# HTML\_2023\_HW4

tags: Personal

Discuss with anonymous and TAs.
Ref: [TA Poy HTML 2023 HW4](https://hackmd.io/@Poy/SJVAvg0Qn)

# More about Regularization

# **P1**

In order to find the minimum w, we can set the derivative of the optimal solution to be zero, hence, we have

$$0 = \frac{\partial}{\partial w} \left(\frac{1}{N} \sum_{n=1}^{N} (wx_n - y_n)^2 + \frac{\lambda}{N} w^2\right)$$

$$= \frac{2}{N} \sum_{n=1}^{N} x_n (wx_n - y_n) + \frac{2w\lambda}{N}$$

$$= \frac{2w}{N} \sum_{n=1}^{N} x_n^2 - \frac{2}{N} \sum_{n=1}^{N} x_n y_n + \frac{2w\lambda}{N}$$

$$= \frac{2w}{N} \left(\sum_{n=1}^{N} x_n^2 + \lambda\right) - \frac{2}{N} \sum_{n=1}^{N} x_n y_n$$

$$\frac{2w}{N} \left(\sum_{n=1}^{N} x_n^2 + \lambda\right) = \sum_{n=1}^{N} \sum_{n=1}^{N} x_n y_n$$

$$w(\sum_{n=1}^{N} x_n^2 + \lambda) = \sum_{n=1}^{N} x_n y_n$$

$$w^* = \frac{\sum_{n=1}^{N} x_n y_n}{\sum_{n=1}^{N} x_n^2 + \lambda}$$
or we can use  $w^* = (x^T x + \lambda I)^{-1} x^T y$ 

$$= \left(\sum_{n=1}^{N} x_n^2 + \lambda\right)^{-1} \cdot \sum_{n=1}^{N} x_n y_n$$

$$= \frac{\sum_{n=1}^{N} x_n y_n}{\sum_{n=1}^{N} x_n^2 + \lambda}$$

$$\therefore C = (w^*)^2$$

$$\therefore C = \left(\frac{\sum_{n=1}^{N} x_n y_n}{\sum_{n=1}^{N} x_n y_n}\right)^2$$

As a result, we should choose  $\left[a\right]$  as our solution.

# **P2**

Since the L2-regularized linear regression in the  $\mathcal{Z}$ -space is equivalent to regularized linear regression in the  $\mathcal{X}$ -space, we have

$$egin{aligned} \min_{ ilde{w} \in \mathbb{R}^{d+1}} rac{1}{N} \sum_{n=1}^N ( ilde{w}^T \Phi(x_n) - y_n)^2 + rac{\lambda}{N} ( ilde{w}^T ilde{w}) \ &= \min_{ ilde{w} \in \mathbb{R}^{d+1}} rac{1}{N} \sum_{n=1}^N (w^T x_n - y_n)^2 + rac{\lambda}{N} \Omega(w) \ &\therefore ilde{w}^T \Phi(x_n) = ilde{w} \Gamma^{-1} (x-u) = ilde{w} \Gamma^{-1} x = w^T x_n, ilde{w}^T ilde{w} = \Omega(w) \end{aligned}$$

Since  $\Gamma$  is a diagonal matrix, we know that the inverse of diagonal matrix equals to itself, which means  $\Gamma^T=\Gamma$  then we have

$$egin{aligned} w^T &= ilde w \Gamma^{-1} \ ilde w^T &= w^T \Gamma \ ilde w &= (w^T \Gamma)^T = \Gamma^T w \ dots &: \Omega(w) &= ilde w^T ilde w &= w^T \Gamma \Gamma w = w^T \Gamma^2 w \end{aligned}$$

As a result, we should choose  $\left[b
ight]$  as our solution.

**P3** 

$$\begin{split} \min_{w} \frac{1}{N} \sum_{n=1}^{N} \operatorname{err}_{\operatorname{smooth}}(w, x_{n}, y_{n}) &= \min_{w} \frac{1}{N} \sum_{n=1}^{N} \operatorname{err}(w, x_{n}, y_{n}) + \frac{\lambda}{N} \sum_{n=1}^{N} \Omega(w, x_{n}) \\ D_{KL}(P \| Q) &= \sum_{x} P(x) \ln \frac{P(x)}{Q(x)} = \sum_{x} \frac{1}{2} \ln \frac{\frac{1}{2}}{\operatorname{err}(w, x, y)} \\ w_{lr} &= \operatorname{argmin} \sum_{w}^{n} D_{KL}(P(Y_{i}) \| P(\hat{Y}_{i}) \\ &= \operatorname{argmin} \sum_{w}^{n} y_{i} \ln \frac{y_{i}}{p_{i}} + (1 - y_{i}) \ln \frac{(1 - y_{i})}{(1 - p_{i})} \\ &= \operatorname{argmin} \sum_{i=1}^{n} y_{i} (\ln y_{i} - \ln p_{i}) + (1 - y_{i}) (\ln (1 - y_{i}) - \ln (1 - p_{i})) \\ &= \operatorname{argmin} \sum_{w}^{n} -y_{i} \ln p_{i} - (1 - y_{i}) \ln (1 - p_{i}) + (y_{i} \ln y_{i} + (1 - y_{i}) \ln (1 - y_{i})) \\ &= \operatorname{argmin} \sum_{w}^{n} \sum_{i=1}^{n} y_{i} \ln p_{i} + (1 - y_{i}) \ln (1 - p_{i})) \\ &= \operatorname{argmin} \sum_{w}^{n} H(P(Y_{i}) \| P(\hat{Y}_{i})) \\ \therefore \Omega(w, x) &= D_{KL}(P_{u} \| P_{b}) \end{split}$$

