# BAYESIAN INFERENCE

Frequentist view    vs    Bayesian view:

Example: You try to evaluate if a coin is fair: $P(H) = P(T) = \frac{1}{2}$

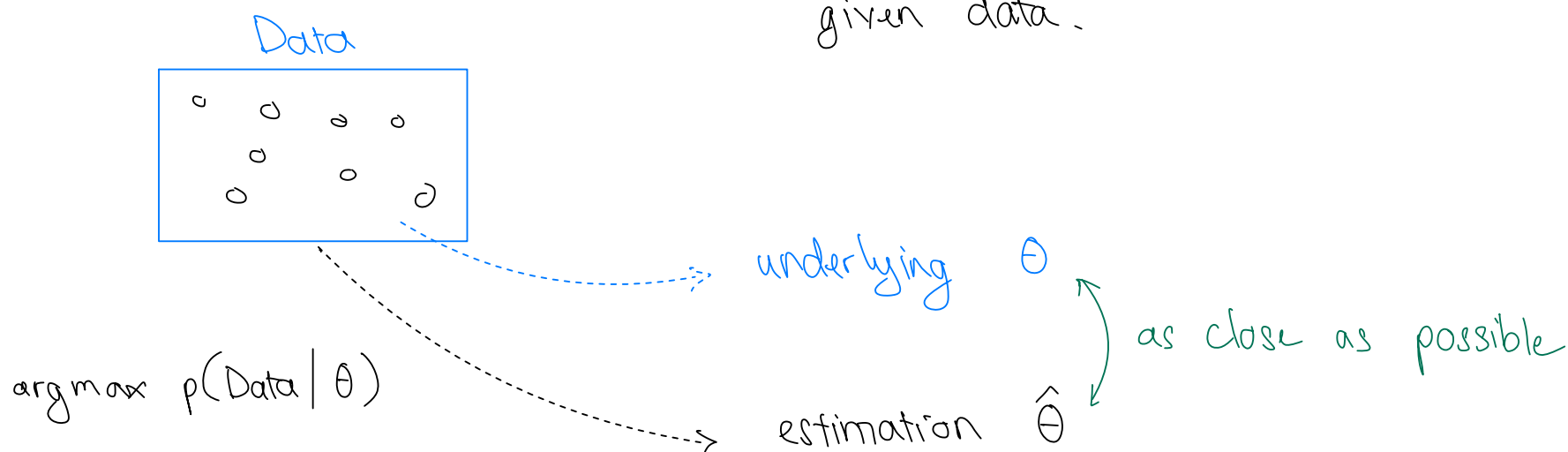| Frequentist | Bayesian |
|---|---|
| ◦ Run experiments: flip the coin 1000 times | ◦ Assume the coin is fair |
| ◦ Collect data | ◦ Collect data, update your assumption |
| ◦ Evaluate if the coin is fair | ◦ By 1000 times, conclude if it is fair or not |

Main Idea of Bayesian Inference:

◦ Recall Maximum Likelihood Estimation, we trying to find the underlying $\Theta$ given the dataset $p(\Theta | D)$ by maximizing the probability of seeing the data given $\Theta$:

$$\hat{\Theta} = \underset{\Theta}{\arg\max} \ p(D | \Theta)$$

estimation        given      unknown

This is the <u>frequentist</u> view: estimating unknown parameter through given data.

Data

$$\arg\max \ p(Data | \Theta)$$

underlying $\Theta$

} as close as possible

estimation $\hat{\Theta}$

◦ Look at the same problem under <u>Bayesian View:</u>

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

prior

posterior        Likelihood

Normalization factor:

$$p(D) = \int p(D | \Theta) \ p(\Theta) \ d\Theta$$

$p(D)$ is treated as constant since we only care ab $\Theta$. So:

$$p(\Theta | D) \propto p(D | \Theta) \cdot p(\Theta)$$

Bayes' rule

"proportional to"

# Example: Why Bayesian Inference is better than MLE (frequentist)

- In this example, we try to evaluate the accuracy of an alarm that goes off when the sun explodes.

