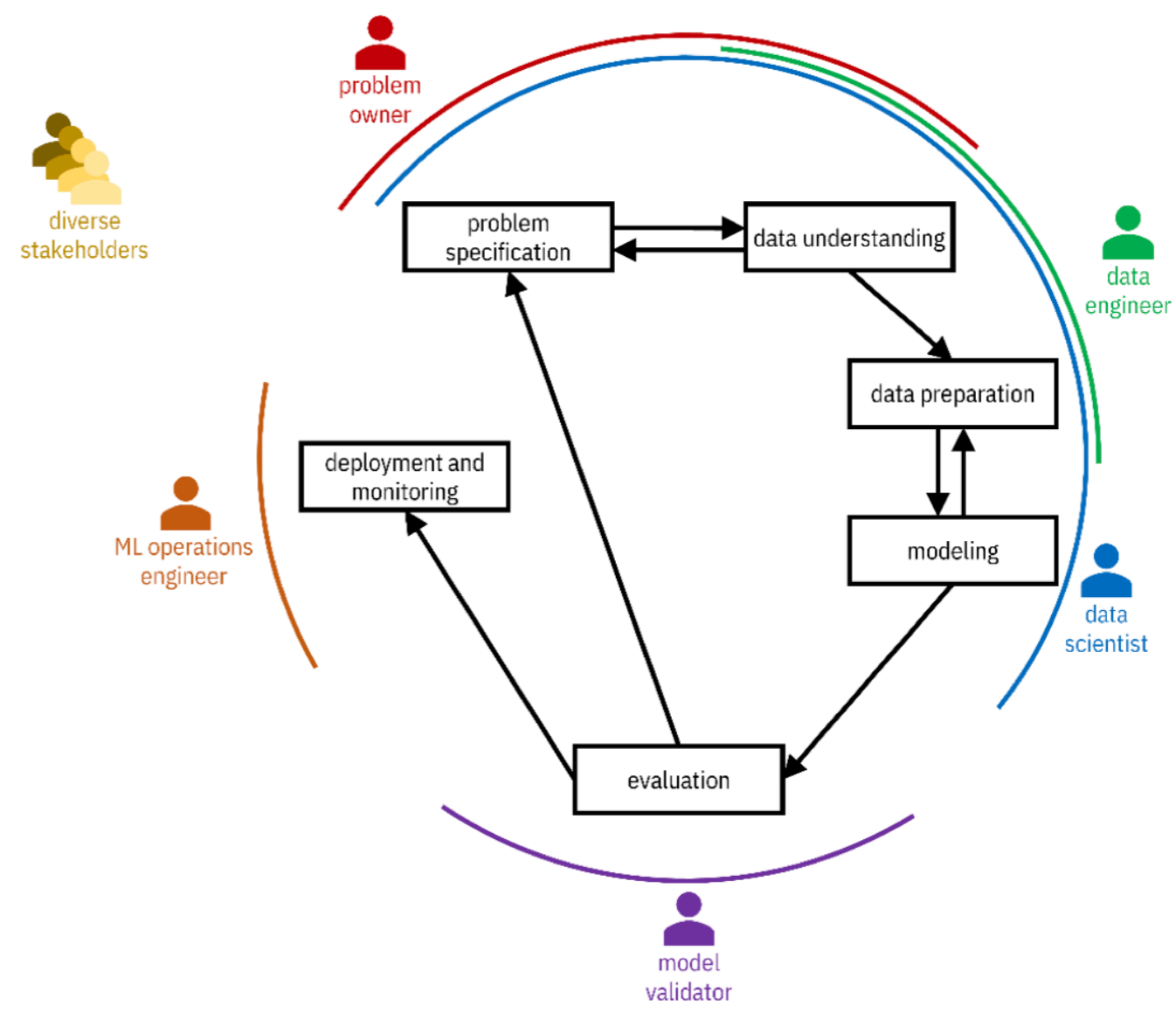# Machine Learning for Design

Lecture 2
Introduction to Machine Learning. *Part 2*

# The Machine Learning Life-Cycle

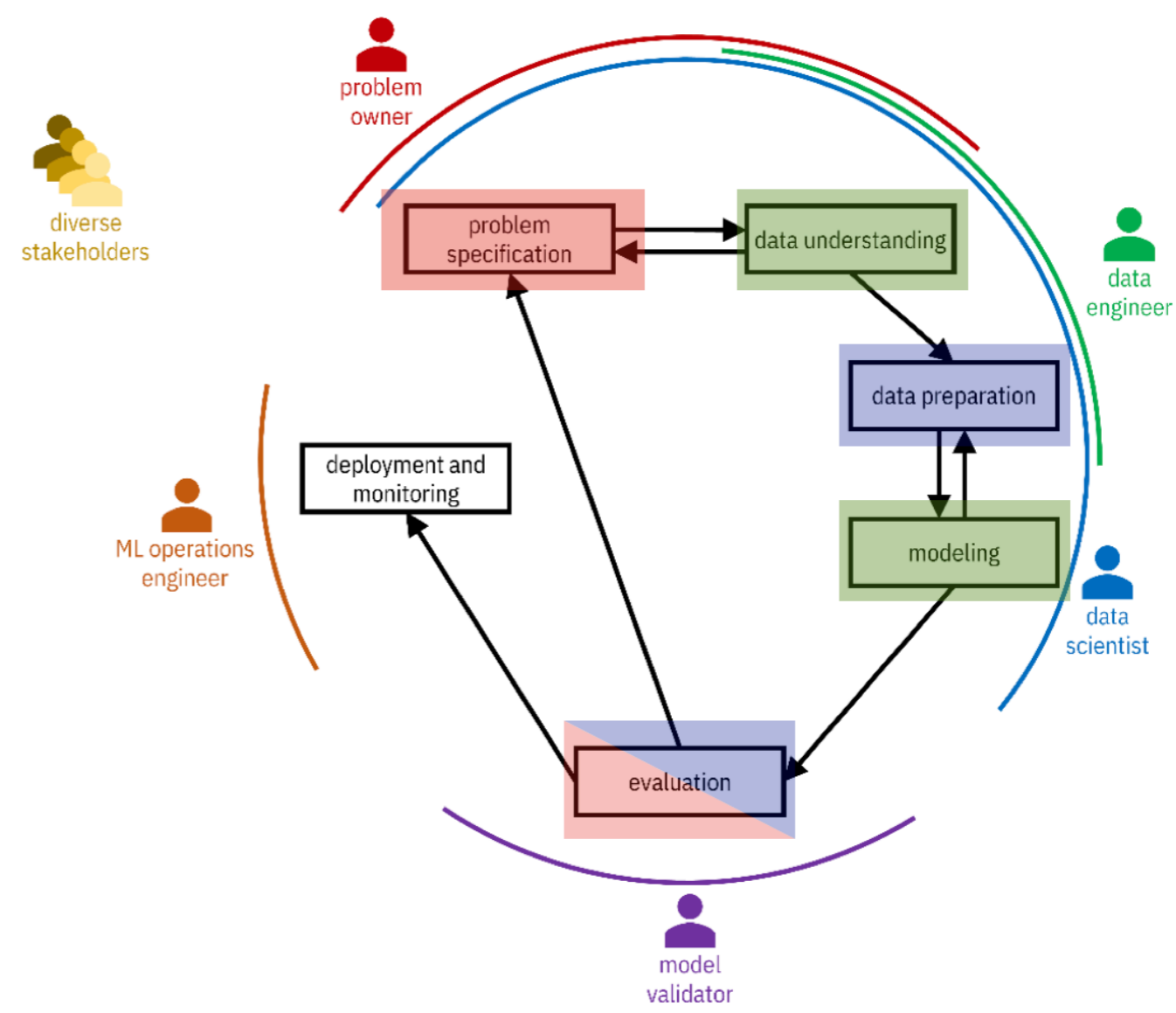# Cross-Industry Standard Process for Data Mining (CRISP-DM) methodology

