

Module 2: Connections between NNs & Kernels.

Review : (1) Kernel Regression

$$x \mapsto \psi(x)$$

$\mathbb{R}^d \quad \mathbb{R}^P$

$$\omega \psi(x^{(i)}) = y^{(i)}$$

$$K(x, \tilde{x}) \in \mathbb{R} \rightarrow \langle \psi(x), \psi(\tilde{x}) \rangle$$



Algorithm : Given data (X, y) , kernel K

$$(1) K_{\text{train}} = K(X, X)$$

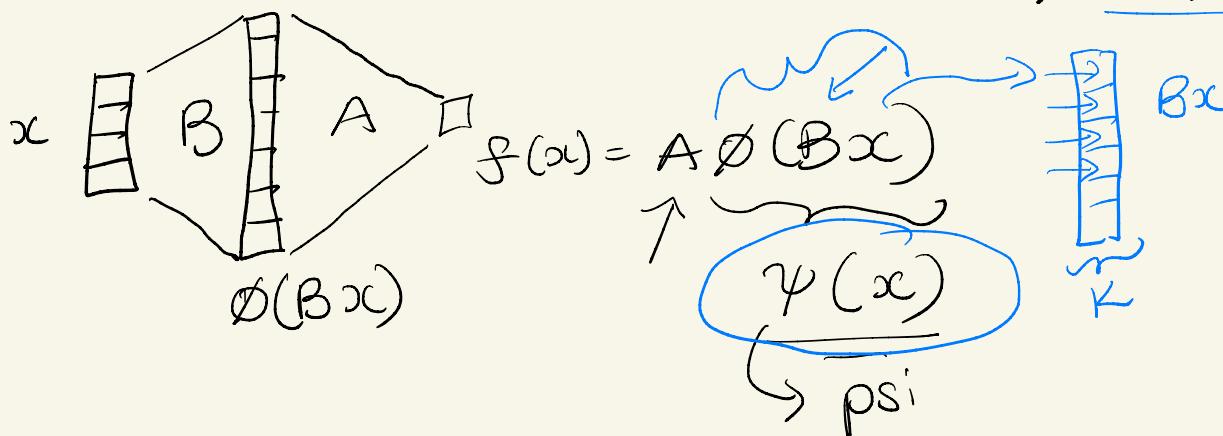
$$(2) \alpha = y K_{\text{train}}^{-1}$$

(2) Neural Networks. $f_w(x) = A \phi(Bx)$

1-hidden layer FC network ; $A \in \mathbb{R}^{1 \times k}$, $B \in \mathbb{R}^{k \times d}$

$$\phi: \mathbb{R} \rightarrow \mathbb{R} ; \phi(z) = \max(z, 0)$$

ReLU.



k : width

$k \rightarrow \infty$

$$f(x) = \underbrace{A \frac{1}{\sqrt{K}} \phi(Bx)}_{\psi(x)}$$

$$\left\{ (x^{(i)}, y^{(i)}) \right\}_{i=1}^n$$

$$\sum_{\text{Kernel}}(x, \tilde{x}) := \langle \psi(x), \psi(\tilde{x}) \rangle = \langle \frac{1}{\sqrt{K}} \phi(Bx), \frac{1}{\sqrt{K}} \phi(B\tilde{x}) \rangle$$

$$= \frac{1}{K} \sum_{i=1}^K \phi(Bx) \cdot \phi(B\tilde{x})$$

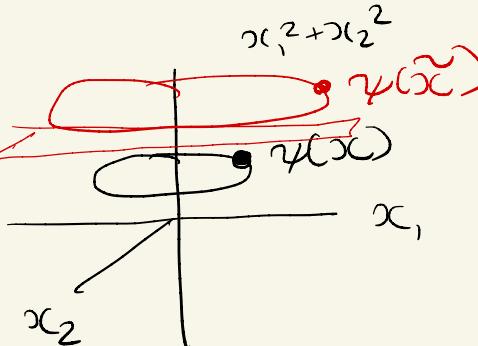
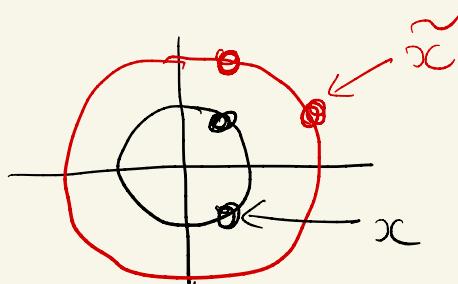
$$= \frac{1}{K} \sum_{i=1}^K \underbrace{\phi(B_i, x)}_{\psi(x)} \underbrace{\phi(B_i, \tilde{x})}_{\psi(\tilde{x})}$$

$$K \begin{bmatrix} B_1, : \\ B_2, : \\ \vdots \\ B_K, : \end{bmatrix} = B$$

As $K \rightarrow \infty$

$$\rightarrow \mathbb{E}_{\omega} [\underbrace{\phi(\omega^T x)}_{\psi(x)} \underbrace{\phi(\omega^T \tilde{x})}_{\psi(\tilde{x})}]$$

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



$$\psi(x) = \underbrace{(x_1, x_2, x_1^2 + x_2^2)}_{\psi(x)}$$

$$\psi(\tilde{x}) = \underbrace{(\tilde{x}_1, \tilde{x}_2, \tilde{x}_1^2 + \tilde{x}_2^2)}_{\psi(\tilde{x})}$$

$$B_{ij} \sim N(0, 1)$$

$$B_{ij} \sim N(0, \frac{1}{K})$$

Aside

Bochner's Thm.

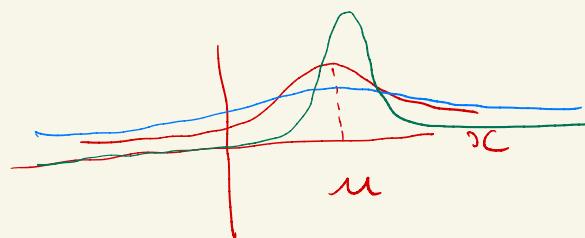
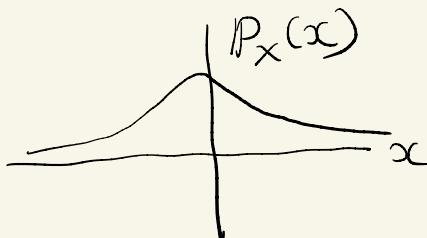
Review of Multivariate Gaussians

$$x \in \mathbb{R}; P_x(x) = \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{1}{2} x^2 \right\} \quad x \sim N(0, 1)$$

