
Introduction to Probability and Statistics (Continued)

Prof. Nicholas Zabaras

Center for Informatics and Computational Science

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

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

Email: nzabaras@gmail.com

URL: <https://www.zabaras.com/>

August 28, 2018



Contents

- Covariance, Uncorrelated Random Variables, Multivariate Random Variables, Independence Vs Uncorrelated Random Variables
- Marginal and Conditional Densities, Conditional Expectation
- The multivariate Gaussian, Multivariate Student t distribution
- Transformations of random variables
- Dirichlet distribution



References

- Following closely [Chris Bishop's 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.



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:

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

- 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:

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



Variance and Covariance

□ Variance

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

□ Covariance

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

It expresses the extent to which X and Y vary (linearly) together.

- If X and Y are independent, $p(X, Y) = p(X)p(Y)$, their covariance vanishes.



Uncorrelated Random Variables

- Consider two random variables $X, Y : \Omega \rightarrow \mathbb{R}$.
- We say that X and Y are uncorrelated when: $\text{cov}(X, Y) = 0 \Rightarrow \mathbb{E}[XY] = \mathbb{E}[X]\mathbb{E}[Y]$
- If X and Y are independent, then they are uncorrelated:

$$\text{cov}(X, Y) = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] = \mathbb{E}[(X - \mathbb{E}[X])]\mathbb{E}[(Y - \mathbb{E}[Y])] = 0$$

The opposite is not true: Uncorrelated random variables are not independent. *Independency affects the whole density, not just the expectation.*

- X and Y are orthogonal if

$$\mathbb{E}[XY] = 0$$

In the last case, the following holds:

$$\mathbb{E}[(X + Y)^2] = \mathbb{E}[X^2] + \mathbb{E}[Y^2]$$



Multivariate Random Variables

- Consider

$$X = \begin{bmatrix} X_1 \\ X_2 \\ .. \\ X_n \end{bmatrix} : \Omega \rightarrow \mathbb{R}^n$$

where each component X_i is an \mathbb{R} - valued variable.

- X is defined by the joint probability density of its components

$$p_X : \mathbb{R}^n \rightarrow \mathbb{R}^+$$

- Define the **cumulative distribution function** is defined as:

$$F(x_1, x_2, \dots, x_n) = \Pr[X_1 < x_1, X_2 < x_2, \dots, X_n < x_n] \in [0, 1]$$

- Then the **probability density function** of X is defined as

$$p(x_1, x_2, \dots, x_n) = \frac{\partial^n F(\mathbf{x})}{\partial x_1 \partial x_2 \dots \partial x_n} \text{ and } \int p(\mathbf{x}) d\mathbf{x} = 1$$



Multivariate Random Variables

- Consider

$$X = \begin{bmatrix} X_1 \\ X_2 \\ .. \\ X_n \end{bmatrix} : \Omega \rightarrow \mathbb{R}^n$$

where each component X_i is an \mathbb{R} -valued variable.

- The expectation is defined as

$$\mathbb{E}[X] = \int_{\mathbb{R}^n} xp(x)dx \in \mathbb{R}^n, \quad or \quad \mathbb{E}[X_i] = \int_{\mathbb{R}^n} x_i p(x)dx = \int_{\mathbb{R}} x_i p(x_i)dx_i \in \mathbb{R}, \quad i = 1, 2, \dots, n$$



Covariance Matrix

- Consider

$$X = \begin{bmatrix} X_1 \\ X_2 \\ .. \\ X_n \end{bmatrix} : \Omega \rightarrow \mathbb{R}^n$$

- The covariance matrix is:

$$\text{cov}[X] = \int_{\mathbb{R}^n} (x - \mathbb{E}[X])(x - \mathbb{E}[X])^T p(x) dx \in \mathbb{R}^{n \times n}$$

or equivalently: $\text{cov}[X]_{ij} = \int_{\mathbb{R}^n} (x_i - \mathbb{E}[X_i])(x_j - \mathbb{E}[X_j]) p(x) dx \in \mathbb{R}, 1 \leq i, j \leq n.$

- The covariance matrix is symmetric and positive semi-definite, i.e.

$$\forall v \in \mathbb{R}^n, v \neq 0, \quad v^T \text{cov}[X] v = \int_{\mathbb{R}^n} [v^T (x - \bar{x})][(x - \bar{x})^T v] p(x) dx = \int_{\mathbb{R}^n} [v^T (x - \bar{x})]^2 p(x) dx \geq 0.$$

- Note that the diagonal of the covariance matrix gives the variances of the individual components:

$$\text{cov}[X]_{ii} = \int_{\mathbb{R}^n} (x_i - \mathbb{E}[X_i])^2 p(x) dx = \int_{\mathbb{R}^n} (x_i - \mathbb{E}[X_i])^2 \int_{\mathbb{R}^{n-1}} p(x_i, x_{i+1}^n) dx_{i+1}^n dx_i = \int_{\mathbb{R}^n} (x_i - \mathbb{E}[X_i])^2 p(x_i) dx_i = \text{var}[X_i]$$



Covariance Matrix

- The covariance matrix of a vector X can be written explicitly:

$$\text{cov}[X] = \begin{pmatrix} \text{var}[X_1] & \text{cov}[X_1, X_2] & \dots & \text{cov}[X_1, X_d] \\ & \text{var}[X_2] & \dots & \\ \text{cov}[X_d, X_1] & \text{cov}[X_d, X_2] & \dots & \text{var}[X_d] \end{pmatrix}$$

- A normalized version of this is *the correlation matrix* (all elements between $[-1,1]$ (diagonal elements = 1)

$$R = \begin{pmatrix} \text{corr}[X_1, X_1] & \text{corr}[X_1, X_2] & \dots & \text{corr}[X_1, X_d] \\ & \text{corr}[X_2, X_2] & \dots & \\ \text{corr}[X_d, X_1] & \text{corr}[X_d, X_2] & \dots & \text{corr}[X_d, X_d] \end{pmatrix}$$

$$\text{corr}[X, Y] = \frac{\text{cov}(X, Y)}{\sqrt{\text{var}(X) \text{var}(Y)}}$$



Correlation Coefficient Between -1 and 1

- Consider two scalar random variables X and Y . We can write the following:

$$\begin{aligned} 0 \leq \text{var}\left[\frac{X}{\sigma_X} + \frac{Y}{\sigma_Y}\right] &= \mathbb{E}\left[\left(\frac{X}{\sigma_X} + \frac{Y}{\sigma_Y}\right)^2\right] - \left(\mathbb{E}\left[\frac{X}{\sigma_X} + \frac{Y}{\sigma_Y}\right]\right)^2 \\ &= \frac{\text{var}[X]}{\sigma_X^2} + \frac{\text{var}[Y]}{\sigma_Y^2} + 2 \frac{\text{cov}[X, Y]}{\sigma_X \sigma_Y} = 1 + 1 + 2\text{Corr}[X, Y] \\ \Rightarrow \text{Corr}[X, Y] &\geq -1 \end{aligned}$$

- Similarly starting with

$$0 \leq \text{var}\left[\frac{X}{\sigma_X} - \frac{Y}{\sigma_Y}\right] \Rightarrow \text{Corr}[X, Y] \leq 1$$

$$\text{corr}[X, Y] = \frac{\text{cov}(X, Y)}{\sqrt{\text{var}(X) \text{var}(Y)}}$$



Variance and Covariance

- The covariance of the vector random variables is:

$$\text{cov}[\mathbf{X}, \mathbf{Y}] = \mathbb{E}_{\mathbf{X}, \mathbf{Y}} \left[(\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{Y} - \mathbb{E}[\mathbf{Y}])^T \right] = \mathbb{E}_{\mathbf{X}, \mathbf{Y}} [\mathbf{XY}^T] - \mathbb{E}[\mathbf{X}] \mathbb{E}[\mathbf{Y}^T]$$

- The covariance between the components of a vector:

$$\text{cov}[\mathbf{X}, \mathbf{X}] = \mathbb{E}_{\mathbf{X}, \mathbf{X}} [\mathbf{XX}^T] - \mathbb{E}[\mathbf{X}] \mathbb{E}[\mathbf{X}^T]$$



Correlation as a Degree of Linearity

- It can be shown that

If $Y = aX + b, a > 0$, then : $\text{corr}[X, Y] = +1$

If $Y = aX + b, a < 0$, then : $\text{corr}[X, Y] = -1$

- The regression coefficient is $a = \text{cov}[X, Y] / \text{var}[X]$.
- *Think of the correlation coefficient as a degree of linearity.*
- If X and Y are independent, $p(X, Y) = p(X)p(Y)$, then $\text{cov}[X, Y] = 0$, and hence $\text{corr}[X, Y] = 0$ so they are uncorrelated.
- The converse is not true: uncorrelated does not imply independence.



