

# **Probability Distributions and Statistical Inference in R**

Instructor: Dr. Azadeh Alimadad  
Department of Statistics & Data Analytics

## **Course Overview**

This lecture introduces probability distributions and statistical inference using R. We'll cover discrete and continuous distributions, hypothesis testing and confidence intervals. Practical R implementations will be demonstrated throughout.

## **Contents**

# 1 Introduction to Probability Distributions in R

## 1.1 R's Distribution Functions

R uses a consistent naming convention for probability distributions:

- **d**: Density/mass function (density)
- **p**: Cumulative distribution function (probability)
- **q**: Quantile function (quantile)
- **r**: Random number generation (random)

## 1.2 Common Distribution Functions

```

1 # General syntax
2 dxxxx(x, parameters)      \# Density/Probability at x
3 pxxxx(q, parameters)      \# P(X $\leq$ q)
4 qxxxx(p, parameters)      \# Quantile for probability p
5 rxxxx(n, parameters)      \# Generate n random values

```

# 2 Discrete Probability Distributions

## 2.1 Binomial Distribution

### 2.1.1 Theory and Formula

The binomial distribution models the number of successes in  $n$  independent trials, each with probability of success  $p$ .

$$P(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}, \quad k = 0, 1, 2, \dots, n \quad (1)$$

Where:

- $n$ : number of trials
- $p$ : probability of success
- $k$ : number of successes

### 2.1.2 Mean and Variance

$$E[X] = np \quad (2)$$

$$Var[X] = np(1 - p) \quad (3)$$

### 2.1.3 R Implementation

```

1 # Binomial distribution functions
2 n <- 10
3 p <- 0.3
4
5 # Probability mass function
6 dbinom(3, size = n, prob = p) # P(X = 3)
7 dbinom(0:10, size = n, prob = p) # All probabilities
8
9 # Cumulative distribution function
10 pbinom(4, size = n, prob = p) # P(X $\leq$ 4)
11 pbinom(4, size = n, prob = p, lower.tail = FALSE) # P(X $>$ 4)
12
13 # Quantile function
14 qbinom(0.95, size = n, prob = p) # 95th percentile
15
16 # Random number generation
17 rbinom(100, size = n, prob = p) # 100 random binomial values
18
19 # Plot binomial distribution
20 x <- 0:n
21 prob <- dbinom(x, size = n, prob = p)
22 plot(x, prob, type = "h", lwd = 2,
23       main = "Binomial Distribution (n=10, p=0.3)",
24       xlab = "Number of Successes",
25       ylab = "Probability",
26       col = "blue")
27 points(x, prob, pch = 16, col = "red")

```

## 2.2 Poisson Distribution

### 2.2.1 Theory and Formula

The Poisson distribution models the number of events occurring in a fixed interval of time or space.

$$P(X = k) = \frac{\lambda^k e^{-\lambda}}{k!}, \quad k = 0, 1, 2, \dots \quad (4)$$

Where  $\lambda$  is the average rate of occurrence.

### 2.2.2 Mean and Variance

$$E[X] = \lambda \quad (5)$$

$$Var[X] = \lambda \quad (6)$$

### 2.2.3 R Implementation

```

1 # Poisson distribution functions
2 lambda <- 3.5
3
4 # Probability mass function
5 dpois(2, lambda)          # P(X = 2)
6 dpois(0:10, lambda)       # First 11 probabilities
7
8 # Cumulative distribution function
9 ppois(5, lambda)          # P(X $\leq$ 5)
10
11 # Quantile function
12 qpois(0.9, lambda)        # 90th percentile
13
14 # Random number generation
15 rpois(100, lambda)        # 100 random Poisson values
16
17 # Plot Poisson distribution
18 x <- 0:15
19 prob <- dpois(x, lambda)
20 plot(x, prob, type = "h", lwd = 2,
21       main = paste("Poisson Distribution ($\lambda$ =", lambda, ")"),
22       xlab = "Number of Events",
23       ylab = "Probability",
24       col = "darkgreen")
25 points(x, prob, pch = 16, col = "red")
26
27 # Example: Customer arrivals
28 # Suppose customers arrive at average rate of 4.2 per hour
29 arrivals <- rpois(100, lambda = 4.2)
30 hist(arrivals, breaks = seq(-0.5, max(arrivals)+0.5, 1),
31       main = "Customer Arrivals per Hour",
32       xlab = "Number of Customers",
33       ylab = "Frequency",
34       col = "lightblue")

```

### 3 Continuous Probability Distributions

#### 3.1 Normal Distribution

##### 3.1.1 Theory and Formula

The probability density function of the normal distribution is:

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2}, \quad -\infty < x < \infty \quad (7)$$

Where:

- $\mu$ : mean

- $\sigma$ : standard deviation

### 3.1.2 Standard Normal Distribution

When  $\mu = 0$  and  $\sigma = 1$ , we have the standard normal distribution  $Z \sim N(0, 1)$ .

