

7. Gradient Descent in Practice.

I. Feature scaling.

feature and parameter values.

$$\hat{\text{price}} = w_1 x_1 + w_2 x_2 + b.$$

x_1 : size (ft)² x_2 : bedrooms
 ↓ ↓
 size bedrooms

$\text{range } 300 - 2000$ $\text{range } 0 - 5$
 large. small.

House: $x_1 = 2000$, $x_2 = 5$, price = 500k.

One training example.
size of parameters w_1, w_2 ?

If we set $w_1 = 50$, $w_2 = 0.1$, $b = 50$ } not good choices.

$$\hat{\text{price}} = 50 \times 2000 + 0.1 \times 5 + 50 = 100,050.5 \text{ k.}$$

} very far.

Say $w_1 = 0.1$, $w_2 = 50$, $b = 50$

$$\hat{\text{price}} = 0.1 \times 2000 + 50 \times 5 + 50 = 500 \text{ k}$$

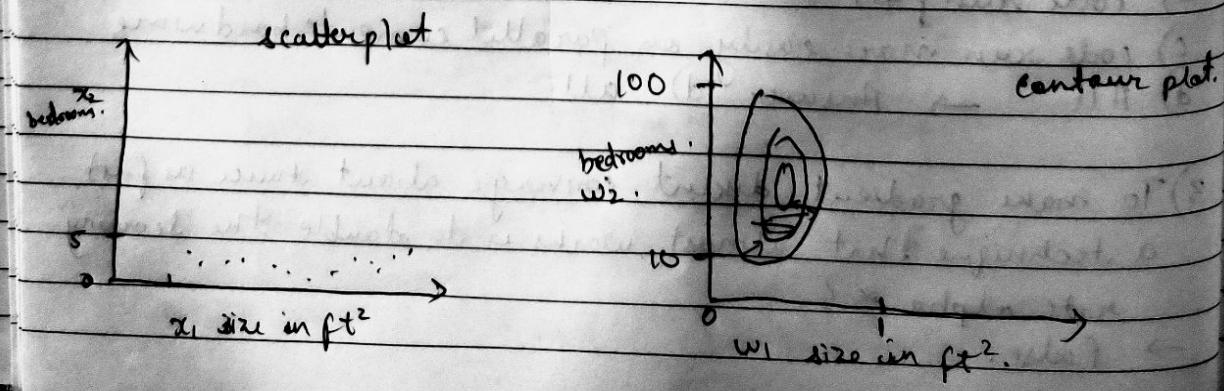
} more reasonable.

When feature value is large, parameter value will be small. If feature value is small, param value might be large.

size of feature x_j size of parameter w_j
 size in ft²
 bedrooms.

features.

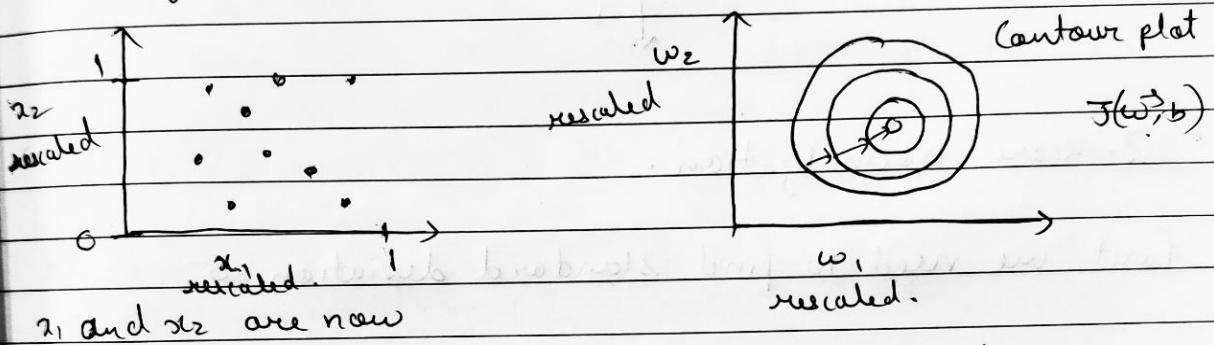
Parameters.



