

# Bayesian Statistics

In the previous lectures, especially on maximum Likelihood method, we focused on  $\text{PCDIO}$  but noted

that, sometimes, when noise effects are large, it may be better not to take the max likelihood estimate. In other words, we can accept some amount of bias to reduce the impact of variance.

<sup>extreme</sup>  
Our example:  $x_i \text{ i.i.d } N(\mu, \sigma^2)$   
 $i=1, \dots, N$

MLE for  $\mu$ :  $\hat{\mu}_{\text{MLE}} = \bar{x}$

$$E(\bar{x} - \mu)^2 = \frac{\sigma^2}{N} \quad \text{MSE}_{\text{MLE}} = \frac{\sigma^2}{N}$$

No bias, all variance.

However, if  $|\mu| < \frac{\sigma}{\sqrt{N}}$  (Sample size too small)

$\hat{\mu}_0 = 0$  has less error.

No variance, but bias.

$$\text{MSE}_0 = \mu^2 \quad \text{MSE}_0 < \text{MSE}_{\text{MLE}}$$

Bayesian approach is one way  
of balancing bias-variance  
trade-off.

Now we allow talking about probability distribution of  $\theta$ 's themselves

Note that  $\log P(\theta | D)$

$$= \log P(D|\theta) + \log P_0(\theta) - \log P(D)$$

log likelihood

additional regularizer constant  
penalty function  
for  $\theta$

If we want to find the  ~~$\hat{\theta}$~~  "most likely"  $\theta$ , we optimize  $\log P(\theta | D)$ . The additional regularization penalty alters  $\hat{\theta}$  from  $\hat{\theta}_{MLE}$ :

This estimate is called the maximum a posteriori (MAP) estimate.

Several non-Bayesian approaches also advocate using additional penalties for estimation in order to control.

(believe)  
Bayesians [redacted] there is more to Bayesian approach. Before we go there, let us look at [redacted] some examples, though.

I will skip examples of discrete choice of hypotheses, as discussed in the beginning of Chapter 3. Let us start with the beta-binomial model.

$$X_i \sim \text{Ber}(\theta)$$

$$\text{Prob}[x_i=1]=\theta$$

$$\text{Prob}[x_i=0]=1-\theta$$

$X_i$ 's [redacted] independent

[Ex: Head-fails for a [redacted] possibly biased coin]

$$P(D|\theta) = \theta^{N_1} (1-\theta)^{N_0}$$

  $N_1 = \sum_{i=1}^N I(x_i=1)$

$$N_0 = \sum_{i=1}^N I(x_i=0)$$

$$N_1 + N_0 = N$$

Note that  $N_1$  is a sufficient statistics for  $\theta$ . We don't need to know which order successes happened to infer  $\theta$ .

$$N_1 \sim \text{Bin}(N, \theta)$$

$$\text{Bin}(k|N, \theta) = \binom{N}{k} \theta^k (1-\theta)^{N-k}$$

Now let us choose a prior. Note that if we choose prior  $p(\theta) \propto \theta^{a-1} (1-\theta)^{b-1}$ , namely  $\theta \sim \text{Beta}(a, b)$ , then,

$$P(\theta|D) \propto P(D|\theta) p(\theta)$$

$$\propto \theta^{N_1+a-1} (1-\theta)^{N_0+b-1}$$

So the posterior distribution follows a beta

distribution as well.

When the posterior and prior belongs to the same family of distributions the prior and the posterior are called conjugate distribution and the prior is called a conjugate prior for the likelihood function.

In this case, the ~~prior~~ prior is effectively adding some pseudo-counts  $\beta$  to  $N$ ,  $N_0$ .

The parameters  $a, b$  are called hyperparameters. We will later talk about how to choose these hyperparameters.

If we set  $a=1, b=1$ , we get an uniform prior.

In principle we could update things sequentially

prior  $\xrightarrow{\text{data } 1}$  Posterior = new prior  $\xrightarrow{\text{data } 2}$  new posterior

is equivalent to

Prior  $\xrightarrow{(\text{data } 1, \text{data } 2)}$  new posterior

In many circumstances, when data arrives sequentially, you can have an ~~online~~ online algorithm updating things sequentially.

