

Machine Learning

Supervised Learning

regression problem: continuous data

classification problem: discrete data

Unsupervised Learning

just give the data, no more.

Cocktail party problem algorithm

$$[W, s, v] = \underline{\text{svd}}((\text{repmat}(\text{sum}(x.*x, 1), \text{size}(x, 1), 1).*x)*x');$$

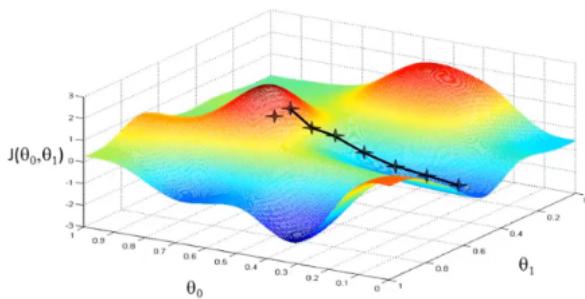
(奇异值分解)

Linear regression with one variable

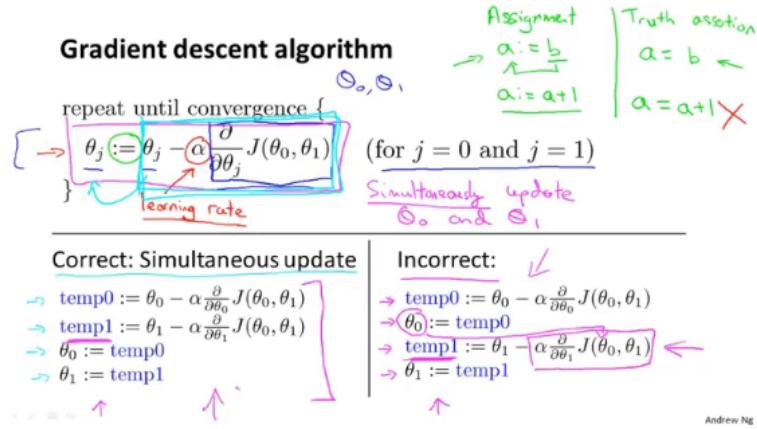
cost function

$$\begin{aligned} & \underset{\theta_0, \theta_1}{\text{minimize}} \quad \frac{1}{2m} \sum_{i=1}^m \left(h_{\theta}(x^{(i)}) - y^{(i)} \right)^2 \\ & \quad \text{#Training examples} \\ & \quad \text{---} \\ & \quad h_{\theta}(x^{(i)}) = \theta_0 + \theta_1 x^{(i)} \\ \\ & \text{or our cost function:} \\ & J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m \left(h_{\theta}(x^{(i)}) - y^{(i)} \right)^2 \end{aligned}$$

Gradient descent



Andrew Ng

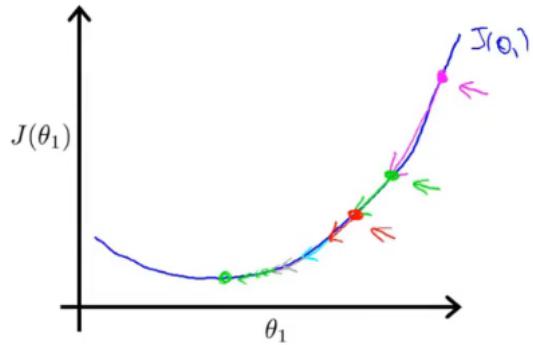


Andrew Ng

Gradient descent can converge to a local minimum, even with the learning rate α fixed.

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

As we approach a local minimum, gradient descent will automatically take smaller steps. So, no need to decrease α over time.



Andrew Ng

"Batch" Gradient Decent

"Batch": Each step of gradient descent uses all the training example

Linear regression with multiple variables

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

Parameters: $\theta_0, \theta_1, \dots, \theta_n$ Θ $n+1$ -dimensional vector

Cost function:

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

Gradient descent:

Repeat {

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

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

Gradient Descent

Previously ($n=1$):

Repeat {

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

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

(simultaneously update θ_0, θ_1)

}

New algorithm ($n \geq 1$):

Repeat {

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

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

}

Andrew Ng

Mean normalization

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

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

$x_2 = \frac{\#\text{bedrooms}-2}{5}$ 1-5 bedrooms

$-0.5 \leq x_1 \leq 0.5$ $-0.5 \leq x_2 \leq 0.5$

$x_1 \leftarrow \frac{x_1 - \mu_1}{\sigma_1}$ avg value of x_1 in training set
range $(\max - \min)$ (or standard deviation)

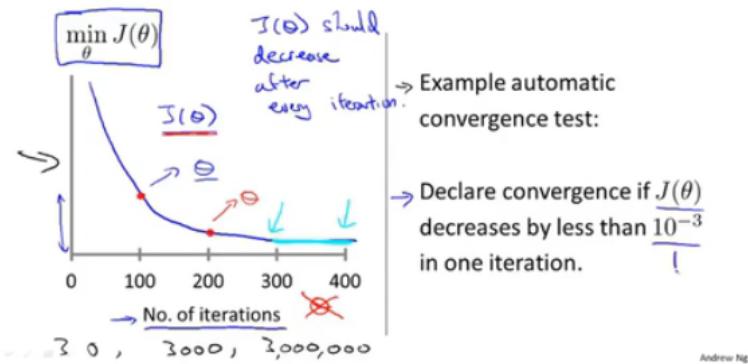
$x_2 \leftarrow \frac{x_2 - \mu_2}{\sigma_2}$

Andrew Ng

Learning rate

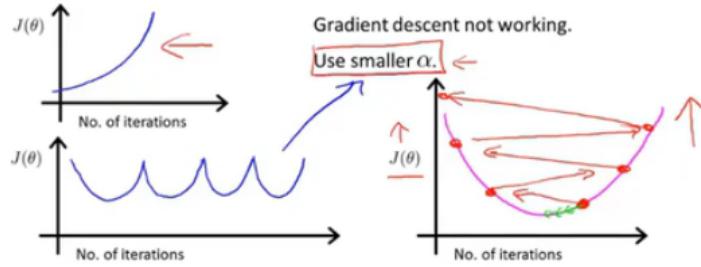
correct

Making sure gradient descent is working correctly.



incorrect

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.

Andrew Ng

Summary

- If α is too small: slow convergence. $\frac{+}{\text{#iters}}$
- If α is too large: $J(\theta)$ may not decrease on every iteration; may not converge. (Slow converge also possible.)

To choose α , try

$$\dots, \underbrace{0.001}_{\approx 2x}, \underbrace{0.003}_{\approx 2x}, \underbrace{0.01}_{\approx 2x}, \underbrace{0.03}_{\approx 2x}, \underbrace{0.1}_{\approx 2x}, \underbrace{0.3}_{\approx 2x}, \underbrace{1}_{\approx 2x}, \dots$$

Andrew Ng

Normal equation

Method to solve for θ analytically

Examples: $m = 4$.

x_0	Size (feet ²)	Number of bedrooms	Number of floors	Age of home (years)	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

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

Andrew Ng

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

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

Set $A = X^T X$

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

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

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

$$\Theta = \text{pinv}(X^T X)^{-1} X^T y$$

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

$$\begin{array}{l|l} X' & X^T \\ \hline \end{array}$$

Feature scaling

$0 \leq x_1 \leq 1$
 $0 \leq x_2 \leq 1000$
 $0 \leq x_3 \leq 10^{-5}$

no need to do feature scaling

Gradient Descent vs Normal Equation

m training examples, n features.

Gradient Descent

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

$$n = 10^6$$

Normal Equation

- No need to choose α .
- Don't need to iterate.
- Need to compute $(X^T X)^{-1}$ $n \times n$ $O(n^3)$
- Slow if n is very large.

$$\begin{array}{l} n = 100 \\ n = 1000 \\ \dots n = 10000 \end{array}$$

Andrew Ng

Non-invertible

What if $X^T X$ is non-invertible?

- Redundant features (linearly dependent).

E.g. $\begin{cases} x_1 = \text{size in feet}^2 \\ x_2 = \text{size in m}^2 \end{cases}$ $1m = 3.28 \text{ feet}$

$$x_1 = (3.28)^2 x_2$$

$$\begin{array}{l} \rightarrow n = 10 \\ \rightarrow h = 100 \\ \Theta \in \mathbb{R}^{101} \end{array}$$

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

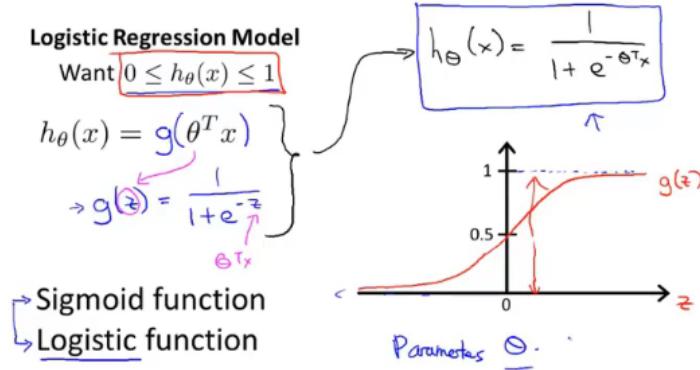
- Delete some features, or use regularization.

↓ later

Andrew Ng

Logistic Regression

Hypothesis Representation



Andrew Ng

Interpretation of Hypothesis Output

$h_\theta(x)$ = estimated probability that $y = 1$ on input x

Example: If $x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumorSize} \end{bmatrix}$

$$h_\theta(x) = 0.7 \quad y=1$$

Tell patient that 70% chance of tumor being malignant

$$h_\theta(x) = P(y=1|x; \theta) \quad \text{"probability that } y=1 \text{, given } x, \text{ parameterized by } \theta"$$

$y = 0 \text{ or } 1$

$$\Rightarrow P(y=0|x; \theta) + P(y=1|x; \theta) = 1$$

$$P(y=0|x; \theta) = 1 - P(y=1|x; \theta)$$

Andrew Ng

Decision boundary

Logistic regression

$$\rightarrow h_\theta(x) = g(\theta^T x) = P(y=1|x; \theta)$$

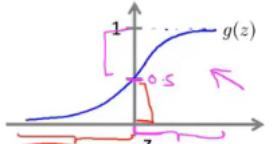
$$\rightarrow g(z) = \frac{1}{1+e^{-z}}$$

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

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

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

$$h_\theta(x) = g(\theta^T x) \quad \rightarrow \theta^T x < 0$$



$$g(z) \geq 0.5$$

when $z \geq 0$

$$h_\theta(x) = g(\theta^T x) \geq 0.5$$

when $\theta^T x \geq 0$

Andrew Ng

Cost function

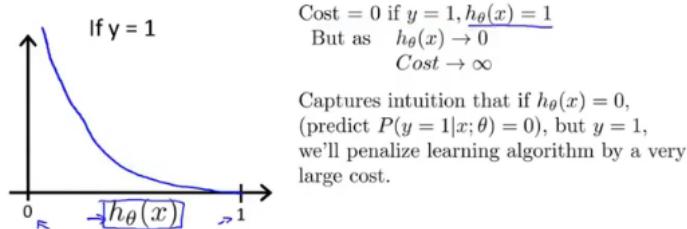
Cost function

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

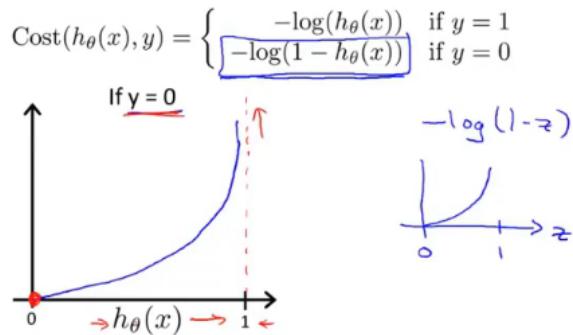
$$\rightarrow \text{Cost}(h_\theta(x), y) = \frac{1}{2} (h_\theta(x) - y)^2$$

Logistic regression cost function

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



Logistic regression cost function



Andrew Ng

Simplified cost function

Logistic regression cost function

$$\begin{aligned} \rightarrow J(\theta) &= \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_\theta(x^{(i)}), y^{(i)}) \\ \rightarrow \text{Cost}(h_\theta(x), y) &= \begin{cases} -\log(h_\theta(x)) & \text{if } y = 1 \\ -\log(1 - h_\theta(x)) & \text{if } y = 0 \end{cases} \\ \text{Note: } y &= 0 \text{ or } 1 \text{ always} \\ \rightarrow \text{Cost}(h_\theta(x), y) &= -y \log(h_\theta(x)) - (1-y) \log(1 - h_\theta(x)) \\ \text{If } y=1: \text{Cost}(h_\theta(x), y) &= -\log(h_\theta(x)) \\ \text{If } y=0: \text{Cost}(h_\theta(x), y) &= -\log(1 - h_\theta(x)) \end{aligned}$$

Andrew Ng

Gradient descent

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

$$\left[\begin{array}{l} \text{Repeat } \{ \\ \quad \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \\ \quad \} \end{array} \right] \quad \text{(simultaneously update all } \theta_j \text{)}$$

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

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

Andrew Ng

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

$$\rightarrow \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix} \quad \theta_0 = S, \theta_1 = S, \theta_2 = S.$$
$$\rightarrow J(\theta) = (\theta_0 - 5)^2 + (\theta_1 - 5)^2 + (\theta_2 - 5)^2$$
$$\rightarrow \frac{\partial}{\partial \theta_0} J(\theta) = 2(\theta_0 - 5)$$
$$\rightarrow \frac{\partial}{\partial \theta_1} J(\theta) = 2(\theta_1 - 5)$$
$$\rightarrow \frac{\partial}{\partial \theta_2} J(\theta) = 2(\theta_2 - 5)$$

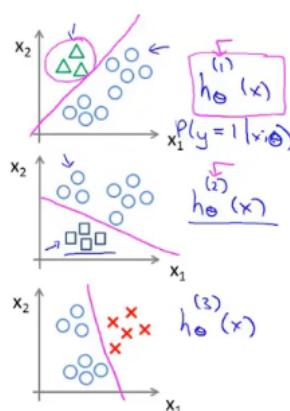
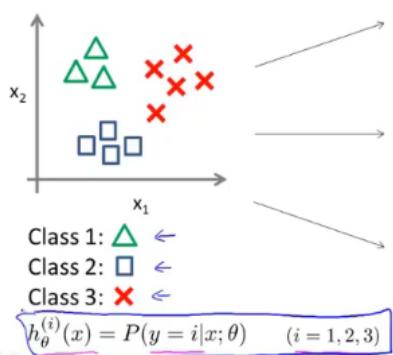
```
options = optimset('GradObj', 'on', 'MaxIter', '100');
initialTheta = zeros(3,1);
[optTheta, functionVal, exitFlag] ...
= fminunc(@costFunction, initialTheta, options);
```

```
function [jVal, gradient] = costFunction(theta)
jVal = (theta(1)-5)^2 + ...  
       (theta(2)-5)^2;  
gradient = zeros(3,1);  
gradient(1) = 2*(theta(1)-5);  
gradient(2) = 2*(theta(2)-5);
```

```
theta = [theta0; theta1; ...; theta_n];
function [jVal, gradient] = costFunction(theta)
jVal = [code to compute J(theta)];
gradient(1) = [code to compute  $\frac{\partial}{\partial \theta_0} J(\theta)$ ];
gradient(2) = [code to compute  $\frac{\partial}{\partial \theta_1} J(\theta)$ ];
...
gradient(n+1) = [code to compute  $\frac{\partial}{\partial \theta_n} J(\theta)$ ];
```

Multi-class classification: One-vs-all

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



Andrew Ng

One-vs-all

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

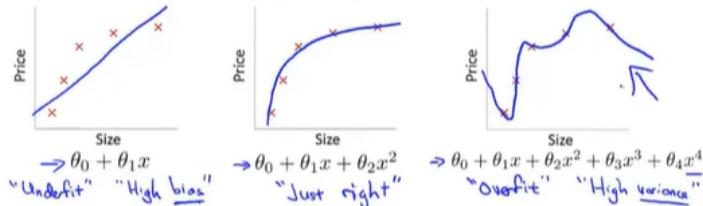
$$\max_i \underline{h_{\theta}^{(i)}(x)}$$

Andrew Ng

Regularization

The problem of overfitting

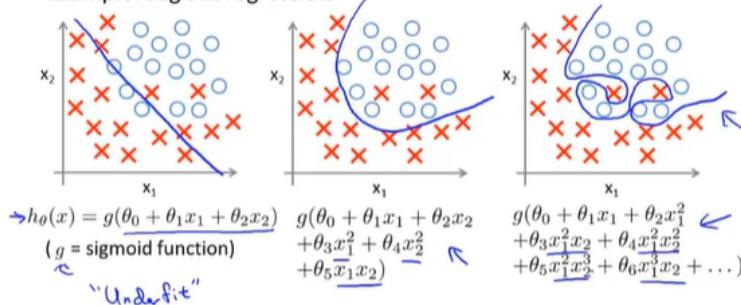
Example: Linear regression (housing prices)



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).

Andrew Ng

Example: Logistic regression



Andrew Ng

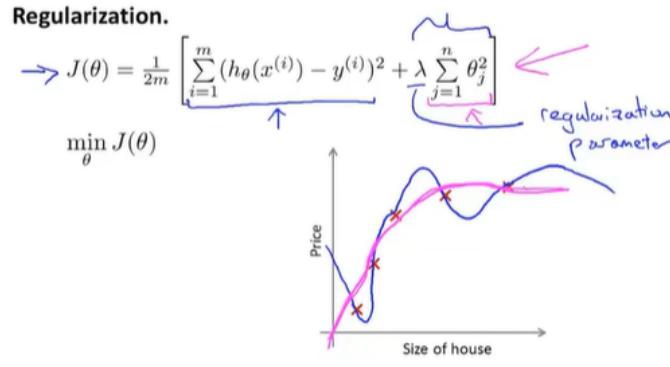
Addressing overfitting:

Options:

1. Reduce number of features.
 - Manually select which features to keep.
 - Model selection algorithm (later in course).
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 .

Andrew Ng

Cost function

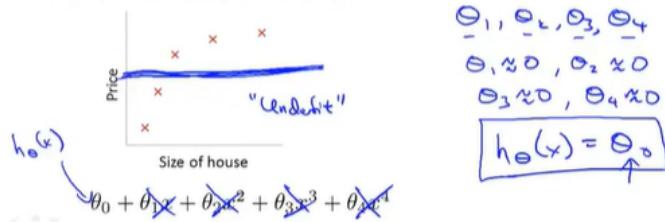


Andrew Ng

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 $\lambda = 10^{10}$)?



Andrew Ng

Regularized linear regression

Gradient descent

Repeat {

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

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

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

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

$\rightarrow J(\theta)$

$\theta_j \times 0.99$

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

Andrew Ng

Non-invertibility (optional/advanced).

Suppose $m \leq n$, \leftarrow

(#examples) (#features)

$$\theta = \underbrace{(X^T X)^{-1} X^T y}_{\text{non-invertible / singular}} \quad \frac{\text{piv}}{\text{inv}} \quad \frac{\text{inv}}{\text{R}}$$

If $\lambda > 0$,

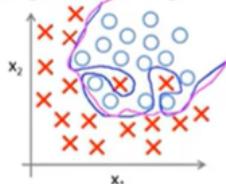
$$\theta = \left(X^T X + \lambda \begin{bmatrix} 0 & 1 & 1 & & \\ & 1 & & & \\ & & 1 & & \\ & & & \ddots & \\ & & & & 1 \end{bmatrix} \right)^{-1} X^T y$$

invertible.

Andrew Ng

Regularized logistic regression

Regularized logistic regression.



$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2^2 + \theta_3 x_1^2 x_2 + \theta_4 x_1^2 x_2^2 + \theta_5 x_1^2 x_2^3 + \dots)$$

Cost function:

$$\rightarrow J(\theta) = - \left[\frac{1}{m} \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\theta}(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2 \quad \boxed{\theta_0, \theta_1, \dots, \theta_n}$$

Andrew Ng

Gradient descent

Repeat {

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

$$\rightarrow \theta_j := \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} \theta_j \right] \quad \begin{array}{l} (j = \textcolor{red}{1}, 2, 3, \dots, n) \\ \theta_0, \dots, \theta_n \end{array}$$

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

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

Andrew Ng

Advanced optimization

```

f minune (Q costfunction)
    theta(1) ←
    theta(i) ←
    theta(n+1) ←
    theta(i+1) ←

function [jVal, gradient] = costFunction(theta)
    jVal = [ code to compute J(theta) ];
    J(theta) = [-\frac{1}{m} \sum_{i=1}^m y^{(i)} \log(h_{\theta}(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)}))] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2
    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)}) + \frac{\lambda}{m} \theta_1 ←
    gradient(3) = [ code to compute \frac{\partial}{\partial \theta_2} J(theta) ];
    : (\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_2^{(i)}) + \frac{\lambda}{m} \theta_2 ←
    gradient(n+1) = [ code to compute \frac{\partial}{\partial \theta_n} J(theta) ];

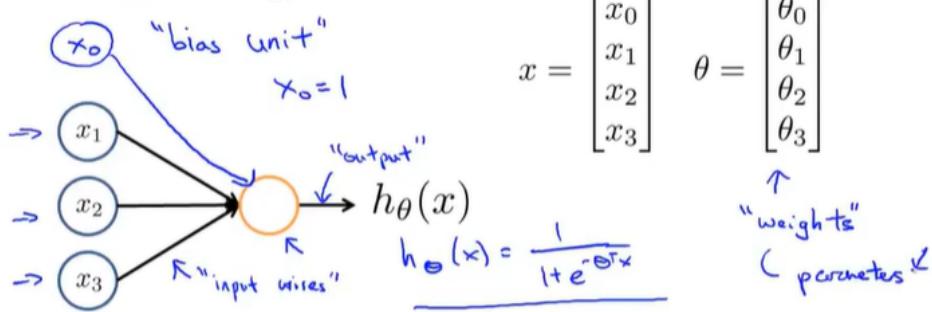
```

Andrew Ng

Neural Networks

Model representation

Neuron model: Logistic unit

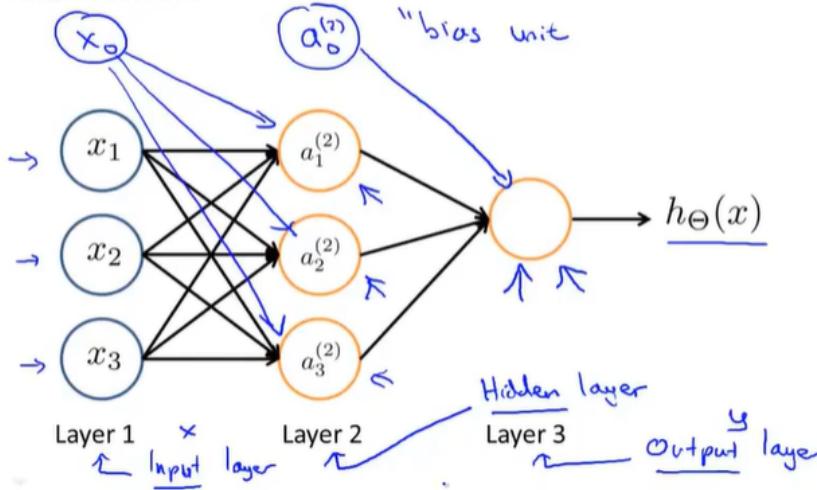


Sigmoid (logistic) activation function.

