

DATA INSIGHT FUNDAMENTALS

Patrick Boily

Data Action Lab | uOttawa | Idlewyld Analytics

pboily@uottawa.ca

[with files from Jen Schellinck | Sysabee]

“Reports that say that something hasn't happened are always interesting to me, because as we know, there are **known knowns**; there are things we know that we know. There are **known unknowns**; that is to say, there are things that we now know we don't know. But there are also **unknown unknowns** – there are things we do not know we don't know.”

Donald Rumsfeld, US Department of Defense News Briefing, 2002

ANALYSIS PLANNING

DATA INSIGHT FUNDAMENTALS

“Plans are nothing. Planning is everything.”

Dwight D. Eisenhower

ANALYSIS PLAN OVERVIEW

Formulate research questions/hypotheses

Identify necessary (and available) datasets

Establish inclusion/exclusion criteria for records/observations

Select variables for use in the analyses

Chose statistical methods and software

DATA 101: BASIC DATA CONCEPTS

DATA INSIGHT FUNDAMENTALS

“You can have data without
information, but you cannot
have information without data.”

Daniel Keys Moran (attributed)

WHAT IS DATA?

4,529

'red'

25.782

'Y'

OBJECTS AND ATTRIBUTES



Object: apple

Shape: spherical

Colour: red

Function: food

Location: fridge

Owner: Jen

Remember: a person or an object is not simply the sum of its attributes!

FROM ATTRIBUTES TO DATASETS

Attributes are **fields** (columns) in a database; objects are **instances** (rows).

Objects are described by their **feature vector**, the collection of attributes associated with value(s) of interest.

ID#	Shape	Colour	Function	Location	Owner
1	spherical	red	food	fridge	Jen
2	rectangle	brown	food	office	Pat
3	round	white	tell time	lounge	School
...

POISONOUS MUSHROOM DATASET

Amanita muscaria

Habitat: woods

Gill Size: narrow

Odor: none

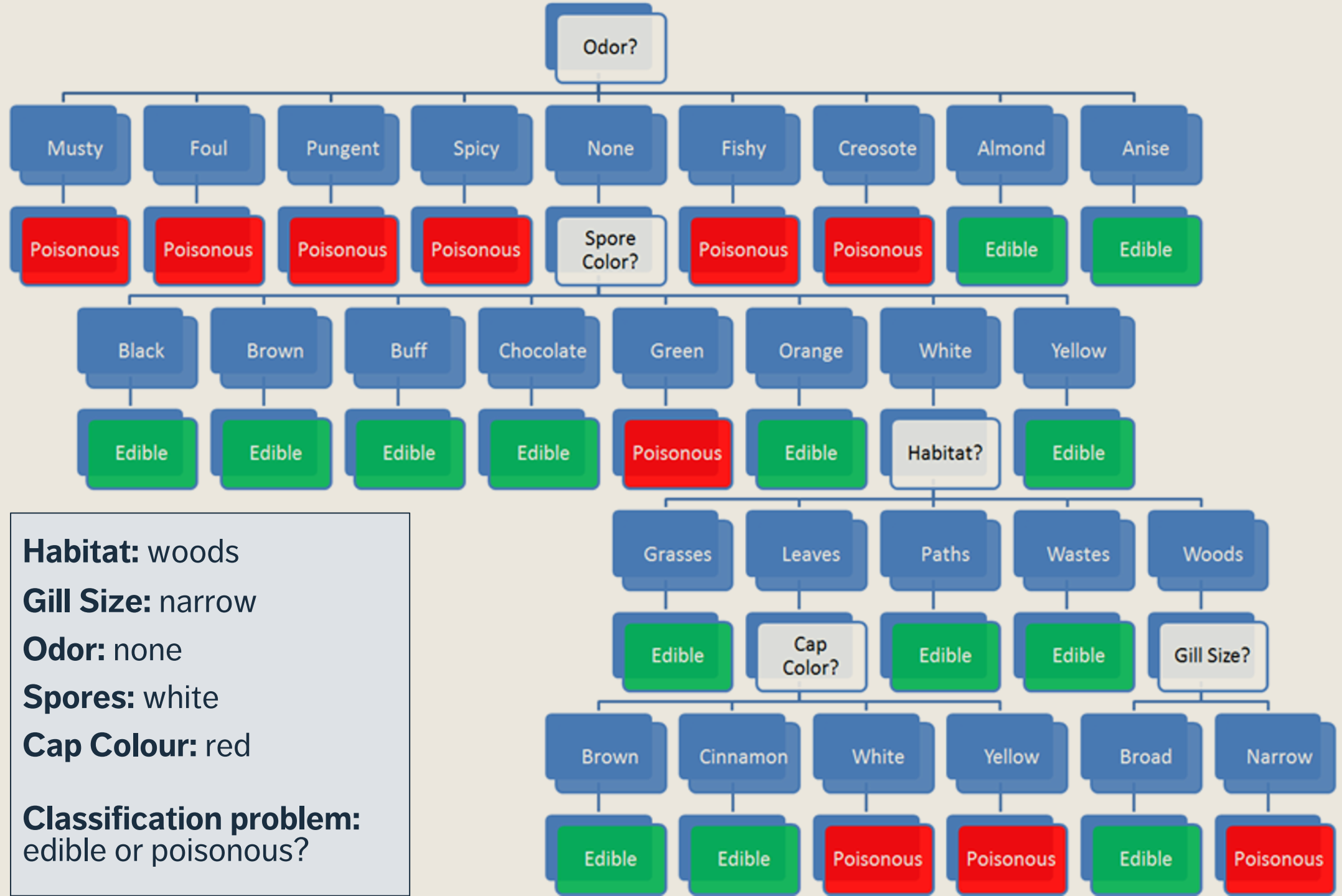
Spores: white

Cap Colour: red

Classification problem:

Is *Amanita muscaria* edible,
or poisonous?





Habitat: woods

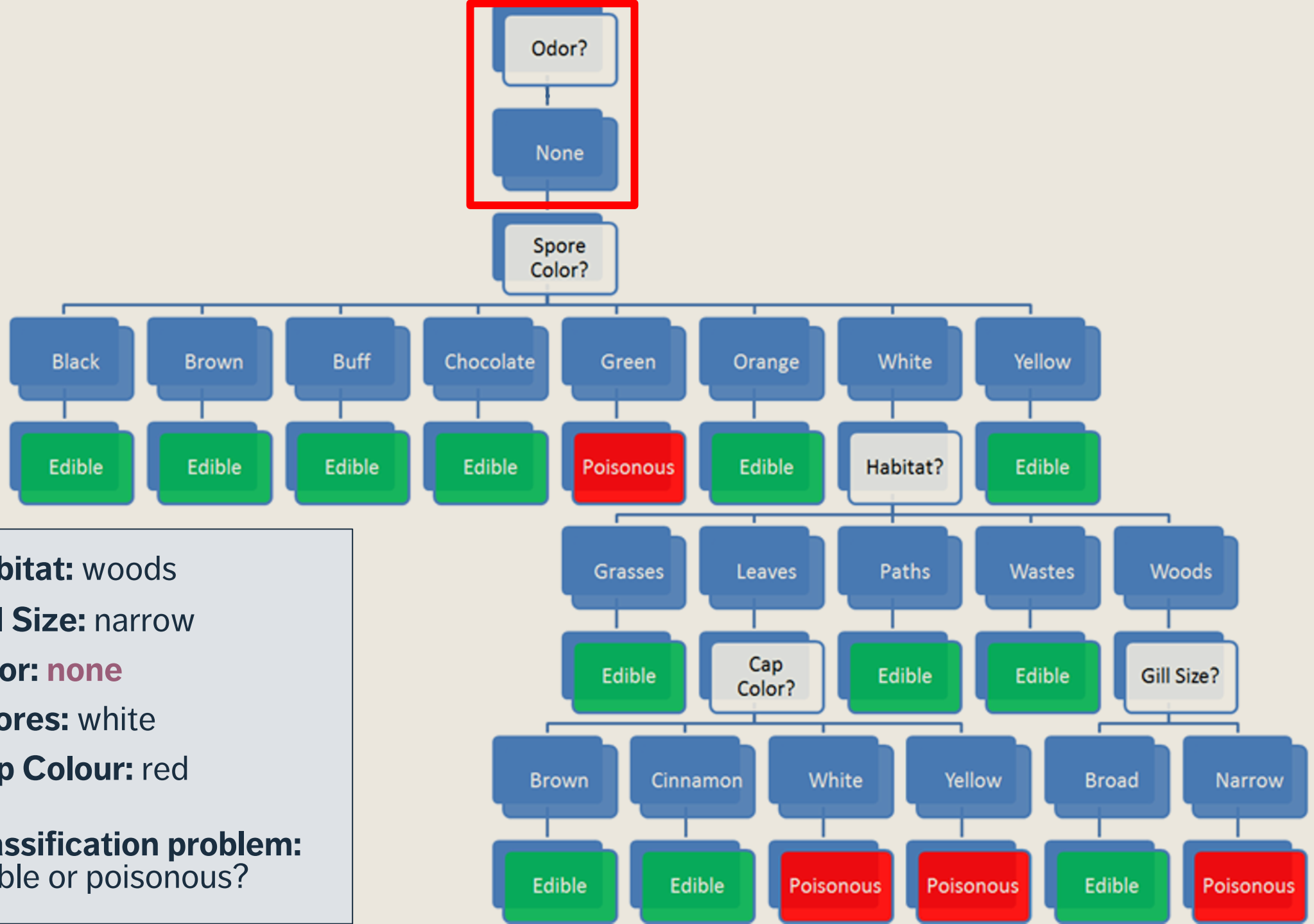
Gill Size: narrow

Odor: none

Spores: white

Cap Colour: red

Classification problem:
edible or poisonous?



