

1. (a).

$$\begin{aligned}
 f(\theta) &= \frac{1}{m} \sum_{i=1}^m \left(\sum_{k=1}^K y_{ik} x_i^T \theta_k - \log \left(\sum_{k=1}^K \exp(x_i^T \theta_k) \right) \right) \\
 \frac{d(f(\theta))}{d\theta_k} &= \frac{1}{m} \sum_{i=1}^m \left(\frac{d \left(\sum_{k=1}^K y_{ik} x_i^T \theta_k \right)}{d\theta_k} - \frac{d \left(\log \left(\sum_{k=1}^K \exp(x_i^T \theta_k) \right) \right)}{d\theta_k} \right) \\
 &= \frac{1}{m} \sum_{i=1}^m \left(y_{ik} x_i - \frac{1}{\sum_{j=1}^K \exp(x_i^T \theta_j)} \cdot \frac{d \left(\sum_{k=1}^K \exp(x_i^T \theta_k) \right)}{d\theta_k} \right) \\
 &= \frac{1}{m} \sum_{i=1}^m \left(y_{ik} x_i - \frac{1}{\sum_{j=1}^K \exp(x_i^T \theta_j)} \cdot \exp(x_i^T \theta_k) \cdot \frac{d(x_i^T \theta_k)}{d\theta_k} \right) \\
 &= \frac{1}{m} \sum_{i=1}^m \left(y_{ik} x_i - \frac{1}{\sum_{j=1}^K \exp(x_i^T \theta_j)} \cdot \exp(x_i^T \theta_k) \cdot x_i \right) \\
 &= \frac{1}{m} \sum_{i=1}^m \left(y_{ik} - \frac{\exp(x_i^T \theta_k)}{\sum_{j=1}^K \exp(x_i^T \theta_j)} \right) x_i
 \end{aligned}$$

(b).

1. Use $D = \max \theta_i$, $\theta_i - D \leq 0$. This makes $\exp(\theta_i - D) \leq 1$, which prevents $\exp(\theta_i)$ from becoming very large due to very large values in θ_i , thus preventing overflow.

2. Use $D = \min \theta_i$, $\theta_i - D \geq 0$. This makes $\exp(\theta_i - D) \geq 1$, which prevents $\exp(\theta_i)$ from becoming very close to zero due to very negative values in θ_i , thus preventing underflow.

