

Факторный анализ многомерных данных

Ортогональная модель с 3 факторами

Елисеев Данила, 2025, ИС

26 декабря 2025 г.

1. Теоретическая часть

Факторный анализ (Factor Analysis, FA) — это метод многомерного статистического анализа, направленный на выявление скрытых (латентных) факторов, объясняющих корреляции между наблюдаемыми переменными.

1.1. Ортогональная факторная модель

Ортогональная факторная модель имеет вид:

$$\mathbf{X} = \mathbf{\Lambda}\mathbf{F} + \mathbf{\Psi}\boldsymbol{\varepsilon}, \quad (1)$$

где:

- $\mathbf{X} \in \mathbb{R}^p$ — вектор наблюдаемых переменных
- $\mathbf{\Lambda} \in \mathbb{R}^{p \times m}$ — матрица факторных нагрузок
- $\mathbf{F} \in \mathbb{R}^m$ — вектор общих факторов (латентные переменные)
- $\mathbf{\Psi} \in \mathbb{R}^{p \times p}$ — диагональная матрица уникальных дисперсий
- $\boldsymbol{\varepsilon} \in \mathbb{R}^p$ — вектор уникальных факторов (ошибки)

Предположения модели:

- $E[\mathbf{F}] = \mathbf{0}$, $\text{Cov}(\mathbf{F}) = \mathbf{I}$ — факторы центрированы и некоррелированы
- $E[\boldsymbol{\varepsilon}] = \mathbf{0}$, $\text{Cov}(\boldsymbol{\varepsilon}) = \mathbf{I}$ — уникальные факторы центрированы и некоррелированы
- $\text{Cov}(\mathbf{F}, \boldsymbol{\varepsilon}) = \mathbf{0}$ — факторы и ошибки некоррелированы

Ковариационная матрица наблюдаемых переменных:

$$\boldsymbol{\Sigma} = \mathbf{\Lambda}\mathbf{\Lambda}^\top + \mathbf{\Psi}. \quad (2)$$

1.2. Интерпретация параметров

- **Факторные нагрузки** λ_{ij} — корреляция между переменной X_i и фактором F_j
- **Общность** (communality) $h_i^2 = \sum_{j=1}^m \lambda_{ij}^2$ — доля дисперсии переменной X_i , объясняемая общими факторами
- **Уникальность** ψ_i — доля дисперсии переменной X_i , не объясняемая общими факторами
- $h_i^2 + \psi_i = 1$ для стандартизированных переменных

1.3. Ротация факторов

Ротация факторов применяется для упрощения интерпретации. Наиболее распространённый метод — **Varimax ротация**, которая максимизирует дисперсию квадратов нагрузок по факторам, что приводит к:

- Факторам с несколькими высокими нагрузками и многими низкими
- Более простой интерпретации факторов
- Сохранению ортогональности факторов

2. Описание данных

Для анализа использован синтетический многомерный набор данных с $p = 8$ переменными и объёмом выборки $n = 300$. Переменные:

1. **GDP_Growth** — рост ВВП (экономический показатель)
2. **Industrial_Output** — промышленное производство (экономический показатель)
3. **Trade_Balance** — торговый баланс (экономический показатель)
4. **Education_Level** — уровень образования (социальный показатель)
5. **Healthcare_Index** — индекс здравоохранения (социальный показатель)
6. **Air_Quality** — качество воздуха (экологический показатель)
7. **Water_Quality** — качество воды (экологический показатель)
8. **Random_Noise** — случайный шум (независимая переменная)

Ожидаемая факторная структура:

- **Фактор 1 (Экономический):** переменные 1–3
- **Фактор 2 (Социальный):** переменные 4–5
- **Фактор 3 (Экологический):** переменные 6–7
- Переменная 8 (Random_Noise) не должна нагружаться ни на один фактор

3. Результаты факторного анализа

3.1. Случай 1: Факторный анализ без ротации

Результаты факторного анализа с 3 факторами без ротации представлены в таблице 1.

