Coursera Machine learning
handwritten notes.

Taught by: Prof. Andrew Ng.

Github:
JaneShu99

Summary of course

Supervised learning
 └ linear regression, logistic regression, neural networks, SVMs.

Unsupervised learning
 └ K-means, PCA, Anomaly detection.

Special applications /special topics
 └ recommender systems, large scale machine learning.

Advice for building a machine learning system.
 └ Bias/variance, regularization, deciding what to work on next
 evaluation of learning algorithms, learning curves, error analysis,
 ceiling analysis.

Coursera Machine learning course
Notes By Jane Shi

**Week 1**

## Introduction

Definition of Machine learning.

A computer program is said to learn from experience E with respect to some task T and some performance measure P if its performance on T as measured by P improves with experience E.

Main two types: (explored in this course)

→ Supervised machine learning vs unsupervised machine learning

→ advice for practical uses of machine learning

→ how to develop ML systems?

## Supervised learning

→ We gave a dataset (where the "right answer" are given) we know what our answer looks like

as a result of relation between input & output

→ regression problem: predict continuous value output.

→ classification problem: predict discrete value output.

→ take account of various number of inputs. / features / infinite many attributes

## Unsupervised learning

→ determine clustering of data, where we have little/ no idea about what result should look like

→ identify cohesive groups of data

→ example: cocktail party problem
      given two recording, with two tracks of different volume, output each sound track   mixed

→ can be written in one line (solution).

→ octave is good! built for lin alg & related programming.

## Model & cost function

## Model representation

↳ linear regression model
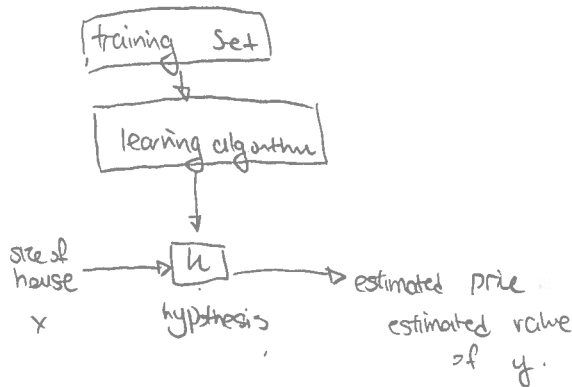
    ↳ training set is the data-set.

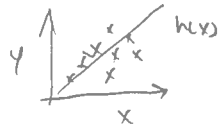    ↳ m = # of training example, x's input var/feature, y's output var

      $(x,y)$ ← training example

      $(x^{(i)}, y^{(i)})$ - ith training example.

    ↳

        | training set |

        | learning algorithm |

    size of house  → [ h ] → estimated price

      x      hypothesis     estimated value of y.

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

    ↳ this is lin. reg w/ 1 variable / univariate lin. reg.

## Cost function

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

                ↑  ↗

              params      ⌐ # of training examples.

    ↳ goal is to minimize $\frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2$

                              $\parallel$

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

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

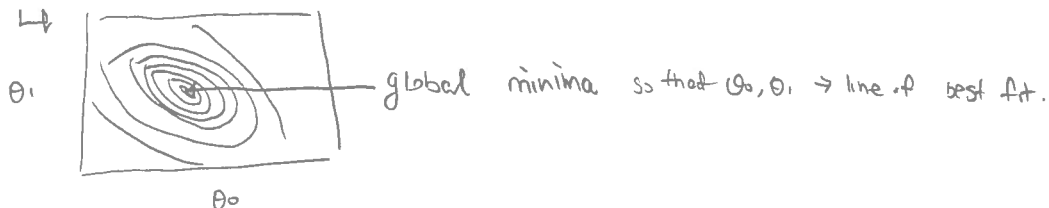       $\theta_0, \theta_1$                      ⌣ Cost function.

    ↳ the square error function.

    ↳ $J(\theta_0, \theta_1)$ is a function in $\theta_0, \theta_1$. Plot $J(\theta_0, \theta_1)$ & find your global minima

    ↳ contour graphs are used for multiple features. (Plot 3D graph)

    ↳

$\theta_1$                — global minima so that $\theta_0, \theta_1$ → line of best fit.

         $\theta_0$

↳ the graphs cannot __always__ be visualized as easily. Thus, we would need some other algo.

↳ Gradient descent algorithm.

  ↳ have function $J(\theta_0, \theta_1)$

  want   $\min\limits_{\theta_0, \theta_1} J(\theta_0, \theta_1)$

  ↳ you start with some $\theta_0, \theta_1$, then keep changing $\theta_0, \theta_1 \rightarrow$ reduce $J(\theta_0, \theta_1)$ each iteration.

  ↳ via calculus

  ↳ you can end up at two different local optimums.

the algorithm:          $(:=)$ ↳ assignment operator

  repeat until converge {

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

  }

  correct __simultaneous__ update:

  $temp0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$          ↳ $\alpha$ is the learning rate.

  $temp1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$            how big step we go down hill?

  $\theta_0 = temp0$                                big step / baby step?

  $\theta_1 = temp1$                      ↳ simultaneously update $\theta_0$ and $\theta_1$, ad same time.

↳ when updating, takes consideration of whether $\frac{\partial}{\partial \theta_0}$ is positive or negative,

  So the new point is closer to x axis. / the absolute value of $\frac{\partial}{\partial \theta_0}$ approach to 0 gradually).

↳ need to choose $\alpha$ so it's not too small, not too large.

  if $\alpha$ too small $\rightarrow$ slow algorithm

  if $\alpha$ too large $\rightarrow$ may even diverge

↳ Putting it altogether:

  Gradient Descent algorithm                          linear regression model.

  repeat until converge {

  $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$                     $h_\theta(x) = \theta_0 + \theta_1 x$

  }                     for $(j = 1, j = 0)$          $J(\theta_0, \theta_1) = \frac{1}{2m} \sum\limits_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2$

                  apply  ↓  to  ← to minimize

Plug in the equation, we obtain

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

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

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

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

It's always a convex function.