(c).

```
import numpy as np
import matplotlib.pyplot as plt
import scipy.io as sio
from sklearn.preprocessing import OneHotEncoder

data = sio.loadmat('mnist.mat')
Xtrain, Xtest = data['trainX'].astype(float), data['testX'].astype(float)
ytrain, ytest = data['trainY'][0], data['testY'][0]
ytrain_onehot = OneHotEncoder(categories='auto').fit_transform(ytrain.reshape(-1, 1))
m, n = Xtrain.shape

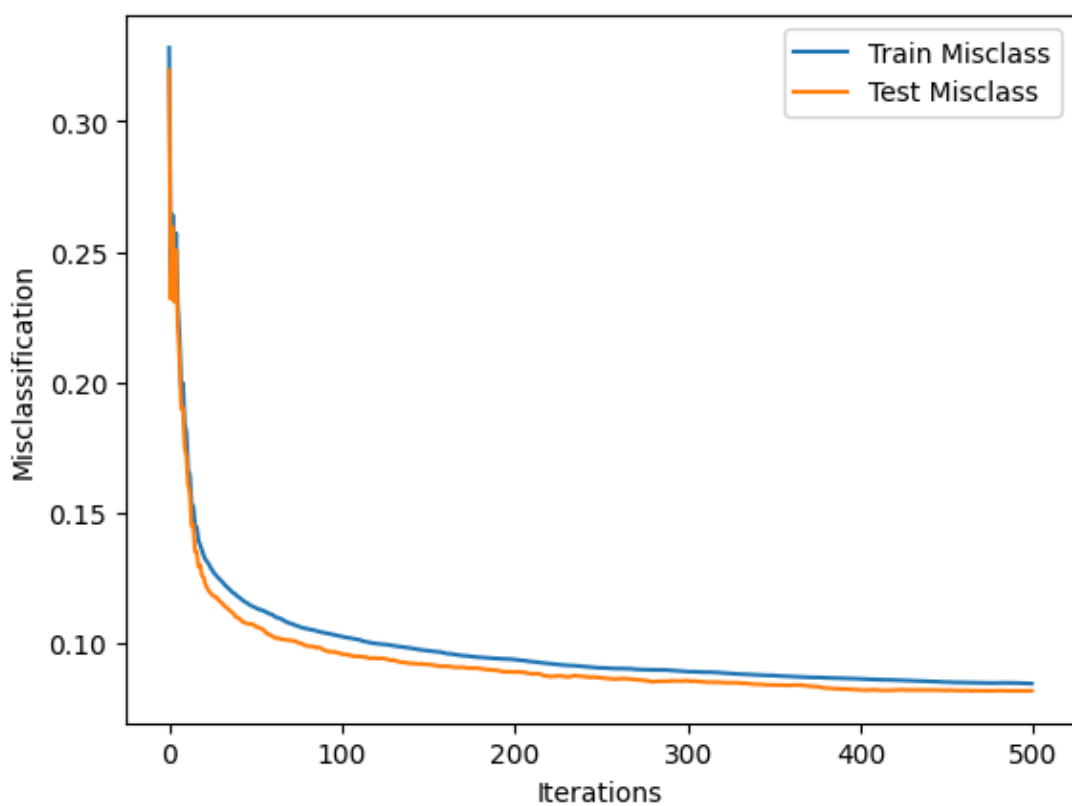
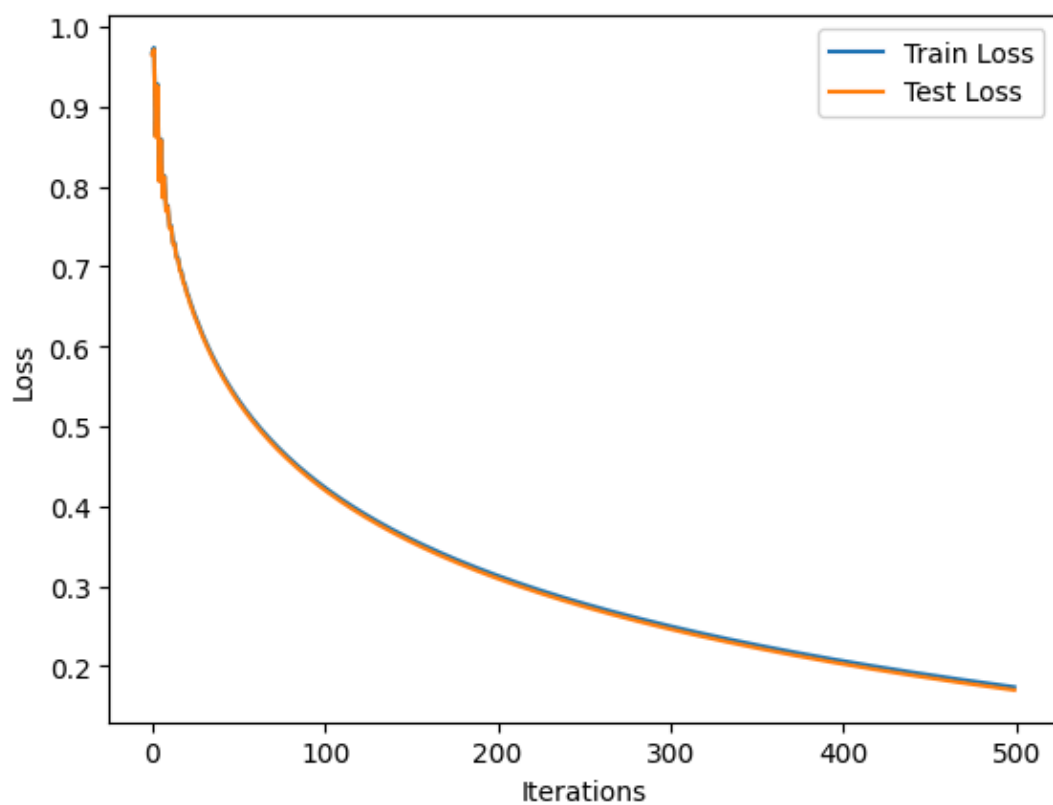
def compute_gradient(X, y, theta):
    s = np.dot(X, theta)
    exp_s = np.exp(s)
    gradient = exp_s.T / np.sum(exp_s, axis=1)
    return np.dot(X.T, (gradient.T - y)) / m

def compute_loss(X, y, theta):
    s = np.dot(X, theta)
    exp_s = np.exp(s)
    gradient = exp_s.T / np.sum(exp_s, axis=1)
    return np.mean(np.log(np.abs(gradient - y)))

def accuracy(X, y, theta):
    preds = np.argmax(np.dot(X, theta), axis=1)
    return np.mean(preds == y)
```

```
stepsize = 1e-5
Theta = np.zeros((n, 10))
train_loss, test_loss, train_misclass, test_misclass = [], [], [], []
for _ in range(500):
    Theta -= stepsize * compute_gradient(Xtrain, ytrain_onehot, Theta)
    train_loss.append(compute_loss(Xtrain, ytrain, Theta))
    test_loss.append(compute_loss(Xtest, ytest, Theta))
    train_misclass.append(1 - accuracy(Xtrain, ytrain, Theta))
    test_misclass.append(1 - accuracy(Xtest, ytest, Theta))

train_loss, test_loss = np.array(train_loss), np.array(test_loss)
train_misclass, test_misclass = np.array(train_misclass), np.array(test_misclass)
plt.plot(train_loss, label="Train Loss")
plt.plot(test_loss, label="Test Loss")
plt.xlabel('Iterations')
plt.ylabel('Loss')
plt.legend()
plt.show()
plt.plot(train_misclass, label="Train Misclass")
plt.plot(test_misclass, label="Test Misclass")
plt.xlabel('Iterations')
plt.ylabel('Misclassification')
plt.legend()
plt.show()
```



