

Chapter 2. Improving Deep Neural Network

1. Setting up your ML application

① Train / dev / Test sets

Data

Training set

- Hold-out
cross validation

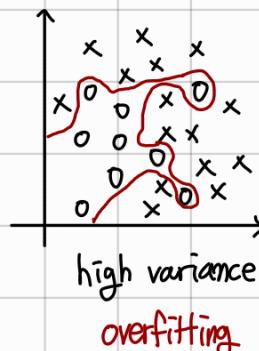
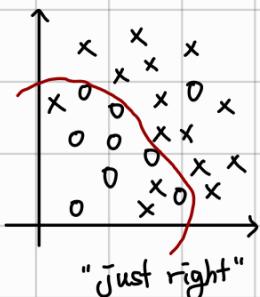
test set

- Development set
"dev"

- Mismatched train/test distribution
 - Training set : Cat pictures from webpages
 - Dev/test sets : Cat pictures from users using your app

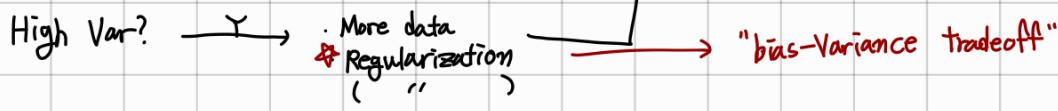
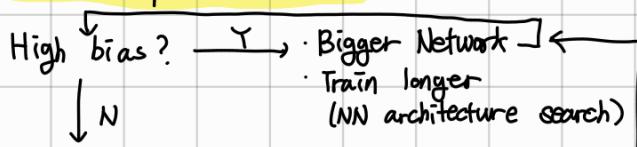
⇒ Make sure dev and test come from same distribution
- Not having a test set might be okay. (Only dev set)

② Bias / Variance



- Train set error : 1% → high var.
- Dev set error : 11% → high bias
- Human : ≈ 0% → Optimal (Bayes) error : ≈ 0%

③ Basic recipe for ML



"bias-Variance tradeoff"

2. Regularizing your neural network

① Frobenius norm

$$\|W^{[l]}\|_F^2 = \sum_{i=1}^{n^{[l]}} \sum_{j=1}^{n^{[l+1]}} (W_{ij}^{[l]})^2$$

행렬의 행 i : 현재 레이어의 뉴런수 $n^{[l]}$

“ 열 j : 이전 레이어의 뉴런수 $n^{[l+1]}$

② Regularization

• Logistic regression

$$\min_{w,b} J(w,b), \quad w \in \mathbb{R}^{n_u}, \quad b \in \mathbb{R}$$

$$J(w,b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \|w\|_F^2 + \frac{\lambda}{2m} b^2, \quad \lambda = \text{regularization parameter}$$

$$L_2 \text{ regularization } \|w\|_F^2 = \sum_{j=1}^{n_u} W_j^2 = W^T W$$

$$L_1 \text{ regularization } \frac{\lambda}{2m} \sum_{i=1}^{n_u} |w_i| = \frac{\lambda}{2m} \|w\|_1, \quad w \text{ will be sparse}$$

• Neural Network

$$J(w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{[l]}\|_F^2$$

$$d w^{[l]} = (\text{from backprop}) + \frac{\lambda}{m} w^{[l]}$$

$$\rightarrow w^{[l]} := w^{[l]} - \alpha d w^{[l]}, \quad \frac{d J}{d w^{[l]}} = d w^{[l]}$$

$$\Rightarrow w^{[l]} - \alpha \left[(\text{from backprop}) + \frac{\lambda}{m} w^{[l]} \right]$$

$$= w^{[l]} - \frac{\alpha \lambda}{m} w^{[l]} - \alpha (\text{from backprop})$$

③ Why regularization reduces overfitting

$$\lambda \uparrow : \|w^{[l]}\| \downarrow$$

$$\hat{y}^{(i)} = w^{[l]} a^{[l-1]} + b^{[l]} \rightarrow \text{조가 상대적으로 작아짐.}$$

$$J(\cdot) = \sum L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum \|w^{[l]}\|_F^2$$

