## A Brief Introduction to Bayesian Inference

Mark Johnson

CG168 notes

## A brief review of *discrete* probability theory

- $\Omega$  is the set of all *elementary events* (c.f. interpretations in logic)
- If  $\omega \in \Omega$ , then  $P(\omega)$  is the probability of event  $\omega$ 
  - $P(\omega) \ge 0$
  - $\sum_{\omega \in \Omega} \mathbf{P}(\omega) = 1$
- A *random variable* X is a function from Ω to some set of values X
  - If  $\mathcal{X}$  is countable then X is a *discrete* random variable
  - If *X* is continuous then *X* is a *continuous* random variable
- If *x* is a possible value for *X*, then

$$P(X = x) = \sum_{\substack{\omega \in \Omega \\ X(\omega) = x}} P(\omega)$$

Independence and conditional distributions

- Two RVs *X* and *Y* are *independent* iff P(X, Y) = P(X)P(Y)
- The *conditional distribution* of *Y* given *X* is:

$$P(Y|X) = \frac{P(Y,X)}{P(X)}$$

so *X* and *Y* are independent iff P(Y|X) = P(Y) (here and below I assume strictly positive distributions)

• We can decompose the joint distribution of a sequence of RVs into a product of conditionals:

 $P(X_1,...,X_n) = P(X_1)P(X_2|X_1)P(X_3|X_2,X_1)...P(X_n|X_{n-1},...,X_1)$ 

i.e., the probability of generating  $X_1, \ldots, X_n$  "at once" is the same as generating them one at a time if each  $X_i$  is conditioned on the  $X_1, \ldots, X_{i-1}$  that preceded it

## Conditional distributions

• It's always possible to factor any distribution over  $X = (X_1, ..., X_n)$  into a product of conditionals

$$\mathbf{P}(\mathbf{X}) = \prod_{i=1}^{n} \mathbf{P}(X_i | X_1, \dots, X_{i-1})$$

• But in many interesting cases, *X<sub>i</sub>* depends only on a subset of *X*<sub>1</sub>,..., *X<sub>i-1</sub>*, i.e.,

$$\mathbf{P}(\mathbf{X}) = \prod_{i} \mathbf{P}(X_{i} | \mathbf{X}_{\mathsf{Pa}(i)})$$

where  $\mathsf{Pa}(i) \subseteq \{1, \ldots, i-1\}$  and  $X_S = \{X_j : j \in S\}$ 

- X and Y are *conditionally independent* given Z iff P(X, Y|Z) = P(X|Z) P(Y|Z) or equivalently, P(X|Y,Z) = P(X|Z)
- Note: the "parents" Pa(*i*) of *X<sub>i</sub>* depend on the order in which the variables are enumerated!



• A Bayes net is a graphical depiction of a factorization of a probability distribution into products of conditional distributions

$$\mathbf{P}(\mathbf{X}) = \prod_{i} \mathbf{P}(X_i | \mathbf{X}_{\mathsf{Pa}(i)})$$

A Bayes net has a node for each variable X<sub>i</sub> and an arc from X<sub>j</sub> to X<sub>i</sub> iff j ∈ Pa(i)

#### Bayes rule

• Bayes theorem:

$$P(Y|X) = \frac{P(X|Y) P(Y)}{P(X)}$$

- Bayes inversion: swap direction of arcs in Bayes net
- Interpreted as a recipe for "belief updating":

$$\underbrace{\frac{P(\text{Hypothesis}|\text{Data})}{\text{Posterior}} \propto \underbrace{\frac{P(\text{Data}|\text{Hypothesis})}{\text{Likelihood}} \underbrace{\frac{P(\text{Hypothesis})}{\text{Prior}}$$

