
Linear Regression with Multiple Variables

Multiple Features

Multiple features (variables)

Size in feet ² (x)	Price (\$) in 1000's (y)
2104	400
1416	232
1534	315
852	178
...	...

$$f_{w,b}(x) = wx + b$$

Multiple features (variables)

	Size in feet ² x_1	Number of bedrooms x_2	Number of floors x_3	Age of home in years x_4	Price (\$) in \$1000's
$i=2$	2104	5	1	45	460
	1416	3	2	40	232
	1534	3	2	30	315
	852	2	1	36	178

$j = 1 \dots 4$
 $n = 4$

$x_j = j^{th}$ feature

n = number of features

$\vec{x}^{(i)}$ = features of i^{th} training example

$x_j^{(i)}$ = value of feature j in i^{th} training example

$$f_{\vec{w},b}(\vec{x}) = w_1x_1 + w_2x_2 + \dots + w_nx_n + b$$

$$\vec{w} = [w_1 \ w_2 \ w_3 \ \dots \ w_n] \quad \text{parameters of the model}$$

$$\text{vector } \vec{x} = [x_1 \ x_2 \ x_3 \ \dots \ x_n]$$

$$f_{\vec{w},b}(\vec{x}) = \vec{w} \cdot \vec{x} + b = w_1x_1 + w_2x_2 + w_3x_3 + \dots + w_nx_n + b$$

↑
dot product

multiple linear regression

(not multivariate regression)

Linear Regression with Multiple Variables

Vectorization

Parameters and features

$$\vec{w} = [w_1 \quad w_2 \quad w_3] \quad n=3$$

b is a number

$$\vec{x} = [x_1 \quad x_2 \quad x_3]$$

linear algebra: count from 1

$w[0] \quad w[1] \quad w[2]$

NumPy 

```
w = np.array([1.0, 2.5, -3.3])
```

```
b = 4
```

$x[0] \quad x[1] \quad x[2]$

```
x = np.array([10, 20, 30])
```

code: count from 0

Without vectorization $n=100,000$

$$f_{\vec{w},b}(\vec{x}) = w_1x_1 + w_2x_2 + w_3x_3 + b$$

```
f = w[0] * x[0] +  
     w[1] * x[1] +  
     w[2] * x[2] + b
```

Without vectorization

$$f_{\vec{w},b}(\vec{x}) = \left(\sum_{j=1}^n w_j x_j \right) + b \quad \sum_{j=1}^n \rightarrow j=1 \dots n$$

1, 2, 3

$\text{range}(0, n) \rightarrow j=0 \dots n-1$

```
f = 0  
for j in range(0, n):  
    f = f + w[j] * x[j]  
f = f + b
```

Vectorization

$$f_{\vec{w},b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$$

```
f = np.dot(w, x) + b
```

Without vectorization

```
for j in range(0,16):  
    f = f + w[j] * x[j]
```

t_0

$$f + w[0] * x[0]$$

t_1

$$f + w[1] * x[1]$$

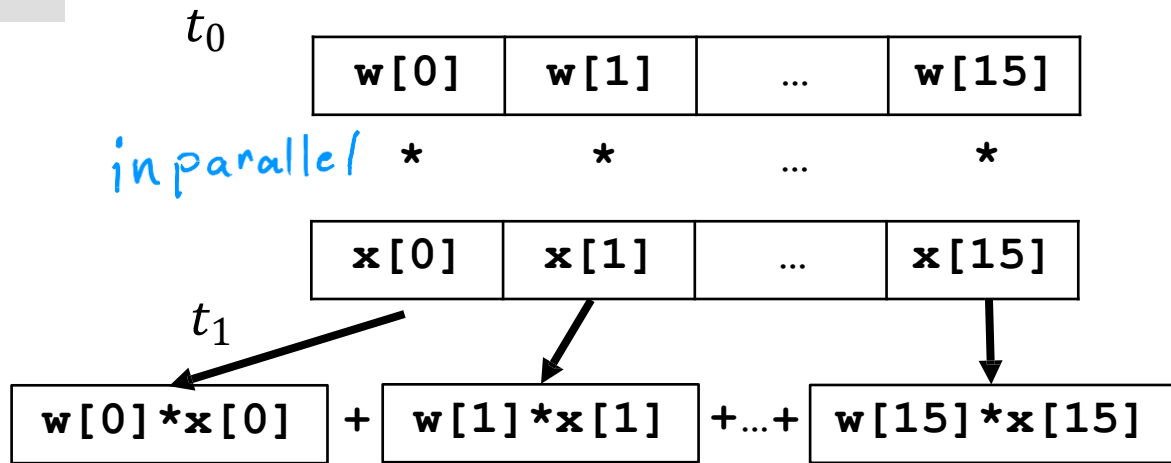
...

t_{15}

$$f + w[15] * x[15]$$

Vectorization

```
np.dot(w,x)
```



efficient → scale to large datasets

Gradient descent $\vec{w} = (w_1 \ w_2 \ \dots \ w_{16})$ ~~b~~ parameters

derivatives $\vec{d} = (d_1 \ d_2 \ \dots \ d_{16})$

```
w = np.array([0.5, 1.3, ... 3.4])
```

```
d = np.array([0.3, 0.2, ... 0.4])
```

compute $w_j = w_j - \underbrace{0.1}_{\text{learning rate } \alpha} d_j$ for $j = 1 \dots 16$