# CRISP-DM In our course

# **Problem Specification**

- What is the problem owner hoping to accomplish and why?
- Why am I (being asked to) solve it?
- Am I the right person to solve this problem?
- What are the (psychological, societal, and environmental) repercussions of building this technology?
- Should this thing be built at all?
- What are the metrics of success?

# Data Understanding

Know your data!

– Data needs to be collected → Datasets

– What data is available?

– What data should be available but isn't?

– What population/system/process is your data representing?

– And what properties of such population/system/process are included (or excluded)?

– What biases (social, population, temporal) are present in your datasets?

# Data Preparation

- **Data integration**
  - Extracting, transforming, and loading (ETL) data from disparate relevant databases and other data sources
  - This step is most challenging when dealing with big data sources
- **Data cleaning**
  - Filling missing values
  - Transforming value types (e.g. binning)
  - Dropping features that should not be considered
- **Feature engineering**
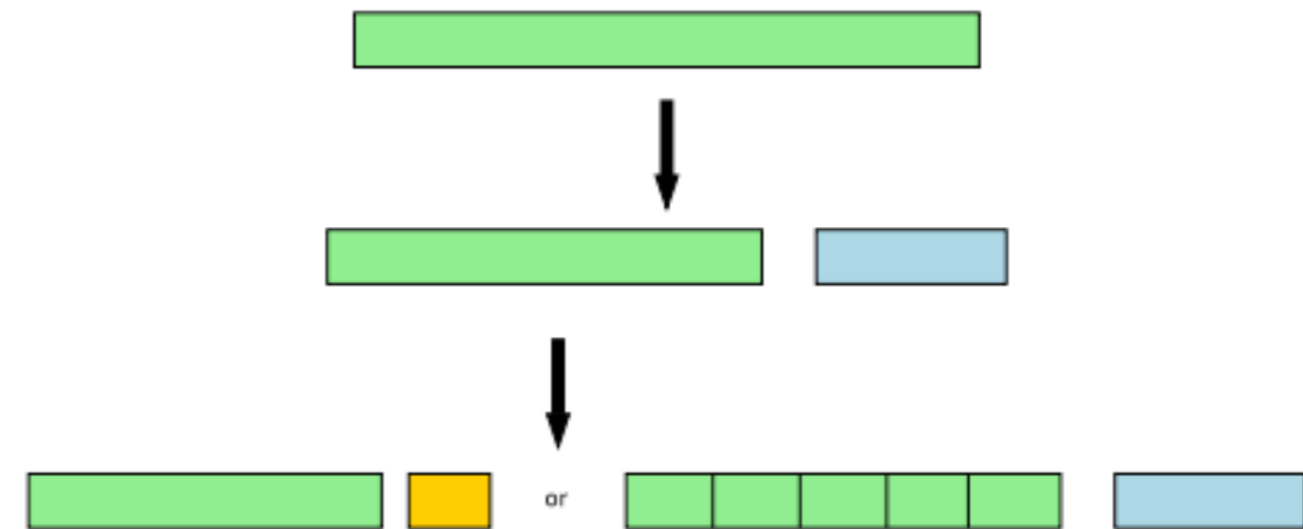  - Transform the data to derive new features

# Modeling

– **Select** a training algorithm

– Use it to **find patterns** in the training dataset

– **Generalize** them to fit a statistical model

– **Enhance** the model to satisfy additional objectives and constraints captured in the problem specification

  – e.g., increase reliability, mitigate biases, generate explanations

– **No free-lunch theorem**

  – There is no one best machine learning algorithm for all problems and datasets

# **Evaluation**

– Testing and validation of the model

    – Also against the problem specification requirements

– Performed on data not used for training

    – Hold out dataset

# Model auditing/risk management



The Pyramid of Criticality for AI Systems

# Deployment and monitoring

– What data infrastructure will bring new data to the model?

– Will predictions be made in batch or one-by-one?

– How much latency is allowed?

– How will the user interact with the system?

– Is there a problem here?

– Tools to monitor the model's performance

– And ensure it is operating as expected

# Data

## The raw material

# Data



| sepal_lenght | sepal_width | petal_lenght | petal_width | Class |
|---|---|---|---|---|
| 5.0 | 3.3 | 1.4 | 0.2 | Iris-setosa |
| 7.0 | 3.2 | 4.7 | 1.4 | Iris-versicolor |
| 5.7 | 2.8 | 4.1 | 1.3 | Iris-versicolor |
| 6.3 | 3.3 | 6.0 | 2.5 | Iris-virginica |

Feature  Feature  Feature  Feature  Label

Setosa   Virginica   Versicolor

petal  sepal

Record / Sample / Data Item

Label Value

Feature Value

Dataset Size

Dataset Dimensionality

# Types of Features / Label Values

- **Categorical**
  - Named Data
  - Can take numerical values, but no mathematical meaning
- **Numerical**
  - -Measurements
  - Take numerical values (discrete or continuous)

