

Week 1 -

20/04/2020

- Youtube 31-08-2016

- conclude WEEK 1 v1 - review, begin (W1)(L2) - inequalities

$$\textcircled{1} A_2 = \{(x, y) \mid g(x, y) \leq z\} \quad (z?)$$

- Notation:-

$$X \perp Y$$

② Independence can be physical fact OR an assumption which we have to eval.

- How to check for independence?

- Joint density factors multiplicatively $p_{X,Y}(x, y) = p_X(x) p_Y(y)$

- easy to check independence against a density/distribution than in terms of events

- same for multiple iid X_i s.

- identical distri - n observations drawn from same population / same marginal distribution

- iid - independent, identically distributed; random sample from population / common distribution

- An assumption

③ iid notation:

$$X_1, \dots, X_n \sim P \quad \text{or} \quad X_1, \dots, X_n \sim F \quad \text{or} \quad X_1, \dots, X_n \sim p$$

(iid drawn from common probab. distribution P) (iid drawn from same CDF) (iid from (prob) same density function)

④ By default in the class - assume iid

• non iid: departing from this is very important e.g. time series

⑤: distributions included for reference

- Normal distribution is a family of distri indexed by 2 parameters μ, σ

$$p(x; \mu, \sigma) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad \mathbb{E}[X] = \mu$$

μ indicates parameter σ^2 indicates r.v. $\text{var}(X) = \sigma^2$

MVN: Random vectors

- vector of random variables $\underline{X} = \begin{pmatrix} X_{(1)} \\ \vdots \\ X_{(d)} \end{pmatrix} \in \mathbb{R}^d$

- Typo ⑤: $p(\underline{x}) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(\underline{x}-\underline{\mu})^T \Sigma^{-1}(\underline{x}-\underline{\mu})\right)$

- $\underline{\mu} = \begin{pmatrix} \mu_{(1)} \\ \vdots \\ \mu_{(d)} \end{pmatrix}$ $\underline{\Sigma}$ = variance covariance matrix = $\begin{bmatrix} & & \\ & & \\ & & \end{bmatrix}$

- very important facts: - ^{MVN (review)}
col prop.

• marginal distri of any set of coordinates is Normal

• condit. distri of any —||— given set of coordinates is also Normal

⑥ χ^2 distri
• $z_1, \dots, z_p \sim N(0, 1)$ - χ^2 distri is a family of distri indexed by parameter p

$$X = \sum_{j=1}^p z_j^2$$

$$\text{then } X \sim \chi_p^2$$

$$\mathbb{E}[X_p^2] = p$$

$$\text{Var}[X_p^2] = 2p$$

• Above a continuous distri.

• onto discrete

Bernoulli $X \sim \text{Bern}(\theta)$ if $P(X=1) = \theta$ $P(X=0) = 1-\theta$

$$p(x; \theta) = \theta^x (1-\theta)^{1-x} \quad x=0,1$$

Binomial $X \sim \text{Binomial}(\theta)$ (sum of n Bernoulli's)

$$p(x; \theta) = P(X=x) = \binom{n}{x} \theta^x (1-\theta)^{n-x} \quad x \in \{0, \dots, n\}$$

• Remind yourself of these!

Sampling distributions

• $X_1, \dots, X_n \sim P$

• sample mean:

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i = g(X_1, \dots, X_n)$$

⑤⑦ "Take a sample of data, take function of sample, which itself defines new r.v." (summarise data)
- A statistic can be viewed as a function, data transform of data

- As X_1, \dots, X_n are random; \bar{X}_n is also random

• \bar{X}_n , as an r.v., also has a distri, the sampling distri

$$G_n(t) = P(\bar{X}_n \leq t)$$

• In general, sampling distri is difficult to find, in particular as P is unknown

However sampling distri has properties:-

- See practice problem for completion the

23:33

⑥: Difference between distribution of X_i s and the \bar{X}_n s

sample variance: $S_n^2 = \frac{1}{n-1} \sum (X_i - \bar{X})^2 \Rightarrow E(S^2) = \sigma^2$

