

Homework 1 Review

Chung-Wei Lee

September 16

Outline

1 Problem 1

- (a)
- (b)
- (c)

2 Problem 2

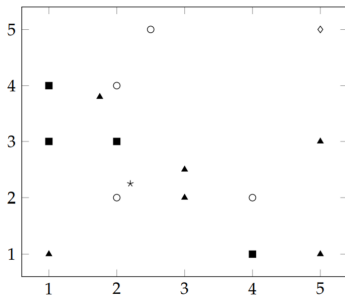
- 2.1
- 2.2.1
- 2.2.2

3 Problem 3

- 3.1
- 3.2
- 3.3
- 3.4

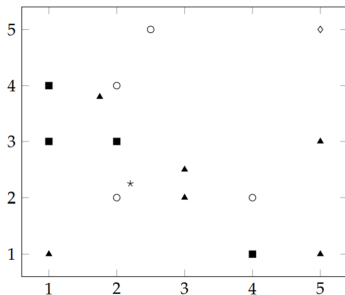
Problem 1

We denote the total number of training points as N and consider K -nearest-neighbor (KNN) classifier with L2 distance.



Problem 1

We denote the total number of training points as N and consider K -nearest-neighbor (KNN) classifier with L2 distance.

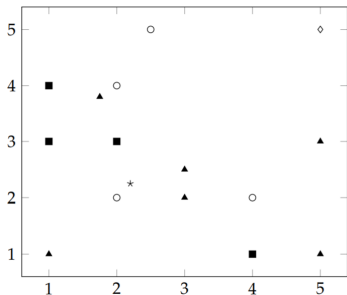


(a)

What is the prediction for the test point star when $K = 4$? Explain why

Problem 1

We denote the total number of training points as N and consider K -nearest-neighbor (KNN) classifier with L2 distance.

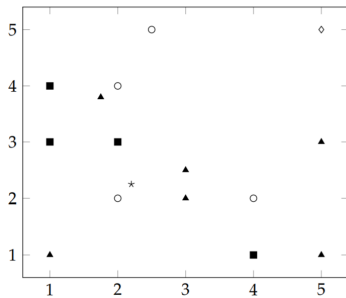


(a)

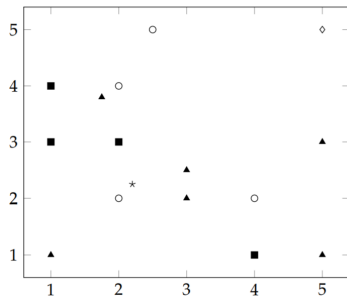
What is the prediction for the test point star when $K = 4$? Explain why

Triangle. The 4 nearest neighbors: 1 square, 1 circle, and 2 triangles.

Problem 1



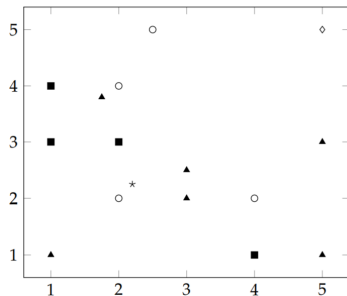
Problem 1



(b)

What is the diamond classified as for $K = N$? Explain why.

Problem 1

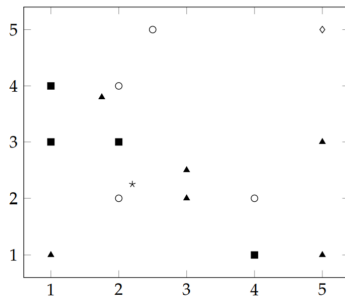


(b)

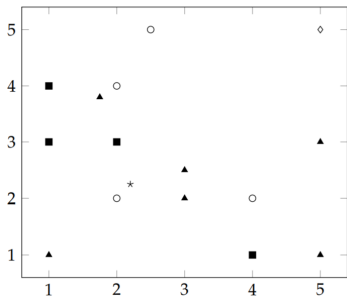
What is the diamond classified as for $K = N$? Explain why.

Triangle. When $K = N$ the prediction is always the majority label of the training set. Triangle is the majority in this example.

Problem 1



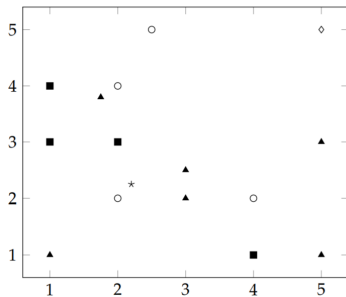
Problem 1



(c)

Suppose one performs leave-one-out validation (that is, N -fold cross validation) to choose the best K . List the triangles that are correctly classified (as a validation point) in this process for the run with $K = 1$.

Problem 1

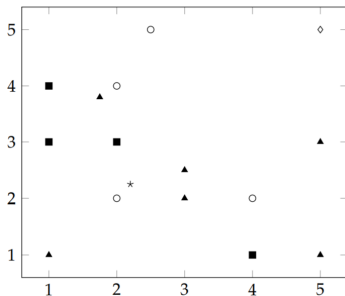


(c)

Suppose one performs leave-one-out validation (that is, N -fold cross validation) to choose the best K . List the triangles that are correctly classified (as a validation point) in this process for the run with $K = 1$.

Recall leave-one-out validation: use each point in turn as a validation dataset and use the others as a training dataset.

Problem 1



(c)

Suppose one performs leave-one-out validation (that is, N -fold cross validation) to choose the best K . List the triangles that are correctly classified (as a validation point) in this process for the run with $K = 1$.

Recall leave-one-out validation: use each point in turn as a validation dataset and use the others as a training dataset. The two triangles with x -value 3 are correctly classified.

Problem 2.1

In this problem we generalize the L2 regularization by

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

where $\mathbf{M} \in \mathbb{R}^{D \times D}$ is a positive definite matrix. Find the closed form of \mathbf{w}'_* .

Problem 2.1

In this problem we generalize the L2 regularization by

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

where $\mathbf{M} \in \mathbb{R}^{D \times D}$ is a positive definite matrix. Find the closed form of \mathbf{w}'_* .

Set $f(\mathbf{w}) = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \mathbf{w}^T \mathbf{M}\mathbf{w}$.

Problem 2.1

In this problem we generalize the L2 regularization by

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

where $\mathbf{M} \in \mathbb{R}^{D \times D}$ is a positive definite matrix. Find the closed form of \mathbf{w}'_* .

Set $f(\mathbf{w}) = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \mathbf{w}^T \mathbf{M}\mathbf{w}$. Then we have

$$\mathbf{w}'_* = \arg \min_{\mathbf{w}} f(\mathbf{w})$$

Problem 2.1

In this problem we generalize the L2 regularization by

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

where $\mathbf{M} \in \mathbb{R}^{D \times D}$ is a positive definite matrix. Find the closed form of \mathbf{w}'_* .

