# Essential Data Science for Business:
# Top 10 Analytics Topics

## What are the key topics that are used in the business?

**NISS Webinar**

**Victor S.Y. Lo**
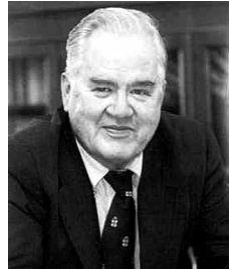
**July 29, 2020**

# Content

► **Three Major Types of Analytics**

► **What Skills Data Scientists Should Have**

► **Top 10 Analytics Topics – Important and Practical**

Disclaimer: The views expressed here are solely those of the speaker and do not in any way represent the views of Fidelity Investments

"**The best thing about being a statistician is that you get to *play in everyone's backyard.*"**

**- John Tukey, decades ago**

"**We no longer simply enjoy the privilege of playing in or cleaning up everyone's backyard. We are <u>now being invited into</u> *everyone's study or living room,* and trusted with the task of being their offspring's first quantitative nanny.**"

**- Xiao-li Meng (2009), Harvard University**

# Three Types of Analytics



**Prescriptive**

**What should we do?**
**What is the Best Decision?**
- Support *decision making* and *proactive* actions

**Predictive**

**What will happen?**
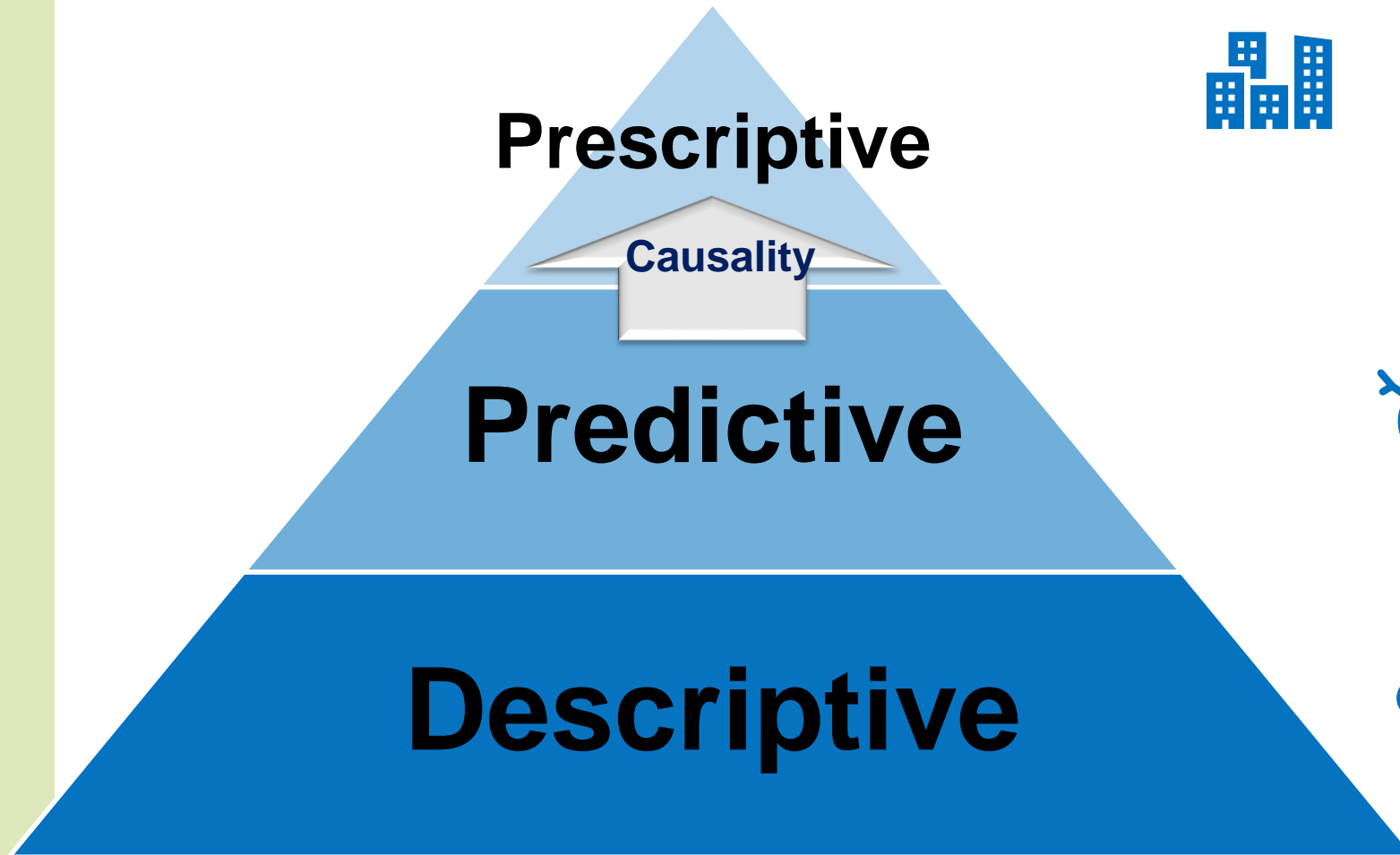- Predict *future* forward-looking behavior, events, probabilities, or trends

**Descriptive**

**What happened?**
- Reports and profiling
- Data visualization

Source: http://www.sas.com/news/sascom/2008q4/column_8levels.html, and https://www.informs.org/Community/Analytics

# Data Science Venn Diagram



**Computer Science**

**Statistics & Math**

**Subject Matter Expertise**

e.g. Marketing, Finance, Insurance, Healthcare, Risk
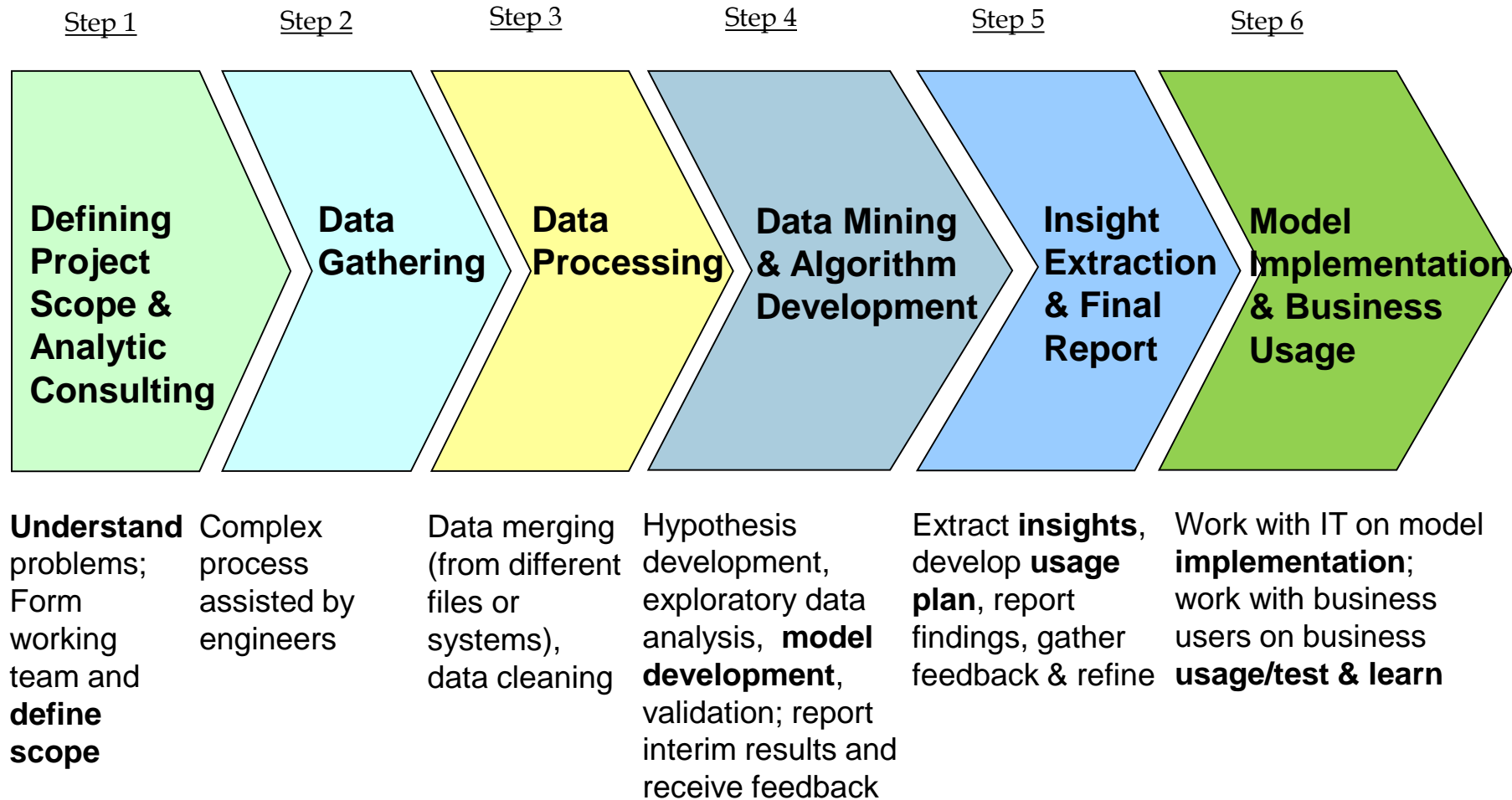
**Soft Skills**

