

WEEK 1

Machine learning definition:

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

— Arthur Samuel

A computer is said to learn from experience E with respect to some class of tasks T and performance P , if its performance at tasks in T , as measured by P , improves with experience E .

— Tom Mitchell

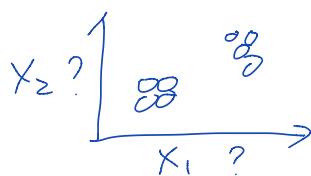
Intro to supervised learning

"right answers" given
categorized as :

{ Regression: Predict continuous valued output
Classification: Discrete valued output

Intro to unsupervised learning

Approach problems with little or no ideas
about what the results should look like

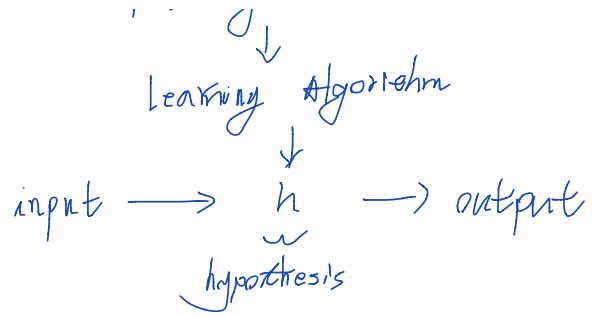


Supervised learning model

Training set

linear regression with one variable:

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



Cost function



θ_0, θ_1 to minimize: $\frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$

m: # training examples

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

y: actual value

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

(Squared error function)

Commonly used for regression problems

Intuition:

hypothesis: $h_\theta(x) = \theta_0 + \theta_1 x$

Parameters: θ_0, θ_1

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

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

Linear regression with one variable — Gradient descent

$$J(\theta_0, \theta_1) \longrightarrow \text{Want } \min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$$

Outline:

start with some θ_0, θ_1

keep changing θ_0, θ_1 to reduce $J(\theta_0, \theta_1)$ until hopefully end up with a minimum

repeat until convergence:

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad \text{for } j=0 \text{ and } j=1 \quad \}$$

α : learning rate

Simultaneously update θ_0, θ_1

Gradient descent for linear regression

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

$$\begin{cases} \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \\ \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)} \end{cases}$$

$$\begin{cases} \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \\ \theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)} \end{cases} \quad \text{update } \theta_0, \theta_1 \text{ simultaneously}$$

Cost function for linear regression will always be a convex function (Bowl-shaped), 1 global/local min

"Batch" Gradient Descent

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

Linear Algebra Review
skip, I know everything

WEEK 2

Multivariate Linear Regression

Notation: n : # of features

$x^{(i)}$: input of i^{th} training example

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

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

define $x_0 = 1$, then $x = \begin{bmatrix} x_0 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1}$

also: $\theta = \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$

then: $h_{\theta}(x) = \theta^T x$

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

Gradient descent.

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

Could be due to learning rate α too large, try use small α

Mathematically, for sufficiently small α , $J(\theta)$ should decrease on every iteration but if α is too small, gradient decent can be slow to converge

Features and Polynomial Regression

$$e.g.: h(\theta) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 = \theta_0 + \theta_1 (\text{size}) + \theta_2 (\text{size})^2 + \theta_3 (\text{size})^3$$

$$\hookrightarrow x_1 = (\text{size}) \quad x_2 = (\text{size})^2 \quad x_3 = (\text{size})^3$$

One should take care of feature scaling, since $x_1 \sim x_3$ takes
every different ranges

Normal Equation:

Normal equation: method solve for θ analytically

$$\frac{\partial J(\theta)}{\partial \theta_j} = 0 \Rightarrow \theta_0 \dots \theta_n$$

x_0	feature 1 x_1	feature 2 x_2	feature 3 x_3	feature 4 x_4	results y
1
1
1
1
1

$$X = \begin{bmatrix} 1 & 1 & \dots & 1 & 1 \\ 1 & 1 & \dots & 1 & 1 \\ 1 & 1 & \dots & 1 & 1 \\ 1 & 1 & \dots & 1 & 1 \\ 1 & 1 & \dots & 1 & 1 \end{bmatrix}_{m \times (n+1)}$$
$$y = \begin{bmatrix} y \\ y \\ y \\ y \\ y \end{bmatrix}_{m \times 1}$$

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

Generally:

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

n features

$$x^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$X = \begin{bmatrix} \cdots & (x^{(1)})^T & \cdots \\ \cdots & \cdots & \cdots \\ \cdots & (x^{(m)})^T & \cdots \end{bmatrix}_{m \times (n+1)}$$

↑
design matrix

It Does NOT matter whether using feature scaling or not in this case

Compare:

Gradient descent

- need to choose α
- needs many iterations
- works well even feature n is longer
- more general

Normal equation

- no need choose α
- No need of iteration
- need to compute $[X^T X]^{-1}$
(slow when n is large)
- linear regression only
- $[X^T X]^{-1}$ may non-invertible

What does $[X^T X]$ is non-invertible mean?

- Redundant features
- Too many features ($m \leq n$)

Vectorization:

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

WEEK 3

Logistic Regression

Regularization

WEEK 3

Classification

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

Logistic regression: $0 \leq h_\theta(x) \leq 1$

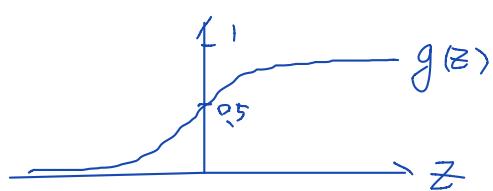
Logistic regression — hypothesis representation

model: want $0 \leq h_\theta(x) \leq 1$

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

$$= \frac{1}{1 + e^{-z}} \leftarrow \text{Si } i \quad \circ \text{ is } -1 \quad - \text{not } 0$$

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



for binary classification: $h_{\theta}(x) = P(y=1|x; \theta)$
 $P(y=0|x; \theta) + P(y=1|x; \theta) = 1$

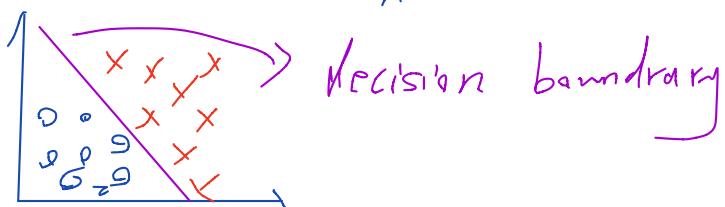
Logistic regression — decision boundary

$$h_{\theta}(x) = P(y=1|\theta; x) \begin{cases} \text{predict } y=1 \text{ if } h_{\theta}(x) \geq 0.5 \\ \dots \quad y=0 \text{ if } h_{\theta}(x) < 0.5 \end{cases}$$

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

$$g(z) \geq 0.5 \text{ when } z \geq 0 \quad \hookrightarrow \theta^T x \geq 0$$

$$g(z) < 0.5 \text{ when } z < 0 \quad \hookrightarrow \theta^T x < 0$$



Logistic regression — cost function

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

$$x = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} \quad x_0 = 1, \quad y \in \{0, 1\}$$

m train sets, n features

$$h_{\theta}(x) = \frac{1}{1+e^{-\theta^T x}} \quad \text{how to choose parameters } \theta ?$$

Cost function

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

* if we use the above same cost function for Logistic regression,
the cost function will be "non-convex" (has many local minimums)



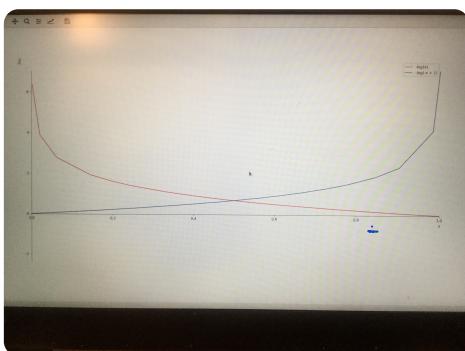
which is not guaranteed to converge to global minimum

Thus we have to modify the cost function for logistic regression

definition:

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_\theta(x_i), y_i)$$

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



Plot of cost

Simplified Cost function and gradient descent

Logistic regression cost function: $J_\theta = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_\theta(x^{(i)}), y^{(i)})$

Simplification:

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

then:

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

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

To fit parameters θ : $\min_{\theta} J(\theta)$

To make a prediction giving new x : $h(\theta) = \frac{1}{1+e^{-\theta^T x}}$

Gradient descent

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

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

then:

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)} \quad [\text{Sim update all } \theta_j]$$

Notice Ng did put $\frac{1}{m}$ here, which I think is a mistake

Looks exactly the same as linear regression!

Except the definition of hypothesis $h_\theta(x)$ is diff.

Feature scaling also works for logistic regression!

