

TIM 245 Lecture 6 (4/19/17)

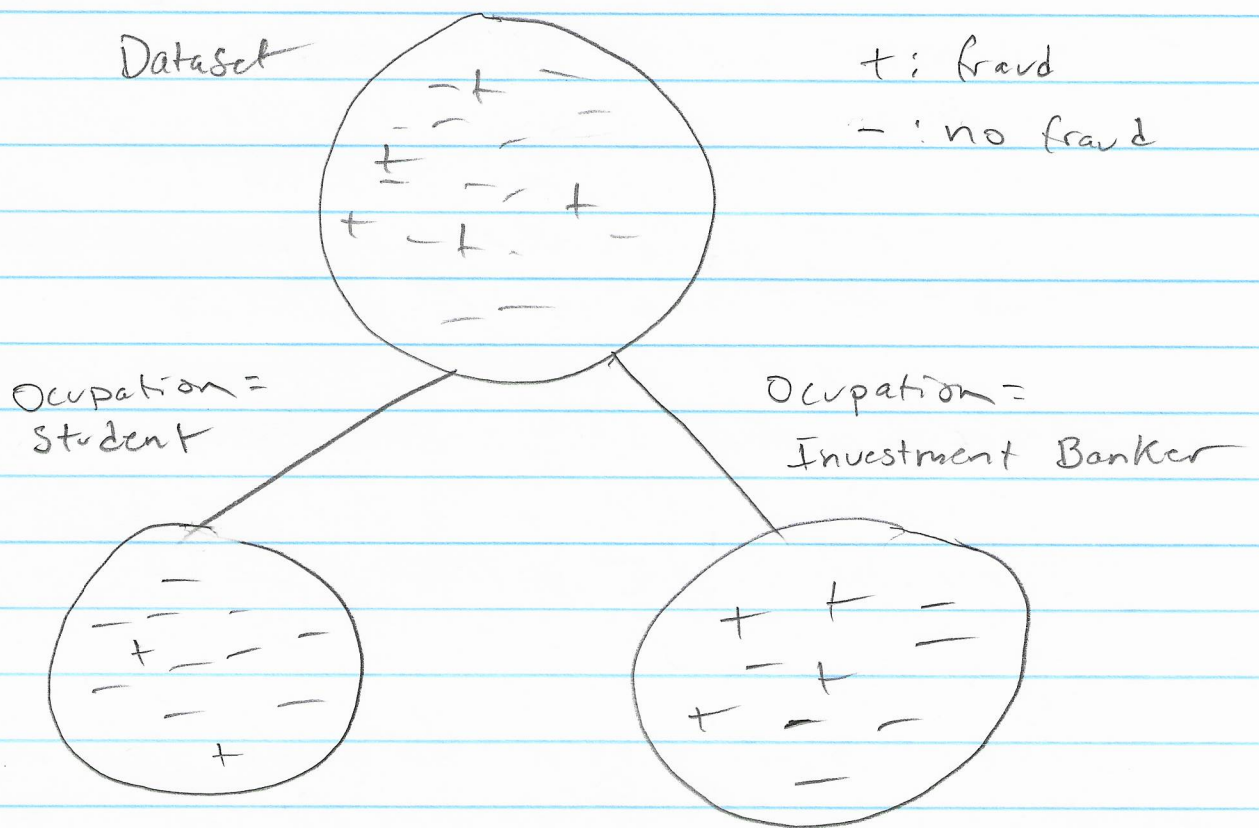
Agenda

- 1) Review Homework 1
- 2) Feature Selection using Information Gain
- 3) Wrapper based methods and general comments on feature selection
- 4) Dimensionality Reduction: PCA
- 5) Project Phase II and roadmap for the course

② Feature Selection using Information Gain

Each attribute gives us information on the target y .

Example : Fraud Detection



How can we select the attributes that provide the most information on y ?

