

$$h_{\theta}(u) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n.$$

So the number of columns in θ would be number of features + 1

$$x = \begin{bmatrix} \dots x_1 \dots \\ \dots x_2 \dots \\ \vdots \\ \dots x_m \dots \end{bmatrix} \quad \begin{array}{l} m \text{ training examples,} \\ n \text{ features,} \end{array}$$

But for calculating $h_{\theta}(u)$ we need $n+1$ columns to add the bias.

$$h_{\theta}(u) = \underbrace{\theta_0}_1 + \boxed{\theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n} \quad \begin{array}{l} + \\ n \text{ columns in } x \text{ vector} \end{array}$$

Formula for gradient descend.

```

Loop {
  for i = 1 to n, {

     $\theta_j := \theta_j + \alpha (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}, \quad (\text{for every } j)$ 

  }
}

```

$$\theta_j = \theta_j - \underbrace{\alpha \cdot \frac{2}{n}}_{\text{Scaler.}} \sum (\underbrace{\vec{x}_b}_{\text{vector}} \underbrace{\theta}_{\text{vector}} - \underbrace{y}_{\text{vector}}) \underbrace{x_j}_{\text{vector}}$$

Vector Dimension

$x_b \rightarrow (m, n+1)$ ($n+1$ to accommodate the bias)

$\theta \rightarrow (n+1, 1)$

$y \rightarrow (m, 1)$

So to apply the vectorization we do the following,

$$x_b \cdot \theta \rightarrow (m, n+1) \times (n+1, 1) \rightarrow (m, 1)$$

$$x_b \cdot \theta - y \rightarrow (m, 1) - (m, 1) \rightarrow (m, 1)$$

So to be able to multiply x_b with the result we need to transpose it.

$$\left. \begin{array}{l} x_b^T \rightarrow (n+1, m) \\ x_b \cdot \theta - y \rightarrow (m, 1) \end{array} \right\} \text{compatible} \rightarrow (n+1, 1)$$

$$\theta_j = \underbrace{\theta_j}_{(n+1,1)} - \underbrace{\alpha \cdot \frac{2}{n} \sum (\vec{x}_b \theta - y) x_j}_{n+1,1}$$

Hence the dimensions are properly matched.

And it would use the parallelism to calculate the gradient.

Further optimizations:

We can stop the G.D after certain steps if the cost reduction is very low,