# ⊙ MULTIPLE FEATURES

→ So far in the previous lectures of linear regression that we discussed, we had linear regression in only 1 variable — the size of the house — with which we wanted to predict the price of the house

→ Now imagine if we had not only the size of the house but more information like this —

→ This would give us a lot more intuition

| Size (feet²) $x_1$ | Number of bedrooms $x_2$ | Number of floors $x_3$ | Age of home (years) $x_4$ | Price ($1000) $y$ |
|---|---|---|---|---|
| 2104 | 5 | 1 | 45 | 460 |
| 1416 | 3 | 2 | 40 | 232 |
| 1534 | 3 | 2 | 30 | 315 |
| 852 | 2 | 1 | 36 | 178 |
| ... | ... | ... | ... | ... |

$$\longrightarrow n = 4$$

$$\longrightarrow x^{(2)} = \begin{pmatrix} 1416 \\ 3 \\ 2 \\ 40 \\ 232 \end{pmatrix}$$

4 Dimensional Vector

$$\longrightarrow x_3^{(2)} = 2$$

→ From now on, we will denote each such variable / feature as $(x_n)$
  ↳ So we will have $x_1$, $x_2$, $x_3$, $x_4$ and so on for features

→ The value being predicted still remains the same — $y$

→ To formalize this approach, now that we have 4 features :

$$n \qquad - \qquad \text{Number of features}$$
$$x^{(i)} \qquad - \qquad \text{Input (features) of } i^{th} \text{ training example}$$
$$x_j^{(i)} \qquad - \qquad \text{Value of feature } j \text{ in } i^{th} \text{ training example}$$

• How does this affect the Hypothesis Function ?

→ Previously, hypothesis was : $h_\theta(x) = \theta_0 + \theta_1 x$ ; but now that we multiple features we are not going to use this simple representation anymore.

→ Instead we are going to use something like this :

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4 \ldots + \theta_n x_4$$

→ So, for example the house pricing problem :

$$h_\theta(x) = 80 + 0.1 x_1 + 0.01 x_2 + 3 x_3 - 2 x_4$$

size    #bedrooms    #floors    age

→ For convenience of notation, let $x_0 = 1$. We will add this feature with value 1 to simplify the equation :

$$\boxed{x_0^{(i)} = 1}$$

$$\Rightarrow \quad x = \begin{pmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{pmatrix} \in \mathbb{R}^{n+1} \qquad \theta = \begin{pmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \\ \vdots \\ \theta_n \end{pmatrix} \in \mathbb{R}^{n+1}$$

These are both (n+1) dimensional vectors now

And,

$$h_\theta(x) = \theta_0 x_0^{=1} + \theta_1 x_1 + \theta_2 x_2 + \ldots + \theta_n x_n$$

$$= \theta^T x$$

## MULTIVARIATE LINEAR REGRESSION

$$\boxed{h_\theta(x) = \theta^T x}$$

◉ GRADIENT DESCENT FOR MULTIVARIATE LINEAR REGRESSION

→ Hypothesis — $h_\theta(x) = \theta^T x = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_3 + \dots + \theta_n x_n$

→ Parameters — $\theta_0, \theta_1, \theta_2 \dots \theta_n$ now replaced by $(n+1)$ dimensional vector $\boxed{\theta}$

→ Cost Function — $J(\theta_0, \theta_1, \theta_2 \dots \theta_n)$ now replaced by $J(\theta)$

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2$$

→ Gradient Descent — Repeat {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

}

(simultaneously update for every $j = 0, 1 \dots n$)

• So, after putting the values from these new definitions :

Gradient Descent —

$$\boxed{\begin{array}{l} \text{Repeat} \quad \{ \\ \qquad \theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right) x_j^{(i)} \\ \} \\ \left(\text{simultaneously update } \theta_j \text{ for } j = 0, 1, 2 \dots n\right) \end{array}}$$

→ New algorithm derivation

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right)^2$$

$$\Rightarrow \frac{1}{2m} \sum_{i=1}^{m} \left( \theta_0 x_0^{(i)} + \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)} + \dots + \theta_n x_n^{(i)} - y^{(i)} \right)^2$$

