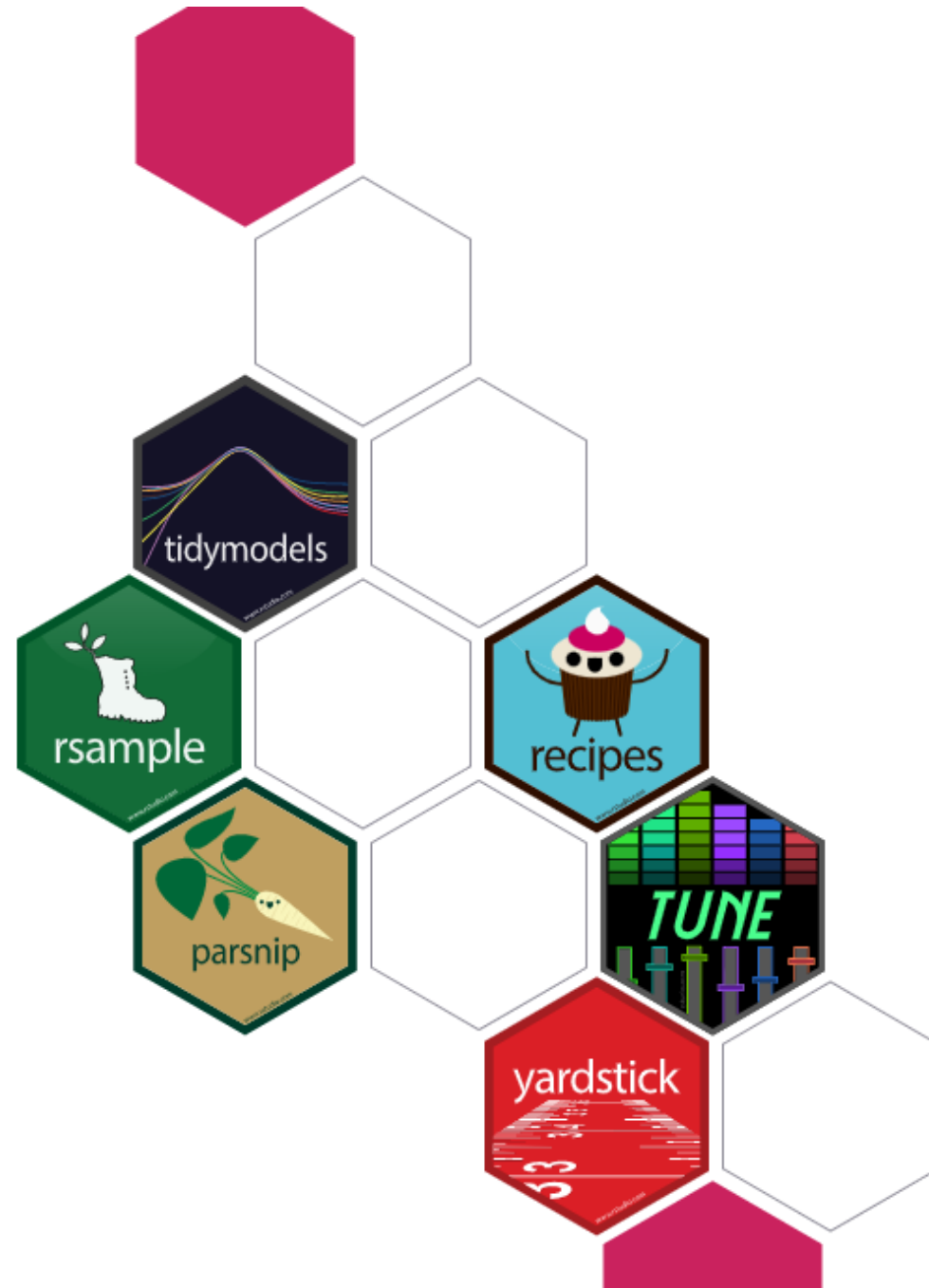


More Tidymodels

Lecture 23

Dr. Colin Rundel



Hotels Data

Original data from Antonio, Almeida, and Nunes (2019), Data dictionary

```
1 hotels = read_csv(  
2   'https://tidymodels.org/start/case-study/hotels.csv'  
3 ) %>%  
4   mutate(  
5     across(where(is.character), as.factor)  
6   )
```

The data

```
1 glimpse(hotels)
```

Rows: 50,000

Columns: 23

```
$ hotel          <fct> City_Hotel, City_Ho...
$ lead_time      <dbl> 217, 2, 95, 143, 13...
$ stays_in_weekend_nights <dbl> 1, 0, 2, 2, 1, 2, 0...
$ stays_in_week_nights  <dbl> 3, 1, 5, 6, 4, 2, 2...
$ adults         <dbl> 2, 2, 2, 2, 2, 2, 2...
$ children       <fct> none, none, none, n...
$ meal           <fct> BB, BB, BB, HB, HB,...
$ country        <fct> DEU, PRT, GBR, ROU,...
$ market_segment <fct> Offline_TA/TO, Dire...
$ distribution_channel <fct> TA/TO, Direct, TA/T...
$ is_repeated_guest <dbl> 0, 0, 0, 0, 0, 0, 0...
$ previous_cancellations <dbl> 0, 0, 0, 0, 0, 0, 0...
$ previous_bookings_not_canceled <dbl> 0, 0, 0, 0, 0, 0, 0...
$ reserved_room_type <fct> A, D, A, A, F, A, C...
$ assigned_room_type <fct> A. K. A. A. F. A. C...
```

