

Data visualization for real-world machine learning

Julia Silge

data visualization informs how we
think, understand,

decide

Computer-based visualization systems provide visual representations of datasets designed to help people carry out tasks more effectively.

— Tamara Munzner in *Visualization Analysis & Design*

The screenshot shows a web browser window with the title bar "Tidymodels". The address bar displays the URL "https://www.tidymodels.org". The main content area features the "Tidymodels" logo in pink at the top left. Below it is a navigation bar with links for "PACKAGES", "GET STARTED", "LEARN", "HELP", "CONTRIBUTE", a search icon, and a GitHub icon. To the left of the text content is a graphic of six hexagonal package icons arranged in a hexagonal pattern. The icons represent: "tidymodels" (dark blue, top), "rsample" (green, second from top-left), "parsnip" (tan, bottom-left), "yardstick" (red, bottom), "TUNE" (black, center), and "recipes" (light blue, second from top-right). The text "TIDYMODELS" is centered above the description. The description text reads: "The tidymodels framework is a collection of packages for modeling and machine learning using **tidyverse** principles." Below this is a section titled "Install tidymodels with:" containing the R code "install.packages("tidymodels")".

Tidymodels

PACKAGES GET STARTED LEARN HELP CONTRIBUTE

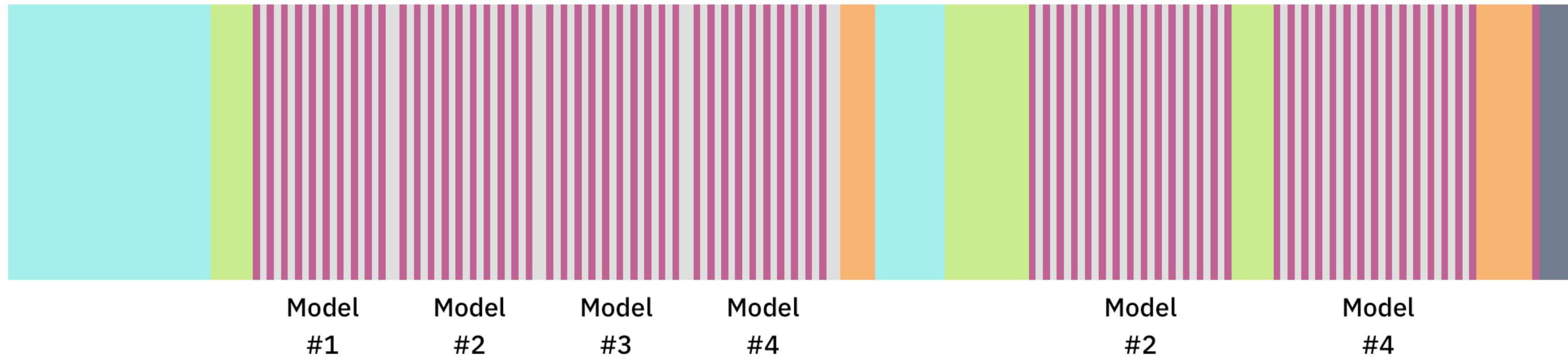
TIDYMODELS

The tidymodels framework is a collection of packages for modeling and machine learning using **tidyverse** principles.

Install tidymodels with:

```
install.packages("tidymodels")
```

When do practitioners build data visualizations?



