

Statistical Processes

Lecture 2

Physics 129AL

Zihang Wang
10/18/2023

Two Schools of Statistics

Recall what we discussed last time,

Frequentist Goal:

Construct procedures with frequency guarantees

Bayesian Goal:

Describe and update degree of belief in propositions

Probabilities

Recall that discrete or continuous event are represented as random variables,

- Random Variable: x
- Outcomes: $S \equiv \{x_1, x_2, \dots\}$

As an example, let consider the probabilities of a dice with the following properties,

$$p(E) : p_{\text{dice}}(1) = \frac{1}{6}$$

$$p(E) \geq 0$$

$$p(A \& B) = p(A) + p(B)$$

$$p(S) = 1 : \sum_i p(x_i) = 1$$

Probabilities

We can define **joint probability** of two random variables, x and y (not necessarily independent),

$$P(x, y) \quad x \in S_x, y \in S_y \quad P(x, y) = P(x)P(y)$$

Or **marginal Probability** of two random variables, x and y

$$P(x = x_i) = \sum_{y \in S_y} P(x = x_i, y)$$


Independent
random variables

And **conditional Probability** of two random variables, x and y

$$P(x = x_i \mid y = y_j) = \frac{P(x = x_i, y = y_j)}{P(y = y_j)} \quad \leftarrow \text{Marginal}$$

Probabilities

Product rule (Chain Rule, giving conditional),

$$P(x, y | \mathcal{H}) = P(x | y, \mathcal{H})P(y | \mathcal{H}) = P(y | x, \mathcal{H})P(x | \mathcal{H}).$$

where \mathcal{H} : Generative Model

Sum rule (giving marginal)

$$\begin{aligned} P(x | \mathcal{H}) &= \sum_y P(x, y | \mathcal{H}) \\ &= \sum_y P(x | y, \mathcal{H})P(y | \mathcal{H}). \end{aligned}$$

And **Bayes Theorem**

$$\begin{aligned} P(y | x, \mathcal{H}) &= \frac{P(x | y, \mathcal{H})P(y | \mathcal{H})}{P(x | \mathcal{H})} \\ &= \frac{P(x | y, \mathcal{H})P(y | \mathcal{H})}{\sum_{y'} P(x | y', \mathcal{H})P(y' | \mathcal{H})}. \end{aligned}$$

Probabilities

Example 2.3. Jo has a test for a nasty disease. We denote Jo's state of health by the variable a and the test result by b .

$$\begin{array}{ll} a = 1 & \text{Jo has the disease} \\ a = 0 & \text{Jo does not have the disease.} \end{array} \quad (2.12)$$

The result of the test is either 'positive' ($b = 1$) or 'negative' ($b = 0$); the test is 95% reliable: in 95% of cases of people who really have the disease, a positive result is returned, and in 95% of cases of people who do not have the disease, a negative result is obtained. The final piece of background information is that 1% of people of Jo's age and background have the disease.

OK – Jo has the test, and the result is positive. What is the probability that Jo has the disease?

Probabilities

Step 1: Write down all probabilities

We are given conditional probabilities

$$\begin{aligned} P(b=1 | a=1) &= 0.95 & P(b=1 | a=0) &= 0.05 \\ P(b=0 | a=1) &= 0.05 & P(b=0 | a=0) &= 0.95; \end{aligned}$$

And marginal probability of a

$$P(a=1) = 0.01 \quad P(a=0) = 0.99.$$

We want $P(a=1|b=1)$

Probabilities

Step 2: Deduce joint probability $P(a,b)$

$$P(a,b)=P(a|b)P(b)=P(b|a)P(a)$$

Step 3: conditional probability $P(a|b)$

$$\begin{aligned}P(a=1 | b=1) &= \frac{P(b=1 | a=1)P(a=1)}{P(b=1 | a=1)P(a=1) + P(b=1 | a=0)P(a=0)} \\&= \frac{0.95 \times 0.01}{0.95 \times 0.01 + 0.05 \times 0.99} \\&= 0.16.\end{aligned}$$

Lots of false positives!

