

Exploratory Data Analysis

S. Arora, J. Harmse, V. Mulholland

April 9, 2018

```
library(tidyverse); theme_set(theme_bw())
library(cowplot)
library(ggjoy)
library(reshape2)
```

Overview

The purpose of this study to measure whether a person is driven by money or not. We found it reasonable to assume that a person who is driven by money would expect to earn more than the average person who has the same skillset and experience.

Our survey has captured the salary of what a participant thinks an average person with their skills and experience should earn, as well as the salary that the participant expects to receive in 1 year's time. Taking inflation and other micro-factors into account, a participant's expected salary in a year's time shouldn't be much higher than the average person with the same skills and experience.

The survey captured the participant's salary in their unique currency. The survey was answered by people from various countries with different currencies. This means that we cannot compare the captured salary values between participants. An easy way of standardising these values is to handle the salary values as a ratio of expected salary over average salary. The ratio should be consistent across different currencies.

For the purpose of this study, social standards will be defined as a person's inclination for a high relative consumption on leisure activities and non-essential expenditure. Our hypothesis relies on the theory that prevailing social conditions will influence one's relationship with money which would translate in whether increase in income is the priority.

Data Pre-processing

Anonymity

In order to maintain user privacy a few manipulations were handled before the raw data was uploaded to the analysis repository. Any confidential information such as IP addresses were omitted, as well as any respondents that did not accept the confidentiality agreement.

Pre-processing Workflow

These were the first steps applied to `surveydata_clean.rds` when the data was downloaded raw from *Survey Monkey*.

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

# remove all traces
rm(survey_results)

```

Once the data is pre-processed, it is reimported and the columns and categories are defined.

```

# import clean data
survey_results <- readRDS(file = '../data/processed/surveydata_clean.rds') # local path - remove iden

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

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

```

A new variable was created as a measurement of relative expected increase in salary. The benefits of using a ratio meant that there would be less extra manipulations and potential confounding variables behind adjustments for foreign currencies.

```

# ensure any NA values are set to 0
survey_results[, spending_cats][is.na(survey_results[,spending_cats])] <- 0

# converting char to numeric
survey_results$salary_base <- as.numeric(as.character(survey_results$salary_base))
survey_results$salary_expect <- as.numeric(as.character(survey_results$salary_expect))

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

```

Outlier Handling

Having chosen to remove outliers on the basis that with a small number of observations applying the statistical method of removing outliers greater than two standard deviations could be erroneous since it cannot be deduced with certainty which distribution is being represented. That being said, a combination of visual

assessments and box-plot/quantile analysis allowed a reasonable upper and lower limit to be chosen.

```
# remove outliers
survey_results <- survey_results %>%
  filter(!ratio %in% boxplot.stats(survey_results$ratio)$out)

# replace NA spendings with 0
survey_results[, spending_cats][is.na(survey_results[, spending_cats])] <- 0
```

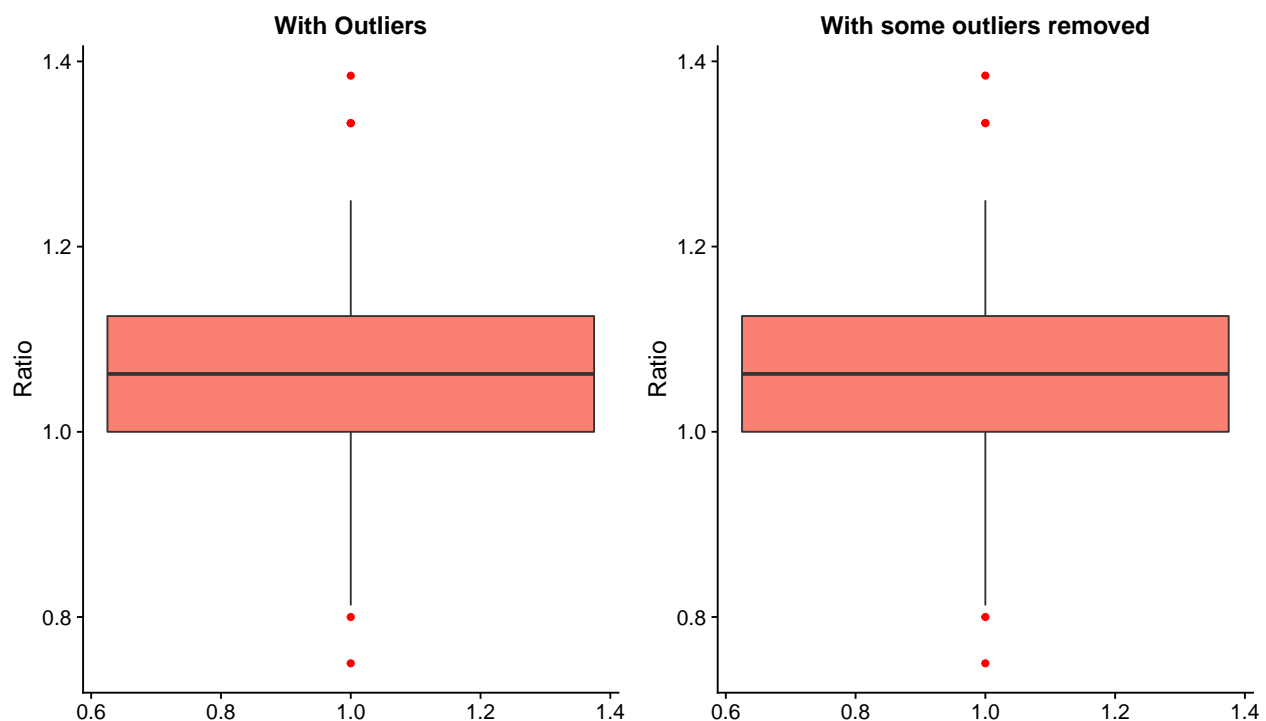
It was decided to remove the values beyond ~95% confidence level. The box-plot method performs a more sophisticated outlier selection than the alternative, the quantile approach, that is more rigid in the 95% threshold. Since we have less observations than ideal, it seemed more appropriate. The visualization below shows the contrast when the most extreme outliers are remove.