Model
#1

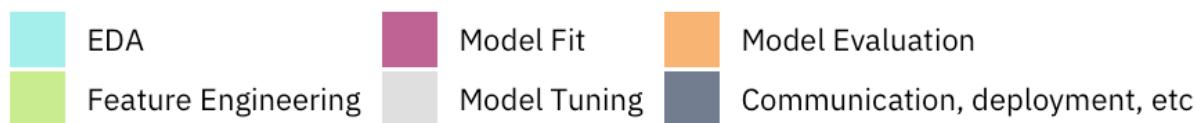
Model
#2

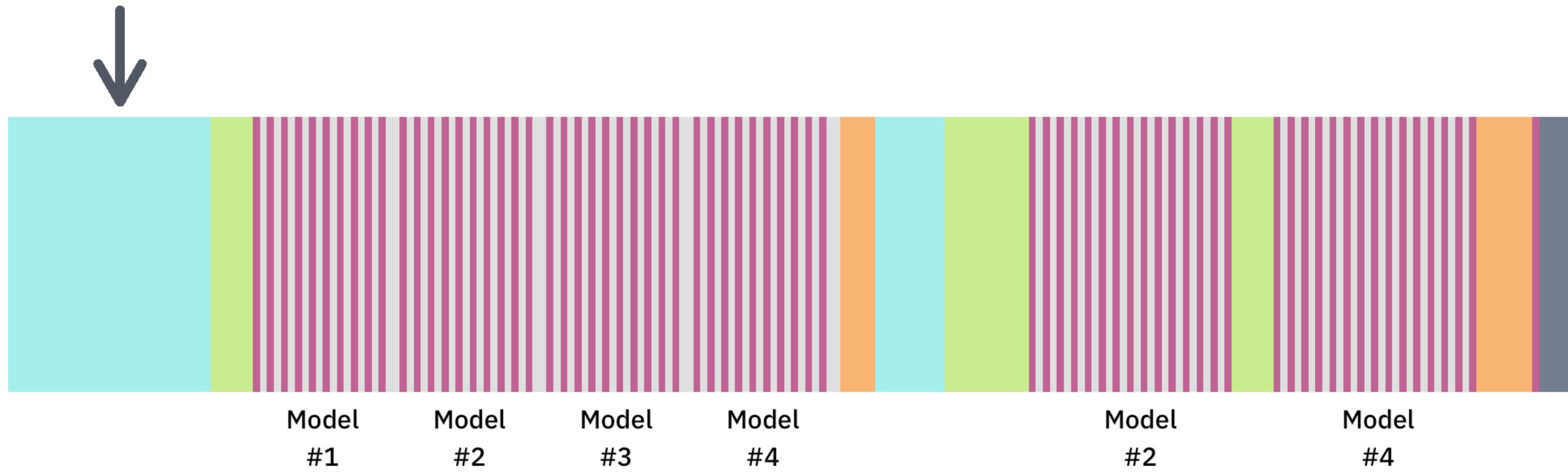
Model
#3

Model
#4

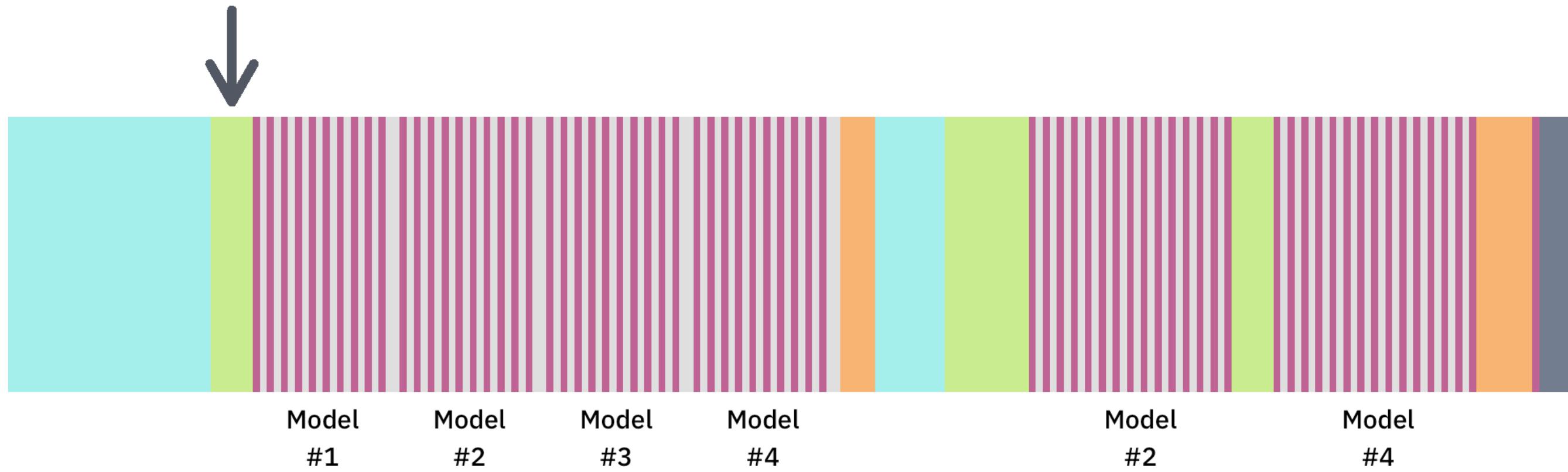
Model
#2

Model
#4

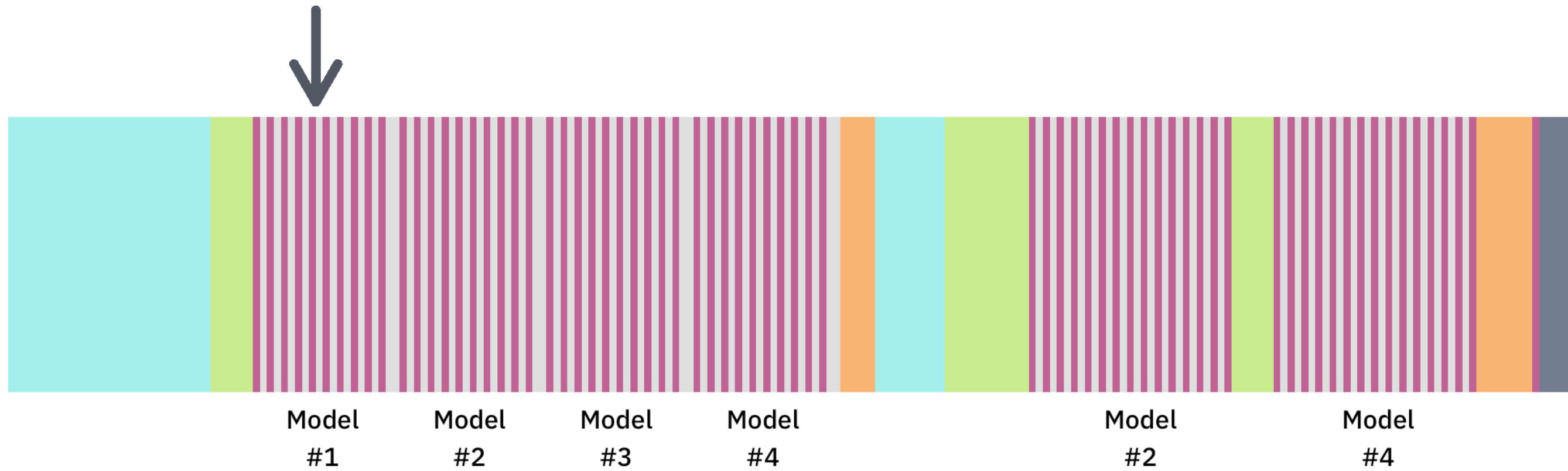




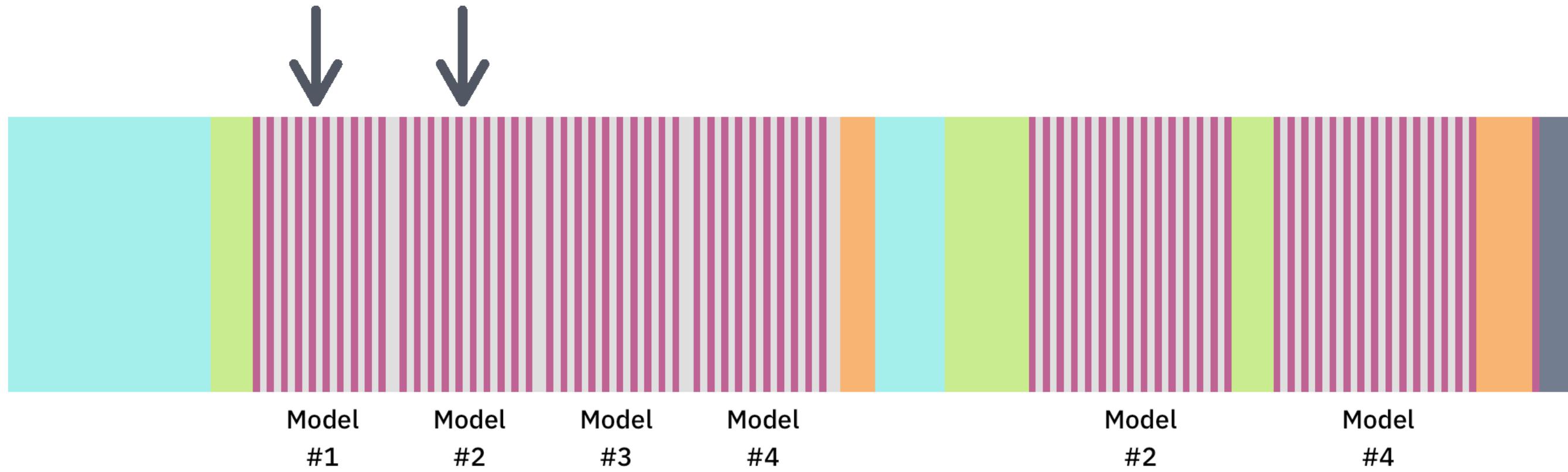
EDA Model Fit Model Evaluation
 Feature Engineering Model Tuning Communication, deployment, etc.



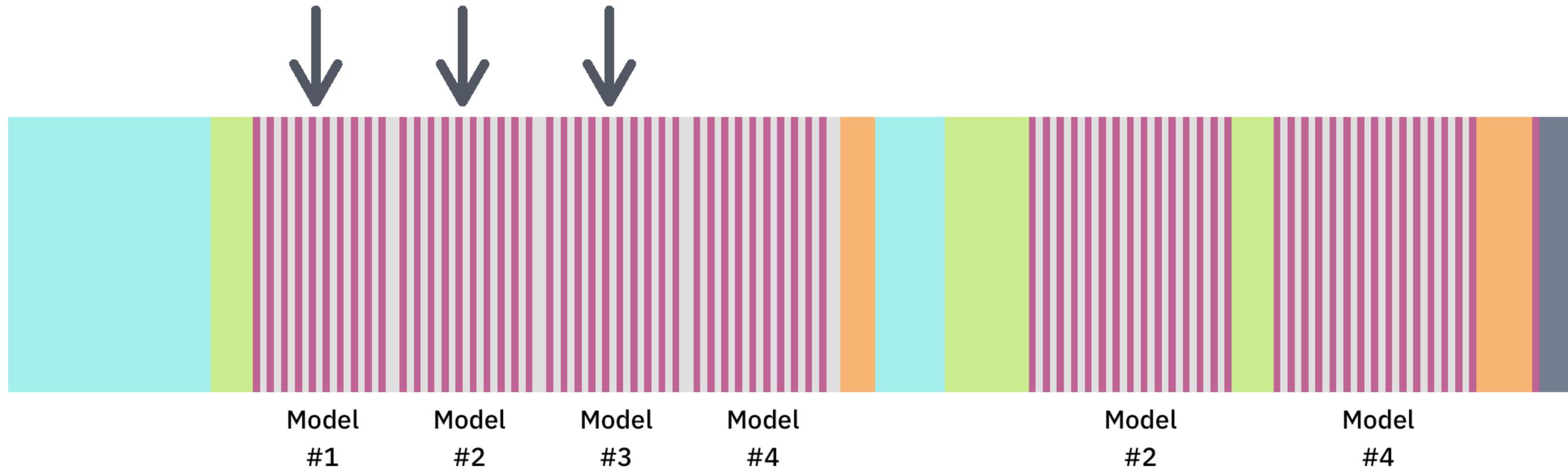
EDA Model Fit Model Evaluation
 Feature Engineering Model Tuning Communication, deployment, etc.



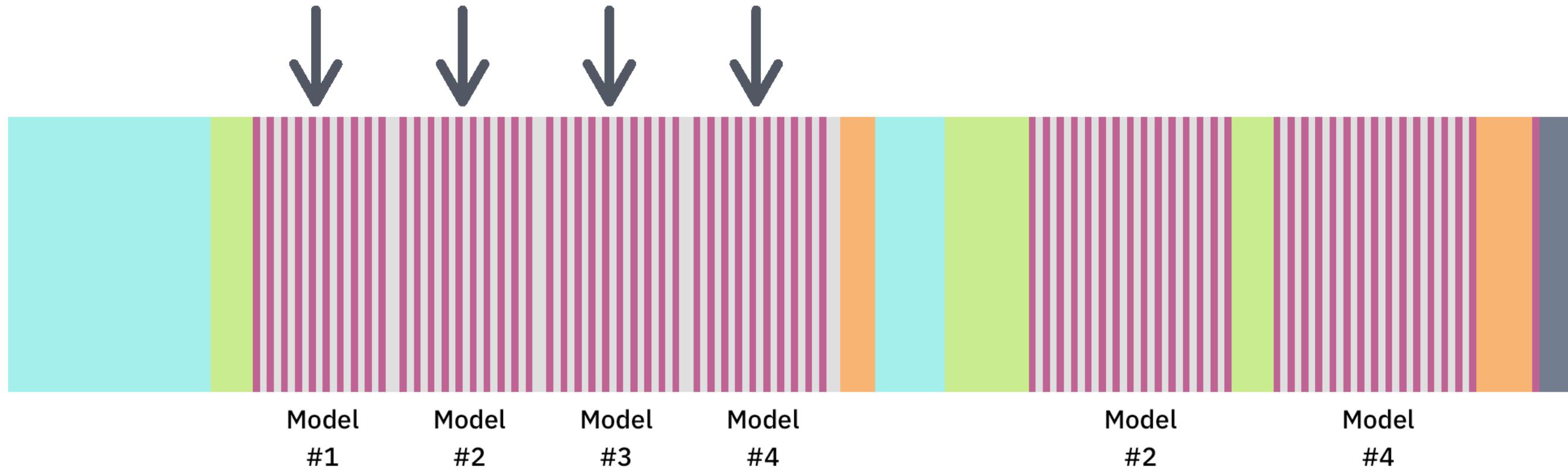
EDA Model Fit Model Evaluation
 Feature Engineering Model Tuning Communication, deployment, etc.