Vectorized implementation:

$$\theta := \theta - \frac{\alpha}{m} X^T \cdot (g(X^T \theta) - \vec{y})$$

Logistic regression — adv. optimization

given θ , we have code that Compute:

- $J(\theta)$
- $-\frac{\partial}{\partial \theta_j} J(\theta)$ for $j = 0, 1, \dots, n$

optimization algorithms

- Gradient decent

- Conjugate gradient
- BFGS
- L-BFGS

advantages:

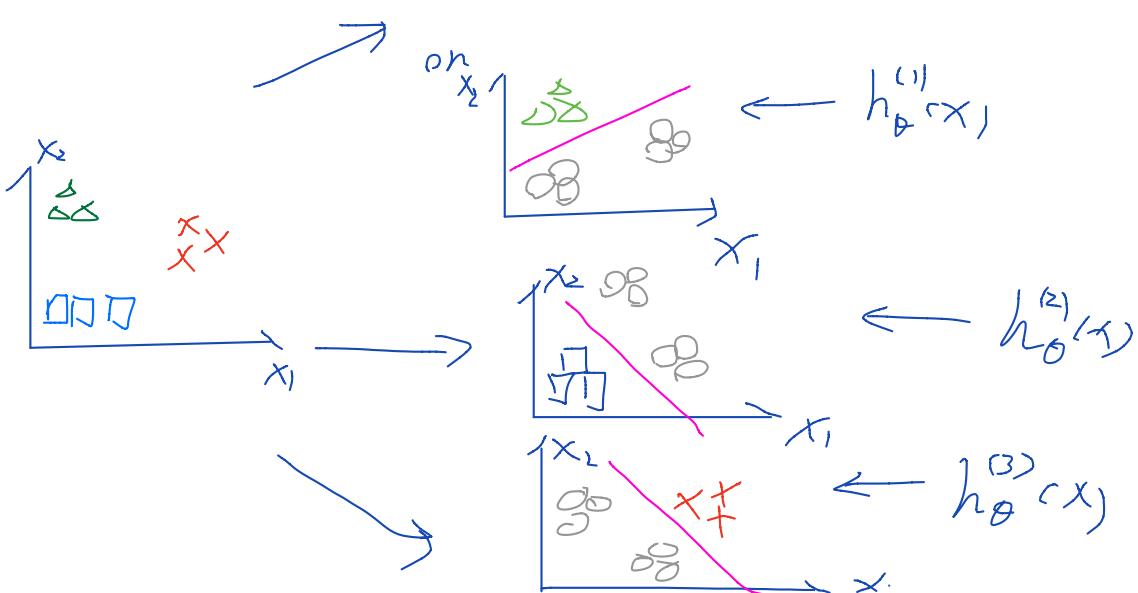
- No need to manually pick α
- often faster than gradient descent

Disadvantages:

- more complex

Multi-class classification

one-vs-all (rest) classification



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

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 output x , to make a prediction, pick the class i that maximizes: $\max_i h_\theta^{(i)}(x)$

Problem of overfitting

Overfitting:

If we have too many features, the learned hypothesis may fit the training set very well ($J(\theta) \approx 0$), but fail to generalize to new example

Addressing overfitting

Options:

1. Reduce number of features

a. Manually select which features to keep

b. Model selection algorithm (later)

2. Regularization

a. keep all features, but reduce magnitude / values of θ_j

↳ works well when having a lot of features, each of which contributes a bit to prediction of

Regularization — cost function

Regularization:

Small values for parameters $\theta_0, \dots, \theta_n$

↪ "Simpler" hypothesis

↪ less prone to over fitting

regularized Cost function

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

omitted θ_0

convention

and θ_0 doesn't make any diff.

regularization parameter

λ : too larger \rightarrow underfit

Regularized linear regression:

① gradient decent:

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) X_0^{(i)}$$

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

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

$$\theta_j := \underbrace{\theta_j}_{\text{---}} \left(1 - \alpha \frac{\lambda}{m} \right) - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) X_j^{(i)}$$

usually α, λ small, m large, so $1 - \alpha \frac{\lambda}{m} < 1$

② Normal Equation

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

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

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

$$\theta = (X^\top X + \lambda I)^{-1} X^\top y$$

Non-invertibility

if $m \leq n$, $[X^T X]$ non-invertible

if $\lambda > 0$, then $[X^T X + \lambda [1, 1]]^{-1}$ is guaranteed to be invertible

Regularized logistic regression

regularize

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

Gradient descent:

Same as linear regression, except that hypothesis func $h_\theta(x)$ is different.

WEEK 4

- Neural Networks: Representation

WEEK 4 Neural Networks

Neural networks motivation

- Complex regression
- Large number of features

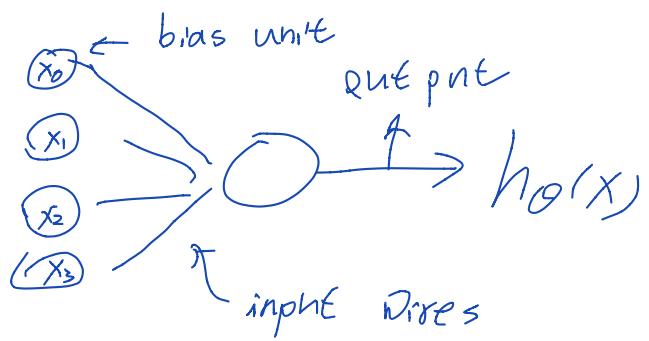
Origins: Algorithms that try to mimic the brain

Somatosensory Cortex "one learning hypothesis"

Neural Networks representation

Neuron model: Logistic unit

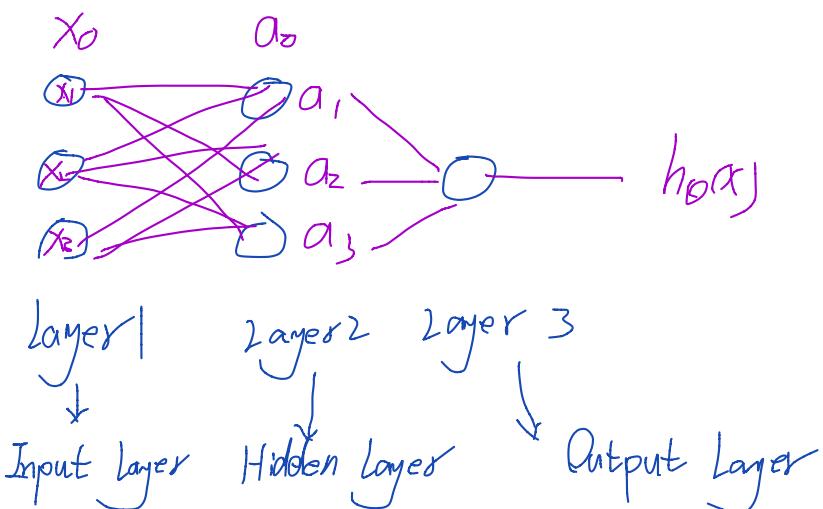
terminology:



Sigmoid (logistic) activation function $g(z) = \frac{1}{1+e^{-z}}$

$$X = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad \Theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix} \leftarrow \text{"weights" = parameters}$$

Neural Network:



$a_i^{(j)}$: "activation of unit i in layer j

$$a_1^{(1)} = g(\Theta_{1,0}^{(1)} x_0 + \Theta_{1,1}^{(1)} x_1 + \Theta_{1,2}^{(1)} x_2 + \Theta_{1,3}^{(1)} x_3)$$

$$a_2^{(1)} = g(\Theta_{2,0}^{(1)} \dots \Theta_{2,1}^{(1)} + \Theta_{2,2}^{(1)} \dots \Theta_{2,3}^{(1)})$$

$$a_3^{(1)} = \dots$$

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

$$z_1^{(2)}$$

$$z_2^{(2)}$$

$$z_3^{(2)}$$

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

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

Forward Propagation: Vectorized Implementation

Vectorize:

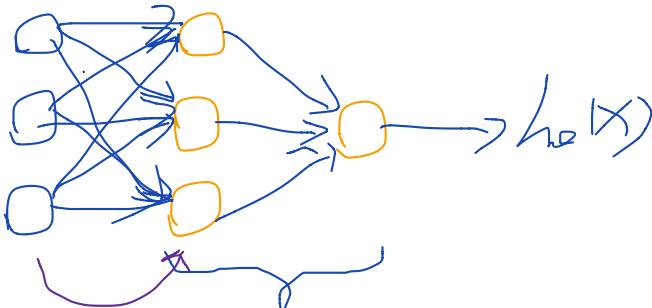
$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix} \quad a^{(2)} = \begin{bmatrix} a_1^{(2)} \\ a_2^{(2)} \\ a_3^{(2)} \end{bmatrix}$$

$$z^{(2)} = \Theta^{(2)} x \quad a^{(2)} = g(z^{(2)})$$

also define $a^{(1)} = x$, add $a_0^{(2)} = 1$

$$z^{(3)} = \Theta^{(3)} a^{(2)} \quad h_{\Theta}(x) = a^{(3)} = g(z^{(3)})$$

Neural Network learning its own features

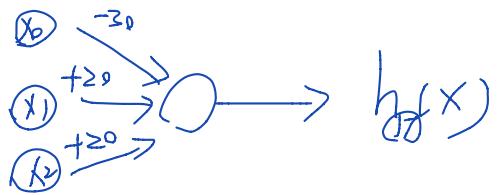


$\Theta^{(2)}$ Logistic regression $h_{\Theta}(x) = g(\Theta_{10}^{(2)} a_0^{(2)} + \Theta_{11}^{(2)} a_1^{(2)} + \Theta_{12}^{(2)} a_2^{(2)} + \Theta_{13}^{(2)} a_3^{(2)})$

Neural network learning its own features

Examples

(1) Logic AND function

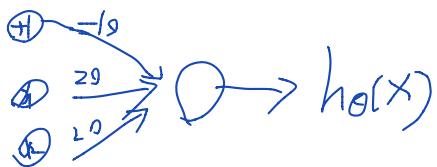


$$h_{\theta}(x) = g(-3x_0 + 2x_1 + 2x_2)$$

g : Sigmoid function

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

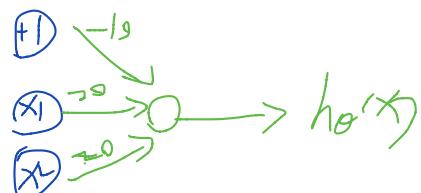
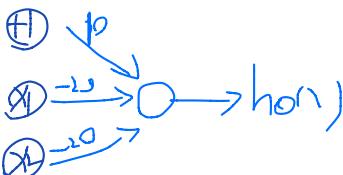
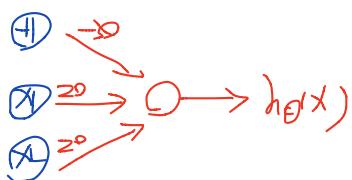
(2) Logic OR function



(3) Logic NOT

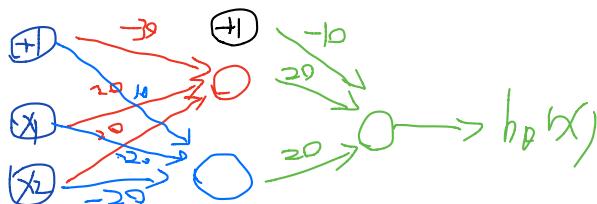


(4) XNOR



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

$x_1 \text{ OR } x_2$



XNOR !!!

