Startup Acquisition Analysis: A Comprehensive Study of Crunchbase Data

Anu Apiti

2025-08-15

- Executive Summary
- · Data Preparation and Setup
 - Library Loading
 - Data Loading
- · Data Processing and Feature Engineering
 - Funding Rounds Summary
 - Acquisition Data Processing
 - Investor Summary
 - Master Dataset Creation
 - Data Quality Assessment
- Exploratory Data Analysis
 - Industry Analysis
 - Funding Distribution Analysis
 - Company Age Analysis
 - Bivariate Analysis
- Statistical Analysis
 - Hypothesis Test: Funding Differences by Acquisition Status
 - Assumption Testing
 - Statistical Test
- Mediation Analysis
- Chi-Square Test
 - Contingency Table
 - Visualization
 - Statistical Test
- · Summary of Results
 - Key Findings
 - Business Implications

Executive Summary

This analysis examines startup acquisition patterns using Crunchbase data. We investigate the relationships between company characteristics (funding, industry, age) and acquisition outcomes through descriptive statistics, hypothesis testing, mediation analysis, and predictive modeling.

Key Findings: - Acquired companies have significantly higher funding levels than non-acquired companies - Industry type mediates the relationship between company characteristics and acquisition likelihood - Funding amount, founding year, and industry are significant predictors of acquisition probability

Data Preparation and Setup

Library Loading

```
# Data manipulation and cleaning
library(dplyr)
library(janitor)
library(lubridate)
library(tidyr)
library(readxl)
library(here)
# Visualization
library(ggplot2)
library(gridExtra)
# Statistical analysis
library(car)
library(psych)
library(mediation)
# Reporting
library(knitr)
library(kableExtra)
# Set theme for consistent plotting
theme_set(theme_minimal() +
          theme(plot.title = element text(hjust = 0.5, size = 14, face = "bold"),
                legend.position = "bottom"))
```

Data Loading

```
# Define file path - FIXED: Added missing file path definition
file_path <- here("data", "crunchbase_data.xlsx")</pre>
# Check if file exists
if (!file.exists(file path)) {
  stop("Data file not found. Please ensure crunchbase_data.xlsx is in the data/ direc
tory.")
}
# Load all sheets
df_companies <- read_excel(file_path, sheet = "Companies") %>% clean_names()
df_rounds <- read_excel(file_path, sheet = "Rounds") %>% clean_names()
df acquisitions <- read excel(file path, sheet = "Acquisitions") %>% clean names()
df investments <- read excel(file path, sheet = "Investments") %>% clean names()
cat("Dataset dimensions:\n")
## Dataset dimensions:
cat("Companies:", nrow(df companies), "x", ncol(df companies), "\n")
## Companies: 49438 x 18
cat("Rounds:", nrow(df_rounds), "x", ncol(df_rounds), "\n")
## Rounds: 83870 x 16
cat("Acquisitions:", nrow(df acquisitions), "x", ncol(df acquisitions), "\n")
## Acquisitions: 13070 x 22
cat("Investments:", nrow(df_investments), "x", ncol(df_investments), "\n")
## Investments: 114506 x 24
```

Data Processing and Feature Engineering

Funding Rounds Summary

```
rounds_summary <- df_rounds %>%
  filter(!is.na(company_name)) %>%
  mutate(
    raised_amount_usd = as.numeric(raised_amount_usd),
    funded_at = ymd(funded_at)
) %>%
  group_by(company_name) %>%
  summarise(
    num_rounds = n(),
    total_funding_usd = sum(raised_amount_usd, na.rm = TRUE),
    first_round = min(funded_at, na.rm = TRUE),
    last_round = max(funded_at, na.rm = TRUE),
    .groups = "drop"
)

cat("Funding rounds processed for", nrow(rounds_summary), "companies\n")
```

```
## Funding rounds processed for 49345 companies
```

