# EECS 545 - Homework 2

## Ding Ding

February 22, 2016

1. Support Vector Machine

**1.**

**a) "$\Rightarrow$"**

For
$$\min_{w,b} \tfrac{1}{2}\|w\|^2 + C\sum_{i=1}^{N}\xi_i$$

$$t^{(i)}(w^Tx^{(i)}+b) \geq 1-\xi_i \quad\Rightarrow\quad \xi_i \geq \max\{0,\ 1-t^{(i)}(w^Tx^{(i)}+b)\}$$
$$\xi_i \geq 0$$

$$L = \tfrac{1}{2}\|w\|^2 + C\sum_{i=1}^{N}\xi_i + \sum_{i=1}^{N}\lambda_i\left(\max\{0,1-t^{(i)}(w^Tx^{(i)}+b)\} - \xi_i\right)$$

$$\min_{w,b}\ \max_{\lambda_i}\ L(w,b,\lambda) = \tfrac{1}{2}\|w\|^2 + C\sum_{i=1}^{N}\xi_i + \sum_{i=1}^{N}\lambda_i\left(\max\{0,1-t^{(i)}(w^Tx^{(i)}+b)\} - \xi_i\right)$$

$$\frac{\partial L}{\partial \lambda} = 0 = \max\{0,\ 1-t^{(i)}(w^Tx^{(i)}+b)\} - \xi_i$$

$$\Rightarrow\quad \xi_i = \max\{0,\ 1-t^{(i)}(w^Tx^{(i)}+b)\}$$

$\Rightarrow$ The problem is equivalent to $\quad \min_{w,b} \tfrac{1}{2}\|w\|^2 + C\sum_{i=1}^{N}\max\{0,\ 1-t^{(i)}(w^Tx^{(i)}+b)\}$

**"$\Leftarrow$"**  let $\xi_i = \max(0,\ 1-t^{(i)}(w^Tx^{(i)}+b)) \quad\Rightarrow\quad \begin{cases}\xi_i \geq 0 \\ \xi_i \geq 1-t^{(i)}(w^Tx^{(i)}+b)\end{cases}$

$\therefore$ this unconstricted problem can be writen as

$$\min_{w,b} \tfrac{1}{2}\|w\|^2 + C\sum_{i=1}^{N}\xi_i$$

$$s.t \quad \begin{cases}\xi_i \geq 0 \\ \xi_i \geq 1 - t^{(i)}(w^Tx^{(i)}+b)\end{cases}$$

**(b)** The Hyperplane is $\quad \{x,\ 1-t^{(i)}((w^*)^Tx + b^*) \neq 0\}$.

we have for $x^{(i)}$ in training data st

$$x^{(i)} = x^{(0)} + r\cdot\frac{w}{\|w\|} \qquad \text{for } x^{(0)} \in H \qquad r \neq 0$$

$$\Rightarrow\ 1-t^{(i)}((w^*)^Tx + b^*) = 1- t^{(i)}\left(w^{*T}(x^{(0)} + r\frac{w}{\|w\|}) + b^*\right)$$

$$= 1- t^{(i)}(w^{*T}x^{(0)}+b) - t^{(i)}(w^{*T}r\frac{w}{\|w\|})$$

$$= -t^{(i)}r\|w\|$$

$$\Rightarrow\quad r = \frac{|t^{(i)}((w^*)^Tx+b^*) - 1|}{\|w\|}$$

since $\xi_i > 0$, $\xi \neq 0$,

For $\mathcal{L} = \frac{1}{2}\|w\|^2 + C\sum_{i=1}^{N}\xi_i + \sum_{i=1}^{N}\alpha_i\xi_i + \sum_{i=1}^{N}\beta_i(1-t^{(i)}(w^Tx^{(i)}+b-\xi_i)$

due to KKT condition,

$\xi_i > 0, \quad \alpha_i = 0,$

$\frac{\partial \mathcal{L}}{\partial \xi_i} = C - \alpha_i - \beta_i = 0.$

$\alpha_i + \beta_i = C \implies \beta_i \neq 0, \quad$ due to KKT, $\quad 1 - t^{(i)}(w^Tx^{(i)}+b) = \xi_i$

$\implies r = +\frac{\xi_i^*}{\|w^*\|} \quad$ which is proportional to $\xi_i^*$.

(c). For $t^{(i)}(w^Tx^{(i)}+b) \neq 1$,

we have

$\nabla_w E(w,b) = w - C\sum_{i=1}^{N}t^{(i)}x^{(i)}\mathbb{1}\{t^{(i)}(w^Tx^{(i)}+b)<1\}$

$\frac{\partial}{\partial b}E(w,b) = -C\sum_{i=1}^{N}t^{(i)}\mathbb{1}\{t^{(i)}(w^Tx^{(i)}+b)<1\}$

$\overset{\text{all (i)}}{\vee}$

For $t^{(i)}(w^Tx^{(i)}+b) = 1$,

$\nabla_w^- E(w,b) = \lim_{w\to w^{*-}}\frac{1-w^{(i)}(w^Tx^{(i)}+b)}{w-w^*} = w - C\sum_{i=1}^{N}t^{(i)}x^{(i)} \quad \frac{\partial E(w,b)^-}{\partial b} = \lim_{w\to w^{*-}}\frac{1-w^{(i)}(w^Tx^{(i)}+b)}{b-b^*} = -C\sum_{i=1}^{N}t^{(i)}.$