```
survey_results_filt <- survey_results %>% filter(ratio > 0.2500000 & ratio < 2)

p1<- ggplot(survey_results)+
  geom_boxplot(aes(x=1, y= ratio), outlier.colour = "red", fill = "salmon") +
  labs(x="",y="Ratio", title = "With Outliers" )

# for visual
p2<- ggplot(survey_results_filt)+
  geom_boxplot(aes(x=1, y = ratio), outlier.colour = "red", fill = "salmon")+
  labs(x="",y="Ratio", title = "With some outliers removed")

# title <- ggdraw() + draw_label("Relative Display of Outliers", fontface='bold')
plot_grid(p1,p2)
```



```
outliers <- boxplot.stats(survey_results$ratio)$out
dh<- data.frame(outliers)
```

```
dh
```

```
## outliers
## 1 0.750000
## 2 1.384615
## 3 1.333333
## 4 0.800000
## 5 1.333333
```

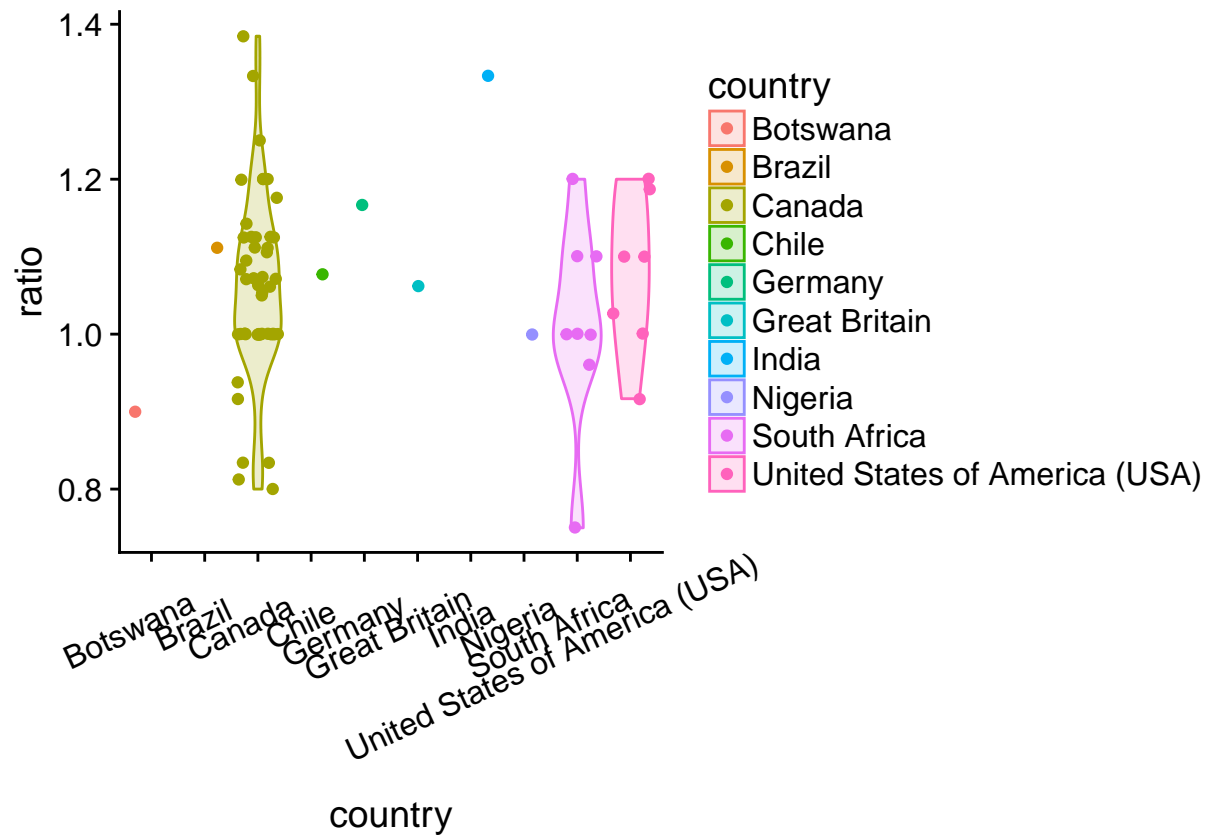
The questions was designed to minimize the potential for entry mistakes when participants entered their responses. A rule was included to ensure that the expenditure percentages summed up to 100 points, but this was not possible with the user salary through the *Survey Monkey* interface. This process of removing outliers will filter out major mistakes in currency where the user entered that they expected a very disproportionate salary increase.

Below each variable is summarized. Since it is difficult to highlight important information from a summary table containing so many variables, a jitter-violin plot was also generated.

