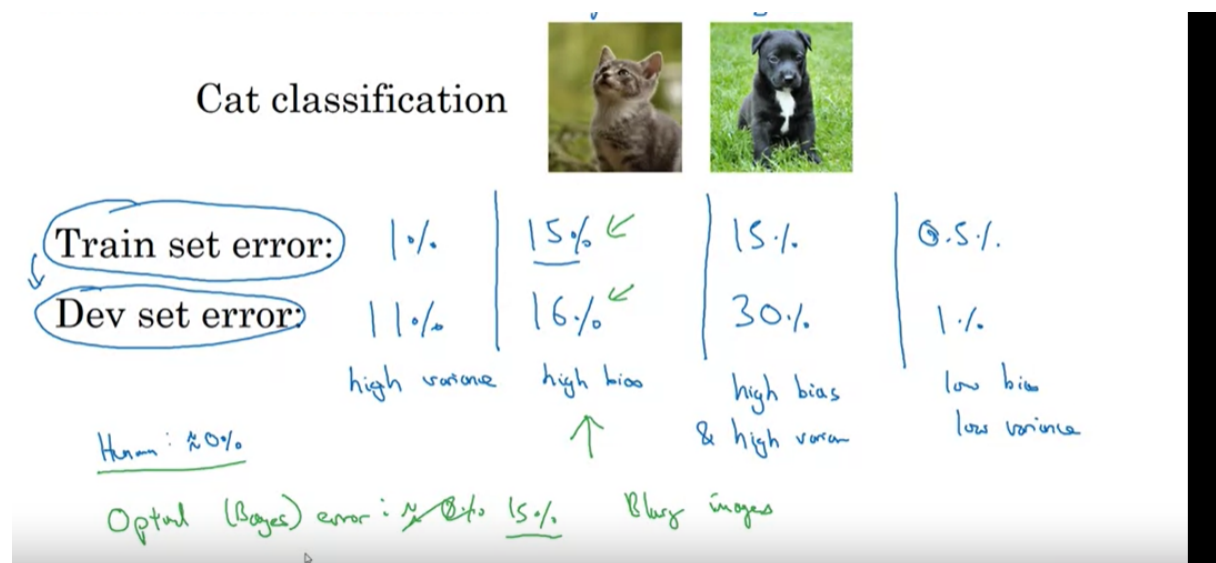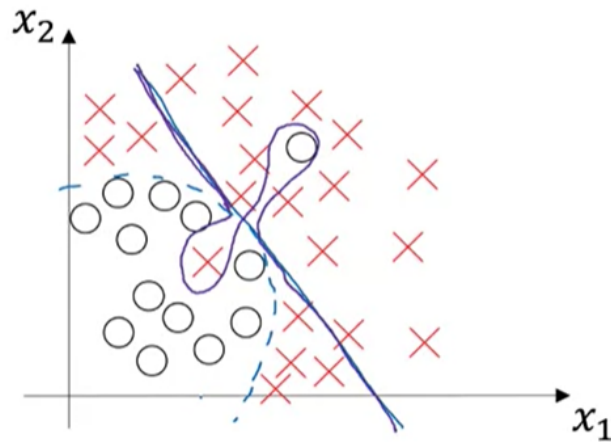# dL.ai_Improvements in Nueral Network

## a)Bias-Variance



- High bias means we assume a lot of things about the classifier, we dont consider eenough dimensions.

- High variance means the data is too sensitive to each individual instance in the training data

- They are not antogonistic to each other

- A classifier having high bias and variance will look like the purple line

# b) Regularization

## i)L1 and L2

- L1 and L2 tries to prevent overfitting by penalizing the model for having high weights..

- We do this by adding weights into the loss function.

- example for logistic loss

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

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

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

- where lambda is the regularization parameter.

