

# Predictive Learning from Data

## LECTURE SET 1

### INTRODUCTION and OVERVIEW

Cherkassky, Vladimir, and Filip M. Mulier. *Learning from data: concepts, theory, and methods*. John Wiley & Sons, 2007.

Source: Dr. Vladimir Cherkassky (revised by Dr. Hsiang-Han Chen)

---

PLEASE DO NOT DISTRIBUTE WITHOUT AUTHOR'S PERMISSION.<sup>1</sup>

# **OUTLINE of Set 1**

1.1 Overview: what is this course about

1.2 Prerequisites and Expected outcomes

1.3 Big Data and Scientific Discovery

1.4 Related Data Modeling Methodologies

1.5 General Experimental Procedure for  
Estimating Models from Data

# 1.1 Overview

## *Uncertainty and Learning*

- Decision making under uncertainty
- Biological learning (adaptation)
- Plausible (uncertain) inference
- Induction in Statistics and Philosophy

Ex. 1: Many old men are bald

Ex. 2: Sun rises on the East every day

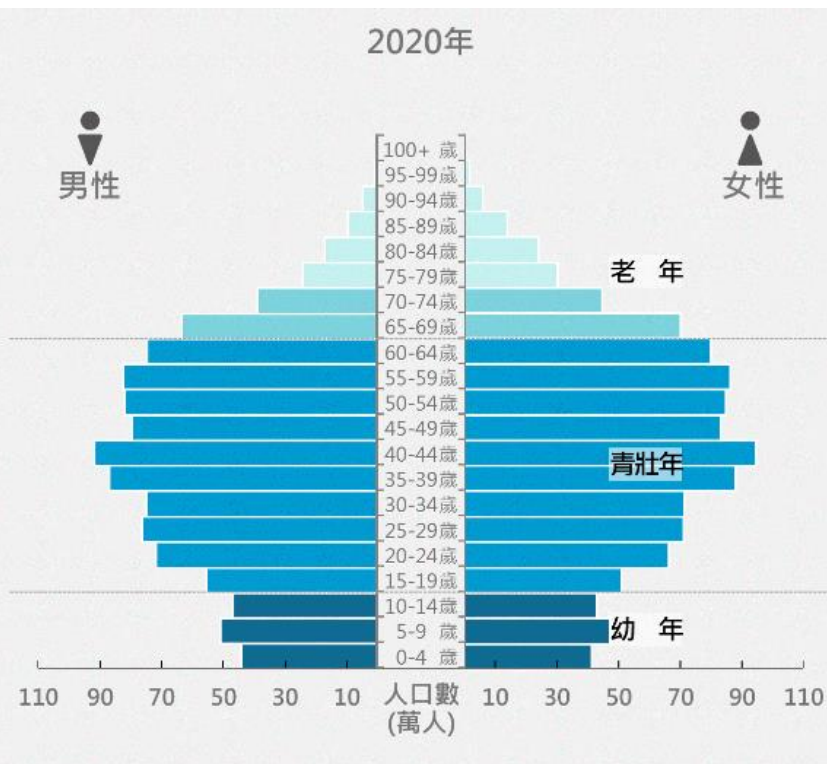
# (cont'd) Many old men are bald

- *Psychological Induction:*
  - inductive statement based on experience
  - also has certain predictive aspect
  - no scientific explanation
- *Statistical View:*
  - the lack of hair = **random variable**
  - estimate its **distribution** (depending on age) from past observations (**training** sample)
- *Philosophy of Science Approach:*
  - find scientific theory to explain the lack of hair
  - explanation itself is not sufficient
  - true theory needs to make **non-trivial predictions**

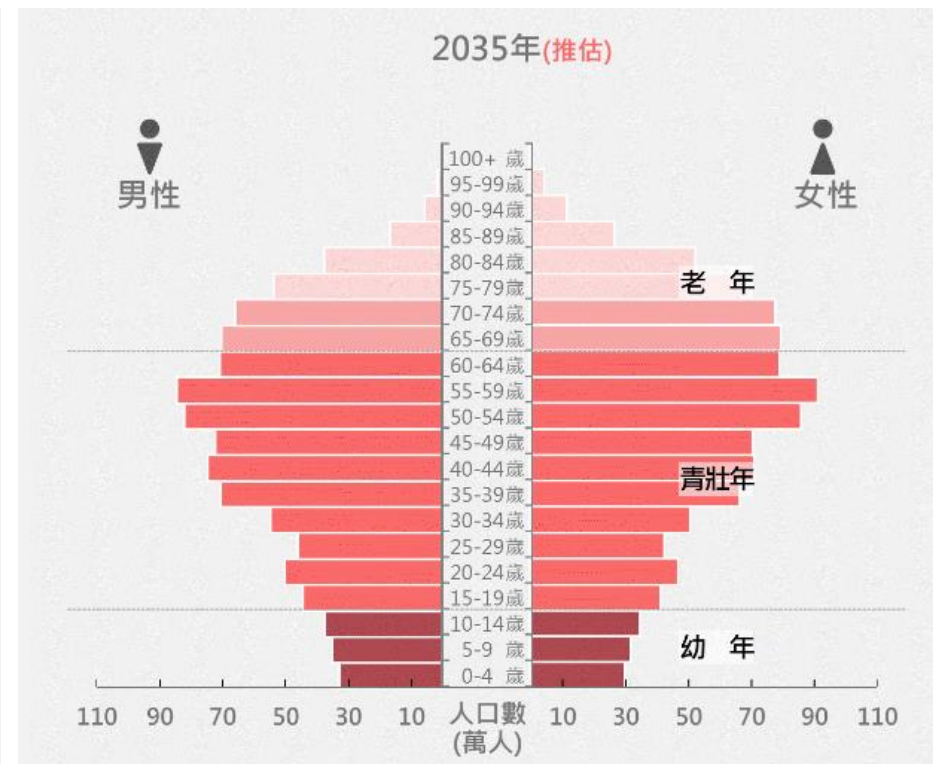
# Conceptual Issues

- Any theory (or model) has two aspects:
  - explanation of past data (observations)
  - prediction of future (unobserved) data

