

EDA

Johannes Harmse

April 9, 2018

```
library(tidyverse)

# removing confidential data
survey_results <- read_csv(file = '../survey_data/Demographic Survey.csv', skip = 1)

## Warning: Missing column names filled in: 'X1' [1], 'X2' [2], 'X3' [3],
## 'X4' [4], 'X5' [5], 'X6' [6], 'X7' [7], 'X8' [8], 'X9' [9]

## Warning: Duplicated column names deduplicated: 'Response' =>
## 'Response_1' [11], 'Response' => 'Response_2' [14]

## Parsed with column specification:
## cols(
##   .default = col_character(),
##   X1 = col_double(),
##   X2 = col_integer(),
##   `Annual Salary (before deductions)` = col_integer(),
##   `Annual salary (before deductions)` = col_integer(),
##   `Living Expenses (utilities, rent, mortgage, transportation, property taxes if owner, etc.)` = col_integer(),
##   `Savings (retirement, investments, emergency funds, etc.)` = col_integer(),
##   `Vacation (lodging, transportation, day trips, etc.)` = col_integer(),
##   `Daily Leisure (eating out, books, movies, self-care, etc.)` = col_integer(),
##   `Consumption Goods (clothing, electronics, other luxury items, etc.)` = col_integer(),
##   `Personal Sports and Hobbies (sporting goods and services, gym, arts and crafts, etc.)` = col_integer(),
##   `Other (health care, taxes, dependent expenses, etc.)` = col_integer()
## )

## See spec(...) for full column specifications.
survey_results <- survey_results[, 10:ncol(survey_results)]

# import data
# survey_results <- read_csv(file = '../survey_data/Demographic Survey.csv') # local path - remove i

# redefine column names
colnames(survey_results) <- c('consent', 'country', 'salary_base', 'salary_expect', 'no_increase_accept',
                             'living_expenses', 'savings', 'vacation', 'daily_leisure', 'consumption_goods',
                             'sports_hobbies', 'other')

# spending categories
spending_cats <- c('living_expenses', 'savings', 'vacation', 'daily_leisure', 'consumption_goods',
                  'sports_hobbies', 'other')

# remove no consent
survey_results <- survey_results %>% filter(consent %in% c('Yes'))

# add observation id
survey_results$id <- 1:nrow(survey_results)
```

```

# save raw clean data
# saveRDS(survey_results, file = '../data/raw/raw_clean.rds')

survey_results %>% head()

## # A tibble: 6 x 13
##   consent country      salary_base salary_expect no_increase_acceptance
##   <chr>   <chr>          <int>         <int> <chr>
## 1 Yes    South Africa    500000        480000 Yes
## 2 Yes    Canada           100000        120000 Yes
## 3 Yes    South Africa    150000        150000 Yes
## 4 Yes    Canada           75000         90000 No
## 5 Yes    <NA>             NA            NA <NA>
## 6 Yes    Nigeria          3000000       3000000 Yes
## # ... with 8 more variables: living_expenses <int>, savings <int>,
## #   vacation <int>, daily_leisure <int>, consumption_goods <int>,
## #   sports_hobbies <int>, other <int>, id <int>

# readRDS(file = '../data/raw/raw_clean.rds')

# get ratio
survey_results <- survey_results %>%
  mutate(ratio = salary_expect/salary_base)

# generic first model
lm_survey <- lm(ratio ~ no_increase_acceptance +
  living_expenses +
  savings +
  vacation +
  daily_leisure +
  consumption_goods +
  sports_hobbies +
  other, data = survey_results)

summary(lm_survey)

##
## Call:
## lm(formula = ratio ~ no_increase_acceptance + living_expenses +
##     savings + vacation + daily_leisure + consumption_goods +
##     sports_hobbies + other, data = survey_results)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.3434 -0.9681 -0.0513  0.3381  7.6359
##
## Coefficients: (1 not defined because of singularities)
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -1.8420404   2.2245253  -0.828   0.411
## no_increase_acceptanceYes  0.3334593   0.4984341   0.669   0.506
## living_expenses    0.0165744   0.0237676   0.697   0.488
## savings           0.0425862   0.0293383   1.452   0.152
## vacation         0.1639607   0.0347074   4.724 1.59e-05 ***
## daily_leisure     0.0023859   0.0333240   0.072   0.943
## consumption_goods  0.0005459   0.0513602   0.011   0.992