• The normalizing constant (which you have to divide Likelihood times Prior by) is:

$$P(\mathsf{Data}) \ = \ \sum_{\mathsf{Hypothesis'}} P(\mathsf{Data}|\mathsf{Hypothesis'}) \ P(\mathsf{Hypothesis'})$$

which is the probability of generating the data under *any* hypothesis

Iterated Bayesian belief updating

• Suppose the data consists of 2 components  $D = (D_1, D_2)$ , and P(H) is our prior over hypotheses H

$$P(H|D_1, D_2) \propto P(D_1, D_2|H) P(H)$$
  
 
$$\propto P(D_2|H, D_1) P(H|D_1)$$

- This means the following are equivalent:
  - ► update the prior P(H) treating (D<sub>1</sub>, D<sub>2</sub>) as a single observation
  - ► update the prior P(H) wrt the first observation D<sub>1</sub> producing posterior P(H|D<sub>1</sub>) ∝ P(D<sub>1</sub>|H) P(H), which serves as the prior for the second observation D<sub>2</sub>

## Incremental Bayesian belief updating

- Consider a "two-part" data set (d<sub>1</sub>, d<sub>2</sub>). We show posterior obtained by Bayesian belief updating on (d<sub>1</sub>, d<sub>2</sub>) together is same as posterior obtained by updating on d<sub>1</sub> and then updating on d<sub>2</sub>.
- Bayesian belief updating on both  $(d_1, d_2)$  using prior P(H) P(H|d\_1, d\_2)  $\propto$  P(d\_1, d\_2|H) P(H) = P(d\_1, d\_2, H)
- Incremental Bayesian belief updating
  - Bayesian belief updating on  $d_1$  using prior P(H)

 $\mathbf{P}(H|d_1) \;\; \propto \;\; \mathbf{P}(d_1|H) \, \mathbf{P}(H) \;=\; \mathbf{P}(d_1,H)$ 

• Bayesian belief updating on  $d_2$  using prior  $P(H|d_1)$ 

$$P(H|d_1, d_2) \propto P(d_2|H, d_1) P(H|d_1)$$
  

$$\propto P(d_2|H, d_1) P(H, d_1)$$
  

$$= P(d_2, d_1, H)$$

### "Distributed according to" notation

- A *probability distribution F* is a non-negative function from some set  $\mathcal{X}$  whose values sum (integrate) to 1
- A random variable *X* is *distributed according* to a distribution *F*, or more simply, *X* has distribution *F*, written *X* ~ *F*, iff:

$$P(X = x) = F(x)$$
 for all  $x$ 

(This is for discrete RVs).

• You'll sometimes see the notion

$$X \mid Y \sim F$$

which means "X is generated conditional on Y with distribution F" (where F usually depends on Y)



#### Dirichlet priors for categorical and multinomial distributions

Comparing discrete and continuous hypotheses

## Continuous hypothesis spaces

• Bayes rule is the same when *H* ranges over a continuous space *except* that P(*H*) and P(*H*|*D*) are *continuous functions* of *H* 



• The normalizing constant is:

$$\mathbf{P}(D) = \int \mathbf{P}(D|H') \mathbf{P}(H') \, dH'$$

- Some of the approaches you can take:
  - Monte Carlo sampling procedures (which we'll talk about later)
  - ► Choose P(H) so that P(H|D) is easy to calculate ⇒ use a prior *conjugate* to the likelihood

## Categorical distributions

- A *categorical distribution* has a finite set of outcomes 1, ..., *m*
- A categorical distribution is parameterized by a vector  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_m)$ , where  $P(X = j | \boldsymbol{\theta}) = \theta_j$  (so  $\sum_{j=1}^m \theta_j = 1$ )
  - Example: An *m*-sided die, where  $\theta_i$  = prob. of face *j*
- Suppose  $X = (X_1, ..., X_n)$  and each  $X_i | \theta \sim CATEGORICAL(\theta)$ . Then:

$$P(\boldsymbol{X}|\boldsymbol{\theta}) = \prod_{i=1}^{n} CATEGORICAL(X_i; \boldsymbol{\theta}) = \prod_{j=1}^{m} \theta_j^{N_j}$$

where  $N_j$  is the number of times *j* occurs in **X**.