- In nueral network $\|w\|$ by the frobenius norm gives

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

- during backpropagation when we find dw an additional term $\frac{\lambda}{m}w$ is also considered this leads to the magnitude of weight decreasing.Henc4e l2 is also called weight decay.
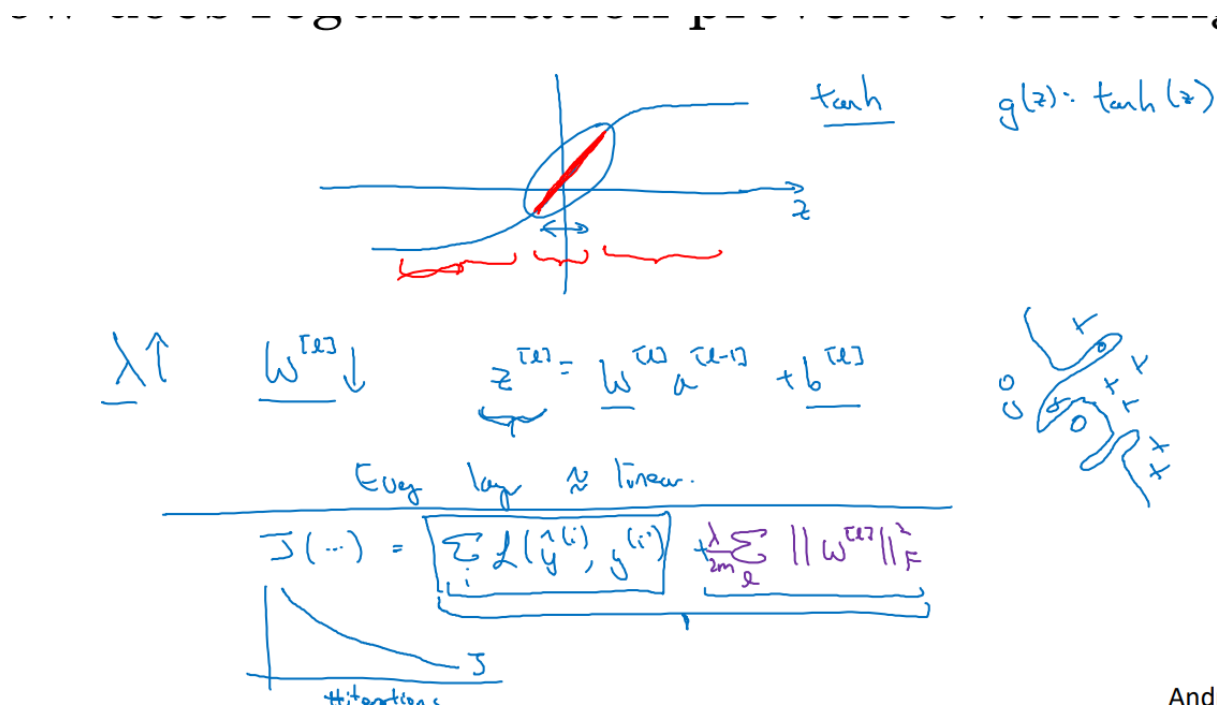


$$dw^{[l]} = \boxed{(\text{from backprop}) + \frac{\lambda}{m}w^{[l]}}$$
$$w^{[l]} := w^{[l]} - \alpha\, dw^{[l]}$$

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

decay

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

Ar

- L1 regularization leads to sparsity in the model.It completely makes some of the weights zero which leads to some nodes being always inactivated leading to a smaller model.I hope this because when we

$$J = CE + \frac{\lambda}{m}|w|$$
$$dw = (bpterm) + \frac{\lambda}{m}$$
$$w^{[l]} - dw^{[l]} = w^{[l]} - a(bpterm) - a\frac{\lambda}{m}$$

- let $w^{[l]}$ be [1,2,3,4,5] then each weight in w will be decreased by the same constant value till they reach zero or no of epochs complete.

- However in L2 since $\frac{\lambda}{m}w$ is proprtional to weight the extreme values are penalised more and the lower values are penalised lesser.It tends to make all the weights smaller but not exactly zero.

