

Copyright Notice

These slides are distributed under the Creative Commons License.

DeepLearning.AI makes these slides available for educational purposes. You may not use or distribute these slides for commercial purposes. You may make copies of these slides and use or distribute them for educational purposes as long as you cite DeepLearning.AI as the source of the slides.

For the rest of the details of the license, see <https://creativecommons.org/licenses/by-sa/2.0/legalcode>

Setting up your ML application

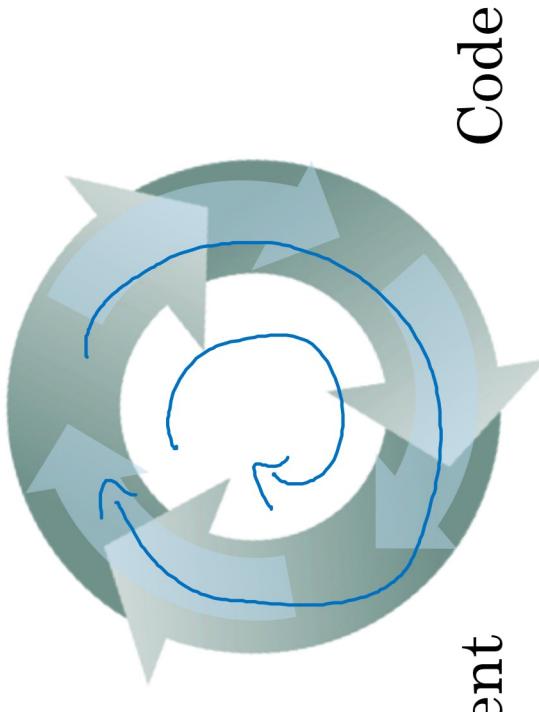
Train/dev/test sets



deeplearning.ai

Applied ML is a highly iterative process

layers
hidden units
learning rates
activation functions
...

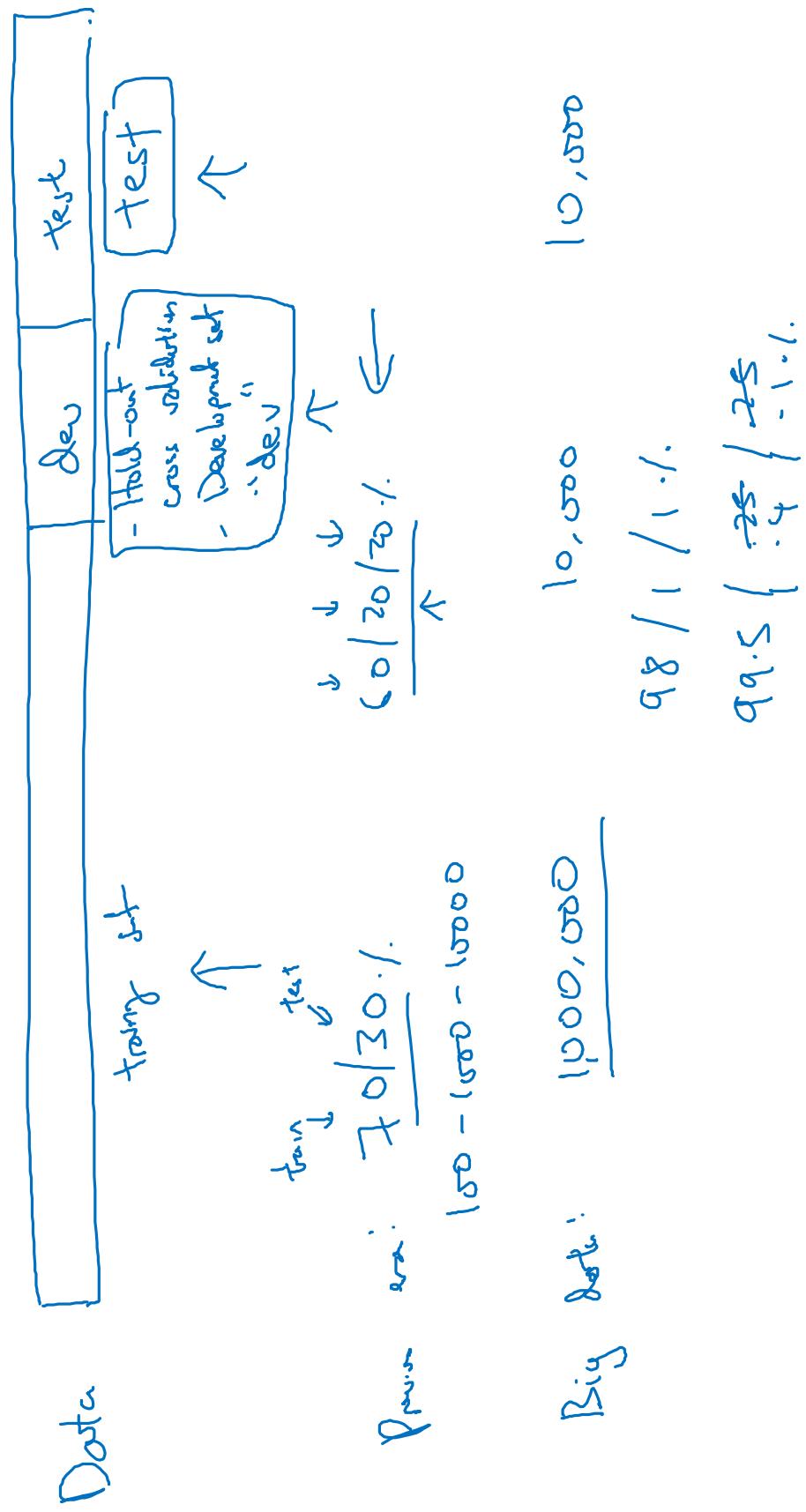


Code
Experiment

NLP, Vision, Speech, Structural Data
ADS
Search
Survey
Logistic ...

Andrew Ng

Train/dev/test sets



Mismatched train/test distribution

Cont'd



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.

"test"