Past data



Future data



# Conceptual Issues

- A model to achieve both goals (explanation & prediction) perfectly is *not possible*
- Important issues to be addressed:
  - quality of explanation and prediction
  - is good prediction possible at all ?
  - if two models explain past data equally well, which one is better?
  - how to measure model complexity?

# Beliefs vs. Scientific Theories

*Men have lower life expectancy than women*

- *Because* they choose to do so
- *Because* they make more money (on average) and experience higher stress managing it
- *Because* they engage in risky activities
- *Because* .....

**Demarcation problem** in philosophy

- How to distinguish between science and non-science?

# Philosophical Connections

- *From Oxford English dictionary:*  
*Induction* is the process of inferring a general law or principle from the observations of particular instances.
- Clearly related to **Predictive Learning**.
- All science and (most of) human knowledge involves (some form of) induction
- How to form ‘good’ inductive theories?
  - inductive principles ~ general rules



# Philosophical Inductive Principles



**William of Ockham:** entities should not be multiplied beyond necessity

**Occam's razor** (also spelled Ockham's razor or Ocham's razor) is the problem-solving principle. It suggests that we should reduce assumptions to their minimum.

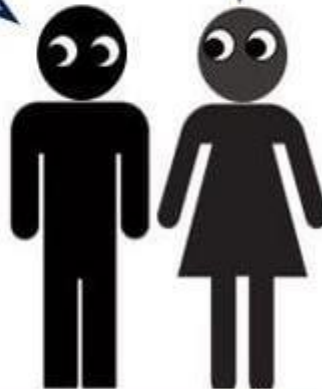
# Ockham's Razor

*Why did the tree fall down?*



*"I agree."*

*"It was the wind. It is the simpler explanation."*



## Two Explanations

1. The wind knocked down the tree.
2. Two meteorites. One hit the tree and knocked it down. Then it hit the other meteorite, thus obliterating evidence of its existence.

***When there are two explanations, choose the simpler one***

# Expected Outcomes

## Scientific / Technical:

- Learning = generalization, concepts and issues
- Math theory: Statistical Learning Theory aka VC-theory
- Conceptual basis for various learning algorithms

## Methodological:

- How to use available statistical/machine learning/ data mining s/w
- How to compare prediction accuracy of different learning algorithms
- Are you getting good modeling results because you are smart or just lucky?

## Practical Applications:

- Financial engineering
- Biomedical + Life Sciences
- Security
- Image recognition etc.

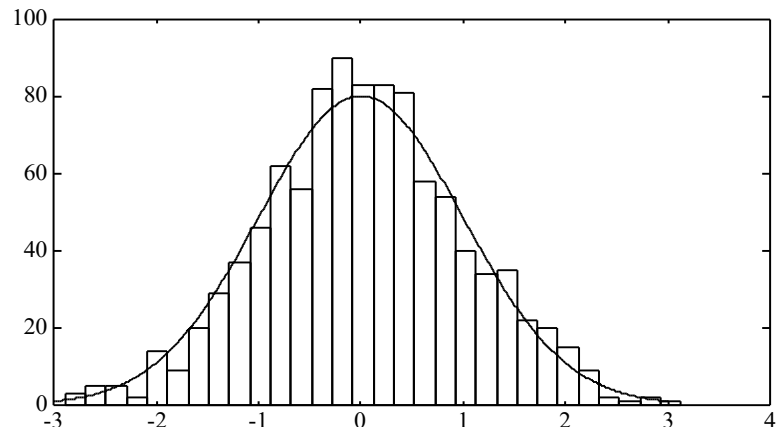
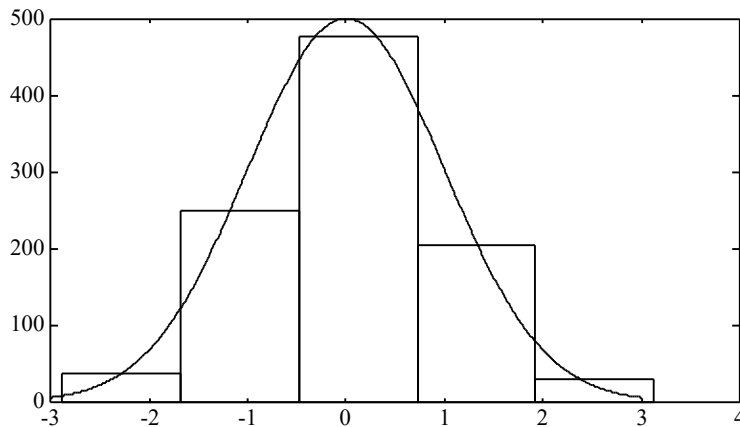
# Modeling financial data on Yahoo! Finance

- Real Data:  $X$ =daily price changes of SP500  
i.e.  $X(t) = \frac{Z(t) - Z(t-1)}{Z(t-1)} * 100\%$  where  $Z(t)$  = closing price
- Is the stock market *truly random*?
- Modeling assumption: price changes  $X$  are i.i.d.  
→ leads to certain analytic relationship that can be verified using empirical data.

