

Applied Machine Learning - Basic Principles

Max Kuhn (RStudio)

Load Packages



Four large, empty light gray rectangular boxes stacked vertically, intended for code or content.

Introduction

In this section, we will introduce concepts that are useful for any type of machine learning model:

- versus the model
- data splitting
- resampling
- tuning parameters and overfitting
- model tuning

Many of these topics will be put into action in later sections.

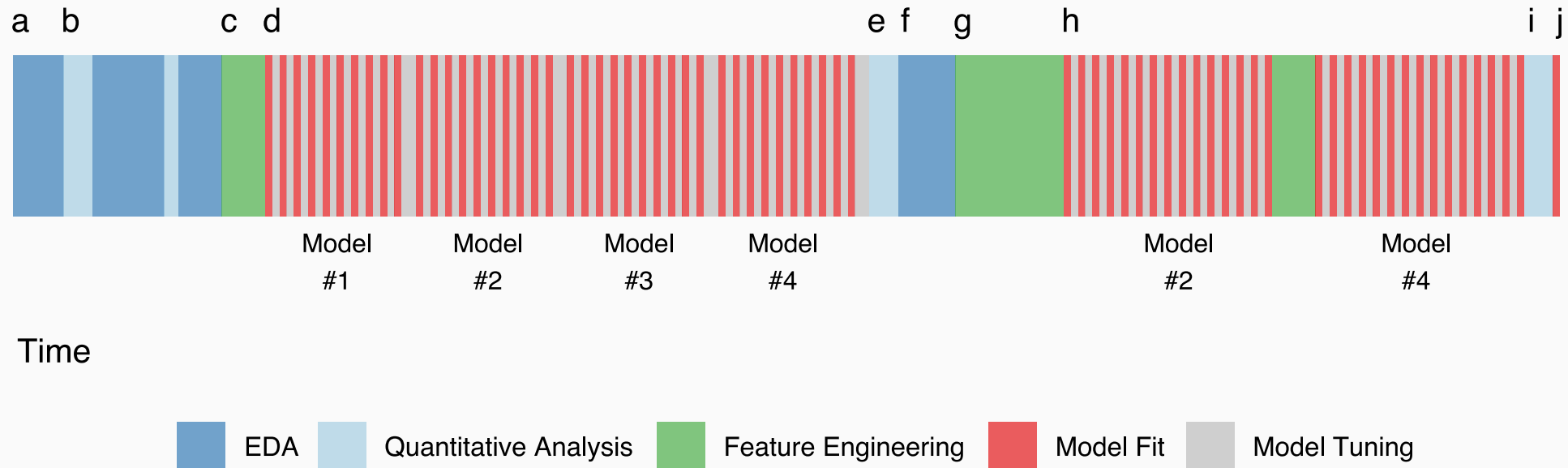
The Modeling

Common steps during model building are:

- estimating model parameters (i.e. training models)
- determining the values of `model.get_params()` that cannot be directly calculated from the data
- model selection (within a model type) and model comparison (between types)
- calculating the performance of the final model that will generalize to new data

Many books and courses portray predictive modeling as a short sprint. A better analogy would be a marathon or campaign (depending on how hard the problem is).

What the Modeling Process Usually Looks Like



Data Usage

Data Splitting and Spending

How do we "spend" the data to find an optimal model?

We split data into training and test data sets:

- : these data are used to estimate model parameters and to pick the values of the complexity parameter(s) for the model.
- : these data can be used to get an independent assessment of model efficacy. They should not be used during model training.

Data Splitting and Spending

The more data we spend, the better estimates we'll get (provided the data is accurate).

Given a fixed amount of data:

- too much spent in training won't allow us to get a good assessment of predictive performance. We may find a model that fits the training data very well, but is not generalizable (overfitting)
- too much spent in testing won't allow us to get a good assessment of model parameters

Statistically, the best course of action would be to use all the data for model building and use statistical methods to get good estimates of error.

From a non-statistical perspective, many consumers of complex models emphasize the need for an untouched set of samples to evaluate performance.

Large Data Sets

When a large amount of data are available, it might seem like a good idea to put a large amount into the training set. , I think that this causes more trouble than it is worth due to diminishing returns on performance and the added cost and complexity of the required infrastructure.

Alternatively, it is probably a better idea to reserve good percentages of the data for specific parts of the modeling process. For example:

- Save a large chunk of data to perform feature selection prior to model building
- Retain data to calibrate class probabilities or determine a cutoff via an ROC curve.

Also, there may be little need for iterative resampling of the data. A single holdout (aka validation set) may be sufficient in some cases if the data are large enough and the data sampling mechanism is solid.

Mechanics of Data Splitting

There are a few different ways to do the split: simple random sampling, by date, or methods that focus on the distribution of the predictors.

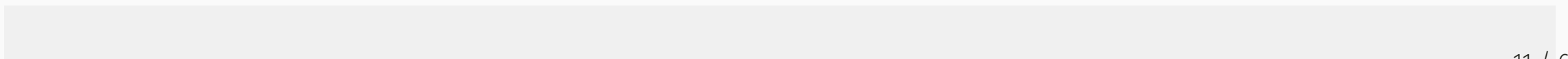
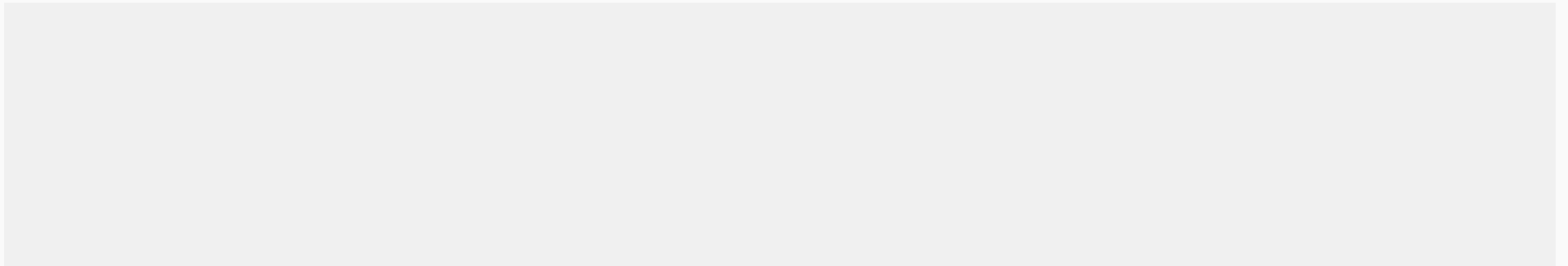
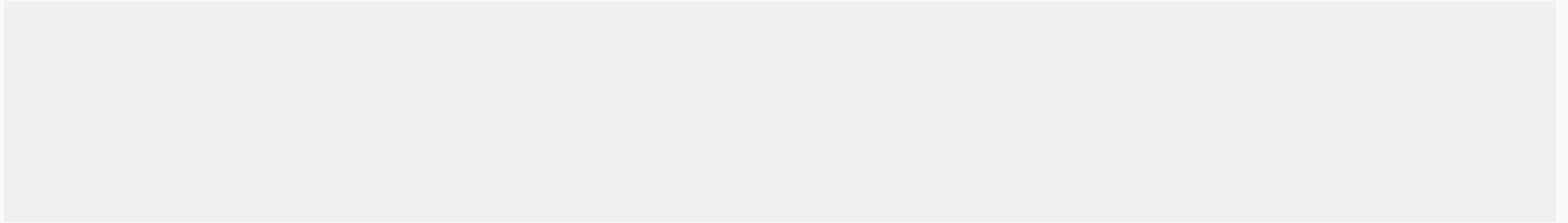
For stratification:

- **Class-based stratification**: this would mean sampling within the classes to preserve the distribution of the outcome in the training and test sets
- **Quantile-based stratification**: determine the quartiles of the data set and sample within those artificial groups

Ames Housing Data



Let's load the example data set and split it. We'll put 75% into training and 25% into testing.



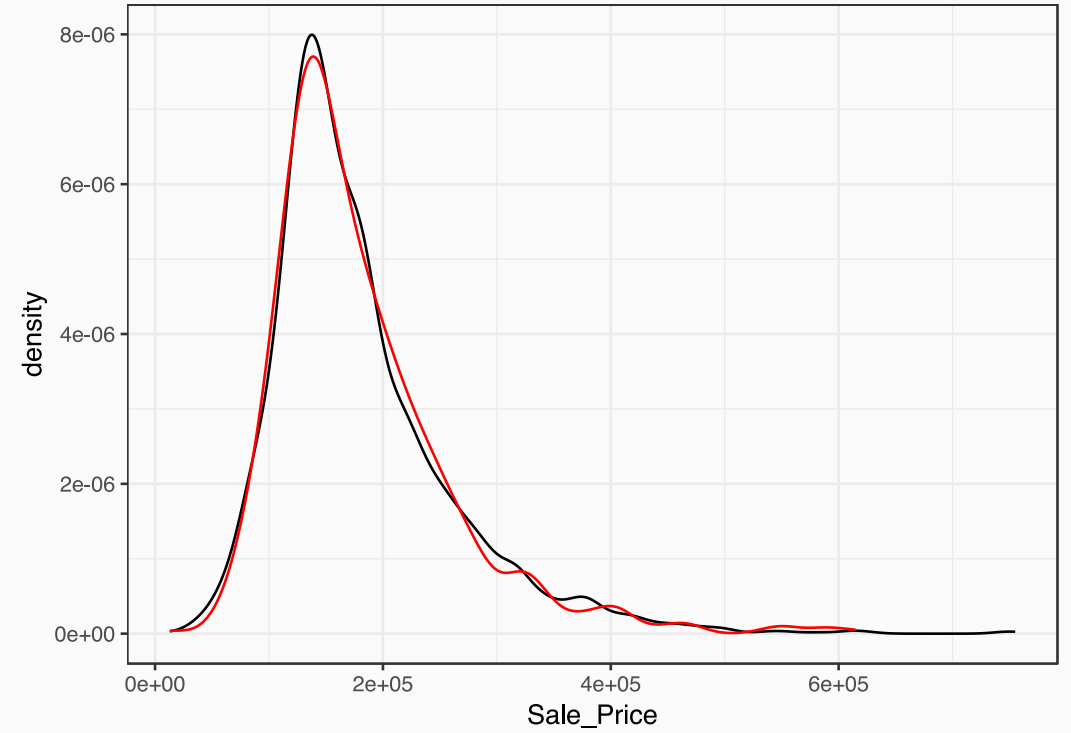
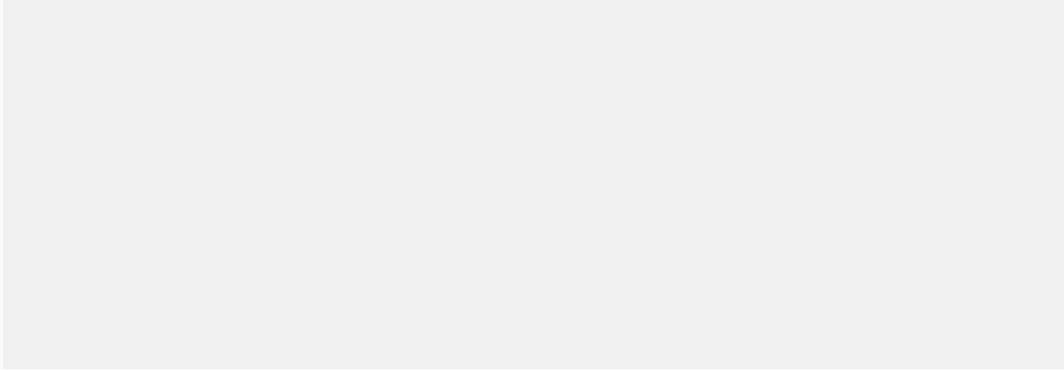
Ames Housing Data



What do these objects look like?

Four large, empty light gray rectangular boxes are stacked vertically, intended for displaying the visual representation of different data objects.

Outcome Distributions



Creating Models in R

Specifying Models in R Using Formulas

To fit a model to the housing data, the model terms must be specified. Historically, there are two main interfaces for doing this.

The `<interface>` using R **formula rules** to specify a `<representation>` of the terms:

Variables + interactions



Shorthand for all predictors



Inline functions / transformations



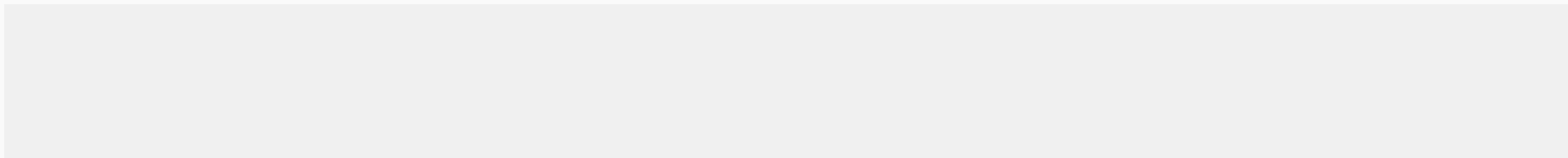
This is very convenient but it has some disadvantages.

Downsides to Formulas

- You can't nest in-line functions such as
.
- All the model matrix calculations happen at once and can't be recycled when used in a model function.
- For very data sets, the formula method can be **extremely inefficient**.
- There are limited that variables can take which has led to several re-implementations of formulas.
- Specifying multivariate outcomes is clunky and inelegant.
- Not all modeling functions have a formula method (consistency!).

Specifying Models Without Formulas

Some modeling functions have a non-formula (XY) interface. This usually has arguments for the predictors and the outcome(s):



This is inconvenient if you have transformations, factor variables, interactions, or any other operations to apply to the data prior to modeling.

Overall, it is difficult to predict if a package has one or both of these interfaces. For example, `glmnet` only has formulas.

There is a `glmnet` package, using `glmnet` that will be discussed later that solves some of these issues.

A Linear Regression Model



Let's start by fitting an ordinary linear regression model to the training set. You can choose the model terms for your model, but I will use a very simple model:

Before looking at coefficients, we should do some model checking to see if there is anything obviously wrong with the model.

To get the statistics on the individual data points, we will use the awesome `car` package:

Hands-On: Some Basic Diagnostics

From these results, let's take 10 minutes and do some visualizations:

- Plot the observed versus fitted values
- Plot the residuals
- Plot the predicted versus residuals

Are there any to this approach?

- A tidy unified interface to models
- `lm()` isn't the only way to perform linear regression
 - `glm()` for regularized regression
 - `brglm2::brglm2()` for Bayesian regression
 - `tf.tensorflow::keras_model_compilation()` for regression using tensorflow
- But...remember the consistency slide?
 - Each interface has its own minutiae to remember
 - `tidy()` standardizes all that!