$\nabla_w^+ E(w,b) = \lim_{w\to w^{*+}}w-\frac{0-0}{w-w^*} = w. \quad \frac{\partial E(w,b)^+}{\partial b} = \lim_{w\to w^{*+}}\frac{1-w^{(i)}(w^Tx^{(i)}+b)}{b-b^*} = 0.$

For $\exists$ some $i$ s.t $t^{(i)}(w^Tx^{(i)}+b)=1$

while other $j \neq i \in [1,N]$ s.t $t^{(j)}(w^Tx^{(j)}+b) \neq 1$,

we have

$\nabla_w^- E(w,b) = w - C\sum_{i=1}^{N}t^{(i)}x^{(i)}\mathbb{1}\{1-t^{(i)}(w^Tx^{(i)}+b)\geq 0\}$

$\nabla_w^+ E(w,b) = w - C\sum_{j\neq i}t^{(j)}x^{(j)}\mathbb{1}\{1-t^{(j)}(w^Tx^{(j)}+b)> 0\}$

$\frac{\partial^- E(w,b)}{\partial b} = -C\sum_{i=1}^{N}t^{(i)}$

$\frac{\partial^+ E(w,b)}{\partial b} = -C\sum_{j\neq i}t^{(j)}.$

(d)
```python
import numpy as np
%matplotlib inline
from matplotlib import pyplot as plt
#import data
traindata=np.loadtxt("digits_training_data.csv",delimiter=",")
trainlabel=np.loadtxt("digits_training_labels.csv",delimiter=",")
#)1 (d)
c=3
ita=0.001
for i in range(trainlabel.size):
    if trainlabel[i]==4:
        trainlabel[i]=1
    else:
        trainlabel[i]=-1
w1=np.zeros(traindata[1,:].size)
b1=0
N1=np.arange(100)
accurate1=np.zeros(100)
m=1
for j in range(1,101):
    wgrab1=w1
    alpha=np.float(ita)/(1+j*(np.float(ita)))
    for k in range(trainlabel.size):
        wgrab1=wgrab1-c*((trainlabel[k]*traindata[k,:]) if trainlabel[k]*(np.dot(traindata[k,:],w1.T
        )+b1)<1 else 0)
        bgrab1=-c*(trainlabel[k] if trainlabel[k]*(np.dot(traindata[k,:],w1.T)+b1)<1 else 0)
    w1=(w1-alpha*wgrab1)
    b1=(b1-alpha*bgrab1)
    a=0
    for k in range(trainlabel.size):
        if trainlabel[k]*(np.dot(traindata[k,:],w1.T)+b1)>=1:
            a=a+1
        else:
            a=a
    accurate1[j-1]=np.float(a)/np.float(trainlabel.size)
plt.savefig('plot1')
```
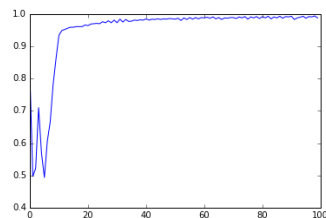
We get the plot as:



Figure 1: Iteration VS Accuracy

(f)
```python
#(f)
w2=np.zeros(traindata[1,:].size)
b2=0
N=np.arange(100)
accurate2=np.zeros(100)
m=1
for j in range(1,101):
    alpha=np.float(ita)/(1+j*(np.float(ita)))
    for k in np.random.permutation(trainlabel.size):
        wgrab2=w2/(trainlabel.size)
        wgrab2=wgrab2-c*((trainlabel[k]*traindata[k,:]) if trainlabel[k]*(np.dot(traindata[k,:],w2.T
)+b2)<1 else 0)
        bgrab2=-c*(trainlabel[k] if trainlabel[k]*(np.dot(traindata[k,:],w2.T)+b2)<1 else 0)
        w2=(w2-alpha*wgrab2)
        b2=(b2-alpha*bgrab2)
    a=0
    for u in range(trainlabel.size):
        if trainlabel[u]*(np.dot(traindata[u,:],w2.T)+b2)>=1:
            a=a+1
        else:
            a=a
    accurate2[j-1]=np.float(a)/np.float(trainlabel.size)
plt.axis([0,100,0.3,1.3])
plt.plot(N,accurate2,label='Batch Gradient Descent')
plt.plot(N,accurate1,label='Stochastic Gradient Descent')
plt.legend()
plt.savefig('plot2')
```
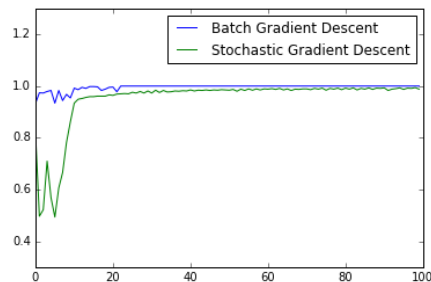
We get the plot as:



Figure 2: Iteration VS Accuracy

(g) We can conclude from the plot that Stochastic Gradient Descent convergent fast than Batch Gradient Descent. Also, we can see that for each iteration, w only update once with Batch Gradient Descent while w update N times in one iteration with Stochastic Gradient Descent.

e). Similarly as c, we have

For $1 - t^{(i)}(w^T x^{(i)} + b) \neq 0$

$$\nabla_w E^{(i)}(w,b) = \frac{1}{N}w - C t^{(i)} x^{(i)} \mathbb{1}_{\{1 - t^{(i)}(w^T x^{(i)} + b) > 0\}}$$

$$\frac{\partial E^{(i)}(w,b)}{\partial b} = -C t^{(i)} \mathbb{1}_{\{1 - t^{(i)}(w^T x^{(i)} + b) > 0\}}$$

