Today: Outline

Piazza Walk Through

- Machine Learning Review 1
 - Cost Functions
 - Regression
 - Classification

 Reminder: Please email me ASAP if you have not received a Piazza enrollment email.



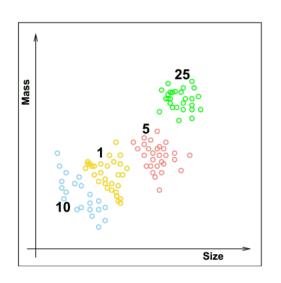
Supervised Learning

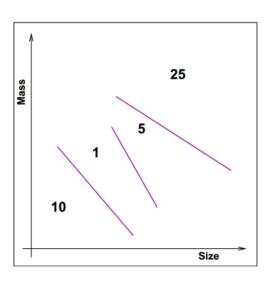
Task Definitions: Classification and Regression Tasks

Example of Supervised Learning:

Classification: recognize coins







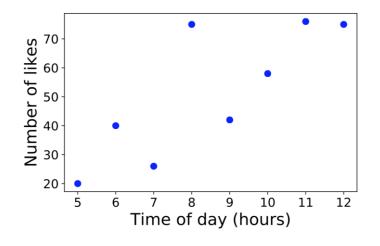
- Given training set consisting of coin denomination (penny, nickel, dime, quarter), mass and size
- Learn to predict denomination
- What is input? Output?

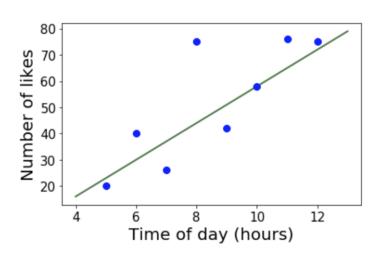
Saenko

Example of Supervised Learning:

Linear Regression: predict number of likes

- When the label is a real number
- Training a model to find a relationship between input and output values
- Learning a line of best fit



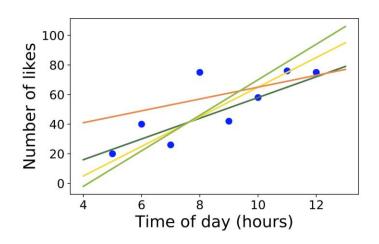


Linear Regression: Model Parameters

• Learning a best fit line means learning parameters (or weights) for our model.

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

 θ_i 's: Parameters

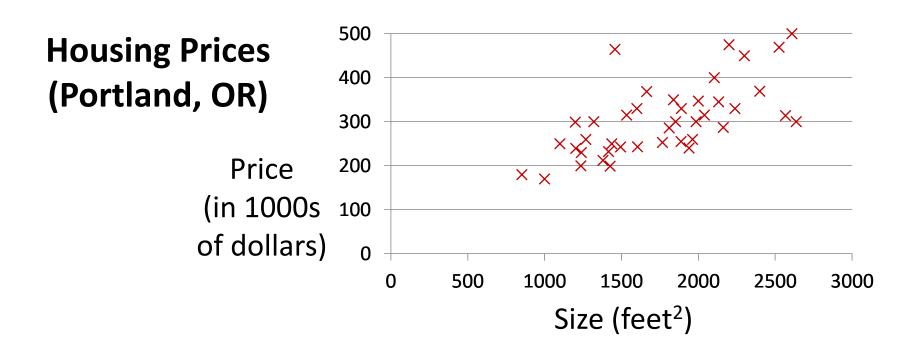




Supervised Learning

Regression

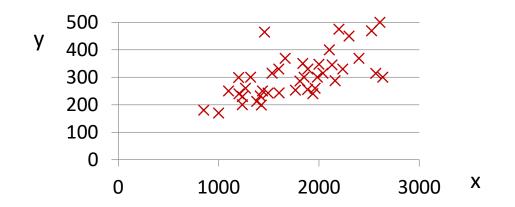
Example: house price prediction



Supervised Learning

What should the learner be??