Our linear regression algorithm turns out to be

repeat until converge {

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

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

} 

"Batch Gradient Descent" each step of gradient descent uses all training examples.

└note: must use the model for $J(\theta_0, \theta_1)$ where there's no other local optima than the global,
   or else it can   end up at   another local min

## Week 2

Multi-feature linear regression
   └having multiple features
       notation:     n = # of feature
                   $x^{(i)}$: input features of $i^{th}$ example (vector)
                   $x_j^{(i)}$: value of feature $j$ in the $i^{th}$ training example.
   └hypothesis:

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \cdots + \theta_n x_n$$

for convenience, $\forall x$, $x_0 = 1$

So $h_\theta(x) = \sum_{i=0}^{n} \theta_i x_i$

4.

$$X = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} \qquad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

Hypothesis:
$$h_\theta(x) = \theta^T x$$

or inner product, $\langle \theta, x \rangle$

Parameter: $\theta$

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

Gradient descent:

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

Simultaneous update $\theta_j$, $j = 0, 1, 2, \cdots n$.
since you're taking derivative with respect to $j$th feature.

Gradient descent in practice:

└ feature scaling.

    └ make sure features are on a similar scale.

    └ $\forall i$, $-1 \leq x_i \leq 1$.

    └ major values around $-3 \vee +3$ ish    not too little as in $\sim 0.1$

└ Mean normalization

    └ replace $x_i$ with $x_i - \mu_i$, to make sure feature have $\sim 0$ mean

    └ do not apply to $x_0 = 1$ though!

    average value of $x_i$
$$x_i \leftarrow \frac{x_i - \mu_i}{s_i} \quad \leftarrow \text{range or std.}$$

└ "debugging" make sure it works properly

└ how to choose your $\alpha$?

↳ "Debugging"  make plot where #iter is x-axis, min J(θ) y,
    ↳ J(θ) should always decrease due to # of iter (every single iter!)
    ↳ if J(θ) error increases, you want to decrease $\alpha$.
↳ convergence test: choose $\varepsilon$ to declare when $J(θ) < \varepsilon$ → converges!
↳ tip: to choose $\alpha$, try 0.001, 0.01, 0.1, 1, ---- try a range of values.
    0.001, 0.003, 0.01, 0.03, 0.1, 0.3, 1, ----

## Features & Polynomial regression

↳ Combine multiple features into 1.

    ↳ combine $x_1$ and $x_2$, by taking $x_3 = x_1 \cdot x_2$
    ↳ Polynomial regression if linear doesn't fit.
        ↳ change the behaviour, so it can be quadratic/cubic etc.
        ↳ ideas = $h_θ(x) = θ_0 + θ_1 x_1 + θ_2 x_1^2 + θ_3 x_1^3$
                    ↑      ↑
                 feature $x_2$  feature $x_3$

$$h_θ(x) = θ_0 + θ_1 x + θ_2 \sqrt{x_1}$$

↳ with this though, keep in mind, feature scaling is very important.

## Normal equation (computing param analytically)

↳ X: design matrix.  $x^{(i)} = \begin{bmatrix} x_0^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \in \mathbb{R}^{n+1}$, then $X = \begin{bmatrix} - (x^{(1)})^T - \\ - (x^{(2)})^T - \\ \vdots \\ - (x^{(n)})^T - \end{bmatrix}$

↳ optimum θ given by

$$θ = (X^T X)^{-1} X^T y$$        octave: print $(x' * x) * x' * y$    $x'$: transpose  * : matrix mult

↳ Note with normal equation, you DONT need feature scaling.

| Gradient descent | VS | Normal equation | |
|---|---|---|---|
| ↳ need to choose $\alpha$ | | ↳ no need to choose $\alpha$ | $\begin{cases} m \# \text{ training example} \\ n \# \text{ of feature} \end{cases}$ |
| ↳ many iteration | | ↳ no iteration needed | |
| ↳ work well even if $n$ is large. | | ↳ $(X^T X)^{-1}$ takes $O(n^3)$ | |
| | | ↳ slow when $n$ is large ($\geq 10,000$) | |

6.

## Normal equation / noninvertibility

↳ what if $X^TX$ is non-invertible?

 ↳ use "pinv" instead of 'inv' (pseudo-inverse)

 ↳ it gives you $\theta$ though $X^TX$ is singular

 ↳ ① happen when there's redundant feature. or ② too many feature: $m < n$, then use regularization) /or delete features.

## Vectorization

helps to compute vectors faster.

Assignment questions include:

 ↳ computing cost for multi / uni variable dataset
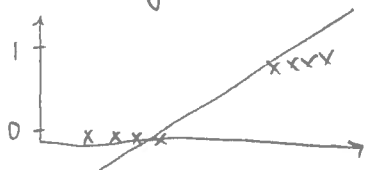
 ↳ computing cost for multiple variable

 ↳ gradient descent for multi / uni variables.

## Week 3

Work with <u>classification problem</u>

↳ $y \in \begin{cases} 0 & \text{Negative class} \\ 1 & \text{Positive class.} \end{cases}$

↳ now: binary class classification.     Does lin-reg work? <u>no</u>. not a good idea.



threshold = 0.5.

If $h_\theta(x) \begin{cases} \geq th \\ < th \end{cases}$ predict $\begin{cases} 1 \\ 0 \end{cases}$

↳     this is the bug! ⤋    mess up your lin reg data! ˙ᗡ˙

↳ So don't use lin reg for classification.

↳ logistic regression: $0 \leq h_\theta(x) \leq 1$
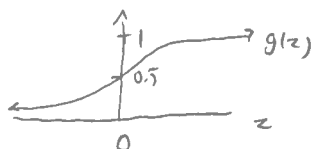
 this is a classification algorithm

## logistic regression

└ want: $0 \le h_\theta(x) \le 1$

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

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

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

└ Interpretation: $h_\theta(x)$ gives you probability that our output is 1. $= P(y=1 \mid x; \theta)$ in probability notation.

