

APPLICATIONS



OF DATA SCIENCE

Modeling in the Tidyverse

Applications of Data Science - Class 5

Giora Simchoni

gsimchoni@gmail.com and add #dsapps in subject

Stat. and OR Department, TAU

2020-03-01

APPLICATIONS



OF DATA SCIENCE

The Problem

APPLICATIONS



OF DATA SCIENCE

Inconsistency, Inextensibility

```
n <- 10000  
x1 <- runif(n)  
x2 <- runif(n)  
t <- 1 + 2 * x1 + 3 * x2  
y <- rbinom(n, 1, 1 / (1 + exp(-t)))
```

```
glm(y ~ x1 + x2, family = "binomial")
```

```
glmnet(as.matrix(cbind(x1, x2)), as.factor(y), family = "binomial")
```

```
randomForest(as.factor(y) ~ x1 + x2)
```

```
gbm(y ~ x1 + x2, data = data.frame(x1 = x1, x2 = x2, y = y))
```



Compare this with sklearn

Detour: A Regression Problem

APPLICATIONS



OF DATA SCIENCE

IPF-Lifts: Predicting Bench Lifting

- Dataset was published as part of the [TidyTuesday](#) initiative
- Comes from [Open Powerlifting](#)
- [Wikipedia](#): Powerlifting is a strength sport that consists of three attempts at maximal weight on three lifts: squat, bench press, and deadlift

The raw data has over 40K rows: for each athlete, for each event, stats about athlete gender, age and weight, and the maximal weight lifted in the 3 types of Powerlifting.

We will be predicting `best3bench_kg` based on a few predictors, no missing values:

```
library(lubridate)

ipf_lifts <- read_csv("https://raw.githubusercontent.com/rfordatascience/tidyverse-tutorial/gh-pages/data/ipf.csv")

ipf_lifts <- ipf_lifts %>%
  drop_na(best3bench_kg, age) %>%
  filter(between(age, 18, 100), best3bench_kg > 0, equipment != "V")
  select(sex, event, equipment, age, division, bodyweight_kg, best3bench_kg)
  drop_na() %>%
  mutate(year = year(date), month = month(date),
        dayofweek = wday(date)) %>%
  select(-date) %>%
  mutate_if(is.character, as.factor)

dim(ipf_lifts)

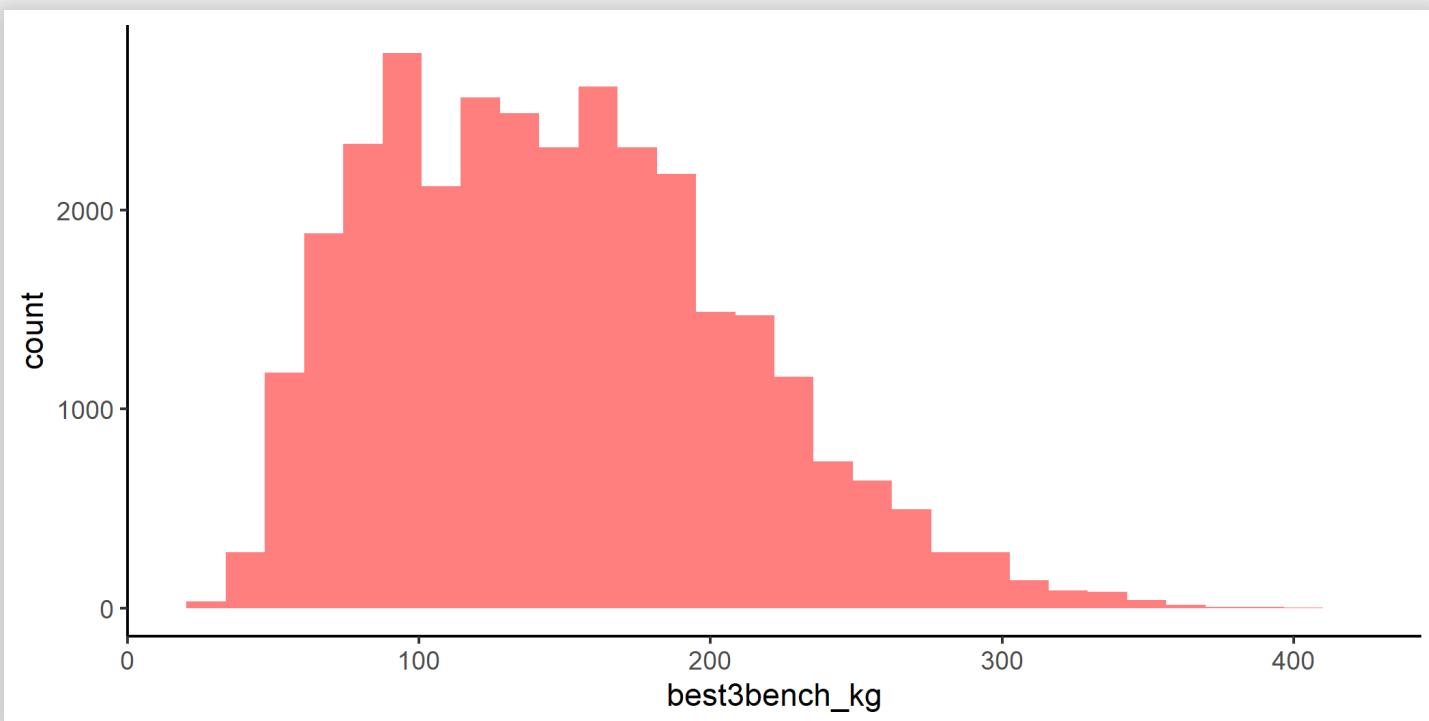
## [1] 32047     11
```

```
glimpse(ipf_lifts)
```

```
## Observations: 32,047
## Variables: 11
## $ sex <fct> F, ...
## $ event <fct> SBD, SBD, SBD, SBD, SBD, SBD, SBD, SBD, SBD...
## $ equipment <fct> Single-ply, Single-ply, Single-ply, Single-ply, ...
## $ age <dbl> 33.5, 34.5, 23.5, 27.5, 37.5, 25.5, 33.5, 26.0, ...
## $ division <fct> Open, Open, Open, Open, Open, Open, Open, Open, ...
## $ bodyweight_kg <dbl> 44, 44, 44, 44, 44, 44, 48, 48, 48, 48, 48, 48, ...
## $ best3bench_kg <dbl> 60.0, 62.5, 62.5, 60.0, 65.0, 45.0, 62.5, 77.5, ...
## $ meet_name <fct> World Powerlifting Championships, World Powerlift...
## $ year <dbl> 1989, 1989, 1989, 1989, 1989, 1989, 1989, 1989, ...
## $ month <dbl> 11, 11, 11, 11, 11, 11, 11, 11, 11, 11, 11, 11, ...
## $ dayofweek <dbl> 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, 2, ...
```

