

---

# *Introduction to Probability and Statistics (Continued)*

*Prof. Nicholas Zabararas  
Center for Informatics and Computational Science*

*<https://cics.nd.edu/>*

*University of Notre Dame  
Notre Dame, Indiana, USA*

*Email: [nzabararas@gmail.com](mailto:nzabararas@gmail.com)*

*URL: <https://www.zabararas.com/>*

*August 27, 2018*



# Contents

---

- The binomial and Bernoulli distributions
- The multinomial and multinoulli distributions
- The Poisson distribution
- Student's T
- Laplace distribution
- Gamma distribution
- Beta distribution
- Pareto distribution
- Introduction to Covariance and Correlation



# References

---

- Following closely [Chris Bishops' PRML book](#), Chapter 2
- Kevin Murphy's, [Machine Learning: A probabilistic perspective](#), Chapter 2
- Jaynes, E. T. (2003). [Probability Theory: The Logic of Science](#). Cambridge University Press.
- Bertsekas, D. and J. Tsitsiklis (2008). [Introduction to Probability](#). Athena Scientific. 2nd Edition
- Wasserman, L. (2004). [All of statistics. A Concise Course in Statistical Inference](#). Springer.



# Binary Variables

- Consider a coin flipping experiment with heads = 1 and tails = 0. With  $\mu \in [0,1]$

$$p(x = 1 \mid \mu) = \mu$$

$$p(x = 0 \mid \mu) = 1 - \mu$$

- This defines the Bernoulli distribution as follows:

$$\mathcal{Bern}(x \mid \mu) = \mu^x (1 - \mu)^{1-x}$$

- Using the indicator function, we can also write this as:

$$\mathcal{Bern}(x \mid \mu) = \mu^{\mathbb{I}(x=1)} (1 - \mu)^{\mathbb{I}(x=0)}$$

# Bernoulli Distribution

- Recall that in general

$$\mathbb{E}[f] = \sum_x p(x) f(x), \quad \mathbb{E}[f] = \int p(x) f(x) dx$$

$$\text{var}[f] = \mathbb{E}[f(x)^2] - \mathbb{E}[f(x)]^2$$

- For the Bernoulli distribution  $\mathcal{Bern}(x | \mu) = \mu^x (1 - \mu)^{1-x}$ , we can easily show from the definitions:

$$\mathbb{E}[x] = \mu$$

$$\text{var}[x] = \mu(1 - \mu)$$

$$\mathbb{H}[x] \equiv - \sum_{x \in \{0,1\}} p(x | \mu) \ln p(x | \mu) = -\mu \ln \mu - (1 - \mu) \ln(1 - \mu)$$

- Here  $\mathbb{H}[x]$  is the “entropy of the distribution”



# Likelihood Function for Bernoulli Distribution

- Consider the data set

$$\mathcal{D} = \{x_1, x_2, \dots, x_N\}$$

in which we have  $m$  heads ( $x = 1$ ), and  $N - m$  tails ( $x = 0$ )

- The **likelihood function** takes the form:

$$p(\mathcal{D} \mid \mu) = \prod_{n=1}^N p(x_n \mid \mu) = \prod_{n=1}^N \mu^{x_n} (1 - \mu)^{1-x_n} = \mu^m (1 - \mu)^{N-m}$$

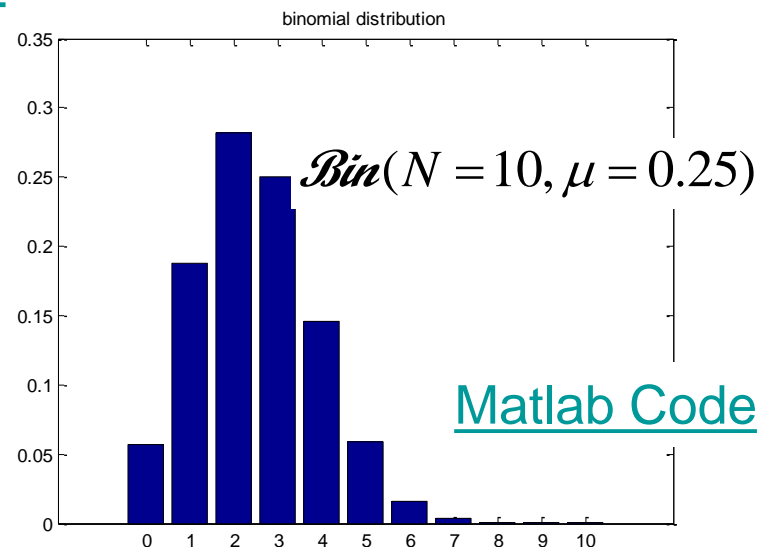
# Binomial Distribution

➤ Consider the discrete random variable  $X \in \{0, 1, 2, \dots, N\}$

➤ We define the Binomial distribution as follows:

$$\mathcal{Bin}(X = m \mid N, \mu) = \binom{N}{m} \mu^m (1 - \mu)^{N-m}$$

➤ In our coin flipping experiment, it gives the probability in  $N$  flips to get  $m$  heads with  $\mu$  being the probability getting heads in one flip.



➤ It can be shown (see S. Ross, Introduction to Probability Models) that the limit of the binomial distribution as

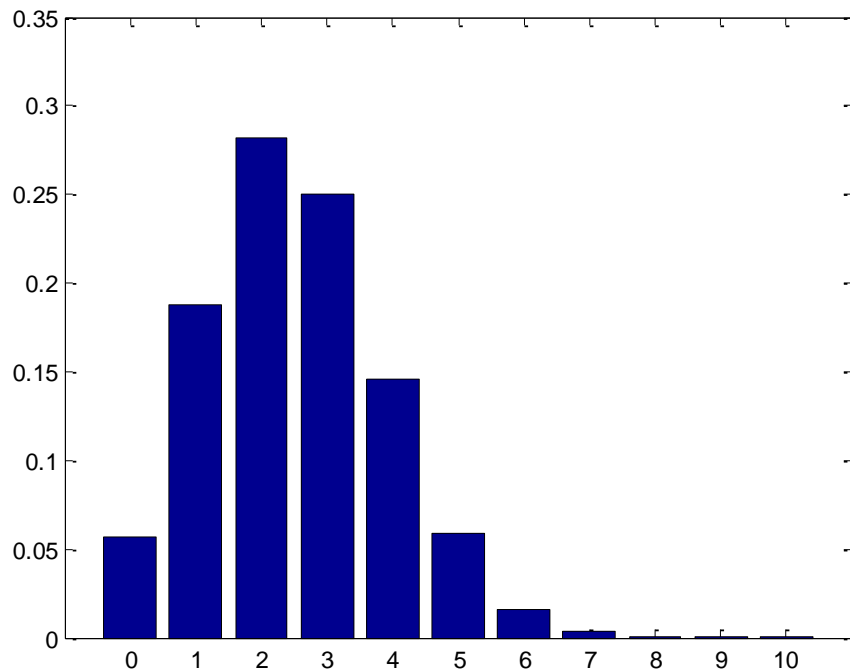
$N \rightarrow \infty, N\mu \rightarrow \lambda$ , is the Poisson( $\lambda$ ) distribution.



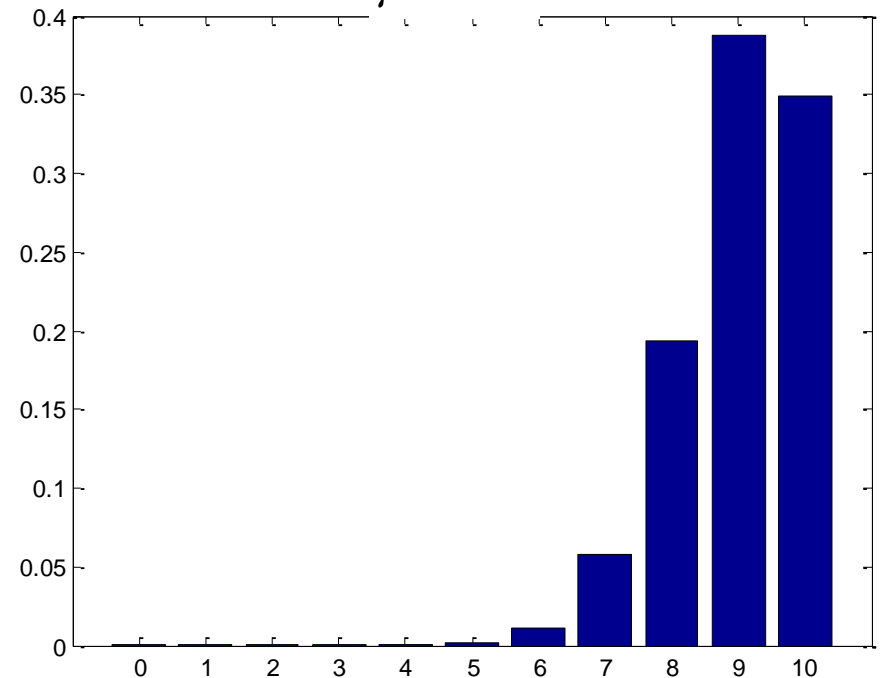
# Binomial Distribution

- The Binomial distribution for  $N = 10$ , and  $\mu \in \{0.25, 0.9\}$  is shown below using MatLab function [binomDistPlot](#) from [Kevin Murphys' PMTK](#)

$\mu = 0.25$



$\mu = 0.9$



$\mathcal{B}in(N, \mu)$



# Mean, Variance of the Binomial Distribution

- Consider for independent events the mean of the sum is the sum of the means, and the variance of the sum is the sum of the variances.

