

APS1070

Foundations of Data Analytics and
Machine Learning

Winter 2022

Week 9:

- *Monte Carlo Simulation*
- *Empirical Risk Minimization*
- *Maximum Likelihood Estimation*
- *Linear Regression*



Slide Attribution

These slides contain materials from various sources. Special thanks to the following authors:

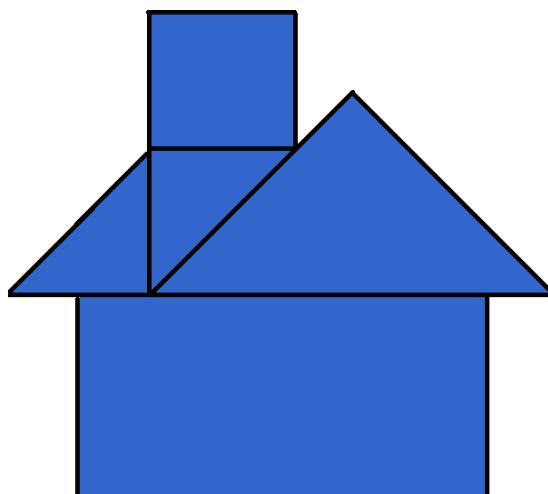
- Roger Grosse
- William Fleshman
- Lisa Zhang
- Andrew Ng
- Jason Riordon

Monte Carlo Simulation

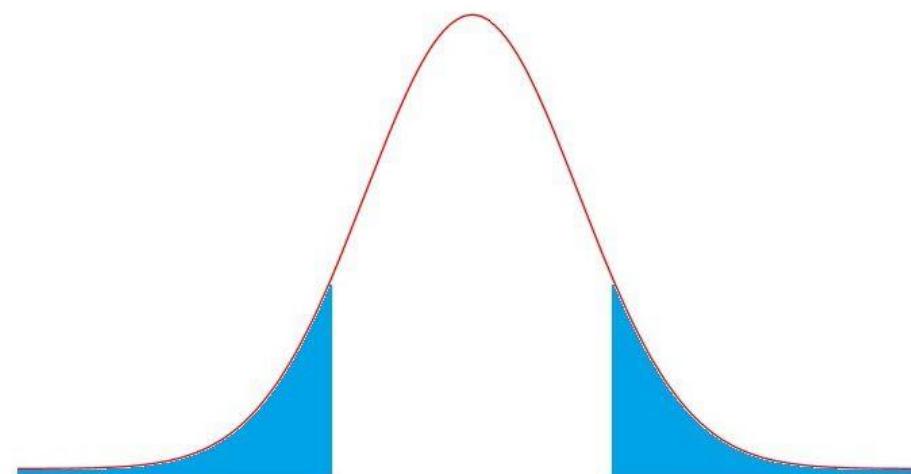
Why Sampling Methods?

➤ Q: How would you solve the following problems?

1. Compute area



2. Compute the blue area



3. Integrate $f(x)$

$$f(x) = x^2$$

$$\int f(x)dx$$

➤ Most of the examples we've seen so far can be solved analytically.

Why Sampling Methods?

- Q: How would you solve the following problems?

Compute Area



Integrate a function

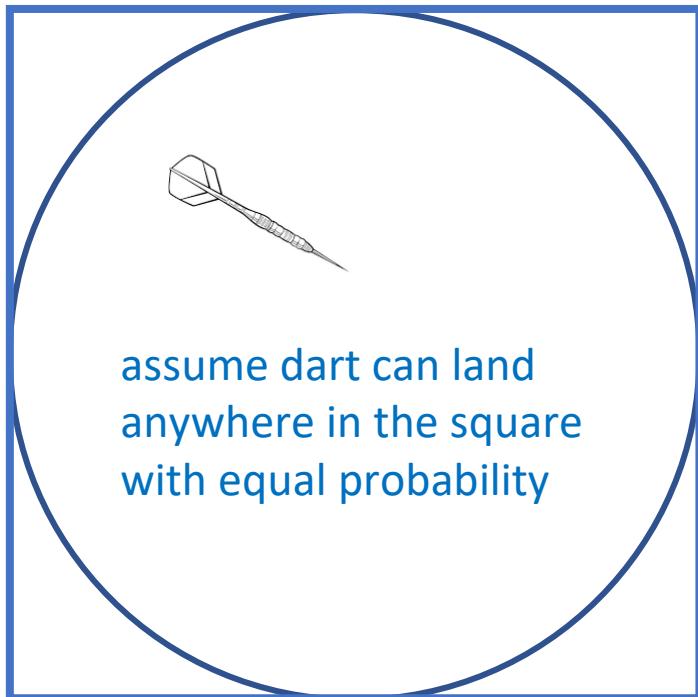
$$\int_0^4 \sqrt[4]{15x^3 + 21x^2 + 41x + 3} \cdot e^{-0.5x} dx$$



- For most **problems of practical interest**, the **analytical solution is intractable**, and so we have to resort to some form of approximation.

Simple Illustration

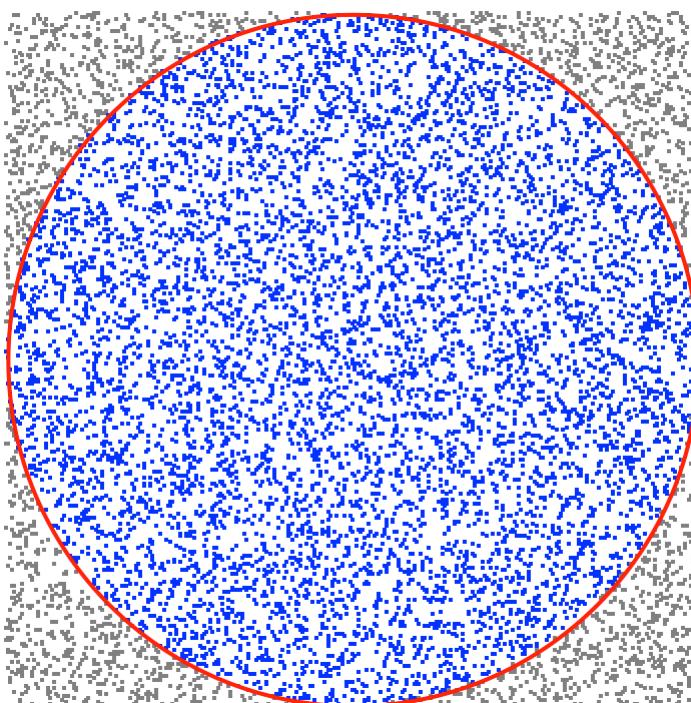
- Q: What is the probability of a dart landing in the circle?



- Simulate many random 2D points
- Use $x^2 + y^2 \leq r^2$
- Then apply $\frac{N_{in}}{N_{tot}}$ to estimate area

Simple Illustration: Simulation

- Q: What is the probability of a dart landing in the circle?



----- Simulation Code -----

```
def throwDarts(num_darts):
    in_circle = 0
    for darts in range(num_darts):
        x = random.random()
        y = random.random()
        if (x*x + y*y)**0.5 <= 1.0:
            in_circle +=1
    return (in_circle/num_darts)
```

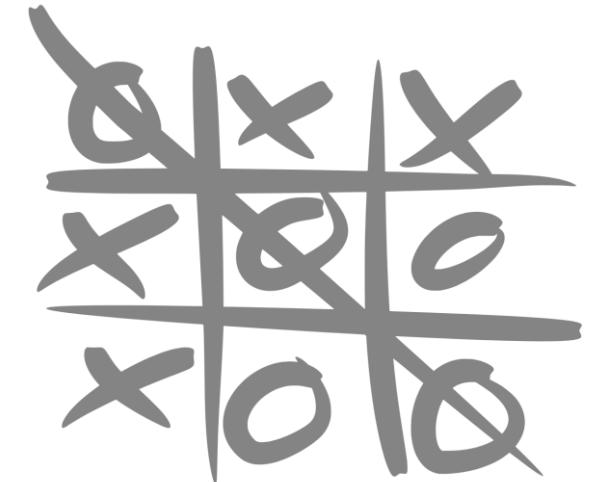
Estimate: 0.78308

Truth: 0.7853981633974483

How can we
obtain π ?

Harder Example

- Q: Assume you are playing tic-tac-toe and you have a $P(\text{win}) = p$. The game ends when you lose 2 times in a row. What is the expected number of rounds we will play the game before losing?
- Coming up with an analytical solution can take a long time. In the process we may make a mistake without realizing it.
- Is there a way we can obtain an estimate?



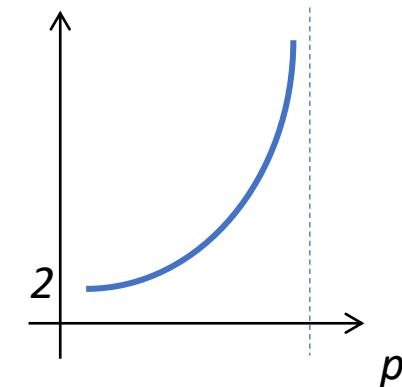
Harder Example: Simulation

----- Simulation Code -----

```
import numpy as np
def mc(n=10**6,p=0.5):
    rounds= []
    for x in range(n):
        r,losses = 0,0
        while losses != 2:
            r+=1
            if np.random.random() <= p:
                losses = 0
            else:
                losses +=1
        rounds.append(r)
    return np.mean(rounds)
```

Analytical Solution

$$\text{rounds} = \frac{2-p}{(1-p)^2}$$



Estimate: 5.996614
Truth: 6

Source: ritvikmath

Monte Carlo Sampling Methods

- It was invented by Stanislaw Ulam during the Manhattan Project.
 - e.g., assess the risk of a runaway chain reaction that could blow up the earth.
- Used extensively in:
 - risk analysis,
 - finance,
 - supply chain logistics,
 - computational physics,
 - robotic



Monte Carlo is not a single method but instead any type of method that relies on **simulations and randomness** to get the solution to a problem.

Law of Large Numbers

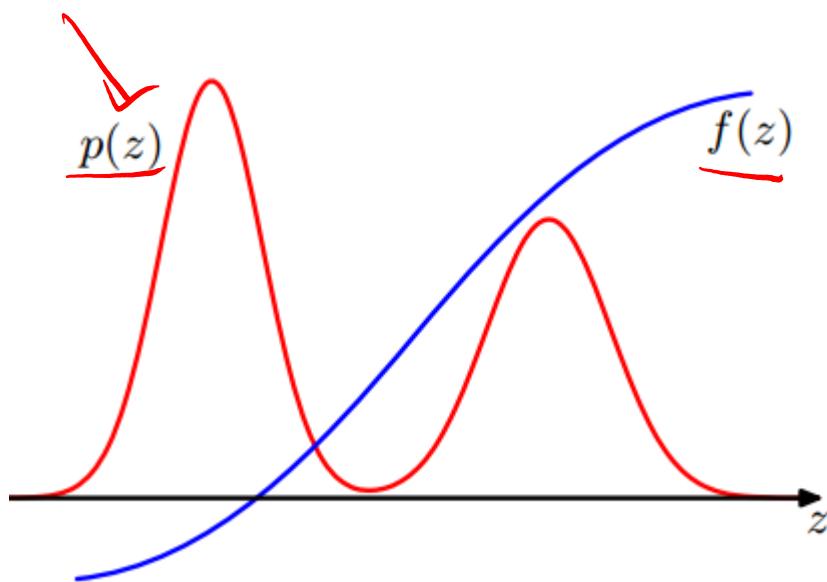
- The Law of Large Numbers describes what happens when performing the same experiment many times.
- After many trials, the average of the results should be close to the expected value and will be more accurate with more trials.
- For Monte Carlo simulation, this means that we can learn properties of a random variable (mean, variance, etc.) simply by simulating it over many trials.

Pros and Cons

- Pros:
 - Monte Carlo simulation is often easy to implement.
 - Doesn't matter if you are subject matter expert. It helps people come to similar results without analytical efforts.
- Cons:
 - It isn't always fast, will depends on the problem (e.g., p=0.999 in our game example will take a long time to finish).
 - Not generalizable, just because you have one result for one set of parameters, it doesn't say anything about other parameters unlike analytical solution.
 - Not interpretable.

Why Sampling Methods?

- For most probabilistic models of practical interest, exact inference is intractable, and so we have to resort to some form of approximation.
- Here we will focus primarily on the problem of estimating expectations for intractable distributions:



$$\mathbb{E}[f] = \int f(\mathbf{z}) p(\mathbf{z}) d\mathbf{z}$$

Comes up frequently in Bayesian modelling

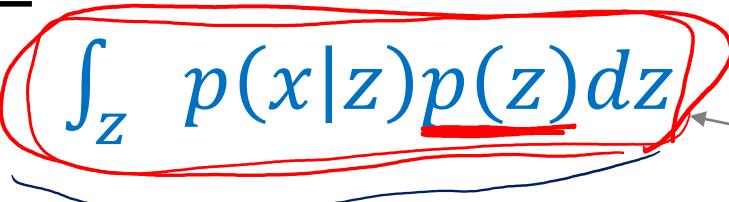
Partition Function

- A common problem in machine learning is that the **partition function** is often intractable (or difficult to compute):

$$p(z|x) = \frac{p(x|z)p(z)}{\int_z p(x|z)p(z)dz}$$

partition function

$p(x)$

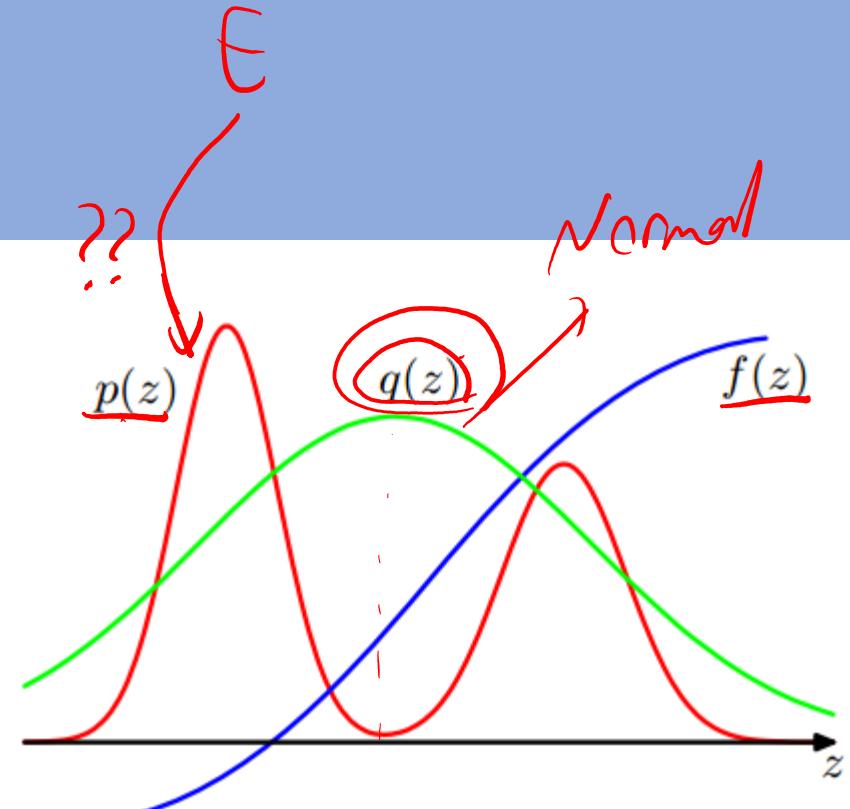


In circumstances where we are unable to sample directly from the distribution $p(z)$ we can use an alternative approach known as importance sampling.

Importance Sampling

- Draws samples from a simpler distribution $q(z)$ to estimate expectation of $f(x)$
- Corresponding terms in the summation are weighted by the ratios $p(x)/q(x)$

$$\begin{aligned}\mathbb{E}_p[f(x)] &= \sum_x p(x)f(x) \\ &= \sum_x q(x) \frac{p(x)f(x)}{q(x)} \\ &= \mathbb{E}_q\left[\frac{p(x)}{q(x)}f(x)\right],\end{aligned}$$



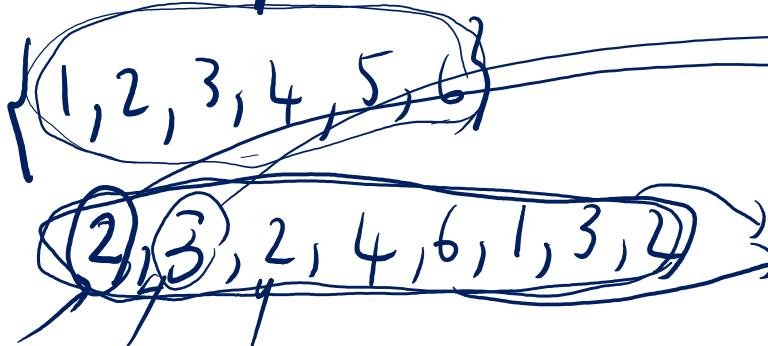
Source: Bishop PRML

Importance Sampling

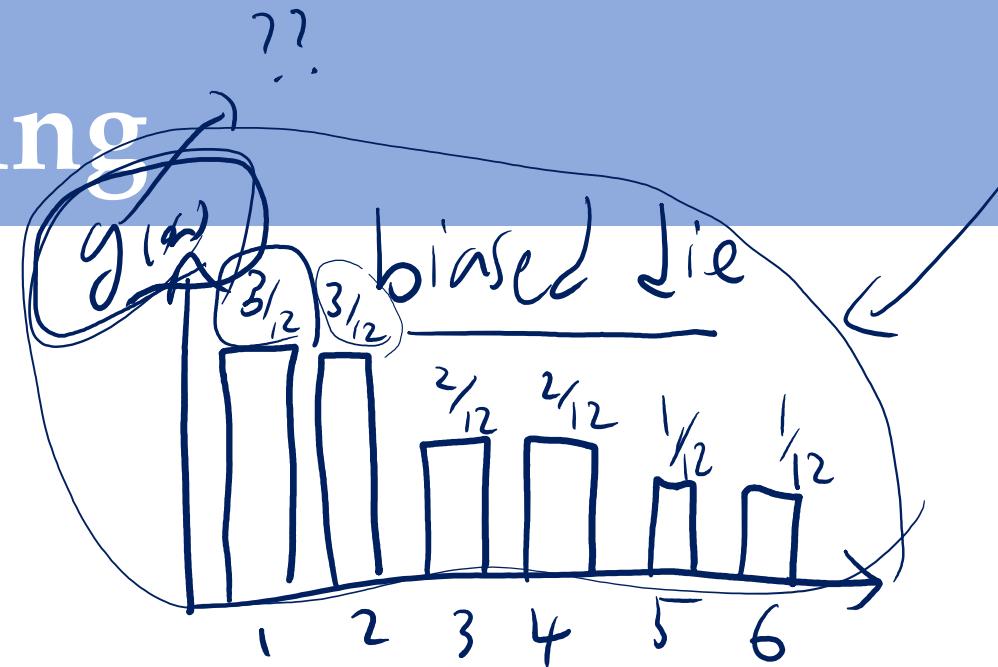


$$E_f[x] = \sum_n x f(n)$$

$$\approx \frac{1}{n} \sum_{i=1}^{1000} x_i \approx 3.5$$



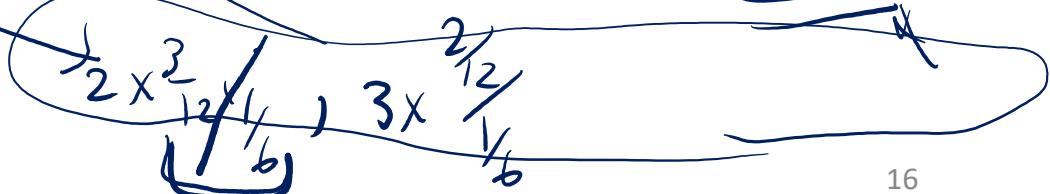
$$\frac{23}{8} = 2.85$$



$$E_g[x] = \sum_n n g(n) = \sum_n x \frac{g(n)}{f(n)} f(n)$$

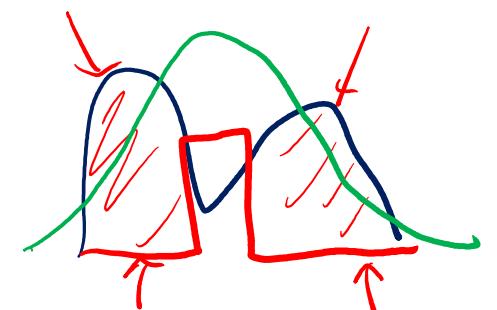
$$= E_f \left[x \frac{g(n)}{f(n)} \right] = \frac{1}{n} \sum_{i=1}^{1000} x_i \frac{g(n)}{f(n)}$$

$$\approx 2.83$$



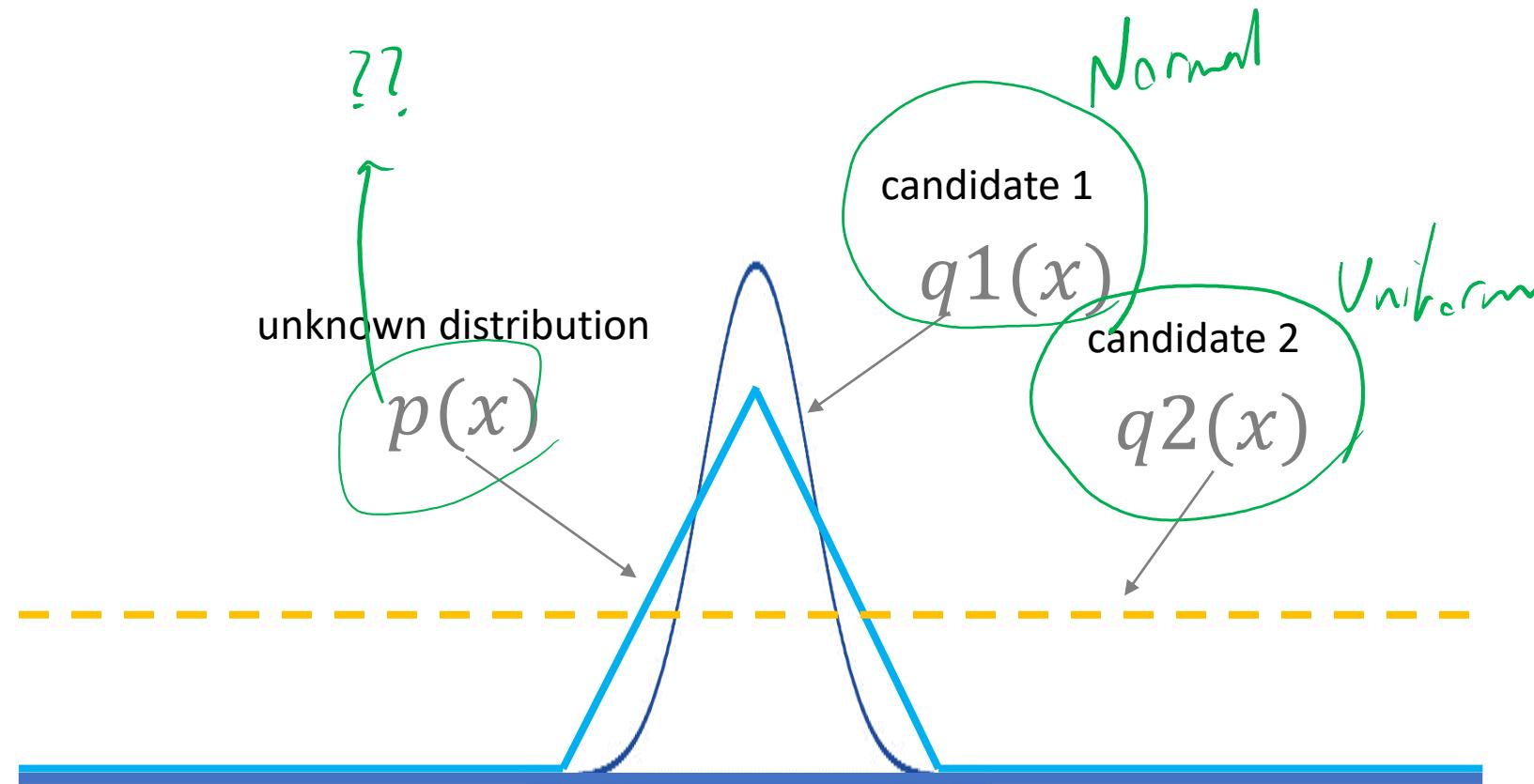
Limitations Con't

- The accuracy depends on how well the sampling distribution $q(z)$ matches the desired distribution $p(z)$.
- **Importance weights may be dominated by a few examples** having large values (remaining weights being small)
- A drawback of the importance sampling method is the potential to produce results that have high error and no diagnostic indication.
- $q(z)$ should not be small or zero where $p(z)$ is substantial.



Source: Bishop PRML

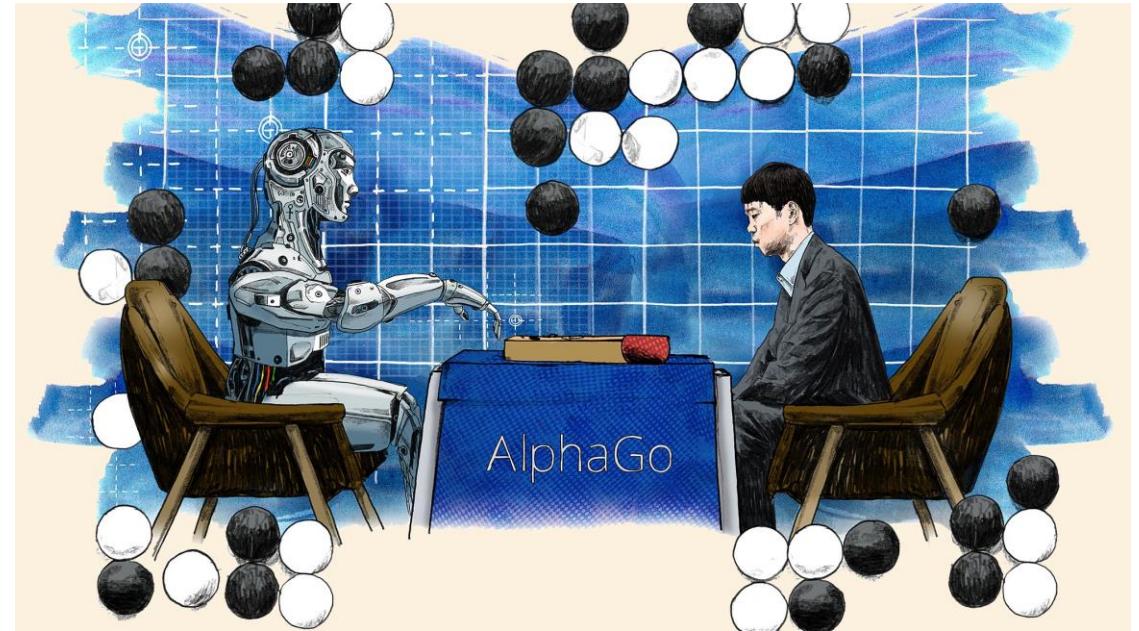
Which candidate distribution is best?