구분	Lasso	Ridge
수식	$\ w\ _1$	$\frac{1}{2} \ w\ ^2$
미분	$\begin{cases} 1 & w < 0 \\ -1 & w > 0 \end{cases}$	w
특성	<ul style="list-style-type: none"> 가중치 값을 정확하게 0으로 만들 중요한 특징을 ‘선택’하는 효과 모델 전반적인 복잡도를 감소시키는 효과 모델에 Sparsity를 가함. 	<ul style="list-style-type: none"> 큰 가중치의 값을 작게 만들 모델 전반적인 복잡도를 감소시키는 효과 가중치의 값이 0이 되게 하자는 뜻함

③ Dropout regularization

- 신경망의 노드와 링크 제거
- Illustrate with layer $l=3$, keep-prob = 0.8

$$d_3 = \text{np.random.rand}(a_3.\text{shape}[0], a_3.\text{shape}[1]) < \text{keep_prob}$$

$$a_3 = \text{np.multiply}(a_3, d_3) \quad \# a_3 *= d_3$$

$$a_3 /= \text{keep_prob}$$

50 units \rightarrow 10 units shut off

$$z^{[4]} = W^{[4]} a^{[3]} + b^{[4]}$$

- No drop-out

$$z^{[1]} = W^{[1]} a^{[0]} + b^{[1]}$$

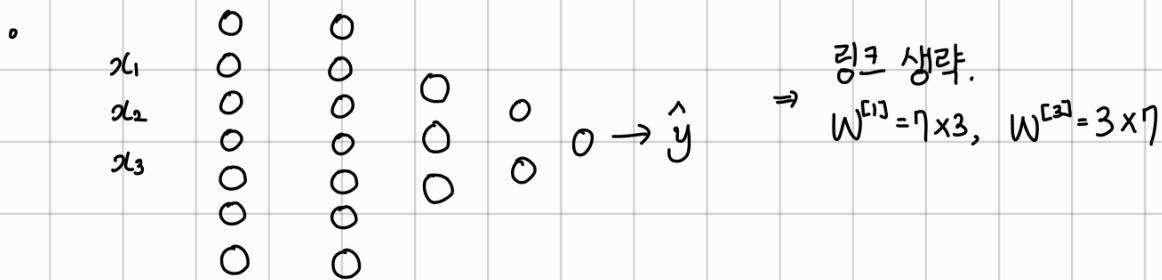
$$a^{[1]} = g^{[1]}(z^{[1]})$$

$$z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}$$

$$a^{[2]} = \dots$$

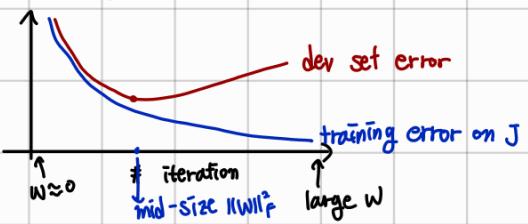
④ Understanding Dropout

- 한 가지 feature에 의존하면 안됨.
- Dropout은 L2 정규화와 유사



⑤ Other regularization methods

- Early Stopping

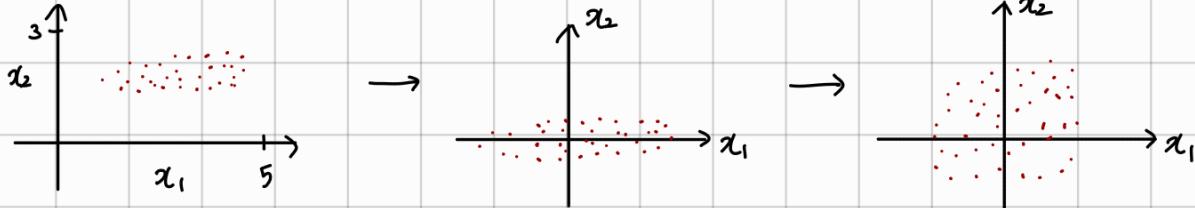


단점 : 비용함수 최적화, 정규화를 한 번에 해결할 수 없음.

