

MACHINE LEARNING

AN INTRODUCTION

Christian Treffs

WHAT IS MACHINE LEARNING?

DEFINITION

Machine Learning is ...

...field of study that gives computers the ability to learn without being explicitly programmed.

~ Arthur Samuel (1959)

... a well-posed learning problem: A computer program is said to learn from experience E with respect to some task T and some performance measure P , if its performance on T , as measured by P , improves with experience E .

~ Tom Mitchell (1998)

... is the scientific study of algorithms and statistical models that computer systems use to effectively perform a specific task without using explicit instructions, relying on patterns and inference instead.

It is seen as a subset of artificial intelligence.

~ Wikipedia

MACHINE LEARNING FIELDS

SUPERVISED LEARNING (CLASSIFICATION, REGRESSION)

- Linear/Polynomial & Logistic regression
- Neural networks (Forward-/Backpropagation)
- Support-Vector Machines (SVMs)

UNSUPERVISED LEARNING

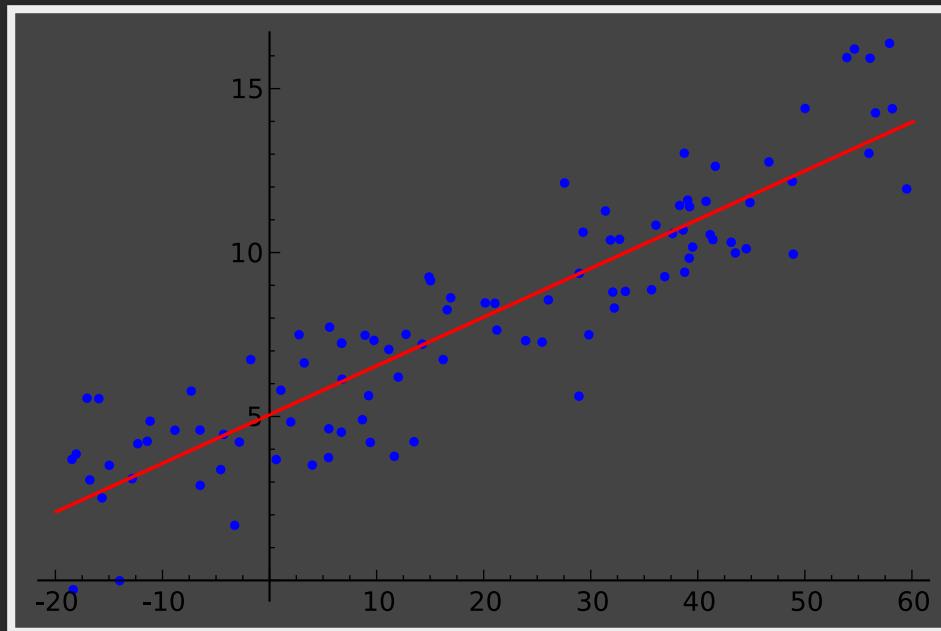
- Clustering (k-Means)
- Dimensionality reduction (Principal Component Analysis (PCA))
- Anomaly detection
- Recommender Systems

SUPERVISED LEARNING

*... is the ML task of learning a function
that maps an input to an output based on
example input-output pairs.*

*It infers a function from labeled training
data consisting of a set of training
examples.*

LINEAR/POLINOMIAL REGRESSION



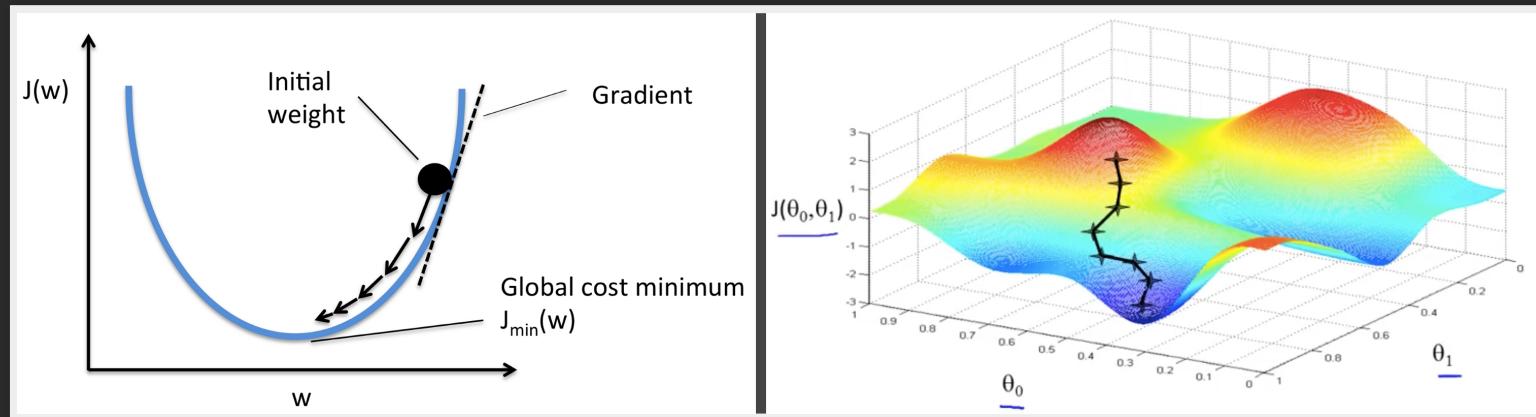
- Housing price prediction
 - linear: size feet => price \$
 - polynomial: (size feet, num rooms, num floors, age of home, ...) => price \$

LINEAR REGRESSION

Hypothesis:	$h_{\theta}(x) = \theta_0 + \theta_1 x$	Gradient descent algorithm	Linear Regression Model
Parameters:	θ_0, θ_1	repeat until convergence {	$h_{\theta}(x) = \theta_0 + \theta_1 x$
Cost Function:	$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$	$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$ (for $j = 1$ and $j = 0$)	$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$
Goal:	$\underset{\theta_0, \theta_1}{\text{minimize}} J(\theta_0, \theta_1)$	}	
	.		

- training set -> learning algorithm -> cost-function w/ hypothesis -> estimation/prediction
- Idea: find θ 's to minimize cost-function
- Learning rate α -> adjust to tune gradient descent

GRADIENT DESCENT & LEARNING RATE α



- Learning rate α
 - too small: slow convergence
 - too large: $J(\theta)$ may not decrease every iteration; may not converge
- Converged: if $J(\theta)$ decreases by less than 10^{-3} in 1 iteration

TECHNIQUE: FEATURE SCALING

- Make sure features are on a similar scale
- Get every feature into approximately $-1 \leq x_i \leq 1$ range.
- Mean normalization:
 - replace x_i with $x_i - \mu_i$ to make features have approx. 0 mean

GRADIENT DESCENT ALTERNATIVE: NORMAL EQUATION

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

$(X^T X)^{-1}$ is inverse of matrix $X^T X$.