The model

Our goal is to develop a predictive model that is able to predict whether a booking will include children or not based on the other characteristics of the booking.

```
1 hotels %>%  
2   count(children) %>%  
3   mutate(prop = n/sum(n))
```

```
# A tibble: 2 × 3
```

	children	n	prop
	<fct>	<int>	<dbl>
1	children	4038	0.0808
2	none	45962	0.919

Clustering the test/train split

```
1 set.seed(123)
2
3 splits = initial_split(
4   hotels, strata = children
5 )
6
7 hotel_train = training(splits)
8 hotel_test = testing(splits)
```

```
1 dim(hotel_train)
```

```
[1] 37500    23
```

```
1 dim(hotel_test)
```

```
[1] 12500    23
```

```
1 hotel_train %>%
2   count(children) %>%
3   mutate(prop = n/sum(n))
```

```
# A tibble: 2 × 3
  children      n  prop
  <fct>      <int> <dbl>
1 children  3027 0.0807
2 none    34473 0.919
```

```
1 hotel_test %>%
2   count(children) %>%
3   mutate(prop = n/sum(n))
```

```
# A tibble: 2 × 3
  children      n  prop
  <fct>      <int> <dbl>
1 children  1011 0.0809
2 none    11489 0.919
```

Logistic Regression model

```
1 show_engines("logistic_reg")
```

```
# A tibble: 7 × 2
  engine    mode
  <chr>    <chr>
1 glm      classification
2 glmnet   classification
3 Liblinear classification
4 spark    classification
5 keras    classification
6 stan     classification
7 brulee   classification
```

```
1 lr_model = logistic_reg() %>%
2   set_engine("glm")
3 translate(lr_model)
```

Logistic Regression Model Specification (classification)

Computational engine: glm

Model fit template:

```
stats::glm(formula = missing_arg(), data = missing_arg(), weights = missing_arg(),
  family = stats::binomial)
```

Recipe

```
1 holidays = c("AllSouls", "AshWednesday", "ChristmasEve", "Easter",  
2             "ChristmasDay", "GoodFriday", "NewYearsDay", "PalmSunday")  
3  
4 lr_recipe = recipe(children ~ ., data = hotel_train) %>%  
5   step_date(arrival_date) %>%  
6   step_holiday(arrival_date, holidays = holidays) %>%  
7   step_rm(arrival_date) %>%  
8   step_rm(country) %>%  
9   step_dummy(all_nominal_predictors()) %>%  
10  step_zv(all_predictors())  
11  
12 lr_recipe
```

Recipe

Inputs:

	role	#variables
outcome		1
predictor		22

Operations:

Date features from arrival_date

Holiday features from arrival_date

Variables removed arrival_date

Variables removed country

Dummy variables from `all_nominal_predictors()`
Zero variance filter on `all_predictors()`

```

1 lr_recipe %>%
2   prep() %>%
3   bake(new_data = hotel_train)

```

A tibble: 37,500 × 76

	lead_time	stays_... ¹	stays_... ²	adults	is_re... ³	previ... ⁴	previ... ⁵	booki... ⁶	days_... ⁷	avera... ⁸	total... ⁹	child... ^x
	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<fct>
1	2	0	1	2	0	0	0	0	0	170	3	none
2	95	2	5	2	0	0	0	2	0	8	2	none
3	67	2	2	2	0	0	0	0	0	49.1	1	none
4	47	0	2	2	0	0	0	0	0	289	1	childr...
5	56	0	3	0	0	0	0	0	0	82.4	1	childr...
6	6	2	2	2	0	0	0	0	0	180	1	childr...
7	130	1	2	2	0	0	0	0	0	71	0	none
8	27	0	1	1	0	0	0	0	0	120.	1	none
9	46	0	2	2	0	0	0	0	0	162	0	none
10	423	1	1	2	0	0	0	0	0	122.	1	none

... with 37,490 more rows, 64 more variables: arrival_date_year <int>, arrival_date_AllSouls <int>,
arrival_date_AshWednesday <int>, arrival_date_ChristmasEve <int>, arrival_date_Easter <int>,
arrival_date_ChristmasDay <int>, arrival_date_GoodFriday <int>, arrival_date_NewYearsDay <int>,
arrival date PalmSundav <int>. hotel Resort Hotel <dbl>. meal FB <dbl>. meal HB <dbl>.

