



Siam Innovation District

DESIGN THINKING FOR BUSINESS INNOVATION

IDEATE: A lightbulb with a person jumping inside it, surrounded by gears, a magnifying glass, and a star.

EMPATHIZE: A brain icon, a question mark, and a person climbing a ladder.

DEFINE: A computer screen with a bar chart, a magnifying glass, and a star.

PROTOTYPE: A person standing next to a cardboard box with a person jumping out of it.

Kaweeuwut Temphuwapat
Team Lead, Innovation Lab and CVC,
PTT Group Design Leader,
MBA, Stanford University d.leader, Stanford d.school

Apply Now – Aug 31, 2017 : Open for Public

Course : SEP 9 – OCT 7, Every Saturday, 09.00 - 12.00
Room Rajakumari 60 Building (Chamchuri 10 Bldg), Fl 4
Selected candidates will be fully-funded towards the course fee, worth 35,000 baht

Contact Us : 093-725-0808, 02-218-3106-7
<https://goo.gl/forms/CQzJLcXWzq>

MACHINE LEARNING

Siam Innovation District Tech Talent

Machine Learning

- Essential tools and libraries: Python, Jupyter Notebook, NumPy, Pandas, SciPy, Scikit-Learn, Matplotlib, and Seaborn
- Data collection through API and web scraping
- Machine Learning Algorithms reviews

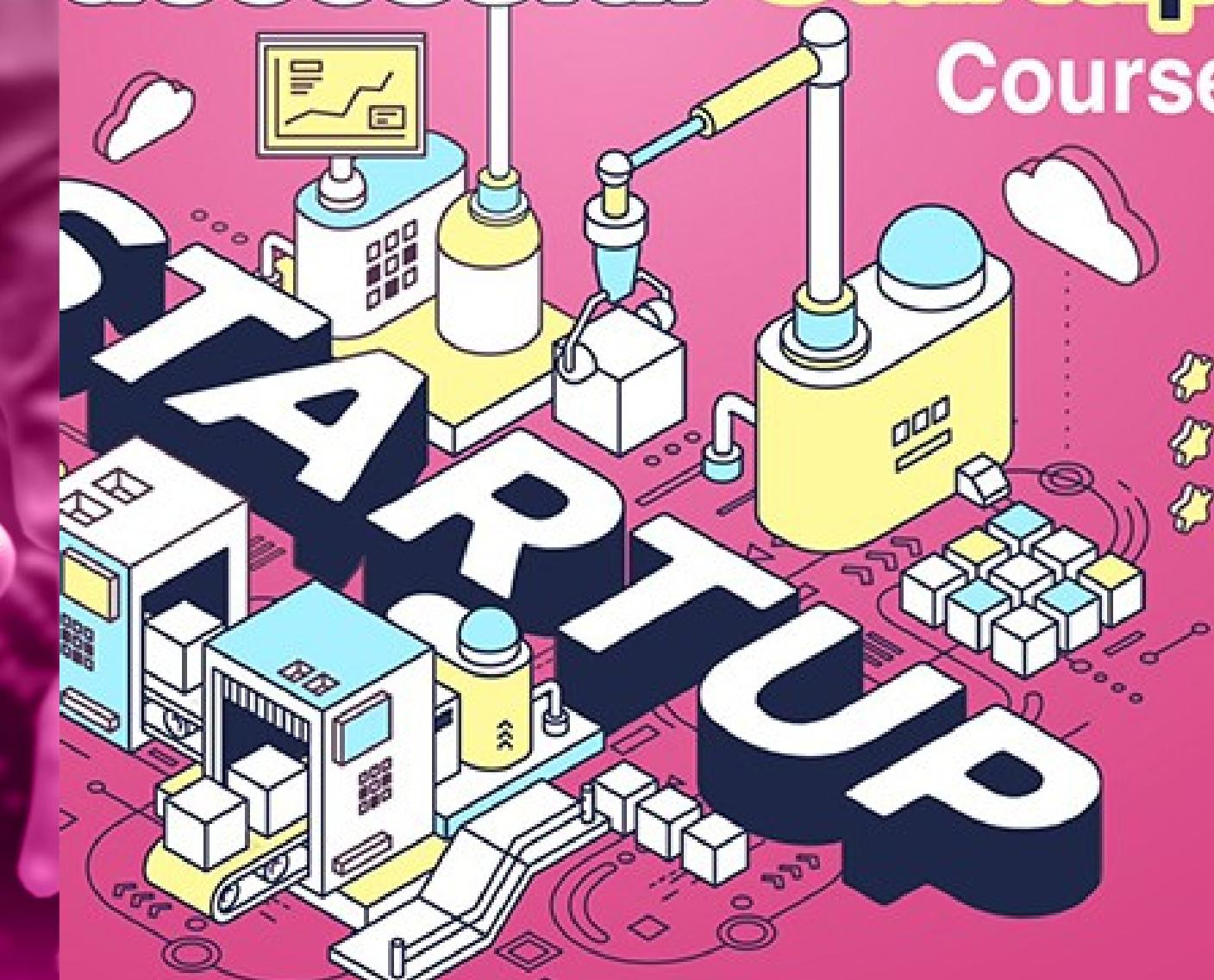
Warodom Khamphanchai, Ph.D.
Bangkok AI Ambassador,
Ex-Software Developer at
Samsung SmartThings in Palo Alto,
CA. Ex-Full Stack Developer at

Sorawit Saengkyongam
Data Scientist at Agoda,
Google Developer Expert in
Machine Learning.

----- Agenda ครับ

1. Meet guest developers from Silicon Valley
2. What you need to know to convince your prospective CTO or developer?
3. How can you communicate your ideas to your CTO or developer?

24 Steps to Successful Startup Course



I want to make happen? Join this comprehensive startup course will take you from idea to product in 6-weeks!

This course to take your startup idea Research through to Building your startup pitch deck and get ready for startup building course from the MIT



Tareef Jaffer
Ex-teaching staff at MIT,
Ex-Goolger, Serial Entrepreneur



Aug 31, 2017 : Open for Public

Course : SEP 9 – OCT 7, Every Saturday, 9.00 - 13.00
Room Rajakumari 60 Building (Chamchuri 10 Bldg), Fl 4
Selected candidates will be fully-funded towards the course fee, worth 35,000 baht

093-725-0808

02-218-3106-7

www.facebook.com/cuinhub

@cuinnovationhub

innovationhub.chula.ac.th

cu.innovation.hub@gmail.com



<http://bit.ly/cueDE>



Warodom Khamphanchai, PhD

- LEAD SMART HOME/BUILDING PLATFORM DEVELOPER @ PEA
- EX-SOFTWARE DEVELOPMENT ENGINEER @ SAMSUNG SMARTTHINGS

Interests: Home/Building Automation, Internet of Things, Smart Grid, Multi-Agent systems, Data Analytics, Machine Learning, Deep Learning, AI, Energy Audit, and Technology Entrepreneurship.

Education:

- [2011-2016] Ph.D. in Electrical and Computer Engineering, Virginia Tech

Dissertation: An Agent-based Platform for Demand Response Implementation in Smart Buildings

- [2009-2011] M.E. in Energy (Area of Specialization: Electric Power System Management), Asian Institute of Technology Thesis: A Multi-Agent Based Power System Restoration

Approach in Distributed Smart Grid

- [2005-2009] B.E. in Electrical Engineering, Chulalongkorn University



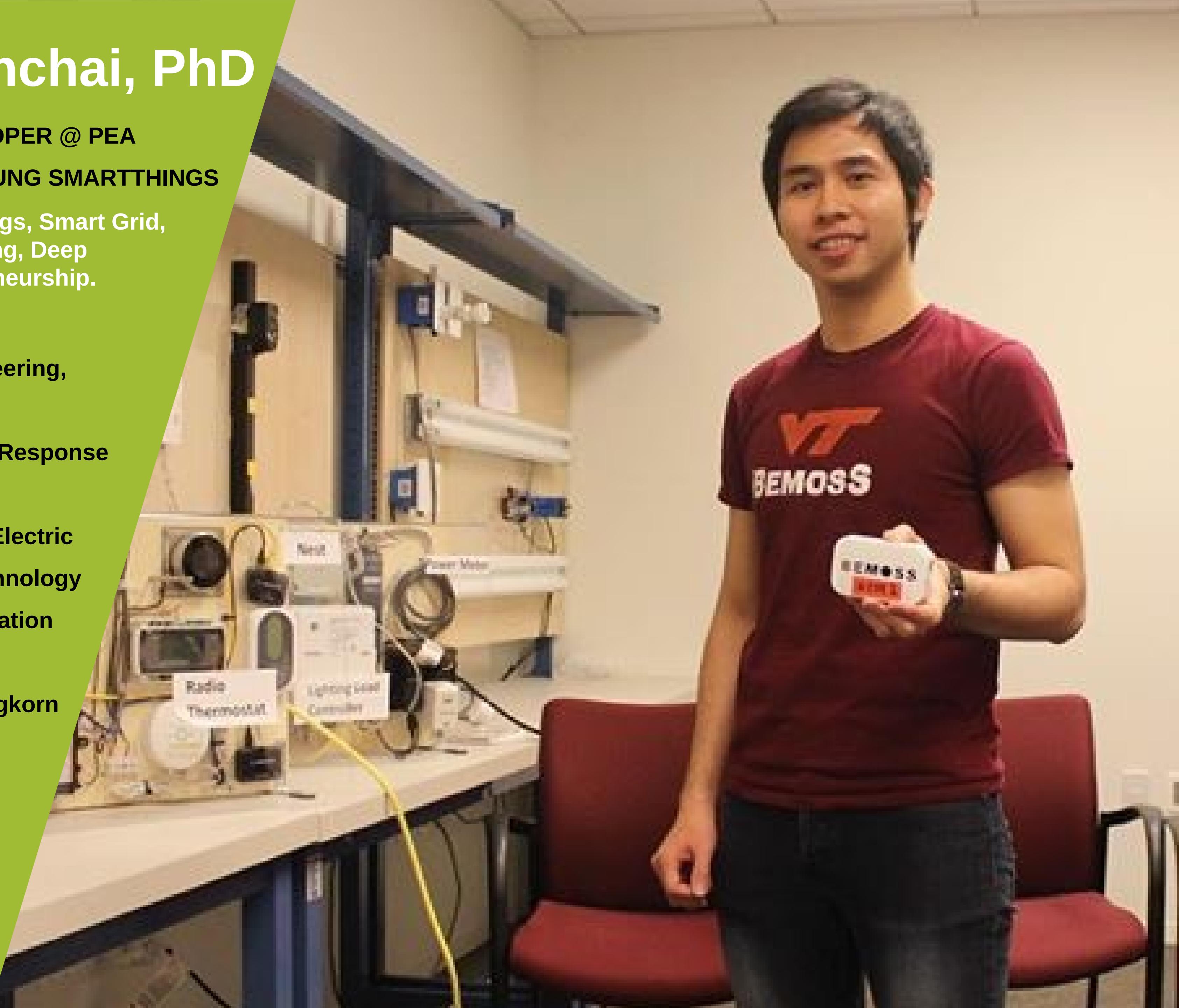
Line: kwarodom

LinkedIn: www.linkedin.com/in/kwarodom

Web: kwarodom.wordpress.com

Email: kwardom@vt.edu

Tel: +6695-161-5011 Github: kwarodom



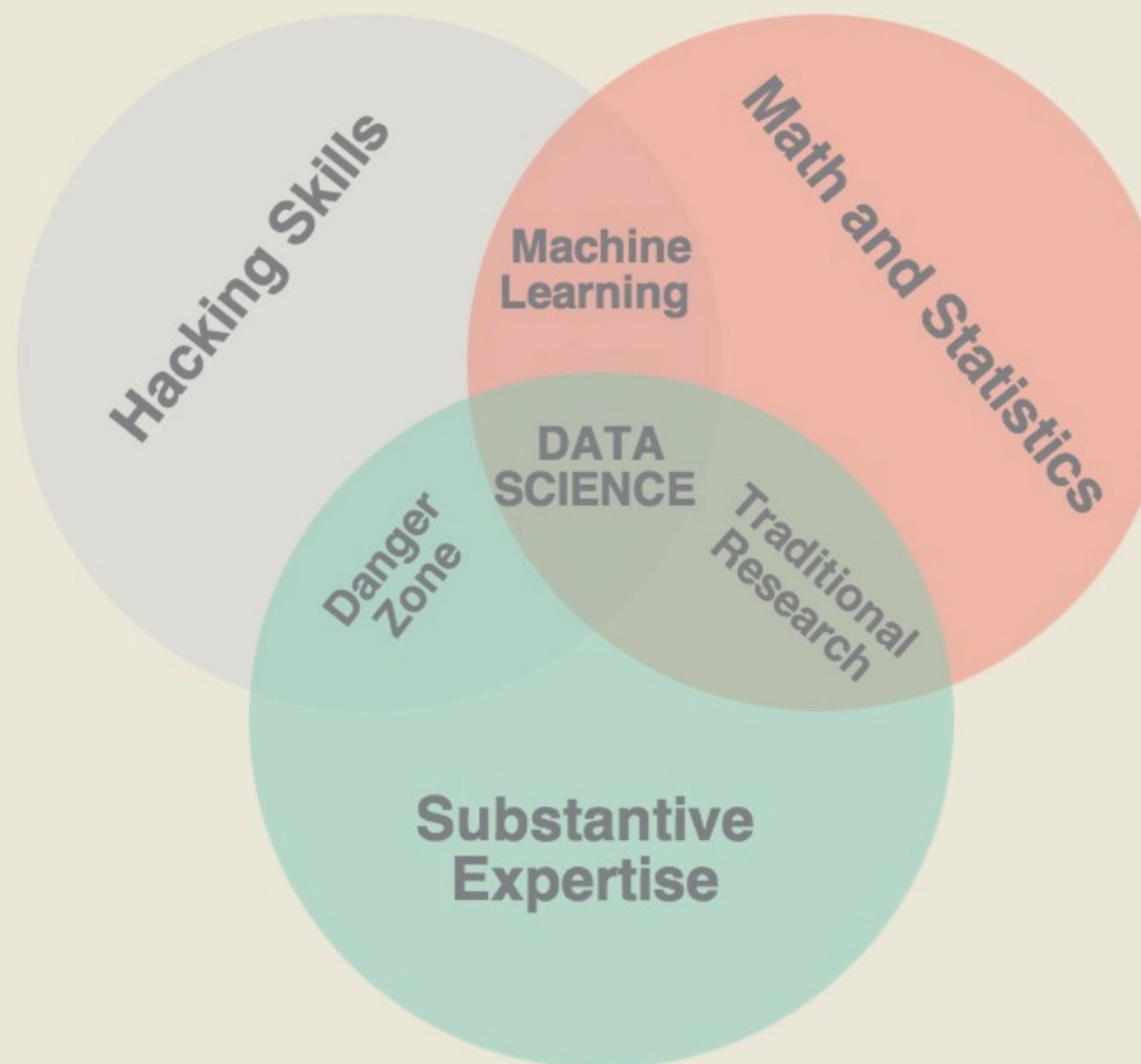


Profile

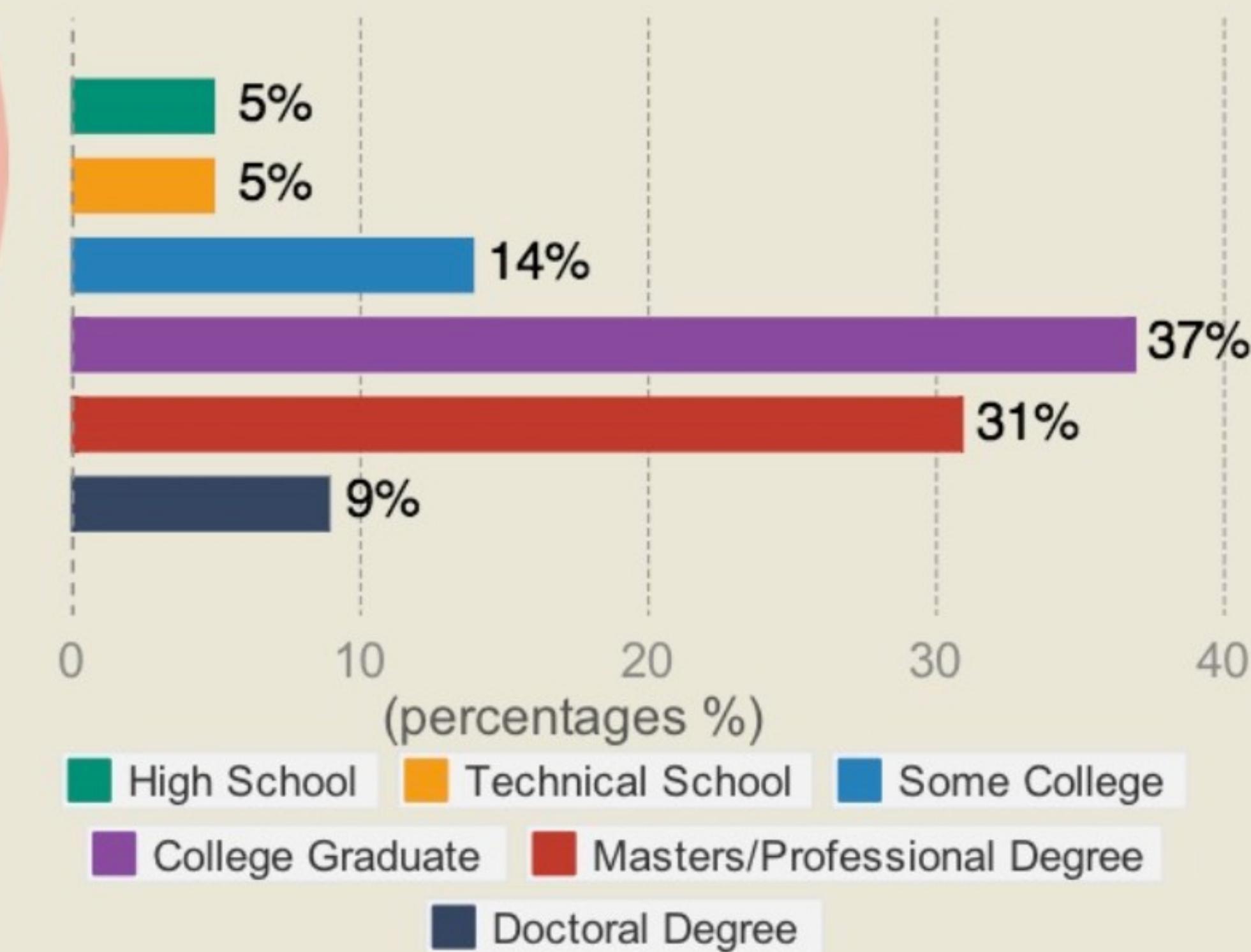
Sorawit Saengkyongam (James)

- Data Scientist at Agoda where he works on research and development of personalization algorithms and recommendation systems
- Google Developer Expert in Machine Learning.
- Organize Bangkok Machine Learning Meetup
- Graduated with a major in Mathematical Statistics from Chulalongkorn University (with first class honours)

What's a data scientist?



Typical Background



A data scientist is someone who is better at statistics than any software engineer and better at software engineering than any statistician.

Gareth James
Daniela Witten
Trevor Hastie
Robert Tibshirani

An Introduction to Statistical Learning

with Applications in R



□ Programming

- R programming language
- Python programming language
- Spreadsheet tools (like Excel)
- JavaScript and HTML
- C/C++

□ Statistics

- Descriptive and Inferential statistics
- Experimental design

□ Mathematics

- College Algebra
- Functions and Graphing
- Multivariable Calculus
- Linear Algebra

□ Machine Learning

- Supervised Learning
- Unsupervised Learning
- Reinforcement Learning

□ Data Wrangling

- Python
- Database Systems
- SQL

□ Communication and Data Visualization

- Visual Encoding
- Data Presentation
- Knowing Your Audience

□ Data Intuition (Thinking like a data scientist)

- Project Management
- Industry Knowledge

05

07

09

10

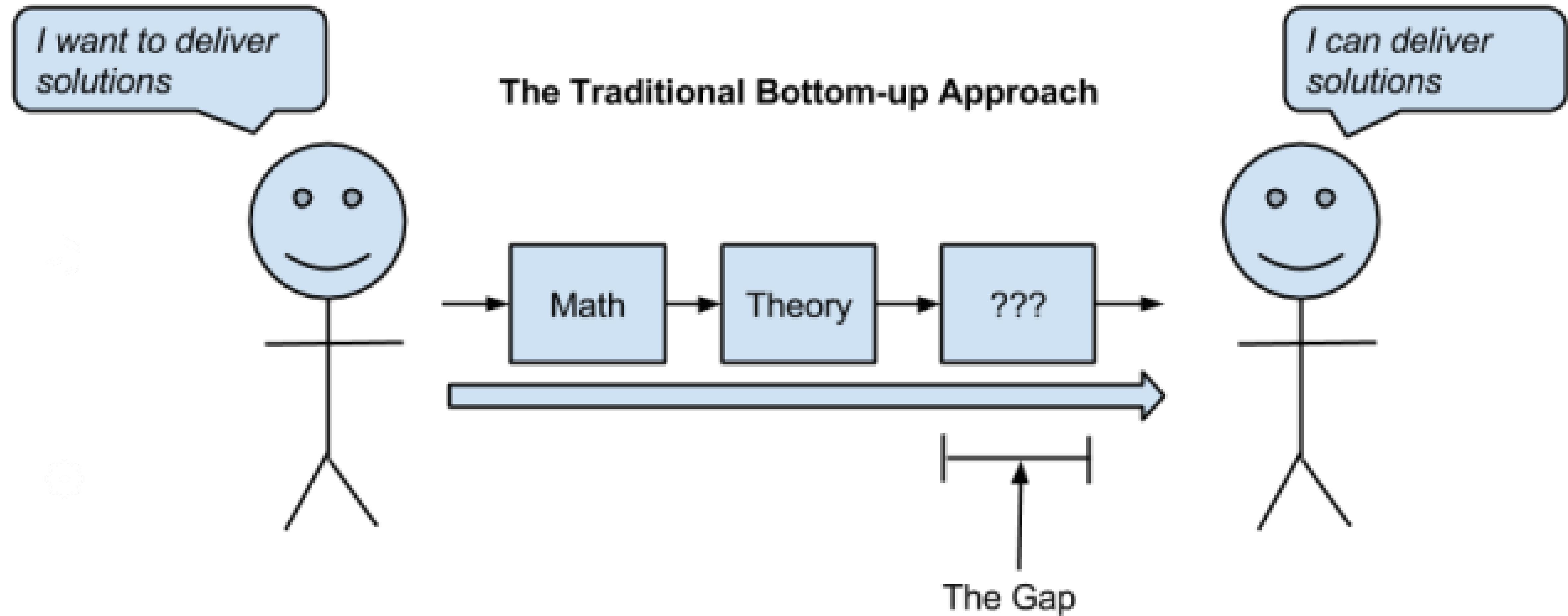
12

13

14



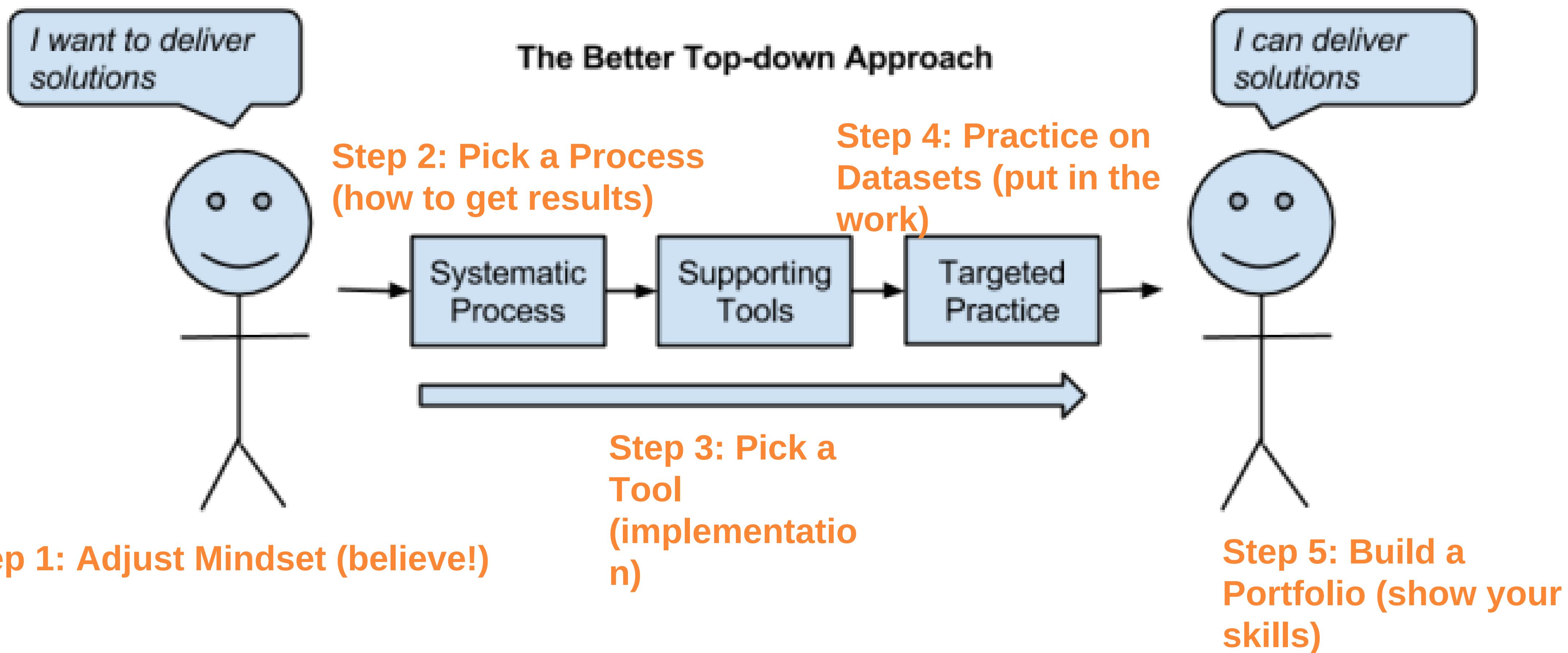
MACHINE LEARNING PRACTITIONER





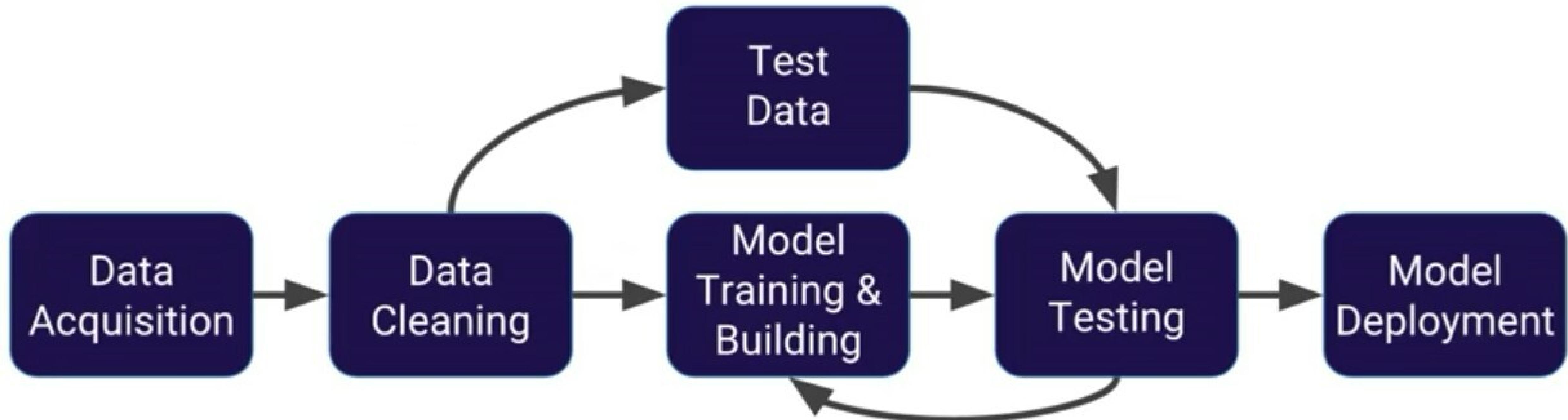
MACHINE LEARNING PRACTITIONER

5-Steps To Get Started and Get Good at Machine Learning



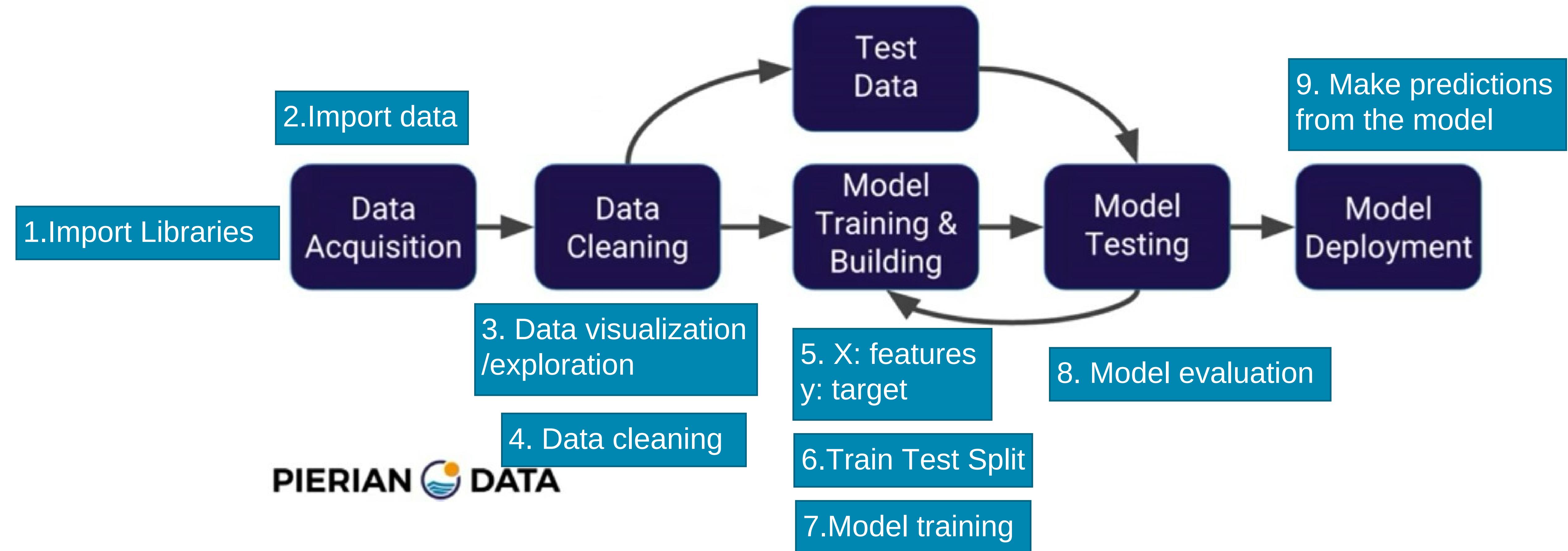


MACHINE LEARNING PROCESS





Machine Learning Process



Machine Learning

What is Machine Learning?

Two definitions of Machine Learning are offered.

Arthur Samuel described it as: "the field of study that gives computers the ability to learn without being explicitly programmed."

Tom Mitchell provides a more modern definition: "A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P, if its performance at tasks in T, as measured by P, improves with experience E."

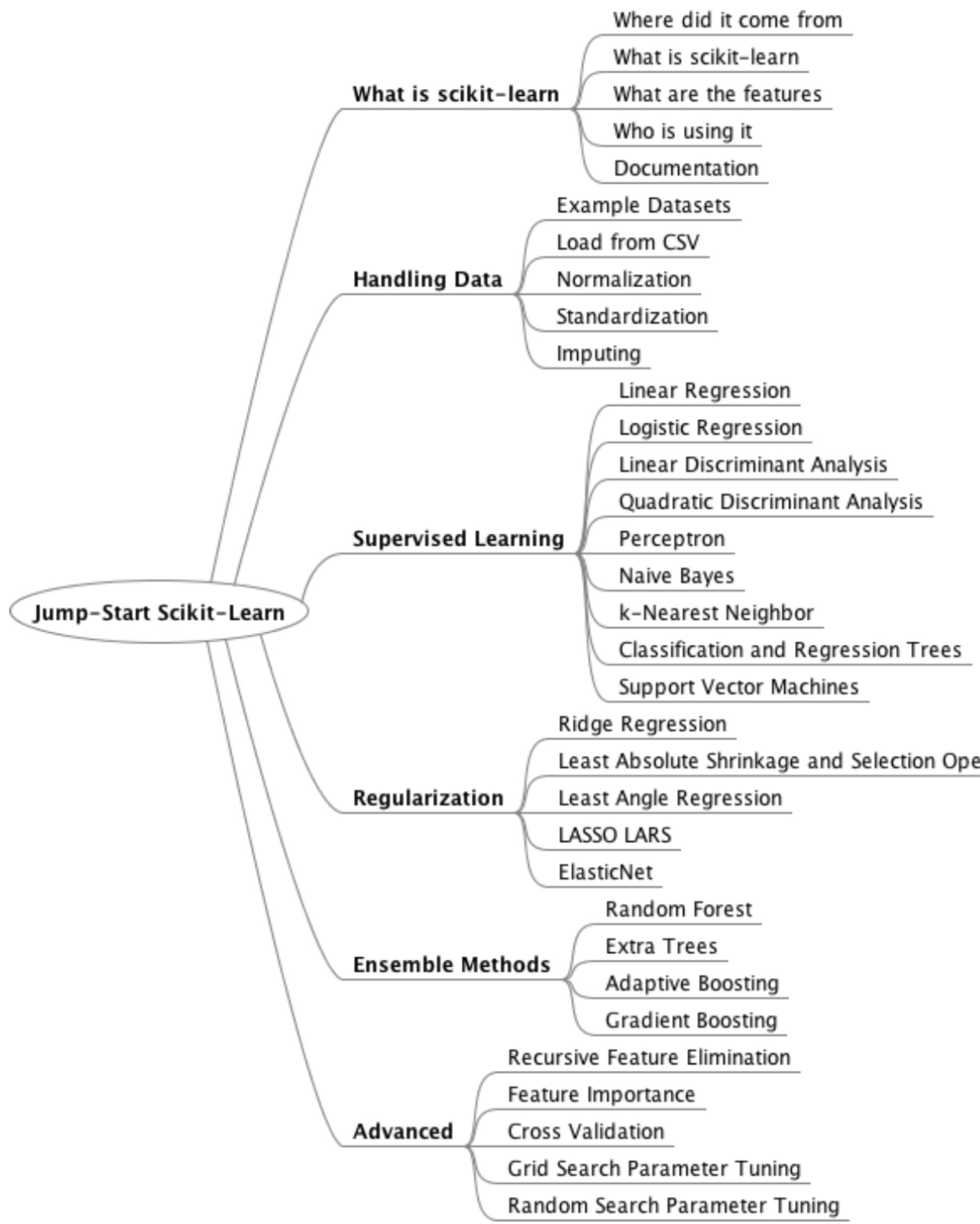
Example: playing checkers.

E = the experience of playing many games of checkers

T = the task of playing checkers.

P = the probability that the program will win the next game.

Scikit Learn



Supervised Learning

In supervised learning, we are given a data set and already know what our correct output should look like, having the idea that there is a relationship between the input and the output.

Supervised learning problems are categorized into "regression" and "classification" problems.

In a regression problem, we are trying to predict results within a continuous output, meaning that we are trying to map input variables to some continuous function.

In a classification problem, we are instead trying to predict results in a discrete output. In other words, we are trying to map input variables into discrete categories.

Data Preparation

Normalization

```
# Normalize the data attributes for the Iris dataset.  
from sklearn.datasets import load_iris  
from sklearn import preprocessing  
# load the iris dataset  
iris = load_iris()  
print(iris.data.shape)  
# separate the data from the target attributes  
X = iris.data  
y = iris.target  
# normalize the data attributes  
normalized_X = preprocessing.normalize(X)
```

Standardisation

```
# Standardize the data attributes for the Iris dataset.  
from sklearn.datasets import load_iris  
from sklearn import preprocessing  
# load the Iris dataset  
iris = load_iris()  
print(iris.data.shape)  
# separate the data and target attributes  
X = iris.data  
y = iris.target  
# standardize the data attributes  
standardized_X = preprocessing.scale(X)
```

Data Preparation

Imputing

Data can have missing values. These are values for attributes where a measurement could not be taken or is corrupt for some reason. It is important to identify and mark this missing data. Once marked, replacement values can be prepared. This is useful because some algorithms are unable to work with or exploit missing data.

```
# Mark 0 values as missing and impute with the mean
import numpy as np
import urllib
from sklearn.preprocessing import Imputer
# Load the Pima Indians Diabetes dataset
url = "http://goo.gl/j0Rvxq"
raw_data = urllib.urlopen(url)
dataset = np.loadtxt(raw_data, delimiter=",")
print(dataset.shape)
# separate the data and target attributes
X = dataset[:,0:7]
y = dataset[:,8]
# Mark all zero values as 0
X[X==0]=np.nan
# Impute all missing values with the mean of the attribute
imp = Imputer(missing_values='NaN', strategy='mean')
imputed_X = imp.fit_transform(X)
```

Regression

Hypothesis:

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

Parameters:

$$\underline{\theta_0, \theta_1}$$

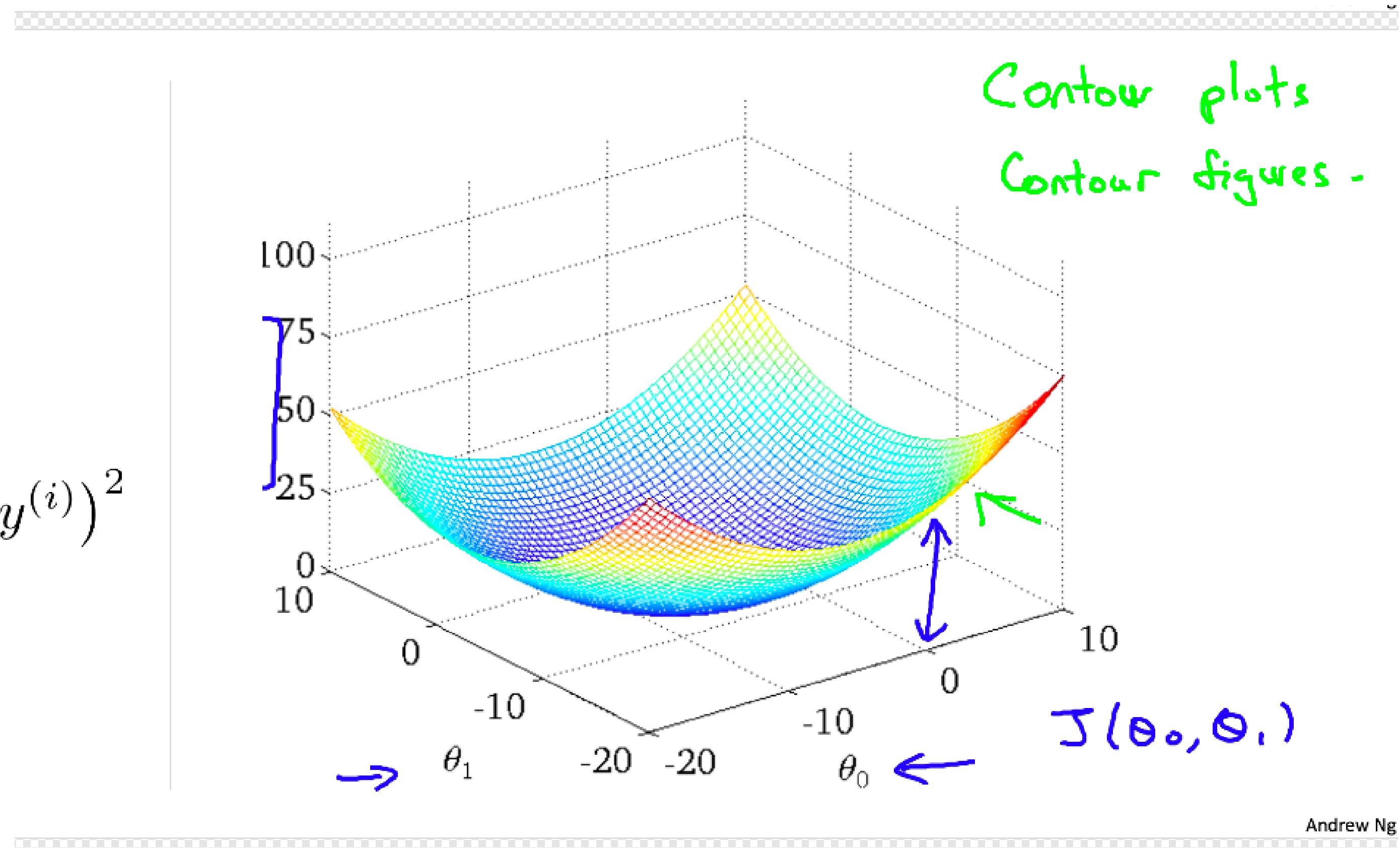
Cost Function:

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

Goal:

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

.



