# Setup and Regularization for training

## Contents

- Setup for training
- Bias & Variance
- Regularization for <span style="color:red">training effectiveness</span>

# Applied ML is a highly iterative process

- Making good choices in how you set up your training, development, and test sets can make a huge difference for you to find a high performance neural network quickly

# Applied ML is a highly iterative process

- When training a neural network, you have to make a lot of decisions:
  - # of layers
  - # of hidden units
  - learning rates
  - activation functions
  - ...

# Applied ML is a highly iterative process

- When training a neural network you have to make a lot of decisions:
  - # of layers
  - # of hidden units
  - learning rates
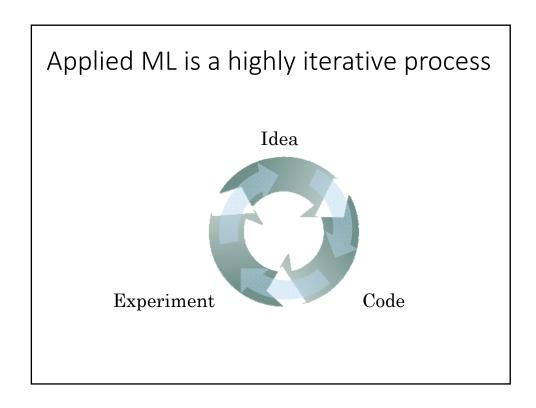  - activation functions
  - …
- It is almost impossible to correctly guess the right values for hyperparameters
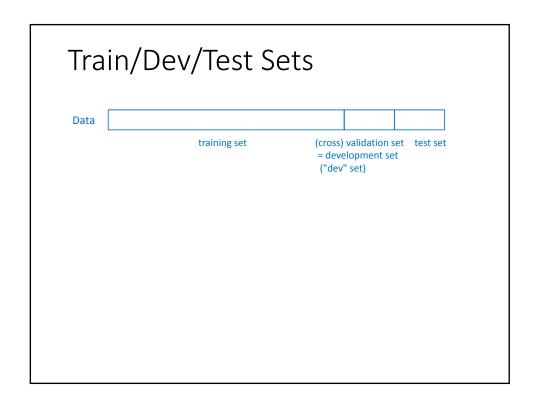
# Applied ML is a highly iterative process

- When training a neural network you have to make a lot of decisions:
  - # of layers
  - # of hidden units
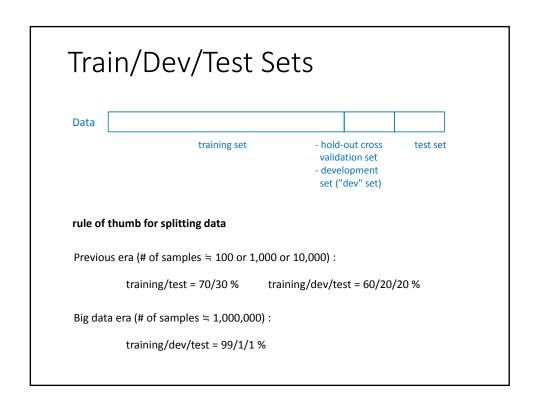  - learning rates
  - activation functions
  - …
- It is almost impossible to correctly guess the right values for hyperparameters
- Therefore, in practice, applied machine learning is <span style="color:red">a highly iterative process</span>

# Applied ML is a highly iterative process

Idea

Experiment · Code

# Train/Dev/Test Sets

Data

training set     (cross) validation set   test set
= development set
("dev" set)

# Train/Dev/Test Sets

Data

training set — - hold-out cross validation set  test set
— - development set ("dev" set)

**rule of thumb for splitting data**

Previous era (# of samples ≒ 100 or 1,000 or 10,000) :

training/test = 70/30 %    training/dev/test = 60/20/20 %

---

# Train/Dev/Test Sets

Data

training set — - hold-out cross validation set  test set
— - development set ("dev" set)

