# Bayesian Sense-Making in Data Science

**Michael L. Thompson**, Ph.D.

https://www.linkedin.com/in/mlthomps/

**6th Annual BayesiaLab Conference in Chicago**

*Nov. 2, 2018*

# Outline

- **Prevalence of Bayesian Applications**

- **Whence Bayesian Analysis?**
  - The *Model Structure* Information Content Diagram
  - Motivation for Bayesian Sense-Making

- **Key Concepts of Bayesian Sense-Making**
  - Use Case: General Recommender/Advisor Systems

- **Future Implications**
  - References to Get Started

# Prevalence of Bayesian Applications
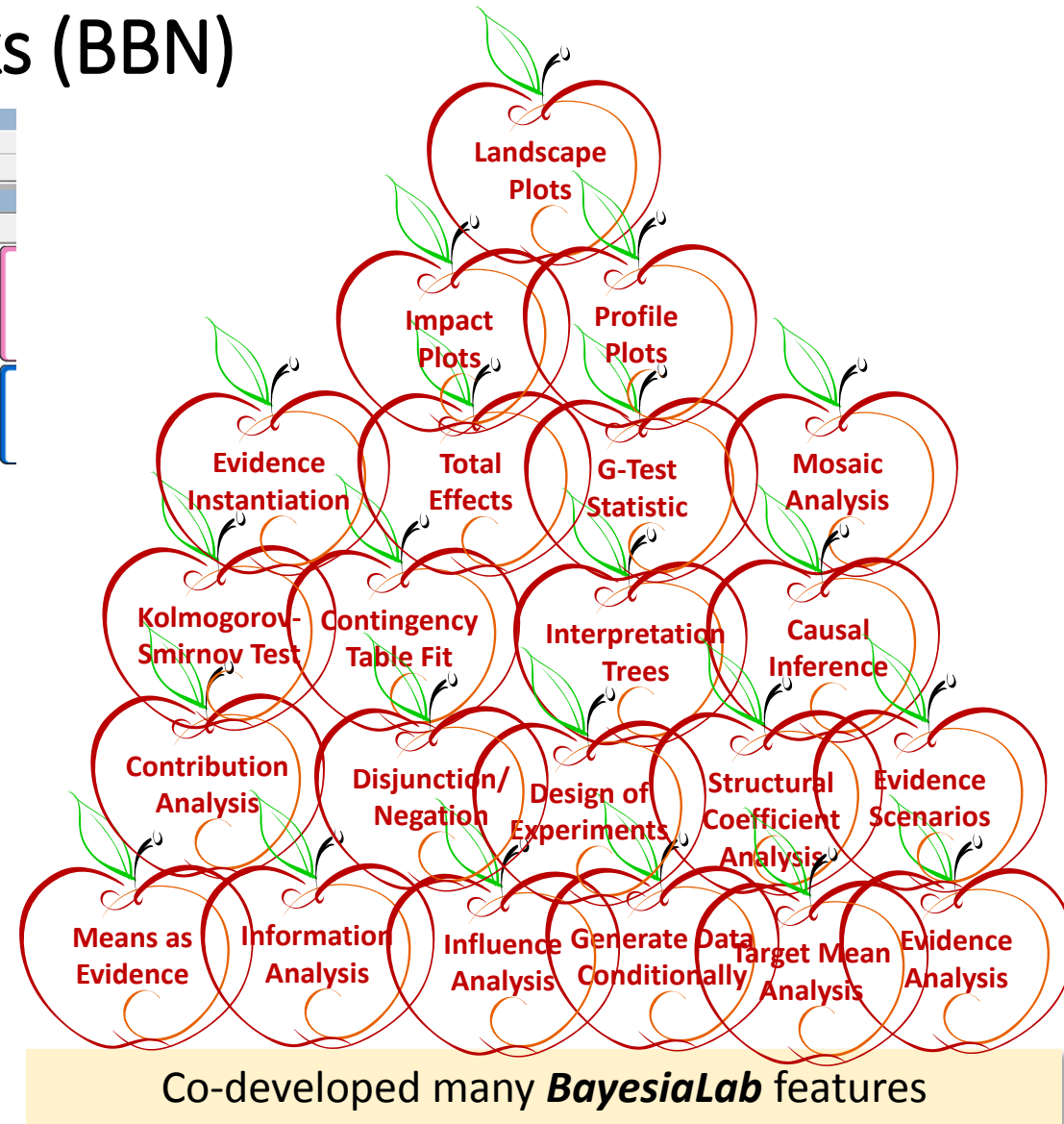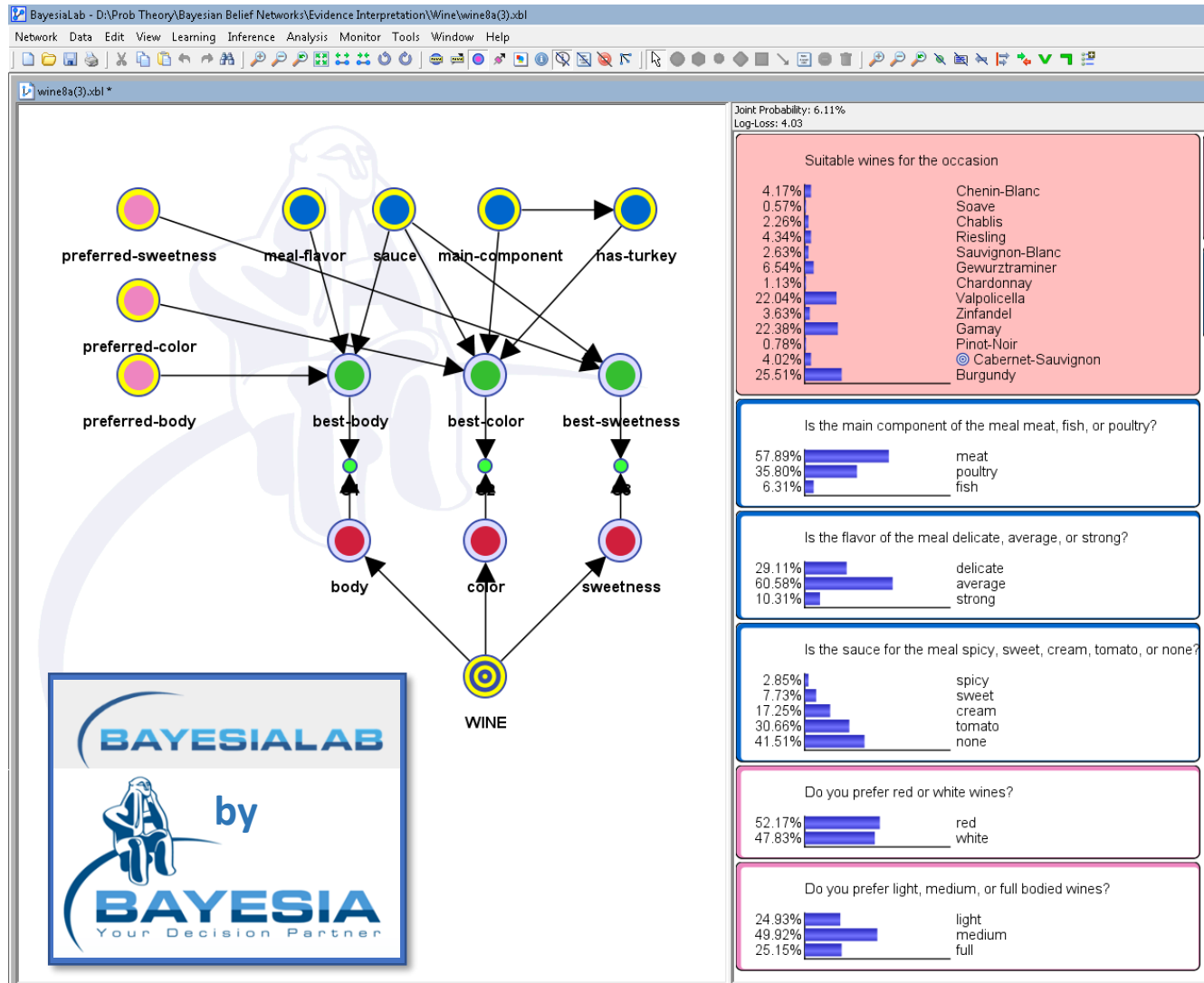# Bayesian Analysis to Delight Consumers at P&G



