

Course 2

Hyperparameter

Tuning

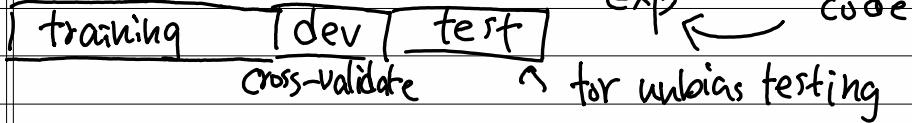
COURSE 2

- Iterative process

for hyperparameter tuning

idea

DATA



Previous

70 / 30 or 60 / 20 / 20.

100 ~ 10000 samples

big data

1,000,000 data, 10,000 test might be enough
10,000

→ 98% / 1% / 1%

Mismatch train/test distribution

training

e.g. cat on web

dev / test

cat on app

- try making

dev / test

with same distribution

Not having test set might be OK

BIAS & VARIANCE

- under fitting — high bias (not fitting well)
- just right
- over fitting — high variance (too flexible)

training set error: 1% 15% 15%

dev set error : 11% 16% 30%

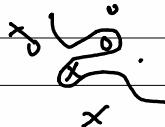
→ overfitting -underfit high bias
high variance high bias var

0.5%

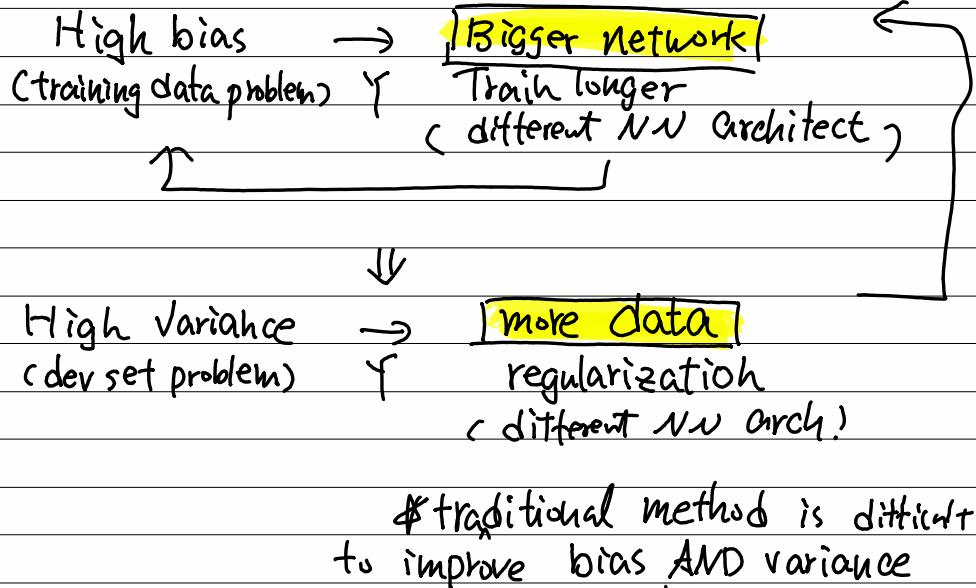
↙ high bias + high variance

1%

↙ low bias var



RECIPE for ML



high bias → more training data is not useful

high variance → regularization
more training data.

Regularization

Logistic Regre..

$$\min_{w, b} J(w, b)$$

$$w \in \mathbb{R}^{k_x}, b \in \mathbb{R}$$

λ : regularization para

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \|w\|_2^2 + \frac{\lambda}{2m} b^2$$

$$\Delta L_2 \text{ regu } \|w\|_2 = \sqrt{\sum_{j=1}^{k_x} w_j^2} = \sqrt{w^T w}$$

omit

$$L_1 \text{ regularization } \|w\|_1, \frac{\lambda}{2m} \sum_{j=1}^{k_x} |w_j| = \frac{1}{2m} \|w\|_1$$

