

# Churn Analysis with logistic regression and random forests

Daniel Moscoe

May, 2021

<https://rpubs.com/dmoscoe/768184>

# Abstract

Churn analysis is a fundamental problem in data science. The investigator obtains information on customer behavior and attributes and uses this information to predict whether the customer will terminate a contract, or not. In this study, I conduct a churn analysis based on simulated cell phone customer data from a Kaggle competition, Customer Churn Prediction 2020. I combine this data with information on Google searches pertaining to each of the four major cell phone carriers in the US. The study addresses three main questions:

Does the mean total daytime charge for customers vary by area code?

Does the mean number of customer service calls vary depending on whether a customer carries an international plan?

How can a customer's characteristics be used to predict whether they will terminate their contract?

The study follows the OSEMN workflow. The main strategy for data exploration, in addition to visualization, is hypothesis testing. In the modeling section, I build a logistic regression as well as a random forest model to predict customer churn. I conclude with a summary of my findings.



# Response variable, churn

```
k3.dat %>%  
  select(churn) %>%  
  table()
```

---

```
## .  
## FALSE TRUE  
## 3652 598
```

# Data collection

Google Trends for wireless carrier names, March 2020. `library(gtrendsR)`

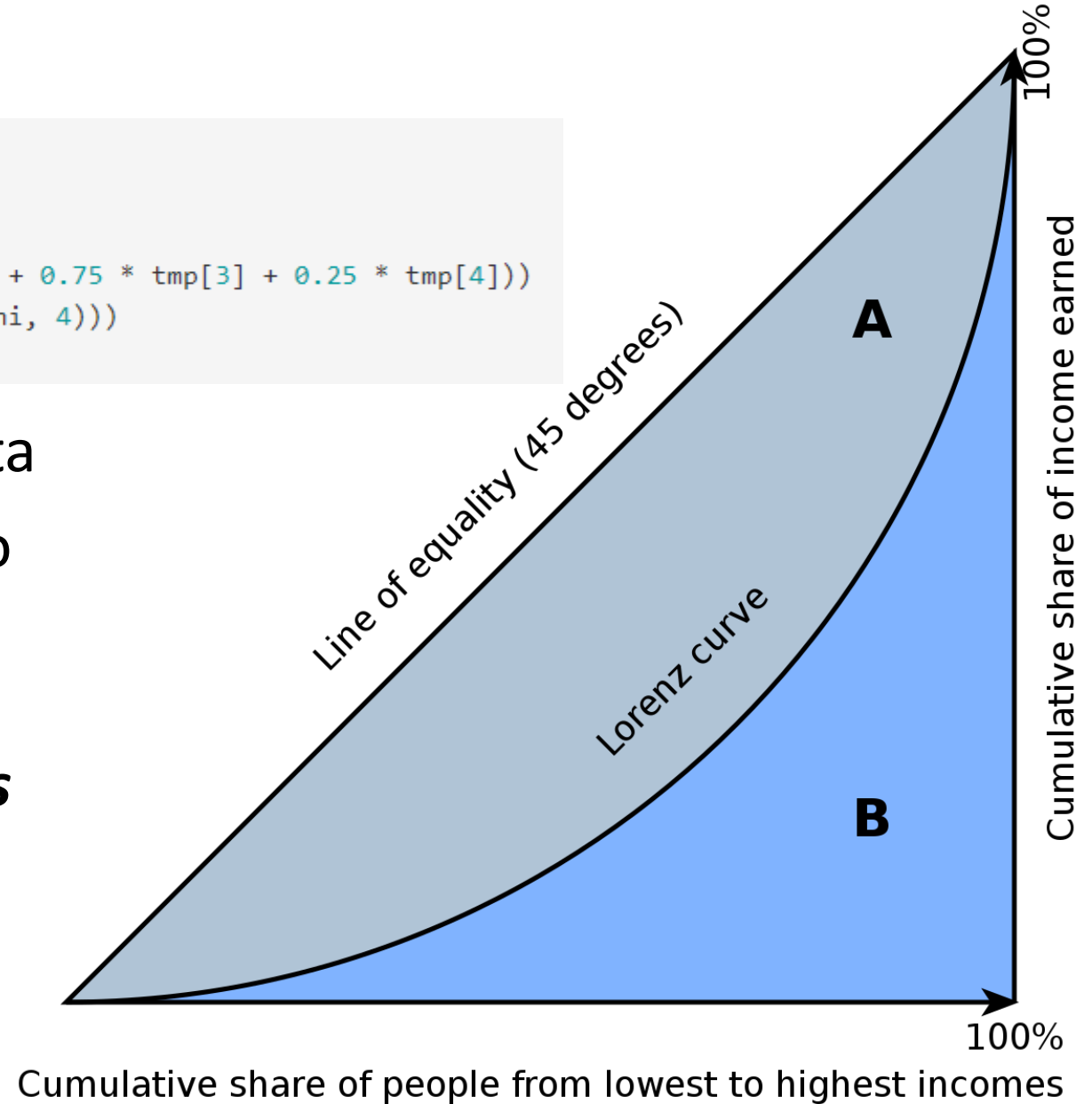
```
gtrends(keyword = gtrends_search_terms[1:4], geo = states[i], time = "2020-03-01 2020-03-30")
```

```
## 'data.frame': 120 obs. of 7 variables:  
## $ date : POSIXct, format: "2020-03-01" "2020-03-02" ...  
## $ hits : int 44 20 84 34 36 76 32 38 69 59 ...  
## $ keyword : chr "att" "att" "att" "att" ...  
## $ geo : chr "US-AL" "US-AL" "US-AL" "US-AL" ...  
## $ time : chr "2020-03-01 2020-03-30" "2020-03-01 2020-03-30" "2020-03-01 2020-03-30" "2020-03-01 2020-03-30" ...  
## $ gprop : chr "web" "web" "web" "web" ...  
## $ category: int 0 0 0 0 0 0 0 0 0 0 ...
```

