

# 1. Linear Algebra

$$\Rightarrow \det A = \sum_{m \times m} \sum_{n \times n} V^* \quad V^* = \text{conjugate transpose}$$

Define  $K_{m \times m} \rightarrow$  Anti-diagonal matrix with all entries 1

Note:  $B = ATK$  will give A rotated by  $90^\circ$  as shown.

$$A = \begin{bmatrix} a_{11} & & a_{1n} \\ a_{21} & \cdots & a_{2n} \\ \vdots & & \vdots \\ a_{m1} & & a_{mn} \end{bmatrix} \quad \begin{bmatrix} a_{11} & a_{21} & \cdots & a_{m1} \\ | & | & \cdots & | \\ a_{1n} & a_{2n} & \cdots & a_{nn} \end{bmatrix} \quad \begin{bmatrix} 0 & & & \\ & \ddots & & \\ & & 0 & \\ & \ddots & & \ddots \end{bmatrix} = \begin{bmatrix} a_{m1} & & a_{21} & a_{1n} \\ 1 & & 1 & \vdots \\ 0_{mm} & & a_{2n} & a_{nn} \end{bmatrix} \quad B$$

$A^T = m \times m \quad K = m \times m$

$$\therefore \text{eigenvalue}(B) = \text{eigenvalues}((U\Sigma V^*)^T K)$$

$$= \text{eigenvalues}(V^T \Sigma U^T K)$$

• Now  $\Sigma^T = \Sigma$  since it's rectangular diagonal matrix (with some abuse of notation.)

$$\bullet \text{ Let } U^T K = U'$$

$$= \text{eigenvalue}(V^T \Sigma U')$$

= diagonal elements of  $\Sigma$

= eigenvalues of A.

Hence, both A & B have same eigen values.

## 1.2 Norms

$x$  (vector),  $A_{m \times n}$  Given

$$(a) \|x\|_\infty \leq \|x\|_2$$

Let  $x = (a_1, a_2, \dots, a_m)$

$$\|x\|_\infty = \max_k |a_k| \quad k \in \{1, 2, \dots, m\} = |a_i| \text{ (let)}$$

Now

$$a_i^2 \leq a_i^2$$

$$\leq a_1^2 + a_2^2 + \dots + a_i^2 + \dots + a_m^2$$

Take square root both sides

$$|a_i| \leq \sqrt{a_1^2 + a_2^2 + \dots + a_m^2}$$

$$\|x\|_\infty \leq \|x\|_2$$

Hence, Proved

Note equality  
holds iff  $a_k = 0$   
 $\forall k \neq i$  i.e  
all vectors  
 $x$  with only  
1 non-zero  
element  
Eg  $x = e_i$

$$(b) \|x\|_2 \leq \sqrt{m} \|x\|_\infty$$

Again let  $\|x\|_\infty = a_i$  for some  $i \in \{1, \dots, m\}$

where  $x = (a_1, a_2, \dots, a_m)$

Now since  $a_i$  has highest magnitude

$$a_1^2 \leq a_i^2 \quad \text{--- (1)}$$

$$a_2^2 \leq a_i^2 \quad \text{--- (2)}$$

|

$$a_m^2 \leq a_i^2 \quad \text{--- (m)}$$

Adding (1), (2) --- (m)

$$a_1^2 + a_2^2 + \dots + a_m^2 \leq m a_i^2$$

$$\Rightarrow \sqrt{a_1^2 + a_2^2 + \dots + a_m^2} \leq \sqrt{m} |a_i|$$

$$\Rightarrow \|x\|_2 \leq \sqrt{m} \|x\|_\infty$$

$$\text{Equality } \|x\|_2 = \sqrt{n} \|x\|_\infty$$

holds iff  $a_1 = a_2 = \dots = a_n$  i.e.  
all elements of  $x$  are equal.

$$\text{Otherwise } \|x\|_2 < \sqrt{n} \|x\|_\infty$$

$$\text{Eg } x = (2, 2, \dots, 2) \quad \lambda \in \mathbb{R}$$

$$(c) \|A\|_\infty \text{ is defined as } A_{m \times n} \quad x_{n \times 1} \quad (Ax)_{m \times 1}$$

$$\max_{x \neq 0} \frac{\|Ax\|_\infty}{\|x\|_\infty} \quad \textcircled{1} - \text{ Using part a) } \quad \|Ax\|_\infty \leq \|Ax\|_2$$

$$\textcircled{2} - \text{ Using part b) } \quad \|x\|_\infty \geq \frac{1}{\sqrt{n}} \|x\|_2$$

