



Assignment 0.

COL 870.

By:-

Ankit Gaug



Section I : Linear Algebra

1.1.

$A \rightarrow m \times n$ matrix

$B \rightarrow n \times m$ matrix

B is obtained by rotating 90° clockwise on paper

By taking examples I observed that A can be converted into B^T by multiplying a matrix H to A .

Here $H \rightarrow n \times n$ matrix

such that $H = \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}_{n \times n}$

$$h_{ij} = \begin{cases} 1 & \text{if } i = n - 1 - j \\ 0 & \text{else.} \end{cases}$$

$$HA = BT$$

Now, According to SVD any matrix \mathbf{Q} can be

Written as $A = U \Sigma V^T$

$$\begin{matrix} A \\ \downarrow m \times n \end{matrix} = \begin{matrix} U \\ \downarrow m \times m \end{matrix} \begin{matrix} \Sigma \\ \downarrow m \times m \end{matrix} \begin{matrix} V^T \\ \downarrow n \times n \end{matrix} \quad - \textcircled{1}$$

$\Sigma \Rightarrow$ diagonal and U & V are orthogonal.

From $\textcircled{1}$, We can write

$$A = U_A \Sigma_A V_A^T \quad - \textcircled{2}$$

Multiply H on both sides

$$HA = HU_A \Sigma_A V_A^T$$

$$B^T = HU_A \Sigma_A V_A^T$$

$$B = (HU_A \Sigma_A V_A^T)^T = V_A \Sigma_A^T U_A^T H^T$$

B =

$$\text{We get } B = V_A \Sigma_A^T U_A^T H^T \quad - \textcircled{3}$$

From $\textcircled{2}$ we already have that V_A is orthogonal, Σ_A is diagonal

Σ_A is diagonal $\Rightarrow \Sigma_A^T$ is also diagonal.

Now, if we prove $U_A^T H^T$ is orthogonal, then we could say that $\textcircled{3}$ is SVD decomposition

$$\text{Now, } H^T = H \text{ and } H^2 = I$$

$$\begin{aligned} (U_A^T H^T) (U_A^T H^T)^T &= U_A^T H^T H U_A \\ &= U_A^T H^2 U_A = U_A^T U_A \\ &= I \end{aligned}$$

because U is orthogonal.

We can say, (3) is a SVD decomposition

$$\underline{U_B} = B = U_B \Sigma_B V_B^T$$

where $U_B = V_A$

$$\Sigma_B = \Sigma_A^T \quad \text{--- (4)}$$

$$V_B^T = U_A^T H^T$$

From (4), we can say that A and B have same singular values.

1.2.

$x \rightarrow$ ~~for~~ $m \times 1$ vector

$A \rightarrow m \times n$ matrix

$$(a) \text{Phone} \quad \|x\|_\infty \leq \|x\|_1$$

$$\|x\|_\infty = \lim_{P \rightarrow \infty} \left(\sum_{i=1}^m |x_i|^P \right)^{\frac{1}{P}}$$

$$\text{Let } |x_{\max}| = \max_{i=1 \rightarrow m} (|x_i|)$$

$$\|x\|_\infty = \lim_{P \rightarrow \infty} |x_{\max}| \left(\sum_{i=1}^m \left| \frac{x_i}{x_{\max}} \right|^P \right)^{\frac{1}{P}}$$

$$= \lim_{P \rightarrow \infty} (|x_{\max}|) \cdot \underbrace{\left(\sum_{i=1}^m \left| \frac{x_i}{x_{\max}} \right|^P \right)^{\frac{1}{P}}} \quad \text{This limit is 1 because if it is of the form}$$

(1)

$$\text{Hence } \|x\|_\infty = |x_{\max}|$$

$$\text{Now } \|x\|_2 = \left(\sum_{i=1}^m x_i^2 \right)^{1/2} = |x_{\max}| \times \left[1 + \underbrace{\frac{\sum_{i=1}^m x_i^2}{|x_{\max}|^2}}_{\leq m} \right]^{1/2}$$

because due to max
to zero

$$\geq |x_{\max}|$$

$$\text{Hence. } \|x\|_2 \geq \|x\|_\infty$$

$$\text{or } \|x\|_\infty \leq \|x\|_2$$

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

$$\begin{aligned} \|x\|_2 &= \left(\sum_{i=1}^m x_i^2 \right)^{1/2} \leq \left(\sum_{i=1}^m x_{\max}^2 \right)^{1/2} \\ &\leq (m x_{\max}^2)^{1/2} \\ &\leq \sqrt{m} |x_{\max}| \\ &\leq \sqrt{m} \|x\|_\infty \end{aligned}$$