# Categorical Nominal

- No order

- No direction

- e.g. marital status, gender, ethnicity

# Categorical Ordinal

- Order

- Direction

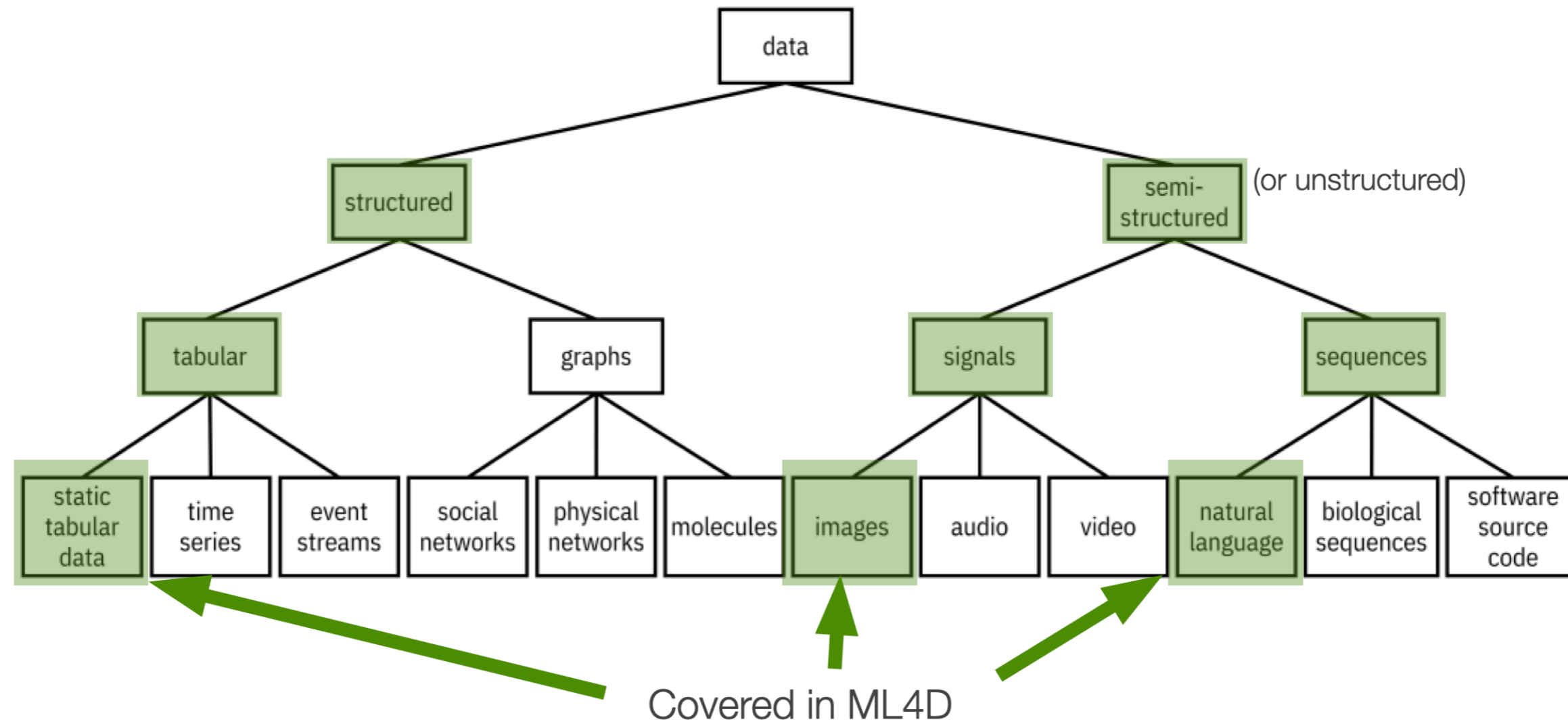- e.g., letter grades ($A$,$B$,$C$,$D$), ratings (*dislike*, *neutral*, *like*)

## Numerical Interval

- Difference between measurements
- No true zero or fixed beginning
- e.g., temperature (C or F), IQ, time, dates

## Numerical Ratio

- Difference between measurements
- True zero exists
- e.g., temperature (K), age, height

# Data Modalities



Covered in ML4D

# Key Dimensions

| Modality | Quantity | Quality | Freshness | Cost |
|---|---|---|---|---|
| Structured | Number of records | Errors | Rate of collection | Acquisition |
| Semi-structured | Number of features | Missing data | | Licensing |
| | | Bias | | Cleaning and integrations |

# Static Tabular Data



| | Feature | Feature | Feature | Feature | Label |
|---|---|---|---|---|---|
| | sepal_lenght | sepal_width | petal_lenght | petal_width | Class |
| | 5.0 | 3.3 | 1.4 | 0.2 | Iris-setosa |
| | 7.0 | 3.2 | 4.7 | 1.4 | Iris-versicolor |
| | 5.7 | 2.8 | 4.1 | 1.3 | Iris-versicolor |
| | 6.3 | 3.3 | 6.0 | 2.5 | Iris-virginica |

Record / Sample / Data Item

Label Value

Feature Value

Dataset Size

Dataset Dimensionality

# Time Series

– tabular data with **time** feature

– For instance

  – Sensor data, Stock market data

– Label is usually associated with a set of records

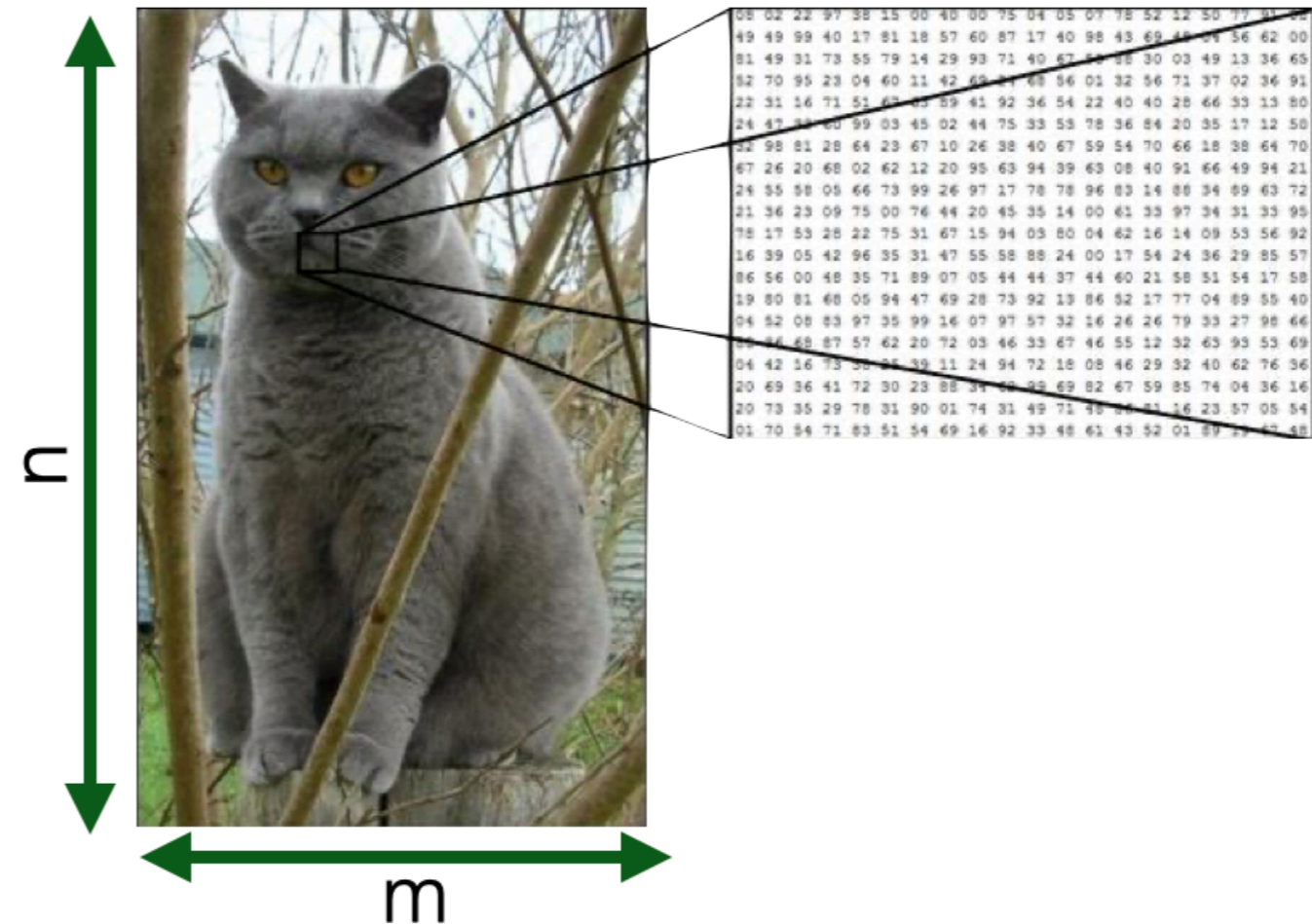  – e.g. a continuous movement of the phone indicating an action

