

1. Load data, get gender, and create app_proc_time column

Load data

```

data_path = "F:/users/adityadasgupta/Documents/ML%<div>CSC672_project_data/"
applications <- read_parquet(paste0(data_path,"app_data_sample.parquet"))
edges <- read_csv(paste0(data_path,"edges_sample.csv"))

## Rows: 32986 Columns: 4
##   = Column specification
##   delimiter: ", "
## chr (1): application_number
## dbl (2): ego_examiner_id, alter_examiner_id
## date (1): advice_date
##
## # Use lapply() to retrieve the full column specification for this data.
## # Specify the column types or stat_ show_col_types = FALSE to quiet this message.

```

applications

application_number	filing_date	examiner_name_last	examiner_name_first	examiner_name_middle
08284457	2000-01-26	HOWARD	JACQUELINE	V
08413193	2000-10-11	YILDIRIM	BEKIR	L
08531853	2000-05-17	HAMILTON	CYNTHIA	NA
08637752	2001-07-20	MOSHER	MARY	NA
09682726	2000-04-10	BARR	MICHAEL	E
08657412	2000-04-28	GRAY	LINDA	LAMEY
08716371	2004-01-26	MC MILLIAN	KARA	RENTA
08765941	2000-06-23	FORD	VANESSA	L
08776818	2000-02-04	STRZELECKA	TERESA	E
08809677	2002-02-20	KIM	SUN	U

1-10 of 10,000 rows
1-5 of 16 columns

Previous

1

2

3

4

5

6

...1000

Next

edges

application_number	advice_date	ego_examiner_id	alter_examiner_id
09402488	2008-11-17	84356	66266
09402488	2008-11-17	84356	63519
09402488	2008-11-17	84356	98531
09451335	2008-08-21	92953	71313
09451335	2008-08-21	92953	93985
09451335	2008-08-21	92953	91818
09479304	2008-12-15	61767	69277
09479304	2008-12-15	61767	92446
09479304	2008-12-15	61767	66805
09479304	2008-12-15	61767	70919

1-10 of 10,000 rows

Previous

1

2

3

4

5

6

...1000

Next

Get gender for examiners

```

# install.packages("Fuzzy") # Only run this line the first time you use the package, to get data for it
# get a list of first names without repetitions
examiner_names <- applications %>%
  distinct(examiner_name) %>%
  <div>show_col_types()</div> # as shown in the example for the package to attach a gender and probability to
# each name and put the results into the table 'examiner_names_gender'
# get a table of names and gender
examiner_names_gender <- examiner_names %>%
  do(results = c(results$examiner_name_first, method = "ssa")) %>%
  uneset(cols = c(results), keep_empty = TRUE) %>%
  select(
    examiner_name_first = name,
    gender,
    proportion_female
  )
examiner_names_gender

```

examiner_name_first

gender

proportion_female

AARON

male

0.0082

ABDEL

male

0.0000

ABDOU

male

0.0000

ABDUL

male

0.0000

ABDULHAKIM

male

0.0000

ABDULLAH

male

0.0000

ABDULLAH

male

0.0000

ABIGAIL

female

0.9982

ABIMBOLA

female

0.9436

ABRAHAM

male

0.0031

1-10 of 1,822 rows

Previous

1

2

3

4

5

6

...

183

Next

```

# Finally, let's join that table back to our original applications data and discard the temporary tables we have
# just created to reduce clutter in our environment.
# remove extra columns from the gender table
examiner_names_gender <- examiner_names_gender %>%
  select(examiner_name_first, gender) %>%
  # joining gender back to the dataset
applications <- applications %>%
  left_join(examiner_names_gender, by = "examiner_name_first")
# cleaning up
rm(examiner_names)
rm(examiner_names_gender)
gc()

## used (Mb) gc trigger (Mb) limit (Mb) max used (Mb)
## Mcells 4848423 258.6 8092228 432.2 NA 5189797 277.7
## Vcells 85594369 385.4 96616289 737.2 16384 96616180 616.7

```

Guess the examiner's race

```

examiner_surnames <- applications %>%
  select(surname = examiner_name_last) %>%
  distinct()
examiner_surnames

```

surname

HOWARD

YILDIRIM

HAMILTON

MOSHER

BARFAR

GRAY

MC MILLIAN

FORD

STRZELECKA

KIM

1-10 of 3,806 rows

Previous

1

2

3

4

5

6

...

381

Next

```

examiner_race <- predict_race(voter_file = examiner_surnames, surname.only = T) %>%
  as_tibble()

## Warning: unknown or uninitialised column: 'state'.

## Proceeding with last name predictions...

## All local files already up-to-date!

## 791 (18.4%) individuals' last names were not matched.