  Given:
  $$\begin{cases} \theta \in \{0,1\} \text{ is indicator if the sun explodes} \\ x \in \{0,1\} \text{ is indicator if the alarm fires} \\ \alpha \in [0,1] \text{ is the error rate of the alarm} \end{cases}$$

  So:
  $$\begin{cases} p(x = \theta \mid \theta) = 1 - \alpha \quad \rightarrow \text{ correctly fired} \\ p(x = 1-\theta \mid \theta) = \alpha \quad \rightarrow \text{ incorrectly fired} \end{cases}$$

- Problem: If the alarm fires, should we believe the sun has exploded or not? Assume that this alarm is very accurate $\alpha = 0.0001$

## 1. Find $\theta$ using MLE

$$\hat{\theta} = \underset{\theta}{\text{argmax}} \; p(x=1 \mid \theta) = \begin{cases} \alpha & \text{if } \theta = 0 \\ 1-\alpha & \text{if } \theta = 1 \end{cases}$$

$$\Rightarrow \quad \hat{\theta} = 1$$

This suggest that we should trust the alarm completely!

## 2. Find $\theta$ using Bayesian Inference

- Step 1: Determine prior (probability of $\theta$ without observing any data)

  $$p(\theta) = \begin{cases} \beta & \text{if } \theta = 1 \\ 1-\beta & \text{if } \theta = 0 \end{cases} \qquad \langle \beta \text{ is very small}, \approx 0 \rangle$$

  prior knowledge is what makes Bayesian better

- Step 2: Determine posterior

  $$p(\theta \mid x = 1) = \frac{p(x=1 \mid \theta) \cdot p(\theta)}{p(x=1)}$$

  $$\propto p(x=1 \mid \theta) \cdot p(\theta)$$

  $$= \begin{cases} (1-\alpha) \cdot \beta & \text{if } \theta = 1 \\ \alpha(1-\beta) & \text{if } \theta = 0 \end{cases}$$

- Step 3: Decide value of $\theta$

- If we predict $\theta = 1$, then:

  $$p(\theta = 1 \mid x = 1) > p(\theta = 0 \mid x = 1)$$

  $$\Leftrightarrow (1-\alpha)\beta > \alpha(1-\beta)$$

  $$\Longleftrightarrow \frac{\beta}{1-\beta} > \frac{\alpha}{1-\alpha}$$

  This is not true since $\beta \approx 0$ and $\alpha = 0.0001$

  So the opposite holds, meaning $\theta = 0$

**Example:** Bayesian Inference on Gaussian Distribution

- You just moved to a new apartment.
- Your friend told you the commute time is $30 \pm 10$ min.
- You drove yourself and recorded the time: $25, 45, 30, 50$.

Problem: How should you predict the commute time?

Solution:

- Let $\theta$ be the time (treat $\theta$ as random variable)

Step 1: Determine prior, we assume $\theta$ is Normally distributed
$$P(\theta) \sim N(\mu_0, \sigma_0^2)$$
where $\begin{cases} \mu_0 = 30 \\ \sigma_0 = 10 \end{cases}$ → assign a reasonable value

Step 2: Observe the data collected.

- Let $D$ be the dataset $\{x_i\}_{i=1}^n$
- Assume the data observed have some noise:
$$x_i = \theta + \sigma_1 \varepsilon_i \qquad \text{where} \begin{cases} \varepsilon_i \sim N(0,1) \\ \sigma_1 : \text{variance} = \text{some value} \end{cases}$$

ground truth (true parameter) ··· noise

- From above assumption, we can say that each data point drawn i.i.d. from Gaussian Distribution:
$$P(x_i | \theta) \sim N(\theta, \sigma_1^2) \Rightarrow P(x_i | \theta) = \frac{1}{\sqrt{2\pi} \, \sigma_1} \exp\left(-\frac{(x_i - \theta)^2}{2\sigma_1^2}\right)$$

Step 3: Determine likelihood function
$$P(D | \theta) = \prod_{i=1}^n P(x_i | \theta)$$

Step 4: Determine posterior
$$P(\theta | D) = \frac{P(D|\theta) \cdot P(\theta)}{P(D)}$$