Set $A = \underline{X^T X}$

$$(X^T X)^{-1} = A^{-1}$$

Octave: $\text{pinv}(X' * X) * X' * y$

$$\text{pinv}(X^T * X) * X^T * y$$

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

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

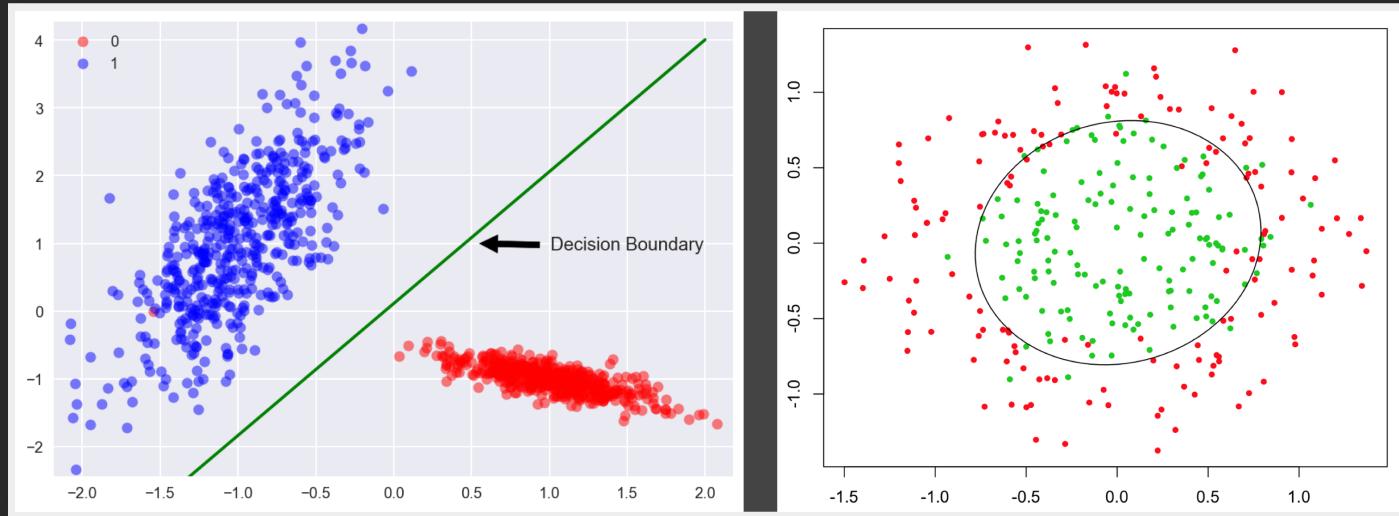
X' X^T
Feature scaling
 $0 \leq x_1 \leq 1$
 $0 \leq x_2 \leq 1000$
 $0 \leq x_3 \leq 10^{-5}$ ✓

m training examples, n features.

<u>Gradient Descent</u> <ul style="list-style-type: none"> → Need to choose α. → Needs many iterations. • Works well even when <u>n</u> is large. <p style="text-align: center;">$n = 10^6$</p>	<u>Normal Equation</u> <ul style="list-style-type: none"> → No need to choose α. → Don't need to iterate. • Need to compute $\boxed{(X^T X)^{-1}}$ $\frac{n \times n}{n \times n} O(n^3)$ • Slow if <u>n</u> is very large. <p style="text-align: center;">$n = 100$ $n = 1000$... $n = 10000$</p>
---	--

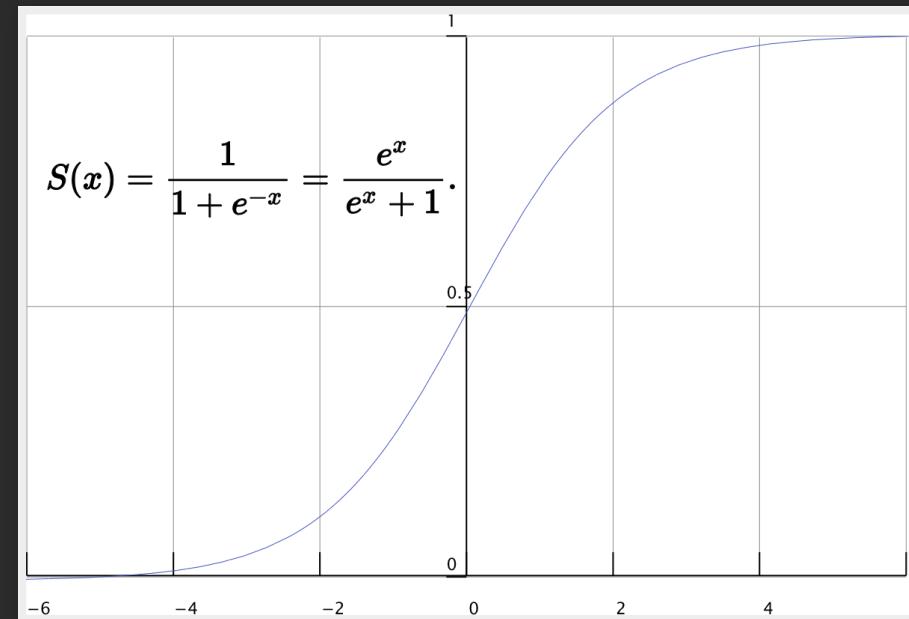
- = Method to solve for α analytically
- + more efficient
- - consumes more memory

LOGISTIC REGRESSION



- Classification
 - Email: Spam / Not Spam?
 - Online Transactions: Fraudulent Yes/No?
 - Tumor: Malignant/Benign?