# The Gini Index

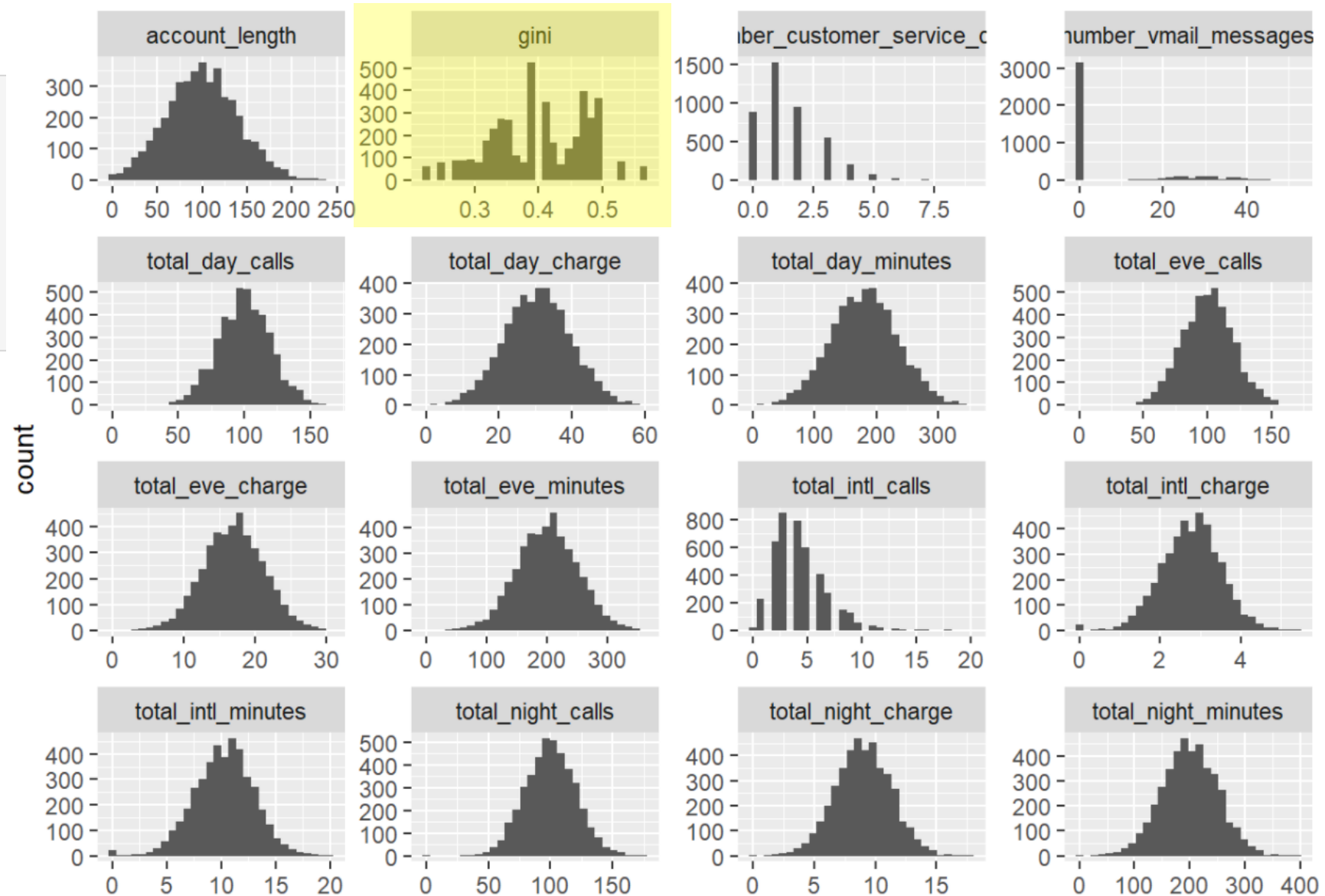
```
ginis <- data.frame("geo" = "x", "gini" = "-1")
for(i in seq(nrow(tmp_query_summaries))) {
  tmp <- sort(as.integer(tmp_query_summaries[i,2:5]))
  tmp.gini <- 1 - ((1/sum(tmp)) * (1.75 * tmp[1] + 1.25 * tmp[2] + 0.75 * tmp[3] + 0.25 * tmp[4]))
  ginis <- rbind(ginis, c(tmp_query_summaries[i,1], round(tmp.gini, 4)))
}
```