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

$\iff P(y=1 \mid x; \theta) + P(y=0 \mid x; \theta) = 1$

## Decision Boundary

└ $h_\theta(x) \ge 0.5 \to y=1$ or $\theta^T x \ge 0 \to y=1$

$h_\theta(x) < 0.5 \to y=0$ or $\theta^T x < 0 \to y=0$

since └ $g(z) \ge 0.5 \iff z \ge 0$.

decision boundary is the line that separate area when $y=0, y=1$ (line where $h_\theta(x) = 0.5$ exactly.)

there are also non linear decision boundaries. then you need more params for higher dim. ie.

$h_\theta(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$

$\to -1 + x_1^2 + x_2^2 \ge 0$ results in $\bigcirc$ boundary

## logistic regression model

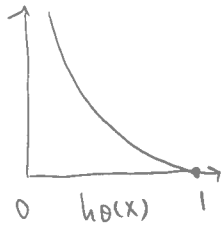training set $= \{(x^{(1)}, y^{(1)}), \dots (x^{(m)}, y^{(m)})\}$

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

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

lin reg won't give a convex, but we <u>want</u> a convex function

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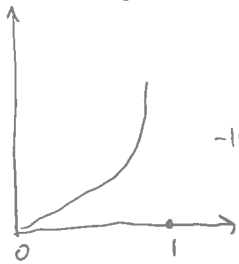
if $y=1$



Cost=0 if $x=1$ but as $h_\theta(x) \to 0$, cost $\to \infty$.

Intuition: if $h_\theta(x)=0$, but you predict it as $1$, you're penalized.

if $y=0$



$-\log (1-z)$

similar as the other intuition

this gives a convex & local optimum free function

note: $y=1$ or $y=0$ always. $\to$ Can combine two equations

the compressed cost function is:

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

total cost J:

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

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

want: $\underset{\theta}{MIN} \, J(\theta)$

gradient descent algorithm:

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

or

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

looks same as lin.reg's grad.des, but! $h_\theta(x)$ refers to $\frac{1}{1+e^{-\theta^T x}}$ now.

9

vectorized implementation

cost
$$h = g(X\theta) \quad \text{this computes quantity } h_\theta(x^{(i)})$$
$$J(\theta) = \tfrac{1}{m} \Sigma \left[ y^{(i)} \log(h_\theta(x^{(i)})) + (1-y^{(i)}) \log(1- h_\theta(x^{(i)})) \right]$$
$$\hookrightarrow J(\theta) = \tfrac{1}{m} \cdot (-y^T \log(h) - (1-y)^T \log(1-h))$$

the gradient descent

Idea: rearrange the vectors until it's easy to type into Matlab. :)
$$\theta := \theta - \alpha \tfrac{1}{m} \sum_{i=1}^{m} \left[ (h_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)}_j \right] \quad \leftarrow \text{column vector}$$

$$\theta := \theta - \tfrac{\alpha}{m} X^T (g(X\theta) - \bar{y})$$

since:  thus,  $x^T$ makes each $x^{(i)}$ a column vector.

$$\begin{bmatrix} \\ \end{bmatrix} \text{ do column wise calculation.}$$
$\theta$ is a $(n \times 1)$ vector.

$X\theta$ returns $x^{(i)} \theta \rightarrow$ vectors corresponding index.

$$X: \begin{matrix} x^{(i)} \end{matrix} \overbrace{\begin{bmatrix} - x^{(1)T} - \\ \vdots \\ - x^{(m)T} - \end{bmatrix}}^{n+1} \Big\} m$$

$\theta = n \times 1$ vec.
$x = m \times n \qquad X\theta = m \times 1$.
$x^T = n \times m$
ans: $n \times 1$.

Advanced Optimization

$\begin{cases} \text{cost function } J(\theta), \text{ want } \min_\theta J(\theta). \\ \text{Given } \theta, \text{ if we can compute } J(\theta), \frac{\partial}{\partial \theta_j} J(\theta) \end{cases}$

as long as you know these two you can use the lib functions
then we can use the following algorithms

$\hookrightarrow$ Conjugate gradient
$\hookrightarrow$ BFGS
$\hookrightarrow$ L-BFGS

$\Big\}$ faster, no need for $\alpha$, but more complex. So we use the library

use function "fminunc()"
Plugin the  $J(\theta)$ & the gradients shall suffice.

logistic optimization for multiple classes    "one vs all classification"

Multiclass classification.
$y = \{0, 1, \dots n\}$ each are category.
assign one class as positive, all other ones, as "the rest"
$y \in \{0, 1, 2, \dots n\}$
$h_\theta^{(0)}(x) = P(y = 0 | x; \theta)$
$h_\theta^{(1)}(x) = P(y = 1 | x; \theta)$
$\vdots$
$h_\theta^{(n)}(x) = P(y = n | x; \theta)$
predictions: $\max h_\theta^{(i)}(x)$

## Problem of over-fitting

under-fitting: hypothesis function maps too poorly to the trend of data. too simple (too little features. ~~function~~

overfitting = not generalized enough. fits available data too well, but might have unnecessary angles/corners i.e. too wiggly (fail to generalize):

to resolve overfitting:
1) reduce # of features. (model selection algorithm to ditch less-important features.)
2) regularization (reduce magnitude of $\theta_j$)

Cost function (the new one with regularization)

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

big $\lambda$ bumps up and forces $\theta_j$ to be small because big $\theta_j$ will be penalized

### Gradient descent

repeat {

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

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

}

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

↑ always less than 1 as it reduce $\theta_j$ each time by a little bit.

### Normal equation

$$\theta = (X^T X + \lambda \cdot L)^{-1} X^T y$$

where $L = \begin{bmatrix} 0 & & & \\ & 1 & & \\ & & 1 & \\ & & & \ddots \\ & & & & 1 \end{bmatrix}$  $\in M_n^{(n+1) \times (n+1)}$

note if m<n, $X^T X$ is non-invertible but adding L makes it invertible. regularization solves non-invertibility as well.

## Regularized logistic regression (advanced optimization works similarly).

regularized cost function for linear regression.