Set $f(\mathbf{w}) = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \mathbf{w}^T \mathbf{M}\mathbf{w}$. Then we have

$$\mathbf{w}'_* = \arg \min_{\mathbf{w}} f(\mathbf{w}) \Rightarrow \nabla f(\mathbf{w}'_*) = 0$$

Problem 2.1

In this problem we generalize the L2 regularization by

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

where $\mathbf{M} \in \mathbb{R}^{D \times D}$ is a positive definite matrix. Find the closed form of \mathbf{w}'_* .

Set $f(\mathbf{w}) = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \mathbf{w}^T \mathbf{M}\mathbf{w}$. Then we have

$$\mathbf{w}'_* = \arg \min_{\mathbf{w}} f(\mathbf{w}) \Rightarrow \nabla f(\mathbf{w}'_*) = 0$$

Since $\nabla f(\mathbf{w}) = 2\mathbf{X}^\top(\mathbf{X}\mathbf{w} - \mathbf{y}) + 2\mathbf{M}\mathbf{w}$, we get

Problem 2.1

In this problem we generalize the L2 regularization by

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

where $\mathbf{M} \in \mathbb{R}^{D \times D}$ is a positive definite matrix. Find the closed form of \mathbf{w}'_* .

Set $f(\mathbf{w}) = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \mathbf{w}^T \mathbf{M}\mathbf{w}$. Then we have

$$\mathbf{w}'_* = \arg \min_{\mathbf{w}} f(\mathbf{w}) \Rightarrow \nabla f(\mathbf{w}'_*) = 0$$

Since $\nabla f(\mathbf{w}) = 2\mathbf{X}^\top(\mathbf{X}\mathbf{w} - \mathbf{y}) + 2\mathbf{M}\mathbf{w}$, we get

$$\nabla f(\mathbf{w}'_*) = 2\mathbf{X}^\top(\mathbf{X}\mathbf{w}'_* - \mathbf{y}) + 2\mathbf{M}\mathbf{w}'_* = 0$$

Problem 2.1

In this problem we generalize the L2 regularization by

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

where $\mathbf{M} \in \mathbb{R}^{D \times D}$ is a positive definite matrix. Find the closed form of \mathbf{w}'_* .

Set $f(\mathbf{w}) = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \mathbf{w}^T \mathbf{M}\mathbf{w}$. Then we have

$$\mathbf{w}'_* = \arg \min_{\mathbf{w}} f(\mathbf{w}) \Rightarrow \nabla f(\mathbf{w}'_*) = 0$$

Since $\nabla f(\mathbf{w}) = 2\mathbf{X}^\top(\mathbf{X}\mathbf{w} - \mathbf{y}) + 2\mathbf{M}\mathbf{w}$, we get

$$\nabla f(\mathbf{w}'_*) = 2\mathbf{X}^\top(\mathbf{X}\mathbf{w}'_* - \mathbf{y}) + 2\mathbf{M}\mathbf{w}'_* = 0 \Rightarrow (\mathbf{X}^\top \mathbf{X} + \mathbf{M}) \mathbf{w}'_* = \mathbf{X}^\top \mathbf{y}$$

Problem 2.1

In this problem we generalize the L2 regularization by

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

where $\mathbf{M} \in \mathbb{R}^{D \times D}$ is a positive definite matrix. Find the closed form of \mathbf{w}'_* .

Set $f(\mathbf{w}) = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \mathbf{w}^T \mathbf{M}\mathbf{w}$. Then we have

$$\mathbf{w}'_* = \arg \min_{\mathbf{w}} f(\mathbf{w}) \Rightarrow \nabla f(\mathbf{w}'_*) = 0$$

Since $\nabla f(\mathbf{w}) = 2\mathbf{X}^\top(\mathbf{X}\mathbf{w} - \mathbf{y}) + 2\mathbf{M}\mathbf{w}$, we get

$$\nabla f(\mathbf{w}'_*) = 2\mathbf{X}^\top(\mathbf{X}\mathbf{w}'_* - \mathbf{y}) + 2\mathbf{M}\mathbf{w}'_* = 0 \Rightarrow (\mathbf{X}^\top \mathbf{X} + \mathbf{M}) \mathbf{w}'_* = \mathbf{X}^\top \mathbf{y}$$

As \mathbf{M} is positive definite, we conclude $\mathbf{w}'_* = (\mathbf{X}^\top \mathbf{X} + \mathbf{M})^{-1} \mathbf{X}^\top \mathbf{y}$

Matrix Calculus Example

Consider $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{A} \mathbf{x}$, where $\mathbf{x} \in \mathbb{R}^D$ and $\mathbf{A} \in \mathbb{R}^{D \times D}$, we have

Matrix Calculus Example

Consider $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{A} \mathbf{x}$, where $\mathbf{x} \in \mathbb{R}^D$ and $\mathbf{A} \in \mathbb{R}^{D \times D}$, we have

$$\nabla f(\mathbf{x}) = \left[\frac{\partial f}{\partial x_1} \quad \frac{\partial f}{\partial x_2} \quad \cdots \quad \frac{\partial f}{\partial x_D} \right]^\top$$

Matrix Calculus Example

Consider $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{A} \mathbf{x}$, where $\mathbf{x} \in \mathbb{R}^D$ and $\mathbf{A} \in \mathbb{R}^{D \times D}$, we have

$$\nabla f(\mathbf{x}) = \left[\frac{\partial f}{\partial x_1} \quad \frac{\partial f}{\partial x_2} \quad \cdots \quad \frac{\partial f}{\partial x_D} \right]^\top$$

Consider entry i , we have

Matrix Calculus Example

Consider $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{A} \mathbf{x}$, where $\mathbf{x} \in \mathbb{R}^D$ and $\mathbf{A} \in \mathbb{R}^{D \times D}$, we have

$$\nabla f(\mathbf{x}) = \left[\frac{\partial f}{\partial x_1} \quad \frac{\partial f}{\partial x_2} \quad \cdots \quad \frac{\partial f}{\partial x_D} \right]^\top$$

Consider entry i , we have

$$\frac{\partial f}{\partial x_i} = \frac{\partial}{\partial x_i} (\mathbf{x}^\top \mathbf{A} \mathbf{x})$$

Matrix Calculus Example

Consider $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{A} \mathbf{x}$, where $\mathbf{x} \in \mathbb{R}^D$ and $\mathbf{A} \in \mathbb{R}^{D \times D}$, we have

$$\nabla f(\mathbf{x}) = \left[\frac{\partial f}{\partial x_1} \quad \frac{\partial f}{\partial x_2} \quad \cdots \quad \frac{\partial f}{\partial x_D} \right]^\top$$

