

# SUMMARY ON STATISTICS

## 1. Useful probability distributions

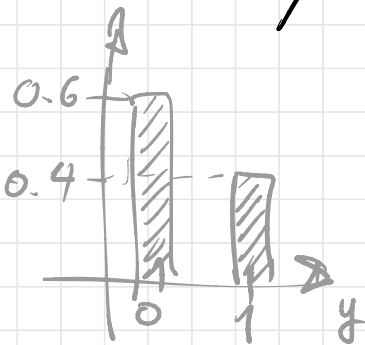
### 1.1) Binary variables

$$y = \{0, 1\}, \{-1, 1\}, \{A, B\}, \dots$$

(The way to treat a binary classification problem)

$$y \sim \text{Bern}(y|\mu) \quad \text{"pipe" notation} \Rightarrow \text{conditioned to ...}$$

$$= \mu^y \cdot (1-\mu)^{1-y} \quad \left. \vphantom{= \mu^y \cdot (1-\mu)^{1-y}} \right\} \text{the likelihood of depending on parameter } \mu$$



$$\begin{aligned} \mu &= P(y=1) \\ P(y=0) &= 1-\mu \end{aligned}$$

Imagine we have a dataset  $D = \{ \vec{x}_i, y_i \}_{i=1}^N$

e.g.  $y_i$ : type of galaxy

(spiral or elliptical)

$\vec{x}_i$ : position in the sky [[given]]

for i.i.d. data, the total likelihood is:

$$P(\vec{y} | X, \vec{\theta}) = \prod_{i=1}^N \text{Bern}(y_i | \mu(\vec{x}_i; \vec{\theta}))$$

$$= \prod_i \mu(\vec{x}_i; \vec{\theta})^{y_i} (1 - \mu(\vec{x}_i; \vec{\theta}))^{1-y_i}$$

$\mu(\vec{x}; \vec{\theta}) \Rightarrow$  This is given by our model  
(ML model, physics model, ...)

Given a model for  $\mu$ , we need to fit it to our data, i.e. to find the optimum parameters,  $\vec{\theta}_{\text{opt}}$ .

e.g.  $\vec{\theta}_{\text{opt}} = \underset{\vec{\theta}}{\text{argmax}} P(\vec{y} | X, \vec{\theta})$

$\hookrightarrow$  Maximum Likelihood Estimator

## 1.2) Categorical variables

e.g. types of jets at LHC, types of events at particle detectors

$$y = \{1, 2, \dots, k\}, \{A, B, C, \dots\}$$

parametrise it in a computer-convenient way: "1-to-k scheme" or "one-hot-encoding".

Suppose  $k=3$  classes

$$y=1 \rightarrow \vec{y} = \{1, 0, 0\}$$

$$y=2 \rightarrow \vec{y} = \{0, 1, 0\}$$

$$y=3 \rightarrow \vec{y} = \{0, 0, 1\}$$

$$\text{i.e. } y=k \rightarrow \vec{y} = \{0, 0, \dots, 1, 0, 0, \dots\}$$

$\hookrightarrow$   $k^{\text{th}}$  position

These variables are distributed according to: the categorical distribution (or generalized Bernoulli):

$$P(\vec{y} | \vec{\mu}) = \prod_{k=1}^k \mu_k^{y_k}$$

$$\mu_k = P(y_k=1)$$

$$\vec{\mu} = \{\mu_k\}$$

indeed for  $k=2$  (binary case)

$$p(\vec{y} | \vec{\mu}) = \mu_1^{y_1} \mu_2^{y_2} = \mu_1^{y_1} (1 - \mu_1)^{1 - y_1}$$

$$\left( \begin{aligned} \mu_1 &= P(y_1=1) = P(y_2=0) \\ &= \text{Bern}(y_1 | \mu_1) \end{aligned} \right)$$

So what is the likelihood of a dataset of categorical variables?

$$D = \{ \vec{x}_i, \vec{y}_i \}_{i=1}^N$$

one-hot-encoded

$$X = \{ \vec{x}_i \}, Y = \{ \vec{y}_i \}$$

$$P(Y|X, \vec{\theta}) = \prod_{i=1}^N \prod_{k=1}^K \mu(\vec{x}_i, \vec{\theta})^{y_{ik}}$$

$$-\ln p(Y|X, \vec{\theta}) = - \sum_i \sum_k y_{ik} \ln \mu(\vec{x}_i, \vec{\theta})$$

