



Statistical Learning II

Lecture 2 - Supervised learning

Bruno Loureiro
@ CSD, DI-ENS & CNRS

brloureiro@gmail.com

Recap of Probability

The butter of statistical learning

Random variable

A **random variable** X mathematically formalises the notion of a “**measurement**” or “**random event**”.

Random variable

A **random variable** X mathematically formalises the notion of a “**measurement**” or “**random event**”. It can be:

- **Discrete**: when the possible outcomes are **countable**.

Examples:

- the outcome of tossing a coin $X \in \{\text{head}, \text{tail}\}$
- rolling a dice $X \in \{1, \dots, 6\}$
- The number of people in France $X \in \mathbb{N}$

Random variable

A **random variable** X mathematically formalises the notion of a “**measurement**” or “**random event**”. It can be:

- **Discrete**: when the possible outcomes are **countable**.

Examples:

- the outcome of tossing a coin $X \in \{\text{head}, \text{tail}\}$
- rolling a dice $X \in \{1, \dots, 6\}$
- The number of people in France $X \in \mathbb{N}$

Discrete r.v.s are described by their probability distribution

$$\mathbb{P}(X = k)$$

A positive “function” that sums to one. $\sum_{k \in \text{supp}(X)} \mathbb{P}(X = k) = 1$

Random variable

A **random variable** X mathematically formalises the notion of a “**measurement**” or “**random event**”. It can be:

- **Continuous**: when the possible outcomes are **uncountable**.

Random variable

A **random variable** X mathematically formalises the notion of a “**measurement**” or “**random event**”. It can be:

- **Continuous**: when the possible outcomes are **uncountable**.

Examples:

- The temperature in the room $X \in \mathbb{R}$
- The GDP of France next year $X \in \mathbb{R}$

Random variable

A **random variable** X mathematically formalises the notion of a “**measurement**” or “**random event**”. It can be:

- **Continuous**: when the possible outcomes are **uncountable**.

Examples:

- The temperature in the room $X \in \mathbb{R}$
- The GDP of France next year $X \in \mathbb{R}$

Continuous r.v.s are described by their probability density function (p.d.f.), which integrates to probabilities:

$$\mathbb{P}(X \in [a, b]) = \int_a^b dx \, p_X(x)$$

A “function” that integrates to one: $\int_{\text{supp}(X)} dx \, p_X(x) = 1$

Random variable

A **random variable** X mathematically formalises the notion of a “**measurement**” or “**random event**”. It can be:

- **Continuous**: when the possible outcomes are **uncountable**.

Examples:

- The temperature in the room $X \in \mathbb{R}$
- The GDP of France next year $X \in \mathbb{R}$

Continuous r.v.s are described by their probability density function (p.d.f.), which integrates to probabilities:

