Lecture Slides for

Machine Learning 2nd Edition



ETHEM ALPAYDIN, modified by Leonardo Bobadilla and some parts from http://www.cs.tau.ac.il/~apartzin/MachineLearning/ © The MIT Press, 2010

alpaydin@boun.edu.tr http://www.cmpe.boun.edu.tr/~ethem/i2m

Outline

Previous class Multivariate Data Parameter Estimation Estimation of Missing Values Multivariate Classification

This class: Ch 5: Multivariate Methods

- Discrete Features
- Multivariate Regression
- Ch 6: Dimensionality reduction

CHAPTER 5: Multivariate Methods

Model Selection





Based on E Alpaydın 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Discrete Features • Binary features: $p_{ij} \equiv p(x_j=1|C_i)$

if x_j are independent (Naive Bayes')

$$p(\mathbf{x} | \mathbf{C}_i) = \prod_{j=1}^{d} p_{ij}^{\mathbf{x}_j} (1 - p_{ij})^{(1 - \mathbf{x}_j)}$$

the discriminant is linear $g_{i}(\mathbf{x}) = \log p(\mathbf{x} | C_{i}) + \log P(C_{i})$ $= \sum_{j} [x_{j} \log p_{ij} + (1 - x_{j}) \log(1 - p_{ij})] + \log P(C_{i})$ Estimated parameters $\hat{p}_{ij} = \frac{\sum_{t} x_{j}^{t} r_{i}^{t}}{\sum_{t} r_{i}^{t}}$

Multivariate Regression

$$r^{t} = g(x^{t} | w_{0}, w_{1}, ..., w_{d}) + \varepsilon$$

Multivariate linear model

•
$$W_0 + W_1 X_1^t + W_2 X_2^t + \dots + W_d X_d^t$$

• $E(W_0, W_1, \dots, W_d \mid X) = \frac{1}{2} \sum_t [r^t - W_0 - W_1 X_1^t - \dots - W_d X_d^t]^2$

Multivariate Regression

$$w_0 + w_1 x_1^t + w_2 x_2^t + \dots + w_d x_d^t$$

 $I E(w_0, w_1, \dots, w_d \mid X) = \frac{1}{2} \sum_t [r^t - w_0 - w_1 x_1^t - \dots - w_d x_d^t]^2$

 $\sum_{t} r^{t} = Nw_{0} + w_{1} \sum_{t} x_{1}^{t} + w_{2} \sum_{t} x_{2}^{t} + \dots + w_{d} \sum_{t} x_{d}^{t}$ $\sum_{t} x_{1}^{t} r^{t} = w_{0} \sum_{t} x_{1}^{t} + w_{1} \sum_{t} (x_{1}^{t})^{2} + w_{2} \sum_{t} x_{1}^{t} x_{2}^{t} + \dots + w_{d} \sum_{t} x_{1}^{t} x_{d}^{t}$ $\sum_{t} x_{2}^{t} r^{t} = w_{0} \sum_{t} x_{2}^{t} + w_{1} \sum_{t} x_{1}^{t} x_{2}^{t} + w_{2} \sum_{t} (x_{2}^{t})^{2} + \dots + w_{d} \sum_{t} x_{2}^{t} x_{d}^{t}$ \vdots $\sum_{t} x_{d}^{t} r^{t} = w_{0} \sum_{t} x_{d}^{t} + w_{1} \sum_{t} x_{d}^{t} x_{1}^{t} + w_{2} \sum_{t} x_{d}^{t} x_{2}^{t} + \dots + w_{d} \sum_{t} (x_{d}^{t})^{2}$

CHAPTER 6: Dimensionality Reduction

Dimensionality of input

- Number of Observables (e.g. age and income)
- If number of observables is increased
 - More time to compute
 - More memory to store inputs and intermediate results
 - More complicated explanations (knowledge from learning)
 - Regression from 100 vs. 2 parameters
 - No simple visualization
 - 2D vs. 10D graph
 - Need much more data (curse of dimensionality)
 - 1M of 1-d inputs is not equal to 1 input of dimension 1M

Dimensionality reduction

- Some features (dimensions) bear little or nor useful information (e.g. color of hair for a car selection)
 - Can drop some features
 - Have to estimate which features can be dropped from data

Several features can be combined together without loss or even with gain of information (e.g. income of all family members for loan application)

- Some features can be combined together
- Have to estimate which features to combine from data