- Even though variance and bias are not perfectly complementary we can say that if a model is closer to a linear regressor it would decrease the variance

- so if value of w is less the value of z would be less, if the value of z is less then sigmoid function will more or less converge to a linear y-x graph. Hence the whole model comes closer to a linear regressor.

$tanh$  $g(z) = tanh(z)$

$\lambda \uparrow$    $W^{[l]} \downarrow$    $z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]}$

Every layer $\approx$ linear.

$J(\cdots) = \left[ \sum_i \mathcal{L}(\hat{y}^{(i)}, y^{(i)}) \right] + \frac{\lambda}{2m} \sum_l \| W^{[l]} \|_F^2$

#iterations    $J$

## ii) Inverted Dropout

- During training in a given layer all the nodes are not taken a fixed ratio of nodes are always dropped out.

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

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

$a3$ = np.multiply(a3, d3)         # a3 *= d3.

$\rightarrow$ $a3 /= \cancel{0.8}$ keep-prob $\leftarrow$

50 units. $\rightsquigarrow$    10 units shut off

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

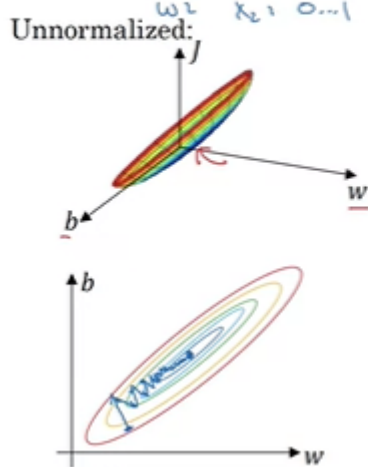$J$         $\uparrow$ reduced by 20%.      Test

$/= 0.8$

- Understand how the value of z doesnt change even though though several a3 i have become zero

- Dropout reduces the dependence of a model on a single node

- Dropout make Loss function more vague

- Drop out ensures that the network cant depend entirely on any one feature this ensures that the weights are spread out more. This has a similar effect to L2 regularization

- During each iteration of training the input data only passes through a smaller nn.
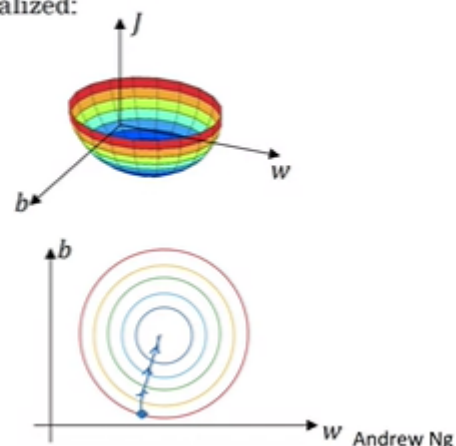
# c)Normalizing



Why normalize inputs?

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

Unnormalized:

Normalized:

Andrew Ng

# d)Vanishing/Exploding Gradient

# Vanishing/exploding gradients



$x_1$

$x_2$

$w^{[1]}$ $w^{[2]}$ $w^{[3]}$ ..... $w^{[L]}$ $\hat{y}$

$L$

$g(z) = z$. $\quad b^{[\ell]} = 0$.

$\hat{y} = w^{[L]} \left( w^{[L-1]} \right) \left( w^{[L-2]} \right) (\cdots) \boxed{w^{[3]} \; w^{[2]} \; w^{[1]} \; x}$ $a^{[3]}$ $\quad 1.5^{L}$

$\quad 0.5^{L}$

$w^{[\ell]} > I$

$w^{[\ell]} < I \quad \begin{bmatrix} 0.9 & \\ & 0.9 \end{bmatrix}$

$z^{[1]} = w^{[1]} x$

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

$a^{[2]} = g(z^{[2]}) = g(w^{[2]} a^{[1]})$

$w^{[\ell]} = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \\ 0.5 \end{bmatrix}$ $\quad \hat{y} = w^{[L]} \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \\ & 0.9 \end{bmatrix}^{L-1} x$ $\quad 1.5^{L-1} x$

$\quad 0.5^{L-1} x$

Andrew Ng