1) Create specification

2) Set the engine

3) Fit the model

```
library(parsnip)
library(gamlss)
spec <- gamlss_formula(
  y ~ 1,
  family = "GAM"
)
```

```
spec <- set_engine(spec, "gamlss")
```

```
model <- fit(spec, data)
```

```
model <- fit(spec, data,
  control = parsnip::control_parsnip(
    method = "fmin"
  )
)
```

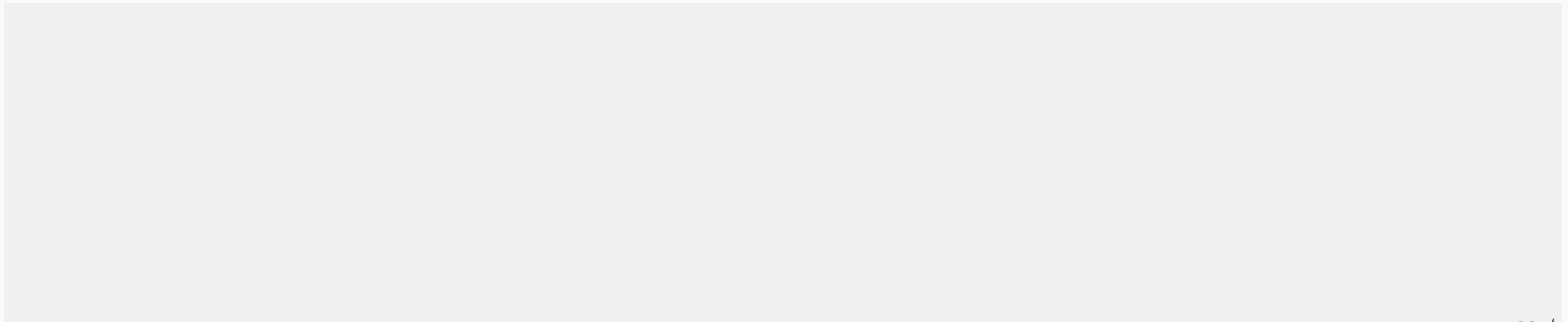
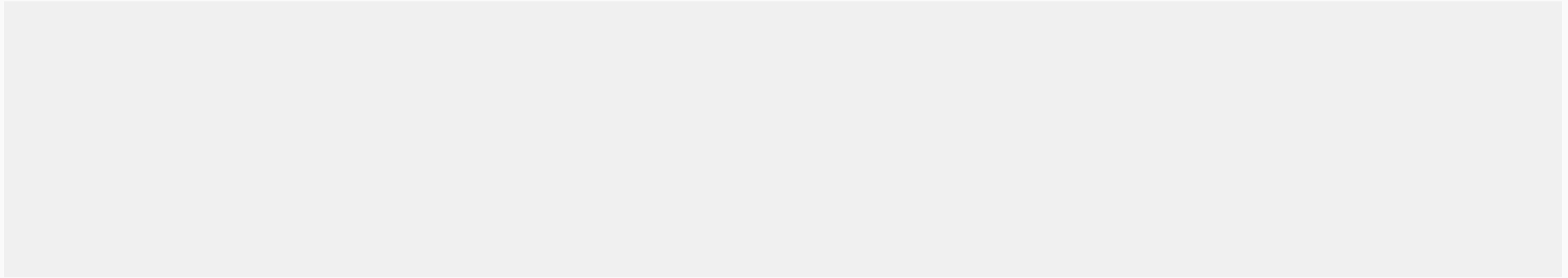
```
library(parsnip)
library(gamlss)
spec <- gamlss_formula(
  y ~ 1,
  family = "GAM"
)
```

```
spec <- set_engine(spec, "gamlss")
model <- fit(spec, data)
model <- fit(spec, data,
  control = parsnip::control_parsnip(
    method = "fmin"
  )
)
```

Different interfaces



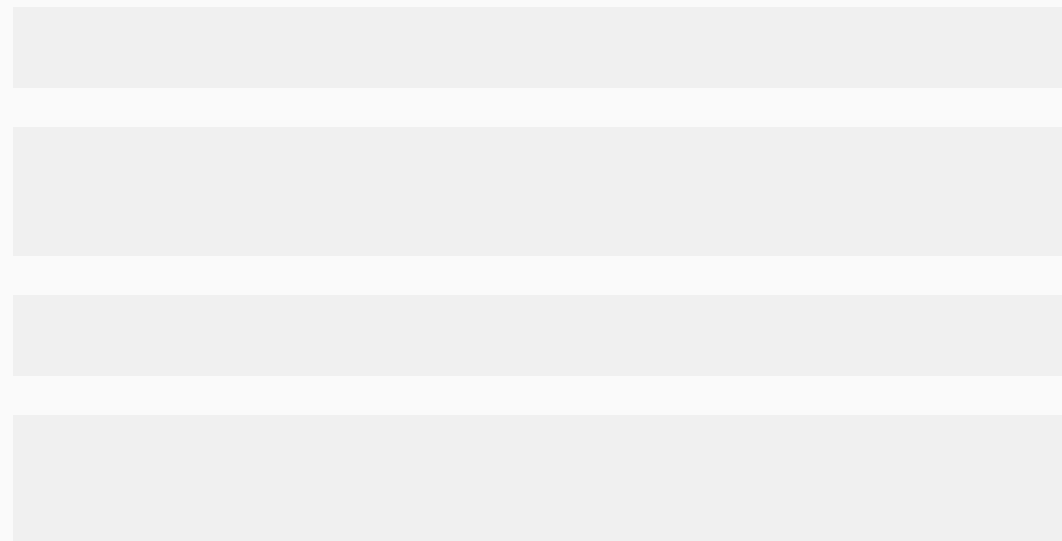
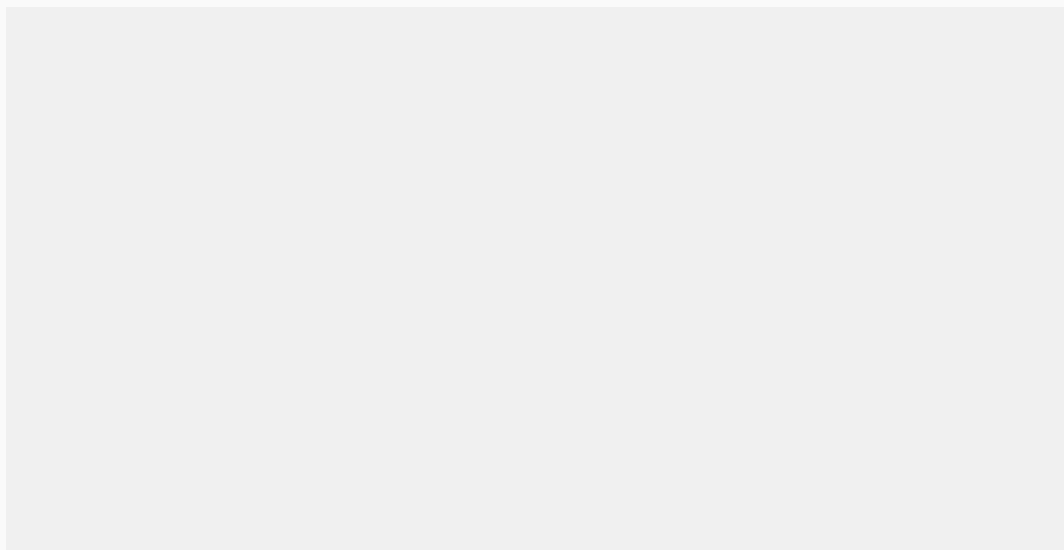
is not picky about the interface used to specify terms. Remember, only allowed the formula interface!



Alternative Engines



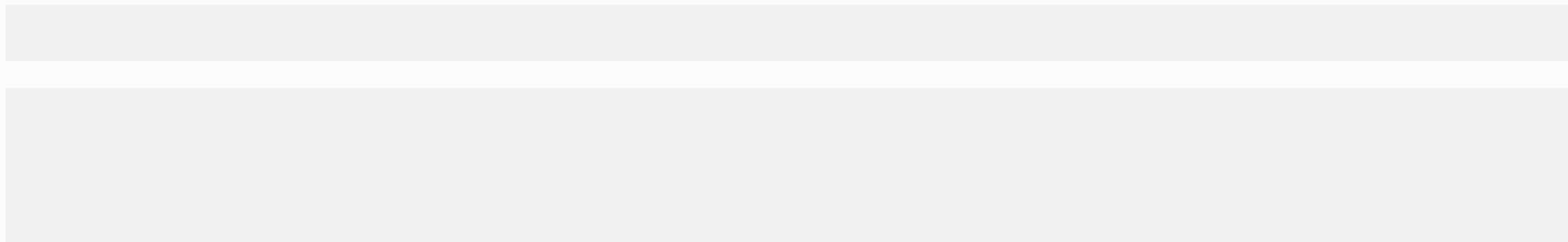
With `engine = "stan"`, it is easy to switch to a different engine, like Stan, to run the same model with alternative backends.



Model Evaluation

Overall Model Statistics

`summary()` holds the actual model object in the `model` slot. If you use the `summary()` method on the underlying `model` object, the bottom shows some statistics:



These statistics are generated from `summary()`. This is problematic because it can lead to optimistic results, especially for flexible models (overfitting).

Overall Model Statistics

holds the actual model object in the slot. If you use the method on the underlying object, the bottom shows some statistics:

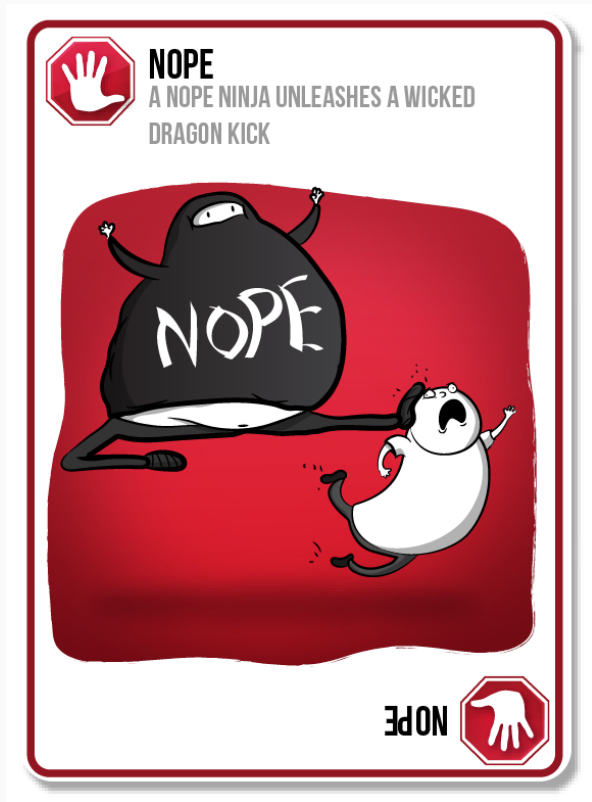
```
##> A tibble: 1 x 10
##>   <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
##> 1  0.78  0.78  0.78  0.78  0.78  0.78  0.78  0.78  0.78  0.78
```

```
##> A tibble: 1 x 10
##>   <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
##> 1  0.78  0.78  0.78  0.78  0.78  0.78  0.78  0.78  0.78  0.78
```

These statistics are generated from . This is problematic because it can lead to optimistic results, especially for flexible models (overfitting).

Idea!

The tests set is used for assessing performance. and use those results to estimate these statistics?



(Matthew Inman/Exploding Kittens)

Assessing Models

until the very end when you have one or two models that are your favorite. We need to use the training set...but how?

Assessing Models

until the very end when you have one or two models that are your favorite. We need to use the training set...but how?

1) For model A, fit on training set, predict on training set

2) For model B, fit on training set, predict on training set

3) Compare performance

Assessing Models

until the very end when you have one or two models that are your favorite. We need to use the training set...but how?

- 1) For model A, fit on training set, predict on training set
- 2) For model B, fit on training set, predict on training set
- 3) Compare performance

For some models, it is possible to get very "good" performance by predicting the training set (it was so flexible you overfit it). That's an issue since we will need to make "honest" comparisons between models before we finalize them and run our final choices on the test set.

If only we had a method for getting honest performance estimates from the ...

Resampling Methods

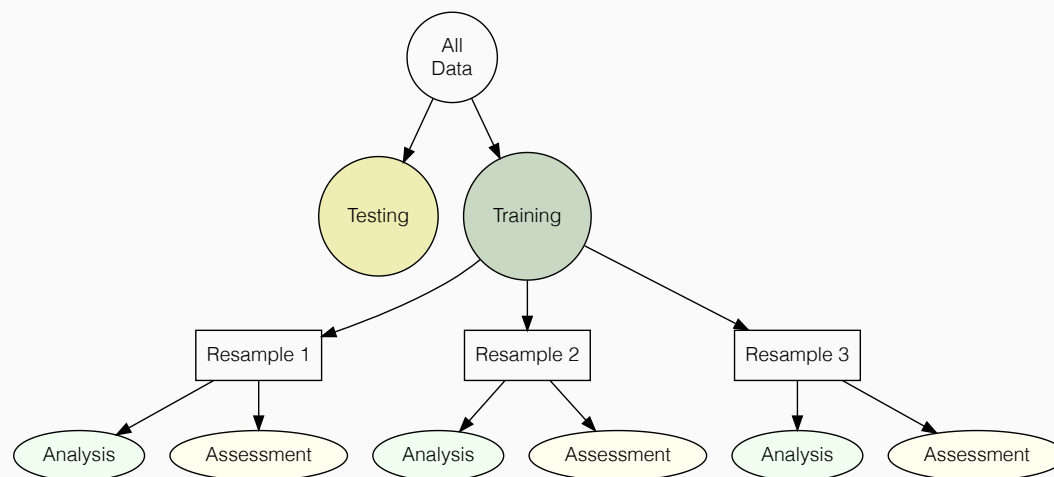
These are additional data splitting schemes that are applied to the `train` set.

They attempt to simulate slightly different versions of the training set. These versions of the original are split into two model subsets:

- The `train` is used to fit the model (analogous to the training set).
- Performance is determined using the `test`.

This process is repeated many times.

There are different flavors of resampling but we will focus on two methods.

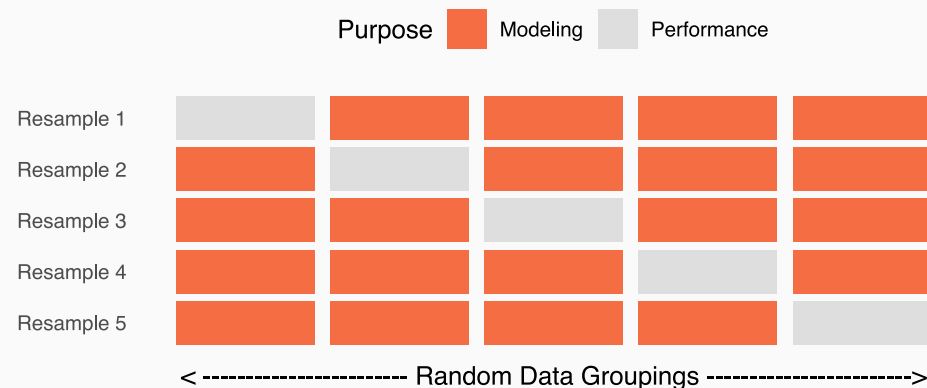


V-Fold Cross-Validation

Here, we randomly split the training data into distinct blocks of roughly equal size.

