deeplearning.ai

# Regularizing your neural network

---

# Regularization

# Logistic regression

$\min\limits_{w,b} J(w, b)$

$\rightarrow w \in \mathbb{R}^{n_x}, \; b \in \mathbb{R}$

- Getting more data can also $\downarrow$ overfitting

$J(w, b) = \frac{1}{M} \sum\limits_{i=1}^{M} L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2M} \|w\|_2^2$

<u>L2 Regularization</u>

$\|w\|_2^2 = \sum\limits_{j=1}^{n_x} w_j^2 = w^T w$

Q why not regularize "b" as well? $\rightarrow \frac{\lambda}{2M} b^2$

well you're doing this to $\downarrow$
overfitting of your params
params are w, b
$\dim(w) = (n_x \times 1)$
$\dim(b) = (1, \times 1)$
$\Rightarrow$ the overfitting would likely
be happening due to w & not
b [more params to tweak in w]

<u>L1 Regularization</u>

This is when we add $\frac{\lambda}{2M} \|w\|_1$

ie, $\frac{\lambda}{M} \sum\limits_{j=1}^{n_x} |w_j|$ to the eqn

If L1 is used, w ends up being
sparse $\Rightarrow$ many 0s in w
$\Rightarrow$ helps in storage of params
as most values are 0s, so you
need to store only non-zeros
& their positions

$\boxed{L_2 \text{ is used more often than } L1}$

$\lambda$ = Regularization param
$\hookrightarrow$ find value via dev set

Andrew Ng

# Neural network

$$J(w^{[1]}, b^{[1]}, \dots w^{[L]}, b^{[L]}) = \frac{1}{M} \sum_{i=1}^{M} L(\hat{y}^{(i)}, y^{(i)}) + \frac{\lambda}{2M} \sum_{j=1}^{L} \| w^{[j]} \|_2^2$$

$\underbrace{\phantom{w}}_{\text{# layers in NN}}$

$$\text{for each} \rightarrow \| \underline{w^{[j]}} \|^2 = \sum_{i=1}^{n^{[L-1]}} \sum_{k=1}^{n^{[L]}} \left( w_{ik}^{[j]} \right)^2 \left.\right\} \text{For each } w \text{ at layer } j, \text{ find } \Sigma \text{ of } Sq, \text{ then do for all layers}$$

for each $j$th layer

aka
"Forbenious Norm"

<span style="color:red">→ why are the limits $n^{[L]}$ & $n^{[L-1]}$

Remember dim(w) from before : $\left( n^{[L]}, n^{[L-1]} \right)$

#hidden units in current layer → #hidden units in prev layer</span>

So now how does Back prop change
w/ the Addition of Regularization?
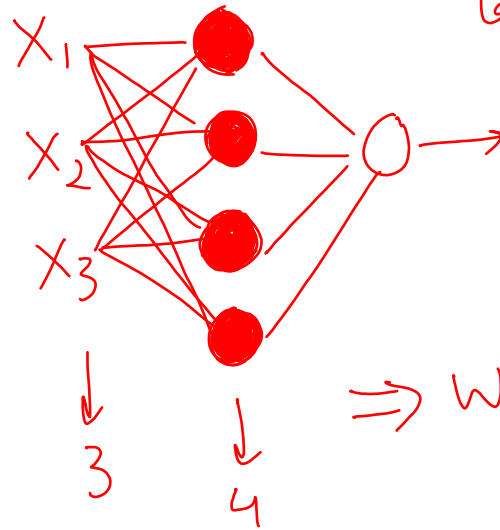
$$\frac{\partial J}{\partial w^{[L]}} = ?$$

Now, $\underline{dw^{[L]}} = \text{old } dw + \frac{\lambda}{M} w^{[L]}$

$\underline{\underline{\text{update}}}$

$$w^{[L]} = w^{[L]} - \alpha \cdot \underline{\underline{dw^{[L]}}}$$

The new $dw^{[L]}$ makes the
update further ↓ by

$$\underline{\underline{\frac{\alpha \lambda}{M}}} w^{[L]}$$

∴ we call gradient descent w/ Regularization "WEIGHT DECAY"

<span style="color:red">$X_1$

$X_2$

$X_3$
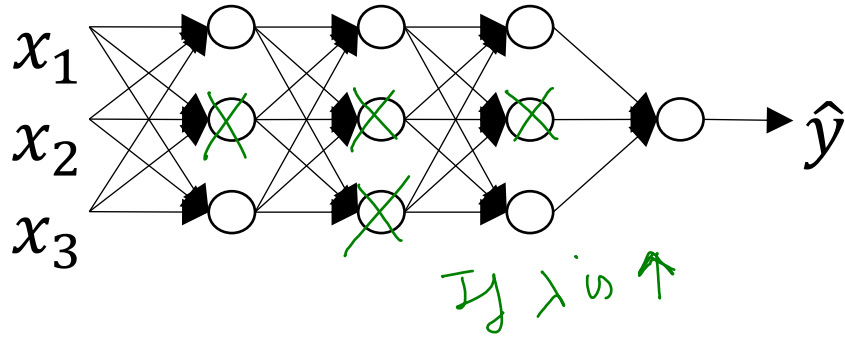
↓       ↓
3       4

$\Rightarrow w^{[1]} = (4 \times 3)$</span>

<span style="color:black">Andrew Ng</span>

# How does regularization prevent overfitting?



$$J\left(w^{[l]}, b^{[l]}\right) = \frac{1}{M} \sum L\left(y^{(i)}, \hat{y}^{(i)}\right) + \frac{\lambda}{2M} \sum_{i=1}^{L} \|w^{[l]}\|_F^2$$

All layers

why does this ↓ overfitting?

If λ is high, then many weights w in diff layers become close to 0 to ↓ cost func ⟹ Simpler model

If λ too high, we get high bias model

If λ is appropriate, we get middle Model

$x_1$

$x_2$

$x_3$

$\hat{y}$

If λ is ↑

High bias

High Variance

Andrew Ng

# How does regularization prevent overfitting?

Assume, we use <u>tanh</u> Activation func



$g(z) = a$

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

Now If we $\uparrow \lambda$, $w^{[l]} \downarrow$

If $w^{[l]} \downarrow$  $z^{[l]} \downarrow$, why?    $z^{[l]} = w^{[l]} \cdot a^{[l-1]} + b^{[l]}$

$$\Rightarrow z \propto w$$

If $z \downarrow$ and lets assume $z$ goes from being
in Range $C \to D$ or Range $A \to B$
to Range $\underline{\underline{B \to C}}$

then we can see derivative of $z$ "$dz$"   $\nearrow g'(z)$
is linear, $\Rightarrow$ every layer is computing a linear
eqn $\Rightarrow$ model has become simpler

$$\Rightarrow \text{If } \lambda \uparrow, \text{ we get simpler}$$
model

Andrew Ng