train / test

train / dev

train / dev

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

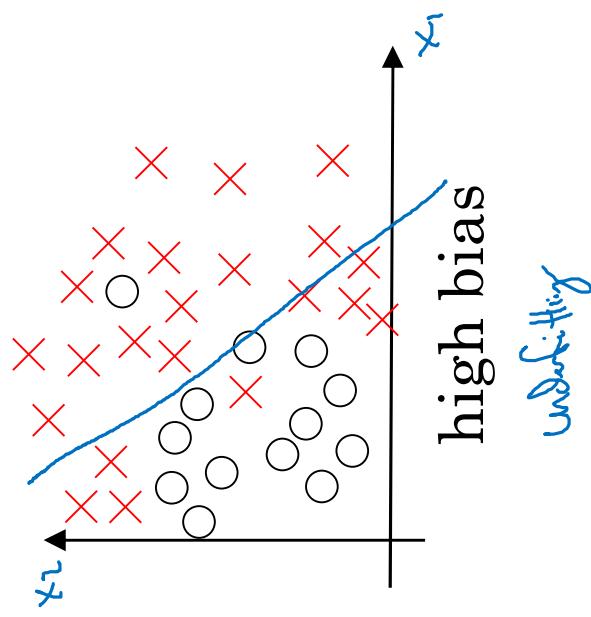
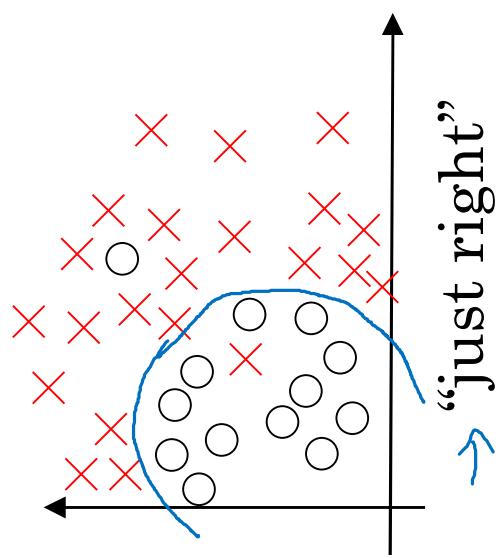
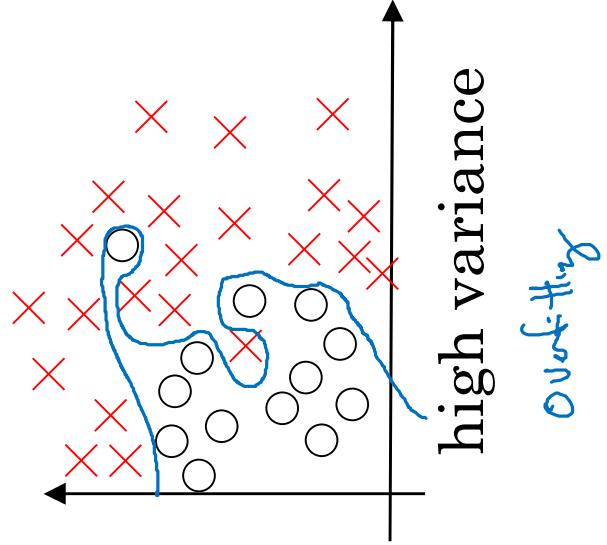
Setting up your
ML application

Bias/Variance



deeplearning.ai

Bias and Variance



Bias and Variance

Cat classification

$$y=0$$



$$y=1$$

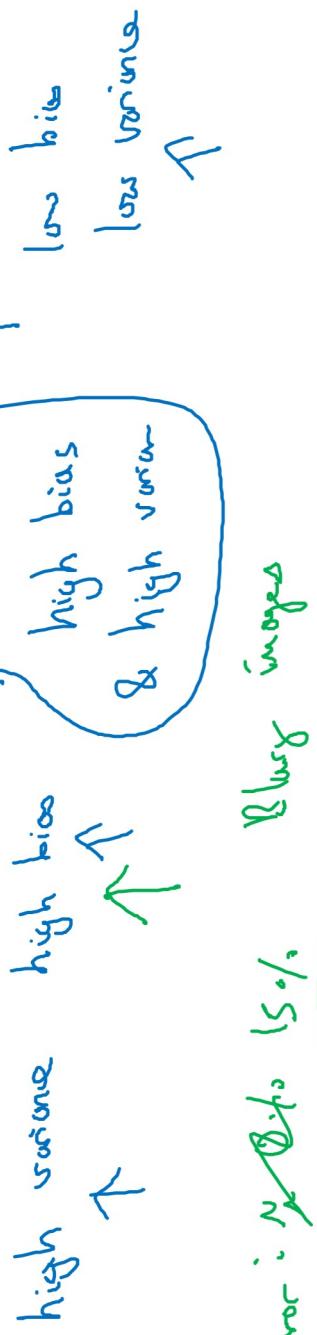


Train set error:

$$\frac{15\% \leftarrow}{16\% \leftarrow}$$

Dev set error:

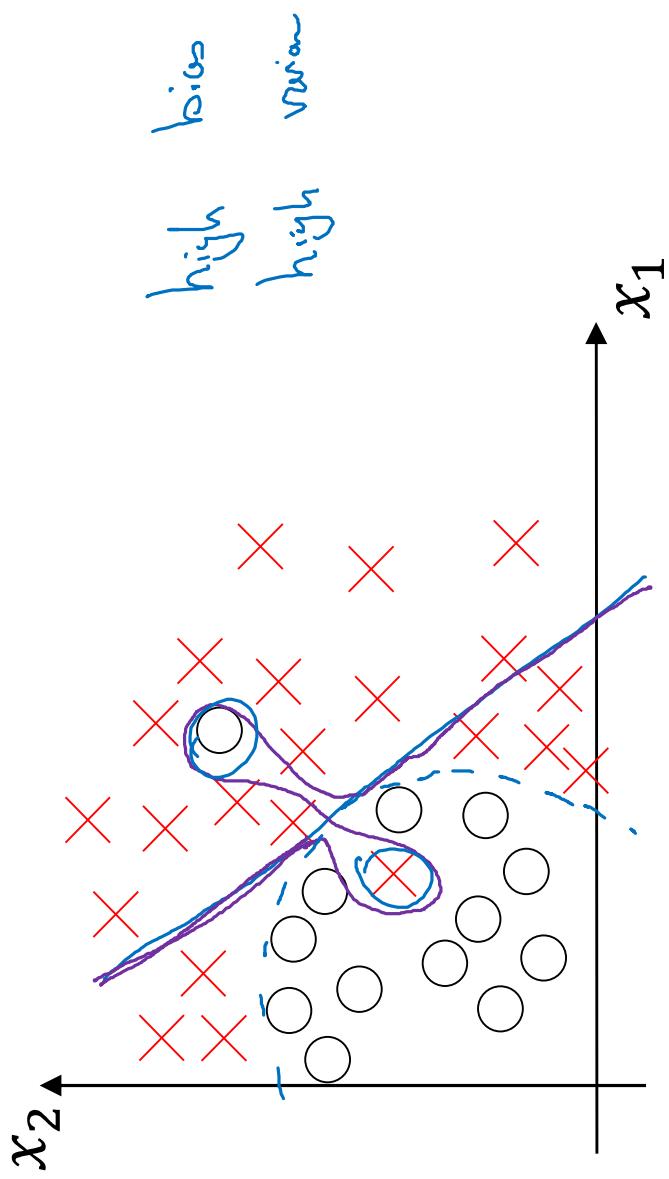
$$30\%$$



$$\text{Human: } 20\%$$

Optimal (bias-variance) error: ~16%
Blury images

High bias and high variance



for machine learning

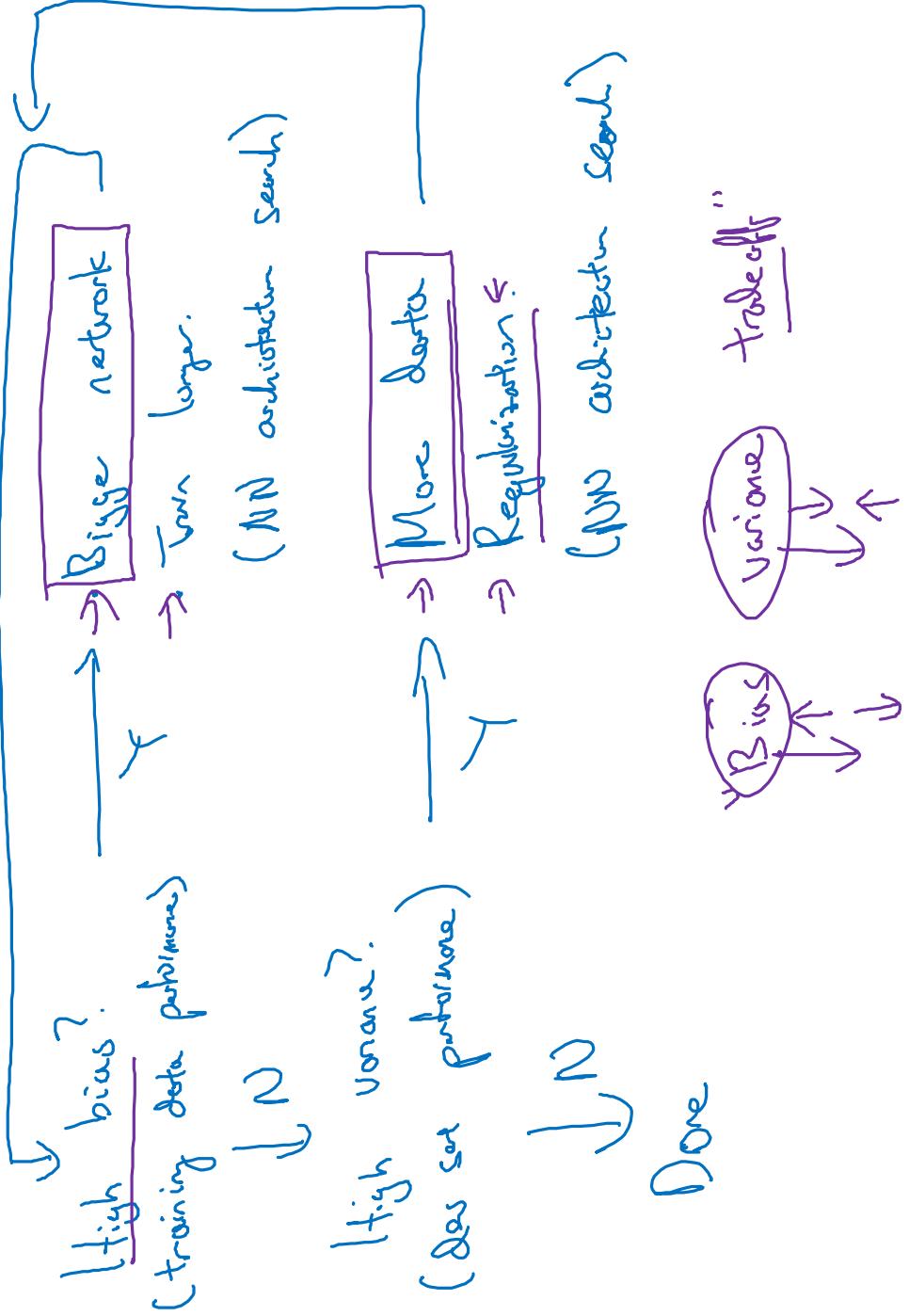
Basic “recipe”

Setting up your
ML application

deeplearning.ai



Basic recipe for machine learning



Regularizing your
neural network

Regularization



deeplearning.ai

Logistic regression

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

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

$\lambda = \text{regularization parameter}$
 λ

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

$$+ \frac{\lambda}{2m} \left(\frac{b^2}{2m} \right)$$

unit

$$\|w\|_2^2 = \sum_{j=1}^n w_j^2 = w^T w \rightarrow$$

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

w will be sparse

Neural network

$$\rightarrow \mathcal{J}(\omega^{(1)}, b^{(1)}, \dots, \omega^{(L)}, b^{(L)}) = \frac{1}{m} \sum_{i=1}^m f(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \sum_{j=1}^{n^{(l-1)}} \|\omega^{(l)}\|_F^2$$

$$\|\omega^{(l)}\|_F^2 = \sum_{i=1}^m \sum_{j=1}^{n^{(l-1)}} (\omega_{ij}^{(l)})^2$$

"Frobenius norm"

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

$$\frac{\partial \mathcal{J}}{\partial \omega^{(l)}} = \frac{\partial \mathcal{J}}{\partial \omega^{(l)}} = \frac{\partial \mathcal{J}}{\partial \omega^{(l)}} = \frac{\partial \mathcal{J}}{\partial \omega^{(l)}}$$

$$\delta \omega^{(l)} = \left[\text{(from backprop)} + \frac{\lambda}{m} \omega^{(l)} \right]$$

$$\Rightarrow \omega^{(l)} := \omega^{(l)} - \alpha \delta \omega^{(l)}$$

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

"Weight decay"

$$\begin{aligned} &= \omega^{(l)} - \frac{\alpha}{m} \omega^{(l)} - \alpha (\text{from backprop}) \\ &= \underbrace{\left(1 - \frac{\alpha \lambda}{m}\right)}_{< 1} \omega^{(l)} - \alpha (\text{from backprop}) \end{aligned}$$

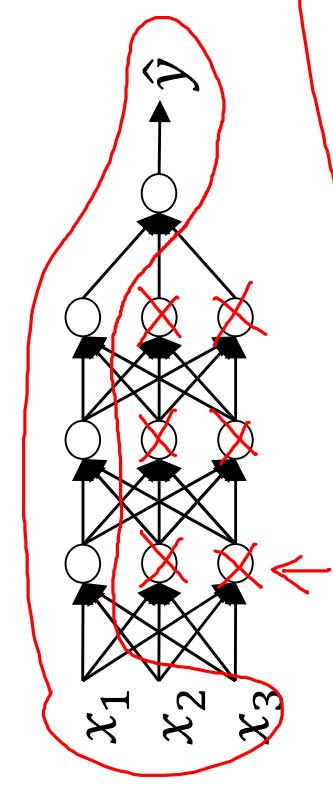
Regularizing your
neural network

Why regularization
reduces overfitting



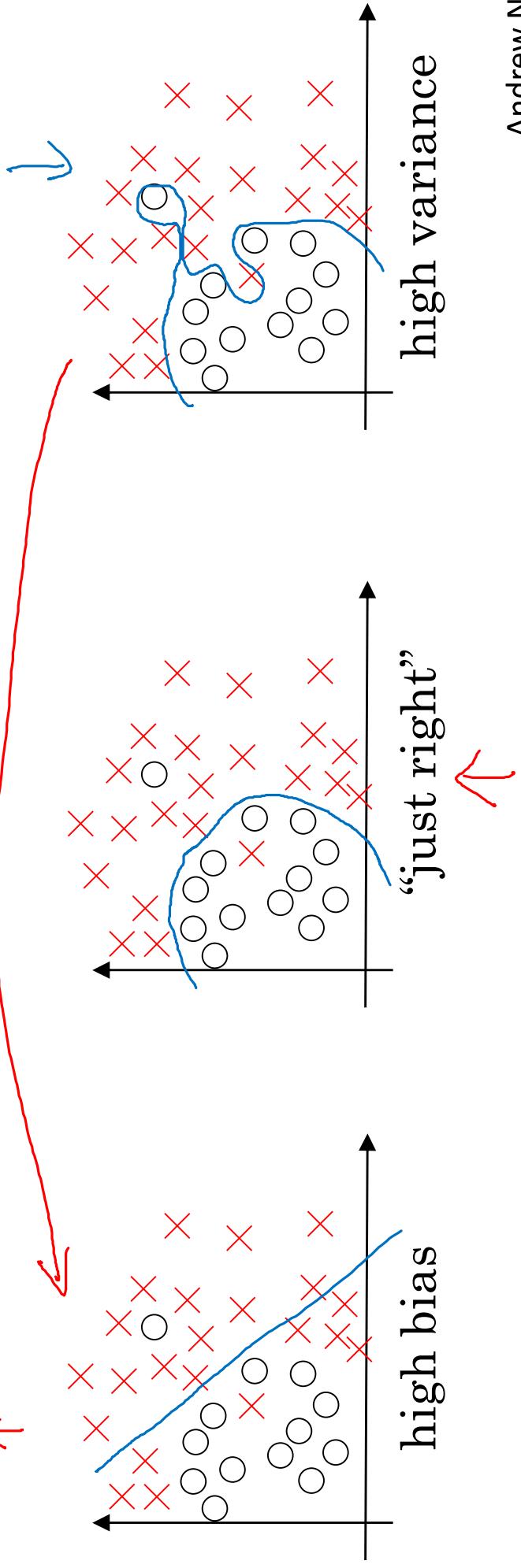
deeplearning.ai

How does regularization prevent overfitting?



