

# CS 340 Lec. 14: Bayesian Statistics

AD

February 2011

# A Bayesian Approach

- In a Bayesian approach, the unknown parameter  $\theta$  is assumed random with an associated prior distribution  $p(\theta)$ .
- Given data  $\{\mathbf{x}^i\}_{i=1}^N$  distributed according to  $p(\mathbf{x}_{1:N}|\theta)$ , inference about  $\theta$  is based on the posterior distribution

$$p(\theta|\mathbf{x}_{1:N}) = \frac{p(\mathbf{x}_{1:N}|\theta)p(\theta)}{p(\mathbf{x}_{1:N})}.$$

- From this posterior, we can obtain various point estimates of  $\theta$ .

# Bernoulli and Binomial Models

- Assume independent  $\{x^i\}$  where  $x^i \in \{0, 1\}$  ( $= \{Tail, Head\}$ ) with

$$p(x|\theta) = \theta^{\mathbb{I}(x=1)} (1-\theta)^{\mathbb{I}(x=0)}$$

so

$$p(x_{1:N}|\theta) = \theta^{n_1} (1-\theta)^{N-n_1}$$

where  $n_1 = \sum_{i=1}^n \mathbb{I}(y^i = 1)$  and  $\hat{\theta}_{MLE} = n_1/N$ .

- $n_1$  is the number of “success” among  $N$  trials, it follows a Binomial distribution

$$p(n_1|\theta) = Bin(n_1; \theta, N) = \binom{N}{n_1} \theta^{n_1} (1-\theta)^{N-n_1}$$

- In a Bayesian framework, we set a prior density  $p(\theta)$  on  $\theta \in [0, 1]$ .
- If you know nothing about  $\theta$  a reasonable prior is the uniform density

$$p(\theta) = \mathbf{1}_{[0,1]}(\theta).$$

# Conjugate Priors

- For simplicity, we will mostly focus on a special kind of prior which has nice mathematical properties.
- A prior  $p(\theta)$  is said to be conjugate to a likelihood  $p(x_{1:N}|\theta)$  (equivalently  $p(n_1|\theta)$ ) if the corresponding posterior  $p(\theta|x_{1:N}) = p(\theta|n_1)$  has the same functional form as  $p(\theta)$ .
- This means the prior family is closed under Bayesian updating.
- So we can recursively apply the rule to update our beliefs as data streams in (online learning).

- Let us introduce the class of Beta densities defined for  $\alpha, \beta > 0$

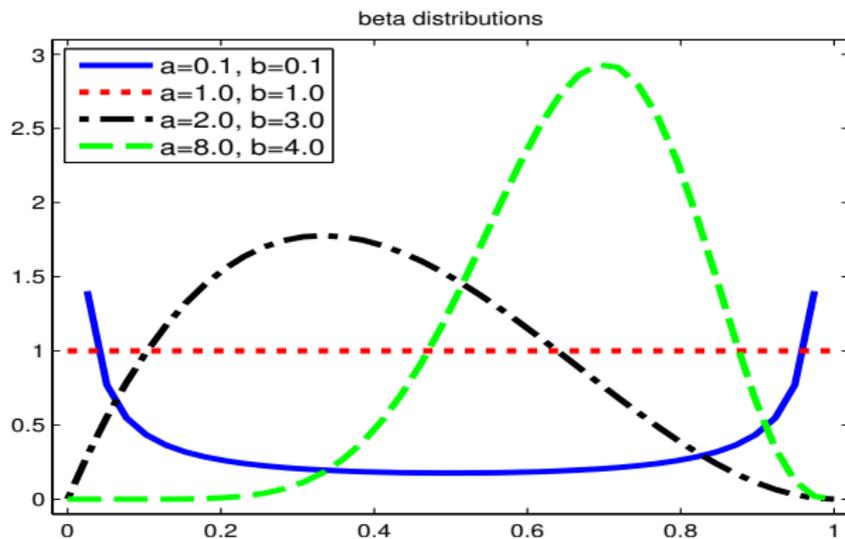
$$\text{Beta}(\theta; \alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1 - \theta)^{\beta-1} \mathbf{1}_{[0,1]}(\theta)$$

where  $\Gamma(u) = \int_0^\infty t^{u-1} e^{-t} dt$ . Note that  $\Gamma(u) = (u-1)!$  for  $u \in \mathbb{N}$ .

