# Classification

## Part 3

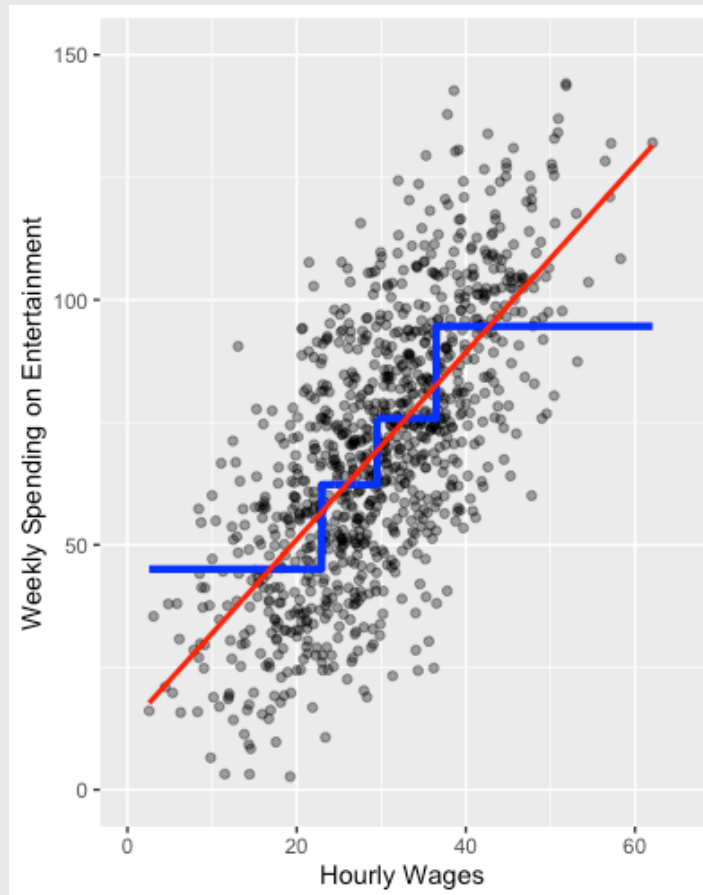Prof. Bisbee

Vanderbilt University

Slides Updated: 2024-03-16

# Agenda

1. Recap of regression and classification

2. Introducing (some) machine learning algorithms
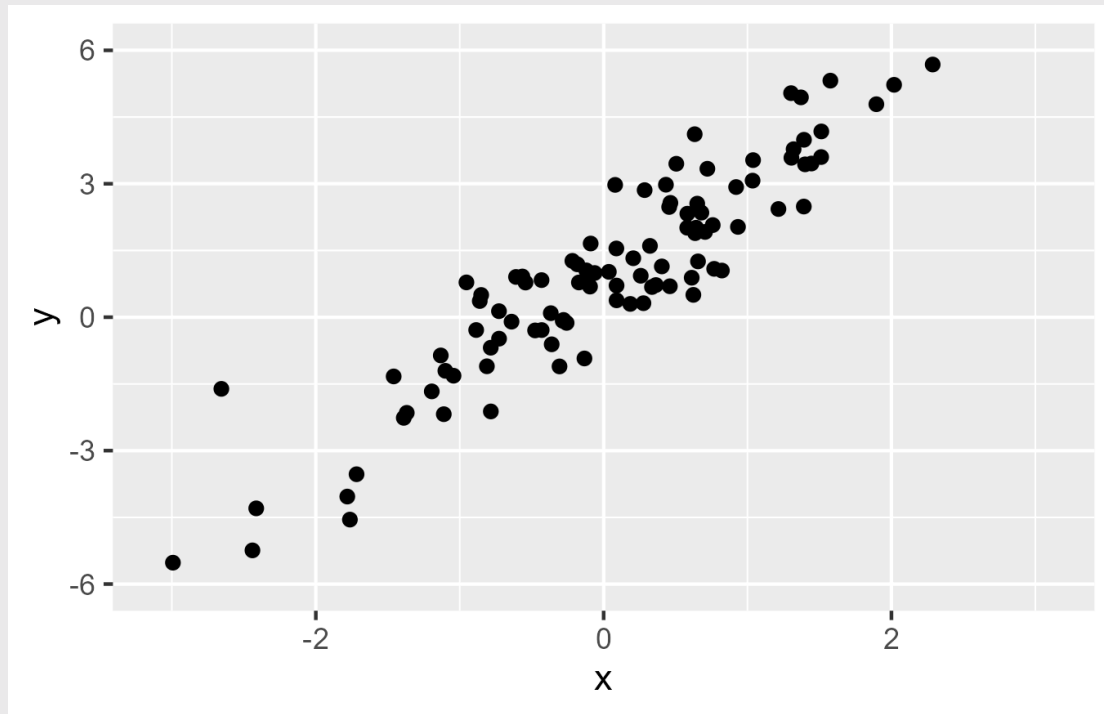
# What is regression?

- Conditional means for continuous data
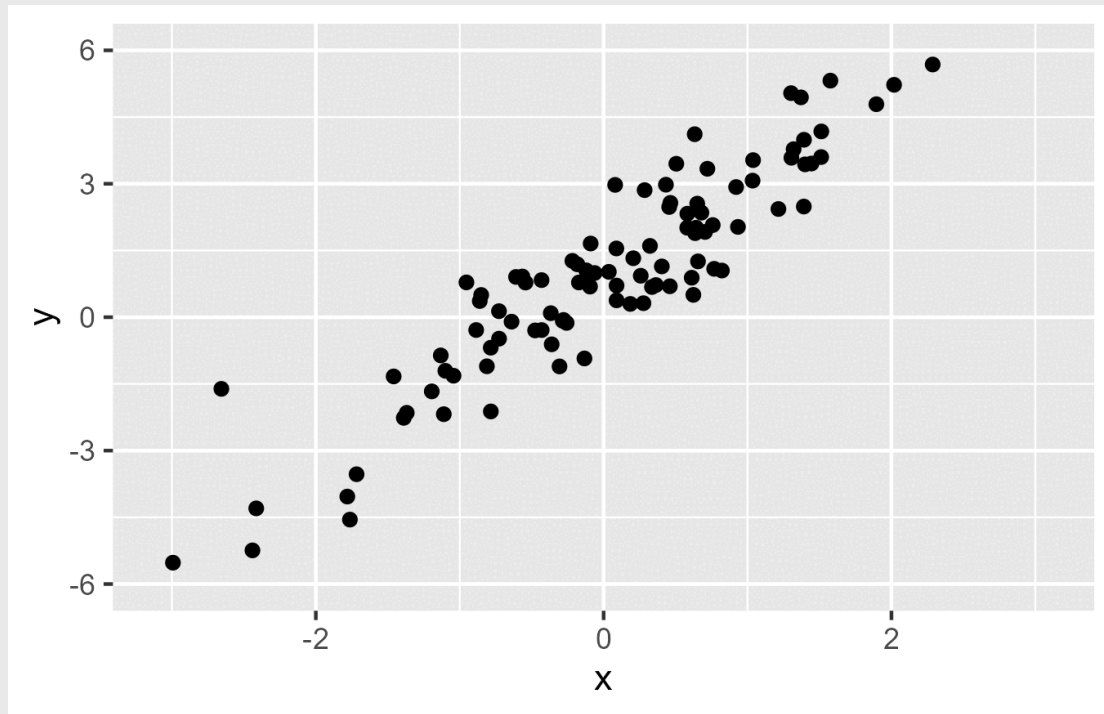
# Regression

- Calculating a **line** that minimizes mistakes *for every observation*

    - NB: could be a curvey line! For now, just assume straight

- Recall from geometry how to graph a straight line

- $Y = a + bX$

    - $a$: the "intercept" (where the line intercepts the y-axis)
    - $b$: the "slope" (how much $Y$ changes for each increase in $X$)

- (Data scientists use $\alpha$ and $\beta$ instead of $a$ and $b$ b/c nerds)

- Regression analysis simply chooses the best line

    - "Best"?

    - The line that minimizes the mistakes (the **line of best fit**)