Time
Feature

| Timestamp | X | y | Z | Class |
|---|---|---|---|---|
| 15060015925 | 2.04 | 3.72 | 8.12 | Device Rotation |
| 15060015943 | 1.96 | 4.73.68 | 7.56 | |
| 15060015980 | 1.63 | 3.56 | 6.53 | |
| 1506001610 | 1.06 | 3.76 | 5.81 | |

# Images

– Visual content acquired through cameras, scanners, etc.

– Each pixel in an image is a feature

  – But spatially and geometrically organised

  – e.g., edges, corners

– Feature values are numerical values across channels

  – e.g., *R,G,B*

– Dimensionality $\rightarrow$ *n x m*

Image →

| P(1,1) | P(2,1) | P(3,1) | ... | P(n,m) | Class |
|--------|--------|--------|-----|--------|-------|
| 255, 0, 0 | 255, 1, 1 | 255, 0, 0 | | R,G,B | Cat |
| 255, 213, 0 | 255, 213, 1 | 255, 213, 4 | | R,G,B | Dog |
| | | | | | Cat |
| | | | | | Duck |

# More in Module 1

# Textual documents

– Sequence of alphanumerical characters

  – Short: e.g. tweets

  – Long: e.g Web documents, interview transcripts

– Features are (set of) words

  – Words are also syntactically and semantically organised

– Feature values are (set of) words occurences

– Dimensionality → at least dictionary size

★★★★☆  **I wear this mask to sing lullabies to my children ...**, 24 May 2015

By **Sir Chubs**

**Verified Purchase** (What is this?)

**This review is from: Overhead Rubber Penguin Mask Happy Feet Animal Fancy Dress (Toy)**

I wear this mask to sing lullabies to my children. They are terrified of the mask. Whenever they protest about their bed time, or ask for too many sweets, I whip on the mask, and they soon know who is the King Penguin.

Document →

| I | Wear | Mask | ... | W(n) | Class |
|---|------|------|-----|------|-------|
| 1 | 1 | 1 | | 0 | Spam |
| 0 | 0 | 1 | | 0 | Not Spam |
| | | | | | Spam |
| | | | | | |

# More in Module 2

# Data Sources

| Purposefully Collected Data | Administrative Data | Social Data | Crowdsourcing |
|---|---|---|---|
| Survey | Call records | Web pages | Distributed sensing |
| Census | Financial transactions | Social Media | Implicit crowd work (e.g. captcha) |
| Economic Indicators | Travel Data | Apps | Micro-work platforms (e.g Amazon Mechanical Turk) |
| Ad-hoc sensing | GPS Data | Search Engines | |

# Data Sources

| Purposefully Collected Data | Administrative Data | Social Data | Crowdsourcing |
|---|---|---|---|
| *Modality*: mostly structured | *Modality*: mostly structured | *Modality*: mostly semi-structured | *Modality*: all |
| *Quantity*: low | *Quantity*: high | *Quantity*: low | *Quantity*: mid-low |
| *Quality*: high | *Quality*: high | *Quality*: low | *Quality*: mid |
| *Freshness*: low | *Freshness*: high | *Freshness*: high | *Freshness*: mid |
| *Cost*: high | *Cost*: high | *Cost*: low | *Cost*: mid-low |

# Categories of Machine Learning

# How do machines learn?

Training Data → Model Training → ML Model → Model Enhancement → Final ML Model

# On Models

A physical, mathematical, logical, or conceptual representation of a system, entity, phenomenon, or process

– A **simple(r)** representation of reality helping us understand how something works or will work.

  – **Not truthful**, just a **useful** one

– The goal of models is to make a particular part or feature of the world more accessible to understand, define, quantify, visualise, or simulate

# Examples of models

Architecture plans

Maps

Music Sheet

Mathematical laws of physics!

