



Probability and Distributions

Mathematics for Machine Learning

San Diego Machine Learning
Liam Barstad

Distributions

Probability distribution – equation that describes likelihood of random variable(s)

Sample space (Ω) – set of all possible outcomes –
 $\{HH, HT, TH, TT\}$

Event space (A) – subset of sample space w/ outcome –
(e.g. at least one head) – $\{HH, HT, TH\}$

Target space (T) – set of values that can result -
(e.g. number of heads) – $\{0, 1, 2\}$

Random Variable (X) – mapping where $X : \Omega \rightarrow T$

Discrete Distribution Functions

Probability mass function (pmf) – describes discrete random variable – $p(x, y)$

	y_1					
Y	y_2			n_{ij}		
	y_3					
		x_1	x_2	x_3	x_4	x_5
		X				

$$P(X = x_i, Y = y_j) = \frac{n_{ij}}{N},$$

$$P(X = x_i) = \frac{c_i}{N} = \frac{\sum_{j=1}^3 n_{ij}}{N}$$

Marginal probability $p(x)$, $p(y)$

$$P(Y = y_j) = \frac{r_j}{N} = \frac{\sum_{i=1}^5 n_{ij}}{N},$$

Conditional probability

$$p(y \mid x) - P(Y = y_j \mid X = x_i) = \frac{n_{ij}}{c_i},$$

$$p(x \mid y) - P(X = x_i \mid Y = y_j) = \frac{n_{ij}}{r_j}.$$

Continuous Distribution Functions

Probability density function (pdf) – describes continuous random variable

Definition 6.1 (Probability Density Function). A function $f : \mathbb{R}^D \rightarrow \mathbb{R}$ is called a *probability density function (pdf)* if

1. $\forall \mathbf{x} \in \mathbb{R}^D : f(\mathbf{x}) \geq 0$
2. Its integral exists and

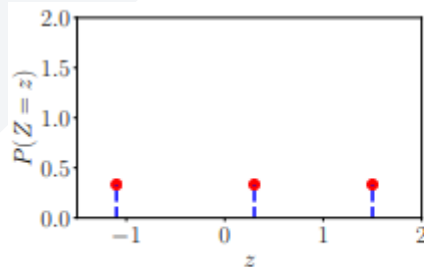
$$\int_{\mathbb{R}^D} f(\mathbf{x}) d\mathbf{x} = 1. \quad (6.15)$$

$$P(a \leq X \leq b) = \int_a^b f(x) dx, \quad (6.16)$$

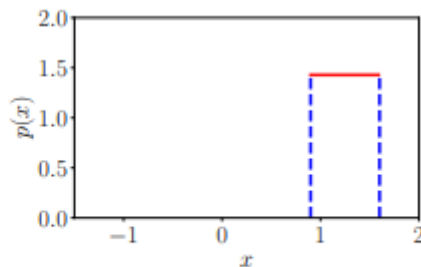
Cumulative distribution function (cdf) – describes likelihood random variable(s) are less

t

$$F_X(\mathbf{x}) = P(X_1 \leq x_1, \dots, X_D \leq x_D), \quad F_X(\mathbf{x}) = \int_{-\infty}^{x_1} \dots \int_{-\infty}^{x_D} f(z_1, \dots, z_D) dz_1 \dots dz_D.$$



(a) Discrete distribution



(b) Continuous distribution

Uniform distribution – all results are

Probability Rules

Sum rule, aka marginalization property

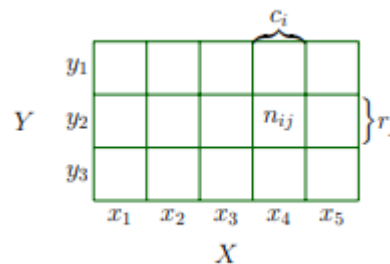
$$p(\mathbf{x}) = \begin{cases} \sum_{\mathbf{y} \in \mathcal{Y}} p(\mathbf{x}, \mathbf{y}) & \text{if } \mathbf{y} \text{ is discrete} \\ \int_{\mathbf{y}} p(\mathbf{x}, \mathbf{y}) d\mathbf{y} & \text{if } \mathbf{y} \text{ is continuous} \end{cases} \quad p(x_i) = \int p(x_1, \dots, x_D) d\mathbf{x}_{\setminus i}$$