SIGMOID FUNCTION



COST FUNCTION & GRADIENT DESCENT

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

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

Repeat {

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

}

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

(simultaneously update all θ_j)

Logistic regression

$$h_\theta(x) = g(\theta^T x)$$

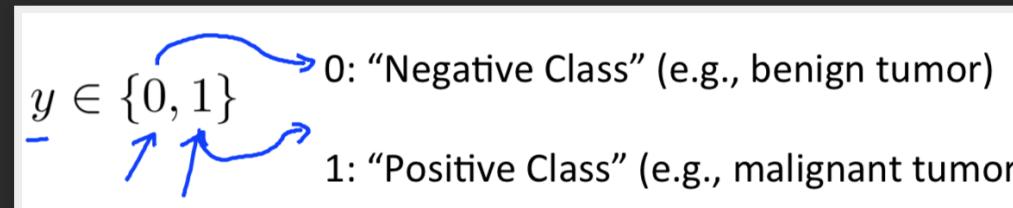
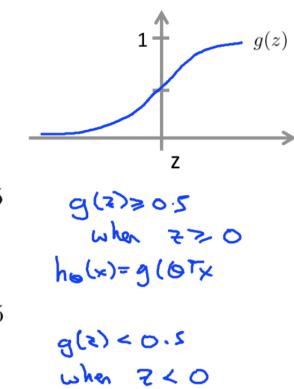
$$g(z) = \frac{1}{1+e^{-z}}$$

Suppose predict " $y = 1$ " if $h_\theta(x) \geq 0.5$

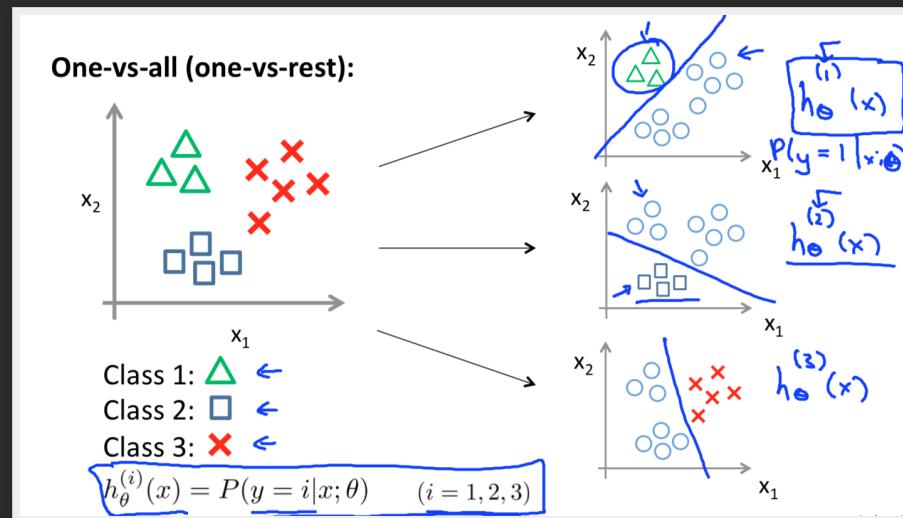
$$\theta^T x \geq 0$$

predict " $y = 0$ " if $h_\theta(x) < 0.5$

$$\theta^T x < 0$$



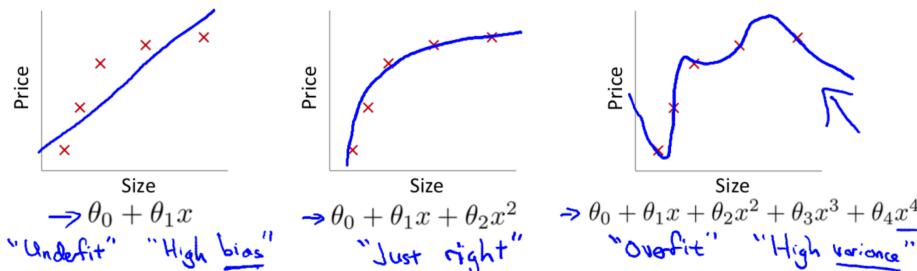
MULTI-CLASS CLASSIFICATION: ONE-VS-ALL



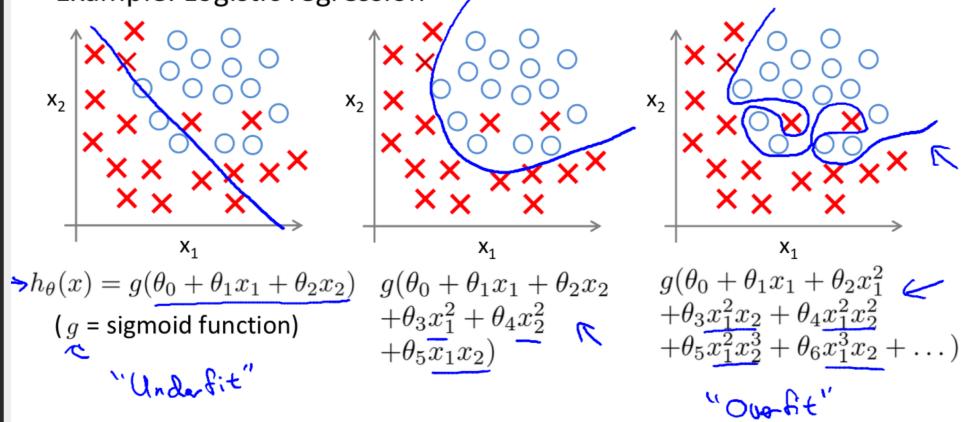
- Classification with multiple segments
 - Email foldering/tagging: Work, Friends, Family, Hobby
 - Medical diagrams: Not ill, Cold, Flu
 - Weather: Sunny, Cloudy, Rain, Snow

REGULARIZATION

Example: Linear regression (housing prices)



Example: Logistic regression



- bias vs. variance
 - high bias - underfit
 - high variance - overfit
- to address overfitting:
 - reduce features
 - regularization

