



deeplearning.ai

Setting up your ML application

Train/dev/test sets

Applied ML is a highly iterative process

layers

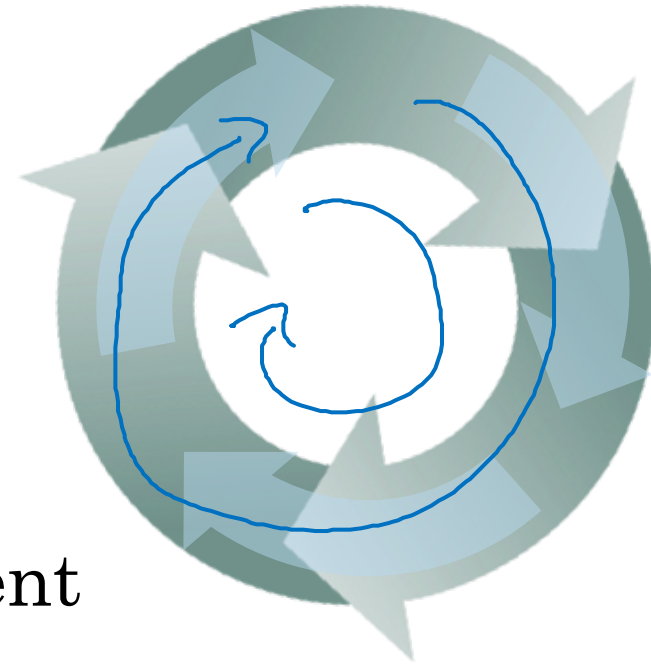
hidden units

learning rates

activation functions

...

Idea



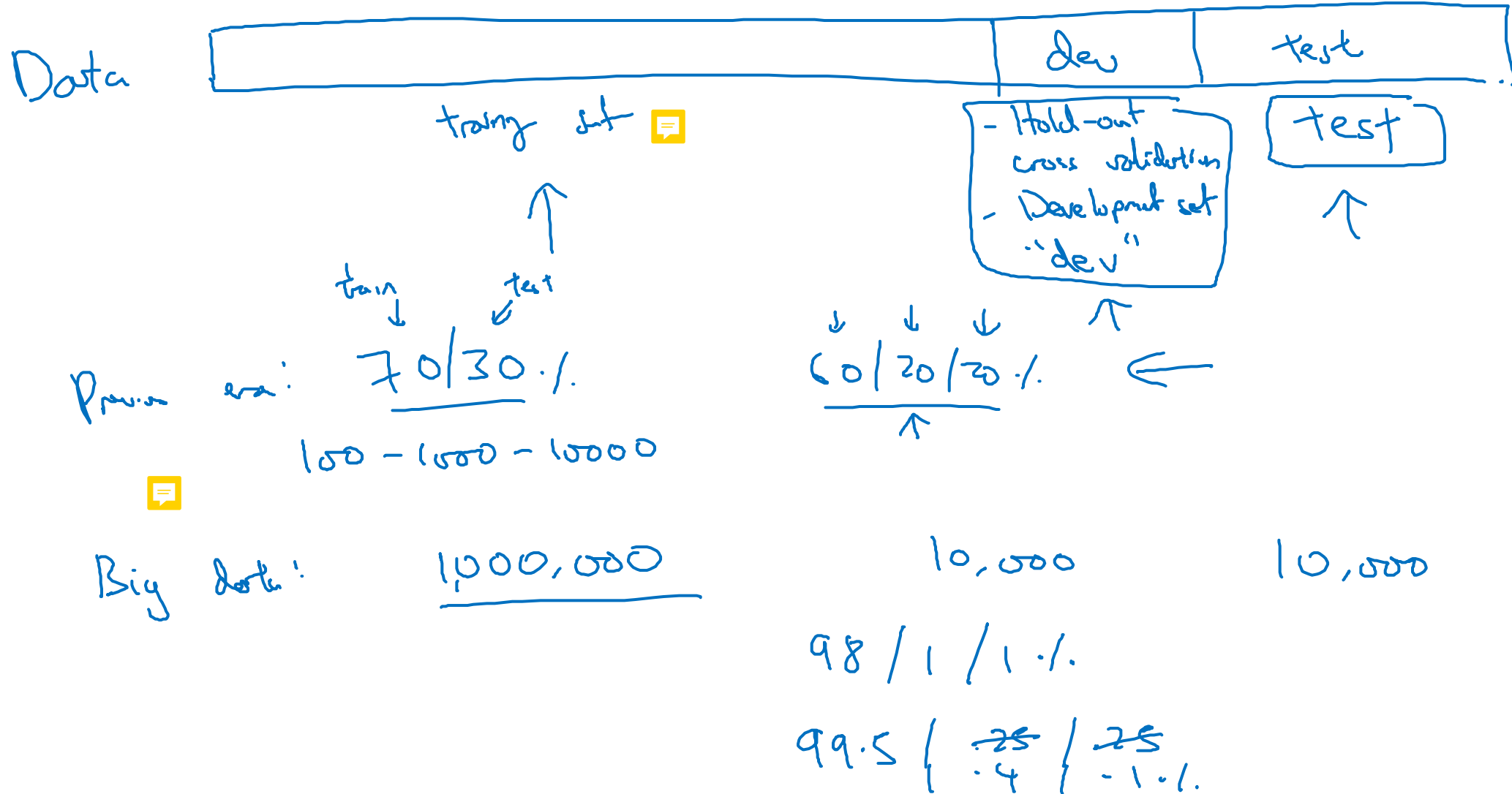
Experiment

Code

NLP, Vision, Speech, Structured Data

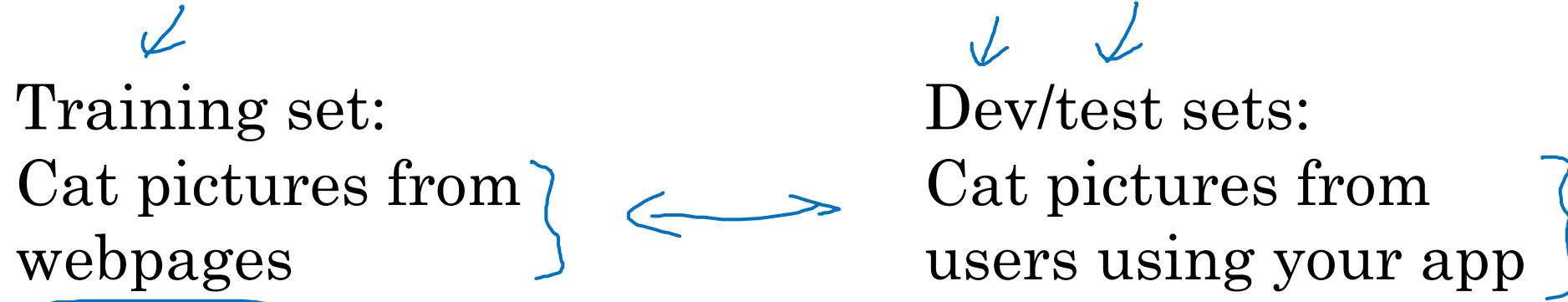
```
graph TD; NLP[NLP] --> Speech[Speech]; StructuredData[Structured Data] --> Search[Search]; StructuredData --> Security[Security]; Speech --> Ads[Ads]; Security --> Logistic[Logistic ...]
```

Train/dev/test sets

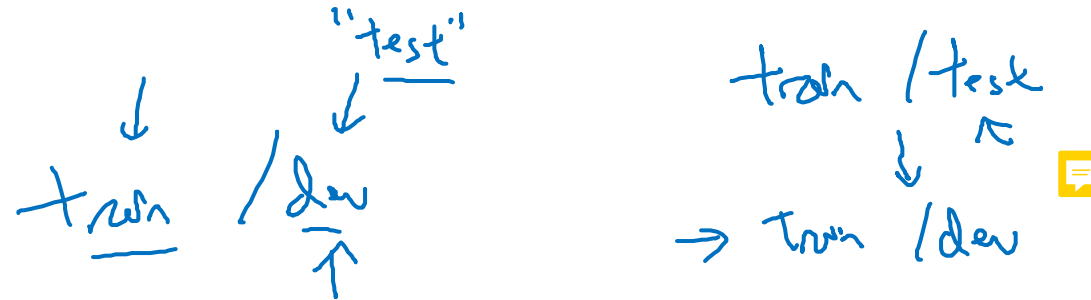


Mismatched train/test distribution

Certs



→ Make sure dev and test come from same distribution.