- Computed from Google Trends data
- $0 \leq A/(A+B) \leq 1$  (perfect equality to perfect inequality)
- ***Inequality in searches implies inequality in number of customers across providers***



# Explanatory variables

```
k3.dat %>%  
  keep(is.numeric) %>%  
  gather() %>%  
  ggplot(aes(value)) +  
  facet_wrap(~ key, scales = "free") +  
  geom_histogram()
```



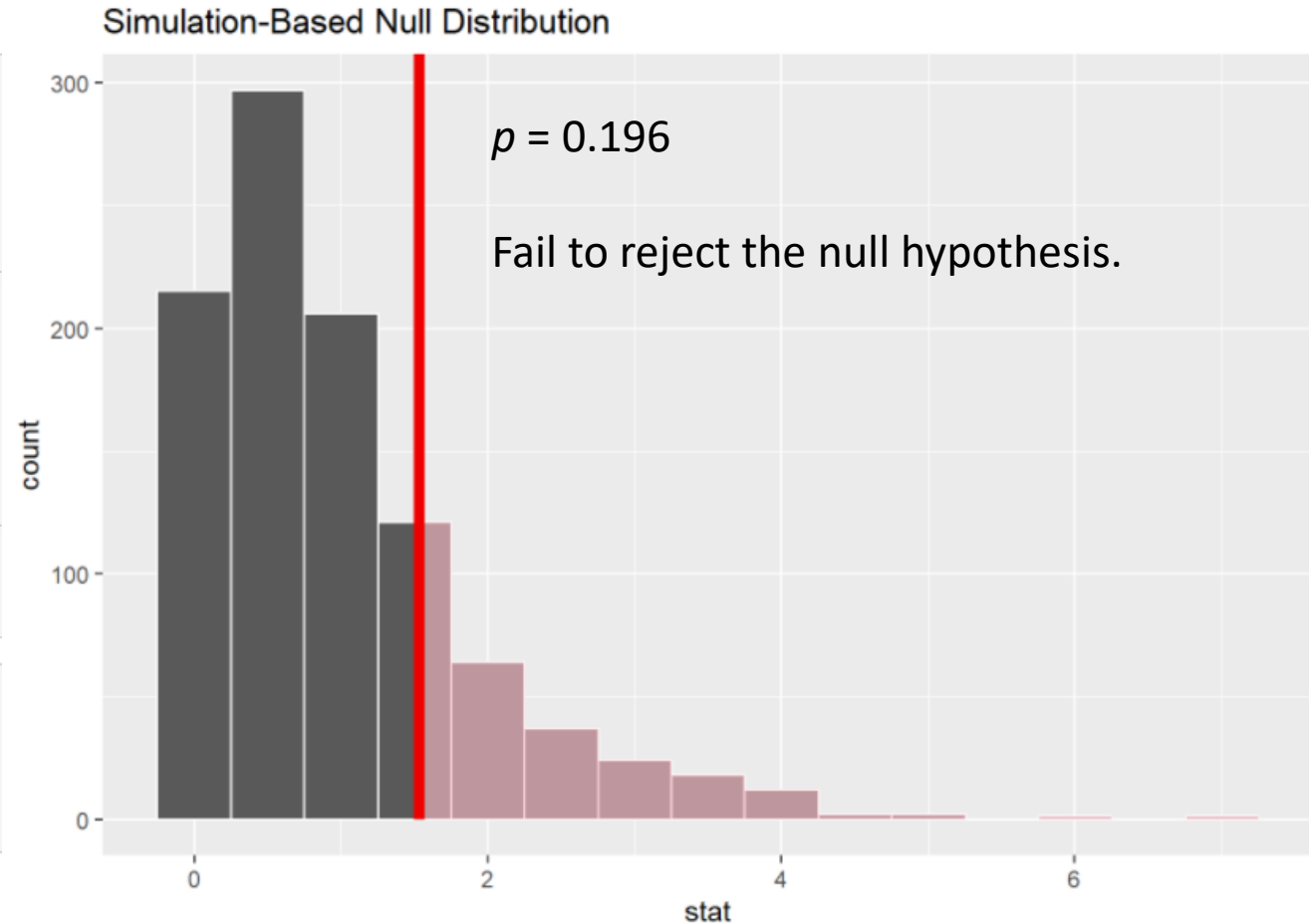
# Does mean total daytime charge vary by area code?

```
#Compute the point estimate
F_hat <- k3.dat %>%
  specify(total_day_charge ~ area_code) %>%
  calculate(stat = "F")

#Generate the null distribution (the sampling distribution of the test stat)
null_f_distn <- k3.dat %>%
  specify(total_day_charge ~ area_code) %>%
  hypothesize(null = "independence") %>%
  generate(reps = 1000, type = "permute") %>%
  calculate(stat = "F")

#Get the p-value for your test statistic with respect to the null distribution
null_f_distn %>%
  get_p_value(obs_stat = F_hat, direction = "greater")
```

```
#Visualize the result
visualize(null_f_distn, method = "simulation") +
  shade_p_value(obs_stat = F_hat, direction = "greater")
```