Independent vs Uncorrelated

- Note that:

$$\begin{aligned}\text{var}[X + Y] &= \mathbb{E}[(X + Y)^2] - (\mathbb{E}[X] + \mathbb{E}[Y])^2 = \\ &= \mathbb{E}[X^2] - (\mathbb{E}[X])^2 + \mathbb{E}[Y^2] - (\mathbb{E}[Y])^2 + 2(\mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]) = \\ &= \text{var}[X] + \text{var}[Y] + 2\text{cov}[X, Y]\end{aligned}$$

- From the above equation, we note that if X, Y are independent then:

$$\text{var}[X + Y] = \text{var}[X] + \text{var}[Y]$$

but note that the linearity of expectation is valid even when the variables are not independent:

$$\mathbb{E}[X + Y] = \mathbb{E}[X] + \mathbb{E}[Y]$$

Correlation and Dependence

- Uncorrelated does not imply independent.
- For example, let $X \sim U(-1, 1)$ and $Y = X^2$. Clearly Y is dependent on X , yet one can show that $\text{corr}[X, Y] = 0$.

$$\mathbb{E}[X] = \frac{-1+1}{2} = 0, \quad \text{var}[X] = \frac{(1-(-1))^2}{12} = \frac{1}{3}$$

$$\mathbb{E}[Y] = \mathbb{E}[X^2] = \text{var}[X] + (\mathbb{E}[X])^2 = \frac{1}{3} + 0^2 = \frac{1}{3}$$

$$\mathbb{E}[XY] = \int_{-1}^1 x^3 p(x) dx = \int_{-1}^1 x^3 \frac{1}{2} dx = 0$$

$$\text{Corr}[X, Y] = \frac{\text{cov}[X, Y]}{\sigma_X \sigma_Y} = \frac{\mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]}{\sigma_X \sigma_Y} = \frac{0 - 0 \times 1/3}{\sigma_X \sigma_Y} = 0$$



Mutual Information

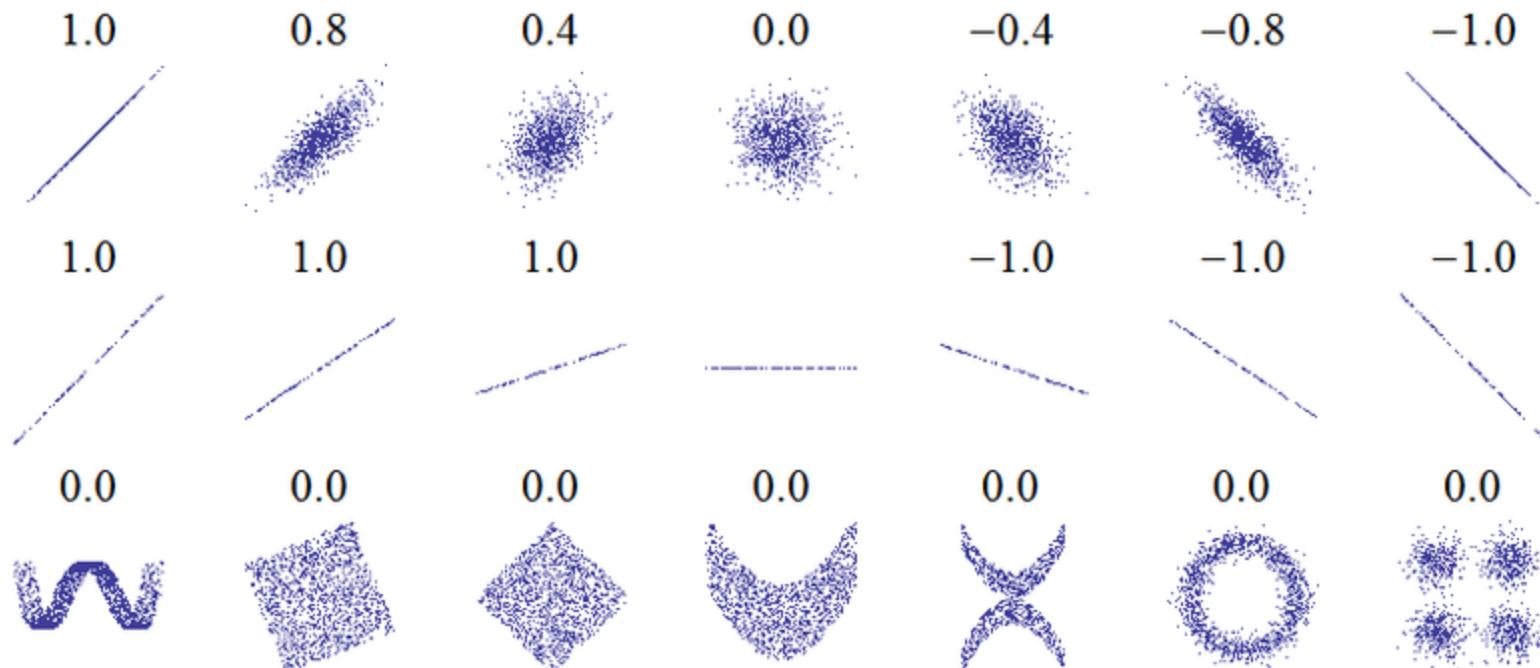
- The Figure given next shows several data sets where there is clear dependence between X and Y , and yet the correlation coefficient is 0.
- A *more general measure of dependence between random variables is mutual information.*

The mutual information is zero if and only if the variables are truly independent.



Uncorrelated Random Variables

- Several sets of (x, y) points, with the correlation coefficient of x and y for each set.
- The correlation reflects the *noisiness and direction of a linear relationship* (top row), but not the slope of that relationship (middle), nor nonlinear relationships (bottom).
- The figure in the center has a slope of 0 but the correlation coefficient is undefined because $\text{var}[Y] = 0$.



Marginal Density

- Consider two random variables $X, Y : \Omega \rightarrow \mathbb{R}$ with joint probability density $p(x, y)$
- The probability density of X when Y can take any value is defined as:

$$p(x) = \int_{\mathbb{R}} p(x, y) dy$$

- Similarly:

$$p(y) = \int_{\mathbb{R}} p(x, y) dx.$$



Conditional Probability Density

- Consider two random variables $X, Y : \Omega \rightarrow \mathbb{R}$ with joint density $p(x, y)$
- The probability density of X assuming that $Y = y$ is defined as

$$p(x | y) = \frac{p(x, y)}{p(y)}, p(y) \neq 0$$

- One can show this by noting the following:

$$P(a \leq X \leq b | y - \varepsilon \leq Y \leq y + \varepsilon) = \frac{\int_{y-\varepsilon}^{y+\varepsilon} \int_a^b p(x, y) dx dy}{\int_{y-\varepsilon}^{y+\varepsilon} p(y) dy} \approx \frac{\int_a^b 2\varepsilon p(x, y) dx}{2\varepsilon p(y)} = \int_a^b \underbrace{\frac{p(x, y)}{p(y)}}_{p(x|Y=y)} dx$$

- From this we derive the following important identity:

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

- **Bayes' rule in terms of densities** now follows as:

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



Conditional Expectations

- Consider two random variables $X, Y : \Omega \rightarrow \mathbb{R}$.
- We define the conditional expectation as: $\mathbb{E}[X | y] = \int_{\mathbb{R}} xp(x | y)dx$
- The expectation of X via conditional expectation can be computed as:

$$\mathbb{E}[X] = \int xp(x)dx = \int x \left(\underbrace{\int p(x, y) dy}_{p(x|y)p(y)} \right) dx \Rightarrow$$

$$\mathbb{E}[X] = \int \left(\int xp(x | y)dx \right) p(y)dy = \int \mathbb{E}[X | y]p(y)dy \Rightarrow$$

$$\boxed{\mathbb{E}[X] = \int \mathbb{E}[X | y]p(y)dy}$$



Linear Transformations

- Suppose $y = f(x) = Ax + b$. You can show that:

$$\mathbb{E}[y] = A\mathbb{E}[x] + b$$

$$\text{cov}[y] = A \text{cov}[x] A^T$$

- For a scalar-valued function $y = f(x) = a^T x + b$:

$$\mathbb{E}[y] = a^T \mathbb{E}[x] + b$$

$$\text{var}[y] = a^T \text{cov}[x] a$$



Multivariate Gaussian

- A random variable $X \in \mathbb{R}$ is Gaussian or normally distributed $X \sim \mathcal{N}(x_0, \sigma^2)$ if:

$$P\{X \leq t\} = \frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^t \exp\left(-\frac{1}{2\sigma^2}(x - x_0)^2\right) dx$$

