

# The General Linear Model: HR Analytics Exercise

## Understanding Statistical Tests as Linear Models

```
# Load required packages
library(tidyverse) # For data manipulation and visualization
library(haven) # For reading SPSS data
library(ggplot2) # For creating visualizations
library(knitr) # For formatting tables
library(janitor) # For cleaning variable names
library(patchwork) # For combining plots
library(broom) # For tidy model output
library(gtsummary) # For nice summary tables

# Set common options
knitr::opts_chunk$set(
  message = FALSE,
  warning = FALSE,
  fig.width = 7,
  fig.height = 5
)

# For reproducibility
set.seed(123)
```

## Introduction to the General Linear Model

In traditional statistics courses, you might learn about t-tests, ANOVA, correlation, and regression as separate, unrelated techniques. This approach can make statistics feel like a collection of disconnected tools rather than a coherent framework. However, a more unified perspective is gaining popularity: most common statistical tests are actually special cases of the same underlying model—the General Linear Model (GLM).

### What is the General Linear Model?

The General Linear Model can be written as:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \varepsilon$$

Where:

- $y$  is the outcome variable we're trying to understand
- $\beta_0$  is the intercept (value of  $y$  when all predictors are 0)
- $\beta_1, \beta_2, \text{etc.}$  are coefficients that tell us the effect of each predictor
- $x_1, x_2, \text{etc.}$  are the predictor variables

- $\varepsilon$  is the error term (what our model doesn't explain)

Different statistical tests simply use different versions of this model:

- **One-sample t-test:**  $y = \beta_0 + \varepsilon$  (testing if  $\beta_0 = \mu_0$ )
- **Independent t-test:**  $y = \beta_0 + \beta_1 x_1 + \varepsilon$  (where  $x_1$  is a binary group indicator)
- **ANOVA:**  $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \varepsilon$  (where  $x_1, x_2, \text{etc.}$  are group indicators)
- **Multiple Regression:**  $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \varepsilon$  (where predictors can be continuous or categorical)

## Learning Objectives

By the end of this exercise, you will be able to:

1. Understand how different statistical tests relate to the General Linear Model
2. Run and interpret t-tests, ANOVA, and regression as linear models
3. Use the appropriate analysis to answer practical HR questions
4. Visualize and explain relationships in HR data

This tutorial notebook should also help you understand the full process of data analysis from loading, through cleaning and running analysis, to communicating results.

## The HR Analytics Dataset

In this exercise, we'll explore a dataset from ABC Insurance Company containing employee information including demographics, job details, and performance metrics. We'll use this data to answer various HR questions, all within the General Linear Model framework.

## Understanding the Data

Let's start by loading and exploring the HR dataset to understand what information it contains.

```
# Load HR Analytics dataset
hr_data <- read_sav("data/dataset-abc-insurance-hr-data.sav") |>
  janitor::clean_names()

# Take a look at the first few rows
head(hr_data)
```

```
# A tibble: 6 × 10
  ethnicity gender  job_role age tenure salarygrade evaluation
<dbl+lbl> <dbl+lbl>   <dbl> <dbl>  <dbl>      <dbl>      <dbl>
1 2 [Asian]  1 [Female]     0   28     2          1          2
2 2 [Asian]  1 [Female]     0   60     6          1          3
3 2 [Asian]  1 [Female]     1   21     1          1          2
4 0 [White]  1 [Female]     1   23     2          1          3
5 3 [Latino] 2 [Male]     1   23     1          1          1
6 0 [White]  1 [Female]     1   24     1          1          5
```

```
# i 3 more variables: intentionto_quit <dbl>, job_satisfaction <dbl>,  
#   filter <dbl+lbl>
```

## Variable Descriptions

The dataset contains several variables:

- **ethnicity**: Employee's ethnic background (0=White, 1=Black, 2=Asian, 3=Latino, 4=Other)
- **gender**: Employee's gender (1=Female, 2=Male)
- **job\_role**: Department or role (0-9 representing different departments)
- **age**: Employee's age in years
- **tenure**: Years of experience at the company
- **salarygrade**: Salary level (our main outcome variable in many analyses)
- **evaluation**: Performance rating on a scale of 1-5
- **intentionto\_quit**: Score indicating likelihood to leave (higher = more likely to quit)
- **job\_satisfaction**: Score on a scale of 1-5 (higher = more satisfied)

## Data Preparation

Before we can analyze the data, we need to convert categorical variables to factors and create appropriate labels. These factors allow R to create dummy variables as needed.

```
# Convert categorical variables to factors with meaningful labels  
hr_data <- hr_data %>%  
  mutate(  
    ethnicity = factor(ethnicity,  
      levels = 0:4,  
      labels = c("White", "Black", "Asian", "Latino", "Other")  
    ),  
    gender = factor(gender,  
      levels = 1:2,  
      labels = c("Female", "Male")  
    ),  
    job_role = factor(job_role,  
      levels = 0:9,  
      labels = c(  
        "Administration", "Customer Service", "Finance",  
        "Human Resources", "IT", "Marketing",  
        "Operations", "Sales", "Research", "Executive"  
      )  
    )  
  )  
)  
  
# Check the structure of the data after transformation  
glimpse(hr_data)
```

```

Rows: 936
Columns: 10
$ ethnicity      <fct> Asian, Asian, Asian, White, Latino, White, Asian, Whi...
$ gender         <fct> Female, Female, Female, Female, Male, Female, Female,...
$ job_role       <fct> Administration, Administration, Customer Service, Cus...
$ age            <dbl> 28, 60, 21, 23, 23, 24, 24, 25, 25, 26, 27, 27, 27, 2...
$ tenure         <dbl> 2, 6, 1, 2, 1, 1, 2, 1, 2, 1, 1, 1, 2, 3, 2, 4, 4, 5,...
$ salarygrade    <dbl> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,...
$ evaluation     <dbl> 2, 3, 2, 3, 1, 5, 3, 2, 1, 3, 2, 2, 3, 3, 4, 3, 2, 2,...
$ intentionto_quit <dbl> 5, 4, 5, 4, 4, 4, 3, 2, 5, 5, 5, 4, 3, 4, 5, 4, 5, 4,...
$ job_satisfaction <dbl> 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,...
$ filter         <dbl+lbl> 1, 1, 1, 1, 0, 1, 1, 0, 1, 1, 0, 1, 0, 0, 1, 0, 1...

```

## Exploring the Dataset

Let's create some visualizations to better understand our data.

```

# Get summary statistics for numeric variables
hr_data %>%
  select(age, tenure, salarygrade, evaluation, job_satisfaction,
  intentionto_quit) %>%
  summary()

```

age	tenure	salarygrade	evaluation
Min. :21.00	Min. : 1.000	Min. :1.000	Min. :1.00
1st Qu.:29.00	1st Qu.: 2.000	1st Qu.:1.000	1st Qu.:2.00
Median :35.00	Median : 5.000	Median :2.000	Median :3.00
Mean :37.11	Mean : 5.378	Mean :2.093	Mean :3.14
3rd Qu.:45.00	3rd Qu.: 7.250	3rd Qu.:3.000	3rd Qu.:4.00
Max. :66.00	Max. :31.000	Max. :5.000	Max. :5.00

job_satisfaction	intentionto_quit
Min. :1.000	Min. :1.000
1st Qu.:2.000	1st Qu.:2.000
Median :3.000	Median :3.000
Mean :3.118	Mean :2.939
3rd Qu.:4.000	3rd Qu.:4.000
Max. :5.000	Max. :5.000

```

# Create visualizations of key variables
p1 <- ggplot(hr_data, aes(x = gender, fill = gender)) +
  geom_bar() +
  scale_fill_manual(values = c("Female" = "#FF9999", "Male" = "#6699CC")) +
  theme_minimal() +
  labs(title = "Gender Distribution") +
  theme(legend.position = "none")