### 3.1.3 R Implementation

```

1 # Normal distribution functions
2 mu <- 100
3 sigma <- 15
4
5 # Density function
6 dnorm(110, mean = mu, sd = sigma)           # f(110)
7 dnorm(seq(70, 130, by = 5), mu, sigma)       # Multiple values
8
9 # Cumulative distribution function
10 pnorm(110, mean = mu, sd = sigma)            # P(X $\leq$ 110)
11 pnorm(110, mean = mu, sd = sigma, lower.tail = FALSE) # P(X $>$ 110)
12
13 # Standard normal
14 pnorm(1.96)                                # P(Z $\leq$ 1.96) =
15             0.975
16 pnorm(-1.96)                               # P(Z $\leq$ -1.96) =
17             0.025
18
19 # Quantile function
20 qnorm(0.95, mean = mu, sd = sigma)          # 95th percentile
21 qnorm(0.975)                                # z-score for 95% CI:
22             1.96
23 qnorm(c(0.025, 0.975))                     # Critical values for
24             95% CI
25
26 # Random number generation
27 rnorm(100, mean = mu, sd = sigma)            # 100 random values
28
29 # Plot normal distribution
30 x <- seq(50, 150, length.out = 200)
31 pdf <- dnorm(x, mean = mu, sd = sigma)
32 plot(x, pdf, type = "l", lwd = 2,
33       main = "Normal Distribution ($\mu=100, \sigma=15)", 
34       xlab = "x",
35       ylab = "Density",
36       col = "blue")
37
38 # Shade areas
39 # P(85 $\leq$ X $\leq$ 115)
40 x_shade <- seq(85, 115, length.out = 100)
41 y_shade <- dnorm(x_shade, mean = mu, sd = sigma)
42 polygon(c(85, x_shade, 115), c(0, y_shade, 0),
43

```

39 |      col = rgb(0, 0, 1, 0.3), border = NA)

## 3.2 t-Distribution

### 3.2.1 Theory and Formula

The t-distribution is used for small sample sizes when population variance is unknown.

$$f(t) = \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\sqrt{\nu\pi}\Gamma\left(\frac{\nu}{2}\right)} \left(1 + \frac{t^2}{\nu}\right)^{-\frac{\nu+1}{2}} \quad (8)$$

Where  $\nu$  is degrees of freedom.

### 3.2.2 R Implementation

```

1 # t-distribution functions
2 df <- 10 # degrees of freedom
3
4 # Density function
5 dt(1.5, df = df) # f(1.5)
6
7 # Cumulative distribution function
8 pt(2.0, df = df) # P(t $\leq$ 2.0)
9 pt(2.0, df = df, lower.tail = FALSE) # P(t $>$ 2.0)
10
11 # Critical values for confidence intervals
12 qt(0.975, df = df) # t-value for 95% CI
13 qt(0.995, df = 20) # t-value for 99% CI with df=20
14
15 # Compare t-distribution with normal distribution
16 x <- seq(-4, 4, length.out = 200)
17 plot(x, dnorm(x), type = "l", lwd = 2, col = "black",
18       main = "t-Distribution vs Normal Distribution",
19       xlab = "x", ylab = "Density")
20 lines(x, dt(x, df = 1), col = "red", lwd = 2)
21 lines(x, dt(x, df = 5), col = "blue", lwd = 2)
22 lines(x, dt(x, df = 30), col = "green", lwd = 2)
23 legend("topright",
24        legend = c("Normal", "t (df=1)", "t (df=5)", "t (df=30)"),
25        col = c("black", "red", "blue", "green"), lwd = 2)
26
27 # t-test example
28 set.seed(123)
29 sample_data <- rnorm(15, mean = 105, sd = 10)
30 t_test_result <- t.test(sample_data, mu = 100)
31 print(t_test_result)
32 cat("t-statistic:", t_test_result$statistic, "\n")
33 cat("p-value:", t_test_result$p.value, "\n")
34 cat("95% CI:", t_test_result$conf.int, "\n")

