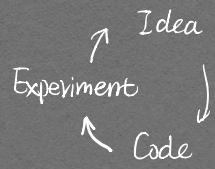


Train / Dev / Test sets

Applied ML is a highly iterative process.



Intuitions are not easily transferrable from NLP to Speech Recognition / CV.

Train / Dev / Test 100 - 1000 - 10000

70% 30%
60% 30% 10%

Train / Dev / Test 1000000

98% 1% 1%
99.5% .4% .1%

Mismatched Train/Test Distribution

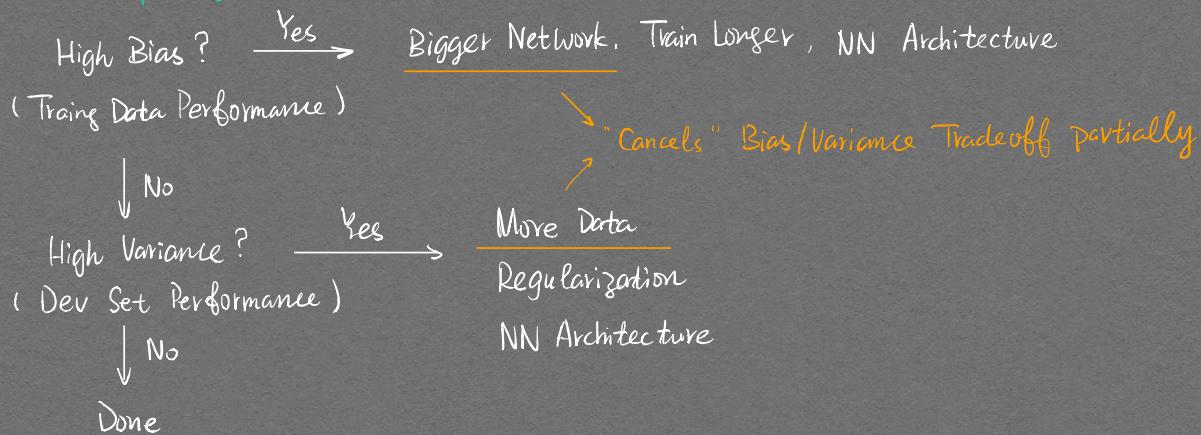
Make sure Dev and Test sets come from the same distribution

Not having a Test set might be okay (Only Dev set) if we don't need an unbiased estimate of model accuracy.

Bias / Variance

	Case 1	Case 2	Case 3	Case 4
Train Set Error	1%	15%	15%	0.5%
Dev Set Error	11%	16%	30%	1%
	High Variance	High Bias	High Bias	Low Bias
Human Error	.5%	1%	High Variance	Low Variance
(Optimal) Bayes Error				

Basic "Recipe" for Machine Learning



Regularizing Your Network

Logistic Regression

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

$$J(w,b) = -\frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \left\{ \begin{array}{l} \boxed{\frac{\lambda}{2m} \|w\|^2} \\ \text{L}_2 \text{ Regularization} \\ \|w\|^2 = \sum_{j=1}^{n_x} w_j^2 = w^T w \end{array} \right.$$

$$J(w,b) = -\frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \left\{ \begin{array}{l} \boxed{\frac{\lambda}{2m} \|w\|_1} \\ \text{L}_1 \text{ Regularization} \\ \|w\|_1 = \sum_{j=1}^{n_x} |w_j| \quad w \text{ will be sparse} \end{array} \right.$$

