# The `ispiranteRanney` Package

*Steven H. Ranney*

## Contents

> *The best experiments require no statistics.*
> *— Al Zale, USGS/Montana Cooperative Fishery Research Unit*

## Introduction

Data analysis is normally undertaken by those with training in statistical analyses (Burnham and Anderson 2002, Anderson 2008). Statistical analyses, and their resulting metrics (e.g., $p$-values, slope nad intercept values, parameter estimates, $R^2$ values) are often misinterpreted by those that posess none or little statistical training. In some cases, reliance on the output from stastical analyses may "cloud" a person's judgement; lack of training may let someone infer from a *significant $p$*-value that a relationship exists when, in reality, the results were a result of spurious data.

I propose to deliver to ispirante the *ispiranteRanney* pckage, written in response to an onine advertisement and subsequent webinar with Stefano Spada. *ispiranteRanney* will have several functions to help ispirante customers visualize, understand, and gain insights into their data. I will write the package—and the functions—with minimal statistical analyses to allow for a deeper, more thorough customer understanding.

Most of the functions provided in *ispiranteRanney* will be simple plotting functions; for example wait time by Day, resolution time by date, satisfaction score by `assignee_name`, etc. Data insights are often generated through simple visualization and help-desk ticket data is no exception. However, the real insight lies in comparing some values to others. In this respect, *ispiranteRanney* will provide the means to make those

visual comparisons. Coupled with simple statistical analyses, users will be able explore how satisfcation score changes as a function of resolution or wait time, for example.

The `ispiranteRanney` package will be written in the most stable release R version 3.1.2 "Pumpkin Helmet" and tested for stability on the development version of R, 3.2.0.

### *A note about R and plots*

This document was produced with the `knitr` package. There will be a combination of text and R data. R data will look like this:

```r
#Create 1000 random values from a normal distribution with a mean of 10 and a SD of 0.8
tmp <- as.data.frame(rnorm(1000, 10, .8))
names(tmp) <- "value"
```

Lines in the R code that begin with a `#` are comments from the author. The output from the R code is preceeded by `##`. Plain text will look like this. Plots will appear in their own window. In some cases plots may not fit at the bottom of a page. In those cases, `knitr` will move the plot to the next page and the bottom of the precedding page may be blank, or nearly so.
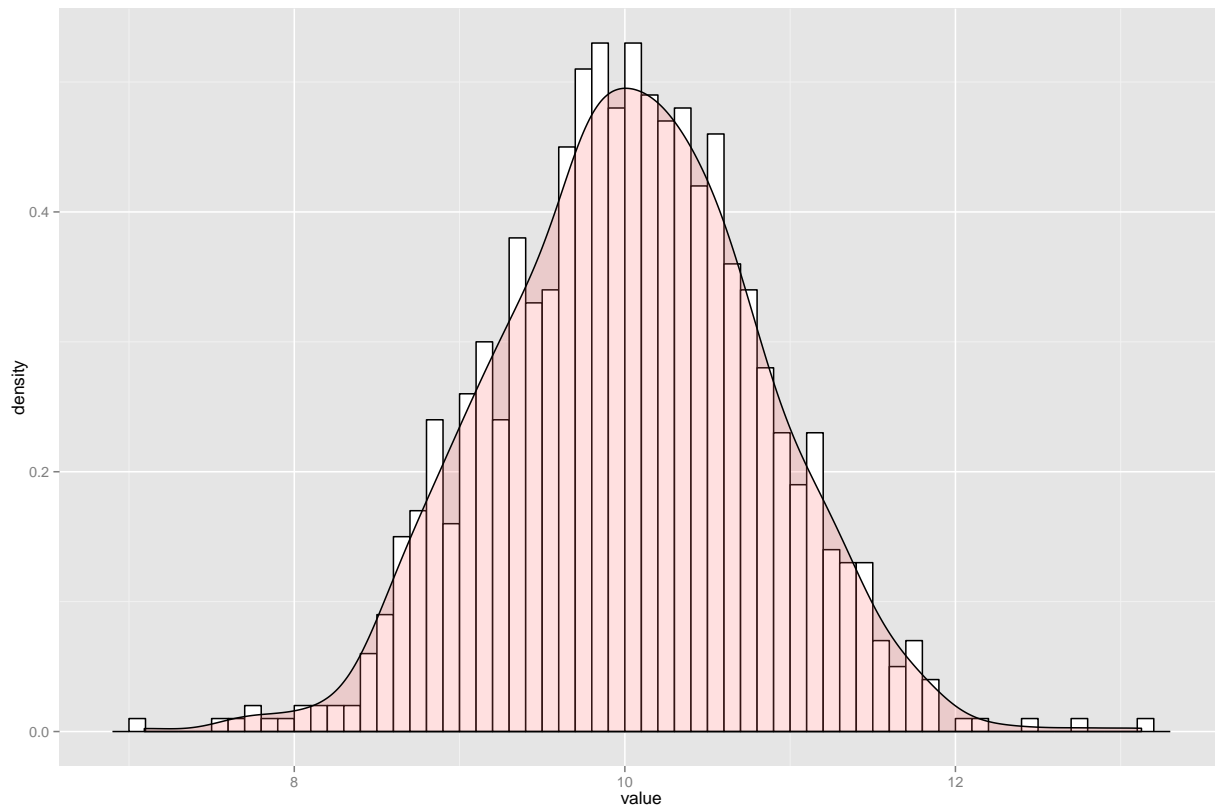