```
sum.tb <- summary(survey_results)
sum.tb
```

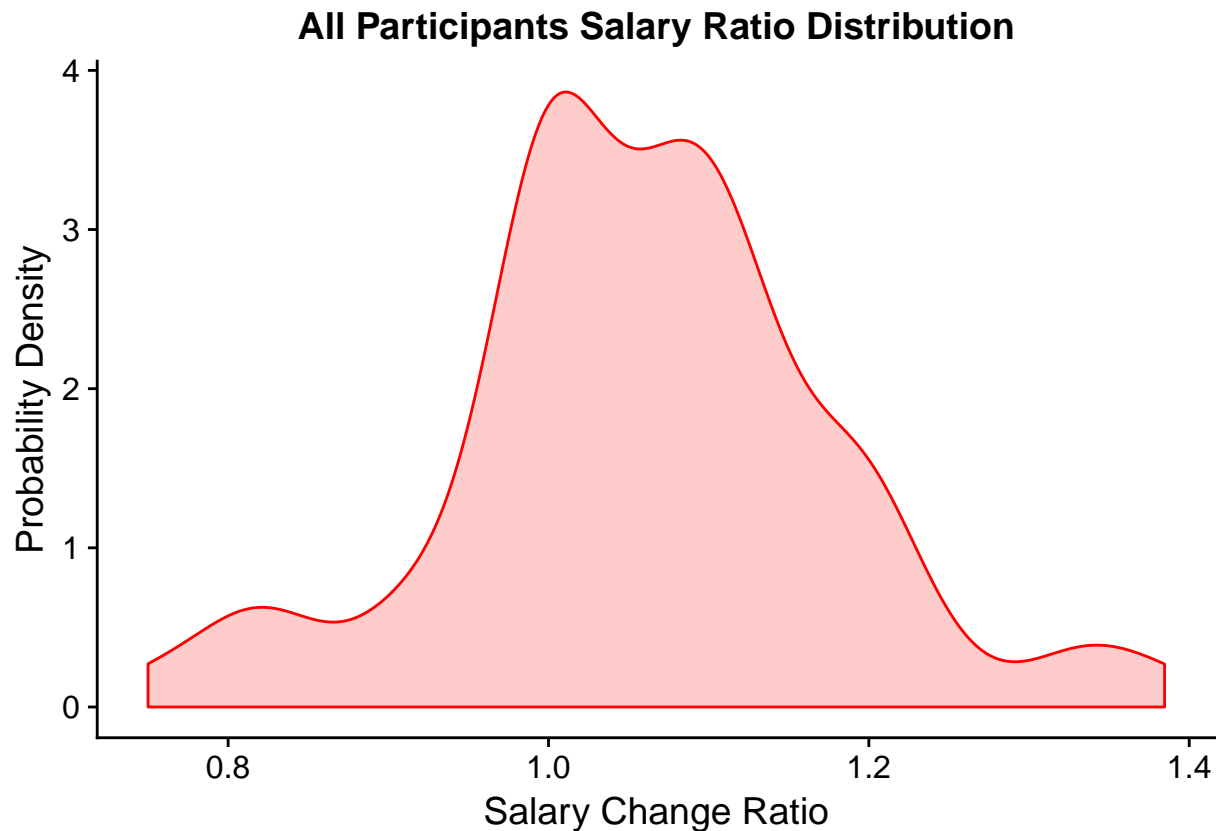
```
## consent          country          salary_base
## Length:71        Length:71        Min.   : 3000
## Class :character  Class :character  1st Qu.: 70000
## Mode  :character  Mode  :character  Median : 80000
##                                     Mean  : 1166042
##                                     3rd Qu.: 117500
##                                     Max.   :65000000
## salary_expect     no_increase_acceptance living_expenses savings
## Min.   : 3000     Length:71        Min.   : 0.00  Min.   : 0.00
## 1st Qu.: 75000     Class :character  1st Qu.:25.00  1st Qu.: 7.00
## Median : 90000     Mode  :character  Median :40.00  Median :10.00
## Mean   : 1265775                    Mean  :38.92  Mean  :14.27
## 3rd Qu.: 120000                    3rd Qu.:50.00  3rd Qu.:20.00
## Max.   :70000000                    Max.   :90.00  Max.   :50.00
## vacation          daily_leisure    consumption_goods sports_hobbies
## Min.   : 0.000     Min.   : 1.00    Min.   : 0.000    Min.   : 0.000
## 1st Qu.: 5.000     1st Qu.: 5.00    1st Qu.: 5.000    1st Qu.: 4.000
## Median :10.000     Median :10.00    Median :10.000    Median : 5.000
## Mean   : 9.394     Mean   :12.45    Mean   : 8.366    Mean   : 6.268
## 3rd Qu.:10.000     3rd Qu.:17.50    3rd Qu.:10.000    3rd Qu.:10.000
## Max.   :30.000     Max.   :60.00    Max.   :30.000    Max.   :25.000
## other             id               ratio
## Min.   : 0.00     Min.   : 2.00    Min.   :0.750
## 1st Qu.: 5.00     1st Qu.:24.00    1st Qu.:1.000
## Median :10.00     Median :45.00    Median :1.062
## Mean   :10.34     Mean   :44.41    Mean   :1.059
## 3rd Qu.:10.00     3rd Qu.:65.50    3rd Qu.:1.125
## Max.   :66.00     Max.   :83.00    Max.   :1.385
```