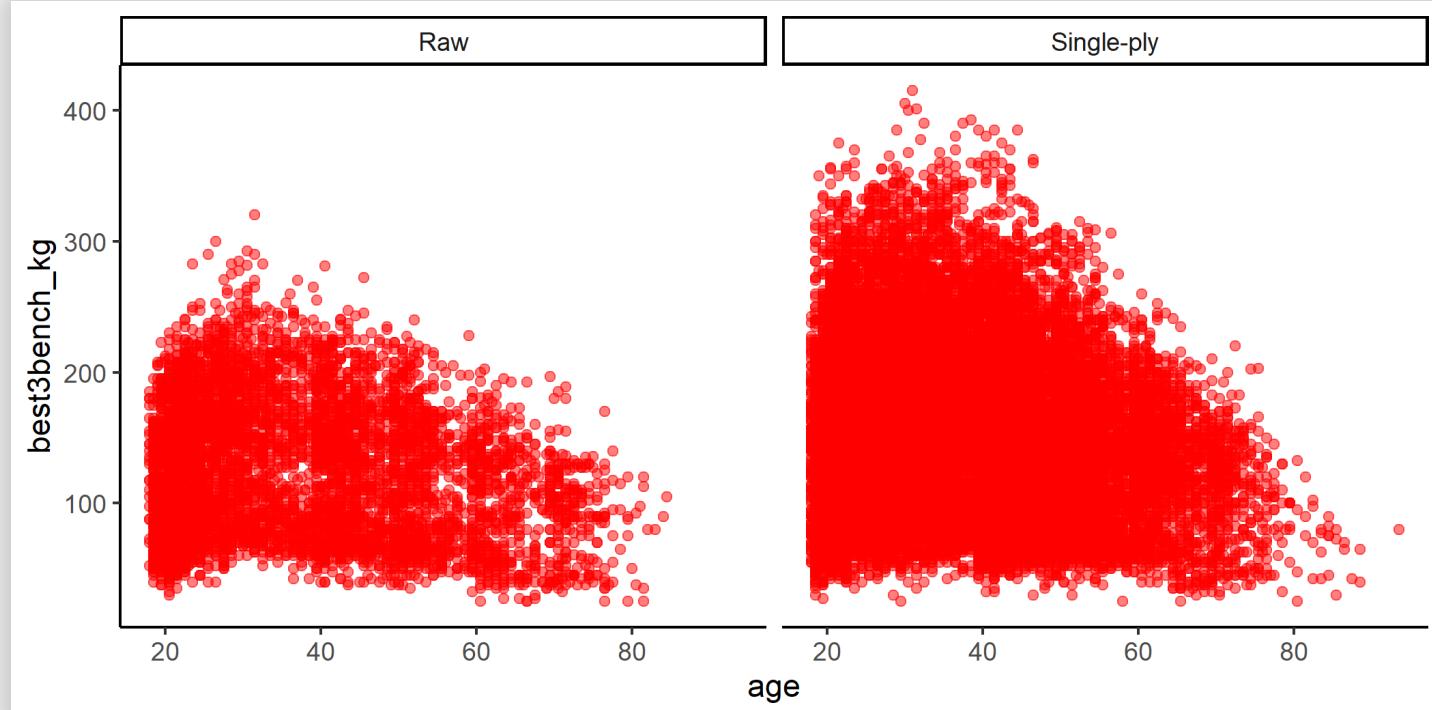
See the dependent variable distribution:

```
ggplot(ipf_lifts, aes(best3bench_kg)) +  
  geom_histogram(fill = "red", alpha = 0.5) +  
  theme_classic()
```



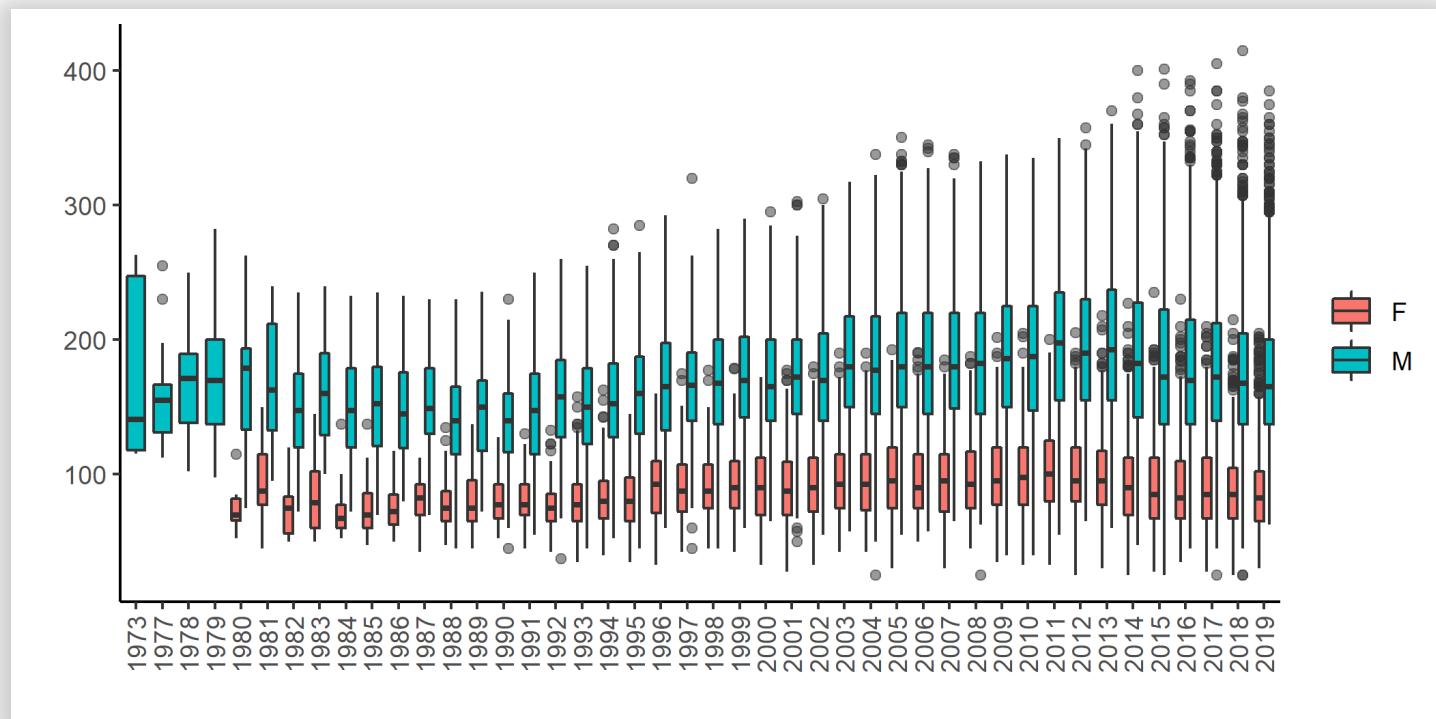
See it vs. say age, faceted by equipment:

```
ggplot(ipf_lifts, aes(age, best3bench_kg)) +  
  geom_point(color = "red", alpha = 0.5) +  
  facet_wrap(~ equipment) +  
  theme_classic()
```



