

Data Pre-Processing

Cleaning, integration, exploration, reduction/transformation, visualization....

- **Readings:**
 - KJ Ch 3, 19

Also see:

- T. Dasu and T. Johnson. Exploratory Data Mining and Data Cleaning. John Wiley & Sons, 2003
- Garcia, Luengo, Herrera, “Data Preprocessing in Data Mining”, Springer 2015.
- Xu Chu, Ihab Ilyas, [Sanjay Krishnan](#), Jiannan Wang, “Data Cleaning: Overview and Emerging Challenges” SIGMOD Tutorial, Jun. 2016.
slides at <https://sites.google.com/site/datacleaningtutorialsigmod16/home/slides>
- **Explore** segmentationOriginal (KJ, 3.1) and German Credit Card datasets
- **Companies:** e.g. <https://www.trifacta.com/>
 - Tons of data warehousing companies

Types of Data

- A data set is a collection of data objects/records
 - Each object is described by several features/attributes
- Data Types
 - Nominal
 - eye color, hobby
 - Binary: special case
 - Ordinal
 - rankings (e.g., taste of potato chips on a scale from 1-10), grades
 - Interval
 - speed, temperature in Celsius
- Categorical: Nominal or Ordinal
- Others: ID, DATE, text/strings, graphs,...

Different data types often need different ways of handling/modeling.

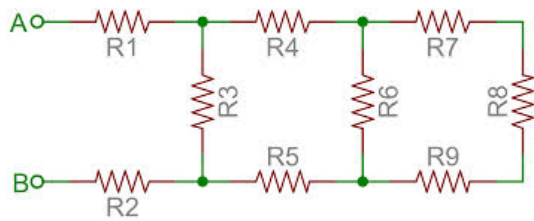
e.g. Proportional odds model for ordinal regression.

Data Issues

- Quantity
- Quality and adequacy
- Acquiring "labels"
- Big Data Issues

How Much Data do You Need? Effort?

- Quantity and quality
- Statistically large vs. computationally large
- Effort?
 - depends on variety of data sources, format, quality e.g. missing/incorrect items,.. Ownership issues,..
 - ETL vs. modern data wrangling



Why Preprocess Data

- **GIGO!**
 - data may be incomplete, inconsistent, noisy; have outliers, or simply too large
- **Why is data dirty?**
 - **Incomplete data** may come from
 - Not available or “Not applicable” data value when collected
 - Thoughtless entry (e.g. 0 vs. missing)
 - **Noisy data** (incorrect values) may come from
 - Faulty data collection instruments
 - Human or computer error at data entry
 - HEB, shoulder surgery, ..
 - Out-of-date
 - **Inconsistent data** may come from
 - Different data sources; formats
 - Inconsistent rules e.g. hotel price on phone vs. internet
 - **Duplicate records** need to be eliminated

Major Preprocessing Steps

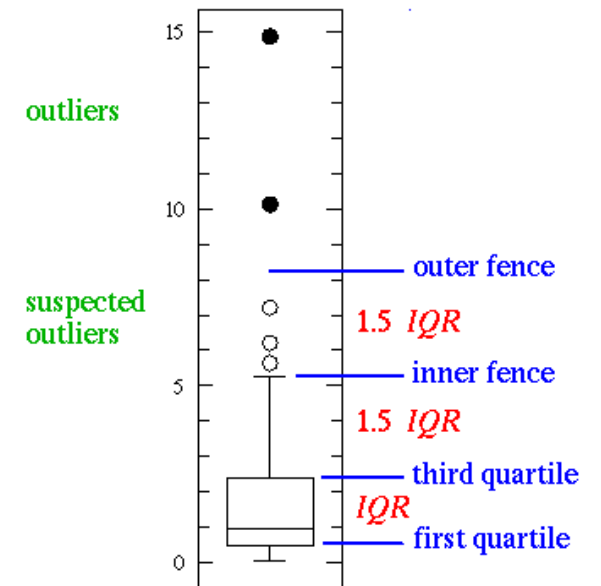
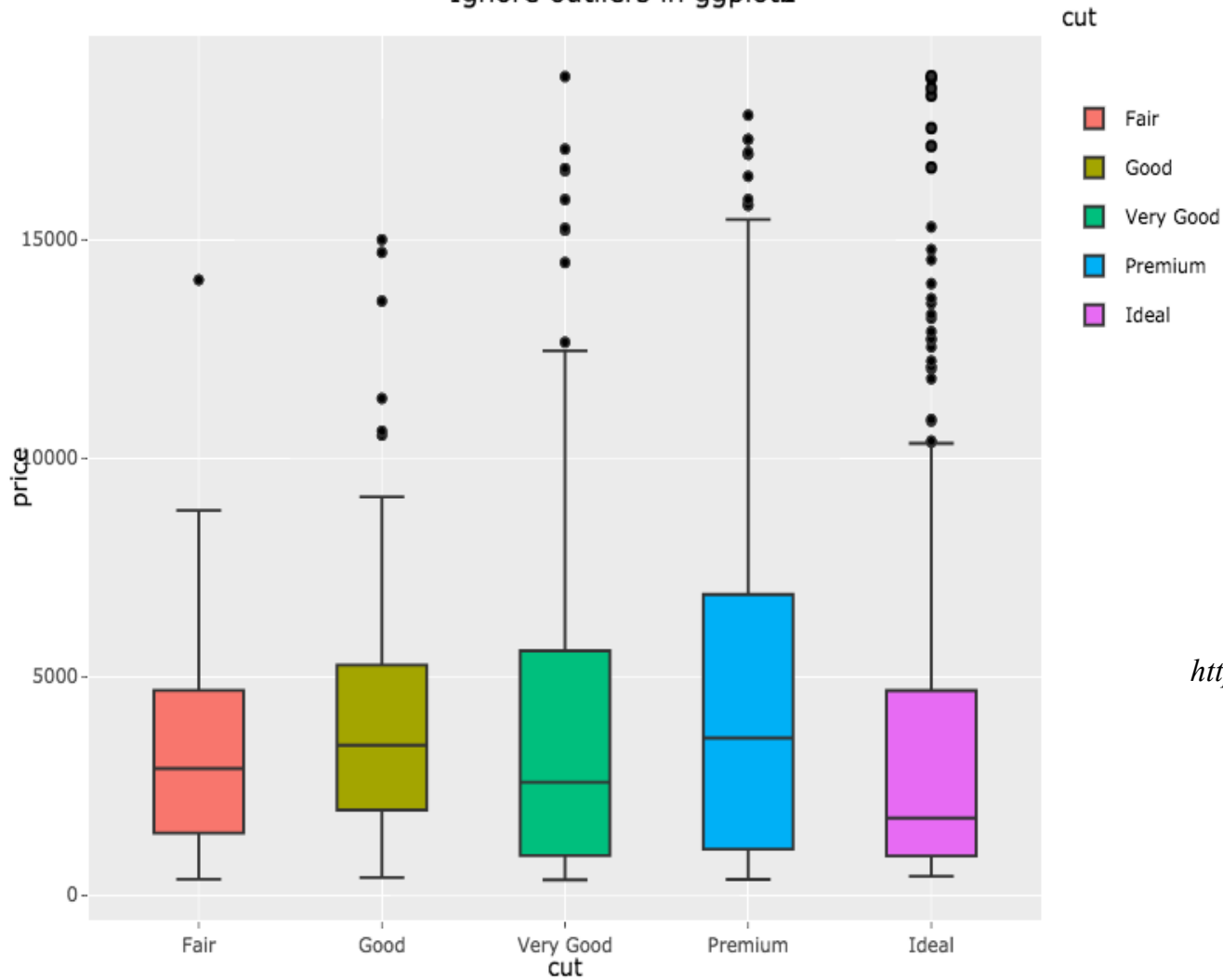
- 1) **Data cleaning**; sanity checks, consistency (already done by ETL tools if data is from warehouse)
 - 2) **Exploratory Data Analysis** (Often based on a sample)
 - 1) Fill missing values, remove noise and outliers
 - 2) transformation/scaling
 - 3) **Data reduction**
 - 1) Of records (sampling)
 - 2) Of attributes (feature selection/extraction)
 - 4) **Visualization**
-
- Often takes over 90% of a project's time!
 - steps 2-4 often revisited after modeling.

Before You Clean the Data..

- .. Do a quick summarization/visualization
 - Single “input” variable summaries
 - Variable type, mean, range, %missing, skewness, histograms, boxplots,
 - Bivariate (X_i vs. Y or X_i vs. X_k) visuals
 - (scatter plots, correlation,..)