REGULARIZED LINEAR REGRESSION

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

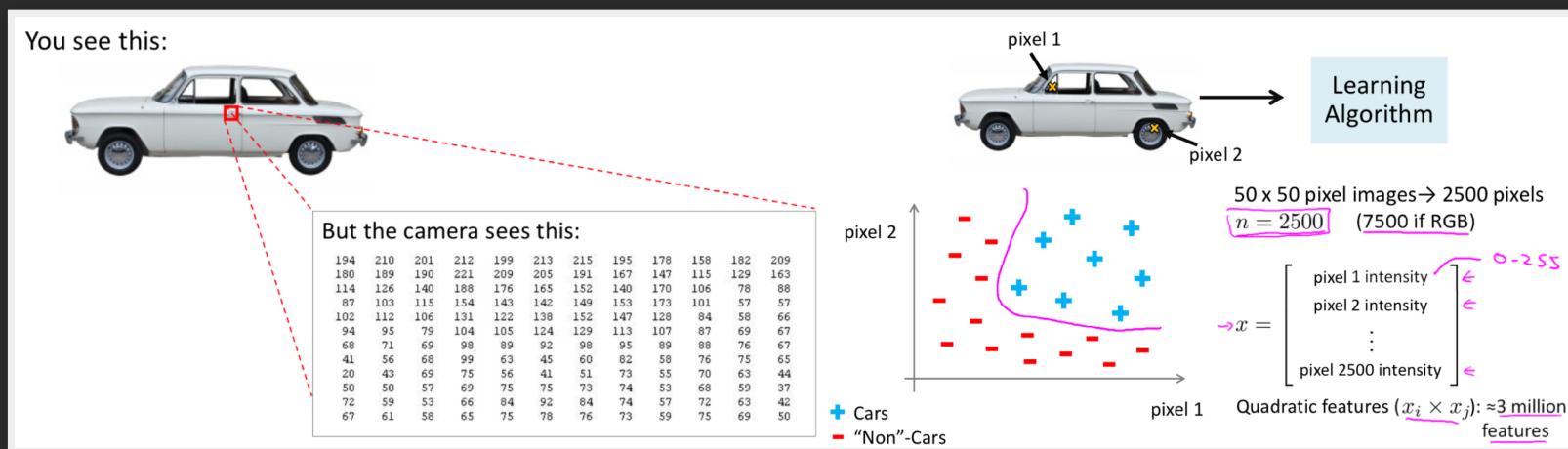
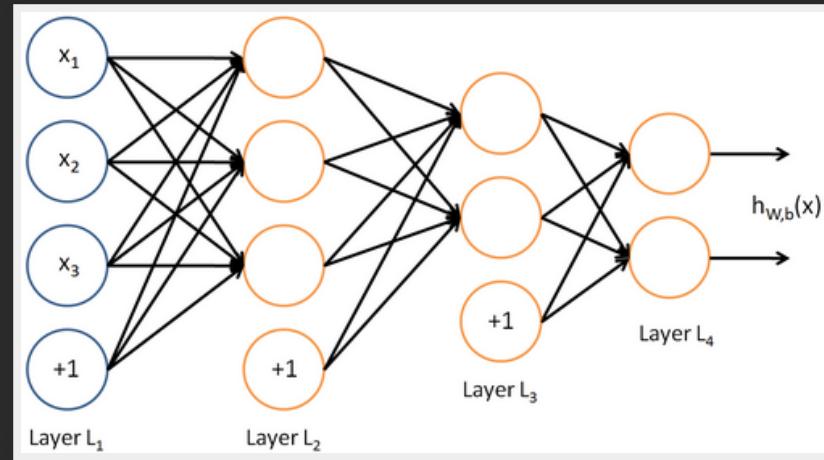
$\min_{\theta} J(\theta)$

REGULARIZED LOGISTIC REGRESSION

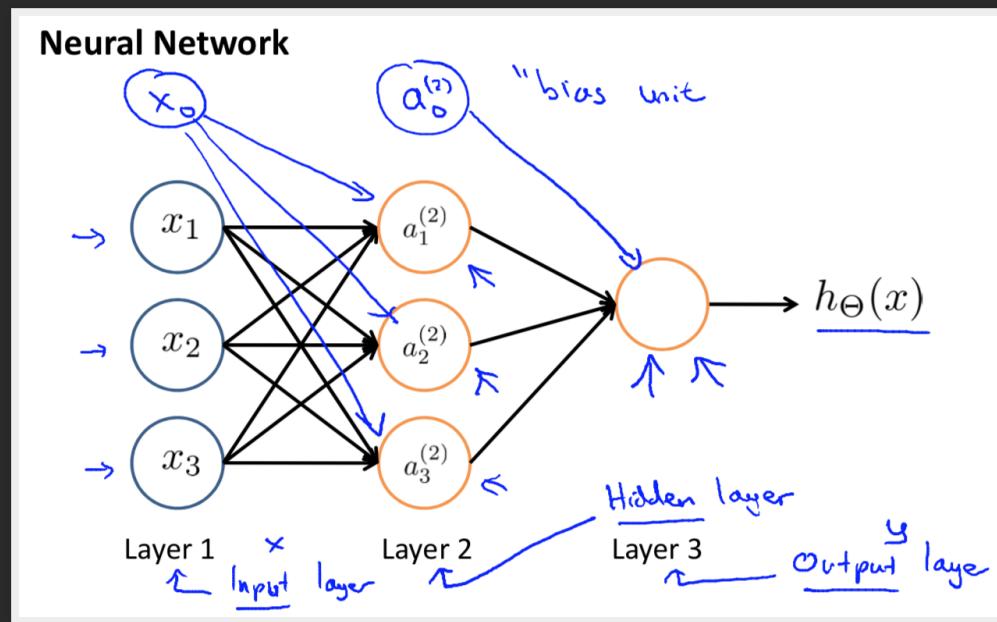
$$\rightarrow J(\theta) = - \left[\frac{1}{m} \sum_{i=1}^m y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

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

NEURAL NETWORKS



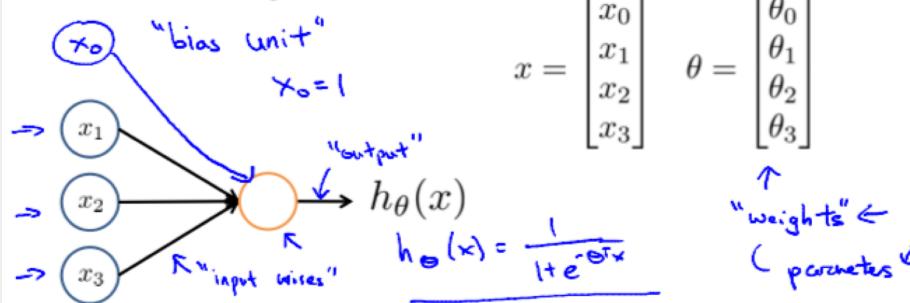
NEURAL NETWORK STRUCTURE



- $\alpha_i^{(j)}$ “activation” of unit i in layer j
- non-linear hypotheses
- computer vision: car detection (car / not a car)
- Autonomous driving

