

Lec 1: Recap of Probability & Statistics

Yanjun Han

Sept. 3, 2024



Probability: mathematical model of random outcomes

(Gaussian, Poisson, Markov chains, Brownian motion, ...)

Statistics: given random outcomes, infer the underlying model

(statistical modeling, parameter estimation, testing, confidence interval, regression, ...)

Central question for this course: given data x_1, \dots, x_n ,

1. How to find a mathematical model of P_θ such that $(x_1, \dots, x_n) \sim P_\theta$?
2. How should we infer the unknown parameter θ ?

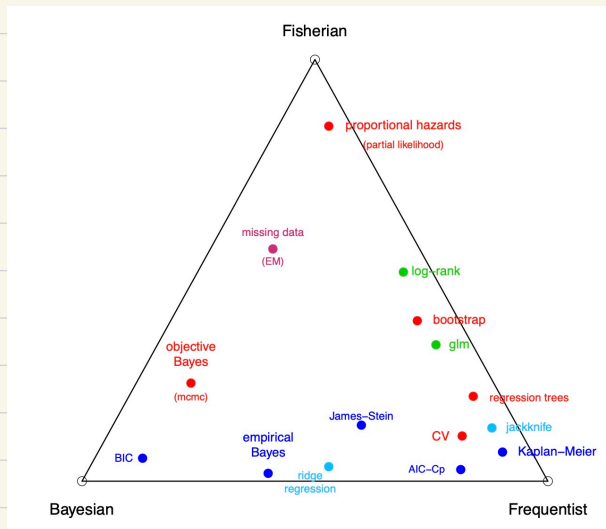
Three classes of models we'll cover:

Parametric: $\theta \in \mathbb{R}^d$ finite-dimensional, typically with a linear structure

Semiparametric: $\theta = (\tau, \eta)$, τ : parameter of interest
 η : nuisance parameter

Nonparametric: $\theta = f$ is parametrized by a function

An (incomplete) list of topics we'll cover:



Limit theorems in probability

Surprising fact (origins of probability theory): sums of independent, identically distributed (i.i.d.) RVs have universal behavior.

Law of large numbers (LLN)

X_1, \dots, X_n i.i.d., $\mathbb{E}[X]$ exists, then

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \rightarrow \mathbb{E}X \quad \text{as } n \rightarrow \infty \quad (\text{in probability \& almost surely})$$

Central limit theorem (CLT)

For iid X_1, \dots, X_n with $\mathbb{E}X = \mu$, $\text{Var}(X) = \sigma^2$,

$$Z_n = \frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \xrightarrow{d} N(0, 1) \quad \text{as } n \rightarrow \infty$$

\nearrow standard normal distribution

\uparrow
converges in distribution, meaning that

$$P(Z_n \leq t) \rightarrow P(N(0, 1) \leq t) \quad \text{for every } t \in \mathbb{R}.$$

Delta method

If g is differentiable & $\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{d} Z \sim N(0, 1)$. Then

$$\sqrt{n}(g(\bar{X}_n) - g(\mu)) \xrightarrow{d} g'(\mu) \cdot Z \sim N(0, g'(\mu)^2)$$

Pf (Taylor expansion)

$$g(\bar{X}_n) = g(\mu) + g'(\mu)(\bar{X}_n - \mu) + o(|\bar{X}_n - \mu|)$$

$$\Rightarrow \sqrt{n}(g(\bar{X}_n) - g(\mu)) = g'(\mu) \underbrace{\sqrt{n}(\bar{X}_n - \mu)}_{\xrightarrow{d} Z} + \underbrace{o(\sqrt{n}|\bar{X}_n - \mu|)}_{\rightarrow 0}$$

$$\rightarrow g'(\mu) Z.$$

□

Example. $g(x) = x^2$, $g'(\mu) = 2\mu$

$$\Rightarrow \sqrt{n}(\bar{X}_n^2 - \mu^2) \xrightarrow{d} N(0, 4\mu^2\sigma^2).$$

Example of how probability is applied to statistics:

Suppose we're watching the US open, where a player is on 1st serve 100 times, and wins 80 of them.

Statistical model: the wins/losses are independent, with an unknown win rate p

LLN: $\frac{1}{n} \# \text{ wins} \rightarrow p$ as $n \rightarrow \infty$.

Since $n=100$ is large enough, $\hat{p} = 0.8$ is a reasonable estimate of p .

