

# Models for Probability Distributions and Density Functions

# General Concepts

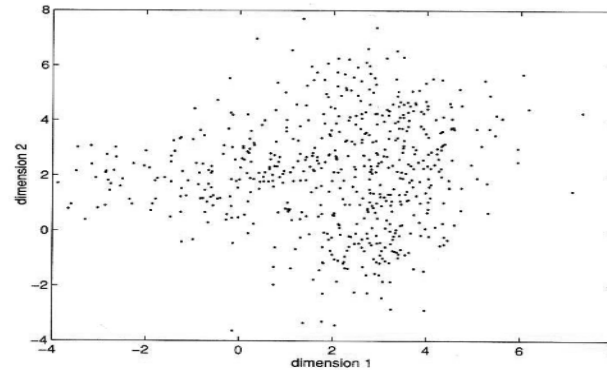
- Parametric:
  - E.g., Gaussian, Gamma, Binomial
- Non-Parametric:
  - E.g., kernel estimates
- Intermediate models: Mixture Models

$$p(x) = \sum_{k=1}^K p_k(x | \theta_k) \pi_k$$

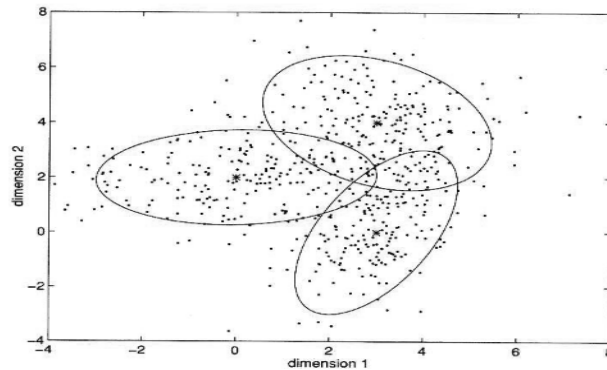
# Gaussian Mixture Model

## Two-dimensional Data Set

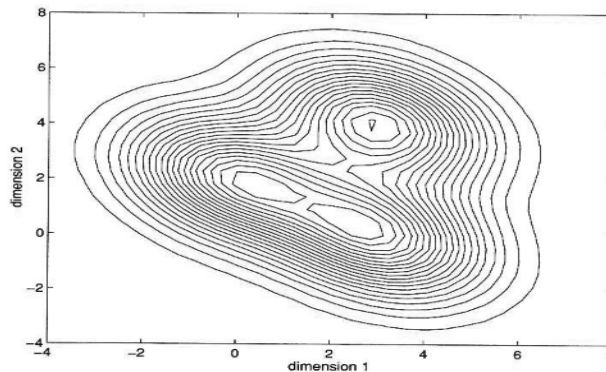
- Mixture models are interpreted as being generated with a hidden variable taking  $K$  values revealed by the data
- EM algorithm is used to learn parameters of mixture models



Data points from three bivariate normal distributions with equal weights



Component Density Contours



Contours of Constant density

# Joint Distributions for Unordered Categorical Variables

Case of Two variables:

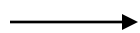
**Variable A:** Dementia: has three possible values

**Variable B:** Smoker: has two possible values

There are **six** possible values for the joint distribution

## Dementia

Contingency  
Table of  
Medical  
Patients  
With  
Dementia



Smoker?	None	Mild	Severe
No	426	66	132
Yes	284	44	88

# Joint Distributions for Unordered Categorical Variables

Variable A:  $\{a_1, a_2, \dots, a_m\}$

Variable B:  $\{b_1, b_2, \dots, b_m\}$

.....p variables

There are  $m^p - 1$  possible independent values for the joint distribution

(to fully specify the model)

The -1 comes from the constraint that they sum to 1  
Contingency tables are impractical when m and p are large  
(e.g., when m=2 and p=20  
impossibly large number of values are needed).

Need systematic techniques for structuring both densities and distribution functions.

