

# Data visualization for real-world machine learning

---

Julia Silge

data visualization informs how we  
think, understand,

decide

---

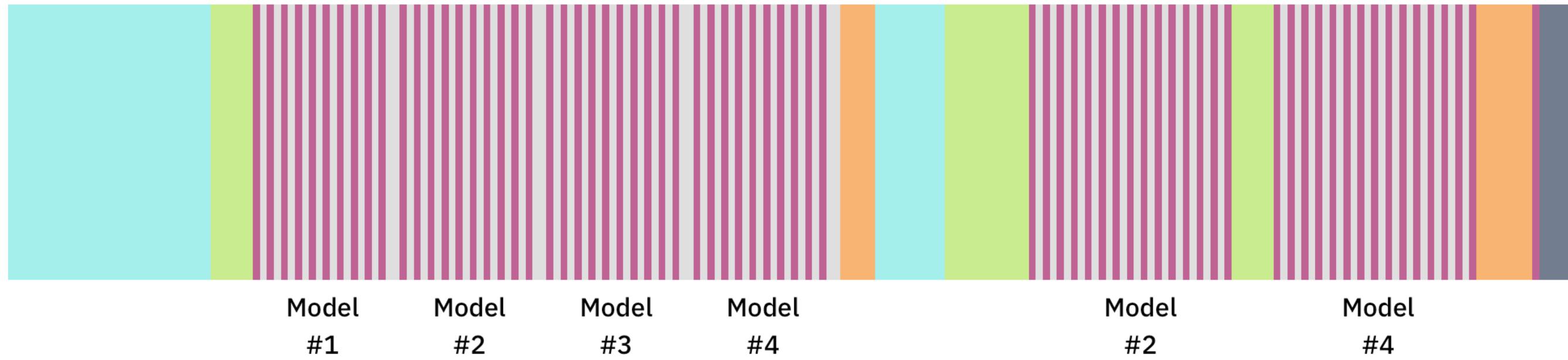
Computer-based visualization systems provide visual representations of datasets designed to help people carry out tasks more effectively.