**rule of thumb for splitting data**

Previous era (# of samples ≒ 100 or 1,000 or 10,000) :

training/test = 70/30 %    training/dev/test = 60/20/20 %

Big data era (# of samples ≒ 1,000,000) :

training/dev/test = 99/1/1 %

# Mismatched train/test distribution

Goal: Cat Detector

Training set:
(high-density) cat pictures
from webpages

Dev/test sets:
(low-density) cat pictures
from users using your app

---

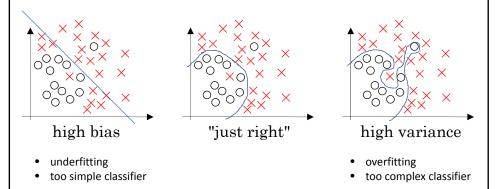# Mismatched train/test distribution

Goal: Cat Detector

Training set:
(high-density) cat pictures
from webpages

Dev/test sets:
(low-density) cat pictures
from users using your app

**rule of thumb for splitting data:**

- Make sure dev and test sets come from same distribution

- Even Not having a test set might be okay (only having dev set) for better training

# Bias and Variance

- Bias and variance tradeoff is a very important problem in traditional machine learning
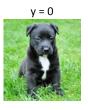- But, in deep learning era, it's importance has been somewhat reduced



| high bias | "just right" | high variance |

- underfitting
- too simple classifier

- overfitting
- too complex classifier

---

# Bias and Variance

Cat classification

y = 1    y = 0



| Train set error: | 1% | 15% | 15% | 0.5% |
|---|---|---|---|---|
| Dev set error: | 11% | 16% | 30% | 1% |
| | high variance | high bias | high bias & high variance | low bias & low variance |

Human's classification error ≈ 0%

# Bias and Variance

Cat classification

y = 1    y = 0

| Train set error: | 1% | 15% | 15% | 0.5% |
|---|---|---|---|---|
| Dev set error: | 11% | 16% | 30% | 1% |
| | high variance | ~~high bias~~ | high bias & high variance | low bias & low variance |

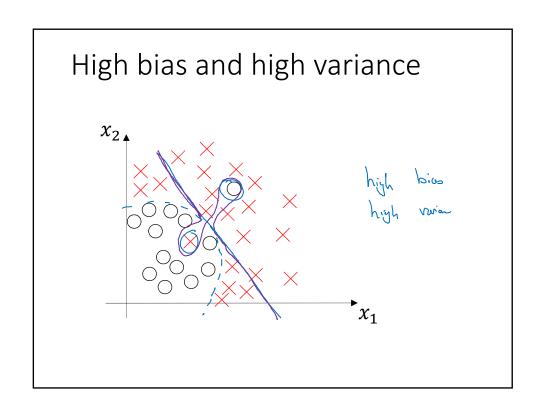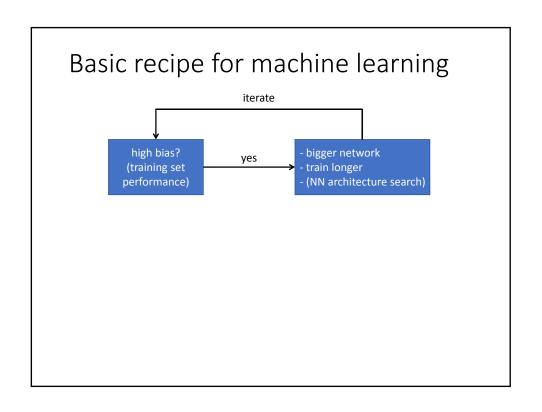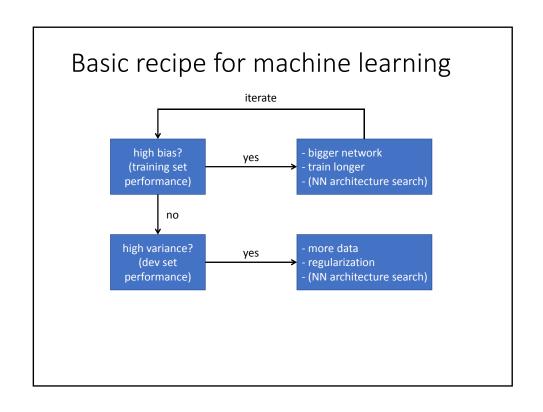- If human's error ≈ 15%, then the second case does not have the high bias problem

