

# Lecture 2: Probability and Statistics (continued)

©Christopher S. Bretherton

Winter 2015

## 2.1 Expectation and moments

**Expectation** of a function  $g(X)$  of a RV  $X$  is

$$\begin{aligned} E[g(X)] &= \sum_{x:p(x)>0} g(x)p(x)dx && \text{discrete RV } X \\ E[g(X)] &= \int_{-\infty}^{\infty} g(x)f(x)dx && \text{continuous RV } X \end{aligned}$$

The expectation of  $X$  is also called its **mean**  $\mu_X$ , sometimes denoted  $\bar{X}$ .

**Variance**  $\text{var}[X] = E[(X - \mu_X)^2] = E[X^2] - (E[X])^2$ , whose square root is the **standard deviation**  $\sigma_X$ , a measure of the spread of  $X$  about its mean.

**$n$ 'th moment**  $E[X^n]$ . The third moment is a measure of **skewness** or asymmetry of the PDF of  $X$  about its mean.

## 2.2 Examples of random variables

**Bernoulli**  $P(X = 1) = p$ ;  $P(X = 0) = q = 1 - p$ .  $\mu_X = p$  and  $\sigma_X = (pq)^{1/2}$ .  
The sum of  $N \geq 1$  independent identically-distributed Bernoulli random variables is a **binomial** distribution with parameters  $N$  and  $p$ .

**Uniform** distribution on  $(\alpha, \beta)$ :

$$f(x) = \frac{1}{\beta - \alpha}, \quad \alpha < x < \beta.$$

$\mu_X = (\alpha + \beta)/2$  and  $\sigma_X = (\beta - \alpha)/\sqrt{12}$ ; note how they scale with  $\alpha, \beta$ .

**Gaussian** (or **normal**) distribution  $n(\mu, \sigma)$ , with mean  $\mu$  and standard deviation  $\sigma$ :

$$\begin{aligned} f(x) &= \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-\mu)^2/2\sigma^2}, && -\infty < x < \infty \\ F(a) &= \int_{-\infty}^a f(x)dx = 0.5 \left( 1 + \text{erf} \left[ \frac{a - \mu}{\sqrt{2}\sigma} \right] \right) \end{aligned}$$

**Lognormal** distribution on  $0 < x < \infty$  with log-mean  $\mu$  and log standard deviation  $\sigma$ :

$$\log(X) = n(\mu, \sigma), \quad \mu_X = \exp\left(\mu + \frac{\sigma^2}{2}\right), \quad \sigma_X = \mu_X \sqrt{\exp(\sigma^2) - 1}.$$

### 2.2.1 Generating random variables in Matlab

**rand(m,n)** returns an  $m \times n$  matrix of random numbers from a uniform distribution on  $(0, 1)$ .

**randn(m,n)** returns an  $m \times n$  matrix of normally-distributed random numbers with mean 0 and standard deviation 1. Fig. 1 shows a histogram of the results of **randn(1,1000)**.

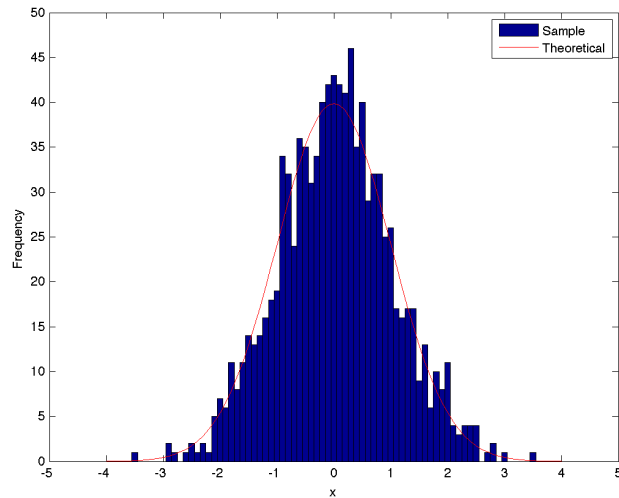


Figure 1: Histogram of 1000 samples of a normal distribution

**random(name,params,[m,n,...])** (Statistics toolbox) returns an  $m \times n \times \dots$  array of random numbers with a pdf described by **name** and **params**, (e.g. 'Binomial',N,p or 'Lognormal',mu,sigma)

## 2.3 Joint distributions

**Joint cumulative distribution** of two RVs  $X$  and  $Y$  can be phrased in terms of their joint CDF

$$F(a, b) = P(X \leq a, Y \leq b)$$

**Joint PDF**  $f(x, y)$  of two continuous RVs:  $f(x, y)dxdy$  is the probability that  $x - dx/2 < X < x + dx/2, y - dy/2 < Y < y + dy/2$ .

**Two RVs are independent** iff

$$F(a, b) = F_X(a)F_Y(b) \quad \forall a, b \quad \text{or} \quad f(x, y) = f_X(x)f_Y(y) \quad \forall x, y$$

**Covariance** of  $X$  and  $Y$ :

$$\text{cov}[X, Y] = E[(X - \bar{X})(Y - \bar{Y})]. \quad (2.3.1)$$

If  $X$  and  $Y$  are independent,  $\text{cov}[X, Y] = 0$  (but not necessarily vice-versa). Note  $\text{cov}[X, X] = \text{var}[X]$  and  $\text{cov}[X, Y + Z] = \text{cov}[X, Y] + \text{cov}[X, Z]$ .

**Correlation coefficient**

$$R_{XY} = \frac{\text{cov}[X, Y]}{\sigma_X \sigma_Y} \quad (2.3.2)$$

$R$  lies between -1 and 1;  $R = 1$  if  $Y = X$  (perfect correlation),  $R = -1$  if  $Y = -X$  (perfect anticorrelation), and  $R = 0$  if  $X$  and  $Y$  are independent. Unlike covariance,  $R$  is *not* additive.