2. (a) Total = 5 + 100 + 30 + 3 + 2 = 140

$$H(X) = -\left(\frac{5}{140} \log_2 \frac{5}{140} + \frac{100}{140} \log_2 \frac{100}{140} + \frac{30}{140} \log_2 \frac{30}{140} + \frac{3}{140} \log_2 \frac{3}{140} + \frac{2}{140} \log_2 \frac{2}{140}\right)$$

$$\approx -(-0.1717 - 0.3467 - 0.4762 - 0.1188 - 0.0876)$$

$$\approx 1.2$$

(b) Top drawer = nails + zip tie = 100 + 30 = 130

Bottom drawer = rejt = 5 + 3 + 2 = 10

$$\log_2 P(X=x | Y=\text{bottom})$$

$$H(X|Y) = -\left[P(\text{top}) \cdot \sum_{x \in X} P(X=x | Y=\text{top}) \log_2 P(X=x | Y=\text{top}) + P(\text{bottom}) \sum_{x \in X} P(X=x | Y=\text{bottom}) \log_2 P(X=x | Y=\text{bottom})\right]$$

$$= -\left[0.5 \left(\frac{100}{130} \log_2 \frac{100}{130} + \frac{30}{130} \log_2 \frac{30}{130}\right) + 0.5 \left(\frac{5}{10} \log_2 \frac{5}{10} + \frac{3}{10} \log_2 \frac{3}{10} + \frac{2}{10} \log_2 \frac{2}{10}\right)\right]$$

$$\approx -[0.5 \cdot (-0.78) + 0.5 \cdot (-1.48)]$$

$$\approx 1.13$$

(c) $I(X;Y) = 1.2 - 1.13 = 0.07$

(d) Top drawer = phillips + flathead + $\frac{1}{2}$ nails = 2 + 3 + $\frac{1}{2} \cdot 100 = 55$

Bottom drawer = rejt = 5 + 3 + $\frac{1}{2} \cdot 100 = 85$

$$\frac{5}{85} \log_2 \frac{5}{85}$$

$$H(X|Y) = -\left[0.5 \left(\frac{2}{55} \log_2 \frac{2}{55} + \frac{3}{55} \log_2 \frac{3}{55} + \frac{50}{55} \log_2 \frac{50}{55}\right) + 0.5 \left(\frac{5}{85} \log_2 \frac{5}{85} + \frac{3}{85} \log_2 \frac{3}{85} + \frac{77}{85} \log_2 \frac{77}{85}\right)\right]$$

$$\approx -[0.5 \cdot (-0.174 - 0.229 - 0.125) + 0.5 \cdot (-0.24 - 0.53 - 0.45)]$$

$$\approx -[0.5 \cdot (-1.748)]$$

$$\approx 0.874$$

$$I(X;Y) = 1.2 - 0.874 = 0.326$$

This one is give the large information gain.