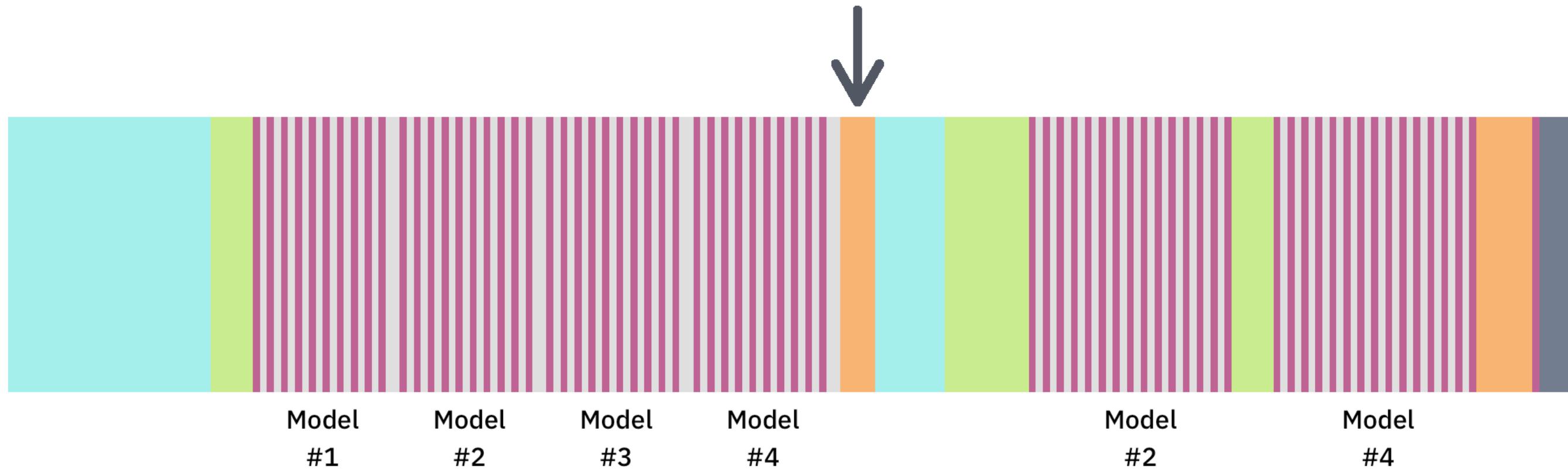
EDA Model Fit Model Evaluation
 Feature Engineering Model Tuning Communication, deployment, etc.



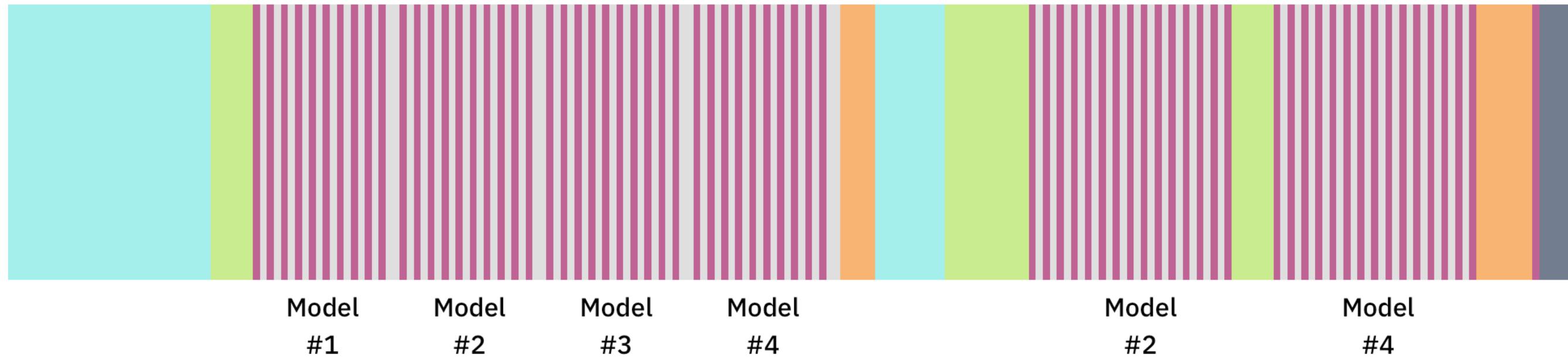
EDA
Feature Engineering
Model Fit
Model Tuning
Model Evaluation
Communication, deployment, etc.



	EDA		Model Fit		Model Evaluation
	Feature Engineering		Model Tuning		Communication, deployment, etc.



EDA Model Fit Model Evaluation
 Feature Engineering Model Tuning Communication, deployment, etc.



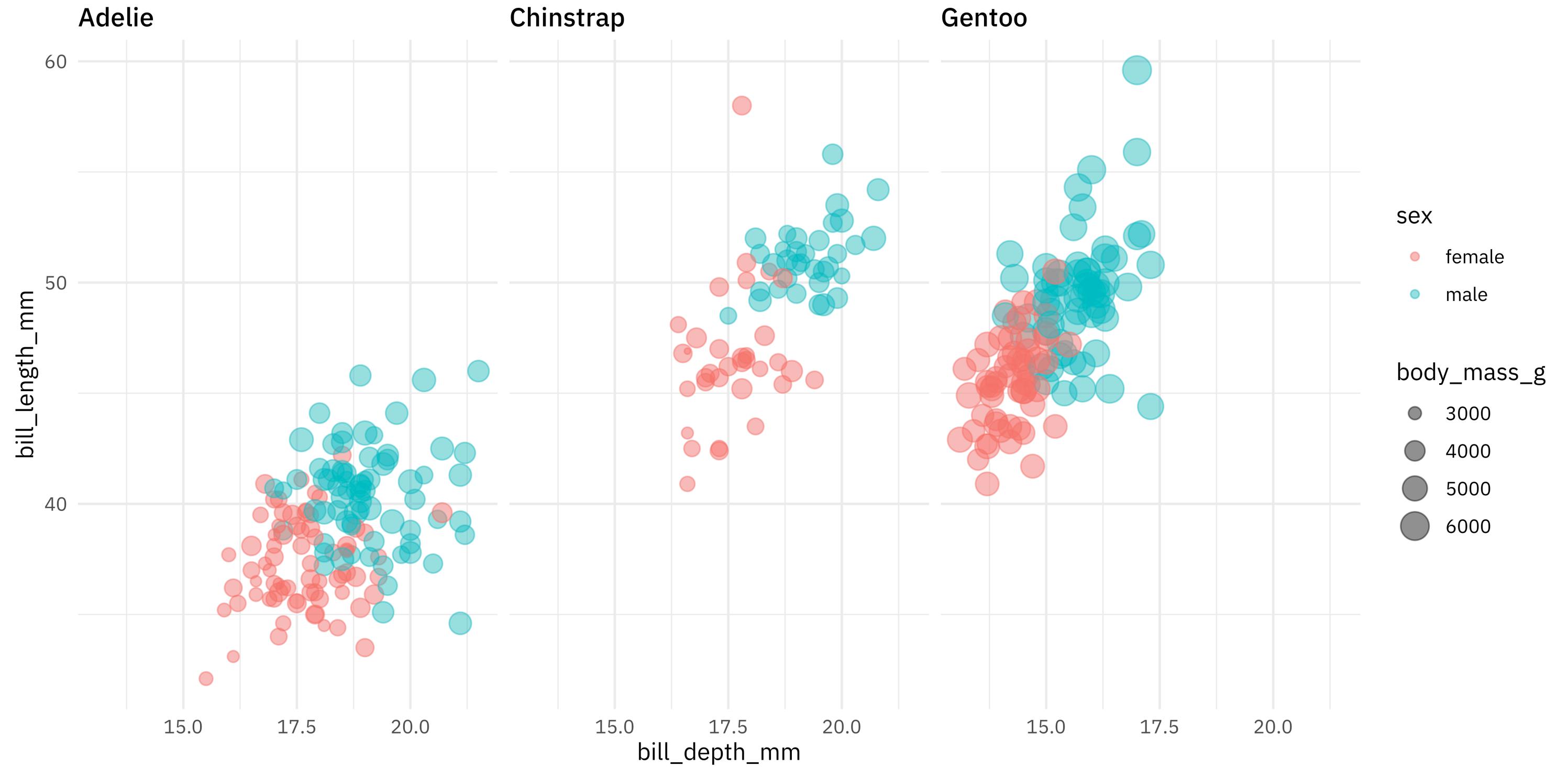
EDA Model Fit Model Evaluation
 Feature Engineering Model Tuning Communication, deployment, etc.

exploratory data analysis

model

evaluation

Why are these plots built?



