# Unsupervised learning for materials informatics

GDS Short Course on Data Science for Physicists I

## Trevor David Rhone

Department of Physics, Applied Physics and Astronomy, Rensselaer Polytechnic Institute

# Google colaboratory excercise preparation

o Navigate to the GDS github page
  o https://github.com/quantum-intelligence/materials-informatics-tutorial
o Open Google colab
  - Load the google colab file associated with this tutorial
o Register with the materials project
  - https://next-gen.materialsproject.org/
  - Record your API key

# Introduction

o Challenges for materials discovery
- Predicting properties of materials with experiments (ab initio)
- Materials design
- Knowledge discovery

o Materials informatics
- Materials science + machine learning

o Challenges for materials informatics
- Lack of labelled data

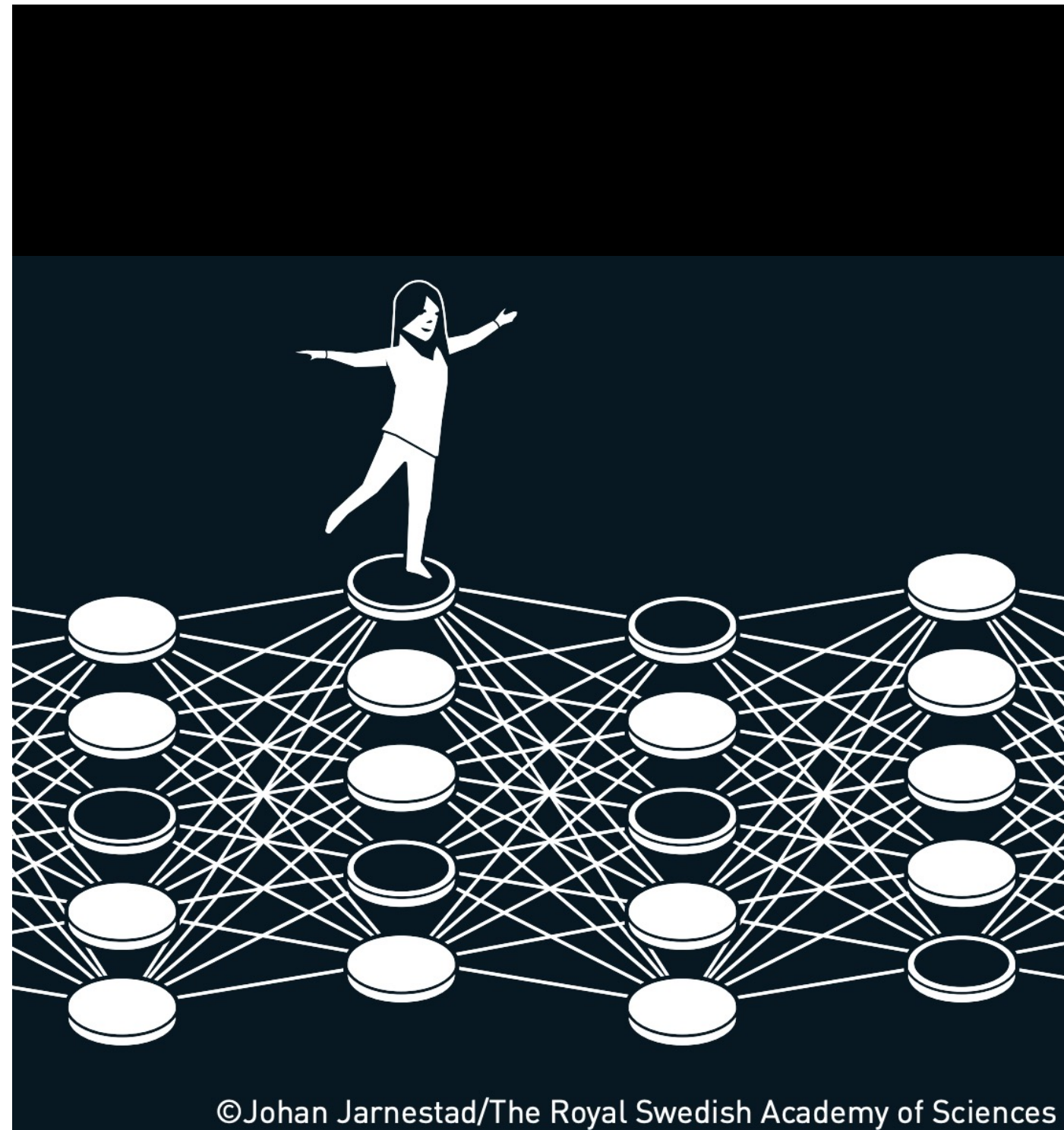# Nobel prize in physics goes to AI

**John J. Hopfield**
Princeton University, NJ, USA
**Geoffrey Hinton**
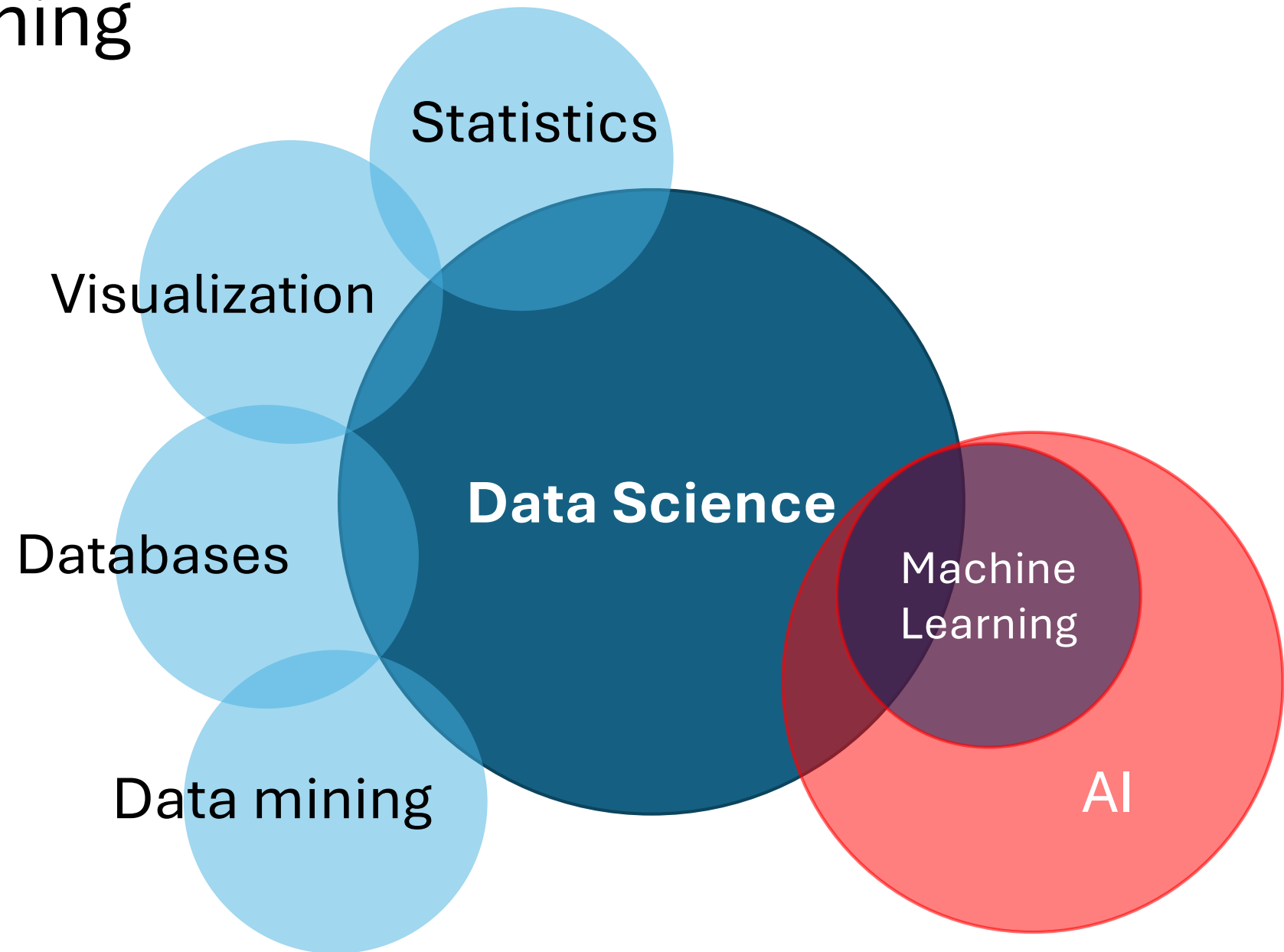University of Toronto, Canada

*"for foundational discoveries and inventions that enable machine learning with artificial neural networks"*



©Johan Jarnestad/The Royal Swedish Academy of Sciences

# Machine learning overview

# Machine learning

# Machine Learning in materials physics



o The rise of the materials databases
  • Data are accessible

o Chemical space descriptors exist
  • Coulomb Kernel[1], Bag of bonds representation[2]

o Datascience tools exist
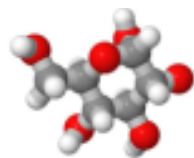  • Scikit learn, Google's TensorFlow, Pytorch

[1] M. Rupp, et al., Phys Rev Lett. 108, 058301 (2012) [2] K. Hansen et al., J. Phys. Chem. Lett. 6, 2326–2331 (2015)