$$\mathbb{P}(X \in [a, b]) = \int_a^b dx \, p_X(x)$$

A “function” that integrates to one: $\int_{\text{supp}(X)} dx \, p_X(x) = 1$

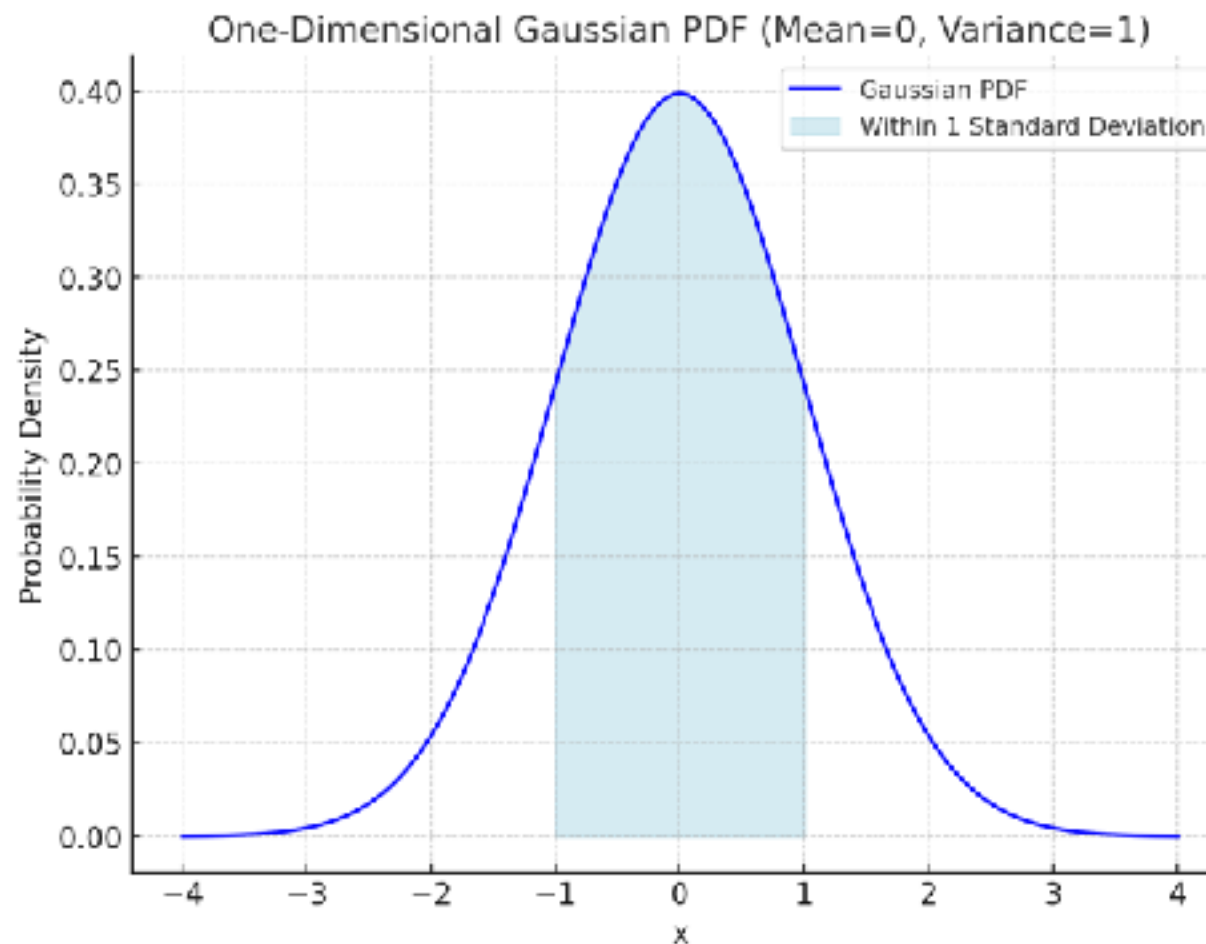


The p.d.f. is NOT a probability. It can be negative.

Normal distribution

A Gaussian r.v. $X \sim \mathcal{N}(\mu, \sigma^2)$ has the following p.d.f.:

$$p_X(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

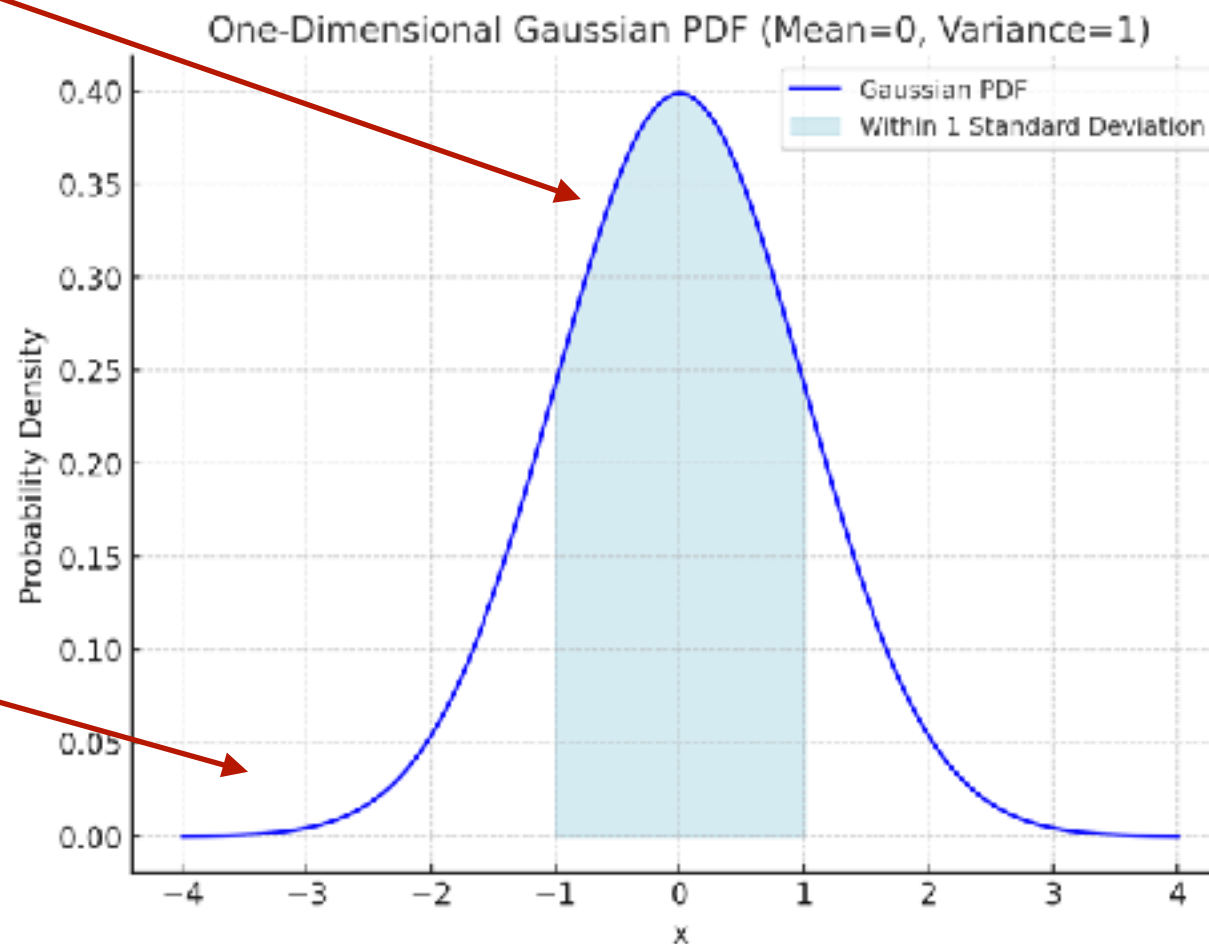


Normal distribution

A Gaussian r.v. $X \sim \mathcal{N}(\mu, \sigma^2)$ has the following p.d.f.:

$$p_X(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

High-probability



Low-probability

Expectation and variance

Let $X \sim p_X$ denote a continuous r.v.

- The **expectation** (or mean) of X is defined as

$$\mathbb{E}[X] = \int dx \, p_X(x) x$$

For example, for $X \sim \mathcal{N}(\mu, \sigma^2)$, we have $\mathbb{E}[X] = \mu$

Expectation and variance

Let $X \sim p_X$ denote a continuous r.v.

