

Linear Regression

Machine learning algorithms

	Supervised Learning	Unsupervised Learning
Discrete	Classification	Clustering
Continuous	Regression	Dimensionality reduction

Nearest neighbor classifier

- **Training data**

$$(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(N)}, y^{(N)})$$

- **Learning**

Do nothing.

- **Testing**

$$h(x) = y^{(k)}, \text{ where } k = \operatorname{argmin}_i D(x, x^{(i)})$$



Instance/Memory-based Learning

1. A distance metric

- Continuous? Discrete? PDF? Gene data? Learn the metric?

2. How many nearby neighbors to look at?

- 1? 3? 5? 15?

3. A weighting function (optional)

- Closer neighbors matter more

4. How to fit with the local points?

- Kernel regression

Validation set

- Splitting training set: *A fake* test set to tune hyper-parameters

```
# assume we have Xtr_rows, Ytr, Xte_rows, Yte as before
# recall Xtr_rows is 50,000 x 3072 matrix
Xval_rows = Xtr_rows[:1000, :] # take first 1000 for validation
Yval = Ytr[:1000]
Xtr_rows = Xtr_rows[1000:, :] # keep last 49,000 for train
Ytr = Ytr[1000:]

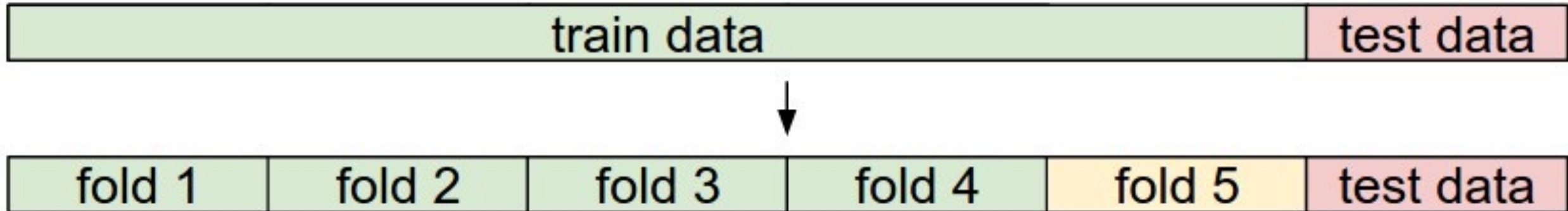
# find hyperparameters that work best on the validation set
validation accuracies = []
for k in [1, 3, 5, 10, 20, 50, 100]:

    # use a particular value of k and evaluation on validation data
    nn = NearestNeighbor()
    nn.train(Xtr_rows, Ytr)
    # here we assume a modified NearestNeighbor class that can take a k as input
    Yval_predict = nn.predict(Xval_rows, k = k)
    acc = np.mean(Yval_predict == Yval)
    print 'accuracy: %f' % (acc,)

# keep track of what works on the validation set
validation accuracies.append((k, acc))
```

Cross-validation

- 5-fold cross-validation -> split the training data into 5 equal folds
- 4 of them for training and 1 for validation



Things to remember

- Supervised Learning
 - Training/testing data; classification/regression; Hypothesis
- k-NN
 - Simplest learning algorithm
 - With sufficient data, very hard to beat “strawman” approach
- Kernel regression/classification
 - Set k to n (number of data points) and chose kernel width
 - Smoother than k-NN
- Problems with k-NN
 - Curse of dimensionality
 - Not robust to irrelevant features
 - Slow NN search: must remember (very large) dataset for prediction



Today's plan: Linear Regression

- Model representation
- Cost function
- Gradient descent
- Features and polynomial regression
- Normal equation

Linear Regression

- **Model representation**
- Cost function
- Gradient descent
- Features and polynomial regression
- Normal equation

Regression

real-valued output

Training set

Learning Algorithm

x

h

y

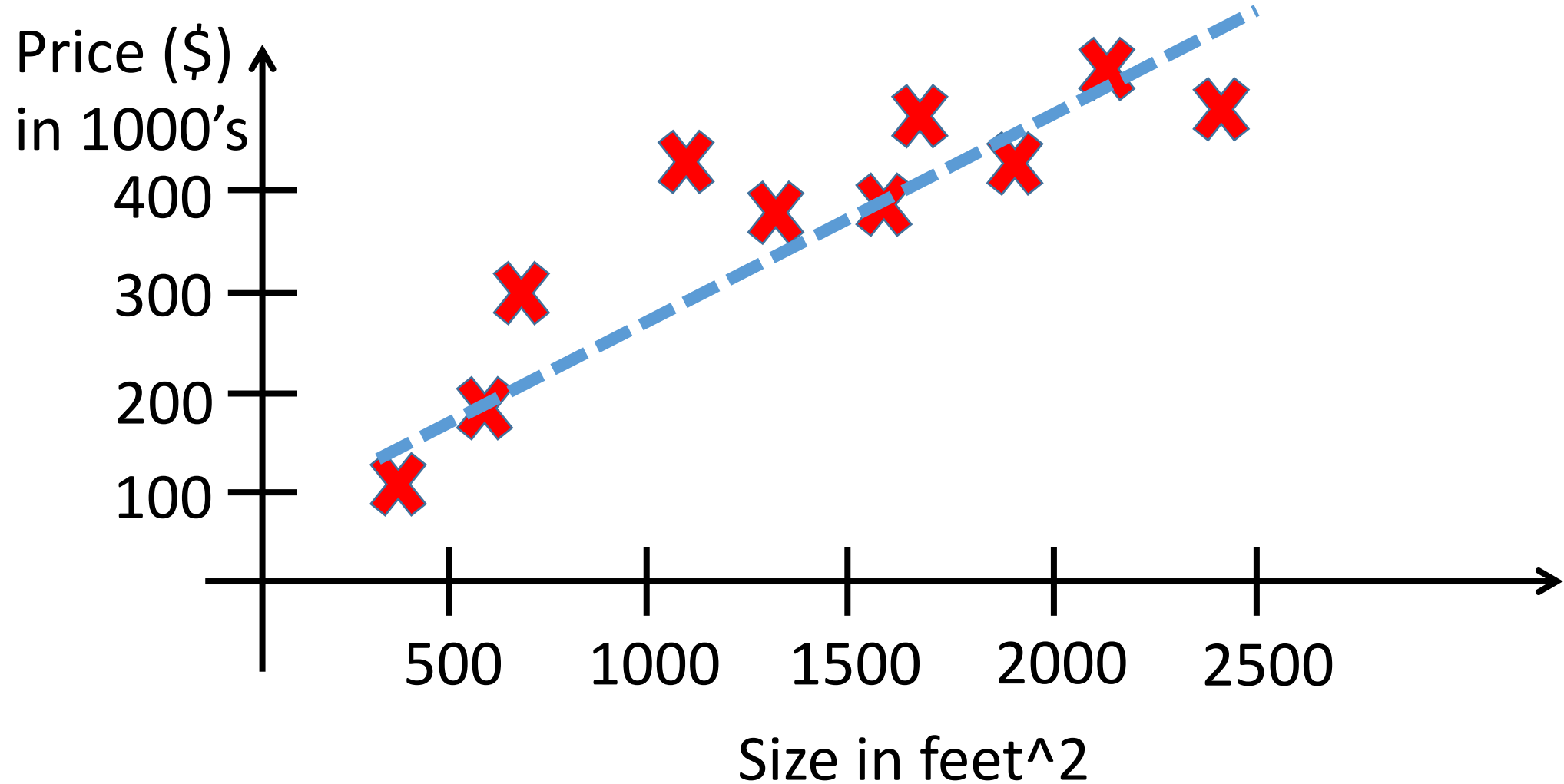
Size of house

Hypothesis

Estimated price

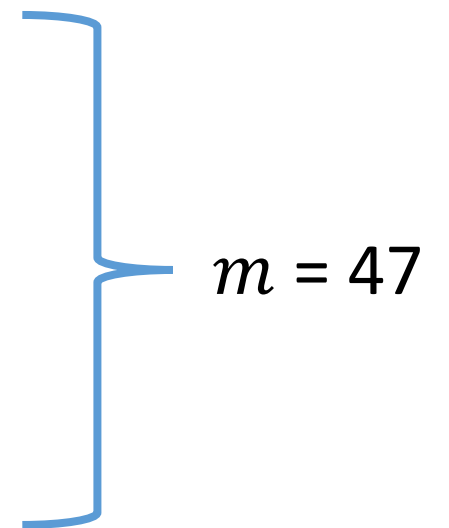


House pricing prediction



Training set

Size in feet ² (x)	Price (\$) in 1000's (y)
2104	460
1416	232
1534	315
852	178
...	...



$m = 47$

- Notation:

- m = Number of training examples
- x = Input variable / features
- y = Output variable / target variable
- (x, y) = One training example
- $(x^{(i)}, y^{(i)}) = i^{th}$ training example

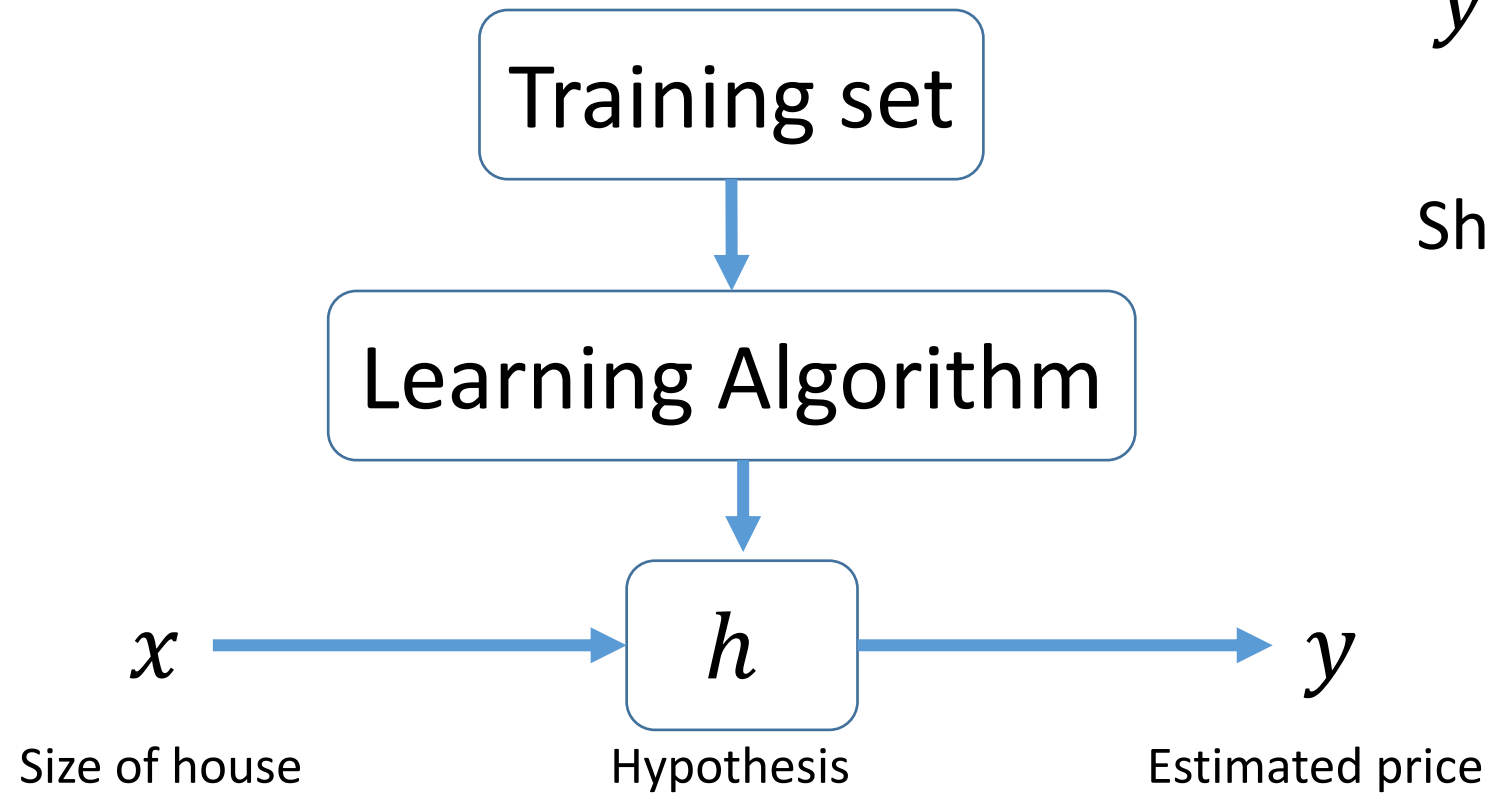
Examples:

