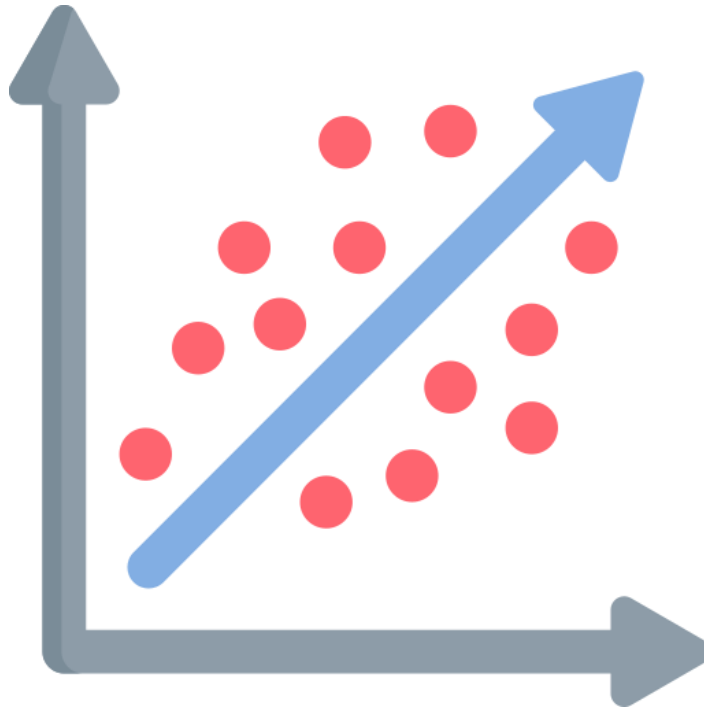# Linear Regression

# Outline

- [Linear Regression with One Variable](#)

- [Multiple Linear Regression](#)

- [Polynomial Regression](#)

- [Regression Practical Considerations](#)

Some of slides are based on the slides from
https://www.deeplearning.ai/courses/machine-learning-specialization/
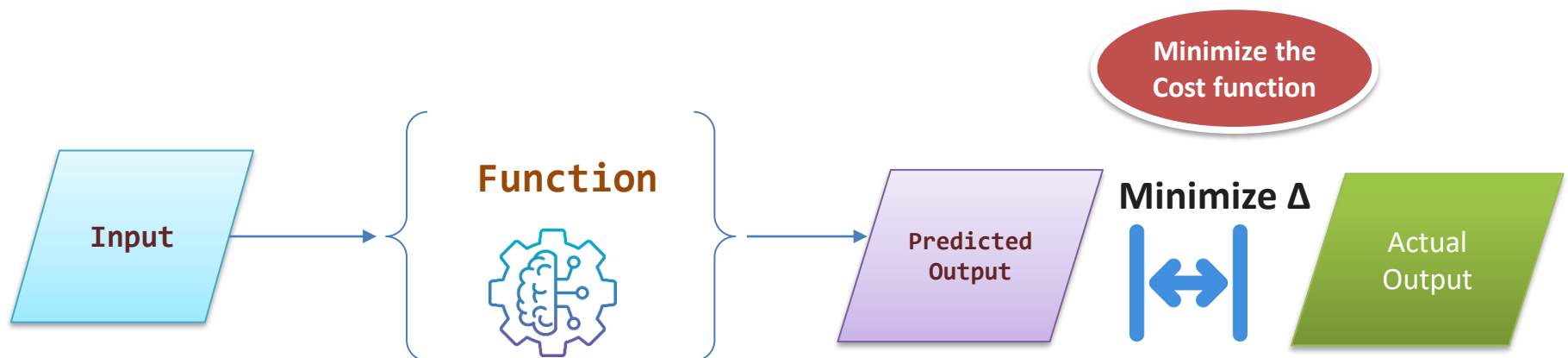
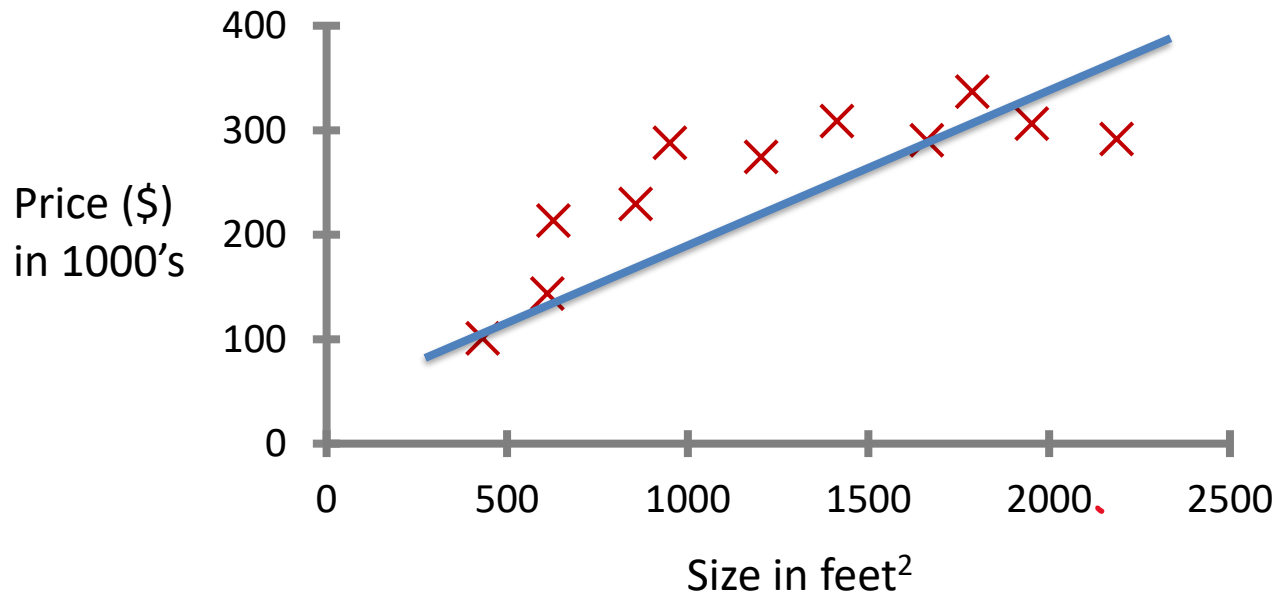# Linear Regression with One Variable

# ML: learn a **Function** that minimizes the cost

- Start with random function parameters
- Repeat intelligent guessing/approximation of the Function parameters such that the difference between the Predicted Output the Actual Output is reduced
  - i.e., minimize a Cost function a.k.a loss, or error function

**Minimize the Cost function**

**Function**

Input

**Minimize Δ**

Predicted Output

Actual Output

# Linear Regression with One Variable

**Housing price prediction**



Linear regression with one variable aka **Univariate linear regression**

- **Regression**: Predict continuous output value (e.g., price)

