

# Introduction to Machine Learning

## Lecture 1: Preliminaries and Clustering

Ali Bereyhi

`ali.bereyhi@utoronto.ca`

Department of Electrical and Computer Engineering  
University of Toronto

Winter 2026

# What is Machine Learning?

It's a hard question to answer *accurately*

Mitchel defines ML as “... the study of computer algorithms that improve automatically through experience...”

and

Goodfellow et al. *informally* define ML as “... a form of applied statistics with increased emphasis on the use of computers to statistically estimate complicated functions...”

and ...

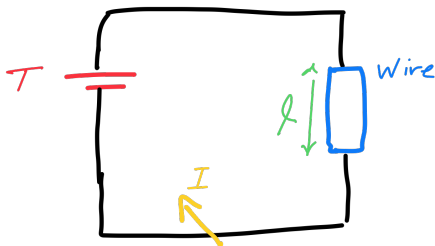
# What is Machine Learning?

But not too hard to answer *practically*

We define ML as *the set of data-driven approaches that help us understand the environment and its behavior, and generalize it!*

*Data-driven* approaches have long been with us in science and engineering!

# Early Example from 1827: *Ohm's Law*



# What Did Ohm Do?

*Georg Ohm did three major steps*

- He saw a pattern and hypothesized some mathematical model
  - ↳ *Electric current increases with voltage*
  - ↳ *The constant changes with the length and material*
  - ↳ ...
- He collected data
  - ↳ *Electric currents and voltages*
- He used mathematical tools to extract the modeled pattern
  - ↳ *Some curve-fitting technique*

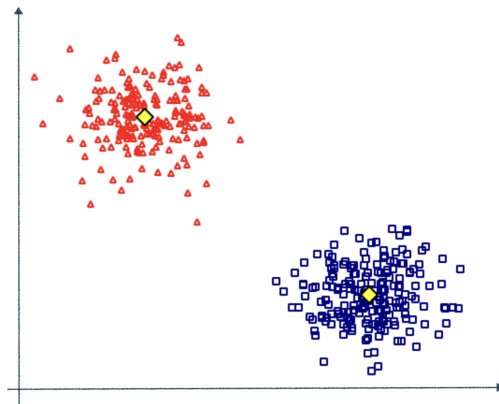
# Learning Task

*Any learning task has three components*

- *Model that captures the Pattern*
- Data
- Learning Algorithm

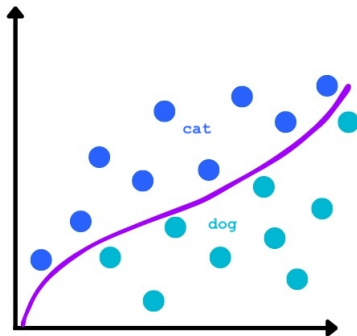
## Example: Clustering

*Monthly amount of transactions* versus *# of transactions per month*



## Example: Classification

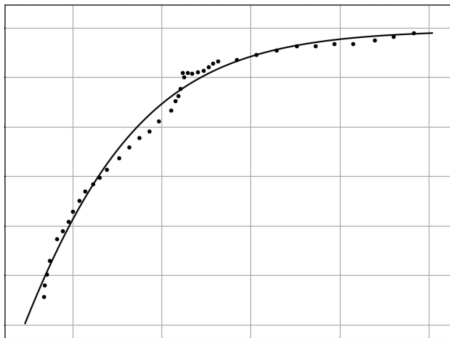
*Sleep time* versus *# of times the pet makes noise*