- Because  $m = x_1 + \dots + x_N$ , and for each observation the mean and variance are known from the Bernoulli distribution:

$$\mathbb{E}[m] \equiv \sum_{m=0}^N m \mathcal{Bin}(m | N, \mu) = \mathbb{E}[x_1 + \dots + x_N] = N\mu$$

$$\text{var}[m] \equiv \sum_{m=0}^N (m - \mathbb{E}[m])^2 \mathcal{Bin}(m | N, \mu) = \text{var}[x_1 + \dots + x_N] = N\mu(1 - \mu)$$

- One can also compute  $\mathbb{E}[m]$ ,  $\mathbb{E}[m^2]$  by differentiating (twice) the identity  $\sum_{m=1}^N \binom{N}{m} \mu^m (1 - \mu)^{N-m} = 1$  wrt  $\mu$ . Try it!

# Binomial Distribution: Normalization

To show that the Binomial is correctly normalized, we use the following identities:

- Can be shown with direct substitution:  $\binom{N}{n} + \binom{N}{n-1} = \binom{N+1}{n} \quad (*)$
- Binomial theorem:  $(1+x)^N = \sum_{m=0}^N \binom{N}{m} x^m \quad (**)$

This theorem is proved by induction using (\*) and noticing:

$$\begin{aligned} (1+x)^{N+1} &= \sum_{m=0}^N \binom{N}{m} x^m (1+x) = \sum_{m=0}^N \binom{N}{m} x^m + \sum_{m=0}^N \binom{N}{m} x^{m+1} = \sum_{m=0}^N \binom{N}{m} x^m + \sum_{m=1}^{N+1} \binom{N}{m-1} x^m = \\ &= \left( 1 + \sum_{m=1}^N \binom{N}{m} x^m \right) + \left( \sum_{m=1}^N \binom{N}{m-1} x^m + x^{N+1} \right) \stackrel{*}{=} 1 + \sum_{m=1}^N \binom{N+1}{m} x^m + x^{N+1} = \sum_{m=0}^{N+1} \binom{N+1}{m} x^m \end{aligned}$$

- To finally show normalization using (\*\*):

$$\sum_{m=0}^N \binom{N}{m} \mu^m (1-\mu)^{N-m} = (1-\mu)^N \sum_{m=0}^N \binom{N}{m} \left( \frac{\mu}{1-\mu} \right)^m = (1-\mu)^N \left( 1 + \frac{\mu}{1-\mu} \right)^N = 1$$



# Generalization of the Bernoulli Distribution

- We are now looking at discrete variables that can take on one of  $K$  possible mutually exclusive states.
- The variable is represented by a  $K$ -dimensional vector  $\mathbf{x}$  in which one of the elements  $x_k$  equals 1, and all remaining elements equal 0:  $\mathbf{x} = (0, 0, \dots, 1, 0, \dots, 0)^T$

These vectors satisfy:  $\sum_{k=1}^K x_k = 1$

- Let the probability of  $x_k = 1$  be denoted as  $\mu_k$ . Then

$$p(\mathbf{x} | \boldsymbol{\mu}) = \prod_{k=1}^K \mu_k^{x_k} = \prod_{k=1}^K \mu_k^{\mathbb{I}(x_k=1)}, \quad \sum_{k=1}^K \mu_k = 1, \mu_k \geq 0$$

where  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_K)^T$ .

# Multinoulli/Categorical Distribution

- The distribution is already normalized:

$$\sum_x p(\mathbf{x} | \boldsymbol{\mu}) = \sum_x \prod_{k=1}^K \mu_k^{x_k} = \sum_{k=1}^K \mu_k = 1$$

- The mean of the distribution is computed as:

$$\mathbb{E}[\mathbf{x} | \boldsymbol{\mu}] = \sum_x \mathbf{x} p(\mathbf{x} | \boldsymbol{\mu}) = (\mu_1, \dots, \mu_K)^T = \boldsymbol{\mu}$$

similar to the result for the Bernoulli distribution.

- The Multinoulli also known as the *Categorical distribution* often denoted as ( $\mathcal{Mu}$  here is the multinomial distribution):

$$\mathcal{Cat}(\mathbf{x} | \boldsymbol{\mu}) = \mathcal{Multinoulli}(\mathbf{x} | \boldsymbol{\mu}) = \mathcal{Mu}(\mathbf{x} | 1, \boldsymbol{\mu})$$

- The parameter 1 stands to emphasize that we roll a  $K$ -sided dice once ( $N = 1$ ) – see next for the multinomial distribution.



# Likelihood: Multinoulli Distribution

- Let us consider a data set  $\mathcal{D} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$ . The likelihood becomes:

$$p(\mathcal{D} | \boldsymbol{\mu}) = \prod_{n=1}^N \prod_{k=1}^K \mu_k^{x_{nk}} = \prod_{k=1}^K \mu_k^{\sum_{n=1}^N x_{nk}} = \prod_{k=1}^K \mu_k^{m_k}, \quad \text{where: } m_k = \sum_{n=1}^N x_{nk}$$

is the # of observations of  $x_k = 1$ .

- $m_k$  is the “**sufficient statistic**” of the distribution.

# MLE Estimate: Multinoulli Distribution

- To compute the maximum likelihood (MLE) estimate of  $\mu$ , we maximize an augmented log-likelihood

$$\ln p(\mathcal{D} | \mu) + \lambda \left( \sum_{k=1}^K \mu_k - 1 \right) = \sum_{k=1}^K m_k \ln \mu_k + \lambda \left( \sum_{k=1}^K \mu_k - 1 \right)$$

- Setting the derivative wrt  $\mu_k$  equal to zero:  $\mu_K = -\frac{m_k}{\lambda}$

- Substitution into the constraint

$$\sum_{k=1}^K \mu_k = 1 \Rightarrow -\frac{\sum_{k=1}^K m_k}{\lambda} = 1 \Rightarrow \lambda = -\sum_{k=1}^K m_k \Rightarrow$$

$$\mu_K = \frac{m_k}{\sum_{k=1}^K m_k} = \frac{m_k}{N}$$

As expected, this is the fraction in the  $N$  observations of  $x_k = 1$

# Multinomial Distribution

- We can also consider *the joint distribution of  $m_1, \dots, m_K$  in  $N$  observations* conditioned on the parameters  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_K)$ .
- From the expression for the likelihood given earlier

$$p(\mathcal{D} | \boldsymbol{\mu}) = \prod_{k=1}^K \mu_k^{m_k}$$

the multinomial distribution  $\mathcal{Mu}(m_1, \dots, m_K | N, \boldsymbol{\mu})$  with parameters  $N$  and  $\boldsymbol{\mu}$  takes the form:

$$p(m_1, m_2, \dots, m_K | N, \mu_1, \mu_2, \dots, \mu_K) = \frac{N!}{m_1! m_2! \dots m_K!} \mu_1^{m_1} \mu_2^{m_2} \dots \mu_K^{m_K} \quad \text{where } \sum_{k=1}^K m_k = N$$

# Example: Biosequence Analysis

- Consider a set of DNA sequences where there are 10 rows (sequences) and 15 columns (locations along the genome).
- Several locations are conserved by evolution (e.g., because they are part

c	g	a	t	a	c	g	g	g	t	c	g	a	a	
c	a	a	t	c	c	g	a	g	a	t	c	g	c	a
c	a	a	t	c	c	g	t	g	t	g	g	g	a	
c	a	a	t	c	g	g	c	a	t	g	c	g	g	g
c	g	a	g	c	c	g	c	g	t	a	c	g	a	a
c	a	t	a	c	g	g	a	g	c	a	c	g	a	a
t	a	a	t	c	c	g	g	g	c	a	t	g	t	a
c	g	a	g	c	c	g	a	g	t	a	c	a	g	a
c	c	a	t	c	c	g	c	g	t	a	a	g	c	a
g	g	a	t	a	c	g	a	g	a	t	g	a	c	a

Sequences ↓  
N=1:10

Location along the genome →

of a gene coding region), since the corresponding columns tend to be pure e.g., column 7 is all *g*'s.

- To visualize the data (**sequence logo**), we plot the letters *A*, *C*, *G* and *T* with a font size proportional to their empirical probability, and with the most probable letter on the top.



