

DSAS105 LECTS

Recap (Hypothesis Space \mathcal{H})

① Linear Basic Function

$$f(x) = \sum_j a_j \phi_j(x) \rightarrow \text{fixed before training}$$

② DT

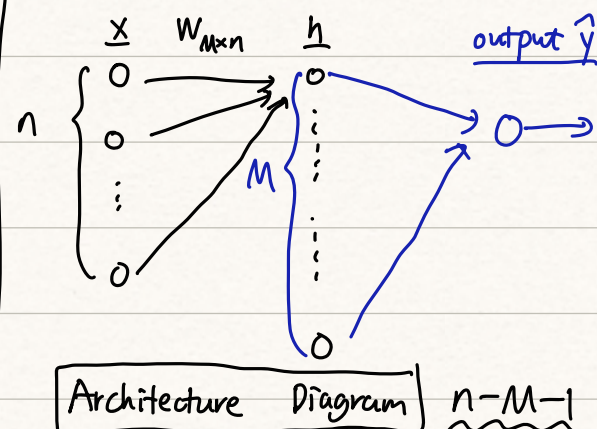
$$f(x) = \sum_j a_j \mathbb{1}_{R_j}(x) \rightarrow \text{learning parameters}$$

Neural Network

$$f(x) = \sum_{j=1}^M v_j f(w_j^T x + b_j)$$

\hat{j} -th hidden unit

only one hidden layer \rightarrow shallow



Activation function

① sigmoid $\in (0, 1)$

$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$

★ ② tanh $\in (-1, 1)$

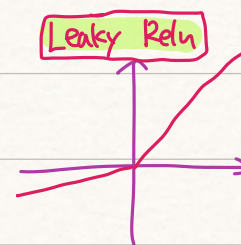
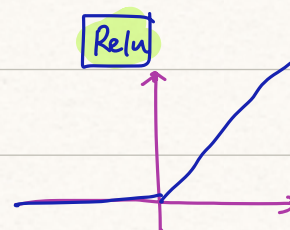
$$\tanh(z)$$

★ ③ ReLU $\in [0, \infty)$

$$\text{Relu}(z) = \max(0, z)$$

④ Leaky Relu

$$\text{Leaky Relu}(z) = \begin{cases} z & z \geq 0 \\ \delta z & z < 0 \end{cases}$$



Universal Approximation Thm

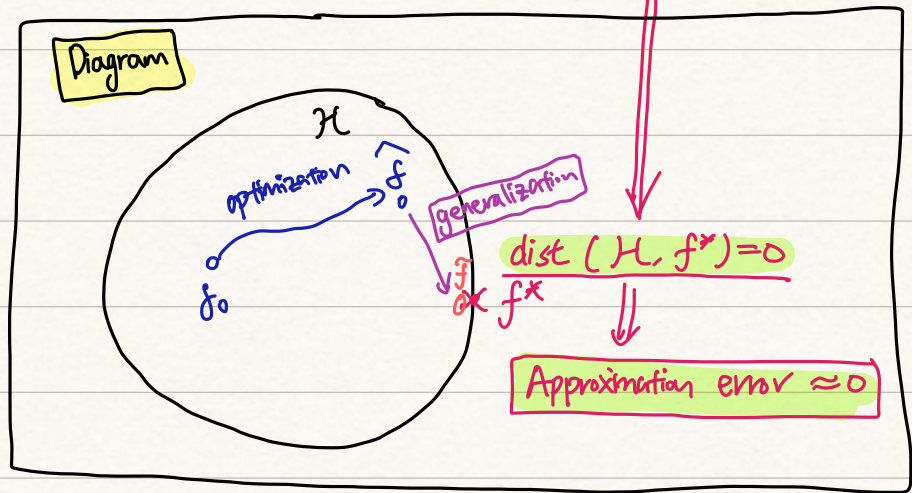
(guarantee)

\Rightarrow NN is an universal approximator

NN can be viewed as a function approximator (any continuous function)

★★★

the hypothesis space for NN is LARGE



Wonderful Part

Question: Will NN suffer from Curse Of Dimensionality?

linear & nonlinear feature map

E.g. $u = (3, 1, 2) = 3e_1 + e_2 + 2e_3$

Basis are fixed



Linear Basis

① what is the optimal approximator of $\hat{u} = a_1 e_1 + a_2 e_2$ (of the form)

→ formulate $\hat{u} = \min \| \hat{u} - u \|_2^2$

$$\Rightarrow \hat{u} = 3e_1 + e_2$$

$$\& \text{Error}(\hat{u}, u) = \| \hat{u} - u \|_2 = 2$$

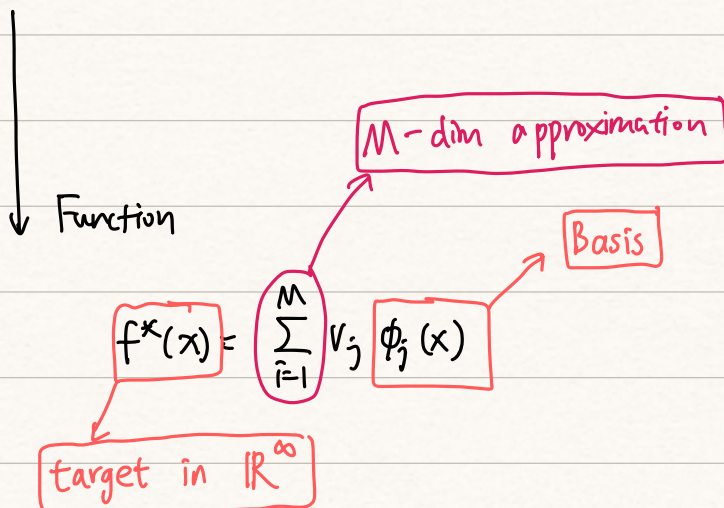
② what is the best 2 basis? → we can arbitrarily choose

2 basis

Neural Network

Answer: $(e_1, e_3) \Rightarrow \text{Error} = 1$

E.g. Assume in \mathbb{R}^{∞} you have a vector v of dim-3



Recap

Will MN suffer from Curse of Dimensionality?

