

Bayesian statistical data science: ^{IP} ① build a probabilistic (BDS) process model for data-generating

process ② use Bayes' Thm. & its corollaries as optimal information-processing algorithm

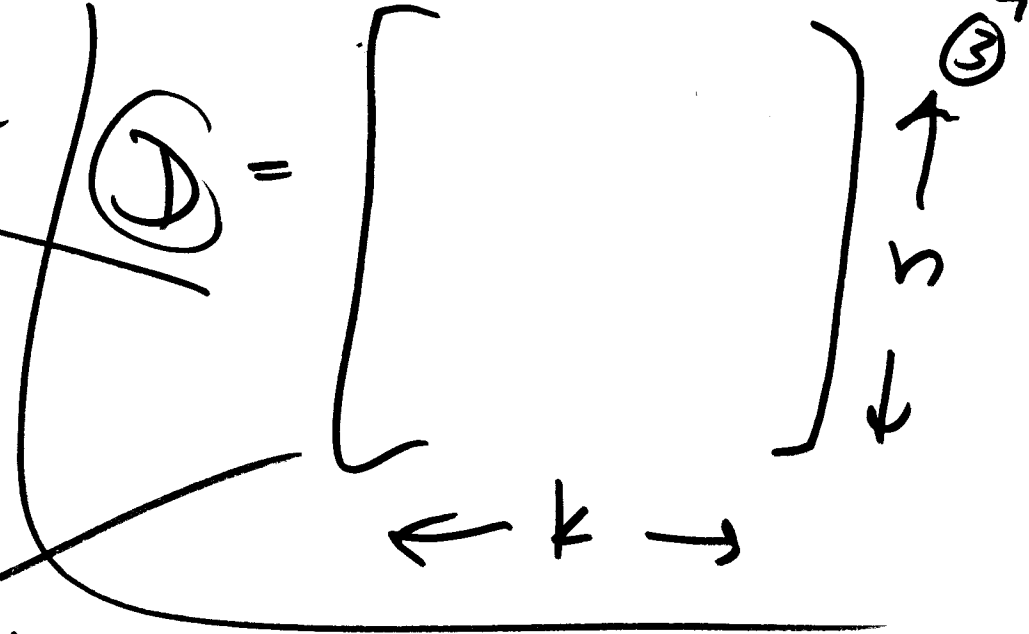
Machine Learning
data science

① Start with $P = (Q, C)$

specify an algorithm that performs the task that answers Q well
(eg., prediction)

Disadvantages of BDS: ① model uncertainty ② Bayesian inference & prediction may well not scale well

to really large
(huge k ,
huge n , or
both)



ML good at scaling their
algorithms to Big Data sizes but
ML usually doesn't try to provide
well-calibrated uncertainty &
assessments (this is a disadvantage
of ML) (to be continued)

Bayesian	vs.	statistical	frequentist
Statistical	vs.	ML	machine learning data science

$\underline{Z} = (Z_1, Z_2, \dots, Z_n, \dots)$ process of getting data

binary observed data vector for $t = \text{me}$ per FE rate

$\underline{Y} = (Y_1, Y_2, \dots, Y_n, \dots)$ before data arrives

de Fine H: start with $\{0, 1\}$

$P(Z_1 = Y_1, Z_2 = Y_2, \dots, Z_n = Y_n | B)$

prior predictive dist. for \underline{Z} (data)

Quiz 2
Case study
 $n = 921$

inference: no feedback loop
prediction: feedback ✓

note: 2^n possible data vectors; $2^{921} = 10^{277}$

if we know nothing about subjects

in the sample, our marginal predictive uncertainty about each person's response is the same! (the Σ_i are marginally ID)

