

Exploratory Data Analysis

S. Arora, J. Harmse, V. Mulholland

April 21, 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

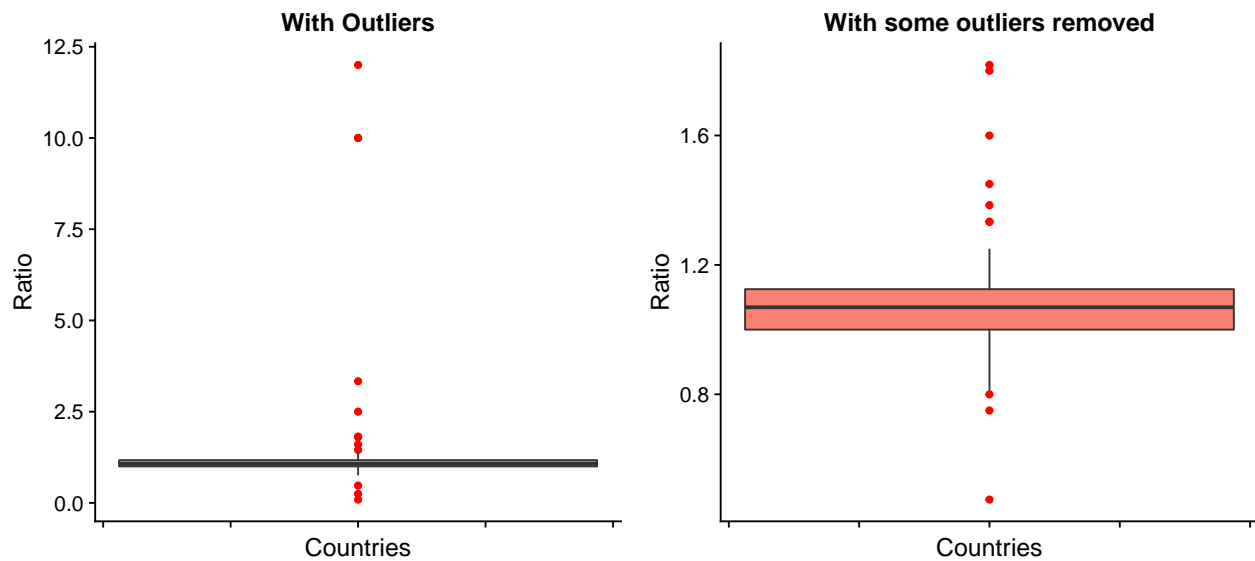


Figure 1: This is a caption

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

```
## saving a copy with all outliers.
survey_results_all_outliers <- survey_results

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

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 removed.

Take a look at Figure 1.

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

The questions were 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.

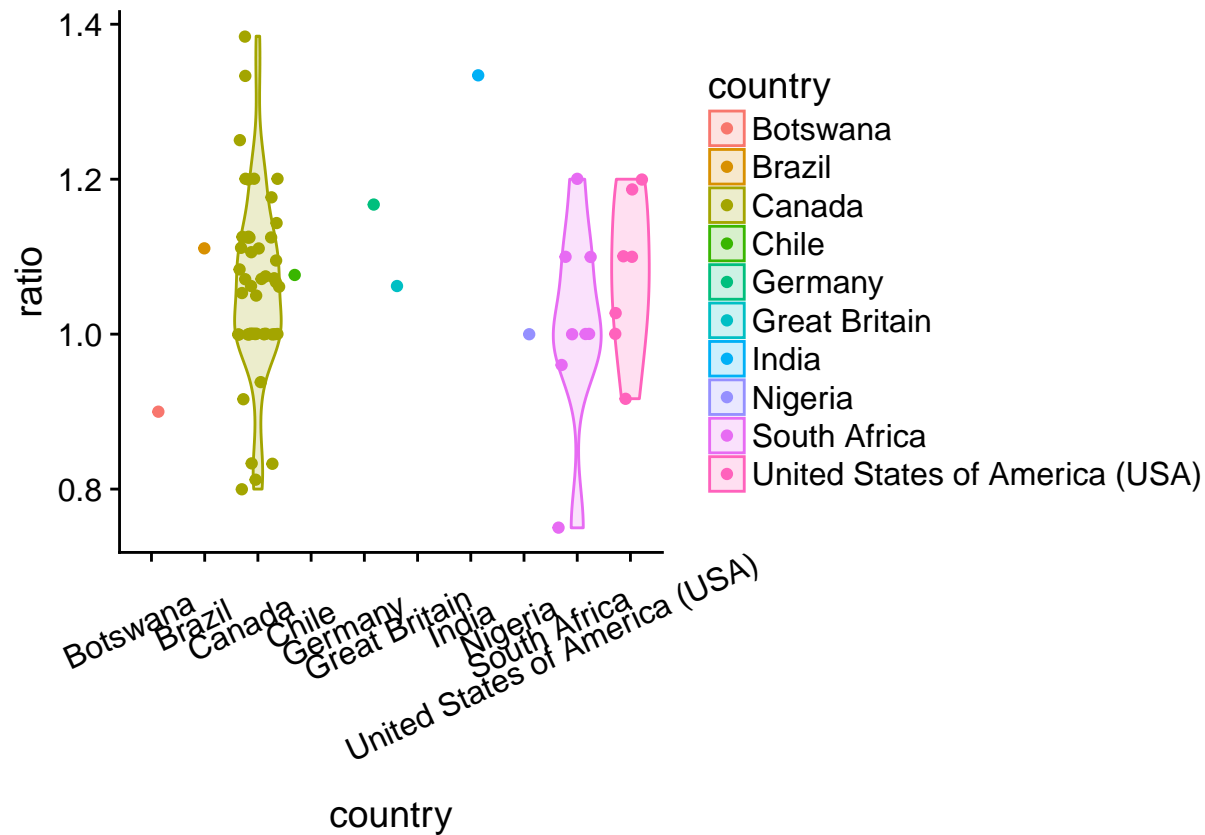
Below each variable is summarized. Since it is difficult to highlight important information from a summary table containing so many variables, a jitter-violin plot was also generated.

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

##	consent	country	salary_base
##	Length:71	Length:71	Min. : 3000

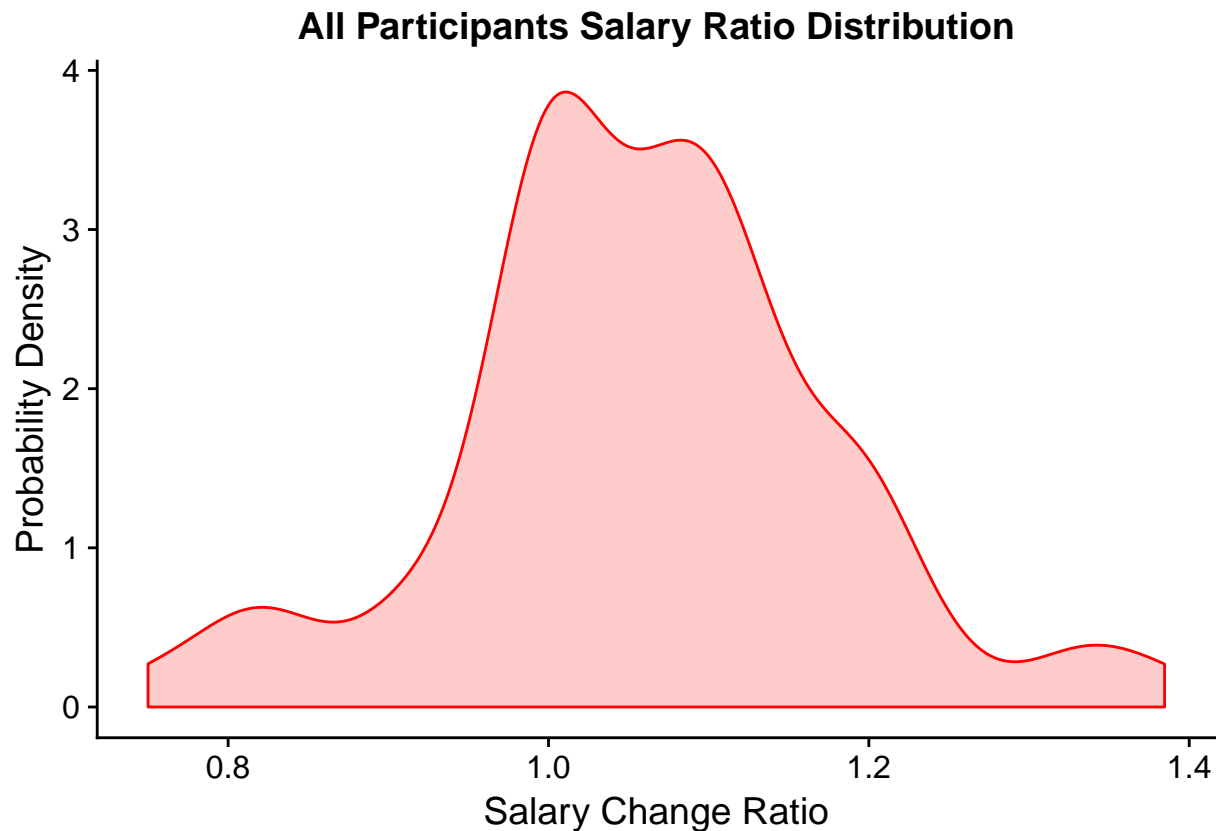
```
## Class :character    Class :character    1st Qu.: 70000
## Mode :character    Mode :character    Median : 80000
##                                     Mean : 1166042
##                                     3rd Qu.: 117500
##                                     Max. :65000000
## salary_expect      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.