BoxPlot

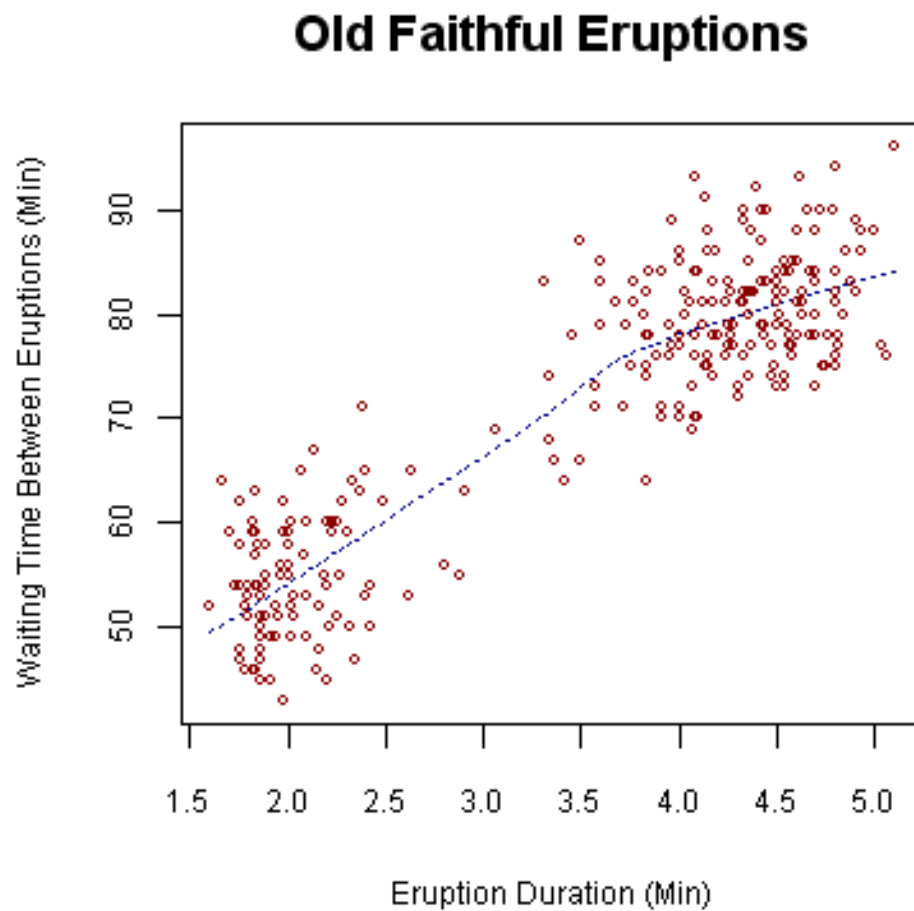
Ignore outliers in ggplot2



Using Plotly (interactive)
<https://plot.ly/ggplot2/box-plots/>

Scatterplot

- Old Faithful Example from Wikipedia

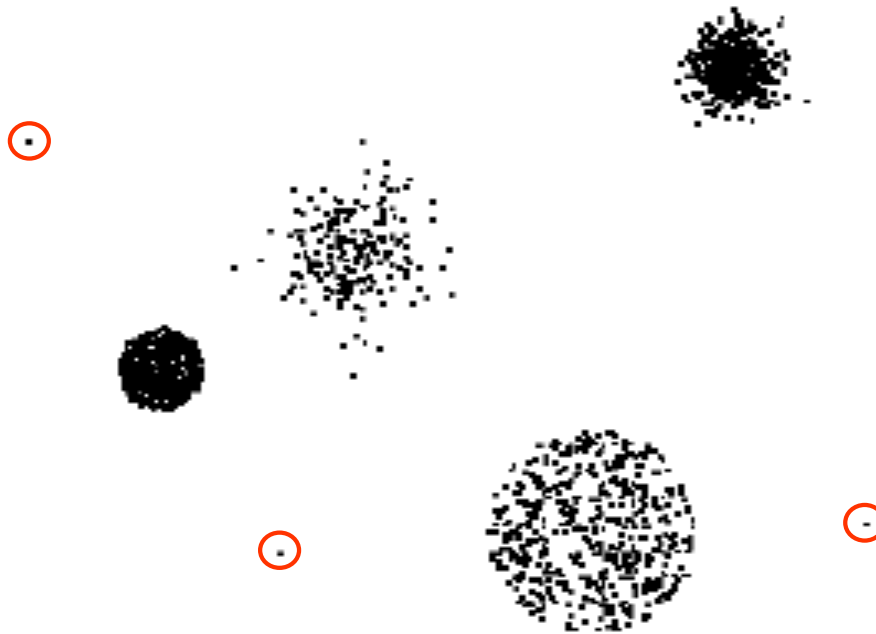


Data Cleaning I

- Dealing with Missing Values (Imputation)
 - Missing Completely at Random (MCAR)?
 - Vs. “informative missingness” (e.g doctor’s choices)
 - ignore record or attribute (often missing values are concentrated in a few instances or attributes)
 - Fill in missing values
 - fill with constant, mean or mode
 - conditional mean/ mode
 - Condition on values of a set of related variable
 - Use K-NN
 - “mathematically optimal” way?

Cleaning II: Handling Outliers

- Outliers are data objects with characteristics that are considerably different than the vast majority of the other data objects in the data set



Dealing with Outliers in “X”

- Probability based (old):
 - Estimate pdf of X, using e.g. Parzen windows or mixture of Gaussians
 - Identify low $p(x)$ points
- Discrimination based
 - Rule based, e.g.
 - » less than 1% for categorical variables
 - » Outside 3 sigma for gaussian looking numeric variables
 - Distance based: see if outlier score is $>$ threshold or not
 - » Score could be av. Distance of k-nearest neighbors; distance to the kth neighbor, etc.

Outliers in Y (robust statistics)

Identify outliers and eliminate
before applying model

OR

Use models that are little
affected by presence of a few
outliers

- trimmed means instead
of means
- alternatives to “squared
error” loss functions
 - e.g. Huber’s loss (quad
→ linear)

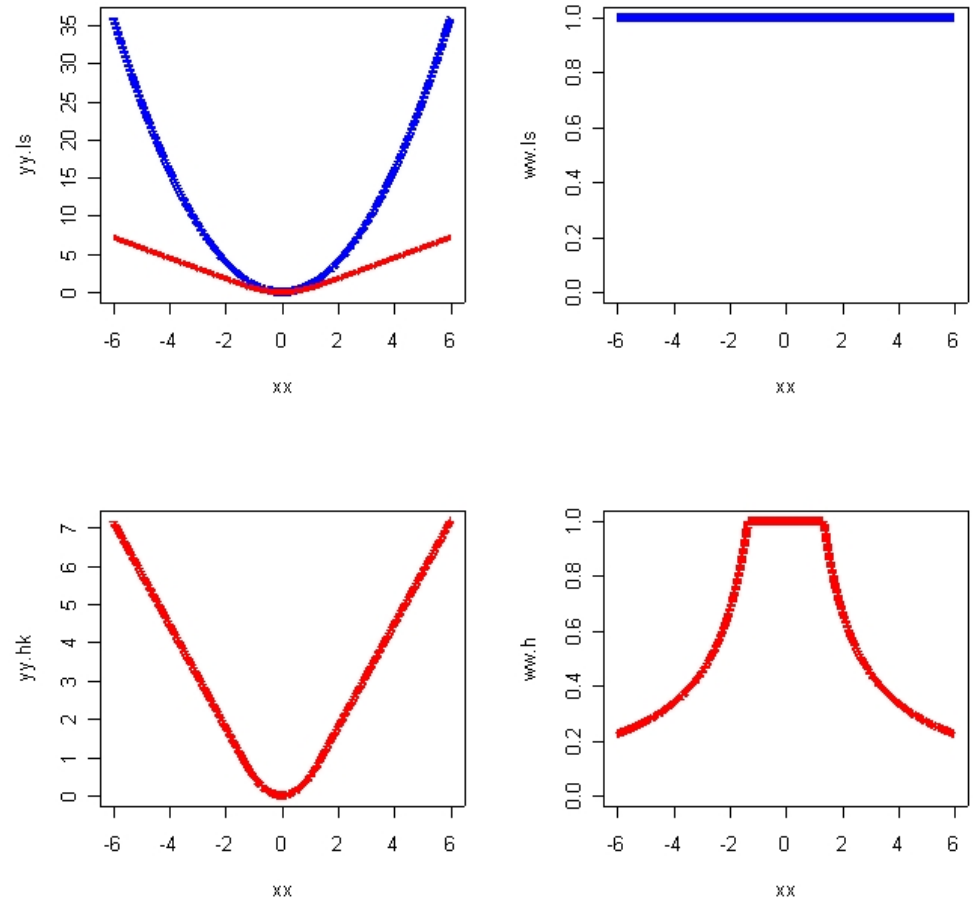


Fig: Plotted as a function of residual ($r = y - \hat{y}$):

Blue: Sq. error loss (left) and “equivalent” weights (right)

Red: Huber loss and equivalent weights if Sq. loss was used (right)

Data Transformation

- Scaling
 - Normalization by Linear scaling
 - Linear $[\min, \max] \rightarrow [0,1]$
 - Centering (e.g. Z-scoring: Normal/Gaussian $\rightarrow N(0, 1)$)
- (non-linear) transformation, e.g. to reduce skew or to show a simpler relationship between x and y (for example a power law shows up as a linear relationship in the log space).
 - Log;
 - square;
 - exponential

Data Reduction Methods

- **Why?**
 - get quicker answers
 - Reducing number of features may (substantially) improve results
!!
 - Reduces **“curse-of-dimensionality”**
 - When dimensionality increases, (randomly distributed) data becomes increasingly sparse in the space that it occupies
 - » Problematic for many types of analysis.
 - Collinearity a problem with MLR
 - Tools, e.g. compute all pairwise correlations (“pairs” in R)
 - Heuristics, e.g. eliminate variables till max pairwise correlation $<$ threshold

How to Reduce Data

- Reduce # of records or instances
 - Reduce # of attributes or features
 - Aggregate (in data cube)
 - Reduce resolution of an attribute e.g. discretization of interval variable.
-
- Note: Data reduction technique will affect quality as well as speed.

Basics of Simple Random Sampling

- Estimating the proportion of a binary choice
- Estimate is based on a sample size n , much smaller than size of underlying population N
- Answer can only be “Probably Approximately Correct”
 - Quantify via ϵ , the **margin of error**, and $1-\alpha$, the **confidence level**
of samples required depends on pre-specified “epsilon” and “alpha”

Estimating Sample Size

- Want: **within ε of mean** with high **probability $(1 - \alpha)$**
 - Normal: 90% of probability within $\pm 1.65 \sigma$ of mean
 - 95% of probability within $\pm 1.96 \sigma$ of mean
 - 99% of probability within $\pm 2.58 \sigma$ of mean
 - **Margin of error is ε ; critical value** (for standardized curve) is denoted by $z_{\alpha/2}$
 - » If $\alpha = 0.05$, then $z_{\alpha/2}$ is 1.96
- Minimum Sample size needed, $n = p(1-p) (z_{\alpha/2} / \varepsilon)^2$
 - **independent of N!!**
 - Use \hat{p} for p in above Eqn; if \hat{p} is unknown, use 0.5 for safe answer.

Sampling

A recent Texas Public Employees Association (TPEA) survey found that 11.7 percent of state employee households received public assistance in the past year. More than 16,000 state employees responded to our survey, and because our sample size was so large, our results can be considered representative of all general state government — approximately 149,000 employees — with a 99 percent confidence level and a 1 percent margin of error.

From AAS, April 24, 2015

Web Resources

- Many good web resources to understanding sampling, confidence intervals, etc.

Understanding confidence intervals:

<http://www.lordsutch.com/pol251/schacht-08-web.pdf>

Introduction to Probability (Undergrad course-notes from MIT).

<http://ocw.mit.edu/OcwWeb/Mathematics/18-05Spring-2005/LectureNotes/index.htm>

Other Sampling Issues

- Very effective when applicable
- good for estimate answer to aggregate query; but not for “needle in haystack” problems
- expensive! Not well supported in databases
- natural choice for **progressive refinement**; hypothesis testing

Reducing # of (Derived) Attributes/Features

- Feature selection (select a subset of original features)

VS

Feature extraction (use derived features, not original ones)

Why is feature selection often preferred to feature extraction?