moreover, any logically-internally-consistent (coherent) predictive

dist. must be invariant to the order in which data vector is given to us

exchangeability

123 { 123 312
132 321
213
231

all possible order permutations

invariant = remain the same under a transformation

→ group theory & other abstract algebra topics

exchangeability is a property of 6

your information about the world,

↓
epistemology

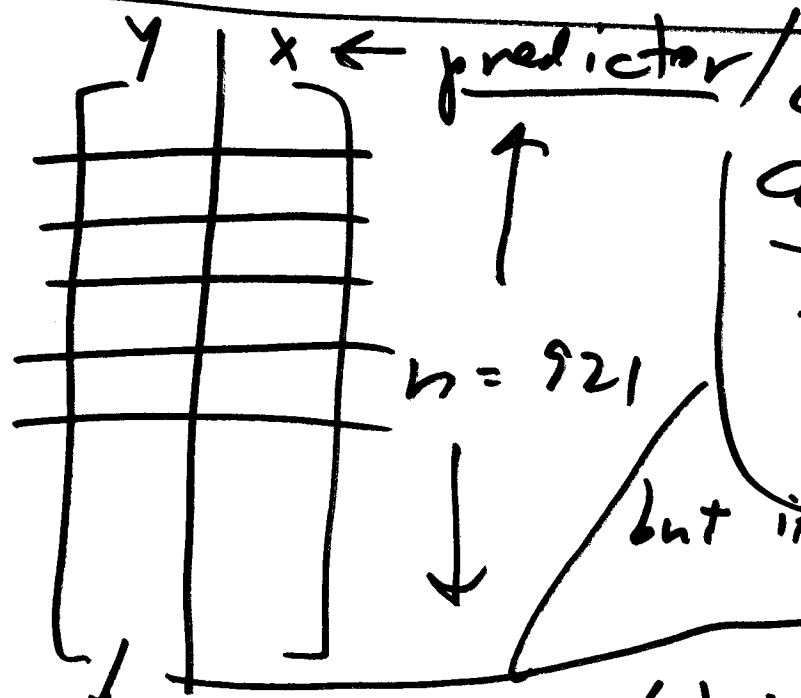
not about the world itself

↓
ontology

The judgment of exchangeability

comes ~~from~~ from $\mathcal{B} \leftarrow \mathcal{C}$

↓ outcome



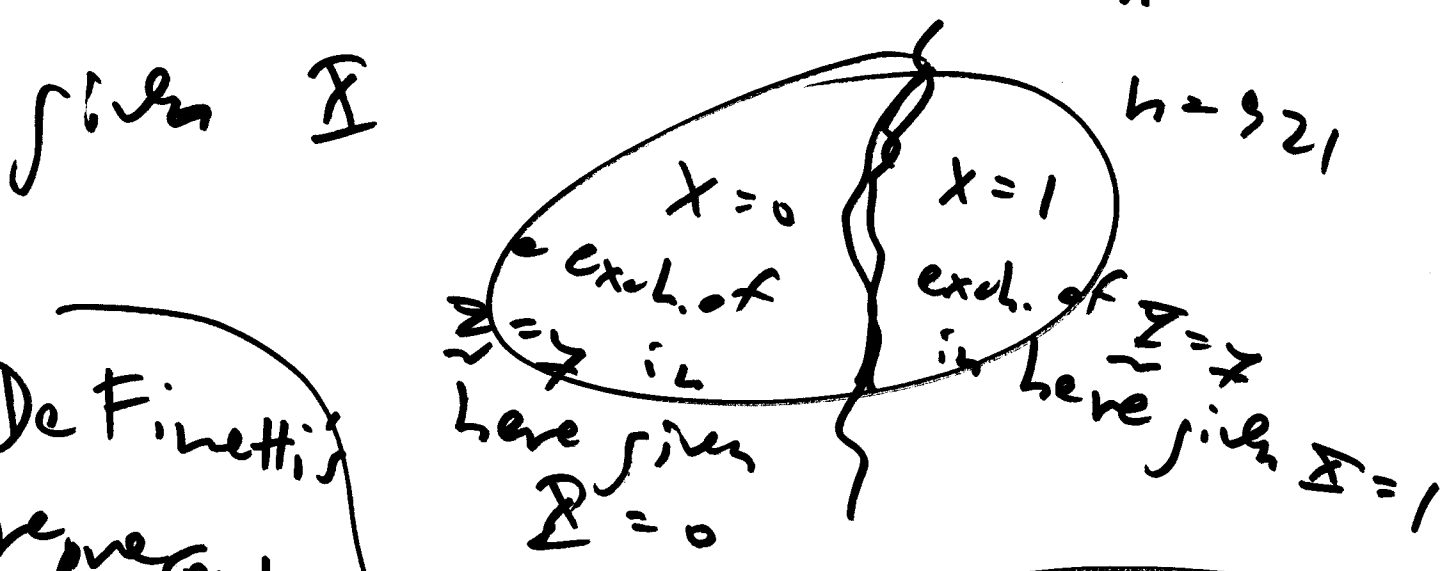
~~Quiz 2~~
Quiz 2
Case study

~~Covariate~~
Covariate / independent variable
our uncertainty about \mathcal{X} is exch. if \mathcal{X} is not known

but if \mathcal{X} is known

$Y = \begin{cases} 1 & \text{if FE} \\ 0 & \text{not} \end{cases}$ $X = \begin{cases} 1 & \text{if known} \\ 0 & \text{not} \end{cases}$

our uncertainty about $\underline{\Sigma}$ is in
layer (unconditionally) exchangeable
(partial)
but becomes conditionally exc h.



