

★ Basic Info

Andrew Ng 👤 lecturer

mixed 👤 level

7/15/2020 👤 Date

By Chaolan Xu

Machine Learning

Week 1 Intro

Introduction

1. what is Machine Learning/(un)supervised
 - 1.1 E,T,P /study gives computer ability to learn without explicitly programmed
 - 1.2 content (supervised, unsupervised, reinforcement learning, recommender system...

Model Representation

1. cost function
2. gradient descent
 - 2.1 intuition
 - 2.2 learning rate alpha (small, large and fixed)
 - 2.3 algorithm (update thetas simultaneously)

Linear Algebra Review

1. scala, singular matrix, etc

Week 3 Logistic Regression

Classification

1. hypothesis
 - 1.1 $g(z)$ /interpretation of hypothesis output
2. decision boundary
 - 2.1 $g(z) \geq 0.5$, z needs to > 0 , i.e. transpose of $\theta^T \cdot x > 0$

Logistic Regression Model

1. cost function and its intuition/gradient descent
2. advanced optimization algorithms (Conjugate Gradient, BFGS, L-BFGS)

Multiclass Classification

1. one-vs-all method

Solving the problem of overfitting

1. reduce number of features
 - 1.1 manually select which features to keep
 - 1.2 model selection algo
2. regularization (reduce magnitude of theta)
 - 2.1 regularize linear regression
 - 2.1.1 gradient descent
 - 2.1.2 normal equation
 - 2.2 regularize logistic regression
 - 2.2.1 gradient descent
 - 2.2.2 advanced optimization

Week 2 Linear Regression

Multivariate Linear Regression

1. cost function, gradient descent
 - 1.1 how to make sure gd is working correctly (declare convergence if cost function decreases by less than 0.001 in one iteration)
2. feature scaling
 - 2.1 about $-3 < x_i < 3$
 - 2.2 mean normalization
 3. polynomial regression
- 3.1 create new features by yourself ($x_2x_3, x_2^2x_3$) (pay attention to feature scaling)

Normal Equation (another way to solve for theta)

1. $\theta = (X^T X)^{-1} X^T y$
2. no need to feature scaling
3. compare with gradient descent
 - 3.1 no need alpha/iteration; slow when n is large
4. noninvertibility
 - 4.1 the 'pinv' function will give you a value of θ even if $X^T X$ is not invertible.
 - 4.2 consider if there are redundant features
 - 4.3 ⚠ may be too many features—delete some features or use regularization

Octave/Matlab

Week 4 Neural Networks 1

Motivation

1. 1. when there are lots of features
2. 2. car recognition (computer vision)
3. 3. back ground: neurons and brain

Neural Networks

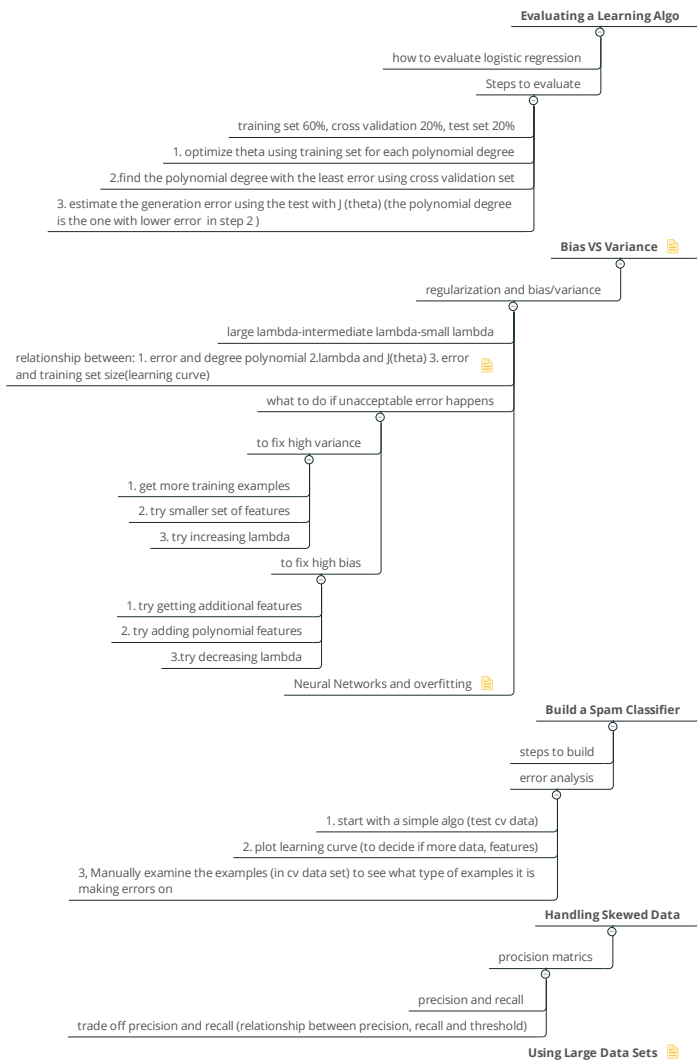
1. Model Representation 📄

Applications

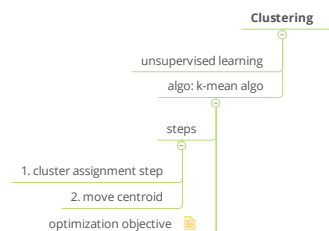
1. x_1 XOR(x_1 NOR) x_2
2. x_1 AND(OR) x_2 (have example)
3. (NOT x_1) AND (NOT x_2) (=1 if and only if $x_1=0=x_2$)
4. Multiclass Classification (with y having more columns)