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

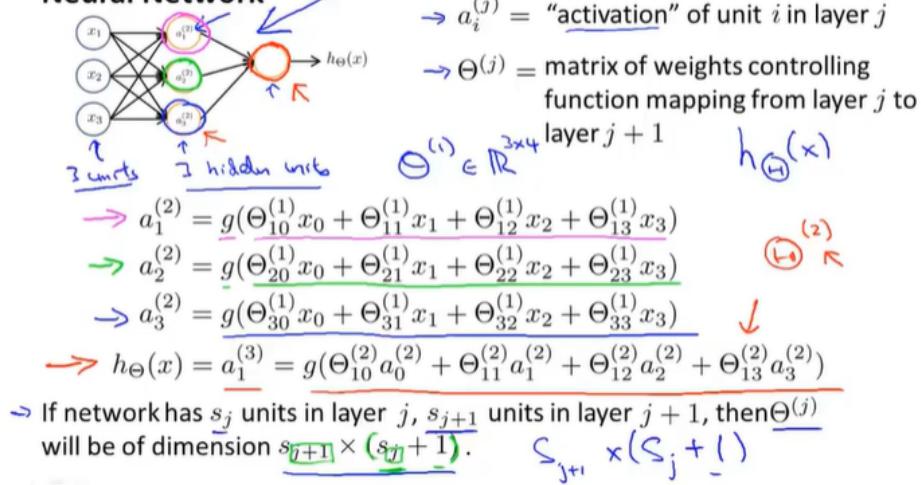
Andrew Ng

Neural Network



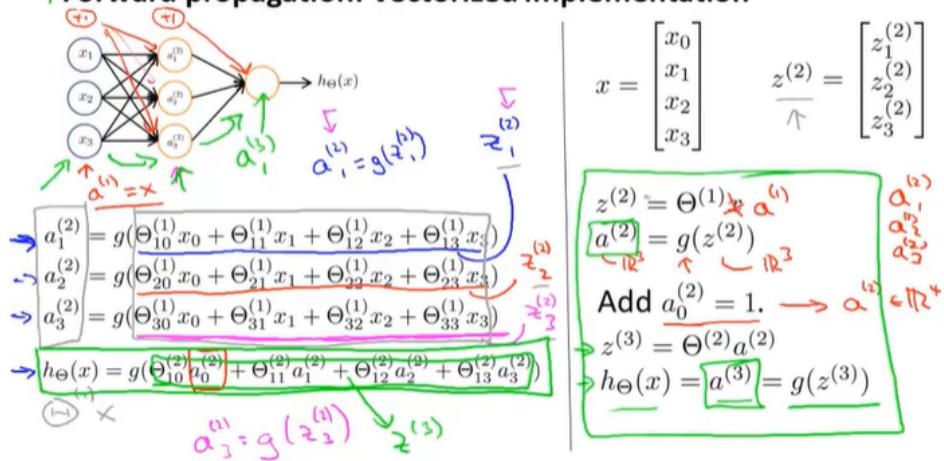
Andrew Ng

Neural Network



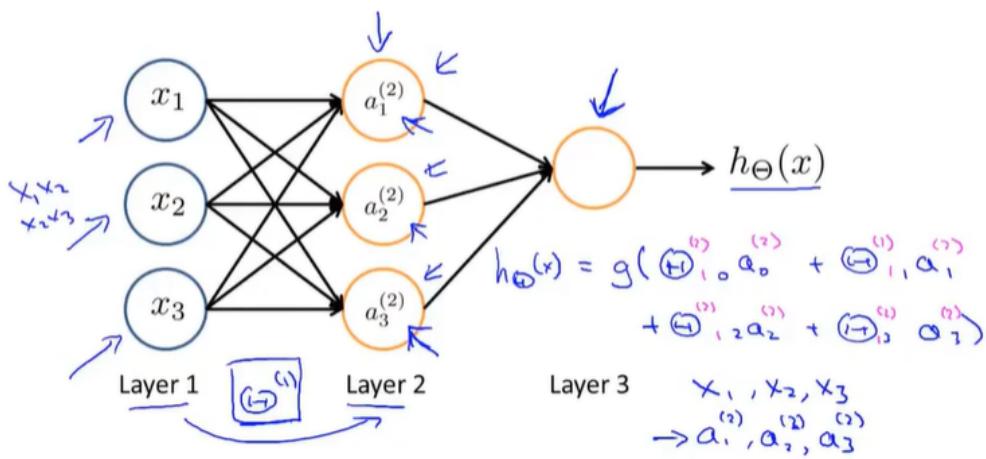
Andrew Ng

Forward propagation: Vectorized implementation



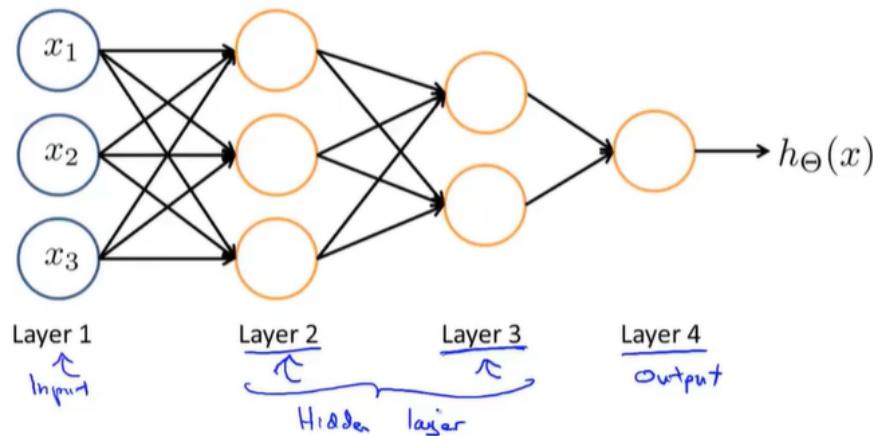
Andrew Ng

Neural Network learning its own features



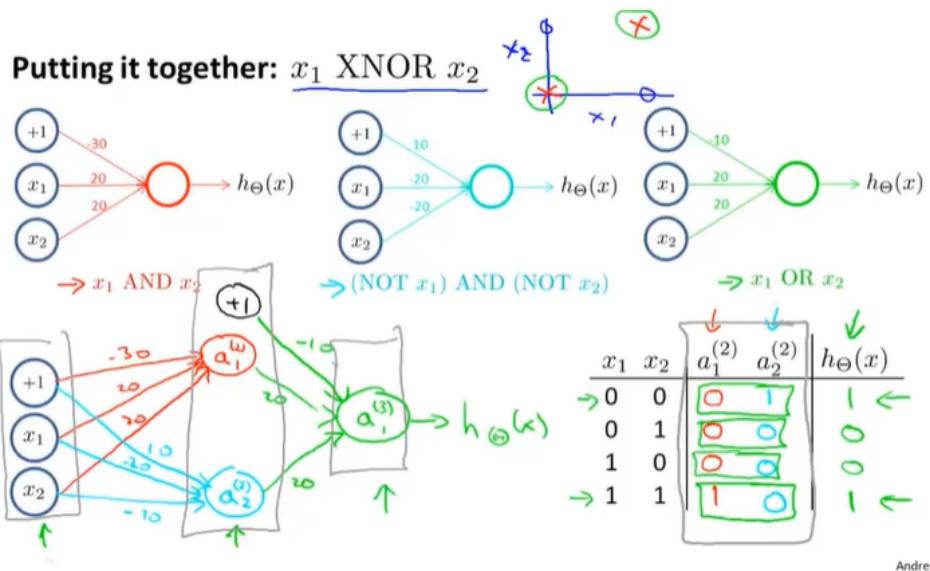
Andrew Ng

Other network architectures



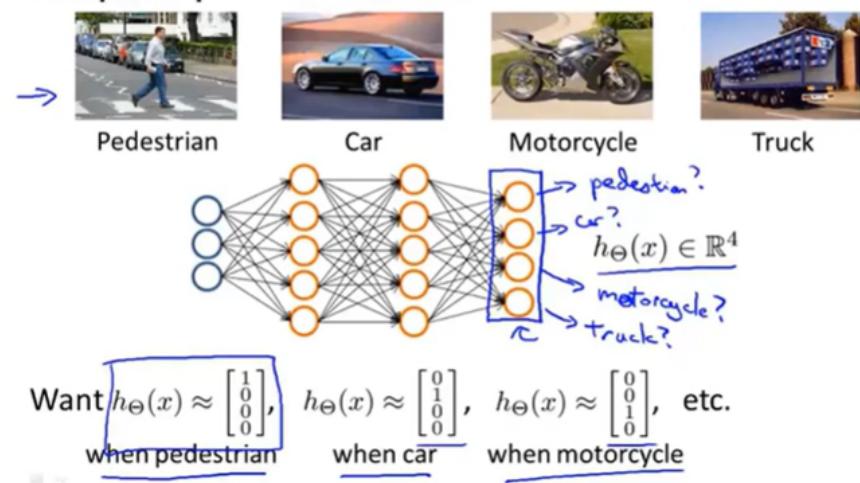
Andrew Ng

Examples and intuitions

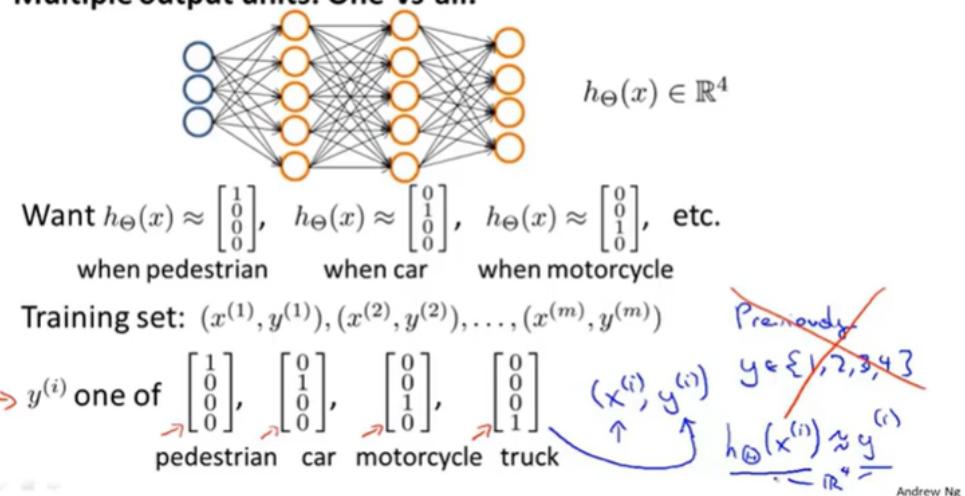


Multi-class classification

Multiple output units: One-vs-all.



Multiple output units: One-vs-all.



Cost function

Cost function

Logistic regression:

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

Neural network:

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

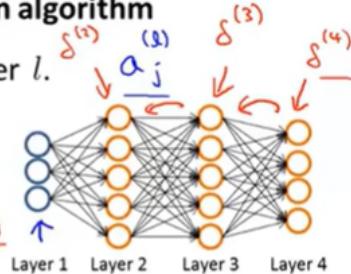
$$\quad \quad \quad \boxed{\frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\Theta_{ji}^{(l)})^2}$$

Andrew Ng

Backpropagation algorithm

Gradient computation: Backpropagation algorithm

Intuition: $\delta_j^{(l)}$ = "error" of node j in layer l .



For each output unit (layer $L = 4$)

$$\delta_j^{(4)} = \underline{a_j^{(4)} - y_j} = (h_\Theta(x))_j - \underline{y_j}$$

$$\rightarrow \delta^{(3)} = (\Theta^{(3)})^T \delta^{(4)} \cdot g'(z^{(3)})$$

$$\rightarrow \delta^{(2)} = (\Theta^{(2)})^T \delta^{(3)} \cdot g'(z^{(2)})$$

(No $\delta^{(1)}$) $\frac{\partial}{\partial \Theta_{ij}^{(2)}} J(\Theta) = \underline{a_j^{(2)} \delta_i^{(3)}} \quad$ (ignoring λ ; if $\lambda = 0$)

Backpropagation algorithm

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

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

For $i = 1$ to $m \leftarrow (x^{(i)}, y^{(i)})$.

- Set $a^{(1)} = x^{(i)}$
 - Perform forward propagation to compute $a^{(l)}$ for $l = 2, 3, \dots, L$
 - Using $y^{(i)}$, compute $\delta^{(L)} = \underline{a^{(L)} - y^{(i)}}$
 - Compute $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$ ~~$\delta^{(1)}$~~
 - $\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$
 - $D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \Theta_{ij}^{(l)}$ if $j \neq 0$
 - $D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)}$ if $j = 0$
- $\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = D_{ij}^{(l)}$

Backpropagation intuition

What is backpropagation doing?

$$J(\Theta) = -\frac{1}{m} \left[\sum_{i=1}^m \underbrace{y^{(i)} \log(h_\Theta(x^{(i)})) + (1-y^{(i)}) \log(1-h_\Theta(x^{(i)}))}_{\text{Cost for one example}} \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\Theta_{ji}^{(l)})^2$$

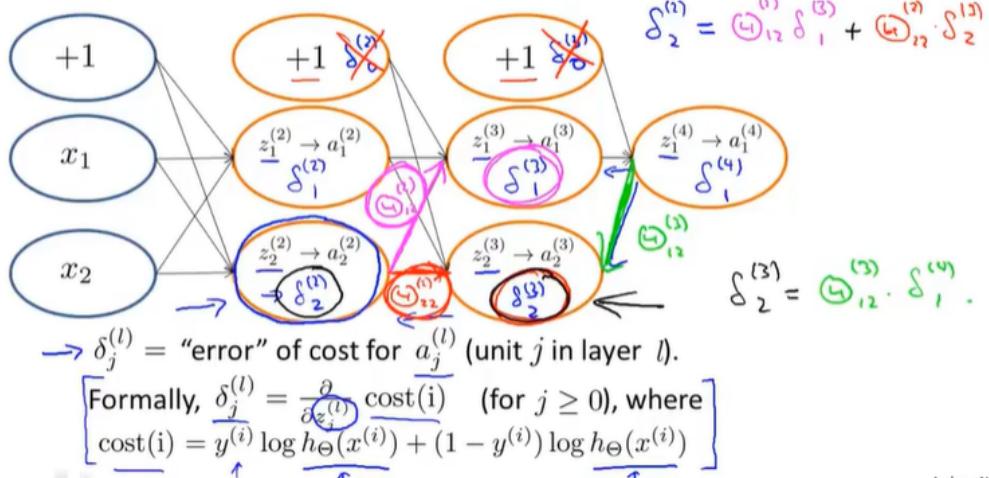
Focusing on a single example $x^{(i)}, y^{(i)}$, the case of 1 output unit, and ignoring regularization ($\lambda = 0$),

$$\text{cost}(i) = y^{(i)} \log h_\Theta(x^{(i)}) + (1-y^{(i)}) \log h_\Theta(x^{(i)})$$

(Think of $\text{cost}(i) \approx (h_\Theta(x^{(i)}) - y^{(i)})^2$)

I.e. how well is the network doing on example i? $y^{(i)}$

Forward Propagation



Andrew Ng

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)

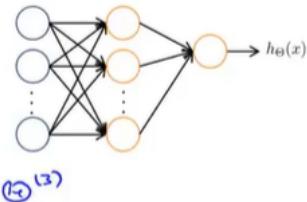
"Unroll" into vectors

Example

```

 $s_1 = 10, s_2 = 10, s_3 = 1$ 
 $\Theta^{(1)} \in \mathbb{R}^{10 \times 11}, \Theta^{(2)} \in \mathbb{R}^{10 \times 11}, \Theta^{(3)} \in \mathbb{R}^{1 \times 11}$ 
 $D^{(1)} \in \mathbb{R}^{10 \times 11}, D^{(2)} \in \mathbb{R}^{10 \times 11}, D^{(3)} \in \mathbb{R}^{1 \times 11}$ 
 $\rightarrow \Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$ 
 $\rightarrow \text{thetaVec} = [\Theta^{(1)}; \Theta^{(2)}; \Theta^{(3)}];$ 
 $\rightarrow \text{DVec} = [D^{(1)}; D^{(2)}; D^{(3)}];$ 
 $\Theta^{(1)} = \text{reshape}(\text{thetaVec}(1:110), 10, 11);$ 
 $\Theta^{(2)} = \text{reshape}(\text{thetaVec}(111:220), 10, 11);$ 
 $\Theta^{(3)} = \text{reshape}(\text{thetaVec}(221:231), 1, 11);$ 

```



Learning Algorithm

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

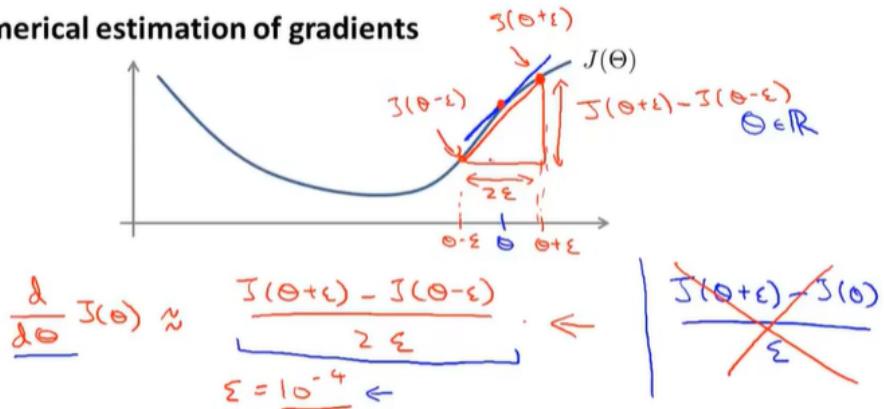
```

function [jval, gradientVec] = costFunction(thetaVec)
     $\rightarrow$  From thetaVec, get  $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$  reshape
     $\rightarrow$  Use forward prop/back prop to compute  $D^{(1)}, D^{(2)}, D^{(3)}$  and  $J(\Theta)$ 
    Unroll  $D^{(1)}, D^{(2)}, D^{(3)}$  to get gradientVec.

```

Gradient checking

Numerical estimation of gradients



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

Andrew Ng

Parameter vector θ

- $\rightarrow \theta \in \mathbb{R}^n$ (E.g. θ is "unrolled" version of $\underline{\Theta^{(1)}}, \underline{\Theta^{(2)}}, \underline{\Theta^{(3)}}$)
- $\rightarrow \theta = [\theta_1, \theta_2, \theta_3, \dots, \theta_n]$
- $\rightarrow \frac{\partial}{\partial \theta_1} J(\theta) \approx \frac{J(\theta_1 + \epsilon, \theta_2, \theta_3, \dots, \theta_n) - J(\theta_1 - \epsilon, \theta_2, \theta_3, \dots, \theta_n)}{2\epsilon}$
- $\rightarrow \frac{\partial}{\partial \theta_2} J(\theta) \approx \frac{J(\theta_1, \theta_2 + \epsilon, \theta_3, \dots, \theta_n) - J(\theta_1, \theta_2 - \epsilon, \theta_3, \dots, \theta_n)}{2\epsilon}$
- \vdots
- $\rightarrow \frac{\partial}{\partial \theta_n} J(\theta) \approx \frac{J(\theta_1, \theta_2, \theta_3, \dots, \theta_n + \epsilon) - J(\theta_1, \theta_2, \theta_3, \dots, \theta_n - \epsilon)}{2\epsilon}$

Andrew Ng

```

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;

```

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

Check that $\text{gradApprox} \approx \text{DVec}$ ←
 From back prop.

Andrew Ng

Implementation Note:

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

Important:

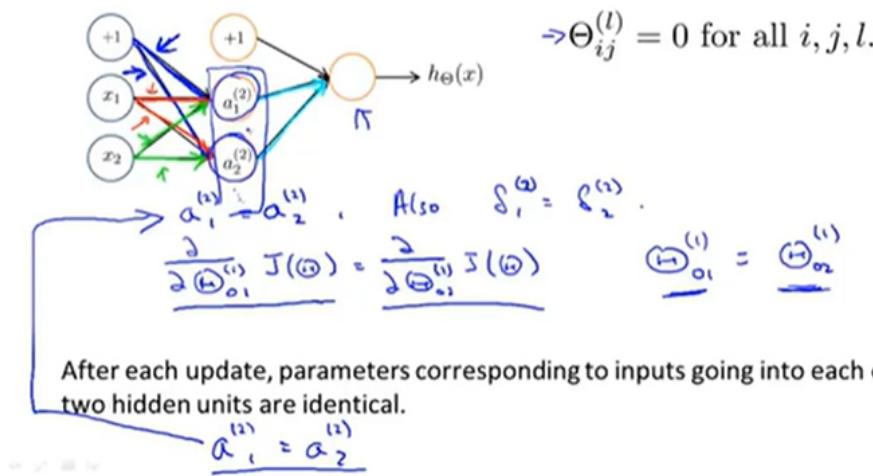
$\begin{array}{c} \xrightarrow{\text{DVec}} \\ g^{(1)}, g^{(2)}, g^{(3)} \end{array}$

- \rightarrow - 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.

Andrew Ng

Random initialization

Zero initialization



Andrew Ng

Random initialization: Symmetry breaking

- Initialize each $\Theta_{ij}^{(l)}$ to a random value in $[-\epsilon, \epsilon]$
 - (i.e. $-\epsilon \leq \Theta_{ij}^{(l)} \leq \epsilon$)
- E.g. Random 10×11 matrix (betw. 0 and 1)
- ```

→ Theta1 = rand(10,11)*(2*INIT_EPSILON)
- INIT_EPSILON; [-\epsilon, \epsilon]

```
- ```

→ Theta2 = rand(1,11)*(2*INIT_EPSILON)
- INIT_EPSILON;

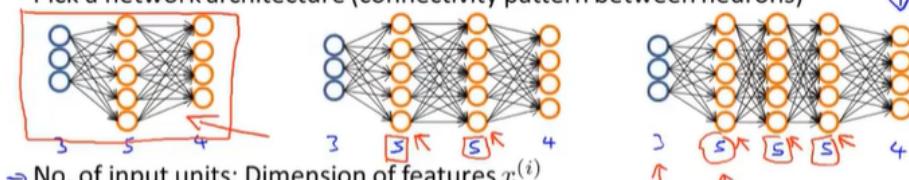
```

Andrew Ng

Putting it together

Training a neural network

Pick a network architecture (connectivity pattern between neurons)



→ No. of input units: Dimension of features $x^{(i)}$

→ No. output units: Number of classes

Reasonable default: 1 hidden layer, or if >1 hidden layer, have same no. of hidden units in every layer (usually the more the better)

↑

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

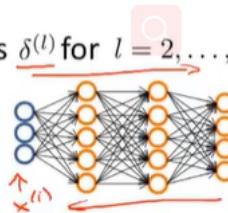
~~y < 5~~

$$y = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \text{ or } \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \leftarrow$$

Andrew Ng

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 $l = 2, \dots, L$.
- $\Delta^{(l)} := \Delta^{(l)} + \delta^{(l)} (a^{(l)})^T$
- ...
- compute $\frac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta)$.