For $1 - t^{(i)}(w^T x^{(i)} + b) = 0$,

$$\nabla_w E^{(i)}(w,b)^- = \frac{1}{N}w - C t^{(i)} x^{(i)}$$

$$\nabla_w E^{(i)}(w,b)^+ = \frac{W}{N}$$

$$\frac{\partial E^{(i)}(w,b)^-}{\partial b} = -C t^{(i)}$$

$$\frac{\partial E^{(i)}(w,b)^+}{\partial b} = 0$$

h. $\min \frac{1}{2}\|w\|^2 + C \sum_{i=1}^{n} \xi_i$

s.t $t_i(w^T x_i + b) \geq 1 - \xi_i$

$\xi_i \geq 0.$

$$\Rightarrow \mathcal{L} = \frac{1}{2}\|w\|^2 + C\sum_{i=1}^{n}\xi_i + \sum_{i=1}^{n}\alpha_i(1 - \xi_i - t_i(w^T x_i + b)) - \sum_{i=1}^{n}\beta_i \xi_i$$

$$\begin{cases} \frac{\partial \mathcal{L}}{\partial w} = w - \sum_{i=1}^{n}\alpha_i t_i x_i = 0 \\ \frac{\partial \mathcal{L}}{\partial b} = \sum_{i=1}^{n}\alpha_i t_i = 0 \\ \frac{\partial \mathcal{L}}{\partial \xi_i} = C - \alpha_i - \beta_i = 0. \end{cases} \Rightarrow \begin{cases} w = \sum_{i=1}^{n}\alpha_i t_i x_i \\ \sum_{i=1}^{n}\alpha_i t_i = 0 \\ C = \alpha_i + \beta_i \end{cases}$$

$$\Rightarrow \mathcal{L} = \frac{1}{2}\|\sum_i \alpha_i t_i x_i\|^2 + \sum_{i=1}^{n}(C - \alpha_i - \beta_i)\xi_i + \sum_{i=1}^{n}\alpha_i(1 - t_i((\sum_{i=1}^{n}\alpha_i t_i x_i)^T x_i + b))$$

$$= \frac{1}{2}\sum_{ij}\alpha_i \alpha_j t_i t_j x_i^T x_j + 0 + \sum_{i=1}^{n}\alpha_i + \sum_{i=1}^{n}(-\alpha_i t_i (\sum_{j=1}^{n}\alpha_j t_j x_j)^T x_i) - b\underbrace{\sum_{i=1}^{n}t_i \alpha_i}_{0}$$

$$= \frac{1}{2}\sum_{ij}\alpha_i \alpha_j t_i t_j x_i^T x_j - \sum_{i=1}^{n}(\sum_{j=1}^{n}\alpha_j t_j x_j^T)(\alpha_i t_i x_i) + \sum_{i=1}^{n}\alpha_i$$

$$= -\frac{1}{2}\sum_{ij}\alpha_i \alpha_j t_i t_j x_i^T x_j + \sum_{i=1}^{n}\alpha_i$$

$\Rightarrow$ The dual problem is

$$\mathcal{L}_D = -\frac{1}{2}\sum_{ij}\alpha_i \alpha_j t_i t_j x_i^T x_j + \sum_{i=1}^{n}\alpha_i$$

s.t $\sum \alpha_i t_i = 0$

$0 \leq \alpha_i \leq C$

Since kernel represent the inner product of $X_i, X_j$,

we change $X^{(i)T}X^{(j)}$ into $k(x^{(i)}, x^{(j)})$

$$\Rightarrow L_D(\alpha_i, \beta_i) = -\frac{1}{2}\sum_{i,j=1}^{N}\alpha_i\alpha_j t^{(i)}t^{(j)} k(x^{(i)}, x^{(j)}) + \sum_i \alpha_i$$

$$s.t \quad \sum \alpha_i t_i = 0$$

$$0 \le \alpha_i \le C.$$

(i) First, let C=3,find an appropriate gamma to do soft-margin SVM with rbf kernel.

```
1  import numpy as np
2  %matplotlib inline
3  from matplotlib import pyplot as plt
4  #(i)
5  import sklearn
6  testdata=np.loadtxt("digits_test_data.csv",delimiter=",")
7  testlabel=np.loadtxt("digits_test_labels.csv",delimiter=",")
8  for i in range(testlabel.size):
9      if testlabel[i]==4:
10         testlabel[i]=1
11     else:
12         testlabel[i]=-1
13 gamma=np.logspace(-9.10,13,num=50)
14 accurate2=np.zeros(50)
15 for i in range(0,50):
16     clf=sklearn.svm.SVC(kernel="rbf",gamma=gamma[i],C=3)
17     clf.fit(traindata,trainlabel)
18     pre=clf.predict(testdata)
19     a=0
20     for u in range(testlabel.size):
21         if testlabel[u]==pre[u]:
22             a=a+1
23         else:
24             a=a
25     accurate2[i]=np.float(a)/np.float(testlabel.size)
26 np.argmax(accurate2)
27 gamma[6]
```

we get gamma=4.0375921650671599e-07.

```
1  import matplotlib.cm as cm
2  clf=svm.SVC(kernel="rbf",gamma=4.0375921650671599e-07,C=3)
3  clf.fit(traindata,trainlabel)
4  pret=clf.predict(traindata)
5  a=0
6  for u in range(trainlabel.size):
7      if trainlabel[u]==pret[u]:
8          a=a+1
9      else:
10         a=a
11 accuratet=np.float(a)/np.float(trainlabel.size)
12 pre=clf.predict(testdata)
13 a=0
14 for u in range(testlabel.size):
15     if testlabel[u]==pre[u]:
16         a=a+1
17     else:
18         a=a
19 accurate=np.float(a)/np.float(testlabel.size)
20 print accuratet,accurate
```

We get training accuracy=1 and test accuracy=0.986.

