

logistics.

[monwenlong.github.io/teaching/sta355f25/](https://monwenlong.github.io/teaching/sta355f25/)

Quercus. Piazza.

Grading.

$$\max \left\{ \begin{array}{l} 30\% \times \text{midterm} + 40\% \times \text{final.} + 30\% \times \text{homework,} \\ \quad \quad \quad (15\% \text{ each}) \\ 70\% \times \text{final} + 30\% \times \text{homework} \end{array} \right\}$$

Homework marks by completion only.

Ptl.

- |    |                        |      |         |
|----|------------------------|------|---------|
| 1. | 50 + 10 + 50 + 10 + 60 | ends | 8pm     |
| 2. | 75 + 10 + 75           | ends | 7:50pm. |

TBD on piazza

Textbook : "All of Statistics" by L. Wasserman.

# Overview.

- Probability review (quick), Statistical models.
- Methodology
  - Empirical distributions & bootstrap.
  - Parameter inference, M-estimators, MLE, etc.
  - Testing: p-values, Likelihood ratio tests.
  - Bayesian inference.
- Theory
  - Asymptotics
  - Basics of decision theory
  - Uniform convergence (\*).
- Models
  - Linear and generalized linear models.
  - Nonparametric models.
  - Classification (\*), causal inference (\*).

## Probability review.

- Basic concepts.

Sample space  $\Omega$ ,  $w \in \Omega$  sample outcome.

$P$ : mapping from subsets of  $\Omega$   
to non-negative reals.

(sometimes some sets are non-measurable  
and we restrict attention to measurable subsets)

## Axioms of probability.

$$(i) P(A) \geq 0 \quad \forall A$$

$$(ii) P(\Omega) = 1.$$

$$(iii). P\left(\bigcup_{i=1}^{+\infty} A_i\right) = \sum_{i=1}^{+\infty} P(A_i) \quad \text{If the subsets } (A_i)_{i=1}^{+\infty} \text{ are disjoint.}$$

(An event is a subset of  $\Omega$ ).

## Independence of events.

$$A, B \text{ independent} \Leftrightarrow P(A \cap B) = P(A) \cdot P(B).$$

## Conditional probability.

$$\text{If } P(B) > 0, \quad P(A|B) = \frac{P(A \cap B)}{P(B)}.$$

For conditioning on zero-probability events.

$$B = \{X = x\} \quad \text{where } X \text{ is a continuous r.v.}$$

Bayes theorem.

let  $(A_i)_{i=1}^{+\infty}$  be a partition of  $\Omega$ ,

$$P(A_i|B) = \frac{P(B|A_i) \cdot P(A_i)}{P(B)} = \frac{P(B|A_i) \cdot P(A_i)}{\sum_j P(B|A_j) \cdot P(A_j)}$$

Random variables.

Sample space  $\Omega$ .

$$X : \Omega \rightarrow \mathbb{R}.$$

$$(\forall \omega \in \Omega, X(\omega) \in \mathbb{R}).$$

- pmf, cdf, and pdf.

Describing the probability distribution of an r.v.

—  $X$  takes values in  $\{x_1, x_2, x_3, \dots\}$

(Can be infinite, but needs to be countable)

Probability mass function defined on  $\{x_1, x_2, \dots\}$

$$p_X(x_i) = P(X = x_i).$$

(By countable summability,  $\sum_{i=1}^{+\infty} p_X(x_i) = 1$ .)

$$p_X(x_i) \geq 0.$$

This does not work for continuous cases.

$$P(X = x) = 0 \text{ for any } x \in \mathbb{R}.$$

— Alternative approach.

