

1. Gradient descent for Kernel regression

Consider given $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^m$

$$\text{and } f(x; \theta) = \frac{1}{\sqrt{m}} \theta^T \phi(x),$$

From these two, we can infer that $x \in \mathbb{R}^d$, $\phi(x) \in \mathbb{R}^m$, $\theta \in \mathbb{R}^m$

$$L(\theta) = \frac{1}{2} \|f(X; \theta) - y\|^2$$

where $f(X; \theta) =$ is a vector of predictions =

$$\begin{bmatrix} f(x_1; \theta) \\ f(x_2; \theta) \\ \vdots \\ f(x_n; \theta) \end{bmatrix}$$

(Assume n data points)

$(n \times 1)$

$$\frac{\partial f(x; \theta)}{\partial \theta} = \frac{1}{\sqrt{m}} \phi(x) \in \mathbb{R}^m \quad \text{--- (1)}$$

$$\frac{d\theta_t}{dt} = - \frac{\partial L(\theta_t)}{\partial \theta_t} = - \underbrace{\frac{\partial f(x; \theta_t)}{\partial \theta}}_{\text{Jacobian matrix of shape } (m \times m)} \cdot \underbrace{(f(x; \theta_t) - y)}_{\text{Vector of shape } (n \times 1)} \quad \text{--- (2)}$$

need to derive, $\frac{d}{dt} f(x; \theta_t)$.

Consider, $\frac{d}{dt} f(x; \theta_t)$

using chain rule,

$$= \sum_{i=1}^M \underbrace{\frac{\partial f(x; \theta_t)}{\partial \theta_i}}_{\text{scalar}} \cdot \underbrace{\frac{d\theta_i}{dt}}_{\text{scalar}},$$

$$= \left\langle \underbrace{\frac{\partial f(x; \theta_t)}{\partial \theta}}_{\text{Vector}}, \underbrace{\frac{d\theta}{dt}}_{\text{Vector}} \right\rangle$$

$$= \left(\frac{\partial f(x; \theta_t)}{\partial \theta} \right)^T \left(\frac{d\theta}{dt} \right)$$

this is replaced by $\textcircled{1}$ replace by $\textcircled{2}$

$$= \frac{1}{\sqrt{m}} \phi(x)^T \cdot \underbrace{\left(-\frac{\partial f(x; \theta_t)}{\partial \theta} \right)}_{m \times n \text{ matrix}} \left(f(x; \theta_t) - y \right) \quad \text{--- } \textcircled{3}$$

need to ~~compute~~ compute, $\frac{\partial f(x; \theta_t)}{\partial \theta}$

$$\text{Entries in } \frac{\partial f(x; \theta_t)}{\partial \theta} = \begin{pmatrix} \frac{\partial f(x_1; \theta_t)}{\partial \theta_1} & \frac{\partial f(x_2; \theta_t)}{\partial \theta_2} & \dots & \frac{\partial f(x_n; \theta_t)}{\partial \theta_1} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f(x_1; \theta_t)}{\partial \theta_m} & \dots & \dots & \frac{\partial f(x_n; \theta_t)}{\partial \theta_m} \end{pmatrix}$$

(3)

~~Each column~~take i^{th} column in the above matrix;

$$\begin{bmatrix} \frac{\partial f(x_i; \theta_t)}{\partial \theta_1} \\ \frac{\partial f(x_i; \theta_t)}{\partial \theta_2} \\ \vdots \\ \frac{\partial f(x_i; \theta_t)}{\partial \theta_m} \end{bmatrix}$$

→ this can be replaced
by, $\frac{1}{\sqrt{m}} \phi(x_i)$
(from (1))

$$\Rightarrow \frac{\partial f(x; \theta_t)}{\partial \theta} = \begin{bmatrix} \phi(x_1) & \dots & \phi(x_i) & \dots & \phi(x_n) \\ | & & | & & | \end{bmatrix} \cdot \frac{1}{\sqrt{m}}$$

(m × n)

define the above matrix is same as $\left(\frac{1}{\sqrt{m}} \phi(X) \right)$

Substituting in (3)

$$\begin{aligned} \Rightarrow \frac{d}{dt} f(x; \theta_t) &= \frac{1}{\sqrt{m}} \phi^T(x) - \frac{1}{\sqrt{m}} \phi(X) \cdot (f(x; \theta_t) - y) \\ &= -\frac{1}{m} \phi^T(x) \phi(X) \cdot (f(x; \theta_t) - y) \end{aligned}$$

B.

