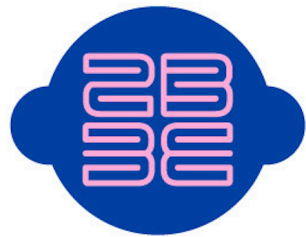


# Linear Regression



**2B3E**

THE 2ND BRAIN, THE 3RD EYE

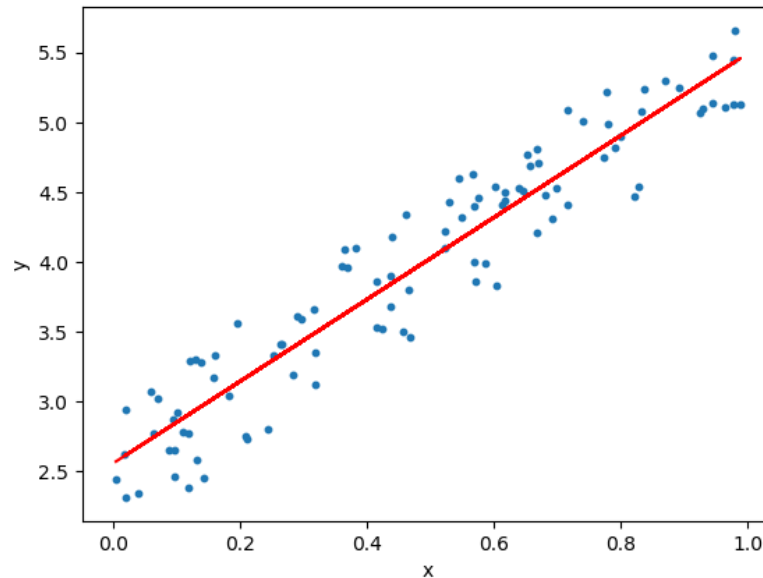
**JP @ 2B3E**

# Contents

- Linear regression
- Multi-variate linear regression
- Gradient descent for linear regression

# Linear regression

- Solving a regression problem, with a linear function



# Linear regression with multiple features

- Multiple variables = multiple features
- House price prediction with one feature
  - $X$  = house size, use this to predict
  - $y$  = house price
- If we have more than one feature (such as number of bedrooms, number floors, age of the home)
  - $x_1, x_2, x_3, x_4$  are the four features
    - $x_1$  - size (feet squared)
    - $x_2$  - Number of bedrooms
    - $x_3$  - Number of floors
    - $x_4$  - Age of home (years)
  - $y$  is the output variable (price)

# Linear regression with multiple features

- More notation **n**
  - number of features ( $n = 4$ )
- **m**
  - number of examples (i.e. number of rows in a table)
- **$x^i$** 
  - vector of the input for an example (so a vector of the four parameters for the  $i^{\text{th}}$  input example)
  - $i$  is an index into the training set
  - So
    - $x$  is an  $n$ -dimensional feature vector
    - $x^3$  is, for example, the 3rd house, and contains the four features associated with that house
- **$x_j^i$** 
  - The value of feature  $j$  in the  $i$ -th training example
  - So
    - $x_2^3$  is, for example, the number of bedrooms in the third house

# Linear regression with multiple features

- hypothesis
  - Hypothesis with one feature
    - $h_{\theta}(x) = \theta_0 + \theta_1 x$ 
      - two parameters (including the bias)
      - One variable  $x$
  - Hypothesis with multiple features
    - $h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4$
  - For example
    - $h_{\theta}(x) = 80 + 0.1x_1 + 0.01x_2 + 3x_3 - 2x_4$ 
      - An example of a hypothesis which is trying to predict the price of a house
      - Parameters are still determined through a cost function

# Linear regression with multiple features

- For convenience of notation,  $x_0 = 1$ 
  - Bias: an additional 0th feature for each example
  - **feature vector**
    - $n + 1$  dimensional feature vector indexed from 0
    - a column **X**: each example has a column vector associated with it
  - **Parameters**
    - are  $n+1$  dimensional vector
    - This is also a column vector called  **$\theta$**
    - This vector is the same for each example
- Considering this, hypothesis can be written
  - $h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4$