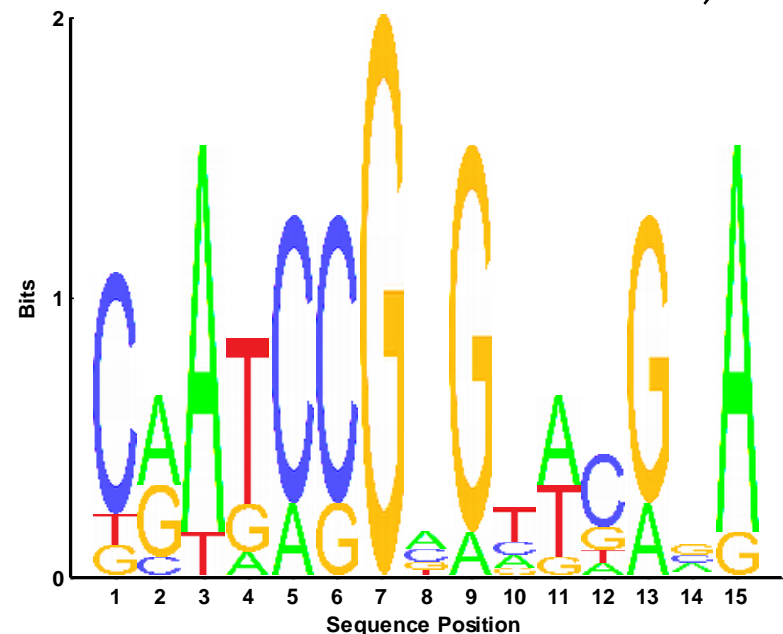


# Example: Biosequence Analysis

- The empirical probability distribution at location  $t$ , is obtained by normalizing the vector of counts ([see MLE estimate](#))

$$\hat{\theta}_t = \frac{1}{N} \left( \sum_{i=1}^N \mathbb{I}(X_{it} = 1), \sum_{i=1}^N \mathbb{I}(X_{it} = 2), \sum_{i=1}^N \mathbb{I}(X_{it} = 3), \sum_{i=1}^N \mathbb{I}(X_{it} = 4) \right)$$

- *This distribution is known as a **motif**.*
- Can also compute the *most probable letter in each location*; this is the *consensus sequence*.



Use MatLab function [seqlogoDemo](#) from [Kevin Murphys' PMTK](#)



# Summary of Discrete Distributions

- A summary of the multinomial and related discrete distributions is summarized below on a Table from [Kevin Murphy's textbook](#)

Name	$n$	$K$	$x$
Multinomial	-	-	$\mathbf{x} \in \{0, 1, \dots, n\}^K, \sum_{k=1}^K x_k = n$
Multinoulli	1	-	$\mathbf{x} \in \{0, 1\}^K, \sum_{k=1}^K x_k = 1$ (1-of- $K$ encoding)
Binomial	-	1	$x \in \{0, 1, \dots, n\}$
Bernoulli	1	1	$x \in \{0, 1\}$

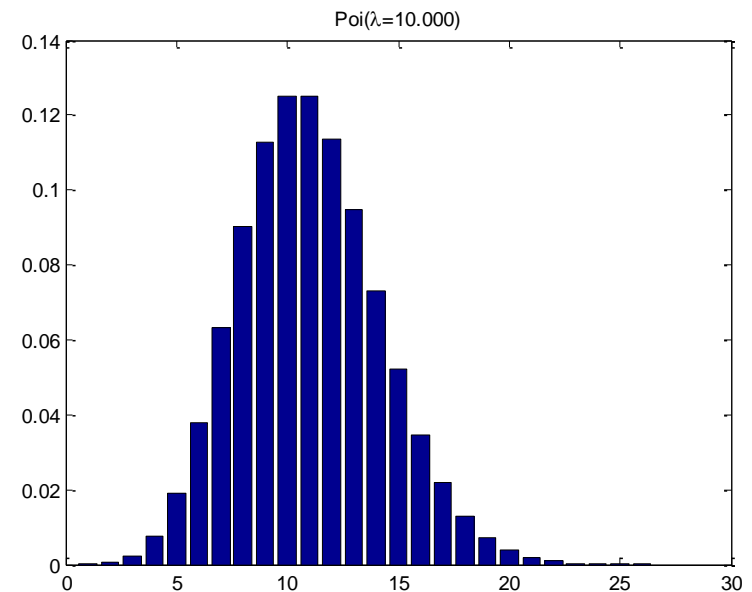
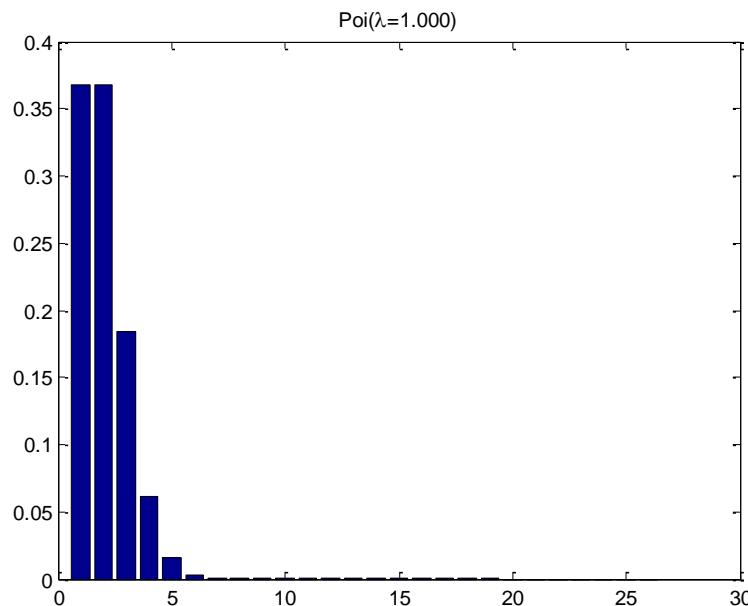
- $n = 1$  (one roll of the dice),  $n = -$  ( $N$  rolls of the dice)
- $K = 1$  (binary variables),  $K = -$  (1-of- $K$  encoding)

# The Poisson Distribution

- We say that  $X \in \{0, 1, 2, 3, \dots\}$  has a Poisson distribution with parameter  $\lambda > 0$ , if its pmf is

$$X \sim \mathcal{Poi}(\lambda) : \mathcal{Poi}(x | \lambda) = e^{-\lambda} \frac{\lambda^x}{x!}$$

- This is a model for **counts of rare events**.



Use MatLab function [poissonPlotDemo](#) from [Kevin Murphys' PMTK](#)

# The Empirical Distribution

- Given data,  $\mathcal{D} = \{x_1, \dots, x_N\}$ , we define the empirical distribution as:

$$p_{emp}(A) = \frac{1}{N} \sum_{i=1}^N \delta_{x_i}(A), \text{ Dirac Measure: } \delta_{x_i}(A) = \begin{cases} 1 & \text{if } x_i \in A \\ 0 & \text{if } x_i \notin A \end{cases}$$

- We can also associate weights with each sample:

$$\text{Generalize } p_{emp}(x) = \frac{1}{N} \sum_{i=1}^N \delta_{x_i}(x) \Rightarrow p_{emp}(x) = \sum_{i=1}^N w_i \delta_{x_i}(x), 0 \leq w_i \leq 1, \sum_{i=1}^N w_i = 1$$

- This corresponds to a histogram with spikes at each sample point with height equal to the corresponding weight. This distribution assigns zero weight to any point not in the dataset.
- Note that the “sample mean of  $f(x)$ ” is the expectation of  $f(x)$  under the empirical distribution:

$$\mathbb{E}[f(x)] = \int f(x) \sum_{i=1}^N \frac{1}{N} \delta_{x_i}(x) dx = \frac{1}{N} \sum_{n=1}^N f(x_i)$$

# Student's $\mathcal{T}$ Distribution

$$p(x | \mu, \lambda, \nu) = \frac{\Gamma(\frac{\nu}{2} + \frac{1}{2})}{\Gamma(\frac{\nu}{2})} \left( \frac{\lambda}{\pi \nu} \right)^{1/2} \left[ 1 + \frac{\lambda (x - \mu)^2}{\nu} \right]^{-\nu/2 - 1/2}$$

- The parameter  $\lambda$  is called the precision of the  $\mathcal{T}$ -distribution, even though it is not in general equal to the inverse of the variance (see below on behavior as  $\nu \rightarrow \infty$ ).
- The parameter  $\nu$  is called the degrees of freedom.
- For the particular case of  $\nu = 1$ , the  $\mathcal{T}$ -distribution reduces to the Cauchy distribution.
- In the limit  $\nu \rightarrow \infty$ , the  $\mathcal{T}$ -distribution  $\mathcal{T}(x|\mu, \lambda, \nu)$  becomes a Gaussian  $\mathcal{N}(x|\mu, \lambda^{-1})$  with mean  $\mu$  and precision  $\lambda$ .