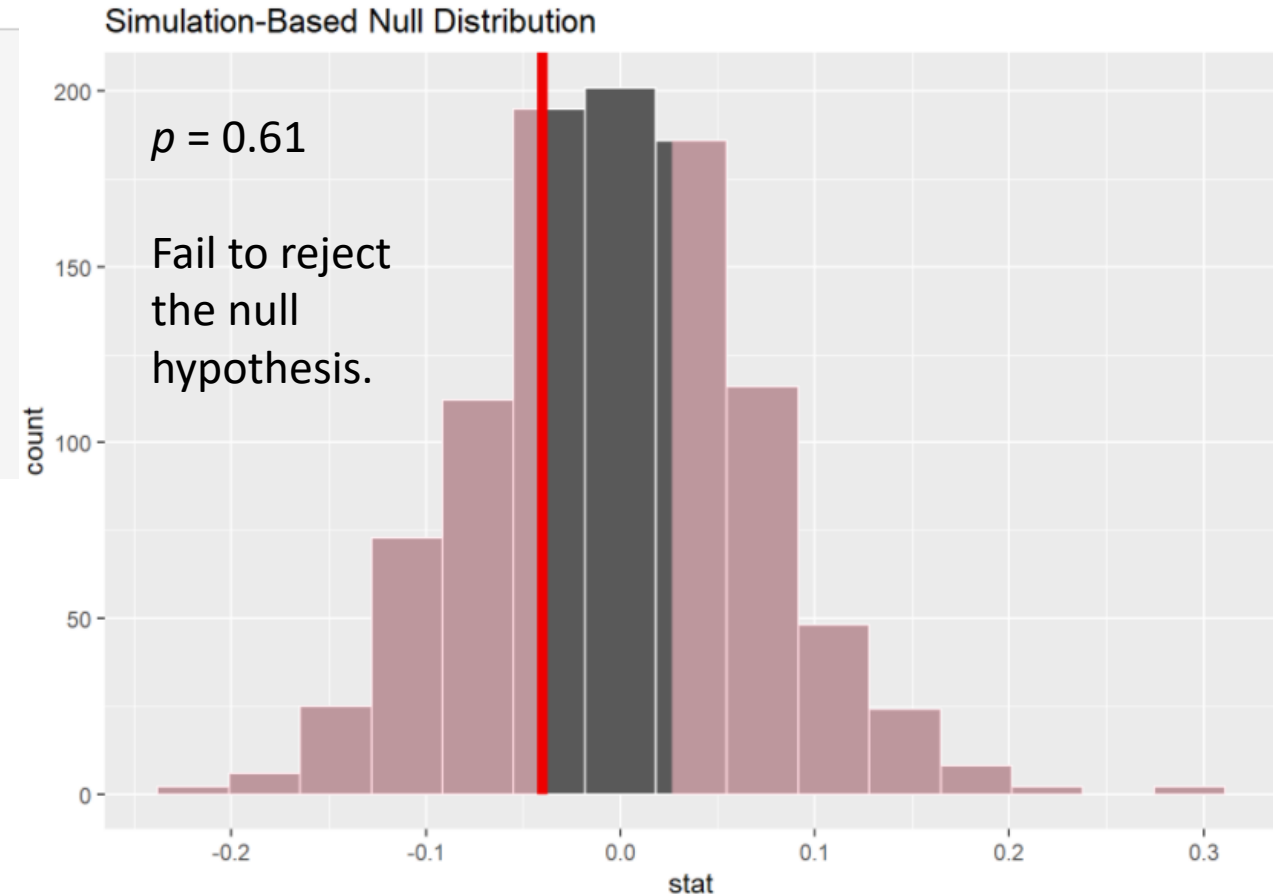




# Does the mean number of customer service calls vary depending on whether a customer carries an international plan?

```
obs_diff <- k3.dat %>%  
  specify(number_customer_service_calls ~ international_plan) %>%  
  calculate(stat = "diff in means", order = c(TRUE, FALSE))  
  
null_t_distn <- k3.dat %>%  
  specify(number_customer_service_calls ~ international_plan) %>%  
  hypothesize(null = "independence") %>%  
  generate(reps = 1000, type = "permute") %>%  
  calculate(stat = "diff in means", order = c(TRUE, FALSE))
```

```
visualize(null_t_distn, method = "simulation") +  
  shade_p_value(obs_stat = obs_diff, direction = "two_sided")
```



# Logistic Regression

- Address class imbalance
- Prune model with backward selection

```
## Call:
## glm(formula = churn ~ total_day_minutes + total_eve_minutes +
##     total_night_minutes + number_customer_service_calls + international_plan +
##     voice_mail_plan, family = binomial(link = "logit"), data = k4_train.dat)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.55164  -0.78148  -0.03148   0.82985   2.67961
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    -6.060681    0.612370  -9.897  < 2e-16 ***
## total_day_minutes    0.014621    0.001472   9.936  < 2e-16 ***
## total_eve_minutes    0.006237    0.001579   3.951  7.79e-05 ***
## total_night_minutes    0.003483    0.001645   2.117   0.0343 *
## number_customer_service_calls  0.630234    0.061077  10.319  < 2e-16 ***
## international_planTRUE    2.487112    0.255333   9.741  < 2e-16 ***
## voice_mail_planTRUE    -1.356577    0.217122  -6.248  4.16e-10 ***
##
##
##
##
```

# Interpreting coefficients of the logistic regression

$$\frac{p}{1-p} = \sum_{i,j} \exp(\beta_i x_{ij})$$

A unit increase in  $x_i$  increases the odds of churn by a *factor* of  $\exp(\beta_i)$ .