---

# Bias and Variance
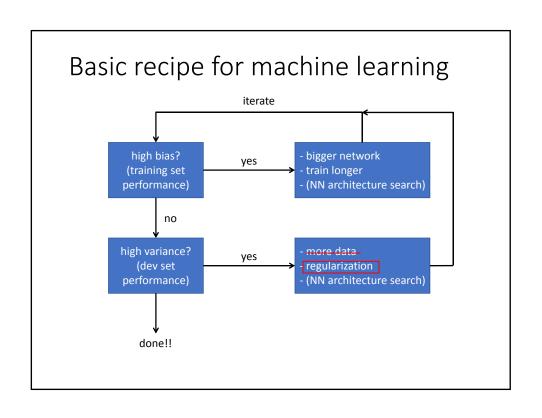
Cat classification

y = 1    y = 0

**bias**

| Train set error: | 1% | 15% | 15% | 0.5% |
|---|---|---|---|---|
| Dev set error: | 11% | 16% | 30% | 1% |
| **variance** = DS error – TS error | high variance | high bias | high bias & high variance | low bias & low variance |

Human's classification error ≈ 0%

# High bias and high variance



# Basic recipe for machine learning

# Basic recipe for machine learning

iterate

| high bias?<br>(training set<br>performance) | →yes→ | - bigger network<br>- train longer<br>- (NN architecture search) |

no

| high variance?<br>(dev set<br>performance) | →yes→ | - more data<br>- regularization<br>- (NN architecture search) |

# Basic recipe for machine learning

iterate

| high bias?<br>(training set<br>performance) | →yes→ | - bigger network<br>- train longer<br>- (NN architecture search) |

no

| high variance?<br>(dev set<br>performance) | →yes→ | - more data<br>- regularization<br>- (NN architecture search) |

done!!

# Regularization

- Regularization helps to prevent overfitting, or to reduce the variance of your network

# Regularization for Logistic Regression

$$\min_{w,b} J(w, b) \qquad w \in \mathbb{R}^n, b \in \mathbb{R}$$

$$J(w, b) = \frac{1}{m} \sum_{i=1}^{m} \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

# Regularization for Logistic Regression

$$\min_{w,b} J(w,b) \qquad w \in \mathbb{R}^n, b \in \mathbb{R}$$

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

usually omitted

$$\text{L}_2 \text{ regularization:} \quad \|w\|_2^2 = \sum_{j=1}^{n} w_j^2 = w^T w$$

---

# Regularization for Logistic Regression

$$\min_{w,b} J(w,b) \qquad w \in \mathbb{R}^n, b \in \mathbb{R}$$

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

usually omitted

$$\text{L}_2 \text{ regularization:} \quad \|w\|_2^2 = \sum_{j=1}^{n} w_j^2 = w^T w$$

$$\text{L}_1 \text{ regularization:} \quad \|w\|_1 = \sum_{j=1}^{n} |w_j| \quad \longrightarrow \quad \text{sparse w}$$

# Regularization for Logistic Regression

$$\min_{w,b} J(w,b) \qquad w \in \mathbb{R}^n, b \in \mathbb{R} \qquad \text{$\lambda$: regularization parameter}$$

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

usually omitted

$$\text{L}_2 \text{ regularization: } \quad \|w\|_2^2 = \sum_{j=1}^{n} w_j^2 = w^T w$$

$$\text{L}_1 \text{ regularization: } \quad \|w\|_1 = \sum_{j=1}^{n} |w_j| \quad \longrightarrow \text{sparse w}$$

# Regularization for Neural Network

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

$$\text{Frobenius norm: } \quad \left\|W^{[l]}\right\|_F^2 = \sum_{i=1}^{n^{[l]}}\sum_{j=1}^{n^{[l-1]}}\left(w_{ij}^{[l]}\right)^2 \qquad W^{[l]}: \left(n^{[l]}, n^{[l-1]}\right)$$

# Regularization for Neural Network

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

Frobenius norm: $\qquad \big\|W^{[l]}\big\|_F^2 = \displaystyle\sum_{i=1}^{n^{[l]}}\sum_{j=1}^{n^{[l-1]}}\big(w_{ij}^{[l]}\big)^2 \qquad W^{[l]}: \big(n^{[l]}, n^{[l-1]}\big)$