# For $v \rightarrow \infty$ , $\mathcal{T}(x|\mu, \lambda, v)$ Becomes a Gaussian

$$p(x|\mu, \lambda, v) = \frac{\Gamma(\frac{v}{2} + \frac{1}{2})}{\Gamma(\frac{v}{2})} \left( \frac{\lambda}{\pi v} \right)^{1/2} \left[ 1 + \frac{\lambda(x - \mu)^2}{v} \right]^{-v/2 - 1/2}$$

➤ We first write the distribution as follows:

$$\mathcal{T}(x|\mu, \lambda, v) \propto \left[ 1 + \frac{\lambda(x - \mu)^2}{v} \right]^{-v/2 - 1/2} = \exp \left\{ -\frac{v+1}{2} \ln \left[ 1 + \frac{\lambda(x - \mu)^2}{v} \right] \right\}$$

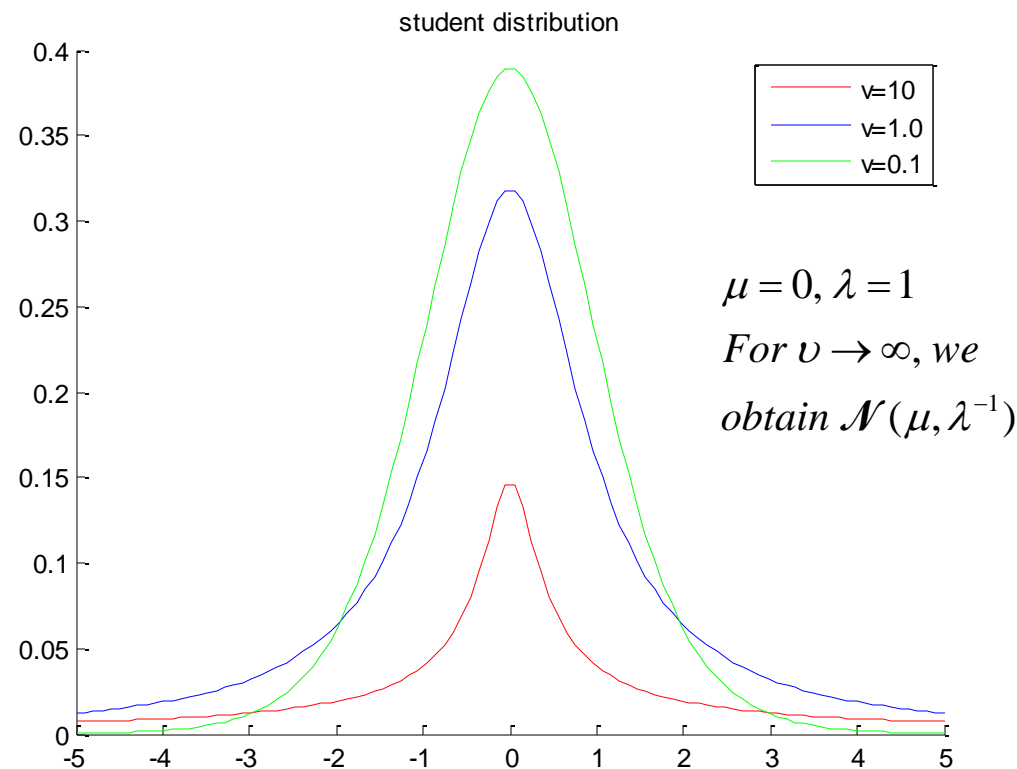
➤ For large  $v$ , we can approximate the log as follows:

$$\mathcal{T}(x|\mu, \lambda, v) \propto \exp \left\{ -\frac{v+1}{2} \left[ \frac{\lambda(x - \mu)^2}{v} + O(v^{-2}) \right] \right\} = \exp \left\{ -\frac{\lambda(x - \mu)^2}{2} + O(v^{-1}) \right\}$$

➤ In the limit  $v \rightarrow \infty$ , the  $\mathcal{T}$ -distribution  $\mathcal{T}(x|\mu, \lambda, v)$  is indeed a Gaussian  $\mathcal{N}(x|\mu, \lambda^{-1})$  with mean  $\mu$  and precision  $\lambda$ . The normalization of the  $\mathcal{T}$  is valid in this limit as well (so the Gaussian obtained is normalized).

# Student's $\mathcal{T}$ Distribution

$$p(x | \mu, \lambda, \nu) = \frac{\Gamma(\frac{\nu}{2} + \frac{1}{2})}{\Gamma(\frac{\nu}{2})} \left( \frac{\lambda}{\pi \nu} \right)^{1/2} \left[ 1 + \frac{\lambda (x - \mu)^2}{\nu} \right]^{-\nu/2 - 1/2}$$



$$\mu = 0, \lambda = 1$$

For  $\nu \rightarrow \infty$ , we

obtain  $\mathcal{N}(\mu, \lambda^{-1})$

MatLab Code

Mean:  $\mu, \nu > 1$   
Mode:  $\mu$   
Var:  $\frac{\nu}{\lambda(\nu - 2)}, \nu > 2$



# Student's $\mathcal{T}$ Vs the Gaussian

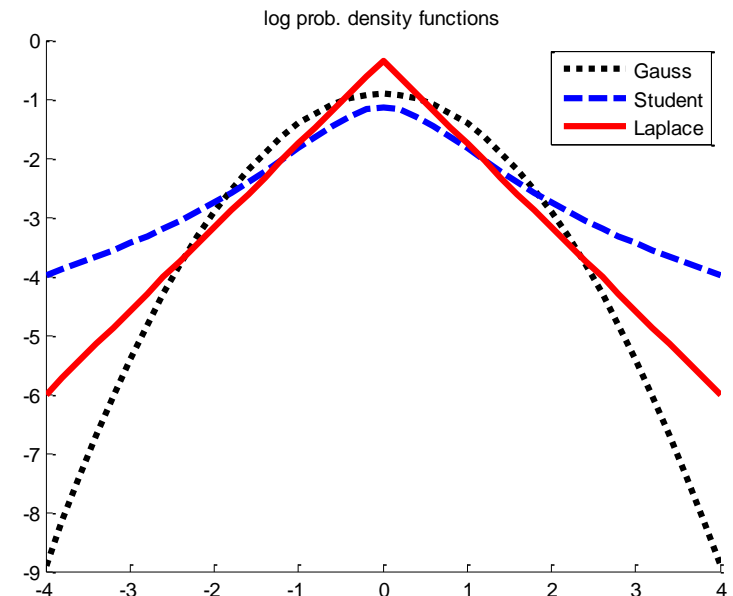
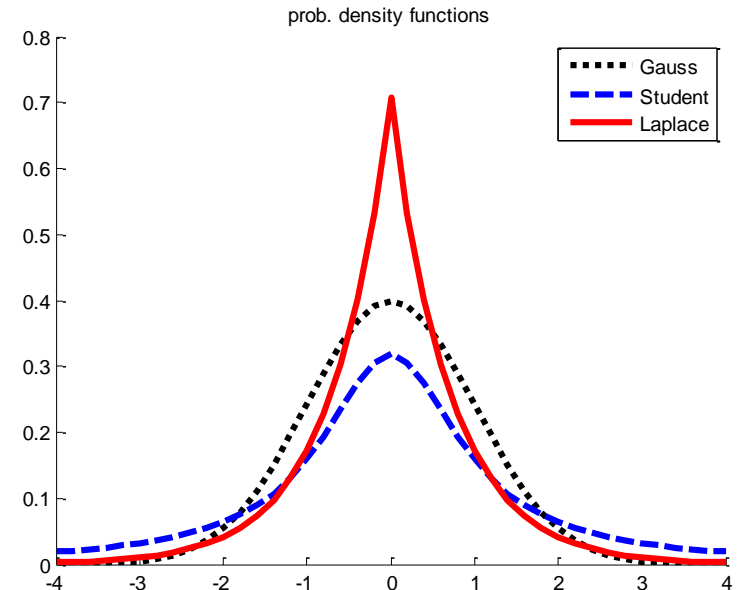
- We plot:

$$\mathcal{N}(x|0,1), \mathcal{T}(x|0,1,1), \mathcal{Lap}(x|0,1/\sqrt{2})$$

- *The mean and variance of the Student's is undefined for  $\nu = 1$ .*
- Logs of the PDFs. *The Student's is NOT log concave.*

Run MatLab function [studentLaplacePdfPlot](#)  
from [Kevin Murphys' PMTK](#)

- *When  $\nu = 1$ , the distribution is known as Cauchy or Lorentz. Due to its heavy tails, the mean does not converge.*
- ***Recommended to use  $\nu = 4$ .***





# Student's $\mathcal{T}$ Distribution

$$\begin{aligned} p(x | \mu, a, b) &= \int_0^{\infty} \mathcal{N}(x | \mu, \tau^{-1}) \textit{Gamma}(\tau | a, b) d\tau = \\ &= \int_0^{\infty} \left( \frac{\tau}{2\pi} \right)^{1/2} \exp\left( -\frac{\tau}{2} (x - \mu)^2 \right) \frac{b^a}{\Gamma(a)} \tau^{a-1} e^{-b\tau} d\tau \end{aligned}$$

- The Student's  $\mathcal{T}$  distribution can be seen from the equation above (see following two slides for proof) as an infinite **mixture of Gaussians each of them with different precision** (governed by a Gamma distribution)
- The result is a distribution that in general has longer 'tails' than a Gaussian.
- This gives the  $\mathcal{T}$ -distribution **robustness**, i.e. the  $\mathcal{T}$ -distribution is much less sensitive than the Gaussian to the presence of outliers.