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

← new term

### Gradient descent:

repeat {

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

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

# Advanced functions (regularization).

Jval: same as previous.

gradient 1: (index 0)

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

gradient (2 ~ n+1) index k (1,2, ... n))

$$\left( \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)}) x_1^{(i)} + \boxed{\frac{\lambda}{m} \theta_j} \right) \quad \swarrow \text{newly added term}$$

watch out following when doing assignment:

↳ display dimensions might be in opposite order.

↳ draw out matrices carefully to visualize the vectorization.

↳ match matrix dimensions always.

## Week 4

### Neural Networks - representation

↳ Computer vision - example.
    ↳ logistic regression would have too many features. (like a few million for images)
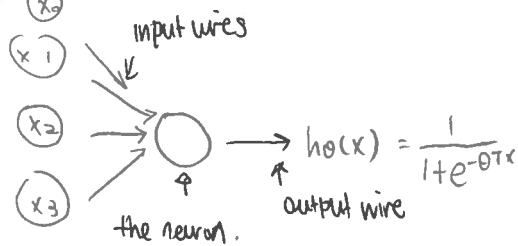↳ mimic the brain.
↳ large scale!
↳ neural-rewiring experiment
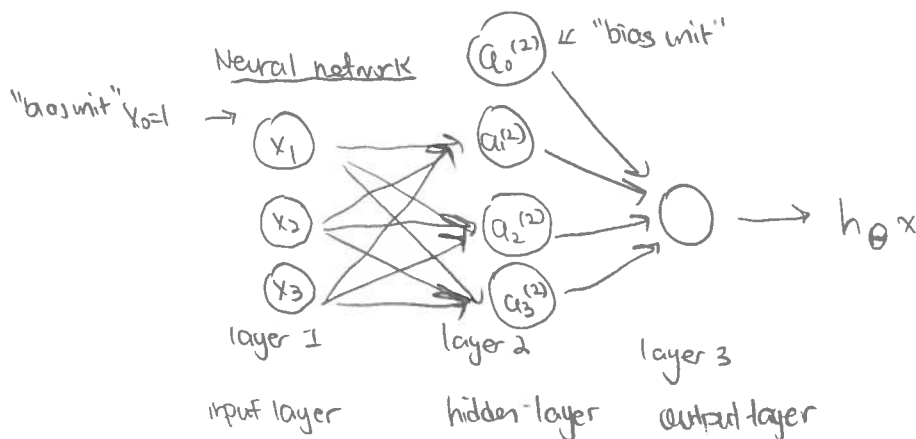↳ adjust/learn the data

### Model representation

Neuron model : logistic unit

=1, "bias unit" → $(x_0)$

$(x_1)$

input wires

$(x_2)$

$(x_3)$

the neuron.

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

output wire

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

"weights"

"params"

" sigmoid activation function."

### Neural network

"bias unit" $x_0 = 1$ →

$(x_1)$

$(x_2)$

$(x_3)$

$(a_0^{(2)})$ ← "bias unit"

$(a_1^{(2)})$

$(a_2^{(2)})$

$(a_3^{(2)})$

$\bigcirc \rightarrow h_\theta x$

layer 1      layer 2      layer 3

input layer      hidden-layer      output layer

$\begin{cases} a_i^{(j)} = \text{"activation" of unit } i \text{ in layer } j \\ \theta^{(j)} = \text{matrix of weights controlling func mapping from layer } j \text{ to } j+1 \end{cases}$

### Vec representation

$\begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \rightarrow \begin{bmatrix} a_1^{(2)} \\ a_2^{(2)} \\ a_3^{(2)} \end{bmatrix} \rightarrow h_\theta(x)$

13.

vector representation "activation nodes"    _example_

layer 1 to layer 2:

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

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

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

layer 2 to layer 3

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

<u>dim ($\theta$)</u>
   └ each layer has its own matrix of weights
      └ if network has $S_j$ layers in level $J$, $S_{j+1}$ layers in level $J+1$, then $\theta^{(j)}$ has dimension
         $$S_{j+1} \times (S_j + 1)$$
            ↑
         comes from bias node
      └ laidout like this, b/c multiply $\theta$, the vector will be on the right.
   intuition: Neural network allows nodes in its hidden layer to "learn" its own features.


<u>Vectorization of computation</u>

   └ $a_1^{(2)} = g(z_1^{(2)})$                    for layer j, node k, z is
   └ $a_2^{(2)} = g(z_2^{(2)})$
   └ $a_3^{(2)} = g(z_3^{(2)})$                    $$z_k^{(j)} = \theta_{k,0}^{(j-1)} x_0 + \theta_{k,1}^{(j-1)} x_1 + \cdots + \theta_{kn}^{(j)} x_n$$

$$X = \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \qquad Z^{(j)} = \begin{bmatrix} z_1^{(j)} \\ z_2^{(j)} \\ \vdots \\ z_n^{(j)} \end{bmatrix}$$

$$Z^{(j)} = \theta^{(j-1)} a^{(j-1)} \qquad \text{note: } \dim(\theta^{(j-1)}) \text{ is } S_j \times (n+1)$$
$$\dim(a^{(j-1)}) \text{ is } (n+1) \times 1.$$

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

adding the bias unit to layer j after computing $a^{(j)}$  i.e. $a_0^{(j)} = 1$.

   to compute final hypothesis, compute z vector
      $z^{(j+1)} = \theta^{(j)} a^{(j)}$   the last matrix $\theta^{(j)}$ has only 1 row, multiplied by   one column vec $a^{(j)}$
      so the result is a real number.

      $$h_\theta(x) = a^{(j+1)} = g(z^{(j+1)})$$

14

Multi-class Classification

    ↳ one-vs-all method

$$h_\theta(x) \in \mathbb{R}^4 \quad \text{if} \quad \text{there are 4 classes}$$

$$h_\theta(x) \approx \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \text{ or } \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix} \text{ or } \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} \text{ or } \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

    with different input x.

    It returns one of the $e_i$'s vector given a particular input

**Week 5**   Goal: learn how to train neural networks

the **cost function** for the neural network.

    ↳ $L$ = total # of layers in the network.