Workflow

```
1 ( lr_work = workflow() %>%  
2   add_model(lr_model) %>%  
3   add_recipe(lr_recipe)  
4 )
```

== Workflow ==

Preprocessor: Recipe

Model: logistic_reg()

— Preprocessor —

6 Recipe Steps

- step_date()
- step_holiday()
- step_rm()
- step_rm()
- step_dummy()
- step_zv()

— Model —

Logistic Regression Model Specification (classification)

Fit

```
1 ( lr_fit = lr_work %>%  
2   fit(data = hotel_train) )
```

== Workflow [trained] ==

Preprocessor: Recipe

Model: logistic_reg()

— Preprocessor —

6 Recipe Steps

- step_date()
- step_holiday()
- step_rm()
- step_rm()
- step_dummy()
- step_zv()

— Model —

Call: stats::glm(formula = ..v ~ .., family = stats::binomial, data = data)

Logistic regression predictions

```
1 ( lr_perf = lr_fit %>%  
2   augment(new_data = hotel_train) %>%  
3   select(children, starts_with(".pred")) )
```

A tibble: 37,500 × 4

	children	.pred_class	.pred_children	.pred_none
	<fct>	<fct>	<dbl>	<dbl>
1	none	none	0.0861	0.914
2	none	none	0.0178	0.982
3	none	none	0.0101	0.990
4	children	children	0.931	0.0693
5	children	none	0.473	0.527
6	children	none	0.144	0.856
7	none	none	0.0710	0.929
8	none	none	0.0596	0.940
9	none	none	0.0252	0.975
10	none	none	0.0735	0.926

... with 37,490 more rows

Performance metrics (within-sample)

```
1 lr_perf %>%
2   conf_mat(children, .pred_class)
```

	Truth	
Prediction	children	none
children	1075	420
none	1952	34053

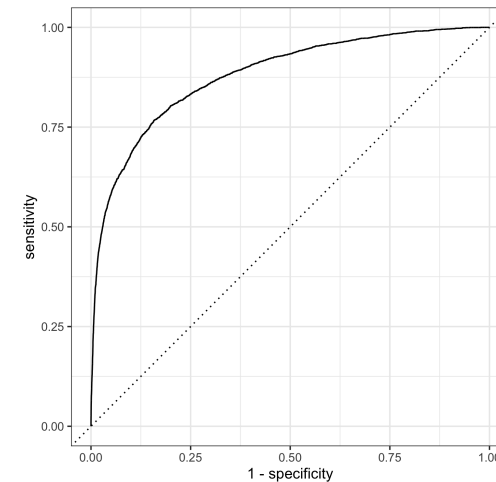
```
1 lr_perf %>%
2   precision(children, .pred_class)
```

```
# A tibble: 1 × 3
  .metric .estimator .estimate
  <chr>    <chr>         <dbl>
1 precision binary      0.719
```

```
1 lr_perf %>%
2   roc_auc(children, .pred_children)
```

```
# A tibble: 1 × 3
  .metric .estimator .estimate
  <chr>    <chr>         <dbl>
1 roc_auc binary      0.881
```

```
1 lr_perf %>%
2   yardstick::roc_curve(
3     children,
4     .pred_children
5   ) %>%
6   autoplot()
```



Performance metrics (out-of-sample)

```
1 lr_test_perf = lr_fit %>%
2   augment(new_data = hotel_test) %>%
3   select(children, starts_with(".pred"))
4
5 lr_test_perf %>%
6   conf_mat(children, .pred_class)
```

	Truth	
Prediction	children	none
children	359	137
none	652	11352

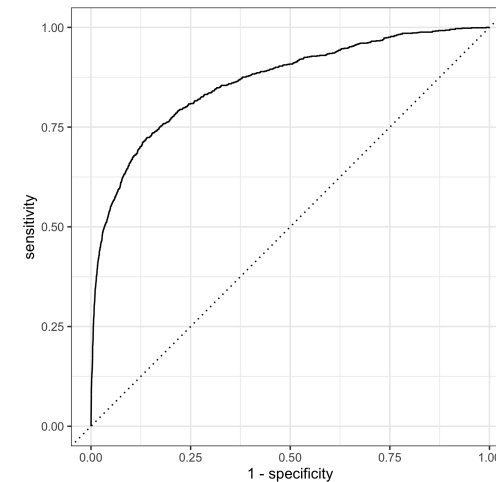
```
1 lr_test_perf %>%
2   precision(children, .pred_class)
```

```
# A tibble: 1 × 3
  .metric .estimator .estimate
  <chr>   <chr>       <dbl>
1 precision binary      0.724
```

```
1 lr_test_perf %>%
2   roc_auc(children, .pred_children)
```

```
# A tibble: 1 × 3
  .metric .estimator .estimate
  <chr>   <chr>       <dbl>
1 roc_auc binary      0.864
```