Statistical Properties: Moments

- **Expectation value** $\langle F(x) \rangle = \int_{-\infty}^{\infty} dx \, p(x) F(x).$
- **Moments** $m_n \equiv \langle x^n \rangle = \int dx p(x) x^n.$
- **Characteristic function** generates moments:

$$\tilde{p}(k) = \langle e^{-ikx} \rangle = \int dx p(x) e^{-ikx}.$$

Fourier Transform of the PDF: contains same information as PDF $p(x)$

Recovering the PDF from the characteristic function through the inverse F.T.:

$$p(x) = \frac{1}{2\pi} \int dk \tilde{p}(k) e^{+ikx}. \quad \int e^{-i(w-w')\tau} d\tau = 2\pi \delta_D(w - w')$$

Moments of the distribution: $\int \frac{dw'}{2\pi} 2\pi \delta_D(w - w') = 1$

$$\tilde{p}(k) = \left\langle \sum_{n=0}^{\infty} \frac{(-ik)^n}{n!} x^n \right\rangle = \sum_{n=0}^{\infty} \frac{(-ik)^n}{n!} \langle x^n \rangle.$$

Infinite series that in principle contains same information as $p(x)$

Statistical Properties: Moments

PDF moments around x_0 :

$$e^{ikx_0} \tilde{p}(k) = \langle e^{-ik(x-x_0)} \rangle = \sum_{n=0}^{\infty} \frac{(-ik)^n}{n!} \langle (x-x_0)^n \rangle.$$

- Cumulant generating function defines cumulants $\langle x^n \rangle_c$

$$\ln \tilde{p}(k) = \sum_{n=1}^{\infty} \frac{(-ik)^n}{n!} \langle x^n \rangle_c.$$

We can obtain relations between moments and cumulants using

$$\ln(1+\epsilon) = \sum_{n=1}^{\infty} (-1)^{n+1} \frac{\epsilon^n}{n}.$$

We expand ϵ into a sum of moments and $\ln(1+\epsilon)$ into a sum of cumulants and match powers of k to obtain

Mean $\langle x \rangle_c = \langle x \rangle,$

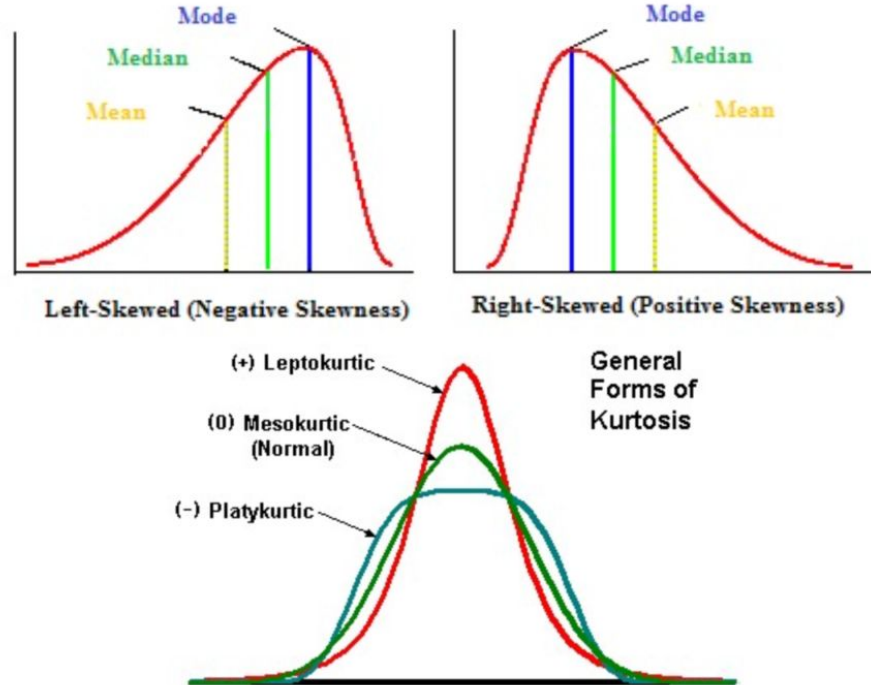
Variance $\langle x^2 \rangle_c = \langle x^2 \rangle - \langle x \rangle^2,$

Skewness $\langle x^3 \rangle_c = \langle x^3 \rangle - 3\langle x^2 \rangle \langle x \rangle + 2\langle x \rangle^3,$

Curtosis $\langle x^4 \rangle_c = \langle x^4 \rangle - 4\langle x^3 \rangle \langle x \rangle - 3\langle x^2 \rangle^2 + 12\langle x^2 \rangle \langle x \rangle^2 - 6\langle x \rangle^4.$

Statistical Properties: Moments

Moments characterize the probability (mass) density function.