Palmer penguins



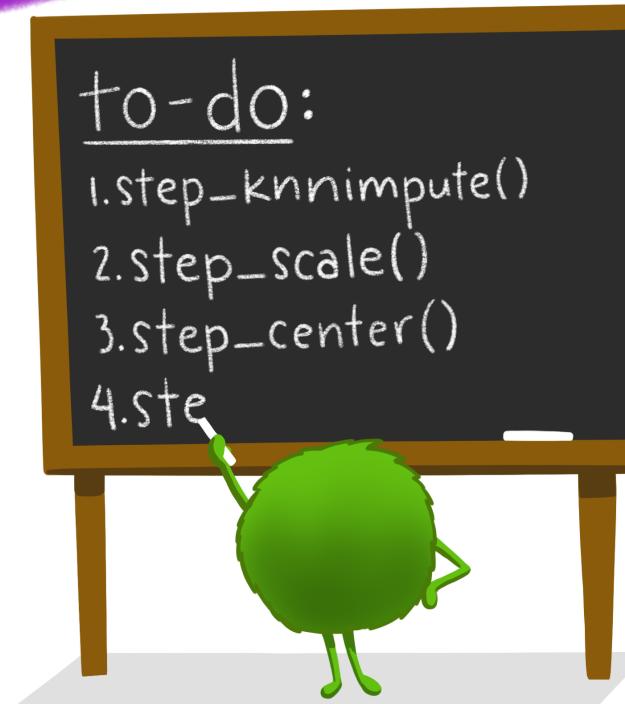
Trees in San Francisco

prioritize efficient iteration
for EDA

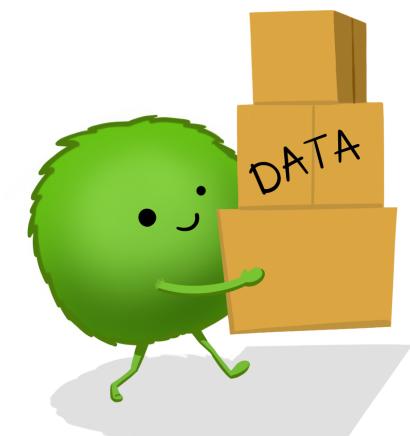


I. SPECIFY VARIABLES
`recipe(y~a+b+..., data=pantry)`

recipes:



STREAMLINED DATA PRE-PROCESSING FOR
STATISTICAL + MACHINE LEARNING MODELS



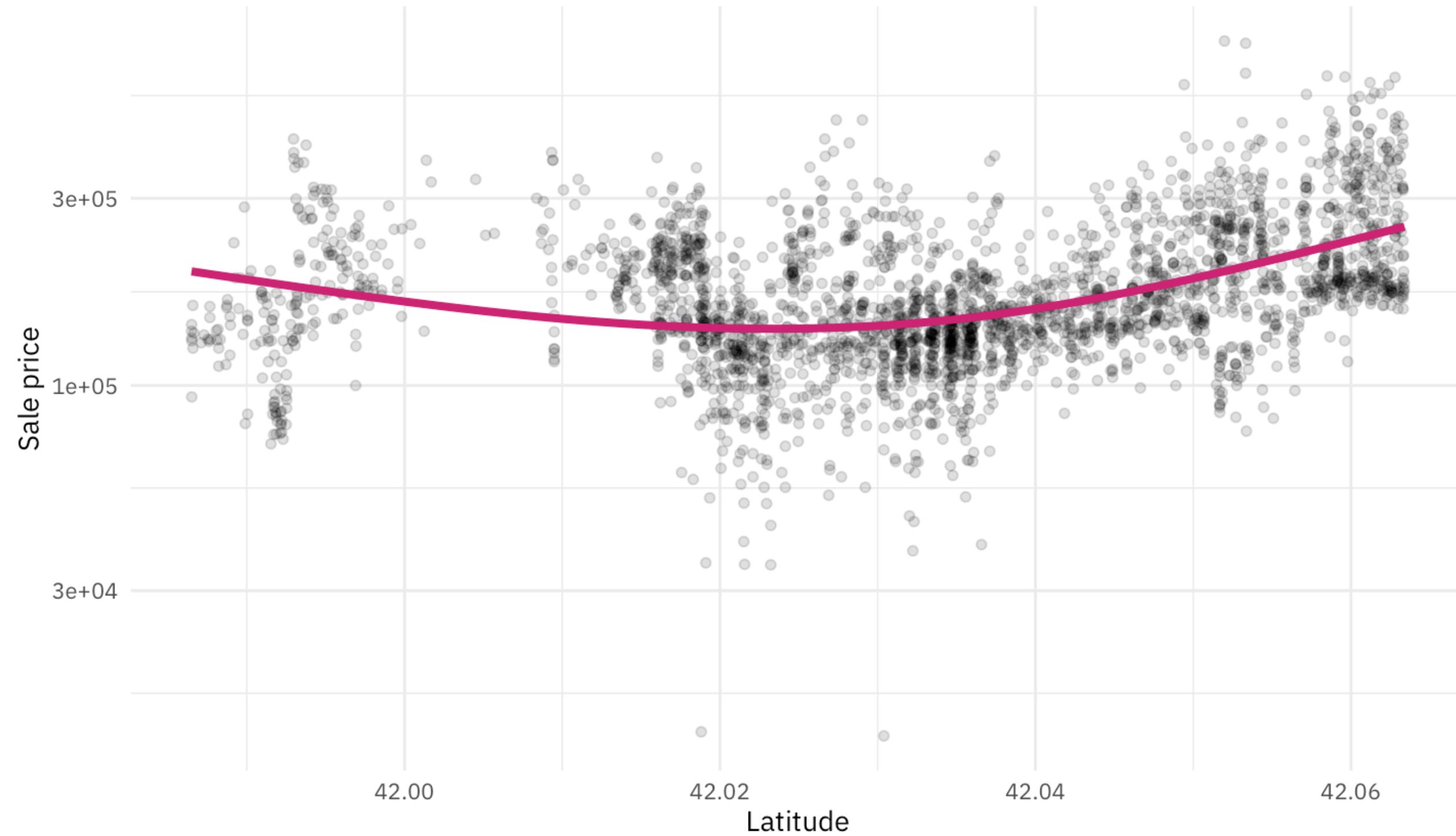
Horst 19

→ 2. DEFINE
PRE-PROCESSING
STEPS (step_*)

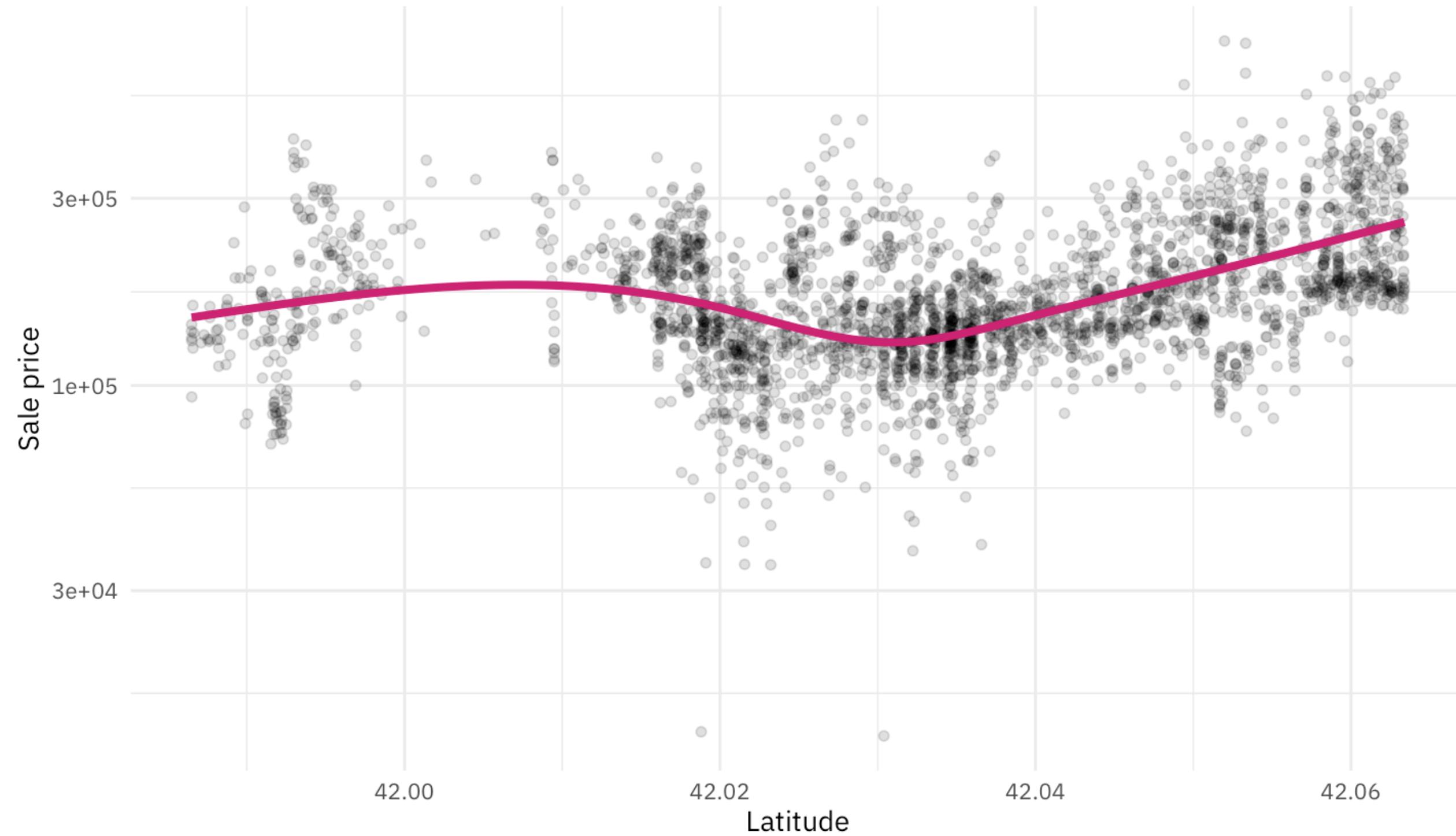
3. PROVIDE
DATASET(S) FOR
RECIPE STEPS
`prep()`