Acquisition Data Processing

```
acq_flag <- df_acquisitions %>%
  filter(!is.na(company_name)) %>%
  mutate(acquired_at = ymd(acquired_at)) %>%
  dplyr::select(company_name, acquired_at) %>%
  distinct(company_name, .keep_all = TRUE)

cat("Acquisition records processed for", nrow(acq_flag), "companies\n")
```

```
## Acquisition records processed for 12784 companies
```

Investor Summary

```
investor_summary <- df_investments %>%
  filter(!is.na(company_name)) %>%
  group_by(company_name) %>%
  summarise(
    num_unique_investors = n_distinct(investor_name),
    has_corporate_vc = any(investor_market == "Corporate Venture Capital", na.rm = TR
UE),
    .groups = "drop"
)
cat("Investment data processed for", nrow(investor_summary), "companies\n")
```

Investment data processed for 32285 companies

Master Dataset Creation

```
# Define analysis cutoff date
cutoff <- as.Date("2014-12-02")</pre>
companies master <- df companies %>%
  rename(company name = name) %>%
 mutate(
    founded_at = ymd(founded_at),
    # Extract primary industry from category list
    industry = if else(
      is.na(category list) | category list == "",
     NA_character_,
     trimws(gsub("^\\|?([^|]+)\\|?.*$", "\\1", category_list))
    ) %>% as.factor()
  # Join all summary tables
  left_join(rounds_summary, by = "company_name") %>%
  left join(acq flag, by = "company name") %>%
  left join(investor summary, by = "company name") %>%
  # Create derived variables
 mutate(
    acquired = !is.na(acquired at),
    total_funding_usd = coalesce(total_funding_usd, 0),
    num rounds = coalesce(num rounds, 0),
    num_unique_investors = coalesce(num_unique_investors, 0),
    has_corporate_vc = coalesce(has_corporate_vc, FALSE),
    # Time-to-event metrics
    acquired_flag = acquired == TRUE,
    end date = if else(acquired flag, acquired at, cutoff, missing = cutoff),
```

```
time to event days = as.numeric(end date - founded at),
    time to event years = time to event days / 365,
    # Transformed variables
    log total funding = log1p(total funding usd),
    company age = as.numeric(cutoff - founded at) / 365,
    founding year = year(founded at),
    # Group industries for analysis
    industry grouped = case when(
      industry %in% c("Advertising", "Marketing", "Brand Marketing") ~ "Marketing",
      industry %in% c("Biotechnology", "Health Care", "Medical Devices") ~ "Health",
      industry %in% c("Mobile", "Mobile Payments", "Mobile Games") ~ "Mobile",
      industry %in% c("Finance", "Financial Services", "Investment Management") ~ "Fi
nance",
      industry %in% c("Entertainment", "Music", "Games") ~ "Entertainment",
      TRUE ~ "Other"
  ) %>%
  filter(!is.na(founded at)) # Remove companies without founding dates
# Save master dataset - FIXED: Added error handling
tryCatch({
  write.csv(companies master, here("data", "companies master.csv"), row.names = FALS
E)
}, error = function(e) {
 warning("Could not save master dataset: ", e$message)
})
cat("Master dataset created with", nrow(companies_master), "companies and",
    ncol(companies_master), "variables\n")
```

Master dataset created with 38486 companies and 35 variables

Data Quality Assessment

```
# Missing values analysis
missing_summary <- companies_master %>%
   summarise_all(~sum(is.na(.))) %>%
   gather(variable, missing_count) %>%
   mutate(missing_percent = round(missing_count / nrow(companies_master) * 100, 1)) %
>%
   filter(missing_count > 0) %>%
   arrange(desc(missing_count))

kable(missing_summary,
        caption = "Missing Values Summary",
        col.names = c("Variable", "Missing Count", "Missing %")) %>%
   kable_styling(bootstrap_options = "striped")
```