Consider entry i , we have

$$\begin{aligned} \frac{\partial f}{\partial x_i} &= \frac{\partial}{\partial x_i} \left(\mathbf{x}^\top \mathbf{A} \mathbf{x} \right) \\ &= \frac{\partial}{\partial x_i} \left(\sum_{j=1}^D x_j \cdot (\mathbf{A} \mathbf{x})_j \right) \end{aligned}$$

Matrix Calculus Example

Consider $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{A} \mathbf{x}$, where $\mathbf{x} \in \mathbb{R}^D$ and $\mathbf{A} \in \mathbb{R}^{D \times D}$, we have

$$\nabla f(\mathbf{x}) = \left[\frac{\partial f}{\partial x_1} \quad \frac{\partial f}{\partial x_2} \quad \cdots \quad \frac{\partial f}{\partial x_D} \right]^\top$$

Consider entry i , we have

$$\begin{aligned} \frac{\partial f}{\partial x_i} &= \frac{\partial}{\partial x_i} \left(\mathbf{x}^\top \mathbf{A} \mathbf{x} \right) \\ &= \frac{\partial}{\partial x_i} \left(\sum_{j=1}^D x_j \cdot (\mathbf{A} \mathbf{x})_j \right) \\ &= \frac{\partial}{\partial x_i} \left(\sum_{j=1}^D x_j \sum_{k=1}^D A_{jk} x_k \right) \end{aligned}$$

Matrix Calculus Example

Consider $f(\mathbf{x}) = \mathbf{x}^\top \mathbf{A} \mathbf{x}$, where $\mathbf{x} \in \mathbb{R}^D$ and $\mathbf{A} \in \mathbb{R}^{D \times D}$, we have

$$\nabla f(\mathbf{x}) = \left[\frac{\partial f}{\partial x_1} \quad \frac{\partial f}{\partial x_2} \quad \cdots \quad \frac{\partial f}{\partial x_D} \right]^\top$$

Consider entry i , we have

$$\begin{aligned} \frac{\partial f}{\partial x_i} &= \frac{\partial}{\partial x_i} \left(\mathbf{x}^\top \mathbf{A} \mathbf{x} \right) \\ &= \frac{\partial}{\partial x_i} \left(\sum_{j=1}^D x_j \cdot (\mathbf{A} \mathbf{x})_j \right) \\ &= \frac{\partial}{\partial x_i} \left(\sum_{j=1}^D x_j \sum_{k=1}^D A_{jk} x_k \right) \\ &= \frac{\partial}{\partial x_i} \left(x_i \sum_{k=1}^D A_{ik} x_k + \sum_{j \neq i} x_j \sum_{k=1}^D A_{jk} x_k \right) \end{aligned}$$

Matrix Calculus Example

$$\frac{\partial}{\partial x_i} (\mathbf{x}^\top \mathbf{A} \mathbf{x}) = \frac{\partial}{\partial x_i} \left(x_i \sum_{k=1}^D A_{ik} x_k + \sum_{j \neq i} x_j \sum_{k=1}^D A_{jk} x_k \right)$$

Matrix Calculus Example

$$\begin{aligned}\frac{\partial}{\partial x_i} (\mathbf{x}^\top \mathbf{A} \mathbf{x}) &= \frac{\partial}{\partial x_i} \left(x_i \sum_{k=1}^D A_{ik} x_k + \sum_{j \neq i} x_j \sum_{k=1}^D A_{jk} x_k \right) \\ &= x_i A_{ii} + \sum_{k=1}^D A_{ik} x_k + \frac{\partial}{\partial x_i} \left(\sum_{j \neq i} x_j \sum_{k=1}^D A_{jk} x_k \right)\end{aligned}$$

Matrix Calculus Example

$$\begin{aligned}\frac{\partial}{\partial x_i} (\mathbf{x}^\top \mathbf{A} \mathbf{x}) &= \frac{\partial}{\partial x_i} \left(x_i \sum_{k=1}^D A_{ik} x_k + \sum_{j \neq i} x_j \sum_{k=1}^D A_{jk} x_k \right) \\ &= x_i A_{ii} + \sum_{k=1}^D A_{ik} x_k + \frac{\partial}{\partial x_i} \left(\sum_{j \neq i} x_j \sum_{k=1}^D A_{jk} x_k \right) \\ &= x_i A_{ii} + \sum_{k=1}^D A_{ik} x_k + \sum_{j \neq i} x_j A_{ji}\end{aligned}$$

Matrix Calculus Example

$$\begin{aligned}\frac{\partial}{\partial x_i} (\mathbf{x}^\top \mathbf{A} \mathbf{x}) &= \frac{\partial}{\partial x_i} \left(x_i \sum_{k=1}^D A_{ik} x_k + \sum_{j \neq i} x_j \sum_{k=1}^D A_{jk} x_k \right) \\ &= x_i A_{ii} + \sum_{k=1}^D A_{ik} x_k + \frac{\partial}{\partial x_i} \left(\sum_{j \neq i} x_j \sum_{k=1}^D A_{jk} x_k \right) \\ &= x_i A_{ii} + \sum_{k=1}^D A_{ik} x_k + \sum_{j \neq i} x_j A_{ji} \\ &= \sum_{k=1}^D A_{ik} x_k + \sum_{j=1}^D x_j A_{ji}\end{aligned}$$

Matrix Calculus Example

$$\begin{aligned}\frac{\partial}{\partial x_i} (\mathbf{x}^\top \mathbf{A} \mathbf{x}) &= \frac{\partial}{\partial x_i} \left(x_i \sum_{k=1}^D A_{ik} x_k + \sum_{j \neq i} x_j \sum_{k=1}^D A_{jk} x_k \right) \\ &= x_i A_{ii} + \sum_{k=1}^D A_{ik} x_k + \frac{\partial}{\partial x_i} \left(\sum_{j \neq i} x_j \sum_{k=1}^D A_{jk} x_k \right) \\ &= x_i A_{ii} + \sum_{k=1}^D A_{ik} x_k + \sum_{j \neq i} x_j A_{ji} \\ &= \sum_{k=1}^D A_{ik} x_k + \sum_{j=1}^D x_j A_{ji} \\ &= (\mathbf{A} \mathbf{x})_i + (\mathbf{A}^\top \mathbf{x})_i\end{aligned}$$

Matrix Calculus Example

Using $\frac{\partial}{\partial x_i} (\mathbf{x}^\top \mathbf{A} \mathbf{x}) = (\mathbf{A} \mathbf{x})_i + (\mathbf{A}^\top \mathbf{x})_i$, we conclude

Matrix Calculus Example

Using $\frac{\partial}{\partial x_i} (\mathbf{x}^\top \mathbf{A} \mathbf{x}) = (\mathbf{A} \mathbf{x})_i + (\mathbf{A}^\top \mathbf{x})_i$, we conclude

$$\nabla f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_D} \end{bmatrix}$$