Variable	Estimate	Odds change by factor of
total_day_minutes	0.015	1.015
total_eve_minutes	0.0062	1.006
total_night_minutes	0.0035	1.004
number_customer_service_calls	0.6302	1.878
international_plan	2.4871	12.026
voicemail_plan	-1.357	0.257

# Logistic Regression Cutoffs

Maximize detection of true positives:

Cutoff = 0.0857

```
k4_train_preds <- ifelse(k4_train_preds >= optimal_cutoff, TRUE, FALSE)
k4_train_preds_table <- table(k3_test.dat$churn, k4_train_preds)
k4_train_preds_table
```

```
##      k4_train_preds
##      FALSE TRUE
## FALSE    95  627
##  TRUE     0  128
```

Maximize accuracy:

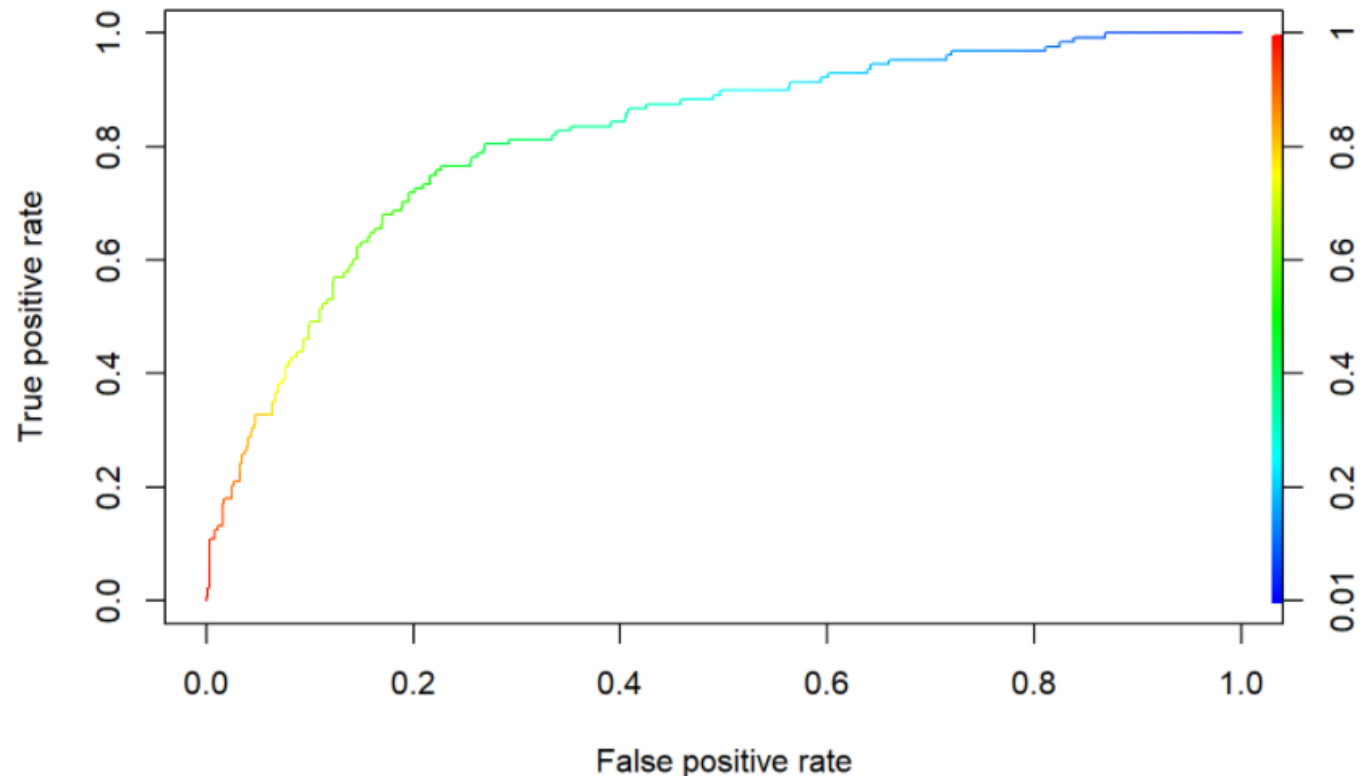
Cutoff = 0.9457

```
k4_train_preds <- ifelse(k4_train_preds >= optimal_cutoff, TRUE, FALSE)
k4_train_preds_table <- table(k3_test.dat$churn, k4_train_preds)
k4_train_preds_table
```

```
##      k4_train_preds
##      FALSE TRUE
## FALSE    719   3
##  TRUE    114  14
```