Linear Regression with One Variable

Model and Cost Function

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

Parameters: $\underline{\theta_0, \theta_1}$

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

Parameter Learning

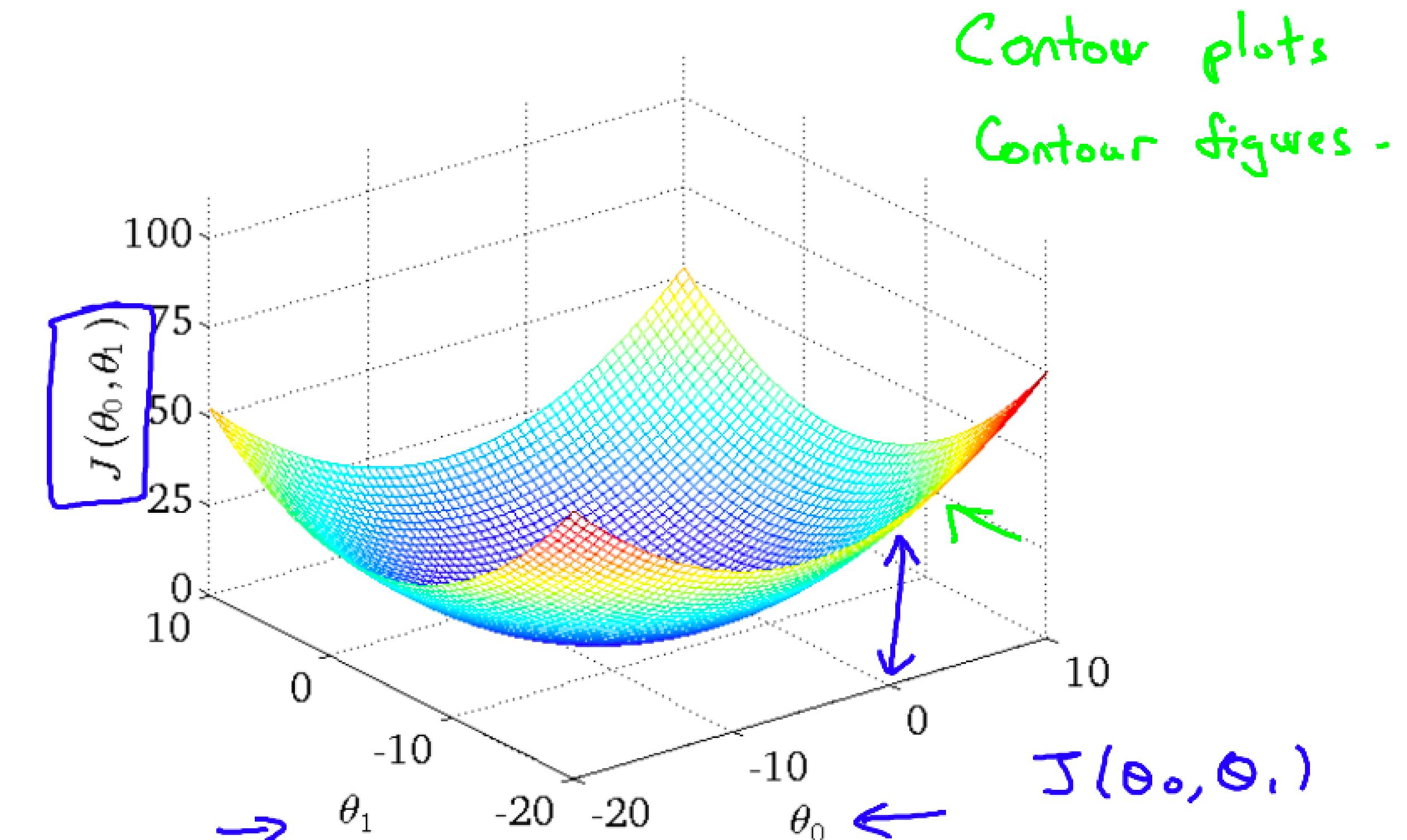
repeat until convergence: {

Gradient Descent

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

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

}



Andrew Ng

Multivariate Linear Regression

Multiple Features

Multiple features (variables).

Size (feet ²) x_1	Number of bedrooms x_2	Number of floors x_3	Age of home (years) x_4	Price (\$1000) 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

$$n = 4$$

→ $x^{(i)}$ = input (features) of i^{th} training example.

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

$$\underline{x}^{(2)} = \begin{bmatrix} 1416 \\ 3 \\ 2 \\ 40 \end{bmatrix}$$

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

Multivariate Linear Regression

Multiple Features

Hypothesis: $\underline{h_\theta(x) = \theta^T x = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n}$ $x_0 = 1$

Parameters: $\underline{\theta_0, \theta_1, \dots, \theta_n}$ $\underbrace{\theta}_{n+1 - \text{dimensional vector}}$

Cost function:

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

$J(\theta)$

Gradient descent:

Repeat {
 $\rightarrow \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \dots, \theta_n)$ $J(\theta)$
}

↑ (simultaneously update for every $j = 0, \dots, n$)

Multivariate Linear Regression

Gradient Descent for Multiple Variables

Gradient Descent

Previously (n=1):

Repeat {

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

$\frac{\partial}{\partial \theta_0} J(\theta)$

$$\rightarrow \theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_1^{(i)}$$

(simultaneously update θ_0, θ_1)

}

New algorithm ($n \geq 1$):

Repeat {

$$\downarrow \frac{\partial}{\partial \theta_j} J(\theta)$$

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

(simultaneously update θ_j for
 $j = 0, \dots, n$)

$$x_0^{(i)} = 1$$

}

$$\rightarrow \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\rightarrow \theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_1^{(i)}$$

$$\rightarrow \theta_2 := \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_2^{(i)}$$

...

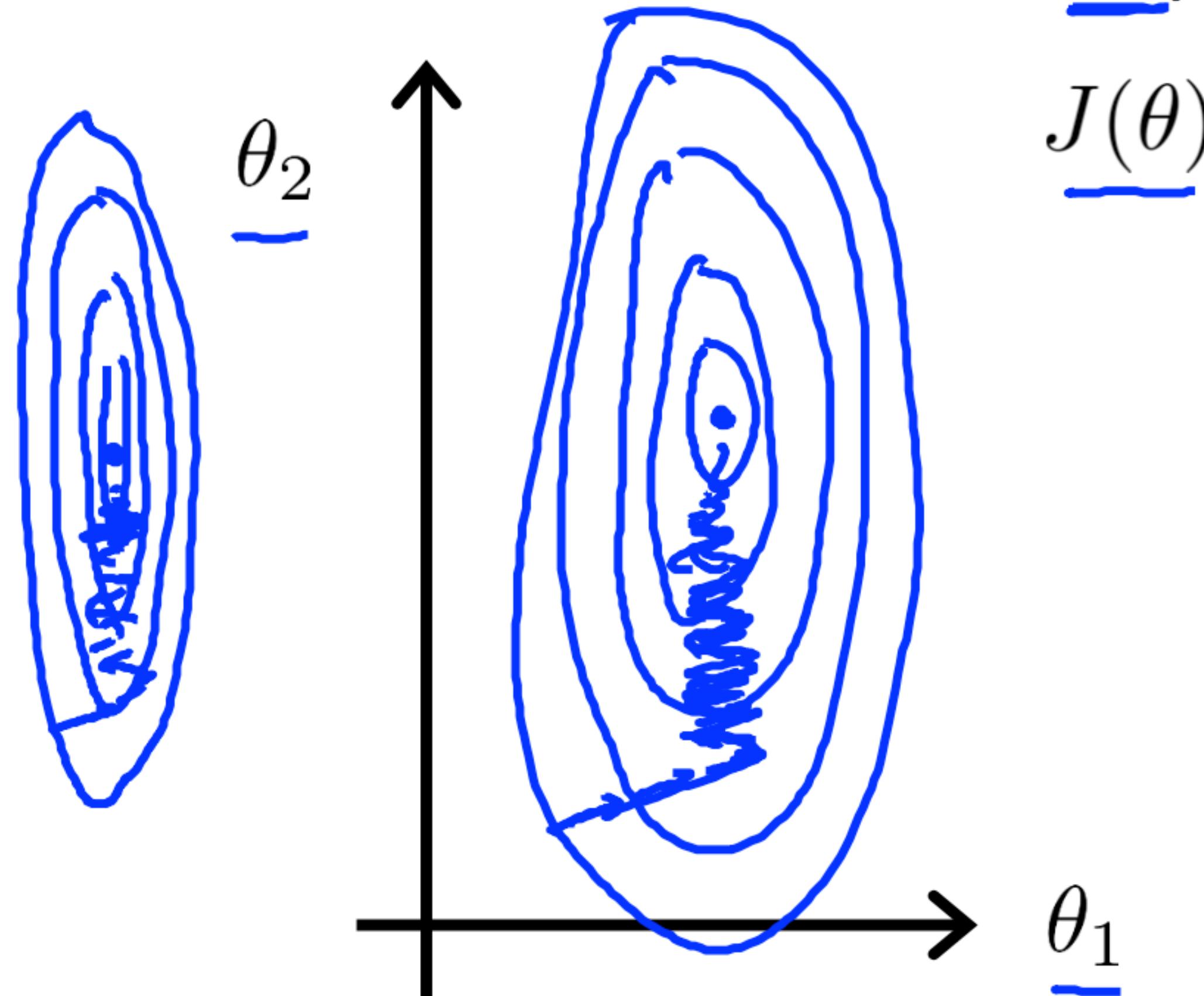
Multivariate Linear Regression

Feature Scaling

Idea: Make sure features are on a similar scale.

$$\text{E.g. } x_1 = \text{size (0-2000 feet}^2)$$

$$x_2 = \text{number of bedrooms (1-5)}$$



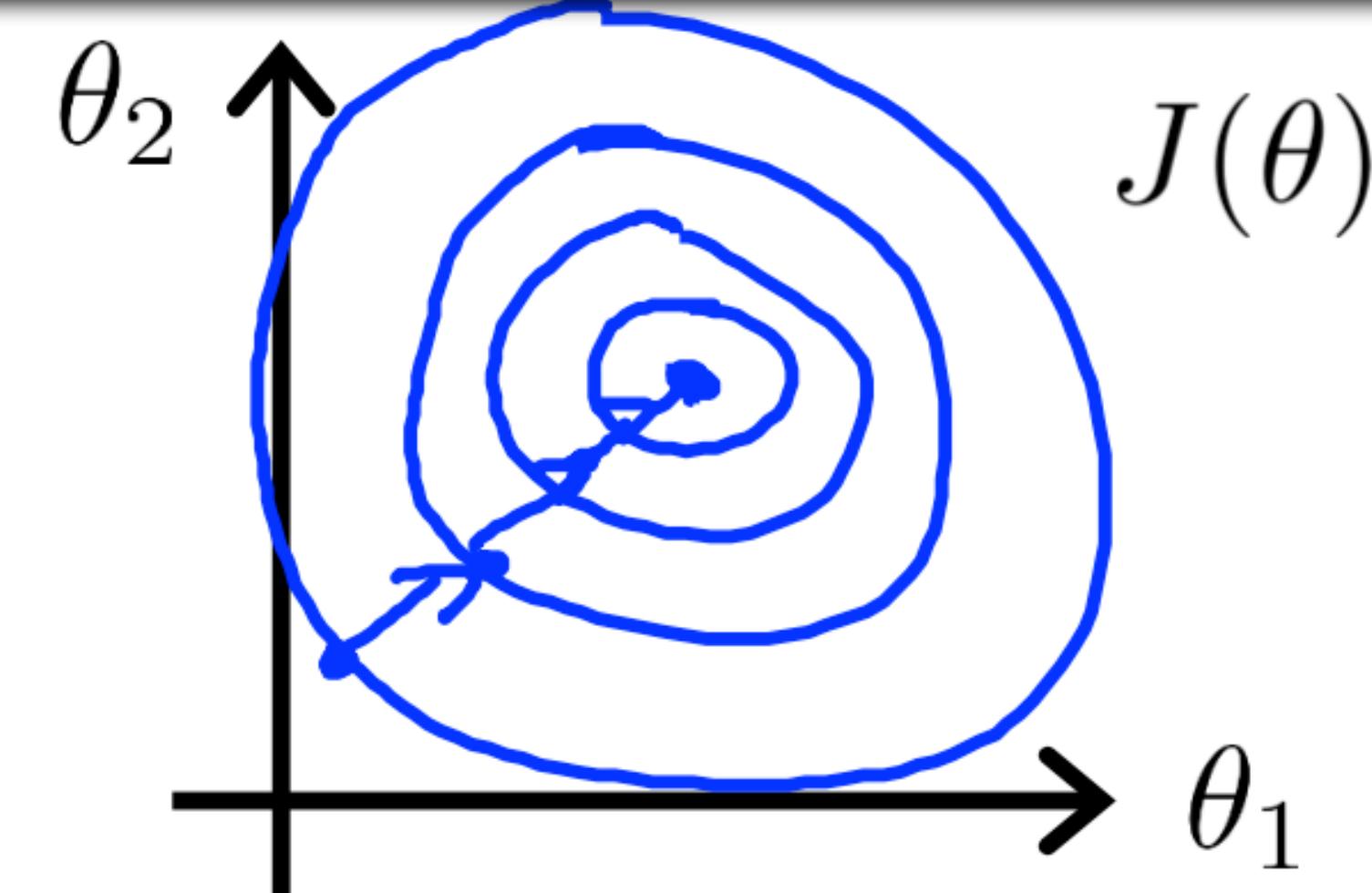
$$\rightarrow x_1 = \frac{\text{size (feet}^2)}{2000}$$

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

$$0 \leq x_1 \leq 1$$

$$0 \leq x_2 \leq 1$$

Get every feature into approximately $-1 \leq x_i \leq 1$



Multivariate Linear Regression

Mean normalization

Replace x_i with $\frac{x_i - \mu_i}{\sigma_i}$ to make features have approximately zero mean
 (Do not apply to $x_0 = 1$).

E.g. $x_1 = \frac{\text{size} - 1000}{2000}$

Avg size = 1000

$$x_2 = \frac{\#\text{bedrooms} - 2}{5 - 4}$$

1-5 bedrooms

$$\rightarrow [-0.5 \leq x_1 \leq 0.5, -0.5 \leq x_2 \leq 0.5]$$

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

avg value of x_1 in training set

range ($\max - \min$)

(or standard deviation)

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

Learning Rate

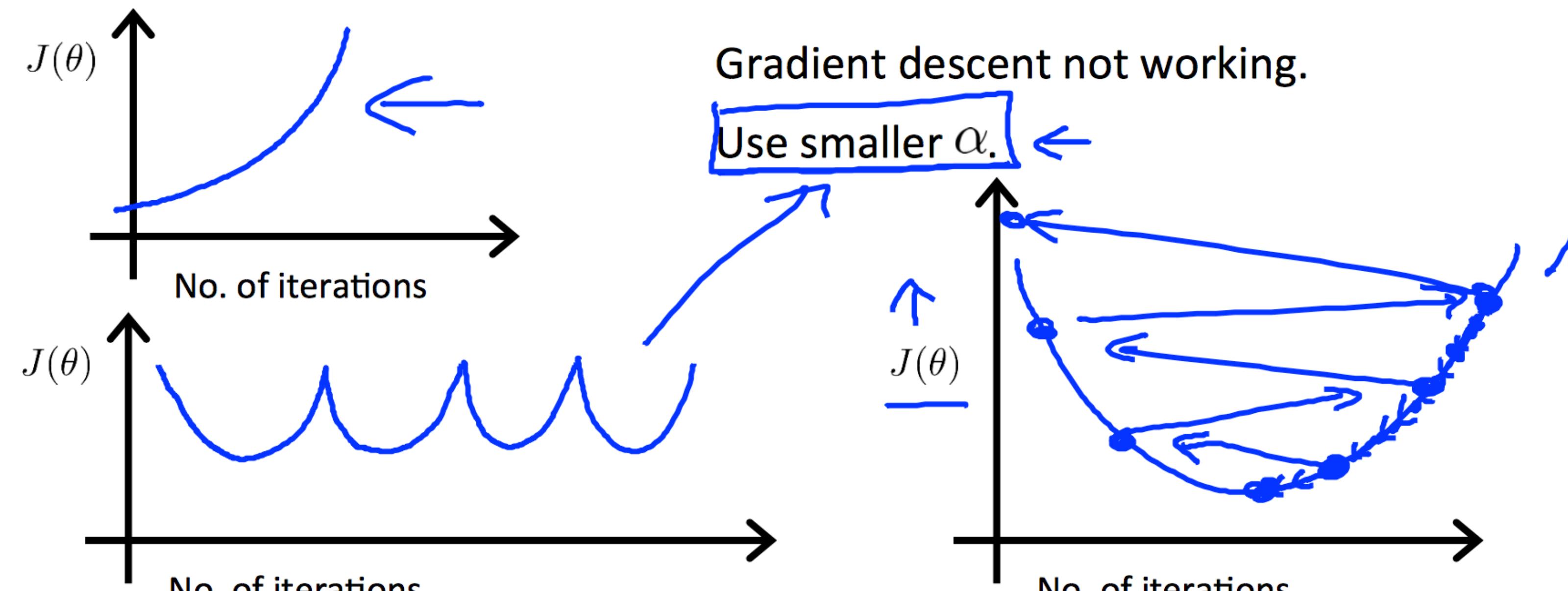
Multivariate Linear Regression

Gradient descent **How to choose learning rate?**

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

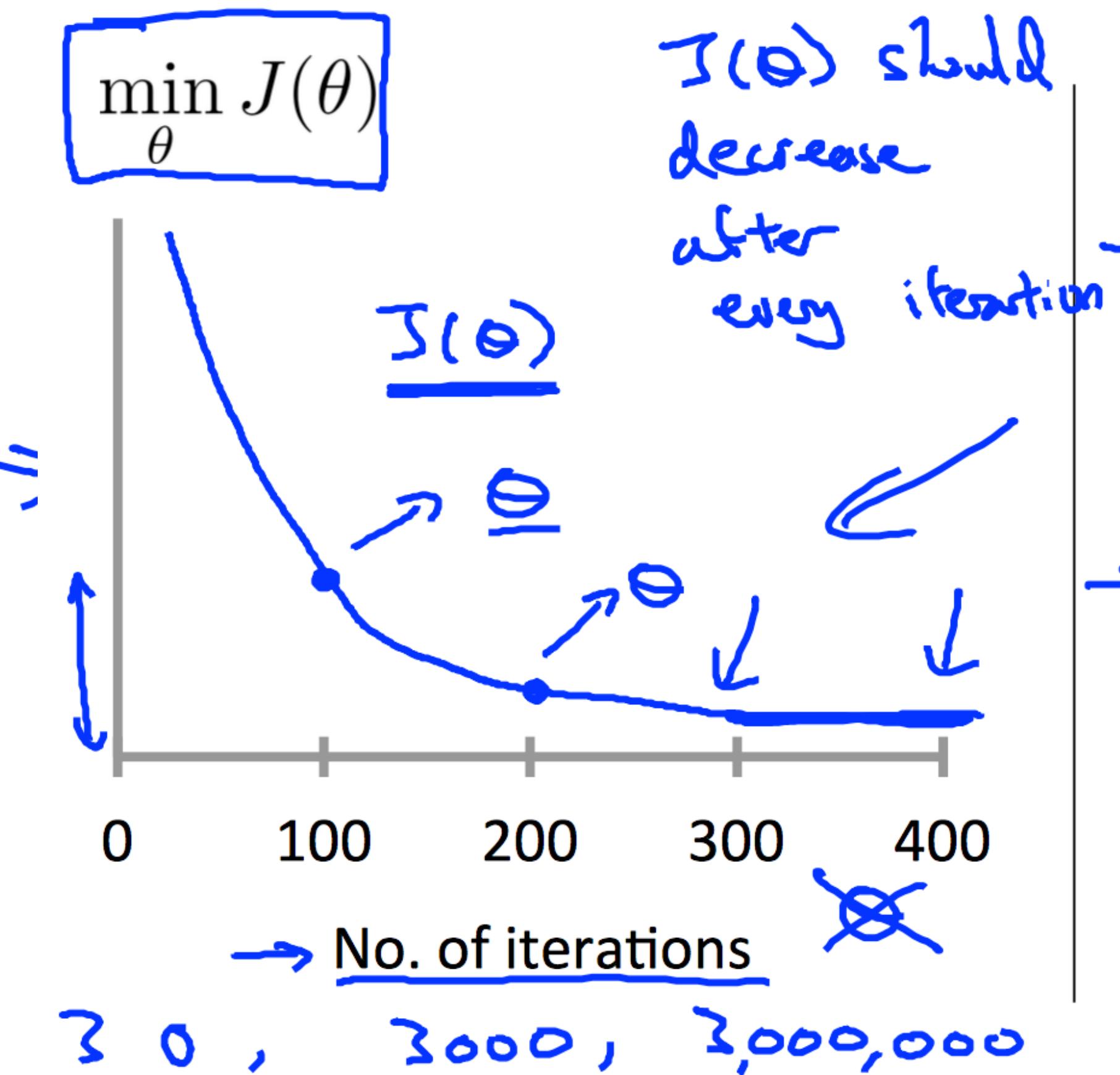
If α is too small, slow convergence
If α is too large, $J(\theta)$ may not decrease on every iteration, may not converge

Making sure gradient descent is working correctly.



- For sufficiently small α , $J(\theta)$ should decrease on every iteration.
- But if α is too small, gradient descent can be slow to converge.

To α choose , try ..., 0.001, 0.003, 0.01, 0.03, 0.1, 0.3, 1, ...



Linear Regression with One Variable

Linear Regression

Linear regression fits a linear model (such as a line in two dimensions) to the data.

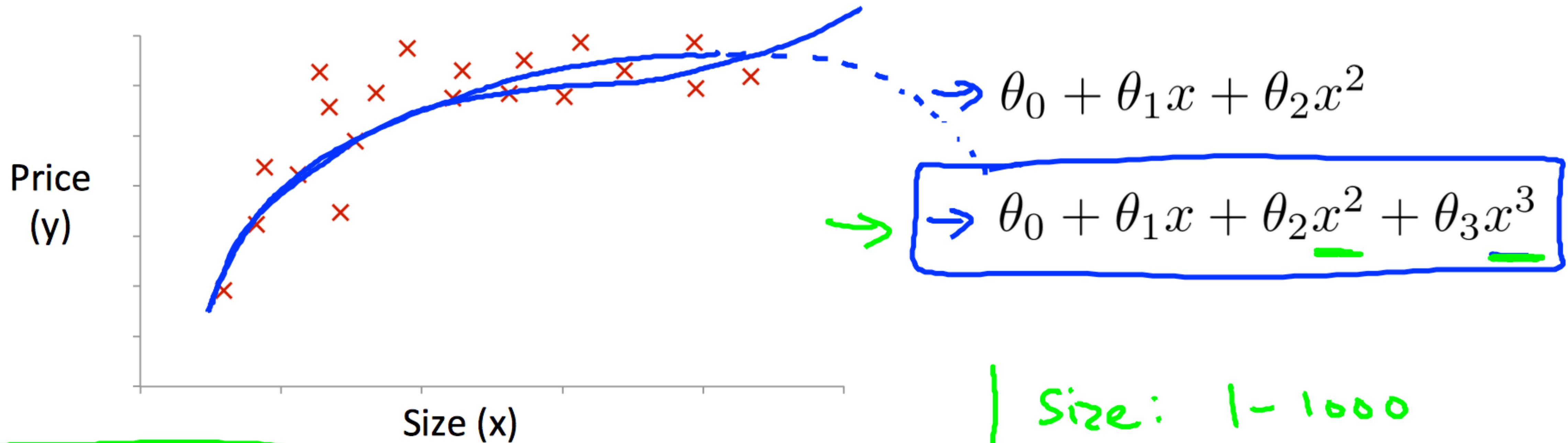
This recipes shows the fitting of a linear regression algorithm on the diabetes dataset.

```
# Linear Regression
import numpy as np
from sklearn import datasets
from sklearn.linear_model import LinearRegression
# load the diabetes datasets
dataset = datasets.load_diabetes()
# fit a linear regression model to the data
model = LinearRegression()
model.fit(dataset.data, dataset.target)
print(model)
# make predictions
expected = dataset.target
predicted = model.predict(dataset.data)
# summarize the fit of the model
mse = np.mean((predicted-expected)**2)
print(mse)
print(model.score(dataset.data, dataset.target))
```

Polynomial Regression

Feature and Polynomial Regression

Polynomial regression



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

$\rightarrow x_1 = (\text{size})$

$\rightarrow x_2 = (\text{size})^2$

$\rightarrow x_3 = (\text{size})^3$

Size: 1 - 1000
Size²: 1 - 1000,000
Size³: 1 - 10⁹

Normal Equation

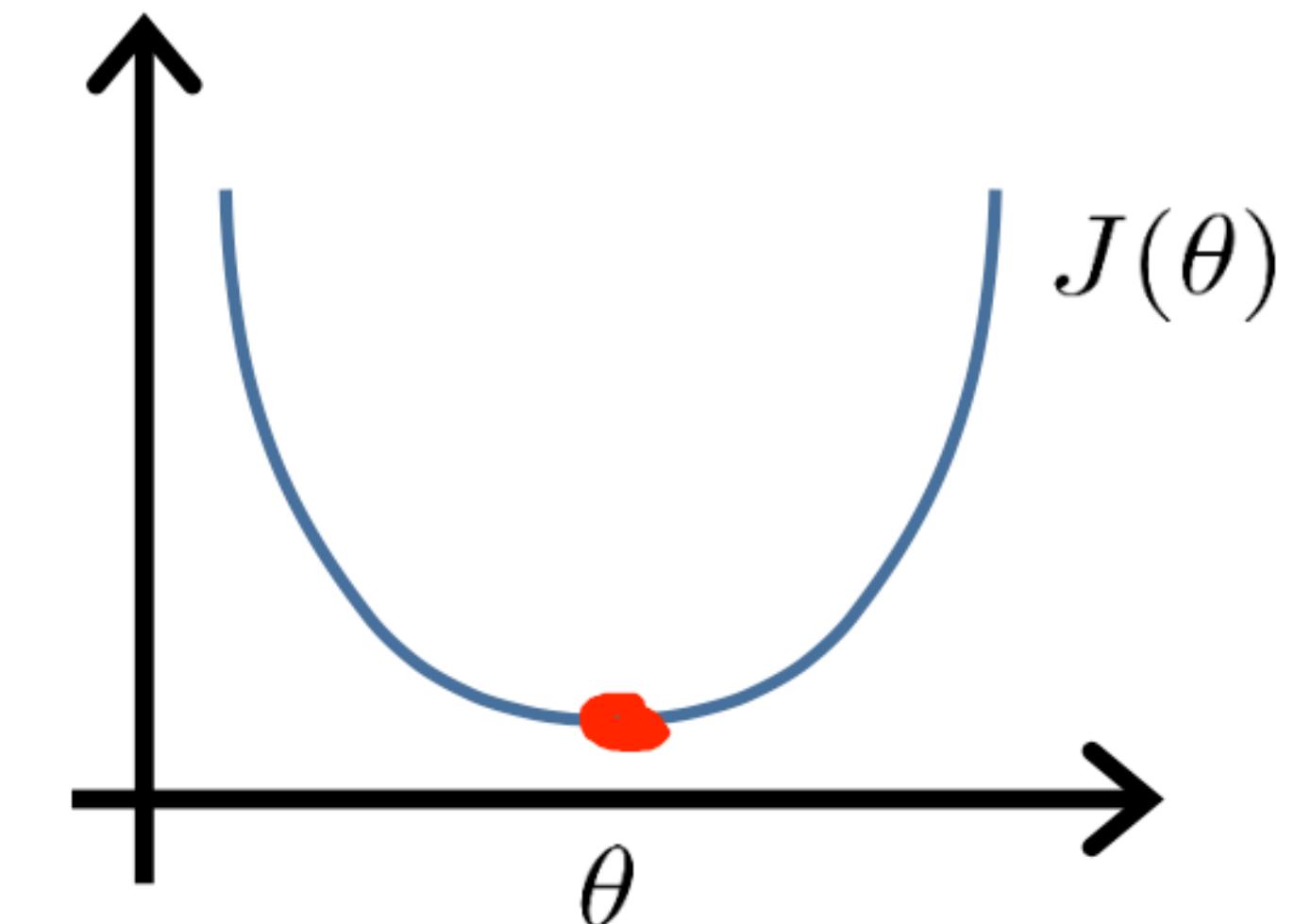
Computing Parameters Analytically

Intuition: If 1D ($\theta \in \mathbb{R}$)

$$\rightarrow J(\theta) = a\theta^2 + b\theta + c$$

$$\frac{\partial}{\partial \theta} J(\theta) = \dots \stackrel{\text{set}}{=} 0$$

Solve for θ



$$\theta \in \mathbb{R}^{n+1}$$

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

$$\frac{\partial}{\partial \theta_j} J(\theta) = \dots \stackrel{\text{set}}{=} 0 \quad (\text{for every } j)$$

Solve for $\theta_0, \theta_1, \dots, \theta_n$

Normal Equation

Computing Parameters Analytically

Examples: $m = 4$.

	Size (feet ²)	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
x_0	x_1	x_2	x_3	x_4	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

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

$m \times (n+1)$

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

m -dimensional vector

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

Computing Parameters Analytically

Normal Equation Noninvertibility

What if $\boxed{X^T X}$ is non-invertible?

- Redundant features (linearly dependent).

E.g. $\underline{x_1} = \text{size in feet}^2$

~~$x_2 = \text{size in m}^2$~~

$$\underline{x_1} = (3.28)^2 \underline{x_2}$$

$$1_m = 3.28 \text{ feet}$$

$$\rightarrow \underline{n = 10} \leftarrow$$

$$\rightarrow \underline{n = 100} \leftarrow$$

$$\Theta \in \mathbb{R}^{101}$$

- Too many features (e.g. $m \leq n$).

- Delete some features, or use regularization.

↓
later

Classification

Logistic regression is a method for classifying data into discrete outcomes. For example, we might use logistic regression to classify an email as spam or not spam. In this module, we introduce the notion of classification, the cost function for logistic regression, and the application of logistic regression to multi-class classification.

We are also covering regularization. Machine learning models need to generalize well to new examples that the model has not seen in practice. We'll introduce regularization, which helps prevent models from overfitting the training data.

Classification and Representation

Classification

Classification

- Email: Spam / Not Spam?
 - Online Transactions: Fraudulent (Yes / No)?
 - Tumor: Malignant / Benign ?
- $y \in \{0, 1\}$
- 0: “Negative Class” (e.g., benign tumor)
- 1: “Positive Class” (e.g., malignant tumor)
- $y \in \{0, 1, 2, 3\}$

Classification and Representation

Hypothesis representation and Cost Function

Logistic regression cost function

$$\begin{aligned} J(\theta) &= \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_\theta(x^{(i)}), y^{(i)}) \\ &= -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_\theta(x^{(i)})) \right] \end{aligned}$$

To fit parameters θ :

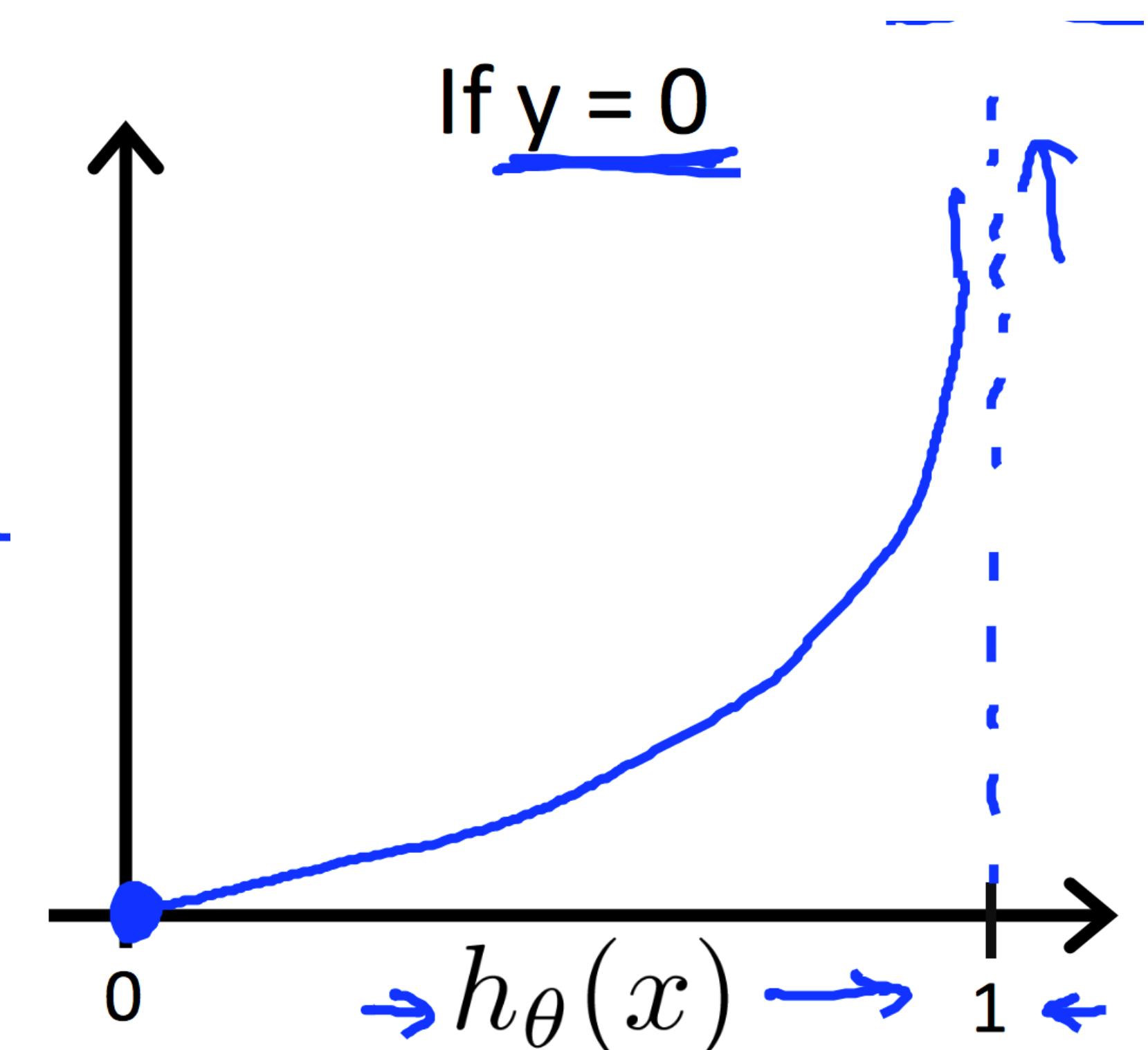
$$\min_{\theta} J(\theta)$$

Gret θ

To make a prediction given new x :

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

$$p(y=1 | x; \theta)$$



Logistic Regression Model

Gradient descent

Gradient Descent

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

Want $\min_{\theta} J(\theta)$:

Repeat {

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

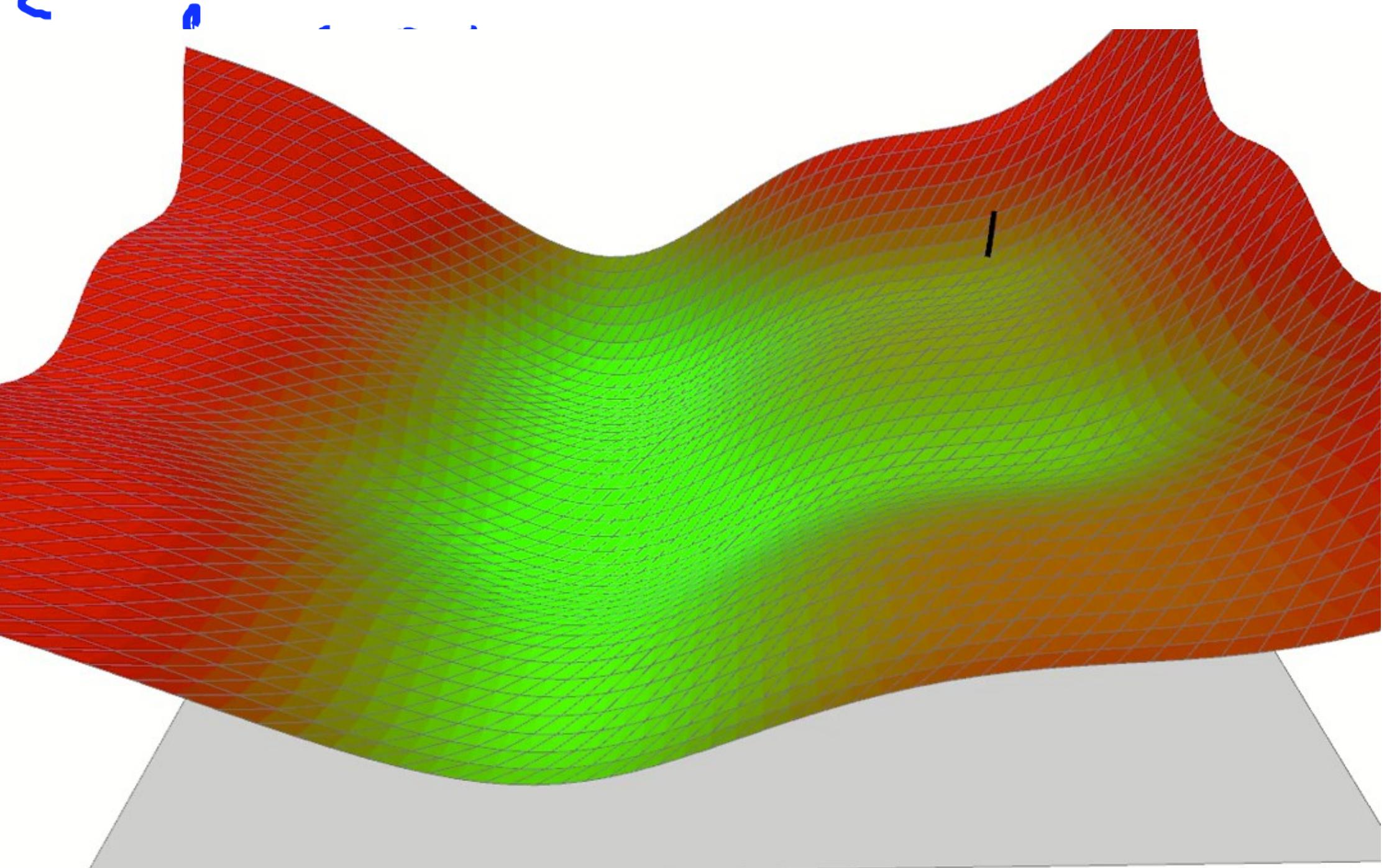
}

(simultaneously update all θ_j)

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

Algorithm looks identical to linear regression!

$$\Theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix}$$



Logistic Regression Model

Advanced optimization

Optimization algorithm

Given θ , we have code that can compute

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

(for $j = 0, 1, \dots, n$)

Optimization algorithms:

- - Gradient descent
- Conjugate gradient
- BFGS
- L-BFGS

Advantages:

- No need to manually pick α
- Often faster than gradient descent.