Posterior mean and mode

$$\text{Mode } \hat{\theta}_{\text{map}} = \arg \max_{\theta} \log \left[ \frac{N_1+a-1}{\theta} \frac{N_0+b-1}{(1-\theta)} \right]$$

$$\theta \text{ deriv} \quad \frac{N_1+a-1}{\theta} - \frac{N_0+b-1}{1-\theta} = 0$$

$$\Rightarrow \frac{\theta}{1-\theta} = \frac{N_1+a-1}{N_0+b-1} \Rightarrow \hat{\theta}_{\text{MAP}} = \frac{N_1+a-1}{N_0+N_1+a+b-2} = \frac{N_1+a-1}{N+a+b-2}$$

$$\text{When } a=b=1 \quad \hat{\theta}_{\text{MAP}} = \frac{N_1}{N} = \theta_{\text{MLE}}$$

What is the average posterior?

$$\bar{\theta} = \frac{\int_0^1 d\theta \theta^{\gamma_1+1} (1-\theta)^{\gamma_2}}{\int_0^1 d\theta \theta^{\gamma_1} (1-\theta)^{\gamma_2}}$$

where  $\gamma_1 = N_1 + a - 1$        $\gamma_2 = N_2 + b - 1$

$$\bar{\theta} = \frac{\int_0^1 d\theta \theta^{\gamma_1+1} (1-\theta)^{\gamma_2}}{\int_0^1 d\theta \theta^{\gamma_1} (1-\theta)^{\gamma_2}} = \frac{B(\gamma_1+1, \gamma_2+1)}{B(\gamma_1+1, \gamma_2+1)}$$

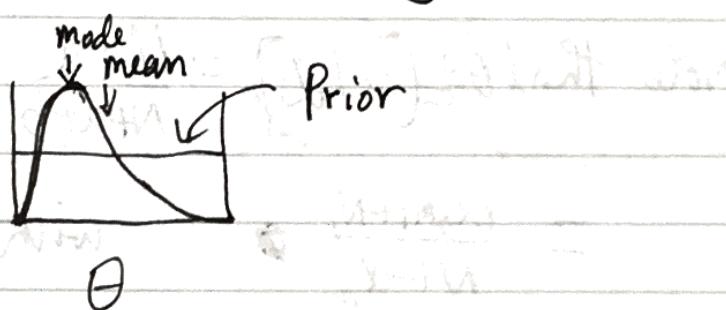
$$= \frac{\Gamma(\gamma_1+2) \Gamma(\gamma_2+1) / \Gamma(\gamma_1+\gamma_2+3)}{\Gamma(\gamma_1+1) \Gamma(\gamma_2+1) / \Gamma(\gamma_1+\gamma_2+2)}$$

$$= \frac{\gamma_1+1}{\gamma_1+\gamma_2+2} = \frac{N_1+a}{N_1+a+b}$$

This is different from  $\hat{\theta}_{MAP}$   
Even if I took  $a=b=1$ ,

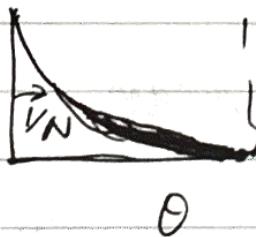
$$\text{I get } \bar{\theta} = \frac{N_1+1}{N+2}$$

The mean posterior seems to have inserted 'pseudocounts' by itself. This is an example of Bayesian averaging 'regularizing' the estimate



Such ~~smooth~~ 'smoothing' is important when data is sparse. Imagine  $N=0$ , no success observed so far.

