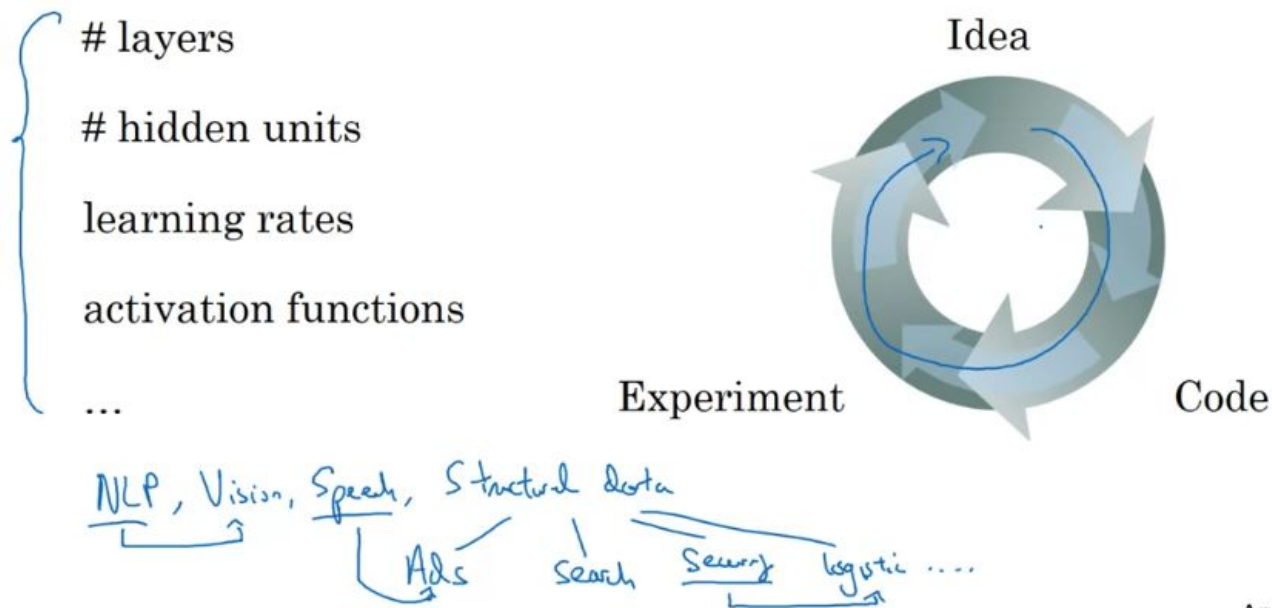# Improving Deep Neural Networks: Hyperparameter Tuning, Regularization and Optimization

27/03/2021 Minsung Kim

# Setting up your machine learning application

## Applied ML is a highly iterative process

# layers

# hidden units

learning rates

activation functions

...

Idea

Experiment

Code

NLP, Vision, Speech, Structured data

Ads    Search    Security    logistic ....

Andrew Ng

# Setting up your machine learning application

Test/dev/test sets

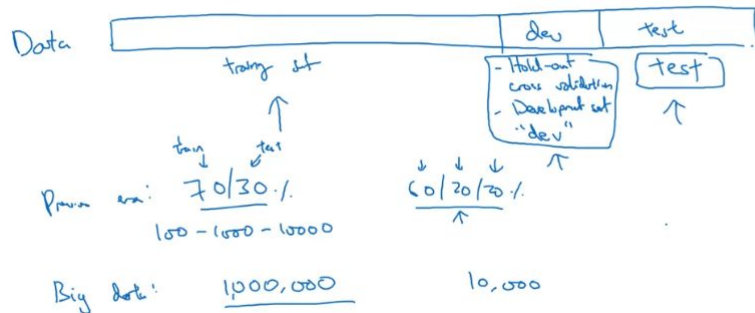Previous era: 70/30 or 60/20/20  - 데이터의 양이 작았기에

Modern era: dev/test -> 데이터가 많아서 굳이 이전처럼 나눌필요 없이 적당히

ex) total 1,000,000  -> Dev 10,000 Test 10,000

98/1/1

99/0.5/0.5



Train/dev/test sets

# Setting up your machine learning application

Mismatched train/test distribution
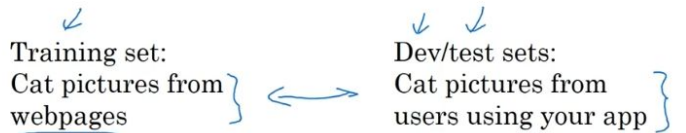
Training vs. Dev/Test

서로 다른 분포에서 나오는 데이터들을 다룸.

Test와 Dev는 같은 분포에서 나오도록 할것.

Training은 상관없나?