Without vectorization

$$w_1 = w_1 - 0.1d_1$$

$$w_2 = w_2 - 0.1d_2$$

\vdots

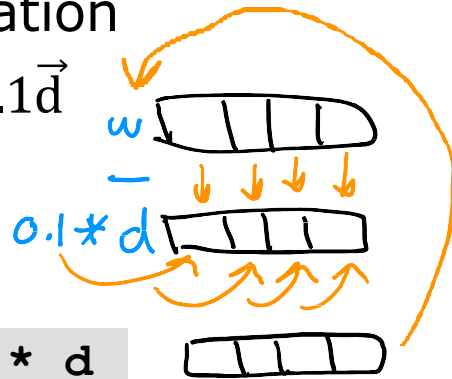
$$w_{16} = w_{16} - 0.1d_{16}$$

```
for j in range(0,16):
```

```
    w[j] = w[j] - 0.1 * d[j]
```

With vectorization

$$\vec{w} = \vec{w} - 0.1\vec{d}$$



```
w = w - 0.1 * d
```


Linear Regression with Multiple Variables

Gradient Descent for Multiple Regression

Previous notation

Parameters

$$w_1, \dots, w_n$$

$$b$$

Model

$$f_{\vec{w},b}(\vec{x}) = w_1 x_1 + \dots + w_n x_n + b$$

Cost function

$$J(w_1, \dots, w_n, b)$$

Gradient descent

repeat {

$$w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(w_1, \dots, w_n, b)$$

$$b = b - \alpha \frac{\partial}{\partial b} J(w_1, \dots, w_n, b)$$

}

Vector notation

← vector of length n

$$\vec{w} = [w_1 \quad \dots \quad w_n]$$

b still a number

$$f_{\vec{w},b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$$

↑ dot product

$$J(\vec{w}, b)$$

repeat {

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

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

}

Gradient descent

One feature

repeat {

$$\underline{w} = w - \alpha \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) x^{(i)}$$

$$\frac{\partial}{\partial w} J(w, b)$$

$$b = b - \alpha \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})$$

simultaneously update w, b

}

n features ($n \geq 2$)

repeat {

$$\overset{j=1}{\underline{w_1}} = w_1 - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\underline{w},b}(\vec{x}^{(i)}) - y^{(i)}) x_1^{(i)}$$

:

$\overset{j=n}{w_n}$

$$w_n = w_n - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\underline{w},b}(\vec{x}^{(i)}) - y^{(i)}) x_n^{(i)}$$

$$b = b - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\underline{w},b}(\vec{x}^{(i)}) - y^{(i)})$$

simultaneously update

w_j (for $j = 1, \dots, n$) and b

}

An alternative to gradient descent

Normal equation

- Only for linear regression
- Solve for w , b without iterations

Disadvantages

- Doesn't generalize to other learning algorithms.
- Slow when the number of features is large ($> 10,000$)

What you need to know

- Normal equation method may be used in machine learning libraries that implement linear regression.
- Gradient descent is the recommended method for finding parameters w, b

Practical Tips for Linear Regression

Feature Scaling



Feature and parameter values

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

\downarrow
size \downarrow
bedrooms

x_1 : size (feet²) x_2 : # bedrooms
range: 300 – 2,000 range: 0 – 5

House: $x_1 = 2000$, $x_2 = 5$, $\text{price} = \$500\text{k}$ one training example

size of the parameters w_1, w_2 ?

$w_1 = 50$, $w_2 = 0.1$, $b = 50$ ←

$$\text{price} = \underbrace{50 * 2000}_{100,000\text{K}} + \underbrace{0.1 * 5}_{0.5\text{K}} + \underbrace{50}_{50\text{K}}$$

$$\text{price} = \$100,050.5\text{k} = \$100,050,500$$

$w_1 = 0.1$, $w_2 = 50$, $b = 50$ →
small

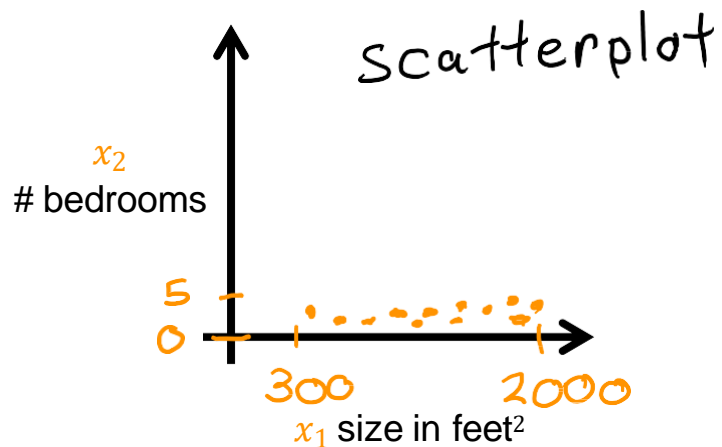
$$\text{price} = \underbrace{0.1 * 2000\text{k}}_{200\text{K}} + \underbrace{50 * 5}_{250\text{K}} + \underbrace{50}_{50\text{K}}$$