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

$$(c) \quad \text{Prove } \|A\|_\infty \leq \sqrt{n} \|A\|_2 \quad ; \quad \begin{matrix} y \rightarrow n \times 1 \text{ vector} \\ A \rightarrow m \times n \text{ matrix} \end{matrix}$$

$$\text{matrix norm} \rightarrow \|A\|_p = \max_{y \neq 0} \frac{\|Ay\|_p}{\|y\|_p} \quad \text{--- (1)}$$

From (1) we can say for any vector y .

$$\|A\|_p \geq \frac{\|Ay\|_p}{\|y\|_p}$$

$$\|A\|_p \|y\|_p \geq \|Ay\|_p - \textcircled{1}$$

$$\frac{\|Ay\|_p}{\|y\|_p} / \leq$$

Now from part (a) of this question we can say

$$\|Ay\|_\infty \leq \|Ay\|_2$$

Using 1 we get

$$\leq \|A\|_2 \|y\|_2$$

from part (b) of this question we can say

$$\leq \|A\|_2 \sqrt{n} \|y\|_\infty$$

$$\|Ay\|_\infty \leq \|A\|_2 \sqrt{n} \|y\|_\infty$$

$$\frac{\|Ay\|_\infty}{\|y\|_\infty} \leq \|A\|_2 \sqrt{n}$$

↳ This is true for any arbitrary y . It must also hold.

for y_s for which you
get $\max_y \frac{\|Ay\|_\infty}{\|y\|_\infty}$

Hence.

$$\|A\|_\infty \leq \sqrt{n} \|A\|_2$$

Hence proved.

$$(d) \text{ Prove } \|A_2\| \leq \sqrt{m} \|A\|_{\infty}$$

By defn:-

$$\|A\|_2 = \max_{\|y\|_2=1} \|Ay\|_2.$$

For any arbitrary y , for which $\|y\|_2=1$

$$\|Ay\|_2. \text{ Now let } A = \begin{bmatrix} \vec{a}_1 \\ \vec{a}_2 \\ \vdots \\ \vec{a}_m \end{bmatrix}$$

* These are now vectors

$$\|Ay\|_2 = \left(\sum_{i=1}^m (\vec{a}_i \cdot \vec{y})^2 \right)^{\frac{1}{2}}$$

For i th term.

$$(\vec{a}_i \cdot \vec{y})^2 = (\sum_{j=1}^n a_{ij} y_j)^2$$

Using Cauchy-Swartz Inequality

$$\leq (\sum_{j=1}^n a_{ij}^2)^{\frac{1}{2}} (\sum_{j=1}^n y_j^2)^{\frac{1}{2}}$$

$$\leq \sum_{j=1}^n a_{ij}^2. \quad [\text{Since } \|y\|_2=1]$$

$$\|Ay\|_2 \leq \left(\sum_{i=1}^m \sum_{j=1}^n a_{ij}^2 \right)^{\frac{1}{2}} \quad \dots \quad (1)$$

Now, $\|Ay\|_{\infty} = \lim_{p \rightarrow \infty} \left(\sum_{i=1}^m (\vec{a}_i \cdot \vec{y})^p \right)^{\frac{1}{p}}$

$$\|A\|_{\infty} = \max_{\|x\|_{\infty}=1} \|Ax\|_{\infty} = \max_{\|x\|_{\infty}=1} \max_{1 \leq i \leq m} |a_i^T x| = \max_{\|x\|_{\infty}=1} \max_{1 \leq i \leq m} \sum_{j=1}^n |a_{ij}| x_j = \max_{\|x\|_{\infty}=1} \sum_{j=1}^n \max_{1 \leq i \leq m} |a_{ij}| x_j = \max_{\|x\|_{\infty}=1} \sum_{j=1}^n \|a_j\|_{\infty} x_j = \max_{\|x\|_{\infty}=1} \sum_{j=1}^n \|a_j\|_{\infty}$$

