

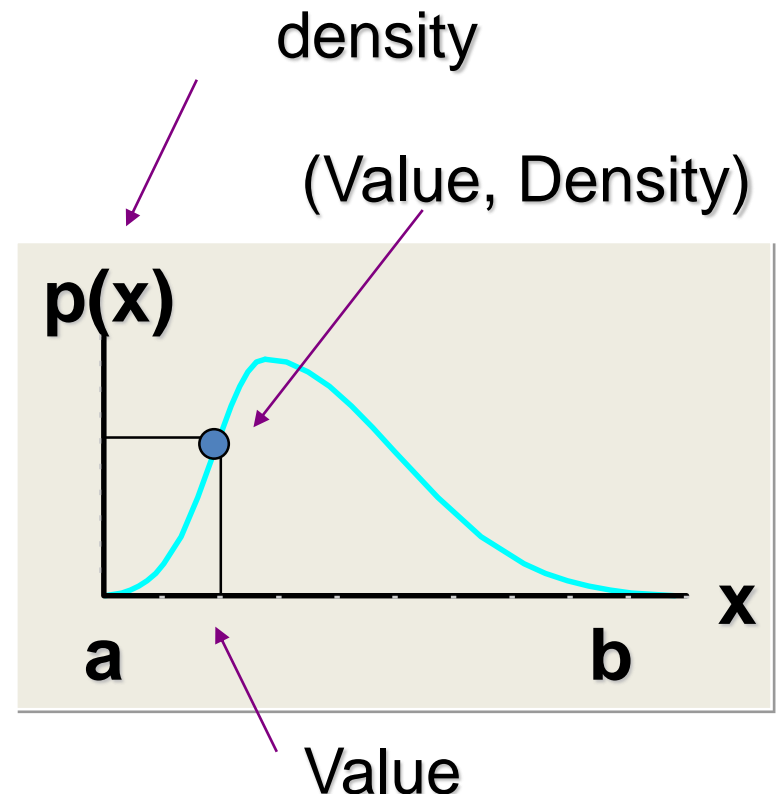
Continuous Probability Distributions

Suresh Manandhar
suresh@cs.york.ac.uk

Continuous Probability Density Function (pdf)

- Shows all values of x in the given interval $[a, b]$, the density $p(x)$
- $p(x)$ is a **probability density function (pdf)**
- Since probabilities need to sum to **1**, the corresponding condition is:

$$F(-\infty \leq X \leq \infty) \\ = \int_{-\infty}^{\infty} p(x) dx = 1$$

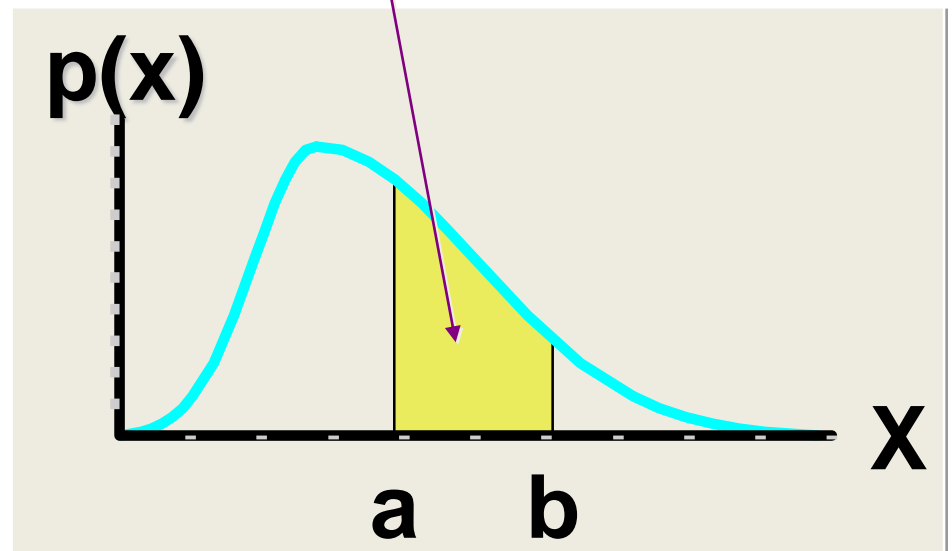


Cumulative density (cdf)

The probability that x lies in the interval $[a, b]$ is given by $F(a \leq X \leq b)$.

**Cumulative probability is
Area Under Curve!**

$$F(a \leq X \leq b) = \int_a^b p(x) dx$$



Some properties

$$F(-\infty \leq X \leq \infty) = \int_{-\infty}^{\infty} p(x) dx = 1$$

$$F(X = a) = F(a \leq X \leq a) = \int_a^a p(x) dx = 0$$

$$F(a) = F(-\infty \leq X \leq a) = \int_{-\infty}^a p(x) dx$$

- or, more generally:

$$F(x) = \int p(x) dx$$

- $F(x)$ is known as the **cumulative distribution function**. (**CDF**).
- The pdf and cdf are related by:

$$\frac{d}{dx} F(x) = p(x)$$

- The following provides an intuitive understanding:

$$F\left(a - \frac{\Delta}{2} \leq X \leq a + \frac{\Delta}{2}\right) = \int_{a - \frac{\Delta}{2}}^{a + \frac{\Delta}{2}} p(x) dx \simeq \Delta p(a)$$

Expectation and Variance

Weighted Average

$$E[X] = \int_{-\infty}^{\infty} x p(x) dx$$

Averaged Squared Distance From Mean

$$Var(X) = \int_{-\infty}^{\infty} (x - E[X])^2 p(x) dx$$

Bayes for Continuous random variables

Law of Total Probability

- For **discrete** random variables:

$$\begin{aligned} p(X = x) &= \sum_i p(X = x, Y = y_i) \\ &= \sum_i p(X = x | Y = y_i) p(Y = y_i) \end{aligned}$$

- For **continuous** random variables:

$$\begin{aligned} p(x) &= \int p(x, y) dy \\ &= \int p(x | y) p(y) dy \end{aligned}$$

Bayes for Continuous random variables

- From the definition of conditional probability:

$$p(y|x) = \frac{p(x, y)}{p(x)} = \frac{p(x|y)p(y)}{\sum_y p(x|y)p(y)}$$

- Substituting for $p(x)$ we derive.

