

# Recipes for Data Processing

Max Kuhn (RStudio)

# recipes and Other Modeling Packages

`recipes` is part of a growing collection of R packages that use the tenets of the `tidyverse`:

1. Reuse existing data structures.
2. Compose simple functions with the pipe.
3. Embrace functional programming.
4. Design for humans.

This approach is exemplified by packages such as: `broom`, `infer`, `rsample`, `recipes`, `tidyposterior`, `tidytext`, and `yardstick`.

The list of packages in the `tidymodels` org continues to grow.

The more *traditional approach* uses high-level syntax and is perhaps the most untidy code that you will encounter.

`caret` is the primary package for *highly untidy* predictive modeling:

1. More traditional R coding style.
2. High-level "I'll do that for you" syntax.
3. More comprehensive (for now) and less modular.
4. Contains many optimizations and is easily parallelized.

Luckily, `recipes` can work with both approaches.

# Example Data Set - House Prices

For regression problems, we will use the Ames IA housing data. There are 2,930 properties in the data.

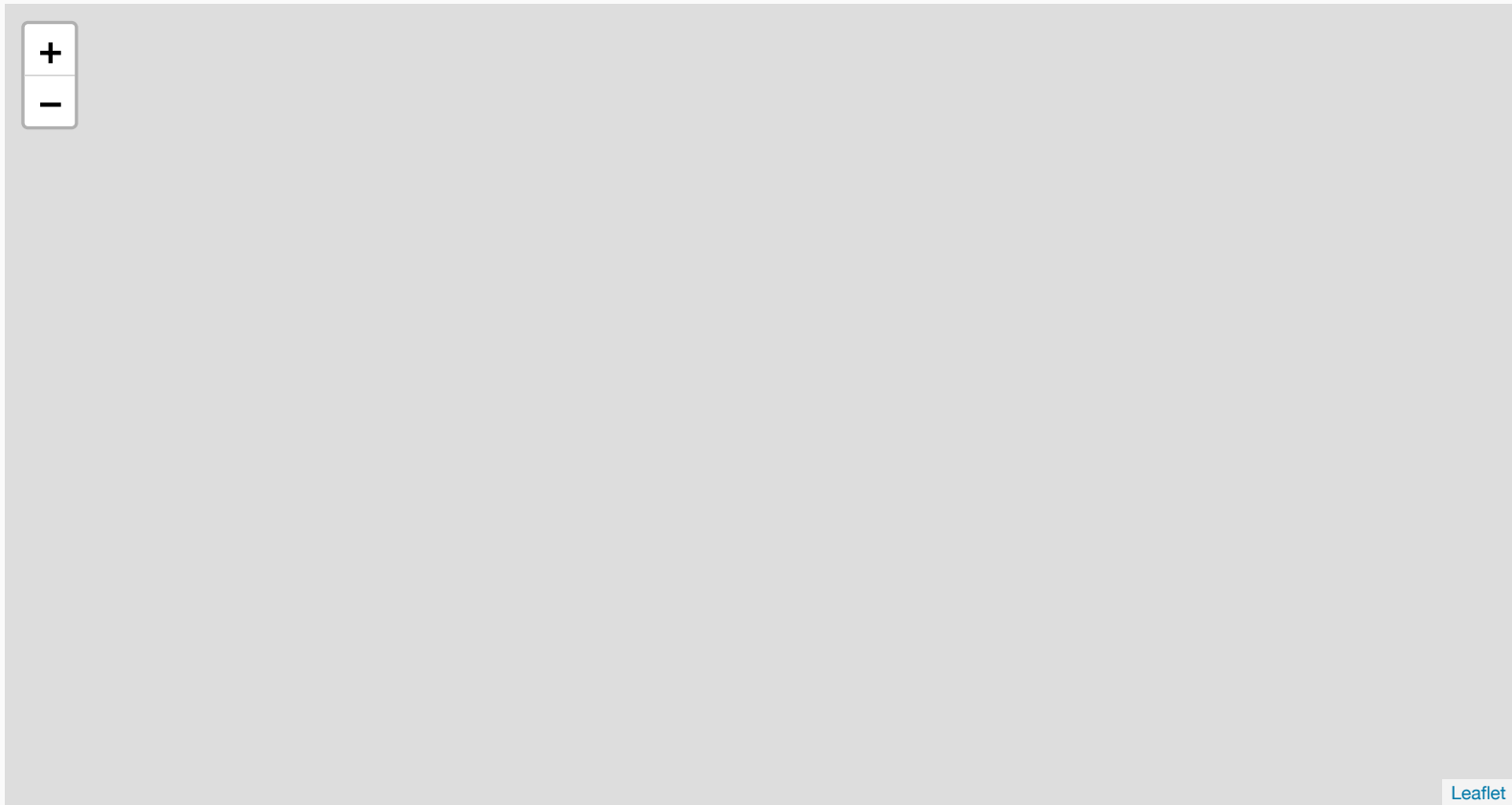
The sale price was recorded along with 81 predictors, including:

- Location (e.g. neighborhood) and lot information.
- House components (garage, fireplace, pool, porch, etc.).
- General assessments such as overall quality and condition.
- Number of bedrooms, baths, and so on.

More details can be found in [De Cock \(2011, Journal of Statistics Education\)](#).

The raw data are at <http://bit.ly/2whgsQM> but we will use a processed version found in the [AmesHousing](#) package.

# Example Data Set - House Prices



# Tidyverse Syntax



Many tidyverse functions have syntax unlike base R code. For example:

- vectors of variable names are eschewed in favor of *functional programming*. For example:

```
contains("Sepal")  
  
# instead of  
  
c("Sepal.Width", "Sepal.Length")
```

- The *pipe* operator is preferred. For example

```
merged <- inner_join(a, b)  
  
# is equal to  
  
merged <- a %>%  
  inner_join(b)
```

- Functions are more *modular* than their traditional analogs (dplyr's `filter` and `select` vs `base::subset`)

# Some Example Data Manipulation Code



```
library(tidyverse)

ames <-
  read_delim("http://bit.ly/2whgsQM", delim = "\t") %>%
  rename_at(vars(contains(' ')), funs(gsub(' ', '_', .))) %>%
  rename(Sale_Price = SalePrice) %>%
  filter(!is.na(Electrical)) %>%
  select(-Order, -PID, -Garage_Yr_Blt)

ames %>%
  group_by(Alley) %>%
  summarize(mean_price = mean(Sale_Price/1000),
            n = sum(!is.na(Sale_Price)))
```

```
## # A tibble: 3 x 3
##   Alley mean_price     n
##   <chr>      <dbl> <int>
## 1 Grvl         124.   120
## 2 Pave         177.    78
## 3 <NA>         183.  2731
```

# Example ggplot2 Code



```
library(ggplot2)

ggplot(ames,
       aes(x = Garage_Type,
           y = Sale_Price)) +
  geom_violin() +
  coord_trans(y = "log10") +
  xlab("Garage Type") +
  ylab("Sale Price")
```



