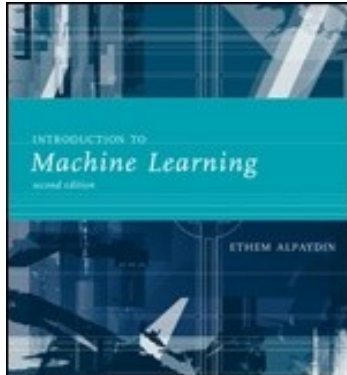


Lecture Slides for

INTRODUCTION TO

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

This class: Ch 5: Multivariate Methods

- Multivariate Data
- Parameter Estimation
- Estimation of Missing Values
- Multivariate Classification

CHAPTER 5:

Multivariate Methods

Multivariate Distribution

- Assume all members of class came from joint distribution
- Can learn distributions from data $P(x|C)$
- Assign new instance for most probable class $P(C|x)$ using Bayes rule
- An instance described by a vector of correlated parameters
- Realm of multivariate distributions
- Multivariate normal

Multivariate Data

- Multiple measurements (sensors)
- d inputs/features/attributes: d -variate
- N instances/observations/examples

$$\mathbf{X} = \begin{bmatrix} X_1^1 & X_2^1 & \dots & X_d^1 \\ X_1^2 & X_2^2 & \dots & X_d^2 \\ \vdots & & & \\ X_1^N & X_2^N & \dots & X_d^N \end{bmatrix}$$

Multivariate Parameters

$$\text{Mean : } E[\mathbf{X}] = \boldsymbol{\mu} = [\mu_1, \dots, \mu_d]^T$$

$$\text{Covariance: } \sigma_{ij} \equiv \text{Cov}(X_i, X_j)$$

$$\text{Correlation : } \text{Corr}(X_i, X_j) \equiv \rho_{ij} = \frac{\sigma_{ij}}{\sigma_i \sigma_j}$$

$$\Sigma \equiv \text{Cov}(\mathbf{X}) = E[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^T] = \begin{bmatrix} \sigma_1^2 & \sigma_{12} & \cdots & \sigma_{1d} \\ \sigma_{21} & \sigma_2^2 & \cdots & \sigma_{2d} \\ \vdots & & & \\ \sigma_{d1} & \sigma_{d2} & \cdots & \sigma_d^2 \end{bmatrix}$$

Parameter Estimation

Sample mean \mathbf{m} : $m_i = \frac{\sum_{t=1}^N x_i^t}{N}, i = 1, \dots, d$

Covariance matrix \mathbf{S} : $s_{ij} = \frac{\sum_{t=1}^N (x_i^t - m_i)(x_j^t - m_j)}{N}$

Correlation matrix \mathbf{R} : $r_{ij} = \frac{s_{ij}}{s_i s_j}$

Estimation of Missing Values

- What to do if certain instances have missing attributes?
- Ignore those instances: not a good idea if the sample is small
- Use 'missing' as an attribute: may give information
- **Imputation**: Fill in the missing value
 - Mean imputation: Use the most likely value (e.g., mean)
 - Imputation by regression: Predict based on other attributes

Multivariate Normal

- Have d-attributes
- Often can assume each one distributed normally
- Attributes might be dependant/correlated
- Joint distribution of correlated several variables
 - $P(X_1=x_1, X_2=x_2, \dots X_d=x_d)=?$
 - X_i is normally distributed with mean μ_i and variance σ_i

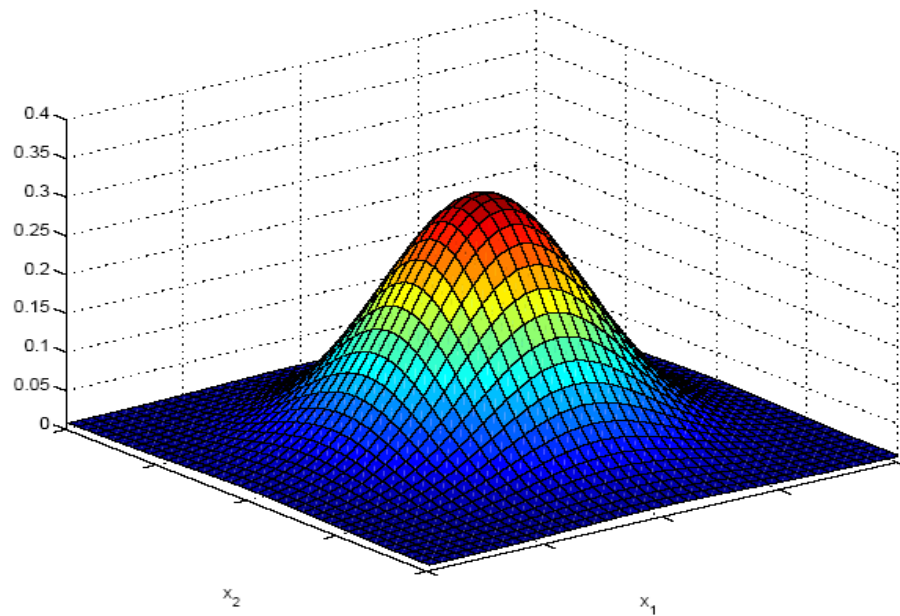
Multivariate Normal

$$\mathbf{x} \sim N_d(\boldsymbol{\mu}, \Sigma)$$

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu})\right]$$

- Mahalanobis distance: $(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})$
- 2 variables are correlated
- Divided by inverse of covariance (large)
- Contribute less to Mahalanobis distance
- Contribute more to the probability

Bivariate Normal



$$\Sigma = \begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix}$$

$$p(x_1, x_2) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \exp \left[-\frac{1}{2(1-\rho^2)} (z_1^2 - 2\rho z_1 z_2 + z_2^2) \right]$$

Multivariate Normal Distribution

- Mahalanobis distance: $(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})$

measures the distance from \mathbf{x} to $\boldsymbol{\mu}$ in terms of $\boldsymbol{\Sigma}$
(normalizes for difference in variances and correlations)

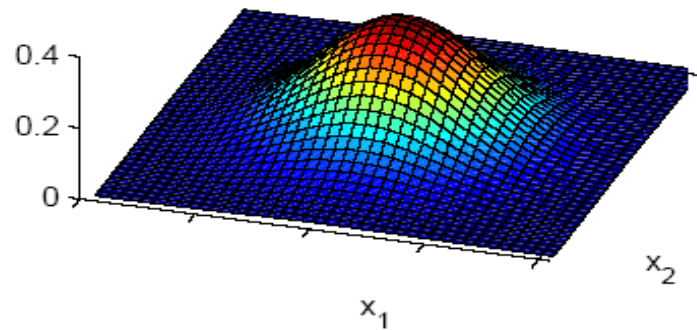
- Bivariate: $d = 2$ $\boldsymbol{\Sigma} = \begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix}$

$$p(x_1, x_2) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \exp\left[-\frac{1}{2(1-\rho^2)}(z_1^2 - 2\rho z_1 z_2 + z_2^2)\right]$$