Let:

$c_i^y \triangleq$ i th ($i=1, 2, \dots, r$) of
the target y

$c_k^{A_j} \triangleq$ k th ($k=1, 2, \dots, s$) of A_j

$p(c_i^y) \triangleq$ probability of c_i^y estimated
as $|c_i^y| / |D|$

Process for computing information gain
of attribute A_j :

- 1) Compute the entropy of the complete dataset D

$$\text{Entropy}(D) = -\sum_{i=1}^r p(c_i^y) \log_2(p(c_i^y))$$

- 2) Partition D into s subsets where partition D_k contains only instances where $A_j = c_k^{A_j}$

$$\text{Entropy}_{A_j}(D) = \sum_{k=1}^s \frac{|D_k|}{|D|} \times \text{Entropy}(D_k)$$

- 3) Compute the information gain

$$\text{Gain}(A_j) = \text{Entropy}(D) - \text{Entropy}_{A_j}(D)$$

Information gain can be biased towards attributes with a large number of values.

Gain ratio is a normalized version of information gain.

$$\text{SplitInfo}_{A_j}(D) = - \sum_{k=1}^s \frac{|D_k|}{|D|} \times \log_2 \left(\frac{|D_k|}{|D|} \right)$$

$$\text{Gain Ratio}(A_j) = \frac{\text{Gain}(A_j)}{\text{SplitInfo}_{A_j}(D)}$$

③ Wrapper Based Methods

Given a particular learning algorithm, e.g. KNN, we can search for the subset of attributes that provide the best performance.

This is called a wrapper based method.

Search can be forwards or backwards through the attributes:

$$\{\emptyset\} \rightarrow \{A_1\} \rightarrow \{A_1, A_2\}$$

$$\{A_1, A_2, A_3\} \rightarrow \{A_1, A_2\}$$

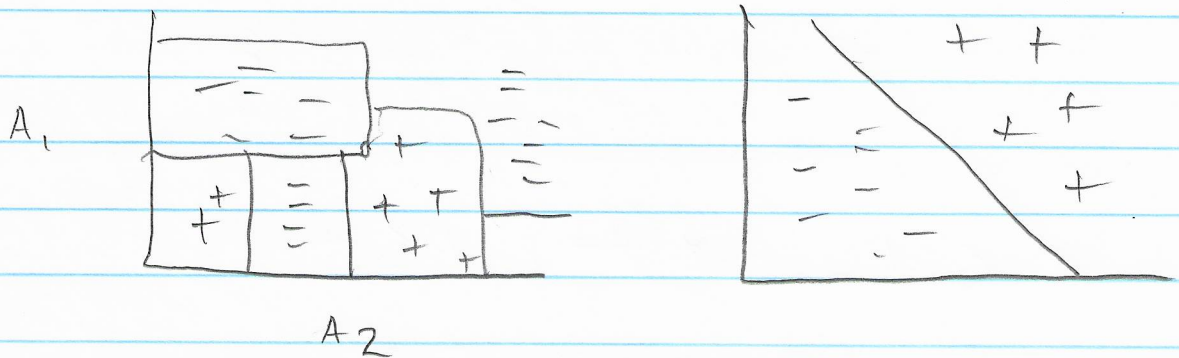
Solved using hill climb or greedy algorithm

Result is the best set of attributes for the specific algorithm.

General Comments on Feature Selection

- 1) Start with fast filter based methods to come up with a preliminary set of attributes

