

Dating Research

Léon Yuan

Table of contents

Import data	1
Conduct logistic regressions separately for male and female	2
Decisions made by females to male when dating	2
Logistic regression only on race	3
Logistic regression includes the race and six attributes	5
Random Forests for females' decisions on more variables	7
Build random forests model for it	8
Load packages	8
Split the data into training and test set	10
Construct a random forest model for this training data	10
Confusion matrix	11
Variable Importance	12
Receiver Operating Characteristic comparing random forests with logistics regression on the same train and test set	14

Import data

```
1 > library(tidyverse)
2 > library(kableExtra)
3 > Speed_Dating_Data <- read_csv("./data/Speed Dating Data.csv")
4 > head(Speed_Dating_Data) |>
5 +   kable(booktabs = TRUE,
6 +       caption = "Speed Dating Data") |>
7 +   kable_styling(latex_options="striped")
```

Table 1: Speed Dating Data

iid	id	gender	idg	condtn	wave	round	position	positin1	order	partner	pid	match	int_cor
1	1	0	1	1	1	10	7	NA	4	1	11	0	0.1
1	1	0	1	1	1	10	7	NA	3	2	12	0	0.5
1	1	0	1	1	1	10	7	NA	10	3	13	1	0.1
1	1	0	1	1	1	10	7	NA	5	4	14	1	0.6
1	1	0	1	1	1	10	7	NA	7	5	15	1	0.2
1	1	0	1	1	1	10	7	NA	6	6	16	0	0.2

Conduct logistic regressions separately for male and female

The reason I build two separate models for females and males is because there are some big differences in dating behaviors between genders and separate models are easier to interpret.

Decisions made by females to male when dating

```

1 > # Filter data when females date males
2 > females_to_males <- Speed_Dating_Data |>
3 +   filter(gender == 0) |>
4 +   select(dec_o, samerace, race_o, age_o, attr_o, sinc_o, intel_o, fun_o, amb_o, shar_o,
5 +         age, race)
6 > # Convert numerical decision into factor type
7 > females_to_males$dec_o <- factor(females_to_males$dec_o,
8 +                                levels = c(0,1),
9 +                                labels = c("No", "Yes"))
10 > # Make the glm predict the Yes as 1
11 > contrasts(females_to_males$dec_o)

```

```

      Yes
No      0
Yes     1

```

```

1 > females_to_males$samerace <- factor(females_to_males$samerace,
2 +                                   levels = c(0,1),
3 +                                   labels = c("no", "yes"))
4 > contrasts(females_to_males$samerace)

```

```

      yes
no      0
yes     1

```

```

1 > females_to_males$race <- factor(females_to_males$race,
2 +                               levels = 1:6,
3 +                               labels = c("Black", "White", "Latino", "Asian", "Native", "Other"))
4 > contrasts(females_to_males$race)

```

	White	Latino	Asian	Native	Other
Black	0	0	0	0	0
White	1	0	0	0	0
Latino	0	1	0	0	0
Asian	0	0	1	0	0
Native	0	0	0	1	0
Other	0	0	0	0	1

```

1 > # delete the missing value rows
2 > females_to_males <- females_to_males |>
3 +   drop_na()

```

Logistic regression only on race

First I only care about how race affects females' decisions to males, only including the **samerace** and **race** columns in this logistic classification model. From the summary of model, we can tell that all females are likely to reject the **Asian** males because **Asian** males has 0.008 p-value which is the most significant in this model. The log odd of saying "yes" to Asian males by all females is -0.48977 given other variables fixed and this is significant negative coefficient meaning that Asian males are very unpopular when dating. Thus, the odd ratio of say "yes" to Asian males is $e^{-0.48977} = 0.6127673$ which means when females date Asian males, they likely decrease 40% probability of saying "yes" to Asian males. Also, **samerace** doesn't show statistical significance because of relatively large p-value 0.06 which is counter intuitive to common sense that females are preferred same race dating.

```

1 > fit <- glm(data = females_to_males,
2 +           formula = dec_o ~ samerace+race,
3 +           family=binomial(link='logit'))
4 > summary(fit)