$P(D)$ → constant

Why remove this even though $\sigma_1$ is a r.v.?
Because the exponent part is much more influential

$$\propto P(D|\theta) \cdot P(\theta)$$
$$= \left[\prod_{i=1}^n P(x_i | \theta)\right] \cdot P(\theta)$$
$$\propto \left[\prod_{i=1}^n \exp\left(-\frac{(x_i - \theta)^2}{2\sigma_1^2}\right)\right] \cdot \exp\left(-\frac{(\theta - \mu_0)^2}{2\sigma_0^2}\right)$$
$$\propto \exp\left[-\sum_{i=1}^n \left(\frac{(\theta - x_i)^2}{2\sigma_1^2}\right) - \frac{(\theta - \mu_0)^2}{2\sigma_0^2}\right]$$
$$= \exp\left[-\frac{1}{2}\left(A\theta^2 - 2B\theta + \text{Const}\right)\right]$$

$$A = \sum_{i}^n \frac{1}{\sigma_1^2} + \frac{1}{\sigma_0^2}$$
$$B = \sum_{i}^n \frac{x_i}{\sigma_1^2} + \frac{\mu_0}{\sigma_0^2}$$

$$= \exp\left(-\frac{1}{2}A\left(\theta - \frac{B}{A}\right)^2 + \text{Const}\right)$$

$$\sim \boxed{N\left(\frac{B}{A}, \frac{1}{A}\right)}$$

So the posterior is Normally Distributed:

$$p(\theta \mid D) \sim N\left(\frac{B}{A}, \frac{1}{A}\right)$$

with:

posterior mean: $\mu_P = \dfrac{B}{A} = \dfrac{\overset{n}{\underset{}{\sum}} \dfrac{x_i}{\sigma_1^2} + \dfrac{\mu_0}{\sigma_0^2}}{\dfrac{n}{\sigma_1^2} + \dfrac{1}{\sigma_0^2}}$

→ weighted sum of data
weight = $\frac{1}{\sigma_1^2}$ which is the "inverse" noise

→ weighted prior
weight = $\frac{1}{\sigma_0^2}$

posterior variance: $\sigma_P^2 = \dfrac{1}{A} = \left(\dfrac{n}{\sigma_1^2} + \dfrac{1}{\sigma_0^2}\right)^{-1}$

→ Some observations:

- If $n = 0$ (no data): $\begin{cases} \mu_P = \mu_0 \\ \sigma_P^2 = \sigma_0^2 \end{cases}$

- If $n > 0$ (have data): $\begin{cases} \mu_P \text{ as stated above} \\ \sigma_P^2 \text{ as stated above} \end{cases}$

- If $n \to \infty$ (infinite data): $\begin{cases} \mu_P \approx \dfrac{\sum x_i}{n} \\ \sigma_P^2 \approx \dfrac{\sigma_1^2}{n} \end{cases}$

Estimate $\theta$ on $\mu_P$ and $\sigma_P$

In Bayesian, what is considered fixed and what is considered random variable?

- Fixed: Observed dataset
  i.e: $\{x_i\}_{i=1}^n$
  $\{x_i, y_i\}_{i=1}^n$
  ...

- Random variable:
  - Parameters: $\theta$
  - Noise: $\sigma, \varepsilon,$
  - Dataset before observed:
    i.e: $x_i \mid \theta$
    $y_i \mid x_i, \theta$

# Example: Bayesian Linear Regression

- Given data points $\{x_i, y_i\}_{i=1}^n$

- Problem: Find $\theta$ such that $y \approx x^T\theta$

$$y = x^T\theta$$
$$= (x_1 \; 1)^T \begin{pmatrix} \theta \\ b \end{pmatrix}$$
$$= x_1\theta + b$$

## Solved with LLS and MLE

- Linear Least Square (LLS): $\hat{\theta} = \underset{\theta}{\arg\min} \sum_{i=1}^n (y_i - x_i^T\theta)^2$

  or $\underset{\theta}{\arg\min} \; \| y - X^T\theta \|_2$

  → vector $\begin{pmatrix} \theta_1 \\ \vdots \\ b \end{pmatrix}$

  → vector

  → data set matrix

  records $\begin{bmatrix} \vdots & \vdots & \vdots & \vdots & 1 \\ & & & & 1 \\ & & & & 1 \end{bmatrix}$ features | bias

- Maximum Likelihood Estimation (MLE): $\hat{\theta} = \underset{\theta}{\arg\max} \; P(D \mid \theta)$

- ⟹ Both methods yield a ==determistic== value $\hat{\theta}$ (meaning fixed)

## Solve with Bayesian Inference

- Bayesian method gives us an ==uncertainty estimation== to see how accurate our estimation is.