Gradient descent: $\quad dW^{[l]} = \text{(from backprop)}$

$\qquad\qquad\qquad W^{[l]} := W^{[l]} - \alpha \cdot dW^{[l]}$

$$\left(dW^{[l]} = \frac{\partial J}{\partial W^{[l]}}\right)$$

---

# Regularization for Neural Network

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

Frobenius norm: $\qquad \big\|W^{[l]}\big\|_F^2 = \displaystyle\sum_{i=1}^{n^{[l]}}\sum_{j=1}^{n^{[l-1]}}\big(w_{ij}^{[l]}\big)^2 \qquad W^{[l]}: \big(n^{[l]}, n^{[l-1]}\big)$

Gradient descent: $\quad dW^{[l]} = \text{(from backprop)} + \dfrac{\lambda}{m}W^{[l]} \quad \left(\because \dfrac{\partial\|W^{[l]}\|_F^2}{\partial W^{[l]}} = 2W^{[l]}\right)$

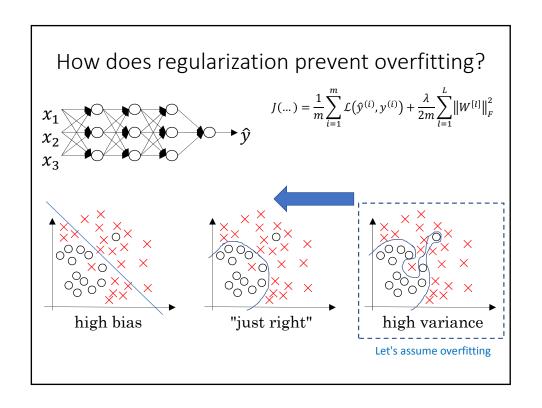$\qquad\qquad\qquad W^{[l]} := W^{[l]} - \alpha \cdot dW^{[l]}$
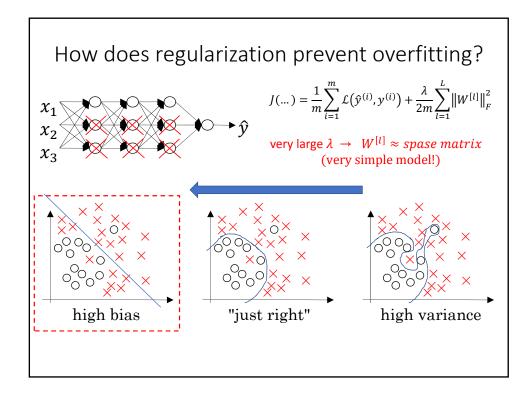
# Regularization for Neural Network

$$J\left(W^{[1]}, b^{[1]}, \ldots, W^{[L]}, b^{[L]}\right) = \frac{1}{m} \sum_{i=1}^{m} \mathcal{L}\left(\hat{y}^{(i)}, y^{(i)}\right) + \frac{\lambda}{2m} \sum_{l=1}^{L} \left\| W^{[l]} \right\|_F^2$$
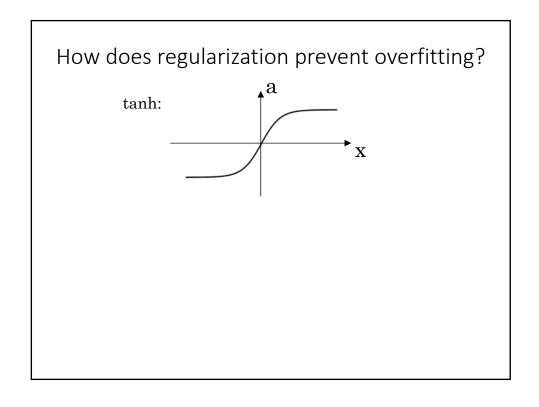
Frobenius norm: $\quad \left\| W^{[l]} \right\|_F^2 = \sum_{i=1}^{n^{[l]}} \sum_{j=1}^{n^{[l-1]}} \left( w_{ij}^{[l]} \right)^2 \quad\quad W^{[l]} : \left( n^{[l]}, n^{[l-1]} \right)$

Gradient descent: $\quad dW^{[l]} = \text{(from backprop)} + \dfrac{\lambda}{m} W^{[l]} \quad \left( \because \dfrac{\partial \left\| W^{[l]} \right\|_F^2}{\partial W^{[l]}} = 2W^{[l]} \right)$