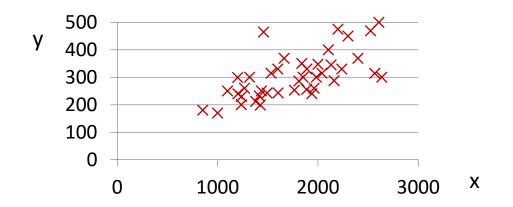




Hypothesis h

h: a function parametrized by ϑ

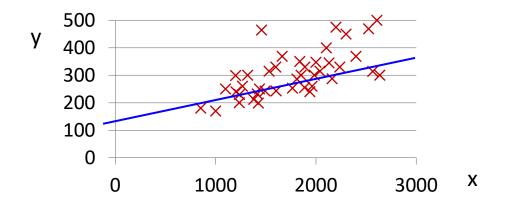
input x
$$\longrightarrow \left[h_{\theta} \right] \longrightarrow \text{output y}$$



How to learn ϑ ?

Given: Training Set $\{x^i, y^i\}$ But what if $y \neq y^i$??

input
$$x^i \longrightarrow h_{\theta} \longrightarrow \text{output } y$$



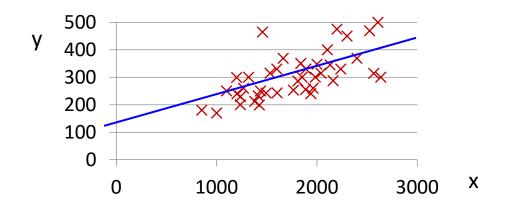
Cost function

Given: Training Set $\{x^i, y^i\}$

Cost function Cost(y, yⁱ)

learning == minimizing cost

input
$$x^i \longrightarrow h_\theta \longrightarrow \text{output } y$$



Supervised Learning

Given: Training Set $\{x^i, y^i\}$

Cost function Cost(y, yⁱ)

learning == minimizing cost

Learn θ^* : min Cost($h_{\theta}(x^i)$, y^i)

want: input $x^i \longrightarrow h_{\theta}^* \longrightarrow$ output y

Training set

Training set:

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

Notation:

```
m = Number of training examples

x^{(i)} = "input" variable / features

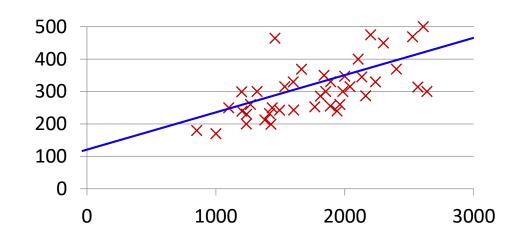
y^{(i)} = "output" variable / "target" variable
```

What should *h* be?

Linear hypothesis:

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

 $heta_i$'s: Parameters



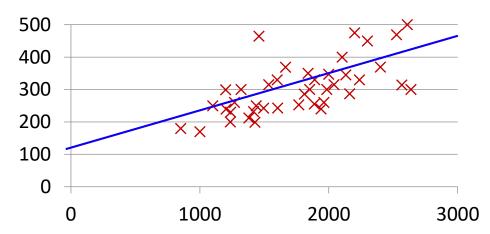
min Cost(
$$h_{\theta}$$
, {xⁱ, yⁱ}) θ

What's a good cost function for this problem?

Hypothesis:

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

 θ_i 's: Parameters



How about "Sum of squared differences"

Cost Function:

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

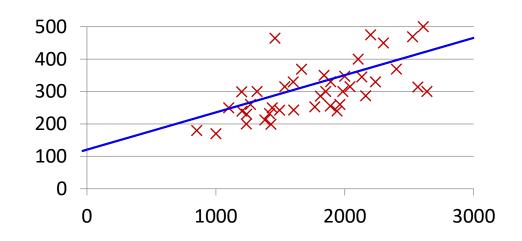
Goal: minimize
$$J(\theta_0, \theta_1)$$

2-dimensional θ

Hypothesis:

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

 θ_i 's: Parameters

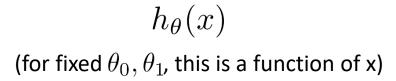


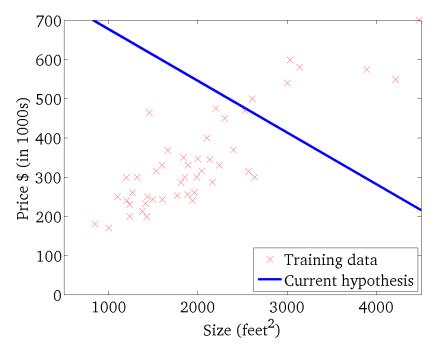
Cost Function:

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

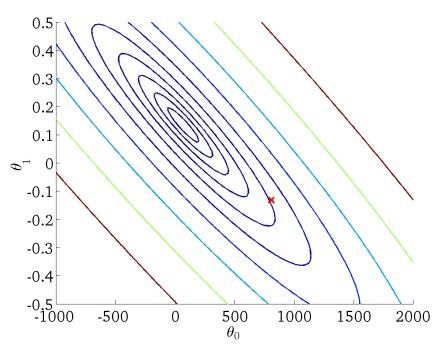
SSD = sum of squared differences, also SSE = sum of squared errors

Plotting cost for 2-dimensional θ

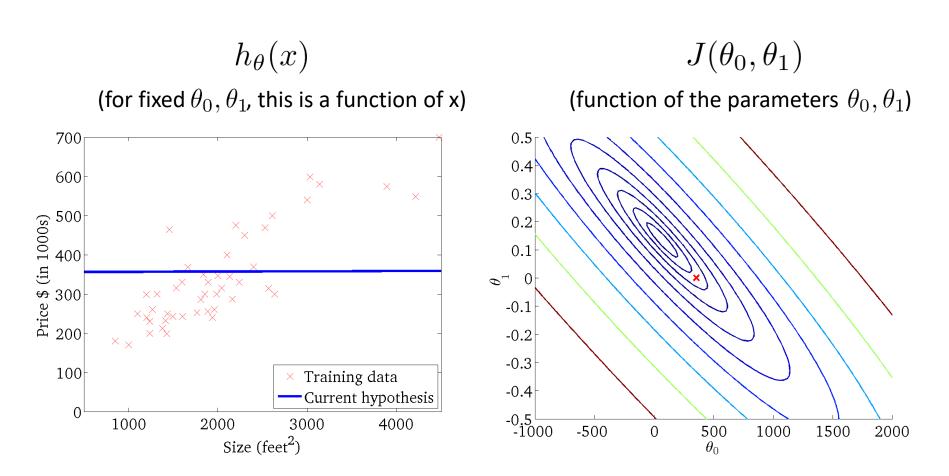




 $J(heta_0, heta_1)$ (function of the parameters $heta_0, heta_1$)