We have shorter on 1 side, longer on another. This is because of value impact.

If we go normal route, our Gradient Descent will go through all circles before reaching global minima.

To rectify this, we should transform our data.



x_1 and x_2 are now take comparable values.

More circular plot.
Gradient descent can find a much more direct path to global minimum.

II. How to scale features?

$$\text{if } 300 \leq x_1 \leq 2000, \quad 0 \leq x_2 \leq 5$$

we can divide x_i by max value, 2000.

$$x_{1, \text{scaled}} = \frac{x_1}{2000} \quad x_{2, \text{scaled}} = \frac{x_2}{5}$$

Now, $0.15 \leq x_{1, \text{scaled}} \leq 1, \quad 0 \leq x_{2, \text{scaled}} \leq 1$.

• Mean Normalization.

$$300 \leq x_1 \leq 2000, \quad 0 \leq x_2 \leq 5$$

they can have both negative and positive values.

let \bar{x}_i be the mean of x_i .

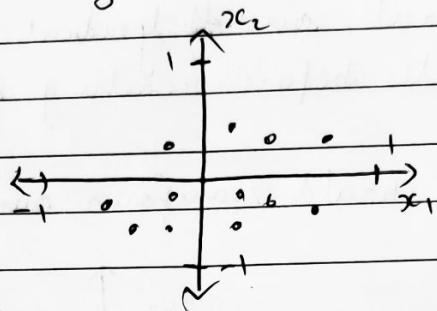
$$\text{let } x_{1, \text{scaled}} = \frac{x_1 - \bar{x}_1}{\frac{2000 - 300}{\max - \min}}$$

$$x_{2, \text{scaled}} = \frac{x_2 - \bar{x}_2}{\frac{5 - 0}{\max - \min}}$$

Date _____
Page 8

$$\text{So, } -0.18 \leq x_1 \leq 0.82, \quad -0.46 \leq x_2 \leq 0.54.$$

Using mean normalized x_1 and x_2 , it will look like this:



• Z-score normalization.

First, we need to find standard deviation σ .

We first find μ_1 , and then find σ .

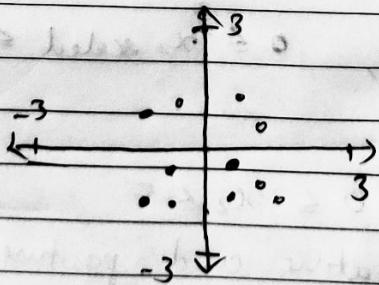
Given, $\mu_1 = 450$, $\mu_2 = 600$ then $\mu_1 = 2.3$, $\mu_2 = 1.4$

then $x_1 = \frac{x_1 - \mu_1}{\sigma_1}$ and $x_2 = \frac{x_2 - \mu_2}{\sigma_2}$

$$-0.67 \leq x_1 \leq 3.1$$

$$-1.6 \leq x_2 \leq 1.9$$

So, the normalized x_1 and x_2 will look like this:



- feature scaling

aim for about $-1 \leq x_j \leq 1$ for each feature x_j :
 $\left. \begin{array}{l} -3 \leq x_j \leq 3 \\ -0.3 \leq x_j \leq 0.3 \end{array} \right\}$ acceptable ranges.

also, $0 \leq x_i \leq 3$ no rescaling

$-2 \leq x_2 \leq 0.5$ no rescaling

$-100 \leq x_3 \leq 100 \rightarrow$ too large, we need to rescale

$-0.001 \leq x_4 \leq 0.001 \rightarrow$ too small, rescale

$98.6 \leq x_5 \leq 105 \rightarrow$ too large, rescale

No harm in carrying out feature scaling.