- Linear regression is used to predict the value of a variable based on the value of another variable(s)

- Linear regression **fits** a straight line that minimizes the discrepancies between predicted and actual output values

**Training set of housing prices**

| Size in feet² ($x$) | Price ($) in 1000's ($y$) |
|:---:|:---:|
| 2104 | 460 |
| 1416 | 232 |
| 1534 | 315 |
| 852 | 178 |
| … | … |

**Notation:**

$m$ = Number of training examples

$x$ = "input" variable / features

$y$ = "output" variable / "target" variable

$x^{(1)}$ = 2104

$x^{(2)}$ = 1416

$y^{(1)}$ = 460

$(x^{(i)}, y^{(i)})$ – the $i^{th}$ training example

**How do we represent $f$ ?**

$$f(x) = wx + b$$

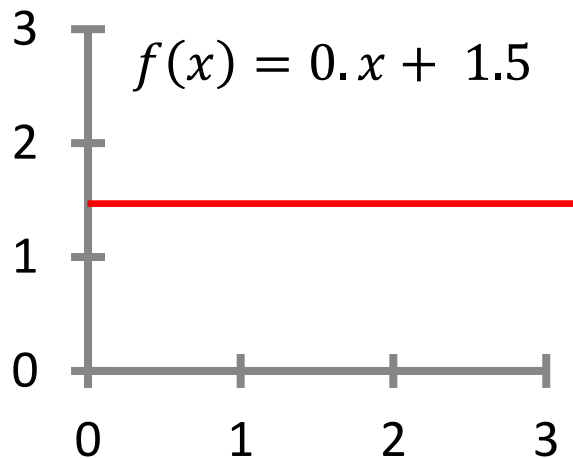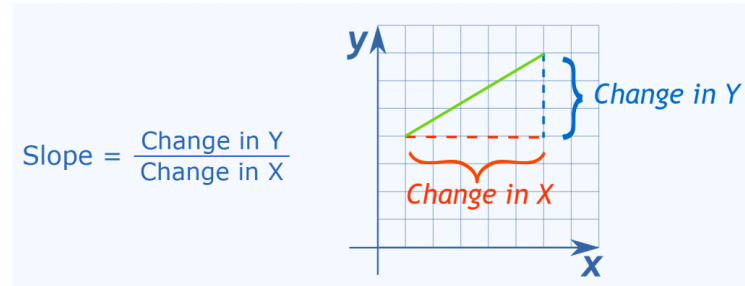$w, b$ are parameters (coefficients) to learn from the training set

- Given a training set, **learn a function $f$** so that $f(x)$ is a "good" predictor for the corresponding value of y

- Find $w, b$ parameters that **minimize** the error between predicted and actual values (i.e., minimize $\frac{1}{m}\sum_{i=1}^{m}(\hat{y}_i - y_i)^2$ for all dataset instances)

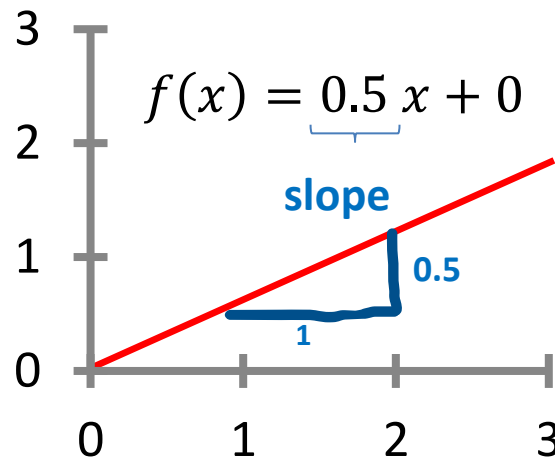# Univariate Linear Regression - Model Representation

$$f(x) = wx + b$$

How to choose $w$ and $b$ ?

- $w$ is the slope of the line
- $b$ is the y-intercept of the line



$$\text{Slope} = \frac{\text{Change in Y}}{\text{Change in X}}$$

$$f(x) = 0.\,x + 1.5$$

$$w = 0$$
$$b = 1.5$$

$$f(x) = 0.5\,x + 0$$

slope

0.5

1

$$w = 0.5$$
$$b = 0$$

$$f(x) = 0.5\,x + 1$$

y intercept (b = 1)

$$w = 0.5$$
$$b = 1$$

**Idea**: Find **w** and **b** so that $f(x)$ is close to $y$ for our training examples $(x, y)$

Find $w, b$:
$\hat{y}^{(i)}$ is close to $y^{(i)}$ for all $(x^{(i)}, y^{(i)})$

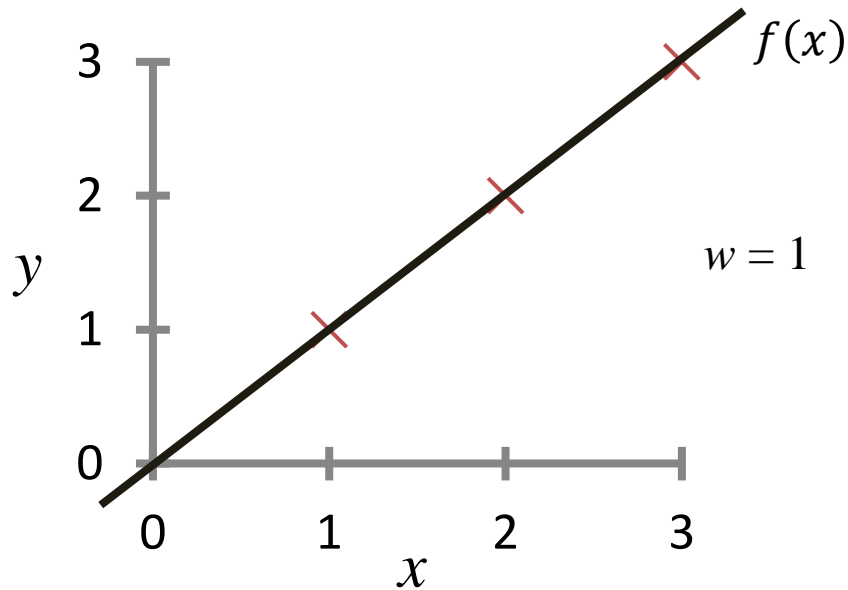**Cost (mean squared error)**
Function:
$$J(w, b) = \frac{1}{2m}\sum_{i=1}^{m}(\hat{y}_i - y_i)^2 = \frac{1}{2m}\sum_{i=1}^{m}(f(x^{(i)}) - y^{(i)})^2$$

**Goal**:     $\underset{w, b}{\text{minimize }} J(w, b)$

With **m** = number of training examples

# Visualizing the Cost function

$$f(x) = wx$$

For simplicity, let us assume our optimization objective is to $\underset{w,b}{\text{minimize}} \, J(w)$, thus, b = 0