Combining \textcircled{1} & \textcircled{2}

$$\Rightarrow \frac{1}{\|x\|_\infty} \leq \frac{\sqrt{n}}{\|x\|_2}$$

$$\|A\|_\infty = \max_{x \neq 0} \frac{\|Ax\|_\infty}{\|x\|_\infty} \leq \max_{x \neq 0} \sqrt{n} \frac{\|Ax\|_2}{\|x\|_2} \leq \sqrt{n} \|A\|_2$$

Hence proved.

\* Equality is achieved if  $\|Ax\|_\infty = \|Ax\|_2$  i.e.  $Ax$  has exactly one non-zero element. Also  $\|x\|_\infty = \frac{\|x\|_2}{\sqrt{2}}$   $\Rightarrow$  All elements of  $x$  are exactly same. Hence  $A$  has exactly one non-zero row than above conditions are satisfied. (This is one such  $A$ )

$$\text{Eg } A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{bmatrix} \begin{bmatrix} \lambda \\ \lambda \\ \vdots \\ \lambda \end{bmatrix} = \begin{bmatrix} \lambda(a_{11} + a_{12} + \dots + a_{1n}) \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

But we can't put restriction on  $x$ . In general if we take with 1 row containing all elements we can show that  $\|A\|_\infty = \|A\|_2$  if rest equal of the rows are zero.

$$\text{Let } A = \begin{bmatrix} \lambda & \lambda & \lambda & \cdots & \lambda \\ 0 & 0 & 0 & \cdots & 0 \\ 0 & 0 & 0 & \cdots & 0 \end{bmatrix}_{m \times n} \quad (\text{Let } \lambda > 0)$$

$$\|A\|_\infty = \max \text{ of all sum} = n\lambda$$

$$\|A\|_2 = \sqrt{\text{highest eigen value } (A^T A)}$$

$$A^T A = \begin{bmatrix} \lambda^2 & \lambda^2 & \cdots & \lambda^2 \\ \lambda^2 & \lambda^2 & \cdots & \lambda^2 \\ \lambda^2 & \lambda^2 & \cdots & \lambda^2 \end{bmatrix}_{n \times n} \quad \text{All value} = \lambda^2$$

Note since rank = 1 all eigenvalues are zero except 1 i.e corresponding to eigen vector  $(1, 1, -1)^T$

$$\begin{bmatrix} \lambda^2 & \lambda^2 & \cdots & \lambda^2 \\ \lambda^2 & \lambda^2 & \cdots & \lambda^2 \\ \lambda^2 & \lambda^2 & \cdots & \lambda^2 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \frac{1}{n} \lambda^2 \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} = \boxed{n\lambda^2}$$

$$\|A\|_2 = \sqrt{n\lambda^2} = \sqrt{n}\lambda$$

$$\left( \begin{array}{l} \|A\|_\infty = \sqrt{n} \|A\|_2 \quad \text{iff} \\ n\lambda = \sqrt{n} (\sqrt{n}\lambda) \\ = n\lambda \end{array} \right)$$

Hence, our matrix A is a general matrix of size  $m \times n$  satisfying

$$\boxed{\|A\|_\infty = \sqrt{n} \|A\|_2}$$

$$(d) \|A\|_2 \leq \sqrt{m} \|A\|_\infty$$

$$\|A\|_2 = \max_{n \neq 0} \frac{\|Ax_n\|_2}{\|x_n\|_2}$$

Using ① & ②

$$\text{Using b)} \|Ax_n\|_2 \leq \sqrt{m} \|Ax_n\|_\infty \quad \text{--- ①}$$

$$\text{Using a)} \|x_n\|_2 \geq \|x_n\|_\infty$$

$$\Rightarrow \frac{1}{\|x_n\|_2} \leq \frac{1}{\|x_n\|_\infty} \quad \text{--- ②}$$

$$\|A\|_2 = \max_{n \neq 0} \frac{\|Ax_n\|_2}{\|x_n\|_2} \leq \max_{n \neq 0} \frac{\sqrt{m} \|Ax_n\|_\infty}{\|x_n\|_\infty} \leq \sqrt{m} \|A\|_\infty$$

