

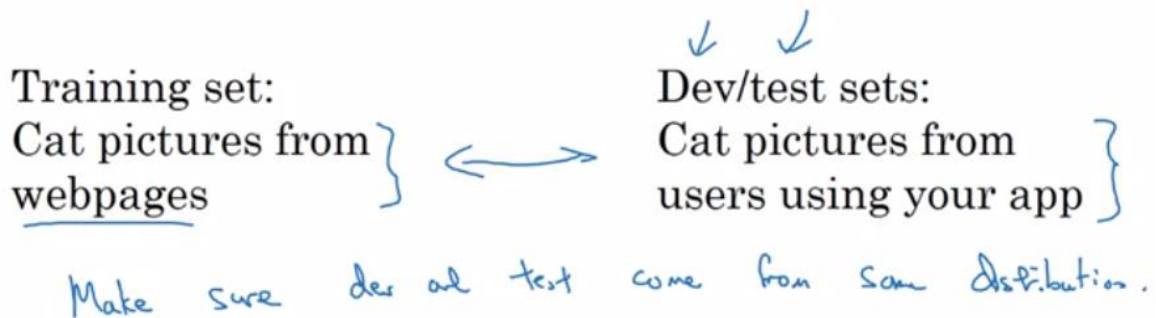
Hyperparameters Tuning

Week 1

Applied ML is a highly iterative process.

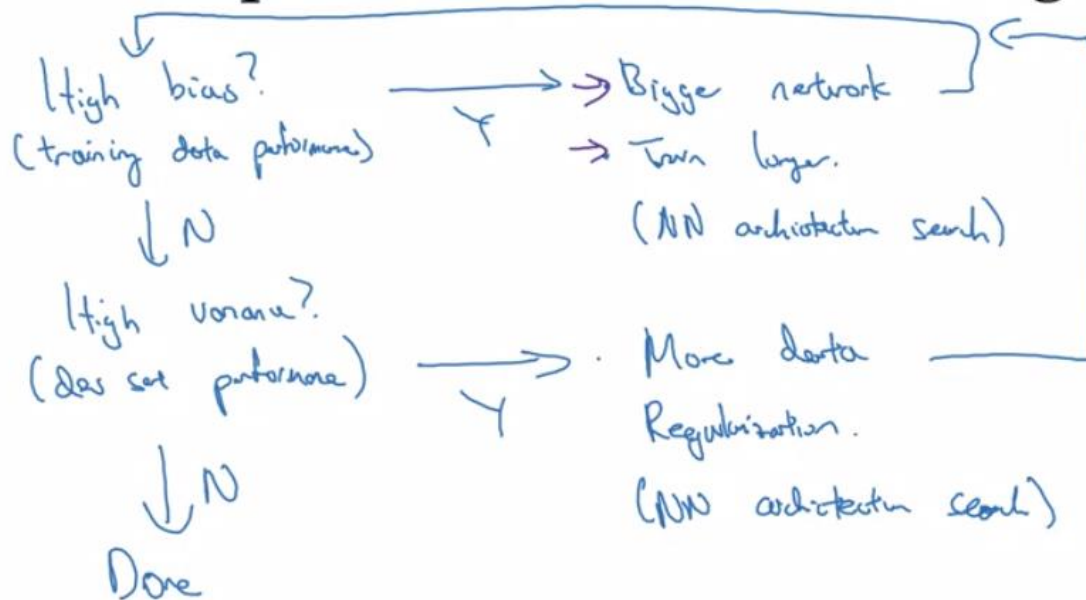
The gathered data is split into training, development(cross-validation) and test set. In the previous era, the split ratio was about 60-20-20. But now, in the Big Data era, the dev and test set require to have a much smaller ratio, much of the data being used in the training set.