$$\text{price} = \$500\text{k} \quad \text{more reasonable}$$

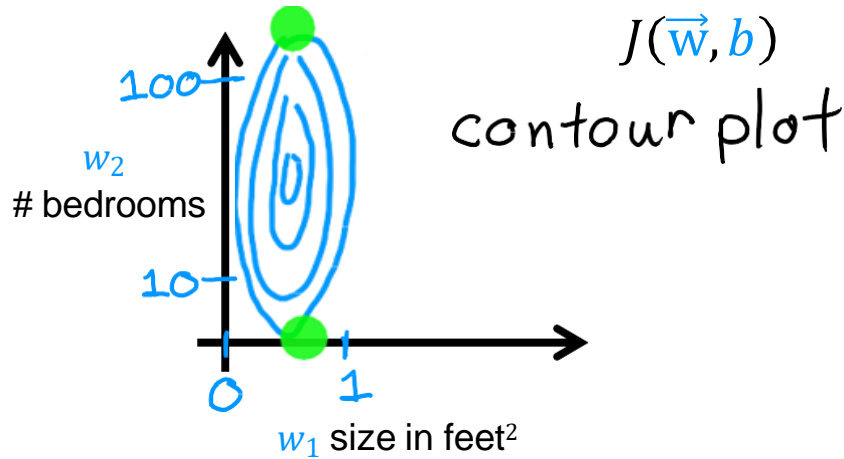
Feature size and parameter size

	size of feature x_j	size of parameter w_j
size in feet ²	←→	←→
#bedrooms	←→	←→

Features



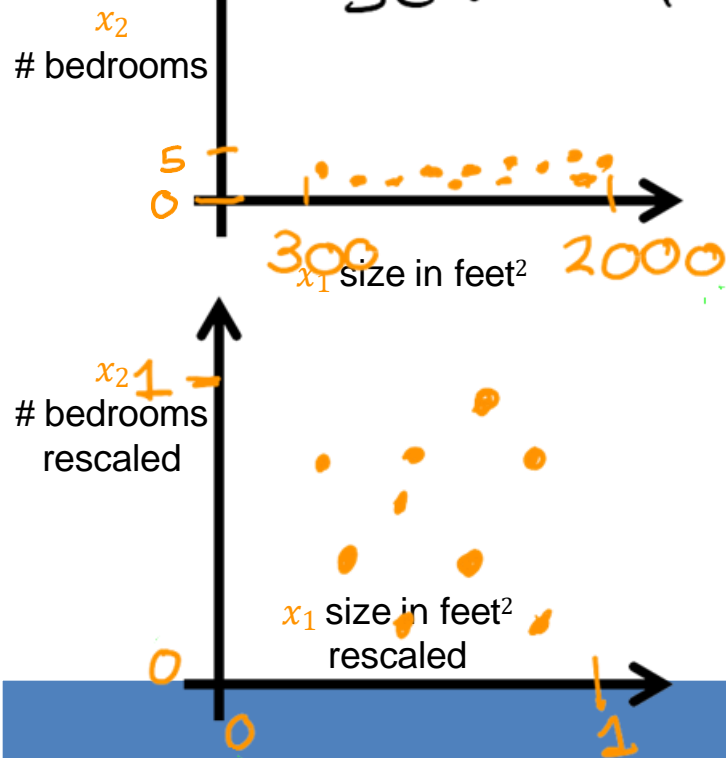
Parameters



Feature size and gradient descent

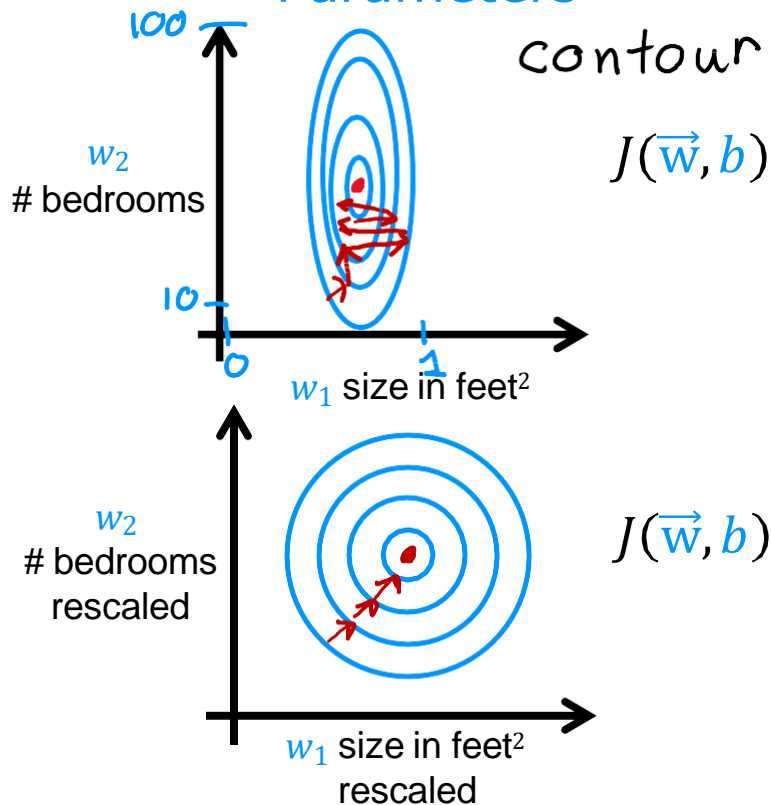
Features

scatterplot



Parameters

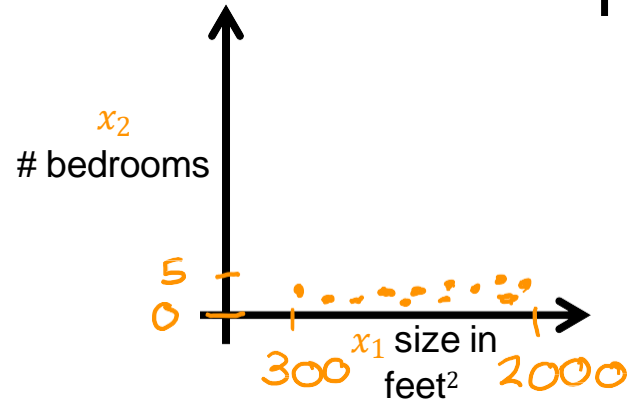
contour plot



$$J(\vec{w}, b)$$

$$J(\vec{w}, b)$$

Feature scaling



$$300 \leq x_1 \leq 2000$$

$$0 \leq x_2 \leq 5$$

$$x_{1,scaled} = \frac{x_1}{2000}$$

max

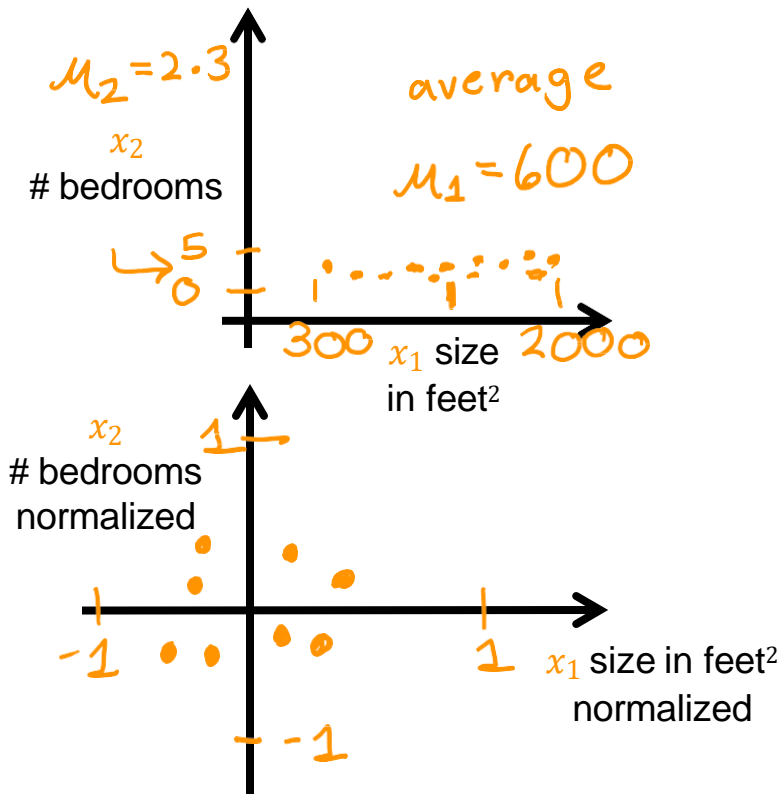
$$x_{2,scaled} = \frac{x_2}{5}$$

max

$$0.15 \leq x_{1,scaled} \leq 1$$

$$0 \leq x_{2,scaled} \leq 1$$

Mean normalization



$$300 \leq x_1 \leq 2000$$

