



## POHANG UNIVERSITY OF SCIENCE AND TECHNOLOGY

77, Cheongam-ro, Nam-gu, Pohang-si, Gyeongsangbuk-do, 37673, Korea

### Lecture 7. Regularization 정제화 Generalization 정제화

$$L_2 \text{ Regularization: } \tilde{J}(w; x, y) = \frac{1}{2} w^T w + J(w; x, y)$$

$$L_1 \text{ Regularization: } \tilde{J}(w; x, y) = \alpha \|w\|_1 + J(w; x, y)$$

Always called "Weight decay"

- Noise → on hidden units : similar with regularization
- Noise → on weights : reflect uncertainty (more stable)
- Noise → on output targets : avoid over-confident

$$\text{Parameter tying: } \Omega(w^{(A)}, w^{(B)}) = \|w^{(A)} - w^{(B)}\|_2^2$$

$$\left\{ \begin{array}{l} \text{Weight decay: } \tilde{J}(\theta; x, y) = J(\theta; x, y) + \alpha \Omega(\theta) \\ \text{Sparse Representation: } \tilde{J}(\theta; x, y) = J(\theta; x, y) + \alpha \Omega(h) \end{array} \right.$$

$$\text{Sparse Representation: } \tilde{J}(\theta; x, y) = J(\theta; x, y) + \alpha \Omega(h)$$

augmentation ↴

Summary: Most Popular Regularizers: L<sub>2</sub>, Early stopping, Dropout, Data

### Lecture 8. Optimization 정제화

{ Learning: Care about Performance (Including Test Set)

Optimization: Care about  $J(\theta)$

Empirical Risk Minimization: Learning Problem  $\rightarrow$  Optimization Problem.

$$\underset{\Delta}{\mathbb{E}}_{x,y \sim p(x,y)} [L(f(x;\theta), y)] \leftarrow \text{minimization}$$

Surrogate Loss: log-likelihood  $\leftrightarrow$  0-1 loss

(Reason: real loss is hard to optimization)

{ Batch Optimization: Use all train set

(also called deterministic)

Minibatch Optimization: Use some train set

Stochastic Optimization: Use single sample to train (also called online)

Accurate ↑

Larger Batch  $\rightarrow$  Efficiency ↓ (low reward)

Memory consuming ↑

Generalization ↓



## POHANG UNIVERSITY OF SCIENCE AND TECHNOLOGY

77, Cheongam-ro, Nam-gu, Pohang-si, Gyeongsangbuk-do, 37673, Korea

Minibatch Select : Unbiased Estimation (By suffling and chunking)

△ Most time : Shuffle once and pass through is okay.

Challenge : ill condition : input  $\rightarrow$  affect output extremely.

Taylor expansion  $\rightarrow x \rightarrow x - \epsilon g$

local minima : many identify model identifiability

saddle points & flat regions :

cliffs & exploding gradients

local and global structure

## Stochastic Gradient Decent (SGD)

N.B. Stochastic, but use mini-batch

Algorithm

$$\frac{k}{T}$$

Requirement : Learning Rate  $\{\varepsilon_k\}$

$$\varepsilon_k = (1-\alpha)\varepsilon_0 + \alpha\varepsilon_T$$

Requirement : Initial Parameter  $\theta$

$T$ : few hundreds

$$k \leftarrow 1$$

$$\varepsilon_T = 1\% \text{ of } \varepsilon_0$$

while Do not meet stopping Condition do.

Small-batch  $\{x^{(i)}, \dots, x^{(m)}\}, y^{(i)}$

Gradient  $\hat{g} \leftarrow \frac{1}{m} \nabla_{\theta} L(f(x^{(i)}; \theta), y^{(i)})$

Update :  $\theta \leftarrow \theta - \varepsilon \hat{g}$

$$k+1$$

end while

Other decay :  $\alpha = \alpha_0 e^{-kt}$  &  $\alpha = \alpha_0 / (1 + kt)$

Momentum :  $v \leftarrow \alpha v - \varepsilon \nabla_{\theta} g$   $\frac{1}{m} \sum_{i=1}^m L(f(x^{(i)}; \theta), y^{(i)})$

poor condition (like ill condition)  
variance in stochastic gradient



## POHANG UNIVERSITY OF SCIENCE AND TECHNOLOGY

77, Cheongam-ro, Nam-gu, Pohang-si, Gyeongsangbuk-do, 37673, Korea

Parameter initialization strategies

(Must break Symmetry)