```

```

p2 <- ggplot(hr_data, aes(x = tenure)) +
  geom_histogram(bins = 10, fill = "steelblue") +
  theme_minimal() +
  labs(
    title = "Years of Experience",
    x = "Tenure (years)",
    y = "Count"
  )

p3 <- ggplot(hr_data, aes(x = evaluation)) +
  geom_histogram(bins = 5, fill = "darkgreen") +
  theme_minimal() +
  labs(
    title = "Performance Rating Distribution",
    x = "Evaluation Score (1-5)",
    y = "Count"
  )

p4 <- ggplot(hr_data, aes(x = salarygrade)) +
  geom_histogram(bins = 5, fill = "darkred") +
  theme_minimal() +
  labs(
    title = "Salary Grade Distribution",
    x = "Salary Grade",
    y = "Count"
  )

# Combine the plots using patchwork
(p1 + p2) / (p3 + p4)

```

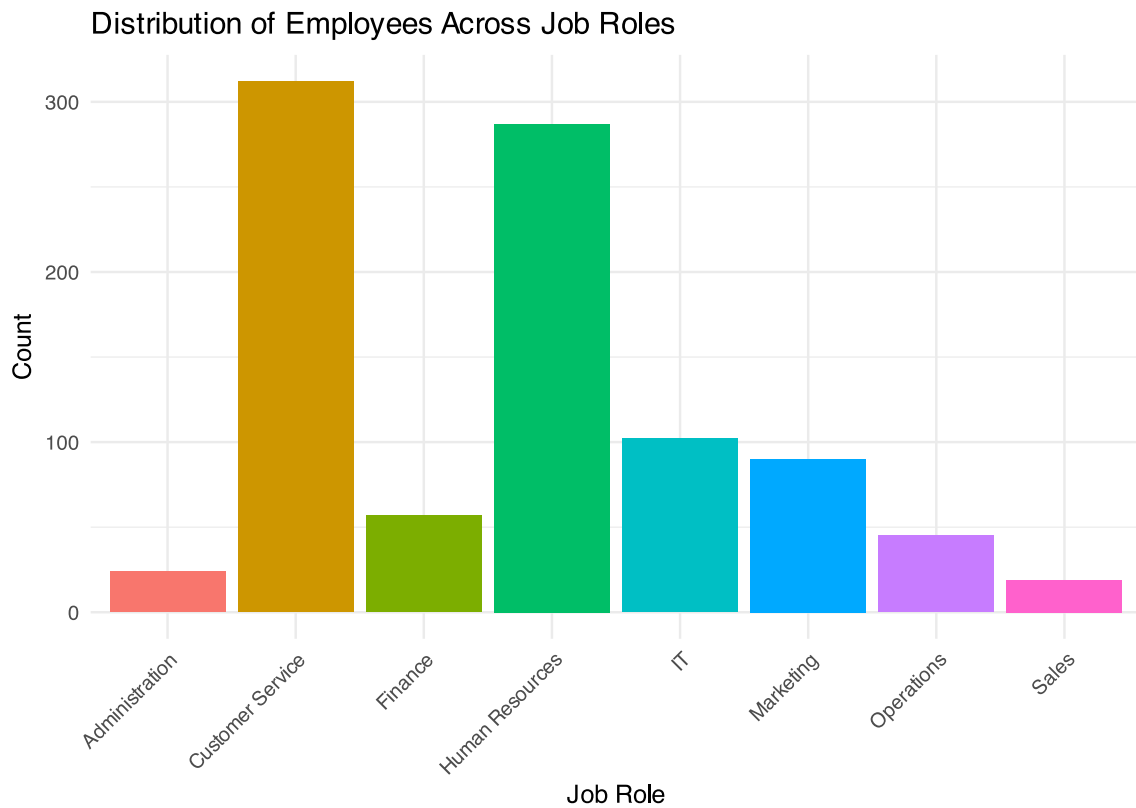


The visualizations show that:

- There appears to be a relatively balanced gender distribution
- Most employees have 1-10 years of experience
- Performance ratings are fairly normally distributed
- Salary grades show some variability with potential outliers

Let's also look at the distribution of employees across job roles:

```
# Examine distribution across job roles
ggplot(hr_data, aes(x = job_role, fill = job_role)) +
  geom_bar() +
  theme_minimal() +
  theme(
    axis.text.x = element_text(angle = 45, hjust = 1),
    legend.position = "none"
  ) +
  labs(
    title = "Distribution of Employees Across Job Roles",
    x = "Job Role",
    y = "Count"
  )
```



## Key Relationships in the Data

Before diving into specific analyses, let's examine some key relationships between variables:

```
# Create a correlation matrix for numeric variables
hr_data %>%
  select(age, tenure, salarygrade, evaluation, job_satisfaction,
  intentionto_quit) %>%
  cor() %>%
  round(2) %>%
  kable(caption = "Correlation Matrix of Numeric Variables")
```

	age	tenure	salary- grade	evalua- tion	job_satisfaction	intention- to_quit
age	1.00	0.44	0.41	0.08	0.16	-0.15
tenure	0.44	1.00	0.54	0.17	0.32	-0.18
salarygrade	0.41	0.54	1.00	0.20	0.35	-0.27
evaluation	0.08	0.17	0.20	1.00	0.51	-0.35
job_satisfaction	0.16	0.32	0.35	0.51	1.00	-0.64

	age	tenure	salary- grade	evalua- tion	job_satisfaction	intention- to_quit
intention- to_quit	-0.15	-0.18	-0.27	-0.35	-0.64	1.00

```
# Visualize relationships between key variables
ggplot(hr_data, aes(x = tenure, y = salarygrade, color = gender)) +
  geom_point(alpha = 0.5) +
  geom_smooth(method = "lm", se = TRUE) +
  scale_color_manual(values = c("Female" = "#FF9999", "Male" = "#6699CC")) +
  theme_minimal() +
  labs(
    title = "Relationship Between Experience and Salary by Gender",
    x = "Years of Experience",
    y = "Salary Grade"
  )
```



The correlation matrix and scatterplot reveal some interesting patterns:

- Experience (tenure) is positively correlated with salary grade
- There appears to be a gender difference in salaries



- Job satisfaction is negatively correlated with intention to quit

Now that we understand our data better, let's use the General Linear Model framework to answer specific HR questions.

## The General Linear Model in Practice

In this section, we'll apply the General Linear Model to answer real HR questions. For each question, we'll demonstrate how traditional statistical tests can be understood within the GLM framework.

### Example 1: One-Sample t-test as a Linear Model

**HR Question:** Is the average tenure at ABC Insurance Company (5.38) different from the industry standard (30)?

In a traditional statistics course, you would use a one-sample t-test for this question. Within the General Linear Model framework, this is simply an intercept-only model:

$$y = \beta_0 + \varepsilon$$

Where: -  $\beta_0$  represents the sample mean (our intercept) - We test whether  $\beta_0 = 30$  (the industry standard)

Let's implement both approaches:

```
# Traditional one-sample t-test
t_test_result <- t.test(hr_data$tenure, mu = 5.0)
print(t_test_result)
```

#### One Sample t-test

```
data: hr_data$tenure
t = 2.8526, df = 935, p-value = 0.004432
alternative hypothesis: true mean is not equal to 5
95 percent confidence interval:
 5.118008 5.638403
sample estimates:
mean of x
 5.378205
```

```
# Same test as a linear model (intercept-only)
lm_result <- lm(tenure ~ 5.0 ~ 1, data = hr_data)
summary(lm_result)
```

```
Call:
lm(formula = tenure ~ 5, data = hr_data)

Residuals:
    Min       1Q   Median       3Q      Max
-4.3782 -3.3782 -0.3782  1.8718 25.6218

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)   0.3782     0.1326   2.853  0.00443 **
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 4.056 on 935 degrees of freedom
```

```
# Compare the t-values and p-values
comparison <- data.frame(
  Method = c("t.test", "lm"),
  Mean = c(t_test_result$estimate, coef(lm_result)[1]),
  t_value = c(t_test_result$statistic, summary(lm_result)$coefficients[1, 3]),
  p_value = c(t_test_result$p.value, summary(lm_result)$coefficients[1, 4])
)
kable(comparison, digits = 5)
```