$$x_1 = \frac{x_1 - \mu_1}{2000 - 300}$$

max-min

$$-0.18 \leq x_1 \leq 0.82$$

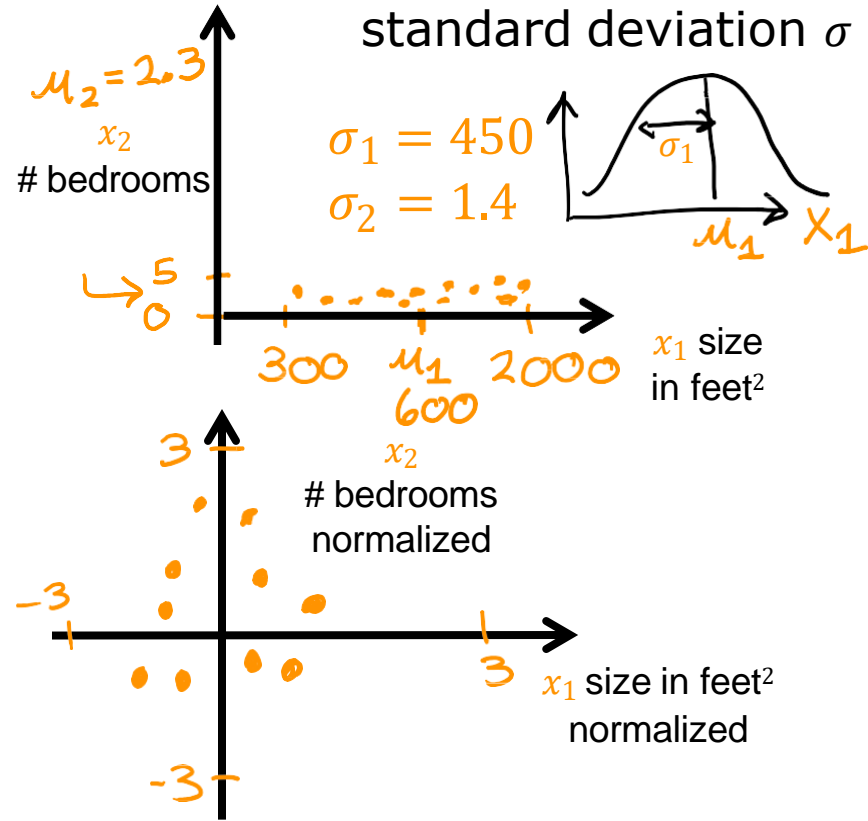
$$0 \leq x_2 \leq 5$$

$$x_2 = \frac{x_2 - \mu_2}{5 - 0}$$

max-min

$$-0.46 \leq x_2 \leq 0.54$$

Z-score normalization



$$300 \leq x_1 \leq 2000$$

$$x_1 = \frac{x_1 - \mu_1}{\sigma_1}$$

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

$$0 \leq x_2 \leq 5$$

$$x_2 = \frac{x_2 - \mu_2}{\sigma_2}$$

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

Feature scaling

aim for about $-1 \leq x_j \leq 1$ for each feature x_j

$$-3 \leq x_j \leq 3$$

$$-0.3 \leq x_j \leq 0.3$$

$$0 \leq x_1 \leq 3$$

$$-2 \leq x_2 \leq 0.5$$

$$-100 \leq x_3 \leq 100$$

$$-0.001 \leq x_4 \leq 0.001$$

$$98.6 \leq x_5 \leq 105$$

Feature scaling