Statistical Properties: Cumulants

Moments are related to clusters of cumulants: useful pictorial that relates cumulants to moments. Cumulants are also called connected moments.

$$\langle x \rangle = \bullet$$

$$\langle x \rangle = \langle x \rangle_c$$

$$\langle x^2 \rangle = \langle \bullet \bullet \rangle + \bullet \bullet$$

$$\langle x^2 \rangle = \langle x^2 \rangle_c + \langle x \rangle_c^2$$

$$\langle x^3 \rangle = \langle \bullet \bullet \bullet \rangle + 3 \langle \bullet \bullet \rangle \bullet + \bullet \bullet \bullet$$

$$\langle x^3 \rangle = \langle x^3 \rangle_c + 3 \langle x^2 \rangle_c \langle x \rangle_c + \langle x \rangle_c^3$$

$$\langle x^4 \rangle = \langle \bullet \bullet \bullet \bullet \rangle + 4 \langle \bullet \bullet \bullet \rangle \bullet + 3 \langle \bullet \bullet \rangle \langle \bullet \bullet \rangle + 6 \langle \bullet \bullet \rangle \bullet \bullet + \bullet \bullet \bullet \bullet$$

$$\langle x^4 \rangle = \langle x^4 \rangle_c + 4 \langle x^3 \rangle_c \langle x \rangle_c + 3 \langle x^2 \rangle_c^2 + 6 \langle x^2 \rangle_c \langle x \rangle_c^2 + \langle x \rangle_c^4$$

Example: Gaussian

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(x-\lambda)^2}{2\sigma^2}\right], \quad -\infty < x < \infty$$

- Characteristic Function

$$\tilde{p}(k) = \int_{-\infty}^{\infty} dx \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(x-\lambda)^2}{2\sigma^2} - ikx\right] = \exp\left[-ik\lambda - \frac{k^2\sigma^2}{2}\right]$$

- Cumulants From $\ln \tilde{p}(k) = -ik\lambda - k^2\sigma^2/2$, we can identify:

$$\langle x \rangle_c = \lambda, \quad \langle x^2 \rangle_c = \sigma^2, \quad \langle x^3 \rangle_c = \langle x^4 \rangle_c = \dots = 0.$$

- Moments from cluster expansion

$$\begin{aligned} \langle x \rangle &= \lambda, & \langle x^3 \rangle &= 3\sigma^2\lambda + \lambda^3, \\ \langle x^2 \rangle &= \sigma^2 + \lambda^2, & \langle x^4 \rangle &= 3\sigma^4 + 6\sigma^2\lambda^2 + \lambda^4, \end{aligned}$$

Many Random Variables

- Joint PDF $p(\mathbf{x})$ where $\mathbf{x} = \{x_1, x_2, \dots, x_N\}$

$$p_{\mathbf{x}}(\mathcal{S}) = \int d^N \mathbf{x} p(\mathbf{x}) = 1. \quad (\mathcal{S} \text{ is the set of all outcomes})$$

- Joint characteristic function $\tilde{p}(\mathbf{k}) = \left\langle \exp \left(-i \sum_{j=1}^N k_j x_j \right) \right\rangle$

Example

$$\langle x_1 x_2 \rangle = \begin{array}{c} \bullet \bullet \\ 1 \ 2 \end{array} + \begin{array}{c} \bullet \bullet \\ \text{---} \\ 1 \ 2 \end{array} \quad \langle x_1 x_2 \rangle = \langle x_1 \rangle_c \langle x_2 \rangle_c + \langle x_1 * x_2 \rangle_c$$

$$\langle x_1^2 x_2 \rangle = \begin{array}{c} 2 \\ \bullet \bullet \\ 1 \ 1 \end{array} + \begin{array}{c} 2 \\ \bullet \bullet \\ \text{---} \\ 1 \ 1 \end{array} + 2 \begin{array}{c} 1 \\ \bullet \bullet \\ 1 \ 2 \end{array} + \begin{array}{c} 2 \\ \bullet \bullet \\ \text{---} \\ 1 \ 1 \end{array}$$

$$\langle x_1^2 x_2 \rangle = \langle x_1 \rangle_c^2 \langle x_2 \rangle_c + \langle x_1^2 \rangle_c \langle x_2 \rangle_c + 2 \langle x_1 * x_2 \rangle_c \langle x_1 \rangle_c + \langle x_1^2 * x_2 \rangle_c$$