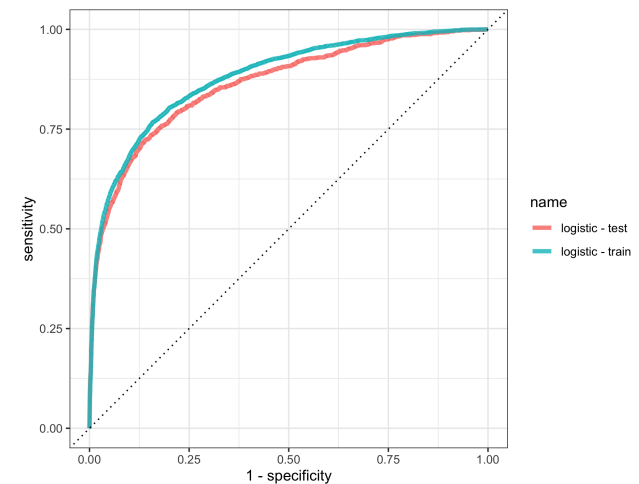
```
1 lr_test_perf %>%
2   yardstick::roc_curve(
3     children,
4     .pred_children
5   ) %>%
6   autoplot()
```



Combining ROC curves

```
1 lr_roc_train = lr_perf %>%  
2   yardstick::roc_curve(children, .pred_children)  
3   mutate(name="logistic - train")  
4  
5 lr_roc_test = lr_test_perf %>%  
6   yardstick::roc_curve(children, .pred_children)  
7   mutate(name="logistic - test")
```

```
1 bind_rows(  
2   lr_roc_train,  
3   lr_roc_test  
4 ) %>%  
5   ggplot(aes(x = 1 - specificity, y = sensitivity,  
6             geom_path(lwd = 1.5, alpha = 0.8) +  
7             geom_abline(lty = 3) +  
8             coord_equal())
```



Lasso

Lasso Model

For this we will be using the `glmnet` package which supports fitting lasso, ridge and elastic net models.

```
1 lasso_model = logistic_reg(penalty = tune(), mixture = 1) %>%  
2   set_engine("glmnet")
```

- `mixture` determines the type of model fit with: `1` -> Lasso, `0` -> Ridge, other -> elastic net.
- `penalty` is λ in the lasso model, scales the penalty for coefficient size.

```
1 lasso_model %>%
2   parameters()
```

Collection of 1 parameters for tuning

identifier	type	object
penalty	penalty	nparam[+]

```
1 lasso_model %>%
2   translate()
```

Logistic Regression Model Specification (classification)

Main Arguments:

```
penalty = tune()
mixture = 1
```

Computational engine: glmnet

Model fit template:

```
glmnet::glmnet(x = missing_arg(), y = missing_arg(), weights = missing_arg(),
  alpha = 1, family = "binomial")
```

Lasso Recipe

Lasso (and Ridge) models are sensitive to the scale of the model features, and so a standard approach is to normalize all features before fitting the model.

```
1 lasso_recipe = lr_recipe %>%  
2   step_normalize(all_predictors())
```

```
1 lasso_recipe %>%  
2   prep() %>%  
3   bake(new_data = hotel_train)
```

```
# A tibble: 37,500 × 76
```

	lead_time	stays_... ¹	stays... ²	adults	is_re... ³	previ... ⁴	previ... ⁵
	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
1	-0.858	-0.938	-0.767	0.337	-0.213	-0.0597	-0.112
2	0.160	1.09	1.32	0.337	-0.213	-0.0597	-0.112
3	-0.146	1.09	-0.245	0.337	-0.213	-0.0597	-0.112
4	-0.365	-0.938	-0.245	0.337	-0.213	-0.0597	-0.112
5	-0.267	-0.938	0.278	-3.59	-0.213	-0.0597	-0.112
6	-0.814	1.09	-0.245	0.337	-0.213	-0.0597	-0.112
7	0.544	0.0735	-0.245	0.337	-0.213	-0.0597	-0.112
8	-0.584	-0.938	-0.767	-1.63	-0.213	-0.0597	-0.112
9	-0.376	-0.938	-0.245	0.337	-0.213	-0.0597	-0.112
10	3.75	0.0735	-0.767	0.337	-0.213	-0.0597	-0.112

```
# ... with 37,490 more rows, 69 more variables:
```

```
#   booking_changes <dbl>, days_in_waiting_list <dbl>,
```

```
#   average_daily_rate <dbl>,
```

Sta 523 - Fall 2022

Lasso workflow

```
1 ( lasso_work = workflow() %>%  
2   add_model(lasso_model) %>%  
3   add_recipe(lasso_recipe)  
4 )
```

== Workflow ==

Preprocessor: Recipe

Model: logistic_reg()

— Preprocessor —

7 Recipe Steps

- step_date()
- step_holiday()
- step_rm()
- step_rm()
- step_dummy()
- step_zv()
- step_normalize()