- A multivariate $X \in \mathbb{R}^D$ is Gaussian if its probability density is

$$p(x) = \left(\frac{1}{(2\pi)^D \det \Sigma} \right)^{1/2} \exp\left(-\frac{1}{2}(x - x_0)^T \Sigma^{-1}(x - x_0)\right)$$

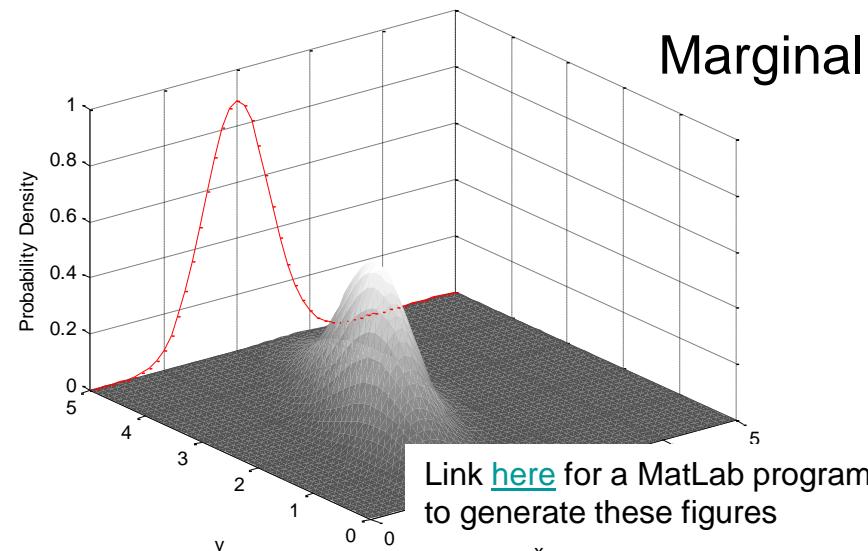
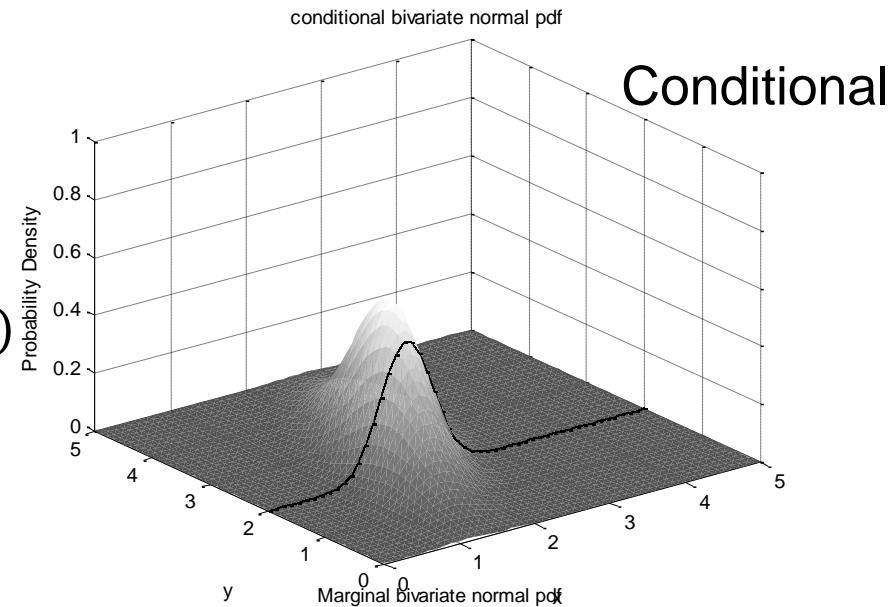
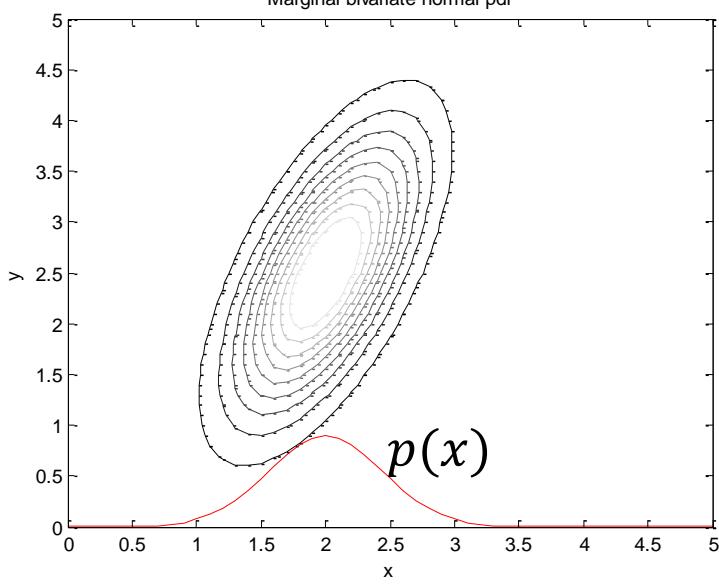
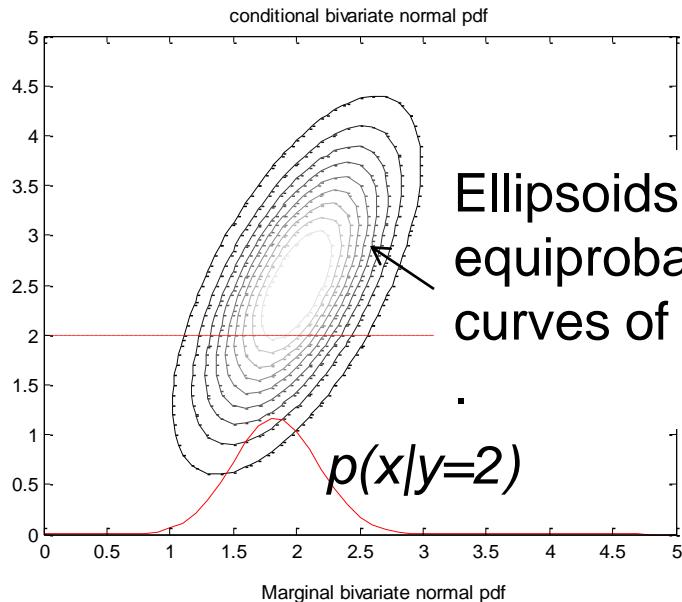
where $x_0 \in \mathbb{R}^D, \Sigma \in \mathbb{R}^{D \times D}$ is *symmetric positive definite (covariance matrix)*.

- The symmetry property of the covariance matrix does not affect the value of $(x - x_0)^T \Sigma^{-1}(x - x_0)$. *However, for symmetric covariance matrices we only need to describe $D(D + 1)/2$ elements rather than D^2 .*
- It is invariant under linear transformations, i.e. for $A, B \in M^{M \times D}, c \in \mathbb{R}^M$

$X_1 \sim \mathcal{N}(\mu_1, \Sigma_1), X_2 \sim \mathcal{N}(\mu_2, \Sigma_2)$
and X_1, X_2 independent \Rightarrow

$$AX_1 + BX_2 + c \sim \mathcal{N}(A\mu_1 + B\mu_2 + c, A\Sigma_1 A^T + B\Sigma_2 B^T)$$

Conditional and Marginal Probability Densities



Link [here](#) for a MatLab program
to generate these figures



Transformation of Probability Density

- A probability density transforms differently from functions.
- Let $x = g(y)$.

$$p_y(y) = p_x(g(y)) \left| \frac{dx}{dy} \right| = p_x(g(y)) |g'(y)| = p_x(g(y)) s g'(y), \quad s \in \{-1, 1\}$$

- This is easily derived by taking observations in the interval $(x, x + dx)$ to be transformed to observations in $(y, y + dy)$, i.e.