— Tamara Munzner in *Visualization Analysis & Design*

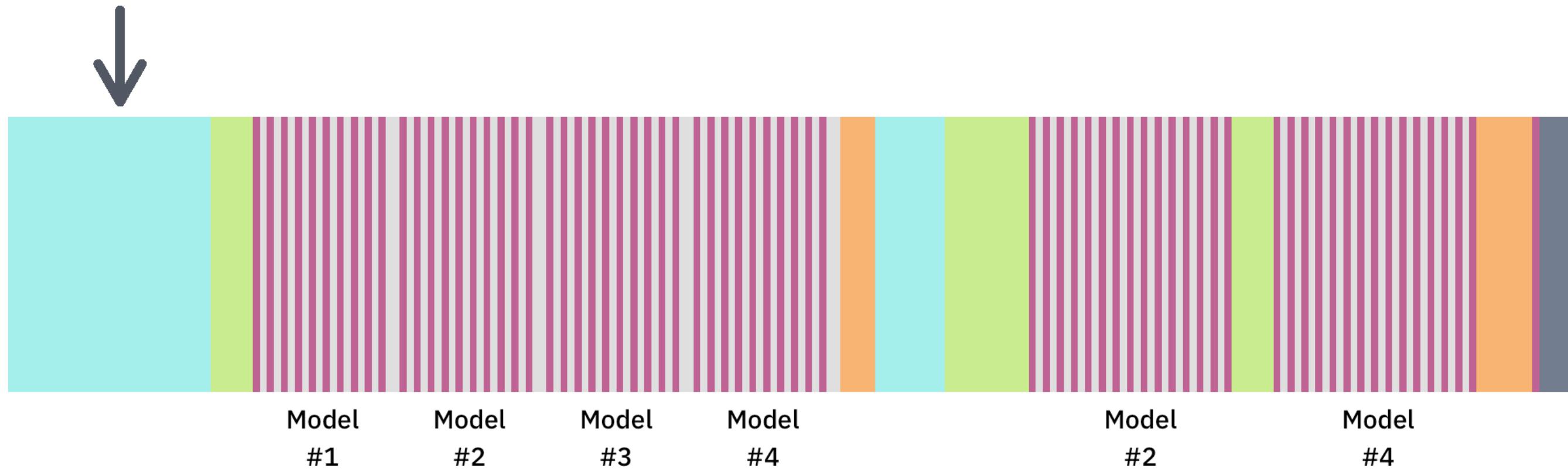
---

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. The central visual element is a hexagonal grid composed of several hexagons, each containing a different package icon from the tidymodels framework. The visible packages include "tidymodels" (dark blue), "rsample" (green, featuring a boot icon), "parsnip" (tan, featuring a plant icon), "recipes" (light blue, featuring a cupcake icon), "TUNE" (black, featuring a colorful grid icon), and "yardstick" (red, featuring a ruler icon). To the right of the grid, the word "TIDYMODELS" is written in bold capital letters. Below this, a paragraph explains the framework: "The tidymodels framework is a collection of packages for modeling and machine learning using **tidyverse** principles." Further down, the text "Install tidymodels with:" is followed by a code block containing the R command "install.packages("tidymodels")".

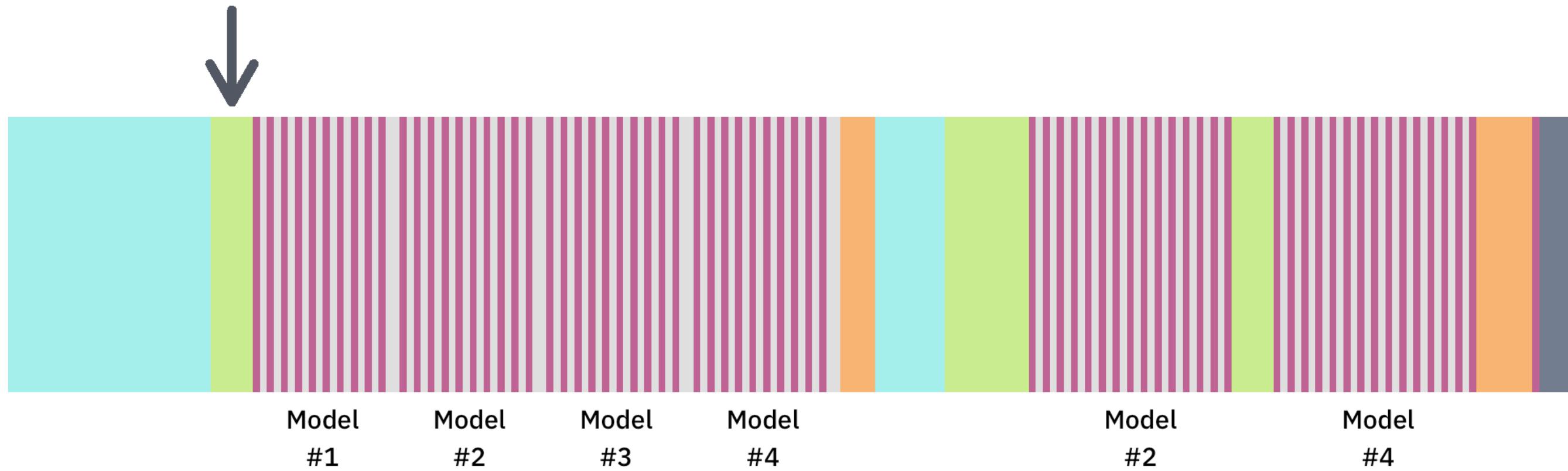
# When do practitioners build data visualizations?