— Model —

Logistic Regression Model Specification (classification)

v-folds for hyperparameter tuning

```
1 ( hotel_vf = rsample::vfold_cv(hotel_train, v=5, strata = children) )
```

```
# 5-fold cross-validation using stratification
```

```
# A tibble: 5 × 2
```

	splits	id
	<list>	<chr>
1	<split [30000/7500]>	Fold1
2	<split [30000/7500]>	Fold2
3	<split [30000/7500]>	Fold3
4	<split [30000/7500]>	Fold4
5	<split [30000/7500]>	Fold5

grid search

```
1 ( lasso_grid = lasso_work %>%
2   tune_grid(
3     hotel_vf,
4     grid = tibble(
5       penalty = 10^seq(-4, -1, length.out = 10)
6     ),
7     control = control_grid(save_pred = TRUE),
8     metrics = metric_set(roc_auc)
9   )
10 )
```

Tuning results

5-fold cross-validation using stratification

A tibble: 5 × 5

	splits	id	.metrics	.notes	.predictions
	<list>	<chr>	<list>	<list>	<list>
1	<split [30000/7500]>	Fold1	<tibble>	<tibble>	<tibble>
2	<split [30000/7500]>	Fold2	<tibble>	<tibble>	<tibble>
3	<split [30000/7500]>	Fold3	<tibble>	<tibble>	<tibble>
4	<split [30000/7500]>	Fold4	<tibble>	<tibble>	<tibble>
5	<split [30000/7500]>	Fold5	<tibble>	<tibble>	<tibble>

Results

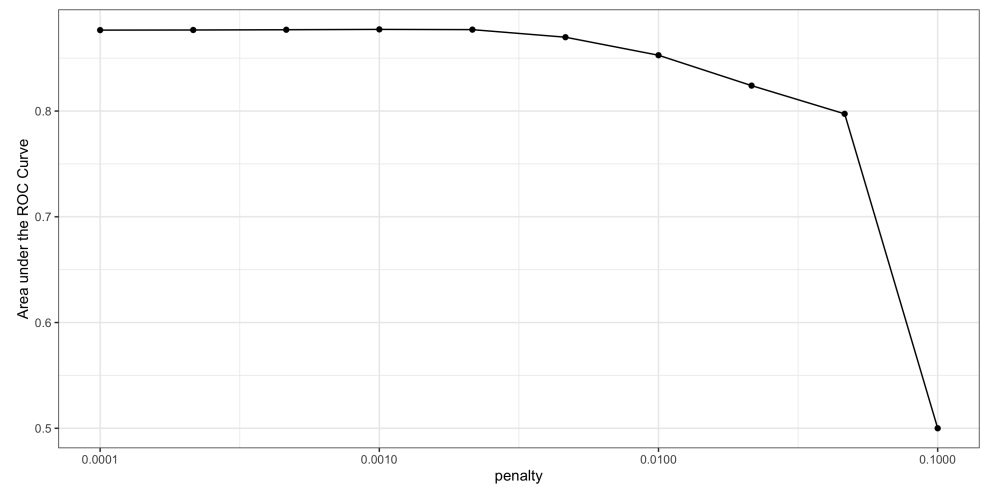
```
1 lasso_grid %>%
2   collect_metrics()
```

A tibble: 10 × 7

	penalty	.metric	.estimator	mean	n	std_err
	<dbl>	<chr>	<chr>	<dbl>	<int>	<dbl>
1	0.0001	roc_auc	binary	0.877	5	0.00318
2	0.000215	roc_auc	binary	0.877	5	0.00316
3	0.000464	roc_auc	binary	0.877	5	0.00314
4	0.001	roc_auc	binary	0.877	5	0.00304
5	0.00215	roc_auc	binary	0.877	5	0.00263
6	0.00464	roc_auc	binary	0.870	5	0.00253
7	0.01	roc_auc	binary	0.853	5	0.00249
8	0.0215	roc_auc	binary	0.824	5	0.00424
9	0.0464	roc_auc	binary	0.797	5	0.00400
10	0.1	roc_auc	binary	0.5	5	0

... with 1 more variable: .config <chr>

```
1 lasso_grid %>%
2   collect_metrics() %>%
3   ggplot(aes(x = penalty, y = mean)) +
4     geom_point() +
5     geom_line() +
6     ylab("Area under the ROC Curve") +
7     scale_x_log10(labels = scales::label_number(
```



“Best” models

```
1 lasso_grid %>%  
2   show_best("roc_auc", n=10)
```

A tibble: 10 × 7

	penalty	.metric	.estimator	mean	n	std_err	.config
	<dbl>	<chr>	<chr>	<dbl>	<int>	<dbl>	<chr>
1	0.001	roc_auc	binary	0.877	5	0.00304	Preproce...
2	0.00215	roc_auc	binary	0.877	5	0.00263	Preproce...
3	0.000464	roc_auc	binary	0.877	5	0.00314	Preproce...
4	0.000215	roc_auc	binary	0.877	5	0.00316	Preproce...
5	0.0001	roc_auc	binary	0.877	5	0.00318	Preproce...
6	0.00464	roc_auc	binary	0.870	5	0.00253	Preproce...
7	0.01	roc_auc	binary	0.853	5	0.00249	Preproce...
8	0.0215	roc_auc	binary	0.824	5	0.00424	Preproce...
9	0.0464	roc_auc	binary	0.797	5	0.00400	Preproce...
10	0.1	roc_auc	binary	0.5	5	0	Preproce...