3. Setting up your optimization problem

① Normalizing inputs

- Mean $\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$, $x := x - \mu$
- Variance $\sigma^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)^2$, $x := \frac{x - \mu}{\sigma}$



- Use same μ, σ to normalize test set

② Vanishing / exploding gradients

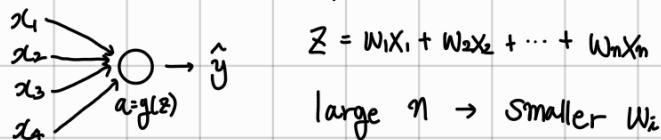
$$y = W^{[L]} W^{[L-1]} W^{[L-2]} \cdots W^{[1]} W^{[2]} W^{[1]} x$$

$$\begin{aligned} z^{[1]} &= W^{[1]} x \\ a^{[1]} &= g(z^{[1]}) = z^{[1]} \\ a^{[2]} &= g(z^{[2]}) = g(W^{[2]} a^{[1]}) \end{aligned}$$

- 해결을 위해 가중치 초기화가 중요!
- 기울기 소실 : 역전파 과정에서 입력층으로 갈수록 기울기가 점차 작아지는 현상
→ 입력층에 가까운 층들에서 가중치들이 업데이트 되지 않음.
- 기울기가 입력층으로 갈수록 점차 커지다가 가중치들이 비정상적으로 큰 값이 되며 발산하는 현상

③ Weight initialization for deep networks

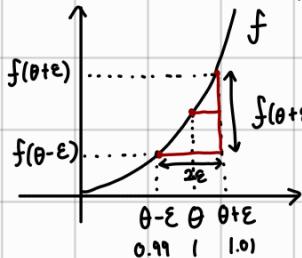
- Single neuron example



$$\text{large } n \rightarrow \text{smaller } w_i$$

- $\text{Var}(w_i) = \frac{2}{n}$, $w^{[l]} = \text{np.random.rand}(\text{slope}) * \text{np.sqrt}(\frac{2}{n^{l+1}}) \rightarrow \text{ReLU}$
- $\tanh \Rightarrow \sqrt{\frac{1}{n^{l+1}}}$: Xavier initialization $\sqrt{\frac{2}{n^{l+1} + n^l}}$

④ Numerical approximation of gradients



- $\frac{f(\theta+\epsilon) - f(\theta-\epsilon)}{2\epsilon} \approx g(\theta)$
- $\frac{(1.01)^3 - (0.99)^3}{2(0.01)} = 3.0001$
- $g(\theta) = 3\theta^2 = 3$

- $f'(\theta) = \lim_{\epsilon \rightarrow 0} \frac{f(\theta+\epsilon) - f(\theta-\epsilon)}{2\epsilon}$
- approx error : 0.0001

⑤ Gradient Checking

- $J(W^{[1]}, b^{[1]}, \dots, W^{[L]}, b^{[L]}) = J(\theta)$
- Is $d\theta$ the gradient of $J(\theta)$?
- Algorithm For each i :

$$d\theta_{\text{approx}}[i] = \frac{J(\theta_1, \theta_2, \dots, \theta_i + \epsilon, \dots) - J(\theta_1, \theta_2, \dots, \theta_i - \epsilon, \dots)}{2\epsilon}$$

$$\approx d\theta[i] = \frac{\partial J}{\partial \theta_i} \quad | \quad d\theta_{\text{approx}} \stackrel{?}{\approx} d\theta$$

Check $\frac{\|d\theta_{\text{approx}} - d\theta\|_2}{\|d\theta_{\text{approx}}\|_2 + \|d\theta\|_2} \approx \begin{cases} 10^{-7} & - \text{great!} \\ 10^{-5} & \\ 10^{-3} & - \text{worry.} \end{cases}$

⑥ Gradient Checking Implementation notes.

- Don't use in training - only to debug.
- If algorithm fails grad check, look at components to try to identify bug.
- Remember regularization
- Doesn't work with dropout.
- Run at random initialization ; perhaps again after some training.

