Lecture 3 (4/10/17)

Agenda

1) General Comments on the Project

2) Formalization of the Data Mining tasks

3) Exploratory Data Analysis (EDA)

4) Outlier analysis

② Formalization of the Data Mining Tasks

Let:

$D \triangleq$ dataset of $n$ instances and $m$ attributes

$d_i \triangleq$ ith instance $(i = 1, 2, \ldots n)$

$A_j \triangleq$ jth attribute $(j = 1, 2, \ldots m)$

$X_{ij} \triangleq$ value of Attribute $A_j$ for instance $d_i$

|        | $A_1$    | $A_2$    | $A_3$    | $\ldots$ | $A_m$    |
|--------|----------|----------|----------|----------|----------|
| $d_1$  | $X_{11}$ | $X_{12}$ | $X_{13}$ |          | $X_{1m}$ |
| $d_2$  |          |          |          |          |          |
| $\vdots$ |        |          |          |          |          |
| $d_n$  | $X_{n1}$ | $X_{n2}$ | $X_{n3}$ |          | $X_{nm}$ |

Predictive Analysis (supervised Learning)

Let attribute, $A_j$, be the target $Y$

$$d_i = (X_{i1}, X_{i2}, \ldots X_{im}, y_i) \equiv (X_i, y_i)$$

Learn function $\Phi : X \rightarrow Y$

classification : $y$ is a nominal attribute

prediction : $y$ is a numerical attribute

# Descriptive Analysis (Unsupervised Learning)

## Cluster Analysis

Organize the dataset $D$ into groups such that all instances in a group are similar to each other and dissimilar to instances in other groups based on some distance metric

## Association Analysis

Let $X$ and $Y$ be sets of attribute values such that $X \cap Y = \emptyset$

Association rule $X \Rightarrow Y$ means that instances containing values from $Y$ also contain values from $X$

③ Exploratory Data Analysis

The purpose of data cleaning is to remove noise and inconsistencies that do not reflect the real-world process that created the data.

In order to do this accurately, we first need to understand the underlying process that created the data.
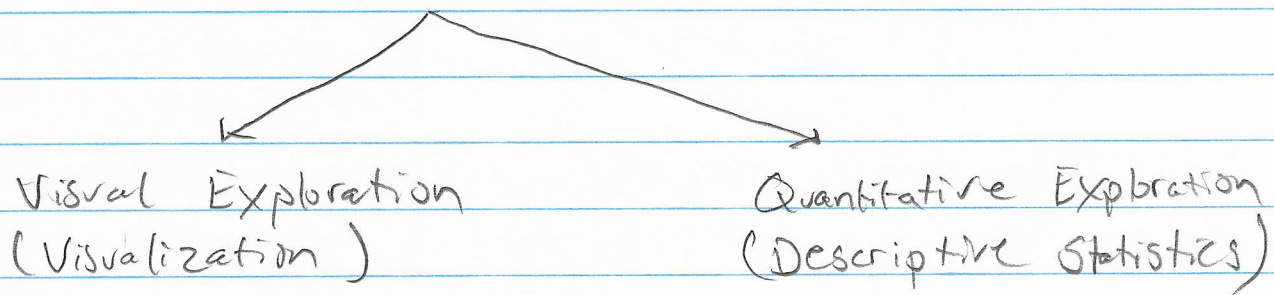
There are two ways to understand a dataset:

1) Domain Knowledge: fit the data to the process ("Inside-Out")

2) Exploratory Data Analysis: fit a process to the data ("Outside-In")

EDA Objectives

1) What is a typical value?

2) What is the uncertainity for a typical value?

3) What is a good distributional fit?

4) Does the attribute affect other attributes?

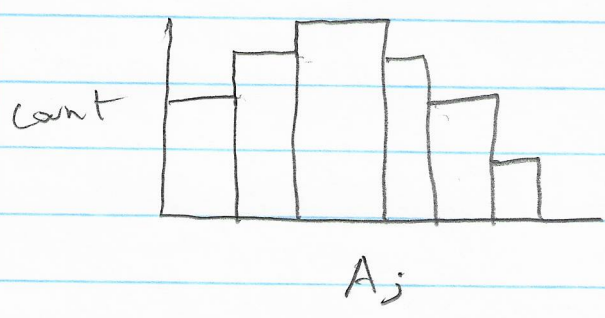5) Does the data contain outliers?