“Best” model

```
1 ( lasso_best = lasso_grid %>%  
2   collect_metrics() %>%  
3   mutate(mean = round(mean, 2)) %>%  
4   arrange(desc(mean), desc(penalty)) %>%  
5   slice(1) )
```

A tibble: 1 × 7

	penalty	.metric	.estimator	mean	n	std_err	.config
	<dbl>	<chr>	<chr>	<dbl>	<int>	<dbl>	<chr>
1	0.00215	roc_auc	binary	0.88	5	0.00263	Preprocess...

Extracting predictions

Since we used `control_grid(save_pred = TRUE)` with `tune_grid()` we can recover the predictions for the out-of-sample values for each fold:

```
1 ( lasso_train_perf = lasso_grid %>%  
2   collect_predictions(parameters = lasso_best) )
```

```
# A tibble: 37,500 × 7  
  id    .pred_child...1 .pred...2  .row penalty child...3 .config  
  <chr>      <dbl>    <dbl> <int>    <dbl> <fct>    <chr>  
1 Fold1      0.366     0.634     5 0.00215 childr... Prepro...  
2 Fold1      0.144     0.856     6 0.00215 childr... Prepro...  
3 Fold1      0.0542    0.946    19 0.00215 none     Prepro...  
4 Fold1      0.0266    0.973    21 0.00215 none     Prepro...  
5 Fold1      0.106     0.894    22 0.00215 childr... Prepro...  
6 Fold1      0.0286    0.971    23 0.00215 none     Prepro...  
7 Fold1      0.0205    0.980    30 0.00215 none     Prepro...  
8 Fold1      0.0192    0.981    31 0.00215 none     Prepro...  
9 Fold1      0.0431    0.957    32 0.00215 none     Prepro...  
10 Fold1     0.0532    0.947    35 0.00215 none     Prepro...  
# ... with 37,490 more rows, and abbreviated variable names  
#   1.pred_children, 2.pred_none, 3children
```

```
1 lasso_train_perf %>%  
2   roc_auc(children, .pred_children)
```

```
# A tibble: 1 × 3  
  .metric .estimator .estimate  
  <chr>   <chr>       <dbl>  
1 roc_auc binary      0.877
```

Re-fitting

Typically with a tuned model we will refit using the complete test data and the “best” parameter value(s),

```
1 lasso_work_tuned = update_model(  
2   lasso_work,  
3   logistic_reg(  
4     mixture = 1,  
5     penalty = lasso_best$penalty  
6   ) %>%  
7     set_engine("glmnet")  
8 )  
9  
10 lasso_fit = lasso_work_tuned %>%  
11   fit(data=hotel_train)
```

Test Performance (out-of-sample)

```
1 lasso_test_perf = lasso_fit %>%
2   augment(new_data = hotel_test) %>%
3   select(children, starts_with(".pred"))
4
5 lasso_test_perf %>%
6   conf_mat(children, .pred_class)
```

	Truth	
Prediction	children	none
children	330	109
none	681	11380

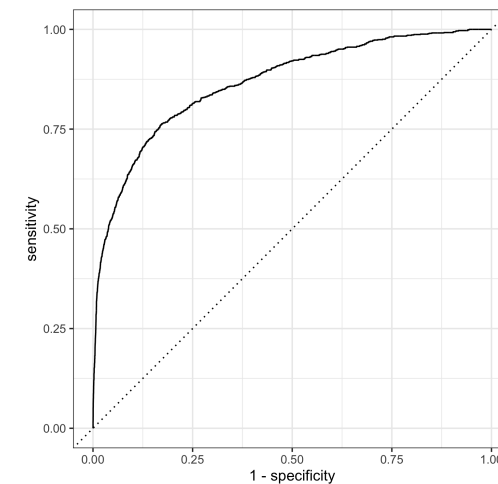
```
1 lasso_test_perf %>%
2   precision(children, .pred_class)
```

```
# A tibble: 1 × 3
  .metric .estimator .estimate
  <chr>    <chr>      <dbl>
1 precision binary      0.752
```