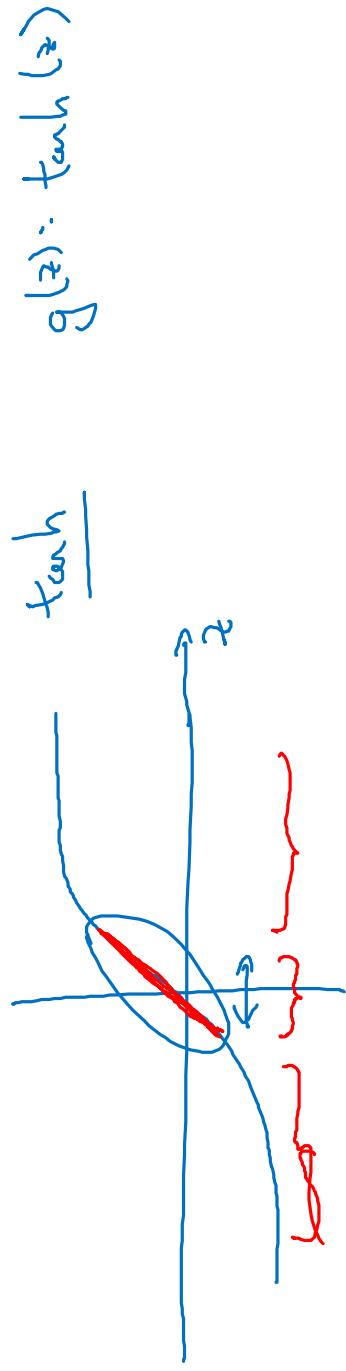
$$J(\omega^{(m)}, b^{(m)}) = \frac{1}{m} \sum_{i=1}^m l(x^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \|\omega^{(m)}\|_F^2$$

$\omega^{(m)} \approx 0$

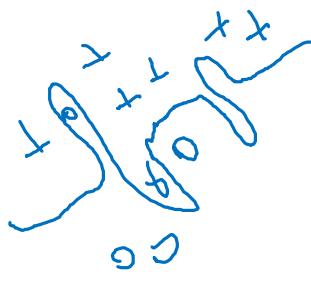


Andrew Ng

How does regularization prevent overfitting?



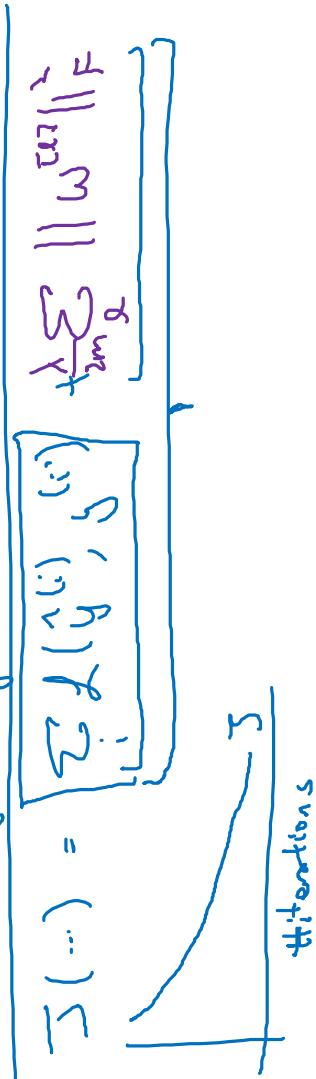
$$g(z) := \tanh(z)$$



$$z^{(i)} = \omega_0 + \omega_1 x_1 + \dots + \omega_n x_n$$

Every layer is linear.

$$\lambda \uparrow \quad \omega^{(k)} \downarrow$$



$$J(\cdot) = \sum \ell(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2m} \sum \|\omega^{(k)}\|_F^2$$

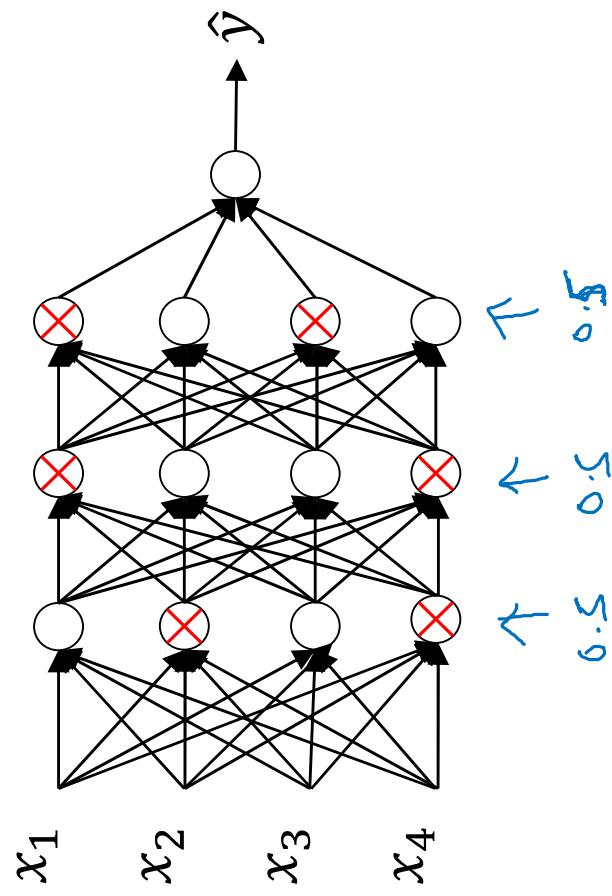
Regularizing your
neural network

Dropout
regularization

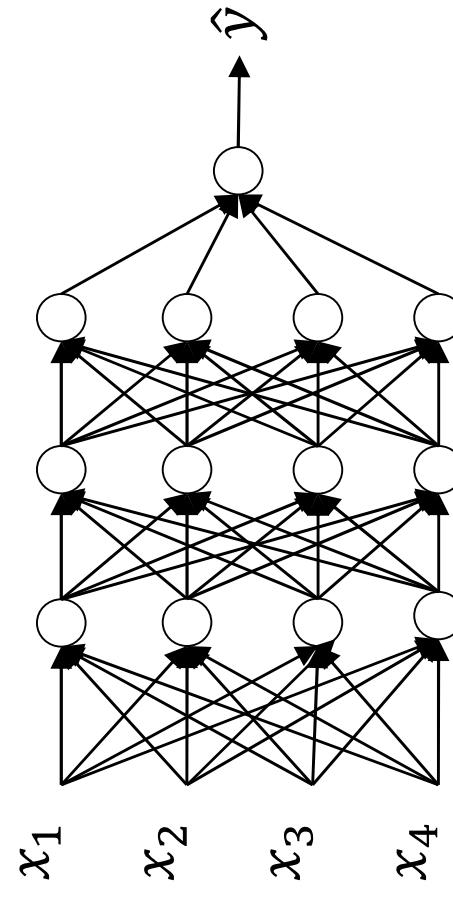
deeplearning.ai



Dropout regularization



Andrew Ng



Implementing dropout (“Inverted dropout”)

Illustrate with layer $\ell=3$. $\text{keep-prob} = \frac{0.8}{\cancel{10}}$

$$\Rightarrow \boxed{\alpha_3} = \text{np.random.rand}(\alpha_3.\text{shape}[0], \alpha_3.\text{shape}[1]) < \text{keep-prob}$$

$$\alpha_3 = \text{np.multiply}(\alpha_3, \delta_3) \quad \# \alpha_3 \neq \delta_3.$$

$$\Rightarrow \boxed{\alpha_3} / = \cancel{\text{keep-prob}} \leftarrow$$

50 units. \rightsquigarrow 10 units shut off

$$z^{(4)} = w^{(4)} \cdot \alpha_3 + b^{(4)}$$

↓

$\cancel{10}$ units reduced by $\cancel{20/1}$.