e.g. Business Consulting, Communication, Writing

**Data Science is a Diversified field with professionals from a variety of disciplines, see Lo (2019)**

# Data Science Project Process

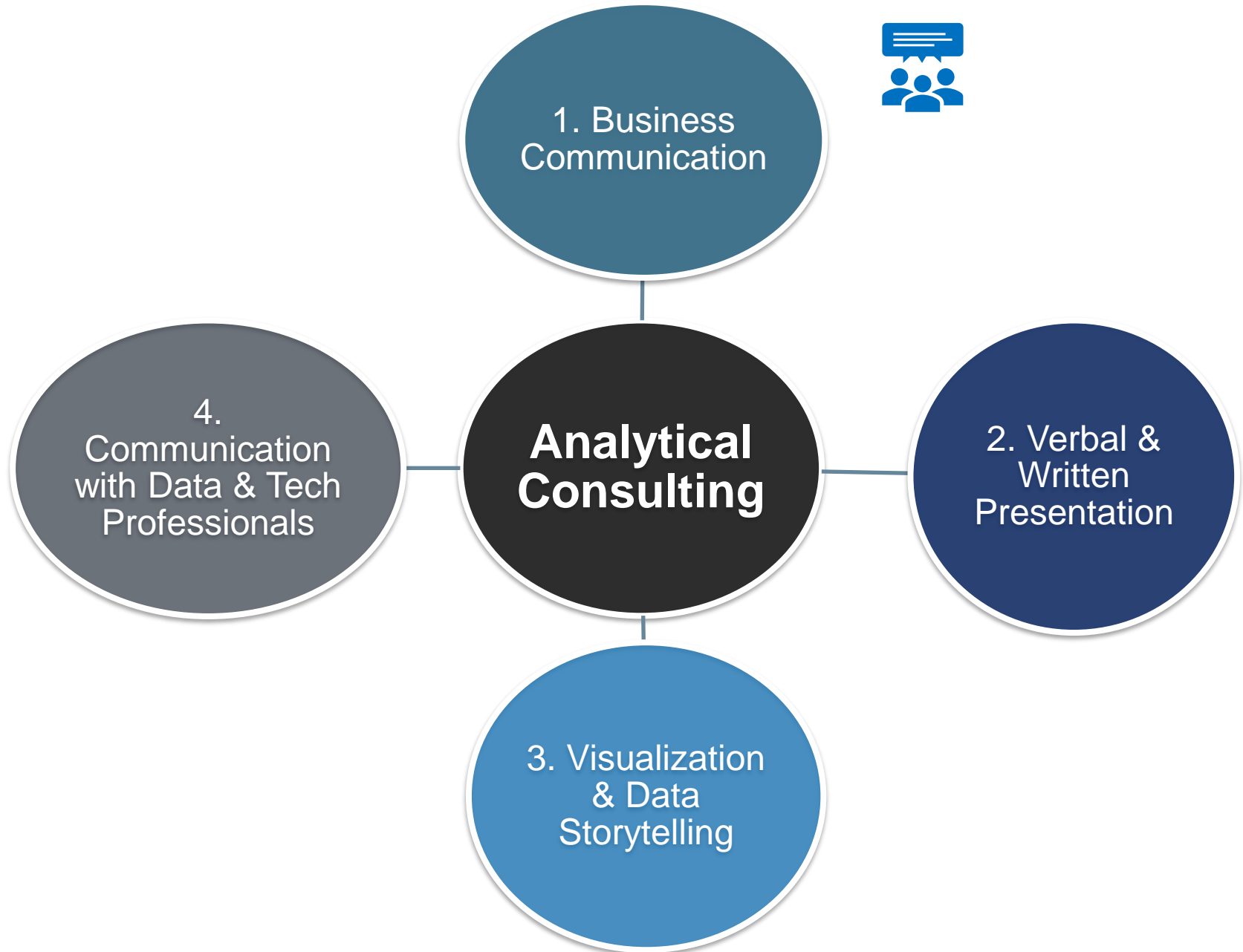| Step 1 | Step 2 | Step 3 | Step 4 | Step 5 | Step 6 |
|---|---|---|---|---|---|
| **Defining Project Scope & Analytic Consulting** | **Data Gathering** | **Data Processing** | **Data Mining & Algorithm Development** | **Insight Extraction & Final Report** | **Model Implementation & Business Usage** |
| **Understand** problems; Form working team and **define scope** | Complex process assisted by engineers | Data merging (from different files or systems), data cleaning | Hypothesis development, exploratory data analysis, **model development**, validation; report interim results and receive feedback | Extract **insights**, develop **usage plan**, report findings, gather feedback & refine | Work with IT on model **implementation**; work with business users on business **usage/test & learn** |

**Top 10 Important and Practical Topics that <u>May NOT Be Covered in Your Education Program</u>…**

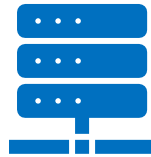# ▶ 1. Analytical Consulting, Communication and Soft Skills

**Analytical Consulting**

1. Business Communication

2. Verbal & Written Presentation

3. Visualization & Data Storytelling

4. Communication with Data & Tech Professionals

9

# Analytic Consulting Process

**1) Initiate Project**

- **Project opportunity identified** by business owner or data scientist

**2) Define Scope**

+ **Ask** right questions and understand **business problem**

+ **Agree** on objective & scope

**3) Propose Approach**

+ **Recommend** initial **approach** and draft **proposal**

+ Establish a feasibility study, if appropriate

**4) Identify Resources and Dependencies**

+ **Determine** data scientist, data engineer, and other **resources** required

+ **Identify** other **dependencies,** e.g., business & data expertise

**5) Leadership during Project Development**

+ **AGILE**

+ Determine **engagement process**

+ Report **Interim analyses** to seek feedback and buy-ins

+ **Present** to multiple levels of the business

**6) Leadership during Deployment**

+ **Deployment** discussion with business and technology stakeholders

+ **Performance** tracking and refinement

# ▶ 2. Computer Science, Programming, and Tools

# ▶ 2. Computer Science, Programming, and Tools

**Demand** for computational power has **dramatically increased** and will need to expand much further:

o  Growth in Structured and **Unstructured** Data

o  Internet of Things (**IoT**)

o  Practical Success of **Deep Learning**

**IDC predicted that the global data size would increase by ~3X from 2019 to 2025** (175 zettabytes), see Reinsel et al (2020)

**1. AI/ML and Statistical Programming**

Python, R, SAS, GPU Programming

**2. AI/ML Tools**

PyTorch, Keras, Tensorflow, MXNET, Caffe, CNTK, SageMaker

**3. Data Knowledge**

**4. ETL**

(Extract, Transform, and Load) skills for Big Data

**5. Model Deployment**

Kubernetes, Docker
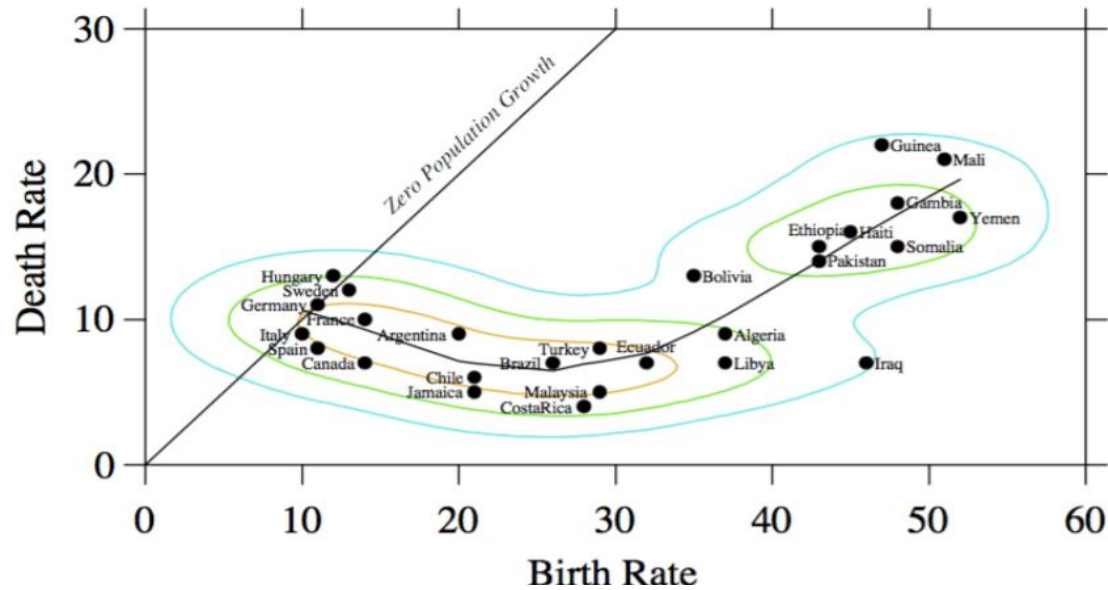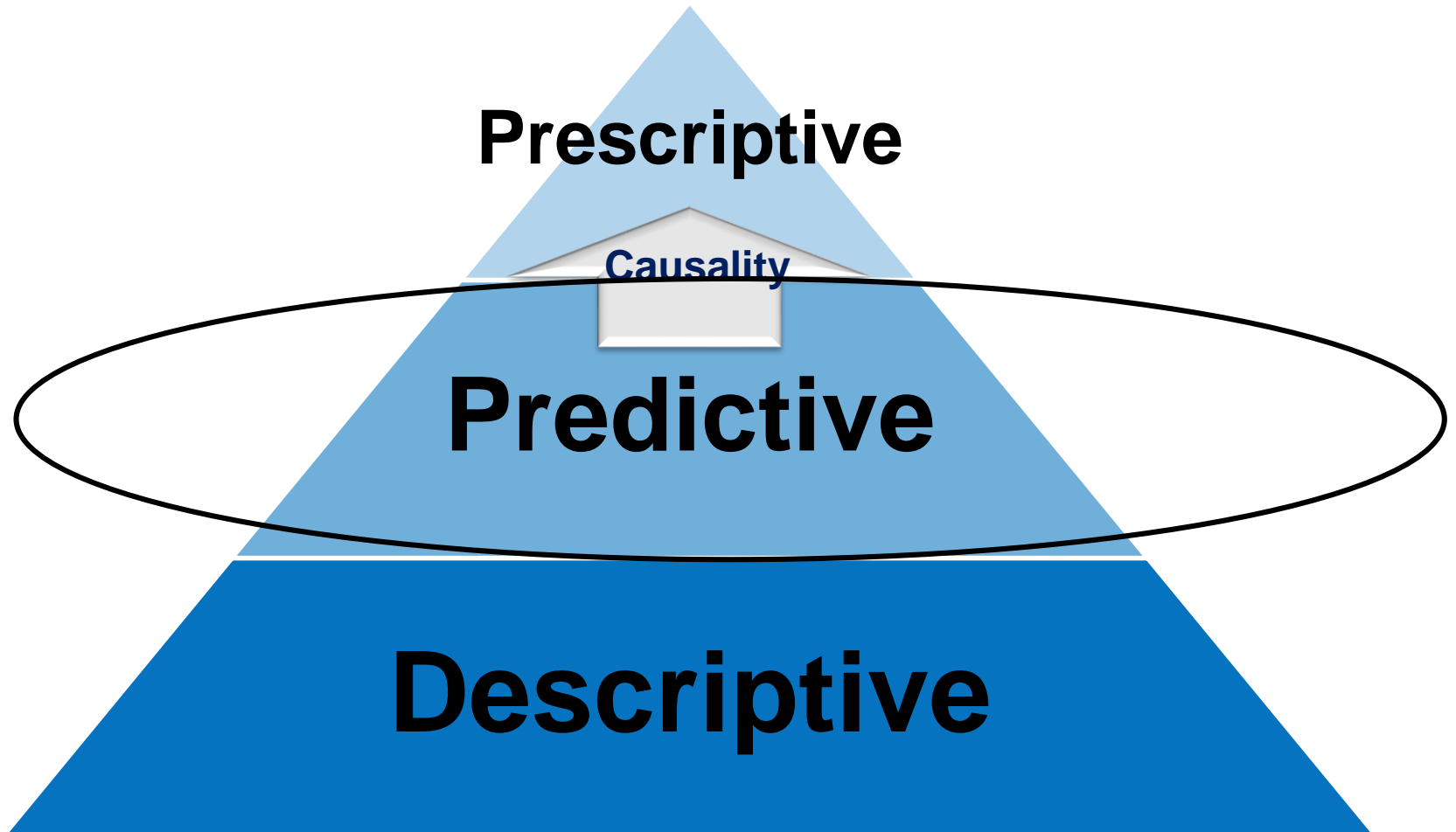
# Foundational Programming Skills and Computer Science