See it vs. year, by gender:

```
ggplot(ipf_lifts, aes(factor(year), best3bench_kg, fill = sex)) +  
  geom_boxplot(outlier.alpha = 0.5) +  
  labs(fill = "", x = "", y = "") +  
  theme_classic() +  
  theme(axis.text.x = element_text(angle = 90, hjust = 1, vjust =
```



Maybe add age^2 and $year^2$ to make Linear Regression's life easier?

```
ipf_lifts <- ipf_lifts %>%
  mutate(age2 = age ^ 2, year2 = year ^ 2)
```

End of Detour

APPLICATIONS



OF DATA SCIENCE

The Present Solution: caret

APPLICATIONS



OF DATA SCIENCE

Split Data

```
library(caret)

train_idx <- createDataPartition(ipf_lifts$best3bench_kg,
                                 p = 0.6, list = FALSE)

ipf_tr <- ipf_lifts[train_idx, ]
ipf_te <- ipf_lifts[-train_idx, ]

library(glue)
glue("train no. of rows: {nrow(ipf_tr)}
      test no. of rows: {nrow(ipf_te)}")

## train no. of rows: 19230
## test no. of rows: 12817
```

Here you might consider some preprocessing.

caret has some nice documentation [here](#).

Tuning and Modeling

Define general methodology, e.g. 10-fold Cross-Validation:

```
fit_control <- trainControl(method = "cv", number = 5)

ridge_grid <- expand.grid(alpha=0, lambda = 10^seq(-3, 1, length =
lasso_grid <- expand.grid(alpha=1, lambda = 10^seq(-3, 1, length =
rf_grid <- expand.grid(splitrule = "variance",
                        min.node.size = seq(10, 30, 10),
                        mtry = seq(2, 10, 2)))

mod_ridge <- train(best3bench_kg ~ ., data = ipf_tr, method = "glm",
                    trControl = fit_control, tuneGrid = ridge_grid,
                    metric = "RMSE")

mod_lasso <- train(best3bench_kg ~ ., data = ipf_tr, method = "glm",
                    trControl = fit_control, tuneGrid = lasso_grid,
                    metric = "RMSE")

mod_rf <- train(best3bench_kg ~ ., data = ipf_tr, method = "ranger",
                  trControl = fit_control, tuneGrid = rf_grid,
                  num.trees = 50, metric = "RMSE")
```

Evaluating Models

```
mod_ridge
```

```
## glmnet
##
## 19230 samples
##     12 predictor
##
## No pre-processing
## Resampling: Cross-Validated (5 fold)
## Summary of sample sizes: 15384, 15383, 15384, 15384, 15385
## Resampling results across tuning parameters:
##
##     lambda      RMSE    Rsquared    MAE
## 0.001000000 26.55832 0.8128887 20.39380
## 0.001206793 26.55832 0.8128887 20.39380
## 0.001456348 26.55832 0.8128887 20.39380
## 0.001757511 26.55832 0.8128887 20.39380
## 0.002120951 26.55832 0.8128887 20.39380
## 0.002559548 26.55832 0.8128887 20.39380
## 0.003088844 26.55832 0.8128887 20.39380
## 0.003727594 26.55832 0.8128887 20.39380
## 0.004498433 26.55832 0.8128887 20.39380
## 0.005428675 26.55832 0.8128887 20.39380
## 0.006551286 26.55832 0.8128887 20.39380
## 0.007000000 26.55832 0.8128887 20.39380
```

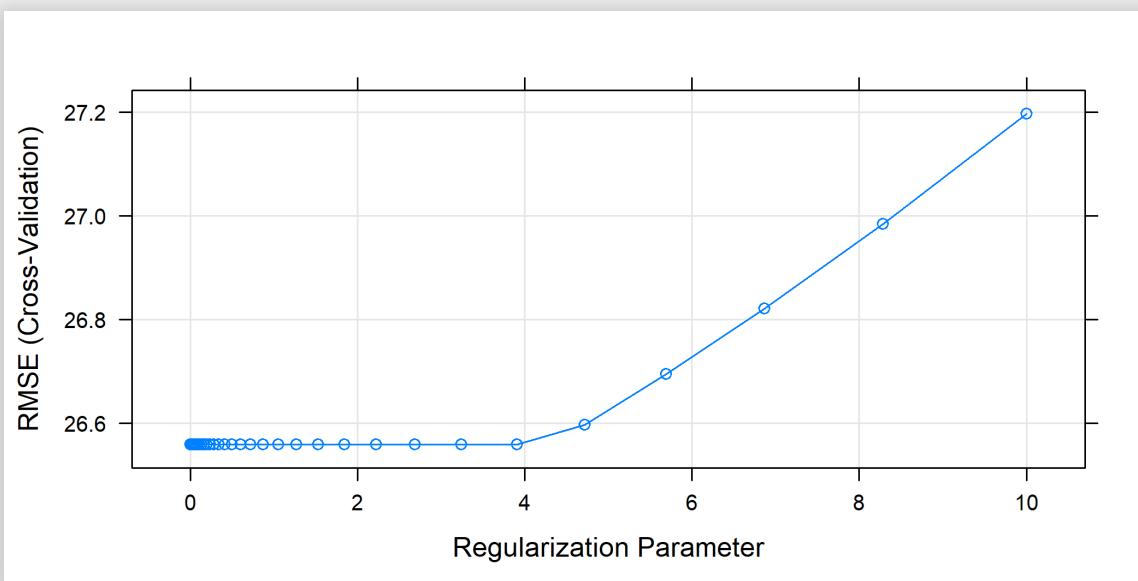
```
mod_lasso
```

```
## glmnet
##
## 19230 samples
##    12 predictor
##
## No pre-processing
## Resampling: Cross-Validated (5 fold)
## Summary of sample sizes: 15383, 15384, 15385, 15384, 15384
## Resampling results across tuning parameters:
##
##     lambda      RMSE    Rsquared     MAE
## 0.001000000 26.11807 0.8171747 20.09521
## 0.001206793 26.11807 0.8171747 20.09521
## 0.001456348 26.11807 0.8171747 20.09521
## 0.001757511 26.11807 0.8171747 20.09521
## 0.002120951 26.11807 0.8171747 20.09521
## 0.002559548 26.11807 0.8171747 20.09521
## 0.003088844 26.11807 0.8171747 20.09521
## 0.003727594 26.11807 0.8171747 20.09521
## 0.004498433 26.11807 0.8171747 20.09521
## 0.005428675 26.11807 0.8171747 20.09521
## 0.006551286 26.11807 0.8171747 20.09521
## 0.007906043 26.11807 0.8171747 20.09521
## 0.009540955 26.11807 0.8171747 20.09521
## 0.011513954 26.11807 0.8171747 20.09521
## 0.013894955 26.11807 0.8171747 20.09521
## 0.016768329 26.11813 0.8171739 20.09523
```