$f(x)$

$w = 1$



$J(w)$

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

$$= \frac{1}{2m} (0^2 + 0^2 + 0^2) = 0$$



slides\regression.xlsx

# Visualizing the Cost function

$$f(x) = \boldsymbol{w}x$$



$f(x)$

$w = 0.5$

For simplicity, let us assume our optimization objective is to $\underset{w,b}{\text{minimize}}\, \boldsymbol{J(w)}$, thus, b = 0



$J(w)$

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

$$= \frac{1}{2m}((0.5-1)^2 + (1\text{-}2)^2 + (1.5\text{-}3)^2)$$

$$= \frac{1}{2\times3}(3.5) = \frac{3.5}{6} = 0.58$$

# Visualizing the Cost function

$$f(x) = wx$$



$w = 0$

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

$$= \frac{1}{2m} (1^2 + 2^2 + 3^2)$$

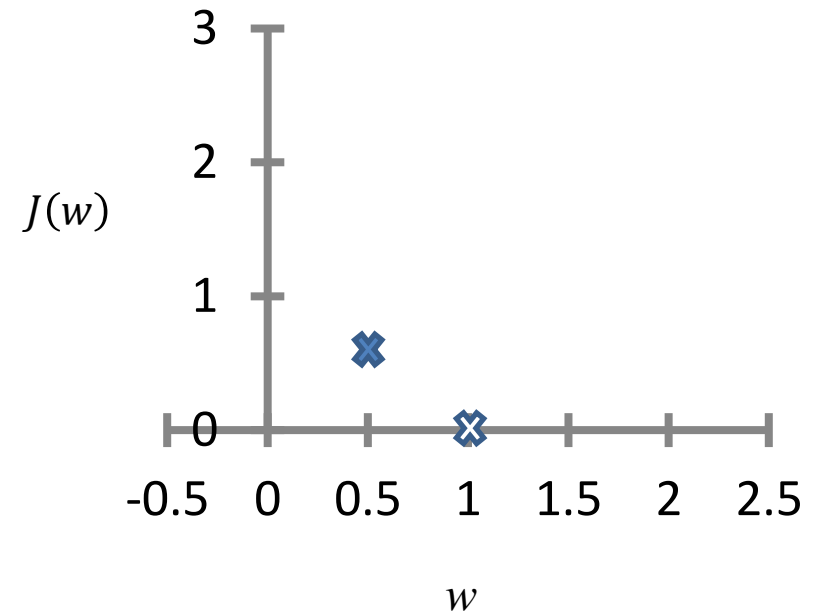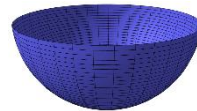$$= \frac{1}{2 \times 3} (14) = \frac{14}{6} = 2.3$$

For simplicity, let us assume our optimization objective is to $\underset{w,b}{\text{minimize}}\, J(w)$, thus, b = 0

$J(w_1)$



$w$



**J(w) cost function has bowl shape Convex function**

$J(w)$

$w$

# Visualizing the Cost function

$$f(x) = wx + b$$

$$J(w, b)$$



**J(w,b) cost function has bowl shape**
**Convex function**

Price ($) in 1000's



Size in feet$^2$ ($x$)

$$f(x) = 0.06x + 50$$

- The fact that the cost function squares the error ensures that the 'error surface' is **convex** like a soup bowl.
- It will always have a minimum that can be reached by following the **gradient** (i.e., the slope)
- Minimizing the cost function yields optimal values of **w** and **b**

`08.regression\02_cost_function.py`

$J(w, b)$

**Bowl shape
Convex function**

**Gradient Descent can be used to find the optimal w and b that minimizes that cost function**

*unsplash.com/photos/3m6vbzY69s4*

**MSE has only 1 global minimum**

$J(w, b)$

**Other cost functions may have multiple local minima**

14

# Gradient descent algorithm

Want to find **w** and **b** that minimize the cost function J $\quad \underset{w,b}{\text{minimize}}\, J(w, b)$

1. Initialize the values of **w** and **b** to some arbitrary values (say 0, 0)

2. Calculate the predicted values of **y** using the current values of **w** and **b**

3. Calculate the gradients of the cost function with respect to **w** and **b**

4. Update the values of **w** and **b** using the gradients and a learning rate

Learning Rate

Derivative of the Cost Function w.r.t **w**

$$w = w - \alpha \frac{d}{dw} J(w, b)$$

$$b = b - \alpha \frac{d}{db} J(w, b)$$

5. Repeat steps 2-4 until convergence (i.e., until the cost function converges to a minimum)

# Gradient descent algorithm

Gradient descent utilizes the partial derivative of the cost function with respect to **w** and $b$ to update **w** and $b$ parameters

**Repeat until convergence** {

$$w = w - \alpha \underbrace{\frac{1}{m} \sum_{i=1}^{m} \left( \hat{y}^{(i)} - y^{(i)} \right) . \, x^{(i)}}_{\frac{\partial J}{\partial w}}$$

$$b = b - \alpha \underbrace{\frac{1}{m} \sum_{i=1}^{m} \left( \hat{y}^{(i)} - y^{(i)} \right)}_{\frac{\partial J}{\partial b}}$$

}

(simultaneously update $w$ and $b$)

# Derivative 101

- Source

- **Derivatives**: it is all about slope!

$$\text{Slope} = \frac{\text{Change in Y}}{\text{Change in X}}$$

We can find an **average** slope between two points

$$\text{average slope} = \frac{24}{15}$$

17

# Derivative = slope at a point

- Fill in this slope formula: $\dfrac{\Delta y}{\Delta x} = \dfrac{f(x+\Delta x) - f(x)}{\Delta x}$

- Simplify it as best we can

- Then make **Δx** shrink towards zero.

**Example**

The slope formula is: $\dfrac{f(x+\Delta x) - f(x)}{\Delta x}$

Use $f(x) = x^2$: $\dfrac{(x+\Delta x)^2 - x^2}{\Delta x}$

Expand $(x+\Delta x)^2$ to $x^2+2x\,\Delta x+(\Delta x)^2$: $\dfrac{x^2 + 2x\,\Delta x + (\Delta x)^2 - x^2}{\Delta x}$

Simplify ($x^2$ and $-x^2$ cancel): $\dfrac{2x\,\Delta x + (\Delta x)^2}{\Delta x}$

Simplify more (divide through by $\Delta x$): $2x + \Delta x$

Then, **as Δx heads towards 0** we get: $2x$

Result: the derivative of **$x^2$** is **2x**

$$\frac{d}{dx}x^2 = 2x$$

In other words, the slope at x is 2x

# Interpretation of Derivative

- So what does $\dfrac{d}{dx}x^2 = 2x$ mean?