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

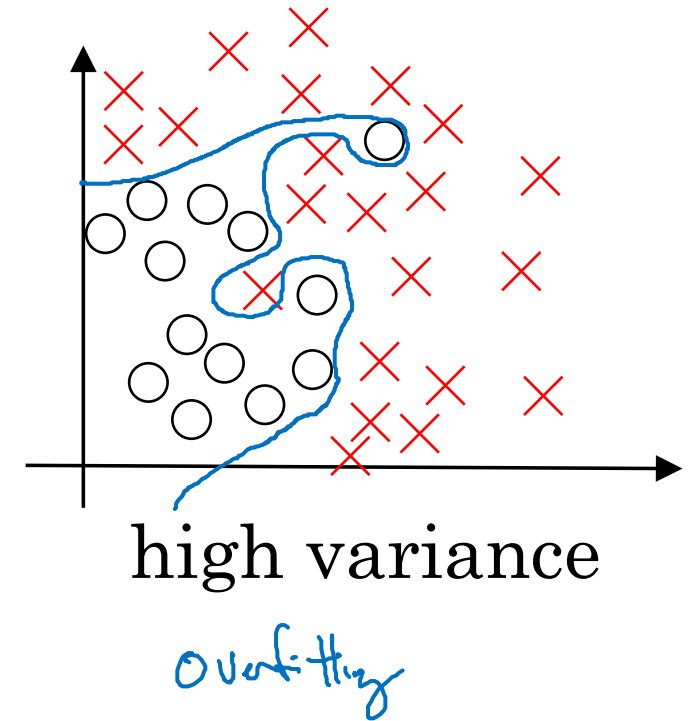
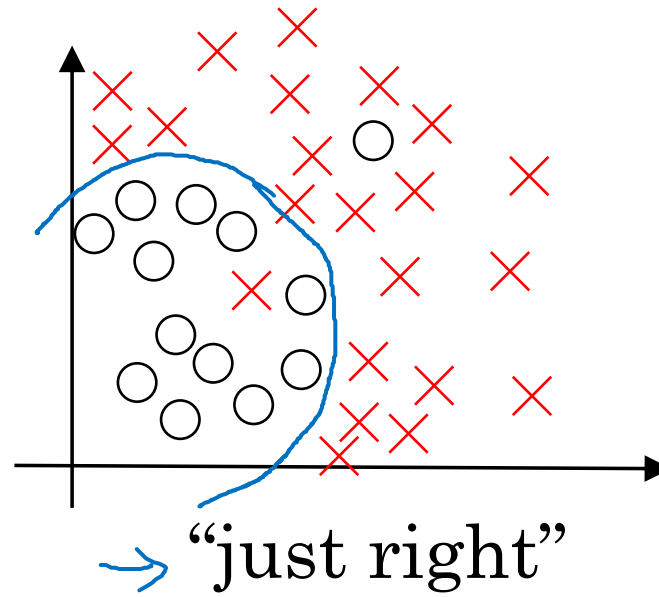
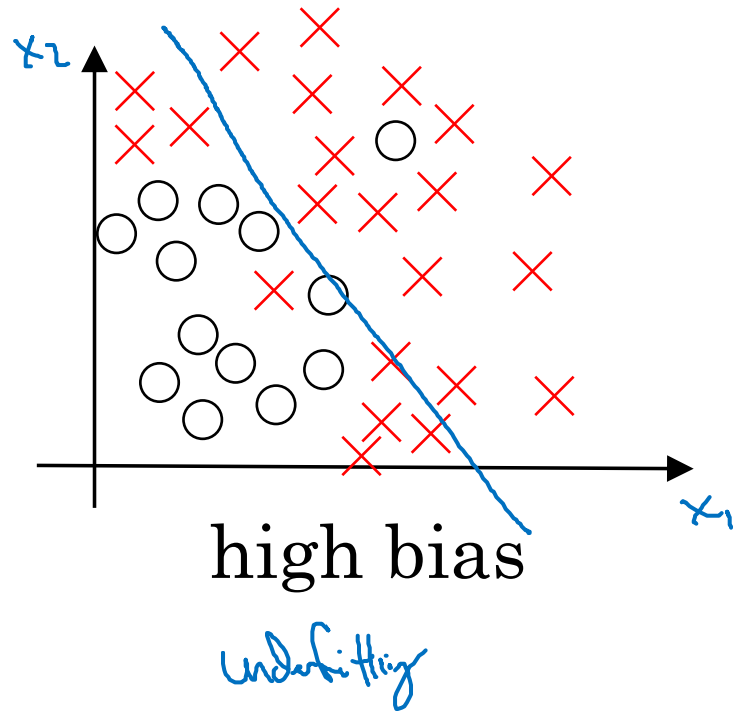


deeplearning.ai

Setting up your ML application

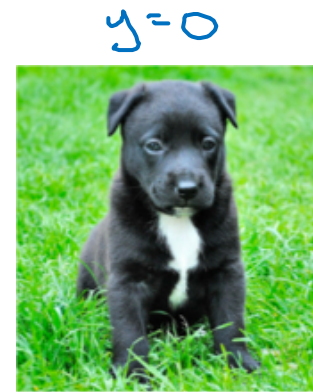
Bias/Variance

Bias and Variance



Bias and Variance

Cat classification



Train set error:

1%



Dev set error:

11%

high variance
↑

15% ←

16% ←

high bias
↑



15%

30%

high bias
& high variance

0.5%

1%

low bias
low variance
↑

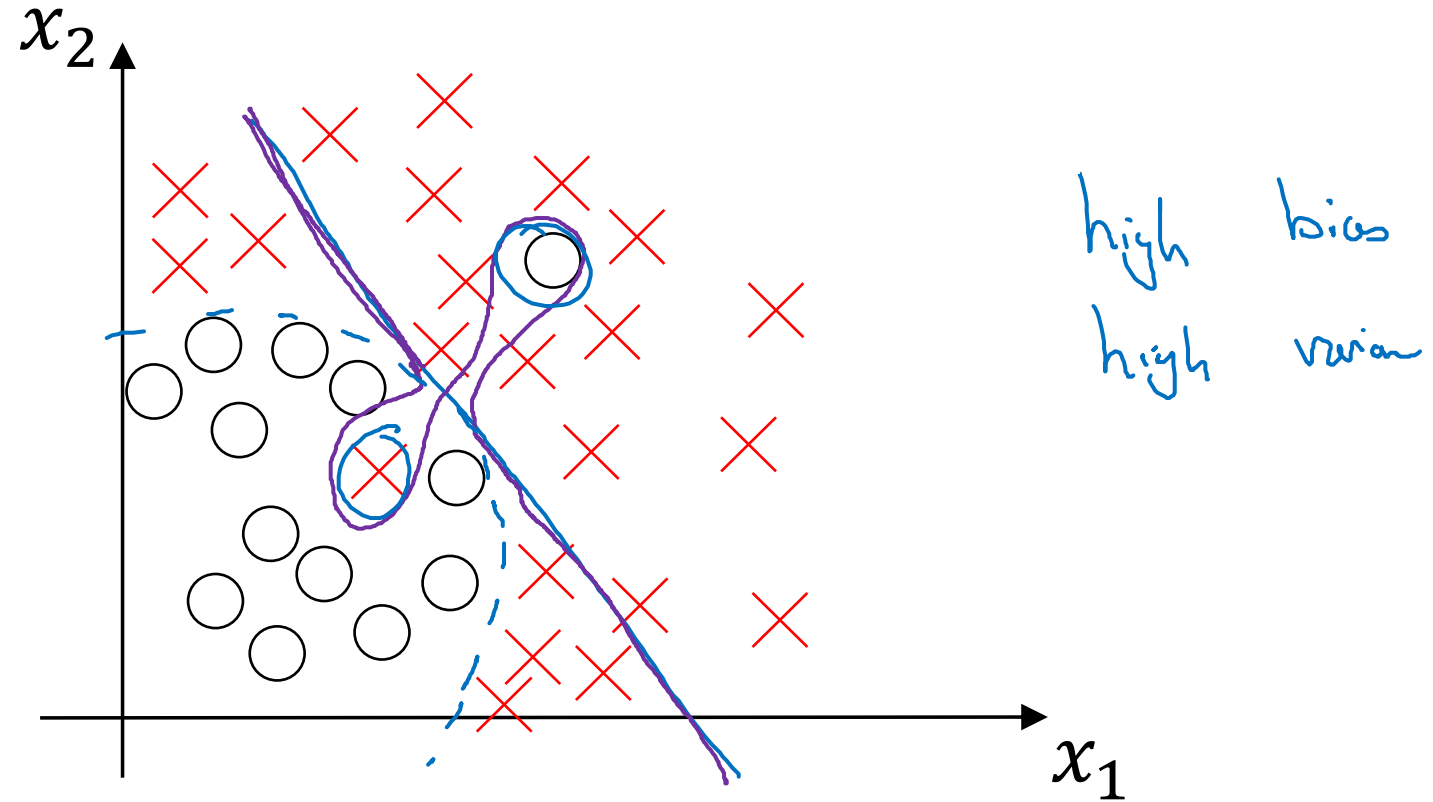


Human: ~0%

Optimal (Bayes) error: ~~~0%~~ 15%

Blurry images

High bias and high variance



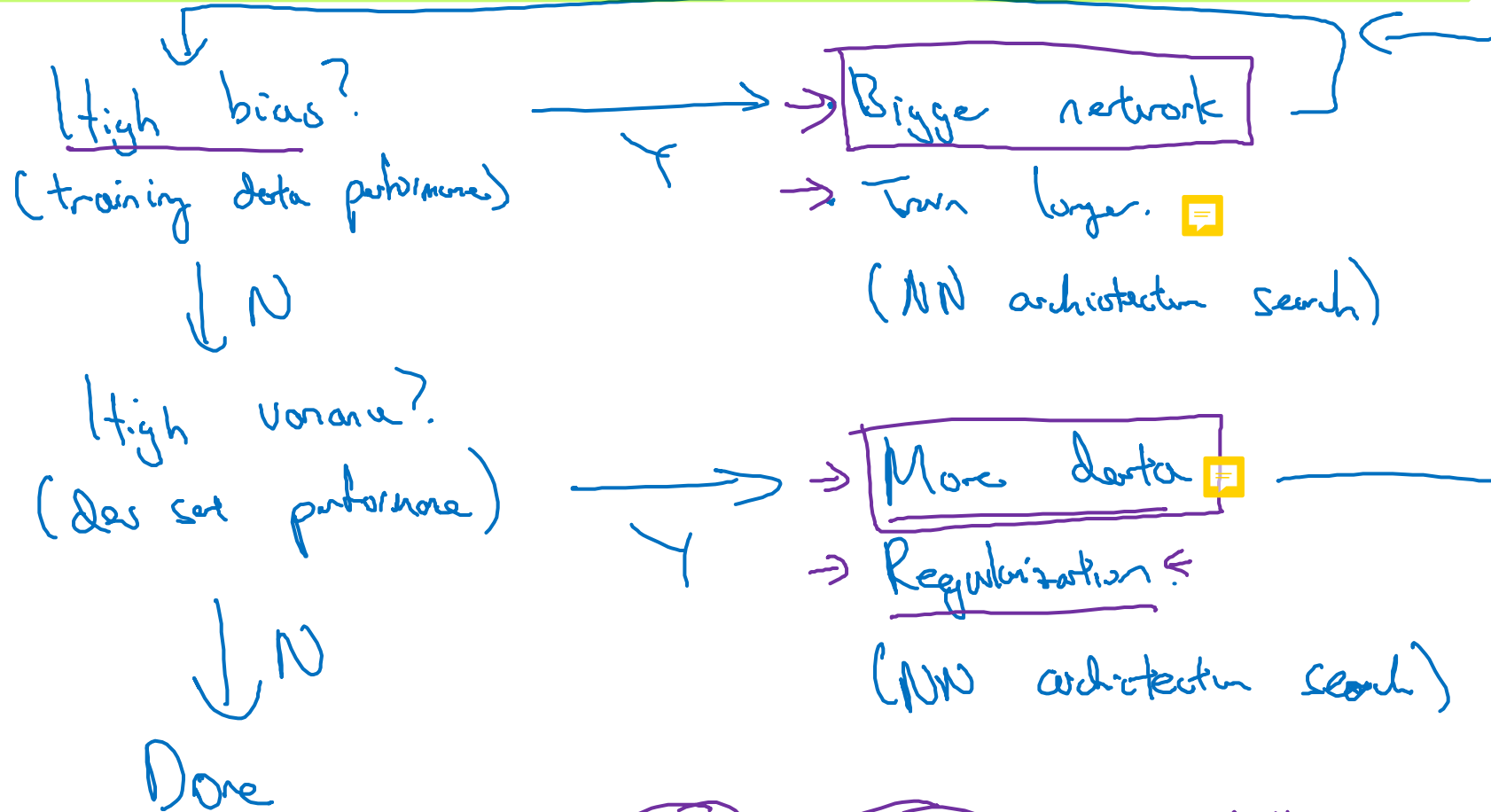


deeplearning.ai

Setting up your ML application

Basic “recipe” for machine learning

Basic recipe for machine learning





deeplearning.ai

Regularizing your neural network

Regularization ?



To reduce variance or prevent overfitting in NN

Logistic regression

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

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

λ = regularization parameter
lambda

$$J(w,b) = \underbrace{\frac{1}{m} \sum_{i=1}^m \ell(y^{(i)}, \hat{y}^{(i)})}_{\text{cost function}} + \frac{\lambda}{2m} \underbrace{\|w\|_2^2}_{\text{L2 regularization}}$$

~~$+\frac{\lambda}{2m} b^2$~~
omit

L_2 regularization $\underbrace{\|w\|_2^2}_{\text{L2 regularization}} = \sum_{j=1}^{n_x} w_j^2 = w^T w \leftarrow$

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

w will be sparse

Neural network

$$\rightarrow J(w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}) = \underbrace{\frac{1}{n} \sum_{i=1}^n \ell(y^{(i)}, \hat{y}^{(i)})}_{\text{loss}} + \underbrace{\frac{\lambda}{2n} \sum_{l=1}^L \|w^{[l]}\|_F^2}_{\text{weight decay}}$$

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

$w^{[l]}: \begin{matrix} n^{[l]} & n^{[l-1]} \\ \uparrow & \uparrow \end{matrix}$

"Frobenius norm"

$$\|\cdot\|_2^2$$

$$\|\cdot\|_F^2$$

$$dw^{[l]} = \left[\text{(from backprop)} + \frac{\lambda}{n} w^{[l]} \right]$$

$$\frac{\partial J}{\partial w^{[l]}} = dw^{[l]}$$

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

"Weight decay"

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

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

$$= \underbrace{\left(1 - \frac{\alpha \lambda}{n}\right)}_{\leq 1} \underbrace{w^{[l]}}_{\text{weight}} - \alpha \text{(from backprop)}$$