(Standard Gaussian)

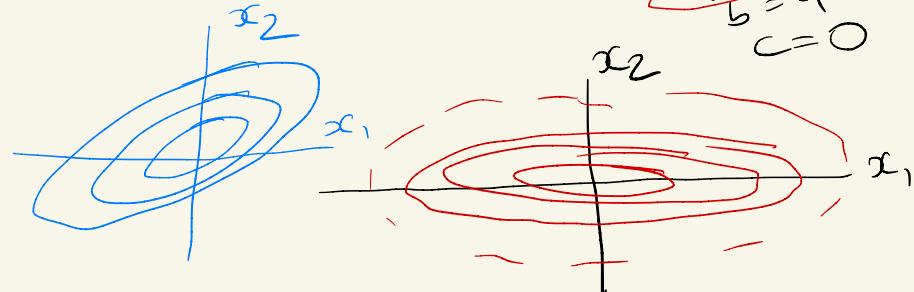
$$P_x(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2\sigma^2}(x-\mu)^2 \right\} \quad x \sim N(\mu, \sigma^2)$$

mean variance



$\sigma \gg 1$

$\sigma \ll 1$



2D - Gaussian Distribution

$$x \in \mathbb{R}^2 \quad x \sim N \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \Sigma \right)$$

Σ symmetric.

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

$$P_x(x) = \frac{1}{2\pi \underbrace{\det(\Sigma)}_{\Sigma^{-1}}} \exp \left\{ -\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu) \right\}$$

$$\Sigma = \begin{bmatrix} a & c \\ c & b \end{bmatrix} ; \det(\Sigma) = ab - c^2 ; \Sigma^{-1} = \frac{1}{ab-c^2} \begin{bmatrix} b & -c \\ -c & a \end{bmatrix}$$

$$f(x) = A \frac{1}{\sqrt{k}} \underbrace{\phi(Bx)}_{\psi(x)}$$

$$\sum(x, \tilde{x}) := \frac{1}{k} \sum_{i=1}^k \phi(B_{i,:}x) \phi(B_{i,:}\tilde{x})$$

$$\underline{\omega^T} [\omega_1, \omega_2, \dots, \omega_d]; \quad \underline{\omega_i \sim N(0, 1)}$$

$B_{ij} \sim N(0, 1)$

Law of Large Numbers \Rightarrow As $k \rightarrow \infty$:

$$\sum(x, \tilde{x}) = \mathbb{E}_{\omega} [\phi(\omega^T x) \phi(\omega^T \tilde{x})]$$

$$u := \underline{\omega^T x}; \quad v := \underline{\omega^T \tilde{x}}$$

$$u \sim N(0, \|x\|_2^2)$$

$$v \sim N(0, \|\tilde{x}\|_2^2)$$

$$uv \sim N(0, x^T \tilde{x})$$

$$u = \sum_{i=1}^d \omega_i x_i$$

$$\mathbb{E}[u] = \sum_{i=1}^k \mathbb{E}[\omega_i x_i] = 0.$$

$$\mathbb{E}[u^2] = \mathbb{E}\left[\left(\sum_{i=1}^k \omega_i x_i\right)^2\right]$$

$$= \sum_{i=1}^d \mathbb{E}[\underline{\omega_i^2 x_i^2}] + \underline{0}$$

$$= \sum_{i=1}^d \underline{x_i^2} = \underline{\|x\|_2^2}$$

$$\sum(x, \tilde{x}) = \mathbb{E}_{(u, v)} [\phi(u) \phi(v)]$$

$$(u, v) \sim N(0, \begin{bmatrix} \|x\|_2^2 & x^T \tilde{x} \\ x^T \tilde{x} & \|\tilde{x}\|_2^2 \end{bmatrix})$$

$$\Sigma(x, \tilde{x}) = \mathbb{E}_{(u, v)} \left[\frac{\phi(u)}{\|x\|_2} \frac{\phi(v)}{\|x^T \tilde{x}\|_2} \right]$$

$$(u, v) \sim N(0, \begin{bmatrix} \|x\|_2^2 & x^T \tilde{x} \\ x^T \tilde{x} & \|x^T \tilde{x}\|_2^2 \end{bmatrix})$$

Neural Net Gaussian Process

NNGP

$$\|x\|_2 = \|\tilde{x}\|_2 = 1.$$

$$\Sigma(x, \tilde{x}) = \check{\phi}(x^T \tilde{x}) \quad \text{"phi-hat": dual activation}$$

$$\phi \rightarrow \check{\phi}$$

$$\phi(z) = z \Rightarrow \check{\phi}(\xi) = \xi$$

$$\phi(z) = h_i(z) \Rightarrow \check{\phi}(\xi) = \xi^i$$

Hermite
polynomials

$$\phi(z) = \max(z, 0) \Rightarrow \check{\phi}(\xi) = \underbrace{\frac{1}{\pi} (\xi(\pi - \cos^{-1}(\xi)) + \sqrt{1-\xi^2})}_{\text{arc cosine kernel.}}$$

Module 2: Connecting NNs and Kernel Machines

Kernel Machines

$$\text{Kernel } K : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$$

Training:

$$(1) \quad K_{\text{train}} = K(x, x) \quad \left. \begin{array}{l} x \mapsto \psi(x) \\ \langle \psi(x), \psi(z) \rangle \end{array} \right\}$$

$$(2) \quad \alpha = \gamma K_{\text{train}}^{-1} \quad \left. \begin{array}{l} " \\ K(x, z) \end{array} \right.$$