Need to compute, $\frac{d}{dt} f(x; \theta)$

$$\text{Consider, } \underbrace{\frac{d}{dt} f(x; \theta)}_{n \times 1 \text{ vector}} = \frac{d}{dt} \begin{bmatrix} f(x_1; \theta) \\ f(x_2; \theta) \\ \vdots \\ f(x_n; \theta) \end{bmatrix} \triangleq \begin{bmatrix} \frac{d}{dt} f(x_1; \theta) \\ \frac{d}{dt} f(x_2; \theta) \\ \vdots \\ \frac{d}{dt} f(x_n; \theta) \end{bmatrix}$$

(5)

we can use the result from part A, which is,

$$\frac{d}{dt} f(x; \theta_t) = -\frac{1}{m} \cdot \phi(x)^T \phi(x) \cdot (f(x; \theta_t) - y)$$

here x could be any data point i . So we substitute in (5)

$$\frac{d}{dt} f(x; \theta) = \begin{bmatrix} -\frac{1}{m} \phi(x_1)^T \phi(x) \cdot (f(x; \theta) - y) \\ -\frac{1}{m} \phi(x_2)^T \phi(x) \cdot (f(x; \theta) - y) \\ \vdots \\ -\frac{1}{m} \phi(x_n)^T \phi(x) \cdot (f(x; \theta) - y) \end{bmatrix}$$

⑤

$$= \begin{bmatrix} -\frac{1}{n} \phi(x_1)^T \phi(x) \\ -\frac{1}{n} \phi(x_2)^T \phi(x) \\ \vdots \\ -\frac{1}{n} \phi(x_n)^T \phi(x) \end{bmatrix} \begin{bmatrix} | \\ p(x; \theta) - y \\ | \end{bmatrix}$$

take i^{th} row of this matrix,

$$\begin{aligned} & -\frac{1}{n} \phi(x_i)^T \phi(x) \\ &= -\frac{1}{n} \phi(x_i)^T \begin{bmatrix} | & & | & & | \\ \phi(x_1) & \dots & \phi(x_j) & \dots & \phi(x_n) \\ | & & | & & | \end{bmatrix} \\ &= -\frac{1}{n} \begin{bmatrix} \phi(x_i)^T \phi(x_1), & \phi(x_i)^T \phi(x_2) & \dots & \phi(x_i)^T \phi(x_j) & \dots & \phi(x_i)^T \phi(x_n) \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix} \end{aligned}$$

(6)

$$\Rightarrow \frac{d}{dt} f(x; \theta) = \frac{1}{n} \begin{bmatrix} \phi(x_1)^T \phi(x_1) & \phi(x_1)^T \phi(x_2) & \dots & \phi(x_1)^T \phi(x_n) \\ \vdots & \vdots & \ddots & \vdots \\ \phi(x_n)^T \phi(x_1) & \dots & \dots & \phi(x_n)^T \phi(x_n) \end{bmatrix}$$

A , $(i, j)^{th}$ element in this matrix,

$$= \frac{1}{n} \phi(x_i)^T \phi(x_j) = k_{i,j}$$

\otimes $\begin{bmatrix} 1 \\ f(x; \theta) - y \\ 1 \end{bmatrix}$

$$\Rightarrow \boxed{\frac{d}{dt} f(x; \theta) = -k (f(x; \theta) - y)}$$

— (6)

$$e) \left. \begin{aligned} \frac{d}{dt} v_{\pm} &= -k v_{\pm} \end{aligned} \right\} v_{\pm} = e^{-\pm k t} v_0$$

to use the above ~~part~~ formula, vector v_{\pm} should be some on left & right,

(7)

denote, $(f(x; \theta) - y) = \underline{a}$

$$\Rightarrow \frac{d}{dt} f(x; \theta) = \frac{da}{dt}$$

the above,
Substitute in equation (6)
 x

$$\Rightarrow \frac{da}{dt} = -k \underline{a}$$

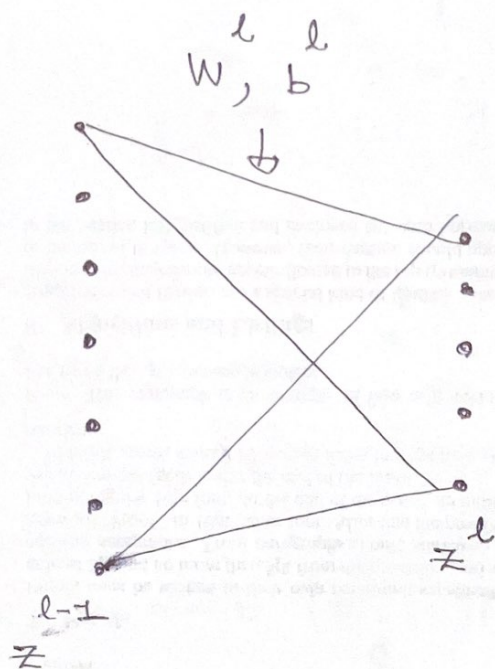
$$\Rightarrow \underline{a} = e^{-tk} a_0$$