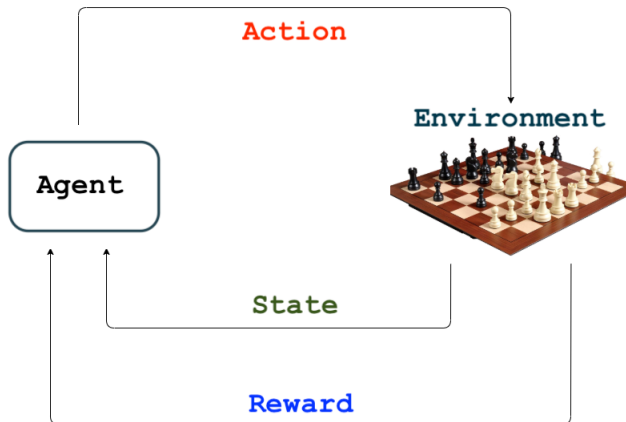


# Example: Regression

*Salary* versus *years of experience*



## Example: *Playing Chess*



# Dataset

A set of data samples

$$\mathbb{D} = \{\mathbf{x}_n : n = 1, \dots, N\}$$

with  $\mathbf{x}_n \in \mathbb{R}^d$

---

*Let's think about data in our examples*

- Clustering
- Classification
- Regression
- Playing Chess

# Model

A pre-assumed function

$$f : x \mapsto y$$

for a data sample  $x$  and **output**  $y$  that *fits the learning task*

---

*Let's formulate model in our examples*

- Clustering
- Classification
- Regression
- Playing Chess

# Learning Algorithm

Algorithm that gets dataset and returns the **exact model**

$$\mathcal{A} : \mathbb{D} \mapsto f^*$$

$f^*$  does the mapping such that we get to the **desired output**

---

*Let's formulate learning algorithm in our examples*

- Clustering
- Classification
- Regression
- Playing Chess

---

? *How can we define a “good” learning algorithm?*

# Unsupervised Learning

Data samples are **not** **labeled**

$$\mathbb{D} = \{\mathbf{x}_n : n = 1, \dots, N\}$$

*Here, we are looking for a pattern in the data*

---

*Other components of an unsupervised task*

- *Model captures the pattern hidden in data*
- Learning Algorithm

---

*Examples of unsupervised learning*

- ✓ Clustering
  - Dimensionality Reduction
  - Distribution Learning

# Supervised Learning

Data samples are **labeled**

$$\mathbb{D} = \{(\mathbf{x}_n, \mathbf{v}_n) : n = 1, \dots, N\}$$

*Here, we are looking for a model that describes the relation*

---

*Other components of a supervised task*

- Model describes the **relation** between data samples and their **labels**
- Learning Algorithm

---

*Examples of supervised learning*

- Classification
- Regression

# Reinforcement Learning

Data samples are **series** of **actions**, **states** and **rewards**

$$\mathbb{D} = \left\{ \left\{ \left( a_n^t, s_n^t, r_n^t \right) : t = 1, \dots \right\} : n = 1, \dots, N \right\}$$

*Here, we are looking for optimal policy, i.e., policy that maximizes future returns*

$$G_t = r^t + r^{t+1} + \dots$$

*Other components of a reinforcement task*

- Model describes a **policy**
- Learning Algorithm

*Examples of reinforcement learning*

- Playing Game, Control Robots, . . .

**Reinforcement learning** is **not** discussed in this course, but you may consider taking **Reinforcement Learning** in **next Fall**



## Further Read

- Bishop
  - ↳ Chapter 1: *Sections 1.1 and 1.3* **Introductory**
- ESL
  - ↳ Chapter 1 **Introductory**
  - ↳ Chapter 2: *Sections 2.1 and 2.2* **Supervised**
  - ↳ Chapter 14: *Sections 14.1 and 14.2* **Unsupervised**
- Mitchell
  - ↳ Chapter 13: *Sections 13.1 and 13.2* **Reinforcement**
- Goodfellow, et al.
  - ↳ Chapter 5: *Sections 5.1 and 5.2* **Introductory**

# Unsupervised Learning

Why do we start with *unsupervised learning*?