    ↳ $S_l$ (# of units not counting bias unit in layer $l$)

    ↳ $K$ = # of output unit / classes.

$$J(\theta) = -\frac{1}{m}\sum_{i=1}^{m}\sum_{k=1}^{K}\left[ y_k^{(i)} \log((h_\theta(x^{(i)}))_k) + (1-y_k^{(i)})(\log(1-h_\theta(x^{(i)}))_k)\right] + \frac{\lambda}{2m}\sum_{l=1}^{L-1}\sum_{i=1}^{S_l}\sum_{j=1}^{S_{l+1}}(\theta_{j,i}^{(i)})^2$$

**Back propagation algorithm**

    ↳ goal is to compute $\min_\theta J(\theta)$

    ↳ look at partial derivative of $J(\theta)$

$$\frac{\partial}{\partial\theta_{i,j}^{(l)}} J(\theta)$$

    the back propagation algorithm works as follows:

    ↳ given training set $\{(x^{(1)}, y^{(1)}), \dots (x^{(m)}, y^{(m)})\}$

    ↳ set $\Delta_{i,j}^{(l)} := 0 \quad \forall i,j$

    ↳ for training example $l = 1 \sim m$

      1. set $a^{(1)} := x^{(l)}$

      2. perform forward propagation to compute $a^{(l)}$, $l = 1, 2, 3, \dots L$

        (i.e. set up $z$ (intermediate, use $g(z)$ to calculate next layer)

      3. $\delta^{(L)} = a^{(L)} - y^{(l)}$

      4. Compute $\delta^{(L-1)}, \delta^{(L-2)}, \dots \delta^{(2)}$ using $\delta^{(l)} = ((\theta^{(l)})^T \delta^{(l+1)}) \underbrace{.* a^{(l)} .* (1-a^{(l)})}_{g'(z^{(l)})}$

        $g'(z^{(l)}) = \underbrace{a^{(l)} .* (1-a^{(l)})}_{\text{calc value.}}$

      5. $\Delta_{i,j}^{(l)} := \Delta_{i,j}^{(l)} + a_j^{(l)} \delta_j^{(l+1)}$ with vectorization, $\quad \Delta^{(l)} := \Delta^{(l)} + \delta^{(l+1)}(a^{(l)})^T$

update: $\begin{cases} D_{i,j}^{(l)} := \frac{1}{m}\Delta_{i,j}^{(l)} + \lambda\theta_{i,j}^{(l)}) \text{ if } j\neq 0. \\ \frac{1}{m}\Delta_{i,j}^{(l)} \text{ if } j=0 \end{cases}$

      $D$ is "accumulator". $\frac{\partial}{\partial\theta_{i,j}^{(l)}} J(\theta) = D_{i,j}^{(l)}$

Implementation details

↳ refer to notes and videos

↳ No need to write code for hand-written notes.

↳ unrolling = you can make / convert between matrix / vector repn of matrices

↳ gradient checking: bug-free impl guarantee

↳ use random to set initial theta

---

Week 6

Evaluating a learning algorithm

Ways to arrive at better hypothesis

↳ more examples

↳ more / less # of features

↳ more / less value of $\lambda$.

to evaluate a hypothesis, we split data into training set & test set.  (70%    (30%

We ↳ learn $\theta$, minimize $J_{train}(\theta)$ using training set

↳ compute test set error $J_{test}(\theta)$

computing test set error

↳ lin. reg. : $J_{test}(\theta) = \frac{1}{2M_{test}} \sum_{i=1}^{M_{test}} (h_\theta(x)_{test}^i - y_{test}^{(i)})^2$

↳ log. reg :

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

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

Model Selection

you can break down data set into three data sets:

training set  ,  cross validation set,  test set

       ↳ 60%       ↳ 20%       ↳ 20%

Idea: test different degree of polynomial, evaluate error function

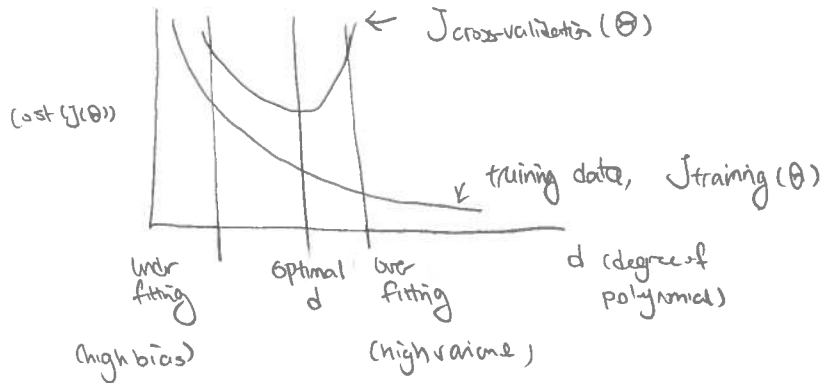1. optimize params in $\theta$ using training set for each degree.

2. find the polynomial degree d that produce least error by cross validation /se.

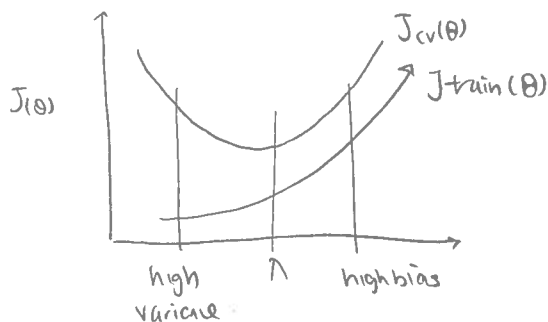3. estimate generalized error with $J_{test}(\theta^{(d)})$ using test set. (d := deg returning /lowest) error

(this way, test set is NOT associated with the param training.

16.

## Diagnosing bias vs Variance.

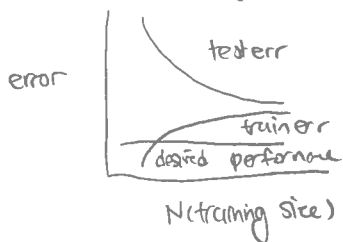If you have bad prediction, you need to figure out whether its high bias or variance.