aim for about $-1 \leq x_j \leq 1$ for each feature x_j

$-3 \leq x_j \leq 3$
 $-0.3 \leq x_j \leq 0.3$ } acceptable ranges

$$0 \leq x_1 \leq 3$$

okay, no rescaling

$$-2 \leq x_2 \leq 0.5$$

okay, no rescaling

$$-100 \leq x_3 \leq 100$$

too large \rightarrow rescale

$$-0.001 \leq x_4 \leq 0.001$$

too small \rightarrow rescale

$$98.6 \leq x_5 \leq 105$$

too large \rightarrow rescale

Practical Tips for Linear Regression

Checking Gradient Descent
for Convergence

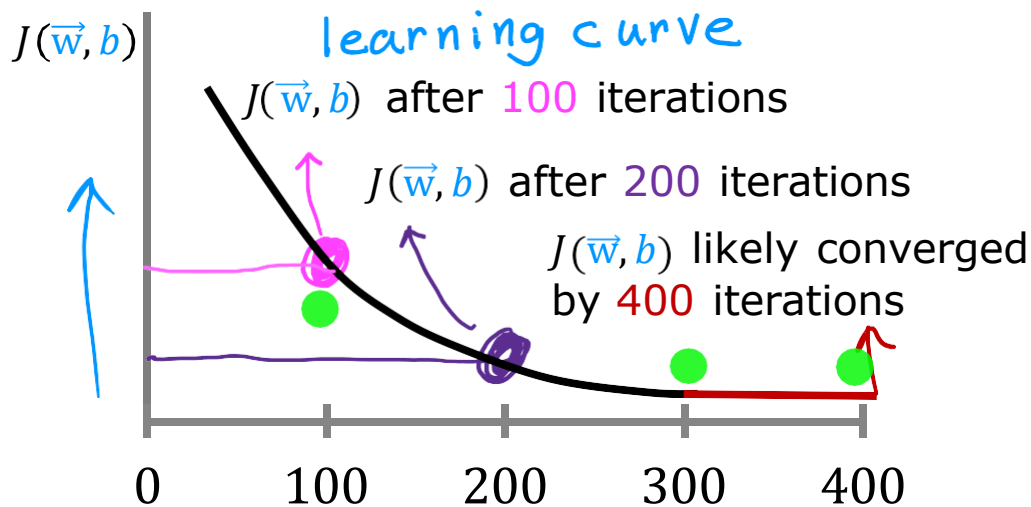
Gradient descent

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

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

Make sure gradient descent is working correctly

objective: $\min_{\vec{w}, b} J(\vec{w}, b)$ $J(\vec{w}, b)$ should **decrease** after every iteration