# 3. Descriptive Analytics, Exploratory Data Analysis, and Data Visualization



Graphics by Leland Wilkinson with permission
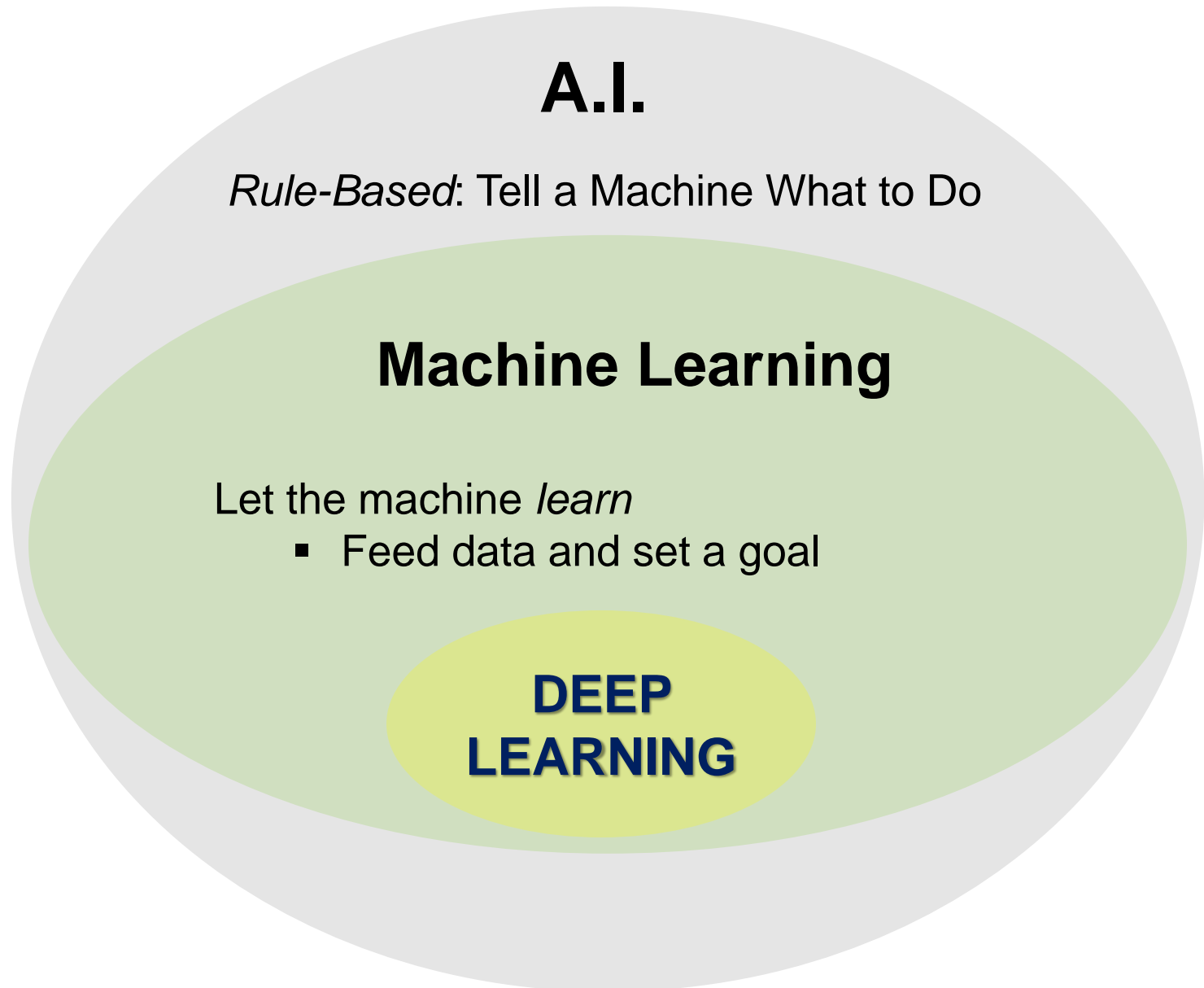
See also a new thought-provoking paper by Efron (2020) and a classic one by Breiman (2001)

# Decision Tree

**How well they are doing financially – Illustrative Only**



Overall — 50% well

Savings:
- Hi → 80% well
- Lo → 25% well

Protection:
- Hi → 90% well
- Lo → 60% well
- Hi → 50% well
- Lo → 15% well

# A.I.

*Rule-Based*: Tell a Machine What to Do

## Machine Learning

Let the machine *learn*
- Feed data and set a goal

### DEEP LEARNING

**Inside a Multi-Layer Perceptron (MLP) neural network, it is a set of nonlinear functions[1].**

The special composite function leads to a **Universal Approximator** to ANY functions.
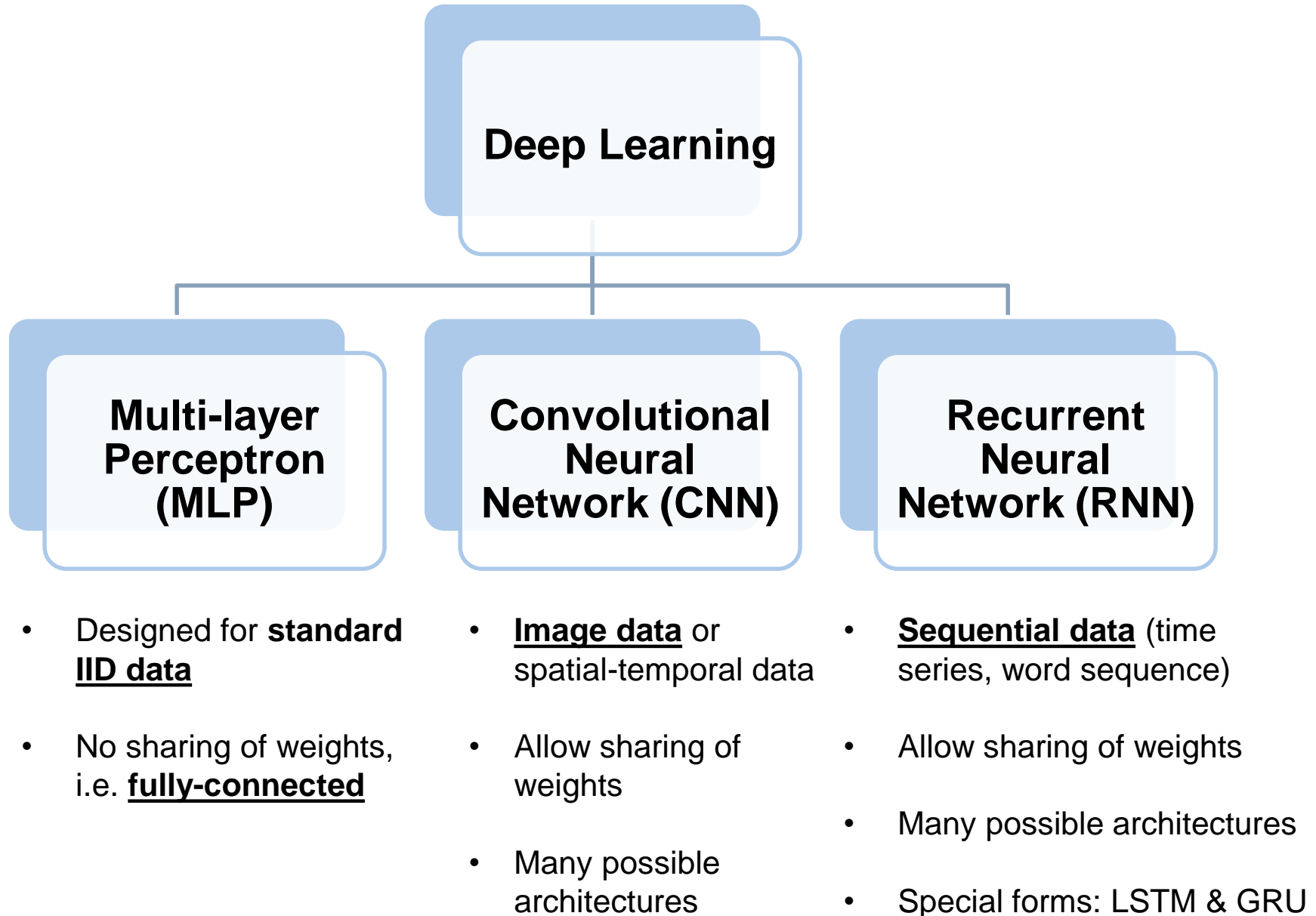


$$Y = w_{2,0} + \Sigma\, w_{2,j}\, I_j \text{ , where } I_j = 1/(1 + \exp(-Z_j)), \; Z_j = \Sigma\, w_{1,kj}\, x_k$$

**Deep Learning: At least 2 Hidden Layers**

[1] Other common activation functions include ReLu (most popular) and Tanh

# Types of Deep Learning

```
                    Deep Learning

   Multi-layer      Convolutional        Recurrent
   Perceptron       Neural               Neural
   (MLP)            Network (CNN)        Network (RNN)
```

- Designed for **standard IID data**

- No sharing of weights, i.e. **fully-connected**

- **Image data** or spatial-temporal data

- Allow sharing of weights

- Many possible architectures

- **Sequential data** (time series, word sequence)

- Allow sharing of weights

- Many possible architectures

- Special forms: LSTM & GRU

# Deep Learning for Medical Use Case

**Goal:** Predict Total Joint Replacement (TJR)

**Data:** De-identified claims data with detailed individual level time series of medical codes (diagnosis, procedure, etc.)

