# Logistic regression

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

$$\min_{w,b} J(w, b)$$

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

omit

$L_2$ regularization

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

$L_1$ regularization

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

$w$ will be sparse.

Andrew Ng

# Neural network

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

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

$$w: (n^{[l]}, n^{[l-1]})$$

"Frobenius norm"

$$\|\cdot\|_2^2 \qquad \|\cdot\|_F^2$$

$$dw^{[l]} = (\text{from backprop})$$

$$w^{[l]} := w^{[l]} - \alpha \, dw^{[l]}$$

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

Andrew Ng

# Neural network

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

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

"Frobenius norm"

$$\|\cdot\|_2^2 \qquad \|\cdot\|_F^2$$

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

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

$$\rightarrow w^{[l]} := w^{[l]} - \alpha \, dw^{[l]}$$

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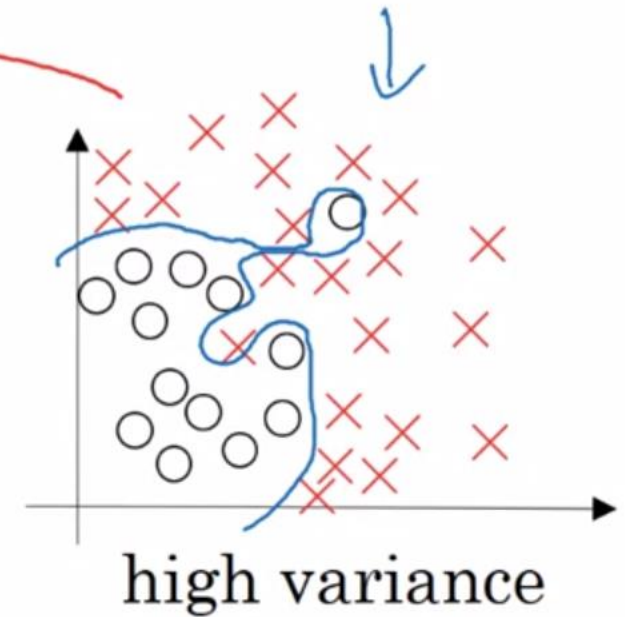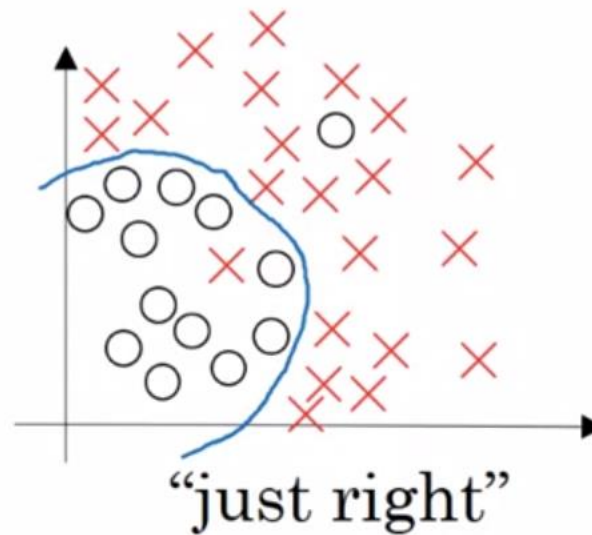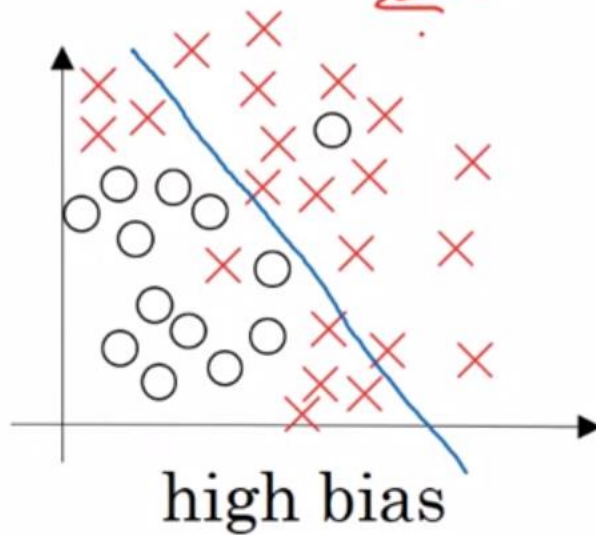
$$w^{[l]} := w^{[l]} - \alpha \left[ (\text{from backprop}) + \frac{\lambda}{m} w^{[l]} \right]$$

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

Andrew Ng

# How does regularization prevent overfitting?



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

$$w^{[l]} \approx 0$$

high bias          "just right"          high variance

# How does regularization prevent overfitting?

$tanh$

$g(z) = tanh(z)$

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

Every layer $\approx$ linear.

