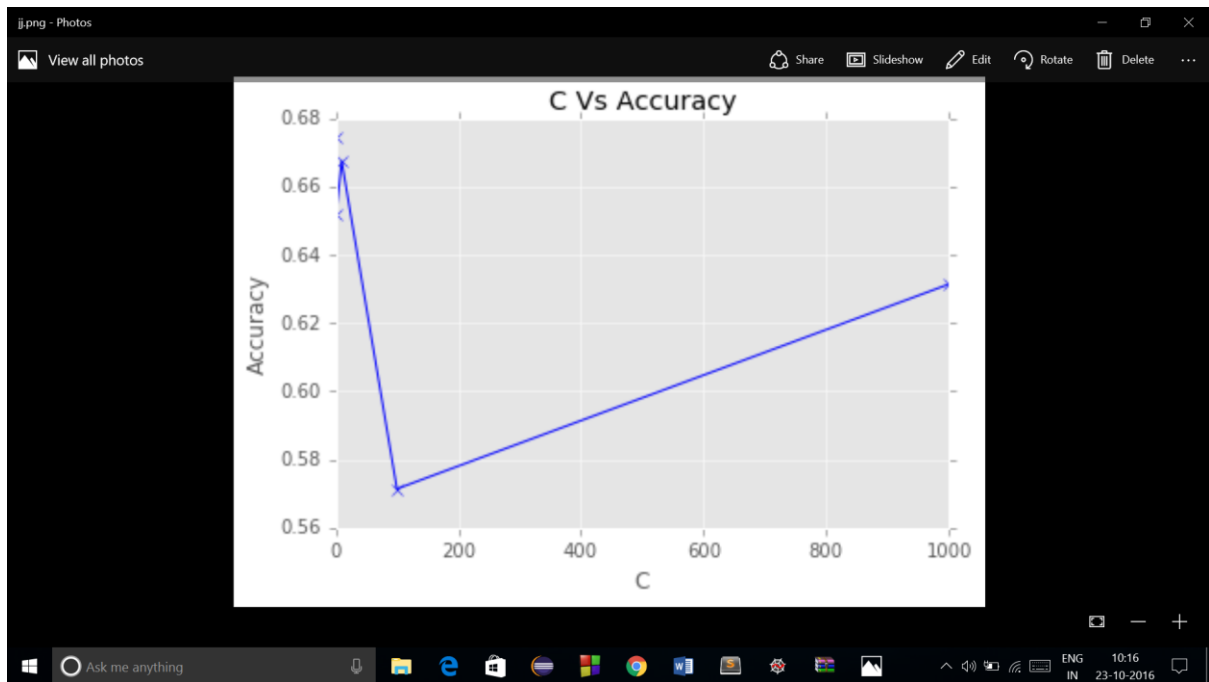
grid.fit(trainimage, trainlabel)
```



```
OneVsRestClassifier(estimator=LinearSVC(C=0.1, class_weight=None, dual=True, fit_intercept=True,  
    intercept_scaling=1, loss='squared_hinge', max_iter=1000,  
    multi_class='ovr', penalty='l2', random_state=None, tol=0.0001,  
    verbose=0),  
    n_jobs=1
```

roc_auc:

0.9679235555555565



C = 0.1, Accuracy = 0.6744

C = 1, Accuracy = 0.6516

C = 10, Accuracy = 0.6674

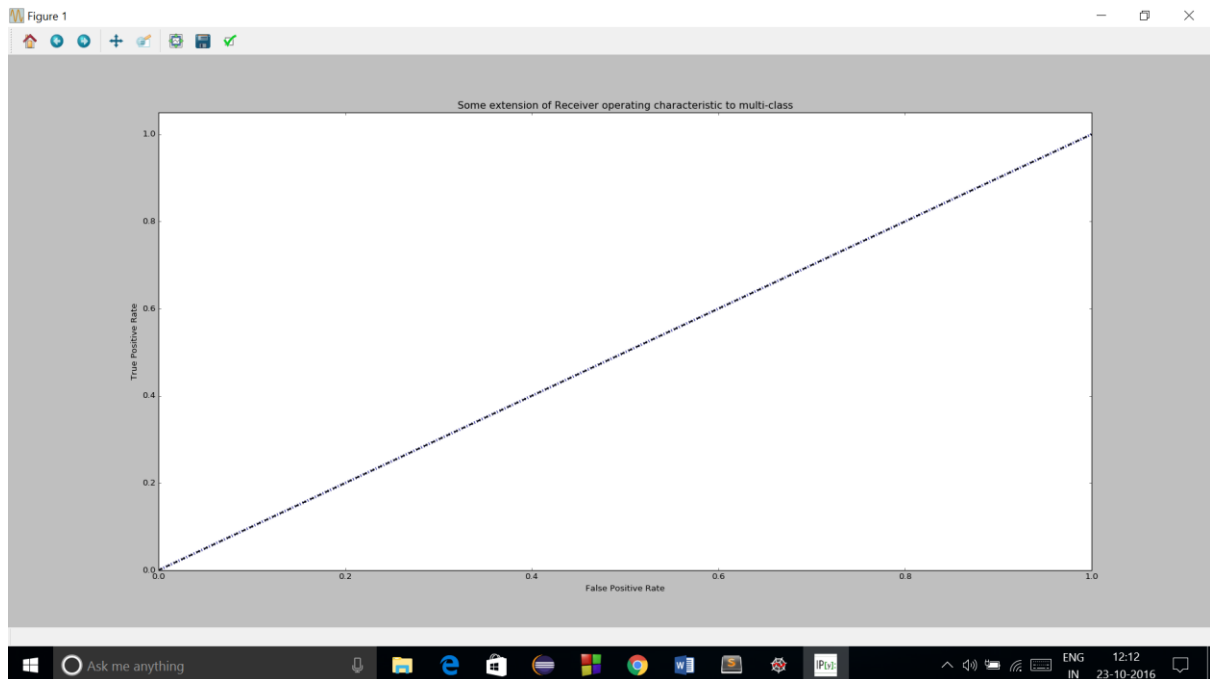
C = 100, Accuracy = 0.5714

C = 1000, Accuracy = 0.6314

Auc:

0: 0.9951613333333323, 1: 0.9949213333333321, 2: 0.9673702222222213, 3: 0.9684706666666659, 4: 0.9884804444444438, 5: 0.9645297777777782, 6: 0.9831453333333343, 7: 0.9783342222222222, 8: 0.8835551111111107, 9: 0.9533431111111114, 'macro': 0.9679235555555565}

2.3).



