



Train / Dev / Test sets

- Applied ML is highly iterative
- Make sure dev and test come from same distribution (this is hard, some people scrap website for more training data, but training distribution may be different)
- Okay not to have test set -> unbiased estimate of model performance (train/dev is okay,)

Bias / Variance

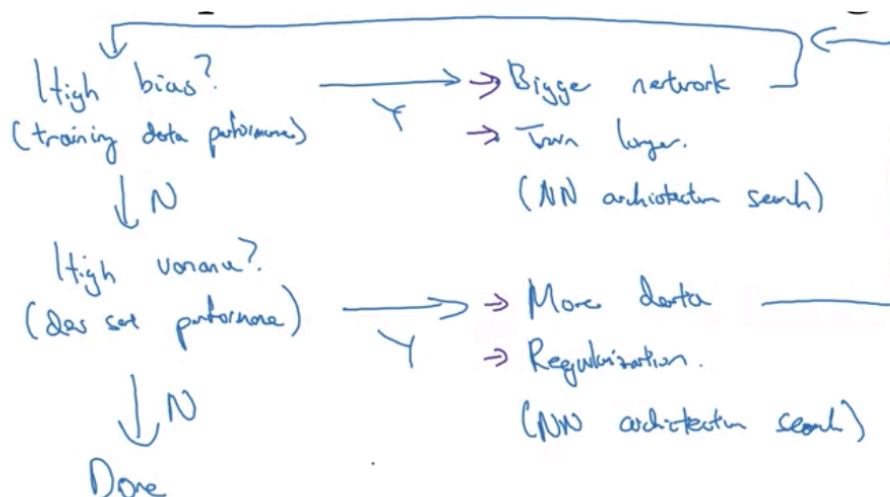
Bias and Variance

Cat classification

	$y=1$	$y=0$	
			
Train set error:	1%	15%	15%
Dev set error:	11%	16%	30%
	high variance	high bias	high bias & high variance
			low bias & low variance
Human: 0%			

- Optimal error: bases error => usually 0%
- The difference between train set error and dev set error give you good estimate on variance

Basic Recipe for machine learning



- Note: Pre DL era, bias variance trade off, now you can use DL decrease either one without hurting the other.

Regularization -> prevent overfitting, or you can just use more data

Logistic regression

$$\min_{w,b} J(w,b) \quad \underline{w \in \mathbb{R}^{n_x}}, \underline{b \in \mathbb{R}}$$

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

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

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

$$L_1 \text{ regularization} \quad \frac{\lambda}{2m} \sum_{i=1}^{n_x} |w_i| = \frac{\lambda}{2m} \|w\|_1$$

w will be sparse

- Lambda - regularization parameter, often set using dev set.
- L2 norm used much more often

Neural network

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

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

$$w: \begin{pmatrix} n^{[l+1]} & n^{[l]} \end{pmatrix}$$

"Frobenius norm"

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

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

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

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

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

"Waggle decay"

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

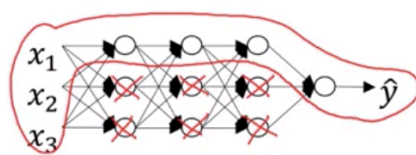
$$\left(1 - \frac{\alpha \lambda}{m}\right) w^{[l]} = \left[w^{[l]} - \left(\frac{\alpha \lambda}{m}\right) w^{[l]} \right] - \alpha (\text{from backprop})$$

Andrew Ng

Initialization

Finally, try "He Initialization"; this is named for the first author of He et al., 2015. (If you have heard of "Xavier initialization", this is similar except Xavier initialization uses a scaling factor for the weights $W^{[l]}$ of $\sqrt{1./\text{layers_dims}[l-1]}$ where He initialization would use $\sqrt{2./\text{layers_dims}[l-1]}$.)

How does regularization prevent overfitting?