Plot 5 mistakes picture:

```
1  diff=np.abs(testlabel-pre)
2  data = np.genfromtxt("digits_test_data.csv", delimiter=',')
3  N=np.zeros(7)
4  i=0
5  for k in range(0,500):
6      if diff[k]!=0:
7          N[i]=k
8          i=i+1
9      else:
10          i=i
11 plt.imshow(data[N[1]].reshape((26,26)), interpolation="nearest", cmap=cm.Greys_r)
12 plt.savefig('wrong1')
13 plt.imshow(data[N[2]].reshape((26,26)), interpolation="nearest", cmap=cm.Greys_r)
14 plt.savefig('wrong2')
15 plt.imshow(data[N[3]].reshape((26,26)), interpolation="nearest", cmap=cm.Greys_r)
16 plt.savefig('wrong3')
17 plt.imshow(data[N[4]].reshape((26,26)), interpolation="nearest", cmap=cm.Greys_r)
18 plt.savefig('wrong4')
19 plt.imshow(data[N[5]].reshape((26,26)), interpolation="nearest", cmap=cm.Greys_r)
20 plt.savefig('wrong5')
```
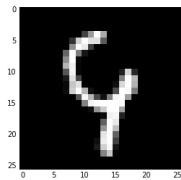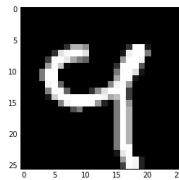


Figure 3: No.122(see 9 as 4)
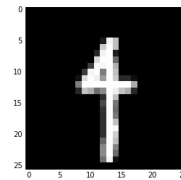


Figure 4: No.163(see 9 as 4)



Figure 5: No.165(see 4 as 9)
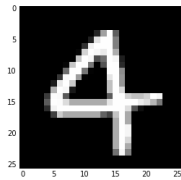


Figure 6: No.209(see 4 as 9)
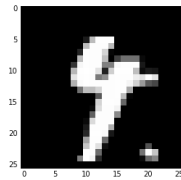


Figure 7: No.329(see 4 as 9)

(j) Apply LDA:

```python
#(j)
import sklearn
testdata=np.loadtxt("digits_test_data.csv",delimiter=",")
testlabel=np.loadtxt("digits_test_labels.csv",delimiter=",")
for i in range(testlabel.size):
    if testlabel[i]==4:
        testlabel[i]=1
    else:
        testlabel[i]=-1
#import data
traindata=np.loadtxt("digits_training_data.csv",delimiter=",")
trainlabel=np.loadtxt("digits_training_labels.csv",delimiter=",")
for i in range(trainlabel.size):
    if trainlabel[i]==4:
        trainlabel[i]=1
    else:
        trainlabel[i]=-1
#get pi
pi1=(trainlabel==1).sum()
pi2=1000-pi1
pi11=np.float(pi1)/np.float(trainlabel.size)
pi22=1-pi11
mu1=np.zeros((traindata[1,:].size,1))
mu2=np.zeros((traindata[1,:].size,1))
sigma1=np.zeros((traindata[:,1].size,traindata[1,:].size))
for i in range(traindata[1,:].size):
    for j in range(trainlabel.size):
        if trainlabel[j]==1:
            mu1[i,0]=mu1[i,0]+(traindata[j,i])
        else:
            mu2[i,0]=mu2[i,0]+(traindata[j,i])
    mu1[i,0]=mu1[i,0]/pi1
    mu2[i,0]=mu2[i,0]/pi2
for j in range(trainlabel.size):
    if trainlabel[j]==1:
        sigma1[j:j+1,:]=traindata[j:j+1,:]-mu1.T
    else:
        sigma1[j:j+1,:]=traindata[j:j+1,:]-mu2.T
sigma=np.dot(sigma1.T,sigma1)/trainlabel.size
gamma1 = -0.5*mu1.T.dot(np.linalg.pinv(sigma)).dot(mu1)+np.log(pi11)
gamma2 = -0.5*mu2.T.dot(np.linalg.pinv(sigma)).dot(mu2)+np.log(pi22)
beta1=np.linalg.pinv(sigma).dot(mu1)
beta2=np.linalg.pinv(sigma).dot(mu2)
y=np.zeros(trainlabel.size)
ytest=np.zeros(testlabel.size)
yita1=np.zeros(trainlabel.size)
yita2=np.zeros(trainlabel.size)
yitat1=np.zeros(testlabel.size)
yitat2=np.zeros(testlabel.size)
for i in range(trainlabel.size):
    yita1[i]=np.exp(np.dot(beta1.T,traindata[i,:])+gamma1)
    yita2[i]=np.exp(np.dot(beta2.T,traindata[i,:])+gamma2)
for i in range(testlabel.size):
```

```
54      yitat1[i]=np.exp(np.dot(beta1.T,testdata[i,:])+gamma1)
55      yitat2[i]=np.exp(np.dot(beta2.T,testdata[i,:])+gamma2)
56  a=0
57  b=0
58  for i in range(trainlabel.size):
59      p1=np.float(yita1[i])/np.float(yita1[i]+yita2[i])
60      p2=np.float(yita2[i])/np.float(yita1[i]+yita2[i])
61      if p1>p2:
62          y[i]=1
63          if y[i]==trainlabel[i]:
64              a=a+1
65          else:
66              a=a
67      else:
68          y[i]=-1
69          if y[i]==trainlabel[i]:
70              a=a+1
71          else:
72              a=a
73  #test data
74  for i in range(testlabel.size):
75      p1=yitat1[i]/np.float(yitat1[i]+yitat2[i])
76      p2=yitat2[i]/np.float(yitat1[i]+yitat2[i])
77      if p1>p2:
78          ytest[i]=1
79          if ytest[i]==testlabel[i]:
80              b=b+1
81          else:
82              b=b
83      else:
84          ytest[i]=-1
85          if ytest[i]==testlabel[i]:
86              b=b+1
87          else:
88              b=b
89  accuratey=np.float(a)/np.float(trainlabel.size)
90  accurateyt=np.float(b)/np.float(testlabel.size)
91  print accuratey
92  print accurateyt
```