How Color Affects Decision

"…We have posed a … **basic question about consumer behavior, and the answer to this question is best captured by a multi-level, dichotomous, logistic regression model … using Bayesian inference by Markov Chain Monte Carlo simulation**,…."
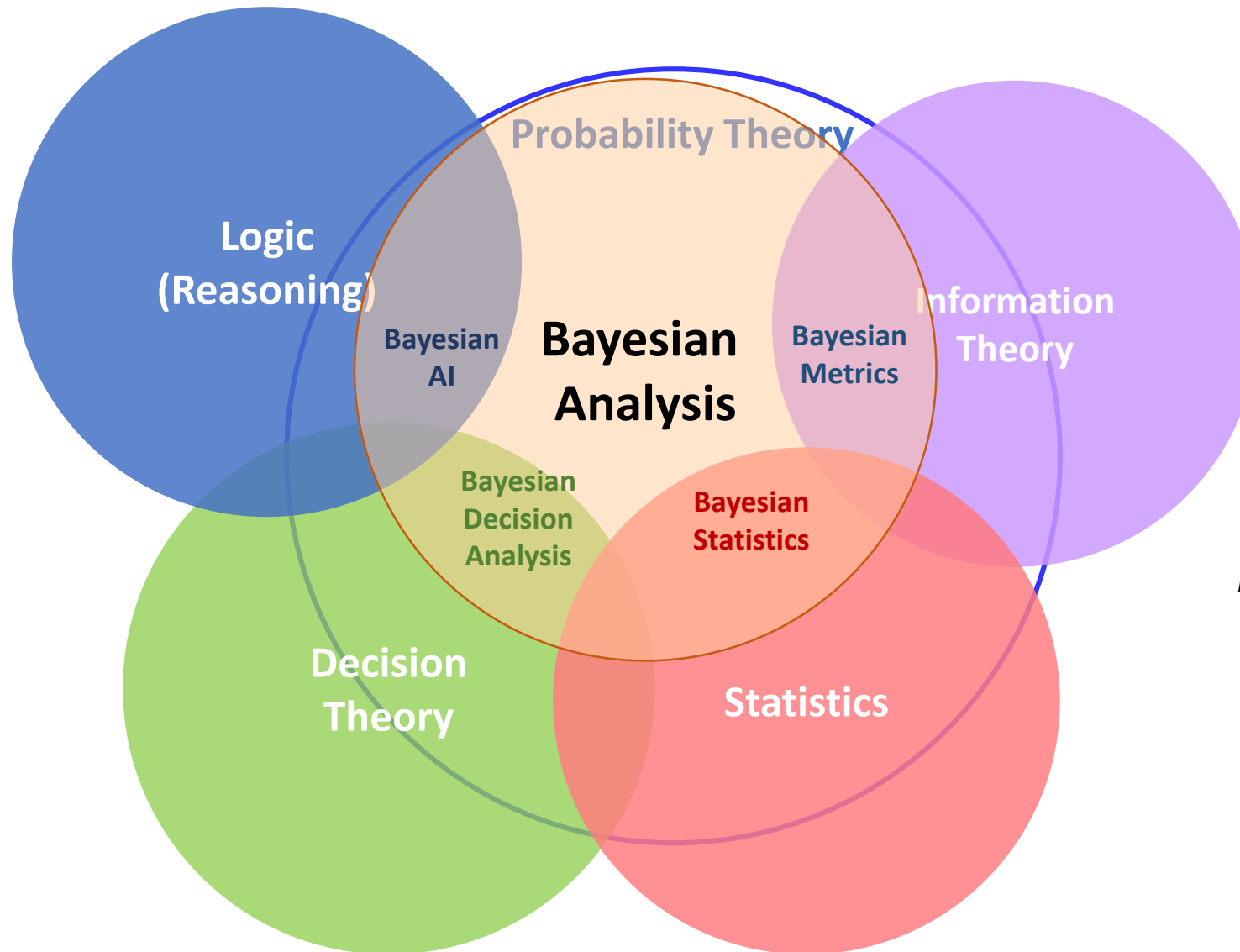
Thompson, Michael L., et al, *P&G Core Technologies J.*, 2000

# P&G/*Bayesia* Strategic Partnership
## *BayesiaLab* for Bayesian Belief Networks (BBN)



Co-developed many *BayesiaLab* features

# Whence Bayesian Analysis?



Probability Theory

Logic (Reasoning)

Information Theory

Bayesian AI

**Bayesian Analysis**

Bayesian Metrics

Bayesian Decision Analysis

Bayesian Statistics

Decision Theory

Statistics

*In short,
Bayesian Analysis is more than just adopting priors to model data.
It's reasoning about the world to learn and to drive decisions, i.e.,
Bayesian Sense-Making!*

# Model Structure: Sources of Information

Models are built by casting the information we have into mathematical functions.

$$V = \frac{4}{3}\pi r^3$$

- **K**nowledge representation
  - Domain knowledge (e.g., physical laws, theories of behavior, etc.)

$$y(x) = \theta_0 + \theta_1 x$$

- **M**athematical approximation
  - Simplifications and canonical functional forms (e.g., linear relationships, response surfaces, etc.)

$$y(x) = \sum_{j=1}^{M} w_j f(x; \Theta)$$

- **D**ata considerations
  - Flexible combinations of basis functions that grow with the data (e.g., nonparametric density estimation, multivariate analyses, etc.)

# Model Structure Information Content Diagram
## A ternary mixture diagram of information sources that dictate model structure