$\langle x_\alpha * x_\beta \rangle_c$ is zero if x_α and x_β are independent (symbol * here means product)

Multi-variate Gaussian

$$p(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^N \det[C]}} \exp \left[-\frac{1}{2} \sum_{mn} (C^{-1})_{mn} (x_m - \lambda_m)(x_n - \lambda_n) \right]$$

where C is the covariance matrix, and C^{-1} is its inverse

$$\tilde{p}(\mathbf{k}) = \exp \left[-ik_m \lambda_m - \frac{1}{2} C_{mn} k_m k_n \right]$$

Note: repeated indices are summed

- Cumulants

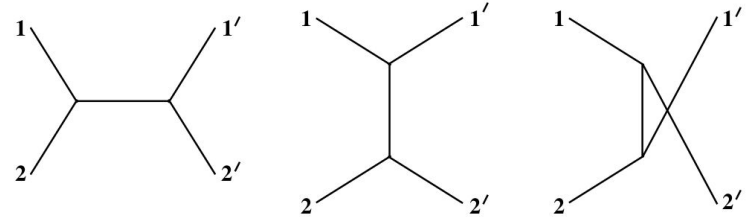
$$\langle x_m \rangle_c = \lambda_m, \quad \langle x_m * x_n \rangle_c = C_{mn}, \quad \text{rest are zero}$$

Wick's Theorem

4 point function for a multi-variate gaussian with zero mean depends only on products of 2 point functions, $\langle x_a x_b x_c x_d \rangle = C_{ab}C_{cd} + C_{ac}C_{bd} + C_{ad}C_{bc}$

The three tree-level Feynman diagrams that contribute to the connected correlation function are related to the Wick's theorem,

$$\begin{aligned} & \langle 0 | T \varphi(x_1) \varphi(x_2) \varphi(x'_1) \varphi(x'_2) | 0 \rangle_C \\ &= (ig)^2 \left(\frac{1}{i} \right)^5 \int d^4y d^4z \Delta(y-z) \\ & \quad \times \left[\Delta(x_1-y) \Delta(x_2-y) \Delta(x'_1-z) \Delta(x'_2-z) \right. \\ & \quad + \Delta(x_1-y) \Delta(x'_1-y) \Delta(x_2-z) \Delta(x'_2-z) \\ & \quad \left. + \Delta(x_1-y) \Delta(x'_2-y) \Delta(x_2-z) \Delta(x'_1-z) \right] \\ & \quad + O(g^4) . \end{aligned}$$



Sum of variables

- **PDF** $p_X(x) = \int d^N \mathbf{x} p(\mathbf{x}) \delta(x - \sum x_i)$

- **Characteristic function**

$$\tilde{p}_X(k) = \left\langle \exp \left(-ik \sum_{j=1}^N x_j \right) \right\rangle = \tilde{p}(k_1 = k_2 = \dots = k_N = k)$$

- **Cumulants**

$$\ln \tilde{p}(k_1 = k_2 = \dots = k_N = k) = -ik \sum_{i_1=1}^N \langle x_{i_1} \rangle_c + \frac{(-ik)^2}{2} \sum_{i_1, i_2}^N \langle x_{i_1} x_{i_2} \rangle_c + \dots$$

$$\langle X \rangle_c = \sum_{i=1}^N \langle x_i \rangle_c, \quad \langle X^2 \rangle_c = \sum_{i,j}^N \langle x_i x_j \rangle_c, \dots$$

If variables are **independent** cross-cumulants vanish ->

$$\langle X^n \rangle_c = \sum_{i=1}^N \langle x_i^n \rangle_c \quad \longrightarrow \quad \langle X^n \rangle_c = N \langle x^n \rangle_c$$