Source: Ben Lambert

Link to Deep Learning

- Leads into more advanced methods such as Markov Chain Monte Carlo (MCMC).
- A whole course could be devoted to this topic.
- DeepMind used a version of Monte Carlo Sampling in AlphaGo.



Source: [Level Up Coding](#)

Last Time

- Matrix decompositions, dimensionality reduction and interpretations.

- SVD
- PCA
- Applications
- Vector Calculus

$$\begin{matrix} n \\ m \end{matrix} X = \begin{matrix} m \\ m \end{matrix} U \begin{matrix} n \\ n \end{matrix} \Sigma \begin{matrix} n \\ m \end{matrix} V^\top z$$

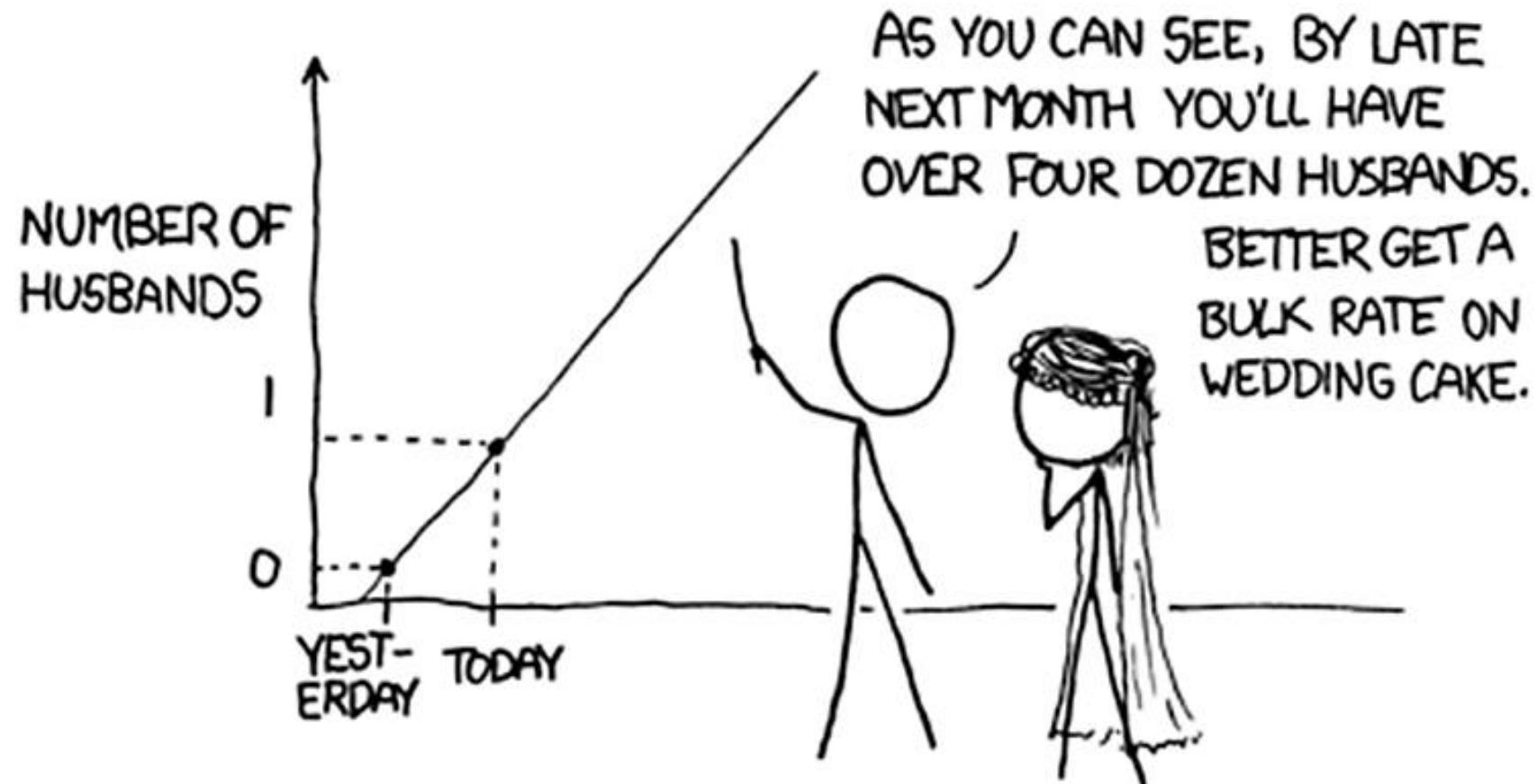
~~10:55~~
9:55 :)

- Today we also discuss learning algorithms, this time focusing on **model-based learning algorithms** starting with linear regression.

Decision Making

- We often need to make some decisions given some data.
- Two common types of decisions that we make are:
 - Classification
 - Discrete outcome
 - Regression
 - Continuous outcome

Linear Regression



Example: House Price Prediction



Recap: Learning Algorithms/Models

- We've seen earlier that there are several approaches to solving regression problems.
- **Instance-Based**
 - The first couple lectures focused on algorithms require long-term storage of data in memory in order to make predictions/decisions.
- **Model-Based**
 - Now we will introduce learning algorithms that replace samples with model parameters for making predictions.

Agenda

- Data
- Empirical Risk Minimization
- Gradient Descent
- Maximum Likelihood Estimation
- Negative Log-Likelihood
- Application of Linear Regression



Theme:
Linear Regression

Linear Regression (Empirical Risk Minimization)



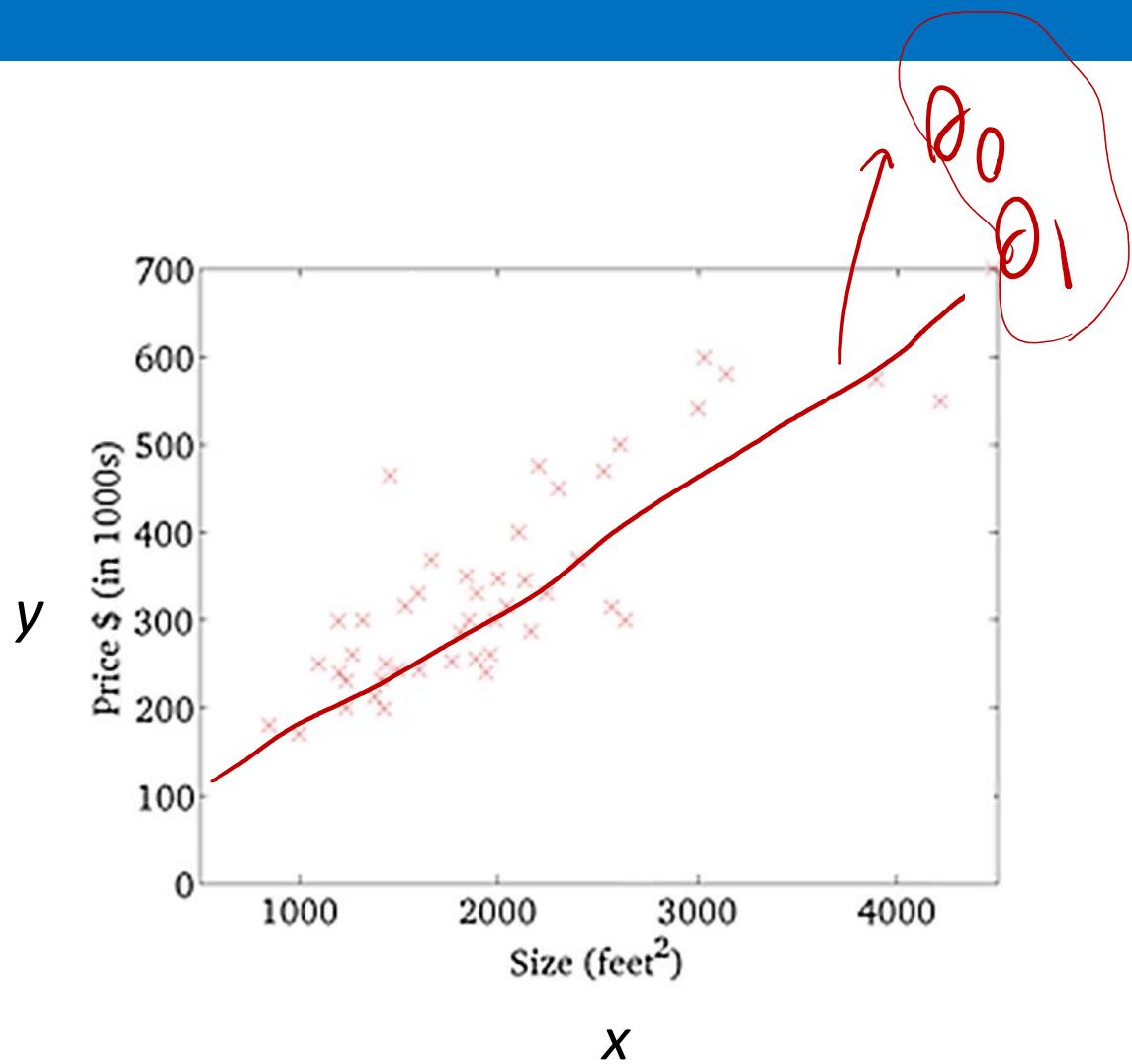
Readings:

- Chapter 7 MML Textbook
- Chapter 8.1-2 MML Textbook

Problem Setup

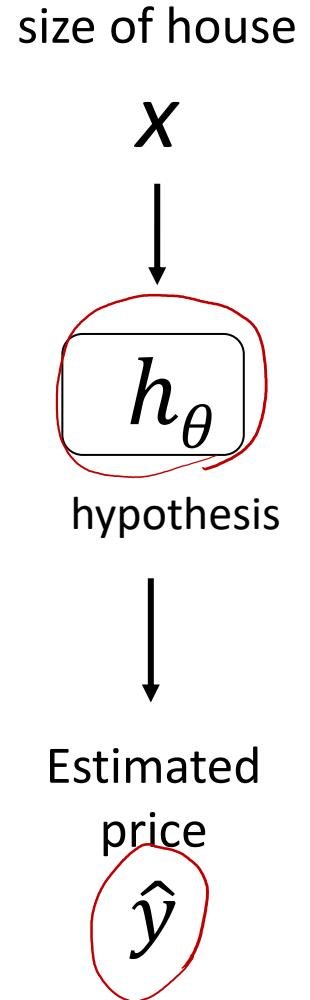
Size in feet ² (x)	Price in 1000's (y)
320	148
450	210
845	362
1043	440
1160	550
...	...

$$\{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$$



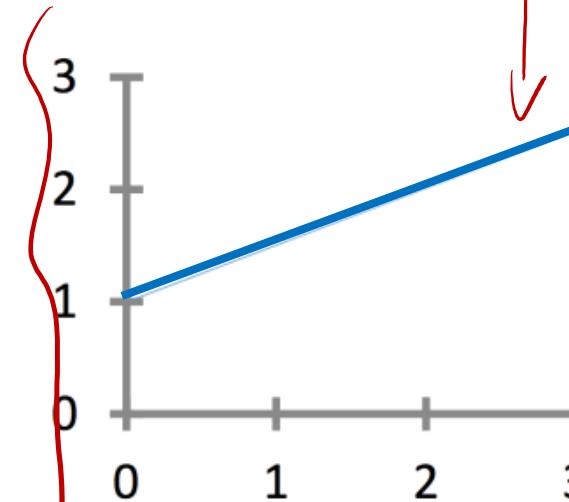
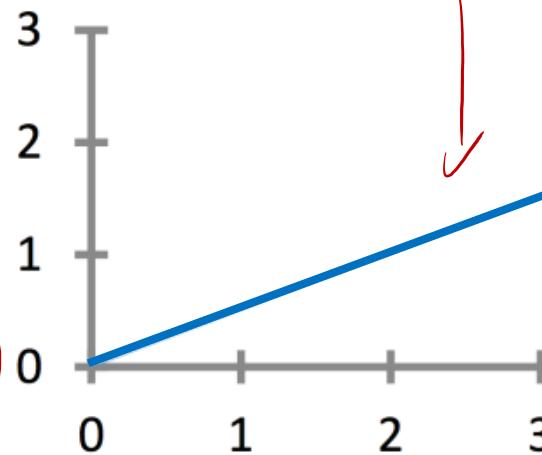
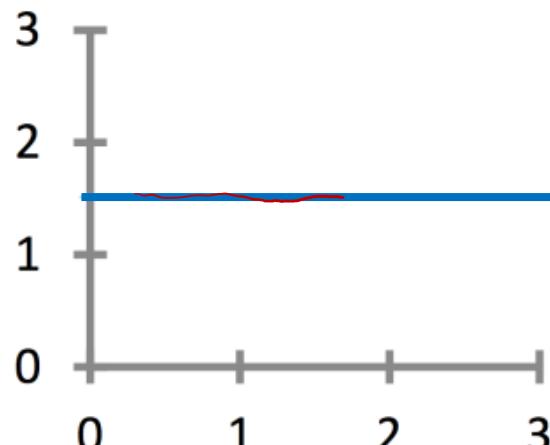
Linear Model

- Parameters $\theta = (\theta_0, \theta_1)$ can be used to test a hypothesis about the data
- If the data is correctly predicted according to the hypothesis h_θ , then $y \sim h_\theta(x) = \theta_0 + \theta_1 x$
- We can then estimate y for new values of x using our h_θ
- If $h_\theta(x)$ is a linear function of a real number x , this procedure is called **linear regression**.



Many Hypotheses to Choose From

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



①

Intercept
slope

$\theta_0 = 1.5$
 $\theta_1 = 0$

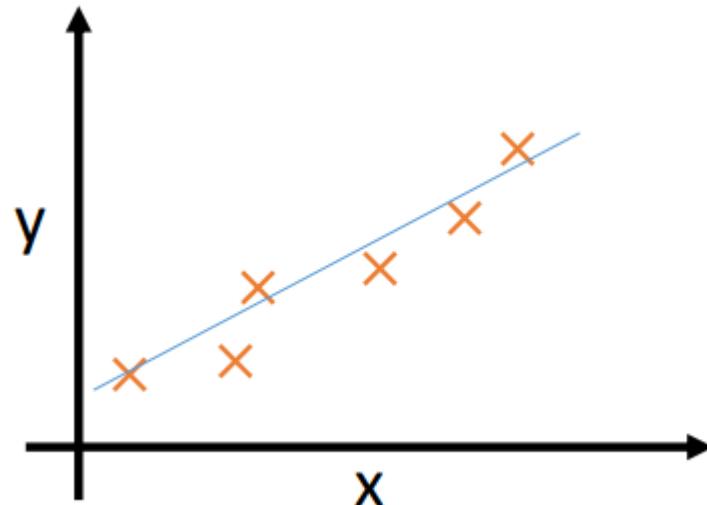
②

$\theta_0 = 0$
 $\theta_1 = 0.5$

③

$\theta_0 = 1$
 $\theta_1 = 0.5$

What is the best Model?



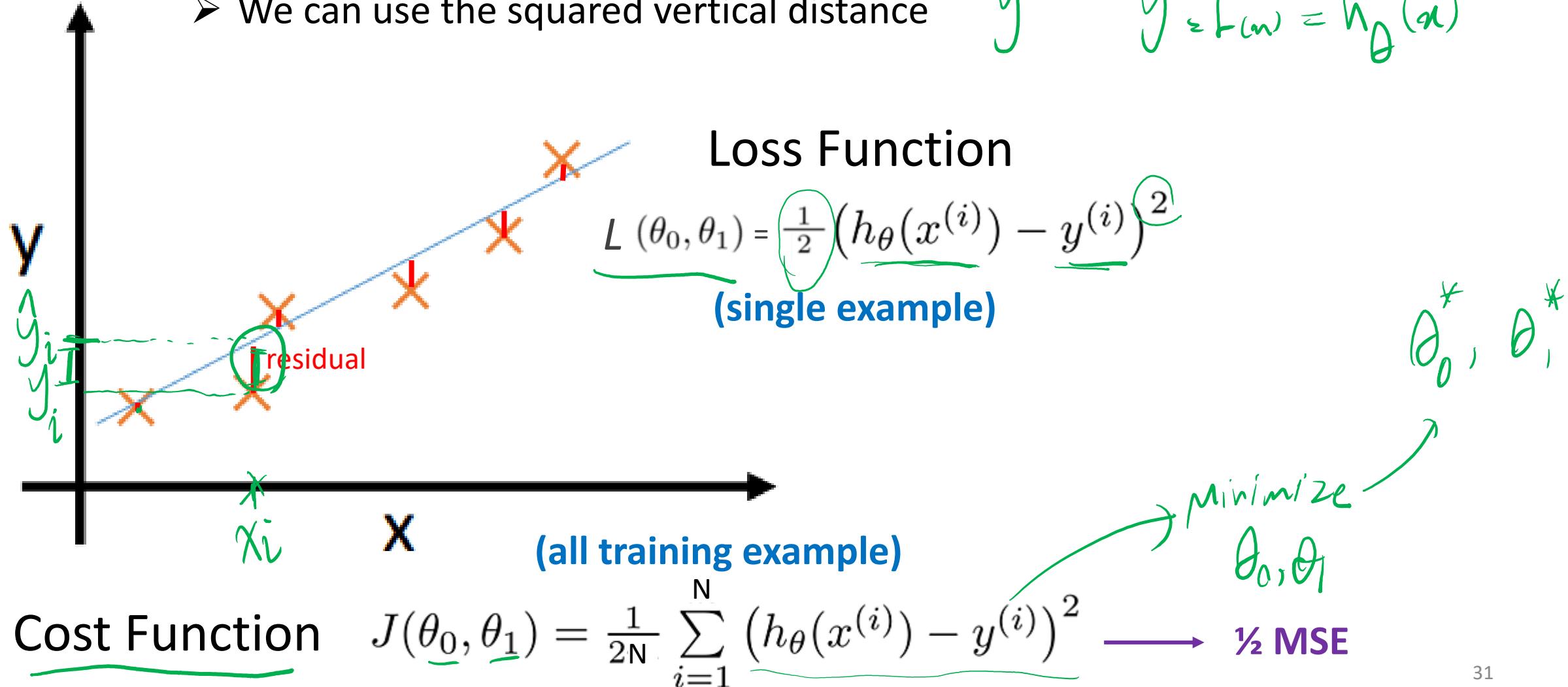
But what does
“close” mean?