Andrew Ng

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

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

$J(\Theta)$ — non-convex.

Andrew Ng

Advice for applying machine learning

Deciding what to try next

Debugging a learning algorithm:

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

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

However, when you test your hypothesis on a new set of houses, you find that it makes unacceptably large errors in its predictions. What should you try next?

- - Get more training examples
- Try smaller sets of features $x_1, x_2, x_3, \dots, x_{100}$
- - Try getting additional features
- Try adding polynomial features ($x_1^2, x_2^2, x_1 x_2$, etc.)
- Try decreasing λ
- Try increasing λ

Andrew Ng

Machine learning diagnostic:

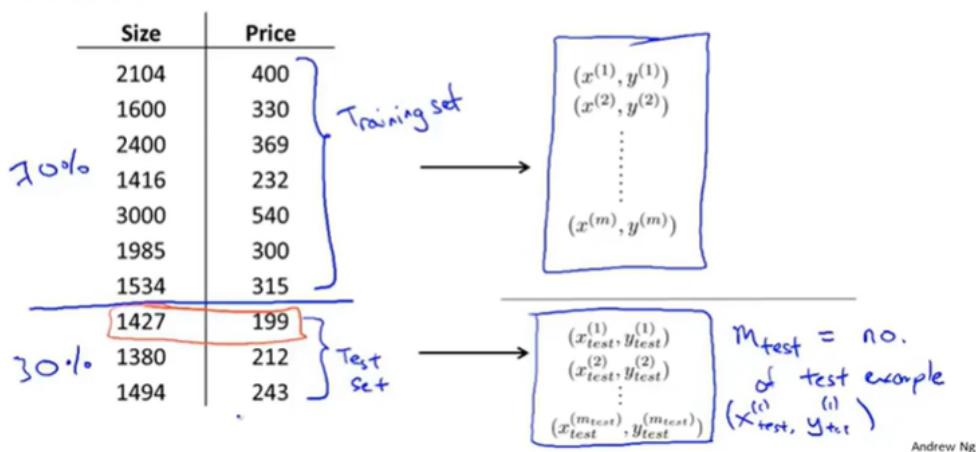
Diagnostic: A test that you can run to gain insight what is/isn't working with a learning algorithm, and gain guidance as to how best to improve its performance.

Diagnostics can take time to implement, but doing so can be a very good use of your time.

Evaluating a hypothesis

Evaluating your hypothesis

Dataset:



Training/testing procedure for linear regression

- - Learn parameter θ from training data (minimizing training error $J(\theta)$)

- Compute test set error:

$$J_{\text{test}}(\theta) = \frac{1}{2m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} (h_{\theta}(x_{\text{test}}^{(i)}) - y_{\text{test}}^{(i)})^2$$

Training/testing procedure for logistic regression

- - Learn parameter θ from training data
- Compute test set error: $J_{\text{test}}(\theta) = -\frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} y_{\text{test}}^{(i)} \log h_{\theta}(x_{\text{test}}^{(i)}) + (1 - y_{\text{test}}^{(i)}) \log h_{\theta}(x_{\text{test}}^{(i)})$
- Misclassification error (0/1 misclassification error):

$$\text{err}(h_{\theta}(x), y) = \begin{cases} 1 & \text{if } h_{\theta}(x) \geq 0.5, y = 0 \\ 0 & \text{otherwise or if } h_{\theta}(x) < 0.5, y = 1 \end{cases} \text{ error}$$

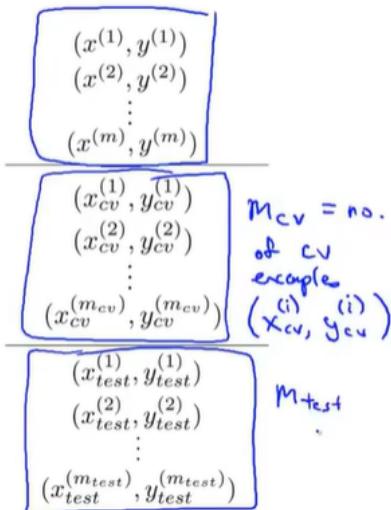
$$\text{Test error} = \frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} \text{err}(h_{\theta}(x_{\text{test}}^{(i)}), y_{\text{test}}^{(i)}).$$

Model selection and training/validation/test sets

Evaluating your hypothesis

Dataset:

Size	Price	
2104	400	60% Training set
1600	330	
2400	369	
1416	232	
3000	540	
1985	300	
1534	315	Cross validation set (cv)
1427	199	
1380	212	
1494	243	
		20% test set



Andrew Ng

Train/validation/test error

Training error:

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

Cross Validation error:

$$\rightarrow J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_{\theta}(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$$

Test error:

$$\rightarrow J_{test}(\theta) = \frac{1}{2m_{test}} \sum_{i=1}^{m_{test}} (h_{\theta}(x_{test}^{(i)}) - y_{test}^{(i)})^2$$

Andrew Ng

Model selection

- δ-1 1. $h_{\theta}(x) = \theta_0 + \theta_1 x \rightarrow \min_{\theta} \mathcal{J}(\theta) \rightarrow \theta^{(1)} \rightarrow \mathcal{J}_{cv}(\theta^{(1)})$
 - δ-2 2. $h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 \rightarrow \theta^{(2)} \rightarrow \mathcal{J}_{cv}(\theta^{(2)})$
 - δ-3 3. $h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_3 x^3 \rightarrow \theta^{(3)} \rightarrow \mathcal{J}_{cv}(\theta^{(3)})$
 - ⋮
 - δ-10 10. $h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_{10} x^{10} \rightarrow \theta^{(10)} \rightarrow \mathcal{J}_{cv}(\theta^{(10)})$
- $\delta = 4$ ↑

Pick $\theta_0 + \theta_1 x_1 + \dots + \theta_4 x^4$ ←

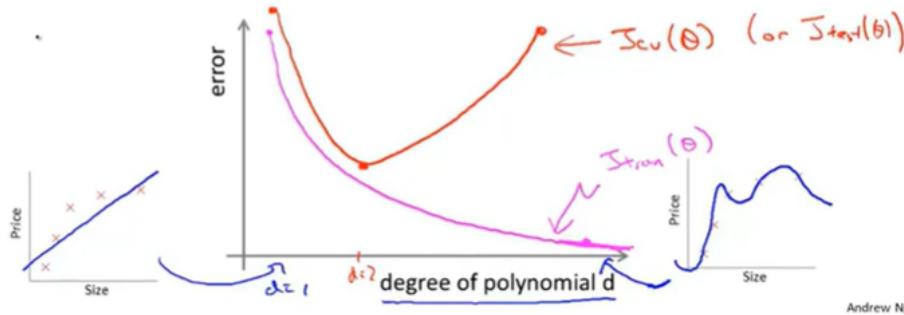
Estimate generalization error for test set $J_{test}(\theta^{(4)})$ ←

Diagnosing bias vs. variance

Bias/variance

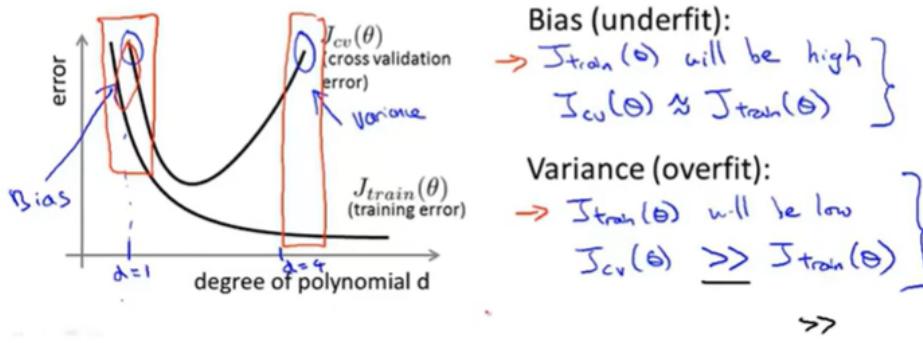
Training error: $J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$

Cross validation error: $J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_\theta(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$ (or $J_{test}(\theta)$)



Diagnosing bias vs. variance

Suppose your learning algorithm is performing less well than you were hoping. ($J_{cv}(\theta)$ or $J_{test}(\theta)$ is high.) Is it a bias problem or a variance problem?



Regularization and bias/variance

Choosing the regularization parameter λ

$$h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4 \quad \leftarrow$$

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^m \theta_j^2 \quad \leftarrow$$

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

$$J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_\theta(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$$

$$J_{test}(\theta) = \frac{1}{2m_{test}} \sum_{i=1}^{m_{test}} (h_\theta(x_{test}^{(i)}) - y_{test}^{(i)})^2 \quad \leftarrow J_{train}, J_{cv}, J_{test}$$

Choosing the regularization parameter λ

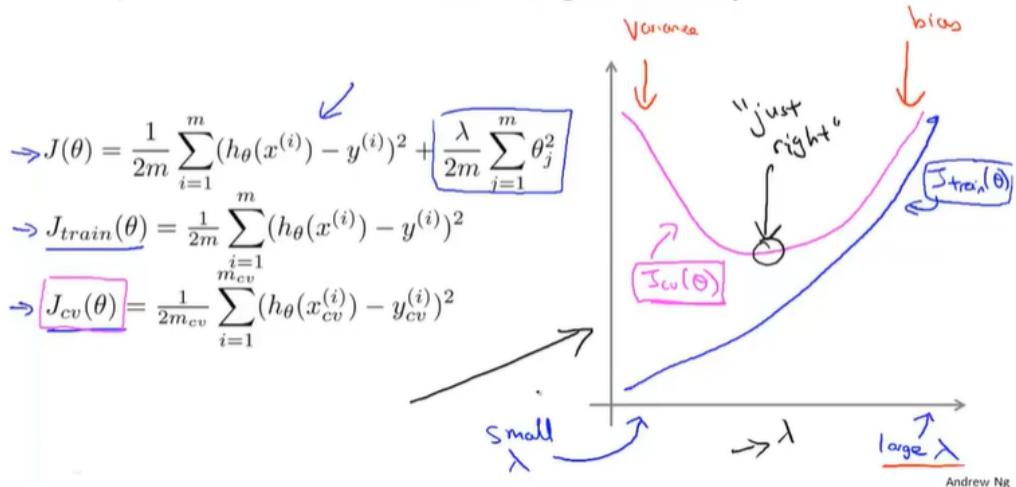
Model: $h_\theta(x) = \theta_0 + \theta_1x + \theta_2x^2 + \theta_3x^3 + \theta_4x^4$

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

- 1. Try $\lambda = 0$ $\rightarrow \min_{\theta} J(\theta) \rightarrow \theta^{(1)} \rightarrow J_{cv}(\theta^{(1)})$
 - 2. Try $\lambda = 0.01$ $\rightarrow \min_{\theta} J(\theta) \rightarrow \theta^{(2)} \rightarrow J_{cv}(\theta^{(2)})$
 - 3. Try $\lambda = 0.02$ $\rightarrow \theta^{(3)} \rightarrow J_{cv}(\theta^{(3)})$
 - 4. Try $\lambda = 0.04$ \vdots
 - 5. Try $\lambda = 0.08$ $\rightarrow \theta^{(5)} \rightarrow J_{cv}(\theta^{(5)})$
 - \vdots
 - 12. Try $\lambda = 10$ $\rightarrow \theta^{(12)} \rightarrow J_{cv}(\theta^{(12)})$
- Pick (say) $\theta^{(5)}$. Test error: $J_{test}(\theta^{(5)})$

Andrew Ng

Bias/variance as a function of the regularization parameter λ



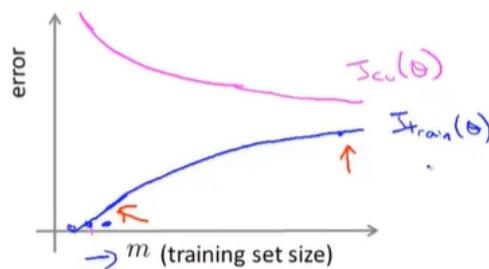
Andrew Ng

Learning curves

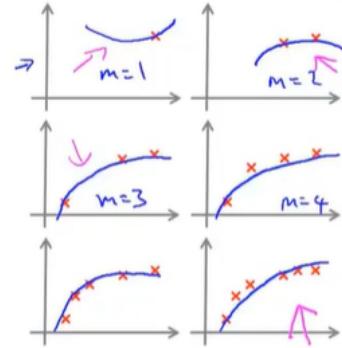
Learning curves

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

$$\Rightarrow J_{cv}(\theta) = \frac{1}{2m_{cv}} \sum_{i=1}^{m_{cv}} (h_\theta(x_{cv}^{(i)}) - y_{cv}^{(i)})^2$$

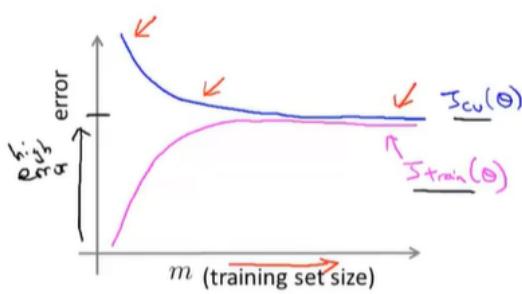


$$h_\theta(x) = \underline{\theta_0 + \theta_1x + \theta_2x^2}$$

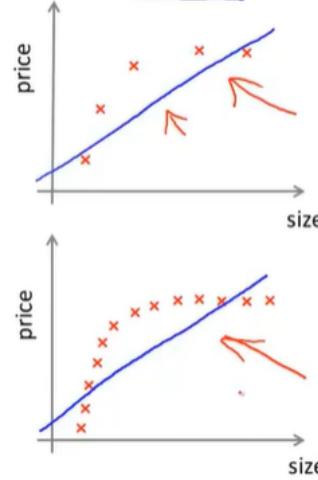


Andrew Ng

High bias

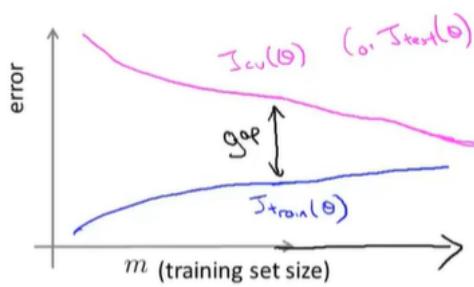


$$h_{\theta}(x) = \theta_0 + \theta_1 x$$



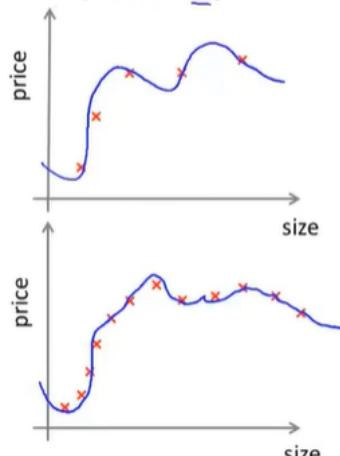
Andrew Ng

High variance



$$h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_{100} x^{100}$$

(and small λ)



Andrew Ng

Deciding what to try next(revisited)

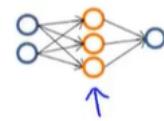
Debugging a learning algorithm:

Suppose you have implemented regularized linear regression to predict housing prices. However, when you test your hypothesis in a new set of houses, you find that it makes unacceptably large errors in its prediction. What should you try next?

- Get more training examples → fixes high variance
- Try smaller sets of features → fixes high variance
- Try getting additional features → fixes high bias
- Try adding polynomial features ($x_1^2, x_2^2, x_1 x_2$, etc) → fixes high bias.
- Try decreasing λ → fixes high bias
- Try increasing λ → fixes high variance

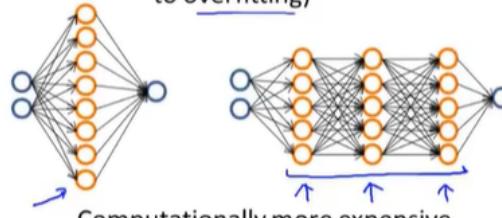
Neural networks and overfitting

- “Small” neural network
(fewer parameters; more prone to underfitting)



Computationally cheaper

- “Large” neural network
(more parameters; more prone to overfitting)



Computationally more expensive.

Use regularization (λ) to address overfitting.

$$J_{\text{reg}}(\theta)$$

Andrew Ng

Machine learning system design

Error analysis

Recommended approach

- - Start with a simple algorithm that you can implement quickly.
Implement it and test it on your cross-validation data.
- - Plot learning curves to decide if more data, more features, etc. are likely to help.
- - Error analysis: Manually examine the examples (in cross validation set) that your algorithm made errors on. See if you spot any systematic trend in what type of examples it is making errors on.

Error Analysis

$m_{CV} = 500$ examples in cross validation set

Algorithm misclassifies 100 emails.

Manually examine the 100 errors, and categorize them based on:

- (i) What type of email it is *pharma, replica, steal passwords, ...*
- (ii) What cues (features) you think would have helped the algorithm classify them correctly.

Pharma: 12

Replica/fake: 4

→ Steal passwords: 53

Other: 31

→ Deliberate misspellings: 5

(m0rgage, med1cine, etc.)

→ Unusual email routing: 16

→ Unusual (spamming) punctuation: 32

← ← ← ← ←

Andrew Ng

The importance of numerical evaluation

Should discount/discounts/discounted/discounting be treated as the same word?

Can use "stemming" software (E.g. "Porter stemmer")
universe/university.

Error analysis may not be helpful for deciding if this is likely to improve performance. Only solution is to try it and see if it works.

Need numerical evaluation (e.g., cross validation error) of algorithm's performance with and without stemming.

Without stemming: 5% error With stemming: 3% error

Distinguish upper vs. lower case (Mom/mom): 3.2%

Error metrics for skewed classes

Cancer classification example

Train logistic regression model $h_\theta(x)$. ($y = 1$ if cancer, $y = 0$ otherwise)

Find that you got 1% error on test set.
(99% correct diagnoses)

Only 0.50% of patients have cancer.

skewed classes.

```
function y = predictCancer(x)
    → y = 0; %ignore x!
    return
```

0.5% error
→ 99.2% accy (0.8% error)
→ 99.5% accy (0.5% error)

Andrew Ng

Precision/Recall

$y = 1$ in presence of rare class that we want to detect

		Actual class	
		1	0
Predicted 1 class	1	True positive	False positive
	0	False negative	True negative

$y = 0$
 $\text{recall} = 0$

→ **Precision**
(Of all patients where we predicted $y = 1$, what fraction actually has cancer?)

$$\frac{\text{True positives}}{\# \text{predicted positive}} = \frac{\text{True positive}}{\text{True pos} + \text{False pos}}$$

→ **Recall**
(Of all patients that actually have cancer, what fraction did we correctly detect as having cancer?)

$$\frac{\text{True positives}}{\# \text{actual positives}} = \frac{\text{True positives}}{\text{True pos} + \text{False neg}}$$

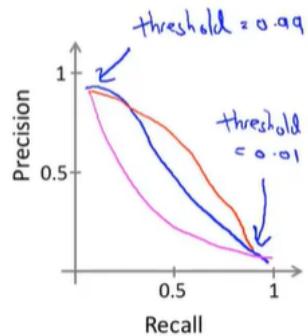
Trading off precision and recall

Trading off precision and recall

- Logistic regression: $0 \leq h_\theta(x) \leq 1$
- Predict 1 if $h_\theta(x) \geq 0.5$ ↗ ↘ ↛ 0.3 ←
- Predict 0 if $h_\theta(x) < 0.5$ ↗ ↘ ↛ 0.3
- Suppose we want to predict $y = 1$ (cancer) only if very confident.
→ Higher precision, lower recall.
- Suppose we want to avoid missing too many cases of cancer (avoid false negatives).
→ Higher recall, lower precision.

$$\rightarrow \text{precision} = \frac{\text{true positives}}{\text{no. of predicted positive}}$$

$$\rightarrow \text{recall} = \frac{\text{true positives}}{\text{no. of actual positive}}$$



More generally: Predict 1 if $h_\theta(x) \geq \text{threshold}$. ←

Andrew Ng

F_1 Score (F score)

How to compare precision/recall numbers?

	Precision(P)	Recall (R)	Average	F_1 Score
→ Algorithm 1	0.5	0.4	0.45	0.444 ←
→ Algorithm 2	0.7	0.1	0.4	0.175 ←
Algorithm 3	0.02	1.0	0.51	0.0392 ←
Average: $\frac{P+R}{2}$			Predict $y=1$ all the time	
F_1 Score: $2 \frac{PR}{P+R}$			$P=0 \text{ or } R=0 \Rightarrow F\text{-score} = 0.$ $P=1 \text{ and } R=1 \Rightarrow F\text{-score} = 1$	

Andrew Ng

Data for machine learning

Large data rationale

→ Assume feature $x \in \mathbb{R}^{n+1}$ has sufficient information to predict y accurately. ↗

Example: For breakfast I ate two eggs. ↗ ↛ ↘ ↙ ↝ ↚

Counterexample: Predict housing price from only size feet^2 and no other features. ↗

Useful test: Given the input x , can a human expert confidently predict y ? ↗

Large data rationale

→ Use a learning algorithm with many parameters (e.g. logistic regression/linear regression with many features; neural network with many hidden units). low bias algorithms. ←

→ $J_{\text{train}}(\theta)$ will be small.

Use a very large training set (unlikely to overfit) low variance ←

→ $J_{\text{train}}(\theta) \approx J_{\text{test}}(\theta)$

→ $J_{\text{test}}(\theta)$ will be small

Support Vector Machines

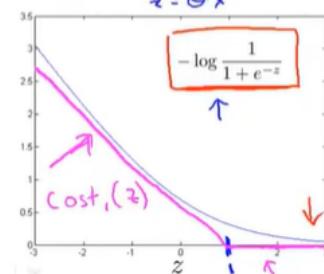
Optimization objective

Alternative view of logistic regression (x, y)

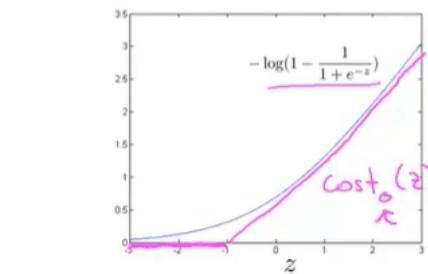
Cost of example: $-(y \log h_{\theta}(x) + (1 - y) \log(1 - h_{\theta}(x)))$ ←

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

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



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



Andrew Ng

Support vector machine

Logistic regression:

$$\min_{\theta} \frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \underbrace{\left(-\log h_{\theta}(x^{(i)}) \right)}_{\text{cost}_1(\theta^T x^{(i)})} + (1 - y^{(i)}) \underbrace{\left(-\log(1 - h_{\theta}(x^{(i)})) \right)}_{\text{cost}_0(\theta^T x^{(i)})} \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

Support vector machine:

$$\min_{\theta} C \sum_{i=1}^m y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) + \frac{1}{2} \sum_{j=0}^n \theta_j^2$$

$$\min_u \frac{(u - S)^2 + 1}{2} \rightarrow u = S$$

$$\min_u \log(u - S)^2 + 1 \rightarrow u = S$$

$$\left| \begin{array}{l} A + \lambda B \\ \rightarrow C = A + B \end{array} \right.$$

$$C = \frac{1}{\lambda}$$

$$\min_{\theta} C \sum_{i=1}^m \left[y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

Andrew Ng

SVM hypothesis

$$\Rightarrow \min_{\theta} C \sum_{i=1}^m [y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)})] + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

Hypothesis:

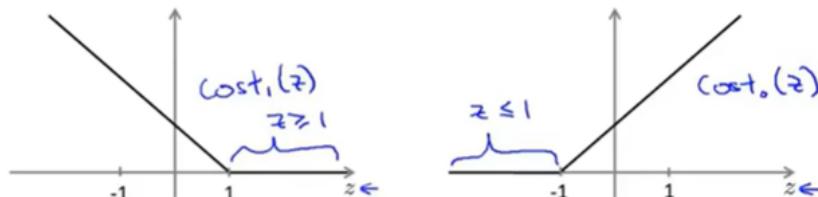
$$h_{\theta}(x) = \begin{cases} 1 & \text{if } \theta^T x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

.