**Approach:**
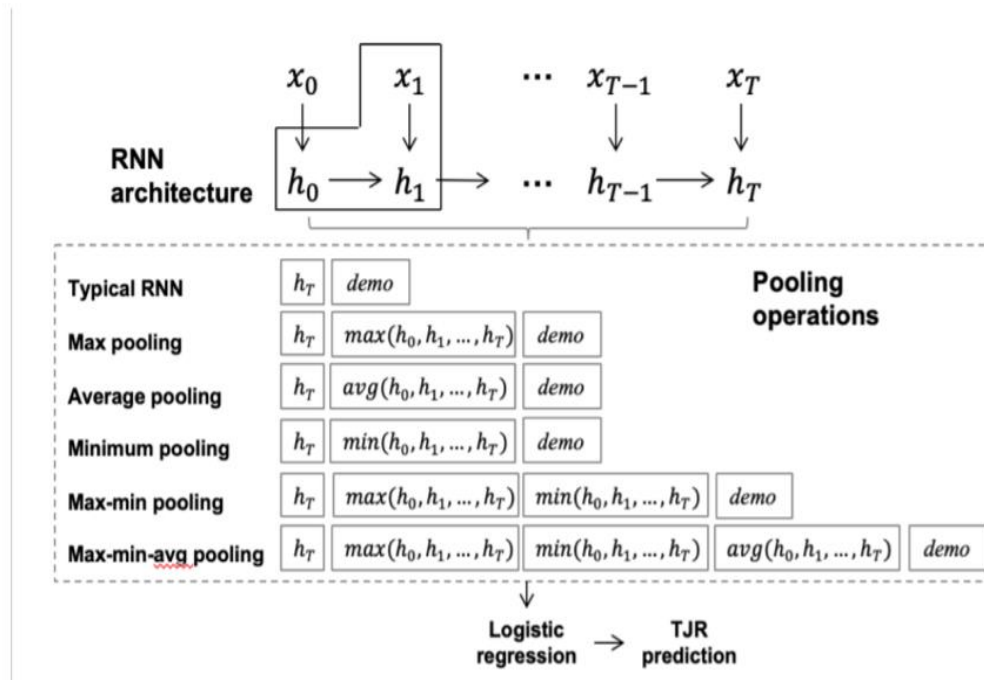Compare various deep learning architectures (CNN, RNN/GRU) with Lasso Logistic and RF
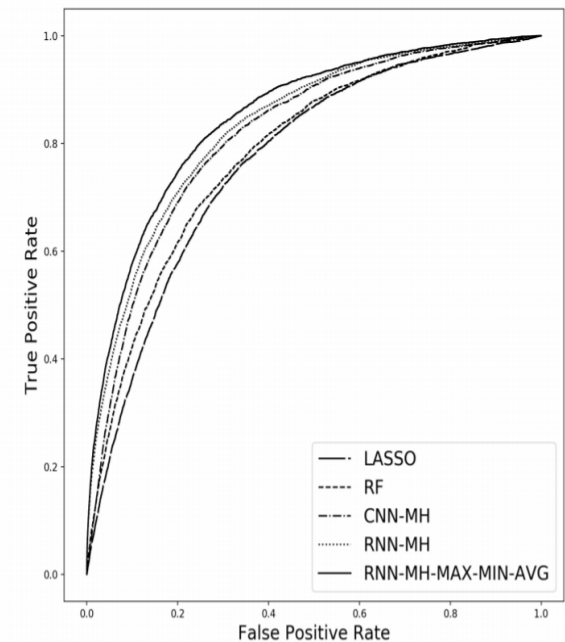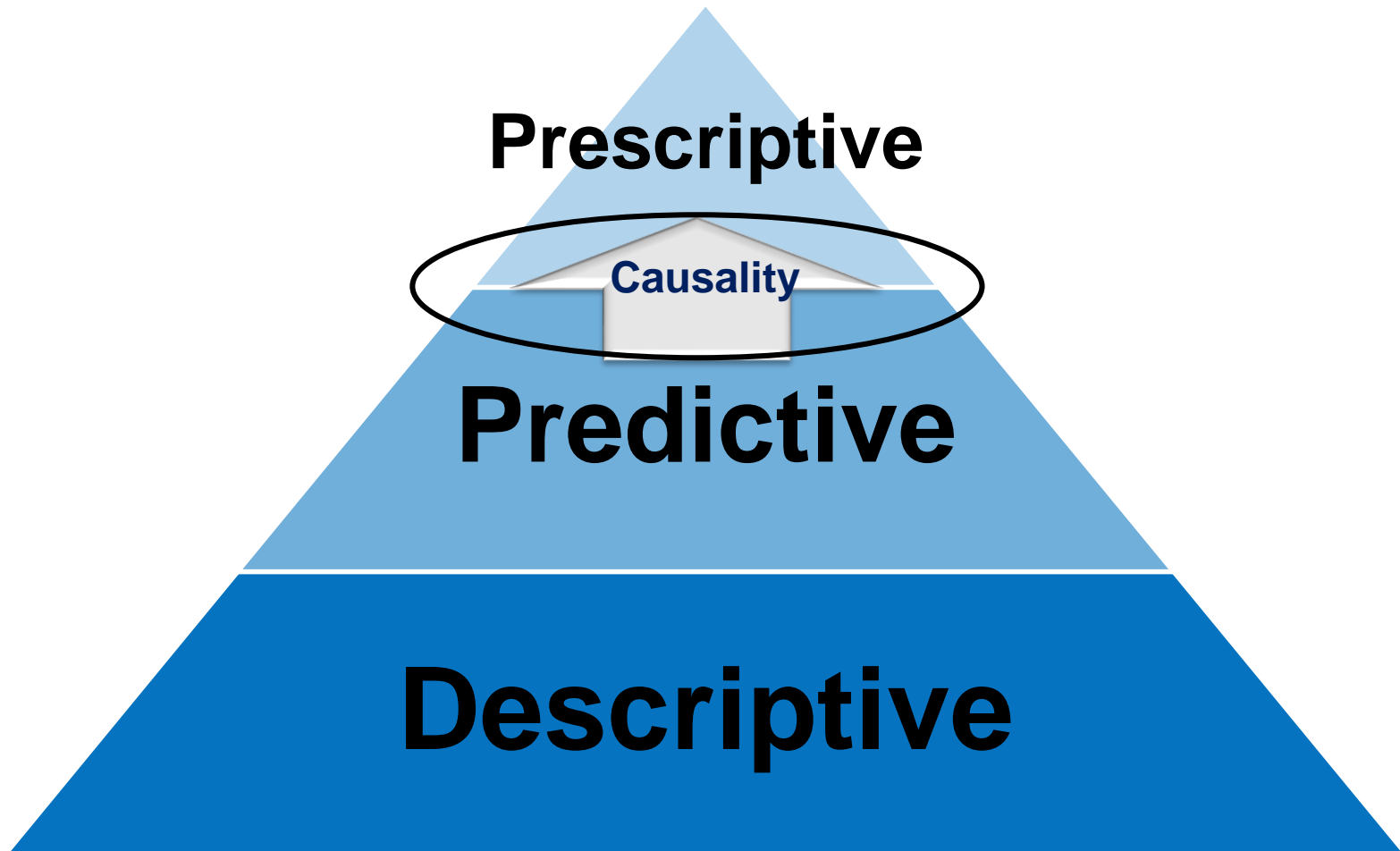
**Figure 1:** The RNN architecture and pooling operations.

**Figure 2:** Comparison of ROC for different models trained with 2014 and 2015

Source: Qiu et al (2019), with permission
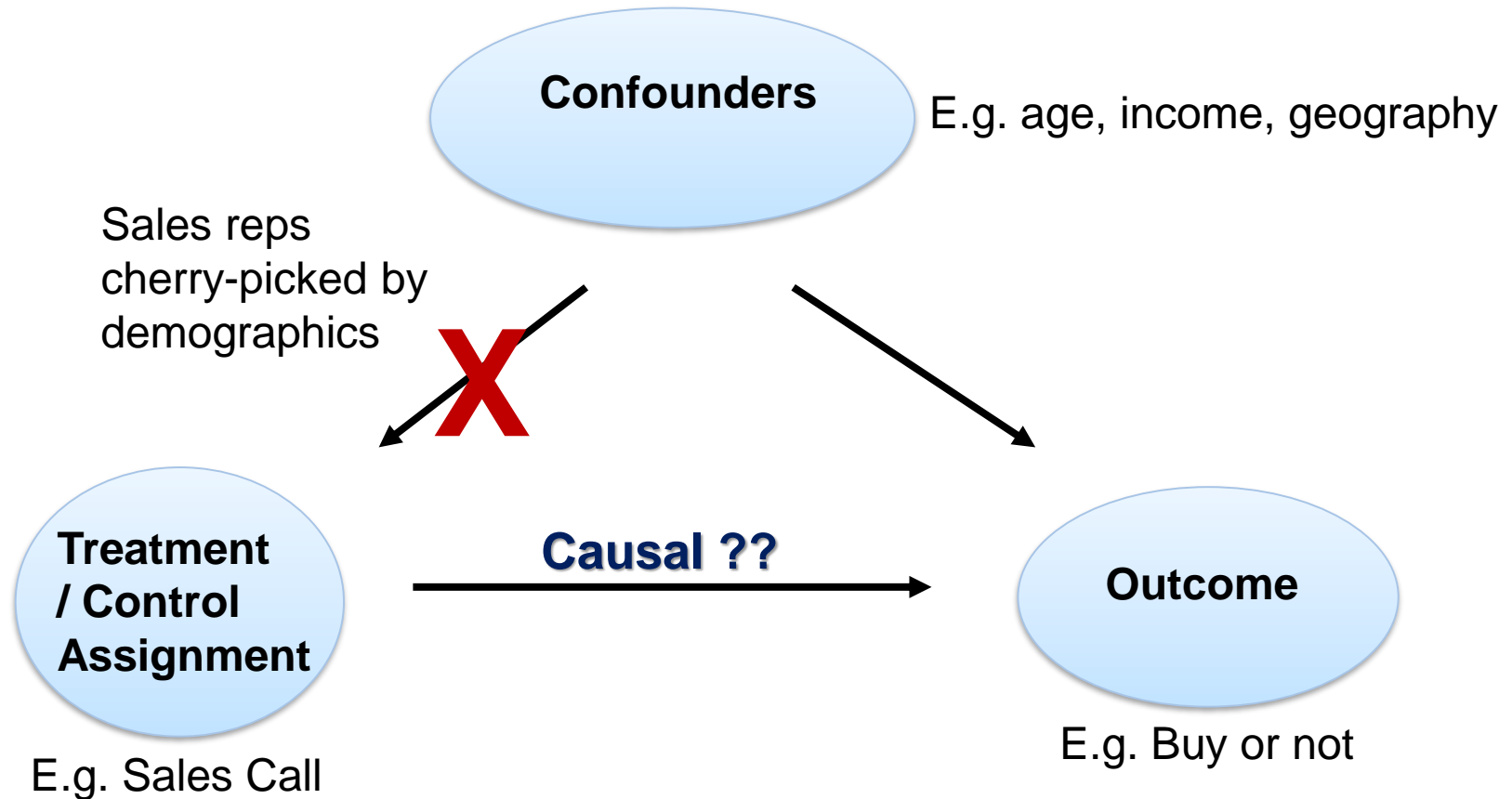
**Common Causality Related Questions in Business**

▶ **Price:** Would a price reduction generate high demand?