Can I conclude  $\theta = 0$ ? With uniform prior:  $P(\theta|D) \propto (1-\theta)^N \approx e^{-N\theta}$



Intuitively that observing is saying that we believe  $\theta \sim O(\frac{1}{N})$

and that is why perhaps we have not observed it.

$$\bar{\theta} = \frac{0+1}{N+2} = \frac{1}{N+2}$$

Posterior mean  $\propto \frac{N+1}{N+2}$  as an estimate gives

Laplace's rule of succession  
and it is also known as  
add-one smoothing!

$$\text{Note that } \bar{\theta} = E[\theta|D] = \frac{N_1 + a}{N + a_0}$$

$$= \frac{\alpha_0 + N_1}{N + a_0} \quad \text{with } a_0 = a + b$$

$$m_1 = \frac{a}{a+b}$$

Prior mean



$$E[\theta|D] = \left( \frac{\alpha_0}{N+a_0} \right) m_1 + \left( \frac{N}{N+a_0} \right) \hat{\theta}_{\text{MLE}}$$

Prior mean                                  MLE

$$= \lambda m_1 + (1-\lambda) \hat{\theta}_{\text{MLE}}$$

?    ↑  
Convex combination

As  $N \rightarrow \infty$   $\bar{\theta}$  goes to  $\hat{\theta}_{\text{MLE}}$

With lots of data, priors do not matter.

## Posterior Variance

$$\text{var}[\theta | D] = E[\theta^2 | D] - E[\theta | D]^2$$

$$= \frac{\Gamma(\gamma_1 + 3)\Gamma(\gamma_2 + 1)/\Gamma(\gamma_1 + \gamma_2 + 4)}{\Gamma(\gamma_1 + 1)\Gamma(\gamma_2 + 1)/\Gamma(\gamma_1 + \gamma_2 + 2)} - \left( \frac{\gamma_1 + 1}{\gamma_1 + \gamma_2 + 2} \right)^2$$

$$= \frac{(\gamma_1 + 2)(\gamma_1 + 1)}{(\gamma_1 + \gamma_2 + 3)(\gamma_1 + \gamma_2 + 2)} - \frac{(\gamma_1 + 1)(\gamma_1 + 1)}{(\gamma_1 + \gamma_2 + 2)(\gamma_1 + \gamma_2 + 2)}$$

$$= \frac{\gamma_1 + 1}{\gamma_1 + \gamma_2 + 2} \left[ \frac{\gamma_1 + 2}{\gamma_1 + \gamma_2 + 3} - \frac{\gamma_1 + 1}{\gamma_1 + \gamma_2 + 2} \right]$$

$$= \frac{(\gamma_1 + 1)[2(\gamma_1 + \gamma_2) + 4 - \{4\gamma_1 + \gamma_2 + 3\}]}{(\gamma_1 + \gamma_2 + 2)^2(\gamma_1 + \gamma_2 + 3)}$$

$$= \frac{(\gamma_1 + 1)(\gamma_2 + 1)}{(\gamma_1 + \gamma_2 + 2)^2 / \gamma_1 + \gamma_2 + 3}$$

$$= \frac{(N_1 + a)(N_0 + b)}{(N_1 + N_0 + a + b)^2 (N_1 + N_0 + a + b + 1)}$$

~~REASON~~ ~~N1 > N0~~

~~N1 > N0~~

when  $N$  is large compared to  $a, b$



$$\text{Var}[\hat{\theta}] \approx \frac{N_1 N_2}{N^3}$$

$$\approx \frac{\hat{\theta}(1-\hat{\theta})}{N}$$

where  $\hat{\theta} = \hat{\theta}_{\text{MLE}}$ .

This is related to  $\hat{\theta}_{\text{MLE}} = \frac{N_1}{N}$

following central limit theorem  
with the std / error bar

$$\sigma = \sqrt{\frac{\hat{\theta}(1-\hat{\theta})}{N}}$$

Note that  $\sigma$  is the largest  
when  $\hat{\theta} = \frac{1}{2}$ , making it

harder to decide whether a  
coin is fair, in comparison to  
establishing large bias.