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

$$\{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$$

Loss Function

- We can use the squared vertical distance



Learning a Hypothesis

Hypothesis: $y \sim h_\theta(x) = \theta_0 + \theta_1 x$

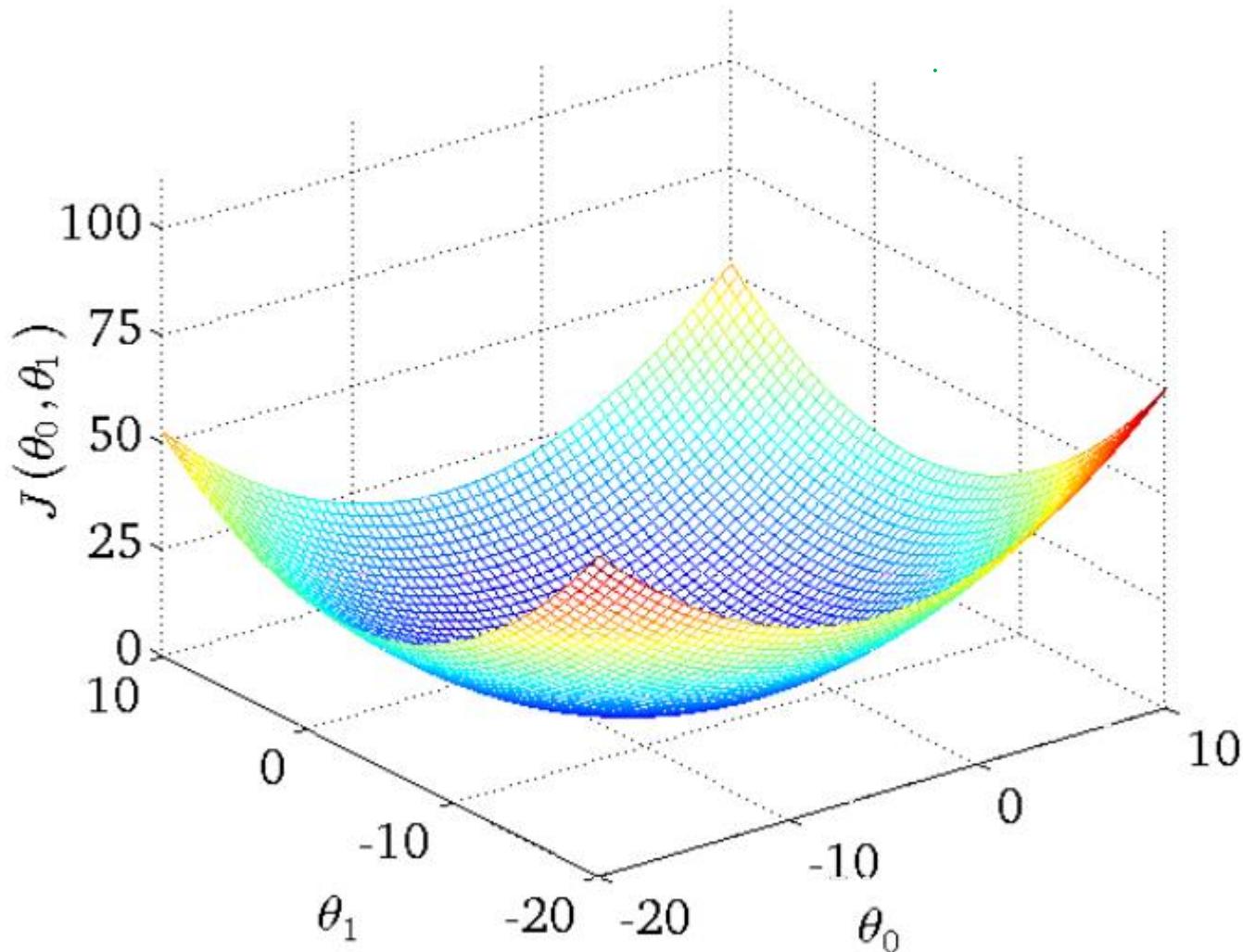
Parameters: θ_0, θ_1

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

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

θ_0, θ_1
model

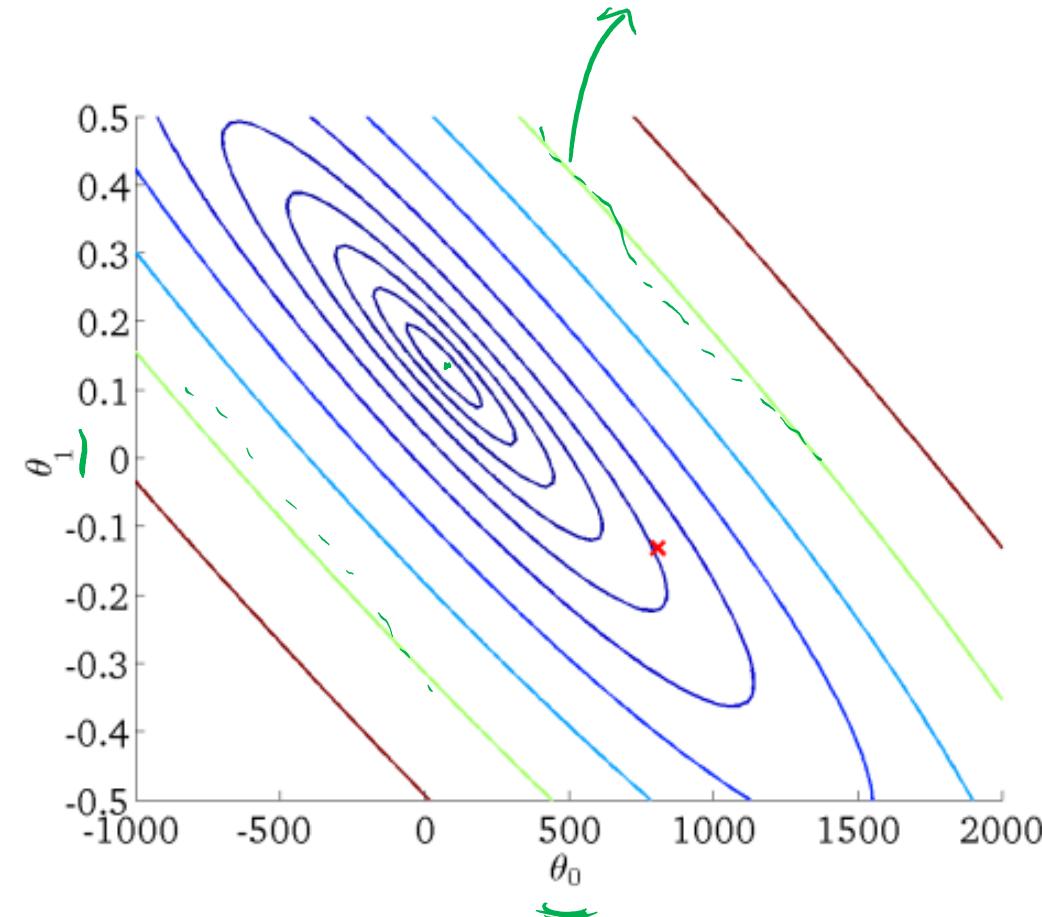
Cost Function Surface Plot



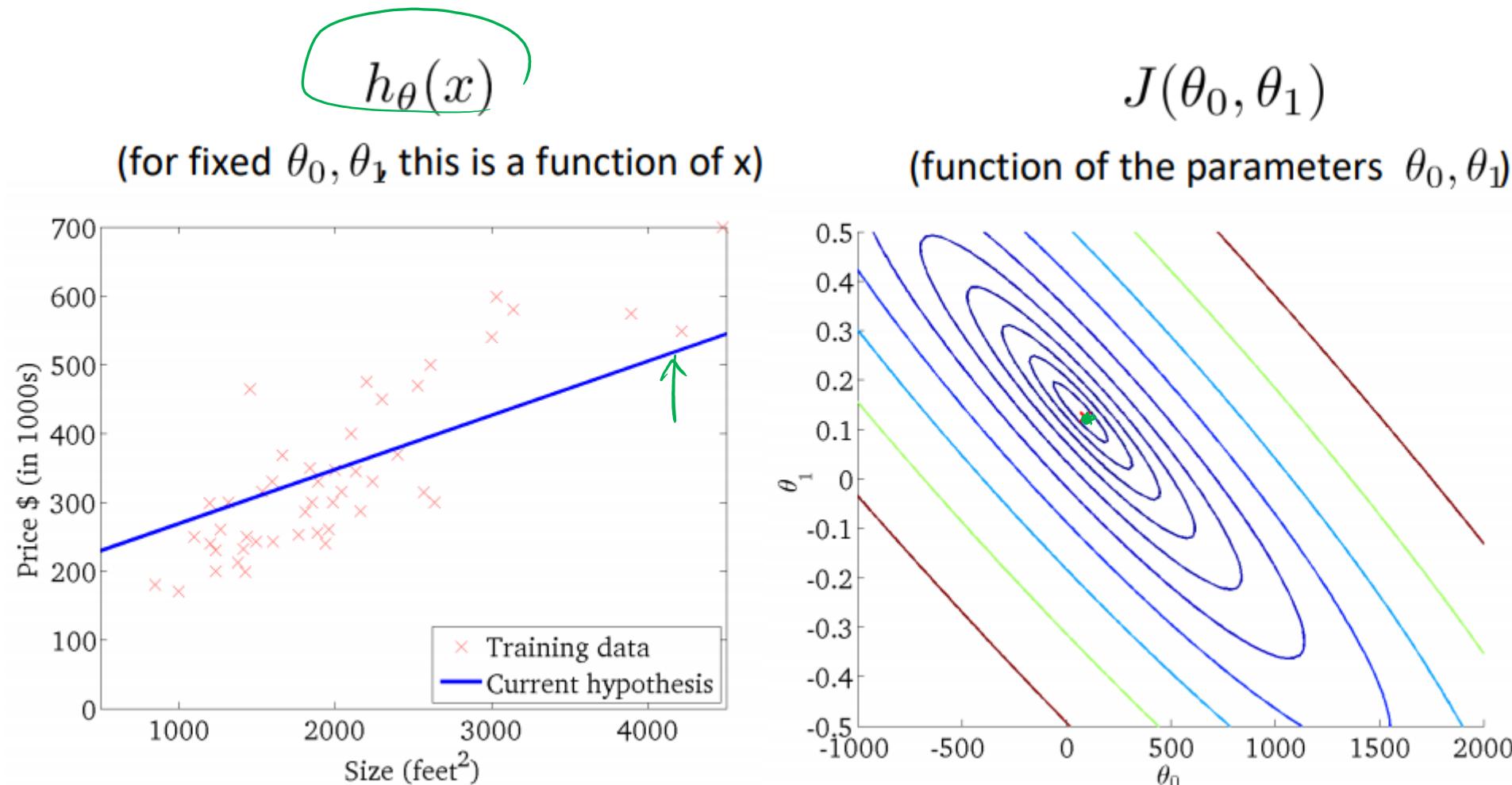
**Notice the
convexity of
the plot**

Cost Function Surface Plot

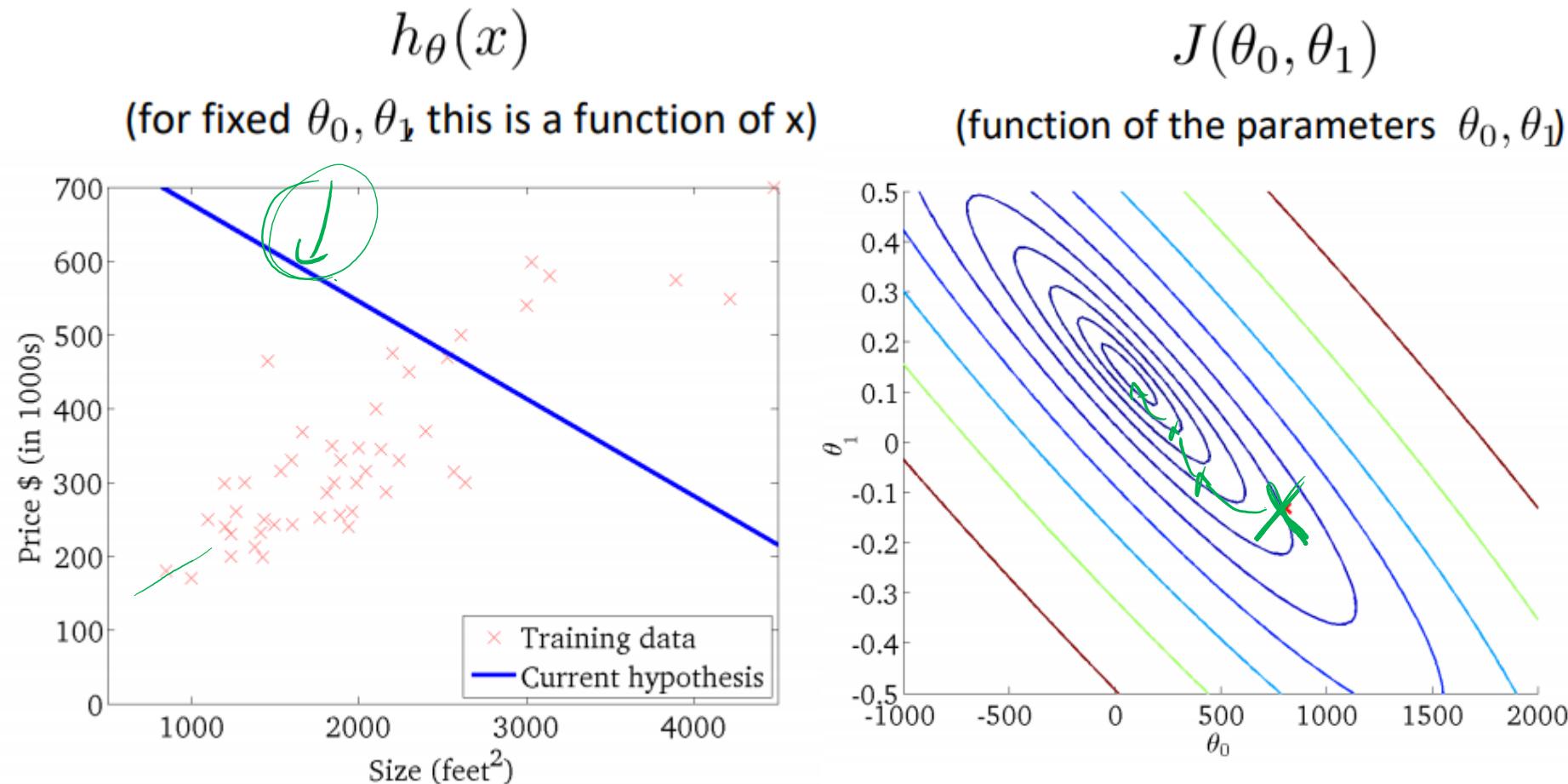
- We can visualize our **Cost Function** $J(\theta_0, \theta_1)$ from above as a contour plot.
- Where contours will be represented as curves in the graph where values of the cost, $J(\theta_0, \theta_1)$ are equal.



Contour Plots



Cost Function Contour Plot



Direct Solution

- The minimum would occur where partial derivatives equal zero:

$$\frac{dJ}{d\theta_i} = 0$$

Linear regression is one of a handful of models that permit direct solutions

- Which results in a **direct solution**:

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



Alternative Approach?

- Q: Great! If we have an analytical solution to finding optimal parameters, what's the issue?

$O(n^3)$

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

- A: In high dimensional spaces **computing matrix inversion is expensive!**
- Instead, we will use a numerical solution known as **gradient descent**

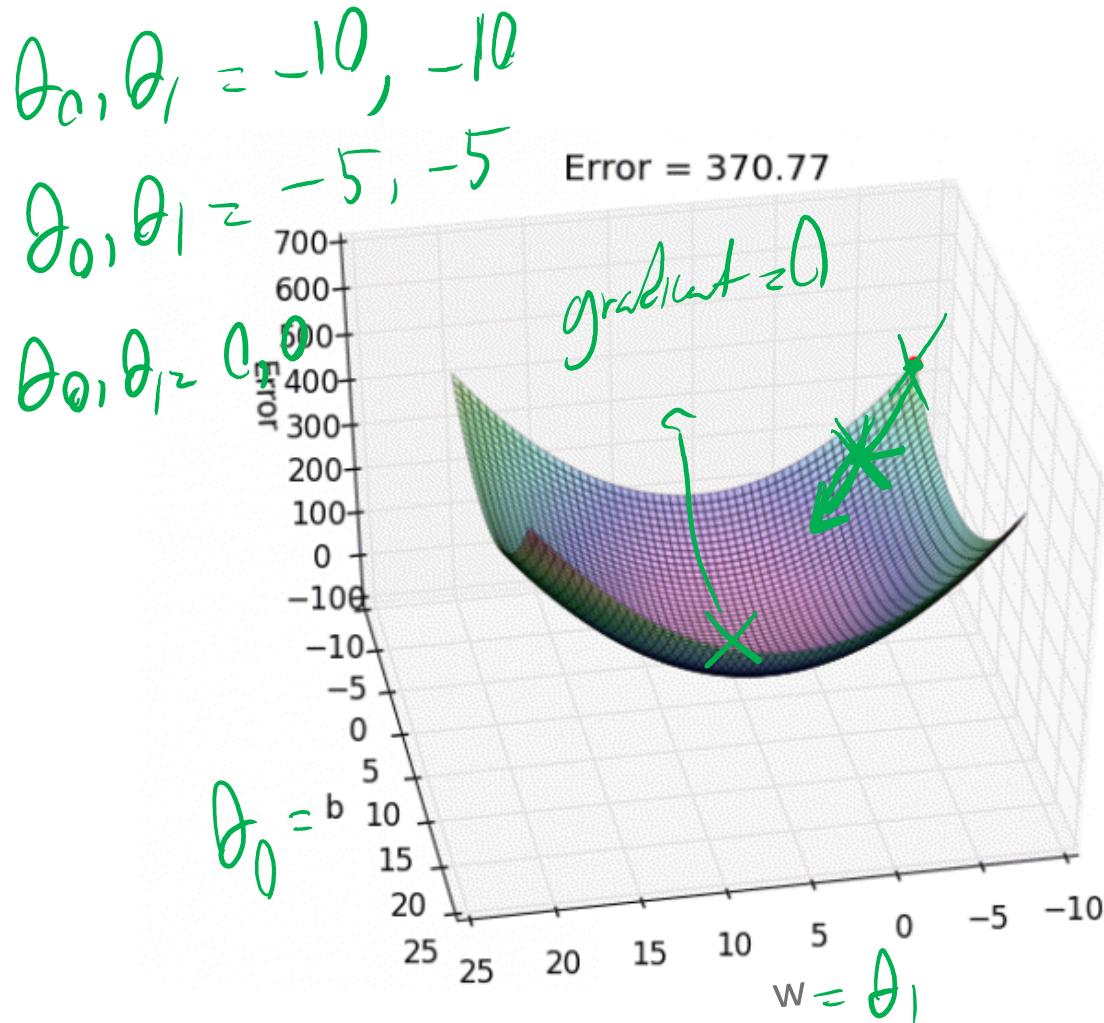
Gradient Descent



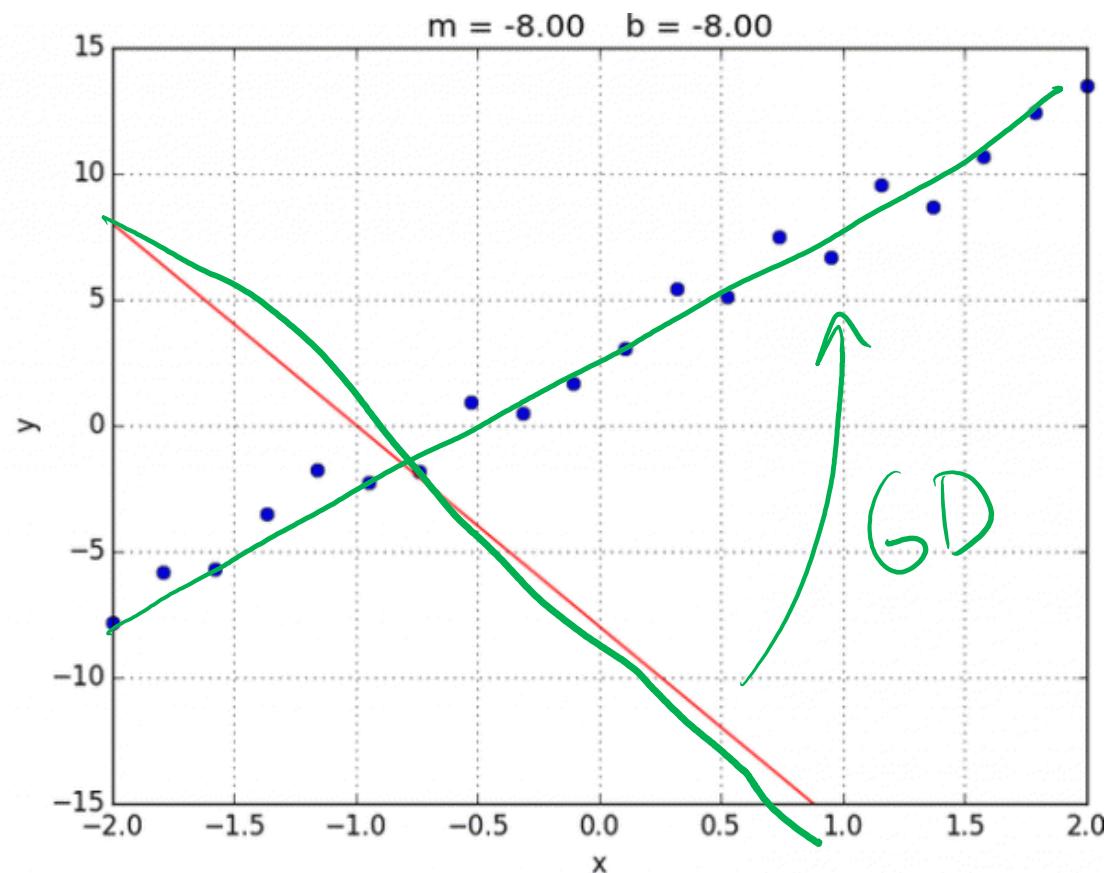
Gradient descent is an iterative algorithm.

We **initialize** a starting point and **repeatedly adjust** based on the direction of **steepest descent**.

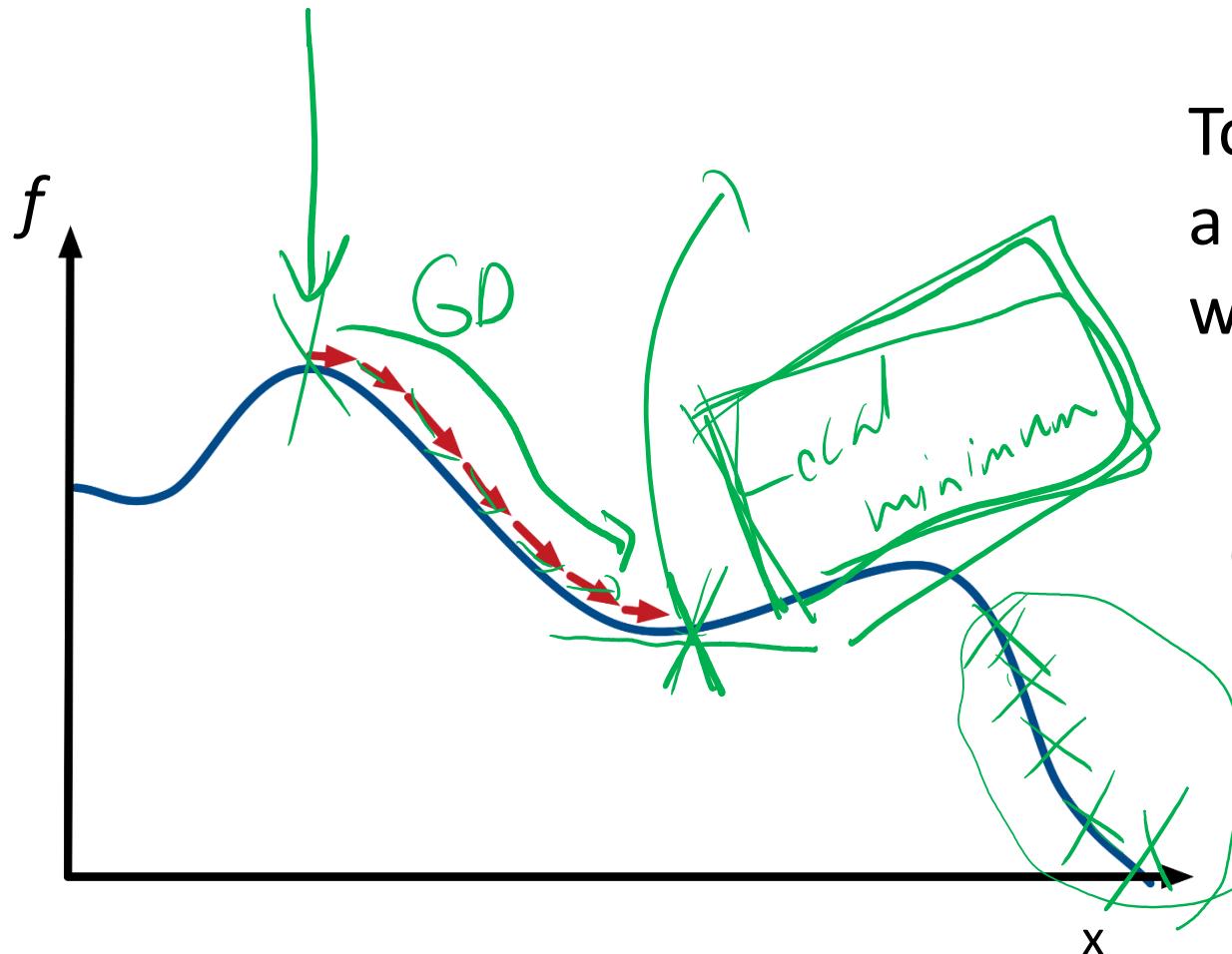
Gradient Descent



https://alykhantejani.github.io/images/gradient_descent_line_graph.gif



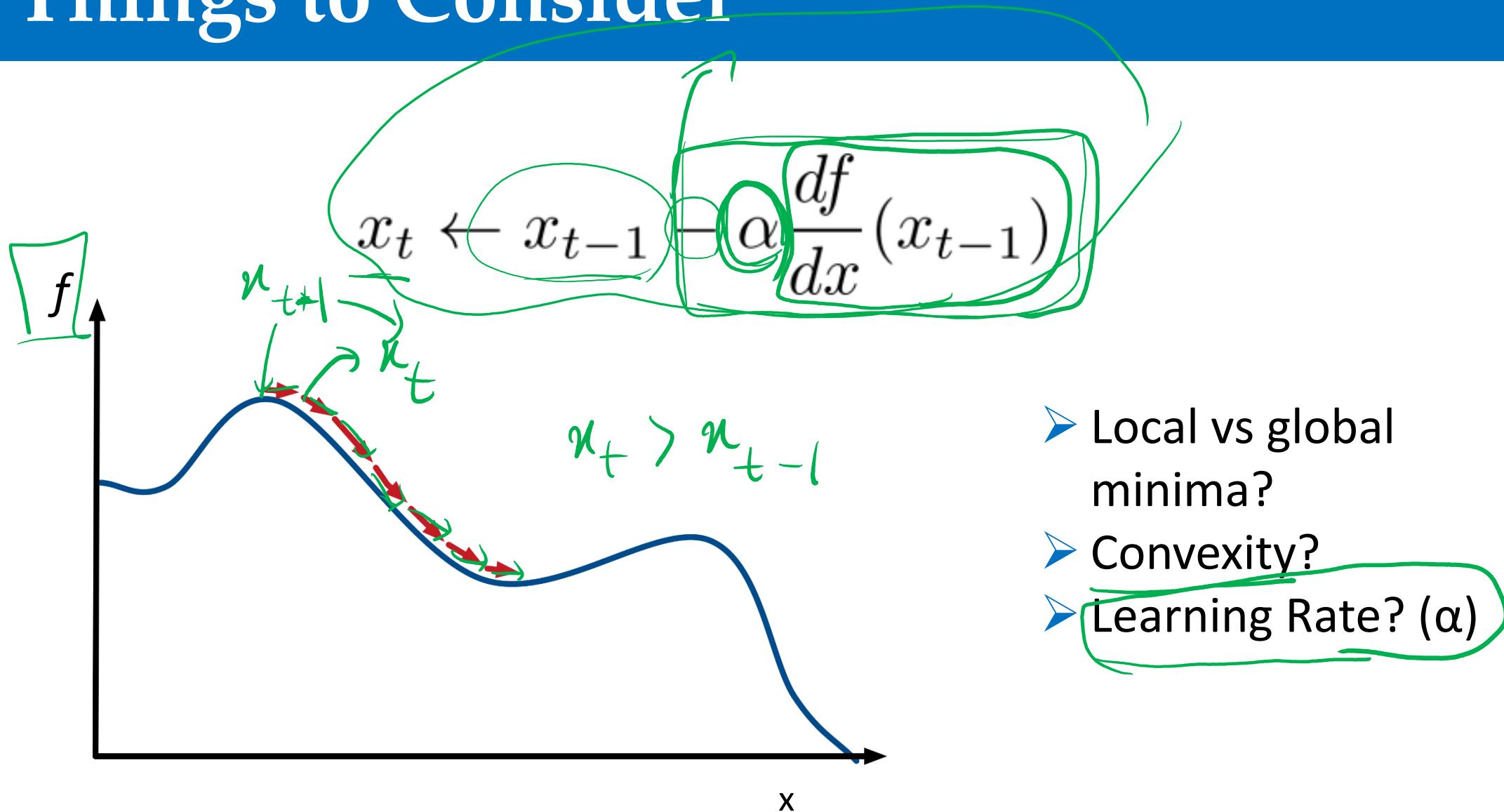
Gradient Descent in 1D



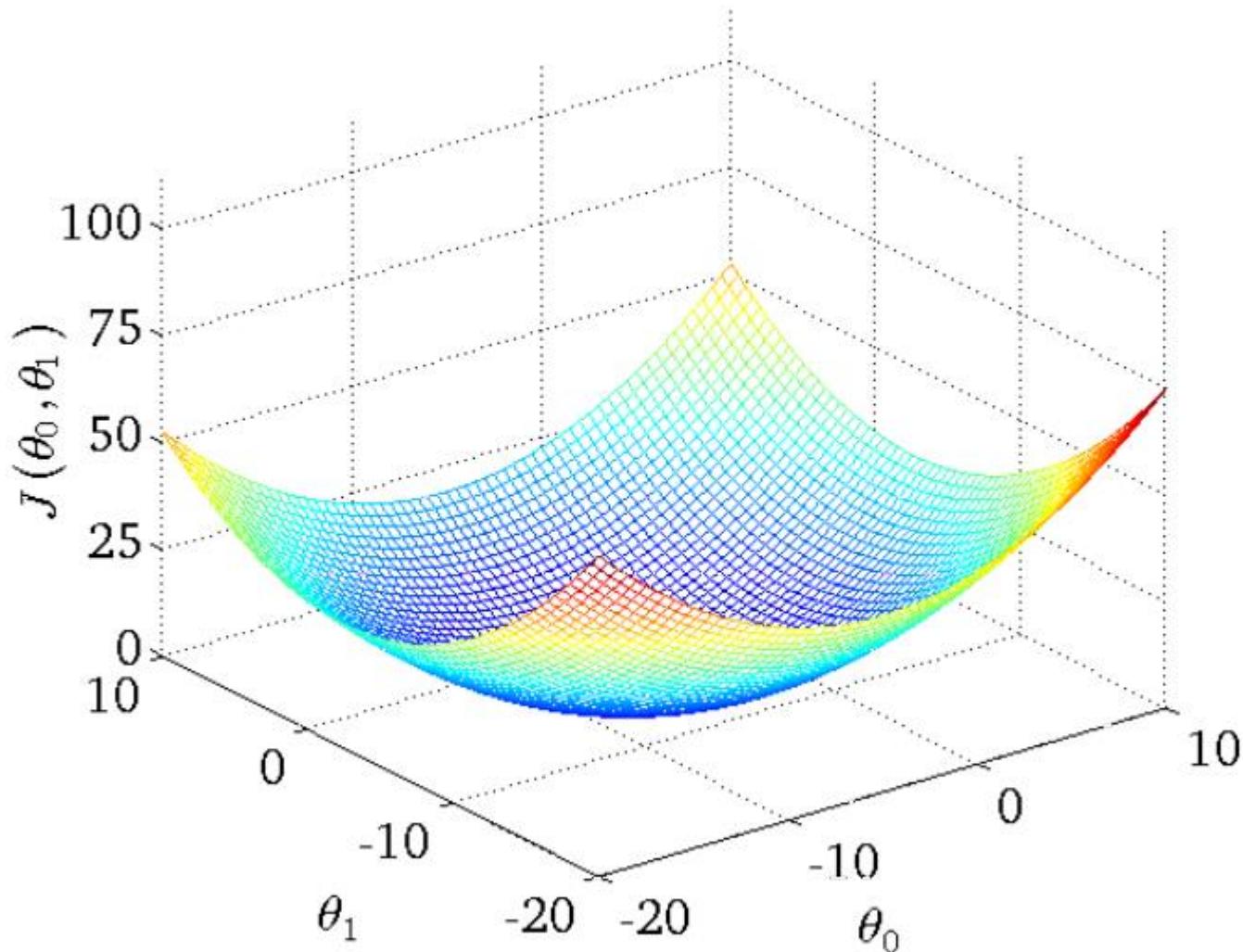
To minimize $f(x)$, we start with a random point and iterate with the update rule:

$$x_t \leftarrow x_{t-1} - \alpha \frac{df}{dx}(x_{t-1})$$

Things to Consider

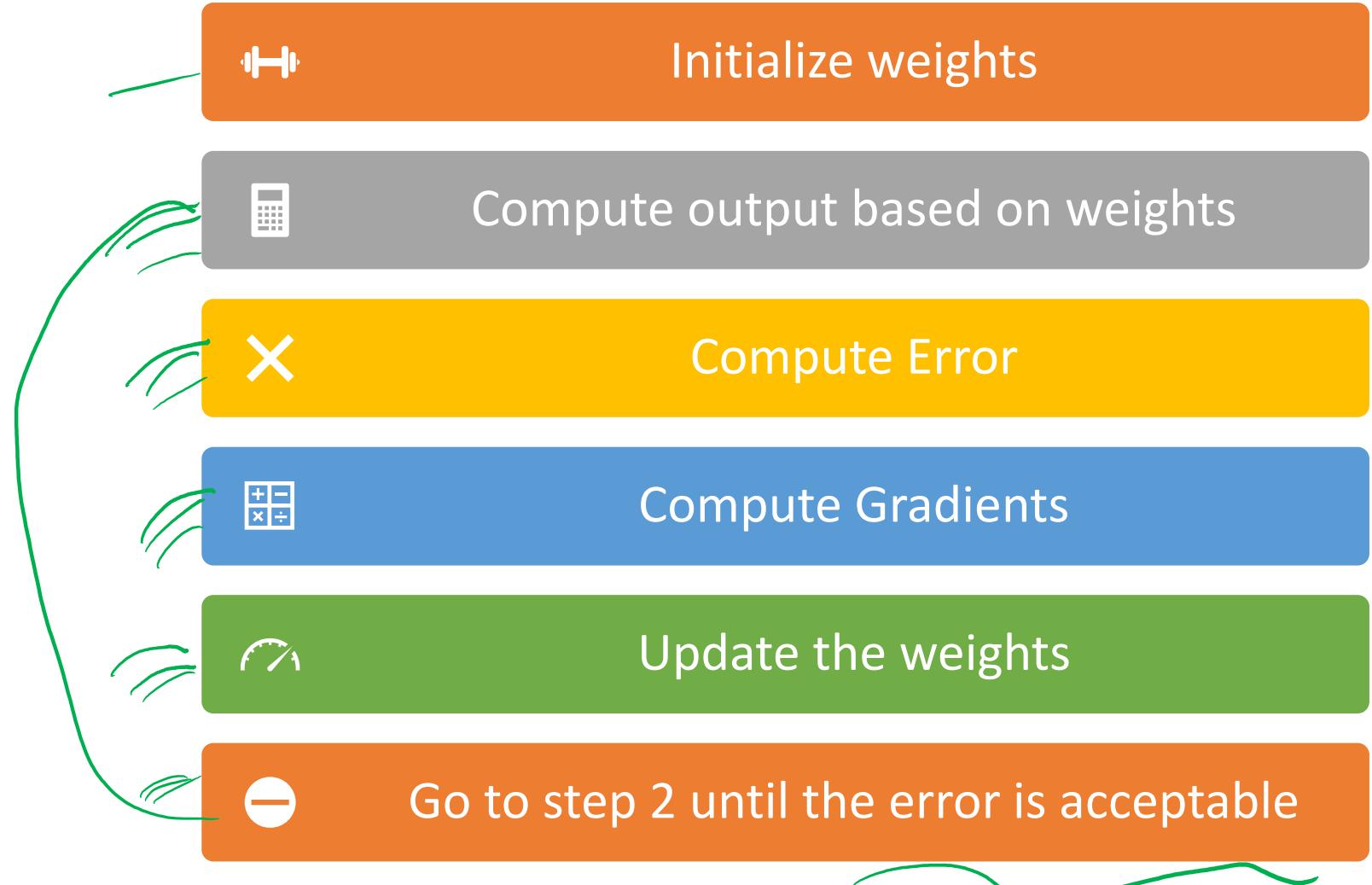


Cost Function Surface Plot

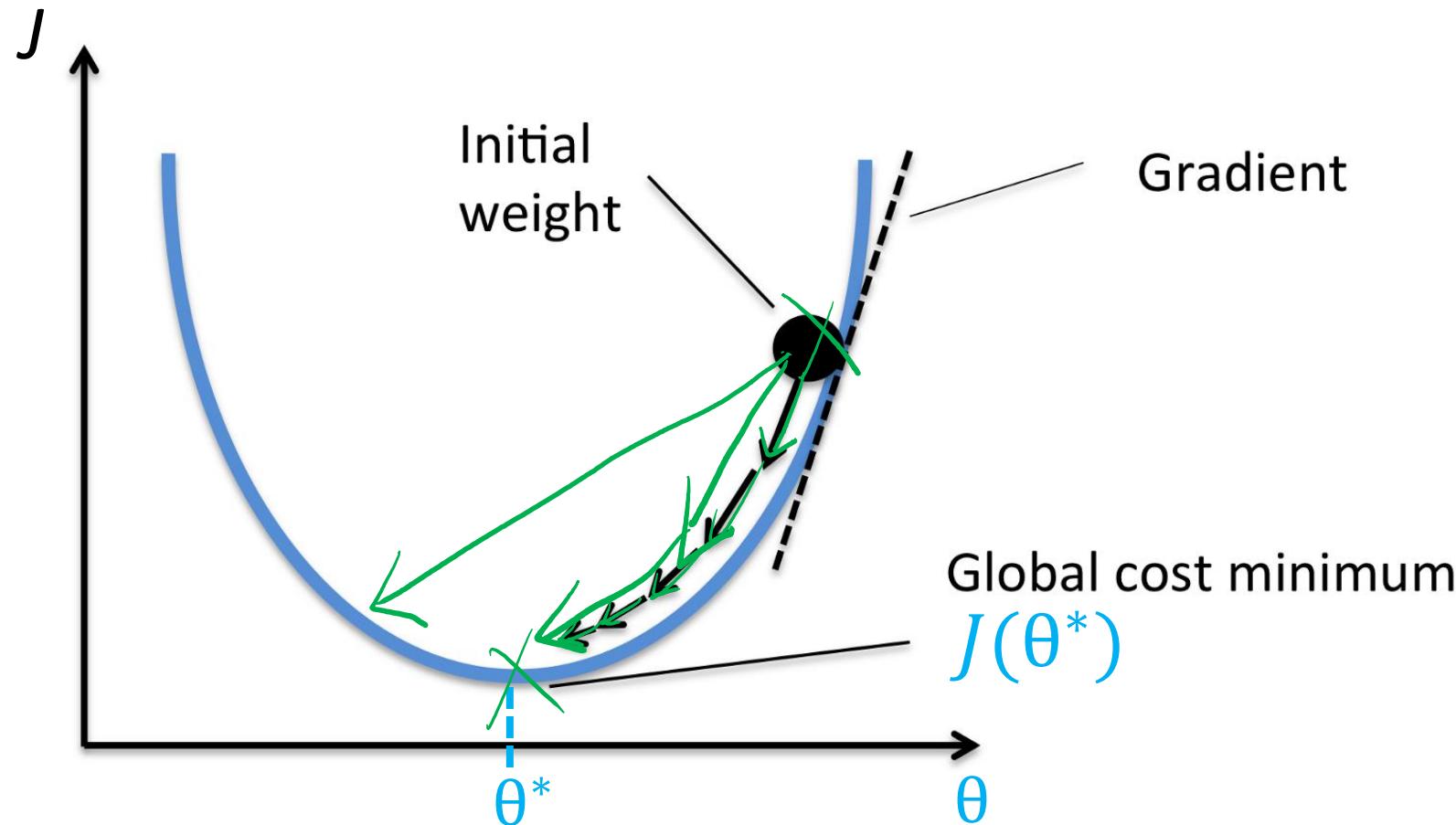


**Notice the
convexity of
the plot**

GD Summary



Role of Learning Rate?

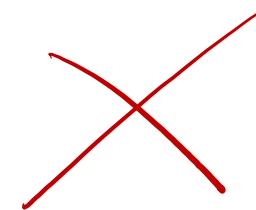
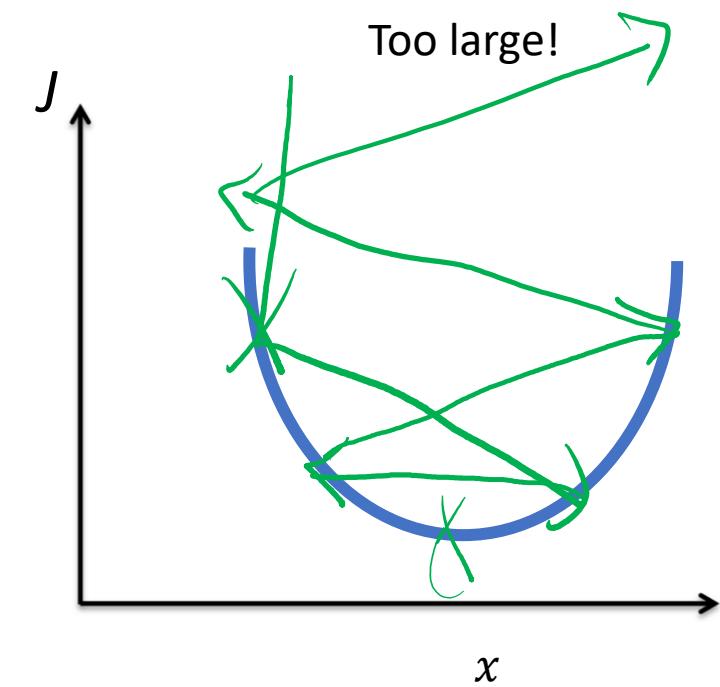
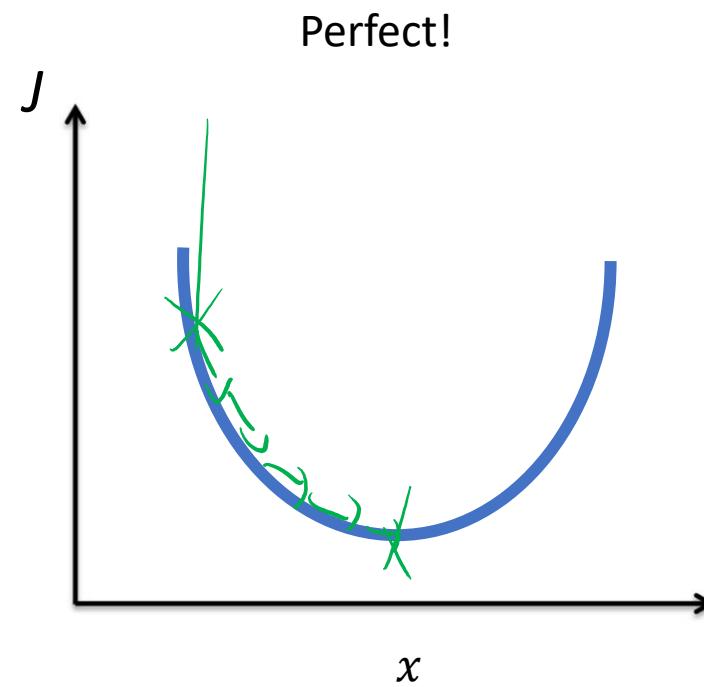
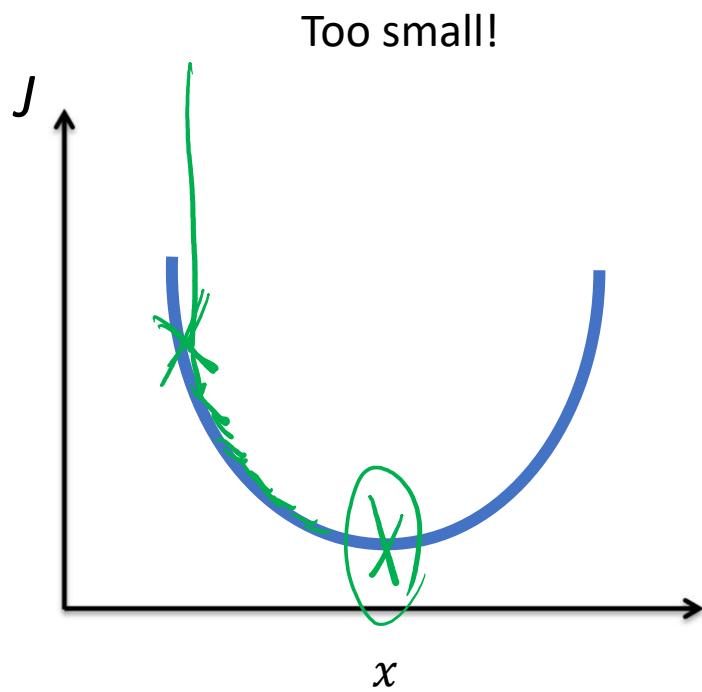


$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m \left(\underbrace{h_\theta(x^{(i)}) - \underline{y}^{(i)}}_{\text{Error}} \right)^2$$

$\theta_j \leftarrow \theta_j - \alpha \frac{\partial J(\theta_j)}{\partial \theta_j}$

- α is a **learning rate**. The larger it is, the faster θ_j changes.
 - We'll see later how to tune the learning rate, but values are typically small, e.g. 0.01 or 0.0001

Role of Learning Rate?



Multivariable Regression

- Suppose we have multiple inputs x_1, \dots, x_D . This is referred to as multivariable regression.
- This is no different than the single input case, just harder to visualize.

➤ Linear model:

$$\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_D x_D$$

$$\hat{y} = \sum_j \theta_j x_j + \theta_0$$

Implementation

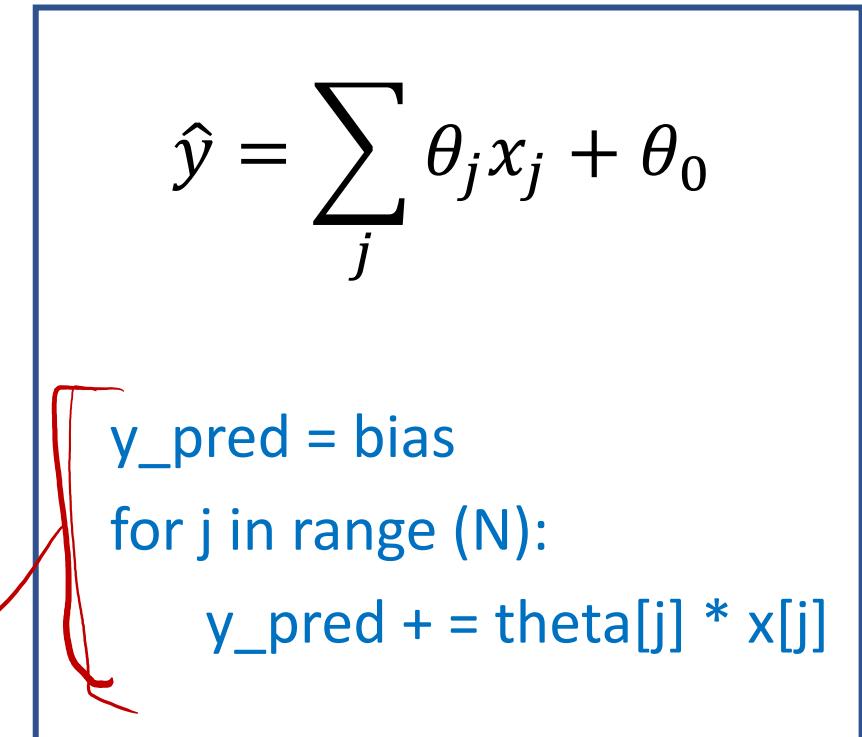
- Computing the prediction using a for loop:
- For-loops in Python are slow, so we vectorize algorithms by expressing them in terms of vectors and matrices.

$$\theta = (\underline{\theta_1, \dots, \theta_D})^T \quad \mathbf{x} = (\underline{x_1, \dots, x_D})^T$$

$$\hat{y} = \underline{\theta^T \mathbf{x}} + \underline{\theta_0}$$

- This is simpler and much faster:

```
y_pred = np.dot(theta, x) + bias
```



Why Vectorize?

- The equations, and the code, will be simpler and more readable.
Gets rid of dummy variables/indices!
- Vectorized code is much faster
 - Cut down on Python interpreter overhead
 - Use highly optimized linear algebra libraries
 - Matrix multiplication is very fast on a Graphical Processing Unit (GPU)

Matrix of Data

- We can take this step further. Organize all the training examples into the **design matrix** \mathbf{X} with one row per training example, and all the targets into the **target vector** \mathbf{y} .

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}^{(1)\top} \\ \mathbf{x}^{(2)\top} \\ \mathbf{x}^{(3)\top} \end{pmatrix} \quad \begin{pmatrix} 1 & \text{---} & 8 & 0 & 3 & 0 \\ 2 & \text{---} & 6 & -1 & 5 & 3 \\ 3 & \text{---} & 2 & 5 & -2 & 8 \end{pmatrix}$$

Annotations:

- Red arrows point from f_1, f_2, f_3, f_4 to the columns of the matrix.
- A green box highlights the first column of the matrix.
- A green box highlights the second column of the matrix.
- A red arrow points to the second column with the text "one feature across all training examples".
- A red arrow points to the second row with the text "Sample 2".
- A red arrow points to the second row with the text "one training example (vector)".

Design Matrix

- Computing the predictions for the entire dataset:

$$\mathbf{x}\theta + \underbrace{\theta_0 \mathbf{1}}_{n \times d} = \begin{pmatrix} \theta^T \mathbf{x}^{(1)} + \theta_0 \\ \vdots \\ \theta^T \mathbf{x}^{(N)} + \theta_0 \end{pmatrix} = \begin{pmatrix} \hat{y}^{(1)} \\ \vdots \\ \hat{y}^{(N)} \end{pmatrix} = \underline{\hat{y}}$$

Loss

- Computing the squared error cost across the entire dataset:

$$\hat{y} = \mathbf{x}\theta + \theta_0 \mathbf{1}$$

$$J = \frac{1}{2N} \|y - \hat{y}\|^2$$

- In Python:

```
y_pred = np.dot(X, theta) + bias
```

```
cost = np.sum((y_pred - y) ** 2) / (2 * N)
```

Compute Gradients

- Partial derivatives: derivatives of a multivariate function with respect to one of its arguments

$$\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

$$\frac{\partial \hat{y}}{\partial \theta_2} = x_2$$

$$\frac{\partial}{\partial x_1} f(x_1, x_2) = \lim_{h \rightarrow 0} \frac{f(x_1 + h, x_2) - f(x_1, x_2)}{h}$$

- To compute, take the single variable derivatives, pretending the other arguments are constant.

$$\hat{y} = \sum_j \theta_j x_j + \theta_0$$

$$\frac{\partial \hat{y}}{\partial \theta_j} = \frac{\partial}{\partial \theta_j} \left[\sum_{j'} \theta_{j'} x_{j'} + \theta_0 \right] \\ = x_j$$

$$\frac{\partial \hat{y}}{\partial \theta_0} = \frac{\partial}{\partial \theta_0} \left[\sum_{j'} \theta_{j'} x_{j'} + \theta_0 \right] \\ = 1$$

Compute Gradients

- Chain rule for derivatives:

$$\mathcal{L}(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2$$

$$\hat{y} = \sum_j \theta_j x_j + \theta_0$$

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial \theta_j} &= \frac{d\mathcal{L}}{d\hat{y}} \frac{\partial \hat{y}}{\partial \theta_j} \\ &= \frac{d}{d\hat{y}} \left[\frac{1}{2} (\hat{y} - y)^2 \right] \cdot x_j \\ &= (\hat{y} - y) \cdot x_j\end{aligned}$$

$\frac{\partial \mathcal{L}}{\partial \theta_0} = (\hat{y} - y)$

- Cost derivative (averages over data points):

$$\frac{\partial \mathcal{J}}{\partial \theta_j} = \frac{1}{N} \sum_{i=1}^N (\hat{y}^{(i)} - y^{(i)}) x_j^{(i)}$$

$$\frac{\partial \mathcal{J}}{\partial \theta_0} = \frac{1}{N} \sum_{i=1}^N (\hat{y}^{(i)} - y^{(i)})$$

Concise Notation

- The bias is often included inside the design matrix \mathbf{X} for convenience as a column of ones.

$$\hat{y} = \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_0 \cancel{1}$$
$$\hat{y} = \underline{\theta_0}(1) + \theta_1 x_1 + \theta_2 x_2 + \cdots$$

$$\hat{y} = \sum_{j=0}^N \theta_j x_j$$

Parameter Update

- Cost derivative (averages over data points):

$$\hat{y} = \sum_{j=0}^N \theta_j x_j$$



$$\frac{\partial J}{\partial \theta_j} = \boxed{\frac{1}{N} \sum_{i=0}^N (\hat{y}^{(i)} - y^{(i)}) x_j^{(i)}}$$

only change starting index and design matrix \mathbf{X}

- Parameter Update:

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial J(\theta_j)}{\partial \theta_j}$$

Vectorize the update

- We can also vectorize this gradient computation:

$$\frac{\partial \mathcal{J}}{\partial \theta} = (\mathbf{y} - \hat{\mathbf{y}})^T \mathbf{X}$$
$$\hat{\mathbf{y}} = \mathbf{X}\theta$$

- and parameter update:

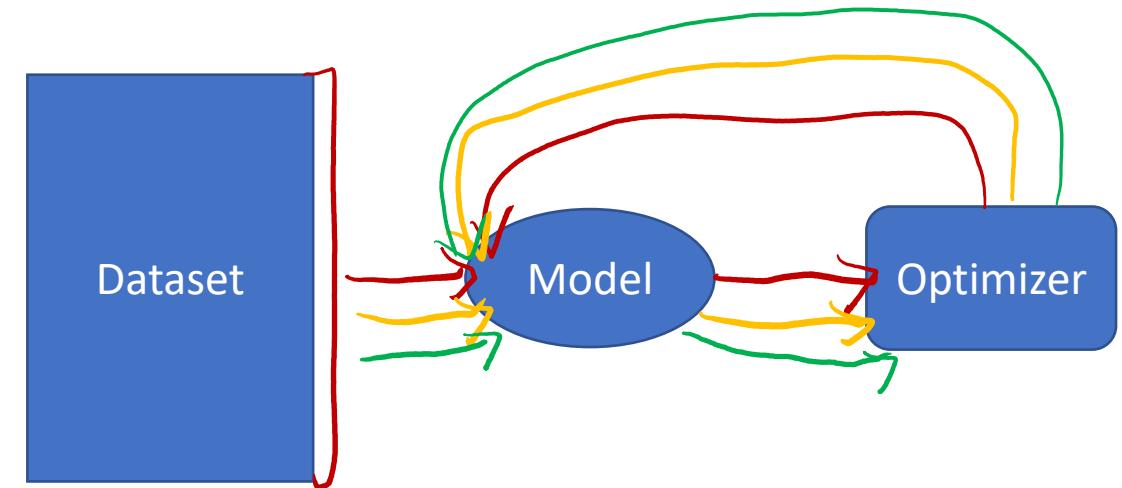
$$\theta \leftarrow \theta - \alpha \frac{\partial \mathcal{J}}{\partial \theta}$$

10:45

Break

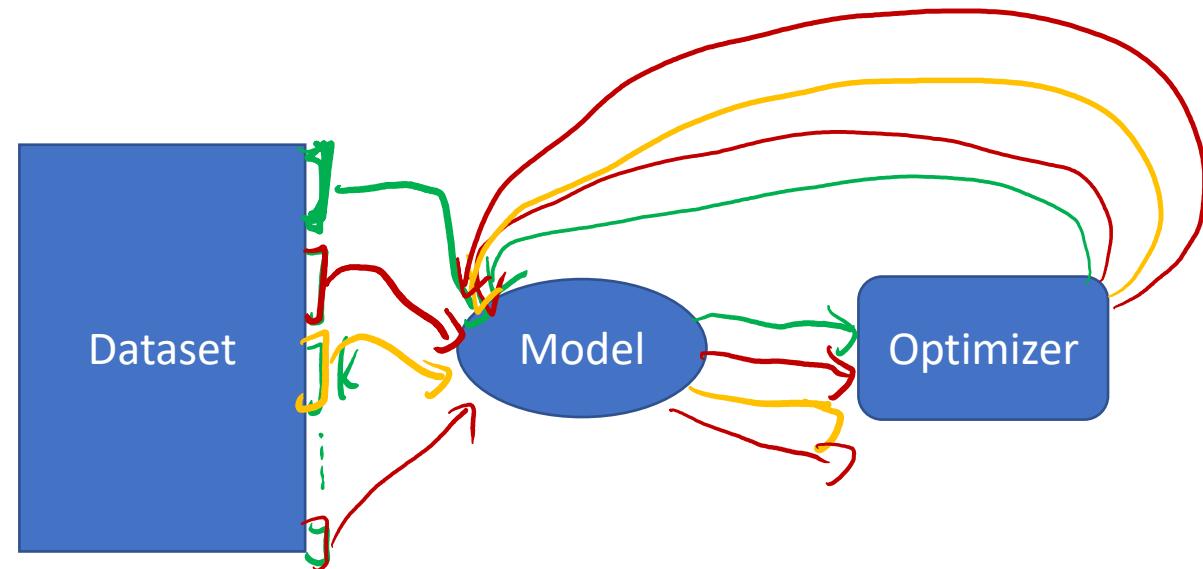
Batching in Gradient Descent

- Iteration:
 - Each time we update the weights is called an iteration
- Epoch:
 - Each time the model sees (learns) the whole dataset.
- Full batch GD:
 - Whole dataset is one batch.
 - One epoch is one iteration.