Two general approaches to EDA

Visual Exploration          Quantitative Exploration
(Visualization)             (Descriptive Statistics)

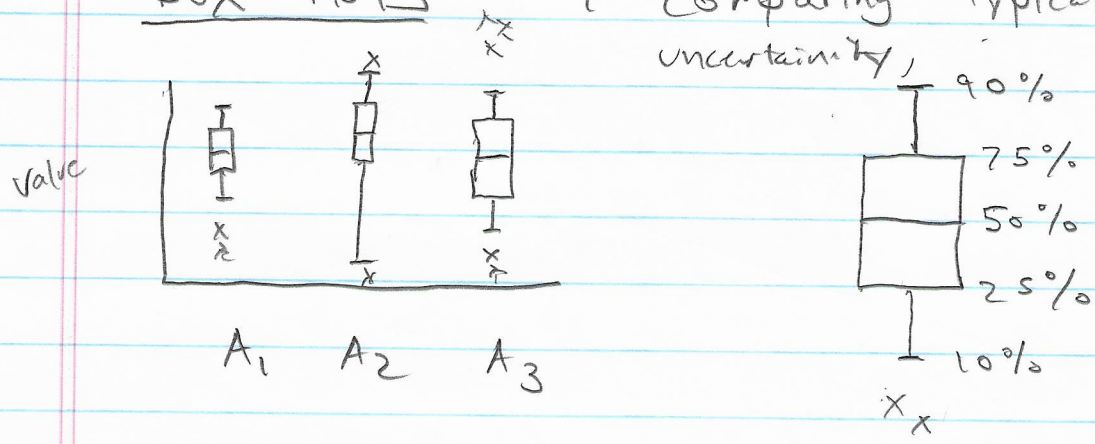Visual and quantitative exploration should be done in parallel
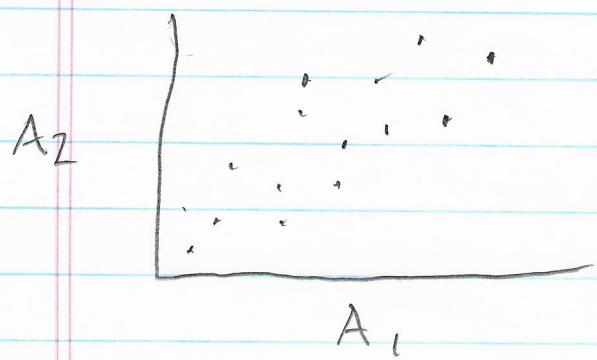
Visual Exploration

Always plot the data!

Histogram → understand typical values, uncertainity, and distribution

count

$A_j$

Box Plots → comparing typical values, uncertainity, outlier detection

value

$A_1$ $A_2$ $A_3$

90%
75%
50%
25%
10%

Scatter Plot → understanding relationship between variables

$A_2$

$A_1$

Quantitative Exploration

Three general ways of describing how attribute, $A_j$, behaves :

1) Central Tendency : What is the typical value of $A_j$

2) Dipersion : What is the spread of $A_j$

3) Correlation : What is the relationship between $A_j$ and another attribute $A_k$

Measures of Central Tendency

Mean (average) $\quad : \frac{1}{n} \sum_{i=1}^{n} x_{ij}$

Median $\qquad :$ middle value of
$\qquad\qquad\qquad x_{1j}, x_{2j}, \ldots x_{nj}$ (sorted)

Mode $\qquad\qquad :$ most frequent value of
$\qquad\qquad\qquad x_{1j}, x_{2j} \ldots x_{nj}$

Mean uses all of the data but
is sensitive to outliers

Median is robust to outliers but can
be sensitive to small changes

Mode useful for nominal data but
rarely used for numerical data

Measures of Dispersion

Variance $(\sigma_{A_j}^2) = \frac{1}{n} \sum_{i=1}^{n} (X_{i,j} - \bar{A}_j)^2$

std dev $(\sigma_{A_j}) = \sqrt{\sigma_{A_j}^2}$

Median Absolute $= $ median $($ 
Deviation (MAD) $|X_{1,j} - \bar{A}_j|, |X_{2,j} - \bar{A}_j|, \ldots$
$|X_{n,j} - \bar{A}_j| )$