Bayes' theorem for continuous random variables:

$$p(y|x) = \frac{p(x, y)}{p(x)} = \frac{p(x|y)p(y)}{\int p(x|y) p(y) dy}$$

Normal Distribution

- The normal (or Gaussian) distribution gives the familiar bell shaped curve
- **Definition:** A random variable X is normally distributed if its probability density function is given by:

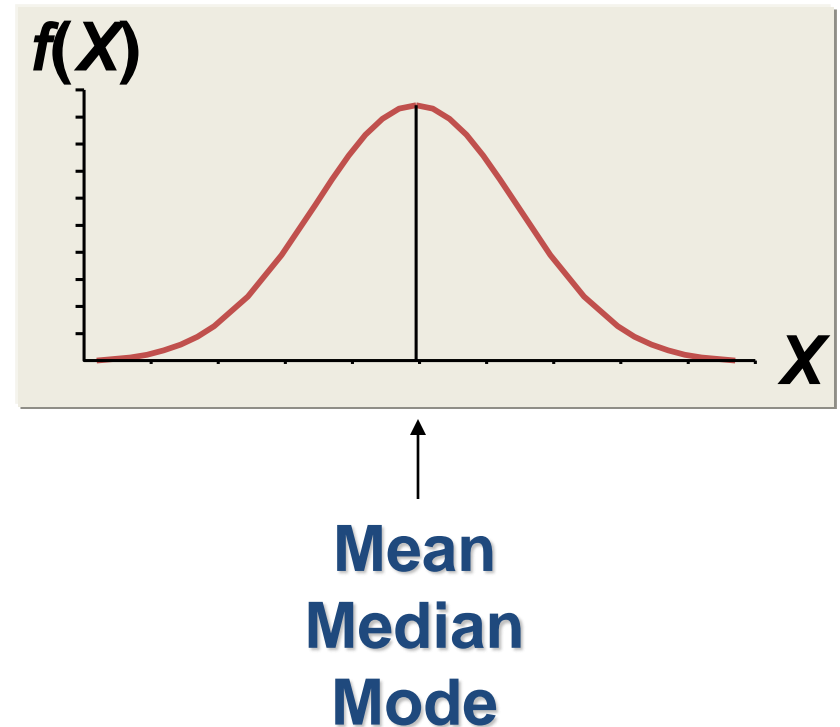
$$p(X = x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

- The parameters μ and σ^2 are known as the mean and variance.
- It turns out that:

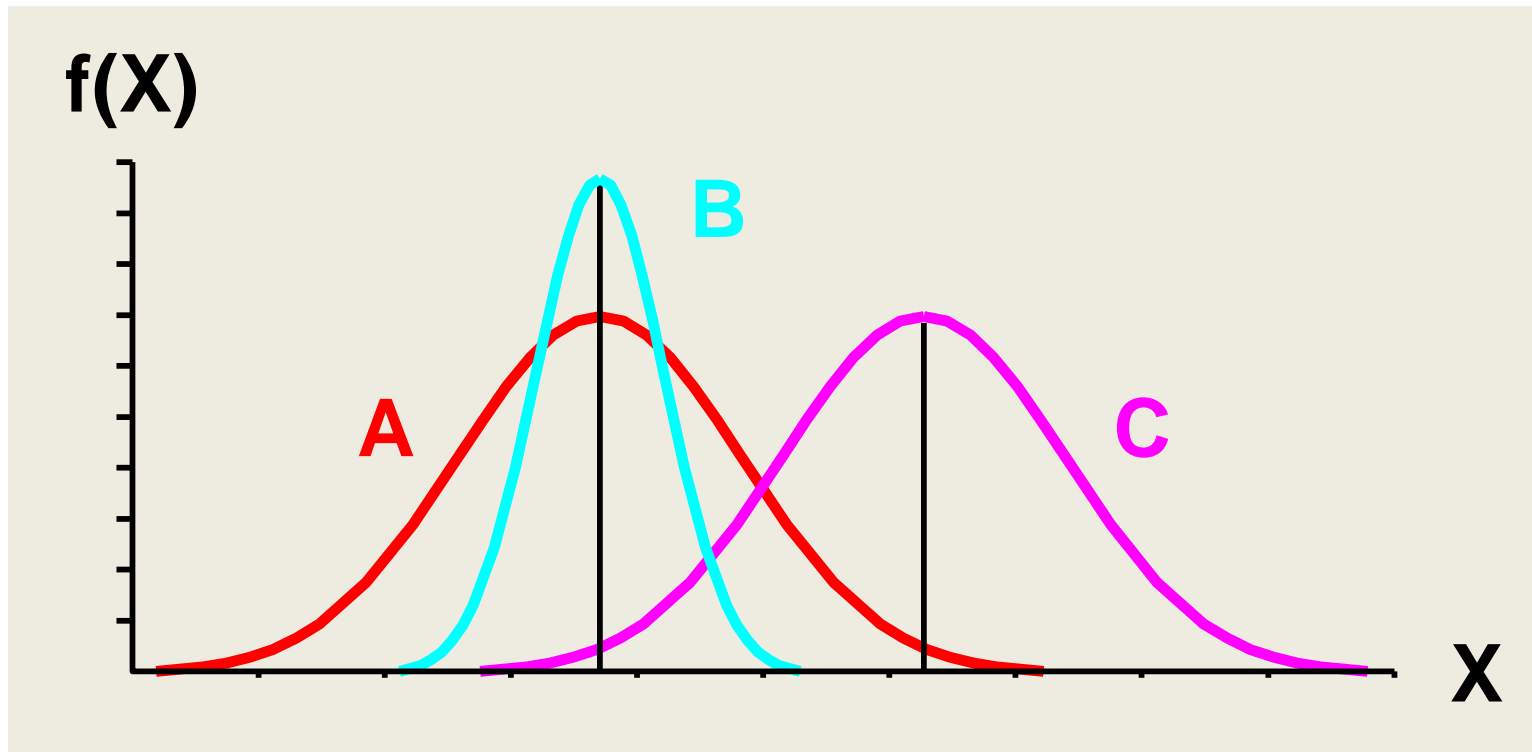
$$\begin{aligned} E[X] &= \mu \\ \text{Var}(X) &= \sigma^2 \end{aligned}$$

The Normal Distribution

- ❑ 1. 'Bell-Shaped' & Symmetrical
- ❑ 2. Mean, median, mode are equal
- ❑ 3. Random variable has infinite range



Effect of varying μ and σ^2



Some properties of Normal distribution

■ Central Limit Theorem

- The sum of a set of independent random variables approaches the Gaussian distribution as the number of variables $\rightarrow \infty$, regardless of the distributions of the individual variables
- (There are generalisations of the CLT.)
- Example: the sum of the face value of randomly drawn playing cards has an approximately Gaussian distribution
- Simplicity: specified by only two intuitive parameters
- Mathematically tractable
 - Many analyses turn out very simply with the Normal distribution

Beta Distribution

- Suppose that the coin factory makes coins that are not perfectly fair all the time
- So, most times the coins are fair but less often the coins are slightly unfair, and less-less often the coins are quite unfair
- So, there is a distribution over the bias of the coins
- **Question:** How to model this distribution?
- The output/sample from this distribution will be the **bias** of a coin i.e. a number between **0** and **1**

Gamma function

- The gamma function $\Gamma(N)$ generalises the factorial function to the reals such that:

$$\Gamma(N + 1) = N! \quad \text{for natural number } N$$

$$\Gamma(1) = 0! = 1$$

$$\Gamma\left(\frac{1}{2}\right) = \left(-\frac{1}{2}\right)! = \sqrt{\pi}$$

$$\Gamma(0) = (-1)! = \frac{\pi}{2}$$

- The gamma function can be used to *interpolate* for values for which the factorial is undefined.

