

### Exercise 1:

1. The methods not working in this case is Logistic regression, the error message is:

```
/Users/houninghi/Library/Python/3.8/lib/python/site-packages/statsmodels/base/model.py:607: ConvergenceWarning: Maximum Likelihood optimization failed to converge. Check mle_retvals
  warnings.warn("Maximum Likelihood optimization failed to converge")
```

The loss function for logit is  $l_w(x_i, y_i) = -y_i \log p_i - (1-y_i) \log(1-p_i)$ , where  $p_i = \frac{1}{1+\exp(-\langle x_i, w \rangle)}$   
when data is linearly separable  $\langle x_i, w \rangle > 0$  when  $y_i=1$ ,  $\langle x_i, w \rangle < 0$  when  $y_i=0$

We can compute:  $l_{2w}(x_i, y_i)$ , If  $y_i=0$ ,  $l_w(x_i, y_i) = -\log(1-p_i)$ ,  $\frac{1}{1+\exp(-\langle x_i, w \rangle)} < \frac{1}{1+\exp(-\langle x_i, 2w \rangle)}$

$l_{2w}(x_i, y_i) > l_w(x_i, y_i)$ , so which means loss gets smaller if  $C \cdot w$ ,  $C$  increases

So the loss never converges to a minimum value, that's exactly what error said.

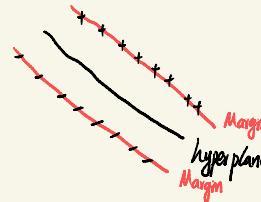
Compare Soft Margin SVC and Logistic regression

(1) Similarity: Both use for classification task, Both Logistic Regression and Soft Margin SVM

(2) Difference: The loss function is different: Logit models the probability that a data point belongs to a class, SVC focus on maximizing margin instead

2. By python code, there are 2000 values  $\leq 1$ , all  $y_i \hat{=} 1$ , so:

For data set A, the sketch is like:



$$W = \sum \alpha_i y_i X_i \quad \text{where } y_i \hat{=} 1: (\text{By complementary slackness, } \alpha_i = 0 \text{ or } y_i \hat{=} 1)$$

By the python results: There are 2 vectors in the SVM, so we require 2 points support

For data set B, the failed model is Hardmargin SVC, because the data set is not linearly separable. I need 192 support vectors

Exercises:

1. The primal form is  $\min_{\substack{w \in \mathbb{R}^n \\ b \in \mathbb{R}}} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \max\{y_i - (w^T x_i + b) - \varepsilon, 0\}$

$$= \min_{\substack{w \in \mathbb{R}^n \\ b \in \mathbb{R}}} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \max\{y_i - (w^T x_i + b) - \varepsilon, (w^T x_i + b) - y_i - \varepsilon, 0\}$$

$$= \min_{\substack{w \in \mathbb{R}^n \\ b \in \mathbb{R}}} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n y_i, \text{ s.t. } y_i - (w^T x_i + b) - \varepsilon \leq y_i, y_i \geq 0, (w^T x_i + b) - y_i - \varepsilon \leq y_i$$

Take Lagrangian dual to remove constraint, and by strong duality.

$$\max_{\substack{\alpha, \beta, d \in \mathbb{R}^n \\ \alpha, \beta, d \geq 0}} \min_{\substack{w \in \mathbb{R}^n \\ b \in \mathbb{R}}} \frac{1}{2} \|w\|^2 + \sum_{i=1}^n (C y_i + \alpha_i (y_i - (w^T x_i + b) - \varepsilon - y_i) - \beta_i \cdot y_i + d_i (w^T x_i + b - y_i - \varepsilon - y_i))$$

$$\text{set gradient to 0: } \frac{\partial}{\partial b} = \sum_i d_i - \alpha_i = 0 \quad \frac{\partial}{\partial w} = \sum_i \alpha_i x_i + \beta_i x_i + d_i x_i = 0 \quad \frac{\partial}{\partial w} = (-\alpha_i - \beta_i - d_i) \cdot x_i = 0 \Rightarrow w = \sum_i (\alpha_i + d_i) \cdot x_i$$

$$\begin{aligned} &= \max_{\substack{\alpha, \beta, d \in \mathbb{R}^n \\ \alpha, \beta, d \geq 0}} \min_{\substack{w \in \mathbb{R}^n \\ b \in \mathbb{R}}} \frac{1}{2} \|w\|^2 + \sum_{i=1}^n (\alpha_i - \alpha_i) (w^T x_i + b - y_i) + \underbrace{\sum_{i=1}^n ((-\alpha_i - \beta_i - d_i) \cdot y_i - \frac{1}{2} (\alpha_i + d_i) \cdot \varepsilon)}_0 \end{aligned}$$