# Factorization and Independence in High Dimensions

- Can construct simpler models for multidimensional data
- If we assume that individual variables are independent, the joint density function can be written as

$$p(x) = p(x_1, \dots, x_p) = \prod_{k=1}^p p_k(x_k) \leftarrow \text{One-dimensional density function}$$

- Simpler to model the one-dimensional densities separately than model them jointly
- Independence model for  $\log p(x)$  has an additive form

# Smoker, Dementia Example

Smoker?	None	Mild	Severe
No	426	66	132
Yes	284	44	88

Smoker?	None	Mild	Severe
P(dementia= , No Smoker)	0.410	0.063	0.126
P(dementia= , Yes Smoker)	0.273	0.042	0.084

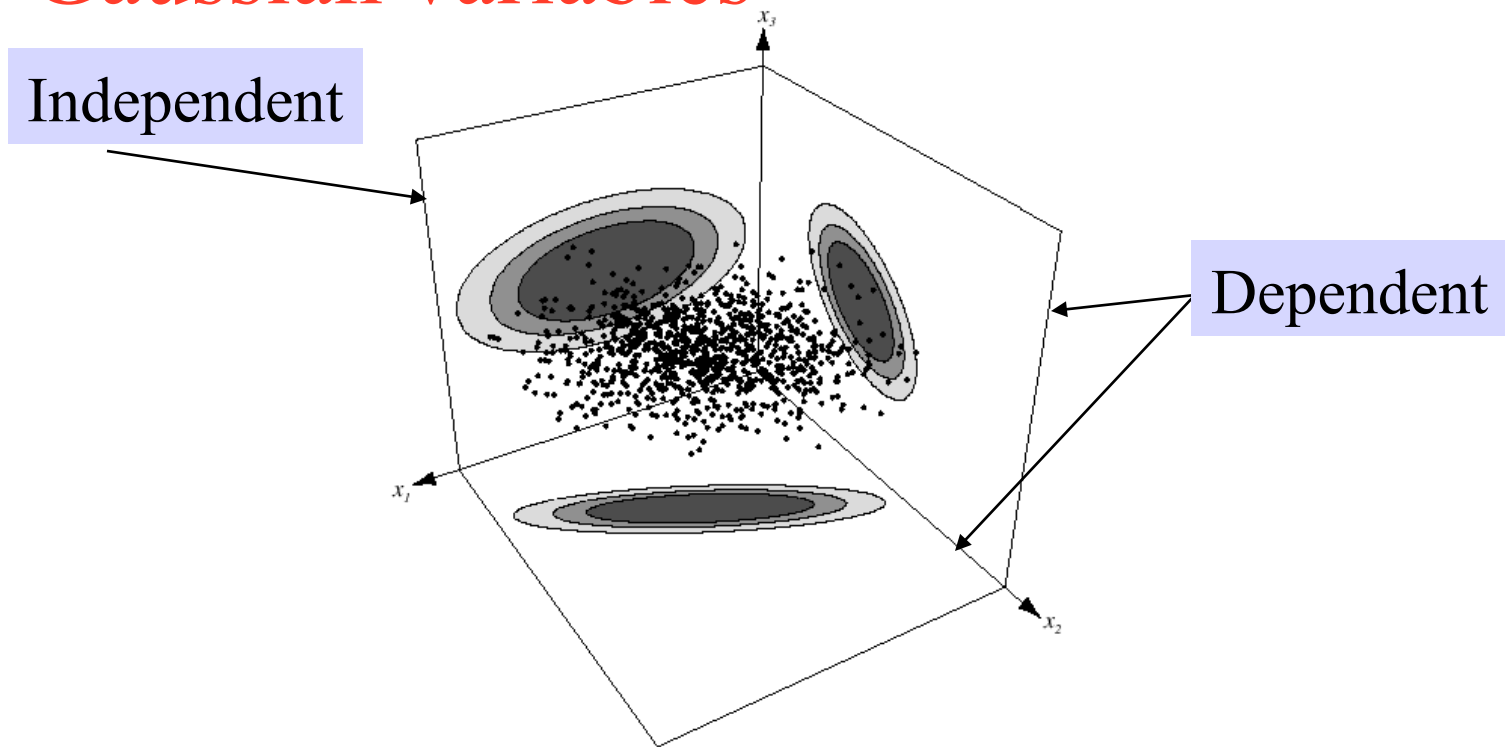
Smoker?	None	Mild	Severe
P(dementia= /No)	0.683	0.105	0.212
P(dementia= /Yes)	0.683	0.105	0.212