We use similar algorithm for testing regularization term $\lambda$.
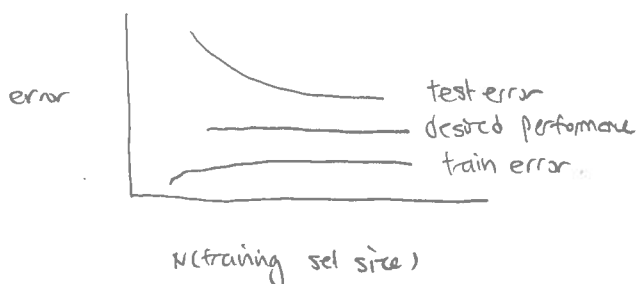
## learning curves
↳ experiencing high bias



↳ low training set size causes $J_{train}(\theta)$ low, $J_{cv}(\theta)$ high.

↳ large training set size causes both $J_{train}(\theta)$, $J_{cv}(\theta)$ high but also $J_{train}(\theta) \approx J_{cv}(\theta)$

↳ getting more data won't help much.

↳ experiencing high variance



↳ low training set size: $J_{train}(\theta)$ low, $J_{cv}(\theta)$ high

↳ large training set size: $J_{train}(\theta)$ increase with set size, and $J_{cv}(\theta)$ continue to decrease without plateauing. $J_{train}(\theta) < J_{cv}(\theta)$, but difference remains significant.

↳ getting more data will likely to help

## Debugging learning algorithm.

| problem | try |
|---------|-----|
| high var | get more training data |
| high var | less features |
| high bias | get more features |
| high bias | add poly features |
| high bias | decrease $\lambda$ |
| high var | increase $\lambda$. |

Small neural network: computationally cheap (prone to underfitting)

large neural network: computationally expensive (prone to overfitting, (use $\lambda$ (regularization) to fix)

## Building a spam classifier

↳ Designing ML system. (building your own system)

↳ Identify features, (X) and classifier (y)

↳ ways to spend more time

    ↳ collect lots of data

    ↳ more sophisticated features

    ↳ algorithms to process input data.

## Error analysis

↳ implement a quick implementation

    ↳ use it to decide how to spend your time

    ↳ plot learning curves, and decide what to do.

    ↳ manually examine errors, analyse

    ↳ implement a metric that returns performance on different change/ideas

### Skewed classes

    ↳ case when one class has very large size, another very little size.

    ↳ different error metric → use tc t, true - & false t, false - to classify

      (precision / recall    precision: true t / pred t  ,* Recall true t / actual t

    can change $h_\theta(X)$ threshold, which trade off precision / recall

    precision metric: $F_1$ score: $PR/(p+R)$

Support vector machine (SVM)

↳ using cost(1), cost(0) similar to $h\theta$, but easier computation wise

↳ ee. if $y=1$



$-\log \frac{1}{1+e^{-z}}$  $\text{cost}_1(\theta)$

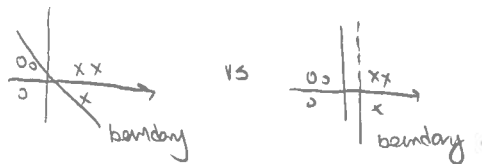$-\log(1-\frac{1}{1+e^{-z}})$  $\text{cost}_0(\theta)$

minimizing $\theta$ given by: optimization projective;

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

↳ large margin classifyer.

↳ allows a decision boundary that stays naturally far apart from dataset.
↳ the perpendicular vector is closest to examples.



boundary    vs    boundary

Kernels

label landmark (defining feature) then use distance measure.

$f_i = \exp\left(-\frac{\|x - \ell^{(i)}\|^2}{2\sigma^2}\right)$    if $\|w\| \approx$ small $f_i \approx 1$
                                                                large $f_i \approx 0$

learn non linear decision boundary



predict 1 if close to $L_1, L_2, L_3$
        0 otherwise.

details (SVM with kernels)

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

choose $\ell^1 = x^1, \ell_2 = x^2, \cdots \ell_m = x^m$

given example $x$, $f_1 = $ similarity $(x, \ell^1)$
                    $f_2 = $ similarity $(x, \ell^2)$

for $i = 1, \cdots m$, $x^{(i)} \Rightarrow \begin{bmatrix} f_1^i \\ f_2^i \\ \vdots \\ f_n^i \end{bmatrix}$    $f = \begin{bmatrix} f_0 \\ f_1 \\ \vdots \\ f_m \end{bmatrix} \to f_0 = 1$

Hypothesis   given $x$, compute features $f \in \mathbb{R}^{m \times 1}$

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

training using $f_{(i)}$'s  similarity metric instead.

$$\min_{\theta} C \sum_{i=1}^{m} y^{(i)} \text{cost}_1 (\theta^T f^{(i)}) + (1-y^{(i)}) \text{cost}_0 (\theta^T f^{(i)}) + \frac{1}{2} \sum_{j=1}^{m} \theta_j^2$$
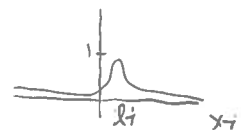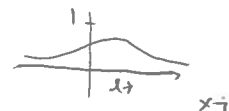
note that since $C = \frac{1}{\lambda}$,  large $C \rightarrow$  low bias, high var   (ie small $\lambda$)

   small $C \rightarrow$  high bias, low var  (ie large $\lambda$)

   large $\sigma^2 \rightarrow$ features $f_{(i)}$ vary smooth,
            high bias, low var

   small $\sigma^2 \rightarrow$ features $f_{(i)}$ vary less smooth,
            low bias, high var

your job is to choose $C$ and $\sigma^2$

## using an SVM

↳ use software libraries: liblinear, libsvm.

   ↳ need to choose: kernel (use or no use) and parameter $C$.
            ↳ linear kernel (no kernel)
            ↳ Gaussian kernel
                 ↳ need to choose $\sigma$
                 ↳ need to do feature scaling
   ↳ for other choices of kernel, it must satisfy Mercer's theorem   so it for sure  do not diverge
↳ multiclassification
      ↳ built in SVM package
      ↳ one-vs-all
↳ logistic regression vs SVM.
   which to choose?
      { $M = $ # features
      { $n = $ # training example.