```

```
Call:
glm(formula = dec_o ~ samerace + race, family = binomial(link = "logit"),
    data = females_to_males)
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-1.232	-1.191	-1.015	1.164	1.349

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-0.13303	0.13970	-0.952	0.3410
sameraceyes	-0.09167	0.07931	-1.156	0.2477
raceWhite	0.25750	0.15447	1.667	0.0955 .
raceLatino	0.26152	0.17670	1.480	0.1389
raceAsian	-0.16904	0.15608	-1.083	0.2788
raceOther	0.03121	0.19477	0.160	0.8727

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 4848.5 on 3498 degrees of freedom
Residual deviance: 4821.1 on 3493 degrees of freedom
AIC: 4833.1

Number of Fisher Scoring iterations: 3

Then the ANOVA table shows that **race** variable has very small p-value which shows it is very significant as I said before.

```
1 > anova(fit, test="Chisq")
```

Analysis of Deviance Table

Model: binomial, link: logit

Response: dec_o

Terms added sequentially (first to last)

	Df	Deviance	Resid. Df	Resid. Dev	Pr(>Chi)
NULL			3498	4848.5	

```

samerace 1 0.0955 3497 4848.4 0.7573
race 4 27.3158 3493 4821.1 1.716e-05 ***
---
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

Logistic regression includes the race and six attributes

Once I included six attribute scores and race together in the logistic regression model, the significance level of race changes radically, because none of race is significant once six attributes are included. This suggests that males' personal attributes can overturn/change the females' impressions or decisions deeply. As we can see from the p-values, all six attributes are statistically significant, especially physical attractiveness, fun, ambitious, shared interests play major roles in making decisions.

The coefficient of *physical attractiveness* is 0.39356 , this means log odds of saying “yes” to males by females increases 0.39356 given other variables fixed, and odd ratios of saying “yes” to males by females increases $e^{0.39356} = 1.482248$ when one more score is given to *attractiveness*.

The coefficient of *fun* is 0.27850 , this means log odds of saying “yes” to males by females increases 0.27850 given other variables fixed, and odd ratios of saying “yes” to males by females increases $e^{0.27850} = 1.321147$ when one more score is given to *fun*.

The coefficient of *shared interests* is 0.27081 , this means log odds of saying “yes” to males by females increases 0.27081 given other variables fixed, and odd ratios of saying “yes” to males by females increases $e^{0.27081} = 1.311026$ when one more score is given to *shared interests*.

All these three most significant attributes have positive coefficient meaning that more scores on these attributes will help females a lot make “yes” decisions to males.

```

1 > fit1 <- glm(data = females_to_males,
2 +   formula = dec_o ~ samerace+race+attr_o+sinc_o+intel_o+fun_o+amb_o+shar_o,
3 +   family=binomial(link='logit'))
4 > summary(fit1)

```

Call:

```

glm(formula = dec_o ~ samerace + race + attr_o + sinc_o + intel_o +
    fun_o + amb_o + shar_o, family = binomial(link = "logit"),
    data = females_to_males)

```

Deviance Residuals:

	Min	1Q	Median	3Q	Max
	-2.6957	-0.8064	-0.1512	0.8302	3.0773

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-5.07251	0.31913	-15.895	< 2e-16 ***
sameraceyes	-0.10644	0.09641	-1.104	0.270
raceWhite	-0.15428	0.19123	-0.807	0.420
raceLatino	-0.33828	0.21824	-1.550	0.121
raceAsian	-0.19103	0.19364	-0.987	0.324
raceOther	-0.16636	0.24082	-0.691	0.490
attr_o	0.69759	0.03381	20.630	< 2e-16 ***
sinc_o	-0.16373	0.03732	-4.388	1.15e-05 ***
intel_o	-0.03924	0.04338	-0.904	0.366
fun_o	0.26326	0.03503	7.514	5.72e-14 ***
amb_o	-0.15977	0.03457	-4.621	3.82e-06 ***
shar_o	0.26955	0.02648	10.178	< 2e-16 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 4848.5 on 3498 degrees of freedom
Residual deviance: 3502.6 on 3487 degrees of freedom
AIC: 3526.6

Number of Fisher Scoring iterations: 5

The ANOVA table also shows that attractiveness, fun and shared interests explain the most deviance residuals by 710.98, 155.36, 106.23 compared to other variables' explained variations which are consistent with our above finding.

```
1 > anova(fit1, test="Chisq")
```