4. Optimization Algorithms

① Mini-batch gradient descent

- 미니배치 : 훈련세트를 다시 나누어 만든 훈련세트

ex) $X = [X^{(1)} \ X^{(2)} \ X^3 \ \dots \ X^{(1000)} \ | \ X^{(1001)} \ X^{(1002)} \ \dots \ X^{(2000)} \ | \ \dots \ | \ \dots \ X^{(m)}]$
mini-batch $t : X^{ftg}, Y^{ftg}$ / X^{ftg} 의 dimension : $(n_x, 1000)$ / $Y^{ftg} : (1, 1000)$

- Mini-batch gradient descent

for $t = 1, \dots, 5000$

Forward prop on X^{ftg}

$$\begin{aligned} Z^{[l]} &= W^{[l]}X^{ftg} + b^{[l]} \\ A^{[l]} &= g^{[l]}(Z^{[l]}) \\ \vdots \\ A^{[L]} &= g^{[L]}(Z^{[L]}) \end{aligned} \quad \left. \begin{array}{l} \\ \\ \end{array} \right\} \text{Vectorized implementation}$$

Compute Cost $J = \frac{1}{1000} \sum_{i=1}^{1000} L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2 \cdot 1000} \sum_l \|W^{[l]}\|_F^2$, l : mini-batch 샘플 수

Backprop to compute gradients cost J^{ftg} (using (X^{ftg}, Y^{ftg}))

$$W^{[l]} := W^{[l]} - \alpha \delta W^{[l]}, \quad b^{[l]} := b^{[l]} - \alpha \delta b^{[l]}$$

- epoch : 훈련세트로 1번 통과

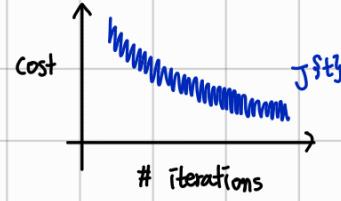
repeat

② Understanding mini-batch gradient descent

- Batch gradient descent



- Mini-batch gradient descent



- Choose mini-batch size

If mini-batch size = m : Batch Gradient descent

If mini-batch size = 1 : Stochastic gradient descent

Batch gradient descent (mini-batch size = m) → Too long per iteration

Stochastic gradient descent → Lose speed up from vectorization.

⇒ In-between (size not too big/small) → fastest learning.

- recommendation

If small training set : Use batch gradient descent. ($m \leq 2000$)

Typical mini-batch sizes : 2의 거듭제곱 (64, 128, 256, 512)

Make sure mini-batch fit in CPU / GPU memory.

③ Exponentially Weighted averages

- $V_t = \beta V_{t-1} + (1-\beta) \theta_t \leftarrow$ 자수 가중 이동 평균

$$\beta = 0.9, V_t \text{ approximately average } \approx \frac{1}{1-\beta}$$

$$\text{ex)} V_{100} = 0.1 \theta_{100} + 0.9 V_{99}$$

$$= 0.1 \theta_{100} + 0.1 \times 0.9 \cdot \theta_{99} + 0.1 (0.9)^2 \theta_{98} + 0.1 (0.9)^3 \theta_{97} + 0.1 (0.9)^4 \theta_{96}$$

$$0.9^{10} \approx 0.35 \approx \frac{1}{e}, \frac{(1-\varepsilon)}{0.9} = \frac{1}{e}$$

- Implementing exponentially weighted averages

$$V_0 = 0$$

$$V_1 = \beta V_0 + (1-\beta) \theta_1$$

$$V_2 = \beta V_1 + (1-\beta) \theta_2$$

$$V_3 = \beta V_2 + (1-\beta) \theta_3$$

$$V_\theta := 0$$

$$V_\theta := \beta V + (1-\beta) \theta,$$

$$V_\theta := \beta V_\theta + (1-\beta) \theta_2$$

$$V_\theta = 0$$

Repeat {

Get next θ_t

$$V_\theta := \beta V_\theta + (1-\beta) \theta_t$$

④ Bias correction in exponentially weighted average

- 처음 데이터에는 이전의 데이터가 존재하지 않으므로 처음 데이터로 0에 대한 가중평균을 구하게 되면 값이 매우 작아짐.