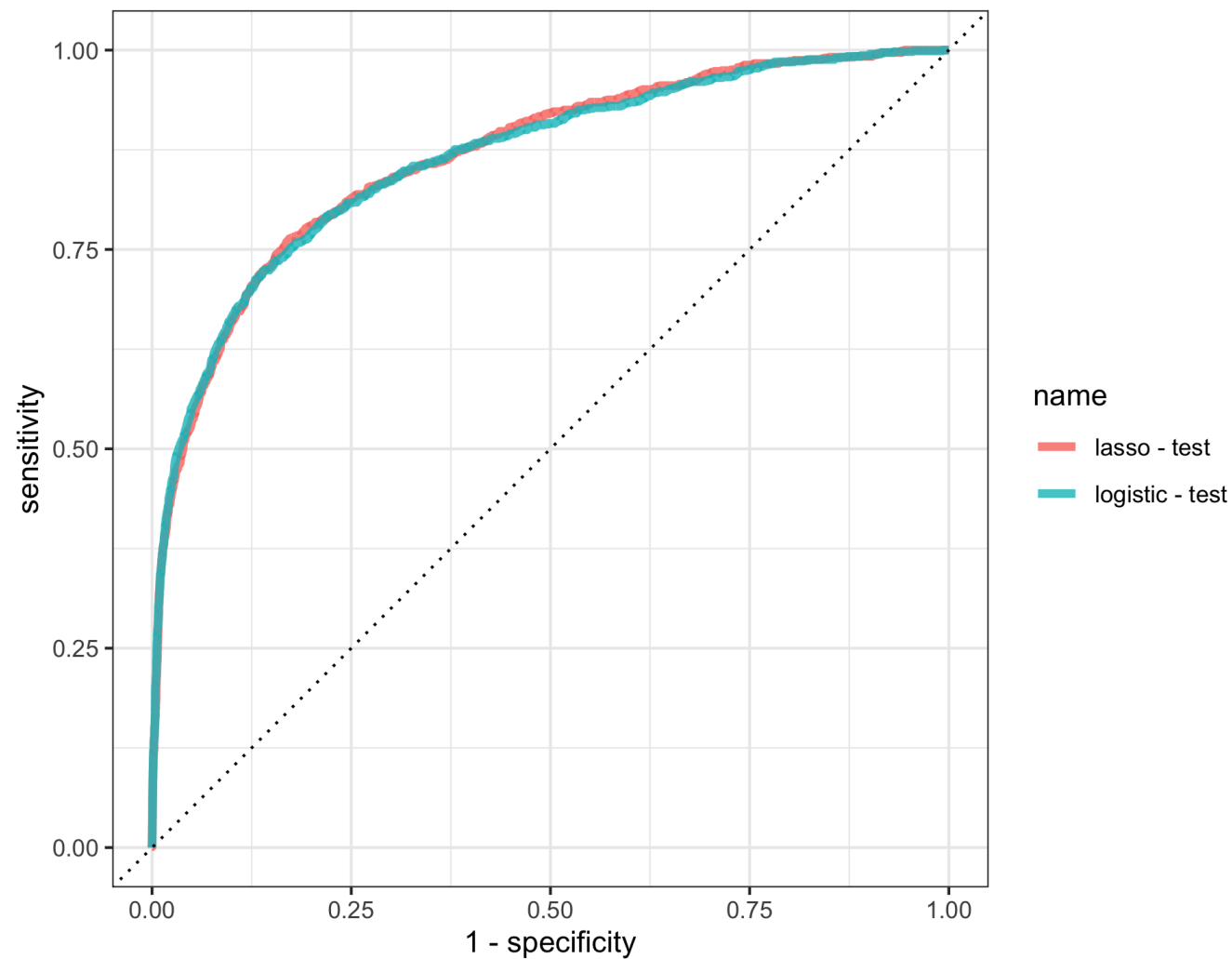
```
1 lasso_test_perf %>%
2   roc_auc(children, .pred_children)
```

```
# A tibble: 1 × 3
  .metric .estimator .estimate
  <chr>    <chr>      <dbl>
1 roc_auc binary      0.866
```

```
1 lasso_roc = lasso_test_perf %>%
2   yardstick::roc_curve(
3     children,
4     .pred_children
5   ) %>%
6   mutate(name = "lasso - test")
7 lasso_roc %>%
8   autoplot()
```



Comparing models



Random Forest

Random forest models

```
1 show_engines("rand_forest")
```

```
# A tibble: 6 × 2
```

	engine	mode
	<chr>	<chr>
1	ranger	classification
2	ranger	regression
3	randomForest	classification
4	randomForest	regression
5	spark	classification
6	spark	regression

```
1 rf_model = rand_forest(mtry = tune(), min_n = tune(), trees = 100) %>%  
2   set_engine("ranger", num.threads = 8) %>%  
3   set_mode("classification")
```

Recipe & workflow

We skip dummy coding in the recipe as it is not needed by ranger,

```
1 rf_recipe = recipe(children ~ ., data = hotel_train) %>%  
2   step_date(arrival_date) %>%  
3   step_holiday(arrival_date, holidays = holidays) %>%  
4   step_rm(arrival_date) %>%  
5   step_rm(country)
```

```
1 rf_work = workflow() %>%  
2   add_model(rf_model) %>%  
3   add_recipe(rf_recipe)
```

Tuning

```
1 rf_work %>%  
2   parameters()
```

Collection of 2 parameters for tuning

identifier	type	object
mtry	mtry	nparam[?]
min_n	min_n	nparam[+]

Model parameters needing finalization:

```
# Randomly Selected Predictors ('mtry')
```

See `?dials::finalize`` or `?``

`dials::update.parameters`` for more information.

```
1 rf_grid = rf_work %>%  
2   tune_grid(  
3     hotel_vf,  
4     grid = 10,  
5     control = control_grid(save_pred = TRUE),  
6     metrics = metric_set(roc_auc)  
7   )
```