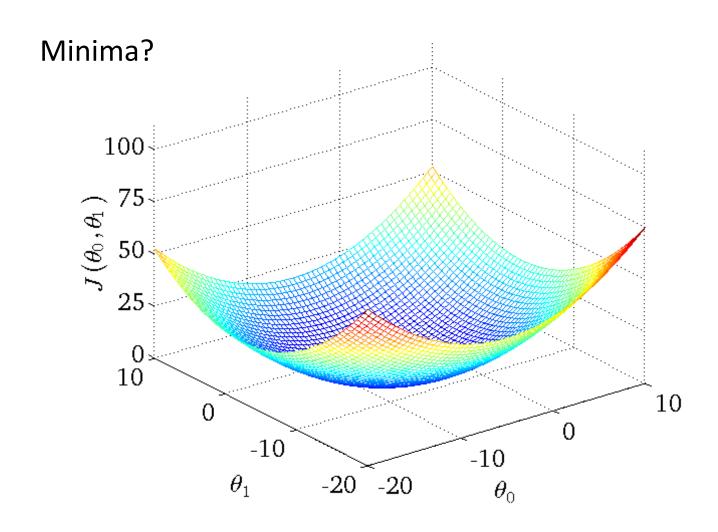


Plotting cost for 2-dimensional θ

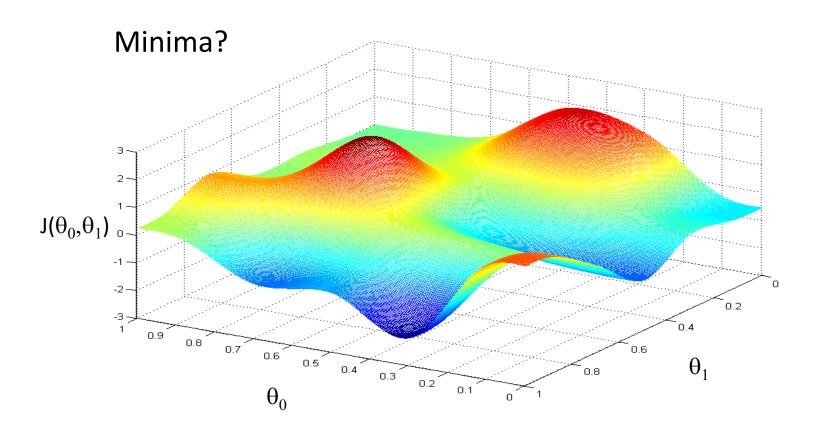


Note, squared loss cost is convex in parameters

SSD cost function is convex



Non-convex cost function



Multivariate Linear Regression

Hypothesis:

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

For convenience of notation, define $x_0 = 1$.

 θ_i 's: Parameters

Cost Function:

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

Goal: minimize
$$J(\theta_0, \theta_1, \dots, \theta_n)$$



Supervised Learning: Classification

Classification

Binary Classification

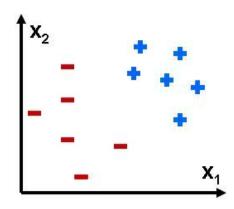
 $y \in \{0,1\}$ 0: "Negative Class" (e.g., No COVID19) 1: "Positive Class" (e.g., COVID19)

Emotion: Smile/ No Smile

Email: Spam / Not Spam?

Video: Viral / Not Viral?

Tumor: Malignant / Benign?



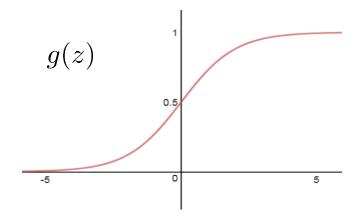
Logistic Regression

$$0 \le h_{\theta}(x) \le 1$$

map to (0, 1) with "sigmoid" function

$$g(z) = \frac{1}{1 + e^{-z}}$$

$$h_{\theta}(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}}$$



$$h_{\theta}(x) = p(y = 1|x)$$
 "probability of class 1 given input"

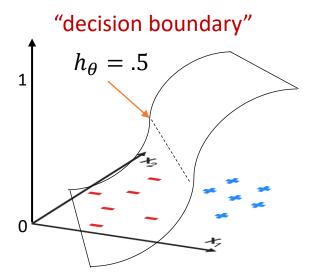
Logistic Regression

Hypothesis:

$$h_{\theta}(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}}$$

Predict "
$$y = 1$$
" if $h_{\theta}(x) \ge 0.5$

Predict "
$$y = 0$$
" if $h_{\theta}(x) < 0.5$



Logistic Regression Cost Function

Hypothesis:

$$h_{\theta}(x) = g(\theta^T x) = \frac{1}{1 + e^{-\theta^T x}}$$

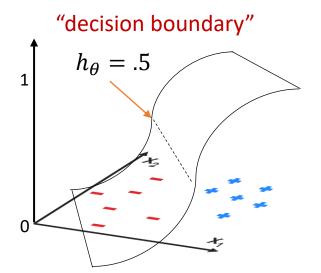
 θ : parameters

$$D = (x^{(i)}, y^{(i)})$$
: data

Cost Function: cross entropy

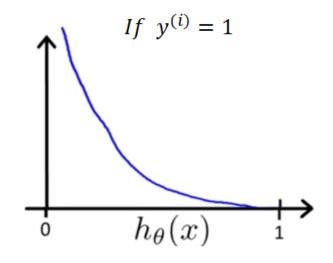
$$J(\theta) = \frac{1}{m} \sum_{i=1}^{m} \operatorname{Cost}(h_{\theta}(x^{(i)}), y^{(i)})$$