We will use the fact that $\|A\|_{\infty} = \text{maximum absolute row sum}$

$$\sqrt{m} \|A\|_{\infty} = \sqrt{m \|A\|_{\infty}^2} \approx \sqrt{m \left(\sum_{j=1}^n |a_{ij}| \right)^2}$$

Let K be the row corresponding to maximum row sum

$$= \sqrt{m \left(\sum_{j=1}^n |a_{kj}| \right)^2} \Rightarrow \sqrt{\sum_{j=1}^m \left(\sum_{i=1}^n |a_{ij}| \right)^2}$$

Now square of sum of positive values is greater than sum of square of those values.

$$\geq \sqrt{\sum_{i=1}^m \sum_{j=1}^n |a_{ij}|^2} \quad \text{--- (2)}$$

Combining (1) & (2) we can say

$$\|Ay\|_2 \leq \sqrt{m} \|A\|_{\infty}.$$

True for any arbitrary y , must be true for the one $\sup_{\|y\|_2=1} \|Ay\|_2$

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

Hence proved

Section 0 : Subgradients.

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

Let n^{th} line be the one, for which.

$$\max_{i=1, \dots, m} (a_i^T x + b_i) = a_n^T x + b_n \quad \text{--- (1)}$$

If there is only one ~~line~~ curve for which (1) is satisfied

the $\boxed{\text{Subgradient } g f(x) = a_n}$

If there are more than one curves, then the curve is non differentiable. We can choose any ~~between left and right derivative~~ point lies between the derivatives around that point

$\nabla f_1 \rightarrow$ be the curve on left.

$$\alpha_1 \nabla f_1 + \alpha_2 \nabla f_2 \dots \alpha_k \nabla f_k ; \alpha_1 + \dots + \alpha_k = 1$$

These k curves pass through the point and contribute to the derivative in some direction.

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

Let n^{th} line be the one for which.

$$\max_i |(a_i^T x + b_i)| = |a_n^T x + b_n|$$

if $a_n^T x + b_n > 0$ then gradient = a_n

if $a_n^T x + b_n < 0$ then gradient = $-a_n$

else $\alpha_1(a_n) + \alpha_2(-a_n)$; where $\alpha_1 + \alpha_2 = 1$.

If more than one curve passes through this point.
Then look at the derivatives at this point, in various directions and take value in-between.

Let f_1, f_2, \dots, f_n be the curves.

$$\text{then } g(x) = \alpha_1 \nabla f_1(x) + \dots + \alpha_n \nabla f_n(x)$$

$$\alpha_1 + \alpha_2 + \dots + \alpha_n = 1$$

$$(c) f(x) = \sup_{0 \leq t \leq 1} p(t) \text{ where } p(t) = x_1 + x_2 t + \dots + x_n t^{n-1}$$

Let t_0 be the value for which we get $\sup_{0 \leq t \leq 1} p(t)$ for given x

$$f(x) = x_1 + x_2 t_0 + \dots + x_n t_0^{n-1}$$

$$\nabla f = [1 \ t_0 \ \dots \ t_0^{n-1}]$$

\Rightarrow subgradient

$$(d) f(x) = x_{[1]} + x_{[2]} + \dots + x_{[k]}$$

$x_{[i]} \rightarrow$ i-th largest element.

$$\nabla f(x) = 1 \quad ; \quad \nabla f(x) = 0$$

$$x_{[i]}, i \leq k \qquad \qquad x_{[l]}, l > n$$

Hence subgradient = $\begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \\ 1 \end{bmatrix} \Rightarrow$ Value is 1, if it is one of the k largest values.
else zero

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

For given x , we have to minimize $\|x - y\|^2$, subject to the constraints $Ay \leq b$.

So this is constrained optimization, with

We can use KKT to solve this

$$\|x - y\|^2 + \lambda(Ay - b)$$

Create dual and then solve. Don't know functions!