3. (a) i.

$$Y_i = x + Z_i = x + N(0,1)$$

$$N(0,1) = Y_i - \hat{x}$$

$$P(Y_i | x) = \frac{1}{\sqrt{2\pi}} \cdot \exp\left[-\frac{(Y_i - x)^2}{2}\right]$$

$$P(Y_1, \dots, Y_m | x) = \frac{1}{\sqrt{2\pi}} \prod_{i=1}^m \exp\left[-\frac{(Y_i - x)^2}{2}\right]$$

$$\begin{aligned} \text{log-likelihood: } \log P(Y_1, \dots, Y_m | x) &= \log\left(\frac{1}{\sqrt{2\pi}}\right)^m + \sum_{i=1}^m \left(-\frac{(Y_i - x)^2}{2}\right) \\ &= \sum_{i=1}^m \log\left(\frac{1}{\sqrt{2\pi}}\right) - \frac{1}{2} \sum_{i=1}^m (Y_i - x)^2 \\ &= -\frac{m}{2} \log(2\pi) - \frac{1}{2} \sum_{i=1}^m (Y_i - x)^2 \end{aligned}$$

$$\text{MLE for } x: \frac{d \log P(Y_1, \dots, Y_m | x)}{dx} = \sum_{i=1}^m (Y_i - x)$$

$$0 = \sum_{i=1}^m (Y_i - x)$$

$$0 = \sum_{i=1}^m Y_i - mx$$

$$\boxed{\hat{x} = \frac{1}{m} \sum_{i=1}^m Y_i}$$

ii. Bias: $E[\hat{x}] = E\left[\frac{1}{m} \sum_{i=1}^m Y_i\right] = \frac{1}{m} \sum_{i=1}^m E[Y_i] = x$, $\boxed{Bias = E[\hat{x}] - x = 0}$

Variance: $Var(\hat{x}) = Var\left(\frac{1}{m} \sum_{i=1}^m Y_i\right) = \frac{1}{m^2} \sum_{i=1}^m Var(Y_i)$

Since $Y_i = x + N(0,1)$, $x \in \mathbb{R}$, $Z \sim N(0,1)$, $Var(Y_i) = 1$

$$Var(\hat{x}) = \frac{1}{m^2} \sum_{i=1}^m 1 = \boxed{\frac{1}{m}}$$

iii. Bias: Bias(\hat{x}) is unbiased, so $\boxed{Bias(\hat{x}) = 0 \text{ as } m \rightarrow +\infty}$

Variance: $\boxed{Var(\hat{x}) = \frac{1}{m} \rightarrow 0 \text{ as } m \rightarrow +\infty}$

(b) i.
$$-\frac{d}{dx} \left(\frac{1}{m} \sum_{i=1}^m (y_i - x)^2 + p(x - \bar{x})^2 \right) = 0$$

$$-\frac{2}{m} \sum_{i=1}^m (y_i - x) + 2p(x - \bar{x}) = 0$$

$$-\frac{2}{m} \sum_{i=1}^m y_i + 2x + 2px - 2p\bar{x} = 0$$

$$(2+2p)x = \frac{2}{m} \sum_{i=1}^m y_i + 2p\bar{x}$$

$$x_{\text{MAP}} = \frac{1}{1+p} \cdot \left(\frac{1}{m} \sum_{i=1}^m y_i + p\bar{x} \right)$$

ii. Bias: $E[x_{\text{MAP}}] = E\left[\frac{1}{1+p} \cdot \left(\frac{1}{m} \sum_{i=1}^m y_i + p\bar{x} \right) \right]$

$$= \frac{1}{1+p} \cdot (E\left[\frac{1}{m} \sum_{i=1}^m y_i \right] + p\bar{x})$$

$$= \frac{1}{1+p} \cdot (E[\bar{x}] + p\bar{x})$$

$$= \frac{1}{1+p} \cdot (x + p\bar{x})$$

$$\text{Bias}(x_{\text{MAP}}) = E[x_{\text{MAP}}] - x = \frac{x + p\bar{x}}{1+p} - \frac{x + px}{1+p} = \frac{p(\bar{x} - x)}{1+p}$$

Variance: $\text{Var}(x_{\text{MAP}}) = \text{Var}\left(\frac{1}{1+p} \cdot \left(\frac{1}{m} \sum_{i=1}^m y_i + p\bar{x} \right) \right)$

$$= \frac{1}{(1+p)^2} \cdot \text{Var}\left(\frac{1}{m} \sum_{i=1}^m y_i \right)$$

$$= \frac{1}{(1+p)^2} \cdot \text{Var}(\bar{x})$$

$$= \frac{1}{(1+p)^2} \cdot \frac{1}{m}$$

$$= \frac{1}{m(1+p)^2}$$

iii. Bias: $\text{Bias}(x_{\text{MAP}}) \neq 0$, it's biased, it's bias

$$\text{Bias}(x_{\text{MAP}}) = \frac{p(\bar{x} - x)}{1+p} \text{ as } m \rightarrow \infty$$