- Be careful:  $(\alpha, \beta)$  are *fixed* quantities. To distinguish them from  $\theta$ , we call them *hyperparameters*. For  $\alpha = \beta = 1$ , the Beta density corresponds to the uniform density.
- The Beta prior is such that

$$\mathbb{E}(\theta) = \frac{\alpha}{\alpha + \beta}, \quad \mathbb{V}(\theta) = \frac{\alpha\beta}{(\alpha + \beta)^2 (\alpha + \beta + 1)}.$$

# Beta Prior



- We obtain

$$\begin{aligned} p(\theta | n_1) &= \frac{p(n_1 | \theta) p(\theta)}{p(n_1)} \\ &\propto p(n_1 | \theta) p(\theta) \\ &\propto \theta^{n_1} (1 - \theta)^{N - n_1} \theta^{\alpha - 1} (1 - \theta)^{\beta - 1} \mathbf{1}_{[0,1]}(\theta) \\ &= \theta^{n_1 + \alpha - 1} (1 - \theta)^{N - n_1 + \beta - 1} \mathbf{1}_{[0,1]}(\theta) \end{aligned}$$

- This implies necessarily that  $p(\theta | x_{1:N}) = \text{Beta}(\theta; n_1 + \alpha, N - n_1 + \beta)$ .
- The prior on  $\theta$  can be conveniently reinterpreted as an imaginary initial sample of size  $(\alpha + \beta - 2)$  with  $\alpha - 1$  observations “1” and  $\beta - 1$  observations “0”. Provided that  $(\alpha + \beta - 2)$  is small with respect to  $n$ , the information carried by the data is prominent.

# Sequential Bayesian Inference with the Binomial-Beta Model

- Assume we first observe at 'time  $t$ '  $n_1^t$  '1' among  $N^t$  trials where  $t = 1, 2, \dots$
- We have

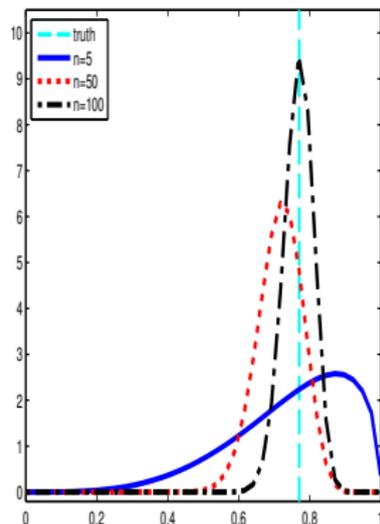
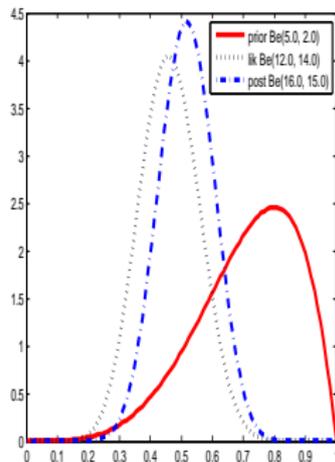
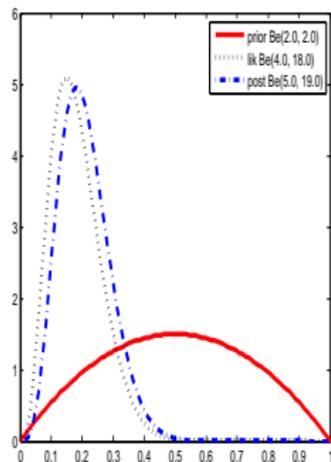
$$p(\theta | n_1^1) = \text{Beta}(\theta; n_1^1 + \alpha, N^1 - n_1^1 + \beta)$$

- At time  $t > 1$ , we use

$$\begin{aligned} p(\theta | n_1^1, \dots, n_1^k) &\propto p(n_1^k | \theta) p(\theta | n_1^1, \dots, n_1^{k-1}) \\ &= \text{Beta}\left(\theta; \alpha + \sum_{i=1}^k n_1^i, \beta + \sum_{i=1}^k (N^i - n_1^i)\right); \end{aligned}$$

i.e. the posterior at time  $k$  can be computed using as a prior the posterior at time  $k - 1$  and the likelihood of the observations at time  $k$ .

# Bayesian Inference with the Binomial-Beta Model



(left) Updating a  $\text{Beta}(2,2)$  prior with a Binomial likelihood with  $n_1 = 3$ ,  $n_0 = 17$  to yield a  $\text{Beta}(5,19)$ ; (center) Updating a  $\text{Beta}(5,2)$  prior with a Binomial likelihood with  $n_1 = 11$ ,  $n_0 = 13$  to yield a  $\text{Beta}(16,15)$  posterior. (c) Sequentially updating a Beta distribution starting with a  $\text{Beta}(1,1)$  and converge to a delta function centered on the true value.