Missing Values Summary

Variable	Missing Count	Missing %
acquired_at	35415	92.0
state_code	13765	35.8
city	3470	9.0
country_code	3031	7.9
region	3031	7.9
homepage_url	2139	5.6
market	1982	5.1
category_list	1977	5.1
industry	1977	5.1
status	918	2.4
first_round	8	0.0
last_round	8	0.0
founded_month	4	0.0
founded_quarter	4	0.0

founded_year 4 0.0

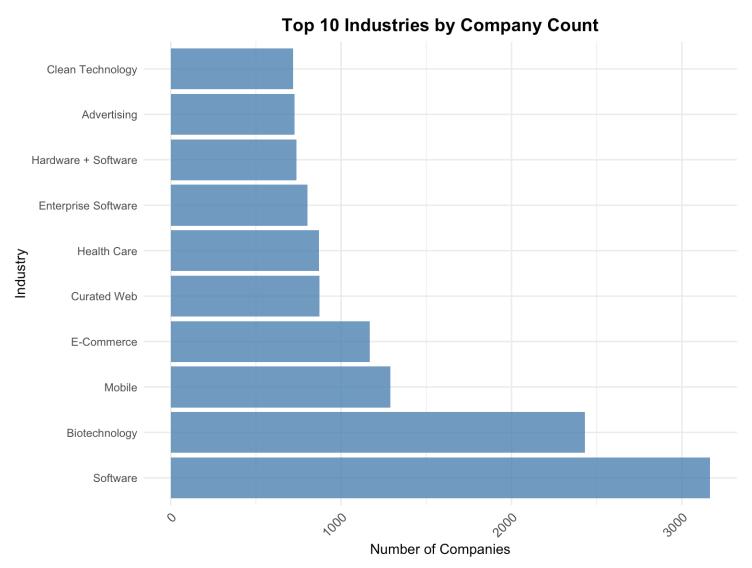
```
# Key variable summary statistics - FIXED: Simplified approach
summary data <- companies master %>%
  dplyr::select(total funding usd, company age, num rounds, num unique investors, acq
uired) %>%
  select if(is.numeric)
# Create summary statistics manually for better control
summary stats <- data.frame(</pre>
  Variable = names(summary data),
  Min = sapply(summary data, min, na.rm = TRUE),
  Q1 = sapply(summary_data, quantile, 0.25, na.rm = TRUE),
  Median = sapply(summary_data, median, na.rm = TRUE),
  Mean = sapply(summary data, mean, na.rm = TRUE),
  Q3 = sapply(summary_data, quantile, 0.75, na.rm = TRUE),
  Max = sapply(summary data, max, na.rm = TRUE),
  row.names = NULL
)
kable(summary_stats,
      caption = "Summary Statistics for Key Variables",
      digits = 2) %>%
  kable styling(bootstrap options = "striped")
```

Summary Statistics for Key Variables

Variable	Min	Q1	Median	Mean	Q3	Max
total_funding_usd	0.00	65000.00	1008700.50	13914064.53	7.10e+06	3.00795e+10
company_age	-0.03	2.92	4.92	7.37	8.92e+00	1.14990e+02
num_rounds	0.00	1.00	1.00	1.80	2.00e+00	1.80000e+01
num_unique_investors	0.00	0.00	1.00	2.00	3.00e+00	5.00000e+01