If all drawn from $p(x)$

Central Limit Theorem

For large N $\langle x \rangle \propto N$, $\langle (x - \langle x \rangle)^2 \rangle \propto N$

$$y = \frac{x - N\langle x \rangle_c}{\sqrt{N}}, \quad \langle y^n \rangle_c \propto N^{1-n/2}$$

$$N \rightarrow \infty, \quad \langle y^n \rangle_c \rightarrow 0 \text{ for } n > 2$$

$$\lim_{N \rightarrow \infty} p\left(y = \frac{\sum_{i=1}^N x_i - N\langle x \rangle_c}{\sqrt{N}}\right) = \frac{1}{\sqrt{2\pi \langle x^2 \rangle_c}} \exp\left(-\frac{y^2}{2 \langle x^2 \rangle_c}\right)$$

Gaussian Distribution

Poisson Distribution

Example: Radioactive decay,

Probability of one and only one event (decay) in $[t, t+dt]$ is proportional to dt as $dt \rightarrow 0$.

Probabilities of events are independent.

Poisson $p(M|T)$ M events in time interval T

- **Limit of binomial:** $N = \frac{T}{dt} \gg 1$

$$p_1 = \alpha dt, \quad p_0 = 1 - \alpha dt, \quad p_2 \propto dt^2 \rightarrow 0$$

- We imagine many binomial events over time interval T

Poisson Distribution

Assume stars randomly distributed around us with density n , what is probability that the nearest star is at distance R ?

$$p(r < R, 0 \text{ star}|n) \sim \frac{(\gamma T)^0}{0!} e^{-n \frac{4\pi}{3} R^3}$$

$$p(r \approx R, 1 \text{ star}|n) \sim \frac{(n4\pi R^2 dR)^1}{1!} e^{-n4\pi R^2 dR}$$

$$\begin{aligned} p(r < R, 0 \text{ star}, r \approx R, 1 \text{ star}|n) &= p(r < R, 0 \text{ star}|n) p(r \approx R, 1 \text{ star}|n) \\ &\sim n4\pi R^2 dR e^{-n4\pi R^2 dR} e^{-n \frac{4\pi}{3} R^3} \approx n4\pi R^2 dR e^{-n \frac{4\pi}{3} R^3} \end{aligned}$$

What do we mean by Probability

- 1) Frequency of outcomes for repeated random experiments
- 2) Degrees of belief in propositions not involving random variables (quantifying uncertainty)