Multi-class classification

$$\text{output} = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \dots$$

Neural Network classification problem

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

L : total no. of layers in the network

s_l : no. of units (exclude bias unit) in layer l

binary classification: $y = 0 \text{ or } 1$ | output unit
 $\begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \dots, \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}$

Multi-class classification (K classes) $\begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \dots, \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}$ K output units

logistic regression

$$J(\Theta) = -\frac{1}{m} \left[\sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_\Theta(x^{(i)}))_k + (1-y_k^{(i)}) \log(1-h_\Theta(x^{(i)}))_k \right] - \frac{1}{m} \sum_{i=1}^m \sum_{j=1}^{s_i} \sum_{k=1}^{s_{i+1}} (h_\Theta^{(i)}(x^{(i)})_{jk})^2$$

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

Neural Networks: Backpropagation algorithm

Gradient Computation: need $J(\Theta)$ $\frac{\partial J(\Theta)}{\partial \Theta_{ij}^k}$

e.g. 1 training example, 4 layers, forward propagation:

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

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

$$a^{(2)} = g(z^{(2)}) \quad (\text{add } a_0^{(2)})$$

$$z^{(3)} = \Theta^{(2)} a^{(2)} \quad a^{(3)} = g(z^{(3)}) \quad (\text{add } a_0^{(3)})$$

$$z^{(4)} = \Theta^{(3)} a^{(3)} \quad a^{(4)} = h_\Theta(x) = g(z^{(4)})$$

Backpropagation algorithm

intuition: $\delta_j^{(l)} = \text{"error" of node } j \text{ on layer } l$

$$\text{e.g. for layer 4: } \delta_j^{(4)} = a_j^{(4)} - y_j$$

$$\delta_j^{(3)} = (\Theta^{(3)})^T \delta_i^{(4)} * g'(z^{(3)}) \quad a_i^{(3)} * (1 - a_i^{(3)})$$

$$\delta_j^{(2)} = (\Theta^{(2)})^T \delta_i^{(3)} * g'(z^{(2)}) \quad a_i^{(2)} * (1 - a_i^{(2)})$$

no $\delta^{(1)}$ term

$$\text{Mathematically: } \frac{\partial J(\Theta)}{\partial \theta_{ij}^{(l)}} = a_j^{(l)} \delta_i^{(l+1)} \quad (\text{ignoring } \lambda \text{ if } \lambda=0)$$

Backpropagation Algorithm

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 :

↳ set $a^{(1)} = X^{(i)}$ (input layer)

↳ perform forward propagation to compute $a^{(l)}$ for $l = 2, 3, \dots, L$

↳ using $y^{(i)}$, compute $\delta^{(L)} = a^{(L)} - y^{(i)}$

↳ Compute $\delta^{(L-1)}, \delta^{(L-2)}, \dots, \delta^{(2)}$

↳ $\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$

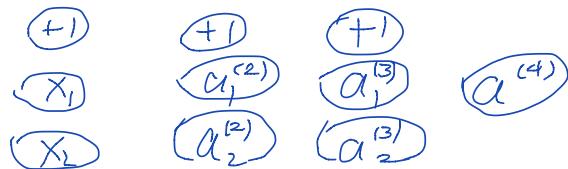
$\rightarrow D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \Theta_{ij}^{(l)} \quad \text{if } j \neq 0$

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

\rightarrow use for gradient descent or other algorithms

Backpropagation intuition

forward propagation



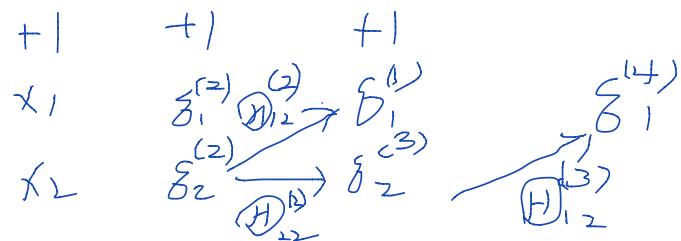
$$z_1^{(3)} = \Theta_{10}^{(4)} x_1 + \Theta_{11}^{(4)} a_1^{(2)} + \Theta_{12}^{(4)} a_2^{(2)}$$

backpropagation

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

$\delta_j^{(l)}$ = "error" of cost for $a_j^{(l)}$ (unit j in layer l)

$$\text{formally: } \delta_j^{(l)} = \frac{\partial}{\partial z_j} \text{cost}(1)$$



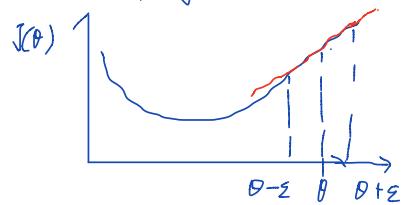
$$\delta^{(4)} = y^{(4)} - a_1^{(4)}$$

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

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

Gradient Checking

Numerical estimation of gradients



$$\frac{d}{d\theta} J(\theta) \approx \frac{J(\theta + \epsilon) - J(\theta - \epsilon)}{2\epsilon}$$

$$\epsilon \sim 10^{-4}$$

Parameter vector θ

$\theta \in \mathbb{R}^n$ (unrolled version of $\theta^{(1)}, \theta^{(2)}, \theta^{(3)}$)

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

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

$$\frac{\partial}{\partial \theta_2} J(\theta) \approx \dots \quad \dots \quad \dots$$

⋮

$$\frac{\partial}{\partial \theta_n} J(\theta) \approx \dots \quad \dots \quad \dots$$

check that above approx $\approx \frac{\partial}{\partial \theta_j} J(\theta)$ from back propagation

Implementation Note

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

Random Initialization

Initial value of θ

- ★. Initialize all parameters to 0 when training a

Neural network is NOT OK !

(after each update, parameters corresponding to inputs going into each of two hidden units are identical
this prevents your neural network from learning anything interesting)

Random Initialization: Symmetry breaking

initialize each Θ_{ij}^w to a random value in $[-\varepsilon, \varepsilon]$

Put it together.

1) Pick a network architecture

No. of input units: Dimension of features $X^{(1)}$

No. of output units: number of classes

No. of hidden layers:

reasonable default: 1 hidden layer

or if > 1 hidden layer, have same number of hidden units in each layer (Usually the more -the better)

2) Randomly initialize weights Θ

3) Implement forward propagation to get $h_\Theta(x^{(1)})$ for any $x^{(1)}$

4) Implement code to compute cost function $J(\Theta)$

5) Implement backprop to compute partial derivatives $\frac{\partial}{\partial \Theta_k} J(\Theta)$

for $i = 1 : m$ $(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$

perform forward prop and backprop using example

$(x^{(1)}, y^{(1)})$, get activations $a^{(1)}$ and delta $\delta^{(1)}$

for $l = 2, \dots, L$

- 6) Use gradient checking to compare $\frac{\partial^2 J(\theta)}{\partial \theta_j \partial \theta_k}$ from backprop and numerical estimation, then disable the gradient checking
- 7) Use gradient descent or adv. method with backprop to try to minimize $J(\theta)$ as a function of parameters θ

decide what to try next

what to do if getting large errors in prediction?

- Get more training examples
- Try smaller sets of features
- Try getting additional features
- Try adding polynomial features ($x_1^2, x_2^2, x_1^3, x_2^3, \dots$)
- Try decreasing/increasing λ

Machine Learning diagnostic

diagnostic: A test 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

Evaluate a hypothesis

Split training data into training set and test set to avoid overfitting

Training / testing procedure for linear regression

- Learn parameter θ from training set (70% of training data, e.g.)
- Compute test set error

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

(logistic regression is similar)

- Misclassification error.

$$\text{err}(h_{\theta}(x), y) = \begin{cases} 1, & \text{if } h_{\theta}(x) \geq 0.5, y=0 \\ & \text{if } h_{\theta}(x) \leq 0.5, y=1 \\ 0, & \text{else} \end{cases}$$

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

Model selection and training/validation/test sets

Model selection:

choose right degree of polynomial:

$$1. h_{\theta}(x) = \theta_0 + \theta_1 x \rightarrow \mathbb{H}^{(1)}$$

$$2. h_{\theta}(x) = \theta_0 + \theta_1 x + \theta_2 x^2 \rightarrow \mathbb{H}^{(2)}$$

$$\vdots \quad | \quad \vdots$$

$$d. h_{\theta}(x) = \theta_0 + \theta_1 x + \dots + \theta_d x^d \rightarrow \mathbb{H}^{(d)}$$

d : degree of polynomial

choose d to minimize the error $J(\theta)$ in test set is NOT right, Since d might just be optimized for this specific test set.,

usually choose:

training data $\left\{ \begin{array}{l} 60\% \text{ Training set } (x^{(i)}, y^{(i)}) \\ 20\% \text{ cross validation set } (CV) \quad (x_{cv}^{(i)}, y_{cv}^{(i)}) \\ 20\% \text{ test set } (x_{test}^{(i)}, y_{test}^{(i)}) \end{array} \right.$

Errors: $J_{train}(\theta)$, $J_{cv}(\theta)$, $J_{test}(\theta)$

choose the degree of polynomial d to get minimum $J_{cv}(\theta)$
then use the θ from same d to calculate $J_{test}(\theta)$
which gives right test set error

Bias vs. Variance

High bias: under fit

High variance: over fit

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

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



high bias ↑ high variance → degree of polynomial
optimal value

High Bias (underfit): $J_{\text{train}}(\theta)$ will be high, $J_{\text{cv}}(\theta) \approx J_{\text{train}}(\theta)$

High variance (over fit): $J_{\text{train}}(\theta)$ will be low, $J_{\text{cv}}(\theta) \gg J_{\text{train}}(\theta)$

Regularization and bias/variance

Linear regression with regularization

$$\begin{aligned} \text{Model: } h_{\theta}(x) &= \theta_0 + \theta_1 x_1 + \theta_2 x_2^2 + \theta_3 x_3^3 + \theta_4 x_4^4 \\ J(\theta) &= \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 + \frac{\lambda}{m} \sum_{j=1}^m \theta_j^2 \end{aligned}$$

Large λ : $\lambda = 1000, \theta_1 \approx 0, \theta_2 \approx 0, \dots, \theta_4 \approx 0 \quad h_{\theta}(x) = \theta_0 \Rightarrow \text{high bias}$