- Gaussian Initialization : By Gaussian or uniform distribution
- Normalized Initialization :  $W_{ij} \sim U(-\sqrt{\frac{6}{m+n}}, \sqrt{\frac{6}{m+n}})$
- Sparse Initialization : fix  $\alpha$

Grad

$$\text{Adam Algorithm: } r \leftarrow r + g \odot g \quad \begin{matrix} \xrightarrow{\text{(time applied element-wise)}} \\ \frac{\varepsilon}{\delta + \sqrt{r}} \odot g \end{matrix} \quad \begin{cases} \text{Scale does matter} \\ \text{Near 0.} \end{cases}$$

$$\text{RMSPProp Algorithm: } r \leftarrow pr + (1-p)g \odot g \quad \begin{matrix} \xrightarrow{\text{(prevent too small learning rate)}} \\ \frac{\varepsilon}{\delta + \sqrt{r}} \odot g \end{matrix}$$

$$\text{Adam (RMSPProp + Momentum)} \quad \begin{cases} s \leftarrow p_1 s + (1-p_1)g \\ r \leftarrow p_2 r + (1-p_2)g \odot g \\ \hat{s} \leftarrow \frac{s}{1-p_1^t} \\ \hat{r} \leftarrow \frac{r}{1-p_2^t} \end{cases} \quad \Delta \theta = -\varepsilon \cdot \frac{\hat{s}}{\sqrt{\hat{r}} + \delta}$$

Batch Normalization



## POHANG UNIVERSITY OF SCIENCE AND TECHNOLOGY

77, Cheongam-ro, Nam-gu, Pohang-si, Gyeongsangbuk-do, 37673, Korea

Lecture 9. Convolutional Network. → all connections / local connections

Convolution Layers

Pooling Layers

Non-Linearity Layers

Convolution.  $f(x,y) \rightarrow h(x,y) \rightarrow g(x,y), g = f * h$

$$g(x,y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x+u, y+v) \cdot h(u,v) \cdot du \cdot dv$$

Pooling : max pooling.

Pooling :  $(F-1)/2$ ; e.g.  $3 \rightarrow 1, 5 \rightarrow 2, 7 \rightarrow 3$ .

Conv Layer

ConvNets: Input → Conv → Nonlinearity → Pooling → Output

e.g. LeNet : Conv → Pooling → Conv → Pooling → Conv → FC.

Lecture 10. Recurrent Networks

Recurrent Neural Network:

$$\begin{aligned} y_{t-1} &\xrightarrow{W_{hy}} y_t \\ h_{t-1} &\xrightarrow{W_{hh}} h_t \\ x_{t-1} &\xrightarrow{W_{xh}} x_t \\ y_t &= W_{hy} \cdot h_t \end{aligned}$$

$$h_t = f_W(h_{t-1}, x_t)$$

$$= \tanh(W_{hh}h_{t-1} + W_{hx}x_t)$$

Origination: Language Models:  $p(w_1, w_2, \dots, w_N)$

$$= p(w_1) \cdot p(w_2 | w_1) \cdots p(w_N | w_1, \dots, w_{N-1})$$

$$= p(w_1) \cdot p(w_2 | w_1) \cdots p(w_N | w_{N-1}, \dots, w_{N-k})$$

Cross function:  $L = -\frac{1}{N} \sum_n w_n \log \hat{p}_n$

$$\begin{array}{c} \hat{p}_n \\ \uparrow \\ h_n \\ \downarrow \\ w_{n-2} \quad w_{n-1} \end{array}$$