Feature Selection vs Extraction

- Feature selection: Choosing k<d important features, ignoring the remaining d k
 - Subset selection algorithms
- Feature extraction: Project the original x_i, i
 =1,...,d dimensions to new k<d dimensions, z_j
 , j =1,...,k
 - Principal Components Analysis (PCA)
 - Factor Analysis (FA)
 - Linear Discriminant Analysis (LDA)

Usage

- Have data of dimension d
- Reduce dimensionality to k<d
 - Discard unimportant features
 - Combine several features in one
- Use resulting k-dimensional data set for
 - Learning for classification problem (e.g. parameters of probabilities P(x|C)
 - Learning for regression problem (e.g. parameters for model $y=g(x|\theta)$

"Goodness" of feature set

- Supervised
 - Train using selected subset
 - Estimate error on validation data set
- Unsupervised
 - Look at input only(e.g. age, income and savings)
 - Select subset of 2 that bear most of the information about the person

Mutual Information

- Have a 3 random variables(features) X,Y,Z and have to select 2 which gives most information
- If X and Y are "correlated" then much of the information about of Y is already in X
- Make sense to select features which are "uncorrelated"

Subset Selection

- There are 2^d subsets of *d* features
- Forward search: Add the best feature at each step
 - Set of features *F* initially Ø.
 - At each iteration, find the best new feature
 - $j = \operatorname{argmin}_i E (F \cup x_i)$

- Add x_j to F if $E(F \cup x_j) < E(F)$

O(d^2) algorithm does not guarantee optimal
 Backward search: Start with all features and remove one at a time, if possible.

Subset Selection

Backward search: Start with all features and remove one at a time, if possible.

- Set of features F
- At each iteration, remove a feature that does not decrease the error
 - $j = \operatorname{argmin} i E (F x_i)$
 - Remove x_j to F if $E(F x_j) < E(F)$

O(d^2) algorithm does not guarantee optimal

Subset-selection (recap)

- Forward search
 - Start from empty set of features
 - Try each of remaining features
 - Estimate classification/regression error for adding specific feature
 - Select feature that gives maximum improvement in validation error
 - Stop when no significant improvement
- Backward search
 - Start with original set of size d
 - Drop features with smallest impact on error

Feature Extraction

- Face recognition problem
 - Training data input: pairs of Image + Label(name)
 - Classifier input: Image
 - Classifier output: Label(Name)
- Image: Matrix of 256X256=65536 values in range 0..256
- Each pixels bear little information so can't select 100 best ones
- Average of pixels around specific positions may give an indication about an eye color. Based on E Alpaydin 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Projection

• Find a projection matrix w from d-dimensional to k-dimensional vectors that keeps error low

$$z = w^T x$$

Based on E Alpaydın 2004 Introduction to Machine Learning © The MIT Press (V1.1)

PCA: Motivation

- Assume that d observables are linear combination of k<d vectors
- We would like to work with basis as it has lesser dimension and have all(almost) required information
- What we expect from such basis
 - Uncorrelated or otherwise can be reduced further
 - Have large variance or otherwise bear no information

PCA: Motivation



Based on E Alpaydın 2004 Introduction to Machine Learning © The MIT Press (V1.1)

PCA: Motivation

- Choose directions such that a total variance of data will be maximum
 - Maximize Total Variance
- Choose directions that are orthogonal
 - Minimize correlation
- Choose k<d orthogonal directions which maximize total variance

Based on E Alpaydın 2004 Introduction to Machine Learning © The MIT Press (V1.1)

PCA

- Choosing only directions: $\|\boldsymbol{w}_1\| = 1$
- $z_1 = \boldsymbol{w}_1^T \boldsymbol{x}$ Cov $(\boldsymbol{x}) = \boldsymbol{\Sigma}$, Var $(z_1) = \boldsymbol{w}_1^T \boldsymbol{\Sigma} \boldsymbol{w}_1$
- Maximize variance subject to a constrain using Lagrange Multipliers

$$\max_{\boldsymbol{w}_1} \boldsymbol{w}_1^T \boldsymbol{\Sigma} \boldsymbol{w}_1 - \boldsymbol{\alpha} (\boldsymbol{w}_1^T \boldsymbol{w}_1 - 1)$$

• Taking Derivatives

 $2\Sigma w_1 - 2\alpha w_1 = 0 \qquad \Sigma w_1 = \alpha w_1$

• Eigenvector. Since want to maximize $w_1^T \Sigma w_1 = \alpha w_1^T w_1 = \alpha$ we should choose an eigenvector with largest eigenvalue

Based on E Alpaydın 2004 Introduction to Machine Learning $\ensuremath{\mathbb{C}}$ The MIT Press (V1.1)

PCA

- d-dimensional feature space
- d by d symmetric covariance matrix estimated from samples $C_{OV}(x) = \Sigma$,
- Select k largest eigenvalue of the covariance matrix and associated k eigenvectors
- The first eigenvector will be a direction with largest variance

What PCA does

 $\boldsymbol{z} = \boldsymbol{W}^{T}(\boldsymbol{x} - \boldsymbol{m})$

where the columns of **W** are the eigenvectors of Σ , and *m* is sample mean

Centers the data at the origin and rotates the axes



How to choose k?