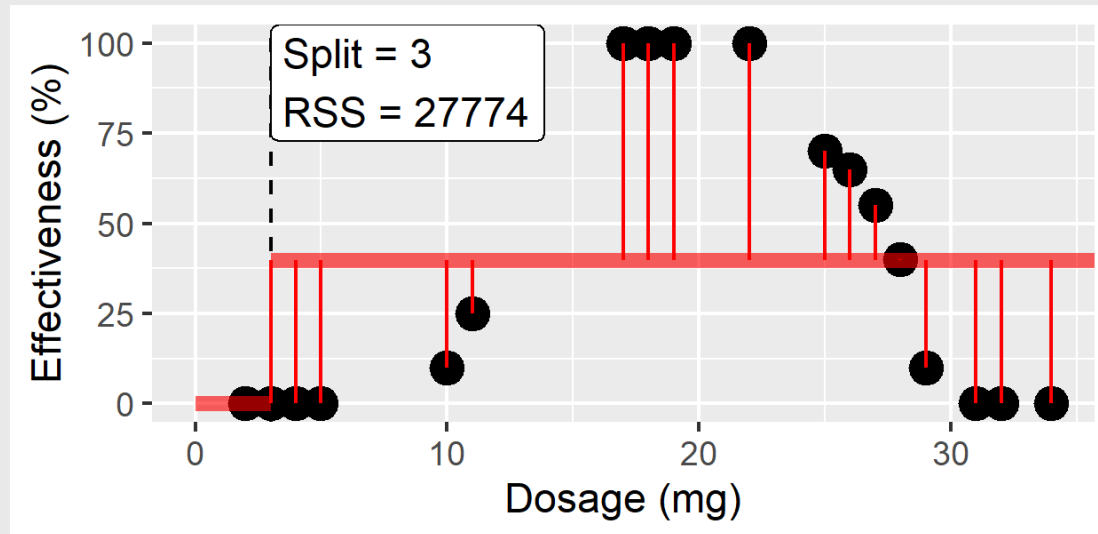# Visual Intuition

# Visual Intuition

# Two Camps Revisited

- Regression is great for **theory testing**

  - Results tell us something **meaningful** about our theory

- But if all we care about is **prediction**...?

  - Want to test every possible predictor (and combinations)

  - Don't care about **relationships**

  - Just care about **accuracy**

- Algorithms can save us time!

  - Random Forests

  - LASSO

# Random Forests

- Identify the best "partition" (split) that divides the data



- In `R`: `ranger`

  - `formula = Y ~ .`

# Random Forests

```r
require(tidyverse)
require(scales)
require(tidymodels)
fn <-
read_rds('C:/Users/bisbeejh/Dropbox/2024_spring/DS1000_S2024/data/fn_cl
```

# Research Question

- What predicts whether you win at Fortnite?

```
form.perf <- 'won ~ hits + assists + accuracy + head_shots +
damage_to_players'

form.games <- 'won ~ eliminations + revives + distance_traveled +
materials_gathered'

form.context <- 'won ~ mental_state + startTime + gameIdSession'

form.full <- 'won ~ hits + assists + accuracy + head_shots +
damage_to_players + eliminations + revives + distance_traveled +
materials_gathered + mental_state + startTime + gameIdSession'
```

# Comparing models

```
m.perf <- lm(as.formula(form.perf),fn)
summary(m.perf)
```

```
##
## Call:
## lm(formula = as.formula(form.perf), data = fn)
##
## Residuals:
##      Min       1Q  Median       3Q      Max
## -0.7905 -0.2756 -0.1563   0.3429   1.0078
##
## Coefficients:
##                      Estimate Std. Error t value Pr(>|t|)
## (Intercept)         8.788e-02  3.768e-02   2.332 0.019893
## hits                6.962e-04  1.001e-03   0.695 0.487053
## assists             3.445e-02  1.020e-02   3.377 0.000764
## accuracy           -4.164e-01  1.081e-01  -3.850 0.000126
## head_shots         -4.808e-03  3.149e-03  -1.527 0.127057
## damage_to_players   4.728e-04  5.713e-05   8.275 4.31e-16
##
## (Intercept)         *
## hits
```

# Comparing models

```
m.games <- lm(as.formula(form.games),fn)
summary(m.games)
```

```
##
## Call:
## lm(formula = as.formula(form.games), data = fn)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -1.0107 -0.2320 -0.1275  0.2028  0.9583
##
## Coefficients:
##                       Estimate Std. Error t value Pr(>|t|)
## (Intercept)          4.171e-02  2.210e-02   1.888 0.059397
## eliminations         1.151e-02  9.765e-03   1.178 0.238921
## revives              6.993e-02  1.809e-02   3.865 0.000119
## distance_traveled    1.805e-04  1.755e-05  10.287  < 2e-16
## materials_gathered  -2.550e-06  3.515e-05  -0.073 0.942186
##
## (Intercept)             .
## eliminations
## revives              ***
```

# Comparing models

```
m.context <- lm(as.formula(form.context),fn)
summary(m.context)
```

```
##
## Call:
## lm(formula = as.formula(form.context), data = fn)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -0.4698 -0.3258 -0.2340  0.5857  0.8553
##
## Coefficients:
##                   Estimate Std. Error t value Pr(>|t|)
## (Intercept)      9.027e+01  2.987e+01   3.022  0.00258
## mental_statesober 1.368e-01  2.933e-02   4.663 3.56e-06
## startTime       -5.672e-08  1.881e-08  -3.015  0.00264
## gameIdSession    1.460e-03  1.463e-03   0.998  0.31863
##
## (Intercept)       **
## mental_statesober ***
## startTime         **
## gameIdSession
```