“Best” parameters

```
1 rf_grid %>%  
2   show_best(metric = "roc_auc")
```

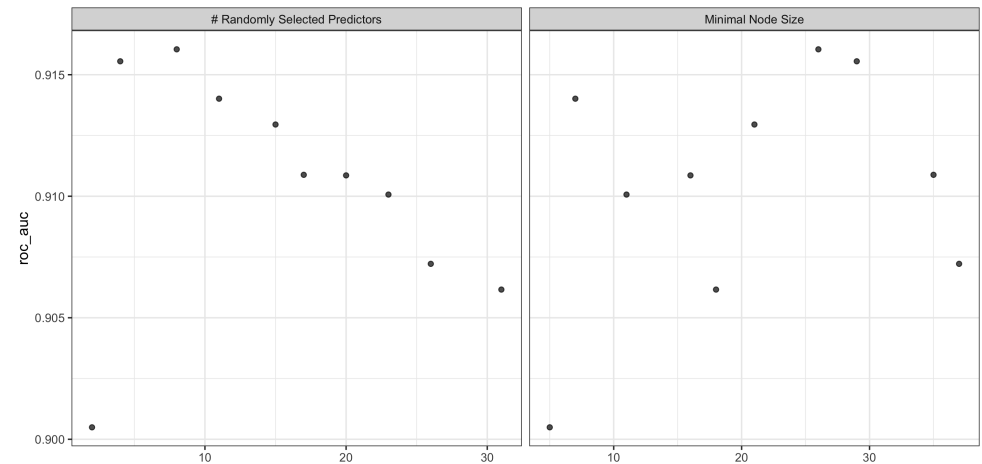
A tibble: 5 × 8

	mtry	min_n	.metric	.estim...	mean	n	std_err
	<int>	<int>	<chr>	<chr>	<dbl>	<int>	<dbl>
1	8	26	roc_auc	binary	0.916	5	0.00172
2	4	29	roc_auc	binary	0.916	5	0.00190
3	11	7	roc_auc	binary	0.914	5	0.00182
4	15	21	roc_auc	binary	0.913	5	0.00118
5	17	35	roc_auc	binary	0.911	5	0.00191

... with 1 more variable: .config <chr>, and

abbreviated variable name ¹.estimator

```
1 autoplot(rf_grid)
```



Refitting

```
1 (rf_best = rf_grid %>%  
2   select_best(metric = "roc_auc"))
```

```
# A tibble: 1 × 3  
  mtry min_n .config  
  <int> <int> <chr>  
1     8    26 Preprocessor1_Model06
```

```
1 rf_work_tuned = update_model(  
2   rf_work,  
3   rand_forest(  
4     trees=100,  
5     mtry = rf_best$mtry,  
6     min_n = rf_best$min_n  
7   ) %>%  
8     set_engine("ranger", num.threads = 8) %>%  
9     set_mode("classification")  
10 )  
11  
12 rf_fit = rf_work_tuned %>%  
13   fit(data=hotel_train)
```

Test Performance (out-of-sample)

```
1 rf_test_perf = rf_fit %>%
2   augment(new_data = hotel_test) %>%
3   select(children, starts_with(".pred"))
4
5 rf_test_perf %>%
6   conf_mat(children, .pred_class)
```

	Truth	
Prediction	children	none
children	402	69
none	609	11420

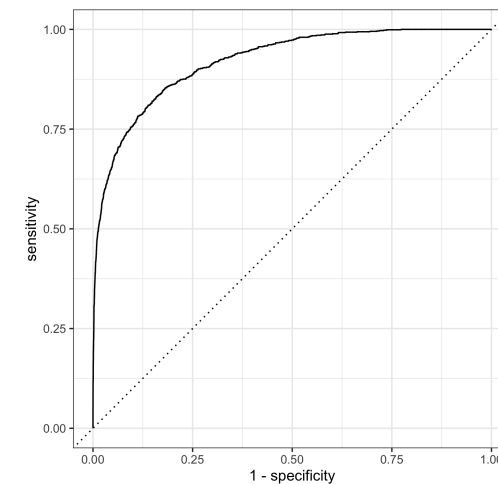
```
1 rf_test_perf %>%
2   precision(children, .pred_class)
```

```
# A tibble: 1 × 3
  .metric .estimator .estimate
  <chr>    <chr>         <dbl>
1 precision binary      0.854
```

```
1 rf_test_perf %>%
2   roc_auc(children, .pred_children)
```

```
# A tibble: 1 × 3
  .metric .estimator .estimate
  <chr>    <chr>         <dbl>
1 roc_auc binary      0.920
```

```
1 rf_roc = rf_test_perf %>%
2   yardstick::roc_curve(
3     children,
4     .pred_children
5   ) %>%
6   mutate(name = "RF - test")
7 rf_roc %>%
8   autoplot()
```



Comparing models