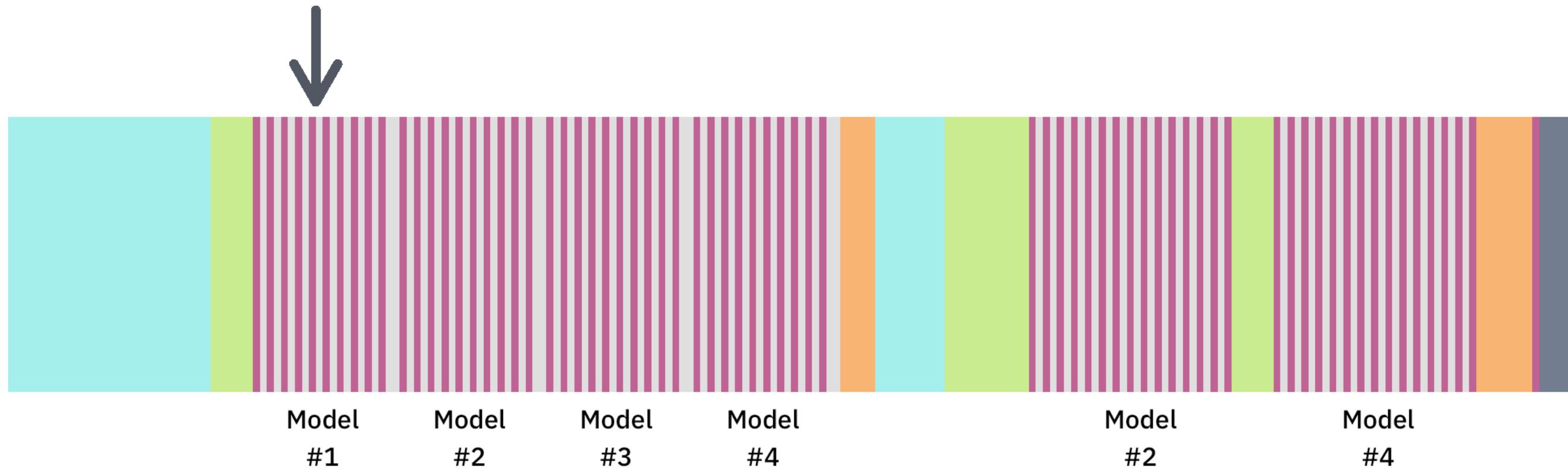
<span style="background-color: cyan; border: 1px solid black; display: inline-block; width: 10px; height: 10px;"></span>	EDA	<span style="background-color: maroon; border: 1px solid black; display: inline-block; width: 10px; height: 10px;"></span>	Model Fit	<span style="background-color: orange; border: 1px solid black; display: inline-block; width: 10px; height: 10px;"></span>	Model Evaluation
<span style="background-color: lightgreen; border: 1px solid black; display: inline-block; width: 10px; height: 10px;"></span>	Feature Engineering	<span style="background-color: grey; border: 1px solid black; display: inline-block; width: 10px; height: 10px;"></span>	Model Tuning	<span style="background-color: darkgrey; border: 1px solid black; display: inline-block; width: 10px; height: 10px;"></span>	Communication, deployment, etc.



