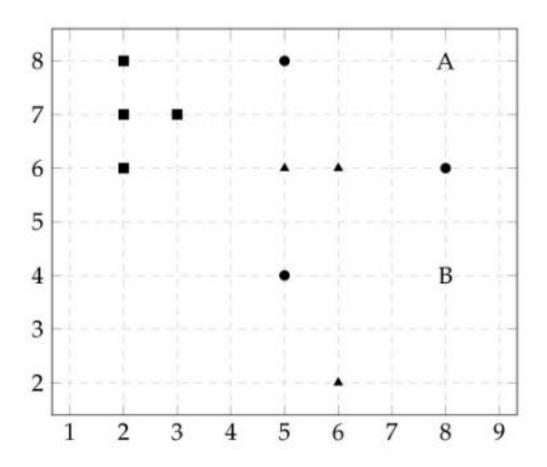
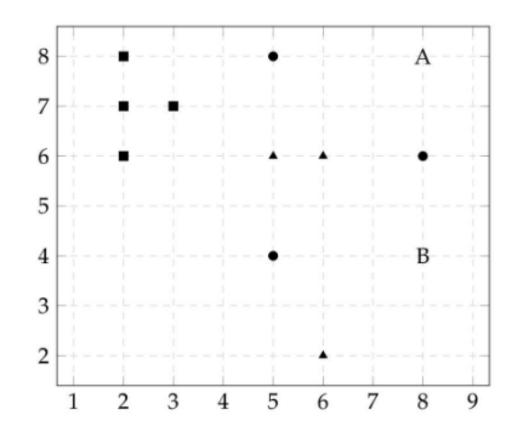
# CSCI 567 Discussion Week 4 - HW1 Review

For the data given below, squares, triangles, and circles are three different classes in the training set, and A and B are two test points with an unknown class. We denote the total number of training points as N (which equals 10) and consider K-nearest-neighbor (KNN) classifier with **L1** distance.





- Three classes:
  - Squares ▲ Triangles Circles
- Two test points:

K-nearest neighbor (KNN) with L1 distance

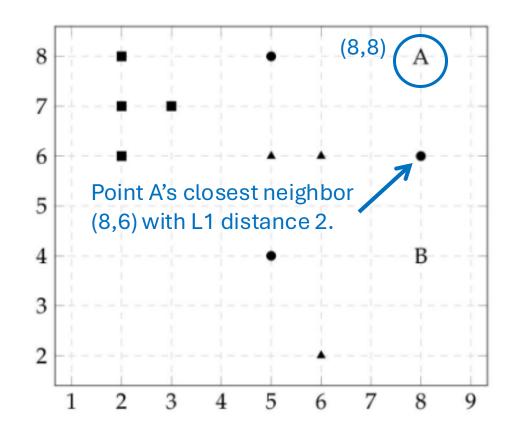
$$d((x1, y1), (x2, y2)) = |x1 - x2| + |y1 - y2|$$

1. What is the test point A classified as for K = 1? Explain briefly.

(2 points)

#### Recall KNN with K=1:

- We only look at the single closest neighbor.
- Whatever class that neighbor belongs to → the predicted class.



- Three classes:
  - Squares ▲ Triangles Circles
- Two test points:

K-nearest neighbor (KNN) with L1 distance

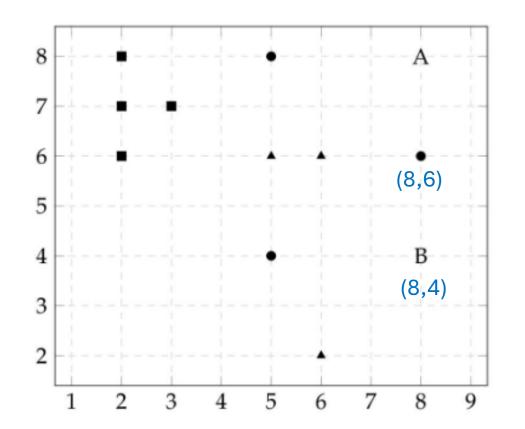
$$d((x1, y1), (x2, y2)) = |x1 - x2| + |y1 - y2|$$

1. What is the test point A classified as for K = 1? Explain briefly.

(2 points)

#### Answer:

- Circle.
- Explanation: The closest point to A is the circle at (8,6) (with L1 distance 2).



- Three classes:
  - Squares ▲ Triangles Circles
- Two test points:

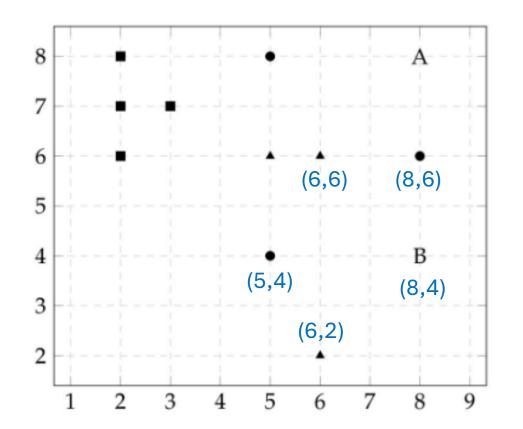
K-nearest neighbor (KNN) with L1 distance

$$d((x1, y1), (x2, y2)) = |x1 - x2| + |y1 - y2|$$

2. What is the smallest odd value of *K* for KNN to predict triangle for the test point B? Explain briefly. (4 points)

K = 1:

- Neighbors:
  - at (8,6) with dist = 2.
- The prediction of B is  $\bullet$ , so K = 1 does not work.