↳ "Cross entropy" cost function

↳ to minimise to obtain  $\vec{\theta}_{\text{opt}}$

### 1.3) Integer-valued variables

$$y = \{0, 1, 2, \dots\} \quad y \in \mathbb{Z}$$

any counting experiment,  
photons from the sky, etc

(not really a categorical variable)

- Under certain assumptions

(independence, non-overlapping, small rate)

$$y \sim \text{Pois}(y | \lambda)$$

Poisson distribution with  
mean  $\lambda$

$$= \frac{\lambda^y e^{-\lambda}}{y!}$$

So assume having data:

$$D = \{ \vec{x}_i, y_i \}_{i=1}^N$$

$$P(\vec{y} | X, \vec{\theta}) = \prod_{i=1}^N \frac{1}{y_i!} \lambda(\vec{x}_i; \vec{\theta})^{y_i} \cdot \exp\{-\lambda(\vec{x}_i; \vec{\theta})\}$$

## 1.4) Continuous variables

$$y \in \mathbb{R}$$

Flux of particles, energy deposition, time of flight, ...

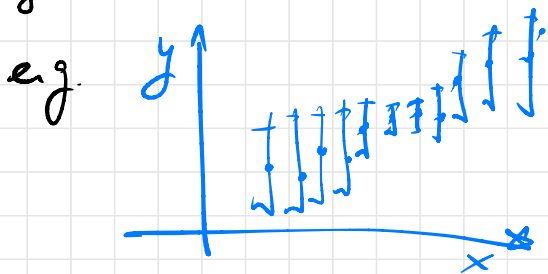
The simplest, most common, is to assume a Gaussian distribution

$$y \sim \mathcal{N}(y | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-\mu)^2}{2\sigma^2}}$$

or its multivariate version

$$\vec{y} \sim \mathcal{N}(\vec{y} | \vec{\mu}, \Sigma) = \frac{1}{(2\pi)^{d/2} (\det \Sigma)^{1/2}} \cdot \exp\left\{-\frac{1}{2} (\vec{y} - \vec{\mu})^T \Sigma^{-1} (\vec{y} - \vec{\mu})\right\}$$

Again the aim is to model  $\mu$ .



⇒ Fit a straight line to these data

$$D = \{x_i, y_i, \sigma_i\} \quad y_i \sim \mathcal{N}(y_i \mid \mu(x_i; \vec{\theta}), \sigma_i^2)$$

$$\mu(x, \vec{\theta}) = \theta_1 x + \theta_2$$

The data likelihood:

$$p(\vec{y} \mid \vec{x}, \vec{\theta}) = \prod_{i=1}^N \mathcal{N}(y_i \mid \mu(x_i; \vec{\theta}), \sigma_i^2)$$

$$\log p(\vec{y} \mid \vec{x}, \vec{\theta}) = \sum_i -\frac{(y_i - \theta_1 x_i - \theta_2)^2}{2\sigma_i^2} + \text{const}$$

$$\vec{\theta}_{\text{opt}} = \underset{\vec{\theta}}{\text{argmin}} \sum_i \frac{(y_i - \theta_1 x_i - \theta_2)^2}{\sigma_i^2} \quad \left. \vphantom{\sum_i} \right\} \begin{array}{l} \text{Least-squares} \\ \text{procedure!} \end{array}$$

The results for  $\theta_1$  and  $\theta_2$  coming from this minimisation coincide with the one in textbooks

All this belongs to the "frequentist" approach, outcome is a priori  $\vec{\theta}_{\text{opt}}$ , i.e. a point-like prediction without uncertainty estimation.