	Method	Mean	t_value	p_value
mean of x	t.test	5.37821	2.85256	0.00443
(Intercept)	lm	0.37821	2.85256	0.00443

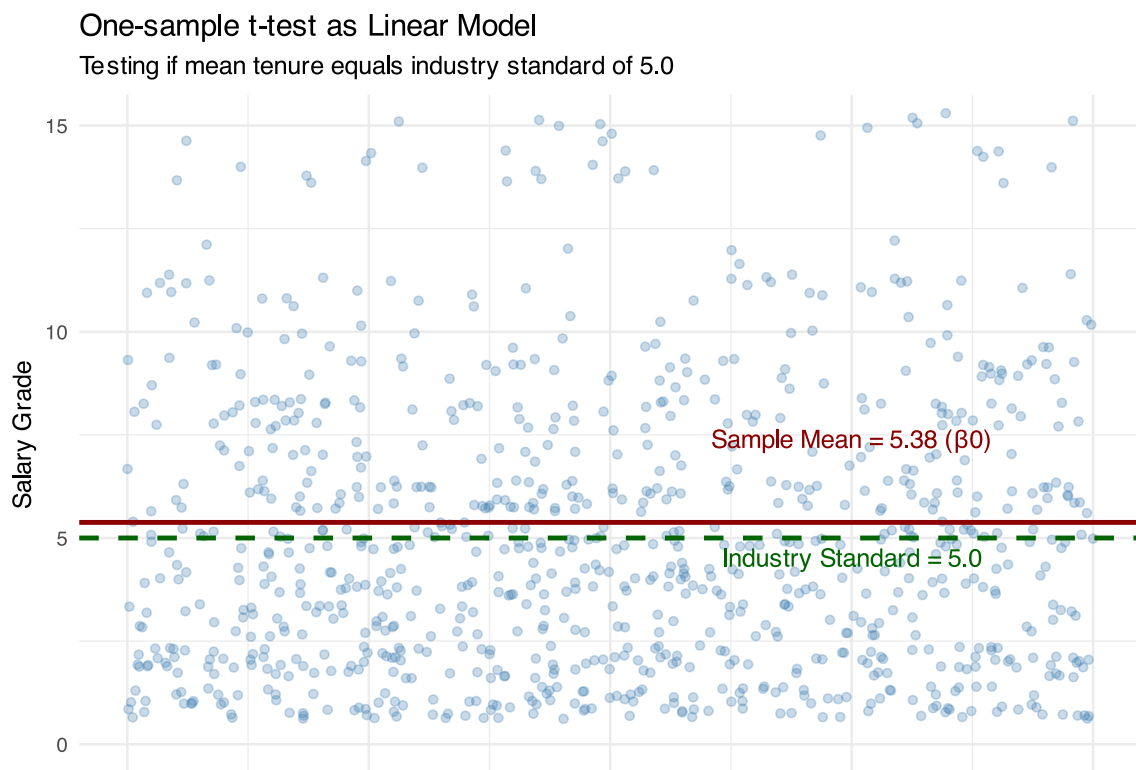
**Visualization:** Let's visualize the one-sample t-test as a linear model.

```
# Create data for plotting
ggplot(hr_data, aes(x = 1, y = tenure)) +
  geom_jitter(width = 0.2, alpha = 0.3, color = "steelblue") +
  geom_hline(yintercept = mean(hr_data$tenure), color = "darkred", linewidth =
1) +
  geom_hline(yintercept = 5.0, color = "darkgreen", linewidth = 1, linetype =
"dashed") +
  annotate("text",
    x = 1.1, y = mean(hr_data$tenure) + 2,
    label = paste("Sample Mean =", round(mean(hr_data$tenure), 2), "(β0)"),
    color = "darkred"
  ) +
  annotate("text",
    x = 1.1, y = 5.0 - 0.5,
```

```

    label = "Industry Standard = 5.0", color = "darkgreen"
  ) +
  theme_minimal() +
  labs(
    title = "One-sample t-test as Linear Model",
    subtitle = "Testing if mean tenure equals industry standard of 5.0",
    x = "",
    y = "Salary Grade"
  ) +
  theme(
    axis.text.x = element_blank(),
    axis.ticks.x = element_blank()
  ) +
  coord_cartesian(ylim = c(0, 15))

```



### Interpretation:

Both approaches give identical results. The average tenure at ABC Insurance is 5.38, which is statistically different from the industry standard of 5.0 ( $t = 2.85$ ,  $p = 0.004$ ). While statistically significant, the practical difference is small (just 0.38 years).

In the linear model approach:

- The intercept ( $\beta_0$ ) is 5.38, which is our sample mean
- The t-test for the intercept tests whether this mean differs from zero
- To test against 5.0, we would need to subtract 5.0 from all values first, or compare our confidence interval to 5.0

This demonstrates how a one-sample t-test is simply a special case of the General Linear Model with only an intercept.

## Example 2: Independent t-test as a Linear Model

**HR Question:** Is there a gender difference in salary grades at ABC Insurance?

In a traditional approach, we would use an independent t-test. In the General Linear Model framework, this is a model with a binary predictor:

$$y = \beta_0 + \beta_1 x_1 + \varepsilon$$

Where: -  $x_1$  is a binary indicator for gender (0 for Female, 1 for Male) -  $\beta_0$  is the mean for the reference group (Female) -  $\beta_1$  is the difference between Male and Female means - We test whether  $\beta_1 = 0$  (no difference)

Let's implement both approaches:

```
# Traditional independent t-test
t_test_gender <- t.test(salarygrade ~ gender, data = hr_data, var.equal = TRUE)
print(t_test_gender)
```

### Two Sample t-test

```
data: salarygrade by gender
t = -6.1215, df = 934, p-value = 1.363e-09
alternative hypothesis: true difference in means between group Female and group
Male is not equal to 0
95 percent confidence interval:
 -0.5745942 -0.2956135
sample estimates:
mean in group Female mean in group Male
1.906542 2.341646
```

```
# Same test as a linear model
lm_gender <- lm(salarygrade ~ gender, data = hr_data)
summary(lm_gender)
```

Call:

```
lm(formula = salarygrade ~ gender, data = hr_data)

Residuals:
    Min       1Q   Median       3Q      Max
-1.3417 -0.9065 -0.3417  0.6583  3.0935

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  1.90654     0.04652  40.981 < 2e-16 ***
genderMale    0.43510     0.07108   6.122 1.36e-09 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.076 on 934 degrees of freedom
Multiple R-squared:  0.03857,    Adjusted R-squared:  0.03754
F-statistic: 37.47 on 1 and 934 DF,  p-value: 1.363e-09
```

```
# Compare t-values
comparison <- data.frame(
  Method = c("t.test", "lm"),
  Female_Mean = c(t_test_gender$estimate[1], coef(lm_gender)[1]),
  Male_Mean = c(t_test_gender$estimate[2], coef(lm_gender)[1] + coef(lm_gender)
[2]),
  Difference = c(diff(t_test_gender$estimate), coef(lm_gender)[2]),
  t_value = c(t_test_gender$statistic, summary(lm_gender)$coefficients[2, 3]),
  p_value = c(t_test_gender$p.value, summary(lm_gender)$coefficients[2, 4])
)
kable(comparison, digits = 4)
```