# Understanding Daily Price Changes

**Histogram** = estimated pdf (from data)

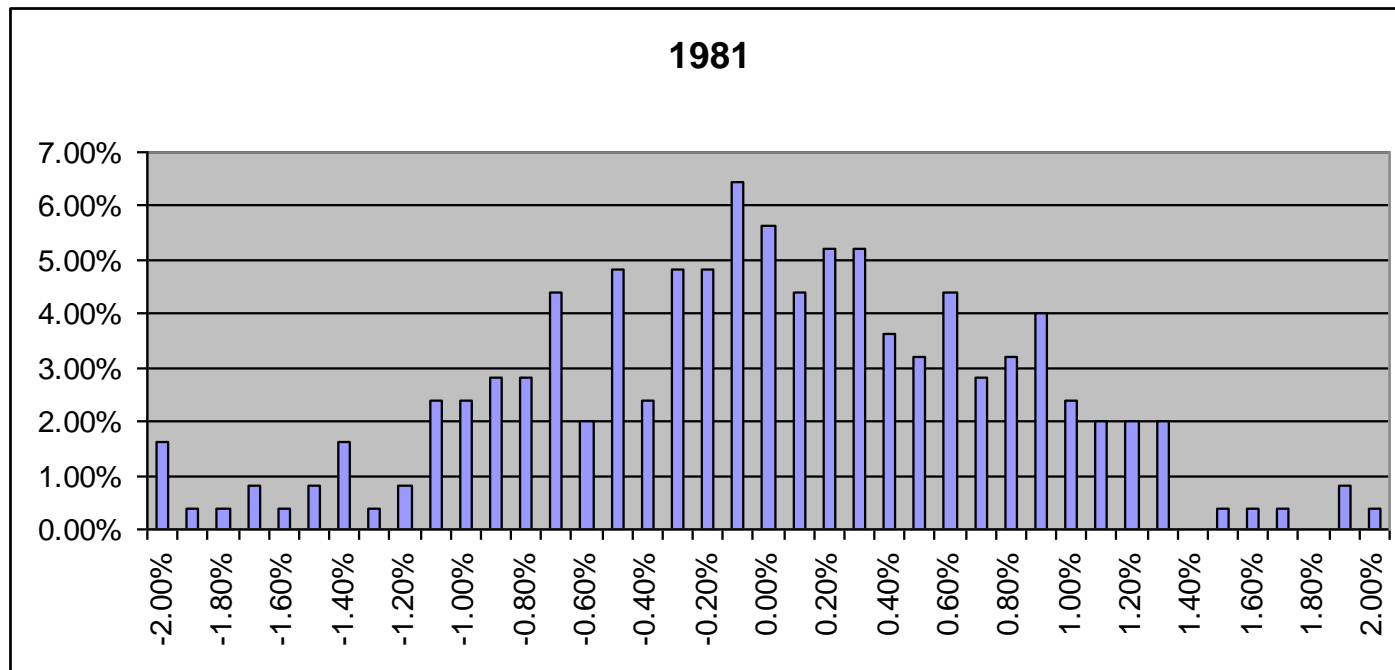
- **Example:** histograms of 5 and 30 bins to model  $N(0, 1)$  also **mean** and **standard deviation** (estimated from data)



# Histogram of daily price changes in 1981

**NOTE:** histogram ~ empirical pdf, i.e. scale of y-axis scale is in % (frequency).

Histogram of SP500 daily price changes in 1981:



# OUTLINE of Set 1

1.1 Overview: what is this course about

1.2 Prerequisites and Expected outcomes

## **1.3 Big Data and Scientific Discovery**

- scientific fairy tales
- promise of Big Data
- characteristics of scientific knowledge
- dealing with uncertainty and risk

1.4 Related Data Modeling Methodologies

1.5 General Experimental Procedure for Estimating Models from Data



# Historical Example

Ulisse Aldrovandi, 16<sup>th</sup> century

## Natural History of Snakes





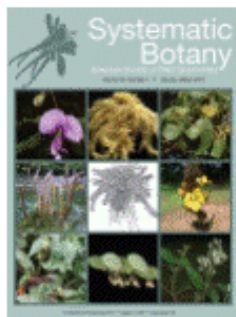
# Promise of Big Data

- **Technical fairy tales in 21<sup>st</sup> century**  
~ marketing + more marketing
- **Promise of Big Data:**  
s/w program + DATA → knowledge  
~ More Data → more knowledge
- **Yes-we-Can !**

# Examples from Life Sciences...

- Duke biologists discovered an **unusual link** btwn the popular singer and a new species of fern, i.e.
  - bisexual reproductive stage of the ferns;
  - the team found the sequence GAGA when analyzing the fern's DNA base pairs





# *Gaga*, a New Fern Genus Segregated from *Cheilanthes* (Pteridaceae)

Buy Article:  
\$20.00 + tax  
(Refund Policy)