III. Checking Gradient Descent for convergence

gradient descent:

$$w_j = w_j - \alpha \frac{\partial J(\vec{w}, b)}{\partial w_j}$$

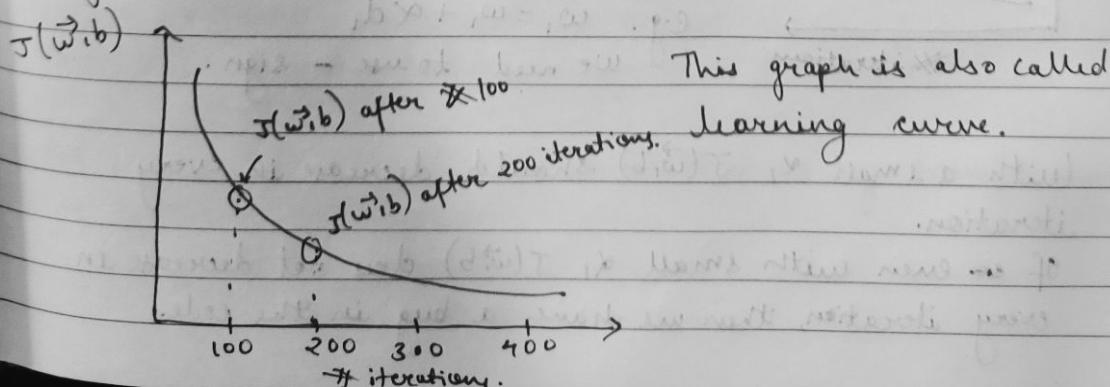
$$b = b - \alpha \frac{\partial J(\vec{w}, b)}{\partial b}$$

a. Make sure gradient descent is working well.

objective: $\min_{\vec{w}, b} J(\vec{w}, b)$, for this, we'll plot the value

of $J(\vec{w}, b)$ w.r.t. no. of iterations we run so far.

Our graph ~~should~~ should look like this:



$J(\vec{w}, b)$ should decrease after every iteration.

at 300, 400 our $J(w, b)$ is almost constant. Our gradient descent is almost converged.

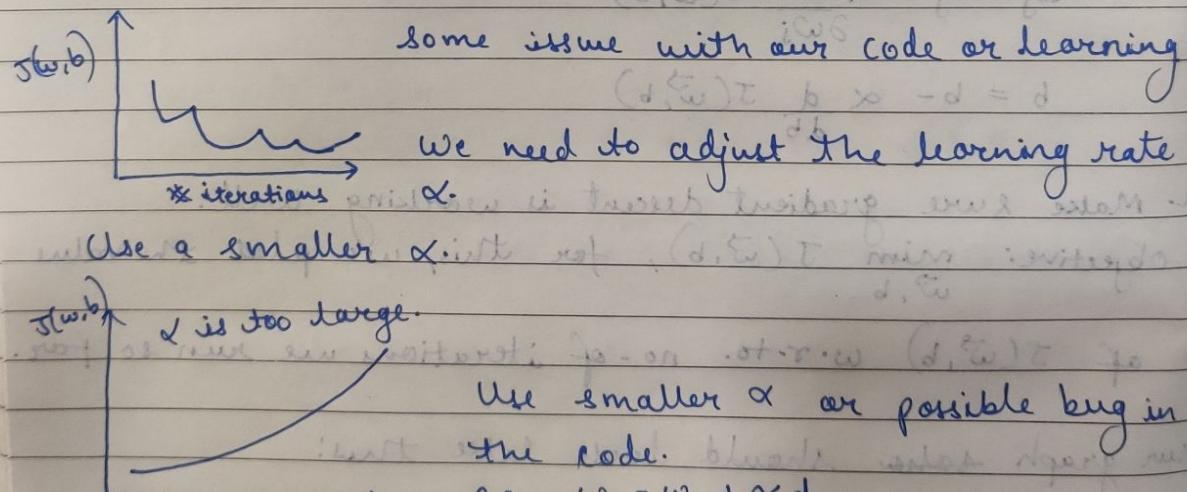
The iterations may vary for each use case.
We can plot this graph to test this out.

b. Automatic convergence test.

Let ϵ "epilen" be 10^{-3}
if $J(\vec{w}, b)$ decreased by $\leq \epsilon$ in one iteration, declare convergence. (found parameters \vec{w}, b to get close to global minimum).

IV. Choosing a learning rate.

a. Identify problems with gradient



(With a small α , $J(\vec{w}, b)$ should decrease in every iteration.)