$\qquad\qquad\qquad\quad W^{[l]} := W^{[l]} - \alpha \cdot dW^{[l]} \quad \longleftarrow$

"Weight decay": $\quad W^{[l]} := W^{[l]} - \alpha \cdot \left[ \text{(from backprop)} + \dfrac{\lambda}{m} W^{[l]} \right]$

$\qquad\qquad\qquad\quad = W^{[l]} - \dfrac{\alpha\lambda}{m} W^{[l]} - \alpha \text{(from backprop)}$

$\qquad\qquad\qquad\quad = \left( 1 - \dfrac{\alpha\lambda}{m} \right) W^{[l]} - \alpha \text{(from backprop)}$

---

# Regularization for Neural Network

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

Frobenius norm: $\quad \left\| W^{[l]} \right\|_F^2 = \sum_{i=1}^{n^{[l]}} \sum_{j=1}^{n^{[l-1]}} \left( w_{ij}^{[l]} \right)^2 \quad\quad W^{[l]} : \left( n^{[l]}, n^{[l-1]} \right)$

Gradient descent: $\quad dW^{[l]} = \text{(from backprop)} + \dfrac{\lambda}{m} W^{[l]} \quad \left( \because \dfrac{\partial \left\| W^{[l]} \right\|_F^2}{\partial W^{[l]}} = 2W^{[l]} \right)$

$\qquad\qquad\qquad\quad W^{[l]} := W^{[l]} - \alpha \cdot dW^{[l]} \quad \longleftarrow$

"Weight decay": $\quad W^{[l]} := W^{[l]} - \alpha \cdot \left[ \text{(from backprop)} + \dfrac{\lambda}{m} W^{[l]} \right]$

$\qquad\qquad\qquad\quad = W^{[l]} - \dfrac{\alpha\lambda}{m} W^{[l]} - \alpha \text{(from backprop)}$

$\qquad\qquad\qquad\quad = \left( 1 - \dfrac{\alpha\lambda}{m} \right) W^{[l]} - \alpha \text{(from backprop)}$

$\qquad\qquad\qquad\qquad\qquad\quad {\scriptstyle <1}$

# How does regularization prevent overfitting?

$$J(\dots) = \frac{1}{m}\sum_{i=1}^{m}\mathcal{L}\big(\hat{y}^{(i)}, y^{(i)}\big)$$

high bias  "just right"  high variance

Let's assume overfitting

# How does regularization prevent overfitting?

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

high bias  "just right"  high variance

Let's assume overfitting

# How does regularization prevent overfitting?

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

very large $\lambda \rightarrow W^{[l]} \approx spase\ matrix$
(very simple model!)

high bias     "just right"     high variance



# How does regularization prevent overfitting?

tanh:

## How does regularization prevent overfitting?

tanh:



almost
linear

$\lambda \uparrow \Rightarrow |W^{[l]}{}_{ij}| \downarrow \Rightarrow |z^{[l]}| \downarrow \quad (\because z^{[l]} = W^{[l]}a^{[l-1]} + b^{[l]})$

$\Rightarrow$ every layers $\approx$ get close to linear $\approx$ get simple

---

# Dropout Regularization



$x_1$

$x_2$

$x_3$

$x_4$

$\hat{y}$

# Dropout Regularization

$x_1$
$x_2$
$x_3$
$x_4$

$\hat{y}$

keeping prob = 0.5



# Dropout Regularization

$x_1$
$x_2$
$x_3$
$x_4$

$\hat{y}$

keeping prob = 0.5

# Dropout Regularization



keeping prob = 0.5

# Dropout Regularization



keeping prob = 0.5

# Why does dropout work?

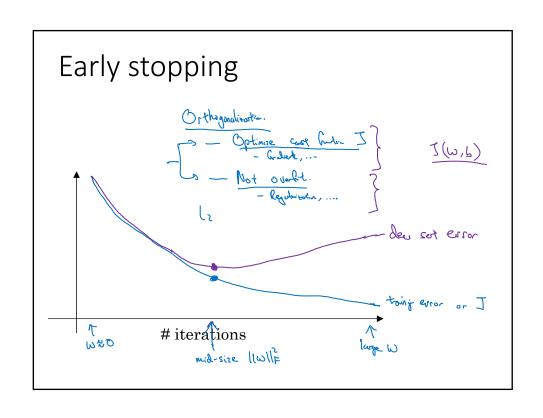Intuition: Can't rely on any one feature, so have to spread out weights.
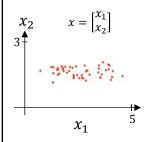


---

# Why does dropout work?

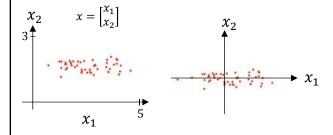Intuition: Can't rely on any one feature, so have to spread out weights.