# Examples of `purrr::map*`



```
library(purrr)
```

```
# summarize via purrr::map
by_alley <- split(ames, ames$Alley)
is_list(by_alley)
```

```
## [1] TRUE
```

```
map(by_alley, nrow)
```

```
## $Grv1
## [1] 120
##
## $Pave
## [1] 78
```

```
# or better yet:
map_int(by_alley, nrow)
```

```
## Grv1 Pave
## 120 78
```

```
# works on non-list vectors too
ames %>%
  mutate(Sale_Price = Sale_Price %>%
    map_dbl(function(x) x / 1000)) %>%
  select(Sale_Price, Yr_Sold) %>%
  head(4)
```

```
## # A tibble: 4 x 2
##   Sale_Price Yr_Sold
##       <dbl>   <int>
## 1       215     2010
## 2       105     2010
## 3       172     2010
## 4       244     2010
```



# Reasons for Modifying the Data

- Some models ( $K$ -NN, SVMs, PLS, neural networks) require that the predictor variables have the same units. **Centering** and **scaling** the predictors can be used for this purpose.
- Other models are very sensitive to correlations between the predictors and **filters** or **PCA signal extraction** can improve the model.
- As we'll see in an example, changing the scale of the predictors using a **transformation** can lead to a big improvement.
- In other cases, the data can be **encoded** in a way that maximizes its effect on the model. Representing the date as the day or the week can be very effective for modeling public transportation data.
- Many models cannot cope with missing data so **imputation** strategies might be necessary.
- Development of new *features* that represent something important to the outcome (e.g. compute distances to public transportation, university buildings, public schools, etc.)

# A Bivariate Example

The plot on the right shows two predictors from a real *test* set where the object is to predict the two classes.

The predictors are strongly correlated and each has a right-skewed distribution.

There appears to be some class separation but only in the bivariate plot; the individual predictors show poor discrimination of the classes.

Some models might be sensitive to highly correlated and/or skewed predictors.

Is there something that we can do to make the predictors *easier for the model to use*?

**Any ideas?**



# A Bivariate Example

We might start by estimating transformations of the predictors to resolve the skewness.

The Box-Cox transformation is a family of transformations originally designed for the outcomes of models. We can use it here for the predictors.

It uses the data to estimate a wide variety of transformations including the inverse, log, sqrt, and polynomial functions.

Using each factor in isolation, both predictors were determined to need inverse transformations (approximately).

The figure on the right shows the data after these transformations have been applied.

A logistic regression model shows a substantial improvement in classifying using the altered data.



# Preprocessing Categorical Predictors

(Chapter **Five**)

# Ames Data

The `AmesHousing` package has a better version of the data that has been cleaned up and includes geocodes for the properties. Let's load that and use the `rsample` package to make a split of the data where 3/4 goes into training set:

```
library(AmesHousing)
ames <- make_ames()
nrow(ames)
```

```
## [1] 2930
```

```
library(rsample)

# Make sure that you get the same random numbers
set.seed(4595)
data_split <- initial_split(ames, strata = "Sale_Price")

ames_train <- training(data_split)
ames_test  <- testing(data_split)
```

# Dummy Variables

One common procedure for modeling is to create numeric representations of categorical data. This is usually done via *dummy variables*: a set of binary 0/1 variables for different levels of an R factor.

For example, the Ames housing data contains a predictor called `Alley` with levels: 'Gravel', 'No\_Alley\_Access', 'Paved'.

Most dummy variable procedures would make two numeric variables from this predictor that are zero for a level and one otherwise:

Data	Dummy Variables	
	No_Alley_Access	Paved
Gravel	0	0
No_Alley_Access	1	0
Paved	0	1

# Dummy Variables

If there are  $C$  levels of the factor, only  $C-1$  dummy variables are created since the last can be inferred from the others. There are different contrast schemes for creating the new variables.

For ordered factors, *polynomial contrasts* are used.

How do you create them in R?

The formula method does this for you<sup>1</sup>. Otherwise, the traditional method is to use the `model.matrix` function to create a matrix. However, there are some caveats to this that can make things difficult.

We'll show another method for making them shortly.

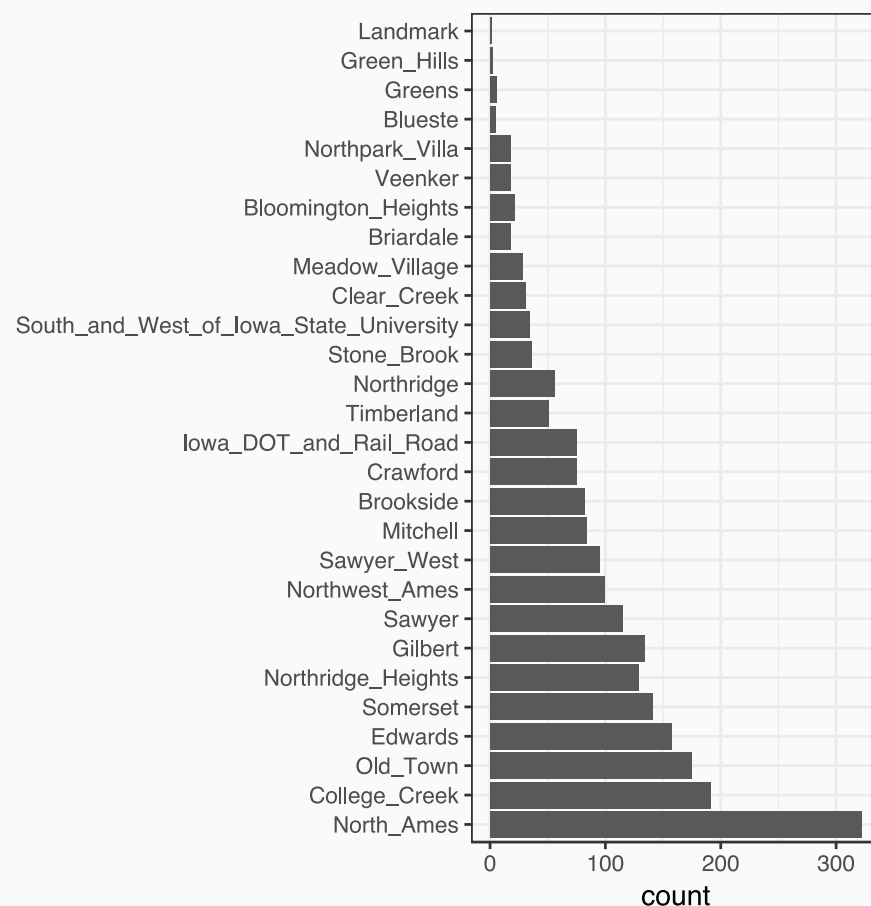
[1] *Almost always* at least. Tree- and rule-based model functions do not. Examples are `randomforest::randomForest`, `ranger::ranger`, `rpart::rpart`, `C50::C5.0`, `Cubist::cubist`, `klaR::NaiveBayes` and others.

