

Tidymodels

Nate Wells

Math 243: Stat Learning

November 22nd, 2021

Outline

In today's class, we will...

- Discuss the `tidymodels` packages for model building in the `tidyverse` framework

Section 1

Intro to tidymodels

Why tidymodels?

- Suppose we plan to classify data with a binary response variable and want predicted probabilities.

Why tidymodels?

- Suppose we plan to classify data with a binary response variable and want predicted probabilities.
 - Several different models are available:

Function	Package	Code
lda	MASS	<code>predict(object)</code>
glm	stats	<code>predict(object, type = "response")</code>
gbm	gbm	<code>predict(object, type = "response", n.trees)</code>
rpart	rpart	<code>predict(object, type = "prob")</code>
kknn	kknn	<code>kknn(...)\$prob</code>

Why tidymodels?

- Suppose we plan to classify data with a binary response variable and want predicted probabilities.
 - Several different models are available:

Function	Package	Code
lda	MASS	<code>predict(object)</code>
glm	stats	<code>predict(object, type = "response")</code>
gbm	gbm	<code>predict(object, type = "response", n.trees)</code>
rpart	rpart	<code>predict(object, type = "prob")</code>
kknn	kknn	<code>kknn(...)\$prob</code>

- Each method has significantly different methods for making class probability predictions

Why tidymodels?

- Suppose we plan to classify data with a binary response variable and want predicted probabilities.
 - Several different models are available:

Function	Package	Code
lda	MASS	<code>predict(object)</code>
glm	stats	<code>predict(object, type = "response")</code>
gbm	gbm	<code>predict(object, type = "response", n.trees)</code>
rpart	rpart	<code>predict(object, type = "prob")</code>
kknn	kknn	<code>kknn(...)\$prob</code>

- Each method has significantly different methods for making class probability predictions
- Additionally, each model takes in different types of data arguments (vectors, model matrices, data frames, model formulas)

tidymodels goals

Broadly, `tidymodels` presents collection of modeling packages that share design philosophy, syntax and data structure to make it easy to move between packages.

tidymodels goals

Broadly, `tidymodels` presents collection of modeling packages that share design philosophy, syntax and data structure to make it easy to move between packages.

Additionally, `tidymodels` fits in the broader `tidyverse` framework:

- Packages and functions should be accessible and easily interpreted
- Outputs should be data frames (or tibbles) whenever possible
- Functions should be compatible with the `%>%` operator and functional programming
- Model objects should be compatible with `ggplot2`

tidymodels goals

Broadly, `tidymodels` presents collection of modeling packages that share design philosophy, syntax and data structure to make it easy to move between packages.

Additionally, `tidymodels` fits in the broader `tidyverse` framework:

- Packages and functions should be accessible and easily interpreted
- Outputs should be data frames (or tibbles) whenever possible
- Functions should be compatible with the `%>%` operator and functional programming
- Model objects should be compatible with `ggplot2`

`tidymodels` takes the mechanics from each individual model package (`mass`, `tree`, `glm` etc.) and unifies the input and output

The tidymodel framework

- 1 Preprocess data using the `recipes` package
- 2 Create training-test data splits using the `rsample` package
- 3 Give a model a functional form and specify fitting method using the `parsnip` package
- 4 Fit the model, tidy the results, and make predictions using the `fit`, `tidy`, and `predict` functions
- 5 Estimate model performance using cross-validation from the `rsample` package
- 6 Tune model parameters by adding model specifications

The tidymodel framework

- 1 Preprocess data using the `recipes` package
- 2 Create training-test data splits using the `rsample` package
- 3 Give a model a functional form and specify fitting method using the `parsnip` package
- 4 Fit the model, tidy the results, and make predictions using the `fit`, `tidy`, and `predict` functions
- 5 Estimate model performance using cross-validation from the `rsample` package
- 6 Tune model parameters by adding model specifications

We'll investigate each of these in-depth (although slightly out of order)

Section 2

Build a Model

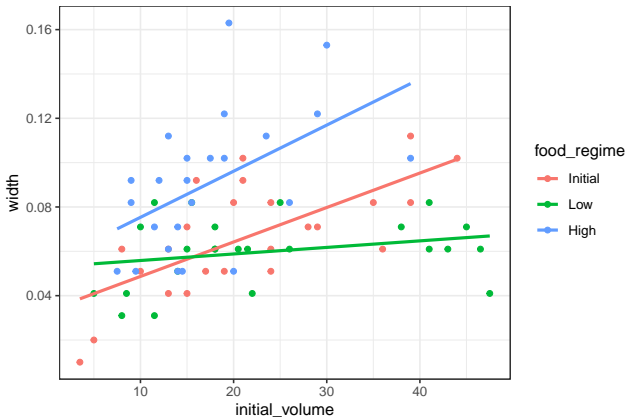
The Data

The `sea_urchins` data set explores the relationship between feeding regimes and size of sea urchins over time:

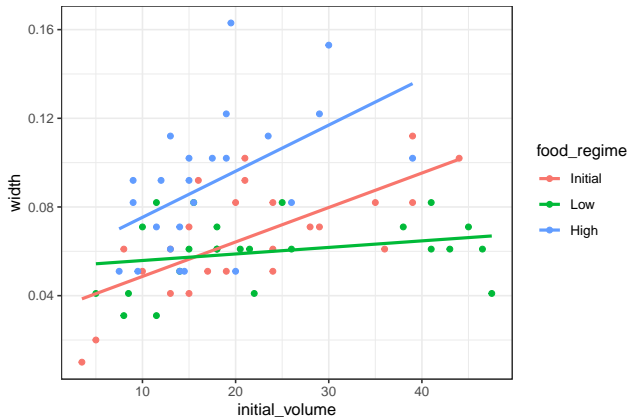
```
sea_urchins<-read_csv("https://tidymodels.org/start/models/urchins.csv") %>%
  setNames(c("food_regime", "initial_volume", "width")) %>%
  mutate(food_regime = factor(food_regime, levels = c("Initial", "Low", "High")))
head(sea_urchins)
```

```
## # A tibble: 6 x 3
##   food_regime initial_volume width
##   <fct>          <dbl> <dbl>
## 1 Initial          3.5  0.01
## 2 Initial          5    0.02
## 3 Initial          8    0.061
## 4 Initial         10    0.051
## 5 Initial         13    0.041
## 6 Initial         13    0.061
```