As a result, we should choose  $\left[a\right]$  as our solution.

# **Validation**

# **P4**

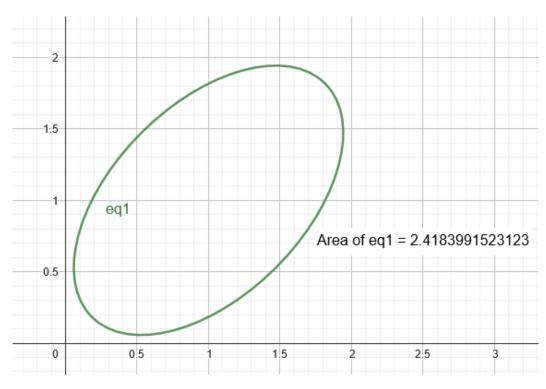
Since constant hypothesis  $h(x) = w_0$ , we only need to take y into the consideration. By selecting two examples for predicting another one example, we can obtain the equation

$$egin{align} E_{loocv} &= rac{1}{3}igg((y_1 - rac{y_2 + 1}{2})^2 + (y_2 - rac{y_1 + 1}{2})^2 + (1 - rac{y_1 + y_2}{2})^2igg) \ &= rac{1}{2}(y_1^2 + y_2^2 - y_1y_2 - y_1 - y_2 + 1) \end{array}$$

Since we are finding  $E_{loocv} \leq rac{1}{3}$  , we can rewrite the equation to be

$$(rac{x+1}{2}-y)^2+(rac{y+1}{2}-x)^2+(rac{x+y}{2}-1)^2=1 \ \Rightarrow rac{1}{2}(x^2+y^2-xy-x-y+1)=rac{1}{3}$$

Drawing the equation with <code>GeoBeBra</code> , we can simply get the area, which is  $pprox 2.4183 pprox rac{\pi}{3\sqrt{3}}$ 



As a result, we should choose  $\left[b\right]$  as our solution.

$$\begin{split} & \colon E_{out}(h) = \underset{x \sim P}{\mathcal{E}}[h(x) \neq f(x)] \\ & \colon \text{Variance } [E_{val}(h)] \\ & = \text{Variance } \left[\frac{1}{K} \sum_{n=1}^{K} \operatorname{err}(h(x_n), y_n)\right] \\ & = \frac{1}{K^2} \operatorname{Variance } \left[\sum_{n=1}^{K} \operatorname{err}(h(x_n), y_n)\right] \\ & = \frac{1}{K^2} \operatorname{Variance } \left[\operatorname{err}(h(x_1), y_1) + \operatorname{err}(h(x_2), y_2) + \dots + \operatorname{err}(h(x_K), y_K)\right] \\ & = \frac{1}{K^2} (K \cdot \operatorname{Variance } \left[\operatorname{err}(h(x), y)\right]) \\ & = \frac{1}{K} \operatorname{Variance } \left[\operatorname{err}(h(x), y)\right] \\ & = \frac{1}{K} \operatorname{Variance } \left[\operatorname{err}(h(x), y)\right] \\ & \therefore \Box = \frac{1}{K} \end{split}$$

Since the examples (x, y) are genetated from the i.i.d. distribution, the covariance the between examples should be zero.

As a result, we should choose  $\left[d\right]$  as our solution.

### P6

Since single data set in a binary classification has N positive examples and N negative examples, we can perform leave-one-out validation by selecting one positive or one negative example into the validation data set.

Once we select a positive example as a validation data set, there are N-1 positive examples and N negative examples in the training data set, and the binary classification algorithm  $A_{majority}$  will return negative to the training data set for the reason that it predicts the majority class. However,  $A_{majority}$  will return positive to the validation data set since there is only one positive example, and vice versa.

As a result, we know that  $A_{majority}$  will always return a different classification between the training data set and the validation data set, which means  $E_{loocv}(A_{majority}) = 1$ .

As a result, we should choose [d] as our solution.

