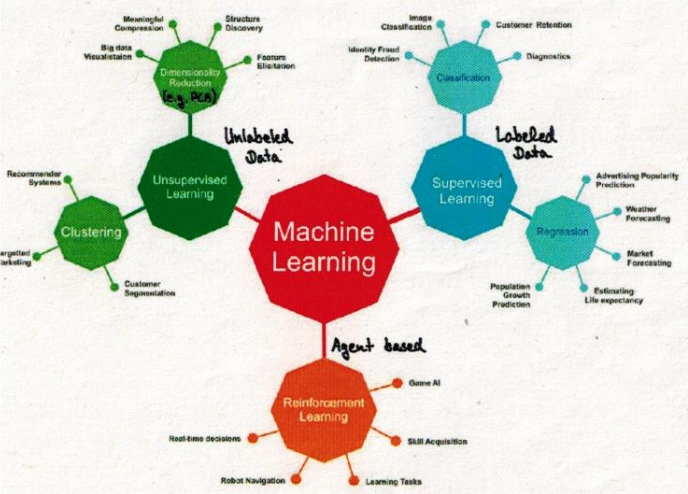# Machine Learning Basics



**Assumption:** Observed data $D$ drawn _i.i.d._ from unknown data-generating process $\mathbb{P}_{xy}$ — independent rows in dataset

(Not always realistic, scenarios like time series are clearly not iid)

## Data in Supervised Learning

ML (typically): prediction more important than explaining — functional ≠ intrinsic pattern

Goal: Discover predictive rule behind some (assumed) relationship between target and features

Row 1 and $y_1$
Feature-Vector 1 $(\underline{x}^{(1)}, y^{(1)})$
labeled data [training]
unlabeled data [predict]

| Features $x$ | | | | Target $y$ |
|---|---|---|---|---|
| SepalLength | SepalWidth | PetalLength | PetalWidth | Species |
| 4.3 | 3.0 | 1.1 | 0.1 | setosa |
| 5.0 | 3.3 | 1.4 | 0.2 | setosa |
| 7.7 | 3.8 | 6.7 | 2.2 | virginica |
| 5.5 | 2.5 | 4.0 | 1.3 | versicolor |
| 5.9 | 3.0 | 5.1 | 1.8 | ? |
| 4.4 | 3.2 | 1.3 | 0.2 | ? |

labels

$n$-Zeilen — $x_2$ Column 2 / Feature 2 — $p$-Merkmale

**Target variable types:**
- Numerical ($\mathbb{R}$)
- Integer ($\mathbb{Z}$)
- Categorical ($\{C_1, ..., C_g\}$)
- Binary ($\{0,1\}$) $\subset \{1, -1\}$

$\left.\begin{array}{l} \text{Regression task} \\ \text{Classification task} \end{array}\right\}$ Supervised task

then Learner called "classifier"

## Models & Parameters

In most Regression tasks $g = 1$. In Classification tasks $g = $ # Categories and $f$ is e.g. a score or class probability

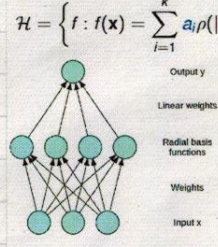A function $f: \mathcal{X} \to \mathbb{R}^g$ is called a **Model** (or **Hypothesis**)

**! ML requires constraining $f$ to a certain type of function which we call structural prior** (e.g. design choice is linear functions) — Not after closely looking at data! choice upfront, usually automatically from Learner choice

The set $\mathcal{H} := \{f \mid f \text{ belongs to the class of the structural prior}\}$ is called **Hypothesis space / Model class**

Usually $\mathcal{H}$ is constructed as a parametrised family of curves with parameters $\underline{\theta} := (\theta_1, ..., \theta_d) \in \Theta$ **Parameter space** $\mathbb{R}^d$