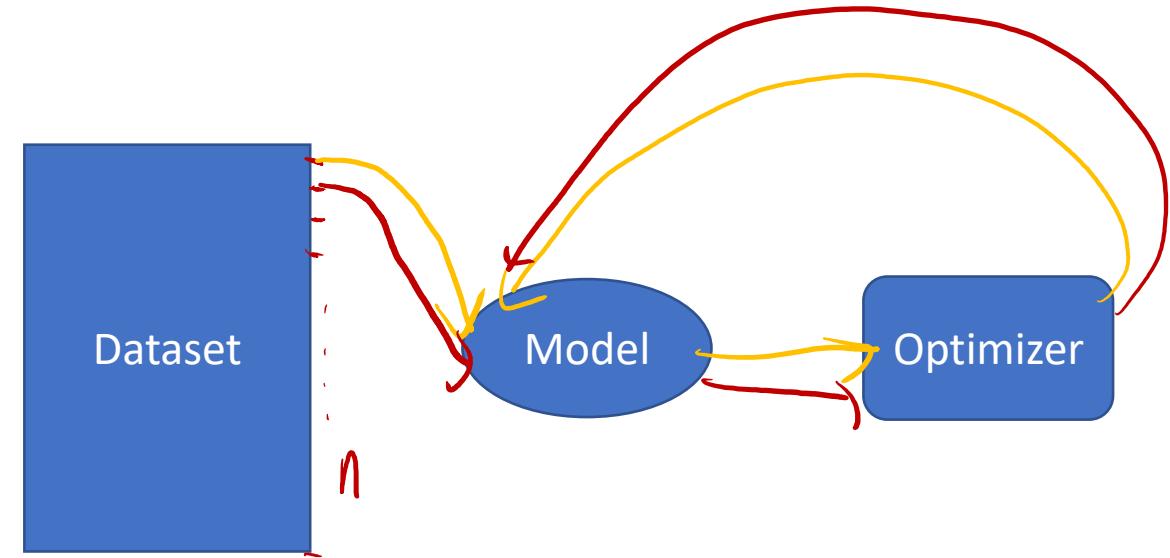
Batching in Gradient Descent

- Iteration:
 - Each time we update the weights is called an iteration
- Epoch:
 - Each time the model sees (learns) the whole dataset.
- Full batch GD:
 - Whole dataset is one batch.
 - One epoch is one iteration.
- Mini batch GD:
 - Break dataset to k smaller batches (mini batch).
 - One epoch takes k iterations.

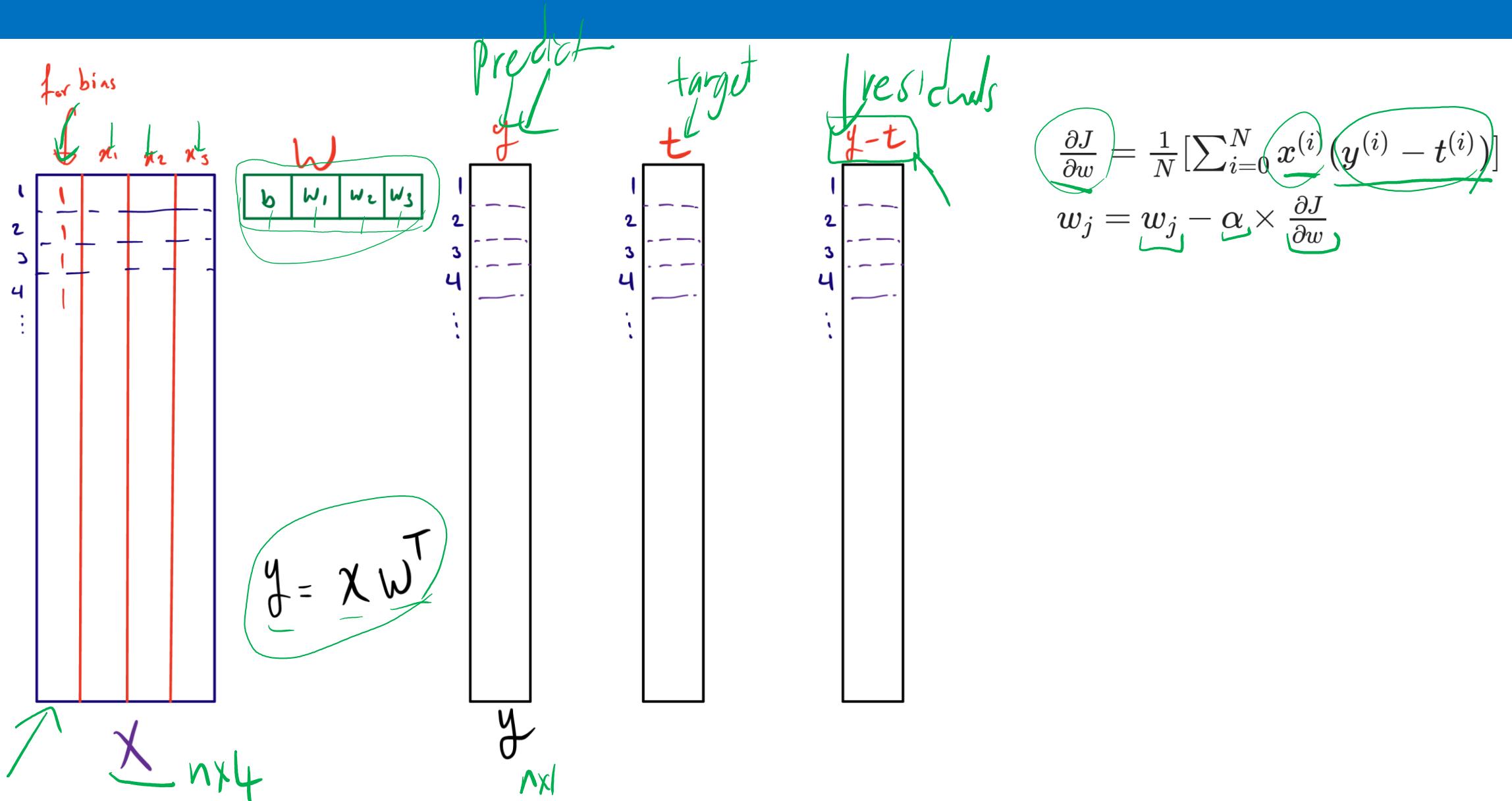


Batching in Gradient Descent

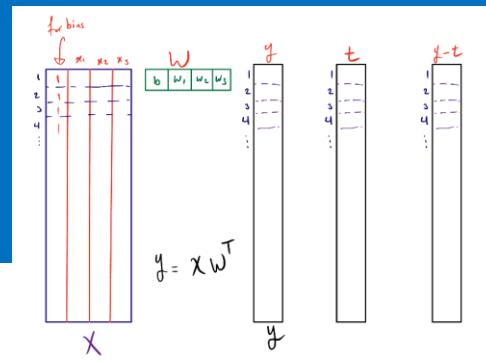
- Iteration:
 - Each time we update the weights is called an iteration
- Epoch:
 - Each time the model sees (learns) the whole dataset.
- Full batch GD:
 - Whole dataset is one batch.
 - One epoch is one iteration.
- Mini batch GD:
 - Break dataset to k smaller batches (mini batch).
 - One epoch takes k iterations.
- Stochastic GD:
 - Each of the n samples is a batch.
 - One epoch takes n iterations.



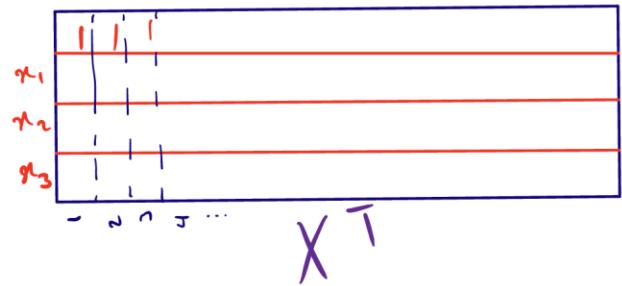
Full Batch



Full Batch

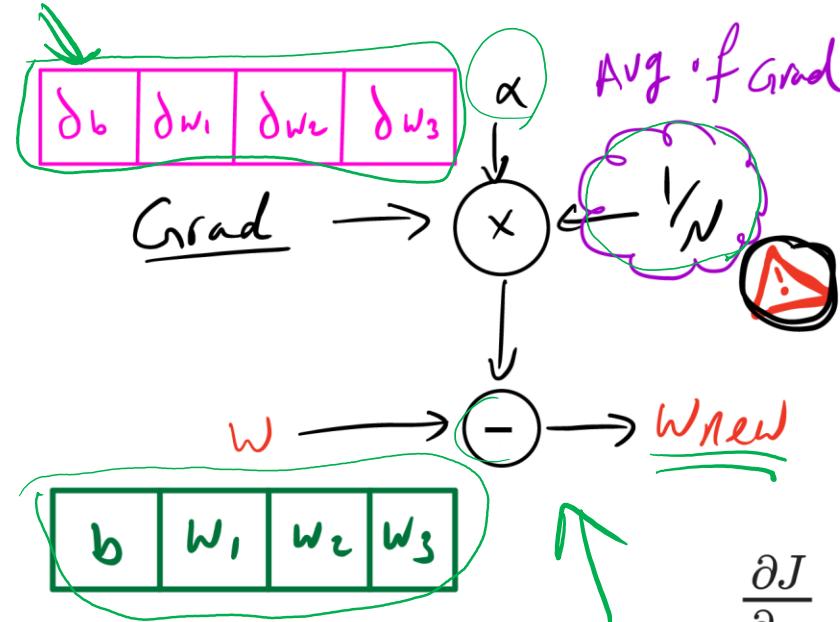


$$y = w_1 u_1 + w_2 u_2 + w_3 u_3 + b$$



$$\text{Grad} = \cancel{x^T(y-t)}$$

$4 \times n \quad n \times 1$



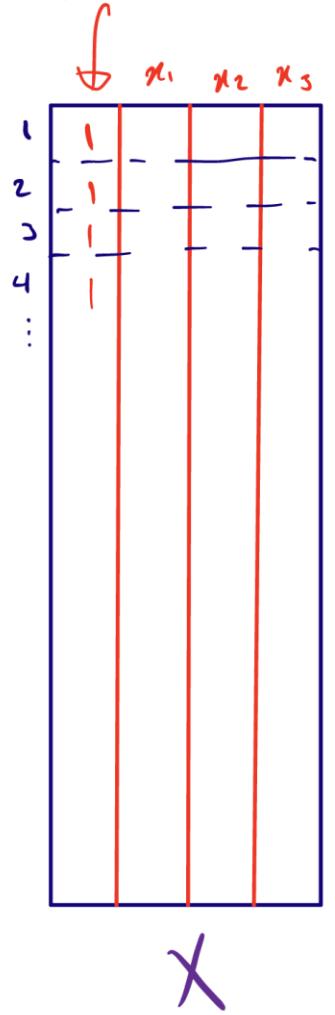
$$\frac{\partial J}{\partial w} = \frac{1}{N} \left[\sum_{i=0}^N x^{(i)} (y^{(i)} - t^{(i)}) \right]$$

$$w_j = w_j - \cancel{\alpha} \times \frac{\partial J}{\partial w}$$

1 iteration Per Epoch

Full Batch

for bins



$$w = \begin{matrix} b \\ w_1 \\ w_2 \\ w_3 \end{matrix}$$

$$y = X w^T$$

$$y = \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ \vdots \end{matrix}$$

$$t = \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ \vdots \end{matrix}$$

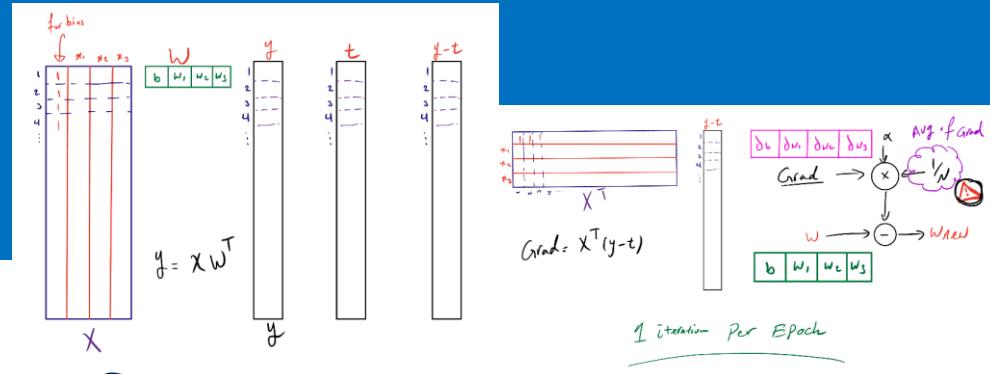
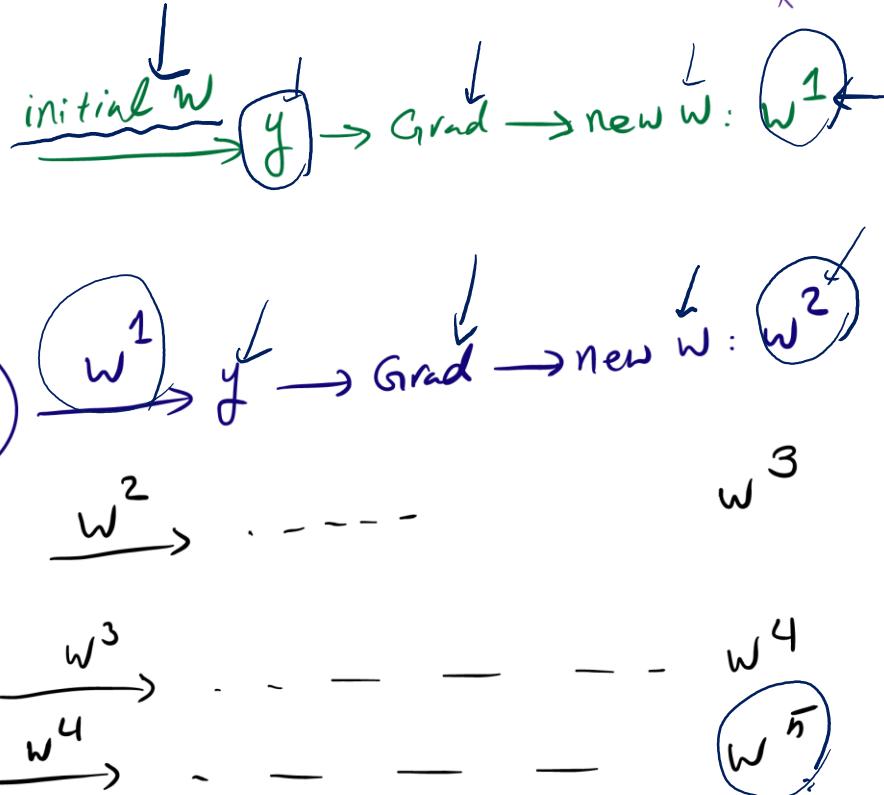
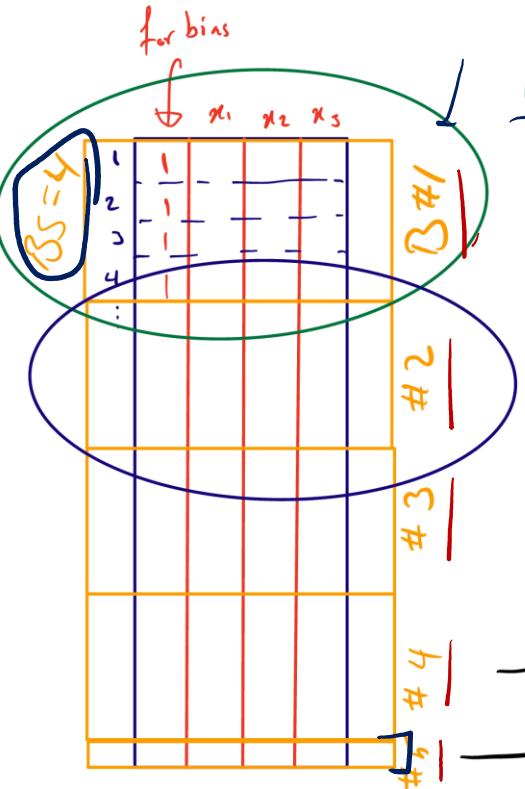
$$y - t = \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ \vdots \end{matrix}$$

$$\frac{\partial J}{\partial w} = \frac{1}{N} [\sum_{i=0}^N x^{(i)} (y^{(i)} - t^{(i)})]$$

$$w_j = w_j - \alpha \times \frac{\partial J}{\partial w}$$

Mini Batch

$n=17$



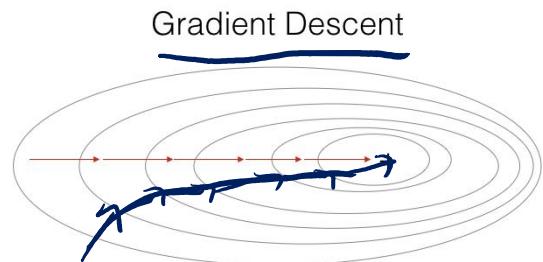
$$-10^{-10} \rightarrow -5, -5 \rightarrow 0, 0$$

X

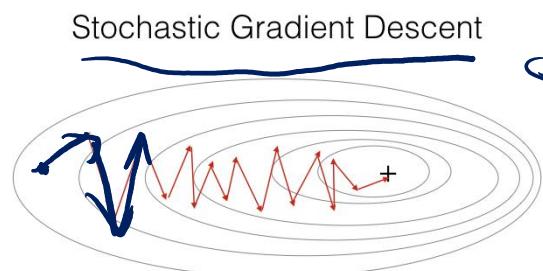
$y = w_1^* n_1 + w_2^* n_2 + w_3^* n_3 + b^*$

5 iteration per Epoch

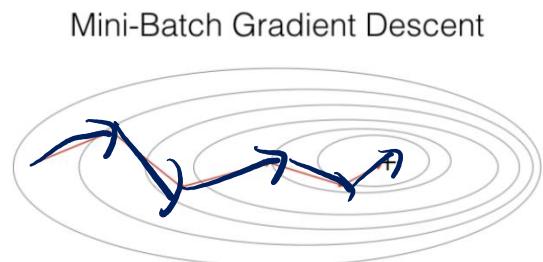
Batching and convergence



Smooth trajectory

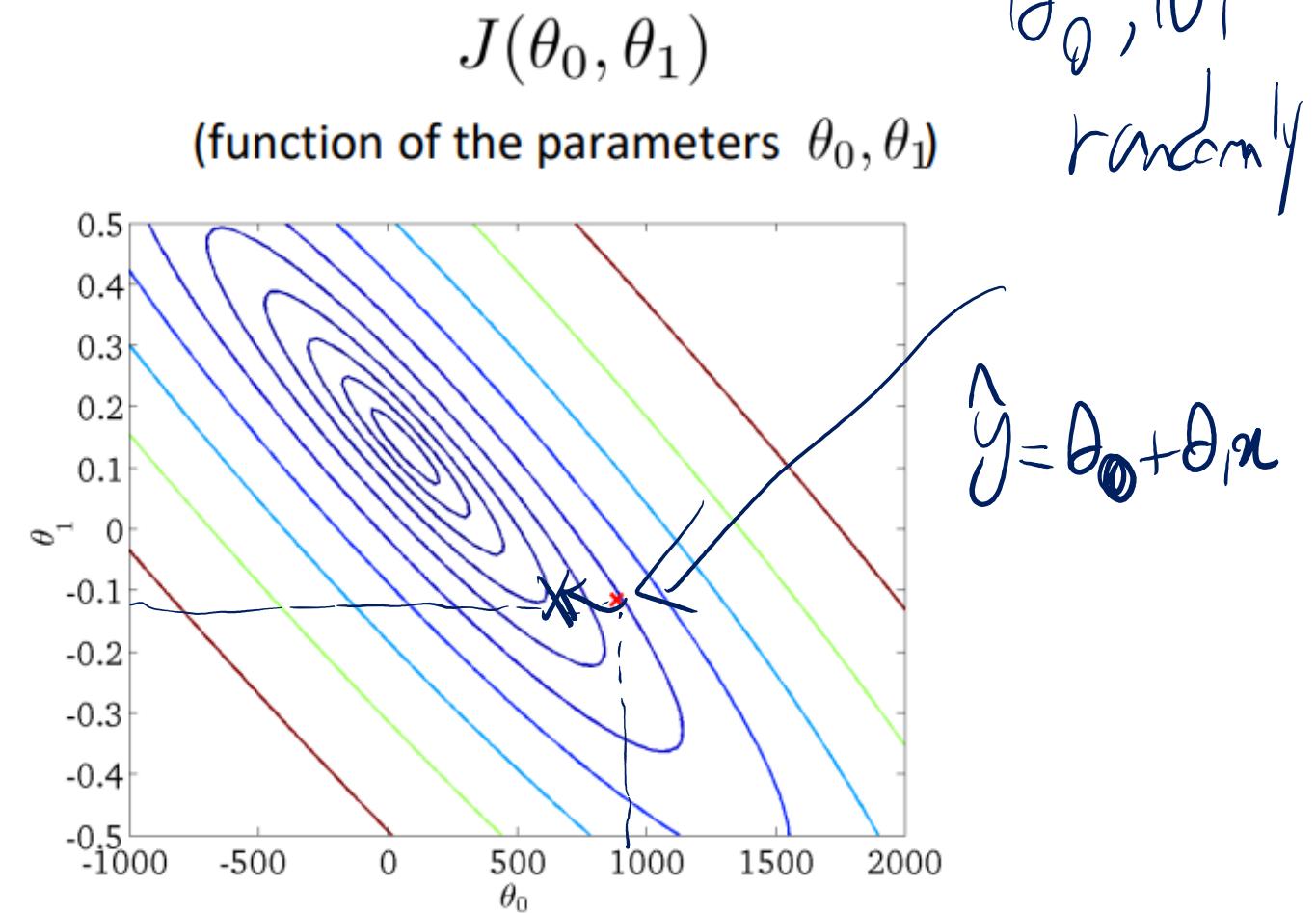
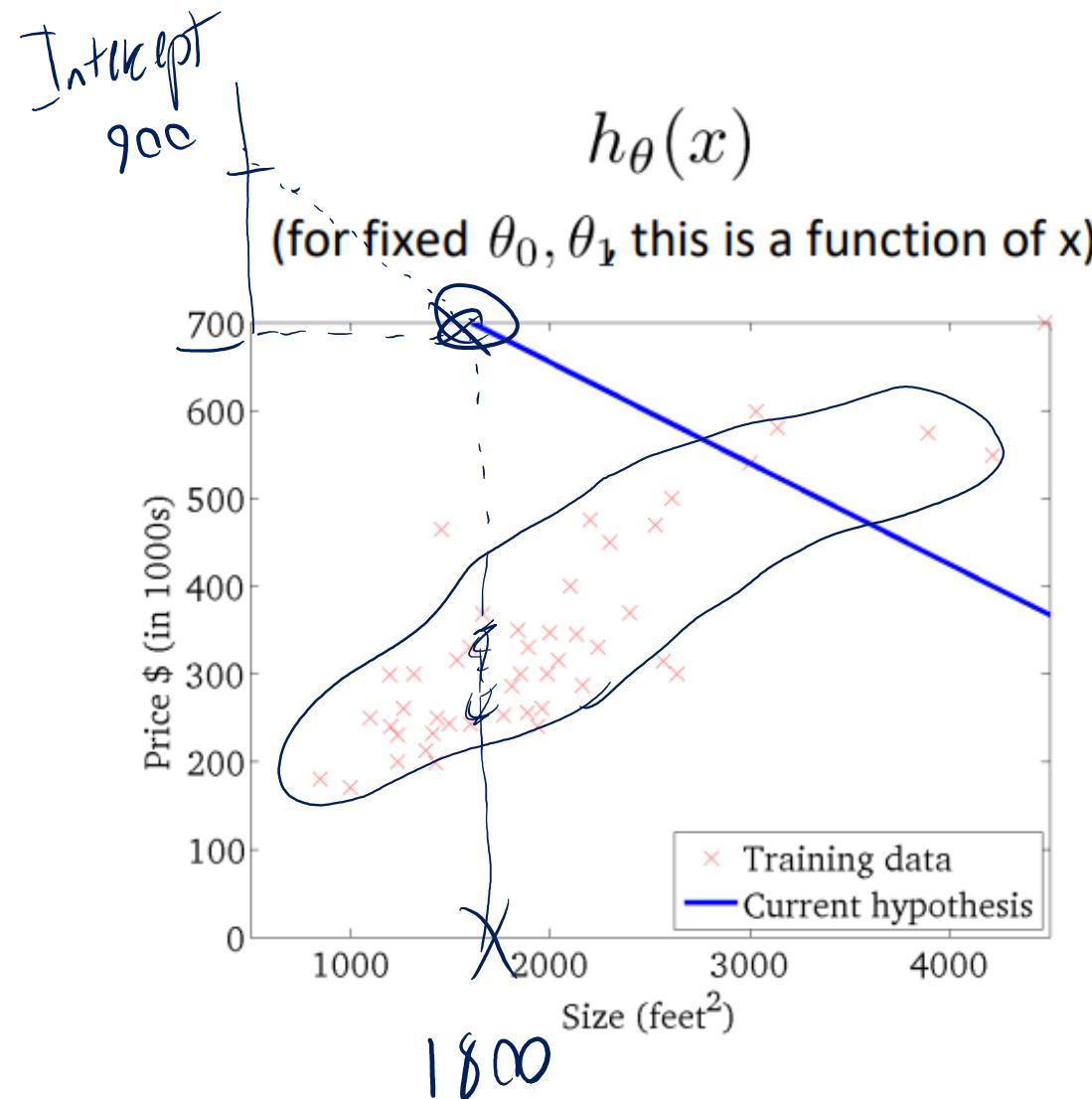