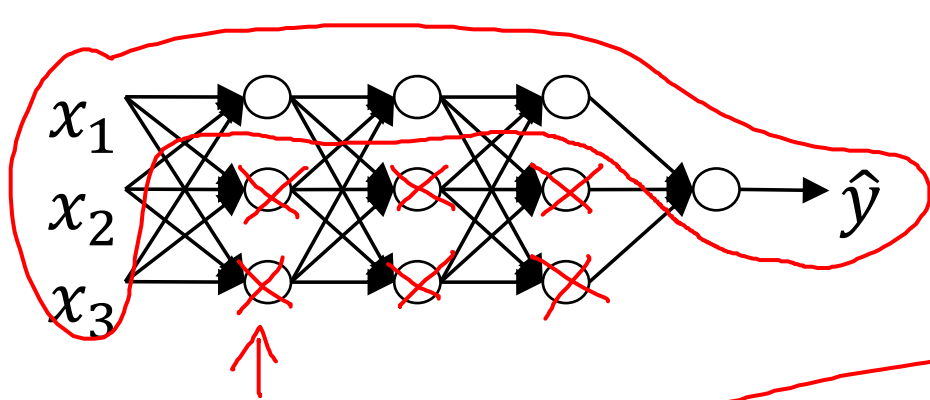


deeplearning.ai

Regularizing your neural network

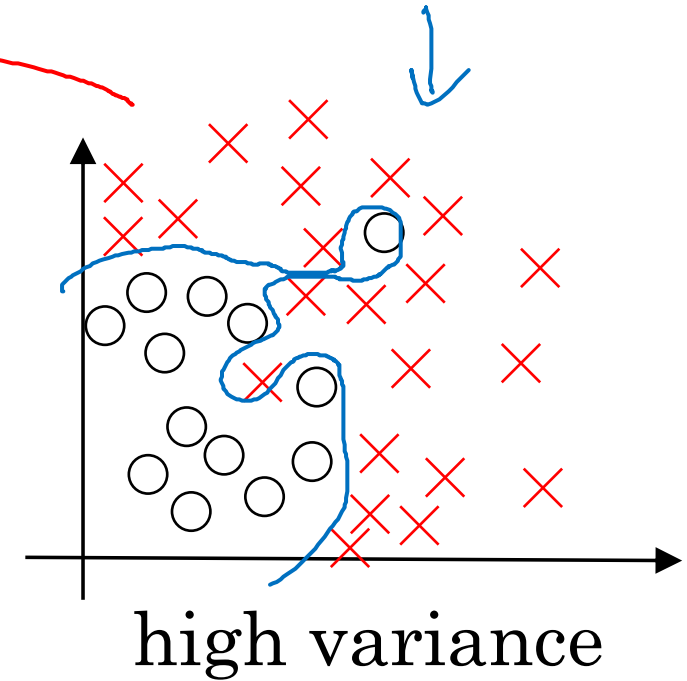
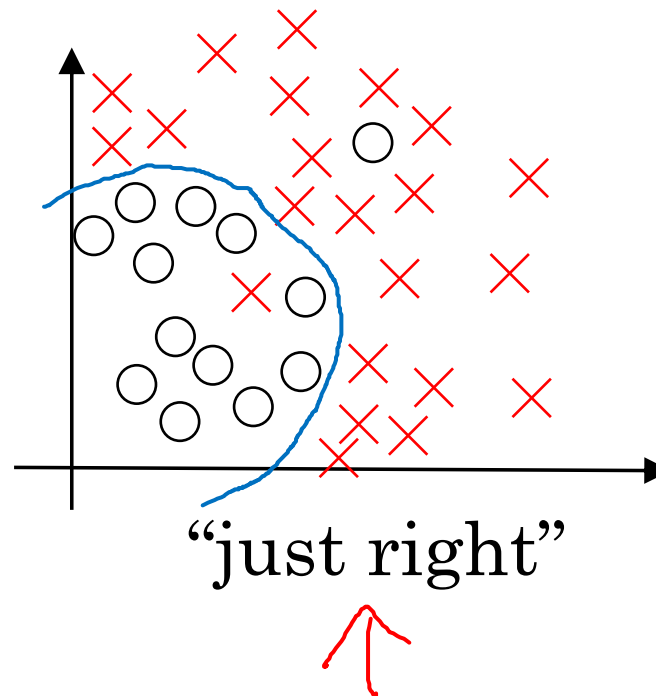
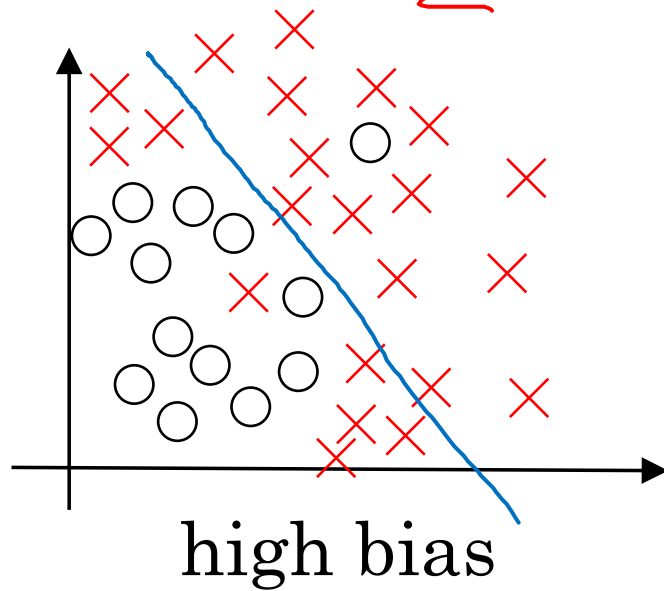
Why regularization reduces overfitting

How does regularization prevent overfitting?

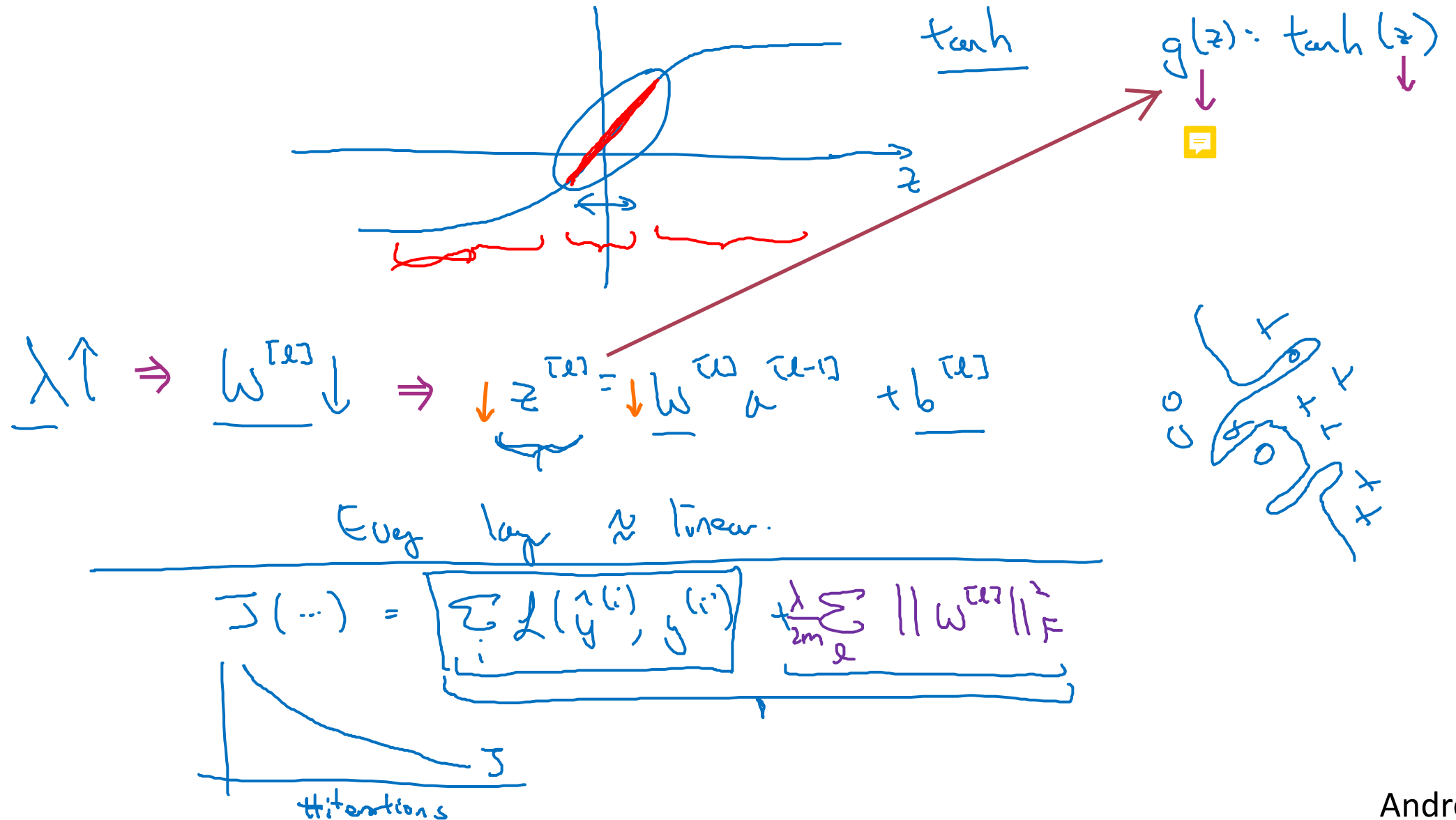


$$J(\mathbf{w}^{(L)}, \mathbf{b}^{(L)}) = \frac{1}{n} \sum_{i=1}^n \ell(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2n} \sum_{l=1}^L \underbrace{\|\mathbf{w}^{(l)}\|_F^2}_{\text{regularization}}$$

□ $\mathbf{w}^{(L)} \approx 0$



How does regularization prevent overfitting?



