

① Noise

② Outliers

$A_1$	$A_2$
7.85	Good
6	Good
5.1	Very bad

$A_{20}$	
Avg	Y
Bad	N
Bad	Y

Senex in corona

Throw out outliers! (Out of scope)

$f_1 \quad f_2 \quad \dots \quad f_{15} \quad \dots \quad f_{50} \quad T(A|N)$





Missing values :- (i) Drop the missing records  
(ii) If too much data is missing  
one fix: eg: If out of 10k, 8k have values  
1 take average of feature

Another fix:

Let's say out of 10k, 6k are positive examples  
and out of 6k, 5.8k have  
values, 200 will get averaged  
6k

Repeat for all classes (eg: positive/negative)  
1, 2, 3, ... - 10

→ Another way:  $f_1 \quad f_2 \quad \dots \quad f_{15} \quad \dots \quad f_{49}$   $\frac{\text{Target}}{Y/N}$   
(1/0)

$D_1$   
 $D_2$   
 $D_3$   
 $D_4$

find nearest neighbours of this point in  $\mathbb{R}^{50}$

• find 5 NN and check their feature if missing  
if not, take average!

### Duplicate data

diff sources may have same data  
→ the data point may skew the data and bias the learning if it is repeated → misleading prediction by kind of saying more data  
(Sometimes good tho, when less features or during random oversampling to clear bias due to dataset)

Similarity & dissimilarity measures

eg: Distance : Euclidean

# Similarity and Dissimilarity Measures

## □ Similarity measure

- Numerical measure of how alike two data objects are.
- Is higher when objects are more alike.
- Often falls in the range [0,1]

## □ Dissimilarity measure

- Numerical measure of how different two data objects are
- Lower when objects are more alike
- Minimum dissimilarity is often 0
- Upper limit varies

□ **Proximity** refers to a similarity or dissimilarity

01/22/2018

Introduction to Data Mining, 2nd Edition

29

## Minkowski Distance

- Minkowski Distance is a generalization of Euclidean Distance

$$d(x, y) = \left( \sum_{k=1}^n |x_k - y_k|^r \right)^{1/r}$$

*(Handwritten:  $n=2$  with a checkmark)*

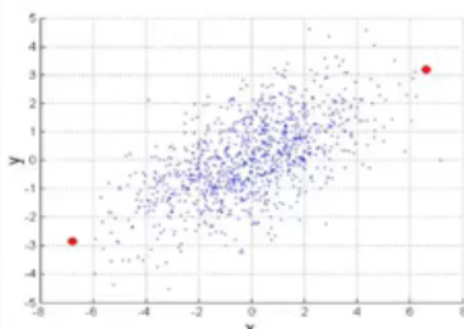
Where  $r$  is a parameter,  $n$  is the number of dimensions (attributes) and  $x_k$  and  $y_k$  are, respectively, the  $k^{\text{th}}$  attributes (components) or data objects  $x$  and  $y$ .

## Minkowski Distance: Examples

- $r = 1$ . City block (Manhattan, taxicab,  $L_1$  norm) distance.
  - A common example of this is the Hamming distance, which is just the number of bits that are different between two binary vectors
- $r = 2$ . Euclidean distance
- $r \rightarrow \infty$ . "supremum" ( $L_{\max}$  norm,  $L_{\infty}$  norm) distance.
  - This is the maximum difference between any component of the vectors
- Do not confuse  $r$  with  $n$ , i.e., all these distances are defined for all numbers of dimensions.

## Mahalanobis Distance

$$\text{mahalanobis}(x, y) = (x - y)^T \Sigma^{-1} (x - y)$$



$\Sigma$  is the covariance matrix

For red points, the Euclidean distance is 14.7, Mahalanobis distance is 6.

many distance metrics from previous class

## Common Properties of a Distance

- Distances, such as the Euclidean distance, have some well known properties.

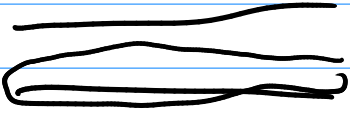
1.  $d(x, y) \geq 0$  for all  $x$  and  $y$  and  $d(x, y) = 0$  only if  $x = y$ . (Positive definiteness)
2.  $d(x, y) = d(y, x)$  for all  $x$  and  $y$ . (Symmetry)
3.  $d(x, z) \leq d(x, y) + d(y, z)$  for all points  $x, y$ , and  $z$ . (Triangle Inequality)

where  $d(x, y)$  is the distance (dissimilarity) between points (data objects),  $x$  and  $y$ .

