

Probabilistic Inference



Bernoulli distribution → La distribución de Bernoulli es una distribución aplicada a una variable aleatoria discreta, la cual solo puede resultar en 2 sucesos posibles: "éxito" y "no éxito" ($\text{(\text{T}, \text{H})}$)

$$p_i(F_i = \text{(\text{T})}) = \theta_i$$

$$p_i(F_i = \text{(\text{H})} | \theta_i) = \text{Ber}(F_i = \text{(\text{T})} | \theta_i) = \theta_i \quad F_i \sim \text{Ber}(\theta_i)$$

↑
probabilidad

Maximum Likelihood

1. The flips are qualitatively the same - identical distribution

$$p_i(F_i = f_i | \theta_i) = p_i(F_i = f_i | \theta)$$

2. The coins flips do not affect each other - independence

$$\begin{aligned} p(\text{(\text{H}, \text{T}, \text{H}, \text{H})} | \theta) &= p(F_1 = \text{(\text{H})} | \theta) \cdot p(F_2 = \text{(\text{T})} | \theta) \cdot p(F_3 = \text{(\text{H})} | \theta) \cdot p(F_4 = \text{(\text{H})} | \theta) \\ &= \prod_{i=1}^4 p(F_i = f_i | \theta) \end{aligned}$$

$$f(\theta) := p(D | \theta)$$

$$\theta_{\text{MLE}} = \underset{\theta \in [0, 1]}{\operatorname{argmax}} f(\theta)$$

↑
is not a probability distribution

Monotonic functions preserve critical points

$$\underset{\theta \in [0, 1]}{\operatorname{argmax}} f(\theta) = \underset{\theta \in [0, 1]}{\operatorname{argmax}} \log f(\theta)$$

$$f(x) = \log(x) \quad \text{ex.: } f(x) = 10 \log(4x)$$

$$f'(x) = \frac{x'}{x} = \frac{1}{x}$$

$$f'(x) = \frac{10}{x}$$

$$f(x) = \theta^t (1-\theta)^{1-t} = \log \theta^t (1-\theta)^{1-t} = \log \theta^t + \log (1-\theta)^{1-t} = t \log \theta + (1-t) \log(1-\theta)$$

$$f'(x) = \frac{t}{\theta} + \frac{1-t}{1-\theta} = \frac{t}{\theta} - \frac{1-t}{1-\theta}$$

$$\text{to find the max } \rightarrow \frac{t}{\theta} - \frac{1-t}{1-\theta} = 0 ; \frac{t}{\theta} = \frac{1-t}{1-\theta} ; (1-\theta)t = \theta(1-t)$$

$$(1-\theta)t = \theta h ; h = \frac{(1-\theta)t}{\theta} ; h = \frac{t - \theta t}{\theta} ; t = \theta h + \theta t ; t = \theta(h+t) ;$$

$$\theta = \frac{t}{h+t}$$

MLE is based on the results of the experiment, but if we have just a few instances, it may seem counter-intuitive.

We have **prior beliefs** that also counts

Bayesian inference

A prior distribution $p(\theta)$ represents our beliefs before we observe any data

1. It must not depend on the data

2. $p(\theta) \geq 0$ for all θ

3. $\int p(\theta) d\theta = 1$

likelihood **prior**

$$\text{Bayes formula: } p(\theta | D) = \frac{\underbrace{p(D|\theta)}_{\text{posterior distribution}} \cdot \underbrace{p(\theta)}_{\text{prior}}}{p(D)} \rightarrow \text{evidence}$$

It acts as a normalizing constant that ensures that the posterior distribution integrates to 1.

$$\text{posterior} \propto \text{likelihood} \cdot \text{prior} \Rightarrow p(\theta | D) \propto p(D|\theta) \cdot p(\theta)$$

$$p(D) = \int p(D, \theta) d\theta = \int p(D|\theta) p(\theta) d\theta$$

if $p(\theta) = 0$ for some particular θ , posterior will always be zero for that particular θ regardless of the likelihood / data.

Maximum a posteriori estimation (MAP)

MLE ignores our prior beliefs and performs poorly if little data is available

$$\theta_{\text{MAP}} = \underset{\theta}{\operatorname{argmax}} p(\theta | D) = \underset{\theta}{\operatorname{argmax}} \frac{p(D|\theta) p(\theta)}{p(D)} = \underset{\theta}{\operatorname{argmax}} p(D|\theta) p(\theta)$$

We can ignore $\frac{1}{p(D)}$ since it's a positive constant independent of θ

we choose Beta distribution for prior

$$\text{Beta}(\theta | a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \theta^{a-1} (1-\theta)^{b-1}, \quad \theta \in [0, 1]$$

where

- $a > 0, b > 0$ are the distribution parameters
- $\Gamma(n) = (n-1)!$ for $n \in \mathbb{N}$ is the gamma function