▶ **Promotion:** What are the impact of direct marketing and advertising?

▶ **Place:** What are the effects of store location and appearance on business outcomes?

▶ **Product:** Would an improvement in product feature be valuable to customers?

Similar questions can apply to other fields

# ► Blocking the "Back-Door" Path

**Goal: Measure Effect of Sales Campaign, using Historical Sales Data**



**Confounders** — E.g. age, income, geography

Sales reps cherry-picked by demographics

**X**

**Treatment / Control Assignment** — E.g. Sales Call

**Causal ??**

**Outcome** — E.g. Buy or not

Estimate **Average Treatment Effect** by breaking the Confounder-Treatment link: **Propensity Score Matching**

Next Level - Prioritize Future Sales Calls: **UPLIFT MODELING,** See Lo (2002, 2008)

# Personalized Medicine:
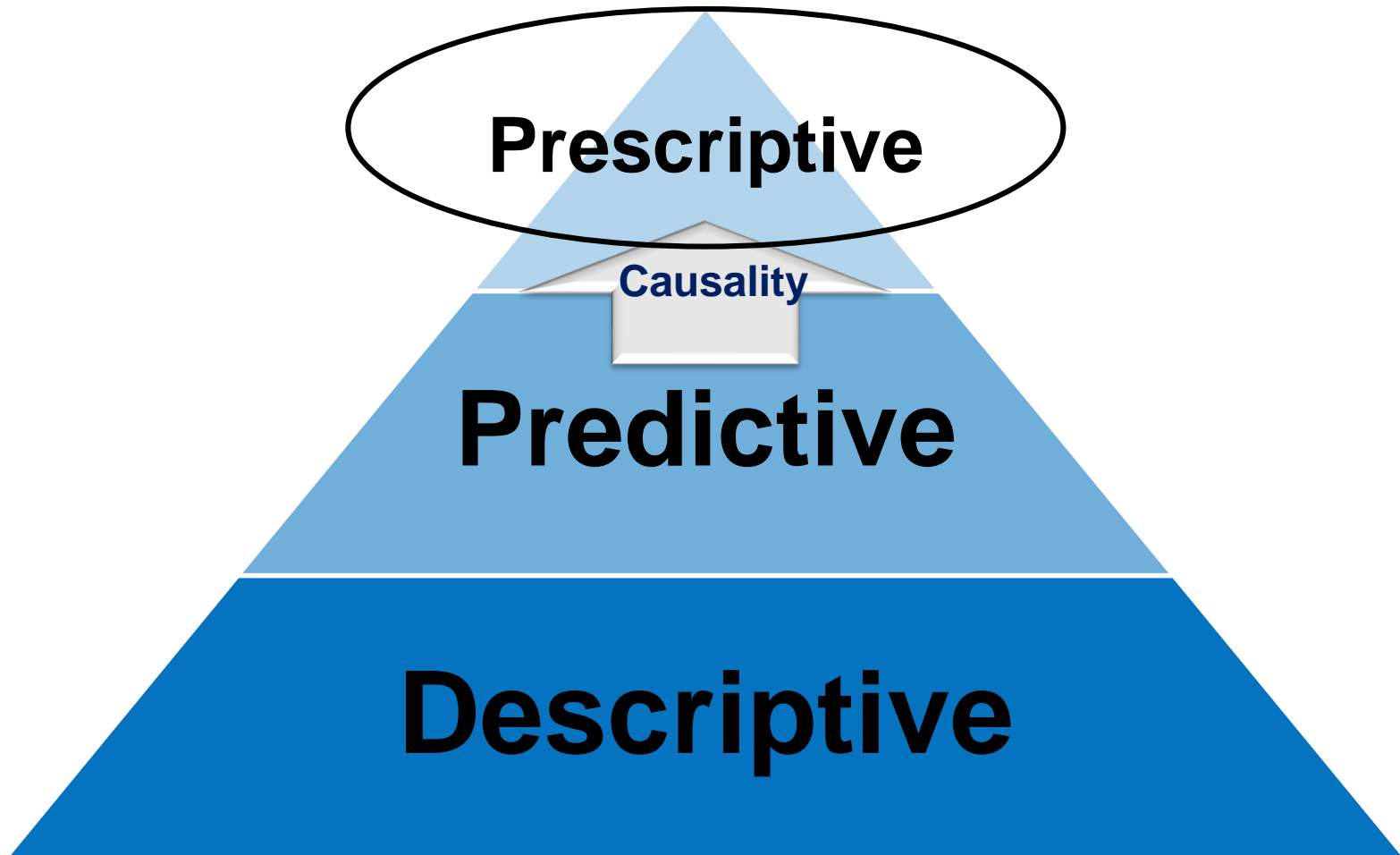## Stratify for more efficient treatment

**Clinical Benefit achieved if Receiving Placebo or no treatment**

|  | | YES | NO |
|---|---|---|---|
| **Clinical Benefit achieved if Receiving Active Treatment** | **YES** | Wasteful [Over-Treat] | Beneficial [Should-Treat] |
| | **NO** | Harmful [Do-Not-Treat] | Futile [Do-Not-Treat] |

Source: Chapter 3 of Yong (2015), with permission

# 7. Prescriptive Analytics and Optimization

| **Causal Inference** | **Prescriptive Analytics/Optimization** |
|---|---|

▶**Price:** **Would a price discount generate high demand?**

▶**Price:** **What is the <u>optimal price</u>?**

▶**Promotion:** **What are the Impact of direct marketing and advertising?**

▶**Promotion:** **How to <u>optimally invest</u> in direct marketing and advertising campaigns?**

▶**Place:** **What are the effects of store location and appearance on business outcomes?**

▶**Place:** **Where to open new stores? How should they look?**

▶**Product:** **Would an improvement in product feature be valuable to customers?**

▶**Product:** **What are <u>best</u> product features?**

# 7. Prescriptive Analytics and Optimization

- **Mindset – Objective Function, Constraints**
- **Mathematical Programming (MP)**
  - LP
  - ILP
  - MIP
  - QP
  - NP
  - DP & MDP
- Heuristics
- Multi-Objective Optimization (MOO)
- **Optimization Under Uncertainty**
  - Stochastic Programming
  - Robust Optimization
  - Mean-Variance Optimization, Nobel Econ 1990
- **Reinforcement Learning** – e.g., Alpha Go
- Stable Marriage and Kidney Exchange, Nobel Econ 2012

# Application: Customer Relationship Management (CRM)

**Individual Characteristics**

**Treatments**

Channel

Message/Offer

Individuals

Best deal…

No risk free trial…

Lowest price ever…

Super benefits…

Unbeatable service…

Top performance…

Scientifically proven…

**Lots of Possible Treatment Combinations**

31

▶ **8. Unstructured Data Analysis**

# 8. Unstructured Data Analysis

## Natural Language Processing (NLP) / Text Analytics

Document Processing
- Contract, legal
- Doctor's notes

Survey Verbatim

Search Engine

Chatbot

## Image Recognition

Radiology

Check Scan

Security & Biometrics

Insurance Claims

## Speech Analytics

Call Center:
- Sentiment Analysis
- Topic Modeling
- Features

# 8. Unstructured Data Analysis

## Natural Language Processing (NLP) / Text Analytics

- Computational Linguistics

- Advanced – **word embedding**, **deep learning** based (esp. **RNN**), Attention, Transformers, ELMO, BERT, etc.

- Specific applications: search, chatbot (QA), topic modeling, sentiment analysis

## Image Recognition

- Convolutional Neural Network (**CNN** or Convnet)
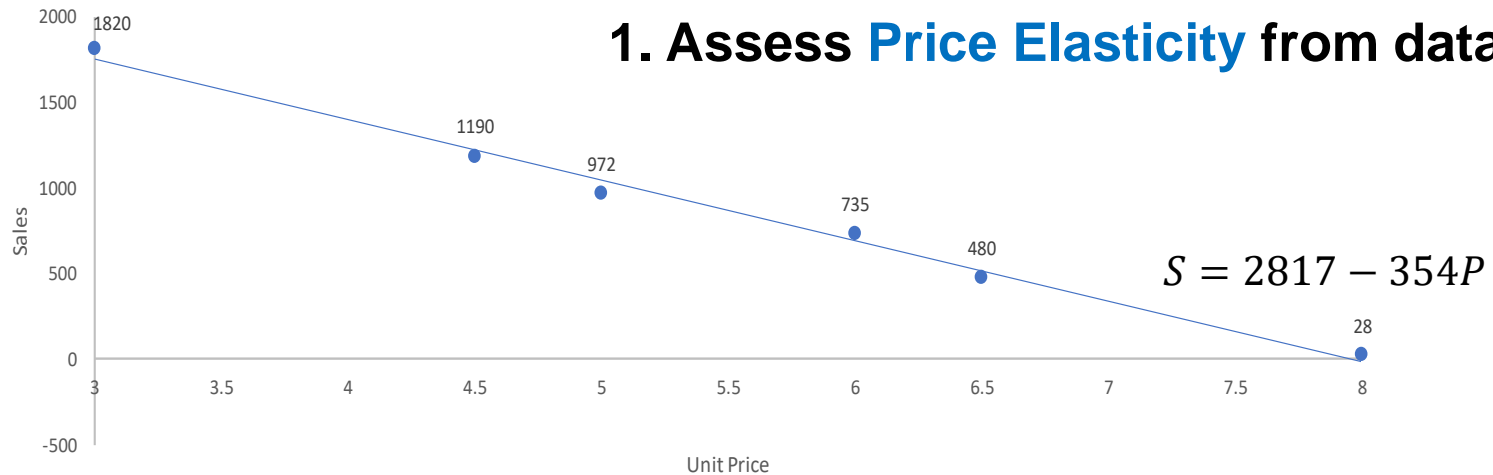- Computer Vision (OCR, R-CNN)