Neural Networks

$$f_w(x) = A \underbrace{\Phi}_{\sqrt{K}} \underbrace{(Bx)}_{\mathbb{R}^d}$$

$$A \in \mathbb{R}^{1 \times K}$$

$$B \in \mathbb{R}^{d \times d}$$

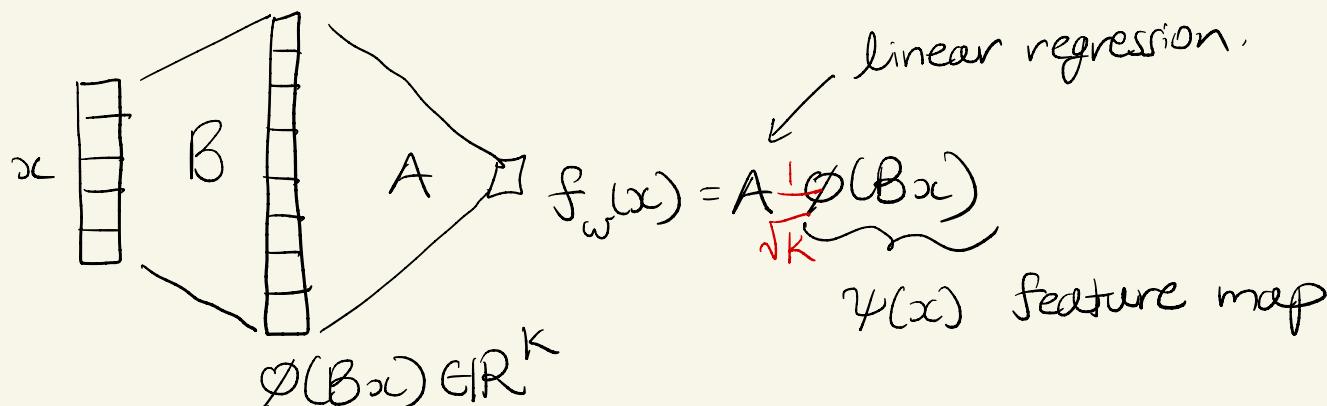
$$\Phi : \mathbb{R} \rightarrow \mathbb{R}$$

elementwise
nonlinearity

Training: Gradient Descent on

$$L(w) = \frac{1}{2} \sum_{p=1}^n (y^{(p)} - f_w(x^{(p)}))^2$$

Last time: NN GP (Neural Net Gaussian Process)



$$K(x, \tilde{x}) = \langle \psi(x), \psi(\tilde{x}) \rangle$$

$$\Sigma(x, \tilde{x}) : \underline{\underline{NNGP}}$$

$$\underline{f_{\omega}(x)} = A \frac{1}{\sqrt{K}} \phi(Bx)$$

$$\omega = \begin{bmatrix} A \\ B_{K \times d} \end{bmatrix}$$

$$B_{ij} \sim N(0, 1)$$

$$\sum(x, \tilde{x}) = \lim_{K \rightarrow \infty} \left\langle \frac{1}{\sqrt{K}} \phi(Bx), \frac{1}{\sqrt{K}} \phi(B\tilde{x}) \right\rangle$$

$$= \lim_{K \rightarrow \infty} \frac{1}{K} \sum_{i=1}^K \underbrace{\phi(B_{:,i}x)}_{\text{random variable.}} \underbrace{\phi(B_{:,i}\tilde{x})}_{\text{random variable.}}$$

$$\begin{aligned} B_{:,i} &= [B_{1,i}, \dots, B_{d,i}] \\ &\Downarrow \\ \omega &\sim N(0, I_d) \end{aligned}$$

$$= \mathbb{E}_{\omega} [\phi(\omega^T x) \phi(\omega^T \tilde{x})]$$

$$\begin{aligned} u &= \omega^T x \\ v &= \omega^T \tilde{x} \end{aligned}$$

$$= \mathbb{E}_{(u,v) \sim N(0, \Delta)} [\phi(u) \phi(v)] ; \quad \Delta = \begin{bmatrix} \|x\|_2^2 & \omega^T \tilde{x} \\ \omega^T x & \|\tilde{x}\|_2^2 \end{bmatrix}$$

$$= \frac{1}{2\pi (\det \Delta)^{1/2}} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \phi(u) \phi(v) \exp \left\{ -\frac{1}{2} [u, v] \Delta^{-1} \begin{bmatrix} u \\ v \end{bmatrix} \right\} du dv$$

$$\text{If } \|x\|_2 = \|\tilde{x}\|_2 = 1, \quad \sum(x, \tilde{x}) = \underbrace{\phi(x^T \tilde{x})}_{\text{dual activation}} = \phi(\xi)$$

$$\xi = x^T \tilde{x}$$

$$\Sigma(x, \tilde{x}) = \check{\phi}(x^T \tilde{x}) := \check{\phi}(\xi) \quad ; \quad \xi = x^T \tilde{x}.$$

Examples: $\underbrace{\phi(x) = x}_{\phi(\tilde{x}) = \tilde{x}} \Rightarrow \underbrace{\check{\phi}(\xi) = \xi}_{\leftarrow h_{\phi}(x)}$

$$f_w(x) = A \underbrace{\frac{1}{\sqrt{k}} B x}$$

$$\text{ReLU: } \phi(x) = \max(x, 0) \Rightarrow \check{\phi}(\xi) = \frac{1}{\pi} \left(\xi (\pi - \cos^{-1}(\xi)) + \sqrt{1 - \xi^2} \right)$$

$$\text{RBF: } \phi(x) = e^{ix} \Rightarrow \Sigma(x, \tilde{x}) = \exp \left\{ -L \|x - \tilde{x}\|_2^2 \right\}$$

$$\hookrightarrow \check{\phi}(\xi) = \exp \left\{ -L (2 - 2\xi) \right\}$$

Hermite Polynomials: $\phi(x) = h_i(x) \Rightarrow \check{\phi}(\xi) = \xi^i$

$$h_0(x) = 1,$$

$$h_1(x) = x$$

$$h_2(x) = \frac{x^2 - 1}{\sqrt{2}}$$

Properties : (1) $\phi(x) \Rightarrow \check{\phi}(\xi)$ $a\phi(x) \Rightarrow \underline{a^2} \check{\phi}(\xi)$ (2) $\phi(x) \Rightarrow \check{\phi}(\xi)$ $\phi'(x) \Rightarrow \check{\phi}'(\xi)$	Ex. $\phi(x) = \frac{x^2 - 1}{\sqrt{2}}$ $\check{\phi}(\xi) = \xi^2$ $\phi'(x) = \frac{2x}{\sqrt{2}}$ $\check{\phi}'(\xi) = 2\xi$
--	--

Neural Tangent Kernel (NTK)

$$f_{\underline{\omega}}(x) = A \frac{1}{\sqrt{K}} \phi(Bx) \iff \boxed{f_x(\omega)} = A \frac{1}{\sqrt{K}} \phi(Bx)$$

↑ changing
↓ fixed

Taylor Series Expansion around initial weights ω_0 .