# Linear regression with multiple features

- Vector form:  $h_{\theta}(x) = \theta^T X$ 
  - $\theta^T$  is an  $[1 \times n+1]$  matrix
  - $\theta$ : a column vector  $\rightarrow \theta^T$  a row vector
  - $\theta$ :  $[(n+1) \times 1]$
  - $\theta^T$ :  $[1 \times (n+1)]$
  - $\theta^T X$ 
    - $[1 \times n+1] * [n+1 \times 1]$
    - $= h_{\theta}(x)$
    - So, in other words, the transpose of our parameter vector \* an input example  $X$  gives you a predicted hypothesis which is  $[1 \times 1]$  dimensions (i.e. a single value)
- This is an example of multivariate linear regression



# Gradient descent for multiple variables

- Fitting parameters for the hypothesis with gradient descent
  - Parameters are  $\theta_0$  to  $\theta_n$
  - Instead of thinking about this as  $n$  separate values, think about the parameters as a single vector ( $\theta$ )
    - Where  $\theta$  is  $n+1$  dimensional
- Our cost function is

$$J(\theta_0, \theta_1, \dots, \theta_n) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

# Gradient descent for multiple variables

- Instead of thinking of  $J$  as a function of the  $n+1$  numbers,  $J()$  is just a function of the parameter vector  $J(\theta)$

$$\begin{array}{l} \text{Repeat } \{ \\ \quad \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \dots, \theta_n) \\ \} \end{array} \quad (\text{simultaneously update for every } j = 0, \dots, n)$$

- This is  $\theta_j = \theta_j - \text{learning rate } (\alpha) \text{ times the partial derivative of } J(\theta) \text{ with respect to } \theta_j$
- We do this through a **simultaneous update** of every  $\theta_j$  value

# Gradient descent for multiple variables

- Implementing this algorithm
  - When  $n = 1$

$$\begin{aligned} &\text{Repeat } \{ \\ &\theta_0 := \theta_0 - \alpha \underbrace{\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})}_{\frac{\partial}{\partial \theta_0} J(\theta)} \\ &\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)} \\ &\quad \text{(simultaneously update } \theta_0, \theta_1 \text{)} \} \end{aligned}$$

- Above, we have slightly different update rules for  $\theta_0$  and  $\theta_1$ 
  - Actually they're the same, except the end has a previously undefined  $x_0^{(i)}$  as 1, so wasn't shown
- $\rightarrow$  almost identical rule for multivariate gradient descent

# Gradient descent for multiple variables

New algorithm ( $n \geq 1$ ):

Repeat {

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

(simultaneously update  $\theta_j$  for  $j = 0, \dots, n$ ) }

$\downarrow \frac{\partial}{\partial \theta_j} J(\theta)$

- We're doing this for each  $j$  (0 until  $n$ ) as a simultaneous update (like when  $n = 1$ )
- So, we re-set  $\theta_j$  to
  - $\theta_j$  minus the learning rate ( $\alpha$ ) times the partial derivative of the  $J$  with respect to  $\theta_j$

# Gradient descent for multiple variables

- In non-calculus words, this means that we do
  - Learning rate
  - Times  $1/m$  (makes the maths easier)
  - Times the sum of
    - The hypothesis taking in the variable vector, minus the actual value, times the  $j$ -th value in that variable vector for EACH example
- It's important to remember that