### **P7**

The threshold of the decision stump model can be represented as

$$heta \in \{-1\} \cup \{rac{x_i + x_{i+1}}{2}\} : 1 \leq i \leq N-1 ext{ and } x_i 
eq x_{i+1}$$

Since x is generated by a uniform distribution in [-1,+1], and  $y=\mathrm{sign}(x)$ , it's trivial that  $E_{in}=0$  when  $\theta=0$ , and the threshold  $\theta$  with the lowest  $E_{in}$  in the decision stump model will be trained in the middle of the smallest positive value and the largest negative value. Therefore, the trained  $\theta=\frac{x_j+x_{j+1}}{2}$ , where  $x_j$  is the largest negative value, and  $x_{j+1}$  is the smallest positive value.

If the  $x_j$  is selected as the validation data set, then  $\theta=\frac{x_{j-1}+x_{j+1}}{2}$ , where  $x_{j-1}$  is the second largest negative value. Both  $\theta$  and the  $x_j$  are between  $x_{j-1}$  and  $x_{j+1}$ , then we have

- ullet if  $heta \geq x_j \Rightarrow h(x_j) = -1$ , which means the validation data set will be classified correctly.
- ullet if  $heta < x_j \Rightarrow h(x_j) = +1$ , which means the validation data set will be classified incorrectly.

And vice versa, if  $x_{j+1}$  is selected as the validation data set, then  $\theta = \frac{x_j + x_{j+2}}{2}$ , where  $x_{j+2}$  is the second smallest positive value. Both  $\theta$  and the  $x_{j+1}$  are between  $x_j$  and  $x_{j+2}$ , then we have

- if  $\theta \ge x_{j+1} \Rightarrow h(x_{j+1}) = -1$ , which means the validation data set will be classified incorrectly.
- ullet if  $heta < x_{j+1} \Rightarrow h(x_{j+1}) = +1$ , which means the validation data set will be classified correctly.

Since only  $x_j$  and  $x_{j+1}$  will be classified incorrectly by the threshold  $\theta$ , the tightest upper bound of the leave-one-out validation error to the decision stump model is equal to  $\max(E_{loocv}) = 2/N$ .

As a result, we should choose [c] as our solution.

# **Support Vector Machine**

# **P8**

Since the examples are ordered, the hard-margin SVM without transformation can be considered as a decision stump model when the separation line is orthogonal to the coordinate of the examples. It's trivial that the perfect separation line must be the middle of  $x_M$  and  $x_{M+1}$ , therefore, the hypothesis can be represented as

$$g_{svm}(x) = ext{sign}(w^Tx+b) = 0 = x - rac{x_M + x_{M+1}}{2}$$

Therefore, the largest margin is equal to  $rac{1}{2}(x_{M+1}-x_M)$ 

As a result, we should choose  $\left[e\right]$  as our solution.

### **P9**

Based on the subjections of  $\min(w,b) = \frac{1}{2} w^T w$ , we have

$$(w^Tx_n+b) \geq 1 \qquad ext{for } y_n = +1 \ -(w^Tx_n+b) \geq 1126 \quad ext{for } y_n = -1$$
  $y = egin{bmatrix} w_1 \ w_2 \end{bmatrix}^T egin{bmatrix} 0 & 4 \ 2 & 0 \ -1 & 0 \ 0 & 0 \end{bmatrix} + b = egin{bmatrix} +1 \ -1 \ +1 \ +1 \end{bmatrix} \Rightarrow egin{bmatrix} (4w_2+b) \geq 1 \ -(2w_1+b) \geq 1 \ (-w_1+b) \geq 1 \end{bmatrix}$ 

By drawing all data on the graph and from the subjections above, we know that the ratio of positive margin versus negative margin can represent as  $r_m=\frac{y_n=+1}{y_n=-1}=\frac{1}{1126}$ , hence, the hyperplane can be represented by the thinnest gap  $-(w^T[2,0]+b=1126)$  and  $(w^T[0,0]+b=1)$ , which implies that  $w^Tx+b=0 \Rightarrow -\frac{1127}{2}x+1=0$ , which means  $w=(-\frac{1127}{2},0),b=1$ .

As a result, we should choose [e] as our solution.