Trigram NN Language Model

→

$$h_n = g(V[w_{n-1}; w_{n-2}] + c)$$

$$\hat{p}_n = \text{softmax}(W h_n + c)$$

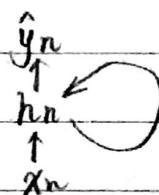


# POHANG UNIVERSITY OF SCIENCE AND TECHNOLOGY

77, Cheongam-ro, Nam-gu, Pohang-si, Gyeongsangbuk-do, 37673, Korea

$$\text{RNN} \quad h_n = g(V[x_n; h_{n-1}] + c)$$

$$\hat{y}_n = Wh_n + b$$



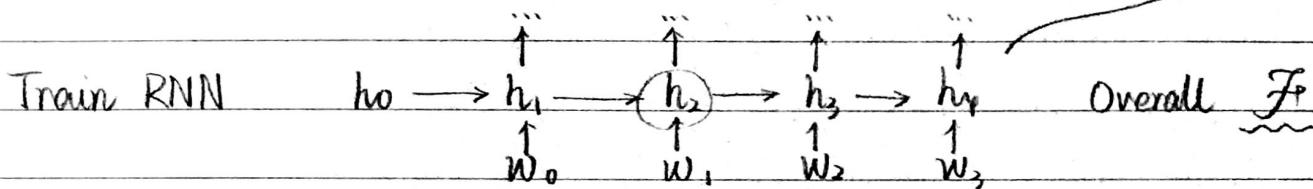
$$\hat{p}_{n+1}$$

$$h_{n+1}$$

$$x_n$$

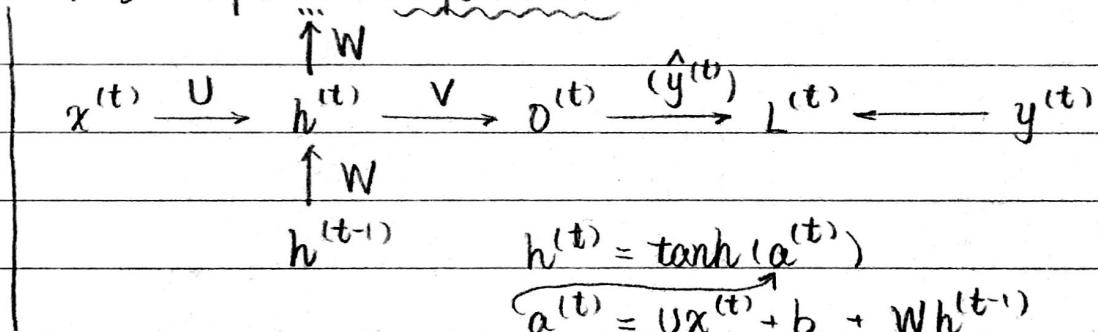
v.s. N-Gram  $\rightarrow$  Limited history

RNN  $\rightarrow$  Entire history (Long range correlations)



$$\text{e.g. } \frac{\partial \hat{F}}{\partial h_2} = \frac{\partial \hat{F}}{\partial \text{Cost}_2} \cdot \frac{\partial \text{Cost}_2}{\partial \hat{p}_2} \frac{\partial \hat{p}_2}{\partial h_2} + \frac{\partial \hat{F}}{\partial \text{Cost}_3} \cdot \frac{\partial \text{Cost}_3}{\partial \hat{p}_3} \frac{\partial \hat{p}_3}{\partial h_2} + \dots$$

N.B. depend on dependencies.



Problems : Short-term dependency.

LSTM Model & Gated Recurrent Unit.

## 1. Mini-batch Gradient Descent

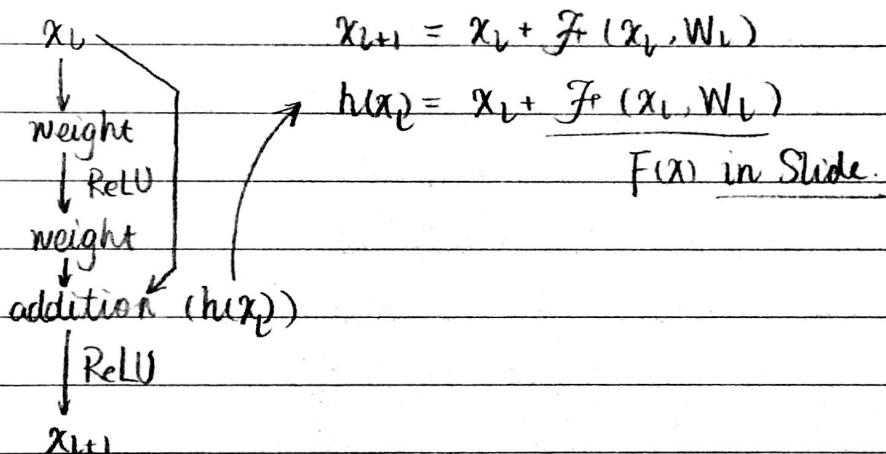
Q. Why the reward of larger samples is less than linear?