Disadvantages:

- More complex

Logistic Regression Model

Advanced optimization

Example:

$$\rightarrow \theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$$

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

$$\theta_1 = 5, \theta_2 = 5.$$

$$\rightarrow J(\theta) = (\theta_1 - 5)^2 + (\theta_2 - 5)^2$$

$$\rightarrow \frac{\partial}{\partial \theta_1} J(\theta) = 2(\theta_1 - 5)$$

$$\rightarrow \frac{\partial}{\partial \theta_2} J(\theta) = 2(\theta_2 - 5)$$

\rightarrow options = optimset('GradObj', 'on', 'MaxIter', '100');

\rightarrow initialTheta = zeros(2,1);

$\boxed{[\text{optTheta}, \text{functionVal}, \text{exitFlag}] \dots}$

$= \text{fminunc}(\text{red box}, \underline{\text{initialTheta}}, \text{options});$

↑ ↑

$\Theta \in \mathbb{R}^d \quad d \geq 2.$

We need code to compute

```
function [jVal, gradient] = costFunction(theta)
    jVal = (theta(1)-5)^2 + ...
            (theta(2)-5)^2;
    gradient = zeros(2,1);
    gradient(1) = 2*(theta(1)-5);
    gradient(2) = 2*(theta(2)-5);
```

Logistic Regression Model

Advanced optimization

```
def fmin_l_bfgs_b(func, x0, fprime=None, args=(),
                   approx_grad=0,
                   bounds=None, m=10, factr=1e7, pgtol=1e-5,
                   epsilon=1e-8,
                   iprint=-1, maxfun=15000, maxiter=15000, disp=None,
                   callback=None):
    """
```

Minimize a function func using the L-BFGS-B algorithm.

https://docs.scipy.org/doc/scipy-0.14.0/reference/generated/scipy.optimize.fmin_l_bfgs_b.html

<https://github.com/scipy/scipy/blob/v0.14.0/scipy/optimize/lbfgsb.py#L47>

Logistic Regression

Logistic Regression

Logistic regression fits a logistic model to data and makes predictions about the probability of an event (between 0 and 1).

This recipe shows the fitting of a logistic regression algorithm to the iris dataset. Because this is a multiclass classification problem and logistic regression makes predictions between 0 and 1, a one-vs-all scheme is used (one model is created per class).

```
# Logistic Regression
from sklearn import datasets
from sklearn import metrics
from sklearn.linear_model import LogisticRegression
# load the iris datasets
dataset = datasets.load_iris()
# fit a logistic regression model to the data
model = LogisticRegression()
model.fit(dataset.data, dataset.target)
print(model)
# make predictions
expected = dataset.target
predicted = model.predict(dataset.data)
# summarize the fit of the model
print(metrics.classification_report(expected, predicted))
print(metrics.confusion_matrix(expected, predicted))
```

Multiclass Classification

Multiclass classification

Multiclass classification

Email foldering/tagging: Work, Friends, Family, Hobby

$y=1$ $y=2$ $y=3$ $y=4$

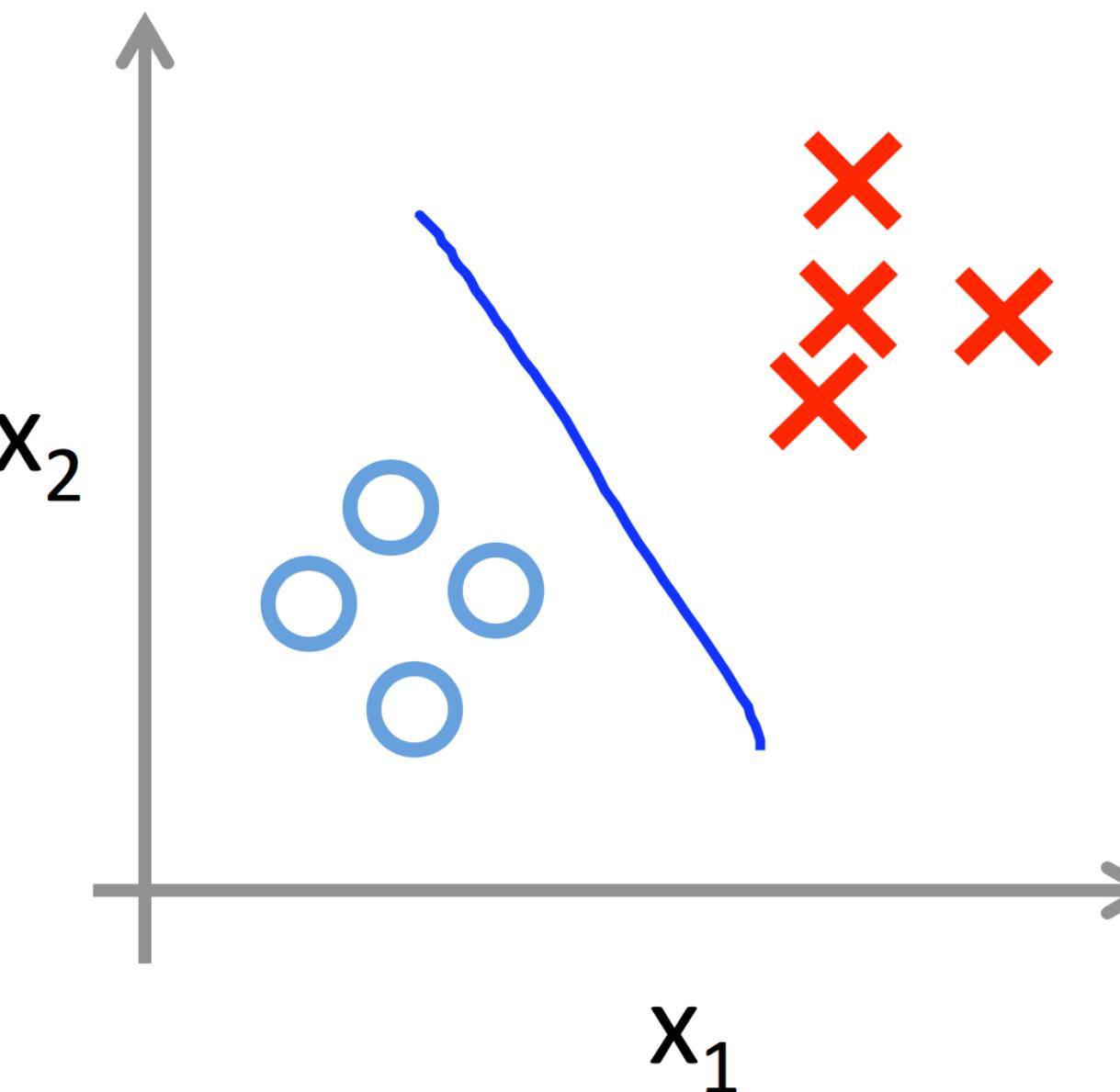
Medical diagrams: Not ill, Cold, Flu

$y=1$ 2 3

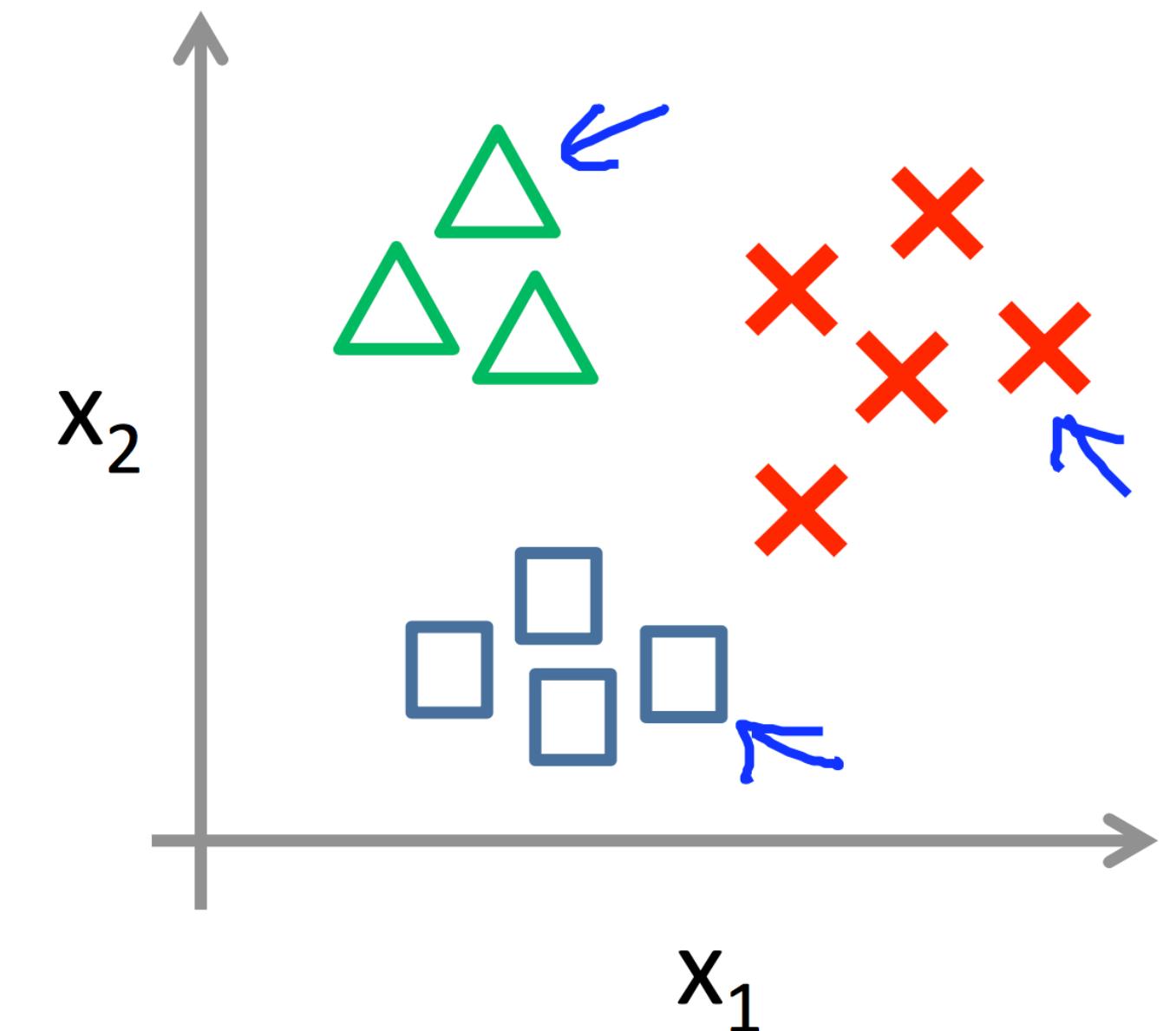
Weather: Sunny, Cloudy, Rain, Snow

$y=1$ 2 3 4

Binary classification:



Multi-class classification:



One vs. All

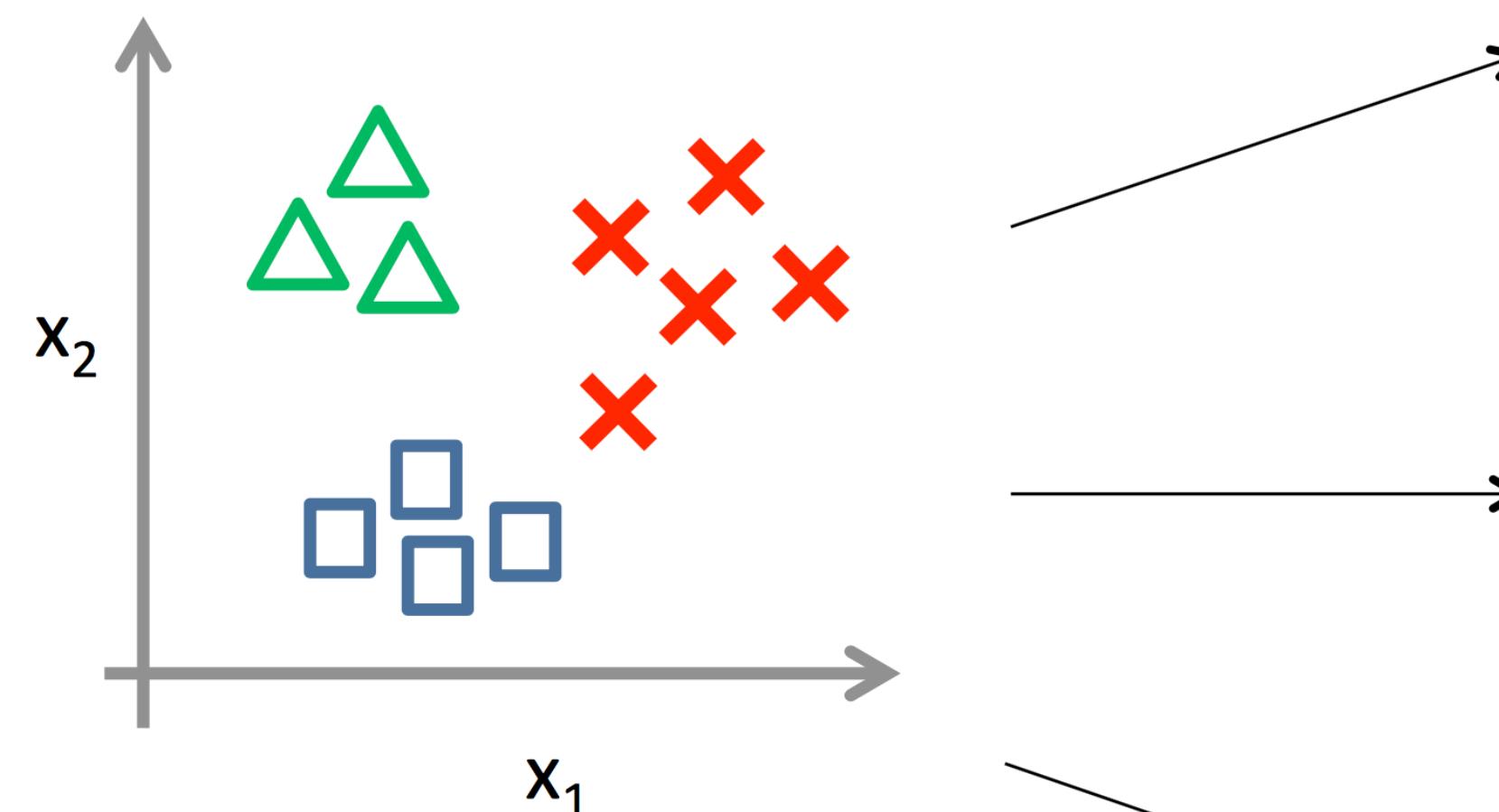
Multiclass Classification

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

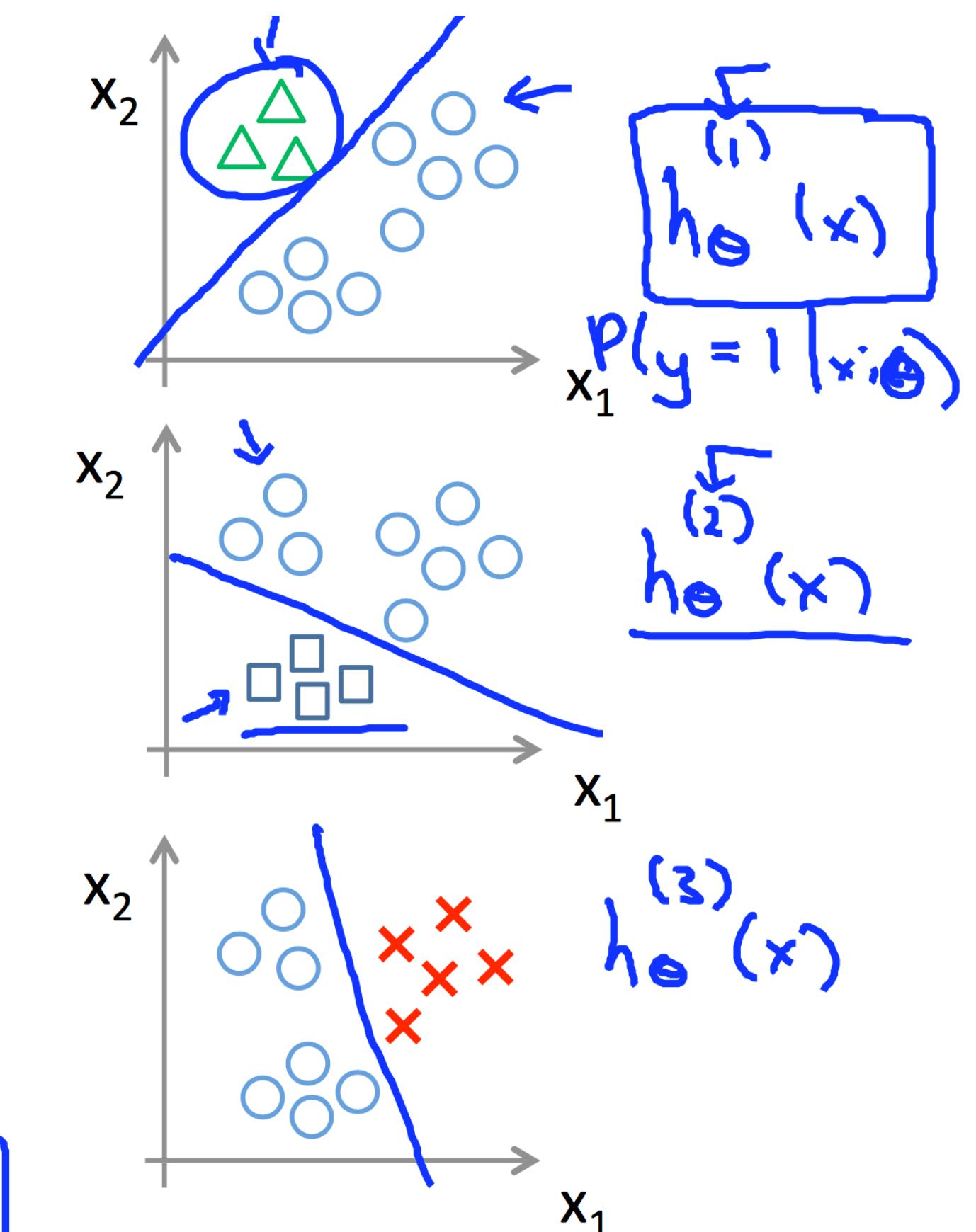
On a new input x , to make a prediction, pick the class i that maximizes

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

One-vs-all (one-vs-rest):



$$h_{\theta}^{(i)}(x) = P(y = i|x; \theta) \quad (i = 1, 2, 3)$$



Linear Discriminant Analysis

Linear Discriminant Analysis

Linear Discriminate Analysis (LDA) fits a conditional probability density function (Gaussian) to each attribute for the class, the discrimination function between the classes is linear.

This recipe shows the fitting of an LDA algorithm to the iris dataset.

```
# Linear Discriminant Analysis
from sklearn import datasets
from sklearn import metrics
from sklearn lda import LDA
# load the iris datasets
dataset = datasets.load_iris()
# fit a LDA model to the data
model = LDA()
model.fit(dataset.data, dataset.target)
print(model)
# make predictions
expected = dataset.target
predicted = model.predict(dataset.data)
# summarize the fit of the model
print(metrics.classification_report(expected, predicted))
print(metrics.confusion_matrix(expected, predicted))
```

Quadratic Discriminant Analysis

Quadratic Discriminant Analysis

Quadratic Discriminant Analysis (QDA) fits a conditional probability density function (Gaussian) to each attribute for the class, the discrimination function between the classes is quadratic.

This recipe shows the fitting of an QDA algorithm to the iris dataset.

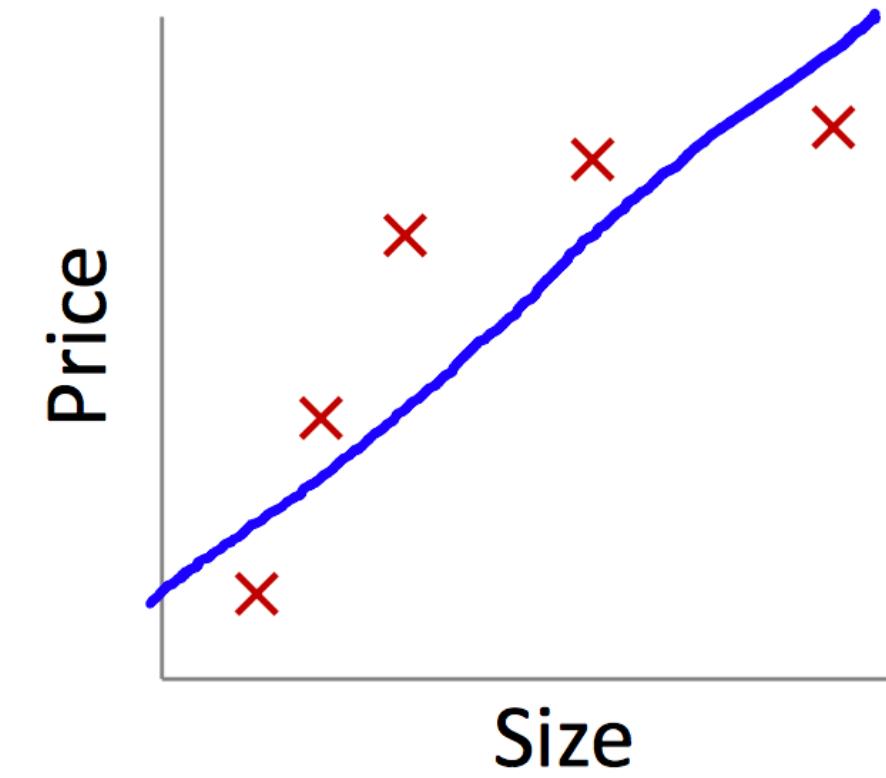
```
# Quadratic Discriminant Analysis
from sklearn import datasets
from sklearn import metrics
from sklearn.qda import QDA
# load the iris datasets
dataset = datasets.load_iris()
# fit a QDA model to the data
model = QDA()
model.fit(dataset.data, dataset.target)
print(model)
# make predictions
expected = dataset.target
predicted = model.predict(dataset.data)
# summarize the fit of the model
print(metrics.classification_report(expected, predicted))
print(metrics.confusion_matrix(expected, predicted))
```

Regularization

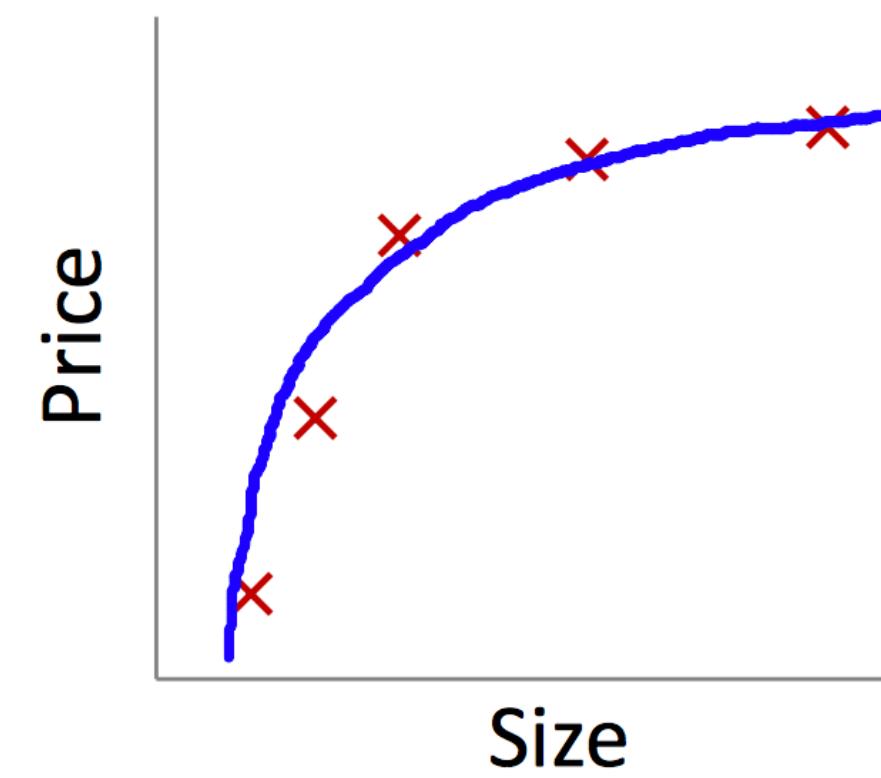
Machine learning models need to generalize well to new examples that the model has not seen in practice. Regularization helps prevent models from overfitting the training data.

The problem of overfitting

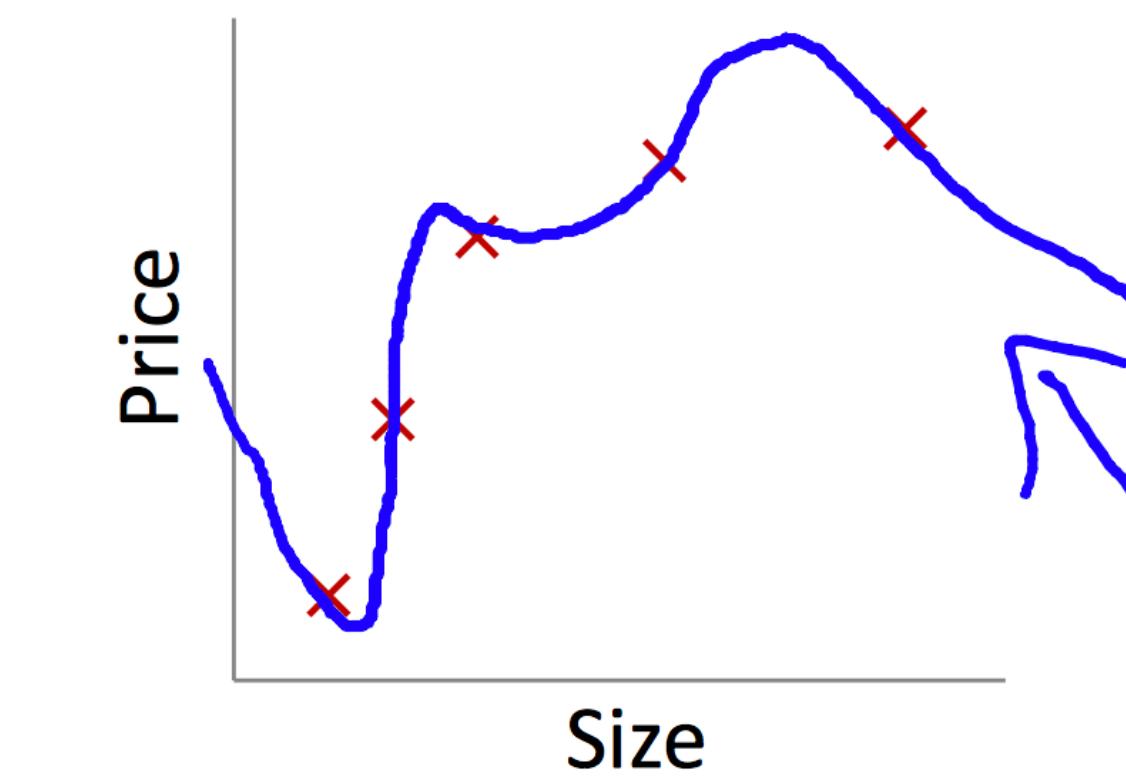
Example: Linear regression (housing prices)



$\rightarrow \theta_0 + \theta_1 x$
"Underfit" "High bias"



$\rightarrow \theta_0 + \theta_1 x + \theta_2 x^2$
"Just right"



$\rightarrow \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$
"Overfit" "High variance"

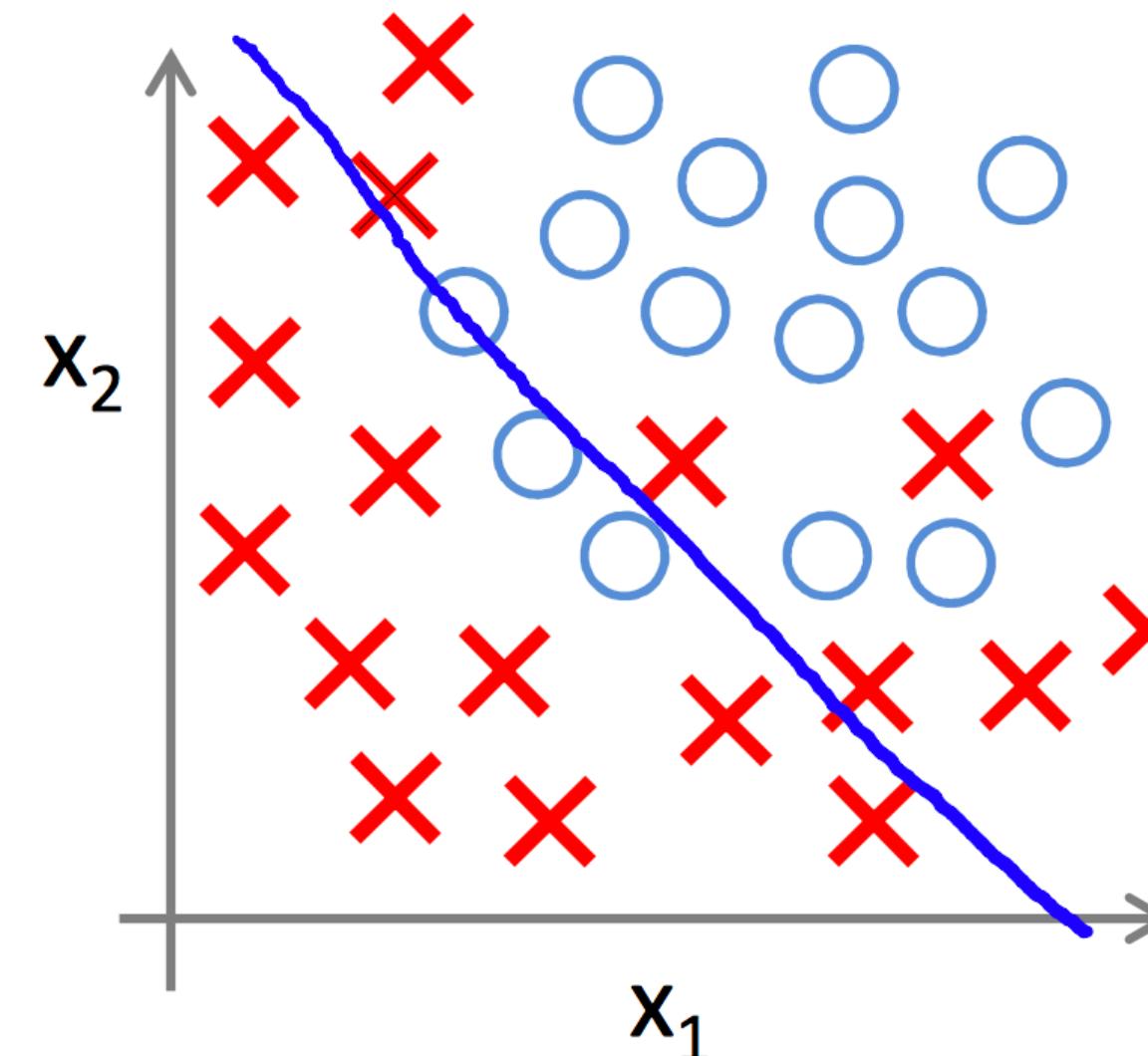
Overfitting: If we have too many features, the learned hypothesis may fit the training set very well ($J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 \approx 0$), but fail to generalize to new examples (predict prices on new examples).

Regularization

Machine learning models need to generalize well to new examples that the model has not seen in practice. Regularization helps prevent models from overfitting the training data.

The problem of overfitting

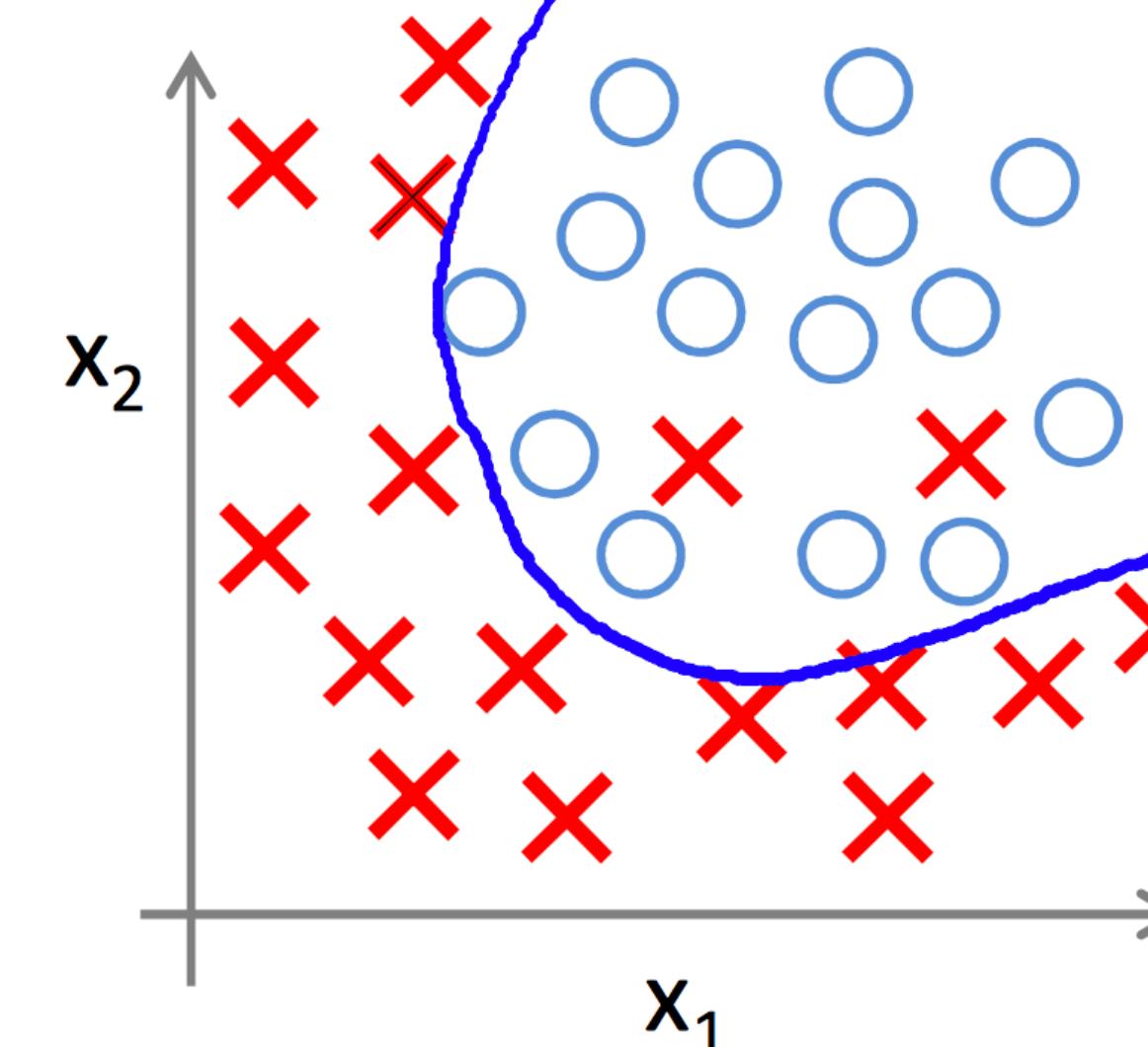
Example: Logistic regression



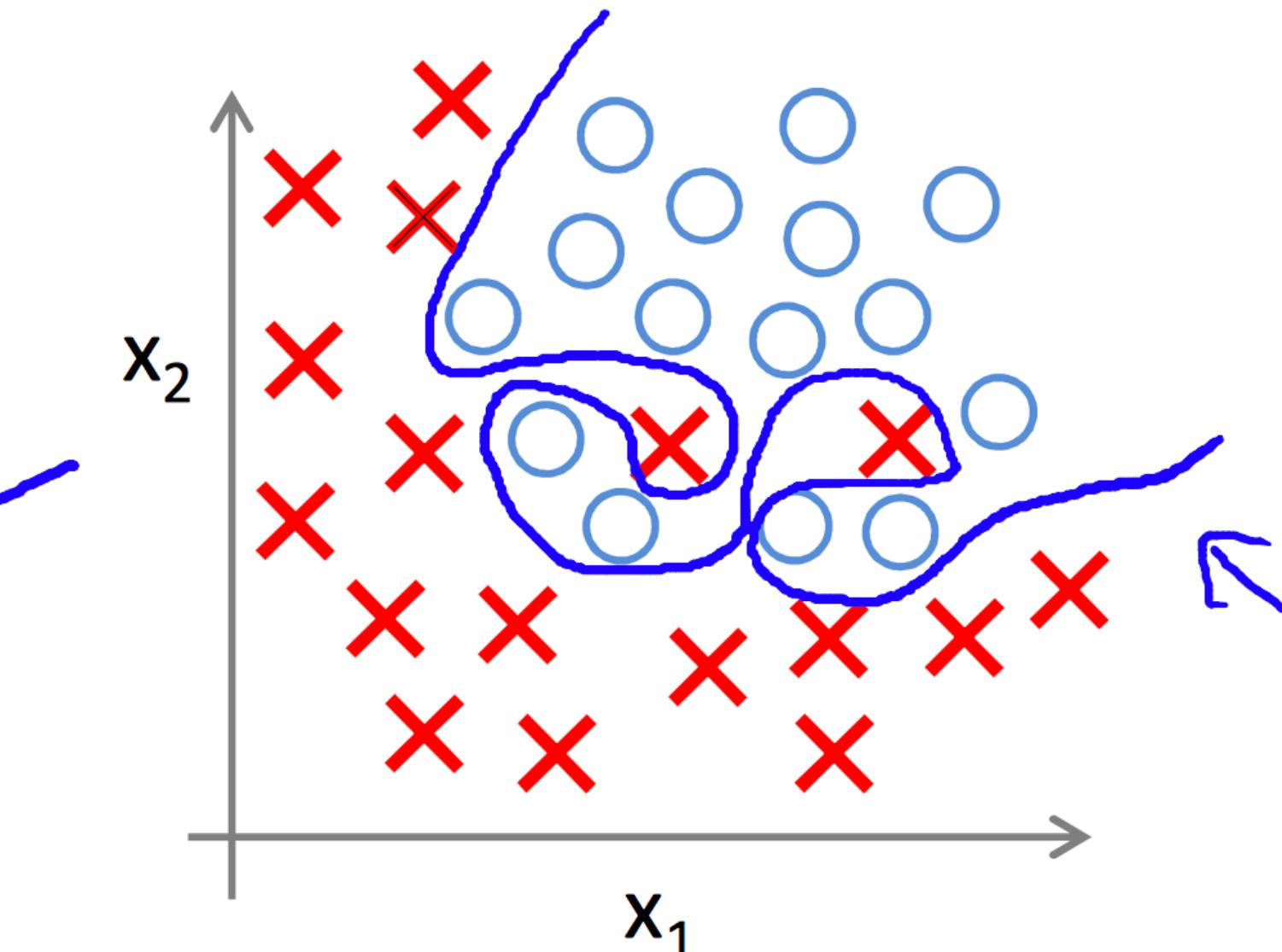
$$h_{\theta}(x) = g(\underline{\theta_0 + \theta_1 x_1 + \theta_2 x_2})$$

(g = sigmoid function)

"Underfit"



$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 \underline{x_1 x_2})$$



$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_1^2 + \theta_3 x_1^2 x_2 + \theta_4 \underline{x_1^2 x_2^2} + \theta_5 \underline{x_1^2 x_2^3} + \theta_6 \underline{x_1^3 x_2} + \dots)$$

"Overfit"