ADD TO CART

BUY NOW

**Authors:** Li, Fay-Wei; Pryer, Kathleen M.; Windham, Michael D.

**Source:** Systematic Botany, Volume 37, Number 4, October-December 2012, pp. 845-860(16)

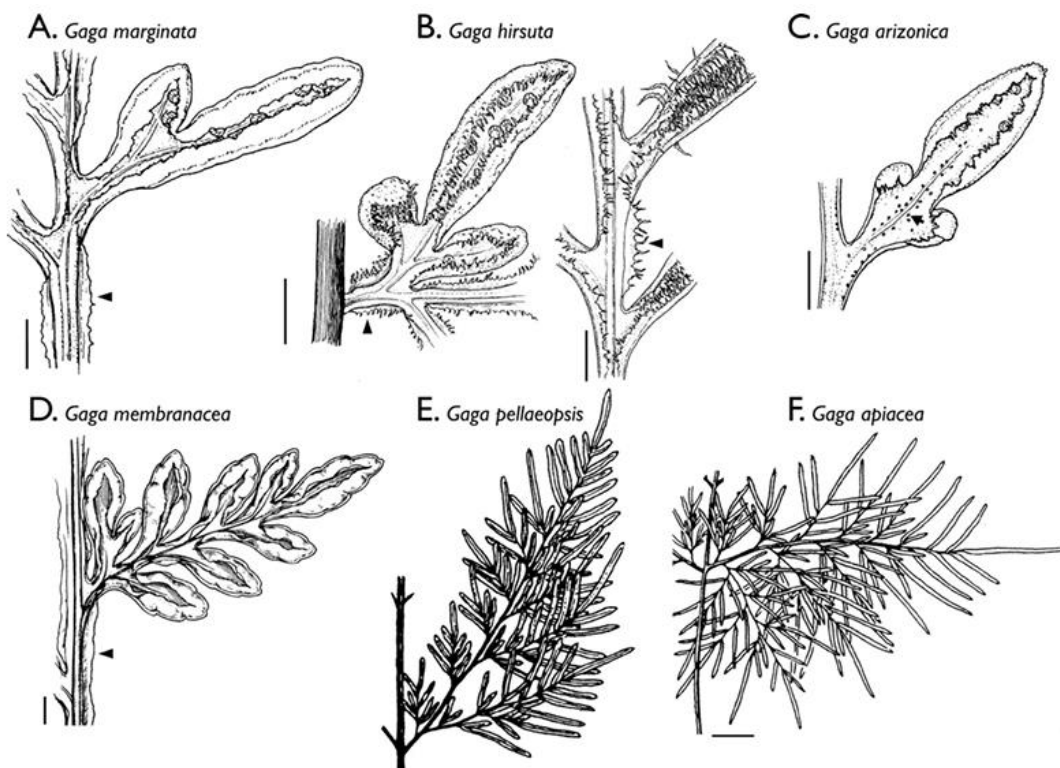
**Publisher:** American Society of Plant Taxonomists

**DOI:** <https://doi.org/10.1600/036364412X656626>



[< previous article](#) | [view table of contents](#) | [next article >](#)

[ADD TO FAVOURITES](#)



	590						600														610	
<i>Cheilanthes micropteris</i>	T	C	G	T	A	T	A	A	G	C	A	G	G	G	G	C	A	A	G	A	G	G
<i>Cheilanthes rufopunctata</i>	T	C	C	T	A	T	A	A	G	G	A	G	G	G	G	C	A	A	G	A	G	G
<i>Hemionitis arifolia</i>	T	C	C	T	A	T	A	A	A	A	A	G	G	G	G	C	A	A	G	A	G	A
<i>Pellaea angulosa</i>	T	C	C	T	A	T	A	A	G	G	A	G	G	G	G	C	A	A	G	A	A	A
<i>Pentagramma triangularis</i>	T	C	C	T	A	T	A	A	G	G	A	G	G	G	G	C	A	A	G	A	G	G
<i>Aspidotis californica</i>	T	C	C	T	A	T	A	A	G	G	A	G	G	G	G	C	A	A	G	A	G	G
<i>Aspidotis densa</i>	T	C	C	T	A	T	A	A	G	G	A	G	G	G	G	C	A	A	G	A	G	G
<i>Aspidotis meifolia</i>	T	C	C	T	A	T	A	A	A	A	A	G	G	G	G	C	A	A	G	A	G	G
<i>Gaga marginata</i>	T	C	C	T	A	T	A	A	G	G	A	G	A	G	G	C	A	A	G	A	G	G
<i>Gaga hirsuta</i>	T	C	C	T	A	T	A	A	G	G	A	G	A	G	G	C	A	A	G	A	G	G
<i>Gaga arizonica</i>	T	C	C	T	A	T	A	A	G	G	A	G	A	G	G	C	A	A	G	A	G	G
<i>Gaga lerstenii</i>	T	C	C	T	A	T	A	A	G	G	A	G	A	G	G	C	A	A	G	A	G	G
<i>Gaga membranacea</i>	T	C	C	T	A	T	A	A	G	G	A	G	A	G	G	C	A	A	G	A	G	G
<i>Gaga germanotta</i>	T	C	C	T	A	T	A	A	G	G	A	G	A	G	G	C	A	A	G	A	G	G
<i>Gaga complanata</i>	T	C	C	T	A	T	A	A	G	G	A	G	A	G	G	C	A	A	G	A	G	G
<i>Gaga harrisii</i>	T	C	C	T	A	T	A	A	G	G	A	G	A	G	G	C	A	A	G	A	G	G
<i>Gaga decurrens</i>	T	C	C	T	A	T	A	A	G	G	A	G	A	G	G	C	A	A	G	A	G	G
<i>Gaga purpusii</i>	T	C	C	T	A	T	A	A	G	G	A	G	A	G	G	C	A	A	G	A	G	G
<i>Gaga apiacea</i>	T	C	C	T	A	T	A	A	G	G	A	G	A	G	G	C	A	A	G	A	G	G
<i>Gaga pellaopsis</i>	T	C	C	T	A	T	A	A	G	G	A	G	A	G	G	C	A	A	G	A	G	G
<i>Gaga hintoniorum</i>	T	C	C	T	A	T	A	A	G	G	A	G	A	G	G	C	A	A	G	A	G	G
<i>Gaga angustifolia</i>	T	C	C	T	A	T	A	A	G	G	A	G	A	G	G	C	A	A	G	A	G	G
<i>Gaga chaerophylla</i>	T	C	C	T	A	T	A	A	G	G	A	G	A	G	G	C	A	A	G	A	G	G
<i>Gaga decomposita</i>	T	C	C	T	A	T	A	A	G	G	A	G	A	G	G	C	A	A	G	A	G	G
<i>Gaga cuneata</i>	T	C	C	T	A	T	A	A	G	G	A	G	A	G	G	C	A	A	G	A	G	G
<i>Gaga monstraparva</i>	T	C	C	T	A	T	A	A	G	G	A	G	A	G	G	C	A	A	G	A	G	G
<i>Gaga kaulfussii</i>	T	C	C	T	A	T	A	A	G	G	A	G	A	G	G	C	A	A	G	A	G	G