- The **expectation** (or mean) of X is defined as

$$\mathbb{E}[X] = \int dx \, p_X(x)x$$

For example, for $X \sim \mathcal{N}(\mu, \sigma^2)$, we have $\mathbb{E}[X] = \mu$

- The **variance** of X is defined as:

$$\text{Var}[X] = \mathbb{E}[(X - \mathbb{E}[X])^2] = \mathbb{E}[X^2] - \mathbb{E}[X]^2$$

For example, for $X \sim \mathcal{N}(\mu, \sigma^2)$, we have $\text{Var}[X] = \sigma^2$

Change of variables

Let $X \sim p_X$ denote a continuous r.v. and $f: \mathbb{R} \rightarrow \mathbb{R}$

Change of variables

Let $X \sim p_X$ denote a continuous r.v. and $f: \mathbb{R} \rightarrow \mathbb{R}$

Then, $Y = f(X)$ is also a random variable, with p.d.f. given by

$$p_Y(y) = \int dx \, p_X(x) \delta(y - f(x))$$

Where $\delta(x)$ is the “Dirac delta function”:

$$\int_{\mathbb{R}} dx \, \delta(x - y) f(x) = f(y)$$

Joint distribution

Two random variables X, Y that concern the same random experiment are characterised by their joint p.d.f.

$$p_{X,Y}(x, y)$$

Joint distribution

Two random variables X, Y that concern the same random experiment are characterised by their joint p.d.f.

$$p_{X,Y}(x, y)$$

The correlation between X, Y is defined by

$$\mathbb{E}[XY] = \int dx \int dy p_{X,Y}(x, y) xy$$

Joint distribution

Two random variables X, Y that concern the same random experiment are characterised by their joint p.d.f.

$$p_{X,Y}(x, y)$$

The correlation between X, Y is defined by

$$\mathbb{E}[XY] = \int dx \int dy p_{X,Y}(x, y) xy$$

We say X, Y are **uncorrelated** if $\mathbb{E}[XY] = \mathbb{E}[X]\mathbb{E}[Y]$

Independence

- Given two r.v.s $X, Y \sim p_{X,Y}$, we define the **marginal distributions**

$$p_X(x) = \int dy \, p_{X,Y}(x, y) \qquad p_Y(y) = \int dx \, p_{X,Y}(x, y)$$

Independence

- Given two r.v.s $X, Y \sim p_{X,Y}$, we define the **marginal distributions**

$$p_X(x) = \int dy \, p_{X,Y}(x, y) \qquad p_Y(y) = \int dx \, p_{X,Y}(x, y)$$

- We say the r.v.s. X, Y are **independent** if

$$p_{X,Y}(x, y) = p_X(x)p_Y(y)$$

Independence

- Given two r.v.s $X, Y \sim p_{X,Y}$, we define the **marginal distributions**

$$p_X(x) = \int dy \, p_{X,Y}(x, y) \qquad p_Y(y) = \int dx \, p_{X,Y}(x, y)$$

- We say the r.v.s. X, Y are **independent** if

$$p_{X,Y}(x, y) = p_X(x)p_Y(y)$$



Note that independence implies uncorrelated, but not the converse!



Exercise: Construct a counter-example

Conditional distribution

- Given two r.v.s $X, Y \sim p_{X,Y}$, we define the conditional p.d.f.

$$p_{X|Y}(x|y) = \frac{p_{X,Y}(x,y)}{p_Y(y)}$$

Conditional distribution

- Given two r.v.s $X, Y \sim p_{X,Y}$, we define the conditional p.d.f.

$$p_{X|Y}(x|y) = \frac{p_{X,Y}(x,y)}{p_Y(y)}$$

Note that $X, Y \sim p_{X,Y}$ are independent if and only if:

$$p_{X|Y}(x|y) = p_X(x)$$

Conditional distribution

- Given two r.v.s $X, Y \sim p_{X,Y}$, we define the conditional p.d.f.

$$p_{X|Y}(x|y) = \frac{p_{X,Y}(x,y)}{p_Y(y)}$$

Note that $X, Y \sim p_{X,Y}$ are independent if and only if:

$$p_{X|Y}(x|y) = p_X(x)$$

Theorem (Bayes theorem)

$$p_{X|Y}(x|y) = \frac{p_{Y|X}(y|x)p_X(x)}{p_Y(y)}$$

Law of large numbers

Let $X_1, \dots, X_n \sim p_X$ denote i.i.d. r.v.s. with mean $\mathbb{E}[X_i] = \mu$

Define the sample mean (note this is itself a r.v.)

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$$

Law of large numbers

Let $X_1, \dots, X_n \sim p_X$ denote i.i.d. r.v.s. with mean $\mathbb{E}[X_i] = \mu$

Define the sample mean (note this is itself a r.v.)

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$$

Theorem (Weak LLN)

$$\bar{X}_n \xrightarrow{P} \mu \quad \text{as} \quad n \rightarrow \infty$$

$$\lim_{n \rightarrow \infty} \mathbb{P}(|\bar{X}_n - \mu| < \epsilon) = 1$$



Be aware there are many variations of the LLN.

Central limit theorem

Let $X_1, \dots, X_n \sim p_X$ denote i.i.d. r.v.s. with mean $\mathbb{E}[X_i] = \mu$ and variance $\text{Var}(X_i) = \sigma^2 < \infty$

Again, consider the sample mean

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$$

Central limit theorem

Let $X_1, \dots, X_n \sim p_X$ denote i.i.d. r.v.s. with mean $\mathbb{E}[X_i] = \mu$ and variance $\text{Var}(X_i) = \sigma^2 < \infty$

Again, consider the sample mean

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$$

Theorem (Lindeberg CLT)

$$\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{d} \mathcal{N}(0, \sigma^2)$$

$$\lim_{n \rightarrow \infty} \mathbb{P}(\sqrt{n}(\bar{X}_n - \mu) \leq z) = \mathbb{P}(Z \leq z/\sigma) \quad Z \sim \mathcal{N}(0, 1)$$



Be aware there are many variations of the CLT.