```

```
## sports_hobbies          0.0434831  0.0545844  0.797  0.429
## other                   NA          NA      NA      NA
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.791 on 56 degrees of freedom
## (23 observations deleted due to missingness)
## Multiple R-squared:  0.3425, Adjusted R-squared:  0.2604
## F-statistic: 4.168 on 7 and 56 DF,  p-value: 0.0009364

# remove outliers
survey_results <- survey_results %>%
  filter(ratio < 10 &
         ratio > 0.1)

# replace NA spendings with 0

survey_results[, spending_cats][is.na(survey_results[, spending_cats])] <- 0

# generic first model (outliers removed and data cleaned)
lm_survey <- lm(ratio ~ no_increase_acceptance +
               living_expenses +
               savings +
               vacation +
               daily_leisure +
               consumption_goods +
               sports_hobbies +
               other, data = survey_results)

summary(lm_survey)

##
## Call:
## lm(formula = ratio ~ no_increase_acceptance + living_expenses +
##     savings + vacation + daily_leisure + consumption_goods +
##     sports_hobbies + other, data = survey_results)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.69315 -0.12927 -0.02927  0.05505  1.81708
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    1.1215946   0.1458869    7.688 9.57e-11 ***
## no_increase_acceptanceYes -0.1434678   0.0935487   -1.534  0.1299
## living_expenses    -0.0007804   0.0018984   -0.411  0.6823
## savings            0.0088407   0.0038760    2.281  0.0258 *
## vacation          -0.0036982   0.0073424   -0.504  0.6162
## daily_leisure     -0.0008830   0.0054082   -0.163  0.8708
## consumption_goods  -0.0015644   0.0087112   -0.180  0.8580
## sports_hobbies     0.0083979   0.0091745    0.915  0.3633
## other             0.0012247   0.0044876    0.273  0.7858
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
```

```

## Residual standard error: 0.3618 on 66 degrees of freedom
## Multiple R-squared:  0.1269, Adjusted R-squared:  0.02109
## F-statistic: 1.199 on 8 and 66 DF,  p-value: 0.3132

survey_tidy <- NULL

non_spending <- colnames(survey_results)[!(colnames(survey_results) %in% spending_cats)]

for (spending in spending_cats){
  temp <- survey_results[, non_spending]
  temp$spending_cat <- spending
  temp$spending_val <- survey_results[[spending]]
  survey_tidy <- rbind(survey_tidy, temp)
}

for (i in unique(survey_tidy$id)){
  temp <- survey_tidy %>% filter(id == i)
  user_living <- as.numeric(temp %>% filter(temp$spending_cat == 'living_expenses') %>% select(spending_val))
  survey_tidy[survey_tidy$id == i, 'spending_ratio'] <- temp$spending_val/user_living
}

p_vals <- data.frame('category' = character(length(spending_cats)), 'slope' = numeric(length(spending_cats)),
                     'p_value' = numeric(length(spending_cats)))

count <- 0

for (i in spending_cats){
  count <- count + 1
  temp <- survey_tidy %>% filter(spending_cat == i)
  temp <- temp %>% filter(!is.na(spending_ratio) & abs(spending_ratio) != Inf)
  temp_lm <- lm(ratio ~ spending_ratio, data = temp)
  lm_summary <- summary(temp_lm)
  p_vals[count, 'category'] <- as.character(i)
  p_vals[count, 'slope'] <- temp_lm$coefficients[2]
  p_vals[count, 'p_value'] <- ifelse(nrow(lm_summary$coefficients) > 1, lm_summary$coefficients[2, 4],
                                     lm_summary$p.value[2])
}

p_vals

##      category      slope      p_value
## 1 living_expenses      NA          NA
## 2      savings 0.24067177 0.000107678
## 3      vacation 0.13913797 0.306369940
## 4  daily_leisure 0.05004306 0.569492133
## 5 consumption_goods 0.10613198 0.236333806
## 6  sports_hobbies 0.25037279 0.069539351
## 7          other 0.01592914 0.669306708

```