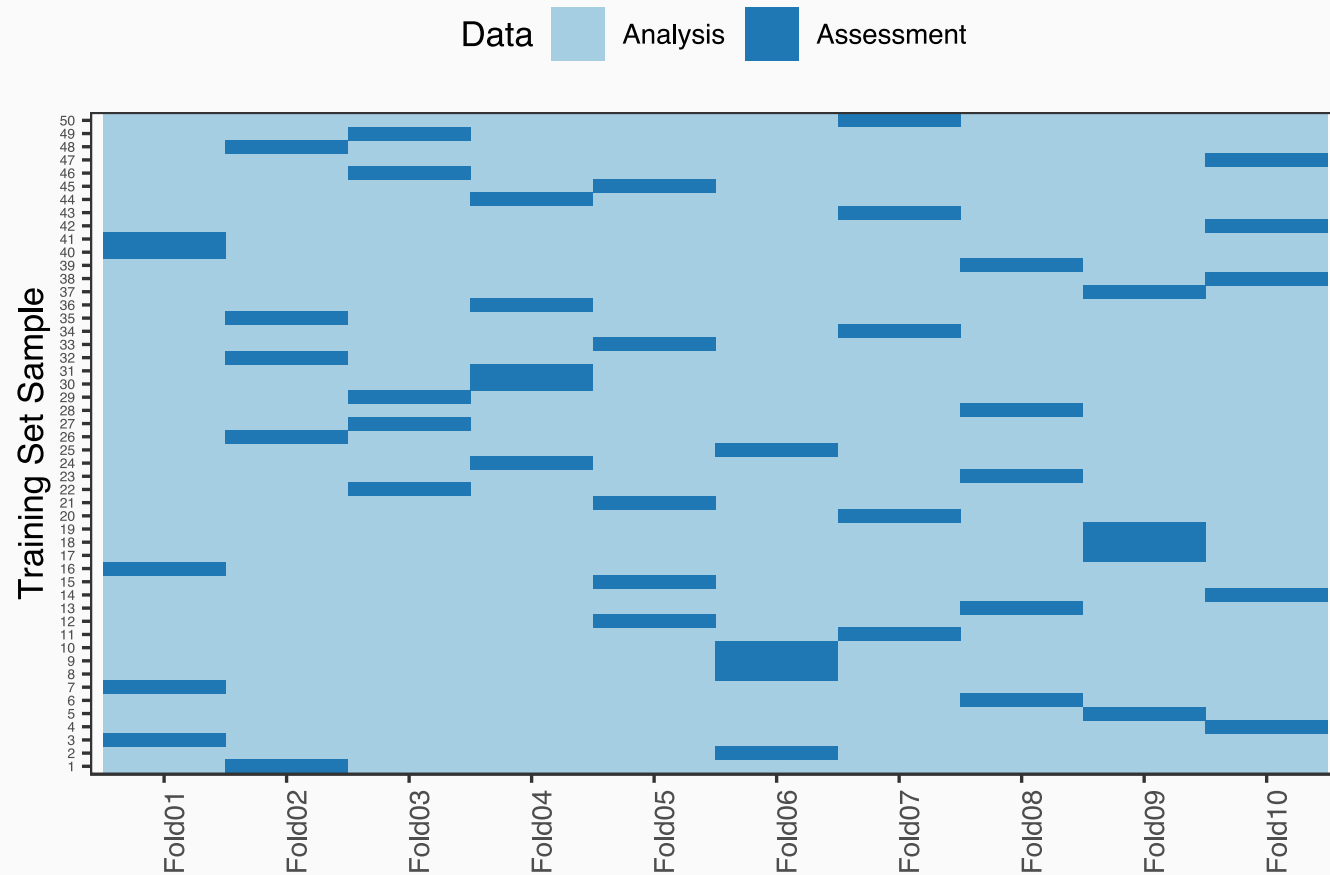
- We leave out the first block of analysis data and fit a model.
- This model is used to predict the held-out block of assessment data.
- We continue this process until we've predicted all assessment blocks

The final performance is based on the hold-out predictions by the statistics from the blocks.



is usually taken to be 5 or 10 and leave one out
cross-validation has each sample as a block.

10-Fold Cross-Validation with $n = 50$



Bootstrapping

A bootstrap sample is the `train` as the training set but each data point is selected `n` times with replacement.

-

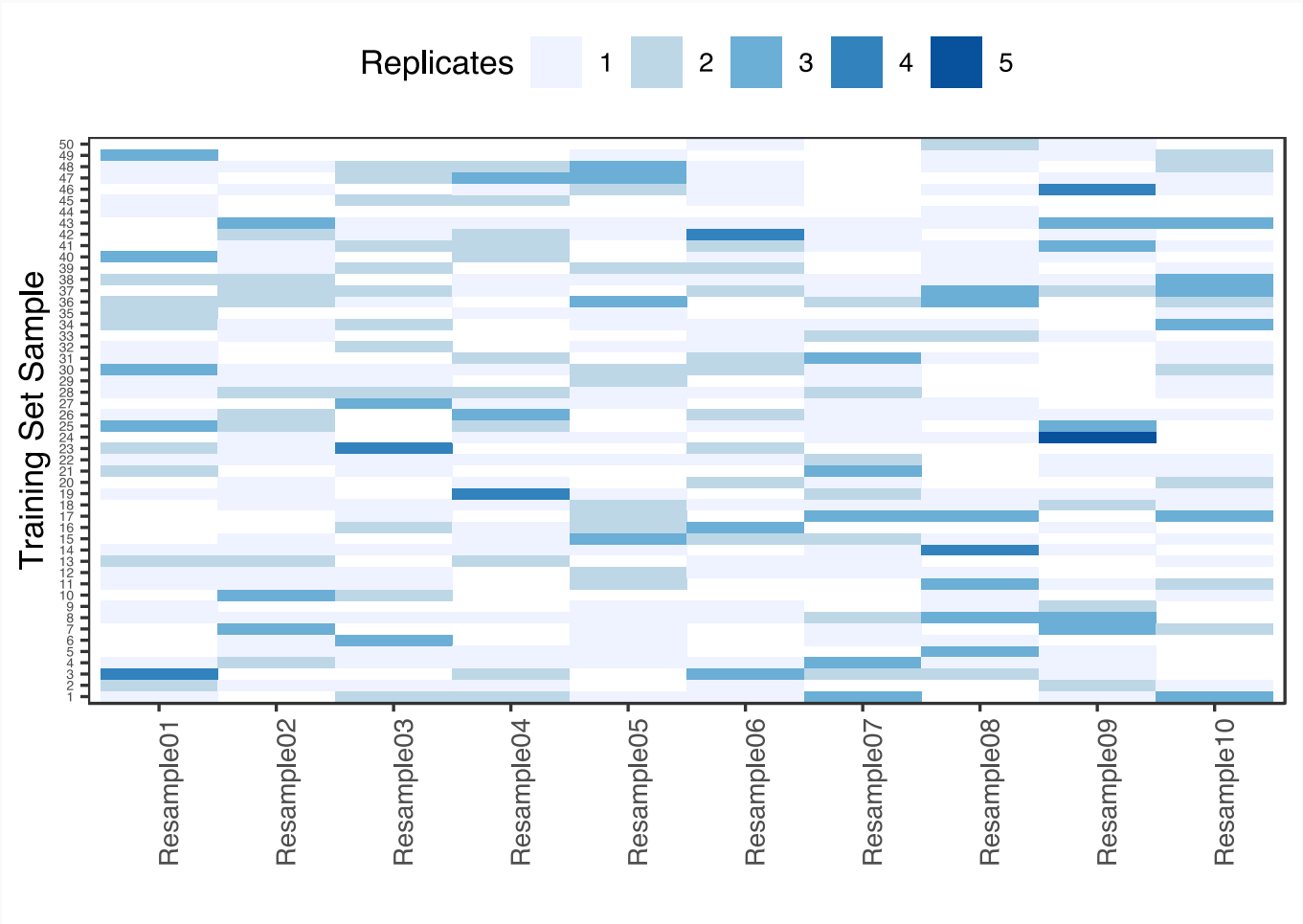
Will contain more than one replicate of a training set instance.

-

Contains all samples that were never included in the corresponding bootstrap set. Often called the "out-of-bag" sample and can vary in size!

On average, 63.2120559% of the training set is contained `train` in the bootstrap sample.

Bootstrapping with $n = 50$



Comparing Resampling Methods

If you think of resampling in the same manner as statistical estimators (e.g. maximum likelihood), this becomes a trade-off between bias and variance:

- Variance is (mostly) driven by the number of resamples (e.g. 5-fold CV has larger variance than 10-fold).
- Bias is (mostly) related to how much data is held back. The bootstrap has large bias compared to 10-fold CV.

There are lengthy blog posts about this subject [here](#) and [here](#).

I tend to favor 5 repeats of 10-fold cross-validation unless the size of the assessment data is "large enough".

For example, 10% of the Ames training set is 219 properties and this is probably good enough to estimate the RMSE and R^2 .

Cross-Validating Using



Each individual split object is similar to the example.

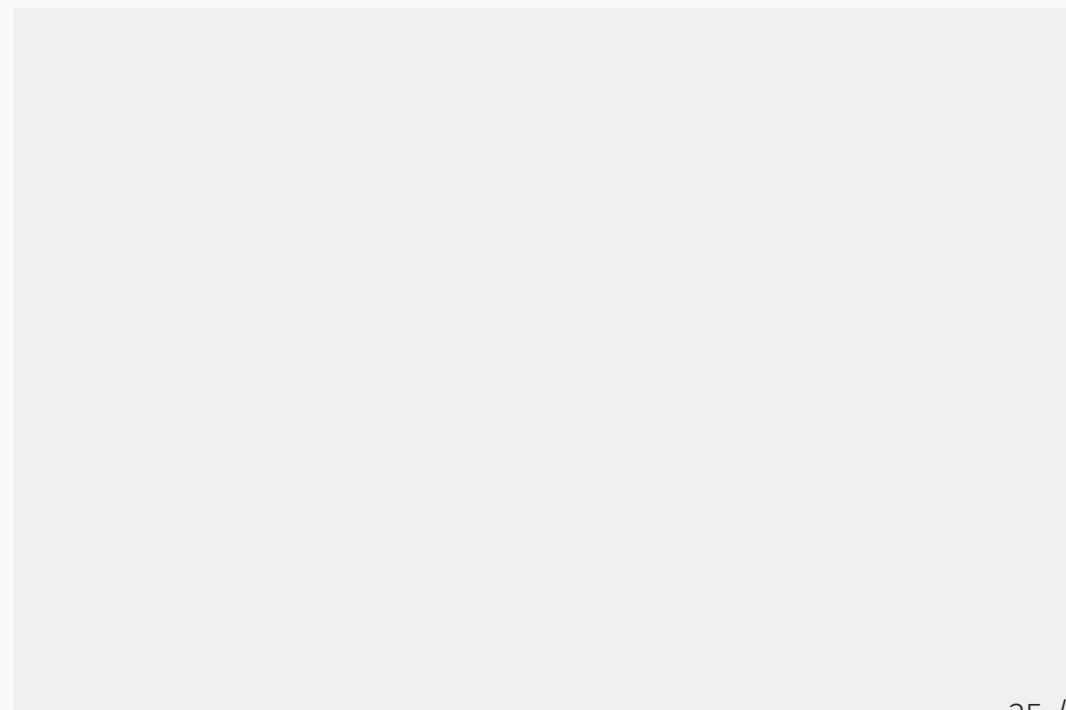
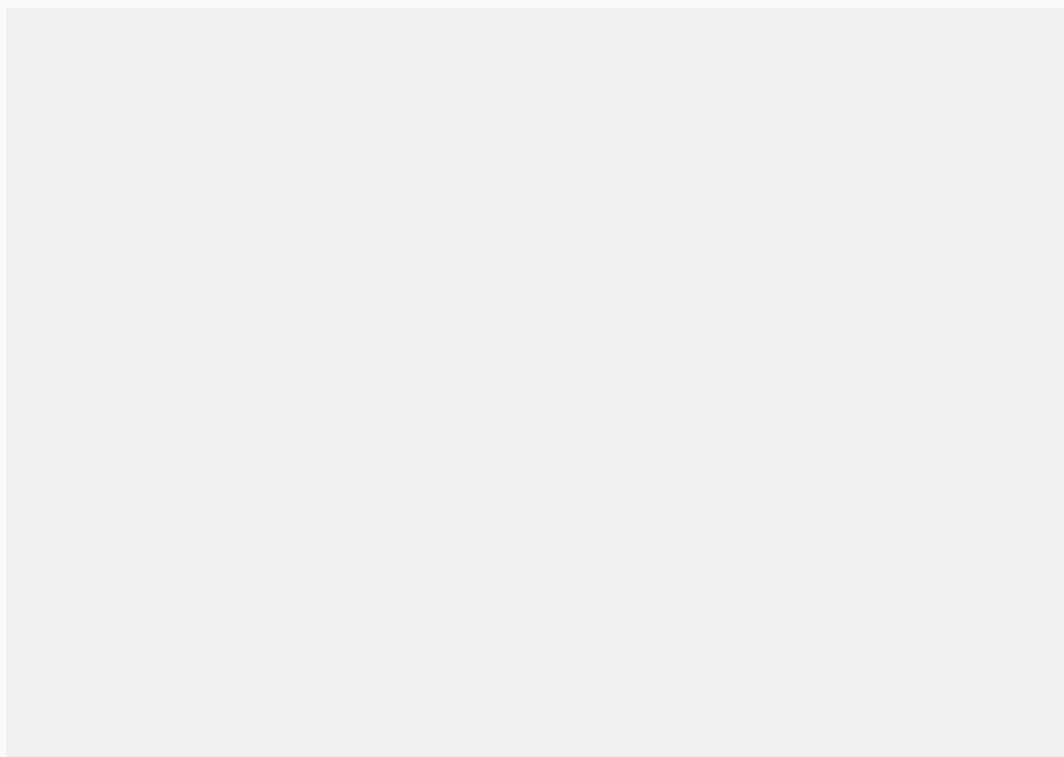


Resampling the Linear Model



Working with resample tibbles generally involves two things:

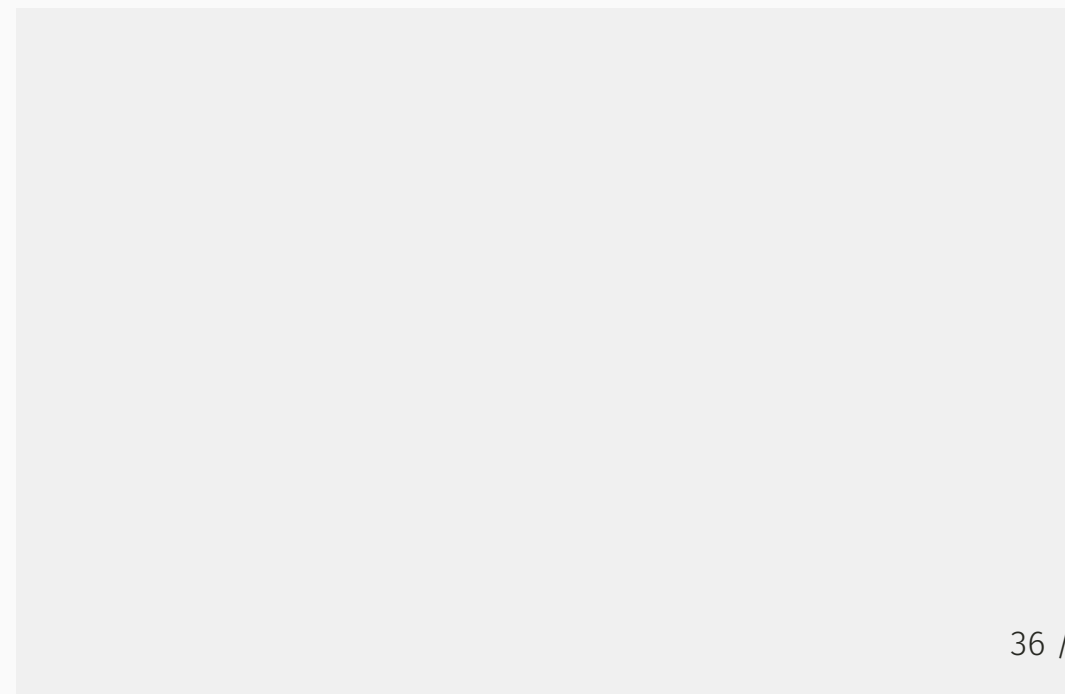
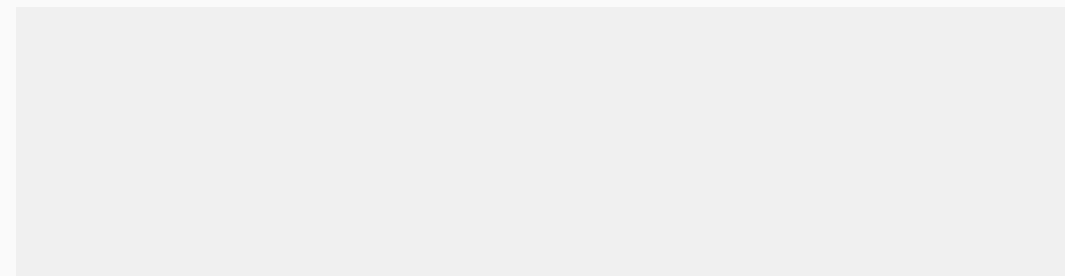
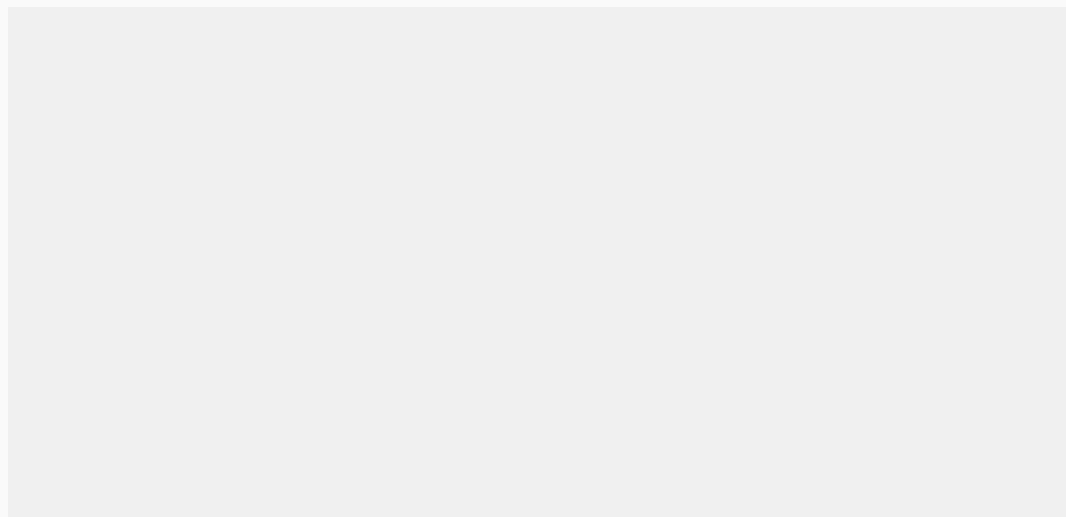
- 1) Small functions that perform an action on a single split.
- 2) The `purrr` package for `map`ping over splits.



Resampling the Linear Model



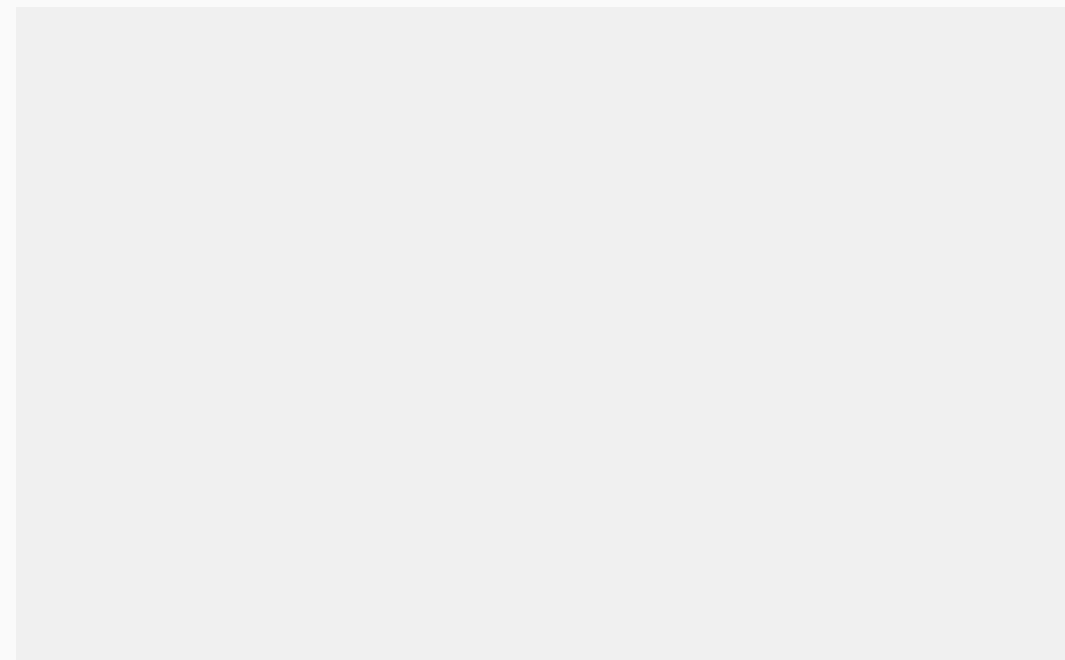
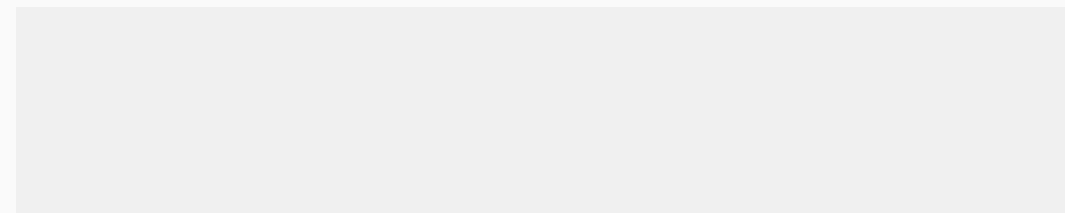
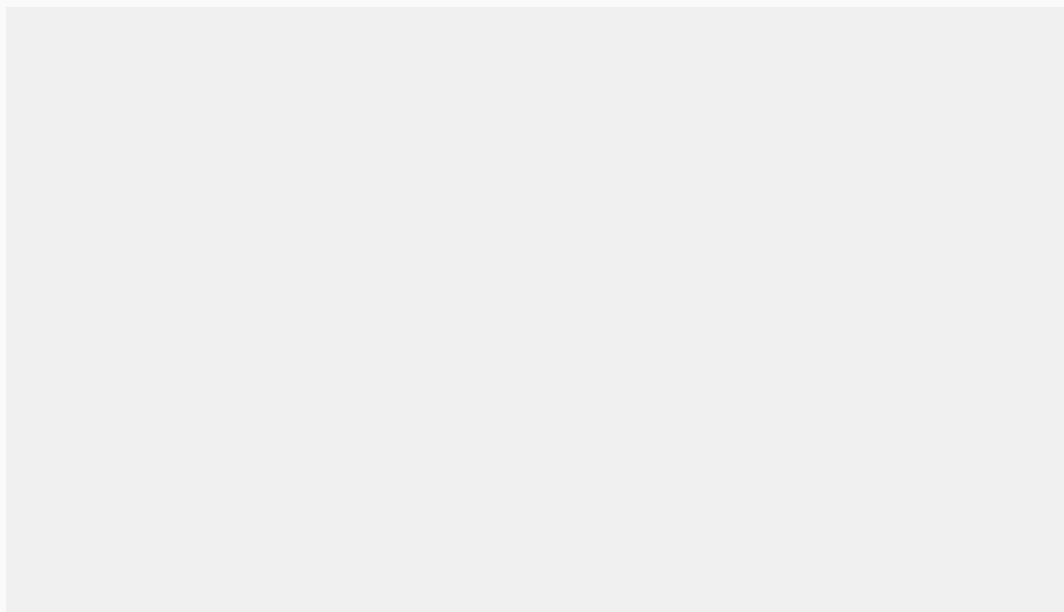
Next, we will attach the predictions for each resample:



Resampling the Linear Model



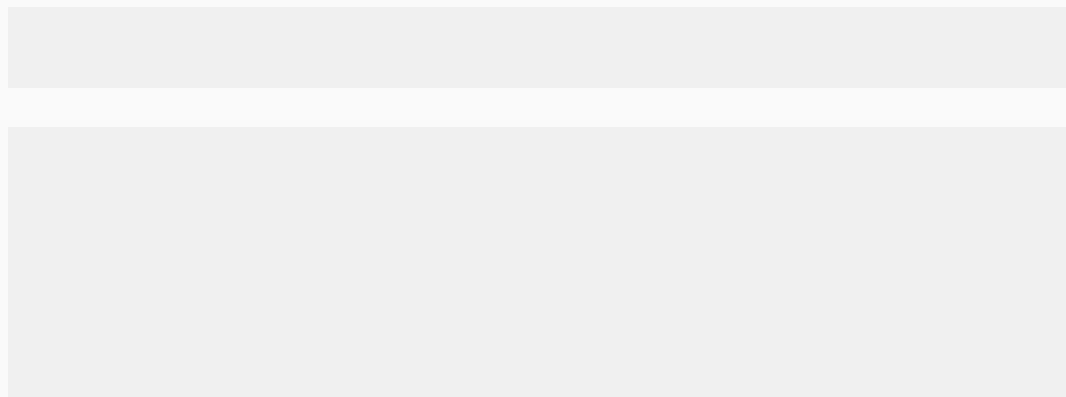
Now, let's compute two performance measures:



Resampling the Linear Model



And finally, let's compute the average of each metric over the resamples:



What Was the Ruckus?

Previously, I mentioned that the performance metrics that were naively calculated from the training set could be optimistic. However, this approach estimates the RMSE to be 0.1614 and cross-validation produced an estimate of 0.1613. What was the big deal?

Linear regression is a *linear* model. This means that it is fairly incapable at being able to adapt the underlying model function (unless it is linear). For this reason, linear regression is unlikely to *generalize* to the training set and our two estimates are likely to be the same.

We'll consider another model shortly that is *non-linear* since it can, theoretically, easily adapt to a wide variety of true model functions.

However, as before, there is also variance to consider. Linear regression is very stable since it leverages all of the data points to estimate parameters. Other methods, such as tree-based models, are not and can drastically change if the training set data is slightly perturbed.

Conclusion: the earlier concern is real but linear regression is less likely to be affected.

Diagnostics Again



Now let's look at diagnostics using the predictions from the assessment sets.



Hands-On: Partial Residual Plots



A partial residual plot is used to diagnose what variables `vars` have been in the model.

We can plot the hold-out residuals versus different variables to understand if they should have been in the model

- If the residuals have no pattern in the data, they are likely to be irrelevant.
- If a pattern is seen, it suggests that the variable should have been in the model.

Take 10 min and use `plot_partial_res()` to investigate the other predictors using the `data` data frame.
`plot_partial_res()` might come in handy.

Tuning Parameters and Overfitting

-Nearest Neighbors Model

Now let's consider a more flexible model that is : -nearest neighbors.

The model stores the training set (including the outcome).

When a new sample is predicted, training set points are found that are most similar to the new sample being predicted.