$$V = \frac{4}{3}\pi r^3$$

**K**  
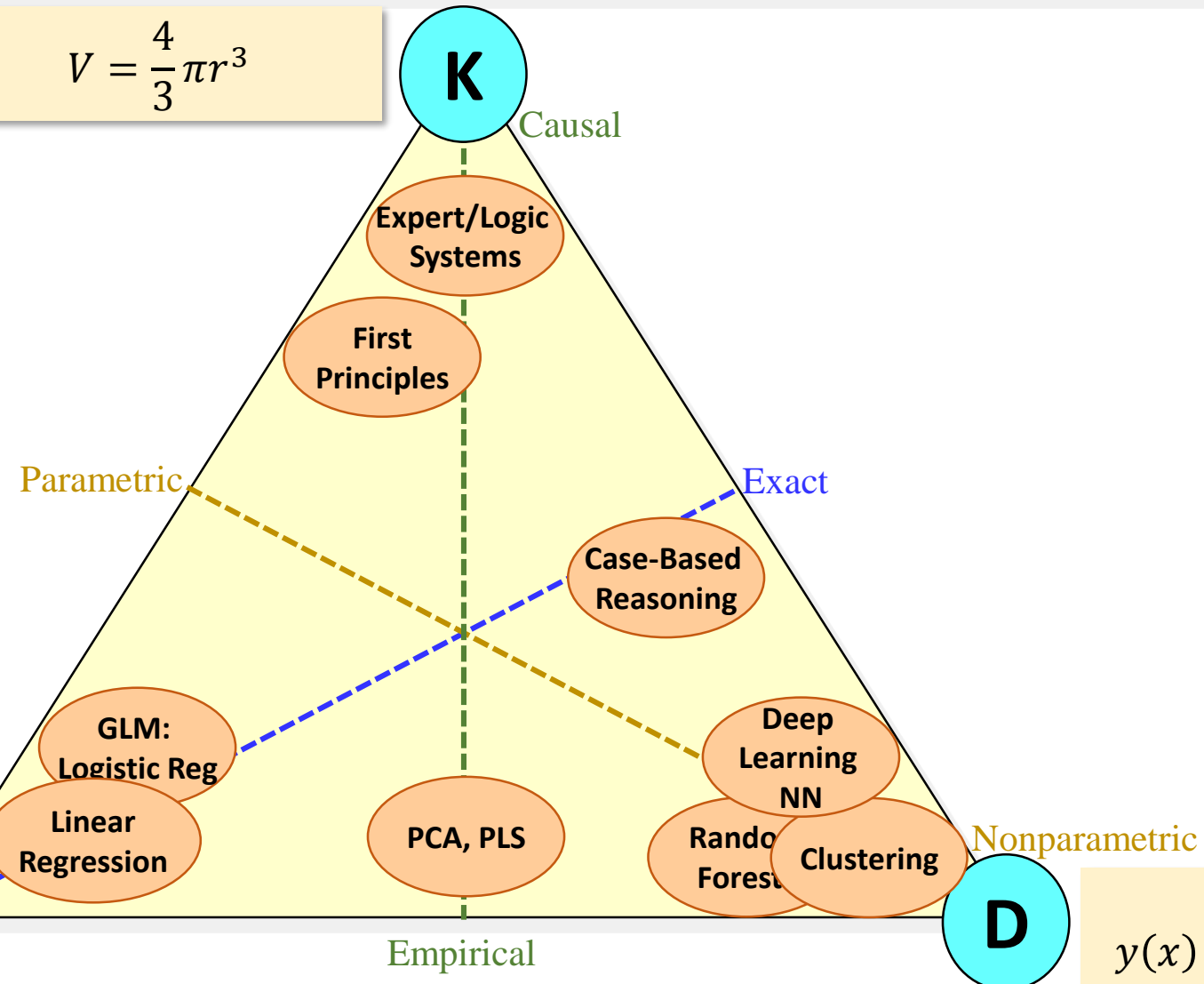Causal

Parametric    Exact

Approximate    Nonparametric

**M**    **D**

Empirical

$$y(x) = \theta_0 + \theta_1 x$$

$$y(x) = \sum_{j=1}^{M} w_j f(x; \Theta)$$

# Common Modeling Paradigms

$$V = \frac{4}{3}\pi r^3$$



K

Causal

Expert/Logic Systems

First Principles

Parametric

Exact

Case-Based Reasoning

GLM: Logistic Reg

Deep Learning NN

Approximate

Linear Regression

PCA, PLS

Random Forest

Clustering

Nonparametric

M

$$y(x) = \theta_0 + \theta_1 x$$

Empirical

D

$$y(x) = \sum_{j=1}^{M} w_j f(x; \Theta)$$

# Motivation for Bayesian Analysis
Fusion of Data Complex, Model Complex, Decision Complex

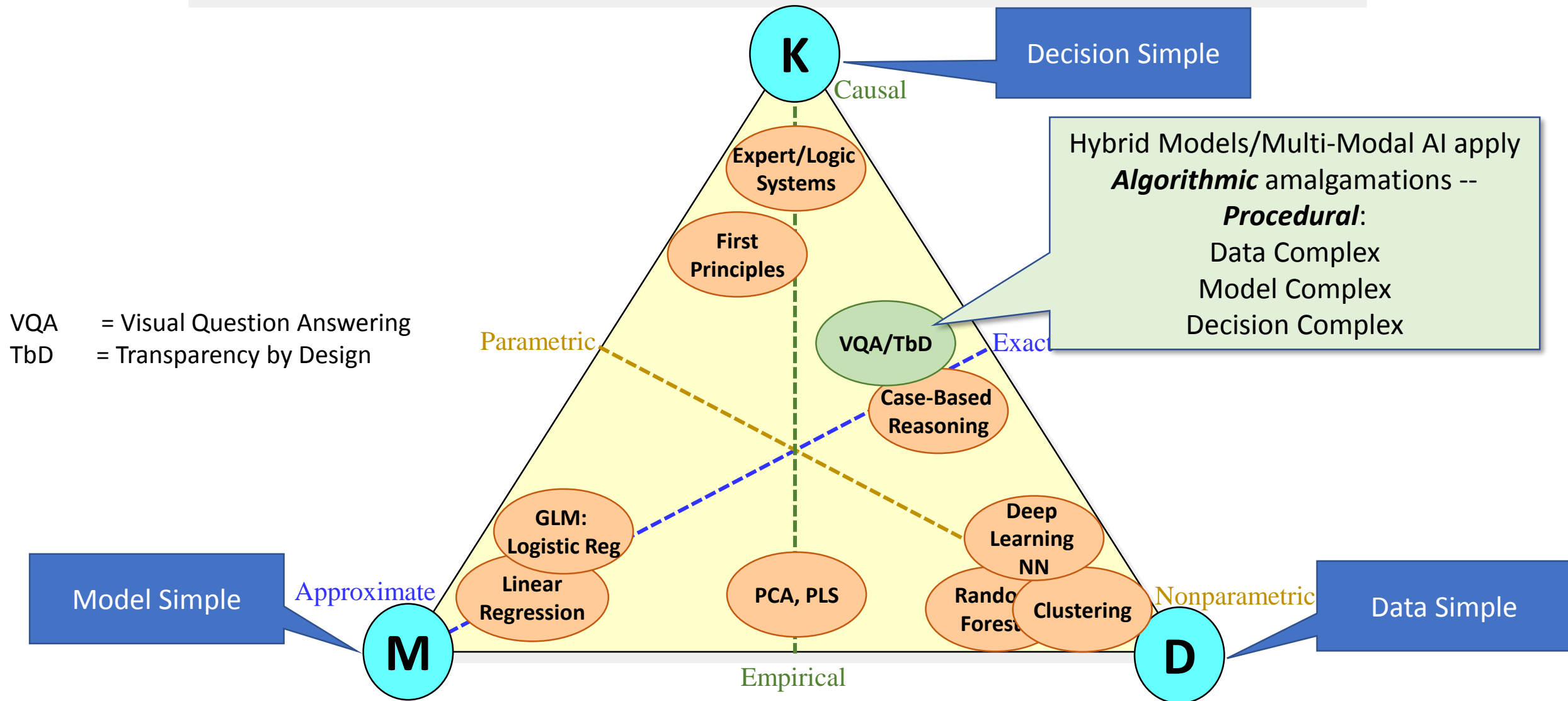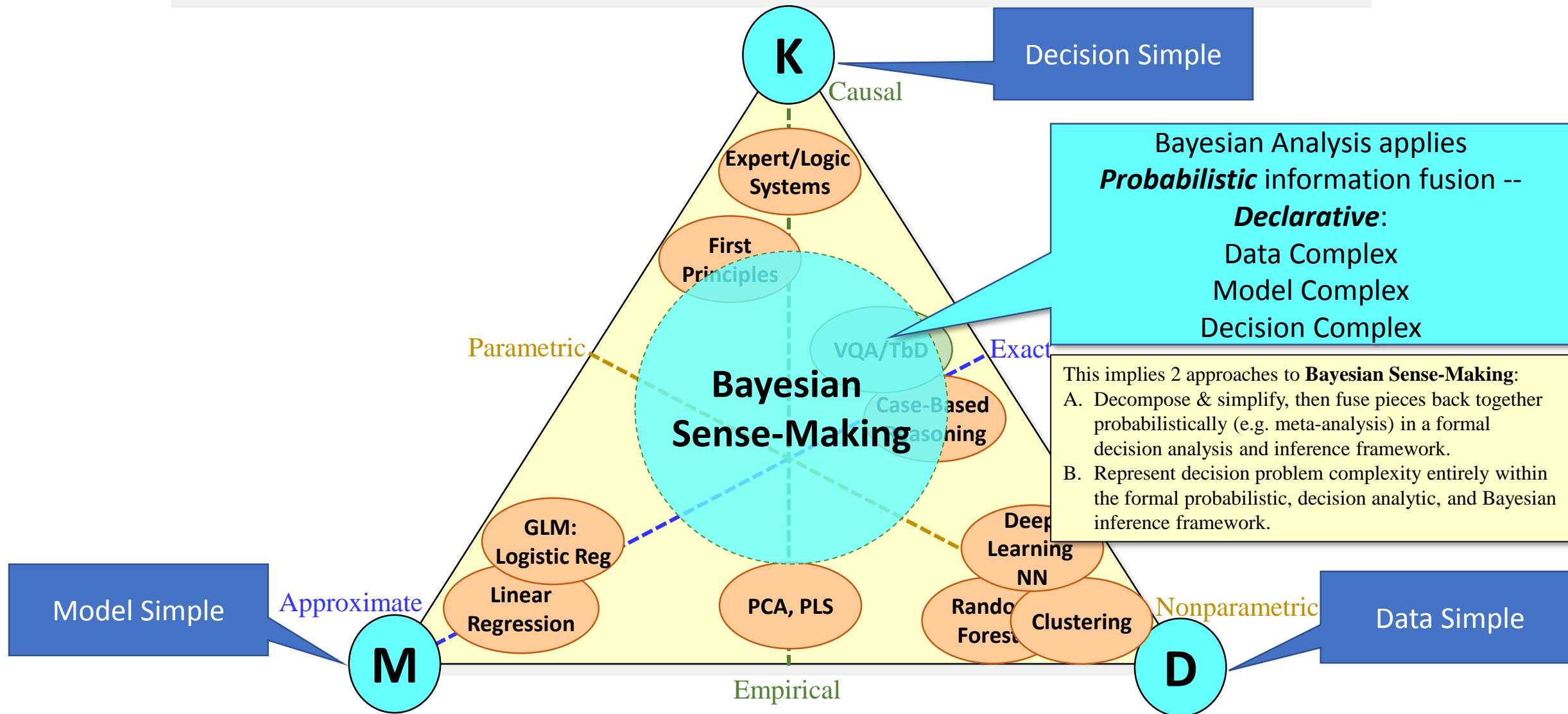| "Simple" | "Complex" |
|---|---|
| **Data** <br> • Single source <br> • Single variable types/distribution families <br> • Tabular & Ample <br>   • Non-missing <br>   • Regular, exchangeable    **Homogeneous** | **Data** <br> • Multiple sources <br> • Multiple variable types/distribution families <br> • Ragged & Sparse <br>   • Missing <br>   • Multigranular aggregation    **Heterogeneous** |
| **Model** <br> • Observations linked to observations (Modeling the Data) <br> • Empirical structure <br>   • Single-level <br>   • Acausal    **Data-to-Data** <br> • Single hypothesis <br> • Component-level estimation; Low-level integration | **Model** <br> • Latent spaces (Modeling the Domain) <br> • Causal structure <br>   • Multi-level <br>   • Mechanisms    **True-to-True** <br> • Mixture phenomena/Multi-Hypothesis <br> • System-level integration |
| **Decision** <br>   **Deterministic, Predictions** <br> • Deterministic assumptions <br> • Modal/point estimate solutions <br> • Predictive inference (What will happen?) <br> • Single objective, Static | **Decision** <br>   **Probabilistic, Explanations** <br> • Reasoning under uncertainty (UQ) <br> • Risk analysis <br> • Explanatory inference (Why did it happen?) <br> • Multi-Objective, Dynamic updating |