```
ggplot(data = survey_results, aes(x = country, y = ratio, colour = country, fill = country)) +
  geom_jitter() +
  geom_violin(alpha = 0.2) +
  theme(axis.text.x = element_text(angle = 25, hjust = 0.7, vjust = 0.8))
```



Our assumption seems to be accurate with regards to countries not varying too greatly in their responses. There is no country that has a significantly higher or lower ratio distribution. As a sanity check, it is a good idea to combine survey answers from all participants to verify that the variance around our mean is somewhat normally distributed (the plot above makes it seem intuitive that this would be the case, but cannot make the assumption). This would verify that we are dealing with a t-distribution.

```
ggplot(data = survey_results, aes(x = ratio)) +
  geom_density(colour = 'red', fill = 'red', alpha = 0.2) +
  labs(x = 'Salary Change Ratio', y = 'Probability Density', title = 'All Participants Salary Ratio Dist.
```



Evaluating the Response

The study is interested in the ratio distribution above. Is there any correlation between the above ratio and social standards? The premise of the study was to develop a metric that would indicate the inclination of individuals to see financial gain as the main driver for success and determine if there is a relationship with the way their income is spent. Three variables were collected that pertain to our model's dependent variable which include:

Dependent Features	Description
<code>salary_base</code>	An indicator meant to be a subjective baseline of what salary a person of their expertise would earn.
<code>salary_expect</code>	The expected salary combined with the base salary provides a relative indicator to the respondents pursuit of monetary gains.
<code>no_increase_acceptance</code>	A binary metric serves as a safety check against false positives, that is respondents that may have over-exaggerated their expected salary skewing the impression of interest in monetary gain while in reality being content with their current situation.
<code>ratio</code>	This is a calculated metric that simplifies handling respondent's country selection.

The survey also captured the percentages of the main expenses of each participant. Each participant had to assign percentages that adds up to 100%. The different expense categories were strategically chosen which are believed to relate to a person's social standards. For example, it is believed that a person who spends a large percentage on vacations and daily leisure most likely has higher social standards than a person who contributes most of their salary to savings. The hypothesis is that a person with higher social standards will

have a higher salary ratio as described above.

In theory this makes sense to simply compare these expense percentages to the salary ratios and look for any significant correlation. But in the real world there are many confounders that have to be accounted for. For example, a person who is close to retirement will most likely not expect an increase in the coming year, but may spend a large portion of their salary on vacations and daily leisure.

It isn't always as clear-cut as to say that the closer you are to retirement, the more you will spend on vacation. Or on the other side of the spectrum, it cannot be assumed that a young person won't spend a large percentage of their income on traveling.

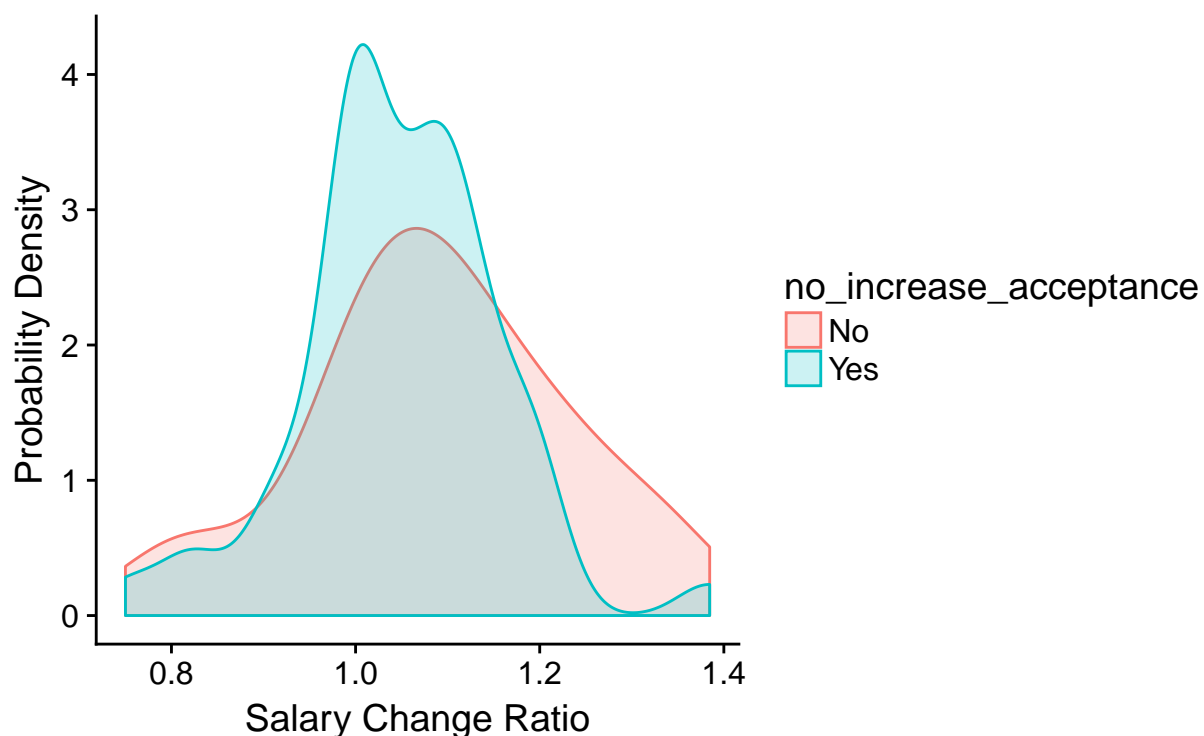
The first confounder that we believe is of importance, is whether a person prefers job satisfaction over an increase in salary. The survey raised the question whether a person would keep their job if they don't receive a salary increase in two years, given high job satisfaction.

A person who spends a lot on vacation and leisure (which can be either the younger or older generation) may strive for a higher salary, but the possibility exists that they don't - possibly depending whether they value job satisfaction over a salary increase.