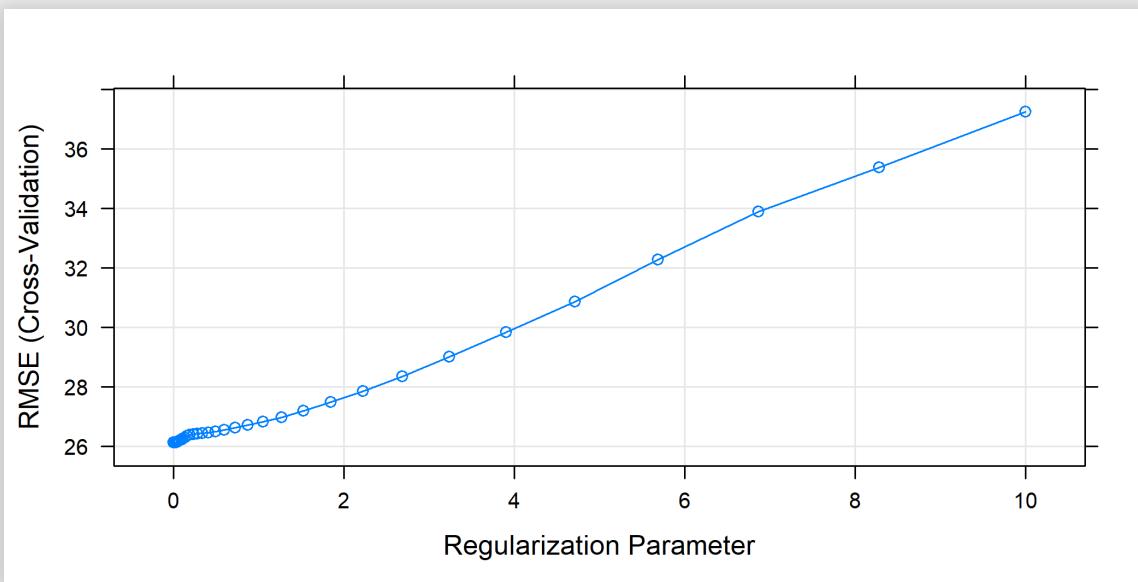
```
mod_rf
```

```
## Random Forest
##
## 19230 samples
##     12 predictor
##
## No pre-processing
## Resampling: Cross-Validated (5 fold)
## Summary of sample sizes: 15383, 15384, 15384, 15383, 15386
## Resampling results across tuning parameters:
##
##     min.node.size  mtry   RMSE      Rsquared    MAE
##     10            2     44.96363  0.6935470  36.14186
##     10            4     33.41223  0.7921971  26.24106
##     10            6     28.19921  0.8220488  21.76780
##     10            8     25.62778  0.8369921  19.65212
##     10           10    24.37503  0.8454751  18.63220
##     20            2     44.07470  0.7191494  35.35157
##     20            4     34.26181  0.7814633  26.98655
##     20            6     28.06338  0.8227133  21.70162
##     20            8     25.68678  0.8370496  19.75861
##     20           10    24.40480  0.8457247  18.67829
##     30            2     44.05877  0.7061486  35.39419
##     30            4     33.32996  0.7882336  26.07677
##     30            6     28.24078  0.8200617  21.81228
##     30            8     25.74167  0.8360185  19.72742
##     30           10    24.39250  0.8459724  18.64230
##
```

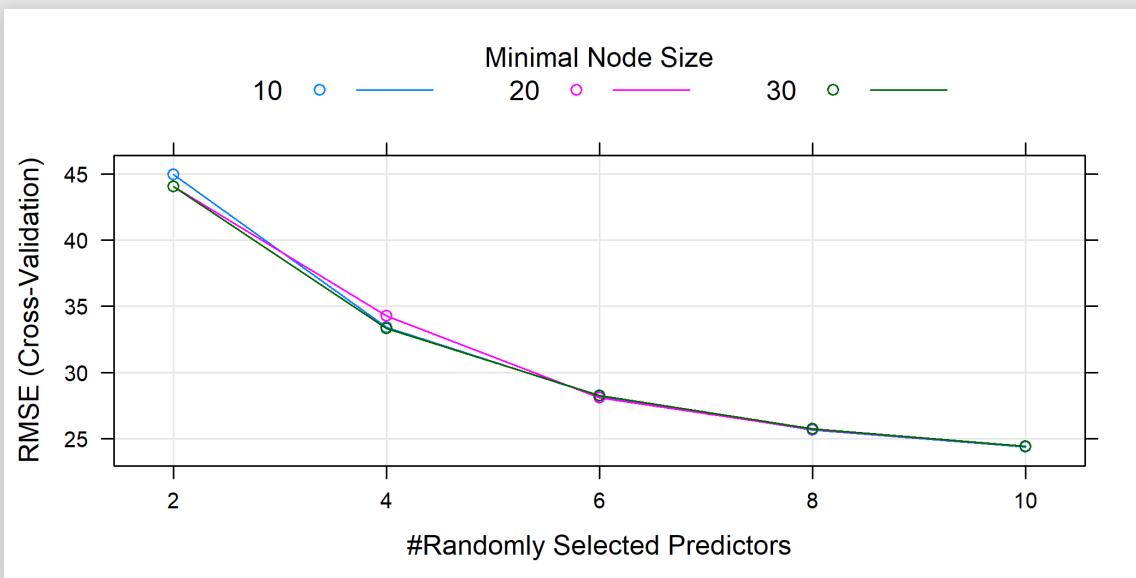
```
plot(mod_ridge)
```



```
plot(mod_lasso)
```



```
plot(mod_rf)
```