• why divide by $(n-1)$ instead of $n \rightarrow$ unbiased

• only assuming data iid, mean and variance exist

$$E(S^2) = \sigma^2$$

• Theorem 10 builds on this by specifying that $X_1, \dots, X_n \sim N(\mu, \sigma^2)$

⑦ we can say $\bar{X}_n \sim N(\mu, \sigma^2/n)$

• proved by mgf of \bar{X}_n ; mgfs completely characterize a distribution

• that is, compute mgf $\Rightarrow \bar{X}_n \sim N(\mu, \sigma^2/n)$

⑧: $\frac{(n-1)S_n^2}{\sigma^2} \sim \chi_{n-1}^2$

• contents of Theorem 10 - only if $X_1, \dots, X_n \sim N(\mu, \sigma^2)$

⑨: 28:20 \rightarrow query about intersections; and clarify on double integral

Week 1 - (12)

Youtube
31-08-2016

20/04/2022

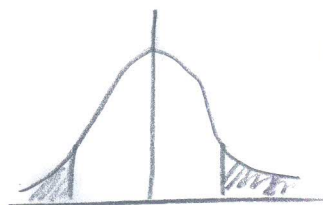
Probability Inequalities

- finding bounds
- Probability bounds useful in ML in particular Hoeffding
- stepping stone to VC Theory (advanced prob. bounds) \rightarrow convergence theory (what happens to \bar{X}_n with more data)
- multitude of bounds; but main:-

- 1) Markov
- 2) Hoeffding
- 3) Gaussian

Theorem 1

$$X \sim N(0, 1)$$



31:50?

(77)

- Tails
- How fast do the tails
- LW uses t (arbitrary positive real t)

- As $t \uparrow$, $P(|X_n| > t)$ gets small (very important)

- 2nd part: intuitively:

$P(|\bar{X}_n| > t)$ goes to 0 quickly in terms of t and sample size n

(W) inequalities tied to thick Gaussian tail inequalities for other distributions

Theorem 1 (Proof)

- Do one side, use symmetry \rightarrow factor 2

(W): $\phi(s)$ is stand normal density

- multiply/divide by s
- check 36:11 little step

Property:

(67)

$$\begin{aligned}\phi'(s) &= -s\phi(s) \\ \Rightarrow s\phi(s) &= -\phi'(s)\end{aligned}$$

(A1) check proof that you understand steps ✓

Theorem 2 - Markov's Inequality

- very weak, but used to prove Hoeffding
- Assume $X > 0$ and that mean exists and is well defined i.e. finite mean $E(X)$
- (W): Cauchy distribution - tails spread out very slowly
- does not have well-defined mean due to fat-tails
- As $X > 0$, integrate in range $[0, \infty)$

A3 check proof

- As $t \uparrow$; $P(X > t) \downarrow$ by $\frac{1}{t} \rightarrow$ qualifying weakness

Theorem 3 - Chebyshev's inequality

- Assuming that variance is well defined and exists (i.e. finite)
- this extra information can be used
- Chebyshev's inequality builds/argues Markov inequality using this info to get a sharper inequality

Hoeffding Inequality

- sharper; assume that X is bounded above and below
- Boundedness \Rightarrow thin tails
- Probability mass is 0 outside interval $[a, b]$; tails cannot be thinner
- improvement over Markov, Chebyshev $P(|\bar{X}_n - \mu| \geq \epsilon) \leq 2e^{-2n\epsilon^2/(b-a)^2}$
- Proof is useful (useful tricks)
- 2 tricks:
 1. All moments of X exists (moment generating function exist?)
 2. How to bound mgf using bound on X

Chernoff's method (going back to epsilon)

$$P(X > \epsilon) \leq \inf ("min") e^{-t\epsilon} E[e^{tx}] \quad \text{--- mgf}$$

$t \geq 0$ --- variational param.

- both are functions of t
- minimise over t

- check 50:06 for chernoff

④ - check you understand chernoff before next lecture