



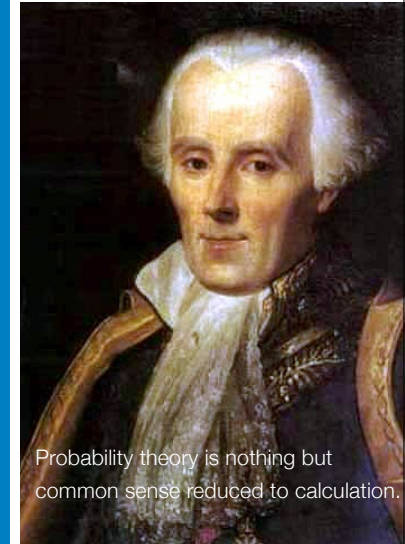
Royal Institute of
Technology

DD2447 STAT. METH. IN CS FALL 2014

- ★ Lecture 2 -
Probability, Bayesian
- ★ Chapter 2, 3

CHAPTER 2

- ★ Known concepts
- ★ Distributions
 - Beta, Gamma, Dirichlet
- ★ Sampling
- ★ Information theory



Probability theory is nothing but
common sense reduced to calculation.

BERNOULLI & BINOMIAL

$$\text{Ber}(x|\theta) = \begin{cases} \theta & \text{if } x = 1 \\ 1 - \theta & \text{if } x = 0 \end{cases} \quad \text{Bin}(k|n, \theta) = \binom{n}{k} \theta^k (1 - \theta)^{n-k}$$

- ★ One or several (unordered) coin tosses

CATEGORICAL & MULTINOMIAL

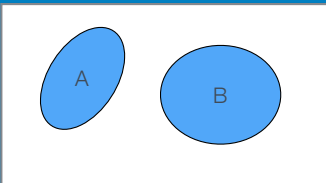
$$\text{Cat}(x|\theta) = \theta_x \quad \text{Mul}(x|n, \theta) = \binom{n}{x_1, \dots, x_K} \prod_{k=1}^K \theta_k^{x_k}$$

- ★ One or several (unordered) coin tosses

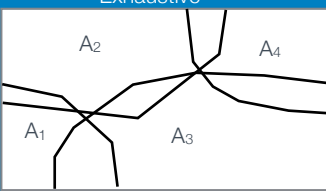
CONDITIONING

$$p(x, y) = p(y)p(x|y) \quad \text{or} \quad p(x|y) = \frac{p(x, y)}{p(y)}$$


Exclusive



Exhaustive



Exclusive & exhaustive



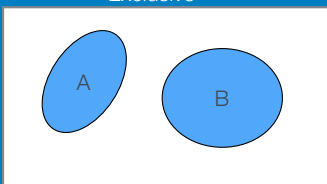
EXCLUSIVE & EXHAUSTIVE

- Exclusive

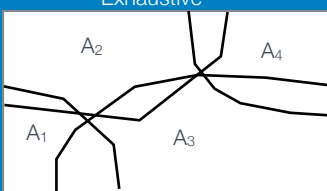
$$p(A \text{ or } B) = p(A) + p(B)$$
- Exclusive & exhaustive

$$\sum_i p(A_i) = 1$$

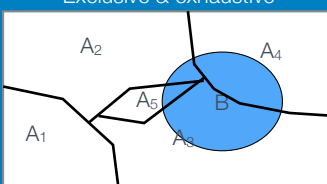
Exclusive



Exhaustive



Exclusive & exhaustive



EXCLUSIVE & EXHAUSTIVE

- Exclusive

$$p(A \text{ or } B) = p(A) + p(B)$$
- Exclusive & exhaustive

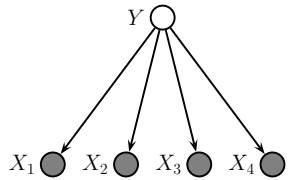
$$p(B) = \sum_i p(B, A_i) = \sum_i p(A_i)p(B|A_i)$$