→ 4. APPLY
PRE-PROCESSING!
`bake()`

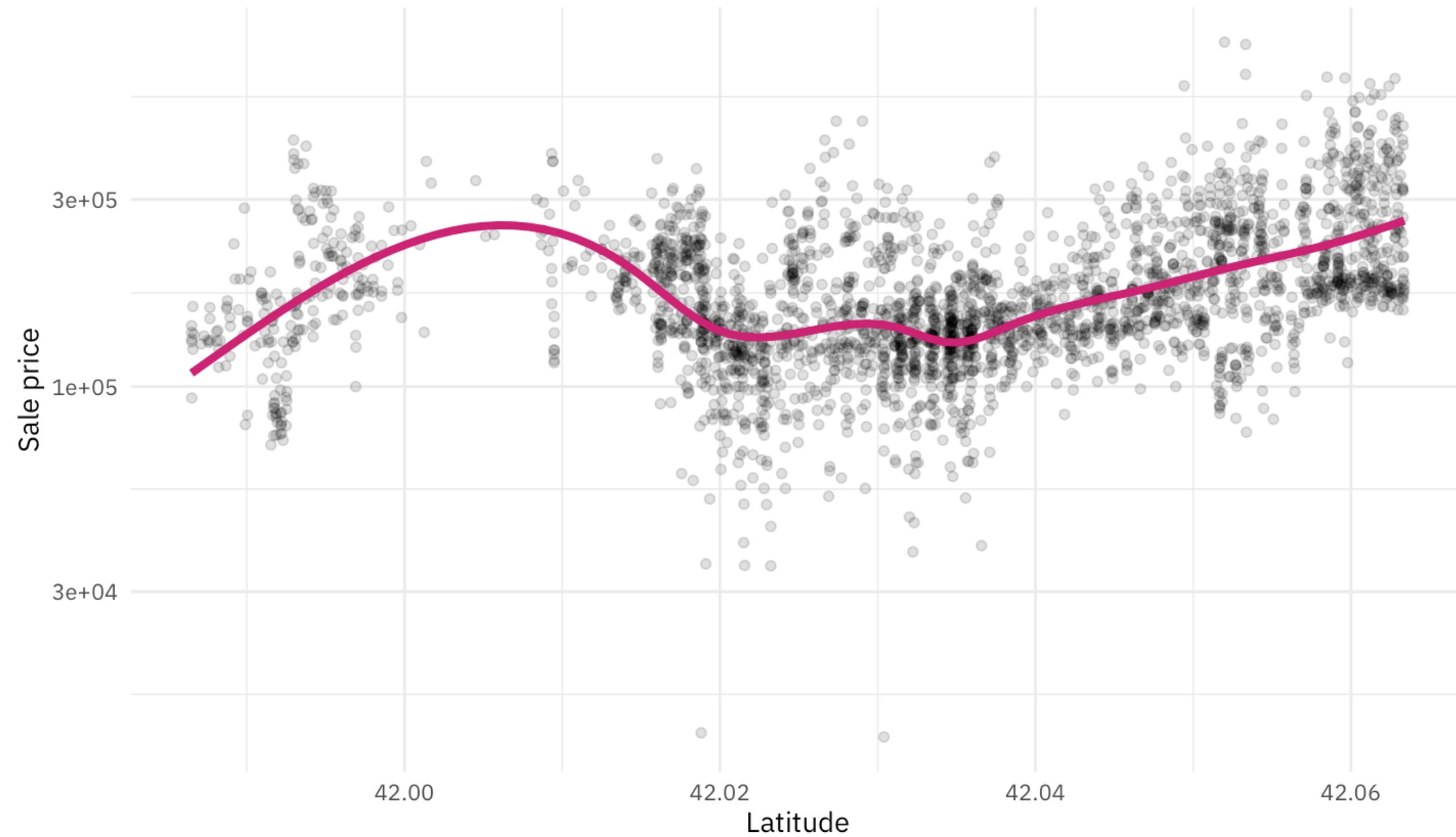
2 spline terms



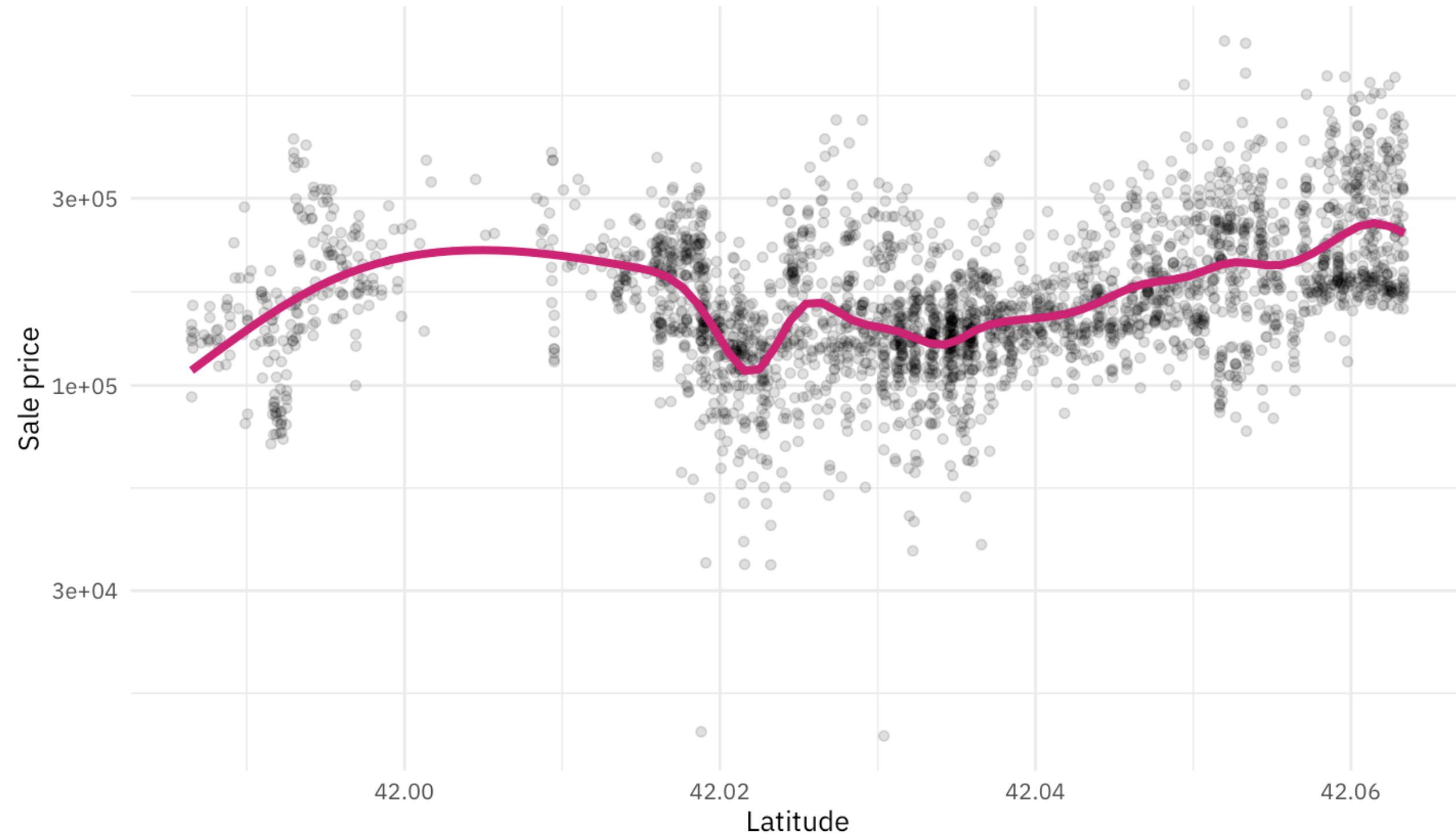
5 spline terms



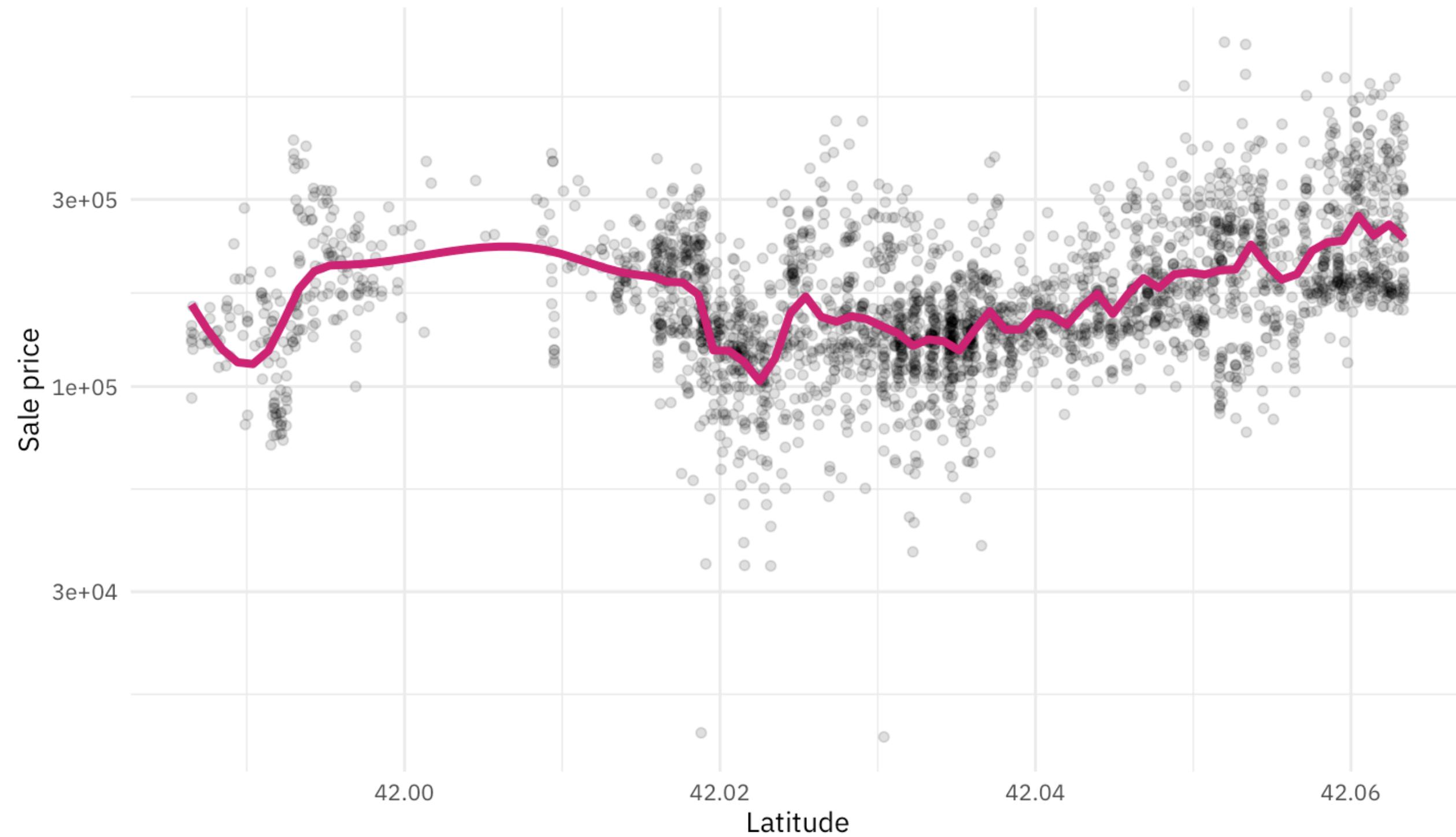
10 spline terms



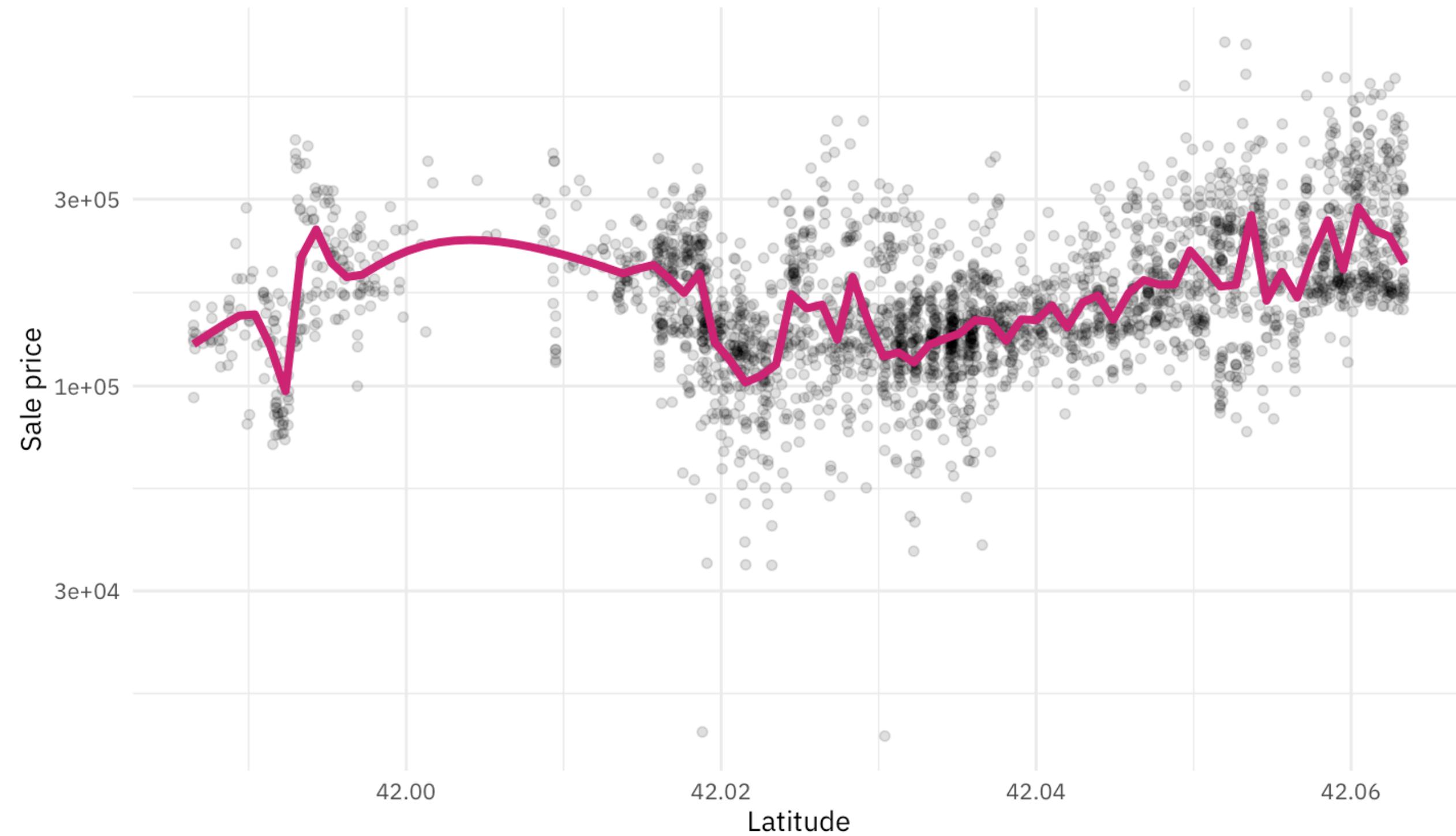
20 spline terms



50 spline terms



100 spline terms



```
recipe(Sale_Price ~ Neighborhood + Gr_Liv_Area + Year_Built +  
    Bldg_Type + Latitude + Longitude,  
    data = ames_train) %>%  
  step_log(Gr_Liv_Area, base = 10) %>%  
  step_other(Neighborhood, threshold = 0.01) %>%  
  step_dummy(all_nominal()) %>%  
  step_ns(Latitude, Longitude, deg_free = 10)
```

```
## Data Recipe  
##  
## Inputs:  
##  
##      role #variables  
##      outcome          1  
##      predictor         6  
##  
## Operations:  
##  
## Log transformation on Gr_Liv_Area  
## Collapsing factor levels for Neighborhood  
## Dummy variables from all_nominal()  
## Natural Splines on Latitude, Longitude
```