P(No)	0.6		
P(Yes)	0.4		

Prob(dementia=none, smoker=No)=0.410

Prob(dementia=none) x Prob(smoker=No)=0.683 x 0.6=0.410 <sup>7</sup>

# Statistically dependent and independent Gaussian variables



3-D distribution which obeys  $p(x1, x3) = p(x1)p(x3)$ ;  $x1$  and  $x3$  are independent but other pairs are not



# Improved Modeling

- Find something in-between independence (low complexity) and complete knowledge (high complexity)
- Factorize into sequence of conditional distributions  
 $p(x_1, x_2) = p(x_1)p(x_2 / x_1)$  from definition of conditional probability

$$p(B/A) = \frac{p(A,B)}{p(A)}$$

$$p(x_1, x_2, x_3) = p(x_1)p(x_2, x_3 / x_1) = p(x_1)p(x_2 / x_1)p(x_3 / x_1, x_2)$$

$$p(x) = p(x_1) \prod_{k=2}^p p(x_k | x_1, \dots, x_{k-1})$$

Some of these can be ignored<sup>9</sup>

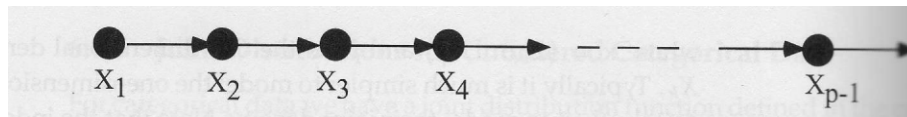
# Graphical Models

- Natural representation of the model as a directed graph
- Nodes correspond to variables
- Edges show dependencies between variables
- Edges directed into node for  $k^{\text{th}}$  variable will come from subset of variables  $X_1, \dots, X_{k-1}$
- Can be used to represent many different structures
  - Markov model
  - Bayesian network
  - Latent variables
  - Naïve Bayes
  - Hidden Markov Model

# Graphical Models

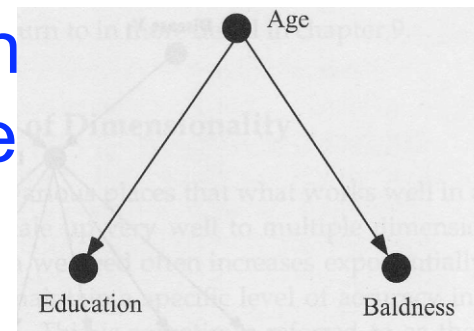
- First order Markov assumption
- Appropriate when the variables represent the same property measured sequentially , e.g., different times

$$p_k(x_k | x_1, \dots, x_{k-1}) = p_k(x_k | x_{k-1})$$



# Bayesian Belief Network

- Variables *age*, *education*, *baldness*
- *Age* cannot depend on *education* or *baldness*
- Conversely *education* and *baldness* depend on *age*
- Given *age*, *education* and *baldness* are not dependent on each other

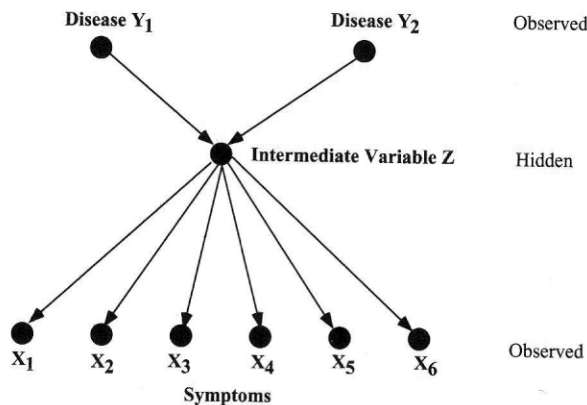


- Two variables *education* and *baldness* that are conditionally independent given *age*

# Latent Variables

- Extension to unobserved hidden variables
  - Two diseases that are conditionally independent
- Simplify relationships in the model structure