We get that training accuray=0.997, test accuracy=0.894;

Plot 5 mistakes picture:

```
1  diff=np.abs(testlabel-ytest)
2  data = np.genfromtxt("digits_test_data.csv", delimiter=',')
3  N=np.zeros(53)
4  i=0
5  for k in range(0,500):
6      if diff[k]!=0:
7          N[i]=k
8          i=i+1
9      else:
10         i=i
11 plt.imshow(data[N[2]].reshape((26,26)), interpolation="nearest", cmap=cm.Greys_r)
```

```
12  plt.savefig('wrong11')
13  plt.imshow(data[N[5]].reshape((26,26)), interpolation="nearest", cmap=cm.Greys_r)
14  plt.savefig('wrong12')
15  plt.imshow(data[N[6]].reshape((26,26)), interpolation="nearest", cmap=cm.Greys_r)
16  plt.savefig('wrong13')
17  plt.imshow(data[N[7]].reshape((26,26)), interpolation="nearest", cmap=cm.Greys_r)
18  plt.savefig('wrong14')
19  plt.imshow(data[N[10]].reshape((26,26)), interpolation="nearest", cmap=cm.Greys_r)
20  plt.savefig('wrong15')
```
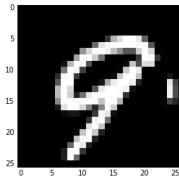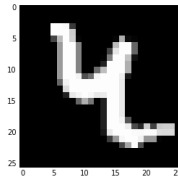
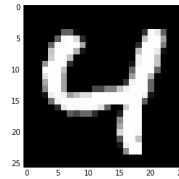Figure 8: No.21(see 9 as 4)          Figure 9: No.34(see 4 as 9)          Figure 10: No.43(see 4 as 9)
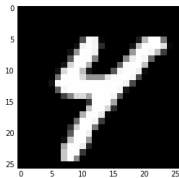
Figure 11: No.50(see 4 as 9)      Figure 12: No.89(see 4 as 9)

There is little different between LDA and SVM.SVM perform better in this case.

2. Open Kaggle Challenge

```
1   import numpy as np
2   %matplotlib inline
3   from matplotlib import pyplot as plt
4   from sklearn import svm
5   from matplotlib import pyplot
6   import matplotlib as mpl
7   trainLabels = np.loadtxt('trainingLabels.gz', dtype=np.uint8, delimiter=',')
8   trainData = np.loadtxt('trainingData.gz', dtype=np.uint8, delimiter=',')
9   testData = np.loadtxt('testData.gz', dtype=np.uint8, delimiter=',')
10  clf=svm.SVC(kernel="rbf",gamma=4.0375921650671599e-07,C=3)
11  clf.fit(trainData,trainLabels)
12  pre=clf.predict(testData)
13  import csv
14  fl = open('test_result.csv', 'w')
15  writer = csv.writer(fl,lineterminator='\n')
16  #writer.writerow(['id', 'Strength']) #if needed
17  for values in pre:
18      writer.writerow([values])
19  fl.close()
```

3. Constructing Kernels

3).

a). 
$$k(u,v) = (\langle u,v \rangle + 1)^4$$
$$= \langle u,v \rangle^4 + 4\langle u,v \rangle^3 + 6\langle u,v \rangle^2 + 4\langle u,v \rangle + 1$$

For $\langle u,v \rangle^4$, since $d=3$, $u=(u_1, u_2, u_3)$, $v=(v_1, v_2, v_3)$

we have

$$\langle u,v \rangle^4 = \left(\sum_{i=1}^{3} u_i v_i\right)^4$$

$$= \sum_{\substack{j_1,j_2,j_3 \\ \sum_k j_k = 4}} \binom{4}{j_1, j_2, j_3}(u_1 v_1)^{j_1}(u_2 v_2)^{j_2}(u_3 v_3)^{j_3}$$