Practitioners use data
visualization during EDA to
inform modeling choices



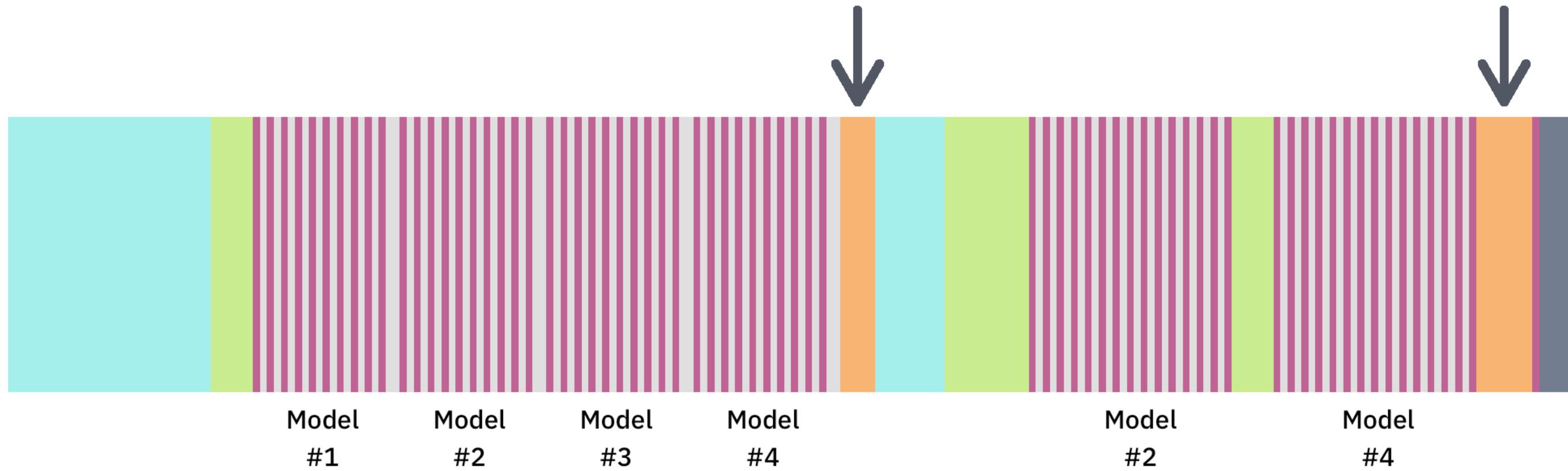
Pinboard
@Pinboard

Machine learning is like a deep-fat fryer. First time you try it you think "Amazing, I bet this will work on anything!"
And it kind of does

10:16 PM · Jul 6, 2016 · [YoruFukurou](#)

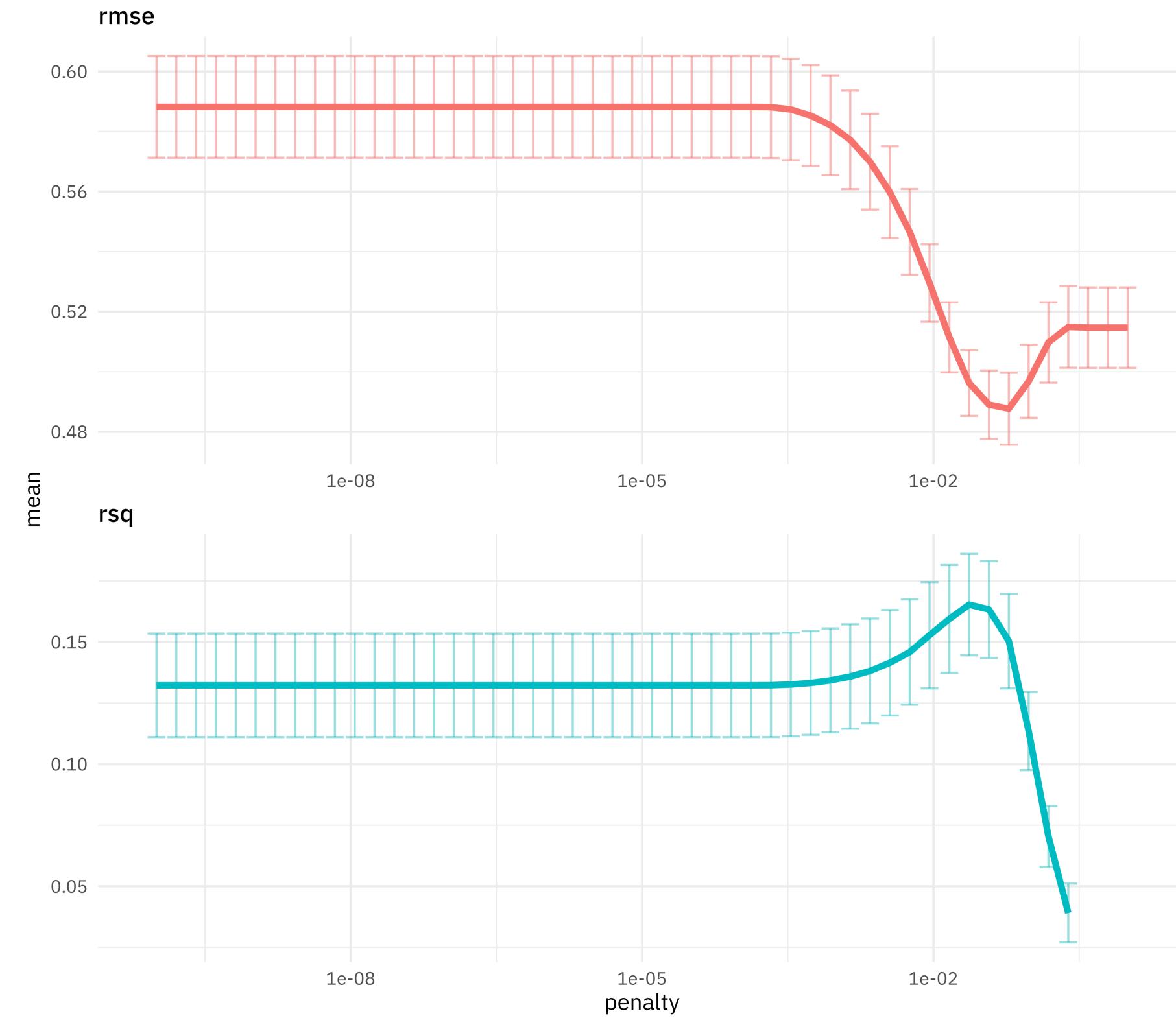
321 Retweets and comments **501** Likes

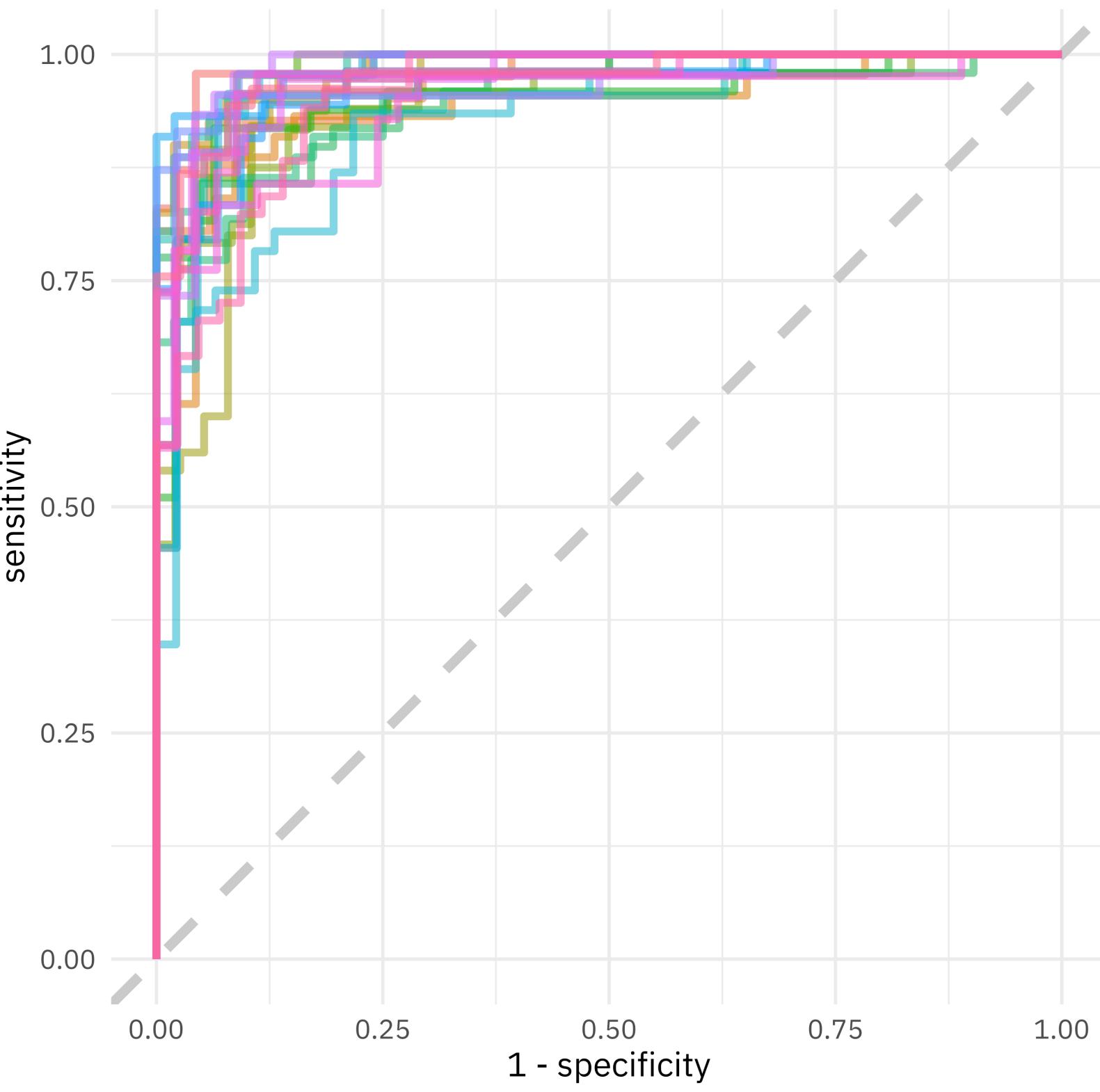




	EDA		Model Fit		Model Evaluation
	Feature Engineering		Model Tuning		Communication, deployment, etc.