# Machine Learning in materials physics
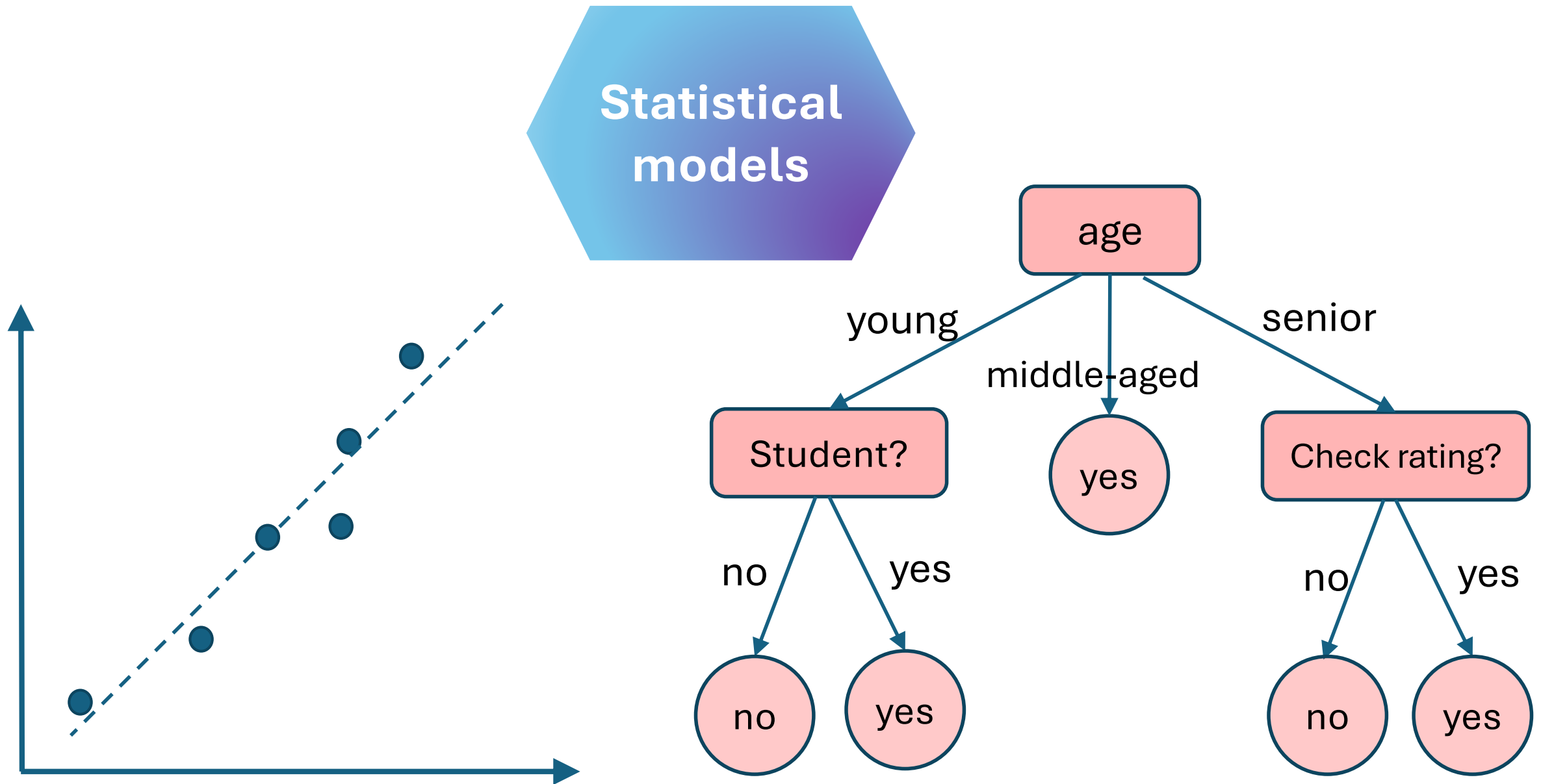
# Machine Learning in materials physics

**Descriptors**

o Description of material
- Atomic properties
- Mathematical representations of crystal structure

o Example:
- Ionic compounds have atoms in different columns of periodic table