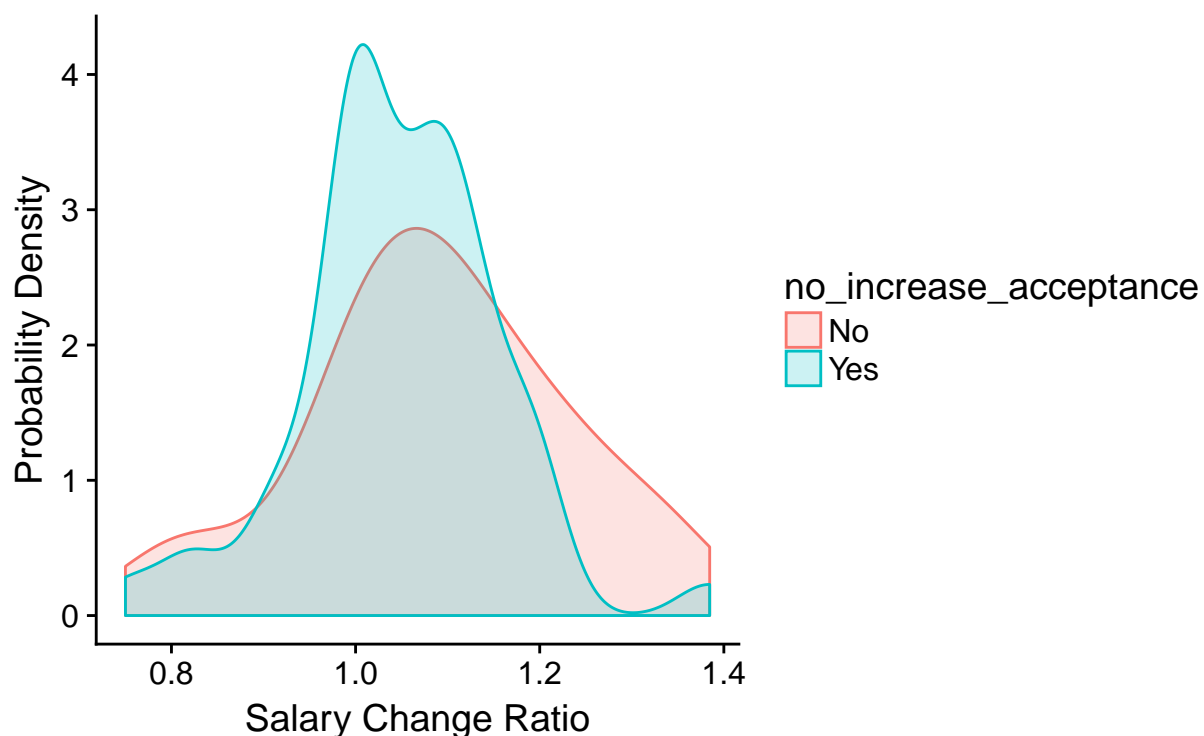
The first confounder that we believe is of importance, is whether a person prefers job satisfaction over an increase in salary. The survey raised the question whether a person would keep their job if they don't receive a salary increase in two years, given high job satisfaction.

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.

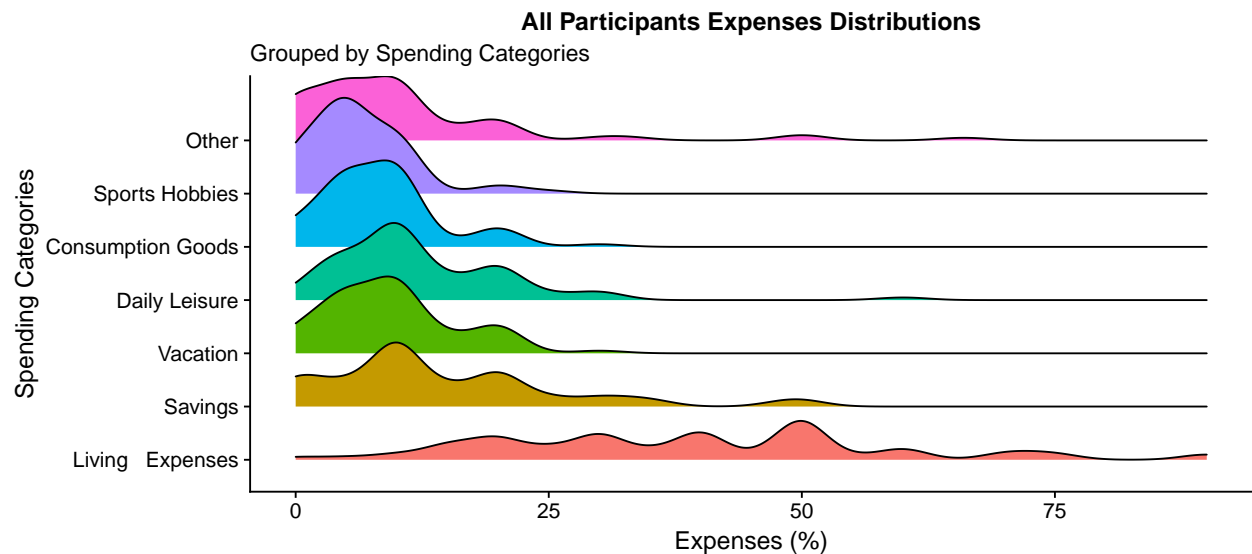
The salary ratio is a continuous variable and from our ratio probability distribution earlier, we saw that the standard deviation is fairly normally distributed around the mean after removing outliers. For this reason a linear regression model seems like a sensible model to fit to our data.

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", bw =2.5)+
  scale_y_discrete(breaks = c("living_expenses", "savings", "vacation", "daily_leisure", "consumption_goods", "sports_hobbies", "other"))+
  theme(legend.position = "None") +
  labs(x = 'Expenses (%)', y = 'Spending Categories', title = 'All Participants Expenses 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)
```

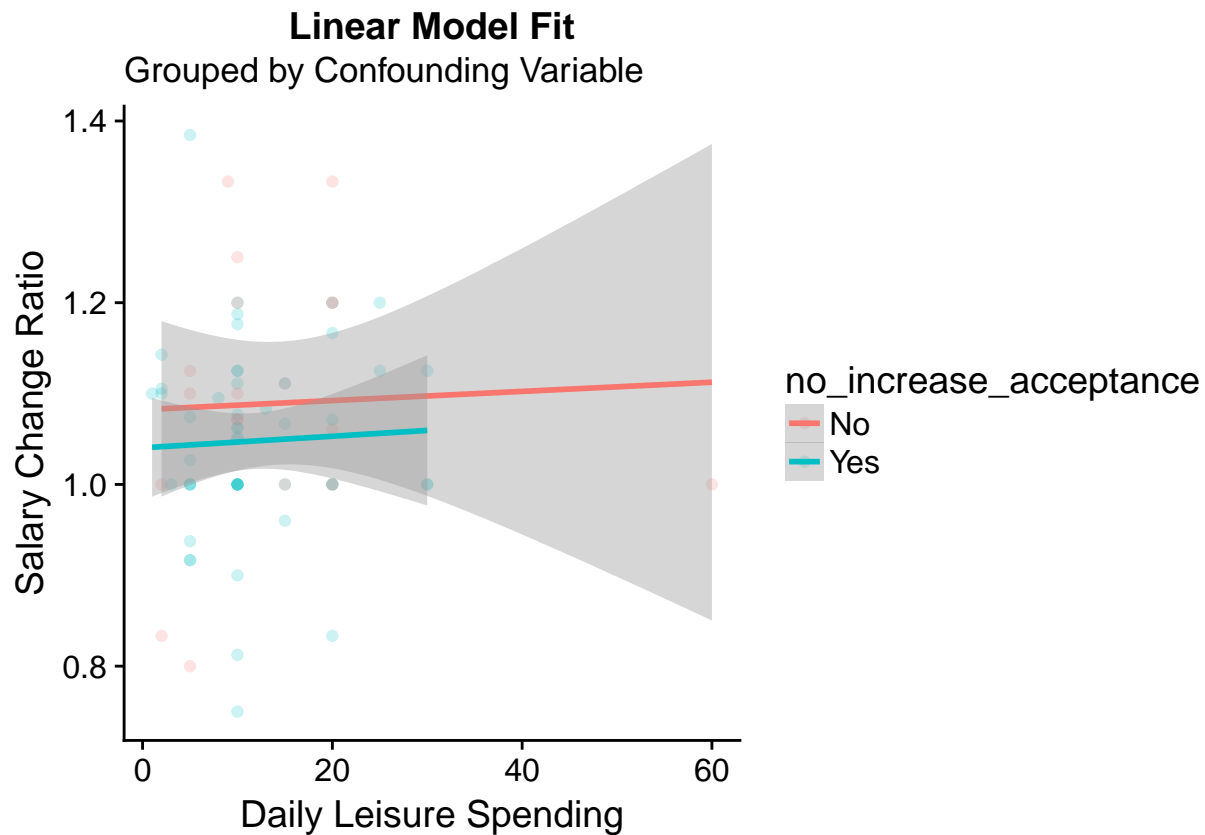
```
##
## Call:
## lm(formula = ratio ~ no_increase_acceptance:(living_expenses +
##      savings + vacation + daily_leisure + consumption_goods +
##      sports_hobbies + other), data = survey_results)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.31909 -0.06406 -0.00684  0.05699  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
```

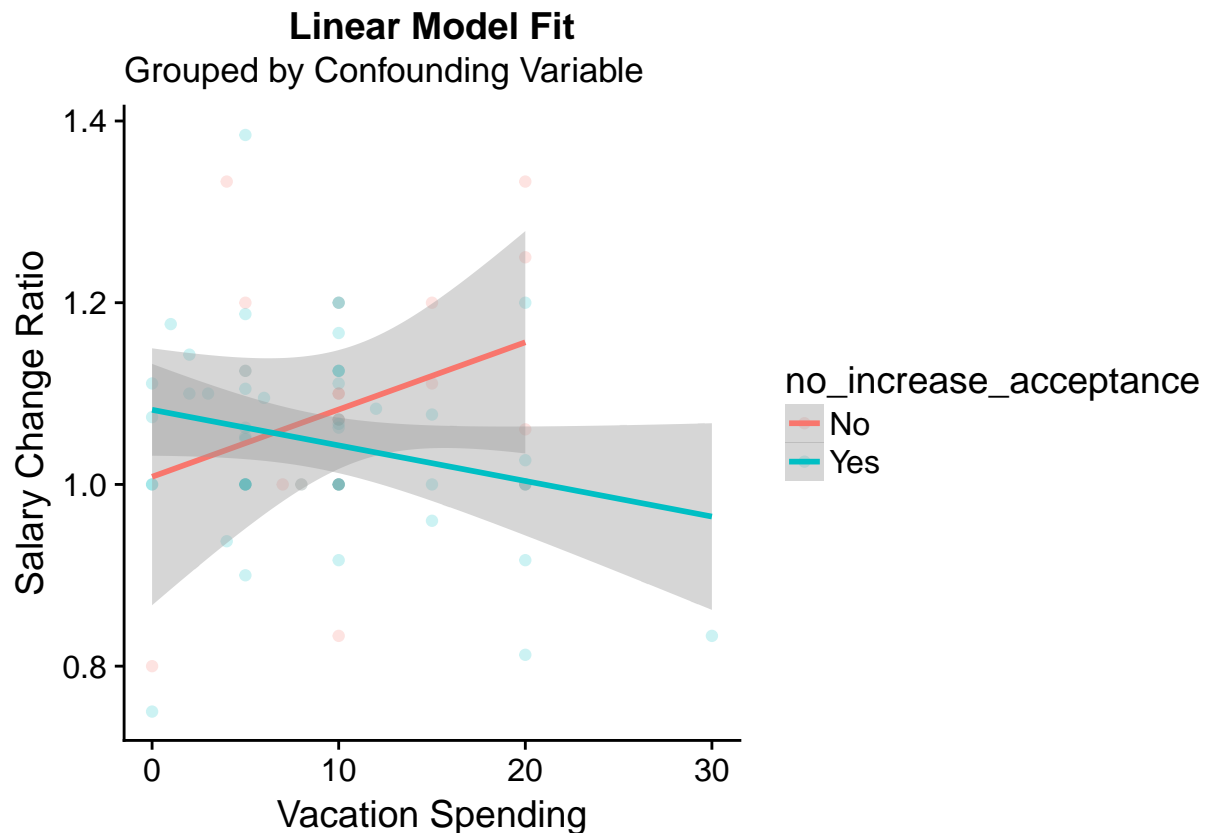


Even though the model found some significance, our visualization seems to disagree to an extent. It might be the daily leisure outlier value that is contributing towards the difference in slopes. The difference in slopes is also quite marginal.

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.

More Analysis