Beta distribution

$$\mathbf{p}(\boldsymbol{\theta}|\boldsymbol{\alpha}, \boldsymbol{\beta}) = \frac{\boldsymbol{\theta}^{\boldsymbol{\alpha}-1}(\mathbf{1} - \boldsymbol{\theta})^{\boldsymbol{\beta}-1}}{B(\boldsymbol{\alpha}, \boldsymbol{\beta})} \sim \textit{Beta}(\boldsymbol{\alpha}, \boldsymbol{\beta})$$

- The $-\mathbf{1}$'s can be thought as mathematical convenience/convention
- To be able to show that such a thing exists, we need to show that:

$$\int_0^1 \mathbf{p}(\boldsymbol{\theta}|\boldsymbol{\alpha}, \boldsymbol{\beta}) d\boldsymbol{\theta} = \int_0^1 \frac{\boldsymbol{\theta}^{\boldsymbol{\alpha}-1}(\mathbf{1} - \boldsymbol{\theta})^{\boldsymbol{\beta}-1}}{B(\boldsymbol{\alpha}, \boldsymbol{\beta})} d\boldsymbol{\theta} = \mathbf{1}$$

- Equivalently, need to show that $B(\boldsymbol{\alpha}, \boldsymbol{\beta})$ is well defined:

$$B(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \int_0^1 \boldsymbol{\theta}^{\boldsymbol{\alpha}-1}(\mathbf{1} - \boldsymbol{\theta})^{\boldsymbol{\beta}-1} d\boldsymbol{\theta}$$

Beta function - derivation

Beta function: $B(\alpha, \beta) = \int_0^1 \theta^{\alpha-1} (1 - \theta)^{\beta-1} d\theta$

- We will apply 'integration by parts', with $u = \theta^{\alpha-1}$, and, $dv = (1 - \theta)^{\beta-1} d\theta$ so the integral becomes

$$B(\alpha, \beta) = \int_0^1 u dv = uv - \int_0^1 v du$$

- $du = (\alpha - 1)\theta^{\alpha-2} d\theta$, and, $v = -\frac{1}{\beta}(1 - \theta)^{\beta}$. Thus:

$$\begin{aligned} &= \theta^{\alpha-1} \left(-\frac{1}{\beta} (1 - \theta)^{\beta} \right) \Big|_0^1 - \int_0^1 \left(-\frac{1}{\beta} (1 - \theta)^{\beta} \right) (\alpha - 1) \theta^{\alpha-2} d\theta \\ &= \frac{(\alpha - 1)}{\beta} \int_0^1 (1 - \theta)^{\beta} \theta^{\alpha-2} d\theta = \frac{(\alpha - 1)}{\beta} B(\alpha - 1, \beta + 1) \\ &= \frac{(\alpha - 1)(\alpha - 2) \dots 1}{\beta(\beta + 1) \dots (\beta + \alpha - 2)} B(1, \beta + \alpha - 1) \\ &= \frac{(\alpha - 1)(\alpha - 2) \dots 1}{\beta(\beta + 1) \dots (\beta + \alpha - 2)} \int_0^1 \theta^{1-1} (1 - \theta)^{\beta + \alpha - 2} d\theta \end{aligned}$$

$$\begin{aligned}
&= \frac{(\alpha - 1)(\alpha - 2) \dots 1}{\beta(\beta + 1) \dots (\beta + \alpha - 2)} \int_0^1 (1 - \theta)^{\beta + \alpha - 2} d\theta \\
&= \frac{(\alpha - 1)(\alpha - 2) \dots 1}{\beta(\beta + 1) \dots (\beta + \alpha - 2)} \left(-\frac{(1 - \theta)^{\beta + \alpha - 1}}{\beta + \alpha - 1} \right) \Big|_0^1 \\
&= \frac{(\alpha - 1)(\alpha - 2) \dots 1}{\beta(\beta + 1) \dots (\beta + \alpha - 2)(\beta + \alpha - 1)} \\
&= \frac{\Gamma(\alpha)}{\beta(\beta + 1) \dots (\beta + \alpha - 2)(\beta + \alpha - 1)} \\
&= \frac{\Gamma(\alpha)}{1 \dots (\beta - 2)(\beta - 1)\Gamma(\beta)} \\
&= \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha + \beta)}
\end{aligned}$$

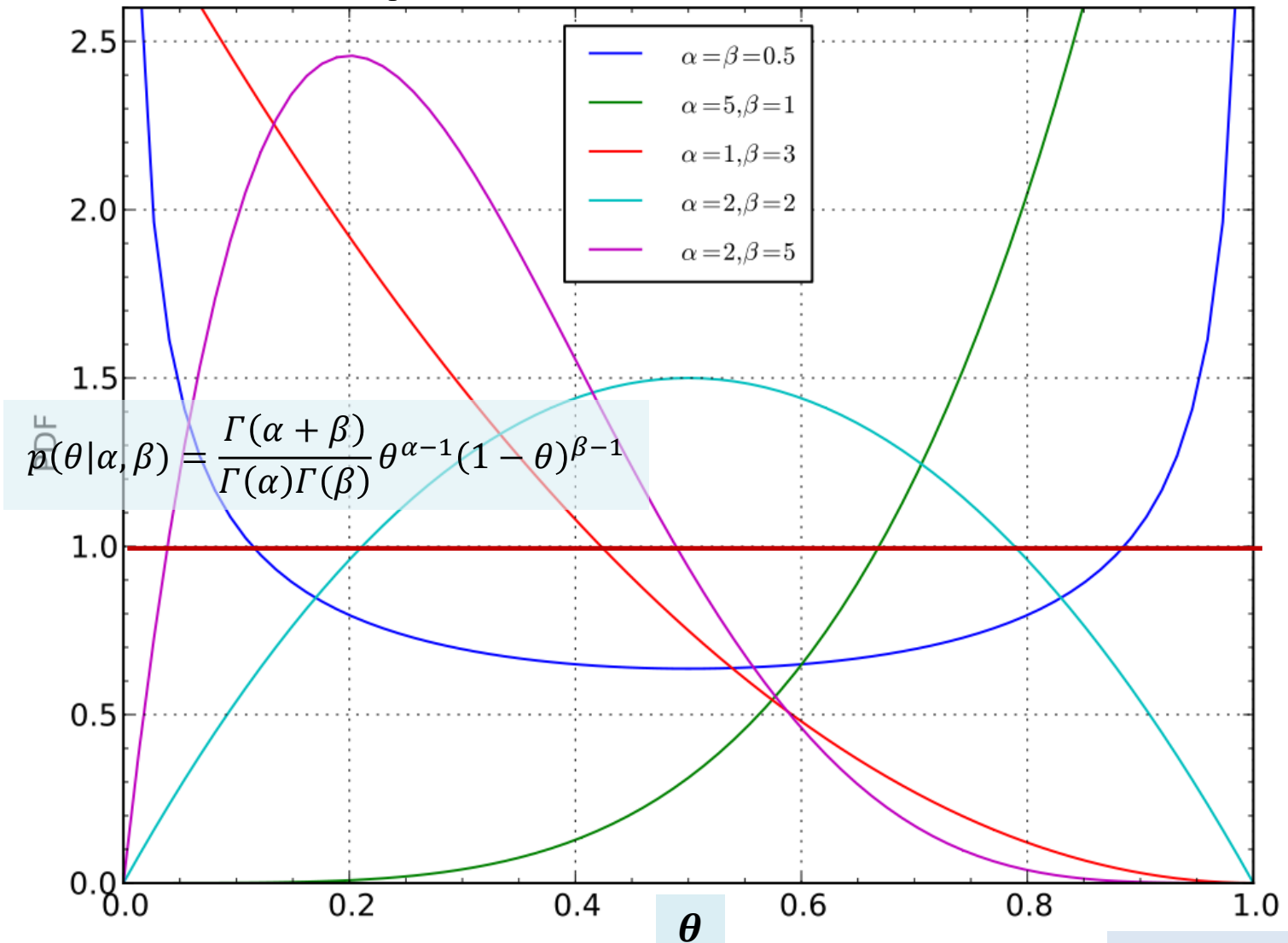
Strictly speaking this proof only applies to positive natural values of α, β