Supervised learning

Supervised Learning



Not grumpy



Grumpy

Supervised Learning



Not grumpy



Grumpy



???

Supervised Learning

Inputs / covariates $x \in \mathcal{X}$



x_1



x_2



x_3



y_1

Not grumpy



y_2

Grumpy



y_3

???

Outputs / Labels / response $y \in \mathcal{Y}$

Supervised Learning

Inputs / covariates $x \in \mathcal{X}$



x_1



x_2



x_3



y_1

Not grumpy

y_2

Grumpy

y_3

???

Outputs / Labels / response $y \in \mathcal{Y}$

Supervised Learning

Inputs / covariates $x \in \mathcal{X}$



x_1



x_2



x_3

$f: \mathcal{X} \rightarrow \mathcal{Y}$



y_1

Not grumpy

y_2

Grumpy

y_3

???

Outputs / Labels / response $y \in \mathcal{Y}$

Supervised Learning

Let $\mathcal{D} = \{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y} : i = 1, \dots, n\}$ denote the training data.

Supervised Learning

Let $\mathcal{D} = \{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y} : i = 1, \dots, n\}$ denote the training data.



- The input space \mathcal{X} is often assumed to be a vector space $\mathcal{X} \subset \mathbb{R}^d$. But keep in mind in real life it can be any data structure (e.g. a pandas.DataFrame)

Supervised Learning

Let $\mathcal{D} = \{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y} : i = 1, \dots, n\}$ denote the **training data**.



- The **input space** \mathcal{X} is often assumed to be a **vector space** $\mathcal{X} \subset \mathbb{R}^d$. But keep in mind in real life it can be any **data structure** (e.g. a pandas.DataFrame)
- The **output space** \mathcal{Y} is often assumed to be a subset $\mathcal{Y} \subset \mathbb{R}$.

Supervised Learning

Let $\mathcal{D} = \{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y} : i = 1, \dots, n\}$ denote the **training data**.



- The **input space** \mathcal{X} is often assumed to be a **vector space** $\mathcal{X} \subset \mathbb{R}^d$. But keep in mind in real life it can be any **data structure** (e.g. a pandas.DataFrame)
- The **output space** \mathcal{Y} is often assumed to be a subset $\mathcal{Y} \subset \mathbb{R}$.
- In particular, if $|\mathcal{Y}| = k$ is a **discrete** set, we say we have a **classification problem**.

Supervised Learning

Let $\mathcal{D} = \{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y} : i = 1, \dots, n\}$ denote the **training data**.



- The **input space** \mathcal{X} is often assumed to be a **vector space** $\mathcal{X} \subset \mathbb{R}^d$. But keep in mind in real life it can be any **data structure** (e.g. a pandas.DataFrame)
- The **output space** \mathcal{Y} is often assumed to be a subset $\mathcal{Y} \subset \mathbb{R}$.
- In particular, if $|\mathcal{Y}| = k$ is a **discrete** set, we say we have a **classification problem**.
- If \mathcal{Y} is a **continuous** set, we say we have a **regression problem**.

Supervised Learning

Let $\mathcal{D} = \{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y} : i = 1, \dots, n\}$ denote the **training data**.



- The **input space** \mathcal{X} is often assumed to be a **vector space** $\mathcal{X} \subset \mathbb{R}^d$. But keep in mind in real life it can be any **data structure** (e.g. a pandas.DataFrame)
- The **output space** \mathcal{Y} is often assumed to be a subset $\mathcal{Y} \subset \mathbb{R}$.
- In particular, if $|\mathcal{Y}| = k$ is a **discrete** set, we say we have a **classification problem**.
- If \mathcal{Y} is a **continuous** set, we say we have a **regression problem**



It is very common to consider a **one-hot encoding** $\mathcal{Y} = \{0, 1\}^k$ in classification.

Supervised Learning

Let $\mathcal{D} = \{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y} : i = 1, \dots, n\}$ denote the training data.

Examples of classification:

- Grumpy vs. Non-grumpy cats

$\mathcal{X} = \{\text{photos of cats}\}, \mathcal{Y} = \{\text{grumpy, not grumpy}\}$

- E-mail spam detection

$\mathcal{X} = \{\text{your inbox}\}, \mathcal{Y} = \{\text{spam, not spam}\}$

- Medical diagnosis

$\mathcal{X} = \{\text{medical data}\}, \mathcal{Y} = \{\text{diseases}\}$

- Sentiment analysis

$\mathcal{X} = \{\text{text}\}, \mathcal{Y} = \{\text{positive, negative, neutral}\}$

Supervised Learning

Let $\mathcal{D} = \{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y} : i = 1, \dots, n\}$ denote the training data.

Examples of regression:

- Temperature prediction

$$\mathcal{X} = \mathbb{R}^3, \mathcal{Y} = \mathbb{R}$$

- Stock price prediction

$$\mathcal{X} = \{\text{list of stocks}\}, \mathcal{Y} = \mathbb{R}_+$$

- Life expectancy

$$\mathcal{X} = \{\text{medical data}\}, \mathcal{Y} = \mathbb{R}_+$$

- Any price, cost, income, etc. prediction.

Supervised Learning

Let $\mathcal{D} = \{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y} : i = 1, \dots, n\}$ denote the training data.

Supervised Learning

Let $\mathcal{D} = \{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y} : i = 1, \dots, n\}$ denote the training data.

In supervised learning, our goal is to use the data to learn a function that correctly assigns the labels to the responses.

$$f : \mathcal{X} \rightarrow \mathcal{Y}$$

Supervised Learning

Let $\mathcal{D} = \{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y} : i = 1, \dots, n\}$ denote the **training data**.