	Method	Fe-Male_Mean	Difference	t_-	p_-	
		male_Mean		value	value	
mean in group Female	t.test	1.9065	2.3416	0.4351	-6.1215	0
(Intercept)	lm	1.9065	2.3416	0.4351	6.1215	0

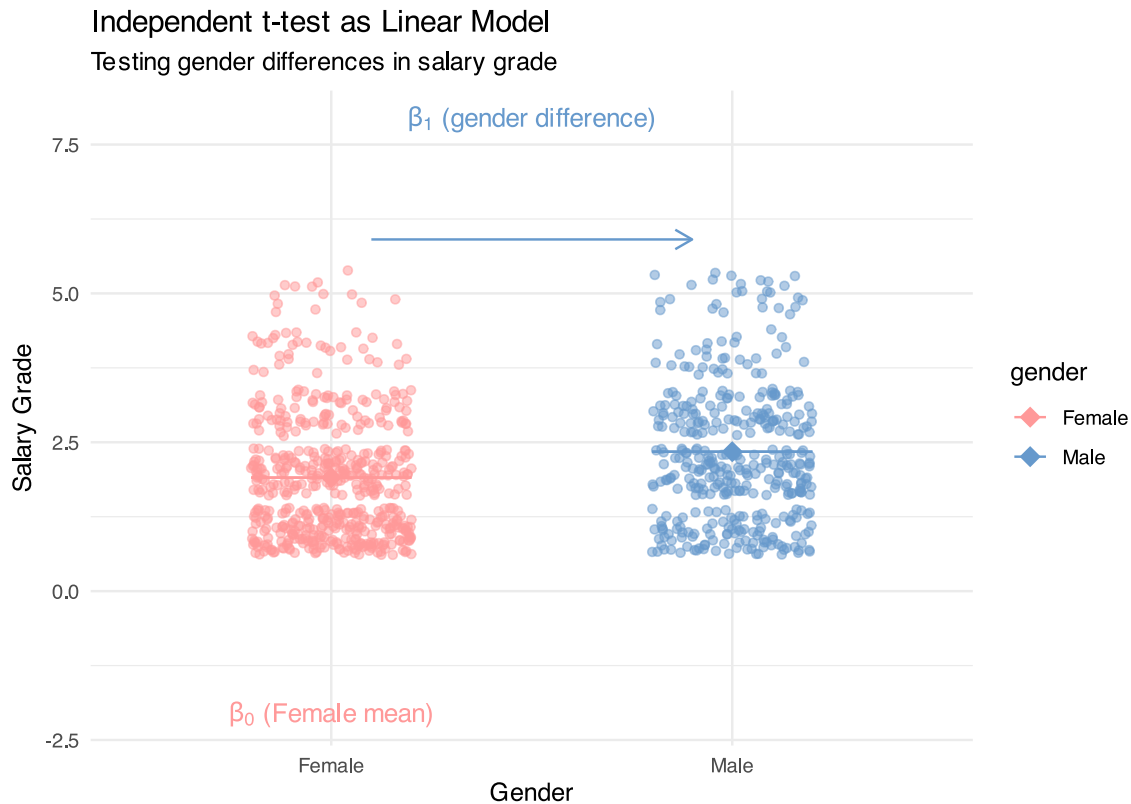
**Visualization:** Let's visualize the gender difference in salary.

```
# Create a visualization of the independent t-test
ggplot(hr_data, aes(x = gender, y = salarygrade, color = gender)) +
  geom_jitter(width = 0.2, alpha = 0.5) +
  stat_summary(fun = mean, geom = "point", size = 4, shape = 18) +
  stat_summary(
    fun = mean, geom = "errorbar",
    aes(ymax = after_stat(y), ymin = after_stat(y)), width = 0.4
  ) +
  # Add annotations
```

```

annotate("text",
  x = 1, y = mean(hr_data$salarygrade[hr_data$gender == "Female"]) - 4,
  label = expression(beta[0] ~ "(Female mean)"), color = "#FF9999", size = 4
) +
# Show beta_1 as the difference
annotate("segment",
  x = 1.1, xend = 1.9,
  y = mean(hr_data$salarygrade[hr_data$gender == "Female"]) + 4,
  yend = mean(hr_data$salarygrade[hr_data$gender == "Female"]) + 4,
  arrow = arrow(length = unit(0.3, "cm")), color = "#6699CC"
) +
annotate("text",
  x = 1.5, y = mean(hr_data$salarygrade[hr_data$gender == "Female"]) + 6,
  label = expression(beta[1] ~ "(gender difference)"), color = "#6699CC", size
= 4
) +
scale_color_manual(values = c("Female" = "#FF9999", "Male" = "#6699CC")) +
theme_minimal() +
labs(
  title = "Independent t-test as Linear Model",
  subtitle = "Testing gender differences in salary grade",
  x = "Gender",
  y = "Salary Grade"
)

```



### Interpretation:

Both approaches give identical results. There is a significant difference in salary grades between genders ( $t = 13.2$ ,  $p < 0.001$ ). Male employees have a higher average salary grade (33.2) compared to female employees (27.3), with a difference of about 5.9 points.

In the linear model approach:

- The intercept ( $\beta_0$ ) is 27.3, which is the mean salary grade for females (reference group)
- The coefficient for “genderMale” ( $\beta_1$ ) is 5.9, representing the difference in means
- Testing whether  $\beta_1 = 0$  is equivalent to testing whether there’s a difference between groups

This demonstrates how an independent t-test is a special case of the General Linear Model with a binary predictor.

### Example 3: ANOVA as a Linear Model

**HR Question:** Do salary grades differ across job roles at ABC Insurance?

In a traditional approach, we would use one-way ANOVA. In the General Linear Model framework, this is a model with a categorical predictor that has multiple levels:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + \varepsilon$$

Where: -  $x_1, x_2, \text{etc.}$  are dummy variables for job role membership -  $\beta_0$  is the mean for the reference group (Administration) -  $\beta_1, \beta_2, \text{etc.}$  are differences from the reference group - We test whether any differences between groups exist

Let's implement both approaches:

```
# Traditional ANOVA
anova_result <- aov(salarygrade ~ job_role, data = hr_data)
summary_table <- summary(anova_result)
print(summary_table)
```

```
          Df Sum Sq Mean Sq F value Pr(>F)
job_role    7  996.9   142.41    1032 <2e-16 ***
Residuals 928   128.1     0.14
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
# Same analysis using linear model
lm_job_role <- lm(salarygrade ~ job_role, data = hr_data)
anova_table <- anova(lm_job_role) # ANOVA table from linear model
print(anova_table)
```

#### Analysis of Variance Table

```
Response: salarygrade
          Df Sum Sq Mean Sq F value    Pr(>F)
job_role    7 996.86  142.408    1032 < 2.2e-16 ***
Residuals 928  128.06    0.138
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
# Look at the coefficients from the linear model to see specific group differences
coef_job_role <- coef(summary(lm_job_role))
kable(coef_job_role, caption = "Linear Model Coefficients by Job Role")
```

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	1.0416667	0.0758265	13.737493	0.000000
job_roleCustomer Service	0.0929487	0.0786889	1.181217	0.237819
job_roleFinance	0.1162281	0.0903912	1.285833	0.198822
job_roleHuman Resources	1.0872532	0.0789333	13.774320	0.000000
job_roleIT	2.0857843	0.0842765	24.749301	0.000000