# Infrequent Levels in Categorical Factors

One issue is: what happens when there are very few values of a level?

Consider the Ames training set and the `Neighborhood` variable.

If these data are resampled, what would happen to Landmark and similar locations when dummy variables are created?





# Infrequent Levels in Categorical Factors

A *zero-variance* predictor that has only a single value (zero) would be the result.

For many models (e.g. linear/logistic regression, etc.) would find this numerically problematic and issue a warning and `NA` values for that coefficient. Trees and similar models would not notice.

There are two main approaches to dealing with this:

- Run a filter on the training set predictors prior to running the model and remove the zero-variance predictors.
- Recode the factor so that infrequently occurring predictors (and possibly new values) are pooled into an "other" category.

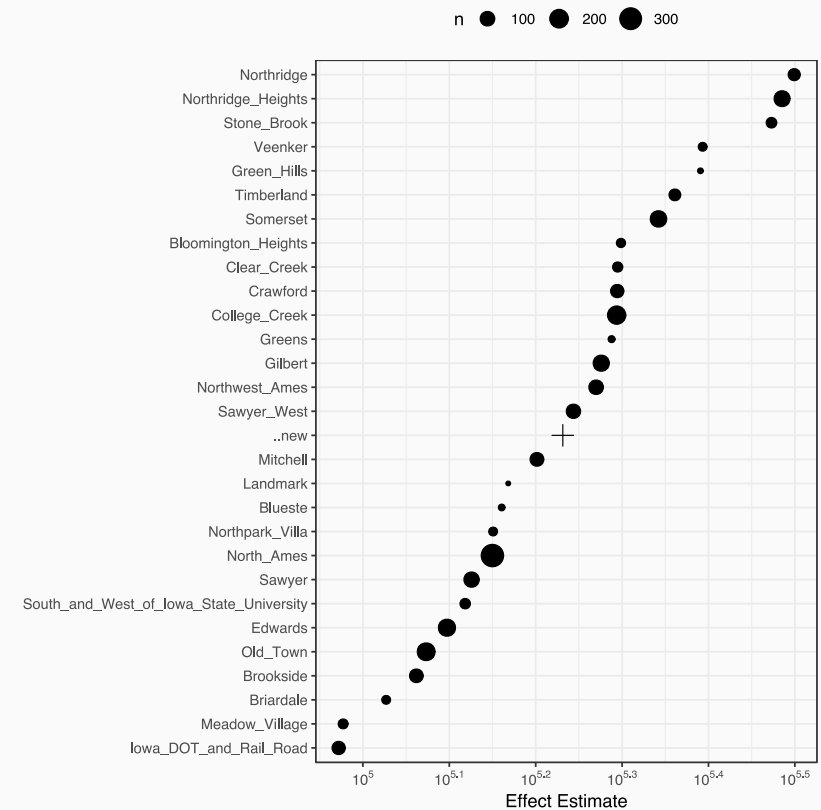
However, `model.matrix` and the formula method are incapable of doing either of these.

# Other Approaches

A few other approaches that might help here:

- *effect* or *likelihood encodings* of categorical predictors estimate the mean effect of the outcome for every factor level in the predictor. These estimates are used in place of the factor levels. Shrinkage/regularization can improve these approaches.
- *entity embeddings* use a neural network to create features that capture the relationships between the categories and the categories and the outcome.

An add-on to `recipes` called `embed` can be used for these encodings.



# Creating Recipes

# Recipes



Recipes are an alternative method for creating the data frame of predictors for a model.

They allow for a sequence of *steps* that define how data should be handled.

Suppose we use the formula `log10(Sale_Price) ~ Longitude + Latitude`? These steps are:

- Assign `Sale_Price` to be the outcome
- Assign `Longitude` and `Latitude` are predictors
- Log transform the outcome

To start using a recipe, the these steps can be done using

```
library(recipes)
mod_rec <- recipe(Sale_Price ~ Longitude + Latitude, data = ames_train) %>%
  step_log(Sale_Price, base = 10)
```

This creates the recipe for data processing (but does not execute it yet)

# Recipes and Categorical Predictors



To deal with the dummy variable issue, we can expand the recipe with more steps:

```
mod_rec <- recipe(Sale_Price ~ Longitude + Latitude + Neighborhood, data = ames_train) %>%  
  step_log(Sale_Price, base = 10) %>%  
  
  # Lump factor levels that occur in <= 5% of data as "other"  
  step_other(Neighborhood, threshold = 0.05) %>%  
  
  # Create dummy variables for _any_ factor variables  
  step_dummy(all_nominal())
```

Note that we can use standard `dplyr` selectors as well as some new ones based on the data type (`all_nominal()`) or by their role in the analysis (`all_predictors()`).

The general process for using recipes is

```
recipe    -->  prepare    -->  bake/juice  
(define) --> (estimate) -->  (apply)
```

# Preparing the Recipe



Now that we have a preprocessing *specification*, let's run it on the training set to *prepare* the recipe:

```
mod_rec_trained <- prep(mod_rec, training = ames_train, retain = TRUE, verbose = TRUE)
```

```
## oper 1 step log [training]  
## oper 2 step other [training]  
## oper 3 step dummy [training]
```

Here, the "training" is to determine which factors to pool and to enumerate the factor levels of the `Neighborhood` variable,

`retain` keeps the processed version of the training set around so we don't have to recompute it.

# Preparing the Recipe

```
mod_rec_trained
```

```
## Data Recipe
##
## Inputs:
##
##      role #variables
## outcome      1
## predictor      3
##
## Training data contained 2199 data points and no missing data.
##
## Operations:
##
## Log transformation on Sale_Price [trained]
## Collapsing factor levels for Neighborhood [trained]
## Dummy variables from Neighborhood [trained]
```

# Getting the Values



Once the recipe is prepared, it can be applied to any data set using `bake`:

```
ames_test_dummies <- bake(mod_rec_trained, newdata = ames_test)
names(ames_test_dummies)
```

```
## [1] "Sale_Price"           "Longitude"           "Latitude"
## [4] "Neighborhood_College_Creek" "Neighborhood_Old_Town" "Neighborhood_Edwards"
## [7] "Neighborhood_Somerset"   "Neighborhood_Northridge_Heights" "Neighborhood_Gilbert"
## [10] "Neighborhood_Sawyer"    "Neighborhood_other"
```

If `retain = TRUE` the training set does not need to be "rebaked". The `juice` function can return the processed version of the training data.

Selectors can be used with `bake` and the default is `everything()`.