Matrix Calculus Example

Using $\frac{\partial}{\partial x_i} (\mathbf{x}^\top \mathbf{A} \mathbf{x}) = (\mathbf{A} \mathbf{x})_i + (\mathbf{A}^\top \mathbf{x})_i$, we conclude

$$\nabla f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_D} \end{bmatrix} = \begin{bmatrix} (\mathbf{A} \mathbf{x})_1 + (\mathbf{A}^\top \mathbf{x})_1 \\ (\mathbf{A} \mathbf{x})_2 + (\mathbf{A}^\top \mathbf{x})_2 \\ \vdots \\ (\mathbf{A} \mathbf{x})_D + (\mathbf{A}^\top \mathbf{x})_D \end{bmatrix}$$

Matrix Calculus Example

Using $\frac{\partial}{\partial x_i} (\mathbf{x}^\top \mathbf{A} \mathbf{x}) = (\mathbf{A} \mathbf{x})_i + (\mathbf{A}^\top \mathbf{x})_i$, we conclude

$$\nabla f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_D} \end{bmatrix} = \begin{bmatrix} (\mathbf{A} \mathbf{x})_1 + (\mathbf{A}^\top \mathbf{x})_1 \\ (\mathbf{A} \mathbf{x})_2 + (\mathbf{A}^\top \mathbf{x})_2 \\ \vdots \\ (\mathbf{A} \mathbf{x})_D + (\mathbf{A}^\top \mathbf{x})_D \end{bmatrix} = (\mathbf{A} + \mathbf{A}^\top) \mathbf{x}$$

Matrix Calculus Example

Using $\frac{\partial}{\partial x_i} (\mathbf{x}^\top \mathbf{A} \mathbf{x}) = (\mathbf{A} \mathbf{x})_i + (\mathbf{A}^\top \mathbf{x})_i$, we conclude

$$\nabla f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f}{\partial x_1} \\ \frac{\partial f}{\partial x_2} \\ \vdots \\ \frac{\partial f}{\partial x_D} \end{bmatrix} = \begin{bmatrix} (\mathbf{A} \mathbf{x})_1 + (\mathbf{A}^\top \mathbf{x})_1 \\ (\mathbf{A} \mathbf{x})_2 + (\mathbf{A}^\top \mathbf{x})_2 \\ \vdots \\ (\mathbf{A} \mathbf{x})_D + (\mathbf{A}^\top \mathbf{x})_D \end{bmatrix} = (\mathbf{A} + \mathbf{A}^\top) \mathbf{x}$$

When A is symmetric, we have $\nabla f(\mathbf{x}) = (\mathbf{A} + \mathbf{A}^\top) \mathbf{x} = 2\mathbf{A} \mathbf{x}$

Problem 2.2

Given a training set $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N) \in \mathbb{R}^D \times \mathbb{R}$, where each outcome $y_n = \mathbf{w}_*^T \mathbf{x}_n + \epsilon_n$ with ϵ_n being an independent Gaussian noise with zero mean and variance σ^2 for some $\sigma > 0$. In other words, the probability of seeing any outcome $y \in \mathbb{R}$ given \mathbf{x}_n is

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

Problem 2.2

Given a training set $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N) \in \mathbb{R}^D \times \mathbb{R}$, where each outcome $y_n = \mathbf{w}_*^T \mathbf{x}_n + \epsilon_n$ with ϵ_n being an independent Gaussian noise with zero mean and variance σ^2 for some $\sigma > 0$. In other words, the probability of seeing any outcome $y \in \mathbb{R}$ given \mathbf{x}_n is

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

Problem 2.2.1

Assume σ is fixed and given and $\mathbf{X}^T \mathbf{X}$ is invertible, find the maximum likelihood estimation for \mathbf{w}_* . In other words, first write down the probability of seeing the outcomes y_1, \dots, y_N given $\mathbf{x}_1, \dots, \mathbf{x}_N$ as a function of the value of \mathbf{w}_* ; then find the value of \mathbf{w}_* that maximizes this probability.

Review of Maximum Likelihood Estimation (MLE)

- The likelihood function $P(\mathbf{w})$: the probability of seeing the input $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$ as a function of \mathbf{w} :

$$P(\mathbf{w}) = \prod_{n=1}^N \mathbb{P}(y_n \mid \mathbf{x}_n; \mathbf{w})$$

Review of Maximum Likelihood Estimation (MLE)

- The likelihood function $P(\mathbf{w})$: the probability of seeing the input $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$ as a function of \mathbf{w} :

$$P(\mathbf{w}) = \prod_{n=1}^N \mathbb{P}(y_n \mid \mathbf{x}_n; \mathbf{w})$$

- Maximum Likelihood Estimation (MLE): find \mathbf{w}^* that maximizes the probability $P(\mathbf{w})$:

Review of Maximum Likelihood Estimation (MLE)

- The likelihood function $P(\mathbf{w})$: the probability of seeing the input $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)$ as a function of \mathbf{w} :

$$P(\mathbf{w}) = \prod_{n=1}^N \mathbb{P}(y_n | \mathbf{x}_n; \mathbf{w})$$

- Maximum Likelihood Estimation (MLE): find \mathbf{w}^* that maximizes the probability $P(\mathbf{w})$:

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmax}} P(\mathbf{w}) = \underset{\mathbf{w}}{\operatorname{argmax}} \prod_{n=1}^N \mathbb{P}(y_n | \mathbf{x}_n; \mathbf{w})$$

In this case, we have

In this case, we have

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

In this case, we have

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

Taking the negative log, this becomes

In this case, we have

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

Taking the negative log, this becomes

$$-\ln(P(\mathbf{w})) = \sum_{n=1}^N -\ln\left(\frac{1}{\sigma\sqrt{2\pi}} \exp\left(\frac{-(y_n - \mathbf{w}^T \mathbf{x}_n)^2}{2\sigma^2}\right)\right)$$

In this case, we have

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

Taking the negative log, this becomes

$$\begin{aligned} -\ln(P(\mathbf{w})) &= \sum_{n=1}^N -\ln\left(\frac{1}{\sigma\sqrt{2\pi}} \exp\left(\frac{-(y_n - \mathbf{w}^T \mathbf{x}_n)^2}{2\sigma^2}\right)\right) \\ &= \sum_{n=1}^N \ln(\sigma\sqrt{2\pi}) + \left(\frac{(y_n - \mathbf{w}^T \mathbf{x}_n)^2}{2\sigma^2}\right) \end{aligned}$$