Not having a test set might be okay. (only dev set)



Mismatched train/test distribution                    Cats

Training set:                Dev/test sets:
Cat pictures from            Cat pictures from
webpages                     users using your app

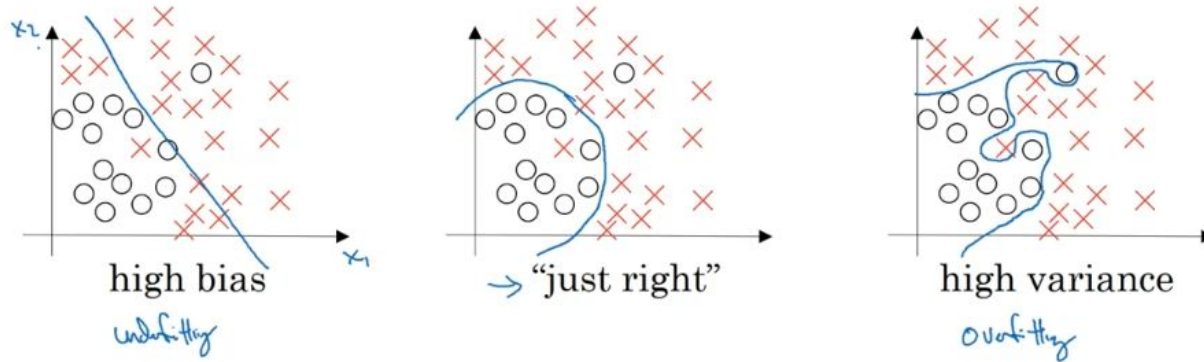→ Make sure dev and test come from same distribution.

train / dev

Not having a test set might be okay. (Only dev set.)

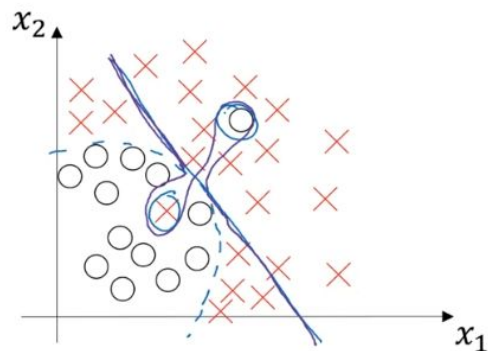# Setting up your machine learning application