- Three classes:
  - Squares ▲ Triangles Circles
- Two test points:

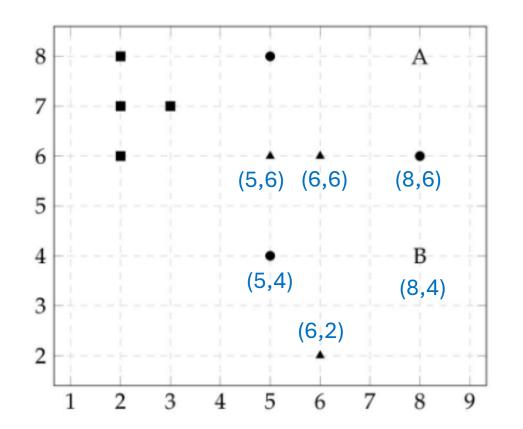
K-nearest neighbor (KNN) with L1 distance

$$d((x1, y1), (x2, y2)) = |x1 - x2| + |y1 - y2|$$

2. What is the smallest odd value of *K* for KNN to predict triangle for the test point B? Explain briefly. (4 points)

K = 3:

- Neighbors: ●, ●, ▲
   at (8,6) with dist = 2; at (5,4) with dist = 3; either ▲ at (6,2) with dist = 4 or ▲ at (6,6) with dist = 4.
- The prediction of B is  $\bullet$ , so K = 3 does not work.



- Three classes:
  - Squares ▲ Triangles Circles
- Two test points:

K-nearest neighbor (KNN) with L1 distance

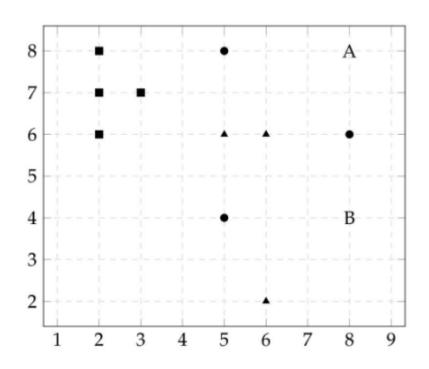
$$d((x1, y1), (x2, y2)) = |x1 - x2| + |y1 - y2|$$

2. What is the smallest odd value of *K* for KNN to predict triangle for the test point B? Explain briefly. (4 points)

K = 5:

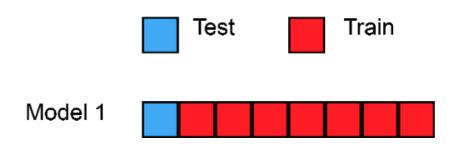
- Neighbors: ●, ●, ▲, ▲, ▲
  - at (8,6) with dist = 2; at (5,4) with dist = 3; at (6,2) with dist = 4; at (6,6) with dist = 4;
  - $\triangle$  at (5,6) with dist = 5.
- The prediction of B is  $\triangle$ , so K = 5 is the answer.

3. Suppose one performs leave-one-out validation (that is, N-fold cross validation) to choose the best hyper-parameter K. List all the points that are misclassified during the N runs when testing the hyper-parameter value K = 1, and report the averaged error rate for this hyper-parameter. (4 points)



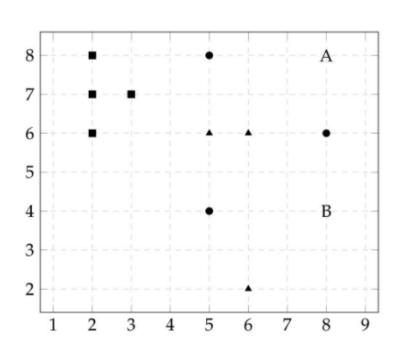
Leave-one-out validation (N-fold cross validation, where N = 10):

- We have N = 10 training points.
- For each **run**  $i = 1 \dots N$ :
  - Leave out 1 point.
  - Classify the point with KNN (K = 1) using the other 9 points.
  - Check if it is correct.



From: https://en.wikipedia.org/wiki/Cross-validation\_(statistics)#/media/File:LOOCV.gif

3. Suppose one performs leave-one-out validation (that is, N-fold cross validation) to choose the best hyper-parameter K. List all the points that are misclassified during the N runs when testing the hyper-parameter value K = 1, and report the averaged error rate for this hyper-parameter. (4 points)



Run	Leave Out	Closest Neighbor (K = 1)	Classification
1	■ at (2,8)	■ at (2,7)	• 📀
2	• at (5,8)	▲ at (5,6)	▲ X
•••	•••		
•••	<ul><li>at (8,6)</li></ul>	▲ at (6,6)	▲ X
•••	• at (5,4)	▲ at (5,6)	▲ X
•••	▲ at (6,2)	<ul><li>at (5,4)</li></ul>	• X
10	•••	•••	

Answer:

The misclassified points:  $\triangle$  at (6,2) and all three  $\bigcirc$  (at (5,8), (8,6), and (5,4)).

The error rate: 4/10 = 0.4

# **Problem 2 Linear Regression**

(24 points)

2.1 (10 points) In the class, we discussed L2 regularized least square solution defined as

$$w_* = \arg\min_{w \in \mathbb{R}^D} \|Xw - y\|_2^2 + \lambda \|w\|_2^2$$
 (1)

where  $X \in \mathbb{R}^{N \times D}$  is the data matrix with each row corresponding to the feature of an example,  $y \in \mathbb{R}^N$  is a vector of all the outcomes,  $\|\cdot\|_2$  stands for the L2 norm, and  $\lambda$  is the regularization coefficient. In this problem, we consider a different regularization method:

$$w'_{*} = \arg\min_{w \in \mathbb{R}^{D}} ||Xw - y||_{2}^{2} + w^{T}Mw$$
 (2)

where  $M \in \mathbb{R}^{D \times D}$  is a sysmetric positive definite matrix.

$$w_* = \arg\min_{w \in \mathbb{R}^D} \|Xw - y\|_2^2 + \lambda \|w\|_2^2$$
 (1)

$$\mathbf{w}_*' = \arg\min_{\mathbf{w} \in \mathbb{R}^D} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \mathbf{w}^T \mathbf{M}\mathbf{w}$$
 (2)

Where,

 $oldsymbol{X} \in \mathbb{R}^{N imes D}$  is the data matrix  $oldsymbol{y} \in \mathbb{R}^N$  is a vector of all the outcomes

 $\lambda$  is the regularization coefficient

 $M \in \mathbb{R}^{D \times D}$  is a sysmetric positive definite matrix.