```
# model with interaction
lm_survey <- glm( as.factor(no_increase_acceptance) ~ (living_expenses +
  savings +
  vacation +
  daily_leisure +
  consumption_goods +
  sports_hobbies +
  other),
  family=binomial(link='logit'),
  data = survey_results)
summary(lm_survey)
```

```
##
```

```
## Call:
## glm(formula = as.factor(no_increase_acceptance) ~ (living_expenses +
##   savings + vacation + daily_leisure + consumption_goods +
##   sports_hobbies + other), family = binomial(link = "logit"),
##   data = survey_results)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.0966  -1.0131   0.5881   0.7992   1.3147
##
## Coefficients: (1 not defined because of singularities)
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)      5.03423     3.53545   1.424  0.1545
## living_expenses  -0.04596     0.03833  -1.199  0.2305
## savings           0.01891     0.04874   0.388  0.6981
## vacation         -0.07792     0.05535  -1.408  0.1592
## daily_leisure    -0.02416     0.04574  -0.528  0.5973
## consumption_goods -0.14560     0.07143  -2.038  0.0415 *
## sports_hobbies   -0.02140     0.07346  -0.291  0.7708
## other              NA           NA      NA      NA
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 84.425  on 70  degrees of freedom
## Residual deviance: 74.970  on 64  degrees of freedom
## AIC: 88.97
##
## Number of Fisher Scoring iterations: 5
```

Propensity Score

The independent variable of interest is the acceptance of no increase (`no_increase_acceptance`) and the salary ratio is the dependent variable.

```
# results are standardized
survey_results %>%
  mutate(ratio_std = (ratio - mean(ratio)) / sd(ratio)) %>% # ratio standardization
  group_by(no_increase_acceptance) %>%
  summarise(mean_ratio = mean(ratio_std))
```

```
## # A tibble: 2 x 2
##   no_increase_acceptance mean_ratio
##   <chr>                  <dbl>
## 1 No                    0.252
## 2 Yes                  -0.0989
```

```
# if data is previously standardized
survey_results %>%
  group_by(no_increase_acceptance) %>%
  summarise(n_participants = n(),
            mean_ratio = mean(ratio),
            std_error = sd(ratio) / sqrt(n_participants))
```

```
## # A tibble: 2 x 4
##   no_increase_acceptance n_participants mean_ratio std_error
##   <chr>                  <int>         <dbl>     <dbl>
## 1 No                     20           1.09     0.0314
## 2 Yes                    51           1.05     0.0152

with(survey_results, t.test(ratio ~ no_increase_acceptance))

##
## Welch Two Sample t-test
##
## data:  ratio by no_increase_acceptance
## t = 1.1971, df = 28.383, p-value = 0.2412
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
##  -0.02964387  0.11314154
## sample estimates:
## mean in group No mean in group Yes
##      1.089479      1.047730

cov <- c('living_expenses', 'savings', 'vacation', 'daily_leisure', 'consumption_goods', 'sports_hobbies')
survey_results %>%
  group_by(no_increase_acceptance) %>%
  select(one_of(cov)) %>%
  summarise_all(funs(mean(., na.rm = T)))

## Adding missing grouping variables: `no_increase_acceptance`

## # A tibble: 2 x 8
##   no_increase_acceptance living_expenses savings vacation daily_leisure
##   <chr>                  <dbl>     <dbl>     <dbl>     <dbl>
## 1 No                     39.9      10.6      11.0      14.4
## 2 Yes                    38.5      15.7       8.78      11.7
## # ... with 3 more variables: consumption_goods <dbl>,
## #   sports_hobbies <dbl>, other <dbl>

T-test is used to evaluate if the difference in means is statistically significant.

lapply(cov, function(v) {
  t.test(unlist(survey_results[, v]) ~ unlist(survey_results[, 'no_increase_acceptance']))
})

## [[1]]
##
## Welch Two Sample t-test
##
## data:  unlist(survey_results[, v]) by unlist(survey_results[, "no_increase_acceptance"])
## t = 0.23582, df = 28.046, p-value = 0.8153
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
##  -10.53408  13.27526
## sample estimates:
## mean in group No mean in group Yes
##      39.90000      38.52941
##
##
## [[2]]
```

```
##
## Welch Two Sample t-test
##
## data: unlist(survey_results[, v]) by unlist(survey_results[, "no_increase_acceptance"])
## t = -2.0147, df = 53.586, p-value = 0.04897
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## -10.18783182 -0.02393288
## sample estimates:
## mean in group No mean in group Yes
## 10.60000 15.70588
##
##
## [[3]]
##
## Welch Two Sample t-test
##
## data: unlist(survey_results[, v]) by unlist(survey_results[, "no_increase_acceptance"])
## t = 1.3695, df = 39.274, p-value = 0.1786
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## -1.032166 5.363539
## sample estimates:
## mean in group No mean in group Yes
## 10.950000 8.784314
##
##
## [[4]]
##
## Welch Two Sample t-test
##
## data: unlist(survey_results[, v]) by unlist(survey_results[, "no_increase_acceptance"])
## t = 0.91007, df = 24.465, p-value = 0.3717
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## -3.434415 8.861866
## sample estimates:
## mean in group No mean in group Yes
## 14.40000 11.68627
##
##
## [[5]]
##
## Welch Two Sample t-test
##
## data: unlist(survey_results[, v]) by unlist(survey_results[, "no_increase_acceptance"])
## t = 1.2158, df = 23.999, p-value = 0.2359
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## -1.635032 6.323267
## sample estimates:
## mean in group No mean in group Yes
## 10.050000 7.705882
##
```

```
##
## [[6]]
##
## Welch Two Sample t-test
##
## data:  unlist(survey_results[, v]) by unlist(survey_results[, "no_increase_acceptance"])
## t = 0.28953, df = 53.867, p-value = 0.7733
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
##  -1.916914  2.563973
## sample estimates:
##  mean in group No mean in group Yes
##      6.500000      6.176471
##
##
## [[7]]
##
## Welch Two Sample t-test
##
## data:  unlist(survey_results[, v]) by unlist(survey_results[, "no_increase_acceptance"])
## t = -1.6326, df = 67.49, p-value = 0.1072
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
##  -8.4714196  0.8478902
## sample estimates:
##  mean in group No mean in group Yes
##      7.60000      11.41176

binom_ps <- glm(as.factor(no_increase_acceptance) ~ (living_expenses +
  savings +
  vacation +
  daily_leisure +
  consumption_goods +
  sports_hobbies +
  other),
  family = binomial(), data = survey_results)
summary(binom_ps)

##
## Call:
## glm(formula = as.factor(no_increase_acceptance) ~ (living_expenses +
##   savings + vacation + daily_leisure + consumption_goods +
##   sports_hobbies + other), family = binomial(), data = survey_results)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.0966  -1.0131   0.5881   0.7992   1.3147
##
## Coefficients: (1 not defined because of singularities)
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    5.03423    3.53545   1.424   0.1545
## living_expenses -0.04596    0.03833  -1.199   0.2305
## savings         0.01891    0.04874   0.388   0.6981
## vacation       -0.07792    0.05535  -1.408   0.1592
## daily_leisure  -0.02416    0.04574  -0.528   0.5973
```



```
## consumption_goods -0.14560    0.07143  -2.038    0.0415 *
## sports_hobbies   -0.02140    0.07346  -0.291    0.7708
## other              NA          NA        NA        NA
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##    Null deviance: 84.425  on 70  degrees of freedom
## Residual deviance: 74.970  on 64  degrees of freedom
## AIC: 88.97
##
## Number of Fisher Scoring iterations: 5

lm_survey <- glm( as.factor(no_increase_acceptance) ~ (living_expenses +
  savings +
  vacation +
  daily_leisure +
  consumption_goods +
  sports_hobbies +
  other),
  family=binomial(link='logit'),
  data = survey_results)
summary(lm_survey)