Habitat: woods

Gill Size: narrow

Odor: none

Spores: white

Cap Colour: red

Classification problem:
edible or poisonous?

Habitat: woods

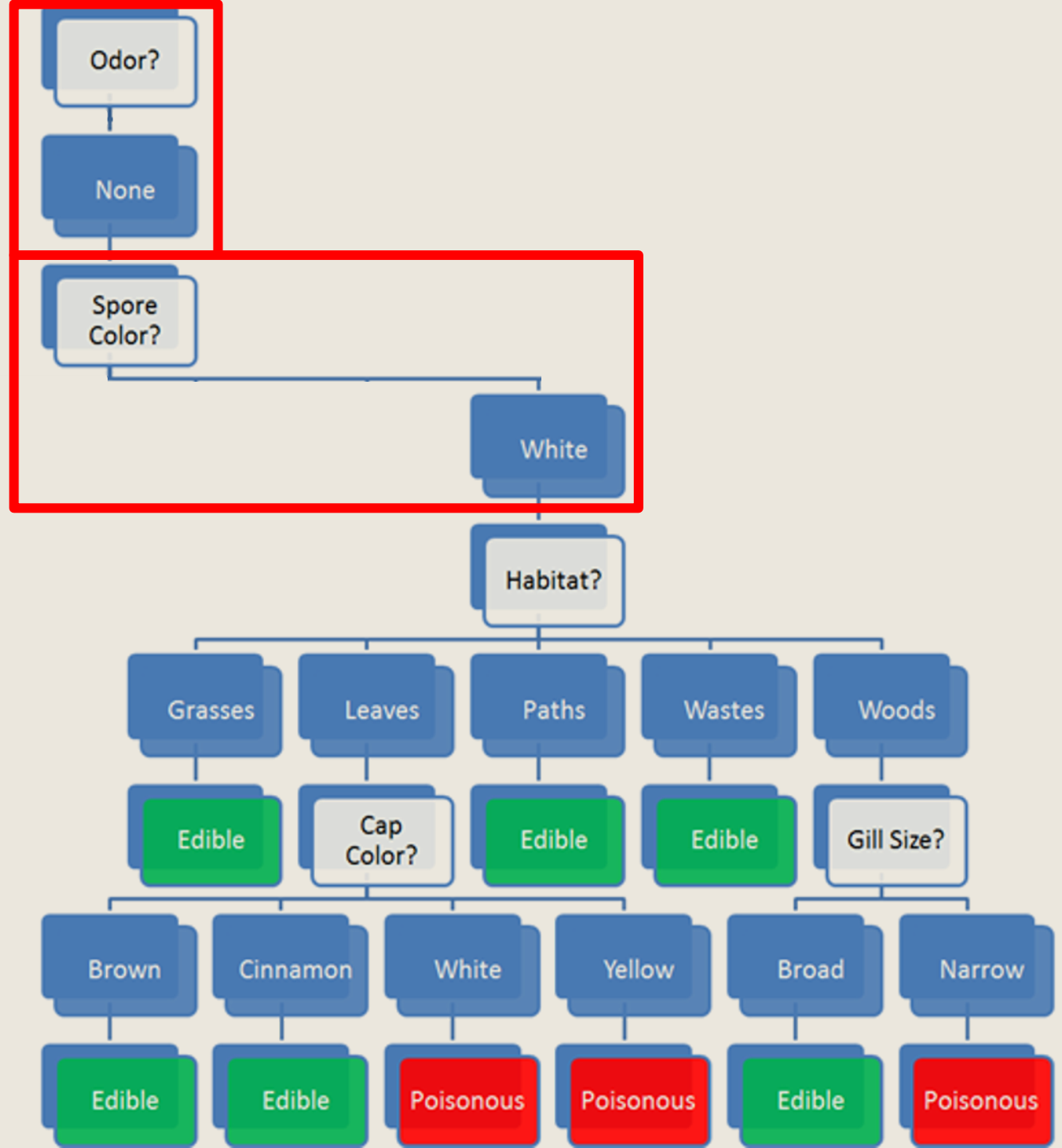
Gill Size: narrow

Odor: none

Spores: white

Cap Colour: red

Classification problem:
edible or poisonous?



Habitat: woods

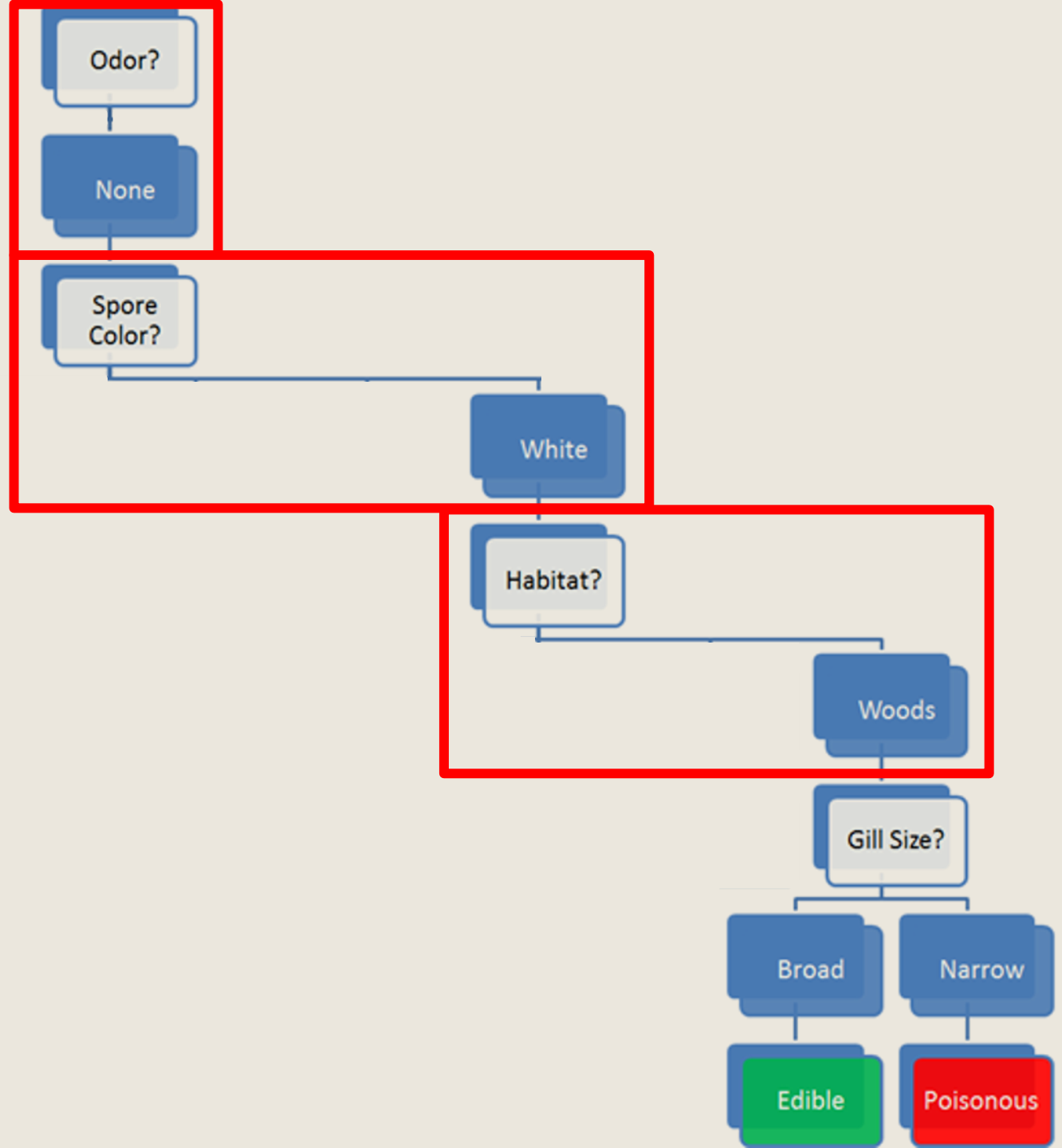
Gill Size: narrow

Odor: none

Spores: white

Cap Colour: red

Classification problem:
edible or poisonous?