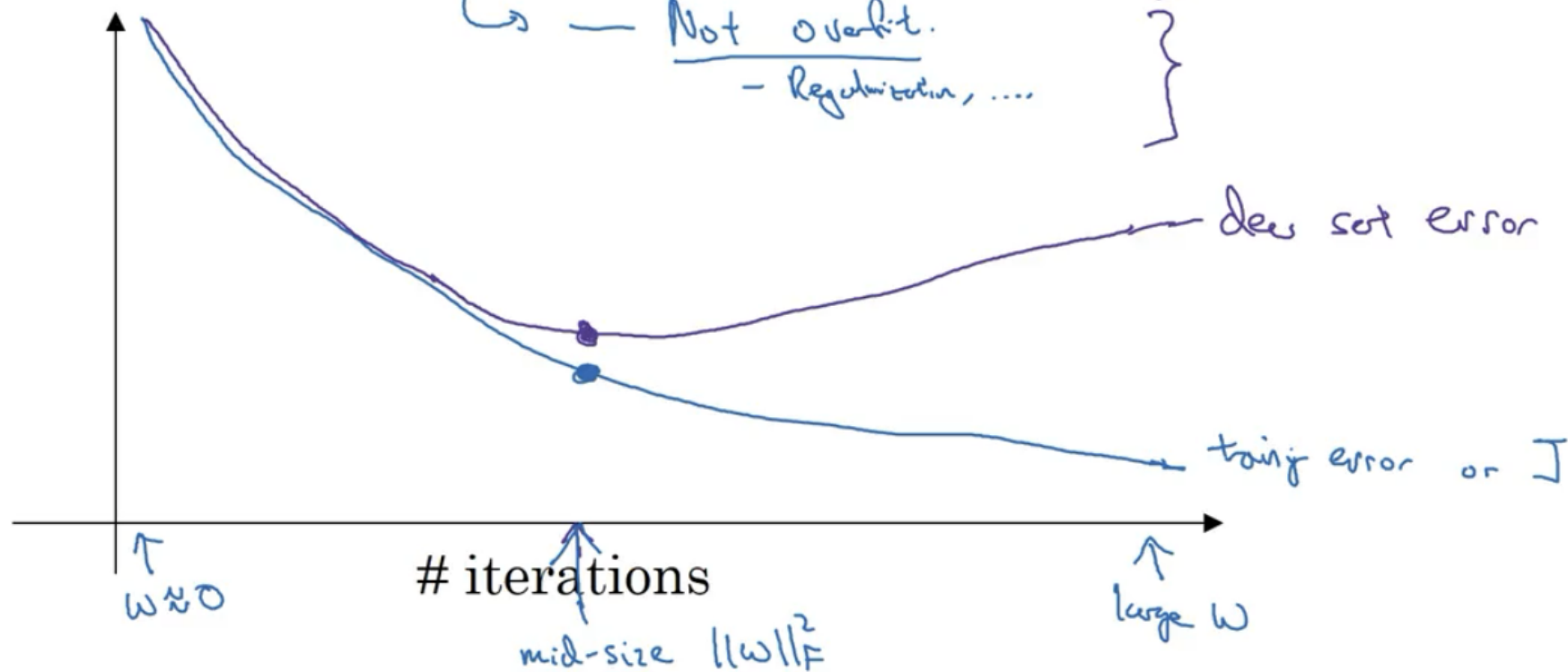# Dropout regularization



Andrew Ng

# Normalizing inputs to speed up learning

$x_1$
$w, b$
$x_2 \longrightarrow \hat{y}$
$x_3$

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

$X = X - \mu$

$\sigma^2 = \frac{1}{m} \sum_i x^{(i)2}$ $\longleftarrow$ elent-wise

$X = X / \sigma^2$

# Normalizing inputs to speed up learning
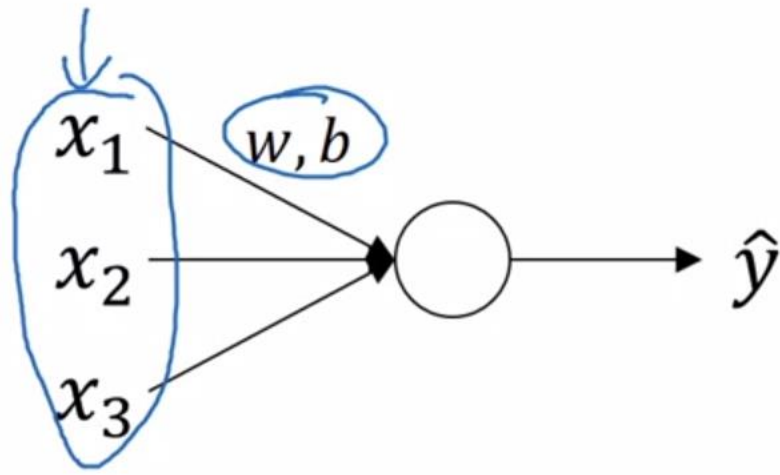