Machine Learning (statistical) Models
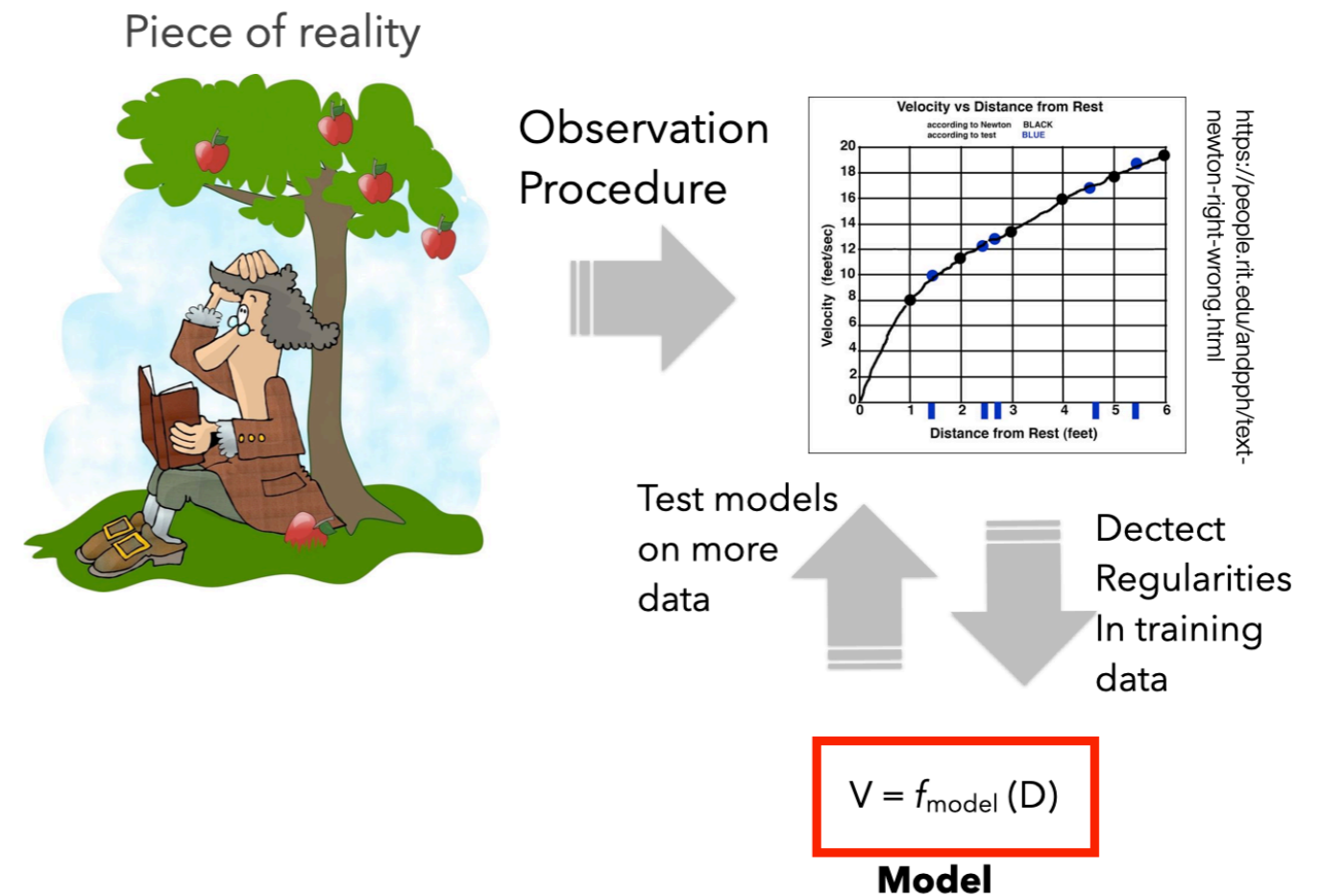


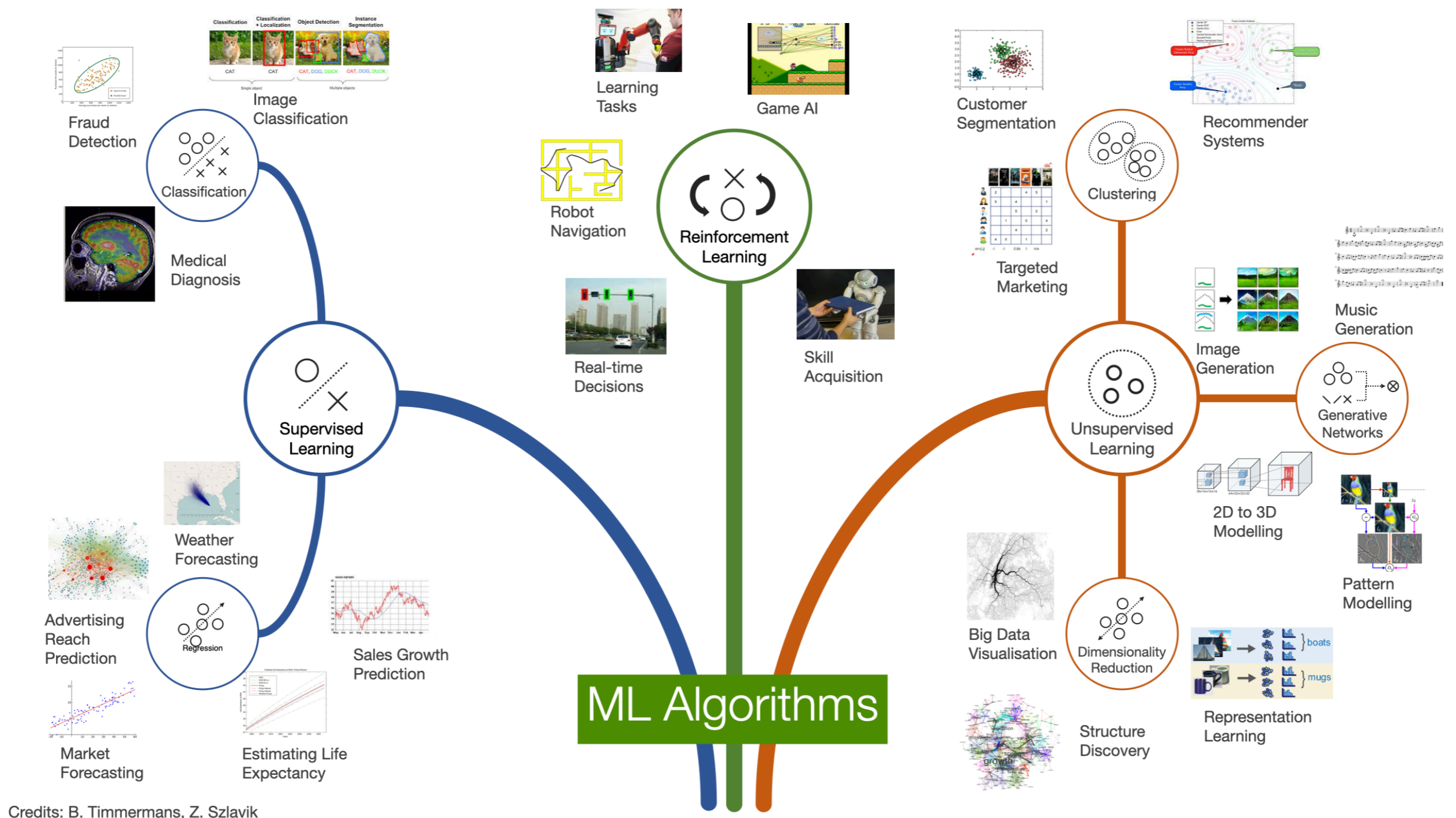$$F_g = G \frac{M_1 M_2}{d^2}$$

## Scientific Models

– GOAL: <mark>explain reality</mark>

– Created to make predictions about the outcomes of future experiments

  – e.g., apples on the moon

– Tested against the **outcome**

– If data from new experiments don't agree, the model has to be modified/extended / refined

  – Falsifiability

– Scientific models should be *small* and *simple*.

– They should generalize phenomena observed in new ways.

Piece of reality

Observation Procedure

https://people.rit.edu/andpph/text-newton-right-wrong.html

Dectect Regularities

Make predictions and test in experiments

## ML Models

- GOAL: <mark>describe the data</mark>

- Designed to capture the *variability* in observational data by exploiting regularities/symmetries/redundancies

- A good ML model doesn't need to explain reality, it **just describe data**

- They don't need to be simple or transparent, or intelligible. Just **accurate**

  - *Black box*

- ML models may be large and complex.