Neural Network

$$J(w^{[0]}, b^{[0]}, \dots, w^{[L]}, b^{[L]}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{[l]}\|^2$$

$$\|w^{[l]}\|^2 = \sum_{i=1}^{n^{[l+1]}} \sum_{j=1}^{n^{[l]}} (w_{ij}^{[l]})^2 \quad w = (w^{[0]}, \dots, w^{[L]}) \quad \text{square norm}$$

Frobenius form sum of square norm of matrix

$$dw^{[l]} = (\text{from back prop}) + \frac{\lambda}{m} w^{[l]}$$

$$w^{[l]} = w^{[l]} - \alpha dw^{[l]}$$

$$= w^{[l]} - \underbrace{\frac{\alpha \lambda}{m} w^{[l]}}_{\text{weight decay}} - \alpha (\text{from back prop})$$

weight decay

$$* J(w^{[0]}, b^{[0]}) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \|w^{[0]}\|^2$$

$\lambda \uparrow \Rightarrow w \downarrow$ (bcz we penalize high cost)
 \Rightarrow simpler NN

Dropout regularization

- drop out random nodes in NN
(hidden unit)

INVERTED DROPOUT

e.g. at layer 3,

$$\text{keep-prop} = 0.8$$

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

$$a_3 = \text{np.multiply}(a_3, d_3) \quad < \text{keep-prop}$$

* $\Rightarrow a_3 \neq \text{keep-prop}$

e.g. total 50 units, start off 10 units (20%)

$$\overset{(20)}{\tilde{z}} = \overset{(20)}{W} \cdot \overset{(20)}{a} + \overset{(20)}{b}$$



bump up by 20%

* no drop-out at test time

why drop-out works?

* intuition \rightarrow can't rely on one feature, so we spread out weights.

* keep-drop can vary in layers \rightarrow shrink weights similar to L2

- Computer Vision uses drop out bcs high dimension

- Downside. J might looks weird in the plot

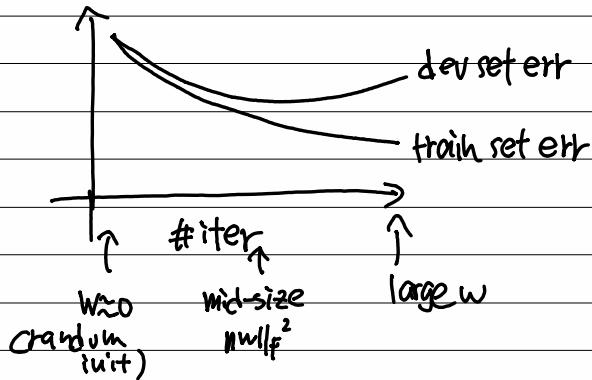


Reguralization methods

Data augmentation

- flipping, random crop images. (flipping cats)
- Adding distortion to characters etc

Early stopping



↳ breaks orthogonalization

Normalizing Input

subtract mean

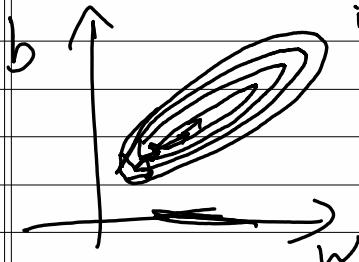
normalize variance

$$X = X - \mu$$

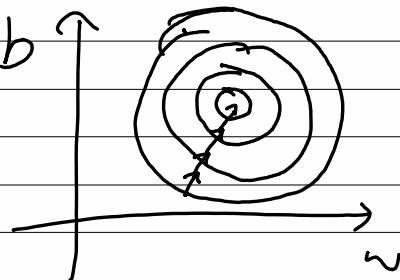
$$X \sim \sigma^2$$

* use same μ & σ^2 for
testing / eval / train sets

unnormalized b might be in very diff range normalized



require small learning rate
for gradient descent



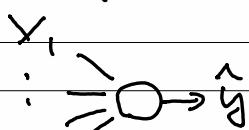
better!
sater to take larger step!