# Bayesian Inference with the Binomial-Beta Model

- We have

$$\mathbb{E}(\theta | n_1) = \frac{n_1 + \alpha}{n_1 + \alpha + n_0 + \beta} = \frac{n_1 + \alpha}{N + \alpha + \beta}$$

- The posterior means behave asymptotically like  $n_1/n$  (the 'frequentist' estimator) and converge to  $\theta^*$ , the 'true' value of  $\theta^*$ .
- We have

$$\begin{aligned} \mathbb{V}(\theta | n_1) &= \frac{(n_1 + \alpha)(n_0 + \beta)}{(n_1 + \alpha + n_0 + \beta)^2 (n_1 + \alpha + n_0 + \beta + 1)} \\ &\approx \frac{\hat{\theta}_{MLE} (1 - \hat{\theta}_{MLE})}{N} \text{ for large } N \end{aligned}$$

- The posterior variance decreases to zero as  $n \rightarrow \infty$ , at rate  $n^{-1}$ : the information you get on  $\theta$  gets more and more precise.
- For  $n$  large enough, the prior is washed out by the data. For a small  $n$ , the prior can have a huge impact.

# Bayesian Inference with the Bernoulli-Beta Model

- We can compute things like

$$\Pr(\theta \in [0.3, 0.7] | n_1) = \int_{0.3}^{0.7} p(\theta | n_1) d\theta$$

- **Be careful:** This has absolutely nothing to do with confidence intervals.
- In classical statistics, and for an univariate problem, the confidence interval at level  $\alpha$  is of the form  $[\hat{\theta} - z_{\alpha/2}\hat{\sigma}, \hat{\theta} + z_{\alpha/2}\hat{\sigma}]$  where  $\hat{\theta}$  is the classical estimator (say MLE) and  $\hat{\sigma}$  is an estimate of its standard deviation.
- In this frequentist perspective, the true value of the parameter is fixed, and the confidence interval is random, having a probability of  $(1 - \alpha)$  to actually contain this true value (when we repeat the same experiment a great number of times) and it is not possible to interpret  $(1 - \alpha)$  as the probability that the parameter lies in the confidence interval for the considered experiment.

- We can also find the maximum a posterior (MAP)

$$\begin{aligned}\hat{\theta}_{MAP} &= \arg \max p(\theta | n_1) \\ &= \arg \max \log p(\theta | n_1) \\ &= \arg \max \log p(n_1 | \theta) + \log p(\theta) \\ &= \frac{n_1 + \alpha - 1}{n_1 + \alpha - 1 + n_0 + \beta - 1} = \frac{n_1 + \alpha - 1}{N + \alpha + \beta - 2}.\end{aligned}$$

- $\hat{\theta}_{MAP} = \hat{\theta}_{MLE}$  when  $\alpha = \beta = 1$  as then  $\log p(\theta)$  is constant over  $[0, 1]$ .

# Prediction: Classical vs Bayesian Approaches

- Assume you have observed  $n_1$  successes among  $N$  trials, we want to use these data to come up with the distribution of the outcome of the next trial.
- Using a Maximum Likelihood approach, we would use the plug-in prediction

$$p(x = 1 | \hat{\theta}_{MLE}) = \hat{\theta}_{MLE} = \frac{n_1}{N}$$

This does not account whatsoever for the uncertainty about  $\hat{\theta}_{MLE}$  (and suffer from Black Swan problem)

- In a Bayesian approach, we will use the predictive distribution

$$\begin{aligned} p(x = 1 | n_1) &= \int p(x = 1 | \theta) p(\theta | n_1) d\theta \\ &= \int \theta p(\theta | n_1) d\theta = \frac{n_1 + \alpha}{N + \alpha + \beta} \end{aligned}$$

so even if  $n_1 = 0$  then  $p(x = 1 | x_{1:N}) > 0$  and our prediction takes into account the uncertainty about  $\theta$ .