$$p_y(y)dy = p_x(x)dx$$



Transformation of Probability Density

- For example consider the Gamma distribution

$$\text{Gamma}(x | a, b) = \frac{b^a}{\Gamma(a)} x^{a-1} e^{-xb}$$

- Let us compute the density of $Y = 1/X$.

$$p_y(y) = p_x(g(y)) \left| \frac{dx}{dy} \right| = \frac{b^a}{\Gamma(a)} y^{-(a-1)} e^{-\frac{b}{y}} \left| -\frac{1}{y^2} \right| = \frac{b^a}{\Gamma(a)} y^{-(a+1)} e^{-\frac{b}{y}},$$

where $\frac{dx}{dy} = -\frac{1}{y^2}$

- This is the Inverse Gamma distribution

$$\text{InuGamma}(y | a, b) = \frac{b^a}{\Gamma(a)} y^{-(a+1)} e^{-\frac{b}{y}}$$

Multivariate Change of Variables

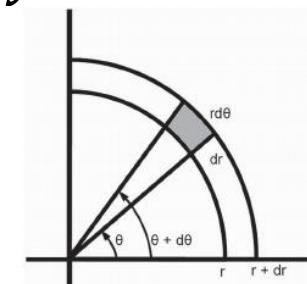
- If f is an invertible mapping, we can define the pdf of the transformed variables using the Jacobian of the inverse mapping $\mathbf{y} \rightarrow \mathbf{x}$:

$$p_y(\mathbf{y}) = p_x(\mathbf{x}) \left| \det \frac{\partial \mathbf{x}}{\partial \mathbf{y}} \right|, \quad \frac{\partial \mathbf{y}}{\partial \mathbf{x}} = \begin{pmatrix} \frac{\partial y_1}{\partial x_1} & \dots & \frac{\partial y_1}{\partial x_n} \\ & \dots & \\ \frac{\partial y_n}{\partial x_1} & \dots & \frac{\partial y_n}{\partial x_n} \end{pmatrix}$$

- As an example, it is trivial to show that transforming a density from Cartesian coordinates $\mathbf{x} = (x, y)$ to polar coordinates $\mathbf{y} = (r, \theta)$, where $x = r \cos \theta$ and $y = r \sin \theta$, $\left| \frac{\partial(x,y)}{\partial(r,\theta)} \right| = r$, gives:

$$p_{r,\theta}(r, \theta) = p_{x,y}(r \cos \theta, r \sin \theta) r$$

$$p_{r,\theta}(r, \theta) dr d\theta = p_{x,y}(r \cos \theta, r \sin \theta) r dr d\theta$$

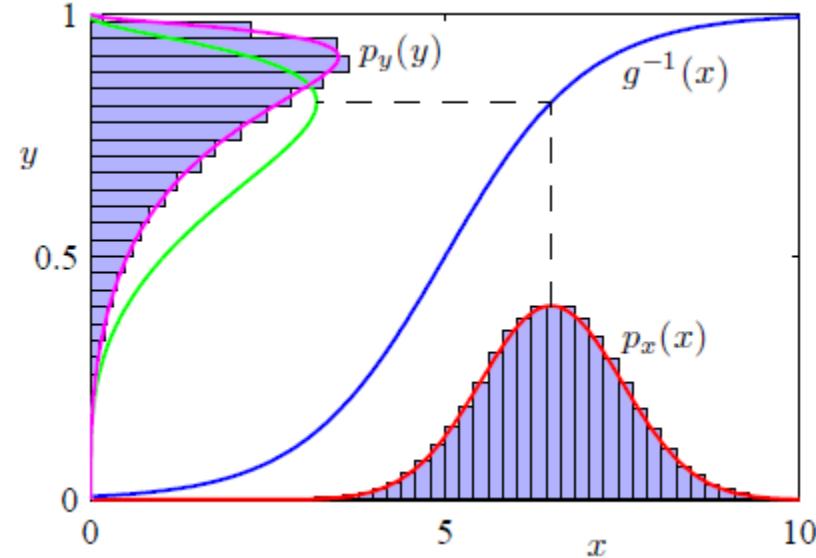


Transformation of Probability Density

- Recall $p_y(y) = p_x(g(y)) \left| \frac{dx}{dy} \right| = p_x(g(y)) s g'(y)$, $s \in \{-1,1\}$
- Using this Eq., note that modes of densities depend on the choice of variables (see 2nd term on the rhs below):

$$p_y(y) = s p_x(x) \{g'(y)\}^2 + s p_x(g(y)) g''(y)$$

- Consider $X \sim \mathcal{N}(6,1)$ and the following



$$x = g(y) = \ln \frac{y}{1-y} + 5, \quad y = g^{-1}(x) = \frac{1}{1+e^{-x+5}}$$

- Transforming $p_x(x)$ as a function gives the same mode for $p_x(g(y))$. The actual mode of $p_y(y)$ is shifted.
- The histogram of $p_y(y)$ is obtained as:

$$y^{(s)} = g^{-1}(x^{(s)}), \text{ where } x^{(s)} \sim p_x(x)$$

Multivariate Student's \mathcal{T} Distribution

$$p(x | \mu, a, b) = \int_0^{\infty} \mathcal{N}(x | \mu, \tau^{-1}) \mathcal{Gamma}(\tau | a, b) d\tau$$

- If we return to the derivation of the univariate Student's \mathcal{T} distribution and substitute $\nu = 2a$, $\lambda = \frac{a}{b}$, $\eta = \tau b / a$, and use

$$\mathcal{Gamma}(\tau | a, b) = \frac{b^a}{\Gamma(a)} \tau^{a-1} e^{-b\tau}$$

we can write the Student's \mathcal{T} distribution as:^{*}

$$\mathcal{T}(x | \mu, \lambda, \nu) = \int_0^{\infty} \mathcal{N}\left(x | \mu, (\eta\lambda)^{-1}\right) \mathcal{Gamma}(\eta | \nu/2, \nu/2) d\eta$$

- This form is useful in *providing generalization to a multivariate Student's \mathcal{T}*

$$\mathcal{T}(x | \mu, \Lambda, \nu) = \int_0^{\infty} \mathcal{N}\left(x | \mu, (\eta\Lambda)^{-1}\right) \mathcal{Gamma}(\eta | \nu/2, \nu/2) d\eta$$

*Use change of variables for distributions, also $d\tau = \lambda d\eta$, and notice that the extra λ -terms that appear cancel out.



Multivariate Student's \mathcal{T} Distribution

$$\mathcal{T}(x | \mu, \Lambda, \nu) = \int_0^\infty \mathcal{N}\left(x | \mu, (\eta\Lambda)^{-1}\right) \text{Gamma}(\eta | \nu/2, \nu/2) d\eta$$

- This integral can be computed analytically as:

$$\mathcal{T}(x | \mu, \Lambda, \nu) = \frac{\Gamma(\frac{D}{2} + \frac{\nu}{2})}{\Gamma(\frac{\nu}{2})} \frac{|\Lambda|^{1/2}}{(\pi\nu)^{D/2}} \left[1 + \frac{\Delta^2}{\nu} \right]^{-\nu/2 - D/2}$$

$$\Delta^2 = (x - \mu)^T \Lambda (x - \mu) \quad (\text{Mahalanobis Distance})$$

- One can derive the above form of the distribution by substitution in the Eq. on the top.

$$\begin{aligned} \mathcal{T}(x | \mu, \Lambda, \nu) &= \frac{(\nu/2)^{\nu/2}}{\Gamma(\nu/2)} \frac{|\Lambda|^{1/2}}{(2\pi)^{D/2}} \int_0^\infty \eta^{D/2} \eta^{\nu/2-1} e^{-\nu\eta/2} e^{-\eta\Delta^2/2} d\eta \quad \text{Use } \tau = \eta(\nu/2 + \Delta^2/2) \\ &= \frac{(\nu/2)^{\nu/2}}{\Gamma(\nu/2)} \frac{|\Lambda|^{1/2}}{(2\pi)^{D/2}} (\nu/2 + \Delta^2/2)^{-D/2-\nu/2} \int_0^\infty \tau^{D/2+\nu/2-1} e^{-\tau} d\tau = \frac{\Gamma(\nu/2 + d/2)}{\Gamma(\nu/2)} \frac{|\Lambda|^{1/2}}{(\pi\nu)^{D/2}} (1 + \Delta^2/\nu)^{-D/2-\nu/2} \end{aligned}$$

- Normalization proof is immediate from the normalization of the normal & Gamma distributions.

