

Lecture 4 - scribbles

Friday, September 11, 2020 13:09

- today:
- finish prob.
 - frequentist vs. Bayesian

binomial distribution:

model n indep. coin flips

sum of n indep. $\text{Bern}(\theta)$ R.V.

let $X_i \overset{\text{iid}}{\sim} \text{Bern}(\theta)$ \rightarrow (mutually) independent and identically distributed

implicitly talking $\rightarrow X_1, \dots, X_n \overset{\text{iid}}{\sim}$

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

"binomial with parameter n & θ "

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

pdf: $p(x; \theta) = \binom{n}{x} \theta^x (1-\theta)^{n-x}$ for $x \in \Omega_X$

of ways to choose x elements out of n
"n choose x"

$$\binom{n}{x} = \frac{n!}{x!(n-x)!}$$

$$\theta^{\sum x_i} (1-\theta)^{\sum (1-x_i)} = \theta^{\sum x_i} (1-\theta)^{n - \sum x_i}$$

$$= \prod_{i=1}^n \text{Bern}(x_i; \theta) = p(x_1, \dots, x_n)$$

mean: $X = \sum_i X_i$

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

indep.

similarly, $\text{Var}(\sum_i X_i) = \sum_i \text{Var}(X_i) = n\theta(1-\theta)$

other distributions: Poisson(1) $\Omega_X = \{0, 1, 2, \dots\} = \mathbb{N}$ [count data]

mean

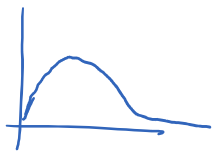
Gaussian in 1D $N(\mu, \sigma^2)$ $\Omega_X = \mathbb{R}$

mean variance

gamma $\Gamma(\alpha, \beta)$ $\Omega_X = \mathbb{R}^+$

shape α inverse scale aka rate β

mean $\frac{\alpha}{\beta}$ variance $\frac{\alpha}{\beta^2}$

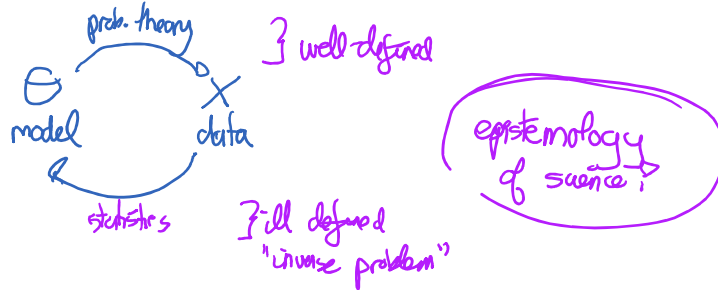


mean $\frac{\alpha}{\beta}$ variance $\frac{\alpha}{\beta^2}$

other: Laplace, Cauchy, exponential, beta, Dirichlet
 $\hookrightarrow \Gamma(1, \beta)$ \hookrightarrow Dirichlet on 2 elements
 $\hookrightarrow \Omega_X = [0, 1]$

Statistical concepts

cartoon



example: model n indep coin flips

prob. theory \rightarrow prob. k heads in a row

statistics: I have observed k heads, what is θ ?
 $n-k$ tails

frequentist vs. Bayesian

Semantic of prob.: meaning of a prob.?

a) (traditional) frequentist semantic

$P\{X=x\}$ represents the limiting frequency of observing $X=x$
if I could repeat ∞ # of icid experiments

b) Bayesian (subjective) semantic

$P\{X=x\}$ encodes an agent "belief" that $X=x$

laws of prob. characterizes a "rational" way to combine "beliefs"
 and "evidence" [observations]

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

operationally:

Bayesian approach: (*) very simple philosophically

• treat all uncertain quantities as R.V.

i.e. encode all knowledge about the system ("beliefs") as a "prior" or probabilistic model and then use law of proba (and Bayes rule) to get updated beliefs and answers. >

justification for frequentist semantics:

• for discrete R.V. X , suppose $P\{X=x\} = \theta$
 $\Rightarrow P\{X \neq x\} = 1 - \theta$

$B \triangleq \mathbb{1}_{\{X=x\}} \Rightarrow B \sim \text{Bern}(\theta)$ R.V.