Bias/Variance

## Bias and Variance



high bias        "just right"        high variance

underfitting                                 overfitting

# Setting up your machine learning application

Bias/Variance



## Bias and Variance

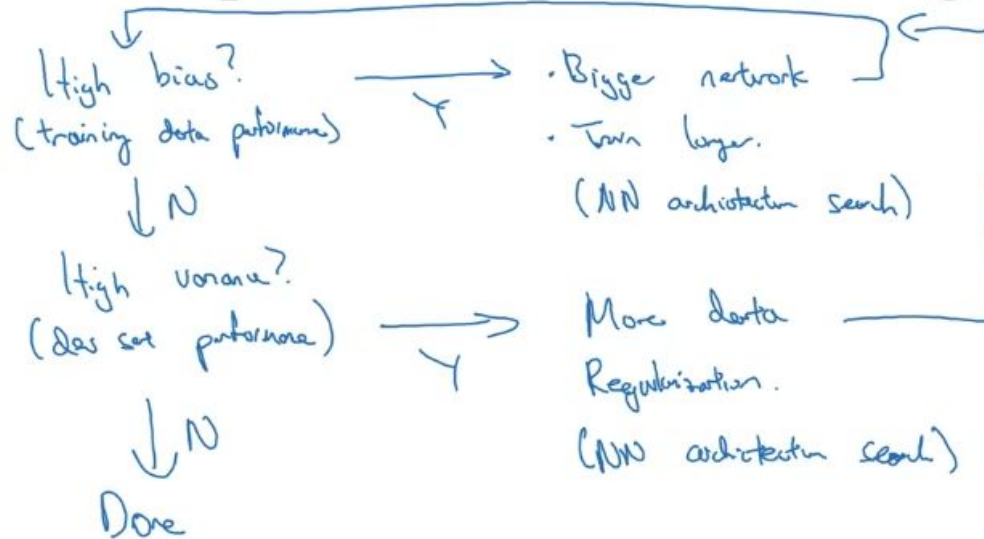Cat classification

y=1    y=0

High bias and high variance

Train set error:    1%    15% ←    15%    0.5%

Dev set error:    11%    16% ←    30%    1%

high variance    high bias    high bias & high varian    low bias low variance

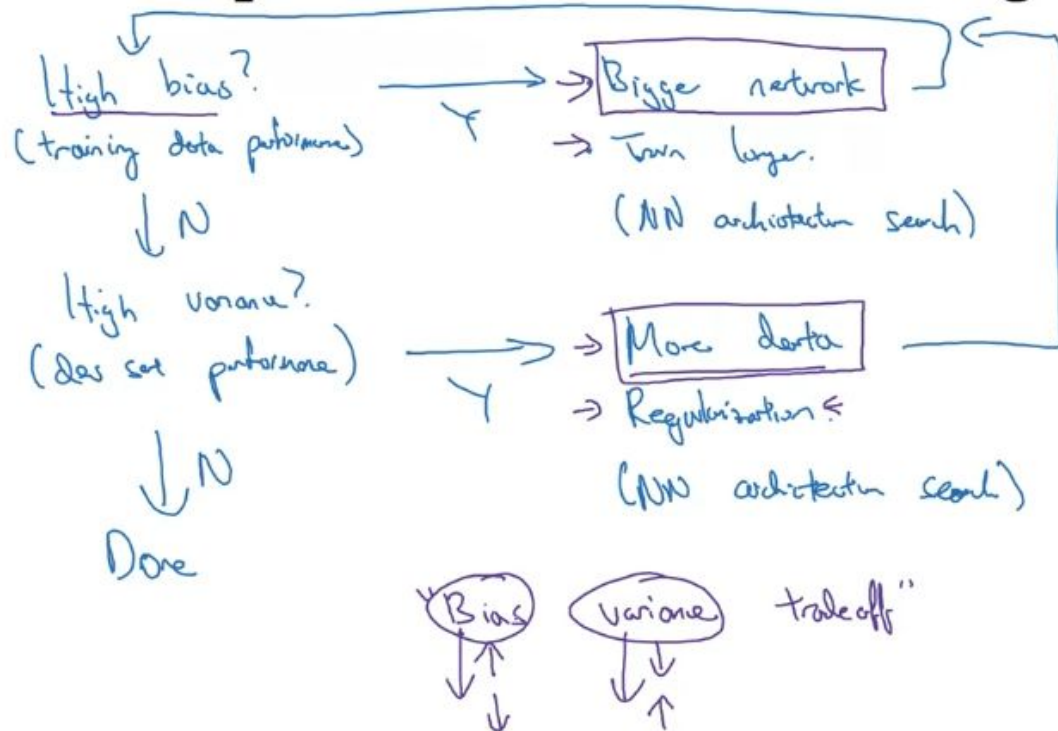Human: ≈0%

Optimal (Bayes) error: ≈ 8% to 15%    Blury images

Andrew Ng

# Setting up your machine learning application



Basic recipe for machine learning

High bias?
(training data performance)
↓ N
High variance?
(dev set performance)
↓ N
Done

• Bigger network
• Train longer
(NN architecture search)

More data
Regularization.
(NN architecture search)

Andrew Ng

# Setting up your machine learning application

## Basic recipe for machine learning



High bias?
(training data performance)

$\downarrow$ N

High variance?
(dev set performance)

$\downarrow$ N

Done

$\rightarrow$ Bigger network
$\rightarrow$ Train longer
(NN architecture search)

$\rightarrow$ More data
$\rightarrow$ Regularization
(NN architecture search)

Bias   Variance   tradeoff"

# Regularizing your neural network

## Logistic regression

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

$$w \in \mathbb{R}^{n_x}, \quad b \in \mathbb{R}$$

$\lambda$ = regularization parameter

lambda         lambd

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

omit

$L_2$ regularization $\quad \|w\|_2^2 = \sum_{j=1}^{n_x} w_j^2 = w^T w \longleftarrow$

$\omega$ will be sparse

$L_1$ regularization $\quad \dfrac{\lambda}{2m} \sum_{j=1}^{n_x} |w_j| = \dfrac{\lambda}{2m} \|w\|_1$

Andrew Ng

# Regularizing your neural network

Neural network

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

"Frobenius norm"

$$\|w^{[l]}\|^2 = \sum_{i=1}^{n^l} \sum_{j=1}^{n^{[l-1]}} (w_{i,j}^{[l]})^2$$

# Regularizing your neural network

Gradient descent?

$$d\omega^{[\ell]} = \boxed{(\text{from backprop}) + \frac{\lambda}{m}\omega^{[\ell]}}$$

$$\rightarrow \omega^{[\ell]} := \omega^{[\ell]} - \alpha\, d\omega^{[\ell]}$$

Before regularization

$$\omega^{[\ell]} := \omega^{[\ell]} - \alpha\left[(\text{from backprop}) + \frac{\lambda}{m}\omega^{[\ell]}\right]$$

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

$$\left(1 - \frac{\alpha\lambda}{m}\right)$$

After regularization

# Regularizing your neural network

Regularization은 왜 오버피팅을 감소시킬까?