- A distance that satisfies these properties is a **metric**

Machine Learning: Feature engineering  
Og. Age height & weight into BMI  
for finding diseases

Sampling Cif too many points

✓  → Sampling with replacement  
w/o replacement

## Data Preprocessing

- Aggregation
- Sampling
- Dimensionality Reduction
- Feature subset selection
- Feature creation
- Discretization and Binarization
- Attribute Transformation

→ get all data in a txn database for a month of January

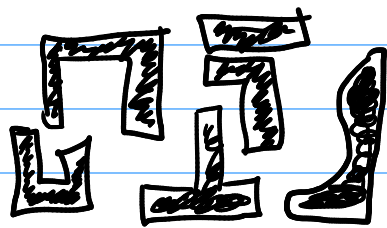
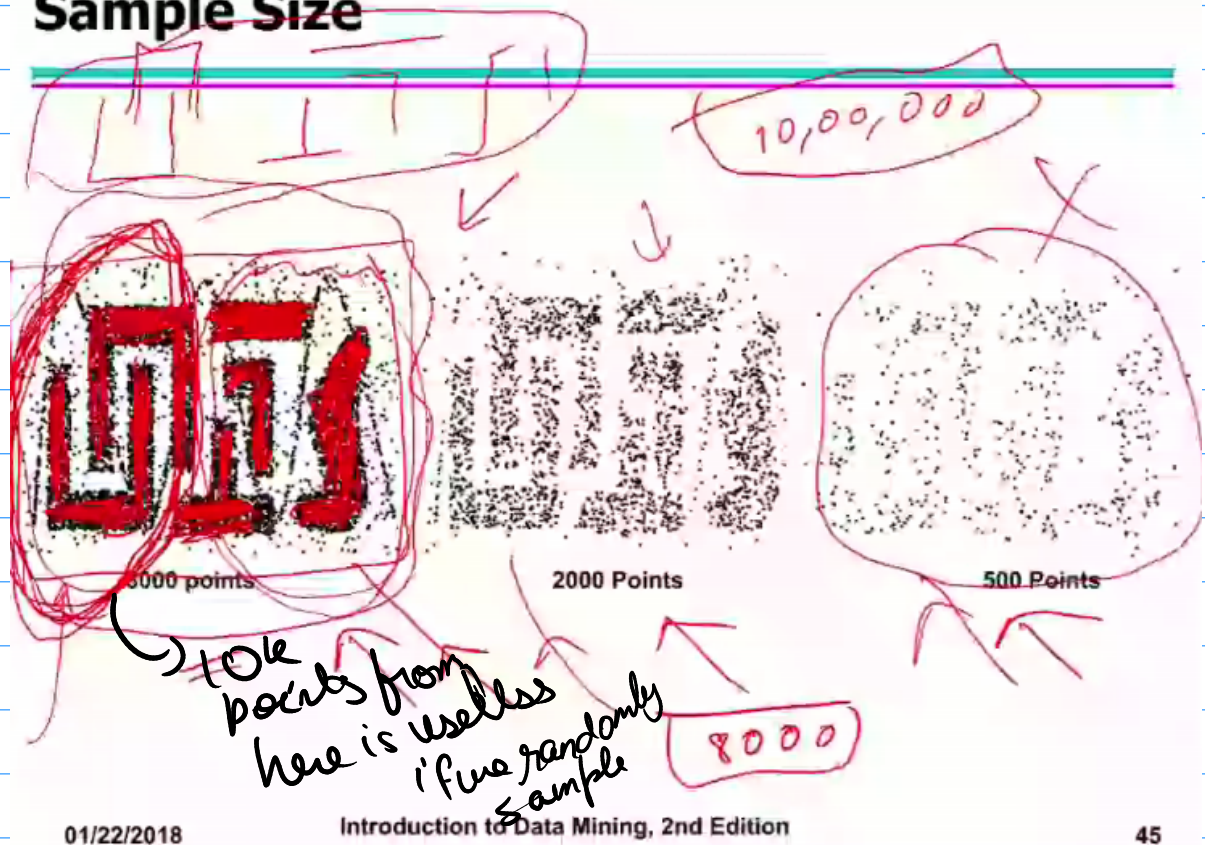
(Aggregation) like concatenation of datapoints

Another way of aggregation  
⇒ involve algebra of features

eg: BMI

We also need to note that <sup>while</sup> sampling points, we can't sample from a single place

## Sample Size



So we do something  
called stratified sampling

You may not know this, you need to have a reasonable

understanding of  
the data beforehand  
for stratified sampling

Stratified sampling is all about dividing your space into partitions and sample from each (maybe equally?)

But you may be unlucky then as well  
since you might sample a part of the partition  
(But if can we do there)



Another thing we should note: PCA

→ Lets say we have 1000 features, and as data scientists, we don't know what feature is useful and what feature is useless  
(eg: in maybe hi-resolution images or sth where the features aren't that very apparent)

★ We need an algorithm to reduce the number of features.