

IRIS DATA SET

Introduced by the statistician Ronald Fisher and one of the most frequently used toy examples.

- Classify iris subspecies based on flower measurements.
- 150 iris flowers: 50 versicolor, 50 virginica, 50 setosa.
- Sepal length / width and petal length / width in [cm].



Source: <https://rpubs.com/vidhividhi/irisdataeda>

Word of warning: "iris" is a small, clean, low-dimensional data set, which is very easy to classify; this is not necessarily true in the wild.

DATA-GENERATING PROCESS

- We assume the observed data \mathcal{D} to be generated by a process that can be characterized by some probability distribution

$$\mathbb{P}_{xy},$$

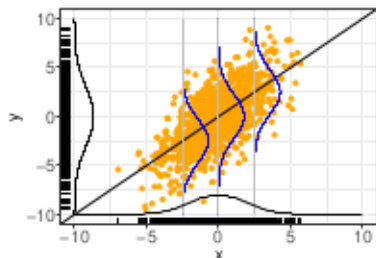
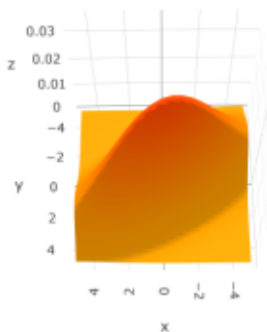
defined on $\mathcal{X} \times \mathcal{Y}$.

- We denote the random variables following this distribution by lowercase \mathbf{x} and y .
- It is important to understand that the true distribution is essentially **unknown** to us. In a certain sense, learning (part of) its structure is what ML is all about.



DATA-GENERATING PROCESS / 2

- We assume data to be drawn *i.i.d.* from the joint probability density function (pdf) / probability mass function (pmf) $p(\mathbf{x}, y)$.
 - *i.i.d.* stands for independent and identically distributed.
 - This means: We assume that all samples are drawn from the same distribution and are mutually independent – the i -th realization does not depend on the other $n - 1$ ones.
 - This is a strong yet crucial assumption that is precondition to most theory in (basic) ML.



Remarks:

- With a slight abuse of notation we write random variables, e.g., \mathbf{x} and y , in lowercase, as normal variables or function arguments. The context will make clear what is meant.
- Often, distributions are characterized by a parameter vector $\theta \in \Theta$. We then write $p(\mathbf{x}, y | \theta)$.
- This lecture mostly takes a frequentist perspective. Distribution parameters θ appear behind the $|$ for improved legibility, not to imply that we condition on them in a probabilistic Bayesian sense. So, strictly speaking, $p(\mathbf{x}|\theta)$ should usually be understood to mean $p_{\theta}(\mathbf{x})$ or $p(\mathbf{x}, \theta)$ or $p(\mathbf{x}; \theta)$. On the other hand, this notation makes it very easy to switch to a Bayesian view.