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

$$W^{[l]} \approx 0$$

# Regularizing your neural network

Regularization은 왜 오버피팅을 감소시킬까?



tanh

$g(z) := \tanh(z)$

$\lambda \uparrow$    $w^{[l]} \downarrow$    $z^{[l]} := w^{[l]} a^{[l-1]} + b^{[l]}$

# Regularizing your neural network



Dropout regularization

Andrew Ng

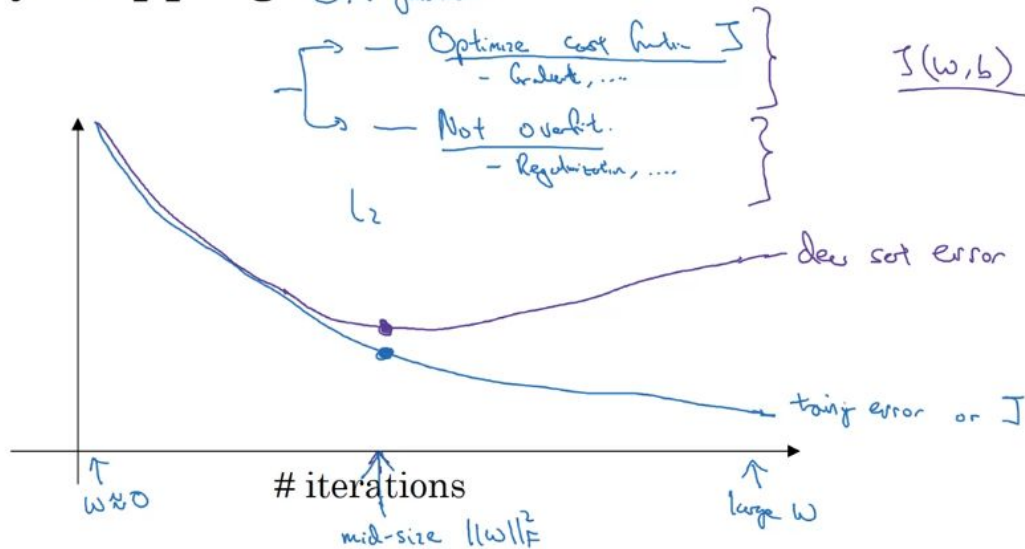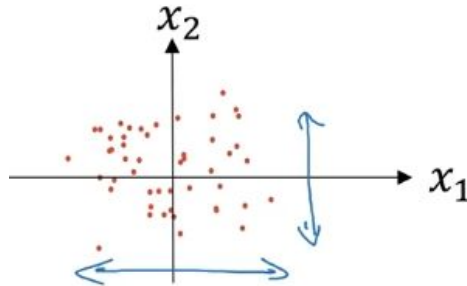# Regularizing your neural network

다른 Regularization 방법들



Data augmentation

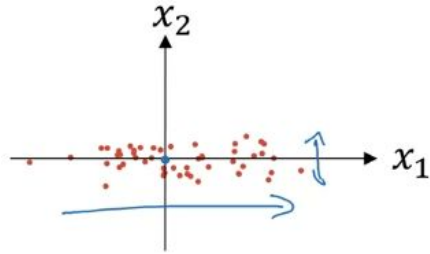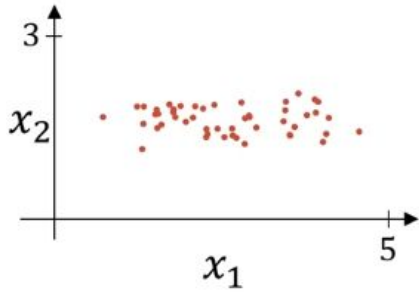# Regularizing your neural network

다른 Regularization 방법들

# Setting up your optimization problem



Subtract mean:

$$\mu = \frac{1}{m} \sum_{i=1}^{m} x^{(i)}$$
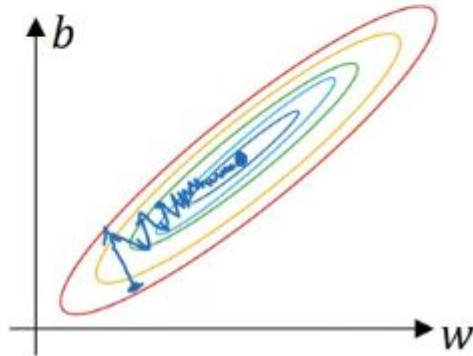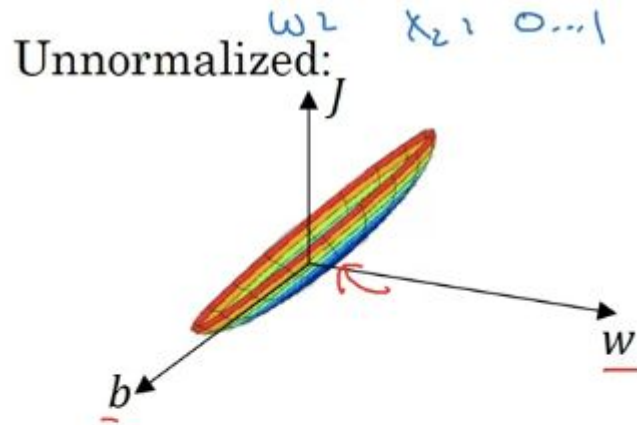
$$x := x - \mu$$