CHAPTER 3 - BEYOND BAYES THEOREM, CONCEPTS



$$p(X|Y) = \frac{p(X, Y)}{p(Y)} = \frac{p(X)p(Y|X)}{\sum_x p(x)p(Y|x)}$$

NAIVE BAYES CLASSIFIER



$$p(\mathbf{x}, y) = p(y) \prod_{t=1}^4 p(x_t|y)$$

BAYES — GENERATIVE CLASSIFIER

- \mathbf{X} given Y the natural direction
- Classifier

$$p(Y = c|\mathbf{x}, \theta) = \frac{p(Y = c|\theta)p(\mathbf{x}|Y = c, \theta)}{\sum_{\mathbf{x}} p(Y = c|\theta)p(\mathbf{x}|Y = c, \theta)}$$

BREAST CANCER TEST

- X test – 1 positive
- Y breast cancer – 1 cancer

THE BRAIN

- ★ Biology
 - the more data the better
- ★ Philosophy of the mind
 - John Searl





- ★ $x_1, \dots, x_N \in [100]$
 - $[M] = \{1, \dots, M\}$
- ★ Ex.
 - 1,4,5,2,2,5,3,5,3,6,...
 - 14,8,28,2,36,...
 - 1,14,16,20,19,...

DATA

- ★ Data $D=14,8,28,2,36, \dots$
- ★ Dice $H = 20$ -sided with even numbers in $[40]$
 - Likelihood $p(D|H) = (1/20)^5$
- ★ Dice $H' = 6$ -sided with numbers 14,8,28,2,36,7
 - Likelihood $p(D|H') = (1/6)^5$
- ★ Least number of sides wins - Occam's razor

LIKELIHOOD

- ★ Data $D=14,8,28,2,36, \dots$
- ★ 6-sided dice H' with numbers 14,8,28,2,36,7
- ★ pretty unnatural
 - we give it prior probability $p(H') = 1/10^6$
- ★ Posterior probability (after observation)
- ★ In general

$$p(H'|D) = \frac{p(D|H')p(H')}{p(D)} = \frac{p(D|H')p(H')}{\sum_{H' \in \mathcal{H}'} p(D|H')}$$
- ★ Here

$$p(H'|D) = \frac{(1/6)^5 10^{-6}}{p(D)}$$

PRIORS

★ Data D=14,8,28,2,36,

- ★ 20-sided Dice H with even numbers in [40]
 - fairly natural
 - we give it prior probability $p(H) = 1/1000$

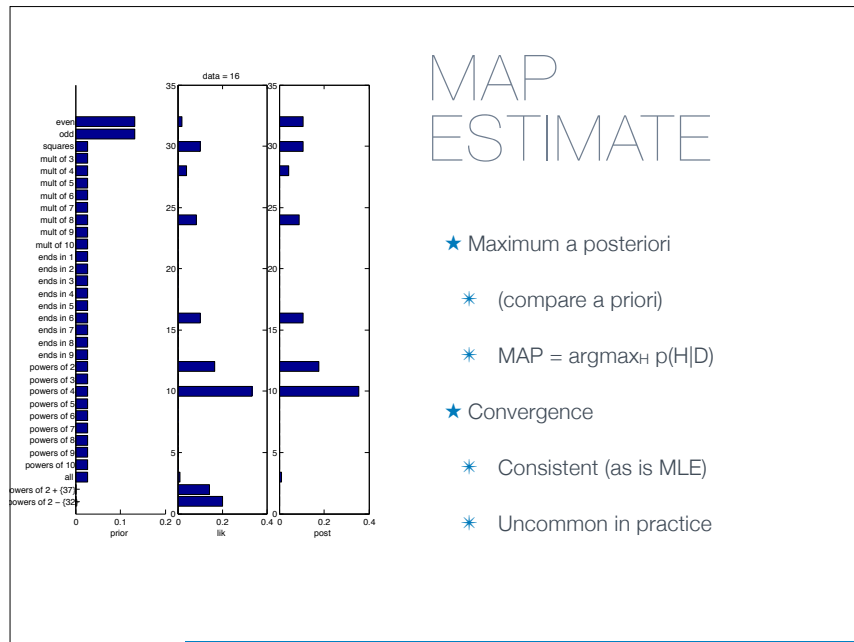
★ Posterior probability
$$p(H|D) = \frac{(1/20)^5 10^{-3}}{p(D)}$$