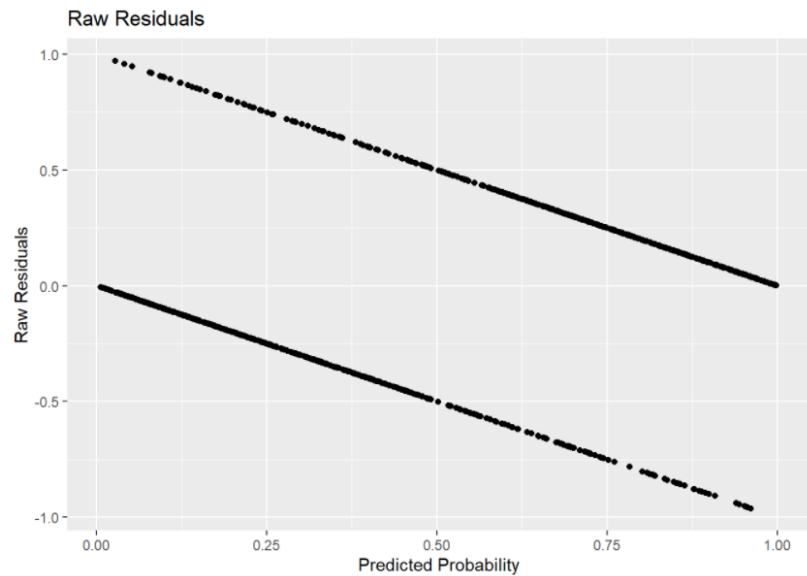
```
k4_train_preds <- predict(k4_train.glm, k3_test.dat, type = "response")
pred <- ROCR::prediction(k4_train_preds, k3_test.dat$churn)
perf <- ROCR::performance(pred, "tpr", "fpr")
plot(perf, colorize = TRUE, main = "ROC curve for logistic regression on churn data")
```