```
ggplot(survey_results, aes(x = ratio, group = no_increase_acceptance, colour = no_increase_acceptance))
  geom_density(aes(fill = no_increase_acceptance), alpha = 0.2) +
  labs(x = 'Salary Change Ratio', y = 'Probability Density', title = 'All Participants Salary Ratio Dist.
```

All Participants Salary Ratio Distribution

Grouped by Accepted/Declined No-increase in Salary



The plot above shows similar salary ratio distributions for participants who prefer high job satisfaction as those who prefer a salary increase. It does seem as if a person who has a higher salary ratio has a higher probability of preferring an increase over job satisfaction, even though this probability is not significant. However, it will be of more importance if the distributions looked different for people with different types of expenses.

It is difficult to visualize the interaction between expenses, salary ratio and job satisfaction versus salary increase preference. It seems more logical and of statistical importance to fit comparative models and observe

whether the confounder variable adds any value to the model.

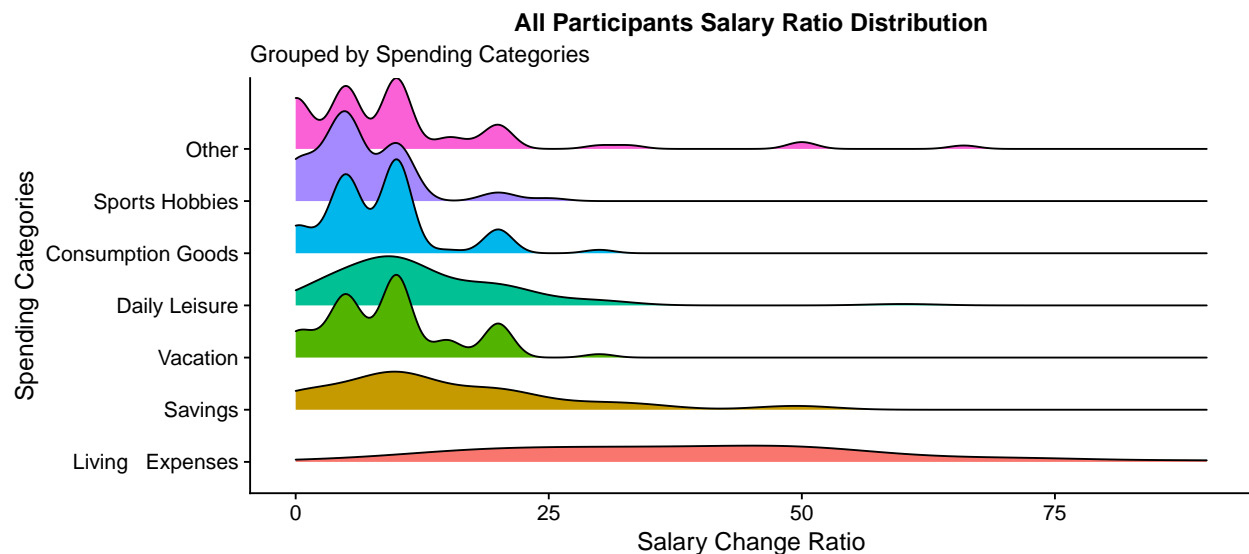
The salary ratio is a continuous variable and from our ratio probability distribution earlier, we saw that the standard deviation is fairly normally distributed around the mean after removing outliers. For this reason a linear regression model seems like a sensible model to fit to our data.

We want to determine whether the preference for job satisfaction interacts with with our explanatory variables. The explanatory variables in our case are the expense categories. We need to compare an additive linear model with a model that considers job satisfaction as a variable that interacts with our expense categories. The following joy plot displays the distribution of participant spendings.