Correlation and information gain find different kinds of relationships



Better suited to
information gain

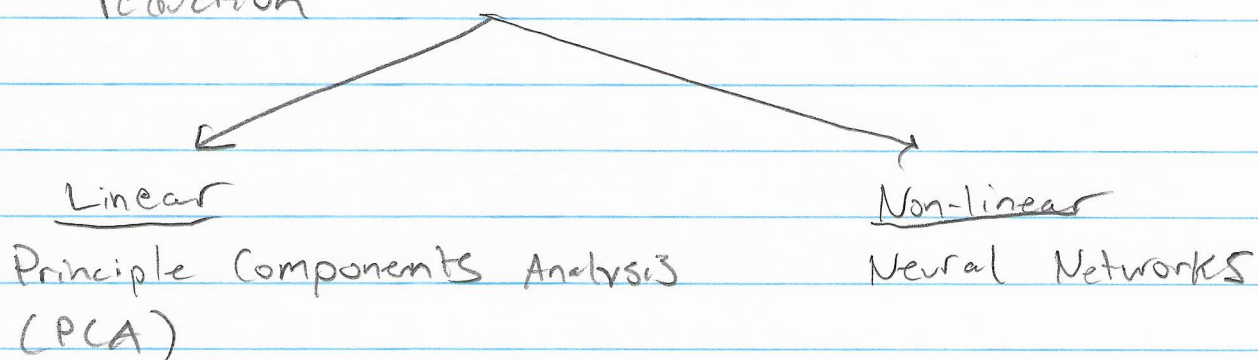
Better suited to
correlation

- 2) Experiment with a variety of different learning algorithms

- 3) Apply wrapper based feature selection using the best learning algorithm

④ Dimensionality Reduction

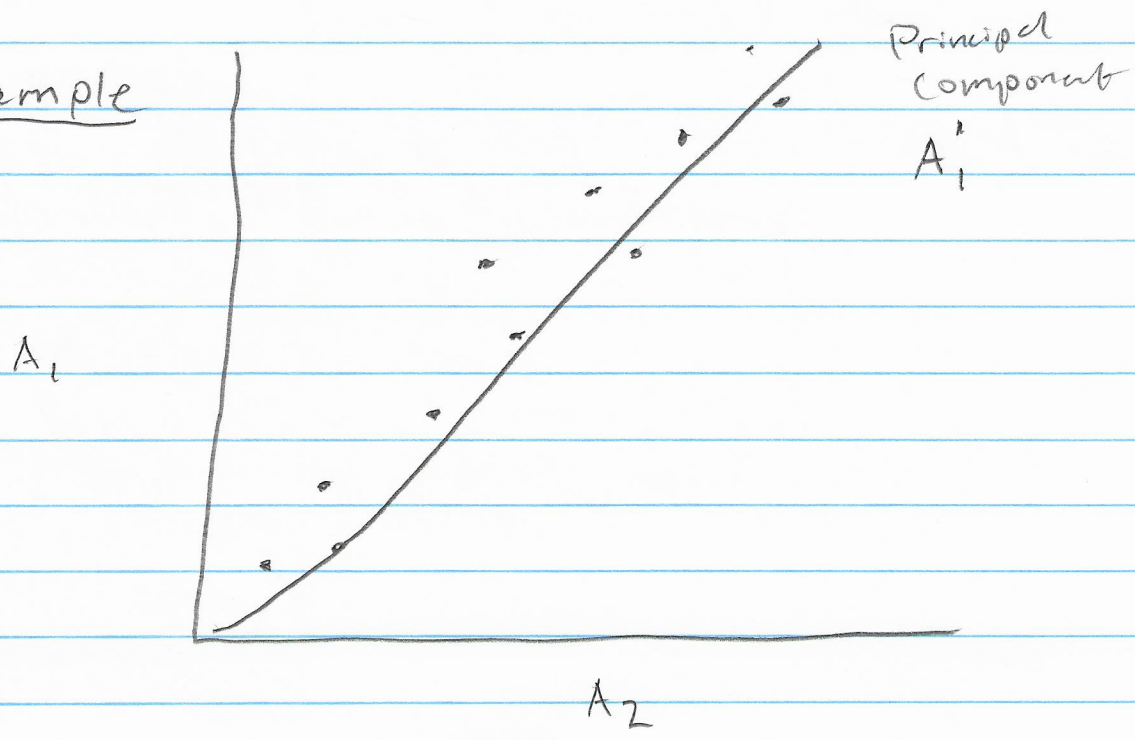
Two general approaches to dimensionality reduction



Principle Components Analysis

Find the internal "axes" of the data set. Each "axes" becomes a new high-level feature

Example



Assumptions:

- 1) Relationship between the variables is linear
- 2) Mean and Covariance is important
- 3) Large Variances have important dynamics