- Thus, $B(\alpha, \beta) = \int_0^1 \theta^{\alpha-1} (1 - \theta)^{\beta-1} d\theta = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}$ hence:

$$p(\theta|\alpha, \beta) = \frac{\theta^{\alpha-1} (1 - \theta)^{\beta-1}}{B(\alpha, \beta)} = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1 - \theta)^{\beta-1}$$

- For $\alpha = \beta = 1$, the distribution is **uniform** with $p(\theta|\alpha, \beta) = 1$

Shape of Beta distribution



from Wikipedia

Modelling the coin factory

- If the factory has a high quality control then choosing $\alpha = \beta = 100$ will give a very peaky pdf
- What would choosing a uniform distribution, $\alpha = \beta = 1$, give? Would that be a good factory?
- Similarly, what about values less than 1?

Modelling typical coins from coin factory

- Suppose, the factory inspector visits the coin factory and picks a coin at random from the factory
- Remember, the factory has beta parameters α, β
- The inspector flips the coin N times and sees

$$\mathbf{c} = (\mathbf{c}_1, \mathbf{c}_2)$$

i.e \mathbf{c}_1 heads and \mathbf{c}_2 tails with $N = \mathbf{c}_1 + \mathbf{c}_2$

- What is the probability of \mathbf{c}_1 heads and \mathbf{c}_2 tails?
- How do we compute this?
- Can we use Bayes?

Modelling typical coins from coin factory

- What we want is: $p(c|\alpha, \beta)$?
- How do we derive this: (*use law of total probability*)

$$\begin{aligned} p(c|\alpha, \beta) &= \int p(c, \theta|\alpha, \beta) d\theta \\ &= \int \underbrace{p(c|\theta)}_{\substack{\text{Binomial} \\ \text{Likelihood}}} \underbrace{p(\theta|\alpha, \beta)}_{\substack{\text{Beta} \\ \text{Prior}}} d\theta \end{aligned}$$

- Integration can be solved analytically (i.e. by hand) if the two distributions have similar form.
- In this case, we say that that two distributions are **conjugate**.

Modelling coin tosses

- Once, we are happy with the choice of α, β for our factory we can plug this in

$$p(c|\alpha, \beta) = \int p(c, \theta|\alpha, \beta) d\theta = \int p(c|\theta) p(\theta|\alpha, \beta) d\theta$$

$$= \int \frac{N!}{c_1! c_2!} \theta^{c_1} (1 - \theta)^{c_2} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1 - \theta)^{\beta-1} d\theta$$

$$= \frac{N!}{c_1! c_2!} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \int \theta^{c_1} (1 - \theta)^{c_2} \theta^{\alpha-1} (1 - \theta)^{\beta-1} d\theta$$

$$= \frac{N!}{c_1! c_2!} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \int \theta^{c_1+\alpha-1} (1 - \theta)^{c_2+\beta-1} d\theta$$

*conjugacy
helps*

$$= \frac{N!}{c_1! c_2!} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \int \theta^{c_1+\alpha-1} (1 - \theta)^{c_2+\beta-1} d\theta$$

$$= \frac{N!}{c_1! c_2!} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \frac{\Gamma(c_1 + \alpha)\Gamma(c_2 + \beta)}{\Gamma(c_1 + c_2 + \alpha + \beta)}$$

*from definition of
Beta function*

N here equals $c_1 + c_2$

The Beta-Binomial Distribution

- What we have just derived is the beta-binomial distribution that gives the (averaged) probability of drawing c_1 heads and c_2 tails from a coin that has been drawn from a beta distribution with parameters α and β .
- The 'averaged' above means that the coin parameter θ has been **integrated out** (and hence no longer appears in the equations):

$$p(c|\alpha, \beta) = \frac{N!}{c_1! c_2!} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \frac{\Gamma(c_1 + \alpha)\Gamma(c_2 + \beta)}{\Gamma(c_1 + c_2 + \alpha + \beta)}$$

- Using the definition of the **beta function**:

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

$$p(c|\alpha, \beta) = \frac{N!}{c_1! c_2!} \frac{B(c_1 + \alpha, c_2 + \beta)}{B(\alpha, \beta)}$$

- $\frac{N!}{c_1! c_2!} \neq \frac{1}{B(C_1+1, C_2+1)}$ since $\frac{1}{B(C_1+1, C_2+1)} = \frac{\Gamma(C_1+C_2+2)}{\Gamma(C_1+1)\Gamma(C_2+1)} = \frac{(N+1)!}{C_1! C_2!}$

Inference

Inference

- The inference problem can be viewed as determining the parameters of your model from observations
- **Example:** If you throw a coin 20 times and you see heads 10 times, is the coin fair?

Inference

- The inference problem can be viewed as determining the parameters of your model from observations
- **Example:** If you throw a coin 20 times and you see heads 10 times, is the coin fair? Or what is the probability that it is fair? How do you even ask such a question?