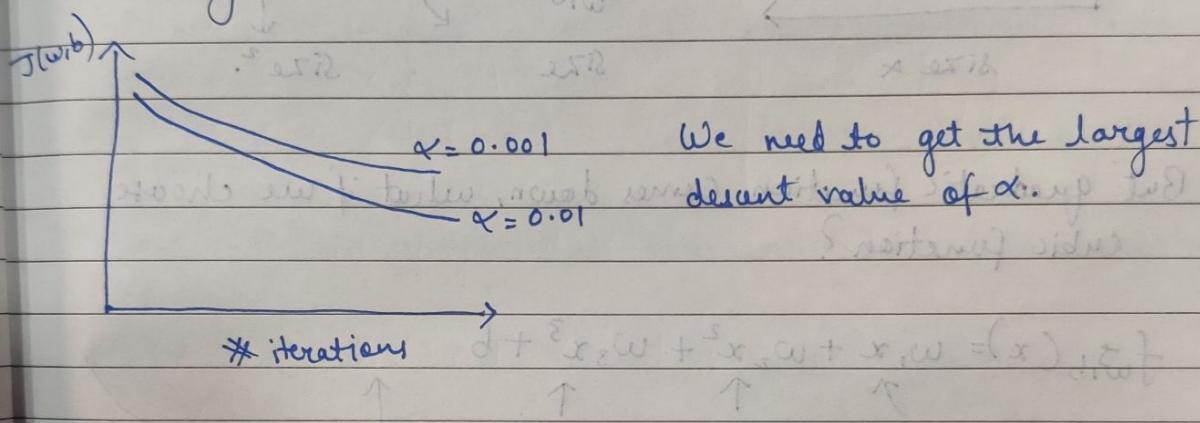
Of even with small α , $J(w, b)$ does not decrease in every iteration, then we have a bug in the code.

If α is too small, it will take a lot of time to converge.

Values of α to try:

0.001 0.01 0.1 1.

We need to plot cost function for these values, and select the value of α where $J(\vec{w}, b)$ decreases rapidly and consistently.



IV. Feature Engineering.

Say we have x_1 as frontage, x_2 as depth.

Let $f_{\vec{w}, b}(\vec{x}) = w_1 x_1 + w_2 x_2 + b$. w_1 and b are removed prior to this.

New, area = frontage \times depth.

$x_3 = x_1 x_2$ (new feature).

$$f_{\vec{w}, b}(\vec{x}) = w_1 x_1 + w_2 x_2 + w_3 x_3 + b. \quad \text{+ } x_1 w_1 + x_2 w_2 + x_3 w_3 = (\vec{w})^T \vec{x} + b$$

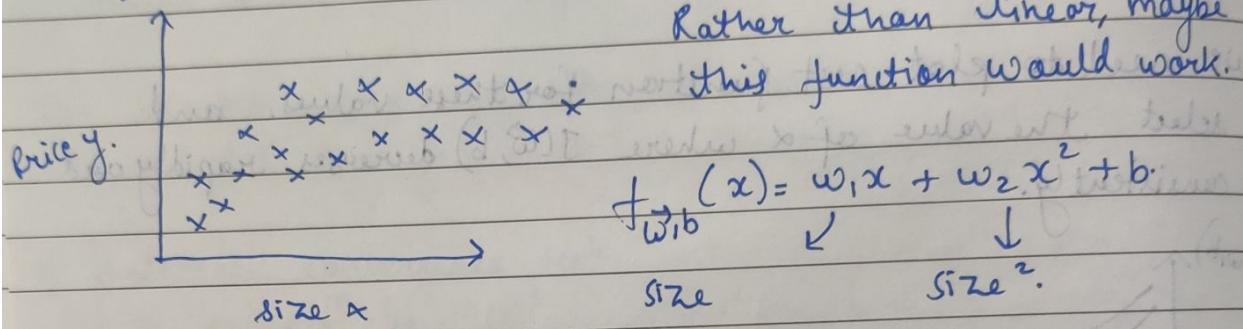
Here, we created a new feature out of existing features.

This is called feature engineering.