NEURAL NETWORK ACTIVATION

Neuron model: Logistic unit



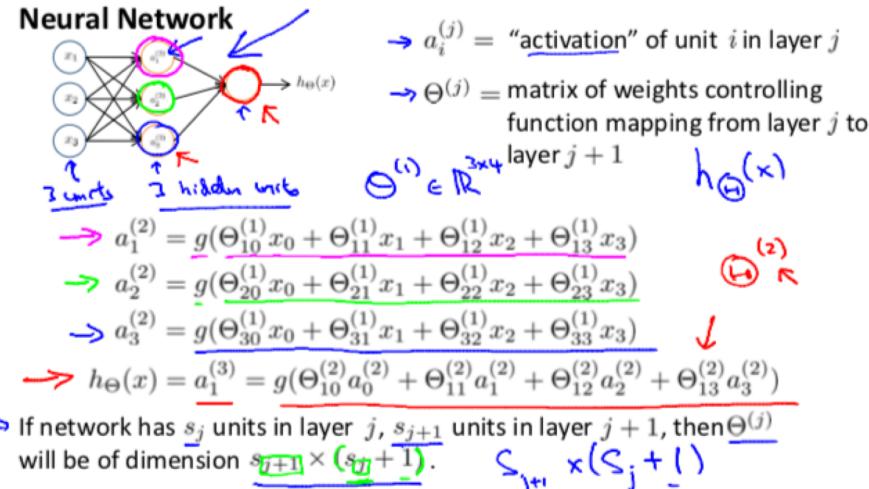
Sigmoid (logistic) activation function.

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

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

"weights" ↗
(parameters) ↗

Neural Network

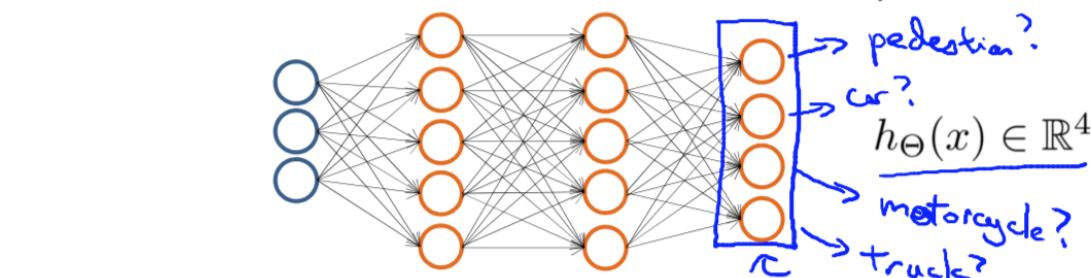


DEMO

- AND / OR

MULTIPLE OUTPUT NEURAL NETWORK

Multiple output units: One-vs-all.



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

COST FUNCTION NEURAL NETWORK

Neural network:

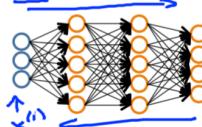
$$h_{\Theta}(x) \in \mathbb{R}^K \quad (h_{\Theta}(x))_i = i^{th} \text{ output}$$

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

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

TRAINING NEURAL NETWORK: FORWARD & BACK-PROPAGATION

Training a neural network

- 1. Randomly initialize weights
 - 2. Implement forward propagation to get $h_{\Theta}(x^{(i)})$ for any $x^{(i)}$
 - 3. Implement code to compute cost function $J(\Theta)$
 - 4. Implement backprop to compute partial derivatives $\frac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta)$
 - for $i = 1:m$ { $(x^{(1)}, y^{(1)})$ $(x^{(2)}, y^{(2)})$, ..., $(x^{(m)}, y^{(m)})$
 - Perform forward propagation and backpropagation using example $(x^{(i)}, y^{(i)})$
 - Get activations $a^{(l)}$ and delta terms $\delta^{(l)}$ for $l = 2, \dots, L$.
 - $\Delta^{(k)} := \Delta^{(k)} + \delta^{(k)} \cdot (a^{(k)})^T$
 - ...
 - compute $\frac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta)$.
- 

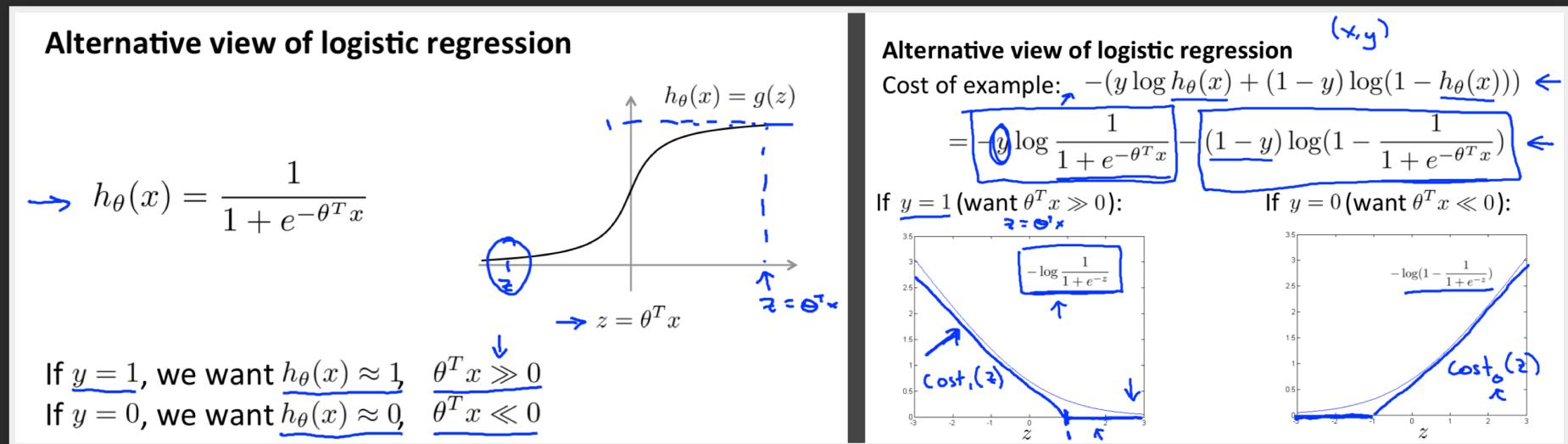
Training a neural network

- 5. Use gradient checking to compare $\frac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta)$ computed using backpropagation vs. using numerical estimate of gradient of $J(\Theta)$.
- Then disable gradient checking code.
- 6. Use gradient descent or advanced optimization method with backpropagation to try to minimize $J(\Theta)$ as a function of parameters Θ

$$\frac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta) \quad \text{↑}$$

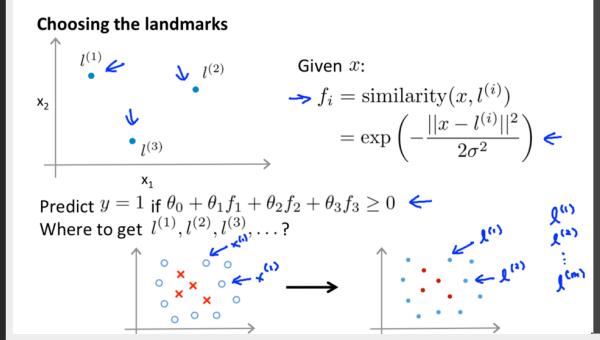
$J(\Theta)$ — non-convex.