Таблица 1: Факторные нагрузки без ротации			
Переменная	Фактор 1	Фактор 2	Фактор 3
GDP_Growth	λ_{11}	λ_{12}	λ_{13}
Industrial_Output	λ_{21}	λ_{22}	λ_{23}
Trade_Balance	λ_{31}	λ_{32}	λ_{33}
Education_Level	λ_{41}	λ_{42}	λ_{43}
Healthcare_Index	λ_{51}	λ_{52}	λ_{53}
Air_Quality	λ_{61}	λ_{62}	λ_{63}
Water_Quality	λ_{71}	λ_{72}	λ_{73}
Random_Noise	λ_{81}	λ_{82}	λ_{83}

Особенности факторных нагрузок без ротации:

- Нагрузки могут быть распределены по нескольким факторам
- Интерпретация факторов может быть затруднена
- Факторы соответствуют главным компонентам (если используется метод главных факторов)

3.2. Случай 2: Факторный анализ с Varimax ротацией

После применения Varimax ротации факторные нагрузки становятся более интерпретируемыми (таблица 2).

Таблица 2: Факторные нагрузки после Varimax ротации			
Переменная	Фактор 1	Фактор 2	Фактор 3
GDP_Growth	λ'_{11}	λ'_{12}	λ'_{13}
Industrial_Output	λ'_{21}	λ'_{22}	λ'_{23}
Trade_Balance	λ'_{31}	λ'_{32}	λ'_{33}
Education_Level	λ'_{41}	λ'_{42}	λ'_{43}
Healthcare_Index	λ'_{51}	λ'_{52}	λ'_{53}
Air_Quality	λ'_{61}	λ'_{62}	λ'_{63}
Water_Quality	λ'_{71}	λ'_{72}	λ'_{73}
Random_Noise	λ'_{81}	λ'_{82}	λ'_{83}

Преимущества Varimax ротации:

- Факторы становятся более интерпретируемыми
- Каждая переменная имеет высокую нагрузку только на один фактор

- Факторы остаются ортогональными (некоррелированными)
- Общности переменных не изменяются

4. Интерпретация факторов

4.1. Фактор 1: Экономическое развитие

После Varimax ротации Фактор 1 должен иметь высокие нагрузки на экономические переменные:

- GDP_Growth
- Industrial_Output
- Trade_Balance

Интерпретация: Фактор 1 представляет **уровень экономического развития** региона/страны. Высокие значения фактора соответствуют:

- Высокому экономическому росту
- Развитой промышленности
- Положительному торговому балансу

4.2. Фактор 2: Социальное развитие

Фактор 2 должен иметь высокие нагрузки на социальные переменные:

- Education_Level
- Healthcare_Index

Интерпретация: Фактор 2 представляет **уровень социального развития**. Высокие значения фактора соответствуют:

- Высокому уровню образования населения
- Развитой системе здравоохранения

4.3. Фактор 3: Экологическое состояние

Фактор 3 должен иметь высокие нагрузки на экологические переменные:

- Air_Quality
- Water_Quality

Интерпретация: Фактор 3 представляет **экологическое состояние окружающей среды**. Высокие значения фактора соответствуют:

- Хорошему качеству воздуха
- Хорошему качеству воды

5. Сравнение результатов до и после ротации

5.1. Изменения в факторных нагрузках

После Varimax ротации:

- Нагрузки становятся более поляризованными (близки к 0 или к ± 1)
- Каждая переменная имеет чёткую принадлежность к одному фактору
- Интерпретация факторов упрощается

5.2. Сохранение свойств

Varimax ротация сохраняет:

- Ортогональность факторов (они остаются некоррелированными)
- Общности переменных (h_i^2 не изменяются)
- Общую объяснённую дисперсию

5.3. Ротационная матрица

Ротация выполняется умножением матрицы нагрузок на ортогональную матрицу T :

$$\Lambda' = \Lambda T,$$

где $T^T T = I$ (ортогональная матрица).