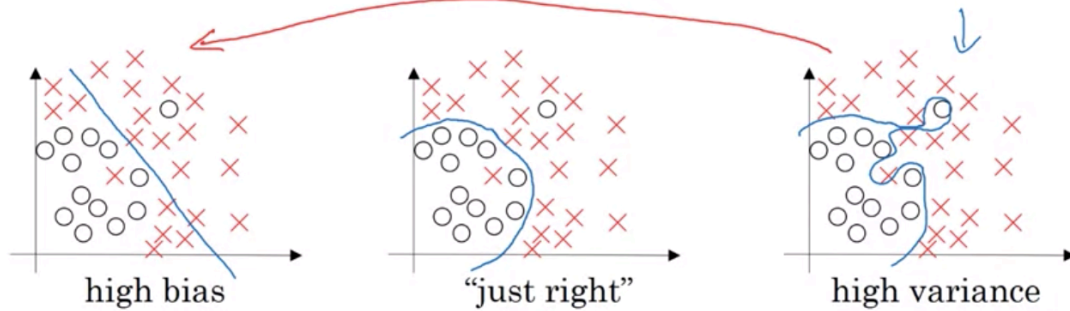


$$J(w^{(l)}, b^{(l)}) = \frac{1}{n} \sum_{i=1}^n L(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{(l)}\|_F^2$$

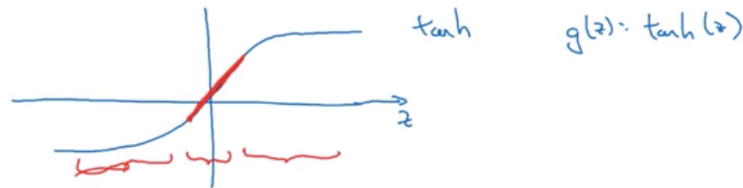
$$w^{(l)} \approx 0$$

Why
Regula
rization
preven
t
overfitti
ng

- Red
uce
varia
nce,
zero
out
the
effec



t of some hidden unit



$$\lambda \uparrow \quad w^{(l)} \downarrow$$

Illustrate with layer $l=3$. keep-prob = $\frac{0.8}{x}$ $\frac{0.2}{x}$

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

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

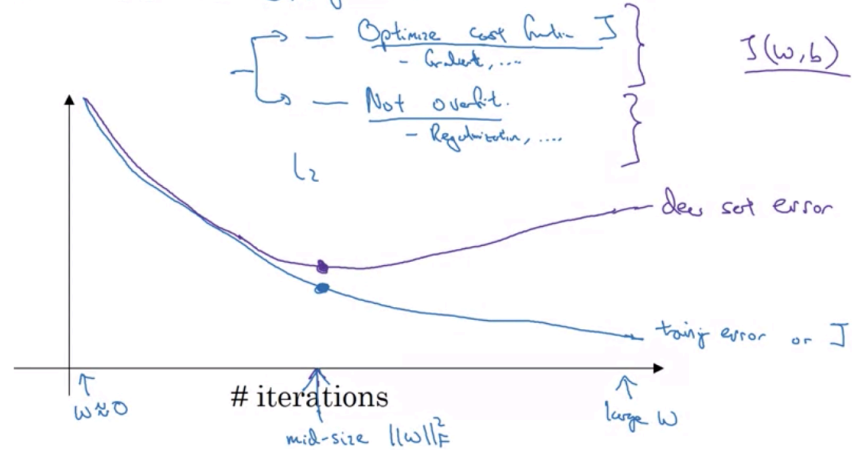
$\rightarrow a3 \neq \text{keep-prob}$

Another intuition

- Restrict the range of activation function (to almost linear)

Early stopping

Dropout



A3 => inverted dropout (remove scaling problem)