Week 5 Neural Networks 2

Week 6 Evaluating



Week 8 Clustering & P C A



Neural Networks' Cost Function and Backpropagation

BackPropagation Algo

Steps

1. set $a(1) = x(i)$
2. perform forward propagation to compute $a(l)$
3. using $y(i)$, compute $\delta(L) = a(L) - y(i)$
4. compute delta of every level
5. get D_{ij} (the derivative of cost function $J(\theta)$)

remarks : to use `fminunc()` in Matlab, need to unroll the theta to one long vector/
theta needs to randomly initialized

math for back propagation (learn beyond this lecture by myself)

Gradient Checking

Six Steps for Neural Networks

1. Randomly initialize weights
2. implement forward propagation to get $h(x_i)$ for any $x(i)$
3. compute cost function
4. implement backpropagation to compute partial derivatives
5. gradient checking
6. use gradient decent or advanced optimization method with backprop to try to minimize cost function (a function of parameters)

Week 7 SVM

Large Margin Classification

Cost Function

inferred from logistic regression

By convention, change the regularized function from $A + \lambda B$ (for logistic regression) to $C * A + B$ (for SVM)

Large Margin Intuition

discuss with setting C (the C in optimization objective $C * A + B$) to be a very large number

Math behind (to know why SVM can choose the best classifier (the one has large margin))

Kernels (for non-linear boundary)

choosing the landmarks

every x_i will get a "new feature" vector, the new features are generated by similarity function $(x(i), \text{landmarks}(x(1) \text{ to } x(m)))$
predict $y=1$ if $(\theta)T f > 0$

SVM parameters

C : large C , lower bias & high variance

sigma for Gaussian Kernels: large sigma, features f vary more smoothly, will cause higher bias and lower variance

Other choices for kernels

polynomial kernel, string kernel, chi-square kernel and etc.

SVM in Practice

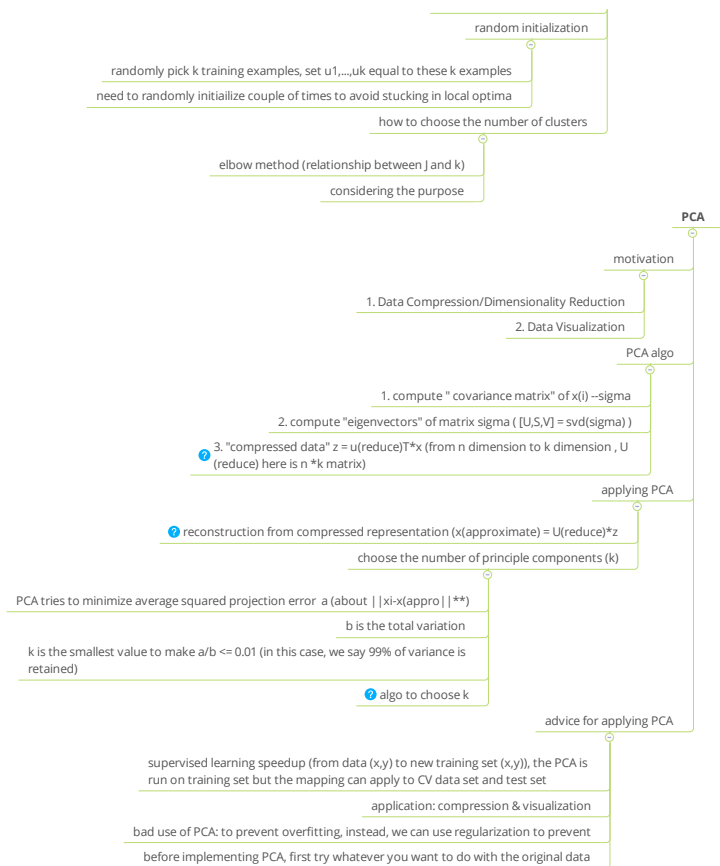
perform feature scaling before using the Gaussian Kernels