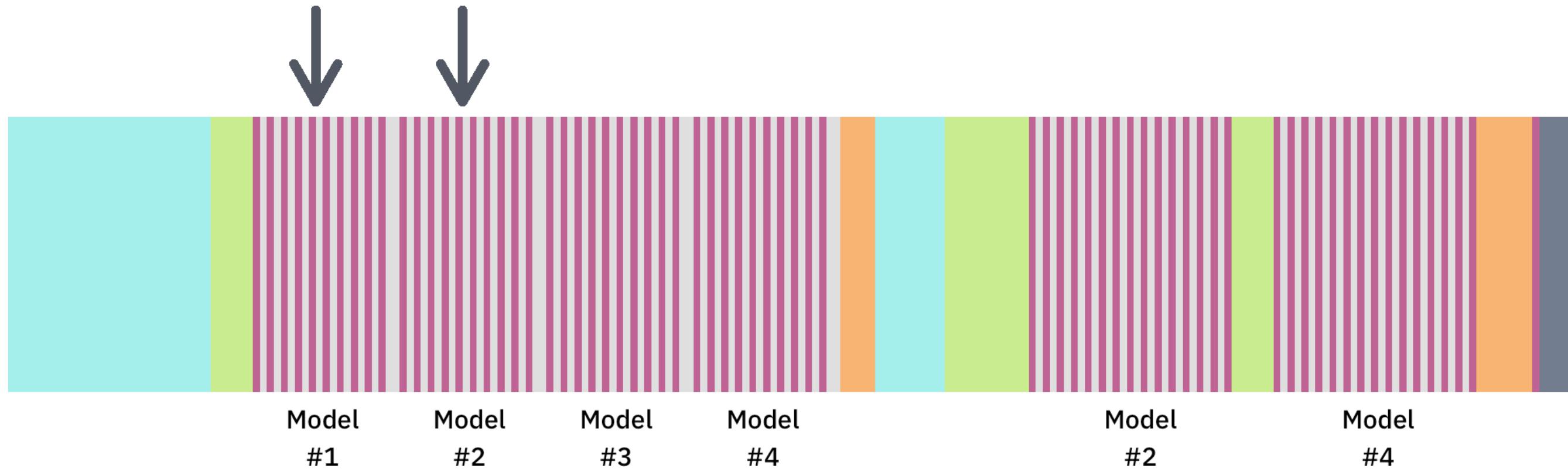
EDA      Model Fit      Model Evaluation  
 Feature Engineering      Model Tuning      Communication, deployment, etc.



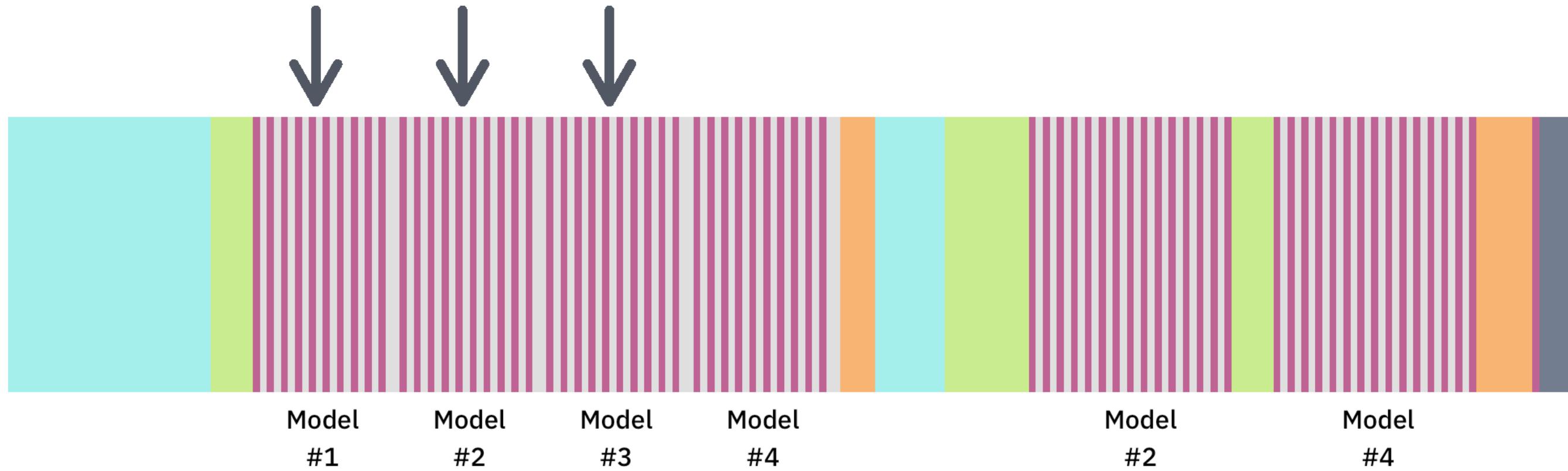
EDA      Model Fit      Model Evaluation  
 Feature Engineering      Model Tuning      Communication, deployment, etc.