Scatterplot

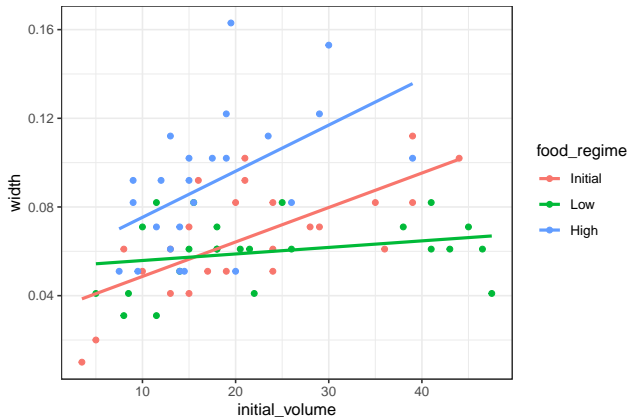


Scatterplot



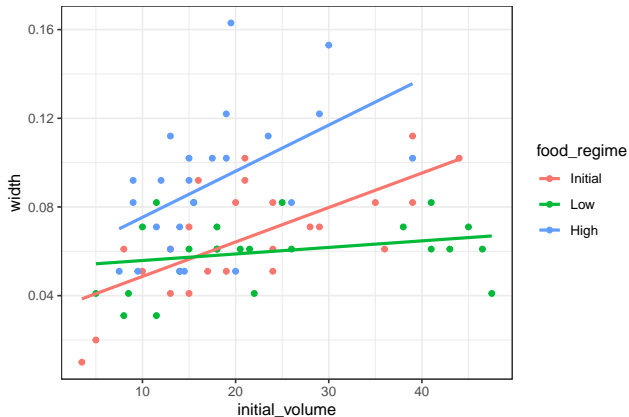
- Goal: Predict `width` as a function of `food_regime` and `initial_volume`.

Scatterplot



- Goal: Predict `width` as a function of `food_regime` and `initial_volume`.
 - Does an additive model seem appropriate?

Scatterplot



- Goal: Predict `width` as a function of `food_regime` and `initial_volume`.
 - Does an additive model seem appropriate?
 - One option might be a linear model with interaction terms.

Build it!

Our model formula takes the form:

```
width ~ initial_volume + food_regime + initial_volume:food_regime
```

Build it!

Our model formula takes the form:

```
width ~ initial_volume + food_regime + initial_volume:food_regime
```

- We need to specify the model's functional form using the `parsnip` package.
- Then specify the method for fitting using `set_engine()`

Build it!

Our model formula takes the form:

```
width ~ initial_volume + food_regime + initial_volume:food_regime
```

- We need to specify the model's functional form using the `parsnip` package.
- Then specify the method for fitting using `set_engine()`

```
library(parsnip)  
linear_reg() %>%  
  set_engine("lm")
```

```
## Linear Regression Model Specification (regression)  
##  
## Computational engine: lm
```

Build it!

Our model formula takes the form:

```
width ~ initial_volume + food_regime + initial_volume:food_regime
```

- We need to specify the model's functional form using the `parsnip` package.
- Then specify the method for fitting using `set_engine()`

```
library(parsnip)  
linear_reg() %>%  
  set_engine("lm")
```

```
## Linear Regression Model Specification (regression)  
##  
## Computational engine: lm
```

- Other engines are possible for `linear_reg()`: `glmnet`, `stan`, and more

Build it!

Our model formula takes the form:

```
width ~ initial_volume + food_regime + initial_volume:food_regime
```

- We need to specify the model's functional form using the `parsnip` package.
- Then specify the method for fitting using `set_engine()`

```
library(parsnip)  
linear_reg() %>%  
  set_engine("lm")
```

```
## Linear Regression Model Specification (regression)  
##  
## Computational engine: lm
```

- Other engines are possible for `linear_reg()`: `glmnet`, `stan`, and more

Now we create the model based on data using the `fit` function:

```
lm_mod<-linear_reg() %>%  
  set_engine("lm")  
  
lm_fit<- lm_mod %>%  
  fit(width ~ initial_volume*food_regime, data = sea_urchins)
```

Results

The output of our `lm_fit` object:

```
lm_fit

## parsnip model object
##
## Fit time: 0ms
##
## Call:
## stats::lm(formula = width ~ initial_volume * food_regime, data = data)
##
## Coefficients:
##              (Intercept)              initial_volume
##              0.0331216                0.0015546
##              food_regimeLow            food_regimeHigh
##              0.0197824                0.0214111
## initial_volume:food_regimeLow  initial_volume:food_regimeHigh
##              -0.0012594                0.0005254
```


Summary Table

To get the traditional summary table:

```
tidy(lm_fit) %>% kable()
```

term	estimate	std.error	statistic	p.value
(Intercept)	0.0331216	0.0096186	3.4434873	0.0010020
initial_volume	0.0015546	0.0003978	3.9077643	0.0002220
food_regimeLow	0.0197824	0.0129883	1.5230864	0.1325145
food_regimeHigh	0.0214111	0.0145318	1.4733993	0.1453970
initial_volume:food_regimeLow	-0.0012594	0.0005102	-2.4685525	0.0161638
initial_volume:food_regimeHigh	0.0005254	0.0007020	0.7484702	0.4568356

Summary Table

To get the traditional summary table:

```
tidy(lm_fit) %>% kable()
```

term	estimate	std.error	statistic	p.value
(Intercept)	0.0331216	0.0096186	3.4434873	0.0010020
initial_volume	0.0015546	0.0003978	3.9077643	0.0002220
food_regimeLow	0.0197824	0.0129883	1.5230864	0.1325145
food_regimeHigh	0.0214111	0.0145318	1.4733993	0.1453970
initial_volume:food_regimeLow	-0.0012594	0.0005102	-2.4685525	0.0161638
initial_volume:food_regimeHigh	0.0005254	0.0007020	0.7484702	0.4568356

Note that the output is a data frame with standard column names

New Data

Suppose we wish to predict the width of 6 sea urchins with `initial_volume` 5 and 30 ml, and with each different `food_regime`.