GAGA

synapomorphy

GAGA  
synapomorphy

# Scientific Discovery

- Combines **ideas/models** and **facts/data**
- **First-principle knowledge:**  
hypothesis → experiment → theory  
~ deterministic, causal, intelligible models
- **Modern data-driven discovery:**  
s/w program + DATA → knowledge  
~ statistical, complex systems
- **Many methodological differences**

# Invariants of Scientific Knowledge

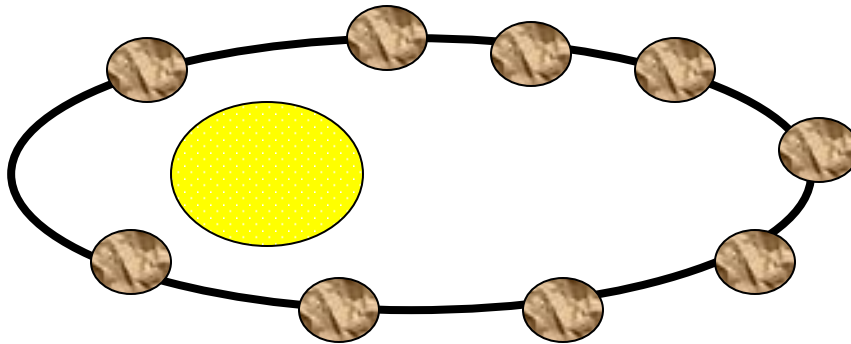
- **Intelligent questions**
- **Non-trivial predictions**
- **Clear limitations/ constraints**
- **All require human intelligence**
  - missing/ lost in Big Data?

# Historical Example: Planetary Motions

- How planets move among the stars?
  - Ptolemaic system (geocentric)
  - Copernican system (heliocentric)
- Tycho Brahe (16 century)
  - measure positions of the planets in the sky
  - use experimental data to support one's view
- Johannes Kepler:
  - used volumes of Tycho's data to discover three remarkably simple laws

# First Kepler's Law

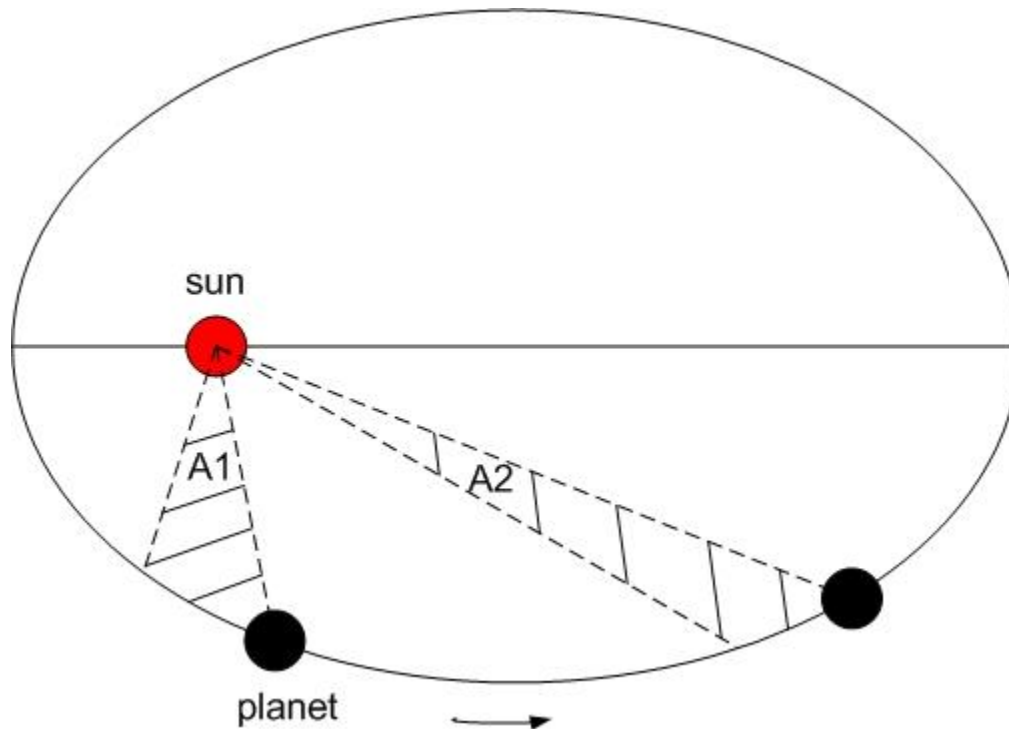
- Sun lies in the plane of orbit, so we can represent positions as (x,y) pairs
- An orbit is an ellipse, with the sun at a focus