# Interaction Effects

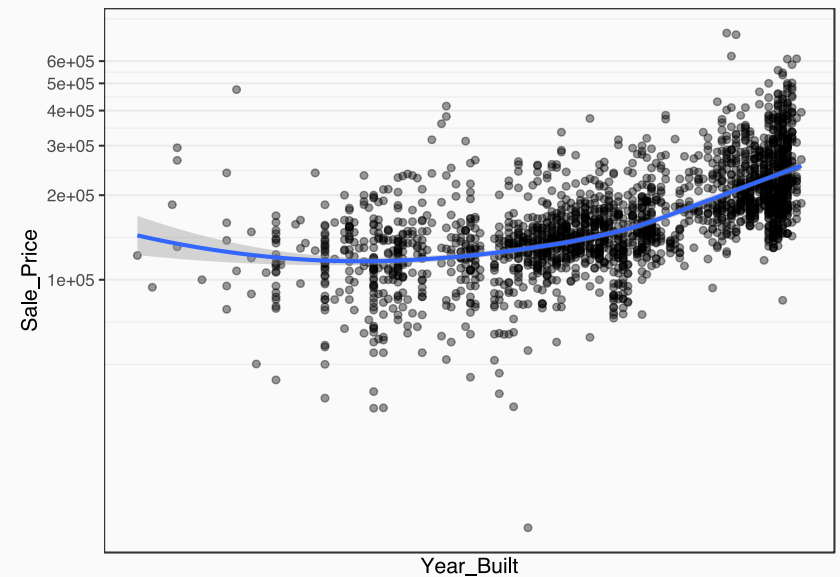
# Interactions

An interaction between two predictors indicates that the relationship between the predictors and the outcome cannot be described using only one of the variables.

For example, let's look at the relationship between the price of a house and the year in which it was built. The relationship appears to be slightly nonlinear, possibly quadratic:

```
price_breaks <- (1:6)*(10^5)

ggplot(ames_train,
       aes(x = Year_Built, y = Sale_Price)) +
  geom_point(alpha = 0.4) +
  scale_x_log10() +
  scale_y_continuous(breaks = price_breaks,
                    trans = "log10") +
  geom_smooth(method = "loess")
```

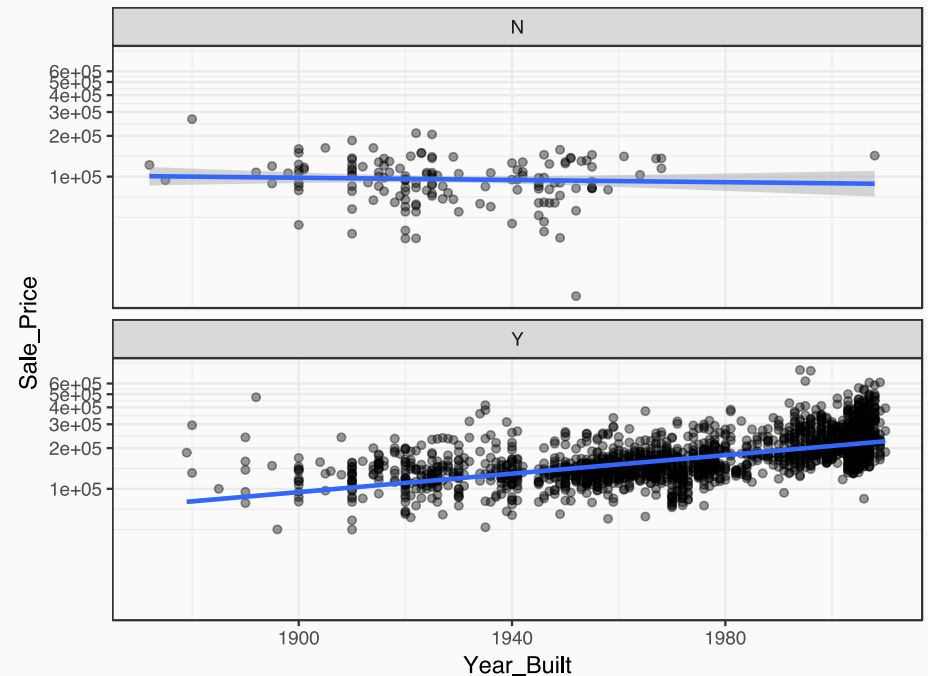


# Interactions

However... what if we separate this trend based on whether the property has air conditioning (93.4% of the training set) or not (6.6%):

```
library(MASS) # to get robust linear regression model

ggplot(ames_train,
      aes(x = Year_Built,
          y = Sale_Price)) +
  geom_point(alpha = 0.4) +
  scale_y_continuous(breaks = price_breaks,
                    trans = "log10") +
  facet_wrap(~ Central_Air, nrow = 2) +
  geom_smooth(method = "rlm")
```



# Interactions

It appears as though the relationship between the year built and the sale price is somewhat *different* for the two groups.

- When there is no AC, the trend is perhaps flat or slightly decreasing
- With AC, there is a linear increasing trend or is perhaps slightly quadratic with some outliers at the low end.

```
mod1 <- lm(log10(Sale_Price) ~ Year_Built + Central_Air, data = ames_train)
mod2 <- lm(log10(Sale_Price) ~ Year_Built + Central_Air + Year_Built: Central_Air, data = ames_train)
anova(mod1, mod2)
```

```
## Analysis of Variance Table
##
## Model 1: log10(Sale_Price) ~ Year_Built + Central_Air
## Model 2: log10(Sale_Price) ~ Year_Built + Central_Air + Year_Built:Central_Air
##   Res.Df  RSS Df Sum of Sq    F Pr(>F)
## 1     2196 41.3
## 2     2195 40.2   1      1.11 60.6 1.1e-14 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

# Interactions in Recipes

We first create the dummy variables for the qualitative predictor (`central_Air`) then use a formula to create the interaction using the `:` operator in an additional step:

```
recipe(Sale_Price ~ Year_Built + Central_Air, data = ames_train) %>%  
  step_log(Sale_Price) %>%  
  step_dummy(Central_Air) %>%  
  step_interact(~ starts_with("Central_Air"):Year_Built) %>%  
  prep(training = ames_train, retain = TRUE) %>%  
  juice %>%  
  # select a few rows with different values  
  slice(153:157)
```