Small λ : $\lambda = 0 \Rightarrow \text{high variance (over fit)}$

Intermediate $\lambda \Rightarrow \text{Just right}$

Choosing the regularization parameter λ

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

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

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

$$J_{\text{cv}}(\theta) = \frac{1}{m_{\text{cv}}} \sum_{i=1}^{m_{\text{cv}}} (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

no regularization

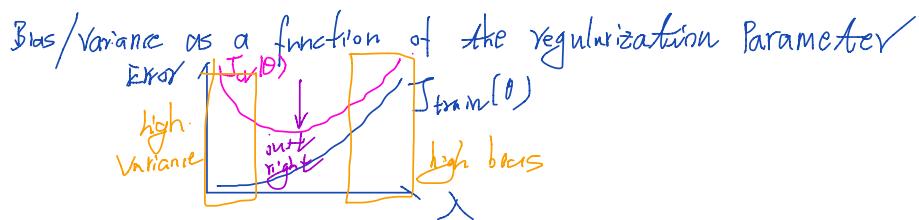
$$J_{\text{test}}(\theta) = \frac{1}{m_{\text{test}}} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

J term involved

1. try $\lambda = 0 \rightarrow \min_{\theta} J(\theta) \rightarrow \theta^{(0)} \rightarrow J_{\text{cv}}(\theta^{(0)})$
2. try $\lambda = 0.01 \rightarrow \min_{\theta} J(\theta) \rightarrow \theta^{(1)} \rightarrow J_{\text{cv}}(\theta^{(1)})$
3. try $\lambda = 0.02 \quad | \quad | \quad |$
4. try $\lambda = 0.04 \quad | \quad | \quad |$
5. try $\lambda = 0.08 \quad | \quad | \quad |$
- ⋮
12. try $\lambda = 10 \quad | \quad | \quad |$

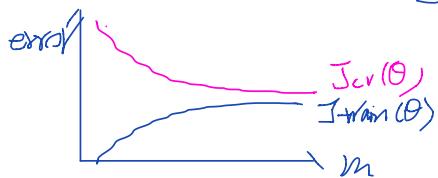
choose λ corresponding to smallest $J_{\text{cv}}(\theta)$, this is the proper value of λ

To see the error of this hypothesis, calculate $J_{\text{test}}(\theta)$

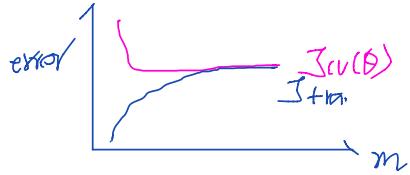


Learning Curves

Plot $J_{\text{train}}(\theta)$ and $J_{\text{cv}}(\theta)$ V.S. training set size m



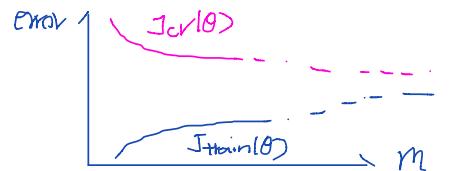
High bias



J_{cv} saturates pretty soon, $J_{train} \approx J_{cv}(\theta)$ by the end

* if a learning algorithm is suffering from high bias,
getting more training data will not help much

High variance



* if a learning algorithm is suffering from high variance (overfitting)
, getting more training data is likely to help

Deciding what to try next - revisit:

debugging a learning algorithm

Suppose you implemented regularized linear regression, however when you apply this hypothesis on a new test set, you find it makes unacceptable large errors. What should you try next?

- Get more training examples
 - Try smaller sets of features
 - Try increasing λ
 - Try adding additional features
- $\left. \begin{array}{l} \\ \\ \end{array} \right\} \text{fix high variance (overfitting) problems}$

- Try adding polynomial features
 - Try decreasing λ
- } fix high bias (underfitting) problems

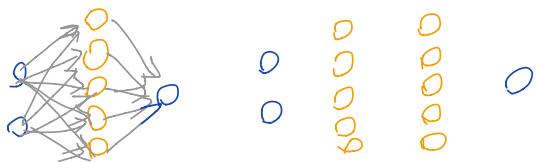
Neural Networks and overfitting

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

PREFERRED !!

Machine Learning System Design

Prioritizing what to work on: Spam classification example

Error Analysis

Recommended Approach

- Start with a simple algorithm that you can implement quick and dirty, 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 CV data set) where your algorithm made errors on, see if you spot any systematic trend in what type of examples it is making errors on.

Error metrics for skewed classes

Cancer classification example

Train logistic regression model $h(x)$, find that you got 1% error on test set (99% correct diagnostics)

BUT if in the training set only 0.5% of people have cancer, prediction $y=0$ (always) will give only 0.5% error, which seems better than "learned" predictions.



Skewed classes

Problem: 99.5% accuracy is better than 99% ??

Precision / Recall.

		<u>Actual class</u>	
		1	0
predicted class	1	True Positive	False Positive
	0	False negative	True negative

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

Precision: of all patients where we predicted $y=1$, what fraction actually has cancer?)

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

Recall: of all patients that actually have cancer, what fraction did we correctly detect as having cancer?

$$\frac{\text{True Pos}}{\# \text{ actual pos}} = \frac{\text{True Pos}}{\text{True pos} + \text{false neg}}$$

Trading off precision and recall

Normally:

Logistical regression: $0 \leq h_0(x) \leq 1$

predict 1 if: $h_0(x) \geq 0.5$

predict 0 if $h_0(x) < 0.5$

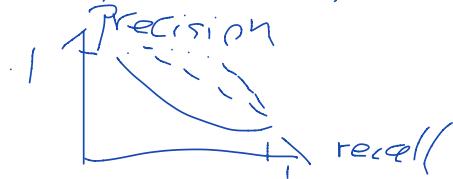
Suppose we want to predict $y=1$ only if very confident

→ higher precision, lower recall: $0.5 \rightarrow 0.7$

Suppose we want to avoid missing too many cases of cancer

→ higher recall, lower precision $0.5 \rightarrow 0.3$

Vary value of threshold \Rightarrow trade off precision vs recall



F₁ Score (on CV-set)

How to compare precision/^(P) recall numbers

Average: $\frac{P+R}{2}$ → not good

F₁ Score: $2 \frac{PR}{P+R}$ $P=0 \text{ or } R=0 \Rightarrow F_1=0$
 $P=1 \text{ & } R=1 \Rightarrow F_1=1$

Data for machine Learning

useful test: given the input x , can a human expert confidently predict y ? (x is features good enough?)

Large data rationale

- └ algorithm with many features \rightarrow low bias algorithms
- └ large training set \rightarrow unlikely to overfit
- └ $J_{\text{train}}(\theta)$ small
- └ $J_{\text{train}}(\theta) \approx J_{\text{test}}(\theta)$

/ - - - - - /

Combine above two:

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

\Rightarrow have many features + Large training set

Large Margin Classification

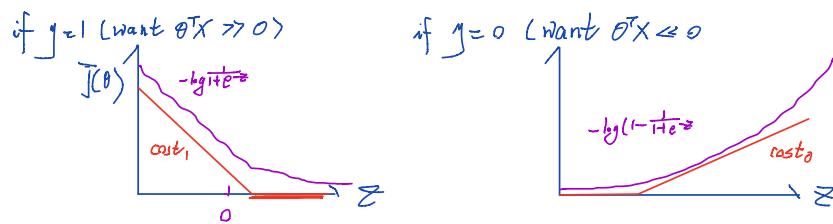
Optimization Objective

Alternative view of logistic regression

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

$$J(\theta) = -y \log h_{\theta}(x) + (1-y) \log(1-h_{\theta}(x))$$

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



Support Vector machine (SVM)

logistic regression $\min_{\theta} \frac{1}{m} \left[\sum_{i=1}^m y^{(i)} (-\log h_{\theta}(x^{(i)})) + (1-y^{(i)}) (\underline{-\log(1-h_{\theta}(x^{(i)}))}) \right] + \frac{\lambda}{2} \sum_{j=0}^n \theta_j^2$

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

In logistic regression $A + \lambda B$

In Support Vector machine $CA + B \quad (C = \frac{1}{\lambda})$

Cost function for SVM:

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

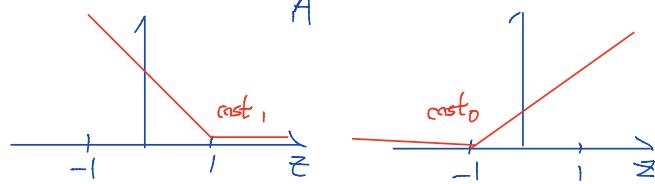
Hypothesis of SVM

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

... | if else

Large Margin intuition

$$\text{SVM: } \min_C \sum_{i=1}^m [\underbrace{y^{(i)} \text{cost}_1(\theta^T x^{(i)})}_{A''} + (1-y^{(i)}) \text{cost}_0(\theta^T x^{(i)})] + \frac{1}{2} \sum_{j=1}^n Q_j^2$$



if $y=1$, we want $\theta^T x \geq 1$ (not just >0)

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

Suppose C very large, $\min_C C A + B \Rightarrow A \rightarrow 0$

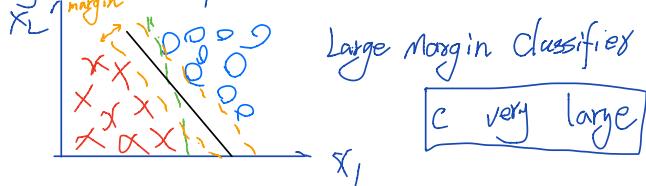
whenever $y^{(i)}=1$:

$$\theta^T x \geq 1 \quad \min C \times 0 + \frac{1}{2} \sum_{j=1}^n Q_j^2$$

whenever $y^{(i)}=0$

$$\theta^T x \leq -1, \quad \min C \times 0 + \frac{1}{2} \sum_{j=1}^n Q_j^2$$

SVM Decision Boundary: Linearly separable case



SVM is a large margin classifier

Mathematics behind Large Margin Classification

SVM decision boundary when C is very large:

(for simplicity, assume $n=2$, $\theta_0=0$)

$$\min \frac{1}{2} \sum_j \theta_j^2 = \frac{1}{2} (\theta_1^2 + \theta_2^2) = \frac{1}{2} \|\theta\|^2 \quad \theta = [\theta_1 \quad \theta_2]$$