New Data

Suppose we wish to predict the width of 6 sea urchins with `initial_volume` 5 and 30 ml, and with each different `food_regime`.

- First, we generate data:

New Data

Suppose we wish to predict the width of 6 sea urchins with `initial_volume` 5 and 30 ml, and with each different `food_regime`.

- First, we generate data:

```
new_urchins <- expand.grid(initial_volume = c(5,30),  
                           food_regime = c("Initial", "Low", "High"))  
new_urchins %>% kable()
```

initial_volume	food_regime
5	Initial
30	Initial
5	Low
30	Low
5	High
30	High

Make predictions

Then we make predictions

```
new_preds <- predict(lm_fit, new_data = new_urchins)
conf_int_preds <- predict(lm_fit, new_data = new_urchins, type = "conf_int")
new_preds %>% kable()
```

<u>.pred</u>
0.0408948
0.0797608
0.0543803
0.0617621
0.0649329
0.1169338

```
conf_int_preds %>% kable()
```

<u>.pred_lower</u>	<u>.pred_upper</u>
0.0251382	0.0566514
0.0688612	0.0906605
0.0396403	0.0691204
0.0522641	0.0712601
0.0483265	0.0815393
0.0999144	0.1339532

Combining Data and Predictions

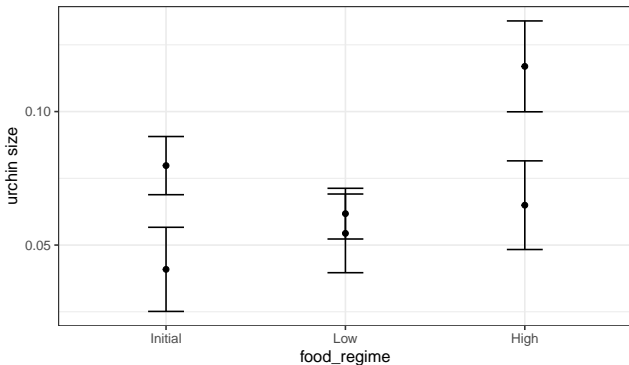
Because the result of `predict()` is tidy, we can easily combine it with the original data:

```
combined_data <- new_urchins %>% cbind(new_preds) %>% cbind(conf_int_preds)
combined_data %>% kable()
```

initial_volume	food_regime	.pred	.pred_lower	.pred_upper
5	Initial	0.0408948	0.0251382	0.0566514
30	Initial	0.0797608	0.0688612	0.0906605
5	Low	0.0543803	0.0396403	0.0691204
30	Low	0.0617621	0.0522641	0.0712601
5	High	0.0649329	0.0483265	0.0815393
30	High	0.1169338	0.0999144	0.1339532

Predictions Plot

```
ggplot(combined_data, aes(x = food_regime)) +  
  geom_point(aes(y = .pred)) +  
  geom_errorbar(aes(ymin = .pred_lower, ymax = .pred_upper), width = .2) +  
  labs(y = "urchin size")+theme_bw()
```



Using a different engine

LASSO?

- With only 3 predictors (`food_regime`, `initial_width` and the interaction term), its unlikely our model will be improved by Penalized Regression. But let's try anyway:

```
glmnet_mod <- linear_reg(penalty = 0.01, mixture = 1) %>% set_engine("glmnet")
```

- `mixture = 1` indicates LASSO (`mixture = 0` is used for Ridge Regression)
- `glmnet` requires us to indicate a value of penalty parameter λ to make predictions.
 - Here, we choose `penalty = 0.01` somewhat arbitrarily (we'll tune later); in any case, `glmnet` will still create models for all λ

Using a different engine

LASSO?

- With only 3 predictors (`food_regime`, `initial_width` and the interaction term), its unlikely our model will be improved by Penalized Regression. But let's try anyway:

```
glmnet_mod <- linear_reg(penalty = 0.01, mixture = 1) %>% set_engine("glmnet")
```

- `mixture = 1` indicates LASSO (`mixture = 0` is used for Ridge Regression)
- `glmnet` requires us to indicate a value of penalty parameter λ to make predictions.
 - Here, we choose `penalty = 0.01` somewhat arbitrarily (we'll tune later); in any case, `glmnet` will still create models for all λ

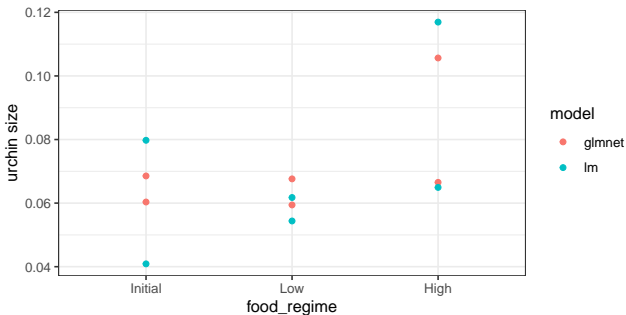
```
glmnet_fit <- glmnet_mod %>% fit(width ~ initial_volume*food_regime, data = sea_urchins)  
tidy(glmnet_fit, penalty = .004) #penalty selects particular value of lambda
```

```
## # A tibble: 6 x 3  
##   term                estimate penalty  
##   <chr>                <dbl>   <dbl>  
## 1 (Intercept)          0.0587  0.004  
## 2 initial_volume      0.000328 0.004  
## 3 food_regimeLow     -0.000918 0.004  
## 4 food_regimeHigh      0         0.004  
## 5 initial_volume:food_regimeLow  0         0.004  
## 6 initial_volume:food_regimeHigh 0.00124  0.004
```