$$\Rightarrow (f(x; \theta) - y) = e^{-tk} a_0$$

at $t=0$ $a = f(x; \theta_0) - y$

$$\Rightarrow f(x; \theta) - y = e^{-tk} (f(x; \theta_0) - y)$$

$$\Rightarrow f(x; \theta) = y + e^{-tk} (f(x; \theta_0) - y)$$



$$z^l(x) = \frac{\sigma_w}{\sqrt{n^{l-1}}} \cdot W^l \cdot z^{l-1}(x) + \sigma_b b^l$$

$$\begin{bmatrix} z_i^{l-1}(x) \\ \vdots \\ z_i^{l-1}(x') \end{bmatrix} \sim N \left(0, \begin{bmatrix} \Sigma^{l-1}(x, x) & \Sigma^{l-1}(x, x') \\ \Sigma^{l-1}(x, x') & \Sigma^{l-1}(x', x') \end{bmatrix} \right)$$

given z_i^{l-1} and z_j^{l-1} are independent for $i \neq j$

Consider j -th component of $z^l(x)$,

$$z_j^l(x) = \frac{\sigma_w}{\sqrt{n^{l-1}}} \cdot \underbrace{W^l}_{\text{first row of } W^l} \cdot z^{l-1}(x) + \sigma_b b_j^l$$

$$z_j^l(x) = \frac{\sigma_w}{\sqrt{n^{l-1}}} \cdot \sum_{i=1}^{n^{l-1}} w_{ij}^l \phi(z_i^{l-1}(x)) + \sigma_b b_j^l$$

similarly,

$$z_j^l(x') = \frac{\sigma_w}{\sqrt{n^{l-1}}} \cdot \sum_{i=1}^{n^{l-1}} w_{ij}^l \phi(z_i^{l-1}(x')) + \sigma_b b_j^l$$

Question is about $(z_j^l(x))(z_j^l(x'))$

So consider $(z_j^l(x))(z_j^l(x'))$

$$= \left(\frac{\sigma_w}{\sqrt{n^{l-1}}} \cdot \sum_{i=1}^{n^{l-1}} w_{ij}^l \phi(z_i^{l-1}(x)) + \sigma_b b_j^l \right)$$

$$\left(\frac{\sigma_w}{\sqrt{n^{l-1}}} \cdot \sum_{i=1}^{n^{l-1}} w_{ij}^l \phi(z_i^{l-1}(x')) + \sigma_b b_j^l \right)$$

$$= \sigma_b^2 b_j^2 + \frac{\sigma_w^2}{n^{l-1}} \left(\sum_{i=1}^{n^{l-1}} w_{ij}^l \phi(z_i^{l-1}(x)) \right) \left(\sum_{i=1}^{n^{l-1}} w_{ij}^l \phi(z_i^{l-1}(x')) \right) + \sigma_b b_j^l \left(\frac{\sigma_w}{\sqrt{n^{l-1}}} \right) \left(\sum_{i=1}^{n^{l-1}} w_{ij}^l \phi(z_i^{l-1}(x)) + \sum_{i=1}^{n^{l-1}} w_{ij}^l \phi(z_i^{l-1}(x')) \right)$$

$$\sigma_b^2 + \left(\frac{\sigma_w^2}{n} \left[\begin{matrix} l & l-1 & l & l-1 & l \\ w_{1j} s(z_1(x)) & w_{1j} s(z_1(x')) & + w_{1j} s(\dots) & w_{2j} s(\dots) \\ \vdots & \vdots & \vdots & \vdots & \vdots \end{matrix} \right] + \dots \right)$$

$$\sigma_b^2 \left(\frac{\sigma_w^2}{\sqrt{n^{l-1}}} \right) \left(\sum_{i=1}^{n^{l-1}} w_{ij} s(z_i^{(l-1)}(x)) + \sum_{i=1}^{n^{l-1}} w_{ij} s(z_i^{(l-1)}(x')) \right)$$

$$E \left[\begin{matrix} l \\ z_j(x), z_j(x') \end{matrix} \right] = E \left[\begin{matrix} l^2 \\ \sigma_b^2 b_j \end{matrix} \right] + E \left[\begin{matrix} \text{ } \\ \text{ } \end{matrix} \right]$$