```
## # A tibble: 5 x 4  
##   Year_Built Sale_Price Central_Air_Y Central_Air_Y_x_Year_Built  
##       <int>     <dbl>         <dbl>                <dbl>  
## 1      1912       12.0           1                1912  
## 2      1930       10.9           0                  0  
## 3      1900       11.8           1                1900  
## 4      1959       12.1           1                1959  
## 5      1917       11.6           0                  0
```

# Adding Recipes to our `rsample` Workflows

# Resampling Methods

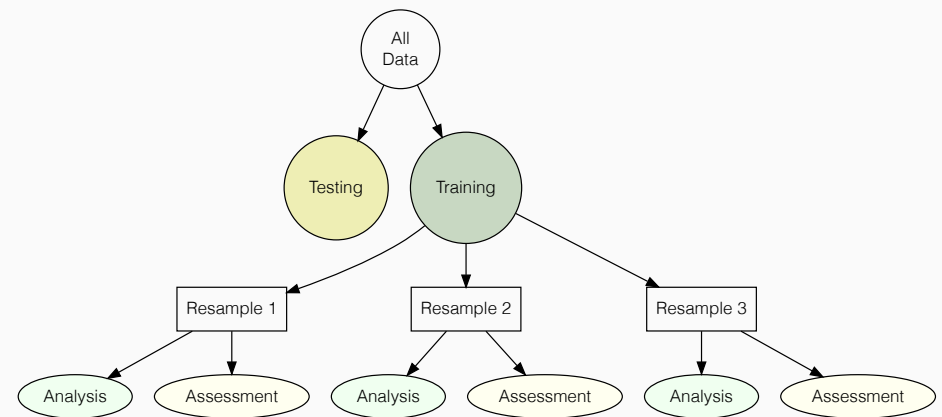
These are **additional data splitting schemes** that are applied to the *training* set.

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

- The *analysis set* is used to fit the model (analogous to the training set).
- Performance is determined using the *assessment set*.

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  $V$  distinct blocks of roughly equal size.

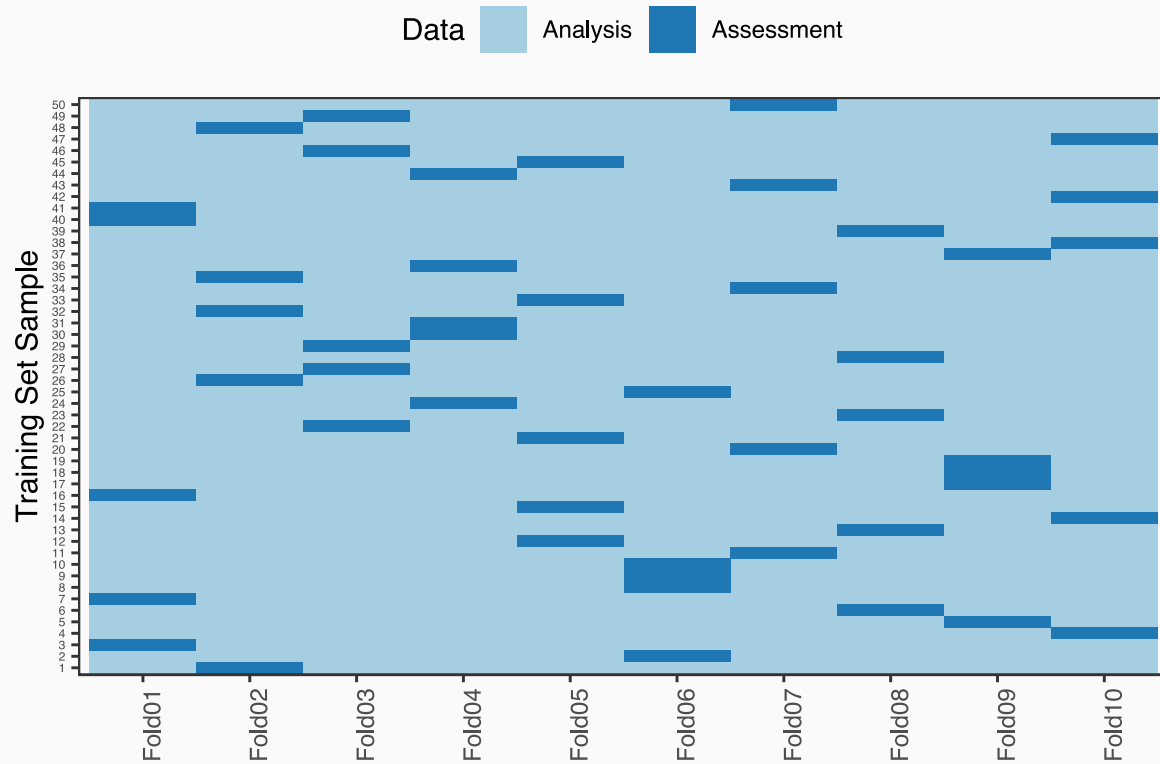
- We leave out the first block of analysis data and fit a model.
- This model is used to predict the held-out block of assessment data.
- We continue this process until we've predicted all  $V$  assessment blocks

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

$V$  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 $n = 50$



# Resampling and Preprocessing

It is important to realize that almost all preprocessing steps that involve estimation should be bundled inside of the resampling process so that the performance estimates are not biased.

- **Bad:** preprocess the data, resample the model
- **Good:** resample the preprocessing and modeling

Also:

- Avoid *information leakage* by having all operations in the modeling process occur only on the training set.
- Do not reestimate anything on the test set. For example, to center new data, the training set mean is used.

# recipes and rsample



Let's go back to the Ames housing data and work on building models with recipes and resampling.

We'll do a 3/4 split of the data with the majority going into the training set. For resampling, 10-fold CV would leave out about 219 houses (which should be enough to compute the RMSE).

```
set.seed(2453)
cv_splits <-
  vfold_cv(ames_train, v = 10, strata = "Sale_Price")
cv_splits[1:3,]
```

```
## # A tibble: 3 x 2
##   splits      id
## * <list>     <chr>
## 1 <S3: rsplit> Fold01
## 2 <S3: rsplit> Fold02
## 3 <S3: rsplit> Fold03
```