vanishing / exploding gradient

(assuming linear activation) $\hat{y} = w \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}^{L-1} x$ - exploding

back-prop $\hat{y} = w \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}^{L-1} x$ - vanishing

solution (initialization)



$$z = w_1 x_1 + \dots + w_n x_n + b$$

large $n \rightarrow$ small w_i

$$\text{Var}(w_i) = \frac{1}{n} \quad \text{if ReLU}, \quad \text{Var}(w_i) = \frac{2}{n}$$

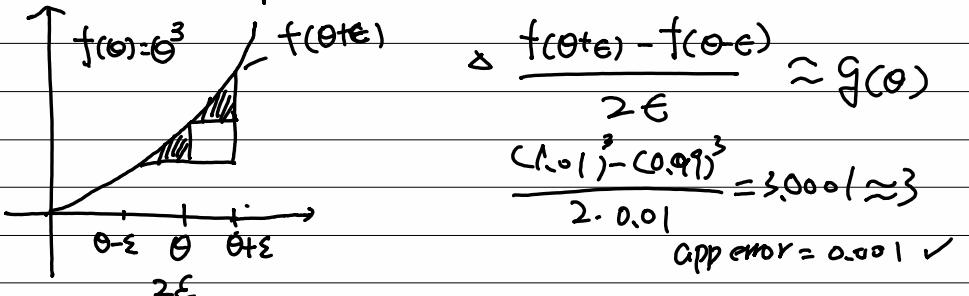
$$w^{[L]} = \text{np.random.rand(slope)} \times \text{np.sqrt}\left(\frac{2}{n^{[L]}}\right)$$

ReLU

$$\tanh = \sqrt{\frac{1}{1+e^{-x}}}$$

Xavier initialization $\sqrt{\frac{2}{n^{[L]} + n^{[L]}}}$

Numeric Approx of Gradients



$$f'(\theta) = \lim_{\epsilon \rightarrow 0} \frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon} \quad O(\epsilon^2)$$

$$f'(\theta) \approx \lim_{\epsilon \rightarrow 0} \frac{f(\theta + \epsilon) - f(\theta)}{\epsilon} \quad \text{err: } O(\epsilon)$$

Δ we use two-sided for gradient checking (more accurate)
instead of one-sided checking

DEBUG TOOL

Gradient Checking (Grad Check)

$$J(w^{(1)}, b^{(1)}, \dots, w^{(m)}, b^{(m)}) = J(\theta)$$

$$J(dw^{(1)}, db^{(1)}, \dots, dw^{(m)}, db^{(m)}) = J(d\theta)$$

for each i:

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

$$\approx d\theta^{(i)} = \frac{\partial J}{\partial \theta_i} \Rightarrow \text{check } d\theta_{\text{app}} \approx d\theta?$$

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

Grad Check Implementation

- Don't use in training
- look at component to try to identify the bug
- remember regularization
- doesn't work with dropout
 - keep-prop=1 for debug
- run at random initialization

mini-batch gradient descent

split [training set] into batches

$$X^{t \in S}, Y^{t \in S}$$

{

△ notation

$$X^{(i)}$$

batch v.s. mini-batch

$$X^{[l]}$$

(all) (small batches)

$$X^{[t \in S]}$$

curly bracket.

training set batch

EXAMPLE

5,000,000 data \rightarrow 5000 mini-batch, each has 1000 data.