Multivariate Student's \mathcal{T} Distribution

$$\mathcal{T}(x | \mu, \Lambda, \nu) = \frac{\Gamma(\frac{D}{2} + \frac{\nu}{2})}{\Gamma(\frac{\nu}{2})} \frac{|\Lambda|^{1/2}}{(\pi\nu)^{D/2}} \left[1 + \frac{\Delta^2}{\nu} \right]^{-\nu/2-D/2}$$

- Some useful results of the multivariate Student's \mathcal{T} are given below:

$$\mathbb{E}[x] = \mu \quad \text{if } \nu > 1, \text{cov}[x] = \frac{\nu}{\nu-2} \Lambda^{-1} \quad \text{if } \nu > 2, \text{mode}[x] = \mu$$

- One can show easily the expression for the mean by using $x = z + \mu$:

$$\mathbb{E}[x] = \int_{-\infty}^{+\infty} \frac{\Gamma(\frac{D}{2} + \frac{\nu}{2})}{\Gamma(\frac{\nu}{2})} \frac{|\Lambda|^{1/2}}{(\pi\nu)^{D/2}} \left[1 + \frac{\Delta^2}{\nu} \right]^{-\nu/2-D/2} (z + \mu) dz$$

- The 1st term drops out since $\mathcal{T}(z | \theta, \Lambda, \nu)$ is even. The 2nd term gives μ from the normalization of the distribution.

- The covariance is computed as:

$$\begin{aligned} \text{cov}[x] &= \int_{\eta=0}^{+\infty} \left[\int_{\mathbf{x}} \mathcal{N}(x | \mu, (\eta\Lambda)^{-1}) (x - \mu)(x - \mu)^T dx \right] \text{Gamma}(\eta | \nu/2, \nu/2) d\eta = \frac{(\nu/2)^{\nu/2}}{\Gamma(\nu/2)} \int_{\eta=0}^{+\infty} (\eta\Lambda)^{-1} \eta^{\nu/2-1} e^{-\nu/2\eta} d\eta \\ &= \Lambda^{-1} \frac{(\nu/2)^{\nu/2}}{\Gamma(\nu/2)} \frac{\Gamma(\nu/2-1)}{(\nu/2)^{\nu/2} (\nu/2)^{\nu/2-1}} = \frac{(\nu/2)\Gamma(\nu/2-1)}{\Gamma(\nu/2)} \Lambda^{-1} = \frac{\nu/2}{\nu/2-1} \Lambda^{-1} = \frac{\nu}{\nu-2} \Lambda^{-1} \end{aligned}$$



Multivariate Student's \mathcal{T} Distribution

$$\mathcal{T}(x | \mu, \Lambda, v) = \frac{\Gamma(\frac{D}{2} + \frac{v}{2})}{\Gamma(\frac{v}{2})} \frac{|\Lambda|^{1/2}}{(\pi v)^{D/2}} \left[1 + \frac{\Delta^2}{v} \right]^{-v/2-D/2}$$

- Differentiation with respect to x also shows the mode being μ :

$$\mathbb{E}[x] = \mu \quad \text{if } v > 1, \text{cov}[x] = \frac{v}{v-2} \Lambda^{-1} \quad \text{if } v > 2, \text{mode}[x] = \mu$$

- The Student's \mathcal{T} has fatter tails than a Gaussian. *The smaller v is the fatter the tails.*
- For $v \rightarrow \infty$, the distribution approaches a Gaussian. Indeed note that:

$$\left[1 + \frac{\Delta^2}{v} \right]^{-v/2-D/2} = \exp\left(-\left(\frac{v}{2} + \frac{D}{2} \right) \ln \left[1 + \frac{\Delta^2}{v} \right] \right)_{v \rightarrow \infty} = \exp\left(-\frac{v}{2} \left(\frac{\Delta^2}{v} - \frac{1}{2} \left(\frac{\Delta^2}{v} \right)^2 \right) \right) = \exp\left(-\frac{\Delta^2}{2} + O(v^{-1}) \right)$$

- The distribution can also be written in terms of $\Sigma = \Lambda^{-1}$ (scale matrix – not the covariance) or $V = v\Sigma$.

Dirichlet Distribution

- We introduce the Dirichlet distribution as a family of “conjugate priors” (*to be formally introduced in a follow up lecture*) for the parameters μ_k in the multinomial distribution.
- The Dirichlet distribution $Dir(\alpha)$, is a family of continuous multivariate probability distributions parametrized by the vector α of positive reals.
- **It is the multivariate generalization of the Beta distribution.**

Dirichlet Distribution

- Its probability density function returns the belief that the probabilities of K rival events are μ_k given that each event has been observed $\alpha_k - 1$ times:

$$p(\boldsymbol{\mu} | \boldsymbol{\alpha}) \propto \prod_{k=1}^K \mu_k^{\alpha_k - 1},$$

$$0 \leq \mu_k \leq 1,$$

$$\sum_{k=1}^K \mu_k = 1$$

- The distribution over the space of μ_k is $K - 1$ dimensional due to the last constraint above.

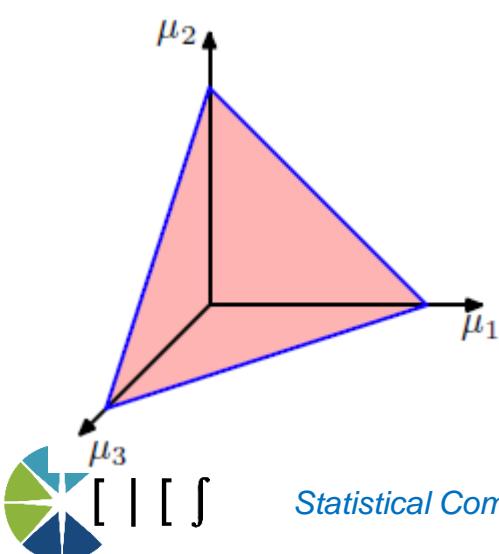


Dirichlet Distribution

- The Dirichlet distribution of order $K \geq 2$ with parameters $\alpha_1, \dots, \alpha_K > 0$ has a PDF with respect to Lebesgue measure on \mathbb{R}^{K-1} given by

$$p(\boldsymbol{\mu} | \boldsymbol{\alpha}) = \frac{1}{Beta(\boldsymbol{\alpha})} \prod_{k=1}^K \mu_k^{\alpha_k - 1}$$

for all $\mu_1, \dots, \mu_{K-1} > 0$ satisfying $\mu_1 + \dots + \mu_{K-1} < 1$, where μ_K is an abbreviation for $1 - \mu_1 - \dots - \mu_{K-1}$. *The normalizing constant is the multinomial Beta function:*



$$Beta(\boldsymbol{\alpha}) = \frac{\prod_{k=1}^K \Gamma(\alpha_k)}{\Gamma\left(\sum_{k=1}^K \alpha_k\right)}, \boldsymbol{\alpha} = (\alpha_1, \alpha_2, \dots, \alpha_K)^T$$

The Dirichlet distribution over (μ_1, μ_2, μ_3) is confined on a plane as shown.

Dirichlet Distribution

➤ We write the Dirichlet distribution as:

$$p(\boldsymbol{\mu} | \boldsymbol{\alpha}) = K(\boldsymbol{\alpha}) \prod_{k=1}^K \mu_k^{\alpha_k - 1}, \quad K(\boldsymbol{\alpha}) = \frac{\Gamma(a_0)}{\Gamma(a_1) \dots \Gamma(a_K)}, \quad a_0 = a_1 + \dots + a_K$$

➤ Note the following useful relation:

$$\frac{\partial}{\partial \alpha_j} \prod_{k=1}^K \mu_k^{\alpha_k - 1} = \frac{\partial}{\partial \alpha_j} \prod_{k=1}^K e^{(\alpha_k - 1) \ln \mu_k} = \ln \mu_j \prod_{k=1}^K \mu_k^{\alpha_k - 1}$$