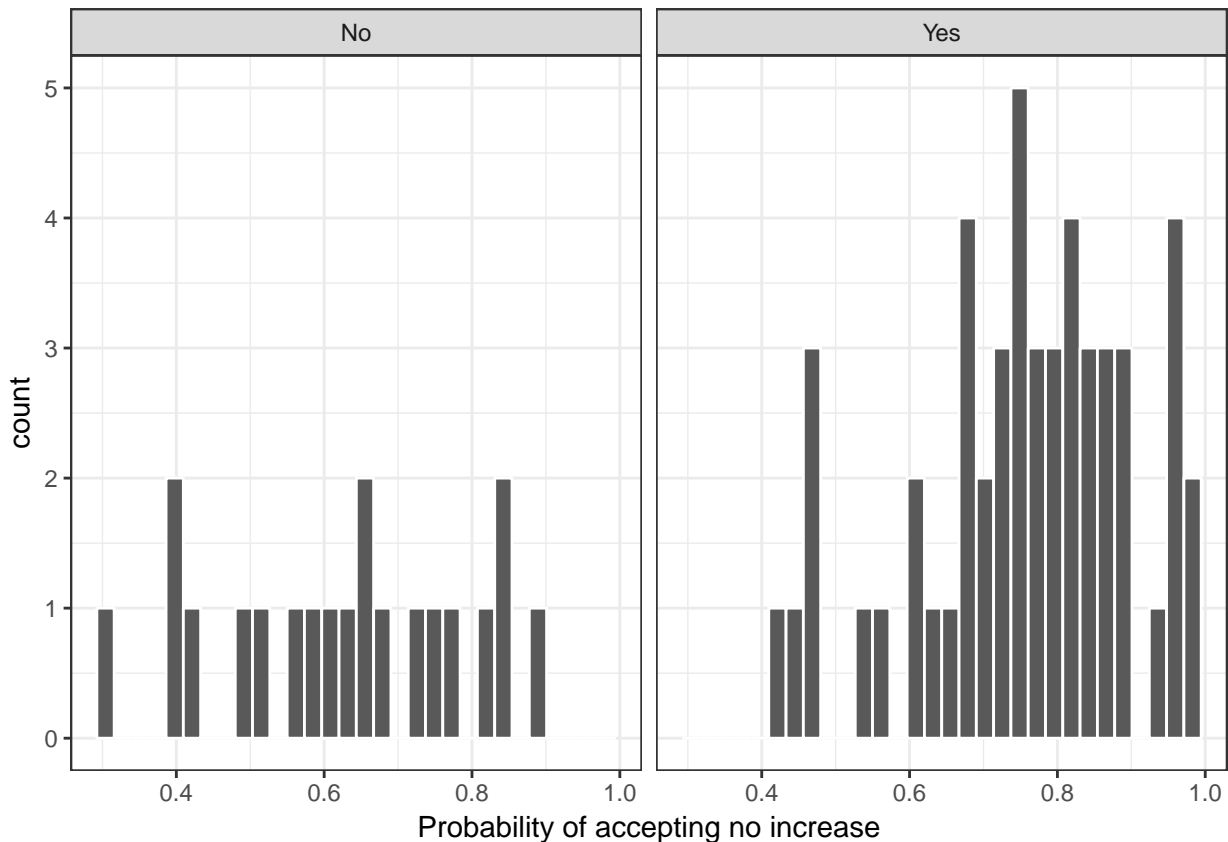
##
## Call:
## glm(formula = as.factor(no_increase_acceptance) ~ (living_expenses +
##   savings + vacation + daily_leisure + consumption_goods +
##   sports_hobbies + other), family = binomial(link = "logit"),
##   data = survey_results)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.0966  -1.0131   0.5881   0.7992   1.3147
##
## Coefficients: (1 not defined because of singularities)
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    5.03423    3.53545   1.424   0.1545
## living_expenses -0.04596    0.03833  -1.199   0.2305
## savings         0.01891    0.04874   0.388   0.6981
## vacation       -0.07792    0.05535  -1.408   0.1592
## daily_leisure  -0.02416    0.04574  -0.528   0.5973
## consumption_goods -0.14560    0.07143  -2.038   0.0415 *
## sports_hobbies  -0.02140    0.07346  -0.291   0.7708
## other              NA          NA        NA        NA
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##    Null deviance: 84.425  on 70  degrees of freedom
## Residual deviance: 74.970  on 64  degrees of freedom
## AIC: 88.97
##
```

```
## Number of Fisher Scoring iterations: 5
prs_df <- data.frame(pr_score = predict(binom_ps, type = "response"),
                     no_increase_accept = binom_ps$model['as.factor(no_increase_acceptance)'])
head(prs_df)

##      pr_score as.factor.no_increase_acceptance.
## 1 0.8335540                                No
## 2 0.7552482                                Yes
## 3 0.8821656                                Yes
## 4 0.9683955                                Yes
## 5 0.7689854                                Yes
## 6 0.8570015                                Yes

labs <- paste("Accepting no salary increase:", c("Yes", "No"))
prs_df %>%
  mutate(no_increase_accept = ifelse(as.factor.no_increase_acceptance. == 1, labs[1], labs[2])) %>%
  ggplot(aes(x = pr_score)) +
  geom_histogram(color = "white") +
  facet_wrap(~as.factor.no_increase_acceptance.) +
  xlab("Probability of accepting no increase") +
  theme_bw()

## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```



```
survey_results

## # A tibble: 71 x 14
##   consent country salary_base salary_expect no_increase_accep~
```

```
##      <chr>      <chr>                <dbl>          <dbl> <chr>
## 1 Yes      Canada                140000.        150000. No
## 2 Yes      Canada                 60000.         65000. Yes
## 3 Yes      South Africa           550000.        550000. Yes
## 4 Yes      United States of ~     100000.        110000. Yes
## 5 Yes      South Africa           160000.        120000. Yes
## 6 Yes      Canada                 30000.         30000. Yes
## 7 Yes      South Africa           3000000.       3300000. Yes
## 8 Yes      United States of ~      75000.         77000. Yes
## 9 Yes      South Africa           500000.        480000. Yes
## 10 Yes     Canada                100000.        120000. Yes
## # ... with 61 more rows, and 9 more variables: living_expenses <int>,
## #   savings <dbl>, vacation <int>, daily_leisure <int>,
## #   consumption_goods <dbl>, sports_hobbies <dbl>, other <dbl>, id <int>,
## #   ratio <dbl>

survey_results_nomiss <- survey_results %>%
  select(ratio, no_increase_acceptance, one_of(cov)) %>%
  na.omit() %>% mutate(no_increase=if_else(no_increase_acceptance == "Yes", 1, 0))
library(MatchIt)
mod_match <- matchit(no_increase ~ living_expenses + savings + vacation + daily_leisure + consumption_goods,
  method = "nearest", data = survey_results_nomiss)

## Warning in matchit2nearest(structure(c(0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, :
## Fewer control than treated units and matching without replacement. Not all
## treated units will receive a match. Treated units will be matched in the
## order specified by m.order: largest

summary(mod_match)

##
## Call:
## matchit(formula = no_increase ~ living_expenses + savings + vacation +
##   daily_leisure + consumption_goods + sports_hobbies + other,
##   data = survey_results_nomiss, method = "nearest", ratio = 1)
##
## Summary of balance for all data:
##           Means Treated Means Control SD Control Mean Diff eQQ Med
## distance           0.7537           0.6281           0.1675           0.1256 0.1247
## living_expenses     38.5294           39.9000           23.4676          -1.3706 5.0000
## savings             15.7059           10.6000            8.1331           5.1059 5.0000
## vacation             8.7843           10.9500            5.7626          -2.1657 3.5000
## daily_leisure       11.6863           14.4000           12.4959          -2.7137 0.0000
## consumption_goods    7.7059           10.0500            8.1206          -2.3441 0.0000
## sports_hobbies        6.1765            6.5000            3.5762          -0.3235 0.0000
## other              11.4118            7.6000            6.1078           3.8118 0.0000
##
##           eQQ Mean eQQ Max
## distance           0.1207 0.1762
## living_expenses     3.8000 15.0000
## savings             4.9000 20.0000
## vacation            3.0000 10.0000
## daily_leisure       3.5000 30.0000
## consumption_goods    2.7000 10.0000
## sports_hobbies       2.5500 15.0000
## other               4.5000 46.0000
```

```
##
##
## Summary of balance for matched data:
##           Means Treated Means Control SD Control Mean Diff eQQ Med
## distance           0.8887           0.6281           0.1675           0.2607 0.2289
## living_expenses     30.6500           39.9000           23.4676          -9.2500 10.0000
## savings             24.5000           10.6000            8.1331          13.9000 12.5000
## vacation            7.3000           10.9500            5.7626          -3.6500  4.0000
## daily_leisure       8.5500           14.4000          12.4959          -5.8500  4.5000
## consumption_goods   6.3500           10.0500            8.1206          -3.7000  1.5000
## sports_hobbies      4.8000            6.5000            3.5762          -1.7000  0.0000
## other              17.8500            7.6000            6.1078          10.2500  5.0000
##           eQQ Mean eQQ Max
## distance           0.2607  0.5093
## living_expenses    11.0500 30.0000
## savings            13.9000 28.0000
## vacation           3.6500 10.0000
## daily_leisure      5.8500 35.0000
## consumption_goods  4.1000 20.0000
## sports_hobbies     1.7000  5.0000
## other             10.4500 46.0000
##
## Percent Balance Improvement:
##           Mean Diff.   eQQ Med   eQQ Mean   eQQ Max
## distance      -107.4685  -83.5707 -115.8871 -188.9601
## living_expenses -574.8927 -100.0000 -190.7895 -100.0000
## savings        -172.2350 -150.0000 -183.6735  -40.0000
## vacation       -68.5378  -14.2857  -21.6667   0.0000
## daily_leisure  -115.5708    -Inf   -67.1429  -16.6667
## consumption_goods -57.8419    -Inf   -51.8519 -100.0000
## sports_hobbies  -425.4545   0.0000   33.3333   66.6667
## other          -168.9043    -Inf  -132.2222   0.0000
##
## Sample sizes:
##           Control Treated
## All           20      51
## Matched       20      20
## Unmatched      0      31
## Discarded      0       0

# trying with ratio of 5 control cases to one treatment
mod_match <- matchit(no_increase ~ living_expenses + savings + vacation + daily_leisure + consumption_g
                      method = "nearest", data = survey_results_nomiss)

## Warning in matchit2nearest(structure(c(0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, :
## Not enough control units for 5 matches for each treated unit when matching
## without replacement. Not all treated units will receive 5 matches

summary(mod_match)

##
## Call:
## matchit(formula = no_increase ~ living_expenses + savings + vacation +
##         daily_leisure + consumption_goods + sports_hobbies + other,
##         data = survey_results_nomiss, method = "nearest", ratio = 5)
```

```

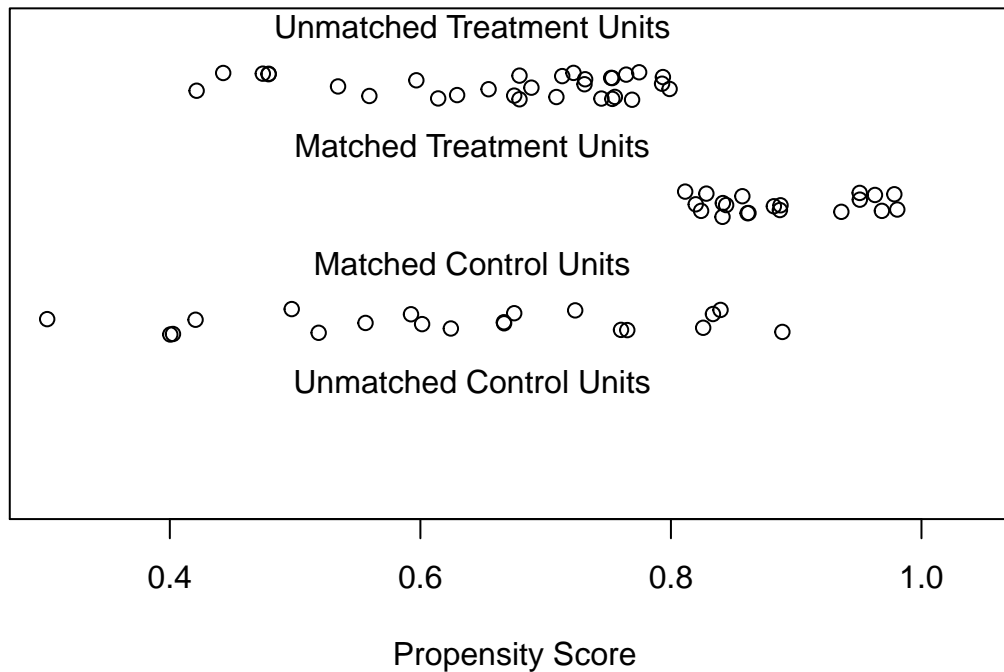
##
## Summary of balance for all data:
##           Means Treated Means Control SD Control Mean Diff eQQ Med
## distance           0.7537           0.6281           0.1675           0.1256  0.1247
## living_expenses    38.5294          39.9000          23.4676          -1.3706  5.0000
## savings            15.7059          10.6000           8.1331           5.1059  5.0000
## vacation           8.7843          10.9500           5.7626          -2.1657  3.5000
## daily_leisure      11.6863          14.4000          12.4959          -2.7137  0.0000
## consumption_goods   7.7059          10.0500           8.1206          -2.3441  0.0000
## sports_hobbies      6.1765           6.5000           3.5762          -0.3235  0.0000
## other              11.4118           7.6000           6.1078           3.8118  0.0000
##           eQQ Mean eQQ Max
## distance           0.1207  0.1762
## living_expenses     3.8000 15.0000
## savings             4.9000 20.0000
## vacation            3.0000 10.0000
## daily_leisure       3.5000 30.0000
## consumption_goods    2.7000 10.0000
## sports_hobbies       2.5500 15.0000
## other               4.5000 46.0000
##
##
## Summary of balance for matched data:
##           Means Treated Means Control SD Control Mean Diff eQQ Med
## distance           0.8887           0.6281           0.1675           0.2607  0.2289
## living_expenses    30.6500          39.9000          23.4676          -9.2500 10.0000
## savings            24.5000          10.6000           8.1331          13.9000 12.5000
## vacation           7.3000          10.9500           5.7626          -3.6500  4.0000
## daily_leisure      8.5500          14.4000          12.4959          -5.8500  4.5000
## consumption_goods   6.3500          10.0500           8.1206          -3.7000  1.5000
## sports_hobbies      4.8000           6.5000           3.5762          -1.7000  0.0000
## other              17.8500           7.6000           6.1078          10.2500  5.0000
##           eQQ Mean eQQ Max
## distance           0.2607  0.5093
## living_expenses    11.0500 30.0000
## savings            13.9000 28.0000
## vacation            3.6500 10.0000
## daily_leisure       5.8500 35.0000
## consumption_goods   4.1000 20.0000
## sports_hobbies       1.7000  5.0000
## other              10.4500 46.0000
##
## Percent Balance Improvement:
##           Mean Diff.   eQQ Med   eQQ Mean   eQQ Max
## distance          -107.4685  -83.5707 -115.8871 -188.9601
## living_expenses   -574.8927 -100.0000 -190.7895 -100.0000
## savings           -172.2350 -150.0000 -183.6735  -40.0000
## vacation          -68.5378  -14.2857  -21.6667   0.0000
## daily_leisure    -115.5708    -Inf    -67.1429  -16.6667
## consumption_goods  -57.8419    -Inf   -51.8519 -100.0000
## sports_hobbies    -425.4545   0.0000   33.3333   66.6667
## other            -168.9043    -Inf  -132.2222   0.0000
##
## Sample sizes:

```

```
##           Control Treated
## All           20      51
## Matched       20      20
## Unmatched      0      31
## Discarded      0       0
```

```
plot(mod_match, type = "jitter")
```

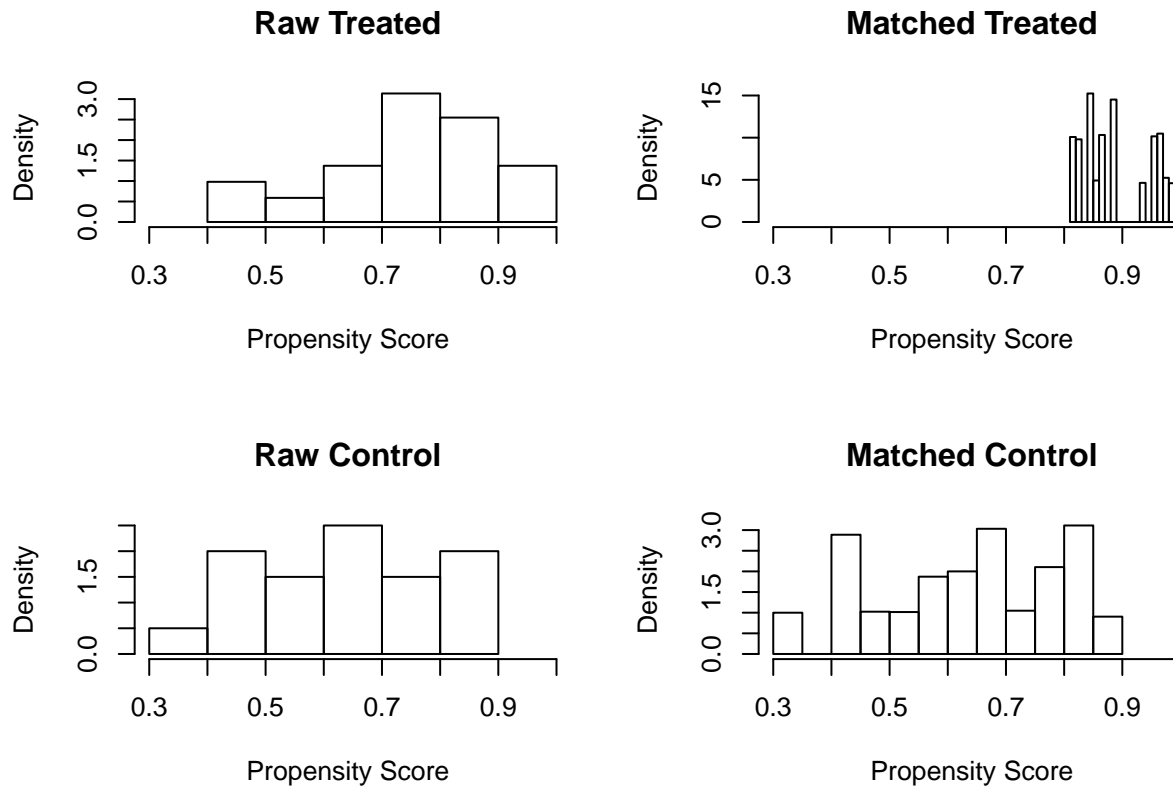
Distribution of Propensity Scores



```
## [1] "To identify the units, use first mouse button; to stop, use second."
```

```
## integer(0)
```

```
plot(mod_match, type = "hist")
```



```
dta_m <- match.data(mod_match)
dim(dta_m)
```

```
## [1] 40 12
```

```
fn_bal <- function(dta, variable) {
  dta$variable <- dta[, variable]
  # if (variable == 'w3income') dta$variable <- dta$variable / 10^3
  # dta$catholic <- as.factor(dta$catholic)
  support <- c(min(dta$variable), max(dta$variable))
  ggplot(dta, aes(x = distance, y = variable, color = no_increase_acceptance)) +
    geom_point(alpha = 0.2, size = 1.3) +
    geom_smooth(method = "loess", se = F) +
    xlab("Propensity score") +
    ylab(variable) +
    theme_bw() +
    ylim(support)
}
```

```
library(gridExtra)
```

```
##
```

```
## Attaching package: 'gridExtra'
```