Figure 1: Histogram and kernel density of 1,000 random numbers with a mean = 10 and a standard deviation = 0.8
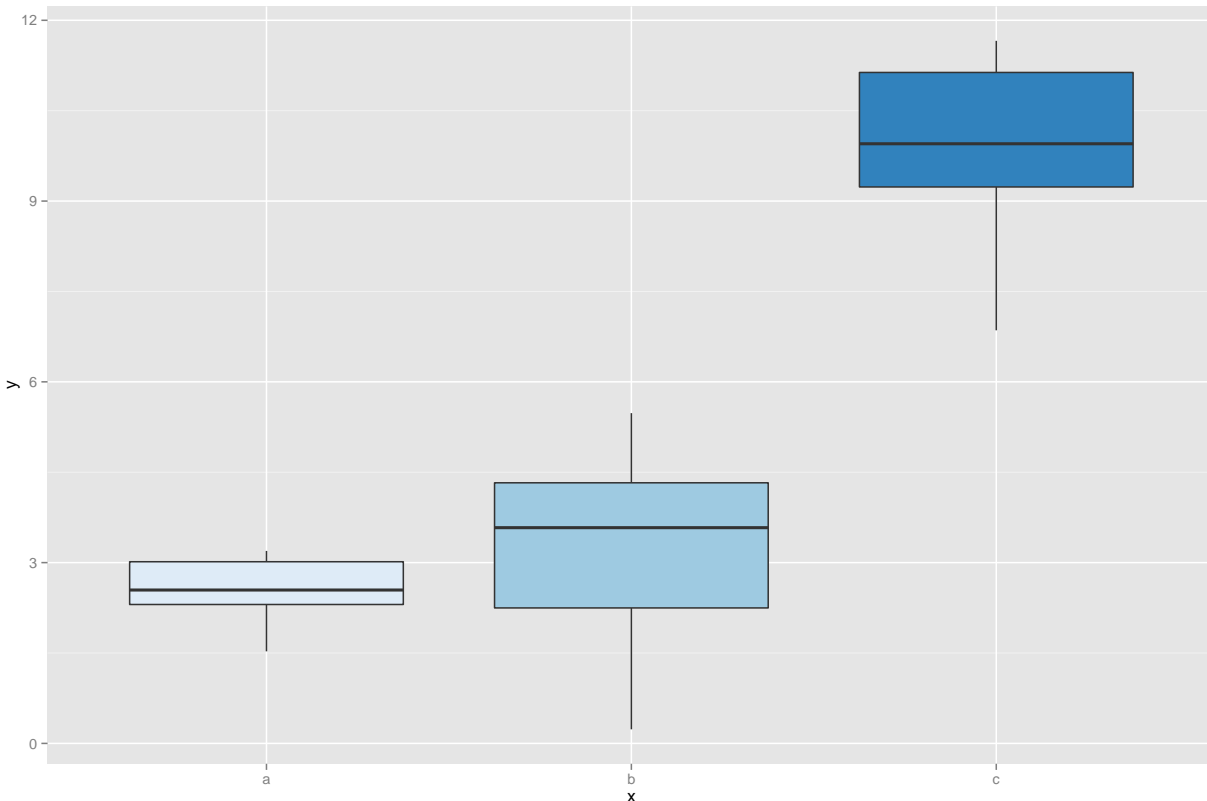
### *A note about boxplots*

Box and whisker plots (boxplots) provide an interesting means of examining differences between categorical variables. To test the difference between categorical variables, a statistician would use an Analysis of Variance

(ANOVA) test. ANOVA is an extention of a linear model and can identify if there are significant differences (with $\alpha$ set to an arbitrary value, normally 0.05) among categories.

```
#Create a vector of numeric values and classify them
y <- c(rnorm(7, 2.5, .5), rnorm(7, 3.5, 2), rnorm(7, 10, 1.5))
x <- c(rep("a", 7), rep("b", 7), rep("c", 7))

#Compile into a dataFrame
tmp <- data.frame(x, y)
names(tmp) <- c("x", "y")
```



In looking at the plot, we can see that there are obvious differences among the groups. For example, group a appears similar to group b and both a and b appear to be quite different from c, but without a statistical test like an ANOVA, we can't know for certain.

```
#Run an ANOVA on the values to look for significant differences
mod <- aov(lm(y~x, data = tmp))
summary(mod)
```

```
##              Df Sum Sq Mean Sq F value   Pr(>F)
## x             2 230.41  115.20   53.78 2.56e-08 ***
## Residuals    18  38.55    2.14
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

The results from the summary(mod) statement show us that the categorical x variable has a significant difference on the y variable. The three asterisks next to the Pr(>F) value in the right column tells us that

the $p$-value is $< 0.05$. (A $p$-value $< 0.05$ is the 'rule of thumb' for a significant value.) Even though we know that x has a significant affect on y, we don't know where the differences lie. A Tukey test will make pair-wise $t$-test comparisons between every x value and every other x value, automatically including the Bonferroni adjustment for multiple $t$-test comparisons.
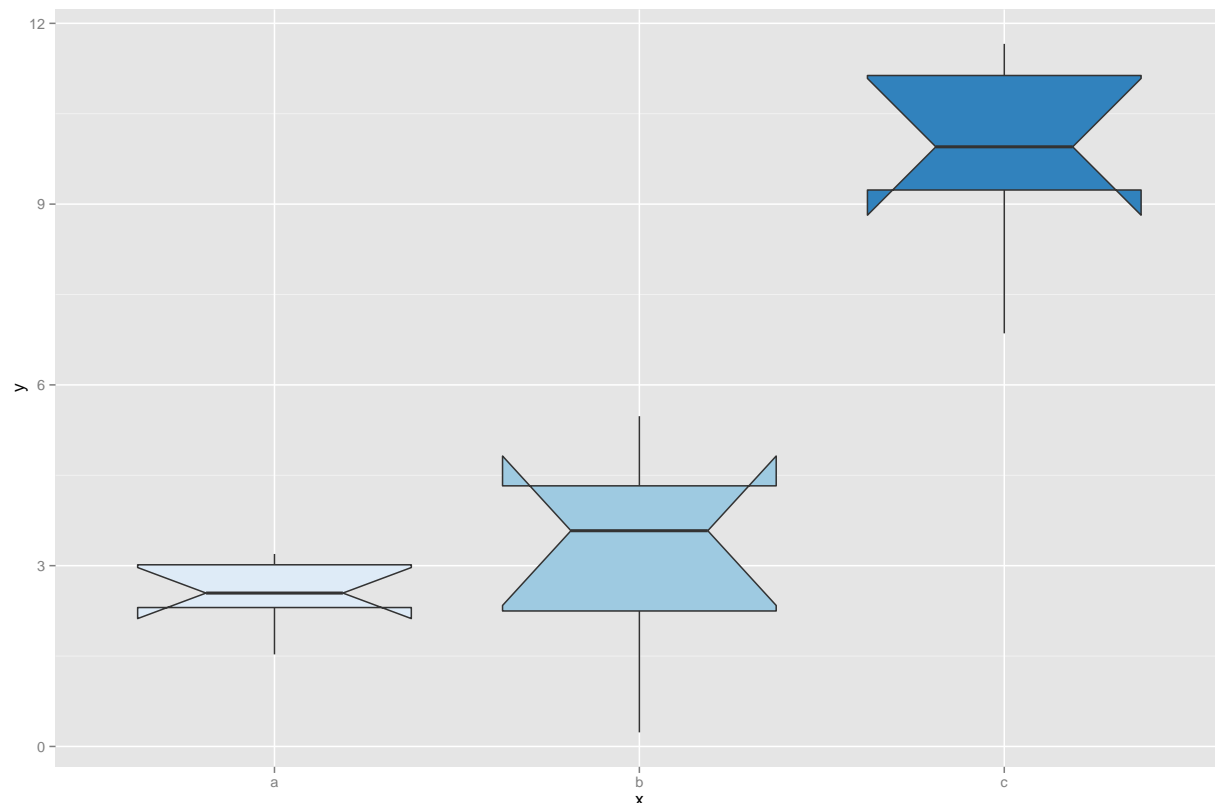
```
#Tukey's test will tell us where the significant (p < 0.05) differences lie
  TukeyHSD(mod)
```