```
cv_splits$splits[[1]]
```

```
## <1977/222/2199>
```

```
cv_splits$splits[[1]] %>% analysis() %>% nrow()
```

```
## [1] 1977
```

```
cv_splits$splits[[1]] %>% assessment() %>% nrow()
```

```
## [1] 222
```

# Linear Models



Let's preprocess the training set prior to adding the data to a linear regression model.

- Two numeric predictors are very skewed and could use a transformation ( `Lot_Area` and `Gr_Liv_Area` ).
- We'll add neighborhood in as well and a few other house features.
- Visualizations suggest that the coordinates can be helpful but probably require a nonlinear representation. We can add these using *B-splines* with 5 degrees of freedom. To evaluate this, we will create two versions of the recipe to evaluate this hypothesis.

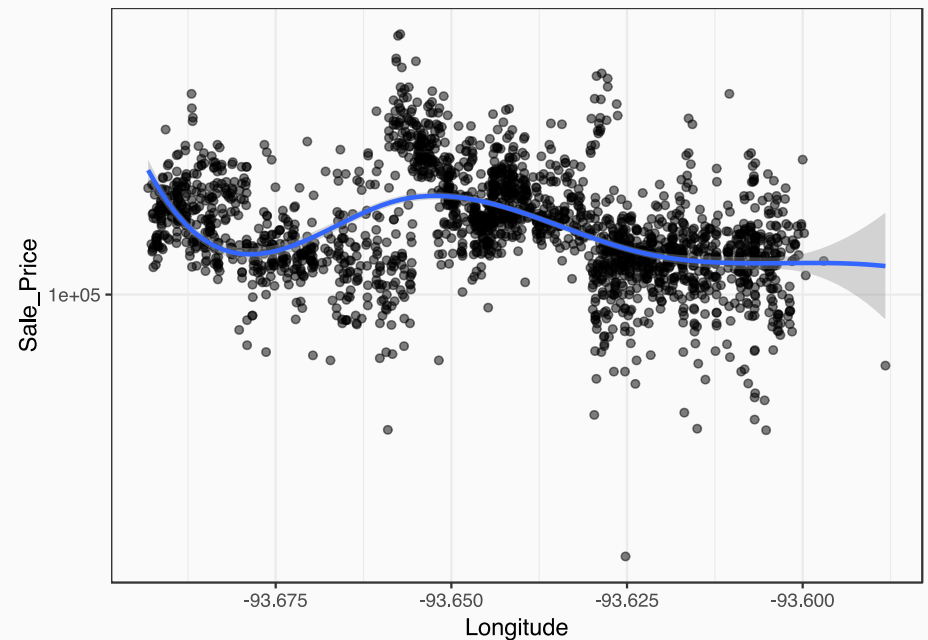
```
ames_rec <- recipe(Sale_Price ~ Bldg_Type + Neighborhood + Year_Built +  
  Gr_Liv_Area + Full_Bath + Year_Sold + Lot_Area +  
  Central_Air + Longitude + Latitude,  
  data = ames_train) %>%  
  step_log(Sale_Price, base = 10) %>%  
  step_YeoJohnson(Lot_Area, Gr_Liv_Area) %>%  
  step_other(Neighborhood, threshold = 0.05) %>%  
  step_dummy(all_nominal()) %>%  
  step_interact(~ starts_with("Central_Air"):Year_Built) %>%  
  step_bs(Longitude, Latitude, options = list(df = 5))
```

# Longitude



```
library(ggplot2)
ggplot(ames_train,
       aes(x = Longitude, y = Sale_Price)) +
  geom_point(alpha = .5) +
  geom_smooth(
    method = "lm",
    formula = y ~ splines::bs(x, 5),
    se = FALSE
  ) +
  scale_y_log10()
```

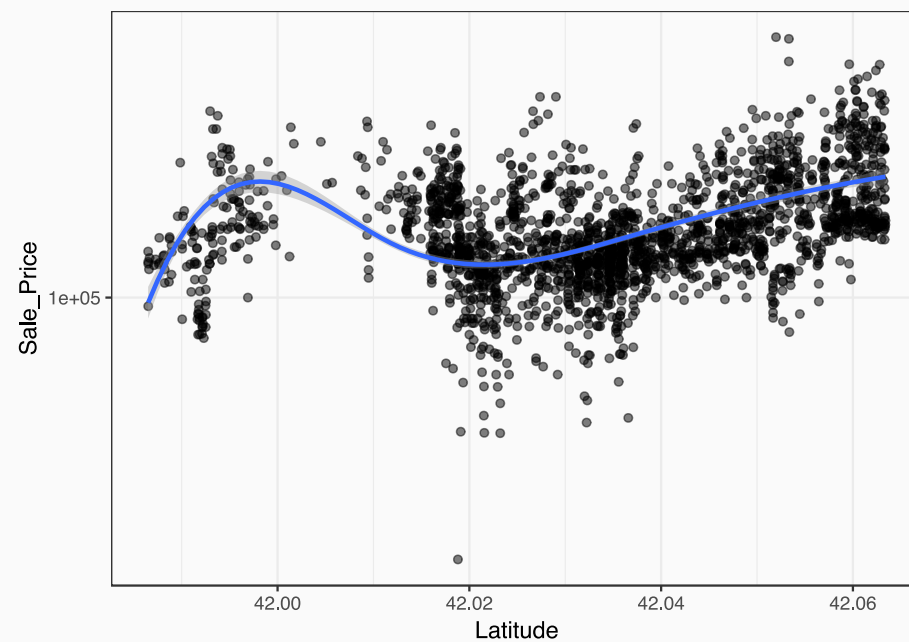
Splines add nonlinear versions of the predictor to a linear model to create smooth and flexible relationships between the predictor and outcome.



# Latitude



```
ggplot(ames_train,  
      aes(x = Latitude, y = Sale_Price)) +  
  geom_point(alpha = .5) +  
  geom_smooth(  
    method = "lm",  
    formula = y ~ splines::bs(x, 5),  
    se = FALSE  
  ) +  
  scale_y_log10()
```



# Preparing the Recipes



Our first step is to run `prep()` on the `ames_rec` recipe but using each of the analysis sets.

`rsample` has a function that is a wrapper around `prep()` that can be used to map over the split objects:

```
prepper
```

```
## function (split_obj, recipe, ...)
## {
##   prep(recipe, training = analysis(split_obj), ...)
## }
## <bytecode: 0x7f80df049348>
## <environment: namespace:rsample>
```

```
library(purrr)
cv_splits <- cv_splits %>%
  mutate(ames_rec = map(splits, prepper, recipe = ames_rec, retain = TRUE))
```

# Fitting the Models



The code below will use the recipe object to get the data. Since each analysis set is used to train the recipe, our previous use of `retain = TRUE` means that the processed version of the data is within the recipe.

This can be returned via the `juice` function.

```
lm_fit_rec <- function(rec_obj, ...)  
  lm(..., data = juice(rec_obj))  
  
cv_splits <- cv_splits %>%  
  mutate(fits = map(ames_rec, lm_fit_rec, Sale_Price ~ .))  
library(broom)  
glance(cv_splits$fits[[1]])
```

```
##   r.squared adj.r.squared  sigma statistic p.value df logLik   AIC    BIC deviance df.residual  
## 1      0.805         0.802 0.0794         276      0 30  2218 -4374 -4200      12.3         1947
```



# Predicting the Assessment Set



This is a little more complex. We need three elements contained in our tibble:

- the split object (to get the assessment data)
- the recipe object (to process the data)
- the linear model (for predictions)

The function is not too bad:

```
assess_predictions <- function(split_obj, rec_obj, mod_obj) {  
  raw_data <- assessment(split_obj)  
  proc_x <- bake(rec_obj, newdata = raw_data, all_predictors())  
  # Now save _all_ of the columns and add predictions.  
  bake(rec_obj, newdata = raw_data, everything()) %>%  
    mutate(  
      .fitted = predict(mod_obj, newdata = proc_x),  
      .resid = Sale_Price - .fitted, # Sale_Price is already logged by the recipe  
      # Save the original row number of the data  
      .row = as.integer(split_obj, data = "assessment")  
    )  
}
```

# Predicting the Assessment Set



Since we have three inputs, we will use `purrr`'s `pmap` function to walk along all three columns in the tibble.