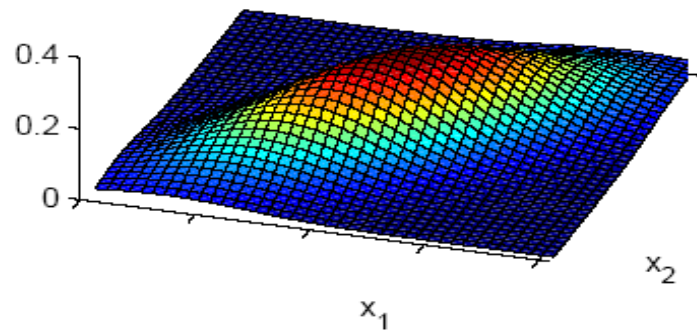
$$z_i = (x_i - \mu_i) / \sigma_i$$

Bivariate Normal

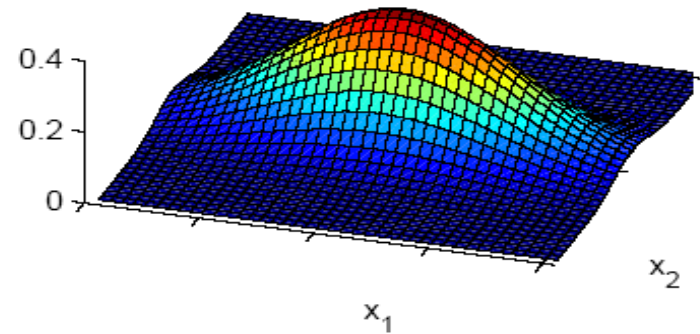
$$\text{Cov}(x_1, x_2) = 0, \text{Var}(x_1) = \text{Var}(x_2)$$



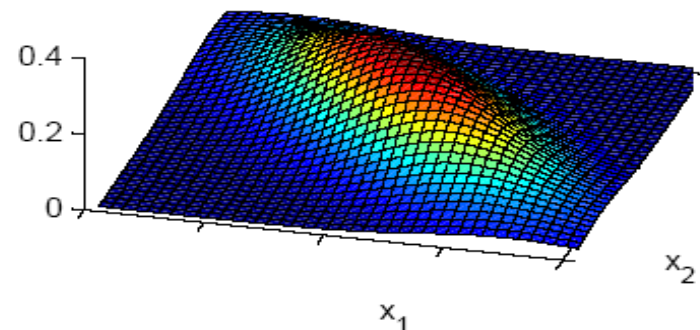
$$\text{Cov}(x_1, x_2) > 0$$



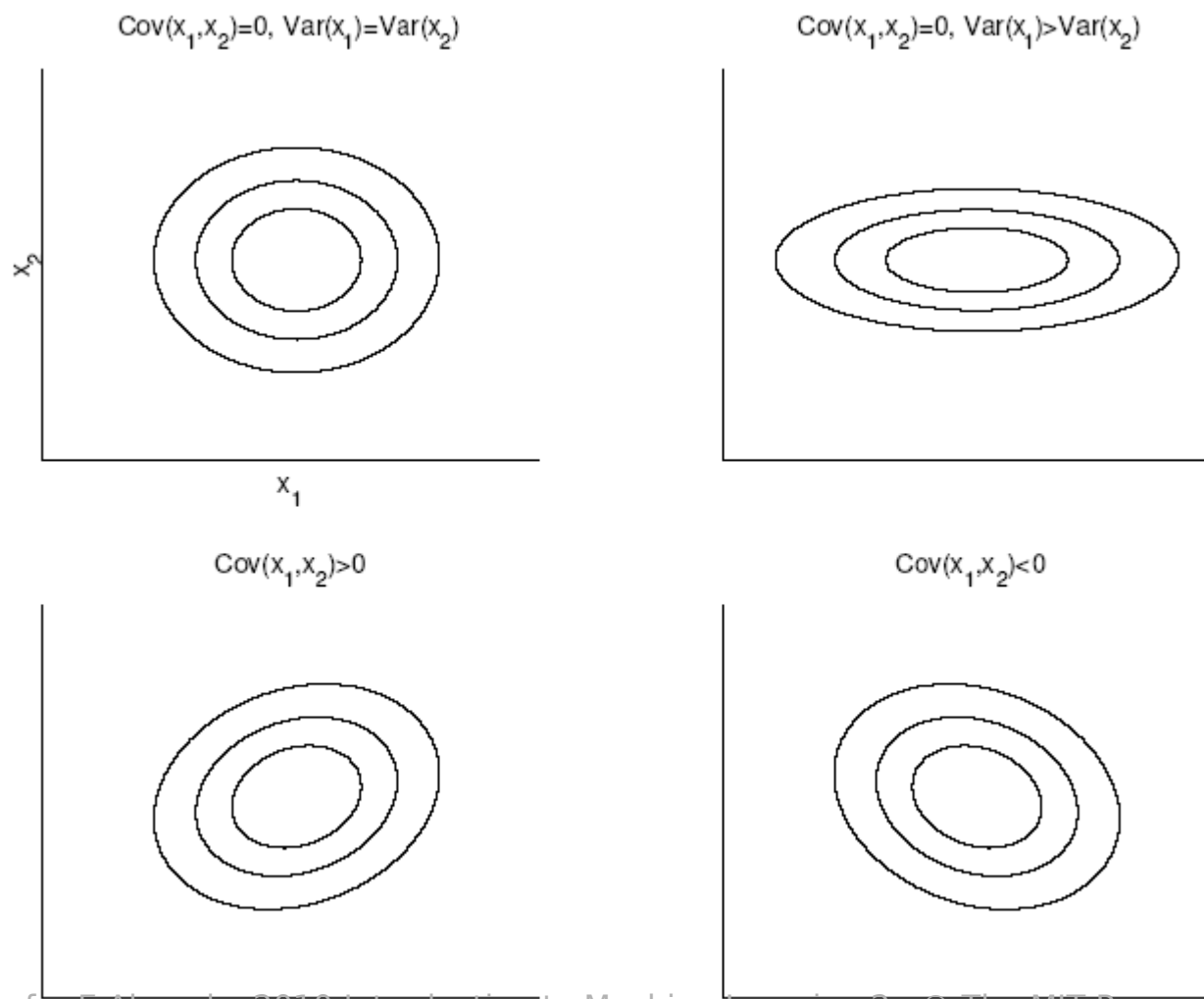
$$\text{Cov}(x_1, x_2) = 0, \text{Var}(x_1) > \text{Var}(x_2)$$



$$\text{Cov}(x_1, x_2) < 0$$



Bivariate Normal



Independent Inputs: Naive Bayes

- If x_i are independent, offdiagonals of Σ are 0, Mahalanobis distance reduces to weighted (by $1/\sigma_i$) Euclidean distance:

$$p(x) = \prod_{i=1}^d p_i(x_i) = \frac{1}{(2\pi)^{d/2} \prod_{i=1}^d \sigma_i} \exp \left[-\frac{1}{2} \sum_{i=1}^d \left(\frac{x_i - \mu_i}{\sigma_i} \right)^2 \right]$$