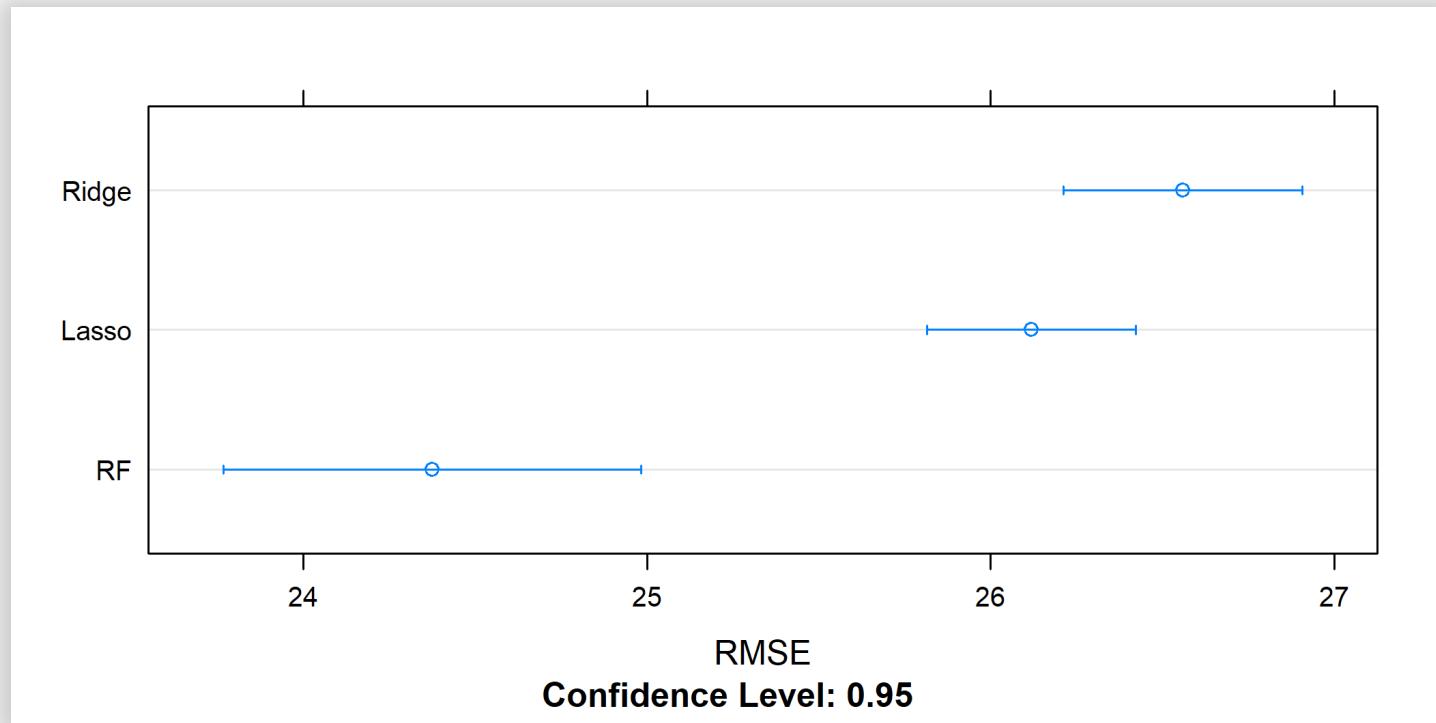


Comparing Models

```
resamps <- resamples(list(Ridge = mod_ridge,
                           Lasso = mod_lasso,
                           RF = mod_rf))
summary(resamps)

##
## Call:
## summary.resamples(object = resamps)
##
## Models: Ridge, Lasso, RF
## Number of resamples: 5
##
## MAE
##      Min. 1st Qu. Median Mean 3rd Qu. Max. NA's
## Ridge 20.16309 20.19137 20.38949 20.39380 20.51260 20.71247 0
## Lasso 19.71590 20.02545 20.12950 20.09521 20.27631 20.32890 0
## RF    18.29181 18.39233 18.67593 18.63220 18.82016 18.98076 0
##
## RMSE
##      Min. 1st Qu. Median Mean 3rd Qu. Max. NA's
## Ridge 26.32026 26.32955 26.47130 26.55832 26.68761 26.98286 0
## Lasso 25.77303 25.96537 26.19022 26.11807 26.31519 26.34653 0
## RF    23.75156 24.08704 24.46089 24.37503 24.52679 25.04886 0
##
```

```
dotplot(resamps, metric = "RMSE")
```