$\hookrightarrow \frac{\partial}{\partial \theta_0} J(\theta) = \frac{1}{2m} \cdot 2 \cdot \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right) \cdot x_0^{(i)}$

$\hookrightarrow \frac{\partial}{\partial \theta_1} J(\theta) = \frac{1}{2m} \cdot 2 \cdot \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right) \cdot x_1^{(i)}$

$\hookrightarrow$ We can clearly see a pattern here :

$$\boxed{\frac{\partial}{\partial \theta_j} J(\theta) = \frac{1}{m} \sum_{i=1}^{m} \left( h_\theta(x^{(i)}) - y^{(i)} \right) x_j^{(i)}} \longrightarrow \text{This is what we have used in the new gradient descent algorithm}$$

◉ FEATURE SCALING

→ If we have a dataset with multiple features and we scale them to a similar scale, then gradient descent will converge more quickly
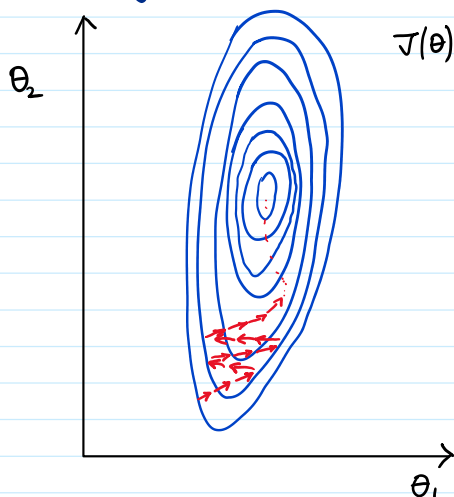
✱ EXAMPLE

Suppose we are dealing with 2 features while predicting the price of a house

$$x_1 = \text{Size of the house} \quad (0 - 2000 \text{ feet}^2)$$

$$x_2 = \text{Number of bedrooms} \quad (1-5)$$

If we are to plot the contour of the cost function for this problem, we will get something like this

$\theta_2$      $J(\theta)$    → When the contour is like this, gradient descent will take a lot of time before it reaches the global minima

→ All the incremental steps will only do very little and hence a LOT of such steps will be required

$\theta_1$

→ In these settings, a useful thing to do is to scale the features. So now, the features will be:

$$x_1 = \frac{\text{size (feet}^2)}{2000}$$

$$x_2 = \frac{\text{Number of bedrooms}}{5}$$

Now, the contours will become much more like circles something like this :

$\theta_2$    $J(\theta)$      Now gradient descent will find a much more direct path to the global minima rather than taking a huge number of iterative steps

$\theta_1$

→ More generally, the aim of feature scaling is to get every feature into approximately a range of $-1 \leq x_i \leq 1$

↳ The numbers -1 and 1 are not too important, it is this range which should be small and comparable with other features

$$0 \leq x_1 \leq 3 \quad \checkmark$$
$$-2 < x_2 < 0.5 \quad \checkmark$$

$$0 \leq x_1 \leq 3 \qquad \checkmark$$
$$-2 \leq x_2 \leq 0.5 \qquad \checkmark$$
$$-100 \leq x_3 \leq 100 \qquad \times \qquad (\text{too large})$$
$$-0.0001 \leq x_4 \leq 0.0001 \qquad \times \qquad (\text{too small})$$

- **MEAN NORMALIZATION**

→ Replace $x_i$ with $\left( x_i - \mu_i \right)$ to make features have approximately ZERO mean
→ Do not apply on $x_0$