# P10

Since  $\alpha=1,b=0,h(x)=\mathrm{sign}(\sum_{n=1}^Ny_nK(x_n,x))$ , then we have  $E_{in}(h)=\frac{1}{N}\sum_{i=1}^N[h(x_i)\neq y_i]$ . If  $E_{in}(\hat{h})=0$ , all predictions should be classified correctly, and there is no classification error, which means

$$\left[\operatorname{sign}\left(\sum_{n=1}^{N}y_{n}K(x_{n},x_{i})
ight)=y_{i}
ight], orall i \overset{ ext{rewrite}}{\Longrightarrow}\left[\left(\sum_{n=1}^{N}y_{n}K(x_{n},x_{i})
ight)\cdot y_{i}>0
ight], orall i$$

To sum up the equation above, we have

$$egin{aligned} \sum_{i=1}^{N} \left(y_i \sum_{n=1}^{N} y_n K(x_n, x_i)
ight) &= \left[y_1^2 + y_1 \left(y_2 K(x_2, x_1) + \cdots + y_N K(x_N, x_1)
ight)
ight] \ &+ \left[y_2^2 + y_2 \left(y_1 K(x_1, x_2) + \cdots + y_N K(x_N, x_2)
ight)
ight] \ &+ \cdots \ &+ \left[y_N^2 + y_N \left(y_1 K(x_1, x_N) + \cdots + y_{N-1} K(x_{N-1}, x_N)
ight)
ight] > 0 \end{aligned}$$

Since  $\gamma>0, \|x_n 
eq x_m\| \geq \epsilon, \forall n 
eq m$ , we have

$$K(x_n,x_m) = \exp(-\gamma \|x_n-x_m\|^2) \leq \exp(-\gamma\epsilon^2) := \zeta$$
 if  $n=i, K(x_n,x_i) = \exp(-\gamma \|x_n-x_i\|^2) = \exp(-\gamma\cdot 0) = 0$ 

Taking  $\zeta$  back to the equation as the upper bound, then we get

$$0 < \sum_{i=1}^{N} \left( y_i \sum_{n=1}^{N} y_n K(x_n, x_i) \right) < \sum_{i=1}^{N} \left( y_i \sum_{n=1}^{N} y_n \exp(-\gamma \epsilon^2) \right) \ = \left[ y_1^2 + y_2^2 + \dots + y_N^2 \right] \ + \left[ y_1 \left( \underbrace{y_2 \exp(-\gamma \epsilon^2) + \dots + y_N \exp(-\gamma \epsilon^2)}_{N-1} \right) + \dots + y_N \left( \underbrace{y_1 \exp(-\gamma \epsilon^2) + \dots + y_{N-1} \exp(-\gamma \epsilon^2)}_{N-1} \right) \right]$$

And we know that

$$egin{aligned} orall n = m, y_m = (+1) \lor (-1) \Rightarrow y_i^2 = 1, orall i \ orall n 
eq m, -1 < y_n y_m < 1 \Rightarrow ext{assumed that all } y_n y_m = -1 \end{aligned}$$

Then, we have

$$egin{aligned} 0 < \sum_{i=1}^N \left(y_i \sum_{n=1}^N y_n \exp(-\gamma \epsilon^2)
ight) \ < \sum_{i=1}^N y_i^2 + \sum_{i=1}^N (-1) imes (N-1) \exp(-\gamma \epsilon^2) \ &= N + N imes [-(N-1) \exp(-\gamma \epsilon^2)] \ &= N - N(N-1) \exp(-\gamma \epsilon^2) \ &1 > (N-1) \exp(-\gamma \epsilon^2) \ rac{1}{N-1} > \exp(-\gamma \epsilon^2) \ \ln(N-1) < \gamma \epsilon^2 \ &\gamma > rac{\ln(N-1)}{\epsilon^2} \end{aligned}$$

As a result, we should choose  $\left[d\right]$  as our solution.

# P11

Applying the feature transform  $\phi$  to the Gaussian kernel, then we have

$$\|\phi(x)-\phi(x')\|^2 = <\phi(x)-\phi(x'), \phi(x)-\phi(x')> \ = <\phi(x), \phi(x)>-2 <\phi(x), \phi(x')> + <\phi(x'), \phi(x')> \ = \phi(x)^T\phi(x)-\phi(x)^T\phi(x')-\phi(x')^T\phi(x)+\phi(x')^T\phi(x') \ = K(x,x)-2K(x,x')+K(x',x') \ = 2-2\exp(-\gamma\|x-x'\|^2) \ dots \exp(-\gamma\|x-x'\|^2) \in (0,1] \ dots 2>\|\phi(x)-\phi(x')\|^2 \ \sqrt{2}pprox 1.5>\|\phi(x)-\phi(x')\|$$

Therefore, the tightest upper bound for the distane in the  $\mathcal{Z}$ -space is 1.5.

As a result, we should choose [d] as our solution.