1. Show that the new method is a generalization of the standard L2 regularization by picking a matrix M such that  $w'_*$  in Eq. (2) equals  $w_*$  in Eq. (1). (2 points)

For example,

$$\|w\|_2^2 = \sum_{j=1}^D w_j^2 = w^ op w_j$$

$$\mathbf{w} = \begin{bmatrix} 2 & 3 & 4 \end{bmatrix}, \mathbf{w}^T = \begin{bmatrix} 2 \\ 3 \\ 4 \end{bmatrix}$$
  
 $\mathbf{w}^T \mathbf{w} = \begin{bmatrix} 2 & 3 & 4 \end{bmatrix} \begin{bmatrix} 2 \\ 3 \\ 4 \end{bmatrix} = 2^2 + 3^2 + 4^2 = \|\mathbf{w}\|_2^2$ 

$$w_* = \arg\min_{w \in \mathbb{R}^D} \|Xw - y\|_2^2 + \lambda \|w\|_2^2$$
 (1) Where,

$$\boldsymbol{w}_{*}' = \arg\min_{\boldsymbol{w} \in \mathbb{R}^{D}} \|\boldsymbol{X}\boldsymbol{w} - \boldsymbol{y}\|_{2}^{2} + \boldsymbol{w}^{T}\boldsymbol{M}\boldsymbol{w}$$
 (2)

 $m{X} \in \mathbb{R}^{N \times D}$  is the data matrix  $m{y} \in \mathbb{R}^N$  is a vector of all the outcomes

 $\lambda$  is the regularization coefficient

 $M \in \mathbb{R}^{D \times D}$  is a sysmetric positive definite matrix.

1. Show that the new method is a generalization of the standard L2 regularization by picking a matrix M such that  $w'_*$  in Eq. (2) equals  $w_*$  in Eq. (1). (2 points)

#### Answer:

To make Eq. (2) equals to Eq. (1), we pick the matrix:

$$M = \lambda I$$

where *I* is the *D* by *D* identify matrix clearly makes  $w^T M w$  equal to  $\lambda ||w||_2^2$ .

$$w^ op M w \ = \ w^ op (\lambda I) w \ = \ \lambda \, w^ op w \ = \ \lambda \|w\|_2^2 \qquad \|w\|_2^2 = \sum_{j=1}^D w_j^2 = w^ op w_j$$

$$w_* = \arg\min_{w \in \mathbb{R}^D} ||Xw - y||_2^2 + \lambda ||w||_2^2$$
 (1)

$$\mathbf{w}_*' = \underset{\mathbf{w} \in \mathbb{R}^D}{\operatorname{arg\,min}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \mathbf{w}^T \mathbf{M}\mathbf{w}$$
 (2)

Where,

 $X \in \mathbb{R}^{N \times D}$  is the data matrix

 $oldsymbol{y} \in \mathbb{R}^N$  is a vector of all the outcomes

 $\lambda$  is the regularization coefficient

 $M \in \mathbb{R}^{D \times D}$  is a sysmetric positive definite matrix.

2. Find the closed form of  $w'_*$  by writing down the gradient of  $F(w) = \|Xw - y\|_2^2 + w^T M w$  and setting it to **0**. (4 points)

$$F(w) = \frac{\|Xw - y\|_2^2 + w^T M w}{\|Xw - y\|_2^2}$$

$$\|Xw-y\|_2^2=(Xw-y)^{ op}(Xw-y)$$
 
$$=\left((Xw)^{ op}-y^{ op}\right)(Xw-y). \qquad \text{using } (AB)^{ op}=B^{ op}A^{ op}$$
 
$$=(Xw)^{ op}(Xw)-(Xw)^{ op}y-y^{ op}(Xw)+y^{ op}y$$

$$(Xw)^{\top}y$$
 is a scalar:  $(N \times D \ D \times 1)^T \ N \times 1 \rightarrow (N \times 1)^T \ N \times 1 \rightarrow 1 \times N \ N \times 1 \rightarrow 1 \times 1$ 

So we have, 
$$(Xw)^{ op}y = \left((Xw)^{ op}y\right)^{ op} = y^{ op}(Xw).$$

$$w_* = \arg\min_{w \in \mathbb{R}^D} \|Xw - y\|_2^2 + \lambda \|w\|_2^2$$
 (1)

$$\mathbf{w}_*' = \underset{\mathbf{w} \in \mathbb{R}^D}{\operatorname{arg\,min}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \mathbf{w}^T \mathbf{M}\mathbf{w}$$
 (2)

Where,

 $m{X} \in \mathbb{R}^{N imes D}$  is the data matrix  $m{y} \in \mathbb{R}^N$  is a vector of all the outcomes

 $\lambda$  is the regularization coefficient

 $M \in \mathbb{R}^{D \times D}$  is a sysmetric positive definite matrix.