  - Descriptor: column # of the periodic table

# Machine Learning in materials physics

# Machine Learning in materials physics

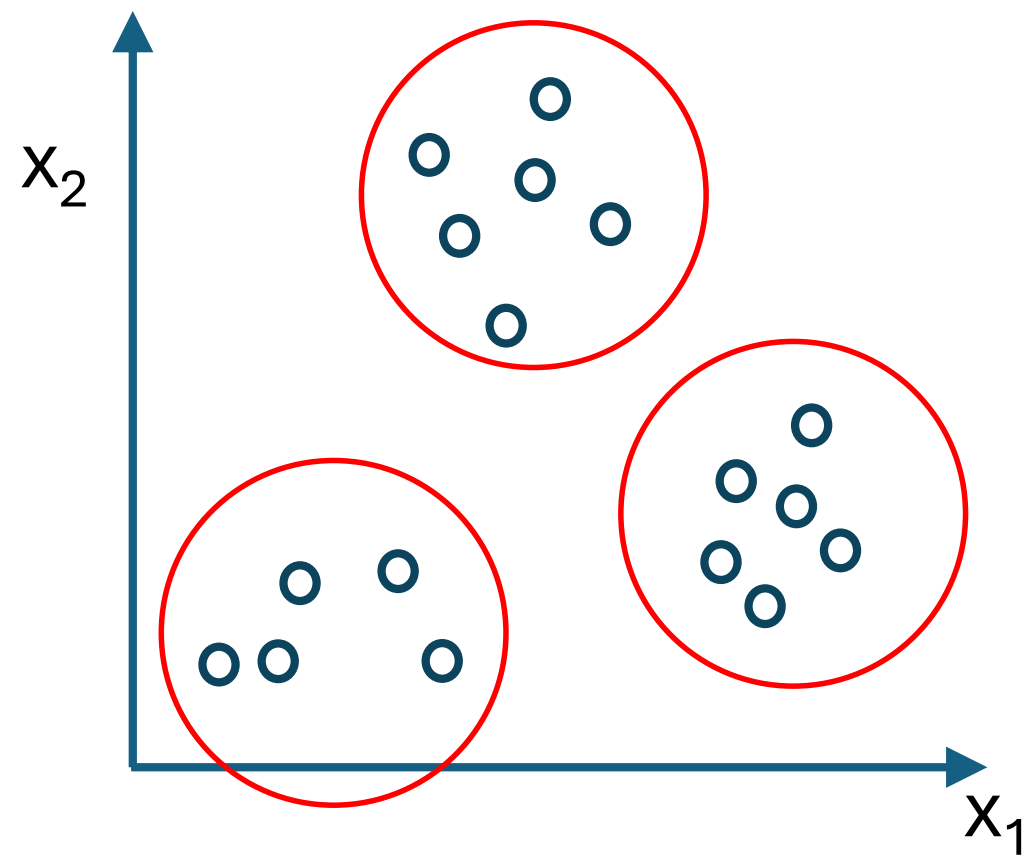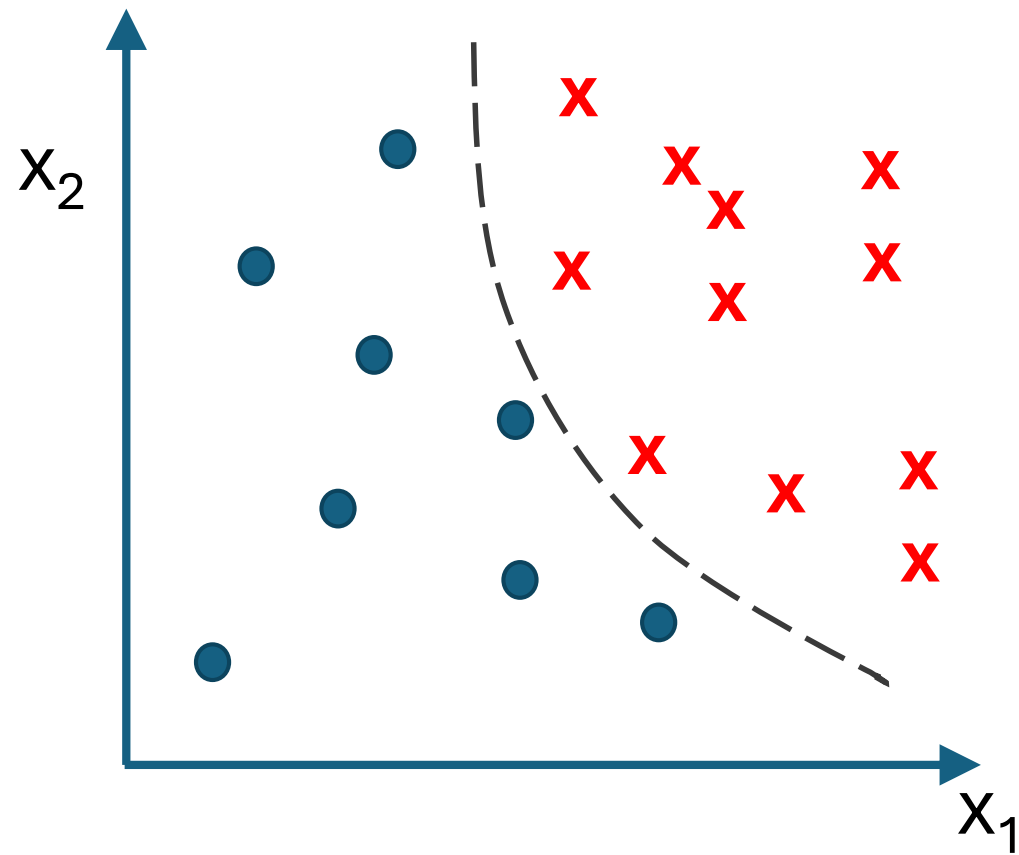Data science ~ data visualization + machine learning