$$x^{(1)} = 2104$$

$$x^{(2)} = 1416$$

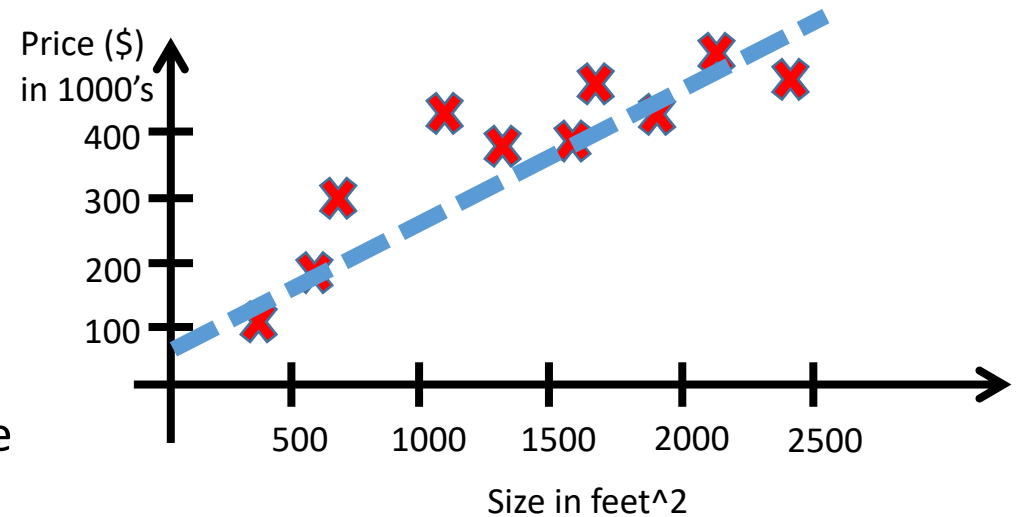
$$y^{(1)} = 460$$

Model representation



$$y = h_{\theta}(x) = \theta_0 + \theta_1 x$$

Shorthand $h(x)$



Univariate linear regression

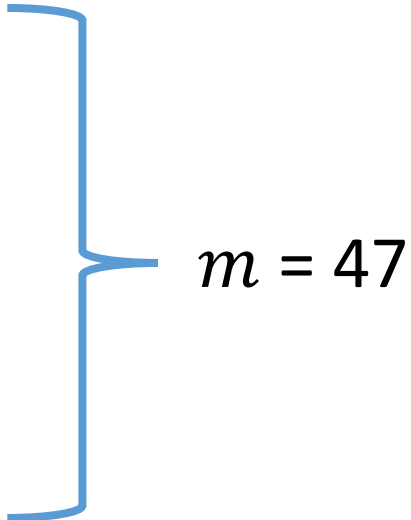
Slide credit: Andrew Ng

Linear Regression

- Model representation
- **Cost function**
- Gradient descent
- Features and polynomial regression
- Normal equation

Training set

Size in feet ² (x)	Price (\$) in 1000's (y)
2104	460
1416	232
1534	315
852	178
...	...



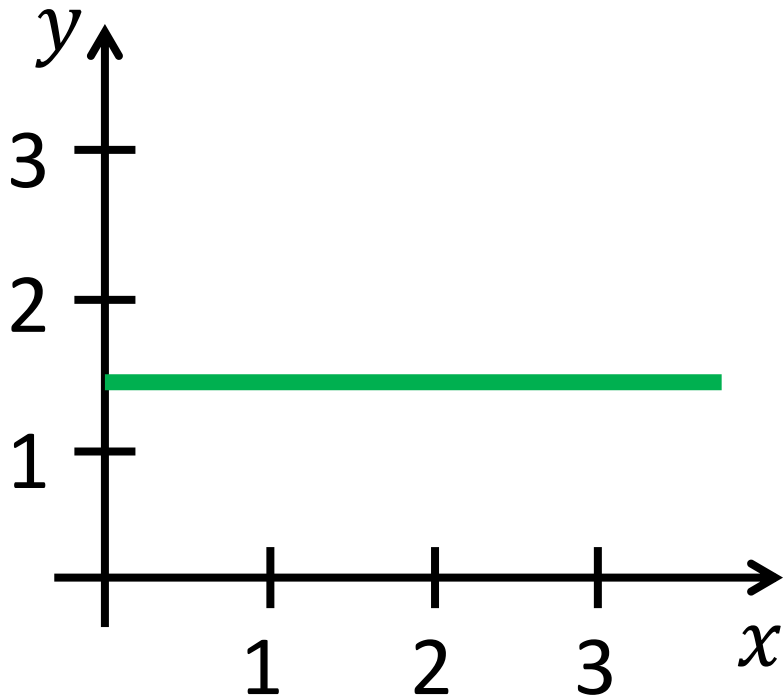
$m = 47$

- Hypothesis $h_{\theta}(x) = \theta_0 + \theta_1 x$

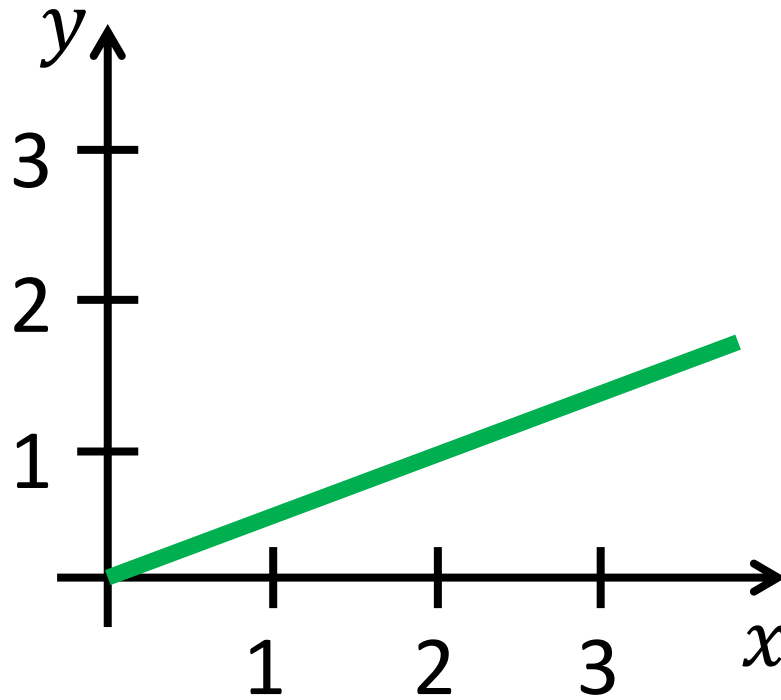
θ_0, θ_1 : parameters/weights

How to choose θ_i 's?

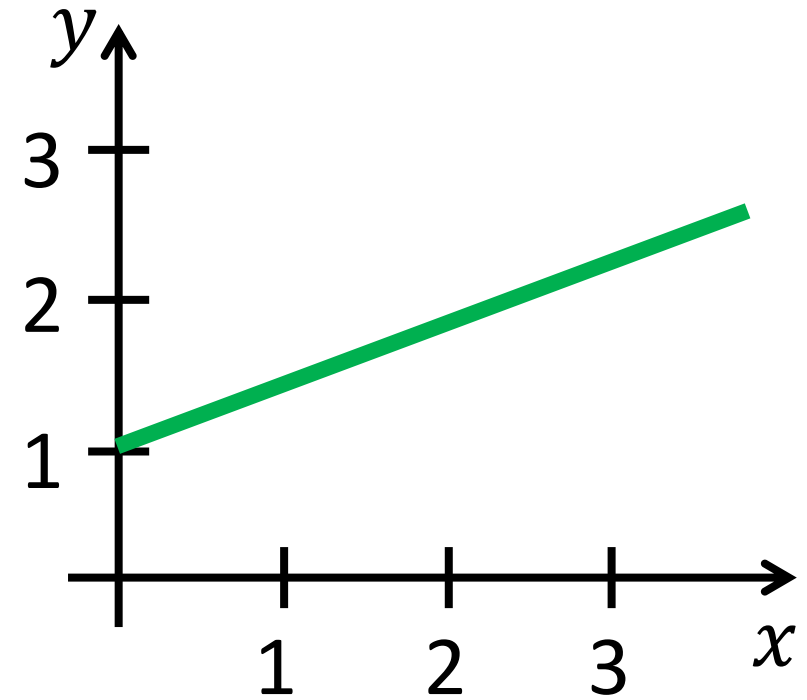
$$h_{\theta}(x) = \theta_0 + \theta_1 x$$



$$\begin{aligned}\theta_0 &= 1.5 \\ \theta_1 &= 0\end{aligned}$$



$$\begin{aligned}\theta_0 &= 0 \\ \theta_1 &= 0.5\end{aligned}$$

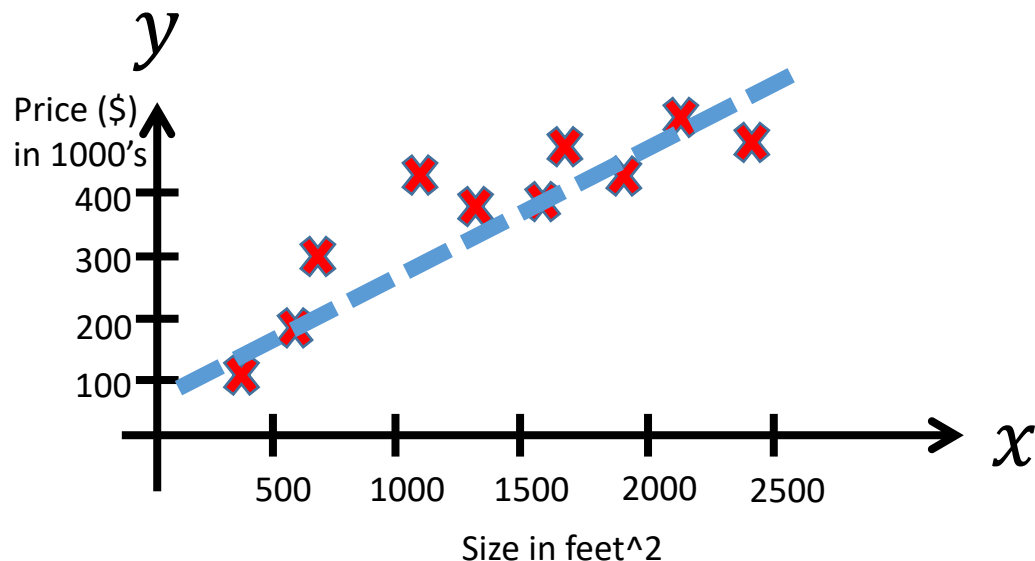


$$\begin{aligned}\theta_0 &= 1 \\ \theta_1 &= 0.5\end{aligned}$$

Cost function

- Idea:

Choose θ_0, θ_1 so that $h_\theta(x)$ is close to y for our training example (x, y)



$$\underset{\theta_0, \theta_1}{\text{minimize}} \quad \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$$

$$h_\theta(x^{(i)}) = \theta_0 + \theta_1 x^{(i)}$$

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

$$\underset{\theta_0, \theta_1}{\text{minimize}} \quad \boxed{J(\theta_0, \theta_1)} \quad \text{Cost function}$$