Test

$$l = \underline{0.8}$$

Making predictions at test time

$$Q^{(t)} = X$$

$$\frac{\text{No. drop out.}}{z^{(t)}} = \frac{w^{(t)} Q^{(t)} + b^{(t)}}{g^{(t)} (z^{(t)})}$$
$$Q^{(t)} = \dots$$
$$z^{(t)} = \frac{w^{(t)} Q^{(t)} + b^{(t)}}{g^{(t)} (z^{(t)})}$$
$$g^{(t)} = \dots$$


$\lambda = \text{keep-pool}$

Regularizing your
neural network

Understanding
dropout

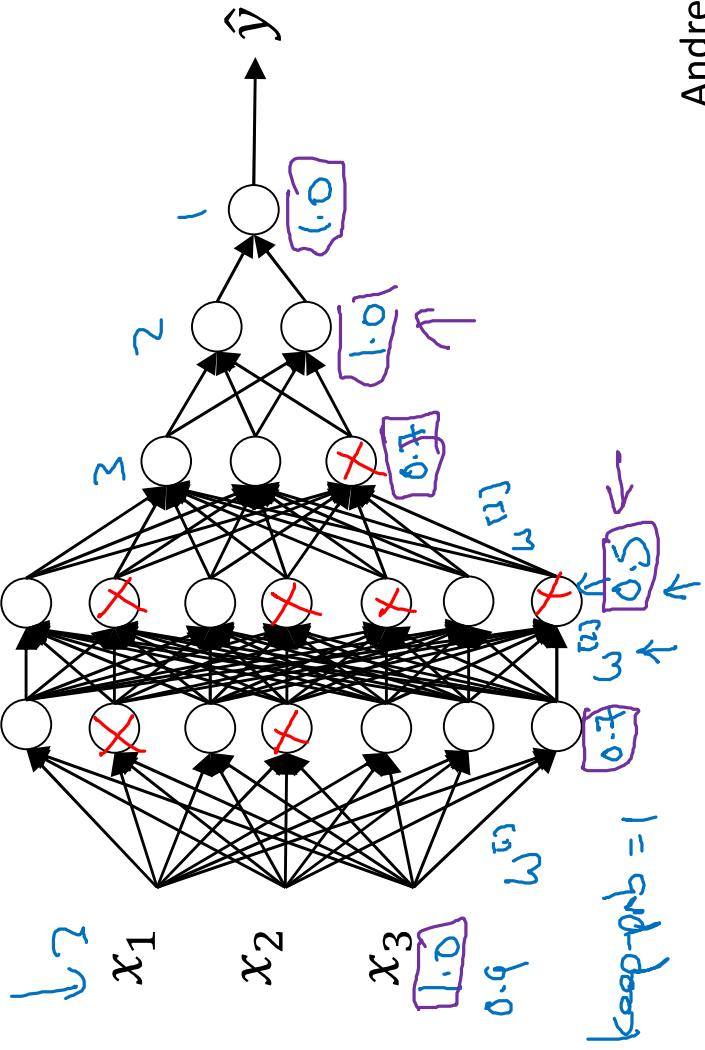
deeplearning.ai



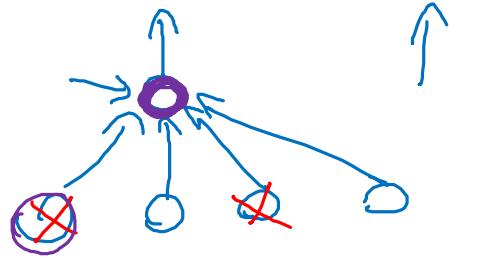
Why does drop-out work?

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

→ Computer Vision
↓
 x_1 x_2 x_3
↓
→ Shrink weights.