Regularization

Addressing Overfitting

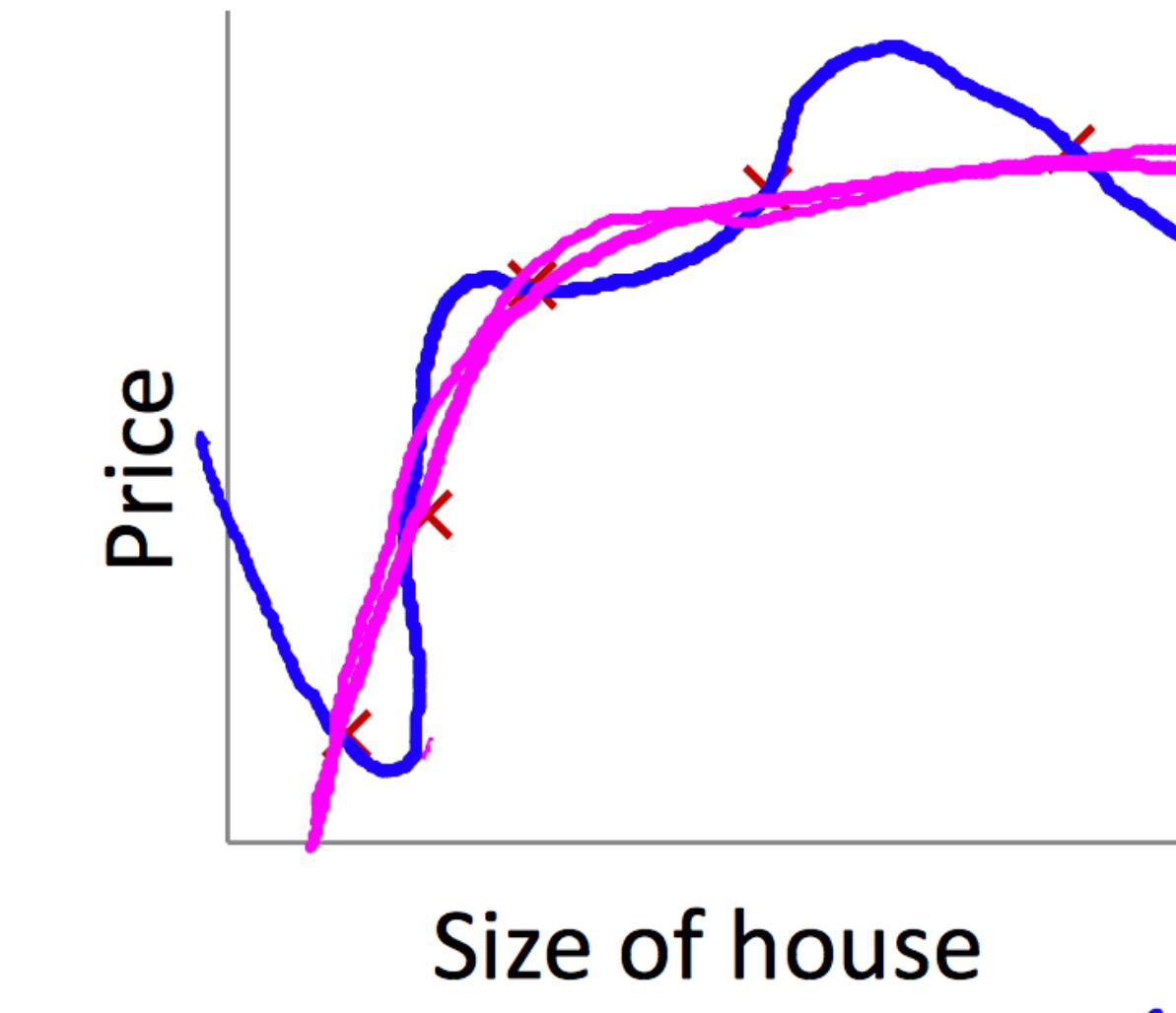
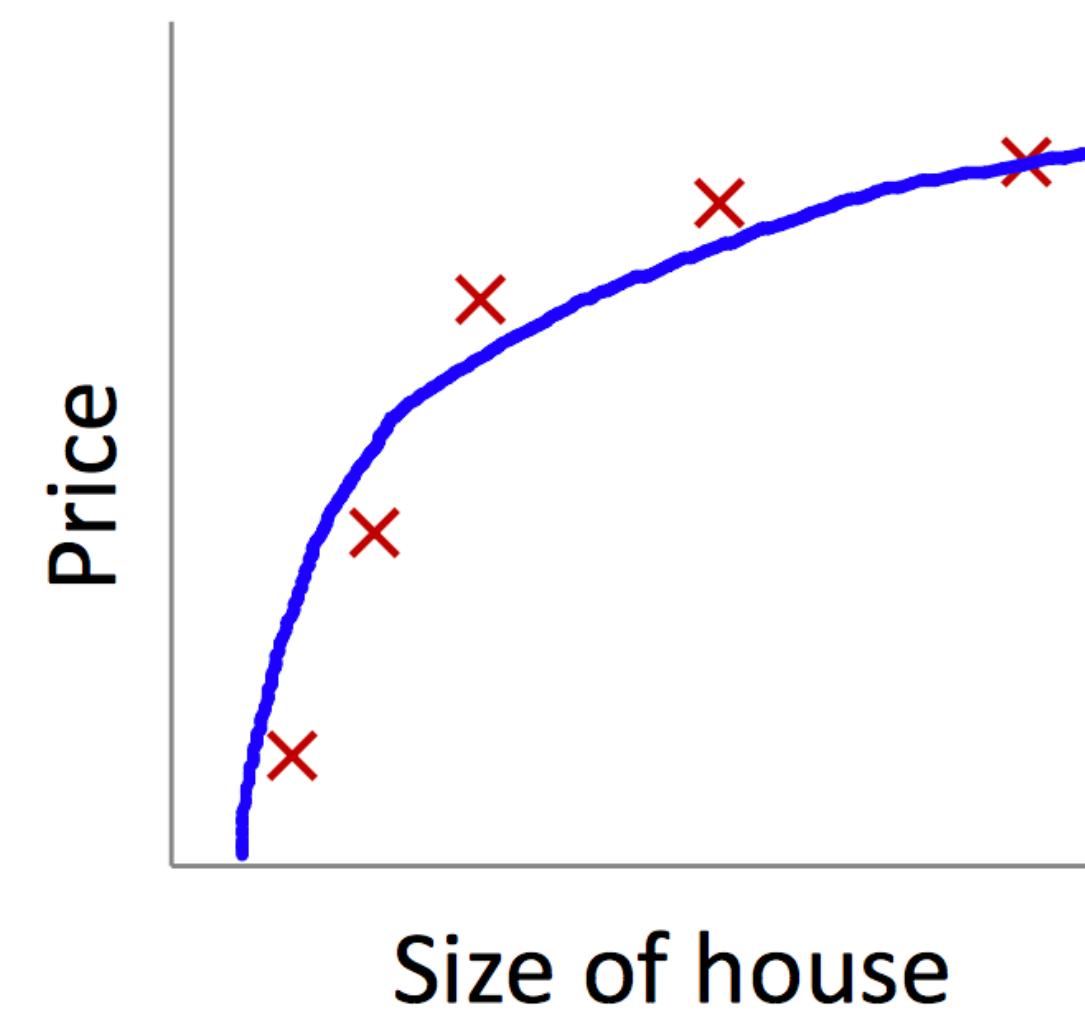
Options:

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

Regularization

Cost function

Intuition



Suppose we penalize and make θ_3, θ_4 really small.

$$\rightarrow \min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + 1000 \underline{\theta_3^2} + 1000 \underline{\theta_4^2}$$

$\theta_3 \approx 0$ $\theta_4 \approx 0$

Regularization

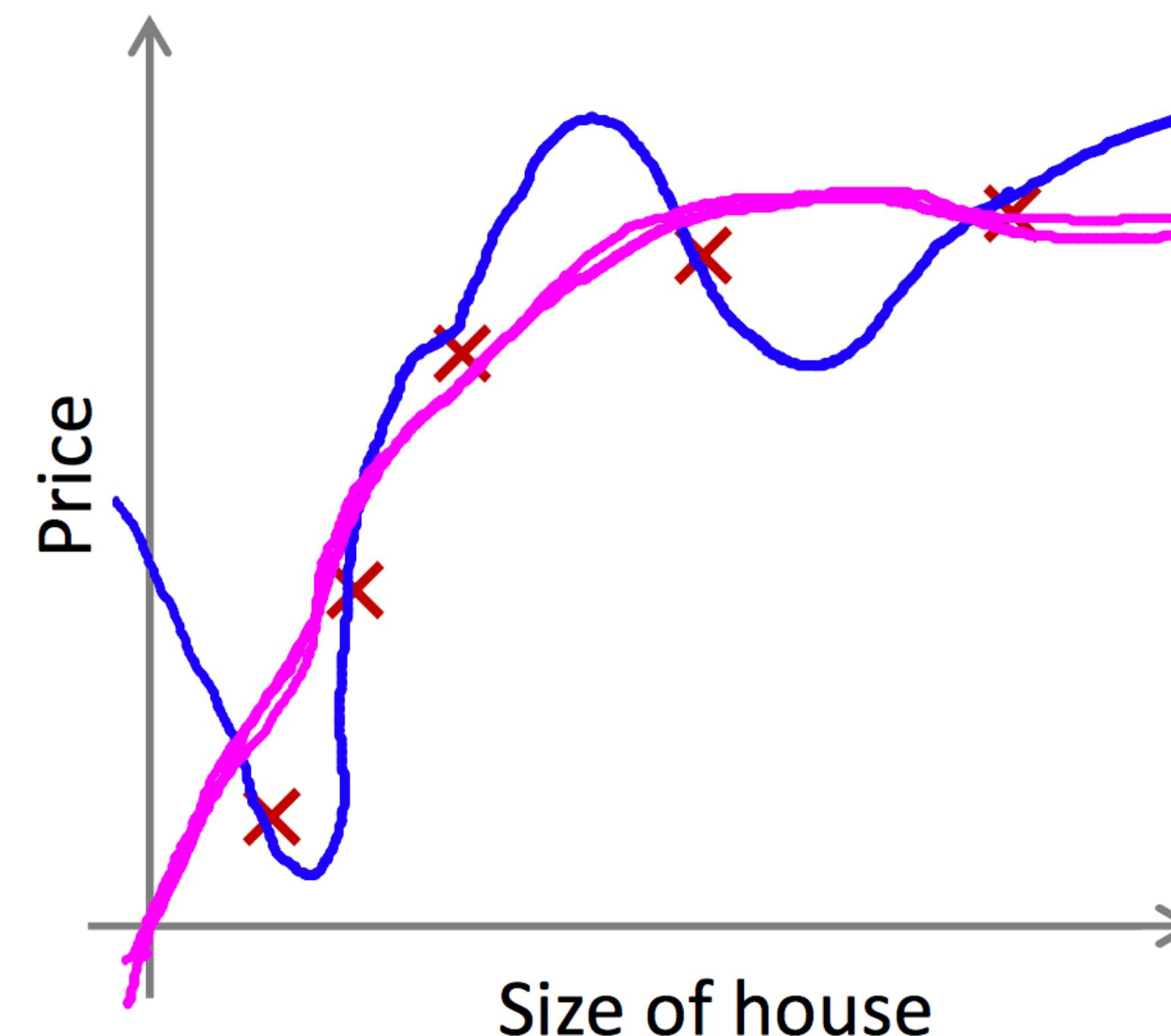
Cost function

Regularization.

$$\rightarrow J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right]$$

min $J(\theta)$

regularization parameter



Regularization

In regularized linear regression, we choose θ to minimize

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

What if λ is set to an extremely large value (perhaps for too large for our problem, say $\lambda = 10^{10}$)?

- Algorithm works fine; setting λ to be very large can't hurt it
- Algorithm fails to eliminate overfitting.
- Algorithm results in underfitting. (Fails to fit even training data well).
- Gradient descent will fail to converge.

Regularization

Regularized Linear Regression

Gradient descent

Repeat {

$$\rightarrow \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\frac{\partial J(\theta)}{\partial \theta_0}$$

$$\theta_1, \theta_2, \dots, \theta_n$$

$$\frac{\partial J(\theta)}{\partial \theta_j}$$

$$\begin{aligned} \theta_j &:= \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right] \\ &\quad (j = \cancel{0}, 1, 2, 3, \dots, n) \\ \theta_j &:= \theta_j (1 - \alpha \frac{\lambda}{m}) - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \end{aligned}$$

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

$$\frac{0.99}{\theta_j \times 0.99}$$

should be invertible

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

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

Normal equation

$$X = \begin{bmatrix} (x^{(1)})^T \\ \vdots \\ (x^{(m)})^T \end{bmatrix} \quad m \times (n+1)$$

$$\rightarrow \min_{\theta} J(\theta)$$

$$\theta = (X^T X + \lambda \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ 0 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{bmatrix})^{-1} X^T y$$

$$y = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix} \quad \mathbb{R}^m$$

** if non-invertible/singular use `>> numpy.linalg.pinv` (Moore-Penrose Matrix Inverse)
<https://docs.scipy.org/doc/numpy-1.13.0/reference/generated/numpy.linalg.pinv.html>

Regularization

Regularized Logistic Regression

Gradient descent

Repeat {

$$\Rightarrow \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\rightarrow \theta_j := \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right] \quad \leftarrow$$

$(j = \cancel{0}, 1, 2, 3, \dots, n)$

$\theta_0, \dots, \theta_n$

$$\frac{\partial}{\partial \theta_j} J(\theta)$$

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

Regularization

We need code to compute

Advanced Optimization for Regularization

Advanced optimization

→ **function** [jVal, gradient] = costFunction(theta)

jVal = [code to compute $J(\theta)$];

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

→ gradient (1) = [code to compute $\frac{\partial}{\partial \theta_0} J(\theta)$];

$$\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

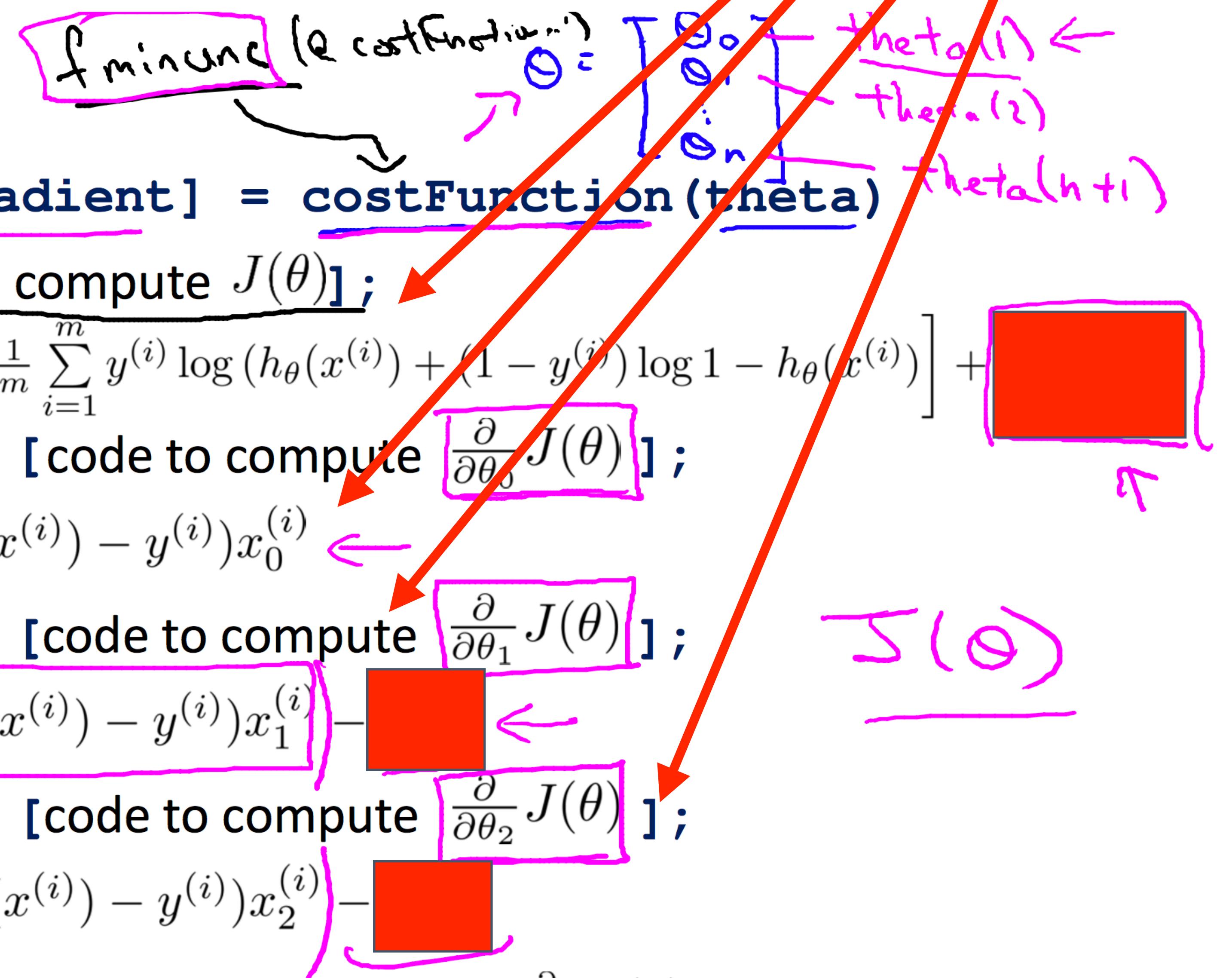
→ gradient (2) = [code to compute $\frac{\partial}{\partial \theta_1} J(\theta)$];

$$\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_1^{(i)}$$

→ gradient (3) = [code to compute $\frac{\partial}{\partial \theta_2} J(\theta)$];

$$\vdots \quad \left(\frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_2^{(i)} \right) - \text{[red box]}$$

gradient (n+1) = [code to compute $\frac{\partial}{\partial \theta_n} J(\theta)$];



Regularization

Ridge Regression

Ridge Regression (also known as Tikhonov regularization) penalizes a least squares regression model on the square of the absolute magnitude of the coefficients (called the L2 norm).

This recipe shows use of the Ridge Regression algorithm to make predictions for the diabetes dataset.

```
# Ridge Regression
import numpy as np
from sklearn import datasets
from sklearn.linear_model import Ridge
# load the diabetes datasets
dataset = datasets.load_diabetes()
# fit a ridge regression model to the data
model = Ridge(alpha=0.1)
model.fit(dataset.data, dataset.target)
print(model)
# make predictions
expected = dataset.target
predicted = model.predict(dataset.data)
# summarize the fit of the model
mse = np.mean((predicted-expected)**2)
print(mse)
print(model.score(dataset.data, dataset.target))
```

L1-norm is also known as least absolute deviations (LAD), least absolute errors (LAE). It is basically minimizing the sum of the absolute differences (S) between the target value (\mathbf{Y}_i) and the estimated values ($f(\mathbf{x}_i)$):

$$S = \sum_{i=1}^n |y_i - f(x_i)|.$$

Regularization

Least Absolute Shrinkage and Selection Operator

The Least Absolute Shrinkage and Selection Operator (LASSO) is a regularization method that penalizes a least squares regression model on the absolute magnitude of the coefficients (called the L1 norm).

This recipe shows use of the LASSO algorithm to make predictions for the diabetes dataset.

```
# Lasso Regression
import numpy as np
from sklearn import datasets
from sklearn.linear_model import Lasso
# load the diabetes datasets
dataset = datasets.load_diabetes()
# fit a LASSO model to the data
model = Lasso(alpha=0.1)
model.fit(dataset.data, dataset.target)
print(model)
# make predictions
expected = dataset.target
predicted = model.predict(dataset.data)
# summarize the fit of the model
mse = np.mean((predicted-expected)**2)
print(mse)
print(model.score(dataset.data, dataset.target))
```



L2-norm is also known as least squares. It is basically minimizing the sum of the square of the differences (S) between the target value (Y_i) and the estimated values ($f(x_i)$):

$$S = \sum_{i=1}^n (y_i - f(x_i))^2$$

Regularization

Least Angle Regression

The Least Angle Regression (LARS) method is a computationally efficient algorithm for fitting a regression model. It is useful for high-dimensional data and is commonly used in conjunction with regularization (such as LASSO).

This recipe shows use of the LARS algorithm to make predictions for the diabetes dataset.

```
# LARS Regression
import numpy as np
from sklearn import datasets
from sklearn.linear_model import Lars
# load the diabetes datasets
dataset = datasets.load_diabetes()
# fit a LARS model to the data
model = Lars()
model.fit(dataset.data, dataset.target)
print(model)
# make predictions
expected = dataset.target
predicted = model.predict(dataset.data)
# summarize the fit of the model
mse = np.mean((predicted-expected)**2)
print(mse)
print(model.score(dataset.data, dataset.target))
```

Regularization

LASSO LARS

The Least Angle Regression (LARS) can be used as an alternative method for calculating Least Absolute Shrinkage and Selection Operator (LASSO) fit.

This recipe shows use of the LASSO with LARS algorithm to make predictions for the diabetes dataset.

```
# LassoLars Regression
import numpy as np
from sklearn import datasets
from sklearn.linear_model import LassoLars
# load the iris datasets
dataset = datasets.load_diabetes()
# fit a LASSO using LARS model to the data
model = LassoLars(alpha=0.1)
model.fit(dataset.data, dataset.target)
print(model)
# make predictions
expected = dataset.target
predicted = model.predict(dataset.data)
# summarize the fit of the model
mse = np.mean((predicted-expected)**2)
print(mse)
print(model.score(dataset.data, dataset.target))
```

Regularization

ElasticNet

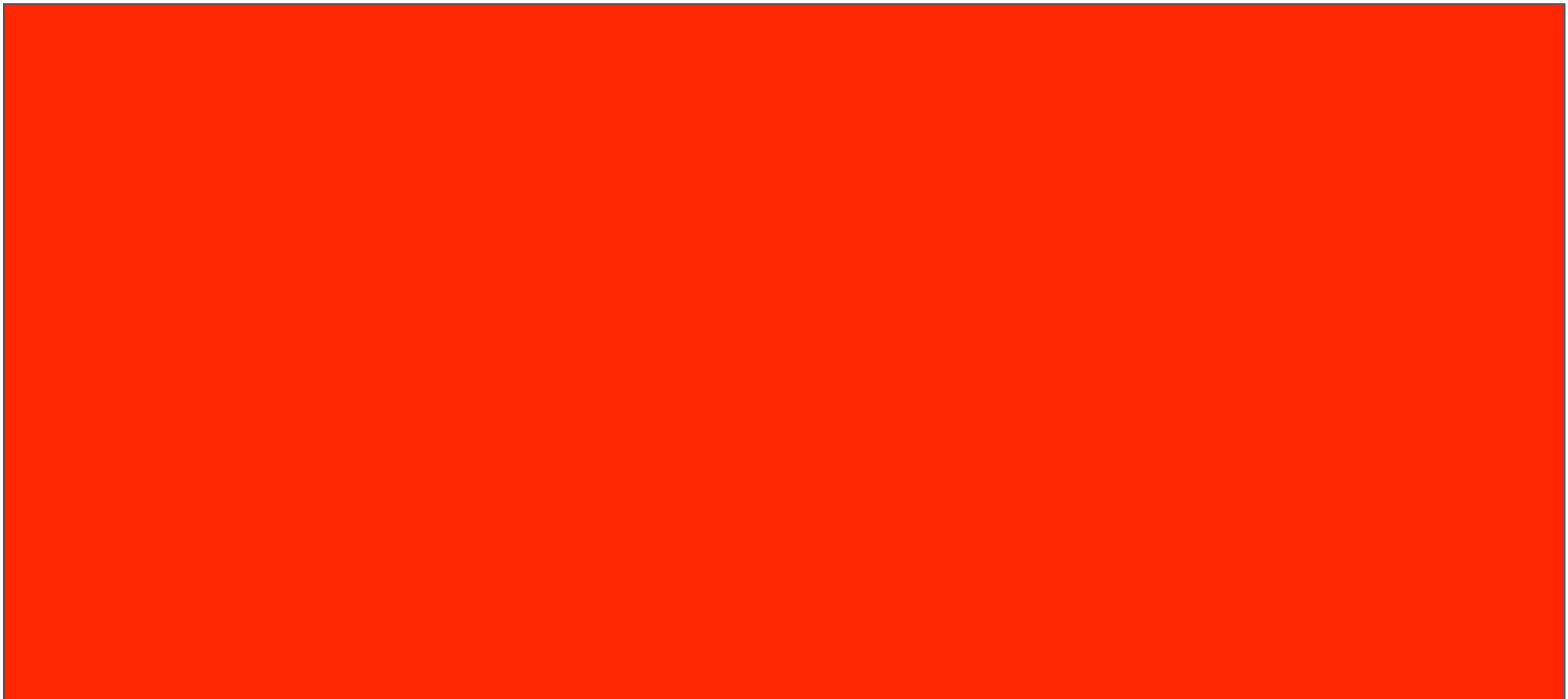
The ElasticNet is a regularized form of regression that penalizes the model with both the L1 norm and the L2 norm.

This recipe shows use of the ElasticNet algorithm to make predictions for the diabetes dataset.

```
# ElasticNet Regression
import numpy as np
from sklearn import datasets
from sklearn.linear_model import ElasticNet
# load the diabetes datasets
dataset = datasets.load_diabetes()
# fit a model to the data
model = ElasticNet(alpha=0.1)
model.fit(dataset.data, dataset.target)
print(model)
# make predictions
expected = dataset.target
predicted = model.predict(dataset.data)
# summarize the fit of the model
mse = np.mean((predicted-expected)**2)
print(mse)
print(model.score(dataset.data, dataset.target))
```

Support Vector Machines

SVMs are considered by many to be the most powerful 'black box' learning algorithm, and by posing a cleverly-chosen optimization objective, one of the most widely used learning algorithms today.



Large Margin Classification

 (x, y)

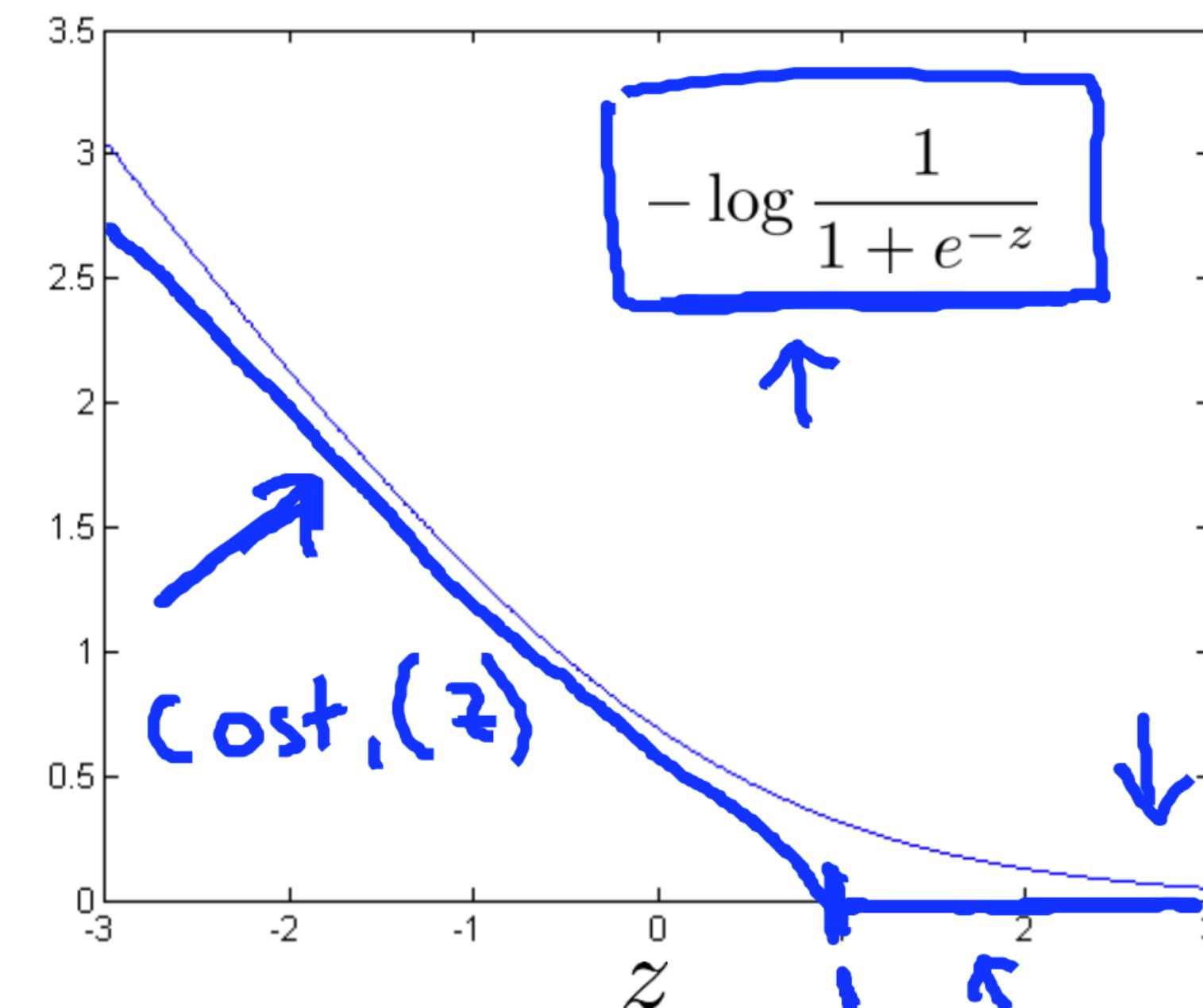
Alternative view of logistic regression

Cost of example: $-(y \log h_\theta(x) + (1 - y) \log(1 - h_\theta(x)))$

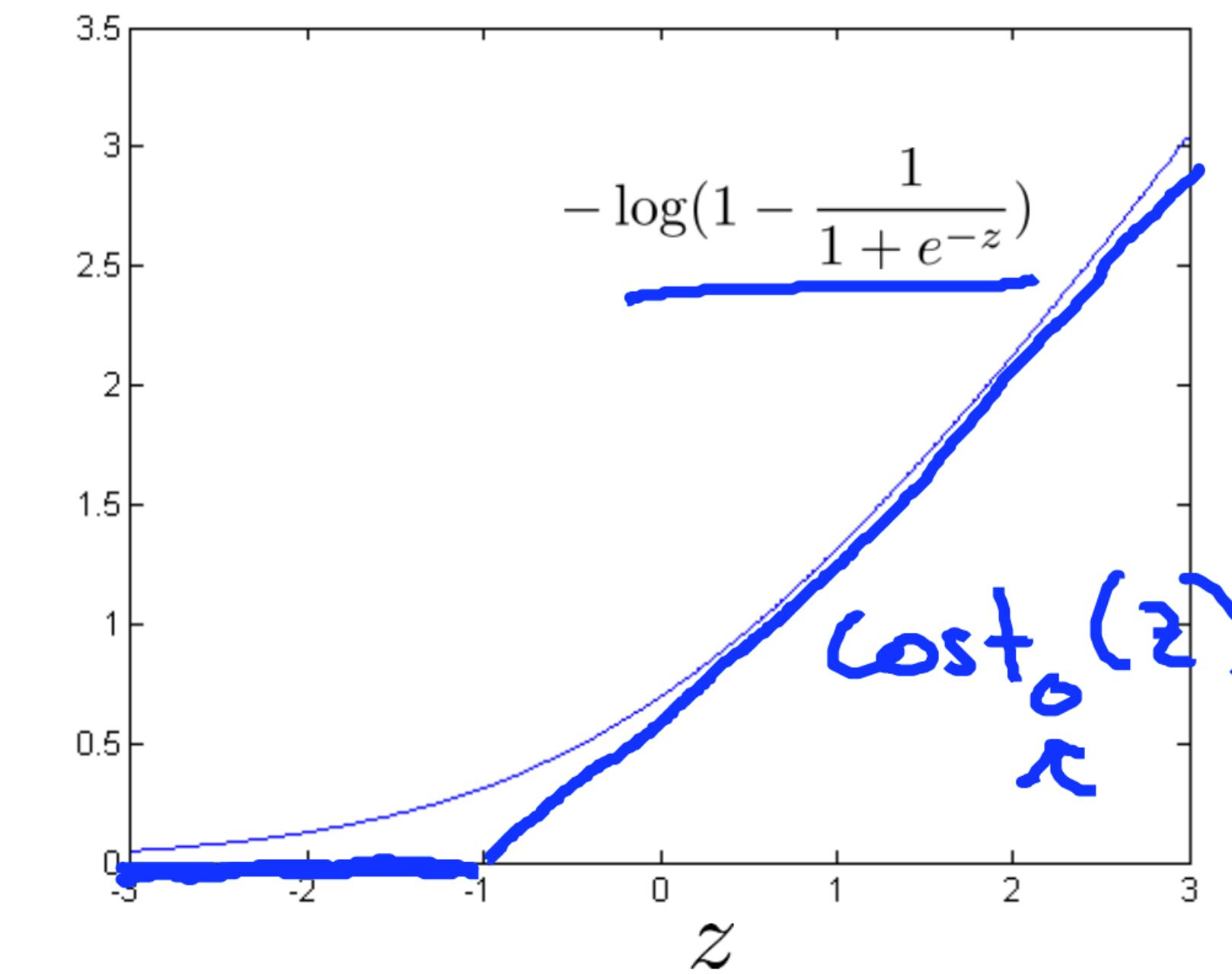
$$= -y \log \frac{1}{1 + e^{-\theta^T x}} - (1 - y) \log \left(1 - \frac{1}{1 + e^{-\theta^T x}}\right)$$

If $y = 1$ (want $\theta^T x \gg 0$):

$$z = \theta^T x$$



If $y = 0$ (want $\theta^T x \ll 0$):



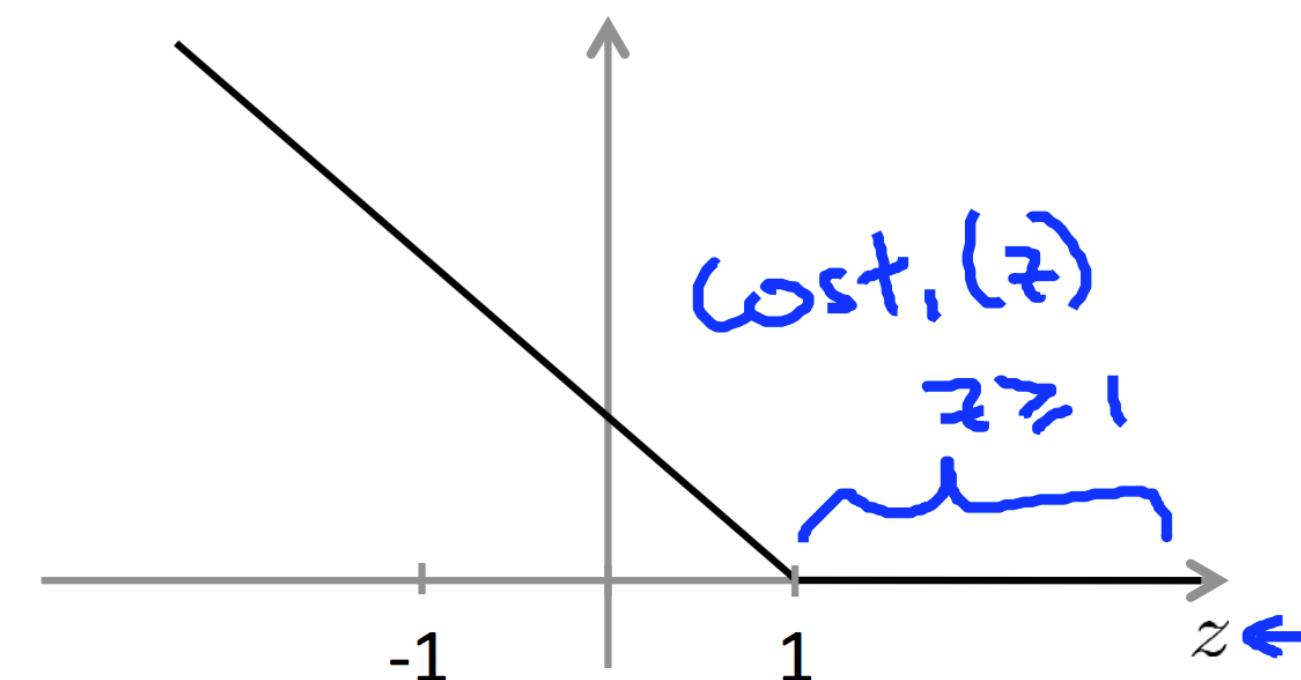
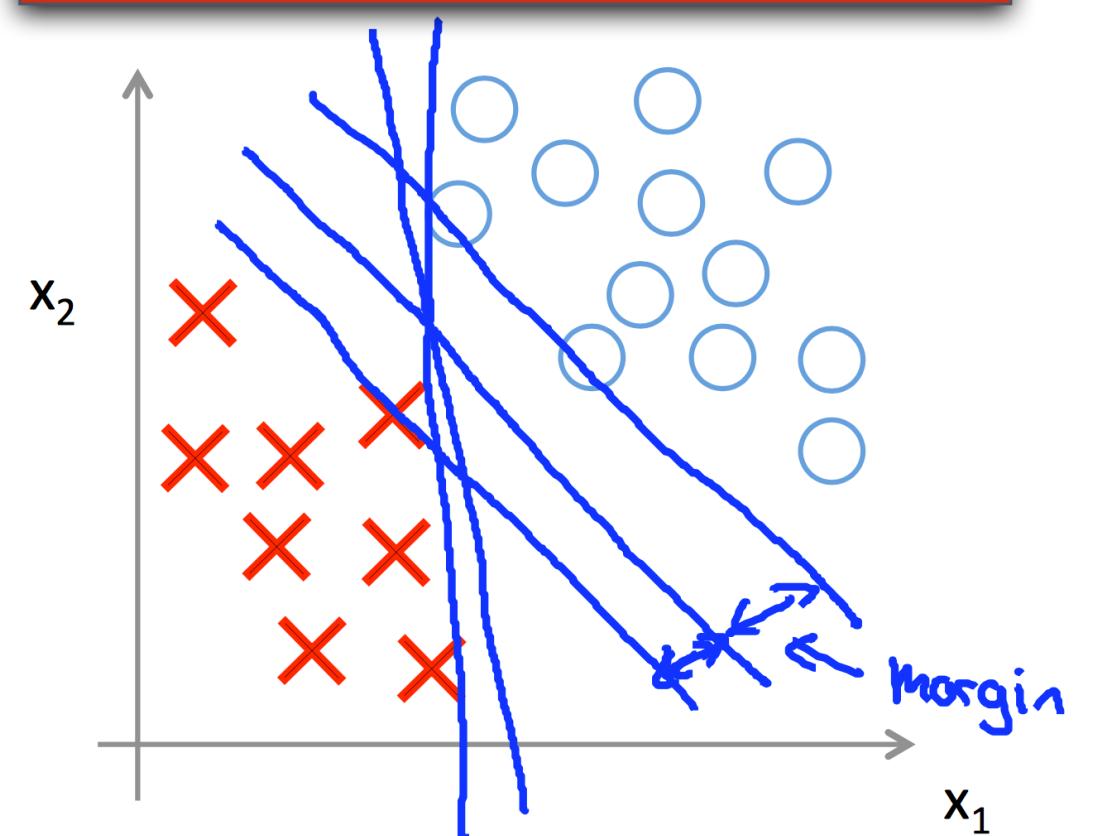
Op

$$\min_{\theta} \frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \left(-\log h_{\theta}(x^{(i)}) \right) + (1 - y^{(i)}) \left((-\log(1 - h_{\theta}(x^{(i)}))) \right) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