Simplified

- **Hypothesis:**

$$h_{\theta}(x) = \theta_0 + \theta_1 x \longrightarrow$$

- **Hypothesis:**

$$h_{\theta}(x) = \theta_1 x \quad \theta_0 = 0$$

- **Parameters:**

$$\theta_0, \theta_1 \longrightarrow$$

- **Parameters:**

$$\theta_1$$

- **Cost function:**

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

- **Cost function:**

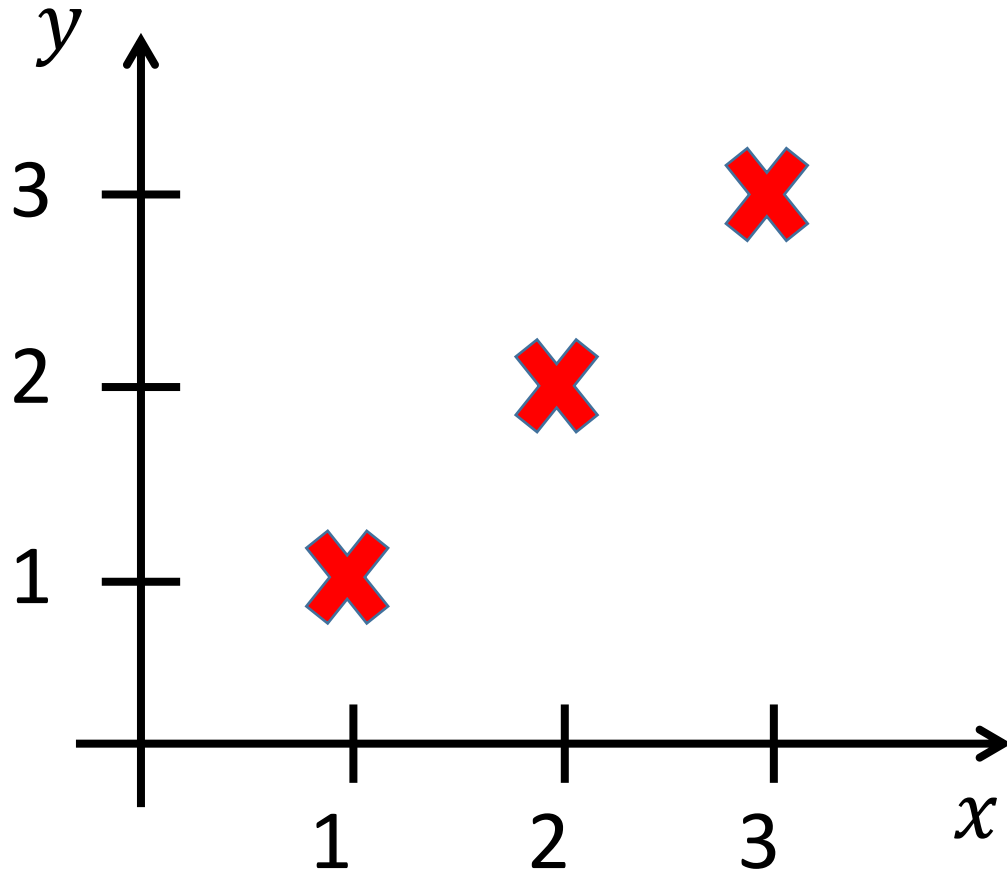
- **Goal:**

$$\underset{\theta_0, \theta_1}{\text{minimize}} J(\theta_0, \theta_1) \longrightarrow$$