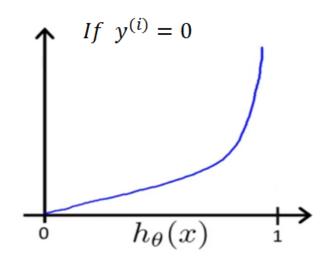
Goal: minimize cost $\min_{\theta} J(\theta)$



Logistic Regression Cost Function

Cost
$$(h_{\theta}(x^{(i)}), y^{(i)}) = \begin{cases} -\log(h_{\theta}(x^{(i)})) & \text{if } y^{(i)} = 1\\ -\log(1 - h_{\theta}(x^{(i)})) & \text{if } y^{(i)} = 0 \end{cases}$$





Cost = 0 if
$$y^{(i)} = 1$$
, $h_{\theta}(x^{(i)}) = 1$

Similarly desirable behavior

But as
$$h_{\theta}(x^{(i)}) \to 0$$

 $\text{Cost} \to \infty$

Cross Entropy Cost

$$\operatorname{Cost}(h_{\theta}(x^{(i)}), y^{(i)}) = \begin{cases} -\log(h_{\theta}(x^{(i)})) & \text{if } y^{(i)} = 1\\ -\log(1 - h_{\theta}(x^{(i)})) & \text{if } y^{(i)} = 0 \end{cases}$$

Can be written more compactly as:

$$\operatorname{Cost}(h_{\theta}(x^{(i)}), y^{(i)}) = -y^{(i)} \log h_{\theta}(x^{(i)}) - (1 - y^{(i)}) \log (1 - h_{\theta}(x^{(i)}))$$

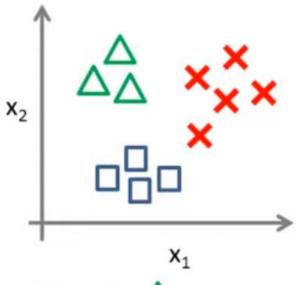
Cross Entropy Cost

Cross entropy compares distribution q to reference p

$$H(p,q) = -\sum_x p(x) \, \log q(x).$$

• Here q is predicted probability of y=1 given x, reference distribution is $p=y^{(i)}$, which is either 1 or 0

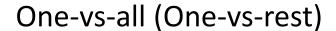
$$J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^{m} y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\theta}(x^{(i)})) \right]$$



Class 1: 🛆

Class 2:

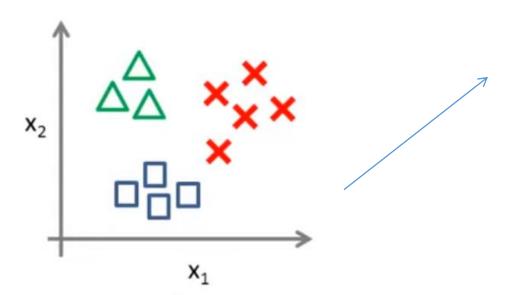
Class 3: X

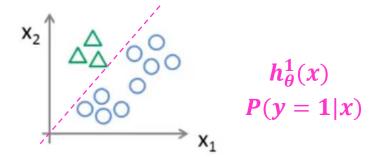


Class 1: \triangle

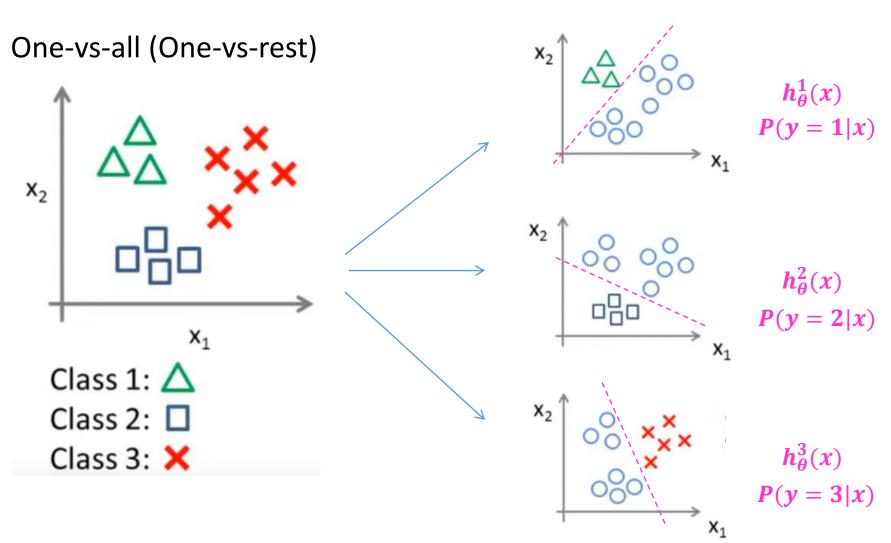
Class 2:

Class 3: X





Andrew Ng



• Trained a logistic regression classifier $h_{\theta}^{i}(x)$ for each class i to predict the probability that y=i.

 On a new input x, to make a prediction, pick the class i that maximizes:

$$\max_{i} h_{\theta}^{i}(x)$$

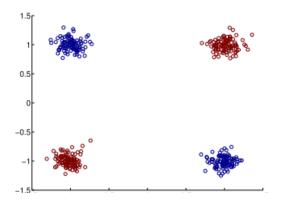


Supervised Learning

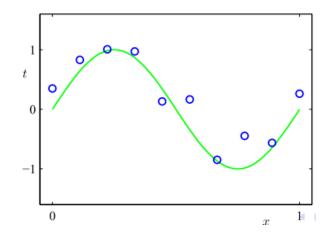
Non-linear Features

What to do if data is nonlinear?

Example of nonlinear classification



Example of nonlinear regression



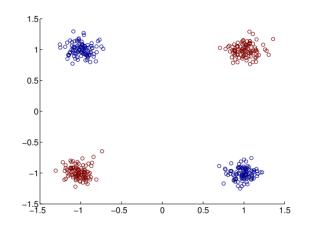
Nonlinear basis functions

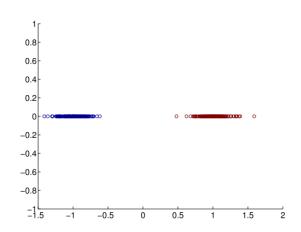
Transform the input/feature

$$\phi(x): x \in R^2 \to z = x_1 \cdot x_2$$

How do we select such features?

Transformed training data: linearly separable!







Supervised Learning

Gradient Descent

Gradient Descent Algorithm

simultaneously for all

 $j = 0, \ldots, n$

Set
$$\theta = 0$$
 or random

Repeat {

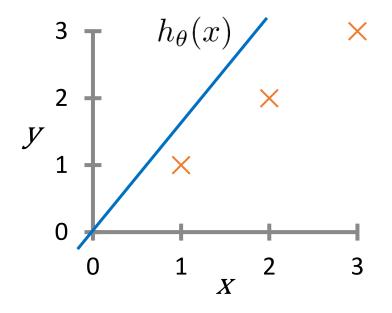
$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