for $t = 1$ to 5000

forward prop on $X^{(t)}$

$$\text{cost } J^{(t)} = \frac{1}{1000} L(Y^{(t)}, Y^{(t)}) + \frac{\lambda}{2 \cdot 1000} \|W^{(t)}\|_F^2$$

backprop using $J^{(t)}$

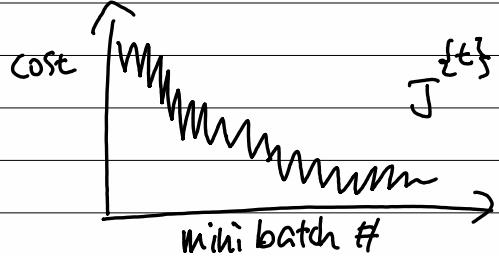
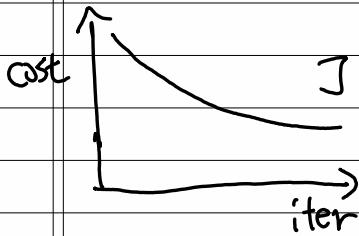
$$w^{(t+1)} = w^{(t)} - \alpha d w^{(t)} \quad b^{(t+1)} = b^{(t)} - \alpha d b^{(t)}$$

1 epoch - a single pass of 'All' data

BATCH - 1 step descent

MINIBATCH - 5000 step descent.

WHY MINI-BATCH WORK?



noisy, but trend down
bcz batch might be difficult

mini-batch size

size = $m \rightarrow$ batch gradient descent

size = 1 \rightarrow stochastic gradient descent.

BATCH - too long per iteration

STOCHASTIC - lose the benefit of vectorization

GUIDE:

Small training set:

(~2000) use batch gradient

typical size

64 ~ 512

make sure mini-batch fits CPU/GPU memory.

Exponentially Weighted Average

$$V_t = \beta V_{t-1} + (1-\beta) \theta_t$$

moving

△

V_t is approx average
over $\frac{1}{1-\beta}$ samples

$$\beta = 0.9$$

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

e.g. $\beta = 0.9 \rightarrow 10$ samples

$$V_{99} = 0.1 \theta_{99} + 0.9 V_{98}$$

$\beta = 0.98 \rightarrow 50$ samples

$$V_{98} = 0.1 \theta_{98} + 0.9 V_{97}$$

:

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

$$= 0.1 \theta_{100} + 0.9 (0.1 \theta_{99} + 0.9 V_{98})$$

$$= 0.1 \theta_{100} + 0.1 \cdot 0.9 \theta_{99} + 0.9^2 (0.1 \theta_{98} + 0.9 V_{97})$$

$$= 0.1 \theta_{100} + 0.1 \cdot 0.9 \theta_{99} + 0.1 \cdot (0.9)^2 \theta_{98} + (0.9)^3 V_{97} + \dots$$

- - -

$(0.9)^{10} \approx 0.35 \approx \frac{1}{e}$ after $\frac{1}{e}$, it drops faster.

$$(0.98)^{50} \approx \frac{1}{e}$$

Bias Correction in EWMA

bias because $V_0 = 0$

$$U_t = \beta V_{t-1} + (1-\beta) \Theta_t$$

incorrect U in the beginning of series

$$V_0 = 0$$

$$U_1 = 0.9V_0 + 0.1\Theta_1 = 0.1\Theta_1$$

$$U_2 = 0.9U_1 + 0.1\Theta_2$$

$$= 0.09\Theta_1 + 0.1\Theta_2$$

⋮

Correction

$$\frac{U_t}{1-\beta^t}, \text{ correct in the beginning}$$

$t \uparrow; \beta^t \sim 0$, correction does not affect U_t .

Gradient Descent with Momentum

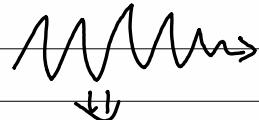
Compute EWA of gradient, then update

on iter t

compute dw db on current mini-batch

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

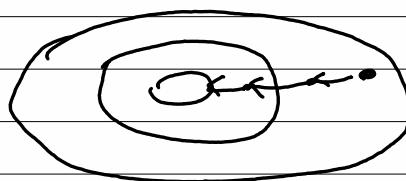
$$V_{db} = \beta V_{db} + (1-\beta)db$$



$$w = w - \alpha V_{dw}$$

$$B = B - \alpha V_{db}$$

dampen oscillation
overshoot.



ball rolling.