   SVM → convex function → return global optima
      ↳ if $n$ large relative to $m$ ($n \gg m$,  $n \lesssim 10,000$  $m \approx 10 \sim 1000$)
            ↳ use L.R. or SVM w/ linear kernel.

      ↳ if $n$ small, $m$ intermediate ($n = 1 \sim 1000$,  $m = 10 \sim 1000$)
            ↳ use SVM w/ Gaussian kernel

      ↳ if $n$ small $m$ large ($n = 1 \sim 1000$,  $m = 50,000+$)
            ↳ add more feature, then  use LR or SVM with  lin' kernel

   ↳ nn works well with these settings, but  is slow to train.

### unsupervised learning

↳ gives input has no labels

↳ algorithm finds clustering data. / structure (as example of unsupervised learning algorithm

↳ used for identifying market segmentation, social network analysis, organize computer clusters.
   galaxy formation / astronomical data analysis

### K-means algorithm

↳ for clustering

↳ step 1: initial cluster centroids randomly

{ step 2: (inum loop) assign each point to centroid that's closer to it (assign $c^i$, to index of closest cluster centroid)

keep on iterating { step 3: (ink loop) move the centroids. to new mean of assigned points.   Mean = $\mu_k$

↳ input:   K ( # of clusters )
   training set   $\{x^1, \cdots x^m\}$
   $x^i \in \mathbb{R}^n$

↳ also useful for non-separated clusters.

   It can still separate out clusters, although may not seem like a obvious separation

### Optimization objective

$c^{(i)}$ = index of cluster (1~K)  that $x^{(i)}$ is currently assigned to

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

$\mu_{c^{(i)}}$ = cluster centroid of cluster that $x^{(i)}$ is assigned to.

objective { 
$$J(c^{(1)}, \cdots c^{(m)}, \mu_1, \cdots \mu_k) = \frac{1}{m} \sum_{i=1}^{m} \|x^{(i)} - \mu_{c^{(i)}}\|^2$$

$$\min_{\substack{c^{(1)}, \cdots, c^{(m)} \\ \mu_1, \cdots \mu_k}} J(c^{(1)}, \cdots c^{(m)}, \mu_1, \cdots \mu_k)$$

note that J does not increase ! as a func of iteration.

### random initialize clustering centroid

steps: ↳ set K < m                                              ↳ k-means may not always
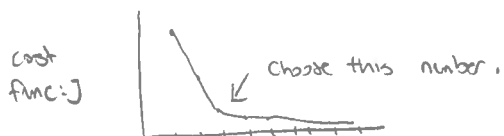   ↳ randomly choose K training examples                        end at global optima.
   ↳ set $\mu_1, \cdots \mu_k$ to these examples.

↳ avoid bad clustering, can initialize many times ! compute J each time & pick one that gives lowest cost.
   when having small # of cluster, multiple/ reinitialization helps the most.

## Pick # of clusters

↳ Choosing # of clusters "elbow method"

cost
func: J



Choose this number.

# clusters.

↳ But often, the 'elbow' does not appear.

↳ don't expect it to work.

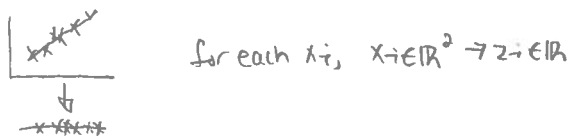↳ better way: pick K from "what you want to do with result of the learning"
   ↳ i.e. how many groups you want? for what purpose are you running this algorithm

## Data compression & dimensionality reduction

↳ example: 2D → 1D

↳ reducing redundant data

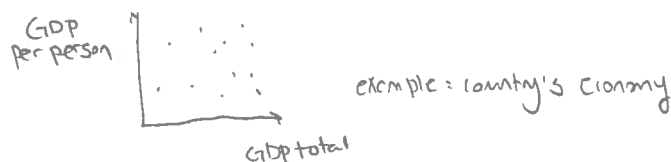↳ represent datapoints in a line by the line (2D→1D) and project each point to the line



for each $x_i$, $x_i \in \mathbb{R}^2 \to z_i \in \mathbb{R}$

example: 3D space roughly on plane → lie on plane

you can generally reduce from high to low dimension given that data roughly lie on like-dimensions

## Visualization & dimensionality reduction.

↳ help with data visualization that looks very complex.

↳ having precision smaller dimension to capture info in more important feature in data

GDP
per person



example: country's economy

GDP total

↳ usually reduce to 2-D or 3-D so it's easy to visualize.

## PCA (Principle component analysis)

↳ goal: find surface of lower dimension that has smallest sum of distance (from actual pt
   ↳ to projected pt)
   ↳ (projection error)

↳ note: PCA is NOT linear regression.
   lin reg minimize vertical distance. (predict y)
   PCA minimize perpendicular distance (nothing to do with y)

## PcA algorithm

↳ training set $= \{x', \cdots x^m\}$.

↳ first, do preprocessing (feature scaling /mean normalization )

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

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

also scale the features

ie. $x_j^{(i)} \leftarrow \frac{x_j^{(i)} - \mu_j}{s_j}$

↳ then, we need to compute the vector $u_i$, and the projections/new representations.

Reduce data from dim $n$ to dim $k$.

covariance matrix: $\Sigma = \frac{1}{m} \sum_{i=1}^{n} (x^{(i)})(x^{(i)})^T$ ⟵ $n \times n$ matrix

eigenvectors of $\Sigma$ : $[U, S, V] = svd(Sigma);$

$U$ will be a $n \times n$ matrix whose columns are $u', u_2, \cdots u^m$.

$$U = \begin{bmatrix} | & | & & | \\ u' & u^2 & \cdots & u^m \\ | & | & & | \end{bmatrix} \in \mathbb{R}^{n \times n}$$

$\underbrace{\qquad}_{k}$

↳ change $x \in \mathbb{R}^n \rightarrow z \in \mathbb{R}^k$

$$U_{redue} = \begin{bmatrix} | & & | \\ u_1 & \cdots & u_k \\ | & & | \end{bmatrix} \qquad z = U_{redue}^T x^{(i)}$$