```
##   Tukey multiple comparisons of means
##     95% family-wise confidence level
##
## Fit: aov(formula = lm(y ~ x, data = tmp))
##
## $x
##          diff       lwr      upr      p adj
## b-a 0.6474949 -1.349041 2.644030 0.6911248
## c-a 7.3279230  5.331387 9.324459 0.0000001
## c-b 6.6804281  4.683893 8.676964 0.0000003
```

From the results of the `TukeyHSD()` test, there is not a significant ($p$-value $> 0.05$) difference between group a and group b but there is a significant difference ($p$-value $< 0.05$) between a and c and b and c.

A Tukey test can be applied to any number of x variables produced by an ANOVA. However, making several comparisons among many different x variables would create a *messy* table and can be visually challenging to interpret (Burnham and Anderson 2002). Box and whisker plots in R allow a user to add a notch into the boxes. When `notch = TRUE`, R will notch the box which provides a *rough* approximation of a 95% confidence interval around the median (the heavy line in the middle of the plot.)

In the above plot, the notch from `a` overlaps with the notch from `b` but the neither the notches from `a` or `b` overlap with `c`. This is a good indicator that `a` and `b` are not significantly different from each other but that `a` and `b` are both significantly different from `c`.

While a statistical test is the only *valid* way an analyst can be sure if there are significant differences among categorical variables, adding the notches to a box plot provide a very quick, easy, and tractable way to identify whether or not differences exist. As a result, rather than providing tables and tables of *p*-values in the *ispiranteRanney* package which would be a quantification of the differences between categorgical variables, boxplots will be notched and users will have an effective visual means of determining significant differences without having to rely on, *magic* numbers, and *p*-values. Always remember Mark Twain's quote about statistics: "*There are three kinds of lies: lies, damn lies, and statistics.*" Although magic *p*-values are handy, decisions should't be based upon magic numbers alone. The human component must be preserved.

# The `ispiranteRanney` Package

## General Information

The `ispiranteRanney` package will rely on other packages for seamless operation. Those packages are:

- ggplot2
- scales
- forecast

In most cases, only a few functions will be needed from the above packages. For example, `ggplot2` (Wickham 2009) and `scales` (Wickham 2014) will be used for all plots and only the `forecast()` (Hyndman 2014) function is used from the `forecast` package (Zeileis and Grothendieck 2005, Hyndman 2014, Team et al. 2014). Without these package dependencies, `ispiranteRanney` will not be able to operate successfully.

The `ispiranteRanney` package will include a number of functions to easily visualise how several metrics change as a function of day, date, priority level, agent, etc. In many cases, these insights may be negligible; for example, if all help-desk agents are of high quality and "know their stuff", it's unlikely that the number of reopens as a function of priority level will provide much insight. However, if customer satisfaction score is routinely low, then investigating how satisfaction score changes as a function of date or agent or a number of other variables may be important.

This section of the documentation will include several parts, one related to each of the major variables that will be used in the `ispiranteRanney` package: `callVolume`, `reopens`, `resolutionTime`, `satisfactionScore`, `waitTime`. In many cases, the `ispiranteRanney` package will plot these values against each other. In others, they will be plotted as a function of `priorityLevel`, `agent`, `dayOfWeek`, and `date`. In all cases, the visualization of the insight is completed by plotting one variable against another.

The functions used to gain insight into the data will be named consistently, normally as `Plot[y]By[x]()`. The naming consistency provides the means by which users–or software engineers–will easily be able to identify what insights will be gained by any `ispiranteRanney` function. For example, a function named `PlotResolutionTimeByPriority.R` will do just that; plot the resolution time of a help-desk ticked as a function of the priority level for that ticket.

As disclosed, a number of the columns in the sample data were `factors` rather than `numeric`. However, values that can be converted to `numeric`–like the `...resolution_time_in_minutes_within_business_hours` values can be redefined with `as.numeric()` and used in regression analyses. Similar to boxplots, the `aov()`, and `TukeyHSD()` test described above, numerical values can be compared against each other and the resulting relationship quantified.

### *A note on linear regressions*

Regression is a stistical process for for estimating the relationships among variables. Like the boxplots discussed above, relationships between variables can be evaluated with certain statistical models. While boxplots evaluate the relationship between numeric and categorical data, regression models allow a statistician to evaluate the relationships among numeric data.