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

Product rule

Bayes rule

$$\underbrace{p(\mathbf{x} | \mathbf{y})}_{\text{posterior}} = \frac{\overbrace{p(\mathbf{y} | \mathbf{x})}^{\text{likelihood}} \overbrace{p(\mathbf{x})}^{\text{prior}}}{\underbrace{p(\mathbf{y})}_{\text{evidence}}}$$



Expected Value

Each value multiplied by its probability

$$\mathbb{E}_{x_d}[x_d] := \begin{cases} \int_{\mathcal{X}} x_d p(x_d) dx_d \\ \sum_{x_i \in \mathcal{X}} x_i p(x_d = x_i) \end{cases} \quad \begin{aligned} \mathbb{E}_X[g(x)] &= \int_{\mathcal{X}} g(x) p(x) dx \\ \mathbb{E}_X[g(x)] &= \sum_{x \in \mathcal{X}} g(x) p(x) \end{aligned}$$

Mean - $\mathbb{E}[x]$

$$\mathbb{E}_X[\mathbf{x}] = \begin{bmatrix} \mathbb{E}_{X_1}[x_1] \\ \vdots \\ \mathbb{E}_{X_D}[x_D] \end{bmatrix} \in \mathbb{R}^D \quad \mathbb{E}_X[g(\mathbf{x})] = \begin{bmatrix} \mathbb{E}_{X_1}[g(x_1)] \\ \vdots \\ \mathbb{E}_{X_D}[g(x_D)] \end{bmatrix} \in \mathbb{R}^D$$

$$\begin{aligned} \mathbb{E}_X[f(\mathbf{x})] &= \int f(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} \\ &= \int [ag(\mathbf{x}) + bh(\mathbf{x})] p(\mathbf{x}) d\mathbf{x} \\ &= a \int g(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} + b \int h(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} \\ &= a\mathbb{E}_X[g(\mathbf{x})] + b\mathbb{E}_X[h(\mathbf{x})]. \end{aligned}$$

Variance + Covariance

Covariance – deviation from central value of 2 random variables

$$\text{Cov}_{X,Y}[x, y] := \mathbb{E}_{X,Y}[(x - \mathbb{E}_X[x])(y - \mathbb{E}_Y[y])]$$

Correlation – normalized covariance

$$\text{Cov}[x, y] = \mathbb{E}[xy] - \mathbb{E}[x]\mathbb{E}[y]$$

$$\text{Cov}[\mathbf{x}, \mathbf{y}] = \mathbb{E}[\mathbf{x}\mathbf{y}^\top] - \mathbb{E}[\mathbf{x}]\mathbb{E}[\mathbf{y}]^\top = \text{Cov}[\mathbf{y}, \mathbf{x}]^\top \in \mathbb{R}^{D \times E}.$$

$$\text{corr}[x, y] = \frac{\text{Cov}[x, y]}{\sqrt{\mathbb{V}[x]\mathbb{V}[y]}} \in [-1, 1]$$

Variance – squared deviation from $\mathbb{E}[x]$, covariance of random variable w/ itself

$$\mathbb{V}_X[\mathbf{x}] = \text{Cov}_X[\mathbf{x}, \mathbf{x}]$$

$$= \mathbb{E}_X[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^\top] = \mathbb{E}_X[\mathbf{x}\mathbf{x}^\top] - \mathbb{E}_X[\mathbf{x}]\mathbb{E}_X[\mathbf{x}]^\top$$

$$= \begin{bmatrix} \text{Cov}[x_1, x_1] & \text{Cov}[x_1, x_2] & \dots & \text{Cov}[x_1, x_D] \\ \text{Cov}[x_2, x_1] & \text{Cov}[x_2, x_2] & \dots & \text{Cov}[x_2, x_D] \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}[x_D, x_1] & \dots & \dots & \text{Cov}[x_D, x_D] \end{bmatrix}.$$

$$\mathbb{V}_X[x] := \mathbb{E}_X[(x - \mu)^2]$$

$$\mathbb{V}_X[x] = \mathbb{E}_X[x^2] - (\mathbb{E}_X[x])^2$$

$$\frac{1}{N^2} \sum_{i,j=1}^N (x_i - x_j)^2 = 2 \left[\frac{1}{N} \sum_{i=1}^N x_i^2 - \left(\frac{1}{N} \sum_{i=1}^N x_i \right)^2 \right]$$