2. Find the closed form of  $w'_*$  by writing down the gradient of  $F(w) = \|Xw - y\|_2^2 + w^T M w$  and setting it to **0**. (4 points)

$$F(\boldsymbol{w}) = \|\boldsymbol{X}\boldsymbol{w} - \boldsymbol{y}\|_2^2 + \boldsymbol{w}^T \boldsymbol{M} \boldsymbol{w}$$

$$egin{aligned} \|Xw-y\|_2^2 &= (Xw-y)^ op (Xw-y) \ &= \left((Xw)^ op - y^ op\right) (Xw-y). & ext{using } (AB)^ op &= B^ op A^ op \ &= (Xw)^ op (Xw) - (Xw)^ op y - y^ op (Xw) + y^ op y \ &= (Xw)^ op (Xw) - 2\,y^ op (Xw) + y^ op y & ext{using } (Xw)^ op y = y^ op (Xw) \ &= w^ op X^ op Xw - 2\,y^ op Xw + y^ op y & ext{using } (AB)^ op &= B^ op A^ op y \end{aligned}$$

$$w_* = \arg\min_{w \in \mathbb{R}^D} \|Xw - y\|_2^2 + \lambda \|w\|_2^2$$

(1) Where,

$$\mathbf{w}_*' = \underset{\mathbf{w} \in \mathbb{R}^D}{\operatorname{arg\,min}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \mathbf{w}^T \mathbf{M}\mathbf{w}$$
 (2)

 $oldsymbol{X} \in \mathbb{R}^{N imes D}$  is the data matrix  $oldsymbol{y} \in \mathbb{R}^N$  is a vector of all the outcomes

 $\lambda$  is the regularization coefficient

 $M \in \mathbb{R}^{D \times D}$  is a sysmetric positive definite matrix.

2. Find the closed form of  $w'_*$  by writing down the gradient of  $F(w) = \|Xw - y\|_2^2 + w^T M w$  and setting it to **0**. (4 points)

$$F(\boldsymbol{w}) = \|\boldsymbol{X}\boldsymbol{w} - \boldsymbol{y}\|_2^2 + \boldsymbol{w}^T \boldsymbol{M} \boldsymbol{w}$$

$$F(w) = w^ op X^ op X w \ - \ 2\,y^ op X w \ + \ y^ op y \ + \ w^ op M w.$$

The gradient:

$$abla_w \left[ y^ op y 
ight] = 0$$

$$abla_wig[w^ op X^ op Xwig] = 2X^ op X\,w$$

$$abla_w \left[ -2\, y^ op X w 
ight] = -2\, X^ op y$$

$$abla_w \left[ w^ op M w 
ight] = \ 2 M w$$

$$\frac{\partial \mathbf{x}^{\top} \mathbf{A} \mathbf{x}}{\partial \mathbf{x}} = 2 \mathbf{A} \mathbf{x}$$

$$rac{\partial \mathbf{a}^ op \mathbf{x}}{\partial \mathbf{x}} = \mathbf{a} \qquad y^ op X w = (X^ op y)^ op w$$

A is not a function of x and A is symmetric

a is not a function of x

(09/05 lecture slides, page 52)

$$w_* = \arg\min_{w \in \mathbb{R}^D} \|Xw - y\|_2^2 + \lambda \|w\|_2^2$$

(1) Where,

$$\mathbf{w}_*' = \underset{\mathbf{w} \in \mathbb{R}^D}{\operatorname{arg\,min}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \mathbf{w}^T \mathbf{M}\mathbf{w}$$
 (2)

 $oldsymbol{X} \in \mathbb{R}^{N imes D}$  is the data matrix  $oldsymbol{y} \in \mathbb{R}^N$  is a vector of all the outcomes

 $\lambda$  is the regularization coefficient

 $M \in \mathbb{R}^{D \times D}$  is a sysmetric positive definite matrix.

2. Find the closed form of  $w'_*$  by writing down the gradient of  $F(w) = \|Xw - y\|_2^2 + w^T M w$  and setting it to **0**. (4 points)

$$F(\boldsymbol{w}) = \|\boldsymbol{X}\boldsymbol{w} - \boldsymbol{y}\|_2^2 + \boldsymbol{w}^T \boldsymbol{M} \boldsymbol{w}$$

$$F(w) = w^ op X^ op X w \ - \ 2\, y^ op X w \ + \ y^ op y \ + \ w^ op M w.$$

The gradient:

$$abla F(w) \ = \ 2X^ op X \, w \ - \ 2X^ op y \ + \ 2Mw$$

$$2X^{T}(Xw - y) + 2Mw.$$

$$w_* = \arg\min_{w \in \mathbb{R}^D} \|Xw - y\|_2^2 + \lambda \|w\|_2^2$$

Where,

(1)

$$\mathbf{w}_*' = \underset{\mathbf{w} \in \mathbb{R}^D}{\operatorname{arg\,min}} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \mathbf{w}^T \mathbf{M}\mathbf{w}$$
 (2)

 $m{X} \in \mathbb{R}^{N imes D}$  is the data matrix  $m{y} \in \mathbb{R}^N$  is a vector of all the outcomes

 $\lambda$  is the regularization coefficient

 $M \in \mathbb{R}^{D \times D}$  is a sysmetric positive definite matrix.

2. Find the closed form of  $w'_*$  by writing down the gradient of  $F(w) = \|Xw - y\|_2^2 + w^T M w$  and setting it to **0**. (4 points)

Set it to 0:

$$2X^{T}(Xw - y) + 2Mw.$$

$$2X^\top X \, w - 2X^\top y + 2Mw = 0$$

$$(X^ op X + M)w = X^ op y$$

$$Aw = b$$

Let 
$$\,A := X^ op X + M\,$$
 and  $\,b := X^ op y\,$ 

$$A^{-1}Aw = A^{-1}b$$

$$oldsymbol{w} = oldsymbol{A}^{-1}oldsymbol{b}$$
 , that is,  $oldsymbol{w}_*' = \left(oldsymbol{X}^{\mathrm{T}}oldsymbol{X} + oldsymbol{M}
ight)^{-1}oldsymbol{X}^{\mathrm{T}}oldsymbol{y}.$ 

3. Recall the Newton method:  $w^{(t+1)} \leftarrow w^{(t)} - H_t^{-1} \nabla F(w^{(t)})$  where  $H_t = \nabla^2 F(w^{(t)})$ . Show that no matter what the initialization  $w^{(0)}$  is, Newton method always takes one step only to find the minimizer  $w'_*$  of  $F(w) = ||Xw - y||_2^2 + w^T M w$ . (4 points)

$$\boldsymbol{w}^{(t+1)} \leftarrow \boldsymbol{w}^{(t)} - \boldsymbol{H}_t^{-1} \nabla F(\boldsymbol{w}^{(t)})$$

$$\nabla F(w) = 2X^{T}(Xw - y) + 2Mw$$

At any point, the Hessian is always  $H:=
abla^2F(w)=2(X^TX+M)$ .