$$\begin{aligned} &= \max_{\substack{\alpha, \beta, d \in \mathbb{R}^n \\ \alpha, \beta, d \geq 0}} \min_{\substack{w \in \mathbb{R}^n \\ b \in \mathbb{R}}} \frac{1}{2} \|w\|^2 + \underbrace{\sum_{i=1}^n (\alpha_i - \alpha_i) b}_0 + \underbrace{w^T \sum_{i=1}^n (\alpha_i + d_i) \cdot x_i - \frac{1}{2} \sum_{i=1}^n (\alpha_i + d_i) \cdot \varepsilon - \frac{1}{2} \sum_{i=1}^n (\alpha_i + d_i) \cdot y_i - \|w\|^2} \end{aligned}$$

$$\begin{aligned} &= \max_{\substack{\alpha, d \in \mathbb{R}^n \\ \alpha, d \geq 0}} \min_{\substack{w \in \mathbb{R}^n \\ b \in \mathbb{R}}} -\frac{1}{2} \|w\|^2 - \sum_{i=1}^n (\alpha_i + d_i) \cdot \varepsilon - \sum_{i=1}^n (\alpha_i + d_i) \cdot y_i \end{aligned}$$

(2) The gradient is:

$$\text{If } \varepsilon + (w^T x_i + b) < y_i, \text{ then } \frac{\partial}{\partial w} \left( C \sum_{i=1}^n \max\{y_i - (w^T x_i + b) - \varepsilon, 0\} \right) = -C x_i$$

$$\frac{\partial}{\partial b} \left( C \sum_{i=1}^n \max\{y_i - (w^T x_i + b) - \varepsilon, 0\} \right) = -C$$

$$\text{If } -\varepsilon < y_i - (w^T x_i + b) \leq \varepsilon, \text{ then } \frac{\partial}{\partial w} \left( C \sum_{i=1}^n \max\{y_i - (w^T x_i + b) - \varepsilon, 0\} \right) = 0$$

$$\frac{\partial}{\partial b} \left( C \sum_{i=1}^n \max\{y_i - (w^T x_i + b) - \varepsilon, 0\} \right) = 0$$

$$\text{if } w^T x_i + b - \varepsilon \geq y_i, \text{ then } \frac{\partial}{\partial w} \left( C \sum_{i=1}^n \max\{(w^T x_i + b) - y_i - \varepsilon, 0\} \right) = C x_i$$

$$\frac{\partial}{\partial b} \left( C \sum_{i=1}^n \max\{(w^T x_i + b) - y_i - \varepsilon, 0\} \right) = C$$

$$(3) P^{\eta}(w) = \arg \min_{z} \frac{1}{2\eta} \|z-w\|_2^2 + \frac{1}{2} \|z\|_2^2$$

$$= \arg \min_{z} \frac{1}{2\eta} \|z\|_2^2 - \frac{1}{\eta} z \cdot w + \frac{1}{2\eta} \|w\|_2^2 + \frac{1}{2} \|z\|_2^2.$$

So that  $\frac{\partial P^{\eta}(w)}{\partial z} = z + \frac{1}{\eta} z - \frac{1}{\eta} w$ , set gradient to 0, then  
 $z = \frac{1}{1+\eta} \cdot w$ .

(4) The output of the training output by selecting appropriate step size

```
the Training loss is 610.554899840419
the Training error is 610.0421934063855
the Test error is 776.4545373565409
0      -0.021302
1      -0.088689
2      -0.005106
3      0.101048
4      0.035666
...
195     0.009747
196    -0.007770
197    -0.079203
198     0.028541
199    -0.007645
Length: 200, dtype: float64 0.006
kous-MacBook-Air:A2 houningzhi$
```

Exercise 3:

$$\begin{aligned} \text{(1) } k(x,y) &= e^{-\alpha(x-y)^2}, \text{ For Taylor expansion, } e^x = \frac{x^0}{0!} + \frac{x^1}{1!} + \frac{x^2}{2!} + \dots \\ &= e^{-\alpha(x^2-2xy+y^2)} \\ &= e^{-\alpha(x^2)y^2} \cdot e^{2\alpha xy} \quad (\text{Taylor expansion on } e^{2\alpha xy}) \\ &= e^{-\alpha(x^2)y^2} \cdot \left[ 1 + \frac{2\alpha xy}{1!} + \frac{(2\alpha xy)^2}{2!} + \dots \right] \\ &= e^{-\alpha x^2} \cdot e^{-\alpha y^2} \left[ 1, \sqrt{\frac{2\alpha}{1!}} x, \sqrt{\frac{2\alpha}{2!}} x^2, \dots \right]^T \left[ 1, \sqrt{\frac{2\alpha}{1!}} y, \sqrt{\frac{2\alpha}{2!}} y^2, \dots \right] \end{aligned}$$

where clearly we can observe that

$$\phi(x) = e^{-\alpha x^2} \left[ 1, \sqrt{\frac{2\alpha}{1!}} x, \sqrt{\frac{2\alpha}{2!}} x^2, \dots \right], \text{ thus that } k(x,y) = \phi(x)^T \phi(y)$$

I'd love to solve the dual representation, because  $\phi(x)$  maps  $x$  to an infinite dimension vector  
It's impossible to solve primal form but in dual form, we only need  $k(x,y)$ .

