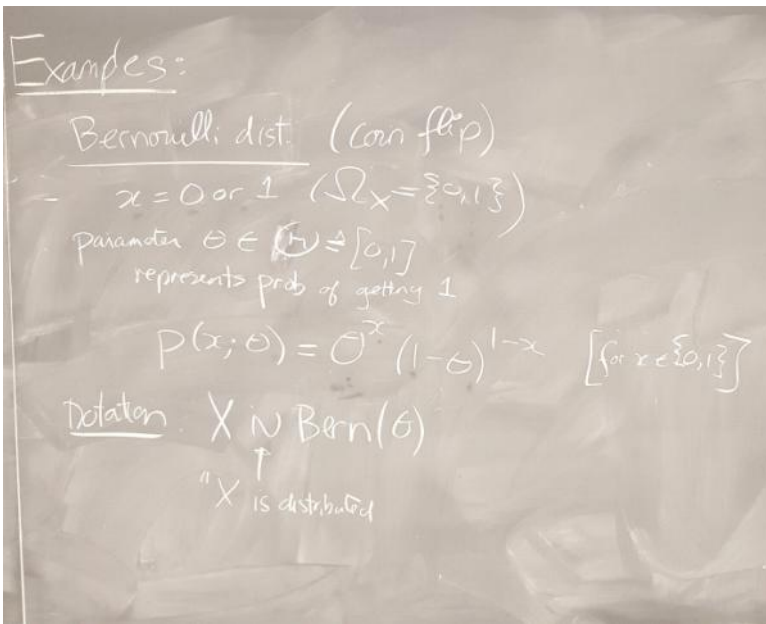
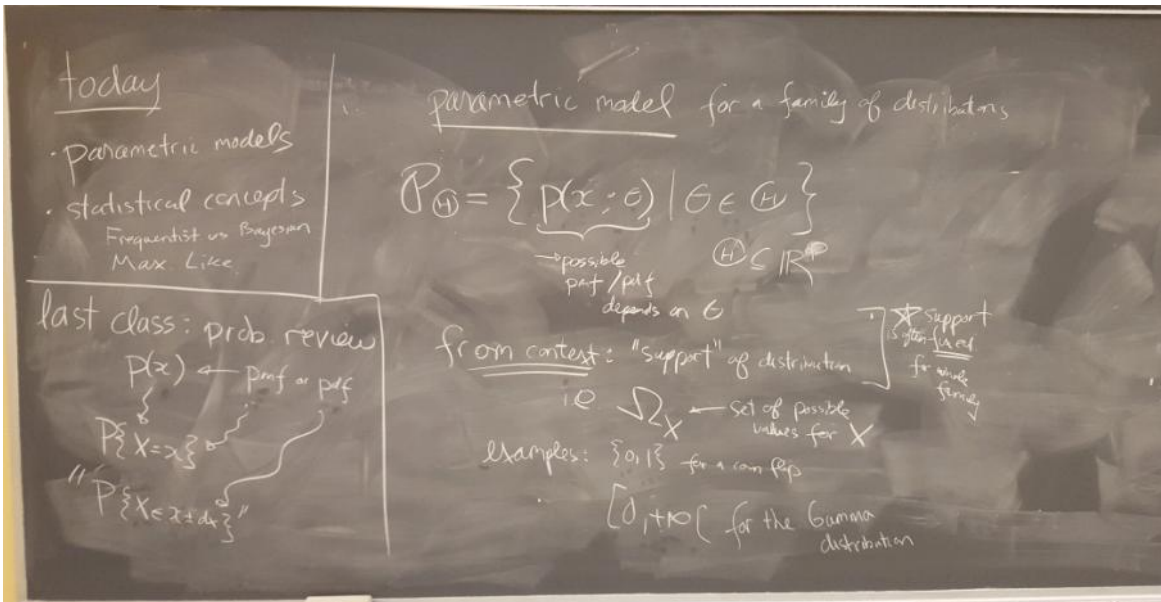


Lecture 3 - scribbles

Friday, September 9, 2016
13:35



notation: $X \sim \text{Bern}(\theta)$

↑
"X is distributed as ..."