SUPPORT VECTOR MACHINES (SVM)



- Alternative view of logistic regression
- large margin classification

KERNELS



SVM with Kernels

- Given $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})$,
- choose $l^{(1)} = x^{(1)}, l^{(2)} = x^{(2)}, \dots, l^{(m)} = x^{(m)}$.

Given example x :

$$\begin{aligned} f_1 &= \text{similarity}(x, l^{(1)}) \\ f_2 &= \text{similarity}(x, l^{(2)}) \\ &\vdots \\ f_m &= \text{similarity}(x, l^{(m)}) \end{aligned}$$

$f = \begin{bmatrix} f_0 \\ f_1 \\ f_2 \\ \vdots \\ f_m \end{bmatrix}$ $f_0 = 1$

For training example $(x^{(i)}, y^{(i)})$:

$$\begin{aligned} x^{(i)} &\rightarrow \begin{bmatrix} f_1^{(i)} \\ f_2^{(i)} \\ \vdots \\ f_m^{(i)} \end{bmatrix} = \begin{bmatrix} \text{sim}(x^{(i)}, l^{(1)}) \\ \text{sim}(x^{(i)}, l^{(2)}) \\ \vdots \\ \text{sim}(x^{(i)}, l^{(m)}) \end{bmatrix} \\ &= \exp\left(-\frac{\|x^{(i)} - l^{(i)}\|^2}{2\sigma^2}\right) = 1 \end{aligned}$$

$f^{(i)} = \begin{bmatrix} f_0^{(i)} \\ f_1^{(i)} \\ f_2^{(i)} \\ \vdots \\ f_m^{(i)} \end{bmatrix} = \begin{bmatrix} 1 \\ f_1^{(i)} \\ f_2^{(i)} \\ \vdots \\ f_m^{(i)} \end{bmatrix}$

SVM with Kernels

Hypothesis: Given x , compute features $f \in \mathbb{R}^{m+1}$ $\Theta \in \mathbb{R}^{n+1}$

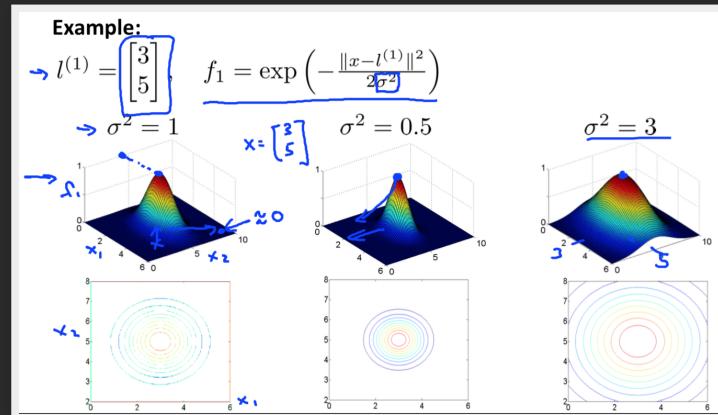
Predict "y=1" if $\theta^T f \geq 0$

Training:

$$\min_{\theta} C \sum_{i=1}^m y^{(i)} \underbrace{\text{cost}_1(\theta^T f^{(i)})}_{\theta^T f^{(i)} < 0} + (1 - y^{(i)}) \underbrace{\text{cost}_0(\theta^T f^{(i)})}_{\theta^T f^{(i)} \geq 0} + \frac{1}{2} \sum_{j=1}^n \theta_j^2 \rightarrow \Theta$$

$\Theta = \Theta^T \Theta \leftarrow \Theta = \begin{bmatrix} \Theta_0 \\ \vdots \\ \Theta_n \end{bmatrix}$ (ignoring Θ_0)