★ So
$$\frac{p(H|D)}{p(H'|D)} = \frac{(1/20)^5 10^{-3}}{(1/6)^5 10^{-6}} \approx 10^3 / 411 > 1$$

PRIORS

INFLUENCE OF PRIOR IN THIS MODEL

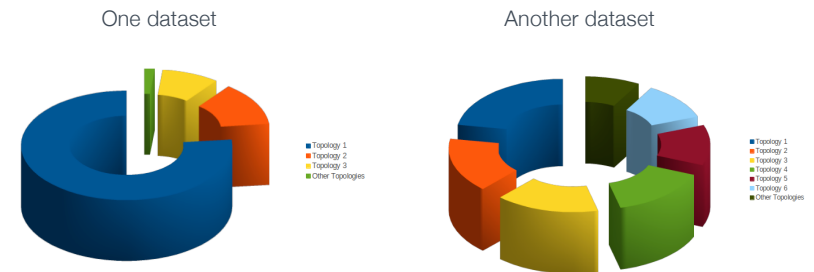
- ★ MAP = $\operatorname{argmax}_H p(H|D) = \operatorname{argmax} \log p(H|D) = \operatorname{argmax} \log p(D|H) + \log p(H)$
- ★ here $p(D|H)=(1/C)^N$ decreases exponentially in $N=\#\text{throws or data points}$
- ★ $p(H)$ is constant
- ★ Conclusions: as $N \rightarrow \infty$ $p(D|H)$ will dominate
- ★ Data overwhelms the prior



MAP ESTIMATE

- ★ Maximum a posteriori
 - * (compare a priori)
 - * $\operatorname{MAP} = \operatorname{argmax}_H p(H|D)$
- ★ Convergence
 - * Consistent (as is MLE)
 - * Uncommon in practice

POSTERIOR IS IMPORTANT (FOR THIS CHOICE)