Hence, Proved

Equality is attained when

- $\sqrt{m} \|A\|_\infty = \|Ax\|_2$  i.e.  $Ax$  has all elements equal.
- $\|x\|_\infty > \|x\|_2$  i.e.  $x$  has only one non zero element

We can choose  $A$  s.t it has all elements in column  $i$  having some values  $a_i$  and all other columns  $0$ , where  $i$  is given by  $\|x\|_\infty = a_i$  where  $x = (0, 0, 0, \dots, a_i, 0)$  i.e. non zero element of  $x$ 's index.

$$\begin{bmatrix} a_i & 0 & 0 & \dots & 0 \\ 0 & a_i & 0 & \dots & 0 \\ 0 & 0 & a_i & \dots & 0 \\ 0 & 0 & 0 & \dots & a_i \\ 0 & 0 & 0 & \dots & 0 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} = \begin{bmatrix} a_i \\ a_i \\ a_i \\ \vdots \\ a_i \end{bmatrix}$$

$A \quad x \quad Ax$

But we can't put a condition on  $x$ . In general we can show that if we take  $A = \begin{bmatrix} 1 & 0 & \dots & 0 \\ 1 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 0 & \dots & 0 \end{bmatrix}$  i.e. only first column non zero with equal values, then  $\|Ax\|_2 = \sqrt{m} \|A\|_\infty$

1st column non zero with equal values

Let  $A = \begin{bmatrix} \lambda & 0 & 0 & \dots & 0 \\ \lambda & 0 & 0 & \dots & 0 \\ \lambda & 0 & 0 & \dots & 0 \\ \vdots & | & | & \dots & | \\ 0 & 0 & 0 & \dots & 0 \end{bmatrix}$   $\lambda > 0$  (det)

$$\|A\|_{\infty} = \lambda \quad (\text{max row sum})$$

$$\|A\|_2 = \sqrt{\text{max eigenvalue}(A^T A)}$$

$$A^T A = \begin{bmatrix} m\lambda^2 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ \vdots & | & | & | & | & | \\ 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

All elements zero except (I, I).

Clearly all eigenvectors are zero except  $(1, 0, 0 \dots)^T$   
and eigen value =  $m\lambda^2$

$$\therefore \|A\|_2 = \sqrt{m\lambda^2} = \sqrt{m}\lambda = \sqrt{m} \|A\|_{\infty}$$

Hence A is an eg. of general  $m \times n$  matrix s.t.

$$\boxed{\|A\|_2 = \sqrt{m} \|A\|_{\infty}}$$

Q2

### Subgradients

(a)  $f(x) = \max_{i=1, 2, \dots, m} (a_i^T x + b_i)$

Find  $k \in 1, 2, \dots, m$  s.t

$$f(x) = a_k^T x + b_k$$

then subgradient at that particular  $x = \frac{\partial f(x)}{\partial x} = a_k$

(b)  $f(x) = \max_{i=1, 2, \dots, m} |a_i^T x + b_i|$

Find  $k \in 1, 2, \dots, m$  s.t

$$f(x) = |a_k^T x + b_k|$$

if  $a_k^T x + b_k \geq 0$   $f(x) = a_k^T x + b_k$  Subgradient  $= \frac{\partial f(x)}{\partial x} = a_k$

if  $a_k^T x + b_k < 0$   $f(x) = -a_k^T x - b_k$  Subgradient  $= \frac{-\partial f(x)}{\partial x} = -a_k$

(c)  $f(x) = \sup_{0 \leq t \leq 1} p(t)$  where

$$p(t) = x_1 + x_2 t + \dots + x_n t^{n-1}$$

Find a  $t'$  such that  $f(x) = p(t')$   $0 \leq t' \leq 1$

Since  $p(t)$  is a polynomial in  $t$  we can plot  $p(t)$  (if  $x$  is given & fixed) in range  $[0, 1]$  and detect  $t'$  at which  $p(t)$  is maximum. Or we can find  $t'$  by finding roots of  $p'(t)$  in  $t \in [0, 1]$ .