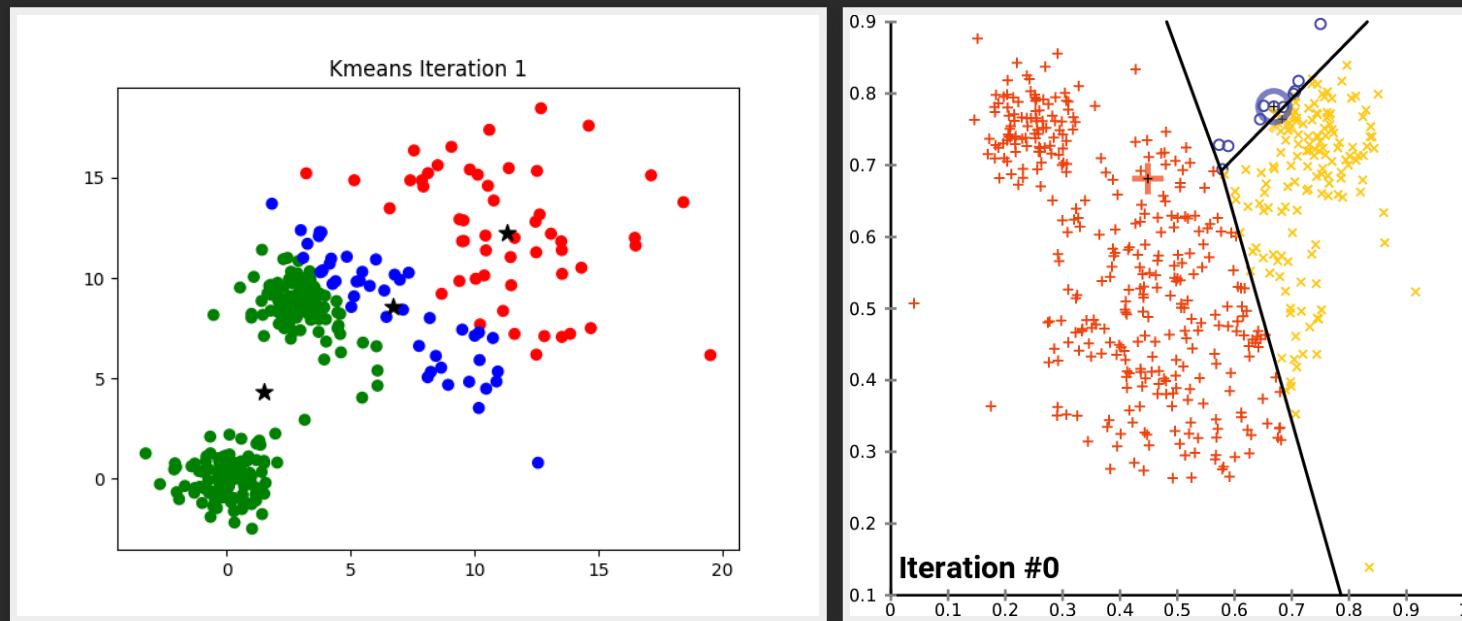
$m = 10,000$



UNSUPERVISED LEARNING

- Organize computing clusters
- Social network analysis
- Market segmentation
- Astronomical data analysis
- aka Clustering

K-MEANS CLUSTERING



Examples:

- T-shirt sizes (S, M, L, XL)

K-MEANS CLUSTERING ALGORITHM

Input:

- K (number of clusters)
- Training set $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$

K-means algorithm



```

Randomly initialize  $K$  cluster centroids  $\mu_1, \mu_2, \dots, \mu_K \in \mathbb{R}^n$ 
Repeat {
    Cluster assignment step:
    for  $i = 1$  to  $m$ 
         $c^{(i)}$  := index (from 1 to  $K$ ) of cluster centroid closest to  $x^{(i)}$ 
    for  $k = 1$  to  $K$ 
         $\rightarrow \mu_k$  := average (mean) of points assigned to cluster  $k$ 
    }
}

```

Note: Handwritten annotations include: "Cluster assignment step" above the first loop, "Move centroid" above the second loop, and a formula for μ_2 : $\mu_2 = \frac{1}{4} [x^{(1)} + x^{(5)} + x^{(6)} + x^{(10)}] \in \mathbb{R}^n$.

K-means optimization objective

$\rightarrow c^{(i)}$ = index of cluster ($1, 2, \dots, K$) to which example $x^{(i)}$ is currently assigned

$\rightarrow \mu_k$ = cluster centroid k ($\mu_k \in \mathbb{R}^n$)

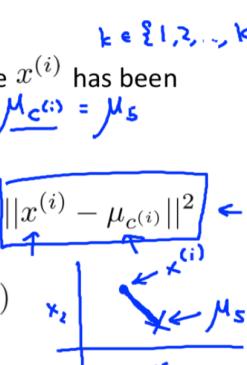
$\mu_{c^{(i)}}$ = cluster centroid of cluster to which example $x^{(i)}$ has been assigned

Optimization objective:

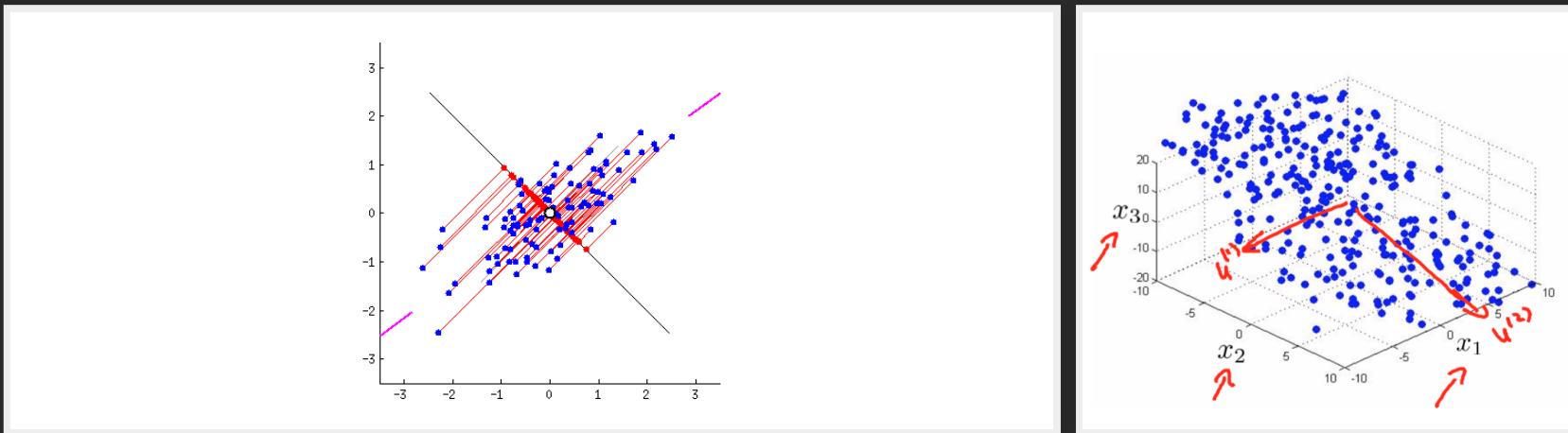
$$\rightarrow J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu_{c^{(i)}}\|^2$$

$$\min_{\substack{c^{(1)}, \dots, c^{(m)}, \\ \mu_1, \dots, \mu_K}} J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$$

Note: Handwritten annotations include: "K" above the range of k, "k in {1, 2, ..., K}" above the mu_k definition, "x^{(i)} -> S" above the mu_{c(i)} definition, and "Distortion" below the second equation.