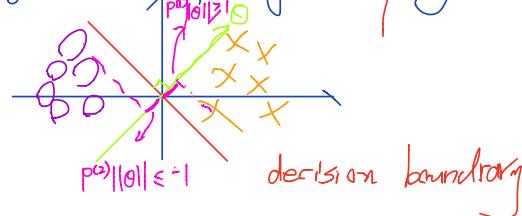
if $y^{(i)}=1$, $\theta^T x^{(i)} \geq 1 \Rightarrow p^{(i)} \cdot \|\theta\| \geq 1$

if $y^{(i)}=0$, $\theta^T x^{(i)} \leq -1 \Rightarrow p^{(i)} \cdot \|\theta\| \leq -1$

where $p^{(i)}$ is $x^{(i)}$ projection on to θ

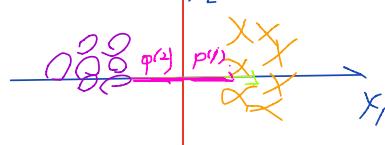
Now prove that SVM with large C will end up with large decision margin (large margin classification)

1) SVM (C large) will not do following small margin:



projection $p^{(i)}$ is small. means $\|\theta\|$ has to be large to make $p^{(i)} \|\theta\| \geq 1$ and $p^{(i)} \|\theta\| \leq -1$, but cost function want $\frac{1}{2} \|\theta\|^2$ to be small, so SVM (large C) will not end up with this situation

2) SVM (large C) will try to do following



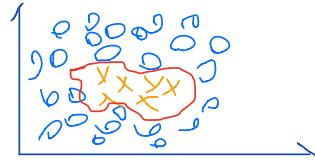
In this case $p^{(i)}$ is larger, to fulfil $p^{(i)} \|\theta\| \geq 1$ and $p^{(i)} \|\theta\| \leq -1$, $\|\theta\|$ doesn't need to be too large, which agrees the minimization of cost function $\min \frac{1}{2} \|\theta\|^2$

thus SVM (with large C) is a Large Margin Classification

Assumption $\theta_0=0$ make θ go through origin of wbar plot

Kernels.

Nonlinear Decision Boundary



predict $y=1$ if $\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 + \dots > 0$

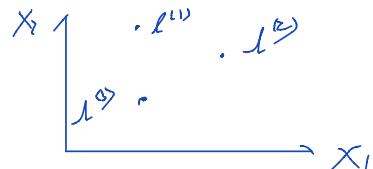
$$h_{\theta}(x) = \begin{cases} 1 & \text{if } \theta_0 + \theta_1 x_1 + \dots > 0 \\ 0 & \text{else} \end{cases}$$

can rewrite: $\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 + \dots$

where: $f_1 = x_1, f_2 = x_2, f_3 = x_1 x_2, \dots$

Q: How to optimize the choice of f_1, f_2, f_3, \dots ?

Kernel



Pick landmarks $l^{(1)}, l^{(2)}, l^{(3)}$, given x , compute new feature depend on proximity to landmarks $l^{(1)}, l^{(2)}, l^{(3)}$

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

$$f_2 = \text{similarity}(x, l^{(2)}) = \exp(\dots)$$

$$f_3 = \dots$$

kernel $k(x, l^{(1)})$ Gaussian kernel

kernels and similarity

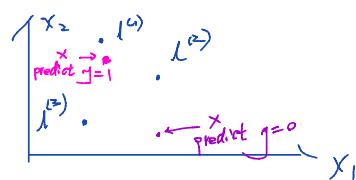
$$f_i = \text{similarity}(x, l^{(i)}) = \exp\left(-\frac{\|x - l^{(i)}\|^2}{2\sigma^2}\right) = \exp\left(-\frac{\sum_{j=1}^n (x_j - l_j^{(i)})^2}{2\sigma^2}\right)$$

if $x \approx l^{(i)}$: $f_i \approx 1$

if x far from $l^{(i)}$, $f_i \approx 0$

each of landmarks defines a new feature

Example:



predict 1 when $\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 + \theta_4 f_4 \geq 0$

say $\theta_0 = -0.5$, $\theta_1 = 1$, $\theta_2 = 1$, $\theta_3 = 0$

case 1: $f_1 \approx 1$, $f_2 \approx 0$, $f_3 \approx 0$

$$\theta_0 + \theta_1 x_1 + \theta_2 x_0 + \theta_3 x_0 = -0.5 + 1 = 0.5 > 0 \quad \text{predict 1}$$

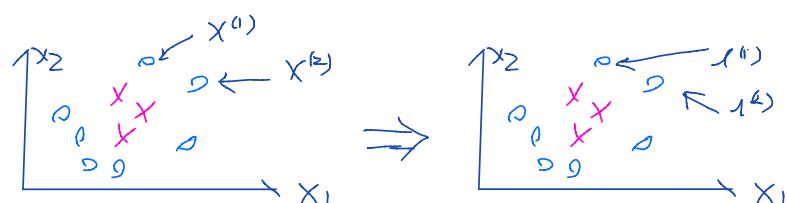
case 2: $f_1, f_2, f_3 \approx 0$

$$\theta_0 + \theta_1 x_0 + \theta_2 x_0 + \theta_3 x_0 = -0.5 < 0 \quad \text{predict 0}$$

Understand: for points close to landmarks, predict 1

for points far away from landmarks, predict 0

choosing the landmarks



choose landmark exactly where the training data is
in training examples \Rightarrow m landmarks

Given $(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$
choose $\ell^{(1)} = x^{(1)}, \dots, \ell^{(m)} = x^{(m)}$

Given example x :

$$f_1 = \text{Similarity}(x, \ell^{(1)})$$

$$f_2 = \dots$$

$$f_m = \dots$$

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

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

$$f_0^{(i)} = \text{Sim}(x^{(i)}, \ell^{(0)})$$

$$f_1^{(i)} = \text{Sim}(x^{(i)}, \ell^{(1)})$$

$$\vdots \quad \vdots \quad \leftarrow f_i^{(i)} = \text{Sim}(x^{(i)}, \ell^{(i)}) = 1$$

$$f_m^{(i)} = \text{Sim}(x^{(i)}, \ell^{(m)})$$

$$\text{then: } x \in \mathbb{R}^{n+1} \rightarrow f^{(i)} = \begin{bmatrix} f_0^{(i)} = 1 \\ f_1^{(i)} \\ \vdots \\ f_m^{(i)} \end{bmatrix} \text{ new feature function}$$

SVM with kernels ($n=m$ or $n=m+1$)

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

predict " $y=1$ " if $\theta^T f \geq 0$ $\theta^T f = \theta_0 f_0 + \theta_1 f_1 + \dots + \theta_m f_m$ $\theta \in \mathbb{R}^{n+1}$

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

In real software/packages, it's actually $\frac{1}{2} \sum_{j=1}^{n+1} \theta_j^2$ instead of $\|\theta\|^2$
where M is some kernel related matrix, in order to let the algorithm scale to much large datasets and computes faster

* SVM with kernel usually runs slow

SVM parameters:

$C (= \frac{1}{\lambda})$: Large C : lower bias, high variance
Small C : higher bias, low variance

σ^2 : large σ^2 : features f_i vary more smoothly
higher bias, lower variance
Small σ^2 : feature f_i vary less smoothly
lower bias, higher variance

Using An SVM

svm software package: liblinear, libsvm

Need to specify:

- choice of parameter C
- choice of kernel

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

predict " $y=1$ " if $\theta^T X \geq 0$

- good for n (# of features) large, m (# of examples) small
to avoid overfitting

e.g. Gaussian kernel

need to choose σ^2

* good for n small, m large

* Note: Do perform feature scaling before using Gaussian kernel, reason is

$$\|x - \mu\|^2 = (x_1 - \mu_1)^2 + (x_2 - \mu_2)^2 + \dots + (x_n - \mu_n)^2$$

$$\sim 10^3 \quad \sim 1$$

then x_2 doesn't show much influence

other choice of kernel:

Note: Not all similarity functions make valid kernels

Need to satisfy "Mercer's Theorem" to make sure SVM packages' optimizations run correctly and do not diverge.

Many off-the-shelf kernels available:

- polynomial kernel: $(x^T x + \text{const})^{\text{degree}} (=2, 3, \dots)$
- string kernel
- chi-square kernel
- histogram intersection kernel

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

Logistic regression vs. SVMs

How to make choice in between these two algorithms?

n : number of features ($x \in \mathbb{R}^{n+1}$)

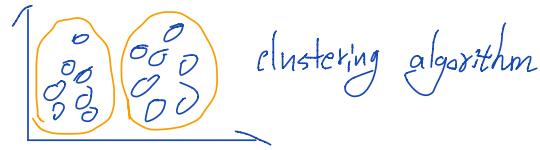
m: number of training examples

- o n is large relative to m (e.g. $n \gg m$, $n=10,000$, $n=10, \dots, 100$)
 - ↳ Use regression, or SVM w/o a kernel ("linear kernel")
- o n is small, m is intermediate (e.g. $n=1-1000$, $m=10-10,000$)
 - ↳ Use SVM w/ Gaussian kernel
- o n is small, m is large ($n=1 \sim 1000$, $M=50,000 + >$)
 - ↳ Create/add more features, then use logistic regression or SVM w/o a kernel ("linear kernel")
- o neural network likely to work well for most of these settings, but may be slower to train

* SVM is a convex problem, can always find ~~the~~ global min

clustering — Unsupervised Learning Intro

Unsupervised Learning



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

Application of clustering.