```r
#Siumulate data to use in a regression example
x <- 1:10
y <- x+rnorm(10, 0, .5)
sim <- data.frame(x, y)
sim$Class <- "1"

x <- 1:10
y <- x+rnorm(10, 0, 30)
sim1 <- data.frame(x, y)
sim1$Class <- "2"

#Bind the two data frames together into one
sim <- rbind(sim, sim1)

modRed <- lm(y~x, data = sim[sim$Class == "1", ])
modBlue <- lm(y~x, data = sim[sim$Class == "2", ])
```
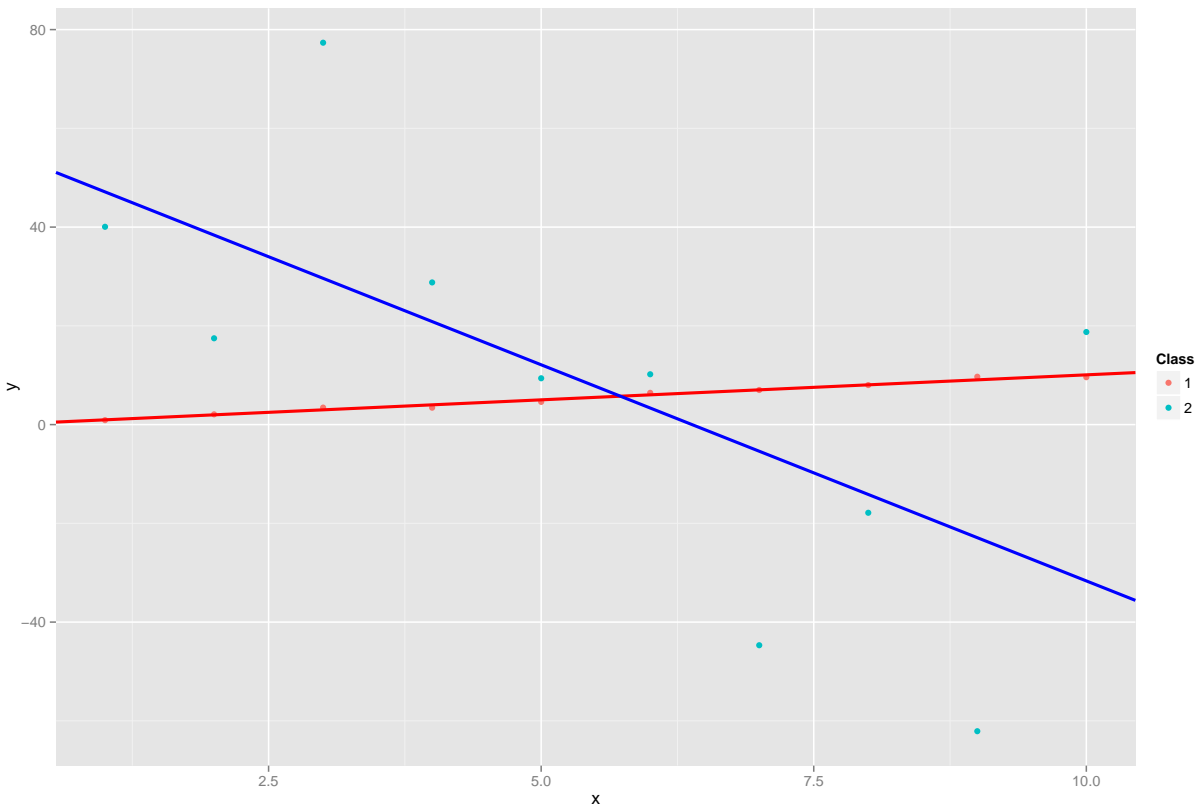


Figure 2: Linear regression of two data sets, one significant with p-value < 0.05 (red) and one with p-value > 0.05 (blue)

The plot above was generated from simulated data. Note how closely the red line follows the data points, especially compared to the blue line. That's a result of the relationship between the x and y variables:

```
summary(modRed)
```

```
##
## Call:
## lm(formula = y ~ x, data = sim[sim$Class == "1", ])
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.57699 -0.32445 -0.03585  0.33252  0.63446
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) -0.04934    0.29813  -0.165    0.873
## x            1.01240    0.04805  21.071  2.7e-08 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.4364 on 8 degrees of freedom
## Multiple R-squared:  0.9823, Adjusted R-squared:  0.9801
## F-statistic:    444 on 1 and 8 DF,  p-value: 2.703e-08
```

```
#P-value of the x variable
summary(modRed)$coefficients[2, 4]
```

```
## [1] 2.703285e-08
```

The *p*-value associated with the x variable the in the model is much less than 0.05, the "rule of thumb" for a significant value. CLearly, the red line follows the red points very closely; even without the *p*-value from the model, it's clear that there is a strong relationship here.

The blue line and blue points appear to be much less closely related. We can infer from the 'scattershot' nature of the blue points that there is unlikely to be a strong impact of the x variable on y:

```
summary(modBlue)
```

```
##
## Call:
## lm(formula = y ~ x, data = sim[sim$Class == "2", ])
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -39.292 -17.440  -3.219   7.651  50.411
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept)   55.886     22.288   2.507   0.0365 *
## x             -8.754      3.592  -2.437   0.0407 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```
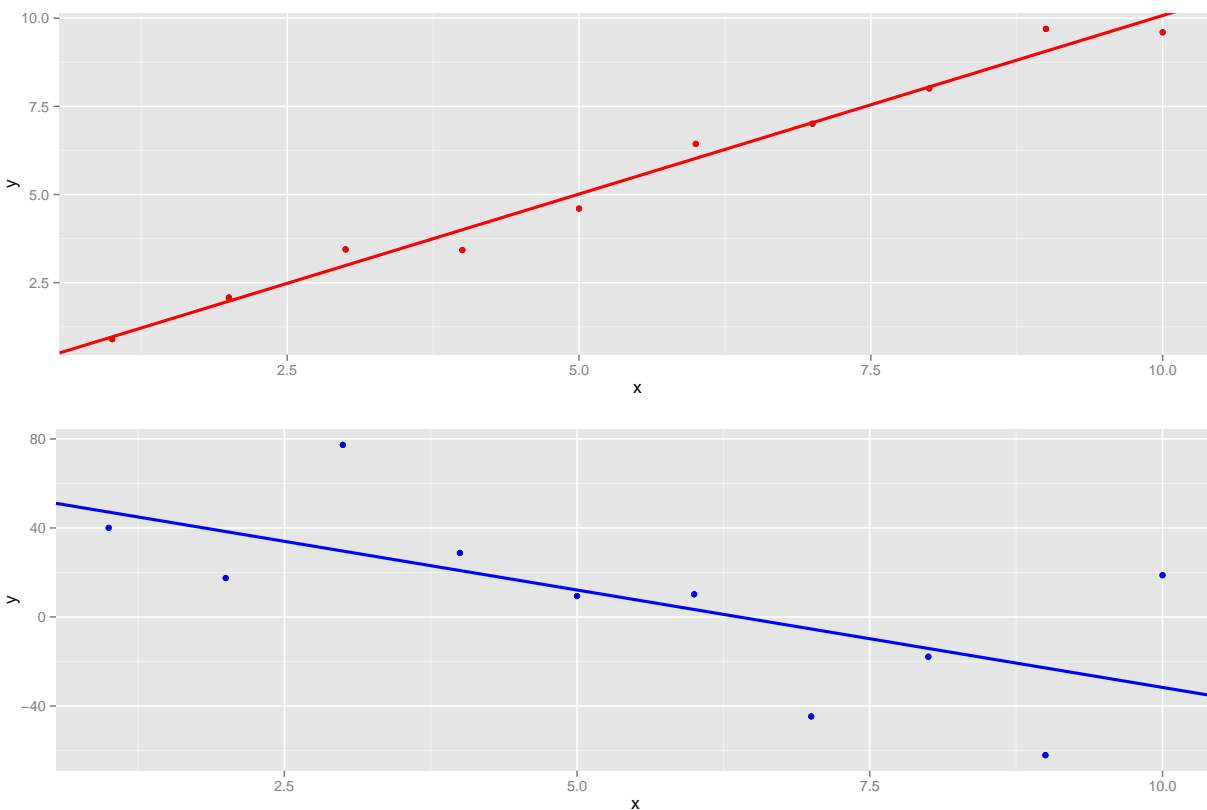
```
##
## Residual standard error: 32.63 on 8 degrees of freedom
## Multiple R-squared:  0.4261, Adjusted R-squared:  0.3544
## F-statistic: 5.939 on 1 and 8 DF,  p-value: 0.04075
```

```
#P-value of the x variable
summary(modBlue)$coefficients[2, 4]
```

```
## [1] 0.0407483
```

the $p$-value of the blue `x` on `y` is $>$ than 0.05, implying that there is not much of a relationship between the two.

Looking at the plots together makes it difficult to see the slope of the red line; the variability in `y` of the blue points condenses the red scale, preventing us from seeing the slope. In looking at them separately:



it's much clearer that there is a strong relationship between the red `x` and `y` variables but not the blue `x` and `y`. While the *significance* can only be calculated by comparing the $p$-value to whatever $\alpha$ we chose, the strength of the relationship can also be visually examined, especially within the context of the randomness of the data points.

The `ispiranteRanney` package will *not* include the actual statistical analyses from linear models presented in these analyses. Unless customers and users have a deep background in statistical analysis, the parameter estimates, standard errors, and $p$-values that result from a `summary()` call to a model can be mis-interpreted. As a result, similar to $p$-values that stem from an `ANOVA` (as identified in the *A Note About Box Plots*), only a red "trend" line will be provided when a numeric value is plotted as a function of another numeric value. Users can clearly see–without a $p$-value to cloud their judgement–the relationship between the `y` variable and the `x` variable. If the line trends up or down in a shallow manner, the relationship can be assumed to be relatively small. Conversely, the steeper the line becomes, the stronger the relationship.

If quantifiable relationships are required, $p$-values, $R^2$ values, and other metrics can be provided.

## Resolution Time

Resolving help-desk ticket quickly is a key component in keeping customer satisfaction levels high. In the help-desk ticket sample data provided by Stefano Spada, there were a numer of values provided for resolution time. The `ispiranteRanney` package will focus on only two values:

- `first_resolution_time_in_minutes_within_business_hours`
- `full_resolution_time_in_minutes_within_business_hours`

The default resolution time value for all functions will be the `first...` value. However, users will be able to select the `full...` value for all visualizations.

There will be five functions that allow users to evaluate resolution time. Functions allow users to visualize resolution time by assignee name, date, day of the week, priority level, and the number of reopens of a ticket.

- `PlotResolutionTimeByAgent()`
- `PlotResolutionTimeByDate()`
- `PlotResolutionTimeByDay()`; day of week will be assigned by a separate, in-package function, `ConvertSatisfactionScoreToNumber()`
- `PlotResolutionTimeByPriority()`
- `PlotResolutionTimeByReopens()`

Resolution time is a key variable in identifying how long it takes a help desk to resolve support tickets. The functions listed here will provide the means for customer service managers to understand how `"full"` or `"first"` resolution time within business hours fluctuates as a function of `agent`, `date`, `day`, `priority`, and the number of `reopens`.

Resolution time likely also plays a key role in customer `satisfaction_score`, another relationship explored in a different proposed function. Regardless of the relationship that resolution time has on any other help-desk ticket item, faster resolution times would increase the number of tickets that help-desk staff could resolve.

## Satisfaction Score

Customer satisfaction score is a powerful measure of how satisfied a customer is with their help-desk experience. In many cases, `satisfaction_score` is not provided by the customer. In cases where satisifaction score *is* returned, it is quite useful.