Then Subgradient is given by

$$\nabla_x f(x) = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \vdots \\ \frac{\partial f(x)}{\partial x_n} \end{bmatrix} = \begin{bmatrix} 1 \\ t' \\ t'^2 \\ \vdots \\ t'^{n-1} \end{bmatrix}$$

(Assuming  $p(t)$  has one, if not maxima is obtained at either  $t=0$  or  $t=1$ )

$$(d) f(x) = x_{[1]} + \dots + x_{[k]} \quad x_{[i]}: i^{\text{th}} \text{ largest element of vector } x.$$

Let  $x \in \mathbb{R}^n$

Let at any given  $x \in \mathbb{R}^n$   $x_{i_1}, x_{i_2}, \dots, x_{i_k}$  are top  $k$  components of vector  $x$ . Define set  $S = \{i_1, \dots, i_k\}$ . Then a gradient at point  $x$  is given by  $g \in \mathbb{R}^n$

St

$$g_i = \begin{cases} 1 & i \in S \quad \forall i \in \{1, 2, \dots, n\} \\ 0 & i \notin S \end{cases}$$

$$(e) f(x) = \inf_{Ay \leq b} \|x - y\|_2^2$$

Let for a given  $x$ ,  $y^*$  be the point in feasible region i.e.  $Ay^* \leq b$  at which  $f(x)$  attains its infimum. Then  $f(x) = (x - y^*)^T(x - y^*)$

$$\text{Subgradient of } f(x) = \nabla_x f(x) = 2x - 2y^* = 2(x - y^*)$$

To find  $y^*$  such we can solve the Dual Problem. Since Slater's condition is satisfied we have "Strict Duality".

~~Primal Problem~~ Primal Problem  $\min \|x - y\|_2^2$  (Note:  $x$  is constant)  
Subject to  $Ay - b \leq 0$

$$L(z, y) = (x - y)^T(x - y) + z^T(Ay - b) = x^T x - x^T y - y^T x + y^T y + z^T A y - z^T b$$

~~Dual~~  $g(z) = \min_y L(z, y)$  To find  $g(z)$  let's find.  $\nabla_y L(z, y) = -2x + 2y + A^T z = 0$   
 $\Rightarrow y = x - \frac{1}{2} A^T z \leftarrow \text{substitute it back into } L \text{ to get } g(z)$

$$g(z) = \left(\frac{1}{2} A^T z\right)^T x - \frac{1}{2} A^T z + z^T (-Ax - \frac{1}{2} A A^T z - b) = -\frac{1}{4} z^T A A^T z + z^T A x - z^T b$$

Dual Problem is  $\max_z g(z)$  subject to  $z \geq 0$ . Let  $z^*$  is soln to dual problem, then optimal value of primal problem i.e.  $f(x) = \left(-\frac{1}{2} z^{*T} A A^T z^* - z^{*T} b\right)$

By weak duality we have  $f(x) \geq -\frac{1}{2} z^{*T} A A^T z^* - z^{*T} A x - z^{*T} b - 2^{*T} b$  (1)

$\therefore f(x') - f(x) \geq z^{*T} A (x - x')$  Hence  $(z^{*T} A)$  is Subgradient at  $x$ .

Now if  $y^*$  &  $z^*$  are optimal soln then by KKT

$$\nabla_y L(y^*, z^*) = 0$$

$$\Rightarrow \boxed{y^* = x - \frac{1}{2} A^T z^*} \rightarrow \text{Hence we found one such } y.$$

Note: Subgradient =  $\boxed{A^T z^* = g(x - y^*)}$

### Q3 Probability

$$\begin{aligned} 1. \quad E(I/x=x) &= P(I=1/x=x) \times 1 + P(I=0/x=x) \times 0 \\ &= P(I=1/x=x) \\ &= P(Z < x/x=x) \\ &= P(Z < x) \quad \text{since } Z \sim N(0, 1) \\ &= \phi(x) \end{aligned}$$

$$\begin{aligned} 2. \quad \textcircled{*} \quad E(I) &= P(I=1) \times 1 + P(I=0) \times 0 = P(I=1) = P(Z < x) - ① \\ \textcircled{+} \quad \text{Using } E[x] &= E[E[x|y]] \\ E[I] &= E(E(I/x)) \quad \text{Using Part 1, } E(I/x) = \phi(x) \\ &= E(\phi(x)) - ② \end{aligned}$$