- If variances are also equal, reduces to Euclidean distance

Projection Distribution

- Example: vector of 3 features
- Multivariate normal distribution
- Projection to 2 dimensional space (e.g. XY plane) Vectors of 2 features
- Projection are also multivariate normal distribution
- Projection of d-dimensional normal to k-dimensional space is k-dimensional normal

$$W^T \mathbf{x} \sim \mathcal{N}_k(W^T \boldsymbol{\mu}, W^T \Sigma W) \quad W \text{ is a } d \times k \text{ matrix}$$

1D projection

$$\mathbf{w}^T \mathbf{x} = w_1 x_1 + w_2 x_2 + \cdots + w_d x_d \sim \mathcal{N}(\mathbf{w}^T \boldsymbol{\mu}, \mathbf{w}^T \boldsymbol{\Sigma} \mathbf{w})$$

$$E[\mathbf{w}^T \mathbf{x}] = \mathbf{w}^T E[\mathbf{x}] = \mathbf{w}^T \boldsymbol{\mu}$$

$$\begin{aligned} \text{Var}(\mathbf{w}^T \mathbf{x}) &= E[(\mathbf{w}^T \mathbf{x} - \mathbf{w}^T \boldsymbol{\mu})^2] = E[(\mathbf{w}^T \mathbf{x} - \mathbf{w}^T \boldsymbol{\mu})(\mathbf{w}^T \mathbf{x} - \mathbf{w}^T \boldsymbol{\mu})] \\ &= E[\mathbf{w}^T (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{w}] = \mathbf{w}^T E[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T] \mathbf{w} \\ &= \mathbf{w}^T \boldsymbol{\Sigma} \mathbf{w} \end{aligned}$$

Multivariate Classification

- Assume members of class from a single multivariate distribution
- Multivariate normal is a good choice
 - Easy to analyze
 - Model many natural phenomena
 - Model a class as having single prototype source (mean) slightly randomly changed

Example

- Matching cars to customers
- Each cat defines a class of matching customers
- Customers described by (age, income)
- There is a correlation between age and income
- Assume each class is multivariate normal
- Need to learn $P(x|C)$ from data
- Use Bayes to compute $P(C|x)$

Parametric Classification

- If $p(\mathbf{x} | C_i) \sim N(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$

$$p(\mathbf{x} | C_i) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}_i|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1}(\mathbf{x} - \boldsymbol{\mu}_i)\right]$$

- Discriminant functions are

$$\begin{aligned} g_i(\mathbf{x}) &= \log P(C_i | \mathbf{x}) = \log \frac{P(\mathbf{x} | C_i) P(C_i)}{P(\mathbf{x})} = \log p(\mathbf{x} | C_i) + \log P(C_i) - \log P(\mathbf{x}) \\ &= -\frac{d}{2} \log 2\pi - \frac{1}{2} \log |\boldsymbol{\Sigma}_i| - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1}(\mathbf{x} - \boldsymbol{\mu}_i) + \log P(C_i) - \log P(\mathbf{x}) \end{aligned}$$

- Need to know Covariance Matrix and mean to compute discriminant functions.
- Can ignore $P(\mathbf{x})$ as the same for all classes

Estimation of Parameters

$$\hat{P}(C_i) = \frac{\sum_t r_i^t}{N}$$

$$\mathbf{m}_i = \frac{\sum_t r_i^t \mathbf{x}^t}{\sum_t r_i^t}$$

$$\mathbf{S}_i = \frac{\sum_t r_i^t (\mathbf{x}^t - \mathbf{m}_i)(\mathbf{x}^t - \mathbf{m}_i)^T}{\sum_t r_i^t}$$

$$g_i(\mathbf{x}) = -\frac{1}{2} \log |\mathbf{S}_i| - \frac{1}{2} (\mathbf{x} - \mathbf{m}_i)^T \mathbf{S}_i^{-1} (\mathbf{x} - \mathbf{m}_i) + \log \hat{P}(C_i)$$

Covariance Matrix per Class

- Quadratic discriminant

$$g_i(\mathbf{x}) = -\frac{1}{2} \log |\mathbf{S}_i| - \frac{1}{2} (\mathbf{x}^T \mathbf{S}_i^{-1} \mathbf{x} - 2\mathbf{x}^T \mathbf{S}_i^{-1} \mathbf{m}_i + \mathbf{m}_i^T \mathbf{S}_i^{-1} \mathbf{m}_i) + \log \hat{P}(C_i)$$

$$= \mathbf{x}^T \mathbf{W}_i \mathbf{x} + \mathbf{w}_i^T \mathbf{x} + w_{i0}$$

where

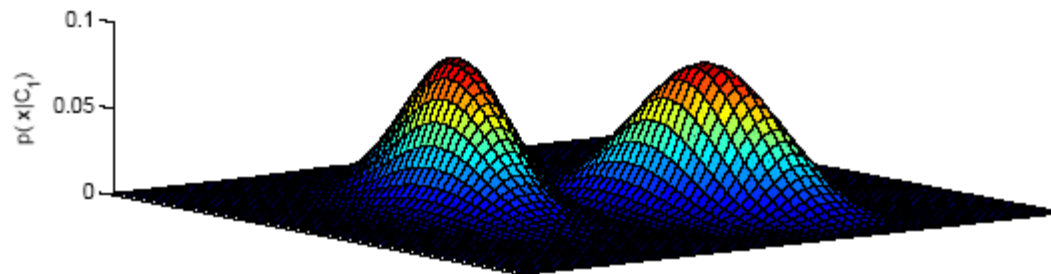
$$\mathbf{W}_i = -\frac{1}{2} \mathbf{S}_i^{-1}$$

$$\mathbf{w}_i = \mathbf{S}_i^{-1} \mathbf{m}_i$$

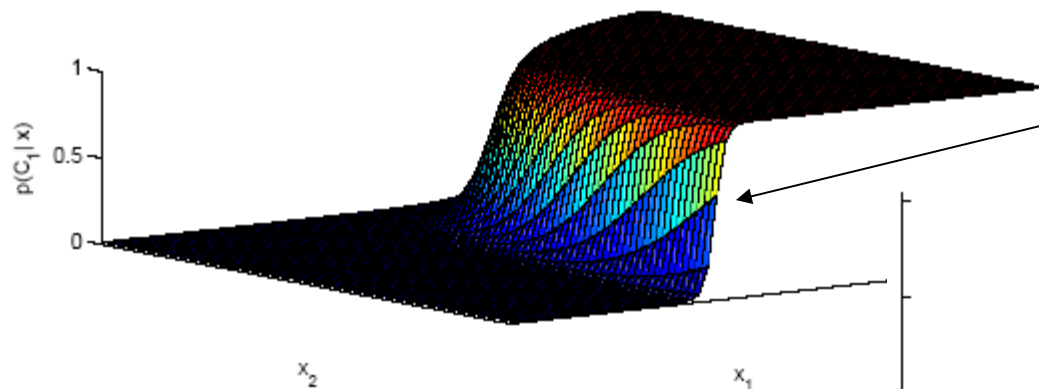
$$w_{i0} = -\frac{1}{2} \mathbf{m}_i^T \mathbf{S}_i^{-1} \mathbf{m}_i - \frac{1}{2} \log |\mathbf{S}_i| + \log \hat{P}(C_i)$$