Results from glmnet

```
new_glmnet_preds <- predict(glmnet_fit, new_data = new_urchins, penalty = 0.004)
combined_glmnet_data <- new_urchins %>% cbind(new_glmnet_preds)
two_models <- rbind(combined_glmnet_data,
                    combined_data %>% select(-.pred_lower, -.pred_upper )) %>%
  mutate(model = rep(c("glmnet", "lm"), each = 6))
```

```
ggplot(two_models, aes(x = food_regime)) +
  geom_point(aes(y = .pred, color = model) ) +
  labs(y = "urchin size")+theme_bw()
```



Section 3

Preprocessing with recipes

Recipes

- The `recipes` package assists with preprocessing before a model is trained

Recipes

- The `recipes` package assists with preprocessing before a model is trained
 - Converts qualitative predictors to dummy variables
 - Transforms data to be on a different scale
 - Transforms several predictors at the same time
 - Extracts features from variable

Recipes

- The `recipes` package assists with preprocessing before a model is trained
 - Converts qualitative predictors to dummy variables
 - Transforms data to be on a different scale
 - Transforms several predictors at the same time
 - Extracts features from variable
- The main advance of `recipes` is that it allows us combine several steps at once, in a reproducible fashion

House Prices

- The house data contains information on 30 predictors for 200 houses in Ames, Iowa

```
names(house)
```

```
## [1] "SalePrice"      "Id"             "Functional"     "BldgType"
## [5] "Foundation"    "LotShape"      "LandSlope"     "SaleCondition"
## [9] "RoofMatl"      "ScreenPorch"   "MSSubClass"    "GarageCars"
## [13] "BedroomAbvGr" "TotalBsmtSF"   "LotArea"       "OpenPorchSF"
## [17] "BsmtFullBath" "WoodDeckSF"    "OverallCond"   "YrSold"
## [21] "GrLivArea"     "MoSold"        "TotRmsAbvGrd"  "PoolArea"
## [25] "YearBuilt"     "GarageArea"    "OverallQual"   "Fireplaces"
## [29] "EnclosedPorch" "FullBath"      "HalfBath"
```


House Prices

- The house data contains information on 30 predictors for 200 houses in Ames, Iowa

```
names(house)
```

```
## [1] "SalePrice"      "Id"              "Functional"      "BldgType"
## [5] "Foundation"    "LotShape"        "LandSlope"      "SaleCondition"
## [9] "RoofMatl"      "ScreenPorch"    "MSSubClass"     "GarageCars"
## [13] "BedroomAbvGr" "TotalBsmtSF"    "LotArea"        "OpenPorchSF"
## [17] "BsmtFullBath"  "WoodDeckSF"     "OverallCond"    "YrSold"
## [21] "GrLivArea"     "MoSold"         "TotRmsAbvGrd"  "PoolArea"
## [25] "YearBuilt"     "GarageArea"     "OverallQual"    "Fireplaces"
## [29] "EnclosedPorch" "FullBath"       "HalfBath"
```

- Note that the variable `Id` is not useful as a predictor, but is useful for referring to houses in the data set.

Investigate Predictors

- Additionally, note that several of the variables are factors, so should be converted to a collection of dummy variables.

Investigate Predictors

- Additionally, note that several of the variables are factors, so should be converted to a collection of dummy variables.
- Moreover, for a few variables, some levels are very underrepresented.

```
house %>% count(RoofMatl)
```

```
##   RoofMatl   n
## 1 CompShg 195
## 2 Membran   1
## 3 Tar&Grv   2
## 4 WdShake   1
## 5 WdShngl   1
```

Data Splitting

- We can use the `rsample` package to create a test-training split

Data Splitting

- We can use the `rsample` package to create a test-training split
 - The `rsample` package allows us to create stratified samples in addition to simple random samples

Data Splitting

- We can use the `rsample` package to create a test-training split
 - The `rsample` package allows us to create stratified samples in addition to simple random samples

```
library(rsample)
set.seed(1221)
data_split <- initial_split(house , prop = 3/4)
train_data <- training(data_split)
test_data <- testing(data_split)
```

Create a recipe and update roles

- We now create a recipe for some data pre-processing

```
library(recipes)
house_rec <-
  recipe(SalePrice ~ ., data = train_data) %>%
  update_role(Id, new_role = "ID")
```

Create a recipe and update roles

- We now create a recipe for some data pre-processing

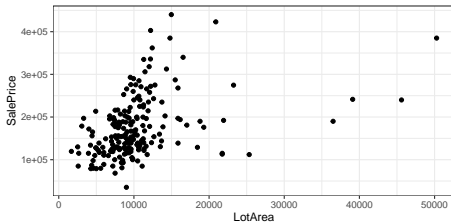
```
library(recipes)
house_rec <-
  recipe(SalePrice ~ ., data = train_data) %>%
  update_role(Id, new_role = "ID")
```

```
summary(house_rec)
```

```
## # A tibble: 31 x 4
##   variable      type      role      source
##   <chr>         <chr>    <chr>    <chr>
## 1 Id            numeric ID      original
## 2 Functional    nominal predictor original
## 3 BldgType      nominal predictor original
## 4 Foundation    nominal predictor original
## 5 LotShape      nominal predictor original
## 6 LandSlope     nominal predictor original
## 7 SaleCondition nominal predictor original
## 8 RoofMatl      nominal predictor original
## 9 ScreenPorch   numeric predictor original
## 10 MSSubClass   numeric predictor original
## # ... with 21 more rows
```

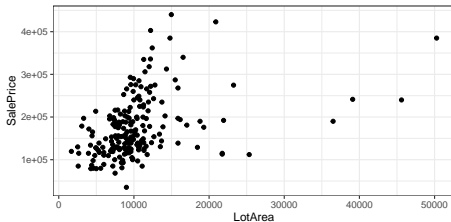

Add steps to recipes

- Consider the relationship between of sale price and lot area:

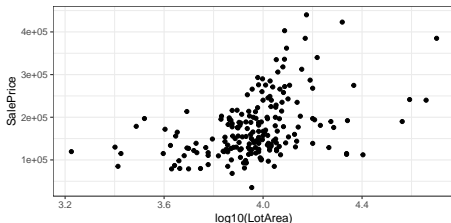


Add steps to recipes

- Consider the relationship between of sale price and lot area:



- Accuracy of a linear model may improve by performing log transformation on LotArea:



Adding steps to recipes

- Let's update our recipe:

```
house_rec <- house_rec %>%  
  step_log(LotArea, base = 10)
```

```
house_rec
```

```
## Recipe
```

```
##
```

```
## Inputs:
```

```
##
```

```
##      role #variables
```

```
##      ID      1
```

```
## outcome      1
```

```
## predictor     29
```

```
##
```

```
## Operations:
```

```
##
```

```
## Log transformation on LotArea
```

Create New Variables from Old

- The original data set contains variables `FullBath` and `HalfBath`. But we want a measure of total number of baths:

$$\text{TotalBath} = \text{FullBath} + \frac{1}{2}\text{HalfBath}$$

Create New Variables from Old

- The original data set contains variables `FullBath` and `HalfBath`. But we want a measure of total number of baths:

$$\text{TotalBath} = \text{FullBath} + \frac{1}{2}\text{HalfBath}$$

- We can also add a mutate step in our recipe to do just this:

```
house_rec <- house_rec %>%  
  step_mutate(TotalBath = FullBath+0.5*HalfBath) %>%  
  step_rm(FullBath, HalfBath)
```

```
house_rec
```

```
## Recipe  
##  
## Inputs:  
##  
##      role #variables  
##      ID      1  
##      outcome  1  
##      predictor 29  
##  
## Operations:  
##  
## Log transformation on LotArea
```

Create Dummy Variables

- Recall that 7 of our variables are factors (Functional, BldgType, Foundation, LotShape, LandSlope, SaleCondition, RoofMatl). To create appropriate dummy variables:

```
house_rec <- house_rec %>% step_dummy(all_nominal(), -all_outcomes())
house_rec
```

```
## Recipe
##
## Inputs:
##
##   role #variables
##   ID      1
##   outcome  1
##   predictor 29
##
## Operations:
##
## Log transformation on LotArea
## Variable mutation
## Delete terms FullBath, HalfBath
## Dummy variables from all_nominal(), -all_outcomes()
```

Create Dummy Variables

- Recall that 7 of our variables are factors (Functional, BldgType, Foundation, LotShape, LandSlope, SaleCondition, RoofMatl). To create appropriate dummy variables:

```
house_rec <- house_rec %>% step_dummy(all_nominal(), -all_outcomes())
house_rec
```

```
## Recipe
##
## Inputs:
##
##   role #variables
##   ID      1
##   outcome  1
##   predictor 29
##
## Operations:
##
## Log transformation on LotArea
## Variable mutation
## Delete terms FullBath, HalfBath
## Dummy variables from all_nominal(), -all_outcomes()
```

- The first argument `all_nominal` selects all variables that are either factors or characters
- The second argument `-all_outcomes` removes any response variables from this step

Create Dummy Variables

- Recall that 7 of our variables are factors (Functional, BldgType, Foundation, LotShape, LandSlope, SaleCondition, RoofMatl). To create appropriate dummy variables:

```
house_rec <- house_rec %>% step_dummy(all_nominal(), -all_outcomes())
house_rec
```

```
## Recipe
##
## Inputs:
##
##   role #variables
##   ID      1
##   outcome  1
##   predictor 29
##
## Operations:
##
## Log transformation on LotArea
## Variable mutation
## Delete terms FullBath, HalfBath
## Dummy variables from all_nominal(), -all_outcomes()
```

- The first argument `all_nominal` selects all variables that are either factors or characters
- The second argument `-all_outcomes` removes any response variables from this step

Remove Problematic Predictors

- Finally, to avoid the situation where an infrequently occurring level doesn't exist in the training or test sets:

```
house_rec <- house_rec %>% step_zv(all_predictors())
house_rec
```

```
## Recipe
##
## Inputs:
##
##       role #variables
##       ID           1
##  outcome           1
## predictor          29
##
## Operations:
##
## Log transformation on LotArea
## Variable mutation
## Delete terms FullBath, HalfBath
## Dummy variables from all_nominal(), -all_outcomes()
## Zero variance filter on all_predictors()
```

Remove Problematic Predictors

- Finally, to avoid the situation where an infrequently occurring level doesn't exist in the training or test sets:

```
house_rec <- house_rec %>% step_zv(all_predictors())  
house_rec
```

```
## Recipe  
##  
## Inputs:  
##  
##      role #variables  
##      ID           1  
##  outcome           1  
## predictor          29  
##  
## Operations:  
##  
## Log transformation on LotArea  
## Variable mutation  
## Delete terms FullBath, HalfBath  
## Dummy variables from all_nominal(), -all_outcomes()  
## Zero variance filter on all_predictors()
```

- The `step_zv` verb removes columns from the training data which have a single value

Workflows

- Why create a recipe when we could just as easily perform the pre-processing steps using `dplyr`?

Workflows

- Why create a recipe when we could just as easily perform the pre-processing steps using `dplyr`?
- ① The recipe allows us to apply the same procedures to both test and training data.
- ② The recipe gives instructions for processing the data **without actually performing that action**

Workflows

- Why create a recipe when we could just as easily perform the pre-processing steps using `dplyr`?
- ① The recipe allows us to apply the same procedures to both test and training data.
- ② The recipe gives instructions for processing the data **without actually performing that action**

To use our recipe across several steps, we will use a *workflow*, which will

- ① Process the recipe using the training set
- ② Apply the recipe to the training set
- ③ Apply the recipe to the test set

Create the workflow

```
house_mod <- linear_reg() %>% set_engine("lm")
```

```
house_wflow <- workflow() %>%  
  add_model(house_mod) %>%  
  add_recipe(house_rec)
```

```
house_wflow
```

```
## == Workflow =====  
## Preprocessor: Recipe  
## Model: linear_reg()  
##  
## -- Preprocessor -----  
## 5 Recipe Steps  
##  
## * step_log()  
## * step_mutate()  
## * step_rm()  
## * step_dummy()  
## * step_zv()  
##  
## -- Model -----  
## Linear Regression Model Specification (regression)  
##  
## Computational engine: lm
```