• Goal of next few slides: compute  $P(\theta|X)$ 

#### Multinomial distributions

- Suppose X<sub>i</sub> ~ CATEGORICAL(θ) for i = 1,..., n, and N<sub>j</sub> is the number of times j occurs in X
- Then  $N|n, \theta \sim MULTI(\theta, n)$ , and

$$\mathbf{P}(\boldsymbol{N}|\boldsymbol{n},\boldsymbol{\theta}) = \frac{\boldsymbol{n}!}{\prod_{j=1}^{m} N_j!} \prod_{j=1}^{m} \theta_j^{N_j}$$

where  $n! / \prod_{j=1}^{m} N_j!$  is the number of sequences of values with occurence counts *N* 

The vector N is known as a *sufficient statistic* for θ because it supplies as much information about θ as the original sequence X does.

## Dirichlet distributions

• *Dirichlet distributions* are probability distributions over multinomial parameter vectors

• called *Beta distributions* when m = 2

Parameterized by a vector *α* = (*α*<sub>1</sub>,..., *α*<sub>m</sub>) where *α*<sub>j</sub> > 0 that determines the shape of the distribution

$$DIR(\boldsymbol{\theta}; \boldsymbol{\alpha}) = \frac{1}{C(\boldsymbol{\alpha})} \prod_{j=1}^{m} \theta_{j}^{\alpha_{j}-1}$$
$$C(\boldsymbol{\alpha}) = \int \prod_{j=1}^{m} \theta_{j}^{\alpha_{j}-1} d\boldsymbol{\theta} = \frac{\prod_{j=1}^{m} \Gamma(\alpha_{j})}{\Gamma(\sum_{j=1}^{m} \alpha_{j})}$$

- Γ is a generalization of the factorial function
- $\Gamma(k) = (k 1)!$  for positive integer *k*
- $\Gamma(x) = (x-1)\Gamma(x-1)$  for all x

#### Plots of the Dirichlet distribution



#### Plots of the Dirichlet distribution (2)



#### Plots of the Dirichlet distribution (3)



#### Dirichlet distributions as priors for $\theta$

• Generative model:

• We can depict this as a Bayes net using *plates*, which indicate *replication* 



Inference for  $\theta$  with Dirichlet priors

- Data  $X = (X_1, ..., X_n)$  generated i.i.d. from CATEGORICAL( $\theta$ )
- Prior is DIR(*α*). By Bayes Rule, posterior is:

$$P(\boldsymbol{\theta}|\boldsymbol{X}) \propto P(\boldsymbol{X}|\boldsymbol{\theta}) P(\boldsymbol{\theta})$$
$$\propto \left(\prod_{j=1}^{m} \theta_{j}^{N_{j}}\right) \left(\prod_{j=1}^{m} \theta_{j}^{\alpha_{j}-1}\right)$$
$$= \prod_{j=1}^{m} \theta_{j}^{N_{j}+\alpha_{j}-1}, \text{ so}$$
$$P(\boldsymbol{\theta}|\boldsymbol{X}) = DIR(\boldsymbol{N}+\boldsymbol{\alpha})$$

- So if prior is Dirichlet with parameters  $\alpha$ , posterior is Dirichlet with parameters  $N + \alpha$
- ⇒ can regard Dirichlet parameters α as "pseudo-counts" from "pseudo-data"

#### Point estimates from Bayesian posteriors

- A "true" Bayesian prefers to use the full P(*H*|*D*), but sometimes we have to choose a "best" hypothesis
- The Maximum a posteriori (MAP) or posterior mode is

$$\widehat{H} = \operatorname{argmax}_{H} \operatorname{P}(H|D) = \operatorname{argmax}_{H} \operatorname{P}(D|H) \operatorname{P}(H)$$

• The *expected value*  $E_P[X]$  of X under distribution P is:

$$\mathbf{E}_{\mathbf{P}}[\mathbf{X}] = \int x \, \mathbf{P}(\mathbf{X} = x) \, dx$$

The expected value is a kind of average, weighted by P(X). The *expected value*  $E[\theta]$  of  $\theta$  is an estimate of  $\theta$ .