Analysis of Deviance Table

Model: binomial, link: logit

Response: dec_o

Terms added sequentially (first to last)

	Df	Deviance	Resid. Df	Resid. Dev	Pr(>Chi)
NULL			3498	4848.5	

samerace	1	0.10	3497	4848.4	0.757305	
race	4	27.32	3493	4821.1	1.716e-05	***
attr_o	1	1082.88	3492	3738.2	< 2.2e-16	***
sinc_o	1	2.40	3491	3735.8	0.121225	
intel_o	1	0.33	3490	3735.5	0.563586	
fun_o	1	117.67	3489	3617.8	< 2.2e-16	***
amb_o	1	7.26	3488	3610.5	0.007047	**
shar_o	1	107.93	3487	3502.6	< 2.2e-16	***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Random Forests for females' decisions on more variables

```

1 > # Filter data when females date males with more variables than logistics regression
2 > tree_females_to_males <- Speed_Dating_Data |>
3 +   filter(gender == 0) |>
4 +   select(dec_o,
5 +         samerace,
6 +         attr_o, sinc_o, intel_o, fun_o, amb_o, shar_o,
7 +         int_corr, age, race, field, from)
8 > # Convert numerical decision into factor type
9 > tree_females_to_males$dec_o <- factor(tree_females_to_males$dec_o,
10 +                                     levels = c(0,1),
11 +                                     labels = c("No", "Yes"))
12 > # Make the glm predict the Yes as 1
13 > contrasts(tree_females_to_males$dec_o)

```

	Yes
No	0
Yes	1

```

1 > tree_females_to_males$samerace <- factor(tree_females_to_males$samerace,
2 +                                         levels = c(0,1),
3 +                                         labels = c("no", "yes"))
4 > contrasts(tree_females_to_males$samerace)

```

	yes
no	0
yes	1

```

1 > tree_females_to_males$race <- factor(tree_females_to_males$race,
2 +                                     levels = 1:6,
3 +                                     labels = c("Black", "White", "Latino", "Asian", "Native", "Other"))
4 > contrasts(tree_females_to_males$race)

```

	White	Latino	Asian	Native	Other
Black	0	0	0	0	0
White	1	0	0	0	0
Latino	0	1	0	0	0
Asian	0	0	1	0	0
Native	0	0	0	1	0
Other	0	0	0	0	1

```

1 > # Drop missing rows
2 > tree_females_to_males <-
3 +   tree_females_to_males |>
4 +   drop_na()

```

Now I checked if the response variable `dec_o` is balanced or not. The ratio of No to Yes is 1.68 which shows relative balanced within the accepted range from 0.5 to 2. Thus, I don't need to make any efforts to balance the dataset.

```

1 > table(tree_females_to_males$dec_o)

```

```

      No  Yes
1786 1705

```

Initially I did want to include `income` variable in the random forest, however, I found there are half of income variables missing, so I have to drop this variable.

```

1 > sum(is.na(Speed_Dating_Data$income))

```

```

[1] 4099

```

Build random forests model for it

Load packages


```
1 > library(randomForest)
```

randomForest 4.7-1.1

Type `rfNews()` to see new features/changes/bug fixes.

Attaching package: 'randomForest'

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

combine

The following object is masked from 'package:ggplot2':

margin

```
1 > library(datasets)
2 > library(caret)
```

Loading required package: lattice

Attaching package: 'caret'

The following object is masked from 'package:purrr':

lift

```
1 > library(pROC)
```

Type `'citation("pROC")'` for a citation.