## Speech Analytics

- Language Model and Acoustic Model
- Hidden Markov Model (HMM)
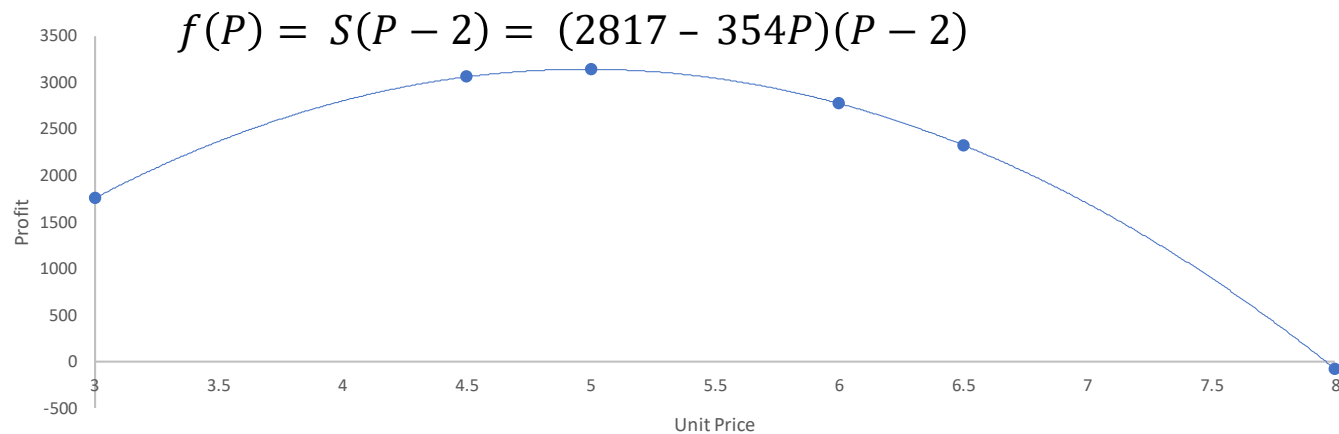- Deep Learning

# ▶9. Social Sciences and Data Science Ethics

**1. Assess Price Elasticity from data**

$$S = 2817 - 354P$$

**2. Determine the Optimal Price**

$$f(P) = S(P-2) = (2817 - 354P)(P-2)$$

# Behavioral Economics = Economics + Psychology

## Prospect Theory



Value

Losses — Gains

Loss Aversion

Daniel Kahneman,
Nobel Econ 2002

See Kahneman (2011)

## Nudge Theory

- Opt-in vs Opt-out

- Choice architecture - # choices

- Language Framing
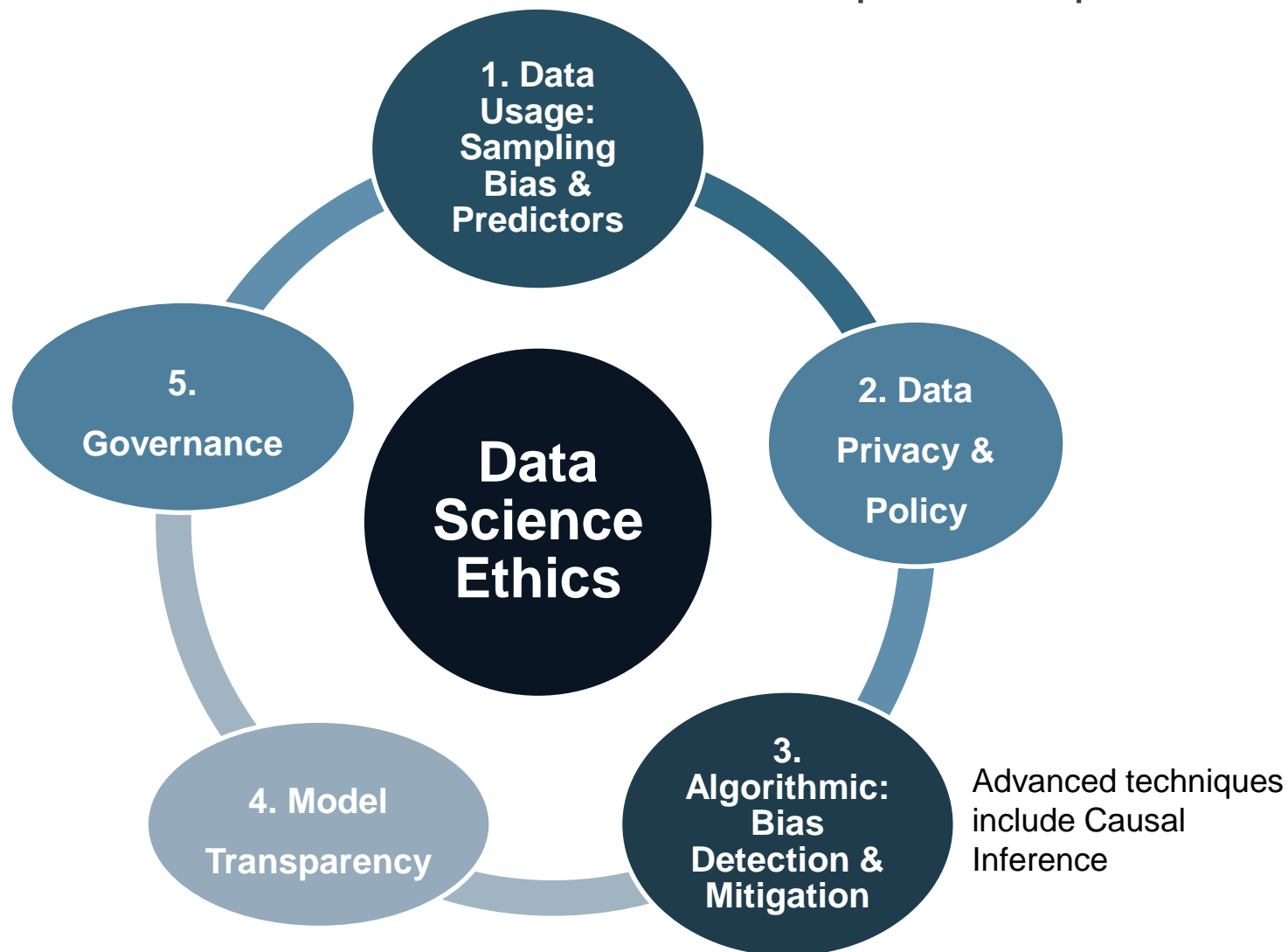
Can be **tested** and **modeled statistically** as input to behavioral change optimization

Richard Thaler,
Nobel Econ 2017

See Thaler & Sunstein (2009)

# Data Science Ethics

## Data Science Ethics involves Multiple Disciplines

1. Data Usage: Sampling Bias & Predictors

5. Governance

Data Science Ethics

2. Data Privacy & Policy

4. Model Transparency

3. Algorithmic: Bias Detection & Mitigation

Advanced techniques include Causal Inference

See O'Neil (2017), Boddington (2017), Lesile (2019), Russell (2019), Sandler & Bast (2019), ASA (2018), IFoA and RSS (2019), and so on.

# ▶ 10. Domain Knowledge and Application Areas

# Common Daily Usage of Data Science

## 1) Marketing & Sales

- Database Marketing / CRM
- Market Research
- Marketing Mix
- Marketing Strategy

## 2) Financial & Risk Management

- Modern Portfolio Theory (MPT)
- Risk Management: Market, Credit, Operational Risks
- Actuarial Science & Insurance

## 3) Operations Management, Supply Chain, and Logistics

- Call Center Analytics
- Logistics & Transportation
- Supply Chain Management
- Intelligent Automation

## 4) Healthcare & Biomedicine

- Health Informatics
- Drug Discovery
- Genomics
- Clinical trials
- Epidemiology

# Future NISS Tutorials, see https://www.niss.org/

1) **Analytical Consulting, Communication and Soft Skills**

2) **Computer Science, Programming, and Tools**

3) **Descriptive Analytics, Exploratory Data Analysis, and Data Visualization**

4) **Predictive Analytics and Machine Learning**

5) **Deep Learning**

6) **Causal Inference and Uplift Modeling**

7) **Predictive Analytics and Optimization**

8) **Unstructured Data Analysis**

9) **Social Sciences and Data Science Ethics**

10) **Domain Knowledge and Application Areas**

# Translation Between Statistics and AI / ML:
## Same or Similar Terminology

| Statistics / Economics / Epidemiology / Math | Data Science / AI / Data Mining |
|---|---|
| Statistical modeling | Machine Learning |
| Dependent Variable / Response Variable | Target Variable / Label |
| Independent Variable | Feature[1] |
| Parameters / coefficients | Weights |
| Intercept | Bias[2] |
| Estimation | Training |
| Out-of-Sample / Holdout Sample | Test Data |
| Regression / Classification | Supervised Learning |
| Cluster Analysis / PCA / Factor Analysis / SVD | Unsupervised Learning |
| Variable Selection | Feature Selection |
| Dimension Reduction | Feature Reduction |
| Data point / observation | Instance / Sample[3] / Example |
| Outlier Detection | Anomaly Detection |
| Log likelihood function of a binary variable | Cross Entropy |
| Logistic function | Sigmoid function |
| Multinomial Logit | Softmax |
| Dummy Coding | One-hot Coding |
| Misclassification Table | Confusion Matrix |
| Bayesian Computation | Probabilistic Programming |
| Approximate Dynamic Programming/Markov Decision Process | Reinforcement Learning |
| Randomized Controlled Trial (RCT) | A/B Testing |
| Factorial Design | Multivariate Testing (MVT) |
| Time series data | Sequential data |
| Classification Matrix | Confusion Matrix |
| Power [P(Reject H0 | H1 is true) or 1-P(Type II error)] | Recall |
| False Discovery Rate (FDR) | $1 -$ Precision |
| Average Treatment Effect (ATE) | Lift (Marketing) |
| Heterogeneous Treatment Effect (Econ.) | Uplift Modeling |
| Or Conditional Average Treatment Effect (CATE; Econ.) | Uplift Modeling |
| Or, Effect Modification (Epidemiology) | Uplift Modeling |
| Or, Impactibility Modeling (Health) | Uplift Modeling |
| Or, Subgroup Analysis (Biostat) | Uplift Modeling |

[1] A feature can also be a function of original variables.