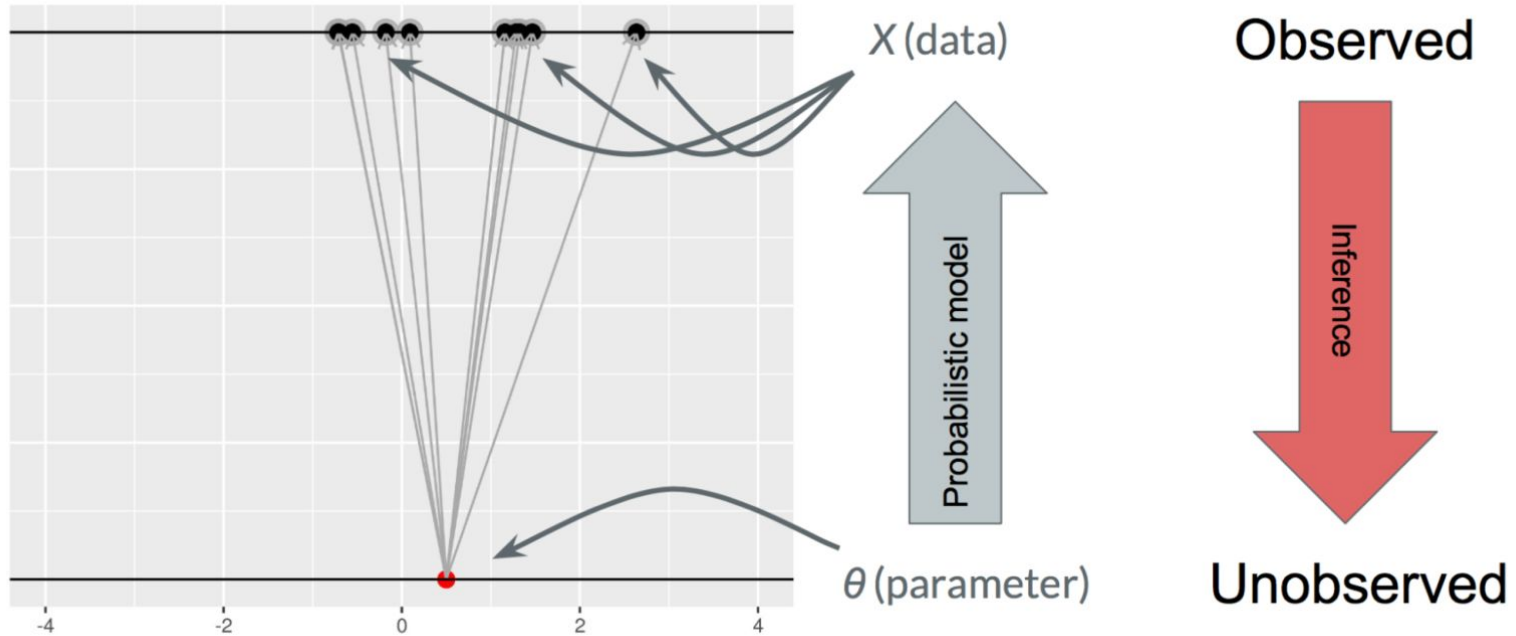
Forward Probability

Generative model describing a process giving rise to some data.

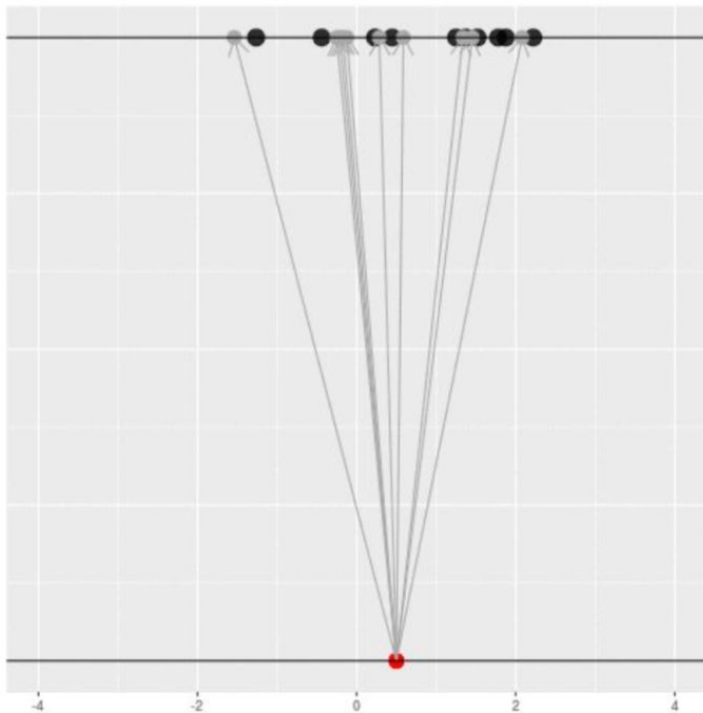
Inverse Probability

We compute probability of some unobserved quantity, given the observed variables, i.e. Bayes theorem.

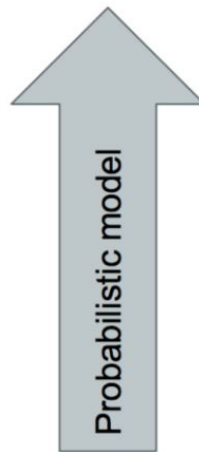
Parameters and Data



Frequentist Inference



X (data)



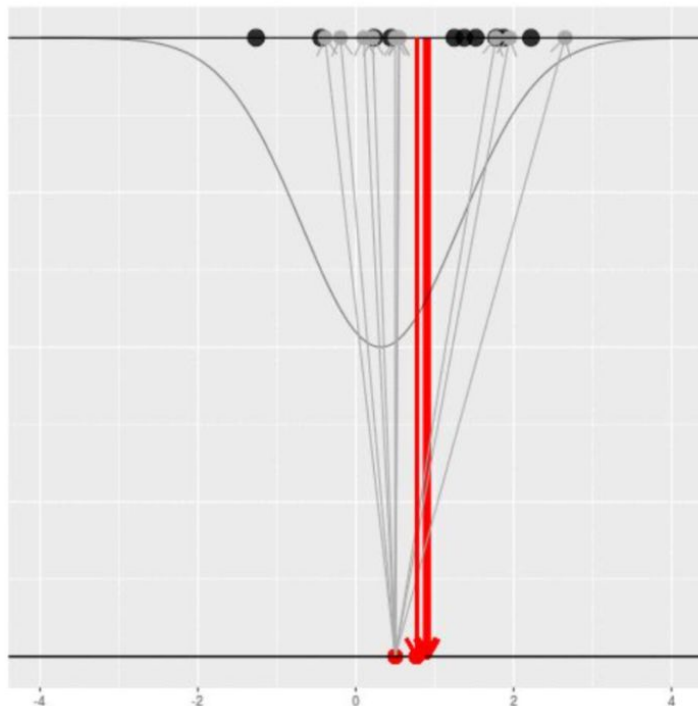
θ (parameter)