```

examiner_race

surname	pred.whi	pred.bla	pred.his	pred.asi	pred.oth
HOWARD	0.596663827	0.2948321643	0.0275207004	0.0068983230	0.074084985
YILDIRIM	0.806879507	0.0273312368	0.0694453842	0.0166026000	0.078841272
HAMILTON	0.655944140	0.2386940681	0.0286452612	0.0074999249	0.069216606
MOSHER	0.914757789	0.0042516658	0.0290908584	0.0091684132	0.042731274
BARFAR	0.783554996	0.1198565790	0.0287955895	0.0083027220	0.061490113
GRAY	0.639787385	0.2522555912	0.0280982391	0.0074809225	0.072377863
MC MILLIAN	0.321548611	0.5544059586	0.0211825127	0.0034011419	0.099461776
FORD	0.576284887	0.3202528684	0.0275472214	0.0062098871	0.5544060
STRZELECKA	0.472364100	0.1708059537	0.2200066888	0.0825296368	0.054293531
KIM	0.016912615	0.0028201203	0.0054594943	0.9429391851	0.031868586

1-10 of 3,806 rows

Previous

1

2

3

4

5

6

...

381

Next

```

examiner_race <- examiner_race %>%
  mutate(max_race_p = pmax(pred.asi, pred.bla, pred.his, pred.oth, pred.whi)) %>%
  mutate(race = case_when(
    max_race_p == pred.asi ~ "Asian",
    max_race_p == pred.bla ~ "black",
    max_race_p == pred.his ~ "Hispanic",
    max_race_p == pred.oth ~ "other",
    max_race_p == pred.whi ~ "White",
    TRUE ~ NA_character_
  ))
examiner_race

```

surname

pred.whi

pred.bla

pred.his

pred.asi

pred.oth

max_race_p

race

HOWARD

0.596663827

0.2948321643

0.0275207004

0.0068983230

0.074084985

0.5966638

white

YILDIRIM

0.806879507

0.0273312368

0.0694453842

0.0166026000

0.078841272

0.8068795

white

HAMILTON

0.655944140

0.2386940681

0.0286452612

0.0074999249

0.069216606

0.6559441

white

MOSHER

0.914757789

0.0042516658

0.0290908584

0.0091684132

0.042731274

0.9147578

white

BARFAR

0.783554996

0.1198565790

0.0287955895

0.0083027220

0.061490113

0.7835550

white

GRAY

0.639787385

0.2522555912

0.0280982391

0.0074809225

0.072377863

0.6397874

white

MC MILLIAN

0.321548611

0.5544059586

0.0211825127

0.0034011419

0.099461776

0.5544060

black

FORD

0.576284887

0.3202528684

0.0275472214

0.0062098871

0.5762849

black

STRZELECKA

0.472364100

0.1708059537

0.2200066888

0.0825296368

0.054293531

0.4723642

white

KIM

0.016912615

0.0028201203

0.0054594943

0.9429391851

0.031868586

0.9429392

Asian

1-10 of 3,806 rows

Previous

1

2

3

4

5

6

...

381

Next

```

# Let's join the data back to the applications table.
# removing extra columns
examiner_race <- examiner_race %>%
  select(surname, race)
applications <- applications %>%
  left_join(examiner_race, by = c("examiner_name_last" = "surname"))
rm(examiner_race)
rm(examiner_surnames)
gc()

## used (Mb) gc trigger (Mb) limit (Mb) max used (Mb)
## Mcells 4915646 262.3 8092228 432.2 NA 7587737 480.5
## Vcells 84895477 418.9 96616289 737.2 16384 96616180 737.2

```

Examiner's tenure

```

examiner_dates <- applications %>%
  select(examiner_id, filing_date, appl_status_date)
examiner_dates

```

examiner_id

filing_date

appl_status_date

96082

2000-01-26

30jan2003 00:00:00

87678

2000-10-11

27sep2010 00:00:00

63213

2000-05-17

30mar2009 00:00:00

73788

2001-07-20

07sep2009 00:00:00

77894

2000-04-10

19sep2001 00:00:00

68606

2000-04-28

16jul2001 00:00:00

89557

2004-01-26

15may2017 00:00:00

97543

2000-06-23

03apw2002 00:00:00

98714

2000-02-04

27nov2002 00:00:00

65530

2002-02-20

23mar2009 00:00:00

1-10 of 10,000 rows

Previous

1

2

3

4

5

6

...