In supervised learning, our goal is to use the data to **learn** a function that correctly **assigns the labels** to the **responses**.

$$f: \mathcal{X} \rightarrow \mathcal{Y}$$



For **classification**, it is common to define instead:

$$f: \mathcal{X} \rightarrow [0,1]^{|\mathcal{Y}|}$$

Where $f(x)$ is a vector of **class probabilities**. In this case, final prediction is given by:

$$\hat{y} = \operatorname{argmax}_{k \in |\mathcal{Y}|} f(x)$$

Supervised Learning

Let $\mathcal{D} = \{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y} : i = 1, \dots, n\}$ denote the training data.

In supervised learning, our goal is to use the data to learn a function that correctly assigns the labels to the responses.

$$f : \mathcal{X} \rightarrow \mathcal{Y}$$

Two key words: correctly and learn.

Loss function

Let $\mathcal{D} = \{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y} : i = 1, \dots, n\}$ denote the **training data**.

In supervised learning, our goal is to use the data to **learn** a function that correctly **assigns the labels** to the **responses**.

$$f : \mathcal{X} \rightarrow \mathcal{Y}$$

Two key words: **correctly** and **learn**. To quantify the first, it is common to introduce a **loss function**:

$$\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$$

Loss function

Let $\mathcal{D} = \{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y} : i = 1, \dots, n\}$ denote the **training data**.

In supervised learning, our goal is to use the data to **learn** a function that correctly **assigns the labels** to the **responses**.

$$f : \mathcal{X} \rightarrow \mathcal{Y}$$

Two key words: **correctly** and **learn**. To quantify the first, it is common to introduce a **loss function**:

$$\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$$

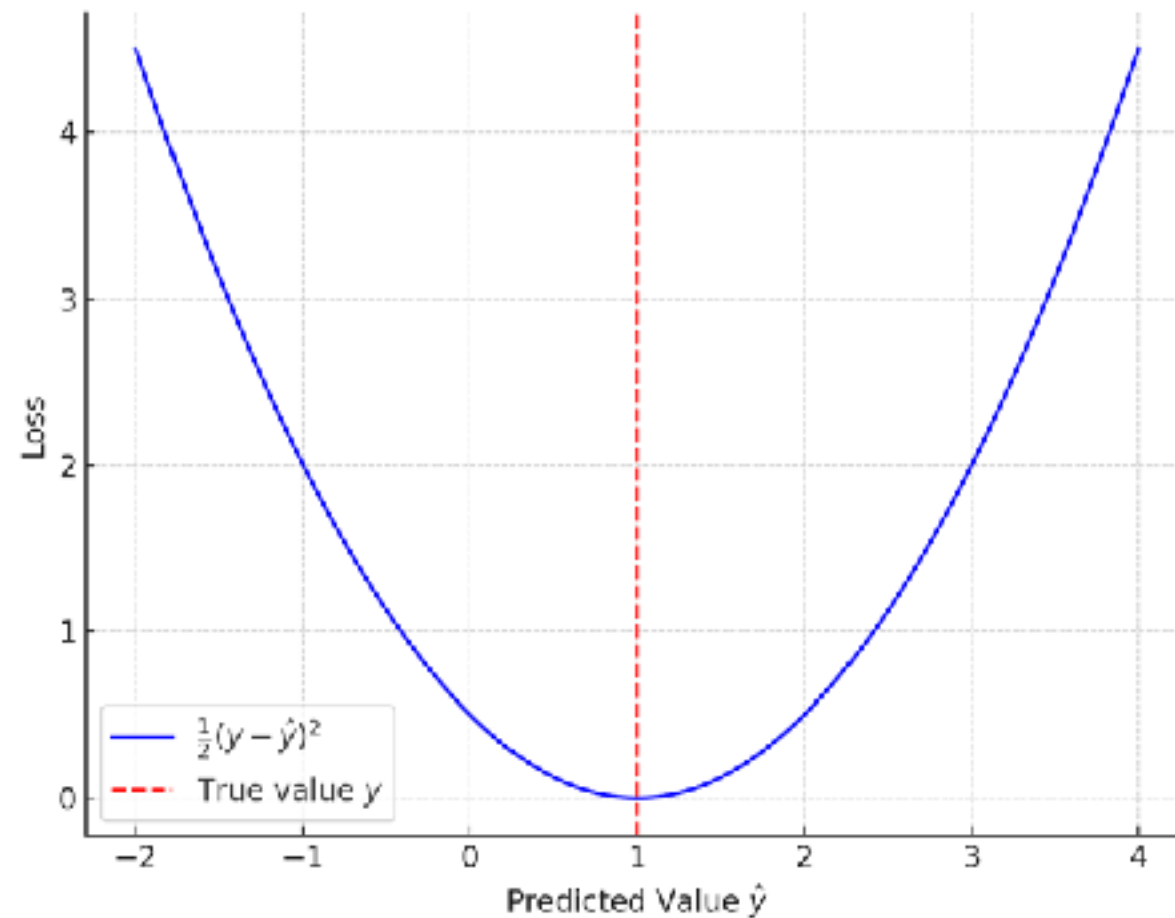


For **classification** this will also depend on the **encoding**.