Large Margin Intuition

Support Vector Machine

$$\rightarrow \min_{\theta} \underbrace{C \sum_{i=1}^m [y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)})]}_{\uparrow} + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$



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

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

$$C = 100,000$$

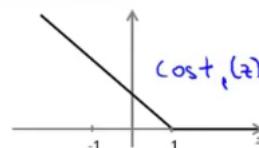
Andrew Ng

SVM Decision Boundary

$$\min_{\theta} \underbrace{C \sum_{i=1}^m [y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T x^{(i)})]}_{\uparrow} + \frac{1}{2} \sum_{j=1}^n \theta_j^2$$

Whenever $y^{(i)} = 1$:

$$\theta^T x^{(i)} \geq 1$$



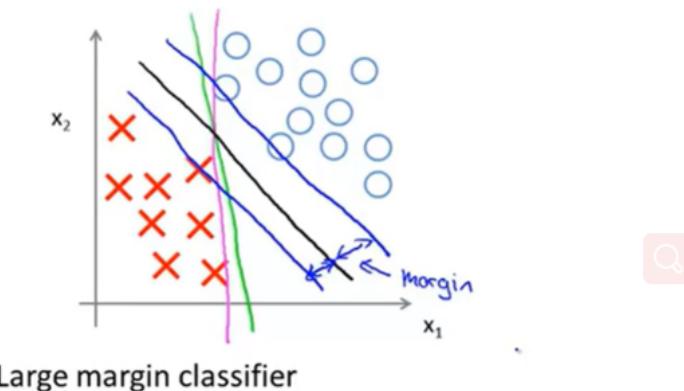
Whenever $y^{(i)} = 0$:

$$\theta^T x^{(i)} \leq -1$$



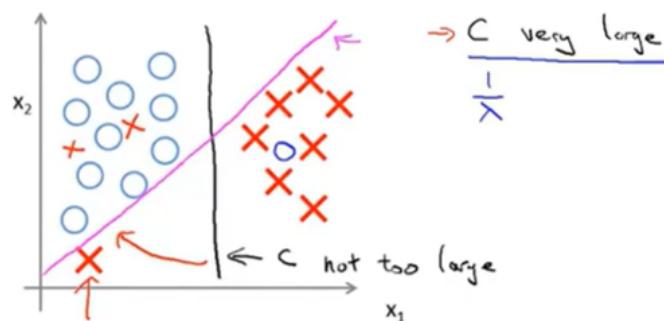
Andrew Ng

SVM Decision Boundary: Linearly separable case



Andrew Ng

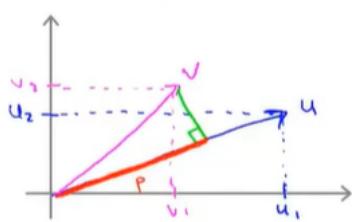
Large margin classifier in presence of outliers



Andrew Ng

The mathematics behind large margin classification

Vector Inner Product



$$\Rightarrow u = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \rightarrow v = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

$$u^T v = ? \quad [u_1 \ u_2] \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

$\|u\|$ = length of vector u

$$= \sqrt{u_1^2 + u_2^2} \in \mathbb{R}$$

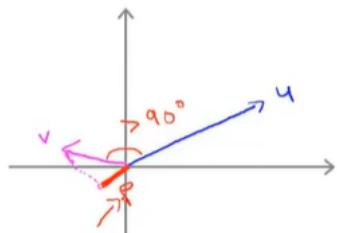
p = length of projection of v onto u .

$$u^T v = p \cdot \|u\| \leftarrow = v^T u$$

$$= u_1 v_1 + u_2 v_2 \leftarrow p \in \mathbb{R}$$

$$u^T v = p \cdot \|u\|$$

$$p < 0$$



Andrew Ng

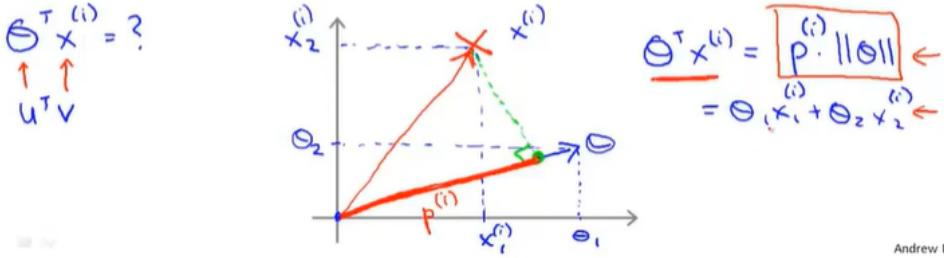
$$\omega = (\sqrt{\omega})^2$$

SVM Decision Boundary

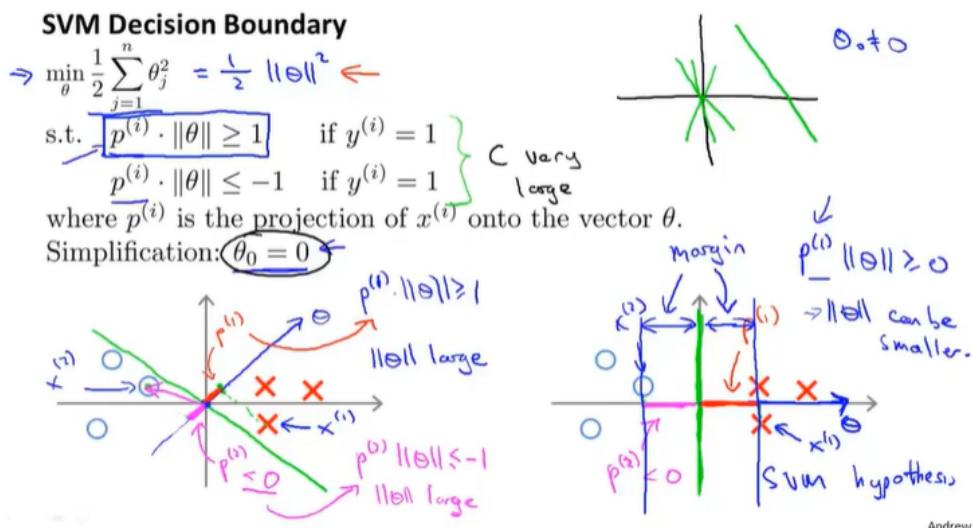
$$\min_{\theta} \frac{1}{2} \sum_{j=1}^n \theta_j^2 = \frac{1}{2} (\Theta_1^2 + \Theta_2^2) = \frac{1}{2} (\underbrace{\Theta_1^2 + \Theta_2^2}_{= \|\theta\|^2}) = \frac{1}{2} \|\theta\|^2$$

s.t. $\theta^T x^{(i)} \geq 1$ if $y^{(i)} = 1$
 $\theta^T x^{(i)} \leq -1$ if $y^{(i)} = 0$

Simplification: $\Theta_0 = 0$, $n=2$

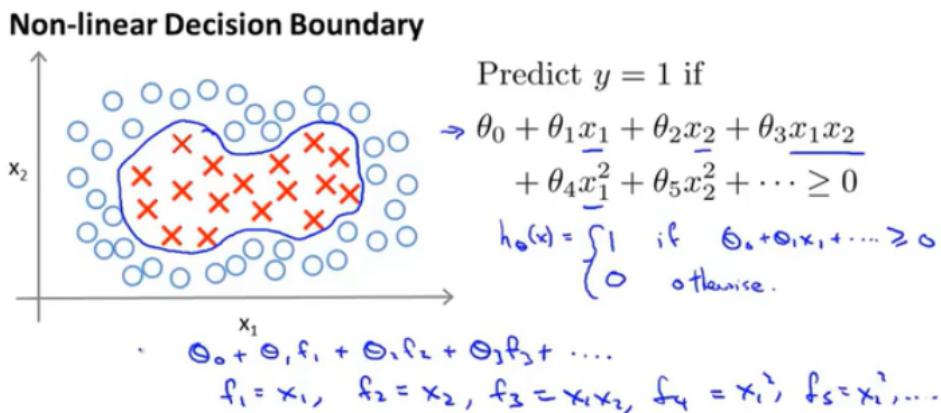


Andrew Ng



Andrew Ng

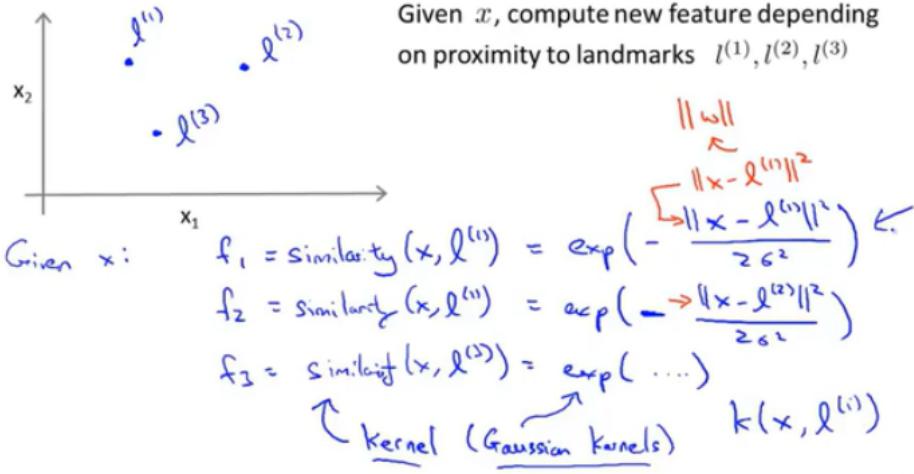
Kernels



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

Andrew Ng

Kernel



Andrew Ng

Kernels and Similarity

$$f_1 = \text{similarity}(x, l^{(1)}) = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right) = \exp\left(-\frac{\sum_{j=1}^n (x_j - l_j^{(1)})^2}{2\sigma^2}\right)$$

If $x \approx l^{(1)}$:

$$f_1 \approx \exp\left(-\frac{0^2}{2\sigma^2}\right) \approx 1$$

$$\begin{aligned} l^{(1)} &\rightarrow f_1 \\ l^{(2)} &\rightarrow f_2 \\ l^{(3)} &\rightarrow f_3. \end{aligned}$$

If x if far from $l^{(1)}$:

$$f_1 = \exp\left(-\frac{(\text{large number})^2}{2\sigma^2}\right) \approx 0.$$

Andrew Ng

Example:

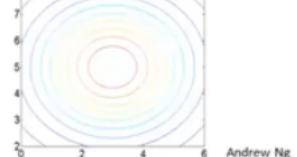
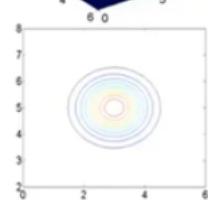
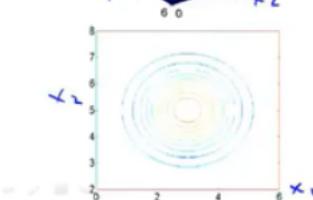
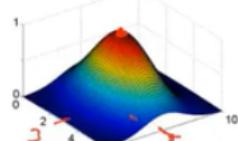
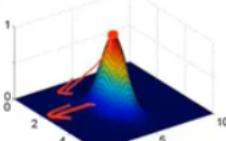
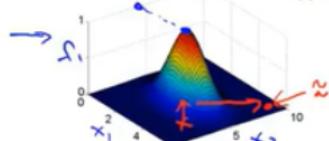
$$\rightarrow l^{(1)} = \begin{bmatrix} 3 \\ 5 \end{bmatrix}, \quad f_1 = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right)$$

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

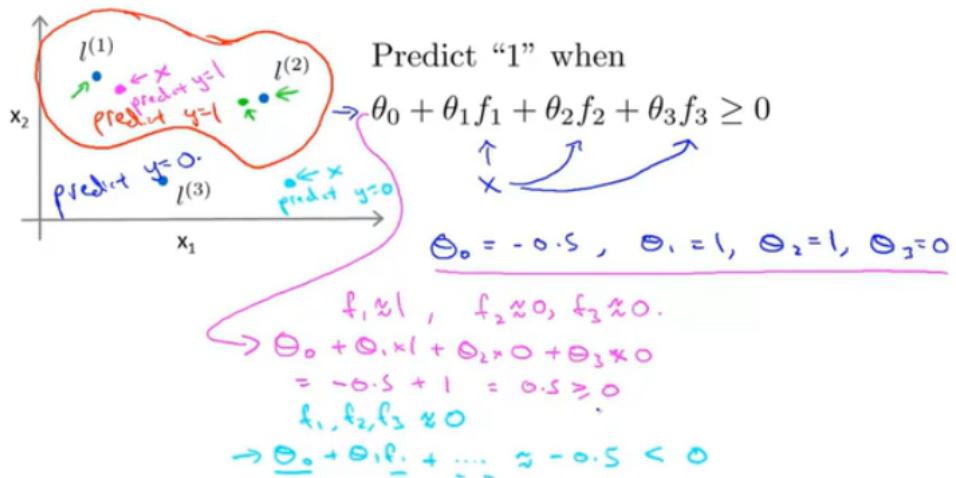
$$x = \begin{bmatrix} 3 \\ 5 \end{bmatrix}$$

$$\sigma^2 = 0.5$$

$$\sigma^2 = 3$$

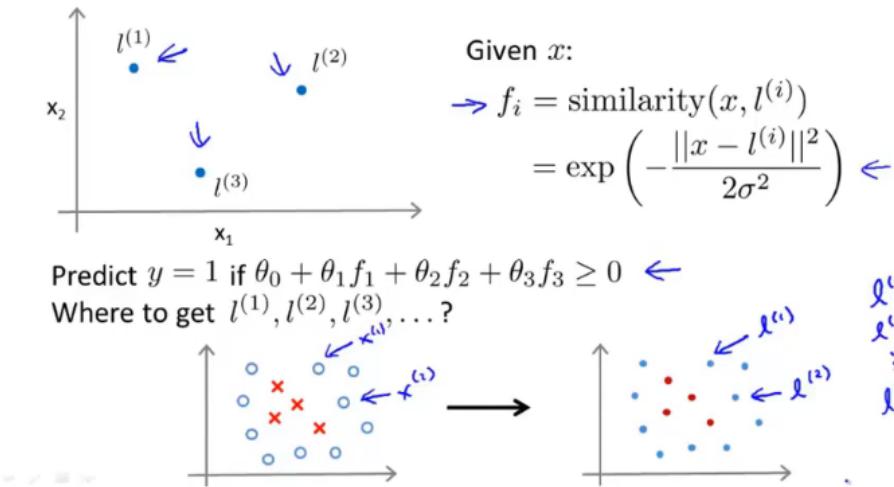


Andrew Ng



Andrew Ng

Choosing the landmarks



Andrew Ng

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 :

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

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

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

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

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

$$\begin{aligned} x^{(i)} \in \mathbb{R}^{n+1} & \quad (\text{or } \mathbb{R}^n) \\ f^{(i)} &= \begin{bmatrix} f_0^{(i)} \\ f_1^{(i)} \\ \vdots \\ f_m^{(i)} \end{bmatrix} \end{aligned}$$

Andrew Ng

SVM with Kernels

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

$$\rightarrow \text{Predict "y=1" if } \theta^T f \geq 0$$

$\theta = \theta_0 + \theta_1 f_1 + \dots + \theta_m f_m$

Training:

$$\min_{\theta} C \sum_{i=1}^m y^{(i)} \text{cost}_1(\theta^T f^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T f^{(i)}) + \frac{1}{2} \sum_{j=1}^m \theta_j^2$$

$\rightarrow \theta_0 = \sum_j \theta_j f_j$

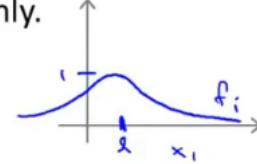
$\rightarrow \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_m \end{bmatrix}$ (ignoring θ_0) $M = 10,000$

Andrew Ng

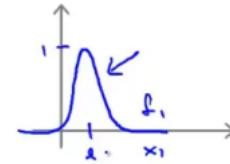
SVM parameters:

$C = \frac{1}{\lambda}$. \rightarrow Large C : Lower bias, high variance. (small λ)
 \rightarrow Small C : Higher bias, low variance. (large λ)

σ^2 Large σ^2 : Features f_i vary more smoothly.
 \rightarrow Higher bias, lower variance.

$$\exp\left(-\frac{\|x - l^{(i)}\|^2}{2\sigma^2}\right)$$


Small σ^2 : Features f_i vary less smoothly.
Lower bias, higher variance.



Andrew Ng

Using an SVM

Use SVM software package (e.g. liblinear, libsvm, ...) to solve for parameters θ .

Need to specify:

\rightarrow Choice of parameter C .

Choice of kernel (similarity function):

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

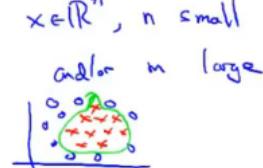
Predict "y = 1" if $\theta^T x \geq 0$

$\theta_0 + \theta_1 x_1 + \dots + \theta_n x_n \geq 0$ $x \in \mathbb{R}^{n+1}$

Gaussian kernel:

$$f_i = \exp\left(-\frac{\|x - l^{(i)}\|^2}{2\sigma^2}\right), \text{ where } l^{(i)} = x^{(i)}$$

Need to choose σ^2 .



Andrew Ng

Kernel (similarity) functions:

```

function f = kernel(x1, x2)
     $f = \exp\left(\frac{-\|x_1 - x_2\|^2}{2\sigma^2}\right)$ 
return

```

$x \rightarrow \begin{matrix} f_1 \\ f_2 \\ \vdots \\ f_m \end{matrix}$

→ Note: Do perform feature scaling before using the Gaussian kernel.

$$\begin{aligned}
& \boxed{\|x - l\|^2} \\
& \quad \begin{aligned}
v &= x - l \\
\|v\|^2 &= v_1^2 + v_2^2 + \dots + v_n^2 \\
&= (x_1 - l_1)^2 + (x_2 - l_2)^2 + \dots + (x_n - l_n)^2
\end{aligned}
\end{aligned}$$

$\underbrace{\quad}_{\text{1000 feet}^2} \quad \underbrace{\quad}_{\text{1-5 bedrooms}}$

Andrew Ng

Other choices of kernel

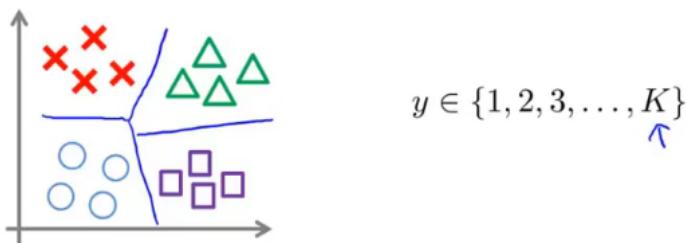
Note: Not all similarity functions $\text{similarity}(x, l)$ make valid kernels.

→ (Need to satisfy technical condition called "Mercer's Theorem" to make sure SVM packages' optimizations run correctly, and do not diverge).

- Many off-the-shelf kernels available:
- Polynomial kernel: $k(x, l) = (x^T l + \text{constant})^{\text{degree}}$
 - More esoteric: String kernel, chi-square kernel, histogram intersection kernel, ...

Andrew Ng

Multi-class classification



Many SVM packages already have built-in multi-class classification functionality.

→ Otherwise, use one-vs.-all method. (Train K SVMs, one to distinguish $y = i$ from the rest, for $i = 1, 2, \dots, K$), get $\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(K)}$. Pick class i with largest $\underline{(\theta^{(i)})^T x}$

Andrew Ng

Logistic regression vs. SVMs

n = number of features ($x \in \mathbb{R}^{n+1}$), m = number of training examples

- If n is large (relative to m): (e.g. $n \geq m$, $n = 10,000$, $m = 10 \dots 1000$)
 - Use logistic regression, or SVM without a kernel ("linear kernel")
- If n is small, m is intermediate:
 $(n = 1-1000, m = 10-10,000)$ ←
 - Use SVM with Gaussian kernel
- If n is small, m is large: $(n = 1-1000, m = 50,000+)$
 - Create/add more features, then use logistic regression or SVM without a kernel
- Neural network likely to work well for most of these settings, but may be slower to train.

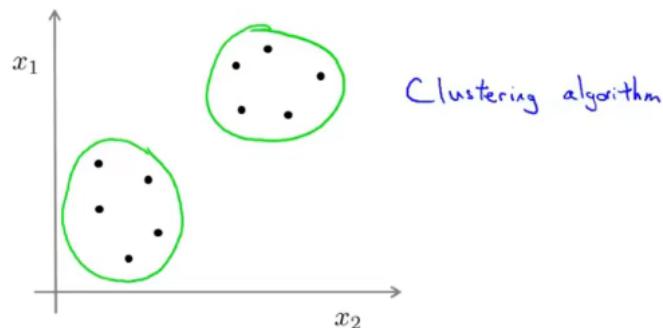


Andrew Ng

Clustering

Unsupervised learning introduction

Unsupervised learning



Training set: $\{x^{(1)}, x^{(2)}, x^{(3)}, \dots, x^{(m)}\}$ ← .