\Rightarrow 이를 수정하는 방법이 Bias correction

- V_t 대신 $\frac{V_t}{1-\beta^t}$ 사용, 추정치 : $\frac{\beta V_{t-1} + (1-\beta) \theta_t}{1-\beta^t}$

⑤ Gradient descent with momentum

- Momentum :

On iteration t :

Compute dW, db on correct mini-batch.

$$V_{dW} = \beta V_{dW} + (1-\beta) dW \quad "V_\theta = \beta V_\theta + (1-\beta) \theta"$$

$$V_{db} = \beta V_{db} + (1-\beta) db$$

$$W = W - \alpha V_{dW}, b := b - \alpha V_{db}$$

⑥ RMS prop (Root Mean Square prop)

On iteration t :

Compute dW, db on correct mini-batch

$$S_{dw} = \beta S_{dw} + (1-\beta) dW^2 \quad \leftarrow \text{element-wise small}$$

$$S_{db} = \beta S_{db} + (1-\beta) db^2 \quad \leftarrow \text{large}$$

$$W := W - \alpha \frac{dW}{\sqrt{S_{dw}}} \quad b := b - \alpha \frac{db}{\sqrt{S_{db}}}$$

⑦ Adam (Adaptive moment estimation) optimization algorithm

- $V_{dw} = 0, S_{dw} = 0, V_{db} = 0, S_{db} = 0$

On iteration t :

Compute dW, db using correct mini-batch

$$V_{dw} = \beta_1 V_{dw} + (1-\beta_1) dW, \quad V_{db} = \beta_1 V_{db} + (1-\beta_1) db \quad \leftarrow \text{"momentum" } \beta_1$$

$$S_{dw} = \beta_2 S_{dw} + (1-\beta_2) dW^2, \quad S_{db} = \beta_2 S_{db} + (1-\beta_2) db^2 \quad \leftarrow \text{"RMS prop" } \beta_2$$

$$V_{dw}^{\text{corrected}} = V_{dw} / (1-\beta_1^t), \quad V_{db}^{\text{corrected}} = V_{db} / (1-\beta_1^t)$$

$$S_{dw}^{\text{corrected}} = S_{dw} / (1-\beta_2^t), \quad S_{db}^{\text{corrected}} = S_{db} / (1-\beta_2^t)$$

$$W := W - \alpha \frac{V_{dw}^{\text{corrected}}}{\sqrt{S_{dw}^{\text{corrected}}} + \epsilon}, \quad b := b - \alpha \frac{V_{db}^{\text{corrected}}}{\sqrt{S_{db}^{\text{corrected}}} + \epsilon}$$

- Hyperparameters choice

α : needs to be tune

$$\beta_1 : 0.9 \quad (dw) \rightarrow 1\text{차 모멘트}$$

$$\beta_2 : 0.999 \quad (dw^2) \rightarrow 2\text{차 모멘트}$$

$$\epsilon : 10^{-8}$$

⑧ Learning Rate Decay

$$\alpha = \frac{1}{1 + \text{decayRate} \times \text{epochNumber}} \alpha_0$$

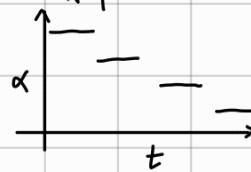
- 1 epoch = 1 pass through data

Epoch	α
1	0.1
2	0.067
3	0.05
4	0.04

- Other learning rate decay methods

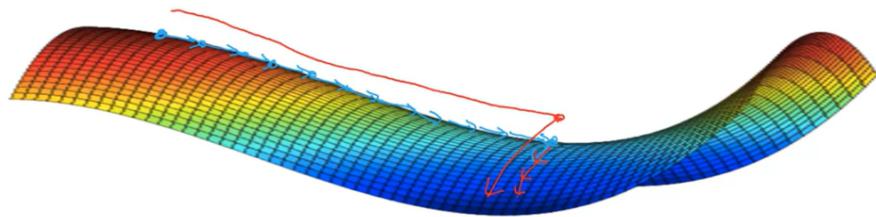
$$\alpha = 0.95^{\text{epoch-num}} \cdot \alpha_0 \quad \text{- exponentially decay}$$

$$\alpha = \frac{k}{\sqrt{\text{epoch-num}}} \cdot \alpha_0 \quad \text{or} \quad \frac{k}{\sqrt{t}} \cdot \alpha_0$$