- Requires estimation of $K \cdot d \cdot (d+1)/2$ parameters for covariance matrix

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

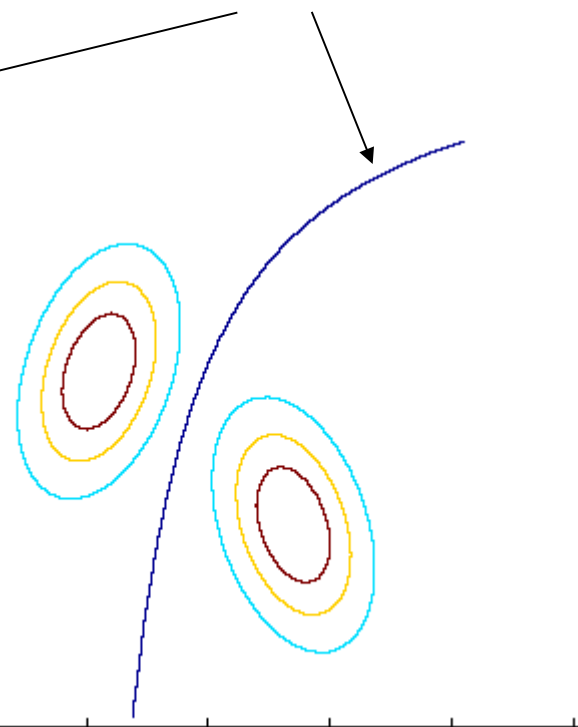


likelihoods



posterior for C_1

discriminant:
 $P(C_1|\mathbf{x}) = 0.5$



Common Covariance Matrix \mathbf{S}

- If not enough data can assume all classes have same common sample covariance matrix \mathbf{S} $\mathbf{S} = \sum_i \hat{P}(C_i) \mathbf{S}_i$

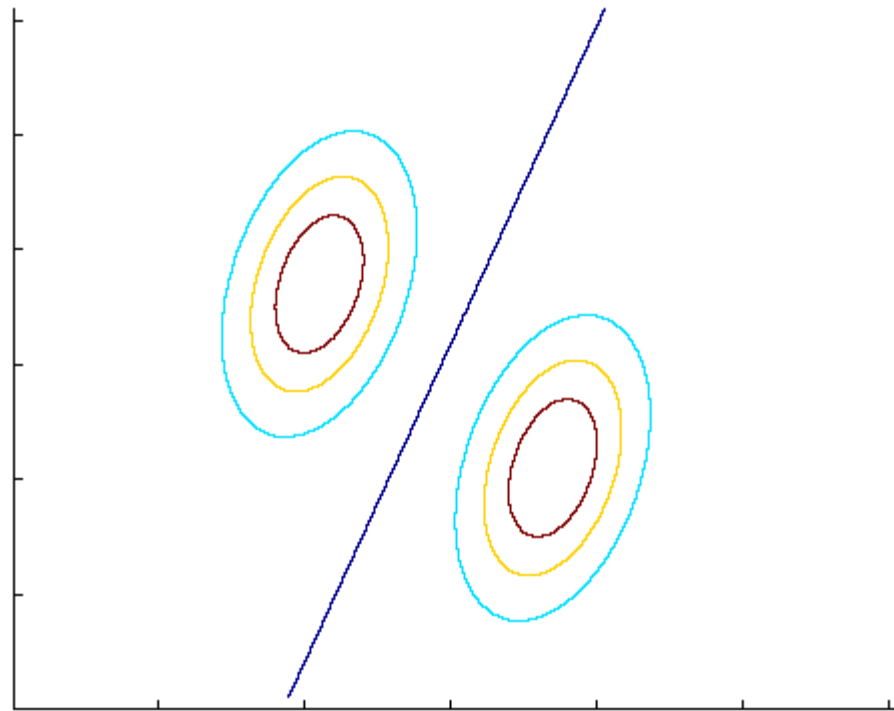
Discriminant reduces to a linear discriminant ($\mathbf{x}^T \mathbf{S}^{-1} \mathbf{x}$ is common to all discriminant and can be removed)

$$g_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \mathbf{m}_i)^T \mathbf{S}^{-1}(\mathbf{x} - \mathbf{m}_i) + \log \hat{P}(C_i)$$

$$g_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + w_{i0}$$

$$\text{where } \mathbf{w}_i = \mathbf{S}^{-1} \mathbf{m}_i \quad w_{i0} = -\frac{1}{2} \mathbf{m}_i^T \mathbf{S}^{-1} \mathbf{m}_i + \log \hat{P}(C_i)$$

Common Covariance Matrix \mathbf{S}



Diagonal Σ

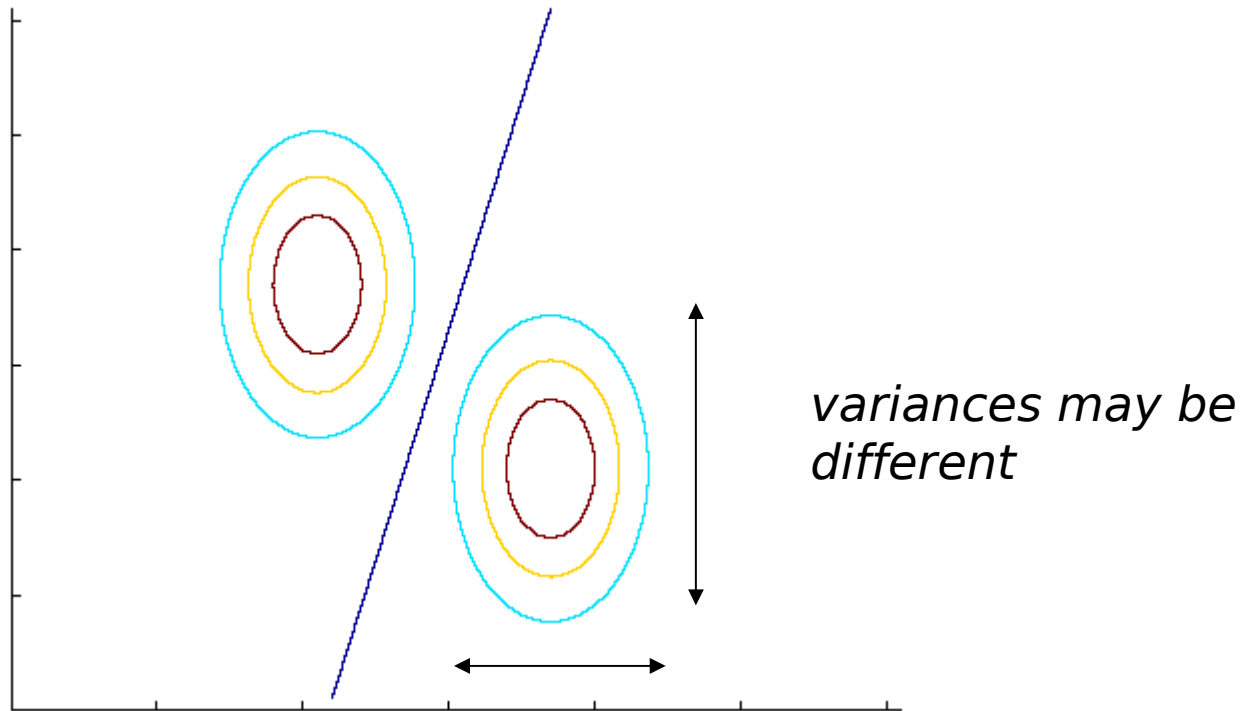
- When $x_j, j = 1, \dots, d$, are independent, Σ is diagonal