ROC curve for logistic regression on churn data

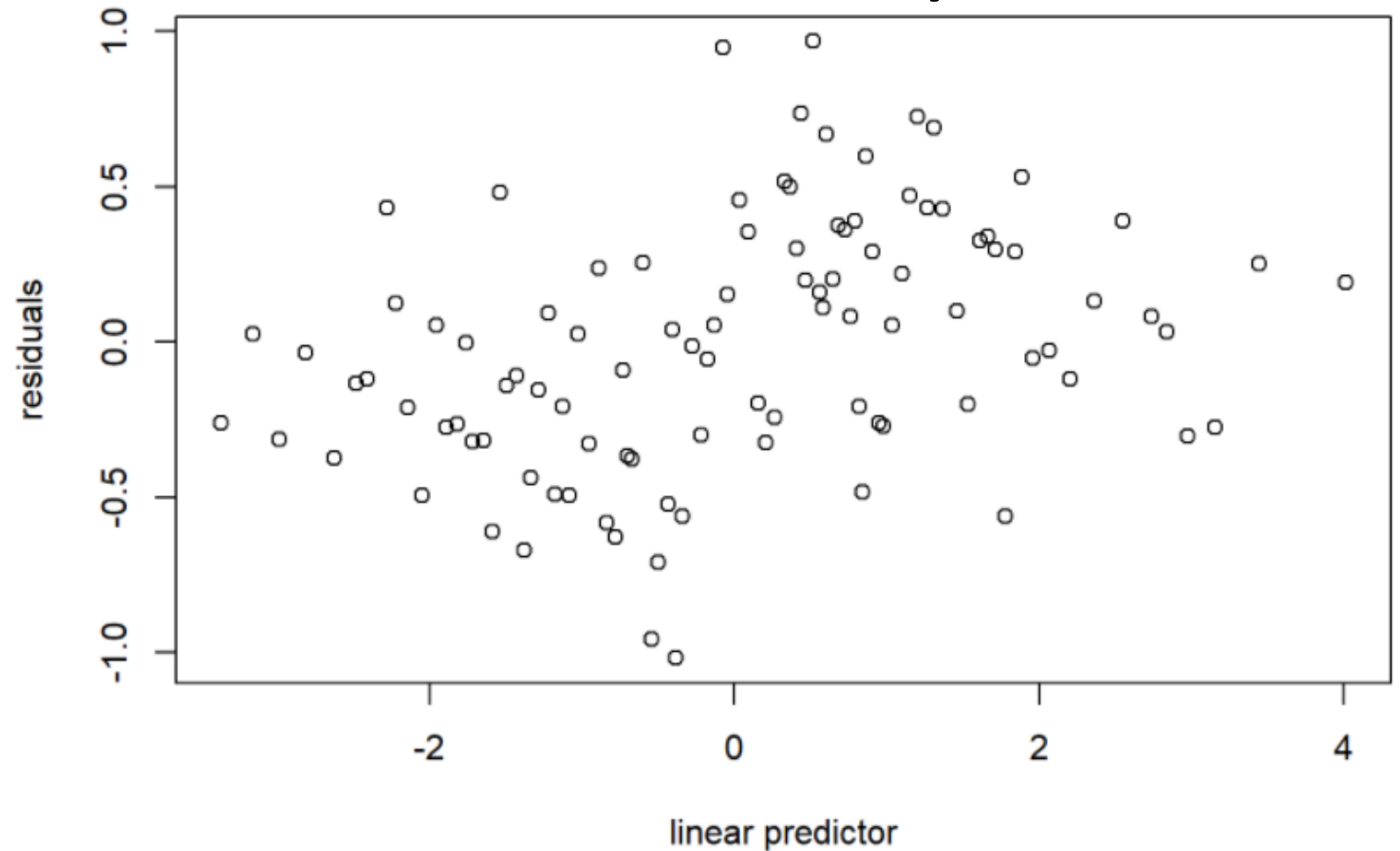


# LR Residuals

```
k5_train.dat <- k4_train.dat %>%  
mutate("residuals" = residuals(k4_train.glm), linpred = predict(k4_train.glm))  
gdf <- group_by(k5_train.dat, cut(linpred, breaks = unique(quantile(linpred, (1:100)/101))))  
diagdf <- summarise(gdf, residuals = mean(residuals), linpred = mean(linpred))  
plot(residuals ~ linpred, diagdf, xlab = "linear predictor")
```



## *Binned residuals plot*



# Tuned Random Forest Model

```
k4_train.for <- randomForest(as.factor(churn) ~ ., k4_train.dat, ntree = 200, mtry = 14, importance = TRUE, proximity = TRUE)
print(k4_train.for)
```

```
##
## Call:
## randomForest(formula = as.factor(churn) ~ ., data = k4_train.dat, ntree = 200, mtry = 14, importance = TRUE, proximity = TRUE)
##              Type of random forest: classification
##              Number of trees: 200
## No. of variables tried at each split: 14
##
##              OOB estimate of error rate: 15.85%
## Confusion matrix:
##      FALSE TRUE class.error
## FALSE  386   84  0.1787234
## TRUE   65  405  0.1382979
```

Applying the new model to the test set:

```
k4_test_preds_for <- predict(k4_train.for, k3_test.dat) #predictions on the k3_test.dat data based on the model trained on k4_train.dat.