#### The posterior mode of a Dirichlet

• The Maximum a posteriori (MAP) or posterior mode is

$$\widehat{H} = \operatorname{argmax}_{H} \operatorname{P}(H|D) = \operatorname{argmax}_{H} \operatorname{P}(D|H) \operatorname{P}(H)$$

• For Dirichlets with parameters *α*, the MAP estimate is:

$$\hat{\theta}_j = \frac{\alpha_j - 1}{\sum_{j'=1}^m (\alpha_{j'} - 1)}$$

so if the posterior is  $DIR(N + \alpha)$ , the MAP estimate for  $\theta$  is:

$$\hat{ heta}_{j} = rac{N_{j} + lpha_{j} - 1}{n + \sum_{j'=1}^{m} (lpha_{j'} - 1)}$$

• If  $\alpha = 1$  then  $\hat{\theta}_j = N_j/n$ , which is also the *maximum likelihood estimate* (MLE) for  $\theta$ 

## The expected value of $\theta$ for a Dirichlet

• The *expected value*  $E_P[X]$  of X under distribution P is:

$$\mathbf{E}_{\mathbf{P}}[X] = \int x \, \mathbf{P}(X=x) \, dx$$

• For Dirichlets with parameters  $\alpha$ , the expected value of  $\theta_j$  is:

$$\mathbf{E}_{\mathrm{DIR}(\boldsymbol{\alpha})}[\theta_j] = \frac{\alpha_j}{\sum_{j'=1}^m \alpha_{j'}}$$

• Thus if the posterior is  $DIR(N + \alpha)$ , the expected value of  $\theta_i$  is:

$$\mathbf{E}_{\mathrm{DIR}(\mathbf{N}+\boldsymbol{\alpha})}[\theta_j] = \frac{N_j + \alpha_j}{n + \sum_{j'=1}^m \alpha_{j'}}$$

• E[*θ*] *smooths* or *regularizes* the MLE by adding pseudo-counts *α* to *N* 

## Sampling from a Dirichlet

$$\boldsymbol{\theta} \mid \boldsymbol{\alpha} \sim \text{DIR}(\boldsymbol{\alpha}) \text{ iff } P(\boldsymbol{\theta} \mid \boldsymbol{\alpha}) = \frac{1}{C(\boldsymbol{\alpha})} \prod_{j=1}^{m} \theta_{j}^{\alpha_{j}-1}, \text{ where:}$$
  
 $C(\boldsymbol{\alpha}) = \frac{\prod_{j=1}^{m} \Gamma(\alpha_{j})}{\Gamma(\sum_{j=1}^{m} \alpha_{j})}$ 

- There are several algorithms for producing samples from DIR(*α*). A simple one relies on the following result:
- If  $V_k \sim \text{GAMMA}(\alpha_k)$  and  $\theta_k = V_k / (\sum_{k'=1}^m V_{k'})$ , then  $\theta \sim \text{DIR}(\alpha)$
- This leads to the following algorithm for producing a sample *θ* from DIR(*α*)
  - Sample  $v_k$  from GAMMA $(\alpha_k)$  for k = 1, ..., m

• Set 
$$\theta_k = v_k / (\sum_{k'=1}^m v_{k'})$$

## **Conjugate priors**

• If prior is  $DIR(\alpha)$  and likelihood is i.i.d.  $CATEGORICAL(\theta)$ , then posterior is  $DIR(N + \alpha)$ 

 $\Rightarrow$  prior parameters  $\alpha$  specify "pseudo-observations"

- A class C of prior distributions P(H) is *conjugate* to a class of likelihood functions P(D|H) iff the posterior P(H|D) is also a member of C
- In general, conjugate priors encode "pseudo-observations"
  - ► the difference between prior P(*H*) and posterior P(*H*|*D*) are the observations in *D*
  - ▶ but P(H|D) belongs to same family as P(H), and can serve as prior for inferences about more data D'
  - $\Rightarrow$  must be possible to encode observations *D* using parameters of prior
- In general, the likelihood functions that have conjugate priors belong to the *exponential family*