```

## 4 Chi-Square Distribution

### 4.1 Theory and Formula

The chi-square distribution with  $k$  degrees of freedom is the sum of squares of  $k$  independent standard normal variables.

$$f(x; k) = \frac{1}{2^{k/2}\Gamma(k/2)}x^{k/2-1}e^{-x/2}, \quad x > 0 \quad (9)$$

### 4.2 Applications

1. Confidence interval for variance
2. Goodness-of-fit tests
3. Test of independence
4. Test of homogeneity

### 4.3 R Implementation

```

1 # Chi-square distribution functions
2 df <- 5 # degrees of freedom
3
4 # Density function
5 dchisq(3, df = df) # f(3)
6
7 # Cumulative distribution function
8 pchisq(10, df = df) # P(chi^2 \leq 10)
9
10 # Critical values
11 qchisq(0.95, df = df) # 95th percentile
12 qchisq(c(0.025, 0.975), df = df) # Critical values for 95% CI
13
14 # Plot chi-square distributions
15 x <- seq(0, 20, length.out = 200)
16 plot(x, dchisq(x, df = 1), type = "l", lwd = 2, col = "red",
17       main = "Chi-Square Distributions",
18       xlab = "x", ylab = "Density",
19       ylim = c(0, 0.5))
20 lines(x, dchisq(x, df = 3), col = "blue", lwd = 2)
21 lines(x, dchisq(x, df = 5), col = "green", lwd = 2)
22 lines(x, dchisq(x, df = 10), col = "orange", lwd = 2)
23 legend("topright",
24        legend = c("df=1", "df=3", "df=5", "df=10"),
25        col = c("red", "blue", "green", "orange"), lwd = 2)

```

## 4.4 Confidence Interval for Variance

### 4.4.1 Theory

For a sample from a normal distribution:

$$\frac{(n-1)s^2}{\sigma^2} \sim \chi^2_{n-1} \quad (10)$$

The  $(1 - \alpha)$  confidence interval for  $\sigma^2$  is:

$$\left[ \frac{(n-1)s^2}{\chi^2_{1-\alpha/2, n-1}}, \frac{(n-1)s^2}{\chi^2_{\alpha/2, n-1}} \right] \quad (11)$$

### 4.4.2 R Implementation

```

1 # Function for confidence interval of variance
2 ci_variance <- function(data, alpha = 0.05) {
3   n <- length(data)
4   s2 <- var(data)
5   df <- n - 1
6
7   # Chi-square critical values
8   chi_lower <- qchisq(alpha/2, df, lower.tail = FALSE)
9   chi_upper <- qchisq(1 - alpha/2, df, lower.tail = FALSE)
10
11  # Confidence interval for variance
12  ci_var <- c((n-1)*s2/chi_lower, (n-1)*s2/chi_upper)
13
14  # Confidence interval for standard deviation
15  ci_sd <- sqrt(ci_var)
16
17  return(list(variance = s2,
18             sd = sd(data),
19             ci_variance = ci_var,
20             ci_sd = ci_sd))
21}
22
23 # Example
24 set.seed(123)
25 sample_data <- rnorm(25, mean = 100, sd = 15)
26 result <- ci_variance(sample_data, alpha = 0.05)
27
28 cat("Sample variance:", result$variance, "\n")
29 cat("Sample standard deviation:", result$sd, "\n")
30 cat("95% CI for variance:", result$ci_variance, "\n")
31 cat("95% CI for standard deviation:", result$ci_sd, "\n")
32
33 # Verify with built-in function
34 # install.packages("DescTools")
35 library(DescTools)
36 VarTest(sample_data)

```

## 4.5 Goodness-of-Fit Test

### 4.5.1 Theory

Tests whether observed frequencies match expected frequencies:

$$\chi^2 = \sum_{i=1}^k \frac{(O_i - E_i)^2}{E_i} \sim \chi^2_{k-1} \quad (12)$$

### 4.5.2 R Implementation

```

1 # Goodness-of-fit test example
2
3 # Example 1: Dice fairness test
4 observed <- c(22, 18, 15, 19, 24, 22) # Observed frequencies for
   6 faces
5 expected <- rep(20, 6)                  # Expected frequencies if
   fair
6
7 chi_test <- chisq.test(observed, p = rep(1/6, 6))
8 print(chi_test)
9
10 # Extract components
11 cat("Chi-square statistic:", chi_test$statistic, "\n")
12 cat("p-value:", chi_test$p.value, "\n")
13 cat("Degrees of freedom:", chi_test$parameter, "\n")
14
15 # Example 2: Mendelian genetics
16 # 9:3:3:1 ratio expected
17 observed_plants <- c(315, 101, 108, 32) # Observed counts
18 expected_ratio <- c(9/16, 3/16, 3/16, 1/16) # Expected ratios
19
20 chi_plants <- chisq.test(observed_plants, p = expected_ratio)
21 print(chi_plants)
22
23 # Manual calculation
24 total <- sum(observed_plants)
25 expected_counts <- total * expected_ratio
26 chi_sq <- sum((observed_plants - expected_counts)^2 / expected_
   counts)
27 p_value <- pchisq(chi_sq, df = 3, lower.tail = FALSE)
28
29 cat("\nManual calculation:\n")
30 cat("Chi-square:", chi_sq, "\n")
31 cat("p-value:", p_value, "\n")

```

## 5 F-Distribution and ANOVA

### 5.1 Theory and Formula

The F-distribution arises as the ratio of two independent chi-square variables divided by their degrees of freedom.

$$F = \frac{U_1/d_1}{U_2/d_2} \sim F(d_1, d_2) \quad (13)$$

Where  $U_1 \sim \chi^2_{d_1}$  and  $U_2 \sim \chi^2_{d_2}$ .

### 5.2 R Implementation

```

1 # F-distribution functions
2 df1 <- 3    # numerator degrees of freedom
3 df2 <- 20   # denominator degrees of freedom
4
5 # Density function
6 df(2.5, df1 = df1, df2 = df2)      # f(2.5)
7
8 # Cumulative distribution function
9 pf(3.0, df1 = df1, df2 = df2)      # P(F $\leq$ 3.0)
10
11 # Critical values
12 qf(0.95, df1 = df1, df2 = df2)      # 95th percentile
13 qf(0.99, df1 = 2, df2 = 15)        # Critical value for $\alpha$ =
14     0.01
15
16 # Plot F-distributions
17 x <- seq(0, 5, length.out = 200)
18 plot(x, df(x, df1 = 1, df2 = 10), type = "l", lwd = 2, col = "red",
19       ,
20       main = "F-Distributions",
21       xlab = "x", ylab = "Density",
22       ylim = c(0, 1))
23 lines(x, df(x, df1 = 3, df2 = 10), col = "blue", lwd = 2)
24 lines(x, df(x, df1 = 5, df2 = 10), col = "green", lwd = 2)
25 lines(x, df(x, df1 = 10, df2 = 10), col = "orange", lwd = 2)
26 legend("topright",
27       legend = c("F(1,10)", "F(3,10)", "F(5,10)", "F(10,10)"),
28       col = c("red", "blue", "green", "orange"), lwd = 2)

```

## 5.3 One-Way ANOVA

### 5.3.1 Theory

Tests whether means of several groups are equal:

$$SS_{total} = \sum_{i=1}^k \sum_{j=1}^{n_i} (x_{ij} - \bar{x})^2 \quad (14)$$

$$SS_{between} = \sum_{i=1}^k n_i (\bar{x}_i - \bar{x})^2 \quad (15)$$

$$SS_{within} = \sum_{i=1}^k \sum_{j=1}^{n_i} (x_{ij} - \bar{x}_i)^2 \quad (16)$$

$$F = \frac{MS_{between}}{MS_{within}} = \frac{SS_{between}/(k-1)}{SS_{within}/(N-k)} \sim F(k-1, N-k) \quad (17)$$

### 5.3.2 R Implementation

```

1 # One-way ANOVA example
2
3 # Create sample data
4 set.seed(123)
5 group1 <- rnorm(20, mean = 100, sd = 10)
6 group2 <- rnorm(20, mean = 110, sd = 10)
7 group3 <- rnorm(20, mean = 105, sd = 10)
8 group4 <- rnorm(20, mean = 115, sd = 10)
9
10 # Combine into data frame
11 data_anova <- data.frame(
12   value = c(group1, group2, group3, group4),
13   group = factor(rep(1:4, each = 20))
14 )
15
16 # Perform ANOVA
17 anova_result <- aov(value ~ group, data = data_anova)
18 summary(anova_result)
19
20 # Extract F-statistic and p-value
21 f_stat <- summary(anova_result)[[1]]$"F value"[1]
22 p_value <- summary(anova_result)[[1]]$"Pr(>F)"[1]
23
24 cat("F-statistic:", f_stat, "\n")
25 cat("p-value:", p_value, "\n")
26
27 # Manual calculation for understanding
28 group_means <- tapply(data_anova$value, data_anova$group, mean)
29 overall_mean <- mean(data_anova$value)
30 n_groups <- length(unique(data_anova$group))
31 n_total <- nrow(data_anova)
32