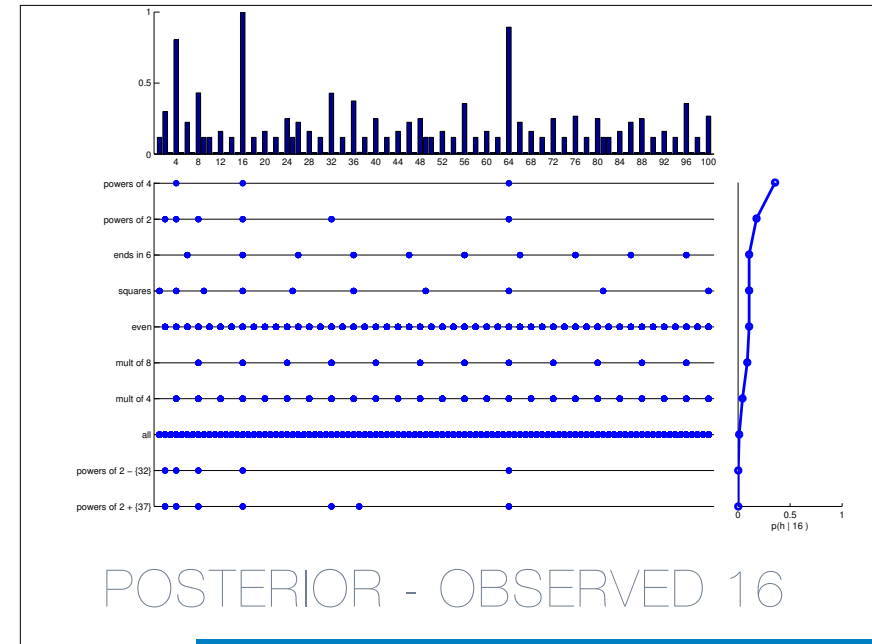
BAYESIAN MODEL AVERAGING

- ★ Posterior - what we believe
- ★ Using the MAP for prediction

$$p(x|D) = p(x|H_{\text{MAP}}^D)$$

- ★ Alternative model averaging

$$p(x|D) = \sum_{H \in \mathcal{H}} p(x, H|D) = \sum_{H \in \mathcal{H}} p(x|H)p(H|D)$$



BINOMIAL - IID BERNOULLI

- ★ Likelihood

- ★ N_1 heads and N_0 tails

- ★ $p(\text{head}) = \theta$

- ★ sequence $p(D) = \theta^{N_1}(1 - \theta)^{N_0}$

- ★ counts

$$p(D) = \binom{N_1 + N_0}{N_1} \theta^{N_1} (1 - \theta)^{N_0}$$

MLE BY FREQUENCIES

- ★ Likelihood is up to a constant $p(D) = \theta^{N_1}(1 - \theta)^{N_0}$

- ★ Log-likelihood $l(D) = N_1 \log \theta + N_0 \log(1 - \theta)$

- ★ Optimized by same θ

- ★ Derivation $l'(D) = \frac{N_1}{\theta} - \frac{N_0}{(1 - \theta)}$

- ★ Setting to zero $\frac{N_1}{\theta} = \frac{N_0}{(1 - \theta)}$

- ★ So

$$N_1 - \theta N_1 = \theta N_0 \quad \text{and} \quad \theta = \frac{N_1}{N_1 + N_0}$$

MLE FOR MULTINOMIAL AND AND CATEGORICAL

- ★ Likelihood $p(D) = \prod_{i \in [k]} \theta_i^{N_i}$

- ★ where $\sum_{i \in [k]} \theta_i = 1$

- ★ as well as loglikelihood $p(D) = \sum_{i \in [k]} N_i \log \theta_i$

- ★ is maximized by $\theta_i = \frac{N_i}{\sum_{i \in [k]} N_i}$

SUFFICIENT STATISTICS

- ★ Maximum Likelihood (ML) estimate maximize the probability of the data
- ★ here $\max_{\theta} p(N_1, N_0 | \theta)$
- ★ The pair N_1, N_0 is a sufficient statistic for our coin model
- ★ i.e., given those ML estimate follows

PRIOR FOR CATEGORICAL AND BERNOULLI - FIRST BERNOULLI

- ★ Assumption: we don't know θ

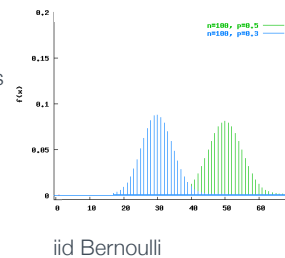
- ★ We need a prior over the outcome probabilities

- ★ N_1 heads and N_0 tails

- ★ p(head) denoted θ

- ★ sequence $p(D) = \theta^{N_1} (1 - \theta)^{N_0}$

- ★ counts $p(D) = \binom{N_1 + N_0}{N_1} \theta^{N_1} (1 - \theta)^{N_0}$ Binomial