$p(\mathbf{x}|C_i) = \prod_j p(x_j|C_i)$ (Naive Bayes' assumption)

$$g_i(\mathbf{x}) = -\frac{1}{2} \sum_{j=1}^d \left(\frac{x_j^t - m_{ji}}{s_j} \right)^2 + \log \hat{P}(C_i)$$

Classify based on weighted Euclidean distance (in s_j units) to the nearest mean

Diagonal S



Diagonal **S**, equal variances

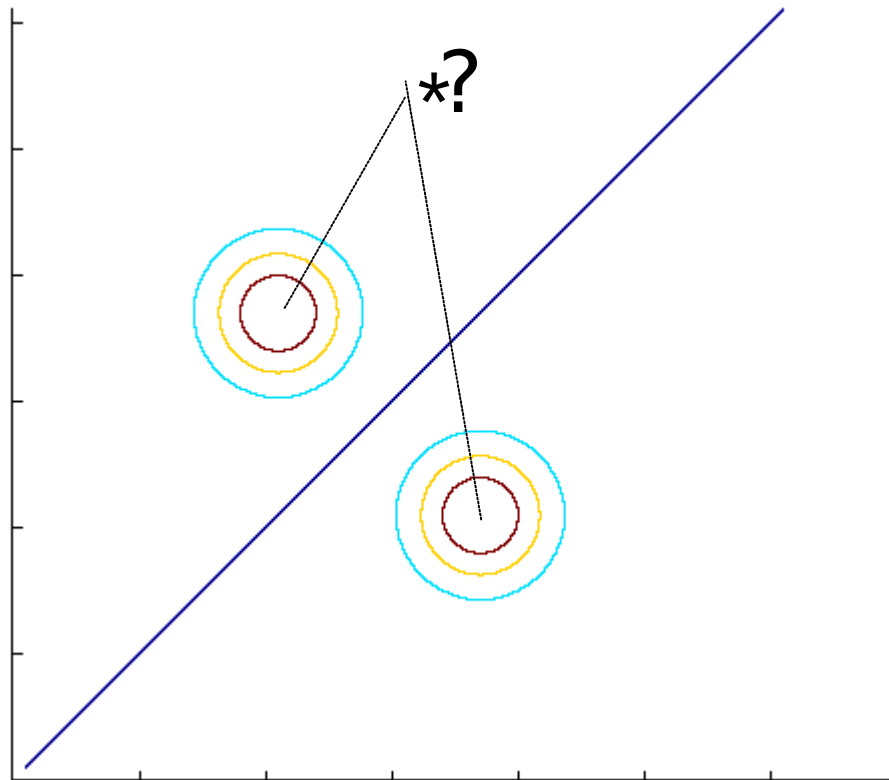
- Nearest mean classifier: Classify based on Euclidean distance to the nearest mean

$$g_i(x) = -\frac{\|x - m_i\|^2}{2s^2} + \log \hat{P}(C_i)$$

$$i = -\frac{1}{2s^2} \sum_{j=1}^d (x_j^t - m_{ij})^2 + \log \hat{P}(C_i)$$

- Each mean can be considered a prototype or template and this is template matching

Diagonal \mathbf{S} , equal variances



Model Selection

<i>Assumption</i>	<i>Covariance matrix</i>	<i>No of parameters</i>
Shared, Hyperspheric	$\mathbf{S}_i = \mathbf{S} = s^2 \mathbf{I}$	1
Shared, Axis-aligned	$\mathbf{S}_i = \mathbf{S}$, with $s_{ij} = 0$	d
Shared, Hyperellipsoidal	$\mathbf{S}_i = \mathbf{S}$	$d(d+1)/2$
Different, Hyperellipsoidal	\mathbf{S}_i	$K d(d+1)/2$

- As we increase complexity (less restricted \mathbf{S}), bias decreases and variance increases
- Assume simple models (allow some bias) to control variance (regularization)

Model Selection

- Different covariance matrix for each class
- Have to estimate many parameters
- Small bias , large variance
- Common covariance matrices, diagonal covariance etc. reduce number of parameters
- Increase bias but control variance
- In-between states?

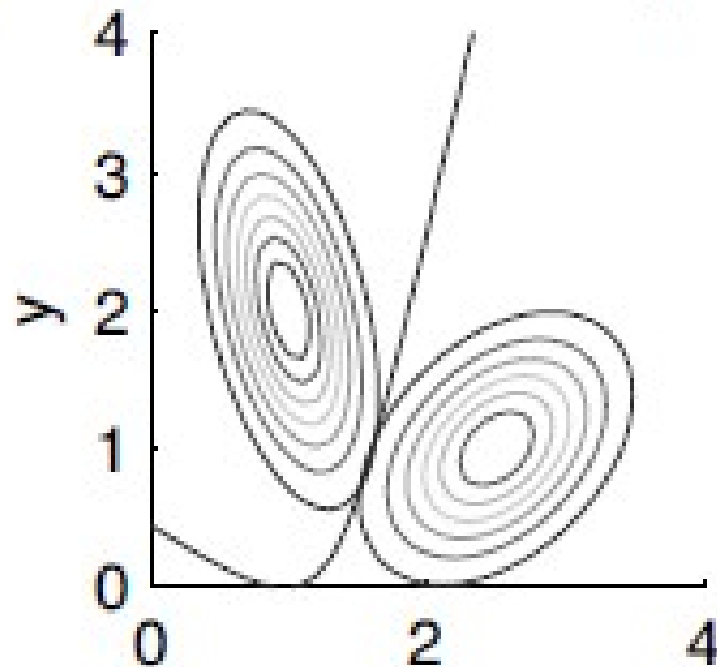
Regularized Discriminant Analysis(RDA)