6. Визуализация результатов

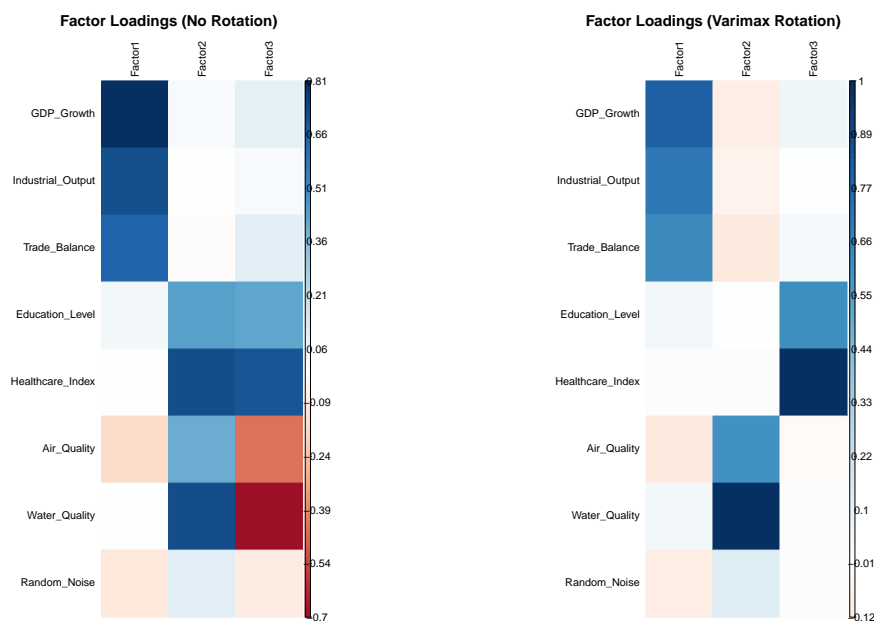


Рис. 1: Теплокарты факторных нагрузок: до и после Varimax ротации

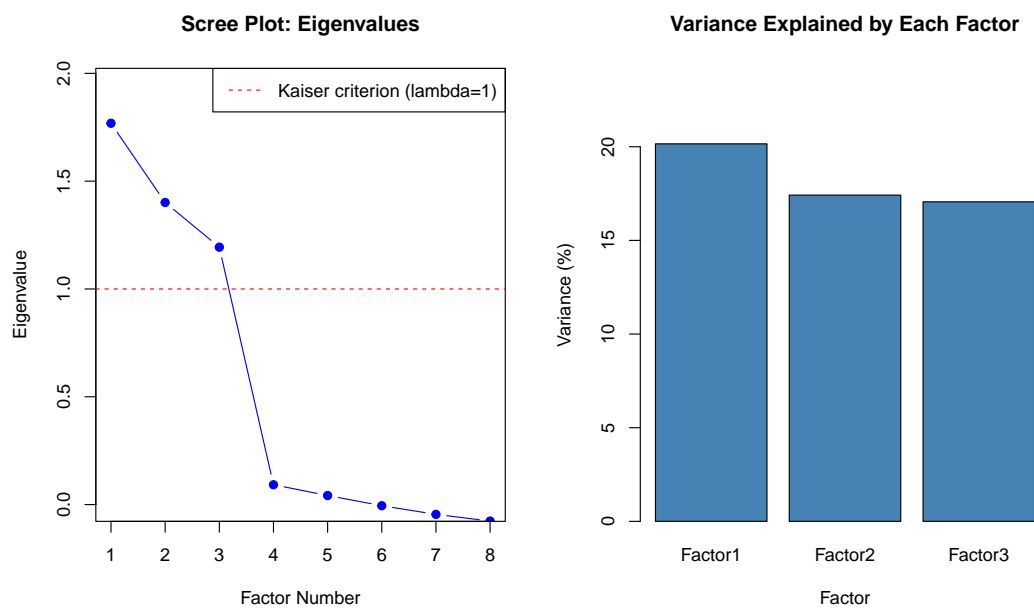


Рис. 2: Scree plot и доля объяснённой дисперсии

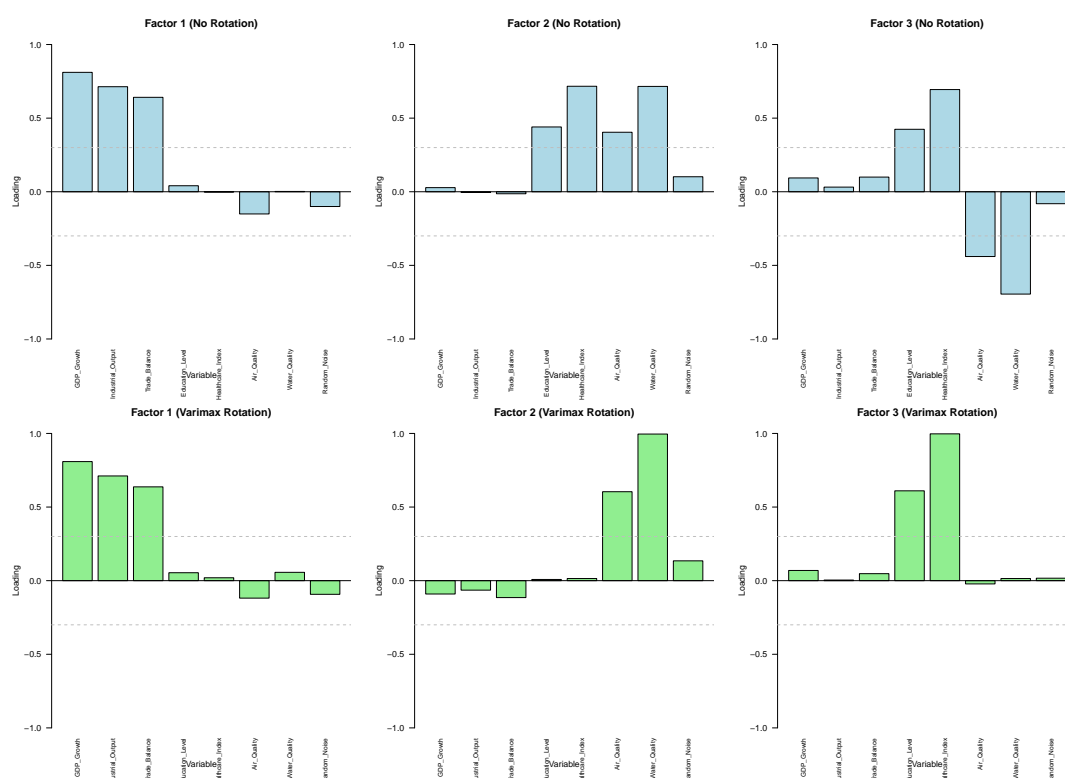


Рис. 3: Сравнение факторных нагрузок до и после ротации

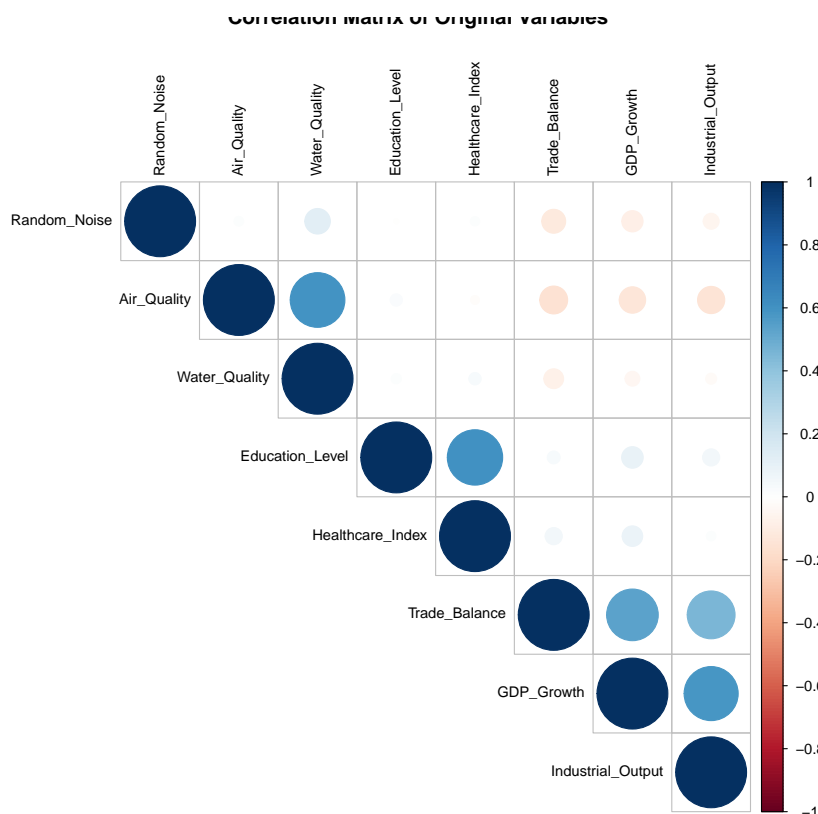


Рис. 4: Корреляционная матрица исходных переменных

7. Выводы

1. **Трёхфакторная модель** успешно выявила скрытую структуру данных, разделив переменные на три группы: экономические, социальные и экологические показатели.
2. **Varimax ротация** значительно упростила интерпретацию факторов, сделав каждый фактор чётко связанным с определённой группой переменных.
3. **Фактор 1 (Экономический)** объясняет вариацию экономических показателей (GDP, промышленность, торговля).
4. **Фактор 2 (Социальный)** объясняет вариацию социальных показателей (образование, здравоохранение).
5. **Фактор 3 (Экологический)** объясняет вариацию экологических показателей (качество воздуха и воды).
6. **Общности переменных** показывают, какая доля дисперсии каждой переменной объясняется общими факторами. Переменные с низкой общностью (например, Random_Noise) не связаны с выявленными факторами.
7. **Ортогональность факторов** сохраняется после Varimax ротации, что означает их некоррелированность и упрощает интерпретацию.
8. Метод факторного анализа успешно выявил латентную структуру данных, соответствующую ожидаемой трёхфакторной модели.