# Evaluate Model Fit

```r
cvRes <- NULL
for(i in 1:100) {
  inds <- sample(1:nrow(fn),size = round(nrow(fn)*.8),replace = F)
  train <- fn %>% slice(inds)
  test <- fn %>% slice(-inds)

  # Train
  mTmp.perf <- lm(as.formula(form.perf),train)
  mTmp.games <- lm(as.formula(form.games),train)
  mTmp.context <- lm(as.formula(form.context),train)

  # Test
  toEval <- test %>%
    mutate(prob.p = predict(mTmp.perf,newdata = test),
           prob.g = predict(mTmp.games,newdata = test),
           prob.c = predict(mTmp.context,newdata = test),
           truth = factor(won,levels = c('1','0')))

  auc.p <- roc_auc(toEval,truth,prob.p) %>%
    mutate(model = 'performance')

  auc.g <- roc_auc(toEval,truth,prob.g) %>%
    mutate(model = 'games')
```
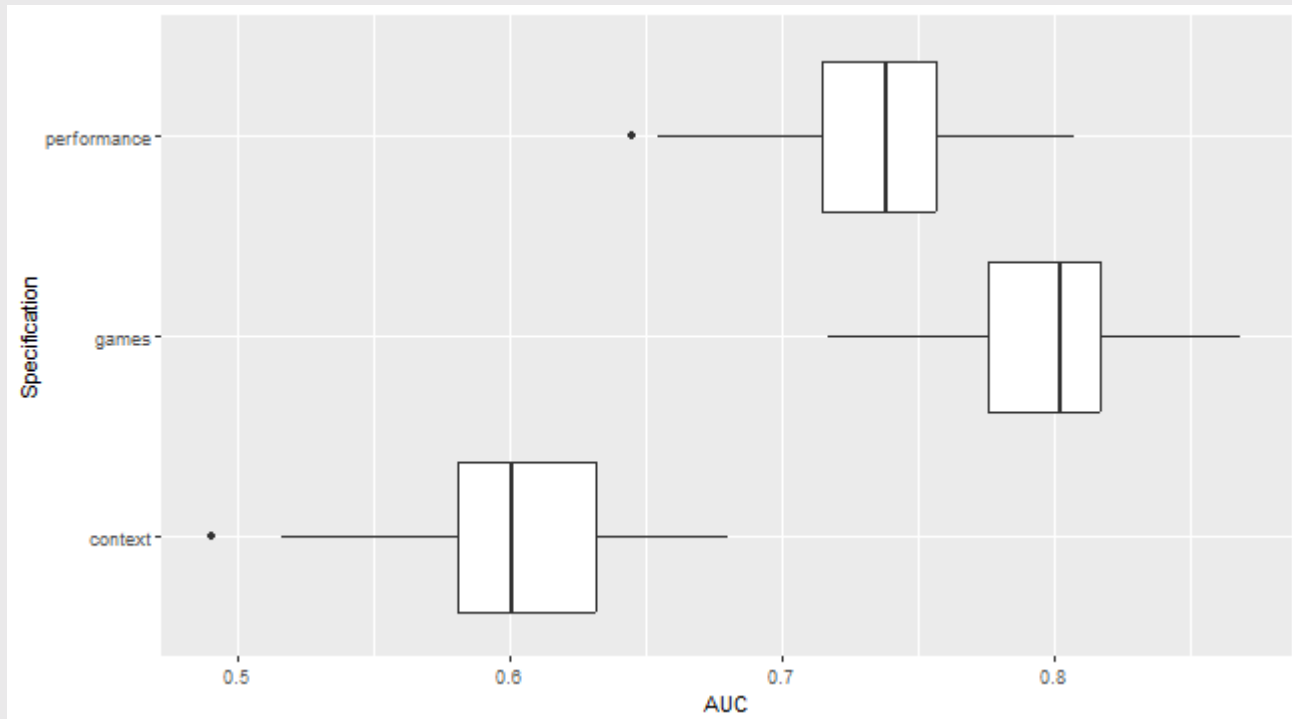
# Evaluate Model Fit

```
cvRes %>%
  ggplot(aes(x = .estimate,y = model)) +
  geom_boxplot() + labs(x = 'AUC',y = 'Specification')
```

# Random Forests

```r
require(ranger) # Fast random forests package
```

```
## Warning: package 'ranger' was built under R version 4.3.3
```

```r
rf.f <- ranger(formula = as.formula(form.full),data = fn)

toEval <- fn %>%
  mutate(prob_won = rf.f$predictions) %>%
  mutate(truth = factor(won,levels = c('1','0')))

roc_auc(toEval,truth,prob_won)
```

```
## # A tibble: 1 × 3
##   .metric .estimator .estimate
##   <chr>   <chr>          <dbl>
## 1 roc_auc binary         0.838
```