$$y = f(x_1, x_2, ..., x_N) + \varepsilon$$

Target property

Inputs of machine learning model

Goal: learn or quantify some relationship

# Supervised versus Unsupervised learning

# Supervised versus Unsupervised learning

**Supervised learning**: For each observation of the predictor measurement(s) $x_i$, i = 1,...,n there can be an associated response measurement $y_i$.

We wish to fit a model that relates the response to the predictors to:
- *accurately predicting the response for future observations* (prediction) or
- better understanding the relationship between the response and the predictors (inference).

**Unsupervised learning** describes the situation in which for every observation i = 1,...,n, we observe a vector of measurements $x_i$ but no associated response $y_i$

# Unsupervised learning

o Dimensionality reduction
  - PCA (principal components analysis) looks to find a low-dimensional representation of the observations that explain a good fraction of the variance
  - tSNE (t-distributed stochastic neighbor embedding)

o Clustering looks to find homogeneous subgroups among the observations
  - K-means clustering
  - Hierarchical clustering

o Associative learning
  - Apriori algorithm
  - FP-growth algorithm

o Generative models
  - Variational autoencoder (VAE)
  - Generative adversarial network (GAN)
  - Transformer based generative models, eg. GPT

# Clustering

o Finds subgroups, or clusters, in a data set.

- Partition observations in a dataset into distinct groups so that
    1. observations within each group are similar, and
    2. observations in different groups are different

o Clustering could be used to find subgroups which can be mapped to subclasses of materials, for example:

- Nonmagnetic materials
- ferromagnets and
- antiferromagnets

o Types of clustering algorithms
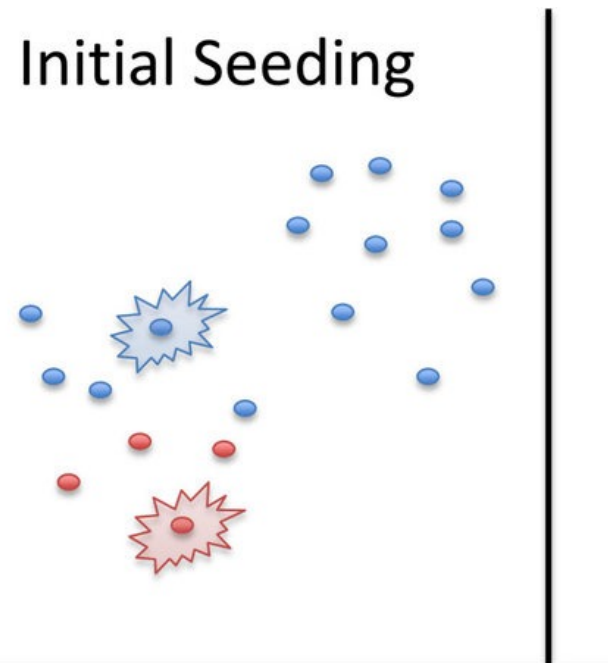
- K-means clustering
- Hierarchical clustering