➤ From this we can derive an interesting expression for $\mathbb{E}[\ln \mu_j]$

$$\begin{aligned} \mathbb{E}[\ln \mu_j] &= K(\boldsymbol{\alpha}) \int_0^1 \dots \int_0^1 \ln \mu_j \prod_{k=1}^K \mu_k^{\alpha_k - 1} d\mu_1 \dots d\mu_K = K(\boldsymbol{\alpha}) \int_0^1 \dots \int_0^1 \frac{\partial}{\partial \alpha_j} \prod_{k=1}^K \mu_k^{\alpha_k - 1} d\mu_1 \dots d\mu_K = \\ K(\boldsymbol{\alpha}) \frac{\partial}{\partial \alpha_j} \int_0^1 \dots \int_0^1 \prod_{k=1}^K \mu_k^{\alpha_k - 1} d\mu_1 \dots d\mu_K &= K(\boldsymbol{\alpha}) \frac{\partial}{\partial \alpha_j} \frac{1}{K(\boldsymbol{\alpha})} = -\frac{\partial \ln K(\boldsymbol{\alpha})}{\partial \alpha_j} = \underbrace{\frac{\partial \ln \Gamma(\alpha_j)}{\partial \alpha_j}}_{\Psi(\alpha_j)} - \underbrace{\frac{\partial \ln \Gamma(\alpha_0)}{\partial \alpha_0}}_{\Psi(\alpha_0)} \end{aligned}$$

where $\Psi(\alpha) = d \ln \Gamma(\alpha) / d\alpha$ is the **digamma function**.

$$\mathbb{E}[\ln \mu_j] = \Psi(\alpha_j) - \Psi(\alpha_0), \quad \Psi(\alpha_j) = \frac{\partial \ln \Gamma(\alpha_j)}{\partial \alpha_j}, \quad \Psi(\alpha_0) = \frac{\partial \ln \Gamma(\alpha_0)}{\partial \alpha_0}$$



Dirichlet Distribution: Normalization

- To show the normalization, we use induction. The case for $M = 2$ was shown earlier for the Beta distribution.
- Assume that the Dirichlet normalization formula is valid for $M - 1$ terms. We will show the formula for M terms:

$$p_M(\mu_1, \dots, \mu_{M-1}) = C_M \prod_{k=1}^{M-1} \mu_k^{\alpha_k-1} \underbrace{\left(1 - \sum_{j=1}^{M-1} \mu_j\right)}_{\mu_M}^{\alpha_M-1}$$

- Let us integrate out μ_{M-1} :

$$\begin{aligned} p_{M-1}(\mu_1, \dots, \mu_{M-2}) &= C_M \int_0^{1 - \sum_{j=1}^{M-2} \mu_j} \left(\prod_{k=1}^{M-2} \mu_k^{\alpha_k-1} \right) \mu_{M-1}^{\alpha_{M-1}-1} \left(1 - \sum_{j=1}^{M-2} \mu_j - \mu_{M-1} \right)^{\alpha_M-1} d\mu_{M-1} = \\ &= C_M \left(\prod_{k=1}^{M-2} \mu_k^{\alpha_k-1} \right) \int_0^1 t^{\alpha_{M-1}-1} \left(1 - \sum_{j=1}^{M-2} \mu_j \right)^{\alpha_{M-1}-1 + \alpha_M-1+1} (1-t)^{\alpha_M-1} dt \end{aligned}$$
$$\mu_{M-1} = t \left(1 - \sum_{j=1}^{M-2} \mu_j \right)$$



Dirichlet Distribution: Normalization

$$\begin{aligned} p_{M-1}(\mu_1, \dots, \mu_{M-2}) &= C_M \left(\prod_{k=1}^{M-2} \mu_k^{\alpha_k - 1} \right) \left(1 - \sum_{j=1}^{M-2} \mu_j \right)^{\alpha_{M-1} + \alpha_M - 1} \int_0^1 t^{\alpha_{M-1}-1} (1-t)^{\alpha_M-1} dt = \\ &= C_M \underbrace{\left(\prod_{k=1}^{M-2} \mu_k^{\alpha_k - 1} \right) \left(1 - \sum_{j=1}^{M-2} \mu_j \right)^{\alpha_{M-1} + \alpha_M - 1}}_{\text{Dirichlet } (M-1)} \frac{\Gamma(\alpha_{M-1})\Gamma(\alpha_M)}{\Gamma(\alpha_{M-1} + \alpha_M)} \end{aligned}$$

- The last step above comes from the normalization of Beta.
- What we have above is an $(M - 1)$ term Dirichlet distribution with coefficients $\alpha_1, \dots, \alpha_{M-2}, \alpha_{M-1} + \alpha_M$. Since we assumed that the normalization formula is valid for $(M - 1)$ terms, we must have:

$$1 = C_M \frac{\Gamma(\alpha_1) \dots \Gamma(\alpha_{M-2}) \Gamma(\alpha_{M-1} + \alpha_M)}{\Gamma(\alpha_1 + \dots + \alpha_M)} \frac{\Gamma(\alpha_{M-1}) \Gamma(\alpha_M)}{\Gamma(\alpha_{M-1} + \alpha_M)} \Rightarrow$$

$$C_M = \frac{\Gamma(\alpha_1 + \dots + \alpha_M)}{\Gamma(\alpha_1) \dots \Gamma(\alpha_{M-2}) \Gamma(\alpha_{M-1}) \Gamma(\alpha_M)}$$



Dirichlet Distribution

- Using the multinomial as a “likelihood” and the Dirichlet as “the conjugate prior”, we arrive at the following “posterior”

$$p(\boldsymbol{\mu} | \mathcal{D}) \propto \underbrace{p(\mathcal{D} | \boldsymbol{\mu})}_{\text{Multinomial}} p(\boldsymbol{\mu}) \Rightarrow p(\boldsymbol{\mu} | \mathcal{D}) \propto \prod_{k=1}^K \mu_k^{\alpha_k + m_k - 1}$$

which is a Dirichlet distribution $\text{Dir}(\boldsymbol{\mu} | \alpha_1 + m_1, \dots, \alpha_K + m_K)$.

- The normalization factor is computed easily from the normalization factor of the Dirichlet as:

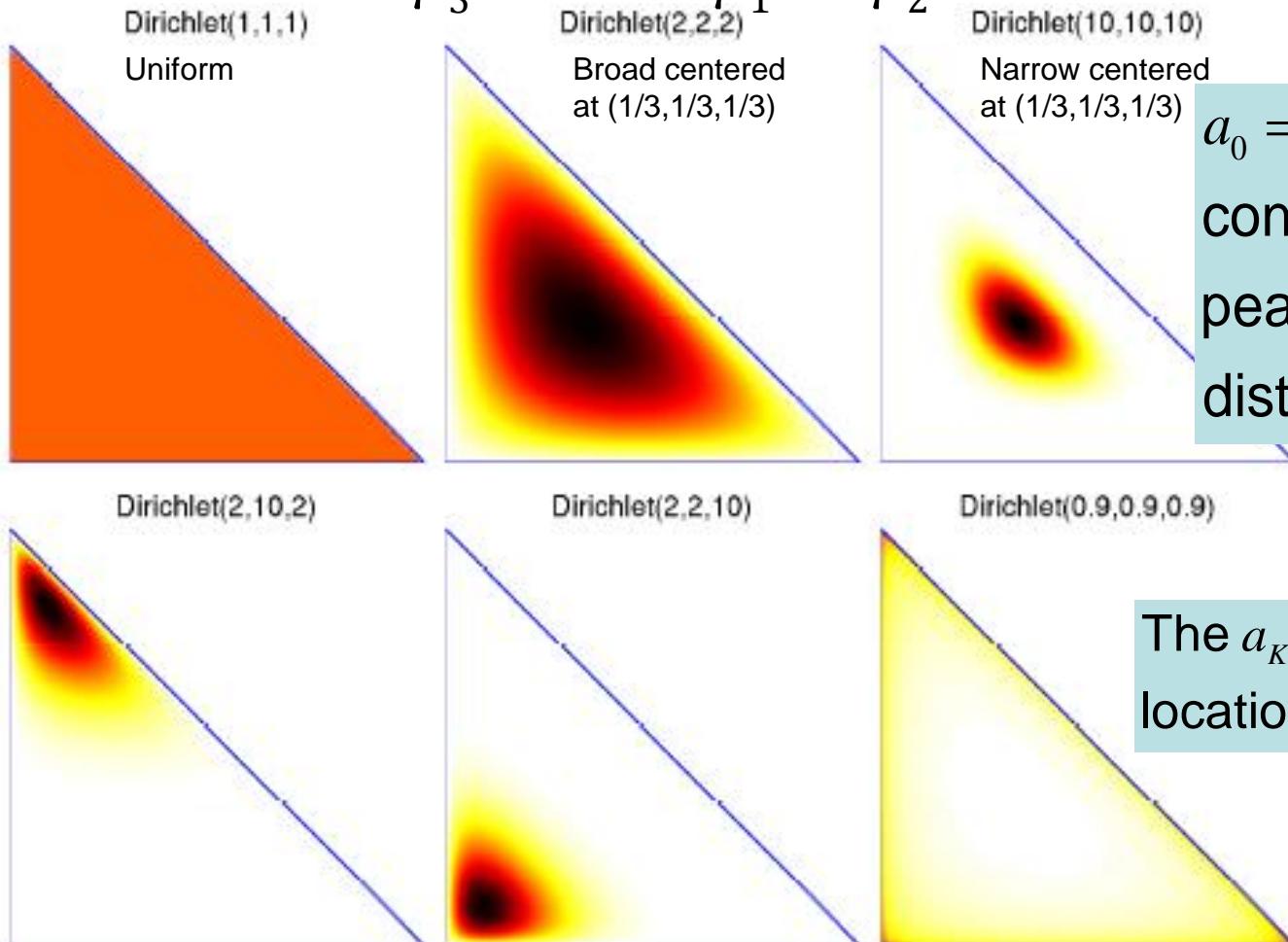
$$p(\boldsymbol{\mu} | \mathcal{D}) = \frac{\Gamma\left(\sum_{k=1}^K \alpha_k + N\right)}{\prod_{k=1}^K \Gamma(\alpha_k + m_k)} \prod_{k=1}^K \mu_k^{\alpha_k + m_k - 1}$$