BETA DISTRIBUTION

- PDF $\text{Beta}(x|a, b) = \frac{1}{B(a, b)} x^{(a-1)} x^{(b-1)}$

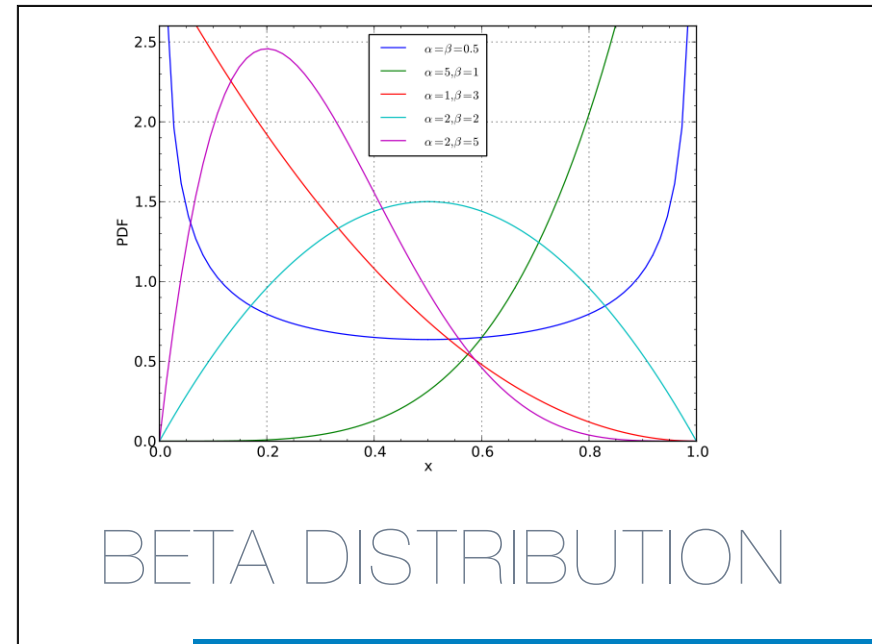
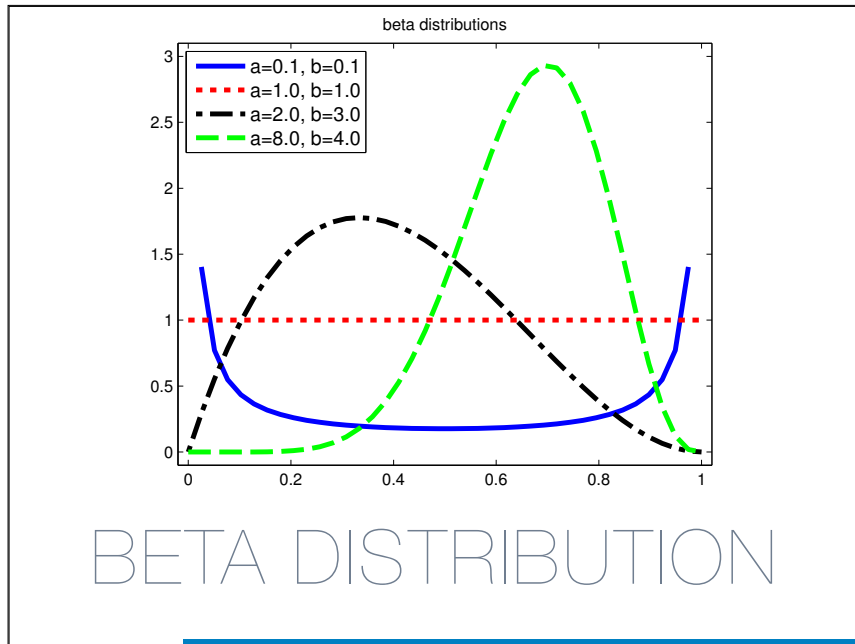
- where $B(a, b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$