Variance: $\text{Var}(x_{\text{MAP}}) = \frac{1}{m(1+p)^2} \rightarrow 0 \text{ as } m \rightarrow \infty$

$$\begin{aligned} \text{C)} \quad \text{MSE} &= E[(x - \hat{x})^2] = B^2 + V \\ &= \left(\frac{p(\bar{x} - x)}{1+p} \right)^2 + \frac{1}{m(1+p)^2} \end{aligned}$$

$$\begin{aligned} B &= \frac{1}{1+p} \quad p = \frac{1}{B} - 1 \\ \text{MSE} &= \left(\frac{1-B}{B} \cdot \frac{1}{1+B-1} \cdot (\bar{x} - x) \right)^2 + \frac{1}{m(1+\frac{1}{B}-1)^2} \\ &= ((1-B) \cdot (\bar{x} - x))^2 + \frac{B^2}{m} \\ &= (1-B)^2 \Delta^2 + \frac{B^2}{m} \end{aligned}$$

Minimize the MSE:

$$\frac{d(\text{MSE})}{dB} = -2(1-B)\Delta^2 + \frac{2B}{m}$$

$$0 = -2(1-B)\Delta^2 + \frac{2B}{m}$$

$$0 = -2\Delta^2 + 2B\Delta^2 + \frac{2B}{m}$$

$$2\Delta^2 = 2B(\Delta^2 + \frac{1}{m})$$

$$B = \frac{2\Delta^2}{2(\Delta^2 + \frac{1}{m})} = \frac{\Delta^2}{\Delta^2 + \frac{1}{m}}$$

$$p = \frac{1}{B} - 1 = \frac{\Delta^2 + \frac{1}{m}}{\Delta^2} - \frac{\Delta^2}{\Delta^2} = \frac{\frac{1}{m}}{\Delta^2} = \boxed{\frac{1}{m\Delta^2}}$$

4. (a)

$$\min_Q \|XQ - Y\|_2^2 + P\|Q\|_2^2$$

$$\frac{d(\|XQ - Y\|_2^2 + P\|Q\|_2^2)}{dQ} = 2X^T(XQ - Y) + 2PQ$$

$$0 = 2X^T(XQ - Y) + 2PQ$$

$$0 = 2X^T XQ - 2X^T Y + 2PQ$$

$$Q(2X^T X + 2P) = 2X^T Y$$

$$Q = (2X^T X + 2P)^{-1} 2X^T Y$$

$$\boxed{Q = (X^T X + P)^{-1} X^T Y}$$

(b)

$$Y_i = X_i^T \bar{Q} + Z_i, \quad Z_i \sim N(0, 1)$$

$$E[Q] = E[(X^T X + P)^{-1} X^T Y]$$

$$= E[(X^T X + P)^{-1} X^T (X^T \bar{Q} + Z)]$$

$$= E[(X^T X + P)^{-1} X^T X \bar{Q} + (X^T X + P)^{-1} X^T Z]$$

$$= (X^T X + P)^{-1} X^T X \bar{Q} + (X^T X + P)^{-1} X^T E[Z]$$

$$\text{Since } Z_i \sim N(0, 1), E[Z] = 0$$

$$\boxed{E[Q] = (X^T X + P)^{-1} X^T X \bar{Q}}$$

(c)

$$\text{Bias}(Q) = E[Q] - \bar{Q}$$

$$\boxed{= ((X^T X + P)^{-1} X^T X - I) \bar{Q}}$$

(d)

$$\Sigma = \text{Cov}(Q) = E[(Q - E[Q])(Q - E[Q])^T]$$

$$= E[(X^T X + P)^{-1} X^T Z (X^T X + P)^{-1} X^T Z^T]$$

$$= E[(X^T X + P)^{-1} X^T Z Z^T X (X^T X + P)^{-1}]$$

$$= (X^T X + P)^{-1} X^T E[Z Z^T] X (X^T X + P)^{-1}$$

$$\text{Since } Z_i \sim N(0, 1), E[Z Z^T] = I$$

$$\boxed{\Sigma = (X^T X + P)^{-1} X^T X (X^T X + P)^{-1}}$$

(e)

$$\text{Bias: } \text{Bias}(Q) = ((X^T X + P)^{-1} X^T X - I) \bar{Q}$$

$$\text{As } P \rightarrow 0: \text{Bias}(Q) \approx (I - I) \bar{Q} \approx 0$$

$$\text{As } P \rightarrow \infty: \text{Bias}(Q) \approx (0 - I) \bar{Q} \approx -\bar{Q}$$

$$\text{As } P \rightarrow \infty: \text{Bias}(Q) \approx (I - I) \bar{Q} \approx 0$$