# Appendix: Student's $\mathcal{T}$ as a Mixture of Gaussians

- If we have a univariate Gaussian  $\mathcal{T}(x|\mu, \tau^{-1})$  together with a prior  $\mathcal{Gamma}(\tau|a, b)$  and we integrate out the precision, we obtain the marginal distribution of  $x$

$$\begin{aligned} p(x|\mu, a, b) &= \int_0^{\infty} \mathcal{N}(x|\mu, \tau^{-1}) \mathcal{Gamma}(\tau|a, b) d\tau = \\ &= \int_0^{\infty} \left(\frac{\tau}{2\pi}\right)^{1/2} \exp\left(-\frac{\tau}{2}(x-\mu)^2\right) \frac{b^a}{\Gamma(a)} \tau^{a-1} e^{-b\tau} d\tau = \end{aligned}$$

- Introduce the transformation  $z = \underbrace{\left[b + \frac{1}{2}(x-\mu)^2\right]}_A \tau$  to simplify as:

$$\begin{aligned} p(x|\mu, a, b) &= \frac{b^a}{\Gamma(a)} \left(\frac{1}{2\pi}\right)^{1/2} \int_0^{\infty} \tau^{1/2} \exp(-z) \tau^{a-1} d\tau = \\ &= \frac{b^a}{\Gamma(a)} \left(\frac{1}{2\pi}\right)^{1/2} \frac{1}{A^{1/2+a-1+1}} \int_0^{\infty} z^{1/2} \exp(-z) z^{a-1} dz \end{aligned}$$

# Appendix: Student's $\mathcal{T}$ as a Mixture of Gaussians

$$\begin{aligned} p(x | \mu, a, b) &= \frac{b^a}{\Gamma(a)} \left( \frac{1}{2\pi} \right)^{1/2} \frac{1}{A^{1/2+a-1+1}} \int_0^\infty z^{1/2} \exp(-z) z^{a-1} dz \\ &= \frac{b^a}{\Gamma(a)} \left( \frac{1}{2\pi} \right)^{1/2} \left[ b + \frac{1}{2} (x - \mu)^2 \right]^{-a-1/2} \int_0^\infty \exp(-z) z^{a-1/2} dz \end{aligned}$$

➤ Recalling the [definition of the Gamma function](#):  $\Gamma(a) = \int_0^\infty \exp(-z) z^{a-1} dz$

$$p(x | \mu, a, b) = \frac{b^a}{\Gamma(a)} \left( \frac{1}{2\pi} \right)^{1/2} \left[ b + \frac{1}{2} (x - \mu)^2 \right]^{-a-1/2} \Gamma(a + \frac{1}{2})$$

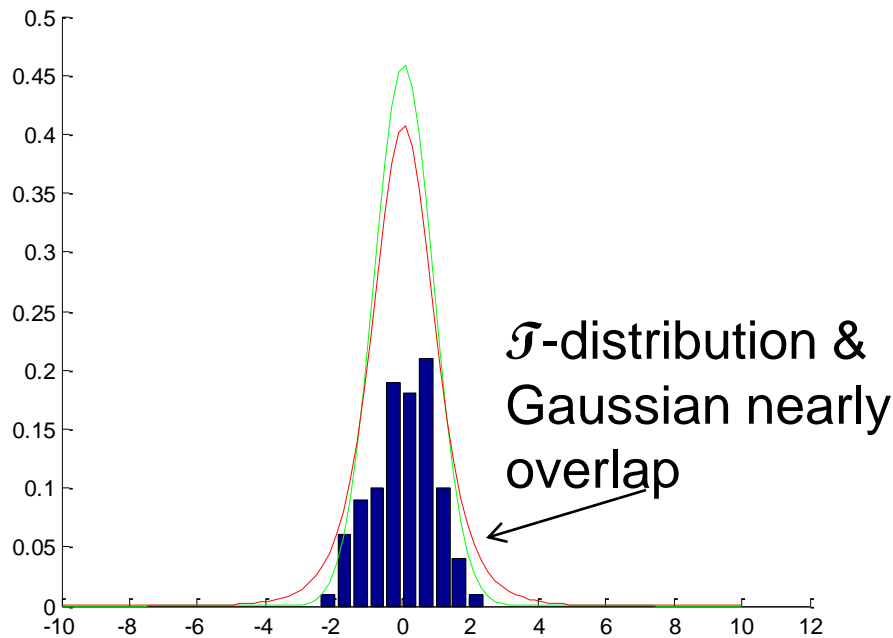
➤ It is common to redefine the parameters in this distribution as:  $\nu = 2a$ ,  $\lambda = \frac{a}{b}$

$$p(x | \mu, \lambda, \nu) = \frac{\Gamma(\frac{\nu}{2} + \frac{1}{2})}{\Gamma(\frac{\nu}{2})} \left( \frac{\lambda}{\pi\nu} \right)^{1/2} \left[ 1 + \frac{\lambda (x - \mu)^2}{\nu} \right]^{-\nu/2-1/2}$$

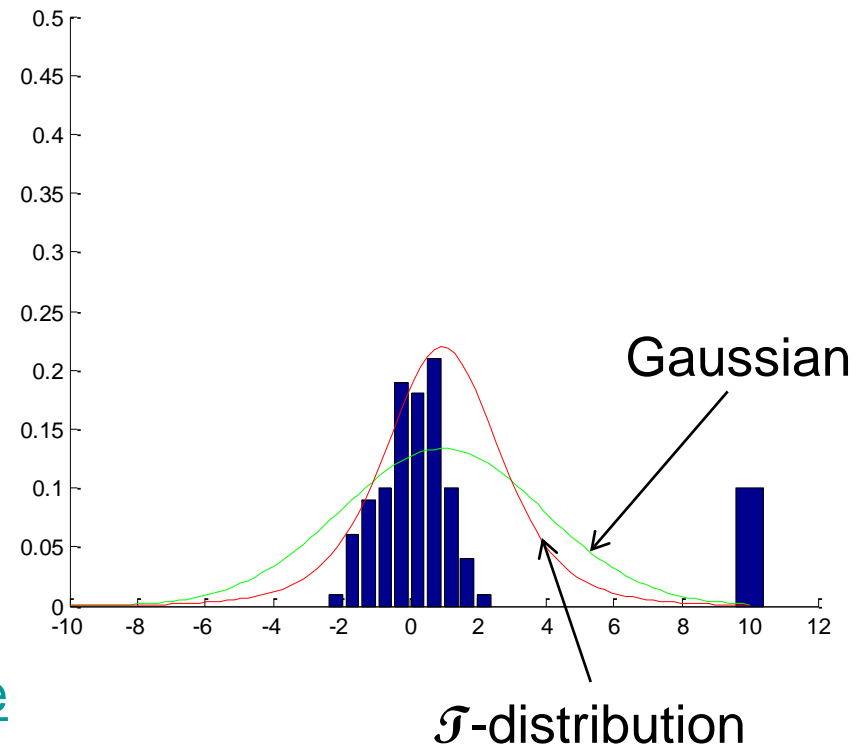


# Robustness of Student's $\mathcal{T}$ Distribution

- The robustness of the  $\mathcal{T}$ -distribution is illustrated here by comparing the “maximum likelihood solutions” for a Gaussian and a  $\mathcal{T}$ -distribution (30 data points from the Gaussian are used).
- The effect of a small number of outliers (Fig. on the right) is less significant for the  $\mathcal{T}$ -distribution than for the Gaussian.

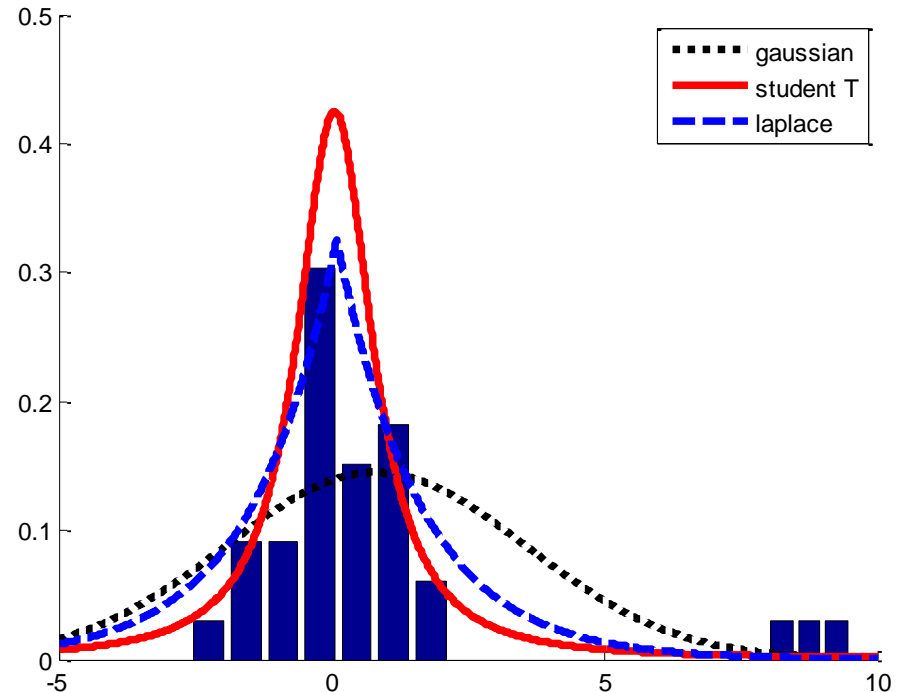
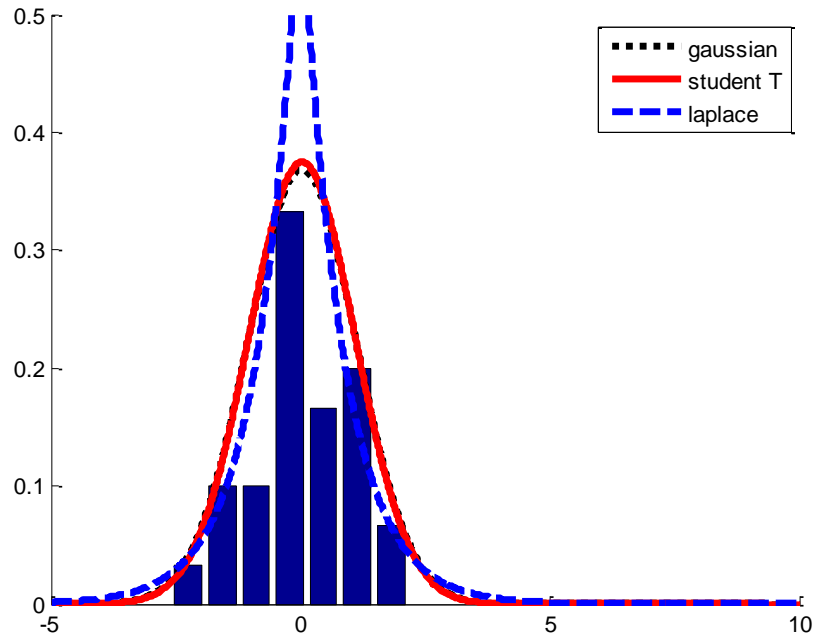


[MatLab Code](#)



# Robustness of Student's $\mathcal{T}$ Distribution

- The earlier simulation is repeated here with the PMTK toolbox.