$h \approx 1/7$

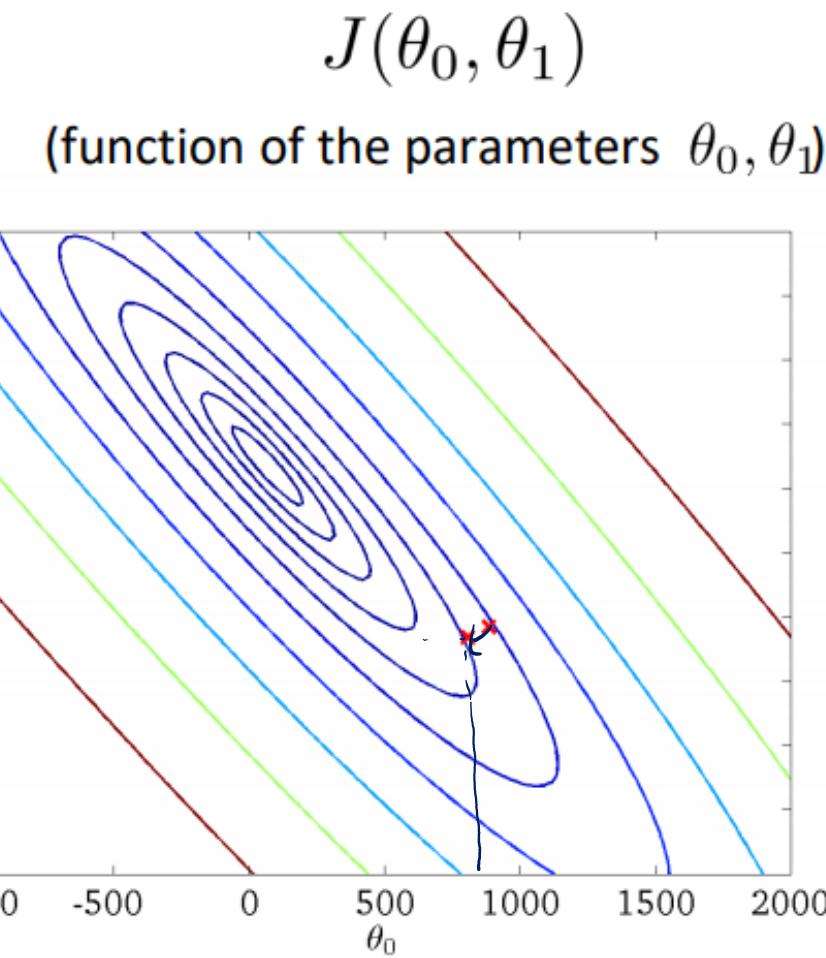
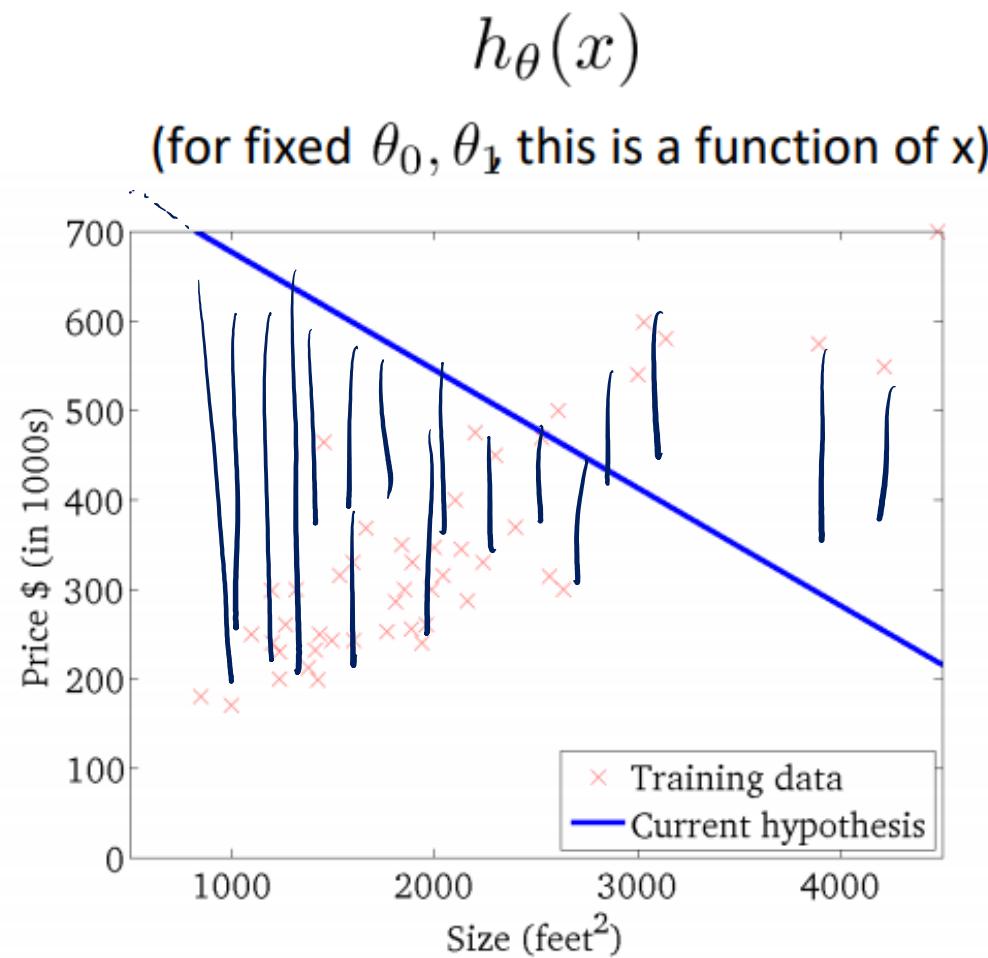


$K \approx 5$

Gradient Descent Example -1



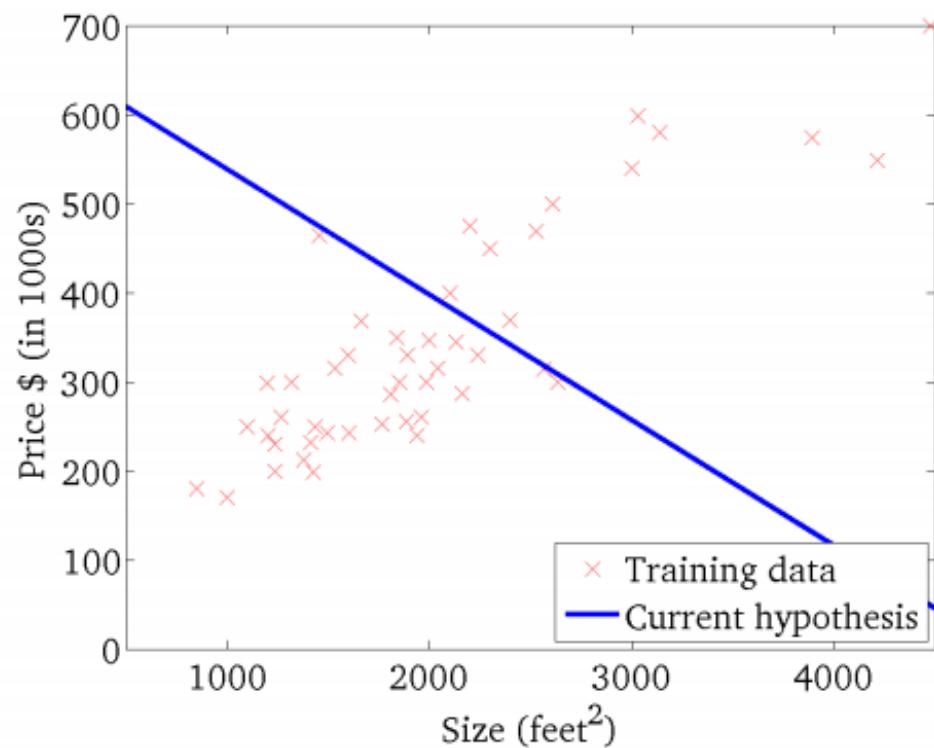
Gradient Descent Example -2



Gradient Descent Example -3

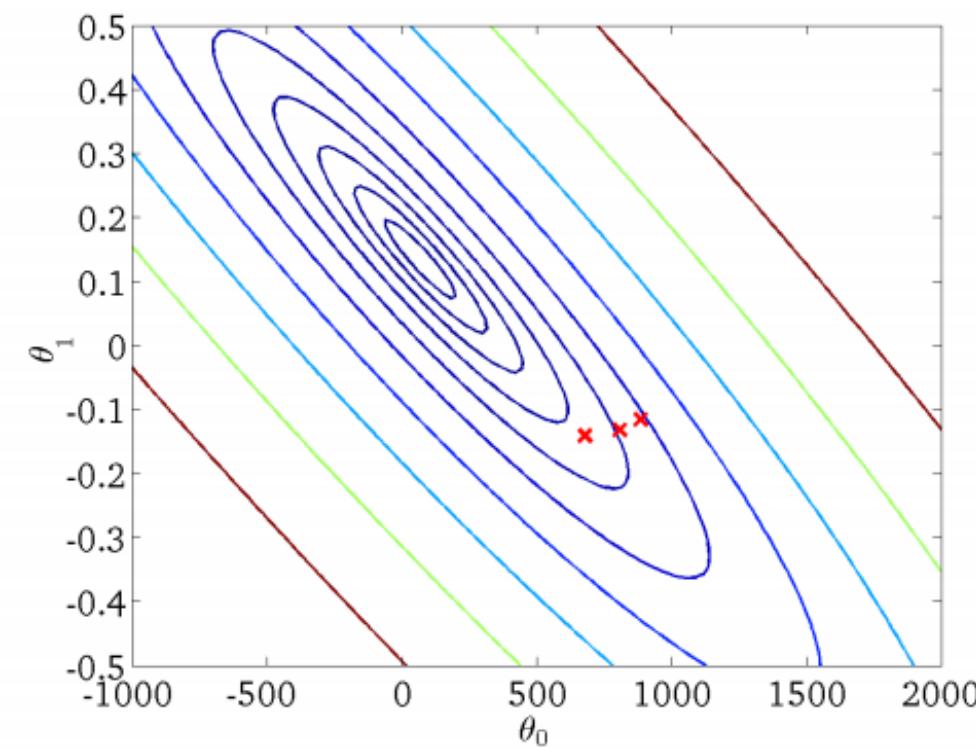
$$h_{\theta}(x)$$

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



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

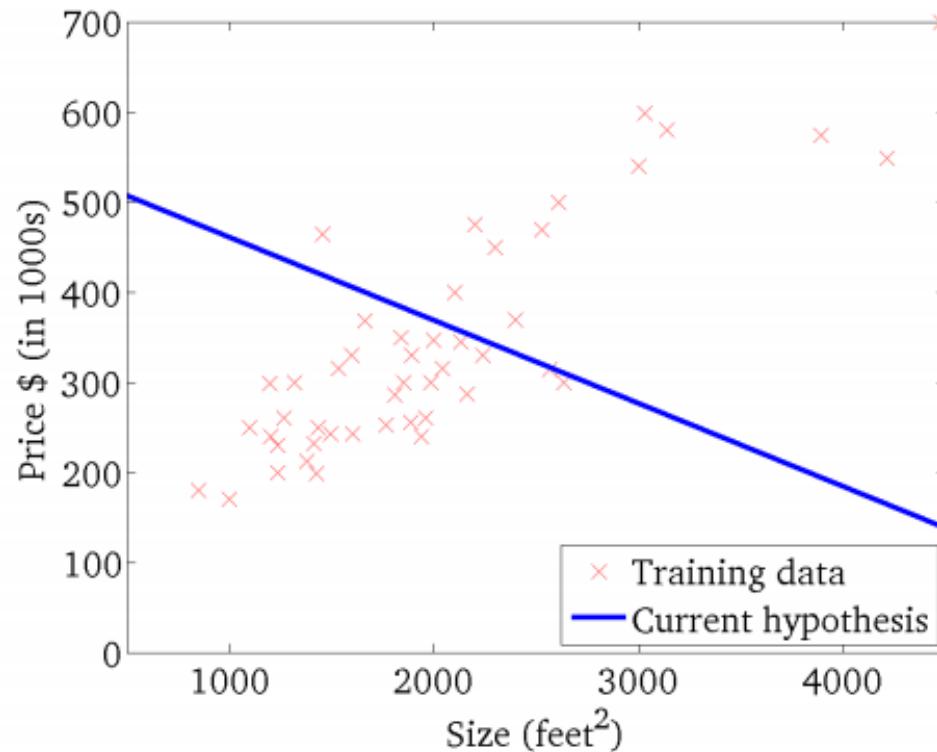
(function of the parameters θ_0, θ_1)



Gradient Descent Example -4

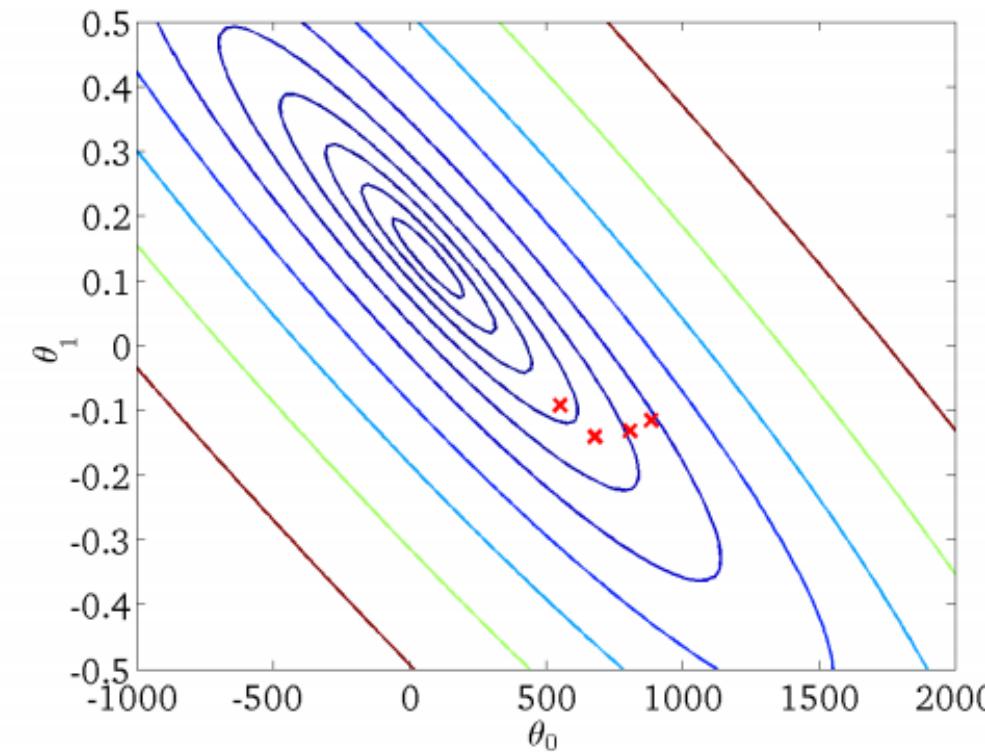
$$h_{\theta}(x)$$

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

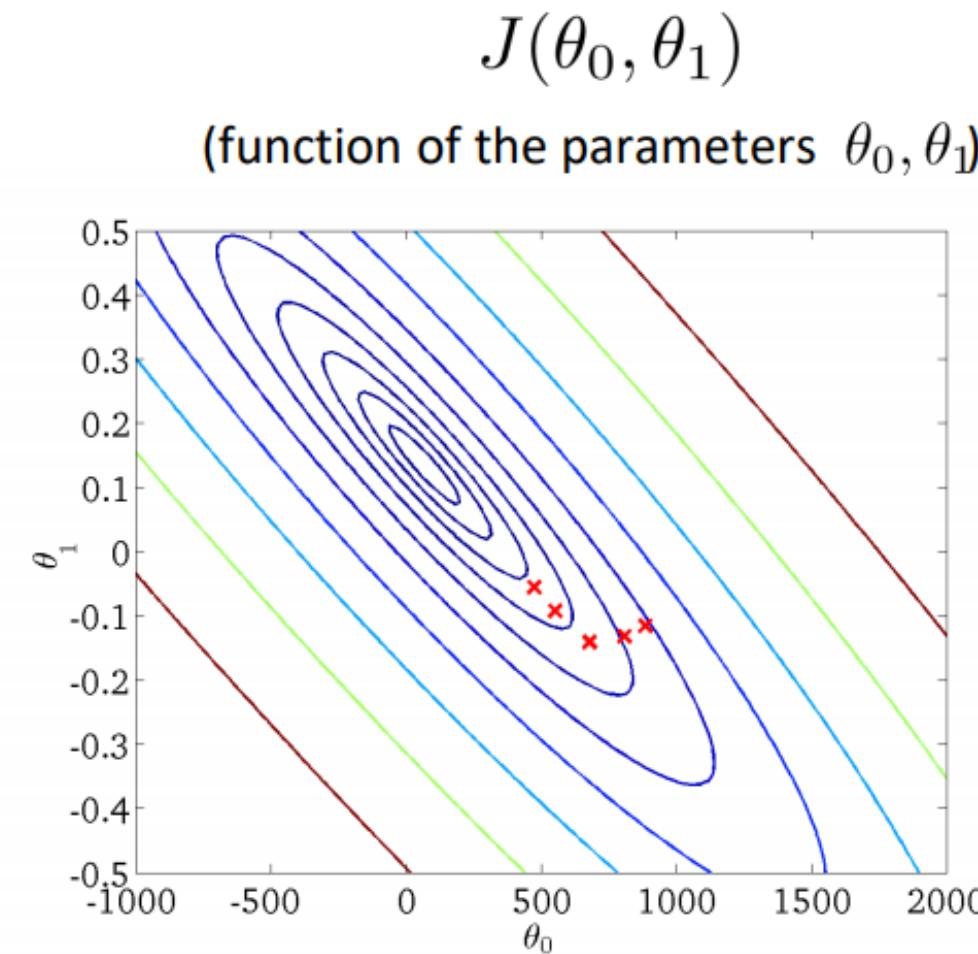
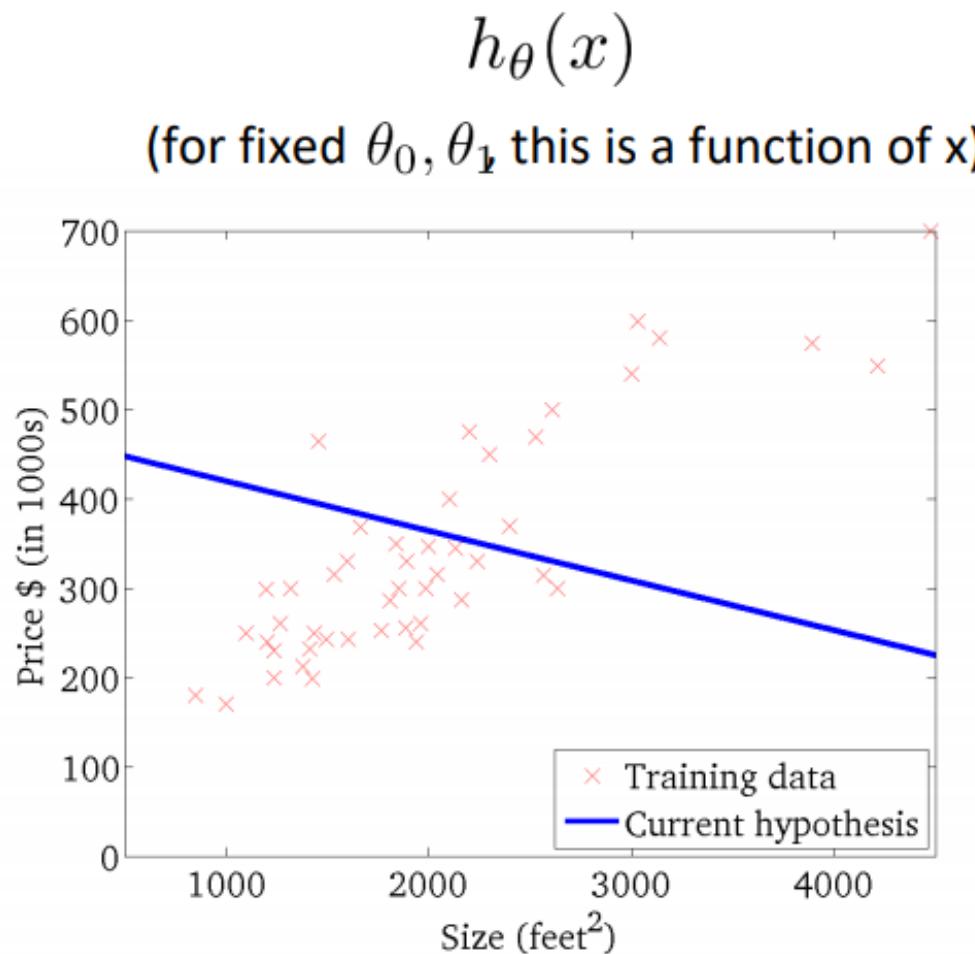


$$J(\theta_0, \theta_1)$$

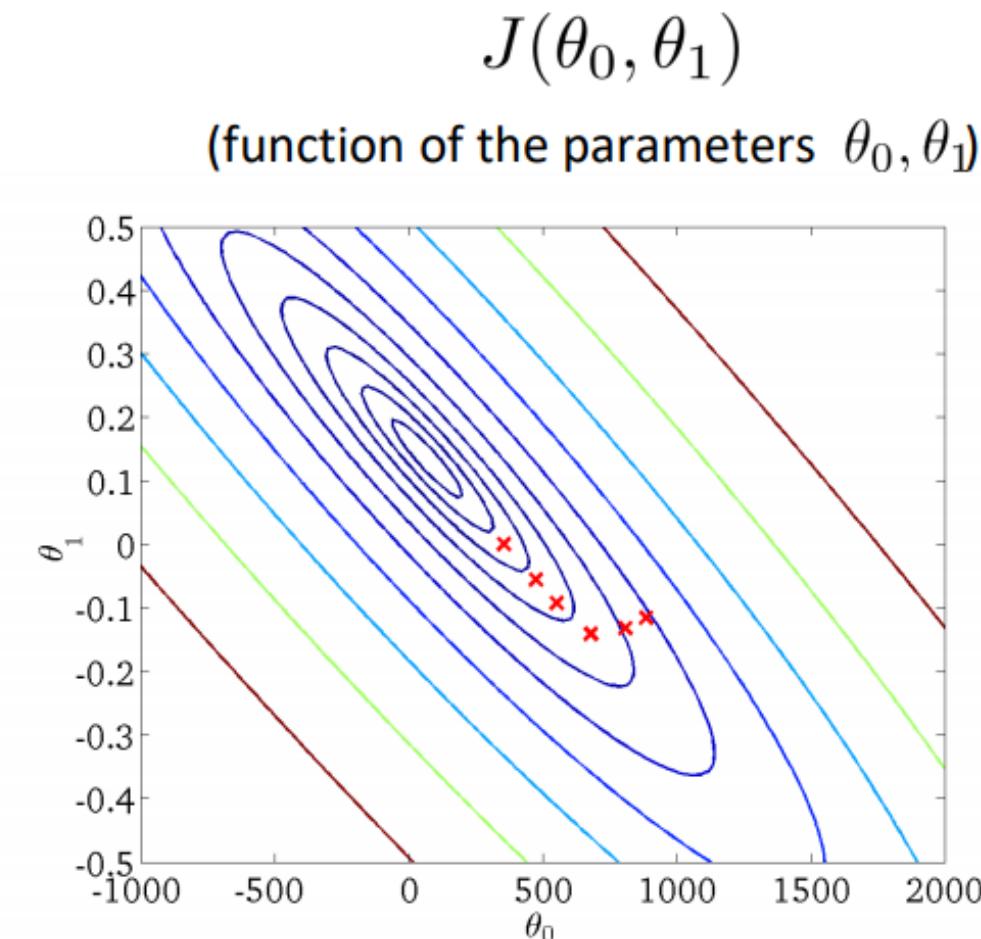
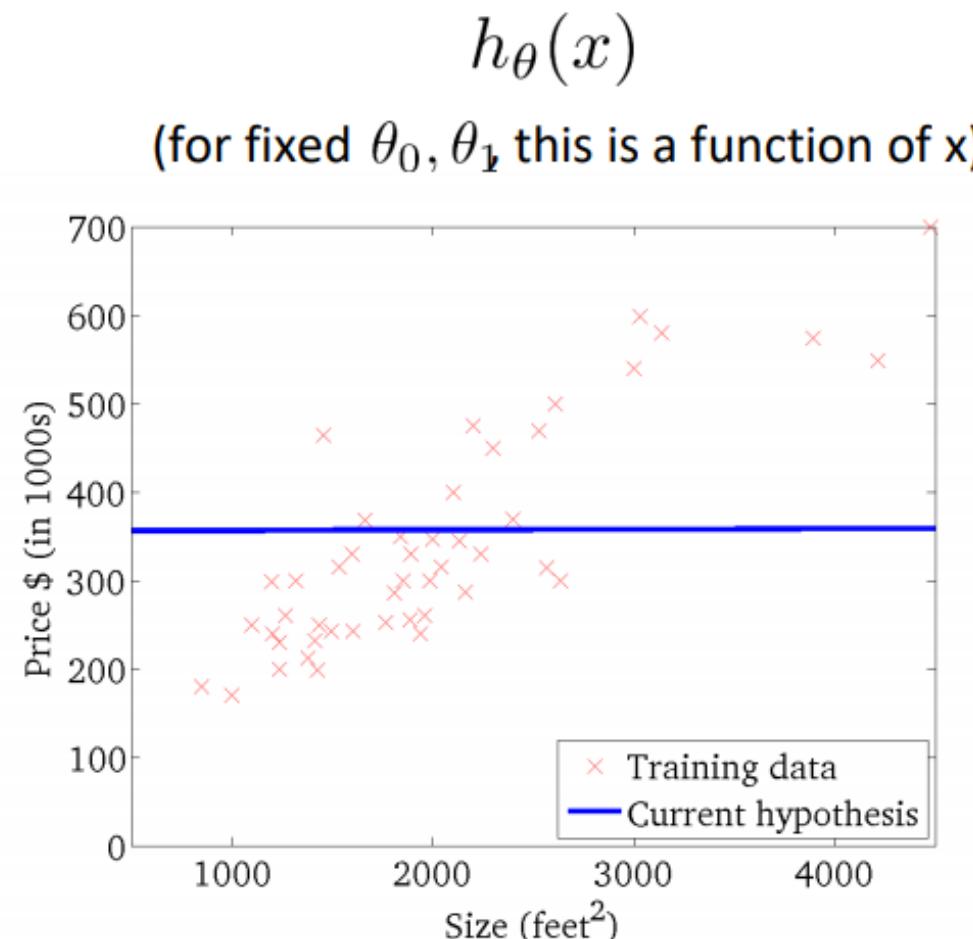
(function of the parameters θ_0, θ_1)