- α_k can be interpreted as “the effective number of prior observations of $x_k = 1$ ”.



Dirichlet Distribution

Examples of Dirichlet distribution over (μ_1, μ_2, μ_3) which can be plotted in 2D since $\mu_3 = 1 - \mu_1 - \mu_2$.

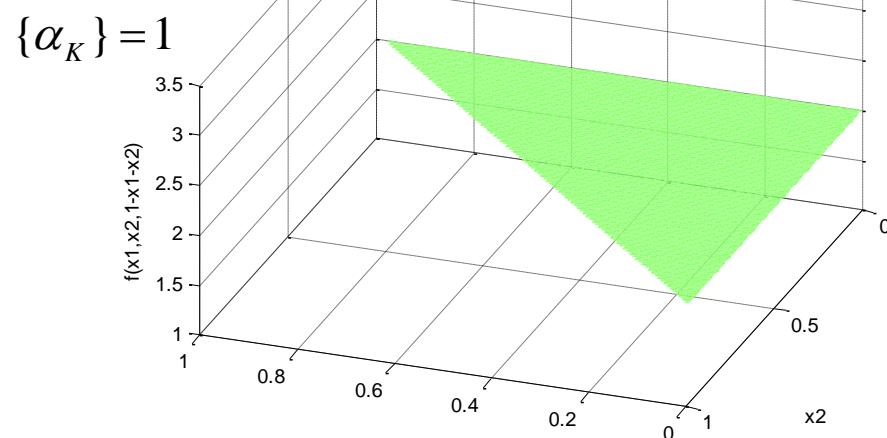
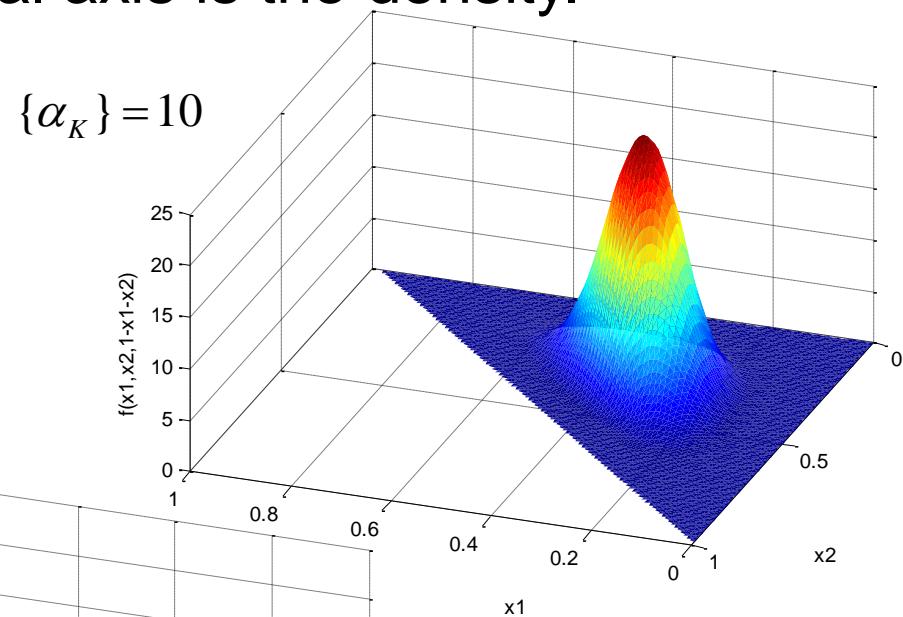
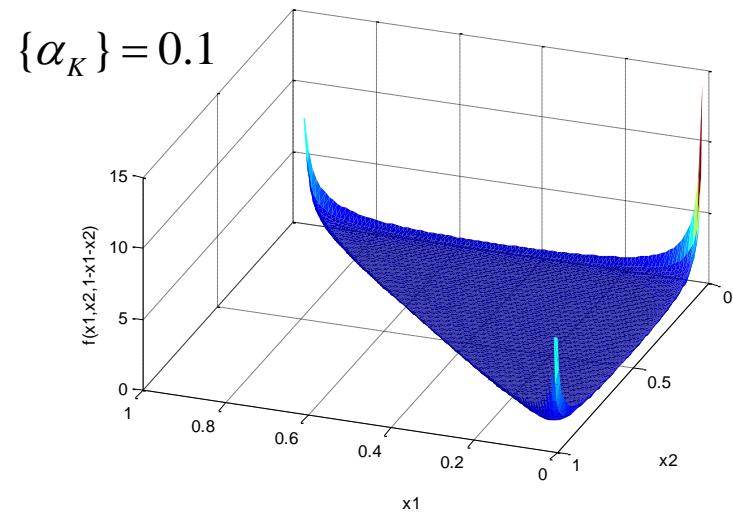


$a_0 = a_1 + \dots + a_K$
controls how
peaked the
distribution is

The a_K control the
location of the peak

Dirichlet Distribution

The Dirichlet distribution over (μ_1, μ_2, μ_3) where the horizontal axes are μ_1 and μ_2 and the vertical axis is the density.