Habitat: woods

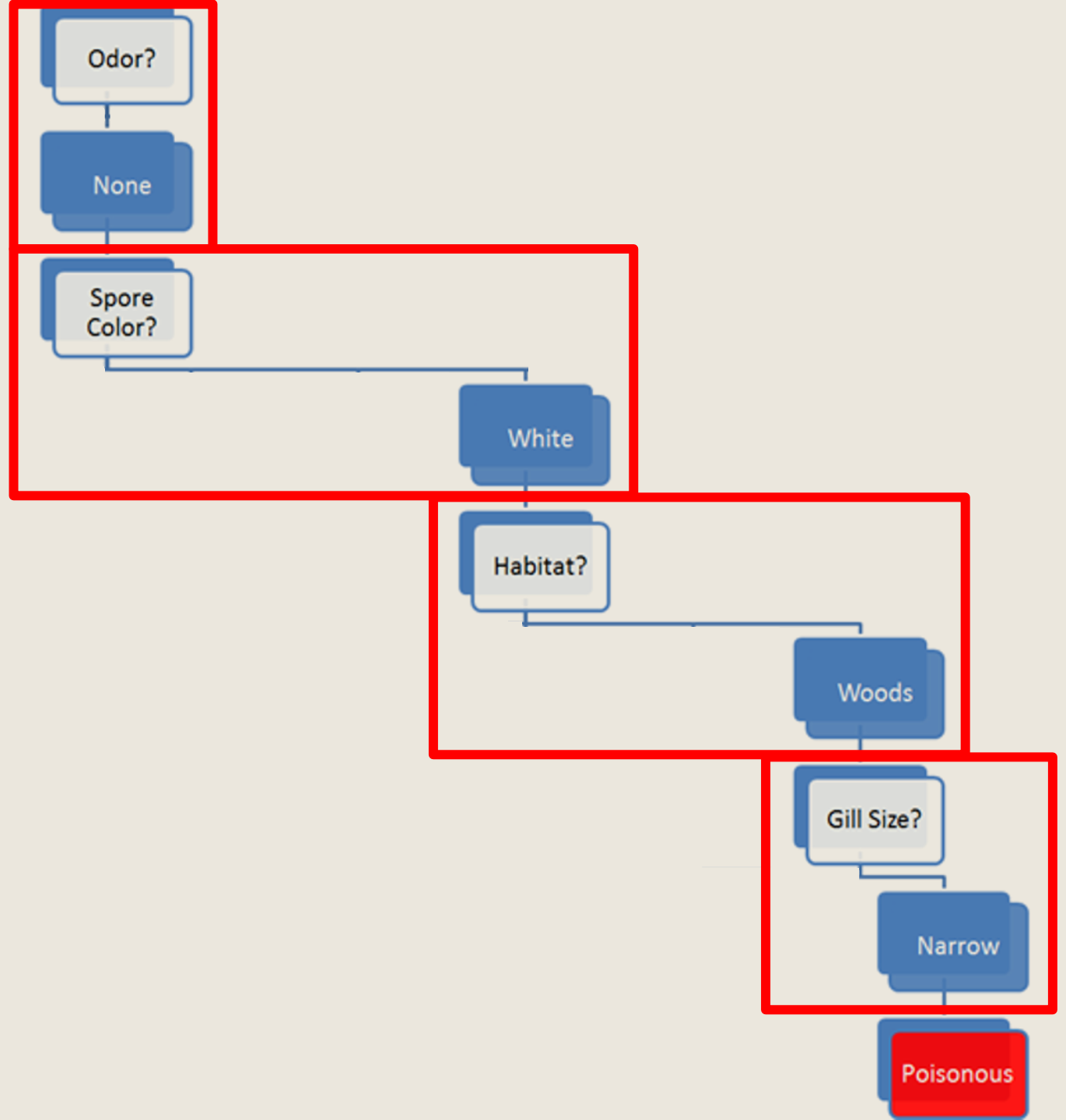
Gill Size: narrow

Odor: none

Spores: white

Cap Colour: red

Classification problem:
edible or poisonous?



Habitat: woods

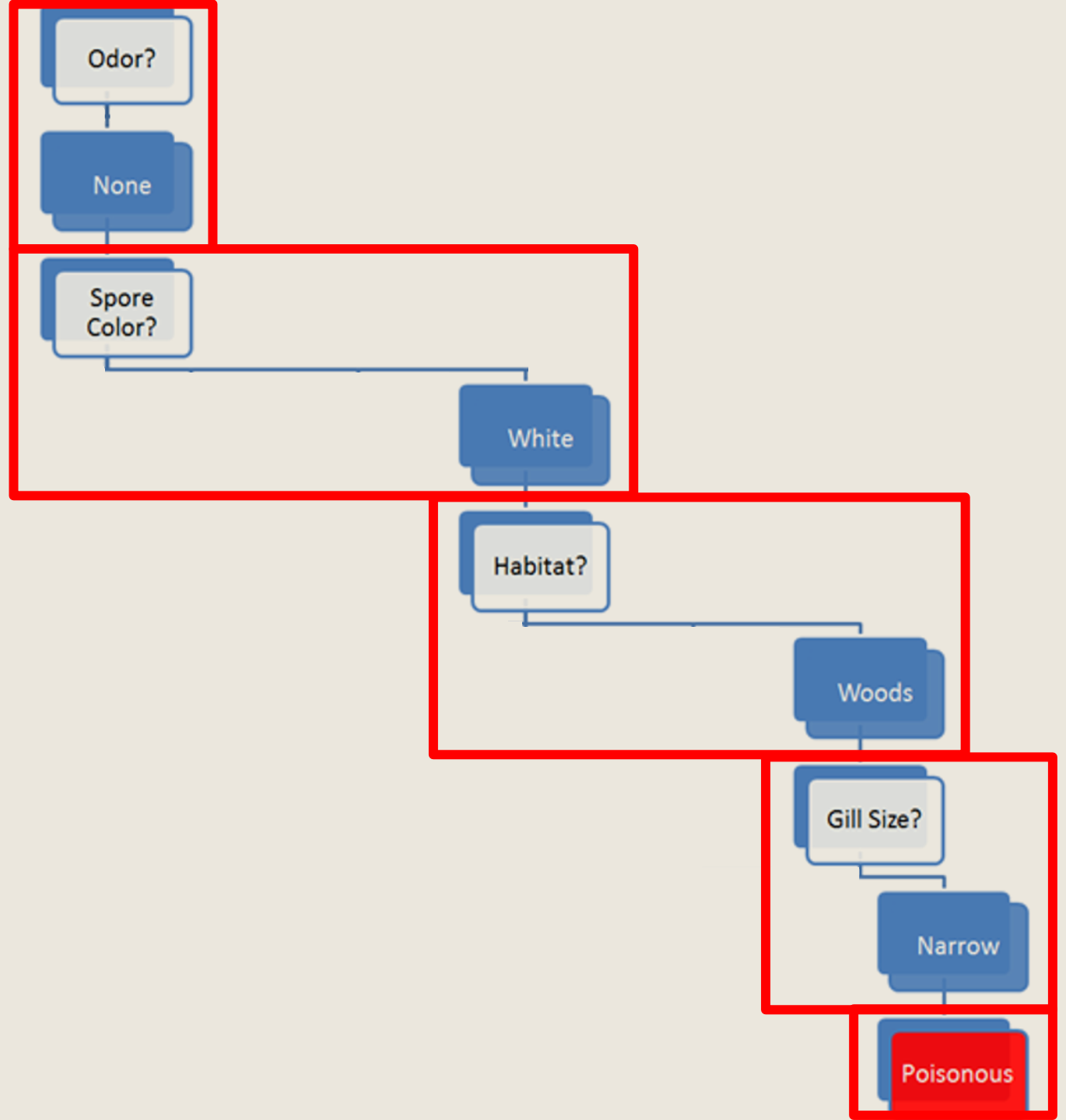
Gill Size: narrow

Odor: none

Spores: white

Cap Colour: red

Classification problem:
edible or **poisonous**



DISCUSSION

Would you have trusted an “**edible**” prediction?

Where is the model coming from?

What would you need to know to trust the model?

What’s the cost of making a classification mistake, in this case?

ASKING THE RIGHT QUESTIONS

Data science is really about asking and answering questions:

- **Analytics:** “How many clicks did this link get?”
- **Data Science:** “Based on this user’s previous purchasing history, can I predict what links they will click on the next time they access the site?”