**The correlation coefficient is useful for describing how strongly  $X$  and  $Y$  are linearly related, but will not perfectly capture non-linear relationships between  $X$  and  $Y$ . In particular, unless  $X$  and  $Y$  are Gaussian, they can be uncorrelated ( $R = 0$ ) yet still be dependent.** For instance, let  $\Theta$  be a uniformly distributed RV over  $[0, 2\pi)$  and let  $X = \cos(\Theta), Y = \sin(\Theta)$  (Fig. 2). Then  $X$  and  $Y$  each have mean zero and they are easily shown to be uncorrelated. However, for any given value  $x$  of  $X$ ,  $Y$  can take only the two values  $\pm(1 - x^2)^{1/2}$  (with equal probability), so  $Y$  is not independent of the value of  $X$ .

**The mean is always additive, and the variance is additive for independent (or uncorrelated) RVs:**

$$E[X + Y] = E[X] + E[Y] \quad (\overline{X + Y} = \bar{X} + \bar{Y}) \quad (2.3.3)$$

$$\begin{aligned} \text{var}[X + Y] &= E[(X + Y - \bar{X} - \bar{Y})^2] \\ &= E[(X - \bar{X})^2] + 2E[(X - \bar{X})(Y - \bar{Y})] + E[(Y - \bar{Y})^2] \\ &= \text{var}[X] + \text{var}[Y] + 2\text{cov}[X, Y] \end{aligned} \quad (2.3.4)$$

$$= \text{var}[X] + \text{var}[Y] \quad \text{if } \text{cov}[X, Y] = 0. \quad (2.3.5)$$

## 2.4 Sample mean and standard deviation

Given  $N$  independent samples  $x_1, x_2, \dots, x_N$  of a random variable  $X$ , we can estimate basic statistics of  $X$ :

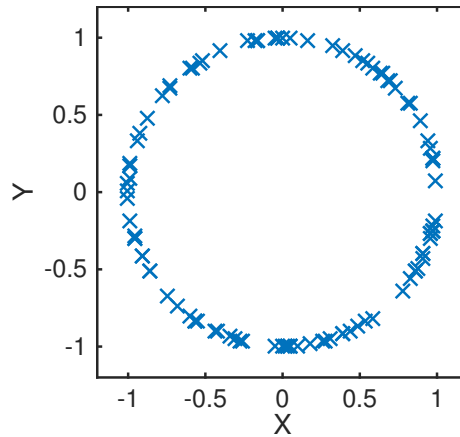


Figure 2: 100 samples of two RVs  $X$  and  $Y$  which are uncorrelated but dependent

### Sample mean

$$\bar{x} = \frac{1}{N} \sum_{j=1}^N x_j \quad (2.4.1)$$

The sample mean is an *estimator* of the true mean  $\bar{X}$  of  $X$ . We will quantify the accuracy of this estimator vs.  $N$  later. For now, we note that the sample mean is an *unbiased* estimator of  $\bar{X}$ , i. e.,  $E[\bar{x}] = \bar{X}$ .

**Sample standard deviation**  $\sigma(x)$  We calculate the variance of the  $x_j$  about the *sample* mean  $\bar{x}$ . Computing the mean from the sample reduces the *effective sample size* (often called the *degrees of freedom* or DOF) by one to  $N - 1$ :

$$\sigma^2(x) = \text{var}(x) = \frac{1}{N-1} \sum_{j=1}^N (x_j - \bar{x})^2 \quad (2.4.2)$$

If the samples are not independent, the effective sample size must be adjusted (Lecture 4). Otherwise  $\sigma^2(x)$  is an unbiased estimator of the true variance  $\sigma_X^2$  of  $X$ .

**Sample covariance and correlation coefficient** between independent samples  $x_j$  of RV  $X$  and corresponding samples  $y_j$  of RV  $Y$ :

$$\text{cov}(x, y) = \frac{1}{N-1} \sum_{j=1}^N (x_j - \bar{x})(y_j - \bar{y}) \quad (2.4.3)$$

(which is an unbiased estimator of the true covariance between  $X$  and  $Y$ );

$$R(x, y) = \frac{\text{cov}(x, y)}{\sigma(x)\sigma(y)} \quad (2.4.4)$$

(*not* usually an unbiased estimator of the true correlation coefficient  $R_{XY}$ .)

### 2.4.1 Matlab for sample statistics

If we arrange the  $x_j$  into a column vector  $\mathbf{x}$ :

**mean(x)** Sample mean.

**std(x), var(x)** Unbiased standard deviation and variance estimators.

For an array  $\mathbf{X}$  these are calculated along the first dimension (the column dimension of a matrix) unless specified otherwise with an optional argument. To get the mean of an array use **mean(X(:))**, i. e. reshape the array into a single vector.

**cov(x,y), corrcoef(x,y)** Given two column data vectors  $\mathbf{x}$  and  $\mathbf{y}$ , these return 2x2 matrices whose off-diagonal (2,1) and (1,2) elements are the sample covariance (or correlation coefficient).

**cov(X), corrcoef(X)** Let  $\mathbf{X}$  be a  $K \times N$  data array whose  $K$  columns  $\mathbf{x}_k$  correspond to different variables, so that  $X_{nk}$  is the  $n$ 'th sample of variable  $k$ . Then these functions return  $K \times K$  matrices whose  $(k,l)$  entry is the sample covariance  $\text{cov}[\mathbf{x}_k, \mathbf{x}_l]$  (or correlation coefficient).