

# Vanilla Bayesian Inference & predictive paradigm

STAT 206  
19 Feb 21

$$P = (\mathcal{Q}, \mathcal{C}) \leftrightarrow (\tilde{\theta}, \mathcal{D}, \mathcal{B})$$

(natural language) (math)

AM disc. sec.

$$\rightarrow M_{IP} = \{ p(\theta | \mathcal{B}), p(\mathcal{D} | \theta, \mathcal{B}) \rightarrow \ell(\theta | \mathcal{D}, \mathcal{B}) \}$$

$\theta \in \Theta$  (parameter space)

input single choice for prior, single choice for sampling distribution

likelihood, use Bayes's Theorem

to get posterior  $p(\tilde{\theta} | \mathcal{D}, \mathcal{B}) =$   
 $c p(\tilde{\theta} | \mathcal{B}) \ell(\tilde{\theta} | \mathcal{D}, \mathcal{B})$

optimal inference given  $M_{IP}$

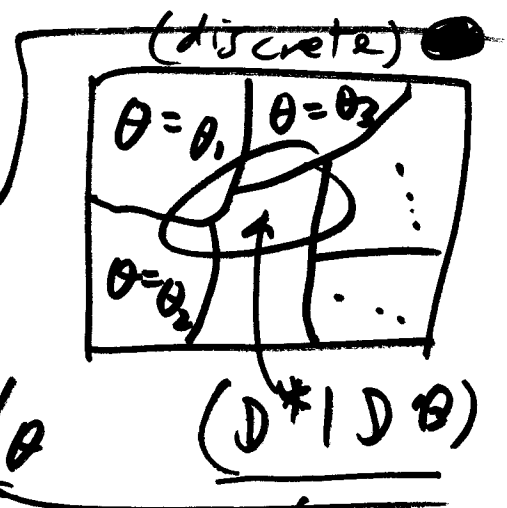
Bayesian prediction: new data  $\mathcal{D}^*$ ,  
 not yet observed: let's predict  $\mathcal{D}^*$   
 based on context, prior info, likelihood  
 info from  $\mathcal{D}$  (already-observed data)

posterior predictive distribution for  $D^*$  given  $D, \mathcal{C},$  prior info

$p(D^* | D, \mathcal{B}) = ?$  (hard) <sup>(prior)  $(M_{\text{pr}})$</sup>  (2)

as with de Finetti's theorem, the problem becomes easier if we know  $\theta$ , so let's bring  $\theta$  into the calculation

by partitioning over it:



$p(D^* | D, \mathcal{B}) = \int p(D^* | \theta, D, \mathcal{B}) d\theta$

$= \int p(D^* | \theta, D, \mathcal{B}) p(\theta | D, \mathcal{B}) d\theta$

this is just the Law of Total Probability for continuous

Now something excellent happens

ex.  $(Z_i | \theta) \stackrel{\text{IID}}{\sim} \text{Bernoulli}(\theta)$  (model  $\text{\textcircled{1}}$ )  $(i = 1, \dots, n)$

If I know  $\theta$ , the  $Z_i$  are conditionally independent given  $\theta$

Bayesian interpretation / definition of ③

independence:

$X_1, X_2$  Independent

if information about  $X_1$  doesn't change your (predictive) distribution for  $X_2$  &

vice versa

in model  $\Theta$ , if we know  $\theta$ , there's no information in  $(Z_1, \dots, Z_n)$  for

predicting  $Z_{n+1}$

but if we don't know ~~the vector~~

$\theta$ ,  $(Z_1, \dots, Z_n)$  has lots of information

in it for predicting  $Z_{n+1}$  (ex. in

first  $n=10$  data values,  $s=8$  were 1s;

you would predict  $Z_{n+1} =$  )

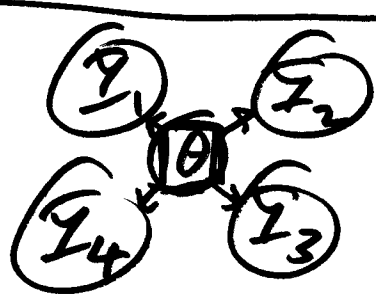
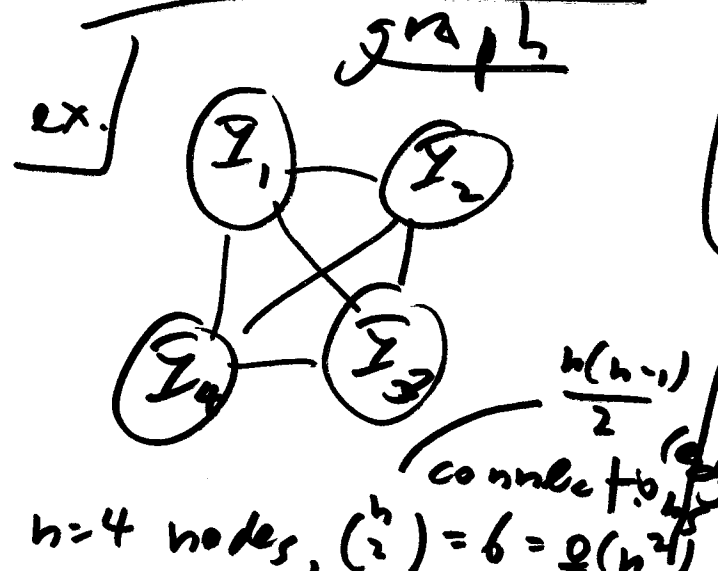
connection