Section ① : Probability.

1. $E[\mathbb{I} | x=x] = ?$

$$Z \sim N(0, 1)$$

$$X \sim N(\mu, 1)$$

X and Z are independent

and $\Phi(z)$ is the cdf of gaussian distribution.

$$\mathbb{I} = \begin{cases} 1 & ; \text{ if } z < x \\ 0 & ; \text{ if } z \geq x \end{cases}$$

$$P(\mathbb{I}=1) = P(Z < x)$$

$$P(\mathbb{I}=0) = P(Z \geq x)$$

$$E(\mathbb{I} | x=x) = 1 \cdot P(\mathbb{I}=1 | x=x) + 0 \cdot P(\mathbb{I}=0 | x=x)$$

$$= P(Z < x) = \Phi(x) \quad \left[\begin{array}{l} \text{By the definition} \\ \text{of cdf} \end{array} \right]$$

Q. $E(\Phi(x)) \neq ?$

Using the rule of conditional expectation

$$= E(E(\Phi(x) | x=x))$$

$$= E(E(P(Z < x)))$$

\rightarrow By the definition
of cdf.

2. In the previous part we derived

$$E(I|X=x) = \Phi(x) \quad \text{--- (1)}$$

We can write

$$E(I) = 1 \cdot P(Z < X) + 0 \cdot P(Z > X)$$

$$\text{Hence } E(I) = P(Z < X)$$

$$E(E(I|X=x)) = P(Z < X)$$

Using the law of conditional expectation

Using (1).

$$E(\Phi(X)) = P(Z < X)$$

Hence proved

3. In the previous part

$$E(\Phi(X)) = P(Z < X) = P(0 < X - Z)$$

$$= P(X - Z > 0)$$

$X \sim N(\mu, 1)$; $Z \sim N(0, 1)$ are independent

their diff're is gaussian $\sim N(\mu - \sigma^2, \sigma^2)$

$$\text{let } Y = X - Z$$

$$P(Y > 0) = P\left(\frac{Y - \mu}{\sigma} > \frac{-\mu}{\sigma}\right) = 1 - \Phi\left(\frac{-\mu}{\sigma}\right)$$

$$= \Phi\left[\frac{\mu}{\sqrt{2}}\right]$$

Property of cdf function

Section P : Machine Learning.

4.1

PCA

$$\max_{\mathbf{u}: \|\mathbf{u}\|_2=1} \tilde{V}[\mathbf{u}^T \mathbf{x}]$$

 $\mathbf{u} \rightarrow p \times 1$ vector. $\mathbf{x} \rightarrow p \times N$ matrix.

$$\text{Use } V_{\mathbf{u}}(\mathbf{x}) = E(x^2) - E(x)^2$$

$$\tilde{V}[\mathbf{u}^T \mathbf{x}] = E(\underline{(\mathbf{u}^T \mathbf{x}_i)^2}) - E(\mathbf{u}^T \mathbf{x}_i)^2$$

$$= \sum_{i=1}^N \frac{(\mathbf{u}^T \mathbf{x}_i)^2}{N} - \left[\frac{\sum_{i=1}^N \mathbf{u}^T \mathbf{x}_i}{N} \right]^2.$$

$$= \sum_{i=1}^N \frac{\mathbf{u}^T \mathbf{x}_i \mathbf{x}_i^T \mathbf{u}}{N} - \underbrace{\left[\mathbf{u}^T \sum_{i=1}^N \frac{\mathbf{x}_i}{N} \right]^2}_{>0, \text{ since data is zero mean.}}$$

$$= \mathbf{u}^T \left(\sum_{i=1}^N \frac{\mathbf{x}_i \mathbf{x}_i^T}{N} \right) \mathbf{u}$$

$$\text{On comparing } \cancel{\sum} = \sum_{i=1}^N \frac{\mathbf{x}_i \mathbf{x}_i^T}{N}$$

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

$$= \min \frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i - \mathbf{u}^T \mathbf{x}_i)^T (\mathbf{x}_i - \mathbf{u}^T \mathbf{x}_i)$$