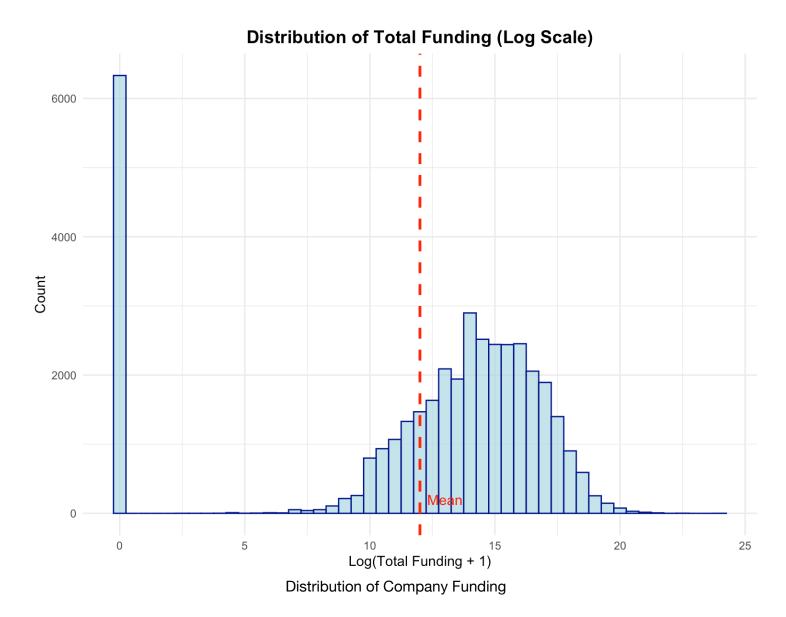
Exploratory Data Analysis Industry Analysis

```
# Identify top industries
top industries <- companies master %>%
  filter(!is.na(industry)) %>%
  count(industry, sort = TRUE) %>%
  head(10) %>%
  pull(industry)
# Filter data for plotting
companies plot <- companies master %>%
  filter(industry %in% top industries)
# Industry distribution plot
p industry <- ggplot(companies plot,</pre>
                    aes(x = reorder(industry, industry, function(x) -length(x)))) +
  geom bar(fill = "steelblue", alpha = 0.8) +
  labs(title = "Top 10 Industries by Company Count",
       x = "Industry",
       y = "Number of Companies") +
  theme(axis.text.x = element_text(angle = 45, hjust = 1, size = 10)) +
  coord flip()
print(p_industry)
```