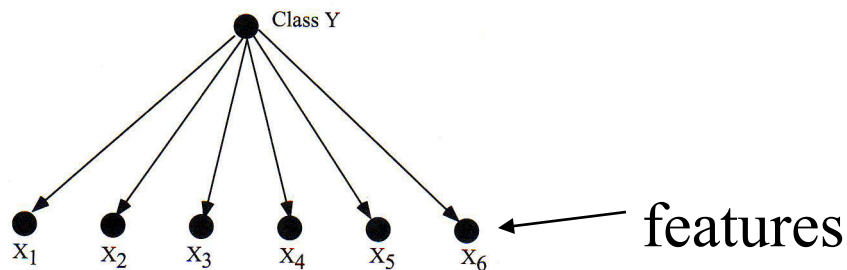
Given the intermediate variable value the symptoms are independent



# First order Bayes graphical model

- Naïve Bayes classifier
- In the context of classification and clustering features are assumed to be independent of each other given the class label  $y$

$$p(x / y) = \prod_{j=1}^p p(x_j / y)$$



# Curse of Dimensionality

- What works well in one dimension may not scale up to multiple dimensions
- Amount of data needed increases exponentially
- Data mining often involves high dimensions

$$E[(\hat{p}(x) - p(x))^2] / p(x)^2 < 0.1 \text{ at } x = 0$$

where  $p(x)$  is the true Normal density and  $\hat{p}(x)$  is a kernel estimate with a normal kernel

- For a 10% relative accuracy
  - In one dimension need 4 points
  - Two dimensions need 19 points
  - Three dimensions 67 points
  - Six dimensions 2790 points
  - 10 dimensions need 842,000 points

# Coping with High Dimensions

- Two basic (obvious) strategies
  1. Use subset of the relevant variables
    - Find a subset  $p'$  of variables where  $p' \ll p$
  2. Transform original  $p$  variables into a new set of  $p'$  variables, with  $p' \ll p$ 
    - Examples are PCA, Projection pursuit, neural networks



# Feature Subset Selection

- Variable selection is a general strategy when dealing with high-dimensional problems
- Consider predicting  $Y$  using  $X_1, \dots, X_p$
- Some may be completely unrelated to predictor variable  $Y$ 
  - Month of person's birth to credit-worthiness
- Others may be redundant
  - Income before tax and income after tax are highly correlated

# Gauging Relevance Quantitatively

- If  $p(y/x_1) = p(y)$  for all values of  $y$  and  $x_1$  then  $Y$  is independent of input variable  $X_1$
- If  $p(y/x_1, x_2) = p(y/x_2)$  then  $Y$  is independent of  $X_1$  if the value of  $X_2$  is already known
- How to estimate this dependence
  - We are not only interested in strict dependence/independence but also in the degree of dependence

# Mutual Information

- Dependence between  $Y$  and  $X$

$$I(Y; X') = \sum_{i,j} p(y_i, x'_j) \log \frac{p(y_i, x'_j)}{p(y_i)p(x'_j)}$$

- Where  $X'$  is a categorical variable (a quantized version of real-valued  $X$ )
- Other measures of the relationship between  $Y$  and  $X$ 's can also be used

# Sets of Variables

- Interaction of individual  $X$  variables does not tell us how sets of variables interact with  $Y$
- Extreme example:
  - $Y$  is a parity function that is 1 if the sum of binary values  $X_1, \dots, X_p$  is even and 0 otherwise
  - $Y$  is independent of any individual  $X$  variable, yet it is a deterministic function of the full set
- $k$  best individual variables (e.g., ranked by correlation) is not the same as the best  $k$  variables
- Since there are  $2^p - 1$  different non-empty subsets of  $p$  variables, exhaustive search is infeasible
- Heuristic search algorithms are used, e.g., greedy selection where one variable at a time is added or deleted

# Transformations for High-Dimensional Data

- Transform the  $X$  variables into  $Z$  variables  $Z_1, \dots, Z_{p'}$
- Called *basis functions, factors, latent variables, principal components*
- Projection Pursuit Regression
- Neural networks use

$$h_j(t) = \frac{1}{1 + e^{-t}}$$

$$\hat{y} = \sum_{j=1}^{p'} w_j h_j(\alpha_j^T x)$$

Projection of  $x$  onto the  $j^{\text{th}}$  weight vector  $\alpha_j$

# Principal Components Analysis

- Linear combinations of the original variables
- Sets of weights are chosen so as to maximize the variance when expressed in terms of the new variables
- PCA may not be ideal when goal is predictive performance
  - For classification and clustering PCA need not emphasize group differences and can hide them