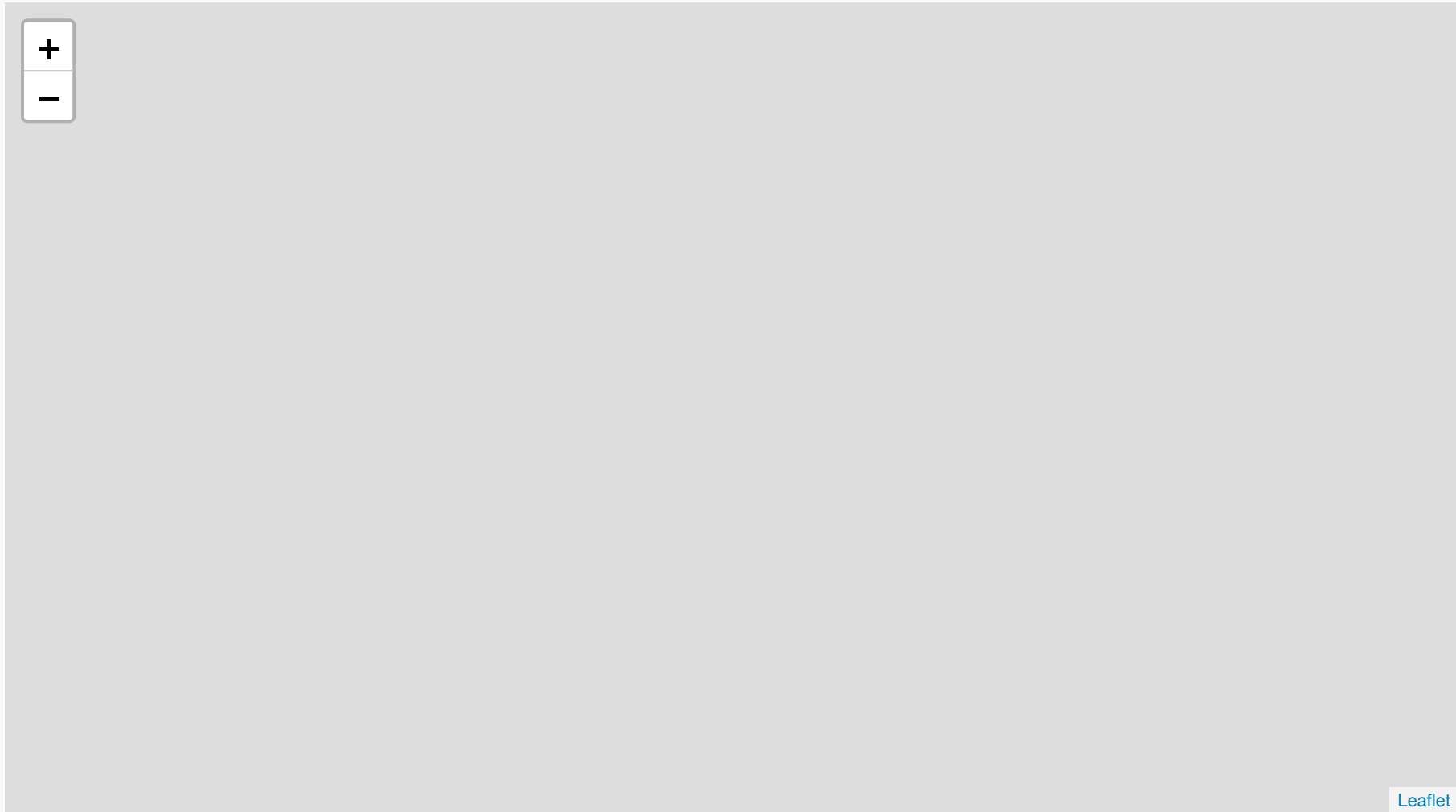
The predicted value for the new sample is some summary statistic of the neighbors, usually:

- the mean for regression, or
- the mode for classification.

When is small, the model might be responsive to the underlying data. When is large, it begins to "over smooth" the neighbors and performance suffers.

Ordinarily, since we are computing a distance, we would want to center and scale the predictors. Our two predictors are already on the same scale so we can skip this step.

-Nearest Neighbors Model



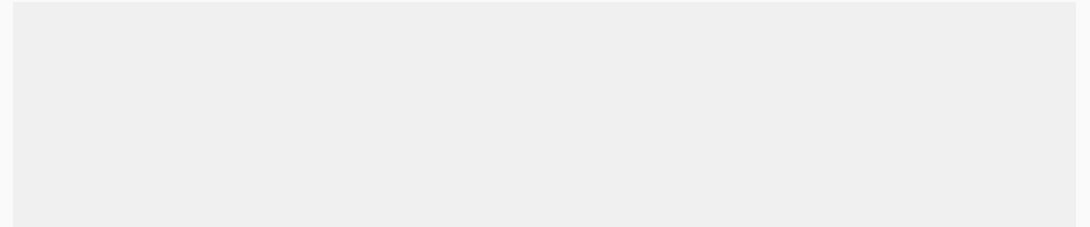
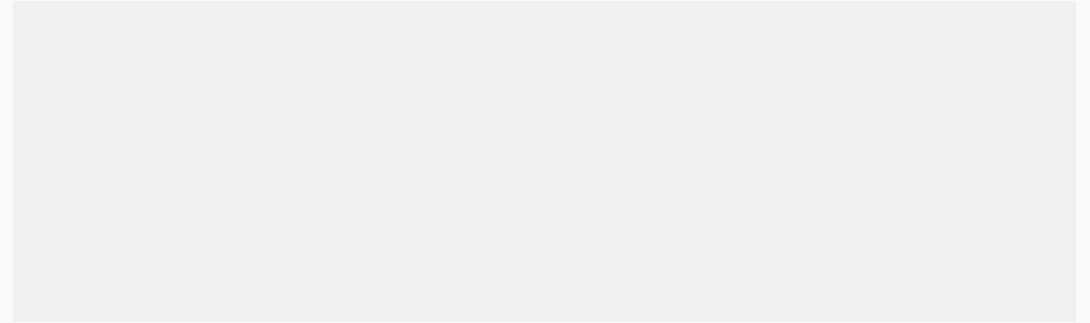
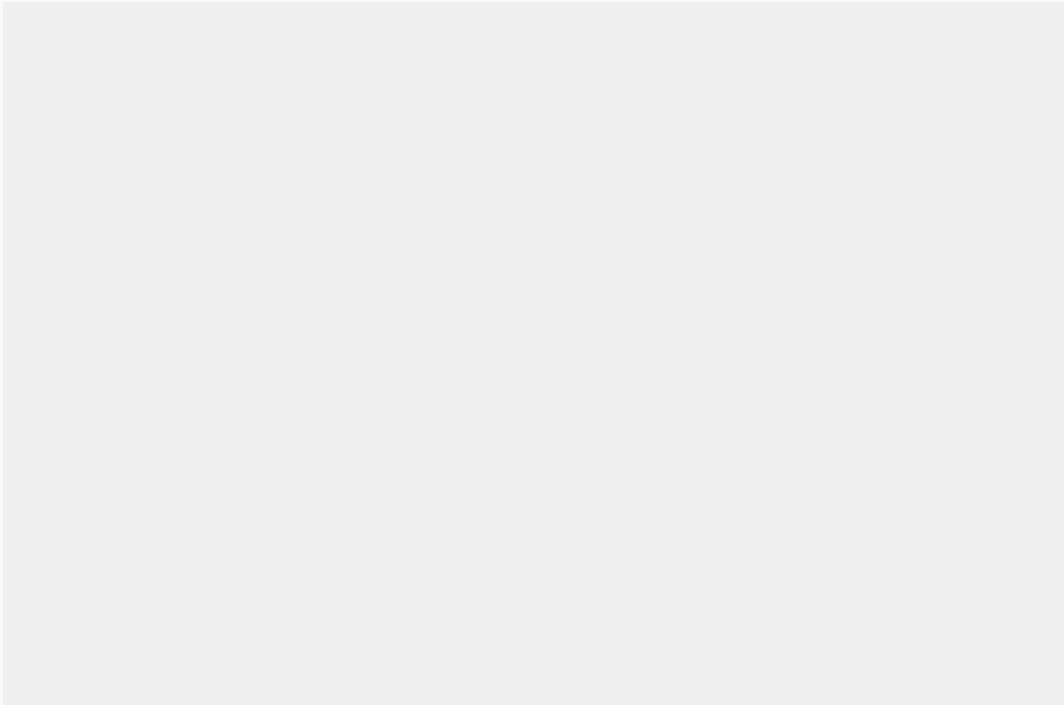
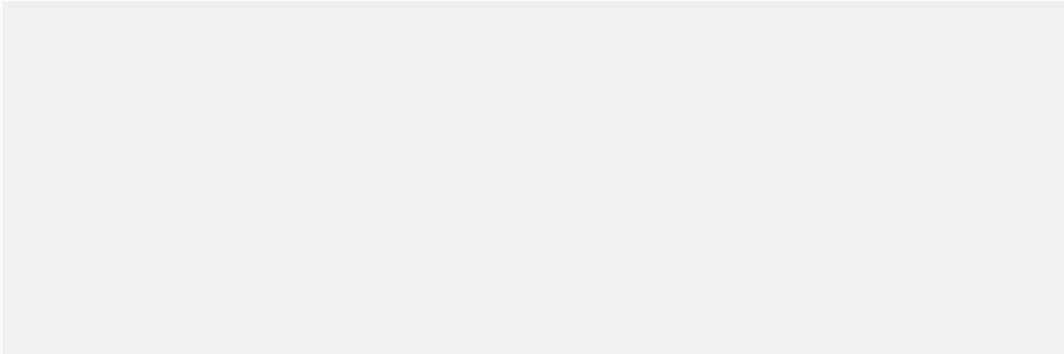
-Nearest Neighbors Model



Consider the 2-nearest neighbor model. Would there be a difference in the estimated model performance between re-prediction and cross-validation?

_____ has a _____ specification that uses the _____ package. _____ is standardized to the name _____ and we will use that going forward.

-Nearest Neighbors Model

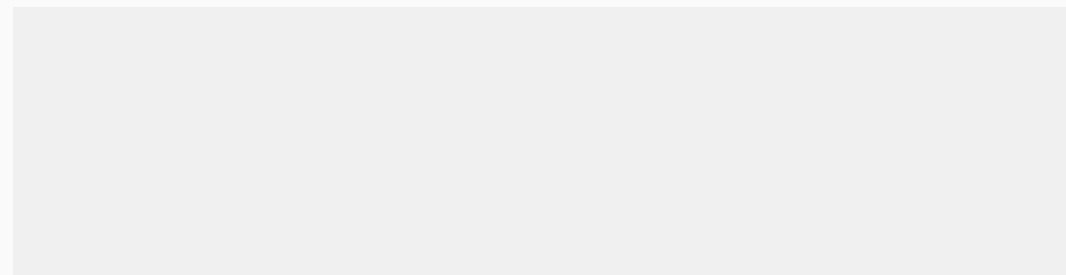
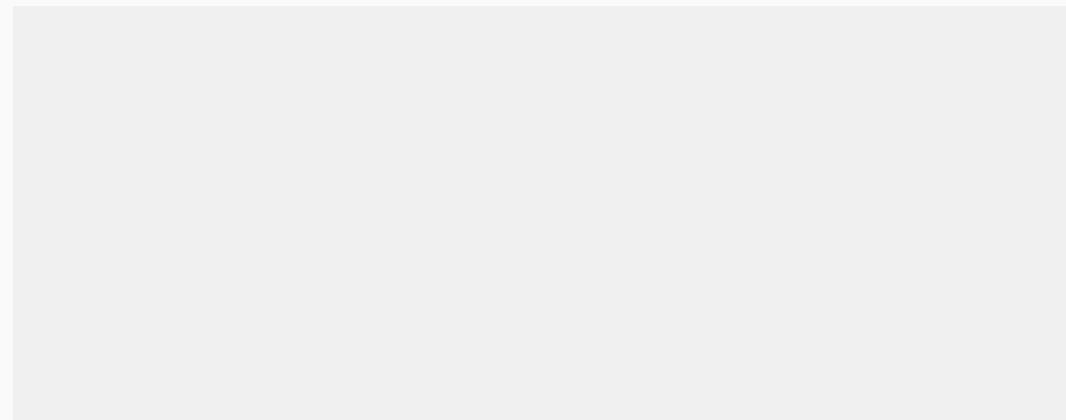
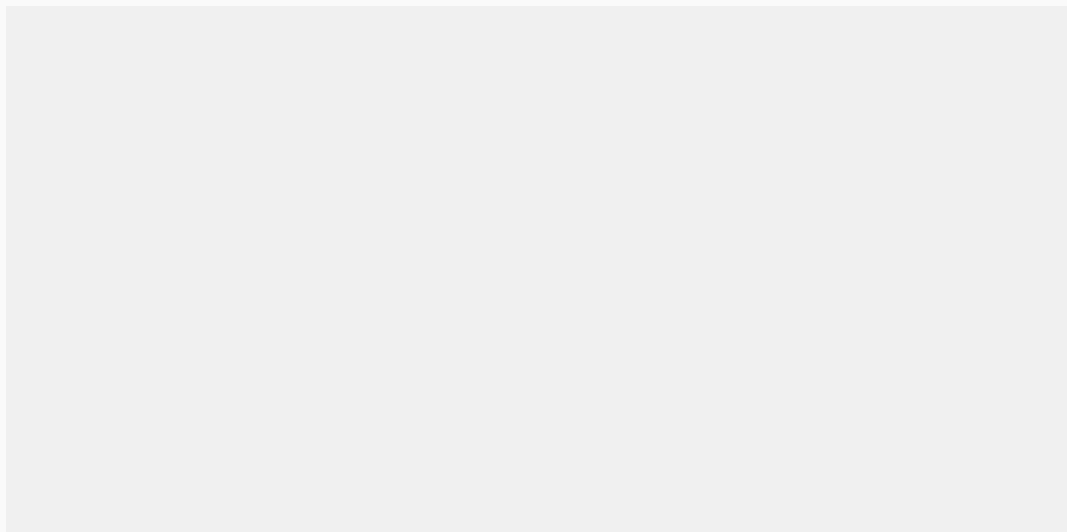


Resampling a 2-Nearest Neighbor Model



That's pretty good but are we tricking ourselves? One of those two neighbors is always itself...

Let's follow the same resampling process as before, reusing some of our other functions for generating the models, predictions, and performance metrics:



Making Formal Comparisons

The model appears to be a drastic improvement over simple linear regression but we are definitely getting highly optimistic results by re-predicting the training set.

We can try to make a more formal assessment of the two current models.

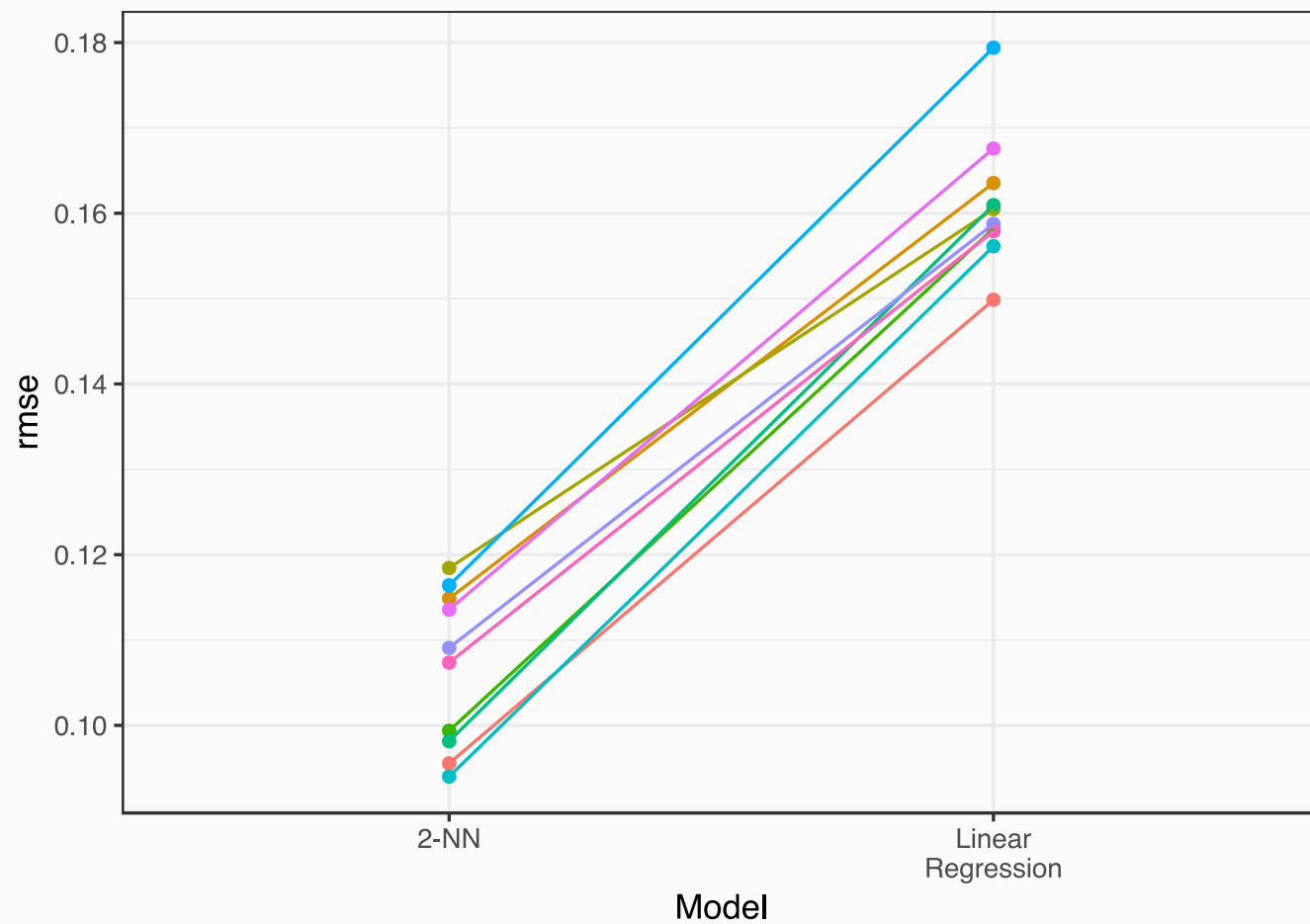
Both models used the `resamples` function, so we have 10 estimates of performance that are matched.