$$\frac{x - \mu}{\sigma}$$

Normalize variance

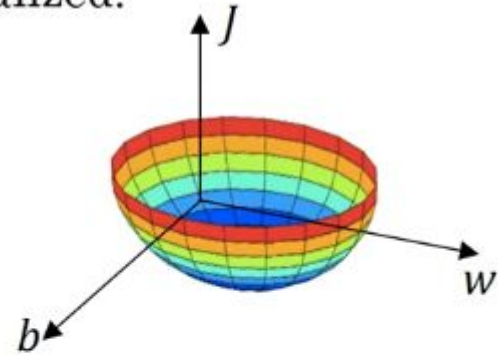$$\sigma^2 = \frac{1}{m} \sum_{i=1}^{m} x^{(i)} ** 2$$

↖ element-wise

# Setting up your optimization problem

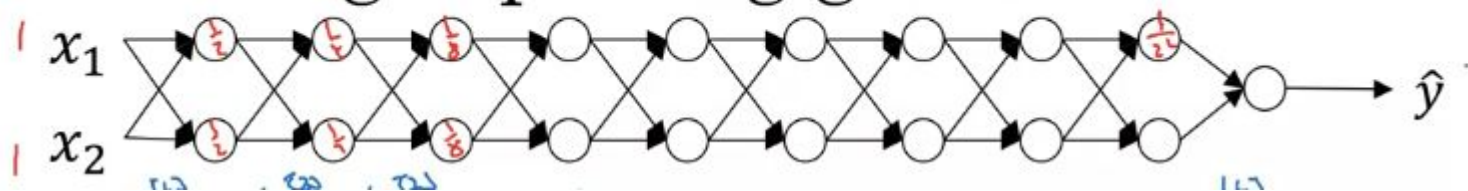# Setting up your optimization problem



Vanishing/exploding gradients

$x_1$, $x_2$ inputs through network to $\hat{y}$

$W^{[n]} > I$  -> W 커지면서 Exploding

$W^{[n]} < I$  -> W가 작아지면서 Vanishing

# Setting up your optimization problem

## Single neuron example

$$a^{[1]}$$

$x_1$

$x_2$ $\qquad w^{[1]}$

$x_3 \longrightarrow \hat{y}$

$x_4 \qquad a = g(z)$

$$z = w_1 x_1 + w_2 x_2 + \cdots + w_n x_n$$

large $n \rightarrow$ smaller $w_i$

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

$$w^{[l]} = \text{np.random.randn(shape)} * \text{np.sqrt}\left(\frac{2}{n^{[l-1]}}\right)$$

ReLU $\qquad g^{[l]}(z) = \text{ReLU}(z)$

Other variants:

tanh $\qquad \sqrt{\dfrac{1}{n^{[l-1]}}}$

Xavier initialization

$$\sqrt{\frac{2}{n^{[l-1]} + n^{[l]}}}$$

# Setting up your optimization problem

## Checking your derivative computation

$$f'(\theta) = \lim_{\varepsilon \to 0} \frac{f(\theta+\varepsilon) - f(\theta-\varepsilon)}{2\varepsilon}$$

Take $W^{[1]}, b^{[1]}, \ldots, W^{[L]}, b^{[L]}$ and reshape into a big vector $\theta$.

Take $dW^{[1]}, db^{[1]}, \ldots, dW^{[L]}, db^{[L]}$ and reshape into a big vector $d\theta$.

Is $d\theta$ the gradient of $J$

# Setting up your optimization problem

for each $i$:

$$d\Theta_{appx}[i] = \frac{J(\Theta_1, \Theta_2, ..., \Theta_i + \varepsilon, ...) - J(\Theta_1, \Theta_2, ..., \Theta_i - \varepsilon, ....)}{2\varepsilon}$$

$$\approx d\Theta[i] = \frac{\partial J}{\partial \Theta_i}$$

$$d\Theta_{appx} \overset{?}{\approx} d\Theta$$

Check

$$\frac{\|d\Theta_{appx} - d\Theta\|_2}{\|d\Theta_{appx}\|_2 + \|d\Theta\|_2}$$

$\varepsilon = 10^{-7}$

$\approx 10^{-7}$ — great!