1000

Next

```

examiner_dates <- examiner_dates %>%
  mutate(start_date = ymd(filing_date), end_date = as.Date(dmy_hms(appl_status_date)))

examiner_dates <- examiner_dates %>%
  group_by(examiner_id) %>%
  summarise(
    earliest_date = min(start_date, na.rm = TRUE),
    latest_date = max(end_date, na.rm = TRUE),
    tenure_days = interval(earliest_date, latest_date) %>% days(1)
  ) %>%
  filter(year(latest_date)<2018)
examiner_dates

```

examiner_id

earliest_date

latest_date

tenure_days

59012

2004-07-28

2015-07-24

4013

59025

2009-10-26

2017-05-18

2761

59030

2005-12-12

2017-05-23

4179

59040

2007-09-11

2017-05-22

3542

59052

2001-08-21

2007-02-28

2017

59054

2000-11-10

2016-12-23

5887

59055

2004-11-02

2007-12-26

1149

59056

2000-03-24

2017-05-22

6268

59074

2000-01-31

2017-03-17

6255

59081

2011-04-21

2017-05-19

2220

1-10 of 5,625 rows

Previous

1

2

3

4

5

6

...

563

Next

```

# Joining back to the applications data.
applications <- applications %>%
  left_join(examiner_dates, by = "examiner_id")
rm(examiner_dates)
gc()

## used (Mb) gc trigger (Mb) limit (Mb) max used (Mb)
## Mcells 4915646 262.3 8092228 432.2 NA 14880674 768.8
## Vcells 85271493 498.8 139383455 1862.9 16384 1393833176 1862.8

```

Application Processing Time

Clean Data

```

# Remove NAs from status date and gender
applications <- applications %>%
  filter(!is.na(appl_status_date) | !is.na(gender) | !is.na(race))
# Clean Date format
edges$edges_backup
edges <- edges %>%
  mutate(from=ego_examiner_id,to=alter_examiner_id) %>%
  select(from, to) %>%
  drop_na()

#Create Nodes From Edges Data
nodes <- as.data.frame(do.call(rbind,append(as.list(edges$from),as.list(edges$to))))
nodes <- nodes %>%
  mutate(id=v1) %>%
  select(id) %>%
  distinct(id) %>%
  drop_na()

```

```

# Remove all the data we will not need based on application status
exclude_list=c("PENDING")
applications <- applications %>%
  filter(disposal_type %in% exclude_list)

#Setting Gender as factor
applications$gender = as.factor(applications$gender)
#Setting ethnicity as factor
applications$race = as.factor(applications$race)
#Setting disposal_type as factor
applications$disposal_type = as.factor(applications$disposal_type)
#Setting the technology center as a factor
applications$tc = as.factor(applications$tc)

```

1. Create 'app_proc_time'

```

#this is the amount of time in days that the applications take
applications$app_proc_time <- applications$date_time - applications$filing_date
applications$app_proc_time <- as.numeric(applications$app_proc_time)

```

```

##Nodes & Edges Creation First we need to create the network data to calculate centrality We will remove any records that contain NAs to avoid
false issues with coding

```

```

#Create the edges from edge data
edges_backup=edges
edges <- edges %>%
  mutate(from=ego_examiner_id,to=alter_examiner_id) %>%
  select(from, to) %>%
  drop_na()

#Create Nodes From Edges Data
nodes <- as.data.frame(do.call(rbind,append(as.list(edges$from),as.list(edges$to))))
nodes <- nodes %>%
  mutate(id=v1) %>%
  select(id) %>%
  distinct(id) %>%
  drop_na()

```

Closeness Measures

```

We will create 3 closeness measures to the nodes data frame:
1 Degree Centrality: The number of connections (or edges) that each node has.
2 Closeness Centrality: A measure that calculates the ability to
spread information efficiently via the edges that node is connected to. It is calculated as the inverse of the average shortest path between nodes.
3 Betweenness Centrality: A measure that detects a node's influence over the flow of information within a graph.

```

```

g <- graph_from_data_frame(edges, directed = nodes) %>% as_tbl_graph(directed=TRUE)
#now run our first test working
edges <- edges, vertices = FALSE)
g <- g %>%
  activate(nodes) %>%
  mutate(degree_cen = centrality_degree(),
    closeness_cen = centrality_closeness(),
    betweenness_cen = centrality_betweenness()) %>%
  activate(edges)
tg_nodes <-
  g %>%
  activate(nodes) %>%
  data.frame() %>%
  mutate(names.integer(name))
nodes <- nodes %>%
  select(id) %>%
  distinct(id) %>%
  drop_na()

```

```

Time to visualise the degree centralities and numeric data
final_data <- applications %>%
  left_join(nodes,by=c("examiner_id"="id"))

plot <- graph_from_data_frame(edges, vertices = nodes) %>% as_tbl_graph(directed=TRUE)
plot(net, edge_arrow.size=4,vertex.label=NA,vertex.size=8,vertex.color="blue",
  edge.color="green")