} until convergence

Gradient Descent: Intuition

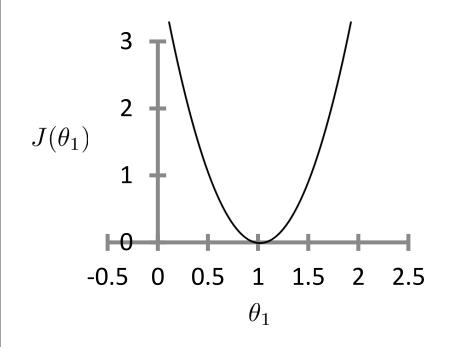
 $h_{\theta}(x)$

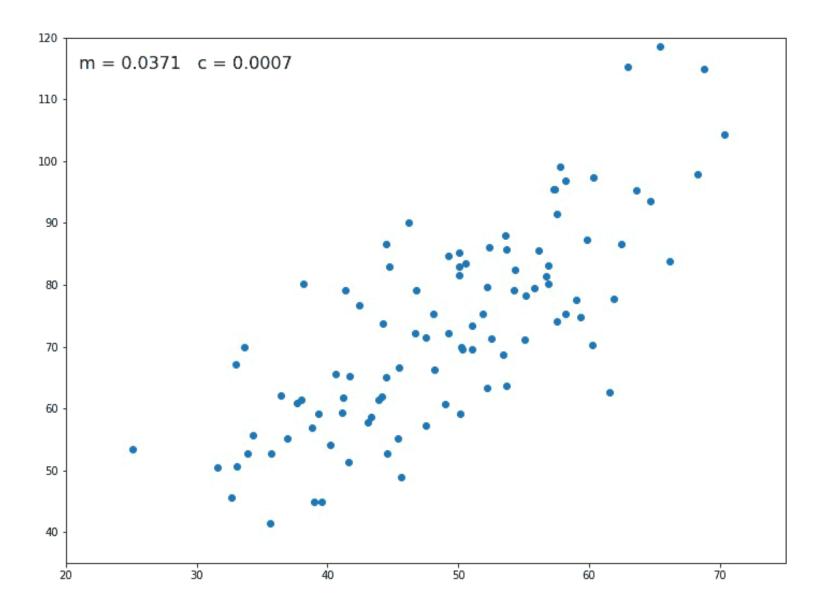
(for fixed θ_1 , this is a function of x)



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

 $J(heta_1)$ (function of the parameter $heta_1$)





Gradient descent illustration (credit: https://towardsdatascience.com/

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

$$\frac{\partial}{\partial \theta_i} J(\theta) =$$

$$\frac{\partial}{\partial \theta_j} J(\theta) = \frac{\partial}{\partial \theta_j} \frac{1}{2} (h_{\theta}(x) - y)^2$$

$$\frac{\partial}{\partial \theta_{j}} J(\theta) = \frac{\partial}{\partial \theta_{j}} \frac{1}{2} (h_{\theta}(x) - y)^{2}$$

$$= 2 \cdot \frac{1}{2} (h_{\theta}(x) - y) \cdot \frac{\partial}{\partial \theta_{j}} (h_{\theta}(x) - y)$$

$$\frac{\partial}{\partial \theta_{j}} J(\theta) = \frac{\partial}{\partial \theta_{j}} \frac{1}{2} (h_{\theta}(x) - y)^{2}$$

$$= 2 \cdot \frac{1}{2} (h_{\theta}(x) - y) \cdot \frac{\partial}{\partial \theta_{j}} (h_{\theta}(x) - y)$$

$$= (h_{\theta}(x) - y) \cdot \frac{\partial}{\partial \theta_{j}} \left(\sum_{i=0}^{n} \theta_{i} x_{i} - y \right)$$

$$\frac{\partial}{\partial \theta_{j}} J(\theta) = \frac{\partial}{\partial \theta_{j}} \frac{1}{2} (h_{\theta}(x) - y)^{2}$$

$$= 2 \cdot \frac{1}{2} (h_{\theta}(x) - y) \cdot \frac{\partial}{\partial \theta_{j}} (h_{\theta}(x) - y)$$

$$= (h_{\theta}(x) - y) \cdot \frac{\partial}{\partial \theta_{j}} \left(\sum_{i=0}^{n} \theta_{i} x_{i} - y \right)$$

$$= (h_{\theta}(x) - y) x_{j}$$

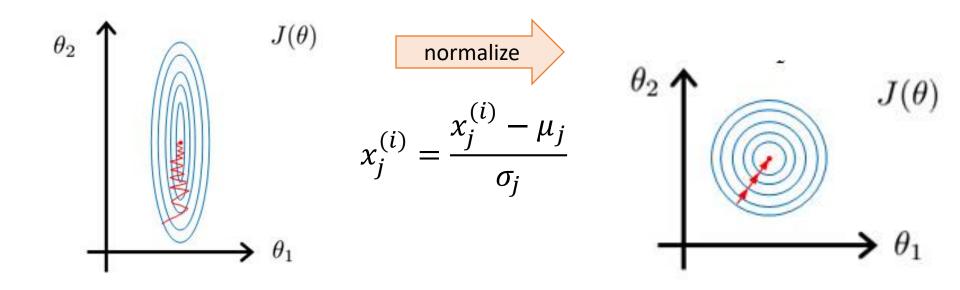
Gradient Descent Algorithm

```
Set m{	heta} = 0

Repeat { m{	heta}_j \coloneqq m{	heta}_j - lpha \frac{1}{m} \sum_{i=1}^m ig( h_{m{	heta}}(m{x}^{(i)}) - m{y}^{(i)} ig) m{x}_j^{(i)} \ j = 0, \dots, n } until convergence
```