- and $\Gamma(a) = (a-1)\Gamma(a-1)$

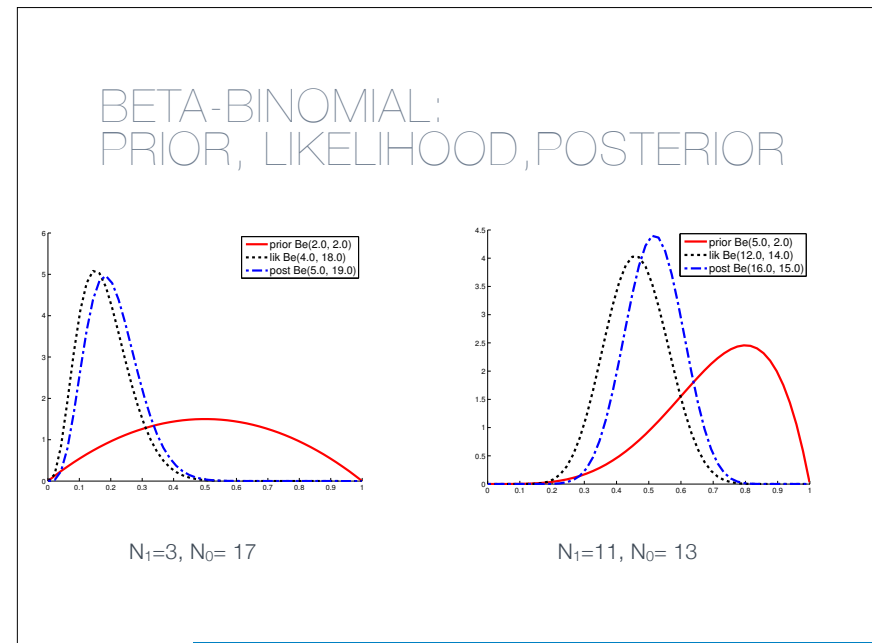
- In particular, for integer n

$$\Gamma(n) = (n-1)!$$

- Proper – the integral is 1.



- ## HYPER-PARAMETERS
- ★ Parameters γ_1 and γ_0 for prior called hyperparameters
 - ★ pseudocounts
 - ★ Prior's effective sample size is $\gamma_1 + \gamma_0$
 - ★ Prior that gives posterior of the same sort is called conjunctive



BETA BINOMIAL

- ★ Beta distribution up to a constant $p(\theta|\gamma_1, \gamma_0) \propto \theta^{\gamma_1-1}(1-\theta)^{\gamma_0-1}$
- ★ Posterior
$$\begin{aligned} p(\theta|D) &= p(D|\theta)p(\theta|\gamma_1, \gamma_0) \\ &\propto \theta^{N_1}(1-\theta)^{N_0}\theta^{\gamma_1-1}(1-\theta)^{\gamma_0-1} \\ &= \theta^{N_1+\gamma_1-1}(1-\theta)^{N_0+\gamma_0-1} \\ &\propto \text{Beta}(\theta|N_1 + \gamma_1, N_0 + \gamma_0) \end{aligned}$$
- ★ Beta is a conjunctive prior for Binomial

$$D = D_a \cup D_b$$

$$S(D_a) = (N_1^a, N_0^a) \text{ and } S(D_b) = (N_1^b, N_0^b)$$

$$\text{so } S(D) = (N_1^a + N_1^b, N_0^a + N_0^b) = (N_1, N_0)$$

$$\text{again } p(\theta|D) = \text{Beta}(\theta|N_1 + \gamma_1, N_0 + \gamma_0)$$

$$2 \text{ steps } p(\theta|D_a, D_b) \propto p(D_b|\theta)p(\theta|D_a) = ?$$

BETA-BINOMIAL: BATCH OR IN TWO STEPS

$$D = D_a \cup D_b$$

$$S(D_a) = (N_1^a, N_0^a) \text{ and } S(D_b) = (N_1^b, N_0^b)$$

$$\text{so } S(D) = (N_1^a + N_1^b, N_0^a + N_0^b) = (N_1, N_0)$$

$$\text{again } p(\theta|D) = \text{Beta}(\theta|N_1 + \gamma_1, N_0 + \gamma_0)$$

$$\begin{aligned} 2 \text{ steps } p(\theta|D_a, D_b) &\propto p(D_b|\theta)p(\theta|D_a) \\ &= \text{Bin}(N_1^b, N_0^b)\text{Beta}(\theta|N_1^a + \gamma_1, N_0^a + \gamma_0) \\ &= \text{Beta}(\theta|N_1^b + N_1^a + \gamma_1, N_0^b + N_0^a + \gamma_0) \\ &= \text{Beta}(\theta|N_1 + \gamma_1, N_0 + \gamma_0) \end{aligned}$$