# K-means Clustering

○ Partition observations into a pre-specified clustering number of clusters

○ Goal is to partition data set into K distinct, non-overlapping clusters
  - first specify the desired number of clusters K
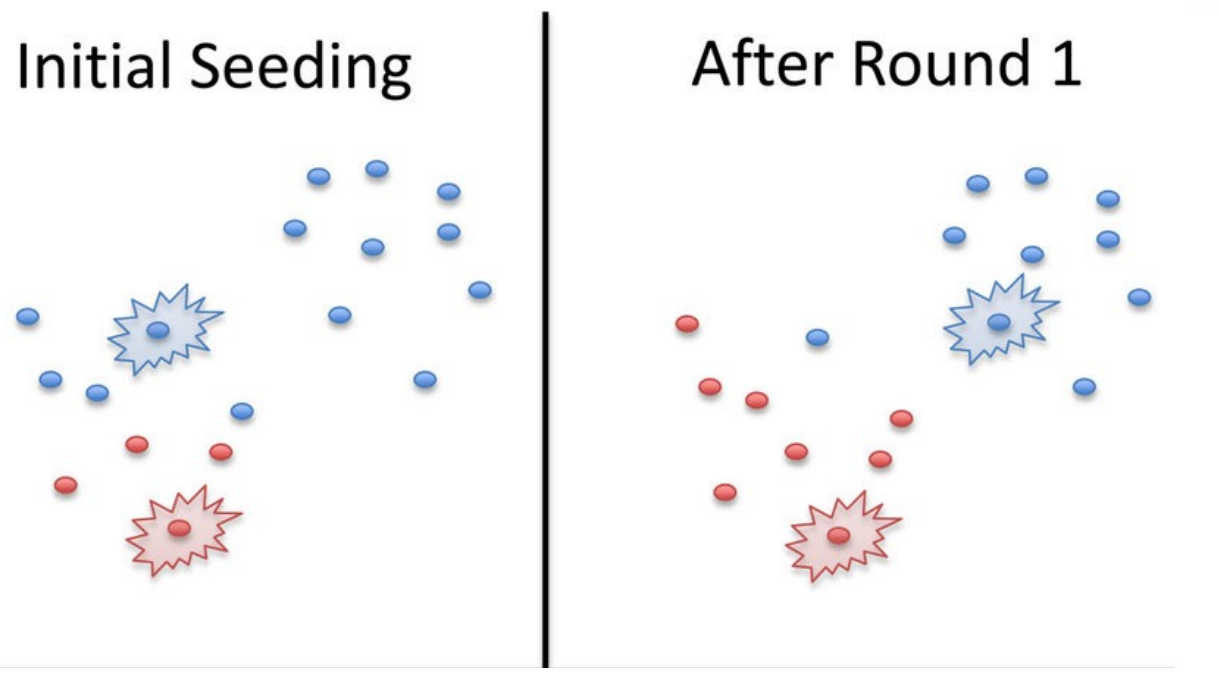  - the K-means algorithm assigns each observation to exactly one of the K clusters

$$\underset{C_1,\ldots,C_K}{\text{minimize}} \left\{ \sum_{k=1}^{K} W(C_k) \right\}$$

$$W(C_k) = \frac{1}{|C_k|} \sum_{i,i' \in C_k} \sum_{j=1}^{p} (x_{ij} - x_{i'j})^2$$

# K-means Clustering



Initial Seeding
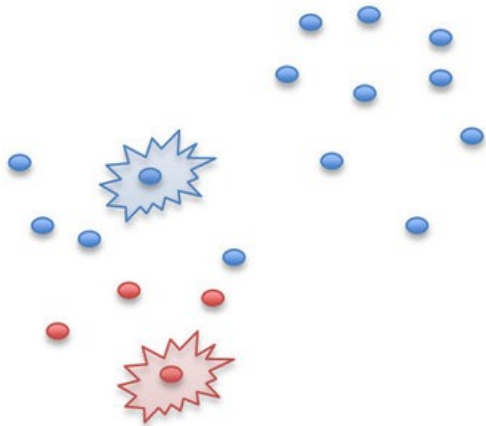
https://devopedia.org/k-means-clustering

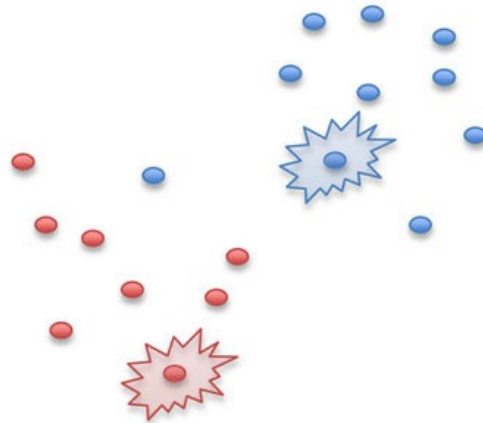# K-means Clustering
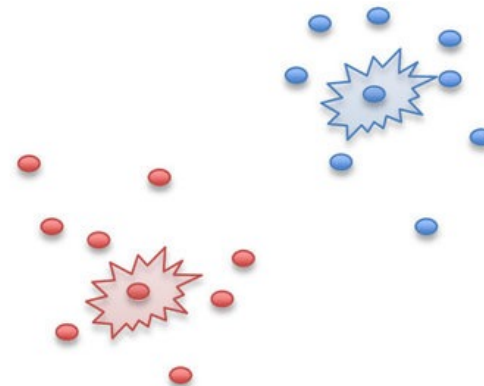
# K-means Clustering



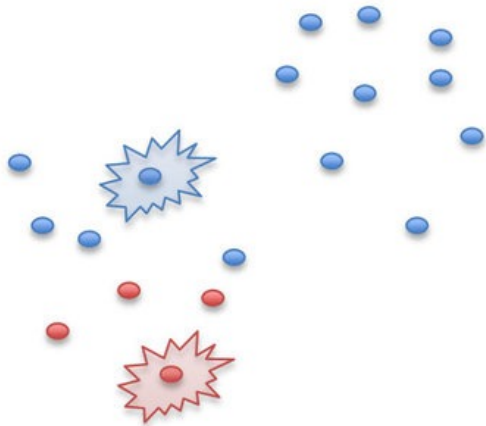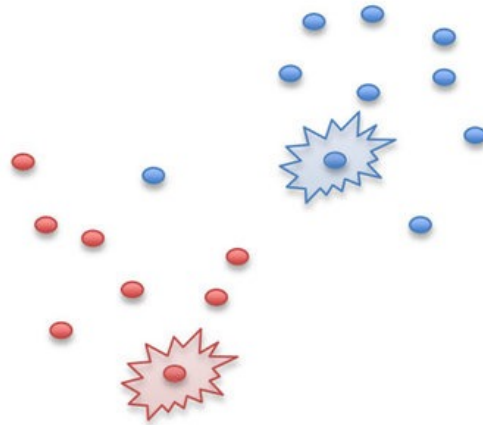Initial Seeding     After Round 1     After Round 2