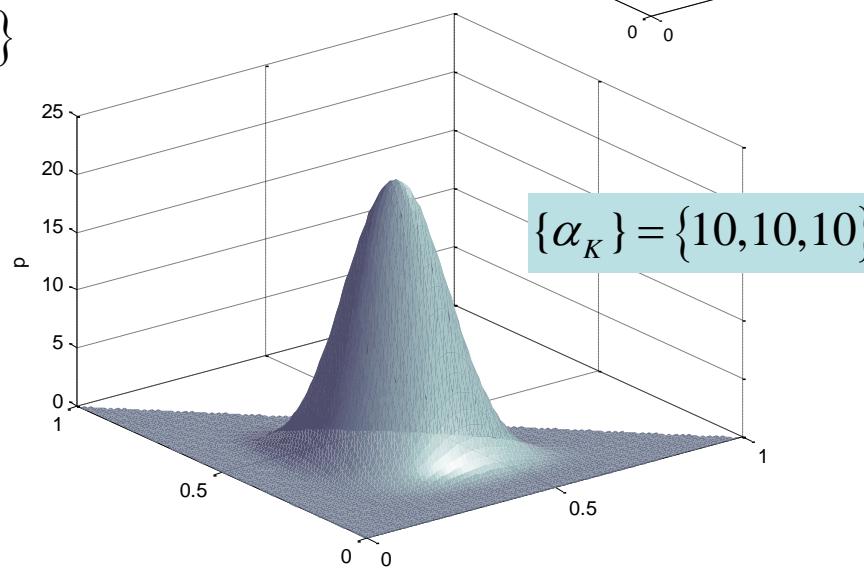
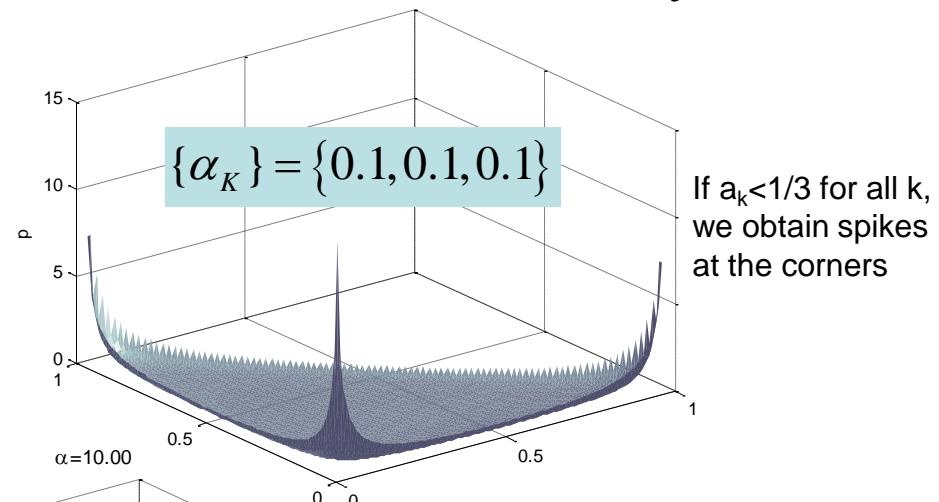
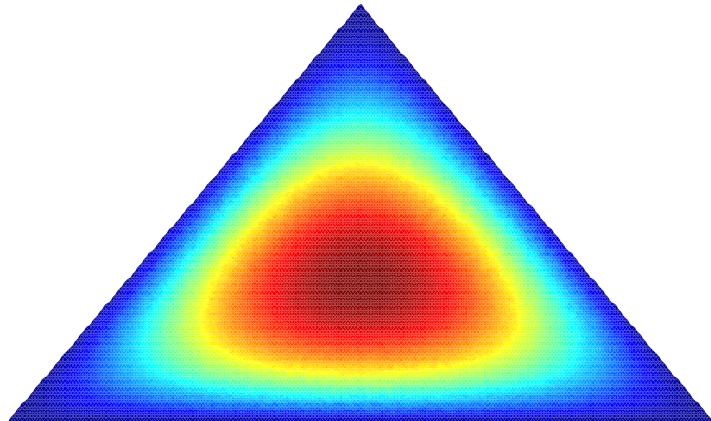


MatLab Code



Dirichlet Distribution

The Dirichlet distribution over (μ_1, μ_2, μ_3) where the horizontal axes are μ_1 and μ_2 and the vertical axis is the density.



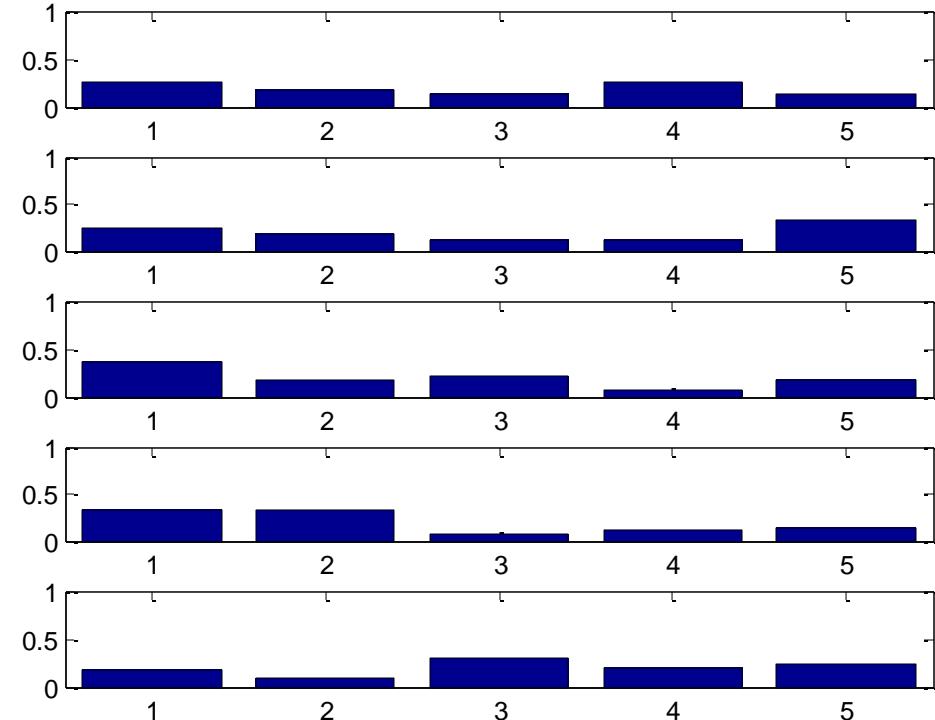
Run [visDirichletGui](#) & [dirichlet3dPlot](#) from [PMTK](#)



Dirichlet Distribution

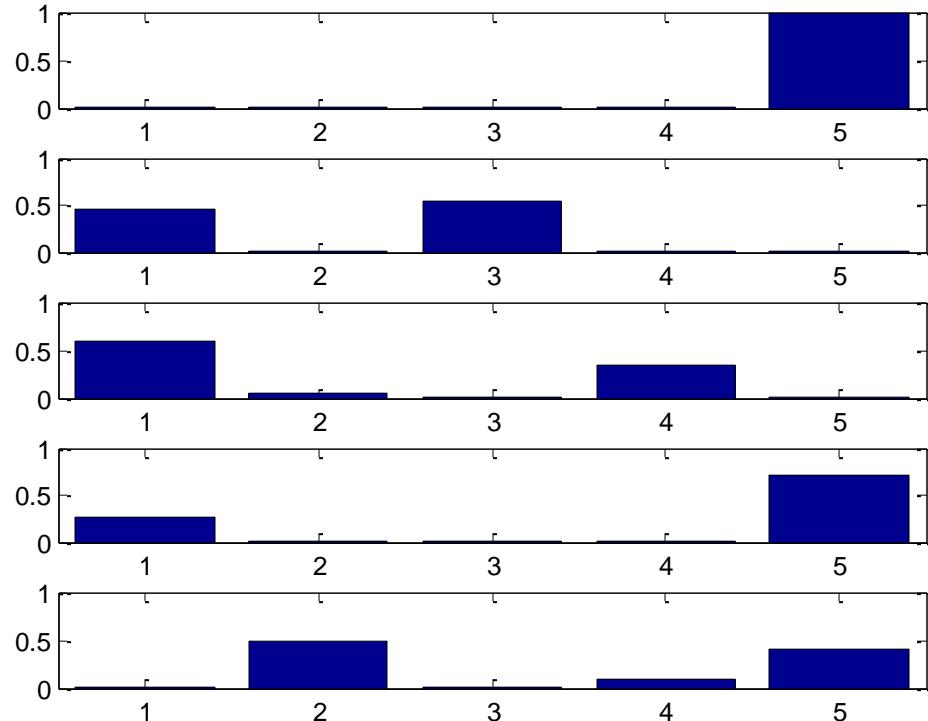
Samples from 5-dimensional symmetric Dirichlet distribution.

Samples from Dir ($\alpha=5$)



$$\{\alpha_K\} = \{5, 5, \dots, 5\}$$

Samples from Dir ($\alpha=0.1$)



$$\{\alpha_K\} = \{0.1, 0.1, \dots, 0.1\}$$

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



Dirichlet Distribution

- In closing, we have the following properties (you only need the normalization $\int \prod_{k=1}^K \mu_k^{\alpha_k-1} d\mu = \frac{\Gamma(a_1) \dots \Gamma(a_K)}{\Gamma(a_0)}$, $a_0 = a_1 + \dots + a_K$ of the Dirichlet distribution and the property $\Gamma(x+1) = x\Gamma(x)$ to prove them):

$$\mathbb{E}[\mu_k] = \frac{\alpha_k}{\alpha_0}, \text{ mode}[\mu_k] = \frac{\alpha_k - 1}{\alpha_0 - 1}, \text{ var}[\mu_k] = \frac{\alpha_k(\alpha_0 - \alpha_k)}{\alpha_0^2(\alpha_0 + 1)}, \text{ cov}[\mu_j \mu_l] = -\frac{\alpha_j \alpha_l}{\alpha_0^2(\alpha_0 + 1)} \quad (j \neq l)$$

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

- Often we use:

$$\alpha_k = \alpha / K$$

In this case:

$$\mathbb{E}[\mu_k] = \frac{1}{K}, \text{ var}[\mu_k] = \frac{K-1}{K^2(\alpha+1)}$$

Increasing α increases the precision of the distribution.