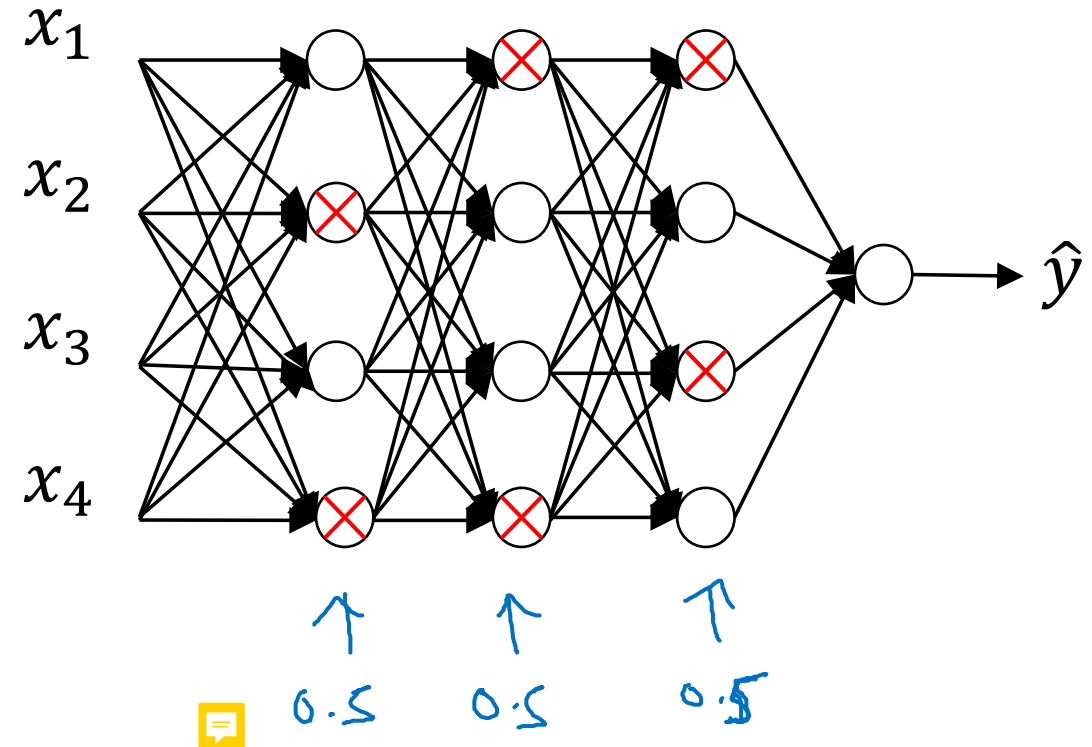
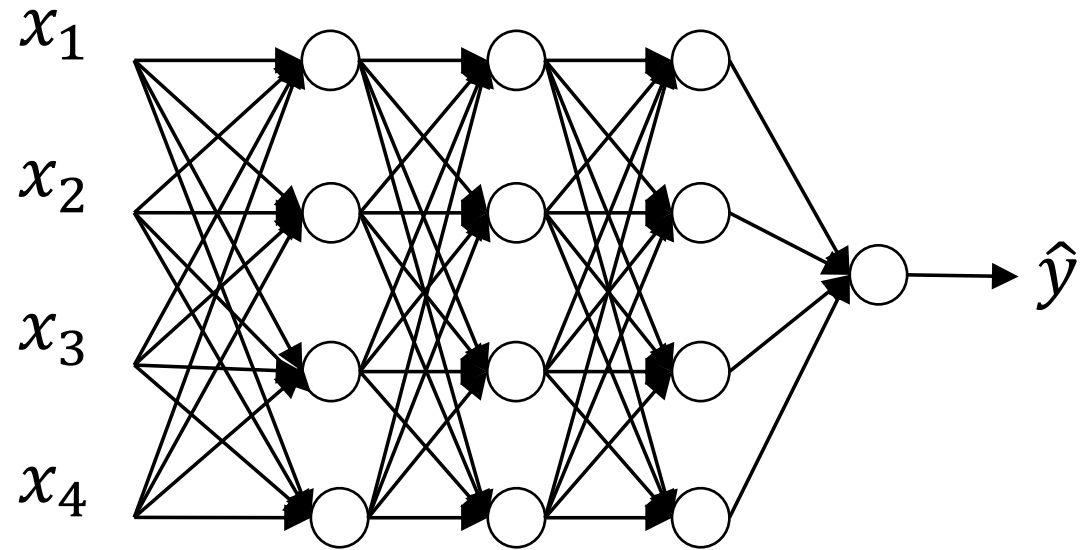


deeplearning.ai

Regularizing your neural network

Dropout regularization

Dropout regularization



Implementing dropout ("Inverted dropout")

Illustrate with layer $l=3$. keep-prob = 0.8 0.2

→ $d3 = \text{np.random.rand}(a3.\text{shape}[0], a3.\text{shape}[1]) < \text{keep-prob}$

$a3 = \text{np.multiply}(a3, d3)$ # $a3 \neq d3$.

→ $a3 /= \text{keep-prob}$ ←

50 units. \leadsto 10 units shut off 

$$z^{[4]} = w^{[4]} \cdot a^{[3]} + b^{[4]}$$

\uparrow

\uparrow 

reduced by 20%.

\uparrow
 $= 0.8$

Test

Making predictions at test time

$$a^{[0]} = X$$

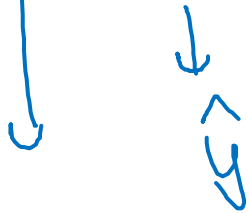
No drop out. 🗨️

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

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

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

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



/= keep-prob



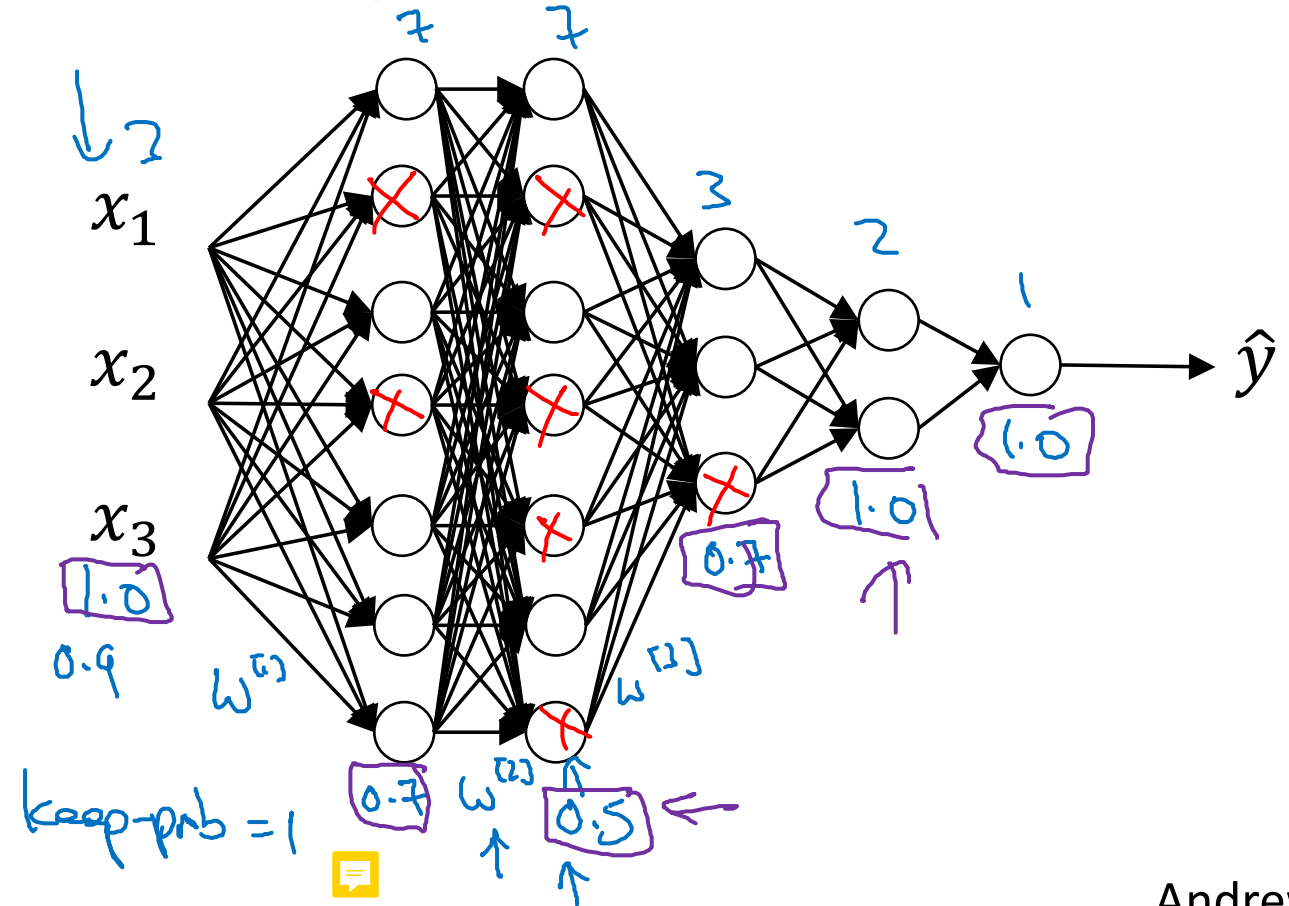
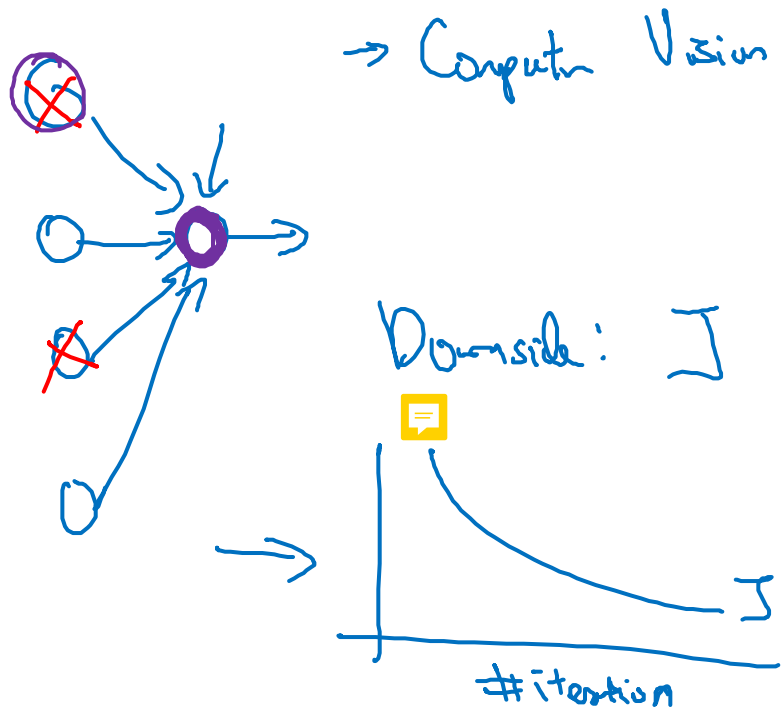
deeplearning.ai

Regularizing your neural network

Understanding dropout

Why does drop-out work?