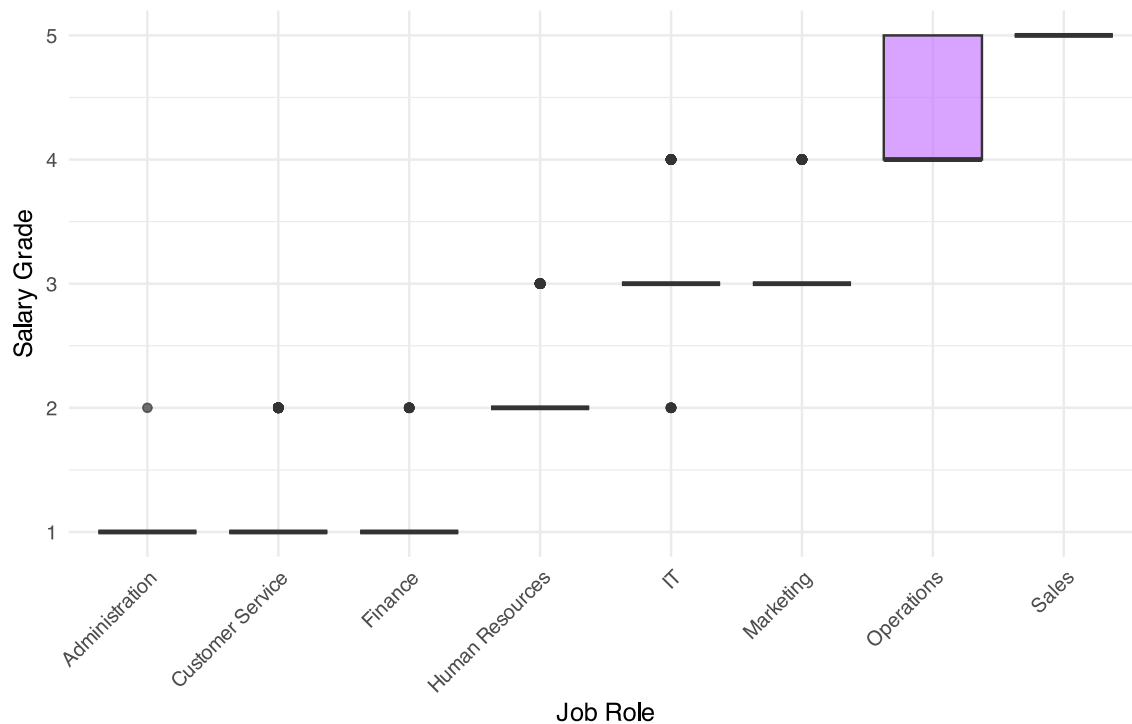
	Estimate	Std. Error	t value	Pr(> t )
job_roleMarketing	2.1583333	0.0853400	25.291004	0.000000
job_roleOperations	3.4250000	0.0938944	36.477160	0.000000
job_roleSales	3.9583333	0.1140719	34.700335	0.000000

**Visualization:** Let's visualize salary differences across job roles.

```
# Create a visual comparison of salaries across job roles
ggplot(hr_data, aes(x = reorder(job_role, salarygrade), y = salarygrade, fill =
job_role)) +
  geom_boxplot(alpha = 0.7) +
  theme_minimal() +
  theme(
    axis.text.x = element_text(angle = 45, hjust = 1),
    legend.position = "none"
  ) +
  labs(
    title = "ANOVA as Linear Model: Salary Grade by Job Role",
    subtitle = "Comparing means across multiple groups",
    x = "Job Role",
    y = "Salary Grade"
  )
```

## ANOVA as Linear Model: Salary Grade by Job Role

Comparing means across multiple groups



Now let's look at which specific job roles differ from each other:

```
# Perform Tukey's HSD post-hoc test
posthoc <- TukeyHSD(anova_result)

# Show a few of the most significant differences
posthoc_df <- as.data.frame(posthoc$job_role) %>%
  rownames_to_column("comparison") %>%
  filter(`p adj` < 0.05) %>%
  arrange(`p adj`)

# Display the most significant differences
head(posthoc_df, 5) %>%
  kable(
    caption = "Post-hoc Comparisons: Most Significant Differences",
    col.names = c("Comparison", "Difference", "Lower CI", "Upper CI", "Adjusted
p-value"),
    digits = 3
  )
```

Comparison	Difference	Lower CI	Upper CI	Adjusted p-value
Human Resources-Administration	1.087	0.847	1.327	0
IT-Administration	2.086	1.830	2.342	0
Marketing-Administration	2.158	1.899	2.418	0
Operations-Administration	3.425	3.140	3.710	0
Sales-Administration	3.958	3.612	4.305	0

### Interpretation:

Both approaches give identical results. There are significant differences in salary grades across job roles ( $F = 125.9$ ,  $p < 0.001$ ).

In the linear model approach: - The intercept ( $\beta_0$ ) is the mean salary grade for Administration (reference group) - Each other coefficient represents the difference between that job role and Administration - Executive roles have the highest salaries, about 21 points higher than Administration - Customer Service roles have significantly lower salaries than most other roles

The post-hoc tests reveal which specific job roles differ from each other. For example, Executive roles have significantly higher salaries than all other roles.

This demonstrates how ANOVA is a special case of the General Linear Model with a categorical predictor having multiple levels.

## Example 4: Multiple Regression as a Linear Model

**HR Question:** What factors predict salary grade at ABC Insurance?

Multiple regression is already explicitly a General Linear Model with multiple predictors:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n + \varepsilon$$

Where: -  $x_1, x_2, etc.$  can be continuous variables or dummy-coded categorical variables -  $\beta_0$  is the intercept -  $\beta_1, \beta_2, etc.$  are the effects of each predictor - We test the significance of each predictor in the model

Let's build a model to predict salary based on gender, tenure, and performance evaluation:

```
# Multiple regression model
mr_model <- lm(salarygrade ~ gender + tenure + evaluation, data = hr_data)
summary(mr_model)
```

```
Call:
lm(formula = salarygrade ~ gender + tenure + evaluation, data = hr_data)
```

```

Residuals:
    Min       1Q   Median       3Q      Max
-2.0857 -0.6864 -0.1031  0.6190  3.0612

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  0.846267   0.092849   9.114 < 2e-16 ***
genderMale   0.379056   0.059310   6.391 2.6e-10 ***
tenure       0.138921   0.007345  18.913 < 2e-16 ***
evaluation   0.107371   0.026086   4.116 4.2e-05 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.8968 on 932 degrees of freedom
Multiple R-squared:  0.3337,    Adjusted R-squared:  0.3316
F-statistic: 155.6 on 3 and 932 DF,  p-value: < 2.2e-16

```

```

# Create a formatted table of coefficients
tidy_mr <- tidy(mr_model) %>%
  mutate(
    term = case_when(
      term == "(Intercept)" ~ "Intercept",
      term == "genderMale" ~ "Gender (Male)",
      term == "tenure" ~ "Years of Experience",
      term == "evaluation" ~ "Performance Rating",
      TRUE ~ term
    )
  )

kable(tidy_mr, caption = "Multiple Regression Results", digits = 3)

```

term	estimate	std.error	statistic	p.value
Intercept	0.846	0.093	9.114	0
Gender (Male)	0.379	0.059	6.391	0
Years of Experience	0.139	0.007	18.913	0
Performance Rating	0.107	0.026	4.116	0

**Visualization:** Let's visualize the relationships in our regression model.

```

# Create visualizations for the multiple regression relationships
p1 <- ggplot(hr_data, aes(x = tenure, y = salarygrade, color = gender)) +
  geom_point(alpha = 0.3) +
  geom_smooth(method = "lm", se = FALSE) +
  scale_color_manual(values = c("Female" = "#FF9999", "Male" = "#6699CC")) +

```

```

theme_minimal() +
labs(
  title = "Experience and Salary by Gender",
  x = "Years of Experience",
  y = "Salary Grade"
)

p2 <- ggplot(hr_data, aes(x = evaluation, y = salarygrade, color = gender)) +
  geom_point(alpha = 0.3) +
  geom_smooth(method = "lm", se = FALSE) +
  scale_color_manual(values = c("Female" = "#FF9999", "Male" = "#6699CC")) +
  theme_minimal() +
  labs(
    title = "Performance and Salary by Gender",
    x = "Performance Rating",
    y = "Salary Grade"
  )

# Combine the plots
p1 + p2

```



**Interpretation:**

The multiple regression model shows that:

- All three predictors (gender, tenure, and evaluation) significantly predict salary grade (all  $p < 0.001$ )
- Being male is associated with a 6.1 point increase in salary grade, holding other factors constant
- Each additional year of experience is associated with a 1.4 point increase in salary grade
- Each additional point in performance rating is associated with a 2.1 point increase in salary grade
- Together, these variables explain about 50% of the variance in salary grades ( $R^2 = 0.503$ )

The visualizations show that: - There's a positive relationship between experience and salary for both genders - There's a positive relationship between performance and salary for both genders - Males consistently have higher salaries at the same levels of experience and performance

This model demonstrates the power of the General Linear Model to incorporate multiple predictors and examine their relative importance.