More than two outcomes: Beta-Binomial  
→ Dirichlet - Multinomial

$N$  'dice' rolls       $\{x_1, \dots, x_N\}$        $x_i \in \{1, \dots, k\}$

$$P(D|\theta) = \prod_k \theta_k^{N_k}$$

$$\sum_k \theta_k = 1$$

$$\sum_k N_k = N$$

For real dice       $\theta_1, \dots, \theta_6$

Honestly, if we were dealing with coin tosses or dice rolls, we can have  $N$  large enough so that all the Bayesian caution may not be too relevant unless  $\theta \sim 1/N$ . In practice one may be dealing with text classification or words in DNA. For words in a natural language  $K \sim 10^5$ . For DNA sequences 6 bases long,  $K = 4^6 \approx 4 \times 10^3$ . Some words may be rare enough not to show up in the training data set, but appear in the test data. We could also have a very common word in the training data be missing in the test set. In all these cases, Bayesian average helps.

$$\text{prior} : \text{Dir}(\theta | \alpha) = \frac{1}{B(\alpha)} \prod_{k=1}^K \theta_k^{\alpha_k - 1}$$

Conjugate prior

$$\text{Posterior } P(\theta | D) \propto P(D | \theta) h(\theta)$$

$$\propto \prod_{k=1}^K \theta_k^{\alpha_k + N_k - 1}$$

$$= \text{Dir}(\theta | \alpha_1 + N_1, \dots, \alpha_K + N_K)$$

Optimize subject to constraint  $\sum \theta_k = 1$

$$\text{Mode } l(\theta, \lambda) = \sum_k (N_k + \alpha_k - 1) \log \theta_k$$

$$+ \lambda \left( 1 - \sum_k \theta_k \right)$$

$$\frac{\partial l}{\partial \theta_k} = 0 \rightarrow \frac{N_k + \alpha_k - 1}{\theta_k} = \lambda \quad \text{or} \quad (N_k + \alpha_k - 1) = \lambda \theta_k$$

$$\text{Sum over } k \quad N + \alpha_0 - k = \lambda$$

$$\alpha_0 = \sum_{k=1}^K \alpha_k$$

$$\text{So } \hat{\theta}_k = \frac{N_k + \alpha_k - 1}{\lambda} = \frac{N_k + \alpha_k - 1}{N + \alpha_0 - K}$$

That is the MLE estimate  $\hat{\theta}_k = \frac{N_k}{N}$  for MLE.

## Posterior predictive

$$p(x=j|D) = \int p(x=j|\theta) p(\theta|D) d\theta$$

$$= \int p(x=j|\theta_j) P(\theta_j|D) d\theta_j$$

$$= \int \theta_j^{\gamma_j} (1-\theta_j)^{r_i} d\theta_j = \bar{\theta}_j$$

$$P(\theta_j|D) = \frac{\int \theta_j^{\gamma_j} (1-\theta_j)^{r_i} d\theta_j}{\int \int \theta_j^{\gamma_j} (1-\theta_j)^{r_i} d\theta_j}$$

$$= \frac{\theta_j^{\gamma_j} (1-\theta_j)^{r_i}}{\int \theta_j^{\gamma_j} (1-\theta_j)^{r_i} d\theta_j}$$

$$\int \theta_j^{\gamma_j} (1-\theta_j)^{r_i} d\theta_j$$

$$\gamma_j = N_j + \alpha_j - 1$$

$$\bar{\theta}_j = \frac{\gamma_j + 1}{\gamma_j + \sum_{i \neq j} \gamma_i + K-1 + 1} = \frac{N_j + \alpha_j}{N + \alpha_0}$$

Once more the posterior predictive avoids the zero count problem [Even when we have a flat prior;  $\alpha_j = 1$ ].

Bag of word model, document classification  
Naïve Bayes.