- Market segmentation
- Social Network analysis
- Organize Computing clusters
- Astronomical data analysis
-

clustering is an example of unsupervised Learning

o Clustering — 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)

- Procedure

- Randomly initialize K cluster centroids $\mu_1, \mu_2, \dots, \mu_K$
 $\in \mathbb{R}^n$
- Repeat {
 - for $i=1$ to m (cluster assignment)

$c^{(i)}$: = index (from 1 to K) of cluster centroid
 closest to $x^{(i)}$

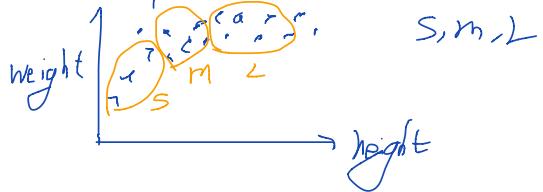
for $k=1$ to K : (update centroid)

μ_k : = average (mean) of points assigned to
cluster k

}

If no points assigned to a certain centroid during a certain step, then eliminate the corresponding centroid.
or randomly re-initialize that centroid

K-means for non-separated clusters



clustering — optimization objective

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

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

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

Optimization objective

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

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

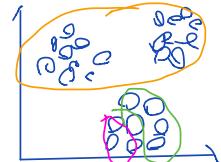
cost function / distortion of k-mean algorithm

- clustering – Random initialization

- should have $k < m$
- randomly pick K training examples
- set μ_1, \dots, μ_K equal to these K examples

Problem of local optima

Assume $K=3$, could end up with:



Solution: Run k-means clustering several times

for $i=1$ to 100 { → 50 - 100

randomly initialize k-means

Run k-means, Get $c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K$

Compute cost function (distortion)

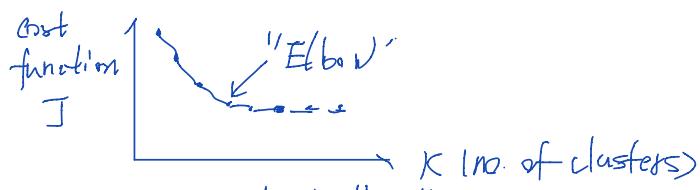
$$J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K)$$

}

Pick clustering that give lowest cost $J(\dots)$

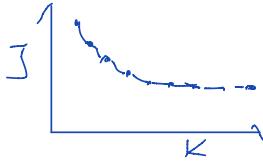
- clustering – choosing the number of clusters

- Elbow method



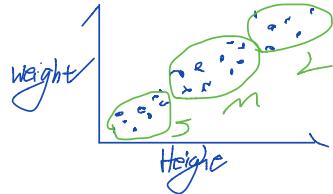
choose K as where corresponding to the "Elbow"

* but "Elbow" is not always obvious, e.g.

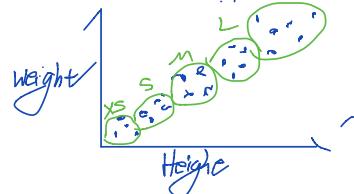


- Base on later purpose

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, e.g. T-shirt size:



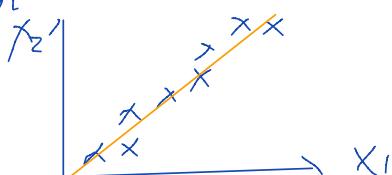
size: S, M, L



size: XS, S, M, L, XL

Dimensionality Reduction

◦ Motivation 1: Data Compression



$$x^{(1)} \in \mathbb{R}^n \rightarrow z^{(1)} \in \mathbb{R}$$

$$x^{(2)} \in \mathbb{R}^n \rightarrow z^{(2)} \in \mathbb{R}$$

$$x^{(3)} \in \mathbb{R}^n \rightarrow z^{(3)} \in \mathbb{R}$$

$$\vdots \quad \vdots \quad \vdots \quad \vdots$$

$$\Rightarrow \longrightarrow 1 \triangleright$$



* This is purely a data projection

$3D \rightarrow 2D$: Project 3D data to 2D plane

o Motivation II: Data visualization

reduce dimensionality to view high dimension data on 2D or 3D or 1D

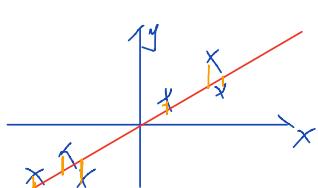
Principal Component Analysis (PCA)

o Problem formulation

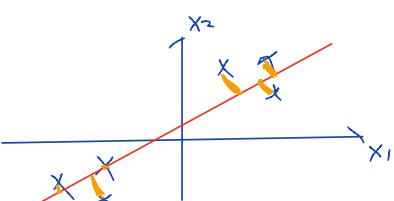
Reduce $2D \rightarrow 1D$: find a direction (a vector $u^{(1)} \in \mathbb{R}^2$) on to which to project the data so as to minimize the projection error.

Reduce from $nD \rightarrow kD$: find k vectors $u^{(1)}, \dots, u^{(k)}$ onto which to project the data, so as to minimize the projection error.

* PCA is not Linear Regression



Linear regression
 $\sum x_i \text{ vs } y$



PCA
 $\sum x_i \text{ vs } x_j, \text{ no } y$

Vertical distance
Projection distance

o PCA Algorithm

1) Data 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, scale features to have comparable ranges

Problem: how to compute $\mu^{(j)}$ and Σ

2) Algorithm: reducing data from $n \rightarrow k \rightarrow$:

- Compute "Covariance matrix"

$$\Sigma = \frac{1}{m} \sum_{i=1}^n \underbrace{(x^{(i)})}_{n \times 1} \underbrace{(x^{(i)})^T}_{1 \times n} \Rightarrow n \times n$$

- Compute "eigenvectors of matrix Σ :

$$[U, S, V] = \text{srd}(\Sigma) \quad (\text{In mathematica})$$

srd: Singular Value Decomposition

$$U(n \times n) = \left[\begin{array}{c|c|c|c} 1 & 1 & \dots & 1 \\ u^{(1)} & u^{(2)} & \dots & u^{(n)} \end{array} \right] \text{ eigenvectors}$$

Take the first k eigenvectors

- Take first k eigenvectors $\dots \rightarrow \begin{bmatrix} 1 & 1 & \dots & 1 \\ u^{(1)} & u^{(2)} & \dots & u^{(k)} \end{bmatrix}$

$$U \quad \text{U reduce} = U_{n \times k}$$

$$\tilde{Z} = U_{\text{reduce}}^T X$$

PCA actually minimizes squared projection error

• Reconstruction from compressed representation

How do you go back from $\tilde{Z}^{(i)}$ to $X^{(i)}$?

$$\tilde{Z} \in \mathbb{R}^k \rightarrow X \in \mathbb{R}^n$$

$$\tilde{Z} = U_{\text{reduce}}^T X \Rightarrow X_{\text{approx}}^{(i)} = U_{\text{reduce}}_{n \times k} \cdot \tilde{Z}^{(i)}$$

• Choosing the number of principal component K

Averaged squared projection error: $\frac{1}{m} \|X^{(i)} - X_{\text{approx}}^{(i)}\|^2$

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

Typically, choose K to be smallest value so that:

$$\frac{\frac{1}{m} \sum_{i=1}^m \|X^{(i)} - X_{\text{approx}}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|X^{(i)}\|^2} \leq 0.01 \quad (1\%)$$

("99% of variance is retained")

Algorithm:

try $k = 1, 2, \dots$

Compute $U_{\text{reduce}}, \tilde{Z}^{(1)}, \dots, \tilde{Z}^{(m)}, X_{\text{approx}}^{(1)}, \dots, X_{\text{approx}}^{(m)}$

check if $\frac{\frac{1}{m} \sum_{i=1}^m \|X^{(i)} - X_{\text{approx}}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|X^{(i)}\|^2} \leq 0.01$

if true, break. find good K value

if True, return \hat{y}

if False, continue.

$$[V, S, Y] = \text{Svd}(\text{Sigma})$$

$$S = \begin{bmatrix} s_{11} & s_{12} & \dots \\ & s_{22} & \\ & & \ddots & s_{nn} \end{bmatrix}$$

for a given k :

$$\frac{\frac{1}{m} \left(\sum_{i=1}^m \|x^{(i)} - x^{(i)}_{\text{approx}}\|^2 \right)^{1/2}}{\frac{1}{m} \sum_{i=1}^m \|x^{(i)}\|^2} = 1 - \frac{\sum_{i=1}^k s_{ii}}{\sum_{i=1}^n s_{ii}} \leq 0.9$$

only need to update numerator for different K

o Advice for applying PCA

Use PCA to speed up learning algorithm

Supervised learning Speed up

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

Extract inputs:

unlabeled dataset: $x^{(1)}, x^{(2)}, \dots, x^{(m)} \in \mathbb{R}^{m \times n}$

\Downarrow PCA

$$z^{(1)}, z^{(2)}, \dots, z^{(m)} \in \mathbb{R}^{m \times k}$$

New training set:

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

New test example

$$x \xrightarrow{\text{PCA}} z$$

$$y = h_n(z) \rightarrow \text{predictions}$$

Note: Mapping $x^{(i)} \rightarrow z^{(i)}$ should be defined by running PCA only on the training set. This mapping can be applied as well to the examples $x_{cv}^{(i)}$ and $x_{test}^{(i)}$ in the cross validation and test sets

Application of PCA:

- Reduce memory/disk needed to store data
- Speed up learning algorithm
(choose K by % of variance retain)
- visualization
 $K=2$ or $K=3$ for visualize high dimensional data

Bad use of PCA I:

Bad thinking: PCA \rightarrow reduce num of features
 \rightarrow prevent overfitting

Actually: This might be ok, but isn't a good way to address overfitting. Use regularization instead

Because: PCA doesn't information of y involved, it might through away some valuable information.
while regularization has all information retained.

Bad use of PCA:

Design of ML system (Bad)

- 1) Get training set $(x^{(1)}, y^{(1)})$, ..., $(x^{(m)}, y^{(m)})$
- 2) Run PCA to reduce $x^{(i)}$ to $z^{(i)}$
- 3) Train
- 4) Test on test set: map $x_{\text{test}}^{(i)}$ to $z_{\text{test}}^{(i)}$, run $h(z^{(i)})$

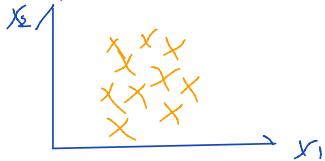
Suggestion:

Before implementing PCA, first try running ML on the raw data. Only if that doesn't do what you want, then implement PCA.