Data mining/science models are usually **predictive** (not **explanatory**): they show connections, but don't reveal why these exist.

Warning: not every situation calls for data science, artificial intelligence, machine learning, statistics, or analytics.

THE WRONG QUESTIONS

Too often, analysts are asking the **wrong questions**:

- questions that are **too broad** or **too narrow**
- questions that **no amount of data could ever answer**
- questions for which **data cannot reasonably be obtained**

The **best-case scenario** is that stakeholders will recognize the answers as irrelevant.

The **worst-case scenario** is that they will erroneously implement policies or make decisions based on answers that have not been identified as misleading and/or useless.

DATA SCIENCE/MACHINE LEARNING/ AUGMENTED INTELLIGENCE TASKS

Classification and class probability estimation: which clients are likely to be repeat customers?

Clustering: do diplomatic missions form natural groups?

Association rule discovery: what books are commonly purchased together?

Others:

profiling and behaviour description; link prediction; value estimation (how much is a client likely to spend in a restaurant); **similarity matching** (which prospective clients are similar to a company's best clients?); **data reduction; influence/causal modeling**, etc.

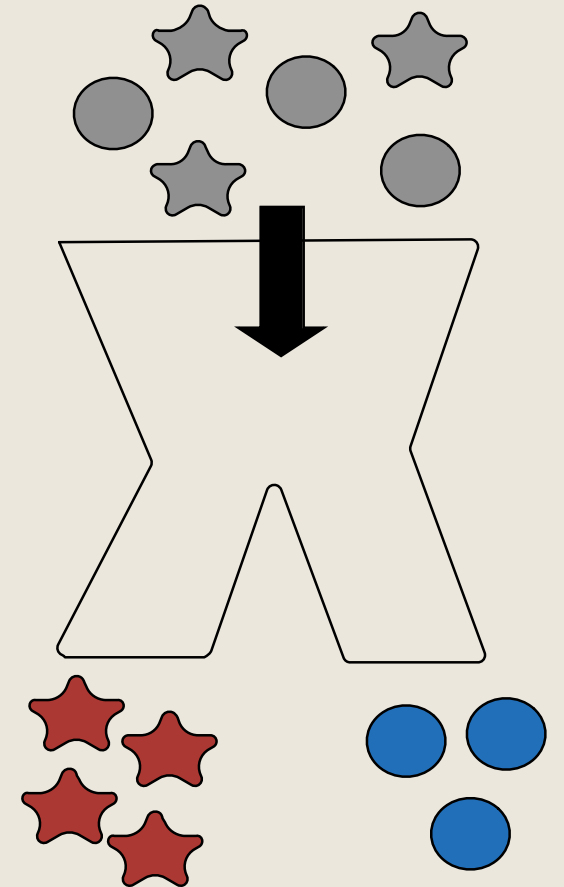
*CLASSIFICATION

Classifier: if presented with an object, can we classify it into one of several **predefined** categories?

There are many different techniques that carry out classification, but the general steps are the same:

- Use a training set to teach the classifier how to classify.
- Test/validate the classifier using new data
- Use the classifier to classify novel instances

Some classifiers (e.g. neural nets) are '**black box**' models. They do a good job, but we don't know what is happening!

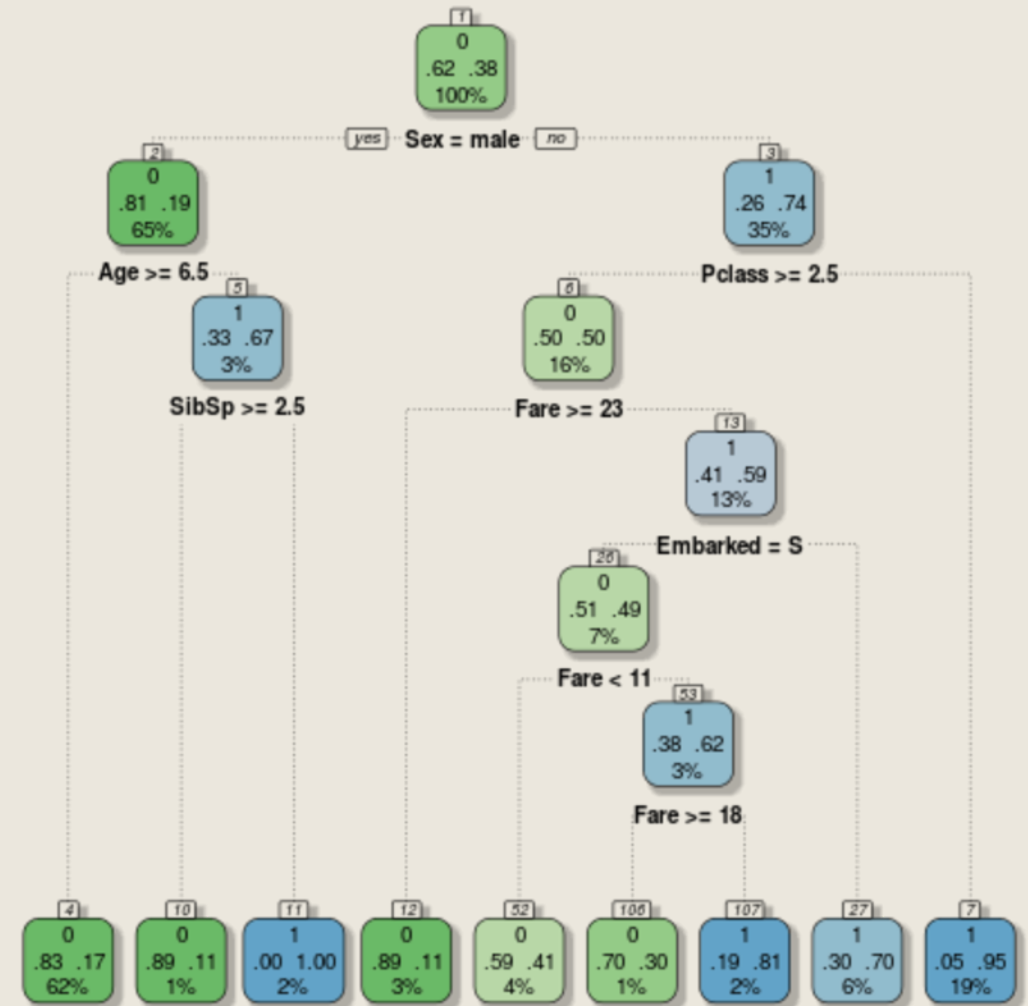


*DECISION TREE CLASSIFIERS

Decision tree: methodically use the available information to classify observations.

These trees are built automatically (**statistical learning**).

Once built, it is easy to see how the tree makes classification decisions (**white box model**).



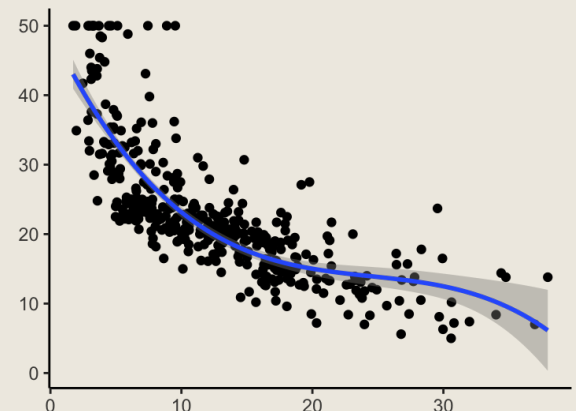
*REGRESSION ANALYSIS

If presented with an object, can we predict the value of its **response variables**?

There are many different techniques to do so, and the methods are somewhat similar (but don't always use the same framework):

- traditional approaches (linear, non-linear, trees, etc.)
- training/testing/validating approaches

As was the case for classification, some regression models are '**black boxes**'.