$x_1$

$x_2$

$x_3$

$\hat{y}$

$W^{[1]}$  $W^{[2]}$  $W^{[3]}$

0.6  0.6  0.8  1.0  1.0

# Data augmentation



# Early stopping



Orthogonalization.

↳ — Optimize cost function $J$
  — Gradient, ...

↳ — Not overfit.
  — Regularization, ...

$L_2$

$J(w,b)$

dev set error

traing error  or  $J$

↑
$w \approx 0$

# iterations

↑
mid-size $\|w\|_F^2$

↑
large $w$

# Normalizing training sets

$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$



# Normalizing training sets

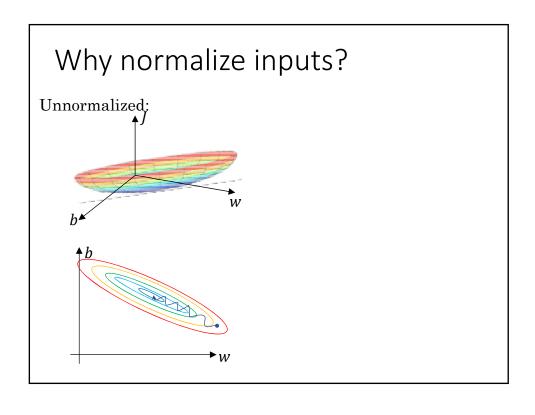$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$



1. subtract mean:

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

$$x^{(i)} := x^{(i)} - \mu$$

# Normalizing training sets

$x_2$   $x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$   $x_2$   $x_2$

3

$x_1$   5   $x_1$   $x_1$

1. subtract mean:

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

$$x^{(i)} := x^{(i)} - \mu$$

2. normalize variance:

$$\sigma_j^2 = \frac{1}{m}\sum_{i=1}^{m} \left\{ x_j^{(i)} \right\}^2$$

$$x_j^{(i)} := x_j^{(i)} / \sigma_j$$

# Normalizing training sets

$x_2$   $x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$   $x_2$   $x_2$

3

$x_1$   5   $x_1$   $x_1$

1. subtract mean:

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

$$x^{(i)} := x^{(i)} - \mu$$

2. normalize variance:

$$\sigma_j^2 = \frac{1}{m}\sum_{i=1}^{m} \left\{ x_j^{(i)} \right\}^2$$

$$x_j^{(i)} := x_j^{(i)} / \sigma_j$$

# Why normalize inputs?

Unnormalized:

---

# Why normalize inputs?

Unnormalized:

# Why normalize inputs?

Unnormalized:

Normalized:

# Why normalize inputs?

Unnormalized:

Normalized:

# Vanishing/exploding gradients

$x_1$

$x_2$

$W^{[1]}$ $\quad W^{[2]}$ $\quad W^{[3]}$ $\quad \dots$ $\qquad\qquad\qquad\qquad\qquad W^{[L]}$

$\hat{y}$

$g(z) = z \qquad b^{[l]} = 0$

$$\hat{y} = W^{[L]}W^{[L-1]}W^{[L-2]} \dots W^{[3]}W^{[2]}W^{[1]}x$$
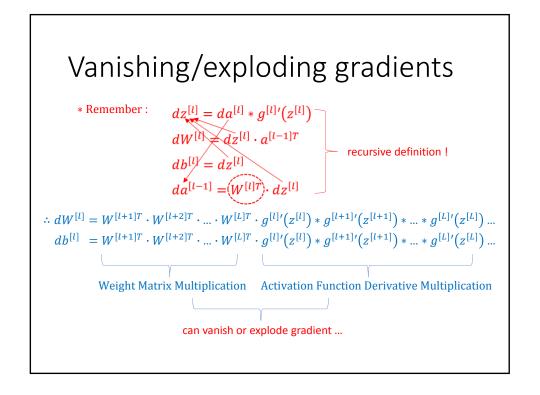
$$W^{[l]} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix} \implies \hat{y} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}^{L} x = \begin{bmatrix} 1.5^{L} & 0 \\ 0 & 1.5^{L} \end{bmatrix} x \implies$$

activations
increase
exponentially!!
"Explode"