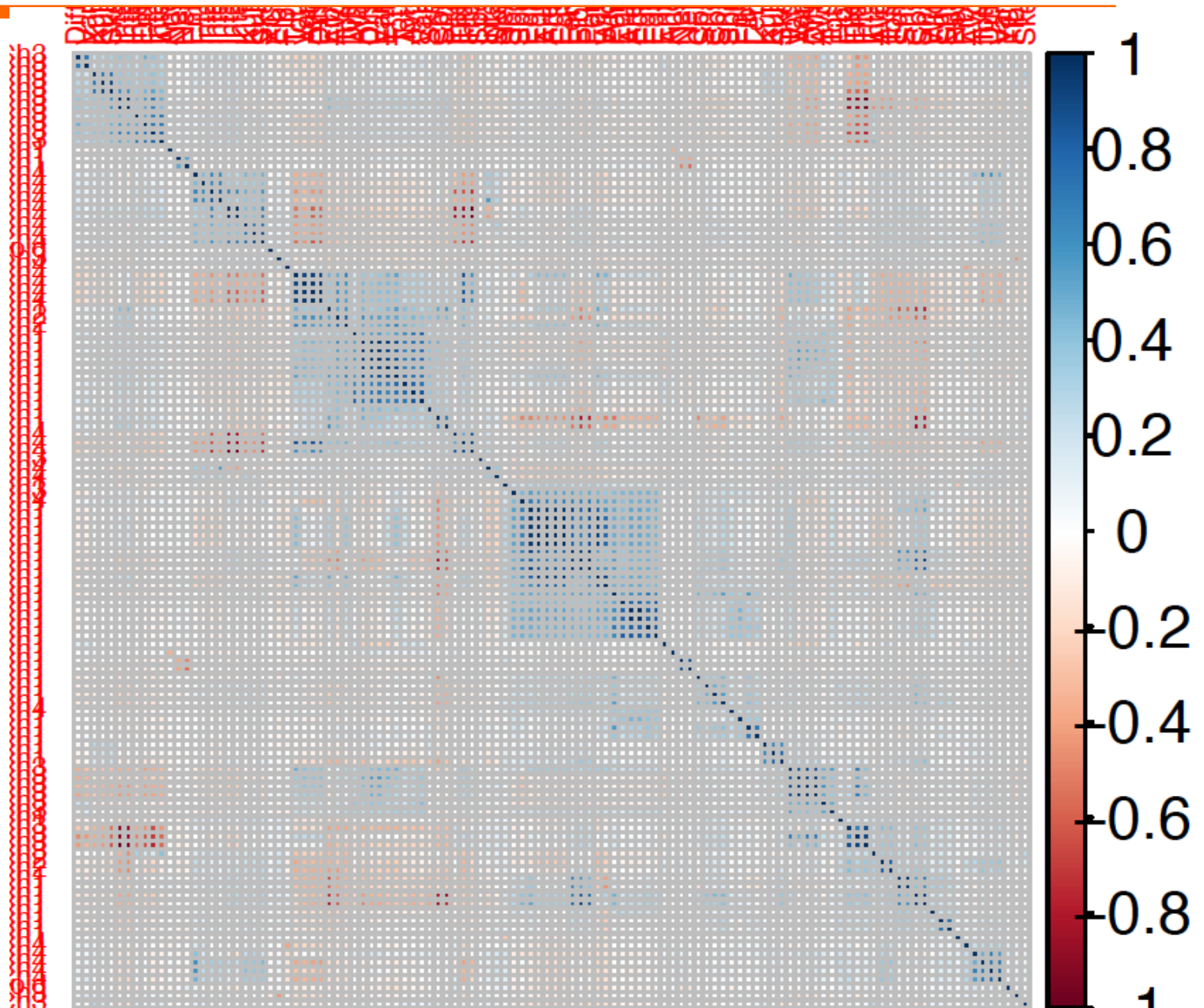
Selecting a Subset of Features

NP-complete, so use heuristics

- **Filter methods** : use intrinsic quality measure
e.g. correlation with other predictors (`cor(data)` in R);
correlation/Chi-sq with target; mutual info with target
- **Wrappers** (extrinsic evaluation)
 - Greedily evaluate using predictive model (e.g. decision tree)
 - **Search strategy for candidate sets to evaluate:**
 - Forward inclusion
 - Backward elimination
 - Stepwise (forward, but may remove predictors that no longer meet criterion)
- **Embedded** (feature selection part of model training, e.g LASSO)
- **Advanced Methods** : <http://featureselection.asu.edu/index.php>

Feature Selection Using Corrplot package

- See KJ
Fig 3.10



Perfect pairings

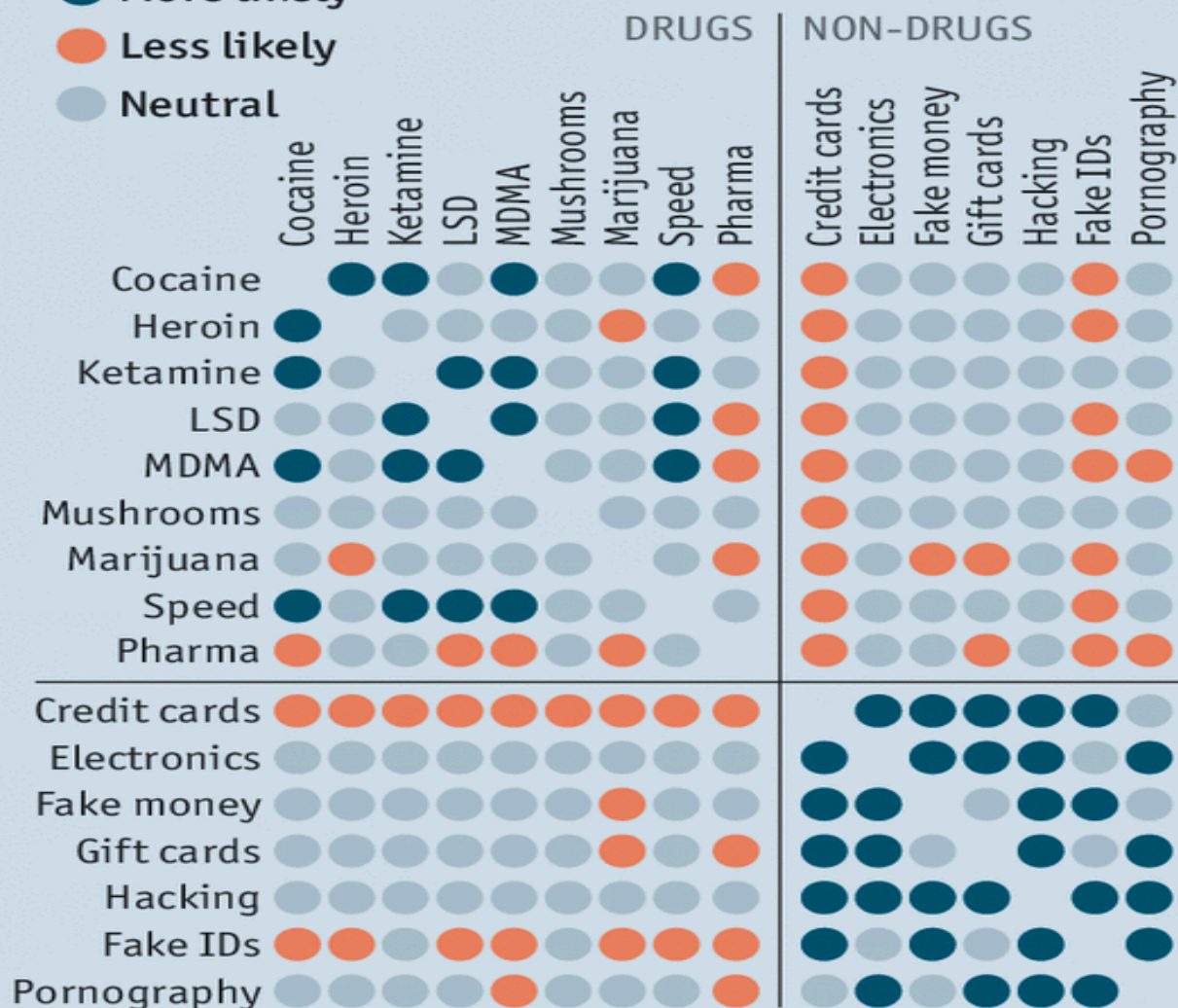
Likelihood that vendors selling one product on dark-web markets will sell another

December 2013–July 2015

● More likely

● Less likely

● Neutral



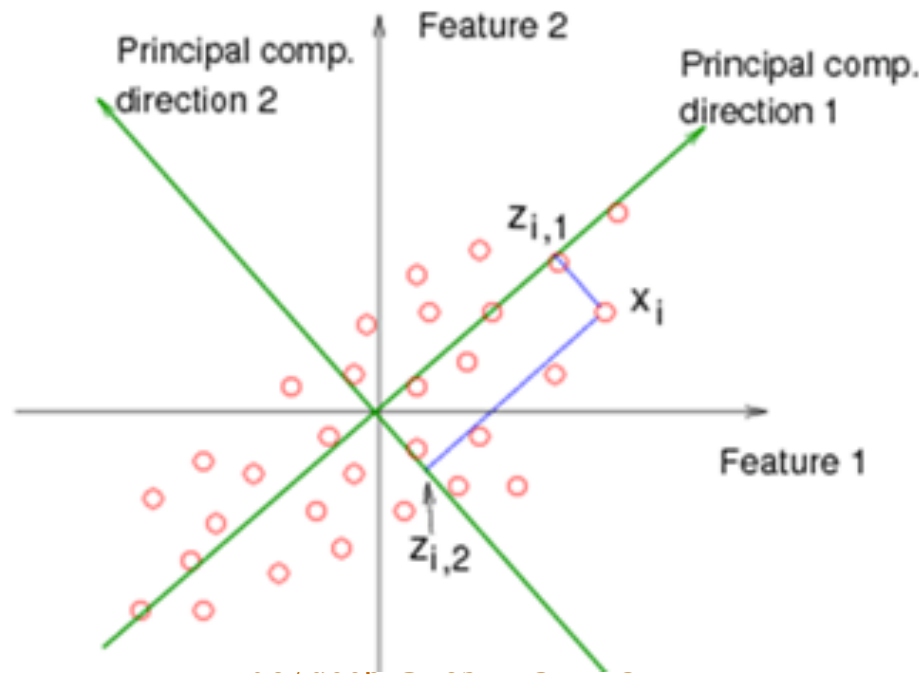
Sources: Gwern Branwen's dark-web archive; *The Economist*

Feature Extraction Choices

- Linear
 - Unsupervised : PCA
 - 5 functions to do Principal Components Analysis in R
<http://gastonsanchez.com/blog/how-to/2012/06/17/PCA-in-R.html>
 - Supervised:
 - Fisher's Linear Discriminant (classification)
 - Canonical Correlation (regression)
- Non-Linear
 - Unsupervised : Principal Curves, Sammon's Map, Kohonen's SOM
 - Supervised: Nonlinear discriminant analysis, e.g. using a multi-layered perceptron.