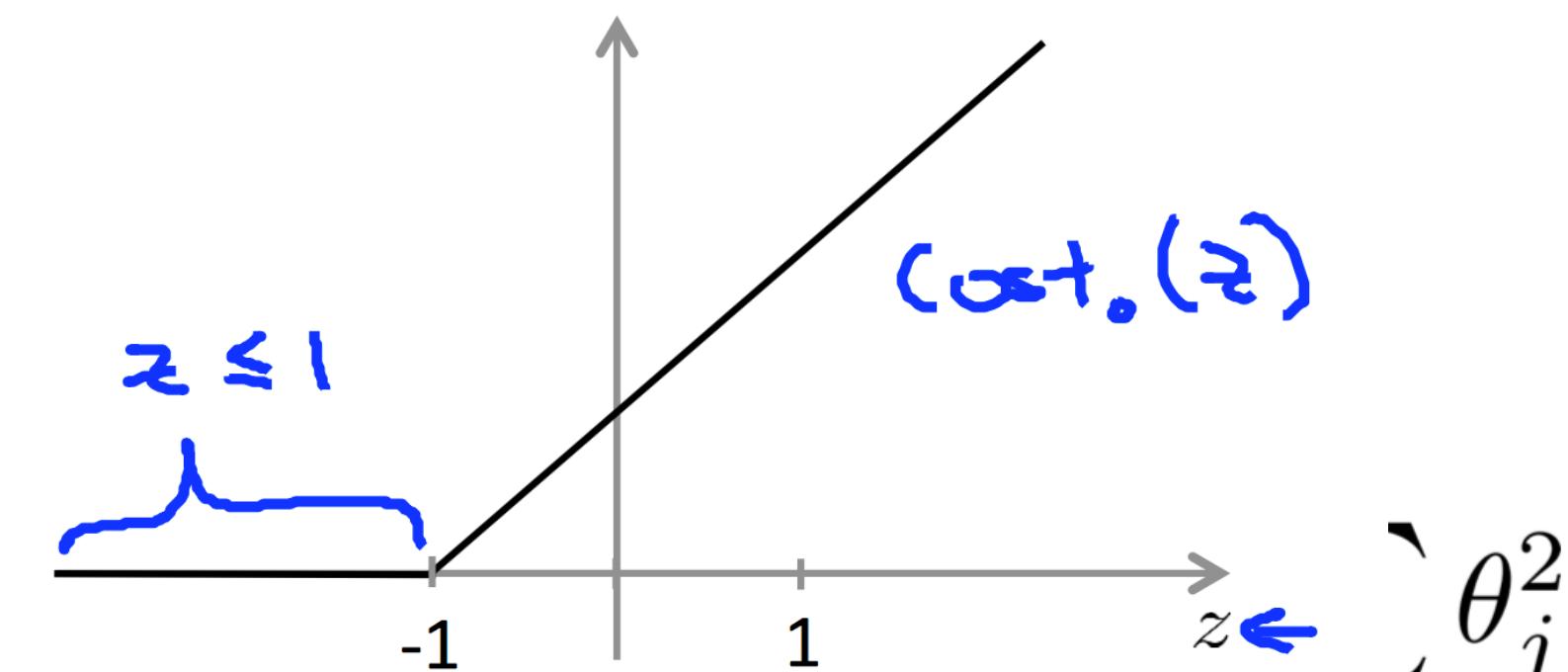
Support vector machine:

Support Vector Machine

Linearly separable cases



If $y = 1$, we want $\theta^T x \geq 1$ (not just ≥ 0)
If $y = 0$, we want $\theta^T x \leq -1$ (not just < 0)



$\theta^T x \geq \alpha - 1$
 $\theta^T x \leq \alpha + 1$

$$\theta_j^2$$

Large Margin Classification

Mathematics behind large margin classification

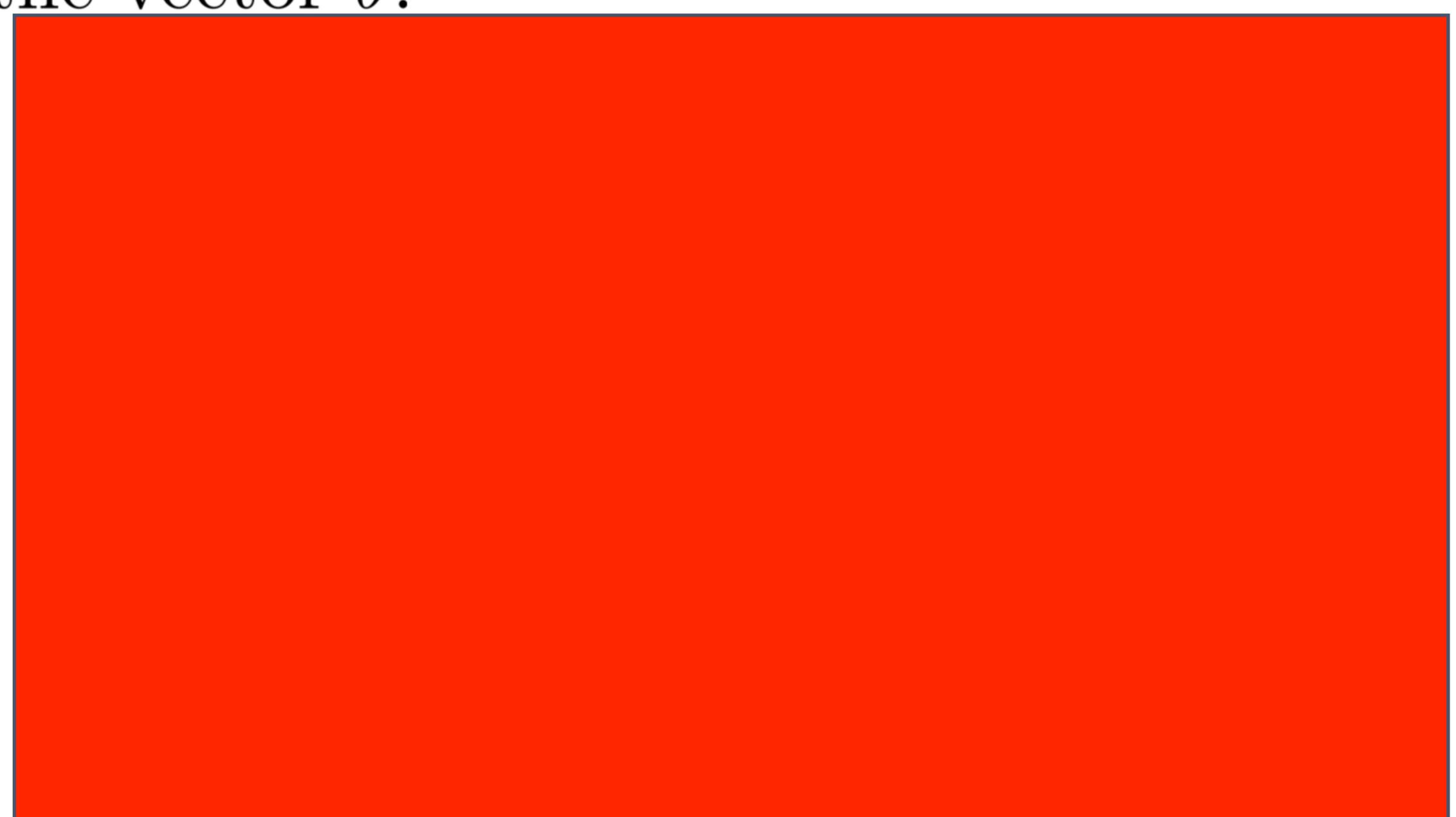
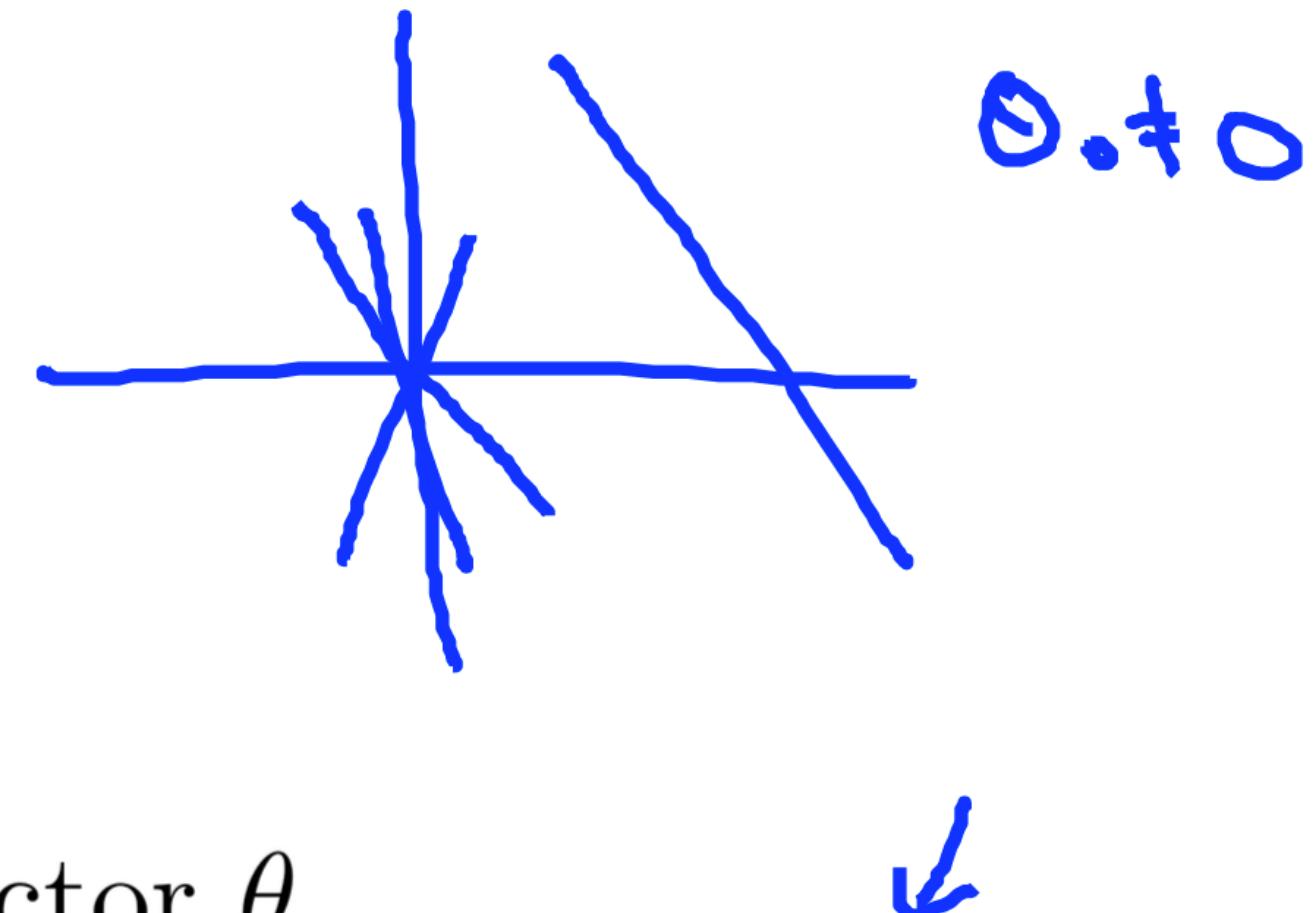
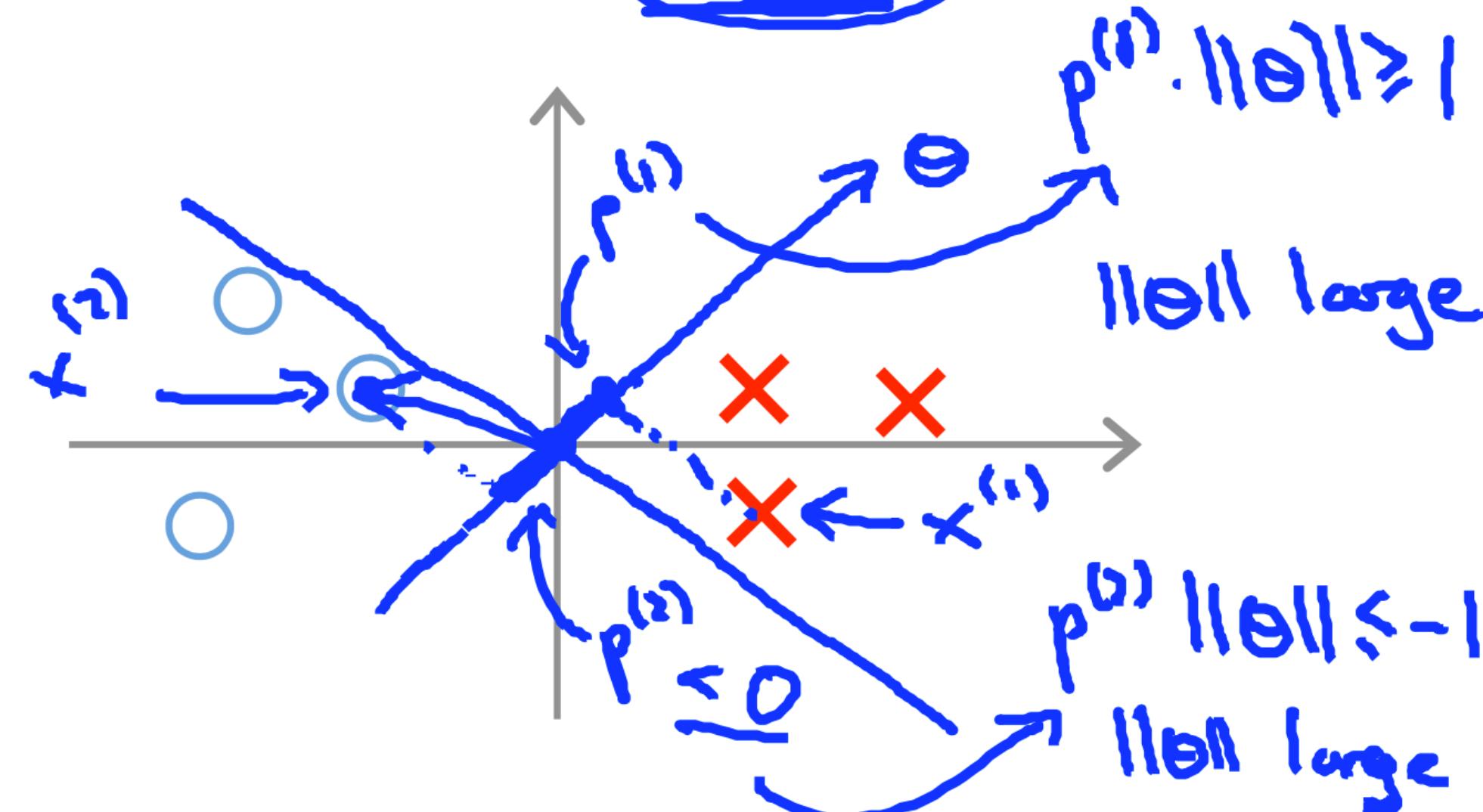
SVM Decision Boundary

$$\min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2 = \frac{1}{2} \|\theta\|^2 \leftarrow$$

$$\text{s.t. } \begin{cases} p^{(i)} \cdot \|\theta\| \geq 1 & \text{if } y^{(i)} = 1 \\ p^{(i)} \cdot \|\theta\| \leq -1 & \text{if } y^{(i)} = -1 \end{cases}$$

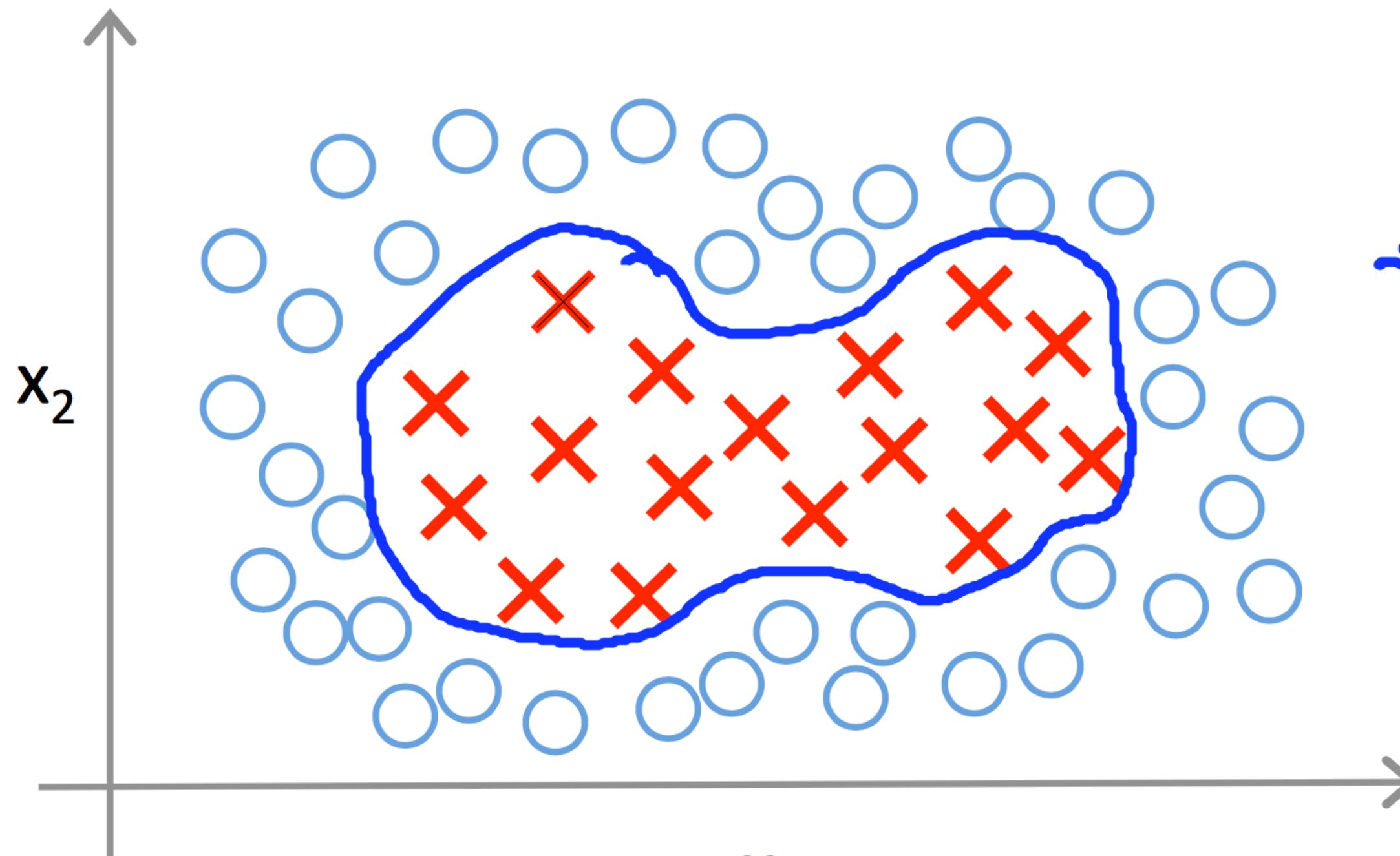
where $p^{(i)}$ is the projection of $x^{(i)}$ onto the vector θ .

Simplification: $\theta_0 = 0$



SVM Kernels

Non-linear Decision Boundary



Predict $y = 1$ if
 $\theta_0 + \theta_1 \underline{x_1} + \theta_2 \underline{x_2} + \theta_3 \underline{x_1 x_2}$
 $+ \theta_4 \underline{x_1^2} + \theta_5 \underline{x_2^2} + \dots \geq 0$

$$h_0(x) = \begin{cases} 1 & \text{if } \theta_0 + \theta_1 x_1 + \dots \geq 0 \\ 0 & \text{otherwise.} \end{cases}$$

$$\Rightarrow \theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 + \dots$$

$$f_1 = x_1, \quad f_2 = x_2, \quad f_3 = x_1 x_2, \quad f_4 = x_1^2, \quad f_5 = x_2^2, \dots$$

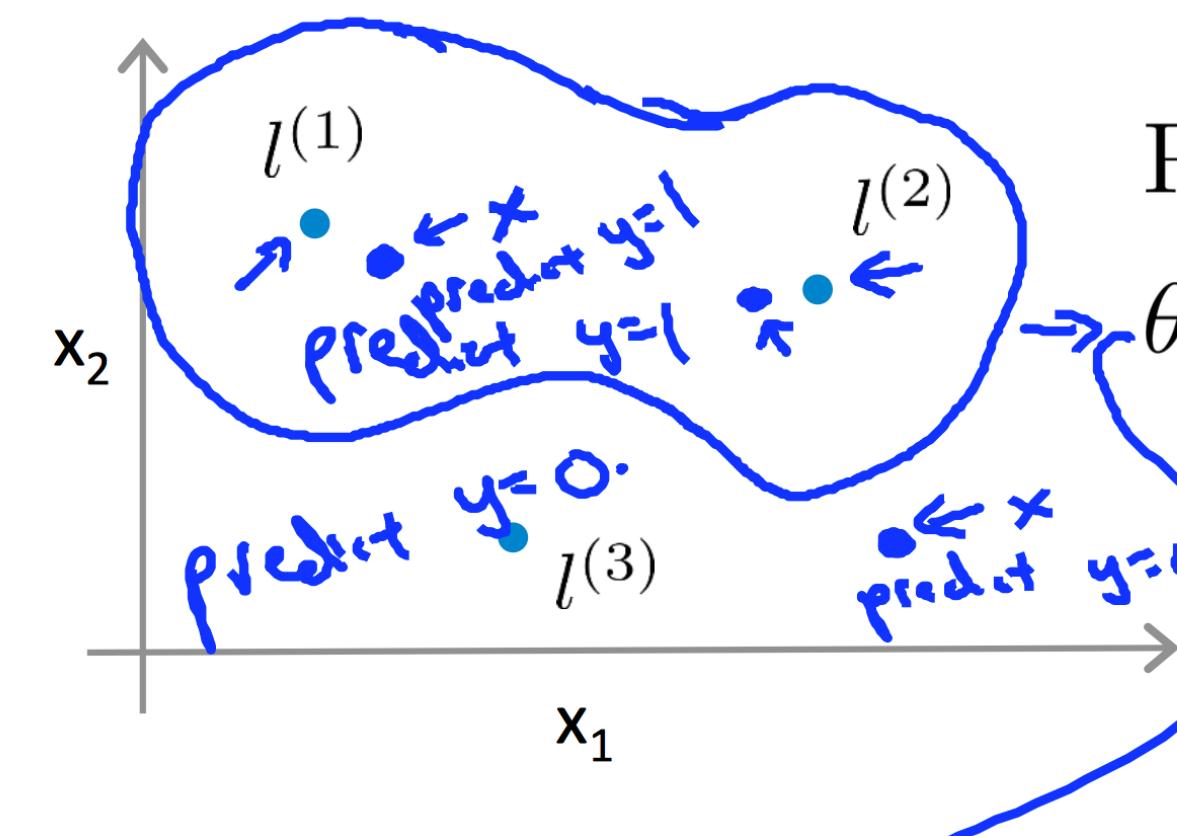
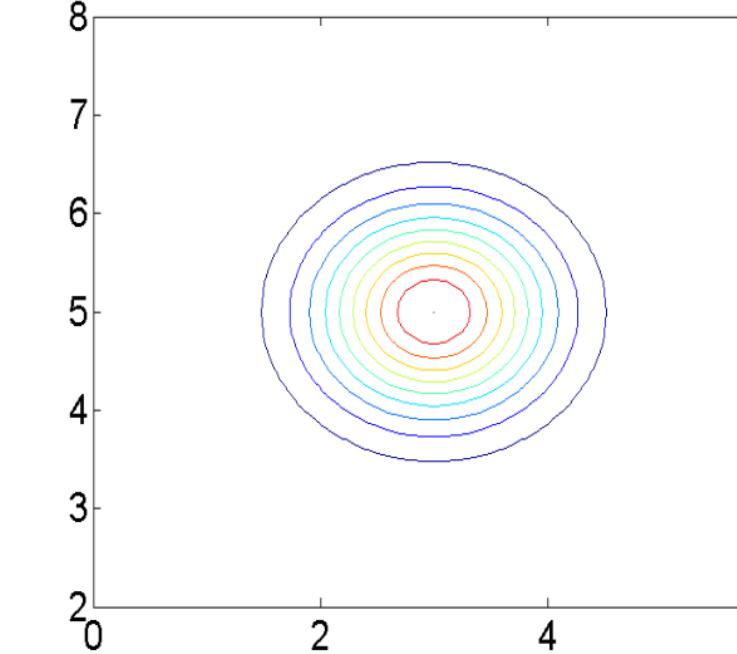
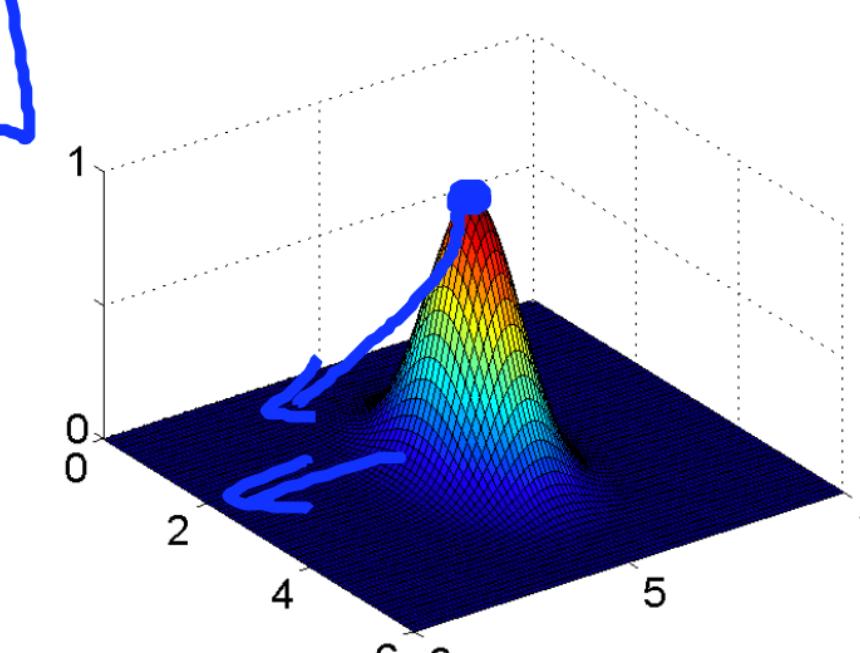
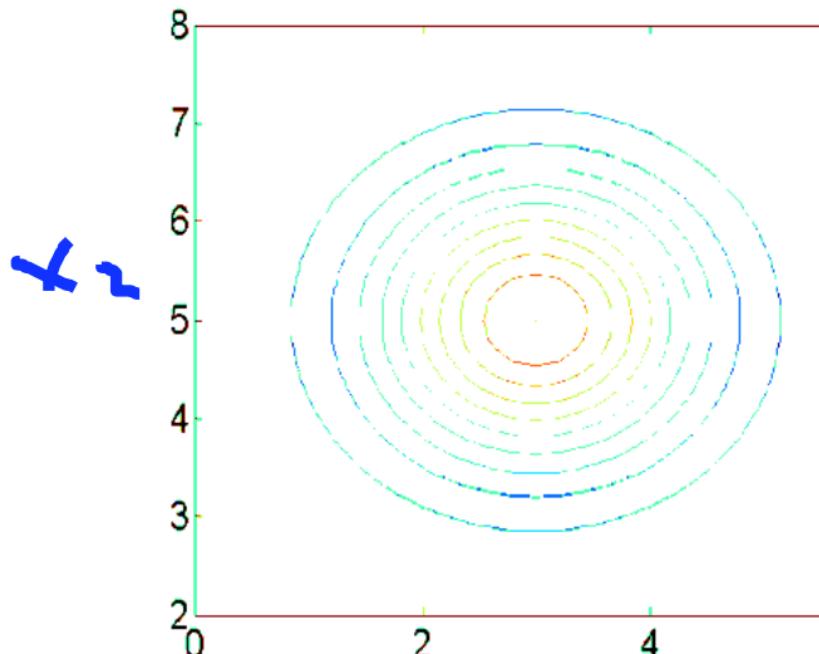
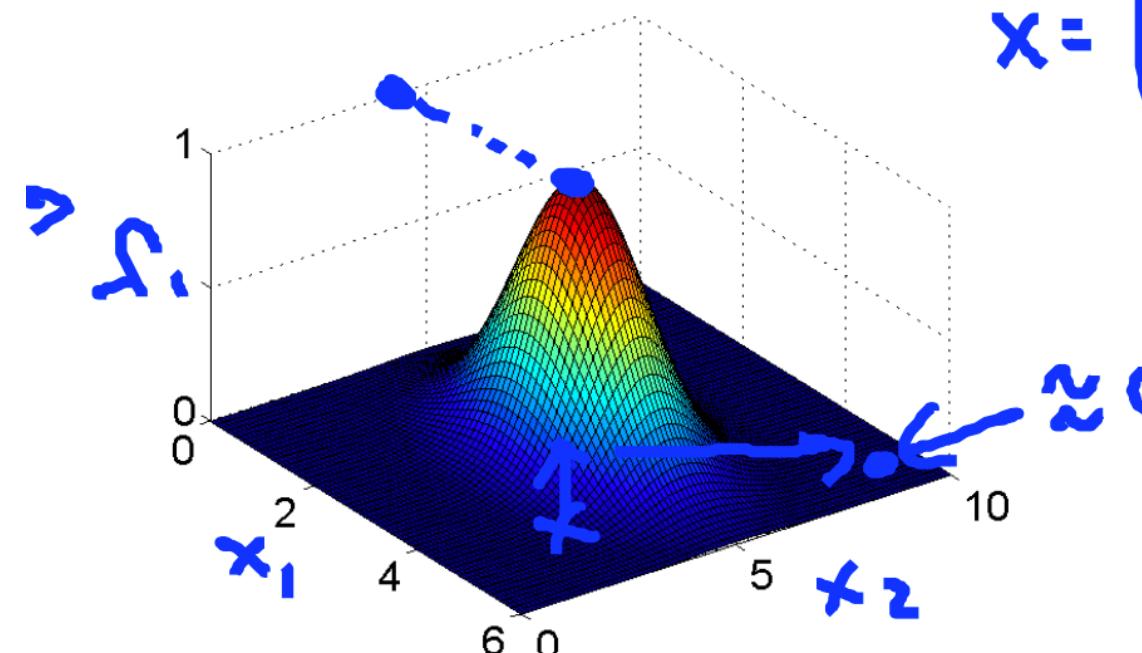
Is there a different / better choice of the features f_1, f_2, f_3, \dots ?

SVM Kernels

Example:

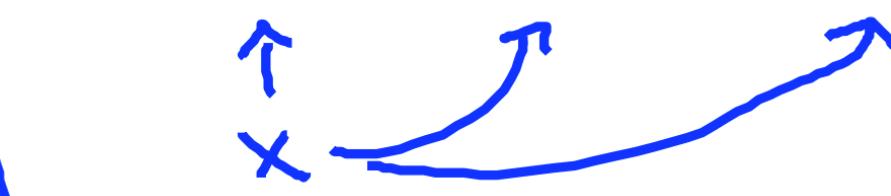
$$l^{(1)} = \begin{bmatrix} 3 \\ 5 \end{bmatrix}, \quad f_1 = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right)$$

$$\rightarrow \sigma^2 = 1$$



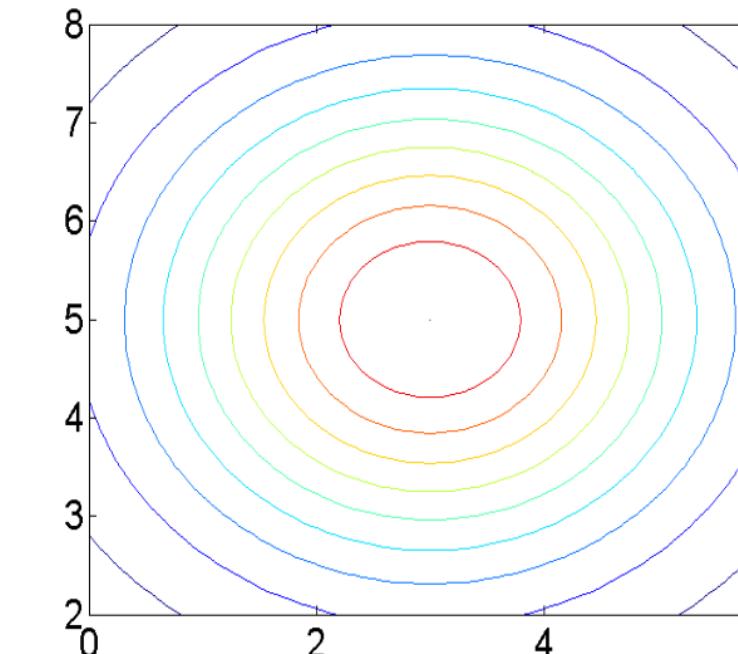
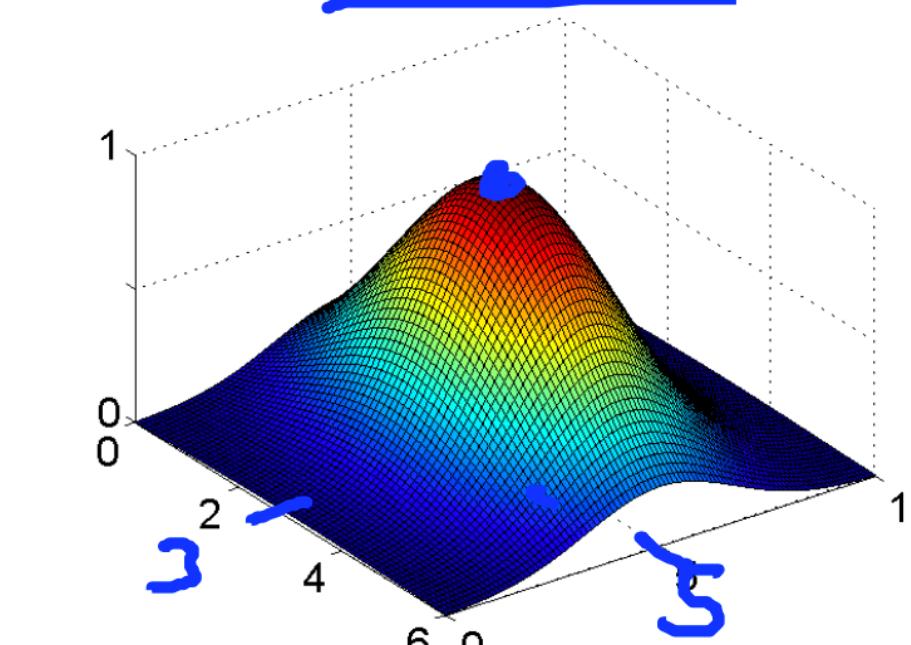
Predict "1" when

$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$$



$$\underline{\theta_0 = -0.5, \theta_1 = 1, \theta_2 = 1, \theta_3 = 0}$$

$$\sigma^2 = 3$$



SVMs in Practice

Need to specify:

→ Choice of parameter C.

Choice of kernel (similarity function):

E.g. No kernel ("linear kernel")

Predict "y = 1" if $\underline{\theta^T x} \geq 0$

$$\Theta_0 + \Theta_1 x_1 + \dots + \Theta_n x_n \geq 0$$

→ n large, m small

$$x \in \mathbb{R}^{n+1}$$

$$x \in \mathbb{R}^n, n \text{ small}$$

and/or m large



SVMs in Practice

Support Vector Machines

Support Vector Machines (SVM) are a method that uses points in a transformed problem space that best separate classes into two groups. Classification for multiple classes is supported by a one-vs-all method. SVM also supports regression by modeling the function with a minimum amount of allowable error.

This recipe shows use of the SVM algorithm to make predictions for the iris dataset (classification).

```
# SVM Classification
from sklearn import datasets
from sklearn import metrics
from sklearn.svm import SVC
# load the iris datasets
dataset = datasets.load_iris()
# fit a SVM model to the data
model = SVC()
model.fit(dataset.data, dataset.target)
print(model)
# make predictions
expected = dataset.target
predicted = model.predict(dataset.data)
# summarize the fit of the model
print(metrics.classification_report(expected, predicted))
print(metrics.confusion_matrix(expected, predicted))
```

This recipe shows use of the SVM algorithm to make predictions for the diabetes dataset (regression).

```
# SVM Regression
import numpy as np
from sklearn import datasets
from sklearn.svm import SVR
# load the diabetes datasets
dataset = datasets.load_diabetes()
# fit a SVM model to the data
model = SVR()
model.fit(dataset.data, dataset.target)
print(model)
# make predictions
expected = dataset.target
predicted = model.predict(dataset.data)
# summarize the fit of the model
mse = np.mean((predicted-expected)**2)
print(mse)
print(model.score(dataset.data, dataset.target))
```

Naive Bayes

Naive Bayes

Naive Bayes uses Bayes Theorem to model the conditional relationship of each attribute to the class variable.

This recipe shows the fitting of an Naive Bayes algorithm to the iris dataset.

```
# Gaussian Naive Bayes
from sklearn import datasets
from sklearn import metrics
from sklearn.naive_bayes import GaussianNB
# load the iris datasets
dataset = datasets.load_iris()
# fit a Naive Bayes model to the data
model = GaussianNB()
model.fit(dataset.data, dataset.target)
print(model)
# make predictions
expected = dataset.target
predicted = model.predict(dataset.data)
# summarize the fit of the model
print(metrics.classification_report(expected, predicted))
print(metrics.confusion_matrix(expected, predicted))
```

K-Nearest Neighbors

k-Nearest Neighbor

The k-Nearest Neighbor (kNN) method makes predictions by locating similar cases to a given data instance (using a similarity function) and returning the average or majority of the most similar data instances. The kNN algorithm can be used for classification or regression.