} Bias(Q) reduce

} Bias(Q) reduce

$$\boxed{\| \text{Bias}(Q) \|_2 = O\left(\frac{P}{P+m}\right)}$$

Σ :

$$\Sigma = (X^T X + P)^{-1} X^T X (X^T X + P)^{-1}$$

$$\text{As } P \rightarrow 0: \Sigma \approx (X^T X)^{-1}$$

$$\text{As } P \rightarrow \infty: \Sigma \approx 0$$

} Σ reduce

$$\text{As } P \rightarrow \infty: \Sigma \approx (X^T X)^{-1} \approx 0$$

} Σ reduce

$$\boxed{\| \Sigma \|_2 = O\left(\frac{1}{P+m}\right)}$$

5. (a) i. $L(\theta, s, \alpha, u) = \frac{1}{2} \|\theta\|_2^2 + p \sum_{i=1}^m s_i - \sum_{i=1}^m u_i (y_i \phi(x_i)^T \theta - 1 + s_i) - \sum_{i=1}^m V_i s_i$

$$= \frac{1}{2} \|\theta\|_2^2 + p \sum_{i=1}^m s_i - \sum_{i=1}^m u_i (y_i \phi(x_i)^T \theta - 1 + s_i) - \sum_{i=1}^m V_i s_i$$

ii. $\frac{dL(\theta, s, \alpha, u)}{d\theta} = \theta - \sum_{i=1}^m u_i y_i \phi(x_i)$

$$0 = \theta - \sum_{i=1}^m u_i y_i \phi(x_i)$$

$$\theta = \sum_{i=1}^m u_i y_i \phi(x_i)$$

$$L(\theta, u, V) = \frac{1}{2} \left\| \sum_{i=1}^m u_i y_i \phi(x_i) \right\|_2^2 - \sum_{i=1}^m u_i y_i \phi(x_i)^T \sum_{j=1}^m u_j y_j \phi(x_j) + \sum_{i=1}^m u_i$$

$$= \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m u_i u_j y_i y_j \phi(x_i)^T \phi(x_j) - \sum_{i=1}^m u_i y_i y_i \phi(x_i)^T \phi(x_i) + \sum_{i=1}^m u_i$$

$$= \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m u_i u_j y_i y_j \phi(x_i)^T \phi(x_j) + u^T \mathbf{1}$$

$$= \frac{1}{2} \left\| \sum_{i=1}^m u_i y_i \phi(x_i) \right\|_2^2 + u^T \mathbf{1}$$

$$\frac{dL(\theta, s, \alpha, u)}{ds} = p - \sum_{i=1}^m u_i - \sum_{i=1}^m V_i$$

$$0 = p - u - V$$

$$p = u + V$$

Since $u \geq 0$ and $V \geq 0$ and $p = u + V$, get $0 \leq u \leq p$

Dual: $\max_{u, V} -\frac{1}{2} \left\| \sum_{i=1}^m u_i y_i \phi(x_i) \right\|_2^2 + u^T \mathbf{1}$
subject to $0 \leq u \leq p$

iii.

Dual: $\max_{u, V} -\frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m u_i u_j y_i y_j K_{ij}$
subject to $0 \leq u \leq p$

iv.

$$\nabla_u \text{Dual} = \left(-\sum_{j=1}^m u_j y_j K_{ij} + 1 \right)_{i=1}^m$$