$x_1$ $x_2$ $x_3$ $\quad w,b \quad \rightarrow \hat{y}$

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

$$X = X - \mu$$

$$\sigma^2 = \frac{1}{m} \sum_i x^{(i)2} \quad \leftarrow \text{element-wise}$$
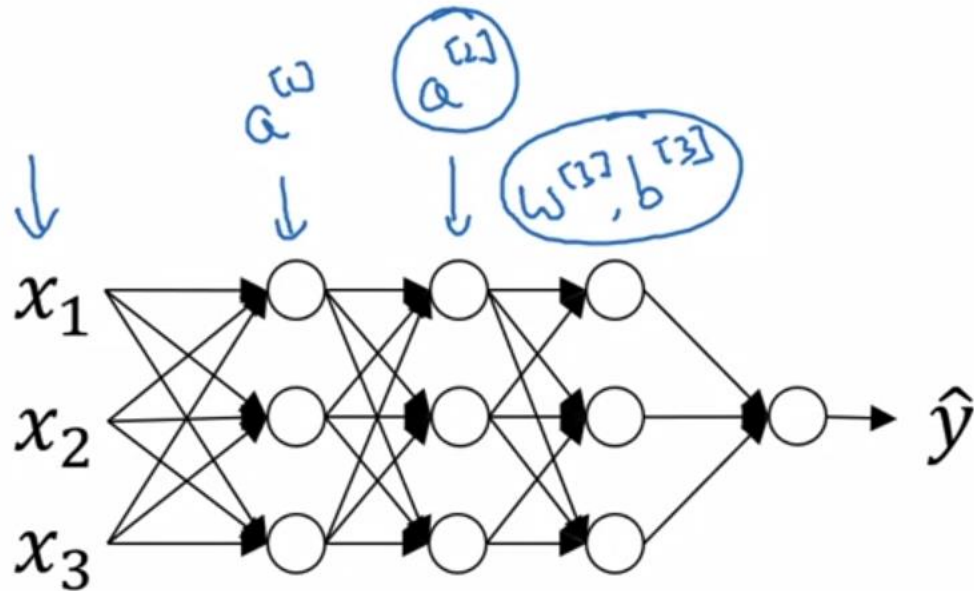
$$X = X / \sigma^2$$

$x_1$ $x_2$ $x_3$ $\quad a^{[1]} \quad a^{[2]} \quad w^{[3]}, b^{[3]} \quad \rightarrow \hat{y}$