`ispiranteRanney` will have a number of ways to visualize how `satisfaction_score` changes with other variables. In the `ispiranteRanney` package, `satisfaction_score` will be converted to a numeric variable by the `ConvertSatisfactionScoreToNumber()`. This function is called "behind the scenes" so a user would never have to use it—or even know it is there. The function converts factor variables (e.g., "Good" and "Bad") into numeric variables:

| Satisfaction Score | Number |
| :---: | :---: |
| Good | 4 |
| Bad | 1 |
| Not Offered | NA |
| Offered | NA |

The sample datasets provided had these unique values for `satisfaction_score`; as a result, these were the only scores that could be converted to a numbers. If other help-desk ticket datasets generate additional score values, they could be incorporated into the functions below simply.

As in the `ResolutionTime` section above, `satisfaction_score` will be compared to other variables to understand how `satisfaction_score` fluctuates:

- `PlotScoreByCallsPerPerson()`
- `PlotSscoreByCallVolume()`
- `PlotScoreByResolutionTime()`
- `PlotScoreByWaitTime()`

These funtions will all be visualizations of `numeric` values. `CallsPerPerson`, `CallVolume`, `ResolutionTime`, and `WaitTime` will all be converted to `numeric` values. `satisfaction_score` will be regressed against the above variables with `lm()`; the resultant plot will be similar to the plots in `A note on linear regressions`.

`satisfaction_score` is an important metric by which to measure customer satisfaction. In the sample data, only one dataset included much `satisfaction_score` data; in other example data, `satisfaction_score` data was sparse. Help-desk managers should stress the importance of customers offering feedback to help-desk employees so that those `satisfaction_score` data can be used in evaluating the effectiveness of the help desk.

In all cases, our statistical hypotheses were that `satisfaction_score` would *decrease* as:

1. the number of calls per person increased;
2. call volume increased;
3. resolution time increased and;
4. wait time increased.

With the sample data provided, bits of these hypotheses play out in the regression lines. In one case, the line was effectively parallel, suggesting no relationship while in another we had a spurious result likely generated from a very unhappy customer (very low wait time and low `satisfaction_score`).

## Number of Reopens

The number of reopens of a ticket is the number of times that a "solved" ticket has been reopened. Reopens likely have a negative affect on a number of other numeric values (e.g., `satisfcation_score` and resolution time, discussed previously). However, simply understanding how the number of reopens changes with different categoriacal variables can help a manager understand how the number of reopens changes by `agent`, `call volume`, `date`, `day`, and `priority` level.

The number of reopens of a help-desk ticket will be visualised with functions:

- `PlotReopensByAgent()`
- `PlotReopensByCallVolume()`
- `PlotReopensByDate()`
- `PlotReopensByPriority()`

Plotting the reopens by `assignee_name` will help customer-service managers understand how efficient each help-desk agent is. Boxplots of the reopens by agent will indicate which agents routinely resolve tickets the first time (i.e., `reopens = 0`) and which agents require more time (i.e,. `reopens > 0`). Visulalizations will provide a horizontal line of the mean number of reopens for the entire dataset as well as the boxplot by agent.

Plotting the number of reopens as a function of call volume will allow `ispiranteRanney` to produce a `lm()` model to quantify the relationship between the two numeric values. Managers will be able to see trends in the data. Questions like *if call volume goes up, what happens to the number of reopens*? By plotting `reopens` as a function of call volume, managers can see for themselves if things get "missed" when help-desk agents are busy. From that, managers can determine staffing levels and, if necessary, have agents "on-call" to help out.

In the help-desk ticket data, `date` is effectively a categorical variable. As a result, while we can't create a linear model of the number of reopens as a function of date, we can plot boxplots and look for differences by date.

Reopens by date is an important relationship to evaluate because by looking at call volumes, agent effectiveness, and reopens by date, a customer service manager can start to put together a full picture of how effective and efficient an agent, a group of agents, or their entire help desk can be.

## Wait Time

Wait time of customers calling into a help desk may have a negative affect on satisfaction score. To understand how wait time is affected by variables in the help-desk data, `waitTime` will be regressed as a function of agent, day of the week, priority level, call volume, and any other variables that assist in understanding how wait time fluctuates.

As a numeric variable plotted against categorical variables, users will be able to visualize and unerstand how the categorical variables affect wait time. For example, if wait time is consitently high on Fridays, customer-service managers can adjust staffing schedules to provide more help-desk agents on Fridays. Alternatively, if wait time is close to or at zero on other days of the week, there may be an over-abundance of help-desk agents available.

In conjunction with other available plots in the `ispiranteRanney` package, customer-service managers will be able to understand how univariate models affect variables important to customer satisfaction.

## Call Volume

Call volume, unlike many of the previous ordinate variables, likely has little or nearly no affect on customer satisfaction score. However, if call volumes are high, help-desk agents may be less likely to resolve customer issues quickly. High call volumes likely:

- increase resolution time;
- lower customer satisfaction scores;
- increase the number of reopens of a given ticket, and;
- increase the wait time.

Call volume will be provided in visualizations by date and by day of week. Because call volume can fluctuate so drastically throughout the year, the sample data provided little ability to determine how call volume is impacted by other variables. Rather, call volume provides the means to generate insight in a number of other issues. For example, when evaluating `ResolutionTimeByAgent`, managers may find it helpful to visualize what the call volume of a given agent was on that day. Because that visualization would be helpful, call volume will be provided when other insights are investigated.