Does the matching mean anything?

Most likely `TRUE`. It is very common to see that there is a resample effect. Similar to repeated measures designs, we can expect a relationship between models and resamples. For example, some resamples will have the worst performance over different models and so on.

In other words, there is usually a within-resample correlation. For the two models, the estimated correlation in RMSE values is 0.671.

The Resample Effect



Model Comparison Accounting for Resampling

With only two models, a paired t -test can be used to estimate the difference in RMSE between the models:

Hothorn (2012) is the [original paper](#) on comparing models using resampling.

We'll do more extensive analyses with [this paper](#) soon.

Overfitting

Overfitting occurs when a model inappropriately picks up on trends in the training set that do not generalize to new samples.

When this occurs, assessments of the model based on the training set can show good performance that does not reproduce in future samples.

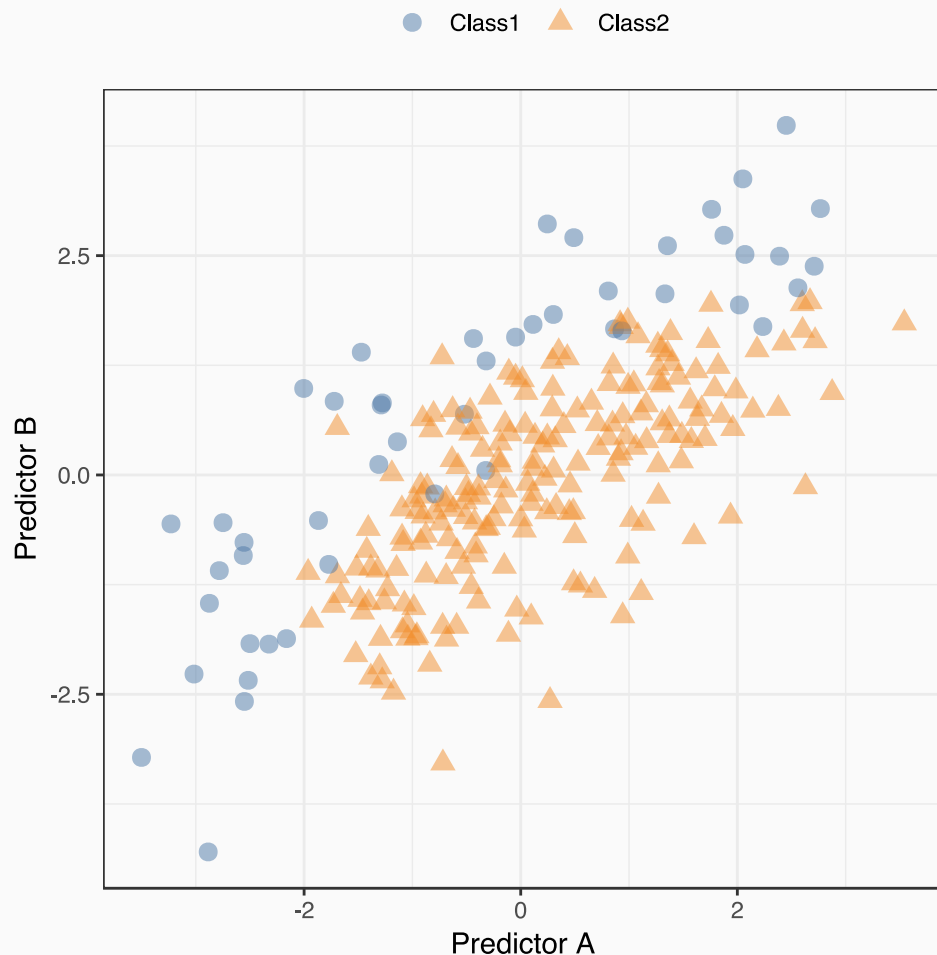
Some models have specific "knobs" to control over-fitting

- neighborhood size in nearest neighbor models is an example
- the number of splits in a tree model

Often, poor choices for these parameters can result in overfitting

For example, the next slide shows a data set with two predictors. We want to be able to produce a line (i.e. decision boundary) that differentiates two classes of data.

Two Class Example

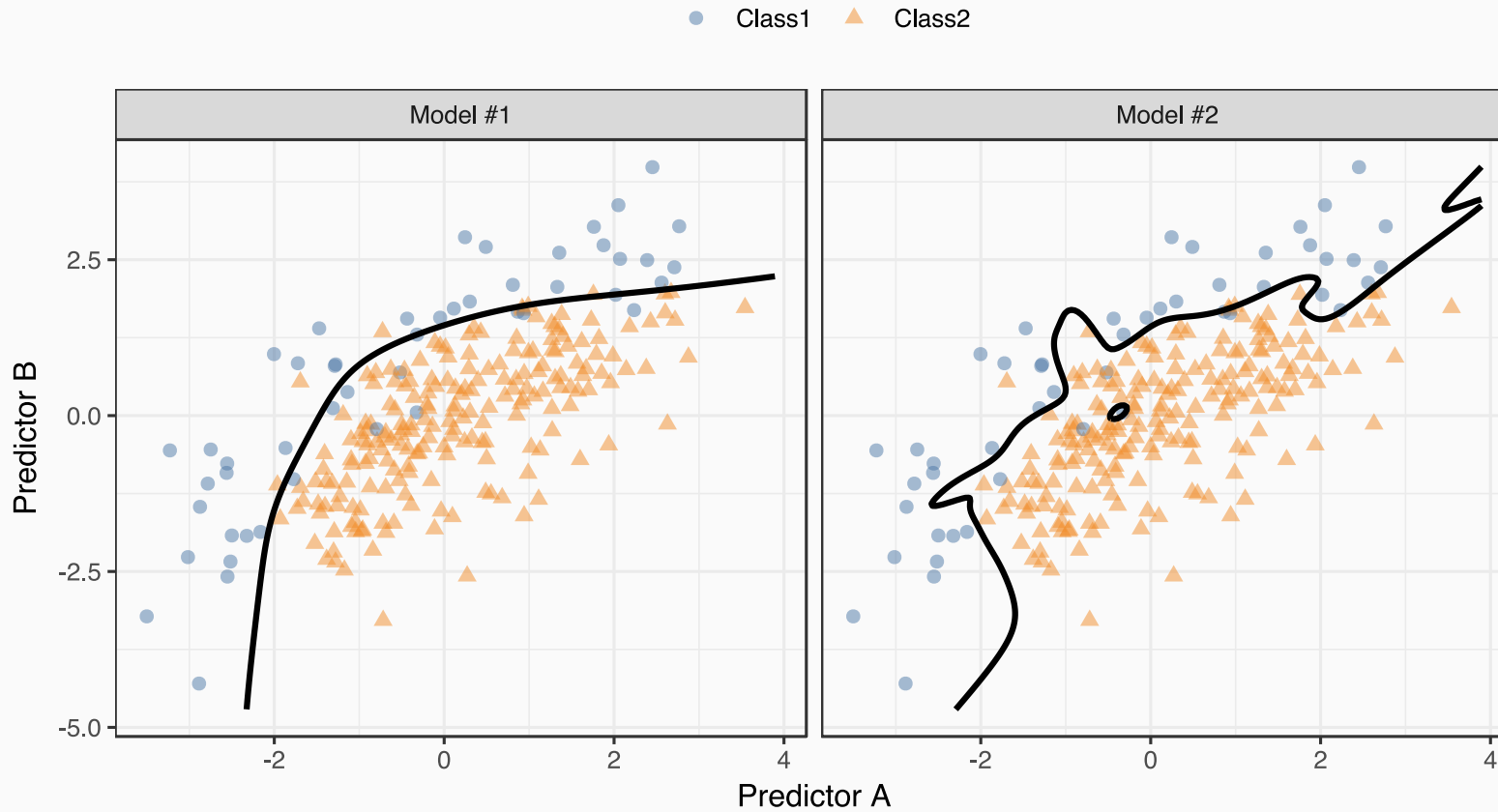


On the next slide, two classification boundaries are shown for a different model type not yet discussed.

The difference in the two panels is solely due to different choices in tuning parameters.

One overfits the training data.

Two Model Fits



Grid Search to Tune Models

We usually don't have two-dimensional data so a quantitative method for under measuring overfitting is needed. `cross_val_score` fits that description. A simple method for tuning a model is to use `GridSearchCV`:

- └─ Create a set of candidate tuning parameter values
- └─ For each resample
 - | └─ Split the data into analysis and assessment sets
 - | └─ [preprocess data]
 - | └─ For each tuning parameter value
 - | | └─ Fit the model using the analysis set
 - | | └─ Compute the performance on the assessment set and save
- └─ For each tuning parameter value, average the performance over resamples
- └─ Determine the best tuning parameter value
- └─ Create the final model with the optimal parameter(s) on the training set

`RandomizedSearchCV` is a similar technique where the candidate set of parameter values are simulated at random across a wide range. Also, an example of `GridSearchCV` can be found [here](#).

Grid Search Computations

- All of the models (except the final model) are discarded.

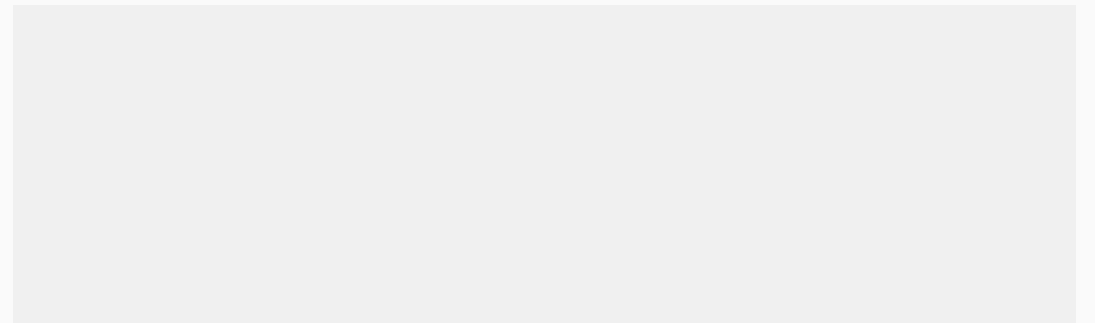
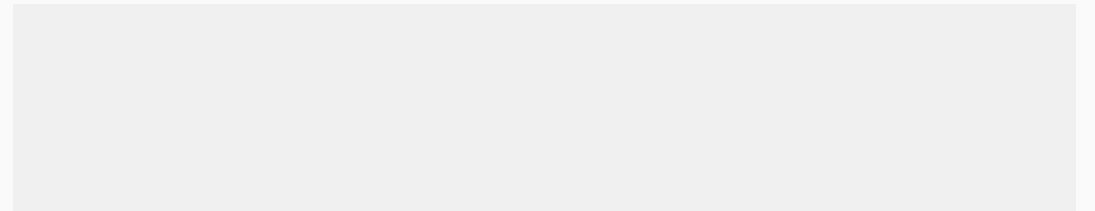
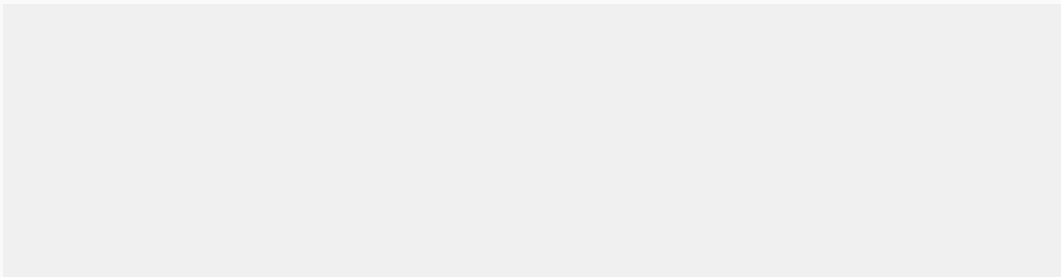
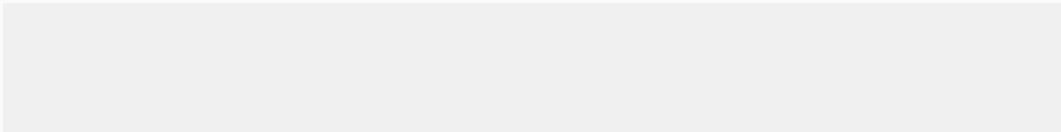
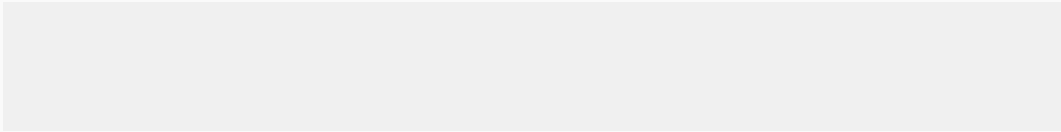
- All of the models (except the final model) can be run in parallel.