indicator fct. $\mathbb{1}_A(u) = \begin{cases} 1 & \text{if } u \in A \\ 0 & \text{o.w.} \end{cases}$

repeat i.i.d. experiments i.e. $B_i \stackrel{\text{i.i.d.}}{\sim} \text{Bern}(\theta)$

by L.L.N. (law of large numbers) $\frac{1}{n} \sum_{i=1}^n B_i \xrightarrow{\text{a.s.}} E[B_i] = \theta$

by CLT $\sqrt{n} \left(\frac{1}{n} \sum_{i=1}^n B_i - \theta \right) \xrightarrow{d} N(0, \frac{\theta(1-\theta)}{n})$ limiting frequency

coin flips - Bayesian approach

biased coin flips

unknown \Rightarrow model it as a R.V.

we believe $X \sim \text{Bin}(n, \theta)$ \Rightarrow need a $p(\theta)$ "prior distribution"

$$\Omega_{\theta} = [0, 1]$$

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

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

$$p(\theta = \sigma | X=x) = \underbrace{p(X=x | \theta)}_{\text{observation model}} \underbrace{p(\theta = \sigma)}_{\text{prior belief}} \underbrace{\Bigg\}}_{\text{normalization}} \text{ "marginal likelihood" }$$

likelihood

[note: $p(z|\theta) \rightarrow$ pmf $p(x, \theta)$ is a "mixed distribution"
 $p(\theta) \rightarrow$ pdf

Example:

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

$p(\theta|x) \propto \underbrace{\theta^x (1-\theta)^{n-x}}_{p(z|\theta) \text{ up to a constant}} \underbrace{\mathbb{1}_{[0,1]}(\theta)}_{p(\theta)}$
 "proportional to"

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

normalization constant $\int_{\theta} p(\theta|x) d\theta = 1$
 $B(a,b) \triangleq \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$ (Beta fct.)
 $\Gamma(a) = \int_0^{\infty} u^{a-1} e^{-u} du$ (Gamma fct.)

here $p(\theta|x)$ is called a "beta distribution"

$B(\theta|\alpha, \beta) \triangleq \frac{\theta^{\alpha-1} (1-\theta)^{\beta-1} \mathbb{1}_{[0,1]}(\theta)}{B(\alpha, \beta)}$
 parameters

• uniform distribution: $B(\theta|1,1)$

• posterior $B(\theta|x+1, n-x+1)$

exercise to the reader: if use $B(\alpha_0, \beta_0)$ as prior

posterior will be $B(x+\alpha_0, n-x+\beta_0)$

(*) posterior $p(\theta|X=x)$ contains all the info from data x that we need to answer new queries about θ
 observation

e.g. Question: what is prob. of head ($F=1$) on the next flip

as a frequentist $P(F=1 | \text{data}) = \hat{\theta}$ (estimate)

as a Bayesian $p(F=1 | X=x) = \int_{\theta} p(F=1, \theta | X=x) d\theta$ (product rule)

$$\overset{\theta}{=} \int \underbrace{p(F=1|\theta, x=x)}_{=\theta \text{ (by our model)}} \underbrace{p(\theta|x=x)}_{\text{posterior}} d\theta$$

$$= \int \theta p(\theta|x=x) d\theta = \mathbb{E}[\theta|x=x] \text{ "posterior mean of } \theta \text{"}$$

* a meaningful "Bayesian" estimator of θ

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

relation: $\hat{\theta} : \text{observation} \rightarrow \theta$

our coin example: $p(\theta|x) = \text{Beta}(\theta | \alpha=x+1, \beta=n-x+1)$

mean of a beta $\propto \frac{\alpha}{\alpha+\beta}$

$$\text{thus } \hat{\theta}_{\text{Bayes}}(x) = \mathbb{E}[\theta|x] = \frac{x+1}{n+2}$$

here, biased estimator $\mathbb{E}_x[\hat{\theta}(x)] \neq \theta$

$$= \mathbb{E}\left[\frac{x+1}{n+2}\right] = \frac{\mathbb{E}[x]+1}{n+2} = \frac{n\theta+1}{n+2}$$

but asymptotically unbiased $\xrightarrow{n \rightarrow \infty} \theta$

compare { contrast with $\hat{\theta}_{\text{MLE}}(x) = \frac{x}{n}$ [unbiased $\mathbb{E}\left[\frac{x}{n}\right] = \frac{\mathbb{E}[x]}{n} = \theta$]

to summarize:

- as a Bayesian: get a posterior + use laws of probability
- in "frequentist statistics"

consider multiple estimators

- MLE
- moment matching
- Bayesian posterior mean
- MAP
- regularized MLE

and then analyze their statistical properties

- biased: ?
- variance: ?
- consistent: ?

Maximum likelihood principle

Setup: given a parametric family $p(x; \theta)$ for $\theta \in \Theta$

we want to estimate/learn θ from x