multi-class classification

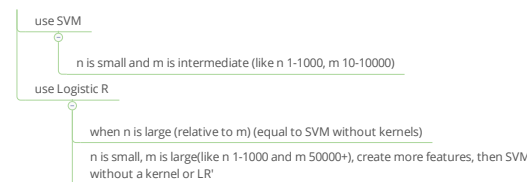
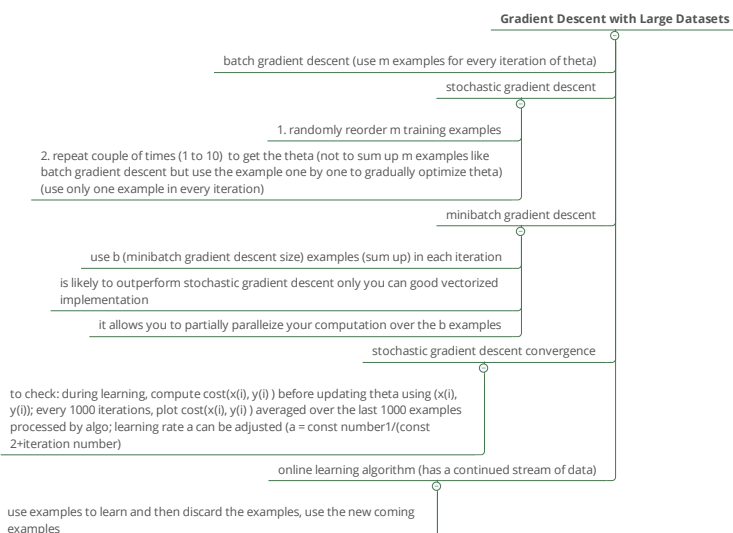
use built-in multi-class classification functionality

use one-vs-all method, train k SVMs to distinguish i from the rest, pick class i with largest $(\theta)T x$

SVM VS Logistic Regression



Week 10 Large Scale Machine Learning



Week 9 Anomaly Detection & Predicting Rating (Collaborative Filtering)

Anomaly Detection

one dimension Gaussian Distribution

algo steps

1. choosing features x_i that you think might be indicative of anomalous examples
2. fit parameters (use data set to find the mean and sigma for every feature)
3. given new example x , compute $p(x)$, anomaly if $p(x) < \text{epsilon}$ (mean and sigma use the ones fitted from data set)

build an anomaly detection system

1. fit model on training set
2. on cross validation/test examples, predict normal(1) or anomaly(0) based on $p(x) < \text{or} \geq \text{epsilon}$ (CV set can also be used to find epsilon)
3. possible evaluation matrix (classification accuracy is not a good indicator--why? --skewed data, always predicting 0 will have high classification accuracy)

anomaly detection VS supervised learning

use anomaly detection when:

- many normal examples
- many different types of anomalies, hard for any algo to learn from positive examples ($y=1$) what the anomalies may look like

use supervised learning when:

- large number of positive and negative examples
- enough positive examples ($y=1$, i.e. anomalies) for algo to get a sense of what positive examples likely to be similar to ones in training set

choosing what feature to use

1. make data more like Gaussian Distri. e.g. use $\log()$
2. most common problem: $p(x)$ is comparable for normal and anomalous examples----> create new feature to capture the anomalies

multivariate Gaussian Distribution

motivation and intuition

anomaly detection using the multivariate Gaussian Distribution

steps

1. fit model with data set
2. give a new example, compute $p(x)$ --flag anomaly if $p(x) < \text{epsilon}$

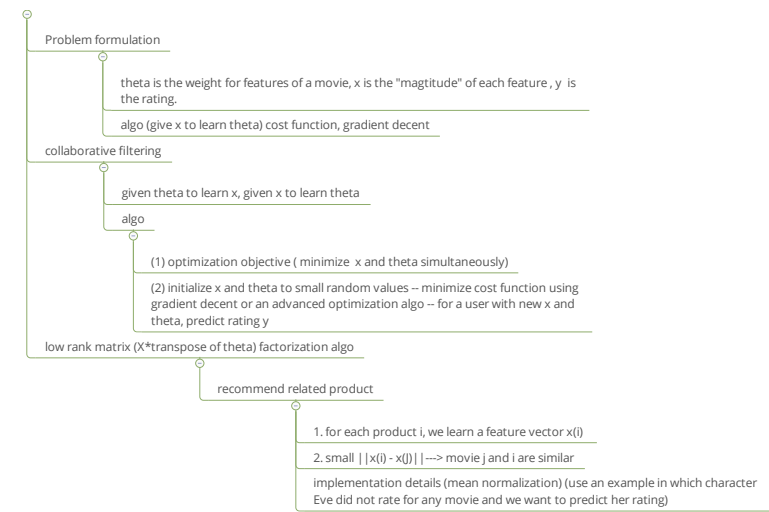
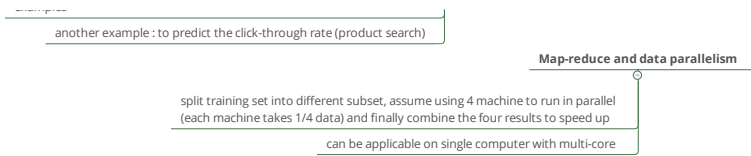
relationship to original model

1. the original model is the special case of multivariate model (with a special sigma matrix)
2. difference

original: need to manually create features to capture anomalies where x_1, x_2 take unusual combinations of values; ok even if m is small; computationally cheaper

multivariate: automatically captures correlations between features; computationally more expensive; must have $m > n$ (otherwise, sigma matrix will be non-invertible) (if sigma matrix is non-invertible, check redundant features)

Predicting Movie Ratings



Week 11 Application Examples