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

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

Acceleration

velocity

friction

hyper parameters

α - learning rate

β

0.9 usually works

RMSprop (root mean square prop)

On iter t.

compute dw, db on mini-batch

$$Sdw = \beta Sdw + (1-\beta)dw^2 \quad \text{- element-wise square}$$

$$Sdb = \beta Sdb + (1-\beta)db^2$$

$$w = w - \alpha \frac{dw}{\sqrt{Sdw} + \epsilon} \quad \text{P.G.}$$

$$b = b - \alpha \frac{db}{\sqrt{Sdb} + \epsilon} \quad \text{avd divided by 0}$$

\nearrow descent

$Sdw \downarrow \quad w \uparrow$

$Sdb \uparrow \quad b \downarrow$

Adam Optimization Algorithm

Adaptive
moment
estimation

Momentum + RMSdrop

$$V_{dw} = 0 \quad S_{dw} = 0 \quad V_{db} = 0, \quad S_{db} = 0$$

On iter t:

compute dw, db of mini-batch

$$V_{dw} = \beta_1 V_{dw} + (1-\beta_1) dw \quad V_{db} = \beta_1 V_{db} + (1-\beta_1) db$$

$$S_{dw} = \beta_2 S_{dw} + \beta_2 dw^2 \quad S_{db} = \beta_2 S_{db} + (1-\beta_2) db^2$$

bias correction

$$V_{dw}^{\text{corrected}} = \frac{V_{dw}}{1-\beta_1^t} \quad V_{db}^{\text{corr}} = \frac{V_{db}}{1-\beta_1^t}$$

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

RMSprop

$$w = w - \alpha \frac{V_{dw}^{\text{corr}}}{\sqrt{S_{dw}^{\text{corr}}} + \epsilon} \quad b = b - \alpha \frac{V_{db}^{\text{corr}}}{\sqrt{S_{db}^{\text{corr}}} + \epsilon}$$

hyperparameter α - tube

$$\beta_1: 0.9 \rightarrow dw$$

first momentum

$$\beta_2: 0.999 \rightarrow dw^2$$

second momentum

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

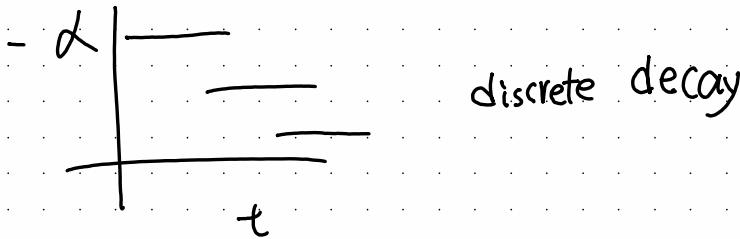
Learning Rate Decay

slowly reduce α (to avoid overshoot when close to solution)

$$\alpha = \frac{1}{\text{lr decay rate} \times \text{epochnum}} \cdot \alpha_0 \Rightarrow \begin{array}{l} \text{para.} \\ \alpha_0 \\ \text{decay rate} \end{array}$$

- $\alpha = 0.95^{\text{epochnum}}$ - α_0 - exponential decay

$$- \alpha = \frac{K}{\text{epochnum}} \cdot \alpha_0 \quad \text{or} \quad \frac{K}{\Delta t} \cdot \alpha_0$$



- manual decay (manually update ...)

Local optima

most local optima happens at saddle point.

- unlikely to stuck in bad local



- plateaus can make learning slow.