$$H^{-1} = \left[2\left(X^ op X + M
ight)
ight]^{-1} = rac{1}{2}\left(X^ op X + M
ight)^{-1}$$

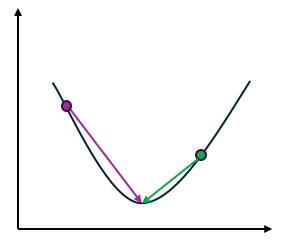
So applying Newton method to any initialization of  $w^{(0)}$  gives,

$$w^{(1)} = w^{(0)} - (X^T X + M)^{-1} (X^T (X w^{(0)} - y) + M w^{(0)})$$

$$= w^{(0)} - (X^T X + M)^{-1} \left( (X^T X + M) w^{(0)} - X^T y \right)$$

$$= \left( X^T X + M \right)^{-1} X^T y = w'_*.$$

Hence Newton method needs only **one step**, regardless of initialization.



Example: Newton method jumps directly to its minimizer in one step from any start.

**2.2 (14 points)** Assume we have a training set  $(x_1, y_1), \ldots, (x_N, y_N) \in \mathbb{R}^D \times \mathbb{R}$ , where each outcome  $y_n$  is generated by a probabilistic model  $\mathbf{w}_*^T \mathbf{x}_n + \epsilon_n$  with  $\epsilon_n$  being an independent Gaussian noise with zero-mean and variance  $\sigma^2$  for some  $\sigma > 0$ . In other words, the probability density of any outcome  $y \in \mathbb{R}$  given  $\mathbf{x}_n$  is

$$\Pr(y \mid \mathbf{x}_n; \mathbf{w}_*, \sigma) = \frac{1}{\sigma \sqrt{2\pi}} \exp\left(\frac{-(y - \mathbf{w}_*^{\mathrm{T}} \mathbf{x}_n)^2}{2\sigma^2}\right).$$

1. Assume  $\sigma$  is fixed and given, find the maximum likelihood estimation for  $w_*$ . In other words, first write down the joint density of the outcomes  $y_1, \ldots, y_N$  given  $x_1, \ldots, x_N$  as a function of the value of  $w_*$ ; then find the value of  $w_*$  that maximizes this density. You can assume  $X^TX$  is invertible where X is the data matrix as used in Problem 2.1.

The probability of seeing label  $y_1, ..., y_n$  given  $x_1, ..., x_n$  is that

$$\mathcal{P}(\boldsymbol{w}) = \prod_{n=1}^{N} \Pr(y_n \mid \boldsymbol{x}_n; \boldsymbol{w}, \sigma) = \prod_{n=1}^{N} \frac{1}{\sigma \sqrt{2\pi}} \exp\left(\frac{-(y_n - \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x}_n)^2}{2\sigma^2}\right).$$

(The joint density for a linear model w)

Now, we need to find a  $w_*$  that maximizes  $\mathcal{P}(w)$ 

1. Assume  $\sigma$  is fixed and given, find the maximum likelihood estimation for  $w_*$ . In other words, first write down the joint density of the outcomes  $y_1, \ldots, y_N$  given  $x_1, \ldots, x_N$  as a function of the value of  $w_*$ ; then find the value of  $w_*$  that maximizes this density. You can assume  $X^TX$  is invertible where X is the data matrix as used in Problem 2.1.

$$\mathcal{P}(\boldsymbol{w}) = \prod_{n=1}^{N} \Pr(y_n \mid \boldsymbol{x}_n; \boldsymbol{w}, \sigma) = \prod_{n=1}^{N} \frac{1}{\sigma \sqrt{2\pi}} \exp\left(\frac{-(y_n - \boldsymbol{w}^T \boldsymbol{x}_n)^2}{2\sigma^2}\right).$$

Taking the negative log, this becomes

From 
$$F(\boldsymbol{w}) = -\ln P(\boldsymbol{w}) = -\sum_{n=1}^{N} \ln \left[ \frac{1}{\sigma \sqrt{2\pi}} \exp\left(-\frac{(y_n - \boldsymbol{w}^{\top} \boldsymbol{x}_n)^2}{2\sigma^2}\right) \right]$$

$$= -\sum_{n=1}^{N} \left[ \ln \left(\frac{1}{\sigma \sqrt{2\pi}}\right) - \frac{(y_n - \boldsymbol{w}^{\top} \boldsymbol{x}_n)^2}{2\sigma^2} \right]$$

$$= -\sum_{n=1}^{N} \left[ -\ln \sigma - \ln \sqrt{2\pi} - \frac{(y_n - \boldsymbol{w}^{\top} \boldsymbol{x}_n)^2}{2\sigma^2} \right]$$

$$=\sum_{n=1}^N ig[\ln\sigma + \ln\sqrt{2\pi}ig] \ + \ rac{1}{2\sigma^2}\sum_{n=1}^N (y_n - w^ op x_n)^2$$

$$\log_b(xy) = \log_b x + \log_b y$$

$$\ln e^t = t$$

$$\log_b \sqrt[p]{x} = rac{\log_b x}{p}$$

$$\log_b \frac{x}{y} = \log_b x - \log_b y$$

$$F(w) = N \ln \sqrt{2\pi} + N \ln \sigma + \frac{1}{2\sigma^2} \sum_{n=1}^{N} (y_n - w^T x_n)^2 = N \ln \sqrt{2\pi} + N \ln \sigma + \frac{1}{2\sigma^2} ||Xw - y||_2^2$$