Regression losses

Examples in regression:

- Square loss: $\ell(y, z) = \frac{1}{2}(y - z)^2$



Regression losses

Examples in regression:

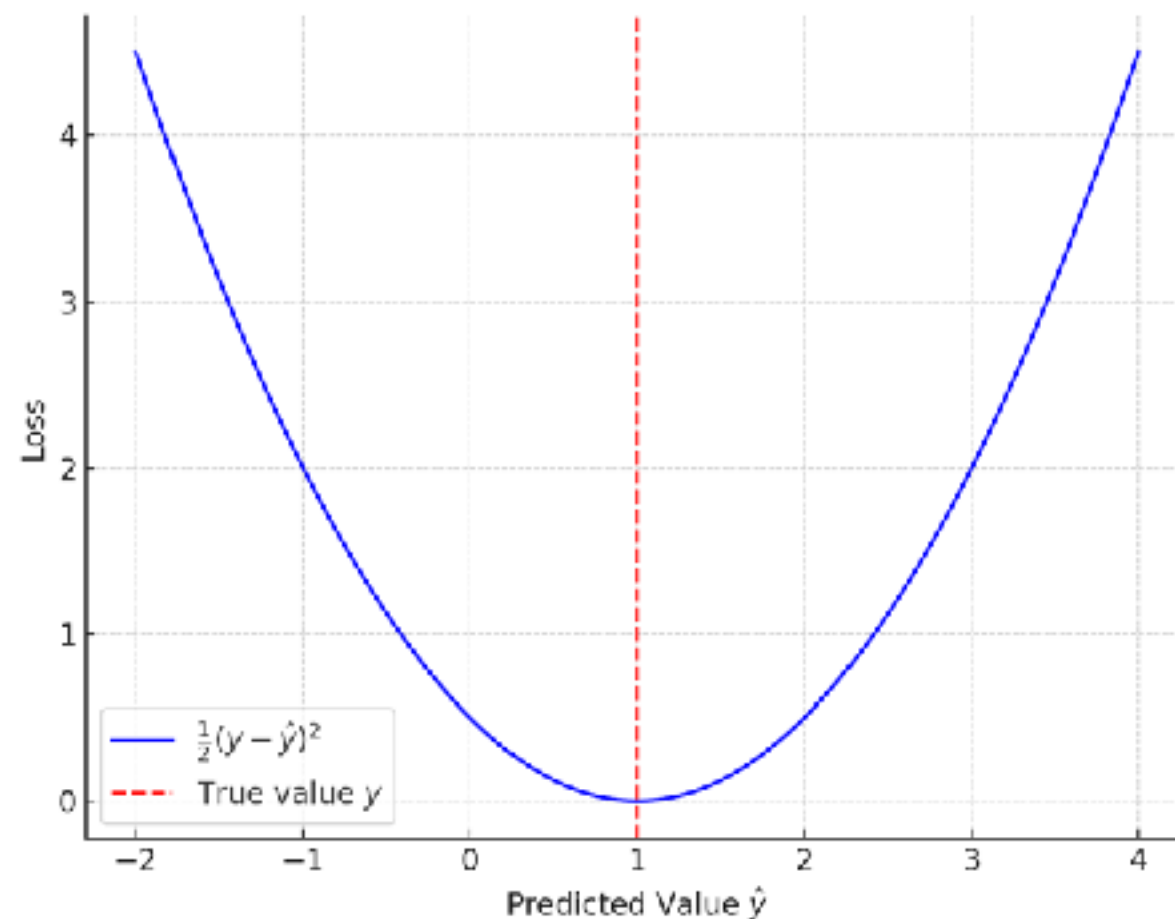
- Square loss: $\ell(y, z) = \frac{1}{2}(y - z)^2$



The square loss is sensitive to outliers



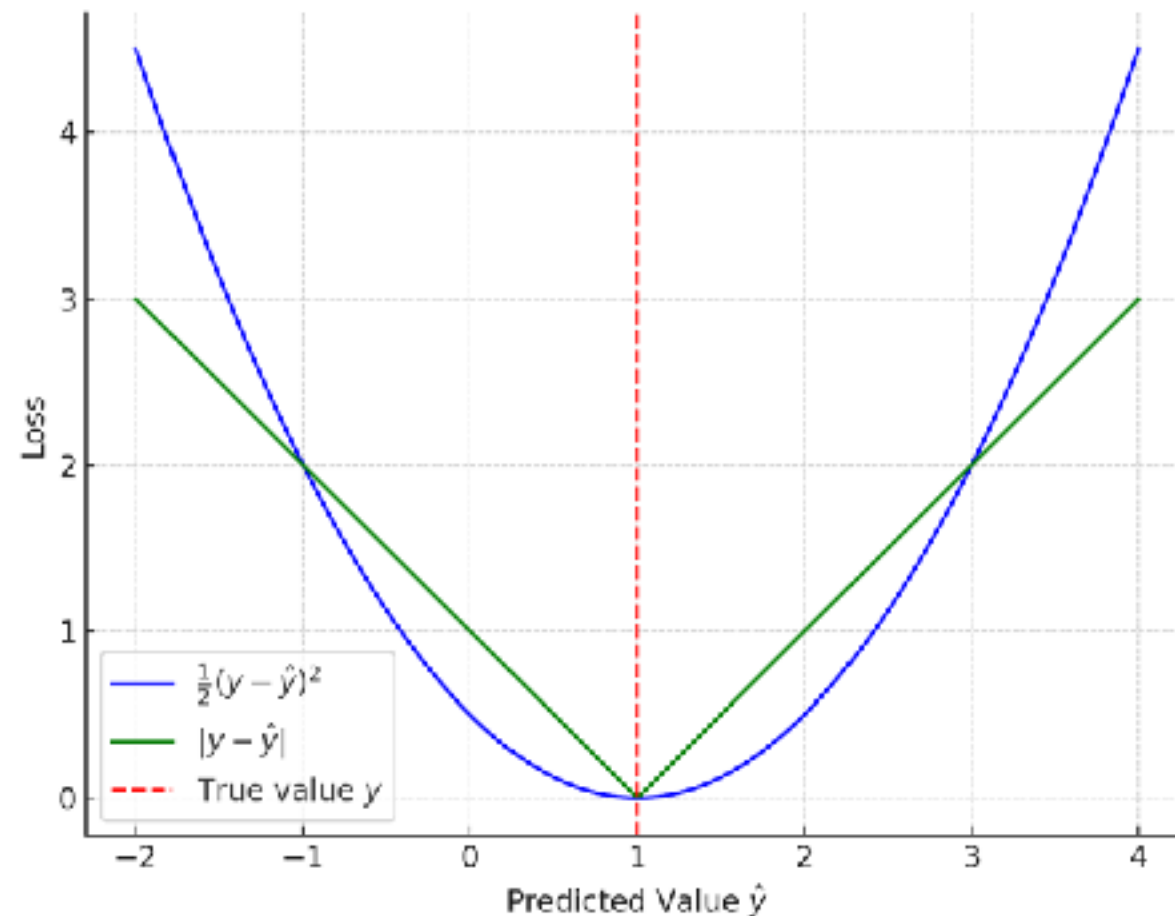
Exercise:
show this.