- It means that, for the function $x^2$, the slope or "rate of change" at any point is **2x**

- So, when **x=2** the slope is **2x = 4**

- Or when **x=5** the slope is **2x = 10**, and so on



08.regression\
03_gradient_decent_x_square_animation.py

# The Impact of Partial Derviative

- For simplicity, let us assume our optimization objective is to $\underset{w,b}{\text{minimize}}\ J(w)$, thus, b = 0

*Learing Rate*

$$w = w - \alpha \frac{d\,J(w)}{dw}$$

$$= w - \alpha\,(Positive\ Number)$$

Decrease $w$ by a certain value

**Positive Derivative**

repeat until convergence {

$$w = w - \alpha \frac{d\,J(w)}{dw}$$

}

New $w$    Old $w$

# The Impact of Partial Derviative

- For simplicity, let us assume our optimization objective is to $\underset{w,b}{\text{minimize}}$ $\textcolor{red}{J(w)}$, thus, b = 0



$$w = w - \alpha \, \frac{d \, \textcolor{red}{J(w)}}{dw}$$

$$= w - \alpha \, (Negative \; Number)$$

Increase $w$ by a certain value

# The Impact of Partial Derviative

- For simplicity, let us assume our optimization objective is to minimize $J(w)$, thus, b = 0
  $w,b$



Derivative = 0

$$w = w - \alpha \frac{d J(w)}{dw}$$

$$= w - \alpha \, (Zero)$$

$w$ remains the same, hence, gradient descent *converges*

# The Impact of Learning Rate
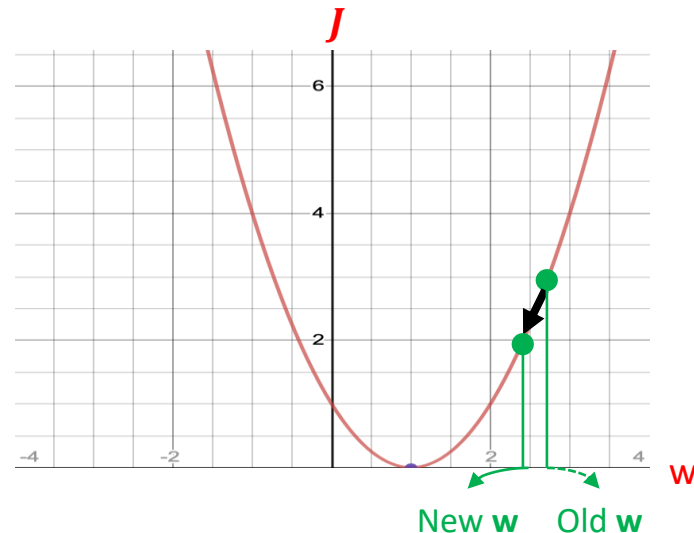


$$w = w - \alpha \, \frac{d\,J(w)}{dw}$$

$$= w - (Too\ Small\ Number)\, \frac{d\,J(w)}{dw}$$

*w* changes only a tiny bit on each step, hence, gradient descent *will render slow (will take more time to converge)*

$$w = w - \alpha \, \frac{d\,J(w)}{dw}$$

$$= w - (Too\ Large\ Number)\, \frac{d\,J(w)}{dw}$$



*w* changes a lot (and probably faster) on each step, hence, gradient descent *will potentially **overshoot the minimum** and, accordingly, fail to converge (or even diverge)*

# The Impact of Learning Rate

We can set α between 0 and 1 (say, 0.1, or a little more or less, hence, not very small or very large)



$$w = w - \alpha \, \frac{d \, J(w)}{dw}$$

$\alpha$ remains fix. As we approach the minimum, gradient descent will automatically start taking smaller steps (i.e., $w$ will start changing at a slower pace because the derivative will become less steep)

# Visualizing gradient descent algorithm

$$f(x) = wx + b$$

$$J(w, b)$$



**Contour Plot**

The optimal values of **w** and **b** are in the center of the inner most 'circle' of the Contour Plot





`08.regression\06_gradient_descent_visualization.py`

# Visualizing gradient descent algorithm

$$f(x) = wx + b$$

$$J(w, b)$$



- The contour plot shows $J(w,b)$ over a range of $w$ and $b$. The cost levels are represented by the rings
- The red crosses is the **path of gradient descent**. The path makes steady progress toward its goal. The initial steps are much larger than the steps near the goal.
- The optimal values of $w$ and $b$ are in the center of the inner most 'circle' of the Contour Plot

# Different modes of gradient descent

- **Batch Gradient Descent (BGD)**: Each iteration of gradient descent uses all the training examples

  – It provides a precise estimate of the gradient but can be computationally expensive, especially for large datasets

- **Stochastic Gradient Descent (SGD):** only a random sample from the dataset is used to compute the gradient in each iteration

  – It is computationally more efficient

  – However, it introduces more noise in the parameter updates, leading to more oscillations in the convergence path

- **Mini-Batch Gradient Descent**: divides the dataset into small mini-batches and computes the gradient using a mini-batch in each iteration

  – This approach combines the efficiency of stochastic gradient descent with the stability of batch gradient descent

# Multiple Linear Regression

# Linear Regression with multiple variables

**Multiple features (variables)**

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

$n = 4$

$w_1$
$w_2$
$w_3$
$w_4$

**Notation:**

$n$   = number of features

$x_j$   = $j^{th}$ feature

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

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

**Univariate Linear Regression:**    $f(x) = wx + b$

**Multiple Linear Regression:**

$$f(x) = w_1 x_1 + w_2 x_2 + \cdots + w_n x_n + b$$

**Example**

$$f(x) = 0.1 x_1 + 4 x_2 + 10 x_3 + -2 x_4 + 80$$

size    #bedrooms    #floors    years    base price

$$f_{\vec{w},b}(\vec{x}) = w_1 x_1 + w_2 x_2 + \cdots + w_n x_n + b$$

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

$$b \text{ is a number}$$

vector $\vec{x} = [x_1 \; x_2 \; x_3 \ldots x_n]$

$$f_{\vec{w},b}(\vec{x}) = \vec{w} \cdot \vec{x} + b = w_1 x_1 + w_2 x_2 + w_3 x_3 + \cdots + w_n x_n + b$$

dot product

# Multiple Linear Regression - Model Representation

## Parameters and features

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

$b$ is a number

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

linear algebra: count from 1

**NumPy**

```
w = np.array([1.0,2.5,-3.3])
b = 4
x = np.array([10,20,30])
```

code: count from 0

## Without vectorization

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

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

## Without vectorization

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

```
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
```