Andrew Ng

K-means algorithm

K-means algorithm

Input:

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

$x^{(i)} \in \mathbb{R}^n$ (drop $x_0 = 1$ convention)

K-means algorithm

$$\mu_1 \quad \mu_2$$

Randomly initialize K cluster centroids $\mu_1, \mu_2, \dots, \mu_K \in \mathbb{R}^n$

Repeat {

Cluster assignment step

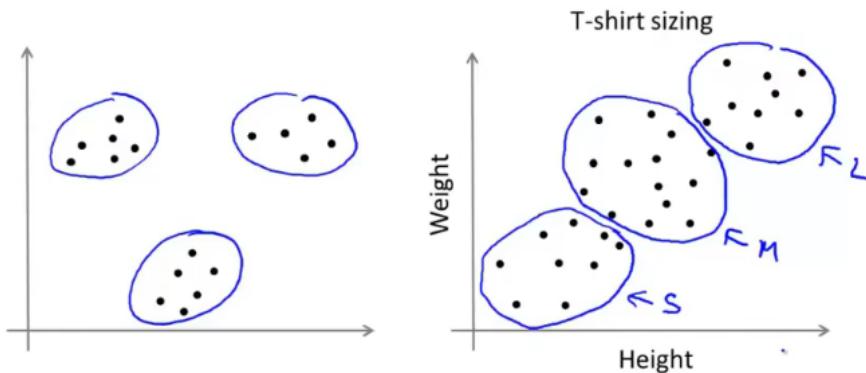
```

for i = 1 to m
    c(i) := index (from 1 to K) of cluster centroid
    closest to x(i)
    min ||x(i) - μk||2
    ↪ c(i)
for k = 1 to K
    → μk := average (mean) of points assigned to cluster k
    x(1), x(2), x(3), x(4) → c(1)=2, c(2)=2, c(3)=2, c(4)=2
    } . μ2 = 1/4 [x(1) + x(2) + x(3) + x(4)] ∈ ℝn
  
```

Andrew Ng

K-means for non-separated clusters

S, M, L



Andrew Ng

Optimization objective

K-means optimization objective

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

→ μ_k = cluster centroid k ($\mu_k \in \mathbb{R}^n$) K $k \in \{1, 2, \dots, K\}$

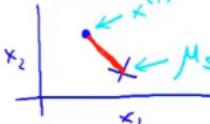
$\mu_{c^{(i)}}$ = cluster centroid of cluster to which example $x^{(i)}$ has been assigned $x^{(i)} \rightarrow S$ $c^{(i)}=S$ $\mu_{c^{(i)}} = \mu_S$

Optimization objective:

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

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

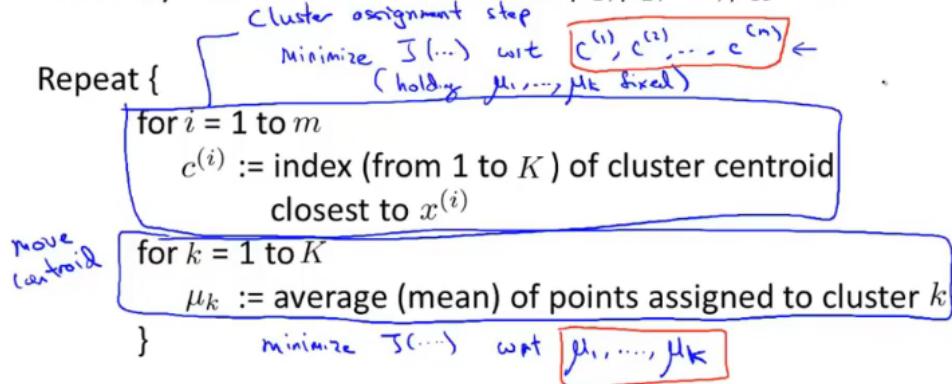
Distortion



Andrew Ng

K-means algorithm

Randomly initialize K cluster centroids $\mu_1, \mu_2, \dots, \mu_K \in \mathbb{R}^n$



Andrew Ng

Random initialization

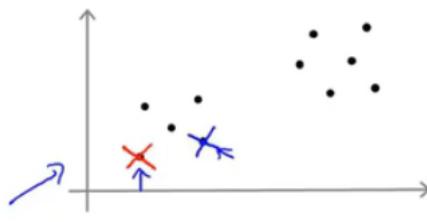
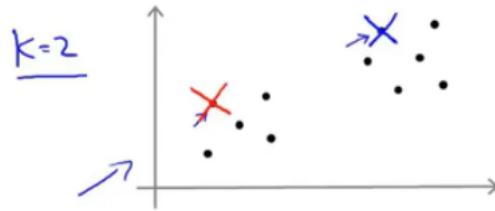
Random initialization

Should have $\underline{K < m}$

Randomly pick \underline{K} training examples.

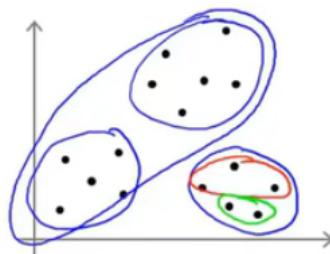
Set μ_1, \dots, μ_K equal to these K examples.

$$\begin{aligned}\mu_1 &= x^{(1)} \\ \mu_2 &= x^{(2)} \\ &\vdots\end{aligned}$$

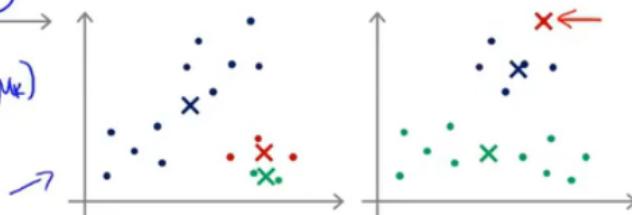
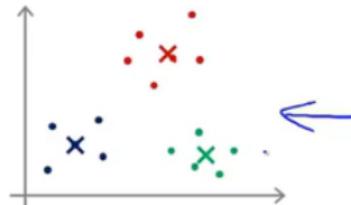


Andrew Ng

Local optima



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



Andrew Ng

Random initialization

For $i = 1$ to $100\}$ $S_0 = 1000$

 → Randomly initialize K-means.
 Run K-means. Get $c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K$.
 Compute cost function (distortion)
 → $J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$

Pick clustering that gave lowest cost $J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$

$K=2-10$

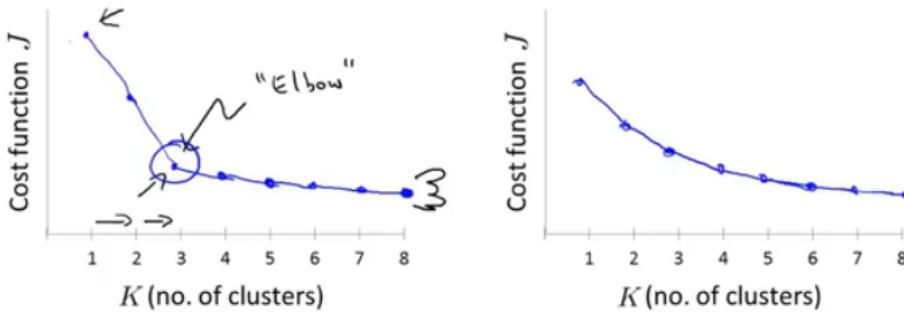
λ

Andrew Ng

Choosing the number of clusters

Choosing the value of K

Elbow method:



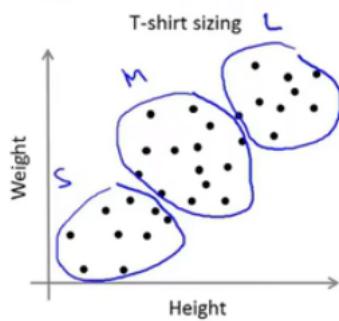
Andrew Ng

Choosing the value of K

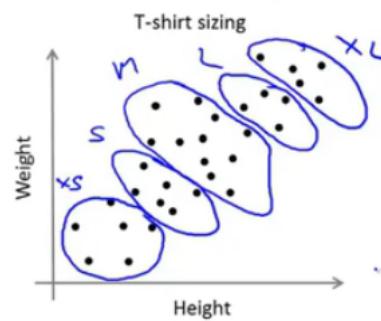
Sometimes, you're running K-means to get clusters to use for some later/downstream purpose. Evaluate K-means based on a metric for how well it performs for that later purpose.

$K=3$ S, M, L

E.g.



$K=5$ XS, S, M, L, XL

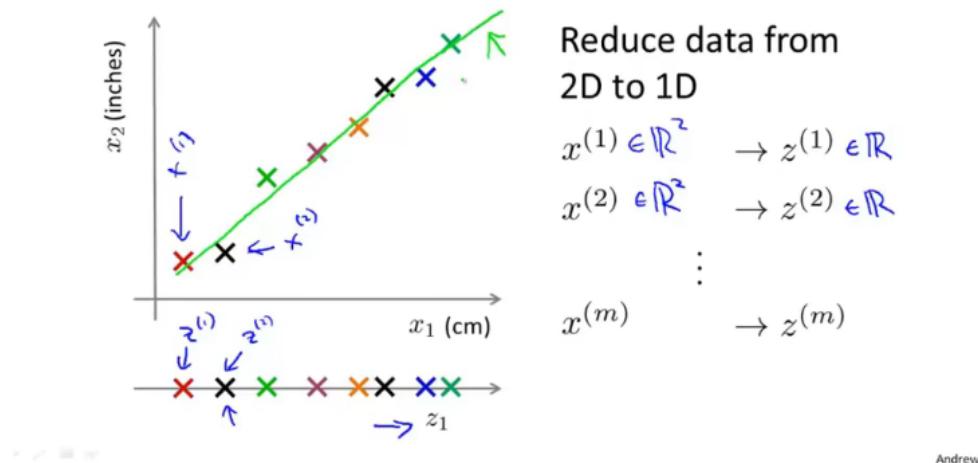


Andrew Ng

Dimensionality Reduction

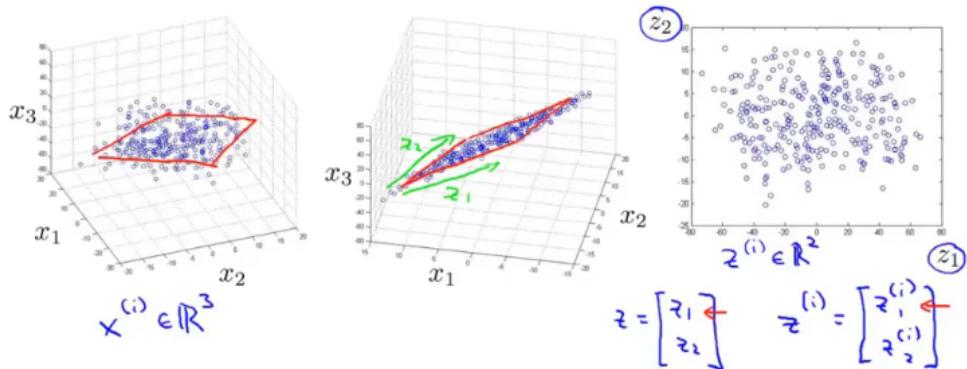
Motivation 1: Data Compression

Data Compression



Data Compression

$10000 \rightarrow 1000$
Reduce data from 3D to 2D



Motivation 2: Data Visualization

Data Visualization

Country	x_1 GDP (trillions of US\$)	x_2 Per capita GDP (thousands of intl. \$)	x_3 Human Development Index	x_4 Life expectancy	x_5 Poverty Index (Gini as percentage)	x_6 Mean household income (thousands of US\$)	...
Canada	1.577	39.17	0.908	80.7	32.6	67.293	...
China	5.878	7.54	0.687	73	46.9	10.22	...
India	1.632	3.41	0.547	64.7	36.8	0.735	...
Russia	1.48	19.84	0.755	65.5	39.9	0.72	...
Singapore	0.223	56.69	0.866	80	42.5	67.1	...
USA	14.527	46.86	0.91	78.3	40.8	84.3	...
...

[resources from en.wikipedia.org]

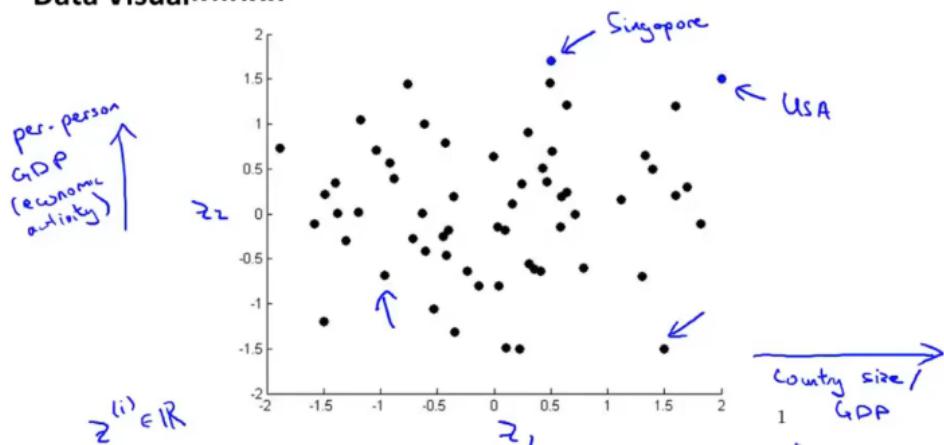
Andrew Ng

Data Visualization

Country	z_1	z_2	$z^{(i)} \in \mathbb{R}^2$
Canada	1.6	1.2	
China	1.7	0.3	
India	1.6	0.2	Reduce data from 500 to 2D
Russia	1.4	0.5	
Singapore	0.5	1.7	
USA	2	1.5	
...	

Andrew Ng

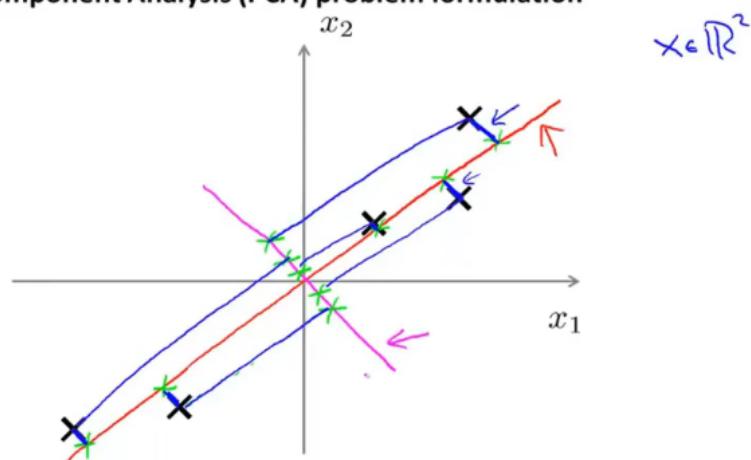
Data Visualization



Andrew Ng

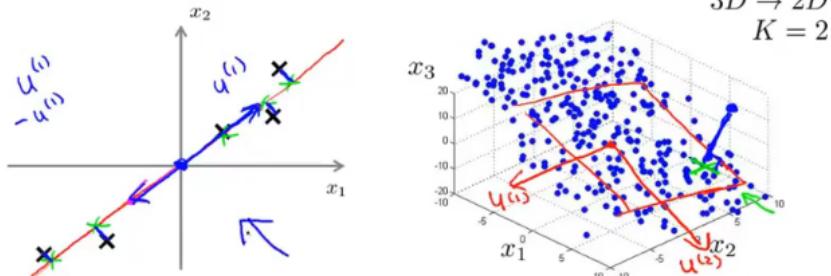
Principal Component Analysis problem formulation

Principal Component Analysis (PCA) problem formulation



Andrew Ng

Principal Component Analysis (PCA) problem formulation

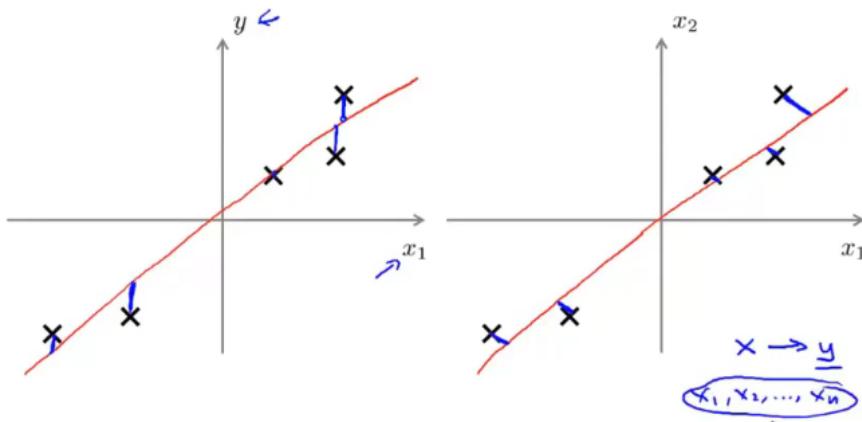


Reduce from 2-dimension to 1-dimension: Find a direction (a vector $u^{(1)} \in \mathbb{R}^n$) onto which to project the data so as to minimize the projection error.

Reduce from n -dimension to k -dimension: Find k vectors $u^{(1)}, u^{(2)}, \dots, u^{(k)}$ onto which to project the data, so as to minimize the projection error.

Andrew Ng

PCA is not linear regression



Andrew Ng

Data preprocessing

Training set: $x^{(1)}, x^{(2)}, \dots, x^{(m)}$

Preprocessing (feature scaling/mean normalization):

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

Replace each $x_j^{(i)}$ with $x_j^{(i)} - \mu_j$.

If different features on different scales (e.g., x_1 = size of house, x_2 = number of bedrooms), scale features to have comparable range of values.

$$x_j^{(i)} \leftarrow \frac{x_j^{(i)} - \mu_j}{\sigma_j}$$

Andrew Ng

Principal Component Analysis (PCA) algorithm

Reduce data from n -dimensions to k -dimensions

Compute "covariance matrix":

$$\Sigma = \frac{1}{m} \sum_{i=1}^n (x^{(i)})(x^{(i)})^T \quad n \times n \quad \text{Sigma}$$

Compute "eigenvectors" of matrix Σ :

$$\rightarrow [U, S, V] = \text{svd}(\text{Sigma}); \quad \rightarrow \begin{array}{l} \text{Singular value decomposition} \\ \text{eig}(\text{Sigma}) \end{array}$$

$$U = \begin{bmatrix} | & | & | & \dots & | \\ u^{(1)} & u^{(2)} & u^{(3)} & \dots & u^{(n)} \\ | & | & | & \dots & | \\ k & & & & \end{bmatrix} \quad U \in \mathbb{R}^{n \times n} \quad u^{(1)}, \dots, u^{(k)}$$

Andrew Ng

Principal Component Analysis (PCA) algorithm

From $[U, S, V] = \text{svd}(\text{Sigma})$, we get:

$$\rightarrow U = \begin{bmatrix} | & | & | & \dots & | \\ u^{(1)} & u^{(2)} & u^{(3)} & \dots & u^{(n)} \\ | & | & | & \dots & | \\ k & & & & \end{bmatrix} \in \mathbb{R}^{n \times n}$$

$$\begin{aligned} x \in \mathbb{R}^n &\rightarrow z \in \mathbb{R}^k \\ \cdot z^{(i)} = \begin{bmatrix} | & | & | & \dots & | \\ u^{(1)} & u^{(2)} & u^{(3)} & \dots & u^{(k)} \\ | & | & | & \dots & | \\ n \times k & & & & \end{bmatrix}^T & x^{(i)} = \begin{bmatrix} -(u^{(1)})^T \\ \vdots \\ -(u^{(k)})^T \\ k \times n \end{bmatrix}_{k \times 1} \\ z \in \mathbb{R}^k && \end{aligned}$$

Andrew Ng

Principal Component Analysis (PCA) algorithm summary

\rightarrow After mean normalization (ensure every feature has zero mean) and optionally feature scaling:

$$\begin{aligned} \text{Sigma} &= \frac{1}{m} \sum_{i=1}^m (x^{(i)})(x^{(i)})^T & X = \begin{bmatrix} -x^{(1)\top} \\ \vdots \\ -x^{(m)\top} \end{bmatrix} \\ \rightarrow [U, S, V] &= \text{svd}(\text{Sigma}); & \text{Sigma} = (1/m) \times X' \times X \\ \rightarrow U_{\text{reduce}} &= U(:, 1:k); & \\ \rightarrow z &= U_{\text{reduce}}' * x; & \begin{array}{c} \uparrow \quad \uparrow \\ x \in \mathbb{R}^n \quad x \neq 1 \end{array} \end{aligned}$$

Andrew Ng

Choosing the number of principal components

Choosing k (number of principal components)

Average squared projection error: $\frac{1}{m} \sum_{i=1}^m \|x^{(i)} - x_{approx}^{(i)}\|^2$

Total variation in the data: $\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2$

Typically, choose k to be smallest value so that

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

$$\rightarrow \frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2 \leq \frac{0.05}{0.10} \quad (50\%)$$

\rightarrow "99% of variance is retained"
goes to 90%

Andrew Ng

Choosing k (number of principal components)

Algorithm:

Try PCA with $k=1, k=2, k=3, k=4$

Compute $U_{reduce}, z^{(1)}, z^{(2)}, \dots, z^{(m)}, x_{approx}^{(1)}, \dots, x_{approx}^{(m)}$

Check if

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

$k=17$

$\rightarrow [U, S, V] = svd(\Sigma)$

$\rightarrow S = \begin{bmatrix} S_{11} & & & \\ & S_{22} & & \\ & & S_{33} & \\ & & & \ddots & S_{nn} \end{bmatrix}$

For given k $k=3$

$$1 - \frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^m S_{ii}} \leq 0.01$$

$$\rightarrow \frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^m S_{ii}} \geq 0.99$$

Andrew Ng

Choosing k (number of principal components)

$\rightarrow [U, S, V] = svd(\Sigma)$

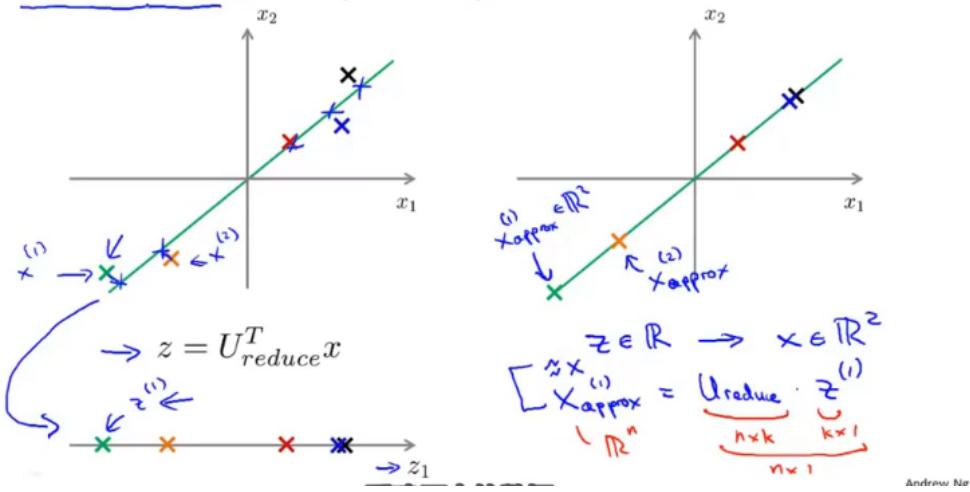
Pick smallest value of k for which

$$\frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^m S_{ii}} \geq 0.99 \quad k=100$$

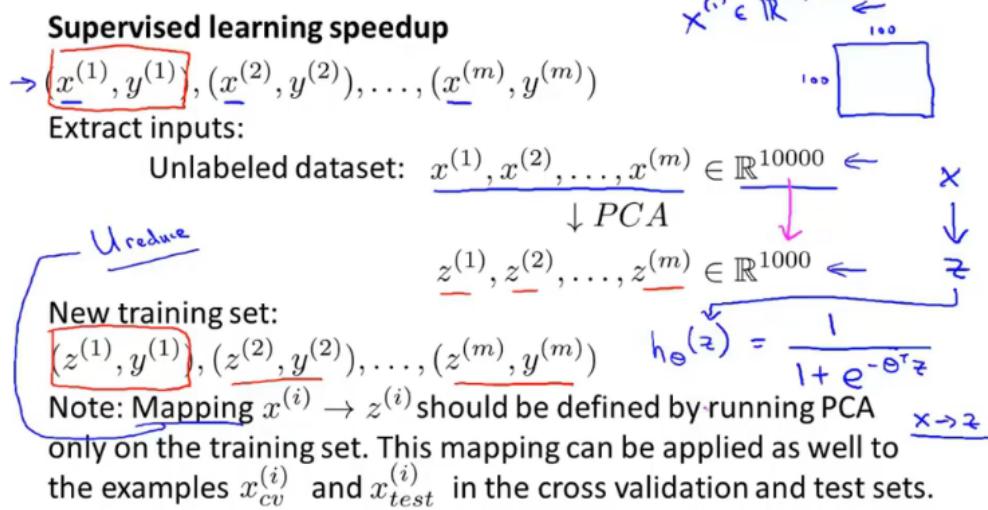
(99% of variance retained)

Reconstruction from compressed representation

Reconstruction from compressed representation



Advice for applying PCA



Application of PCA

- Compression

- Reduce memory/disk needed to store data
- Speed up learning algorithm ←

Choose k by % of variance retain

← - Visualization

$k=2$ or $k=3$

Andrew Ng

Bad use of PCA: To prevent overfitting

- Use $\underline{z}^{(i)}$ instead of $\underline{x}^{(i)}$ to reduce the number of features to $\underline{k} < \underline{n}$.
Thus, fewer features, less likely to overfit.