Feature normalization

- If features have very different scale, GD can get "stuck" since x_j affects size of gradient in the direction of $j^{\rm th}$ dimension
- Normalizing features to be zero-mean (μ) and same-variance (σ) helps gradient descent converge faster

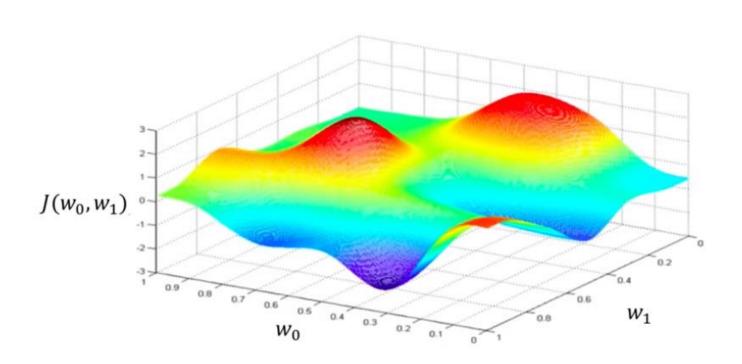




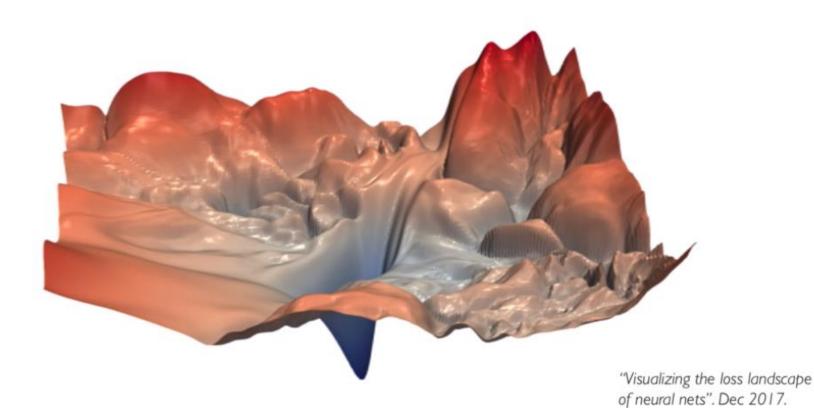
Supervised Learning

Loss Function and Learning Rate

Loss/Cost Function

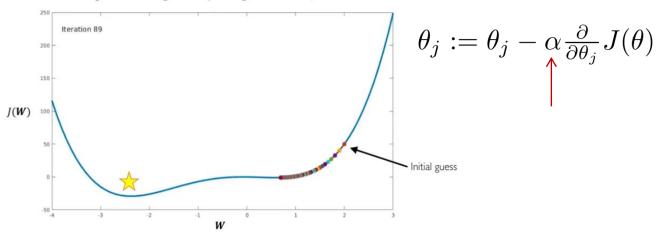


Landscape Visualization



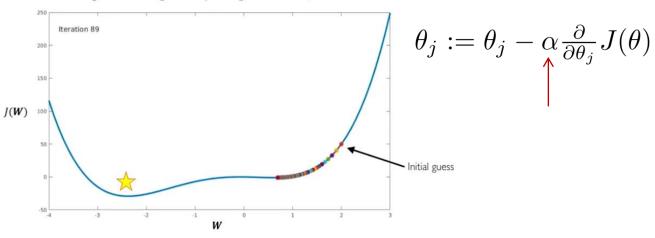
Setting the Learning Rate

Small learning rate converges slowly and gets stuck in false local minima



Setting the Learning Rate

Small learning rate converges slowly and gets stuck in false local minima



Large learning rates overshoot, become unstable and diverge