- Step 1: Determine prior
  $$P(\theta) \sim N(\mu_0, \sigma_0^2) \quad,$$
  where $\begin{cases} \mu_0 = 0 \\ \sigma_0^2 = \text{some large number} \end{cases}$

  → default value of no prior knowledge

- Step 2: Determine likelihood function
  - Let $D$ be data set $\{x_i, y_i\}_{i=1}^n$
  - There should be some noise in data observed:
    $$y_i = x_i^T\theta + \sigma_1^2 \epsilon_i \quad \text{where} \begin{cases} \epsilon_1 \sim N(0,1) \\ \sigma_1^2 : \text{variance} \end{cases}$$

    noise

  - From above assumption, the likelihood function for each data point is:
    $$P(\{y_i, x_i\} \mid \theta) = P(y_i \mid x_i, \theta) \cdot P(x_i)$$
    $$\propto P(y_i \mid x_i, \theta)$$

    → fixed, dont care

Because
$y_i \mid x_i, \theta \sim N(x_i^T\theta, \sigma_1^2)$

$$= \boxed{\frac{1}{\sqrt{2\pi}\,\sigma_1} \exp\left(-\frac{(y_i - x_i^T\theta)^2}{2\sigma_1^2}\right)}$$

Because exponent part
is more influential

$$\propto \exp\left(-\frac{(y_i - x_i^T\theta)^2}{2\sigma_1^2}\right)$$

· So the likelihood function for whole dataset is

$$P(D \mid \theta) = \prod_{i=1}^{n} P(y_i \mid x_i, \theta)$$

$$\propto \prod_{i=1}^{n} \exp\left(-\frac{(y_i - x_i^T\theta)^2}{2\sigma_1^2}\right)$$

Step 4: Determine posterior

$$P(\theta \mid D) \propto P(D \mid \theta) \cdot P(\theta)$$

$$\propto \left[\prod_{i=1}^{n} \exp\left(-\frac{(y_i - x_i^T\theta)^2}{2\sigma_1^2}\right)\right] \cdot \exp\left(-\frac{(\theta - \mu_0)^2}{2\sigma_0^2}\right)$$

$$= \exp\left[\sum_{i=1}^{n}\left(-\frac{(y_i - x_i^T\theta)^2}{2\sigma_1^2}\right) - \frac{(\theta - \mu_0)^2}{2\sigma_0^2}\right]$$

$$= \exp\left[-\frac{1}{2}\left(\theta^T A\,\theta - 2B^T\theta + \text{Const}\right)\right]$$

$$\sim N\left(A^{-1}B, \ A^{-1}\right)$$

$$A = \sum^{n} \frac{x_i x_i^T}{\sigma_1^2} + \frac{I}{\sigma_0^2}$$

$$B = \sum^{n} \frac{y_i x_i}{\sigma_1^2} + \frac{\mu_0}{\sigma_2^2}$$

So the posterior is Normally Distributed

$$P(\theta \mid D) \sim N\left(A^{-1}B, A^{-1}\right)$$

with:

· posterior mean: $\mu_p = \left(\sum^{n} \frac{x_i x_i^T}{\sigma_1^2} + \frac{I}{\sigma_0^2}\right)^{-1}\left(\sum^{n} \frac{y_i x_i}{\sigma_1^2} + \frac{\mu_0}{\sigma_0^2}\right)$

· posterior variance: $\sigma_p^2 = \left(\sum^{n} \frac{x_i x_i^T}{\sigma_1^2} + \frac{I}{\sigma_0^2}\right)^{-1}$

→ Some observations:

· If $n = 0$ (no data): $\begin{cases} \mu_p = \mu_0 \\ \sigma_p^2 = \sigma_0^2 \end{cases}$

· If $n > 0$ (have data): as stated above

· If $n \to \infty$ (infinite data): $\begin{cases} \mu_p = \left(\sum^{\infty} \frac{x_i x_i^T}{\sigma_1^2}\right)^{-1}\left(\sum^{\infty} \frac{y_i x_i}{\sigma_1^2}\right) \\ \sigma_p^2 = \left(\sum^{\infty} \frac{x_i x_i^T}{\sigma_1^2}\right)^{-1} \end{cases}$

The distribution of prior, likelihood and posterior can be different.
In total, there are 3 distributions:

- Dataset distribution ( likelihood )
- Parameter distribution ( prior )
- Parameter | Dataset distribution ( posterior )