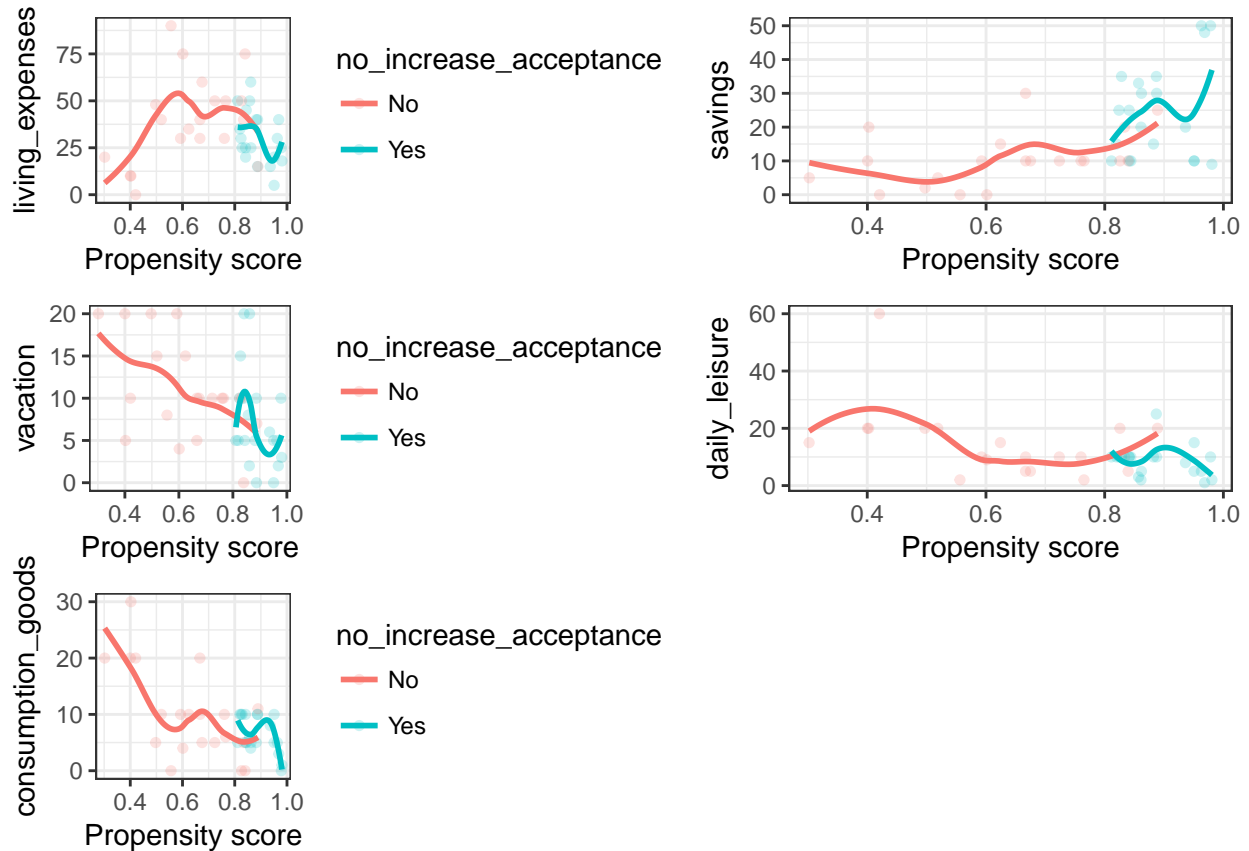
```
## The following object is masked from 'package:dplyr':
```

```
##
```

```
## combine
```

```
grid.arrange(
  fn_bal(dta_m, "living_expenses"),
  fn_bal(dta_m, "savings") + theme(legend.position = "none"),
```

```
fn_bal(dta_m, "vacation"),
fn_bal(dta_m, "daily_leisure") + theme(legend.position = "none"),
fn_bal(dta_m, "consumption_goods"),
nrow = 3, widths = c(1, 0.8)
)
```



Our grouping variable is `no_increase_acceptance` and our matching variables are all the spending categories (`living_expenses`, `savings`, `vacation`, `daily_leisure`, `consumption_goods`, `sports_hobbies`, `other`) where we aim to equalize the groups on.

```
# trying with ratio of 5 control cases to one treatment
mod_match <- matchit(no_increase ~ living_expenses + savings + vacation + daily_leisure + consumption_g
                      method = "subclass", data = survey_results_nomiss)

summary(mod_match)
```

```
## Warning: Not enough control units in subclass 4
## Warning: Not enough control units in subclass 6
## Warning: Not enough control units in subclass 4
## Warning: Not enough control units in subclass 6
## Warning: Not enough control units in subclass 4
## Warning: Not enough control units in subclass 6
## Warning: Not enough control units in subclass 4
## Warning: Not enough control units in subclass 6
```



```

## Warning: Not enough control units in subclass 4
## Warning: Not enough control units in subclass 6
## Warning: Not enough control units in subclass 4
## Warning: Not enough control units in subclass 6
## Warning: Not enough control units in subclass 4
## Warning: Not enough control units in subclass 6
## Warning: Not enough control units in subclass 4
## Warning: Not enough control units in subclass 6

##
## Call:
## matchit(formula = no_increase ~ living_expenses + savings + vacation +
##         daily_leisure + consumption_goods + sports_hobbies + other,
##         data = survey_results_nomiss, method = "subclass", sub.by = "treat")
## Summary of balance for all data:
##               Means Treated Means Control Mean Diff eQQ Med eQQ Mean
## distance           0.7537         0.6281    0.1256  0.1247  0.1207
## living_expenses    38.5294        39.9000   -1.3706  5.0000  3.8000
## savings            15.7059        10.6000    5.1059  5.0000  4.9000
## vacation           8.7843        10.9500   -2.1657  3.5000  3.0000
## daily_leisure     11.6863        14.4000   -2.7137  0.0000  3.5000
## consumption_goods  7.7059        10.0500   -2.3441  0.0000  2.7000
## sports_hobbies     6.1765         6.5000   -0.3235  0.0000  2.5500
## other             11.4118         7.6000    3.8118  0.0000  4.5000
##               eQQ Max
## distance           0.1762
## living_expenses    15.0000
## savings            20.0000
## vacation           10.0000
## daily_leisure     30.0000
## consumption_goods  10.0000
## sports_hobbies     15.0000
## other             46.0000
##
##
## Summary of balance by subclasses:
## , , Subclass 1
##
##               Means Treated Means Control Mean Diff eQQ Med eQQ Mean
## distance           0.5112         0.4768    0.0343  0.0184  0.0384
## living_expenses    47.2222        35.8889   11.3333 10.0000 12.4444
## savings            6.2222         5.7778    0.4444  0.0000  0.6667
## vacation           11.6667        13.5556   -1.8889  3.0000  4.1111
## daily_leisure     13.3333        19.5556   -6.2222  4.0000  6.8889
## consumption_goods  10.8889        13.2222   -2.3333  0.0000  3.0000
## sports_hobbies     5.5556         6.4444   -0.8889  3.0000  3.1111
## other             5.1111         5.5556   -0.4444  0.0000  2.6667
##               eQQ Max
## distance           0.1193
## living_expenses    30.0000
## savings            5.0000

```

```

## vacation          10.0000
## daily_leisure     35.0000
## consumption_goods 10.0000
## sports_hobbies    10.0000
## other              10.0000
##
## , , Subclass 2
##
##               Means Treated Means Control Mean Diff eQQ Med eQQ Mean
## distance          0.6784         0.6582    0.0202   0.0103   0.0160
## living_expenses   40.6250        41.2500   -0.6250   7.5000   7.5000
## savings           10.0000        16.2500  -6.2500   7.5000   7.5000
## vacation          11.8750        10.0000   1.8750   0.0000   1.2500
## daily_leisure     15.0000         8.7500   6.2500   5.0000   6.2500
## consumption_goods  8.1250        11.2500  -3.1250   2.5000   2.5000
## sports_hobbies     6.3750         7.5000  -1.1250   2.0000   2.2500
## other              8.0000         5.0000   3.0000   2.5000   2.5000
##               eQQ Max
## distance          0.0383
## living_expenses   15.0000
## savings           10.0000
## vacation          5.0000
## daily_leisure     15.0000
## consumption_goods  5.0000
## sports_hobbies     5.0000
## other              5.0000
##
## , , Subclass 3
##
##               Means Treated Means Control Mean Diff eQQ Med eQQ Mean
## distance          0.7429         0.7419   0.0010   0.0031   0.0031
## living_expenses   47.7500        40.0000   7.7500  11.0000  11.0000
## savings           10.0000        10.0000   0.0000  10.0000  10.0000
## vacation          7.7500        10.0000  -2.2500   7.5000   7.5000
## daily_leisure     14.3750        10.0000   4.3750  14.0000  14.0000
## consumption_goods  6.0000         7.5000  -1.5000   2.5000   2.5000
## sports_hobbies     6.6250        10.0000  -3.3750  12.5000  12.5000
## other              7.5000        12.5000  -5.0000   5.0000   5.0000
##               eQQ Max
## distance          0.0049
## living_expenses   22.0000
## savings           10.0000
## vacation          10.0000
## daily_leisure     20.0000
## consumption_goods  5.0000
## sports_hobbies    15.0000
## other              5.0000
##
## , , Subclass 4
##
##               Means Treated Means Control Mean Diff eQQ Med eQQ Mean
## distance          0.7943         0.7652
## living_expenses   37.2222        50.0000
## savings           16.1111        10.0000

```

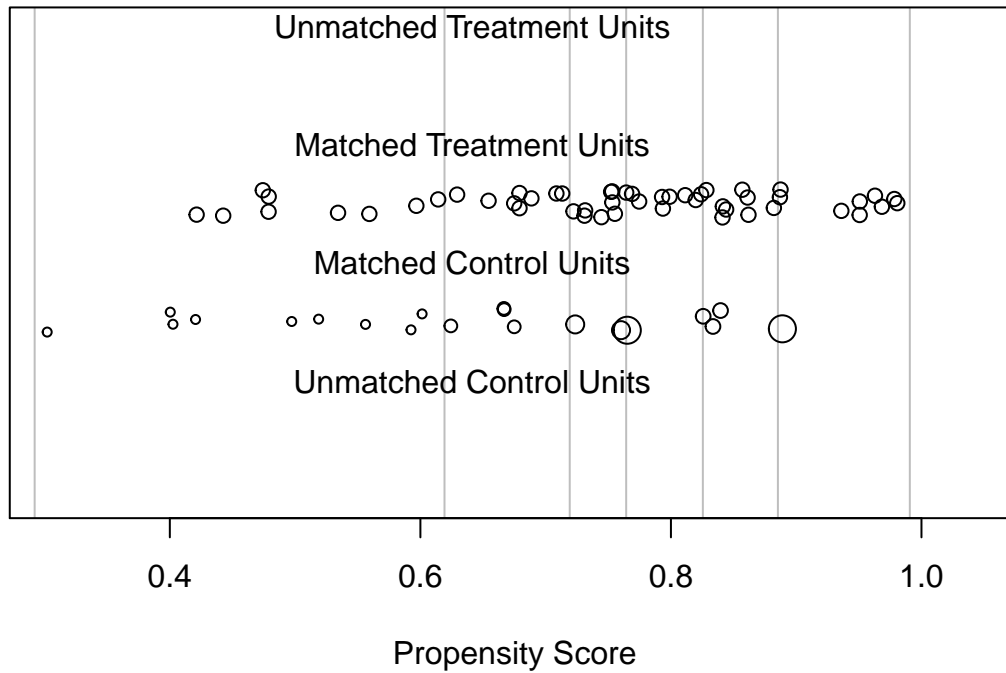
```

## vacation          6.6667      10.0000
## daily_leisure     11.1111       2.0000
## consumption_goods  8.8889       6.0000
## sports_hobbies     8.8889       2.0000
## other             11.1111      20.0000
##                  eQQ Max
## distance
## living_expenses
## savings
## vacation
## daily_leisure
## consumption_goods
## sports_hobbies
## other
##
## , , Subclass 5
##
##               Means Treated Means Control Mean Diff  eQQ Med eQQ Mean
## distance          0.8522      0.8330    0.0192   0.0105   0.0185
## living_expenses    36.2500     55.0000   -18.7500  20.0000  20.0000
## savings            22.2500     13.3333    8.9167  10.0000   8.3333
## vacation          10.6250      6.6667    3.9583   2.0000   4.6667
## daily_leisure       7.5000     11.6667   -4.1667   3.0000   4.3333
## consumption_goods   6.2500      1.6667    4.5833   5.0000   4.6667
## sports_hobbies      5.2500      5.0000    0.2500   0.0000   0.0000
## other             11.8750      6.6667    5.2083   5.0000  10.0000
##                  eQQ Max
## distance          0.0426
## living_expenses    25.0000
## savings            15.0000
## vacation          10.0000
## daily_leisure     10.0000
## consumption_goods   5.0000
## sports_hobbies      0.0000
## other             20.0000
##
## , , Subclass 6
##
##               Means Treated Means Control Mean Diff  eQQ Med eQQ Mean
## distance          0.9447      0.8890    0.0557   0.0000   0.0000
## living_expenses    23.1111     15.0000    8.1111  20.0000  20.0000
## savings            29.1111     25.0000    4.1111  20.0000  20.0000
## vacation          4.5556      7.0000   -2.4444   0.0000   0.0000
## daily_leisure       9.0000     20.0000  -11.0000   0.0000   0.0000
## consumption_goods   5.7778     11.0000  -5.2222   0.0000   0.0000
## sports_hobbies      4.3333      5.0000  -0.6667   0.0000   0.0000
## other             24.1111     17.0000    7.1111  20.0000  20.0000
##                  eQQ Max
## distance
## living_expenses
## savings
## vacation
## daily_leisure
## consumption_goods