In this case, we have

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

Taking the negative log, this becomes

$$\begin{aligned} -\ln(P(\mathbf{w})) &= \sum_{n=1}^N -\ln\left(\frac{1}{\sigma\sqrt{2\pi}} \exp\left(\frac{-(y_n - \mathbf{w}^T \mathbf{x}_n)^2}{2\sigma^2}\right)\right) \\ &= \sum_{n=1}^N \ln(\sigma\sqrt{2\pi}) + \left(\frac{(y_n - \mathbf{w}^T \mathbf{x}_n)^2}{2\sigma^2}\right) \\ &= N \ln \sqrt{2\pi} + N \ln \sigma + \frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \mathbf{w}^T \mathbf{x}_n)^2 \end{aligned}$$

In this case, we have

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

Taking the negative log, this becomes

$$\begin{aligned} -\ln(P(\mathbf{w})) &= \sum_{n=1}^N -\ln\left(\frac{1}{\sigma\sqrt{2\pi}} \exp\left(\frac{-(y_n - \mathbf{w}^T \mathbf{x}_n)^2}{2\sigma^2}\right)\right) \\ &= \sum_{n=1}^N \ln(\sigma\sqrt{2\pi}) + \left(\frac{(y_n - \mathbf{w}^T \mathbf{x}_n)^2}{2\sigma^2}\right) \\ &= N \ln \sqrt{2\pi} + N \ln \sigma + \frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \mathbf{w}^T \mathbf{x}_n)^2 \\ &= N \ln \sqrt{2\pi} + N \ln \sigma + \frac{1}{2\sigma^2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 \end{aligned}$$

Therefore, maximizing $P(\mathbf{w})$ is the same minimizing

$$-\ln(P(\mathbf{w})) = N \ln \sqrt{2\pi} + N \ln \sigma + \frac{1}{2\sigma^2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2$$

Therefore, maximizing $P(\mathbf{w})$ is the same minimizing

$$-\ln(P(\mathbf{w})) = N \ln \sqrt{2\pi} + N \ln \sigma + \frac{1}{2\sigma^2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2$$

As \mathbf{w} only appears in the last term, the problem is equivalent to minimizing $\|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2$, the same objective as for least square regression.

Therefore, maximizing $P(\mathbf{w})$ is the same minimizing

$$-\ln(P(\mathbf{w})) = N \ln \sqrt{2\pi} + N \ln \sigma + \frac{1}{2\sigma^2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2$$

As \mathbf{w} only appears in the last term, the problem is equivalent to minimizing $\|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2$, the same objective as for least square regression. Therefore, the MLE for is exactly the same as the least square solution:

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

Problem 2.2.2

Now consider σ as a parameter of the probabilistic model too, that is, the model is specified by both w_* and σ . Find the maximum likelihood estimation for w_* and σ .

Problem 2.2.2

Now consider σ as a parameter of the probabilistic model too, that is, the model is specified by both \mathbf{w}_* and σ . Find the maximum likelihood estimation for \mathbf{w}_* and σ .

Similar to 1., maximizing $P(\mathbf{w}, \sigma)$ is the same as minimizing

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

Problem 2.2.2

Now consider σ as a parameter of the probabilistic model too, that is, the model is specified by both \mathbf{w}_* and σ . Find the maximum likelihood estimation for \mathbf{w}_* and σ .

Similar to 1., maximizing $P(\mathbf{w}, \sigma)$ is the same minimizing

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

We first fix σ and minimize over \mathbf{w} , which leads to the same MLE for \mathbf{w}_* :

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

Problem 2.2.2

Now consider σ as a parameter of the probabilistic model too, that is, the model is specified by both \mathbf{w}_* and σ . Find the maximum likelihood estimation for \mathbf{w}_* and σ .

Similar to 1., maximizing $P(\mathbf{w}, \sigma)$ is the same minimizing

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

We first fix σ and minimize over \mathbf{w} , which leads to the same MLE for \mathbf{w}_* :

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

Next we minimize $-\ln(P(\mathbf{w}_*, \sigma))$ as function of σ by setting the derivative w.r.t. σ to be 0 :

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

Similar to 1., maximizing $P(\mathbf{w}, \sigma)$ is the same minimizing

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

We first fix σ and minimize over w , which leads to the same MLE for \mathbf{w}_* :

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

Next we minimize $-\ln(P(\mathbf{w}_*, \sigma))$ as function of σ by setting the derivative w.r.t. σ to be 0 :

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

Similar to 1., maximizing $P(\mathbf{w}, \sigma)$ is the same as minimizing

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

We first fix σ and minimize over w , which leads to the same MLE for \mathbf{w}_* :

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

Next we minimize $-\ln(P(\mathbf{w}_*, \sigma))$ as a function of σ by setting the derivative w.r.t. σ to be 0 :

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

Solving for σ and plugging $\mathbf{w}_* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ gives the MLE estimate:

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

Problem 3

For a binary classification dataset of N data points, where every \mathbf{x}_i has a corresponding label $y_i \in \{-1, 1\}$ and is normalized: $\|\mathbf{x}_i\| = \sqrt{\mathbf{x}_i^T \mathbf{x}_i} = 1, \forall i \in \{1, 2, \dots, N\}$, the perceptron algorithm proceeds as below: Suppose there exists an unknown hyperplane \mathbf{w}_{opt} such that $\|\mathbf{w}_{\text{opt}}\| = 1$ and $y_i \mathbf{w}_{\text{opt}}^T \mathbf{x}_i \geq \gamma, \forall i \in [N]$ for some positive value γ .

Problem 3

For a binary classification dataset of N data points, where every \mathbf{x}_i has a corresponding label $y_i \in \{-1, 1\}$ and is normalized: $\|\mathbf{x}_i\| = \sqrt{\mathbf{x}_i^T \mathbf{x}_i} = 1, \forall i \in \{1, 2, \dots, N\}$, the perceptron algorithm proceeds as below: Suppose there exists an unknown hyperplane \mathbf{w}_{opt} such that $\|\mathbf{w}_{\text{opt}}\| = 1$ and $y_i \mathbf{w}_{\text{opt}}^T \mathbf{x}_i \geq \gamma, \forall i \in [N]$ for some positive value γ .

Algorithm 1 Perceptron

Initialize $\mathbf{w}_1 = \mathbf{0}$

for $k = 1, 2, \dots$, **do**

 Pick a data point \mathbf{x}_i randomly

 Make a prediction $y = \text{sign}(\mathbf{w}_k^T \mathbf{x}_i)$ using current \mathbf{w}_k

if $y \neq y_i$ **then**

$\mathbf{w}_{k+1} \leftarrow \mathbf{w}_k + y_i \mathbf{x}_i$

else

$\mathbf{w}_{k+1} \leftarrow \mathbf{w}_k$

Problem 3

For a binary classification dataset of N data points, where every \mathbf{x}_i has a corresponding label $y_i \in \{-1, 1\}$ and is normalized: $\|\mathbf{x}_i\| = \sqrt{\mathbf{x}_i^T \mathbf{x}_i} = 1, \forall i \in \{1, 2, \dots, N\}$, the perceptron algorithm proceeds as below: Suppose there exists an unknown hyperplane \mathbf{w}_{opt} such that $\|\mathbf{w}_{opt}\| = 1$ and $y_i \mathbf{w}_{opt}^T \mathbf{x}_i \geq \gamma, \forall i \in [N]$ for some positive value γ .