In evaluating resolution time by agent, customer-service managers would likely find it helpful to evaluate the call volume of a specific agent. Because that would be useful, `ispiranteRanney` will provide both plots:
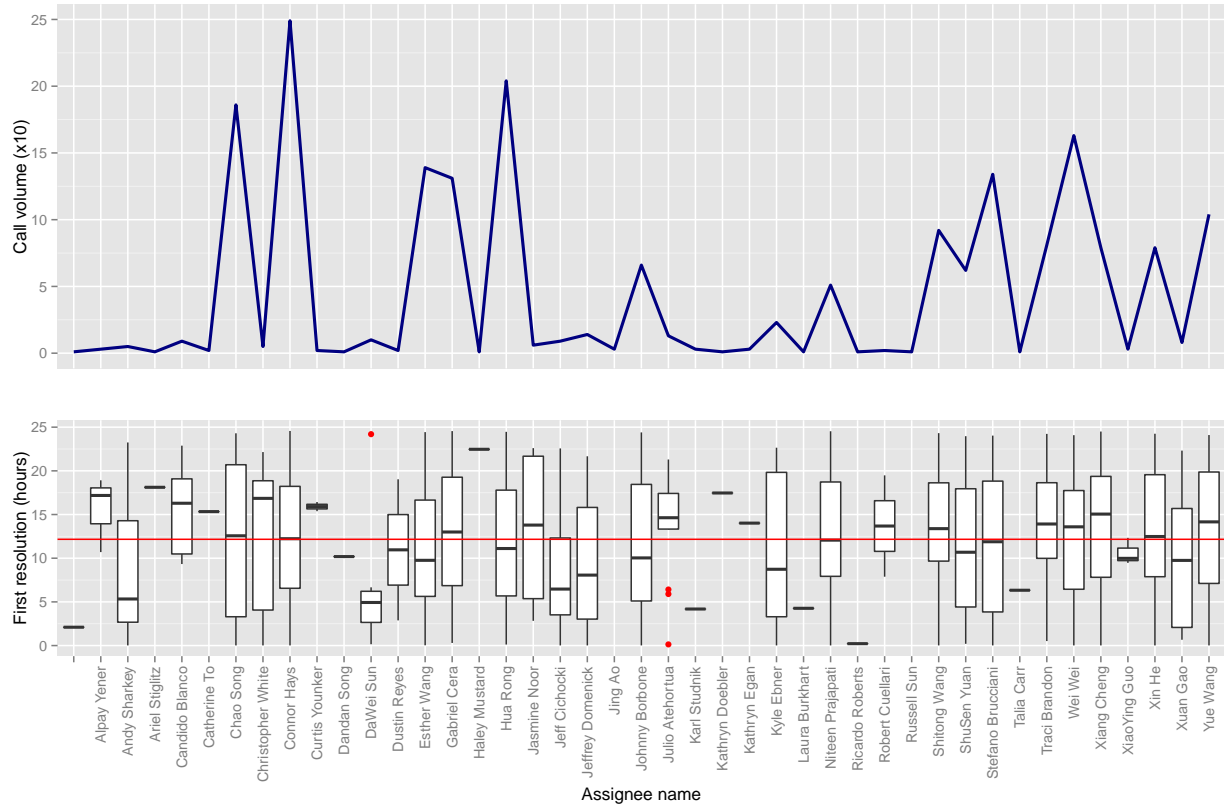
```
## NULL
```

Figure 3: Upper plot of call volume by agent and lower plot of resolution time by agent. From these data, customer-service managers can determine if the resolution time of an `assignee` fluctuates as a result of call volume. For example, the assignee `Connor Hays` has an average resolution time value but his call volume was particularly high. This may be an indication that `Connor Hays` could be used in a training capacity for other agents

In this example, a customer-service manager can see that some agents have very low resolution times while others are much higher. In an effort to identify if high resolution times are a function of call volume, it's easy to see that some agents have very high call volumes with average resolution times (i.e., `Assignee_name ==` `Connor Hays`). This kind of insight is useful for managers to identify who are their top performing agents.

## Forecasting Help-Desk Calls

Predicting the future is difficult. However, with timeseries and the `forecast` package (Hyndman 2014), valid statistical methodology can be applied to approximate how many calls a help-desk may receive $x$ days in the future. For example, in some of the sample data provided, the data was yearly, ranging from `2014-01-01` through `2014-11-29`. The `ispiranteRaney` package will use an `ARIMA` model (Autoregressive Inegrated Moving Average, (Hyndman 2014)) to predict the number of calls a help-desk will receive for as many days into the future as a manager would like.
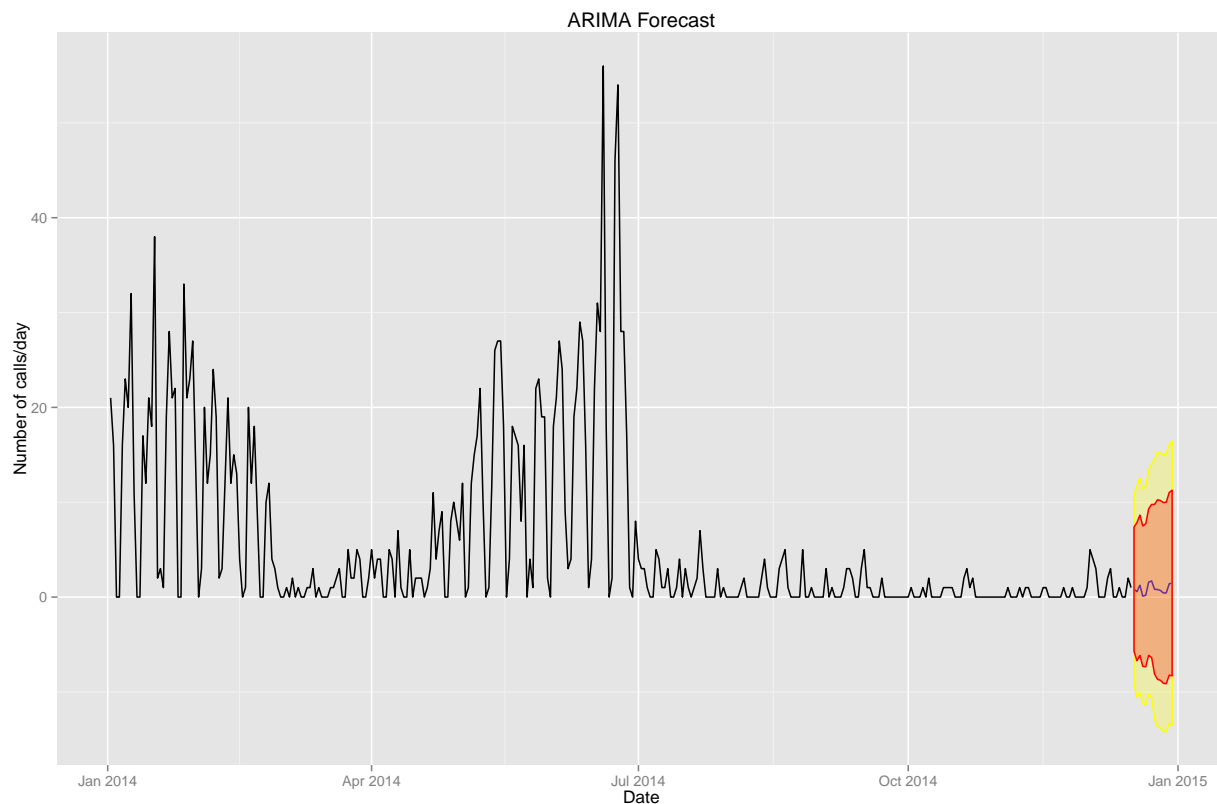


Figure 4: Forecast values for the number of help-desk calls received two weeks from the end of the dataset.