$10^{-5}$

$\rightarrow 10^{-3}$ — worry.

# Setting up your optimization problem

## Gradient checking implementation notes

- Don't use in training – only to debug

$$d\Theta_{approx}[i] \longleftrightarrow \frac{d\Theta[i]}{}$$

$\uparrow \quad \uparrow \qquad\qquad \uparrow$

- If algorithm fails grad check, look at components to try to identify bug.

$$\underline{db^{[l]}_r} \qquad \underline{dw^{[l]}_r}$$

- Remember regularization.

$$J(\theta) = \frac{1}{m} \sum_i \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_\ell \|w^{[\ell]}\|^2_F$$

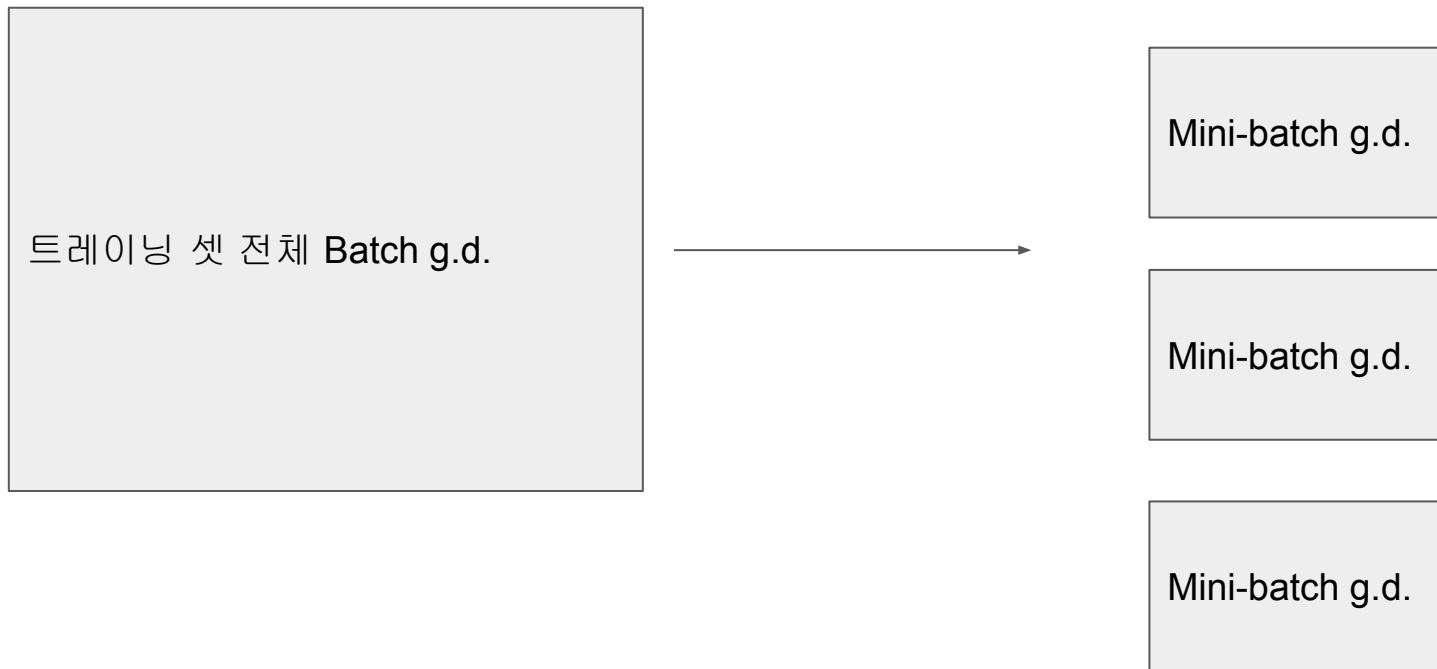$$d\theta = \text{grad of } J \text{ w.r.t. } \Theta$$

- Doesn't work with dropout.

$$J \qquad \underline{keep\text{-}prob = 1.0}$$

- Run at random initialization; perhaps again after some training.

# Optimization algorithms

Mini-batch gradient descent

트레이닝 셋 전체 Batch g.d.

Mini-batch g.d.

Mini-batch g.d.

Mini-batch g.d.

Mini-batch g.d.