```

```
## sports_hobbies
## other
##
##
## Sample sizes by subclasses:
##      Subclass 1 Subclass 2 Subclass 3 Subclass 4 Subclass 5 Subclass 6
## Treated      9      8      8      9      8      9
## Control      9      4      2      1      3      1
## Total       18     12     10     10     11     10
##
## Summary of balance across subclasses
##      Means Treated Means Control Mean Diff eQQ Med eQQ Mean
## distance      0.7537      0.7263      NA      NA      NA
## living_expenses 38.5294     39.1765      NA      NA      NA
## savings       15.7059     13.4052      NA      NA      NA
## vacation       8.7843      9.5752      NA      NA      NA
## daily_leisure  11.6863     12.1046      NA      NA      NA
## consumption_goods 7.7059      8.5359      NA      NA      NA
## sports_hobbies  6.1765      5.9020      NA      NA      NA
## other        11.4118     11.3007      NA      NA      NA
##      eQQ Max
## distance      NA
## living_expenses NA
## savings        NA
## vacation        NA
## daily_leisure  NA
## consumption_goods NA
## sports_hobbies  NA
## other          NA
##
## Percent Balance Improvement:
##      Mean Diff. eQQ Med eQQ Mean eQQ Max
## distance      78.2151      NA      NA      NA
## living_expenses 52.7897      NA      NA      NA
## savings        54.9411      NA      NA      NA
## vacation       63.4827      NA      NA      NA
## daily_leisure  84.5857      NA      NA      NA
## consumption_goods 64.5894      NA      NA      NA
## sports_hobbies 15.1515      NA      NA      NA
## other         97.0850      NA      NA      NA
plot(mod_match, type = "jitter")
```

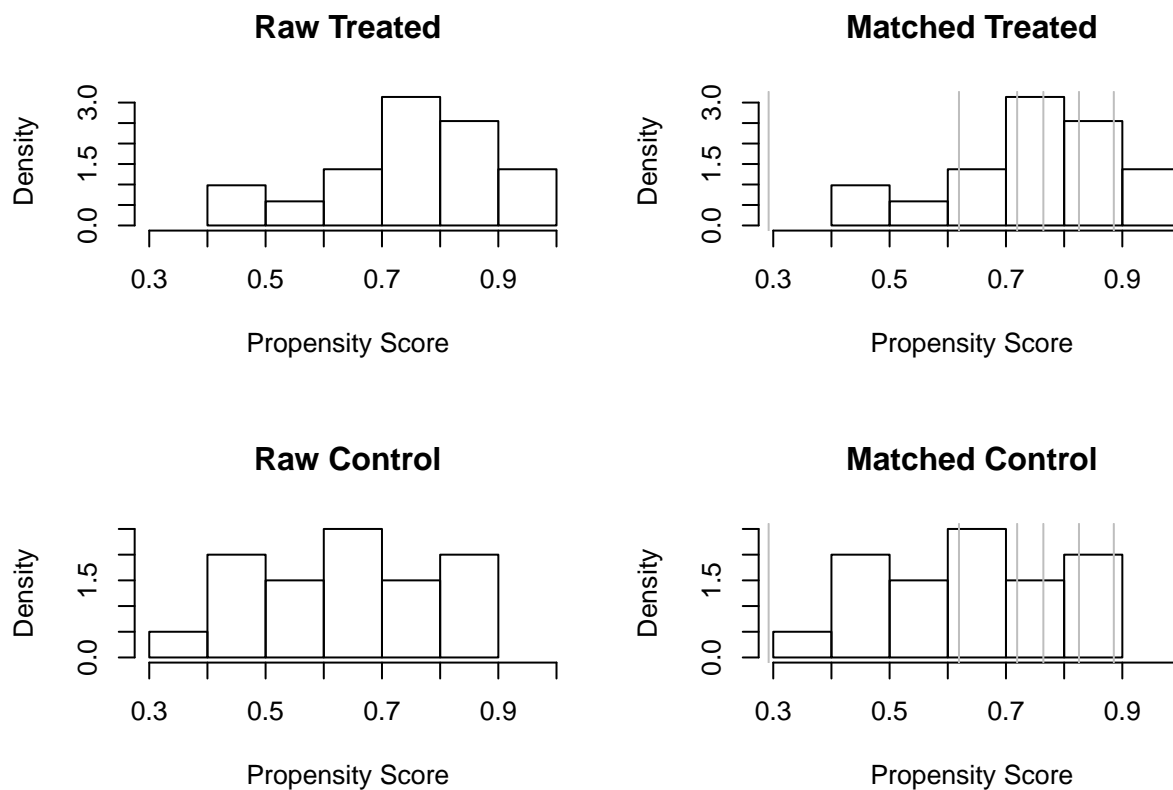
Distribution of Propensity Scores



```
## [1] "To identify the units, use first mouse button; to stop, use second."
```

```
## integer(0)
```

```
plot(mod_match, type = "hist")
```

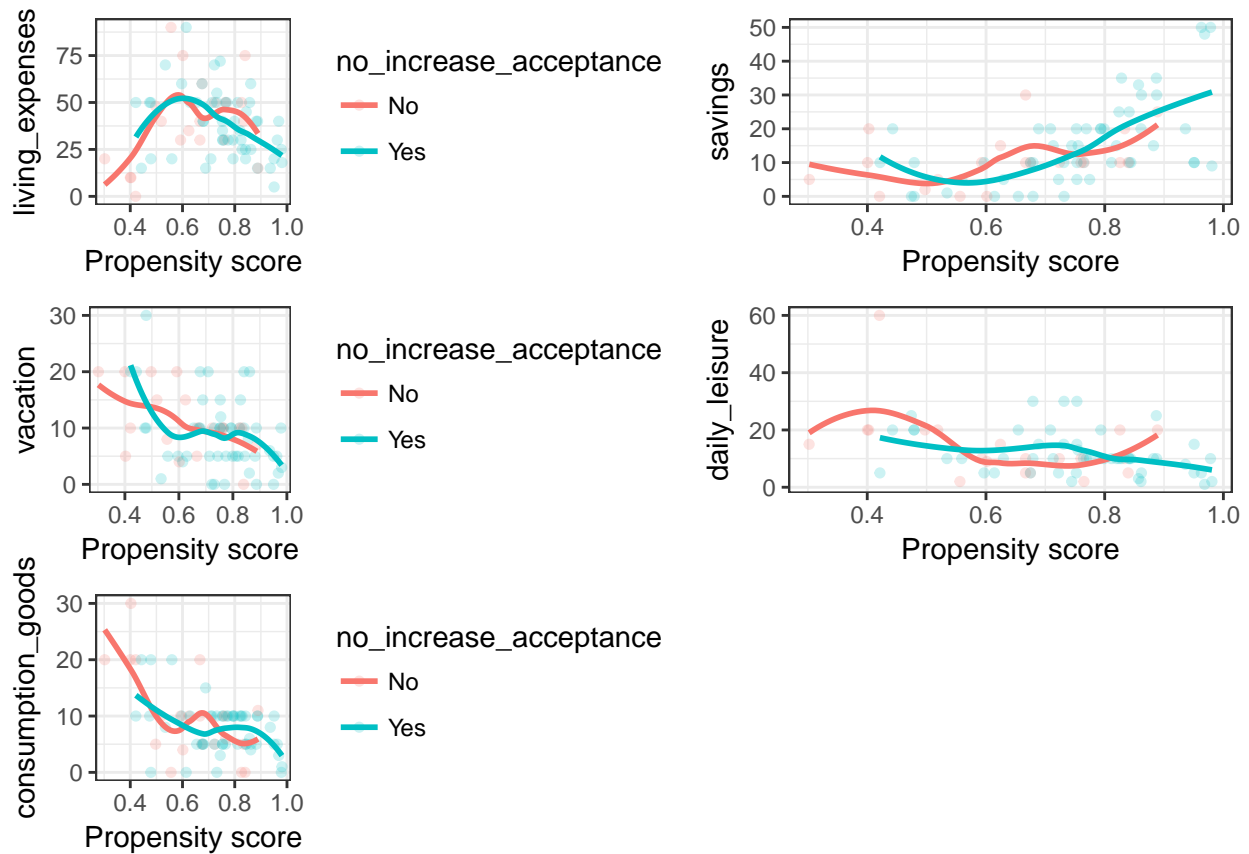


```
dta_m <- match.data(mod_match)
dim(dta_m)

## [1] 71 13

fn_bal <- function(dta, variable) {
  dta$variable <- dta[, variable]
  # if (variable == 'w3income') dta$variable <- dta$variable / 10^3
  # dta$catholic <- as.factor(dta$catholic)
  support <- c(min(dta$variable), max(dta$variable))
  ggplot(dta, aes(x = distance, y = variable, color = no_increase_acceptance)) +
    geom_point(alpha = 0.2, size = 1.3) +
    geom_smooth(method = "loess", se = F) +
    xlab("Propensity score") +
    ylab(variable) +
    theme_bw() +
    ylim(support)
}

library(gridExtra)
grid.arrange(
  fn_bal(dta_m, "living_expenses"),
  fn_bal(dta_m, "savings") + theme(legend.position = "none"),
  fn_bal(dta_m, "vacation"),
  fn_bal(dta_m, "daily_leisure") + theme(legend.position = "none"),
  fn_bal(dta_m, "consumption_goods"),
  nrow = 3, widths = c(1, 0.8)
)
```



```
# trying with ratio of 5 control cases to one treatment
mod_match <- matchit(no_increase ~ living_expenses + savings + vacation + daily_leisure + consumption_g
                      method = "cem", data = survey_results_nomiss)
```

```
##
## Using 'treat'='1' as baseline group
```

```
summary(mod_match)
```

```
##
## Call:
## matchit(formula = no_increase ~ living_expenses + savings + vacation +
##       daily_leisure + consumption_goods + sports_hobbies + other,
##       data = survey_results_nomiss, method = "cem")
##
## Summary of balance for all data:
```