## Prediction: Classical vs Bayesian Approaches

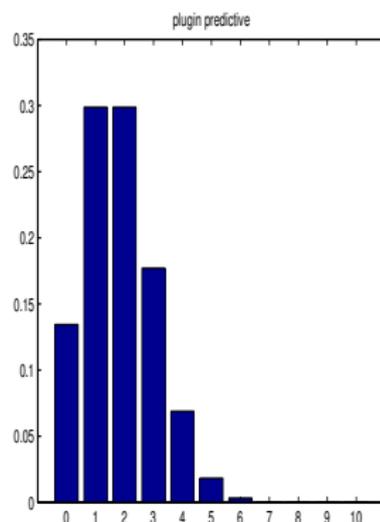
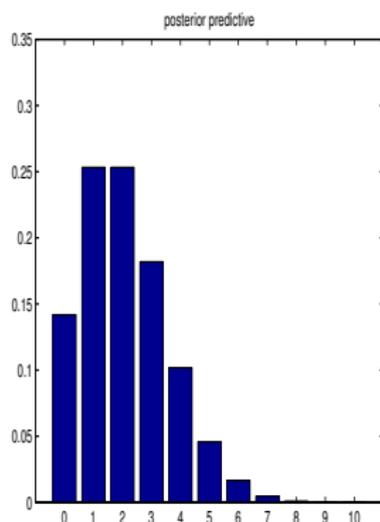
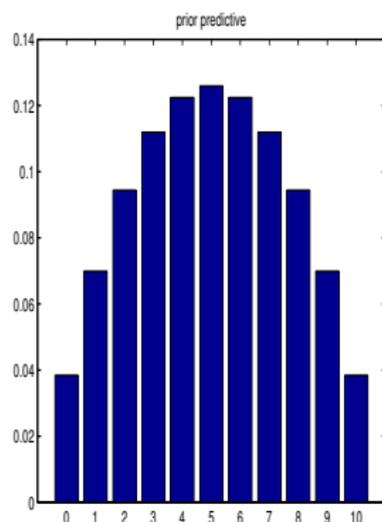
- Suppose now we want to predict the number  $m_1$  of heads in  $M$  future trials.
- The standard MLE approach would give use

$$p(m_1 | \hat{\theta}_{MLE}) = \text{Bin}(m_1; \hat{\theta}_{MLE}, M) = \binom{M}{m_1} \hat{\theta}_{MLE}^{m_1} (1 - \hat{\theta}_{MLE})^{M-m_1}$$

- The Bayesian approach yields

$$\begin{aligned} p(m_1 | n_1) &= \int p(m_1 | \theta) p(\theta | n_1) d\theta \\ &= \binom{M}{m_1} \frac{\Gamma(N+\alpha+\beta)}{\Gamma(n_1+\alpha)\Gamma(N-n_1+\beta)} \int \theta^{m_1+n_1-1} (1-\theta)^{N+M-m_1-n_1-1} d\theta \\ &= \binom{M}{m_1} \frac{\Gamma(N+\alpha+\beta)}{\Gamma(n_1+\alpha)\Gamma(N-n_1+\beta)} \frac{\Gamma(m_1+n_1+\alpha)\Gamma(N+M-m_1-n_1+\beta)}{\Gamma(N+M+\alpha+\beta)} \end{aligned}$$

# Prediction: Classical vs Bayesian Approaches



(left) Prior predictive dist. for a Binomial likelihood with  $M = 10$  and a  $\text{Beta}(2,2)$  prior. (center) Posterior predictive after having seen  $n_1 = 3$ ,  $N = 20$ . (right) Plug-in approximation using  $\hat{\theta}_{MLE}$ .

# From Coins to Dice: Multinomial

- Assume you have independent observations  $\{\mathbf{x}^i\}_{i=1}^M$  such that

$$p(\mathbf{x}|\theta) = \frac{P!}{\prod_{k=1}^d x_k!} \prod_{k=1}^d \theta_k^{x_k}$$

for  $\theta_k > 0$ ,  $\sum_{k=1}^d \theta_k = 1$  and  $x_k = 0, 1, 2, \dots, P$  with  $\sum_k x_k = P$ .

- We have seen that

$$\hat{\theta}_{k,MLE} = \frac{\sum_{i=1}^M x_k^i}{\sum_{i=1}^M \sum_{k=1}^d x_k^i} = \frac{N_k}{N}$$

- We want now to perform a Bayesian analysis

$$p(\theta|\mathbf{x}^{1:M}) = \frac{p(\mathbf{x}^{1:M}|\theta) p(\theta)}{p(\mathbf{x}^{1:M})}$$