Putting everything together because $p(D)$ is constant

$$p(\theta | D) = \frac{p(D|\theta) \cdot p(\theta)}{p(D)} \propto p(D|\theta) \cdot p(\theta)$$

↑ ↑
likelihood prior
↓
posterior

We know

$$p(D|\theta) = \theta^{|T|} (1-\theta)^{|H|},$$

$$p(\theta) = p(\theta | a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \theta^{a-1} (1-\theta)^{b-1}$$

So:

$$p(\theta | D) \propto \theta^{|T|} (1-\theta)^{|H|} \cdot \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \theta^{a-1} (1-\theta)^{b-1}$$

$$\propto \theta^{|T|+a-1} (1-\theta)^{|H|+b-1}$$

We are looking for $\theta_{\text{MAP}} = \underset{\theta}{\operatorname{argmax}} p(\theta | D) = \underset{\theta}{\operatorname{argmax}} \theta^{|T|+a-1} (1-\theta)^{|H|+b-1} =$

$$\underset{\theta}{\operatorname{argmax}} \log p(\theta | D) = \underset{\theta}{\operatorname{argmax}} (|T|+a-1) \log \theta + (|H|+b-1) \log(1-\theta)$$

We obtain $\theta_{\text{MAP}} = \frac{|T|+a-1}{|H|+|T|+a+b-2}$

Estimating the posterior distribution

We need to consider the entire posterior distribution $p(\theta|D)$, not just its mode θ_{MAP}

$$p(\theta|D) \propto \theta^{|T|+\alpha-1} (1-\theta)^{|H|+b-1}$$

Finding the true posterior $p(\theta|D)$ boils down to finding the normalization constant, such that the distribution integrates to 1.

Sol. → Pattern matching

$$\text{Unnormalized posterior} \rightarrow p(\theta|D) \propto \theta^{|T|+\alpha-1} (1-\theta)^{|H|+b-1}$$

$$\text{Beta distribution} \rightarrow \text{Beta}(\theta|\alpha, \beta) = \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1-\theta)^{\beta-1}$$

Thus, we can conclude that

$$\text{the appropriate normalizing constant is} \rightarrow \frac{\Gamma(|T|+\alpha+|H|+b)}{\Gamma(|T|+\alpha)\Gamma(|H|+b)}$$

and the posterior is a Beta distribution $p(\theta|D) = \text{Beta}(\theta|\alpha+|T|, \beta+|H|)$

The prior is also a Beta distribution $p(\theta) = \text{Beta}(\theta|\alpha, \beta)$

This is not a coincidence, this is an instance of a more general principle.

Beta distribution is a conjugate prior for the Bernoulli likelihood.

If a prior is conjugate for the given likelihood, then the posterior will be of the same family as the prior.

Normal distribution (Gaussian)

$$N(\mu, \sigma^2) = \sqrt{\frac{\beta}{2\alpha}} \exp\left(-\frac{\beta}{2} (\mu - \bar{x})^2\right) \propto \exp\left(-\frac{\beta}{2} \bar{x}^2 + \beta \bar{x}\mu\right)$$

$$\text{mode of } p(\theta | D) = \text{Beta}(\alpha, \beta) = \frac{\alpha-1}{\alpha+\beta-2} \quad \begin{array}{l} \text{Relation between MAP, MLE and} \\ \text{posterior distribution} \end{array}$$

MAP is the mode of the posterior distribution

And if we choose a uniform prior ($\alpha = \beta = 1$) we obtain MLE

All this is due to we have chosen a conjugate prior. Had we chosen a non-conjugate prior, $p(\theta | D)$ and θ_{MAP} could not have a closed form

Predicting the next flip

We want to compute the probability that the next coin flip is H , given observations D and prior belief a, b :

$$p(F = \text{H} | D, a, b)$$

This distribution is called the **posterior predictive distribution**

$$f \in \{0, 1\}$$

$$p(F = f | D, a, b) = p(f | D, a, b) = \int_0^1 p(f | \theta) p(\theta | D, a, b) d\theta$$

Recall

$$p(f | \theta) = \text{Ber}(f | \theta) = \theta^f (1-\theta)^{1-f}$$

$$p(\theta | D, a, b) = \frac{\Gamma(|T|+a+|H|+b)}{\Gamma(|T|+a)\Gamma(|H|+b)} \theta^{|T|+a-1} (1-\theta)^{|H|+b-1}$$

$$= \frac{(|T|+a)^f (|H|+b)^{1-f}}{|T|+a+|H|+b} = \text{Ber}\left(f \mid \frac{|T|+a}{|T|+a+|H|+b}\right)$$

it does not contain θ

fully Bayesian analysis