Neural Network

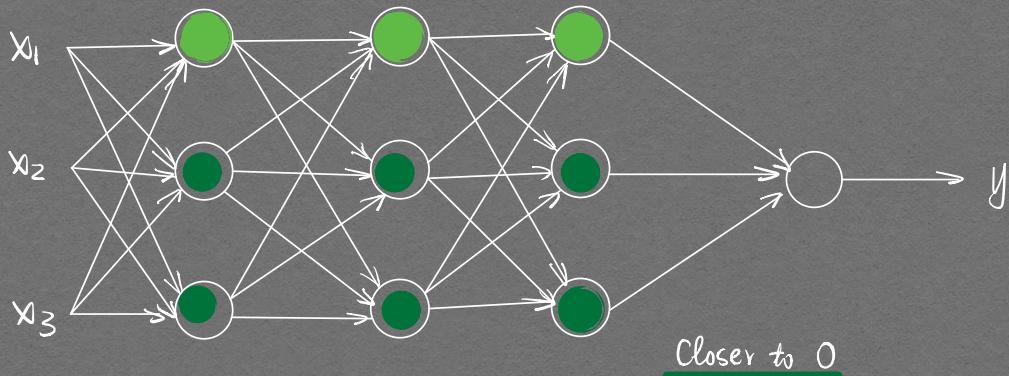
$$J(w^{[0]}, b^{[0]}, \dots, w^{[L]}, b^{[L]}) = -\frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w^{[l]}\|^2$$

of neurons in layer l \rightarrow # neurons in the previous layer $l-1$

$$\text{Frobenius Norm: } \|w^{[l]}\|_F^2 = \sum_{i=1}^{n_l} \sum_{j=1}^{n_{l-1}} (w_{i,j}^{[l]})^2 \quad w^{[l]} \sim n^{[l]} \times n^{[l-1]}$$

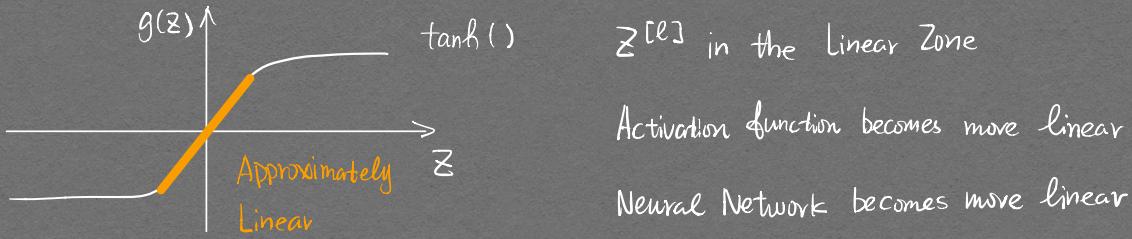
Calculate Gradient: $d w^{[l]} = (\text{from backprop}) + \frac{\lambda}{m} w^{[l]}$

Update Step: $w^{[l]} = w^{[l]} - \alpha d w^{[l]}$ "Weight Decay"