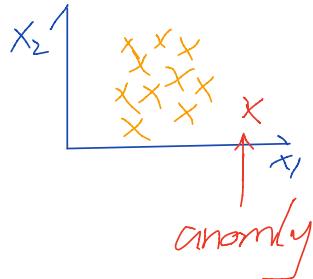
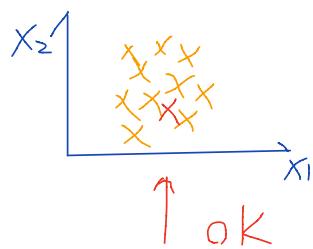
Anomaly Detection — Problem motivation

• An Example

Features $\{x_1^{(1)}, x_2^{(1)}\}, \dots, \{x_1^{(m)}, x_2^{(m)}\}$ has distribution:



If now having a new data set, it could be of following cases:



• Density Estimation

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

Is x_{test} anomalous?

Model: density $P(x)$

$$\begin{cases} P(x_{\text{test}}) < \varepsilon & \rightarrow \text{anomaly} \\ P(x_{\text{test}}) \geq \varepsilon & \rightarrow \text{OK} \end{cases}$$

Example:

— Fraud detection

- $x^{(i)}$ features of user i's activities
- Model $P(x)$ from data

- Identify unusual users by checking which have $p(x) < \epsilon$

- Manufacturing

- Monitoring computers in a data center

Anomaly Detection — Gaussian distribution

Anomaly Detection — Algorithm

o Density Estimation

Training set $\{x^{(1)}, \dots, x^{(m)}\} \quad x \in \mathbb{R}^n$

$$p(x) = P(x_1; \mu_1, \sigma_1^2) P(x_2; \mu_2, \sigma_2^2) \cdots P(x_n; \mu_n, \sigma_n^2)$$

Assume each feature follows Gaussian distribution, and they are independent to each other. (In computation, even if they are not independent, it's fine ?? Andrew said so, I'm not sold)

$$\begin{aligned} p(x) &= P(x_1; \mu_1, \sigma_1^2) \cdots P(x_n; \mu_n, \sigma_n^2) \\ &= \prod_{j=1}^n P(x_j; \mu_j, \sigma_j^2) \end{aligned}$$

Algorithm:

1. choose features x_i that you think might be indicative of anomalous examples

2. Fit parameters $\mu_1, \dots, \mu_n, \sigma_1^2, \dots, \sigma_n^2$

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

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

3. Given next 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 $p(x) < \epsilon$

Building an anomaly detection system

→ Develop anomaly system

The importance of real-number evaluation

When developing a learning algorithm, making decisions is much easier if we have a way of evaluating our algorithm

Assume we have some labeled data, of anomalies and non-anomalous examples ($y=0$ if normal, $y=1$ if anomalous)

Training set: $x^{(1)}, x^{(2)}, \dots, x^{(m)}$, assume they are all normal examples (it's ok if a few anomalous examples are included.)

Cross validation set: $(x_{cv}^{(1)}, y_{cv}^{(1)}), \dots, (x_{cv}^{(m_{cv})}, y_{cv}^{(m_{cv})})$

Test set: $(x_{test}^{(1)}, y_{test}^{(1)}), \dots, (x_{test}^{(m_{test})}, y_{test}^{(m_{test})})$

are labeled data sets (meaning we know exactly if they are anomalous or non-anomalous)

Aircraft engines motivating example

10000 good (normal) engines, 20 flawed engines (anomalies)

Training set: 6000 good engines ($y=0$), 10 anomalies ($y=1$)

CV set: 2000 good engines ($y=0$), 10 anomalies ($y=1$)

Test set: 2000 good engines ($y=0$), 10 anomalies ($y=1$)

Algorithm:

1) Fit model $p(x)$ on training set $\{x^{(1)}, \dots, x^{(n)}\}$

2) on a cv/test set example x , predict

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

Possible evaluation metrics:

- True positive, false positive, false negative

true negative

- Precision / Recall

- F₁-Score

Can also use CV set to choose ε value

Note CV/test set is very skewed !! so classification accuracy is not a good way to evaluate the algorithm.

• Anomaly detection VS supervised learning

Anomaly detection

- very small num of positive examples ($y=1$)
- large num of negative examples ($y=0$)
- Many different "types" of anomaly examples
hard for any algorithm to learn from it
- future anomalies may look nothing like
any of the anomalous examples seen before

Supervised learning

- Large number of positive and negative examples
- Enough positive examples for algorithm to get sense of what positive examples are look like
- Future positive examples likely to be similar to ones in training set

Application

Anomaly detection

- fraud detection
- manufacturing
- Monitoring machines in data center

Supervised learning

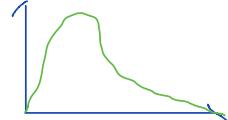
- Email spam classification
- weather prediction
- cancer classification

If have large number of positive examples ($y=1$), anomaly detection can be transferred to supervised learning

- o choosing what features to use

non-gaussian features transform

- Plot the feature histogram to exam if it's Gaussian
- If its not a good gaussian
do some transformation to the feature
to make it more gaussian, e.g. $\log(x + c)$, \sqrt{x} , x^t , ...



Error analysis for anomaly detection

- want: $P(x)$ large for normal ($y=0$) examples \times
 $P(x)$ Small for anomalous ($y=1$) examples \times .

- most common problem:

$P(x)$ is comparable for normal/anomalous examples

- Solution

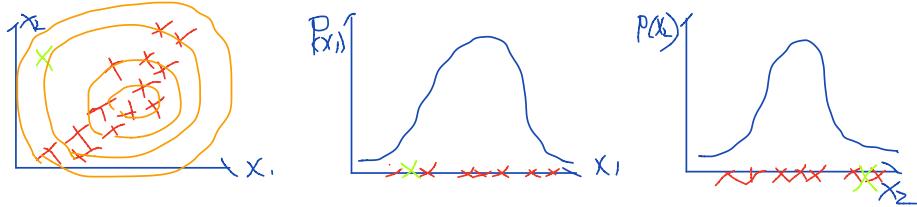
Inspect mistakes and try to create new features

- experience

- choose features that might take on unusually large or small values in an event of an anomaly

Multivariate Gaussian Distribution

- when anomaly detection fails:



Fail to detect $\text{green } x$ as positive

Solution:

$x \in \mathbb{R}^n$, don't model $p(x_1), \dots, p(x_n)$ separately

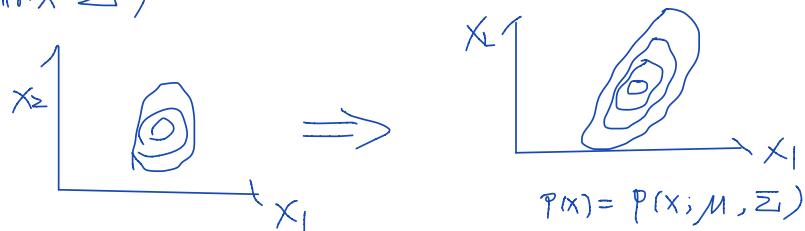
Model $p(x)$ all in one go

Parameters: $\mu \in \mathbb{R}^n$, $\Sigma \in \mathbb{R}^{n \times n}$ (covariance matrix)

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

$$|\Sigma| = \det(\Sigma)$$

By doing this, actually we're allowed to "rotate" the Gaussian distribution (due to off-diagonal value of covariance matrix Σ)



$$p(x) = p(x_1)p(x_2)\dots p(x_n)$$

Algorithm

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

Recommender Systems

- Problem formulation

Example: predicting movie ratings

MOVIE	USER 1	USER 2	USER 3	USER 4	...
movie 1	5	?	3	4	
movie 2	1	2	?	4	
⋮	⋮	⋮	⋮	⋮	

$n_u = \text{No. users}$ $n_m = \text{No. movies}$

$r(i, j) = 1$ if user j has rated movie i

$y^{(u,i)} = \text{rating given by user } j \text{ to movie } i$ (defined only if
when $r(i, j) = 1$)

Question: how to predict movie ratings that hasn't been
rated by a certain user?

- Content based recommendations

$n_u=4, n_m=5$
 x_1 romance x_2 fiction

MOVIE	USER 1	USER 2	USER 3	USER 4		
movie 1	5	?	3	4	0.9	0
movie 2	1	2	?	4	1.0	0.01
...

$$x_0 = 1 \quad x^{(1)} = \begin{bmatrix} 1 \\ 0.9 \\ 0 \end{bmatrix} \quad x^{(2)} = \begin{bmatrix} 1 \\ 1.0 \\ 0.01 \end{bmatrix} \dots$$

num of features : 3 $x_0 = 1, x_1, x_2$

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

problem formulation:

- $r(i, j) = 1$ if user j has rated movie i , else 0
- $y^{(i,j)}$ = rating by user j on movie i (if defined)
- $\theta^{(j)}$ parameter vector for user j
- $X^{(i)}$ feature vector for movie i
- for user j movie i , predicted rating $(\theta^{(j)})^T X^{(i)}$

To learn $\theta^{(j)}$, $m^{(j)} = \text{no. of movies rated by user } j$

$$\min_{\theta^{(j)}} \frac{1}{2m^{(j)}} \sum_{i: r(i,j)=1} ((\theta^{(j)})^T (X^{(i)}) - y^{(i,j)})^2 + \frac{\lambda}{2m^{(j)}} \sum \theta^{(j)}$$

just like linear regression, use gradient descent to find θ

Collaborative Filtering (feature learning)

MOVIE	USER1	USER2	USER3	USER4	x_1	x_2
movie 1	5	?	3	4	?	?
movie 2	1	2	?	4	?	1
...	?	?
...	?	?

If the users can tell us how much they like each kind of movies (e.g. $\theta^{(1)} = \begin{bmatrix} 0 \\ 2 \\ 1 \end{bmatrix}$, $\theta^{(2)} = \begin{bmatrix} 0 \\ 0 \\ 5 \end{bmatrix} \dots \dots \dots$)
 then we can learn the features for each rated movies

The problem turned out to be; what feature value $x_i^{(1)}$ for movie i should be, so that $(\theta^{(1)})^T x^{(1)} \approx 5$, $(\theta^{(2)})^T x^{(1)} \approx 0$, ...

Optimization algorithm.