12

**Hybrid Models**
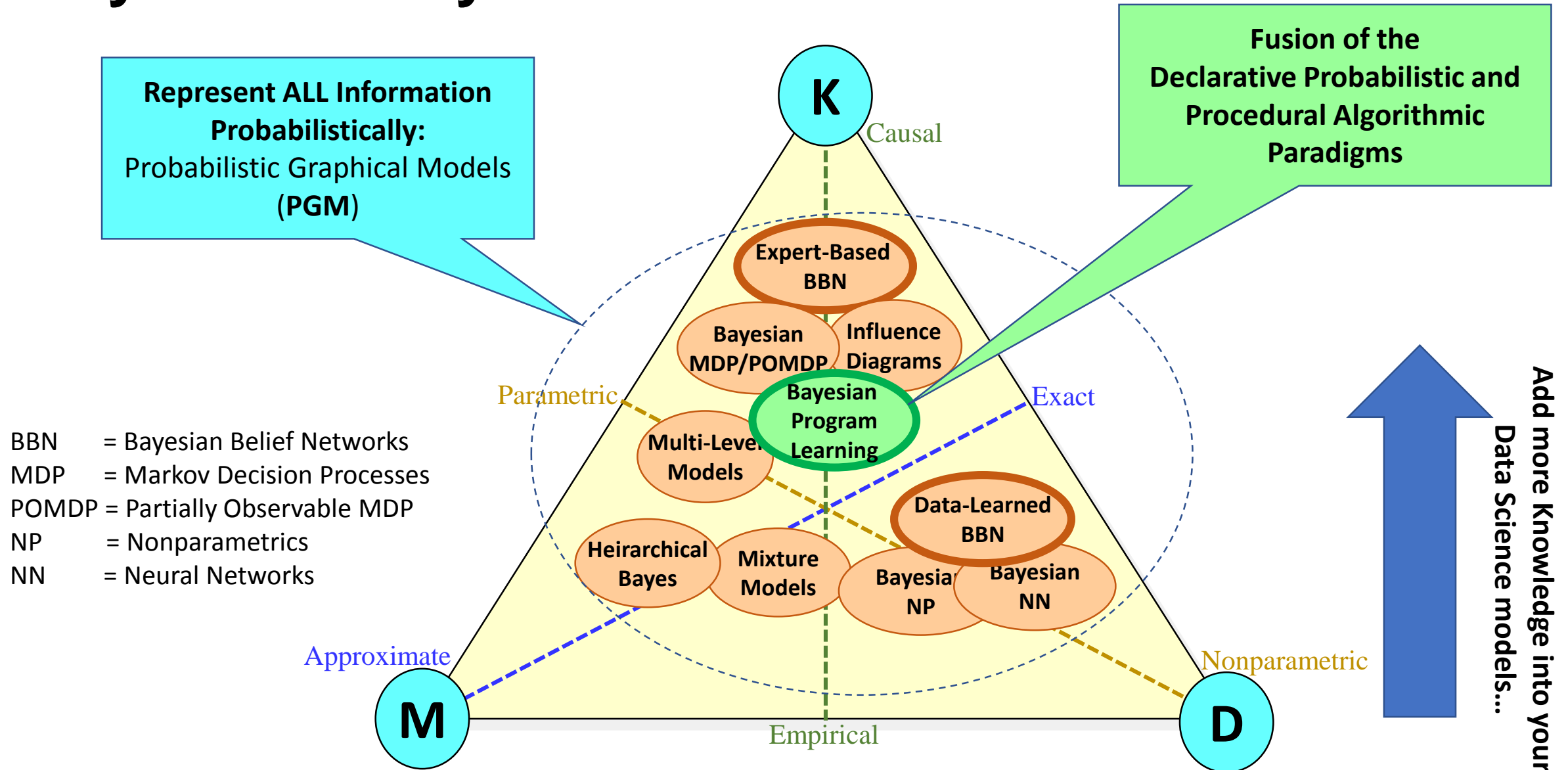Combine Components to Make Sense in the Face of Real-World Complexity

K

Decision Simple

Causal

Expert/Logic Systems

First Principles

Hybrid Models/Multi-Modal AI apply
*Algorithmic* amalgamations --
*Procedural*:
Data Complex
Model Complex
Decision Complex

VQA = Visual Question Answering
TbD = Transparency by Design

Parametric

VQA/TbD

Exact

Case-Based Reasoning

GLM: Logistic Reg

Model Simple

Approximate

Linear Regression

PCA, PLS

Deep Learning NN

Random Forest

Clustering

Nonparametric

Data Simple

M

Empirical

D

13

# Bayesian Analysis Formalizes Sense-Making



**K** — Causal

Decision Simple

Expert/Logic Systems

First Principles

Parametric

VQA/TbD

Exact

**Bayesian Sense-Making**

Case-Based Reasoning

Bayesian Analysis applies *Probabilistic* information fusion --
*Declarative*:
Data Complex
Model Complex
Decision Complex

This implies 2 approaches to **Bayesian Sense-Making**:
A. Decompose & simplify, then fuse pieces back together probabilistically (e.g. meta-analysis) in a formal decision analysis and inference framework.
B. Represent decision problem complexity entirely within the formal probabilistic, decision analytic, and Bayesian inference framework.

GLM: Logistic Reg

Deep Learning NN

Model Simple

Approximate

Linear Regression

PCA, PLS

Random Forest

Clustering

Nonparametric

Data Simple

**M**

**D**

Empirical

# Bayesian Analysis

Use Case: Background
A daughter picks which colleges to visit & apply to...

# Bayesian Sense-Making: Key Concepts

**I.   Domain-Relevant Model Structure**

- Generative Probabilistic Graphical Models (PGM)
- Latent Spaces, Mixture & Multi-Level Models Causal Structural Models

**II.   Information Theoretic Principles**

- Informative but Least Committal Probability Distributions
- Probability-Based Metrics for Association, Goodness, and Discrepancy

**III.   Bayesian Inference within Probabilistic Programming Languages**

- Model-Based Machine Learning
- Declarative Probabilistic Programs

**IV.   Explanatory and Causal Inference**

- Most Relevant Explanations
- Simulation & Implications of Interventions &Counterfactuals

**V.   Risk Analysis and Decision Analysis**

- Uncertainty Quantification (UQ) & Optimization – Maximum Expected Utility
- Value of Information

**VI.   Optimal Learning**

- Sequential & Adaptive Design of Experiments: Optimal Exploration & Exploitation

# I. Domain-Relevant Structure
## Generative Probabilistic Graphical Models (PGM)

- Declarative specification of data generation process
  - Exploit **Conditional Independence**
  - **Probabilistic Programming Languages**
  - **Model-Based Machine Learning**

- Explicitly represent Latent Spaces
  - **Model the System, *NOT* the Data**



$$P(Responses \mid LatentSpace, Stimuli, \theta_m)$$

$P(\theta_m)$

$P(\theta)$

$P(LatentSpace \mid Stimuli, \theta)$

$N$

The Latent Space is the link that fuses together observations from many different contexts.

$$P(Responses, LatentSpace, \Theta \mid Stimuli)$$

$$= P(\Theta) \prod^{N} P(Responses \mid LatentSpace, Stimuli, \Theta) \, P(LatentSpace \mid Stimuli, \Theta)$$

$$\rightarrow \prod_{Contexts} P(Responses = D \mid LatentSpace = Z, Stimuli) \prod_j P(LatentSpace = Z_j \mid Par(Z_j), Stimuli)$$

Measurement Models over multiple Contexts        Causal Structural Model

# Knowledge Elicitation

## Capturing and representing domain knowledge

**BEKEE, *Bayesia* Expert Knowledge Elicitation Environment**

Seminar: Knowledge Elicitation & Reasoning with Bayesian Networks (video)



Source: Bayesia S.A.S

# II. Information Theoretic Concepts
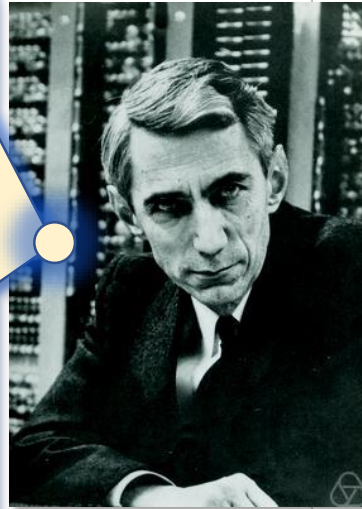## Basis for prior distributions & discrepancy/association metrics

- Basics
  - Surprisal, $S(x) = \log(1/P(x))$
  - Entropy , $H(x) = \Sigma_x P(x)S(x)$
  - Information, $I(x|y) = H(x)-H(x|y)$

> "My greatest concern was what to call it. **I thought of calling it 'information,' but the word was overly used, so I decided to call it 'uncertainty.'** When I discussed it with John von Neumann, he had a better idea. <mark>Von Neumann</mark> **told me, 'You should call it entropy, for two reasons.** In the first place your uncertainty function has been used in statistical mechanics under that name, so it already has a name. **In the second place, and more important, no one really knows what entropy really is, so in a debate you will always have the advantage**."
> <mark>Claude E. Shannon</mark>,
> *Scientific American,* 1971, v225, p180

- **Distribution Derivations**
  - **MaxEnt**: Given moments, quantiles, and/or bounds, derive the probability distribution that satisfies these constraints while admitting no other information.

## Fitness Function Metrics
  - **MDL(p(x,D,$\Theta$))**: Measure of information content of a probabilistic model p(x,D,$\Theta$)
  - **KLD(p||q)**: Measure of discrepancy between a probability distribution p and a reference distribution q.