- The Dirichlet density is given by

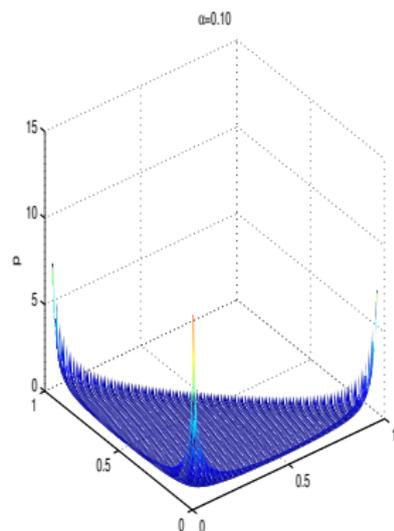
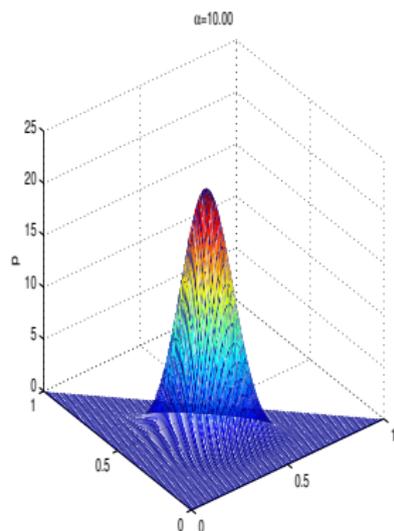
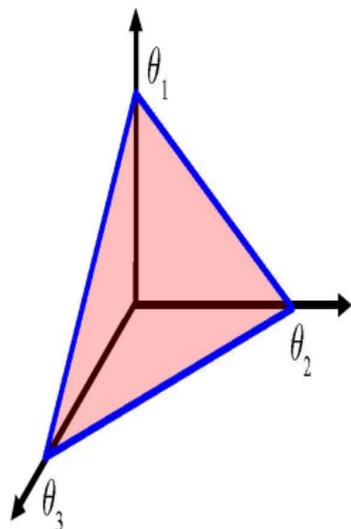
$$\text{Dir}(\theta; (\alpha_1, \dots, \alpha_d)) = \frac{\Gamma\left(\sum_{k=1}^d \alpha_k\right)}{\prod_{k=1}^d \Gamma(\alpha_k)} \prod_{k=1}^d \theta_k^{\alpha_k - 1}$$

for  $\alpha_k > 0$  and corresponds to a Beta density for  $d = 2$ . It is defined on  $\left\{ \theta : \theta_k > 0 \text{ and } \sum_{k=1}^d \theta_k = 1 \right\}$ .

- $\alpha_0 = \sum_{k=1}^d \alpha_k$  controls how peaky the distribution is and the  $\alpha_k$  controls where the peak is located.
- We have

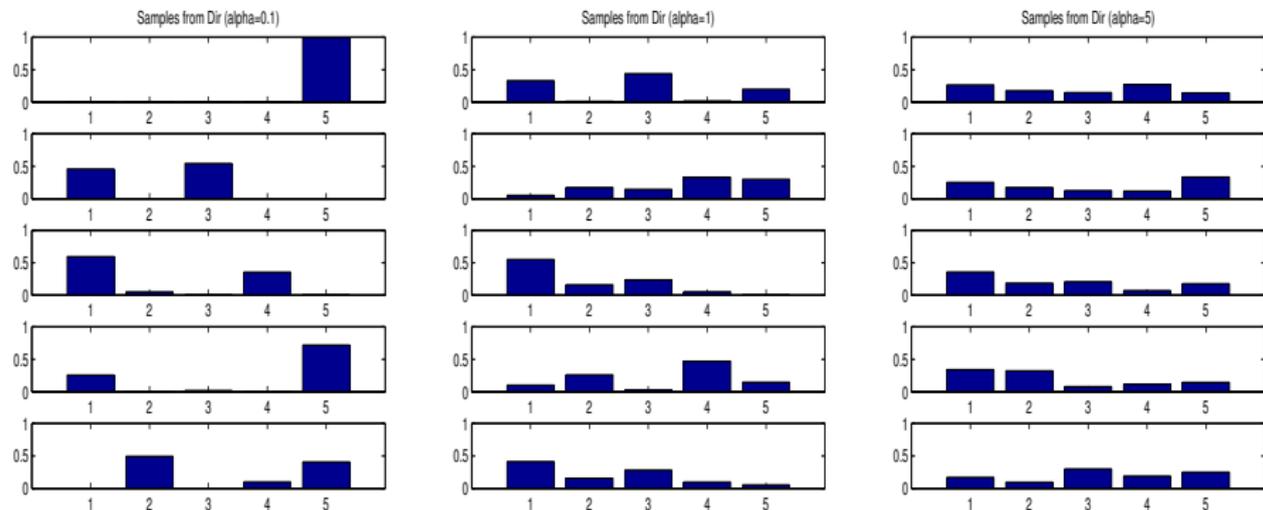
$$\mathbb{E}(\theta_k) = \frac{\alpha_k}{\alpha_0}, \quad \text{mode}(\theta_k) = \frac{\alpha_k - 1}{\alpha_0 - d}, \quad \mathbb{V}(\theta_k) = \frac{\alpha_k (\alpha_0 - \alpha_k)}{\alpha_0^2 (\alpha_0 + 1)}$$