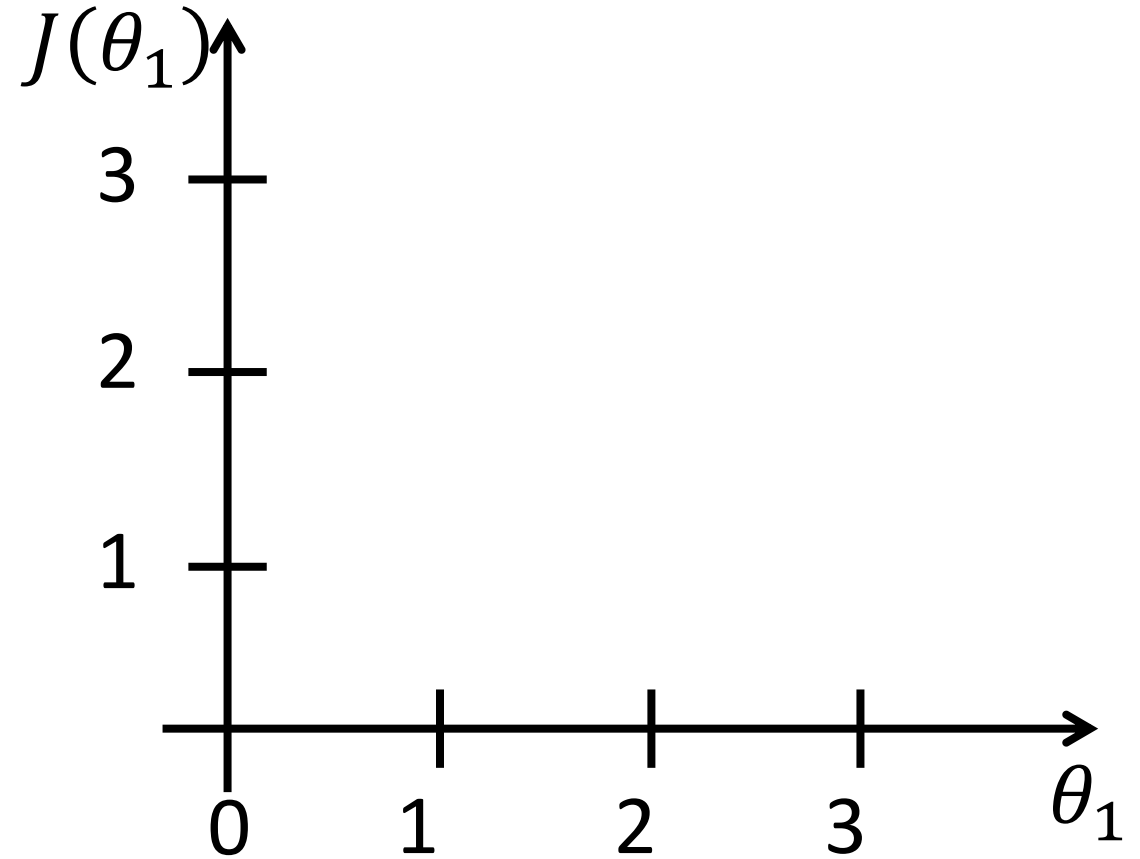
- **Goal:**

$$\underset{\theta_0, \theta_1}{\text{minimize}} J(\theta_1)$$

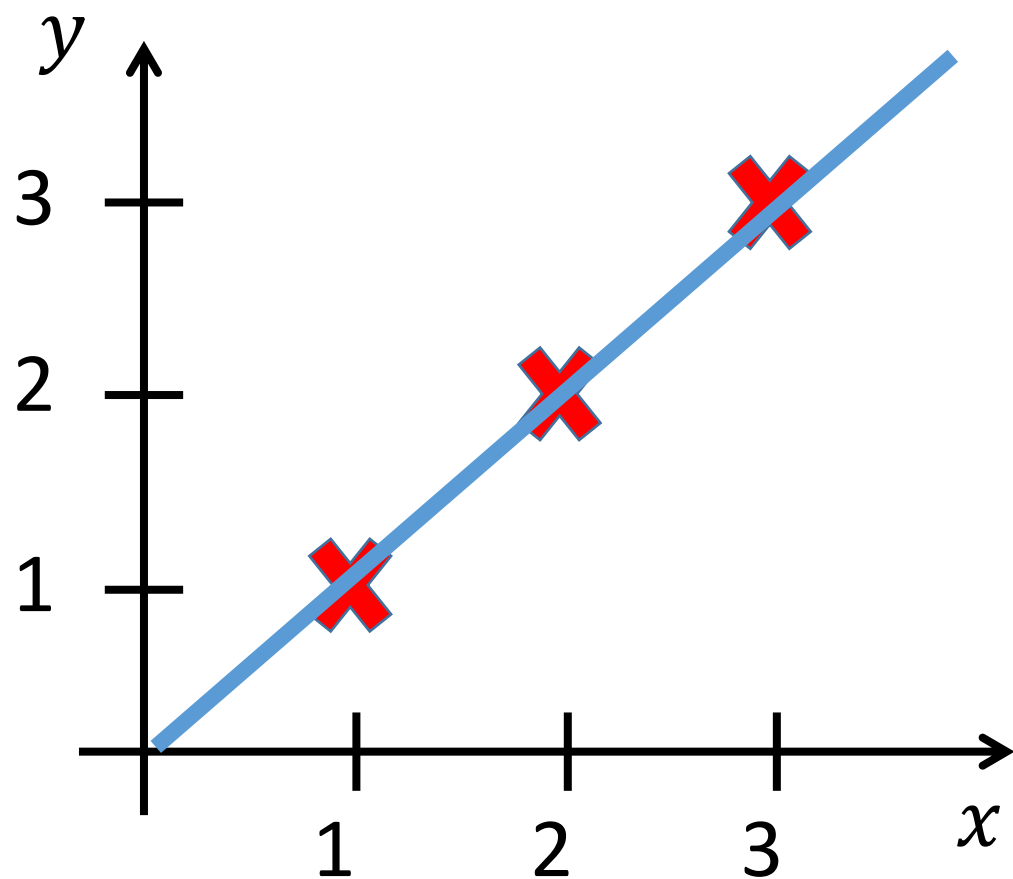
$h_{\theta}(x)$, function of x



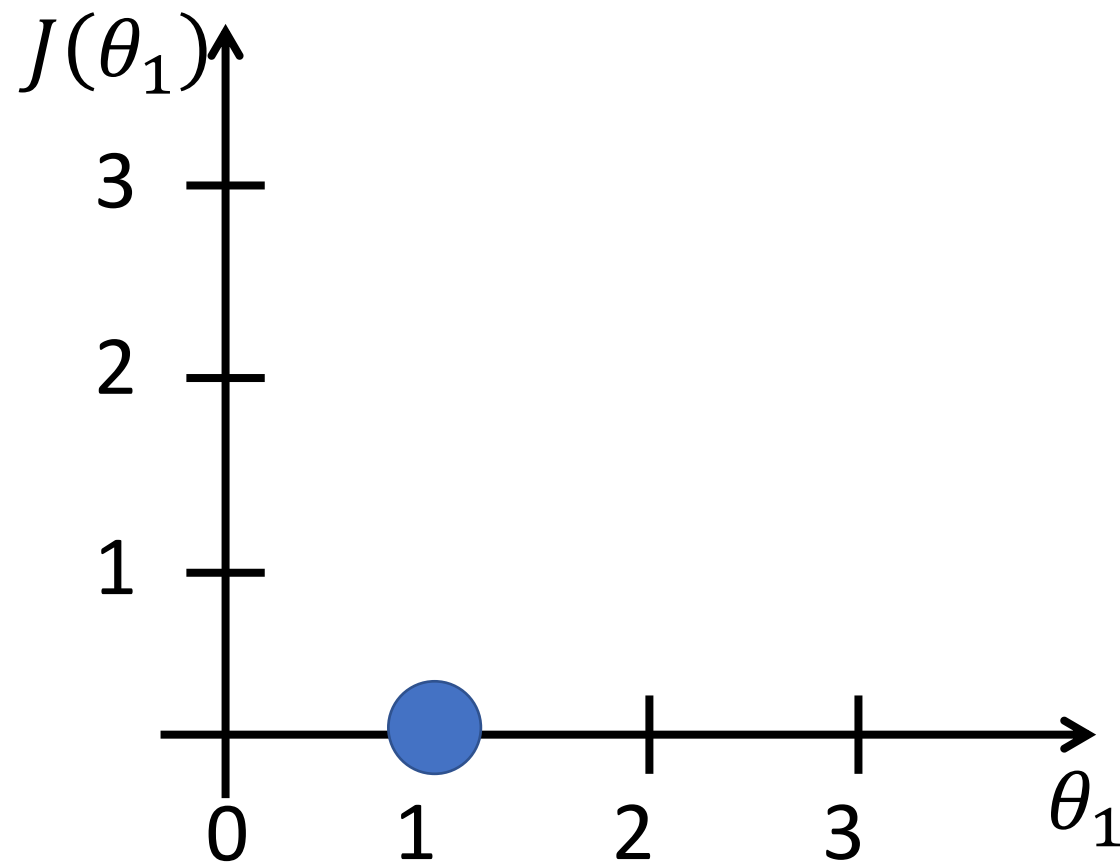
$J(\theta_1)$, function of θ_1



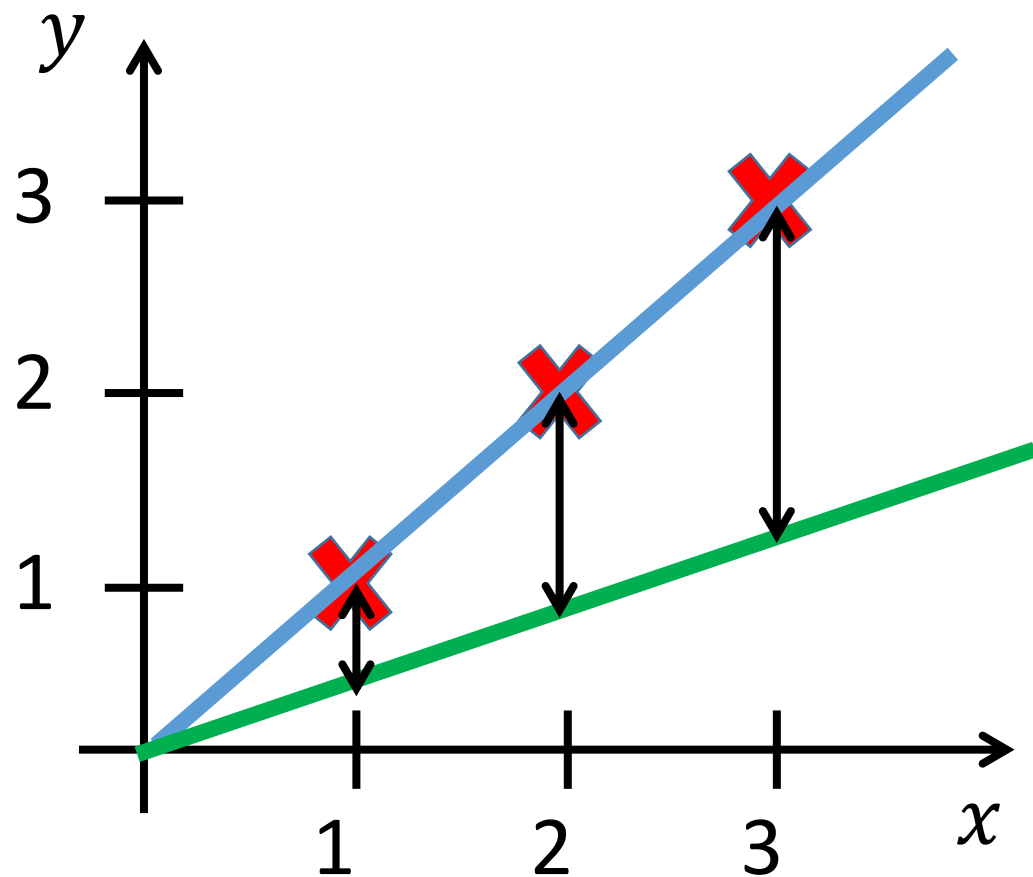
$h_{\theta}(x)$, function of x



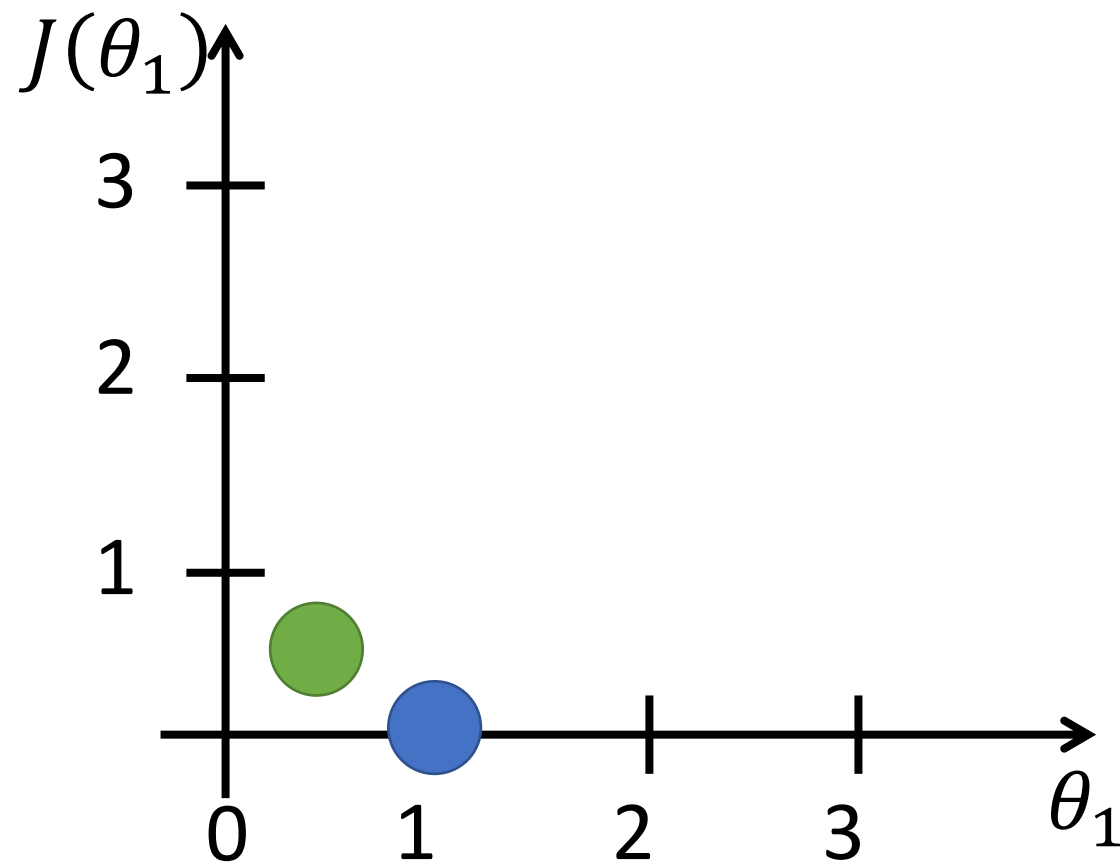
$J(\theta_1)$, function of θ_1



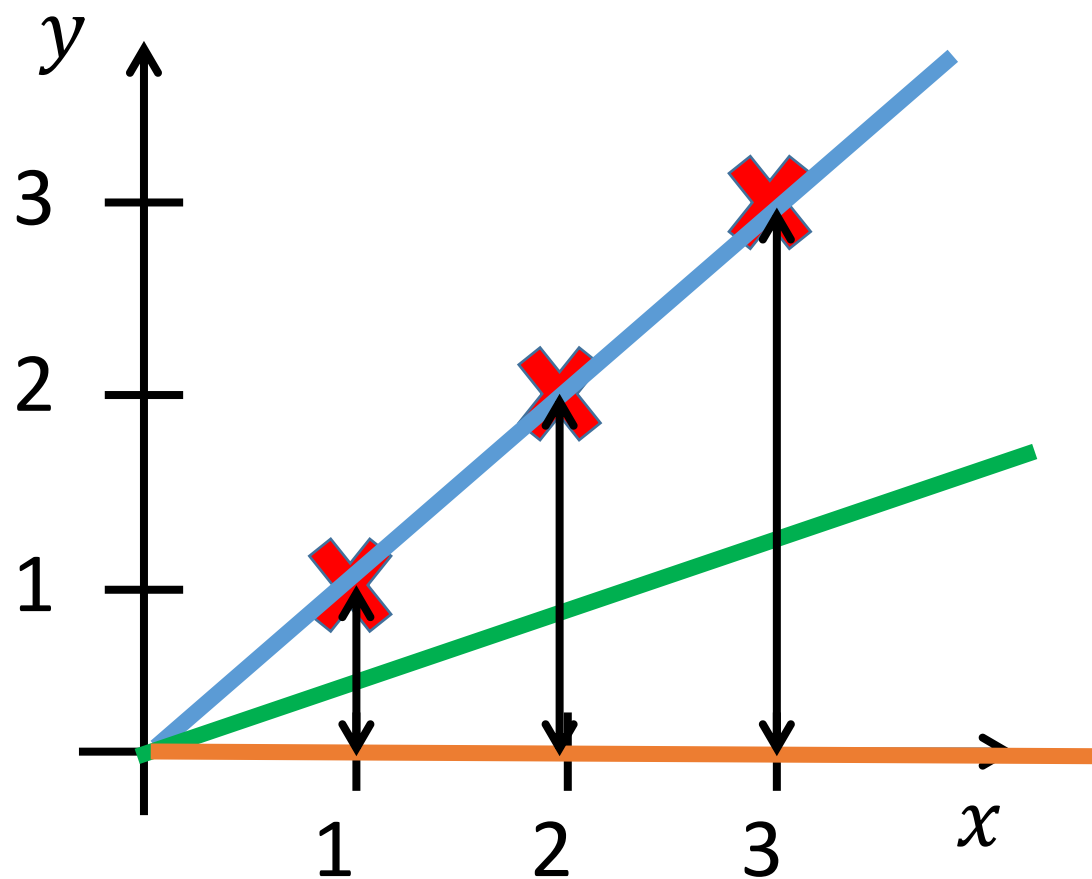
$h_{\theta}(x)$, function of x



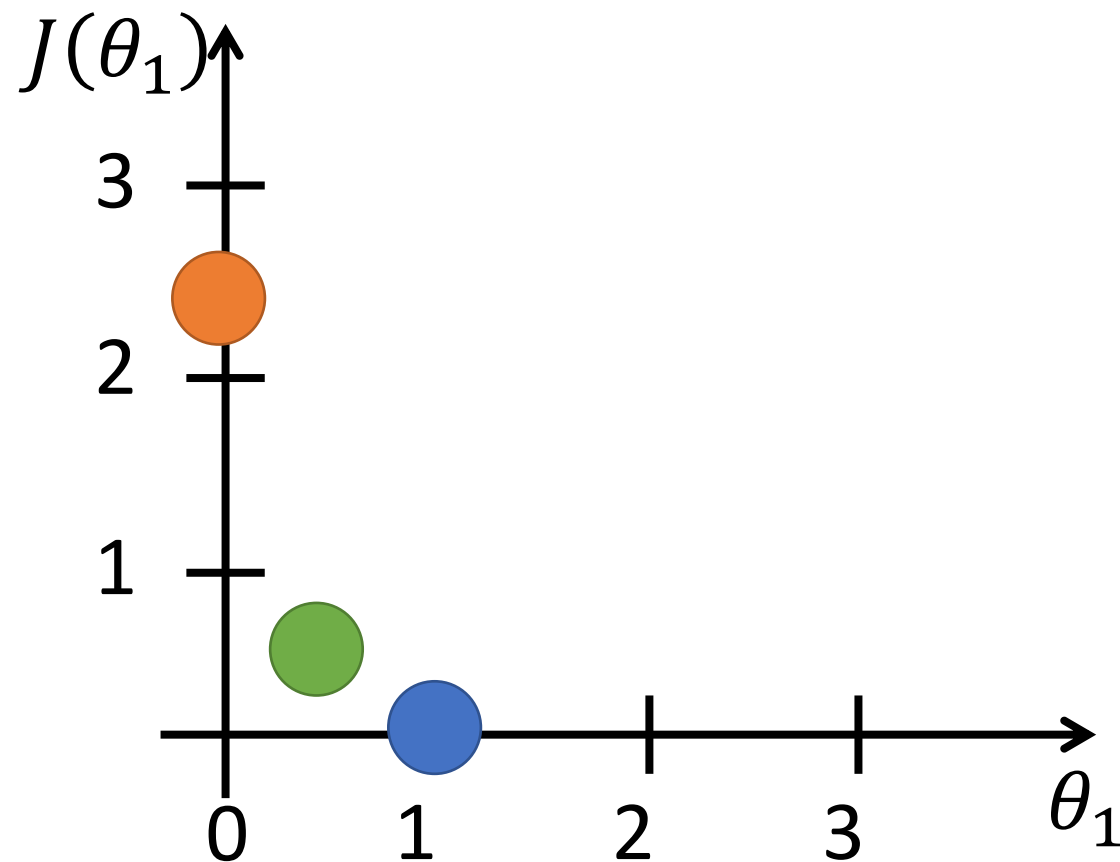