Using ① & ②

$$E[\phi(x)] = P(Z < x)$$

$$3. \quad E(\phi(x)) = P(Z < x) = P(x - Z > 0)$$

$$\begin{aligned} \text{Now } X &\sim N(\mu, 1) \\ Z &\sim N(0, 1) \\ -Z &\sim N(0, 1) \end{aligned} \Rightarrow \begin{aligned} X - Z &\sim N(\mu, 2) \\ \frac{(x-z)-\mu}{\sqrt{2}} &\sim N(0, 1) \end{aligned}$$

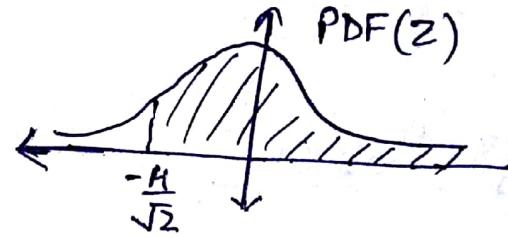
$$\therefore E(\phi(x)) = P(x - \mu > 0)$$

$$= P\left(\frac{(x-\mu)}{\sqrt{2}} > \frac{-\mu}{\sqrt{2}}\right)$$

It follows standard normal distribution

$$= P\left(Z > \frac{-\mu}{\sqrt{2}}\right)$$

Using Symmetry



$$= P\left(Z < \frac{\mu}{\sqrt{2}}\right)$$

$$= \Phi\left(\frac{\mu}{\sqrt{2}}\right)$$

## Q4 Machine Learning

### 4.1 PCA

$$\underset{u: \|u\|_2=1}{\text{max}} \tilde{V}[u^T x]$$

$$\text{Given } \frac{1}{N} \sum_{i=1}^N x_i = 0 \quad i=1, 2, \dots, N$$

④ Projected points are  $(u^T x_1, u^T x_2, \dots, u^T x_N)$

⑤ Avg of projected points  $\frac{1}{N}(u^T \sum_{i=1}^N x_i) = 0$

$$\therefore \tilde{V}[u^T x] = \frac{\sum_{i=1}^N (u^T x_i - 0)^2}{N}$$

$$\max_{u: \|u\|_2=1} \frac{\sum_{i=1}^N (u^T x_i)(u^T x_i)}{N}$$

$$\max_{u: \|u\|_2=1} u^T \left( \frac{\sum_{i=1}^N x_i x_i^T}{N} \right) u$$

$$\therefore \sum = \frac{1}{N} \sum_{i=1}^N x_i x_i^T$$

# 4.1 PCA

2.

$$\min_{\mathbf{u}: \|\mathbf{u}\|_2=1} \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - \mathbf{u} \mathbf{u}^T \mathbf{x}_i\|_2^2$$

$$= \min_{\mathbf{u}: \|\mathbf{u}\|_2=1} \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \mathbf{u} \mathbf{u}^T \mathbf{x}_i)^T (\mathbf{x}_i - \mathbf{u} \mathbf{u}^T \mathbf{x}_i)$$

$$= \min_{\mathbf{u}: \|\mathbf{u}\|_2=1} \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i^T \mathbf{x}_i - \mathbf{x}_i^T \mathbf{u} \mathbf{u}^T \mathbf{x}_i - \mathbf{x}_i^T \mathbf{u} \mathbf{u}^T \mathbf{x}_i + \mathbf{x}_i^T \mathbf{u} \mathbf{u}^T \mathbf{u}^T \mathbf{x}_i$$

Using Fact that  $\mathbf{u}^T \mathbf{u} = 1$

and  $\mathbf{x}_i^T \mathbf{u} = \mathbf{u}^T \mathbf{x}_i$  is a scalar

$$= \min_{\mathbf{u}: \|\mathbf{u}\|_2=1} \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i^T \mathbf{x}_i - (\mathbf{u}^T \mathbf{x}_i)^2$$