A.  $n$ : batch size.  $\sigma$ : Standard variance

batch  $\rightarrow$  average  $\rightarrow$  Standard variance:  $\frac{\sigma}{\sqrt{n}}$  standard deviation

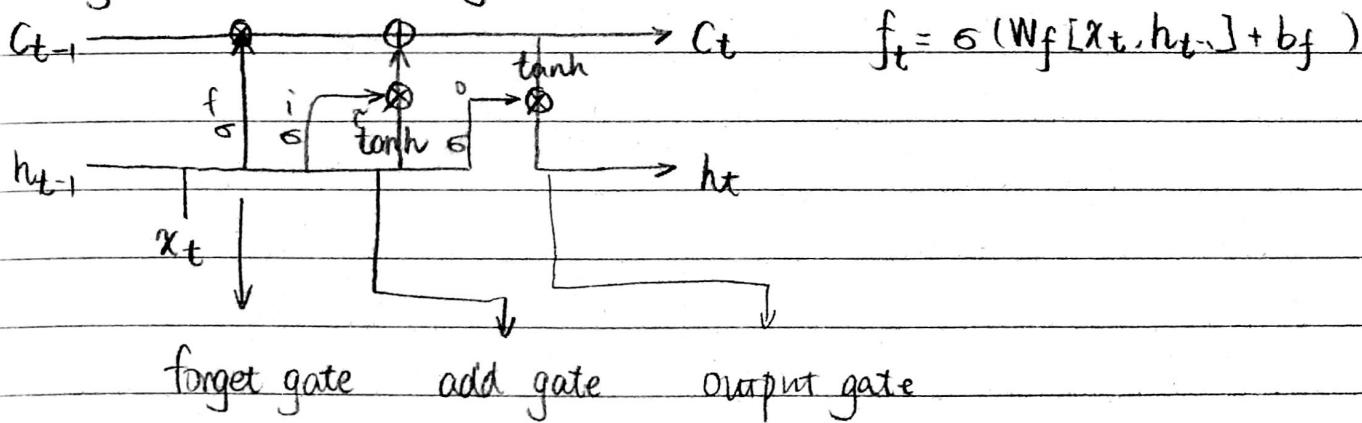
$$\text{i.e. } \hat{\mu} : \text{estimate mean} \quad SE(\hat{\mu}) = \sqrt{\text{Var}\left[\frac{1}{m} \sum_i x^{(i)}\right]} = \frac{\sigma}{\sqrt{m}} \quad \Delta$$

## 2. Residual neural network



## 3. Long short-term memory (LSTM)

fico



## 4. Back propagation through time (BPTT)

(next page.)



# POHANG UNIVERSITY OF SCIENCE AND TECHNOLOGY

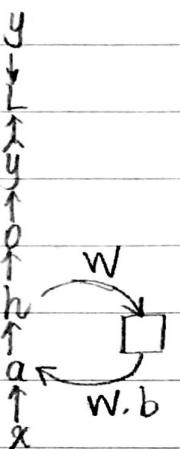
77, Cheongam-ro, Nam-gu, Pohang-si, Gyeongsangbuk-do, 37673, Korea

$$a = Ux + Wh + b \quad \Theta = \{U, V, W, b, c\}$$

$$h = \tanh(a)$$

$$o = Vh + c$$

$$\hat{y} = \text{Softmax}(o)$$



$$\nabla_{\Theta} L = \sum_t \left( \frac{\partial \Theta}{\partial \Theta} \right)^T \nabla_{\Theta} L = \sum_t \nabla_{\Theta} L$$

$$\nabla_b L = \sum_t \left( \frac{\partial \Theta}{\partial b} \right)^T \nabla_{\Theta} L = \sum_t \underbrace{\text{diag}(1-h^2)}_{\text{Remember}} \nabla_h L$$

$$\nabla_V L = \sum_t \sum_i \left( \frac{\partial \Theta}{\partial V} \right)^T \nabla_{\Theta} L = \sum_t h^T \nabla_{\Theta} L$$

$$\sum_t \sum_i \left( \frac{\partial L}{\partial \Theta_i} \right)^T \nabla_{\Theta} L$$

$$\nabla_W L = \sum_t \sum_i \left( \frac{\partial L}{\partial \Theta_i} \right)^T \nabla_{\Theta} L = \sum_t D_{hi} - \text{diag}(1-h^2) \cdot h$$

$\Delta \Theta, h$  都是 3D

✓

$$5. \text{Graph Convolution} \quad h_v^k = g \left( W_k \cdot \sum_{u \in N(v)} \frac{h_u^{k-1}}{|N(v)|} + B_k h_v^{k-1} \right)$$

$N(v)$  : v's neighbor's list