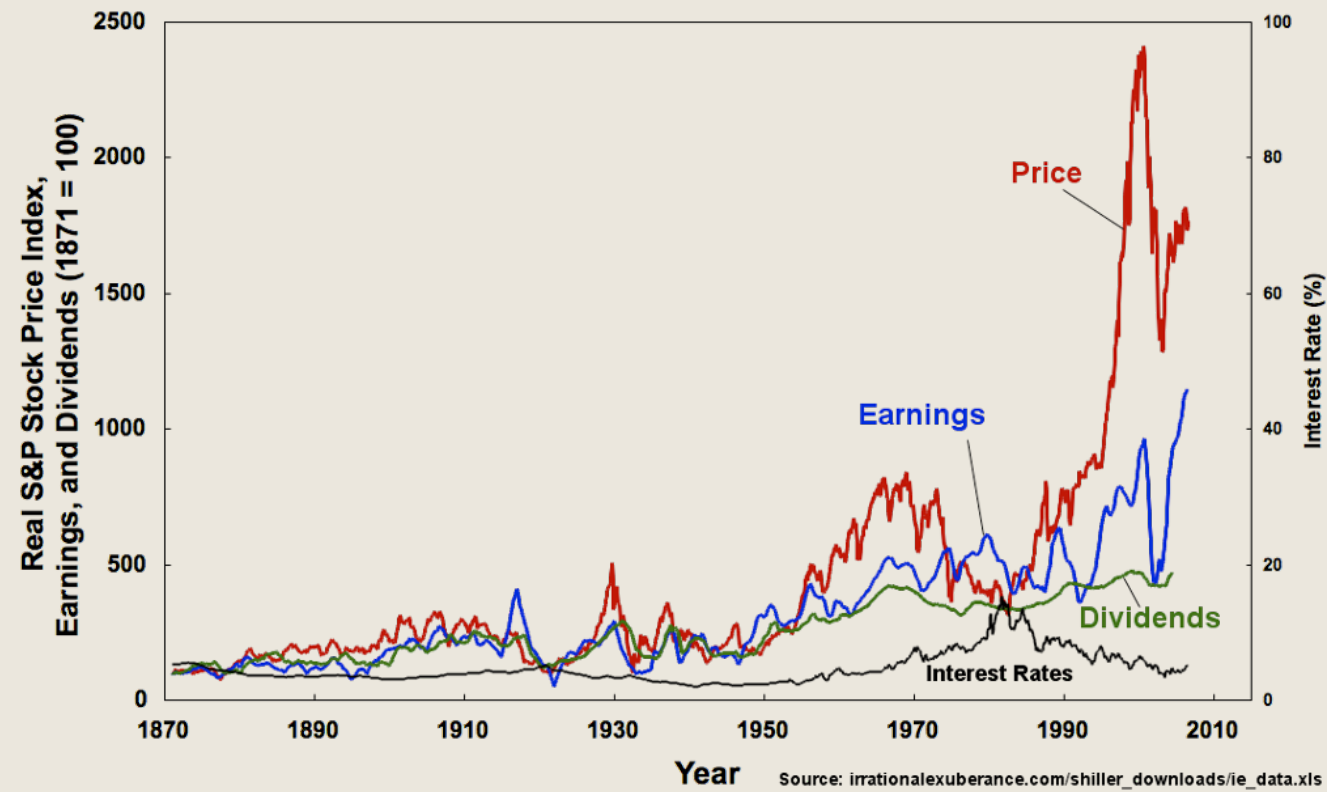


*TIME SERIES ANALYSIS

A simple time series has two variables:
time + 2nd variable.

What is the pattern of behaviour of this
second variable over time? Relative to
other variables?

Can we use this information to forecast
the behaviour of the variable in the
future?



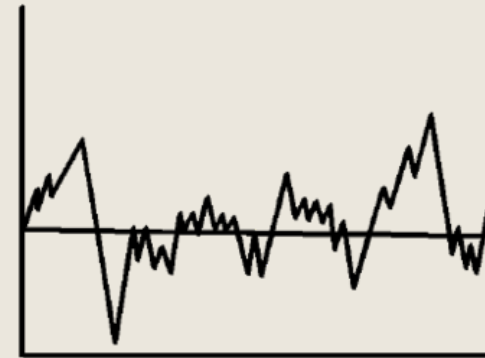
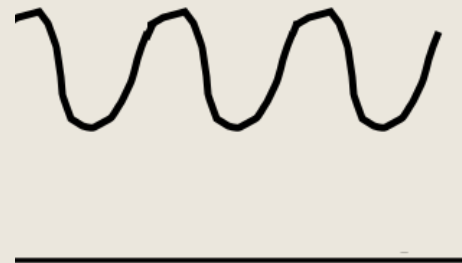
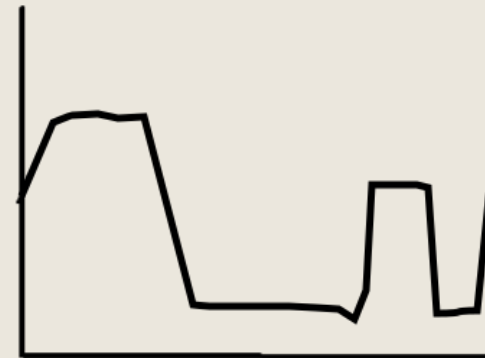
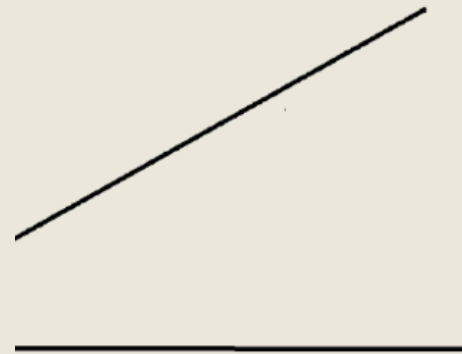
*TEMPORAL PATTERNS

The goals are familiar:

- find patterns in the data
- create a mathematical model that captures the essence of these patterns

The patterns can be quite complex – some fancy analysis typically required!

The overall series can often be broken down into multiple **component models**.



*SUPERVISED/UNSUPERVISED LEARNING

Automated behaviours vs intelligent behaviours

Supervised: examples are given (training set), algorithm learns from those

Unsupervised: learning happens based on what is seen in the data

Unsupervised techniques:

- association rules
- recommender engines
- novel categories (clustering)

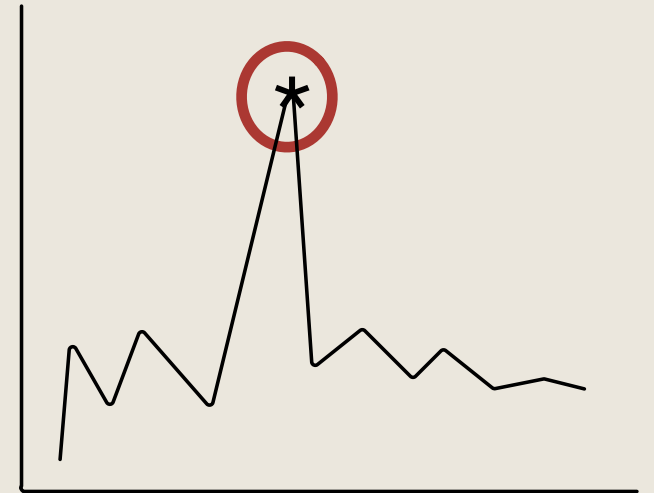
*ANOMALY DETECTION

An **anomaly** is an unexpected, unusual, atypical, or statistically unlikely event.

A fair number of data analysis pipeline are built in order to alert users when things are **out of the ordinary**.

Analytical approaches include:

- supervised (classification), unsupervised (clustering)
- association rules deviation
- ensemble techniques
- directed



SOME PRACTICAL DEFINITIONS

DATA INSIGHT FUNDAMENTALS

“What’s in a name? That which we call a rose
By any other name would smell as sweet.”

W. Shakespeare, Romeo and Juliet, Act II, Scene 2

WHAT IS DATA ANALYSIS?

Finding **patterns** in data

Using data to do something (answer a question, assist in decision-making, predict a future occurrence, draw a conclusion)

Creating models of the data

Describing or explaining a situation (the **system**)

(Testing (scientific) hypotheses?)

(Carrying out calculations on data?)

WHAT IS DATA SCIENCE?

Data science is the collection of processes by which we extract useful and **actionable insights** from data.

T. Kwartler (paraphrased)

Data science is the **working intersection** of statistics, engineering, computer science, domain expertise, and “hacking.” It involves two main thrusts: **analytics** (counting things) and **inventing new techniques** to draw insights from data.

H. Mason (paraphrased)

WHAT IS MACHINE LEARNING?

Starting around the 1940s, researchers began the earnest study of how to **teach machines to learn**.

The goal of **machine learning** was (is?) to create machines that can **learn, adapt, and respond** to novel situations

A wide variety of techniques, accompanied by a great deal of theoretical underpinning, was created to achieve this goal.

WHAT IS ARTIFICIAL/AUGMENTED INTELLIGENCE?

Artificial Intelligence (A.I.) is non-human intelligence that has been engineered rather than one that has evolved naturally.

A.I. research is research carried out in pursuit of this goal.

Pragmatically speaking, A.I. is “computers carrying out tasks that only humans can usually do”.

Augmented Intelligence is human intelligence that is supported or enhanced by machine intelligence.

WORKFLOWS AND PIPELINES

DATA INSIGHT FUNDAMENTALS

“All models are wrong.
Some models are useful.”