- Many **basic** problems are **unsupervised**
  - ↳ We naturally **cluster** everything around us
  - ↳ We get sense about quantities by understanding its **statistical behavior**
- It helps us to recap some basics we need later
  - ↳ *Linear Algebra*
  - ↳ *Probability Theory*

# Problem of Clustering

This is a basic sign of **intelligence**

- *We cluster everything around us*
  - ↳ *Trees, flowers, animals, . . .*
- *We often start with simple clustering and extend hierarchically*
  - ↳ *Plants and animals*
  - ↳ *Plants could be trees, flowers, . . .*
  - ↳ *Animals could be mammals, birds, . . .*
- *The further we go, the more intelligent we get!*

## Basic Clustering Task: *Data*

Data samples are points in  $d$ -dimensional space

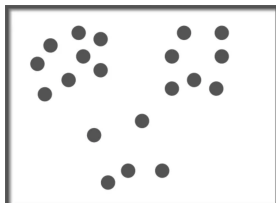
$$\mathbb{D} = \{\mathbf{x}_n : n = 1, \dots, N\}$$

with  $\mathbf{x}_n \in \mathbb{R}^d$

### In Examples

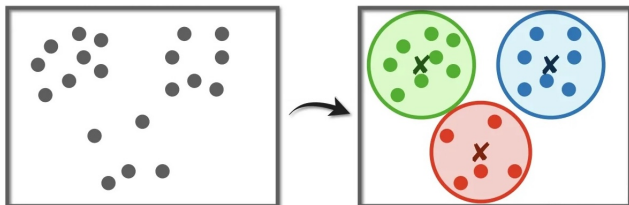
*In examples, we always think of two dimensions for sake of simplicity*

*Recall our bank record example*



## Basic Clustering Task: *Pattern*

We assume that the samples can be grouped into clusters



*Recall our bank record example*

## Basic Clustering Task: *Model*

We use a model to capture the clustering pattern

$$f(\mathbf{x}) \rightarrow k \in \{1, \dots, K\}$$

for some integer number of clusters  $K$

---

Some definitions and assumptions

- *Cluster subspace  $k$*

$$\mathbb{C}_k = \{\mathbf{x} : f(\mathbf{x}) = k\}$$

- *Cluster subspaces partition the data space*

$$\mathbb{C}_1 \cup \dots \cup \mathbb{C}_K = \mathbb{X} \rightsquigarrow \text{all possible samples}$$

$$\mathbb{C}_j \cap \mathbb{C}_k = \emptyset \rightsquigarrow \forall j \neq k$$

# Basic Clustering Task: *Learning Algorithm*

The learning algorithm gets a dataset and finds a **good**  $f$

$$\mathcal{A} : \mathbb{D} \mapsto f^*$$

? What is a “**good**” model?

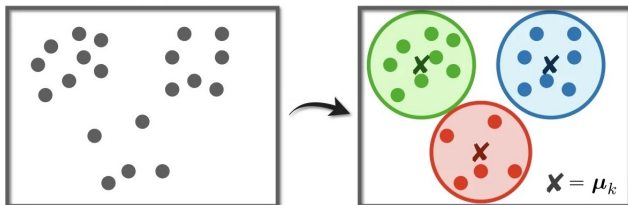
! We'll answer it!

## An Intuitive Model: $K$ Centroids

Let's use a simple and intuitive model

$$f(\mathbf{x}) = \operatorname{argmin}_{k \in \{1, \dots, K\}} \|\mathbf{x} - \boldsymbol{\mu}_k\|$$

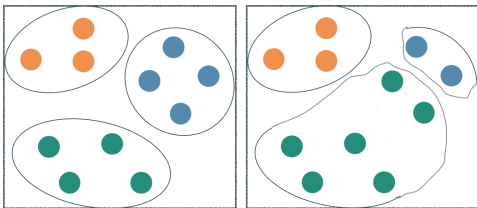
for  $K$  centroids  $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K \in \mathbb{R}^d$





## $K$ Centroids: Learning Algorithm

*This model is valid for any set of centroids!*



The learning algorithm is to start from  $\mathbb{D}$  and **learn good** centroids

$$\mathcal{A} : \mathbb{D} \mapsto \mu_1^*, \dots, \mu_K^*$$

? What is a “good” set of centroids?

! We'll answer it!