$$S'_i = \alpha \sigma^2 \mathbf{I} + \beta S + (1 - \alpha - \beta) S_i$$

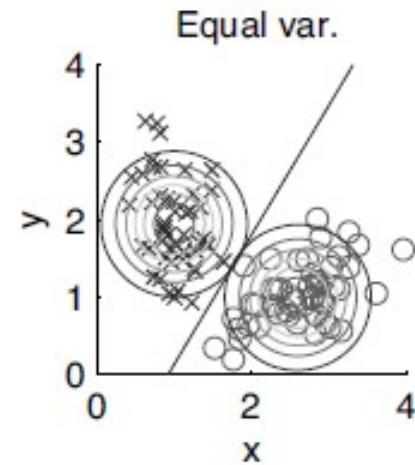
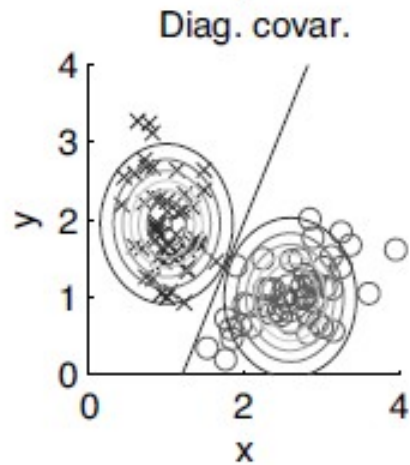
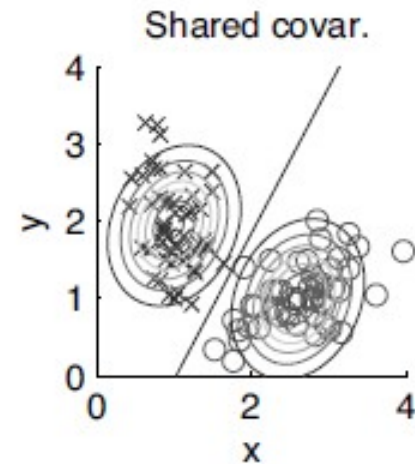
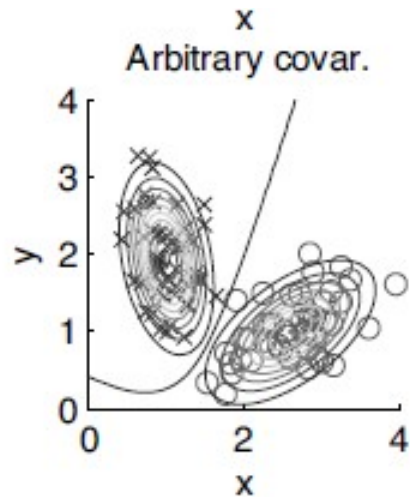
- $a=b=0$: Quadratic classifier
- $a=0, b=1$: Shared Covariance, linear classifier
- $a=1, b=0$: Diagonal Covariance
- Choose best a, b by cross validation

Model Selection: Example

Population likelihoods and posteriors



Model Selection



Discrete Features

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

if x_j are independent (Naive Bayes')

$$p(\mathbf{x}|C_i) = \prod_{j=1}^d p_{ij}^{x_j} (1-p_{ij})^{(1-x_j)}$$

the discriminant is linear

$$\begin{aligned} g_i(\mathbf{x}) &= \log p(\mathbf{x}|C_i) + \log P(C_i) \\ &= \sum_j [x_j \log p_{ij} + (1-x_j) \log(1-p_{ij})] + \log P(C_i) \end{aligned}$$

Estimated parameters $\hat{p}_{ij} = \frac{\sum_t x_j^t r_i^t}{\sum_t r_i^t}$

Multivariate Regression

$$r^t = g(x^t | w_0, w_1, \dots, w_d) + \varepsilon$$

- Multivariate linear model

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

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

Multivariate Regression

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

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

$$\sum_t r^t = Nw_0 + w_1 \sum_t x_1^t + w_2 \sum_t x_2^t + \cdots + w_d \sum_t x_d^t$$

$$\sum_t x_1^t r^t = w_0 \sum_t x_1^t + w_1 \sum_t (x_1^t)^2 + w_2 \sum_t x_1^t x_2^t + \cdots + w_d \sum_t x_1^t x_d^t$$

$$\sum_t x_2^t r^t = w_0 \sum_t x_2^t + w_1 \sum_t x_1^t x_2^t + w_2 \sum_t (x_2^t)^2 + \cdots + w_d \sum_t x_2^t x_d^t$$

\vdots

$$\sum_t x_d^t r^t = w_0 \sum_t x_d^t + w_1 \sum_t x_d^t x_1^t + w_2 \sum_t x_d^t x_2^t + \cdots + w_d \sum_t (x_d^t)^2$$

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

Based on E. Alpaydm, 2004, Introduction to Machine Learning © The MIT Press (v1.2)

Dimensionality reduction

- Some features (dimensions) bear little or no 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 , $i = 1, \dots, d$ dimensions to new $k < d$ dimensions, z_j , $j = 1, \dots, k$
 - Principal Components Analysis (PCA)
 - Linear Discriminant Analysis (LDA)
 - Factor Analysis (FA)

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 $y=g(x|\Theta)$)

Subset selection

- Have initial set of features of size d
- There are 2^d possible subsets
- Need a criteria to decide which subset is the best
- A way to search over the possible subsets
- Can't go over all 2^d possibilities
- Need some heuristics

“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”
- Mutual Information (Kullback–Leibler Divergence) is more general measure of “mutual information”
- Can be extended to n variables (information variables x_1, \dots, x_n have about variable x_{n+1})

Subset-selection

- 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

Subset Selection

- There are 2^d subsets of d features
- Forward search: Add the best feature at each step
 - Set of features F initially \emptyset .
 - At each iteration, find the best new feature
$$j = \operatorname{argmin}_i E (F \cup x_i)$$
 - Add x_j to F if $E (F \cup x_j) < E (F)$
- Hill-climbing $O(d^2)$ algorithm
- Backward search: Start with all features and remove one at a time, if possible.
- Floating search (Add k , remove l)

Floating Search

- Forward and backward search are “greedy” algorithms
 - Select best options at single step
 - Do not always achieve optimum value
- Floating search
 - Two types of steps: Add k , remove l
 - *More computations*

Feature Extraction

- Face recognition problem
 - Training data input: pairs of Image + Label(name)
 - Classifier input: Image
 - Classifier output: Label(Name)
- Image: Matrix of $256 \times 256 = 65536$ 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.