George Box



Supported by a foundation of stewardship, metadata, standards and quality

THE DATA SCIENCE WORKFLOW

Objective/
Rationale

Data
Collection

Data
Exploration

Utilization and
Decision
Support

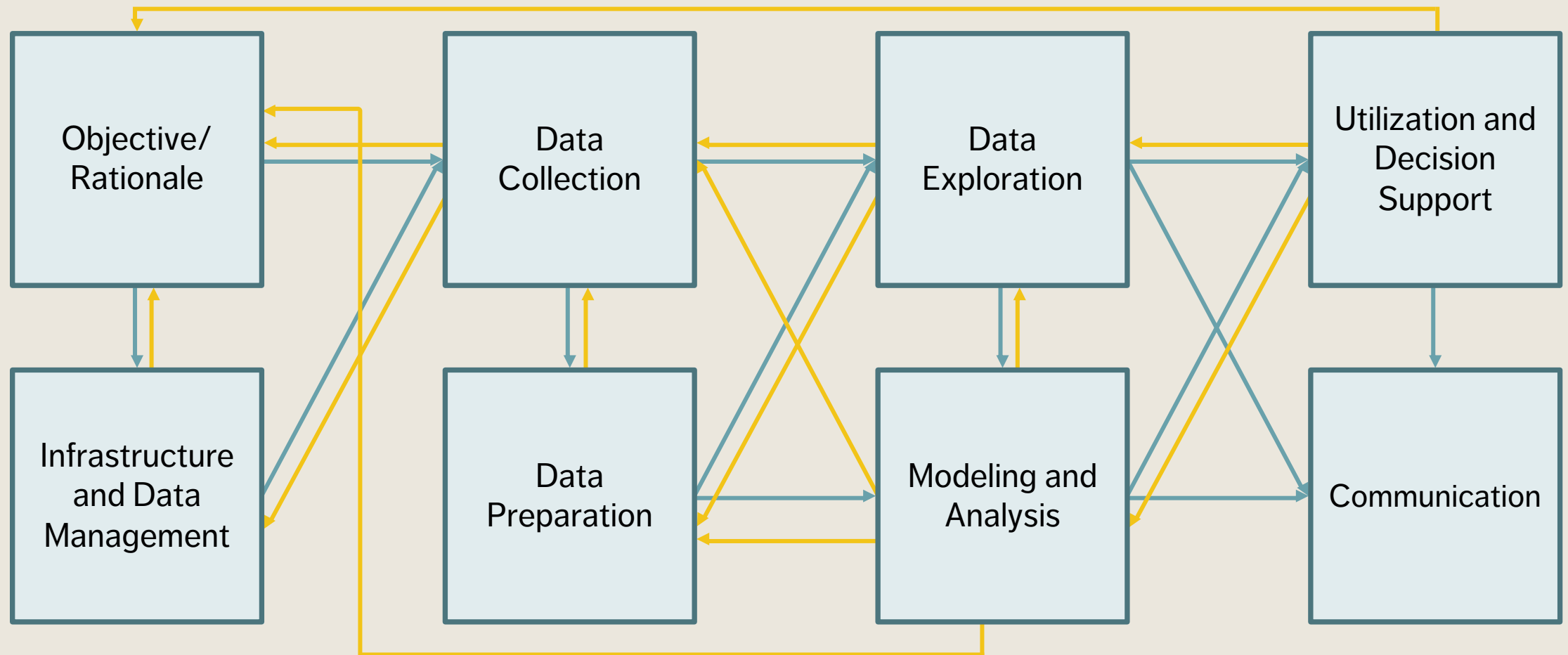
Infrastructure
and Data
Management

Data
Preparation

Modeling and
Analysis

Communication

THE DATA SCIENCE WORKFLOW



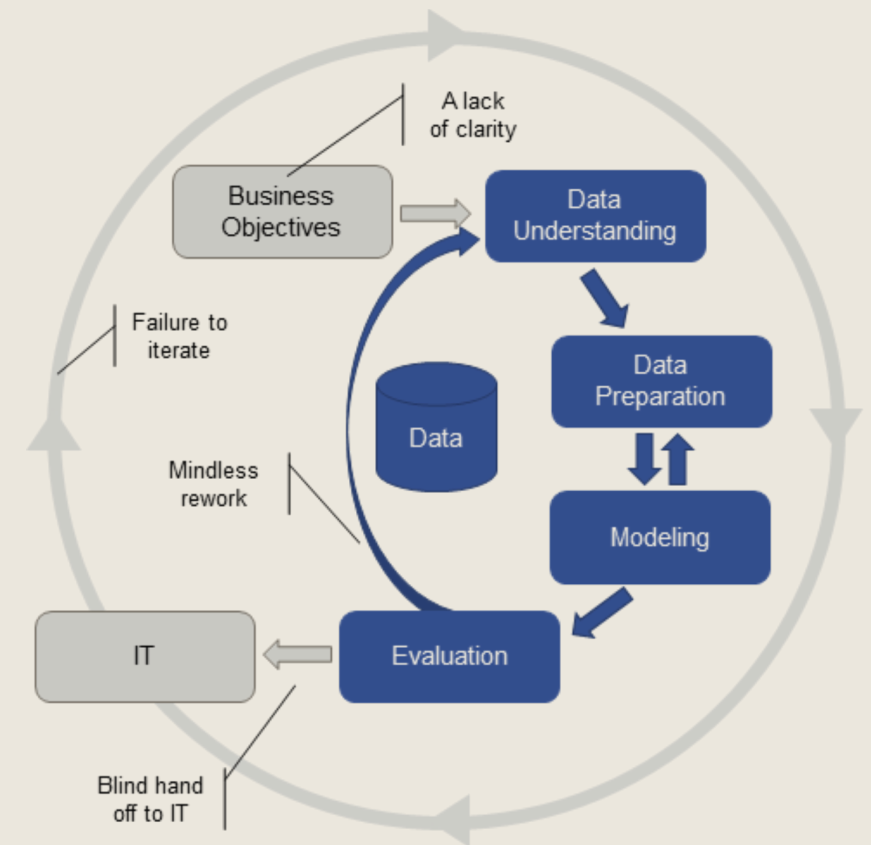
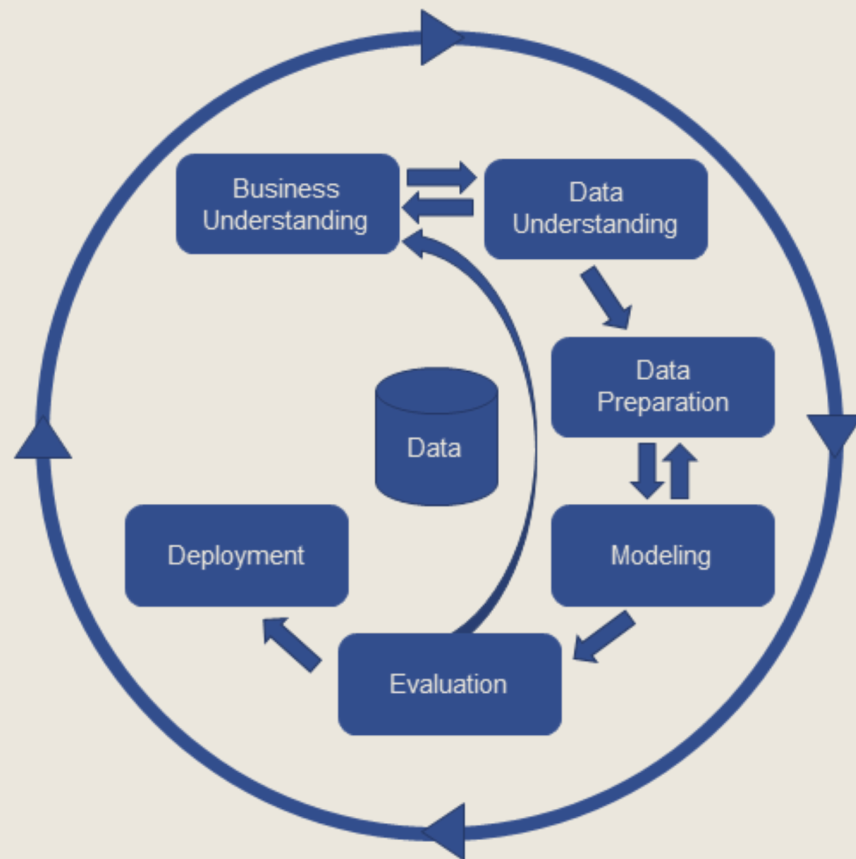
THE DATA ANALYSIS PROCESS

A **large number of analytical models** have to be generated before a final selection can be made.

Iterative process: feature selection and data reduction may require numerous visits to domain experts before models start yielding promising results.

Domain-specific knowledge has to be integrated in the models in order to beat random classifiers and clustering schemes, **on average**.