## *K*-Means Clustering Algorithm: *Intuitive Derivation*

*Given the centroids, we can easily assign each  $x_n \in \mathbb{D}$  to a cluster-set*

```
Cluster_Assignment( $\mu_1, \dots, \mu_K$ ):
```

```
  # we want to find  $\mathcal{C}_1 \cup \dots \cup \mathcal{C}_K = \mathbb{D}$ 
```

```
1: for  $n = 1 : N$  do
```

```
2:   Find the index of the closest centroid
```

$$k^* \leftarrow \operatorname{argmin}_{k \in \{1, \dots, K\}} \|x_n - \mu_k\|$$

```
3:   Assign  $x_n$  to cluster-set  $\mathcal{C}_{k^*}$ 
```

```
4: end for
```

```
5: Return  $\mathcal{C}_1, \dots, \mathcal{C}_K$ 
```

# *K*-Means Clustering Algorithm: *Intuitive Derivation*

*Given the cluster sets, we can move centroids to the center of cluster-sets*

Centroid\_Update( $\mathcal{C}_1, \dots, \mathcal{C}_K$ ):

# we want to find  $\mu_1, \dots, \mu_K$

1: **for**  $k = 1 : K$  **do**

2:   **if**  $\mathcal{C}_k \neq \emptyset$  **then**

3:     Move  $\mu_k$  to the center of cluster  $\mathcal{C}_k$ , i.e.,

$$\mu_k \leftarrow \frac{1}{|\mathcal{C}_k|} \sum_{\mathbf{x}_n \in \mathcal{C}_k} \mathbf{x}_n$$

4:   **else**

5:     Leave  $\mu_k$  unchanged

6:   **end if**

7: **end for**

8: Return  $\mu_1, \dots, \mu_K$

# K-Means Clustering Algorithm

*We could iterate till we converge*

*K-Means()* :

- 1: Initiate  $\mu_1, \dots, \mu_K$
- 2: **while**  $\mu_1, \dots, \mu_K$  changing **do**
- 3:   Set  $\mathcal{C}_1, \dots, \mathcal{C}_K \leftarrow \text{Cluster\_Assignment}(\mu_1, \dots, \mu_K)$
- 4:   Update  $\mu_1, \dots, \mu_K \leftarrow \text{Centroid\_Update}(\mathcal{C}_1, \dots, \mathcal{C}_K)$
- 5: **end while**
- 6: Return  $\mu_1, \dots, \mu_K$

## Example: 2-Means Clustering<sup>1</sup>

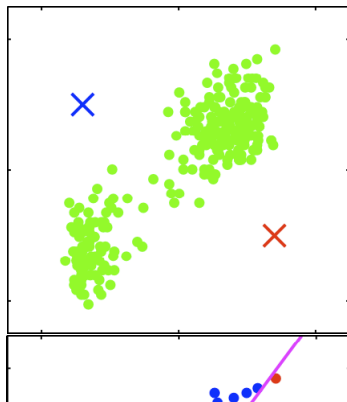
*Initial centroids*

*Iteration 1*

*Iteration 2*

*Iteration 3*

*Iteration 4 – Converged*



# K-Means Clustering Algorithm: Alternative Formulation

Cluster\_Assignment( $\mu_1, \dots, \mu_K$ ):

1: **for**  $n = 1 : N$  **do**

2:   Assign  $K$  weights  $r_{n,1}, \dots, r_{n,K}$  to sample  $x_n$  as

$$r_{n,k} = \begin{cases} 1 & \text{if } \mu_k = \text{closest centroid to } x_n \\ 0 & \text{otherwise} \end{cases}$$

3: **end for**

4: Return  $r_{n,k}$  for  $k = 1 : K$  and  $n = 1 : N$

*Properties of  $r_{n,k}$*

$$\sum_{k=1}^K r_{n,k} = 1 \quad \text{and} \quad \sum_{n=1}^N r_{n,k} = |\mathcal{C}_k|$$