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

•

# Gradient Decent in practice

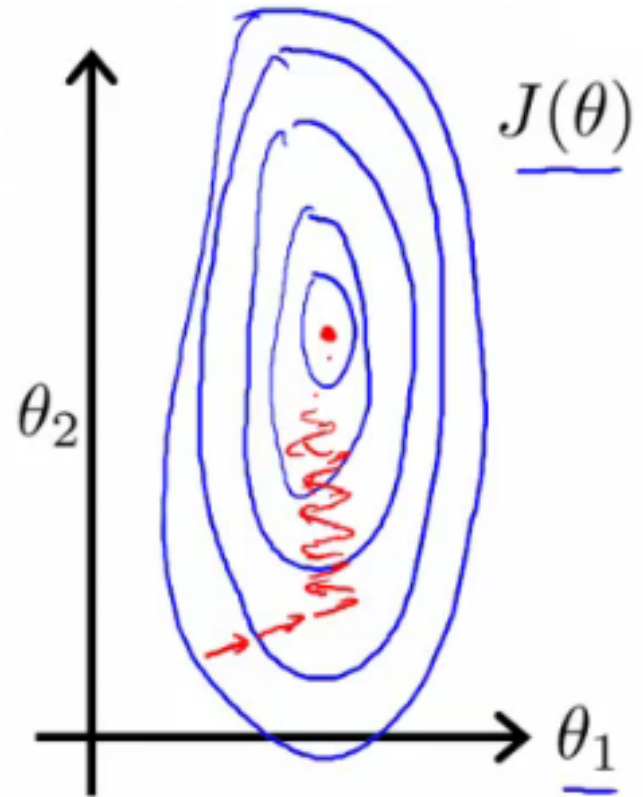
- Feature scaling
- Learning rate

# Feature Scaling

- Having covered the theory
  - —> move on to learn about some of the practical tricks
- Feature scaling
  - If you have a problem with multiple features
  - You should make sure those features have a similar scale
    - Means gradient descent will converge more quickly
  - e.g.
    - $x_1$  = size (0 - 2000 feet)
    - $x_2$  = number of bedrooms (1-5)
    - Means the contours generated if we plot  $\theta_1$  vs.  $\theta_2$  give a very tall and thin shape due to the huge range difference
  - Running gradient descent on this kind of cost function can take a long time to find the global minimum

# Feature Scaling

- Pathological input to gradient descent
- So **we need to rescale** this input so it's more effective
- So, if you define each value from  $x_1$  and  $x_2$  by dividing by the max for each feature
- Contours become more like circles (as **scaled between 0 and 1**)





# Feature Scaling

- May want to get everything into -1 to +1 range (approximately)
  - Want to avoid large ranges, small ranges or very different ranges from one another
  - Rule a thumb regarding acceptable ranges
    - -3 to +3 is generally fine - any bigger bad
    - -1/3 to +1/3 is ok - any smaller bad
- Can do **mean normalization**
  - Take a feature  $x_i$ 
    - Replace it by  $(x_i - \text{mean})/\text{max}$
    - So your values all have an average about 0

# Feature Scaling

$$x_1 \leftarrow \frac{x_1 - \mu_1}{s_1}$$

Annotations:

- $\mu_1$ : avg value of  $x_1$  in training set
- $s_1$ : range (max-min) (or standard deviation)

- Instead of max can also use **standard deviation**

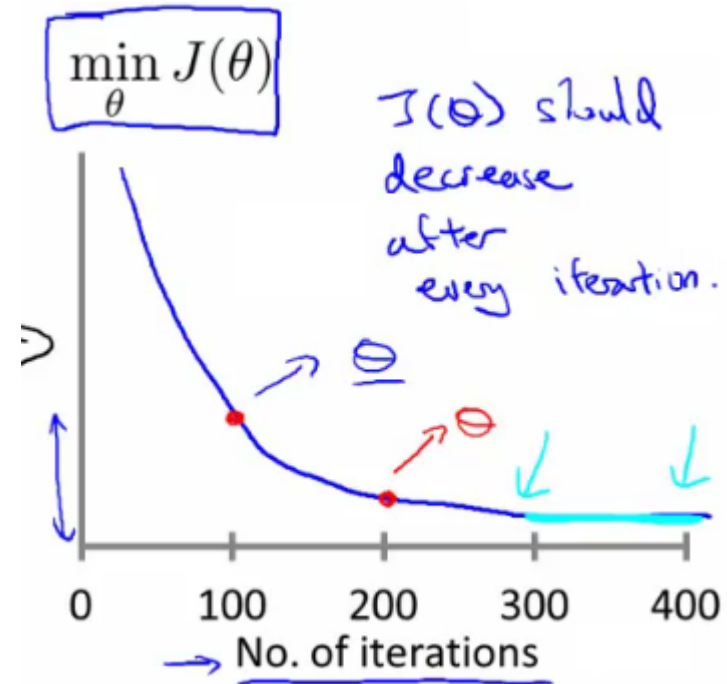
# Learning Rate $\alpha$

- Focus on the learning rate ( $\alpha$ )
- Topics
  - Update rule
  - Debugging
  - How to chose  $\alpha$