Inference

- The inference problem can be viewed as determining the parameters of your model from observations
- **Example:** If you throw a coin 20 times and you see heads 10 times, is the coin fair? Or what is the probability that it is fair? How do you even ask such a question?
- The key to Bayesian inference is the mechanism to integrate our prior beliefs into the modelling process and provide mathematically grounded answers to the above questions.

Inference

- **Example:** If you throw a coin 20 times and you see heads 10 times, is the coin fair?

$$p(\mathbf{c} = (\mathbf{c}_1 = 10, \mathbf{c}_2 = 10) | \theta) = \binom{20}{10} \theta^{10} (1 - \theta)^{10}$$

- What is $p(\theta | \mathbf{c} = (\mathbf{c}_1 = 10, \mathbf{c}_2 = 10))$?
- Using the definition of conditional probability:

$$p(\theta, \mathbf{c}) = p(\theta | \mathbf{c}) p(\mathbf{c}) = p(\mathbf{c} | \theta) p(\theta)$$

- Using Law of total probability:

$$p(\mathbf{c}) = \int_0^1 p(\theta, \mathbf{c}) d\theta = \int_0^1 p(\mathbf{c} | \theta) p(\theta) d\theta$$

- Substituting, we get:

$$p(\theta | \mathbf{c}) = \frac{p(\mathbf{c} | \theta) p(\theta)}{\int_0^1 p(\mathbf{c} | \theta) p(\theta) d\theta}$$

Inference

- Thus to find out $p(\theta | c = (c_1 = 10, c_2 = 10))$ we can use:

$$p(\theta | c) = \frac{p(c | \theta)p(\theta)}{\int_0^1 p(c | \theta)p(\theta)d\theta}$$

- However, we do not know what $p(\theta)$ is?
- Hmm... What might this mean?

Inference

- Thus to find out $p(\theta | c = (c_1 = 10, c_2 = 10))$ we can use:

$$p(\theta | c) = \frac{p(c | \theta) p(\theta)}{\int_0^1 p(c | \theta) p(\theta) d\theta}$$

- However, we do not know what $p(\theta)$ is?
- Hmm... What might this mean?
- $p(\theta)$ is our prior belief about the coin with bias θ
- What does this mean?

Inference

Data:

- c is the *data*
- or *observations*

$$p(\theta|c) = \frac{p(c|\theta)p(\theta)}{\int p(c|\theta)p(\theta)d\theta}$$

Posterior:

- $p(\theta|c)$ is the *posterior distribution*
- Gives the probability of the parameter given data

Prior:

- $p(\theta)$ is the *prior distribution*
- Gives the probability for different values of θ and quantifies our belief regarding θ

Likelihood:

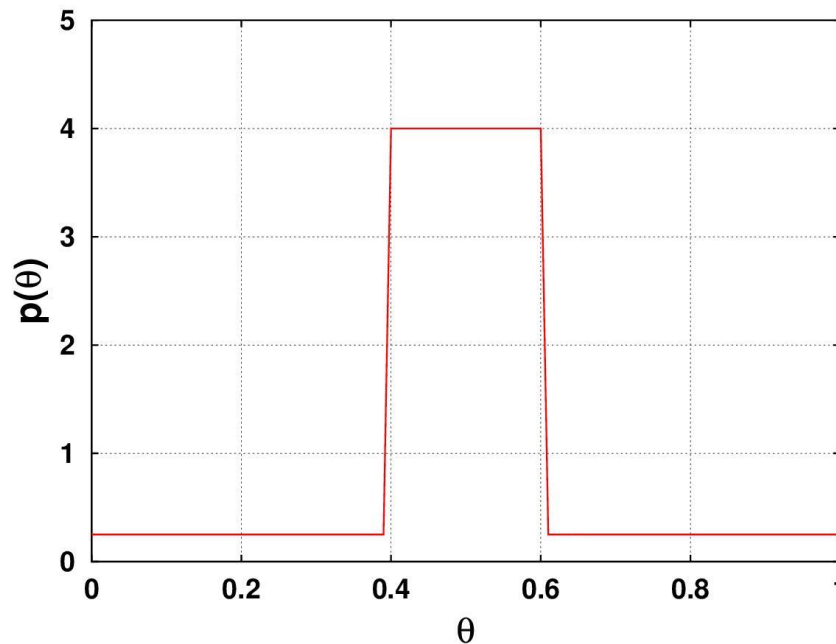
- $p(c|\theta)$ is the *likelihood* of the data given by the model
- Gives the probability of the data being generated by the model

Partition function/Normalising constant/Evidence:

- $\int_0^1 p(c|\theta)p(\theta)d\theta$ is the partition function/normalising constant/evidence.
- This is a constant needed to ensure that the probabilities sum to 1.

Inference

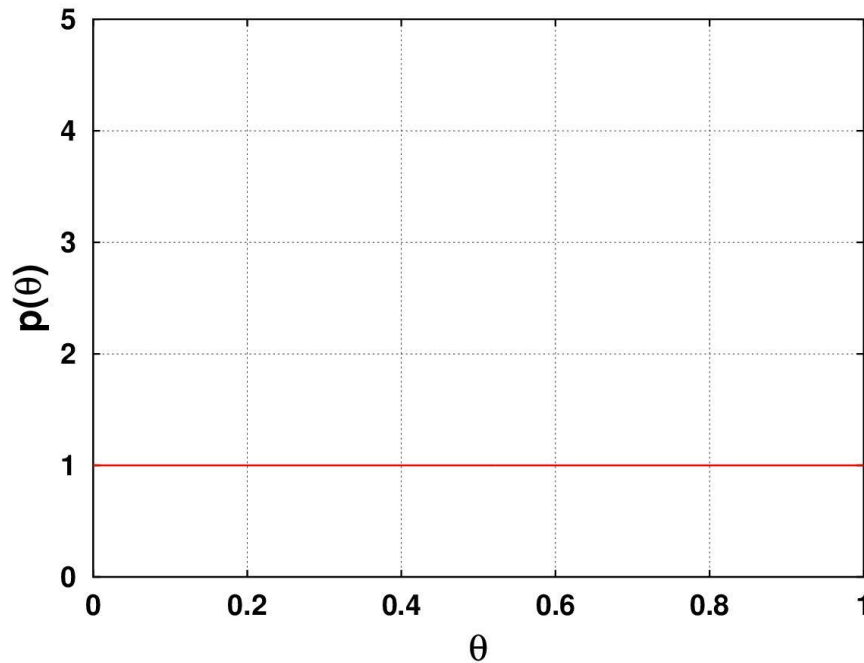
- $p(\theta)$ is our prior belief about the coin with bias θ
- Suppose, we say that, $p(\theta)$ is **piecewise uniform** with:
 - $F(0.4 \leq \theta \leq 0.6) = 0.8, F(0 \leq \theta < 0.4) = 0.1, F(0.6 < \theta \leq 1) = 0.1$