confusionMatrix(k4_test_preds_for, as.numeric(k3_test.dat$churn)) #A confusion matrix comparing predicted values for k3_test.dat with actual values.
```

```
##  FALSE TRUE
## 0    603  119
## 1     20  108
```

# Choosing a model

Act.↓ Pred. →	Churn	Remain	Pro	Con
<b>LR max true +</b>			<b>Detects about 100% of churn. Use if cost of churn is much greater than cost of promo.</b>	<b>Low accuracy. Detects many false positives. Many non-churners will receive promo.</b> acc = 0.26, prec = 0.17, rec = 1, spec = 0.13
Churn	0.15	0		
Remain	0.74	0.11		
<b>LR max acc.</b>			<b>High accuracy.</b>	<b>Essentially equivalent to predicting majority class for all obs's.</b> acc = 0.87, prec = 0.83, Rec = 0.13, spec = 0.99
Churn	0.02	0.13		
Remain	0.004	0.85		
<b>Random Forest</b>			<b>Detects almost all churn. Detects majority of non-churn. High accuracy.</b>	<b>Some false positives.</b> acc = 0.84, prec = 0.48, Rec = 0.87, spec = 0.84
Churn	0.13	0.02		
Remain	0.14	0.71		

Submission and Description

Private Score

Public Score

[sub\\_set.csv](#)

0.81333

0.80444

2 minutes ago by [Daniel Moscoe](#)

# Conclusion

- Finding value doesn't always mean improving accuracy.
- The LR yields “collateral insights.”