Bad!

This might work OK, but isn't a good way to address overfitting. Use regularization instead.

$$\rightarrow \min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \boxed{\frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2} \quad \leftarrow$$

Andrew Ng

PCA is sometimes used where it shouldn't be

Design of ML system:

- - Get training set $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$
- - ~~Run PCA to reduce $x^{(i)}$ in dimension to get $\underline{z}^{(i)}$~~
- - Train logistic regression on $\{(z^{(1)}, y^{(1)}), \dots, (z^{(m)}, y^{(m)})\}$
- - Test on test set: Map $x_{test}^{(i)}$ to $z_{test}^{(i)}$. Run $h_{\theta}(z)$ on $\{(z_{test}^{(1)}, y_{test}^{(1)}), \dots, (z_{test}^{(m)}, y_{test}^{(m)})\}$
- How about doing the whole thing without using PCA?
- Before implementing PCA, first try running whatever you want to do with the original/raw data $\underline{x}^{(i)}$. Only if that doesn't do what you want, then implement PCA and consider using $\underline{z}^{(i)}$.

Andrew Ng

Anomaly detection

Problem motivation

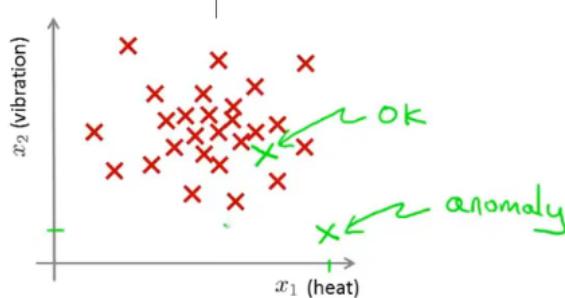
Anomaly detection example

Aircraft engine features:

- x_1 = heat generated
- x_2 = vibration intensity
- ...

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

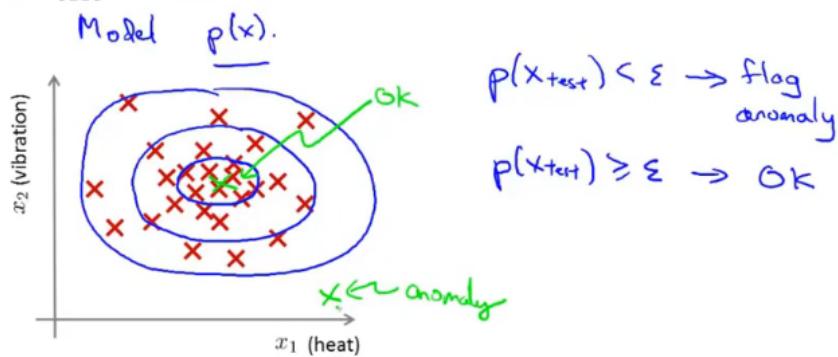
New engine: \underline{x}_{test}



Andrew Ng

Density estimation

- Dataset: $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$
- Is x_{test} anomalous?



Andrew Ng

Anomaly detection example

- Fraud detection:
 - $x^{(i)}$ = features of user i 's activities
 - Model $p(x)$ from data.
 - Identify unusual users by checking which have $p(x) < \varepsilon$
- Manufacturing
- Monitoring computers in a data center.
 - $x^{(i)}$ = features of machine i
 - x_1 = memory use, x_2 = number of disk accesses/sec,
 - x_3 = CPU load, x_4 = CPU load/network traffic.
 - ... $p(x) < \varepsilon$

$$\begin{matrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{matrix} \quad p(x)$$

Andrew Ng

Gaussian distribution

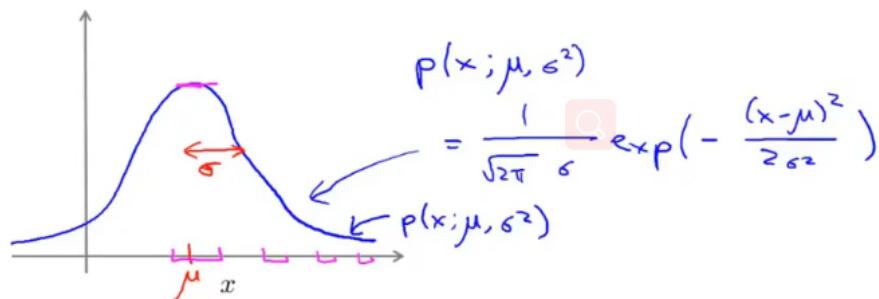
Gaussian (Normal) distribution

Say $x \in \mathbb{R}$. If x is a distributed Gaussian with mean μ , variance σ^2 .

$$x \sim \mathcal{N}(\mu, \sigma^2)$$

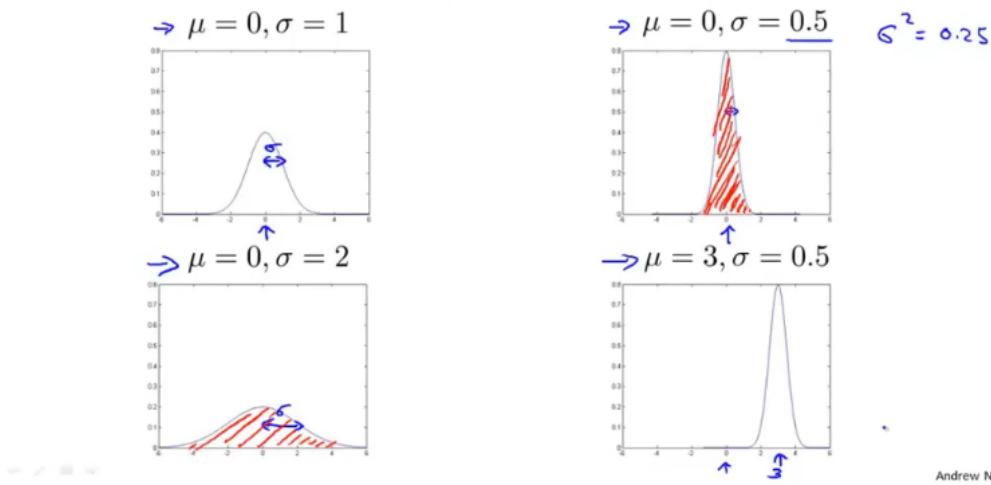
"distributed as"

σ standard deviation



Andrew Ng

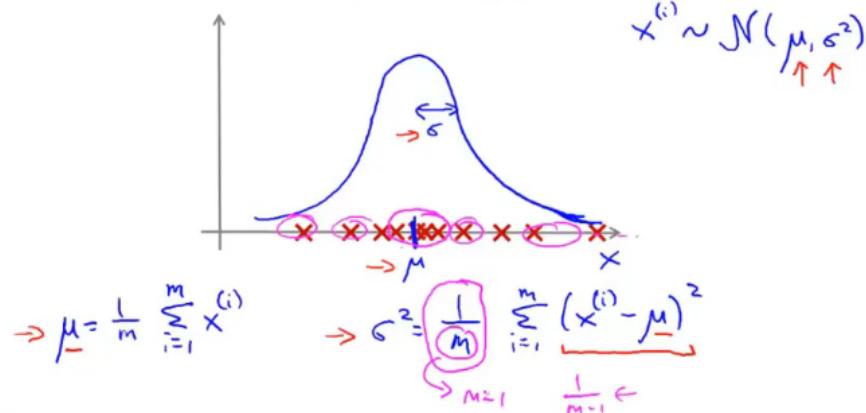
Gaussian distribution example



Andrew Ng

Parameter estimation

→ Dataset: $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\}$ $x^{(i)} \in \mathbb{R}$



Andrew Ng

Algorithm

→ Density estimation

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

Each example is $x \in \mathbb{R}^n$

→ $p(x)$

$$= p(x_1; \mu_1, \sigma_1^2) p(x_2; \mu_2, \sigma_2^2) p(x_3; \mu_3, \sigma_3^2) \cdots p(x_n; \mu_n, \sigma_n^2) \quad \leftarrow$$

$$= \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2)$$

$$x_1 \sim \mathcal{N}(\mu_1, \sigma_1^2)$$

$$x_2 \sim \mathcal{N}(\mu_2, \sigma_2^2)$$

$$x_3 \sim \mathcal{N}(\mu_3, \sigma_3^2)$$

$$\sum_{i=1}^n i = 1+2+3+\dots+n$$

$$\prod_{i=1}^n i = 1 \times 2 \times 3 \times \dots \times n$$

Andrew Ng

Anomaly detection algorithm

- 1. Choose features x_i that you think might be indicative of anomalous examples. $\{x^{(1)}, \dots, x^{(m)}\}$
- 2. Fit parameters $\mu_1, \dots, \mu_n, \sigma_1^2, \dots, \sigma_n^2$

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

$$\rightarrow \sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$$

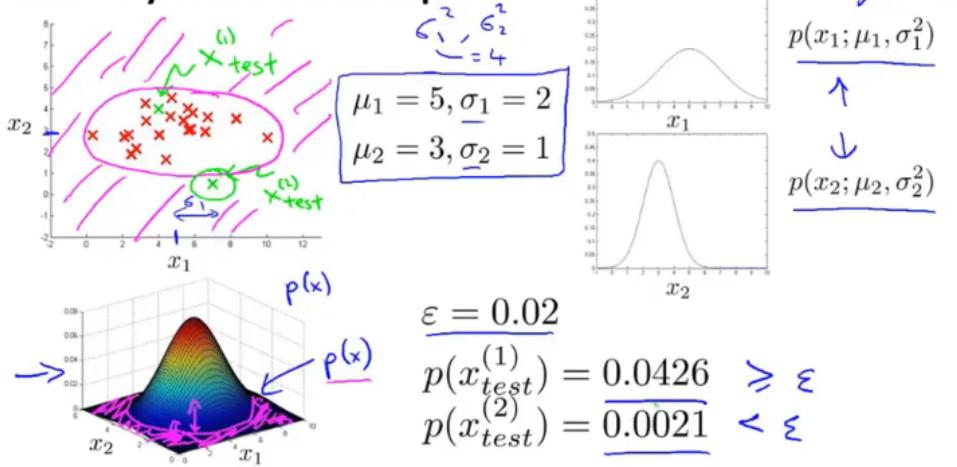
$$p(x_j; \mu_j, \sigma_j^2)$$

$$\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{bmatrix} = \frac{1}{m} \sum_{i=1}^m x^{(i)}$$
- 3. Given new example x , compute $p(x)$:
$$p(x) = \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi\sigma_j^2}} \exp\left(-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}\right)$$

Anomaly if $p(x) < \varepsilon$

Andrew Ng

Anomaly detection example



Andrew Ng

Developing and evaluating an anomaly detection system

The importance of real-number evaluation

When developing a learning algorithm (choosing features, etc.), making decisions is much easier if we have a way of evaluating our learning algorithm.

- Assume we have some labeled data, of anomalous and non-anomalous examples. ($y = 0$ if normal, $y = 1$ if anomalous).
- Training set: $x^{(1)}, x^{(2)}, \dots, x^{(m)}$ (assume normal examples/not anomalous)
- Cross validation set: $(x_{cv}^{(1)}, y_{cv}^{(1)}), \dots, (x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})})$
- Test set: $(x_{test}^{(1)}, y_{test}^{(1)}), \dots, (x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$

$y=1$

Andrew Ng

Aircraft engines motivating example

- 10000 good (normal) engines
 - 20 flawed engines (anomalous) $\frac{2-50}{y=1}$
 - Training set: 6000 good engines ($y=0$) $p(x) = p(x_i; \mu_1, \sigma^2_1) \dots p(x_i; \mu_n, \sigma^2_n)$
CV: 2000 good engines ($y=0$), 10 anomalous ($y=1$)
Test: 2000 good engines ($y=0$), 10 anomalous ($y=1$)
- Alternative:
Training set: 6000 good engines
CV: 4000 good engines ($y=0$), 10 anomalous ($y=1$)
Test: 4000 good engines ($y=0$), 10 anomalous ($y=1$)

Andrew Ng

Algorithm evaluation

- Fit model $p(x)$ on training set $\{x^{(1)}, \dots, x^{(m)}\}$ $(x_{\text{test}}^{(i)}, y_{\text{test}}^{(i)})$
- On a cross validation/test example x , predict

$$y = \begin{cases} 1 & \text{if } p(x) < \epsilon \text{ (anomaly)} \\ 0 & \text{if } p(x) \geq \epsilon \text{ (normal)} \end{cases} \quad y=0$$

Possible evaluation metrics:

- - True positive, false positive, false negative, true negative
 - - Precision/Recall
 - - F_1 -score
- \uparrow
CV
Test set.

Can also use cross validation set to choose parameter ϵ ↵

Andrew Ng

Anomaly detection vs. supervised learning

Anomaly detection

- Very small number of positive examples ($y=1$). (0-20 is common).
- Large number of negative ($y=0$) examples. $p(x)$ ↵
- Many different "types" of anomalies. Hard for any algorithm to learn from positive examples what the anomalies look like;
- future anomalies may look nothing like any of the anomalous examples we've seen so far.

vs.

Supervised learning

Large number of positive and negative examples. ↵

Enough positive examples for algorithm to get a sense of what positive examples are like, future positive examples likely to be similar to ones in training set. ↵

Spam ↵

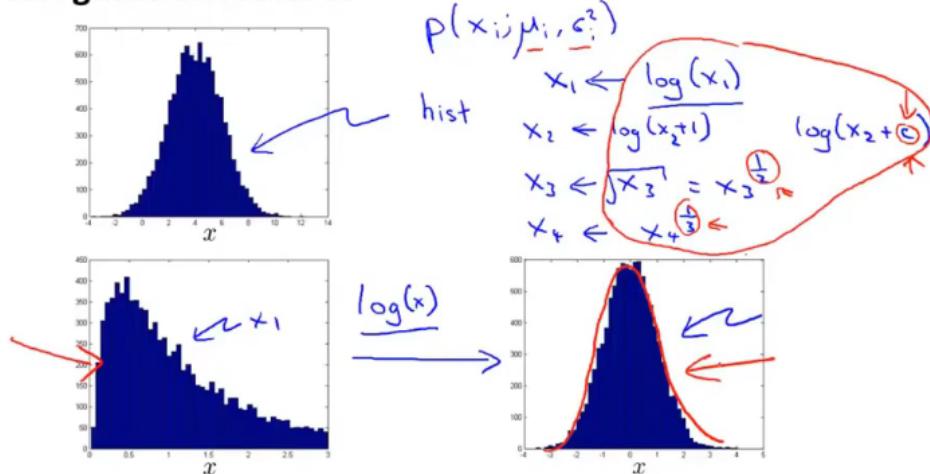
Andrew Ng

Anomaly detection	vs.	Supervised learning
→ • Fraud detection $y=1$		• Email spam classification \leftarrow
→ • Manufacturing (e.g. aircraft engines)		• Weather prediction \leftarrow (sunny/rainy/etc).
→ • Monitoring machines in a data center		• Cancer classification \leftarrow
⋮		⋮

Andrew Ng

Choosing what features to use

Non-gaussian features

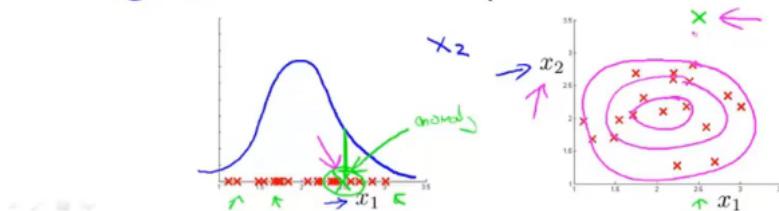


→ Error analysis for anomaly detection

[Want $p(x)$ large for normal examples x .
p(x) small for anomalous examples x .

Most common problem:

[$p(x)$ is comparable (say, both large) for normal and anomalous examples



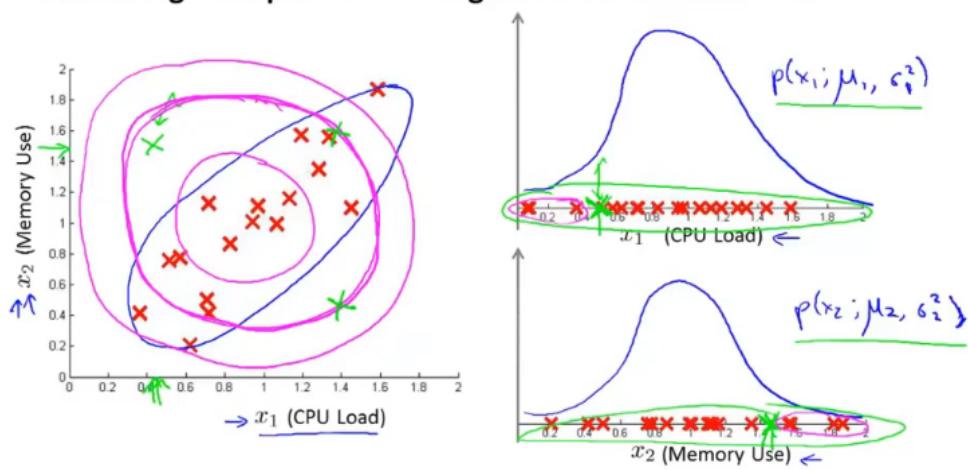
- Monitoring computers in a data center
- Choose features that might take on unusually large or small values in the event of an anomaly.
- x_1 = memory use of computer
- x_2 = number of disk accesses/sec
- x_3 = CPU load ←
- x_4 = network traffic ←

$$x_5 = \frac{\text{CPU load}}{\text{network traffic}}$$

$$x_6 = \frac{(\text{CPU load})^2}{\text{network traffic}}$$

Multivariate Gaussian distribution

Motivating example: Monitoring machines in a data center



Andrew Ng

Multivariate Gaussian (Normal) distribution

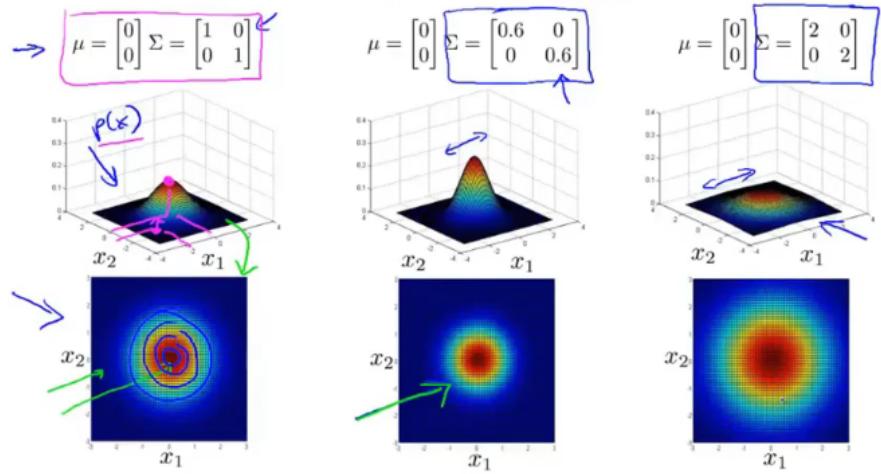
- $x \in \mathbb{R}^n$. Don't model $p(x_1), p(x_2), \dots$, etc. separately.
- Model $p(x)$ all in one go.
- Parameters: $\mu \in \mathbb{R}^n, \Sigma \in \mathbb{R}^{n \times n}$ (covariance matrix)

$$p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} \exp(-\frac{1}{2} (x - \mu)^\top \Sigma^{-1} (x - \mu))$$