Inference

- Alternatively suppose that, $p(\theta)$ is **uniform** with:

$$p(\theta) = \begin{cases} 1 & \text{if } 0 < \theta < 1 \\ 0 & \text{otherwise} \end{cases}$$



Inference

- Alternatively still, we can assume that $p(\theta)$ is given by a **beta distribution** with parameters α, β

$$p(\theta|\alpha, \beta) = \frac{\theta^{\alpha-1}(1-\theta)^{\beta-1}}{B(\alpha, \beta)} = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1}(1-\theta)^{\beta-1}$$

- The uniform distribution can be recovered by $\alpha = 1, \beta = 1$

$$p(\theta|1, 1) = \frac{\theta^{1-1}(1-\theta)^{1-1}}{B(\alpha, \beta)} = \frac{\Gamma(1 + 1)}{\Gamma(1)\Gamma(1)} = 1$$

- The piecewise uniform distribution may also be approximated by suitable choices of α, β
- Choosing the **beta distribution** as a **prior** means that we can take advantage of **conjugacy**.

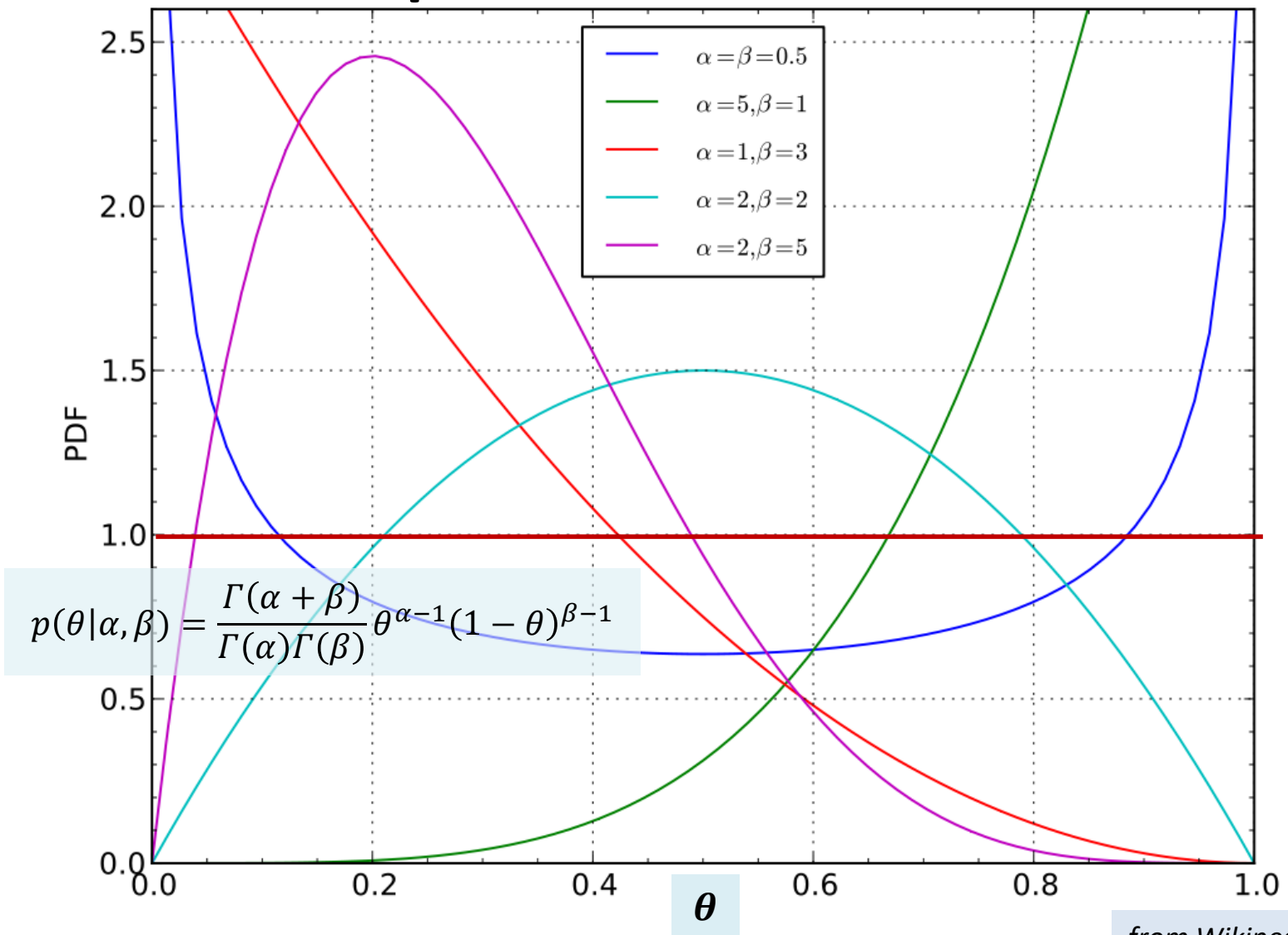
Inference: Binomial Posterior Distribution

- For the case when $p(\theta)$ is given by the **beta distribution** with parameters α, β :

$$\begin{aligned}
 p(\theta|c, \alpha, \beta) &= \frac{p(c|\theta, \alpha, \beta) p(\theta|\alpha, \beta)}{\int_0^1 p(c|\theta, \alpha, \beta) p(\theta|\alpha, \beta) d\theta} && \text{Remember } c = (c_1, c_2) \\
 &= \frac{\frac{N!}{c_1! c_2!} \theta^{c_1} (1 - \theta)^{c_2} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1 - \theta)^{\beta-1}}{\int_0^1 \frac{N!}{c_1! c_2!} \theta^{c_1} (1 - \theta)^{c_2} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1 - \theta)^{\beta-1} d\theta} \\
 &= \frac{\theta^{c_1 + \alpha - 1} (1 - \theta)^{c_2 + \beta - 1}}{\frac{\Gamma(c_1 + \alpha)\Gamma(c_2 + \beta)}{\Gamma(c_1 + c_2 + \alpha + \beta)}} \\
 &= \frac{\Gamma(c_1 + c_2 + \alpha + \beta)}{\Gamma(c_1 + \alpha)\Gamma(c_2 + \beta)} \theta^{c_1 + \alpha - 1} (1 - \theta)^{c_2 + \beta - 1} \\
 &= \frac{1}{B(c_1 + \alpha, c_2 + \beta)} \theta^{c_1 + \alpha - 1} (1 - \theta)^{c_2 + \beta - 1} \\
 &= p(\theta|\alpha + c_1, \beta + c_2) \sim \text{Beta}(\alpha + c_1, \beta + c_2)
 \end{aligned}$$

- Thus, posterior of binomial (under beta) is beta with **counts added into the parameters**