$J(\theta_1)$, function of θ_1



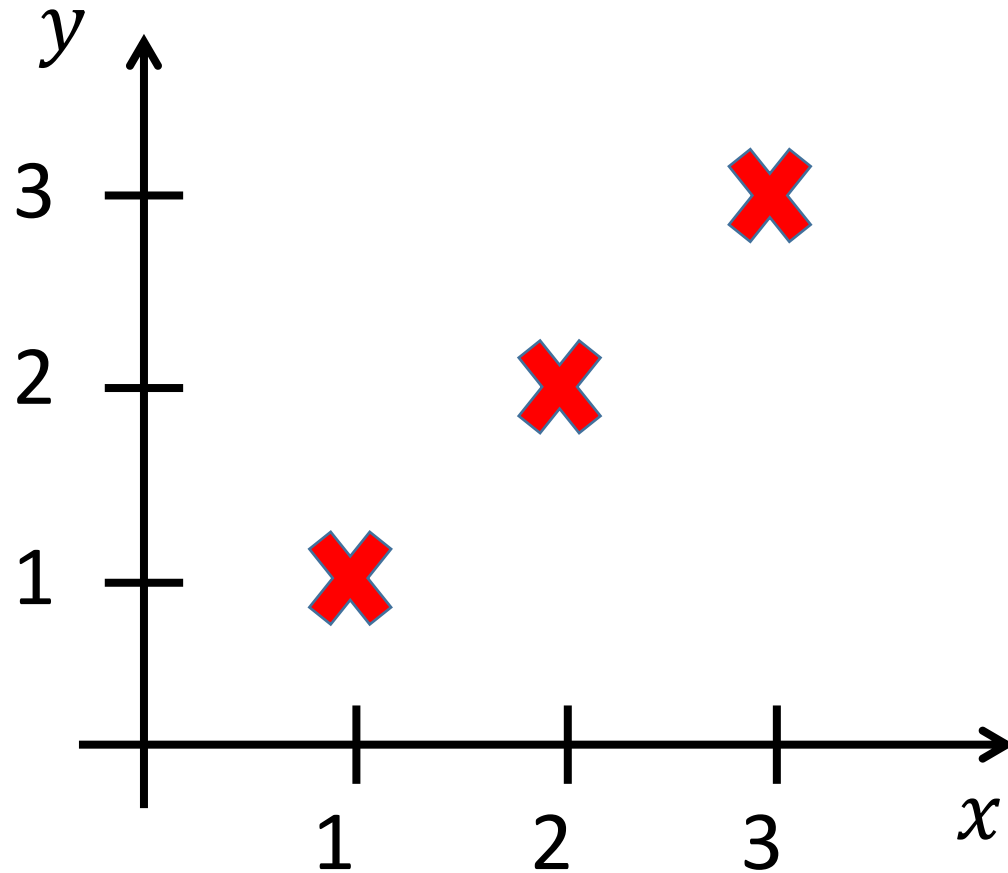
$h_{\theta}(x)$, function of x



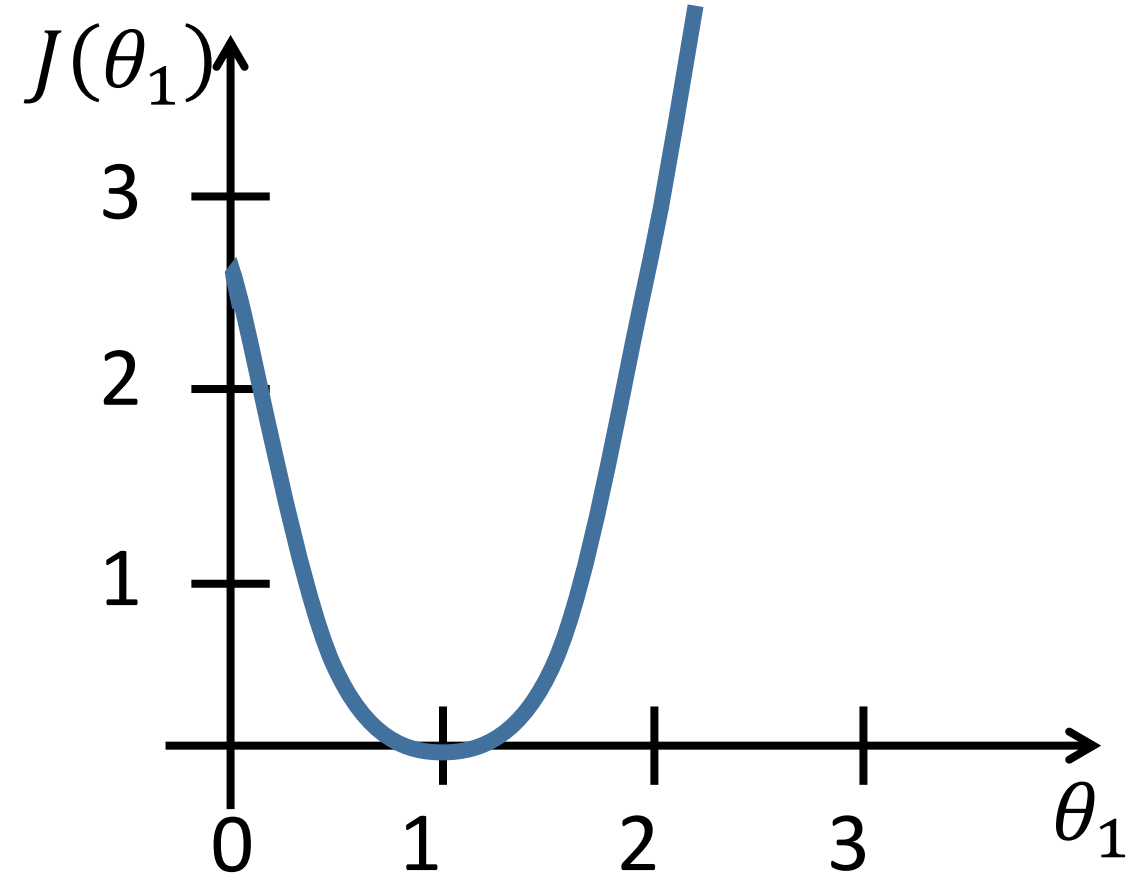
$J(\theta_1)$, function of θ_1



$h_{\theta}(x)$, function of x

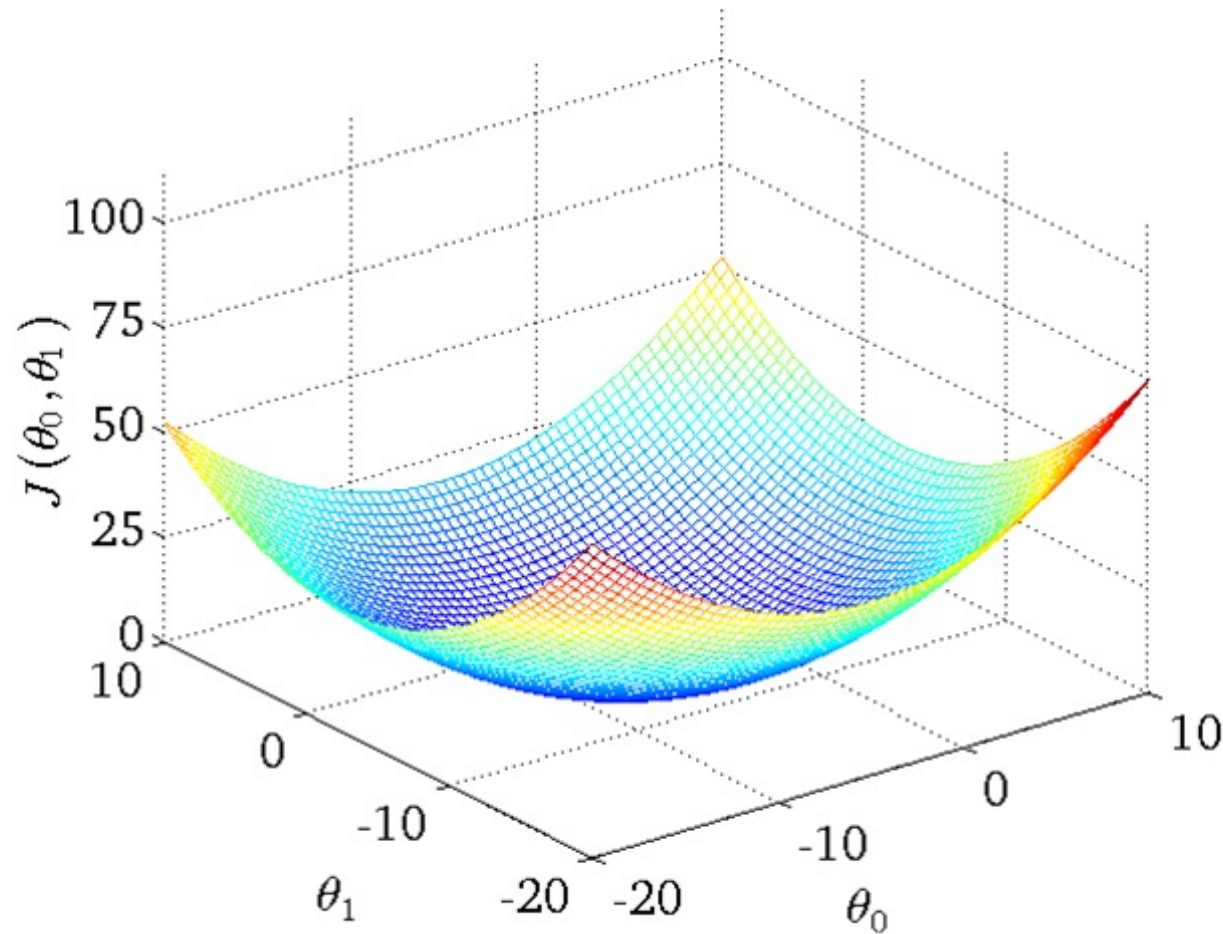


$J(\theta_1)$, function of θ_1



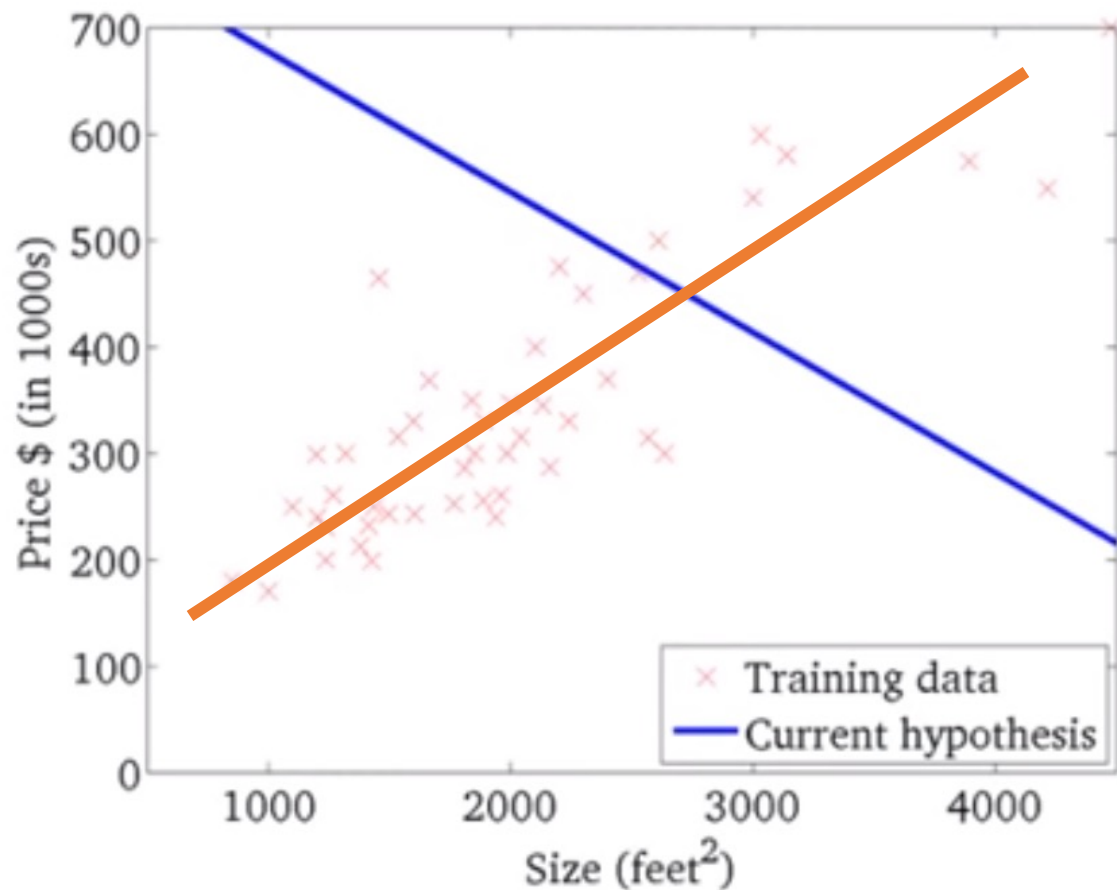
- **Hypothesis:** $h_{\theta}(x) = \theta_0 + \theta_1 x$
- **Parameters:** θ_0, θ_1
- **Cost function:** $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$
- **Goal:** minimize $J(\theta_0, \theta_1)$
 θ_0, θ_1