# Dirichlet Prior



(left) Support of the Dirichlet density for  $d = 3$  (center) Dirichlet density for  $\alpha_k = 10$  (right) Dirichlet density for  $\alpha_k = 0.1$ .

# Dirichlet Prior



Samples from a Dirichlet distribution for  $d = 5$  when  $\alpha_k = \alpha_l$  for  $k \neq l$ .

- We obtain

$$\begin{aligned} p(\theta | \mathbf{x}^{1:M}) &= \frac{p(\mathbf{x}^{1:M} | \theta) p(\theta)}{p(\mathbf{x}^{1:M})} \\ &\propto \prod_{k=1}^d \theta_k^{N_k} \prod_{k=1}^d \theta_k^{\alpha_k - 1} \\ &\propto \prod_{k=1}^d \theta_k^{\alpha_k + N_k - 1} \end{aligned}$$

- This implies necessarily that

$$p(\theta | \mathbf{x}_{1:M}) = \text{Dir}(\theta; \alpha_1 + N_1, \dots, \alpha_d + N_d).$$

# Predictive Distribution with the Multinomial-Dirichlet Model

- We have for a single categorical variable

$$\begin{aligned}\Pr(x = k | \mathbf{x}^{1:M}) &= \int \Pr(x = k | \theta) p(\theta | \mathbf{x}^{1:M}) d\theta \\ &= \int \theta_k p(\theta | \mathbf{x}^{1:M}) d\theta \\ &= \int \theta_k p(\theta_k | \mathbf{x}^{1:M}) d\theta_k \\ &= \frac{\alpha_k + N_k}{\alpha_0 + N}.\end{aligned}$$

- Once more this avoids the black-swan problem.

# Bayesian Naive Bayes for Multinomial Data

- We assume that we have  $M$  data  $(\mathbf{x}^i, y^i) \in \mathbb{N}^d \times \{0, 1\}^C$  and we use the model

$$p(\mathbf{x}, y = c | \theta) = \pi_c \frac{P!}{\prod_{i=1}^d x_k!} \prod_{k=1}^d \theta_{k,c}^{x_k}$$

where  $(\pi_1, \dots, \pi_C, \theta_{1,1}, \dots, \theta_{d,1}, \dots, \theta_{1,C}, \dots, \theta_{d,C})$  are the unknown parameters.

- If we do MLE, then

$$\hat{\pi}_{c,MLE} = \frac{M_c}{M}, \quad \hat{\theta}_{k,c,MLE} = \frac{N_{k,c}}{N_c}$$

where  $M_c$  = nb. documents class  $c$ ,  $N_{k,c}$  = nb. occurrences word  $k$  in class  $c$ ,  $M = \sum_{k=1}^C M_c$ ,  $N_c = \sum_{k=1}^d N_{k,c}$ .

# Bayesian Naive Bayes for Multinomial Data

- In a Bayesian context, we can set independent Dirichlet priors

$$\begin{aligned}p(\boldsymbol{\pi}) &= \text{Dir}((\pi_1, \dots, \pi_C); \beta_1, \dots, \beta_C), \\p(\boldsymbol{\theta}_c) &= \text{Dir}((\theta_{1,c}, \dots, \theta_{d,c}); \alpha_{1,c}, \dots, \alpha_{d,c}), \quad c = 1, \dots, C\end{aligned}$$

and obtain

$$\begin{aligned}p(\pi_1, \dots, \pi_C | D) &= \text{Dir}((\pi_1, \dots, \pi_C); \beta_1 + M_1, \dots, \beta_C + M_C), \\p(\boldsymbol{\theta}_c | D) &= \text{Dir}((\theta_{1,c}, \dots, \theta_{d,c}); \alpha_{1,c} + N_{1,c}, \dots, \alpha_{d,c} + N_{d,c}).\end{aligned}$$