# Multiple Linear Regression - Model Representation

|  | Previous notation | Vector notation |
|---|---|---|

**Parameters**

$$w_1, \cdots, w_n$$
$$b$$

$$\vec{\mathrm{w}} = [w_1 \quad \cdots \quad w_n]$$
$$b$$

**Model**

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

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

*dot product*

**Cost function**

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

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

**Gradient descent**

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

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

# Vector representation of X, y and w

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

$$\vec{w} \cdot \vec{x} = \begin{bmatrix} w_0 \\ w_1 \\ \cdots \\ w_n \end{bmatrix} \begin{bmatrix} x_0^{(0)} & x_1^{(0)} & \cdots & x_n^{(0)} \\ x_0^{(1)} & x_1^{(1)} & \cdots & x_n^{(1)} \\ \cdots & \cdots & \cdots & \cdots \\ x_0^{(m-1)} & x_1^{(m-1)} & \cdots & x_n^{(m-1)} \end{bmatrix} = \begin{bmatrix} w_0 x_0^{(0)} + w_1 x_1^{(0)} + \cdots + w_n x_n^{(0)} \\ w_0 x_0^{(1)} + w_1 x_1^{(1)} + \cdots + w_n x_n^{(1)} \\ \cdots \\ w_0 x_0^{(m-1)} + w_1 x_1^{(m-1)} + \cdots + w_n x_n^{(m-1)} \end{bmatrix}$$

## Dot product operation

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{bmatrix} \begin{bmatrix} 2 \\ 3 \end{bmatrix} = \begin{bmatrix} \square & \square \\ \square & \square \\ \square & \square \end{bmatrix} = \begin{bmatrix} \square \\ \square \\ \square \end{bmatrix}$$

# Dot product operation

$$\boldsymbol{a} \cdot \boldsymbol{b} = \boldsymbol{a}^T \boldsymbol{b} = \sum_{i=0}^{n-1} a_i b_i$$

$$\boldsymbol{a} \cdot \boldsymbol{b} = \begin{bmatrix} a_0 \\ a_1 \\ a_2 \\ a_3 \end{bmatrix} \begin{bmatrix} b_0 \\ b_1 \\ b_2 \\ b_3 \end{bmatrix} = [a_0 b_0 + a_1 b_1 + a_2 b_2 + a_3 b_4]$$

$$\boldsymbol{a}^T \boldsymbol{b} = \begin{bmatrix} a_0 \\ a_1 \\ a_2 \\ a_3 \end{bmatrix} \begin{bmatrix} a_0 & a_1 & a_2 & a_3 \end{bmatrix} \quad \text{transform} \quad \begin{bmatrix} b_0 \\ b_1 \\ b_2 \\ b_3 \end{bmatrix} = [a_0 b_0 + a_1 b_1 + a_2 b_2 + a_3 b_4]$$

# Gradient Descent

### One feature

repeat {

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

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

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

simultaneously update $w$, $b$

}

### $n$ features ($n \geq 2$)

repeat {

$j=1$

$$w_1 = w_1 - \alpha \frac{1}{m} \sum_{i=1}^{m} \left(f_{\vec{w},b}\left(\vec{X}^{(i)}\right) - y^{(i)}\right) x_1^{(i)}$$

$\vdots$

$j=n$

$$\frac{\partial}{\partial w_1} J(\vec{w}, b)$$

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

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

simultaneously update
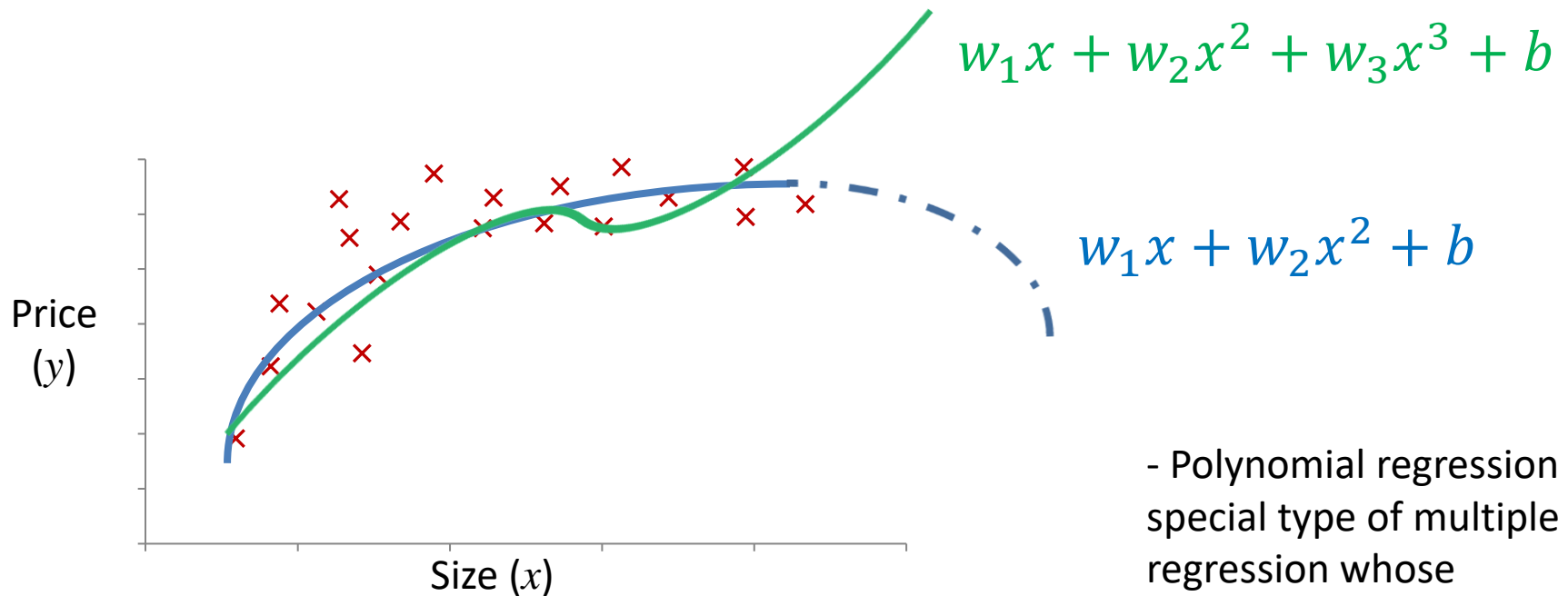$w_j$ (for $j = 1, \cdots, n$) and $b$

}

# Polynomial Regression

# Polynomial Regression

- When the relationship between the set of features and the target variable is not linear then we need to use polynomial regression of some degree