# Optimization algorithms

## Mini-batch gradient descent

repeat {

for $t = 1, \cdots, 5000$ {

    Forward prop on $X^{\{t\}}$.

$$Z^{[1]} = W^{[1]} X^{\{t\}} + b^{[1]}$$
$$A^{[1]} = g^{[1]}(Z^{[1]})$$
$$\vdots$$
$$A^{[L]} = g^{[L]}(Z^{[L]})$$

Vectorized implementation (1000 examples)

    Compute cost $J^{\{t\}} = \frac{1}{1000} \sum_{i=1}^{l} \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2 \cdot 1000} \sum_{l} \|W^{[l]}\|_F^2$.

    $\sqrt{}$ from $X^{\{t\}}, Y^{\{t\}}$.

    Backprop to compute gradients wrt $J^{\{t\}}$ (using $(X^{\{t\}}, Y^{\{t\}})$)

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

}

}

"1 epoch" — pass through training set.

1 step of gradient descent

using $X^{\{t\}}, Y^{\{t\}}$.

(as if $m = 1000$)

$X, Y$

Andrew Ng

# Optimization algorithms

## Mini-batch gradient descent

repeat {
for $t = 1, \ldots, 5000$ {

Forward prop on $X^{\{t\}}$.

$$Z^{[1]} = W^{[1]} X^{\{t\}} + b^{[1]}$$
$$A^{[1]} = g^{[1]}(Z^{[1]})$$
$$\vdots$$
$$A^{[L]} = g^{[L]}(Z^{[L]})$$

$\}$ Vec

Compute cost $J^{\{t\}} = \frac{1}{1000} \sum_{i=1}^{\ell} \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2 \cdot 1000} \sum_{\ell} \|W^{[\ell]}\|_F^2$.

Backprop to compute gradients wrt $J^{\{t\}}$ (using $(X^{\{t\}}, Y^{\{t\}})$)

$W^{[\ell]} := W^{[\ell]} - \alpha dW^{[\ell]}$, $b^{[\ell]} := b^{[\ell]} - \alpha db^{[\ell]}$

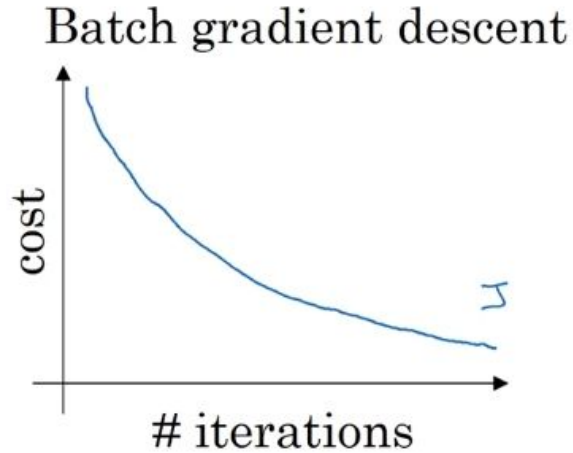}
}

"1 epoch"
↳ pass through training set.

1 step of gradient descent
using $X^{\{t\}}, Y^{\{t\}}$.
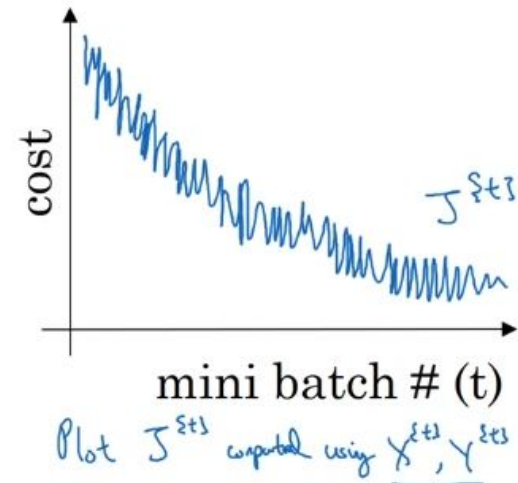(as if $m = 1000$)

$X, Y$

- Epoch means one pass over the full training set
- Batch means that you use all your data to compute the gradient during one iteration.
- Mini-batch means you only take a subset of all your data during one iteration.

Andrew Ng

# Optimization algorithms

## Training with mini batch gradient descent

# Optimization algorithms



— Batch gradient descent
— Mini-batch gradient Descent
— Stochastic gradient descent

Stochastic
gradient
descent

{

Lose speedup
from vectorization

In-between
(mini-batch size
not too big/small)

}

Fastest learning.
• Vectorization.
  (~ 1,000)

Batch
gradient descent
(mini-batch size = $m$)

↓

Too long
per iteration

Typical mini-batch sizes:

$$64 \quad , \quad 128, \quad 256, \quad 512 \qquad \frac{1024}{2^{10}}$$
$$2^6 \qquad 2^7 \qquad 2^8 \qquad 2^9$$

Make sure mini-batch fit in CPU/GPU memory.
$X^{\{t\}}, Y^{\{t\}}$.