CLT: $\frac{1}{\sqrt{np(1-p)}} (\# \text{ wins} - np) \xrightarrow{d} N(0,1)$ as $n \rightarrow \infty$

\Rightarrow the estimation error $\hat{p} - p \approx N(0, 0.04^2)$, so a 95% confidence interval for p is $p \in [\hat{p} - 2 \cdot 0.04, \hat{p} + 2 \cdot 0.04] = [0.72, 0.86]$

Estimation: given $X_1, \dots, X_n \stackrel{\text{i.i.d.}}{\sim} P_\theta$ with a known distribution family $\theta \mapsto P_\theta$ but an unknown parameter θ , the target of estimation is to find θ .

Approach I: estimating equation

Suppose one can find functions F_1, F_2, \dots, F_p s.t.

$$\mathbb{E}_{X \sim P_\theta} [F_j(\theta, X)] = 0 \quad \forall j = 1, \dots, p.$$

Then a reasonable estimator $\hat{\theta}_n$ is defined as the solution to

$$\begin{cases} \frac{1}{n} \sum_{i=1}^n F_1(\theta, X_i) = 0 \\ \vdots \\ \frac{1}{n} \sum_{i=1}^n F_p(\theta, X_i) = 0 \end{cases}$$

Analysis: ① By LLN, $0 = \frac{1}{n} \sum_{i=1}^n F_j(\hat{\theta}_n, X_i) \approx \mathbb{E}_{X \sim P_\theta} [F_j(\hat{\theta}_n, X)]$
 $\Rightarrow \hat{\theta}_n \approx \theta$ (i.e. "consistency")

② CLT can also be used to establish the asymptotic normality of $\sqrt{n}(\hat{\theta}_n - \theta)$

Example: $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2)$ with unknown (μ, σ^2)

$$\begin{aligned} \text{then } \mathbb{E}_{\mu, \sigma^2}[X - \mu] &= 0 & \Rightarrow F_1(X, (\mu, \sigma^2)) &= X - \mu \\ \mathbb{E}_{\mu, \sigma^2}[X^2 - \mu^2 - \sigma^2] &= 0 & F_2(X, (\mu, \sigma^2)) &= X^2 - \mu^2 - \sigma^2 \end{aligned}$$

The estimator $(\hat{\mu}, \hat{\sigma}^2)$ solves

$$\begin{cases} \frac{1}{n} \sum_{i=1}^n (X_i - \hat{\mu}) = 0 \\ \frac{1}{n} \sum_{i=1}^n (X_i^2 - \hat{\mu}^2 - \hat{\sigma}^2) = 0 \end{cases} \Rightarrow \begin{cases} \hat{\mu} = \bar{X} \\ \hat{\sigma}^2 = \overline{X^2} - \bar{X}^2 \end{cases}$$

Approach II: MLE (maximum likelihood estimator)

Def (likelihood).

Suppose X_1, \dots, X_n have joint pdf $f_\theta(X_1, \dots, X_n)$ (or pmf $p_\theta(X_1, \dots, X_n)$).

The likelihood function is

$$L_n(\theta) = f_\theta(X_1, \dots, X_n) \quad (\text{pdf viewed as function of } \theta)$$

The log likelihood function is

$$l_n(\theta) = \log L_n(\theta) = \log f_\theta(X_1, \dots, X_n).$$

$$\underline{\text{MLE}}: \quad \hat{\theta}_n = \arg\max_{\theta} L_n(\theta) = \arg\max_{\theta} l_n(\theta)$$

$$\underline{\text{Example (cont'd)}}: f_{\mu, \sigma^2}(X_1, \dots, X_n) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(X_i - \mu)^2}{2\sigma^2}\right)$$

$$\Rightarrow l_n(\mu, \sigma^2) = -\frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu)^2 - \frac{n}{2} \log(2\pi\sigma^2)$$

$$\text{F.O.C.: } \begin{cases} \frac{\partial l_n}{\partial \mu} = 0 \\ \frac{\partial l_n}{\partial \sigma^2} = 0 \end{cases} \Rightarrow \begin{cases} \hat{\mu} = \bar{X} \\ \hat{\sigma}^2 = \overline{X^2} - \bar{X}^2 \end{cases}$$

(more on both topics in future lectures)

Testing: given $X_1, \dots, X_n \sim p_\theta$, we'd like to test between

$$H_0: \theta \in \Theta_0 \quad \text{vs.} \quad H_1: \theta \in \Theta_1 \quad (\Theta_0 \cap \Theta_1 = \emptyset)$$