Shape of Beta distribution



Aside: Inference

- If θ is *piecewise uniform* with, $F(0.4 \leq \theta \leq 0.6) = 0.8$,
 $F(0 \leq \theta < 0.4) = 0.1$, $F(0.6 < \theta \leq 1) = 0.1$
- We can calculate, $p(\theta)$:

$$F(0.4 \leq \theta \leq 0.6) = 0.8 \text{ implies } p(\theta) = \frac{0.8}{0.2} = 4$$

$$F(0 \leq \theta < 0.4) = 0.1 \text{ implies } p(\theta) = \frac{0.1}{0.4} = \frac{1}{4}$$

$$F(0.6 < \theta \leq 1) = 0.1 \text{ implies } p(\theta) = \frac{0.1}{0.4} = \frac{1}{4}$$

Aside: Inference

$$\begin{aligned} p(\theta|c) &= \frac{p(c|\theta) p(\theta)}{\int_0^1 p(c|\theta) p(\theta) d\theta} \\ &= \frac{p(c|\theta) p(\theta)}{\int_0^{0.4} p(c|\theta) \cdot \frac{1}{4} d\theta + \int_{0.4}^{0.6} p(c|\theta) \cdot 4 d\theta + \int_{0.6}^1 p(c|\theta) \cdot \frac{1}{4} d\theta} \\ &= \frac{\theta^{c_1}(1-\theta)^{c_2} p(\theta)}{\frac{1}{4} IB(.4, c_1, c_2) + 4[IB(.6, c_1, c_2) - IB(.4, c_1, c_2)] + \frac{1}{4} [B(c_1, c_2) - IB(.6, c_1, c_2)]} \end{aligned}$$

- Where ***IB*** is the ***incomplete beta function*** defined by:

$$IB(a, c_1, c_2) = \int_0^a \theta^{c_1-1} (1-\theta)^{c_2-1} d\theta \quad 0 \leq a \leq 1$$

- Hence: $IB(1, c_1, c_2) = B(c_1, c_2)$
- Most math packages provide implementations of the incomplete beta function.
- Thus, posterior probability for any θ and c can be calculated.

Bayesian Inference

$$p(\theta|c) = \frac{p(c|\theta)p(\theta)}{\int p(c|\theta)p(\theta)d\theta}$$

- ***Being Bayesian*** typically means that we treat $p(\theta|c)$ as a distribution
- This means that we do not get a single value for the model parameter θ
- However, sometimes, we may want to find out what the '**best**' value for the model parameter θ would be.
- How do define what '**best**' means ?

MLE Inference

- This is the maximum likelihood estimate (where c is the data)

$$\theta_{MLE} = \arg \max_{\theta} p(c|\theta)$$

- Thus MLE is simply the maximum/**mode** of the model likelihood.
- We don't even need to define a prior, so (apparently) more “objective”.
- Typically deal with situations where there is only one mode, so can talk about **the** MLE.

MAP Inference

$$p(\theta|c) = \frac{p(c|\theta)p(\theta)}{\int p(c|\theta)p(\theta)d\theta}$$

MAP (Maximum a Posteriori):

- Assume that $p(\theta)$ is distributed as per **some given distribution**

$$\theta_{MAP} = \arg \max_{\theta} p(\theta|c) = \arg \max_{\theta} \frac{p(c|\theta)p(\theta)}{\int_0^1 p(c|\theta)p(\theta)d\theta}$$

- Since the denominator is a constant:

$$\theta_{MAP} = \arg \max_{\theta} p(c|\theta)p(\theta)$$

- Thus MAP estimate is the **mode** of the posterior distribution

Multinomial Distribution

- Probability of observing $\mathbf{c} = (c_1, \dots, c_k)$ heads in all possible ways out of $C = \sum_i c_i$ throws from a k-headed dice with probability of heads $\boldsymbol{\theta} = (\theta_1, \dots, \theta_k)$ s.t. $\sum_i \theta_i = 1$

$$p(\text{heads} = \mathbf{c} | \boldsymbol{\theta}) = \frac{C!}{\prod_i c_i!} \prod_i \theta_i^{c_i}$$

- Points to note:
 - Generalises the binomial distribution to $k > 2$
 - Equivalent to the binomial for $k = 2$

The Dirichlet Distribution

- Dirichlet distribution generalises the Beta distribution to the $k - 1$ probability simplex

- The Beta distribution:

$$p(\theta|\alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \theta^{\alpha-1} (1 - \theta)^{\beta-1} \sim \text{Beta}(\alpha, \beta)$$

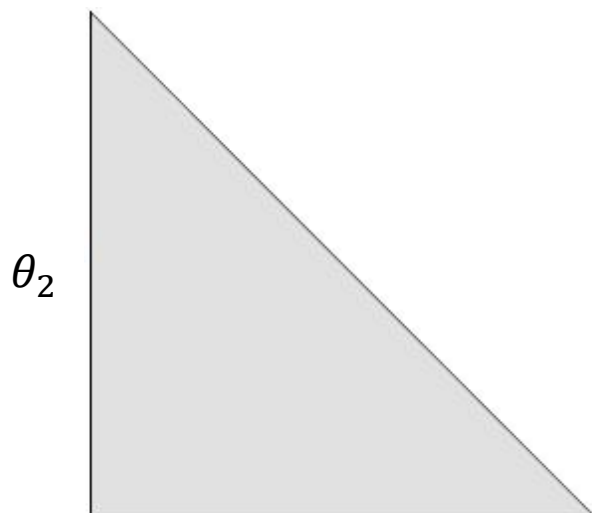
- The Dirichlet distribution:

$$p(\theta|\alpha) = \frac{\Gamma(A)}{\prod_i \Gamma(\alpha_i)} \prod_i \theta_i^{\alpha_i-1} \sim \text{Dir}(\alpha_1, \dots, \alpha_k)$$

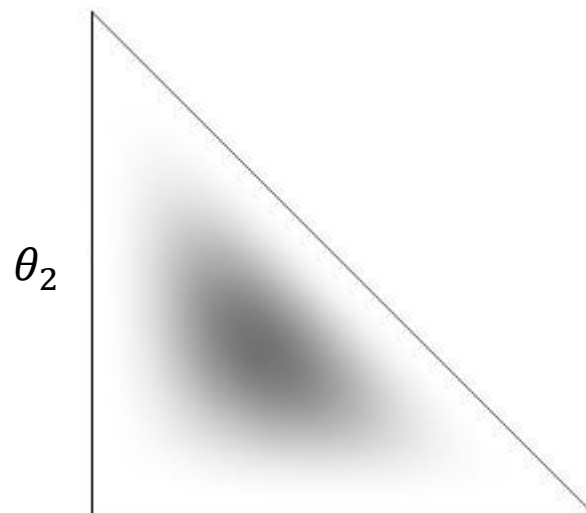
with:

$$\alpha = (\alpha_1, \dots, \alpha_k), \theta = (\theta_1, \dots, \theta_k), \sum_i \theta_i = 1, A = \sum_i \alpha_i$$