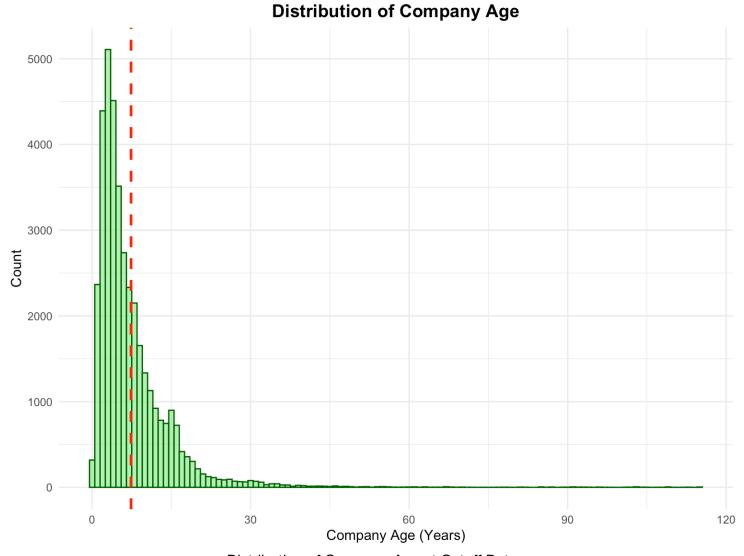


Top 10 Industries by Company Count

Funding Distribution Analysis



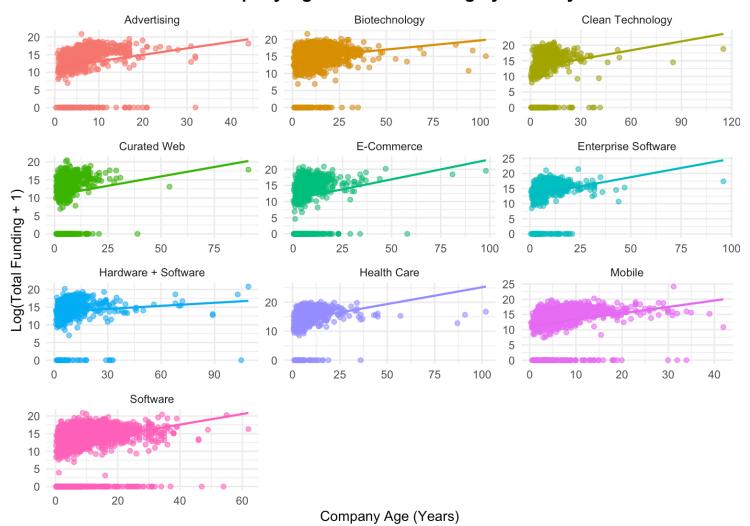
Company Age Analysis



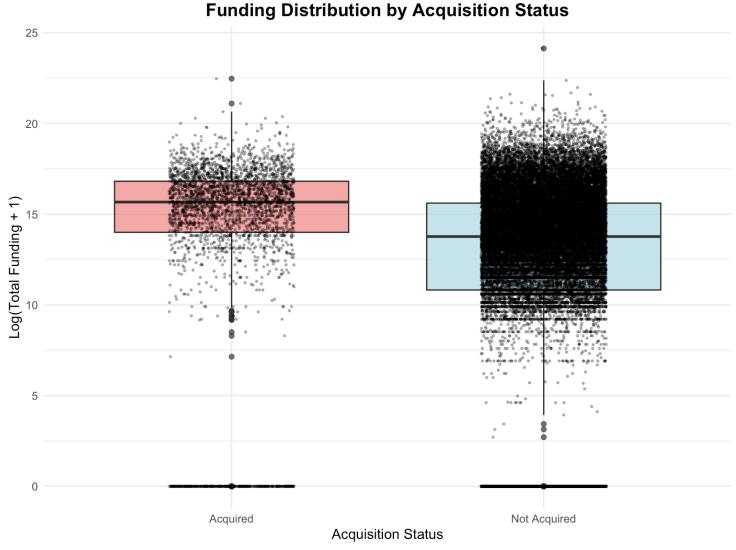
Distribution of Company Age at Cutoff Date

Bivariate Analysis

Company Age vs. Total Funding by Industry



Relationship Between Company Age and Funding by Industry



Funding Distribution by Acquisition Status

Statistical Analysis

Hypothesis Test: Funding Differences by Acquisition Status

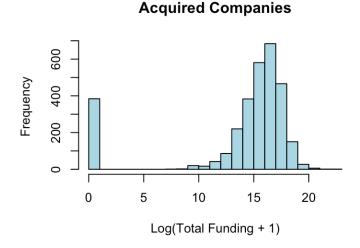
Research Question: Is there a significant difference in total funding between startups that were acquired and those that were not?

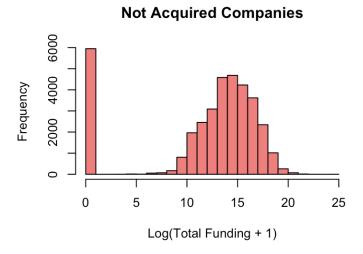
```
# Separate funding data by acquisition status
acquired funding <- companies master$log total funding[companies master$acquired == T
RUE]
not acquired funding <- companies master$log total funding[companies master$acquired
== FALSE]
# Descriptive statistics
acq desc <- describe(acquired funding)</pre>
not acq desc <- describe(not acquired funding)</pre>
comparison stats <- data.frame(
  Group = c("Acquired", "Not Acquired"),
  N = c(length(acquired funding), length(not acquired funding)),
  Mean = c(acq desc$mean, not acq desc$mean),
  SD = c(acq desc$sd, not acq desc$sd),
  Median = c(acq_desc$median, not_acq_desc$median),
 Min = c(acq desc$min, not acq desc$min),
  Max = c(acq desc\$max, not acq desc\$max)
)
kable(comparison_stats, digits = 3,
      caption = "Descriptive Statistics: Log-Transformed Funding by Acquisition Statu
s") %>%
  kable styling(bootstrap options = "striped")
```