BETA-BINOMIAL: BATCH OR IN TWO STEPS

BETA-BINOMIAL - MLE, MAP, AND PM

$$\text{Let } \gamma = \gamma_1 + \gamma_0 \text{ and } N = N_1 + N_0$$

Then

$$\theta_{\text{MLE}} = \frac{N_1}{N}, \quad \theta_{\text{MAP}} = \frac{N_1 + \gamma_1 - 1}{N + \gamma - 2}, \text{ and } \theta_{\text{PM}} = \frac{N_1 + \gamma_1}{N + \gamma}$$

PM \nearrow Posterior mean

$$\text{Let } f_1 = \gamma_1/\gamma$$

$$\text{Then } E[\theta|D] = \frac{\gamma f_1 + N_1}{N + \gamma} = \frac{\gamma}{N + \gamma} f_1 + \frac{N}{N + \gamma} \theta_{\text{MLE}}$$

Posterior to our observation what is the probability of a specific outcome?

$$\begin{aligned} p(x = 1|D) &= \int_0^1 p(x = 1|\theta)p(\theta|D)d\theta \\ &= \int_0^1 \theta \text{Beta}(\theta|N_1 + \gamma_1, N_0 + \gamma_0)d\theta \\ &= \frac{N_1 + \gamma_1}{N + \gamma} \end{aligned}$$

POSTERIOR
PREDICTIVE

Uniform prior, i.e., $\gamma_1 = \gamma_0 = 1$, gives

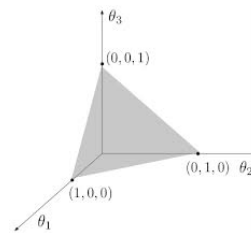
$$p(x = 1|D) = \frac{N_1 + 1}{N + 2}$$

Black Swan “paradox”

LAPLACE’S RULE OF
SUCCESSION



DIRICHLET



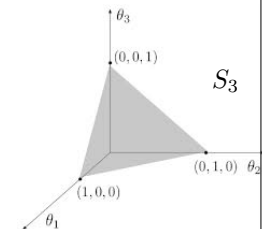
$$D = \{x_1, \dots, x_N\} \quad \text{where} \quad x_i \in \{1, \dots, K\}$$

$$\theta = (\theta_1, \dots, \theta_K) \in S_K \text{ the } K\text{-dim. probability simplex, i.e., pos. } \& \sum_{i \in [K]} \theta_k = 1$$

$$\alpha = (\alpha_1, \dots, \alpha_K) \text{ hyperparameters}$$

$$\text{Prior } Dir(\theta|\alpha) = \frac{1}{B(\alpha)} \prod_{k \in [K]} \theta_k^{\alpha_k - 1}$$

$$\text{Likelihood } p(\theta|D) \propto \prod_{k \in [K]} \theta_k^{N_k}$$



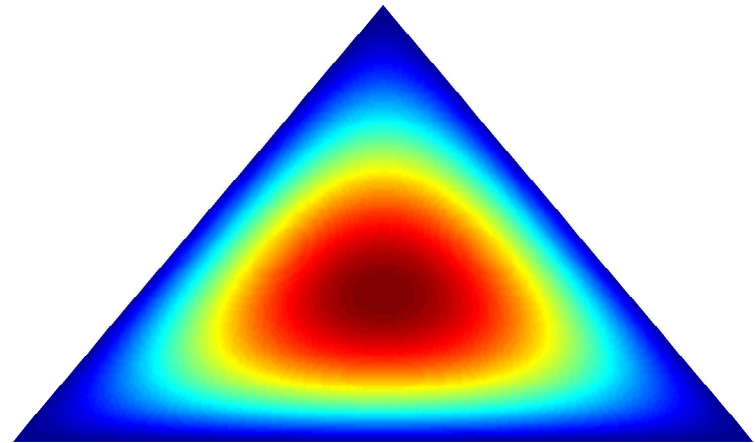
BACK TO THE DICE –
DIRICHLET-MULTINOMIAL