Run MatLab function [\*robustDemo\*](#)  
from [Kevin Murphys' PMTK](#)



# The Laplace Distribution

---

- Another distribution with heavy tails is the [Laplace distribution](#), also known as the *double sided exponential distribution*. It has the following pdf:

$$\mathcal{Lap}(x | \mu, b) = \frac{1}{2b} e^{-\frac{|x-\mu|}{b}}$$

- $\mu$  is a location parameter and  $b > 0$  is a scale parameter

$$\text{Mean} = \mu, \text{Mode} = \mu, \text{Var} = 2b^2$$

- Its robust to outliers (see [earlier demonstration](#)).
- *It puts more probability density at 0 than the Gaussian. This property is a useful way to encourage sparsity in a model.*

# Beta Distribution

- The  $\text{Beta}(\alpha, \beta)$  distribution with  $x \in [0, 1], \alpha, \beta > 0$  is defined as follows:

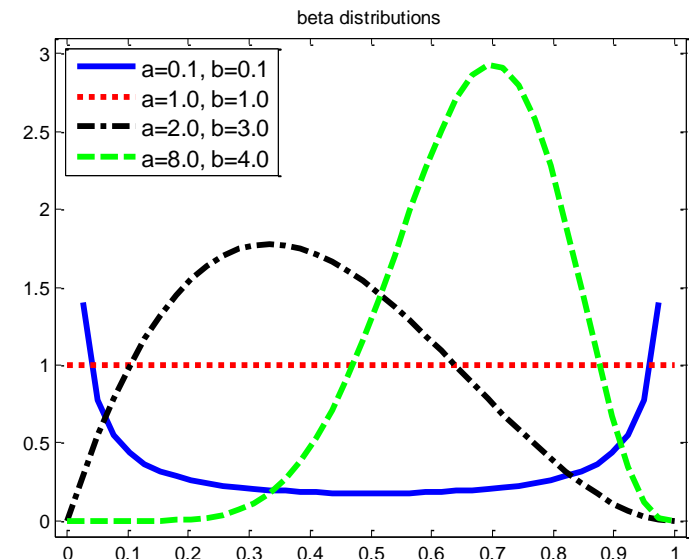
$$\text{Beta}(x) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1} = \underbrace{x^{\alpha-1} (1-x)^{\beta-1}}_{\text{Normalizing factor}}, \text{beta}(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha + \beta)}$$

- The expected value, mode and variance of a  $\text{Beta}$  random variable  $x$  with (hyper-) parameters  $\alpha$  and  $\beta$  :

$$\mathbb{E}[x] = \frac{\alpha}{\alpha + \beta}, \quad \text{mode}[x] = \frac{a - 1}{a + b - 2}$$

$$\text{var}[x] = \frac{\alpha\beta}{(\alpha + \beta)^2 (\alpha + \beta + 1)}$$

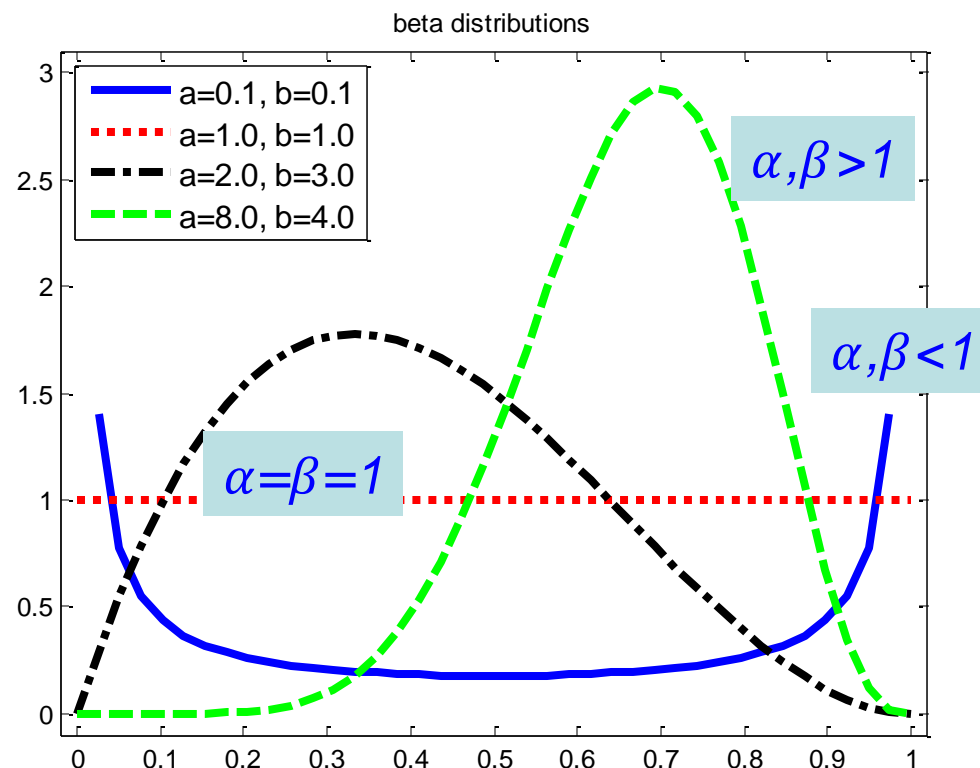
- For more information visit [this link](#).



# Beta Distribution

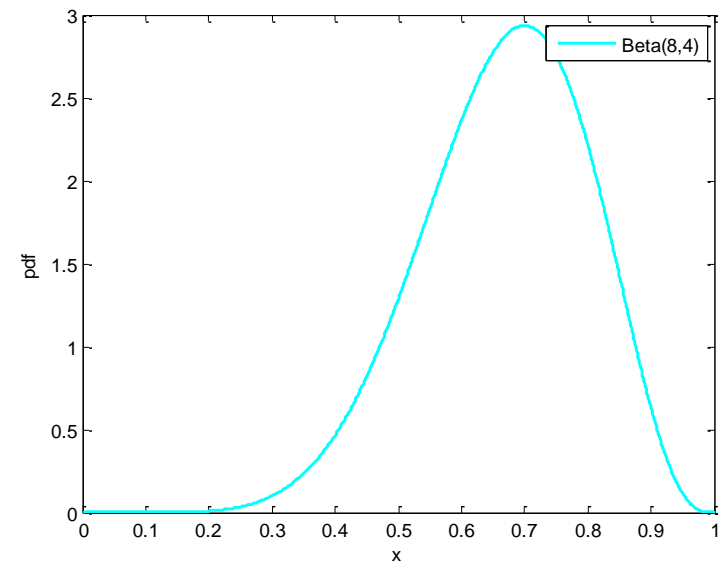
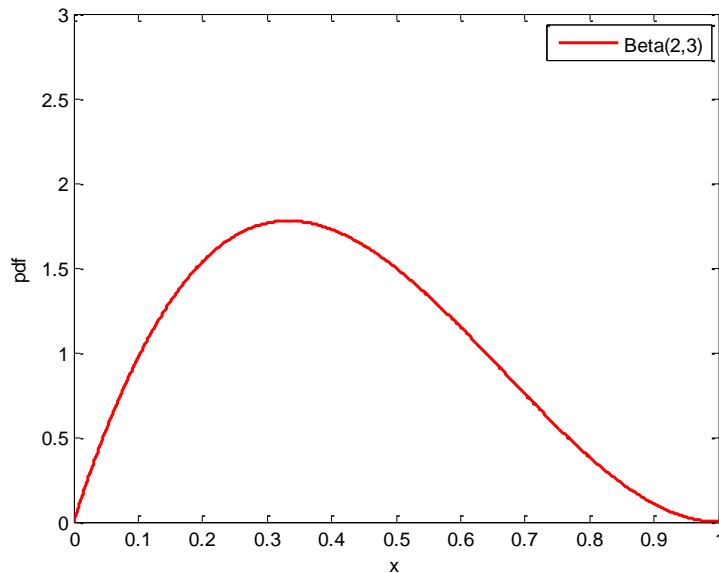
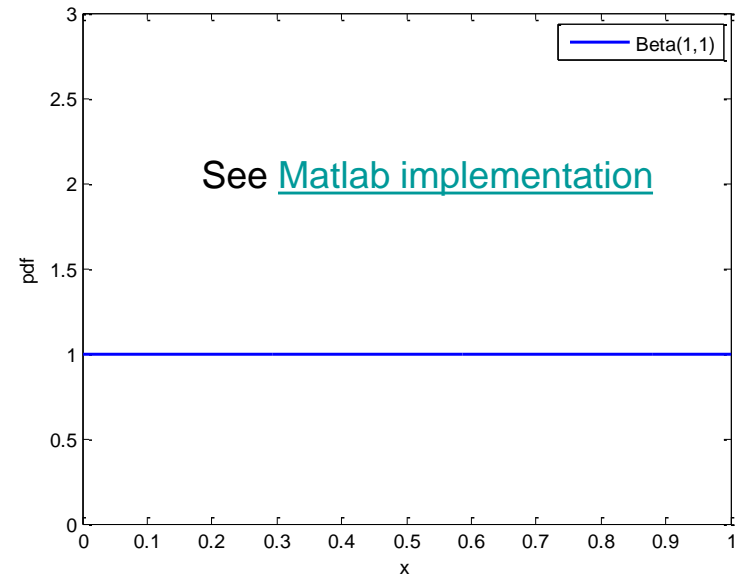
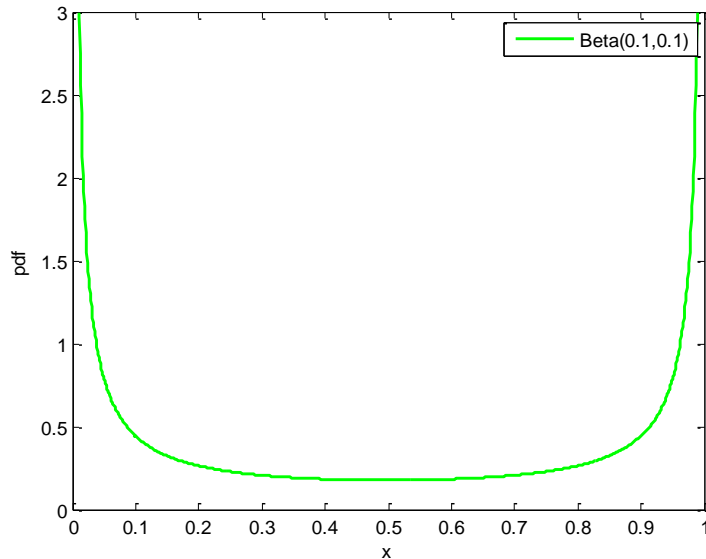
- If  $\alpha = \beta = 1$ , we obtain a uniform distribution.
- If  $\alpha$  and  $\beta$  are both less than 1, we get a bimodal distribution with spikes at 0 and 1.
- If  $\alpha$  and  $\beta$  are both greater than 1, the distribution is unimodal.