	Means Treated	Means Control	SD Control	Mean Diff	eQQ Med
distance	0.7537	0.6281	0.1675	0.1256	0.1247
living_expenses	38.5294	39.9000	23.4676	-1.3706	5.0000
savings	15.7059	10.6000	8.1331	5.1059	5.0000
vacation	8.7843	10.9500	5.7626	-2.1657	3.5000
daily_leisure	11.6863	14.4000	12.4959	-2.7137	0.0000
consumption_goods	7.7059	10.0500	8.1206	-2.3441	0.0000
sports_hobbies	6.1765	6.5000	3.5762	-0.3235	0.0000
other	11.4118	7.6000	6.1078	3.8118	0.0000

```
##
##           eQQ Mean eQQ Max
## distance      0.1207  0.1762
## living_expenses 3.8000 15.0000
```

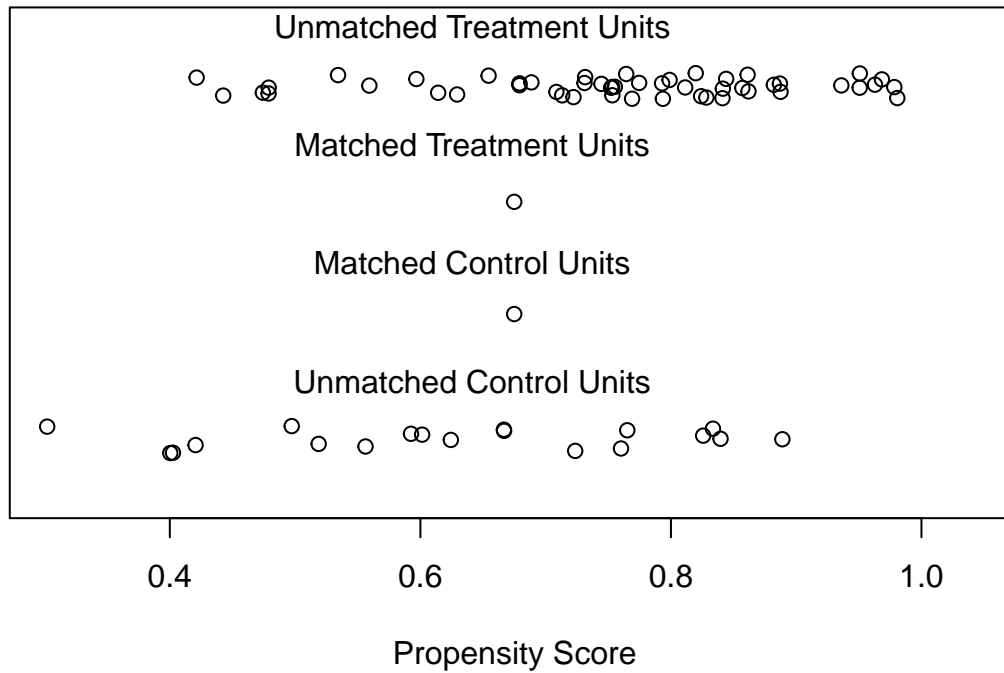
```

## savings          4.9000 20.0000
## vacation         3.0000 10.0000
## daily_leisure    3.5000 30.0000
## consumption_goods 2.7000 10.0000
## sports_hobbies   2.5500 15.0000
## other            4.5000 46.0000
##
##
## Summary of balance for matched data:
##           Means Treated Means Control SD Control Mean Diff eQQ Med
## distance          0.6749      0.6749    NaN          0      0
## living_expenses    60.0000     60.0000    NaN          0      0
## savings            10.0000     10.0000    NaN          0      0
## vacation           10.0000     10.0000    NaN          0      0
## daily_leisure       5.0000      5.0000    NaN          0      0
## consumption_goods    5.0000      5.0000    NaN          0      0
## sports_hobbies       5.0000      5.0000    NaN          0      0
## other              5.0000      5.0000    NaN          0      0
##           eQQ Mean eQQ Max
## distance          0      0
## living_expenses    0      0
## savings            0      0
## vacation           0      0
## daily_leisure      0      0
## consumption_goods   0      0
## sports_hobbies      0      0
## other              0      0
##
## Percent Balance Improvement:
##           Mean Diff. eQQ Med eQQ Mean eQQ Max
## distance          100    100    100    100
## living_expenses    100    100    100    100
## savings            100    100    100    100
## vacation           100    100    100    100
## daily_leisure      100      0    100    100
## consumption_goods   100      0    100    100
## sports_hobbies      100      0    100    100
## other              100      0    100    100
##
## Sample sizes:
##           Control Treated
## All           20      51
## Matched        1       1
## Unmatched      19     50
## Discarded       0       0

```

```
plot(mod_match, type = "jitter")
```

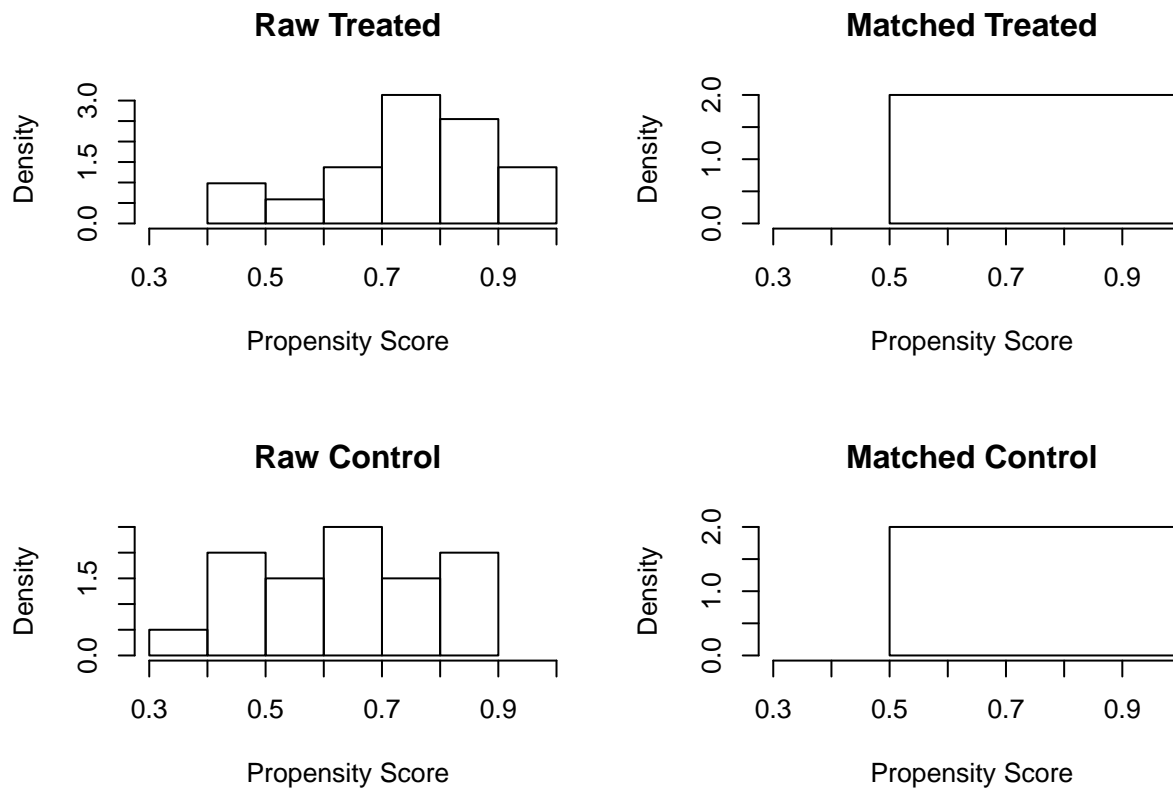

Distribution of Propensity Scores



```
## [1] "To identify the units, use first mouse button; to stop, use second."
```

```
## integer(0)
```

```
plot(mod_match, type = "hist")
```



```

dta_m <- match.data(mod_match)
dim(dta_m)

## [1] 2 13

fn_bal <- function(dta, variable) {
  dta$variable <- dta[, variable]
  # if (variable == 'w3income') dta$variable <- dta$variable / 10^3
  # dta$catholic <- as.factor(dta$catholic)
  support <- c(min(dta$variable), max(dta$variable))
  ggplot(dta, aes(x = distance, y = variable, color = no_increase_acceptance)) +
    geom_point(alpha = 0.2, size = 1.3) +
    geom_smooth(method = "loess", se = F) +
    xlab("Propensity score") +
    ylab(variable) +
    theme_bw() +
    ylim(support)
}

library(gridExtra)
grid.arrange(
  fn_bal(dta_m, "living_expenses"),
  fn_bal(dta_m, "savings") + theme(legend.position = "none"),
  fn_bal(dta_m, "vacation"),
  fn_bal(dta_m, "daily_leisure") + theme(legend.position = "none"),
  fn_bal(dta_m, "consumption_goods"),
  nrow = 3, widths = c(1, 0.8)
)

```