- They should generalize to new data obtained in the same way as the training data

  - Same application context and data acquisition process

Piece of reality

Observation Procedure
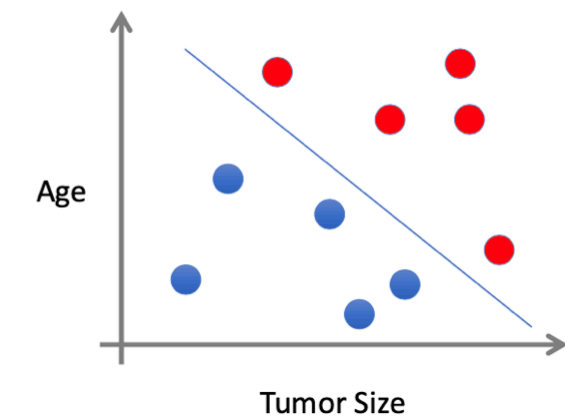
https://people.rit.edu/andpph/text-newton-right-wrong.html

Test models on more data

Dectect Regularities In training data

$V = f_{model}(D)$

**Model**

# ML Algorithms

## Supervised Learning

### Classification
- Fraud Detection
- Image Classification
- Medical Diagnosis

### Regression
- Advertising Reach Prediction
- Weather Forecasting
- Sales Growth Prediction
- Market Forecasting
- Estimating Life Expectancy

## Reinforcement Learning
- Learning Tasks
- Game AI
- Robot Navigation
- Real-time Decisions
- Skill Acquisition

## Unsupervised Learning

### Clustering
- Customer Segmentation
- Recommender Systems
- Targeted Marketing

### Generative Networks
- Image Generation
- Music Generation
- 2D to 3D Modelling
- Pattern Modelling

### Dimensionality Reduction
- Big Data Visualisation
- Structure Discovery
- Representation Learning

Credits: B. Timmermans, Z. Szlavik

# Supervised Learning

- Input: **labeled** data
  - Data + expected prediction
- During training, labels are used to associate patterns with outputs
- Learns how to make input-output **predictions**

- *Classification*
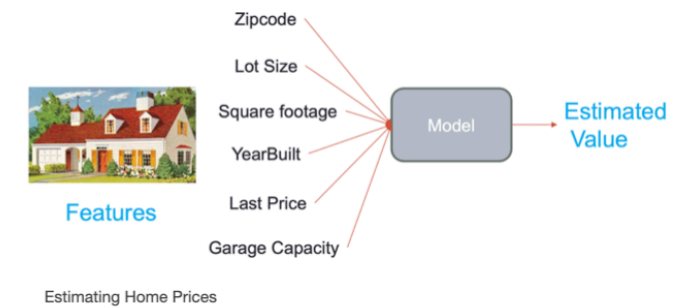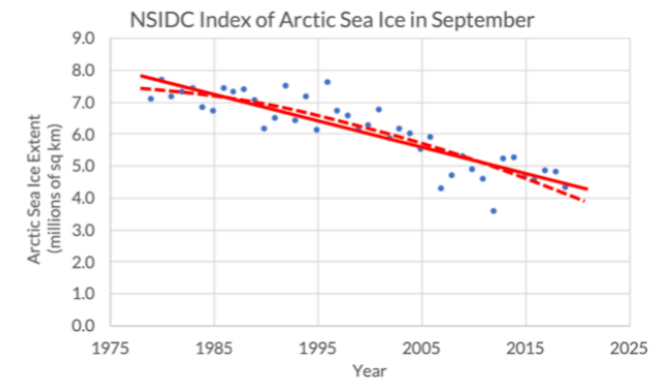- *Regression*
- *Ranking*
- *Recommendation*

# Classification

– Learn to output a **category label**

– Binary

   – e.g. *Spam / not Spam*, *Cat / not cat*
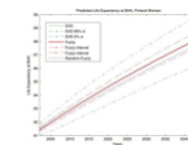
– Multi-class

   – e.g. *cat*, *dog*, *bird*

Ocular Tumor (Malignant / Benign)

$f(x)$

Predict Benign     Predict Malignant

Tumor Size

Age

Tumor Size

Fraud Detection

Medical Diagnosis

| Classification | Classification + Localization | Object Detection | Instance Segmentation |
|---|---|---|---|
| CAT | CAT | CAT, DOG, DUCK | CAT, DOG, DUCK |

Single object    Multiple objects

Image Classification

# **Regression**

– Learn to output one or more **numbers**

– e.g., value of a share, number of stars in a review
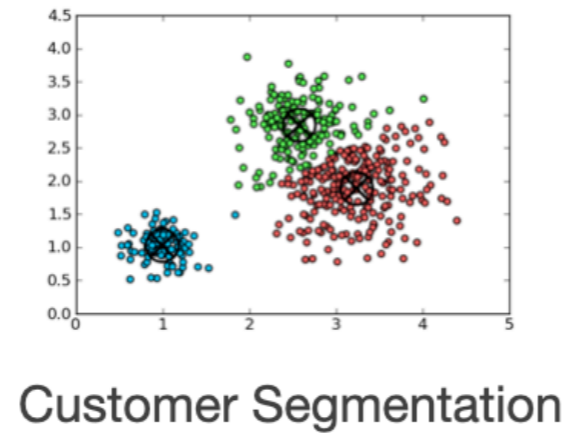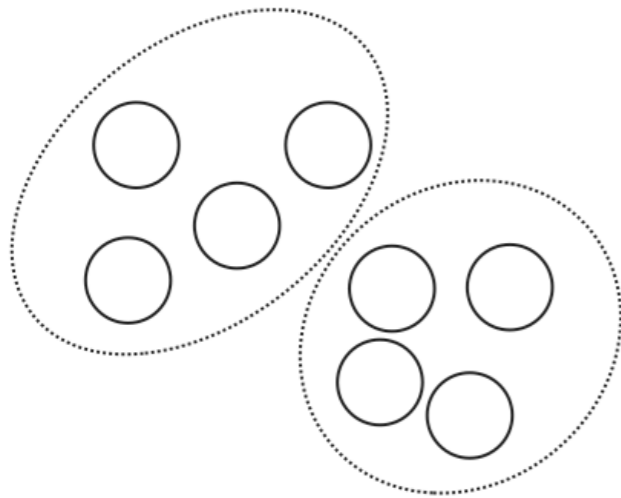
# **Unsupervised Learning**

– Input: **unlabeled** data

– The machine learns structures (patterns) from the data without human guidance

– *Clustering*

– *Dimensionality Reduction* (e.g. Large Language Models)
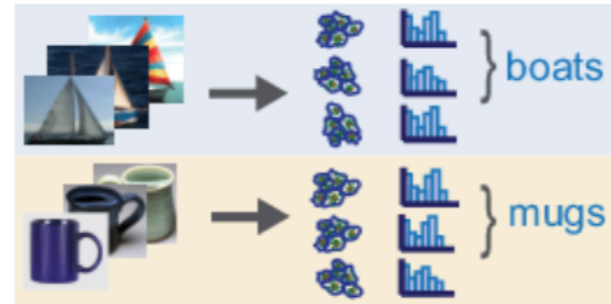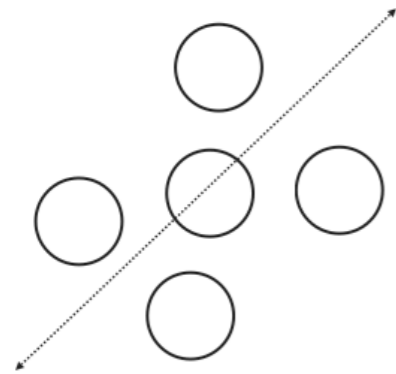
– *Anomaly detection*