#### Outline

Dirichlet priors for categorical and multinomial distributions

Comparing discrete and continuous hypotheses

# Categorical and continuous hypotheses about coin flips

- Data: A sequence of coin flips  $X = (X_1, \ldots, X_n)$
- Hypothesis  $h_1$ : *X* is generated from a fair coin, i.e.,  $\theta_H = 0.5$
- Hypothesis  $h_2$ : X is generated from a biased coin with unknown bias, i.e.,  $\theta_H \sim \text{DIR}(\alpha)$

$$P(H|X) = P(X|H) P(H)$$

- Assume  $P(h_1) = P(h_2) = 0.5$
- $P(X|h_1) = 2^{-n}$ , but *what is*  $P(X|h_2)$ ?
- P(X|h<sub>2</sub>) is the probability of generating θ from DIR(α) and then generating X from CATEGORICAL(θ). But we don't care about the value of θ, so we *marginalize* or *integrate out* θ

$$P(\boldsymbol{X}|\boldsymbol{\alpha},h_2) = \int P(\boldsymbol{X},\boldsymbol{\theta}|\boldsymbol{\alpha}) d\boldsymbol{\theta}$$

#### Posterior with Dirichlet priors

• *Integrate out*  $\theta$  to calculate posterior probability of *X* 

$$P(\mathbf{X}|\boldsymbol{\alpha}) = \int P(\mathbf{X}, \boldsymbol{\theta}|\boldsymbol{\alpha}) d\boldsymbol{\theta} = \int P(\mathbf{X}|\boldsymbol{\theta}) P(\boldsymbol{\theta}|\boldsymbol{\alpha}) d\boldsymbol{\theta}$$
$$= \int \left(\prod_{j=1}^{m} \theta_{j}^{N_{j}}\right) \left(\frac{1}{C(\boldsymbol{\alpha})} \prod_{j=1}^{m} \theta_{j}^{\alpha_{j}-1}\right) d\boldsymbol{\theta}$$
$$= \frac{1}{C(\boldsymbol{\alpha})} \int \prod_{j=1}^{m} \theta_{j}^{N_{j}+\alpha_{j}-1} d\boldsymbol{\theta}$$
$$= \frac{C(N+\boldsymbol{\alpha})}{C(\boldsymbol{\alpha})}, \text{ where } C(\boldsymbol{\alpha}) = \frac{\prod_{j=1}^{m} \Gamma(\alpha_{j})}{\Gamma(\sum_{j=1}^{m} \alpha_{j})}$$

 Collapsed Gibbs samplers and the Chinese Restaurant Process rely on this result

#### Posteriors under $h_1$ and $h_2$



## Understanding the posterior

$$P(\mathbf{X}|\boldsymbol{\alpha}) = \frac{C(N+\boldsymbol{\alpha})}{C(\boldsymbol{\alpha})} \text{ where } C(\boldsymbol{\alpha}) = \frac{\prod_{j=1}^{m} \Gamma(\alpha_{j})}{\Gamma(\boldsymbol{\alpha} \bullet)} \text{ and } \boldsymbol{\alpha} \bullet = \sum_{j=1}^{m} \alpha_{j}$$

$$P(\mathbf{X}|\boldsymbol{\alpha}) = \left(\frac{\prod_{j=1}^{m} \Gamma(N_{j} + \alpha_{j})}{\Gamma(n + \boldsymbol{\alpha} \bullet)}\right) \left(\frac{\Gamma(\boldsymbol{\alpha} \bullet)}{\prod_{j=1}^{m} \Gamma(\alpha_{j})}\right)$$

$$= \left(\prod_{j=1}^{m} \frac{\Gamma(N_{j} + \alpha_{j})}{\Gamma(\alpha_{j})}\right) \left(\frac{\Gamma(\boldsymbol{\alpha} \bullet)}{\Gamma(n + \boldsymbol{\alpha} \bullet)}\right)$$

$$= \frac{\alpha_{1}}{\alpha_{\bullet}} \times \frac{\alpha_{1} + 1}{\alpha_{\bullet} + 1} \times \ldots \times \frac{\alpha_{1} + N_{1} - 1}{\alpha_{\bullet} + N_{1} - 1}$$

$$\times \frac{\alpha_{2}}{\alpha_{\bullet} + N_{1}} \times \frac{\alpha_{2} + 1}{\alpha_{\bullet} + N_{1} + 1} \times \ldots \times \frac{\alpha_{2} + N_{2} - 1}{\alpha_{\bullet} + N_{1} + N_{2} - 1}$$

$$\times \ldots$$

$$\times \frac{\alpha_{m}}{\alpha_{\bullet} + n - N_{m} - 1} \times \frac{\alpha_{m} + 1}{\alpha_{\bullet} + n - N_{m}} \times \ldots \times \frac{\alpha_{m} + N_{m} - 1}{\alpha_{\bullet} + n - 1}$$