$$p(x; \theta) = \text{Bern}(x; \theta) \quad \text{Bern}(n; \theta) \dots$$

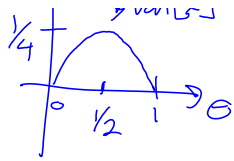
↑ parameter
variable for density

$$E[X] = \theta$$

$$\text{Var}[X] = \theta(1-\theta)$$

$$E[X] = \theta$$

$$\text{Var}[X] = \theta(1-\theta)$$



• binomial distribution : \rightarrow n independent coin flips
 \rightarrow sum of n indep. Bern(θ) R.V.
 \uparrow often denoted p

let $X_i \overset{\text{iid.}}{\sim} \text{Bern}(\theta)$ \rightarrow independent & identically distributed parameter

let $X = \sum_{i=1}^n X_i$ we have $X \sim \text{Bin}(n, \theta)$
context

$$\Omega_X = \{0, 1, \dots, n\}$$

$$P(x; \theta) = \binom{n}{x} \theta^x (1-\theta)^{n-x} \text{ for } x \in \Omega_X$$

$\underbrace{\binom{n}{x}}_{\substack{\# \text{ of ways} \\ \text{to get } x \text{ successes} \\ \text{out of } n \text{ throws}}} \underbrace{\theta^x (1-\theta)^{n-x}}_{\substack{= \prod_{i=1}^n \text{Bern}(x_i; \theta) \\ \text{indep coin flips}}}$

mean $X = \sum_{i=1}^n X_i$ θ

$$E[X] = \sum_i E[X_i] = n\theta$$

similarly $\text{Var}[X] = n\theta(1-\theta)$

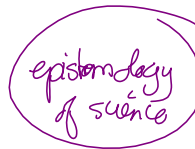
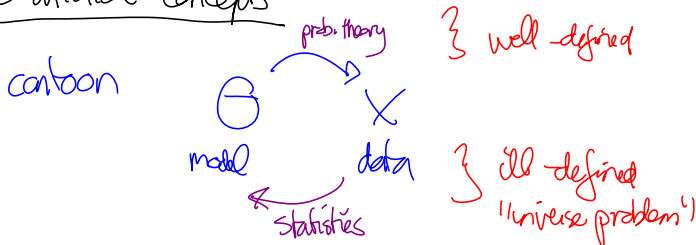
other distributions: Poisson (λ) $\Omega_X = \{0, 1, \dots\}$ [count data] \uparrow mean

Gaussian (1D) $N(\mu, \sigma^2)$ $\Omega_X = \mathbb{R}$ \uparrow variance

• gamma $\Gamma(\alpha, \beta)$ $\Omega_X = \mathbb{R}_+$ \uparrow shape scale

others: Laplace, Cauchy, exponential, beta & Dirichlet dist. . .

Statistical concepts



example: model of n indep coin flips

example: model of n indep coin flips

prob theory: proba of k heads in a row

statistics: I have observed n_b tails what is θ ?
 n_h heads

Frequentist vs. Bayesian

philos: meaning of a probability??

a) (traditional) frequentist semantic

$p(X=x)$ represents the relative frequency
of observing $X=x$
if could repeat ∞ # of i.i.d. experiments

b) Bayesian (subjective) semantic

$p(X=x)$ encodes our "belief" that $X=x$

laws of probability characterizes a "rational"
way to combine "beliefs" and "evidence"

[has motivation in terms gambling,
utility/decision theory, etc...]

operationally:

• for a discrete R.V. suppose that $P(X=x) = \theta$

$$\Rightarrow P(X \neq x) = 1 - \theta$$

$$B = \mathbb{1}_{\{X=x\}} \sim \text{Bern}(\theta) \text{ p.v.}$$

$$\text{by L.L.N. } \frac{1}{n} \sum_{i=1}^n B_i \xrightarrow{\text{a.s.}} \mathbb{E}[B_i] = \theta$$

↳ limiting frequency interpretation

$$[\text{notes: } \sum_{i=1}^n B_i \sim \text{Bin}(n, \theta) \quad \frac{1}{n} \mathbb{E} \left[\sum_{i=1}^n B_i \right] = \frac{n\theta}{n} = \theta$$

$$\text{Var} \left[\frac{1}{n} \sum_{i=1}^n B_i \right] = \frac{1}{n^2} \text{Var} [\text{Bin}(n, \theta)] = \frac{1}{n^2} n \theta (1-\theta)$$

$$\text{Var}\left[\frac{1}{n} \sum_{i=1}^n B_i\right] = \frac{1}{n^2} \text{Var}[\text{Bin}(n, \theta)] = \frac{1}{n^2} n \theta(1-\theta)$$

$$\text{(use Var}(aX) = a^2 \text{Var}(X)) \quad = \frac{1}{n} \theta(1-\theta)$$

$$\sqrt{n} \left(\frac{1}{n} \text{Bin}(n, \theta) - \theta \right) \xrightarrow{d} N(0, \theta(1-\theta))$$

CLT central limit theorem

all works?