Let's look at the Ames K-NN model and evaluate $k = 1, 2, \dots, 20$ using the same 10-fold cross-validation as before.

dials

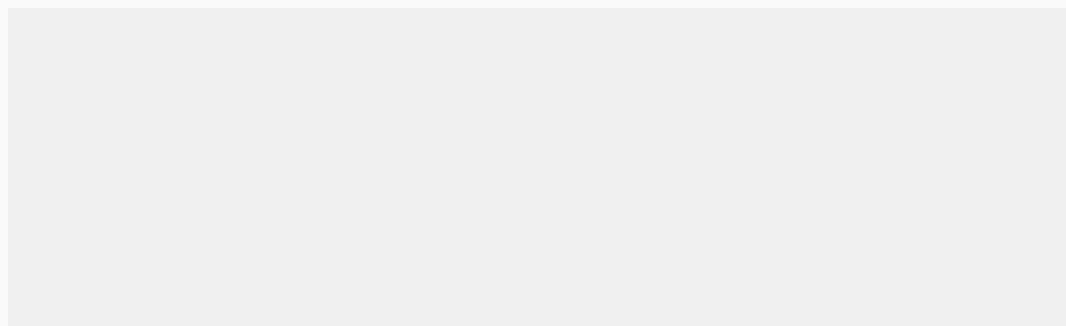
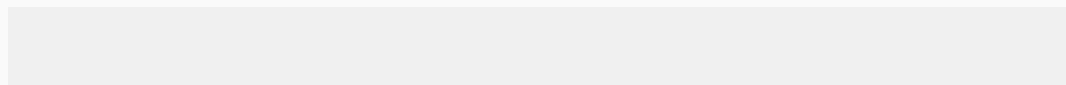
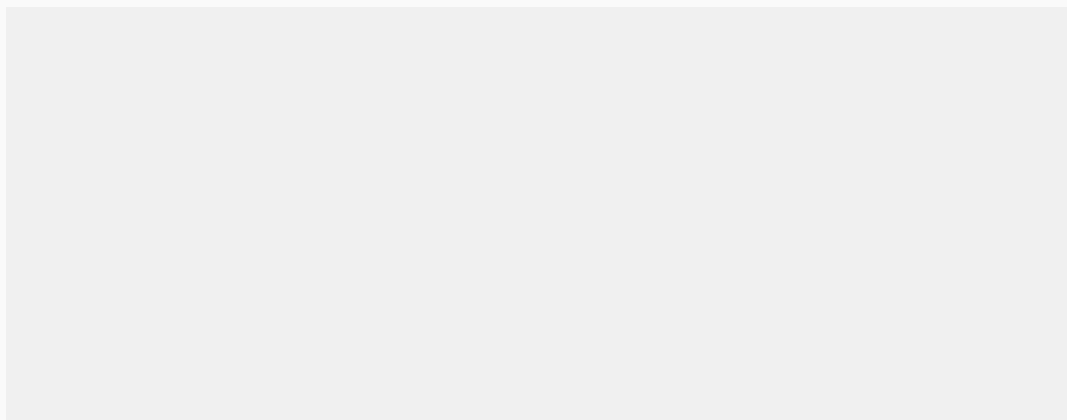
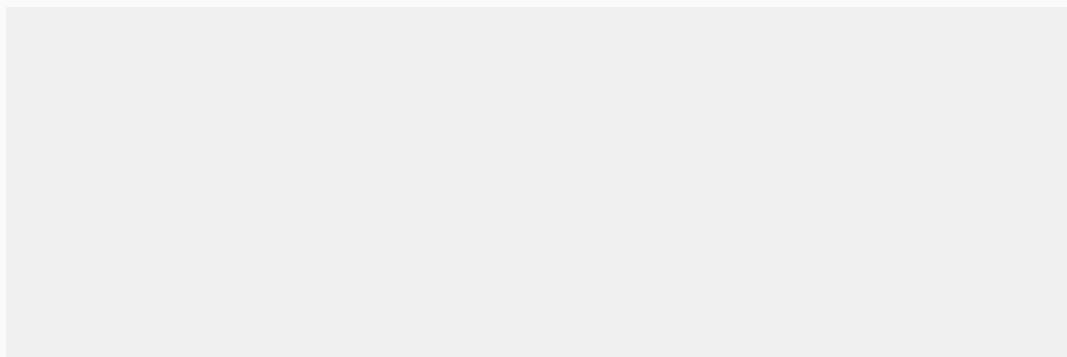
To tune `brglm2` models, two concepts are important to understand:

- 1) `parameters` are those that you want to tune over.
- 2) `grid` is a package for filling in those varying parameters with a tuning grid.



dials

can be used join a parameter grid to a specification. This fills in the field of the specification with the 20 different values.



Now that we have the specification grid, we will start coding this algorithm from the inside out...

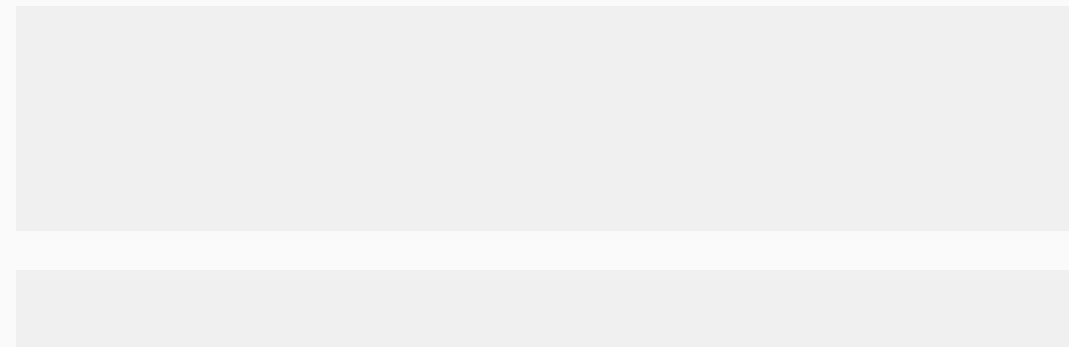
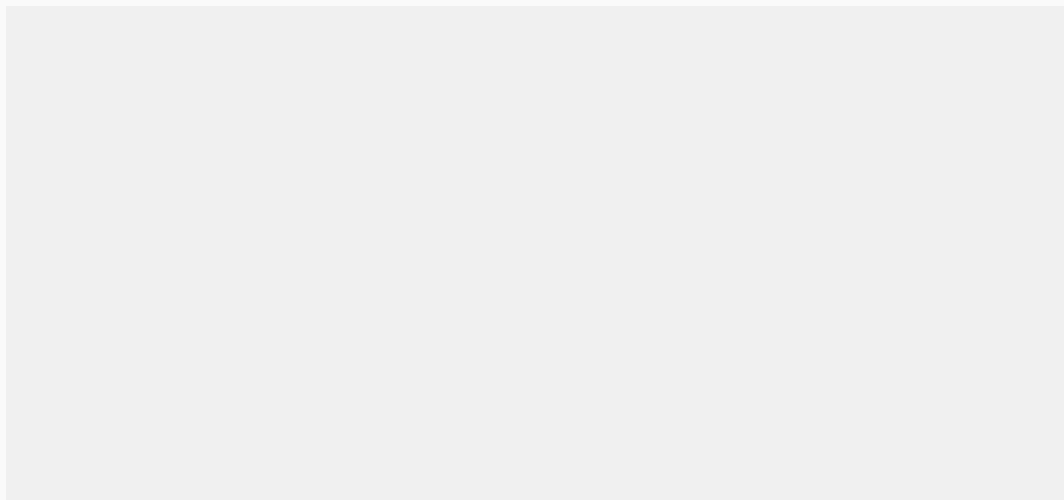
A One Split + One Spec Combo



These steps are:

- └─ Fit the model using the analysis set
- └─ Compute the performance on the assessment set and save

In the code below, `fit_resamples` will be one of the elements of `fit_resamples` and `perf` is one of `perf`. We can reuse many of the functions we've already created.

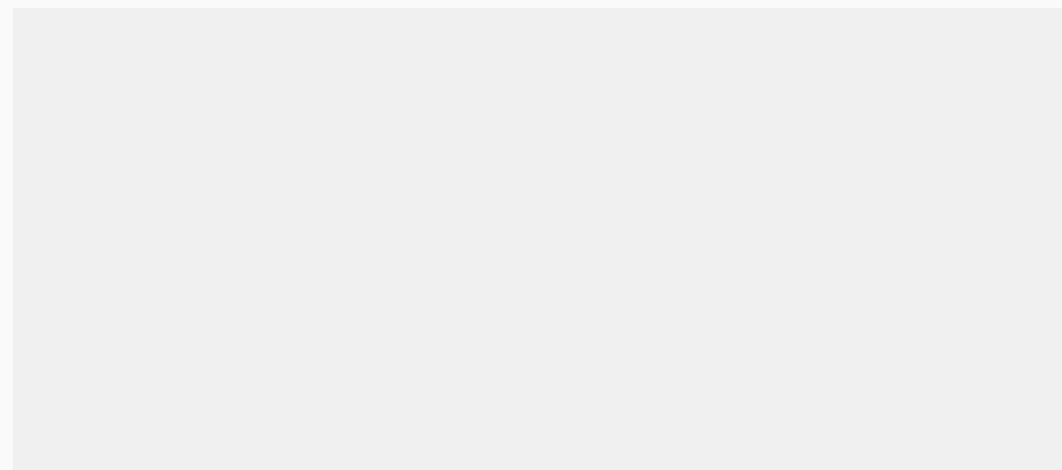
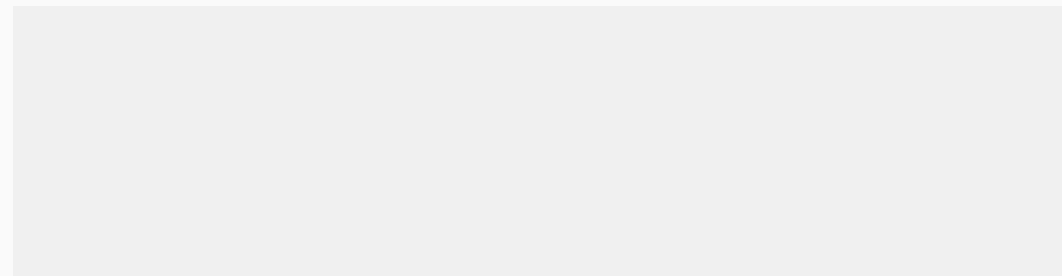
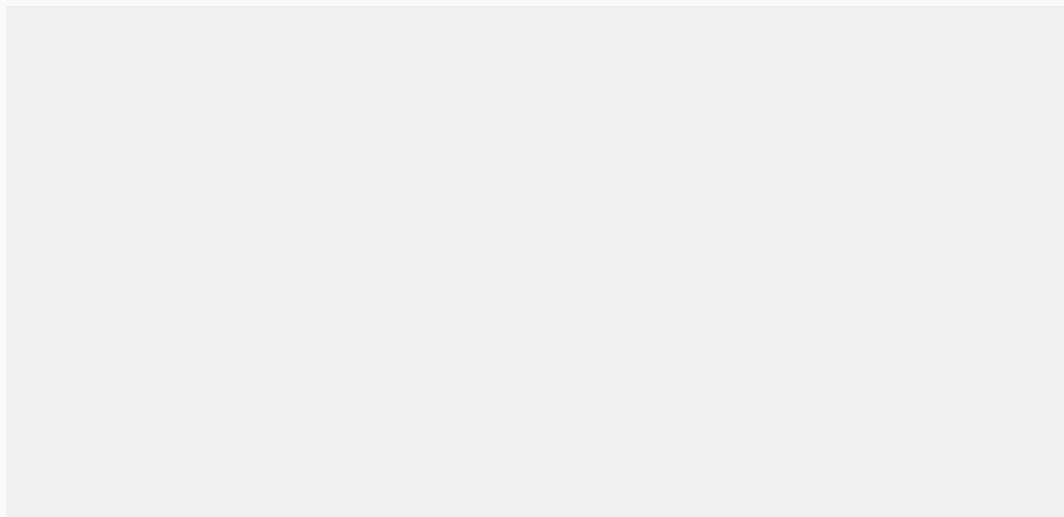


One Split + All Specs



Now we apply `fit_one_spec_one_split()` to every parameter combination.

```
| └─ For each tuning parameter value  
|   └─ Run `fit_one_spec_one_split()`
```