Mismatched train/test distribution



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: 20%				

Basic recipe for machine learning



L_2 regularization $\|w\|_2^2 = \sum_{j=1}^n w_j^2 = w^T w$

OMIT

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

w will be sparse

Lambda = Regularization Parameter (given by Dev Set)

Neural network

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

Frobenius Norm of a Matrix

$$dW^{[L]} = \left[\text{(from backprop)} + \frac{\lambda}{n} W^{[L]} \right]$$

$$\rightarrow W^{[L]} := W^{[L]} - \alpha dW^{[L]}$$

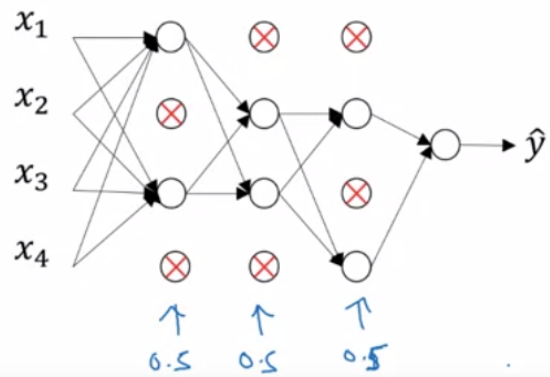
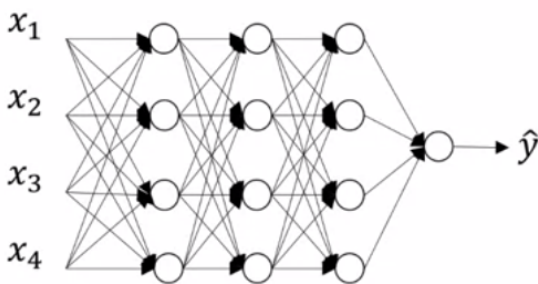
"Weight decay"

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

$$= \left(1 - \frac{\alpha \lambda}{n}\right) W^{[L]} - \alpha \text{(from backprop)}$$

9:42

Dropout regularization



Implementing dropout ("Inverted dropout")