## Association Metrics
  - **I(X,Y)**: Mutual Information is the KLD of the true joint probability distribution $P(X,Y)$ from the joint under independence $P(X)P(Y)$

# III. Bayesian Inference in Probabilistic Programming Languages
## Declarative Probabilistic Specification Distinct from Inference Algorithms

**Model-Based Machine Learning**
E.g., Microsoft's C. Bishop ([PDF](#))
Environment:
      **BayesiaLab** by *Bayesia*: solely
Bayesian Belief Networks & Influence Diagrams
**Probabilistic Programming Languages:**
(see https://github.com/topics/bayesian-inference )
**Stan** (esp. R), PyMC3 (Python)
Google: TensorFlow Probability;
Uber AI: Pyro
Microsoft: Infer.NET

- Natively encode probability distributions
- Syntax for conditioning upon evidence
- Make available a variety of inference algorithms for any model: e.g. Hamiltonian Monte *Carlo-No-U-Turn* Sampling (HMC-NUTS); Automatic Differentiation Variational Inference (ADVI); and robust optimizers

Stan

```
1   // Bayesian Plackett
2   data {
3     int N;
4     int M;
5     int<lower=0> K[N];
6     int Kmax;
7     int D;
8     int Nsource;
9     int Ncountry;
10    int w[N,Kmax];
11    row_vector[D] X[M]
12    int country[M];
13    int source[N];
14    // HYPERPARAMETERS
15    vector[D] mu0;
16    cov_matrix[D] Sign
17  }
18  parameters {
19    vector[D] beta;
20    vector[M] etaraw;
21    vector[Ncountry] b
22    vector[Nsource] bs
23    real<lower=0> Smag
    strengths)
24  }
25  transformed paramete
26    vector[M] eta;
27    vector[Ncountry] b
28    vector[Nsource] bs
29    vector[M] strength
30    real etamean;
31    real<lower=0> etas
32    etamean = mean(etaraw);
33    etastdv =   sd(etaraw);
34    eta   = (etaraw   - etamean)/etastdv;
35    bCntry = (bCntryraw - etamean)/etastdv;
36    bSrc  = (bSrcraw  - etamean)/etastdv;
37    for (j in 1:M) {
38      strength[j] = Smag * ( eta[j] + X[j] * beta + bCntry[country[j]] );
39    }
40  }]
```

```
41  model {
42    vector[M] pinstitution;  // conditional probability of
43    vector[M] vi;             // strength of each institutio
44    // Priors for the coefficients/random effects.
45    beta      ~ multi_normal(mu0,Sigma);
46    Smag      ~ exponential(1);
47    etaraw    ~ normal(0,1);
48    bCntryraw ~ normal(0,1);
49    bSrcraw   ~ normal(0,1);
50    // Compute conditional probabilities Bayes of ranks for each
51    for ( i in 1:N ) { // for each of N competitions/ranki
52      // Strength of each institution adjusted for source.
53      vi = exp( strength + Smag * bSrc[source[i]]  ); // (
54      for ( r in 1:K[i] ) { // for each position, i.e. ran
55        for ( j in 1:M ) {
56          pinstitution[j] = 0.0;
57        }
58        for ( rj in r:K[i] ) { // for institutions with ra
59          pinstitution[w[i,rj]] = vi[w[i,rj]];
60        }
61        pinstitution = pinstitution / sum(pinstitution);
62        // Likelihood based on the ordering data: Sample f
63        w[i,r] ~ categorical(pinstitution);
64      }
65    }
66  }
67
```

Source: "Bayesian Plackett-Luce Rankings Model", M.L.Thompson, Kaggle.com kernel, 2016, Apache 2.0 license

# IV. Explanatory and Causal Inference
## Deriving Insights & Reliable Policies by Explaining Why

- Most Relevant (Representative) Explanations: **Generalized Bayes Factor, GBF(H;E)**
  - *Which hypothesis, H, best explains given evidence, E?*

- Implications of Interventions and Decision Policies: Causal inference

- $\text{GBF}(H; E) = \dfrac{\text{P}(Evidence=E|Hypothesis=H)}{\text{P}(Evidence=E|Hypothesis\neq H)}$

  $= \dfrac{\text{Odds}(Hypothesis=H|Evidence=E)}{\text{Odds}(Hypothesis=H)}$

- Weight of Evidence, $\text{WE}(H; E) \triangleq \log \dfrac{\text{P}(E|=H)}{\text{P}(E|\neq H)}$

where

$$\text{Odds}(X = x) \equiv \frac{\text{P}(X=x)}{\text{P}(X\neq x)} = \frac{\text{P}(X=x)}{1-\text{P}(X=x)}$$

Yuan, C., et al., Most relevant explanations in Bayesian networks, *J. AI Research*, 2011

Good, I.J., Weight of Evidence: A Brief Survey, Bayesian Statistics 2, 1985

# IV. Explanatory and Causal Inference
## Deriving Insights & Reliable Policies by Explaining Why

"**It is therefore natural to call it *the factor in favour of H provided by E* and this was the name given to it by** A.M. Turing **in a vital cryptanalytic application in WWII in 1941.** He did not mention Bayes's theorem, with which it is of course closely related, because he always liked to work out everything for himself. **When I said to him that the concept was essentially an application of Bayes's theorem he said 'I suppose so'.**

... **Thus *weight of evidence* is equal to the logarithm of the Bayes factor.**"

"Weight of Evidence: A Brief Survey",
**Good, I.J.,** *Bayesian Statistics 2;*
Bernardo, et al. (eds), 1985

- $\text{GBF}(H;E) = \dfrac{\text{P}(Evidence=E|Hypothesis=H)}{\text{P}(Evidence=E|Hypothesis\neq H)}$

$\qquad = \dfrac{\text{Odds}(Hypothesis=H|Evidence=E)}{\text{Odds}(Hypothesis=H)}$

- Weight of Evidence, $\text{WE}(H;E) \triangleq \log \dfrac{\text{P}(E|=H)}{\text{P}(E|\neq H)}$