⑨ The problem of local optima

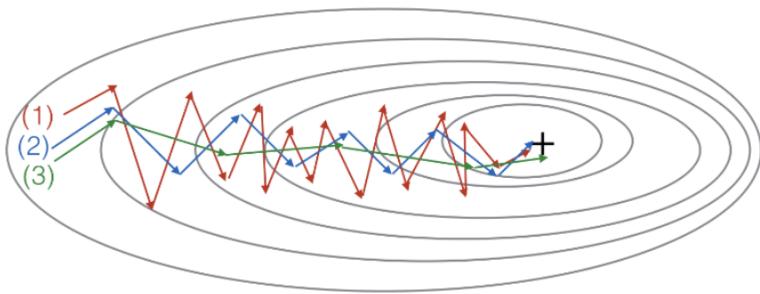
- 안장점 (Saddle point) : 비용함수에서 기울기가 0인 지점
- problem of plateaus



- Unlikely to get stuck in a bad local optima
- Plateaus can make learning slow
⇒ Momentum, RMSprop, Adam이 도움이 될 수 있음.

* 퀴즈

- $\alpha^{[2]} \nabla J^{(3)}$: 네 번째 미니배치, 세 번째 예일 때, 두 번째 레이어
- $V_t = \beta V_{t-1} + (1-\beta) \theta_t \Rightarrow$ 평형 보정 : $\frac{V_t}{1-\beta^t}$
- 곡선의 noise가 심한 경우 β 값이 더 작음.



- (1) : 경사 하강
 - (2) : 모멘텀이 있는 경사하강 (small β)
 - (3) : 모멘텀이 있는 경사 하강 (large β)
- NN에서 J 가 작은 값에 도달하도록 할 수 있는 기법
 - 1) 모멘텀이 있는 경사 하강 사용
 - 2) 입력 데이터 정규화
 - 3) 가중치에 대한 더 나은 무작위 초기화 시도.

5. Hyperparameter Tuning

① Tuning process

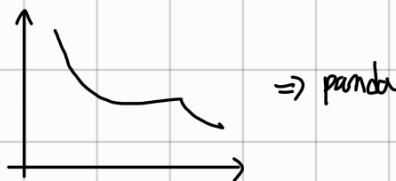
- Hyper parameters
 - 1) α
 - 2) β (default: 0.9)
 - 3) $\beta_1, \beta_2, \epsilon \rightarrow$ Adam 사용하는 경우 0.9, 0.999, 10^{-8}
 - 4) # layers
 - 5) # hidden units
 - 6) learning rate decay
 - 7) mini-batch size
- grid를 사용하지 말고, 랜덤하게 배치하여 하이퍼파라미터 선택

② Using an appropriate scale to pick hyperparameters.

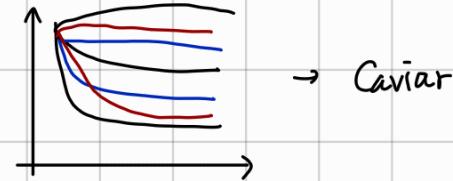
- $\eta^{[L]} = 50, \dots, 100$
- # layers L : 2-4
- $\alpha = 0.0001, \dots, 1 \Rightarrow r = -4 * np.random.random(), r \in [-4, 0], \alpha = 10^r$
- $\beta = 0.9, \dots, 0.999$