Input Data

**ML Algorithm**

# Clustering



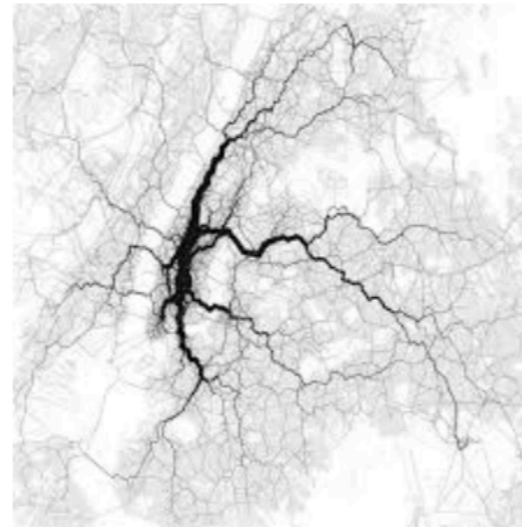Customer Segmentation

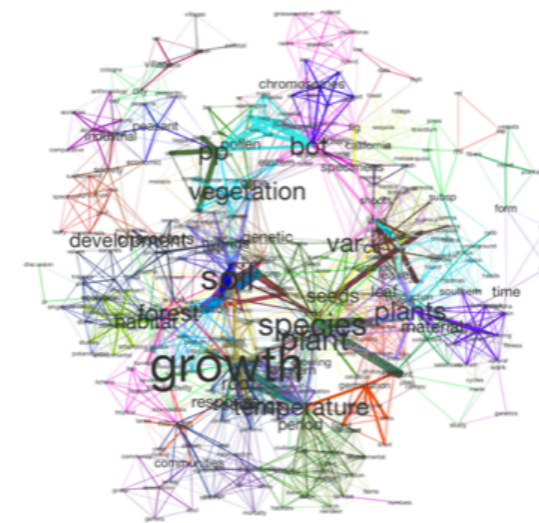Targeted Marketing

Recommender Systems

# Dimensionality Reduction
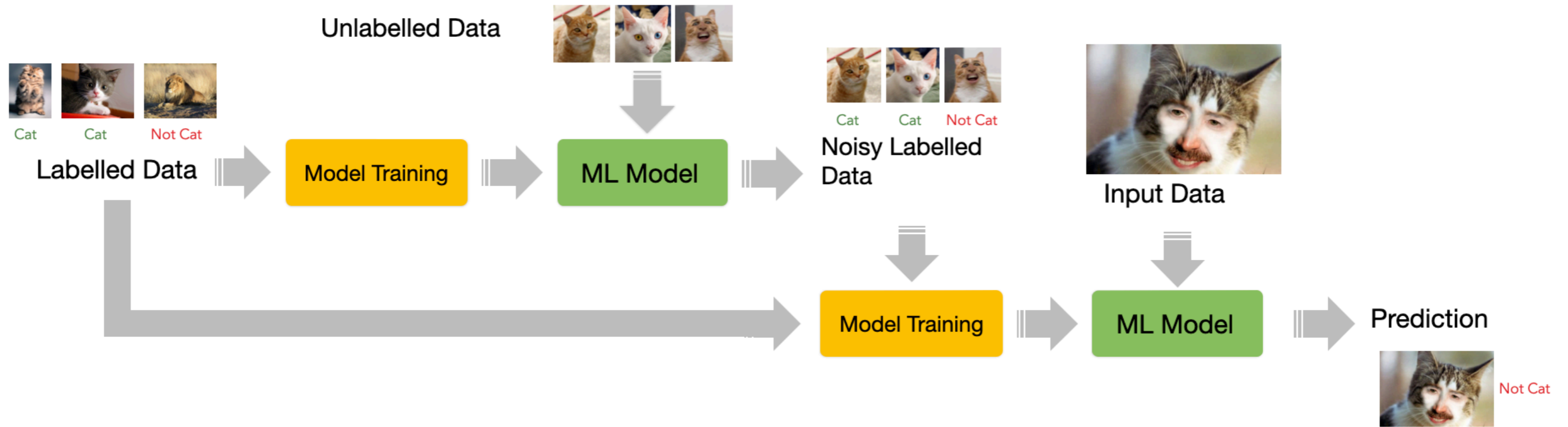


Foundational Models
For Transfer Learning



Big Data Visualisation



Structure Discovery

# Semi-Supervised Learning

– Combination of **supervised** and **unsupervised** learning

– Few **labeled** data in the input are used to create **noisy labeled data**

– With more labeled data, the machine learns how to make input-output **predictions**