$$f_x(\omega); \quad \omega \in \mathbb{R}, \quad \omega_0: \quad f_x(\omega) = f_x(\omega_0) + f'_x(\omega_0)(\omega - \omega_0)$$

$$\omega \in \mathbb{R}^P, \quad \omega_0: \quad f_x(\omega) = f_x(\omega_0) + \underbrace{\nabla f_x(\omega_0)^T}_{\mathbb{R}^{1 \times P}} \underbrace{(\omega - \omega_0)}_{\mathbb{R}^P}$$

$$f_x(\omega) \approx \cancel{f_x(\omega_0)} + \underbrace{\nabla f_x(\omega_0)^T}_{\mathbb{R}^{1 \times P}} (\omega - \omega_0)$$

$$L(\omega) = \frac{1}{2} \sum_{i=1}^n \left(y^{(i)} - \underbrace{\nabla f_{x^{(i)}}(\omega_0)^T}_{\mathcal{W}(x^{(i)})} \omega \right)^2 \quad \left. \right\} \text{Kernel Regression}$$

$$\gamma: x \longmapsto \nabla f_x(\omega_0)$$

$$\langle \gamma(x), \gamma(\tilde{x}) \rangle = \boxed{\langle \nabla f_x(\omega_0), \nabla f_{\tilde{x}}(\omega_0) \rangle} \xrightarrow{\text{NTK}}$$

Example: $f_x(\omega) = A \frac{1}{\sqrt{K}} \phi(Bx) = \frac{1}{\sqrt{K}} \sum_{i=1}^K A_i \phi(B_{i,:}x)$

$A \in \mathbb{R}^{1 \times K}, B \in \mathbb{R}^{K \times d}, \phi: \mathbb{R} \rightarrow \mathbb{R}$.

$$\nabla f_x(\omega_0) = \begin{bmatrix} \frac{\partial f}{\partial A_1} \\ \vdots \\ \frac{\partial f}{\partial A_K} \\ \frac{\partial f}{\partial B_{11}} \\ \vdots \\ \frac{\partial f}{\partial B_{Kd}} \end{bmatrix} \Big|_{\omega=\omega_0}$$

$$\frac{\partial f}{\partial A_i} = \frac{1}{\sqrt{K}} \phi(B_{i,:}x)$$

$$\frac{\partial f}{\partial B_{ij}} = \frac{1}{\sqrt{K}} A_i \phi'(B_{i,:}x) x_j$$

NTK

$$K(x, \tilde{x}) = \langle \nabla f_x(\omega_0), \nabla f_{\tilde{x}}(\omega_0) \rangle = \sum_{i=1}^K \frac{\partial f_x}{\partial A_i} \frac{\partial f_{\tilde{x}}}{\partial A_i} + \sum_{i=1}^K \sum_{j=1}^d \frac{\partial f_x}{\partial B_{ij}} \frac{\partial f_{\tilde{x}}}{\partial B_{ij}}$$

NTK

$$K(x, \tilde{x}) = \langle \nabla f_x(\omega_0), \nabla f_{\tilde{x}}(\omega_0) \rangle =$$

$$\frac{\partial f}{\partial A_i} = \frac{1}{\sqrt{K}} \phi'(B_i; x)$$

$$+ \sum_{i=1}^K \frac{\partial f_x}{\partial A_i} \frac{\partial f_{\tilde{x}}}{\partial A_i} + \sum_{i=1}^K \sum_{j=1}^d \frac{\partial f_x}{\partial B_{ij}} \frac{\partial f_{\tilde{x}}}{\partial B_{ij}}$$

$$\frac{\partial f}{\partial B_{ij}} = \frac{1}{\sqrt{K}} A_i \phi'(B_i; x) x_j \rightarrow \phi'(\xi)$$

$$K(x, \tilde{x}) = \frac{1}{K} \sum_{i=1}^K \phi'(B_i; x) \phi'(B_i; \tilde{x}) \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{NNGP.}$$
$$+ \frac{1}{K} \sum_{i=1}^K \sum_{j=1}^d \boxed{A_i^2 \phi'(B_i; x) \phi'(B_i; \tilde{x})} \quad \boxed{x_j \tilde{x}_j}$$

$$\text{Red: } \frac{1}{K} \sum_{i=1}^K \overline{A_i^2} \overline{\phi'(B_i; x)} \overline{\phi'(B_i; \tilde{x})} \quad \frac{\overline{\phi'(\xi)}}{x^\top \tilde{x}}$$

$$\boxed{\sum_{j=1}^d x_j \tilde{x}_j}$$

$$K(x, \tilde{x}) = \frac{1}{K} \sum_{i=1}^K \phi(B_i; x) \phi(B_i; \tilde{x})$$

$$+ \frac{1}{K} \sum_{i=1}^K \sum_{j=1}^K A_i^2 \phi'(B_i; x) \phi'(B_i; \tilde{x})$$

$x_j \quad \tilde{x}_j$

$$\lim_{K \rightarrow \infty} K(x, \tilde{x}) = \check{\phi}(\xi) + \lim_{K \rightarrow \infty} \left(\frac{1}{K} \sum_{i=1}^K A_i^2 \right) \check{\phi}'(\xi) \quad x^T \tilde{x}$$

$A_i \sim N(0, 1)$
 $B_{ij} \sim N(0, 1)$

$$= \boxed{\check{\phi}(\xi)} + \underbrace{\check{\phi}'(\xi) \xi}_{\text{Connection}}$$

$\mathbb{E}[A_i^2] = 1$
 $A \sim N(0, 1)$

Depth L net:

$$K_L(x, \tilde{x}) = \boxed{\check{\phi}^{(L)}(\xi)} + \boxed{\check{\phi}'(\check{\phi}^{(L)}(\xi)) \quad K_{L-1}(x, \tilde{x})}$$

$\check{\phi}^{(L)} = \check{\phi}(\check{\phi}(\dots \check{\phi}(\xi)))$

connection.

$L \rightarrow \infty$.

$$\xi = x^T \tilde{x}$$

Neural Tangents

$$\frac{1}{K} \sum_{i=1}^K \phi(B_i, x) \phi'(B_i, \tilde{x}) \xrightarrow{K \rightarrow \infty} \check{\phi}(\xi)$$

$$f_x(\omega) \approx f_x(\omega_0) + \underbrace{\nabla f_x(\omega_0)^T (\omega - \omega_0)}_{\omega_0 \sim N(0, I_p)}$$

$$\left. \begin{aligned} f_{x^{(i)}}(\omega_0) + \nabla f_{x^{(i)}}(\omega_0)^T \underline{\omega} &= y^{(i)} \\ \nabla f_{x^{(i)}}(\omega_0)^T \underline{\omega} &= y^{(i)} - \underline{f_{x^{(i)}}(\omega_0)} \end{aligned} \right\}$$

$$f_x(\omega) = \underbrace{A \frac{1}{\sqrt{K}} \phi(Bx)}_{f_x(\omega)}$$

Review:

Kernel Regression

$x \mapsto \psi(x)$ feature map

$$(1) K_T = K(x, x)$$

$$(2) \alpha = y K_T^{-1}$$

FC-NN

$$f_w(x) = A \frac{1}{\sqrt{K}} \phi(Bx)$$

$$A \in \mathbb{R}^{1 \times K}, B \in \mathbb{R}^{K \times d}$$

$\phi: \mathbb{R} \rightarrow \mathbb{R}$ elementwise activation

$$\text{Ex: } \phi(x) = \max(x, 0) \text{ ReLU}$$

Training: G.D. on $\mathcal{L}(w) = \frac{1}{2} \|y - f(x)\|_2^2$

NTK

$$A_i, B_{ij} \sim N(0, 1)$$

$$\text{NTK: } K(x, \tilde{x}) = \langle \nabla f_x(\omega_0), \nabla f_{\tilde{x}}(\omega_0) \rangle ; \quad \|x\|_2 = \|\tilde{x}\|_2 = 1$$

$$\lim_{K \rightarrow \infty} K(x, \tilde{x}) = \underbrace{\phi(x^T \tilde{x})}_{\text{NNGP} \leq (x, \tilde{x})} + \underbrace{\phi'(x^T \tilde{x}) x^T \tilde{x}}_{\text{correction term}}$$

ϕ : dual activation; closed form for many activations ϕ .

Motivating Problems

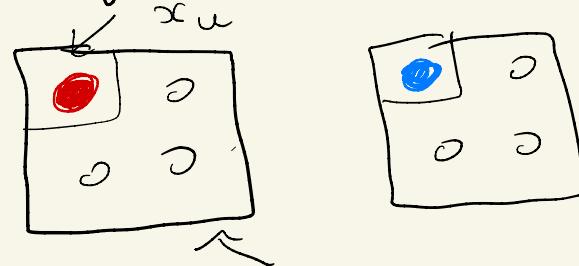
Tabular Datasets ;

Ex. 1 Housing Price Prediction

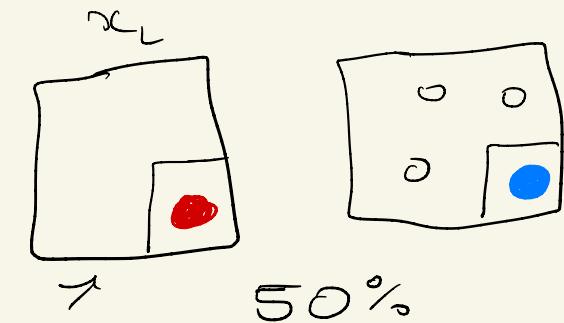
$$x \begin{bmatrix} \# \text{ beds} \\ \# \text{ baths} \\ \text{sq. ft} \end{bmatrix} \rightarrow \$ \text{ House } y$$

Ex. 2 Image Classification

Training:



Test:



$$\text{NN GP: } x^T \tilde{x} = 0.$$

$$x_u =$$

$$x_u = \begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$x_L =$$

$$x_L = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 1 \end{bmatrix}$$

$$x_u^T x_L = 0.$$

$$[1, 1, 1, 0, 0, 0] \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 1 \end{bmatrix} = 0.$$

1D - Convolutional Nets.

Training Test

Ex. 1 $x_R: \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ $\begin{bmatrix} 0 & 0 & 0 & 1 \end{bmatrix} \tilde{x}_R$ $\sum_{i=1}^d (A * x_R)_i = \frac{0+0+1+1}{4} = \frac{1}{2}$

$x_B: \begin{bmatrix} -1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ $\begin{bmatrix} 0 & 0 & 0 & -1 \end{bmatrix} \tilde{x}_B$ $\sum_{i=1}^d (A * x_B)_i = \frac{-1+0+0+0}{4} = \frac{1}{2}$

Filter: $\begin{bmatrix} 1 & 0 & 1 \end{bmatrix} : A$

$$(A * x_R) = \left[\underbrace{\begin{bmatrix} 1 & 0 & 1 \end{bmatrix}}_1 \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}, \underbrace{\begin{bmatrix} 1 & 0 & 1 \end{bmatrix}}_1 \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}, \underbrace{\begin{bmatrix} 1 & 0 & 1 \end{bmatrix}}_0 \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}, \underbrace{\begin{bmatrix} 1 & 0 & 1 \end{bmatrix}}_0 \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \right]$$

$$= [1, 1, 0, 0]$$

$$(A * x_B) = [-1, -1, 0, 0]$$

$$(A * x)_i = \sum_{\alpha=-\frac{q-1}{2}}^{\frac{q-1}{2}} A_\alpha x_{i+\alpha}$$

$$A = [A_{-1}, A_0, A_1]$$

$$q=3$$

$$(A * \underline{\underline{x}}_R)_i = \sum_{\alpha=-1}^1 A_\alpha x_{i+\alpha} = A_{-1} \cdot x_{i-1} + A_0 x_i + A_1 x_{i+1}$$

$$i=0 \Rightarrow A_{-1} \cdot 0 + A_0 \cdot 1 + A_1 \cdot 0 = A_0 = 1$$

Def: 1D conv. * ; $A \in \mathbb{R}^g$ $A = [A_{-\frac{g-1}{2}} \dots A_0 \dots A_{\frac{g-1}{2}}]$

$x \in \mathbb{R}^d \Rightarrow (A * x) \in \mathbb{R}^d$; $(A * x)_i := \sum_{\alpha=-\frac{g-1}{2}}^{\frac{g+1}{2}} A_\alpha x_{i+\alpha}$.

Ex. $A = [1 \ 1 \ 1]$; $x = [\underline{0} \ 0 \ 0 \ 1]$, $g=3$, $\alpha \in \{-1, 0, 1\}$

$$\begin{aligned}(A * x)_0 &= A_{-1} \cdot \underline{x_{-1}} + A_0 x_0 + A_1 x_1 \\ &= 1 \cdot 0 + 1 \cdot 0 + 1 \cdot 0 = 0\end{aligned}$$

$$(A * x) = [0 \ 0 \ 1 \ 1]$$

1D-CNN: (1) $f_\omega(x) = \frac{1}{d} \sum_{i=1}^d (A * x)_i$; Linear 1-filter CNN

\rightarrow (2) $f_\omega(x) = \frac{1}{d} \sum_{i=1}^d \phi(A * x)_i$; Nonlinear 1-filter CNN.