Predicting

```
pred_ridge <- predict(mod_ridge, newdata = ipf_te)
pred_lasso <- predict(mod_lasso, newdata = ipf_te)
pred_rf <- predict(mod_rf, newdata = ipf_te)

rmse_ridge <- RMSE(pred_ridge, ipf_te$best3bench_kg)
rmse_lasso <- RMSE(pred_lasso, ipf_te$best3bench_kg)
rmse_rf <- RMSE(pred_rf, ipf_te$best3bench_kg)

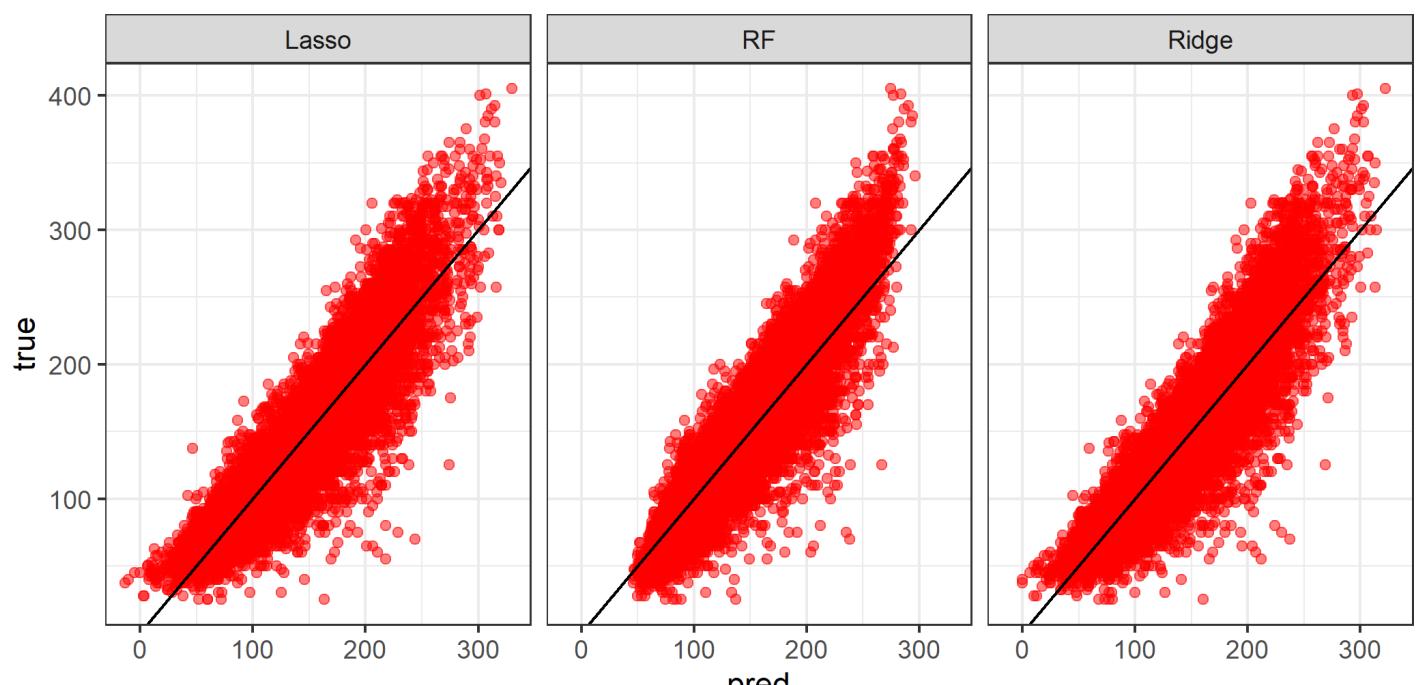
glue("Test RMSE Ridge: {format(rmse_ridge, digits = 3)}
      Test RMSE Lassoe: {format(rmse_lasso, digits = 3)}
      Test RMSE RF: {format(rmse_rf, digits = 3)}")
```

```
## Test RMSE Ridge: 26.9
## Test RMSE Lassoe: 26.4
## Test RMSE RF: 25.1
```

⚠️ Is using a "regular" regression model the natural approach for these data?

Ask yourself what is this model good for, if at all 🤔

```
bind_rows(  
  tibble(method = "Ridge", pred = pred_ridge, true = ipf_te$best3k)  
  tibble(method = "Lasso", pred = pred_lasso, true = ipf_te$best3k)  
  tibble(method = "RF", pred = pred_rf, true = ipf_te$best3bench_1)  
  ggplot(aes(pred, true)) +  
  geom_point(color = "red", alpha = 0.5) +  
  geom_abline(slope = 1, intercept = 0) +  
  facet_wrap(~ method) +  
  theme_bw()
```



The Future Solution: tidymodels

Inspired by [Julia Silge](#)

APPLICATIONS



OF DATA SCIENCE

Packages under `tidymodels`

- `parsnip`: `tidy` `caret`
- `dials` and `tune`: specifying and tuning model parameters
- `rsample`: sampling, data partitioning
- `recipes` and `embed`: preprocessing and creating model matrices
- `infer`: `tidy` statistics
- `yardstick`: measuring models performance
- `broom`: convert models output into tidy tibbles

And [more](#).



All `tidymodels` packages are under development!

Split Data

The `initial_split()` function is from the `rsample` package:

```
library(tidymodels)

ipf_split_obj <- ipf_lifts %>%
  initial_split(prop = 0.6, strata = equipment)

ipf_tr <- training(ipf_split_obj)
ipf_te <- testing(ipf_split_obj)

glue("train no. of rows: {nrow(ipf_tr)}\n"
     "test no. of rows: {nrow(ipf_te)}")
```



```
## train no. of rows: 19229
## test no. of rows: 12818
```



```
print(ipf_split_obj)
```



```
## <19229/12818/32047>
```

Preprocess (but we're not gonna use it)

The `recipe()` function is from the `recipes` package. It allows you to specify a python-like pipe you can later apply to any dataset, including all preprocessing steps:

```
ipf_rec <- recipe(best3bench_kg ~ ., data = ipf_tr)  
ipf_rec
```

```
## Data Recipe  
##  
## Inputs:  
##  
##      role #variables  
##      outcome          1  
##      predictor         12
```

`recipes` contains more preprocessing step_s than you imagine:

```
ipf_rec <- ipf_rec %>%  
  step_normalize(all_numeric())
```

After you have your recipe you need to prep() materials...

```
ipf_rec <- ipf_rec %>% prep(ipf_tr)  
  
ipf_rec
```

```
## Data Recipe  
##  
## Inputs:  
##  
##      role #variables  
##      outcome          1  
## predictor          12  
##  
## Training data contained 19229 data points and no missing data.  
##  
## Operations:  
##  
## Centering and scaling for age, bodyweight_kg, year, month, ... [trained]
```

At this point our recipe has all necessary sd and means for numeric variables.

And then we bake ():

```
ipf_tr2 <- ipf_rec %>% bake(ipf_tr)
ipf_te2 <- ipf_rec %>% bake(ipf_te)

glue("mean of age in orig training: {format(mean(ipf_tr$age), digits = 1)}  
      mean of age in baked training: {format(mean(ipf_tr2$age), digits = 1)}
```



```
## mean of age in orig training: 36.6, sd: 14.3
## mean of age in baked training: 0, sd: 1

glue("mean of age in orig testing: {format(mean(ipf_te$age), digits = 1)}  
      mean of age in baked testing: {format(mean(ipf_te2$age), digits = 1)}
```



```
## mean of age in orig testing: 36.7, sd: 14.3
## mean of age in baked testing: 0, sd: 0.998
```

Or you can do it all in a single pipe:

```
ipf_rec <- recipe(best3bench_kg ~ ., data = ipf_tr) %>%
  step_normalize(all_numeric()) %>%
  prep(ipf_tr)

ipf_tr2 <- ipf_rec %>% bake(ipf_tr)
ipf_te2 <- ipf_rec %>% bake(ipf_te)

glue("mean of age in orig training: {format(mean(ipf_tr$age), digits = 1)}")
     mean of age in baked training: {format(mean(ipf_tr2$age), digits = 1)}

## mean of age in orig training: 36.6, sd: 14.3
## mean of age in baked training: 0, sd: 1

glue("mean of age in orig testing: {format(mean(ipf_te$age), digits = 1)}")
     mean of age in baked testing: {format(mean(ipf_te2$age), digits = 1)}

## mean of age in orig testing: 36.7, sd: 14.3
## mean of age in baked testing: 0, sd: 0.998
```

Modeling

For now let us use the original `ipf_tr` data.