Cost function



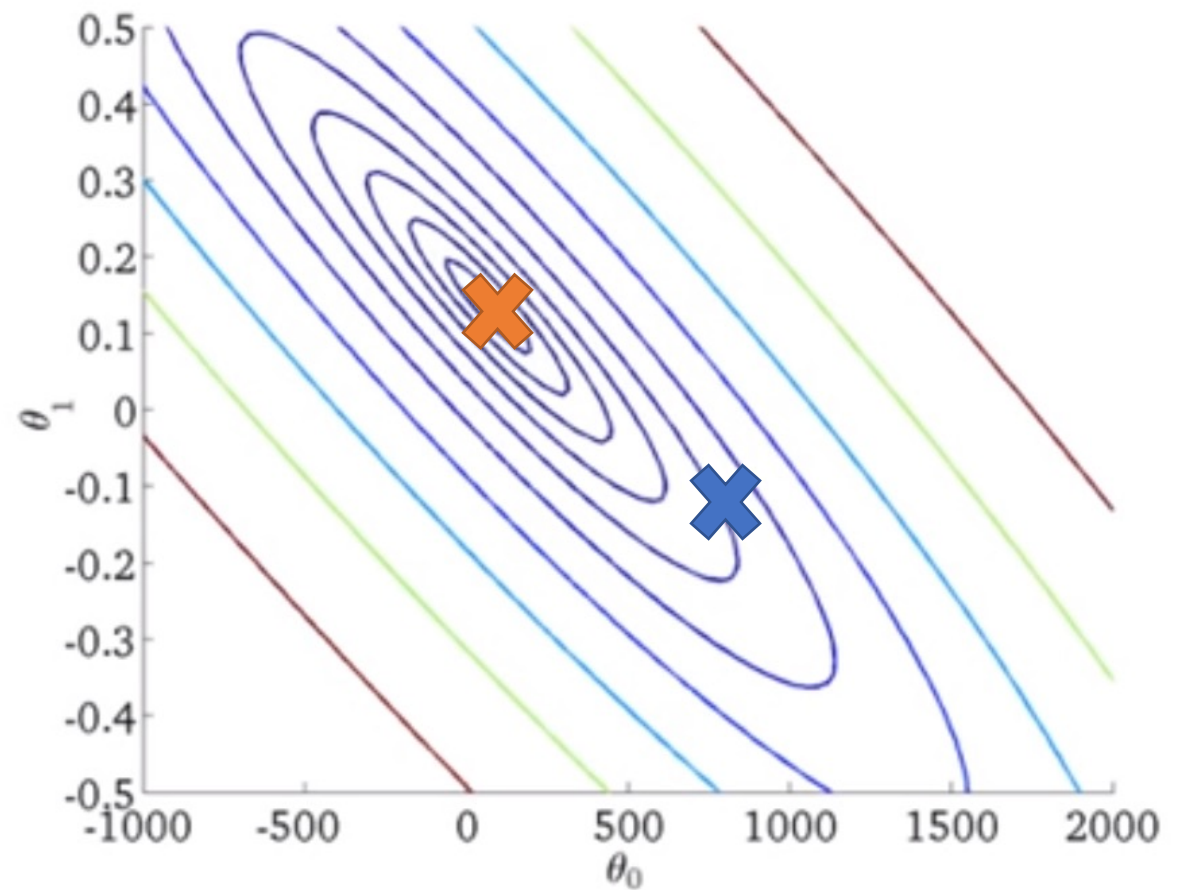
$$h_{\theta}(x)$$

(for fixed θ_0, θ_1 , this is a function of x)



$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



How do we find good θ_0, θ_1 that minimize $J(\theta_0, \theta_1)$?

Linear Regression

- Model representation
- Cost function
- **Gradient descent**
- Features and polynomial regression
- Normal equation

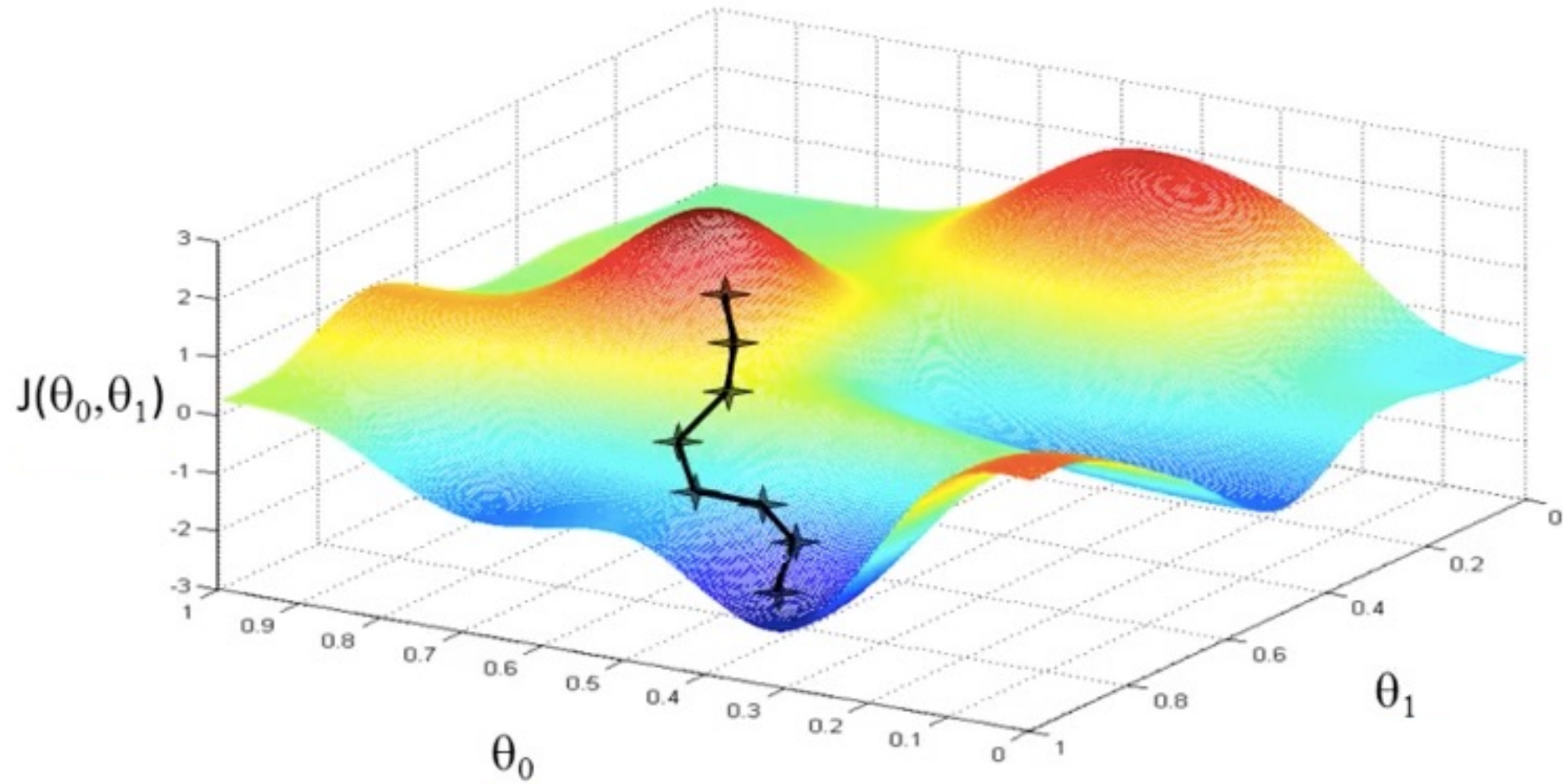
Gradient descent

Have some function $J(\theta_0, \theta_1)$

Want $\operatorname{argmin}_{\theta_0, \theta_1} J(\theta_0, \theta_1)$

Outline:

- Start with some θ_0, θ_1
- Keep changing θ_0, θ_1 to reduce $J(\theta_0, \theta_1)$
until we hopefully end up at minimum



Gradient descent

Repeat until convergence{

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad (\text{for } j = 0 \text{ and } j = 1)$$

}

α : Learning rate (step size)

$\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$: derivative (rate of change)

Gradient descent

Correct: simultaneous update

$$\text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

$$\theta_0 := \text{temp0}$$

$$\theta_1 := \text{temp1}$$

Incorrect:

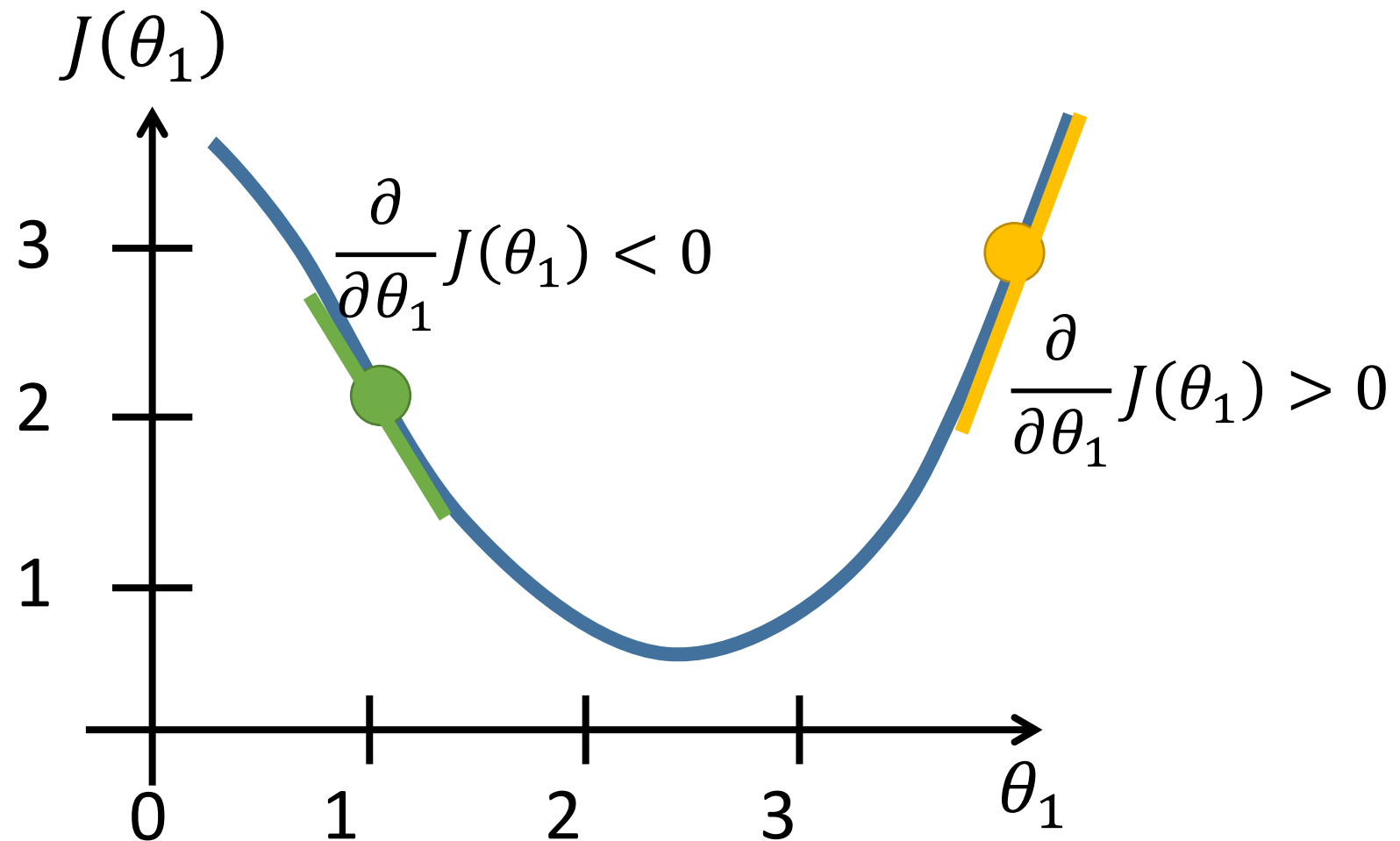
$$\text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