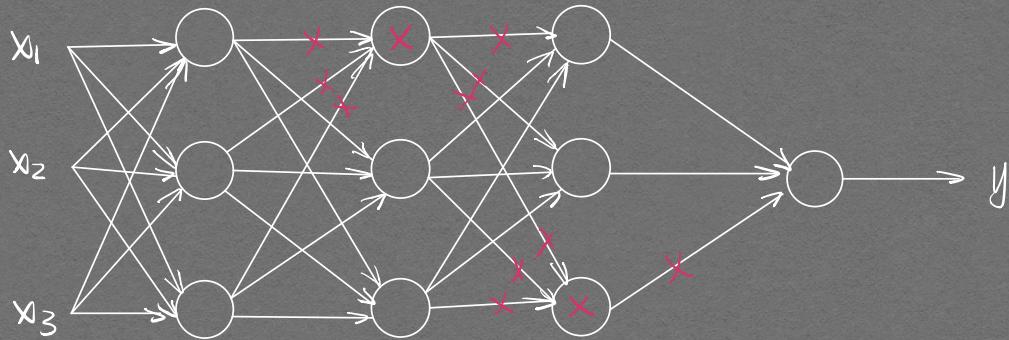


$\lambda \rightarrow \text{Inf}$

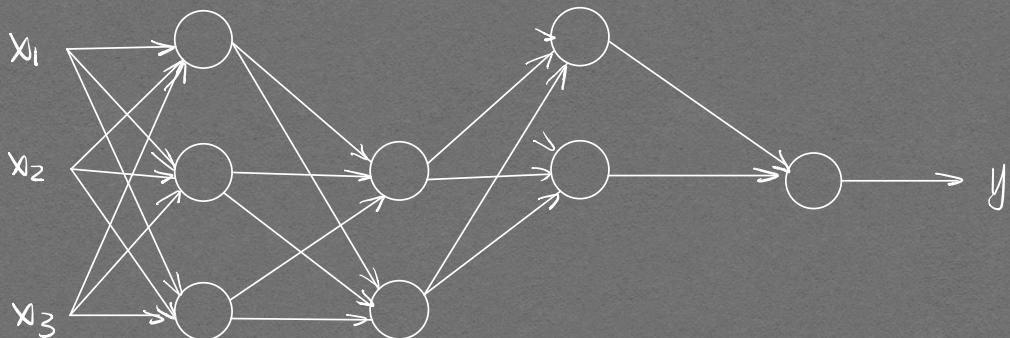
$w \rightarrow 0 \quad z^{[l]} : w^{[l]} a^{[l-1]} + b^{[l]} \rightarrow 0$



Dropout Regularization



kill neurons and their input/output edges



Training Time (Inverted Dropout) : $\text{keep_prob} = .8$

At layer 3 : $d_3 = \text{np.random.rand}(a_3.\text{shape}[0], a_3.\text{shape}[1])$

$< \text{keep_prob}$ A True/False Mask

$a_3 = \text{np.multiply}(a_3, d_3)$ Element-wise Multiplication

$a_3 /= \text{keep_prob}$ Inverted Dropout

to keep expected value of a_3 unchanged

Test Time: No Dropout

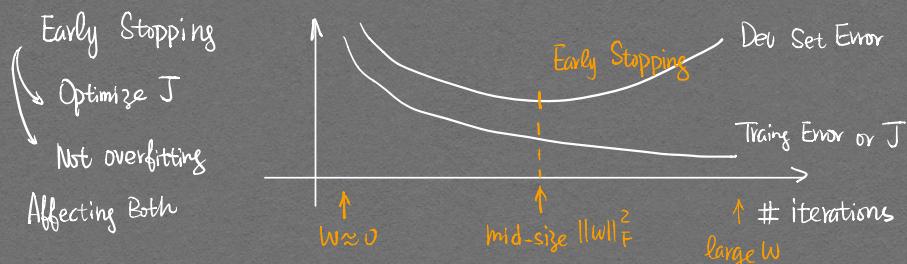
(We don't want random outputs)

Understanding Dropout

- 1) Smaller Network
- 2) Input Features are killed randomly, thus weights are encouraged to spread out
- 3) Different Dropout / keep-prob Ratio in Different layers

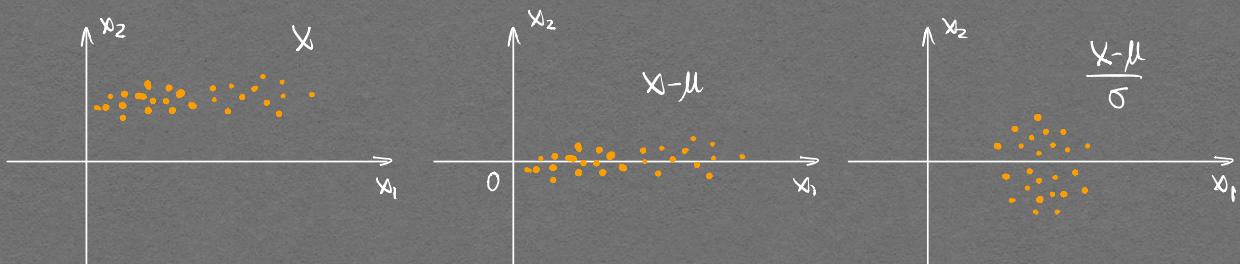
Other Regularization Methods

Data Augmentation (Mirroring Images, Random Cropping Images, ...)

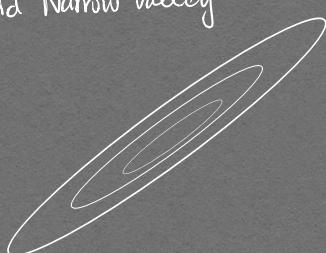


Setting Up Your Optimization Problem

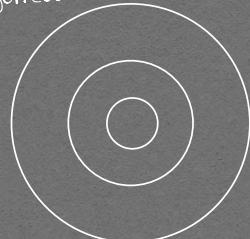
Normalizing Inputs Calculate μ, σ with training set.
Use same μ, σ for test set.