Variance + Covariance Useful Properties

Given affine transformation $y = Ax + b$

$$\mathbb{E}_Y[\mathbf{y}] = \mathbb{E}_X[\mathbf{Ax} + \mathbf{b}] = \mathbf{A}\mathbb{E}_X[\mathbf{x}] + \mathbf{b} = \mathbf{A}\boldsymbol{\mu} + \mathbf{b}$$

$$\mathbb{V}_Y[\mathbf{y}] = \mathbb{V}_X[\mathbf{Ax} + \mathbf{b}] = \mathbb{V}_X[\mathbf{Ax}] = \mathbf{A}\mathbb{V}_X[\mathbf{x}]\mathbf{A}^\top = \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top$$

$$\begin{aligned}\text{Cov}[\mathbf{x}, \mathbf{y}] &= \mathbb{E}[\mathbf{x}(\mathbf{Ax} + \mathbf{b})^\top] - \mathbb{E}[\mathbf{x}]\mathbb{E}[\mathbf{Ax} + \mathbf{b}]^\top \\ &= \mathbb{E}[\mathbf{x}]\mathbf{b}^\top + \mathbb{E}[\mathbf{x}\mathbf{x}^\top]\mathbf{A}^\top - \boldsymbol{\mu}\mathbf{b}^\top - \boldsymbol{\mu}\boldsymbol{\mu}^\top\mathbf{A}^\top \\ &= \boldsymbol{\mu}\mathbf{b}^\top - \boldsymbol{\mu}\mathbf{b}^\top + (\mathbb{E}[\mathbf{x}\mathbf{x}^\top] - \boldsymbol{\mu}\boldsymbol{\mu}^\top)\mathbf{A}^\top \\ &\stackrel{(6.38b)}{=} \boldsymbol{\Sigma}\mathbf{A}^\top,\end{aligned}$$

where $\boldsymbol{\Sigma} = \mathbb{E}[\mathbf{x}\mathbf{x}^\top] - \boldsymbol{\mu}\boldsymbol{\mu}^\top$ is the covariance of X .

Geometric Interpretation

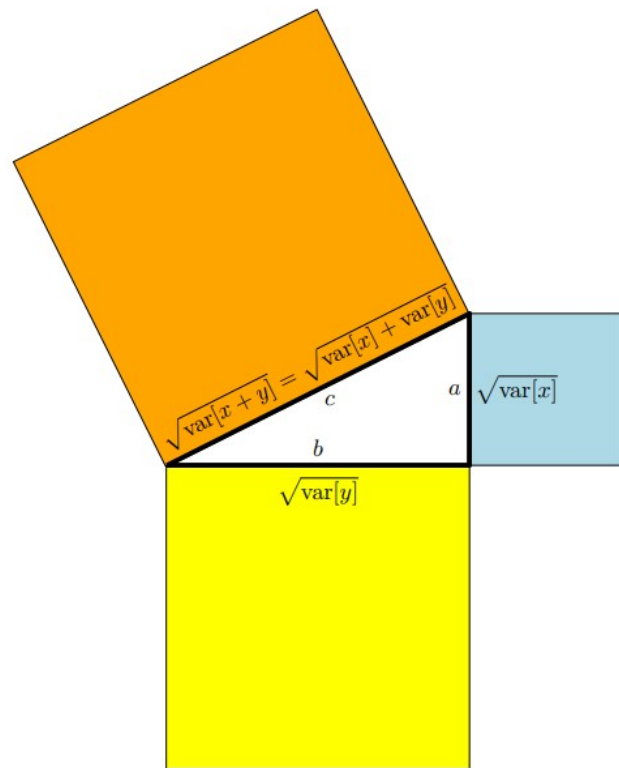
Covariance can be thought of as inner product

$$\langle X, Y \rangle := \text{Cov}[x, y]$$

If inner product = 0, a and b are orthogonal, $V[x + y] = V[x] + V[y]$

$$\cos \theta = \frac{\langle X, Y \rangle}{\|X\| \|Y\|} = \frac{\text{Cov}[x, y]}{\sqrt{V[x]V[y]}}$$