"**...**the terminology of Bayes factors and weights of evidence has more intuitive appeal [than log-likelihood ratio]. This intuitive appeal persists in the general case when the weight of evidence is not the logarithm of a likelihood ratio.
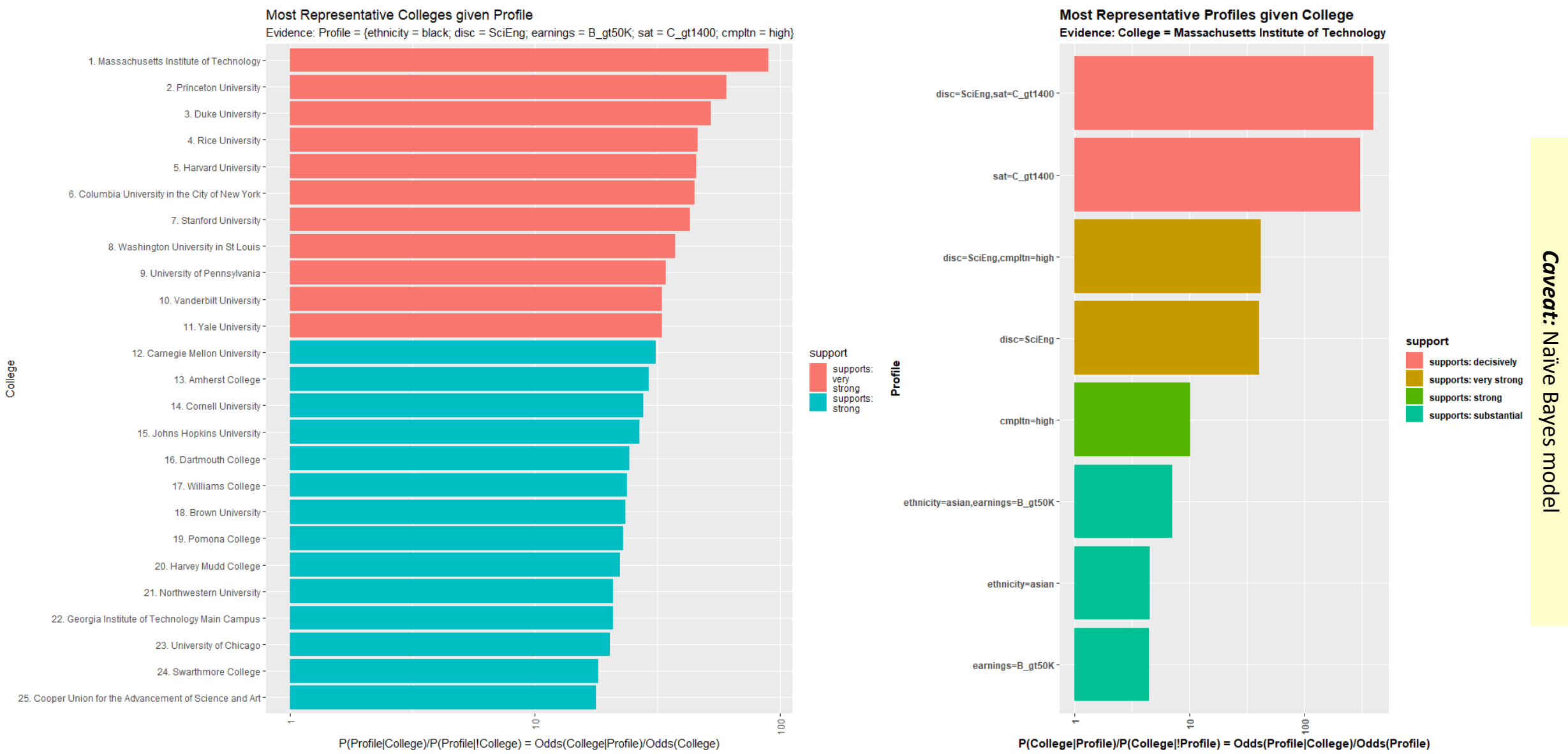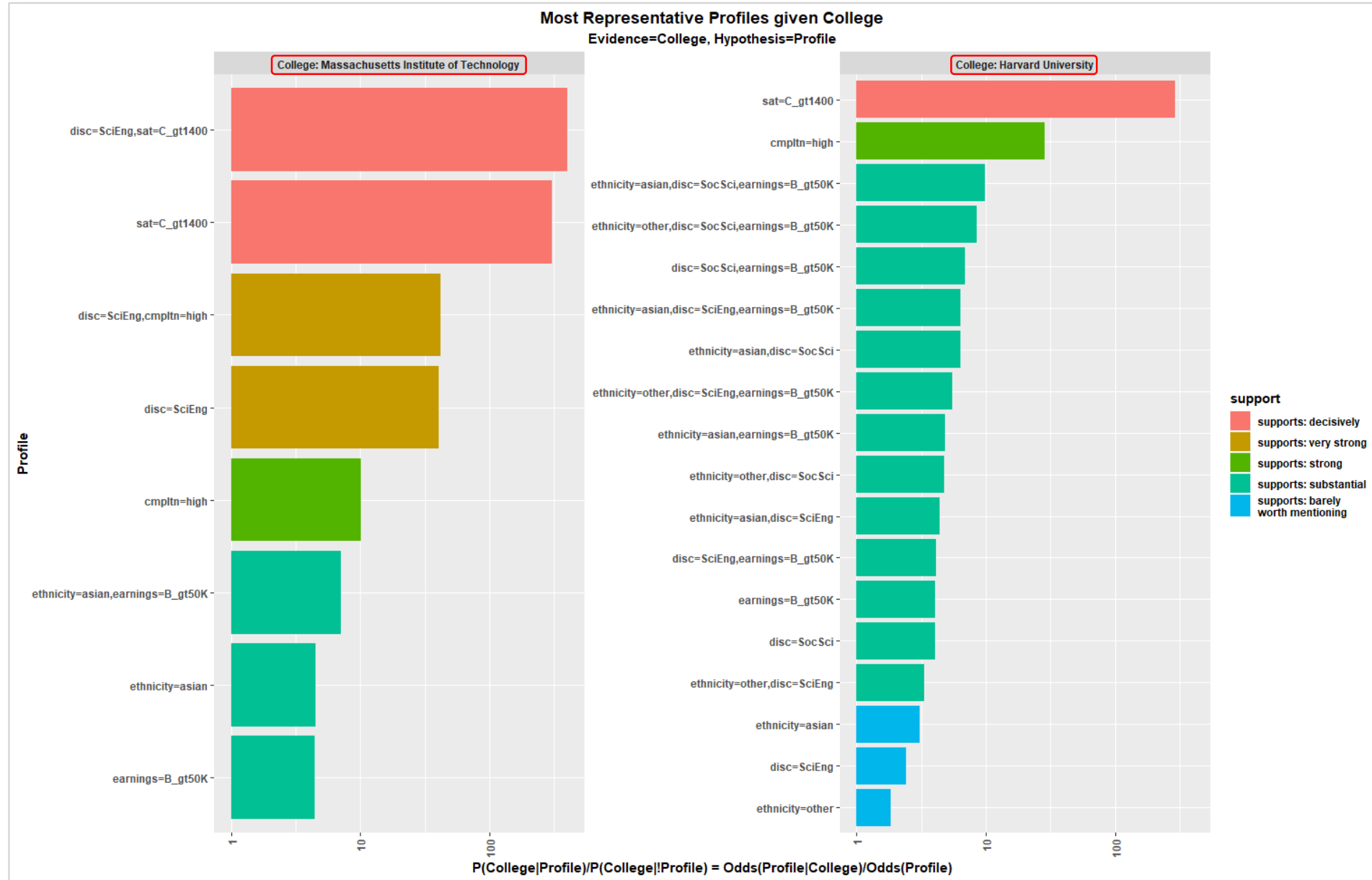**I conjecture that juries, detectives, doctors, and perhaps most educated citizens, will eventually express their judgments in these intuitive terms.",** *ibid*

Yuan, C., et al., Most relevant explanations in Bayesian networks, *J. AI Research*, 2011
Good, I.J., Weight of Evidence: A Brief Survey, Bayesian Statistics 2, 1985

# Ranking Colleges as hypotheses given Student Profiles as evidence, & vice versa



**Most Representative Colleges given Profile**
Evidence: Profile = {ethnicity = black; disc = SciEng; earnings = B_gt50K; sat = C_gt1400; cmpltn = high}

**Most Representative Profiles given College**
Evidence: College = Massachusetts Institute of Technology

*Caveat:* Naïve Bayes model

P(Profile|College)/P(Profile|!College) = Odds(College|Profile)/Odds(College)

P(College|Profile)/P(College|!Profile) = Odds(Profile|College)/Odds(Profile)

# Contrasting Colleges as evidence scenarios



*Caveat:* Naïve Bayes model

# Causal Inference: Climbing Pearl's Causal Ladder

- Motivates imposing a
  ***Causal Structural Model of the Latent Space***

- Predictive ***Simulations***: Implications of Interventions/Decision Policies

- Pearl, J. and Mackenzie, D., *The Book of Why: The New Science of Cause and Effect*, 2018
  - Downloadable Chapter 1

- **DAGitty** [http://www.dagitty.net/] ...
  - "... is a browser-based environment for creating, editing, and analyzing causal models (also known as directed acyclic graphs or causal Bayesian networks). **The focus is on the use of causal diagrams for minimizing bias in empirical studies** in epidemiology and other disciplines."
  - Developed & maintained by Johannes Textor (Tumor Immunology Lab and Institute for Computing and Information Sciences, Radboud University Nijmegen)
  - Textor, J., et al., "Robust causal inference using directed acyclic graphs: the R package 'dagitty'", *Intl. J. Epidemiology*, 45, 6, 1 Dec. 2016, 1887–1894

Leads to plausible reasoning about a person's underlying motives ….
Hence, we go beyond measured data and into latent constructs.



**3. COUNTERFACTUALS**

ACTIVITY: Imagining, Retrospection, Understanding

QUESTIONS: *What if I had done …? Why?*
(Was it X that caused Y? What if X had not occurred? What if I had acted differently?)

EXAMPLES: Was it the aspirin that stopped my headache? Would Kennedy be alive if Oswald had not killed him? What if I had not smoked for the last 2 years?

**2. INTERVENTION**

ACTIVITY: Doing, Intervening

QUESTIONS: *What if I do …? How?*
(What would Y be if I do X? How can I make Y happen?)

EXAMPLES: If I take aspirin, will my headache be cured? What if we ban cigarettes?

**1. ASSOCIATION**

ACTIVITY: Seeing, Observing

QUESTIONS: *What if I see …?*
(How are the variables related? How would seeing X change my belief in Y?)

EXAMPLES: What does a symptom tell me about a disease? What does a survey tell us about the election results?

IMAGINING

DOING

SEEING

JUDEA PEARL
WINNER OF THE TURING AWARD
AND DANA MACKENZIE

THE BOOK OF WHY

THE NEW SCIENCE OF CAUSE AND EFFECT

# Latent Motivations of Students
## Manifest in behavioral theories & data and expressed attitudes

# V. Risk Analysis & Decision Analysis
## Quantifying the Uncertainty, Risk & Value of Decisions and Policies

- **Sensitivity Analysis**

- Uncertainty Quantification

- Optimization: Maximum Expected Utility
  - Influence Diagrams

- Learning Optimal Policies
  - Bayesian Reinforcement Learning
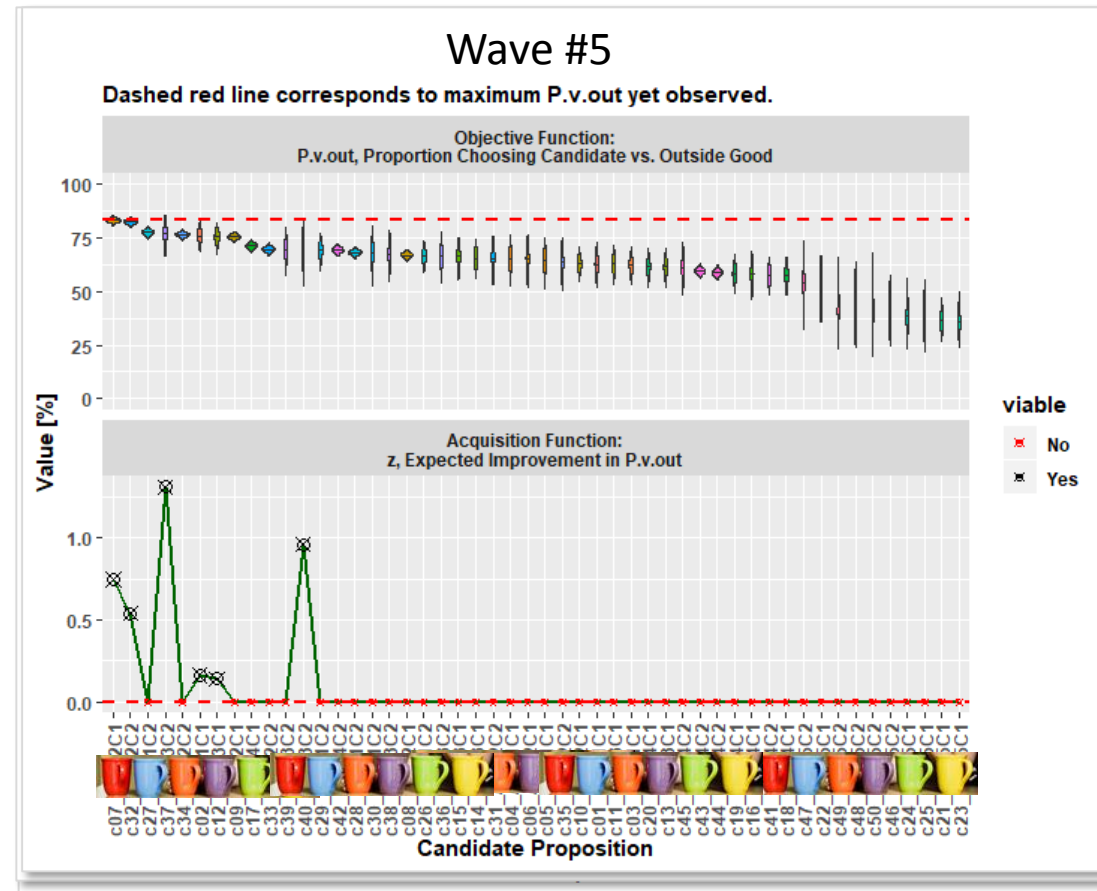    - Markov Decision Processes (MDP)
    - Partially-Observable MDP (POMDP)