[2] The standard statistical definition of Bias is the discrepancy between the actual value of an unknown parameter and the expected value of its estimator. Such definition is also used in machine learning, which is totally different from the Intercept-equivalent meaning in neural networks.

[3] The traditional definition of a sample refers to a subset of the population, which is a collection of observations. In some AI/ML literature, a single observation is sometimes called a sample.

# ►References

American Statistical Association (2018), "Ethical Guidelines for Statistical Practice."

Boddington, Paula (2017) *Towards a Code of Ethics for Artificial Intelligence*. Springer.

Breiman, Leo (2001), "Statistical Modeling: The Two Cultures," *Statistical Science*, 16, p.199-231.

Leslie, David (2019) "Understanding Artificial Intelligence Ethics and Safety." The Alan Turing Institute.

Efron, Bradley (2020), "Prediction, Estimation, and Attribution," *Journal of the American Statistical Association*, v.115, no.530, p.646-655.

Freedman, D. (2010). *Statistical Methods and Causal Inference*. Cambridge.

Hamburg, M.A. and Collins, F.S. (2010). "The path to personalized medicine." *The New England Journal of Medicine*, 363;4, p.301-304.

Hastie, T., R. Tibshirani, and J. Friedman (2009), *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, 2nd ed., Springer.

Haughton, D., Haughton, J., and Lo, V.S.Y. (2020, expected) *Cause-and-Effect Business Analytics*, CRC/Chapman & Hall.

Holland, C. (2005). *Breakthrough Business Results with MVT*, Wiley.

Independent High-Level Expert group on Artificial Intelligence, set up by the European Commission (8 April 2019) "Ethics Guidelines for Trustworthy AI" .Retrieved from: https://ai.bsa.org/wp-content/uploads/2019/09/AIHLEG_EthicsGuidelinesforTrustworthyAI-ENpdf.pdf

Kahneman, Daniel (2011). *Thinking, Fast and Slow*, FSG.

Kane, K., Lo, V.S.Y., and Zheng, J. (2014) "Mining for the Truly Responsive Customers and Prospects Using True-Lift Modeling: Comparison of New and Existing Methods." *Journal of Marketing Analytics*, v.2, Issue 4, p.218-238.

Kearns, Michael and Aaron Roth (2019) *The Ethical Algorithm*. Oxford University Press.

Lai, Lilly Y.-T. (2006) Influential Marketing: A New Direct Marketing Strategy Addressing the Existence of Voluntary Buyers. Master of Science thesis, Simon Fraser University School of Computing Science, Burnaby, BC, Canada.

Leamer, Edward E. *Macroeconomic Patterns and Stories: A Guide for MBAs,* Springer.

Lo, V.S.Y. (2002) "The True Lift Model – A Novel Data Mining Approach to Response Modeling in Database Marketing." *SIGKDD Explorations* 4, Issue 2, p.78-86, at: http://www.acm.org/sigs/sigkdd/explorations/issues/4-2-2002-12/lo.pdf

Lo, V.S.Y. (2008), "New Opportunities in Marketing Data Mining," in *Encyclopedia of Data Warehousing and Mining*, Wang (2008) ed., 2nd edition, Idea Group Publishing.

Lo, Victor S.Y. (2019), "Searching for the Perfect Unicorn," *Analytics Magazine.*

Lo, V.S.Y. and D. Pachamanova (2015), "A Practical Approach to Treatment Optimization While Accounting for Estimation Risk," *Journal of Marketing Analytics*, v.3, Issue 2, p.79-95.

43

Lynch, S. (2017), "Andrew Ng: Why AI is the New Electricity," *Insights by Stanford Business.*

Morgan, S.L. and Winship C. (2007). *Counterfactuals and Causal Inference.* Cambridge University Press.

O'Neil, Cathy (2017). *Weapons of Math Destruction: How Big Data Increases Inequality and Threatens Democracy*, Broadway Books.

Oxelheim, Lars and Clas Wihlborg (2008), *Corporate Decision-Making with Macroeconomic Uncertainty*, Oxford University Press.

Pachamanova, D., V.S.Y. Lo, and N. Gulpinar (2020) "Uncertainty Representation and Risk Management in Direct Segmented Market," *Journal of Marketing Management,* v.36, Issue 1-2*.*

Pearl, Judea (2000), *Causality*. Cambridge University Press.

Pearl, Judea and Dana MacKenzie (2018), *The Book of Why: The New Science of Cause and Effect.* Basic Books.

Porter, Daniel (2013) Pinpointing the Persuadables: Convincing the Right Voters to Support Barack Obama. Presented at Predictive Analytics World; Oct, Boston, MA; http://www.predictiveanalyticsworld.com/patimes/pinpointing-the-persuadables-convincing-the-right-voters-to-support-barack-obama/ (available with free subscription).

Qiu, R., Y. Jia, F. Wang, P. Divakarmurthy, S. Vinod, B. Sabir and, M. Hadzikadic (2019), "Predictive Modeling of the Total Joint Replacement Surgery Risk: a Deep Learning Based Approach with Claims Data," *AMIA Summits on Translational Science Proceedings,* American Medical Informatics Association.

Radcliffe, N.J. and Surry, P. (1999). "Differential response analysis: modeling true response by isolating the effect of a single action," *Proceedings of Credit Scoring and Credit Control VI*, Credit Research Centre, U. of Edinburgh Management School.

Radcliffe, N.J. (2007). "Using Control Groups to Target on Predicted Lift," *DMA Analytics Annual Journal*, Spring, p.14-21.

Reinsel, D., J. Gantz, and J. Rydning (2020), "The Digitization of the World: From Edge to Core," *IDC White Paper*.

Rubin, D.B. (2006), *Matched Sampling for Causal Effects*. Cambridge University Press.

Rubin, D.B. and Waterman, R.P. (2006), "Estimating the Causal Effects of Marketing Interventions Using Propensity Score Methodology," *Statistical Science*, p.206-222.

Russell, Stuart (2019) *Human Compatible: Artificial Intelligence and the Problem of Control*. Viking.

Sandler, Ronald and John Basl (2019) "Building Data and AI Ethics Committees." Northeastern University Ethics Institute and Accenture.

Thaler, Richard H. and Cass R. Sunstein. *Nudge: Improving Decisions About Health, Wealth, and Happiness.* Penguin Group.

The Institute and Faculty of Actuaries (IFoA) and the Royal Statistical Society (RSS) (2019) "A Guide for Ethical Data Science."

Wikipedia (2010), "Uplift Modeling," at http://en.wikipedia.org/wiki/Uplift_modelling

Yong, Florence H. (2015), "Quantitative Methods for Stratified Medicine," *PhD Dissertation*, Department of Biostatistics, Harvard T.H. Chan School of Public Health, Harvard University.
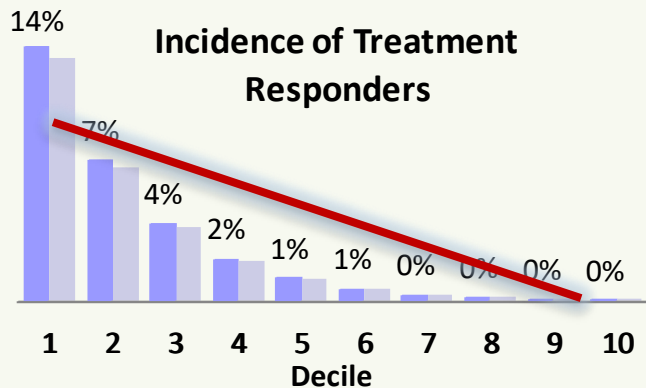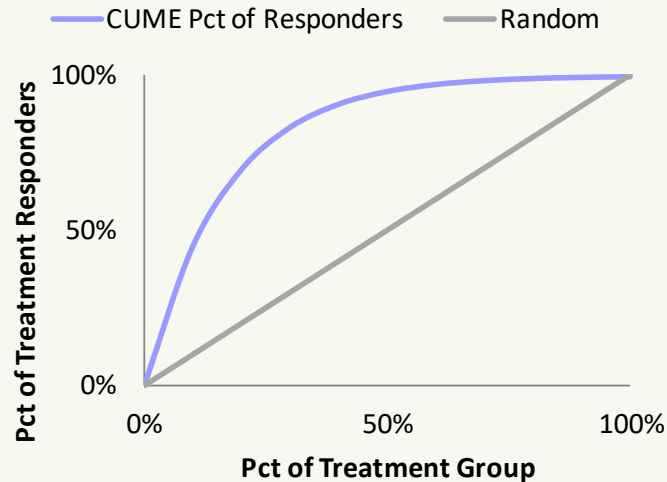
# APPENDIX

# History of Data Science

- Wu 1997, proposed:
  - Statistics → Data Science
  - Statistician → Data Scientist

- Cleveland 2001, proposed:
  - Enlarge the major areas of Statistics → Data Science

Source:
https://www.forbes.com/sites/gilpress/2013/05/28/a-very-short-history-of-data-science/#5a5a13cd55cf
https://course.ccs.neu.edu/cs7280sp16/CS7280-Spring16_files/50YearsOfDataScience.pdf
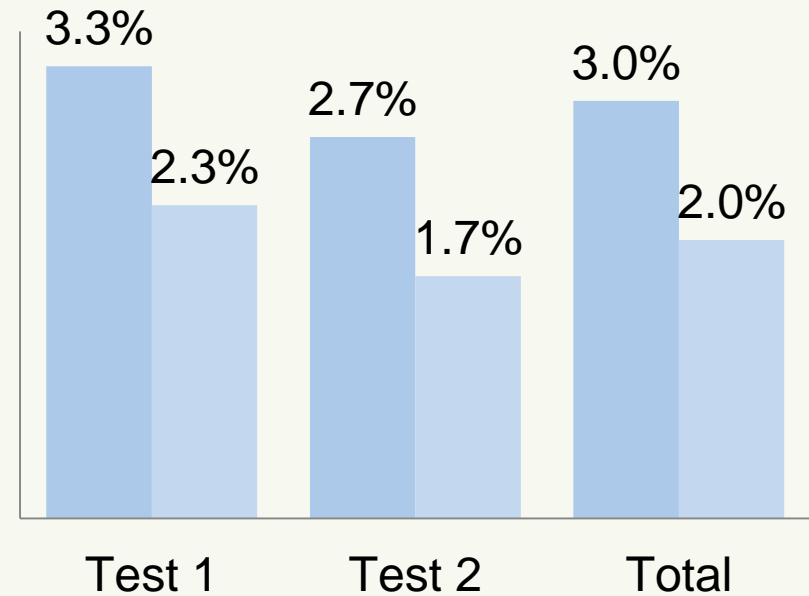
# What is the Right Way to Measure Lift?