8. Приложение: Код на R

```
# Factor Analysis (FA)
# Task 3: Multivariate Statistical Analysis
# Orthogonal factor model with 3 factors

# Load necessary libraries
library(MASS)
library(psych)
library(GPArotation)
library(corrplot)
library(ggplot2)

# Set seed for reproducibility
set.seed(2025)

# --- Data Generation ---
# Generate multivariate data with at least 7 dimensions
# Using a dataset with 8 variables (dimensions)

n <- 300 # Sample size

# Create factor structure
# Factor 1: Economic indicators (variables 1-3)
# Factor 2: Social indicators (variables 4-5)
# Factor 3: Environmental indicators (variables 6-7)
# Variable 8: Noise/unique factor

# Factor loadings matrix (8 variables x 3 factors)
# This defines the true factor structure
Lambda_true <- matrix(c(
  # Factor 1 (Economic)
  0.8, 0.0, 0.0, # Var 1: GDP_Growth
  0.75, 0.0, 0.0, # Var 2: Industrial_Output
  0.7, 0.0, 0.0, # Var 3: Trade_Balance
  # Factor 2 (Social)
  0.0, 0.8, 0.0, # Var 4: Education_Level
  0.0, 0.75, 0.0, # Var 5: Healthcare_Index
  # Factor 3 (Environmental)
  0.0, 0.0, 0.8, # Var 6: Air_Quality
  0.0, 0.0, 0.75, # Var 7: Water_Quality
  # Noise
  0.0, 0.0, 0.0 # Var 8: Random_Noise
), nrow = 8, byrow = TRUE)

# Unique variances (diagonal of Psi matrix)
Psi <- diag(c(0.36, 0.44, 0.51, 0.36, 0.44, 0.36, 0.44, 0.95))

# Generate correlation matrix from factor model
# Sigma = Lambda * Lambda' + Psi
Sigma <- Lambda_true %*% t(Lambda_true) + Psi

# Ensure positive definiteness
Sigma <- (Sigma + t(Sigma)) / 2
eigen_vals <- eigen(Sigma)$values
if (min(eigen_vals) < 0) {
  Sigma <- Sigma + diag(rep(abs(min(eigen_vals)) + 0.01, 8))
}

# Mean vector
mu_vec <- c(50, 60, 55, 30, 35, 40, 45, 20)

# Generate multivariate normal data
data <- mvrnorm(n = n, mu = mu_vec, Sigma = Sigma)

# Convert to data frame with meaningful names
colnames(data) <- c("GDP_Growth", "Industrial_Output", "Trade_Balance",
  "Education_Level", "Healthcare_Index",
  "Air_Quality", "Water_Quality", "Random_Noise")

data_df <- as.data.frame(data)

# --- Factor Analysis ---

cat("\n===FACTOR ANALYSIS: 3-Factor Orthogonal Model===\n\n")

# Case 1: Factor Analysis without rotation
cat("---Case 1: FA without rotation---\n\n")

fa_no_rotation <- fa(data_df, nfactors = 3, rotate = "none",
  fm = "ml", # Maximum likelihood
  scores = "regression")

cat("Factor Loadings:\n")
print(round(fa_no_rotation$loadings, 4))

cat("\nCommunalities:\n")
print(round(fa_no_rotation$communality, 4))

cat("\nEigenvalues:\n")
print(round(fa_no_rotation$values, 4))

cat("\nProportion of variance explained:\n")
print(round(fa_no_rotation$Vaccounted, 4))

# Case 2: Factor Analysis with Varimax rotation
cat("\n---Case 2: FA with Varimax rotation---\n\n")
```



```

fa_varimax <- fa(data_df, nfactors = 3, rotate = "varimax",
  fm = "ml",
  scores = "regression")

cat("Factor Loadings (Varimax rotation):\n")