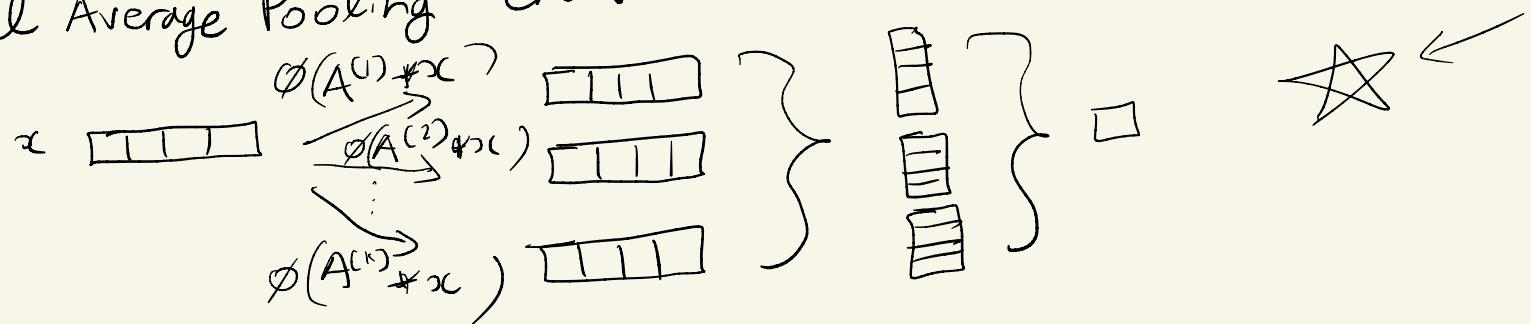
\rightarrow (3) $f_\omega(x) = \frac{1}{d} \sum_{i=1}^d \sum_{l=1}^k \phi(A^{(l)} * x)_i$; k -filter Nonlinear CNN.

(4) $f_\omega(x) = \sum_{l=1}^k \underbrace{B^{(l)}}_{\mathbb{R}^{1 \times d}} \underbrace{\phi(A^{(l)} * x)}_{\mathbb{R}^d}$; k -filter Nonlinear FC NN

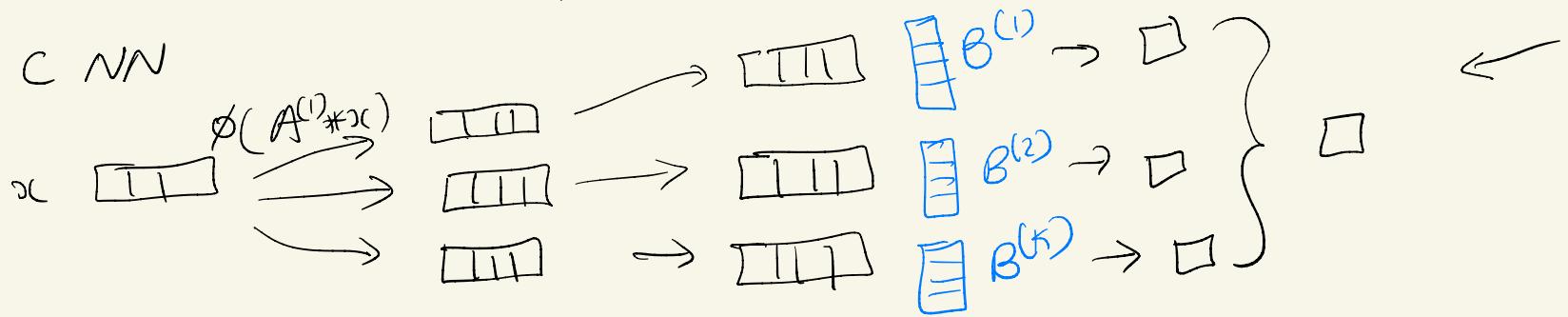
- 1D-CNN : (1) $f_{\omega}(x) = \frac{1}{d} \sum_{i=1}^d (A * x)_i$; Linear 1-filter CNN
- \rightarrow (2) $f_{\omega}(x) = \frac{1}{d} \sum_{i=1}^d \phi(A * x)_i$; Nonlinear 1-filter CNN.
- $\rightarrow \left\{ \rightarrow$ (3) $f_{\omega}(x) = \frac{1}{d} \sum_{i=1}^d \sum_{l=1}^k \phi(A^{(l)} * x)_i$; k -filter nonlinear CNN.
- $\rightarrow \left\{ \text{(4) } f_{\omega}(x) = \sum_{l=1}^k \underbrace{B^{(l)}}_{R^{1 \times d}} \underbrace{\phi(A^{(l)} * x)}_{R^d}$; k -filter Nonlinear FC NN

Diagrams:

(1) Global Average Pooling CNN : GAP-CNN.



(2) FC NN



$(A * x) \in \mathbb{R}^d$, $x \in \mathbb{R}^d$

Ex: $A = \begin{bmatrix} A_{-1} & A_0 & A_1 \end{bmatrix}$ $x \in \mathbb{R}^4$

$$\begin{bmatrix} A_0 & A_1 & 0 & 0 \\ A_{-1} & A_0 & A_1 & 0 \\ 0 & A_{-1} & A_0 & A_1 \\ 0 & 0 & A_{-1} & A_0 \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} = A * x$$

Toepiltz.

FC.

$$\begin{bmatrix} A_{11} & A_{12} & A_{13} & A_{14} \\ \vdots & \vdots & \vdots & \vdots \\ A_{41} & \ddots & A_{44} \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

CNTK:

$$\text{GIAP-CNN} = f_{\omega}(x) = \frac{1}{d} \sum_{i=1}^d \sum_{l=1}^k \frac{1}{\sqrt{K}} \phi(A^{(l)} * x);$$
$$= \boxed{\frac{1}{d} \sum_{i=1}^d \sum_{l=1}^k \frac{1}{\sqrt{K}} \phi\left(\sum_{\alpha} A_{\alpha}^{(l)} x_{i+\alpha}\right)} *$$

$$\nabla f_x(\omega_0) = \left[\frac{\partial f(\omega_0)}{\partial A_{-\frac{g-1}{2}}^{(1)}} , \dots , \frac{\partial f(\omega_0)}{\partial A_{\frac{g-1}{2}}^{(K)}} \right] \quad A_{\alpha}^{(l)} \sim N(0, 1)$$

$$\frac{\partial f(\omega_0)}{\partial A_j^{(l)}} = \frac{1}{d} \sum_{i=1}^d \frac{1}{\sqrt{K}} \phi' \left(\sum_{\alpha} A_{\alpha}^{(l)} x_{i+\alpha} \right) x_{i+j}$$