Frequentist idea:

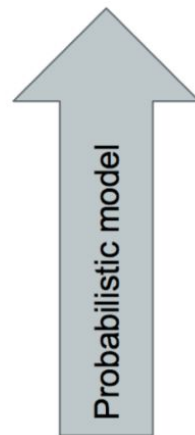
We got the parameter indicated by the red dot and saw the dataset in black.

But the same parameter could have given us lots of other datasets.

Frequentist Inference



X (data)



θ (parameter)

Frequentist idea:

For each dataset, we might pick some summary function and call it an “estimate”.

It will be different each time because the data will be different each time.

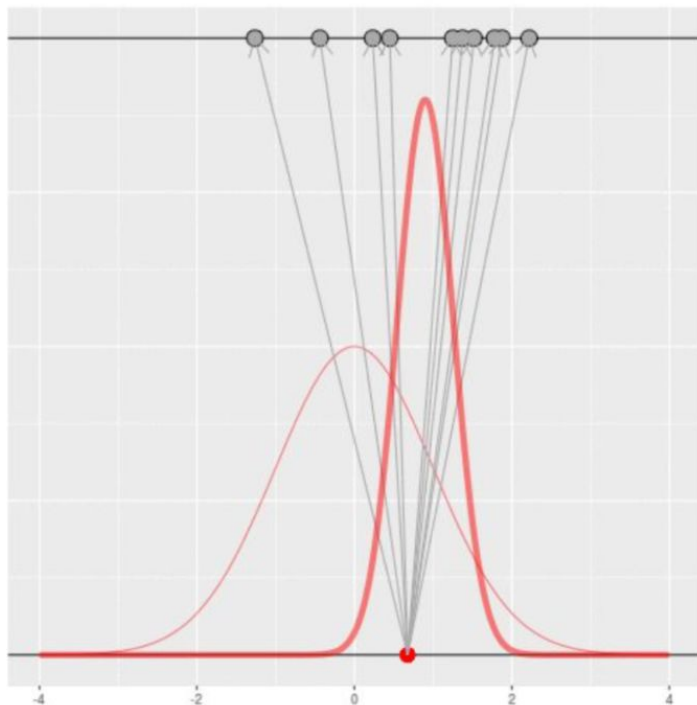
A typical estimate -- but not the only one -- is the value that maximizes the likelihood of the data.

We hope the estimate is usually near the true parameter in some sense.

Frequentist Inference

The frequentist uses the concept of frequency or repeated sampling. It focuses on estimating statistical parameters, e.g. means and variances and making statistical inferences based on the given data. In frequentist statistics, parameters governing the underlying distribution are treated as **fixed, unknown values**, and the goal is to estimate these parameters using **point estimates** (e.g., maximum likelihood) or confidence intervals. The frequentist approach does not incorporate prior beliefs on the parameter p (it is given and fixed), and **it relies solely on the data at hand**.

Parameters and Data



X (data)



θ (parameter)