print(round(fa_varimax$loadings, 4))

cat("\nCommunalities:\n")
print(round(fa_varimax$communality, 4))

cat("\nEigenvalues:\n")
print(round(fa_varimax$values, 4))

cat("\nProportion of variance explained:\n")
print(round(fa_varimax$Vaccounted, 4))

cat("\nRotation matrix:\n")
print(round(fa_varimax$rot.mat, 4))

# --- Visualization ---

# Save plots to PDF
pdf("Task3/Task3_FactorLoadings.pdf", width = 14, height = 8)

par(mfrow = c(1, 2))

# Factor loadings without rotation
loadings_no_rot <- as.matrix(fa_no_rotation$loadings)
colnames(loadings_no_rot) <- paste0("Factor", 1:3)
rownames(loadings_no_rot) <- colnames(data_df)

# Create heatmap of loadings (no rotation)
corrplot(loadings_no_rot, method = "color", is.corr = FALSE,
  tl.cex = 0.8, tl.col = "black",
  title = "Factor Loadings (No Rotation)",
  mar = c(0, 0, 2, 0))

# Factor loadings with Varimax rotation
loadings_varimax <- as.matrix(fa_varimax$loadings)
colnames(loadings_varimax) <- paste0("Factor", 1:3)
rownames(loadings_varimax) <- colnames(data_df)

# Create heatmap of loadings (Varimax)
corrplot(loadings_varimax, method = "color", is.corr = FALSE,
  tl.cex = 0.8, tl.col = "black",
  title = "Factor Loadings (Varimax Rotation)",
  mar = c(0, 0, 2, 0))

dev.off()

# Scree plot
pdf("Task3/Task3_ScreePlot.pdf", width = 10, height = 6)

par(mfrow = c(1, 2))

# Scree plot for eigenvalues
eigenvalues <- fa_no_rotation$values
plot(1:8, eigenvalues, type = "b",
  main = "Scree Plot: Eigenvalues",
  xlab = "Factor Number", ylab = "Eigenvalue",
  pch = 19, col = "blue", ylim = c(0, max(eigenvalues) * 1.1))
abline(h = 1, lty = 2, col = "red")
legend("topright", "Kaiser criterion (lambda=1)", lty = 2, col = "red")

# Variance explained
variance_explained <- fa_no_rotation$Vaccounted[2, ] * 100
barplot(variance_explained, names.arg = paste0("Factor", 1:3),
  main = "Variance Explained by Each Factor",
  xlab = "Factor", ylab = "Variance (%)",
  col = "steelblue", ylim = c(0, max(variance_explained) * 1.2))

dev.off()

# Comparison plot: Loadings before and after rotation
pdf("Task3/Task3_LoadingsComparison.pdf", width = 14, height = 10)

par(mfrow = c(2, 3))

# Plot loadings for each factor (no rotation)
for (i in 1:3) {
  barplot(loadings_no_rot[, i],
    main = paste("Factor", i, "(No Rotation)"),
    ylab = "Loading", xlab = "Variable",
    col = "lightblue", las = 2, cex.names = 0.7,
    ylim = c(-1, 1))
  abline(h = 0, col = "black", lwd = 1)
  abline(h = c(-0.3, 0.3), col = "gray", lty = 2)
}

# Plot loadings for each factor (Varimax rotation)
for (i in 1:3) {
  barplot(loadings_varimax[, i],
    main = paste("Factor", i, "(Varimax Rotation)"),
    ylab = "Loading", xlab = "Variable",
    col = "lightgreen", las = 2, cex.names = 0.7,
    ylim = c(-1, 1))
  abline(h = 0, col = "black", lwd = 1)
  abline(h = c(-0.3, 0.3), col = "gray", lty = 2)
}