$|\Sigma| = \text{determinant of } \Sigma \quad |\det(\Sigma)|$

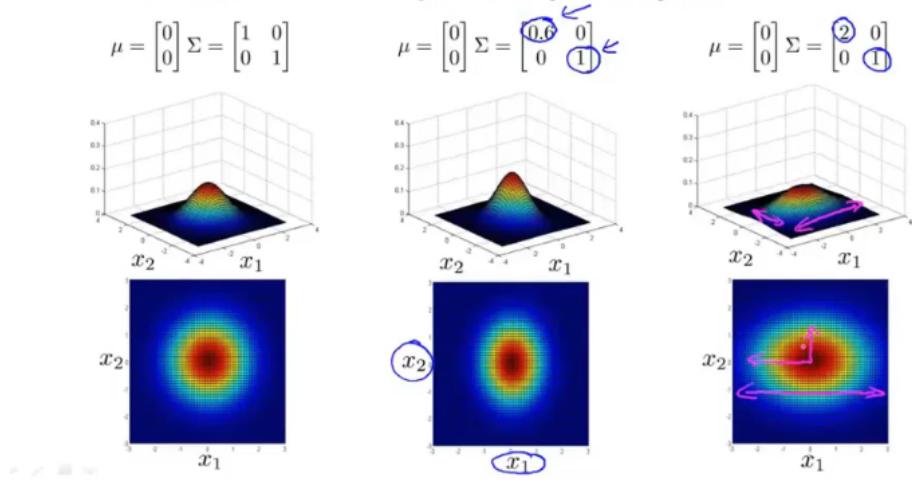
Andrew Ng

Multivariate Gaussian (Normal) examples



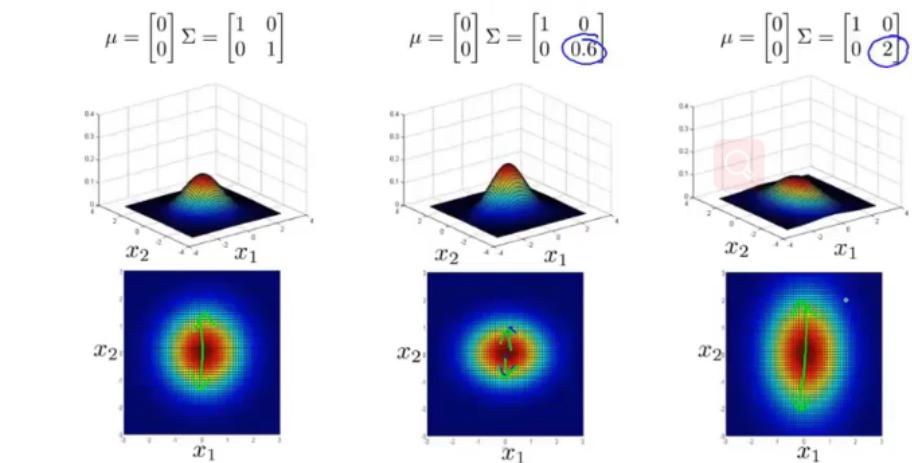
Andrew Ng

Multivariate Gaussian (Normal) examples



Andrew Ng

Multivariate Gaussian (Normal) examples



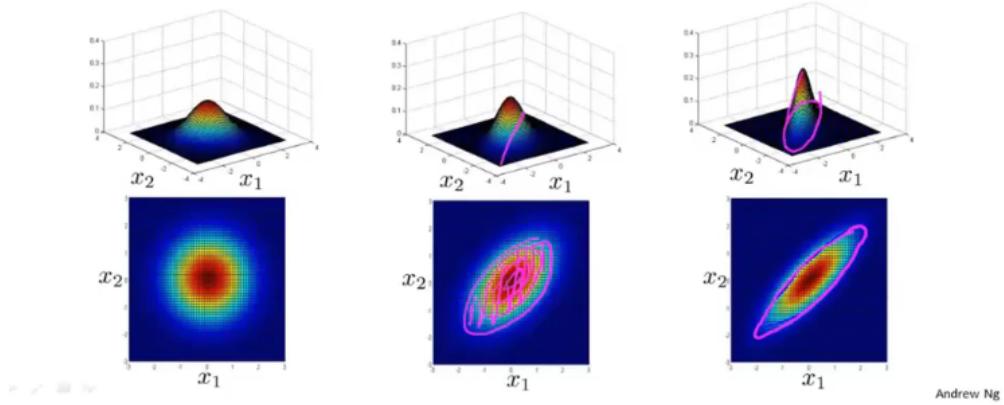
Andrew Ng

Multivariate Gaussian (Normal) examples

$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$$

$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$$

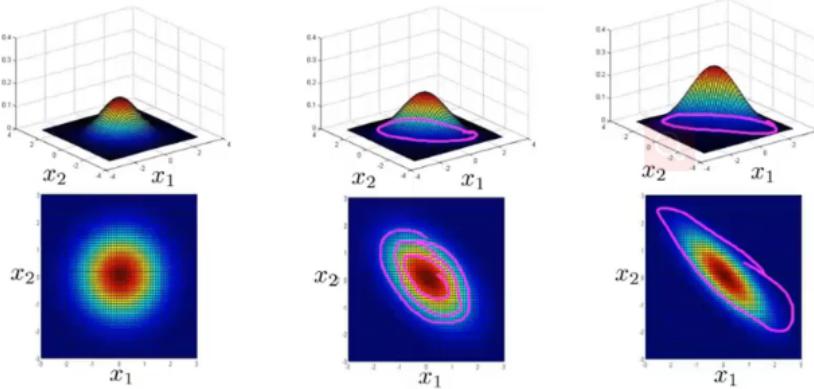


Multivariate Gaussian (Normal) examples

$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ * & 1 \end{bmatrix}$$

$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & -0.5 \\ -0.5 & 1 \end{bmatrix}$$

$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & -0.8 \\ -0.8 & 1 \end{bmatrix}$$

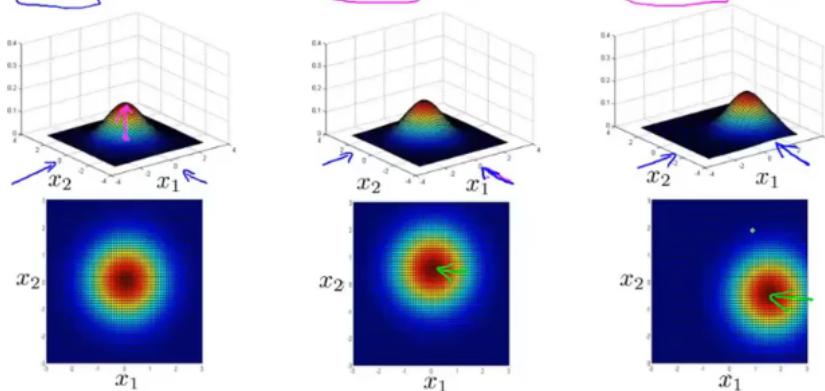


Multivariate Gaussian (Normal) examples

$$\mu = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\mu = \begin{bmatrix} 0 \\ 0.5 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\mu = \begin{bmatrix} 1.5 \\ -0.5 \end{bmatrix} \Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



Anomaly detection using the multivariate Gaussian distribution

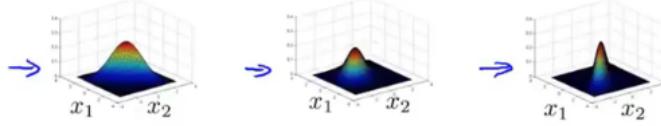
Multivariate Gaussian (Normal) distribution

Parameters μ, Σ

$$\mu \in \mathbb{R}^n$$

$$\Sigma \in \mathbb{R}^{n \times n}$$

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



Parameter fitting:

Given training set $\{x^{(1)}, x^{(2)}, \dots, x^{(m)}\} \leftarrow x \in \mathbb{R}^n$

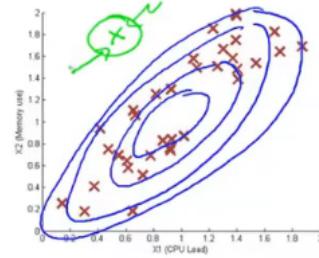
$$\boxed{\mu} = \frac{1}{m} \sum_{i=1}^m x^{(i)} \quad \boxed{\Sigma} = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)(x^{(i)} - \mu)^T$$

Andrew Ng

Anomaly detection with the multivariate Gaussian

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

$$\begin{cases} \mu = \frac{1}{m} \sum_{i=1}^m x^{(i)} \\ \Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu)(x^{(i)} - \mu)^T \end{cases}$$



2. Given a new example x , compute

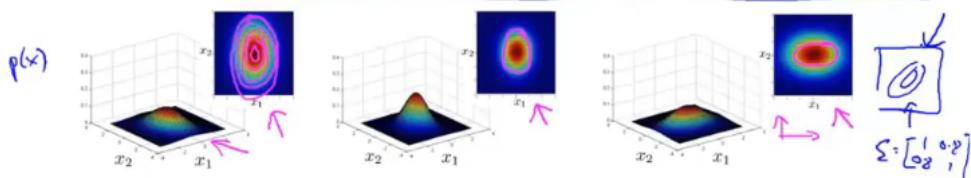
$$p(x) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1} (x - \mu) \right)$$

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

Andrew Ng

Relationship to original model

Original model: $p(x) = p(x_1; \mu_1, \sigma_1^2) \times p(x_2; \mu_2, \sigma_2^2) \times \dots \times p(x_n; \mu_n, \sigma_n^2)$



Corresponds to multivariate Gaussian

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

where

$$\Sigma = \begin{bmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_n^2 \end{bmatrix}$$

Andrew Ng

→ Original model

$$p(x_1; \mu_1, \sigma_1^2) \times \cdots \times p(x_n; \mu_n, \sigma_n^2)$$

Manually create features to capture anomalies where x_1, x_2 take unusual combinations of values.

$$x_3 = \frac{x_1}{x_2} = \frac{\text{CPU load}}{\text{memory}}$$

→ Computationally cheaper (alternatively, scales better to large n) $n=10,000, m=100,000$

OK even if m (training set size) is small

→ Multivariate Gaussian

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

→ Automatically captures correlations between features

$$\Sigma \in \mathbb{R}^{n \times n} \quad \Sigma^{-1}$$

Computationally more expensive

$$\rightarrow \Sigma \sim \frac{n^2}{2}$$

Must have $m > n$ or else Σ is non-invertible.

Andrew Ng

Recommender Systems

Problem formulation

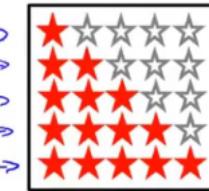
Example: Predicting movie ratings

→ User rates movies using one to five stars

~~zero~~

Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)
Love at last	5	5	0	6
Romance forever	5	4.5	0	0
Cute puppies of love	5	4	0	0
Nonstop car chases	0	0	5	4
Swords vs. karate	0	0	5	4

$$n_u = 4 \quad n_m = 5$$



$\rightarrow n_u = \text{no. users}$
 $\rightarrow n_m = \text{no. movies}$
 $\rightarrow r(i, j) = 1$ if user j has rated movie i
 $\rightarrow y^{(i,j)}$ = rating given by user j to movie i (defined only if $r(i, j) = 1$)

Andrew Ng

Content-based recommendations

Content-based recommender systems

Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)	$\rightarrow n_u = 4, n_m = 5$
Love at last	1	5	5	0	$\rightarrow x_0 = 1$
Romance forever	2	5	?	0	$\downarrow x_1 \text{ (romance)}$
Cute puppies of love	3	4	?	0	$\downarrow x_2 \text{ (action)}$
Nonstop car chases	4	0	5	4	$\rightarrow 0.9 \rightarrow 0$
Swords vs. karate	5	0	5	?	$\rightarrow 1.0 \rightarrow 0.01$

→ For each user j , learn a parameter $\theta^{(j)} \in \mathbb{R}^3$. Predict user j as rating movie i with $(\theta^{(j)})^T x^{(i)}$ stars.

$$x^{(3)} = \begin{bmatrix} 1 \\ 0.99 \\ 0 \end{bmatrix} \Leftrightarrow \theta^{(1)} = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix} \quad (\theta^{(1)})^T x^{(3)} = 5 \times 0.99 = 4.95$$

Andrew Ng

Problem formulation

$\rightarrow r(i, j) = 1$ if user j has rated movie i (0 otherwise)

$\rightarrow y^{(i,j)}$ = rating by user j on movie i (if defined)

$\rightarrow \theta^{(j)}$ = parameter vector for user j

$\rightarrow x^{(i)}$ = feature vector for movie i

\rightarrow For user j , movie i , predicted rating: $\underline{(\theta^{(j)})^T(x^{(i)})}$ $\theta^{(j)} \in \mathbb{R}^{n+1}$

$\rightarrow m^{(j)}$ = no. of movies rated by user j

To learn $\underline{\theta^{(j)}}$:

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

Andrew Ng

Optimization objective:

To learn $\underline{\theta^{(j)}}$ (parameter for user j):

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

To learn $\underline{\theta^{(1)}, \theta^{(2)}, \dots, \theta^{(n_u)}}$:

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

Andrew Ng

Optimization algorithm:

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

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

Gradient descent update:

$$\theta_k^{(j)} := \theta_k^{(j)} - \alpha \sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} \quad (\text{for } k=0)$$

$$\theta_k^{(j)} := \theta_k^{(j)} - \alpha \left(\sum_{i:r(i,j)=1} ((\theta^{(j)})^T x^{(i)} - y^{(i,j)}) x_k^{(i)} + \frac{\lambda \theta_k^{(j)}}{\frac{\partial}{\partial \theta_k^{(j)}} J(\theta^{(1)}, \dots, \theta^{(n_u)})} \right) \quad (\text{for } k \neq 0)$$

Andrew Ng

Collaborative filtering

Problem motivation

Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)	x_1 (romance)	x_2 (action)	$x_3 = 1$
Love at first	5	5	0	0	1.0	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	?	?	?	
	$\theta^{(1)} = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix}$	$\theta^{(2)} = \begin{bmatrix} 0 \\ 5 \\ 0 \end{bmatrix}$	$\theta^{(3)} = \begin{bmatrix} 0 \\ 0 \\ 5 \end{bmatrix}$	$\theta^{(4)} = \begin{bmatrix} 0 \\ 0 \\ 5 \end{bmatrix}$	$\Theta^{(j)}$	$(\Theta^{(1)})^T x^{(1)} \approx 5$ $(\Theta^{(2)})^T x^{(1)} \approx 5$ $(\Theta^{(3)})^T x^{(1)} \approx 0$ $(\Theta^{(4)})^T x^{(1)} \approx 0$	$\times^{(1)}$

Andrew Ng

Optimization algorithm

Given $\theta^{(1)}, \dots, \theta^{(n_u)}$, to learn $x^{(i)}$:

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

Given $\theta^{(1)}, \dots, \theta^{(n_u)}$, to learn $x^{(1)}, \dots, x^{(n_m)}$:

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

Andrew Ng

Collaborative filtering

[Given $x^{(1)}, \dots, x^{(n_m)}$ (and movie ratings), can estimate $\theta^{(1)}, \dots, \theta^{(n_u)}$]

$$r^{(i,j)}$$

 $y^{(i,j)}$

[Given $\theta^{(1)}, \dots, \theta^{(n_u)}$, can estimate $x^{(1)}, \dots, x^{(n_m)}$]

Guess $\Theta \rightarrow x \rightarrow \Theta \rightarrow x \rightarrow \Theta \rightarrow x \rightarrow \dots$

Andrew Ng

Collaborative filtering algorithm

Collaborative filtering optimization objective

Given $x^{(1)}, \dots, x^{(n_m)}$, estimate $\theta^{(1)}, \dots, \theta^{(n_u)}$:

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

Given $\theta^{(1)}, \dots, \theta^{(n_u)}$, estimate $x^{(1)}, \dots, x^{(n_m)}$:

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

Minimizing $x^{(1)}, \dots, x^{(n_m)}$ and $\theta^{(1)}, \dots, \theta^{(n_u)}$ simultaneously:

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

$$\min_{\substack{x^{(1)}, \dots, x^{(n_m)} \\ \theta^{(1)}, \dots, \theta^{(n_u)}}} J(x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)})$$

Andrew Ng

Collaborative filtering algorithm

- 1. Initialize $x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)}$ to small random values.
- 2. Minimize $J(x^{(1)}, \dots, x^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)})$ using gradient descent (or an advanced optimization algorithm). E.g. for every $j = 1, \dots, n_u, i = 1, \dots, n_m$:

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

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

- 3. For a user with parameters θ and a movie with (learned) features x , predict a star rating of $\theta^T x$.

$$(\theta^{(i)})^T (x^{(i)})$$

Andrew Ng

Vectorization: Low rank matrix factorization

Collaborative filtering

Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)
Love at last	5	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	?

$$n_m = 5$$

$$n_u = 4$$

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 \\ 5 & ? & ? & 0 \\ ? & 4 & 0 & ? \\ 0 & 0 & 5 & 4 \\ 0 & 0 & 5 & 0 \end{bmatrix}$$

$$y^{(i,j)}$$

Andrew Ng

Collaborative filtering

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 \\ 5 & ? & ? & 0 \\ ? & 4 & 0 & ? \\ 0 & 0 & 5 & 4 \\ 0 & 0 & 5 & 0 \end{bmatrix}$$

Predicted ratings: $\underline{\underline{X}}(\underline{\underline{\Theta}})^T \leftarrow (\underline{\underline{\Theta}}^{(i)})^T(x^{(i)})$

$$\begin{bmatrix} (\theta^{(1)})^T(x^{(1)}) & (\theta^{(2)})^T(x^{(1)}) & \dots & (\theta^{(n_u)})^T(x^{(1)}) \\ (\theta^{(1)})^T(x^{(2)}) & (\theta^{(2)})^T(x^{(2)}) & \dots & (\theta^{(n_u)})^T(x^{(2)}) \\ \vdots & \vdots & \vdots & \vdots \\ (\theta^{(1)})^T(x^{(n_m)}) & (\theta^{(2)})^T(x^{(n_m)}) & \dots & (\theta^{(n_u)})^T(x^{(n_m)}) \end{bmatrix}$$

$$\underline{\underline{X}} = \begin{bmatrix} -(x^{(1)})^T \\ -(x^{(2)})^T \\ \vdots \\ -(x^{(n_m)})^T \end{bmatrix} \quad \underline{\underline{\Theta}} = \begin{bmatrix} -(\underline{\underline{\Theta}}^{(1)})^T \\ -(\underline{\underline{\Theta}}^{(2)})^T \\ \vdots \\ -(\underline{\underline{\Theta}}^{(n_u)})^T \end{bmatrix}$$

→ Low rank matrix factorization

Andrew Ng

Finding related movies

For each product i , we learn a feature vector $x^{(i)} \in \mathbb{R}^n$.

→ $x_1 = \text{romance}$, $x_2 = \text{action}$, $x_3 = \text{comedy}$, $x_4 = \dots$

How to find movies j related to movie i ?

small $\|x^{(i)} - x^{(j)}\| \rightarrow$ movie j and i are "similar"

5 most similar movies to movie i :

→ Find the 5 movies j with the smallest $\|x^{(i)} - x^{(j)}\|$.

Andrew Ng

Implementational detail: Mean normalization

Users who have not rated any movies

Movie	Alice (1)	Bob (2)	Carol (3)	Dave (4)	Eve (5)
Love at last	5	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	?	?

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 & ? \\ 5 & ? & ? & 0 & ? \\ ? & 4 & 0 & ? & ? \\ 0 & 0 & 5 & 4 & ? \\ 0 & 0 & 5 & 0 & ? \end{bmatrix}$$

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

$$\underline{\underline{\Theta}}^{(s)} \in \mathbb{R}^2 \quad \underline{\underline{\Theta}}^{(s)} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\frac{\lambda}{2} [\underline{\underline{\Theta}}_1^{(s)}]^2 + [\underline{\underline{\Theta}}_2^{(s)}]^2 \leftarrow$$

$$(\underline{\underline{\Theta}}^{(s)})^T \underline{\underline{X}}^{(i)} = 0$$

Andrew Ng

Mean Normalization:

$$Y = \begin{bmatrix} 5 & 5 & 0 & 0 & ? \\ 5 & ? & ? & 0 & ? \\ ? & 4 & 0 & ? & ? \\ 0 & 0 & 5 & 4 & ? \\ 0 & 0 & 5 & 0 & ? \end{bmatrix} \xrightarrow{\text{mean}} \mu = \begin{bmatrix} 2.5 \\ 2.5 \\ 2 \\ 2.25 \\ 1.25 \end{bmatrix} \rightarrow Y = \begin{bmatrix} 2.5 & 2.5 & -2.5 & -2.5 & ? \\ 2.5 & ? & ? & -2.5 & ? \\ ? & 2 & -2 & ? & ? \\ -2.25 & -2.25 & 2.75 & 1.75 & ? \\ -1.25 & -1.25 & 3.75 & -1.25 & ? \end{bmatrix}$$

For user j , on movie i predict:

$$\rightarrow (\Theta^{(s)})^T (X^{(i)}) + \mu_i$$

learn $\underline{\Theta^{(s)}}, \underline{X^{(i)}}$

User 5 (Eve):

$$\underline{\Theta^{(s)}} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$(\Theta^{(s)})^T (X^{(i)}) + \boxed{\mu_i}$$

Andrew Ng

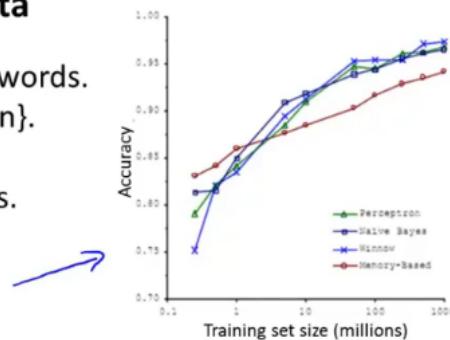
Large scale machine learning

Learning with large datasets

Machine learning and data

Classify between confusable words.
E.g., {to, two, too}, {then, than}.

For breakfast I ate two eggs.

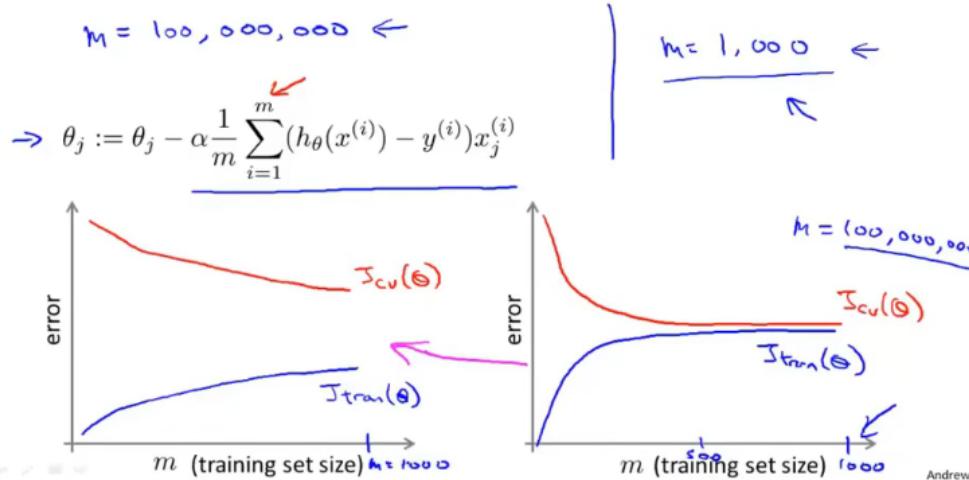


→ “It’s not who has the best algorithm that wins.
It’s who has the most data.”

[Figure from Banko and Brill, 2001]

Andrew Ng

Learning with large datasets



Andrew Ng

Stochastic gradient descent

Linear regression with gradient descent

$$h_{\theta}(x) = \sum_{j=0}^n \theta_j x_j$$

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

Repeat {

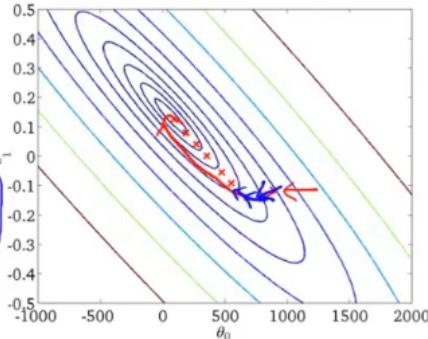
$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}$$

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

}

$$m = 300,000,000$$

Batch gradient descent



Andrew Ng

Batch gradient descent

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

Repeat {

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

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

}

$$m = 300,000,000$$

Stochastic gradient descent

$$\rightarrow \text{cost}(\theta, (x^{(i)}, y^{(i)})) = \frac{1}{2} (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$\rightarrow J_{train}(\theta) = \frac{1}{m} \sum_{i=1}^m \text{cost}(\theta, (x^{(i)}, y^{(i)}))$$

1. Randomly shuffle dataset. ←

2. Repeat {

[for $i := 1, \dots, m$ {

$$\theta_j := \theta_j - \alpha \frac{\partial \text{cost}(\theta, (x^{(i)}, y^{(i)}))}{\partial \theta_j}$$

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

}]

$$\rightarrow (x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), (x^{(3)}, y^{(3)}), \dots$$

Andrew Ng

Stochastic gradient descent

→ 1. Randomly shuffle (reorder) training examples

→ 2. Repeat {

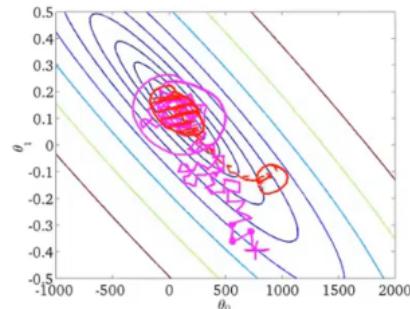
for $i := 1, \dots, m$ {

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

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

}

$$\rightarrow m = 300,000,000$$



Andrew Ng

Mini-batch gradient descent

Mini-batch gradient descent

- Batch gradient descent: Use all m examples in each iteration
- Stochastic gradient descent: Use 1 example in each iteration

Mini-batch gradient descent: Use b examples in each iteration

$$b = \text{mini-batch size.} \quad b = 10. \quad \frac{2 - 100}{2 - 100}$$

Get $b = 10$ examples $(x^{(i)}, y^{(i)}), \dots (x^{(i+9)}, y^{(i+9)})$

$$\Rightarrow \theta_j := \theta_j - \alpha \frac{1}{10} \sum_{k=i}^{i+9} (h_\theta(x^{(k)}) - y^{(k)}) x_j^{(k)}$$

$i := i + 10$

Andrew Ng

Mini-batch gradient descent

Say $b = 10, m = 1000$.