This recipe shows use of the kNN algorithm to make predictions for the iris dataset (classification).

```
# k-Nearest Neighbor Classification
from sklearn import datasets
from sklearn import metrics
from sklearn.neighbors import KNeighborsClassifier
# load iris the datasets
dataset = datasets.load_iris()
# fit a k-nearest neighbor model to the data
model = KNeighborsClassifier()
model.fit(dataset.data, dataset.target)
print(model)
# make predictions
expected = dataset.target
predicted = model.predict(dataset.data)
# summarize the fit of the model
print(metrics.classification_report(expected, predicted))
print(metrics.confusion_matrix(expected, predicted))
```

This recipe shows use of the kNN algorithm to make predictions for the diabetes dataset (regression).

```
# k-Nearest Neighbor Regression
import numpy as np
from sklearn import datasets
from sklearn.neighbors import KNeighborsRegressor
# load the diabetes datasets
dataset = datasets.load_diabetes()
# fit a model to the data
model = KNeighborsRegressor()
model.fit(dataset.data, dataset.target)
print(model)
# make predictions
expected = dataset.target
predicted = model.predict(dataset.data)
# summarize the fit of the model
mse = np.mean((predicted-expected)**2)
print(mse)
print(model.score(dataset.data, dataset.target))
```

Ensemble Methods



Random Forest

This recipe shows use of the Random Forest algorithm to make predictions for the iris dataset (classification).

```
# Random Forest Classification
from sklearn import datasets
from sklearn import metrics
from sklearn.ensemble import RandomForestClassifier
# load iris the datasets
dataset = datasets.load_iris()
# fit a random forest model to the data
model = RandomForestClassifier()
model.fit(dataset.data, dataset.target)
print(model)
# make predictions
expected = dataset.target
predicted = model.predict(dataset.data)
# summarize the fit of the model
print(metrics.classification_report(expected, predicted))
print(metrics.confusion_matrix(expected, predicted))
```

This recipe shows use of the Random Forest model to make predictions for the diabetes dataset (regression).

```
# Random Forest Regression
import numpy as np
from sklearn import datasets
from sklearn.ensemble import RandomForestRegressor
# load the diabetes datasets
dataset = datasets.load_diabetes()
# fit a random forest model to the data
model = RandomForestRegressor()
model.fit(dataset.data, dataset.target)
print(model)
# make predictions
expected = dataset.target
predicted = model.predict(dataset.data)
# summarize the fit of the model
mse = np.mean((predicted-expected)**2)
print(mse)
print(model.score(dataset.data, dataset.target))
```

Extra Trees



This recipe shows use of the Extra Trees algorithm to make predictions for the iris dataset (classification).

```
# Extra Trees Classification
from sklearn import datasets
from sklearn import metrics
from sklearn.ensemble import ExtraTreesClassifier
# load the iris datasets
dataset = datasets.load_iris()
# fit a Extra Trees model to the data
model = ExtraTreesClassifier()
model.fit(dataset.data, dataset.target)
print(model)
# make predictions
expected = dataset.target
predicted = model.predict(dataset.data)
# summarize the fit of the model
print(metrics.classification_report(expected, predicted))
print(metrics.confusion_matrix(expected, predicted))
```

This recipe shows use of the Extra Trees algorithm to make predictions for the diabetes dataset (regression).

```
# Extra Trees Regression
import numpy as np
from sklearn import datasets
from sklearn.ensemble import ExtraTreesRegressor
# load the diabetes datasets
dataset = datasets.load_diabetes()
# fit a extra trees model to the data
model = ExtraTreesRegressor()
model.fit(dataset.data, dataset.target)
print(model)
# make predictions
expected = dataset.target
predicted = model.predict(dataset.data)
# summarize the fit of the model
mse = np.mean((predicted-expected)**2)
print(mse)
print(model.score(dataset.data, dataset.target))
```

Adaptive Boosting

This recipe shows use of the AdaBoost algorithm to make predictions for the iris dataset (classification).

```
# AdaBoost Classification
from sklearn import datasets
from sklearn import metrics
from sklearn.ensemble import AdaBoostClassifier
# load the iris datasets
dataset = datasets.load_iris()
# fit an AdaBoost model to the data
model = AdaBoostClassifier()
model.fit(dataset.data, dataset.target)
print(model)
# make predictions
expected = dataset.target
predicted = model.predict(dataset.data)
# summarize the fit of the model
print(metrics.classification_report(expected, predicted))
print(metrics.confusion_matrix(expected, predicted))
```

This recipe shows use of the AdaBoost algorithm to make predictions for the diabetes dataset (regression).

```
# Random Forest Regression
import numpy as np
from sklearn import datasets
from sklearn.ensemble import AdaBoostRegressor
# load the diabetes datasets
dataset = datasets.load_diabetes()
# fit an AdaBoost model to the data
model = AdaBoostRegressor()
model.fit(dataset.data, dataset.target)
print(model)
# make predictions
expected = dataset.target
predicted = model.predict(dataset.data)
# summarize the fit of the model
mse = np.mean((predicted-expected)**2)
print(mse)
print(model.score(dataset.data, dataset.target))
```

Gradient Boosting

This recipe shows use of the Gradient Boosting algorithm to make predictions for the iris dataset (classification).

```
# Gradient Boosting Classification
from sklearn import datasets
from sklearn import metrics
from sklearn.ensemble import GradientBoostingClassifier
# load the iris dataset
dataset = datasets.load_iris()
# fit a Gradient Boosting model to the data
model = GradientBoostingClassifier()
model.fit(dataset.data, dataset.target)
print(model)
# make predictions
expected = dataset.target
predicted = model.predict(dataset.data)
# summarize the fit of the model
print(metrics.classification_report(expected, predicted))
print(metrics.confusion_matrix(expected, predicted))
```

This recipe shows use of the Gradient Boosting algorithm to make predictions for the diabetes dataset (regression).

```
# Gradient Boosting Regression
import numpy as np
from sklearn import datasets
from sklearn.ensemble import GradientBoostingRegressor
# load the diabetes dataset
dataset = datasets.load_diabetes()
# fit a Gradient Boosting model to the data
model = GradientBoostingRegressor()
model.fit(dataset.data, dataset.target)
print(model)
# make predictions
expected = dataset.target
predicted = model.predict(dataset.data)
# summarize the fit of the model
mse = np.mean((predicted-expected)**2)
print(mse)
print(model.score(dataset.data, dataset.target))
```

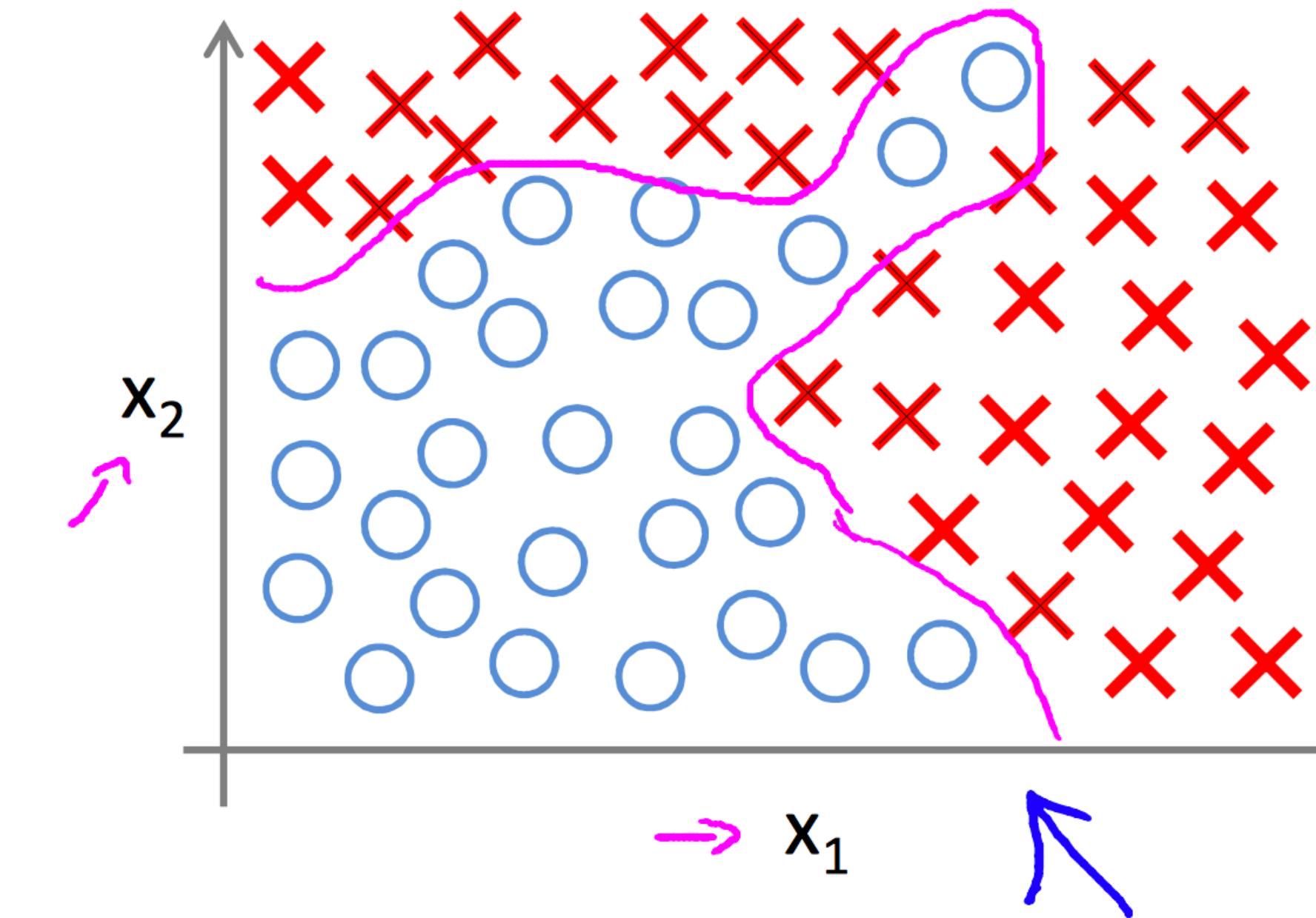
Neural Networks

Neural networks is a model inspired by how the brain works. It is widely used today in many applications: when your phone interprets and understand your voice commands, it is likely that a neural network is helping to understand your speech; when you cash a check, the machines that automatically read the digits also use neural networks.

Motivations

Issues with Non-linear Classification

Non-linear Classification



$\rightarrow \underline{x_1} = \text{size}$
 $\underline{x_2} = \# \text{bedrooms}$
 $\underline{x_3} = \# \text{floors}$
 $x_4 = \text{age}$
 \dots
 $x_{100} -$

$$n=100$$

$\underbrace{g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 x_2 + \theta_5 x_1^3 x_2 + \theta_6 x_1 x_2^2 + \dots)}$

$\rightarrow \underline{x_1^2}, \underline{x_1 x_2}, \underline{x_1 x_3}, \underline{x_1 x_4} \dots \underline{x_1 x_{100}}$
 $\underline{x_2^2}, \underline{x_2 x_3} \dots$

$\approx \underline{5000 \text{ feature}}$

$\underline{\underline{x_1^2, x_2^2, x_3^2, \dots, x_{100}^2}}$

$\rightarrow \underline{x_1 x_2 x_3}, \underline{x_1^2 x_2}, \underline{x_{10} x_{11} x_{12}}, \dots$

$\mathcal{O}(n^2)$

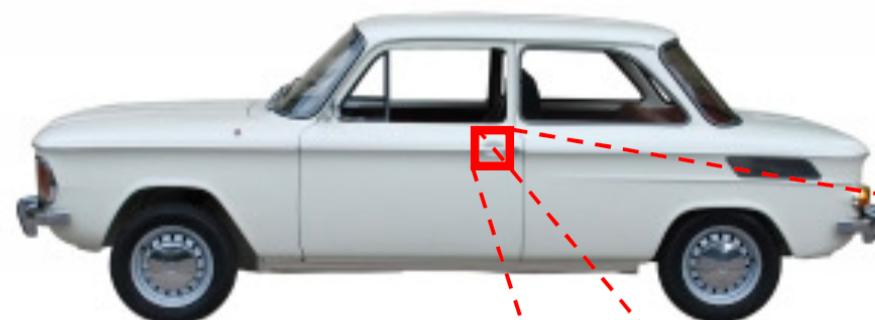
$\frac{n^2}{2}$

$\cancel{10}$

Motivations

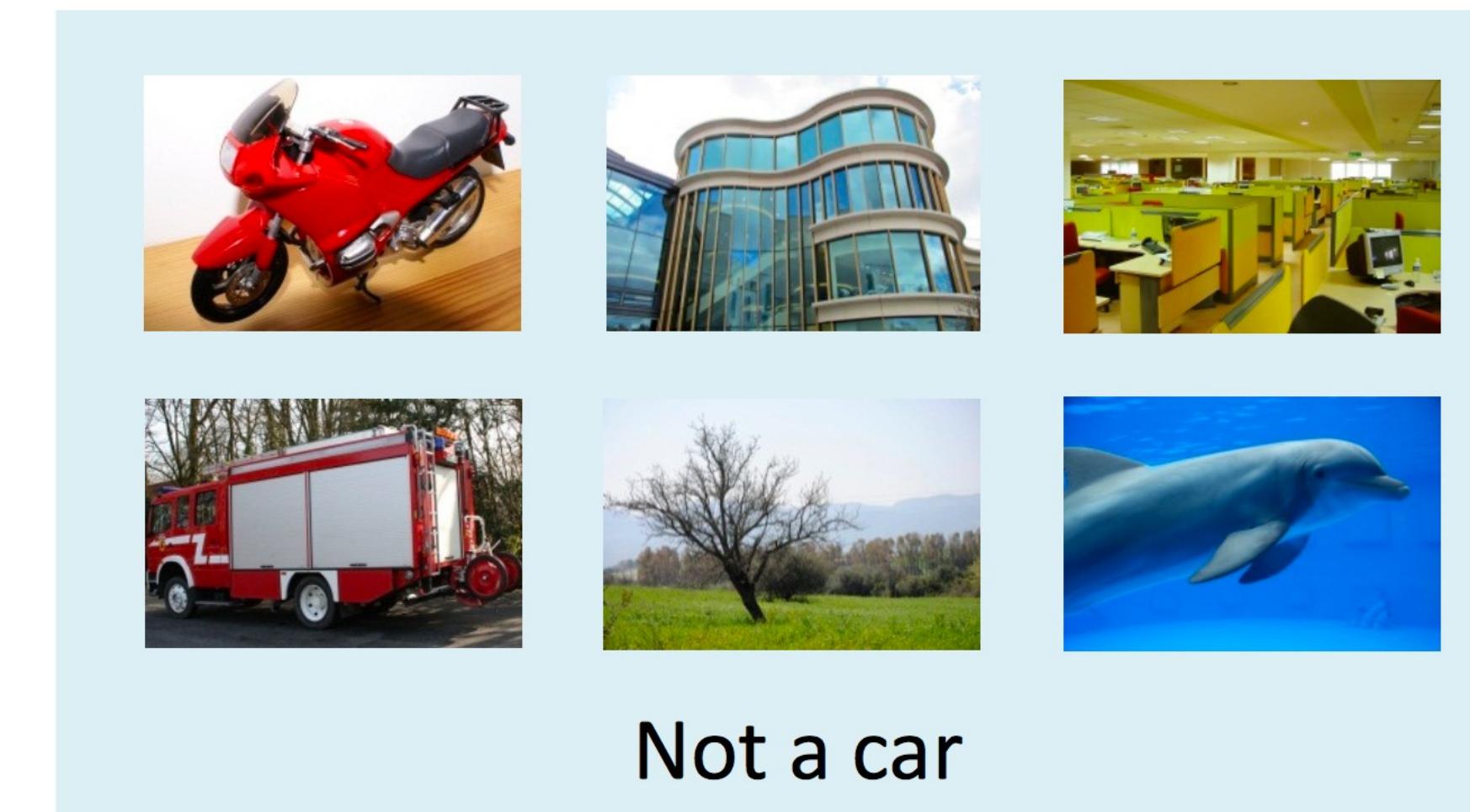
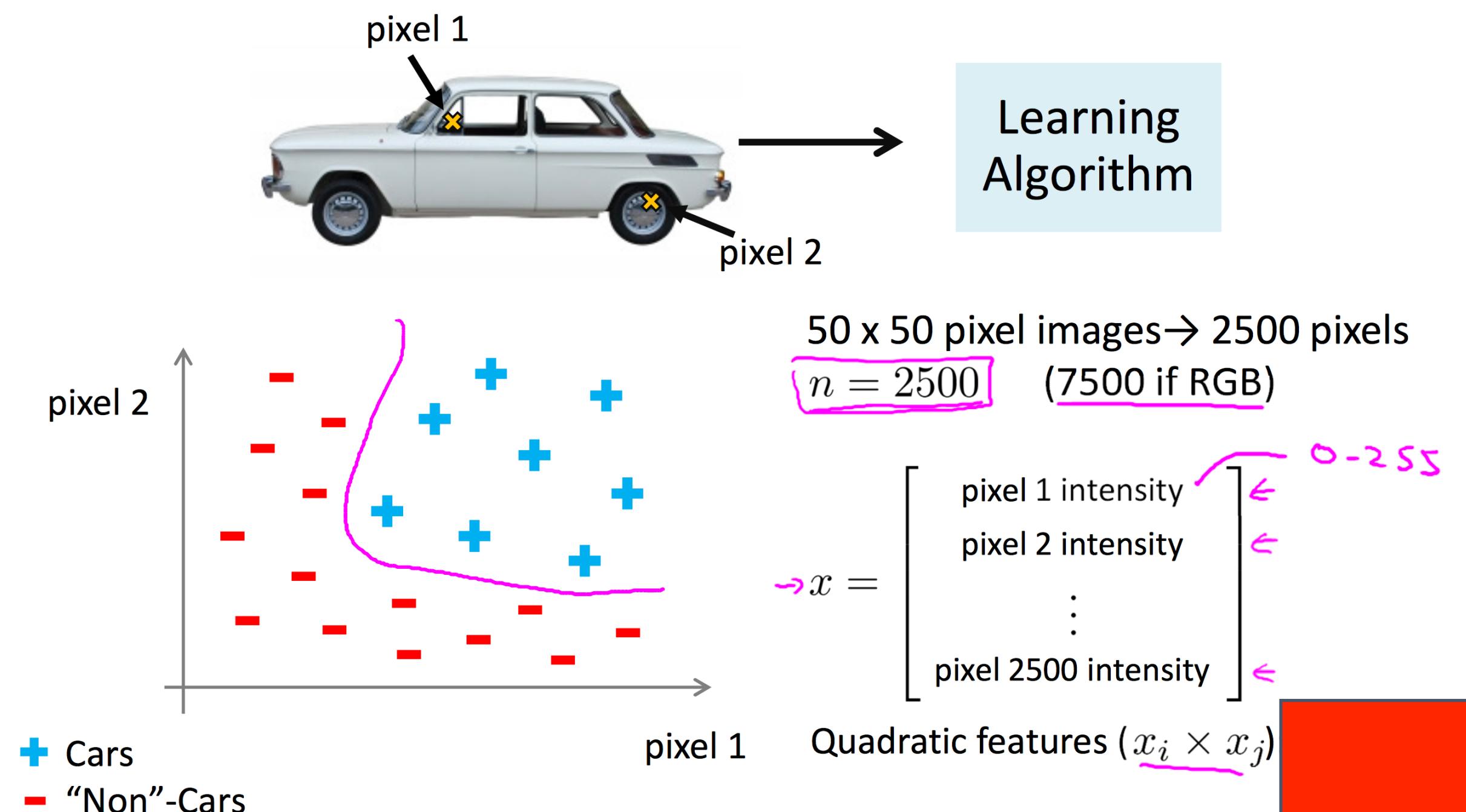
What is this?

You see this:



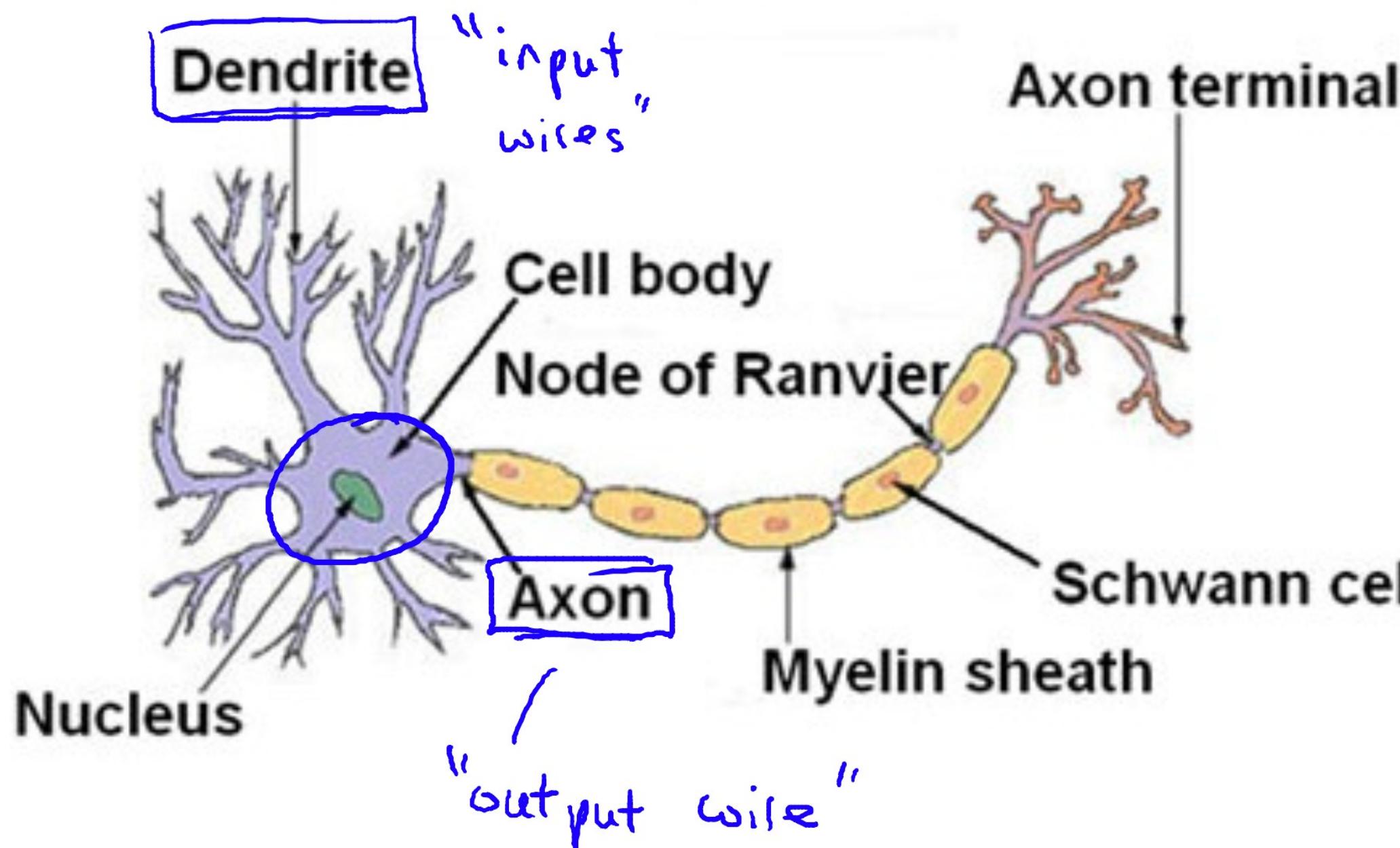
But the camera sees this:

194	210	201	212	199	213	215	195	178	158	182	209
180	189	190	221	209	205	191	167	147	115	129	163
114	126	140	188	176	165	152	140	170	106	78	88
87	103	115	154	143	142	149	153	173	101	57	57
102	112	106	131	122	138	152	147	128	84	58	66
94	95	79	104	105	124	129	113	107	87	69	67
68	71	69	98	89	92	98	95	89	88	76	67
41	56	68	99	63	45	60	82	58	76	75	65
20	43	69	75	56	41	51	73	55	70	63	44
50	50	57	69	75	75	73	74	53	68	59	37
72	59	53	66	84	92	84	74	57	72	63	42
67	61	58	65	75	78	76	73	59	75	69	50

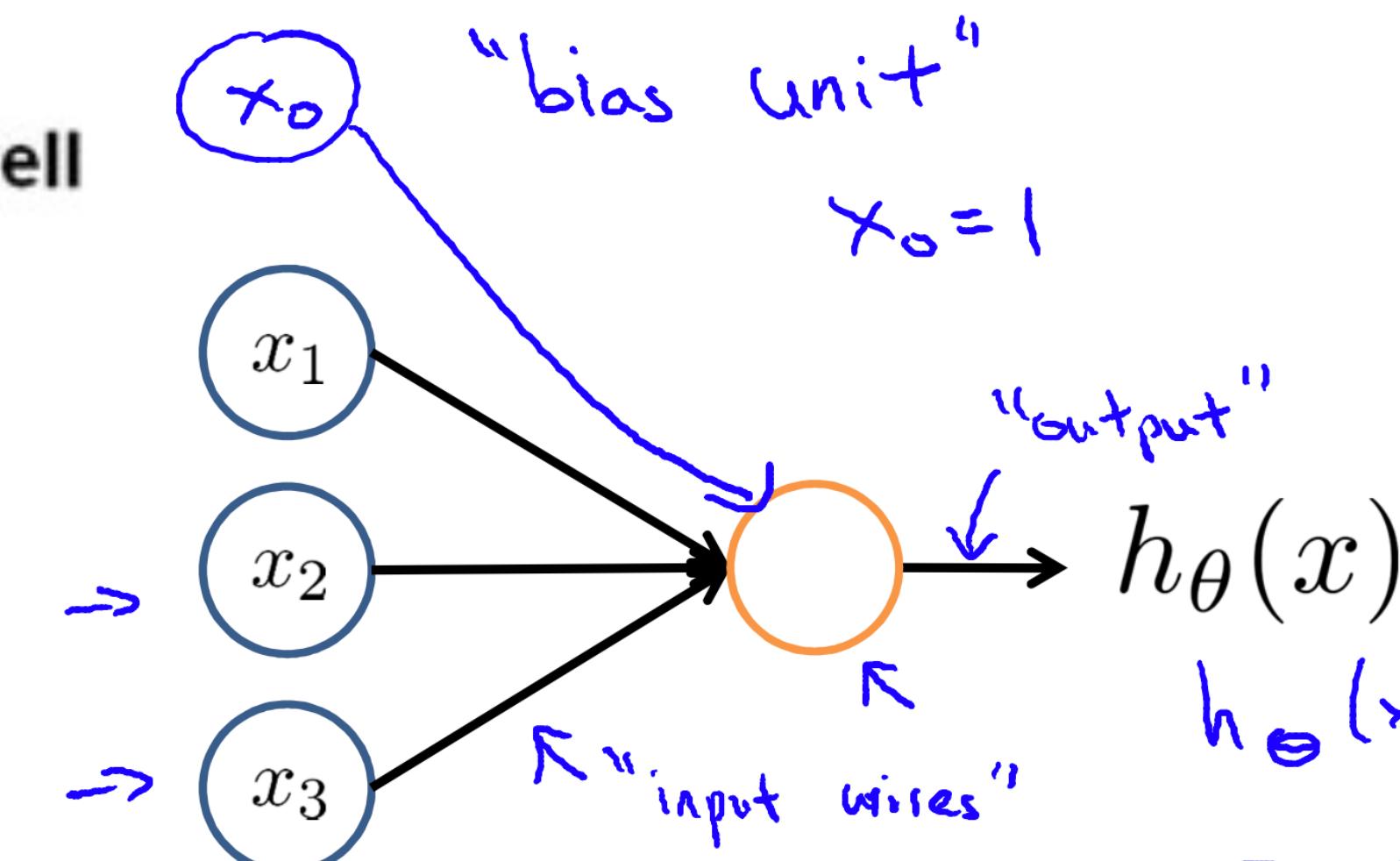


Motivations

Neural Networks perform better for non-linear data with number of features



Non-linear Hypotheses
neuron model: Logistic unit



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

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$$

↑
"weights" ←
(parameters)

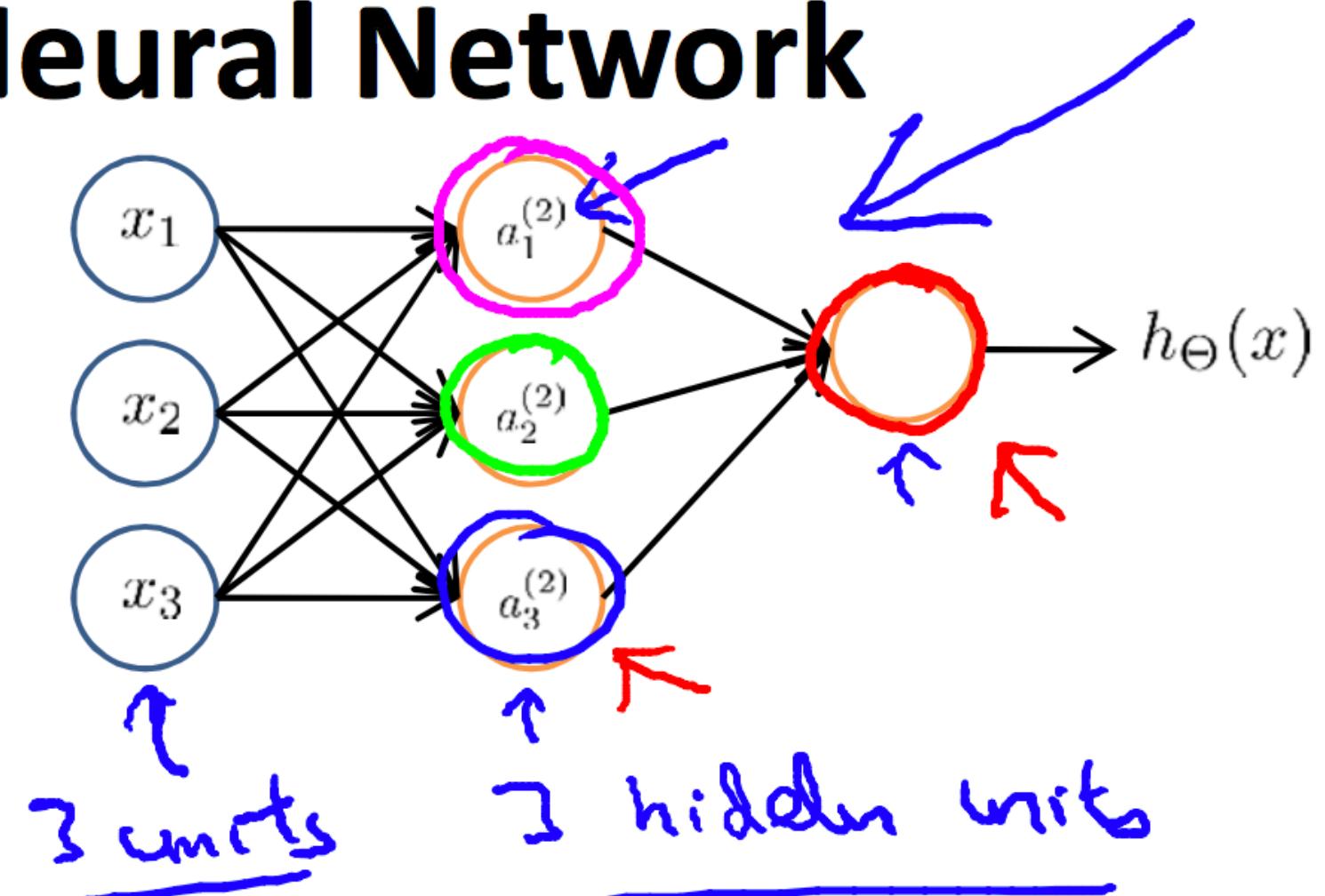
Sigmoid (logistic) activation function.

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

Neural Networks

Model representation I

Neural Network



$\rightarrow a_i^{(j)} = \text{"activation" of unit } i \text{ in layer } j$

$\rightarrow \Theta^{(j)} = \text{matrix of weights controlling function mapping from layer } j \text{ to layer } j + 1$

$$\Theta^{(1)} \in \mathbb{R}^{3 \times 4}$$

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

$$\rightarrow a_1^{(2)} = g(\underline{\Theta_{10}^{(1)} x_0 + \Theta_{11}^{(1)} x_1 + \Theta_{12}^{(1)} x_2 + \Theta_{13}^{(1)} x_3})$$

$$\rightarrow a_2^{(2)} = g(\underline{\Theta_{20}^{(1)} x_0 + \Theta_{21}^{(1)} x_1 + \Theta_{22}^{(1)} x_2 + \Theta_{23}^{(1)} x_3})$$

$$\rightarrow a_3^{(2)} = g(\underline{\Theta_{30}^{(1)} x_0 + \Theta_{31}^{(1)} x_1 + \Theta_{32}^{(1)} x_2 + \Theta_{33}^{(1)} x_3})$$

$$\rightarrow h_{\Theta}(x) = \underline{a_1^{(3)}} = g(\underline{\Theta_{10}^{(2)} a_0^{(2)} + \Theta_{11}^{(2)} a_1^{(2)} + \Theta_{12}^{(2)} a_2^{(2)} + \Theta_{13}^{(2)} a_3^{(2)}})$$

$$h_{\Theta}^{(2)}$$

↓

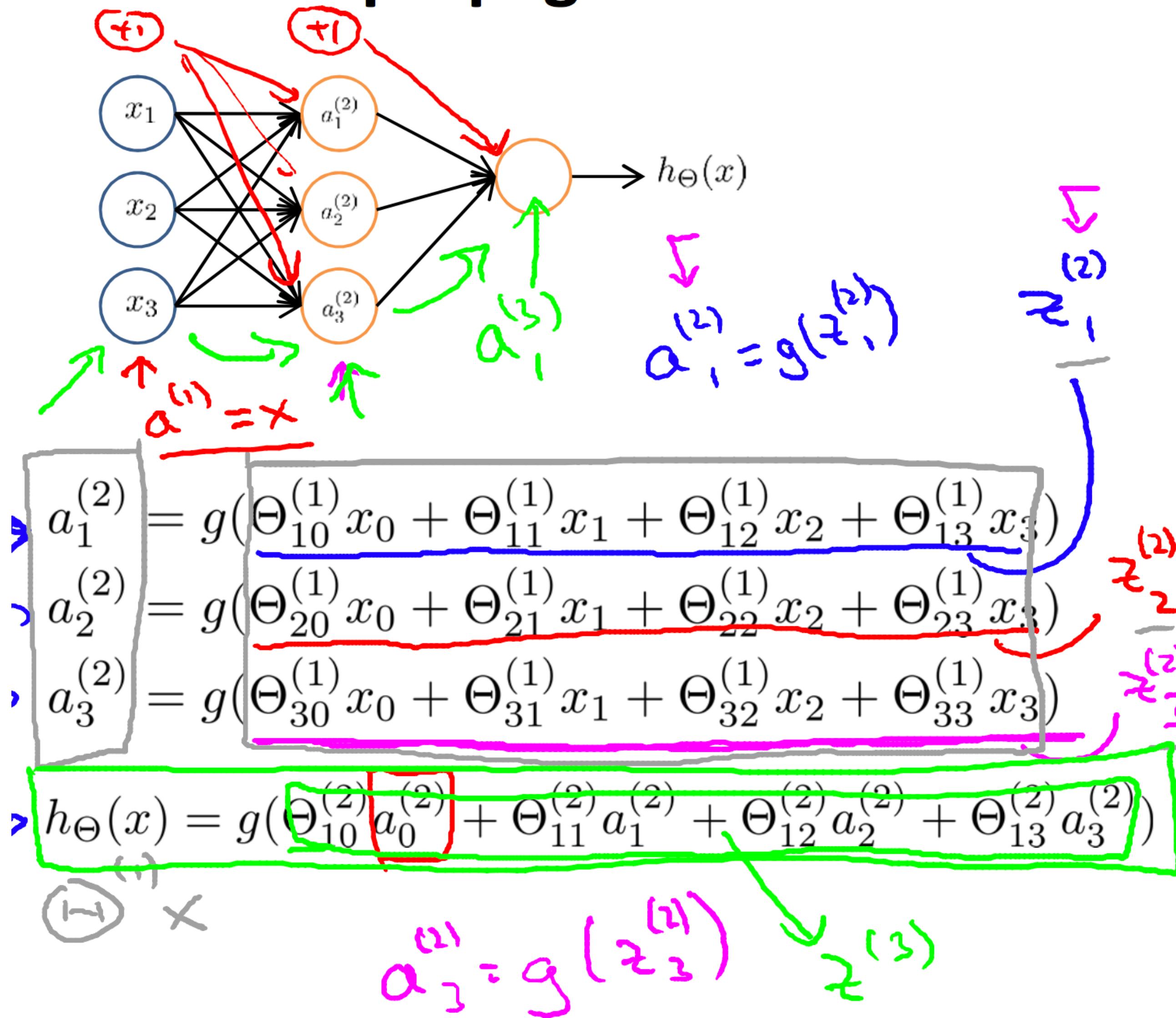
If network has s_j units in layer j , s_{j+1} units in layer $j+1$, then $\underline{\Theta^{(j)}}$ will be of dimension $\underline{s_{j+1}} \times (\underline{s_j} + 1)$.

$$s_{j+1} \times (s_j + 1)$$

Neural Networks

Model representation II

Forward propagation: Vectorized implementation



$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix}$$

$z^{(2)} = \Theta^{(1)} \times a^{(1)}$

$a^{(2)} = g(z^{(2)})$

Add $a_0^{(2)} = 1.$ $\rightarrow a^{(2)} \in \mathbb{R}^4$

$z^{(3)} = \Theta^{(2)} a^{(2)}$

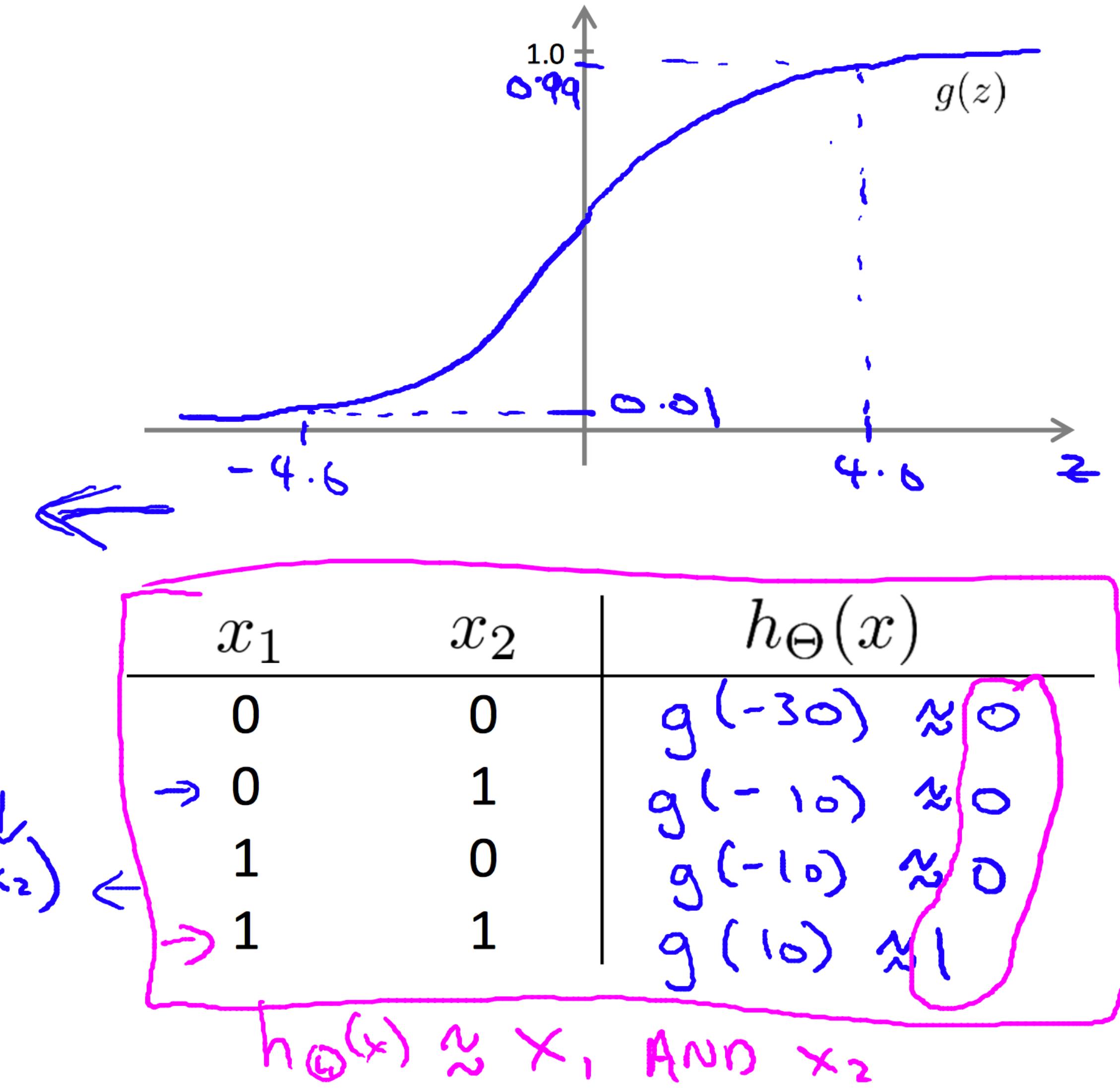
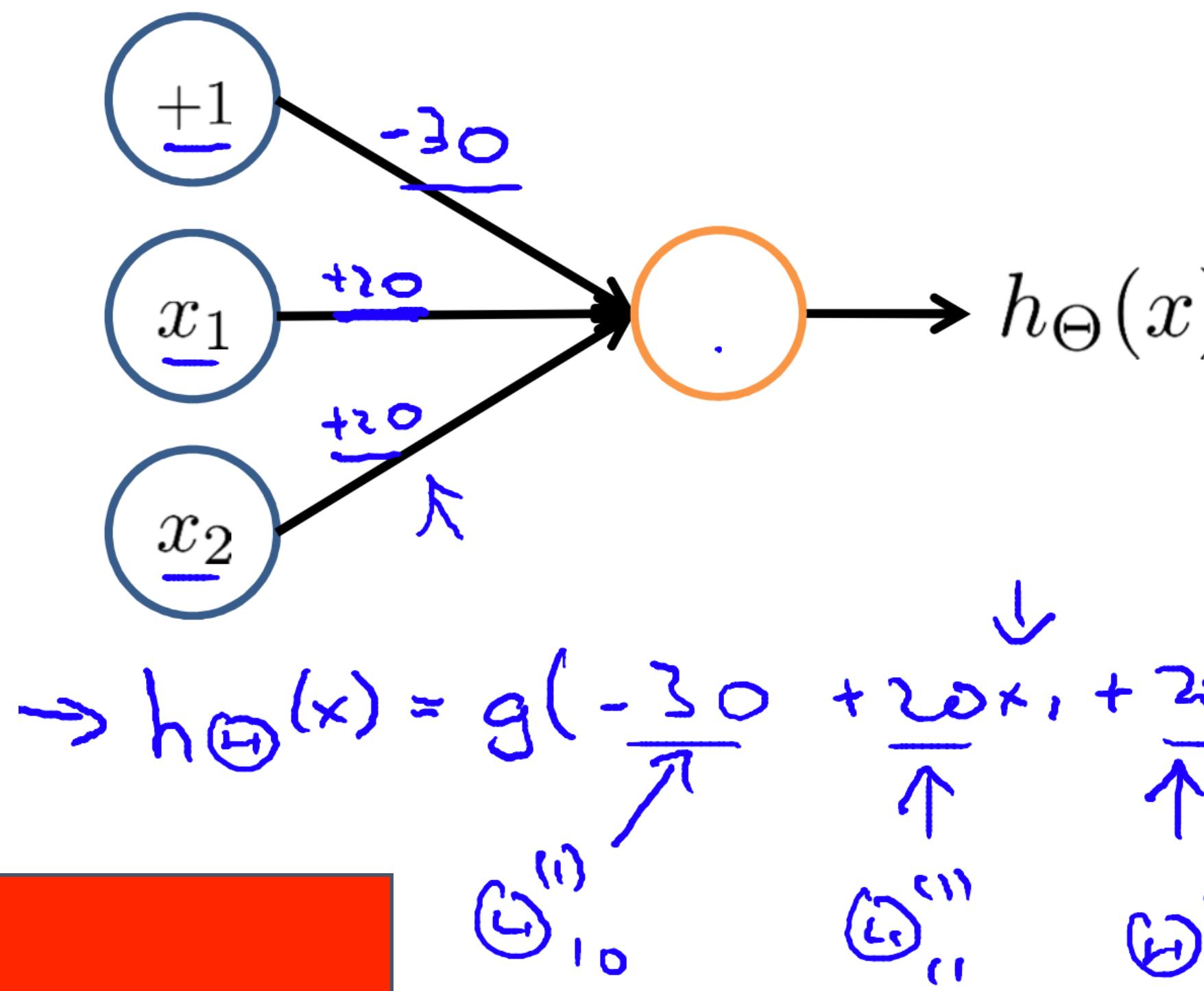
$h_{\Theta}(x) = a^{(3)} = g(z^{(3)})$