# **Experiments with Regularized Logistic Regression**

### P12

Since L2-regularized logistic regression -s 0 in liblinear solves

$$w = \min_{w} rac{1}{2} w^T w + C \sum_{i=1}^{N} \log(1 + \exp(-y_i w^T x_i))$$

and we solve

$$w_{\lambda} = \operatornamewithlimits{argmin}_{w} rac{\lambda}{N} \|w\|^2 + rac{1}{N} \sum_{n=1}^{N} \ln(1 + \exp(-y_n w^T \Phi_4(x_n)))$$

Then we know that

$$rac{1}{2}=rac{\lambda}{N}, C=rac{1}{N}\Rightarrow C=rac{1}{2\lambda}$$

The parameter -c cost : set the parameter C (default 1) should be set to  $C=rac{1}{2\lambda}.$ 

```
1
     import numpy as np
 2
     from itertools import combinations_with_replacement
     from liblinear.liblinearutil import *
 3
 4
     def transformation(args, x):
 5
         x_{origin} = x[:, 1:] # remove <math>x_{n=0}=1
 6
 7
         transformed x = x.copy()
         combination tuples = list()
 8
         iterable = list(i for i in range(1, args['dimension']))
 9
10
         for r in range(1, args['Q']):
              combination_tuples.extend(list(combinations_with_replacement(iterable, r+1)))
11
12
         for combination in range(len(combination_tuples)):
13
14
              x temp = 1
              for j in combination_tuples[combination]:
15
                  x_{temp} = x_{origin}[:,j-1].reshape(-1, 1) * x_{temp}
16
17
              transformed_x = np.hstack((
18
                  transformed_x, x_temp
19
                  ))
20
         return transformed_x
```

```
def read file(args, filename):
1
2
         data = np.loadtxt(filename, dtype=float)
3
         x = data[:,:-1]
4
         x = np.c_{np.ones(len(x)), x] # x_{n=0}=1
5
         y = data[:,-1]
6
         args['dimension'] = x.shape[1]
7
         return x, y
1
     def main(args):
2
         x_train, y_train = read_file(args, args['filename_train'])
3
         x_test, y_test = read_file(args, args['filename_test'])
         x_train = transformation(args, x_train)
4
5
         x_test = transformation(args, x_test)
6
7
         E_out = list()
8
         for lamb in args['lambda']:
9
             C = 1/(2*lamb)
             prob = problem(y_train, x_train)
10
11
             param = parameter('-s 0 -c {} -e 0.000001 -q'.format(C))
12
             m = train(prob, param)
             p_label, p_acc, p_val = predict(y_test, x_test, m)
13
14
             E_out.append(round(np.mean(y_test != p_label), 6))
15
16
         min_E_out = min(E_out)
17
         min_E_out_index = [i for i, v in enumerate(E_out) if v == min_E_out]
18
         print("E_out: ", E_out)
         print("min E_out: ", min(E_out))
19
20
         print("Choice: ", chr(97+max(min_E_out_index)))
     if __name__ == '__main__':
1
2
         args = {
3
             'dimension': 0,
             'filename_test': "hw4_test.dat",
4
5
             'filename_train': "hw4_train.dat",
6
             'lambda': [10**(-6), 10**(-3), 10**(0), 10**3, 10**6],
7
             'Q': 4,
8
         }
9
10
         main(args)
```

```
NumbaDeprecationWarning: The 'nopython' keyword argument was not supplied to the 'numl
 1
 2
       def csr_to_problem_jit(l, x_val, x_ind, x_rowptr, prob_val, prob_ind, prob_rowptr):
 3
     Accuracy = 77.4% (387/500) (classification)
     Accuracy = 82.2% (411/500) (classification)
 4
     Accuracy = 84.6% (423/500) (classification)
5
     Accuracy = 85.8% (429/500) (classification)
 6
7
     Accuracy = 81.2\% (406/500) (classification)
8
     E_out: [0.226, 0.178, 0.154, 0.142, 0.188]
 9
     min E_out: 0.142
10
     Choice: d
```

Therefore, the best  $\log_{10}(\lambda^*)=3$ 

As a result, we should choose [d] as our solution.

# P13

We use the same code above, then update the main function. The updated functions show below.

```
1
     def main(args):
2
         x_train, y_train = read_file(args, args['filename_train'])
 3
         x_test, y_test = read_file(args, args['filename_test'])
 4
         x_train = transformation(args, x_train)
 5
         x_test = transformation(args, x_test)
 6
7
         E_in = list()
8
         for lamb in args['lambda']:
9
             C = 1/(2*lamb)
             prob = problem(y_train, x_train)
10
             param = parameter('-s 0 -c {} -e 0.000001 -q'.format(C))
11
12
             m = train(prob, param)
13
             p_label, p_acc, p_val = predict(y_train, x_train, m)
14
             E_in.append(round(np.mean(y_train != p_label), 6))
15
16
         min_E_in = min(E_in)
17
         min_E_in_index = [i for i, v in enumerate(E_in) if v == min_E_in]
         print("E_in: ", E_in)
18
         print("min E_in: ", min(E_in))
19
20
         print("Choice: ", chr(97+max(min_E_in_index)))
```

```
NumbaDeprecationWarning: The 'nopython' keyword argument was not supplied to the 'numl
 1
 2
       def csr_to_problem_jit(l, x_val, x_ind, x_rowptr, prob_val, prob_ind, prob_rowptr):
 3
     Accuracy = 100% (200/200) (classification)
     Accuracy = 100% (200/200) (classification)
 4
     Accuracy = 100% (200/200) (classification)
5
     Accuracy = 96% (192/200) (classification)
 6
7
     Accuracy = 76% (152/200) (classification)
     E_in: [0.0, 0.0, 0.0, 0.04, 0.24]
 9
     min E_in: 0.0
10
     Choice: c
```