Intuition: Can't rely on any one feature, so have to spread out weights. \rightsquigarrow Shrink weights. b_2





deeplearning.ai

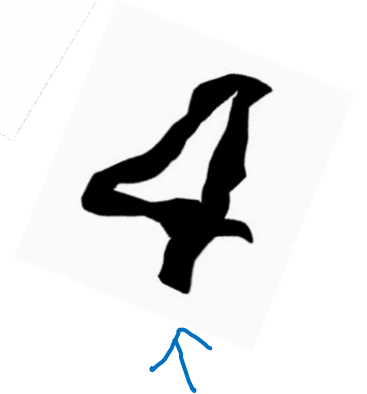
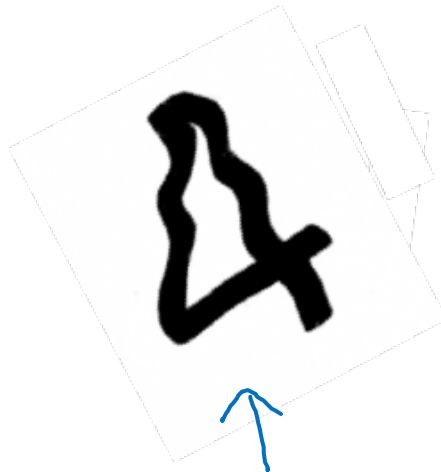
Regularizing your neural network

Other regularization methods

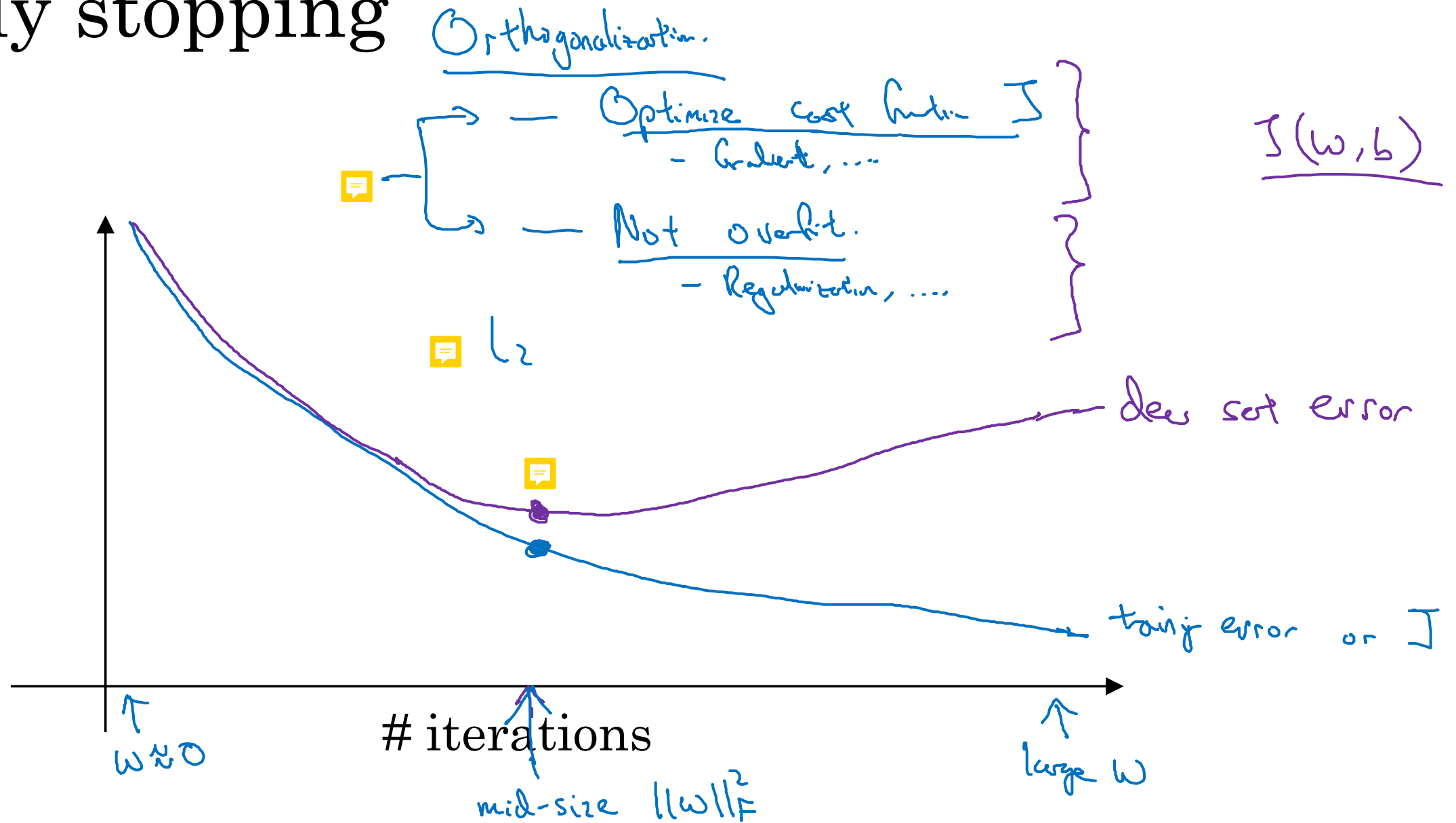
Data augmentation



4



Early stopping





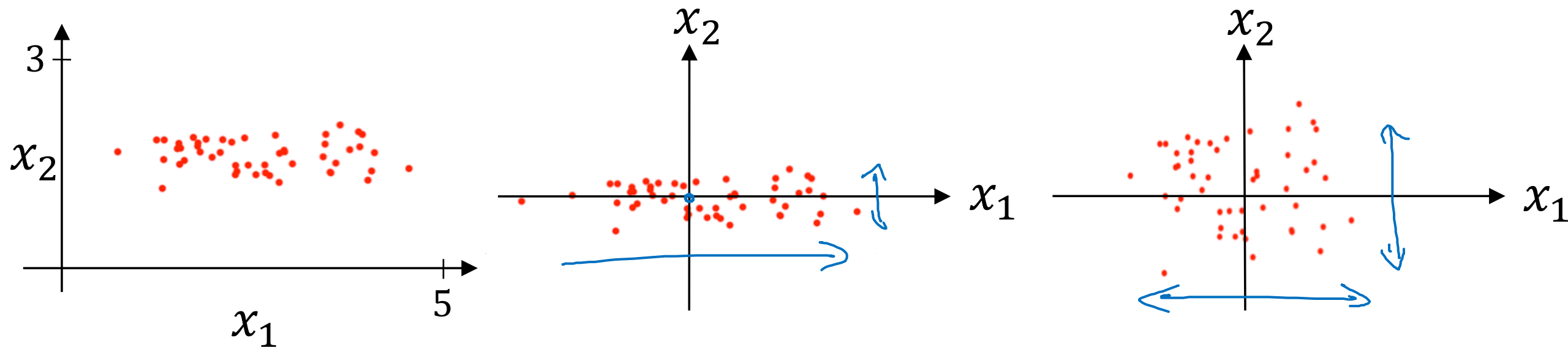
deeplearning.ai

Setting up your
optimization problem

Normalizing inputs

Normalizing training sets

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



Subtract mean:

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

$$x := x - \mu$$

Normalize variance

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

← element-wise

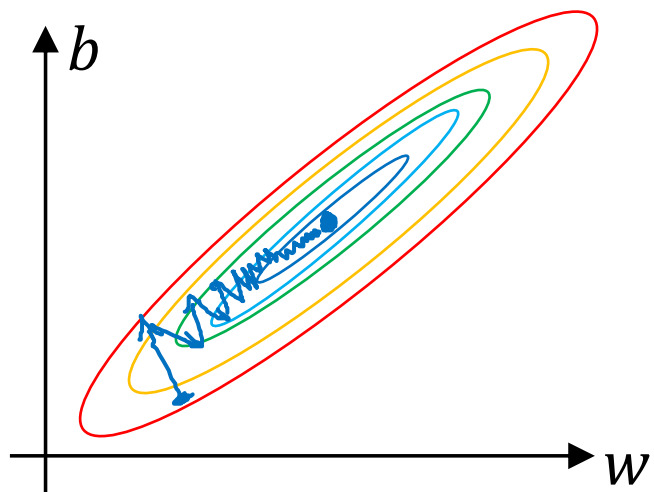
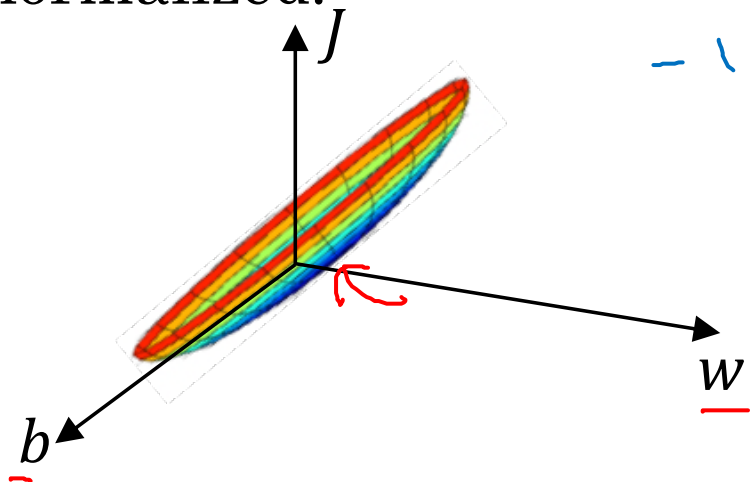
$$x /= \sigma^2$$

* * Use same μ σ^2 to normalize test set. * *

Why normalize inputs?

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

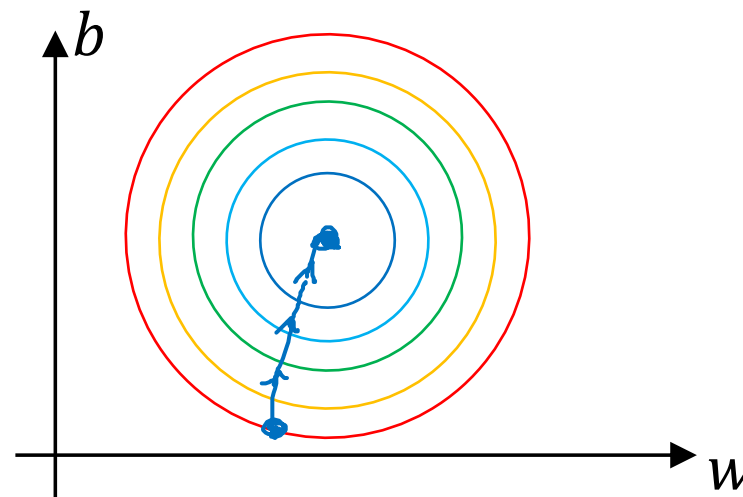
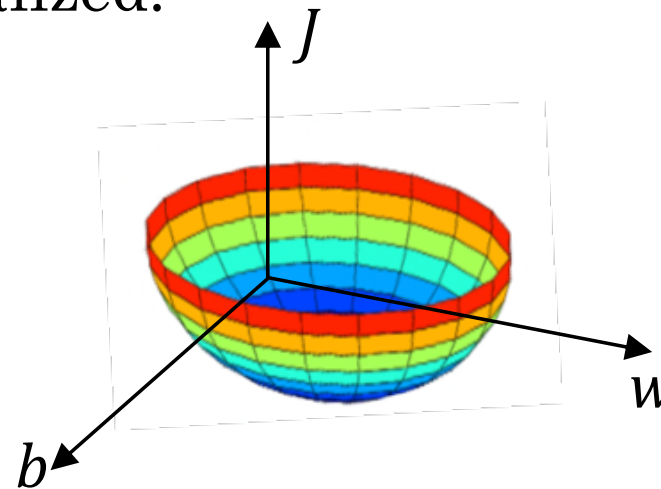
Unnormalized:



$w_1: x_1: \underline{1 \dots 1000} \leftarrow$
 $w_2: x_2: \underline{0 \dots 1} \leftarrow$
 $\quad \quad \quad -1 \dots 1$

$x_1: 0 \dots 1$
 $x_2: -1 \dots 1$
 $x_3: 1 \dots 2$

Normalized:



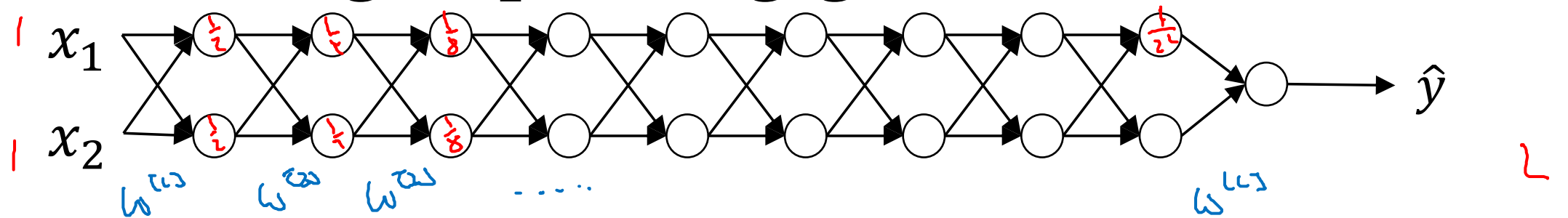


deeplearning.ai

Setting up your
optimization problem

Vanishing/exploding
gradients

Vanishing/exploding gradients



Handwritten notes: $g(z) = z$ and $b^{(L)} = 0$.

Diagram illustrating the forward pass of the RNN. The output \hat{y} is calculated as $\hat{y} = w^{(L)} a^{(L)}$, where $a^{(L)}$ is the hidden state at time step L . The hidden state is updated recursively: $a^{(t)} = g(w^{(t)} x + a^{(t-1)})$. The diagram shows the sequence of hidden states $a^{(1)}, a^{(2)}, \dots, a^{(L)}$ and the weights $w^{(1)}, w^{(2)}, \dots, w^{(L)}$.

Handwritten notes: $w^{(1)} > I$

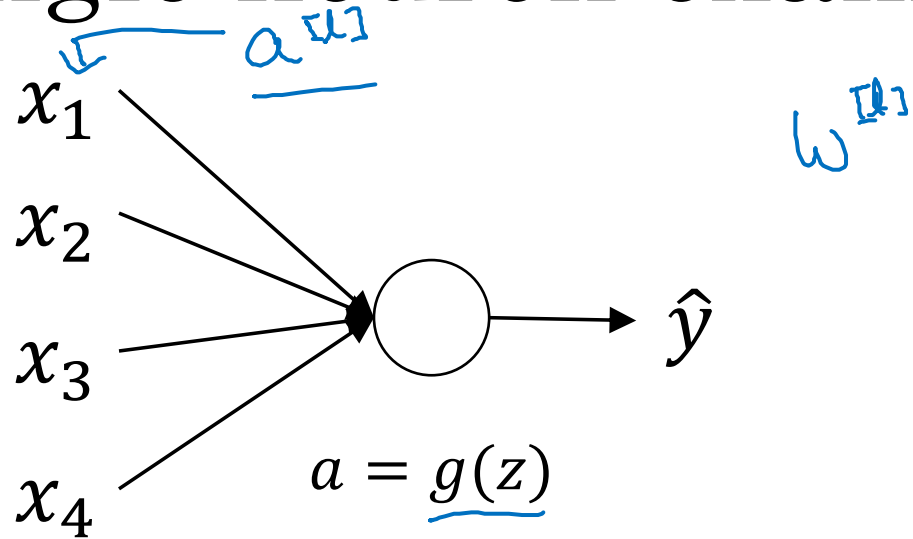
Handwritten notes: $w^{(2)} < I$ and a matrix $\begin{bmatrix} 0.9 & 0.9 \\ 0 & 0 \end{bmatrix}$.

Handwritten notes: A matrix $w^{(2)} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}$ with 0.5 and 0.5 written above the diagonal elements.