$$\phi_4(u) = \left(\binom{4}{0,0,4}^{1/2} u_3^4, \binom{4}{4,0,0}^{1/2} u_1^4, \binom{4}{0,4,0}^{1/2} u_2^4, \binom{4}{1,0,3}^{1/2} u_3^3 u_1, \binom{4}{1,3,0}^{1/2} u_1 u_2^3, \binom{4}{0,1,3}^{1/2} u_2 u_3^3, \binom{4}{0,3,1}^{1/2} u_2^3 u_3, \binom{4}{3,0,1}^{1/2} u_1^3 u_3, \right.$$
$$\binom{4}{3,1,0}^{1/2} u_1^3 u_2, \binom{4}{2,1,1}^{1/2} u_1^2 u_2 u_3, \binom{4}{1,2,1}^{1/2} u_1 u_2^2 u_3, \binom{4}{1,1,2}^{1/2} u_1 u_2 u_3^2, \binom{4}{2,2,0}^{1/2} u_1^2 u_2^2, \binom{4}{2,0,2}^{1/2} u_1^2 u_3^2, \left.\binom{4}{0,2,2}^{1/2} u_2^2 u_3^2\right)$$

$$= \left(u_1^4, u_2^4, u_3^4, 2u_3^2 u_3, 2u_1^3 u_2, 2u_3^3 u_1, 2u_2^3 u_3, 2u_2^3 u_1, 2u_1^3 u_3, 2\sqrt{3}u_1^2 u_2 u_3, 2\sqrt{3}u_1 u_2^2 u_3, 2\sqrt{3}u_1 u_2 u_3^2, \sqrt{6}u_1^2 u_2^2, \sqrt{6}u_1^2 u_3^2, \sqrt{6}u_2^2 u_3^2\right)$$

$$\langle u,v \rangle^3 = \left(\sum_{i=1}^{3} u_i v_i\right)^3 = \sum_{\substack{j_1,j_2,j_3 \\ \sum_k j_k = 3}} \binom{3}{j_1, j_2, j_3}(u_1 v_1)^{j_1}(u_2 v_2)^{j_2}(u_3 v_3)^{j_3}$$

$$\phi_3(u) = \left(\binom{3}{0,0,3}^{1/2} u_3^3, \binom{3}{0,3,0}^{1/2} u_2^3, \binom{3}{3,0,0}^{1/2} u_1^3, \binom{3}{2,1,0}^{1/2} u_1^2 u_2, \binom{3}{2,0,1}^{1/2} u_1^2 u_3, \binom{3}{0,2,1}^{1/2} u_2^2 u_3, \binom{3}{0,1,2}^{1/2} u_2 u_3^2, \binom{3}{1,0,2}^{1/2} u_1 u_3^2, \right.$$
$$\left.\binom{3}{1,2,0}^{1/2} u_1 u_2^2, \binom{3}{1,1,1}^{1/2} u_1 u_2 u_3\right)$$

$$= \left(u_1^3, u_2^3, u_3^3, \sqrt{3}u_1^2 u_2, \sqrt{3}u_1^2 u_3, \sqrt{3}u_2^2 u_3, \sqrt{3}u_2^2 u_1, \sqrt{3}u_3^2 u_1, \sqrt{3}u_3^2 u_2, \sqrt{6}u_1 u_2 u_3\right)$$

$$\langle u,v \rangle^2 = \left(\sum_{i=1}^{3} u_i v_i\right)^2$$

$$= \sum_{\substack{j_1,j_2,j_3 \\ \sum_k j_k = 2}} \binom{2}{j_1, j_2, j_3}(u_1 v_1)^{j_1}(u_2 v_2)^{j_2}(u_3 v_3)^{j_3}$$

$$\phi_2(u) = \left(\binom{2}{2,0,0}^{1/2} u_1^2, \binom{2}{0,2,0}^{1/2} u_2^2, \binom{2}{0,0,2}^{1/2} u_3^2, \binom{2}{1,1,0} u_1 u_2, \binom{2}{1,0,1} u_1 u_3, \binom{2}{0,1,1} u_2 u_3\right)$$

$$= \left(u_1^2, u_2^2, u_3^2, \sqrt{2}u_1 u_2, \sqrt{2}u_2 u_3, \sqrt{2}u_1 u_3\right)$$

$$\langle u,v \rangle = \sum_{i=1}^{3} u_i v_i$$

$$\phi_1(u) = (u_1, u_2, u_3)$$

Then, we have $\phi(u) = (\phi_4(u), 2\phi_3(u), \sqrt{6}\phi_2(u), 2\phi_1(u), 1)$

$$\Phi(v) = (\phi_4(v), \phi_3(v), \phi_2(v), \phi_1(v), 1)$$

which satisfy $k(u,v) = \phi(u)^T \phi(v)$

For any $d$ dimension $u$, & $v$, we can still get $\phi_4(u), \phi_3(u), \phi_2(u), \phi_1(u), 1$.

$$\langle u,v\rangle^4 = \left(\sum_{i=1}^{d} u_i v_i\right)^4 = \sum_{\substack{j_1\cdots j_d \\ \Sigma_k j_k=4}} \binom{4}{j_1\cdots j_d}(u_1 v_1)^{j_1}\cdots(u_d v_d)^{j_d}$$

$$\phi_4(u) = \left(\cdots, \binom{4}{j_1\cdots j_d}^{1/2} u_1^{j_1}\cdots u_d^{j_d}, \cdots\right)_{\Sigma_k j_k=4}$$

similarly

$$\phi_3(u) = \left(\cdots, \binom{3}{j_1\cdots j_d}^{1/2} u_1^{j_1}\cdots u_d^{j_d}, \cdots\right)_{\Sigma_k j_k=3}$$

$$\phi_2(u) = \left(\cdots, \binom{2}{j_1\cdots j_d}^{1/2} u_1^{j_1}\cdots u_d^{j_d}, \cdots\right)_{\Sigma_k j_k=2}$$

$$\phi_1(u) = (u_1, \cdots\cdots, u_d)$$

Then we have $\phi(u) = (\phi_4(u), \phi_3(u), \phi_2(u), \phi_1(u), 1)$

$k(u,v) = \phi(u)^T \phi(v)$.

b). i). we already know that $k_1$, $k_2$ are positive-definite kernel functions.