Regression losses

Examples in regression:

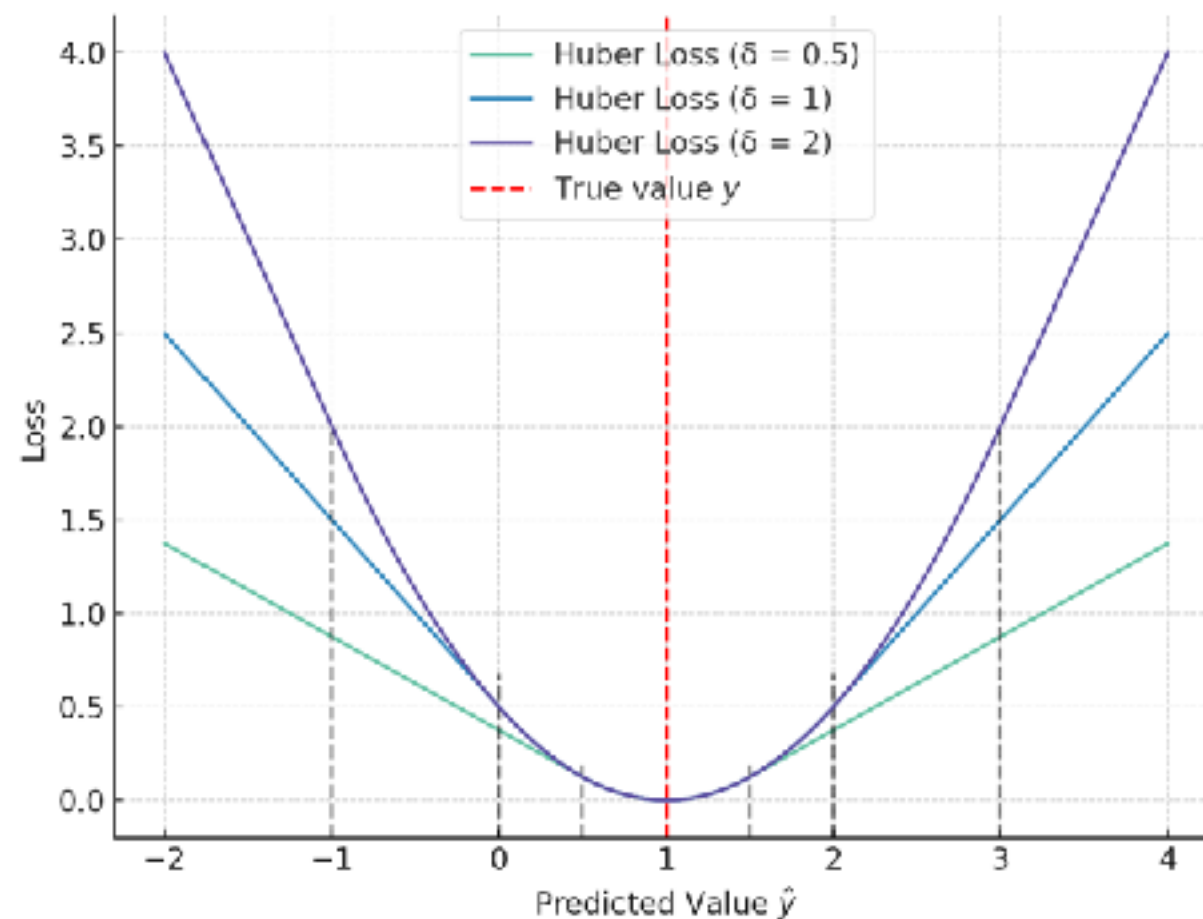
- Square loss: $\ell(y, z) = \frac{1}{2}(y - z)^2$
- Absolute loss: $\ell(y, z) = |y - z|$



Regression losses

Examples in regression:

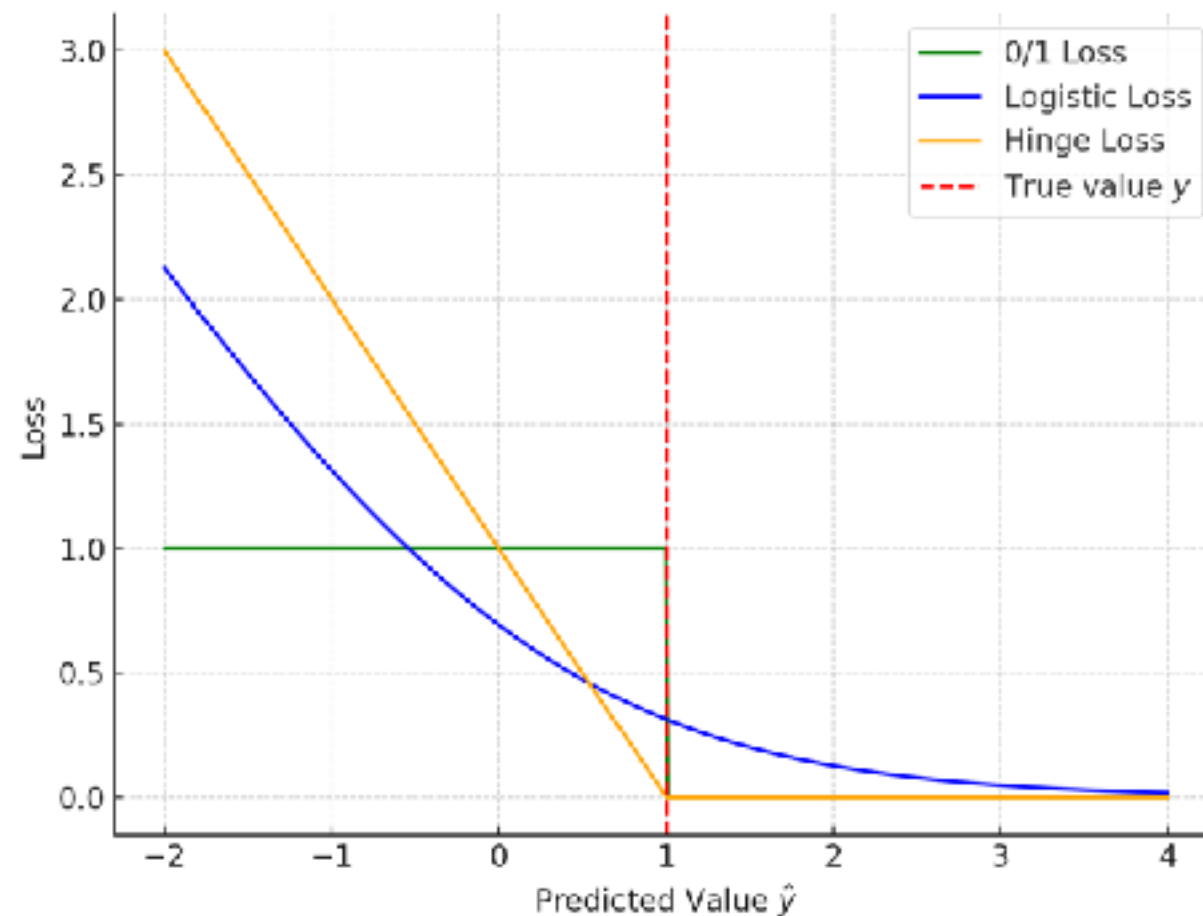
- Huber loss: $\ell_{\delta}(y, z) = \begin{cases} \frac{1}{2}(y - z)^2 & \text{if } |y - z| \leq \delta \\ \delta(|y - z| - \frac{1}{2}\delta) & \text{if } |y - z| > \delta \end{cases}$



Classification losses

Examples in binary classification $\mathcal{Y} = \{-1, +1\}$:

- 0/1 loss: $\ell(y, z) = \delta_{yz}$ (or $\ell(y, z) = \theta(y - z) = \begin{cases} 1 & \text{if } y - z \leq 0 \\ 0 & \text{otherwise} \end{cases}$)
- Logistic loss: $\ell(y, z) = \log(1 + e^{-yz})$
- Hinge loss: $\ell(y, z) = \max(0, 1 - yz)$



Empirical risk

Let $\mathcal{D} = \{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y} : i = 1, \dots, n\}$ denote the **training data**.

Given a loss function $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$, and a predictor $f : \mathcal{X} \rightarrow \mathcal{Y}$ define the **empirical risk**:

$$\hat{\mathcal{R}}_n(f) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i))$$

Also known as the **training loss**.

Empirical risk

Let $\mathcal{D} = \{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y} : i = 1, \dots, n\}$ denote the **training data**.

Given a loss function $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_+$, and a predictor $f : \mathcal{X} \rightarrow \mathcal{Y}$ define the **empirical risk**:

$$\hat{\mathcal{R}}_n(f) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, f(x_i))$$

Also known as the **training loss**. This quantifies how well we fit the data. But is this a good notion of learning?

$$f(x) = \begin{cases} y_i & \text{if } x \in \mathcal{D} \\ 0 & \text{otherwise} \end{cases} \quad \Rightarrow \quad \hat{\mathcal{R}}_n = 0$$

Probabilistic framework

Instead, we would like predictors that do well on **unseen data**.

Probabilistic framework

Instead, we would like predictors that do well on **unseen data**.

Assume there is an underlying data distribution p over $\mathcal{X} \times \mathcal{Y}$:

$$(x_i, y_i) \sim p \quad \text{i.i.d.}$$

Probabilistic framework

Instead, we would like predictors that do well on **unseen data**.

Assume there is an underlying data distribution p over $\mathcal{X} \times \mathcal{Y}$:

$$(x_i, y_i) \sim p \quad \text{i.i.d.}$$



- The “i.i.d.” assumption might not always hold. (Sampling bias, distribution shift, etc.)
- Under this assumption, $\hat{\mathcal{R}}_n$ is a random function.

Population risk

Instead, we would like predictors that do well on **unseen data**.

Assume there is an underlying data distribution p over $\mathcal{X} \times \mathcal{Y}$:

$$(x_i, y_i) \sim p \quad \text{i.i.d.}$$

Define the notion of population risk of a predictor $f: \mathcal{X} \rightarrow \mathcal{Y}$:

$$\mathcal{R}(f) = \mathbb{E} [\ell(y, f(x))]$$

Also known as the **generalisation** or **test error**.

Population risk

Instead, we would like predictors that do well on **unseen data**.

Assume there is an underlying data distribution p over $\mathcal{X} \times \mathcal{Y}$:

$$(x_i, y_i) \sim p \quad \text{i.i.d.}$$

Define the notion of population risk of a predictor $f: \mathcal{X} \rightarrow \mathcal{Y}$:

$$\mathcal{R}(f) = \mathbb{E} [\ell(y, f(x))]$$

Also known as the **generalisation** or **test error**.



\mathcal{R} is a deterministic function of the predictor f