Why are these plots built?



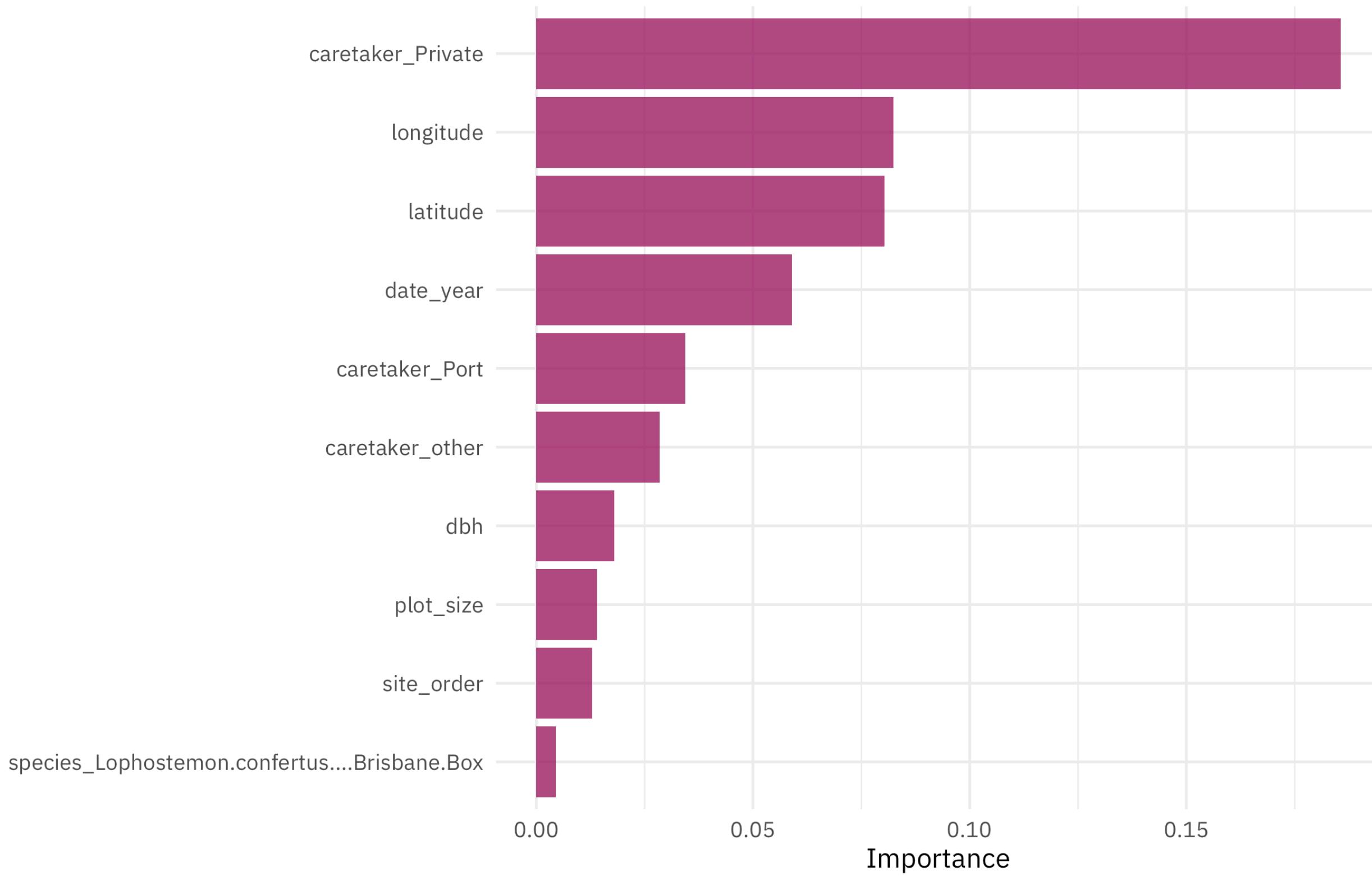


automatic plotting methods

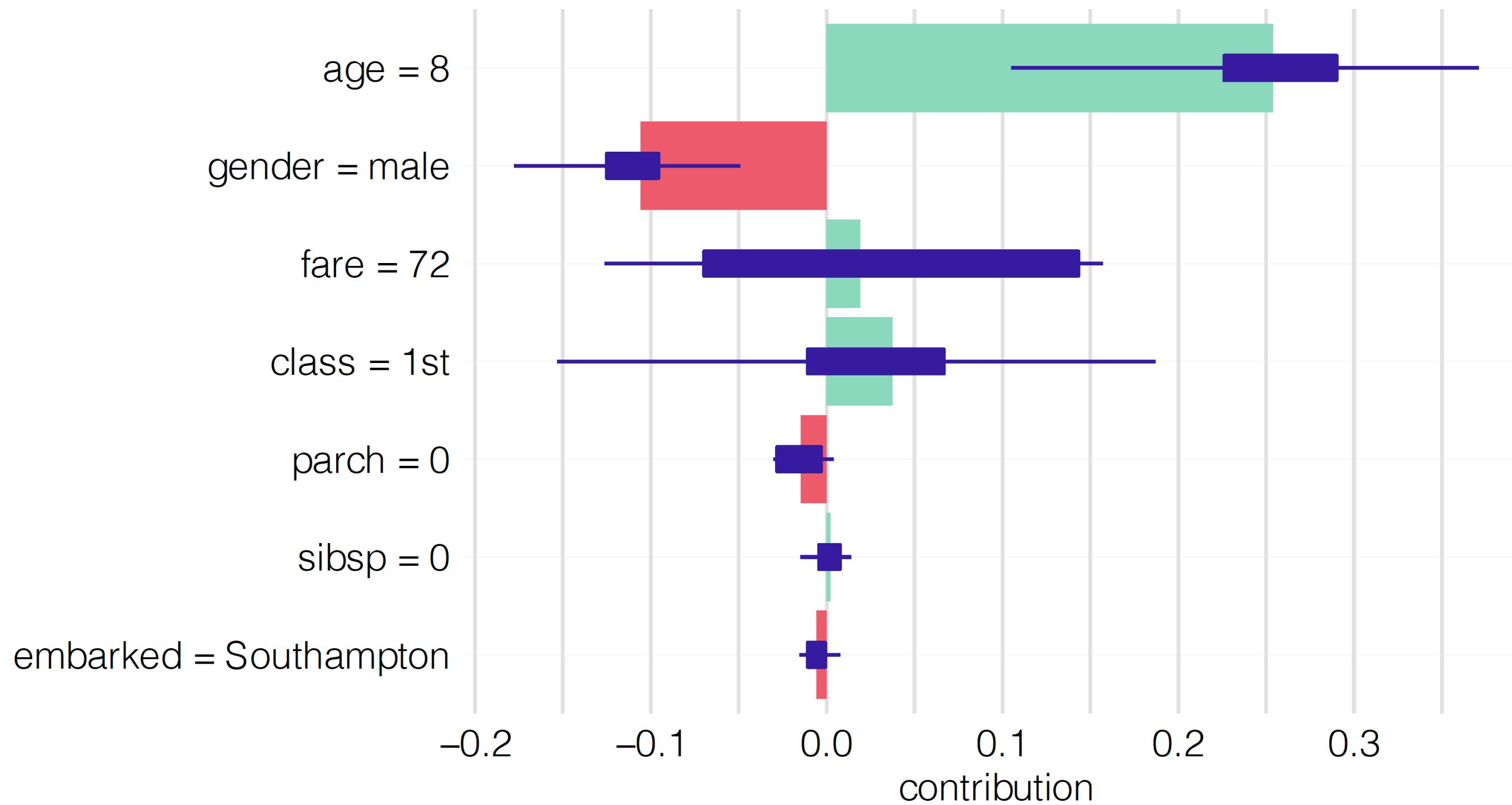


Volcano classification

explaining models



Average attributions for Johnny D



Who are these plots for?

understand
1
gatua

understand
models

Thank you!

Julia Silge