# Random Forest Comparison

```r
cvRes <- NULL
for(i in 1:100) {
  inds <- sample(1:nrow(fn),size = round(nrow(fn)*.8),replace = F)
  train <- fn %>% slice(inds)
  test <- fn %>% slice(-inds)

  # Train
  mLM.f <- lm(as.formula(form.full),train)
  mRF.f <- ranger(as.formula(form.full),train)

  # Test
  # NEED TO RUN PREDICTION ON RF FIRST
  tmpPred <- predict(mRF.f,test)

  toEval <- test %>%
    mutate(prob.lm = predict(mLM.f,newdata = test),
           prob.rf = tmpPred$predictions,
           truth = factor(won,levels = c('1','0')))

  auc.lm <- roc_auc(toEval,truth,prob.lm) %>%
    mutate(model = 'linear')

  auc.rf <- roc_auc(toEval,truth,prob.rf) %>%
```
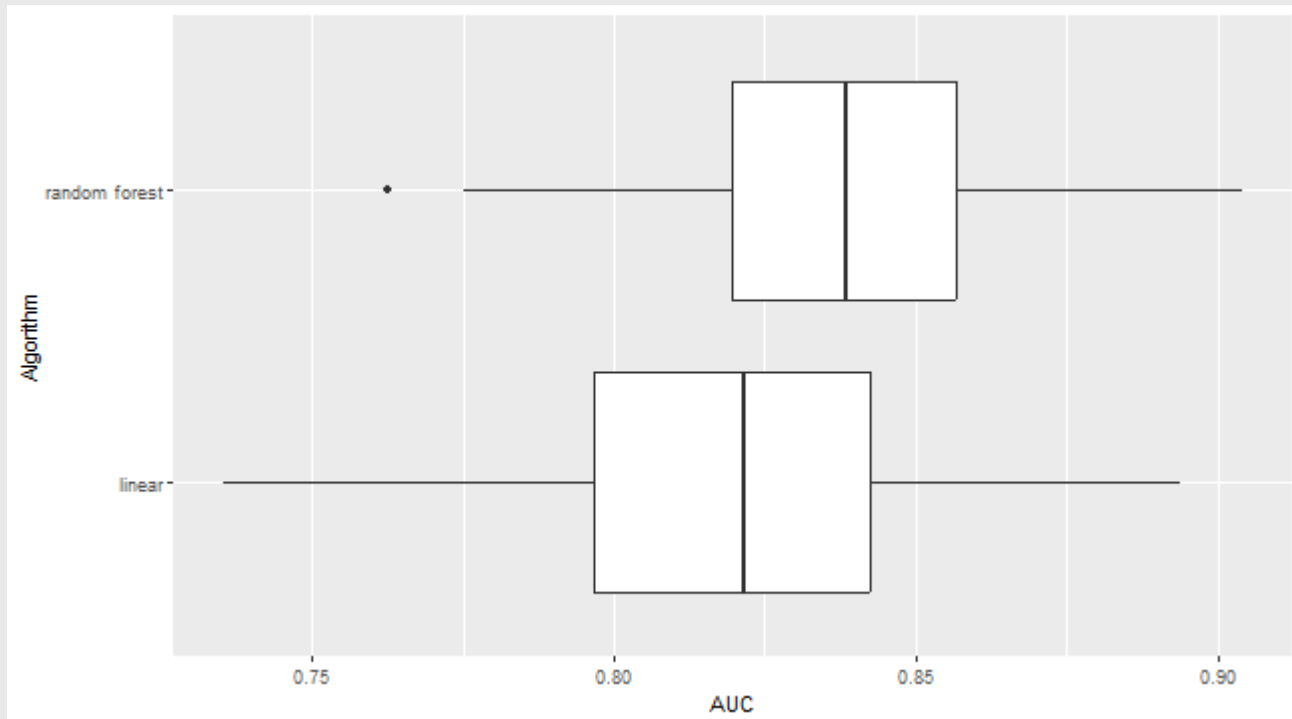
# Random Forest Comparison

```
cvRes %>%
  ggplot(aes(x = .estimate,y = model)) +
  geom_boxplot() + labs(x = 'AUC',y = 'Algorithm')
```

# What matters most?

- Random Forests are particularly suitable for investigating **variable importance**

    - I.e., which $X$ predictors are most helpful?

- A few options, but we rely on **permutation tests**

    - Idea: run the best model you have, then re-run it after "permuting" one of the variables