Demand: \int
Iteration #



Regularizing your
neural network

Other regularization
methods

deeplearning.ai



Data augmentation

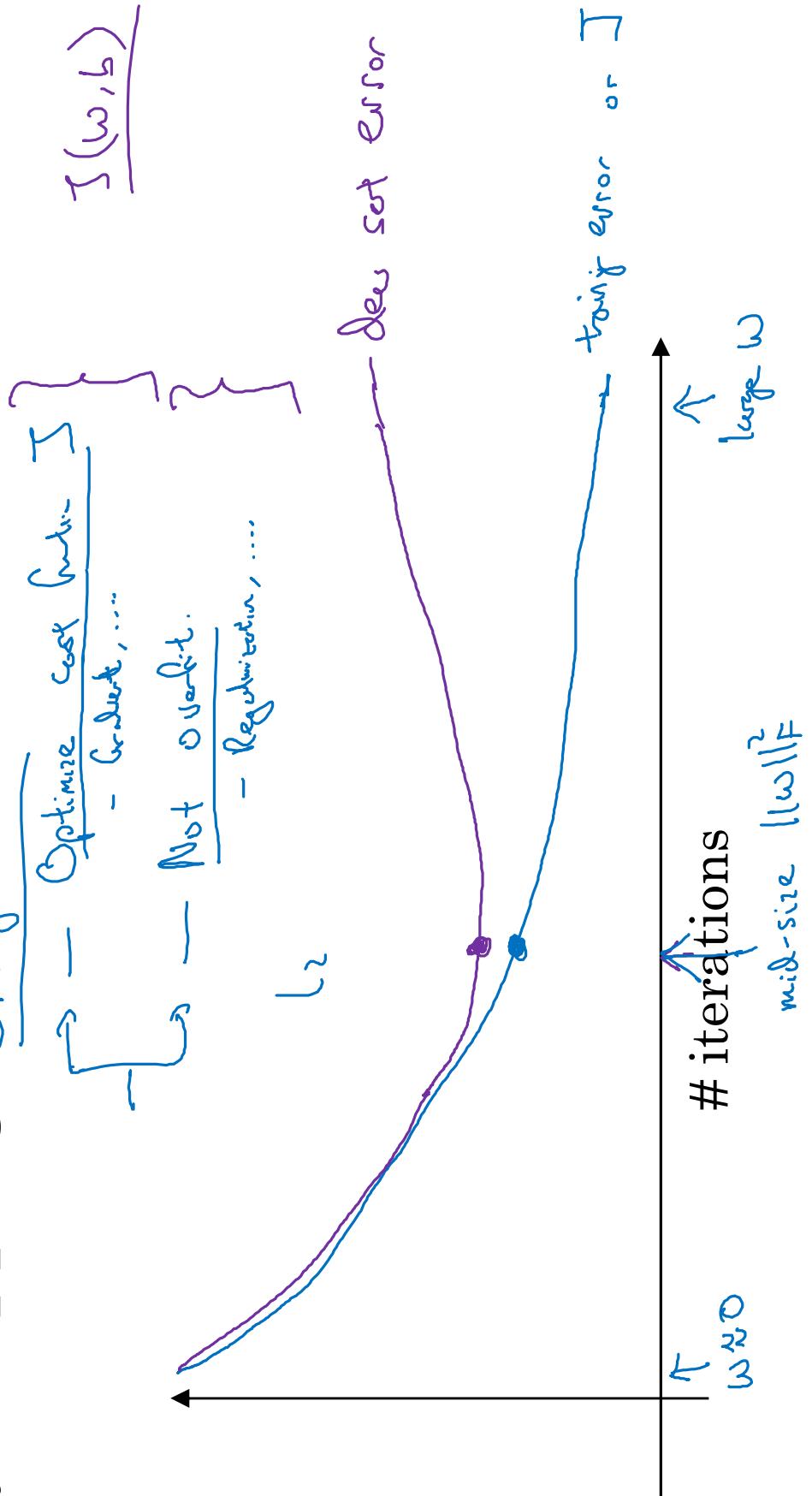


4



Early stopping

Orthogonalization.



Setting up your
optimization problem

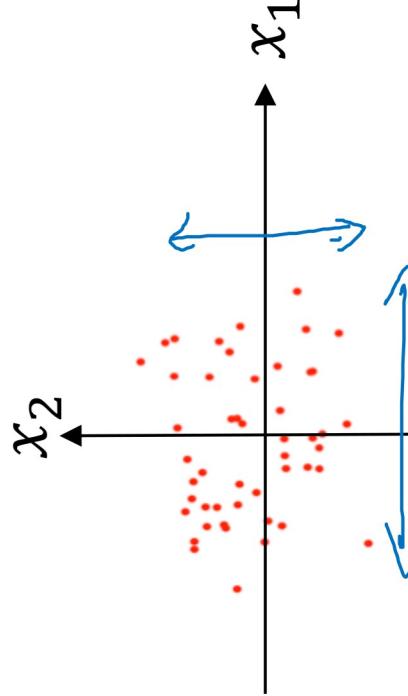
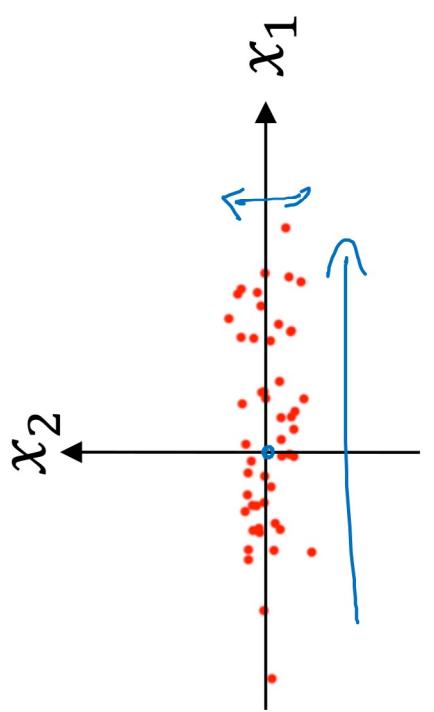
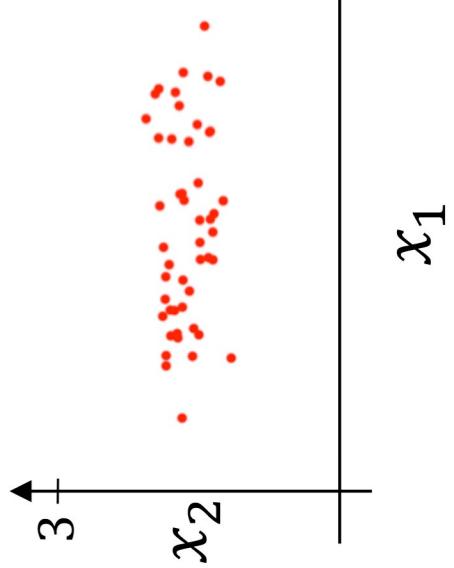
Normalizing inputs



deeplearning.ai

Normalizing training sets

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



Subtract mean:

$$\hat{x}_i = \frac{1}{n} \sum_{i=1}^n x^{(i)}$$

Normalize variance:

$$\frac{\sum_{i=1}^n (x^{(i)} - \hat{x})^2}{n}$$

↑ element-wise

Use $\hat{x}_i = \frac{x_i - \mu}{\sigma}$ to normalize test set.

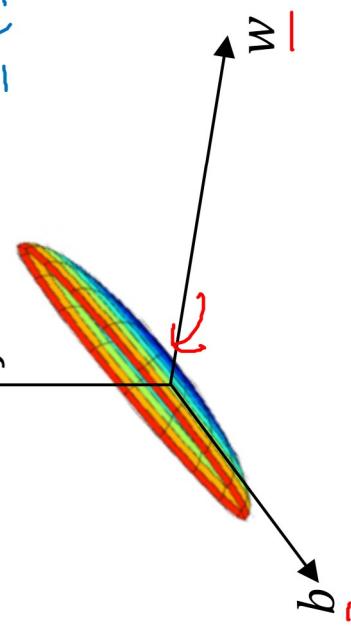
Andrew Ng

Why normalize inputs?

$$\omega_1: x_1: \frac{1 \dots 1000}{0 \dots 1} \leftarrow$$

$$\omega_2: x_2: \frac{0 \dots 1}{-1 \dots 1} \leftarrow$$

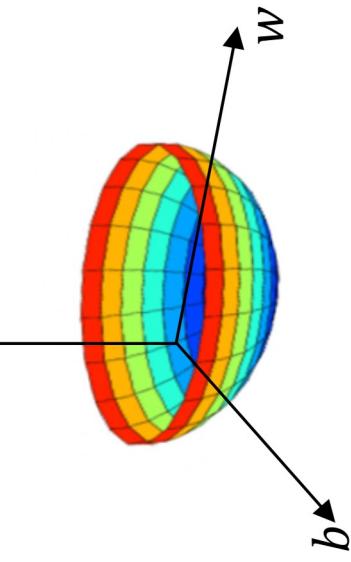
Unnormalized:
 $x_1: 6 \dots 1$
 $x_2: -1 \dots 1$
 $x_3: 1 \dots 2$



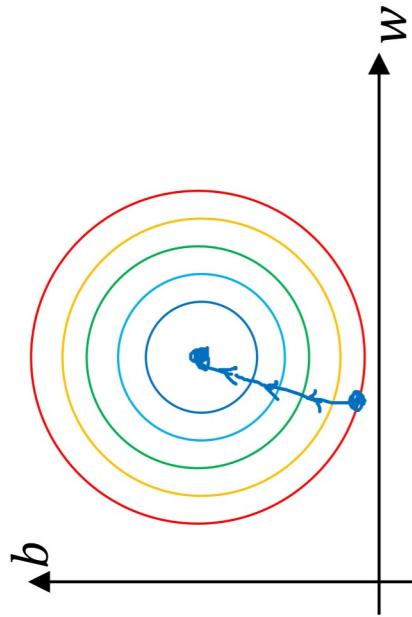
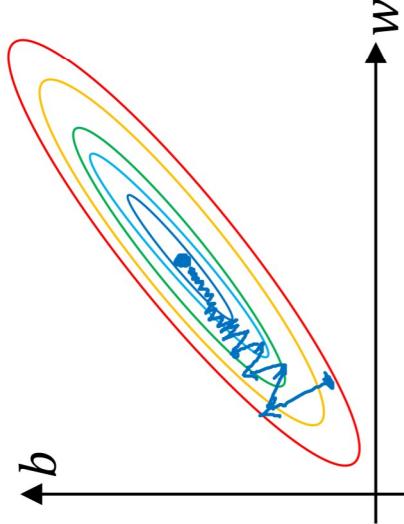
$$J(w, b) = \frac{1}{m} \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