# K-means Clustering



Initial Seeding     After Round 1     After Round 2     Final

https://devopedia.org/k-means-clustering

# K-means Clustering

- Hands-on example

# Dimensionality reduction

o Dimension Reduction
- Project p predictors into an M-dimensional subspace, where M < p
- These M projections can be used as predictors to fit a regression model
- M=2, 3 projections can be used to create visualizations that create insight into patterns in the data

- Dimensionality reduction methods
  - Principal components analysis (PCA)
  - t-distributed stochastic neighbor embedding (tSNE )
  - Autoencoders

# Dimension Reduction Methods

Principal Components Analysis

o Principal components summarize a data set with a smaller number of variables that collectively explain most of the variability in the original set

o PCA is an unsupervised approach, since it involves only a set of features $X_1, X_2, . . . , X_p$, and no associated response Y

o Uses:

    o Produces variables for use in supervised learning

    o Serves as a data visualization tool

# Principal Component Analysis



Two-dimensions

Three-dimensions

# Principal Component Analysis

# Considerations in High Dimensions

o Traditional statistical techniques for regression/classification are intended for the low-dimensional setting in which n (the # of observations) >> p (the # of features)

o In the past 20 years, new technologies have changed the way that data are collected. We now collect a large number of features, p.

o The number of observations n is often limited due to sample availability, etc.

o E.g., Rather than predicting blood pressure on the basis of age, gender, and BMI, one might *also* collect measurements for half a million low-dimensional single nucleotide polymorphisms (i.e. DNA mutations)

  - Then n ≈ 200 and p ≈ 500,000.

o Data sets containing more features than observations are often referred to as high-dimensional.

# Principal Component Analysis

- Hands-on example

# Associative rule learning

o Association rule learning is a rule-based ML method for discovering interesting relations between variables in large databases

o It identifies strong rules discovered in databases using several measures

o In any given transaction with a variety of items, association rules are meant to discover the rules that determine how certain items are connected

# Associative rule learning

| Algorithm | Use case | Efficiency | Cons |
|---|---|---|---|
| Apriori | Small the medium datasets | Slow on large datasets | High computational cost due to candidate generation |
| FP-Growth | Large datasets with complex relationships | Faster than Apriori | More complex |

# Associative rule learning

Apriori algorithm: identifies frequent itemsets and generates association rules from them

o Support
  • The proportion of materials that contain an item or itemset

o Confidence
  • The likelihood that item Y appears given that item X is already present

o Lift
  • Measures how much more likely Y is to appear with X compared to chance

$$\text{Support}(X) = \frac{\text{Transactions containing } X}{\text{Total transactions}}$$

$$\text{Confidence}(X \Rightarrow Y) = \frac{\text{Support}(X \cup Y)}{\text{Support}(X)}$$

$$\text{Lift}(X \Rightarrow Y) = \frac{\text{Confidence}(X \Rightarrow Y)}{\text{Support}(Y)}$$

Lift > 1: X and Y are positively correlated.
Lift = 1: No correlation.
Lift < 1: X and Y are negatively correlated.

# Associative rule learning

o Step 1: Find Frequent Itemsets

- Scan the dataset and find items that appear in at least a given support threshold
- Example: If we set **support ≥ 10%**, only elements appearing in **at least 10%** of materials are considered

o Step 2: Generate Candidate Itemsets

- Combine frequent single elements into pairs, then triplets, and so on
- Example:
  - Frequent **single elements**: {Ti}, {O}, {Fe}, {Ni}
  - Generate **pairs**: {Ti, O}, {Fe, O}, {Ni, O}, etc.

# Associative rule learning

o Step 3: Prune Non-Frequent Itemsets

•If a set is not frequent, its supersets cannot be frequent either.

•Example:
- If {Ti, O} is frequent, we check {Ti, O, Fe}.
- If {Ti, Fe} is not frequent, we discard {Ti, O, Fe}.

o Step 4: Generate Association Rules
- From frequent itemsets, we generate rules with high confidence.
- Example:
  - Rule: {Ti, O} → {High Band Gap}
  - Confidence: 80% (If a material contains Ti and O, it has a high band gap in 80% of cases).
  - Lift: 1.5 (Having Ti and O increases the likelihood of a high band gap).

# Associative rule learning

o Example for oxide materials and band gap

| Rule | Support | Confidence | Lift |
|---|---|---|---|
| {Ti, O} → {High band gap} | 0.25 | 0.80 | 1.5 |
| {Fe, O} → {Low band gap} | 0.20 | 0.70 | 1.3 |
| {Ti, O} → {High chemical stability} | 0.30 | 0.82 | 1.8 |

o Rule #1: If a materials contains Ti and O, there's an 80% probability that is has a high band gap

o Rule #2: If a material contains Fe and O. it's more likely to have a low band gap

# Associative rule learning

- Hands-on example

# Generative models

o Generative models
  - learn the underlying data distribution and generate new data samples that resemble the original dataset

o Discriminative models
  - predict labels, generative models create new data

| | Generative models | Discriminative models |
|---|---|---|
| Goal | Learn data distribution and generate materials | Learn decision boundary for classification |
| Examples | VAEs, GANs, Diffusion models | Decision Trees, SVM, Neural networks |
| Output | New materials | Labels or class probabilities |

# Types of Generative Models

- Probabilistic Generative Models
- Deep Learning-Based Generative Models
- Transformer-Based Generative Models