Gradient Descent Example -5



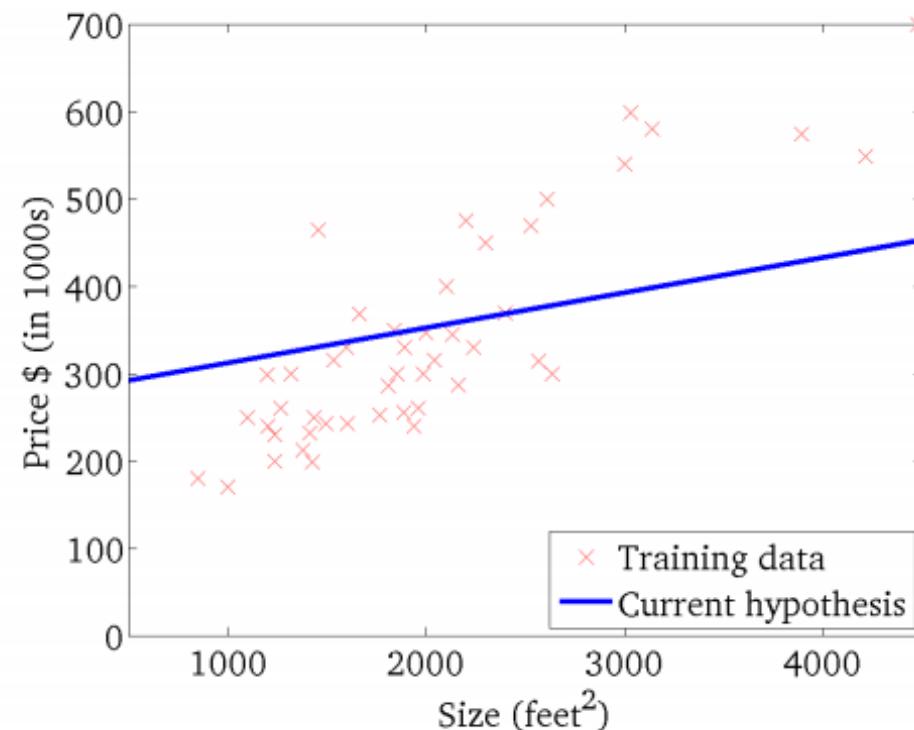
Gradient Descent Example -6



Gradient Descent Example -7

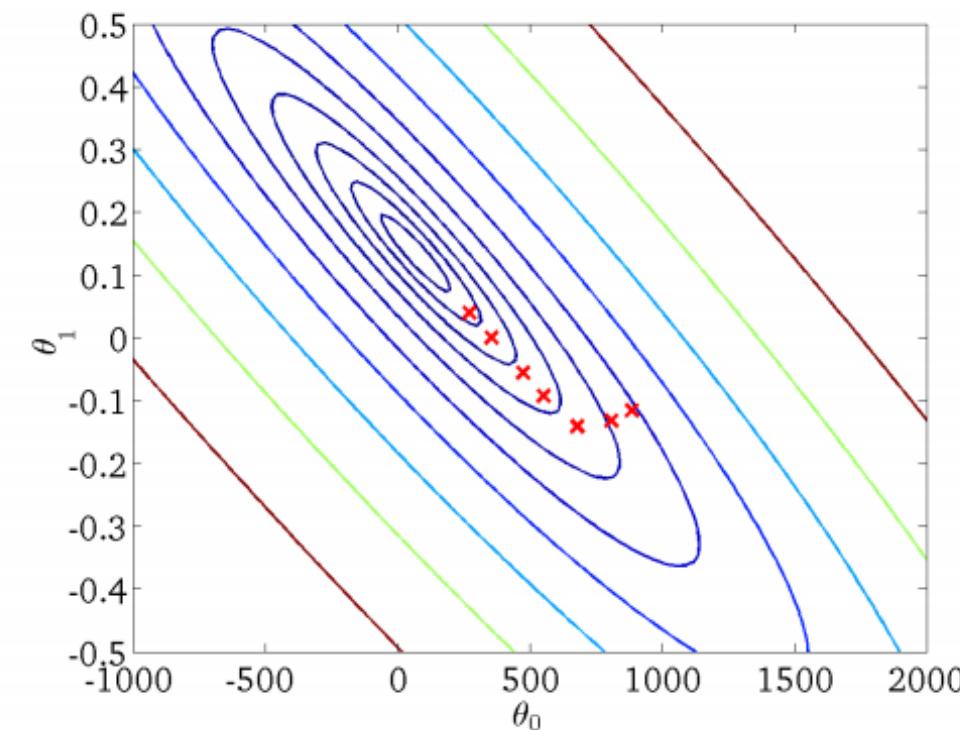
$$h_{\theta}(x)$$

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



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

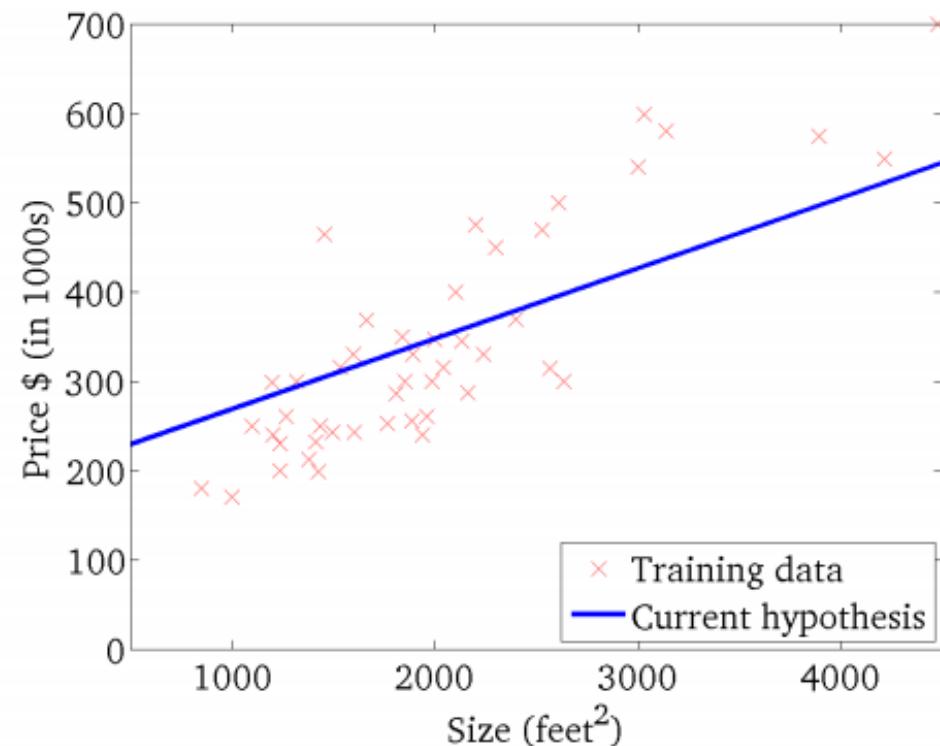
(function of the parameters θ_0, θ_1)



Gradient Descent Example -8

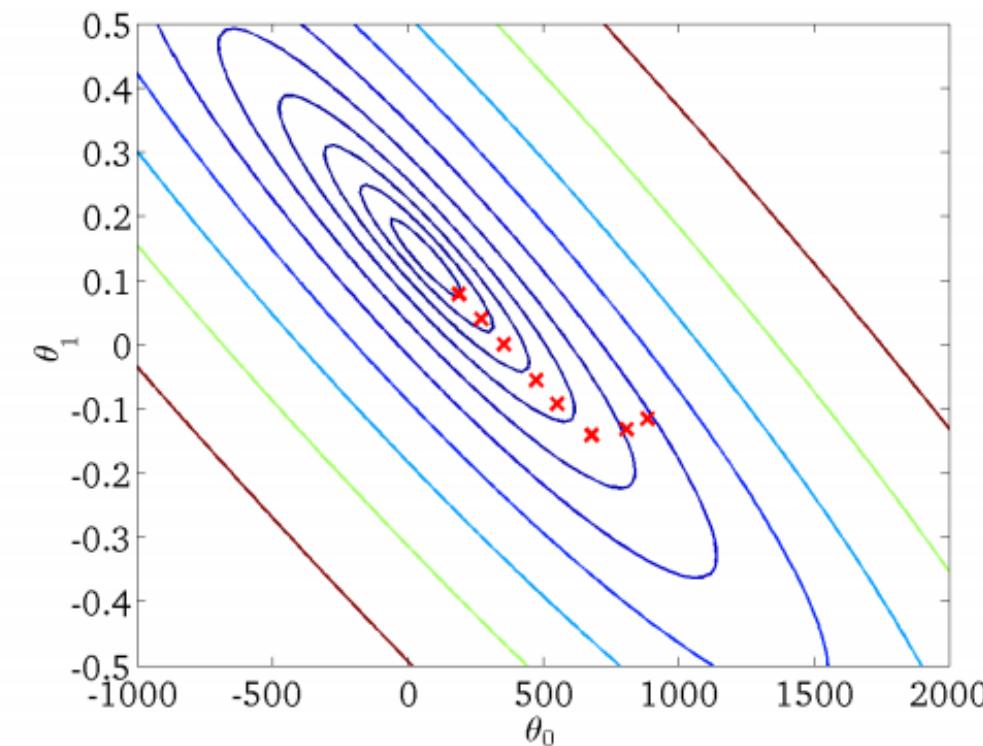
$$h_{\theta}(x)$$

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

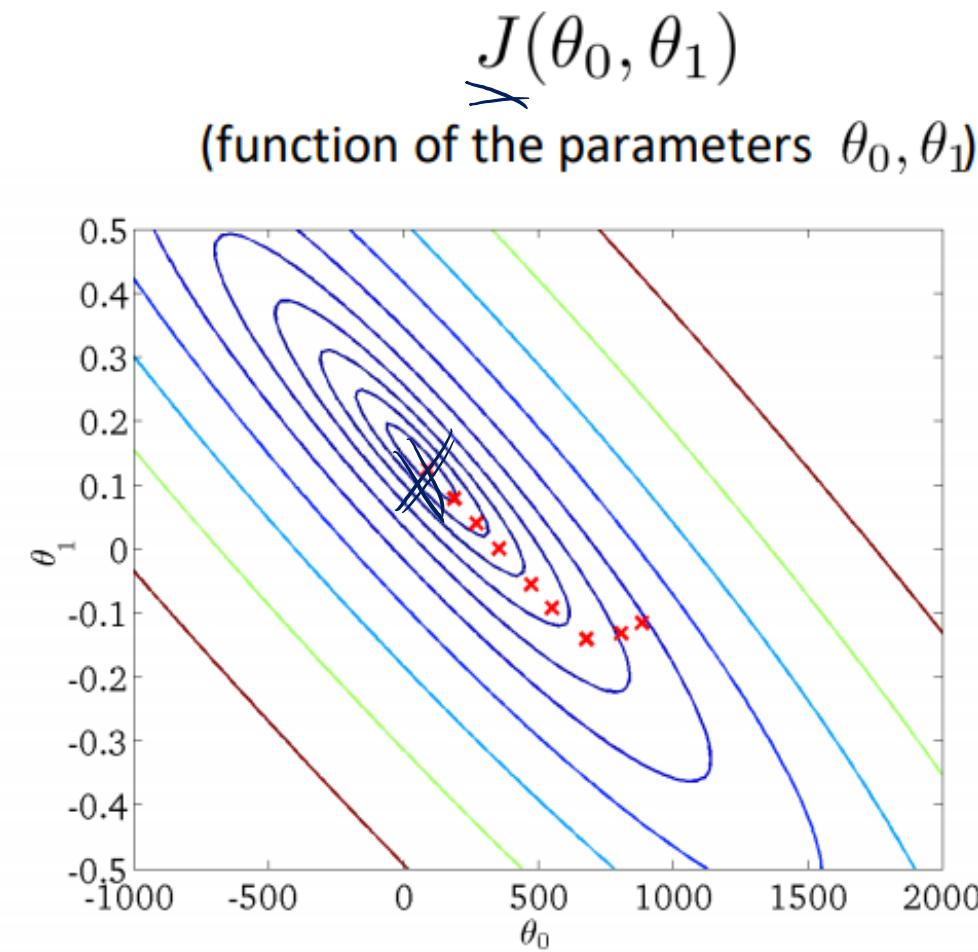
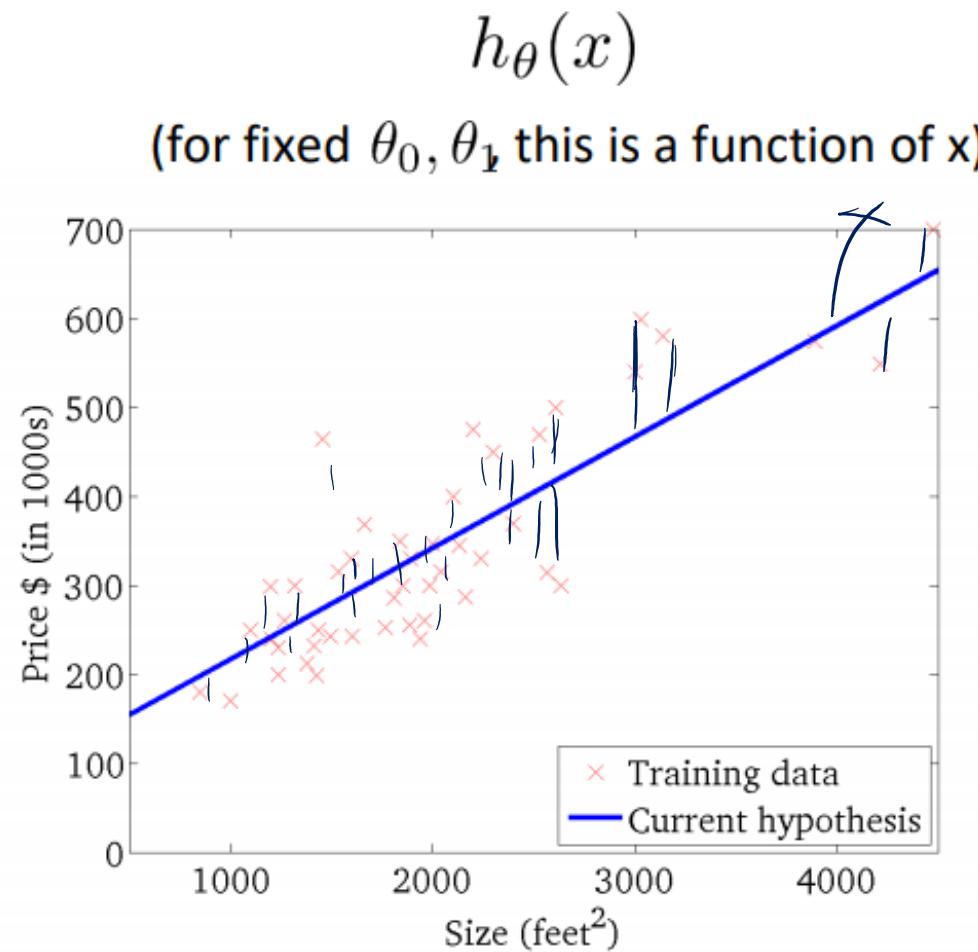


$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



Gradient Descent Example -9



Maximum Likelihood Estimation

Readings:

- Chapter 8.3 MML Textbook

Maximum Likelihood Estimation (MLE)

- A frequentist approach for estimating the parameters of a model given some observed data.
- Uses **probability distributions to model the uncertainty**.
- The general approach for using MLE is:
 1. Observe some data.
 2. Write down a **model for how we believe the data was generated**.
 3. **Set the parameters** of our model to values which **maximize the likelihood of observing what we have observed**.

Models

- A model is a formal representation of our beliefs, assumptions, and simplifications surrounding some event or process.
- **Example:** Building a model for flipping a coin
 - The coin has two faces and an edge
 - The faces have different designs
 - The coin can sit on either face or the edge
 - The weight of the coin
 - The diameter and thickness of the coin



Models Con't

➤ What assumptions can we make?

- The different designs probably cause the coin's center of mass to slightly favor one side over another.
- There's no way to measure the force or angle exerted on the coin when it's flipped.

Models Con't

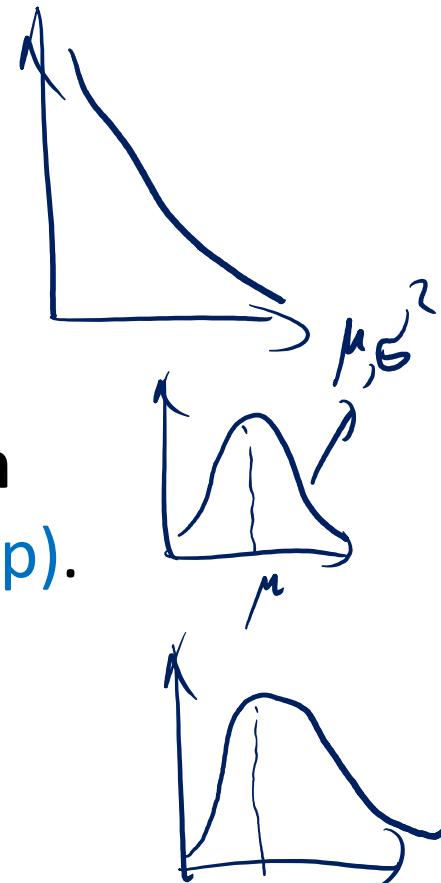
➤ First attempt at a model (without simplifications):

- The initial position of the coin is drawn from a Bernoulli distribution (i.e. flipper's preference).
- The force exerted on the coin is drawn from an exponential distribution.
- The angle in which the force is exerted is drawn from a truncated normal distribution on the interval $[-\pi, \pi]$.
- The center of mass of the coin is at some coordinate (x,y,z) in a system (center of the coin is the origin).
- The force of gravity is ... OK, I think you get the picture.



Models Con't

- The real world can be complicated. A simplified model can often do just as well or better!
- Let's make a simplified model:
 - The outcome of the flip is drawn from a **Bernoulli distribution** with the probability of heads p , and the probability of tails $(1-p)$.
- Our simplified model only has a **single parameter!**



MLE Example: Coin Flips -1

- Step 1: To start let us assume we observed the following sequence of coin flips:



- $X = \text{heads, heads, tails, heads, tails, tails, tails, heads, tails, tails}$

MLE Example: Coin Flips -2

- Step 2: Write down a model for how we believe the data was generated (we'll start with a single flip):

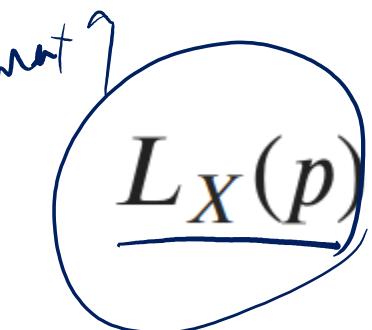
$$L_x(p) = P(x|p) = p^x(1-p)^{1-x}$$

likelihood function
(for a single flip)

Recall that we're **modeling the outcome** of a coin flip by a **Bernoulli distribution**, where the parameter p represents the probability of getting a heads.

MLE Example: Coin Flips -3

- Step 2: Write down a model for how we believe the data was generated:


$$L_X(p) = P(X|p) = \prod_{x \in X} p^x (1-p)^{1-x}$$

Since the **coin flips are iid**, we can write the likelihood of seeing a particular sequence as the **product of each individual flip**

$$\hat{p} = \operatorname{Argmax}(L(p))$$

likelihood function
(for all flips)

MLE Example: Coin Flips -3

- Step 2: We can generalize this further...

$$\underline{L(p)} = \underline{p^h} \cdot (1 - p)^{n-h}$$

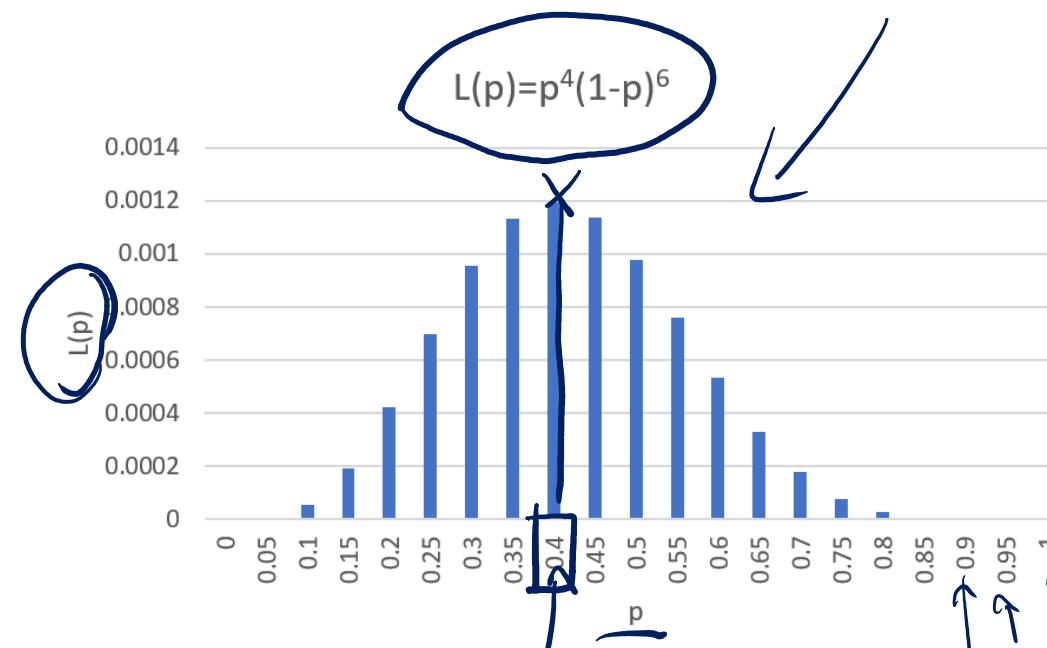
$$n=10 \\ h=4$$

where n is the number of coin flips with h the number of heads that were recorded

In our data (X):

We had observed

$h=4, n=10$



MLE Example: Coin Flips -3

- Step 2: We can generalize this further...

$$L(p) = p^h \cdot (1 - p)^{n-h}$$

where n is the number of coin flips with h the number of heads that were recorded

- Log-Likelihood Function

$$l(p) = h \cdot \log(p) + (n - h) \cdot \log(1 - p)$$

MLE Example: Coin Flips -4

- Step 3: Set the parameters of our model to values which maximize the likelihood of the parameters given the data

$$l(p) = h \cdot \log(p) + (n - h) \cdot \log(1 - p)$$

- Maximum => take derivative of function $l(p)$ with respect to p

$$l'(p) = \frac{h}{p} - \frac{n-h}{1-p} = 0$$

~~$\frac{h - hp - n + np}{p(1-p)} = 0$~~ → $np = h \rightarrow p^* = \frac{h}{n}$

Setting $l'(p)$ to 0 gives us:

$$p = h/n$$

11 : 33

Linear Regression (Maximum Likelihood Estimation)

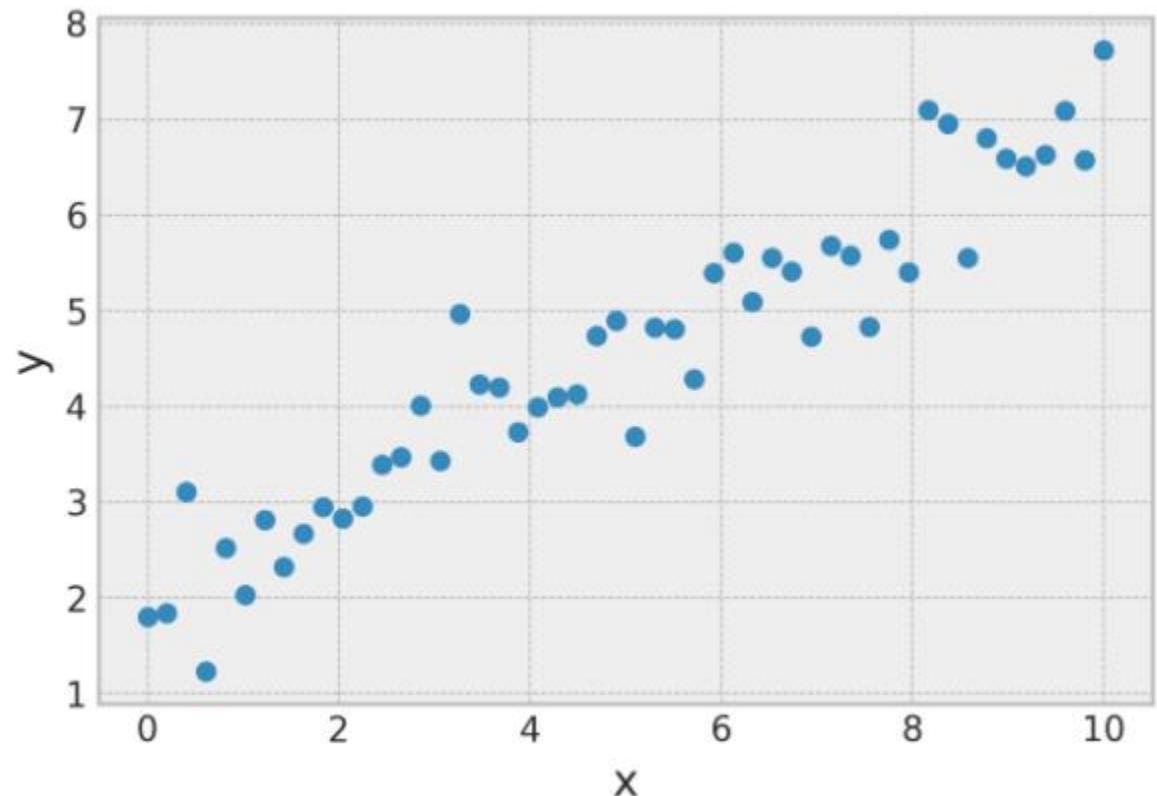
Readings:

- Chapter 9.1-2 MML Textbook

Linear Model

- We'd like to build a model of the data in order to predict new values of y given x.
- The data almost looks like a line, so let's start with that as our model.