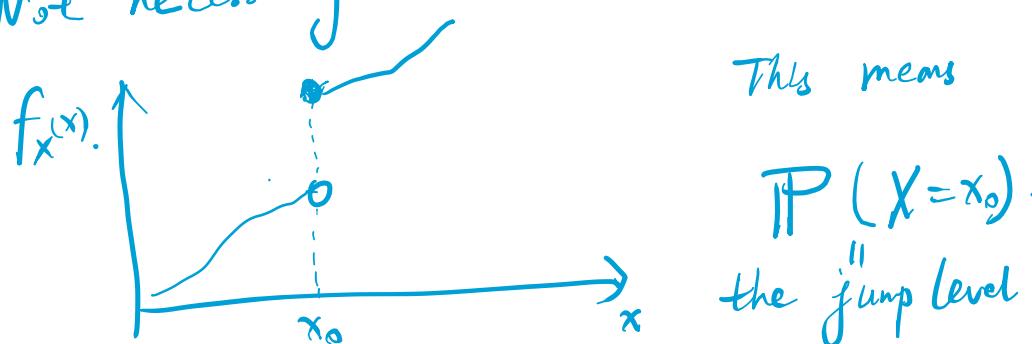
$$\forall x \in \mathbb{R}, \text{ define } f_X(x) = P(X \leq x).$$

"cumulative distribution function".

- Always well-defined.
- Monotonically non-decreasing.

- $\lim_{x \rightarrow +\infty} F_X(x) = 1, \quad \lim_{x \rightarrow -\infty} F_X(x) = 0.$

- Not necessarily continuous.



- Probability density function (pdf).

Suppose  $f_X$  satisfies

$$\int_a^b P_X(x) dx = f_X(b) - f_X(a)$$

||(

$$P(X \in (a, b))$$

then we call  $P_X$  the pdf of  $X$ .

- May not exist for certain r.v.

e.g.  $X \sim \text{Unif}(0, 1)$  w.p.  $\frac{1}{2}$

$$= \frac{1}{2}$$
 w.p.  $\frac{1}{2}$

- "Almost" requires  $f_X$  to be differentiable, but it still exists under weaker conditions.

- pdf can be extended to multivariate cases.  
If there exists a function  $P_X$ . s.t.  

$$P(X \in A) = \int_A P_X(x) dx$$
  
 d-dimensional. for any (measurable) subset  $A \subseteq \mathbb{R}^d$   
 Then we call  $P_X$  to be the pdf of  $X$ .

For this class: focus on discrete and continuous r.v.s  
 (probably also their mixtures).

- Conditional distribution

$(X, Y)$  jointly distributed

| discrete  
continuous.

Discrete 
$$P(Y=y_i | X=x_i) = \frac{P(X=x_i, Y=y_i)}{P(X=x_i)}$$
  
 (Assuming  $P(X=x_i) > 0$ ).      //

$$\frac{P(X=x_i, Y=y_i)}{\sum_{j=1}^{\infty} P(X=x_i, Y=y_j)}$$

$P(\cdot | X=x_i)$  defines  
 a probability distribution  
 on the discrete set  $\{y_1, y_2, \dots\}$ .

Cts Case:

$$\underbrace{P_{Y|X}(y|x)}_{\text{Conditional density of } Y \text{ at } y, \text{ given that } X=x} = \frac{P_{XY}(x,y)}{P_X(x)} = \frac{P_{X,Y}(x,y)}{\int P_{X,Y}(x,y') dy'}.$$

"Conditional density of  $Y$  at  $y$ , given that  $X=x$ ".

e.g.  $\begin{bmatrix} X_1 \\ X_2 \end{bmatrix} \sim \mathcal{N}\left(\begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix}\right)$   
 $\rho \in (-1, 1).$

$$P_{X_1, X_2}(x_1, x_2) = \frac{1}{2\pi\sqrt{\det(\Sigma)}} \exp\left(-\frac{1}{2} \begin{bmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{bmatrix}^\top \Sigma^{-1} \begin{bmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{bmatrix}\right)$$

where  $\Sigma = \begin{bmatrix} \sigma_1^2 & \rho\sigma_1\sigma_2 \\ \rho\sigma_1\sigma_2 & \sigma_2^2 \end{bmatrix}$  is the covariance matrix.