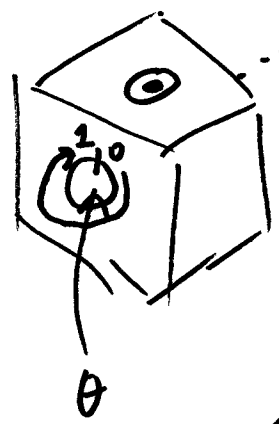
with Mr. de Finetti / before we see

the data  $(\bar{I} = \tau) : (\tau_1 = \tau, \tau_2 = \tau, \dots, \tau_n = \tau)$ , our uncertainty is exchangeable, but exchangeability  $\neq$   $(\text{II})$ : we just agree that there's information in any of the  $\tau_i$  that's helpful in predicting any of the other  $\tau_i$ , so exchangeable  $\tau_i$  are not independent; but in this example they become independent

if we know  $\theta$ :

exchangeability = conditionally IID given  $\theta$





$(Y_i | \theta) \stackrel{i.i.d.}{\sim} \text{Bernoulli}(\theta)$   
 $(i = 1, \dots, n)$

815, 205 ( $n=10$ )

if I don't know  $\theta$ , there

is useful information in

$(y_1, \dots, y_n)$  for predicting  $Y_{n+1}$ :

with  $(815, 205)$  in  $(y_1, \dots, y_n)$  ~~i.i.d.~~ ✓

you'd predict  $Y_{n+1} = 1$   
 if  $\theta$  is unknown  
 the  $(Y_1, \dots, Y_n)$  are dependent

but exchangeable

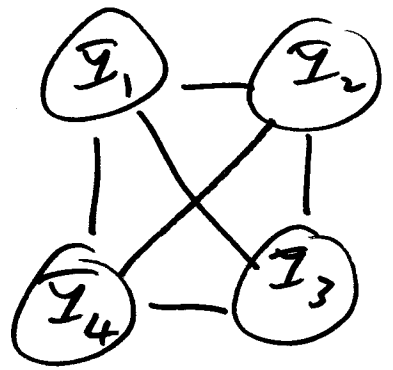
now instead,  
 $\theta$  known = 0.87

$(915, 905)$ : if I know  $\theta$ ,

there's no info in that helps  
 we predict  $Y_{n+1}$ ; if  $\theta$  is known,  
 $(Y_1, \dots, Y_n)$  are now (conditionally independent)  
 given  $\theta$

graph theory

(unknown  $\theta$ )



(nodes)  
 ○ (unknown)  
 □ (known)

$n = 4 = O(n)$   
 nodes

$\binom{n}{2} = \binom{4}{2}$   
 $= \frac{4 \cdot 3}{2 \cdot 1}$   
 $= 6$   
 edges

○ — ○ dependent

$\theta \rightarrow I_1$  "  $\theta$  causes  $I_1$  "  $= \frac{n(n-1)}{2}$

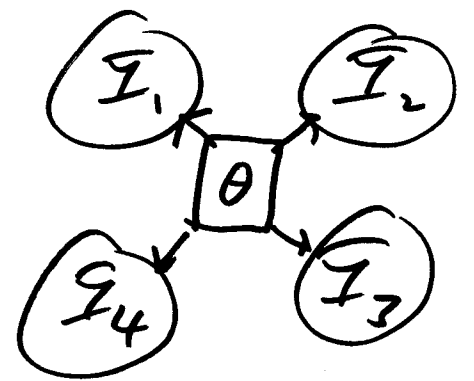
$= O(n^2)$

(complexity of the graph)

exchangeability

known  $\theta$

⑥



$\binom{n+1}{2} = 5$   
 ~~$\binom{n}{2}$~~   $= O(n)$   
 nodes

$\binom{n}{1} = n = O(n)$   
 edges

conditional independence given  $\theta$  dramatically reduces probabilistic model complexity

optimal prediction

$$p(D^* | D, \mathcal{B}) \quad \left[ P = (\theta, \mathcal{B}) \rightarrow (\theta, D, \mathcal{B}) \right]$$

posterior predictive distribution

given  $D, \text{prior } \pi_0, M_{EP}, \mathcal{B}$

$$p(D^* | D, \mathcal{B}) = \int p(D^* | \theta, \mathcal{B}) \cdot p(\theta | D, \mathcal{B}) d\theta$$

6.5  
case  
study

Lecture notes part 6 pp. 77-120

$(Z_i | \lambda, \mathcal{P}, \mathcal{B}) \stackrel{i.i.d.}{\sim} \text{poisson}(\lambda)$   
discrete  $(i=1, \dots, 4)$   
 $\lambda > 0$

PMF for  $Z_i$

$$P(Z_i = y_i | \lambda, \mathcal{P}, \mathcal{B}) = \frac{\lambda^{y_i} e^{-\lambda}}{y_i!} \quad I(y_i = 0 \text{ or } 1 \text{ or } \dots)$$

PDF for  $\lambda$

$$L(\lambda | z, \mathcal{P}, \mathcal{B}) = c \lambda^s e^{-4\lambda}, \quad s = \sum_{i=1}^4 z_i$$

(conjugate prior)

$$p(\lambda | \Gamma \mathcal{B}) = \Gamma(\alpha, \beta)$$

gamma  $\begin{pmatrix} \alpha > 0 \\ \beta > 0 \end{pmatrix}$

$$p(\lambda | \Gamma \mathcal{B}) = c \lambda^{\alpha-1} e^{-\beta \lambda} \mathbb{I}(\lambda > 0)$$

prior sample size =  $\beta$

prior mean (estimate of  $\lambda$ )

$$= \frac{\alpha}{\beta} = \mu$$

LI prior:  $\beta = \epsilon > 0$   
small (eg.) positive (0.001)

$$\alpha = \beta \mu$$

class:  $\mu = 5$

2 main uses

4 to h:  $\mu = 16$

for predictive dist.

① genuine interest in predicting new data (Loo)

② model-checking / leave one out = jackknife