# *K*-Means Clustering Algorithm: Alternative Formulation

Centroid\_Update( $\{r_{n,k}\}$ ):

- 1: **for**  $k = 1 : K$  **do**
- 2:   **if**  $\sum_n r_{n,k} > 0$  **then**
- 3:     Move  $\mu_k$  to the center of cluster  $k$ , i.e.,

$$\mu_k = \frac{\sum_{n=1}^N r_{n,k} \mathbf{x}_n}{\sum_{n=1}^N r_{n,k}}$$

- 4:   **else**
- 5:     Leave  $\mu_k$  unchanged
- 6:   **end if**
- 7: **end for**
- 8: Return  $\mu_1, \dots, \mu_K$

# K-Means Clustering Algorithm: Alternative Formulation

*We could iterate till we converge*

*K-Means()* :

- 1: Initiate  $\mu_1, \dots, \mu_K$
- 2: **while**  $\mu_1, \dots, \mu_K$  changing **do**
- 3:   Set  $\{r_{n,k}\} \leftarrow \text{Cluster\_Assignment}(\mu_1, \dots, \mu_K)$
- 4:   Update  $\mu_1, \dots, \mu_K \leftarrow \text{Centroid\_Update}(\{r_{n,k}\})$
- 5: **end while**
- 6: Return  $\mu_1, \dots, \mu_K$

*This is a better form to extend K-means clustering to a **soft** format*



# Defining Objective: *Risk*

? What is a “good” set of centroids?

! We'll answer it!

We may define a metric to evaluate how our model performs

$$\mathcal{J}(\{r_{n,k}\}, \{\mu_k\}) = \frac{1}{N} \sum_{k=1}^K \sum_{n=1}^N r_{n,k} \|\mathbf{x}_n - \mu_k\|^2$$

This specifies the **risk** we take with this model

# Notion of Optimality

? What is a “good” set of centroids?

! We'll answer it!

## Optimal Clustering

Optimal assignments  $\{r_{n,k}^*\}$  and centroids  $\{\mu_k^*\}$  minimize the *risk*

$$\{r_{n,k}^*\}, \{\mu_k^*\} = \operatorname{argmin}_{\{r_{n,k}\}, \{\mu_k\}} \mathcal{J}(\{r_{n,k}\}, \{\mu_k\})$$

# *K*-Means Clustering: Risk Minimization

*Risk minimization for clustering is hard, so we can use alternating optimization*

Risk\_Minimization():

- 1: Initiate  $\mu_1^*, \dots, \mu_K^*$
- 2: **while**  $\mu_1^*, \dots, \mu_K^*$  changing **do**
- 3:   Minimize the risk for fixed centroids  $\{\mu_k^*\}$

$$\{r_{n,k}^*\} = \operatorname{argmin}_{\{r_{n,k}\}} \mathcal{J}(\{r_{n,k}\}, \{\mu_k^*\})$$

- 4:   Minimize the risk for fixed assignments  $\{r_{n,k}^*\}$

$$\{\mu_k^*\} = \operatorname{argmin}_{\{\mu_k\}} \mathcal{J}(\{r_{n,k}^*\}, \{\mu_k\})$$

- 5: **end while**
- 6: Return  $\mu_1^*, \dots, \mu_K^* \approx \mu_1^*, \dots, \mu_K^*$

# $K$ -Means Clustering: Risk Minimization

Minimize the risk for fixed centroids  $\{\mu_k^*\}$

$$\{r_{n,k}^*\} = \underset{\{r_{n,k}\}}{\operatorname{argmin}} \mathcal{J}(\{r_{n,k}\}, \{\mu_k^*\})$$

which is done by

Cluster\_Assignment( $\mu_1, \dots, \mu_K$ ):

1: **for**  $n = 1 : N$  **do**

2: Assign  $K$  weights  $r_{n,1}, \dots, r_{n,K}$  to sample  $x_n$  as

$$r_{n,k} = \begin{cases} 1 & \text{if } k = \operatorname{argmin}_{j \in \{1, \dots, K\}} \|x_n - \mu_j\| \\ 0 & \text{otherwise} \end{cases}$$