$$c_1x^2 + c_2y^2 + c_3xy + c_4x + c_5y + c_6 = 0$$

# Second Kepler's Law

- The radius vector from the sun to the planet sweeps out equal areas in the same time intervals





# Third Kepler's Law

	P	D	P <sup>2</sup>	D <sup>3</sup>
Mercury	0.24	0.39	0.058	0.059
Venus	0.62	0.72	0.38	0.39
Earth	1.00	1.00	1.00	1.00
Mars	1.88	1.53	3.53	3.58
Jupiter	11.90	5.31	142.0	141.00
Saturn	29.30	9.55	870.0	871.00

P = orbit period    D = orbit size (half-diameter)

**For any planet:  $P^2 \sim D^3$**

# Empirical Scientific Theory

- Kepler's Laws can
  - explain experimental data
  - predict new data (i.e., other planets)
  - *BUT* do not explain *why planets move*.
- Popular explanation
  - planets move because there are invisible angels beating the wings behind them
- **First-principle scientific explanation**

Galileo and Newton discovered laws of motion and gravity that explain Kepler's laws.

# OUTLINE of Set 1

1.1 Overview: what is this course about

1.2 Prerequisites and Expected outcomes

1.3 Big Data and Scientific Discovery

## **1.4 Related Data Modeling Methodologies**

- growth of empirical knowledge
- empirical vs first-principle knowledge
- handling uncertainty and risk
- related data modeling methodologies

1.5 General Experimental Procedure.

# Scientific knowledge

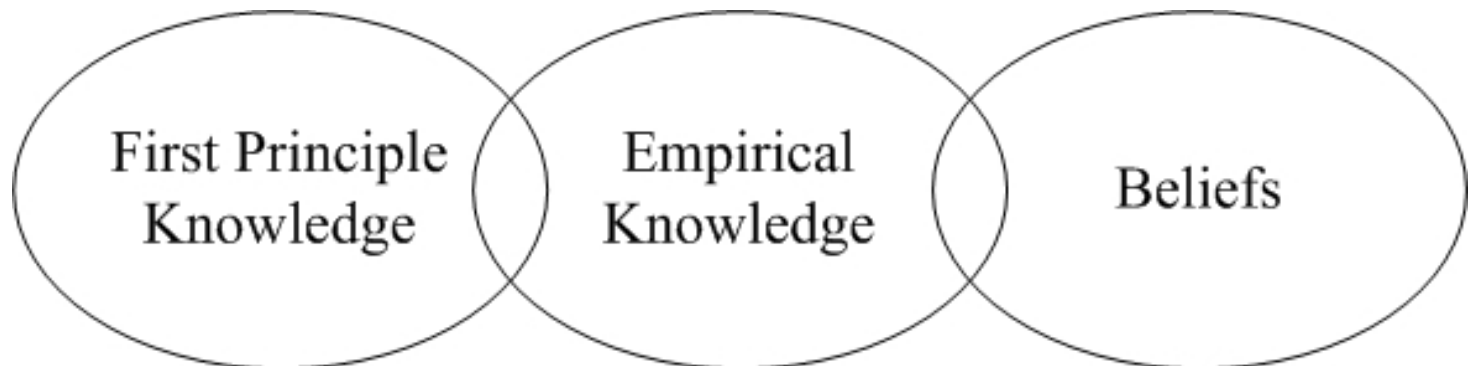
- **Knowledge**
  - ~ stable relationships between facts and ideas (mental constructs)
- **Classical first-principle knowledge:**
  - rich in ideas
  - relatively **few facts** (amount of data)
  - **simple relationships**

# Growth of empirical knowledge

- Huge growth of the amount of data in 20<sup>th</sup> century (computers and sensors)
- Complex systems (engineering, life sciences and social)
- Classical first-principles science is inadequate for empirical knowledge
- Need for new Methodology:  
How to estimate good predictive models from noisy data?