    - "Permute" means randomly reshuffle...breaks relationship

    - How much **worse** is the model when you break a variable?
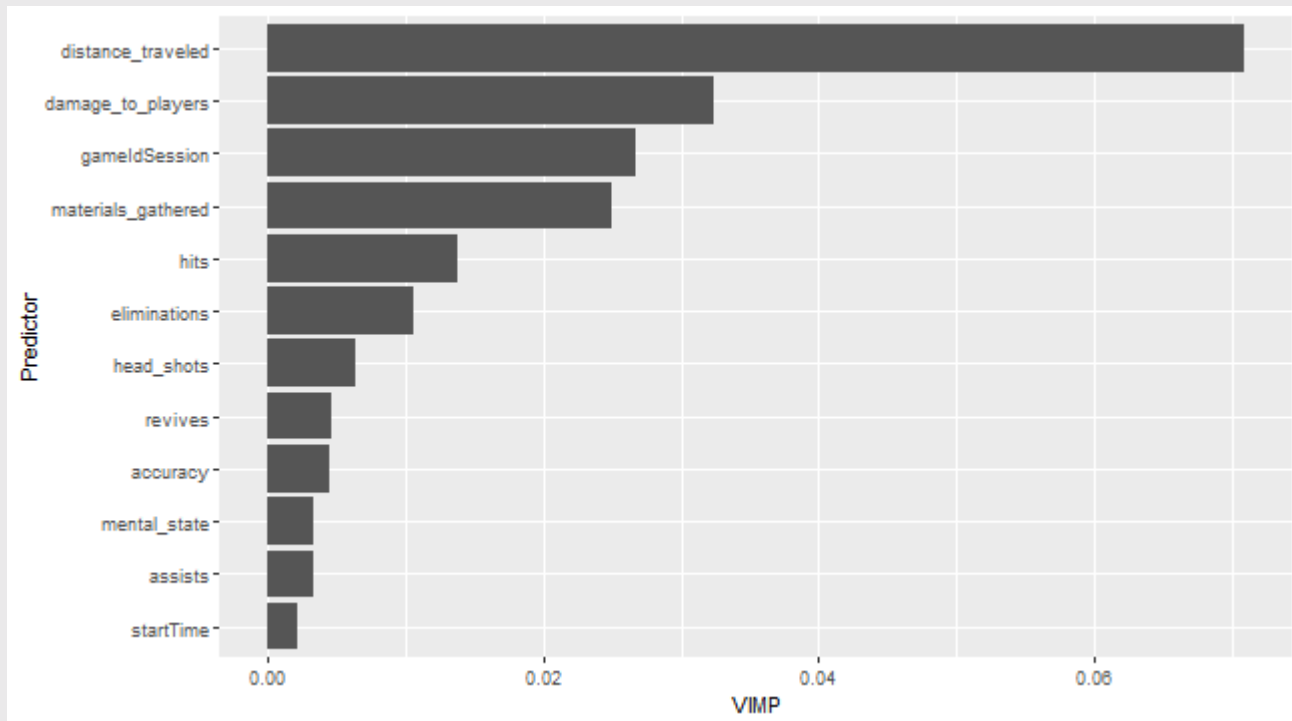
# Variable Importance

- In `ranger()`, use `importance = "permutation"`

```
rf.full <- ranger(formula = as.formula(form.full),data = fn %>%
                  mutate(mental_state = ifelse(mental_state ==
'sober',1,0)),importance = 'permutation')

rf.full$variable.importance
```

```
##                 hits              assists             accuracy
##          0.013744542          0.003357181          0.004465332
##           head_shots    damage_to_players         eliminations
##          0.006363476          0.032384850          0.010595106
##               revives     distance_traveled   materials_gathered
##          0.004715412          0.070833478          0.025052451
##         mental_state            startTime        gameIdSession
##          0.003392765          0.002135140          0.026763760
```

# Variable Importance

```
toplot <- data.frame(vimp = rf.full$variable.importance,
                     vars = names(rf.full$variable.importance))

toplot %>%
  ggplot(aes(x = vimp,y = reorder(vars,vimp))) +
  geom_bar(stat = 'identity') + labs(x = 'VIMP',y = 'Predictor')
```

# LASSO

- "Least Absolute Shrinkage and Selection Operator"

- Concept: Make it hard for predictors to matter

  - Practice: $\lambda$ penalizes how many variables you can include

  - $\sum_{i=1}^{n}(y_i - \sum_j x_{ij}\beta_j)^2 + \lambda \sum_{j=1}^{p} |\beta_j|$

  - Minimize the errors, but penalize for each additional predictor

  - You *could* kitchen-sink a regression and get super low errors

  - LASSO penalizes you from throwing everything into the kitchen sink

- In `R`, need to install a new package! `install.packages('glmnet')`

```
require(glmnet)
```

```
## Warning: package 'glmnet' was built under R version 4.3.2
```