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

$$\rightarrow \min_{x^{(1)}} \frac{1}{2} \sum_{j: x_{ij} \neq 0} ((\theta_j^{(1)})^T x^{(1)} - y^{(1,j)})^2 + \frac{\lambda}{2} \sum_{k=1}^n (x_k^{(1)})^2$$

Collaborative filtering

- Given $x^{(1)}, \dots, x^{(n)}$ (and movie ratings >

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

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

can estimate $x^{(1)}, \dots, x^{(n)}$

\Rightarrow Guess $\theta \rightarrow x \rightarrow \theta \rightarrow x \rightarrow \theta \rightarrow \dots$

This is awesome !!!

Collaborative filtering algorithm

non-optimized algorithm

(i) Given $x^{(1)}, \dots, x^{(n)}$, estimate $\theta^{(1)}, \dots, \theta^{(n)}$

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

(ii) Given $\theta^{(1)}, \dots, \theta^{(n)}$, estimate $x^{(1)}, \dots, x^{(n)}$

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

Guess $x^{(1)}, \dots, x^{(n)}$ or $\theta^{(1)}, \dots, \theta^{(n)}$, then go back and forth (i) and (ii) to optimize the learning

— Optimized algorithm

Do (i) and (ii) at the same time

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

$$\min_{x^{(1)}, \dots, x^{(n)}, \theta^{(1)}, \dots, \theta^{(n)}} J(x^{(1)}, \dots, x^{(n)}, \theta^{(1)}, \dots, \theta^{(n)})$$

$\text{no } x^0 = 1$ in this case

1) Initialize $x^{(1)}, \dots, x^{(m)}, \theta^{(1)}, \dots, \theta^{(n)}$ to small random values

2) Minimize $J(x^{(1)}, \dots, \theta^{(1)}, \dots)$ using gradient descent, e.g.

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

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

3) for a user with parameters θ and a movie with features x , predict a star rating of $\theta^T x$

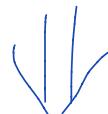
Low rank Matrix Factorization

MOVIE	USER 1	USER 2	USER 3	USER 4	...
-------	--------	--------	--------	--------	-----

movie 1	5	?	3	4	
---------	---	---	---	---	--

movie 2	1	2	?	4	
---------	---	---	---	---	--

...	
-----	----	----	----	----	--



$$Y = \begin{bmatrix} 5 & ? & 3 & 4 \\ 1 & 2 & ? & 4 \\ \dots & \dots & \dots & \dots \\ y[1, j] \end{bmatrix}$$

Predicted ratings

$$\left[(\theta^{(1)})^T x^{(1)} \quad (\theta^{(2)})^T x^{(1)} \quad \dots \right]$$

↓

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

$$= X \circ H^T$$

$$X = \begin{bmatrix} - (x^{(1)})^T & - \\ - (x^{(2)})^T & - \\ \vdots & \end{bmatrix}$$

$$H = \begin{bmatrix} - (\theta^{(1)})^T & - \\ - & \vdots & - \\ - & \vdots & - \end{bmatrix}$$

$$H = \begin{bmatrix} - & & \\ & \ddots & \\ & & - \end{bmatrix}$$

Low rank matrix factorization

- Finding related movies

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

? How to find movies j related to movie i ?

Small $\|x^{(i)} - x^{(j)}\| \rightarrow$ more i, j are "similar"

Implementation detail: Mean Normalization

for user j hasn't rated any movie, $\min J(\theta)$ \rightarrow will give
 $\theta = \begin{bmatrix} ? \\ ? \end{bmatrix}$ so predict this user will always rate movie 0
NOT RIGHT!

Solution: Mean Normalization

$$X = \begin{bmatrix} 5 & 5 & 0 & 0 & ? \\ \dots & \dots & \dots & \dots & \dots \\ - & - & - & - & - \end{bmatrix}$$
$$\Rightarrow \mu = \begin{bmatrix} 2.5 \\ 2.5 \\ \vdots \\ i \end{bmatrix} \rightarrow X' = \begin{bmatrix} 2.5 & 2.5 & -2.5 & 2.5 & ? \\ \sim & \sim & \sim & \sim & \sim \\ \sim & \sim & \sim & \sim & \sim \end{bmatrix}$$

For user j on movie i , predict $(\theta^{(i)})^T(X'^{(i)}) + \mu_i$
In this case the above problem can be avoided.

Learning with large datasets

Stochastic Gradient Descent

- Linear regression with gradient descent

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

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

Repeat {

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m [h_{\theta}(x^{(i)}) - y^{(i)}] X^{(i)} \quad \frac{\partial J_{\text{train}}(\theta)}{\partial \theta_j}$$

$$j = 0, 1, \dots, n$$

}

Batch gradient descent

Computational expensive since need to load all m examples to do one small step scanning through the entire dataset periodically

- Stochastic gradient descent

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

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

1) Randomly shuffle dataset

2) Repeat { (usually 1 ~ 10 times repeating, depends on m) }

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

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

$$\text{for } i = 0, 1, \dots, n > \quad = \frac{\partial}{\partial \theta_i} \text{cost}(\theta, (x^{(i)}, y^{(i)}))$$

why stochastic gradient descent is much faster than gradient descent?

- for large dataset, gradient descent has to run through all training set to make one baby step of descent
- Stochastic gradient descent take one descent step only when run through one training example. Especially when dataset is larger, running through all training examples once is enough to get good training result.

o Mini-batch gradient descent

Batch Gradient Descent uses

Stochastic Gradient Descent uses

Mini Gradient Descent uses

$$\begin{bmatrix} \text{all } m \\ 1 \\ b \end{bmatrix}$$

Examples in each iteration

o Mini-batch gradient descent

1) choose b .

2) Repeat {

for $i = 1, 1+b, 1+2b, \dots$ {

$$\theta_j := \theta_j - \alpha \frac{1}{b} \sum_{k=1}^{b-1} (\hat{h}_\theta(x^{(k)}) - y^{(k)}) X_j^{(k)}$$

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

γ

γ

- Mini-batch gradient descent can be faster than stochastic gradient descent

* vectorization!

Mini-gradient descent use b examples for one step, can be vectorized, meaning partially parallelize !!!

Stochastic Gradient Descent Convergence

- Checking for convergence

$$\cdot \text{cost}(\theta, (x^{(i)}, y^{(i)})) = \frac{1}{n} (\hat{y}_\theta(x^{(i)}) - y^{(i)})$$

• Putting the learning compute $\text{cost}(\theta, (x^{(i)}, y^{(i)}))$ before updating using $(x^{(i)}, y^{(i)})$

• Every 100 iterations (say), print cost averaged over the last 100 examples processed by the algorithm.

* Stochastic gradient descent does not converge to the global minimum, but only oscillates around the global minimum

* Still need to choose the proper learning rate α .

* if want SGD to converge, can slowly decrease learning rate α over time

$$\text{e.g. } \alpha = \frac{\text{const 1}}{\text{iteration} + \text{const 2}}$$

Online learning

Eg1: Shipping service website. users sometimes choose your shipping service ($y=1$)

sometimes not ($y=0$)

Features x capture the properties of users, want to learn $P(y=1|x;\theta)$ to optimize the prize to get maximum users

Background: have large stream of user's data

Repeat {

Get (x, y) corresponding to user

update θ using (x, y)

$$\theta_j := \theta_j - \alpha (h_\theta(x) - y) \cdot x_j \quad (j=0, 1, \dots, n)$$

Discard (x, y) after using it once

}

*Advantage of online learning

- Can adapt to change of users' preferences

Eg2: Product Search (learning to search)

User searches for "Android phone 10mp camera"

Have 100 phones in store. Will return 10 results

x : features of phone that matches the query

$y = 1$ if user clicked the button

0 otherwise

learn $P(y=1 | x; \theta)$ [predicted CTR]

CTR: click through rate

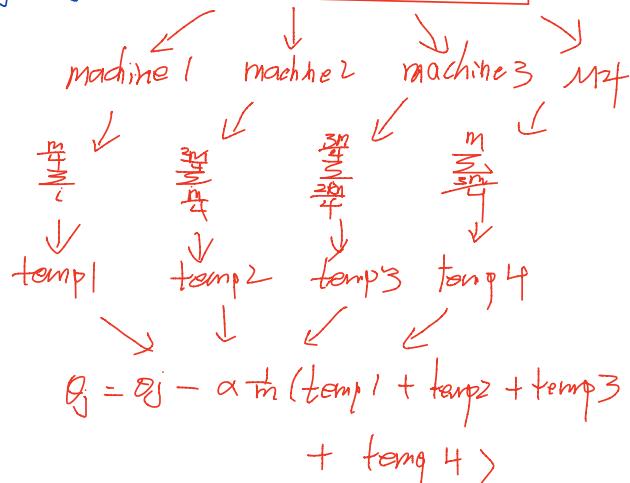
other examples: choosing special offers to users; customized selection of news articles; product recommendation

These examples can also be transferred to classical ML problem

Map-reduce and data parallelism

Map-reduce

$$\text{Batch-gradient descent } \theta_j = \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (\hat{h}_\theta(x^{(i)}) - y^{(i)}) x_j^{(i)}$$



training set



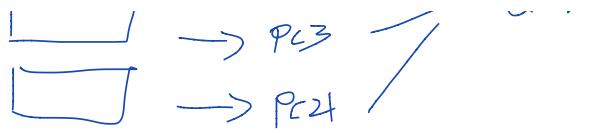
$\rightarrow PC_1$



$\rightarrow PC_2$



Combine results



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

- Multi-moves machines

PC_i → core:

Photo OCR

- Problem description and pipeline

photo OCR : photo Optical Character Recognition

- Photo OCR Pipeline

1. Text detection

2. character segmentation

3. character classification

Each above steps may be a machine learning problem

Pipeline : A system with many components / stages, several of which may use machine learning

- Sliding Windows

Text detection

Getting lots of data : Artificial data synthesis

- Creating new data from scratch.

Take real data and add background / noise / distortion

Distortion should represent the type of noise in the test set.

- Discussion on getting more data

1. Make sure you have low bias classifier before expanding the effort.

2. How much work would it be to get 10x much data as we

currently have?

- Artificial data synthesis?
- Collect/label it yourself?
- Crowd source (e.g. Amazon mechanical Turk)

Ceiling Analysis: What part of pipeline to work on next?

Image → Text detection → character segmentation → character recognition

Ceiling analysis: estimating the error due to each component