$\Rightarrow \mathcal{H} = \{f_\theta \mid f_\theta \text{ belongs to a family of curves parametrised by } \theta\}$

training dataset — hyperparameter control setting

Inducer, Learner (-Algorithm) $\mathcal{I}: \mathbb{D} \times \Lambda \to \mathcal{H}$ bzw. $\Theta$, $(D, \lambda) \mapsto \hat{f}_{D, \lambda}$

Therefore: **! Finding the optimal model is equivalent to finding the optimal parameters** $\Rightarrow$ picks best element of hypothesis space for given training data

**! The parameter-to-model mapping can be non-injective, i.e. one model is described by different parameter vectors.**

↳ evaluates $f \in \mathcal{H} \mid \theta$ on $D$ with risk function $R_{emp}(f)$ or $R_{emp}(\theta)$

↳ uses optimization method to find $\hat{f}$ or $\hat{\theta}$, the model with the lowest risk

## Example for Hypothesis / Parameter Spaces

Learning components: $\mathcal{H}$, risk, optimization

### Bivariate quadratic function:

$\mathcal{H} = \{f : f(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2 + \theta_5 x_1 x_2, \theta \in \mathbb{R}^6\}$,



### RBF Network:

$\mathcal{H} = \left\{ f : f(\mathbf{x}) = \sum_{i=1}^{k} a_i \rho(\|\mathbf{x} - c_i\|) \right\}$,



Output y
Linear weights
Radial basis functions
Weights
Input x

$a_i :=$ weight of $i$-th neuron

$c_i :=$ its center vector

$\rho(\|\vec{x} - c_i\|) := \exp(-\beta \cdot \|\vec{x} - c_i\|^2)$ is the $i$-th **radial basis** function with bandwidth $\beta \in \mathbb{R}$.

## Notation

usual $\mathcal{X} \subset \mathbb{R}^p$

Input/Feature Space: $\mathcal{X}$ ] All data sets of size $n$: $\mathbb{D}_n := (\mathcal{X} \times \mathcal{Y})^n$

Output/Target Space: $\mathcal{Y}$ ] All finite data sets: $\mathbb{D} := \bigcup_{n \in \mathbb{N}} (\mathcal{X} \times \mathcal{Y})^n$

For the actually observed data we denote:

- $i$-th observation: $(\underline{x}^{(i)}, y^{(i)}) \in \mathcal{X} \times \mathcal{Y}$ $i$-th dataset row
- $j$-th feature vector: $\underline{x}_j = (x_j^{(1)}, ..., x_j^{(n)})^T$ $j$-th dataset column
- observed data set: $D := ((\underline{x}^{(1)}, y^{(1)}), ..., (\underline{x}^{(n)}, y^{(n)})) \in \mathbb{D}_n$
- data generating process: $\mathbb{P}_{xy}: \mathcal{X} \times \mathcal{Y} \to [0, 1]$

## Encoding of Categorical Features

Why? Most learning algos can only handle numerical features

Let $x$ be a nominal-scaled feature, i.e. $x \in \{C_1, ..., C_k\}$.

We transform the scalar $x$ into a vector: $o(x) = \begin{pmatrix} \mathbb{1}_{C_1, x}(x) \\ \vdots \\ \mathbb{1}_{C_k}(x) \end{pmatrix} \in \{0,1\}^k$.

Each entry of $o(x)$ is treated as a seperate (binary) feature. $\Rightarrow \forall x, o(x) = (0 ... 1 ... 0)$ Length $k$

**One-hot-encoding**: $\tilde{k} = k$ dummies. So exactly one entry of $o(x)$ is 1 ("hot").

**Dummy-encoding**: $\tilde{k} = k - 1$ dummies. So at most one entry of $o(x)$ is 1.

↳ cuts redundancy of one-hot-encoding.

Necessary if non-singular-matrix is required (e.g. Lin.Reg.)

For an ordinal-scaled feature $x$ we use an encoding that reflects the ... e.g. a sequence of integer values.