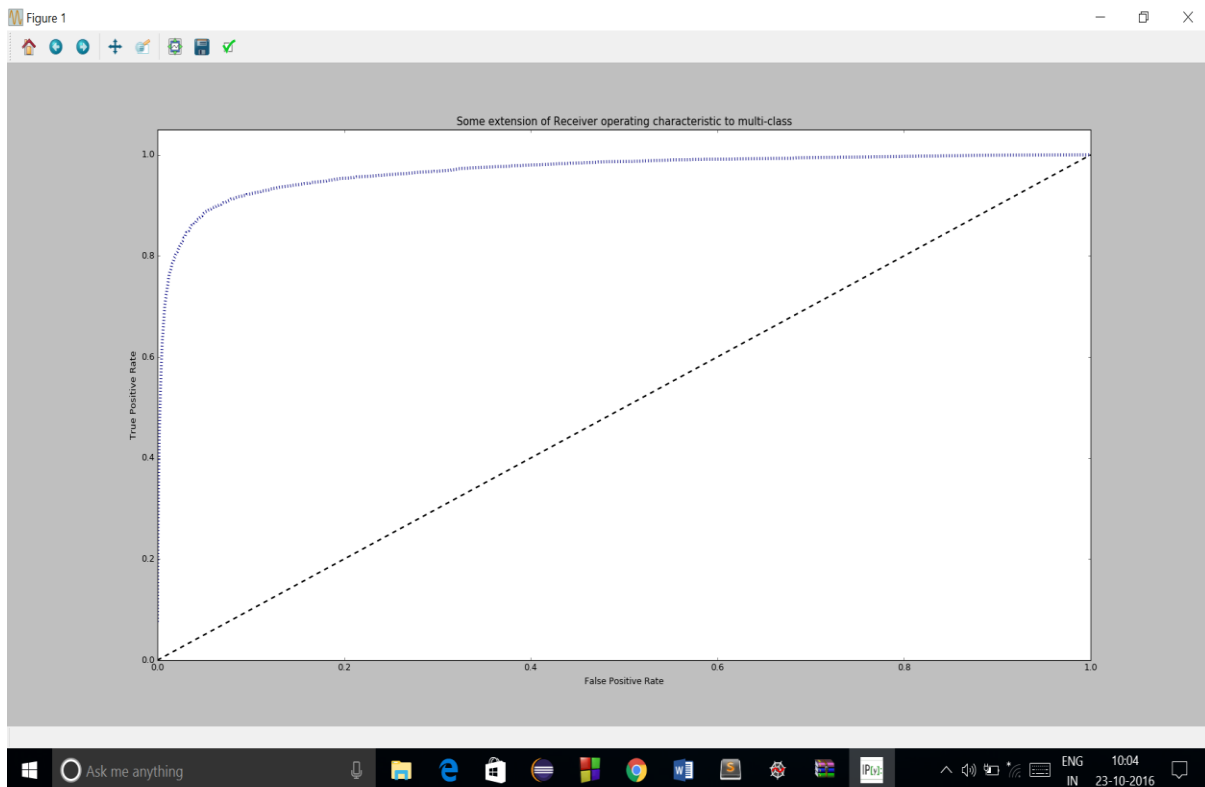
0: 0.5, 1: 0.5, 2: 0.5, 3: 0.5,

4: 0.5, 5: 0.5, 6: 0.5, 7: 0.5, 8: 0.5, 9: 0.5, 'macro': 0.51

$C = 1$, $\text{Gamma} = 0.000001$

I have chosen this C and gamma as others were giving me lower accuracy through `function estimator.score()`.

Accuracy of 2 and 3 together:



Ans 3).

Chose 150 randomly selected samples per class from the training set. Then a test set of 100 samples per class was selected. Then KPCA and KNN was performed with $k = 3$. Five-fold cross-validation was performed with grid search to estimate the gamma. Code is provided in the source section.

gamma: [0.1, 1, 10, 100, 1000], then gamma = 0.1 as the best one.

0.290666666667

Predict_proba returns this:

```
array([[ 0.33333333,  0.33333333,  0.        , ...,  0.        ,
         0.        ,  0.        ],
       [ 0.33333333,  0.33333333,  0.        , ...,  0.        ,
         0.        ,  0.        ],
       [ 0.66666667,  0.33333333,  0.        , ...,  0.        ,
         0.        ,  0.        ],
       ...,
       [ 0.33333333,  0.33333333,  0.        , ...,  0.        ,
         0.        ,  0.33333333],
```


[0. , 0.33333333, 0. , ..., 0. ,
0. , 0.33333333],
[0. , 0.33333333, 0. , ..., 0. ,
0. , 0.33333333]])