# Different types of knowledge

- Three types of knowledge
  - scientific (first-principles, deterministic)
  - empirical (uncertain, statistical)
  - metaphysical (beliefs)



- Boundaries are poorly understood

# Handling Uncertainty and Risk(1)

- Ancient times
- Probability for quantifying **uncertainty**
  - degree-of-belief
  - frequentist (Cardano-1525, Pascale, Fermat)
- Newton and **causal determinism**
- Probability theory and statistics (20<sup>th</sup> century)
- **Modern classical science** (A. Einstein)
  - Goal of science: estimating a **true model** or **system identification**

# Handling Uncertainty and Risk(2)

- Making decisions under uncertainty  
~ *risk management, adaptation, intelligence...*
- **Probabilistic approach:**
  - estimate probabilities (of future events)
  - assign costs and minimize expected risk
- **Risk minimization** approach:
  - apply decisions to known past events
  - select one minimizing expected risk



# Summary

- **First-principles** knowledge:  
deterministic relationships between a few concepts (variables)
- *Importance of empirical knowledge:*
  - statistical in nature
  - (usually) many input variables
- **Goal of modeling:** to act/perform well, rather than system identification

# Other Related Methodologies

- **Estimation of empirical dependencies** is commonly addressed many fields
  - *statistics, data mining, machine learning, neural networks, signal processing* etc.
  - each field has its own methodological bias and terminology → confusion
- Quotations from popular textbooks:

The field of *Pattern Recognition* is concerned with the automatic discovery of regularities in data.

*Data Mining* is the process of automatically discovering useful information in large data repositories.

*Statistical Learning* is about learning from data.
- All these fields are concerned with estimating predictive models from data.

# Other Methodologies (cont'd)

- **Generic Problem**

Estimate (learn) useful models from available data

- **Methodologies differ** in terms of:

- what is **useful**
- (assumptions about) **available data**
- **goals of learning**

- Often these important notions are not well-defined.

# Common Goals of Modeling

- **Prediction (Generalization)**
- **Interpretation ~ descriptive model**
- **Human decision-making** using both above
- **Information retrieval**, i.e. predictive or descriptive modeling of unspecified subset of available data

## *Note:*

- These goals usually ill-defined
- Formalization of these goals in the context of application requirements is THE MOST IMPORTANT aspect of 'data mining'

# Three Distinct Methodologies (section 1.5)

- **Statistical Estimation**

- from classical statistics and fct approximation

- **Predictive Learning (~ machine learning)**

- practitioners in machine learning /neural networks
- Vapnik-Chervonenkis (VC) theory for estimating predictive models from empirical(finite) data samples

- **Data Mining**

- exploratory data analysis, i.e. selecting a subset of available (large) data set with **interesting properties**

# OUTLINE of Set 1

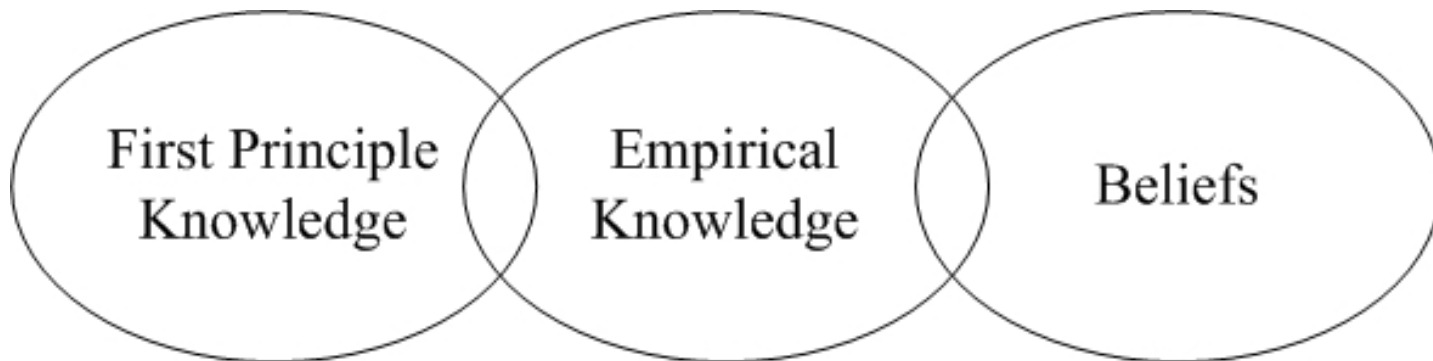
1.1 Overview: what is this course about

1.2 Prerequisites and Expected outcomes

1.3 Big Data and Scientific Discovery

1.4 Related Data Modeling Methodologies

## **1.5 General Experimental Procedure for Estimating Models from Data**



# 1.5 General Experimental Procedure

1. Statement of the Problem
2. Hypothesis Formulation (Problem Formalization) –  
*different from classical statistics*
3. Data Generation/ Experiment Design
4. Data Collection and Preprocessing
5. Model Estimation (learning)
6. Model Interpretation, Model Assessment and Drawing Conclusions

## Note:

- each step is complex and requires several iterations
- estimated model depends on all previous steps
- **observational data** (*not experimental\_design*)