# Deep Learning-Based Generative Models

Neural networks used to learn complex data distributions

o Variational Autoencoders (VAEs)
  - Encodes data into a latent space, then reconstructs it.
  - Used in drug discovery, and materials discovery

o Generative Adversarial Networks (GANs)
  - Comprises a Generator (creates fake data) and a Discriminator (distinguishes real from fake)
  - Used for microscopy image synthesis, material microstructure simulation, and molecule generation

o Diffusion Models
  - Adds noise to data, then learns to denoise the data
  - Used to generate novel molecules and crystal structures
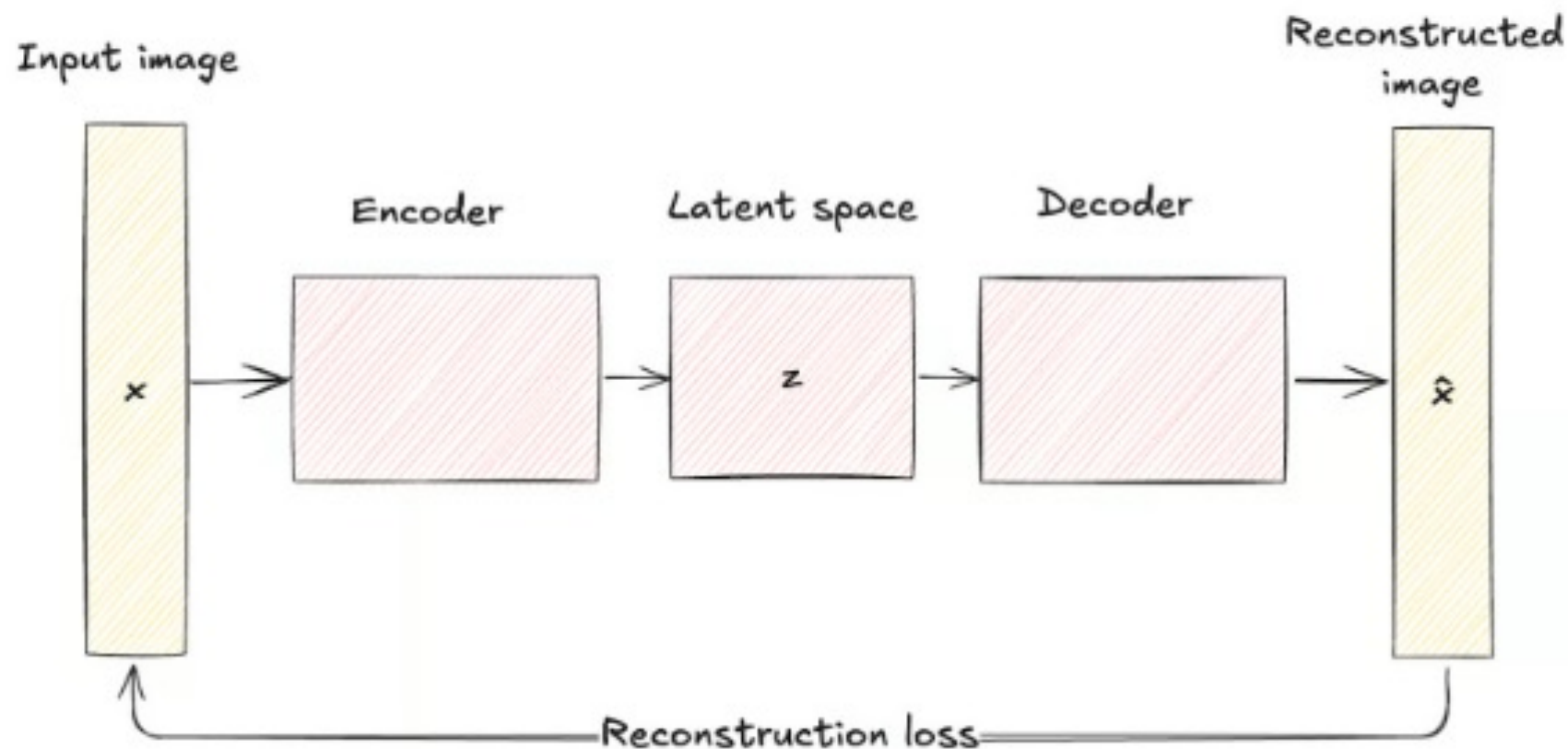
# Transformer-Based Generative Models