```

```

33 # Calculate sums of squares
34 ss_between <- sum(table(data_anova$group) * (group_means -
35   overall_mean)^2)
35 ss_within <- sum(tapply(data_anova$value, data_anova$group,
36     function(x) sum((x - mean(x))^2)))
37 ss_total <- sum((data_anova$value - overall_mean)^2)
38
39 # Calculate mean squares
40 ms_between <- ss_between / (n_groups - 1)
41 ms_within <- ss_within / (n_total - n_groups)
42
43 # Calculate F-statistic
44 f_calc <- ms_between / ms_within
45 p_calc <- pf(f_calc, df1 = n_groups - 1, df2 = n_total - n_groups
46   ,
47   lower.tail = FALSE)
48
48 cat("\nManual calculation:\n")
49 cat("SS_between:", ss_between, "\n")
50 cat("SS_within:", ss_within, "\n")
51 cat("MS_between:", ms_between, "\n")
52 cat("MS_within:", ms_within, "\n")
53 cat("F-statistic:", f_calc, "\n")
54 cat("p-value:", p_calc, "\n")
55
56 # Post-hoc test (Tukey's HSD)
57 if (p_value < 0.05) {
58   tukey_result <- TukeyHSD(anova_result)
59   print(tukey_result)
60   plot(tukey_result)
61 }
```

## 6 Exercises

### 6.1 Basic Exercises

1. Binomial Distribution:

- Simulate 1000 coin flips (fair coin) and plot the distribution
- Calculate probability of getting 7 or more heads in 10 flips
- Find the number of successes corresponding to the 95th percentile

2. Normal Distribution:

- Generate 500 values from  $N(100, 15)$  and verify empirical vs theoretical probabilities
- Calculate the probability that a value is between 85 and 115
- Find the 99th percentile

### 3. t-Distribution:

- Compare t-distribution with 5, 15, and 30 df to standard normal
- Calculate critical values for 90%, 95%, and 99% confidence intervals with  $df=10$

## 6.2 Intermediate Exercises

### 1. Chi-Square Applications:

- Test if a 6-sided die is fair using 60 rolls with observed frequencies (9, 11, 10, 8, 12, 10)
- Calculate 95% CI for variance of a sample from  $N(0,1)$  with  $n=25$

### 2. ANOVA:

- Simulate data for 3 groups with different means and perform ANOVA
- Check ANOVA assumptions (normality, homogeneity of variance)
- Perform post-hoc tests if ANOVA is significant

## 6.3 Advanced Exercises

### 1. Comprehensive Analysis:

- Design a complete study with multiple groups and covariates
- Perform all appropriate statistical tests
- Create a complete report with interpretations

### 2. Power Analysis:

- Calculate required sample size for detecting effect sizes in t-tests and ANOVA
- Simulate power curves for different effect sizes

## 7 Resources and References

### 7.1 Books

- "Introduction to Probability and Statistics Using R" by G. Jay Kerns
- "The R Book" by Michael J. Crawley
- "Statistical Inference via Data Science" by Chester Ismay and Albert Y. Kim

### 7.2 Online Resources

- R Documentation: <https://www.rdocumentation.org/>
- Quick-R: <https://www.statmethods.net/>
- CRAN Task Views: <https://cran.r-project.org/web/views/>
- R-bloggers: <https://www.r-bloggers.com/>

### 7.3 Packages for Advanced Analysis

```
1 # Useful packages for probability and statistics
2 install.packages(c(
3   "ggplot2",          # Advanced plotting
4   "dplyr",            # Data manipulation
5   "tidyverse",         # Data tidying
6   "broom",             # Tidy model outputs
7   "infer",             # Tidy statistical inference
8   "effect",            # Effect displays
9   "lsr",                # Learning statistics with R
10  "pwr"                 # Power analysis
11 ))
```

## Conclusion

This lecture covered essential probability distributions and statistical inference methods in R. Understanding these concepts is crucial for data analysis and interpretation. Remember to:

1. Choose appropriate distributions for your data
2. Check assumptions before applying tests
3. Interpret results in context
4. Report estimates with confidence intervals
5. Visualize your data and results

Practice is key to mastering these concepts. Start with simple examples and gradually work on more complex analyses.

**Remember: Statistics is about understanding uncertainty, not eliminating it.**