## Exchangability

- The individual  $X_i$  in a Dirichlet-multinomial distribution  $P(X|\alpha) = C(N + \alpha)/C(\alpha)$  are *not independent* 
  - the probability of  $X_i$  depends on  $X_1, \ldots, X_{i-1}$

$$P(X_n = k | X_1, \dots, X_{n-1}, \boldsymbol{\alpha}) = \frac{P(X_1, \dots, X_n | \boldsymbol{\alpha})}{P(X_1, \dots, X_{n-1} | \boldsymbol{\alpha})}$$
$$= \frac{\alpha_k + N_k(X_1, \dots, X_{n-1})}{\alpha_{\bullet} + n - 1}$$

- but  $X_1, \ldots, X_n$  are *exchangable* 
  - $P(X|\alpha)$  depends only on *N*
  - $\Rightarrow$  doesn't depend on *the order* in which the *X* occur
- A distribution over a sequence of random variables is *exchangable* iff *the probability of all permutations of the random variables are equal*

## Summary so far

- Bayesian inference can compare models of *differing complexity* (assuming we can *calculate posterior probability*)
  - ▶ Hypothesis *h*<sup>1</sup> has no free parameters
  - Hypothesis  $h_2$  has one free parameter  $\theta_H$
- *Bayesian Occam's Razor:* "A more complex hypothesis is only prefered if its greater complexity consistently provides a better account of the data"
- But:  $h_1$  makes every sequence equally likely.  $h_2$  seems to *dislike*  $\theta_H \approx 0.5$ What's going on here?

Posteriors with n = 10,  $\alpha = 10$ 



Posteriors with  $n = 20, \alpha = 1$ 



#### Dirichlet-Multinomial distributions

- Only one sequence of 10 heads out of 10 coin flips
- but 252 different sequences of 5 heads out of 10 coin flips
- Each particular sequence of 5 heads out of 10 flips is unlikely, but there are so many of them that *the group is very likely*
- The number of ways of picking *N* outcomes out of *n* trials is:

$$\frac{n!}{\prod_{j=1}^m N_j!}$$

• The probability of observing N given  $\theta$  is:

$$P(\boldsymbol{N}|\boldsymbol{\theta}) = \frac{n}{\prod_{j=1}^{m} N_j!} \prod_{j=1}^{m} \theta_j^{N_j}$$

• The probability of observing N given  $\alpha$  is:

$$P(\boldsymbol{N}|\boldsymbol{\alpha}) = \frac{n}{\prod_{j=1}^{m} N_j!} \frac{C(\boldsymbol{N}+\boldsymbol{\alpha})}{C(\boldsymbol{\alpha})}$$

Dirichlet-**multinomial** posteriors with  $n = 10, \alpha = 1$ 



Dirichlet-**multinomial** posteriors with n = 10, varying  $\alpha$ 



Dirichlet-**multinomial** posteriors with n = 20, varying  $\alpha$ 



Dirichlet-**multinomial** posteriors with n = 50, varying  $\alpha$ 



Entropy vs. "rich get richer"

• Notation: If  $X = (X_1, ..., X_n)$ , then  $X_{-j} = (X_1, ..., X_{j-1}, X_{j+1}, ..., X_n)$ 

$$P(X_n = k | \boldsymbol{\alpha}, \boldsymbol{X}_{-n}) = \frac{N_k(\boldsymbol{X}_{-n}) + \alpha_k}{\alpha_{\bullet} + n - 1}$$

- The probability of generating an outcome is proportional to the number of times it has been seen before (including prior)
- ⇒ Next outcome is most likely to be most frequently generated previous outcome ⇒ *sparse outcomes* 
  - But there are far fewer sparse outcomes than non-sparse outcomes ⇒ entropy "prefers" non-sparse outcomes
  - If α > 1 then most likely outcomes are not sparse i.e., entropy is stronger than prior
  - If α < 1 then most likely outcomes are sparse i.e., prior is stronger than entropy