Conditional distribution.

$$X_2 | X_1 = x_1 \sim \mathcal{N}\left(\mu_2 + \frac{\rho\sigma_2}{\sigma_1}(x_1 - \mu_1), \sigma_2^2 - \rho^2\sigma_2^2\right).$$

Expectation.

$$X \left\{ \begin{array}{l} \text{Discrete r.v.} \\ \text{Continuous r.v.} \end{array} \right. \quad \mathbb{E}[X] := \sum_{i=1}^{+\infty} x_i P_X(x_i)$$

(Assuming that  $\sum_{i=1}^{+\infty} |x_i| \cdot P_X(x_i) < +\infty$ )

$$\mathbb{E}[X] := \int_{\mathbb{R}} x p_X(x) dx$$

(Assuming that  $\int |x| \cdot p_X(x) dx < +\infty$ ).

- Expectations are linear.

$$E[aX+bY] = aE[X] + bE[Y].$$

$(a, b \in \mathbb{R} \text{ are deterministic quantities})$ .

does not require (in)dependence between  $X, Y$ .

- There exists r.v. s.t. expectations do not exist.

e.g. "Cauchy distribution".

$$p(x) = \frac{1}{\pi(1+x^2)} \quad (x \in \mathbb{R})$$

$$\begin{aligned} \int_0^{+\infty} |x| p(x) dx &= \frac{1}{\pi} \int_0^{+\infty} \frac{x dx}{1+x^2} \\ &= \frac{1}{2\pi} \left[ \log(1+x^2) \right]_0^{+\infty} = +\infty. \end{aligned}$$

## Inequalities.

From expectation to tail behaviour.

- Markov ineq.  $X \geq 0, \text{r.v. } E[X] < +\infty$

then  $P(X \geq t) \leq \frac{E[X]}{t} \quad (\forall t > 0)$ .

Tail probability decaying at  $\gamma t$  rate (or faster).

Proof:

$$\begin{aligned} P(X \geq t) &= E[1_{X \geq t}] \leq E\left[\frac{X}{t} \cdot 1_{X \geq t}\right] \\ &= \frac{1}{t} E[X \cdot 1_{X \geq t}] \leq \frac{1}{t} E[X]. \end{aligned}$$

Chebyshev ineq. Suppose  $E[X^2] < \infty$

(stronger requirement than Markov ineq.).

$$P(|X - E(X)| \geq t) \leq \frac{\text{var}(X)}{t^2}$$

( $\gamma t^2$ , faster decaying tail).

Proof: Apply Markov ineq to  $Y = |X - E(X)|^2$ .

$$P(|X - E(X)| \geq t) = P(Y \geq t^2) \leq \frac{E(Y)}{t^2} = \frac{\text{var}(X)}{t^2}.$$

In general, if  $E[|X - E(X)|^p] < \infty$ .

then  $P(|X - E(X)| \geq t) \leq \frac{E[|X - E(X)|^p]}{t^p}$ .

Proof: same as Chebyshev.

Rmk: tighter tail bound (eg. exponentially decaying)

is possible, assuming  $\exists$  moment generating functions.

Convergence of random variables.

Notions of convergence.

- almost sure convergence

$$X_1, X_2, \dots, X_n, \dots$$

we say  $X_n \xrightarrow{\text{a.s.}} X$  if  $P(\lim_{n \rightarrow \infty} X_n = X) = 1$ .

- convergence in probability.

we say  $X_n \xrightarrow{P} X$  when

$$\forall \varepsilon > 0, \quad \lim_{n \rightarrow +\infty} P(|X_n - X| > \varepsilon) = 0.$$

- $L^p$  convergence.

we say  $X_n \xrightarrow{L^p} X$  when

$$\lim_{n \rightarrow \infty} \mathbb{E}[|X_n - X|^p] = 0.$$