Handwritten notes: The output \hat{y} is calculated as $\hat{y} = w^{(L)} a^{(L)}$, where $a^{(L)}$ is the hidden state at time step L . The diagram shows the sequence of hidden states $a^{(1)}, a^{(2)}, \dots, a^{(L)}$ and the weights $w^{(1)}, w^{(2)}, \dots, w^{(L)}$.

Handwritten notes: $1.5^{L-1} \times$ and $0.5^{L-1} \times$.

Single neuron example



$$z = \underline{w_1} x_1 + \underline{w_2} x_2 + \dots + \underline{w_n} x_n \quad \text{to}$$

large $n \rightarrow$ Smaller w_i

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

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

ReLU $g^{[2]}(z) = \text{ReLU}(z)$

Other vars:

$$\text{tanh} \left(\frac{1}{n^{[l-1]}} \right)$$

Xavier initialization ↑

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

↑



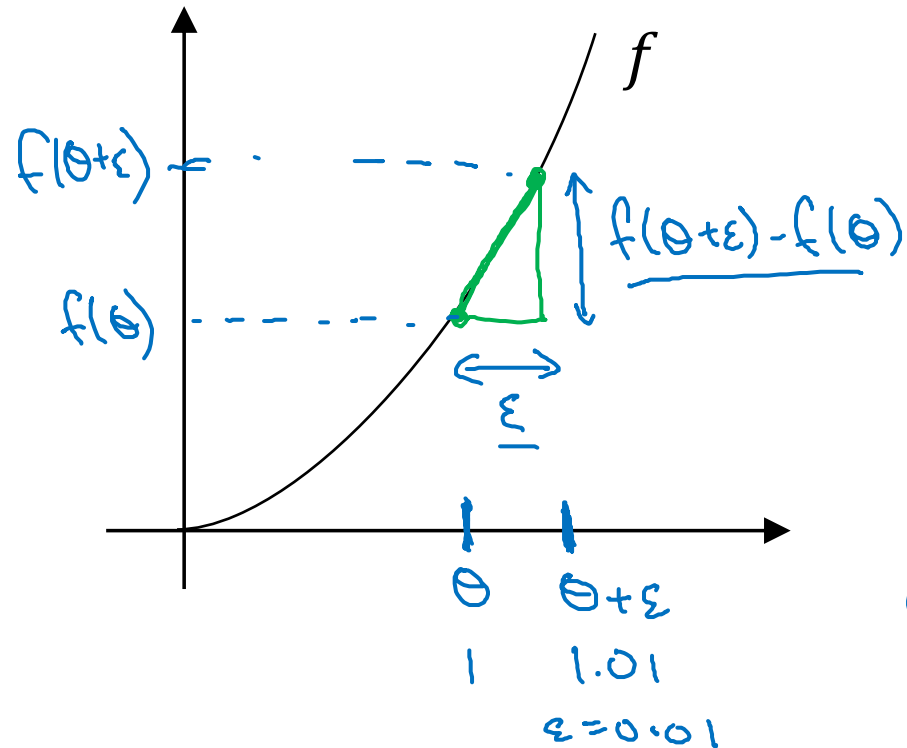
deeplearning.ai

Setting up your optimization problem

Numerical approximation of gradients

Checking your derivative computation

I $f(\theta) = \theta^3$
 $\theta \in \mathbb{R}.$



$$g(\theta) = \frac{d}{d\theta} f(\theta) = f'(\theta)$$

$\frac{dw}{db}$ \rightarrow $g(\theta) = 3\theta^2$

$g(\theta) = 3 \cdot (1)^2 = 3$
 when $\theta = 1$

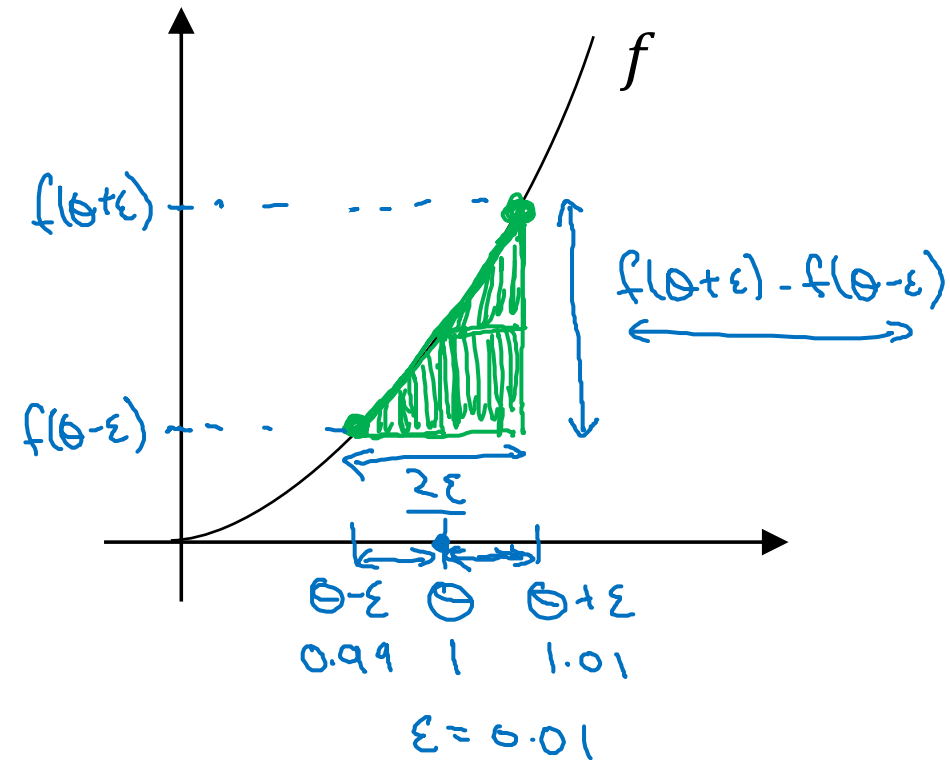
$$\frac{f(\theta + \epsilon) - f(\theta)}{\epsilon} \approx g(\theta)$$

$$\frac{(1.01)^3 - 1^3}{0.01} = \frac{1.030301 - 1}{0.01} = \frac{0.0301}{0.01} = 3.0301 \approx 3$$

Annotations: $\theta = 1$, $\theta + \epsilon = 1.01$, $\epsilon = 0.01$, 3.1 , 3.2

Checking your derivative computation

$$\underline{f(\theta) = \theta^3}$$



$$\left[\frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon} \approx \underline{g(\theta)} \right]$$

$$\frac{(1.01)^3 - (0.99)^3}{2(0.01)} = 3.0001 \approx 3$$

$$g(\theta) = 3\theta^2 = 3$$

approx error: 0.0001

(prev slide: 3.0301. error: 0.03)

$$\left\{ \begin{array}{l} f'(\theta) = \lim_{\epsilon \rightarrow 0} \frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon} \end{array} \right.$$

$$\frac{O(\epsilon^2)}{O(\epsilon)} = O(\epsilon)$$

0.01
0.0001

$$\frac{f(\theta + \epsilon) - f(\theta)}{\epsilon} \quad \text{error: } O(\epsilon)$$

0.01



deeplearning.ai

Setting up your
optimization problem

Gradient Checking

Gradient check for a neural network

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

$$J(w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}) = J(\theta)$$

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

Is $d\theta$ the gradient of $J(\theta)$?

Gradient checking (Grad check)

$$J(\theta) = J(\theta_1, \theta_2, \theta_3, \dots)$$

for each i :

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

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

Checks

$$\rightarrow \frac{\|d\theta_{\text{approx}} - d\theta\|_2}{\|d\theta_{\text{approx}}\|_2 + \|d\theta\|_2}$$
$$\epsilon = 10^{-7}$$

$$\approx \frac{10^{-7}}{10^{-5}} - \text{great!} \leftarrow$$
$$\rightarrow 10^{-3} - \text{worry.} \leftarrow$$



deeplearning.ai

Setting up your
optimization problem

Gradient Checking
implementation notes

Gradient checking implementation notes

- Don't use in training – only to debug

$$\frac{d\theta_{\text{approx}}[\vec{i}]}{\uparrow \uparrow} \longleftrightarrow \frac{d\theta[\vec{i}]}{\uparrow}$$

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

$$\frac{db^{[L]}}{\uparrow} \quad \frac{dW^{[L]}}{\uparrow}$$

- Remember regularization. 🗨

$$\underline{J(\theta)} = \frac{1}{n} \sum_i \ell(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2n} \sum_l \|W^{[l]}\|_F^2$$

$d\theta = \text{gradient of } J \text{ wrt. } \theta$

- Doesn't work with dropout. 🗨 J

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

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

$$\underline{W, b \approx 0}$$