Functions `linear_reg()` and `set_engine()` are from the `parsnip` package:

```
mod_ridge_spec <- linear_reg(mixture = 0, penalty = 0.001) %>%
  set_engine(engine = "glmnet")  
  
mod_ridge_spec
```

```
## Linear Regression Model Specification (regression)
## 
## Main Arguments:
##   penalty = 0.001
##   mixture = 0
## 
## Computational engine: glmnet
```

```
mod_ridge <- mod_ridge_spec %>%
  fit(best3bench_kg ~ ., data = ipf_tr)
```

```
mod_ridge
```

```
## parsnip model object
##
## Fit in: 40ms
## Call: glmnet::glmnet(x = as.matrix(x), y = y, family = "gaussian",
##
##          Df    %Dev Lambda
## 1    51 0.00000 43050
## 2    51 0.00394 39230
## 3    51 0.00433 35740
## 4    51 0.00474 32570
## 5    51 0.00521 29670
## 6    51 0.00571 27040
## 7    51 0.00626 24640
## 8    51 0.00687 22450
## 9    51 0.00754 20450
## 10   51 0.00826 18640
## 11   51 0.00906 16980
## 12   51 0.00994 15470
## 13   51 0.01089 14100
## 14   51 0.01194 12840
## 15   51 0.01309 11700
## 16   51 0.01435 10660
## 17   51 0.01573  9717
## 18   51 0.01723  8853
```

In a single pipe:

```
mod_lasso <- linear_reg(mixture = 1, penalty = 0.001) %>%
  set_engine(engine = "glmnet") %>%
  fit(best3bench_kg ~ ., data = ipf_tr)
```

```
mod_lasso
```

```
## parsnip model object
##
## Fit in: 40ms
## Call: glmnet::glmnet(x = as.matrix(x), y = y, family = "gaussian",
##
##          Df      %Dev Lambda
## 1    0 0.00000 43.050
## 2    1 0.08403 39.230
## 3    2 0.16040 35.740
## 4    2 0.24020 32.570
## 5    2 0.30640 29.670
## 6    2 0.36140 27.040
## 7    2 0.40710 24.640
## 8    2 0.44500 22.450
## 9    2 0.47640 20.450
## 10   2 0.50260 18.640
## 11   2 0.52420 16.980
## 12   2 0.54220 15.470
## 13   2 0.55720 14.100
## 14   3 0.57120 12.840
```

Can also use `fit_xy()` a-la `sklearn`:

```
mod_rf <- rand_forest(mode = "regression", mtry = 4, trees = 50, n  
  set_engine("ranger") %>%  
  fit_xy(x = ipf_tr[, -7],  
         y = ipf_tr$best3bench_kg)  
  
mod_rf
```

```
## parsnip model object  
##  
## Fit in: 671msRanger result  
##  
## Call:  
##   ranger::ranger(formula = formula, data = data, mtry = ~4, num.trees = ~  
##  
##   Type:                           Regression  
##   Number of trees:                 50  
##   Sample size:                     19229  
##   Number of independent variables: 12  
##   Mtry:                            4  
##   Target node size:                30  
##   Variable importance mode:       none  
##   Splitrule:                       variance  
##   OOB prediction error (MSE):     572.301  
##   R squared (OOB):                 0.8471802
```

Notice how easy it is to get the model's results in a tidy way using the `tidy()` function:

```
tidy(mod_ridge)
```

```
## # A tibble: 5,200 x 5
##   term                step estimate lambda dev.ratio
##   <chr>              <dbl>    <dbl>  <dbl>     <dbl>
## 1 (Intercept)          1  1.49e+ 2 43051.  2.56e-36
## 2 sexM                 1  8.47e-35 43051.  2.56e-36
## 3 eventSB               1  2.53e-36 43051.  2.56e-36
## 4 eventSBD               1  -2.59e-35 43051.  2.56e-36
## 5 equipmentSingle-ply      1  2.85e-35 43051.  2.56e-36
## 6 age                   1  -5.30e-37 43051.  2.56e-36
## 7 divisionJuniors        1  -4.55e-36 43051.  2.56e-36
## 8 divisionLight            1  -1.68e-35 43051.  2.56e-36
## 9 divisionMasters 1       1  -2.96e-36 43051.  2.56e-36
## 10 divisionMasters 2      1  -1.39e-35 43051.  2.56e-36
## # ... with 5,190 more rows
```

Predicting

```
results_test <- mod_ridge %>%
  predict(new_data = ipf_te, penalty = 0.001) %>%
  mutate(
    truth = ipf_te$best3bench_kg,
    method = "Ridge"
  ) %>%
  bind_rows(mod_lasso %>%
    predict(new_data = ipf_te) %>%
    mutate(
      truth = ipf_te$best3bench_kg,
      method = "Lasso"
    )) %>%
  bind_rows(mod_rf %>%
    predict(new_data = ipf_te) %>%
    mutate(
      truth = ipf_te$best3bench_kg,
      method = "RF"
    ))
dim(results_test)
```

```
## [1] 38454      3
```