→ b examples

→ 1 example

Repeat {

Vectorization

- for $i = 1, 11, 21, 31, \dots, 991$ {
- $\theta_j := \theta_j - \alpha \frac{1}{10} \sum_{k=i}^{i+9} (h_\theta(x^{(k)}) - y^{(k)}) x_j^{(k)}$
- (for every $j = 0, \dots, n$)
- }
- }

$$m = \underbrace{300,000,000}_{\uparrow}$$

$$b = \underbrace{10}_{\uparrow}$$

Andrew Ng

Stochastic gradient descent convergence

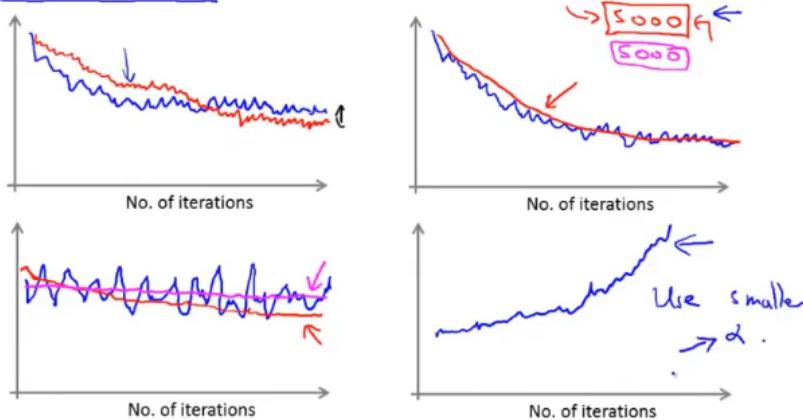
Checking for convergence

- Batch gradient descent:
 - Plot $J_{train}(\theta)$ as a function of the number of iterations of gradient descent.
 - $J_{train}(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$
- Stochastic gradient descent:
 - $cost(\theta, (x^{(i)}, y^{(i)})) = \frac{1}{2} (h_\theta(x^{(i)}) - y^{(i)})^2$
 - During learning, compute $cost(\theta, (x^{(i)}, y^{(i)}))$ before updating θ using $(x^{(i)}, y^{(i)})$.
 - Every 1000 iterations (say), plot $cost(\theta, (x^{(i)}, y^{(i)}))$ averaged over the last 1000 examples processed by algorithm.

Andrew Ng

Checking for convergence

Plot $\text{cost}(\theta, (x^{(i)}, y^{(i)}))$, averaged over the last 1000 (say) examples



Andrew Ng

Stochastic gradient descent

$$\text{cost}(\theta, (x^{(i)}, y^{(i)})) = \frac{1}{2}(h_\theta(x^{(i)}) - y^{(i)})^2$$

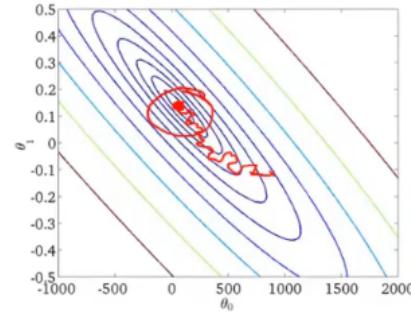
$$J_{\text{train}}(\theta) = \frac{1}{2m} \sum_{i=1}^m \text{cost}(\theta, (x^{(i)}, y^{(i)}))$$

1. Randomly shuffle dataset.

2. Repeat {

```

for i := 1, ..., m {
     $\theta_j := \theta_j - \alpha(h_\theta(x^{(i)}) - y^{(i)})x_j^{(i)}$ 
    (for j = 0, ..., n)
}
}
```



Learning rate α is typically held constant. Can slowly decrease α

over time if we want θ to converge. (E.g. $\alpha = \frac{\text{const1}}{\text{iterationNumber} + \text{const2}}$) $\alpha \rightarrow 0$

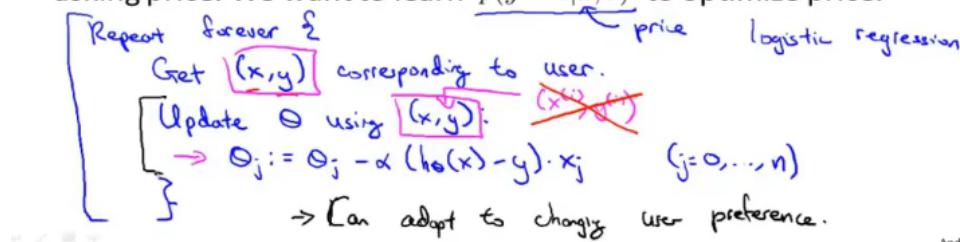
Andrew Ng

Online learning

Online learning

Shipping service website where user comes, specifies origin and destination, you offer to ship their package for some asking price, and users sometimes choose to use your shipping service ($y = 1$), sometimes not ($y = 0$).

Features x capture properties of user, of origin/destination and asking price. We want to learn $p(y = 1|x; \theta)$ to optimize price.



Andrew Ng

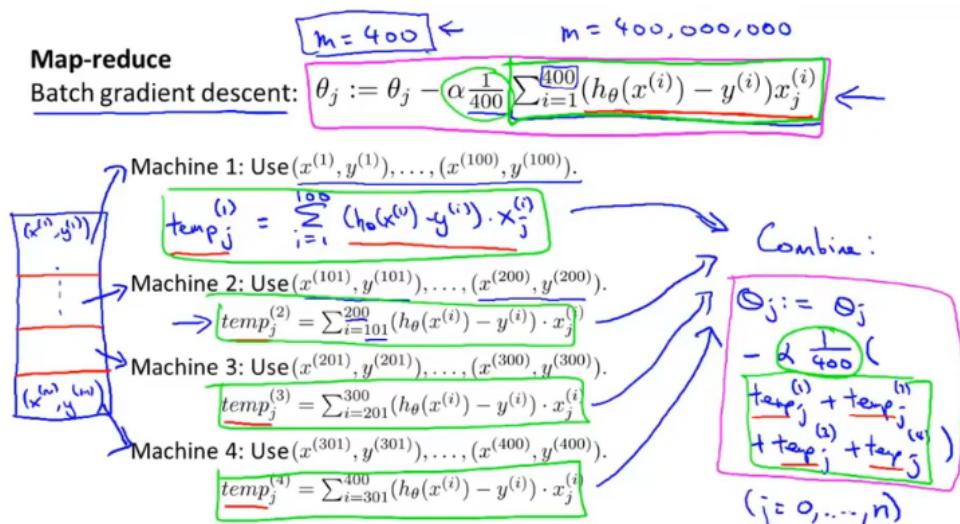
Other online learning example:

Product search (learning to search)

- User searches for "Android phone 1080p camera" ←
- Have 100 phones in store. Will return 10 results.
- $x = \text{features of phone, how many words in user query match name of phone, how many words in query match description of phone, etc.}$
- $y = 1$ if user clicks on link. $y = 0$ otherwise. $(x, y) \leftarrow$
- Learn $p(y = 1|x; \theta)$. ← predicted CTR
- Use to show user the 10 phones they're most likely to click on.
- Other examples: Choosing special offers to show user; customized selection of news articles; product recommendation; ...

Andrew Ng

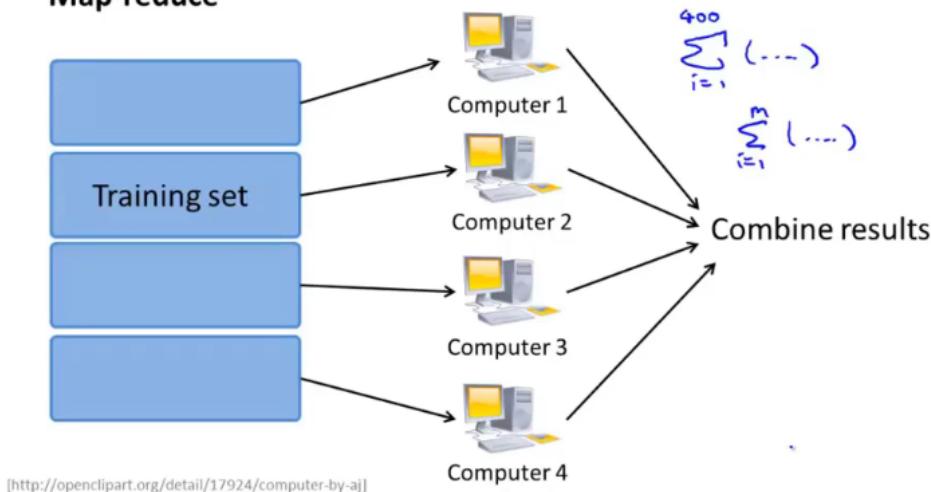
Map-reduce and data parallelism



[Jeffrey Dean and Sanjay Ghemawat] ←

Andrew Ng

Map-reduce



[<http://openclipart.org/detail/17924/computer-by-aj>] ←

Andrew Ng

Map-reduce and summation over the training set

Many learning algorithms can be expressed as computing sums of functions over the training set.

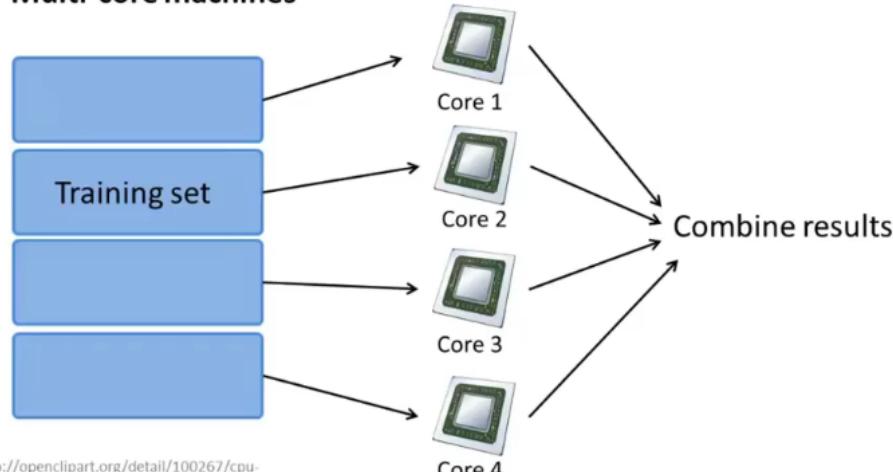
E.g. for advanced optimization, with logistic regression, need:

$$\rightarrow \underline{J_{train}(\theta)} = -\frac{1}{m} \sum_{i=1}^m \underline{y^{(i)} \log h_{\theta}(x^{(i)})} - \underline{(1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)}))}$$
$$\rightarrow \underline{\frac{\partial}{\partial \theta_j} J_{train}(\theta)} = \frac{1}{m} \sum_{i=1}^m \underline{(h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x_j^{(i)}}$$

\uparrow
 $\text{temp}^{(i)}$ $\text{temp}_j^{(i)} \leftarrow$

Andrew Ng

Multi-core machines



Andrew Ng

Application example: Photo OCR

Problem description and pipeline

Photo OCR pipeline

- 1. Text detection



- 2. Character segmentation

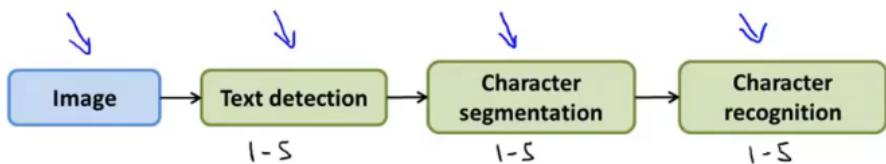


- 3. Character classification



Andrew Ng

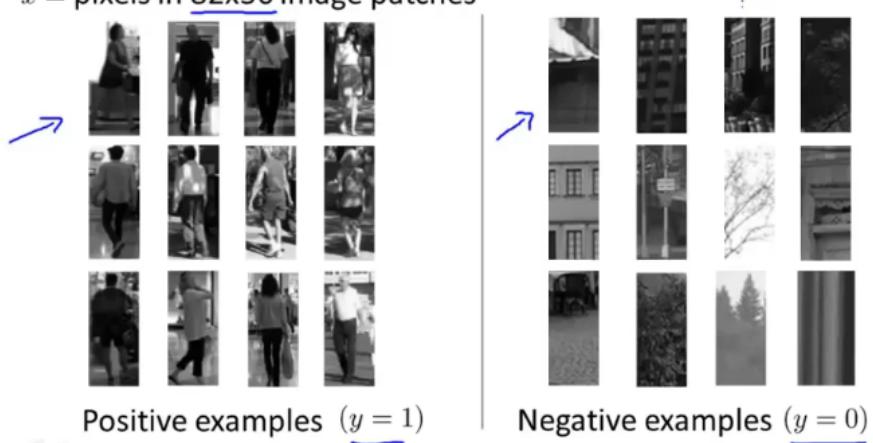
Photo OCR pipeline



Sliding windows

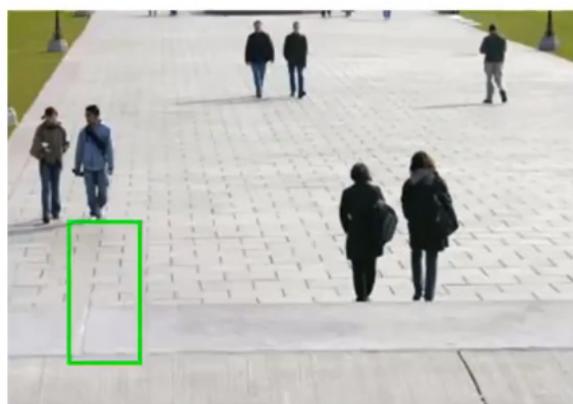
Supervised learning for pedestrian detection

x = pixels in 82x36 image patches



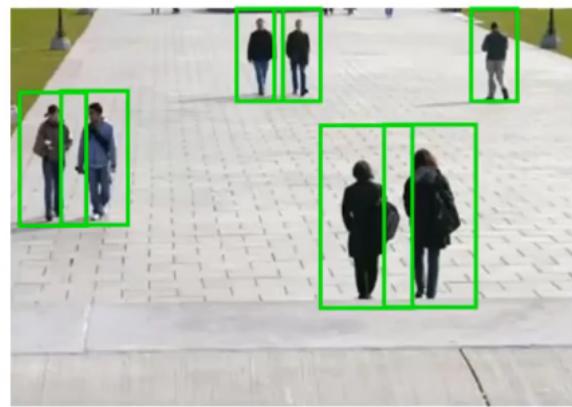
Andrew Ng

Sliding window detection



Andrew Ng

Sliding window detection



Andrew Ng

Text detection

PATRIOTDNF

OUTONICU

Positive examples ($y = 1$)



Negative examples ($y = 0$)

Text detection



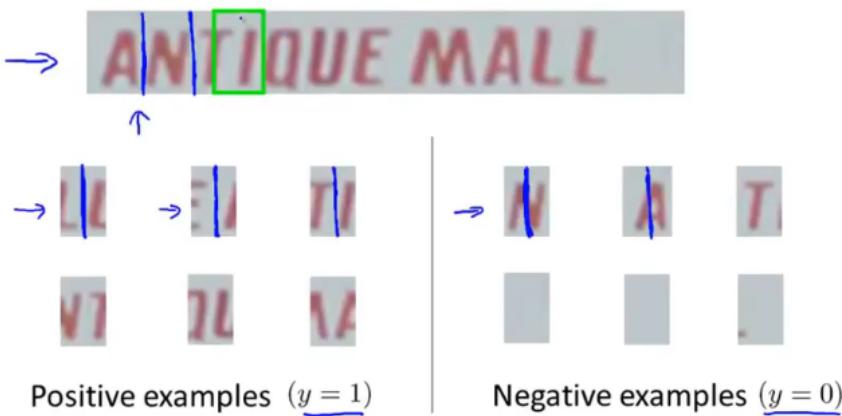
"expansion"



[David Wu]

Andrew Ng

1D Sliding window for character segmentation



Andrew Ng

Photo OCR pipeline

- 1. Text detection



- 2. Character segmentation



- 3. Character classification



Andrew Ng

Getting lots of data: Artificial data synthesis

Artificial data synthesis for photo OCR



Real data

[Adam Coates and Tao Wang]

A
B
c
d
e
f
g
A
b
c
d
e
f
g
A
B
c
d
e
f
g
A
B
c
d
e
f
g

Andrew Ng

Artificial data synthesis for photo OCR



Real data

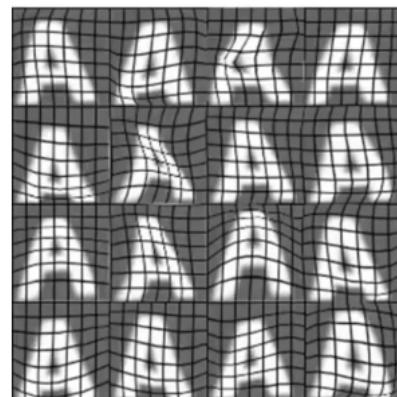
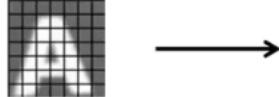


Synthetic data

[Adam Coates and Tao Wang]

Andrew Ng

Synthesizing data by introducing distortions



[Adam Coates and Tao Wang]

Andrew Ng

Synthesizing data by introducing distortions: Speech recognition

🔊 Original audio: ↪

🔊 Audio on bad cellphone connection

🔊 Noisy background: Crowd

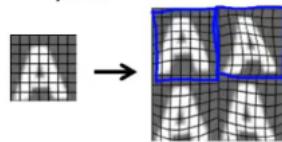
🔊 Noisy background: Machinery

[www.pdsounds.org]

Andrew Ng

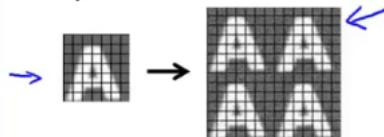
Synthesizing data by introducing distortions

- Distortion introduced should be representation of the type of noise/distortions in the test set.



→ Audio:
Background noise,
bad cellphone connection

- Usually does not help to add purely random/meaningless noise to your data.



→ x_i = intensity (brightness) of pixel i
→ $x_i \leftarrow x_i + \text{random noise}$

[Adam Coates and Tao Wang]

Andrew Ng

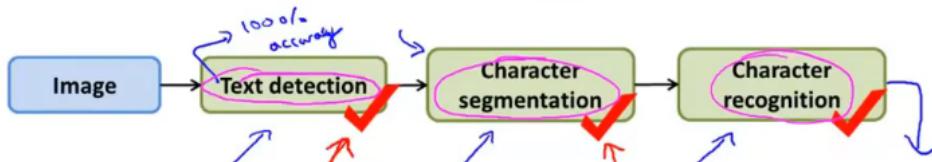
Discussion on getting more data

1. Make sure you have a low bias classifier before expending the effort. (Plot learning curves). E.g. keep increasing the number of features/number of hidden units in neural network until you have a low bias classifier.
2. "How much work would it be to get 10x as much data as we currently have?"
 - Artificial data synthesis
 - Collect/label it yourself
 - "Crowd source" (E.g. Amazon Mechanical Turk)

Andrew Ng

Ceiling analysis: What part of the pipeline to work on next

Estimating the errors due to each component (ceiling analysis)



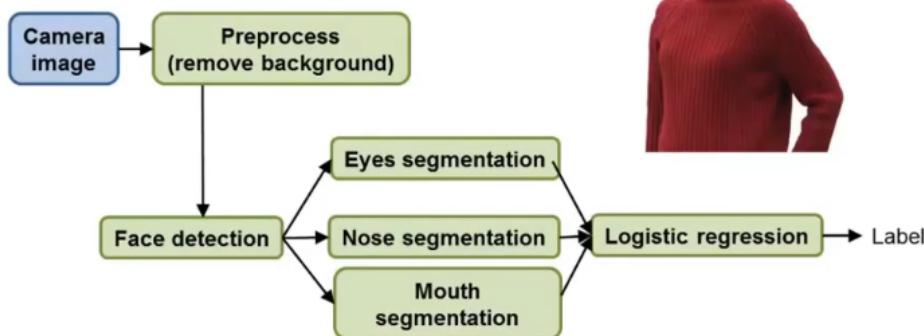
What part of the pipeline should you spend the most time trying to improve?

Component	Accuracy
Overall system	72%
→ Text detection	89% $\downarrow 17\%$
Character segmentation	90% $\downarrow 1\%$
Character recognition	100% $\downarrow 10\%$

Andrew Ng

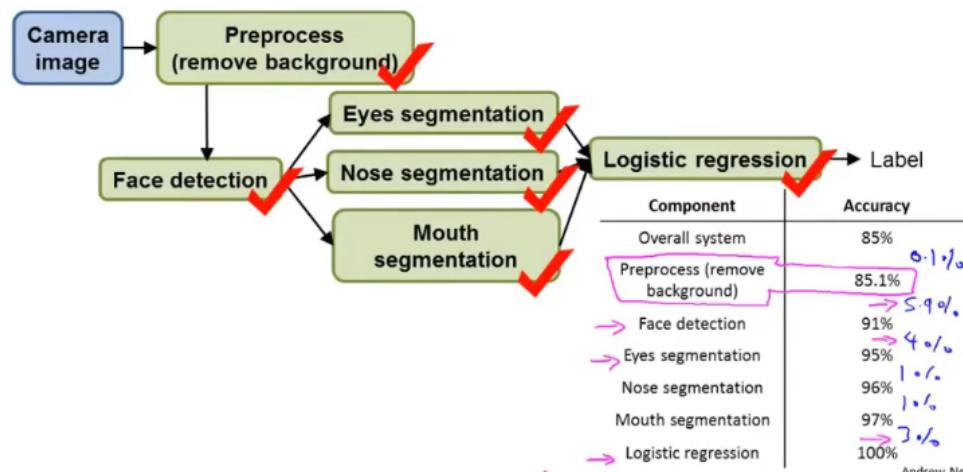
Another ceiling analysis example

Face recognition from images
(Artificial example)



Andrew Ng

Another ceiling analysis example



Andrew Ng

Summary

Summary: Main topics

- [Supervised Learning $(x^{(i)}, y^{(i)})$
- Linear regression, logistic regression, neural networks, SVMs]
- [Unsupervised Learning $x^{(i)}$
- K-means, PCA, Anomaly detection]
- [Special applications/special topics
- Recommender systems, large scale machine learning.]
- [Advice on building a machine learning system
- Bias/variance, regularization; deciding what to work on next: evaluation of learning algorithms, learning curves, error analysis, ceiling analysis.]

Andrew Ng