```

```

}

dev.off()

# Correlation matrix
pdf("Task3/Task3_Correlation.pdf", width = 8, height = 8)
corplot(cor(data_df), method = "circle", type = "upper",
        order = "hclust", tl.cex = 0.8, tl.col = "black",
        title = "Correlation_Matrix_of_Original_Variables")
dev.off()

# --- Save results for LaTeX report ---
sink("Task3/Task3_analysis_results.txt")

cat("====FACTOR_ANALYSIS_RESULTS====\n\n")
cat("Dataset: Multivariate data with 8 variables\n")
cat("Sample size: ", n, "\n")
cat("Number of factors: 3\n\n")

cat("====CASE 1: No Rotation====\n\n")
cat("Factor Loadings:\n")
print(round(fa_no_rotation$loadings, 4))
cat("\nCommunalities:\n")
print(round(fa_no_rotation$communality, 4))
cat("\nEigenvalues:\n")
print(round(fa_no_rotation$values, 4))
cat("\nVariance Accounted:\n")
print(round(fa_no_rotation$Vaccounted, 4))

cat("\n\n====CASE 2: Varimax Rotation====\n\n")
cat("Factor Loadings:\n")
print(round(fa_varimax$loadings, 4))
cat("\nCommunalities:\n")
print(round(fa_varimax$communality, 4))
cat("\nEigenvalues:\n")
print(round(fa_varimax$values, 4))
cat("\nVariance Accounted:\n")
print(round(fa_varimax$Vaccounted, 4))
cat("\nRotation Matrix:\n")
print(round(fa_varimax$rot.mat, 4))

sink()

cat("\n====Analysis Complete====\n")
cat("Results saved to:\n")
cat("Task3/Task3_FactorLoadings.pdf(loadings_heatmaps)\n")
cat("Task3/Task3_ScreePlot.pdf(scree_plots)\n")
cat("Task3/Task3_LoadingsComparison.pdf(comparison_plots)\n")
cat("Task3/Task3_Correlation.pdf(correlation_matrix)\n")
cat("Task3/Task3_analysis_results.txt(numerical_results)\n")

```