# Learning Rate $\alpha$

- **Make sure gradient descent is working**
  - Plot  $\min J(\theta)$  vs. # of iterations
    - (i.e. plotting  $J(\theta)$  over the course of gradient descent)
- If gradient descent is working then  $J(\theta)$  should decrease after (almost) every iteration
- Can also show if you're not making huge gains after a certain number
  - Can apply heuristics to reduce number of iterations
  - If, for example, after 1000 iterations you reduce the parameters by nearly nothing you could choose to only run 1000 iterations in the future
  - Make sure you don't accidentally hard-code thresholds

# Learning Rate $\alpha$



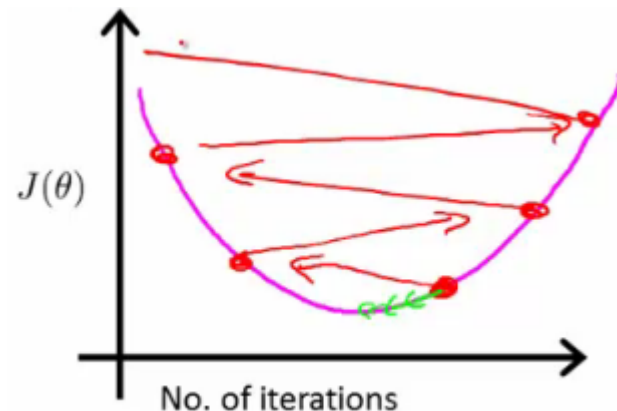
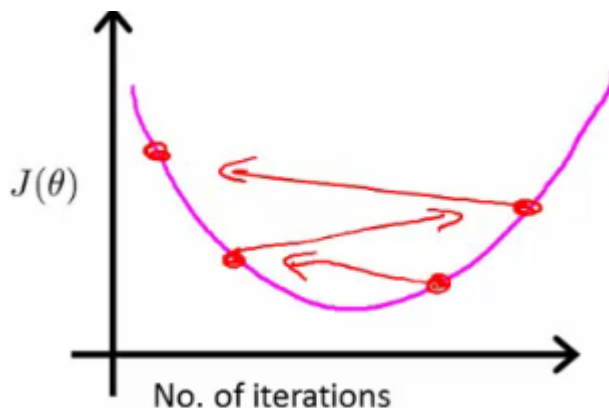
- Number of iterations varies a lot
  - 30 iterations
  - 3000 iterations
  - 3000 000 iterations
- Very hard to tell in advance how many iterations will be needed
- Can often make a guess based a plot like this after the first 100 or so iterations

# Learning Rate $\alpha$

- Automatic convergence tests
  - Check if  $J(\theta)$  changes by a small threshold or less
    - Choosing this threshold is hard
    - So often **easier to check for a straight line**
      - Why? - Because we're seeing the straightness in the context of the whole algorithm

# Learning Rate $\alpha$

- Checking its working
  - If you plot  $J(\theta)$  vs iterations and see the value is increasing - means you probably need a smaller  $\alpha$ 
    - Cause is because your minimizing a function which looks like this



But you overshoot, so reduce learning rate so you actually reach the minimum (green line)  $\rightarrow$  So, use a smaller  $\alpha$

# Learning Rate $\alpha$

- Another problem might be if  $J(\theta)$  looks like a series of waves
  - Here again, you need a smaller  $\alpha$
- However
  - If  $\alpha$  is small enough,  $J(\theta)$  will decrease on every iteration
  - BUT, if  $\alpha$  is too small then rate is too slow
    - A less steep incline is indicative of a slow convergence, because we're decreasing by less on each iteration than a steeper slope
- Typically
  - Try a range of alpha values
  - Plot  $J(\theta)$  vs number of iterations for each version of alpha
  - Go for roughly threefold (or fivefold) increases
    - 0.001, 0.003, 0.01, 0.03, 0.1, 0.3