Run [betaPlotDemo](#)  
from [PMTK](#)





# Beta Distribution



# Gamma Function

---

$$\text{Beta}(x) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1}$$

- The gamma function extends the factorial to real numbers:

$$\Gamma(x) = \int_0^{\infty} u^{x-1} e^{-u} du$$

- With integration by parts:

$$\Gamma(x+1) = x\Gamma(x)$$

- For integer  $n$ :

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

- For more information visit [this link](#).

# Beta Distribution: Normalization

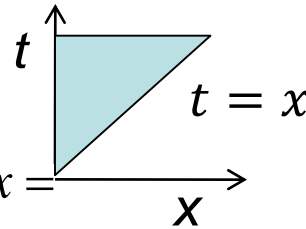
- Showing that the  $\mathcal{Beta}(\alpha, \beta)$  distribution is normalized correctly is a bit tricky. We need to prove that:

$$\Gamma(\alpha)\Gamma(\beta) = \Gamma(\alpha + \beta) \int_0^1 \mu^{\alpha-1} (1-\mu)^{\beta-1} d\mu$$

Follow the steps:

- ✓ (a) change the variable  $y$  below to  $y = t - x$ ;
- ✓ (b) change the order of integration in the shaded triangular region;
- ✓ and (c) change  $x$  to  $\mu$  via  $x = t\mu$ :

$$\begin{aligned} \Gamma(\alpha)\Gamma(\beta) &= \left( \int_0^\infty x^{\alpha-1} e^{-x} dx \right) \left( \int_0^\infty y^{\beta-1} e^{-y} dy \right) \underset{y=t-x}{=} \int_0^\infty x^{\alpha-1} \left( \int_x^\infty e^{-t} (t-x)^{\beta-1} dt \right) dx \\ &= \int_0^\infty \left( \int_0^t x^{\alpha-1} e^{-t} (t-x)^{\beta-1} dx \right) dt = \int_0^\infty \underset{t}{t}^{\alpha-1} e^{-\underset{t}{t}} \underset{t}{t}^{\beta-1} \underset{t}{t} dt \int_0^1 \underset{\mu}{\mu}^{\alpha-1} (1-\underset{\mu}{\mu})^{\beta-1} d\mu = \\ &= \Gamma(\alpha + \beta) \int_0^1 \mu^{\alpha-1} (1-\mu)^{\beta-1} d\mu \end{aligned}$$



# Gamma Distribution- Rate Parametrization

- It is frequently a model for waiting times. For important properties [see here](#).
- It is more often parameterized in terms of a shape parameter  $a$  and an inverse scale parameter  $b = 1/\theta$ , called a *rate parameter*:

$$p(x | a, b) = \frac{b^a}{\Gamma(a)} x^{a-1} e^{-bx}, x \in [0, \infty], \Gamma(a) = \int_0^{\infty} u^{a-1} e^{-u} du$$

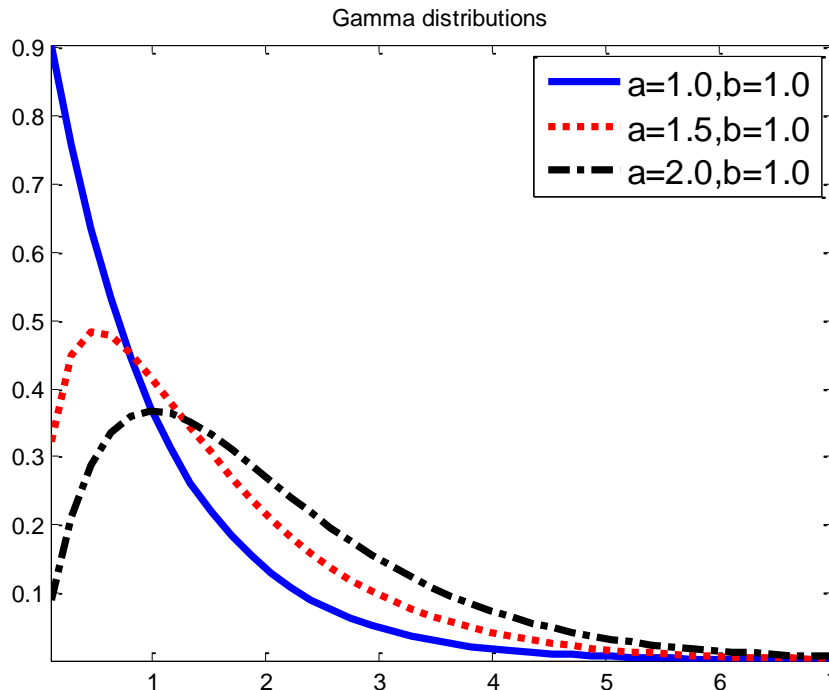
- The mean, mode and variance with this parametrization are:

$$\mathbb{E}[x] = \frac{\alpha}{b} \quad \text{mode}[x] = \begin{cases} \frac{a-1}{b}, & \text{for } a > 1 \\ 0 & \text{otherwise} \end{cases} \quad \text{var}[x] = \frac{\alpha}{b^2}$$

# Gamma Distribution

➤ Plots of  $\text{Gamma}(X | a, b) = \frac{b^a}{\Gamma(a)} x^{a-1} \exp(-xb), b = 1$

➤ As we increase the rate  $b$ , the distribution squeezes leftwards and upwards. For  $a < 1$ , the mode is at zero.