Algorithm 1 Perceptron

Initialize $\mathbf{w}_1 = \mathbf{0}$

for $k = 1, 2, \dots$, **do**

 Pick a data point \mathbf{x}_i randomly

 Make a prediction $y = \text{sign}(\mathbf{w}_k^T \mathbf{x}_i)$ using current \mathbf{w}_k

if $y \neq y_i$ **then**

$\mathbf{w}_{k+1} \leftarrow \mathbf{w}_k + y_i \mathbf{x}_i$

else

$\mathbf{w}_{k+1} \leftarrow \mathbf{w}_k$

Problem 3.1

In iteration k such that $y_i \neq \text{sign}(\mathbf{w}_k^T \mathbf{x}_i)$, prove $\mathbf{w}_{k+1}^T \mathbf{w}_{opt} \geq \mathbf{w}_k^T \mathbf{w}_{opt} + \gamma$

Problem 3.1

In iteration k such that $y_i \neq \text{sign}(\mathbf{w}_k^T \mathbf{x}_i)$, prove $\mathbf{w}_{k+1}^T \mathbf{w}_{opt} \geq \mathbf{w}_k^T \mathbf{w}_{opt} + \gamma$

Problem 3.1

In iteration k such that $y_i \neq \text{sign}(\mathbf{w}_k^T \mathbf{x}_i)$, prove $\mathbf{w}_{k+1}^T \mathbf{w}_{opt} \geq \mathbf{w}_k^T \mathbf{w}_{opt} + \gamma$

When $y_i \neq \text{sign}(\mathbf{w}_k^T \mathbf{x}_i)$, we have

$$\mathbf{w}_{k+1} = \mathbf{w}_k + y_i \mathbf{x}_i$$

Problem 3.1

In iteration k such that $y_i \neq \text{sign}(\mathbf{w}_k^T \mathbf{x}_i)$, prove $\mathbf{w}_{k+1}^T \mathbf{w}_{opt} \geq \mathbf{w}_k^T \mathbf{w}_{opt} + \gamma$

When $y_i \neq \text{sign}(\mathbf{w}_k^T \mathbf{x}_i)$, we have

$$\mathbf{w}_{k+1} = \mathbf{w}_k + y_i \mathbf{x}_i$$

Taking inner product with \mathbf{w}_{opt} on the both sides, we get

$$\mathbf{w}_{k+1}^T \mathbf{w}_{opt} = \mathbf{w}_k^T \mathbf{w}_{opt} + y_i \mathbf{x}_i^T \mathbf{w}_{opt}$$

Problem 3.1

In iteration k such that $y_i \neq \text{sign}(\mathbf{w}_k^T \mathbf{x}_i)$, prove $\mathbf{w}_{k+1}^T \mathbf{w}_{opt} \geq \mathbf{w}_k^T \mathbf{w}_{opt} + \gamma$

When $y_i \neq \text{sign}(\mathbf{w}_k^T \mathbf{x}_i)$, we have

$$\mathbf{w}_{k+1} = \mathbf{w}_k + y_i \mathbf{x}_i$$

Taking inner product with \mathbf{w}_{opt} on the both sides, we get

$$\mathbf{w}_{k+1}^T \mathbf{w}_{opt} = \mathbf{w}_k^T \mathbf{w}_{opt} + y_i \mathbf{x}_i^T \mathbf{w}_{opt}$$

Since $y_i \mathbf{x}_i^T \mathbf{w}_{opt} \geq \gamma$ guaranteed by the problem, we conclude

$$\mathbf{w}_{k+1}^T \mathbf{w}_{opt} = \mathbf{w}_k^T \mathbf{w}_{opt} + y_i \mathbf{x}_i^T \mathbf{w}_{opt} \geq \mathbf{w}_k^T \mathbf{w}_{opt} + \gamma$$

Problem 3.2

In iteration k such that $y_i \neq \text{sign}(\mathbf{w}_k^T \mathbf{x}_i)$, prove $\|\mathbf{w}_{k+1}\|^2 \leq \|\mathbf{w}_k\|^2 + 1$

Problem 3.2

In iteration k such that $y_i \neq \text{sign}(\mathbf{w}_k^T \mathbf{x}_i)$, prove $\|\mathbf{w}_{k+1}\|^2 \leq \|\mathbf{w}_k\|^2 + 1$

When $y_i \neq \text{sign}(\mathbf{w}_k^T \mathbf{x}_i)$, we have

$$\mathbf{w}_{k+1} = \mathbf{w}_k + y_i \mathbf{x}_i$$

Problem 3.2

In iteration k such that $y_i \neq \text{sign}(\mathbf{w}_k^T \mathbf{x}_i)$, prove $\|\mathbf{w}_{k+1}\|^2 \leq \|\mathbf{w}_k\|^2 + 1$

When $y_i \neq \text{sign}(\mathbf{w}_k^T \mathbf{x}_i)$, we have

$$\mathbf{w}_{k+1} = \mathbf{w}_k + y_i \mathbf{x}_i$$

Therefore,

$$\|\mathbf{w}_{k+1}\|^2 = \mathbf{w}_{k+1}^T \mathbf{w}_{k+1} = (\mathbf{w}_k + y_i \mathbf{x}_i)^T (\mathbf{w}_k + y_i \mathbf{x}_i)$$

Problem 3.2

In iteration k such that $y_i \neq \text{sign}(\mathbf{w}_k^T \mathbf{x}_i)$, prove $\|\mathbf{w}_{k+1}\|^2 \leq \|\mathbf{w}_k\|^2 + 1$

When $y_i \neq \text{sign}(\mathbf{w}_k^T \mathbf{x}_i)$, we have

$$\mathbf{w}_{k+1} = \mathbf{w}_k + y_i \mathbf{x}_i$$

Therefore,

$$\begin{aligned}\|\mathbf{w}_{k+1}\|^2 &= \mathbf{w}_{k+1}^T \mathbf{w}_{k+1} = (\mathbf{w}_k + y_i \mathbf{x}_i)^T (\mathbf{w}_k + y_i \mathbf{x}_i) \\ &= \|\mathbf{w}_k\|^2 + 2y_i \mathbf{w}_k^T \mathbf{x}_i + y_i^2 \mathbf{x}_i^T \mathbf{x}_i\end{aligned}$$