3: **end for**

4: Return  $r_{n,k}$  for  $k = 1 : K$  and  $n = 1 : N$

## $K$ -Means Clustering: Risk Minimization

Minimize the risk for fixed assignments  $\{r_{n,k}^*\}$

$$\{\mu_k^*\} = \underset{\{\mu_k\}}{\operatorname{argmin}} \mathcal{J}(\{r_{n,k}^*\}, \{\mu_k\})$$

which is done by

Centroid\_Update( $\{r_{n,k}\}$ ):

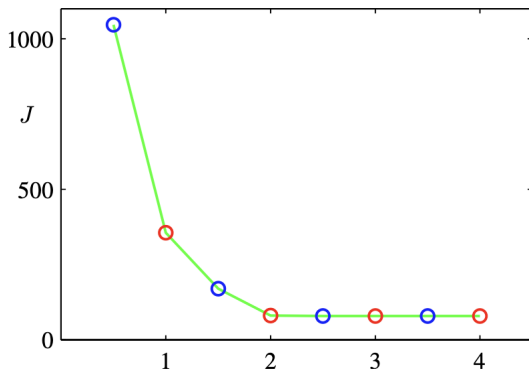
```
1: for  $k = 1 : K$  do
2:   if  $\sum_n r_{n,k} > 0$  then
3:     Move  $\mu_k$  to the center of cluster  $k$ 
4:   else
5:     Leave  $\mu_k$  unchanged
6:   end if
7: end for
8: Return  $\mu_1, \dots, \mu_K$ 
```

## $K$ -Means Clustering $\equiv$ Risk Minimization

So we conclude

$$\text{Risk\_Minimization}() \equiv K\text{-Means}()$$

*Back to our binary example*



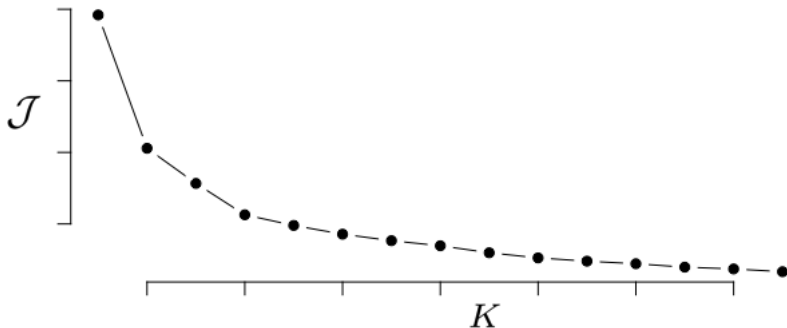
## More Sophisticated Example: Segmentation<sup>2</sup>

Each RGB pixel is a sample  $x_n \in \mathbb{R}^3$ : we cluster with  $K = 10$   $K = 3$   $K = 2$



## Choice of *Hyperparameter*

- ? How do we know  $K$ ?
- ! This is a *hyperparameter*



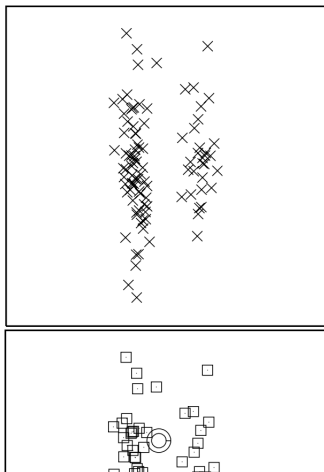


# $K$ -Means Clustering Always Converge

? Does  $K$ -means clustering always **converge** to a **stable state**?

! **Yes!** You can show it!

However, it does **not necessary** end with **what we want!**<sup>3</sup>



## Further Read

- MacKay
  - ↳ Chapter 20 *K-means*
- Bishop
  - ↳ Chapter 9: *Section 9.1* *K-means*
- ESL
  - ↳ Chapter 14: *Section 14.3* *Clustering*