**Data Mining for Business Analytics**

**Lecture 1: Introduction to Data Mining**

**Stern School of Business**
**New York University**
**Spring 2019**

G. Valkanas — New York University

---

## The Good



G. Valkanas — New York University

---

## The Good



G. Valkanas — New York University

---

## The Good



G. Valkanas — New York University

---

## The Good

- Data Mining is pervasive

G. Valkanas — New York University

---

## The "Bad"

- No Free Lunch

G. Valkanas — New York University

## The "Bad"

- Effective Data Science requires / builds on a *SET* of skills:

  - Analytical thinking

  - Technical skills

  - Creativity

  - Communication

  - Domain Knowledge **(!)**

## The "Ugly"

- We will be doing some math

## The "Ugly"

- We will be doing some math

- We will be doing some Programming

  - Highly sought-after skill !

  - **BUT REMEMBER:**

    - Data Mining is *not* (just) about coding, especially in business settings!

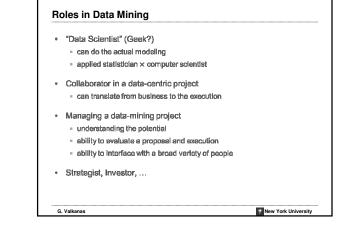    - We **will also** be focusing on several non-technical areas

## Let's play a game…

## Data Mining Approach

**"If we have data, let's look at data. If all we have are opinions, let's go with mine."**

-- *James Love Barksdale*

*Former CEO of Netscape*

## Data Mining

- A set of principles, concepts, and techniques that **structure thinking** and **analysis** of data

- Extracts **useful information** and **knowledge** from large volumes of data by **following a process** with reasonably well defined steps

- Changes the way you think about data and its role in business
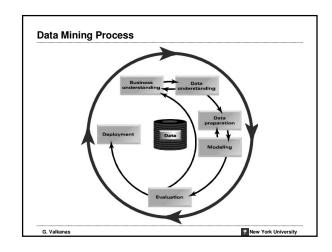
## Learning Goals

- Approach business problems data-analytically

- Interact competently on the topic of data mining for business intelligence
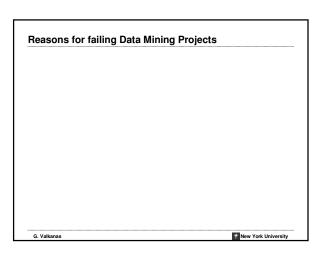
- Hands-on experience mining data

---

## Roles in Data Mining

- "Data Scientist" (Geek?)
  - can do the actual modeling
  - applied statistician × computer scientist

- Collaborator in a data-centric project
  - can translate from business to the execution

- Managing a data-mining project
  - understanding the potential
  - ability to evaluate a proposal and execution
  - ability to interface with a broad variety of people

- Strategist, Investor, …

---

## Business data mining is a process

science + craft + creativity + common sense

a *process*

---

## Data Mining Process

---

## Outline

- Business Understanding

- Data Understanding

- Data Preparation

- Modeling

- Evaluation

- Deployment

---

## Reasons for failing Data Mining Projects

**This is NOT a course about…**

- Statistics

- Database Querying
  - SQL

- Data Warehousing

- Regression Analysis
  - Explanatory vs Predictive modeling

- Big Data

---

**Data Mining versus…**

- Data Warehousing / Storage
  - Data warehouses coalesce data from across an enterprise, often from multiple transaction-processing systems

- Querying / Reporting (SQL, Excel, QBE, other GUI-based querying)
  - Very flexible interface to ask factual questions about data
  - No modeling or sophisticated pattern finding
  - Most of the cool visualizations

- OLAP – On-line Analytical Processing
  - OLAP provides easy-to-use GUI to explore large data collections
  - Exploration is <u>manual</u>; no modeling
  - Dimensions of analysis preprogrammed into OLAP system

---

**Data Mining versus…**

- Traditional statistical analysis
  - Mainly based on hypothesis testing or estimation / quantification of uncertainty
  - Should be used to follow-up on data mining's <u>hypothesis generation</u>

- **Automated statistical modeling** (e.g., advanced regression)
  - This <u>is</u> data mining, one type – usually based on linear models
  - Massive databases allow non-linear alternatives

---

**Answering business questions with these techniques..**

- Who are the most profitable customers?
  - **Database querying**

- Is there really a difference between profitable customers and the average customer?
  - **Statistical hypothesis testing**

- But who really are these customers? Can I characterize them?
  - **OLAP** (manual search), **Data mining** (automated pattern finding)

- Will some particular new customer be profitable? How much revenue should I expect this customer to generate?
  - **Data mining** (predictive modeling)

---

# Thanks!

---

# Questions?