```
cv_splits <- cv_splits %>%  
  mutate(  
    pred =  
      pmap(  
        lst(split_obj = cv_splits$splits, rec_obj = cv_splits$ames_rec, mod_obj = cv_splits$fits),  
        assess_predictions  
      )  
  )
```

We do get some warnings that the assessment data are outside the range of the analysis set values:

```
## Warning in bs(x = c(-93.6235954, -93.636372, -93.627536, -93.65332,  
## -93.649447, : some 'x' values beyond boundary knots may cause ill-  
## conditioned bases
```

# Predicting the Assessment Set



`yardstick::metrics` will compute a small set of summary statistics for metrics model based on the type of outcome (e.g. regression, classification, etc).

```
library(yardstick)

# Compute the summary statistics
map_df(cv_splits$pred, metrics, truth = Sale_Price, estimate = .fitted) %>%
  colMeans
```

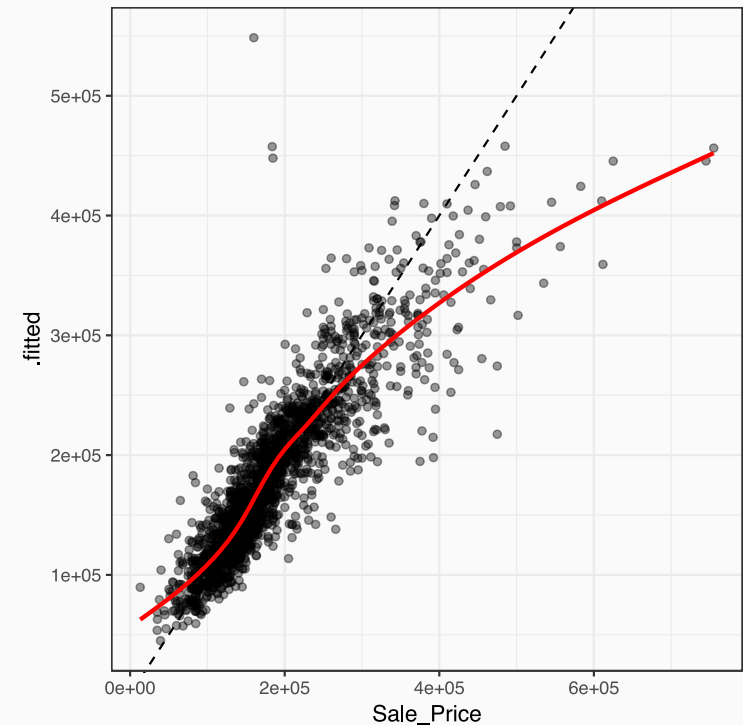
```
##   rmse    rsq   mae
## 0.0798 0.8005 0.0569
```

These results are good but *"the only way to be comfortable with your results is to never look at them"*.

# Graphical Checks for Fit



```
assess_pred <- bind_rows(cv_splits$pred) %>%  
  mutate(Sale_Price = 10^Sale_Price,  
         .fitted = 10^.fitted)  
  
ggplot(assess_pred,  
       aes(x = Sale_Price, y = .fitted)) +  
  geom_abline(lty = 2) +  
  geom_point(alpha = .4) +  
  geom_smooth(se = FALSE, col = "red")
```



# caret and recipes

To facilitate model tuning and other conveniences, `recipes` can be used as inputs into `train` (and, soon, the various feature selection routines). This allows for:

- more extensive feature engineering and preprocessing tools.
- the ability to estimate performance using columns in the data set.

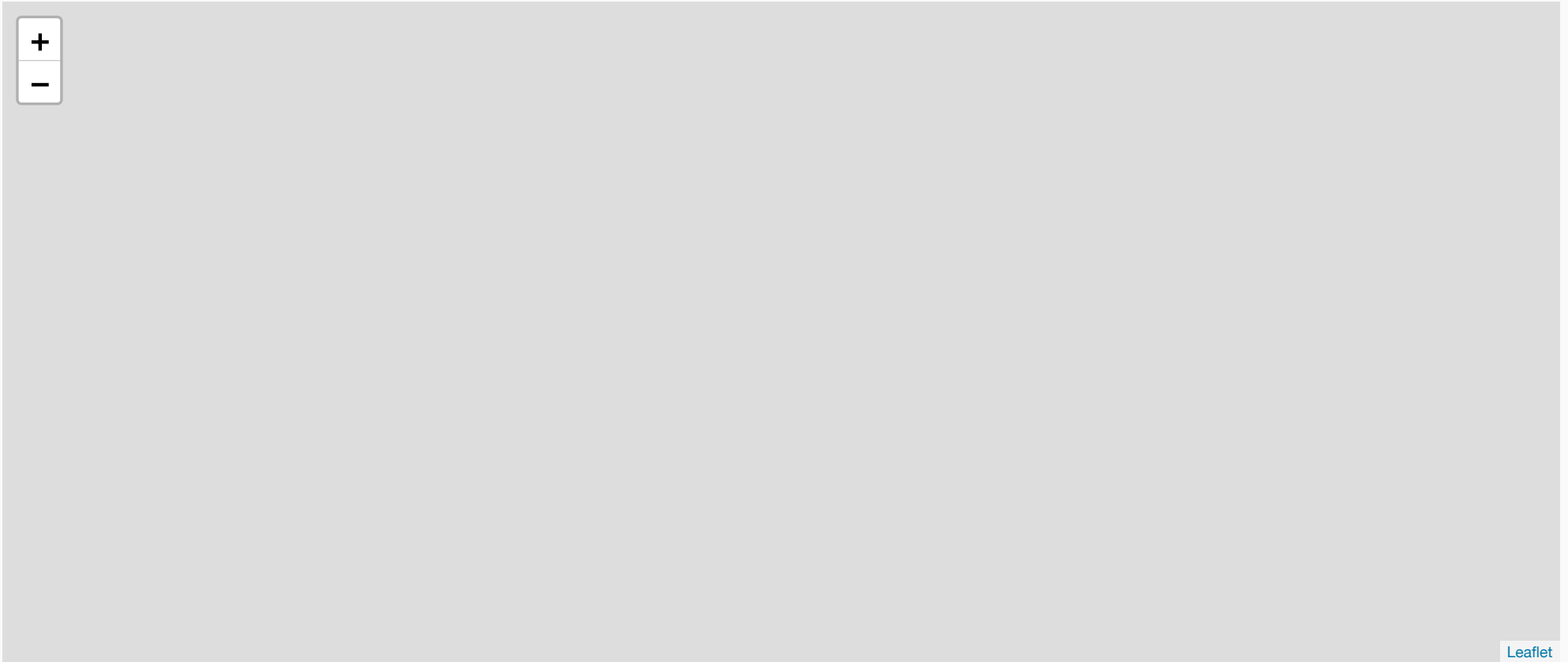
(Unfortunately, the recipe parameters can't be included as part of the tuning process.)

The general syntax is

```
train(recipe, data, method, ...)
```

We'll use a *K*-nearest neighbors model to illustrate.

## 8-Nearest Neighbors (Geographically)



# Recipe for a K-NN Model

Since a *distance* will be computed, we'll need to add two intermediary steps that puts the predictors on the same scale:

```
ames_rec_norm <-  
  recipe(Sale_Price ~ Bldg_Type + Neighborhood + Year_Built +  
    Gr_Liv_Area + Full_Bath + Year_Sold + Lot_Area +  
    Central_Air + Longitude + Latitude,  
    data = ames_train) %>%  
  step_log(Sale_Price, base = 10) %>%  
  step_YeoJohnson(Lot_Area, Gr_Liv_Area) %>%  
  step_other(Neighborhood, threshold = 0.05) %>%  
  step_dummy(all_nominal()) %>%  
  step_center(all_predictors()) %>%  
  step_scale(all_predictors()) %>%  
  step_interact(~ starts_with("Central_Air"):Year_Built) %>%  
  step_bs(Longitude, Latitude, options = list(df = 5))
```

We will tune the model over the number of neighbors as well as how to weight the distances to the neighbors.

# K-NN Model Code

```
converted_resamples <-  
  rsample2caret(cv_splits)  
  
ctrl <- trainControl(method = "cv",  
                     savePredictions = "final")  
ctrl$index <- converted_resamples$index  
ctrl$indexOut <- converted_resamples$indexOut  
  
knn_grid <- expand.grid(  
  kmax = 1:9,  
  distance = 2,  
  kernel = c("rectangular", "triangular", "gaussian")  
)  
  
knn_fit <- train(  
  ames_rec_norm, data = ames_train,  
  method = "kkn",  
  tuneGrid = knn_grid,  
  trControl = ctrl  
)
```

<= Standardize the resamples format to go between `rsample` and `caret`

<= Define the grid of parameters to tune over.

<= Fit the model

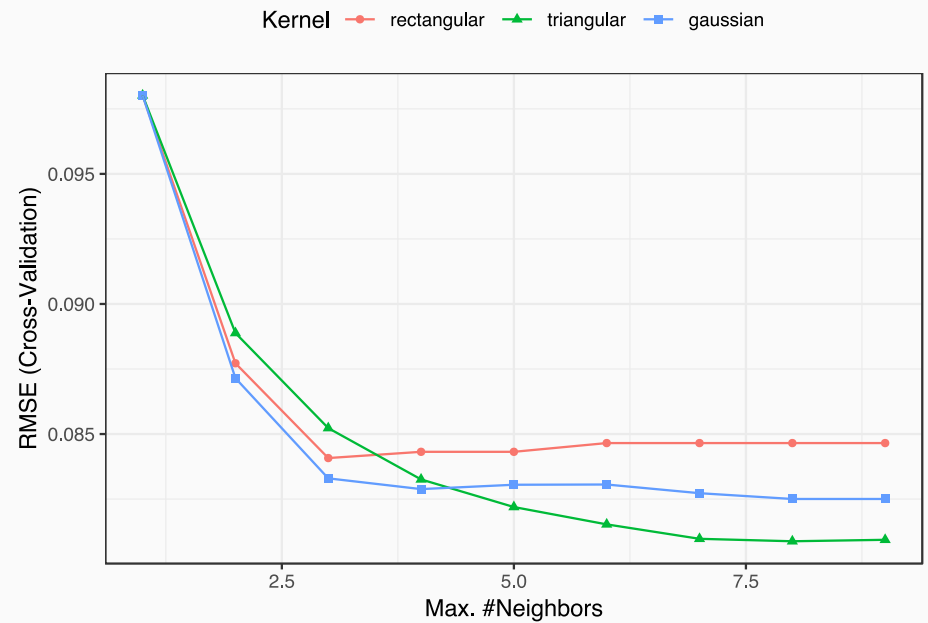


# K-NN Model

```
ggplot(knn_fit)
```

```
getTrainPerf(knn_fit)
```

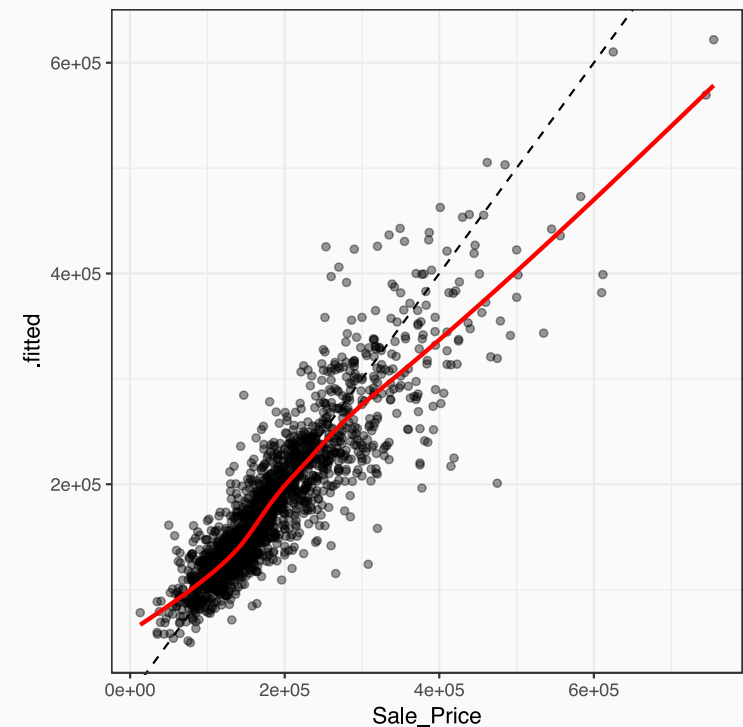
##	TrainRMSE	TrainRsquared	TrainMAE	method
## 1	0.0809	0.794	0.0571	kkn



# Graphical Check for Fit



```
knn_pred <- knn_fit$pred %>%  
  mutate(Sale_Price = 10^obs,  
         .fitted = 10^pred)  
  
ggplot(knn_pred,  
       aes(x = Sale_Price, y = .fitted)) +  
  geom_abline(lty = 2) +  
  geom_point(alpha = .4) +  
  geom_smooth(se = FALSE, col = "red")
```



# About Some of These Houses

From Dmytro Perepolkin via [GitHub/topepo/AmesHousing](#):

I did a little bit of research on [Kaggle](#) regarding those five houses (three "partial sale" houses in Edwards and two upscale in Northridge):

*In fact all three of Edwards properties were sold within three months by the same company. The lots are located next to each other in Cochrane Pkwy Ames, IA. I guess the story was that developer was either in trouble or wanted to boost sales, so they were signing sales contracts on half-finished houses in a new development area at a deep discount. Those few houses are likely to have been built first out of the whole residential block that followed.*

*Two of the upscale houses in Northridge (sold at 745 and 755 thousand, respectively) are located next to each other. Houses are, of course, eyecandies (at least to my taste). They are outliers with regards to size in their own neighborhoods!*

Sometimes it really pays off to be a [forensic statistician](#).