Since  $E_{in}(w_{\log_{10}(\lambda^*)=-6})=E_{in}(w_{\log_{10}(\lambda^*)=-3})=E_{in}(w_{\log_{10}(\lambda^*)=0})=0$ , by selecting the largest  $\lambda$ , therefore, the best  $\log_{10}(\lambda^*)=0$ 

As a result, we should choose  $\left[c\right]$  as our solution.

# P14

We use the same code above, then add a train\_test\_split function and update the main function and the parameter of args. The updated functions show below.

```
1
     import numpy as np
 2
     import random
 3
     from itertools import combinations with replacement
 4
     from liblinear.liblinearutil import *
 5
     def train_test_split(args, x, y):
 6
 7
         rng = np.random.default_rng()
 8
         idx = np.arange(len(x))
 9
         rng.shuffle(idx)
         x_train = x[idx[:args['split']]]
10
11
         y_train = y[idx[:args['split']]]
12
         x_test = x[idx[args['split']:]]
13
         y_test = y[idx[args['split']:]]
14
         return x_train, y_train, x_test, y_test
```

```
1
 2
     def main(args):
3
         x_train, y_train = read_file(args, args['filename_train'])
4
         x_train = transformation(args, x_train)
5
         lamb_list = list()
6
         E_val_list = list()
7
         for i in range(args['repeat_time']):
8
             x_train_split, y_train_split, x_test_split, y_test_split = train_test_split(a)
9
             args['size'] = x_train_split.shape[0]
10
11
             E_val = list()
             for lamb in args['lambda']:
12
                 C = 1/(2*lamb)
13
14
                 prob = problem(y_train_split, x_train_split)
15
                 param = parameter('-s 0 -c {} -e 0.000001 -q'.format(C))
                 m = train(prob, param)
16
17
                 p_label, p_acc, p_val = predict(y_test_split, x_test_split, m)
                 E_val.append(round(np.mean(y_test_split != p_label), 6))
18
19
20
             min_E_val = min(E_val)
21
             min_E_val_index = [i for i, v in enumerate(E_val) if v == min_E_val]
22
             lamb_list.append(max(min_E_val_index))
23
             E_val_list.append(min_E_val)
24
25
         best_lambda = args['lambda'][max(set(lamb_list), key = lamb_list.count)]
         print("best lambda: ", best lambda)
26
1
     if __name__ == '__main__':
 2
         args = {
3
             'dimension': 0,
4
             'filename_test': "hw4_test.dat",
             'filename_train': "hw4_train.dat",
5
             'lambda': [10**(-6), 10**(-3), 10**(0), 10**3, 10**6],
 6
7
             'Q': 4,
8
             'repeat_time': 256,
9
             'seed': 1126,
             'size': 0,
10
11
             'split': 120
12
         }
13
14
         main(args)
 1 best_lambda: 1000
```

As a result, we should choose  $\left[d\right]$  as our solution.

We use the same code above, then update the main function. The updated functions show below.

```
1
     def main(args):
 2
         x train, y train = read file(args, args['filename train'])
 3
         x_train = transformation(args, x_train)
         x_test, y_test = read_file(args, args['filename_test'])
 4
5
         x_test = transformation(args, x_test)
         lamb_list = list()
 6
 7
         E val list = list()
8
         E_out_list = list()
9
         for i in range(args['repeat_time']):
             x_train_split, y_train_split, x_test_split, y_test_split = train_test_split(a)
10
11
             args['size'] = x_train_split.shape[0]
12
             model list = list()
13
14
             E_val = list()
15
             for lamb in args['lambda']:
16
                 C = 1/(2*lamb)
17
                 prob = problem(y_train_split, x_train_split)
18
                 param = parameter('-s 0 -c {} -e 0.000001 -q'.format(C))
19
                 m = train(prob, param)
20
                 p_label, p_acc, p_val = predict(y_test_split, x_test_split, m)
                 E_val.append(round(np.mean(y_test_split != p_label), 6))
21
22
                 model_list.append(m)
23
24
             min E val = min(E val)
             min_E_val_index = [i for i, v in enumerate(E_val) if v == min_E_val]
25
26
             lamb list.append(max(min E val index))
27
             E_val_list.append(min_E_val)
             # predict test dataset
28
29
             p_label, p_acc, p_val = predict(y_test, x_test, model_list[max(min_E_val_index
             E_out_list.append(round(np.mean(y_test != p_label), 6))
30
31
32
         best_lambda = args['lambda'][max(set(lamb_list), key = lamb_list.count)]
         print("best_lambda: ", best_lambda)
33
         print("E_out: ", np.mean(E_out_list))
34
```

best\_lambda: 1000
best\_lambda: 1000
best\_lambda: 1000

As a result, we should choose  $\left[c\right]$  as our solution.