Using intuition to design new features, by transforming or combining new features.

VI. Polynomial Regression.

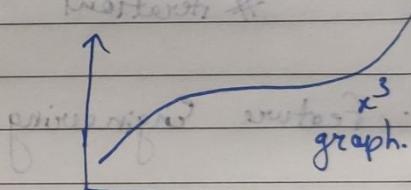
a. Polynomial Regression.



But quadratic function comes down, what if we choose cubic function?

$$f_{w,b}(x) = w_1 x + w_2 x^2 + w_3 x^3 + b$$

↑ ↑ ↑
size size² size³.



When we create x^2 or x^3 values of our input features, feature scaling becomes extremely important.

b. Choices of features.

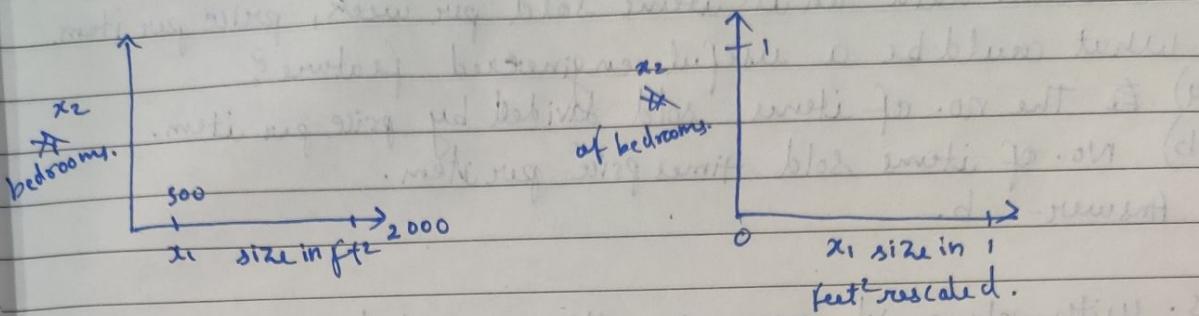
We can also use \sqrt{x} , so our model can look like this.

$$f_{w,b}(x) = w_1 x + w_2 \sqrt{x} + b$$

↑ ↑
size $\sqrt{\text{size}}$

What features to use?

VII. Quiz.

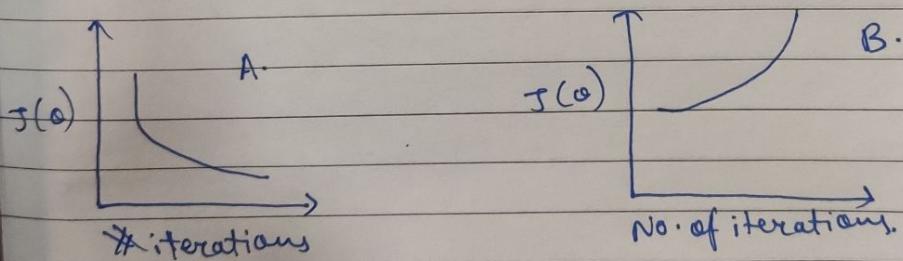


1. Which of the following is a valid step used during feature scaling?

- Add the mean from each value and then divide by the (max - min).
- Subtract the mean from each value and then divide by (max - min).

Answer - b.

2. Suppose a friend ran gradient three separate times with three choices of the learning rate α , and plotted the learning curves for each (cost J for each iteration).



for which case, α was too large.

Answer - for B.

3. Of the circumstances below, for which one is feature scaling particularly helpful?

- When a feature is much larger or smaller than others.
- When all features in the original data range from 0 to 1.

Answer - a.

4. You are helping a grocery store predict its revenue, and have data on its items sold per week, price per item, what could be a useful engineered feature?

- a) The no. of items sold divided by price per item.
- b) No. of items sold times price per item.

Answer - b.

5. With polynomial regression, the predicted values $f_{w,b}(x)$ does not necessarily have to be a straight line (or linear) function of the input feature?

→ True. bivariate least squares does not mean we fit a straight line.

→ It's about minimizing the sum of squared errors ($(y_i - \hat{y}_i)^2$)