③ Hyperparameters tuning in practice : Pandas VS Caviar

- Babysitting one model vs Training many models in parallel



\Rightarrow panda



\rightarrow Caviar

6. Batch Normalization

① Normalizing activations in a network

- $\mu = \frac{1}{m} \sum_i z^{(i)}$, $x = x - \mu$

$$\sigma^2 = \frac{1}{m} \sum_i (z^{(i)} - \mu)^2, \quad x = x / \sigma$$

- Batch Norm

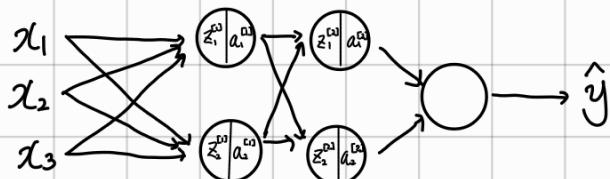
$$\mu = \frac{1}{m} \sum_i z^{(i)}, \quad \sigma^2 = \frac{1}{m} \sum_i (z^{(i)} - \mu)^2, \quad z_{\text{norm}}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$\tilde{z}^{(i)} = \gamma z_{\text{norm}}^{(i)} + \beta \quad (\gamma, \beta : \text{learnable parameters of Model})$$

$$\gamma = \sqrt{\sigma^2 + \epsilon}, \quad \beta = \mu \quad \text{then } \tilde{z}^{(i)} = z^{(i)}$$

② Fitting Batch Norm into a neural network.

- Adding Batch Norm to a network



$$X \xrightarrow{W^{[l]}, b^{[l]}} Z^{[l]} \xrightarrow[\text{BatchNorm(BN)}]{\gamma^{[l]}, \beta^{[l]}} \tilde{Z}^{[l]} \rightarrow a^{[l]} = g^{[l]}(\tilde{Z}^{[l]}) \xrightarrow{W^{[l+1]}, b^{[l+1]}} Z^{[l+1]} \xrightarrow[\text{BN}]{\gamma^{[l+1]}, \beta^{[l+1]}} \tilde{Z}^{[l+1]} \rightarrow a^{[l+1]} \rightarrow \dots$$

Parameters: $W^{[l]}, b^{[l]}, W^{[l+1]}, b^{[l+1]}, \dots, W^{[L]}, b^{[L]} + \beta^{[l]}, \gamma^{[l]}, \beta^{[l+1]}, \gamma^{[l+1]}, \dots, \beta^{[L]}, \gamma^{[L]}$

$$\rightarrow d\beta^{[l]}, \quad \beta^{[l+1]} = \beta^{[l]} - \alpha d\beta^{[l]}$$

- Working with mini-batches

$$X^{[1]} \xrightarrow{W^{[l]}, b^{[l]}} Z^{[1]} \xrightarrow[\text{BN}]{\gamma^{[l]}, \beta^{[l]}} \tilde{Z}^{[1]} \rightarrow g^{[l]}(\tilde{Z}^{[1]}) = a^{[1]} \xrightarrow{W^{[l+1]}, b^{[l+1]}} Z^{[2]} \rightarrow \dots$$

$$X^{[2]} \rightarrow Z^{[2]} \xrightarrow[\text{BN}]{\gamma^{[2]}, \beta^{[2]}} \tilde{Z}^{[2]} \rightarrow \dots$$

parameters: $W^{[l]}, b^{[l]}, \beta^{[l]}, \gamma^{[l]}$

$$\tilde{Z}^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]}$$

$$= W^{[l]} a^{[l-1]}$$

$$Z^{[l]} \rightarrow (\eta^{[l]}, 1)$$

$$\tilde{Z}^{[l]} = \gamma^{[l]} Z_{\text{norm}}^{[l]} + \beta^{[l]}$$

- Implementing gradient descent

for $t=1, \dots$, # Mini-Batches

compute forward prop on $X^{t, j}$

In each hidden layer, use BN to $\hat{z}^{[l]}$ with $\tilde{z}^{[l]}$

Use backprop to compute $dW^{[l]}, db^{[l]}, d\beta^{[l]}, d\gamma^{[l]}$

Update parameters $W^{[l]} := W^{[l]} - \alpha dW^{[l]}$

$$\beta^{[l]} := \beta^{[l]} - \alpha d\beta^{[l]}$$

$$\gamma^{[l]} := \dots$$

Work w/ momentum, RMS prop, Adam, ...