Ignoring  $\sum_{i=1}^N \mathbf{x}_i^T \mathbf{x}_i$  since it doesn't depend on  $\mathbf{u}$

$$= \min_{\mathbf{u}: \|\mathbf{u}\|_2=1} -\frac{1}{N} \sum_{i=1}^N (\mathbf{u}^T \mathbf{x}_i)^2$$

$$= \max_{\mathbf{u}: \|\mathbf{u}\|_2=1} \frac{1}{N} \sum_{i=1}^N (\mathbf{u}^T \mathbf{x}_i)^2 \quad \leftarrow \text{Same as Part 1 i.e variance term}$$

$$= \max_{\mathbf{u}: \|\mathbf{u}_2\|=1} \mathbf{u}^T \Sigma \mathbf{u} \quad \text{where } \boxed{\Sigma = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T}$$

## 4.2 Bias Variance Trade Off

$$E_D(J(\theta)) = ?$$

$$\text{Bias} = \bar{h}_\theta(x) - y \quad \text{where } \bar{h}_\theta(x) = E_D[h_\theta(x)]$$

$$\begin{aligned}\text{Variance} &= E_D((h_\theta(x) - \bar{h}_\theta(x))^2) \\ &= E_D(h_\theta^2(x) - 2h_\theta(x)\bar{h}_\theta(x) + \bar{h}_\theta^2(x)) \\ &= E_D(h_\theta^2(x)) - (\bar{h}_\theta(x))^2\end{aligned}$$

a)  $E_D(J(\theta)) = E_D((y - h_\theta(x))^2)$

$$\begin{aligned}&= E_D(y^2 + h_\theta^2(x) - 2yh_\theta(x)) \\ &= y^2 + E_D(h_\theta^2(x)) - 2yh_\theta(x) \\ &\quad \text{Add & Subtract } (\bar{h}_\theta(x))^2 \\ &= \underbrace{[y^2 - 2yh_\theta(x) + (\bar{h}_\theta(x))^2]}_{\text{Bias}^2} + \underbrace{E_D(h_\theta^2(x)) - (\bar{h}_\theta(x))^2}_{\text{Variance}} \\ &= \text{Bias}^2 + \text{Variance}\end{aligned}$$

b)  $\text{Bias} = \bar{h}_\theta(x) - f(x) \quad \text{Variance} = E_D(h_\theta^2(x)) - (\bar{h}_\theta(x))^2$

$$\begin{aligned}E(J(\theta)) &= E_D((f(x) + \varepsilon - h_\theta(x))^2) = E_D((f(x) + \varepsilon)^2 + h_\theta^2(x) - 2h_\theta(x)f(x) + \varepsilon^2) \\ &= E_D(f(x)^2 + \varepsilon^2 + 2f(x)\varepsilon) + E_D(h_\theta^2(x)) - 2E(h_\theta(x)f(x)) - 2E(\varepsilon h_\theta(x))\end{aligned}$$

Note  $f(x)$  is deterministic & not random  $\varepsilon \sim \mathcal{N}(\theta, \sigma^2)$

$$\therefore E(\varepsilon) = 0 \quad \text{Variance} = \sigma^2 = E((\varepsilon - 0)^2) = E(\varepsilon^2)$$

$$\begin{aligned}&= f(x)^2 + E(\varepsilon^2) + 2f(x)E(\varepsilon) + E_D(h_\theta^2(x)) - 2f(x)\bar{h}_\theta(x) - 2E(\varepsilon)\bar{h}_\theta(x) \\ &\quad \text{Add & subtract } (\bar{h}_\theta(x))^2 \\ &= [f(x)^2 - 2f(x)\bar{h}_\theta(x) + (\bar{h}_\theta(x))^2] + (E_D(h_\theta^2(x)) - (\bar{h}_\theta(x))^2) + E(\varepsilon^2) \\ &= \text{Bias}^2 + \text{Variance} + \sigma^2\end{aligned}$$

### 4.3 Kernelizing the Perceptron

$$\Theta^{t+1} = \Theta^t + \alpha [y^{t+1} - h_{\Theta^t}(x^{t+1})] x^{t+1}$$

Note: we have found a way to update our parameters as follows

$$\Theta^{t+1} = \Theta^t + \alpha [y^{t+1} - h_{\Theta^t}(\phi(x^{t+1}))] \phi(x^{t+1})$$