Attaching package: 'pROC'

The following objects are masked from 'package:stats':

cov, smooth, var

```
1 > library(glmnet)
```

Loading required package: Matrix

Attaching package: 'Matrix'

The following objects are masked from 'package:tidyr':

expand, pack, unpack

Loaded glmnet 4.1-4

Split the data into training and test set

I randomly splited data into 80% training and 20% test set.

```
1 > set.seed(222)
2 > ind <- sample(2, nrow(tree_females_to_males), replace = TRUE, prob = c(0.8, 0.2))
3 > train <- tree_females_to_males[ind==1,]
4 > test <- tree_females_to_males[ind==2,]
```

Check how many observations in training and how many in test set. There are 3366 rows in training and 826 in the test set.

```
1 > dim(train)
```

```
[1] 2808  13
```

```
1 > dim(test)
```

```
[1] 683  13
```

Construct a random forest model for this training data

I chose 500 trees and 4 random predictors at each split.

```

1 > rf <- randomForest(x = train[-1],
2 +                   y = train$dec_o,
3 +                   xtest = test[-1],
4 +                   ytest = test$dec_o,
5 +                   ntree = 500,
6 +                   mtry = 4,
7 +                   proximity = TRUE)

```

Print out the random forests. The OOB estimate of error rate is 24.99% which has 75% accuracy on the training set while on the test set, this RF has roughly 73% test accuracy which is not bad on this dating data.

```

1 > print(rf)

```

Call:

```

randomForest(x = train[-1], y = train$dec_o, xtest = test[-1],      ytest = test$dec_o, ntr
              Type of random forest: classification
              Number of trees: 500

```

No. of variables tried at each split: 4

OOB estimate of error rate: 23.43%

Confusion matrix:

	No	Yes	class.error
No	1131	325	0.2232143
Yes	333	1019	0.2463018

Test set error rate: 25.77%

Confusion matrix:

	No	Yes	class.error
No	248	82	0.2484848
Yes	94	259	0.2662890

Confusion matrix

Print out the confusion matrix and other statistical measures on this classification results. The whole accuracy on the test set is 72.71%. The true “Yes” rate is $138/(138+118) = 53.90$ which is a little bit over 50% random guess rate. However, the true “No” rate is $347/(347+64) = 84.43$ which is a better prediction rate on the test set because training set has more “No” classes than “Yes”.

```

1 > confusionMatrix(data = rf$test$predicted,
2 +               reference = test$dec_o)

```

Confusion Matrix and Statistics

```

              Reference
Prediction   No  Yes
   No      248   94
   Yes     82  259

      Accuracy : 0.7423
      95% CI : (0.7078, 0.7747)
No Information Rate : 0.5168
P-Value [Acc > NIR] : <2e-16

      Kappa : 0.4847

McNemar's Test P-Value : 0.407

      Sensitivity : 0.7515
      Specificity : 0.7337
   Pos Pred Value : 0.7251
   Neg Pred Value : 0.7595
      Prevalence : 0.4832
   Detection Rate : 0.3631
Detection Prevalence : 0.5007
   Balanced Accuracy : 0.7426

      'Positive' Class : No

```

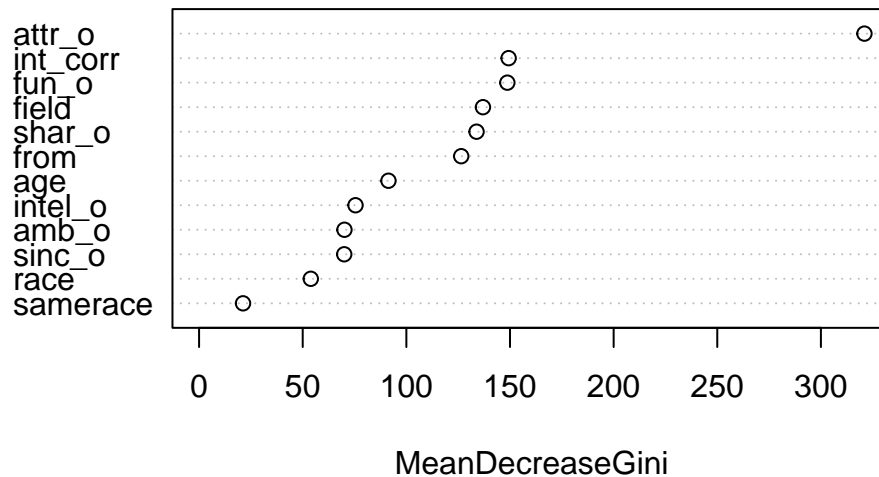
Variable Importance

From the variable importance plot, I roughly classify the top 6 predictors into three classes. The first top class only has one predictor, physical attractiveness. This is consistent with my logistic regression. *Physical Attractiveness* has 205.911 mean decrease of Gini that is measurement of building trees. This Gini decrease is almost twice of other variables. Thus, *Physical attractiveness* is the most significant factor when females make decision to males. The second top class has *shared interests* and the correlation between participant's and partner's ratings of interests. These two variables actually are highly correlated however random forest is robust to the highly correlated predictors because of its ability of randomly selecting a