**CROSS INDUSTRY STANDARD PROCESS DATA MINING (CRISP-DM)



LIFE AFTER ANALYSIS

When an analysis or model is 'released into the wild', it can take on a life of its own.

Analysts may eventually have to relinquish control over dissemination. Results may be misappropriated, misunderstood, or shelved. What can the analyst do to prevent this?

Finally, because of **analytic decay**, it's important to view the last analytical step NOT as a static dead end, but rather as an invitation to return to the beginning of the process.

DATA SCIENCE ECOSYSTEM

Data analysis is a **team sport**, with team members needing a good understanding of both **data** and **context**

- data management
- data preparation
- analysis
- communications

Even slight improvements over a current approach can find a useful place in an organization – **data science is not solely about Big Data and disruption!**

*MODEL ASSESSMENT AND VALIDITY

Models should be **current**, **useful**, and **valid**.

Data can be used in conjunction with existing models to come to some conclusions or can be used to update the model itself.

At what point does one determine that the current data model is **out-of-date** or is **not useful anymore**?

Past successes can lead to **reluctance** to re-assess/re-evaluate a model.

MODELS AND SYSTEMS THINKING

DATA INSIGHT FUNDAMENTALS

“What if the only valid model of the
Universe is the Universe itself?”

Unknown

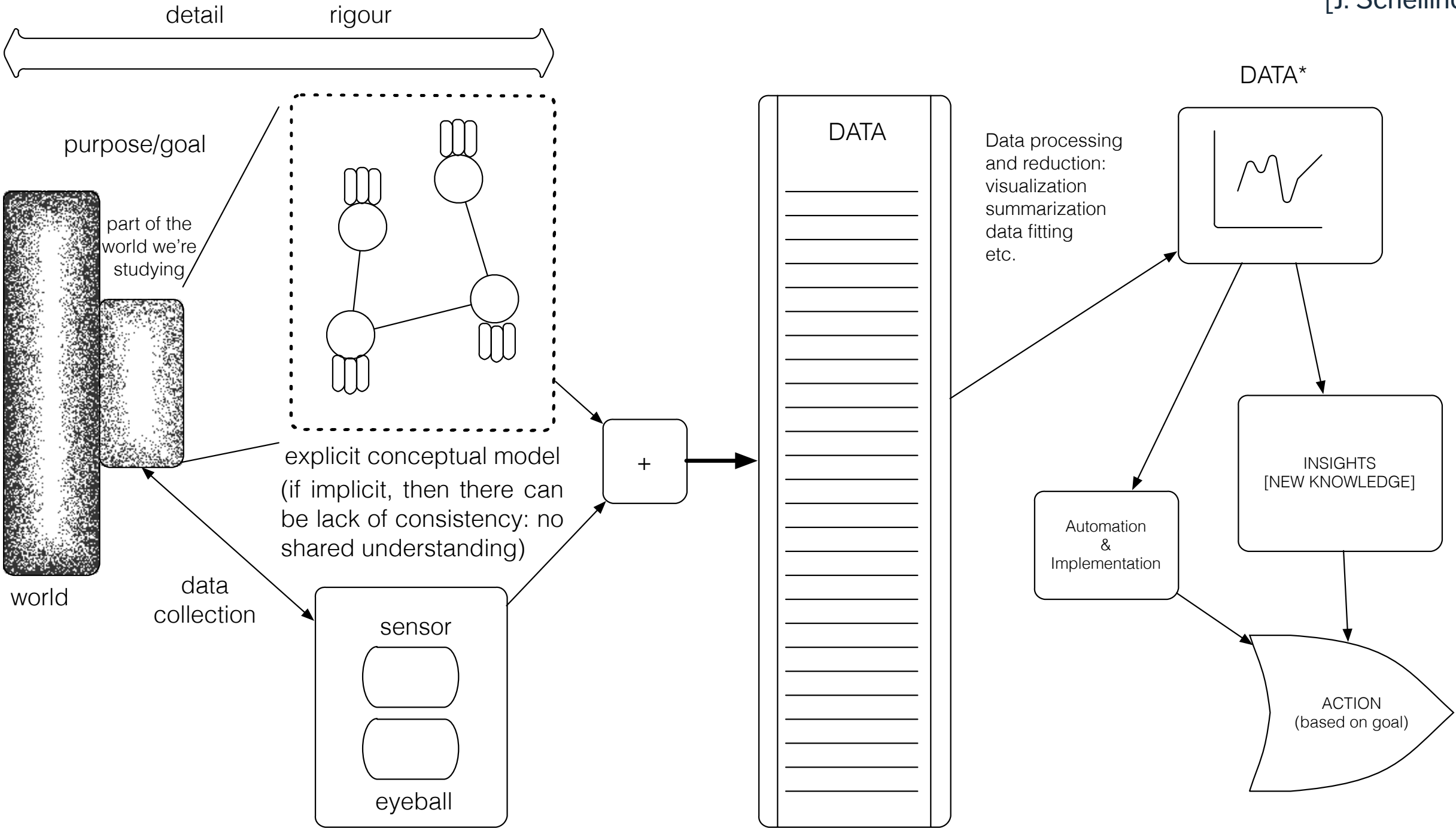
REPRESENTATIONS

A **representation** is an object that stands in for another object.

A representation may or may not physically resemble the object it represents.

Representations of the world help us to **understand**, **navigate**, and **manipulate** the world.





THINKING IN SYSTEMS TERMS

In order to understand how various aspects of the World interact with one another, we need to **carve out chunks** corresponding to the aspects and define their **boundaries**.

Working with other intelligences requires **shared understanding** of what is being studied.

A **system** is made up of **objects** with **properties** that potentially change over time. Within the system we perceive **actions** and **evolving** properties leading us to think in terms of **processes**.

THINKING IN SYSTEMS TERMS

Objects themselves have various properties. Natural processes generate (or destroy) objects and may change the properties of these objects over time.

We **observe**, **quantify**, and **record** particular values of these properties at particular points in time.

This generates data points, capturing the **underlying reality** to some degree of **accuracy** and **error** (biased or unbiased).

IDENTIFYING GAPS IN KNOWLEDGE

A **gap in knowledge** is identified when we realize that what we thought we knew about a system proves incomplete (or false).

This might happen repeatedly, at any moment in the process:

- data cleaning
- data consolidation
- data analysis

The solution is to be flexible. When faced with such a gap, **go back, ask questions, and modify the system representation.**

CONCEPTUAL MODELING

Exercise:

- an acquaintance has just set foot in your living space for the first time
- you are on the phone with them but not currently at home
- explain to them how to go about changing a breaker
- (how would you do so if the acquaintance was visually impaired?)

Conceptual models are built using methodical investigation tools

- diagrams
- structured interviews
- structured descriptions
- etc.

RELATING THE DATA TO THE SYSTEM

Is the data which has been collected and analyzed going to be of any use when it comes to understanding the system?

This question can only be answered if we understand:

- how the data is **collected**
- the **approximate nature** of both data and system
- what the data **represents** (observations and features)

Is the combination of system and data **sufficient** to understand the aspects of the world under consideration?

Real World



Model



→
Theory
→

Identification of
details relevant to
description and
translation of real-
world objects into
model variables

MODELS IN GENERAL

First principles modeling

- examine a system
- write down a set of rules/equations that describe the essence of the system
- ignore complicating details that are “less” important

Statistical modeling

- typically a set of equations with parameters
- parameters are learned (model is “trained”) using multiple data observations
- data sample vs. population

**MODELING HEURISTICS

In a sense, modeling is a **straightforward** (and **formulaic**?) process, guided by **intuition** and **experience** at each step.

Basic steps in building a statistical model:

1. **defining the goals**

- what are we trying to achieve?

- under what situations will the model be used and what is the outcome we are trying to predict?

2. **gathering data**

- what data is available?

- how many records will we have?

- generally, modelers want as much data as possible

****MODELING HEURISTICS**

Basic steps in building a statistical model: (continued)

3. deciding on the model structure

should we run a linear regression, logistic regression, or a nonlinear model? Which kind?

choices of model structure require experience and deep knowledge of the strength and weaknesses of each technique

4. preparing the data

assemble data into appropriate form for the model

encode the data into inputs, using expert knowledge as much as possible

separate the data into the desired training, testing, and validation sets

****MODELING HEURISTICS**

Basic steps in building a statistical model: (continued)

5. selecting and removing features

variables are examined for model importance and selected or eliminated
a list of candidate appropriate variables are ordered by importance

6. building candidate models

begin with baseline linear models and try to improve using more complex nonlinear models
keep in mind the environment in which the model will be implemented

7. finalizing the model

select among the candidates the most appropriate model to be implemented

8. implementing and monitoring

embed the model into necessary system process; implement monitoring steps to examine the model performance

****MODELING PITFALLS**

Common pitfalls surrounding the modeling process:

- 1. defining the goals**

- lack of clarity around problem definition

- lack of understanding of how and where the model will be used

- 2. gathering data**

- using data that is too old or otherwise not relevant going forward

- not considering additional key data sources or data sets that might be available

- 3. deciding on the model structure**

- using a modeling methodology that is not appropriate for the nature of the data (sizes, dimensions, noise...)