(null) (alternative)

Test: a function $\Omega \longrightarrow \{H_0, H_1\}$

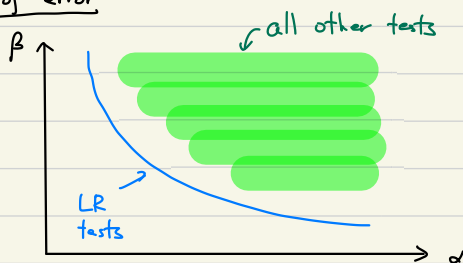
\uparrow set of outcomes \uparrow accept H_0 \uparrow reject H_0

Test \ Truth	H_0	H_1
H_0	Correct	Type I error (false positive)
H_1	Type II error (false negative)	Correct

Probability \ Truth	$P(\text{output } H_0)$	$P(\text{output } H_1)$
H_0	$1 - \alpha$	α (prob of type I error, or (significance) level)
H_1	β (prob. of type II error)	$1 - \beta$ (power of the test)

Fundamental tradeoff between two types of error

making smaller type I error
 \Rightarrow setting higher bar for rejecting H_0
 \Rightarrow larger type II error



Simple hypothesis testing: $H_0: \theta = \theta_0$ vs. $H_1: \theta = \theta_1$

Neyman-Pearson lemma: the likelihood ratio test is optimal

$$\text{output} \begin{cases} H_0 & \text{if } \frac{P(X_1, \dots, X_n | H_0)}{P(X_1, \dots, X_n | H_1)} > c \\ H_1 & \text{if } \frac{P(X_1, \dots, X_n | H_0)}{P(X_1, \dots, X_n | H_1)} \leq c \end{cases}$$

Composite hypothesis testing: $|\Theta_0| > 1$ and/or $|\Theta_1| > 1$

Unfortunately, no complete picture here. So statisticians have made:

- ① compromise I: focus only on significance level
- ② compromise II: focus only on asymptotic tests ($n \rightarrow \infty$)

Idea: find a function $F(\theta, X)$ s.t. $F(\theta, X)$ (asymptotically) follows a known distribution P for every $\theta \in \Theta_0$. (e.g. by CLT),
given a significant level α , find A s.t. $P(A) = 1 - \alpha$;

level- α test: reject $H_0: \theta = \theta_0$ if $F(\theta_0, X) \notin A$

$(1-\alpha)$ -confidence interval for θ : $C = \{\theta: F(\theta, X) \in A\}$.

Example. Given $X \sim B(n, p)$, then CLT gives

$$\frac{X - np}{\sqrt{np(1-p)}} \xrightarrow{d} N(0,1) \text{ as } n \rightarrow \infty.$$

level- α test for $H_0: p = p_0$: reject H_0 iff $\left| \frac{X - np_0}{\sqrt{np_0(1-p_0)}} \right| > z_{\alpha/2}$

$(1-\alpha)$ -confidence interval for p : $C = \{p: \left| \frac{X - np}{\sqrt{np(1-p)}} \right| \leq z_{\alpha/2}\}$.

p-value: every test can be equivalently represented by a p-value $\in [0,1]$.

and rejects H_0 iff p-value $\leq \alpha$

(p-value contains more information than yes/no answers)

Generalized likelihood ratio test: $LR = \frac{\max_{\theta \in \Theta_0} L_n(\theta)}{\max_{\theta \in \Theta_0 \cup \Theta_1} L_n(\theta)}$ ↖ likelihood function

Wilk's thm, under mild conditions, $-2 \log LR \xrightarrow{d} \chi_d^2$ under H_0
with $d = \dim(\Theta_0 \cup \Theta_1) - \dim(\Theta_0)$

Regression (one of the great ideas in statistics)

Goal: prediction (Given independent/explanatory/predictor variables x_1, \dots, x_{p-1} ,
predict dependent/response/outcome variable y)

(p : feature dimension; we use $p-1$ as there's an additional intercept term)

Idea: least-squares: $\operatorname{argmin}_{\beta_0, \dots, \beta_{p-1}} \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_{i,1} + \dots + \beta_{p-1} x_{i,p-1}))^2$

Much more general than it appears by data transformation:

Ex. 1 Suppose $y \approx \sum_{k=0}^d a_k x^k$ has a polynomial relationship with x , then
for each data point x_i we can define

$$z_{i,1} = x_i, \quad z_{i,2} = x_i^2, \quad \dots, \quad z_{i,d} = x_i^d.$$

Then

$$y \approx \sum_{k=0}^d a_k x^k = a_0 + \sum_{k=1}^d a_k z_k$$

Ex. 2 Suppose $y \approx C_0 e^{C_1 x}$, then by defining $w_i = \log y_i$,

$$y_i \approx C_0 e^{C_1 x_i} \Rightarrow w_i \approx \log C_0 + C_1 x_i$$

So (x_i, w_i) has a linear relationship $w_i \approx \beta_0 + \beta_1 x_i$, with

$$\beta_0 = \log C_0, \quad \beta_1 = C_1.$$

Solution:

$$X = \begin{bmatrix} 1 & x_{1,1} & \dots & x_{1,p-1} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n,1} & \dots & x_{n,p-1} \end{bmatrix} \in \mathbb{R}^{n \times p}, \quad \beta = \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_{p-1} \end{bmatrix} \in \mathbb{R}^p$$

Least squares solution:

$$\min_{\beta_0, \dots, \beta_{p-1}} \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_{i,1} + \dots + \beta_{p-1} x_{i,p-1}))^2 \rightarrow \min_{\beta \in \mathbb{R}^p} \|y - X\beta\|^2$$

Theorem. The minimizer $\hat{\beta}$ is

$$\hat{\beta} = (X^T X)^{-1} X^T y,$$

provided that $X^T X \in \mathbb{R}^{p \times p}$ is invertible (requiring $n \geq p$)

Statistical analysis

Statistical modeling: $y = X\beta + e$, $\mathbb{E}[e] = 0$, $\text{Cov}(e) = \sigma^2 I_n$

Then: • $\mathbb{E}[\hat{\beta}] = \beta$

• $\text{Cov}(\hat{\beta}) = \sigma^2 (X^T X)^{-1}$

• $\mathbb{E} \left[\frac{\|y - X\hat{\beta}\|^2}{n-p} \right] = \sigma^2$ (motivating the estimator $\hat{\sigma}^2 = \frac{\|y - X\hat{\beta}\|^2}{n-p}$ for σ^2)

• $\frac{\hat{\beta}_j - \beta_j}{\underbrace{\sqrt{\hat{\sigma}^2 [(X^T X)^{-1}]_{jj}}}_{\text{estimated variance of } \hat{\beta}_j}} \xrightarrow{d} N(0,1) \text{ as } n \rightarrow \infty.$

Pf. $\hat{e} = y - X\hat{\beta} = y - X(\beta + (X^T X)^{-1} X^T e)$
 $= y - X\beta - X(X^T X)^{-1} X^T e$
 $= (I - \underbrace{X(X^T X)^{-1} X^T}_{\text{call it } P})e.$

Note that P is a projection matrix: $P^T = P$, $P^2 = P$, since

$$P^T = [X(X^T X)^{-1} X^T]^T = X(X^T X)^{-1} X^T = P$$

$$P^2 = \underbrace{X(X^T X)^{-1} X^T}_{=I_p} X(X^T X)^{-1} X^T = X(X^T X)^{-1} X^T = P.$$

$$\Rightarrow (I_n - P)^T = I_n - 2P + P^2 = I_n - 2P + P = I_n - P$$

$$\Rightarrow \|\hat{e}\|^2 = \hat{e}^T \hat{e} = e^T (I_n - P)^T (I_n - P) e = e^T (I_n - P)^T e = e^T (I_n - P) e$$

$$\Rightarrow \mathbb{E}[\|\hat{e}\|^2] = \text{Tr}(\text{Cov}(e)(I_n - P)) = \sigma^2 \text{Tr}(I_n - P)$$

$$(\mathbb{E}[\hat{e}^T A \hat{e}] = \text{Tr}(\text{Cov}(\hat{e}) A))$$

Since $\text{Tr}(I_n) = n$ check $\text{Tr}(AB) = \text{Tr}(BA)!$

$$\text{Tr}(P) = \text{Tr}(X(X^T X)^{-1} X^T) = \text{Tr}(X^T X(X^T X)^{-1}) = \text{Tr}(I_p) = p$$

$$\Rightarrow \text{Tr}(I_n - P) = \text{Tr}(I_n) - \text{Tr}(P) = n - p$$

$$\Rightarrow \mathbb{E}[\|\hat{e}\|^2] = (n-p) \sigma^2.$$