Projection

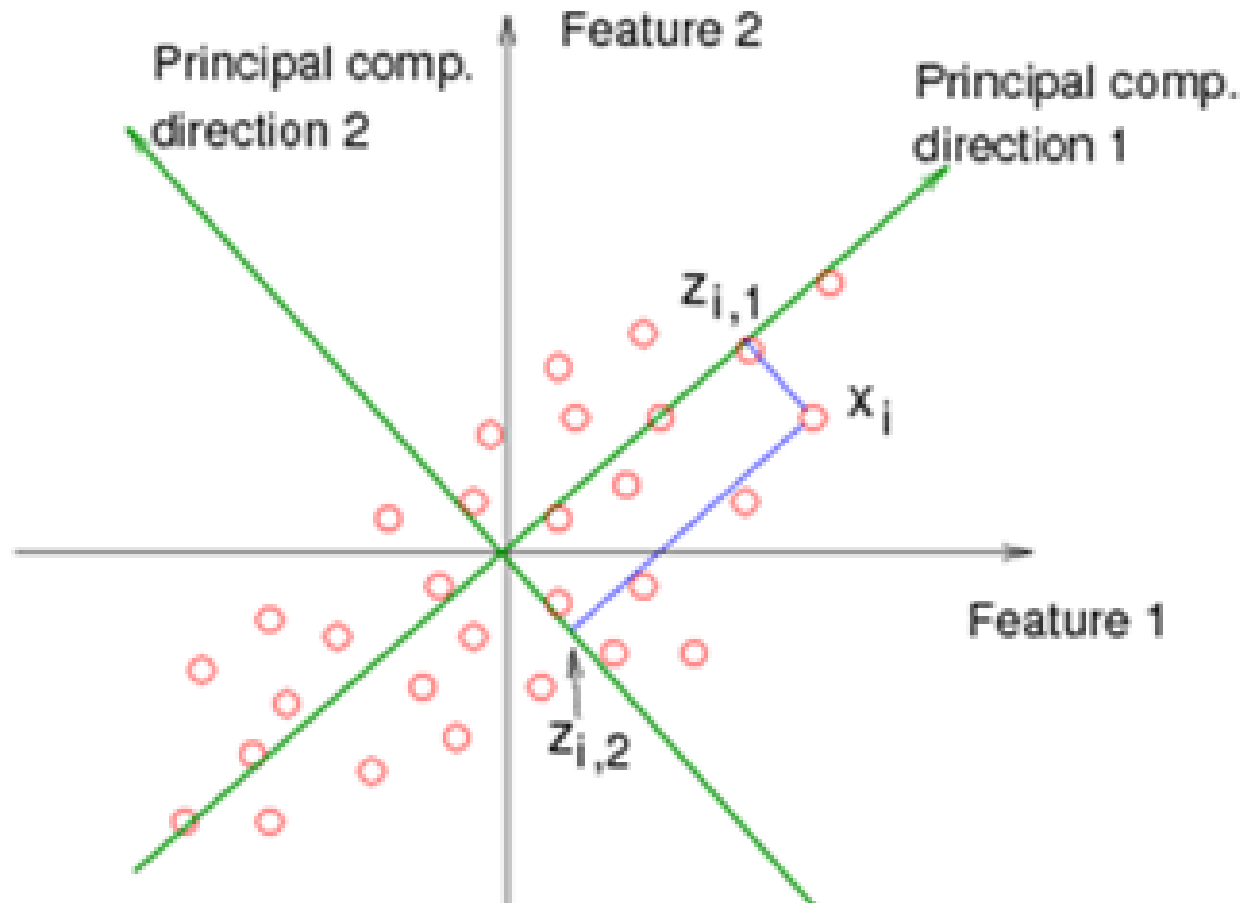
- Find a projection matrix w from d -dimensional to k -dimensional vectors that keeps error low

$$z = w^T x$$

PCA: Motivation

- Assume that d observables are linear combination of $k < d$ vectors
- $Z_i = W_{i1}X_{i1} + \dots + W_{ik}X_{id}$
- 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 (e.g. w_{i1} have large variation) or otherwise bear no information

PCA: Motivation



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

PCA

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

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

- Taking Derivatives

$$2\Sigma\mathbf{w}_1 - 2\alpha\mathbf{w}_1 = 0 \quad \Sigma\mathbf{w}_1 = \alpha\mathbf{w}_1$$

- Eigenvector. Since want to maximize $\mathbf{w}_1^T \Sigma \mathbf{w}_1 = \alpha \mathbf{w}_1^T \mathbf{w}_1 = \alpha$ we should choose an eigenvector with largest eigenvalue

PCA

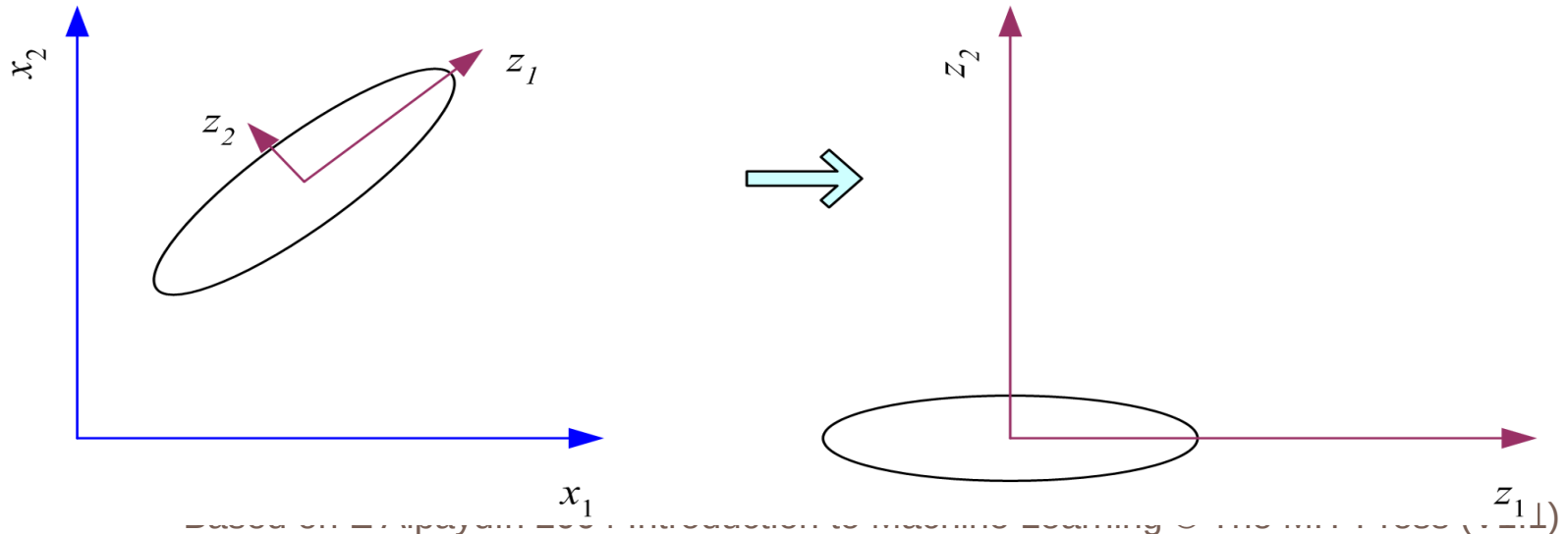
- d-dimensional feature space
- d by d symmetric covariance matrix estimated from samples
$$\text{Cov}(\mathbf{x}) = \Sigma,$$
- Select k largest eigenvalue of the covariance matrix and associated k eigenvectors
- The first eigenvector will be a direction with largest variance