- From this posterior, you can compute  $\hat{\pi}_{MAP}, \hat{\boldsymbol{\theta}}_{c,MAP}$  or  $\hat{\pi}_{MMSE} = \mathbb{E}(\pi_M | D), \hat{\boldsymbol{\theta}}_{c,MMSE} = \mathbb{E}(\boldsymbol{\theta}_c | D)$  and use

$$p(y = c | \mathbf{x}, \hat{\boldsymbol{\pi}}, \hat{\boldsymbol{\theta}}) \propto p(y = c | \hat{\boldsymbol{\pi}}) p(\mathbf{x} | y = c, \hat{\boldsymbol{\theta}})$$

# Bayesian Naive Bayes for Multinomial Data

- A better way to do it is to
- Given a new input  $\mathbf{x}$ , we compute using

$$p(y = c | \mathbf{x}, D) \propto p(y = c | D) p(\mathbf{x} | y = c, D)$$

where

$$p(y = c | D) = \int \underbrace{p(y = c | D, \pi)}_{=\pi_c} p(\pi | D) d\pi = \frac{\beta_c + M_c + 1}{\beta_0 + M + 1},$$

$$p(\mathbf{x} | y = c, D) = \int p(\mathbf{x} | D, \theta_c) p(\theta_c | D) d\theta_c = \dots$$

- Assume you have independent data  $\{x^i\}_{i=1}^N$  such that

$$p(x|\theta) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

where  $\theta = (\mu, \sigma^2)$ .

- We have seen that

$$\hat{\mu}_{ML} = \frac{1}{N} \sum_{i=1}^N x^i, \quad \hat{\sigma}^2_{ML} = \frac{1}{N} \sum_{i=1}^N (x^i - \hat{\mu}_{ML})^2.$$

# Bayesian Inference for Normal Data

- In a Bayesian framework, the conjugate prior is  $p(\mu, \sigma^2) = p(\mu | \sigma^2) p(\sigma^2)$  where

$$p(\mu | \sigma^2) = \mathcal{N}\left(\mu; \mu_0, \frac{\sigma^2}{\kappa_0}\right) = \frac{1}{\sqrt{2\pi(\sigma^2/\kappa_0)}} \exp\left(-\frac{\kappa_0(\mu - \mu_0)^2}{2\sigma^2}\right)$$
$$p(\sigma^2) = \mathcal{IG}(\sigma^2; \alpha, \beta) = \frac{\beta^\alpha}{\Gamma(\alpha)} (\sigma^2)^{-\alpha-1} \exp(-\beta/\sigma^2) \mathbf{1}_{(0, \infty)}(\sigma^2).$$

- The posterior is given by  $p(\theta | x^{1:n}) = p(\sigma^2 | x^{1:n}) p(\mu | x^{1:n}, \sigma^2)$  where

$$p(\mu | x^{1:n}, \sigma^2) = \mathcal{N}\left(\mu; \frac{\kappa_0 \mu_0 + N \hat{\mu}_{ML}}{\kappa_0 + N}, \frac{\sigma^2}{\kappa_0 + N}\right)$$
$$p(\sigma^2 | x^{1:n}) = \mathcal{IG}\left(\sigma^2; \alpha + N/2, \beta + \frac{N}{2} \hat{\sigma}_{ML}^2 + \frac{N \kappa_0}{2(N + \kappa_0)} (\hat{\mu}_{ML} - \mu_0)^2\right)$$

- Once more we see clearly the influence of the prior on the posterior and, as  $N \rightarrow \infty$ , the posterior concentrates around  $\hat{\mu}_{ML}$  and  $\hat{\sigma}_{ML}^2$ .

# Bayesian Model Selection

- Suppose we have  $K$  different models for the data  $D$ ; each model being associated to some parameters  $\theta_i$ .
- Using a Bayesian approach, we can compute

$$p(M = i | D) = \frac{p(M = i) p(D | M = i)}{P(D)}$$

where

$$p(D) = \sum_{i=1}^K p(M = i) p(D | M = i)$$

- The *marginal likelihood* or *evidence*  $p(D | M = i)$  is given by

$$p(D | M = i) = \int p(D | \theta_i) p(\theta_i) d\theta_i$$

which is the normalizing constant of

$$p(\theta_i | D) = \frac{p(\theta_i) p(D | \theta_i)}{p(D)}.$$

- To compare two models, we can use posterior odds of Bayes factors

$$\underbrace{\frac{p(M = i | D)}{p(M = j | D)}}_{\text{posterior odds}} = \underbrace{\frac{p(D | M = i)}{p(D | M = j)}}_{\text{Bayes factor}} \underbrace{\frac{p(M = i)}{p(M = j)}}_{\text{prior odds}}$$

- The Bayes factor is a Bayesian version of a likelihood ratio test, that can be used to compare models of different complexity.
- Bayes factors and posterior odds tell you whether one should prefer  $M = i$  to  $M = j$ : it does NOT tell you whether these models are sensible!