$$K(x, \tilde{x}) = \sum_{l=1}^k \sum_{j=-\frac{g-1}{2}}^{\frac{g+1}{2}} \frac{\partial f(\omega_0)}{\partial A_j^{(l)}} \frac{\partial f_{\tilde{x}}(\omega_0)}{\partial A_j^{(l)}}$$
$$= \sum_{l=1}^k \sum_{i=1}^d \frac{1}{d^2} \frac{1}{K} \sum_{i'=1}^d \phi'(A^{(l)} * x) ; x_{i+j} \cdot \phi'(A^{(l)} * \tilde{x}) ; \tilde{x}_{i'+j}$$

$$\text{As } K \rightarrow \infty = \frac{1}{d^2} \sum_{i=1}^d \sum_{i'=1}^d \phi'(x[i], \tilde{x}[i']) \langle x[i], \tilde{x}[i] \rangle \}$$
$$x[i] : [x_{i-1}, x_i, x_{i+1}]$$

$$K(x, \tilde{x}) = \sum_{\ell=1}^K \sum_{j=-\frac{q-1}{2}}^{q+1} \frac{\partial f_\ell(\omega_j)}{\partial A_j^{(\ell)}} \frac{\partial \tilde{f}_x^{(j)}(\omega_j)}{\partial A_j^{(\ell)}}$$

$$= \sum_{\ell=1}^K \sum_0^d \frac{1}{d^2} \sum_{i=1}^d \sum_{i'=1}^d \phi'(A^{(\ell)} * x); x_{i+j} \cdot \phi'(A^{(\ell)} * \tilde{x}); \tilde{x}_{i'+j}$$

$$= \frac{1}{K} \sum_{\ell=1}^K \frac{1}{d^2} \sum_{i=1}^d \sum_{i'=1}^d \phi'(A^{(\ell)} * x); \phi'(A^{(\ell)} * \tilde{x}); \sum_{j=-\frac{q-1}{2}}^{q+1} x_{i+j} \tilde{x}_{i'+j}$$

$$\underbrace{[x_{i-\frac{q-1}{2}}, \dots, x_{i+\frac{q+1}{2}}]}_{x[i]}$$

$$\underbrace{[\tilde{x}_{i-\frac{q-1}{2}}, \dots, \tilde{x}_{i+\frac{q+1}{2}}]}_{\tilde{x}[i']}$$

$$= \frac{1}{d^2} \sum_{i=1}^d \sum_{i'=1}^d \frac{1}{K} \sum_{\ell=1}^K \phi'(A^{(\ell)} * x); \phi'(A^{(\ell)} * \tilde{x}); \langle x[i], \tilde{x}[i'] \rangle$$

As $K \rightarrow \infty$. $\phi'(x[i], \tilde{x}[i'])$

$$= \frac{1}{d^2} \sum_{i=1}^d \sum_{i'=1}^d \phi'(x[i], \tilde{x}[i']) \langle x[i], \tilde{x}[i'] \rangle.$$

$$K_{\text{GAP}}(x, \tilde{x}) = \frac{1}{d^2} \sum_{i=1}^d \sum_{i'=1}^d \phi'(x[i], \tilde{x}[i']) \langle x[i], \tilde{x}[i'] \rangle.$$

$$K_{\text{FC}}(x, \tilde{x}) = \frac{1}{d} \sum_{i=1}^d \phi'(x[i], \tilde{x}[i]) \langle x[i], \tilde{x}[i] \rangle \}$$

Role of Increasing Depth in CNNs

$$f_w(x) = W_1 * (W_2 * (W_3 * x))$$

$$W_3 = \begin{bmatrix} & & & \\ & \ddots & & \\ & & \ddots & 0 \\ & & & \ddots \\ 0 & & & & \ddots \\ & & & & & \ddots \end{bmatrix}$$

$$\underbrace{W_1 \ W_2 \ W_3}_{=} = \underbrace{\begin{bmatrix} & & & \\ & \ddots & & \\ & & \ddots & 0 \\ & & & \ddots \\ 0 & & & & \ddots \\ & & & & & \ddots \end{bmatrix}}$$

Review :

(1) Linear Regression: $\omega X = y$ $\| \omega X - y \|_2^2$ $\omega = y X^+$

(2) Kernel Regression: $x \mapsto \psi(x)$ (1) $K_T = K(X, X)$
 $\mathbb{R}^d \quad \mathbb{R}^P$ $K(x, z) = \langle \psi(x), \psi(z) \rangle$

(3) Neural Nets : $f_\omega(x) = A \phi(Bx)$ $\| y - f_\omega(x) \|_2^2$ using G.D.

(4) NNGP & NTK connected in wide neural nets
with kernel methods.

NTK \Rightarrow $K(x, z) = \langle \nabla f_x(\omega_0), \nabla f_z(\omega_0) \rangle$
 $\| x \|_2 = \| z \|_2 = 1 \Rightarrow K(x, z) = \phi(\xi) + \xi \phi'(\xi)$
where $\xi = x^T z$.

Course challenge:

NTK : .2, .3

Test R²
.98

Laplace : <.3

$$f^*(x) = \underbrace{x_1 x_2}_{\dots}$$

$$M = W_1^\top W_1 = \begin{bmatrix} 1 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 100 \end{bmatrix} \quad \left. \right\}$$

Expected Gradient Outer Product

$$\nabla f^*(x) = \begin{bmatrix} \frac{\partial f^*}{\partial x_1} \\ \frac{\partial f^*}{\partial x_2} \\ \vdots \\ 0 \end{bmatrix} = \begin{bmatrix} x_2 \\ x_1 \\ \vdots \\ 0 \end{bmatrix}$$

$$\nabla f^*(x) \nabla f^*(x)^\top = \begin{bmatrix} x_2 \\ x_1 \\ \vdots \\ 0 \end{bmatrix} \begin{bmatrix} x_2 & x_1 & 0 & \dots & 0 \end{bmatrix} = \begin{bmatrix} x_2^2 & x_2 x_1 & 0 & \dots & 0 \\ x_2 x_1 & x_1^2 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & \dots & 0 & 0 \end{bmatrix}$$

$$x_i \sim \mathcal{N}(0, 1)$$

$$G(f^*) = \mathbb{E}_x \left[\nabla f^*(x) \nabla f^*(x)^\top \right] = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$f^* \approx \hat{f}; G(f^*) \approx G(\hat{f})$