```

```

# Degree centrality linear regression model
model_degree <- lm(app_proc_time ~ degree_cen + gender + race + tenure_days, data = final_data)

# Betweenness centrality linear regression model
model_betweenness <- lm(app_proc_time ~ betweenness_cen + gender + race + tenure_days, data = final_data)

# Closeness centrality linear regression model
model_closeness <- lm(app_proc_time ~ closeness_cen + gender + race + tenure_days, data = final_data)

```

```

# Display the model summaries
summary(model_degree)

##
## Call:
## lm(formula = app_proc_time ~ degree_cen + gender + race + tenure_days, data = final_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1575.7   -662.8   -288.1    327.6   4727.8
##
## Coefficients:
## (Intercept)      777.94239      9.78831      88.132 < 2e-16 ***
## degree_cen      -0.38643      0.04804   -7.578 3.51e-14 ***
## gendermale      -4.01289      2.32959   -1.722 0.08903
## raceblack      -56.75969      6.80807   -8.459 < 2e-16 ***
## racehispanic    64.06730      7.64724      8.373 < 2e-16 ***
## raceother      122.21068    31.55291      3.857 0.00015 ***
## racewhite      -6.13950      2.43558   -2.521 0.01219 *
## tenure_days     0.12773      0.00157    81.357 < 2e-16 ***
##
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1617 on 986468 degrees of freedom
## (782248 observations deleted due to missingness)
## Multiple R-squared:  0.007785, Adjusted R-squared:  0.007761
## F-statistic: 1804.19 on 7 and 986468 DF, p-value: < 2.2e-16

```

```

summary(model_betweenness)

##
## Call:
## lm(formula = app_proc_time ~ betweenness_cen + gender + race +
##   tenure_days, data = final_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1767.4   -662.3   -288.1    326.8   4722.5
##
## Coefficients:
## (Intercept)      848.624951      12.813499      69.973 < 2e-16 ***
## degree_cen      -4.054e+00      2.338e+00      2.083 0.0371 *
## raceblack      -5.345e+01      5.999e+00      8.911 < 2e-16 ***
## racehispanic    6.729e+01      7.648e+00      8.834 < 2e-16 ***
## raceother      1.273e+02      2.155e+01      4.205 < 2e-16 ***
## racewhite      -4.723e+00      2.435e+00      1.939 0.0525 .
## tenure_days     0.1279e-01      1.564e-03      81.755 < 2e-16 ***
##
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1617 on 986468 degrees of freedom
## (782248 observations deleted due to missingness)
## Multiple R-squared:  0.008015, Adjusted R-squared:  0.008007
## F-statistic: 1846 on 7 and 986468 DF, p-value: < 2.2e-16

```

```

summary(model_closeness)

##
## Call:
## lm(formula = app_proc_time ~ closeness_cen + gender + race +
##   tenure_days, data = final_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1575.7   -662.8   -288.1    326.8   4722.5
##
## Coefficients:
## (Intercept)      813.204054      12.207231      67.436 < 2e-16 ***
## degree_cen      -13.385628      6.192129      2.162 0.03864 *
## gendermale      11.433802      3.621863      3.157 0.00159 **
## raceblack      -12.323887      7.121369      1.728 0.08548
## racehispanic    21.847666      8.709093      2.417 0.01566 *
## raceother      11.174069      54.142668      0.206 0.83649
## racewhite      -35.839249      2.904477   -12.338 < 2e-16 ***
## tenure_days     0.123189      0.001987      61.976 < 2e-16 ***
## degree_cen:gendermale -0.46938      0.09538   -4.921 8.61e-07 ***
##
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1617 on 986468 degrees of freedom
## (782248 observations deleted due to missingness)
## Multiple R-squared:  0.007795, Adjusted R-squared:  0.007786
## F-statistic: 890.2 on 8 and 986468 DF, p-value: < 2.2e-16

```

```

Get the summary of the linear regressions!
model_degree_interaction <- lm(app_proc_time ~ degree_cen + gender + race + tenure_days, data = final_data)
model_betweenness_interaction <- lm(app_proc_time ~ betweenness_cen + gender + race + tenure_days, data = final_d
ata)
model_closeness_interaction <- lm(app_proc_time ~ closeness_cen + gender + race + tenure_days, data = final_data)
summary(model_degree_interaction)

```

```

##
## Call:
## lm(formula = app_proc_time ~ degree_cen + gender + race + tenure_days,
##   data = final_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1575.7   -662.8   -288.1    326.8   4722.5
##
## Coefficients:
## (Intercept)      813.204054      12.207231      67.436 < 2e-16 ***
## degree_cen      -13.385628      6.192129      2.162 0.03864 *
## gendermale      11.433802      3.621863      3.157 0.00159 **
## raceblack      -12.323887      7.121369      1.728 0.08548
## racehispanic    21.847666      8.709093      2.417 0.0156
```