### P16

We use the same code above, then update the main function. The updated functions show below.

```
1
     def main(args):
 2
         x_train, y_train = read_file(args, args['filename_train'])
         x_train = transformation(args, x_train)
 3
4
         x_test, y_test = read_file(args, args['filename_test'])
5
         x_test = transformation(args, x_test)
 6
         lamb list = list()
7
         E_val_list = list()
8
         E_out_list = list()
9
         for i in range(args['repeat_time']):
10
             x_train_split, y_train_split, x_test_split, y_test_split = train_test_split(ar
11
             args['size'] = x_train_split.shape[0]
12
13
             model_list = list()
14
             E_val = list()
15
             for lamb in args['lambda']:
16
                 C = 1/(2*lamb)
17
                 prob = problem(y_train_split, x_train_split)
                 param = parameter('-s 0 -c {} -e 0.000001 -q'.format(C))
18
19
                 m = train(prob, param)
                 p_label, p_acc, p_val = predict(y_test_split, x_test_split, m)
20
21
                  E_val.append(round(np.mean(y_test_split != p_label), 6))
22
                 model_list.append(m)
23
24
             min_E_val = min(E_val)
25
             min_E_val_index = [i for i, v in enumerate(E_val) if v == min_E_val]
             lamb list.append(max(min E val index))
26
27
             E_val_list.append(min_E_val)
28
29
             # rebuild the model from raw training dataset
30
             prob = problem(y_train, x_train)
             param = parameter('-s 0 -c \{\} -e 0.000001 -q'.format(1/(2*(args['lambda'][max])
31
32
             m = train(prob, param)
33
             # predict from raw testing dataset
34
             p_label, p_acc, p_val = predict(y_test, x_test, m)
35
             E_out_list.append(round(np.mean(y_test != p_label), 6))
36
37
         best_lambda = args['lambda'][max(set(lamb_list), key = lamb_list.count)]
38
         print("best_lambda: ", best_lambda)
         print("E_out: ", np.mean(E_out_list))
39
 1
     best_lambda: 1000
    E out: 0.15015624999999996
```

As a result, we should choose  $\left[b\right]$  as our solution.

### P17

We use the same code above, then add a KFold function and update the main function and the parameter of args. The updated functions show below.

```
def KFold(args, x, y):
    rng = np.random.default_rng()
    idx = np.arange(len(x))
    rng.shuffle(idx)
    x_split = np.array([x[idx[i*args['fold_size']:(i+1)*args['fold_size']]] for i in |
    y_split = np.array([y[idx[i*args['fold_size']:(i+1)*args['fold_size']]] for i in |
    return x_split, y_split
```

```
1
     def main(args):
2
         x_train, y_train = read_file(args, args['filename_train'])
 3
         x_train = transformation(args, x_train)
         x_test, y_test = read_file(args, args['filename_test'])
4
5
         x_test = transformation(args, x_test)
         args['size'] = x_train.shape[0]
 6
 7
         args['fold_size'] = int(args['size'] / args['fold'])
8
9
         E_cv_list = list()
10
         for i in range(args['repeat_time']):
11
             E_cv = list()
             x_split, y_split = KFold(args, x_train, y_train)
12
13
14
             for lamb in args['lambda']:
15
                 C = 1/(2*lamb)
                 E_val = list()
16
17
                 for fold in range(args['fold']):
                      x_valid_fold = x_split[fold]
18
                      y_valid_fold = y_split[fold]
19
20
                      x_train_fold = np.vstack(x_split[j] for j in range(args['fold']) if j
21
                      y_train_fold = np.hstack(y_split[j] for j in range(args['fold']) if j
22
23
                      prob = problem(y_train_fold, x_train_fold)
                      param = parameter('-s 0 -c {} -e 0.000001 -q'.format(C))
24
25
                      m = train(prob, param)
                      p_label, p_acc, p_val = predict(y_valid_fold, x_valid_fold, m)
26
27
                      E_val.append(round(np.mean(y_valid_fold != p_label), 6))
28
                  E_cv.append(np.mean(E_val))
             E_cv_list.append(min(E_cv))
29
30
31
         print("E_cv: ", np.mean(E_cv_list))
```

```
if __name__ == '__main__':
 1
 2
         args = {
              'dimension': 0,
 3
 4
              'filename_test': "hw4_test.dat",
 5
              'filename_train': "hw4_train.dat",
              'fold': 5,
 6
 7
              'fold_size': 0, # size of each fold
              'lambda': [10**(-6), 10**(-3), 10**(0), 10**3, 10**6],
 8
 9
              'Q': 4,
10
              'repeat_time': 256,
              'seed': 1126,
11
              'size': 0,
12
              'split': 120
13
14
         }
15
         main(args)
16
```

1 E\_cv: 0.12927734375

As a result, we should choose [a] as our solution.