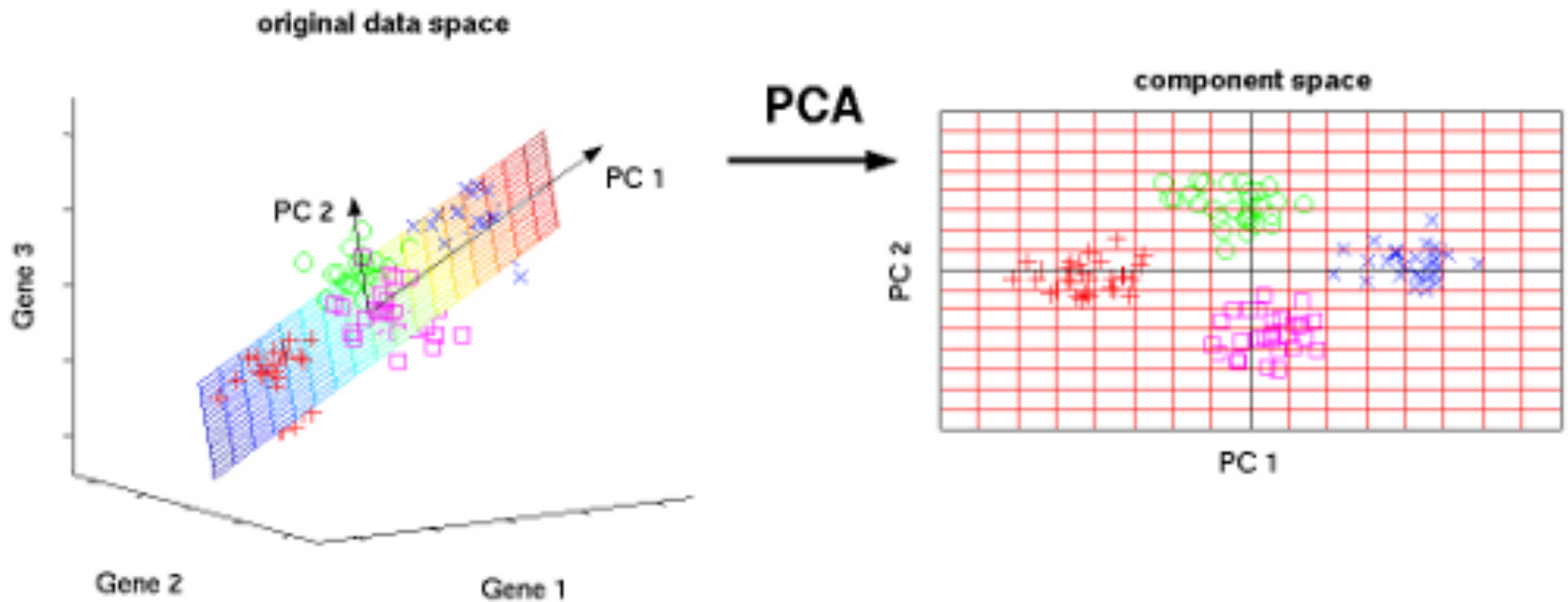
PCA

- Principal Components Analysis:
 - Reduce dimensions while retaining info about original data
 - PCA finds the best “subspace” that captures as much data variance as possible
 - optimal linear projection/reconstruction in MSE sense
 - Based on eigen-decomposition of data covariance matrix; can also be obtained sequentially



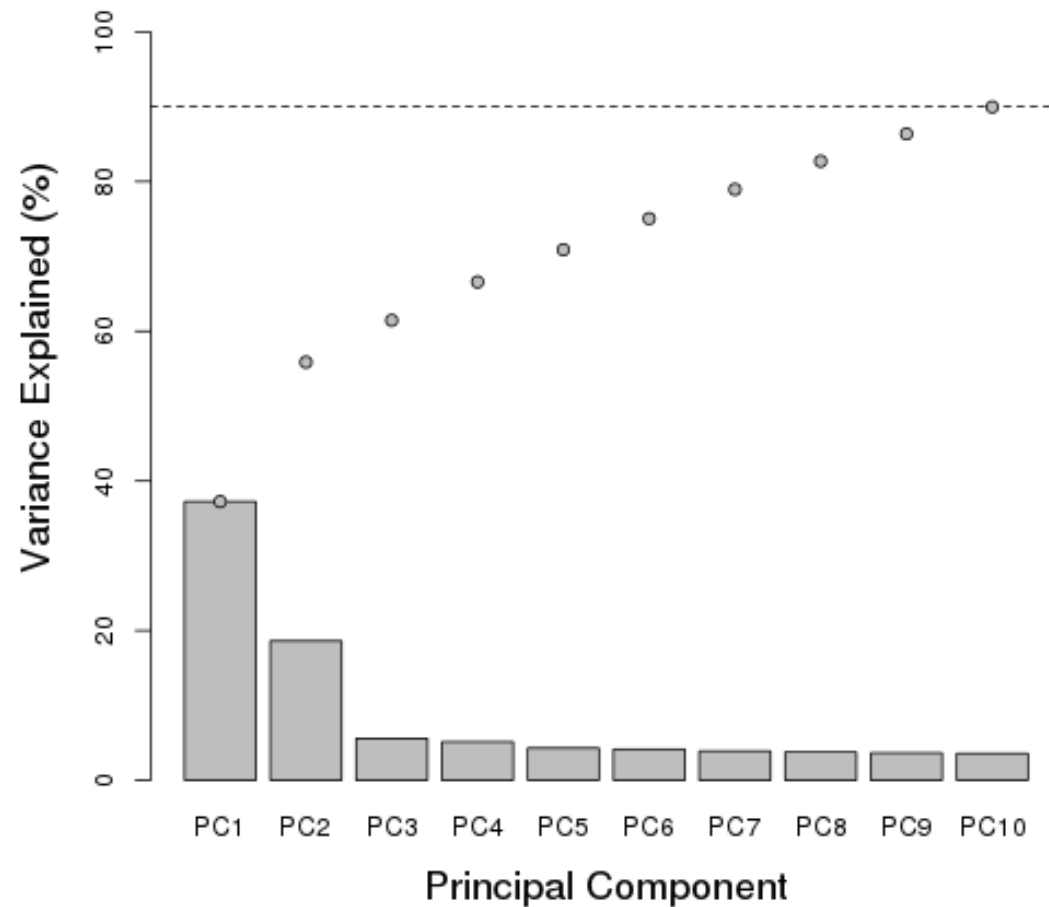
Another Example

From http://www.nlpca.org/pca_principal_component_analysis.html



PCs sorted by amount of variance explained

Example scree plot



Tensor Board Visualization

- PCA and t-SNE for MNIST (using Tensorboard)
 - 28x28 images
- <https://www.youtube.com/watch?v=eBbEDRsCmv4> 19:11 onwards
(Part of [Tensorflow dev summit 2017](#))

Visualize

- <http://setosa.io/ev/principal-component-analysis/>
- <http://setosa.io/ev/eigenvectors-and-eigenvalues/>

Singular Value Decomposition (SVD)

- Practical way of obtaining Principal components

| | day | We | Th | Fr | Sa | Su |
|----------|-----|---------|---------|---------|---------|---------|
| customer | | 7/10/96 | 7/11/96 | 7/12/96 | 7/13/96 | 7/14/96 |
| ABC Inc. | | 1 | 1 | 1 | 0 | 0 |
| DEF Ltd. | | 2 | 2 | 2 | 0 | 0 |
| GHI Inc. | | 1 | 1 | 1 | 0 | 0 |
| KLM Co. | | 5 | 5 | 5 | 0 | 0 |
| Smith | | 0 | 0 | 0 | 2 | 2 |
| Johnson | | 0 | 0 | 0 | 3 | 3 |
| Thompson | | 0 | 0 | 0 | 1 | 1 |

$$\mathbf{A} = \begin{bmatrix} 0.18 & 0 \\ 0.36 & 0 \\ 0.18 & 0 \\ 0.90 & 0 \\ 0 & 0.53 \\ 0 & 0.80 \\ 0 & 0.27 \end{bmatrix} \times \begin{bmatrix} 9.64 & 0 \\ 0 & 5.29 \end{bmatrix} \times \begin{bmatrix} 0.58 & 0.58 & 0.58 & 0 & 0 \\ 0 & 0 & 0 & 0.71 & 0.71 \end{bmatrix}$$