Neural Networks Applications

Example and intuition I

Simple example: AND

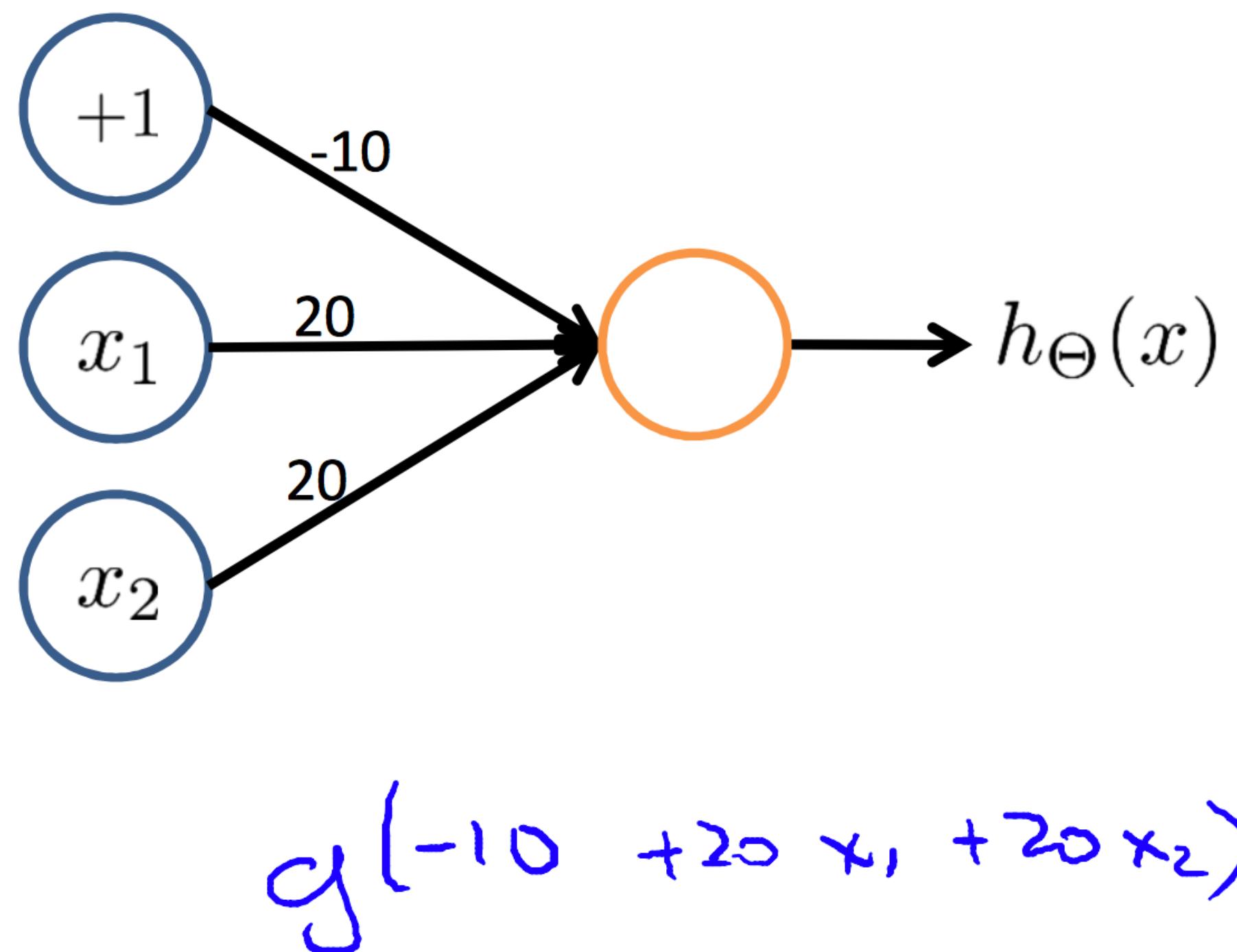
- $x_1, x_2 \in \{0, 1\}$
- $y = x_1 \text{ AND } x_2$



Neural Networks Applications

Example and intuition I

Example: OR function



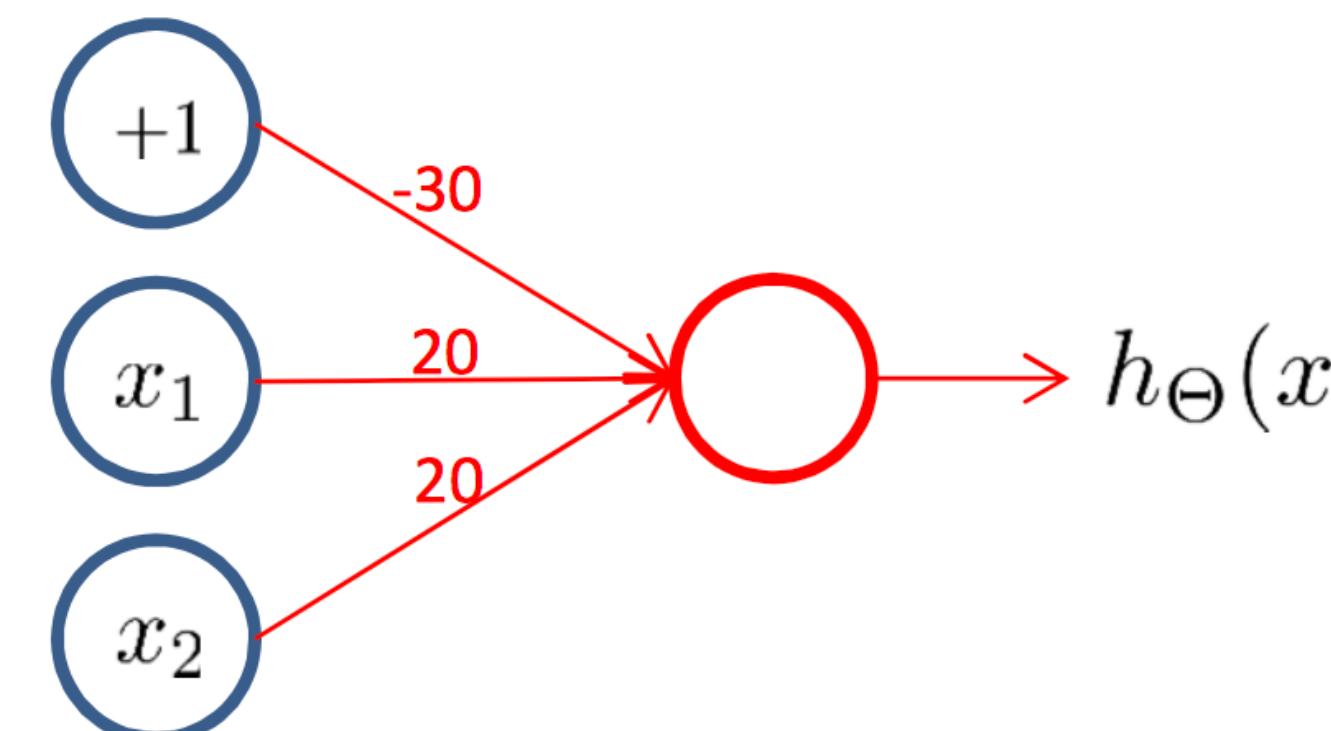
x_1	x_2	$h_{\Theta}(x)$
0	0	0
0	1	1
1	0	1
1	1	1

Handwritten annotations in blue and pink are present on the right side of the table:

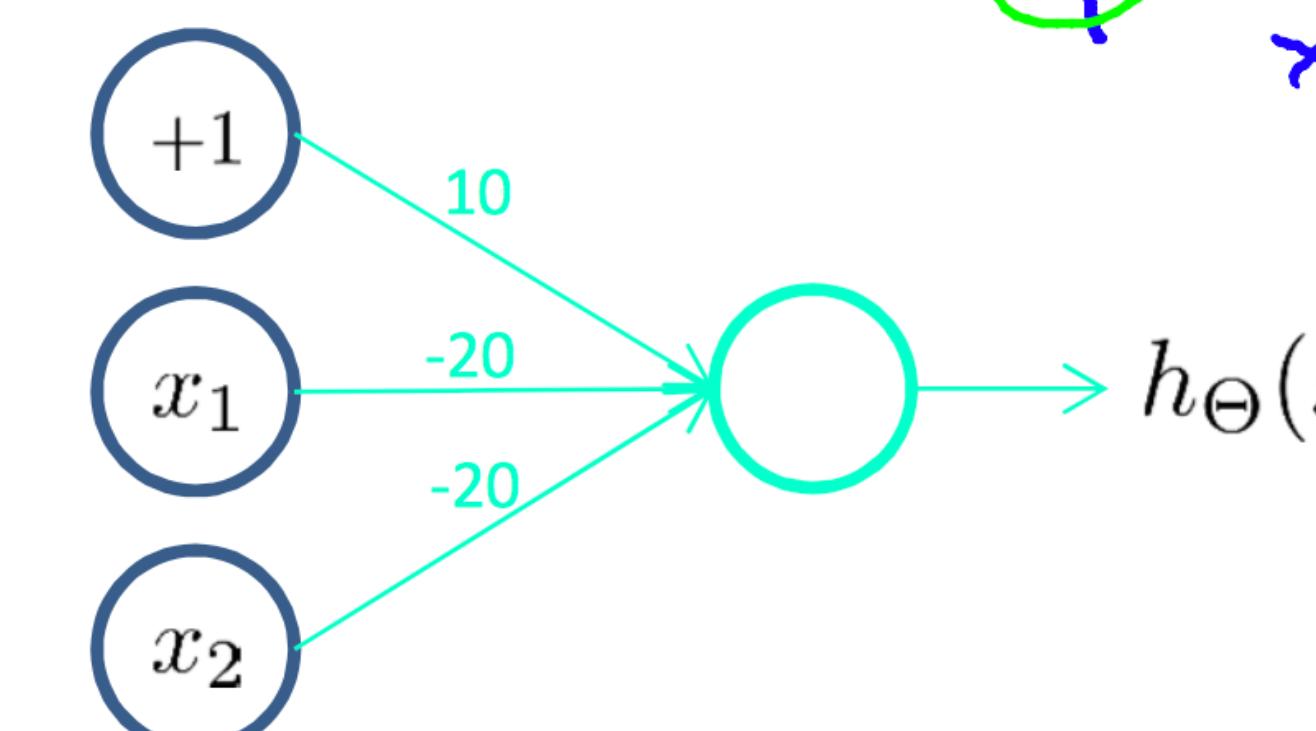
- $g(-10) \approx 0$ (near $x_1=0$)
- $g(10) \approx 1$ (near $x_1=1$)
- ≈ 1 (near $x_2=1$)
- ≈ 1 (near $x_1=1, x_2=1$)

Neural Networks Applications

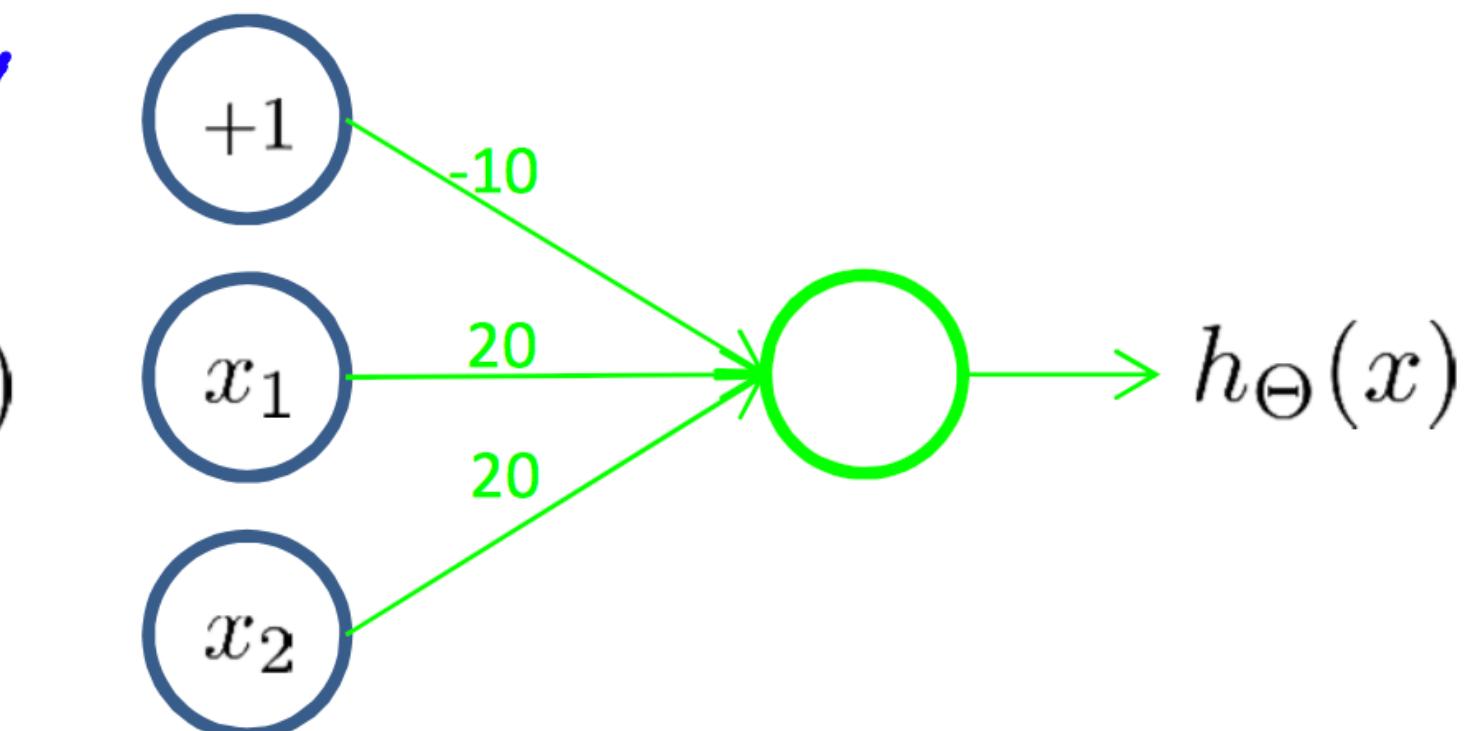
Putting it together:



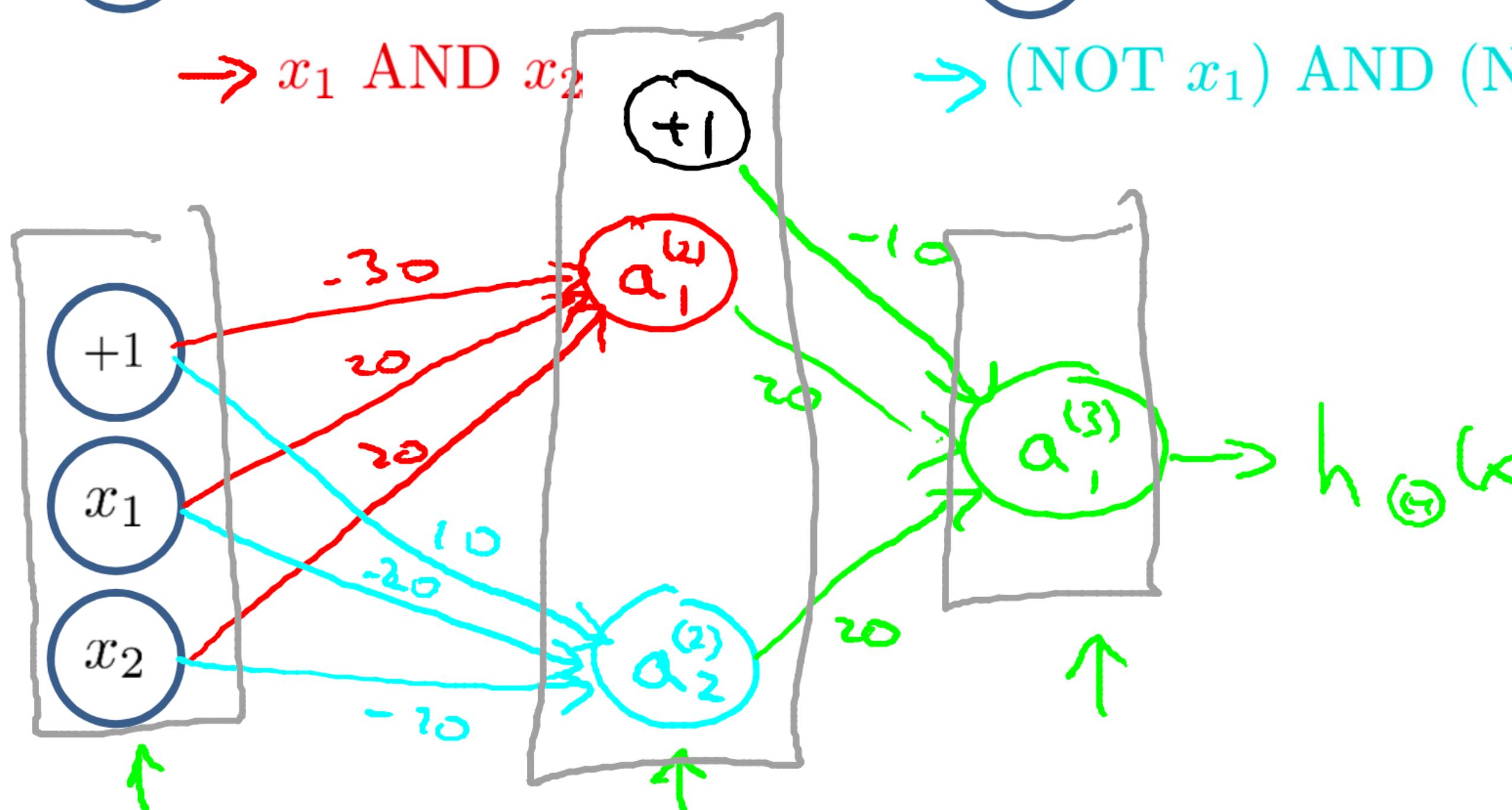
$\rightarrow x_1 \text{ AND } x_2$



$\rightarrow (\text{NOT } x_1) \text{ AND } (\text{NOT } x_2)$



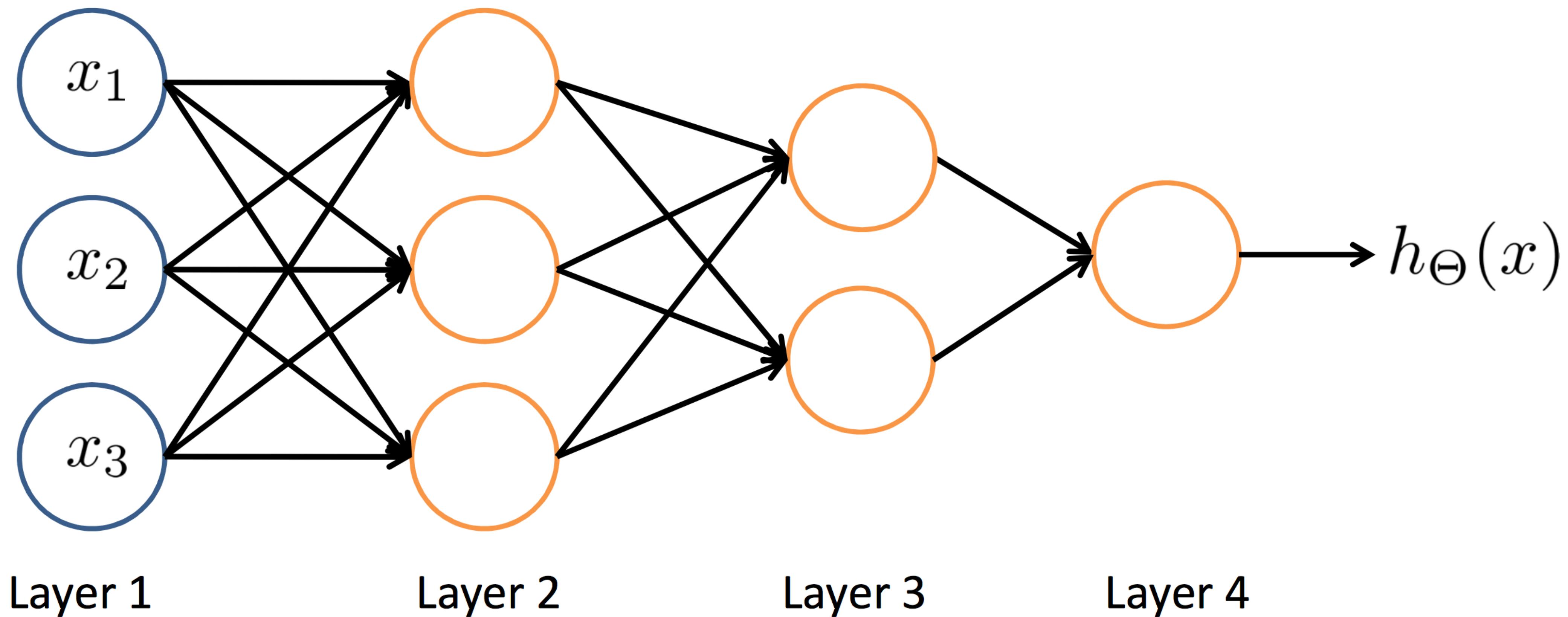
$\rightarrow x_1 \text{ OR } x_2$



x_1	x_2	$a_1^{(2)}$	$a_2^{(2)}$	$h_{\Theta}(x)$
0	0	0	1	1
0	1	0	0	0
1	0	0	0	0
1	1	1	1	1

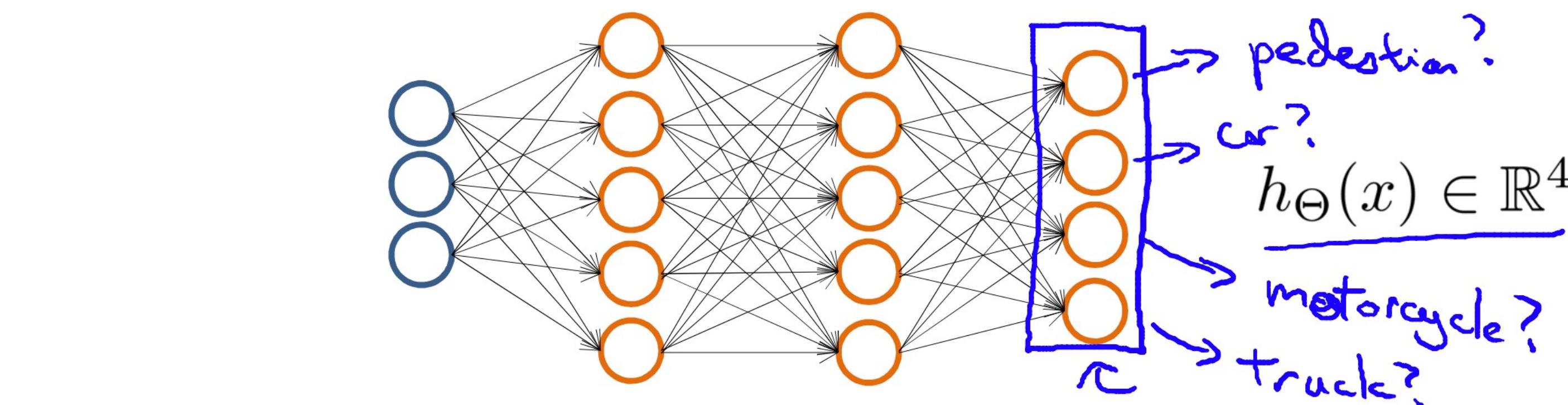
Neural Networks Applications

Neural Network intuition



Neural Networks Applications

Multiple output units: One-vs-all.



Want $h_{\Theta}(x) \approx \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$,
when pedestrian $h_{\Theta}(x) \approx \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$, when car $h_{\Theta}(x) \approx \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$, etc.
when motorcycle

$y^{(i)}$ one of $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$, $\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$, $\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$, $\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$
 pedestrian car motorcycle truck

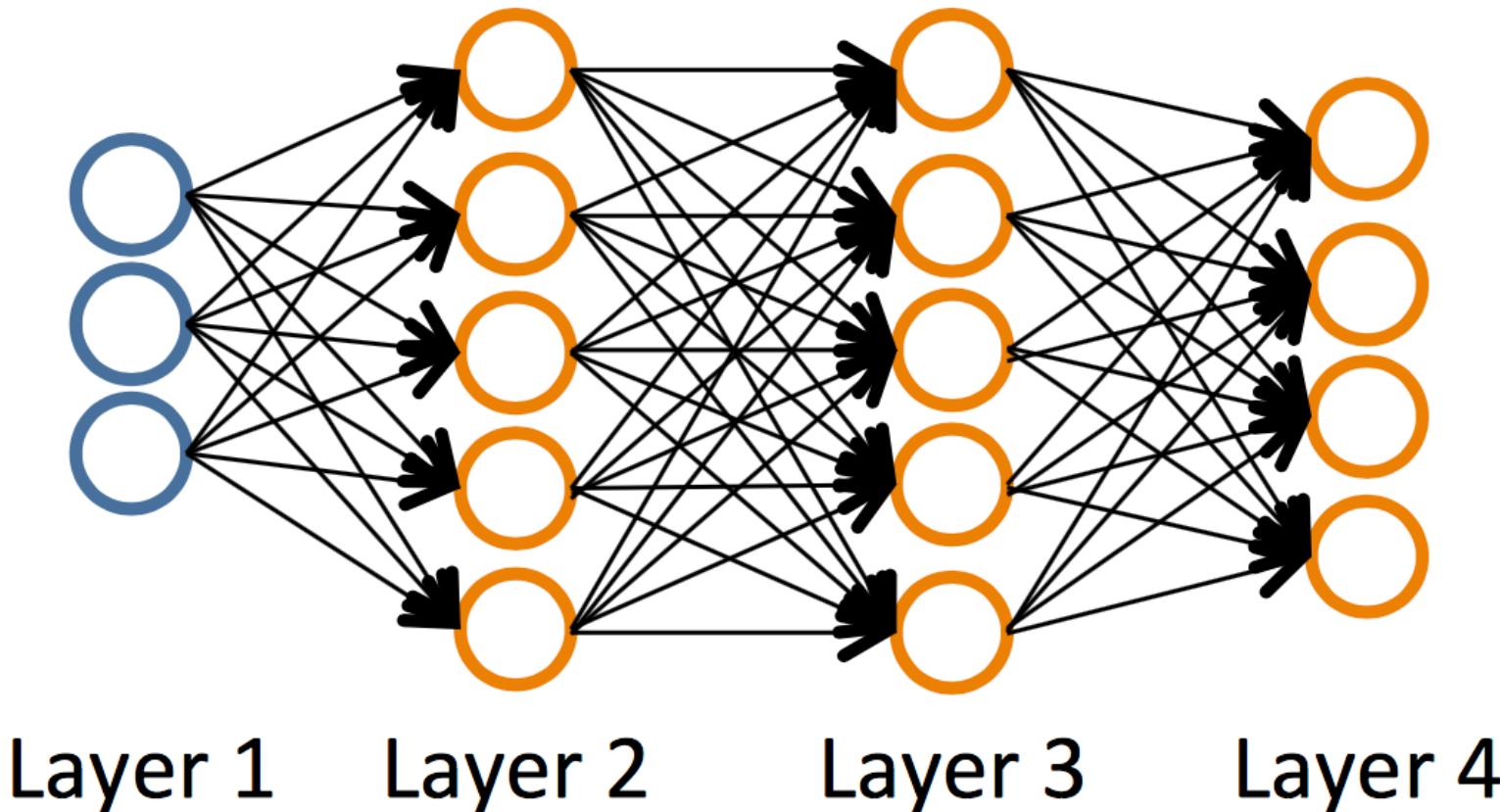
$(x^{(i)}, y^{(i)})$ $y \in \{1, 2, 3, 4\}$

$h_{\Theta}(x^{(i)}) \approx y^{(i)}$

Neural Networks: Learning

The Neural Network is one of the most powerful learning algorithms
linear classifier doesn't work)

Neural Network (Classification)



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

L = total no. of layers in network

s_l = no. of units (not counting bias unit) in layer l

Binary classification

$$y = 0 \text{ or } 1$$

1 output unit

Multi-class classification (K classes)

$$y \in \mathbb{R}^K \quad \text{E.g. } \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

pedestrian car motorcycle truck

K output units

Cost function

Logistic regression:

$$J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$



Cost Function and Backpropagation

Gradient computation

$$\Rightarrow \underline{J(\Theta)} = -\frac{1}{m} \left[\sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log h_\theta(x^{(i)})_k + (1 - y_k^{(i)}) \log(1 - h_\theta(x^{(i)})_k) \right]$$

$$+ \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\Theta_j^{(l)})^2$$

$$\Rightarrow \min_{\Theta} J(\Theta)$$

Need code to compute:



Cost Function and Backpropagation

Gradient computation

Given one training example $(\underline{x}, \underline{y})$:

Forward propagation:

$$\underline{a}^{(1)} = \underline{x}$$

$$\rightarrow z^{(2)} = \Theta^{(1)} \underline{a}^{(1)}$$

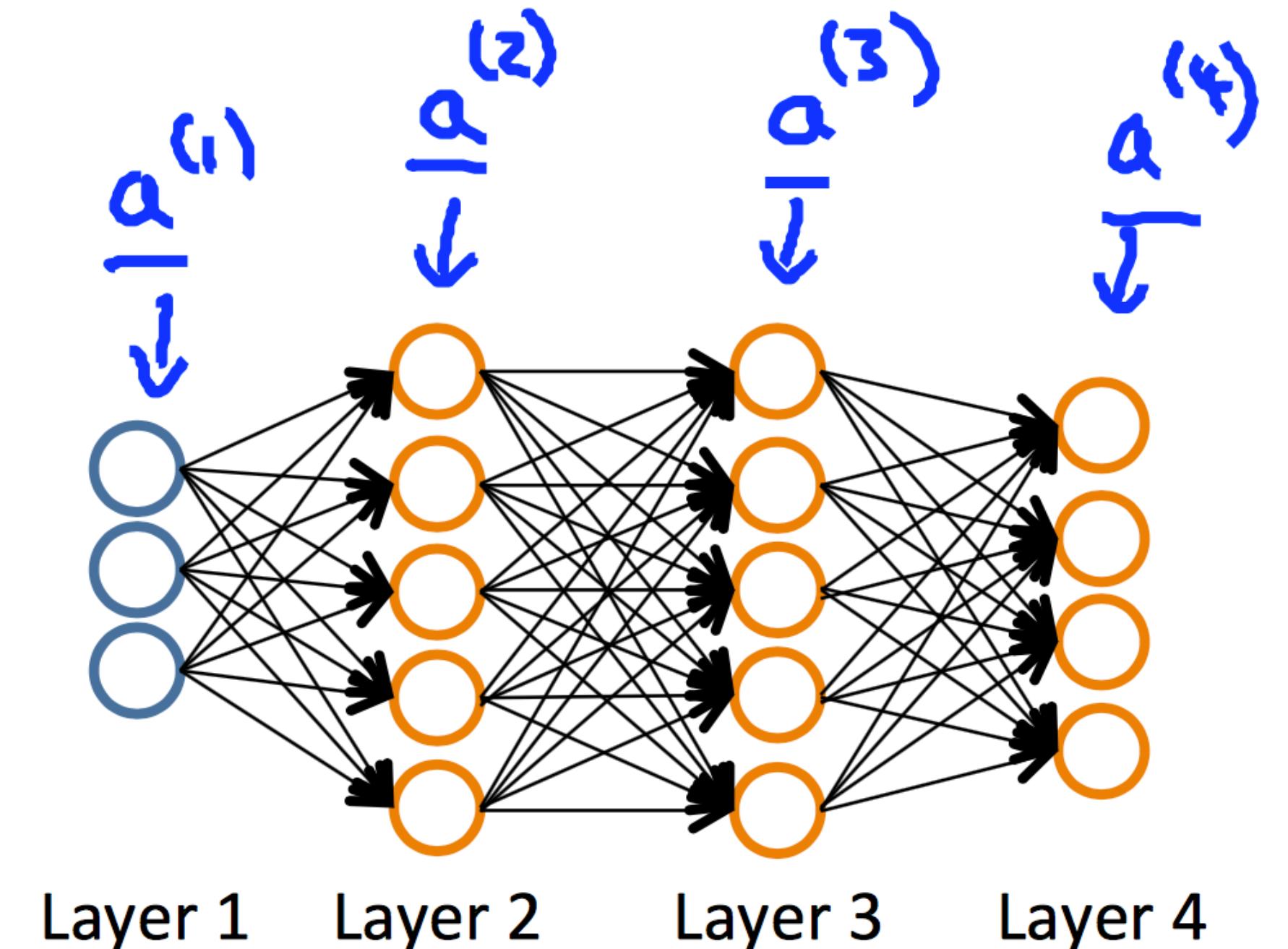
$$\rightarrow \underline{a}^{(2)} = g(z^{(2)}) \quad (\text{add } \underline{a}_0^{(2)})$$

$$\rightarrow z^{(3)} = \Theta^{(2)} \underline{a}^{(2)}$$

$$\rightarrow \underline{a}^{(3)} = g(z^{(3)}) \quad (\text{add } \underline{a}_0^{(3)})$$

$$\rightarrow z^{(4)} = \Theta^{(3)} \underline{a}^{(3)}$$

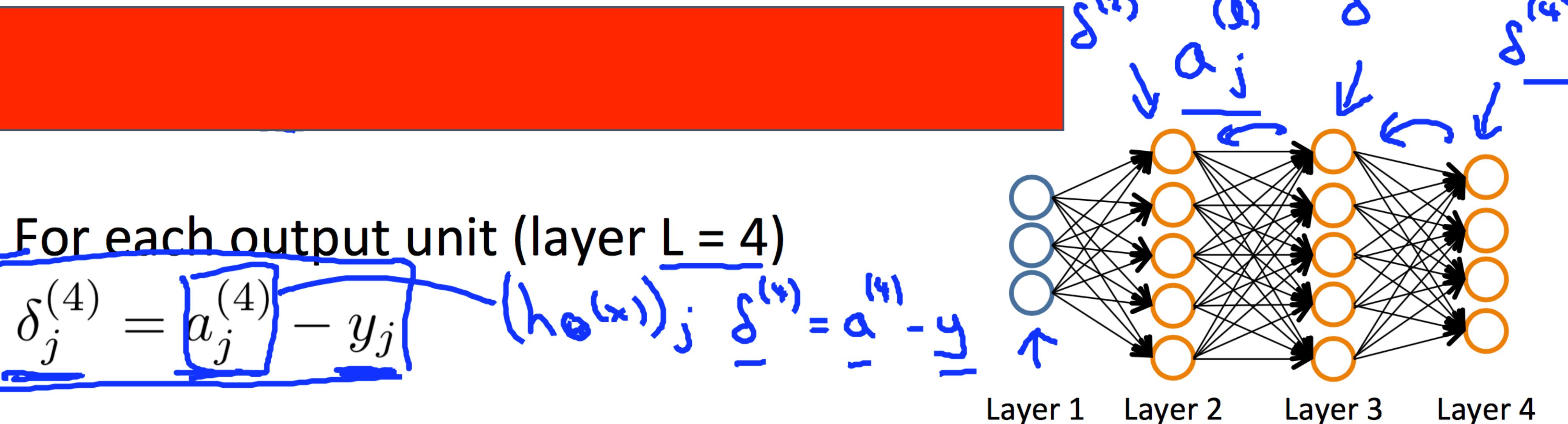
$$\rightarrow \underline{a}^{(4)} = h_{\Theta}(\underline{x}) = g(z^{(4)})$$



Cost Function and Backpropagation

Backpropagation Algorithm

Gradient computation: Backpropagation algorithm



$$\frac{\partial}{\partial \Theta_j^{(l)}} J(\Theta) = a_j^{(l)} \delta_i^{(l+1)}$$

(ignoring λ ; if $\lambda = 0$)

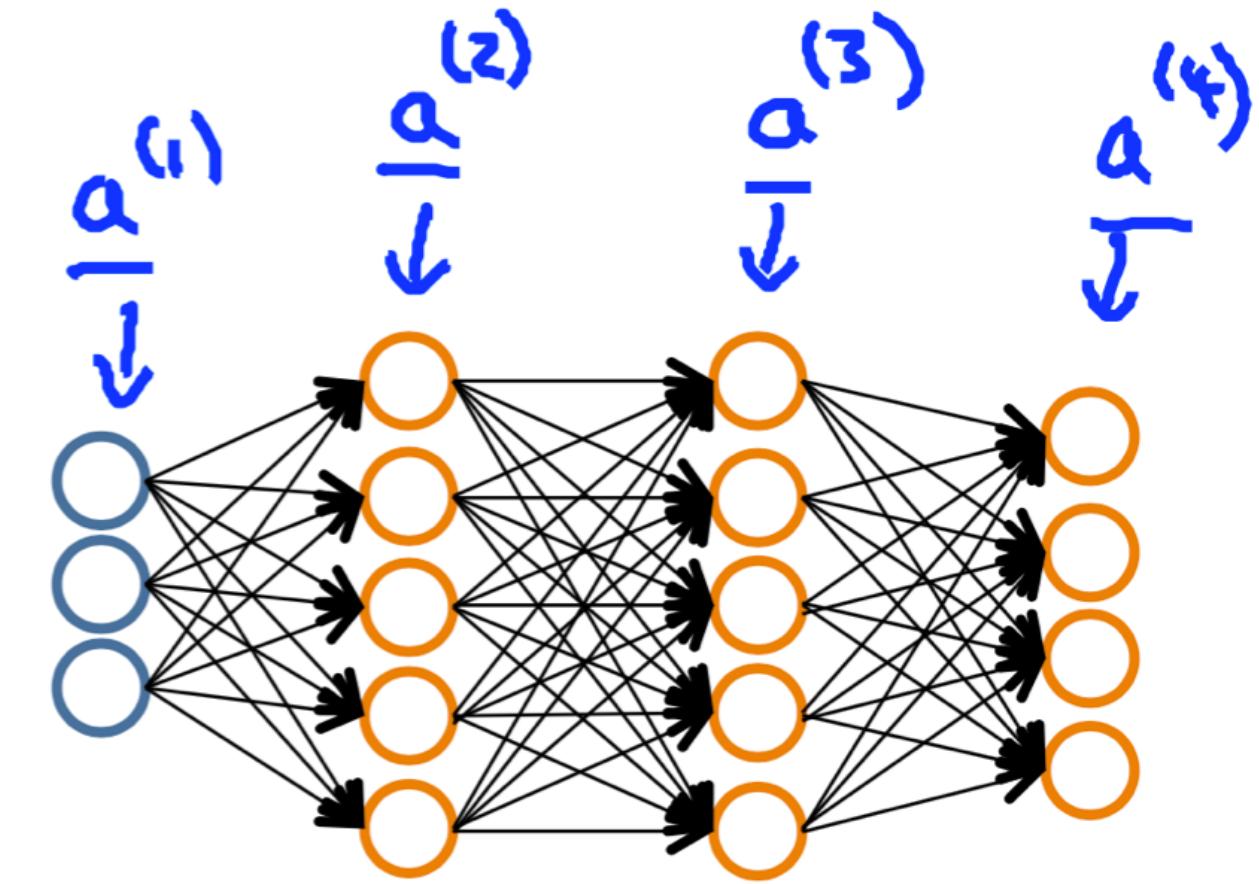
Cost Function and Backpropagation

Backpropagation Algorithm

Backpropagation algorithm

Training set $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$

Set $\Delta_{ij}^{(l)} = 0$ (for all l, i, j).
 $\Delta_{ij}^{(l)}$ (use to compute $\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta)$)

