

## Statistics Formulae

### Mean and Variance

A random variable  $X$  has a distribution over some subset  $x$  of the real numbers. When the distribution of  $X$  is discrete, the probability that  $X = x_i$  is  $P_i$ . When the distribution is continuous, the probability that  $X$  lies in an interval  $\delta x$  is  $f(x)\delta x$ , where  $f(x)$  is the probability density function.

$$\text{Mean } \mu = E(X) = \sum P_i x_i \text{ or } \int x f(x) dx.$$

$$\text{Variance } \sigma^2 = V(X) = E[(X - \mu)^2] = \sum P_i (x_i - \mu)^2 \text{ or } \int (x - \mu)^2 f(x) dx.$$

### Probability distributions

Error function:  $\text{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x e^{-y^2} dy$

Binomial:  $f(x) = \binom{n}{x} p^x q^{n-x}$  where  $q = (1 - p)$ ,  $\mu = np$ ,  $\sigma^2 = npq$ ,  $p < 1$ .

Poisson:  $f(x) = \frac{\mu^x}{x!} e^{-\mu}$ , and  $\sigma^2 = \mu$

Normal:  $f(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{(x - \mu)^2}{2\sigma^2}\right]$

### Weighted sums of random variables

If  $W = aX + bY$  then  $E(W) = aE(X) + bE(Y)$ . If  $X$  and  $Y$  are independent then  $V(W) = a^2V(X) + b^2V(Y)$ .

### Statistics of a data sample $x_1, \dots, x_n$

$$\text{Sample mean } \bar{x} = \frac{1}{n} \sum x_i$$

$$\text{Sample variance } s^2 = \frac{1}{n} \sum (x_i - \bar{x})^2 = \left(\frac{1}{n} \sum x_i^2\right) - \bar{x}^2 = E(x^2) - [E(x)]^2$$

### Regression (least squares fitting)

To fit a straight line by least squares to  $n$  pairs of points  $(x_i, y_i)$ , model the observations by  $y_i = \alpha + \beta(x_i - \bar{x}) + \epsilon_i$ , where the  $\epsilon_i$  are independent samples of a random variable with zero mean and variance  $\sigma^2$ .

$$\text{Sample statistics: } s_x^2 = \frac{1}{n} \sum (x_i - \bar{x})^2, \quad s_y^2 = \frac{1}{n} \sum (y_i - \bar{y})^2, \quad s_{xy}^2 = \frac{1}{n} \sum (x_i - \bar{x})(y_i - \bar{y}).$$

$$\text{Estimators: } \hat{\alpha} = \bar{y}, \hat{\beta} = \frac{s_{xy}^2}{s_x^2}; E(Y \text{ at } x) = \hat{\alpha} + \hat{\beta}(x - \bar{x}); \hat{\sigma}^2 = \frac{n}{n-2} (\text{residual variance}),$$

$$\text{where residual variance} = \frac{1}{n} \sum \{y_i - \hat{\alpha} - \hat{\beta}(x_i - \bar{x})\}^2 = s_y^2 - \frac{s_{xy}^4}{s_x^2}.$$

$$\text{Estimates for the variances of } \hat{\alpha} \text{ and } \hat{\beta} \text{ are } \frac{\hat{\sigma}^2}{n} \text{ and } \frac{\hat{\sigma}^2}{ns_x^2}.$$

Correlation coefficient:  $\hat{\rho} = r = \frac{s_{xy}^2}{s_x s_y}$ .