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

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

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

$$\rightarrow \boxed{a3 /= \text{keep-prob}}$$

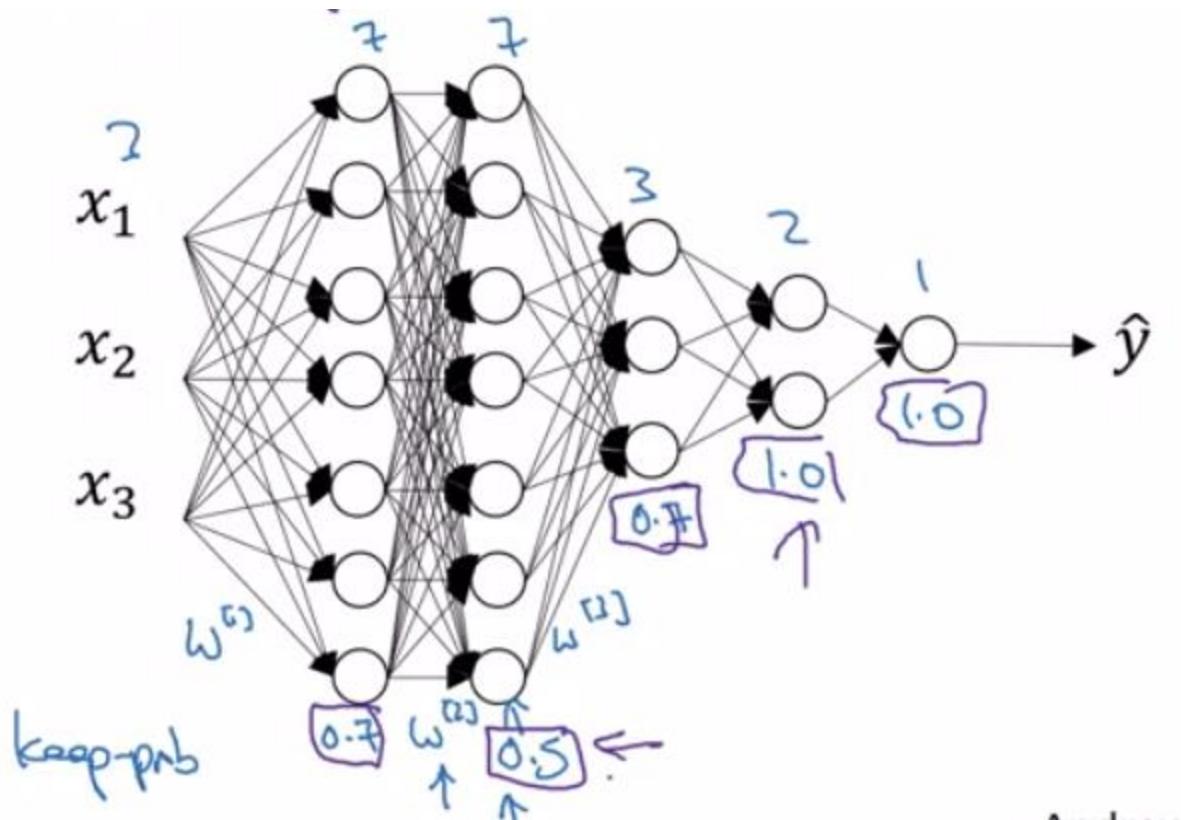
50 units. \leadsto 10 units shut off

$$z^{[4]} = w^{[4]} \cdot \underbrace{a^{[3]}}_{\substack{\uparrow \\ \text{reduced by } 20\%}} + b^{[4]}$$

$$/= \underline{0.8}$$

Why does drop-out work?

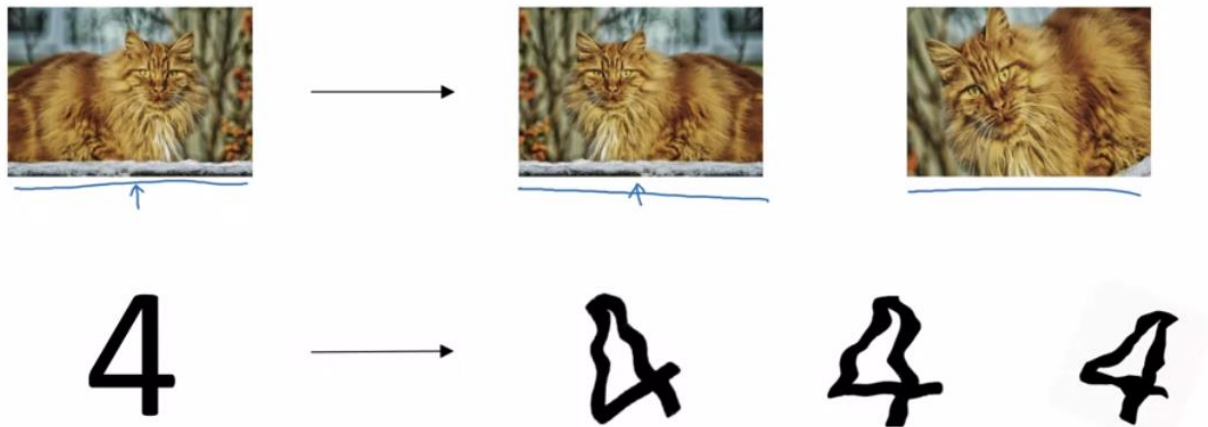
Intuition: Can't rely on any one feature, so have to spread out weights. \leadsto Shrink weights. $\frac{1}{2}$



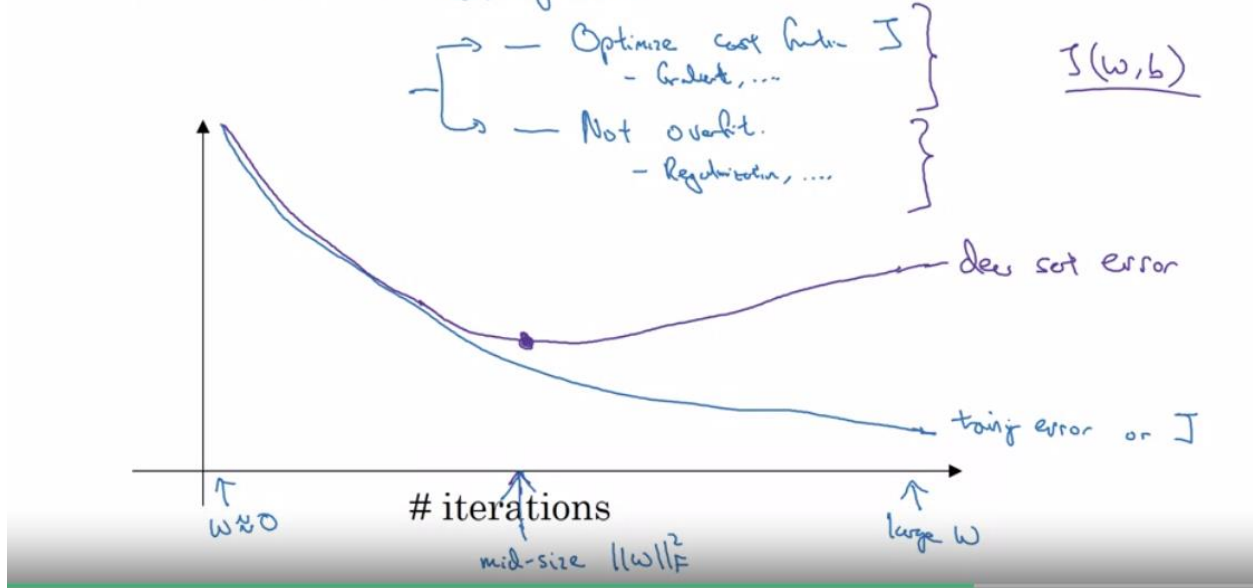
This technique is usually used in Computer Vision.

