# Bioinformatics Programming: A Comprehensive Review Guide

## Table of Contents

## Introduction to R Programming

### Setting Up Your Environment

RStudio is an Integrated Development Environment (IDE) for R that provides a user-friendly interface with multiple panes: - Script editor (top left) - Console (bottom left) - Environment/History (top right) - Files/Plots/Packages/Help (bottom right)

### Understanding Data Types in R

R has several fundamental data types that form the building blocks of data analysis:

```r
# 1. Numeric (includes both integers and floating-point numbers)
x <- 42        # Integer
y <- 3.14      # Floating-point
class(x)       # Shows the type of x
class(y)       # Shows the type of y

# 2. Character (strings)
name <- "Gene"
class(name)    # "character"

# 3. Logical (boolean)
is_active <- TRUE
class(is_active)   # "logical"

# 4. Factor (categorical data)
blood_types <- factor(c("A", "B", "O", "AB"))
class(blood_types) # "factor"
```

### Data Structures in R

R provides several key data structures for organizing information:

```r
# 1. Vectors (1-dimensional, same type)
numeric_vector <- c(1, 2, 3, 4, 5)
```

```r
char_vector <- c("gene1", "gene2", "gene3")

# 2. Lists (1-dimensional, different types)
my_list <- list(
    numbers = c(1, 2, 3),
    name = "Sample",
    is_valid = TRUE
)

# 3. Matrices (2-dimensional, same type)
my_matrix <- matrix(
    1:9,            # Data
    nrow = 3,       # Number of rows
    ncol = 3        # Number of columns
)

# 4. Data Frames (2-dimensional, different types per column)
df <- data.frame(
    gene_id = c("g1", "g2", "g3"),
    expression = c(100, 150, 200),
    is_significant = c(TRUE, FALSE, TRUE)
)
```

## Control Structures

### If Statements

```r
# Basic if-else structure
expression_value <- 150

if (expression_value > 100) {
    print("High expression")
} else if (expression_value > 50) {
    print("Medium expression")
} else {
    print("Low expression")
}

# Vectorized if-else with ifelse()
expression_levels <- c(80, 120, 40, 160)
categories <- ifelse(expression_levels > 100,
                     "High", "Low")
```

### For Loops

```r
# Basic for loop
genes <- c("BRCA1", "TP53", "EGFR")
```

```r
for (gene in genes) {
    cat("Processing gene:", gene, "\n")
}

# Loop with index
for (i in 1:length(genes)) {
    cat("Gene", i, "is", genes[i], "\n")
}
```

**File Input/Output**

```r
# Reading data
# CSV files
data <- read.csv("expression_data.csv",
                 header = TRUE,
                 row.names = 1)

# Writing data
write.csv(data,
          "processed_data.csv",
          row.names = TRUE)

# Reading tab-delimited files
tab_data <- read.delim("data.txt",
                       sep = "\t")

# Saving R objects
save(data, file = "analysis.RData")
load("analysis.RData")
```

# Advanced R Programming and Statistics

**Working with Logical Values**

```r
# Creating logical vectors
high_expression <- expression_levels > 100
is_significant <- p_values < 0.05

# Combining logical conditions
interesting_genes <- high_expression & is_significant

# Subsetting data based on logical vectors
significant_data <- data[is_significant, ]
```

### Exploratory Data Analysis (EDA)

```r
# Basic summary statistics
summary(data)

# Visual exploration with base R
hist(expression_values,
     main = "Distribution of Expression Values",
     xlab = "Expression Level")

boxplot(expression_values ~ conditions,
        main = "Expression by Condition")

# Using ggplot2 for advanced visualization
library(ggplot2)

ggplot(data, aes(x = condition, y = expression)) +
    geom_boxplot() +
    theme_minimal() +
    labs(title = "Gene Expression by Condition",
         x = "Experimental Condition",
         y = "Expression Level")
```

### Correlation Analysis

```r
# Computing correlation
correlation <- cor(gene1_expr, gene2_expr)

# Testing correlation significance
cor_test <- cor.test(gene1_expr, gene2_expr)

# Visualizing correlation
plot(gene1_expr, gene2_expr,
     main = "Gene Expression Correlation",
     xlab = "Gene 1",
     ylab = "Gene 2")
```

### Linear Models

```r
# Simple linear regression
model <- lm(expression ~ treatment, data = exp_data)
summary(model)

# Extracting model components
coefficients <- coef(model)
p_values <- summary(model)$coefficients[,4]
```

**Statistical Tests**

```r
# T-test
t_test_result <- t.test(group1, group2)

# ANOVA
anova_result <- aov(expression ~ condition,
                    data = exp_data)
summary(anova_result)
```

## Python for Bioinformatics

Since you're proficient in Python, here's a brief overview of key bioinformatics-related modules:

```python
# Key libraries for bioinformatics
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

# Reading biological data
expression_data = pd.read_csv("expression.csv")
```

## Bioinformatics Applications

### Bulk RNA-seq Analysis in R

```r
# Loading required packages
library(DESeq2)
library(edgeR)
library(ggplot2)

# Normalizing count data
dds <- DESeqDataSetFromMatrix(
    countData = counts,
    colData = sample_info,
    design = ~ condition
)

# Differential expression analysis
dds <- DESeq(dds)
results <- results(dds)
```

### Integration with Python

```python
# Reading RNA-seq results
bulk_results = pd.read_csv("bulk_results.csv")
```

```python
sc_results = pd.read_csv("sc_results.csv")

# Integration analysis
integrated_data = pd.merge(
    bulk_results,
    sc_results,
    on="gene_id",
    how="inner"
)
```

## Additional Resources

- CRAN: The Comprehensive R Archive Network
- Bioconductor: Repository for bioinformatics packages
- R Markdown: For reproducible research
- RStudio Cheatsheets
- Bioinformatics-focused R packages:
    - DESeq2
    - edgeR
    - limma
    - Seurat
    - biomaRt

## Exam Preparation Tips

1. Practice with real datasets
2. Review basic R syntax and data structures
3. Understand statistical concepts
4. Focus on RNA-seq analysis workflow
5. Practice data visualization
6. Understand integration between R and Python

Remember to: - Run example code - Modify parameters - Create your own test cases - Document your analysis steps