Problem 3.2

In iteration k such that $y_i \neq \text{sign}(\mathbf{w}_k^T \mathbf{x}_i)$, prove $\|\mathbf{w}_{k+1}\|^2 \leq \|\mathbf{w}_k\|^2 + 1$

When $y_i \neq \text{sign}(\mathbf{w}_k^T \mathbf{x}_i)$, we have

$$\mathbf{w}_{k+1} = \mathbf{w}_k + y_i \mathbf{x}_i$$

Therefore,

$$\begin{aligned}\|\mathbf{w}_{k+1}\|^2 &= \mathbf{w}_{k+1}^T \mathbf{w}_{k+1} = (\mathbf{w}_k + y_i \mathbf{x}_i)^T (\mathbf{w}_k + y_i \mathbf{x}_i) \\ &= \|\mathbf{w}_k\|^2 + 2y_i \mathbf{w}_k^T \mathbf{x}_i + y_i^2 \mathbf{x}_i^T \mathbf{x}_i\end{aligned}$$

Since $y_i \neq \text{sign}(\mathbf{w}_k^T \mathbf{x}_i)$, we have $2y_i \mathbf{w}_k^T \mathbf{x}_i \leq 0$.

Problem 3.2

In iteration k such that $y_i \neq \text{sign}(\mathbf{w}_k^T \mathbf{x}_i)$, prove $\|\mathbf{w}_{k+1}\|^2 \leq \|\mathbf{w}_k\|^2 + 1$

When $y_i \neq \text{sign}(\mathbf{w}_k^T \mathbf{x}_i)$, we have

$$\mathbf{w}_{k+1} = \mathbf{w}_k + y_i \mathbf{x}_i$$

Therefore,

$$\begin{aligned}\|\mathbf{w}_{k+1}\|^2 &= \mathbf{w}_{k+1}^T \mathbf{w}_{k+1} = (\mathbf{w}_k + y_i \mathbf{x}_i)^T (\mathbf{w}_k + y_i \mathbf{x}_i) \\ &= \|\mathbf{w}_k\|^2 + 2y_i \mathbf{w}_k^T \mathbf{x}_i + y_i^2 \mathbf{x}_i^T \mathbf{x}_i\end{aligned}$$

Since $y_i \neq \text{sign}(\mathbf{w}_k^T \mathbf{x}_i)$, we have $2y_i \mathbf{w}_k^T \mathbf{x}_i \leq 0$. Moreover, by the assumption on input, we have $y_i^2 \mathbf{x}_i^T \mathbf{x}_i = y_i^2 \|\mathbf{x}_i\|^2 = 1$, so

Problem 3.2

In iteration k such that $y_i \neq \text{sign}(\mathbf{w}_k^T \mathbf{x}_i)$, prove $\|\mathbf{w}_{k+1}\|^2 \leq \|\mathbf{w}_k\|^2 + 1$

When $y_i \neq \text{sign}(\mathbf{w}_k^T \mathbf{x}_i)$, we have

$$\mathbf{w}_{k+1} = \mathbf{w}_k + y_i \mathbf{x}_i$$

Therefore,

$$\begin{aligned}\|\mathbf{w}_{k+1}\|^2 &= \mathbf{w}_{k+1}^T \mathbf{w}_{k+1} = (\mathbf{w}_k + y_i \mathbf{x}_i)^T (\mathbf{w}_k + y_i \mathbf{x}_i) \\ &= \|\mathbf{w}_k\|^2 + 2y_i \mathbf{w}_k^T \mathbf{x}_i + y_i^2 \mathbf{x}_i^T \mathbf{x}_i\end{aligned}$$

Since $y_i \neq \text{sign}(\mathbf{w}_k^T \mathbf{x}_i)$, we have $2y_i \mathbf{w}_k^T \mathbf{x}_i \leq 0$. Moreover, by the assumption on input, we have $y_i^2 \mathbf{x}_i^T \mathbf{x}_i = y_i^2 \|\mathbf{x}_i\|^2 = 1$, so

$$\|\mathbf{w}_{k+1}\|^2 = \|\mathbf{w}_k\|^2 + 2y_i \mathbf{w}_k^T \mathbf{x}_i + y_i^2 \mathbf{x}_i^T \mathbf{x}_i$$

Problem 3.2

In iteration k such that $y_i \neq \text{sign}(\mathbf{w}_k^T \mathbf{x}_i)$, prove $\|\mathbf{w}_{k+1}\|^2 \leq \|\mathbf{w}_k\|^2 + 1$

When $y_i \neq \text{sign}(\mathbf{w}_k^T \mathbf{x}_i)$, we have

$$\mathbf{w}_{k+1} = \mathbf{w}_k + y_i \mathbf{x}_i$$

Therefore,

$$\begin{aligned}\|\mathbf{w}_{k+1}\|^2 &= \mathbf{w}_{k+1}^T \mathbf{w}_{k+1} = (\mathbf{w}_k + y_i \mathbf{x}_i)^T (\mathbf{w}_k + y_i \mathbf{x}_i) \\ &= \|\mathbf{w}_k\|^2 + 2y_i \mathbf{w}_k^T \mathbf{x}_i + y_i^2 \mathbf{x}_i^T \mathbf{x}_i\end{aligned}$$

Since $y_i \neq \text{sign}(\mathbf{w}_k^T \mathbf{x}_i)$, we have $2y_i \mathbf{w}_k^T \mathbf{x}_i \leq 0$. Moreover, by the assumption on input, we have $y_i^2 \mathbf{x}_i^T \mathbf{x}_i = y_i^2 \|\mathbf{x}_i\|^2 = 1$, so

$$\begin{aligned}\|\mathbf{w}_{k+1}\|^2 &= \|\mathbf{w}_k\|^2 + 2y_i \mathbf{w}_k^T \mathbf{x}_i + y_i^2 \mathbf{x}_i^T \mathbf{x}_i \\ &\leq \|\mathbf{w}_k\|^2 + y_i^2 \mathbf{x}_i^T \mathbf{x}_i = \|\mathbf{w}_k\|^2 + 1\end{aligned}$$

Problem 3.3

For any iteration $k + 1$, with M being the total number of mistakes the algorithm has made so far for the first k iterations, we have

$$M\gamma \leq \|\mathbf{w}_{k+1}\| \leq \sqrt{M}$$

Problem 3.3

For any iteration $k + 1$, with M being the total number of mistakes the algorithm has made so far for the first k iterations, we have

$$M\gamma \leq \|\mathbf{w}_{k+1}\| \leq \sqrt{M}$$

By repeatedly applying results from Problem 3.1 we have

$$\mathbf{w}_{k+1}^T \mathbf{w}_{opt} \geq \mathbf{w}_k^T \mathbf{w}_{opt} + \gamma$$