    – The degree of the polynomial is usually a hyperparameter of the model

# Polynomial Regression



$$w_1 x + w_2 x^2 + w_3 x^3 + b$$

$$w_1 x + w_2 x^2 + b$$

Price $(y)$

Size $(x)$

$$f(x) = w_1 x_1 + w_2 x_2 + w_3 x_3 + b$$
$$= w_1 (size) + w_2 (size)^2 + w_3 (size)^3 + b$$
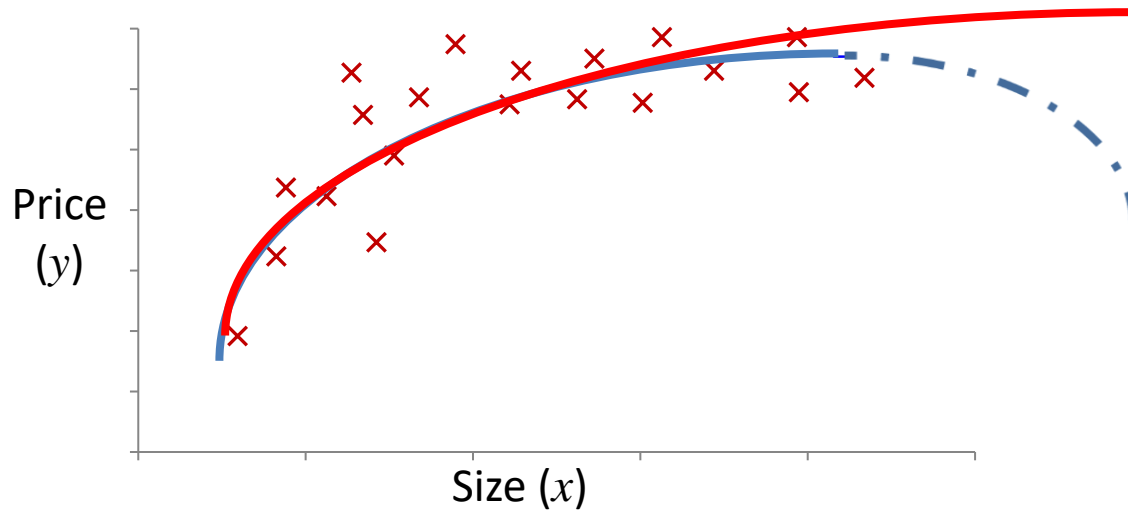
$$x_1 = (size)$$
$$x_2 = (size)^2$$
$$x_3 = (size)^3$$

- Polynomial regression is a special type of multiple regression whose independent variables are **powers** of variable X
- It is used to **approximate a curve** with unknown functional form
- For each additional power of X added to the model, the regression line will have one more bend

# Polynomial Regression



$$f(x) = w_1(size) + w_2(size)^2 + b$$

$$f(x) = w_1(size) + w_2\sqrt{(size)} + b$$

# Training Polynomial Regression Model

- Add **Polynomial Features** then train a Linear Regression Model

- For example, if we have two features then our new features will be: $x_1$, $x_2$, $x_1x_2$, $x_1^2$, $x_2^2$, and the features matrix will be:

$$X = \begin{pmatrix} 1 & x_{11} & x_{12} & x_{11}x_{12} & x_{11}^2 & x_{12}^2 \\ 1 & x_{21} & x_{22} & x_{21}x_{22} & x_{21}^2 & x_{22}^2 \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n1} & x_{n2} & x_{n1}x_{n2} & x_{n1}^2 & x_{n2}^2 \end{pmatrix}$$
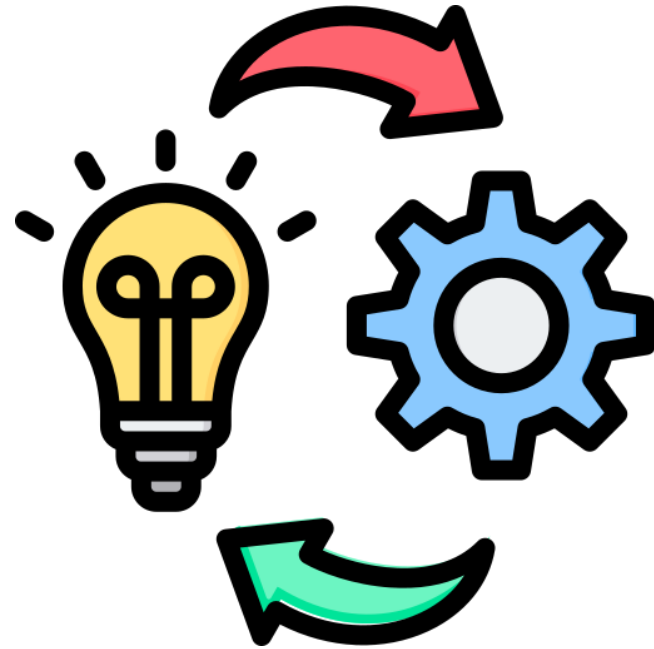
- Scikit-learn provides the transformer [PolynomialFeatures](#) that creates the features matrix consisting of all the polynomial combinations of the features up to a **specified degree**

- Then we can use the created features matrix to train a linear regression model

# Polynomial Regression

- We can treat the powers of x: x, $x^2$, …, $x^d$, as distinct independent variables
  - Then, polynomial regression becomes a special case of multiple linear regression, since the model is still linear in the parameters that need to be estimated

- Therefore, we can find the optimal parameters $\vec{w}$ using gradient descent
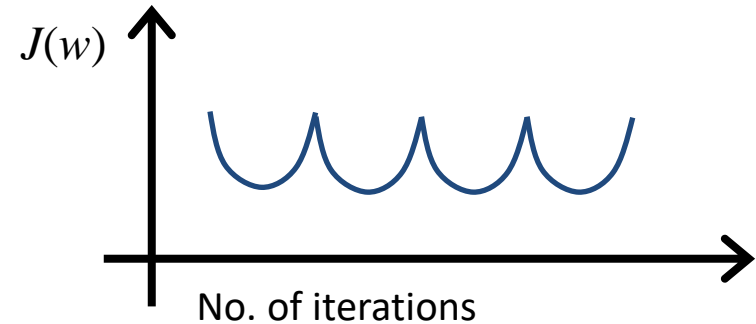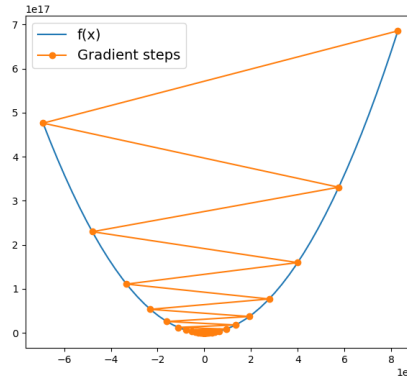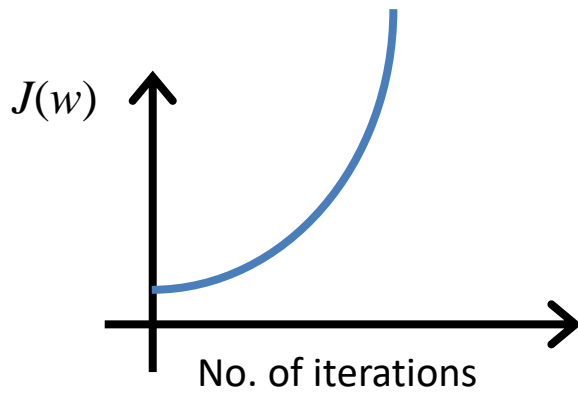
# Regression Practical Considerations