③ Why does Batch Norm work?

- Internal Covariance Shift

네트워크의 각 레이어나 Activation마다 입력값의 분산이 달라지는 현상

\Rightarrow Covariate Shift : 이전 레이어의 파라미터 변화로 인하여 현재 레이어의 입력의 분포가 바뀌는 현상

\Rightarrow Internal Covariate Shift : 레이어를 통과할 때마다 Covariate Shift가 일어나면서 입력의 분포가
약간씩 변하는 현상

- 각 미니 배치가 하나의 미니 배치에서 계산된 평균과 분산으로 조정됨. \rightarrow 약간의 노이즈 존재
- 미니 배치에서 $\hat{z}^{[l]}$ 로 조정 시 노이즈가 생김. dropout과 비슷하게 각 hidden layer's activation의 noise가 생생됨.
- 약간의 정규화 효과가 존재함.

④ Batch Norm at test time

- μ, σ^2 : estimate using exponentially weighted average (across mini-batch)

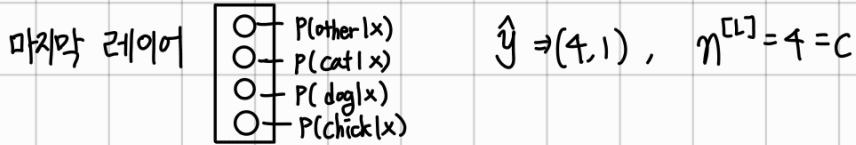
$$X^{1, j}, X^{2, j}, X^{3, j}, \dots \downarrow \mu^{1, j}, \mu^{2, j}, \dots \longrightarrow \mu$$

$$z_{norm} = \frac{z - \mu}{\sqrt{\sigma^2 + \epsilon}}, \quad \tilde{z} = \gamma z_{norm} + \beta$$

7. Multi-class classification

① Softmax regression

Example) 클래스 4개 (0-어느곳에도 속하지 않는 경우, 1-고양이, 2-개, 3-꿩아리)



$$z^{[L]} = W^{[L]} a^{[L-1]} + b^{[L]}$$

$$\text{Activation function : } t = e^{(z^{[L]})}, \quad a^{[L]} = \frac{e^{z^{[L]}}}{\sum_i t_i}$$

$$\text{ex)} z^{[L]} = \begin{bmatrix} 5 \\ 2 \\ -1 \\ 3 \end{bmatrix}, \quad t = \begin{bmatrix} e^5 \\ e^2 \\ e^{-1} \\ e^3 \end{bmatrix} = \begin{bmatrix} 148.4 \\ 7.4 \\ 0.4 \\ 20.1 \end{bmatrix}, \quad \frac{4}{176.3} t_i = [0.002, 0.042, 0.842, 0.114]$$

$$\text{layer L} \Rightarrow 1) \frac{e^5}{176.3} = 0.842, \quad 2) \frac{e^2}{176.3} = 0.042, \quad 3) \frac{e^{-1}}{176.3} = 0.002, \quad 4) \frac{e^3}{176.3} = 0.114$$

- 신경망 결과값 : $a^{[L]} = g^{[L]}(z^{[L]})$ or \hat{y}

② Training a softmax classifier

- Softmax 회귀는 로지스틱 회귀를 C개의 class로 일반화함.

- Loss function $L(\hat{y}, y) = -\sum_j y_j \log \hat{y}_j$

$$J(W^{[L]}, b^{[L]}, \dots) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

- $z^{[L]} \rightarrow a^{[L]} = \hat{y} \rightarrow L(\hat{y}, y)$

- Backprop : $dz^{[L]} = \hat{y} - y$

8. Programming frameworks

① Deep Learning frameworks

- Caffe / Caffe2
 - CNTK
 - DL4J
 - Keras
 - Lasagne
 - mxnet
 - PaddlePaddle
 - TensorFlow
 - Theano
 - Torch
- Choosing deep learning frameworks
 - Ease of programming (development and deployment)
 - Running Speed
 - Truly open (open source with good governance)