Fitting Models with Workflows

```
house_fit <- house_wflow %>% fit(data = train_data)
```

```
house_fit %>% pull_workflow_fit() %>% tidy()
```

```
## # A tibble: 47 x 5
```

```
##   term          estimate std.error statistic p.value
##   <chr>          <dbl>     <dbl>     <dbl>  <dbl>
## 1 (Intercept) 1920452. 3087160.    0.622 0.535
## 2 ScreenPorch    59.7      69.1      0.863 0.390
## 3 MSSubClass   -323.      129.     -2.50 0.0138
## 4 GarageCars   3602.     7179.    0.502 0.617
## 5 BedroomAbvGr 1509.     3933.    0.384 0.702
## 6 TotalBsmtSF   12.4      8.82     1.41 0.162
## 7 LotArea     20971.    20280.    1.03 0.304
## 8 OpenPorchSF  -17.0     37.0     -0.461 0.646
## 9 BsmtFullBath 16888.    5087.    3.32 0.00124
## 10 WoodDeckSF   18.5     17.9     1.03 0.305
## # ... with 37 more rows
```

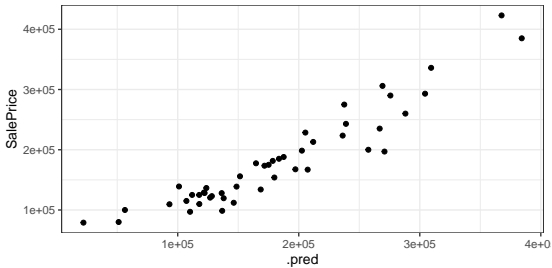
Making predictions with workflow

```
house_preds<- predict(house_fit, test_data)
house_preds
```

```
## # A tibble: 50 x 1
##   .pred
##   <dbl>
## 1 178531.
## 2 236275.
## 3 257502.
## 4 269208.
## 5  51212.
## 6 123668.
## 7 136742.
## 8 136343.
## 9 238991.
## 10      NA
## # ... with 40 more rows
```


Evaluate performance

```
house_results <- house_preds %>% cbind(test_data)
```



```
rbind(  
  rmse(house_results, truth = SalePrice, estimate = .pred),  
  rsq(house_results, truth = SalePrice, estimate = .pred)  
)
```

```
## # A tibble: 2 x 3  
##   .metric .estimator .estimate  
##   <chr>   <chr>         <dbl>  
## 1 rmse    standard      27049.  
## 2 rsq     standard        0.885
```

Section 4

Resampling

Resampling with `rsample`

- We previously built a linear model for `SalePrice` as a function of predictors in the house data and found the following accuracy measures on **test** data:

```
## # A tibble: 2 x 3
##   .metric .estimator .estimate
##   <chr>   <chr>         <dbl>
## 1 rmse    standard    27049.
## 2 rsq     standard     0.885
```

Resampling with `rsample`

- We previously built a linear model for `SalePrice` as a function of predictors in the house data and found the following accuracy measures on **test** data:

```
## # A tibble: 2 x 3
##   .metric .estimator .estimate
##   <chr>   <chr>         <dbl>
## 1 rmse    standard    27049.
## 2 rsq     standard     0.885
```

- But how typical are these estimates? Let's perform cross-validation.

```
set.seed(271)
library(rsample)
folds <- vfold_cv(train_data, v = 10, stattra = RoofMatl)
```

Delving Deeper

- Which observations are in each fold?

```
folds$splits[[1]]
```

```
## <Analysis/Assess/Total>
```

```
## <135/15/150>
```

```
folds$splits[[1]] %>% analysis() %>% head() %>% select(1:5)
```

```
##      SalePrice   Id Functional BldgType Foundation
## 58      81000    387          Typ      1Fam      PConc
## 37     113000    240          Typ      1Fam      CBlock
## 85     124000    631          Typ      1Fam      BrkTil
## 108    183000   821          Typ      1Fam      PConc
## 192    340000  1418          Typ      1Fam      PConc
## 165     93000   1180          Min2     1Fam      Slab
```

```
folds$splits[[1]] %>% assessment() %>% head() %>% select(1:5)
```

```
##      SalePrice   Id Functional BldgType Foundation
## 193    271000  1427          Typ      1Fam      PConc
## 74     130000   499          Typ      1Fam      PConc
## 131    136500   997          Typ      1Fam      CBlock
## 123    169990   923          Typ      1Fam      PConc
## 82     240000   622          Typ      1Fam      CBlock
## 194    119000  1429          Typ      1Fam      CBlock
```

Adding resampling to workflow

```
house_fit_resamples <- house_wflow %>% fit_resamples(folds)
house_fit_resamples
```

```
## # Resampling results
## # 10-fold cross-validation
## # A tibble: 10 x 4
##   splits          id      .metrics      .notes
##   <list>         <chr> <list>      <list>
## 1 <split [135/15]> Fold01 <tibble [2 x 4]> <tibble [1 x 1]>
## 2 <split [135/15]> Fold02 <tibble [2 x 4]> <tibble [1 x 1]>
## 3 <split [135/15]> Fold03 <tibble [2 x 4]> <tibble [1 x 1]>
## 4 <split [135/15]> Fold04 <tibble [2 x 4]> <tibble [1 x 1]>
## 5 <split [135/15]> Fold05 <tibble [2 x 4]> <tibble [1 x 1]>
## 6 <split [135/15]> Fold06 <tibble [2 x 4]> <tibble [1 x 1]>
## 7 <split [135/15]> Fold07 <tibble [2 x 4]> <tibble [1 x 1]>
## 8 <split [135/15]> Fold08 <tibble [2 x 4]> <tibble [1 x 1]>
## 9 <split [135/15]> Fold09 <tibble [2 x 4]> <tibble [1 x 1]>
## 10 <split [135/15]> Fold10 <tibble [2 x 4]> <tibble [1 x 1]>
```