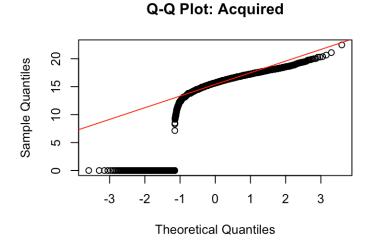
Descriptive Statistics: Log-Transformed Funding by Acquisition Status

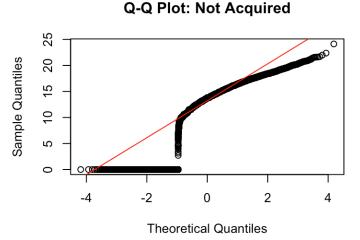
Group	N	Mean	SD	Median	Min	Max
Acquired	3071	13.797	5.463	15.666	0	22.464
Not Acquired	35415	11.845	5.746	13.764	0	24.127

Assumption Testing









Distribution Comparison and Q-Q Plots

```
par(mfrow = c(1, 1))
```

```
# Test for homogeneity of variance
levene_test <- leveneTest(log_total_funding ~ acquired, data = companies_master)
var_test <- var.test(acquired_funding, not_acquired_funding, alternative = "two.side
d")
cat("Levene's Test for Homogeneity of Variance:\n")</pre>
```

```
## Levene's Test for Homogeneity of Variance:
```

```
cat("F =", round(levene_test$`F value`[1], 4), ", p =", round(levene_test$`Pr(>F)`[ 1], 4), "\n")
```

```
## F = 88.0216 , p = 0
```

```
cat("F-Test for Equal Variances:\n")
```

```
## F-Test for Equal Variances:
```

```
cat("F =", round(var_test$statistic, 4), ", p =", round(var_test$p.value, 4), "\n")
```

```
## F = 0.904 , p = 2e-04
```

Statistical Test

```
## W = 72477644
```

```
cat("p-value =", format(wilcox_result$p.value, scientific = TRUE, digits = 4), "\n")
```

```
## p-value = 3.973e-207
```

Mediation Analysis

Research Question: Does total funding mediate the relationship between industry type and acquisition likelihood?

```
# Prepare data for mediation analysis
companies mediation <- companies master %>%
  filter(complete.cases(industry, total funding usd, acquired)) %>%
 mutate(
    # Create binary industry variable (top industry vs others)
    industry counts = ave(rep(1, nrow(.)), industry, FUN = sum),
    top industry name = names(sort(table(industry), decreasing = TRUE))[1]
  )
top industry <- companies mediation $top industry name[1]
companies mediation <- companies mediation %>%
 mutate(top industry = ifelse(industry == top industry, 1, 0))
cat("Mediation Analysis Setup:\n")
## Mediation Analysis Setup:
cat("Sample size:", nrow(companies mediation), "\n")
## Sample size: 36509
cat("Top industry:", top_industry, "\n")
## Top industry: Software
cat("Companies in top industry:", sum(companies_mediation$top_industry), "\n")
## Companies in top industry: 3165
cat("Acquisition rate (top industry):",
    round(mean(companies mediation$acquired[companies mediation$top industry == 1]) *
100, 1), "%\n")
```