Answer: NO



[Baron, 1993]

the reason is intuitively from previous example.

$$f^*: [0,1]^d \rightarrow \mathbb{R}$$

we can choose M -basis (hidden neurons), s.t

$$\|f^* - f_M\|^2 \leq O(M^{-1})$$

→ Bound is not related to $d!!!$



overcome the curse of dim!

we have the freedom to choose
feature map (Basis)

[Comparison] for linear basis function

$$\|f^* - f_M\|^2 \leq O(M^{-\frac{\alpha}{d}}) \rightarrow \text{Curse of Dimensionality!}$$

[Optimization]

Oracle

→ Universal Approximation Thm \implies there is a good approximator of f^* in \mathcal{H}

① Hypothesis Space $\mathcal{H} = \{f: f(x) = f_\theta(x); \theta \in \Theta\}$

② ERM Framework

$$\hat{\theta} = \underset{\theta \in \Theta}{\operatorname{argmin}} \Phi(\theta) \rightarrow \text{empirical loss}$$

unconstrained
optimization

$$= \underset{\theta \in \Theta}{\operatorname{argmin}} \frac{1}{N} \sum_{n=1}^N L(f_\theta(x_n), y_n)$$

(too much non-linearity)

③ Gradient Descent (SGD)

a) Necessary Condition :

$$\nabla_{\theta} \Phi(\theta) \big|_{\theta=\hat{\theta}} = 0$$

Non-linear equation

No-closed form

(not solvable)

b) Iterative Method

1) Gradient Descent (GD)

update: $\theta_{k+1} = \theta_k - \alpha_k \nabla_{\theta} \Phi(\theta) \big|_{\theta=\theta_k}$ ($\alpha_k > 0$)

(to achieve $\min_{\theta} \Phi(\theta)$)

Approximation Model

Taylor Approximation Order-1

$$\Phi(\theta + \Delta\theta) \approx \Phi(\theta) + \nabla_{\theta} \Phi(\theta)^T \Delta\theta + o(\|\Delta\theta\|^2)$$

a choice of $\Delta\theta_{GD} = -\nabla_{\theta} \Phi(\theta)$

can only guarantee

$$\|\nabla_{\theta} \Phi(\theta) \big|_{\theta=\theta_k}\| \rightarrow 0$$

stationary point / local minima

→ depend on the landscape of $\Phi(\cdot)$

we cannot have

$$\theta_k \rightarrow \theta^* \text{ as } k \rightarrow \infty$$

global minima

(Special Case)

→ for CONVEX f^* , global minima \iff local minima

$$\iff \nabla_{\theta} \Phi(\hat{\theta}) = 0$$

(or sub-gradient $0 \in \partial \Phi(\hat{\theta})$)

rigorously speaking

★

GD can guarantee converge to global minima

[E.g.]

① linear basis model

$$y = \Phi w$$

convex

$$\begin{aligned} \rightarrow \Phi(w) &= (y - \Phi w)^T (y - \Phi w) \\ &= w^T \Phi^T \Phi w - 2y^T \Phi w + \text{const} \end{aligned}$$

$$\text{since } \nabla \Phi(w) = 2\Phi^T \Phi w - 2\Phi^T y$$

$$\nabla^2 \Phi(w) = 2\Phi^T \Phi \rightarrow (\text{PSD})$$

implies

$\Phi(w)$ is convex

② SVM \rightarrow convex

③ NN \rightarrow highly non-convex

2) SGD \rightarrow works very well \rightarrow we may lose the advantages of convergence

$$\Phi(\theta) = \frac{1}{N} \sum_{n=1}^N \text{Loss}(f_{\theta}(x_n), y_n)$$

of convergence

But we can be close to the place we want

$$\nabla_{\theta} \Phi(\theta) = \frac{1}{N} \sum_{n=1}^N \nabla_{\theta} \text{Loss}(f_{\theta}(x_n), y_n) \quad (*)$$

① when N is large, $(*)$ is computationally expensive

② use $\frac{1}{|B|} \sum_{n \in B} \nabla_{\theta} \text{Loss}(f_{\theta}(x_n), y_n)$ sample B randomly

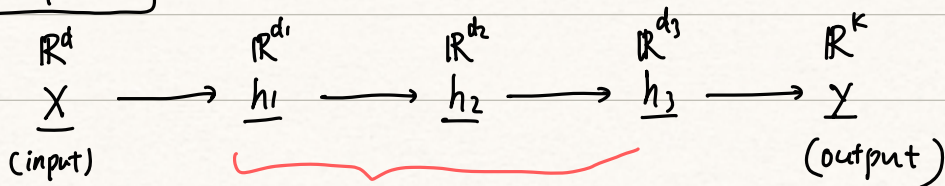
to treat as one APPROXIMATOR of $\nabla_{\theta} \Phi(\theta)$

$|B| \rightarrow$ Batch Size
 \Downarrow
Small

a noise approximator

★ help to go out of local minimum

Deep NN \rightarrow more hidden layers



hidden layers (3)

How to calculate gradient?

→ through BP algorithm (chain rule)