subset of variables at each split. *Shared Interests* has 150.324 decrease Gini of mean which is as three times as other less importance variables. Females secondary emphasize shared interests with males. The third top class has three variables: *fun*, *from*, *field*. They have very close Gini decrease mean about 125 which is as twice as other less important variables. Females put equally emphasis on the fun, where males are from, and which field males' careers belong to. Overall females are likely to date males who are very physical attractive then have common/shared interests as they do while males' career fields and where they're from play a secondary role in dating.

```
1 > varImpPlot(rf,
2 +           sort = T,
3 +           n.var = 12,
4 +           main = "Top 12 - Variable Importance")
```

Top 12 – Variable Importance



```
1 > importance(rf)
```

```

      MeanDecreaseGini
samerace      21.20644
attr_o       320.99262
sinc_o        70.02605
intel_o       75.47259
fun_o        148.72444
amb_o         70.13057
shar_o        133.86015
```

int_corr	149.34420
age	91.37031
race	53.89209
field	136.94824
from	126.47190

Receiver Operating Characteristic comparing random forests with logistics regression on the same train and test set

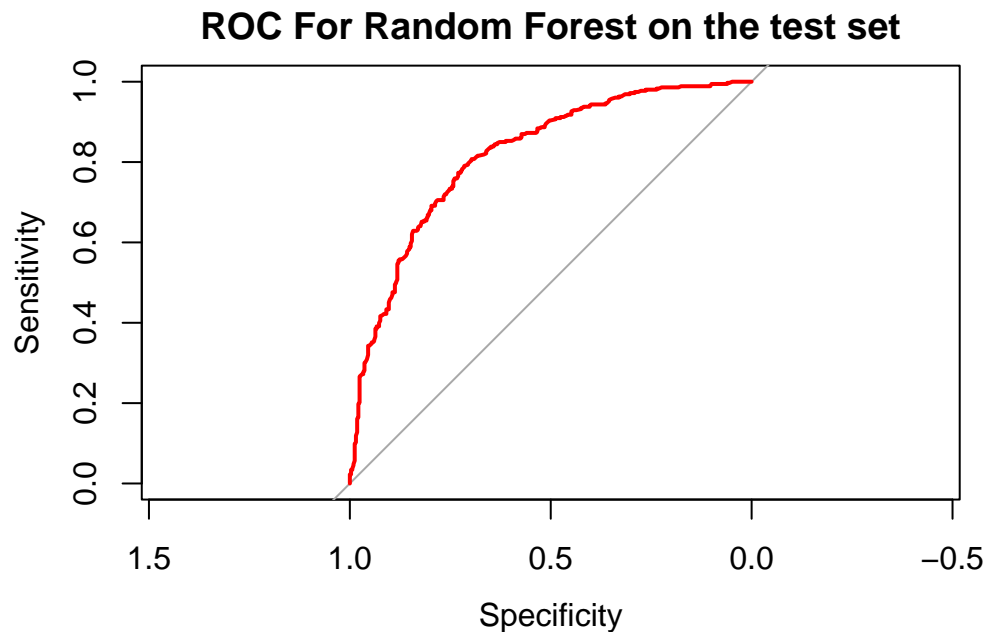
First I built the same random forest model and plot the ROC curve:

```
1 > # Build a random forest model
2 > rf <- randomForest(x = train[-1],
3 +                   y = train$dec_o,
4 +                   ntree = 500,
5 +                   mtry = 4,
6 +                   proximity = TRUE)
7 > # Make "Yes" as positive classes
8 > pred_roc <- predict(rf, newdata = test, type = "prob")[,2]
9 > ROC_rf <- roc(response = test$dec_o,
10 +              predictor = pred_roc)
```

Setting levels: control = No, case = Yes

Setting direction: controls < cases

```
1 > plot(ROC_rf, col = "red", main = "ROC For Random Forest on the test set")
```