29

# VI. Optimal Learning
## Trading Off Exploration & Exploitation

- **Bayesian Optimization** for Adaptive/Sequential Experimental Design and Active Learning

  - **Maximum Expected Improvement** to rank order new stimuli

Joo, Mingyu,
Thompson, Michael L.,
Allenby, Greg M.,
Optimal Product Design by Sequential Experiments in High Dimensions,
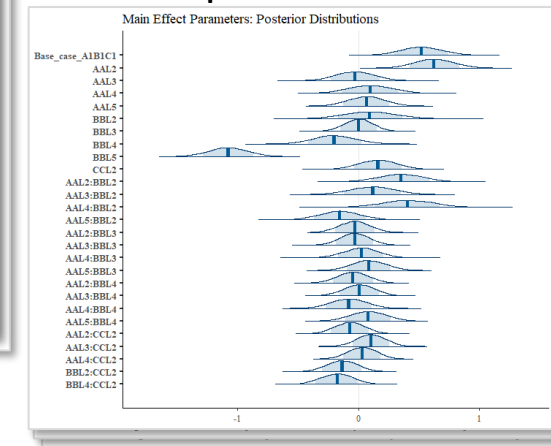Management Science (INFORMS),
Oct. 8, 2018
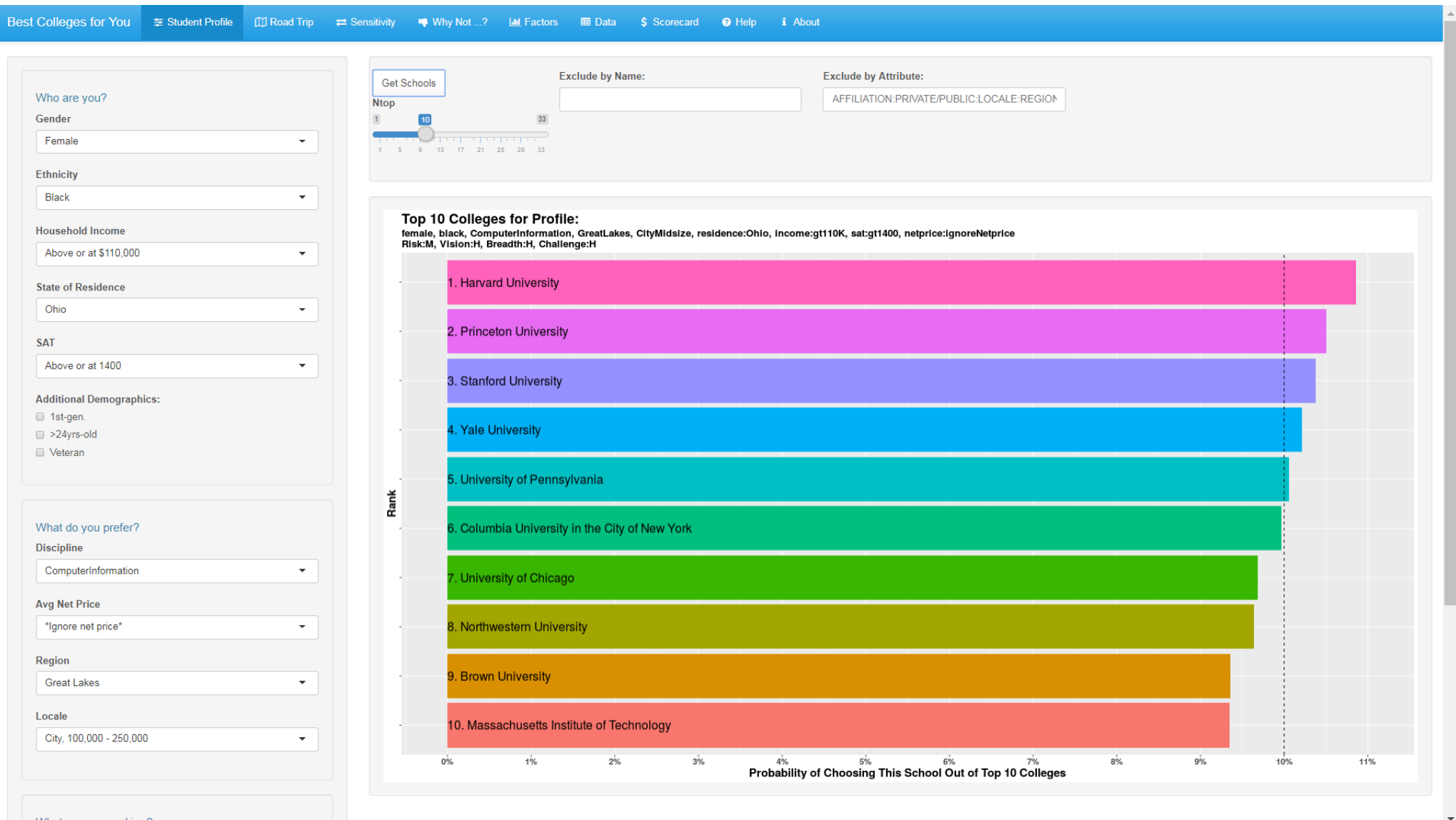
**1. Evaluate & Pick Stimuli**

**2. Perform Experiment**

**3. Update Model**

# And so, the *"Best Colleges for You"* App was born!

https://thompsonml.shinyapps.io/BestCollegeApp/
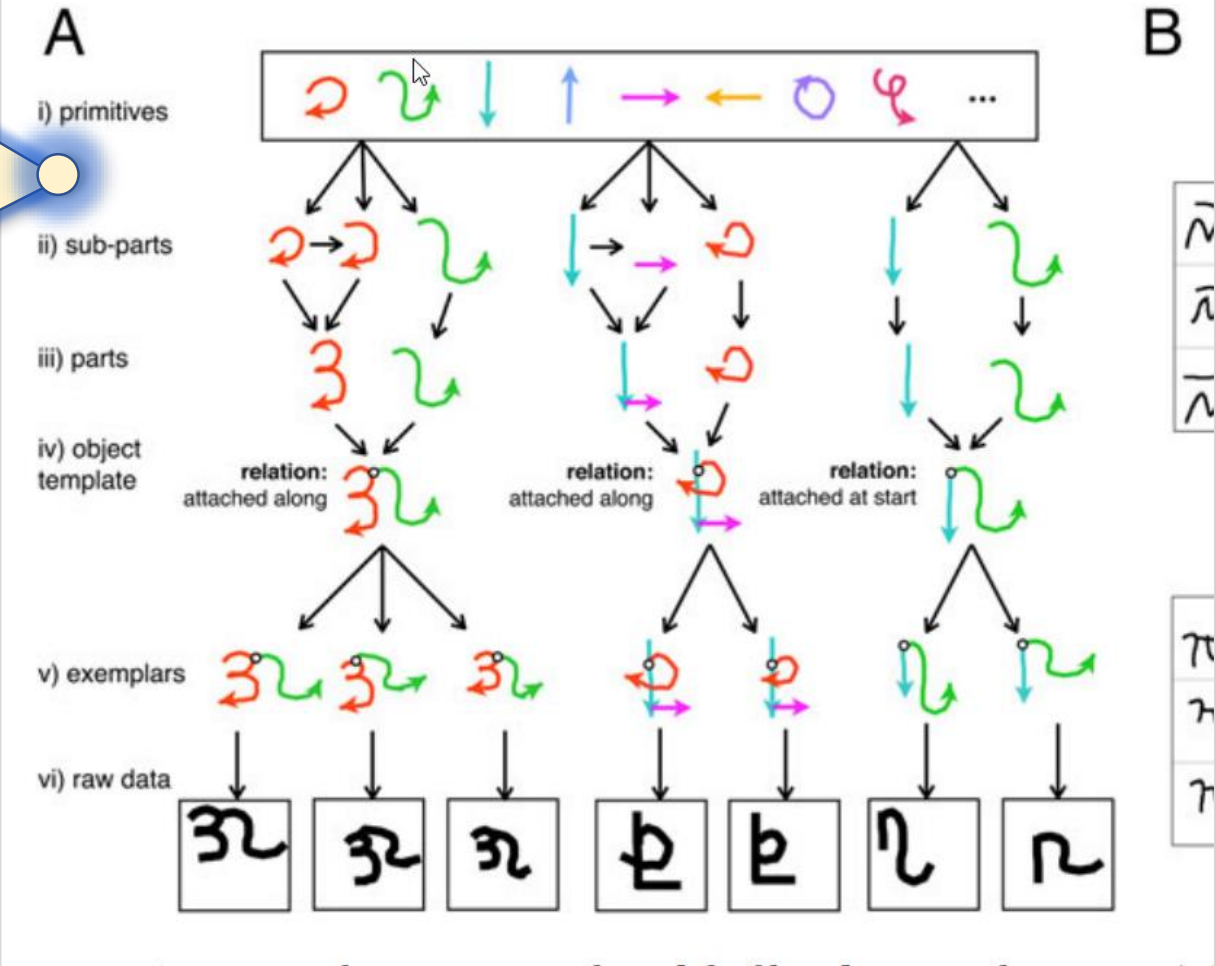
# Future Implications
# Sense-Making Systems