Now $\mu^T x_i = x_i^T \mu = k$ [scalar quantity]

$$= \min \frac{1}{N} \sum_{i=1}^N (x_i - \mu k)^T (x_i - \mu k)$$

$$= \min \frac{1}{N} \sum_{i=1}^N (x_i^T x_i - \mu^T x_i k - x_i^T \mu k + \mu^T \mu)$$

$$\mu^T x_i = k$$

$$= \min \frac{1}{N} \sum_{i=1}^N (x_i^T x_i - \underbrace{\mu^2 - \mu^2}_{\text{independent of } \mu} + \mu^2)$$

$$= \min \frac{1}{N} \sum_{i=1}^N (-k^2) = \max \frac{1}{N} (\mu^2)$$

$$= \max \frac{1}{N} \sum_{i=1}^N (\mu^T x_i x_i^T \mu)$$

$$= \mu^T \left(\frac{\sum_{i=1}^N x_i x_i^T}{N} \right) \mu$$

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

Hence proved.

4.2.

Bias - Variance Trade-off.

(a) $(x, y) \rightarrow$ new data point

$$J(\theta) = (y - h_{\theta}(x))^2$$

$h_{\theta}(x) \rightarrow$ classifier trained over a dataset D .

D which is dataset is random variable in this question

$$E_D [J(\theta)] = E_D [(y - h_{\theta}(x))^2]$$

$$= E_D [(y - h_{\theta}(x) + h_{\bar{\theta}}(x) - h_{\theta}(x))^2]$$

$$= E_D [(y - h_{\theta}(x))^2 + (h_{\bar{\theta}}(x) - h_{\theta}(x))^2 + 2(y - h_{\theta}(x))(h_{\bar{\theta}}(x) - h_{\theta}(x))]$$

$$= E_D [(y - h_{\theta}(x))^2] + E_D [(h_{\bar{\theta}}(x) - h_{\theta}(x))^2] + E_D [2(y - h_{\theta}(x))(h_{\bar{\theta}}(x) - h_{\theta}(x))]$$

$$= (y - h_{\theta}(x))^2 + E_D [(h_{\bar{\theta}}(x) - h_{\theta}(x))^2] + 2(y - h_{\theta}(x))(h_{\bar{\theta}}(x) - h_{\theta}(x))$$

$$= 0$$

$$= \text{Bias}^2 + \text{Variance}$$

(b) Now for noise

$$J(\theta) = (y - h_{\theta}(x))^2 = (f(x) + \epsilon - h_{\theta}(x))^2$$

$$= (f(x) - h_{\theta}(x))^2 + \epsilon^2 + 2\epsilon(f(x) - h_{\theta}(x))$$

$$E_D [J(\theta)] = E_D [(f(x) - h_{\theta}(x))^2] + E_D (\epsilon^2)$$

$$+ E_D (2\epsilon (f(x) - h_{\theta}(x)))$$

Using results from the previous part

$$= \text{Bias}^2 + \text{Variance} + \text{Var}(\epsilon) - E_D (\epsilon)$$

$$+ 2 E_D (\epsilon) E (f(x) - h_{\theta}(x))$$

$$\text{since } \epsilon \sim N(0, \sigma^2) \Rightarrow 0$$

$$= \text{Bias}^2 + \text{Variance} + \sigma^2 - 0 + 0$$

$$= \text{Bias}^2 + \text{Variance} + \sigma^2$$

4.3.

let's make some observation, then I will describe the 3 parts.

We are doing only one epoch and we take one example at a time.

Therefore we have

$\theta^{(0)}, \theta^{(1)}, \dots, \theta^{(m)}$ parameters
 [assuming dataset has m examples]

According to perceptron

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

Note that this term is either 0, α or $-\alpha$.

Depending upon whether $(i+1)$ th example was misclassified, and what was its actual label.