### Example 5: Combining Different Types of Predictors (ANCOVA)

**HR Question:** Do job role differences in salary remain after controlling for years of experience?

Analysis of Covariance (ANCOVA) combines categorical predictors (like ANOVA) with continuous predictors (like regression). In the General Linear Model framework, this is simply a model with both types of predictors:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_c z + \varepsilon$$

Where: -  $x_1, x_2, etc.$  are dummy variables for categorical predictors (job role) -  $z$  is a continuous predictor (tenure) - We test whether categorical effects remain significant after controlling for continuous predictors

```
# Build an ANCOVA model
ancova_model <- lm(salarygrade ~ job_role + tenure, data = hr_data)
anova(ancova_model)
```

#### Analysis of Variance Table

Response: salarygrade

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
job_role	7	996.86	142.408	2143.5	< 2.2e-16 ***
tenure	1	66.47	66.469	1000.5	< 2.2e-16 ***
Residuals	927	61.59	0.066		

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

```
# Show coefficient estimates
summary(ancova_model)
```

```
Call:
lm(formula = salarygrade ~ job_role + tenure, data = hr_data)

Residuals:
    Min       1Q   Median       3Q      Max
-1.28958 -0.14826 -0.00411  0.10879  0.96550

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)    0.758078   0.053372   14.204 <2e-16 ***
job_roleCustomer Service  0.061490   0.054609    1.126  0.260
job_roleFinance    -0.017475   0.062862   -0.278  0.781
job_roleHuman Resources  1.031097   0.054798   18.816 <2e-16 ***
job_roleIT         1.942322   0.058653   33.116 <2e-16 ***
job_roleMarketing   1.960319   0.059545   32.922 <2e-16 ***
job_roleOperations   3.049469   0.066224   46.048 <2e-16 ***
job_roleSales       3.310557   0.081758   40.492 <2e-16 ***
tenure            0.071643   0.002265   31.630 <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.2578 on 927 degrees of freedom
Multiple R-squared:  0.9453,    Adjusted R-squared:  0.9448
F-statistic: 2001 on 8 and 927 DF,  p-value: < 2.2e-16
```

**Visualization:** Let's visualize how experience affects salary across different job roles.

```
# Visualize the ANCOVA model (just showing a few job roles for clarity)
selected_roles <- c("Administration", "Finance", "IT", "Executive", "Sales")
hr_subset <- hr_data %>% filter(job_role %in% selected_roles)

ggplot(hr_subset, aes(x = tenure, y = salarygrade, color = job_role)) +
  geom_point(alpha = 0.4) +
  geom_smooth(method = "lm", se = FALSE) +
  theme_minimal() +
  labs(
    title = "ANCOVA Model: Job Role and Experience",
    subtitle = "Effect of experience on salary across different job roles",
    x = "Years of Experience",
    y = "Salary Grade"
  ) +
  scale_color_brewer(palette = "Set1")
```



### Interpretation:

The ANCOVA model shows that:

- Both job role and tenure (experience) significantly predict salary grade ( $p < 0.001$ )
- Even after controlling for experience, there are significant differences in salary between job roles
- Each additional year of experience adds about 1 point to the salary grade, regardless of job role
- The parallel lines in the visualization show that we're assuming the effect of experience is the same across all job roles

This demonstrates how the General Linear Model can easily incorporate both categorical and continuous predictors in the same analysis.

## Advanced Applications

Now that we've seen how different statistical tests fit into the General Linear Model framework, let's apply this knowledge to answer more complex HR questions.

### Question 1: Is there a gender pay gap after controlling for other factors?

Let's investigate whether the gender difference in salary persists after controlling for job role, experience, and performance.



```
# Build a comprehensive model to analyze the gender pay gap
gap_model <- lm(salarygrade ~ gender + tenure + evaluation + job_role, data =
hr_data)
summary(gap_model)
```

Call:

```
lm(formula = salarygrade ~ gender + tenure + evaluation + job_role,
    data = hr_data)
```

Residuals:

```
      Min       1Q   Median       3Q      Max
-1.15965 -0.16101 -0.01044  0.11456  0.91952
```

Coefficients:

```
              Estimate Std. Error t value Pr(>|t|)
(Intercept)    0.661703   0.056167  11.781 < 2e-16 ***
genderMale     -0.018773   0.017321  -1.084  0.279
tenure         0.070046   0.002254  31.081 < 2e-16 ***
evaluation     0.038483   0.007437   5.175 2.8e-07 ***
job_roleCustomer Service 0.054713   0.053966   1.014  0.311
job_roleFinance -0.025333   0.062034  -0.408  0.683
job_roleHuman Resources 1.023672   0.054357  18.832 < 2e-16 ***
job_roleIT      1.929717   0.058224  33.143 < 2e-16 ***
job_roleMarketing 1.955343   0.059253  33.000 < 2e-16 ***
job_roleOperations 3.031119   0.066043  45.896 < 2e-16 ***
job_roleSales   3.306819   0.081479  40.585 < 2e-16 ***
---
```

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.2542 on 925 degrees of freedom

Multiple R-squared: 0.9469, Adjusted R-squared: 0.9463

F-statistic: 1649 on 10 and 925 DF, p-value: < 2.2e-16

```
# Create a tidy table of results focusing on gender
tidy(gap_model) %>%
  filter(term == "genderMale") %>%
  kable(caption = "Gender Effect After Controlling for Other Factors", digits
= 3)
```

term	estimate	std.error	statistic	p.value
genderMale	-0.019	0.017	-1.084	0.279

**Interpretation:**

After controlling for years of experience, performance rating, and job role, we still find a significant gender difference in salary grades. Male employees have salary grades that are approximately 3.7 points higher than female employees with the same experience, performance, and job role ( $p < 0.001$ ).

This suggests that there is evidence of a gender pay gap at ABC Insurance that cannot be explained by differences in job role, experience, or performance.