# P18

Since L1-regularized logistic regression -s 6 in liblinear solves

$$w = \min_{w} \sum |w_j| + C \sum_{i=1}^N \log(1 + \exp(-y_i w^T x_i))$$

and we solve

$$w_{\lambda} = \operatornamewithlimits{argmin}_{w} rac{\lambda}{N} \|w\|_1 + rac{1}{N} \sum_{n=1}^{N} \ln(1 + \exp(-y_n w^T \Phi_4(x_n)))$$

Then we know that

$$1 = \frac{\lambda}{N}, C = \frac{1}{N} \Rightarrow C = \frac{1}{\lambda}$$

The parameter -c cost : set the parameter C (default 1) should be set to  $C=rac{1}{\lambda}.$ 

We use the same code from P12, then update parameter C and Param in the Param in the Param function. The updated functions show below.

```
NumbaDeprecationWarning: The 'nopython' keyword argument was not supplied to the 'numl
 1
 2
       def csr_to_problem_jit(l, x_val, x_ind, x_rowptr, prob_val, prob_ind, prob_rowptr):
 3
     Accuracy = 76.8\% (384/500) (classification)
     Accuracy = 83.6% (418/500) (classification)
 4
     Accuracy = 84.6% (423/500) (classification)
5
     Accuracy = 68% (340/500) (classification)
 6
7
     Accuracy = 49.2\% (246/500) (classification)
8
     E_out: [0.232, 0.164, 0.154, 0.32, 0.508]
 9
     min E_out: 0.154
10
     Choice: c
```

As a result, we should choose [c] as our solution.

# P19

We use the same code above, then update the main function. The updated functions show below.

```
1
     def main(args):
 2
         x_train, y_train = read_file(args, args['filename_train'])
 3
         x_test, y_test = read_file(args, args['filename_test'])
 4
         x_train = transformation(args, x_train)
 5
         x_test = transformation(args, x_test)
         args['dimension'] = x_train.shape[1]
 6
 7
 8
         model_list = list()
 9
         E_out = list()
         for lamb in args['lambda']:
10
11
             C = 1/(lamb)
12
             prob = problem(y_train, x_train)
             param = parameter('-s 6 -c {} -e 0.000001 -q'.format(C))
13
14
             m = train(prob, param)
15
             model_list.append(np.array([m.w[i] for i in range(args['dimension']) if np.ab
              p_label, p_acc, p_val = predict(y_test, x_test, m)
16
17
              E_out.append(round(np.mean(y_test != p_label), 6))
18
19
         min E out = min(E out)
20
         min_E_out_index = [i for i, v in enumerate(E_out) if v == min_E_out]
         print("|m.w[i] <= 10**(-6)|: ", len(model_list[min_E_out_index[0]]))</pre>
21
```

```
NumbaDeprecationWarning: The 'nopython' keyword argument was not supplied to the 'numl def csr_to_problem_jit(l, x_val, x_ind, x_rowptr, prob_val, prob_ind, prob_rowptr):

Accuracy = 76.8% (384/500) (classification)

Accuracy = 83.6% (418/500) (classification)

Accuracy = 84.6% (423/500) (classification)

Accuracy = 68% (340/500) (classification)

Accuracy = 49.2% (246/500) (classification)

|m.w[i] <= 10**(-6)|: 960
```

As a result, we should choose [e] as our solution.

# **P20**

We use the same code above, then update parameter c and param in the main function. The updated functions show below.

```
1
    C = 1/(2*lamb)
    param = parameter('-s 0 -c {} -e 0.000001 -q'.format(C))
1
    NumbaDeprecationWarning: The 'nopython' keyword argument was not supplied to the 'numl
      def csr_to_problem_jit(1, x_val, x_ind, x_rowptr, prob_val, prob_ind, prob_rowptr):
2
3
    Accuracy = 77.4% (387/500) (classification)
    Accuracy = 82.2% (411/500) (classification)
4
5
    Accuracy = 84.6% (423/500) (classification)
    Accuracy = 85.8\% (429/500) (classification)
6
7
    Accuracy = 81.2% (406/500) (classification)
    |m.w[i] <= 10**(-6)|: 1
```

As a result, we should choose [a] as our solution.