Let this value be $\gamma^{(i+1)}$ corresponds to $y^{(i+1)}$

$$\text{Then } \theta^{(i+1)} = \theta^{(i)} + \gamma^{(i+1)} x^{(i+1)} \quad \text{--- (1)}$$

where $\gamma^{(i+1)} \in \{-\alpha, 0, \alpha\}$

From (1)

for any $i \in \{0, N\}$ we can say are of the form

$$\theta^{(i)} = \sum_{j=1}^i \gamma^{(j)} x^{(j)} + \theta^{(0)} \rightarrow \vec{\theta}$$

So $\theta^{(i)}$ can be represented as in the basis formed by examples i.e. $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

For example

$$\theta^{(i)} = \begin{bmatrix} \alpha \\ -\alpha \\ 0 \\ \vdots \\ \alpha \end{bmatrix}_{m \times 1} \rightarrow \text{is the representation.}$$

Notice that this representation is independent of the dimension of x or θ .

1. Now, in transformed space, from the above arguments.

$$\theta^{(i)} = \sum_{j=1}^i y_j^{(i)} \phi(x^{(j)})$$

$\theta^{(i)}$ can be represented on the basis

$$\text{Basis} = [\phi(x^{(1)}), \phi(x^{(2)}), \dots, \phi(x^{(m)})] \quad \hookrightarrow \textcircled{1}$$

$\theta^{(i)}$ will be represented as ($M \times 1$) vector.

Note that, the representation is independent of the dimension of $\theta^{(i)}$ or even $\phi(x^{(i)})$.

For example

$$\theta^{(i)} = [0 \ 0 \ 0 \ 0 \ \dots \ 0]_{1 \times m}$$

$$\theta^{(i)} = [y_1^{(i)} \ 0 \ 0 \ \dots \ 0]_{1 \times m}$$

$$\theta^{(i)} = [y_1^{(i)} \ y_2^{(i)} \ 0 \ 0 \ \dots \ 0]_{1 \times m}$$

$$\theta^{(i)} = [y_1^{(i)} \ y_2^{(i)} \ \dots \ y_m^{(i)}]_{1 \times m}$$

The above representation holds when basis is in Eq. ①.

2. Now, we have already represented $\theta^{(i)}$ in our own defined basis.

For prediction on $\theta^{(i+1)}$ example.

$$h_{\theta^{(i)}}(x^{(i+1)}) = g(\theta^{(i)T} \phi(x^{i+1})) \quad \text{--- (2)}.$$

Now $\theta^{(i)} = \sum_{j=1}^i \gamma_j^{(i)} \phi(x^{i+1})$ in our basis.

(2) becomes.

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

all the γ^1 to γ^i have been calculated till now

$$\text{So. } \phi^T(x^j) \phi(x^{i+1}) = K(x^j, x^{i+1})$$

\rightarrow Mercer's Kernel
for this basis function

$$= g\left(\sum_{j=1}^i \gamma_j^{(i)} K(x^j, x^{i+1})\right)$$

We can compute this.

and hence can make our prediction.

3.

When at $(i+1)$ example, we have already calculated $\gamma^{(i)}$ to $\gamma^{(i)}$, and we need to calculate $\gamma^{(i+1)}$ at this step.

$$\theta^{(i+1)} = \theta^{(i)} + \alpha (y^{(i+1)} - g(\theta^{(i)T} \phi(x^{(i+1)})) \phi(x^{(i)})$$

$$\sum_{j=1}^{i+1} \gamma^j \phi(x^{(j)}) = \sum_{j=1}^i \gamma^j \phi(x^{(j)}) + \underbrace{\alpha (y^{(i+1)} - g(\theta^{(i)T} \phi(x^{(i+1)})) \phi(x^{(i)}))}_{\text{We say how to compute this in the last part}}$$

$$\sum_{j=1}^{i+1} \gamma^j \phi(x^{(j)}) - \sum_{j=1}^i \gamma^j \phi(x^{(j)}) = \alpha (y^{(i+1)} - g(\theta^{(i)T} \phi(x^{(i+1)})) \phi(x^{(i+1)}))$$

$$\gamma^{(i+1)} \phi(x^{(i+1)}) = \alpha (y^{(i+1)} - g(\theta^{(i)T} \phi(x^{(i+1)})) \phi(x^{(i+1)}))$$

$$\boxed{\gamma^{(i+1)} = \alpha (y^{(i+1)} - g(\theta^{(i)T} \phi(x^{(i+1)})))}$$

Update Rule.

Calculation of prediction has been discussed in the previous part!