## Example: Is the Euro coin biased?

- Suppose we toss a coin  $N = 250$  times and observe  $n_1 = 141$  heads and  $n_0 = 109$  tails:

$$p(x^{1:N} | \theta) = \theta^{n_1} (1 - \theta)^{n_0}$$

- Consider two models/hypotheses:  $M_1 =$  coin unbiased, that is  $\theta_1 = 0.5$  and  $M_2 =$  coin biased and  $p(\theta_2) = \text{Beta}(\theta_2; \alpha_1, \alpha_0)$ .
- We have

$$p(D|M_1) = 0.5^{n_1} (1 - 0.5)^{n_0} = 0.5^N$$

and

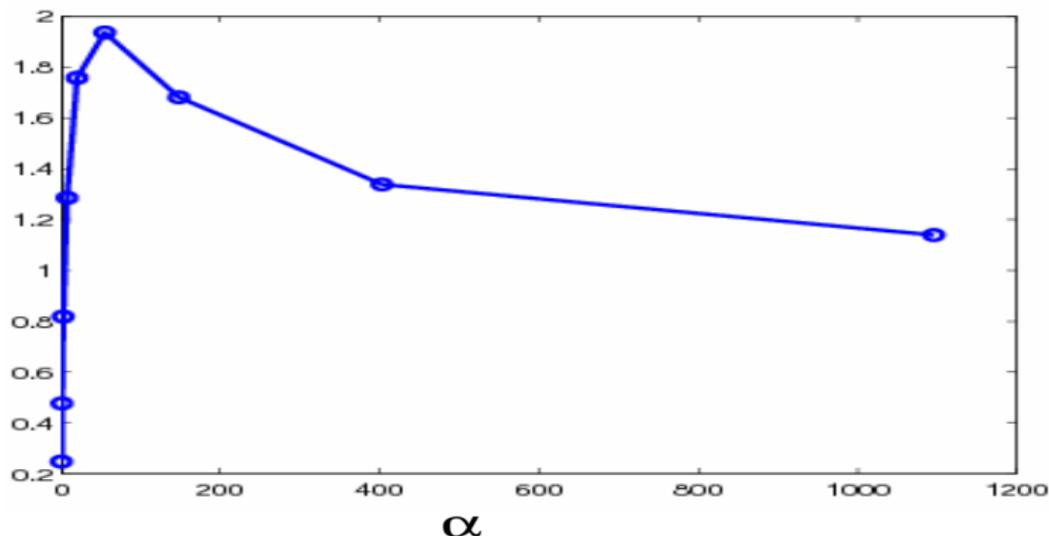
$$p(D|M_2) = \frac{\Gamma(\alpha_0 + n_0) \Gamma(\alpha_1 + n_1)}{\Gamma(\alpha_0 + \alpha_1 + N)} \frac{\Gamma(\alpha_0 + \alpha_1)}{\Gamma(\alpha_0) \Gamma(\alpha_1)}$$

so

$$\frac{p(D|M_2)}{p(D|M_1)} = \frac{\Gamma(\alpha_0 + n_0) \Gamma(\alpha_1 + n_1)}{\Gamma(\alpha_0 + \alpha_1 + N)} \frac{\Gamma(\alpha_0 + \alpha_1)}{\Gamma(\alpha_0) \Gamma(\alpha_1)} 0.5^{-N}$$

# Computation of Bayes Factors

- Let  $\alpha = \alpha_0 = \alpha_1$  varying over 0 to 1000.



Bayes factor  $p(D|M_2) / p(D|M_1)$  as a function of  $\alpha$ .

- The largest BF in favor of  $M_2$  (biased coin) is only 2.0, which is very weak evidence of bias.

- For complex Bayesian models, we cannot compute the posterior and marginal likelihood analytically.
- In such cases, analytical (Laplace, variational) and Monte Carlo methods approximations are necessary.
- For example, a crude approximation of the marginal likelihood is provided by the Bayesian Information Criterion

$$\log p(D|M_i) = \log p\left(D|\theta_i^{MLE}\right) - \frac{d}{2} \log n$$

where  $n$  is the number of data and  $d$  is the dimension/number of free parameters