```
# additional wrangling for plotting purposes
survey_results_spendings <- survey_results %>% select(spending_cats)
survey_results_spendings <- map_df(survey_results_spendings, as.numeric)
survey_results_spendings<- melt(survey_results_spendings)

## No id variables; using all as measure variables

# joy plot per participant
ggplot(survey_results_spendings, aes(x = value, y = variable, height = ..density.., fill = variable ))+
  geom_joy(stat = "density")+
  scale_y_discrete(breaks = c("living_expenses", "savings", "vacation", "daily_leisure", "consumption_goods", "sports_hobbies", "other"))+
  theme(legend.position = "None") +
  labs(x = 'Salary Change Ratio', y = 'Spending Categories', title = 'All Participants Salary Ratio Distribution')
```



Additional Modelling and Exploratory Analysis

```
# model without interaction
lm_survey <- lm(ratio ~ living_expenses +
  savings +
  vacation +
  daily_leisure +
  consumption_goods +
  sports_hobbies +
  other, data = survey_results)

summary(lm_survey)
```

```
##
```



```
## Call:
## lm(formula = ratio ~ living_expenses + savings + vacation + daily_leisure +
##      consumption_goods + sports_hobbies + other, data = survey_results)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.32290 -0.05802 -0.00767  0.06133  0.30502
##
## Coefficients: (1 not defined because of singularities)
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    1.0961144  0.1122856   9.762  2.7e-14 ***
## living_expenses -0.0007987  0.0012893  -0.619   0.538
## savings        -0.0008328  0.0016622  -0.501   0.618
## vacation       -0.0017813  0.0024973  -0.713   0.478
## daily_leisure  -0.0003918  0.0020366  -0.192   0.848
## consumption_goods 0.0021812  0.0031719   0.688   0.494
## sports_hobbies   0.0015482  0.0033636   0.460   0.647
## other              NA           NA      NA      NA
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.1209 on 64 degrees of freedom
## Multiple R-squared:  0.05402,    Adjusted R-squared:  -0.03467
## F-statistic: 0.6091 on 6 and 64 DF,  p-value: 0.7221
```

Without any interaction, none of the expenses carry any statistical significance. Below we allow the job satisfaction versus salary increase preference to interact with the expense explanatory variables.

```
# model with interaction
lm_survey <- lm(ratio ~ no_increase_acceptance:(living_expenses +
              savings +
              vacation +
              daily_leisure +
              consumption_goods +
              sports_hobbies +
              other), data = survey_results)