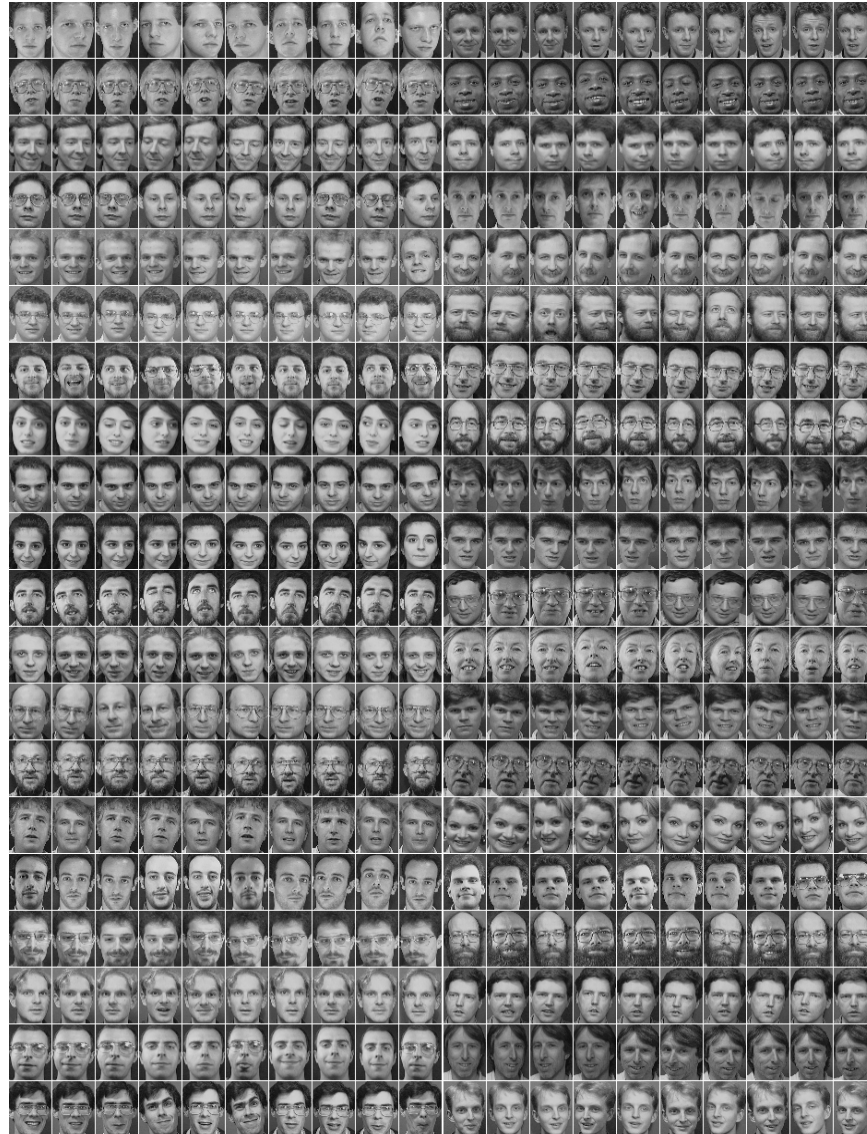
SVD

- Singular Value Decomposition (SVD)

$$A = U \times \Lambda \times V^T$$

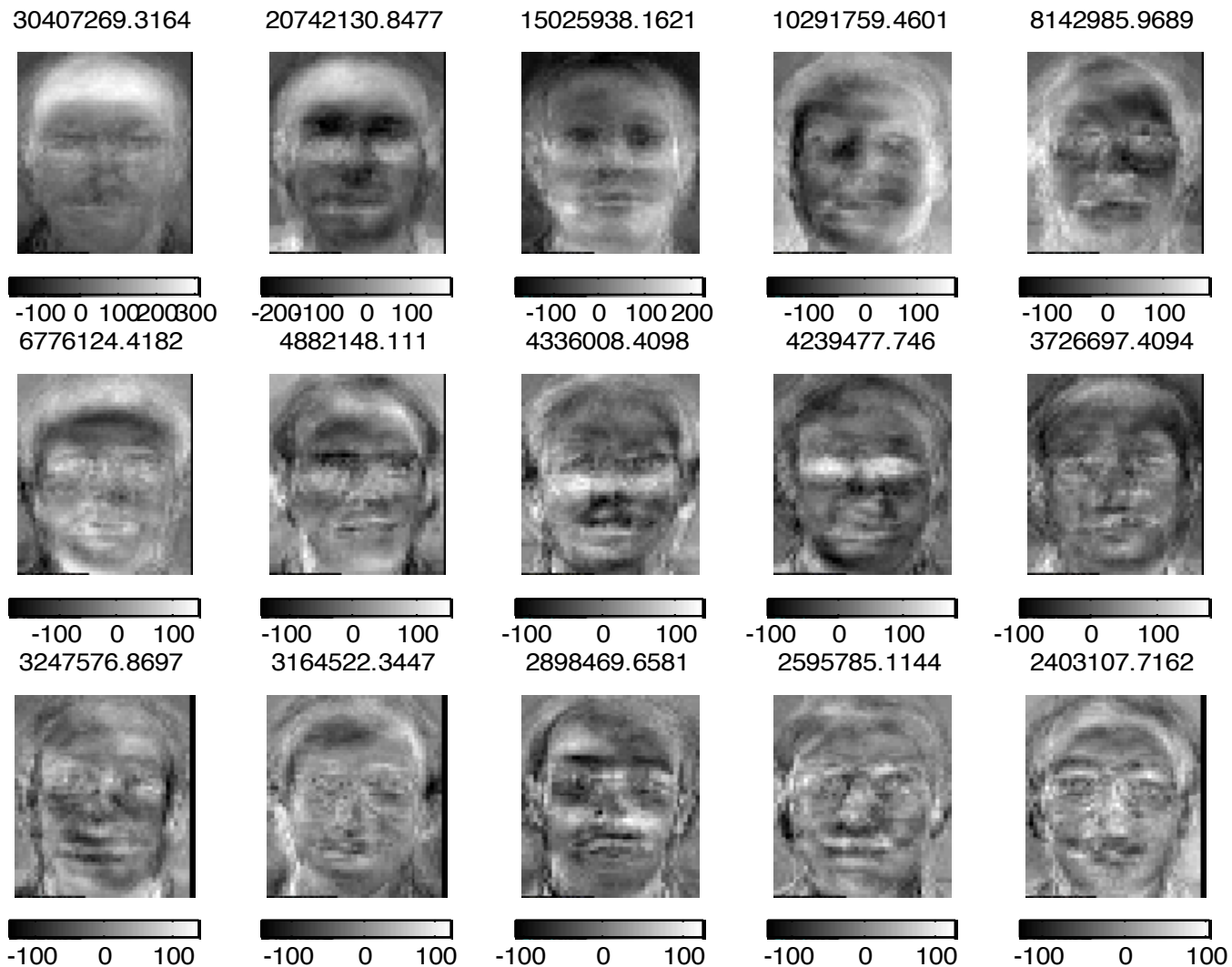
- for A = customer -day matrix, interpret
- U as customer-to-pattern similarity matrix
 - Columns of U are (orthonormal) eigen-“days”
 - Eigenvectors of AA^T
- V as day-to-pattern similarity matrix
 - Rows of V are (orthonormal) eigen-“customers”
 - Eigenvectors of $A^T A$
- is diagonal matrix of singular values (sorted)
 - (sq. root of eigen-values of AA^T Or $A^T A$)

Image Database



EigenFaces

<http://www.cs.princeton.edu/~cdecoreo/eigenfaces/>
<https://www.youtube.com/watch?v=jQOZrXZTXcw>



Reconstruction

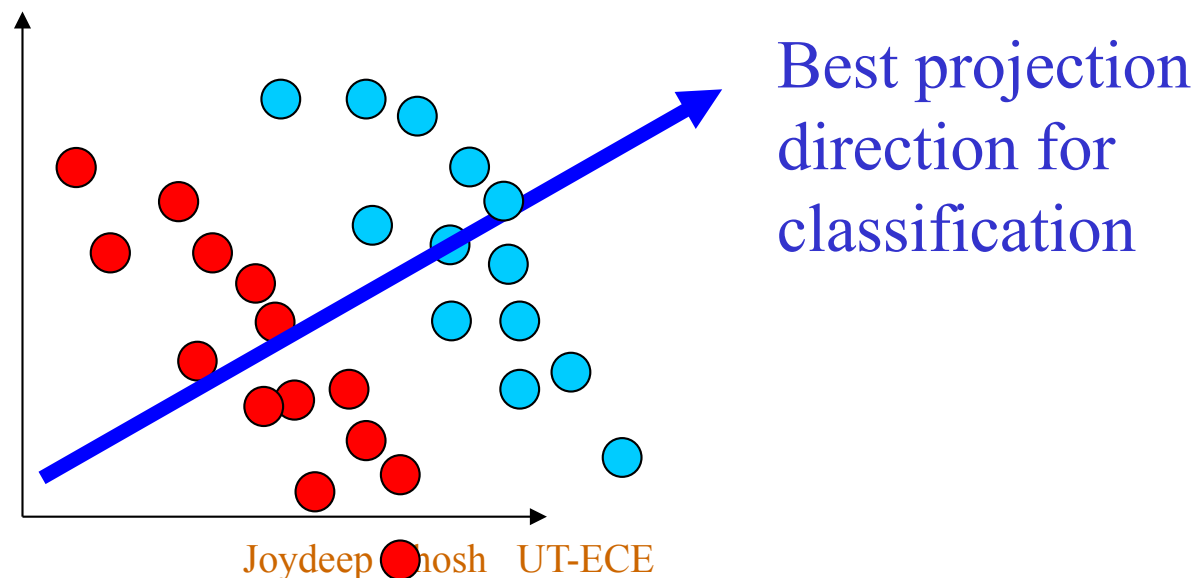
an example of reconstruction [from Turk and Pentland, 91] using different numbers of bases.



Effect of using more bases: each consecutive image uses 8 more bases

Linear Supervised Method: Fisher's Linear Discriminant (FLD)

- FLD finds the projection direction that best separates the two classes
- Multiple discriminant analysis (MDA) extends LDA to multiple classes
- For fun: Fisherfaces vs. Eigenfaces <https://www.youtube.com/watch?v=cKEF6iZAIMc> (David Mumford at 6:30)
- <https://www.youtube.com/watch?v=jQOZrXZTXcw>



Multi-dimensional Scaling (MDS);

Also see [Perceptual Mapping](#)

- When only pairwise distances (or similarities) are known
 - Objects may not be in Euclidean space
- Minimize a distortion measure (“stress”)

Table 1 Flying Mileages Between 10 American Cities

| Atlanta | Chicago | Denver | Houston | Los Angeles | Miami | New York | San Francisco | Seattle | Washington, DC | |
|---------|---------|--------|---------|-------------|-------|----------|---------------|---------|----------------|----------------|
| 0 | 587 | 1212 | 701 | 1936 | 604 | 748 | 2139 | 2182 | 543 | Atlanta |
| 587 | 0 | 920 | 940 | 1745 | 1188 | 713 | 1858 | 1737 | 597 | Chicago |
| 1212 | 920 | 0 | 879 | 831 | 1726 | 1631 | 949 | 1021 | 1494 | Denver |
| 701 | 940 | 879 | 0 | 1374 | 968 | 1420 | 1645 | 1891 | 1220 | Houston |
| 1936 | 1745 | 831 | 1374 | 0 | 2339 | 2451 | 347 | 959 | 2300 | Los Angeles |
| 604 | 1188 | 1726 | 968 | 2339 | 0 | 1092 | 2594 | 2734 | 923 | Miami |
| 748 | 713 | 1631 | 1420 | 2451 | 1092 | 0 | 2571 | 2408 | 205 | New York |
| 2139 | 1858 | 949 | 1645 | 347 | 2594 | 2571 | 0 | 678 | 2442 | San Francisco |
| 2182 | 1737 | 1021 | 1891 | 959 | 2734 | 2408 | 678 | 0 | 2329 | Seattle |
| 543 | 597 | 1494 | 1220 | 2300 | 923 | 205 | 2442 | 2329 | 0 | Washington, DC |