(b).

```
def fit_data(X,y,Kfun,rho,maxiter=1000):
    K = Kfun(X,X)
    def get_grad(u):
        return np.ones_like(u) - (y * np.dot(u * y, K))

    u = np.ones(X.shape[0])
    for _ in range(maxiter):
        u = np.clip(u - 0.01 * get_grad(u), 0, rho)

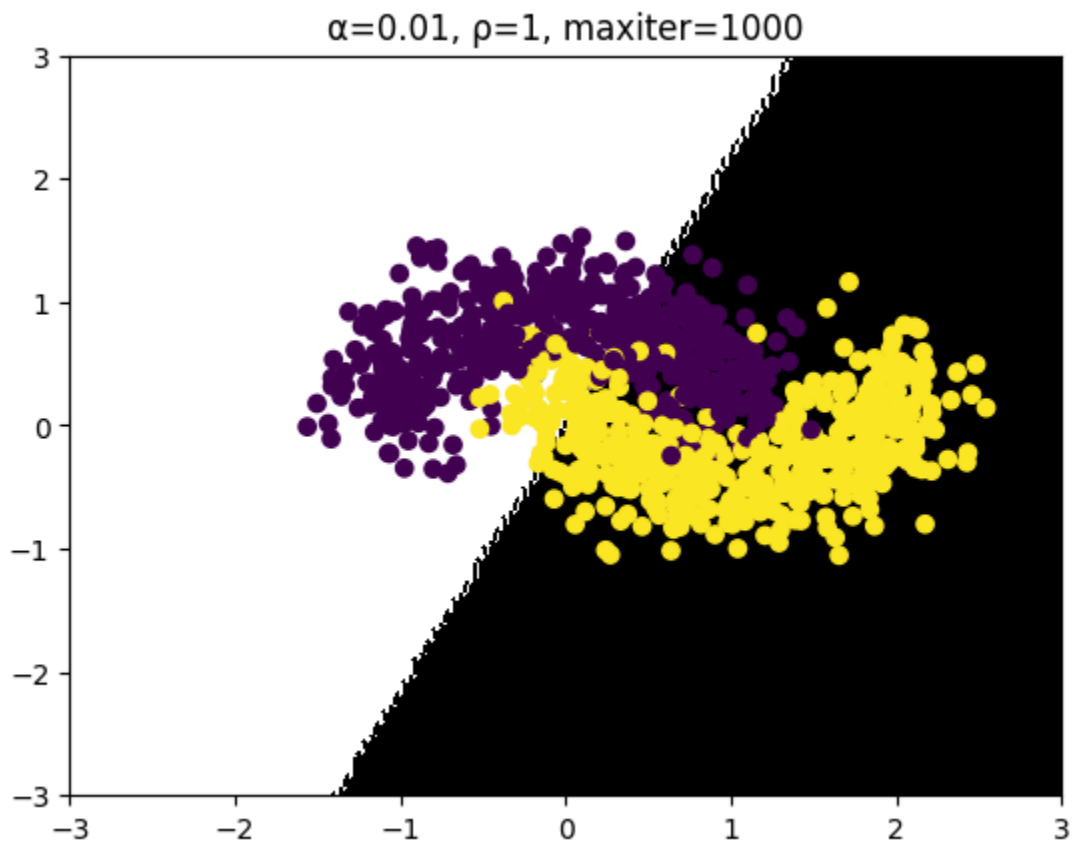
    def inference(x):
        Kt = Kfun(x,X)
        return np.sign(np.dot(y * u, Kt.T))
    return inference

def LinearKfun(X1,X2):
    return np.dot(X1,X2.T)

rhos = [0.01, 0.1, 1, 10]
for rho in rhos:
    inference = fit_data(X, y, LinearKfun, rho)
    print(f"rho = {rho}, Train error: {np.mean(inference(X) != y):.4f},
Test error: {np.mean(inference(Xt) != yt):.4f}")
rho = 0.01, Train error: 0.2580, Test error: 0.2900
rho = 0.1, Train error: 0.2540, Test error: 0.2900
rho = 1, Train error: 0.2530, Test error: 0.2900
rho = 10, Train error: 0.2530, Test error: 0.2900
```

(c).

```
inference = fit_data(X, y, LinearKfun, 1)
Xp, Yp, Zp = plot_model_contour(inference)
plt.contourf(Xp, Yp, Zp, colors=['#FFFFFF', '#000000'])
plt.scatter(*X.T, c=y)
plt.title(" $\alpha=0.01$ ,  $\rho=1$ , maxiter=1000")
plt.show()
```