What PCA does

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

where the columns of \mathbf{W} are the eigenvectors of Σ , and \mathbf{m} is sample mean

Centers the data at the origin and rotates the axes



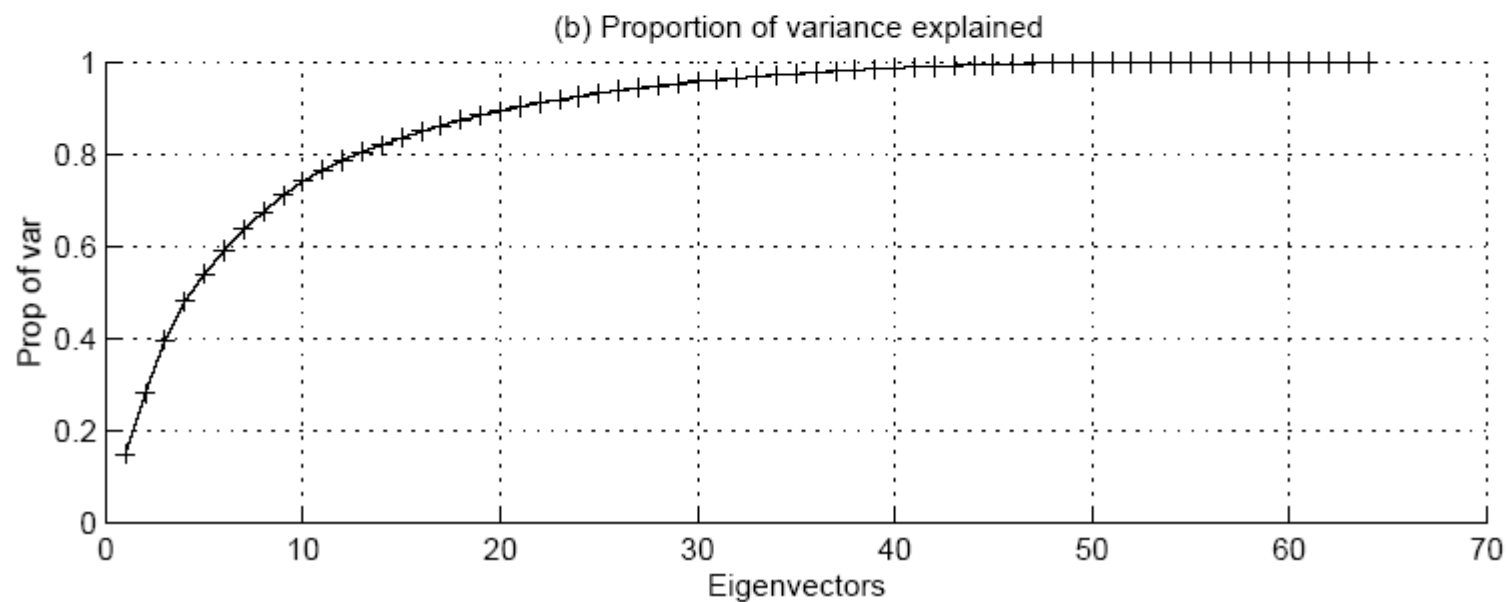
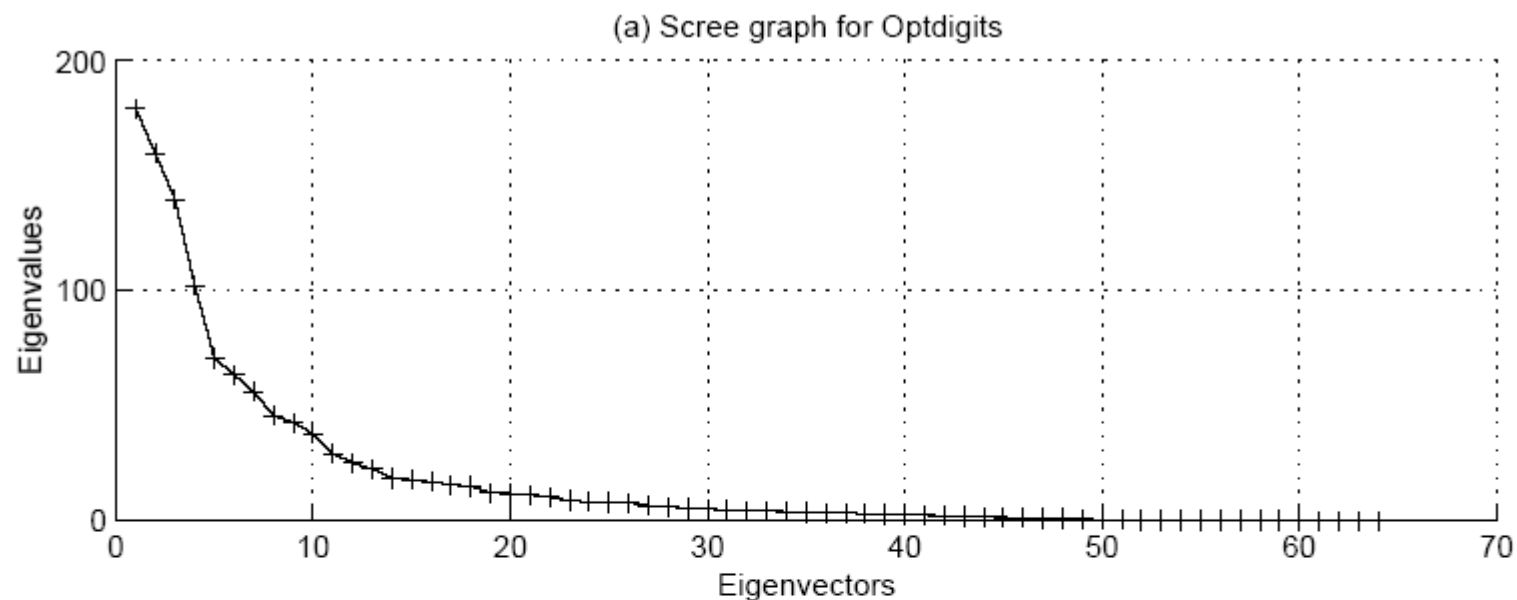
How to choose k ?

- Proportion of Variance (PoV) explained

$$\frac{\lambda_1 + \lambda_2 + \cdots + \lambda_k}{\lambda_1 + \lambda_2 + \cdots + \lambda_k + \cdots + \lambda_d}$$

when λ_i are sorted in descending order

- Typically, stop at $\text{PoV} > 0.9$
- Scree graph plots of PoV vs k , stop at “elbow”



PCA

- PCA is unsupervised (does not take into account class information)
- Can take into account classes : Karhunen-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

- 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