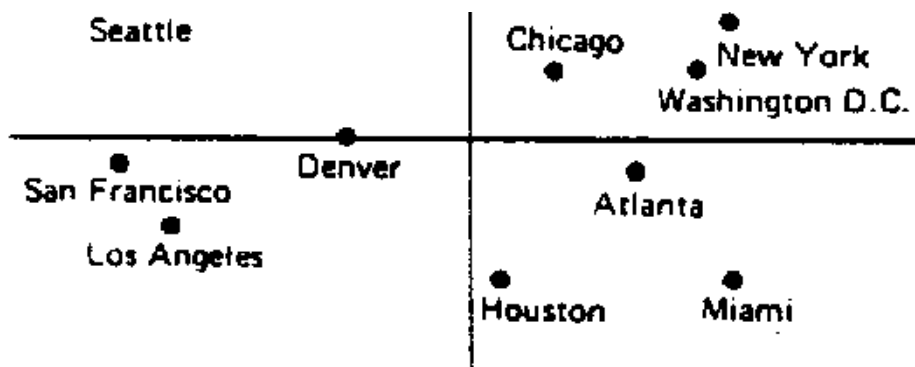
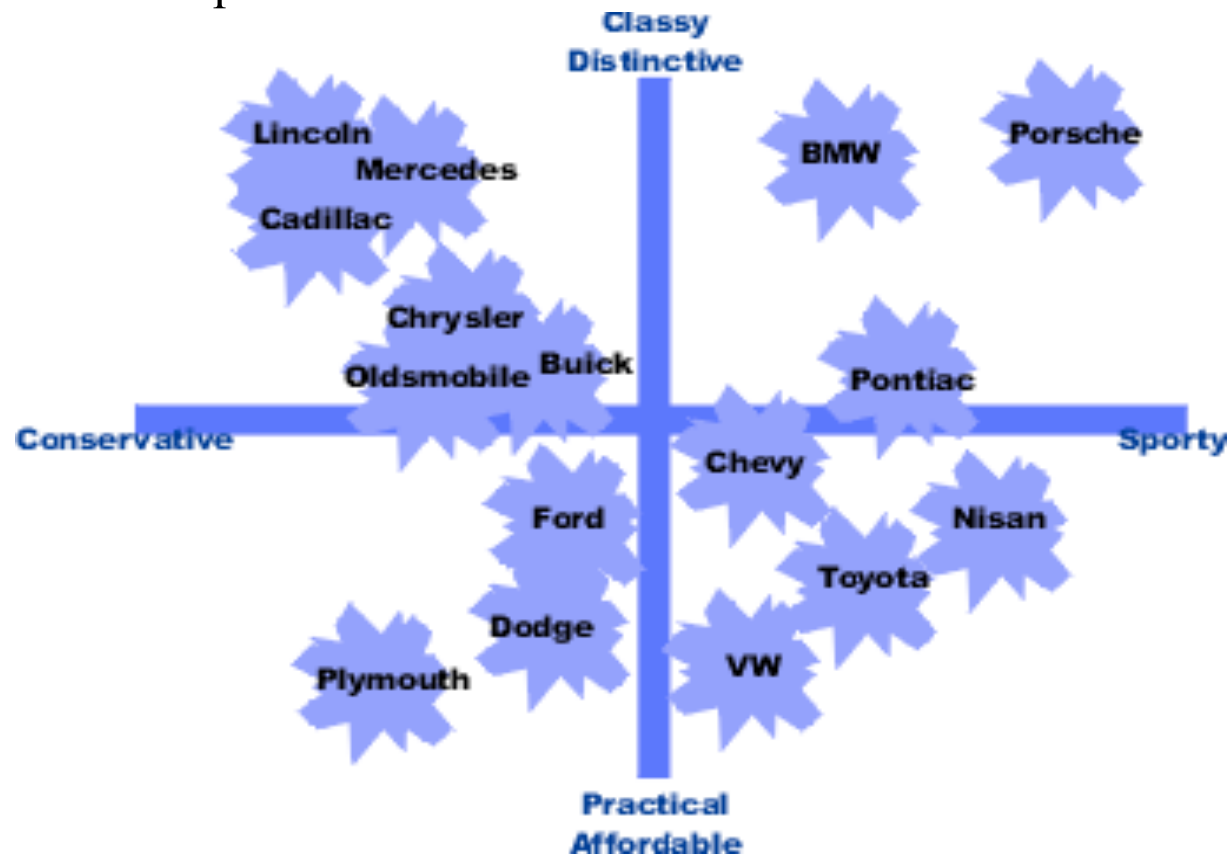


Figure 1 CMDS of flying mileages between 10 American cities.

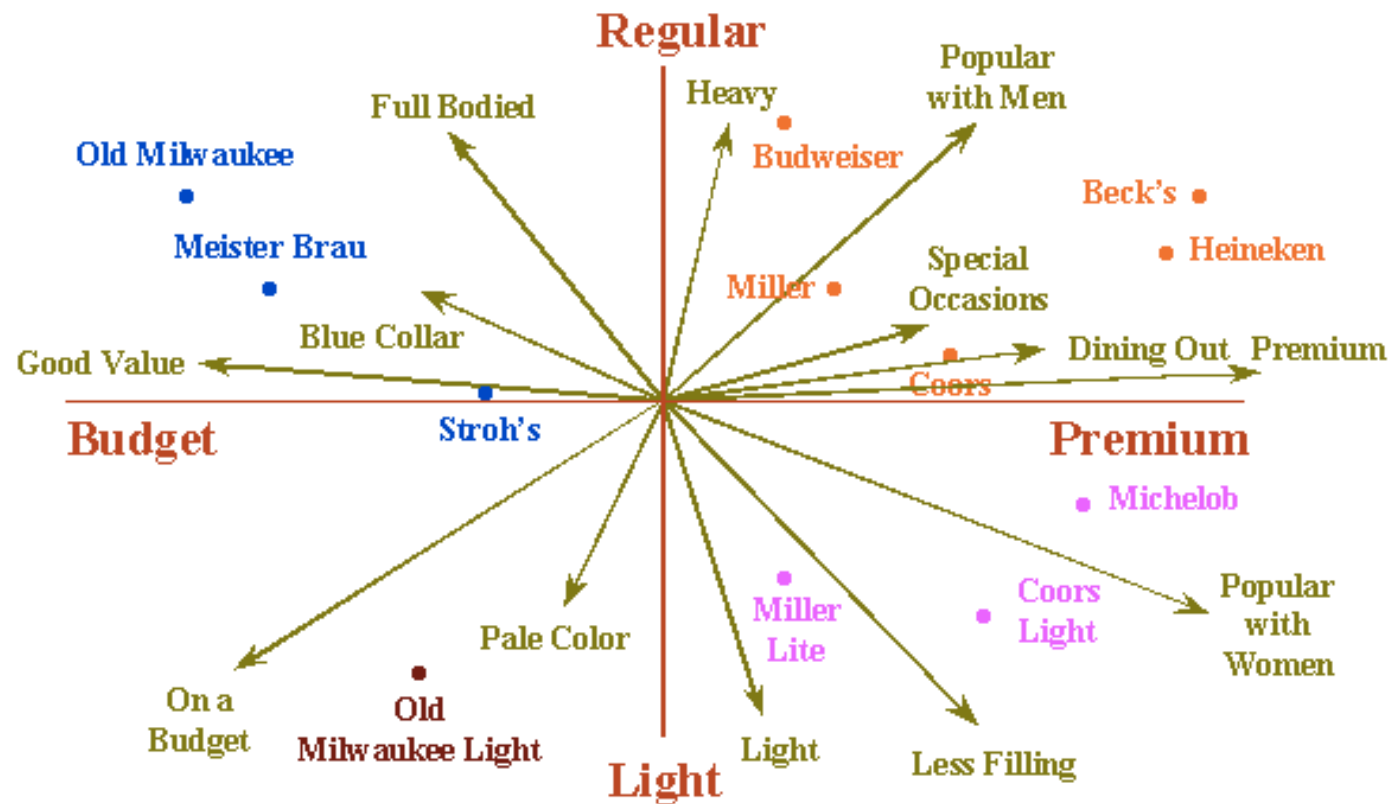
MDS also called Perceptual Mapping

- In the marketing literature
 - Used for brand positioning etc.
 - Wikipedia example below



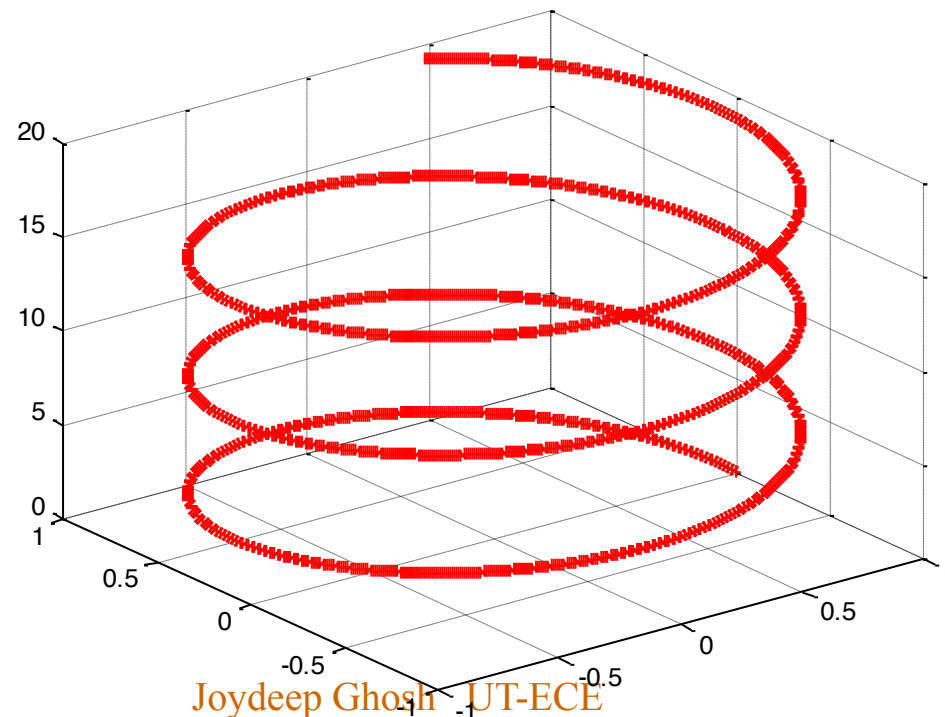
Beer Market

Perceptual Mapping



Deficiencies of Linear Methods

- Data may not be best summarized by linear combination of features
 - Example: PCA cannot discover 1D structure of a helix