Normalized:

$\uparrow J$



$$\begin{aligned}x_1 &: 6 \dots 1 \\x_2 &: -1 \dots 1 \\x_3 &: 1 \dots 2\end{aligned}$$



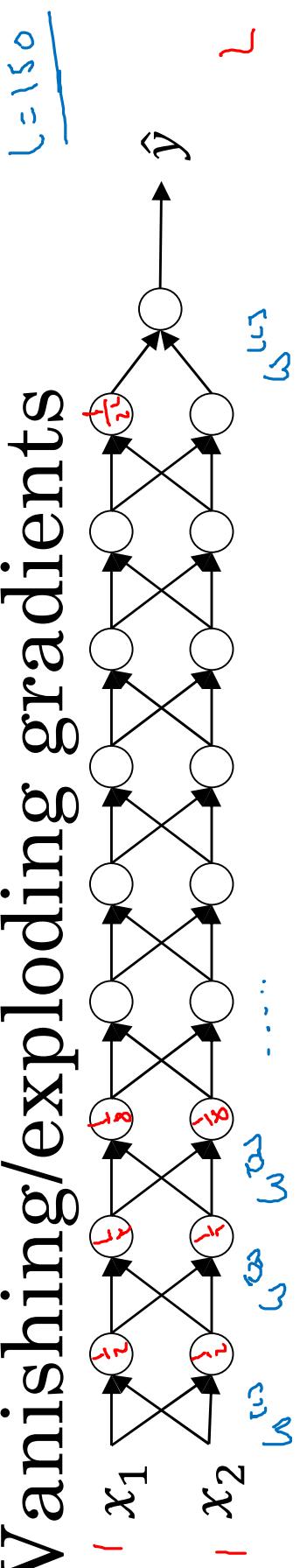
W Andrew Ng

Vanishing/exploding gradients

Setting up your
optimization problem



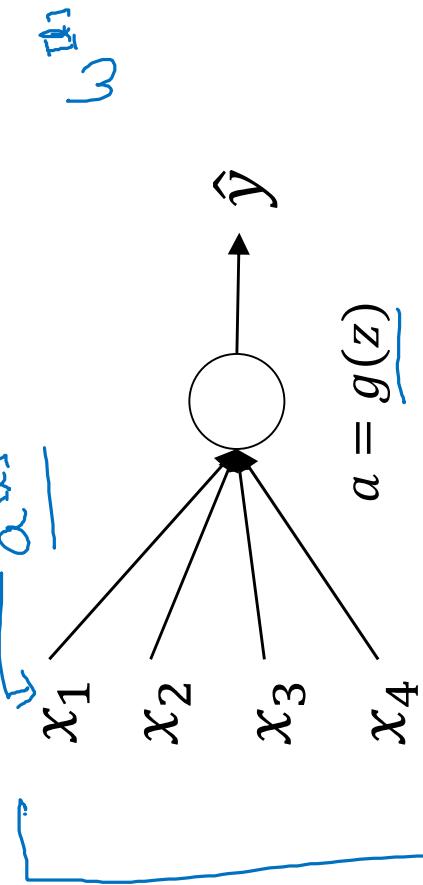
Vanishing/exploding gradients



$$\begin{aligned}
 & \text{Input: } x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \\
 & \text{Hidden Layer 1: } z^{(1)} = \begin{bmatrix} z^{(1)}_1 \\ z^{(1)}_2 \\ z^{(1)}_3 \end{bmatrix}, \quad w^{(1)} = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{bmatrix}, \quad a^{(1)} = g(z^{(1)}) = \begin{bmatrix} a^{(1)}_1 \\ a^{(1)}_2 \\ a^{(1)}_3 \end{bmatrix} \\
 & \text{Hidden Layer 2: } z^{(2)} = \begin{bmatrix} z^{(2)}_1 \\ z^{(2)}_2 \\ z^{(2)}_3 \end{bmatrix}, \quad w^{(2)} = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{bmatrix}, \quad a^{(2)} = g(z^{(2)}) = \begin{bmatrix} a^{(2)}_1 \\ a^{(2)}_2 \\ a^{(2)}_3 \end{bmatrix} \\
 & \text{Output: } \hat{y} = \hat{w}^T a^{(2)} + b
 \end{aligned}$$

Below the equations, there are two purple ovals. The left oval contains a matrix $w^{(1)}$ with values 0.5 , 0.5 , 0.5 and 0.5 , 0.5 , 0.5 . The right oval contains a matrix $w^{(2)}$ with values 0.5 , 0.5 , 0.5 and 0.5 , 0.5 , 0.5 .

Single neuron example



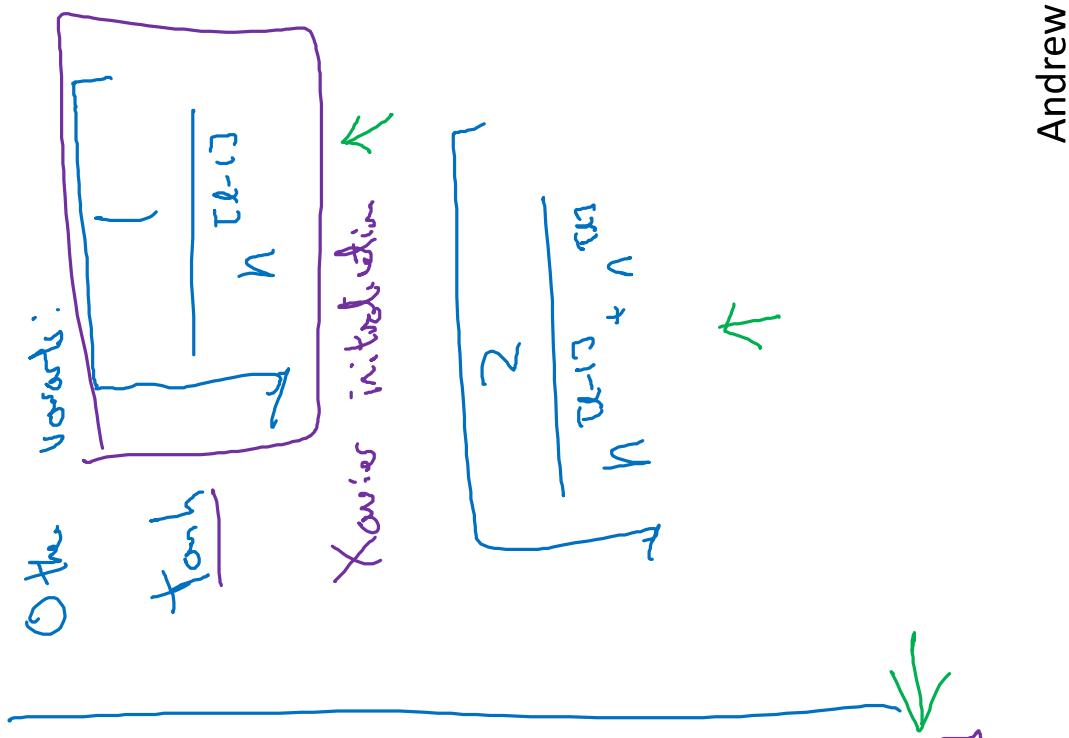
$$z = \sum_{i=1}^n w_i x_i + b$$

~~Large $n \rightarrow$ smaller w_i~~

$$\overbrace{w_i}^{(1)} = \frac{2}{n}$$

$$\overbrace{w_i}^{(2)} = \frac{n}{\sqrt{n}} \cdot \text{random_val} \quad * \overbrace{\text{np.sqrt} \left(\frac{2}{n} \right)}^{\text{ReLU slope}}$$