∴ $\forall x \in \mathbb{R}^d, \neq 0$, we have for it's gram Matrix $K_1^M, K_2^M$, $x^T(k_1^M)x>0$, $x^T(k_2^M)x>0$.

$x^T(k^M)x = x^T(k_1^M+k_2^M)x > 0$. $\Rightarrow$ the gram Matrix of $k$, is positive-def

Set $k_1(x,z) = \phi_1(x)\phi_1(z)$

$k_2(x,z) = \phi_2^T(x)\phi_2(z)$

$k = \phi_1^T(x)\phi_1(z) + \phi_2^T(x)\phi_2(z) = \langle(\phi_1(x), \phi_2(x)), (\phi_1(z), \phi_2(z))\rangle$.

is also positive-definite kernel functions

ii). Counterexample:

simply let $k_1(x,z) = k_2(x,z)$, we have

$k = k_1(x,z) - k_2(x,z) = 0$

∴ the gram matrix of $k$ is a $D\times D$ matrix with all value be 0 which is not a positive-definite gram matrix.

$\Rightarrow$ $k$ is not a positive definite kernel function.

iii). since $a \in \mathbb{R}^+$, $x^T(k_1^M)x > 0$ for the gram matrix of $k_1$

we have $x^T k^M x = x^T a k^M x = a(x^T k^M x) > 0$, which is also positive-definite matrix.

$k = \langle\sqrt{a}\phi_1(x), \sqrt{a}\phi_1(z)\rangle$ is a positive-definite kernel function.

(iv). $k(x,z) = k_1(x,z) k_2(x,z)$

$$= \phi_1^T(x) \phi_1(z) \phi_2^T(x) \phi_2(z) = \sum_{i=1}^{D} \phi_{1i}(x) \phi_{1i}(z) \sum_{i=1}^{D} \phi_{2i}(x) \phi_{2i}(z) = \sum_{i,j=1}^{D} \phi_{1i}(x) \phi_{1i}(z) \phi_{2j}(x) \phi_{2j}(z)$$

$$= \sum_{i,j=1}^{D} \phi_{1i}(x) \phi_{2j}(x) \phi_{1i}(z) \phi_{2j}(z)$$

$$\therefore k(x,z) = \langle (\phi_{11}(x)\phi_{21}(x), \phi_{11}(x)\phi_{22}(x), \cdots \phi_{11}(x)\phi_{2D}(x), \cdots, \phi_{1D}(x)\phi_{21}(x), \cdots, \phi_{1D}(x)\phi_{2D}(x)),$$

$$(\phi_{11}(z)\phi_{21}(z), \phi_{11}(z)\phi_{22}(z) \cdots, \phi_{11}(x)\phi_{2D}(x), \cdots, \phi_{1D}(z)\phi_{21}(z), \cdots, \phi_{1D}(z)\phi_{2D}(z)) \rangle$$

which is a positive-definite kernel function.

(V). Counterexample:

set $f(x): \mathbb{R}^D \to 0$, which means for $\forall x \in \mathbb{R}^D$, we have $f(x)=0$.

$\Rightarrow$. $k(x,z) = f(x)f(z) = 0$ for $\forall x, z$,

the gram matrix of $k$ is a matrix with all value $=0$. which is not positive-definite matrix.

$\Rightarrow$ $k$ is not a positive-definite kernel function at this time.

(Vi). $k(x,z) = P(k_1(x,z)) = \sum_{i=1}^{N} a_i (k_1(x,z))^i$ where $a_i > 0$

we only need to prove $k_1(x,z)^i$ is a positive-definit kernel, then due from i) & iii),

$k(x,z)$ is also a positive-definit kernel function.

as we proved in (i,V), $k(x,z) = k_1(x,z) k_2(x,z)$ is also a postive-definit kernel function.

$k(x,z) = (k_1(x,z))^2$ also satisfy if let $k_2 = k_1$

similarly, $k(x,z) = (k_1(x,z))^i$ is positive-def for $\forall i \geq 1$

$\Rightarrow$. $k(x,z) = P(k_1(x,z))$ is positive-definit kernel function.