Goal: $f(x) = A \frac{1}{\sqrt{K}} Bx$; $A \in \mathbb{R}^{1 \times K}$, $B \in \mathbb{R}^{K \times d}$

$$B^{(1)\top} B^{(1)} = \nabla f^{(1)}(x) \nabla f^{(1)}(x)^\top$$

$$A^{(0)} = A^{(0)} + n \left(y - A^{(0)} \frac{1}{\sqrt{K}} B^{(0)} x \right) x^\top B^{(0)\top} \frac{1}{\sqrt{K}} \quad (1)$$

$$B^{(0)} = \underline{B^{(0)}} + n A^{(0)\top} \left(y - A^{(0)} \frac{1}{\sqrt{K}} B^{(0)} x \right) x^\top \frac{1}{\sqrt{K}} \quad (2)$$

$$B^{(0)} = 0, \quad A^{(0)} \sim N(0, I)$$

$$(1) \quad A^{(1)} = A^{(0)}$$

$$(2) \quad B^{(1)} = n \underbrace{A^{(0)\top} y}_{\frac{1}{\sqrt{K}}} x^\top \frac{1}{\sqrt{K}}$$

$$B^{(0)\top} B^{(0)} = n x y^\top A^{(0)} \frac{1}{\sqrt{K}} \cdot \frac{1}{\sqrt{K}} n A^{(0)\top} y x^\top$$

$$= n^2 x y^\top \underbrace{\frac{1}{K} A^{(0)} A^{(0)\top}}_{\sum_{i=1}^K A^{(0)i}^2} y x^\top$$

$$\lim_{K \rightarrow \infty} B^{(0)\top} B^{(0)} = n^2 x y^\top y x^\top$$

$$\textcircled{1} \quad A^{(1)} = A^{(0)}$$

$$\textcircled{2} \quad B^{(1)} = n \underbrace{A^{(0)T} y}_{\frac{1}{\sqrt{k}}} x^T$$



$$\lim_{K \rightarrow \infty} B^{(1)T} B^{(1)} = n^2 \times y^T y \times x^T$$

$$f^{(1)}(x) = A^{(1)} \frac{1}{\sqrt{k}} B^{(1)} x ; \quad \nabla f^{(1)}(x) = B^{(1)T} \frac{1}{\sqrt{k}} A^{(1)T}$$

$$\nabla f^{(1)}(x) \nabla f^{(1)}(x)^T = B^{(1)T} \frac{1}{\sqrt{k}} A^{(1)T} A^{(1)} \frac{1}{\sqrt{k}} B^{(1)}$$

$$= \underbrace{n \times y^T A^{(0)} \frac{1}{\sqrt{k}}}_{B^{(1)T}} \cdot \frac{1}{\sqrt{k}} A^{(0)T} \cdot A^{(0)} \frac{1}{\sqrt{k}} \underbrace{n A^{(0)T} y x^T}_{B^{(1)}}$$

$$= n^2 \times y^T \left(\frac{1}{k} A^{(0)T} A^{(0)} \right) \left(\frac{1}{k} A^{(0)T} A^{(0)} \right) y x^T$$

$\underbrace{\qquad}_{\substack{\rightarrow 1 \\ \text{as } k \rightarrow \infty}}$ $\underbrace{\qquad}_{\substack{\rightarrow 1 \\ \text{as } k \rightarrow \infty}}$

$$\lim_{K \rightarrow \infty} \nabla f^{(1)}(x) \nabla f^{(1)}(x)^T = n^2 \times y^T y \times x^T$$

After 1-step : $\underline{f^{(1)}(x)}$ ←

$$\underline{\underline{B^{(1)T} B^{(1)}}} = \nabla f^{(1)}(x) \nabla f^{(1)}(x)^T$$

Recursive Feature Machines (RFMs)

$$(1) \quad \hat{f}(M_t^{\frac{1}{2}} x) \approx f^* \leftarrow \text{kernels for } \hat{f}, M^{(0)} = I.$$

$$(2) \quad M_t = \mathbb{E}_x \left[\nabla \hat{f}(x) \nabla \hat{f}(x)^T \right]$$

Euclidean dist.

$$\text{Algorithm: } K(x, z) = \exp \left\{ -L \underbrace{\|x - z\|_2}_{\text{Euclidean dist.}} \right\}$$

$$K_M(x, z) = K(M^{\frac{1}{2}} x, M^{\frac{1}{2}} z)$$
$$= \exp \left\{ -L \|M^{\frac{1}{2}} x - M^{\frac{1}{2}} z\|_2 \right\}$$

$$= \exp \left\{ -L \sqrt{\underbrace{(x - z)^T M (x - z)}_{\text{Mahalanobis distance}}} \right\}$$

Mahalanobis distance.

$$\text{Algorithm : } k_M(x, z) = \exp \left\{ -L \sqrt{(x-z)^T M (x-z)} \right\}$$

(1) Solve kernel Regression with Km7

$$K_{\text{train}} = K_m(x, x)$$

$$\alpha = y^k_{\text{train}}^{-1}$$

$$\hat{f}(x) = \alpha K_m(x, x)$$

$$\rightarrow (2) M = \frac{1}{n} \sum_{x^{(p)} \in X} \left\{ \hat{f}(x^{(p)}) \circ \hat{f}(x^{(p)})^T \right\} G(\hat{f})$$

(3) Iterate steps 1 & 2 -

$$x \in \mathbb{R}^{100} \quad f^*(x) = x, x_2 \quad \}$$

$$K(x, z) = \exp \left\{ -\frac{\|x - z\|_2^2}{2} \right\}$$

$$f^*(x) = g(u^T x)$$

$$\nabla f^*(x) = g'(u^T x) u$$

$$\begin{aligned} \mathbb{E}_x [\nabla f^*(x) \nabla f^*(x)^T] &= \mathbb{E}_x [g'(u^T x)^2 u u^T] \\ &= \underbrace{\mathbb{E}_x [g'(u^T x)^2]}_C u u^T \end{aligned}$$

$$u \in \mathbb{R}^{50}$$

$$g(u^T x) \rightarrow \mathcal{O}_{\underline{u u^T}}$$

$$= \underbrace{C}_{\mathbb{E}} \boxed{u u^T} \leftarrow$$

$$\Sigma = \begin{bmatrix} 1 & \xi \\ \bar{\xi} & 1 \end{bmatrix}$$

$$\exp \left\{ -\frac{1}{2} \frac{[u \quad v]}{(1-\xi^2)} \begin{bmatrix} 1 & -\xi \\ -\bar{\xi} & 1 \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} \right\}$$

$$\begin{bmatrix} u - v\xi & -u\xi + v \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix}$$

$$u^2 - uv\xi - uv\xi + v^2 = u^2 + v^2 - 2uv\xi$$