Q(3)

$$X = \begin{bmatrix} \underline{\phantom{---}X^{(1)}\phantom{---}} \\ \underline{\phantom{--}X^{(2)}\phantom{--}} \\ \underline{\phantom{--}X^{(3)}\phantom{--}} \\ \vdots \\ \underline{\phantom{--}X^{(m)}\phantom{--}} \end{bmatrix} \quad , \quad Y = \begin{bmatrix} Y^{(1)} \\ \vdots \\ Y^{(m)} \end{bmatrix}$$

$\underbrace{\phantom{----}}_{m \times (n+1)}$ $\longrightarrow (n+1)$ indicates we have an extra feature containing only $1$'s.

Parameters $\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$

$\theta$. ~~this~~ ~~removes the~~

This ensures that $\hat{Y} = X\theta$ works (instead of $Y = W^T X + b$)

## In BGD

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j}, \text{ where } J(\theta) = \frac{1}{2} \sum_{i=1}^{m} (X_\theta^{(i)} \theta - Y^{(i)})$$

Repeat

$$\Rightarrow \theta_j \leftarrow \theta_j - \alpha \sum_{i=1}^{m} (X^{(i)}\theta - Y^{(i)})X_j^{(i)}, \quad j=1,2,\dots,n+1$$

until $J(\theta)$ stops ~~improving~~ reducing.

## SGD

Repeat
$$\theta_j \leftarrow \theta_j - \alpha(X^{(i)}\theta - Y^{(i)}), \quad j=1,2,\dots,n+1$$
until $\mathcal{L}(\theta) = (X^{(i)}\theta - Y^{(i)})^2$ stops reducing.

Repeat this
for each training
example
i.e $i=1,2,\dots,m$.

For complex (large) datasets, SGD is better as it provides ~~better~~ faster convergence.

~~Math~~     ~~A mathematical proof is~~
~~Intuit~~
Intuition

In BGD, we have to calculate ~~at~~ $\sum_{i=1}^{m} X^{(i)}\theta - Y^{(i)}$ by iterating over all training examples before updating $\theta_j$.

Whereas in ~~❸~~ SGD, an ~~update~~ optimization is made for each training example. ~~This makes sure It is highly bi~~ It is highly likely that ~~the~~ optimizing $\theta$ for each training example is ~~the~~ "close" to the dir^n when we optimize entire cost fn.