(vii). $k(x,z) = \exp(-\frac{\|x-z\|^2}{2\sigma^2}) = \exp(-\frac{x^Tx}{2\sigma^2}) \exp(-\frac{x^Tz}{\sigma^2}) \exp(-\frac{z^Tz}{2\sigma^2})$

it can be write as $k(x,z) = f(x) \exp(-\frac{x^Tz}{\sigma^2}) f(z)$ where $f(x) = \exp(-\frac{x^Tx}{2\sigma^2})$, $f(z) = \exp(-\frac{z^Tz}{2\sigma^2})$

$\exp(-\frac{x^Tz}{\sigma^2}) = \sum_{n=1}^{\infty} \frac{(\frac{x^Tz}{\sigma^2})^n}{n!}$ by Taylor expansion.

set $k_0(x,z) = x^Tz = \langle x, z \rangle$ which is a kernel function.

due from (iii) $\frac{k_0(x,z)}{\sigma^2}$ is a kernel function, due from (vi), $\sum_{n=1}^{\infty} \frac{(\frac{x^Tz}{\sigma^2})^n}{n!}$ is also a kernel function,

with infinit term, we have $\sum_{n=1}^{\infty} \frac{(\frac{x^Tz}{\sigma^2})^n}{n!}$ can be write as $\langle \phi(x), \phi(z) \rangle$ where $\phi(x), \phi(z)$ are infinite dimension

$\therefore k(x,z) = f(x) \langle \phi(x), \phi(z) \rangle f(z) = \langle f(x)\phi(x), f(z)\phi(z) \rangle = \langle \phi'(x), \phi'(z) \rangle$.    Q.E.D

**4)**

a).

$$P^{-1} - P^{-1}Q(R^{-1} + SP^{-1}Q)^{-1}SP^{-1} \quad (\ast)$$

$$= P^{-1}Q(R^{-1} + SP^{-1}Q)^{-1}(R^{-1} + SP^{-1}Q)Q^{-1} - P^{-1}Q(R^{-1} + SP^{-1}Q)^{-1}SP^{-1}$$

$$= P^{-1}Q(R^{-1} + SP^{-1}Q)^{-1}(R^{-1} + SP^{-1}Q - SP^{-1}Q)Q^{-1}$$

$$= P^{-1}Q(R^{-1} + SP^{-1}Q)^{-1}R^{-1}Q^{-1}$$

Since $(\ast) = (P + QRS)^{-1} = P^{-1}Q(R^{-1} + SP^{-1}Q)^{-1}R^{-1}Q^{-1}$

with $\hat{w} = (\Phi^T\Phi + \lambda I)^{-1}\Phi^T t$

set $P = \lambda I$, $Q = \Phi^T$, $R = I$ $S = \Phi$, we have

$$(\lambda I + \Phi^T\Phi)^{-1} = (\lambda I)^{-1}\Phi^T((I)^{-1} + \Phi(\lambda I)^{-1}\Phi^T)^{-1}(I)^{-1}(\Phi^T)^{-1}$$

$$(\lambda I + \Phi^T\Phi)^{-1} = \frac{1}{\lambda}\Phi^T(I + \frac{1}{\lambda}\Phi\Phi^T)^{-1}(\Phi^T)^{-1}$$

$$(\lambda I + \Phi^T\Phi)^{-1}\Phi^T = \Phi^T(\lambda I + \Phi\Phi^T)^{-1}$$

$$\Rightarrow \hat{w} = (\Phi^T\Phi + \lambda I)^{-1}\Phi^T t = \Phi^T(\lambda I + \Phi\Phi^T)^{-1}t = \Phi^T a \quad \text{where } a = (\lambda I + \Phi\Phi^T)^{-1}t$$

$$\hat{f}(x) = w^T\phi(x) = (\Phi^T a)^T\phi(x) = a^T\Phi\phi(x) = a^T k(x) \quad \text{where } k(x) = \Phi\phi(x) = [k(x,x), \cdots, k(x_1,x)]^T$$

$$E(w) = (\Phi w - t)^T(\Phi w - t) + \lambda w^T w$$

$$= w^T\Phi^T\Phi w - 2t^T\Phi w + t^T t + \lambda w^T w$$

$$= a^T\Phi\Phi^T\Phi\Phi^T a - 2t^T\Phi\Phi^T a + t^T t + \lambda a^T\Phi\Phi^T a$$

$$= a^T KKa - 2t^T ka + t^T t + \lambda a^T ka$$

4. Kernelized Ridge Regression

(b)
```
import numpy as np
#import data
data=np.loadtxt("steel_composition_train.csv",delimiter=",", skiprows=1,usecols=(1,2,3,4,5,6,7,8,9))
X=data[:,:-1] #X is a feature set
Y=data[:,-1] #Y is an array with target value
#normalized
X=(X-np.mean(X,0))/np.std(X,0)
#get y with kernel
k0=np.dot(X,X.T)
I=np.identity(k0[:,1].size)
E=np.zeros(4)
for i in range(3):
    k=np.power((k0+1),i+2)
    a=np.dot(np.linalg.inv((I+k)),Y)
    Ypre=np.dot(a.T,k)
    E[i]=np.dot(Ypre,Ypre.T)-2*np.dot(Y.T,Ypre.T)+np.dot(Y.T,Y)+np.dot(Ypre,a)
k4=np.zeros((k0[:,1].size,k0[:,1].size))
for i in range(k0[:,1].size):
    for j in range(k0[:,1].size):
        k4[i,j]=np.exp(-np.power(np.linalg.norm(X[i,:]-X[j,:]),2)/2)
a=np.dot(np.linalg.inv((I+k4)),Y)
Ypre=np.dot(a.T,k4)
E[3]=np.dot(Ypre,Ypre.T)-2*np.dot(Y.T,Ypre.T)+np.dot(Y.T,Y)+np.dot(Ypre,a)
RMSE=np.sqrt(E/Y.size)
```

We get RMSE with four different kernel be:

| kernel | poly with degree 2 | poly with degree 3 | poly with degree 4 | Gaussian |
|--------|--------------------|--------------------|--------------------|----------|
| RMSE   | 7.41955281         | 4.83144123         | 2.96214783         | 13.18899638 |