$$+ E \left[\begin{matrix} \text{ } \\ \text{ } \end{matrix} \right]$$

Given $b_j^l \sim N(0, 1) \Rightarrow E[b_j^l] = 0, E[b_j^{l^2}] = 1$

$w_{ij}^l \sim N(0, 1) \Rightarrow E[w_{ij}^l] = 0, E[w_{ij}^{l^2}] = 1$

$$b_j^l + \left(\frac{\sigma_w^2}{n^{l-1}} \left[w_{1j}^l s(z_1^l(w)) \cdot w_{1j}^{l-1} s(z_1^{l-1}(x)) + w_{1j}^l s(\dots) w_{2j}^{l-1} s(\dots) \right. \right. \\ \left. \left. + \dots \right] \right) + \sigma_b^l b_j^l \left(\frac{\sigma_w}{\sqrt{n^{l-1}}} \right) \left(\sum_{i=1}^{n^{l-1}} w_{ij}^l s(z_i^{l-1}(w)) \right. \\ \left. + \sum_{i=1}^{n^{l-1}} w_{ij}^l s(z_i^{l-1}(x)) \right)$$

$$\equiv \begin{bmatrix} z_j^l(x), z_j^l(x') \end{bmatrix} = E \begin{bmatrix} \sigma_b^2 b_j^{l^2} \end{bmatrix} + E \begin{bmatrix} \leftarrow \right] \\ + E \left[\leftarrow \right]$$

Given $b_j^l \sim N(0, 1) \Rightarrow E[b_j^l] = 0, E[b_j^{l^2}] = 1$

$w_{ij}^l \sim N(0, 1) \Rightarrow E[w_{ij}^l] = 0, E[w_{ij}^{l^2}] = 1$

$$\sigma_b^2 b_j^l + \left(\frac{\sigma_w^2}{n^{l-1}} \left[w_{1j}^l s(z_1^l(w)) \cdot w_{1j}^{l-1} s(z_1^{l-1}(x)) + w_{1j}^l s(\dots) w_{2j}^{l-1} s(\dots) + \dots \right] \right) + \sigma_b^2 b_j^l \left(\frac{\sigma_w}{\sqrt{n^{l-1}}} \right) \left(\sum_{i=1}^{n^{l-1}} w_{ij}^l s(z_i^{l-1}(w)) + \sum_{i=1}^{n^{l-1}} w_{ij}^l s(z_i^{l-1}(x)) \right)$$

$$E \left[z_j^l(x) \cdot z_j^l(x') \right] = E \left[\sigma_b^2 b_j^{l2} \right] + E \left[\begin{matrix} \text{ } \\ \text{ } \end{matrix} \right] + E \left[\begin{matrix} \leftarrow \end{matrix} \right]$$

Given $b_j^l \sim N(0, 1) \Rightarrow E[b_j^l] = 0, E[b_j^{l2}] = 1$

$w_{ij}^l \sim N(0, 1) \Rightarrow E[w_{ij}^l] = 0, E[w_{ij}^{l2}] = 1$

In the

and w_{ij}^l and z_i^{l-1} are independent

$$\Rightarrow E[w_{ij}^l z_i^{l-1}] = 0$$

or stronger statement $E[f(w_{ij}^l) \cdot g(z_i^{l-1})] = 0$

first term $= E[\sigma_b^2 b_j^{l^2}] = \sigma_b^2 E[b_j^{l^2}] = \sigma_b^2$

Second term $= E\left[\frac{\sigma_w^2}{n^{l-1}} \left[w_{1j}^l s(z_1^{l-1}(x)), w_{1j}^l s(z_1^{l-1}(x')) \right. \right.$
 $\left. + w_{1j}^l s(z_1^{l-1}(x)), w_{2j}^l s(z_2^{l-1}(x')) \right.$
 $\left. + \dots \right]$

$$= \frac{\sigma_w^2}{n^{l-1}} \left[E\left[\sum_{i=1}^{n^{l-1}} w_{ij}^l s(z_i^{l-1}(x)) \cdot s(z_i^{l-1}(x')) \right] + \left\{ w_{1j}^l s(z_1^{l-1}(x)) w_{2j}^l s(z_2^{l-1}(x')) \right. \right.$$

 $\left. + \dots \right\}]$

all these are cross terms

$$\frac{\sigma_w^2}{n^{l-1}} \cdot n \cdot \sum_{i=1}^n E \left[s(z_i^{l-1}(x)) \cdot s(z_i^{l-1}(x')) \right]$$

$$= \sigma_w^2 \cdot E \left[s(z_i^{l-1}(x))^T s(z_i^{l-1}(x')) \right]$$

third term :- Expectation is zero because of independence.

$$E \left[z_j^l(x) z_j^l(x') \right] = \text{first term} + \text{second term} + \text{third}$$

$$= \sigma_B^2 + \sigma_w^2 \cdot E \left[s(z_i^{l-1}(x))^T s(z_i^{l-1}(x')) \right]$$