Hyperparameter tuning

α , β , $\beta_1, \beta_2, \epsilon$, α

layers

1st / # hidden units second important

- learning rate decay

mini-batch

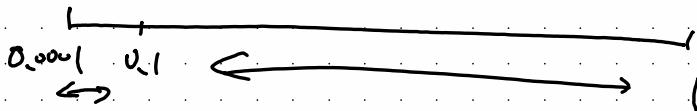
- use "random" search, not grid

- coarse to fine

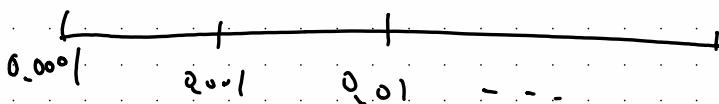
appropriate scale to pick hyper-parameter.

$n^{[l]}$, L uniform sample works

α .



in log scale



$$\gamma = -4 * \text{np.random.rand} \leftarrow \gamma \in [-4, 0]$$

$$\alpha = 10^\gamma \leftarrow \gamma \in [-4, 0]$$

Hyperparameter for EWA (exponentially weighted averages)

$$\beta = 0.9 \dots 0.999$$

$$r \in [-1, 1] \quad \overbrace{0.1 \quad 0.01 \quad 0.001}$$

$$1-\beta = 0.1 \dots 0.001$$

$$1-\beta = 10^r$$

$$\frac{1}{1-\beta} \leftarrow \text{when } \beta \rightarrow 1.$$

$$\beta = 1 - 10^r$$

Small changes have big

impact.

$$\beta = 1 - 10^{(r-1)}$$

Hyperparameter tuning

Babysitting one model \leftarrow fewer resource
Training many models in parallel

large model

Panda
Cavier.

Batch Normalization (single layer)

can we normalize $a^{(l)}$ so that training $w^{(l+1)}, b^{(l+1)}$ faster

$$\mu = \frac{1}{m} \sum z^{(l)}$$

$$\sigma^2 = \frac{1}{m} \sum_i (z_i - \mu)^2$$

$$z_{\text{norm}}^{(l)} = \frac{z^{(l)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$\tilde{z}^{(l)} = \gamma z_{\text{norm}}^{(l)} + \beta$$

$$\text{if } \gamma = \sqrt{\sigma^2 + \epsilon}$$

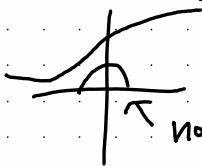
$$\beta = \mu$$

$$\text{then } \tilde{z}^{(l)} = z^{(l)}$$

$$\underbrace{z^{(l)} \dots z^{(n)}}_{z^{(l), c_i}} \xrightarrow{\gamma, \beta}$$

learnable hyperparameter

\Rightarrow use $\tilde{z}^{(l)}$ instead of $z^{(l)}$
activation

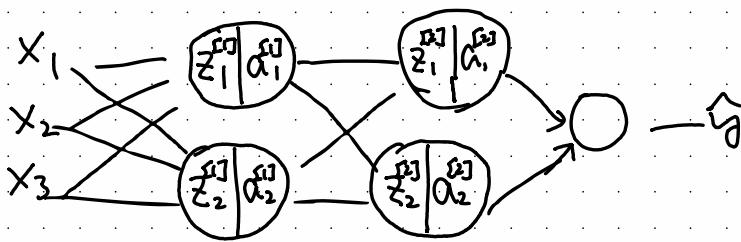


normalization

avoids value squeeze in small area

(by spreading the value)

Batch norm in a NN



$$X \xrightarrow{w^{[l]}, b^{[l]}} z^{[l]} \xrightarrow{\text{BatchNorm}, \mu^{[l]}, \sigma^2} \tilde{z}^{[l]} \xrightarrow{\alpha^{[l]}, \beta^{[l]}} a^{[l]} = g^{[l]}(\tilde{z}^{[l]})$$

$$w^{[l+1]}, b^{[l+1]} \xrightarrow{z^{[l]}} \tilde{z}^{[l]} \xrightarrow{\text{BatchNorm}, \mu^{[l]}, \sigma^2} \tilde{z}^{[l]} \xrightarrow{\alpha^{[l+1]}, \beta^{[l+1]}} a^{[l+1]}$$