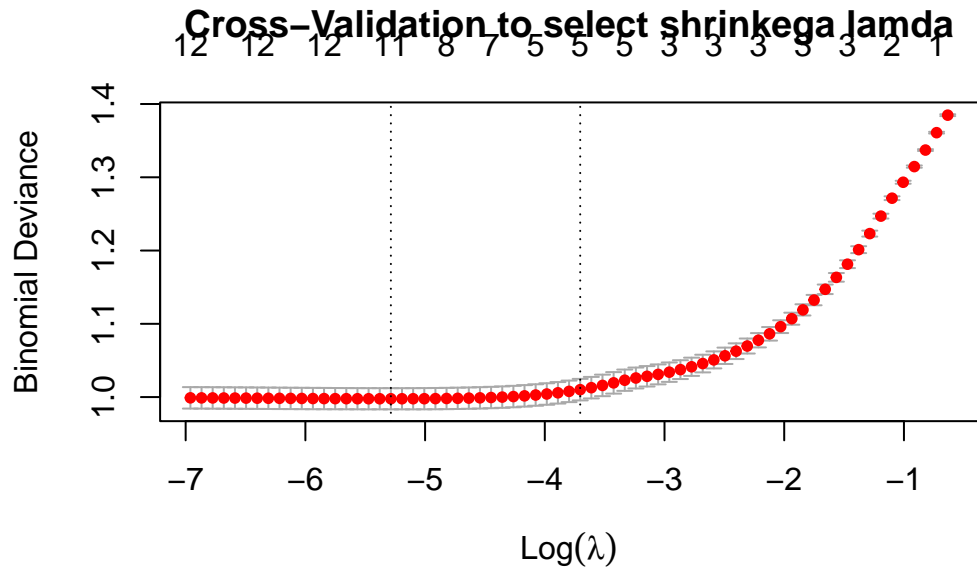
Second I built a Penalized Logistic Regression because there are many predictors, thus selecting and shrinking variables are necessary to build a good logistic regression.

First I encoded train data frame into a form of dummy variables for all categorical variables and I deleted the `from` column because it has 164 unique values which produces huge number of variables and also `field`. Then I use Elastic net with logistics regression on this train matrix with $\alpha = 0.5$. I build a final model with the lambda that gives the simplest model but also lies within one standard error of the optimal lambda selected by cross validation measured by Binomial Deviance. Finally I plot the cross-validation plot when selecting lambda.

```

1 > set.seed(123)
2 > # Encode matrix into dummy variable forms for all categorical variables
3 > x.train <- model.matrix(dec_o~., train[, -c(12,13)]), [-1]
4 > # Use Elastic net with logistics regression on this train dataset with alpha = 0.5
5 > cv.elastic <- cv.glmnet(x = x.train, y = train$dec_o,
6 +                         alpha = 0.5, family = "binomial")
7 > # Build a final model with the best lambda selected by cross validation measured by Binomial
8 > best_elastic <- glmnet(x.train, y = train$dec_o, alpha = 0.5, family = "binomial",
9 +                         lambda = cv.elastic$lambda.1se)
10 > plot(cv.elastic, main = "Cross-Validation to select shrinkage lambda")

```



Next I make predictions on the test set by this best Elastic net model as following:

```
1 > x.test <- model.matrix(dec_o ~., test[, -c(12,13)]), [-1]
2 > prob_elastic <- predict(best_elastic, newx = x.test, type = "response")
```

From this plot, I can tell that Random Forest and Penalized Logistic Regression perform almost equally while Penalized Logistic regression performs slightly better than Random Forest.

```
1 > ROC_lr <- roc(response = test$dec_o,
2 +             predictor = prob_elastic)
```

Setting levels: control = No, case = Yes

Warning in roc.default(response = test\$dec_o, predictor = prob_elastic):
 Deprecated use a matrix as predictor. Unexpected results may be produced, please
 pass a numeric vector.

Setting direction: controls < cases

```
1 > plot(ROC_rf, col = "red", main = "Compare ROC of Random Forest and Penalized Logistic Regression")
2 > lines(ROC_lr, col = "blue")
```


Compare ROC of Random Forest and Penalized Logistic Regression