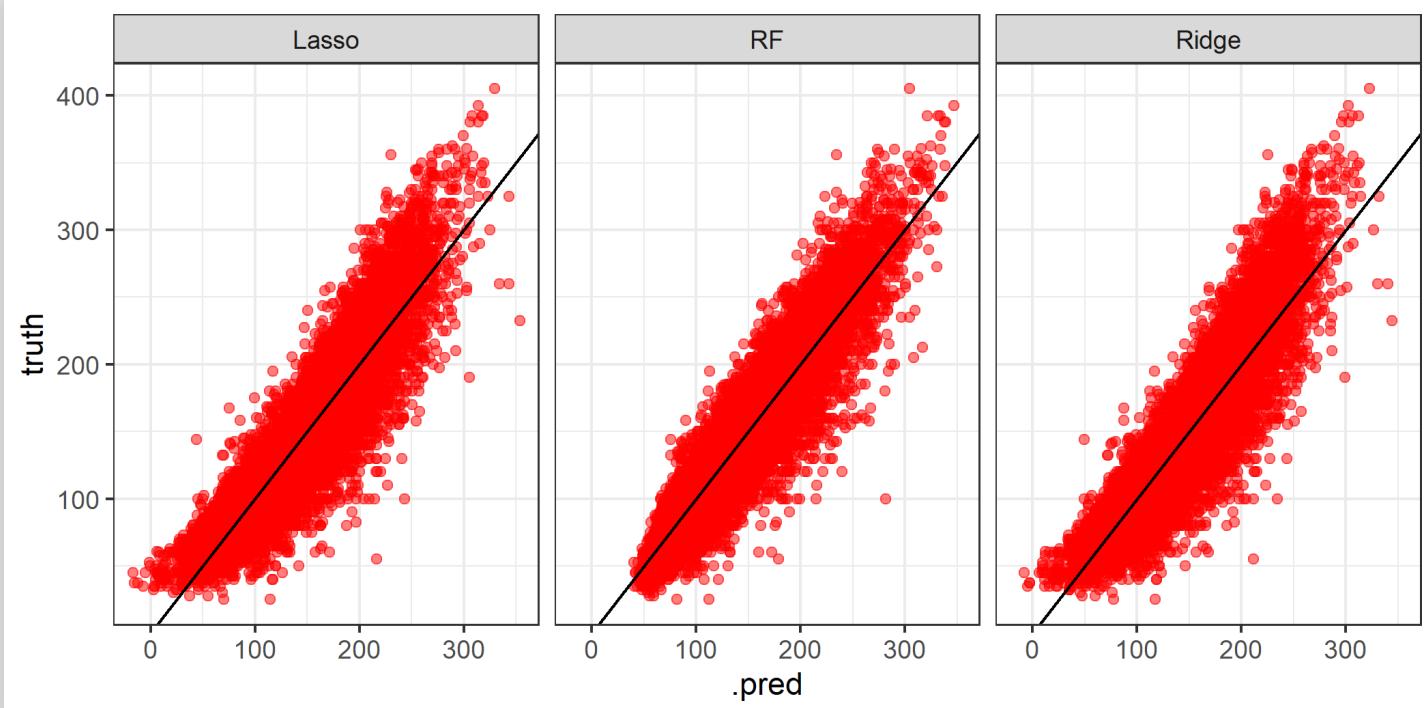
Comparing Models

The package `yardstick` has tons of performance metrics:

```
results_test %>%
  group_by(method) %>%
  yardstick::rmse(truth = truth, estimate = .pred)
```

```
## # A tibble: 3 x 4
##   method .metric .estimator .estimate
##   <chr>   <chr>    <chr>        <dbl>
## 1 Lasso   rmse     standard     26.1
## 2 RF      rmse     standard     23.3
## 3 Ridge   rmse     standard     26.5
```

```
results_test %>%
  ggplot(aes(.pred, truth)) +
  geom_point(color = "red", alpha = 0.5) +
  geom_abline(slope = 1, intercept = 0) +
  facet_wrap(~ method) +
  theme_bw()
```



Tuning

This isn't completely clear to me, but it seems to work:

Define your model spec, using `tune()` from the `tune` package (needs to be installed separately) for a parameter you wish to tune:

```
library(tune)

mod_rf_spec <- rand_forest(mode = "regression",
                           mtry = tune("mtry"),
                           min_n = tune("min_n")) %>%
  set_engine("ranger")
```

Define the grid on which you train your params, with the `dials` package:

```
rf_grid <- grid_regular(mtry(range(2, 10)), min_n(range(10, 30)),  
                         levels = c(5, 3))
```

```
rf_grid
```

```
## # A tibble: 15 x 2  
##       mtry   min_n  
##   <int> <int>  
## 1     2     10  
## 2     4     10  
## 3     6     10  
## 4     8     10  
## 5    10     10  
## 6     2     20  
## 7     4     20  
## 8     6     20  
## 9     8     20  
## 10    10    20  
## 11    2     30  
## 12    4     30  
## 13    6     30  
## 14    8     30  
## 15    10    30
```

Split your data into a few folds for Cross Validation with `vfold_cv()` from the `rsample` package:

```
cv_splits <- vfold_cv(ipf_tr, v = 5)

cv_splits
```

```
## # 5-fold cross-validation
## # A tibble: 5 x 2
##   splits          id
##   <named list>    <chr>
## 1 <split [15.4K/3.8K]> Fold1
## 2 <split [15.4K/3.8K]> Fold2
## 3 <split [15.4K/3.8K]> Fold3
## 4 <split [15.4K/3.8K]> Fold4
## 5 <split [15.4K/3.8K]> Fold5
```

Now perform cross validation with `tune_grid()` from the `tune` package:

```
tune_res <- tune_grid(recipe(best3bench_kg ~ ., data = ipf_tr),
                       model = mod_rf_spec,
                       resamples = cv_splits,
                       grid = rf_grid,
                       metrics = metric_set(rmse))

tune_res
```

```
## # 5-fold cross-validation
## # A tibble: 5 x 4
##   splits              id    .metrics      .notes
##   * <list>            <chr> <list>        <list>
## 1 <split [15.4K/3.8K]> Fold1 <tibble [15 x 5]> <tibble [0 x 1]>
## 2 <split [15.4K/3.8K]> Fold2 <tibble [15 x 5]> <tibble [0 x 1]>
## 3 <split [15.4K/3.8K]> Fold3 <tibble [15 x 5]> <tibble [0 x 1]>
## 4 <split [15.4K/3.8K]> Fold4 <tibble [15 x 5]> <tibble [0 x 1]>
## 5 <split [15.4K/3.8K]> Fold5 <tibble [15 x 5]> <tibble [0 x 1]>
```

Collect the mean metric across folds:

```
estimates <- collect_metrics(tune_res)
```

```
estimates
```

```
## # A tibble: 15 x 7
##       mtry min_n .metric .estimator   mean     n std_err
##       <int> <int> <chr>   <chr>     <dbl> <int>  <dbl>
## 1       2     10  rmse    standard  24.7     5  0.197
## 2       2     20  rmse    standard  24.7     5  0.190
## 3       2     30  rmse    standard  24.8     5  0.192
## 4       4     10  rmse    standard  23.7     5  0.186
## 5       4     20  rmse    standard  23.6     5  0.192
## 6       4     30  rmse    standard  23.5     5  0.181
## 7       6     10  rmse    standard  23.9     5  0.191
## 8       6     20  rmse    standard  23.7     5  0.186
## 9       6     30  rmse    standard  23.6     5  0.186
## 10      8     10  rmse    standard  24.0     5  0.203
## 11      8     20  rmse    standard  23.7     5  0.185
## 12      8     30  rmse    standard  23.6     5  0.185
## 13     10     10  rmse    standard  24.1     5  0.206
## 14     10     20  rmse    standard  23.8     5  0.181
## 15     10     30  rmse    standard  23.7     5  0.183
```

Choose best parameter:

```
estimates %>%
  mutate(min_n = factor(min_n)) %>%
  ggplot(aes(x = mtry, y = mean, color = min_n)) +
  geom_point() +
  geom_line() +
  labs(y = "Mean RMSE") +
  theme_classic()
```