Problem 3.3

For any iteration $k + 1$, with M being the total number of mistakes the algorithm has made so far for the first k iterations, we have

$$M\gamma \leq \|\mathbf{w}_{k+1}\| \leq \sqrt{M}$$

By repeatedly applying results from Problem 3.1 we have

$$\mathbf{w}_{k+1}^T \mathbf{w}_{opt} \geq \mathbf{w}_k^T \mathbf{w}_{opt} + \gamma \geq \cdots \geq \mathbf{w}_1^T \mathbf{w}_{opt} + M\gamma$$

Problem 3.3

For any iteration $k + 1$, with M being the total number of mistakes the algorithm has made so far for the first k iterations, we have

$$M\gamma \leq \|\mathbf{w}_{k+1}\| \leq \sqrt{M}$$

By repeatedly applying results from Problem 3.1 we have

$$\mathbf{w}_{k+1}^T \mathbf{w}_{opt} \geq \mathbf{w}_k^T \mathbf{w}_{opt} + \gamma \geq \cdots \geq \mathbf{w}_1^T \mathbf{w}_{opt} + M\gamma = M\gamma$$

Problem 3.3

For any iteration $k + 1$, with M being the total number of mistakes the algorithm has made so far for the first k iterations, we have

$$M\gamma \leq \|\mathbf{w}_{k+1}\| \leq \sqrt{M}$$

By repeatedly applying results from Problem 3.1 we have

$$\mathbf{w}_{k+1}^T \mathbf{w}_{opt} \geq \mathbf{w}_k^T \mathbf{w}_{opt} + \gamma \geq \dots \geq \mathbf{w}_1^T \mathbf{w}_{opt} + M\gamma = M\gamma$$

Using $\mathbf{a}^T \mathbf{b} \leq \|\mathbf{a}\| \|\mathbf{b}\|$ and $\|\mathbf{w}_{opt}\| = 1$, we have

$$M\gamma \leq \mathbf{w}_{k+1}^T \mathbf{w}_{opt} \leq \|\mathbf{w}_{k+1}\| \|\mathbf{w}_{opt}\| = \|\mathbf{w}_{k+1}\|$$

Problem 3.3

For any iteration $k + 1$, with M being the total number of mistakes the algorithm has made so far for the first k iterations, we have

$$M\gamma \leq \|\mathbf{w}_{k+1}\| \leq \sqrt{M}$$

By repeatedly applying results from Problem 3.1 we have

$$\mathbf{w}_{k+1}^T \mathbf{w}_{opt} \geq \mathbf{w}_k^T \mathbf{w}_{opt} + \gamma \geq \dots \geq \mathbf{w}_1^T \mathbf{w}_{opt} + M\gamma = M\gamma$$

Using $\mathbf{a}^T \mathbf{b} \leq \|\mathbf{a}\| \|\mathbf{b}\|$ and $\|\mathbf{w}_{opt}\| = 1$, we have

$$M\gamma \leq \mathbf{w}_{k+1}^T \mathbf{w}_{opt} \leq \|\mathbf{w}_{k+1}\| \|\mathbf{w}_{opt}\| = \|\mathbf{w}_{k+1}\|$$

Similarly, using results of Problem 3.2 repeatedly, we have

$$\|\mathbf{w}_{k+1}\|^2 \leq \|\mathbf{w}_k\|^2 + 1$$

Problem 3.3

For any iteration $k + 1$, with M being the total number of mistakes the algorithm has made so far for the first k iterations, we have

$$M\gamma \leq \|\mathbf{w}_{k+1}\| \leq \sqrt{M}$$

By repeatedly applying results from Problem 3.1 we have

$$\mathbf{w}_{k+1}^T \mathbf{w}_{opt} \geq \mathbf{w}_k^T \mathbf{w}_{opt} + \gamma \geq \dots \geq \mathbf{w}_1^T \mathbf{w}_{opt} + M\gamma = M\gamma$$

Using $\mathbf{a}^T \mathbf{b} \leq \|\mathbf{a}\| \|\mathbf{b}\|$ and $\|\mathbf{w}_{opt}\| = 1$, we have

$$M\gamma \leq \mathbf{w}_{k+1}^T \mathbf{w}_{opt} \leq \|\mathbf{w}_{k+1}\| \|\mathbf{w}_{opt}\| = \|\mathbf{w}_{k+1}\|$$

Similarly, using results of Problem 3.2 repeatedly, we have

$$\|\mathbf{w}_{k+1}\|^2 \leq \|\mathbf{w}_k\|^2 + 1 \leq \dots \leq \|\mathbf{w}_1\|^2 + M$$

Problem 3.3

For any iteration $k + 1$, with M being the total number of mistakes the algorithm has made so far for the first k iterations, we have

$$M\gamma \leq \|\mathbf{w}_{k+1}\| \leq \sqrt{M}$$

By repeatedly applying results from Problem 3.1 we have

$$\mathbf{w}_{k+1}^T \mathbf{w}_{opt} \geq \mathbf{w}_k^T \mathbf{w}_{opt} + \gamma \geq \dots \geq \mathbf{w}_1^T \mathbf{w}_{opt} + M\gamma = M\gamma$$

Using $\mathbf{a}^T \mathbf{b} \leq \|\mathbf{a}\| \|\mathbf{b}\|$ and $\|\mathbf{w}_{opt}\| = 1$, we have

$$M\gamma \leq \mathbf{w}_{k+1}^T \mathbf{w}_{opt} \leq \|\mathbf{w}_{k+1}\| \|\mathbf{w}_{opt}\| = \|\mathbf{w}_{k+1}\|$$

Similarly, using results of Problem 3.2 repeatedly, we have

$$\|\mathbf{w}_{k+1}\|^2 \leq \|\mathbf{w}_k\|^2 + 1 \leq \dots \leq \|\mathbf{w}_1\|^2 + M = M$$

Problem 3.4

Using result of Problem 3.3, conclude $M \leq \gamma^{-2}$

Problem 3.4

Using result of Problem 3.3, conclude $M \leq \gamma^{-2}$

From Problem 3.3, we have

$$M\gamma \leq \sqrt{M}$$

Problem 3.4

Using result of Problem 3.3, conclude $M \leq \gamma^{-2}$

From Problem 3.3, we have

$$M\gamma \leq \sqrt{M}$$

This implies

$$\sqrt{M} \leq \frac{1}{\gamma}$$

Problem 3.4

Using result of Problem 3.3, conclude $M \leq \gamma^{-2}$

From Problem 3.3, we have

$$M\gamma \leq \sqrt{M}$$

This implies

$$\sqrt{M} \leq \frac{1}{\gamma} \Rightarrow M \leq \frac{1}{\gamma^2}$$