De Finetti's representation

Theorem for binary observables

(a) your uncertainty about $\underline{\Sigma} = (\Sigma_1, \dots, \Sigma_n)$ is

If (b) you're willing to extend this in (a) exch. and

judgment of exc h. from

$\underline{\Sigma} = (\Sigma_1, \Sigma_2, \dots, \Sigma_n)$ to generalizing

to population
 $\underline{\Sigma}^* = (\Sigma_1, \Sigma_2, \dots, \Sigma_n, \dots)$ then

1 (form $\bar{Y}_n = \frac{1}{n} \sum_{i=1}^n Y_i$) then

$\lim_{n \rightarrow \infty} \bar{Y}_n$ exists & $= \theta$ ($0 \leq \theta \leq 1$)

2 $P(Y_i = y_i | \mathcal{B}) = \begin{cases} \theta & \text{for } y_i = 1 \\ 1 - \theta & \text{for } y_i = 0 \\ 0 & \text{else} \end{cases}$
(i.e.) $(Y_i | \mathcal{B}) \sim \text{Bernoulli}(\theta)$
and 3

all coherent predictive dist. are expressible in following way:

$P(Y_1 = y_1, \dots, Y_n = y_n | \mathcal{B}) = \int_0^1 \theta^s (1-\theta)^{n-s} p(\theta) d\theta$
where $s = \sum_{i=1}^n y_i$
Likelihood from IID Bernoulli model
 $p(\theta)$ prior PDF

$$P(\underline{Y} = \gamma | \mathcal{B}) = \int_0^1 P(\underline{Y} = \gamma, \theta | \mathcal{B}) d\theta$$

STAT 131
LTP

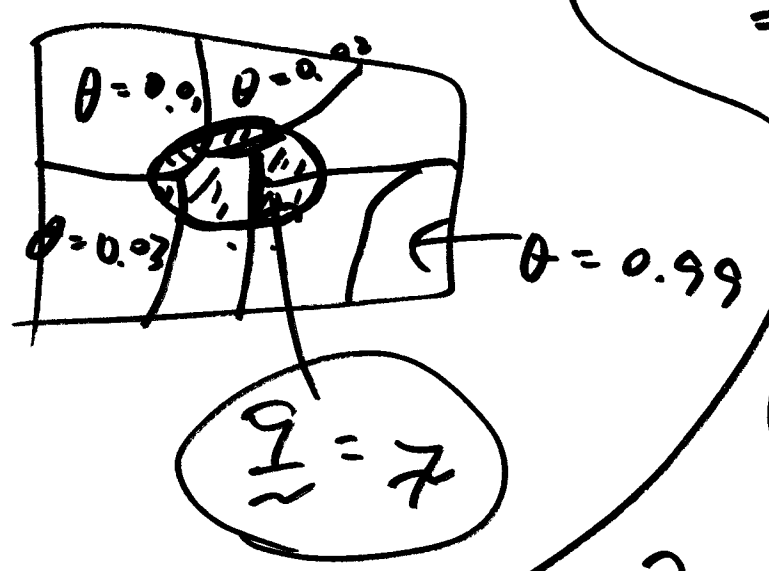
hard
but much easier if I knew θ ("truth")

$$= \int_0^1 P(\underline{Y} = \gamma | \theta, \mathcal{B}) \cdot$$

$$\boxed{p(\theta | \mathcal{B})} d\theta$$

$$= \int_0^1 \theta^s (1-\theta)^{h-s} \underline{p(\theta | \mathcal{B})} d\theta$$

$$\approx f = s$$



Bayesian story: if

λ is unknown to you,

model your information/uncertainty

about λ with a PDF/PMF

ie. treat λ as if it were a "random" variable $p(\lambda | \dots)$

If your uncertainty about binary observables $(\Sigma_1, \dots, \Sigma_n)$ not yet observed is ⁽¹⁰⁾ exchangeable, the only coherent modeling must look like this:

$$\begin{cases}
 \theta \in (0, 1) & \leftarrow \text{prior} \\
 (\theta | \mathcal{B}) \sim p(\theta | \mathcal{B}) \\
 (\Sigma_i | \theta \mathcal{B}) \stackrel{\text{i.i.d.}}{\sim} \text{Bernoulli}(\theta) & \leftarrow \text{freq.} \\
 (i = 1, \dots, n)
 \end{cases}$$

Bayesian hierarchical model

ie. Bayesians in this model are frequentists with a prior

ie. in this model, exch. \rightarrow sampling dist. ^{no} uncertainty