summary(lm_survey)
```

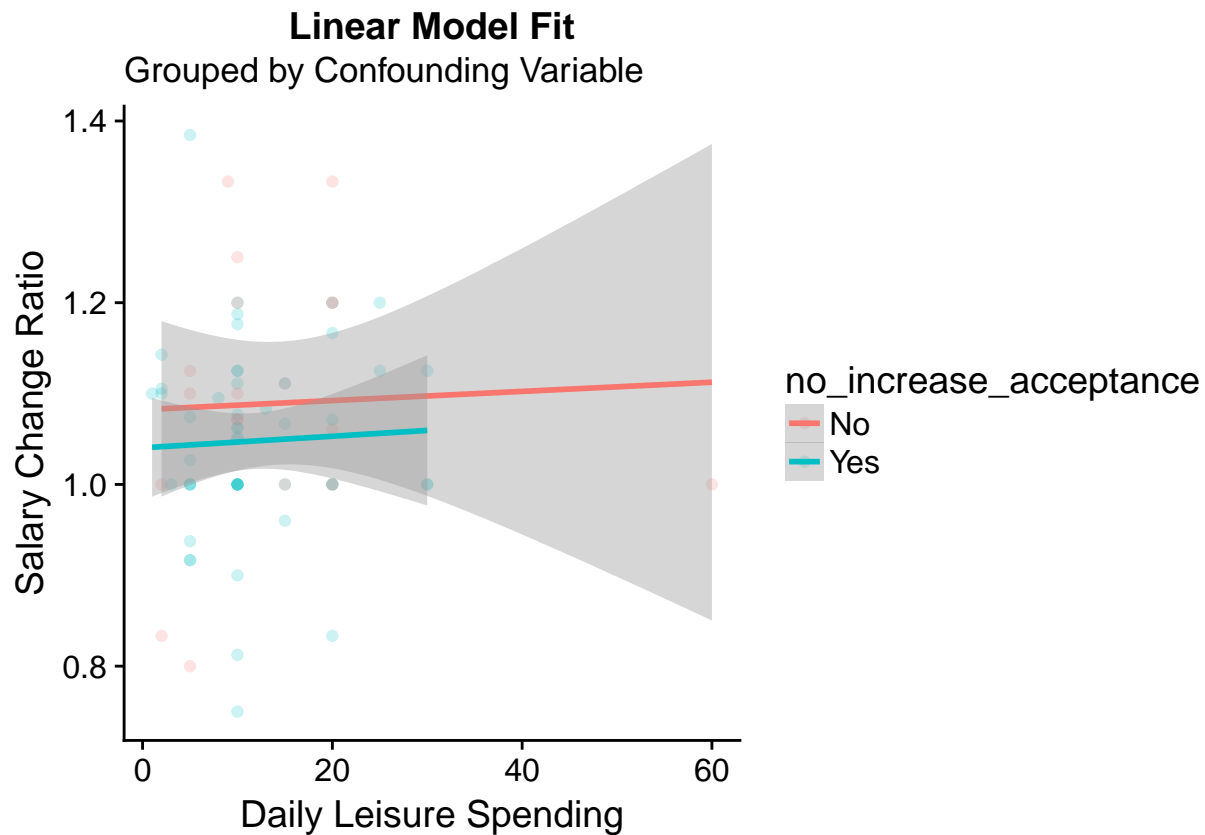
```
##
## Call:
## lm(formula = ratio ~ no_increase_acceptance:(living_expenses +
##      savings + vacation + daily_leisure + consumption_goods +
##      sports_hobbies + other), data = survey_results)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.31909 -0.06406 -0.00684  0.05699  0.28757
##
## Coefficients: (1 not defined because of singularities)
##              Estimate Std. Error t value
## (Intercept)    1.166e+00  1.116e-01  10.450
## no_increase_acceptanceNo:living_expenses -8.975e-04  1.448e-03  -0.620
## no_increase_acceptanceYes:living_expenses -1.526e-03  1.312e-03  -1.163
## no_increase_acceptanceNo:savings        -4.790e-03  4.066e-03  -1.178
```

```
## no_increase_acceptanceYes:savings      -4.291e-04  1.662e-03  -0.258
## no_increase_acceptanceNo:vacation       -1.371e-03  5.880e-03  -0.233
## no_increase_acceptanceYes:vacation      -5.661e-03  2.805e-03  -2.018
## no_increase_acceptanceNo:daily_leisure  -6.113e-03  2.830e-03  -2.160
## no_increase_acceptanceYes:daily_leisure  1.533e-03  2.629e-03   0.583
## no_increase_acceptanceNo:consumption_goods  7.915e-03  4.884e-03   1.621
## no_increase_acceptanceYes:consumption_goods -2.720e-03  4.083e-03  -0.666
## no_increase_acceptanceNo:sports_hobbies   1.477e-02  1.053e-02   1.402
## no_increase_acceptanceYes:sports_hobbies  -4.979e-05  3.474e-03  -0.014
## no_increase_acceptanceNo:other           -8.263e-03  4.640e-03  -1.781
## no_increase_acceptanceYes:other          NA          NA      NA
##                                         Pr(>|t|)
## (Intercept)                          7.32e-15 ***
## no_increase_acceptanceNo:living_expenses  0.5378
## no_increase_acceptanceYes:living_expenses  0.2496
## no_increase_acceptanceNo:savings          0.2436
## no_increase_acceptanceYes:savings         0.7972
## no_increase_acceptanceNo:vacation         0.8164
## no_increase_acceptanceYes:vacation        0.0483 *
## no_increase_acceptanceNo:daily_leisure    0.0350 *
## no_increase_acceptanceYes:daily_leisure   0.5622
## no_increase_acceptanceNo:consumption_goods 0.1106
## no_increase_acceptanceYes:consumption_goods 0.5080
## no_increase_acceptanceNo:sports_hobbies   0.1663
## no_increase_acceptanceYes:sports_hobbies   0.9886
## no_increase_acceptanceNo:other            0.0803 .
## no_increase_acceptanceYes:other          NA
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.1154 on 57 degrees of freedom
## Multiple R-squared:  0.2328, Adjusted R-squared:  0.05782
## F-statistic:  1.33 on 13 and 57 DF,  p-value: 0.2231
```

Above we see that that the job satisfaction confounder variable does contribute towards the correlation between daily leisure, vacation and salary ratio.

Below we visualize daily leisure while accounting for our confounder variable.