There are several methods to estimate the uncertainty on  $\vec{\theta}_{\text{opt}}$  in the frequentist framework

{ Bootstrap, Analytical method, Fisher information...  
or just Neyman-Pearson intervals

On the other hand, uncertainties can be estimated "from first principles" using the Bayesian approach:

$$\mathcal{D} = \{ \vec{x}_i, y_i \}_{i=1}^N$$

$$p(\vec{\theta} | \mathcal{D}) = \frac{p(\mathcal{D} | \vec{\theta}) p(\vec{\theta})}{p(\mathcal{D})}$$

$$p(\vec{y} | \mathcal{X}, \vec{\theta})$$

is the same likelihood as before

"Evidence" or "marginal likelihood"

prior distribution of the parameters of the model "m"

$$p(\mathcal{D}) = p(\mathcal{D} | \mathcal{M}) = \int d\vec{\theta} p(\mathcal{D} | \vec{\theta}) p(\vec{\theta})$$



Since  $p(D)$  doesn't depend on  $\vec{\theta}$ , it is not needed for computing  $\vec{\theta}_{opt}$  given a dataset (however it is needed for "model comparison", as we will see later)

- So  $\vec{\theta}_{opt}$  is obtained in this framework by maximising the (un-normalised) posterior

$$\vec{\theta}_{opt} = \vec{\theta}_{MAP} = \underset{\vec{\theta}}{\operatorname{argmax}} p(D|\vec{\theta})p(\vec{\theta})$$

So if we have a way to estimate  $p(\vec{\theta}|D)$ , then it naturally gives a way to estimate the uncertainties

$$1 - \beta = \int_{\vec{\theta}_{min}}^{\vec{\theta}_{max}} d\vec{\theta} p(\vec{\theta}|D) = 0.95 \quad \text{e.g.}$$

$$\vec{\theta} = [\vec{\theta}_{min}, \vec{\theta}_{max}]$$

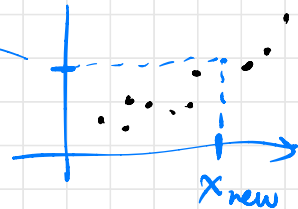
e.g. 95%  
credible interval  
for  $\beta = 0.05$

- In many situations we need to go a step beyond  $\vec{\theta}_{\text{opt}}$  and actually make a prediction for a new input  $\vec{x}_{\text{new}}$

- Frequentist approach:

Simply substitute  $\vec{\theta}_{\text{opt}}$  inside the expected value of " $y$ ":

$$\mu(\vec{x}_{\text{new}}; \vec{\theta}_{\text{opt}})$$



- Bayesian approach:

Distribution of  $\vec{\theta}$   $p(\vec{\theta} | D)$  translates into a distribution of the prediction, a.k.a. the "predictive distribution"

$$P(y_{\text{new}} | \vec{x}_{\text{new}}, D) = \int d\vec{\theta} p(\vec{\theta} | D) p(y_{\text{new}} | \vec{x}_{\text{new}}, \vec{\theta})$$

likelihood of  $y_{\text{new}}$

In very few situations  $p(\vec{\theta} | D)$  is analytical (or even tractable!).

The prototypical example is in the following situation:

- Prior:  $p(\vec{\theta}) = \mathcal{N}(\vec{\theta} | \vec{\mu}_0, \Sigma_0)$

- Likelihood:  $p(\vec{y} | \vec{\theta}) = \mathcal{N}(\vec{y} | \Phi(x) \cdot \vec{\theta} + \vec{b}, C)$

i.e. Both prior & likelihood are Gaussian, but the mean of the likelihood has to be linear in the parameters!

(e.g. a polynomial:

$$f(x; \vec{\theta}) = \theta_0 + \theta_1 x + \theta_2 x^2 + \dots,$$

but also a linear combination of non-linear functions of  $x$ :

$$f(x; \vec{\theta}) = \sum_{k=0}^M \theta_k \cdot \phi_k(x) \equiv \vec{\theta}^T \cdot \vec{\phi}(x)$$

$$\vec{\phi}(x) = \{1, \phi_1(x), \phi_2(x), \dots, \phi_M\}$$

$\{\phi_k(x)\}_{k=1}^M$  is a given catalog of non-linear functions (whose params. are fixed)

So if we have  $N$  input values  $\{x_i\}$  we can build a matrix

$$\Phi(X) = \begin{pmatrix} 1 & \phi_1(x_1) & \phi_2(x_1) & \dots & \phi_M(x_1) \\ 1 & \phi_1(x_2) & \phi_2(x_2) & \dots & \phi_M(x_2) \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & \phi_1(x_N) & \phi_2(x_N) & \dots & \phi_M(x_N) \end{pmatrix}$$

Note that in spite of having a simple structure, this model is actually a "universal approximator", the same way a neural network is! with the advantage that the solution is analytical!