Bayesian approach

⊗ very simple philosophically: treat all uncertain quantities as R.V.s

→ encode all your knowledge about the system ("beliefs") as a "prior" or probabilistic models

(later we'll see more with model selection)

simplest example: biased coin we do not know θ

we believe $X \sim \text{Bin}(n, \theta)$

→ θ is a R.V. for the Bayesian ⇒ need a $p(\theta)$

prior distribution
 $\mathcal{I}_\theta = [0, 1]$

suppose observe $X=x$ (in n flips) result of flips

then we can "update" our belief about θ using Bayes rule

$$p(\theta | X=x) = \frac{\underbrace{p(x|\theta)}_{\text{observation model / likelihood}} \underbrace{p(\theta)}_{\text{prior}}}{\underbrace{p(x)}_{\text{normalization}}}$$

posterior belief

[note $p(x|\theta) \rightarrow$ is pdf $p(x, \theta)$ is a mixed distribution]
 $p(\theta) \rightarrow$ is pdf

example: suppose $p(\theta)$ is uniform on $[0, 1]$
 "no specific preference"

$$p(\theta | x) \propto \theta^x (1-\theta)^{n-x}$$

$\{0, 1, \dots, n\}$
 \downarrow
 $[x \in 0:n]$

scaling: $\int_0^1 \theta^x (1-\theta)^{n-x} d\theta = B(x+1, n-x+1)$ (proportional)

$$B(a,b) \triangleq \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)} \quad \Gamma(a) \triangleq \int_0^\infty u^{a-1} e^{-u} du$$

→ this is called a Beta function "Beta distribution"

$$\text{Beta}(\theta; \alpha, \beta) = \frac{\theta^{\alpha-1} (1-\theta)^{\beta-1}}{B(\alpha, \beta)} \mathbb{1}_{[0,1]}(\theta)$$

parameters

so $p(\theta|x) = \text{Beta}(\theta; x+1, n-x+1)$

look at $E[\theta|x]$ (posterior mean)

as an estimator for θ

notation: $\hat{\theta}(x)$ estimator function for θ observation / "training data"

$$\hat{\theta}_{\text{Bayes}}(x) = E[\theta|x]$$

mean for a $\text{Beta}(\alpha, \beta) = \frac{\alpha}{\alpha+\beta}$ so here $E[\theta|x] = \frac{x+1}{n+2}$ (this is result of experiment)

$$\hat{\theta}(x) = E[\theta|X] = \frac{X+1}{n+2}$$

$$\theta \sim \text{Bin}(n, \theta) + 1$$

by LLN, if repeat m times independent experiment (with θ fixed)

we have $\frac{1}{m} \sum_{i=1}^m \hat{\theta}(X_i) \xrightarrow{\text{a.s.}} E[\hat{\theta}] = \frac{n\theta+1}{n+2}$

$$= \left(\frac{n}{n+2}\right)\theta + \frac{2}{n+2} \left(\frac{1}{2}\right)$$

↑ true parameter ↑ prior mean

convex combination

note that for finite n , $E[\hat{\theta}(x)] \neq \theta$

→ "biased method" ← from the prior belief

but as $n \rightarrow \infty$, $E[\hat{\theta}_n(x)] \rightarrow \theta$

"asymptotically unbiased"

$\text{Beta}(\alpha, \beta)$ s.t. posterior mean of $\text{Beta}(\alpha+x, \beta+n-x)$
has mean of θ with respect to R.V. X

exercise?

summary:

frequentist considers various "estimation" rules

[e.g. max. likelihood is most common]
but there are others

- method of moments
- M-estimation
- etc.

then they analyze "frequentist" properties of these estimators

- biased?
- variance?
- consistent?
- etc.

to get guarantees and compare them together

→ but no best rule [vs. Bayes rule for Bayesians]

sometimes MLE is "inadmissible" [e.g. James-Stein estimator]