$$\theta_0 := \text{temp0}$$

$$\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$

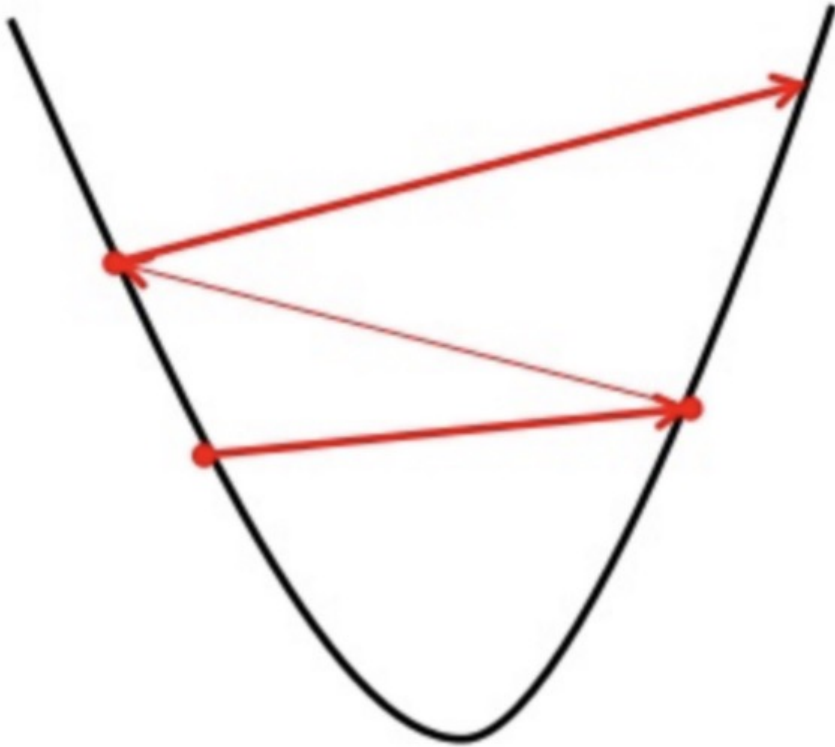
$$\theta_1 := \text{temp1}$$

$$\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_1)$$

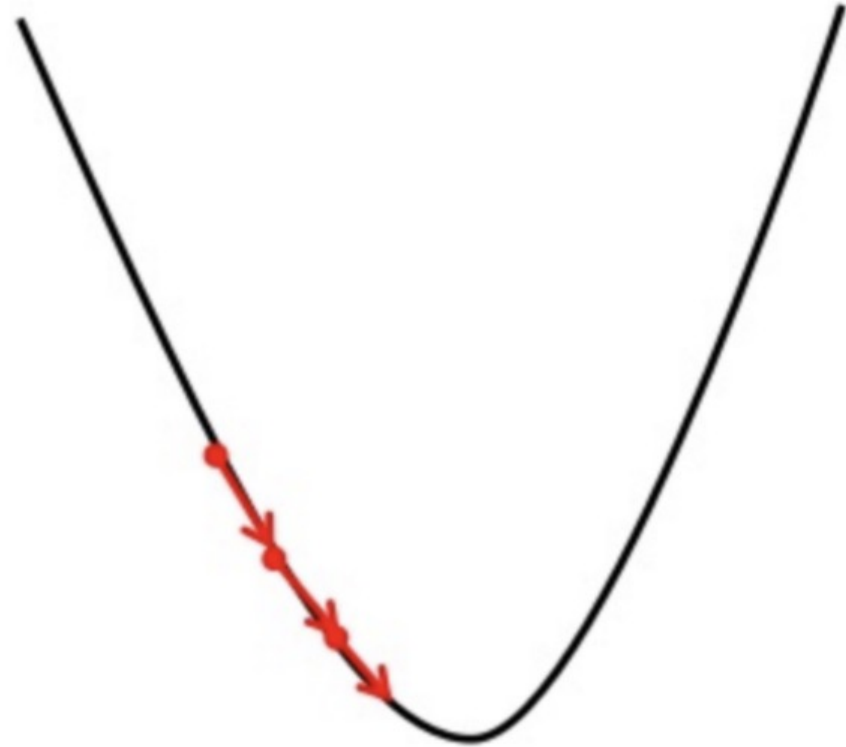


Learning rate

Big learning rate



Small learning rate



Gradient descent for linear regression

Repeat until convergence{

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad (\text{for } j = 0 \text{ and } j = 1)$$

}

- Linear regression model

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

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

Computing partial derivative

- $$\begin{aligned}\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) &= \frac{\partial}{\partial \theta_j} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \\ &= \frac{\partial}{\partial \theta_j} \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2\end{aligned}$$
- $j = 0: \quad \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})$
- $j = 1: \quad \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)}$

Gradient descent for linear regression

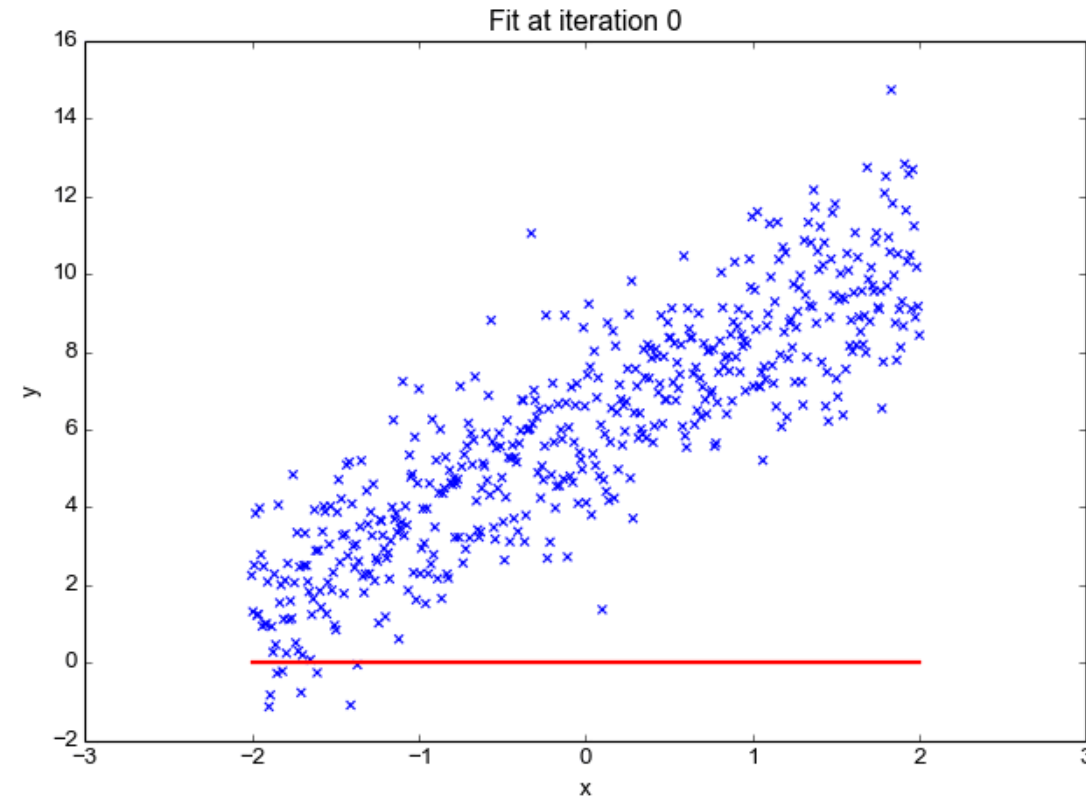
Repeat until convergence{

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)}$$

}

Update θ_0 and θ_1 simultaneously



Batch gradient descent

- “Batch”: Each step of gradient descent uses all the training examples

Repeat until convergence{

m : Number of training examples

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)}$$

}

Linear Regression

- Model representation
- Cost function
- Gradient descent
- **Features and polynomial regression**
- Normal equation

Training dataset

Size in feet ² (x)	Price (\$) in 1000's (y)
2104	460
1416	232
1534	315
852	178
...	...

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

Multiple features (input variables)

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

Notation:

n = Number of features

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

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

$$x_3^{(2)} = ?$$

$$x_3^{(4)} = ?$$

Hypothesis

Previously:

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

Now:

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

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n$$

- For convenience of notation, define $x_0 = 1$
($x_0^{(i)} = 1$ for all examples)

$$\bullet \mathbf{x} = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in R^{n+1} \qquad \boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \in R^{n+1}$$

- $$\begin{aligned} h_{\theta}(x) &= \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n \\ &= \boldsymbol{\theta}^{\top} \mathbf{x} \end{aligned}$$

Gradient descent

- Previously ($n = 1$)

Repeat until convergence{

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)}$$

}

- New algorithm ($n \geq 1$)

Repeat until convergence{

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

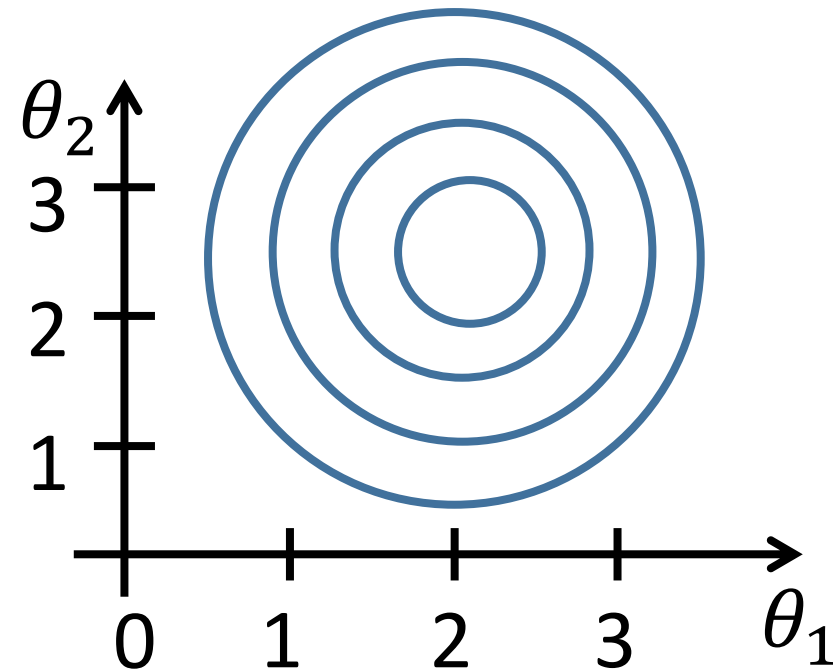
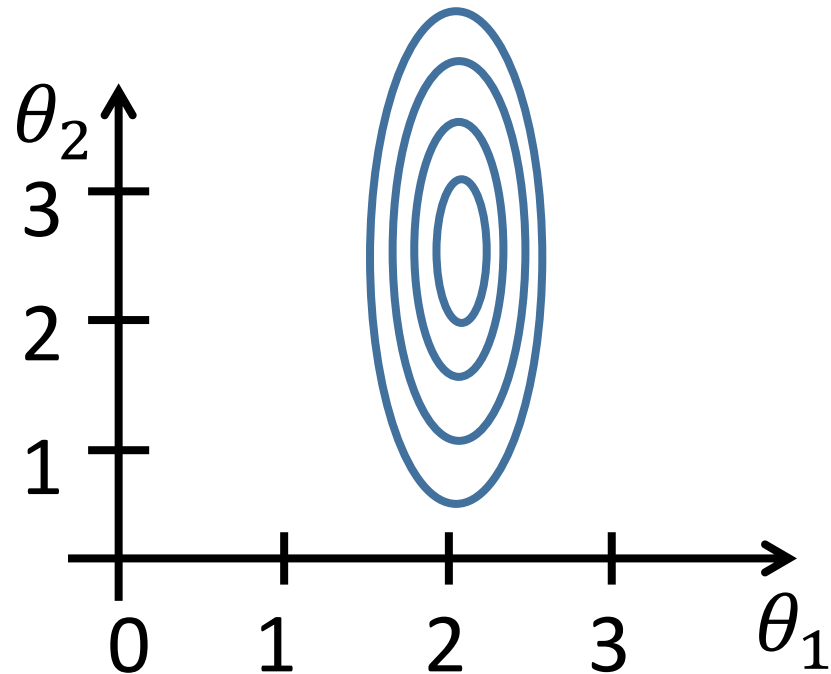
}

Simultaneously update

θ_j , for $j = 0, 1, \dots, n$

Gradient descent in practice: Feature scaling

- Idea: Make sure features are on a similar scale (e.g., $-1 \leq x_i \leq 1$)
- E.g. x_1 = size (0-2000 feat²)
 x_2 = number of bedrooms (1-5)

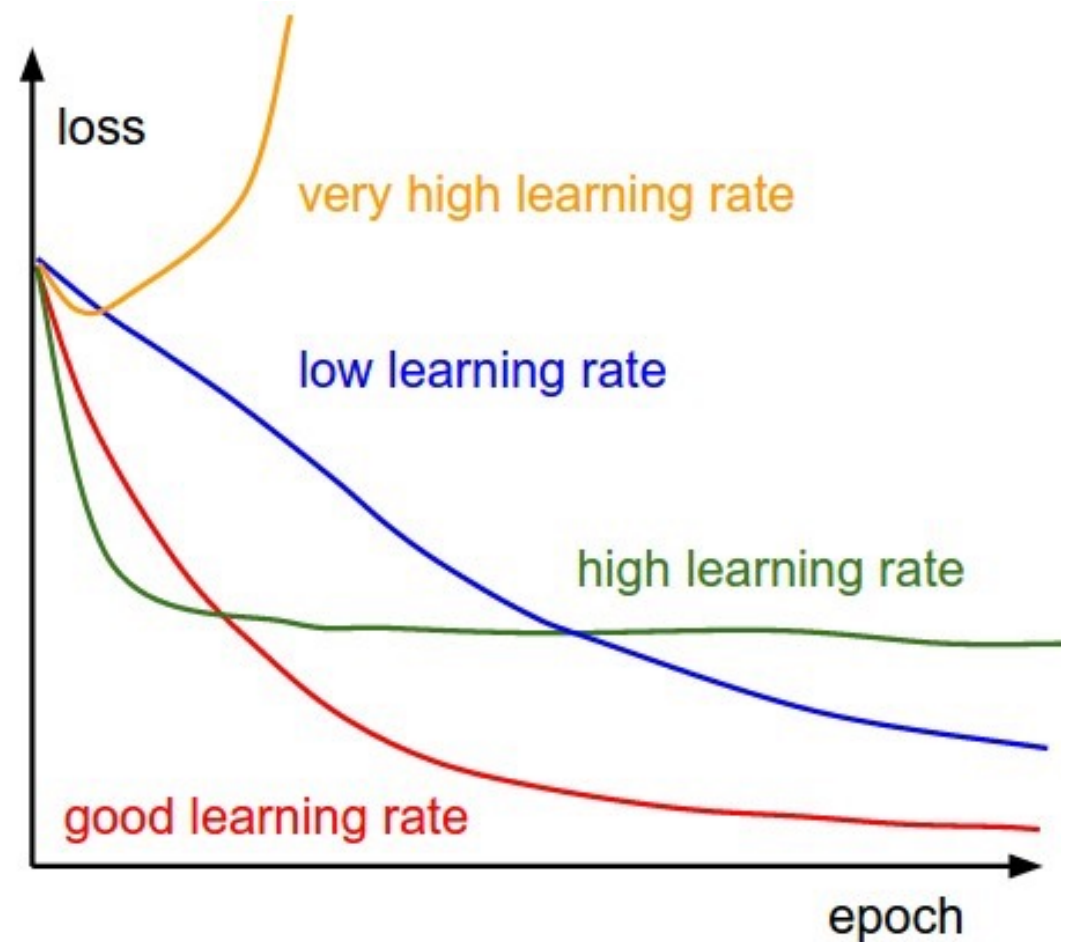


Gradient descent in practice: Learning rate

- Automatic convergence test
- α too small: slow convergence
- α too large: may not converge