`ForecastHelpDeskCalls()` will use a 14-order model (e.g., a 14-day moving average) to predict the number of calls received. Dates that do not appear in the dataset will be assumed to have zero calls. With the function, users could easily predict out to whatever date they are interested, though, if data is sparse towards the end of the dataset, predictions will be of little value:

```
#Predicting for the next quarter, from whatever is the final date of the dataset
ForecastHelpDeskCalls(dF, 365/4)
```
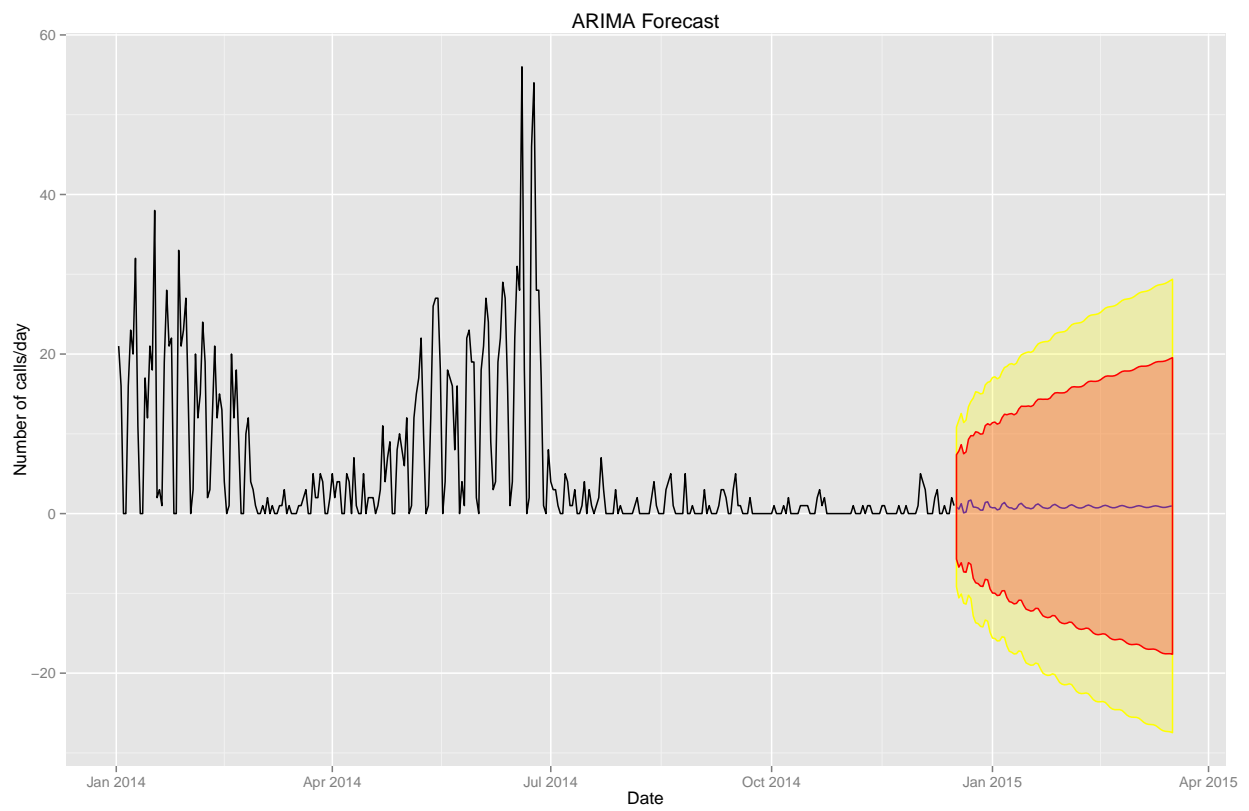
Figure 5: Forecast values for the number of help-desk calls received in the next quarter from the end of the dataset.

## Total Estimated Cost

The estimated cost for creating the `ispiranteRanney` functions, detailed help files, and final packaging into a `.zip` file is:

| Item | Hours | Cost |
|---|---|---|
| Writing functions and documentation | 50.00 | $4250.00 |
| **Total** | **50.00** | **$4250.00** |

Estimated costs include data exploration, authoring and testing of the functions necessary to make analysis and visualization easy, and writing help files. Packaging of the final `R` file can be completed both through the standard `R CMD` procedure, installation through a GitHub repository, or both.

## References

Anderson, D. R. 2008. Model based inference in the life science: A primer on evidence. Springer, New York, NY.

Burnham, K. P., and D. R. Anderson. 2002. Model selection and multimodel inference: A practical information-theoretic approach, 2nd editions. Springer, New York, NY.

Hyndman, R. J. 2014. Forecast: Forecasting functions for time series and linear models.

Team, R. C., D. Wuertz, T. Setz, Y. Chalabi, M. Maechler, and J. W. Byers. 2014. TimeDate: Rmetrics - chronological and calendar objects.

Wickham, H. 2009. Ggplot2: Elegant graphics for data analysis. Springer New York.

Wickham, H. 2014. Scales: Scale functions for graphics.

Zeileis, A., and G. Grothendieck. 2005. Zoo: S3 infrastructure for regular and irregular time series. Journal of Statistical Software 14(6):1–27.