- **Compositionality**
- **Causality**
- **Learning-to-Learn**

> "This trend is an avenue of potential integration of deep learning models with probabilistic models and probabilistic programming: **Training neural networks to help perform probabilistic inference in a generative model or a probabilistic program."**
> "Building machines that learn and think like people",
> Lake, Brenden, et al. *Behavioral & Brain Sciences*, 40, E253. 2017

- **One-Shot Learning & General AI**
  - PGM over programs & complex schema (Brenden Lake, NYU; Josh Tenenbaum, MIT)
- **Explanatory AI**
  - DLNN as sensory apparatus fused with PGM answering "Why?" for diagnostic/advisory systems
- **Federated Learning**
  - "Computation at the Edges", e.g., Mobile Phone Deep Learning with Bayesian multi-level models



Lake et al.: Building machines that learn and think like people

# Future Implications
## Organized to Innovate

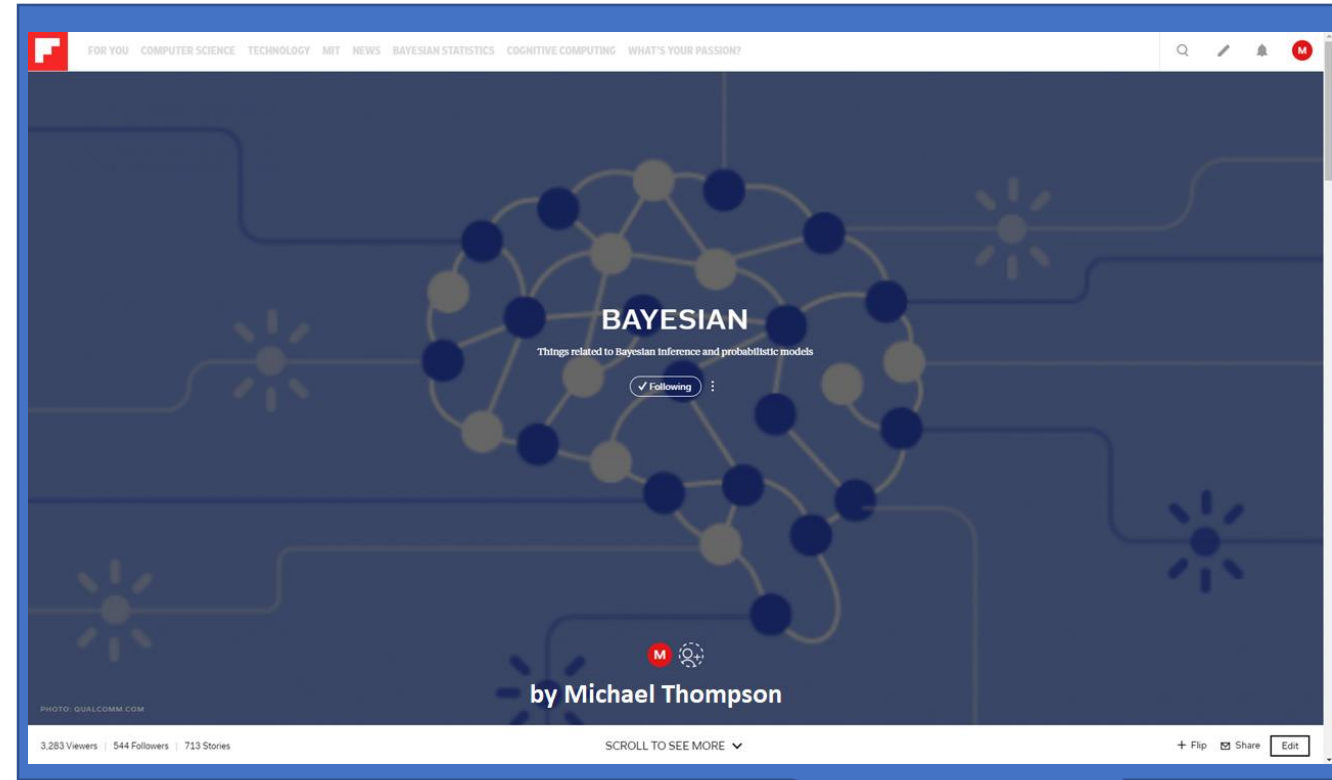- **The Procter & Gamble Company**

- **MIT Quest for Intelligence**

- **Notable Business Models**
  - Gamalon (customer intelligence)

  - Zighra (AI-powered continuous authentication): Decentralized AI through Bayesian Learning

Stay informed: Flipboard magazine "Bayesian"

# References to Get Started

**I.      Basics**
- Kurt, Will, "Count Bayesie" blog series, "A Guide to Bayesian Statistics", May 2, 2016

**II.     Domain-Relevant Model Structure**
- Tenenbaum, Joshua, et al., How to Grow a Mind: Statistics, Structure, and Abstraction (PDF), *Science*, 2011; lecture video (must-see!)
- Koller, Daphne & Friedman, Nir, Probabilistic Graphical Models, 2009, downloadable excerpt Introduction (must-read!)
- Murphy, Kevin P., Introduction to Graphical Models, 2001
- Daly, Ronan, et al., Learning Bayesian Networks: Approaches and Issues, *The Knowledge Engineering Review*, 2011
- Van Horn, Kevin S., Constructing a Logic of Plausible Inference: A Guide to Cox's Theorem, *Intl. J. Approximate Reasoning*, 2003

**III.    Information Theoretic Principles**
- Jaynes, E.T., The Relation of Bayesian and Maximum Entropy Methods, *Maximum-Entropy and Bayesian Methods in Sci. and Eng.*, 1988
- MacKay, David J.C., *Information Theory, Inference, and Learning Algorithms*, 2003; esp. Ch. 2 "Probability, Entropy, and Inference" and Ch. 3 "More About Inference".

**IV.    Bayesian Inference within Probabilistic Programming Languages**
- Bishop, Christopher, Model-Based Machine Learning, *Phil. Trans. Roy. Soc. A*, 2013
- Davidson-Pilon, Cameron, Probabilistic Programming and Bayesian Methods for Hackers (PyMC3), 2015
- Conrady, Stefan & Jouffe, Lionel (Bayesia S.A.S), Bayesian Networks and BayesiaLab: A Practical Introduction for Researchers, 2015
  - Ch. 8, Probabilistic Structural Equation Models with Bayesian Networks for Key Drivers Analysis and Product Optimization, 2015
- Stan Development Team, Modeling Language User's Guide and Reference Manual, Version 2.17.0, 2018; esp. Section III. Example Models

**V.     Explanatory and Causal Inference**
- Yuan, Changhe, et al., Most relevant explanations in Bayesian networks, *J. Artificial Intelligence Research*, 2011
- Pacer, Michael, et al., Evaluating computational models of explanation using human judgment, *Proc. 29th Conf. on Uncertainty in AI (UAI2013)*, 2013
- Rydall, Michael D. and Bramson, Aaron L., *Inference and Intervention: Causal Models for Business Analysis*, 2013
- Pearl, Judea and Mackenzie, Dana, *The Book of Why: The New Science of Cause & Effect*, 2018; downloadable excerpts Introduction, Chapter 1, Chapter 2

**VI.    Risk Analysis and Decision Analysis**
- Fenton, Norman and Neil, Martin, Managing Risk in the Real World: Applications of Bayesian Networks, 2007
- Barry Matthew and Horvitz, Eric, Vista Goes Online: Decision Analytic Systems for Real-Time Decision-Making in Mission Control, 1994

**VII.   Optimal Learning**
- Shahrari, Bobak, et al., Taking the Human Out of the Loop: A Review of Bayesian Optimization, *Proc. IEEE*, 2016
- Joo, Mingyu, et al., Optimal Product Design by Sequential Experiments in High Dimensions, *Management Science* (INFORMS), Oct. 8, 2018