- To choose α , try

0.001, ... 0.01, ..., 0.1, ... , 1



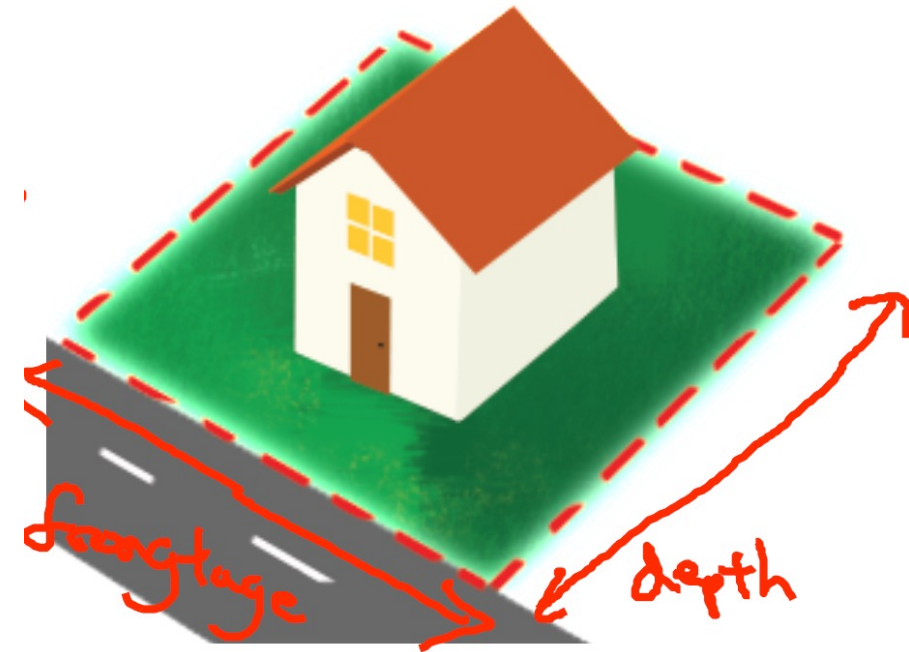
House prices prediction

- $h_{\theta}(x) = \theta_0 + \theta_1 \times \text{frontage} + \theta_2 \times \text{depth}$

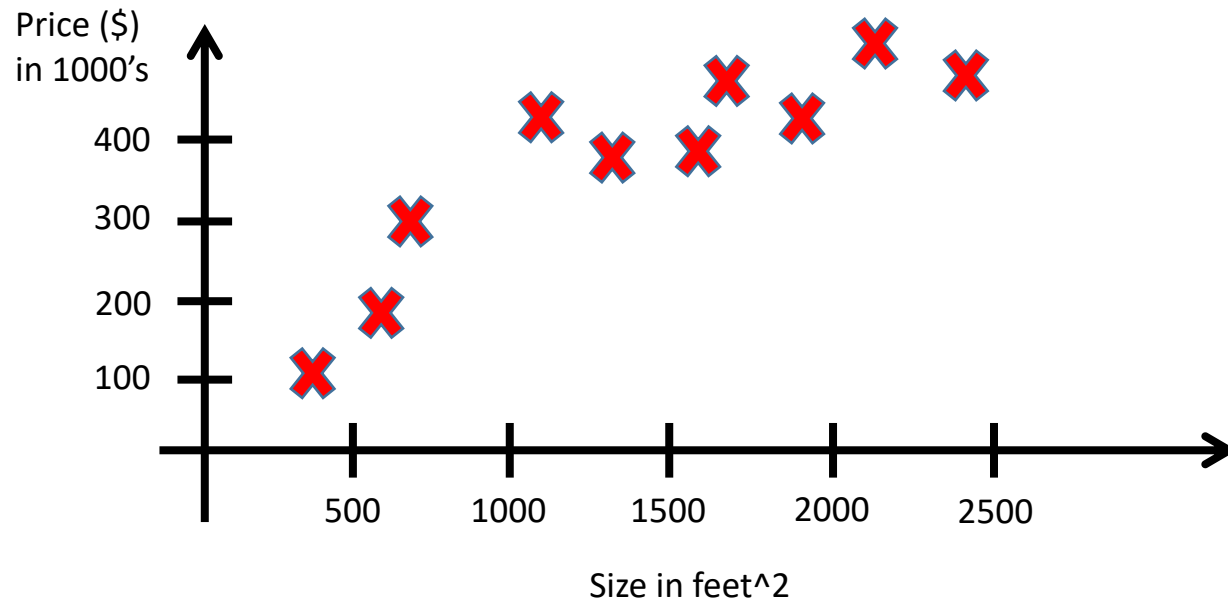
- Area

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

- $h_{\theta}(x) = \theta_0 + \theta_1 x$



Polynomial regression



$$\begin{aligned}x_1 &= (\text{size}) \\x_2 &= (\text{size})^2 \\x_3 &= (\text{size})^3\end{aligned}$$

- $$\begin{aligned}h_{\theta}(x) &= \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 \\&= \theta_0 + \theta_1 (\text{size}) + \theta_2 (\text{size})^2 + \theta_3 (\text{size})^3\end{aligned}$$

Linear Regression

- Model representation
- Cost function
- Gradient descent
- Features and polynomial regression
- **Normal equation**

(x_0)	Size in feet ² (x_1)	Number of bedrooms (x_2)	Number of floors (x_3)	Age of home (years) (x_4)	Price (\$) in 1000's (y)
1	2104	5	1	45	460
1	1416	3	2	40	232
1	1534	3	2	30	315
1	852	2	1	36	178

$$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix}$$

$$y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix}$$

$$\theta = (X^T X)^{-1} X^T y$$

Least square solution

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

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

$$= \frac{1}{2m} \|X\theta - y\|_2^2$$

- $\frac{\partial}{\partial \theta} J(\theta) = 0$

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

Justification/interpretation 1

- **Loss minimization**

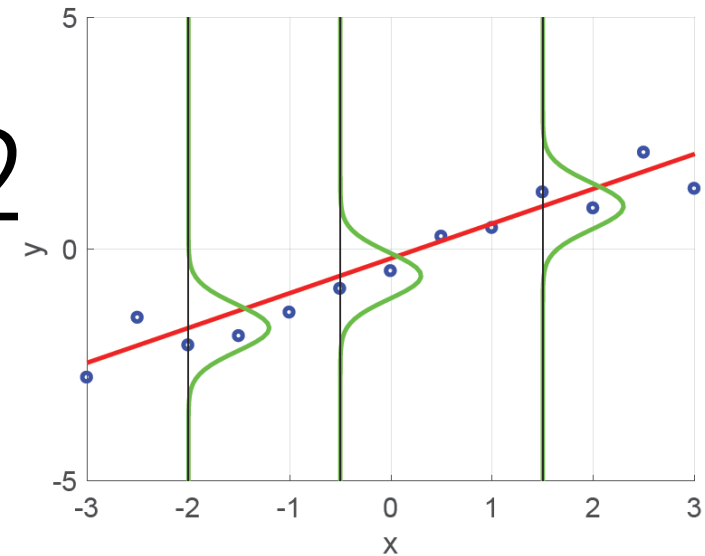
$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 = \frac{1}{m} \sum_{i=1}^m L_{ls}(h_{\theta}(x^{(i)}), y^{(i)})$$

- $L_{ls}(y, \hat{y}) = \frac{1}{2} \|y - \hat{y}\|_2^2$: Least squares loss

- Empirical Risk Minimization (ERM)

$$\frac{1}{m} \sum_{i=1}^m L_{ls}(y^{(i)}, \hat{y})$$

Justification/interpretation 2



- **Probabilistic model**

- Assume linear model with Gaussian errors

$$p_{\theta}(y^{(i)}|x^{(i)}) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} (y^{(i)} - \theta^{\top}x^{(i)})^2\right)$$

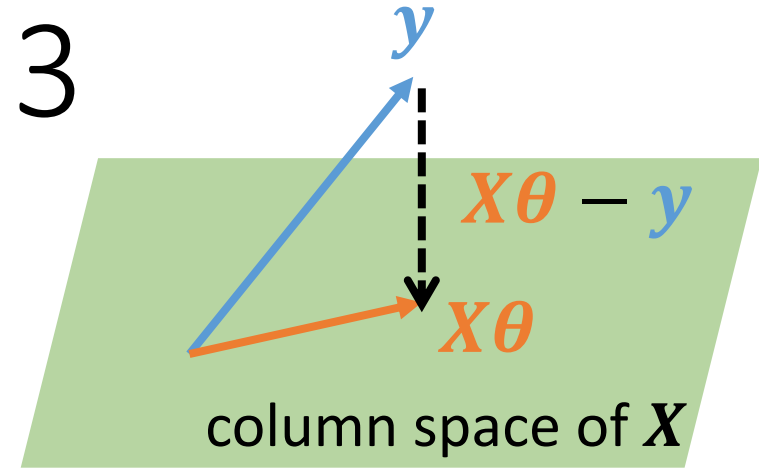
- Solving maximum likelihood

$$\operatorname{argmin}_{\theta} \prod_{i=1}^m p_{\theta}(y^{(i)}|x^{(i)})$$

$$\operatorname{argmin}_{\theta} \log\left(\prod_{i=1}^m p(y^{(i)}|x^{(i)})\right) = \operatorname{argmin}_{\theta} \frac{1}{2\sigma^2} \sum_{i=1}^m \frac{1}{2} (\theta^{\top}x^{(i)} - y^{(i)})^2$$

Justification/interpretation 3

- Geometric interpretation



$$X = \begin{bmatrix} 1 & \leftarrow \mathbf{x}^{(1)} \rightarrow \\ 1 & \leftarrow \mathbf{x}^{(2)} \rightarrow \\ \vdots & \vdots \\ 1 & \leftarrow \mathbf{x}^{(m)} \rightarrow \end{bmatrix} = \begin{bmatrix} \uparrow & \uparrow & \uparrow & \cdots & \uparrow \\ \mathbf{z}_0 & \mathbf{z}_1 & \mathbf{z}_2 & \cdots & \mathbf{z}_n \\ \downarrow & \downarrow & \downarrow & \cdots & \downarrow \end{bmatrix}$$

- $X\theta$: column space of X or $\text{span}(\{\mathbf{z}_0, \mathbf{z}_1, \dots, \mathbf{z}_n\})$
- Residual $X\theta - \mathbf{y}$ is orthogonal to the column space of X
- $X^\top (X\theta - \mathbf{y}) = 0 \rightarrow (X^\top X)\theta = X^\top \mathbf{y}$

m training examples, n features

Gradient Descent

- Need to choose α
- Need many iterations
- Works well even when n is large

Normal Equation

- No need to choose α
- Don't need to iterate
- Need to compute $(X^T X)^{-1}$
- Slow if n is very large

Things to remember

- **Model representation**

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n = \theta^{\top} x$$

- **Cost function**

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

- **Gradient descent for linear regression**

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

- **Features and polynomial regression**

Can combine features; can use different functions to generate features (e.g., polynomial)

- **Normal equation** $\theta = (X^{\top} X)^{-1} X^{\top} y$