# LASSO

- Function doesn't use formulas

- Give it the raw data instead, divided into Y (outcome) and X (predictors)

```
rhsVars <- str_split(gsub('won ~ ','',form.full),' \\+ ')[[1]]
fn <- fn %>%
  drop_na(all_of(rhsVars))
Y <- fn$won
X <- fn %>% select(all_of(rhsVars)) %>%
  mutate(mental_state = ifelse(mental_state == 'sober',1,0),
         startTime = as.numeric(startTime))
```
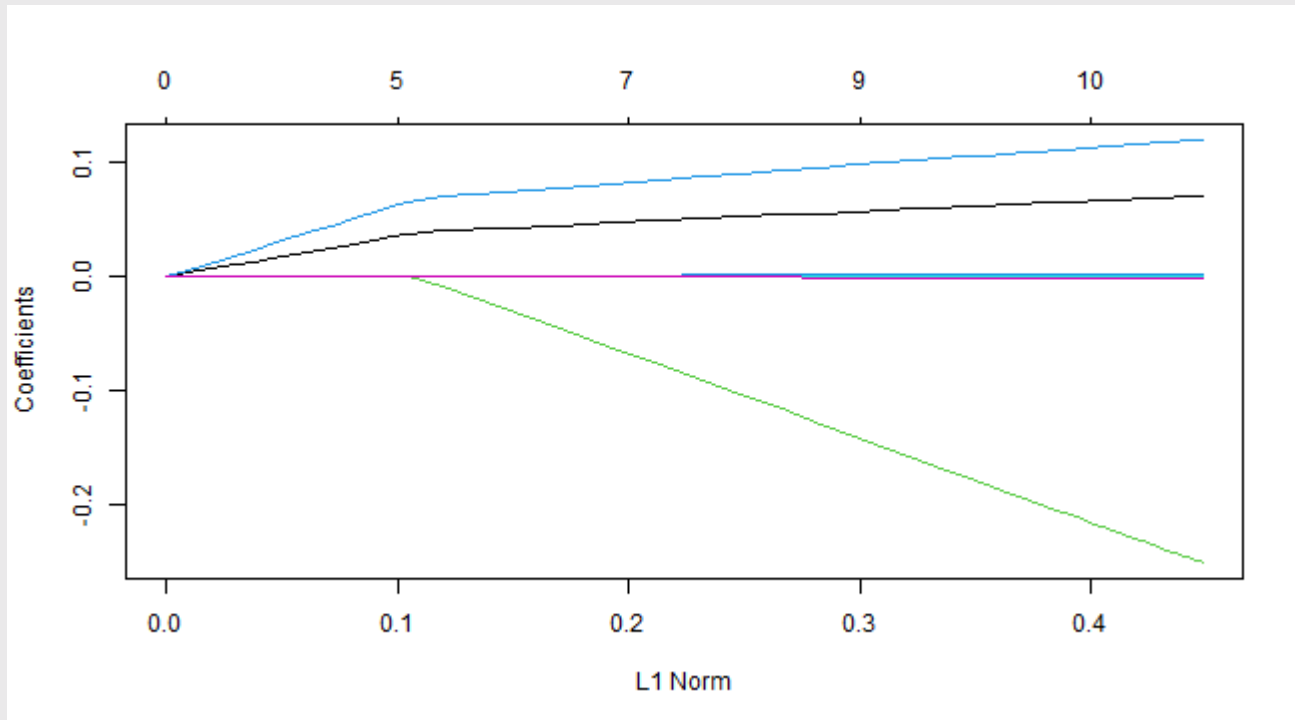
# LASSO