20

1. Assume  $\sigma$  is fixed and given, find the maximum likelihood estimation for  $w_*$ . In other words, first write down the joint density of the outcomes  $y_1, \ldots, y_N$  given  $x_1, \ldots, x_N$  as a function of the value of  $w_*$ ; then find the value of  $w_*$  that maximizes this density. You can assume  $X^TX$  is invertible where X is the data matrix as used in Problem 2.1. (6 points)

$$\mathcal{P}(\boldsymbol{w}) = \prod_{n=1}^{N} \Pr(y_n \mid \boldsymbol{x}_n; \boldsymbol{w}, \sigma) = \prod_{n=1}^{N} \frac{1}{\sigma \sqrt{2\pi}} \exp\left(\frac{-(y_n - \boldsymbol{w}^T \boldsymbol{x}_n)^2}{2\sigma^2}\right).$$

$$F(w) = N \ln \sqrt{2\pi} + N \ln \sigma + \frac{1}{2\sigma^2} \sum_{n=1}^{N} (y_n - w^T x_n)^2 = N \ln \sqrt{2\pi} + N \ln \sigma + \frac{1}{2\sigma^2} ||Xw - y||_2^2$$

Maximizing  $P(\mathbf{w})$  is the same as minimizing  $F(\mathbf{w})$ , the negative log of  $P(\mathbf{w})$ :

•  $\sigma$  is fixed and given, so it becomes minimizing the last part of  $F(\mathbf{w})$ :  $\sum_{n=1}^{N} (y_n - \mathbf{w}^T x_n)^2$ 

(the same objective as for least square regression)

Similarly,

$$egin{aligned} 
abla_w \|y-Xw\|_2^2 &= 2X^ op(Xw-y) = 0 \ X^ op X w &= X^ op y \ oldsymbol{w}_* &= (oldsymbol{X}^ op oldsymbol{X})^{-1} oldsymbol{X}^ op oldsymbol{y}. \end{aligned}$$

2. Now consider  $\sigma$  as a parameter of the probabilistic model too, that is, the model is specified by both  $w_*$  and  $\sigma$ . Find the maximum likelihood estimation for  $w_*$  and  $\sigma$ . (8 points)

Follow the same steps from the previous question, we can get  $F(\mathbf{w}, \sigma)$ :

$$F(\boldsymbol{w}, \sigma) = N \ln \sqrt{2\pi} + N \ln \sigma + \frac{1}{2\sigma^2} ||\boldsymbol{X}\boldsymbol{w} - \boldsymbol{y}||_2^2.$$

MLE for w\* and  $\sigma$  is the minimizer of the function  $F(w, \sigma)$ .

Take the derivative over w and set it to 0 to find  $w_*$ :

$$\frac{\partial F(w,\sigma)}{\partial w} = \frac{1}{\sigma^2} X^{\top} (Xw - y) = \frac{1}{\sigma^2} (X^{\top} X w - X^{\top} y) \qquad \qquad \boldsymbol{w}_* = (\boldsymbol{X}^{\top} \boldsymbol{X})^{-1} \boldsymbol{X}^{\top} \boldsymbol{y}.$$
(Does not depend on  $\sigma$ )

Take the derivative over  $\sigma$ :

$$rac{\partial F(w,\sigma)}{\partial \sigma} \ = \ rac{N}{\sigma} - rac{\|Xw-y\|^2}{\sigma^3}.$$

$$N\sigma^2 - \|Xw - y\|^2 = 0$$
  $\sigma > 0.$   $\sigma = \sqrt{\frac{\|Xw - y\|^2}{N}}$ 

This depends on w, so we need to find  $w_*$  first.

2. Now consider  $\sigma$  as a parameter of the probabilistic model too, that is, the model is specified by both  $w_*$  and  $\sigma$ . Find the maximum likelihood estimation for  $w_*$  and  $\sigma$ . (8 points)

$$F(\boldsymbol{w},\sigma) = N \ln \sqrt{2\pi} + N \ln \sigma + \frac{1}{2\sigma^2} \|\boldsymbol{X}\boldsymbol{w} - \boldsymbol{y}\|_2^2. \qquad \sigma > 0.$$

We first fix  $\sigma$  and minimize over w (the same MLE from the previous question):

$$\boldsymbol{w}_* = (\boldsymbol{X}^{\mathrm{T}} \boldsymbol{X})^{-1} \boldsymbol{X}^{\mathrm{T}} \boldsymbol{y}.$$

For any  $\sigma > 0$ ,

$$N \ln \sqrt{2\pi} + N \ln \sigma + \frac{1}{2\sigma^2} |Xw - y|^2 \ge N \ln \sqrt{2\pi} + N \ln \sigma + \frac{1}{2\sigma^2} |Xw - y|^2$$

Then we plug in  $\mathbf{w}_*$  and minimize  $F(\mathbf{w}_*, \sigma)$  to find  $\sigma$ :

$$\frac{\partial F(\boldsymbol{w}_*, \sigma)}{\partial \sigma} = \frac{N}{\sigma} - \frac{1}{\sigma^3} \|\boldsymbol{X} \boldsymbol{w}_* - \boldsymbol{y}\|_2^2 = 0.$$

$$\sigma = \frac{1}{\sqrt{N}} \|\boldsymbol{X} \boldsymbol{w}_* - \boldsymbol{y}\|_2 = \frac{1}{\sqrt{N}} \|\boldsymbol{X} (\boldsymbol{X}^T \boldsymbol{X})^{-1} \boldsymbol{X}^T \boldsymbol{y} - \boldsymbol{y}\|_2.$$