• Proportion of Variance (PoV) explained

$$\lambda_1 + \lambda_2 + \dots + \lambda_k$$
$$\lambda_1 + \lambda_2 + \dots + \lambda_k + \dots + \lambda_d$$

when λ_i are sorted in descending order

- Typically, stop at PoV>0.9
- Scree graph plots of PoV vs k, stop at "elbow"

Lecture Notes for E Alpaydın 2004 Introduction to Machine Learning © The MIT Press (V1.1)





PCA

ullet

- Can take into account classes : Karhuned-Loeve Expansion
 - Estimate Covariance Per Class
 - Take average weighted by prior
- Common Principle Components
 - Assume all classes have same eigenvectors (directions) but different variances

PCA

- PCA is unsupervised (does not take into account class information)
- Does not try to explain noise
 - Large noise can become new dimension/largest
 PC
- Interested in resulting uncorrelated variables which explain large portion of **total** sample variance

Sometimes interested in explained shared variance (common factors) that affect data

Factor Analysis

- Assume set of unobservable ("latent") variables
- Goal: Characterize dependency among observables using latent variables
- Suppose group of variables having large correlation among themselves and small correlation with other variables
- Single factor?

Factor Analysis

- Assume k input factors (latent unobservable) variables generating d observables
- Assume all variations in observable variables are due to latent or noise (with unknown variance)
- Find transformation from unobservable to observables which explain the data

Based on E Alpaydın 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Factor Analysis

• Find a small number of factors *z*, which when combined generate *x* :

 $X_i - \mu_i = V_{i1}Z_1 + V_{i2}Z_2 + \dots + V_{ik}Z_k + \varepsilon_i$

where z_j , j = 1, ..., k are the latent factors with $E[z_i]=0$, $Var(z_i)=1$, $Cov(z_i, z_i)=0$, $i \neq j$,

 ε_i are the noise sources E[ε_i]= ψ_i , Cov(ε_i , ε_j) =0, $i \neq j$, Cov(ε_i , z_j) =0

,

and v_{ij} are the factor loadings $x - \mu = Vz + \epsilon$

Lecture Notes for E Alpaydın 2004 Introduction to Machine Learning © The MIT Press (V1.1)



Factor Analysis

 In FA, factors z_j are stretched, rotated and translated to generate x



Lecture Notes for E Alpaydın 2004 Introduction to Machine Learning © The MIT Press (V1.1)

FA Usage

- Speech is a function of position of small number of articulators (lungs, lips, tongue)
- Factor analysis: go from signal space (4000 points for 500ms) to articulation space (20 points)
- Classify speech (assign text label) by 20 points
- Speech Compression: send 20 values

Linear Discriminant Analysis

 Find a low-dimensional space such that when x is projected, classes are well-separated



Based on E Alpaydın 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Means and Scatter after projection

$$m_1 = \frac{\sum_t w^T x^t r^t}{\sum_t r^t} = w^T m_1$$

$$m_2 = \frac{\sum_t w^T x^t (1 - r^t)}{\sum_t (1 - r^t)} = w^T m_2$$

$$s_{1}^{2} = \sum_{t} (\mathbf{w}^{T} \mathbf{x}^{t} - m_{1})^{2} r^{t}$$

$$s_{2}^{2} = \sum_{t} (\mathbf{w}^{T} \mathbf{x}^{t} - m_{2})^{2} (1 - r^{t})$$

Based on E Alpaydın 2004 Introduction to Machine Learning © The MIT Press (V1.1)

Good Projection

- Means are far away as possible
- Scatter is small as possible
- Fisher Linear Discriminant

$$J(\mathbf{w}) = \frac{(m_1 - m_2)^2}{s_1^2 + s_2^2}$$





Summary

- Feature selection
 - Supervised: drop features which don't introduce large errors (validation set)
 - Unsupervised: keep only uncorrelated features (drop features that don't add much information)
- Feature extraction
 - Linearly combine feature into smaller set of features
 - Supervised
 - PCA: explain most of the total variability
 - FA: explain most of the common variability
 - Unsupervised
 - LDA: best separate class instances