(d).

```
def fit_data(X,y,Kfun,sigma,rho,maxiter=1000):
    K = Kfun(X,X,sigma)
    def get_grad(u):
        return np.ones_like(u) - (y * np.dot(u * y, K))

    u = np.ones(X.shape[0])
    for _ in range(maxiter):
        u = np.clip(u - 0.01 * get_grad(u), 0, rho)

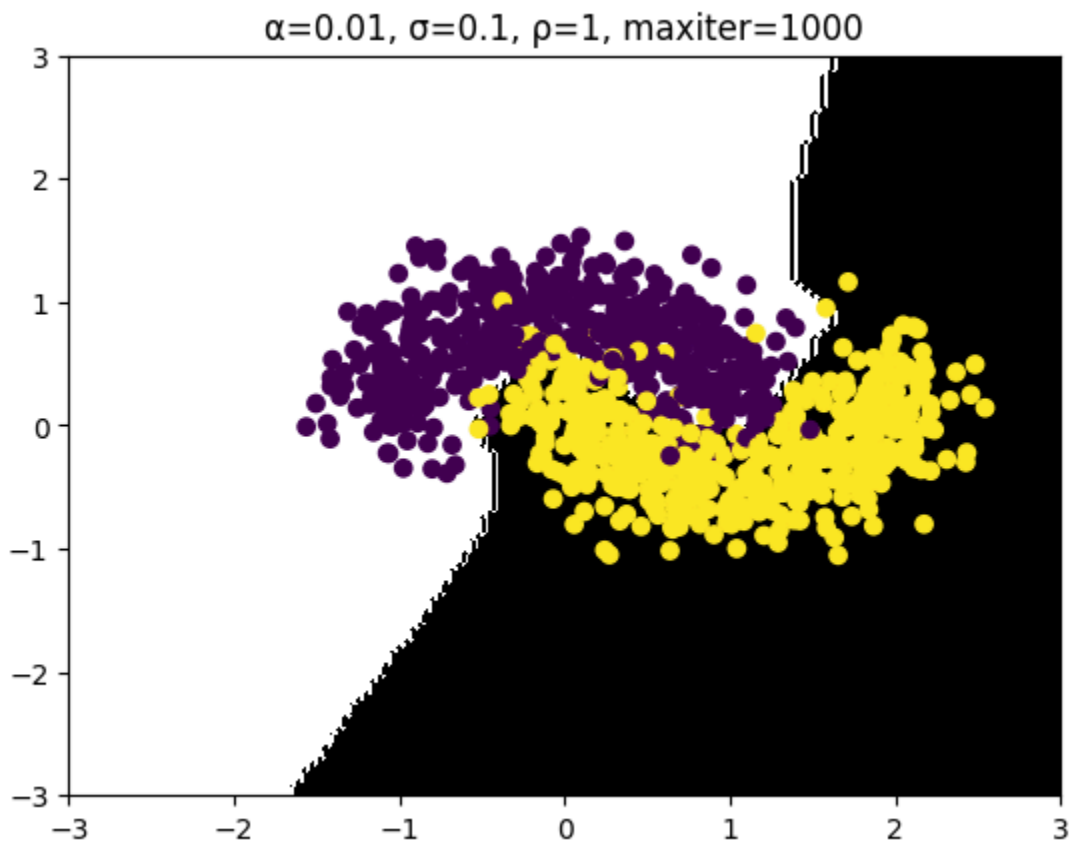
    def inference(x):
        Kt = Kfun(x,X,sigma)
        return np.sign(np.dot(y * u, Kt.T))
    return inference

def RBFKfun(X1,X2,sigma):
    sq = np.sum(X1**2, axis=1).reshape(-1, 1) - 2 * np.dot(X1, X2.T) +
    np.sum(X2**2, axis=1)
    return np.exp(- sq / (2 * sigma**2))

sigmas = [1, 0.1, 0.001, 0.0001]
for sigma in sigmas:
    inference = fit_data(X, y, RBFKfun, sigma, 1)
    print(f"Sigma: {sigma}, Train error: {np.mean(inference(X) != y):.4f},
Test error: {np.mean(inference(Xt) != yt):.4f}")
Sigma: 1, Train error: 0.1990, Test error: 0.2200
Sigma: 0.1, Train error: 0.0520, Test error: 0.0400
Sigma: 0.001, Train error: 0.0020, Test error: 0.5200
Sigma: 0.0001, Train error: 0.0000, Test error: 0.9700
```


(e).

```
inference = fit_data(X, y, RBFKfun, 0.1, 1)
Xp, Yp, Zp = plot_model_contour(inference)
plt.contourf(Xp, Yp, Zp, colors=['#FFFFFF', '#000000'])
plt.scatter(*X.T, c=y)
plt.title(" $\alpha=0.01$ ,  $\sigma=0.1$ ,  $\rho=1$ , maxiter=1000")
plt.show()
```



(f).

When $\sigma \rightarrow 0$, SVM is too sensitive to a single data point, causing the model to fit the training data very closely but lose generalization ability, leading to overfitting. This is why although Train error $\rightarrow 0$, the Test error is getting bigger and bigger.