- Now estimate!

```
lassFit <- glmnet(x = as.matrix(X),
                  y = as.matrix(Y))
```
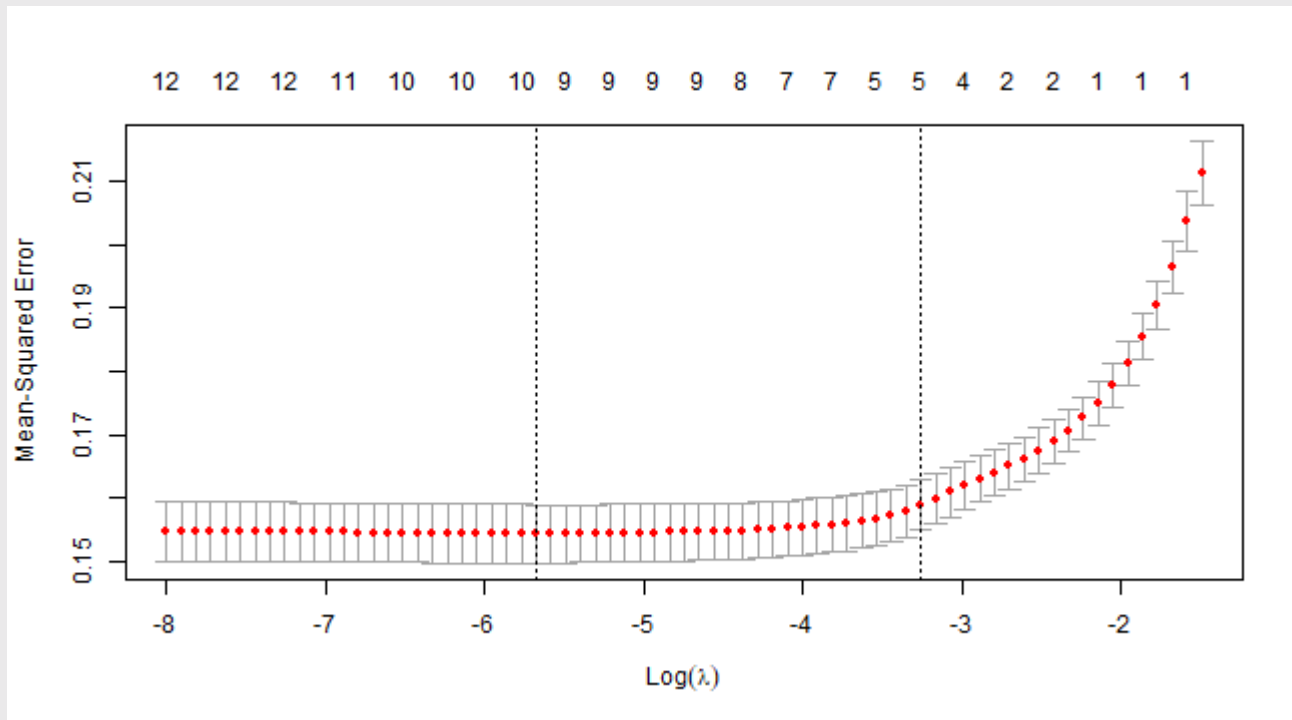
# LASSO
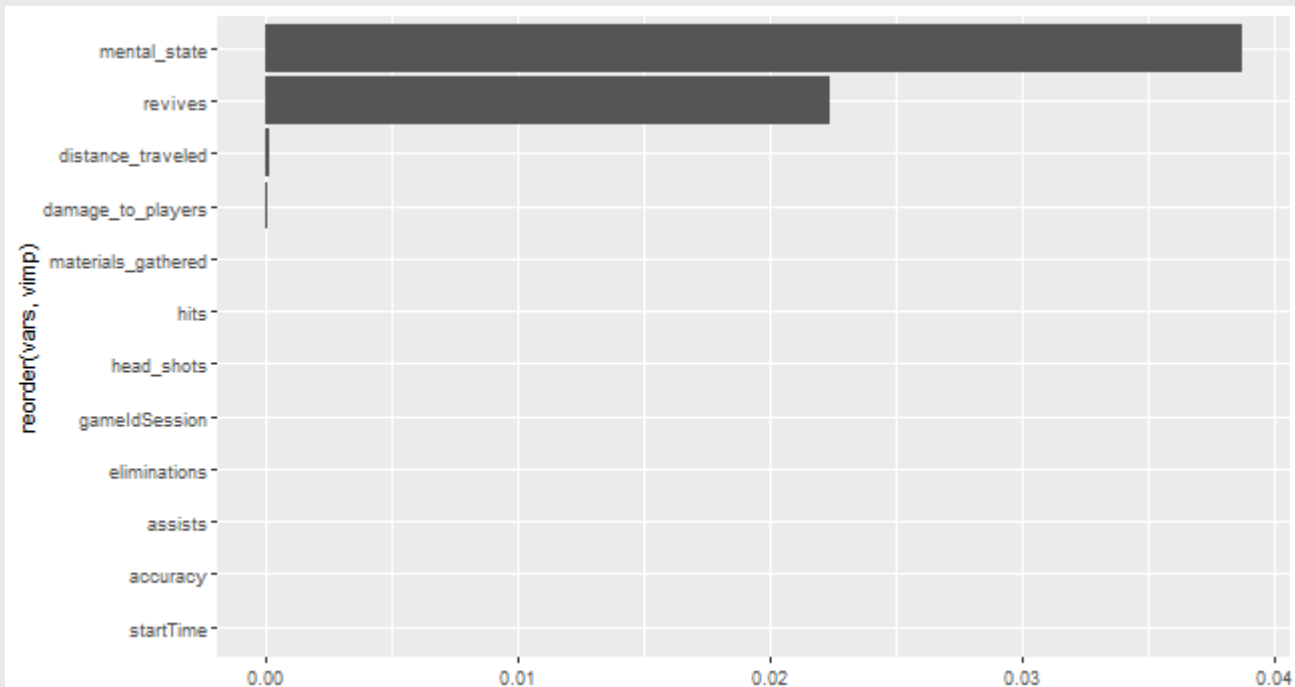
```
plot(lassFit)
```

# Has its own CV!

```
cv.lassFit <- cv.glmnet(x = as.matrix(X),y = as.matrix(Y))

plot(cv.lassFit)
```

# Variable Importance

```
best <- cv.lassFit$glmnet.fit$beta[,cv.lassFit$index[2,]]
vimpLass <- data.frame(vimp = best,
                       vars = names(best))
vimpLass %>%
  ggplot(aes(x = vimp,y = reorder(vars,vimp))) +
  geom_bar(stat = 'identity')
```

# Conclusion

- Lots of powerful tools out there!

- Make sure to take **more classes** on these topics!

- Go to Brightspace and take the **16th** quiz

- **Homework:**

  - Problem Set 9 (due 2024-03-29 by 11:59PM)

  - HW 17