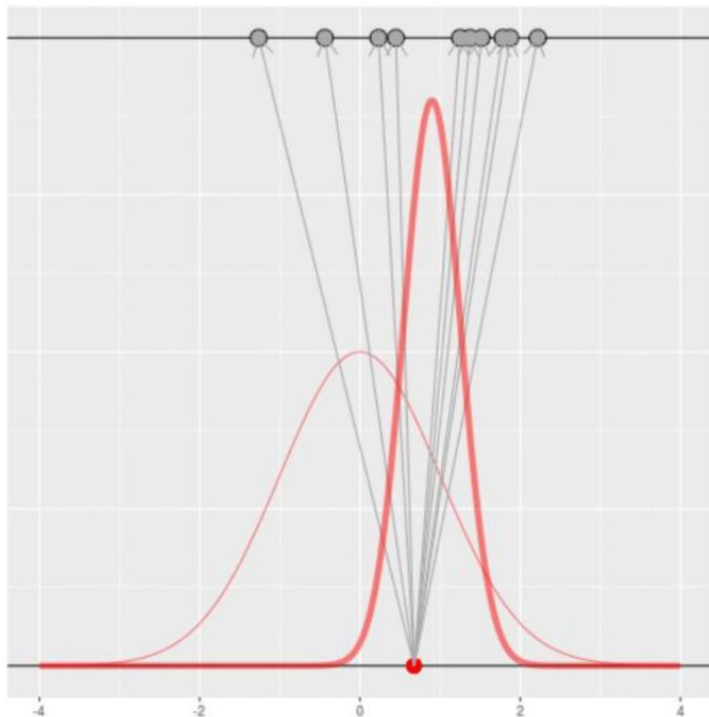
Bayesian idea:

Suppose we draw a bunch of parameters and datasets, and then throw out every pair where the data doesn't match what we observed.

The distribution of the parameters that are left represents which parameters could have given us the dataset we saw.

We hope the prior is reasonable and the model is accurate.

Parameters and Data



Of course, in practice you don't usually generate parameters and data hoping to get your original dataset.

Instead, you use Bayes' rule:

$$p(\theta|x) = \frac{p(x|\theta) p(\theta)}{p(x)}$$

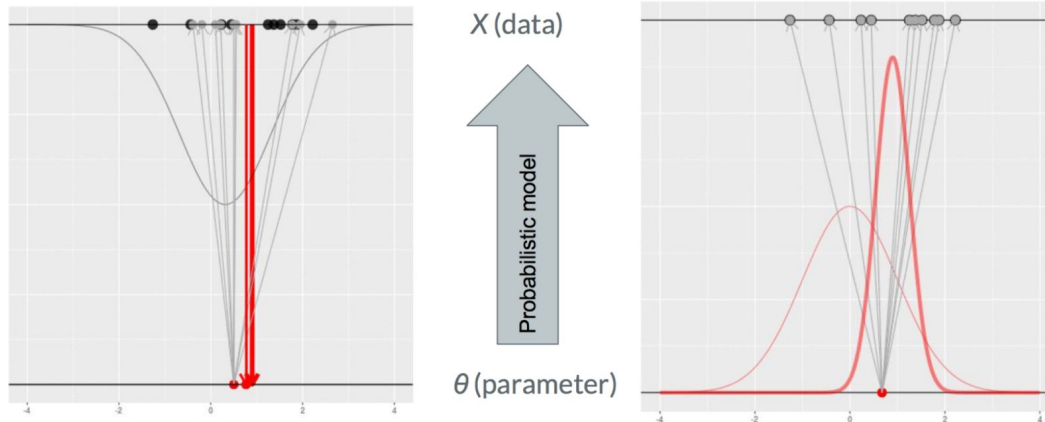
This is intractable in general (the denominator is a problem). Turn to approximations schemes like MCMC, variational Bayes, &c.

Bayesian Inference

Bayesian inference treats both observed data and parameters as random variables with probability distributions, in contrast to the frequentist approach. Bayesian inference typically begins with **prior beliefs or knowledge about the controlling parameters**. The likelihood function quantifies the probability of observing the data with a specific set of controlling parameters and often relies on certain assumptions. The posterior probability represents the **updated beliefs about controlling parameters after integrating the observed data**. It can then be employed as the new prior for subsequent observations, enabling the continuous refinement of beliefs in the presence of additional data.

Frequentist v.s. Bayesian Inference

Bayesian inference incorporates subjective prior beliefs and provides posterior probability distributions of the controlling parameters, while frequentist inference focuses on objective measures based solely on observed data. The choice between the two approaches often depends on the specific problem and available data.



Acknowledgement

The slides are partially developed or inspired by Professor Uros Seljak at UC Berkeley. For more information, please visit the github page.

A short story: When I was doing my undergraduate, I took his class (the exact same one, linked below).

<https://phy151-ucb.github.io/seljak-phy151-fall-2018/#course-syllabus>



Uros Seljak (Berkeley)