One disadvantage of the Dropout Regularization is that now the Cost Function is not “well-defined”, and plotting J over Iterations may not always result in a smooth decreasing graph.

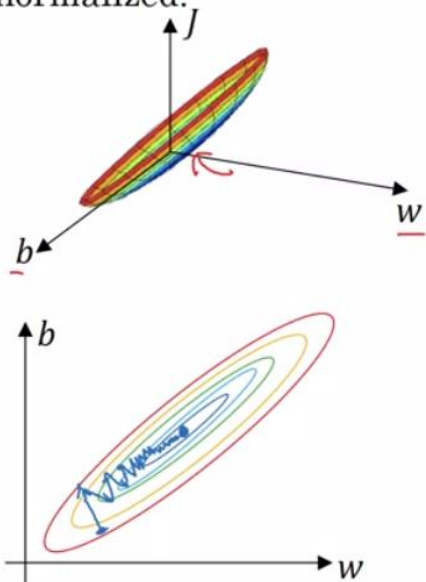
Data augmentation



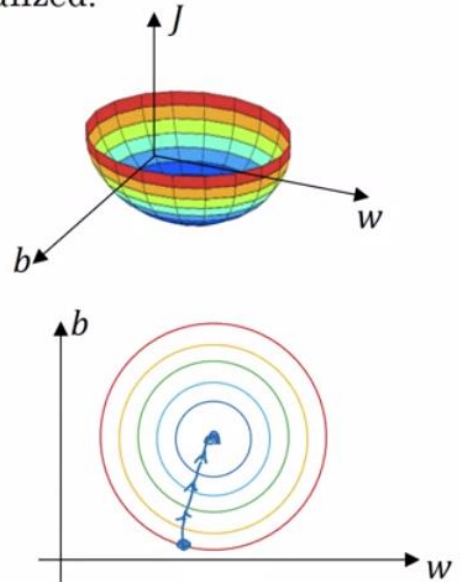
Early stopping



Unnormalized: $w_2, x_2, 0, \dots, 1$



Normalized:



Vanishing/Exploding Gradients Problem

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

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

Other variants:

tanh

Xavier initialization ↑

↑

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

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

$\underline{\Sigma = 10^{-7}}$

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 $\text{sqrt}(1./\text{layers_dims}[1-1])$ where He initialization would use $\text{sqrt}(2./\text{layers_dims}[1-1])$.)

He initialization works well for networks with RELU activations.

With dropout, your neurons thus become less sensitive to the activation of one other specific neuron, because that other neuron might be shut down at any time.

- A **common mistake** when using dropout is to use it both in training and testing. You should use dropout (randomly eliminate nodes) only in training.
- Deep learning frameworks like [tensorflow](#), [PaddlePaddle](#), [keras](#) or [caffe](#) come with a dropout layer implementation. Don't stress - you will soon learn some of these frameworks.

Note that regularization hurts training set performance! This is because it limits the ability of the network to overfit to the training set. But since it ultimately gives better test accuracy, it is helping your system.