# Gradient descent in practice:

- Debugging Learning rate
- Feature Scaling
- Regularization
- Regression Evaluation

**Learning rate:**  Gradient descent not working

Use smaller $\alpha$

$J(w)$

No. of iterations

1e17

- f(x)
- Gradient steps

$J(w)$

No. of iterations

- For sufficiently small $\alpha$, $J(w)$ should decrease on every iteration
- If $\alpha$ is too small, gradient descent can be slow to converge
- If $\alpha$ is too large: $J(w)$ may not decrease on every iteration; may not converge
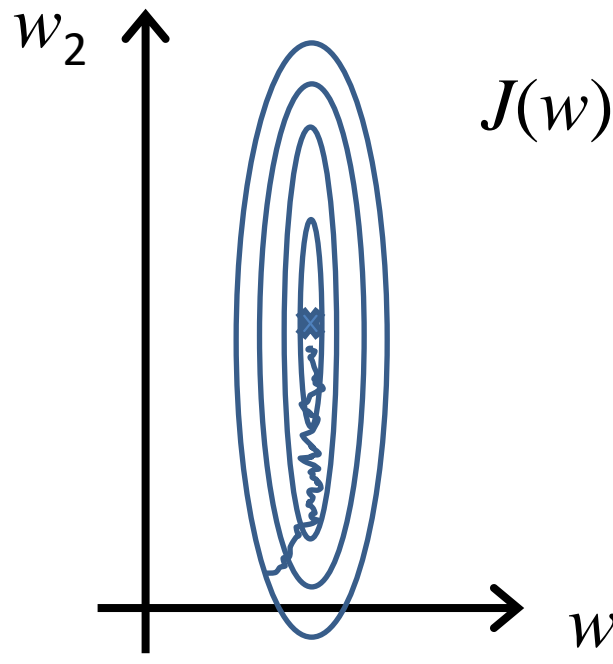
To choose $\alpha$, try

..., 0.001, 0.003, 0.01, 0.03, 0.1, 0.3, 1, ...

**Feature Scaling:** divide the input values by the range (i.e. the maximum value minus the minimum value) of the input variable, resulting in a new range of just 1.

The idea: Make sure features are on a similar scale. So that the gradient descent converges faster.
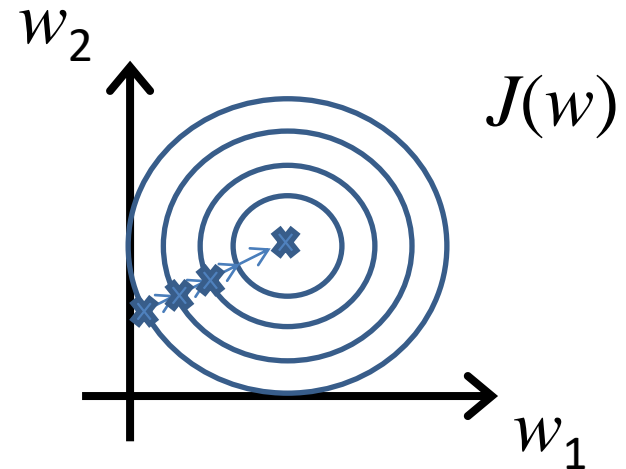
E.g. $x_1$ = size (0-2000 feet²)

$x_2$ = number of bedrooms (1-5)

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

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



Rule-of-thumb: Get every feature into approximately $-1 \leq x_i \leq 1$ range, $-0.5 \leq x_i \leq 0.5$, or other similar small ranges.

# Mean normalization

- Replace $x_i$ to make features have approximately zero mean: $\qquad x_i = \dfrac{x_i - \mu_i}{s_i}$

  Where $\mu_i$ is the **average** of all the values for feature ($i$) (**in the training set**) and $s_i$ is the standard deviation.
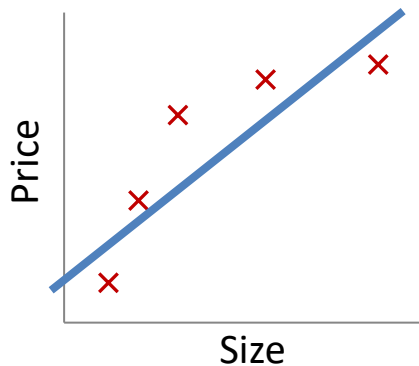
  - e.g.,

  $x_1 = \dfrac{size - 1000}{400}$ (average size of the houses is 1000, and standard deviation is 400)

  $x_2 = \dfrac{\#bedrooms - 2}{3}$ (average # of bedrooms is 2, and the standard deviation 3)
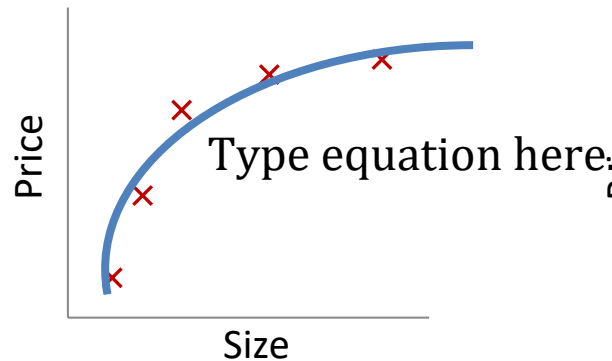
# Regularization

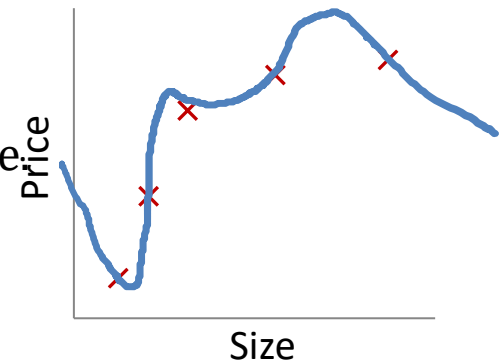- The problem of overfitting

Example: Linear regression (housing prices)



$$w_0 + w_1 x$$

"underfit/high bias"

$$w_0 + w_1 x + w_2 x^2$$

"just right"

$$w_0 + w_1 x + w_2 x^2 + w_3 x^3 + w_4 x^4$$

"overfit/high variance"

**Overfitting:** If we have too many features, the learned function may fit the training set very well ( $J(w) = \frac{1}{2m} \sum_{i=1}^{m} (h_w(x^{(i)}) - y^{(i)})^2 \approx 0$ ) but fail to generalize to new examples (predict prices on new examples)