# Vanishing/exploding gradients

$x_1$
$x_2$
$W^{[1]}$   $W^{[2]}$   $W^{[3]}$   ...           $W^{[L]}$    $\hat{y}$

$$g(z) = z \qquad b^{[l]} = 0$$

$$\hat{y} = W^{[L]}W^{[L-1]}W^{[L-2]} \dots W^{[3]}W^{[2]}W^{[1]}x$$

$W^{[l]} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}$ ➡ $\hat{y} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}^L x = \begin{bmatrix} 1.5^L & 0 \\ 0 & 1.5^L \end{bmatrix}x$ ➡ activations increase exponentially!!
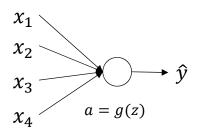
$W^{[l]} = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}$ ➡ $\hat{y} = \begin{bmatrix} 0.5 & 0 \\ 0 & 0.5 \end{bmatrix}^L x = \begin{bmatrix} 0.5^L & 0 \\ 0 & 0.5^L \end{bmatrix}x$ ➡ activations decrease exponentially!! "Vanish"

---

# Vanishing/exploding gradients

$*$ Remember :

$$dz^{[l]} = da^{[l]} * g^{[l]\prime}(z^{[l]})$$
$$dW^{[l]} = dz^{[l]} \cdot a^{[l-1]T}$$
$$db^{[l]} = dz^{[l]}$$
$$da^{[l-1]} = W^{[l]T} \cdot dz^{[l]}$$

recursive definition !

$$\therefore dW^{[l]} = W^{[l+1]T} \cdot W^{[l+2]T} \cdot \dots \cdot W^{[L]T} \cdot g^{[l]\prime}(z^{[l]}) * g^{[l+1]\prime}(z^{[l+1]}) * \dots * g^{[L]\prime}(z^{[L]}) \dots$$
$$db^{[l]} = W^{[l+1]T} \cdot W^{[l+2]T} \cdot \dots \cdot W^{[L]T} \cdot g^{[l]\prime}(z^{[l]}) * g^{[l+1]\prime}(z^{[l+1]}) * \dots * g^{[L]\prime}(z^{[L]}) \dots$$

Weight Matrix Multiplication     Activation Function Derivative Multiplication

can vanish or explode gradient …

# Weight Variance initialization for deep network
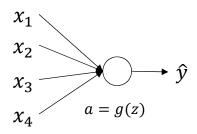
$x_1$
$x_2$
$x_3$
$x_4$

$\hat{y}$

$a = g(z)$

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

if large $n \rightarrow$ smaller variance of $w_i$

$\therefore \text{Var}(w_i) = \frac{1}{n}$ when b exists for general case
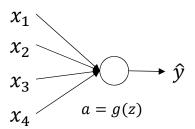
# Weight Variance initialization for deep network

$x_1$
$x_2$
$x_3$
$x_4$

$\hat{y}$

$a = g(z)$

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

if large $n \rightarrow$ smaller variance of $w_i$

$\therefore \text{Var}(w_i) = \frac{2}{n}$ for ReLU activation, which has been claimed by some experiments

# Weight Variance initialization for deep network

$x_1$
$x_2$
$x_3$
$x_4$

$a = g(z)$

$\hat{y}$

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

if large $n \rightarrow$ smaller variance of $w_i$

$\therefore \text{Var}(w_i) = \dfrac{2}{n}$ for ReLU activation

Other variants for tanh activation

$\text{Var}(w_i) = \dfrac{1}{n^{[l-1]}}$

: Xavier initialization

$\text{Var}(w_i) = \dfrac{2}{n^{[l-1]} + n^{[l]}}$

: proposed by Bengio