# Data Preprocessing and Scaling

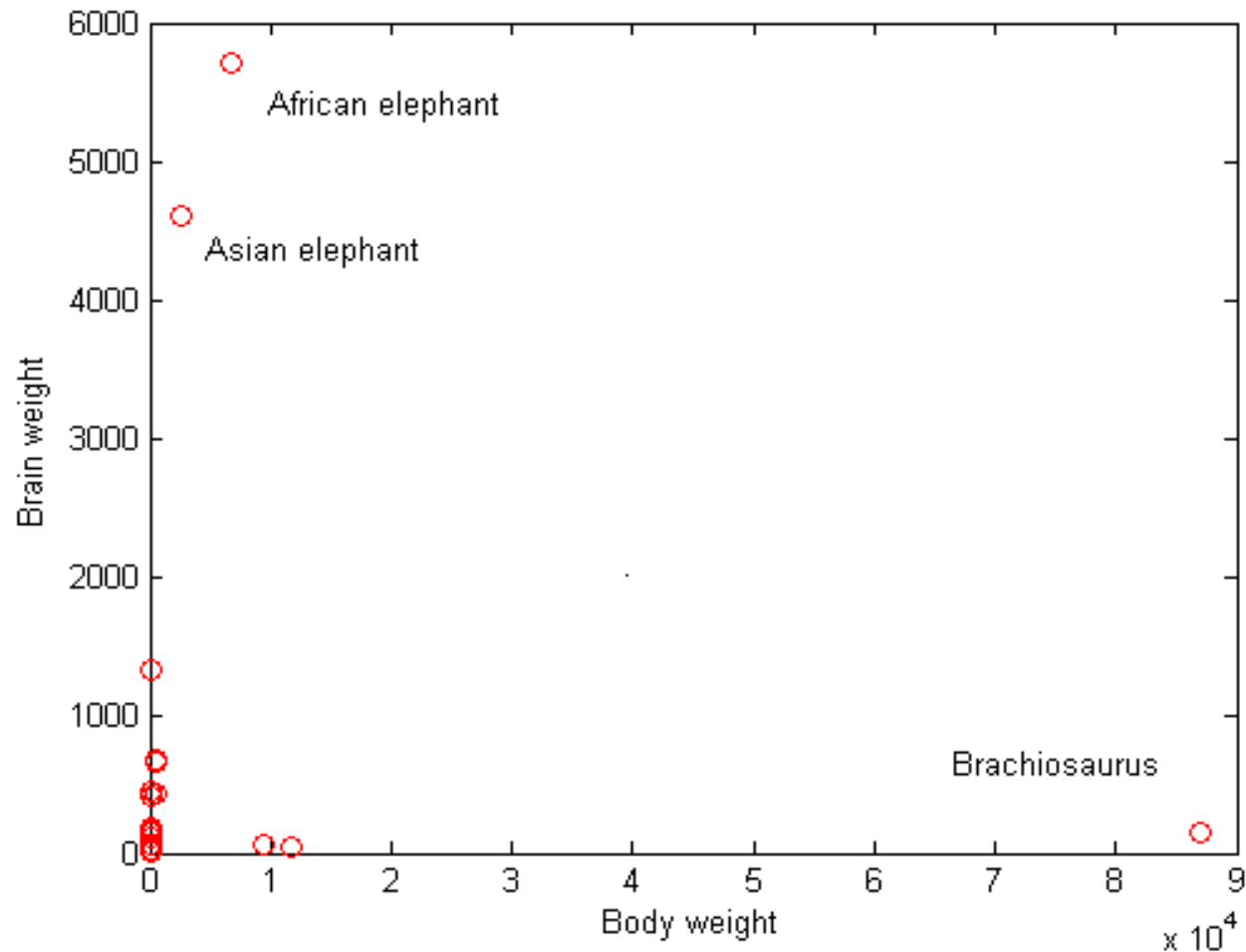
- Preprocessing is required with observational data (**step 4** in general experimental procedure)

*Examples: ...*

- Basic preprocessing includes
  - summary univariate statistics: *mean, st. deviation, min + max value, range, boxplot* performed independently for each input/output
  - *detection (removal) of outliers*
  - *scaling* of input/output variables (may be **necessary** for some learning algorithms)
- Visual inspection of data is tedious but useful



# Original Unscaled Animal Data



# Cultural + Ethical Aspects

- **Cultural and business aspects** usually affect:
  - problem formalization
  - data access/ sharing (i.e., in life sciences)
  - model interpretation
- **Possible (idealistic) solution approach**
  - to adopt common methodology
  - critical for interdisciplinary projects

# Honest Disclosure of Results

- **Modern drug studies**

Review of studies submitted to FDA

- Of 74 studies reviewed, 38 were judged to be positive by the FDA.

*All but one were published.*

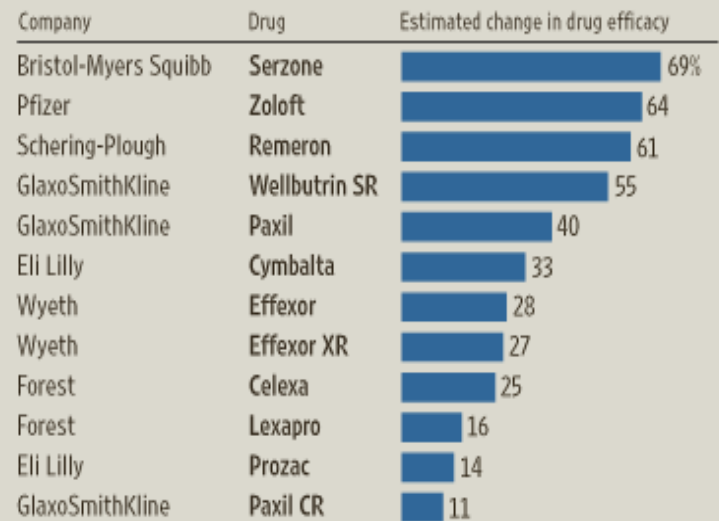
- Most of the studies found to have negative or questionable results *were not published.*

**Source:** The New England Journal of Medicine, WSJ Jan 17, 2008)

**Publication bias:** common in modern research

## Under Wraps

Estimate of how much the impression of each drug's effectiveness was inflated by not publishing unfavorable studies



Source: New England Journal of Medicine