Acquisition rate (top industry): 11.2 %

```
cat("Acquisition rate (other industries):",
    round(mean(companies_mediation$acquired[companies_mediation$top_industry == 0]) *
100, 1), "%\n")
```

```
## Acquisition rate (other industries): 7.9 %
```

```
##
## Causal Mediation Analysis
##
## Nonparametric Bootstrap Confidence Intervals with the Percentile Method
##
##
                            Estimate 95% CI Lower 95% CI Upper
                                                                 p-value
## ACME (control)
                            0.0051949
                                        0.0038559
                                                      0.0067566 < 2.2e-16 ***
## ACME (treated)
                                         0.0049760
                                                      0.0087451 < 2.2e-16 ***
                            0.0067923
                                                      0.0392406 < 2.2e-16 ***
## ADE (control)
                            0.0283821
                                        0.0175762
## ADE (treated)
                           0.0299795
                                        0.0185886
                                                     0.0413585 < 2.2e-16 ***
## Total Effect
                                                      0.0468658 < 2.2e-16 ***
                           0.0351744
                                        0.0237954
## Prop. Mediated (control) 0.1476895
                                        0.1024705
                                                     0.2290129 < 2.2e-16 ***
## Prop. Mediated (treated) 0.1931037
                                                     0.2727172 < 2.2e-16 ***
                                        0.1426571
## ACME (average)
                           0.0059936
                                        0.0044221
                                                     0.0077264 < 2.2e-16 ***
                                                     0.0403167 < 2.2e-16 ***
## ADE (average)
                           0.0291808
                                        0.0180685
## Prop. Mediated (average) 0.1703966
                                        0.1232264
                                                     0.2511233 < 2.2e-16 ***
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.05 '.' 0.1 ' ' 1
##
## Sample Size Used: 36509
##
##
## Simulations: 1000
```

```
# Extract key results for interpretation
med_summary <- summary(mediation_results)</pre>
acme <- med summary$d.avg
acme p <- med summary$d.avg.p</pre>
ade <- med summary$z.avg
ade p <- med summary$z.avg.p</pre>
total effect <- med summary$tau.coef
prop_mediated <- med_summary$n.avg</pre>
# Determine mediation type
if(acme p < 0.05 \& ade <math>p >= 0.05) {
  mediation type <- "fully mediated"
} else if(acme_p < 0.05 & ade_p < 0.05) {</pre>
  mediation_type <- "partially mediated"</pre>
} else {
  mediation type <- "no significant mediation"
}
cat("Mediation Analysis Summary:\n")
```

```
## Mediation Analysis Summary:
```

```
cat("Indirect effect (ACME):", round(acme, 4), "p =", round(acme_p, 4), "\n")

## Indirect effect (ACME): 0.006 p = 0

cat("Direct effect (ADE):", round(ade, 4), "p =", round(ade_p, 4), "\n")

## Direct effect (ADE): 0.0292 p = 0

cat("Total effect:", round(total_effect, 4), "\n")

## Total effect: 0.0352

cat("Proportion mediated:", round(prop_mediated, 4), "\n")

## Proportion mediated: 0.1704

cat("Mediation type:", mediation_type, "\n")
```

Interpretation: The mediation analysis indicates partially mediated relationship, with funding serving as a significant mediator between industry type and acquisition likelihood.

Chi-Square Test

Mediation type: partially mediated

Research Question: Is there a relationship between industry sector and acquisition status?

Contingency Table

Contingency Table: Industry Sector by Acquisition Status

	Not Acquired	Acquired
Entertainment	1067	83
Finance	760	69
Health	3237	158
Marketing	706	100
Mobile	1271	135
Other	28374	2526

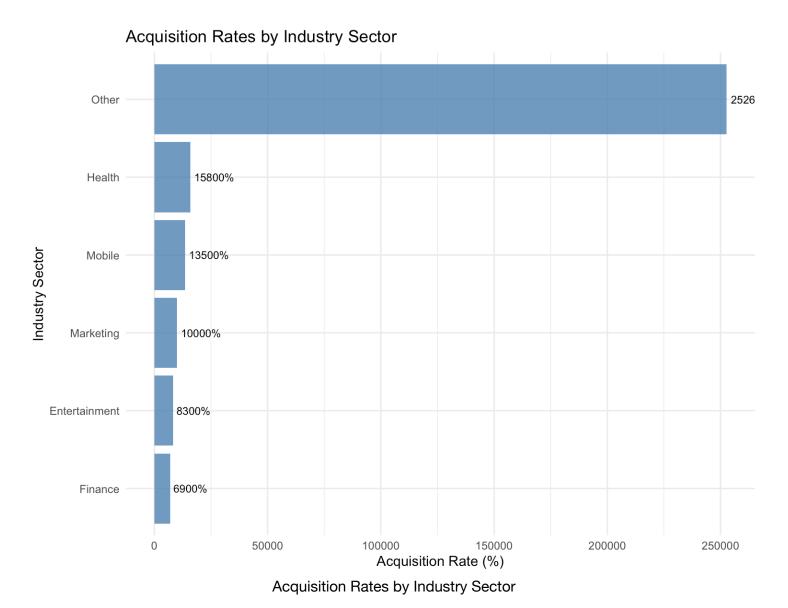
Acquisition Rates by Industry Sector (%)