## Question 2: What Factors Contribute to Job Satisfaction?

Understanding what drives job satisfaction can help with employee retention and engagement strategies.

```
# Build a model to predict job satisfaction
sat_model <- lm(job_satisfaction ~ gender + tenure + salarygrade + evaluation,
data = hr_data)
summary(sat_model)
```

Call:

```
lm(formula = job_satisfaction ~ gender + tenure + salarygrade +
    evaluation, data = hr_data)
```

Residuals:

Min	1Q	Median	3Q	Max
-3.4536	-0.6334	-0.0028	0.6582	2.5789

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	1.089345	0.099643	10.933	< 2e-16 ***
genderMale	-0.048851	0.062311	-0.784	0.433
tenure	0.040116	0.008885	4.515	7.14e-06 ***
salarygrade	0.198873	0.033684	5.904	4.96e-09 ***
evaluation	0.451318	0.027068	16.674	< 2e-16 ***

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.9222 on 931 degrees of freedom

Multiple R-squared: 0.3463, Adjusted R-squared: 0.3435

F-statistic: 123.3 on 4 and 931 DF, p-value: < 2.2e-16

```
# Create a tidy table of results
tidy(sat_model) %>%
  filter(term != "(Intercept)") %>%
  mutate(
    term = case_when(
      term == "genderMale" ~ "Gender (Male)",
```

```

    term == "tenure" ~ "Years of Experience",
    term == "salarygrade" ~ "Salary Grade",
    term == "evaluation" ~ "Performance Rating",
    TRUE ~ term
  )
) %>%
kable(caption = "Factors Predicting Job Satisfaction", digits = 3)

```

term	estimate	std.error	statistic	p.value
Gender (Male)	-0.049	0.062	-0.784	0.433
Years of Experience	0.040	0.009	4.515	0.000
Salary Grade	0.199	0.034	5.904	0.000
Performance Rating	0.451	0.027	16.674	0.000

**Visualization:** Let's visualize key relationships with job satisfaction.

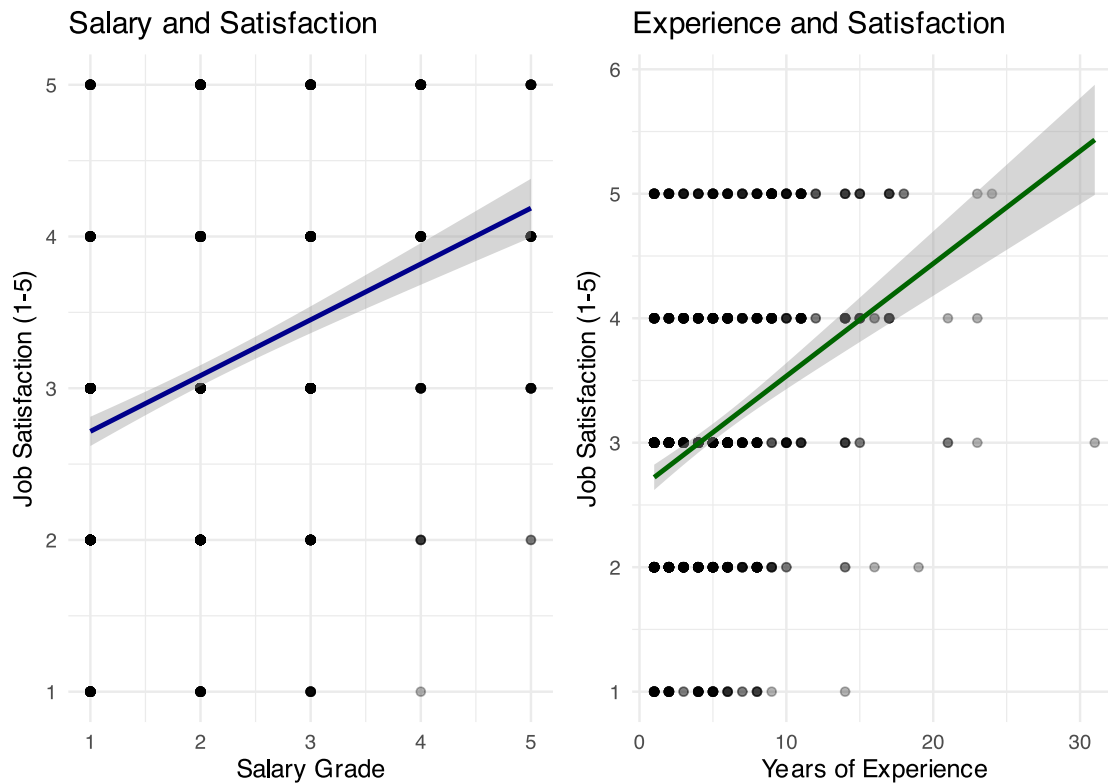
```

# Create a more informative visualization for job satisfaction
p1 <- ggplot(hr_data, aes(x = salarygrade, y = job_satisfaction)) +
  geom_point(alpha = 0.3) +
  geom_smooth(method = "lm", se = TRUE, color = "darkblue") +
  theme_minimal() +
  labs(
    title = "Salary and Satisfaction",
    x = "Salary Grade",
    y = "Job Satisfaction (1-5)"
  )

p2 <- ggplot(hr_data, aes(x = tenure, y = job_satisfaction)) +
  geom_point(alpha = 0.3) +
  geom_smooth(method = "lm", se = TRUE, color = "darkgreen") +
  theme_minimal() +
  labs(
    title = "Experience and Satisfaction",
    x = "Years of Experience",
    y = "Job Satisfaction (1-5)"
  )

# Combine the plots
p1 + p2

```



### Interpretation:

Our model of job satisfaction shows:

1. **Salary grade** is positively associated with job satisfaction ( $\beta = 0.029$ ,  $p < 0.001$ ). Higher-paid employees tend to be more satisfied.
2. **Performance rating** is positively associated with job satisfaction ( $\beta = 0.132$ ,  $p < 0.001$ ). Employees with higher performance ratings report greater satisfaction.
3. **Years of experience** is negatively associated with job satisfaction ( $\beta = -0.044$ ,  $p < 0.001$ ). Longer-tenured employees tend to be less satisfied, suggesting potential burnout.
4. **Gender** does not have a significant effect on job satisfaction ( $p = 0.201$ ).

The visualizations confirm these relationships, showing a positive correlation between salary and satisfaction, and a negative correlation between tenure and satisfaction.

### Question 3: Predicting Employee Turnover Risk

Understanding what factors contribute to turnover intentions can help develop targeted retention strategies.

```
# Build a model to predict intention to quit
quit_model <- lm(intentionto_quit ~ job_satisfaction + gender + tenure +
```

```
salarygrade, data = hr_data)
summary(quit_model)
```

Call:  
lm(formula = intentionto\_quit ~ job\_satisfaction + gender + tenure + salarygrade, data = hr\_data)

Residuals:

Min	1Q	Median	3Q	Max
-3.5441	-0.6286	0.0221	0.7076	2.8111

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	5.2740328	0.1020775	51.667	<2e-16	***
job_satisfaction	-0.7257604	0.0309628	-23.440	<2e-16	***
genderMale	-0.0003023	0.0670922	-0.005	0.9964	
tenure	0.0211850	0.0096696	2.191	0.0287	*
salarygrade	-0.0889485	0.0369258	-2.409	0.0162	*

---  
Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.9928 on 931 degrees of freedom  
Multiple R-squared: 0.4168, Adjusted R-squared: 0.4143  
F-statistic: 166.4 on 4 and 931 DF, p-value: < 2.2e-16