**A successful response model**

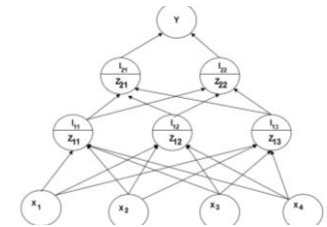**A successful treatment (e.g. marketing) program**

# ►5. Deep Learning

## History of A.I.: From Programming to Deep Learning



| 1840's | 1950's | 1970's-1980's | 2000's | Present |
|---|---|---|---|---|

**1. Birth of Programming**

- Ada Lovelace: *computers can never be as intelligent as humans*

**2. Birth of A.I.**

- Alan Turing: *a machine can possibly think for itself*

- Participants at Dartmouth workshop called it **Artificial Intelligence**

**3. Rule-Based Expert Systems (Classical A.I.)**

- *Rule-based*: hard-code with human expertise

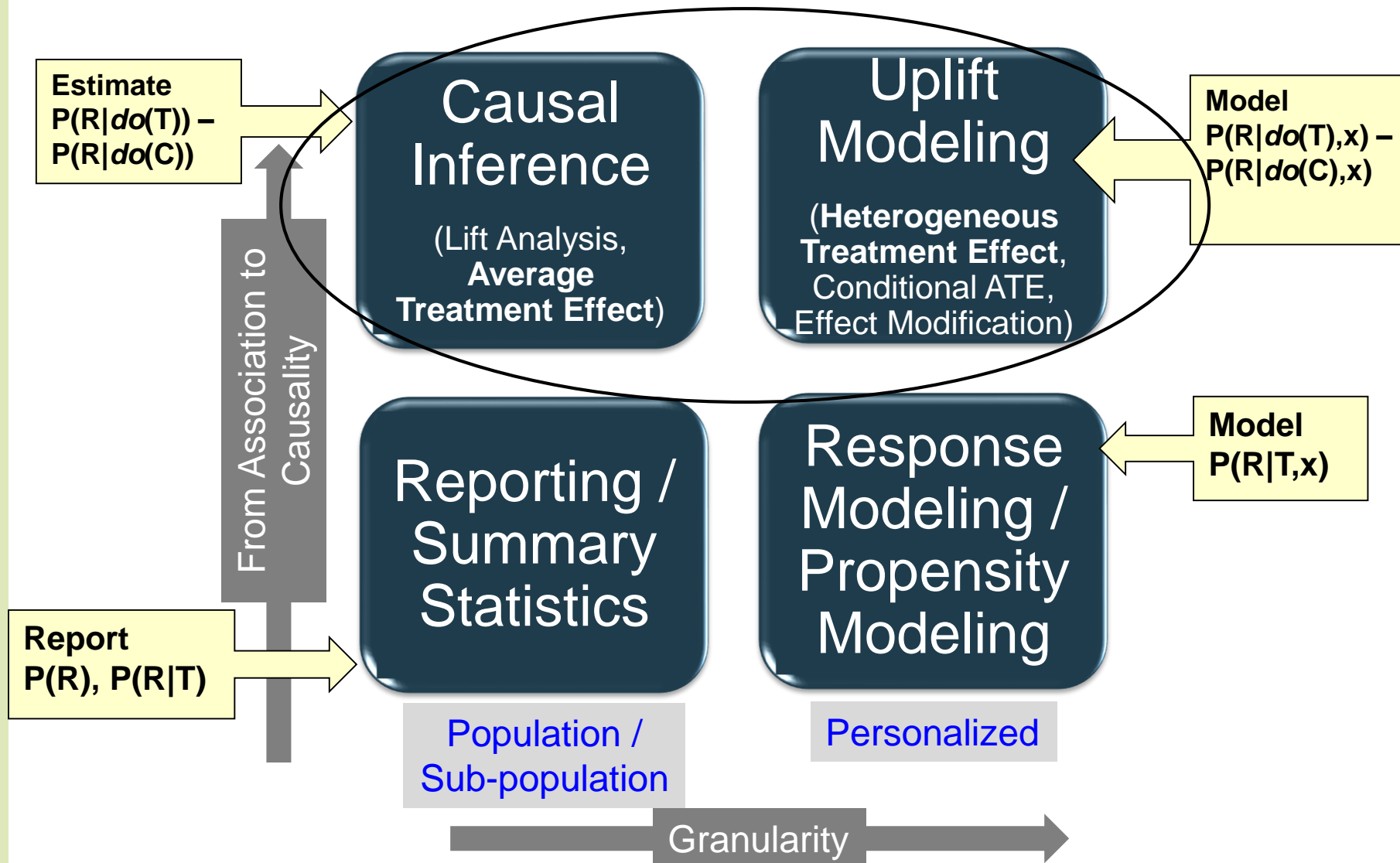- Work well on limited applications

**4. Machine Learning (Modern A.I.)**

- Feed *Big Data* data to an algorithm and set a goal

- Wide applications in medicine, marketing, finance, logistics, operations, and beyond

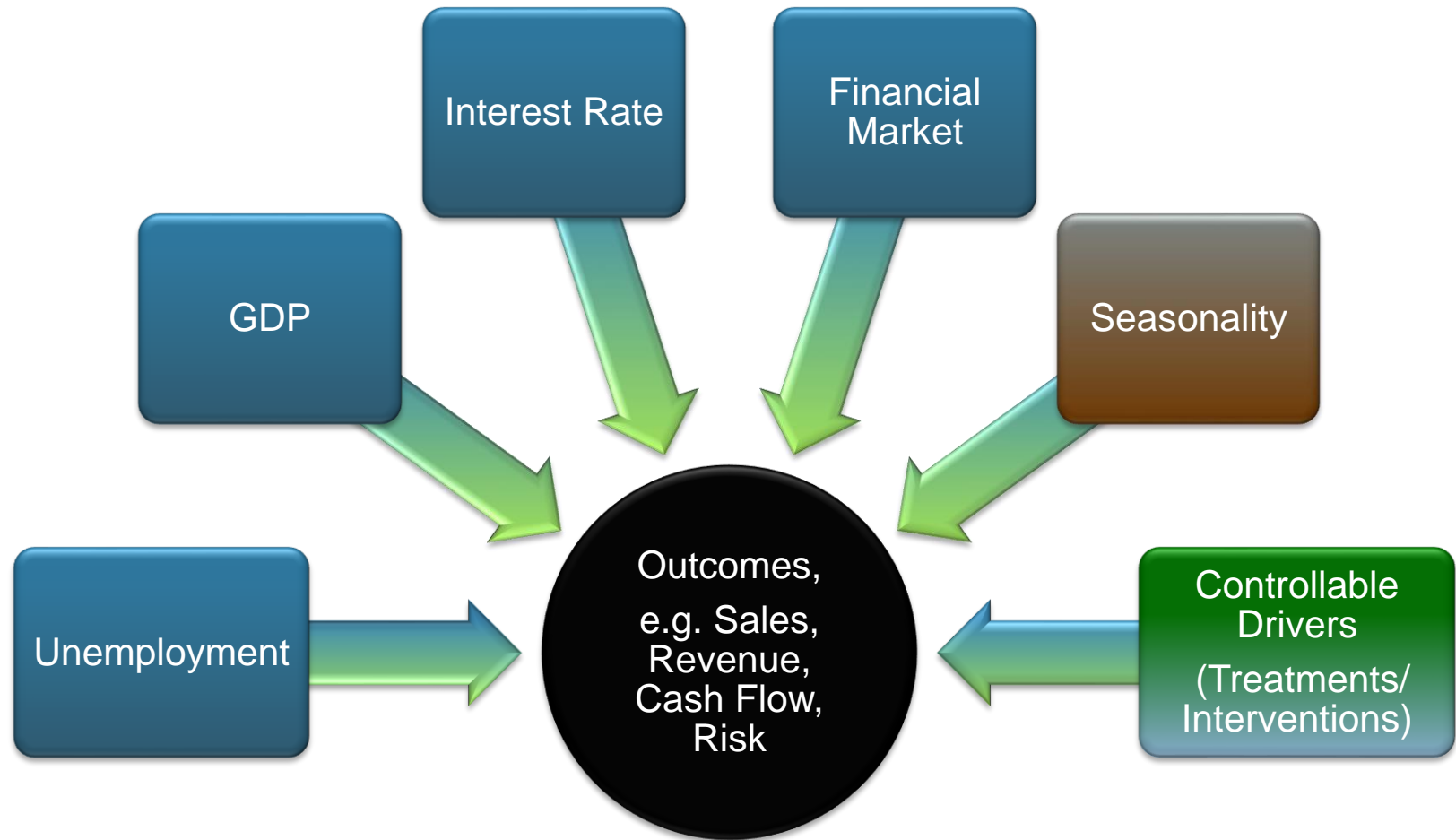**5. Deep Learning (Latest Machine Learning)**

- Deep Learning became widely used *for image processing, natural language processing (NLP),* and so on

- **Hinton, Bengio, LeCun won 2018 Turing Award**

---

- Most A.I.'s are designed to do a single task: **Narrow AI**

- *Can Machines Think? It depends…*

# Framework for Causal and Association Analysis

Estimate
P(R|*do*(T)) −
P(R|*do*(C))

Model
P(R|*do*(T),x) −
P(R|*do*(C),x)

## Causal Inference

(Lift Analysis, **Average Treatment Effect**)

## Uplift Modeling

(**Heterogeneous Treatment Effect**, Conditional ATE, Effect Modification)

From Association to Causality

## Reporting / Summary Statistics

## Response Modeling / Propensity Modeling

Model
P(R|T,x)

Report
P(R), P(R|T)

Population / Sub-population

Personalized

Granularity

For the *do*-operator above, see **Pearl (2000)** and **Pearl and MacKenzie (2018).**

49

# **Macroeconomics**: Sensitivity to Macroecon Factors



See Oxelheim and Wihlborg (2008) and Leamer (2009)