Information Geometry



Gaussian (Normal) Distribution

Univariate density

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

Multivariate density

$$p(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-\frac{D}{2}} |\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

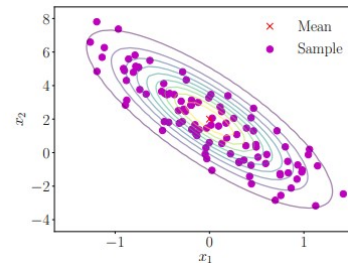
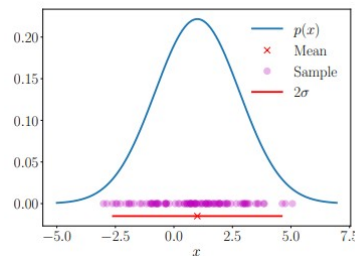
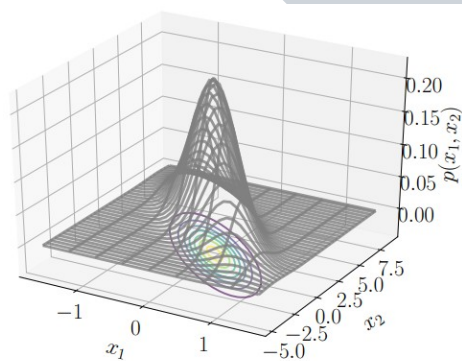
Both marginals and conditionals of gaussians are gaussian

$$p(\mathbf{x}, \mathbf{y}) = \mathcal{N}\left(\begin{bmatrix} \boldsymbol{\mu}_x \\ \boldsymbol{\mu}_y \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{xx} & \boldsymbol{\Sigma}_{xy} \\ \boldsymbol{\Sigma}_{yx} & \boldsymbol{\Sigma}_{yy} \end{bmatrix}\right)$$

$$p(\mathbf{x} | \mathbf{y}) = \mathcal{N}(\boldsymbol{\mu}_{x|y}, \boldsymbol{\Sigma}_{x|y})$$

$$\boldsymbol{\mu}_{x|y} = \boldsymbol{\mu}_x + \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_{yy}^{-1}(\mathbf{y} - \boldsymbol{\mu}_y)$$

$$\boldsymbol{\Sigma}_{x|y} = \boldsymbol{\Sigma}_{xx} - \boldsymbol{\Sigma}_{xy} \boldsymbol{\Sigma}_{yy}^{-1} \boldsymbol{\Sigma}_{yx}$$



Properties of Gaussian Distributions

Products of Gaussians – $N(x | a, A)$, $N(x | b, B)$ is gaussian
scaled by C , $cN(x | c, C)$

$$c = \mathcal{N}(a | b, A + B) = \mathcal{N}(b | a, A + B)$$

Sums of Gaussians

$$p(\mathbf{x} + \mathbf{y}) = \mathcal{N}(\boldsymbol{\mu}_x + \boldsymbol{\mu}_y, \boldsymbol{\Sigma}_x + \boldsymbol{\Sigma}_y)$$

$$p(a\mathbf{x} + b\mathbf{y}) = \mathcal{N}(a\boldsymbol{\mu}_x + b\boldsymbol{\mu}_y, a^2\boldsymbol{\Sigma}_x + b^2\boldsymbol{\Sigma}_y)$$