→ # iterations

~~w, b~~

iterations needed varies 30 1,000 100,000

Automatic convergence test

Let ε "epsilon" be 10^{-3} .
0.001

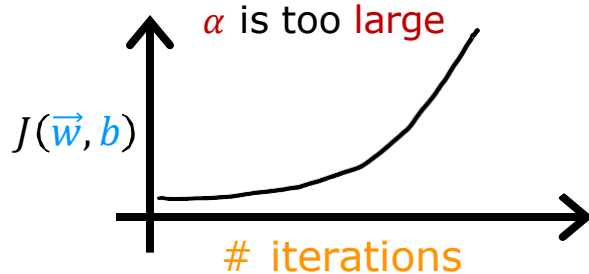
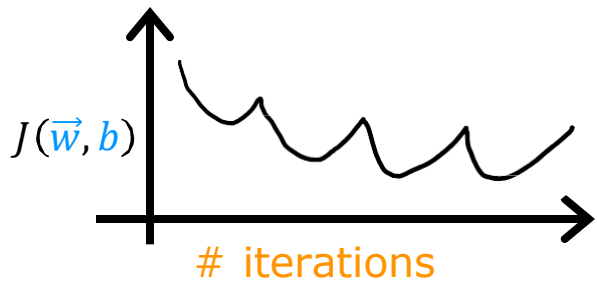
If $J(\vec{w}, b)$ decreases by $\leq \varepsilon$ in one iteration, declare **convergence**.

(found parameters \vec{w}, b to get close to global minimum)

Practical Tips for Linear Regression

Choosing the
Learning Rate

Identify problem with gradient descent



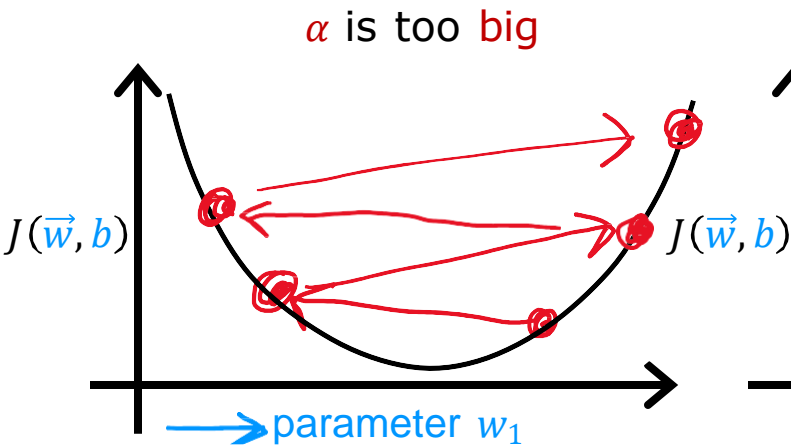
or learning rate is too large

$$w_1 = w_1 + \alpha d_1$$

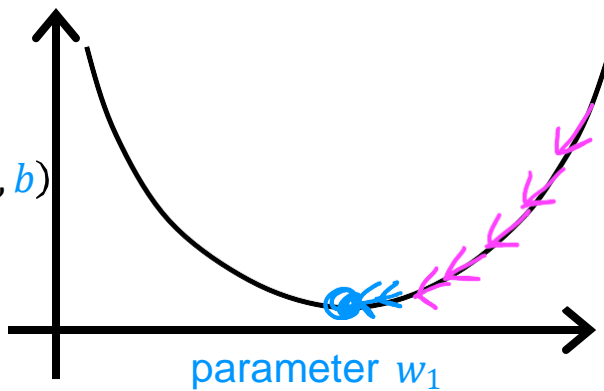
use a minus sign

$$w_1 = w_1 - \alpha d_1$$

Adjust learning rate



Use smaller α

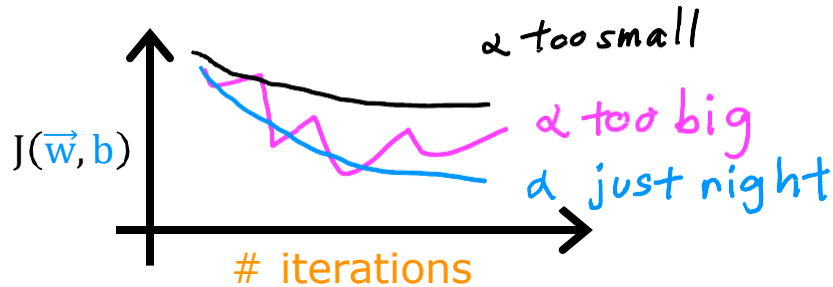
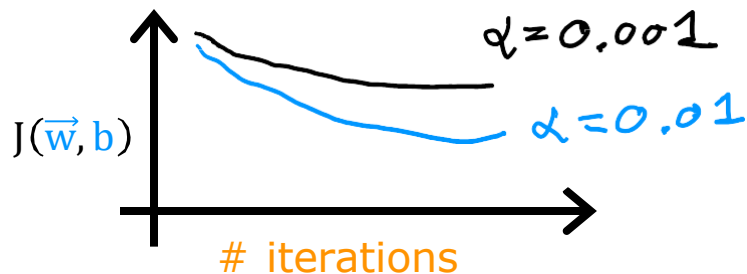


With a small enough α , $J(\vec{w}, b)$ should **decrease** on every iteration

If α is too small, gradient descent takes a lot more iterations to **converge**

Values of α to try:

... 0.001 0.003 0.01 0.03 0.1 0.3 1 ...
 \nearrow \nearrow \nearrow \nearrow \nearrow \nearrow
 $3\times$ $\approx 3\times$ $3\times$ $\approx 3\times$ $3\times$ $\approx 3\times$