(2) It's a kernel, we know that:

$$\begin{aligned} k(x,y) &= \frac{1}{1-xy}, \\ &= 1+xy+x^2y^2+\dots, \text{ by Taylor expansion} \\ &= [1, x, x^2, \dots]^T [1, y, y^2, \dots] \\ &= \phi(x)^T \phi(y), \text{ So:} \\ \text{The feature map } \phi(x) &= [1, x, x^2, \dots]^T. \end{aligned}$$

(3). It's not a kernel, we can give a counter example:

$x = [0.1 \ 10 \ 100]$  So the matrix  $k$  looks like:

$$k = \begin{bmatrix} \log(0.01) & \log(0.1) & \log(1) \\ \log(0.1) & \log(10) & \log(100) \\ \log(1) & \log(100) & \log(1000) \end{bmatrix}, \text{ If it's a kernel, } k \text{ must be PSD. given } D = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$$

$$D^T k D = \{(\log(0.01)+\log(0.1)-\log(1)), (\log(0.1)+\log(10)+\log(100)), (\log(1)+\log(100)+\log(1000))\}$$

$$\begin{aligned} D^T k D &= (\log(0.01)+\log(0.1)+\log(1000))-2\log(1)-2\log(100) \\ &= -147291 \dots < 0 \end{aligned}$$

So  $k$  is not PSD, so  $k$  is not a kernel

(4) It's not a kernel, we can give a counter example:

$x = [\frac{\pi}{3} \ \frac{2\pi}{3}]$ , so that  $k = \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}$ . let  $D = \begin{bmatrix} 1 \\ -1 \end{bmatrix}$ , so we compute  $D^T k D$ :

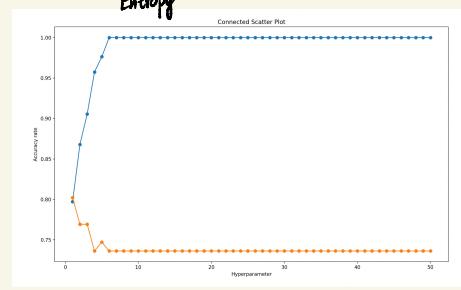
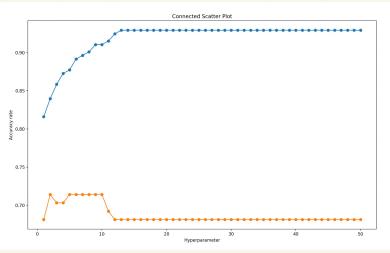
$$D^T k D = -4 < 0$$

So that  $k$  is not PSD, contradiction; so  $k(x,y)$  is not a kernel.

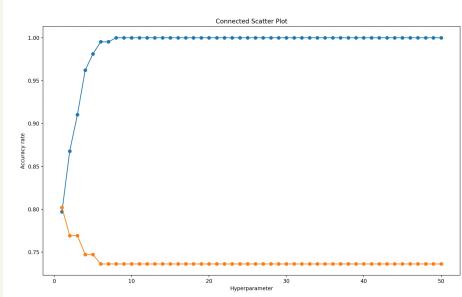
Exercise 4: Hyperparameter: depth ranges from 1 to 50

(a)

Misclassification error



Gini-Index



Analysis: These three plots shows how accuracy changes for test data when maximum depth for tree increases.

Generally: The training accuracy can reach 1 finally, which makes sense. If we have many leaves, we can classify 1 data on one leave to be 100% accurate on Training data, which is also overfitting.

In contrast, the test accuracy decreases gradually due to the overfitting

(b) The result for bagging is better than the result without (73%)

```
101
101
the median is 0.8131868131868132, the maximum is 0.8351648351648352, the minimum is 0.7912087912087912
kous-MacBook-Air:A2 houningzhi
```

The result for random forest and bagging is:

```
101
101
the median is 0.8241758241758241, the maximum is 0.8461538461538461, the minimum is 0.8021978021978022
kous-MacBook-Air:A2 houningzhi
```

It's better than any result before!