DIMENSIONALITY REDUCTION: PRINCIPAL COMPONENT ANALYSIS (PCA)



- Data compression (Reduce memory needed to store data)
- Speed up other learning algorithms
- Data visualization (reduce to 2D or 3D)
 - > Reduce data from n -dimensions to k -dimensions

PCA - STEPS

1. Data preprocessing: feature scaling/mean normalization

2. Calculate covariance matrix

3. Compute eigenvectors of matrix Sigma

([U, __, __] = svd (Sigma) ;)

4. z = U' * x where

\$x\$:= original data \$z\$:= dim. reduced data

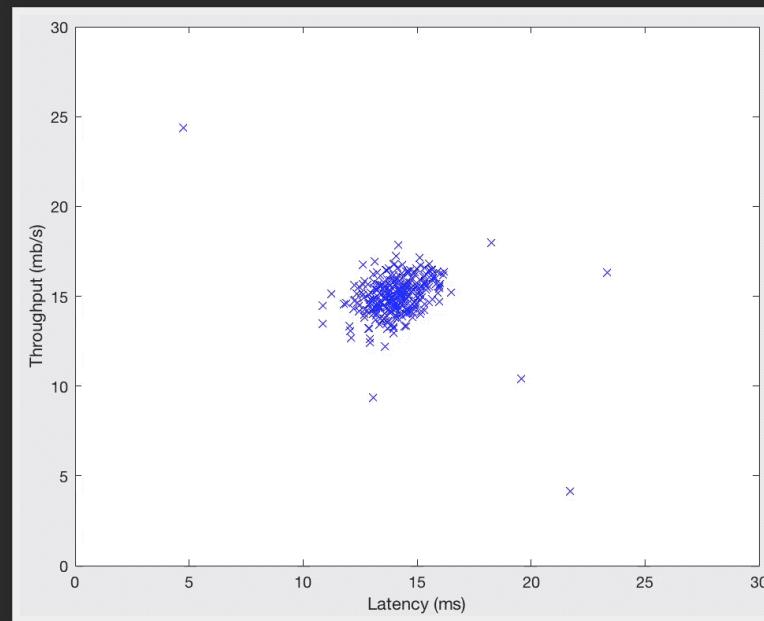
5. principal components \$k\$ (~ new axes in the data)

Choose \$k\$ to be smallest value so that:

$$\frac{\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{approx}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} \leq \frac{0.01}{0.05}$$



ANOMALY DETECTION



- Fraud detection
- Manufacturing (e.g. engines)
- Monitoring machines in data center

ANOMALY DETECTION - GAUSSIAN DISTRIBUTION

Anomaly detection algorithm

- 1. Choose features x_j that you think might be indicative of anomalous examples.
- 2. Fit parameters $\mu_1, \dots, \mu_n, \sigma_1^2, \dots, \sigma_n^2$

$$\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$$

$$\sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$$
- 3. Given new example x , compute $p(x)$:

$$p(x) = \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi}\sigma_j} \exp\left(-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}\right)$$
 Anomaly if $p(x) < \varepsilon$

Anomaly detection example

$p(x) = p(x_1; \mu_1, \sigma_1^2) \times p(x_2; \mu_2, \sigma_2^2)$

$\mu_1 = 5, \sigma_1 = 2$
 $\mu_2 = 3, \sigma_2 = 1$

$\varepsilon = 0.02$
 $p(x_{test}^{(1)}) = 0.0426 \geq \varepsilon$
 $p(x_{test}^{(2)}) = 0.0021 < \varepsilon$

Anomaly detection with the multivariate Gaussian

1. Fit model $p(x)$ by setting

$$\mu = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)(x^{(i)} - \mu)^T$$
2. Given a new example x , compute

$$p(x) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right)$$

Flag an anomaly if $p(x) < \varepsilon$

RECOMMENDER SYSTEMS



- predict movie ratings (netflix, hulu, ...)
- predict music by style (spotify radio, last.fm, pandora)
- predict products that customers would buy (amazon)

CONTENT-BASED FILTERING

CONTENT-BASED FILTERING

Read by user → Similar articles → Recommended to user

Content-based recommender systems

Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)
Love at last	1	5	0	0
Romance forever	5	?	?	0
Cute puppies of love	?	4	0	?
Nonstop car chases	0	0	5	4
Swords vs. karate	0	0	5	?

$n_u = 4, n_m = 5$

$x^{(1)} = \begin{bmatrix} 1 \\ 0.9 \\ 0 \end{bmatrix}$

$\theta^{(1)} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$

$(\theta^{(1)})^T x^{(1)} = 5 \times 0.99 = 4.95$

Optimization algorithm:

$$\min_{\theta^{(1)}, \dots, \theta^{(n_u)}} \frac{1}{2} \sum_{j=1}^{n_u} \sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (\theta_k^{(j)})^2$$

$J(\theta^{(1)}, \dots, \theta^{(n_u)})$

Gradient descent update:

$$\theta_k^{(j)} := \theta_k^{(j)} - \alpha \sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)}$$

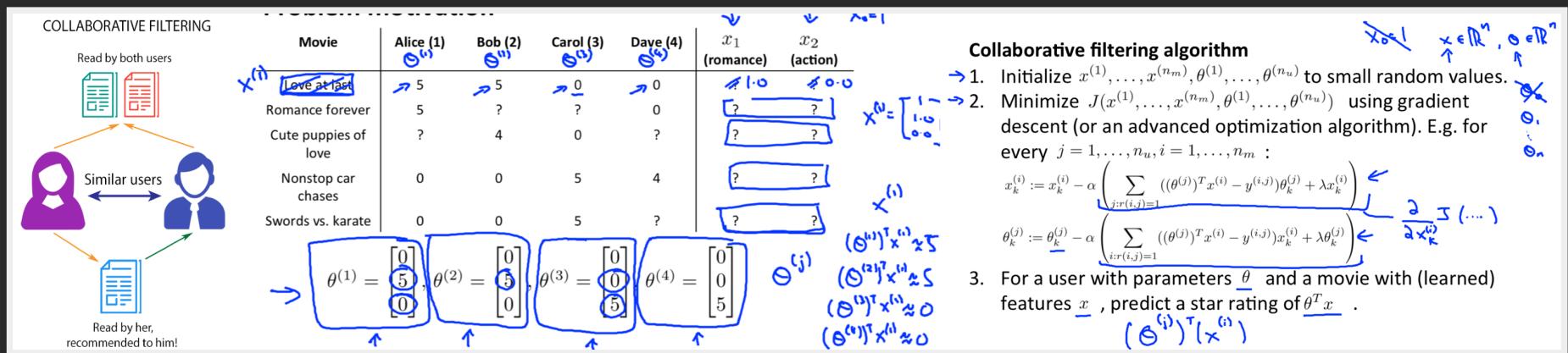
(for $k = 0$)

$$\theta_k^{(j)} := \theta_k^{(j)} - \alpha \left(\sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} + \lambda \theta_k^{(j)} \right)$$

(for $k \neq 0$)

$\frac{\partial}{\partial \theta_k^{(j)}} J(\theta^{(1)}, \dots, \theta^{(n_u)})$

COLLABORATIVE FILTERING



QUESTIONS?

SOURCES

- <https://www.coursera.org/learn/machine-learning>
- https://en.wikipedia.org/wiki/Machine_learning