## Bias / Variance

High Bias ?
(training data performance) $\quad$ Y $\longrightarrow$ Bigger DataSet
$\quad\quad\quad\quad$ $\longrightarrow$ train longer
$\quad\quad\quad$ ↓ N

High Variance ?
( validation set performance) $\rightarrow$ More data
$\quad\quad\quad\quad\quad\quad\quad$ $\rightarrow$ Regularization
$\quad\quad\quad\quad$ ↓ N $\quad\quad\quad$ Y

Done

## Regularization

logistic Regression

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

$L_2$ Regularization: $\|w\|_2^2 = w^T w$

$L_1$ Regularization: $\frac{\lambda}{2m} \sum_{j=1}^{n_x} |w_j| = \frac{\lambda}{2m} \|w\|_1$ $\quad$ } w will be sparse

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$ ⇓
$\quad\quad\quad\quad\quad\quad\quad\quad\quad$ help a little to compress model

## Neural Network

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

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

$w \quad\quad n^{[l-1]} \times n^{[l]}$

"Frobenius norm" $\quad\quad$ ↑ # of hidden ↑

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

$$\rightarrow w^{[l]} = w^{[l]} - \alpha \, dw^{[l]} \qquad \Rightarrow \text{from backprop}$$

$$\Rightarrow \text{Weight Decay}: \quad w^{[l]} = w^{[l]} - \alpha \, [(\text{from backprop})$$
$$+ \frac{\lambda}{m} w^{[l]} \,]$$
$$= w^{[l]} - \frac{\partial \lambda}{m} w^{[l]} - \partial (\text{from back prop})$$
$$= (1 - \frac{\partial \lambda}{m}) w - \partial (\text{from backprop})$$

<u>always gets</u> smaller

## Why Regularization Works?



tan h

$z^{[l]}$

$$x \uparrow \qquad w^{[l]} \downarrow \qquad z^{[l]} = w^{[l]} a^{[l-1]} + b^{[l]}$$

force the weight to be closer to linear.
therefore reduce the complexity of Network.

## DropOut Regularization (No Dropout at test time)
randomly turning off hidden unit during training.

Implementation: "inverted dropout"

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

$$\nwarrow \text{ zero out } 20\%$$
$$\text{then } /= 0.8$$

## Why Dropout Works?

Intuition: Can't rely on one feature , so have
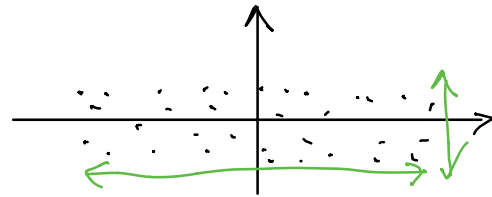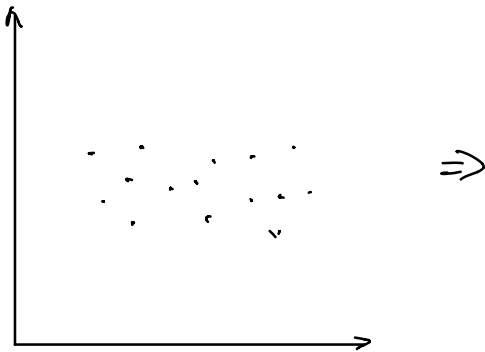to spread out weights. ⟿ shrink
weights
Similar to L2.

[ When debugging gradient descent: turn off
dropout and see if gradient monotonically
decreasing )

Other method:
⎰ Data augmentation
⎱ Early Stopping → orthognolization

# Normalizing Input

subtract mean

$$\mu = \frac{1}{m} \sum_{i=1}^{m} x^{(i)}$$

$$x -= \mu$$

Normalize Variance

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^{m} x^{(i)} ** 2$$

$$x /= \sigma^2$$