Parameters:

$$(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, \dots, w^{[L]}, b^{[L]}) \\ (\beta^{[1]}, \gamma^{[1]}, \dots, \beta^{[L]}, \gamma^{[L]})$$

this β is different
from β in Momentum

ADAM
etc

(tf.nn.batch-normalization)

WORKING with MINI-BATCH

$$X^{[1]} \\ X^{[2]} \\ \vdots \\ X^{[n]}$$

parameter: $w^{[l]}, b^{[l]}, \beta^{[l]}, \gamma^{[l]}$ → $z^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]}$

Batch-Normalization

$(n^{[l]}, h^{[l]})$ eliminate

$$\Rightarrow \tilde{z}^{[l]} = \frac{z^{[l]} - \mu^{[l]}}{\sigma^{[l]}}$$

$$\tilde{z}^{[l]} = \gamma^{[l]} \tilde{z}^{[l]} + \beta^{[l]}$$

IMPLEMENTATION.

for $t = 1 \dots \# \text{MINI BATCH}$

compute forward prop on $X^{[t]}$

in each layer, use BN to replace $Z^{[e]}$ with $\tilde{Z}^{[e]}$

use backprop to compute $dW^{[e]}$, $b^{[e]}$, $ds^{[e]}$, $dr^{[e]}$

update para $W^{[e]} = W^{[e]} - \alpha dW^{[e]}$

$B^{[e]} = B^{[e]} - \alpha ds^{[e]}$

work with momentum, RMSprop, ADAM etc

Batch norm work bcz.

- learning on shifting of distribution

'covariance shift' e.g. training on black cat

✗ Batch norm reduce "cor shift" hard to generalize to yellow cat.

△ Batch norm as regularization

- mean/var of each mini-batch differ

- Add noise to value $Z^{[e]}$

- similar to dropout.

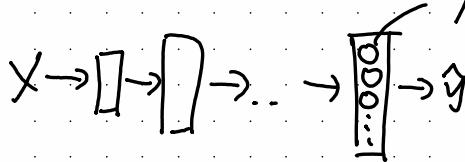
if add noise to each hidden layer's activation

→ slight regulation

Batch norm at test time

M, S^2 for testing: exponentially weighted average across all mini-batch

softmax regression



$$p(\text{class}_i | x)$$

$$N^{\text{class}} = \# \text{ class}$$

e.g. 4 class

$$z^{(l)} = w^{(l)} a^{(l-1)} + b^{(l)} \quad (4, 1)$$

Activation Function

$$t = e^{z^{(l)}}$$

$$a^{(l)} = \frac{e^{z^{(l)}}}{\sum_{j=1}^4 t_j}, \quad a_i^{(l)} = \frac{t_i}{\sum_{j=1}^4 t_j}$$

hard max

$$\begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix} \rightarrow \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

$$a^{(l)} = g^{(l)}(z^{(l)})$$

$$(4, 1) \quad (4, 1)$$

softmax activation.

vector \rightarrow vector

others output real number.

$C=2 \rightarrow$ softmax reduced to logistic regression

Loss function

$$C=4, \quad y = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \quad \hat{y} = \begin{bmatrix} 0.3 \\ 0.2 \\ 0.1 \\ 0.4 \end{bmatrix}$$
$$L(y, \hat{y}) = -\sum_{j=1}^4 y_j \log \hat{y}_j$$

in sample

maximum likelihood estimation

$$\mathbb{J}(w^{(i)}, b^{(i)}, \dots) = \frac{1}{m} \sum L(y^{(i)}, \hat{y}^{(i)})$$

$$\text{e.g. } -y_2 \log \hat{y}_2$$

reduce loss $\rightarrow \hat{y}_2 \uparrow$

Deep Learning Frameworks

TensorFlow

$$\leftarrow \{w=s\}, w-s=0$$

$$J(w) = w^2 - 10w + 25$$

↖ only need to calculate forward prop.

Gradient Tape for back prop