Bag of word model

Mary had [1] little lamb, little lamb, little lamb,  
Mary had [2] little lamb, its fleece [3] white [4] snow

Drop 'stop words'

Vocab list

mary lamb little big fleece white black snow

rain unk.  
q ↑ unknown

mary unk little lamb little lamb little lamb

1 10 3 2 3 2 3 2

mary unk little lamb unk fleece white snow

1 10 3 2 10 5 8 8

$$N = 16 \quad \alpha_j = 1$$

$$\bar{\theta}_j = \frac{N_j + 1}{16 + 10}$$

$$\bar{\theta}_{\text{mary}} = \frac{2+1}{16+10} = \frac{3}{26}, \quad \bar{\theta}_{\text{big}} = \frac{1}{16+10} = \frac{1}{26}$$

big never shows up

Naive Bayes!  $(x_1, \dots, x_D)$

'Naive' assumption feature conditionally independent  
(conditioned upon class).

$$p(x|y=c, \theta) = \prod_{j=1}^D p(x_j|y=c, \theta_{jc})$$

Example: 1) For real-valued ~~real-valued~~ features  $p(x_j|y=c, \theta) = \frac{1}{\sqrt{2\pi}} N(x_j; \mu_c, \sigma_c^2)$

2) Binary features  $x_j \in \{0, 1\}$ ,  $p(x_j|y=c, \theta) = \prod_{j=1}^D \text{Ber}(x_j; \theta_c)$

Let us see an application to document classification.

Bernoulli document model:

$D$  'feature' words

|                  |                   |         |                   |
|------------------|-------------------|---------|-------------------|
| $x_1$            | $x_2$             | $\dots$ | $x_D$             |
| 0                | 1                 | $\dots$ | 0                 |
| ↑                | ↑                 | $\dots$ | ↑                 |
| Word 1<br>absent | Word 2<br>present | $\dots$ | Word D<br>present |

We throw away the details of word abundance

Bayesian priors

$\pi = (\pi_1, \dots, \pi_C) \rightarrow$  Probs of belonging to different classes

$$p(\pi) \rightarrow \text{Dir}(\alpha)$$

$$p(\theta) = \prod_{j=1}^D \prod_{c=1}^C p(\theta_{jc})$$

$$\rightarrow \text{Beta}(\beta_0, \beta_1)$$

With training data  $\mathcal{D} \rightarrow N_1, \dots, N_c$

$\frac{1}{\text{numbers}}$

$\frac{1}{\text{different classes}}$

$N_{jc}$  = No. of word  $j$  in docs  
of class  $c$

$$P(\theta | D) = P(\pi | D) \prod_{j=1}^J \prod_{c=1}^C P(\theta_{jc} | D)$$

$$P(\pi | D) = \text{Dir}(N_1 + \alpha_1, \dots, N_c + \alpha_c) \propto \prod_{c=1}^C \bar{\pi}_c^{N_c + \alpha_c}$$

$$P(\theta_{jc} | D) = \text{Beta}(N_c - N_{jc} + \beta_j, N_{jc} + \beta_j) \propto \theta_{jc}^{N_{jc} + \beta_j} (1 - \theta_{jc})^{\frac{N_c - N_{jc} + \beta_j}{N_c + \beta_j}}$$

For posterior predictive; integrate out parameters

$$\begin{aligned} P(y=c | x, D) &\propto \int_D \text{Cat}(y=c | \pi) P(\pi | D) d\pi \\ &\propto \prod_{j=1}^J \left[ \int \text{Ber}(x_j | y=c, \theta_{jc}) P(\theta_{jc} | D) d\theta_{jc} \right] \\ &= \bar{\pi}_c \prod_{j=1}^J \bar{\theta}_{jc}^{I(x_j=1)} (1 - \bar{\theta}_{jc})^{I(x_j=0)} \end{aligned}$$

$$\text{with } \bar{\pi}_c = \frac{N_c + \alpha_c}{N + \alpha_0} \quad \bar{\theta}_{jc} = \frac{N_{jc} + \beta_j}{N_c + \beta_0 + \beta_j}$$

$$\alpha_0 = \sum_c \alpha_c$$

Note that  $\bar{\theta}_{jc}$  is being used for evaluating the class rather than  $(\hat{\theta}_{jc})_{\text{MAP}}$

One can use ~~the~~ multinomial model using word frequency, but it could be problematic if certain words are 'bursty': They appear many times when they do in a document, but does not at all appear in most documents.