Prior $Dir(\boldsymbol{\theta}|\boldsymbol{\alpha}) = \frac{1}{B(\boldsymbol{\alpha})} \prod_{k \in [K]} \theta_k^{\alpha_k - 1}$

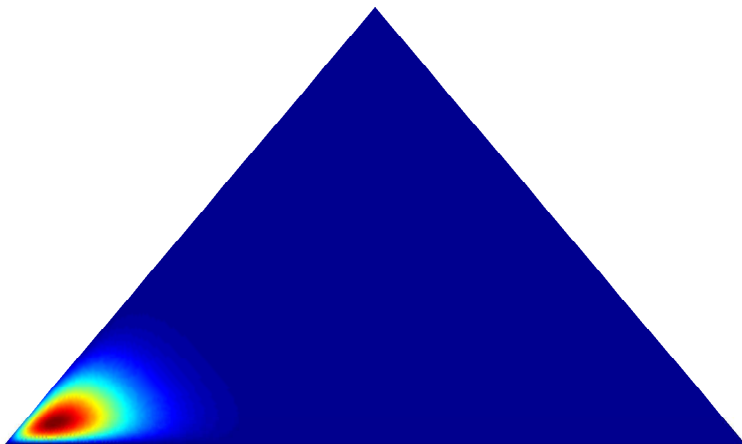
Likelihood $p(\boldsymbol{\theta}|D) \propto \prod_{k \in [K]} \theta_k^{N_k}$

Posterior $p(\boldsymbol{\theta}|D) = p(D|\boldsymbol{\theta})p(\boldsymbol{\theta})$
 $\propto \prod_{k \in [K]} \theta_k^{N_k} \prod_{k \in [K]} \theta_k^{\alpha_k - 1}$
 $= \prod_{k \in [K]} \theta_k^{N_k + \alpha_k - 1}$
 $\propto Dir(\boldsymbol{\theta}|N_1 + \alpha_1, \dots, N_K + \alpha_K)$

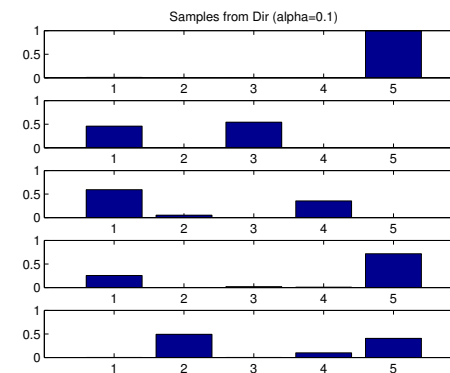
POSTERIOR FOR DIRICHLET-MULTINOMIAL



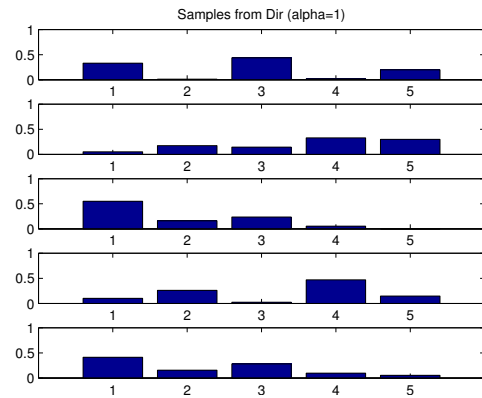
$\boldsymbol{\alpha} = (2, 2, 2)$



$\boldsymbol{\alpha} = (20, 2, 2)$



SAMPLES $\boldsymbol{\alpha} = (0.1, 0.1, 0.1)$



SAMPLES $\alpha = (1, 1, 1)$

The end