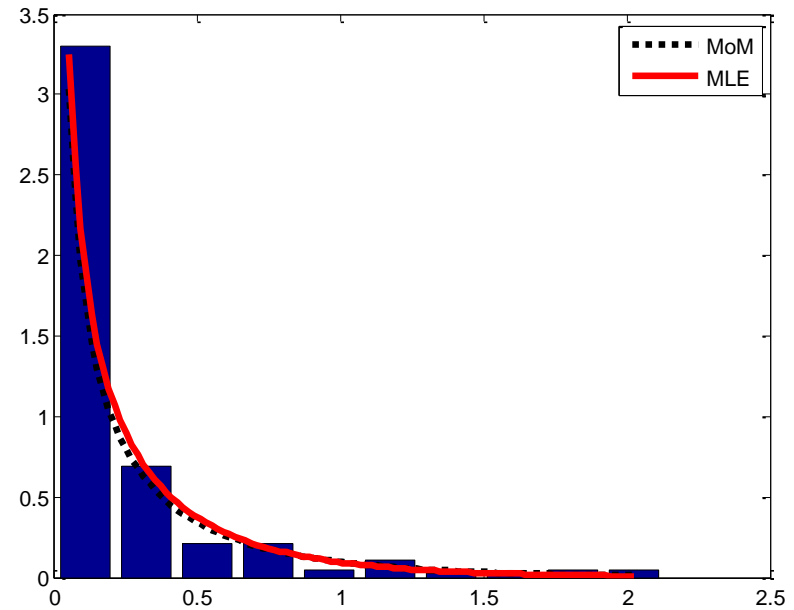
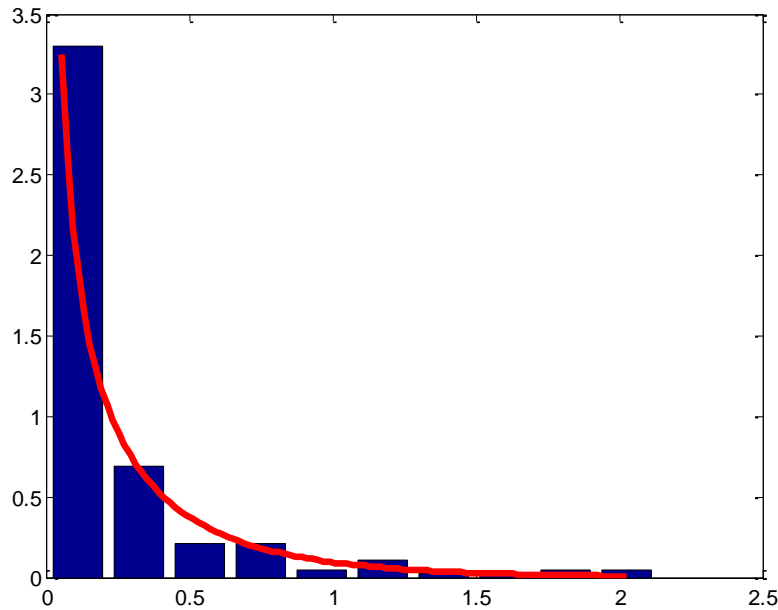


Run [gammaPlotDemo](#)  
from [PMTK](#)



# Gamma Distribution

- An empirical PDF of rainfall data fitted with a Gamma distribution.



Run MatLab function [gammaRainfallDemo](#)  
from [PMTK](#)

# Exponential Distribution

---

- This is defined as

$$\textit{Expon}(X \mid \lambda) = \textit{Gamma}(X \mid 1, \lambda) = \lambda \exp(-x\lambda), x \in [0, \infty]$$



- Here  $\lambda$  is the rate parameter.

- This *distribution describes the times between events in a Poisson process, i.e. a process in which events occur continuously and independently at a constant average rate  $\lambda$ .*

# Chi-Squared Distribution

- This is defined as

$$\chi^2(X | \nu) = \textit{Gamma}(X | \frac{\nu}{2}, \frac{1}{2}) = \frac{\left(\frac{1}{2}\right)^{\frac{\nu}{2}}}{\Gamma\left(\frac{\nu}{2}\right)} x^{\frac{\nu}{2}-1} \exp\left(-\frac{x}{2}\right), x \in [0, \infty]$$

- *This is the distribution of the sum of squared Gaussian random variables.*
- More precisely,

$$\textit{Let } Z_i \sim \mathcal{N}(0,1) \textit{ and } S = \sum_{i=1}^{\nu} Z_i^2, \textit{ then : } S \sim \chi_{\nu}^2$$



# Inverse Gamma Distribution

- This is defined as follows:

$$\text{If } X \sim \text{Gamma}(X | a, b) \Rightarrow X^{-1} \sim \text{InvGamma}(X | a, b)$$

where:

$$\text{InvGamma}(X | a, b) = \frac{b^a}{\Gamma(a)} x^{-(a+1)} \exp(-b / x), x \in [0, \infty]$$

- $a$  is the shape and  $b$  the scale parameters.

- Note that  $b$  is a scale parameter since:

$$\text{InvGamma}(X | a, b) = \frac{\text{InvGamma}(\frac{X}{b} | a, 1)}{b}$$

- It can be shown that:

$$\text{Mean} = \frac{b}{a-1} \text{ (exists for } a > 1), \text{ Mode} = \frac{b}{a+1}, \text{ var} = \frac{b^2}{(a-1)^2(a-2)} \text{ (exists for } a > 2)$$



# The Pareto Distribution

- Used to model the distribution of quantities that exhibit long tails (heavy tails)

$$\text{Pareto}(X|k, m) = km^k x^{-(k+1)} \mathbb{I}(x \geq m)$$

- This density asserts that  $x$  must be greater than some constant  $m$ , but not too much greater,  $k$  controls what is “too much”.
- Modeling the frequency of words vs. their rank (e.g. “the”, “of”, etc.) or the wealth of people.\*
- As  $k \rightarrow \infty$ , the distribution approaches  $\delta(x - m)$ .
- On a log-log scale, the pdf forms a straight line of the form  $\log p(x) = a \log x + c$  for some constants  $a$  and  $c$  (power law, Zipf's law).

\* Basis of the distribution: a high proportion of a population has low income and only few have very high incomes.



# The Pareto Distribution

- Applications: Modeling the frequency of words vs their rank, distribution of wealth ( $k$  =Pareto Index), etc.

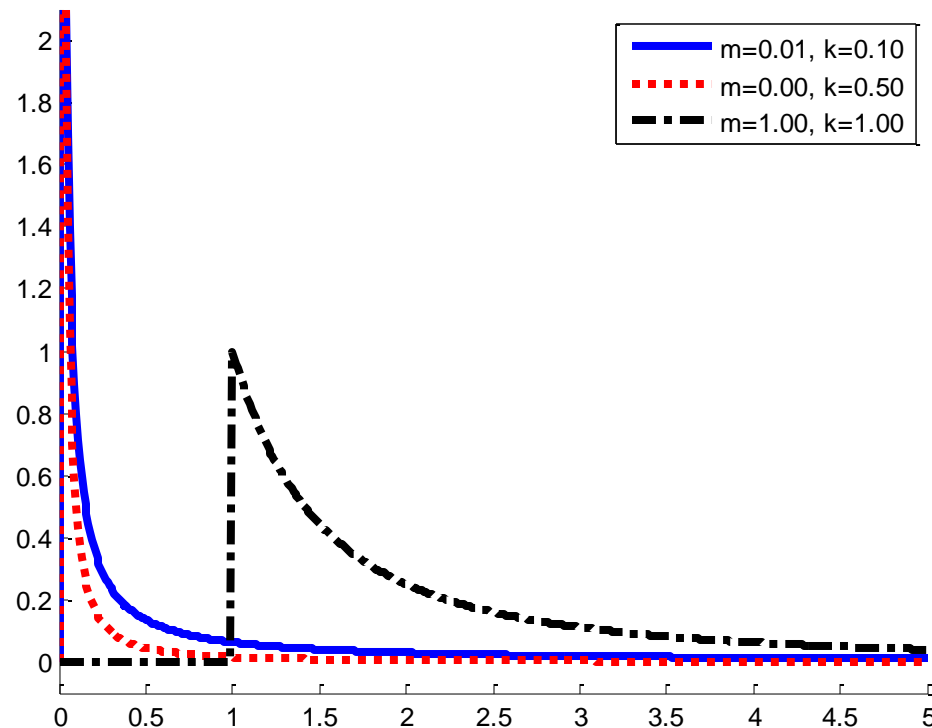
$$\mathcal{P}areto(X | k, m) = km^k x^{-(k+1)} \mathbb{I}(x \geq m),$$

$$Mean = \frac{km}{k-1} \text{ (if } k > 1),$$

$$Mode = m,$$

$$var = \frac{m^2 k}{(k-1)^2 (k-2)} \text{ (if } k > 2)$$

Pareto distribution



ParetoPlot from PMTK



# Covariance

➤ Consider two random variables  $X, Y : \Omega \rightarrow \mathbb{R}$ .

➤ The joint probability distribution is defined as:

$$P\{X \in A, Y \in B\} = P\{X^{-1}(A) \cap Y^{-1}(B)\} = \iint_{A \times B} p(x, y) dx dy$$

➤ Two random variables are independent if

$$p(x, y) = p(x)p(y)$$

➤ The covariance of  $X$  and  $Y$  is defined as:

$$\text{cov}(X, Y) = \mathbb{E}[(X - \mathbb{E}(X))(Y - \mathbb{E}(Y))]$$

➤ It is straight forward to verify that:  $\text{cov}(X, Y) = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]$

# Correlation, Center Normalized Random Variables

- Consider two random variables  $X, Y : \Omega \rightarrow \mathbb{R}$ .
- The correlation coefficient of  $X$  and  $Y$  is defined as:

$$\text{corr}(X, Y) = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y}$$

where the standard deviations of  $X$  and  $Y$  are

$$\sigma_X = \sqrt{\text{cov}(X)}, \sigma_Y = \sqrt{\text{cov}(Y)}$$

- The center normalized random variables are defined as: 
$$\begin{cases} \tilde{X} = \frac{X - \mathbb{E}[X]}{\sigma_X} \\ \tilde{Y} = \frac{Y - \mathbb{E}[Y]}{\sigma_Y} \end{cases}$$
- It is straight forward to verify that:

$$\mathbb{E}[\tilde{X}] = \mathbb{E}[\tilde{Y}] = 0 \quad \text{var}[\tilde{X}] = \text{var}[\tilde{Y}] = 1$$