Deficiencies of Linear Methods

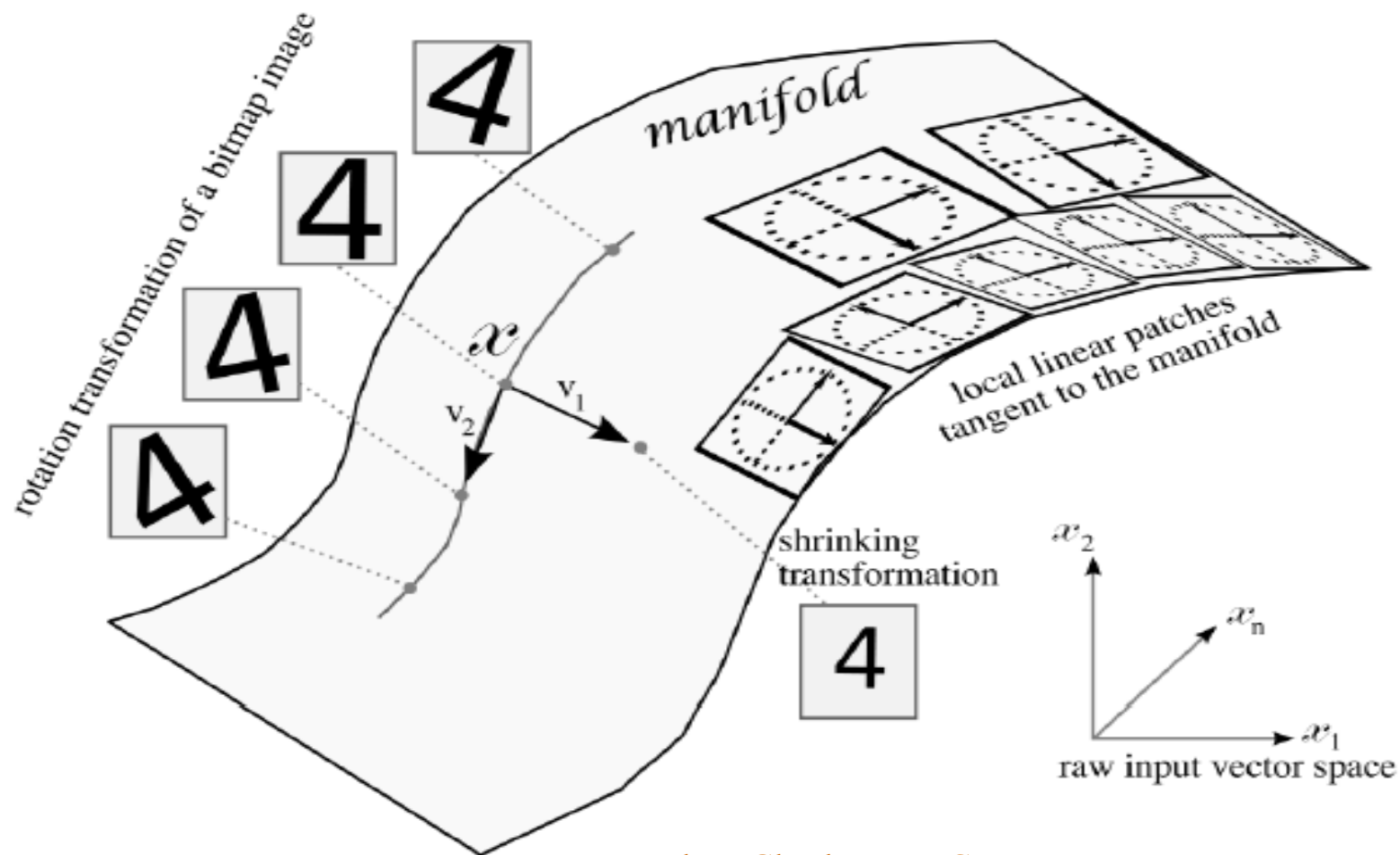
- Useful characteristics for real world data are often not linear combination of features
 - Example: poses of faces



- Example: when shall I get hit (from motion data)

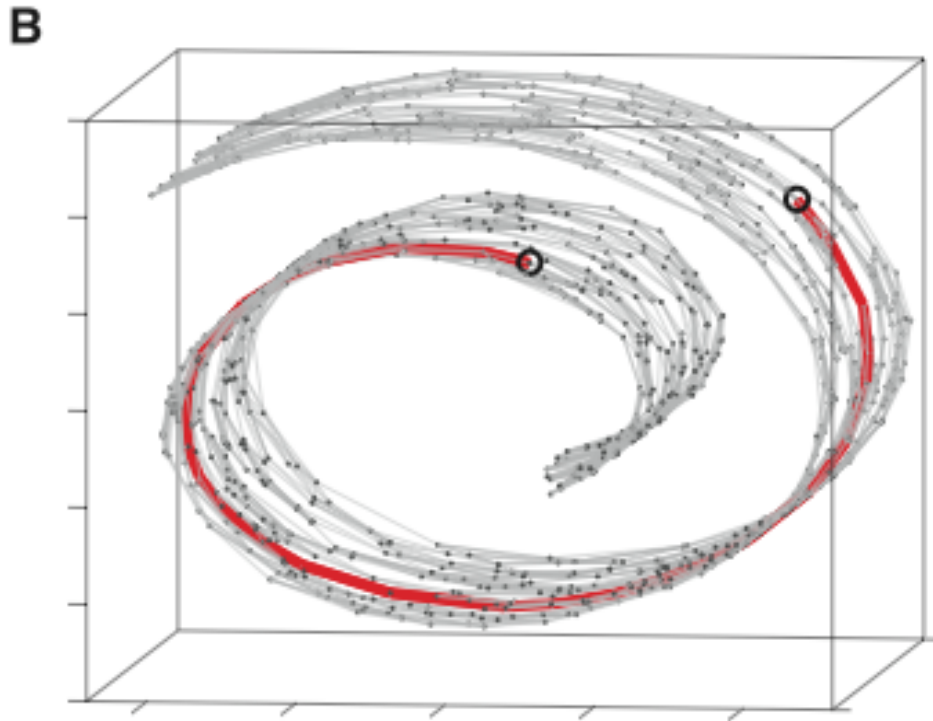


Handwritten Digits lie near a low-D manifold



Non-Linear Dimensionality Reduction

- Manifold based (ISOMAP, SOM,...)
 - The “swiss roll” below is an example of a manifold
 - Distance should be Measured on the Manifold and not in original space
- Non-linear versions of Multi-dimensional Scaling



UCL

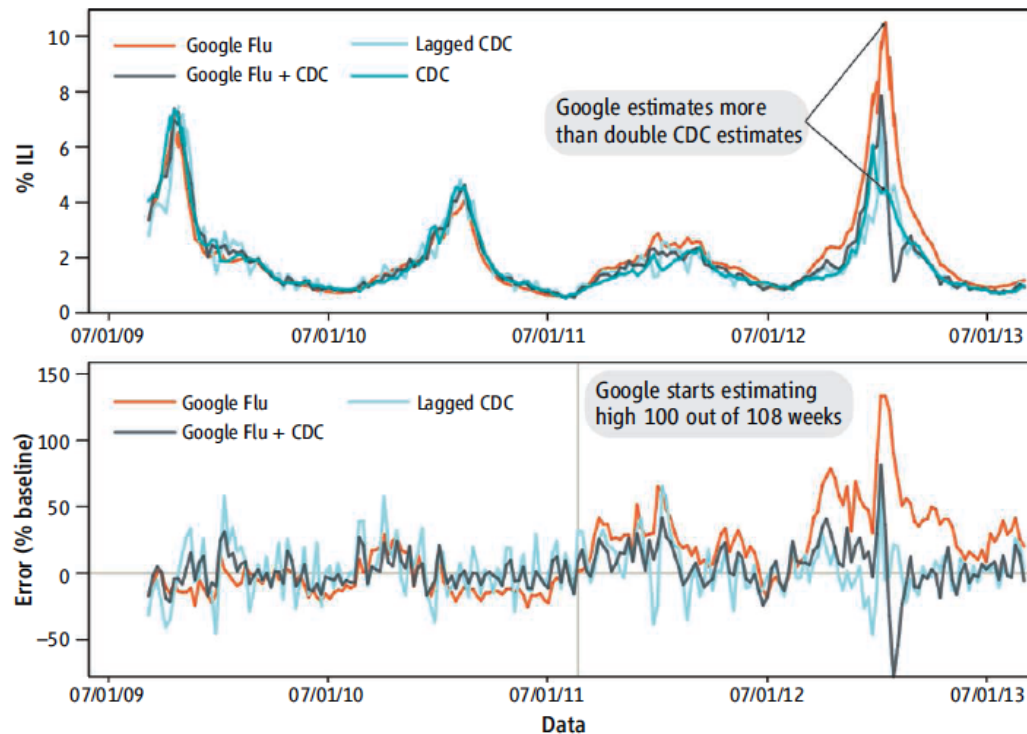
Caution: Validate Data and Results!