Common Distributions

Bernoulli Distribution – single binary random variable, represents probability of $X = 1$

$$p(x | \mu) = \mu^x (1 - \mu)^{1-x}$$

Binomial Distribution – probability of observing m occurrences of $X = 1$ in a set of N samples

$$p(m | N, \mu) = \binom{N}{m} \mu^m (1 - \mu)^{N-m}$$

Beta Distribution – models uncertainty over a continuous random variable, α is number of successes, β is number of failures (i.e. there's a 95% chance the success rate is between 3% and 7%)

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

$$\Gamma(t) := \int_0^\infty x^{t-1} \exp(-x) dx, \quad t > 0$$
$$\Gamma(t + 1) = t\Gamma(t).$$

Conjugacy

Prior is **conjugate** for the likelihood function if the posterior is the same form as the prior (retain same distance structure geometrically)

Example 1: Beta-Binomial, x is number of heads, μ is prob. of heads (prior)

$$p(x | N, \mu) = \binom{N}{x} \mu^x (1 - \mu)^{N-x}$$

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

$$\begin{aligned} p(\mu | x = h, N, \alpha, \beta) &\propto p(x | N, \mu) p(\mu | \alpha, \beta) \\ &\propto \mu^h (1 - \mu)^{(N-h)} \mu^{\alpha-1} (1 - \mu)^{\beta-1} \\ &= \mu^{h+\alpha-1} (1 - \mu)^{(N-h)+\beta-1} \end{aligned}$$

Example 2: Beta-Bernoulli

$$p(\theta | x, \alpha, \beta) = p(x | \theta) p(\theta | \alpha, \beta)$$

$$\underbrace{p(\mathbf{x} | \mathbf{y})}_{\text{posterior}} = \frac{\overbrace{p(\mathbf{y} | \mathbf{x})}^{\text{likelihood}} \overbrace{p(\mathbf{x})}^{\text{prior}}}{\underbrace{p(\mathbf{y})}_{\text{evidence}}}$$

Exponential Family

Only family where the number of sufficient statistics used to describe data has finite dimensions

$$p(\mathbf{x} | \boldsymbol{\theta}) = h(\mathbf{x}) \exp (\langle \boldsymbol{\theta}, \boldsymbol{\phi}(\mathbf{x}) \rangle - A(\boldsymbol{\theta}))$$

- $h(\mathbf{x})$ can be absorbed into sufficient statistics by adding $\log(h(\mathbf{x}))$ to $\boldsymbol{\phi}(\mathbf{x})$
- $A(\boldsymbol{\theta})$ is log-partition function which makes sure the sum is 1

$$p(\mathbf{x} | \boldsymbol{\theta}) \propto \exp (\boldsymbol{\theta}^\top \boldsymbol{\phi}(\mathbf{x}))$$

For a Gaussian distribution:

$$p(x | \boldsymbol{\theta}) \propto \exp(\theta_1 x + \theta_2 x^2)$$

$$\boldsymbol{\theta} = \left[\frac{\mu}{\sigma^2}, -\frac{1}{2\sigma^2} \right]^\top$$

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

Exponential Family - Conjugates

Every exponential distribution has a conjugate prior

$$p(\boldsymbol{\theta} | \boldsymbol{\gamma}) = h_c(\boldsymbol{\theta}) \exp \left(\left\langle \begin{bmatrix} \gamma_1 \\ \gamma_2 \end{bmatrix}, \begin{bmatrix} \boldsymbol{\theta} \\ -A(\boldsymbol{\theta}) \end{bmatrix} \right\rangle - A_c(\boldsymbol{\gamma}) \right)$$

For Bernoulli/Beta distribution

$$p(x | \mu) = \exp \left[x \log \frac{\mu}{1 - \mu} + \log(1 - \mu) \right] \quad \boldsymbol{\gamma} := [\alpha, \beta + \alpha]^\top \text{ and } h_c(\mu) := \mu / (1 - \mu)$$

$$p(\mu | \alpha, \beta) = \frac{\mu}{1 - \mu} \exp \left[\alpha \log \frac{\mu}{1 - \mu} + (\beta + \alpha) \log(1 - \mu) - A_c(\boldsymbol{\gamma}) \right]$$

$$p(\mu | \alpha, \beta) \propto \mu^{\alpha-1} (1 - \mu)^{\beta-1}$$