Avoid Narrow Valley



Round Shape is preferred

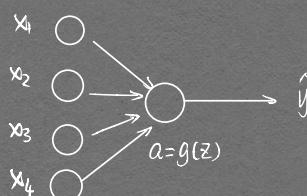


Vanishing / Exploding Gradients

Weight Initialization
Partial Solution

=

Single Neuron Example



$b=0$
 $z = w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4$
 $\text{large } n \rightarrow \text{small } w_i$
 $\text{Var}(w_i) = \frac{1}{n} \left(\frac{2}{n} \text{ works better} \right)$

Other variants $W^{[l]} = \text{np.random.randn(shape)} * \text{np.sqrt}(c)$

If $\tanh()$

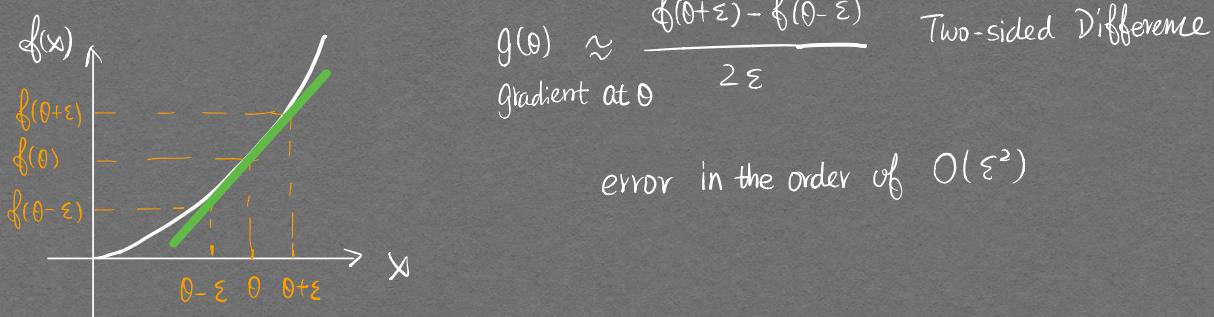
$$\text{use Xavier activation, } c = \sqrt{\frac{1}{n^{[l-1]}}}$$

If $\text{relu}()$

$$\text{use He activation, } c = \sqrt{\frac{2}{n^{[l-1]}}}$$

Numerical Approximation of Gradients

Checking your derivative computation



Gradient Checking

$$\Theta = (w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}, \dots, w^{[l]}, b^{[l]})$$

$$d\Theta = (dw^{[1]}, db^{[1]}, dw^{[2]}, db^{[2]}, \dots, dw^{[l]}, db^{[l]})$$

for each i in $1: \text{len}(\Theta)$:

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

$$\approx d\theta_{[i]} = \frac{\partial J}{\partial \theta_i} \quad \text{pick } \epsilon = 10^{-7}$$

$$\text{check } \frac{\|d\Theta_{\text{approx}} - d\Theta\|_2}{\|d\Theta_{\text{approx}}\|_2 + \|d\Theta\|_2} \approx 10^{-7} - \text{Great!}$$

$10^{-3} - \text{Wrong}$

Implement Gradient Checking

Don't use grad check in training - only to debug

If fails, look at components

Remember regularization term in the cost function when calculating gradients
Doesn't work with dropout. (keep-prob = 1, pass grad check, turn on dropout)
Run at random initialization, perhaps after some training.

Notes from Jupyter Notebook

Zero Initialization \rightarrow Fail to break symmetry

Every neuron in each layer will learn the same thing.
We end up training a neural network with $n^{[l]} = 1$
for every layer. The network is a lot less powerful.

Random Initialization \rightarrow Initialize with very large values
 \rightarrow Vanishing / Exploding Gradients
Slows optimization algorithm

He Initialization \rightarrow Works well for ReLU activation