Can we normalize $a^{[2]}$ so as to train $w^{[3]}, b^{[3]}$ faster

Normalize $z^{[2]}$

# Implementing Batch Norm

Given some intermediate values in NN $\quad z^{(1)}, \ldots, z^{(m)}$

$$z^{[l](i)}$$

$$\mu = \frac{1}{m} \sum_i z^{(i)}$$

$$\sigma^2 = \frac{1}{m} \sum_i (z_i - \mu)^2$$

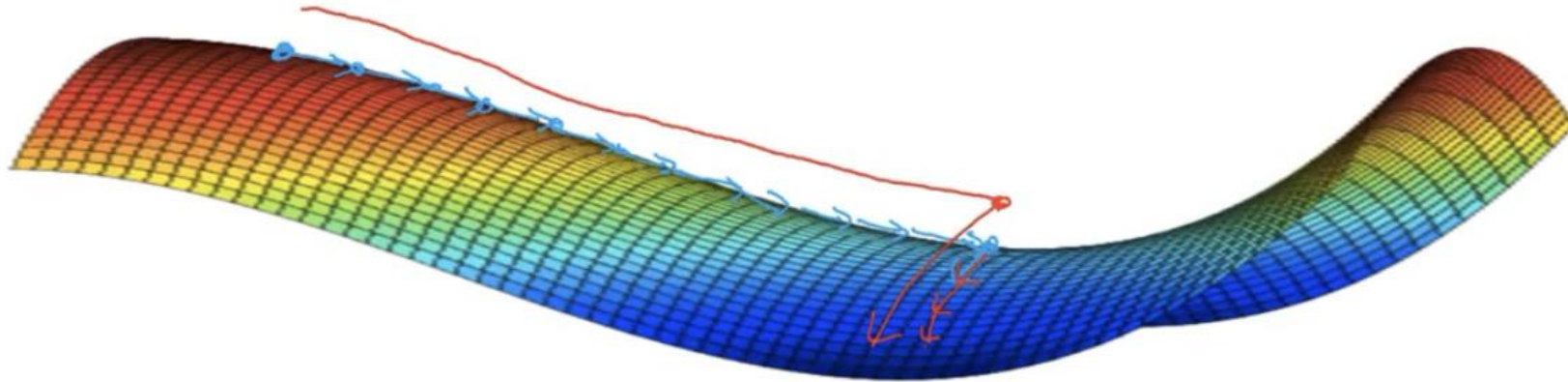$$z_{norm}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \varepsilon}}$$

$$\gamma = \sqrt{\sigma^2 + \varepsilon}$$

$$\beta = \mu$$

$$\tilde{z}^{(i)} = \gamma \, z_{norm}^{(i)} + \beta$$

learnable parameters of model.

# Problem of plateaus



- Unlikely to get stuck in a bad local optima
- Plateaus can make learning slow

Andrew Ng

```html
<html>
<head>
<script src="https://cdn.jsdelivr.net/npm/@tensorflow/tfjs@latest"></script>
<script>
    async function run(){
        const MODEL_URL = 'http://127.0.0.1:8887/model.json';
        const model = await tf.loadLayersModel(MODEL_URL);
        console.log(model.summary());
        const input = tf.tensor2d([10.0], [1, 1]);
        const result = model.predict(input);
        alert(result);
    }
    run();
</script>
<body>
</body>
</html>
```
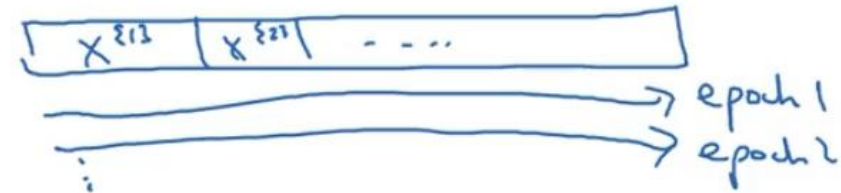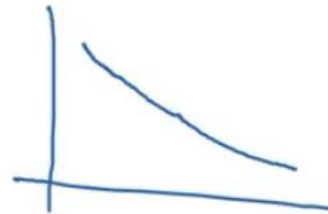
# Learning rate decay

1 epoch = 1 pass through data.

$$\alpha = \frac{1}{1 + decay\text{-}rate * epoch\text{-}num} \underline{\alpha_0} \leftarrow$$

| Epoch | $\alpha$ |
|---|---|
| 1 | 0.1 |
| 2 | 0.67 |
| 3 | 0.5 |
| 4 | 0.4 |
| $\vdots$ | $\vdots$ |

$X^{\{1\}}$ $X^{\{2\}}$ $\cdots$

$\rightarrow$ epoch 1
$\rightarrow$ epoch 2

$\alpha_0 = 0.2$
decay-rate = 1

# Learning rate decay

Slowly reduce $\alpha$