In Lecture 3 we have seen the hinge loss  $\ell(z) = \max\{0, 1-z\}$ , which is non-differentiable at z=1. To avoid this issue, we can consider the square of hinge loss  $\ell(z)^2$ , which is differentiable everywhere. More specifically, given a binary dataset  $(x_1, y_1), \ldots, (x_N, y_N) \in \mathbb{R}^D \times \{-1, 1\}$ , we define the following new loss function for a linear model  $w \in \mathbb{R}^D$ :

$$F(\boldsymbol{w}) = \frac{1}{N} \sum_{n=1}^{N} F_n(\boldsymbol{w}), \quad \text{where } F_n(\boldsymbol{w}) = \left( \max \left\{ 0, 1 - y_n \boldsymbol{w}^{\mathsf{T}} \boldsymbol{x}_n \right\} \right)^2. \tag{3}$$

1. For a fixed n, write down the gradient  $\nabla F_n(w)$  (show your derivation), then fill in the missing details in the repeat-loop of the algorithm below which applies SGD to minimize F. (6 points)

When 
$$1 - y_n w^T x_n < 0$$
,  $F_n(w) = \left( \max \left\{ 0, 1 - y_n w^T x_n \right\} \right)^2 = 0$  The gradient of it is 0.

When 
$$1 - y_n w^T x_n > 0$$
,  $F_n(w) = (1 - y_n w^T x_n)^2$ 

Take the derivative by using the chain rule:

$$egin{aligned} rac{d}{dz}z^2 &= 2z & 
abla_w F_n(w) &= 2(1-y_n w^ op x_n) \cdot 
abla_w (1-y_n w^ op x_n) &= 2(1-y_n w^ op x_n)(-y_n x_n) \ 
abla F_n(w) &= -2y_n (1-y_n w^ op x_n) x_n, \end{aligned}$$

## **Problem 3** Linear Classifiers

# (16 points)

In Lecture 3 we have seen the hinge loss  $\ell(z) = \max\{0, 1-z\}$ , which is non-differentiable at z = 1. To avoid this issue, we can consider the square of hinge loss  $\ell(z)^2$ , which is differentiable everywhere. More specifically, given a binary dataset  $(x_1, y_1), \ldots, (x_N, y_N) \in \mathbb{R}^D \times \{-1, 1\}$ , we define the following new loss function for a linear model  $w \in \mathbb{R}^D$ :

$$F(\boldsymbol{w}) = \frac{1}{N} \sum_{n=1}^{N} F_n(\boldsymbol{w}), \quad \text{where } F_n(\boldsymbol{w}) = \left( \max \left\{ 0, 1 - y_n \boldsymbol{w}^{\mathsf{T}} \boldsymbol{x}_n \right\} \right)^2. \tag{3}$$

1. For a fixed n, write down the gradient  $\nabla F_n(w)$  (show your derivation), then fill in the missing details in the repeat-loop of the algorithm below which applies SGD to minimize F. (6 points)

When 
$$1 - y_n w^T x_n < 0$$
,  $F_n(w) = \left( \max \left\{ 0, 1 - y_n w^T x_n \right\} \right)^2 = 0$  The gradient of it is 0.

When 
$$1 - y_n w^T x_n > 0$$
,  $F_n(w) = (1 - y_n w^T x_n)^2$   $\max\{0, 1 - y_n w^T x_n\} = \begin{cases} 0 & \text{if } 1 - y_n w^T x_n \leq 0, \\ 1 - y_n w^T x_n & \text{if } 1 - y_n w^T x_n > 0. \end{cases}$   $\nabla F_n(w) = -2y_n(1 - y_n w^T x_n) x_n$ 

which is also **0** when  $y_n w^T x_n$  approaches 1

$$\nabla F_n(\boldsymbol{w}) = -2y_n \max\{0, 1 - y_n \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x}_n\} \boldsymbol{x}_n.$$

## **Problem 3** Linear Classifiers

# (16 points)

In Lecture 3 we have seen the hinge loss  $\ell(z) = \max\{0, 1-z\}$ , which is non-differentiable at z = 1. To avoid this issue, we can consider the square of hinge loss  $\ell(z)^2$ , which is differentiable everywhere. More specifically, given a binary dataset  $(x_1, y_1), \ldots, (x_N, y_N) \in \mathbb{R}^D \times \{-1, 1\}$ , we define the following new loss function for a linear model  $w \in \mathbb{R}^D$ :

$$F(\boldsymbol{w}) = \frac{1}{N} \sum_{n=1}^{N} F_n(\boldsymbol{w}), \quad \text{where } F_n(\boldsymbol{w}) = \left( \max \left\{ 0, 1 - y_n \boldsymbol{w}^{\mathsf{T}} \boldsymbol{x}_n \right\} \right)^2. \tag{3}$$

1. For a fixed n, write down the gradient  $\nabla F_n(w)$  (show your derivation), then fill in the missing details in the repeat-loop of the algorithm below which applies SGD to minimize F. (6 points)

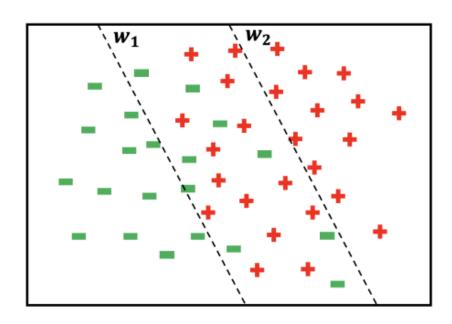
## **Algorithm 1:** SGD for minimizing Eq. (3)

- 1 **Input:** A training set  $(x_1, y_1), \ldots, (x_N, y_N) \in \mathbb{R}^D \times \{-1, 1\}$ , learning rate  $\eta > 0$
- 2 Initialization: w = 0
- 3 Repeat:
- 4 | randomly pick an example  $(x_n, y_n)$
- 5 update  $\mathbf{w} \leftarrow \mathbf{w} + 2\eta y_n \max\{0, 1 y_n \mathbf{w}^{\mathrm{T}} \mathbf{x}_n\} \mathbf{x}_n$

$$\nabla F_n(\boldsymbol{w}) = -2y_n \max\{0, 1 - y_n \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x}_n\} \boldsymbol{x}_n.$$

$$F_n(\boldsymbol{w}) = \begin{cases} \left( \max\left\{0, 1 - \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x}_n \right\} \right)^2, & \text{if } y_n = 1, \\ \underline{0.1} \left( \max\left\{0, 1 + \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x}_n \right\} \right)^2, & \text{else.} \end{cases}$$
(4)