$\underbrace{\qquad}_{k}$

## reconstruction from compressed representation

$z = U_{redue}^T * x$

$x_{approx} = U_{redue} * z$

## Applying PCA

choosing $k$ (# of principle components)

↳ average square projection error $= \frac{1}{m} \sum_{i=1}^{m} \|x^{(i)} - x_{approx}^{(i)}\|^2$

↳ total variance in data $= \frac{1}{m} \sum_{i=1}^{m} \|x^{(i)}\|^2$

choose $k$ to be smallest value, s.t.

$$\frac{\frac{1}{m} \sum_{i=1}^{m} \|x^{(i)} - x_{approx}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^{m} \|x^{(i)}\|^2} \leq 0.01 \qquad \text{`` 99\% variance is retained ''}$$

## algorithm

$[U, S, V] = svd(sigma)$        S is a diagnal matrix

Pick smallest K s.t.

$$\frac{\sum_{i=1}^{K} S_{ii}}{\sum_{i=1}^{m} S_{ii}} \geq 0.99$$        (0.99 variance retained).

## use PCA to speed up learning algorithm

When running PCA, only run it on the training set.

to speed up:

input: $\{(x^1, y^1), (x^2, y^2), \cdots (x^m, y^m)\}$

extracted input:    $x^1, x^2, \cdots x^m \in \mathbb{R}^{10000}$

$\downarrow$ PCA

$z^1, z^2, \cdots z^m \in \mathbb{R}^{1000}$

then train $\{(z^1, y^1), \cdots (z^m, y^m)\}$    by $h_\theta(z)$

## Summary of Applications of PCA

 └ compression
        └ reduce storage to store data
        └ speed up learning algos.          } chose K by % variance

 └ visualization
        └ K=2 or K=3

 └ Note: DONT use PCA to prevent overfitting.
        └ It's a bad way of using PCA.
        └ Use regularization instead.
 └ note: when designing system
        └ first try without PCA.
        └ only if it does not work, utilize PCA
                └ if you need to save storage space or time.

Anomaly detection

↳ detect abnormally behaving data

↳ use Gaussian distribution (probability distribution)

↳ Density estimation
  ↳ training set $= \{x^1, x^2, \cdots x^m\}$  each training  $X \in \mathbb{R}^n$
  ↳ $P(x) = \prod\limits_{j=1}^{n} P(x_j; u_j; \sigma_j^2)$

Anomaly detection algorithm

↳ 1. choose features $X_i$ that you think might be anomaly.

↳ 2. fit params $u_1, \cdots u_n, \sigma_1^2, \cdots \sigma_n^2, \Rightarrow \mu_j = \frac{1}{m} \sum x_j^{(i)}$   $\sigma_j^2 = \frac{1}{m} \sum (x_j^{(i)} - u_j)^2$

↳ 3. given new example $x$, compute $P(x)$

  $P(x) = \prod\limits_{j=1}^{n} P(x_j; u_j; \sigma_j^2)$      $P(x) < \epsilon$


developing & evaluating an anomaly detection system

↳ use the training set, cross validation set, and test set.

↳ algorithm evaluation:

  fit model $P(x)$ on test set $x^1, \cdots x^m$

  con cv / test example)

  predict $= \begin{cases} 1 & \text{if } P(x) < \epsilon \text{ (anomaly)} \\ 0 & \text{if } P(x) > \epsilon \text{ (normal)} \end{cases}$

↳ can use cv set to choose $\epsilon$

↳ use precision /recall   & $F_1$-score

When to use anomaly detection, when to use supervised learning?

| AD | SL |
|---|---|
| ↳ very small # of positive examples. (y=1) (0~20) large number of negative examples. | ↳ large # of positive /negative examples. |
| | ↳ future anomaly look like previous ones. |
| ↳ many types of anomaly ( difficult for learning) | ↳ enough positive examples for the algorithm to learn. |
| ↳ possible anomaly in future that we've never seen. | |

↳ Multivariate Gaussian distribution can detect abnormal relationship between features.

    ↳ it's more computationally expensive.

    ↳ must have $m > n$ or else $\Sigma$ is singular.

## Recommender systems

    ↳ give recommendations to subscribers to a service.

    ↳ is a high priority to many companies.

    ↳ example- predicting movies.

    ↳ includes:

        ↳ feature learning

        ↳ & content based recommendation

        both are included in videos. (implementation)

## Collaborative filtering algorithm

1. initialize $x^1, \cdots x^{nm}, \theta^1, \cdots \theta^{hu} \in \mathbb{R}^n$ to small values.

2. minimize $J(x^1, \cdots x^{nm}, \theta^1, \cdots \theta^{hu})$ using gradient descent or other advanced optimization algorithm, i.e.

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

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

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

vectorized implementation is in videos.

    "low rank matrix factorization"

Use mean normalization to predict entries with no previous data

# Week 10

## Note
More data, the better algorithm it is.

## Stochastic gradient descent
↳ updating data as data come along.

| Batch gradient descent | vs. | Stochastic gradient descent |

**Batch gradient descent**

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

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

Repeat {

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

for every $j=0, \cdots M$

}

**Stochastic gradient descent**

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

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

① randomly shuffle set
② repeat {

for $i=1, \cdots M$ {

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

for $j=0, \cdots n$

} }

## Mini-batch gradient descent
↳ batch grad. des. use all $m$ examples each iteration.
↳ stochastic grad. des. use 1 example each time (iteration).
↳ mini batch grad. des. use $\underline{b}$ examples each iteration

## Checking for convergence
↳ check trend / convergence over time.

## large scale machine learning
Map reduce & data parallelism.
Map reduce is to split batch gradient descent over diffrent computers for computation then the central machine combines the result.

Example: photo OCR problem.
↳ design a <u>pipeline</u> that goes through multiple step of machine learning algorithm.

Getting lots of data: Artificial data analysis.
↳ distort the data to generate more test examples
↳ thus more data and better ML algorithm.

Ceiling analysis
↳ for each process in the pipeline, take its maximum performance and see how it improves the algorithm. then component wise improve the process by their priorities.