```
ggplot(survey_results, aes(y = ratio, x = daily_leisure, group = no_increase_acceptance, colour = no_in
  geom_point(aes(fill = no_increase_acceptance), alpha = 0.2) +
  geom_smooth(method = "lm") +
  labs(x = 'Daily Leisure Spending', y = 'Salary Change Ratio', title = 'Linear Model Fit', subtitle = "G
```

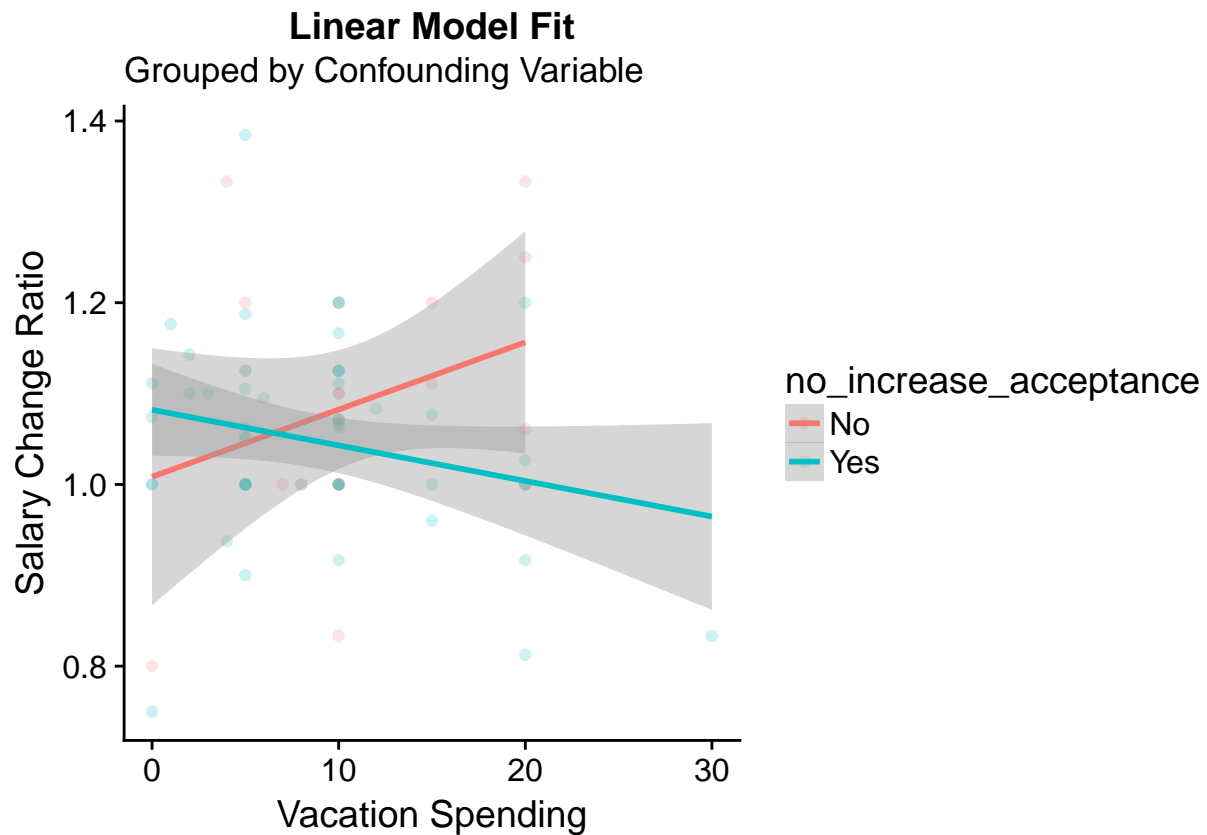


Even though the model found some significance, our visualization seems to disagree to an extent. It might be the daily leisure outlier value that is contributing towards the difference in slopes. The difference in slopes is also quite marginal.

We aren't directly interested in a person's preference between job satisfaction and salary increase, but we do need to take into account how this variable is influencing our study. There are various ways of dealing with confounding variables, but given our dataset size, our options are limited. For now, including this interaction in our model should be sufficient to maintain awareness of its effect. We should also strongly consider removing higher leverage outliers for the different expense categories which may eliminate the effect of the confounding variable, especially in the case above as linear regression model are highly susceptible to outliers.

Below we visualize vacation while taking our confounding variable into account.

```
ggplot(survey_results, aes(y = ratio, x = vacation, group = no_increase_acceptance, colour = no_increas
  geom_point(aes(fill = no_increase_acceptance), alpha = 0.2) +
  geom_smooth(method = "lm") +
  labs(x = 'Vacation Spending', y = 'Salary Change Ratio', title = 'Linear Model Fit', subtitle= "Grouped
```



The difference in slopes is more radical in this case. It would appear that people who spend a larger percentage on vacation have a larger salary ratio **only** if they prefer a salary increase. The confidence intervals are fairly wide, but there might be some truth in the finding. It could contribute towards our hypothesis - people who spend a large percentage on vacation may be the people who are driven by money. In this case, it seems as if our confounding variable interaction could support our hypothesis - people who prefer a salary increase above job satisfaction are those with (possibly) higher social standards (we should be careful to assume that vacation is a direct indication of social standards) and are the same people who expect a higher salary ratio. However, the lack of statistical significance (we aren't yet considering adjusted p-values) and small number of observations mean that we cannot draw any conclusions. However, it is important to differentiate between the people who prefer job satisfaction and those who prefer an increase.