Metrics

- Let's look at the results:

```
house_fit_resamples$.metrics[[1]]
```

```
## # A tibble: 2 x 4
##   .metric .estimator .estimate .config
##   <chr>   <chr>         <dbl> <chr>
## 1 rmse    standard    28540. Preprocessor1_Model1
## 2 rsq     standard     0.874 Preprocessor1_Model1
```

```
house_fit_resamples$.metrics[[2]]
```

```
## # A tibble: 2 x 4
##   .metric .estimator .estimate .config
##   <chr>   <chr>         <dbl> <chr>
## 1 rmse    standard    22605. Preprocessor1_Model1
## 2 rsq     standard     0.884 Preprocessor1_Model1
```

```
house_fit_resamples$.metrics[[3]]
```

```
## # A tibble: 2 x 4
##   .metric .estimator .estimate .config
##   <chr>   <chr>         <dbl> <chr>
## 1 rmse    standard    22019. Preprocessor1_Model1
## 2 rsq     standard     0.876 Preprocessor1_Model1
```

CV Performance

- How do the models do overall?

```
#Baseline
```

```
rbind(  
  rmse(house_results, truth = SalePrice, estimate = .pred),  
  rsq(house_results, truth = SalePrice, estimate = .pred)  
)
```

```
## # A tibble: 2 x 3  
##   .metric .estimator .estimate  
##   <chr>   <chr>         <dbl>  
## 1 rmse    standard    27049.  
## 2 rsq     standard     0.885
```


CV Performance

- How do the models do overall?

#Baseline

```
rbind(  
  rmse(house_results, truth = SalePrice, estimate = .pred),  
  rsq(house_results, truth = SalePrice, estimate = .pred)  
)
```

```
## # A tibble: 2 x 3  
##   .metric .estimator .estimate  
##   <chr>   <chr>         <dbl>  
## 1 rmse    standard    27049.  
## 2 rsq     standard     0.885
```

- Cross-validation:

```
collect_metrics(house_fit_resamples)
```

```
## # A tibble: 2 x 6  
##   .metric .estimator   mean     n  std_err .config  
##   <chr>   <chr>         <dbl> <int>   <dbl> <chr>  
## 1 rmse    standard    24994.    10  1939.   Preprocessor1_Model1  
## 2 rsq     standard     0.862    10   0.0237 Preprocessor1_Model1
```

Section 5

Tuning Hyperparameters

Building a LASSO model

- The linear model did fine. But can we improve our results using penalized regression?

Building a LASSO model

- The linear model did fine. But can we improve our results using penalized regression?
 - Note that our data pre-processing recipe `house_rec` is still valid (although we could change it)

Building a LASSO model

- The linear model did fine. But can we improve our results using penalized regression?
 - Note that our data pre-processing recipe `house_rec` is still valid (although we could change it)
- If we wanted a LASSO model with particular penalty (say $\lambda = 4$) we could use

```
house_lasso_mod <- linear_reg(penalty =4 ) %>% set_engine("glmnet")
```

Building a LASSO model

- The linear model did fine. But can we improve our results using penalized regression?
 - Note that our data pre-processing recipe `house_rec` is still valid (although we could change it)
- If we wanted a LASSO model with particular penalty (say $\lambda = 4$) we could use

```
house_lasso_mod <- linear_reg(penalty = 4 ) %>% set_engine("glmnet")
```

- But we are really interested in finding the **BEST** value of λ . So instead

```
house_lasso_mod <- linear_reg(penalty = tune() ) %>% set_engine("glmnet")
```

Building a LASSO model

- The linear model did fine. But can we improve our results using penalized regression?
 - Note that our data pre-processing recipe `house_rec` is still valid (although we could change it)
- If we wanted a LASSO model with particular penalty (say $\lambda = 4$) we could use

```
house_lasso_mod <- linear_reg(penalty = 4 ) %>% set_engine("glmnet")
```

- But we are really interested in finding the **BEST** value of λ . So instead

```
house_lasso_mod <- linear_reg(penalty = tune() ) %>% set_engine("glmnet")
```

- Let's fit the model and tune

```
lasso_grid <- grid_regular(penalty() %>% range_set(c(-5,5)), levels = 10)  
lasso_wf <- workflow() %>% add_model(house_lasso_mod) %>% add_recipe(house_rec)  
lasso_res <- lasso_wf %>% tune_grid(grid = lasso_grid, resamples = folds)
```

Results

```
collect_metrics(lasso_res)
```

```
## # A tibble: 20 x 7
##       penalty .metric .estimator      mean      n  std_err .config
##       <dbl> <chr>   <chr>         <dbl> <int>    <dbl> <chr>
## 1  0.00001 rmse    standard  24577.     10  1966.   Preprocessor1_Mod~
## 2  0.00001 rsq     standard    0.867     10   0.0209 Preprocessor1_Mod~
## 3  0.000129 rmse    standard  24577.     10  1966.   Preprocessor1_Mod~
## 4  0.000129 rsq     standard    0.867     10   0.0209 Preprocessor1_Mod~
## 5  0.00167 rmse    standard  24577.     10  1966.   Preprocessor1_Mod~
## 6  0.00167 rsq     standard    0.867     10   0.0209 Preprocessor1_Mod~
## 7  0.0215  rmse    standard  24577.     10  1966.   Preprocessor1_Mod~
## 8  0.0215  rsq     standard    0.867     10   0.0209 Preprocessor1_Mod~
## 9  0.278   rmse    standard  24577.     10  1966.   Preprocessor1_Mod~
## 10 0.278   rsq     standard    0.867     10   0.0209 Preprocessor1_Mod~
## 11 3.59    rmse    standard  24577.     10  1966.   Preprocessor1_Mod~
## 12 3.59    rsq     standard    0.867     10   0.0209 Preprocessor1_Mod~
## 13 46.4   rmse    standard  24499.     10  1981.   Preprocessor1_Mod~
## 14 46.4   rsq     standard    0.868     10   0.0209 Preprocessor1_Mod~
## 15 599.   rmse    standard  23612.     10  2259.   Preprocessor1_Mod~
## 16 599.   rsq     standard    0.878     10   0.0111 Preprocessor1_Mod~
## 17 7743.  rmse    standard  28677.     10  3090.   Preprocessor1_Mod~
## 18 7743.  rsq     standard    0.844     10   0.0281 Preprocessor1_Mod~
## 19 100000 rmse    standard  67654.     10  5580.   Preprocessor1_Mod~
## 20 100000 rsq     standard    NaN        0    NA     Preprocessor1_Mod~
```