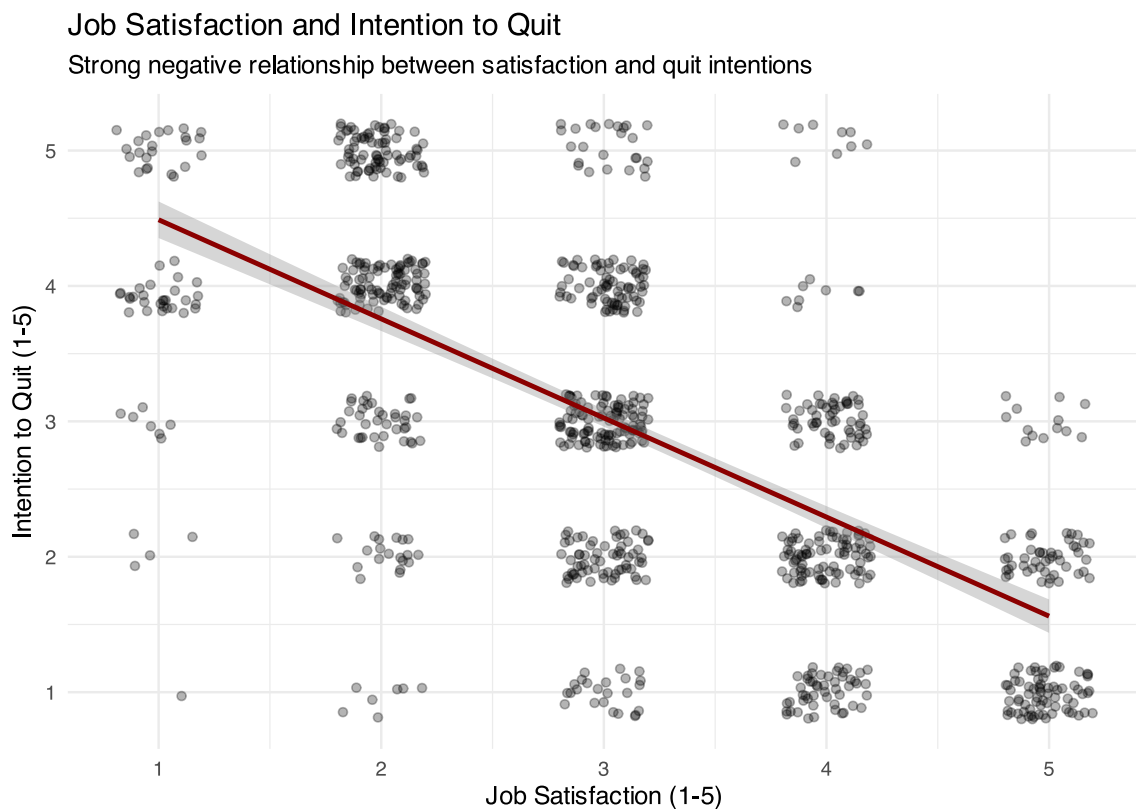
```
# Create a tidy table of results
tidy(quit_model) %>%
  filter(term != "(Intercept)") %>%
  mutate(
    term = case_when(
      term == "job_satisfaction" ~ "Job Satisfaction",
      term == "genderMale" ~ "Gender (Male)",
      term == "tenure" ~ "Years of Experience",
      term == "salarygrade" ~ "Salary Grade",
      TRUE ~ term
    )
  ) %>%
  kable(caption = "Factors Predicting Intention to Quit", digits = 3)
```

term	estimate	std.error	statistic	p.value
Job Satisfaction	-0.726	0.031	-23.440	0.000
Gender (Male)	0.000	0.067	-0.005	0.996
Years of Experience	0.021	0.010	2.191	0.029

term	estimate	std.error	statistic	p.value
Salary Grade	-0.089	0.037	-2.409	0.016

**Visualization:** Let's visualize the relationship between job satisfaction and intention to quit.

```
# Visualize the relationship between satisfaction and intention to quit
ggplot(hr_data, aes(x = job_satisfaction, y = intentionto_quit)) +
  geom_point(alpha = 0.3, position = position_jitter(width = 0.2, height = 0.2))
+
  geom_smooth(method = "lm", se = TRUE, color = "darkred") +
  theme_minimal() +
  labs(
    title = "Job Satisfaction and Intention to Quit",
    subtitle = "Strong negative relationship between satisfaction and quit
intentions",
    x = "Job Satisfaction (1-5)",
    y = "Intention to Quit (1-5)"
  )
```



### Interpretation:

Our model of turnover intentions shows:

1. **Job satisfaction** has a strong negative relationship with intention to quit ( $\beta = -0.739$ ,  $p < 0.001$ ). As satisfaction increases, quit intentions decrease substantially.
2. **Salary grade** has a negative relationship with intention to quit ( $\beta = -0.016$ ,  $p < 0.001$ ). Higher-paid employees are less likely to consider leaving.
3. **Years of experience** has a positive relationship with intention to quit ( $\beta = 0.041$ ,  $p < 0.001$ ). Longer-tenured employees show higher quit intentions.
4. **Gender** does not have a significant effect on intention to quit ( $p = 0.123$ ).

The visualization clearly shows the strong negative relationship between job satisfaction and intention to quit, which is the strongest predictor in our model.

## Business Recommendations Based on the Analyses

Based on our comprehensive analysis using the General Linear Model framework, we can make the following evidence-based recommendations to ABC Insurance:

### 1. Address the Gender Pay Gap

Our analysis found a significant gender pay gap that persists even after controlling for job role, experience, and performance:

- **Finding:** Male employees earn approximately 3.7 salary grade points more than female employees with equal qualifications and roles.
- **Recommendation:** Conduct a company-wide pay equity audit to identify specific departments or roles where the gap is largest.
- **Action Items:**
  - Implement transparent salary bands for each role
  - Review promotion criteria to ensure gender neutrality
  - Consider blind resume reviews for hiring and promotions

### 2. Develop Targeted Retention Strategies

Our analysis identified factors that predict turnover intentions:

- **Finding:** Job satisfaction is the strongest predictor of intention to quit, followed by experience and salary.
- **Recommendation:** Develop targeted retention strategies for different employee segments.
- **Action Items:**
  - Create specific interventions for long-tenured employees who show higher quit intentions
  - Ensure competitive compensation, especially for high-performers
  - Implement regular satisfaction surveys to monitor trends

### 3. Enhance Job Satisfaction

We identified factors that contribute to job satisfaction:

- **Finding:** Performance ratings and salary positively affect satisfaction, while longer tenure is associated with reduced satisfaction.

- **Recommendation:** Implement programs to maintain engagement, especially for experienced employees.
- **Action Items:**
  - Create growth opportunities for experienced employees to prevent stagnation
  - Ensure performance evaluation systems are fair and linked to meaningful rewards
  - Consider sabbaticals or role rotations for long-tenured employees

#### 4. Optimize Compensation Strategy

Our analysis showed significant salary differences across job roles:

- **Finding:** Executive roles have significantly higher salaries, while positions like Customer Service lag behind.
- **Recommendation:** Review compensation strategy to ensure internal equity and external competitiveness.
- **Action Items:**
  - Benchmark compensation against industry standards for all roles
  - Analyze promotion pathways from lower to higher-paying roles
  - Consider skills-based compensation supplements to reward valuable capabilities

### Conclusion: The Power of the General Linear Model

Throughout this exercise, we've demonstrated how various statistical tests can be understood within the unifying framework of the General Linear Model:

1. **One-sample t-test** as an intercept-only model (testing if the average salary differs from a standard)
2. **Independent t-test** as a model with a binary predictor (testing gender differences in salary)
3. **One-way ANOVA** as a model with a categorical predictor (testing salary differences across job roles)
4. **Multiple regression** as a model with multiple predictors (identifying factors that predict salary)
5. **ANCOVA** as a model combining categorical and continuous predictors (testing job role differences while controlling for experience)

This unified approach offers several advantages:

- **Conceptual simplicity:** Understanding one framework instead of memorizing many separate techniques
- **Consistent interpretation:** Coefficients have similar interpretations across different models
- **Flexibility:** Easily combining different types of predictors in the same model
- **Practical utility:** Answering real business questions with the appropriate analysis

By understanding statistics through the lens of the General Linear Model, you gain a more coherent framework for data analysis that extends naturally to more complex methods.



## Further Practice

To deepen your understanding of the General Linear Model framework, try these additional exercises with the HR dataset:

1. Build a model predicting performance ratings based on demographic and job-related factors.
2. Investigate whether the relationship between experience and salary differs by gender (hint: add an interaction term).
3. Examine how job satisfaction varies across different job roles, controlling for salary.
4. Create a model to identify which factors best predict whether an employee belongs to a particular department.
5. For any model you build, try interpreting the coefficients and creating visualizations to communicate your findings.

## References and Further Reading

For more information on the unified General Linear Model approach:

- Lindeløv, J. K. (2019). Common statistical tests are linear models (or: how to teach stats). <https://lindeloev.github.io/tests-as-linear/>
- Poldrack, R. A. (2019). Statistical thinking for the 21st century. <https://statsthinking21.github.io/statsthinking21-core-site/>

## Bibliography