

Biostatistics

Lecture 4

Probability Distributions

Continuous Probability Distributions

- Continuous distribution has an infinite number of values between any two values assumed by the continuous variable
- As the number of observations, n , approaches infinity, the the width of the class intervals approaches zero, the graph of the frequencies (frequency polygon) approaches a smooth curve

Continuous Probability Distributions

- As with other probability distributions, the total area under the curve equals 1
- Relative frequency (probability) of occurrence of values between any two points on the x-axis is equal to the total area bounded by the curve, the x-axis, and perpendicular lines erected at the two points on the x-axis
- Probability of any specific value of the random variable is 0

Area under a smooth curve

- Integration of the density function over the range a to b
 - Density function is a formula used to represent the distribution of a continuous random variable
- A nonnegative function $f(x)$ is called a probability distribution or probability density function of the continuous random variable X if the total area bounded by its curve and the x -axis is equal 1 and if the subarea under the curve bounded by the curve, the x -axis, and perpendiculars erected at any two points a and b gives the probability that X is between the points a and b

Normal distribution

- Most important distribution in statistics
- Also called the **Gaussian distribution**
- Density given by

$$f(x) = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

- for $-\infty < x < \infty$
- where μ is the mean and σ the standard deviation

Characteristics of Normal Distribution

- Symmetrical about mean, μ
- Mean, median, and mode are equal
- Total area under the curve above the x-axis is one square unit
- 1 standard deviation on both sides of the mean includes approximately 68% of the total area
 - 2 standard deviations includes approximately 95%
 - 3 standard deviations includes approximately 99%

Characteristics of the Normal Distribution

- Normal distribution is completely determined by the parameters μ and σ
 - Different values of μ shift the distribution along the x-axis
 - Different values of σ determine degree of flatness or peakedness of the graph

Standard Normal Distribution

- Normal distribution is really family of curves determined by μ and σ
- Standard normal distribution is one with a $\mu = 0$ and $\sigma = 1$
 - Standard normal density given by:

$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-z^2/2}$$

- for $-\infty < x < \infty$
- where $z = (x - \mu) / \sigma$

Standard Normal Distribution

- To find probability that z takes on a value between any two points on the z -axis, need to find area bounded by perpendiculars erected at these points, the curve, and the z -axis
 - Values are tabled (Appendix Table 3)
 - Standard normal distribution is symmetric

Applications of Normal Distribution

- Frequently, data are normally distributed
 - Essential for some statistical procedures
 - Due to the combination of many random factors/influences having different effects on the outcome
 - If not, possible to transform to a more normal form
- Approximations for other distributions
- Because of the frequent occurrence of the **normal distribution** in nature, much statistical theory has been developed for it

Examples of Standard Normal Distribution

- Height and weight
 - Calculate z-statistics
 - $\Pr(X < x)$
 - $\Pr(X > x)$
 - $\Pr(x_1 < X < x_2)$
- Why?
 - Determine percentiles
 - Comparisons between different distributions

Linear Combination of Random Variables

- Frequently, we work with combinations of variables; ie., variables ‘added’ or combined together in some way
- Definition:
 - $L = c_1X_1 + c_2X_2 + \dots + c_nX_n$
- also known as **linear contrast**
- Sum = $1X_1 + 1X_2$
- Difference = $1X_1 + (-1)X_2$
- Mean = $0.5X_1 + 0.5X_2$

Linear Combination of Random Variables

- Expected value (typically, the mean) of a linear combination is the sum of the expected values of the variables
 - $E(L) = c_1 E(X_1) + c_2 E(X_2) + \dots$
 - $E(L) = \sum_i c_i X_i$
- Variance of a linear combination of independent random variables is
 - $\text{Var}(L) = \sum_i c_i^2 \text{var}(X_i)$
- For normally distributed random variables,
 - $E(L) = \sum_i c_i \mu_i$
 - $\text{Var}(L) = \sum_i c_i^2 \sigma_i^2$

Linear Combination of Random Variables

- Variables that are not independent have a **covariance** between each pair
 - $\text{Cov}(X, Y) = E[(X - \mu_x)(Y - \mu_y)]$
- Covariance hard to interpret directly, so we calculate the **correlation coefficient** for descriptive purposes
 - $\rho = \text{Corr}(X, Y) = \text{Cov}(X, Y) / (\sigma_x \sigma_y)$
- Correlation coefficient ranges from -1 to 1 and is without units
 - Frequently used to describe strength of relationship between variables

Linear Combination of Random Variables

- Variance of a linear combination takes into account the covariance
 - $\text{Var}(c_1X_1 + c_2X_2) = c_1^2 \text{var}(X_1) + c_2^2 \text{var}(X_2) + 2c_1c_2\text{Cov}(X_1, X_2)$
 - Typically, we calculate the variance-covariance matrix for these calculations
- Note that the mean of a linear combination is not affected by the covariance 'structure' of the combination, only the variance
 - Important for regression analysis and analysis of correlated analysis

Normal Approximation to the Binomial Distribution

- When n is moderate in size ($n \geq 25$) and p is not too extreme, then the normal distribution is a good approximation to the binomial distribution with $\mu = np$ and $\sigma^2 = npq$
- $\Pr(a \leq X \leq b)$ is approximately the area under the $N(np, npq)$ curve between $a-1/2$ and $b+1/2$

Normal Approximation to the Poisson Distribution

- When μ is moderate in size ($\mu \geq 10$), the Poisson distribution is cumbersome to use and the normal distribution is a good approximation with $\mu = \mu$ and $\sigma^2 = \mu$
- $\Pr(X = x)$ is approximately the area under the $N(\mu, \mu)$ curve between $x-1/2$ and $x+1/2$

Summary

- Laid the ground work for 'normal theory' estimation and hypothesis testing
- Showed how the binomial and Poisson distribution generalize to the normal with large samples
- Investigated the properties of linear contrasts, important for ANOVA and regression