****MODELING PITFALLS**

Common pitfalls surrounding the modeling process: (continued)

4. preparing the data

- not cleaning or considering outliers

- not properly scaling data

- not giving enough thought to building special expert variables

- not having data from important categories of records

5. selecting and eliminating features

- keeping too many variables, making it hard for modeling, interpretation, implementation, or model maintenance

- too much reliance on simply eliminating correlated variables

****MODELING PITFALLS**

Common pitfalls surrounding the modeling process: (continued)

6. building candidate models

- overfitting

- not doing proper training/testing as one examines candidate models

- not doing a simpler linear regression to use as baseline

7. finalizing the model

- not rebuilding the final model optimally using all the appropriate data

- improperly selecting the final model without consideration to some implementation constraints

8. implementing and monitoring

- errors in implementation process: data input streams, variable encodings, algorithm mistakes

- not monitoring model performance

TAKE-AWAYS

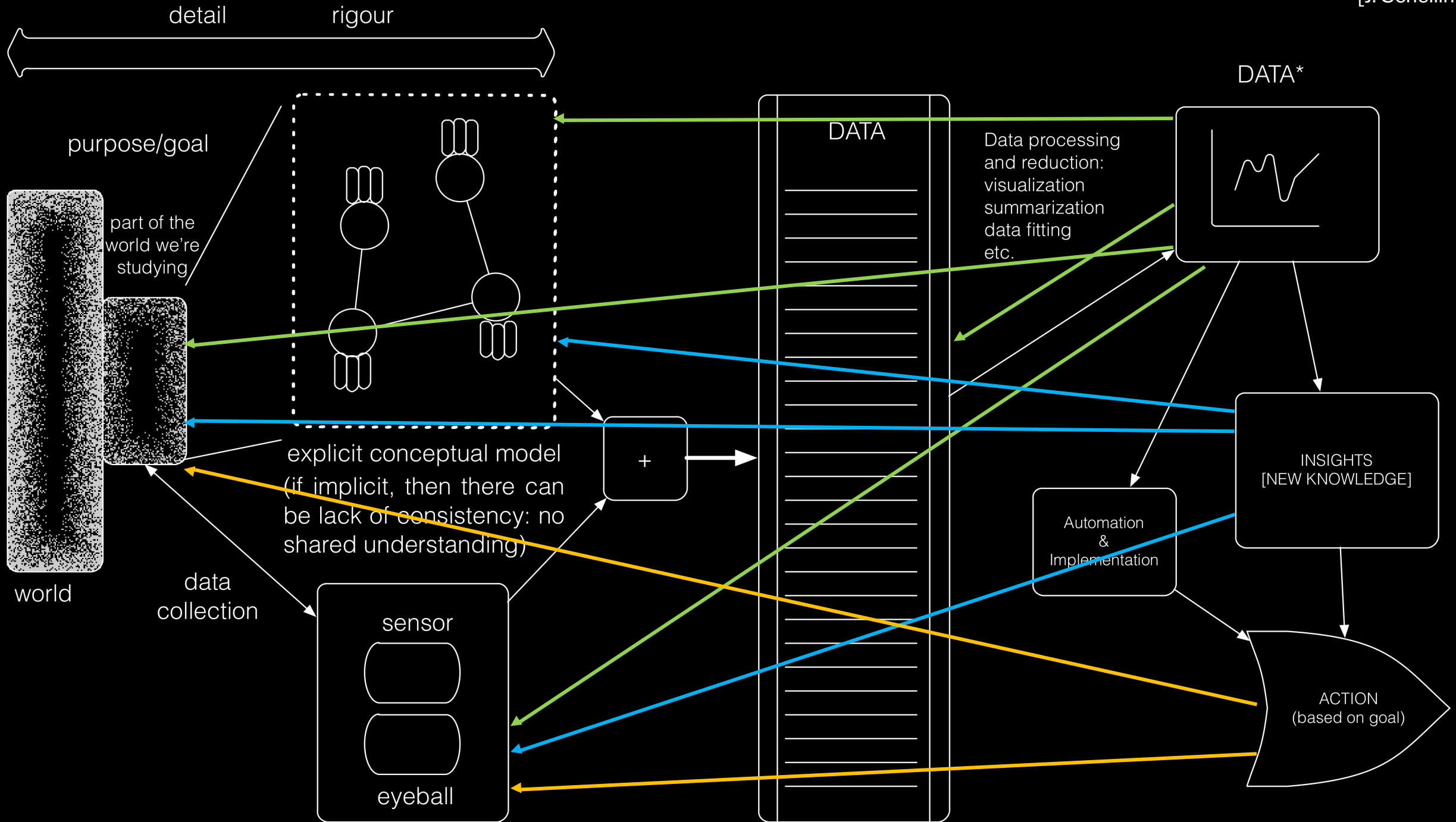
Systems can approximate certain aspects of the Universe.

System models provide the basis under which data is identified and collected, but data itself is approximate and selective.

Knowledge gaps happen – be ready to re-visit your set-up regularly.

Implicit conceptual modeling can lead to problematic situations.

If the data, the system, and the world are out of alignment, data analysis insights might ultimately prove useless.



ETHICAL CONSIDERATIONS & BEST PRACTICES

DATA INSIGHT FUNDAMENTALS

“We have flown the air like birds and swum
the sea like fish but have yet to learn the
simple act of walking the Earth like brothers.”

Martin Luther King, Jr.

What harm can come from data?

THE NEED FOR ETHICS

Formerly: “**Wild West**” mentality to data collection (and use). Whatever wasn’t technologically forbidden was allowed.

Now: professional codes of conduct are being devised for data scientists (outline responsible ways to practice data science).

Additional responsibility for data scientists; but also, **protection** against being hired to carry out questionable analyses.

Does your organization have a code of ethics for its data scientists? For its employees?

WHAT ARE ETHICS?

Broadly speaking, ethics refers to the **study** and **definition** of **right and wrong conducts**:

- “not [...] social convention, religious beliefs, or laws”. (R.W. Paul, L. Elder)

Influential ethical theories:

- Kant's **golden rule** (do onto others...), **consequentialism** (the ends justify the means), **utilitarianism** (act in order to maximize positive effect), etc.
- **Confucianism**, **Taoism**, **Buddhism** (?), etc.
- **Ubuntu**, **Maori**, etc.

WHAT ARE ETHICS?

First Nations Principles of OCAP®:

- **Ownership**
cultural knowledge, data, and information is owned by communities
- **Control**
communities have the right to control all aspects of research and information management that impact them
- **Access**
communities must have access to information and data about themselves no matter where it is held
- **Possession**
communities must have physical control of relevant data

ETHICS IN THE DATA CONTEXT

Data ethics questions:

- **Who**, if anyone, owns data?
- Are there **limits** to how data can be used?
- Are there **value-biases** built into certain analytics?
- Are there categories that should **not** be used in analyzing personal data?
- Should some data be **publicly available** to **all** researchers?

Analytically, the **general** is preferred to the **anecdotal**, but decisions made on the basis of machine learning and A.I. (security, financial, marketing, etc.) may affect real beings in **unpredictable ways**.

BEST PRACTICES

“Do No Harm”: data collected from an individual **should not be used to harm** the individual.

Informed Consent:

- Individuals must **agree to the collection and use** of their data
- Individuals must have a **real understanding of what they are consenting to**, and of **possible consequences** for them and others

Respect “Privacy”: excessively hard to maintain in the age of constant trawling of the Internet for personal data.

BEST PRACTICES

Keep Data Public: data should be kept **public** (all? most? any?).

Opt-In/Opt-Out: Informed consent requires the ability to **opt out**.

Anonymize Data: removal of id fields from data prior to analysis.

“Let the Data Speak”:

- no cherry picking
- importance of validation (more on this later)
- correlation and causation (more on this later, too)
- repeatability

GBA+

Gender-Based Analysis Plus is an analytical process used to assess how different gendered people may experience policies, programs and initiatives.

Example: [Work interruptions and financial vulnerability](#), D. Messacar, R. Morrisette

- If the data had not been collected and/or analyzed in a GBA+ manner, it would be harder to see how financial vulnerability affects different groups (if the analysis had looked only at age groups and gender, for example, instead of also including family composition).

Policies and events **impact real people in real way**, and not always in the same manner. Data analysis methods are typically used to predict and/or describe **average** (or central) outcomes, but it is often those who are far from the centre who are most affected.

Let's be smart about this.