$\text{ReLU}(z) = \max(0, z)$



Setting up your
optimization problem

Numerical approximation
of gradients

deeplearning.ai



Checking your derivative computation

$$\frac{f(\theta)}{\theta} = \frac{\theta^3}{\theta} = \theta^2$$

Σ

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

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

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

when $\theta = 1$

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

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

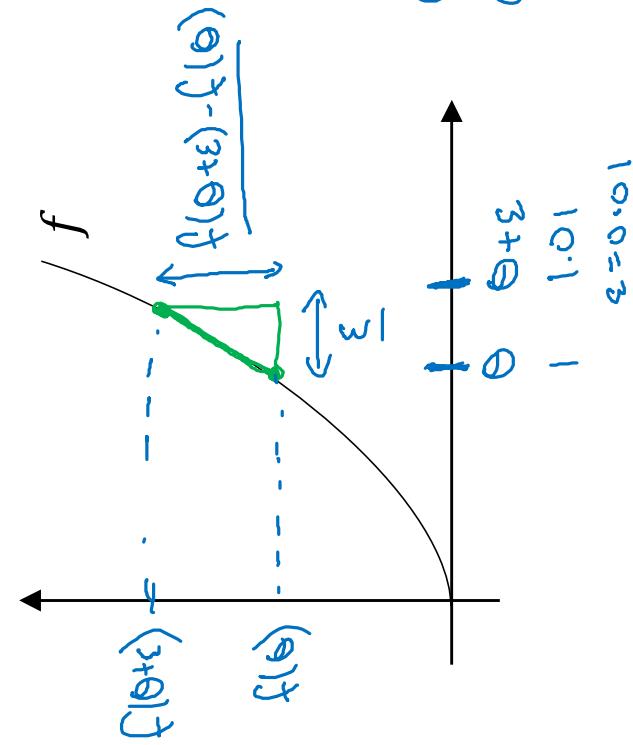
$$3.0301$$

Σ

$$3.1$$

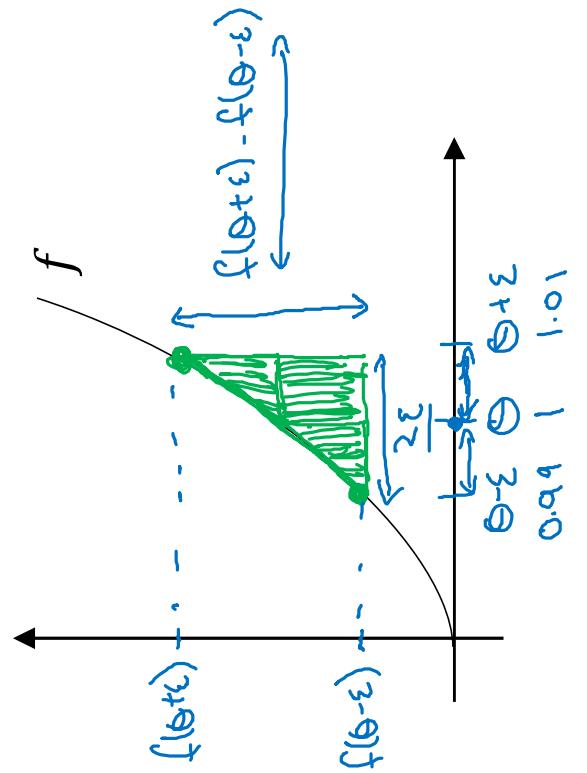
Σ

$$3.2$$



Checking your derivative computation

$$f'(\theta) = \frac{\theta^2}{2}$$



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

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

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

approx error: 0.0001
(prev slide: 3.0301. error: 0.03)

$f(\theta) = \frac{H(\theta + \epsilon) - H(\theta - \epsilon)}{2\epsilon}$

$$\approx \frac{0.01}{0.0001} = 100$$

err: 0.01
err: 0.0001

Andrew Ng

Setting up your
optimization problem

Gradient Checking



deeplearning.ai

Gradient check for a neural network

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

concatenate
 $\mathcal{J}(\underbrace{w^{[1]}, b^{[1]}, \dots, w^{[L]}, b^{[L]}}_{\theta}) = \underline{\mathcal{J}(\theta)}$

Take $\underbrace{dW^{[1]}, db^{[1]}, \dots, dW^{[L]}, db^{[L]}}$ and reshape into a big vector $\underline{\mathcal{J}(\theta)}$.

concatenate
Is $\mathcal{J}(\theta)$ the gradient of $\mathcal{J}(\theta)$?

Gradient checking (Grad check)

for each i :

$$\rightarrow \frac{\partial \theta_{\text{approx}}[i]}{\rightarrow}$$

$$\begin{aligned} \text{for each } i: & \quad \uparrow \\ & \quad J(\theta_1, \theta_2, \dots, \theta_i + \varepsilon, \dots) - J(\theta_1, \theta_2, \dots, \theta_i - \varepsilon, \dots) \\ \rightarrow \frac{\partial \theta_{\text{approx}}[i]}{\rightarrow} = & \quad \downarrow \\ & \quad \approx \varepsilon \end{aligned}$$

$$\approx \frac{\partial \theta[i]}{\partial \theta^i} = \frac{\partial J}{\partial \theta^i}$$

$$\approx \frac{\partial \theta_{\text{approx}}}{\partial \theta^i} \approx \frac{\partial \theta}{\partial \theta^i}$$

Check

$$\rightarrow \frac{\| \theta_{\text{approx}} - \theta \|^2}{\| \theta_{\text{approx}} \|_2 + \| \theta \|^2} \approx$$

$$\boxed{\begin{array}{l} \approx 10^{-7} - \text{rest!} \\ \downarrow \\ 10^{-5} \end{array}}$$

$$\varepsilon = 10^{-7}$$

$$\rightarrow 10^{-3} - \text{worse.} \quad \downarrow$$

Implementation notes

Gradient Checking

Setting up your
optimization problem

deeplearning.ai



Gradient checking implementation notes

- Don't use in training – only to debug
 $\frac{\partial \text{output}^{[l]}}{\partial \theta}$ ← → $\frac{\partial \text{out}}$
- If algorithm fails grad check, look at components to try to identify bug.
 $\frac{\partial \text{out}}{\partial w^l}$
- Remember regularization.
 $J(\theta) = \frac{1}{m} \sum_i f(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2m} \sum_l \|\omega^{(l)}\|_F^2$
 $\frac{\partial \theta}{\partial \theta} = \text{grad of } J \text{ wrt. } \theta$
- Doesn't work with dropout.
 $\frac{\text{keep-prob}}{1.0}$
- Run at random initialization; perhaps again after some training.
 $\frac{w, b \approx 0}{}$