(a) Consider a binary classification dataset of points in two dimensions as shown in Figure 1, where the red, plus signs denote samples with label +1, and the green, minus signs denote samples with label -1. When training a linear classifier with the modified loss in Eq. (4), which of  $w_1$  or  $w_2$  in Figure 1 do you think is more likely the resulting decision boundary? Explain briefly. (2 points)

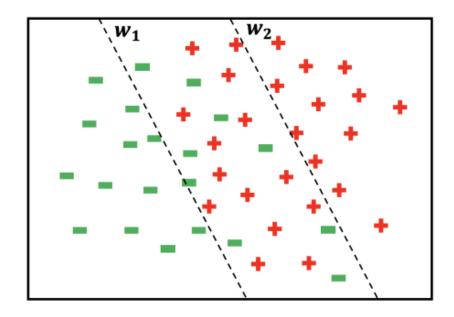


- The loss function now does not penalize misclassifying negative points as heavily as it penalizes misclassifying positive points.
- => The model does not want to misclassify + points.

Figure 1: A binary classification task

$$F_n(\boldsymbol{w}) = \begin{cases} \left( \max\left\{0, 1 - \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x}_n \right\} \right)^2, & \text{if } y_n = 1, \\ \underline{0.1} \left( \max\left\{0, 1 + \boldsymbol{w}^{\mathrm{T}} \boldsymbol{x}_n \right\} \right)^2, & \text{else.} \end{cases}$$
(4)

(a) Consider a binary classification dataset of points in two dimensions as shown in Figure 1, where the red, plus signs denote samples with label +1, and the green, minus signs denote samples with label -1. When training a linear classifier with the modified loss in Eq. (4), which of  $w_1$  or  $w_2$  in Figure 1 do you think is more likely the resulting decision boundary? Explain briefly. (2 points)



- The dataset is not linearly separable.
  - Any linear classifier will make mistakes.
- Therefore,  $w_1$  is more likely.

Figure 1: A binary classification task

$$F_n(w) = \begin{cases} \left( \max \left\{ 0, 1 - w^{\mathsf{T}} x_n \right\} \right)^2, & \text{if } y_n = 1, \\ 0.1 \left( \max \left\{ 0, 1 + w^{\mathsf{T}} x_n \right\} \right)^2, & \text{else.} \end{cases}$$
(4)

(b) Based on your answer from the last question, give an example where one would want to modify the loss function in such a way. (2 points)

For example, in fraud detection (+1 represents a fraud), it is far more important to make sure that an actual fraud is not missed, and thus the loss function should give more weights to positive labels.

We accept any reasonable example where the one label represents a critical or high-stakes outcome and is far more important than the other label.

$$F_n(w) = \begin{cases} \left( \max \left\{ 0, 1 - w^{\mathsf{T}} x_n \right\} \right)^2, & \text{if } y_n = 1, \\ 0.1 \left( \max \left\{ 0, 1 + w^{\mathsf{T}} x_n \right\} \right)^2, & \text{else.} \end{cases}$$
(4)

(c) Similarly to Question 3.1, write down the gradient of this modified loss  $F_n$ , then fill in the missing details in the repeat-loop of the algorithm below which applies SGD to minimize F. (6 points)

Similar steps as in Question 3.1,

When  $y_n = 1$ ,

- If  $1-w^{\top}x_n \leq 0$  , the gradient is 0.
- If  $1-w^{\top}x_n>0$  ,

$$abla F_n(w) = 2(1-w^ op x_n) \, 
abla (1-w^ op x_n) = 2(1-w^ op x_n)(-x_n) = -2 \, \max\{0, 1-w^ op x_n\} \, x_n$$

Else,

- If  $1+w^{\top}x_n\leq 0$ , the gradient is 0.
- If  $1+w^{ op}x_n>0$  ,

$$abla F_n(w) = 0.1 \cdot 2(1 + w^ op x_n) \, 
abla (1 + w^ op x_n) = 0.2 \, \max\{0, 1 + w^ op x_n\} \, x_n.$$

$$F_n(w) = \begin{cases} \left( \max \left\{ 0, 1 - w^{\mathsf{T}} x_n \right\} \right)^2, & \text{if } y_n = 1, \\ 0.1 \left( \max \left\{ 0, 1 + w^{\mathsf{T}} x_n \right\} \right)^2, & \text{else.} \end{cases}$$
(4)

(c) Similarly to Question 3.1, write down the gradient of this modified loss  $F_n$ , then fill in the missing details in the repeat-loop of the algorithm below which applies SGD to minimize F. (6 points)

$$\nabla F_n(w) = \begin{cases} -2 \max\{0, 1 - w^{\mathsf{T}} x_n\} x_n, & \text{if } y_n = 1, \\ 0.2 \max\{0, 1 + w^{\mathsf{T}} x_n\} x_n & \text{else.} \end{cases}$$

### **Algorithm 2:** SGD for minimizing modified loss Eq. (4)

- 1 **Input:** A training set  $(x_1, y_1), \ldots, (x_N, y_N) \in \mathbb{R}^D \times \{-1, 1\}$ , learning rate  $\eta > 0$
- 2 Initialization: w = 0
- 3 Repeat:
- randomly pick an example  $(x_n, y_n)$ update  $w \leftarrow w - \begin{cases} -2\eta \max\{0, 1 - w^T x_n\} x_n, & \text{if } y_n = 1, \\ 0.2\eta \max\{0, 1 + w^T x_n\} x_n & \text{else.} \end{cases}$