Variance and std dev are sensitive to outliers

Median Absolute Deviation is useful for data sets with long-tails and/or outliers

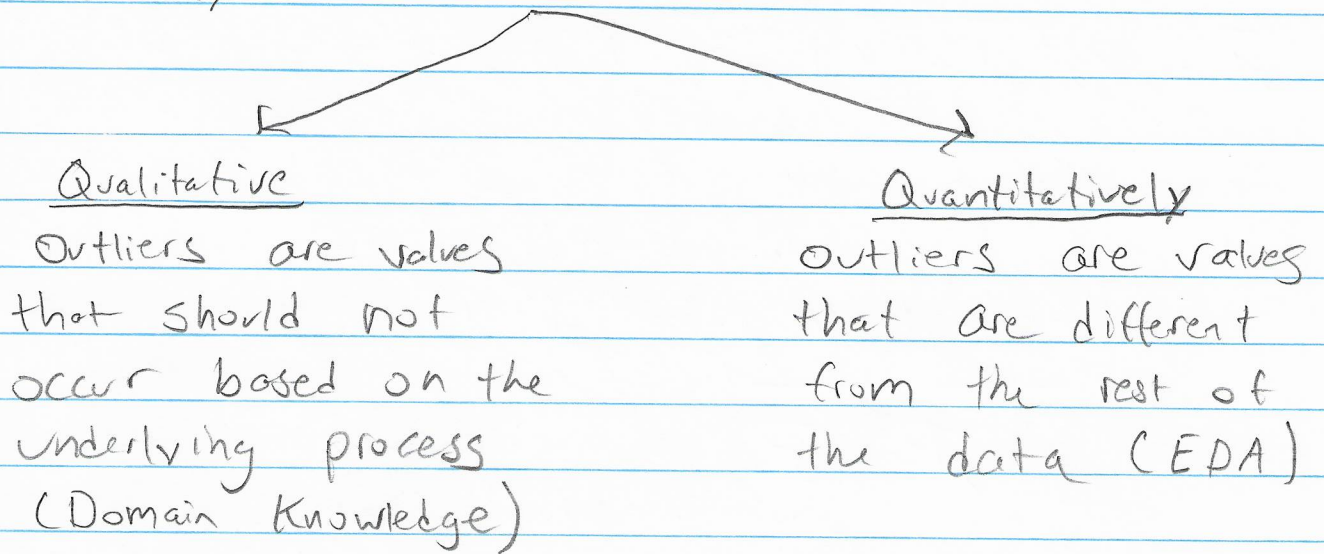Measures of Correlation

Pearson Correlation Coefficient $(R^2)$

$$\text{Covariance}(A_j, A_k) = \frac{\sum_{i=1}^{n} (X_{ij} - \bar{A_j})(X_{ik} - \bar{A_k})}{n-1}$$

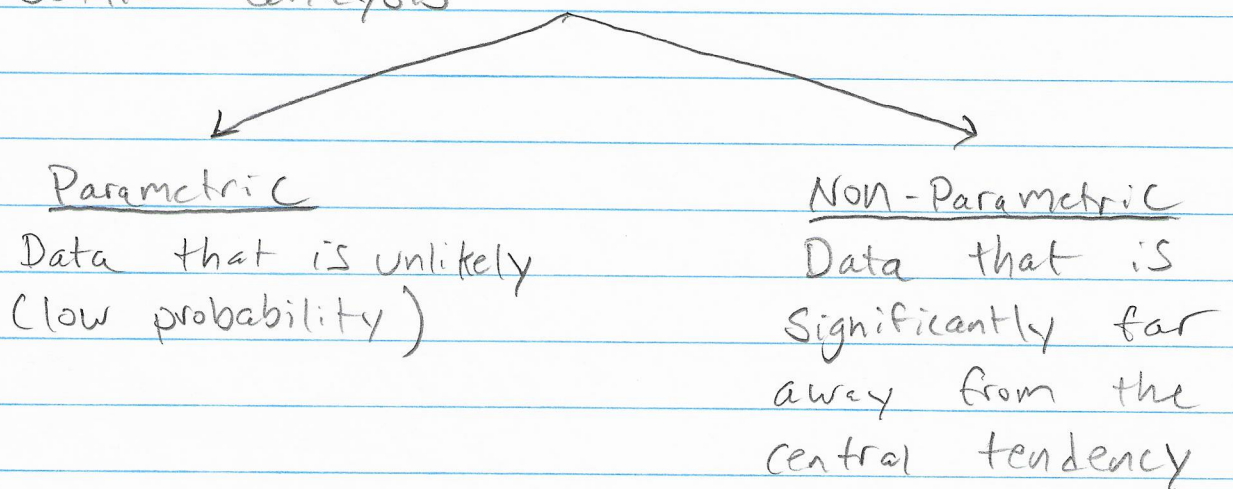$$\text{Correlation}(A_j, A_k) = \frac{\text{Cov}(A_j, A_k)}{\sigma_{A_j} \sigma_{A_k}}$$