These use self-attention mechanisms to model long-range dependencies

o GPT (Generative Pre-trained Transformer)
- Used for text generation, code generation, and language modeling
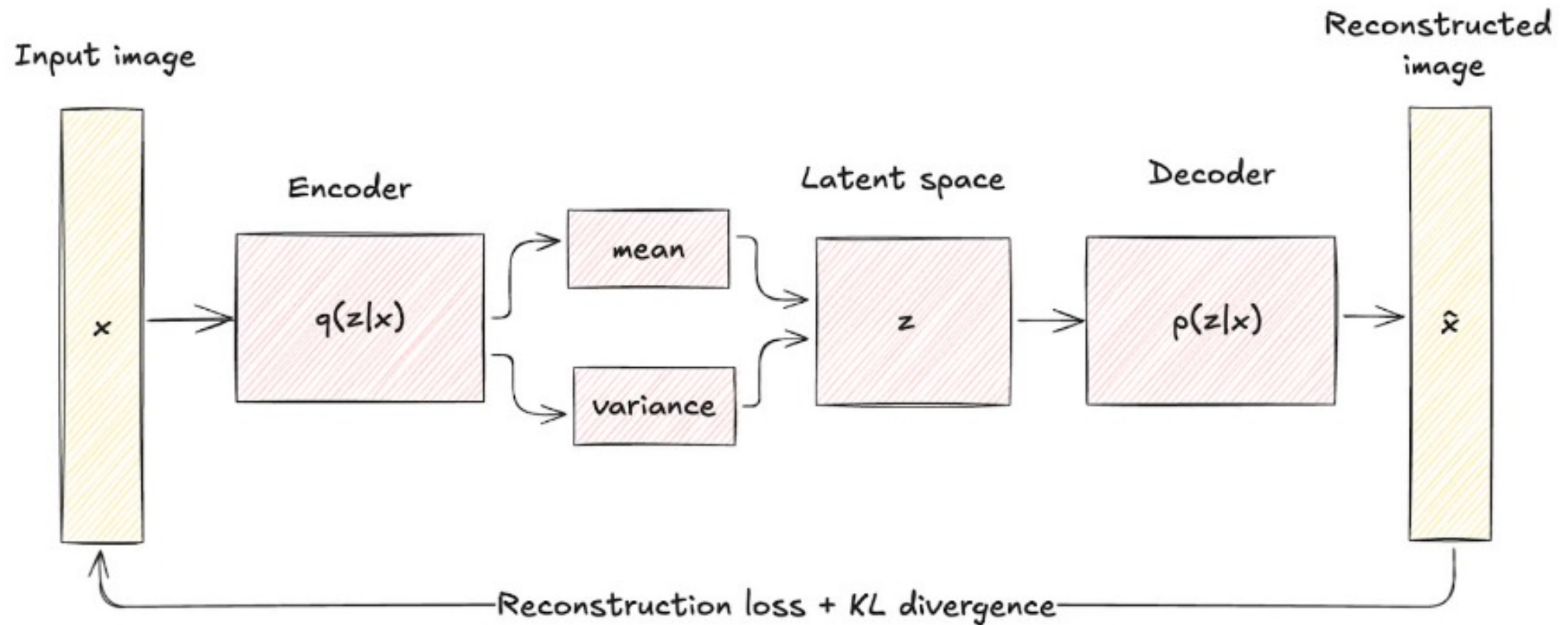- Example: data-mining materials science research papers

# Variational autoencoder (VAE)

o The VAE architecture
- Similar to the autoencoder architecture



https://www.datacamp.com/tutorial/variational-autoencoders

# Variational autoencoder (VAE)

o The VAE architecture

# Variational autoencoder (VAE)

o The VAE architecture

The VAE introduces a probabilistic element into the encoding process
- The VAE encoder maps the input data to a probability distribution over the latent variables
- This is modeled as a Gaussian distribution with mean $\mu$ and variance $\sigma^2$

- VAEs incorporate regularization through the KL divergence
- The KL divergence makes the latent space continuous and well-structured
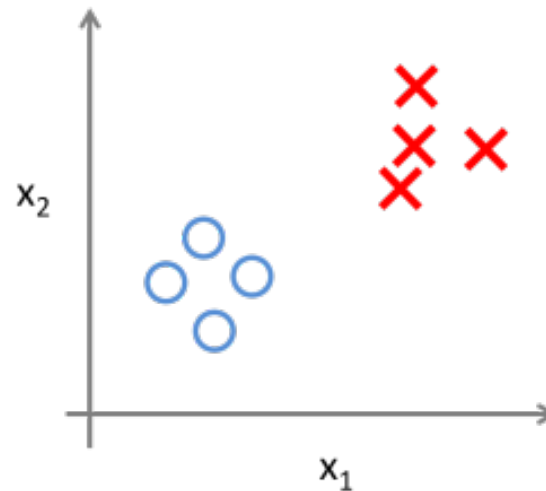
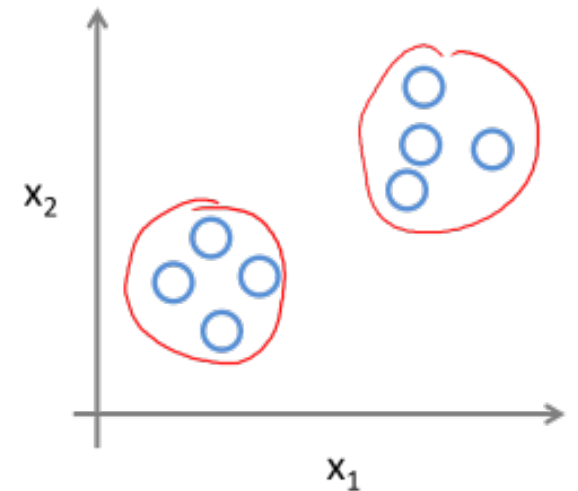# Generative Models

o Hands-on example

# Summary

Unsupervised learning can identify patterns in data

- Clustering
  - K-means clustering
- Dimensionality reduction
  - PCA
- Associative rule learning
  - Apriori
- Generative models
  - Variational autoencoders

# Outlook

○ Semi-supervised learning

- Unsupervised + supervised learning
- Leverage unlabeled data to improve supervised learning tasks

**Artificial Intelligence Guided Studies of van der Waals Magnets**

T. D. Rhone *et al.*, "Artificial Intelligence Guided Studies of van der Waals Magnets," *Adv. Theory Simulations*, vol. 6, no. 6, p. 2300019, 2023.

# Hands-on project

o $CrGeTe_3$-type vdW magnets dataset

o Identify patterns in data using unsupervised learning

o $CrGeTe_3$ materials in this dataset can be ferromagnetic (FM) or antiferromagnetic (AFM)

o Can you find clusters of FM and AFM materials?

- Data exploration
  - Data visualization
  - PCA
- K-means clustering
- How to choose materials descriptors?

T. D. Rhone *et al.*, "Data-driven studies of magnetic two-dimensional materials," *Sci. Rep.*, vol. 10, no. 1, p. 15795, 2020.

# Want more information?

www.materials-informatics.com

- o GDS Webinar series on YouTube
- o Jarvis Tutorials by NIST
- o Andrew Ng's Machine Learning course on Coursera

# Apply to Rensselaer Polytechnic Institute (RPI)

Patterns emerge free,
atoms whisper their secrets,
no guide, yet they learn.