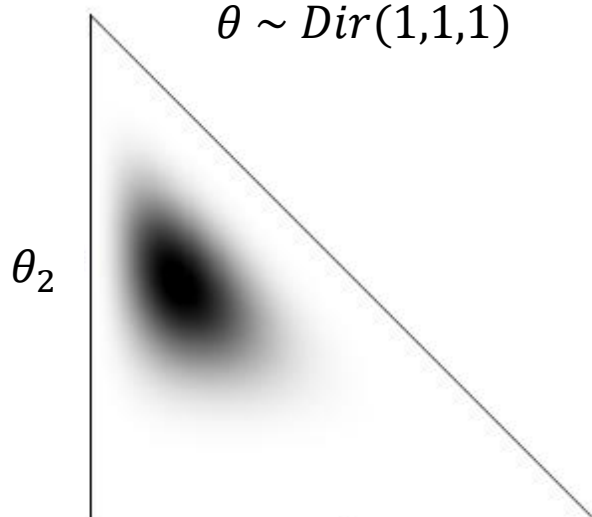
- So, the samples from the Dirichlet distribution can be used to model the bias in a k -sided dice.



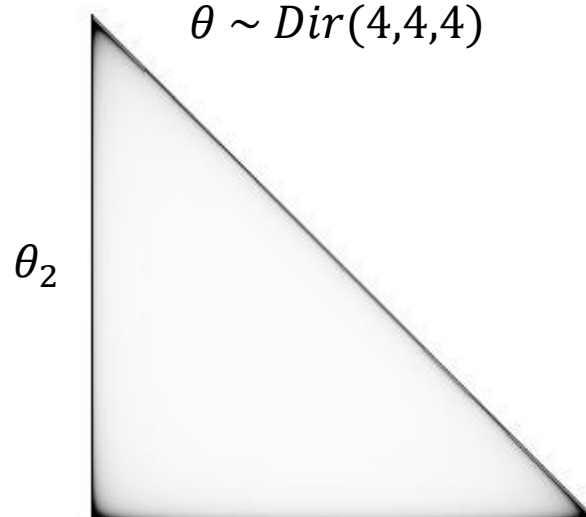
$$\theta \sim \text{Dir}(1,1,1)$$



$$\theta \sim \text{Dir}(4,4,4)$$



$$\theta \sim \text{Dir}(4,9,7)$$



$$\theta \sim \text{Dir}(0.2, 0.2, 0.2)$$

The Dirichlet-Multinomial

- Like for the Beta-Binomial distribution we can integrate out the Multinomial parameters
- Here, we are modelling the case where we need to predict the outcome $\mathbf{c} = (c_1, \dots, c_k)$ of a dice throw from a dice sampled from a factory with parameter $\alpha = (\alpha_1, \dots, \alpha_k)$

$$\begin{aligned}
 p(\mathbf{c}|\alpha) &= \int p(\mathbf{c}, \boldsymbol{\theta}|\alpha) d\boldsymbol{\theta} = \int p(\mathbf{c}|\boldsymbol{\theta}) p(\boldsymbol{\theta}|\alpha) d\boldsymbol{\theta} \\
 &= \int \left[\frac{\mathbf{C}!}{\prod_i c_i!} \prod_i \theta^{c_i} \right] \frac{\Gamma(A)}{\prod_i \Gamma(\alpha_i)} \prod_i \theta^{\alpha_i-1} d\boldsymbol{\theta} \\
 &= \frac{\mathbf{C}!}{\prod_i c_i!} \frac{\Gamma(A)}{\prod_i \Gamma(\alpha_i)} \int \prod_i \theta^{c_i+\alpha_i-1} d\boldsymbol{\theta}
 \end{aligned}$$

- From earlier: $B(\alpha, \beta) = \int_0^1 \theta^{\alpha-1} (1-\theta)^{\beta-1} d\theta = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}$

- This generalises to: $\Delta(\alpha) = \int_0^1 \prod_i \theta^{\alpha_i-1} d\boldsymbol{\theta} = \frac{\prod_i \Gamma(\alpha_i)}{\Gamma(A)}$

$$p(\mathbf{c}|\alpha) = \frac{\mathbf{C}!}{\prod_i c_i!} \frac{\Gamma(A)}{\prod_i \Gamma(\alpha_i)} \int \prod_i \theta^{c_i+\alpha_i-1} d\boldsymbol{\theta} = \frac{\mathbf{C}!}{\prod_i c_i!} \frac{\Gamma(A)}{\prod_i \Gamma(\alpha_i)} \frac{\prod_i \Gamma(c_i + \alpha_i)}{\Gamma(\mathbf{C} + A)}$$

Posterior Multinomial under Dirichlet prior

- Suppose we observe counts $\mathbf{c} = (c_1, \dots, c_k)$ from a dice sampled from our factory.
- We would like to predict the most likely parameters $\boldsymbol{\theta}$ for this dice.

$$\begin{aligned} p(\boldsymbol{\theta}|\mathbf{c}, \boldsymbol{\alpha}) &= \frac{p(\mathbf{c}, \boldsymbol{\theta}|\boldsymbol{\alpha})}{p(\mathbf{c}|\boldsymbol{\alpha})} = \frac{p(\mathbf{c}, \boldsymbol{\theta}|\boldsymbol{\alpha})}{\int p(\mathbf{c}, \boldsymbol{\theta}|\boldsymbol{\alpha}) d\boldsymbol{\theta}} = \frac{p(\mathbf{c}|\boldsymbol{\theta})p(\boldsymbol{\theta}|\boldsymbol{\alpha})}{\int p(\mathbf{c}|\boldsymbol{\theta})p(\boldsymbol{\theta}|\boldsymbol{\alpha}) d\boldsymbol{\theta}} \\ &= \frac{\frac{\mathbf{C}!}{\prod_i c_i!} \prod_i \theta_i^{c_i} \frac{\Gamma(A)}{\prod_i \Gamma(\alpha_i)} \prod_i \theta_i^{\alpha_i-1}}{\int \frac{\mathbf{C}!}{\prod_i c_i!} \frac{\Gamma(A)}{\prod_i \Gamma(\alpha_i)} \prod_i \theta_i^{c_i+\alpha_i-1} d\boldsymbol{\theta}} = \frac{\prod_i \theta_i^{c_i+\alpha_i-1}}{\frac{\Gamma(\mathbf{C} + A)}{\prod_i \Gamma(c_i + \alpha_i)}} \sim \text{Dir}(c_1 + \alpha_1, \dots, c_k + \alpha_k) \end{aligned}$$

- So, the shape of the posterior is exactly like that of the prior with **counts added into the parameters**