☆ **EXAMPLE**

$$x_1 = \frac{size - 1000}{2000} \qquad (\text{Mean size} = 1000 \text{ feet}^2)$$

$$x_2 = \frac{\#bedrooms - 2}{5} \qquad (\text{Mean } \#bedrooms = 2)$$

$$\Rightarrow \qquad -0.5 \leq x_1 \leq 0.5$$
$$-0.5 \leq x_2 \leq 0.5$$

L

→ As a general rule, if $\mu \rightarrow$ mean and $S \rightarrow$ range $(max - min)$

$$x_1 \leftarrow \frac{x_1 - \mu_1}{S_1}$$
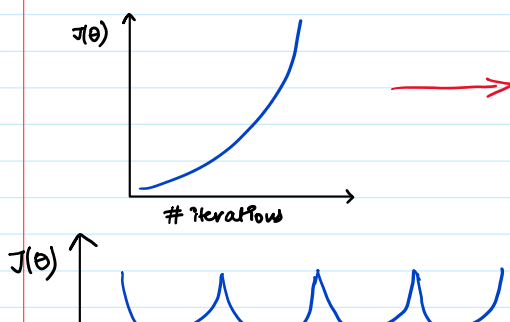
▣ **LEARNING RATE** $- \boxed{\alpha}$

Ⓠ How to make sure gradient descent is working?

→ One way to find out if gradient descent is working properly is to plot the cost function against the number of iterations of gradient descent
→ If gradient descent is running properly, the value of the cost function should decrease with each iteration
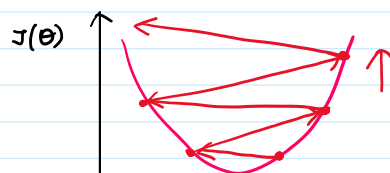


→ Also, if we look after 400 iterations the value of $J(\theta)$ more or less remains the same
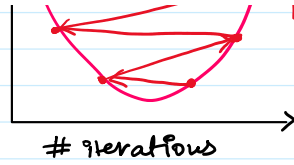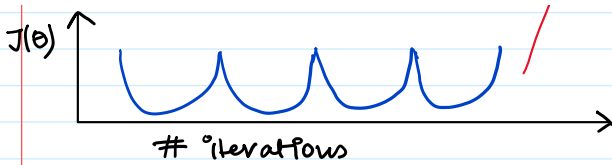→ This suggests that after 400 iterations the cost function converged
→

→ The plot cost vs. # iterations also tells us when the algorithm is not working



→ This suggests that we might be using larger than required value of the learning rate
→ We need to use a smaller learning rate

→ If the learning rate is high, the algorithm overshoots and gradient descent fails to converge

$J(\theta)$



\# iterations



\# iterations

... ~~high~~, the ~~algorithm~~
overshoots and gradient
descent fails to converge
→ It might even diverge

→ For a sufficiently small $\boxed{\alpha}$, $J(\theta)$ should decrease on every iteration
→ But if $\alpha$ is too small, gradient descent can be slow to converge

## ⊡ FEATURE GENERATION

→ Suppose we have 2 features for a house :
  ○ frontage
  ○ depth

→ When applying linear regression we can have :

$$h_\theta(x) = \theta_0 + \theta_1 \times \underbrace{\text{frontage}}_{x_1} + \theta_2 \times \underbrace{\text{depth}}_{x_2}$$



frontage          depth

→ But we don't necessarily have to use the features $x_1$ and $x_2$ that we are given
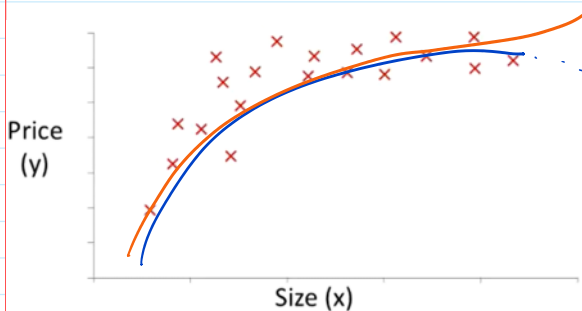→ What we can do b is to create new features from the given ones if we feel they are more appropriate

## ★ EXAMPLE

→ If we have frontage and depth, we can come up with a new feature :

$$\text{Area} = \text{frontage} \times \text{depth}$$

$$\Rightarrow \quad h_\theta(x) = \theta_0 + \theta_1 \times \text{Area}$$

## ⊡ POLYNOMIAL REGRESSION

→ Instead of a linear model, which will not be the most appropriate choice all the time, we can have a polynomial model
→ This will enable us to generate a line which better fits the data



Price (y)

Size (x)

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2^2$$

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2^2 + \theta_3 x_3^3$$