**Addressing overfitting:**

$x_1 =$ size of house
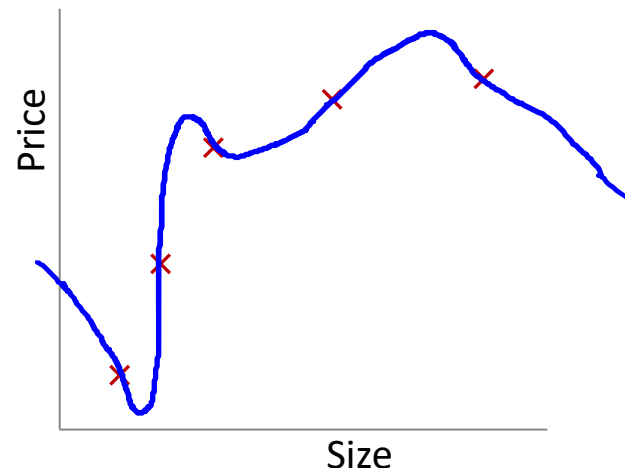$x_2 =$ no. of bedrooms
$x_3 =$ no. of floors
$x_4 =$ age of house
$x_5 =$ average income in neighborhood
$x_6 =$ kitchen size
$\vdots$

$x_{100}$

**Addressing overfitting:**

Options:

1. Reduce number of features.
   - Manually select which features to keep
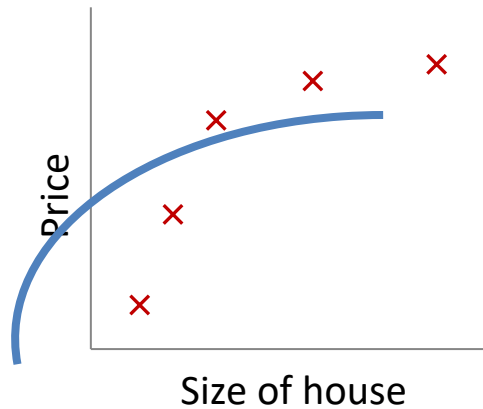   - Use feature selection algorithm
2. Regularization.
   - Keep all the features, but reduce magnitude/values of parameters $w_j$
   - Works well when we have a lot of features, each of which contributes a bit to predicting $y$
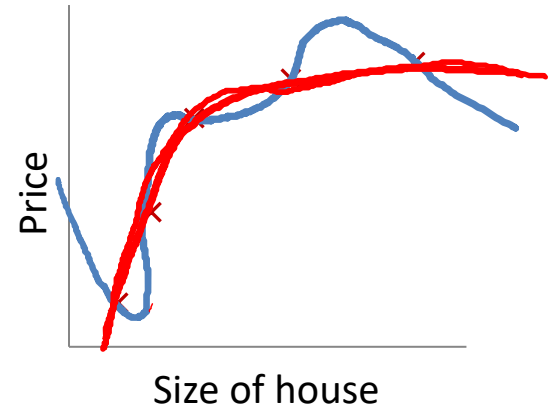
# Regularization

- Cost function

**Intuition**



$$w_0 + w_1 x + w_2 x^2$$

$$w_0 + w_1 x + w_2 x^2 + w_3 x^3 + w_4 x^4$$

Suppose we penalize and make $w_3, w_4$ really small

$$J(w) = \frac{1}{2m} \sum_{i=1}^{m} (f_w(x^{(i)}) - y^{(i)})^2 + 1000\, w_3^2 + 1000\, w_4^2$$

$$w_3 \approx 0, \quad w_4 \approx 0$$

# Regularization

Small values for parameters $w_1, w_2, \ldots, w_n$
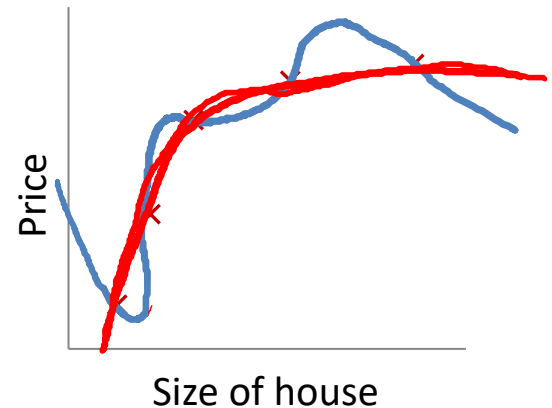- "Simpler/smoother" function
- Less prone to overfitting

Housing:
- Features: $x_1, x_2, \ldots, x_{100}$
- Parameters: $w_1, w_2, \ldots, w_{100}$

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

$$J(w) = \frac{1}{2m} \left[ \sum_{i=1}^{m} (f_w(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^{n} w_j^2 \right]$$

$$\min_{w} J(w)$$

Price

Size of house

# Lasso (L1) Regularization

$$\frac{1}{2m} \sum_{i=1}^{m} (y - Xw)^2 + alpha \sum_{j=1}^{p} |w_j|$$

# Ridge (L2) Regularization

$$\sum_{i=1}^{n} (y - Xw)^2 + alpha \sum_{j=1}^{p} w_j^2$$

# Regularized linear regression

$$J(w) = \frac{1}{2m}\left[\sum_{i=1}^{m}(f_w(x^{(i)}) - y^{(i)})^2 + \lambda\sum_{j=1}^{n} w_j^2\right]$$

$$\min_{w} J(w)$$

**Gradient descent**

Repeat {

$$\overbrace{\phantom{xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx}}^{\frac{\partial}{\partial w_0}J(w)}$$

$$b = b - \alpha\frac{1}{m}\sum_{i=1}^{m}(f_w(x^{(i)}) - y^{(i)})\, x_0^{(i)}$$

$$w_j = w_j - \alpha\left[\frac{1}{m}\sum_{i=1}^{m}(f_w(x^{(i)}) - y^{(i)})\, x_j^{(i)} - \frac{\lambda}{m}w_j\right]$$

$$\underbrace{\phantom{xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx}}_{\frac{\partial}{\partial w_j}J(w)}$$ "Regularized"

} $(j = 1, 2, 3, \ldots, n)$

$$w_j = w_j\left(1 - \alpha\frac{\lambda}{m}\right) - \alpha\frac{1}{m}\sum_{i=1}^{m}(f_w(x^{(i)}) - y^{(i)})\, x_j^{(i)}$$

$$1 - \alpha\frac{\lambda}{m} < 1$$

# Regression Evaluation

- Performance measured by

  - Mean Squared Error (MSE)          $MSE = \frac{1}{n}\sum(y - \hat{y})^2$

  - Root-Mean-Squared-Error (RMSE)    $RMSE = \sqrt{\frac{(y-\hat{y})^2}{n}}$

  - Mean-Absolute-Error (MAE)         $MAE = \frac{1}{n}\sum|y - \hat{y}|$

  - …others