# Optimization algorithms

## Exponentially weighted averages

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

$\beta = 0.9$ : $\approx$ 10 days' temperature.

$\beta = 0.98$ : $\approx$ 50 days

$V_t$ as approximately average over

$\approx \frac{1}{1-\beta}$ days' temperature.

$$\frac{1}{1-0.98} = 50$$

$$SMA = \frac{A_1 + A_2 + \ldots + A_n}{n}$$

where:

$A$ = Average in period $n$

$n$ = Number of time periods

# Optimization algorithms

## Exponentially weighted averages

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

$$\beta = 0.9 \quad : \quad \approx 10 \text{ days' temperature}$$

$$\beta = 0.98 \quad : \quad \approx 50 \text{ days}$$

$V_t$ as approximately average over

$$\approx \frac{1}{1-\beta} \text{ days' temperature.}$$

$$\frac{1}{1-0.98} = 50$$

$$SMA = \frac{A_1 + A_2 + \ldots + A_n}{n}$$

**where:**

$A =$ Average in period $n$

$n =$ Number of time periods



temperature

## What Is an Exponential Moving Average (EMA)?

An exponential moving average (EMA) is a type of moving average (MA) that places a greater weight and significance on the most recent data points. The exponential moving average is also referred to as the exponentially weighted moving average.

# Optimization algorithms

## Exponentially weighted averages

$$v_t = \beta v_{t-1} + (1-\beta)\theta_t$$
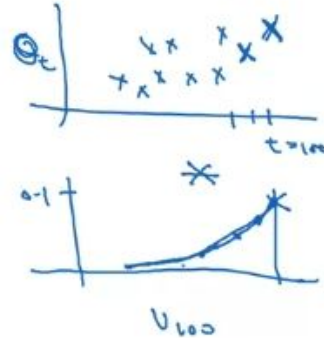
$$v_{100} = 0.9v_{99} + 0.1\theta_{100}$$

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

$$v_{98} = 0.9v_{97} + 0.1\theta_{98}$$
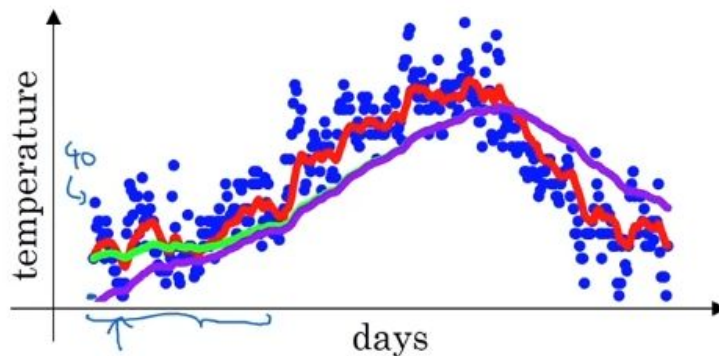
...

$\rightarrow V_{100} = 0.1\,\theta_{100} + 0.9 \times (0.1\,\theta_{99} + 0.9 \times v_{98})$

$0.1\,\theta_{98} + 0.9\,v_{97}$

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

$+ \cdots$

$\theta_t$

$t=100$

$v_{100}$

# Optimization algorithms

## Bias correction



$\beta = 0.98$

$\to v_t = \beta v_{t-1} + (1 - \beta)\theta_t$

$V_0 = 0$

$V_1 = \cancel{0.98 V_0} + 0.02\, \theta_1$

$V_2 = 0.98\, V_1 + 0.02\, \theta_2$

$\quad = 0.98 \times 0.02 \times \theta_1 + 0.02\, \theta_2$

$\quad = 0.0196\, \theta_1 + 0.02\, \theta_2$

$\dfrac{V_t}{1 - \beta^t}$

$t = 2: \quad 1 - \beta^t = 1 - (0.98)^2 = 0.0396$

$\dfrac{V_2}{0.0396} = \dfrac{0.0196\,\theta_1 + 0.02\,\theta_2}{0.0396}$

Andrew Ng

# Optimization algorithms

Gradient descent example



↑ slower learning

↔ faster learning.

Momentum:

On iteration $t$:

  Compute $dW, db$ on current mini-batch.

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

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

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

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

Andrew Ng

# Optimization algorithms

## Implementation details

On iteration $t$:

    Compute $dW, db$ on the current mini-batch

$$v_{dW} = \beta v_{dW} + (1 - \beta)dW$$

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

$$W = W - \alpha v_{dW}, \quad b = b - \alpha v_{db}$$

Hyperparameters: $\alpha, \beta$          $\beta = 0.9$