Nominal Attributes $\rightarrow$ Chi Squared Test $(x^2)$

④ Outlier Analysis

Two general approaches to outlier analysis

<u>Qualitative</u>
Outliers are values that should not occur based on the underlying process (Domain Knowledge)

<u>Quantitatively</u>
Outliers are values that are different from the rest of the data (EDA)

Two general approaches to quantitative outlier analysis

<u>Parametric</u>
Data that is unlikely (low probability)

<u>Non-Parametric</u>
Data that is significantly far away from the central tendency
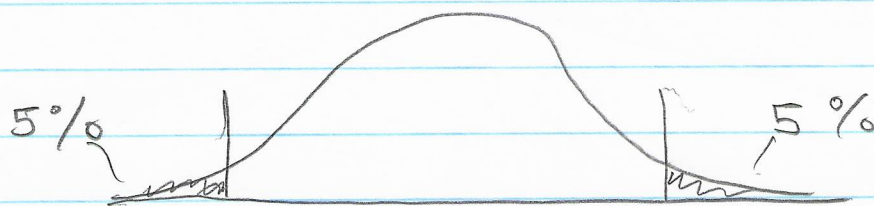
# Process for Parametric Outlier Analysis

1) Fit the data to an appropriate theoretical distribution

   e.g. Normal, Binomial, Exponential

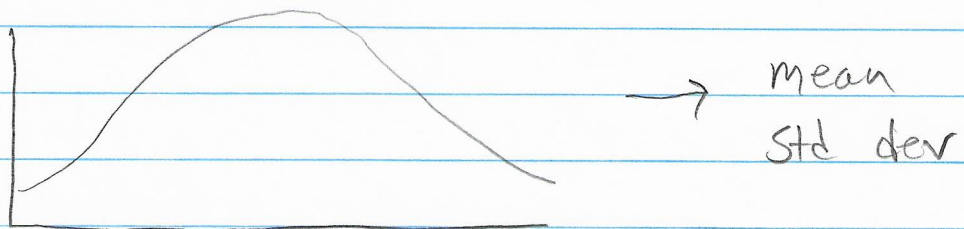2) Select a probability cutoff for outliers

   e.g. 10%

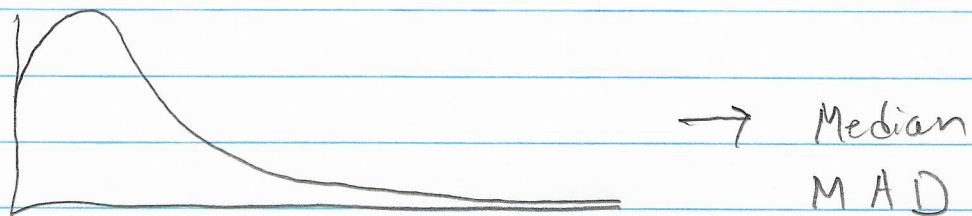3) Determine the cutoff value based on the PDF/CDF

Process for Non-Parametric Outlier Analysis

1) Determine the appropriate measure
   for central tendency and dispersion

Symmetric Histogram



→ mean
  std dev

Skewed / long-tail histogram



→ Median
  MAD

2) Determine the appropriate threshold
   based on visual exploration and
   domain knowledge

Cut off = central tendency $\pm$ K spread