

Recap

Shallow NN \rightarrow adaptive basis model
linear basis model

\rightarrow universal approximator

advantage: do not suffer from Curse of Dimensionality
 dis-advantage: difficult to train

Problem

gradient calculation \Rightarrow expensive
 how to calculate gradient in practice

DNN \rightarrow FCNN (simplest one)
 extension of GD-method for large scale problems {SGD, Momentum}
 BP Algo \Leftrightarrow Computational Graph

① Shallow NN \rightarrow Deep NN (DNN)

Regression

$$\begin{aligned} h^{(0)} &= x \\ h^{(i)} &= g^{(i)}(W^{(i)}h^{(i-1)} + b^{(i)}) \\ h^{(i+1)} &= g^{(i+1)}(W^{(i+1)}h^{(i)} + b^{(i+1)}) \\ \hat{y} &= w^T h^{(l)} + c \end{aligned}$$

$$i = 0, 1, \dots, l-1$$

Q: How many layers is appropriate?

A: larger # of layers may not result in improvements!

\downarrow

Trade-off { model capacity
 difficulty in training process

② Advantages for FCNN Architecture

1. sequential model can be viewed as one PRIOR

↕
 ✱ built-in property for FCNN

↖
inductive bias

2. { better approximation property
 better generalization property

✱ ③ Issues

1. Optimization related ↔ non-convexity

2. Prone to over-fitting with inappropriate { initialization
 regularization

④ GD → SGD

GD :

$$\hat{\theta} = \min_{\theta} R_{\text{emp}}(\theta)$$

Mini-batch SGD

(combine)

SGD

$$\Rightarrow \theta^{(k+1)} = \theta^{(k)} - \epsilon \nabla_{\theta} R_{\text{emp}}(\theta^{(k)})$$

→ what we use in real application

unbiased estimator with large variance

$$R_{\gamma_k}(\theta) \rightarrow \gamma_k \sim \text{Unif}(1, 2, \dots, K)$$

✱ $\nabla_{\theta} R_{\gamma_k}(\theta^{(k)})$ → estimator of $\nabla_{\theta} R_{\text{emp}}(\theta^{(k)})$

$$\Rightarrow \theta^{(k+1)} = \theta^{(k)} - \epsilon \nabla_{\theta} R_{\gamma_k}(\theta^{(k)})$$

estimator of $\nabla_{\theta} R_{\text{pop}}(\theta^{(k)})$

$$\begin{cases} R_{\text{emp}}(f) = \frac{1}{N} \sum_{i=1}^N L(f(x_i), y_i) \\ R_{\text{pop}}(f) = \mathbb{E}_{(x,y) \sim \mathcal{D}} [L(f(x), y)] \end{cases}$$

★

Dynamics of SGD [Example] \leftrightarrow hand analysis

$$\Rightarrow R(\theta) = \frac{1}{N} \sum_{i=1}^N \frac{1}{2} (\theta - \theta^{(i)})^2$$

$$\{\theta^{(i)}\}_{i=1}^N \text{ satisfies that } \begin{cases} \frac{1}{N} \sum_{i=1}^N \theta^{(i)} = 0 \\ \frac{1}{N} \sum_{i=1}^N (\theta^{(i)})^2 = 1 \end{cases}$$

Data Point

then $R(\theta) = \frac{1}{2} \theta^2 + \frac{1}{2}$

GD iterates : $\theta_R^{\text{GD}} = \theta_{R-1}^{\text{GD}} - \varepsilon \theta_{R-1}^{\text{GD}}$
 $= (1-\varepsilon)^k \theta_0$

SGD iterates : $\theta_R^{\text{SGD}} = \theta_{R-1}^{\text{SGD}} - \varepsilon (\theta_{R-1}^{\text{SGD}} - \theta^{(\gamma_{R-1})})$

$$= (1-\varepsilon) \theta_{R-1}^{\text{SGD}} + \varepsilon \theta^{(\gamma_{R-1})}$$

$$\gamma_{R-1} \sim \text{Unif}([1, 2, \dots, N])$$



Convergence? $\Rightarrow \theta_R^{\text{SGD}} = (1-\varepsilon)^2 \theta_{R-2}^{\text{SGD}} + (1-\varepsilon) \cdot \varepsilon \theta^{(\gamma_{R-2})} + \varepsilon \theta^{(\gamma_{R-1})}$

= ...

$$= \underbrace{(1-\varepsilon)^k \theta_0}_{\text{deterministic}} + \varepsilon \underbrace{\sum_{j=1}^k (1-\varepsilon)^{j-1} \theta^{(\gamma_{R-j})}}_{\text{randomness}}$$

deterministic

randomness

Note that

$$\mathbb{E}_{\gamma_k} [\theta^{(\gamma_k)}] = \frac{1}{N} \sum_{i=1}^N \theta^{(i)} = 0$$

⇒ Property :