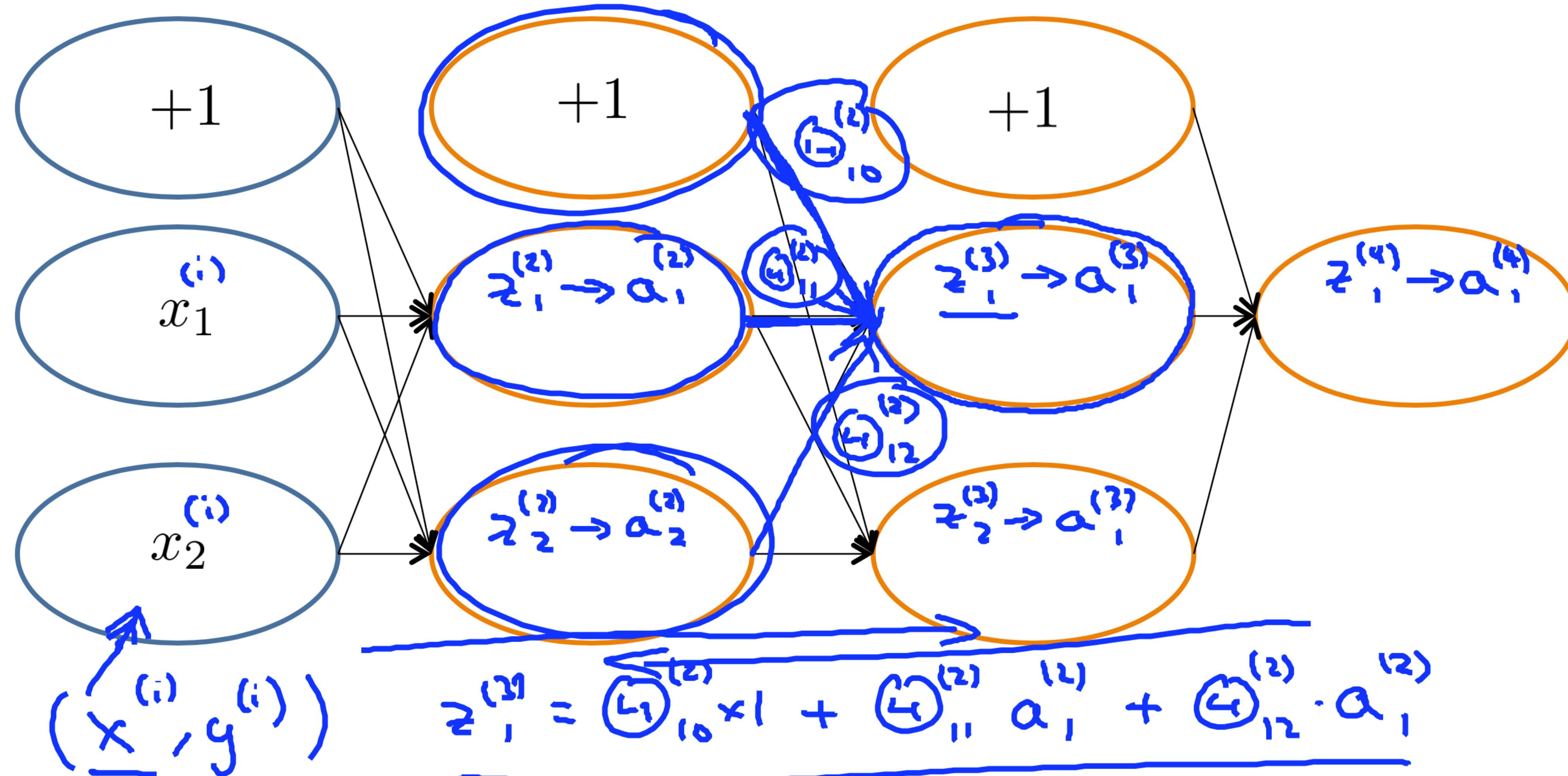


$$D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \Theta_{ij}^{(l)} \quad \text{if } j \neq 0$$
$$D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} \quad \text{if } j = 0$$

Cost Function and Backpropagation

Backpropagation Intuition

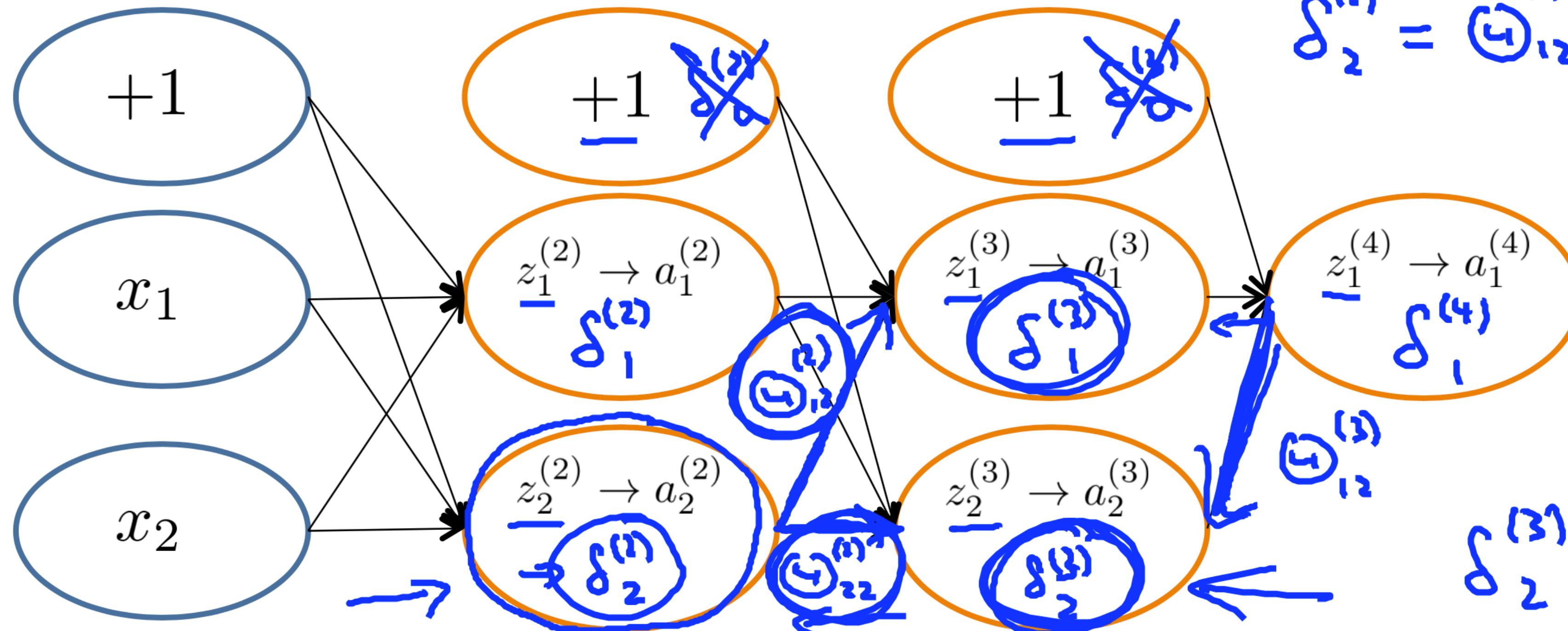
Forward Propagation



Cost Function and Backpropagation

Backpropagation Intuition

Forward Propagation



$$\delta_1^{(4)} = \underline{y_1^{(i)}} - \underline{a_1^{(4)}}$$

$$\delta_2^{(2)} = (\Theta_{12}^{(2)} \delta_1^{(3)}) + (\Theta_{22}^{(2)} \delta_2^{(3)})$$

$$\delta_2^{(3)} = \Theta_{12}^{(3)} \cdot \delta_1^{(4)}$$

$$\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = D_{ij}^{(l)}$$

Neural Networks: Learning

Implementation Note: Unrolling Parameters

Advanced optimization

```
[function [jVal, gradient] = costFunction(theta)
...
optTheta = fminunc(@costFunction, initialTheta, options)
```

Neural Network (L=4):

→ $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$ - matrices (Theta1, Theta2, Theta3)

→ $D^{(1)}, D^{(2)}, D^{(3)}$ - matrices (D1, D2, D3)

To be suitable for optimization algorithm

Neural Networks: Learning

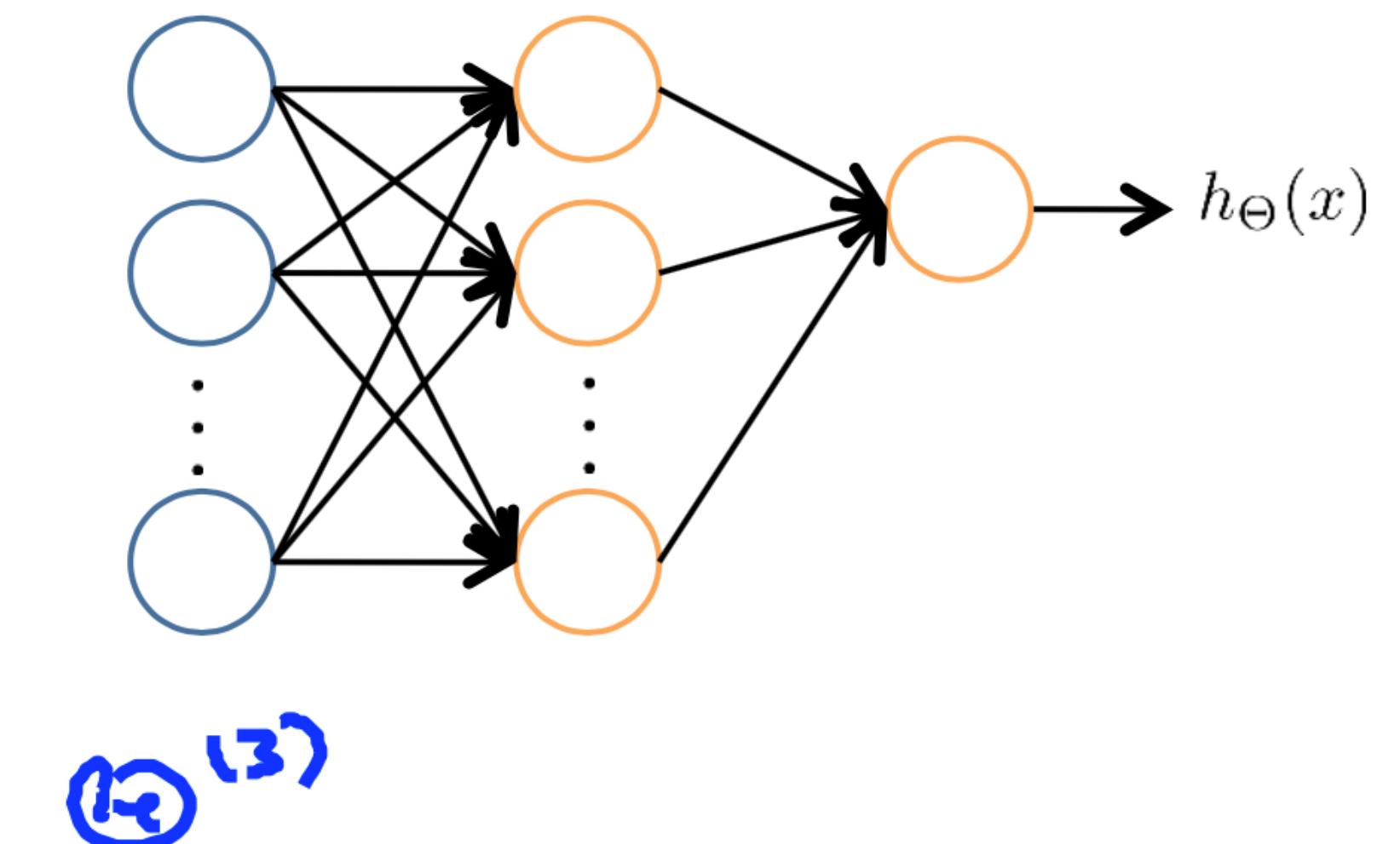
Implementation Note: Unrolling Parameters

Example

$s_1 = 10, s_2 = 10, s_3 = 1$
→ $\Theta^{(1)} \in \mathbb{R}^{10 \times 11}, \Theta^{(2)} \in \mathbb{R}^{10 \times 11}, \Theta^{(3)} \in \mathbb{R}^{1 \times 11}$

→ $D^{(1)} \in \mathbb{R}^{10 \times 11}, D^{(2)} \in \mathbb{R}^{10 \times 11}, D^{(3)} \in \mathbb{R}^{1 \times 11}$

→ **thetaVec** = [Theta1(:); Theta2(:); Theta3(:)];
→ **DVec** = [D1(:); D2(:); D3(:)];



Neural Networks: Learning

Implementation Note: Unrolling Parameters

Learning Algorithm

Learning Algorithm

- Have initial parameters $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$.
- Unroll to get `initialTheta` to pass to
`fminunc (@costFunction, initialTheta, options)`

Learning Algorithm

- Have initial parameters $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$.
- Unroll to get `initialTheta` to pass to
`fminunc (@costFunction, initialTheta, options)`

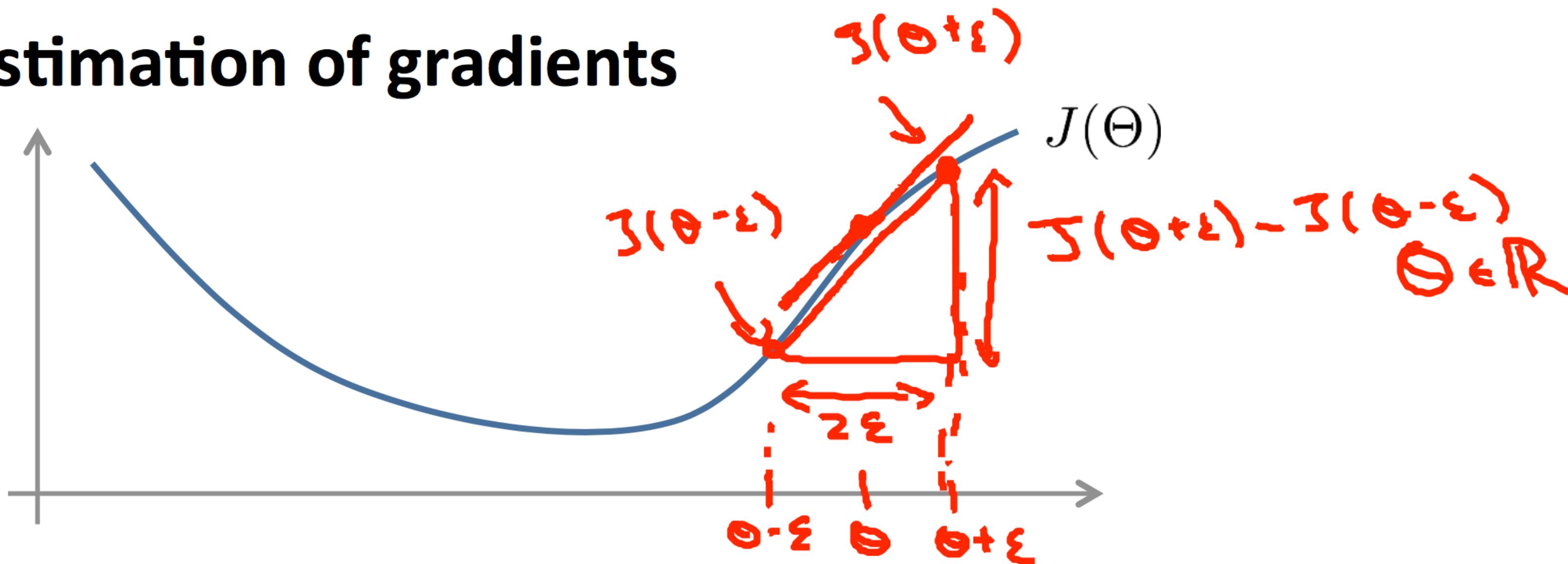
on the right side of the slide.

to get `gradientvec`.

Neural Networks: Learning

Gradient checking

Numerical estimation of gradients



$$\frac{\partial}{\partial \Theta} J(\Theta) \approx$$

$$\frac{J(\Theta + \epsilon) - J(\Theta - \epsilon)}{2\epsilon}$$

$\epsilon = 10^{-4}$

$$\frac{J(\Theta + \epsilon) - J(\Theta)}{\epsilon}$$

To verify forward and backward propagation are moving to the right direction!!!

Implement: gradApprox =
$$\frac{(J(\text{theta} + \text{EPSILON}) - J(\text{theta} - \text{EPSILON}))}{(2 * \text{EPSILON})}$$

Backpropagation in Practice

Parameter vector θ

$\rightarrow \theta \in \mathbb{R}^n$ (E.g. θ is “unrolled” version of $\underline{\Theta^{(1)}}, \underline{\Theta^{(2)}}, \underline{\Theta^{(3)}}$)

$$\rightarrow \theta = [\theta_1, \theta_2, \theta_3, \dots, \theta_n]$$

$$\rightarrow \frac{\partial}{\partial \underline{\theta_1}} J(\theta) \approx \frac{J(\underline{\theta_1 + \epsilon}, \theta_2, \theta_3, \dots, \theta_n) - J(\underline{\theta_1 - \epsilon}, \theta_2, \theta_3, \dots, \theta_n)}{2\epsilon}$$

$$\rightarrow \frac{\partial}{\partial \underline{\theta_2}} J(\theta) \approx \frac{J(\theta_1, \underline{\theta_2 + \epsilon}, \theta_3, \dots, \theta_n) - J(\theta_1, \underline{\theta_2 - \epsilon}, \theta_3, \dots, \theta_n)}{2\epsilon}$$

⋮
⋮
⋮

$$\rightarrow \frac{\partial}{\partial \underline{\theta_n}} J(\theta) \approx \frac{J(\theta_1, \theta_2, \theta_3, \dots, \underline{\theta_n + \epsilon}) - J(\theta_1, \theta_2, \theta_3, \dots, \underline{\theta_n - \epsilon})}{2\epsilon}$$

Backpropagation in Practice

Gradient checking

```
for i = 1:n, ←  
    thetaPlus = theta;  
    thetaPlus(i) = thetaPlus(i) + EPSILON;  
    thetaMinus = theta;  
    thetaMinus(i) = thetaMinus(i) - EPSILON;  
    gradApprox(i) = (J(thetaPlus) - J(thetaMinus))  
                    / (2*EPSILON);  
end;
```

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

Check that gradApprox \approx DVec \leftarrow

Derivative we got

From back prop.

$$\begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_i + \epsilon \\ \vdots \\ \theta_n \end{bmatrix} \rightarrow \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_i - \epsilon \\ \vdots \\ \theta_n \end{bmatrix}$$

Backpropagation in Practice

Gradient checking

Implementation Note:

- Implement backprop to compute DVec (unrolled $D^{(1)}, D^{(2)}, D^{(3)}$).
- Implement numerical gradient check to compute gradApprox.
- Make sure they give similar values. Backprop is more computationally efficient
- Turn off gradient checking. Using backprop code for learning.

Once we are sure that our backpropagation implementation is correct turn it off!!!



Important:

- Be sure to disable your gradient checking code before training your classifier. If you run numerical gradient computation on every iteration of gradient descent (or in the inner loop of **costFunction(...)**) your code will be very slow.

Neural Networks: Learning

Random initialization

Initial value of Θ

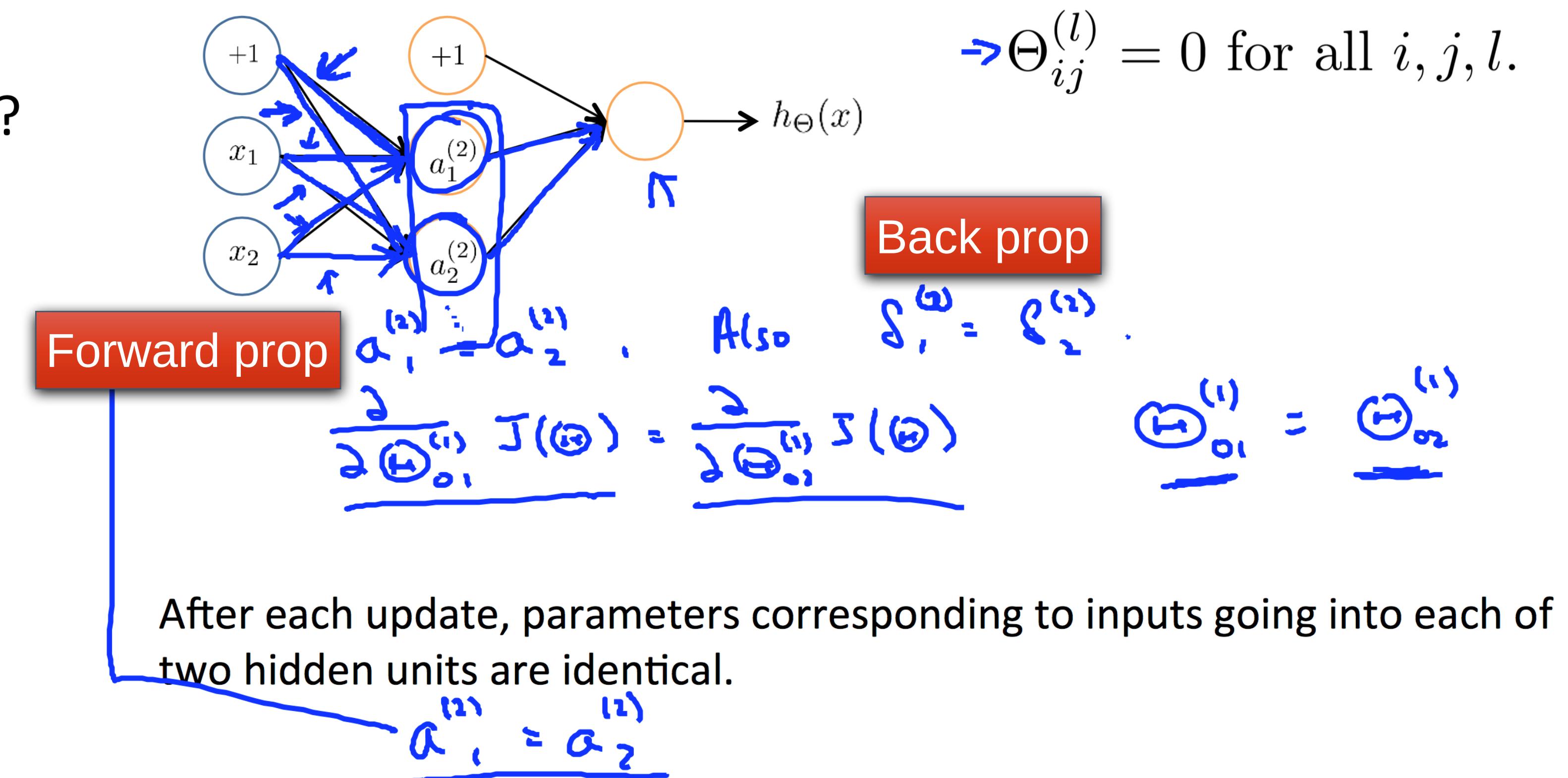
For gradient descent and advanced optimization method, need initial value for Θ .

```
optTheta = fminunc(@costFunction,  
initialTheta, options)
```

Zero initialization

Consider gradient descent

Set initialTheta = zeros(n, 1)?



Neural Networks: Learning

Random initialization

Random initialization: Symmetry breaking

Initialize each $\Theta_{ij}^{(l)}$ to a random value in $[-\epsilon, \epsilon]$
(i.e. $-\epsilon \leq \Theta_{ij}^{(l)} \leq \epsilon$)

*** not the same epsilon
as gradient checking!!!!

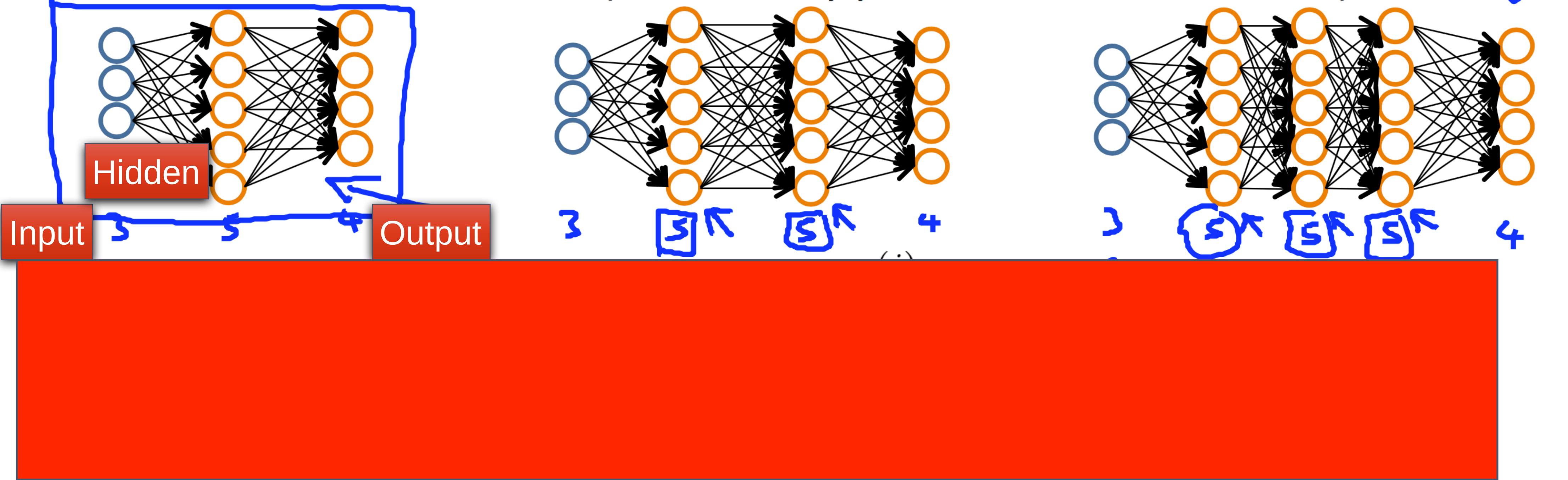
E.g.

Random 10x11 matrix (between 0 and 1)

Neural Networks: Learning

Training a neural network

Pick a network architecture (connectivity pattern between neurons)



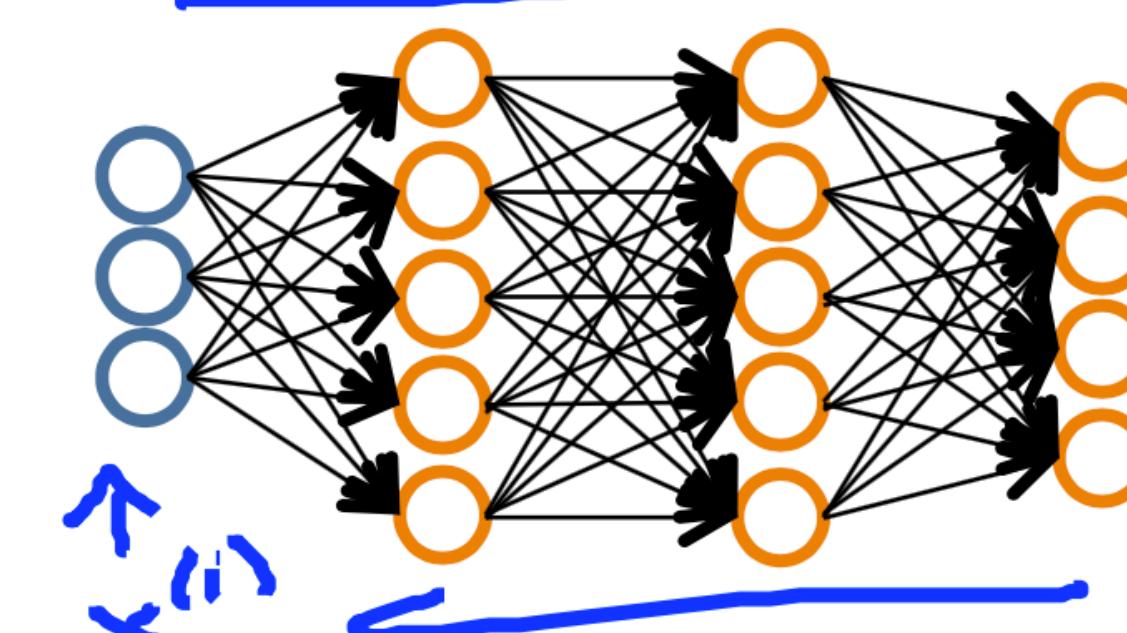
$$y \in \{1, 2, 3, \dots, 10\}$$

~~$y \in \{1, 2, 3, \dots, 10\}$~~

$$y = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \text{ or } \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \begin{bmatrix} 1 & 0 & 0 & \cdots & 0 \end{bmatrix}$$

More hidden units, more computationally expensive!!!!

Neural Networks: Learning

- ▶ **for** $i = 1:m$ { $(x^{(i)}, y^{(i)})$ $(x^{(i)}, y^{(i)})$, ..., $(x^{(m)}, y^{(m)})$
 - Perform forward propagation and backpropagation using example $(x^{(i)}, y^{(i)})$
 - Try to find efficient ways to do for loop over training examples!!
 - (Get activations $a^{(l)}$ and delta terms $\delta^{(l)}$ for $l = 2, \dots, L$).
 - $\Delta^{(l)} := \Delta^{(l)} - \delta^{(l)}(a^{(l)})^T$
 - }
 - ...
 - compute $\frac{\partial}{\partial \Theta^{(2)}} J(\Theta)$.
- 

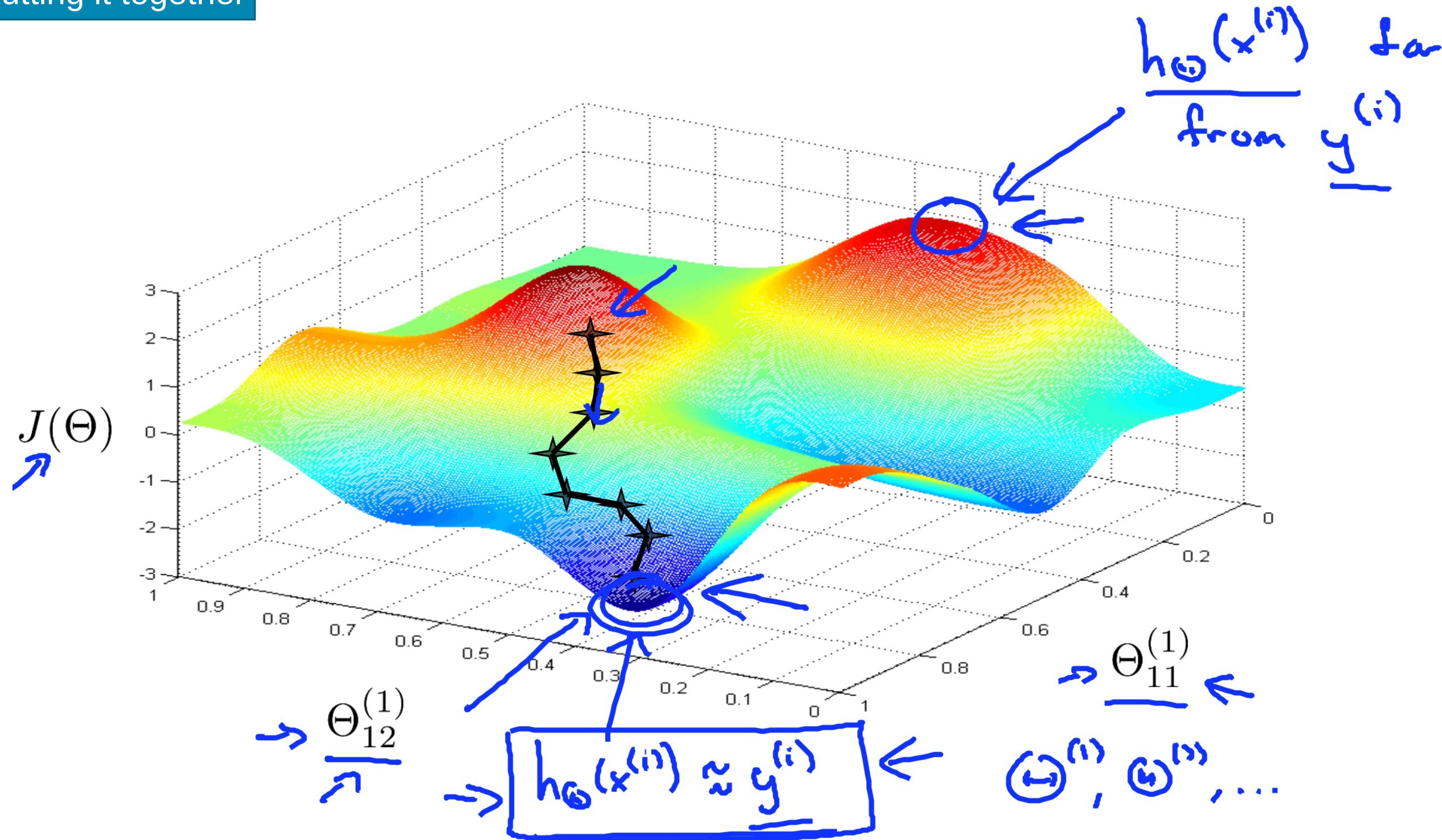
Neural Networks: Learning

$$\nabla \frac{\partial J(\Theta)}{\partial \theta_j}$$

$J(\Theta)$ — non-convex.

Neural Networks: Learning

Putting it together



Neural Networks Applications

Autonomous driving

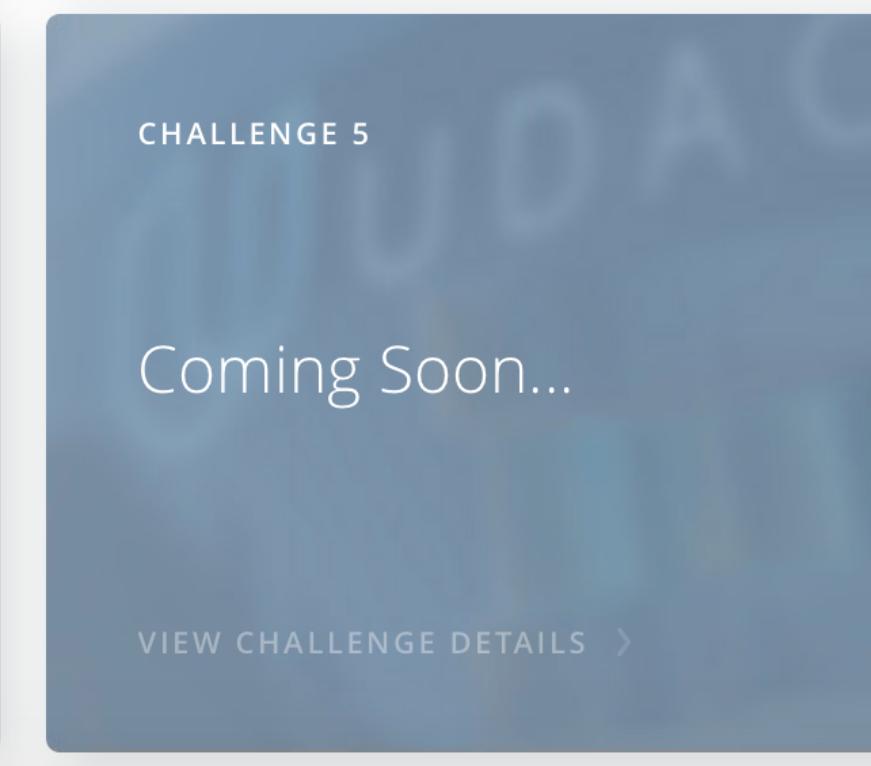
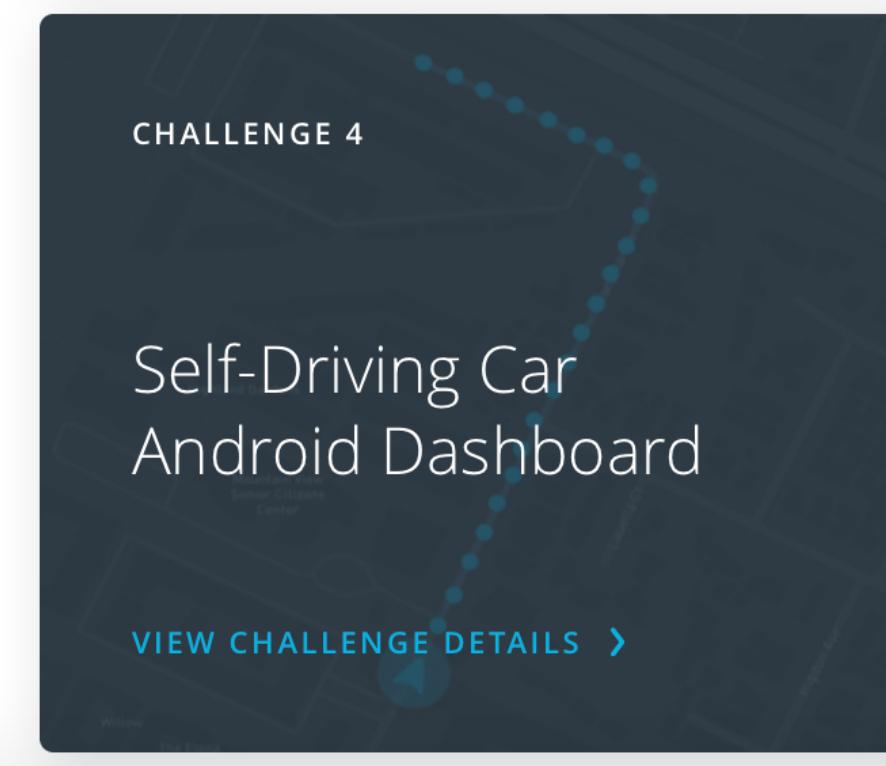
An Open Source Self-Driving Car

Udacity is building an open source self-driving car, and we want your help! Join the effort to create the world's first open source autonomous vehicle. We've broken down the problem into multiple complex challenges, and you or a team can compete to have your solution run in a real self-driving car.

LEARN MORE

JOIN SLACK

GITHUB



Perceptron

Perceptron

The Perceptron is a primitive type of neural network that learns weights for input attributes and transfers the weighted inputs into a network output or prediction.

This recipes shows the fitting of a Perceptron algorithm on the diabetes dataset.

```
# Perceptron
import numpy as np
from sklearn import datasets
from sklearn.linear_model import Perceptron
# load the diabetes datasets
dataset = datasets.load_diabetes()
# fit a Perceptron model to the data
model = Perceptron()
model.fit(dataset.data, dataset.target)
print(model)
# make predictions
expected = dataset.target
predicted = model.predict(dataset.data)
# summarize the fit of the model
mse = np.mean((predicted-expected)**2)
print(mse)
print(model.score(dataset.data, dataset.target))
```

Advice for Applying Machine Learning

how to tell when a learning algorithm is doing poorly?

describe the 'best practices' for how to 'debug' your learning algorithm and go about improving its performance.

We will also be covering machine learning system design.

you'll need to first understand where the biggest improvements can be made.

In these lessons, we discuss how to understand the performance of a machine learning system with multiple parts, and also how to deal with skewed data.

When you're applying machine learning to real problems, a solid grasp of this week's content will easily save you a large amount of work.

Evaluating a Learning Algorithm

Debugging a learning algorithm:

Suppose you have implemented regularized linear regression to predict housing prices.

$$\rightarrow J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^m \theta_j^2 \right]$$