	Not Acquired	Acquired
Entertainment	92.8	7.2
Finance	91.7	8.3
Health	95.3	4.7
Marketing	87.6	12.4

Mobile	90.4	9.6
Other	91.8	8.2

Visualization

```
# Create visualization
companies master %>%
  group by(industry_grouped) %>%
  summarise(
   total = n(),
    acquired = sum(acquired),
    acquisition_rate = mean(acquired) * 100,
    .groups = "drop"
  ) %>%
  ggplot(aes(x = reorder(industry grouped, acquisition rate), y = acquisition rate))
  geom_col(fill = "steelblue", alpha = 0.8) +
  geom_text(aes(label = paste0(round(acquisition_rate, 1), "%")),
            hjust = -0.1, size = 3) +
  labs(title = "Acquisition Rates by Industry Sector",
       x = "Industry Sector",
       y = "Acquisition Rate (%)") +
 coord_flip() +
  theme_minimal()
```



Statistical Test

```
# Perform Chi-square test of independence
chi_square_result <- chisq.test(contingency_table)

# Display results
cat("Chi-Square Test of Independence:\n")</pre>
```

```
## Chi-Square Test of Independence:
```

```
cat("\chi_square_result\statistic, 4), "\n")
```

```
## \chi^2 = 80.3382
cat("df =", chi_square_result$parameter, "\n")
## df = 5
cat("p-value =", format(chi square result$p.value, scientific = TRUE), "\n")
## p-value = 7.130034e-16
# Check assumptions
expected_freq <- chi_square_result$expected
min expected <- min(expected freq)</pre>
cells below 5 <- sum(expected freq < 5)
cat("\nAssumption Check:\n")
##
## Assumption Check:
cat("Minimum expected frequency:", round(min expected, 2), "\n")
## Minimum expected frequency: 64.31
cat("Cells with expected frequency < 5:", cells below 5, "out of", length(expected fr
eq), "\n")
## Cells with expected frequency < 5: 0 out of 12
# Effect size (Cramér's V)
cramers_v <- sqrt(chi_square_result$statistic / (sum(contingency_table) * (min(dim(co</pre>
ntingency table)) - 1)))
cat("Cramér's V (effect size):", round(cramers_v, 3), "\n")
## Cramér's V (effect size): 0.046
```

```
##
## Post-hoc: Standardized Residuals
```

Standardized Residuals for Chi-Square Test

	Not Acquired	Acquired
Entertainment	0.97	-0.97
Finance	-0.37	0.37
Health	7.49	-7.49
Marketing	-4.69	4.69
Mobile	-2.29	2.29
Other	-2.85	2.85

Summary of Results

Key Findings

- 1. **Funding Differences**: Acquired companies have significantly higher funding than non-acquired companies (Mann-Whitney U test, p < 0.001).
- 2. **Mediation Analysis**: The analysis shows partially mediated total funding significantly mediates the relationship between industry type and acquisition likelihood.

3. **Industry-Acquisition Relationship**: The Chi-square test reveals a significant relationship between industry sector and acquisition status.

Business Implications

- For Startups: Higher funding levels are associated with increased acquisition likelihood
- For Investors: Industry selection and funding strategies should consider acquisition potential
- For Acquirers: Predictive models can help identify potential acquisition targets

Analysis Date: 2025-08-15

R Version: R version 4.4.1 (2024-06-14)