Back to general issues with Bayesian method. We should take advantage of the fact that Bayesian approach gives us a whole distribution of  $\theta$ .

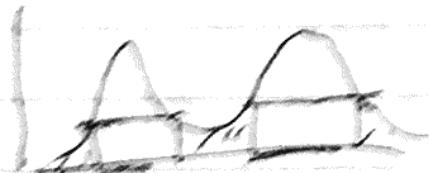
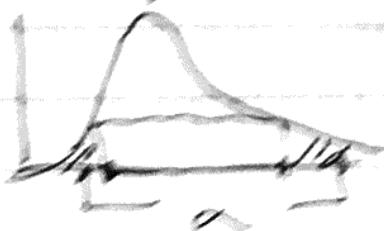
We could do confidence intervals  $\rightarrow$  Credible interval



If the dist<sup>n</sup> is skewed



Highest posterior density region



If the posterior distribution is not analytically tractable, one can use Monte Carlo to sample from it and get values.

Another important aspect of Bayesian approach is its take of model selection.

## Bayesian model Selection

Remember the polynomial fit exercise where choosing different degrees led to models with different complexities. How does Bayesian approach choose among them?

As opposed to the cross validation approach, which requires doing many fits, we compute, for a model  $m$

$$P(m | \mathcal{D}) = \frac{P(\mathcal{D} | m) p(m)}{\sum_{m'} P(\mathcal{D} | m') p(m')}$$

Biggest  $P(m | \mathcal{D})$  wins.

For uniform prior ( $p(m)$  constant) we just take largest  $P(\mathcal{D} | m)$

$$P(\mathcal{D} | m) = \int P(\mathcal{D} | \theta) P(\theta | m) d\theta$$

integrate over  
all choices of  
parameters

Evidence,

Marginal likelihood,  
Integrated likelihood

## Bayesian Occam's razor

Occam's razor: One should pick the simplest model that adequately explains the data.

Use  $P(D|\hat{\theta}_m)$  to choose between models

Too complex models

$$P(D) = P(y_1) P(y_2|y_1) P(y_3|y_1, y_2) \cdots P(y_n|y_1, \dots, y_{n-1})$$

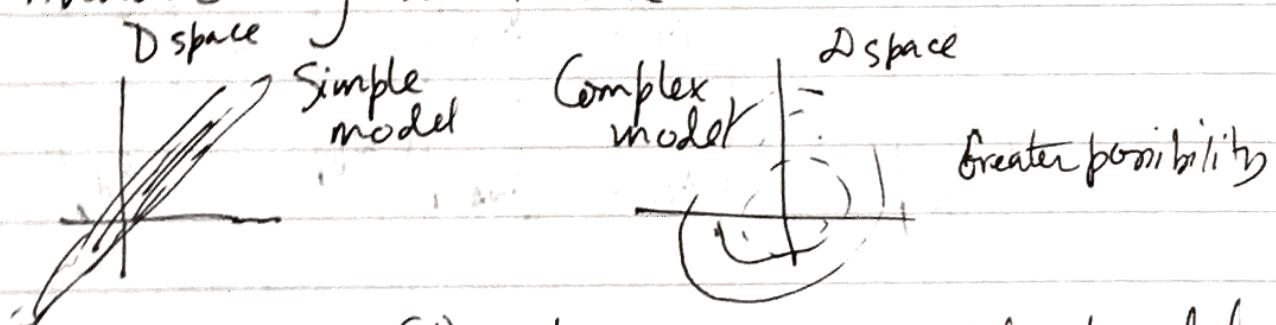
good likelihood by overfitting initial data

Bad likelihood from later pieces.

Like training  
good

validation  
bad.

Another way is to see it is that



$$\text{Since } \sum_{D'} P(D') = 1$$

more spread out model usually produces smaller  $P(D)$

Example  $x_i \sim N(\mu, \sigma^2)$

$$P(\mu) \propto e^{-\frac{\mu^2}{2\sigma_m^2}}$$

Model 1  $\mu_1 = \mu_2$

Model 2  $\mu_1, \mu_2$  independent