1. SGD iteration $\xleftrightarrow{\text{equivalent}}$ GD iteration w.r.t expectation

$$\mathbb{E}_{\gamma} [\theta_k^{\text{SGD}}] = (1-\epsilon)^k \theta_0 = \theta_k^{\text{GD}}$$

2. What about variance? (Second - Moment)

$$\begin{aligned} \rightarrow \mathbb{E}_{\gamma} [(\theta_k^{\text{SGD}})^2] &= (1-\epsilon)^{2k} \theta_0^2 + 2\epsilon(1-\epsilon)^k \theta_0 \sum_{j=1}^k (1-\epsilon)^{j-1} \mathbb{E}_{\gamma_j} [\theta^{(2k-j)}] \\ &\quad + \epsilon^2 \sum_{j,l=1}^k (1-\epsilon)^{j+l-2} \mathbb{E}_{\gamma_{kj}, \gamma_{kl}} [\theta^{(2k-j)} \theta^{(2k-l)}] \end{aligned}$$

(if $\gamma_1, \dots, \gamma_k$ i.i.d.) $= (1-\epsilon)^{2k} \theta_0^2$

$$+ \epsilon^2 \sum_{j=1}^k (1-\epsilon)^{2j-2}$$

$$\begin{aligned} &\mathbb{E}_{\gamma_i, \gamma_j} [\theta^{(2i)} \theta^{(2j)}] \\ &= \begin{cases} 0, & i \neq j \\ 1, & i = j \end{cases} \end{aligned}$$

$$= (1-\epsilon)^{2k} \theta_0^2 + \frac{\epsilon}{2-\epsilon} [1 - (1-\epsilon)^{2k}]$$

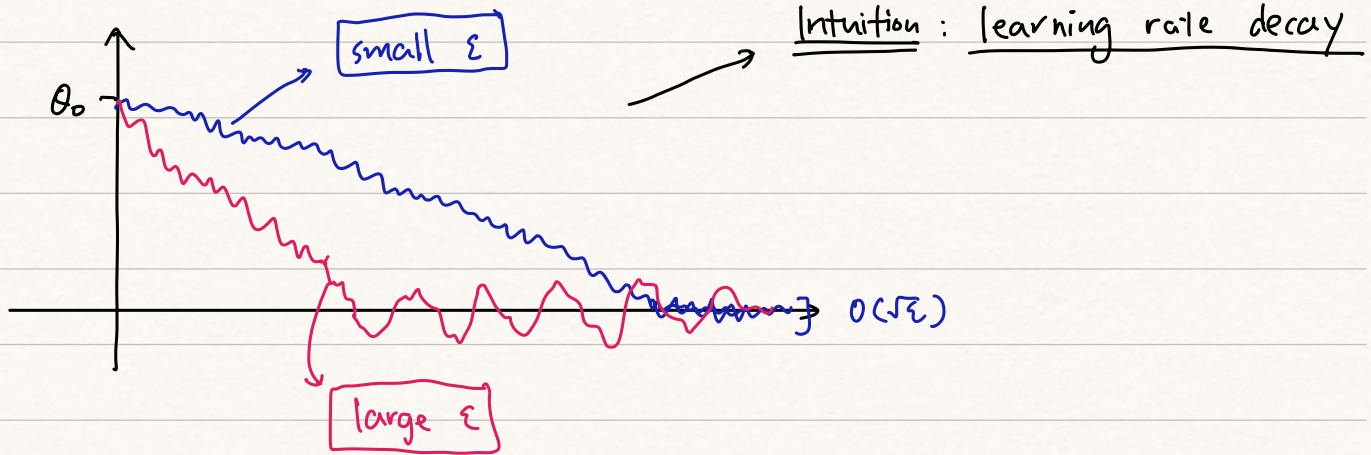
$$\rightarrow \text{Var}_{\gamma} [\theta_k^{\text{SGD}}] = \mathbb{E}_{\gamma} [(\theta_k^{\text{SGD}})^2] - (\mathbb{E}_{\gamma} [\theta_k^{\text{SGD}}])^2$$

$$= \frac{\epsilon}{2-\epsilon} [1 - (1-\epsilon)^{2k}] \xrightarrow{k \rightarrow \infty} \frac{\epsilon}{2-\epsilon} > 0$$

→ when $\begin{cases} \epsilon \rightarrow 0 \\ k \rightarrow \infty \end{cases}$, then $\text{Var}_{\gamma} [\theta_k^{\text{SGD}}] \approx \frac{\epsilon}{2}$

std $\rightarrow O(\sqrt{\epsilon})$ as $k \rightarrow \infty$

Diagram



⑤ SGD with varying learning rate

↓
Learning Rate Decay

Condition

$$\begin{cases} \sum \epsilon_k = \infty \\ \sum \epsilon_k^2 < \infty \end{cases}$$

$$\rightarrow \epsilon_k = \frac{1}{k+1}$$

↓
Sufficient Condition for convergence

⑥ Variants of SGD

↓
Momentum method ↔ reduce zig-zag behaviour

GD with momentum

$$\begin{cases} v_{k+1} = \alpha v_k - \epsilon \nabla_{\theta} \text{Loss}(\theta_k) \\ \theta_{k+1} = \theta_k + v_{k+1} \end{cases}$$