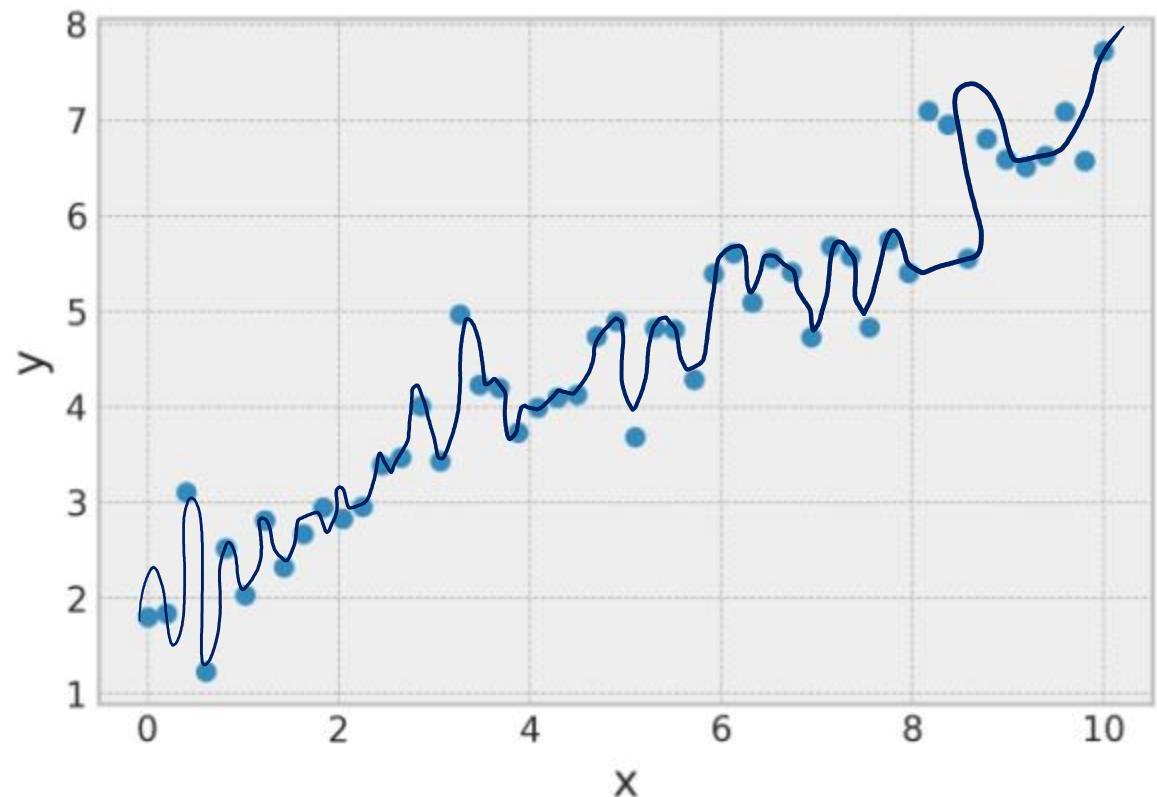
$$y = \underline{\theta_1 x + \theta_0}$$



Linear Model

- Q: How do we account for the deviations we're observing?

- A: Imagine we're using a sensor to collect this data. Most sensors have some amount of error in their measurements. We can think of the deviations from our model as being caused by an error prone sensor.



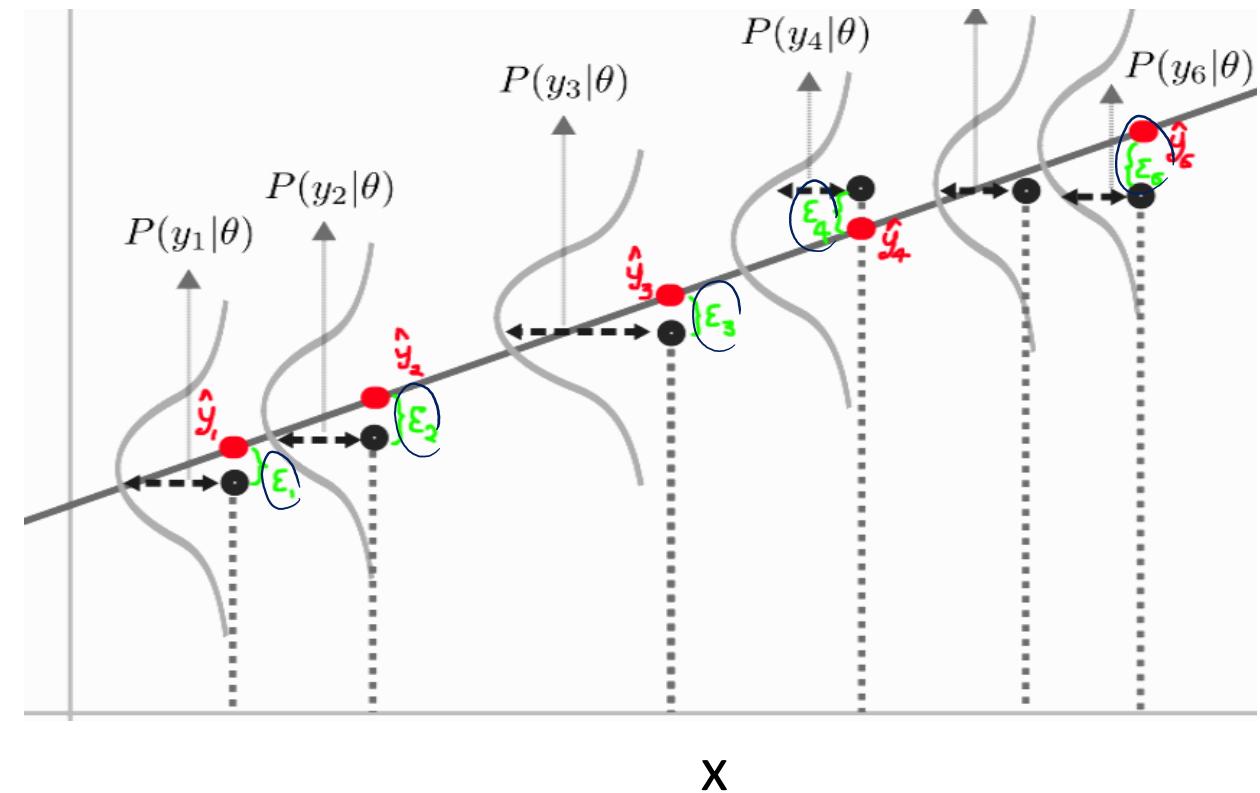
Linear Model

- It's very common to model the error as being drawn from a Gaussian distribution with mean zero and variance σ^2 .

$$\epsilon \sim N(0, \sigma^2)$$
$$y = \theta_1 x + \theta_0 + \epsilon$$

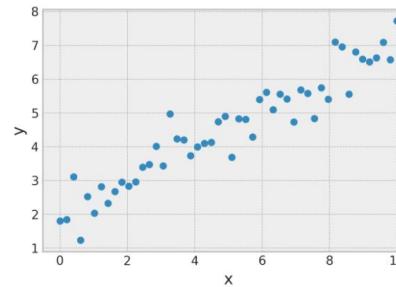
Signal y noise

$$y \sim N(\theta_1 x + \theta_0, \sigma^2)$$



MLE Example: Linear Regression -1

- Step 1: Obtain some sample data:



$$\{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$$

- Step 2: Write down a model for how we believe the data was generated (we'll start with a single sample):

$$f(y|x, \theta_0, \theta_1, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-(\theta_1 x + \theta_0))^2}{2\sigma^2}}$$



likelihood function
(for a single point)

This time we are **modeling the outcome with a Gaussian distribution.**

MLE Example: Linear Regression -2

- Step 2: Write down a model for how we believe the data was generated:

$$L(\theta_0, \theta_1, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \prod_{n=1}^N e^{\frac{-(y_n - (\theta_1 x_n + \theta_0))^2}{2\sigma^2}}$$

likelihood function
(for all points)

since each point is iid, we can write the likelihood function with respect to all the observed points as a product of each point.

MLE Example: Linear Regression -3

➤ Step 2: Rewrite in terms of log-likelihood

$$\begin{aligned}\mathcal{L}(\theta_0, \theta_1, \sigma^2) &= \log\left[\frac{1}{\sqrt{2\pi\sigma^2}} \prod_{n=1}^N e^{-\frac{(y_n - (\theta_1 x_n + \theta_0))^2}{2\sigma^2}}\right] \\ &= \log\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) + \sum_{n=1}^N \frac{-(y_n - (\theta_1 x_n + \theta_0))^2}{2\sigma^2} \\ &= \cancel{\log(1)} - \log(\sqrt{2\pi\sigma^2}) - \frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - (\theta_1 x_n + \theta_0))^2 \\ &= -\log(\sqrt{2\pi\sigma^2}) \left[-\frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - (\theta_1 x_n + \theta_0))^2 \right]\end{aligned}$$

Recall $\log(ab) = \underline{\log(a)} + \underline{\log(b)}$

MLE Example: Linear Regression -4

- Step 2: Rewrite in terms of log-likelihood

$$\mathcal{L}(\theta_0, \theta_1, \sigma^2) = -\underbrace{\log(\sqrt{2\pi\sigma^2})}_{\text{constant}} - \frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - (\theta_1 x_n + \theta_0))^2$$

- To clean things up a bit more, let's write the output of our line as a single value:

$$\hat{y} = \theta_1 x + \theta_0$$

- Now our log-likelihood can be written as:

$$\mathcal{L}(\theta_0, \theta_1, \sigma^2) = -\log(\sqrt{2\pi\sigma^2}) - \frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \hat{y}_n)^2$$

multiply by -1 to make it a negative log-likelihood

Sum of Squared Errors

- Step 3: Parameters which maximize the log-likelihood are the same as the ones that minimize the negative log-likelihood.

$$\mathcal{L}(\theta_0, \theta_1, \sigma^2) = \underbrace{\log(\sqrt{2\pi\sigma^2})}_{\text{constant}} + \frac{1}{2\sigma^2} \sum_{n=1}^N (y - \hat{y}_n)^2$$

- Removing constants which don't include our θ s won't alter the solution. We can simplify what we're trying to minimize:

$$\sum (y - \hat{y})^2$$

Sum of
Squared
Errors

Solving for Parameters

- The maximum likelihood estimates for our slope and intercept (θ parameters) can be found by **minimizing the sum of squared errors.**

This is the same as our empirical risk minimization where we assumed a sum of squared error loss function!

$$\begin{aligned}\underline{\mathcal{L}(\theta)} &= \sum_{n=1}^N (y_n - \hat{y}_n)^2 \\ &= \sum_{n=1}^N (y_n - (\theta_1 x_n + \theta_0))^2 \\ &= (\mathbf{y} - \mathbf{X}\boldsymbol{\theta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\theta}) \\ &= \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|^2\end{aligned}$$

vector notation

Solving for the Parameters (θ)

- We start by taking the **partial derivative** with respect to **parameters θ** :

$$\begin{aligned}\frac{d\mathcal{L}}{d\theta} &= \frac{d}{d\theta} \left((\mathbf{y} - \mathbf{X}\theta)^\top (\mathbf{y} - \mathbf{X}\theta) \right) \\ &= \frac{d}{d\theta} \left(\mathbf{y}^\top \mathbf{y} - 2\mathbf{y}^\top \mathbf{X}\theta + \theta^\top \mathbf{X}^\top \mathbf{X}\theta \right)\end{aligned}$$

$$= \underbrace{-\mathbf{y}^\top \mathbf{X} + \theta^\top \mathbf{X}^\top \mathbf{X}}$$

Q: How do we find our parameters (θ)?

Solving for the Parameters (θ)

➤ At this point we have two options for finding the parameters θ :

$$\frac{d\mathcal{L}}{d\theta} = -y^\top X + \theta^\top X^\top X$$

(1) iterative solution for θ :

$$\theta_i \leftarrow \theta_i - \alpha \frac{\partial \mathcal{L}(\theta_i)}{\partial \theta_i}$$

(2) direct solution for θ : ✓

$$\frac{d\mathcal{L}}{d\theta} = 0$$

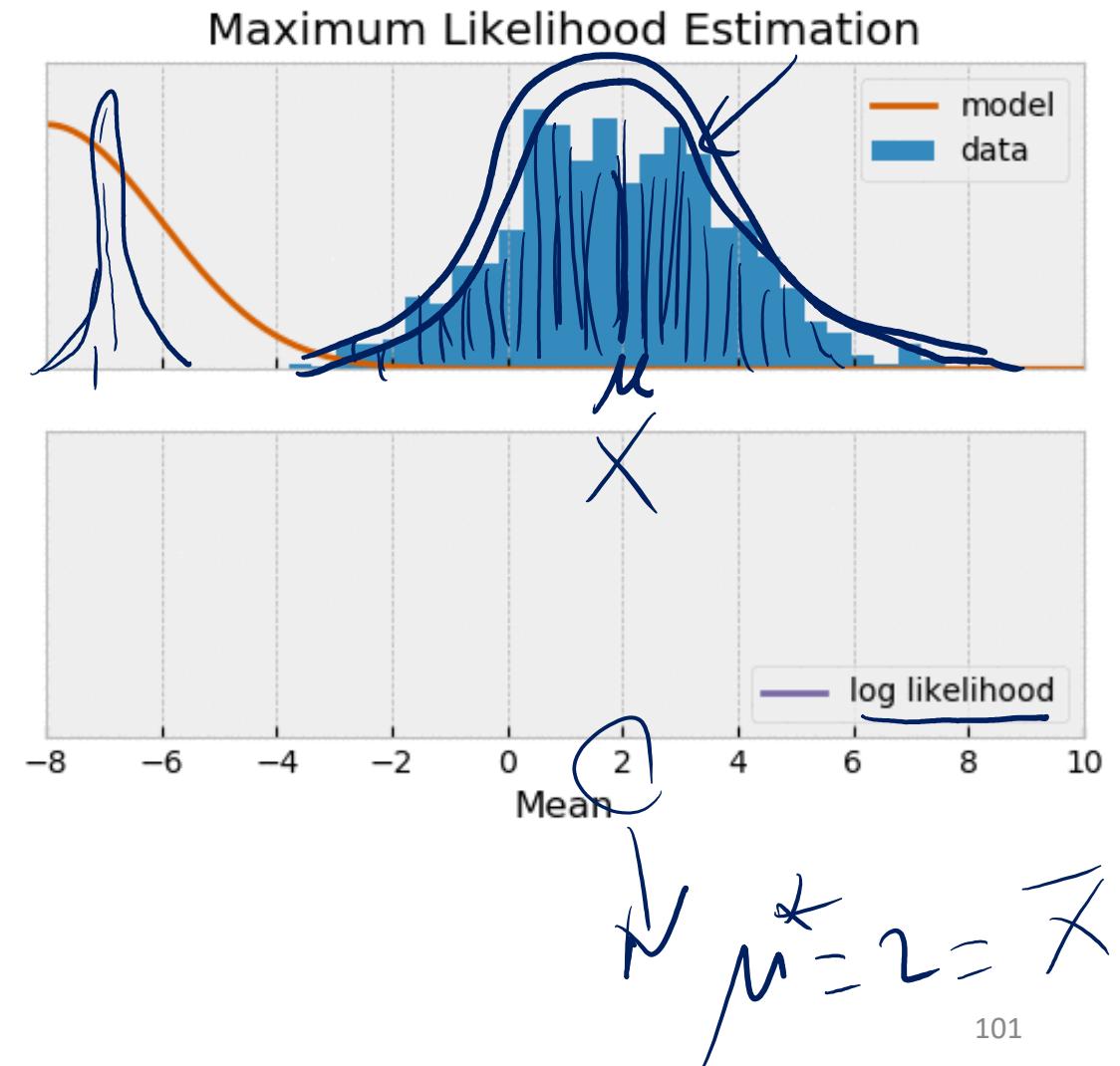
$$\theta_{ML}^\top X^\top X = y^\top X$$

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

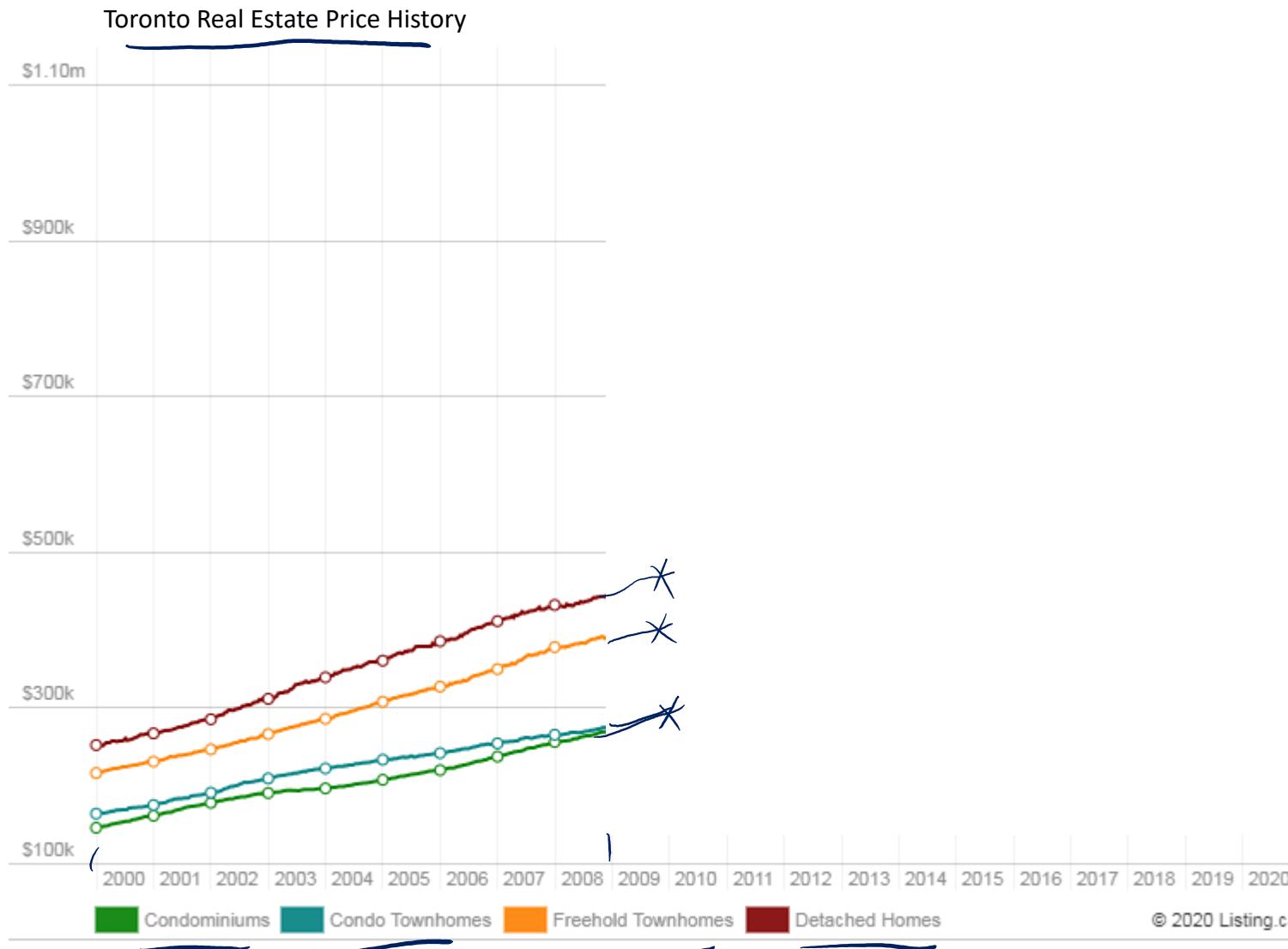
$$\theta_{ML} = (X^\top X)^{-1} X^\top y$$

Summary: Maximum Likelihood

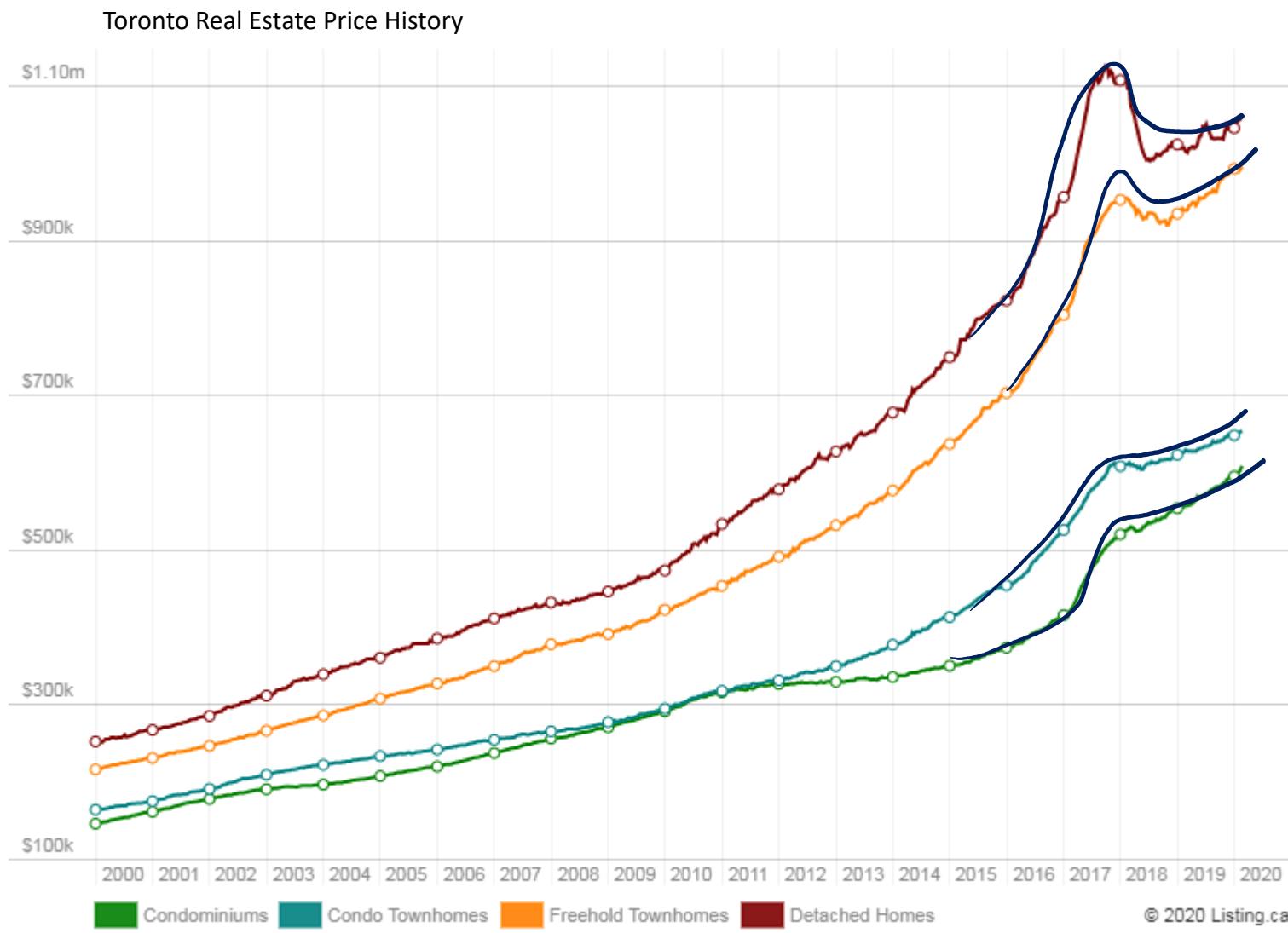
- Imagine we have some data generated from a Gaussian distribution with a known variance, but we don't know the mean.
- You can think of MLE as taking the Gaussian, sliding it over all possible means, and choosing the mean which causes the model to fit the data best.



Nonlinear Regression



Nonlinear Regression



LR Capturing
non-linear effect

Next Time

- Week 9 Lab session: Tutorial 4 on Regression
- Week 10 Lecture – Nonlinear Regression
 - Polynomial Regression
 - Optimization and Convexity
 - Regression with Regularization
 - Classification
 - Neural Networks