Results

```
collect_metrics(lasso_res)
```

```
## # A tibble: 20 x 7
##       penalty .metric .estimator      mean     n  std_err .config
##       <dbl> <chr> <chr>      <dbl> <int>   <dbl> <chr>
## 1  0.00001 rmse  standard  24577.    10  1966.  Preprocessor1_Mod~
## 2  0.00001 rsq   standard    0.867    10   0.0209 Preprocessor1_Mod~
## 3  0.000129 rmse  standard  24577.    10  1966.  Preprocessor1_Mod~
## 4  0.000129 rsq   standard    0.867    10   0.0209 Preprocessor1_Mod~
## 5  0.00167 rmse  standard  24577.    10  1966.  Preprocessor1_Mod~
## 6  0.00167 rsq   standard    0.867    10   0.0209 Preprocessor1_Mod~
## 7  0.0215  rmse  standard  24577.    10  1966.  Preprocessor1_Mod~
## 8  0.0215  rsq   standard    0.867    10   0.0209 Preprocessor1_Mod~
## 9  0.278   rmse  standard  24577.    10  1966.  Preprocessor1_Mod~
## 10 0.278   rsq   standard    0.867    10   0.0209 Preprocessor1_Mod~
## 11 3.59    rmse  standard  24577.    10  1966.  Preprocessor1_Mod~
## 12 3.59    rsq   standard    0.867    10   0.0209 Preprocessor1_Mod~
## 13 46.4   rmse  standard  24499.    10  1981.  Preprocessor1_Mod~
## 14 46.4   rsq   standard    0.868    10   0.0209 Preprocessor1_Mod~
## 15 599.   rmse  standard  23612.    10  2259.  Preprocessor1_Mod~
## 16 599.   rsq   standard    0.878    10   0.0111 Preprocessor1_Mod~
## 17 7743.  rmse  standard  28677.    10  3090.  Preprocessor1_Mod~
## 18 7743.  rsq   standard    0.844    10   0.0281 Preprocessor1_Mod~
## 19 100000 rmse  standard  67654.    10  5580.  Preprocessor1_Mod~
## 20 100000 rsq   standard    NaN      0    NA     Preprocessor1_Mod~
```

Which penalties?

- Focus just on optimal penalties for rmse:

```
lasso_res %>%  
  show_best("rmse")
```

```
## # A tibble: 5 x 7  
##   penalty .metric .estimator   mean     n std_err .config  
##   <dbl> <chr>   <chr>     <dbl> <int>  <dbl> <chr>  
## 1 599.     rmse     standard 23612.   10  2259. Preprocessor1_Model108  
## 2 46.4     rmse     standard 24499.   10  1981. Preprocessor1_Model107  
## 3 0.00001  rmse     standard 24577.   10  1966. Preprocessor1_Model101  
## 4 0.000129 rmse     standard 24577.   10  1966. Preprocessor1_Model102  
## 5 0.00167  rmse     standard 24577.   10  1966. Preprocessor1_Model103
```

Which penalties?

- Focus just on optimal penalties for rmse:

```
lasso_res %>%  
  show_best("rmse")
```

```
## # A tibble: 5 x 7  
##   penalty .metric .estimator   mean     n std_err .config  
##   <dbl> <chr>   <chr>     <dbl> <int>  <dbl> <chr>  
## 1 599.     rmse     standard 23612.   10   2259. Preprocessor1_Model108  
## 2 46.4     rmse     standard 24499.   10   1981. Preprocessor1_Model107  
## 3 0.00001  rmse     standard 24577.   10   1966. Preprocessor1_Model101  
## 4 0.000129 rmse     standard 24577.   10   1966. Preprocessor1_Model102  
## 5 0.00167  rmse     standard 24577.   10   1966. Preprocessor1_Model103
```

- Let's collect the best model:

```
best_lasso <- lasso_res %>% select_best(metric = "rmse")  
best_lasso
```

```
## # A tibble: 1 x 2  
##   penalty .config  
##   <dbl> <chr>  
## 1 599. Preprocessor1_Model108
```

Finalize the model

- We update or finalize our workflow with the values from `select_best`:

```
final_lasso_wf <- lasso_wf %>% finalize_workflow(best_lasso)
final_lasso_wf
```

```
## == Workflow =====
## Preprocessor: Recipe
## Model: linear_reg()
##
## -- Preprocessor -----
## 5 Recipe Steps
##
## * step_log()
## * step_mutate()
## * step_rm()
## * step_dummy()
## * step_zv()
##
## -- Model -----
## Linear Regression Model Specification (regression)
##
## Main Arguments:
##   penalty = 599.484250318942
##
## Computational engine: glmnet
```

Fit the Best Model

- Thus far, we've just focused on finding the best parameter. But we haven't actually built a LASSO model on training data. Let's do that:

Fit the Best Model

- Thus far, we've just focused on finding the best parameter. But we haven't actually built a LASSO model on training data. Let's do that:

```
final_lasso_fit<-final_lasso_wf %>% last_fit(data_split )
```

```
final_lasso_fit$.metrics
```

```
## [[1]]
## # A tibble: 2 x 4
##   .metric .estimator .estimate .config
##   <chr>   <chr>         <dbl> <chr>
## 1 rmse    standard    26889. Preprocessor1_Model1
## 2 rsq     standard     0.883 Preprocessor1_Model1
```

```
final_lasso_fit$.predictions
```

```
## [[1]]
## # A tibble: 50 x 4
##   .pred .row SalePrice .config
##   <dbl> <int>   <int> <chr>
## 1 178628.     1   181500 Preprocessor1_Model1
## 2 231755.     2   223500 Preprocessor1_Model1
## 3 247905.     3   200000 Preprocessor1_Model1
## 4 265349.     7   306000 Preprocessor1_Model1
## 5  58165.    12    80000 Preprocessor1_Model1
## 6 129233.    14   136500 Preprocessor1_Model1
## 7 135097.    16    98600 Preprocessor1_Model1
```