Note: if $\alpha=0$, then we go back to conventional GD

⑦ Back Propagation Algo (BP)



Question: How to calculate $\nabla_{\theta} \text{Loss}(\theta)$

$$= \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} L(f_{\theta}(x_i), y_i)$$

TensorFlow Architecture

Computational Graph \longleftrightarrow formulate the network architecture

↓
DAG Directed Acyclic Graph

↓
visualize the parameter dependency

Chain Rule

1. scalar function

$$h = g \circ f$$

$$\frac{dh}{dx} = \frac{dg}{df} \cdot \frac{df}{dx}$$

2. vector function

$$h = g \circ f$$

$$h: \mathbb{R}^m \rightarrow \mathbb{R}$$

$$f: \mathbb{R}^m \rightarrow \mathbb{R}^n \quad g: \mathbb{R}^n \rightarrow \mathbb{R}$$

$$\nabla_x h(x) = \nabla_x f(x) \nabla_f g(f)$$

$$\begin{cases} \nabla_x h(x) \in \mathbb{R}^m \\ \nabla_x f(x) \in \mathbb{R}^{m \times n} \quad \nabla_f g(f) \in \mathbb{R}^n \end{cases}$$

$$\frac{dh}{dx} = \frac{dg}{df} \cdot \frac{df}{dx} \longrightarrow \text{Jacobian} \quad \frac{dh}{dx} = \nabla_x h(x)^T$$

Note: Jacobian = Gradient^T

Toy Example :

$$\bullet h^{(1)} = \omega^{(1)} x$$

$$\bullet h^{(2)} = \omega^{(2)} h^{(1)}$$

$$\bullet h^{(3)} = \omega^{(3)} h^{(2)}$$

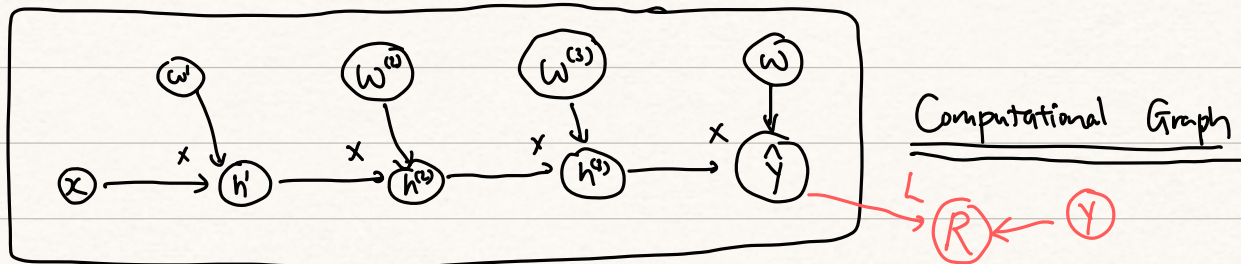
$$\bullet \hat{y} = \omega h^{(3)}$$

$$\Theta = \{\omega^{(1)}, \omega^{(2)}, \omega^{(3)}, \omega\}$$

$$\mathcal{D} = \{(x, y)\}$$

$$R(\Theta) = L(\hat{y}, y)$$

$$= L(\hat{y}(\Theta), y)$$



→ Forward Propagation

given $x, \omega^{(1)}, \omega$, calculate $h^{(1)}, h^{(2)}, h^{(3)}, \hat{y}$

→ Backward Propagation

calculate $\frac{dR}{d\omega}, \frac{dR}{d\omega^{(3)}}, \frac{dR}{d\omega^{(2)}}, \frac{dR}{d\omega^{(1)}}$

$$a) \frac{\partial R}{\partial \hat{y}} = \frac{\partial L}{\partial \hat{y}} \rightarrow \textcircled{P} \quad (\text{store})$$

$$b) \frac{\partial R}{\partial h^{(3)}} = \frac{\partial R}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial h^{(3)}} = P \cdot \omega \rightarrow \textcircled{p^{(3)}} \quad (\text{store})$$

$$\Rightarrow \frac{\partial R}{\partial \omega^{(3)}} = p^{(3)} \cdot h^{(2)}$$

$$c) \frac{\partial R}{\partial h^{(2)}} = p^{(3)} \cdot \omega^{(3)} \rightarrow \textcircled{p^{(2)}} \quad (\text{store})$$

$$\Rightarrow \frac{\partial R}{\partial \omega^{(2)}} = p^{(2)} h^{(1)}$$

$$d) \frac{\partial R}{\partial h^{(1)}} = p^{(2)} \omega^{(2)} \rightarrow \textcircled{p^{(1)}} \quad (\text{store})$$

$$\Rightarrow \frac{\partial R}{\partial \omega^{(1)}} = p^{(1)} h^{(0)} = p^{(1)} x$$