Practical Tips for Linear Regression

Feature Engineering



Feature engineering

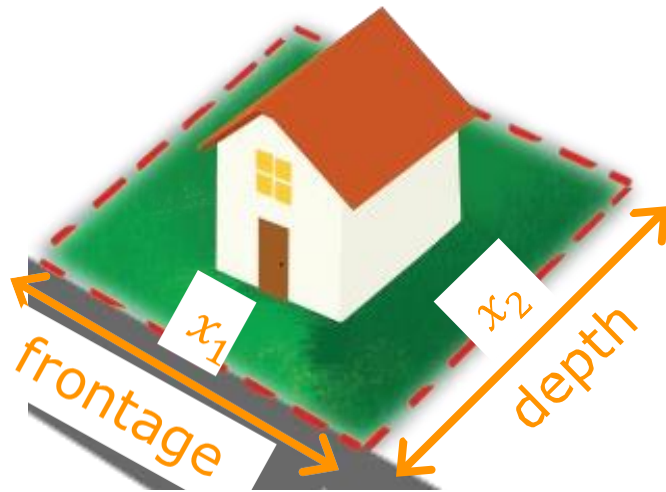
$$f_{\vec{w},b}(\vec{x}) = \underbrace{w_1}_{\text{frontage}} \underbrace{x_1}_{\text{depth}} + \underbrace{w_2}_{\text{depth}} \underbrace{x_2}_{\text{frontage}} + b$$

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

$$x_3 = x_1 x_2$$

new feature

$$f_{\vec{w},b}(\vec{x}) = \underbrace{w_1}_{\text{frontage}} x_1 + \underbrace{w_2}_{\text{depth}} x_2 + \underbrace{w_3}_{\text{area}} x_3 + b$$



Feature engineering:
Using **intuition** to design
new features, by
transforming or combining
original features.

Important notes

- Sci-kit learn library `LinearRegression()` model uses Ordinary Least Squares (OLS) Method with no involvement of optimization
- Stochastic Gradient Descent (SGD) regressor in sci-kit learning uses a stochastic gradient descent algorithm to optimize the cost function.

Ordinary Least Squares (OLS) Method

The ordinary Least Squares (OLS) estimator for linear regression is derived by minimizing the sum of squared errors between the observed values and the predicted values of a linear model.

The goal of OLS is to minimize the sum of the squared residuals (errors) between the actual values \mathbf{y} and the predicted values $\hat{\mathbf{y}}$:

$$\mathbf{y} - \hat{\mathbf{y}} = \mathbf{y} - \mathbf{X}\boldsymbol{\beta}$$

We define the sum of squared errors (SSE) as:

$$S(\boldsymbol{\beta}) = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

This function, $S(\boldsymbol{\beta})$, is the **objective function** we aim to minimize with respect to $\boldsymbol{\beta}$.

Expanding the Objective Function:

Let's expand the SSE to make it easier to differentiate:

$$S(\boldsymbol{\beta}) = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

Using the distributive property of the transpose operation, we get:

$$S(\boldsymbol{\beta}) = \mathbf{y}^\top \mathbf{y} - 2\mathbf{y}^\top \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{X}\boldsymbol{\beta}$$

This expression represents the sum of squared errors that we need to minimize with respect to $\boldsymbol{\beta}$.

Minimization by Taking the Derivative:

To find the value of β that minimizes $S(\beta)$, we take the derivative of $S(\beta)$ with respect to β and set it equal to zero.

The derivative of $S(\beta)$ with respect to β is:

$$\frac{\partial S(\beta)}{\partial \beta} = -2\mathbf{X}^\top \mathbf{y} + 2\mathbf{X}^\top \mathbf{X}\beta$$

Now, set the derivative equal to zero to find the minimum:

$$-2\mathbf{X}^\top \mathbf{y} + 2\mathbf{X}^\top \mathbf{X}\beta = 0$$

Simplifying the equation:

$$\mathbf{X}^\top \mathbf{X}\beta = \mathbf{X}^\top \mathbf{y}$$

OLS Estimator

To solve for β , we multiply both sides of the equation by $(\mathbf{X}^\top \mathbf{X})^{-1}$, assuming $\mathbf{X}^\top \mathbf{X}$ is invertible:

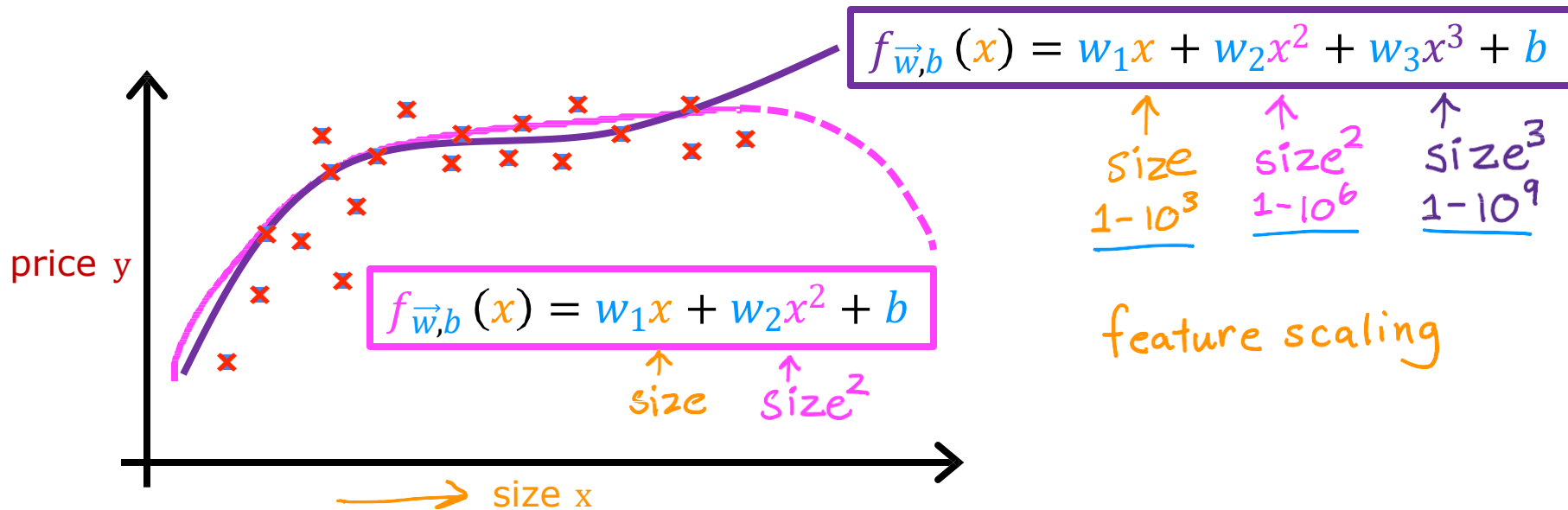
$$\beta = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

The above expression is the **OLS estimator** for β , which gives the coefficients that minimize the sum of squared residuals:

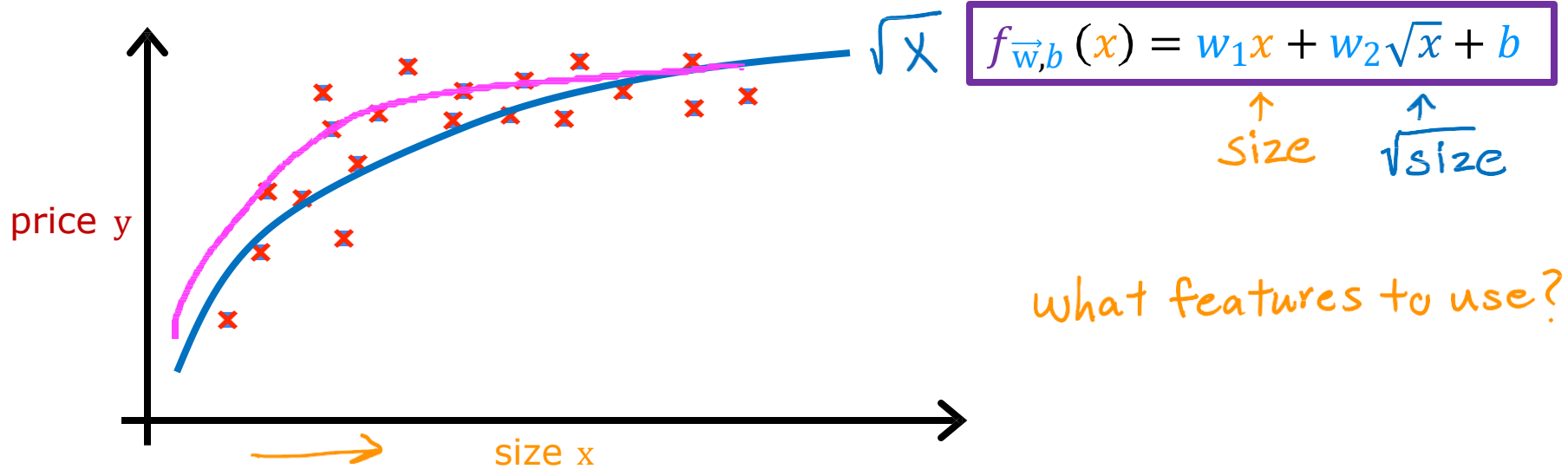
$$\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

This is the closed-form solution for the OLS estimator.

Polynomial regression



Choice of features



Assumptions for Lines Regression

1. Linear relationship:

1. Little or no multi-collinearity:

- It is assumed that there is little or no multicollinearity in the data.
- Multicollinearity occurs when the features (or independent variables) are not independent of each other.

2. Little or no autocorrelation:

- Another assumption is that there is little or no autocorrelation in the data.
- Autocorrelation occurs when the residual errors are not independent of each other.

Assumptions for LR

- **No outliers:**

- We assume that there are no outliers in the data. Outliers are data points that are far away from the rest of the data. Outliers can affect the results of the analysis.

- **Homoscedasticity:**

- Homoscedasticity describes a situation in which the error term (that is, the “noise” or random disturbance in the relationship between the independent variables and the dependent variable) is the same across all values of the independent variables.