Note:  $y^{t+1} \in \{0, 1\}$

$$h_{\Theta^t}(\phi(x^{t+1})) = g(\Theta^t \cdot \phi(x^{t+1})), \in \{0, 1\}$$

$$\therefore \underbrace{y^{t+1} - h_{\Theta^t}(\phi(x^{t+1}))}_{\text{error}} \in \{-1, 0, 1\}$$

Let's call it error

$$\therefore \Theta^{t+1} = \begin{cases} \Theta^t - \alpha \phi(x^{t+1}) \text{ error} = -1 \\ \Theta^t \quad \text{error} = 0 \\ \Theta^t + \alpha \phi(x^{t+1}) \text{ error} = 1 \end{cases}$$

$$\text{Thus } \Theta^{t+1} = \Theta^t + \beta_{t+1} \phi(x^{t+1}) \text{ where } \beta \in \{-1, 0, 1\}$$

$$\Theta^t = \Theta^{t-1} + \beta_t \phi(x^t)$$

$$\Theta^1 = \Theta^0 + \beta_1 \phi(x_1)$$

Adding we get

$$\Theta^{t+1} = \beta_1 \phi(x_1) + \beta_2 \phi(x_2) + \dots + \beta_{t+1} \phi(x^{t+1})$$

where  $\beta_i \in \{-1, 0, 1\}$

(2)

From derivation above we can clearly see that

$\Theta^{(i)}$  can be represented as (implicitly) the linear combinations of  $\phi(x^1), \phi(x^2), \dots, \phi(x^i)$  i.e

$$\Theta^{(i)} = \sum_{k=1}^i \beta_k \phi(x^k) \text{ where } \beta_k \in \{0, 1, -1\}$$

Hence  $\Theta$  can be represented as list of  $\beta$ 's.  
The value  $\Theta^{(0)}$  is initialised to  $\vec{0}$  and can be thought as base case where the summation ends as there are no summation terms included in  $\Theta^{(0)}$ .

(Note: this was helpful to end recursion / act as base case for  $\Theta^{(1)} = \Theta^{(0)} + \beta_1 \phi(x^1)$ )

→ We can say the list of  $\beta$ 's for  $\Theta^{(0)}$  is empty.

(2)

How to predict  $h_{\Theta^i}(x^{i+1})$

$$= g\left((\Theta^i)^T \phi(x^{i+1})\right)$$

$$\text{Since } \Theta^i = \sum_{k=1}^i \beta_k \phi(x^k) \quad (\Theta^i)^T = \sum_{k=1}^i \beta_k \phi^T(x^k)$$

$$= g\left(\underbrace{\sum_{k=1}^i \beta_k \phi^T(x^k)}_{\text{This dot product}} \phi(x^{i+1})\right) = g\left(\sum_{k=1}^i \beta_k K(x^k, x^{i+1})\right)$$

This dot product in high dimensional space can be computed as sum function of dot products on  $x^k, x^{i+1}$  in original space. Let that function be  $K$ .

Note: for prediction on any new point we only need  $\beta_k$  &  $K$ . i.e  $\Theta^i$  &  $K$ .

(3)

Since we are representing  $\Theta$  as list of  $\beta$ 's  
i.e.

$$\Theta^i = \{ \beta_1, \beta_2, \dots, \beta_i \} \quad \text{where } \beta_k \in \{-1, 0, 1\}$$

Updating  $\Theta$  to  $\Theta^{i+1}$  means getting new value of  
 $\beta$  i.e.  $\beta_{i+1}^o$

Now we know that

$$\Theta^{i+1} = \Theta^i + \alpha \left[ y^{i+1} - h_{\Theta^i} (g \phi(x^{i+1})) \right] \phi(x^{i+1})$$

↓  
This gives us  $\beta$

$$\therefore \beta^{i+1} = y^{i+1} - h_{\Theta^i} (g \phi(x^{i+1}))$$

↓  
where 2nd term is calculated  
using method in part 2  
using  $\{ \alpha, \beta_2, \dots, \beta_i \}$  & Kernel  
function  $K$ .

$$\therefore \Theta^{i+1} = \{ \alpha, \beta_2, \dots, \beta_i, \beta_{i+1}^o \}$$