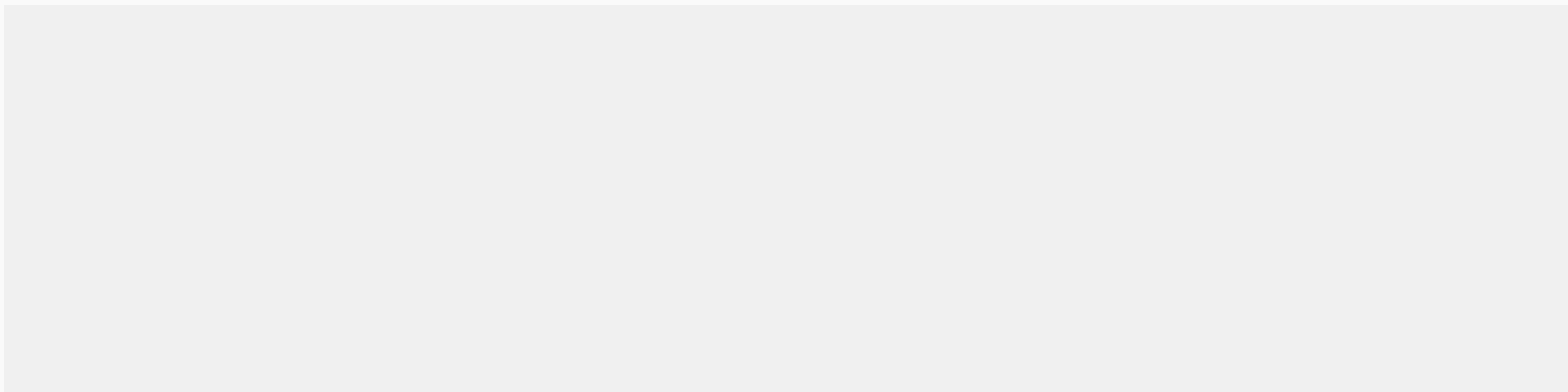


All Splits + All Specs



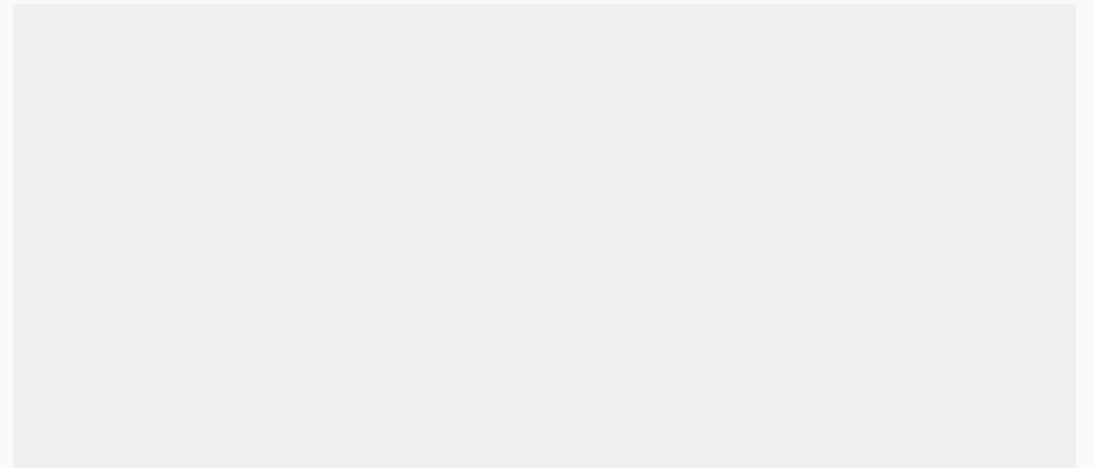
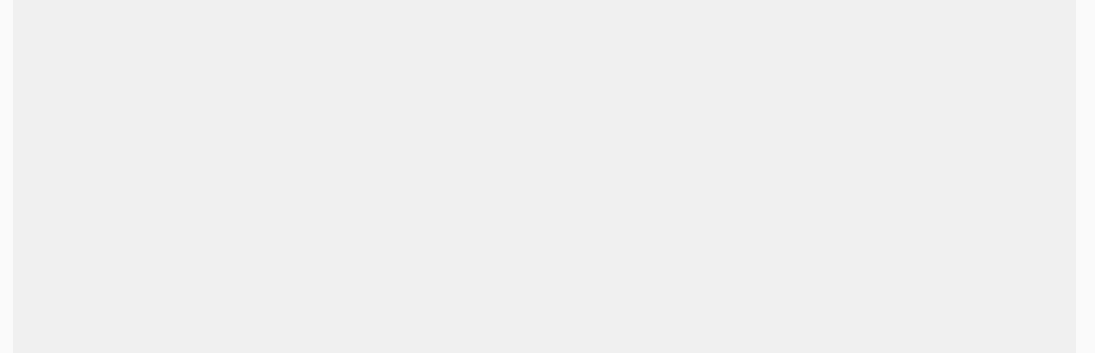
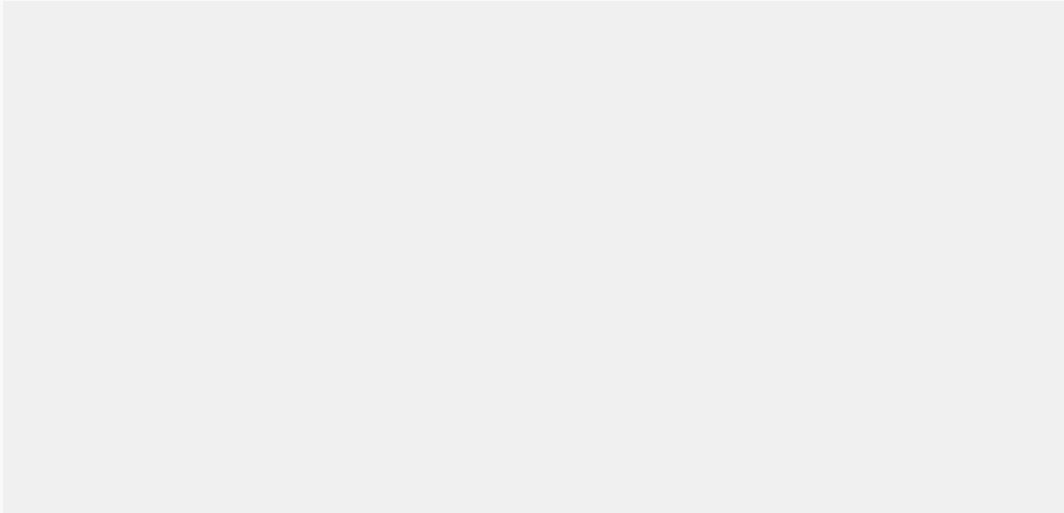
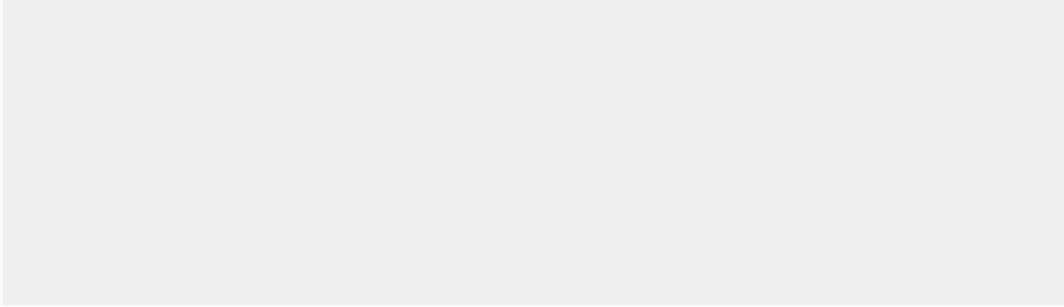
```
└─ For each resample  
  └─ Run `fit_all_specs_one_split()`
```

Here, `resample` is the resample object and `fit_resample` is the function.



This outputs a tibble with columns for the resample, the resample id (e.g., `resample_id`), and a list column of the performance of each tuning parameter combination for that resample.

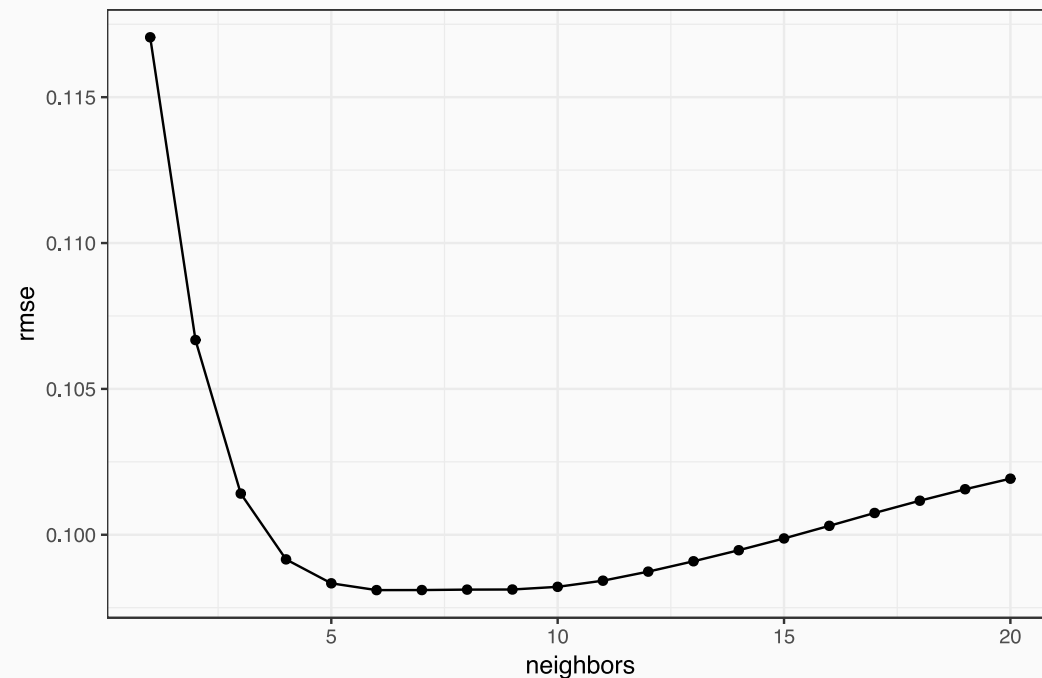
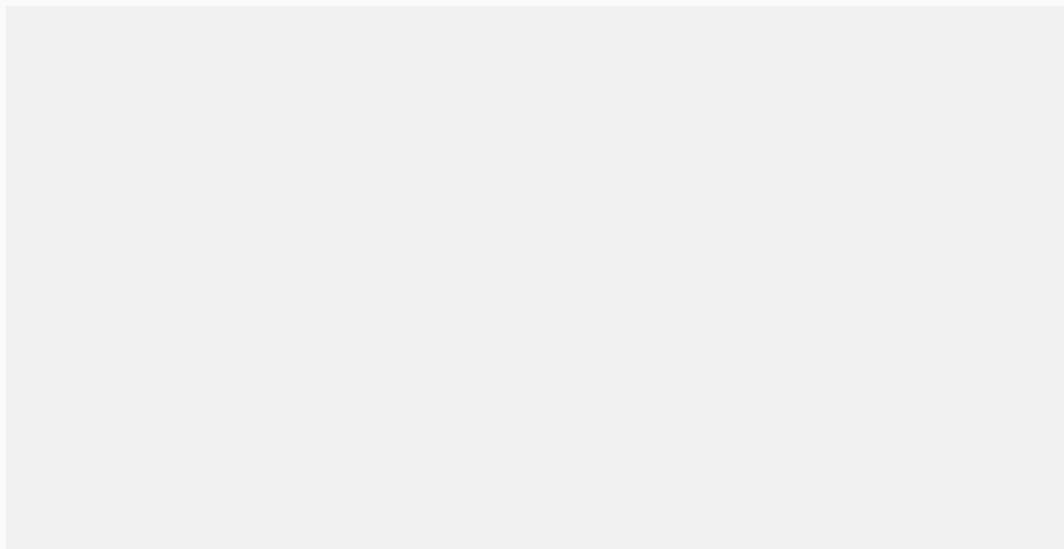
Running the Code



The Performance Profile

To summarize the results for each value of

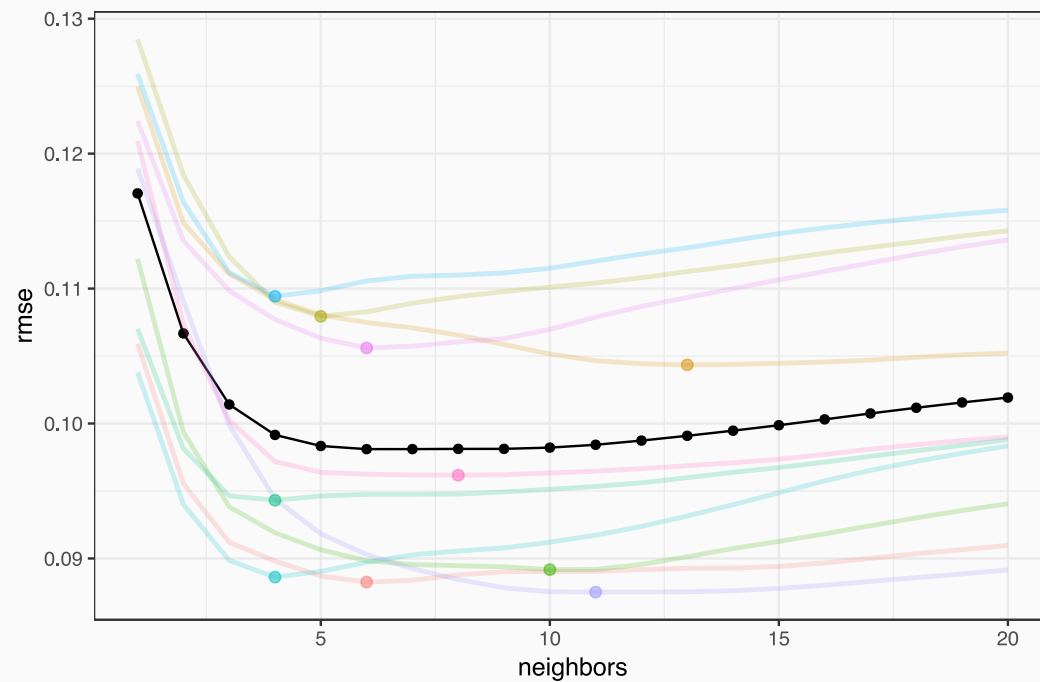
:



Although it is numerically optimal, we are not required to use a value of 6 neighbors for the final model.

Resampling Variation

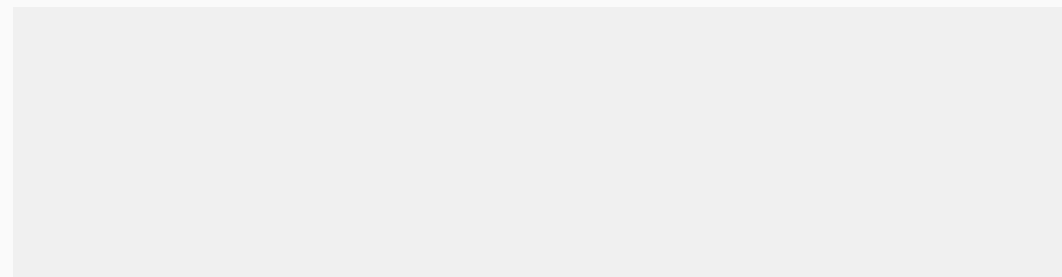
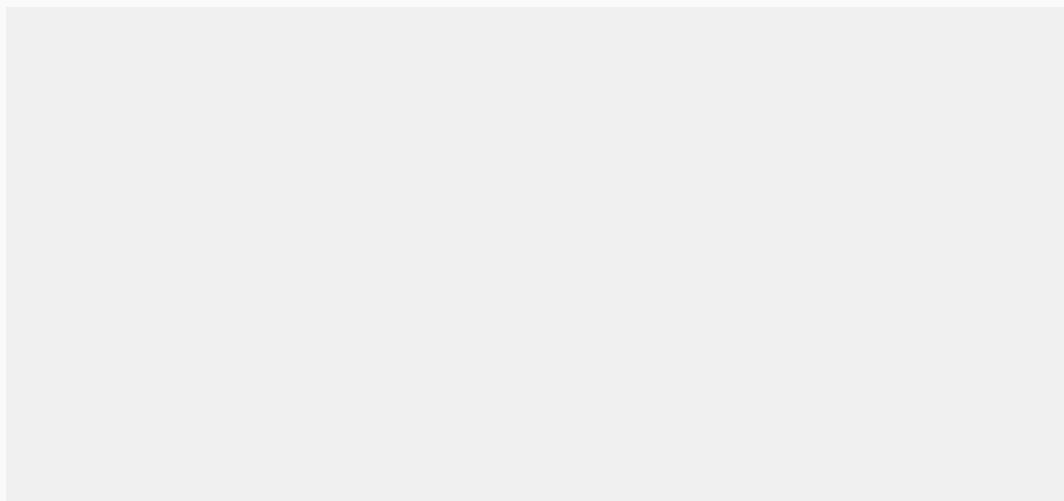
How stable is this? We can also plot the individual curves and their minimums.



Next Steps



At this point, we would decide on a good value for `best_model`, fit that model, and use it going forward:



To reiterate: the previous 10 models created during the grid search are not used once `best_model` is set.

Later, we will look at a high-level API in `tidymodels` that streamlines almost all of this process for many different models. A similar process is being created for the modular tidy packages you see here. We'll talk about this later.