- Making prediction => no dropout

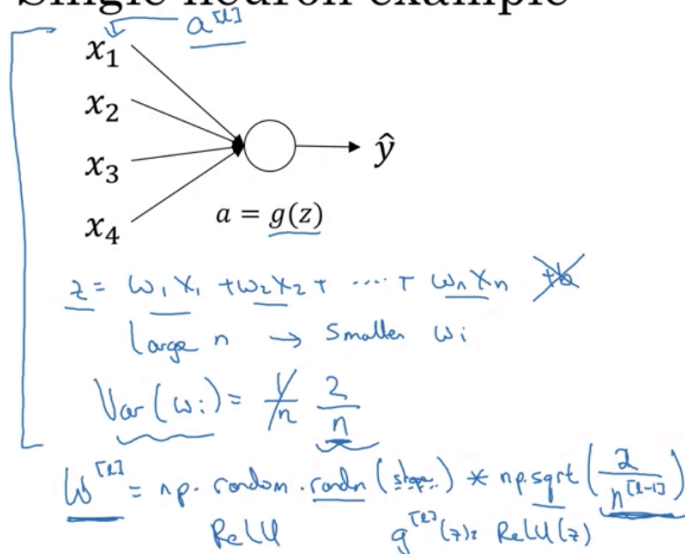
Intuition of Dropout

- Can't rely on any one feature, so have to spread out weights. ~> shrink weights
- Kinda similar to L2 norm but not exactly

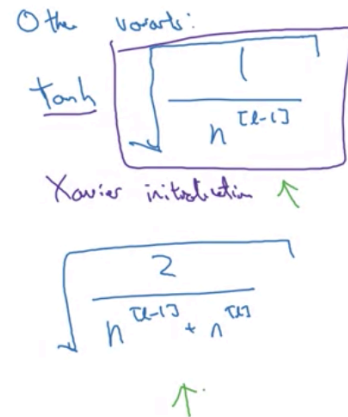
Other methods

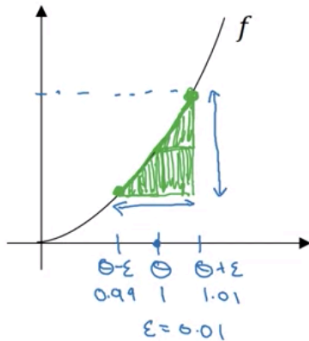
- Augmentation
- Early stopping

Single neuron example



(Do not need too tune L2 norm)





Normalizing input features

- Why?
 - Less elongated space
 - Gradient descent less oscillating

$$\frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon} \approx g(\theta)$$

Weight

- Red gradient (relu example)
-

Initialization
use vanishing

^^^^ tanh -> Xavier initialization

$$f'(\theta) = \lim_{\epsilon \rightarrow 0} \frac{f(\theta + \epsilon) - f(\theta - \epsilon)}{2\epsilon} \quad \begin{matrix} \text{O}(\epsilon^2) \\ 0.01 \\ 0.0001 \end{matrix} \quad \left| \quad \frac{f(\theta + \epsilon) - f(\theta)}{\epsilon} \quad \begin{matrix} \text{error: O}(\epsilon) \\ 0.01 \end{matrix}$$

Andrew Ng

Numerical

Approx of Gradients

- Big green angle gives better approximation (2-sided difference)
- Twice as slow, but worth it for accuracy

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

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

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

Is $d\theta$ the gradient of $\mathcal{J}(\theta)$.

Definiton of

limit: with relate to limit in numerical approximation (think about CS370)

Why 2 sided (on the left) is much more accurate

Gradient checking (Grad check)

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

for each i :

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

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

Gradient
Checking ->
debugging
gradient

Goal?

Check $\frac{\|d\theta_{approx} - d\theta\|_2}{\|d\theta_{approx}\|_2 + \|d\theta\|_2} \approx \frac{10^{-7}}{10^{-5}} - \text{great!}$

$\epsilon = 10^{-7}$ $10^{-3} - \text{worry.}$

Gradient checking implementation notes

- Don't use in training – only to debug

$$\frac{d\theta_{approx}[i]}{\uparrow \uparrow} \leftrightarrow \frac{d\theta[i]}{\uparrow}$$

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

$$\underline{db_F^{[n]}} \quad \underline{dw_F^{[n]}}$$

- Remember regularization.

$$J(\theta) = \frac{1}{n} \sum_i f(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2n} \sum_k \|w^{(k)}\|_2^2$$

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

- Doesn't work with dropout.

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

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

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

How
to
check?
Compute
the
Euclidean

Andrew Ng

distance of the 2