Unlabelled Data

Cat    Cat    Not Cat

Labelled Data

Model Training

ML Model

Cat    Cat    Not Cat

Noisy Labelled Data

Input Data

Model Training
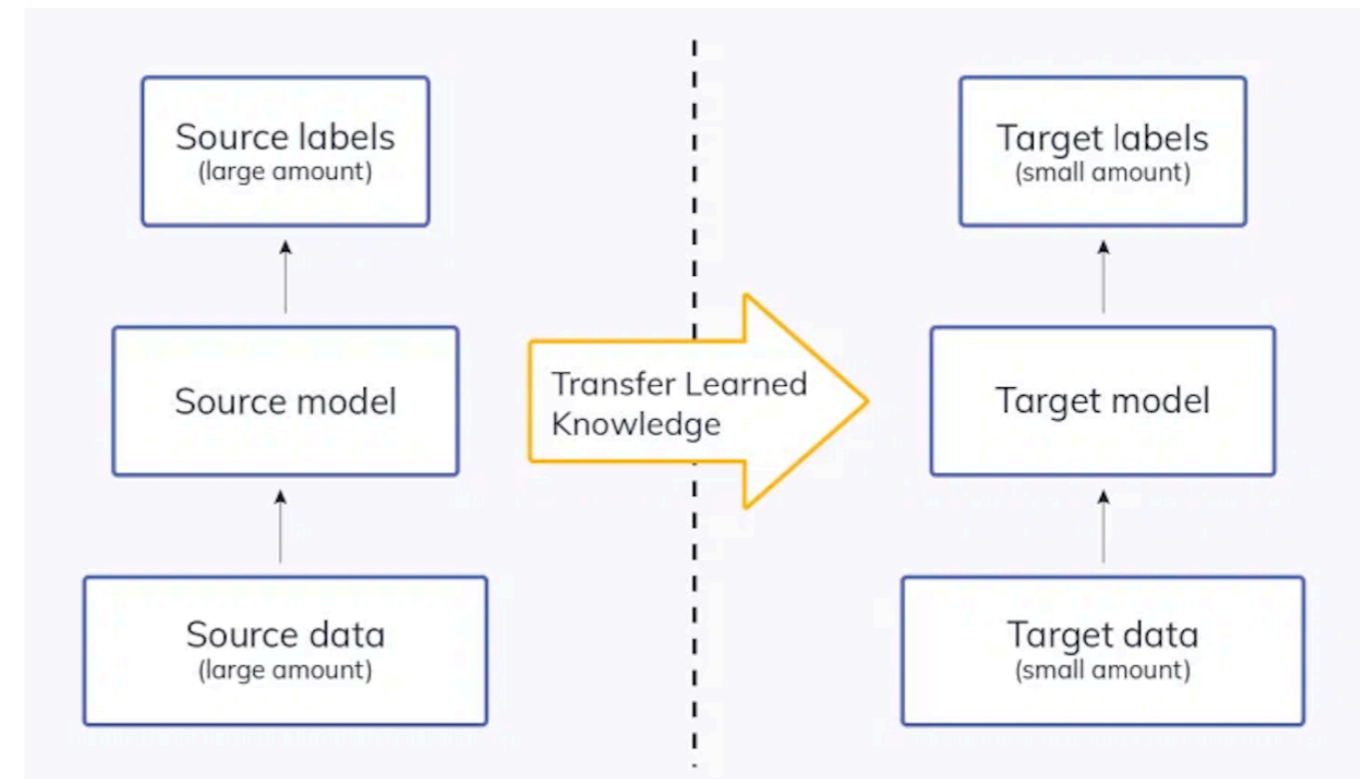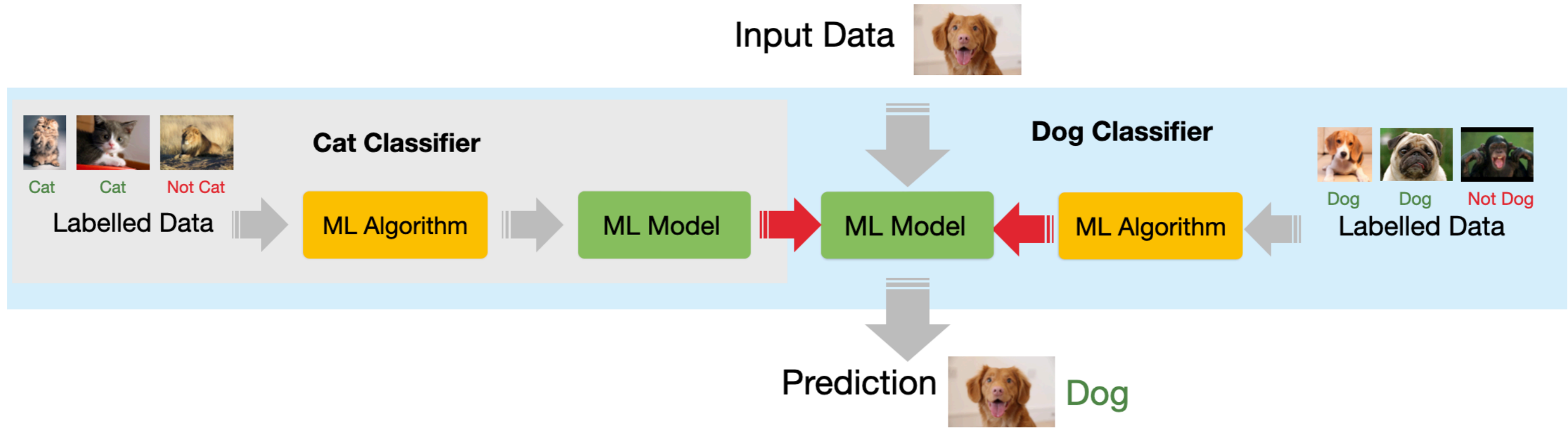
ML Model

Prediction

Not Cat

# Transfer Learning

Often called *fine-tuning*

Reuse a model trained for one task is **re-purposed** (tuned) on a different but related task

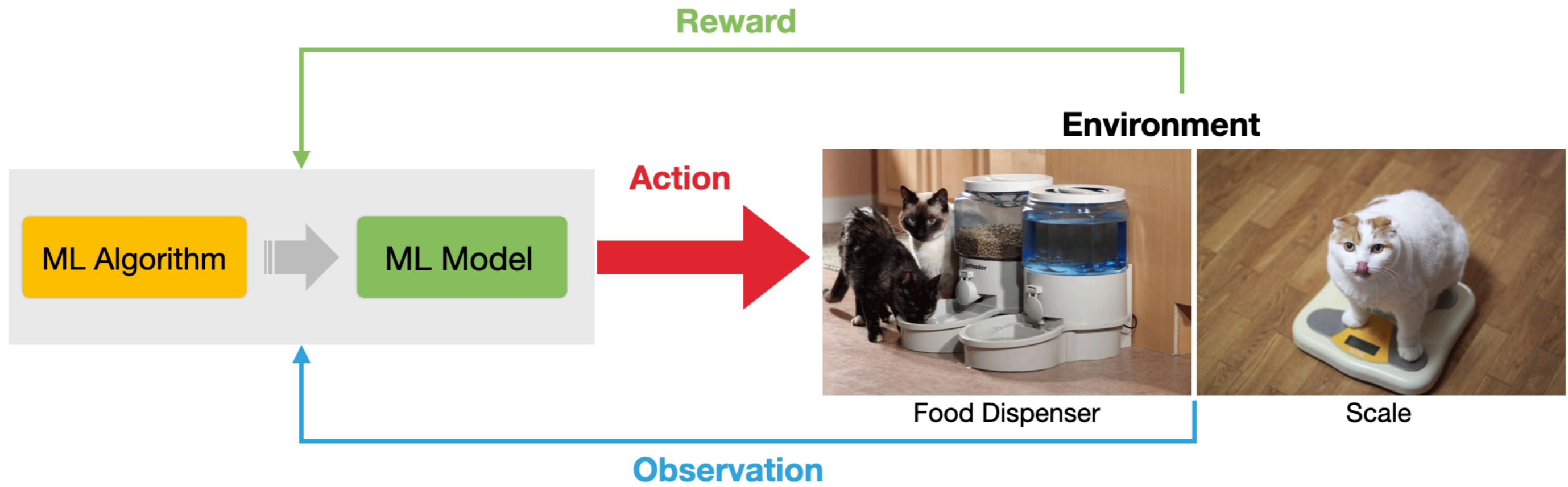Useful in tasks lacking abundant data

Input Data

**Cat Classifier**

Cat    Cat    Not Cat

Labelled Data → ML Algorithm → ML Model → ML Model ← ML Algorithm ← Labelled Data

**Dog Classifier**

Dog    Dog    Not Dog

Prediction    Dog

# Reinforcement Learning

– Data about the **environment** and **reward function** as input

– The machine can perform **actions** influencing the environment

– The machine learns behaviours that result in **greater reward**

Reward

Environment

Action

ML Algorithm → ML Model

Food Dispenser

Scale

Observation

# Don't forget domain expertise

– ML makes some tasks automatic, but we still need our brains

– More in Module 3 and Module 4

– Defining the prediction task

– Define the evaluation metrics

– Designing features

– Designing inclusions and exclusion criteria for the data

– Annotating (hand-labelling) training (and testing) data

– Select the right model

– Error analysis

# Machine Learning for Design

Lecture 2
Introduction to Machine Learning. *Part 2*