- **Bonferroni's Theorem:** if there are too many possible conclusions to draw, some will be true for purely statistical reasons, with no physical validity
- If possible, see that the entries make sense and data was collected properly
 - Ex: milk study at Lanarkshire, Scotland
- Data is often observational and not experimental
- Results validation vs. data dredging, snooping, fishing
 - E.g. S&P index almost perfectly predicted by butter, cheese production and sheep population in US and Bangladesh
 - “parapsychologist” David Rhine found (1950's) found about .1% guessed all 10 card colors correctly, but failed in next round.
 - Concluded that “telling people they have ESP causes them to lose it”!
 - www.tylervigen.com

The Google Flu-Trends Fiasco

The Parable of Google Flu: Traps in Big Data Analysis

<http://science.sciencemag.org/content/343/6176/1203>



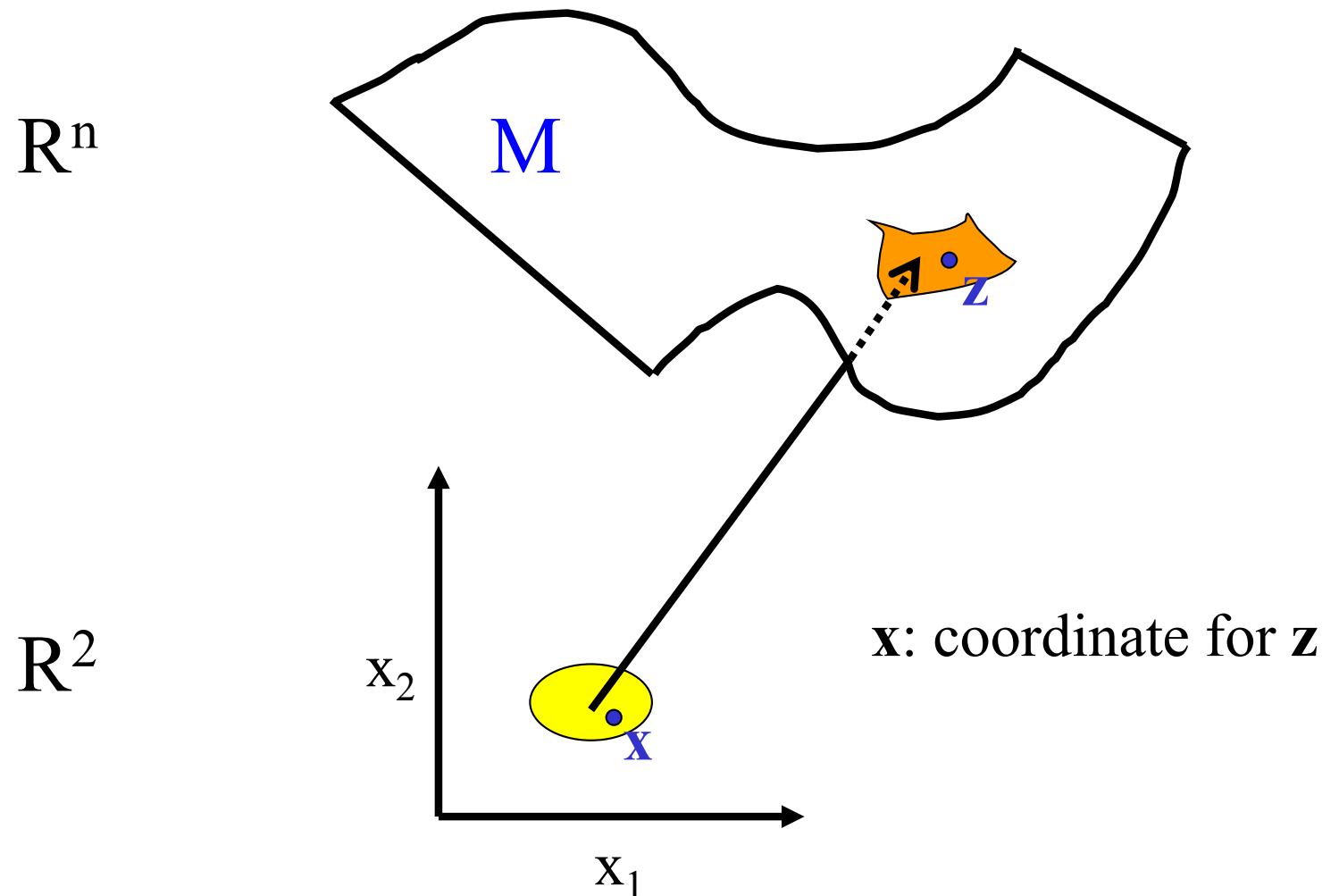
GFT overestimation. GFT overestimated the prevalence of flu in the 2012–2013 season and overshot the actual level in 2011–2012 by more than 50%. From 21 August 2011 to 1 September 2013, GFT reported overly high flu prevalence 100 out of 108 weeks. (**Top**) Estimates of doctor visits for ILI. "Lagged CDC" incorporates 52-week seasonality variables with lagged CDC data. "Google Flu + CDC" combines GFT, lagged CDC estimates, lagged error of GFT estimates, and 52-week seasonality variables. (**Bottom**) Error [as a percentage $[(\text{Non-CDC estimate}) - (\text{CDC estimate})] / (\text{CDC estimate})$]. Both alternative models have much less error than GFT alone. Mean absolute error (MAE) during the out-of-sample period is 0.486 for GFT, 0.311 for lagged CDC, and 0.232 for combined GFT and CDC. All of these differences are statistically significant at $P < 0.05$. See SM.

Extras

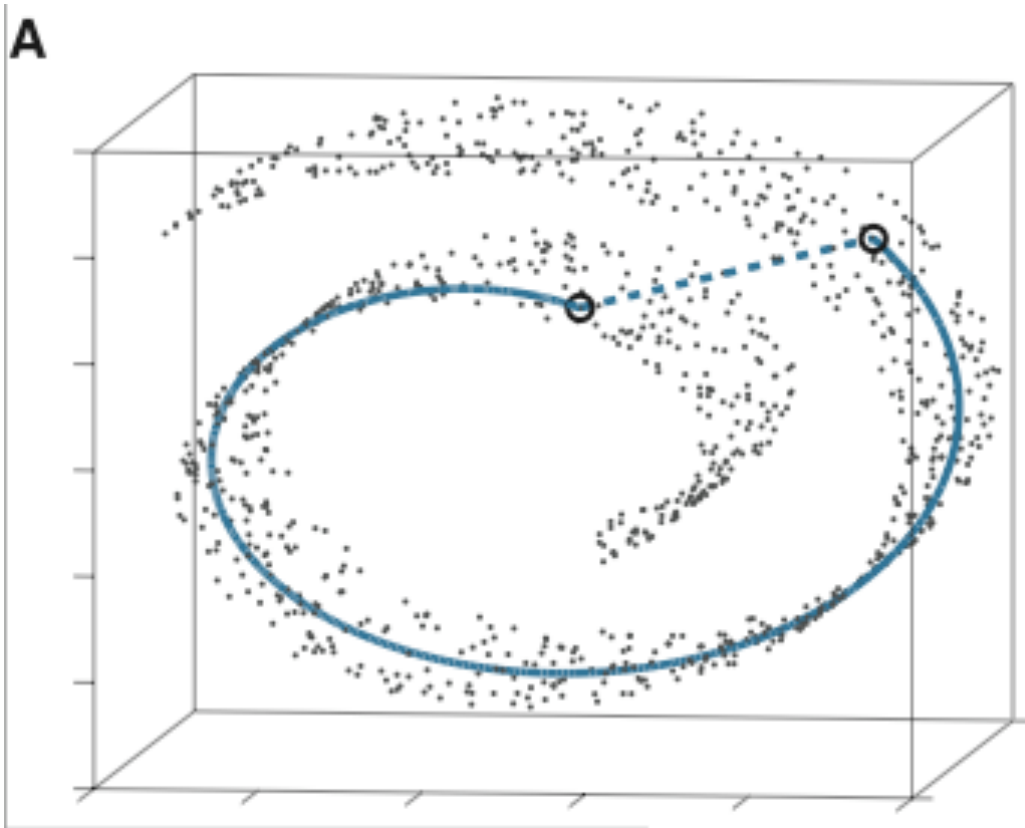
Manifold and Dimensionality Reduction

- Manifold: generalized “subspace” in \mathbb{R}^n ($n \gg 1$)
- Points in a *local* region on a manifold can be indexed by a subset of \mathbb{R}^k
 - The value of k is usually small
 - Thus map n -dim space into local k -dim coordinates.
 - Neural approaches include SOM and GTM

Example of a Manifold



Example: Manifold in Swiss Roll



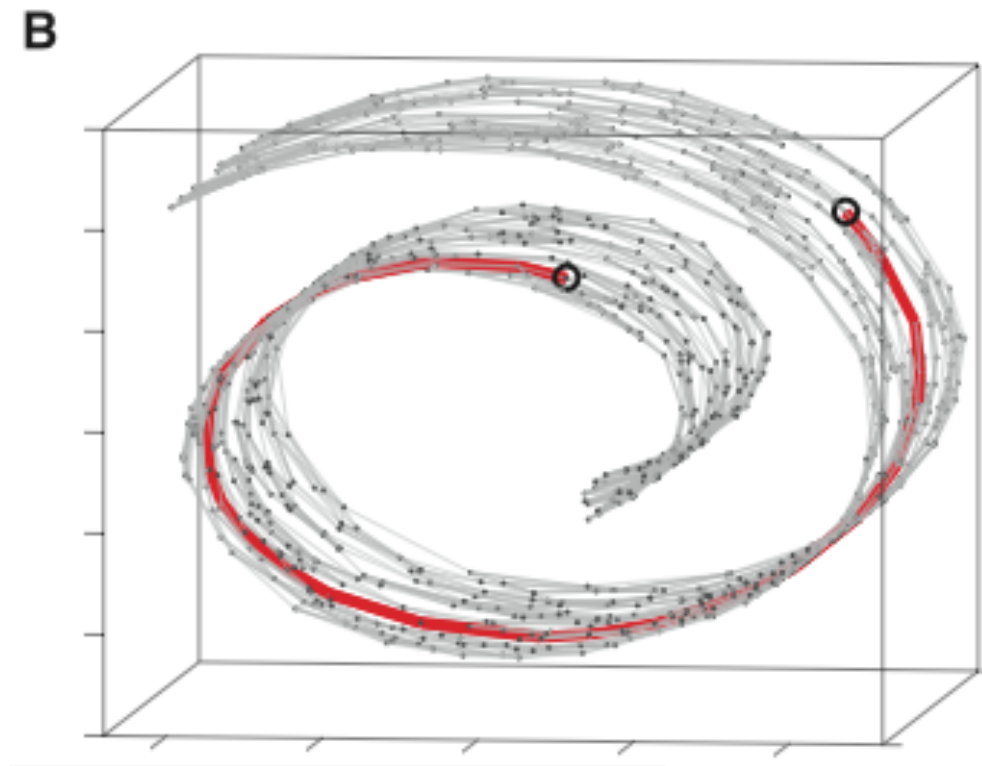
ISOMAP Algorithm

- Goal: preserve intrinsic geometry of data as captured via distances (**along the manifold**) between pairs of data points

Steps:

1. Determine which points are neighbors
2. Estimate geodesic distances and compute shortest path
 - For near points, measuring input space distances provides good enough approximation
 - For distant points, add up a sequence of short hops between neighboring points
3. Apply MDS to matrix of distances

-
- Estimating geodesic distances via shortest path



Step 3

- Apply classical MDS to matrix of distances