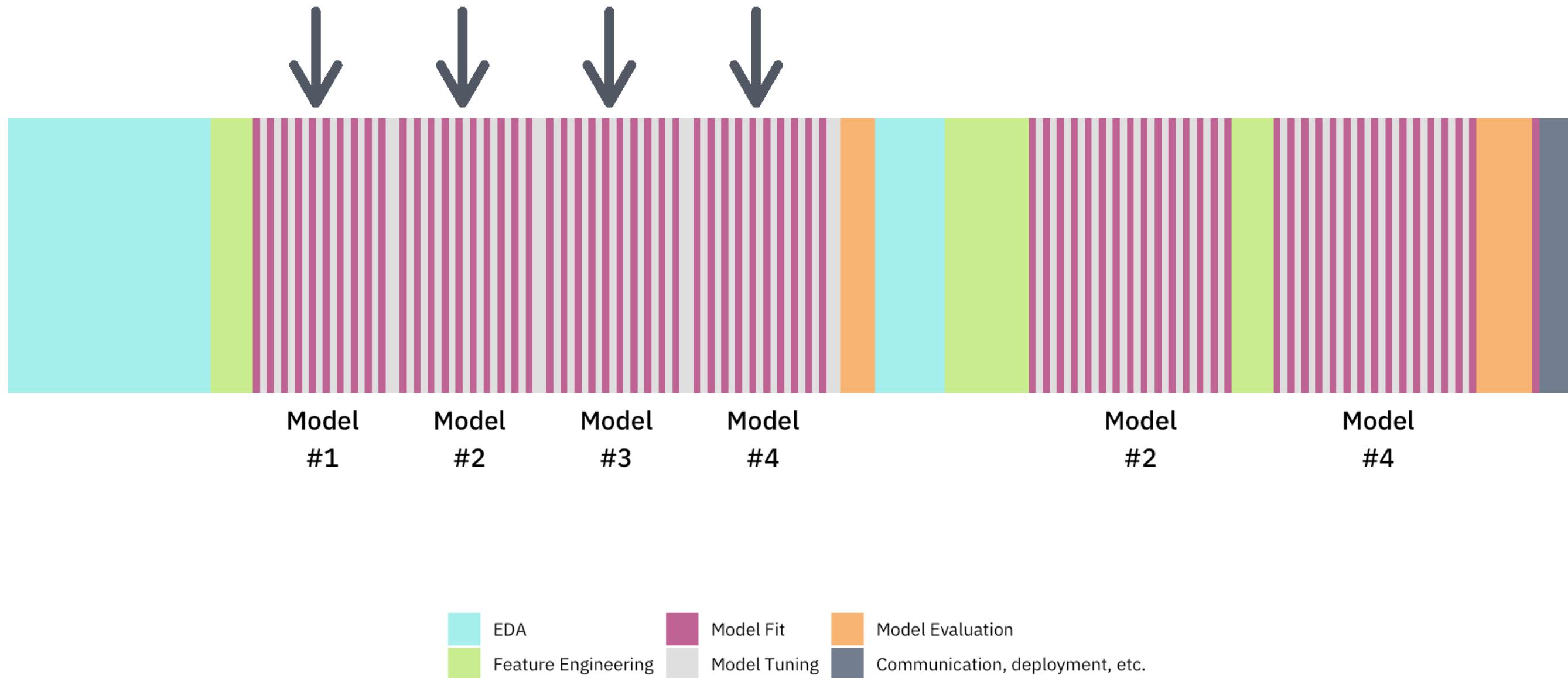
EDA      Model Fit      Model Evaluation  
 Feature Engineering      Model Tuning      Communication, deployment, etc.

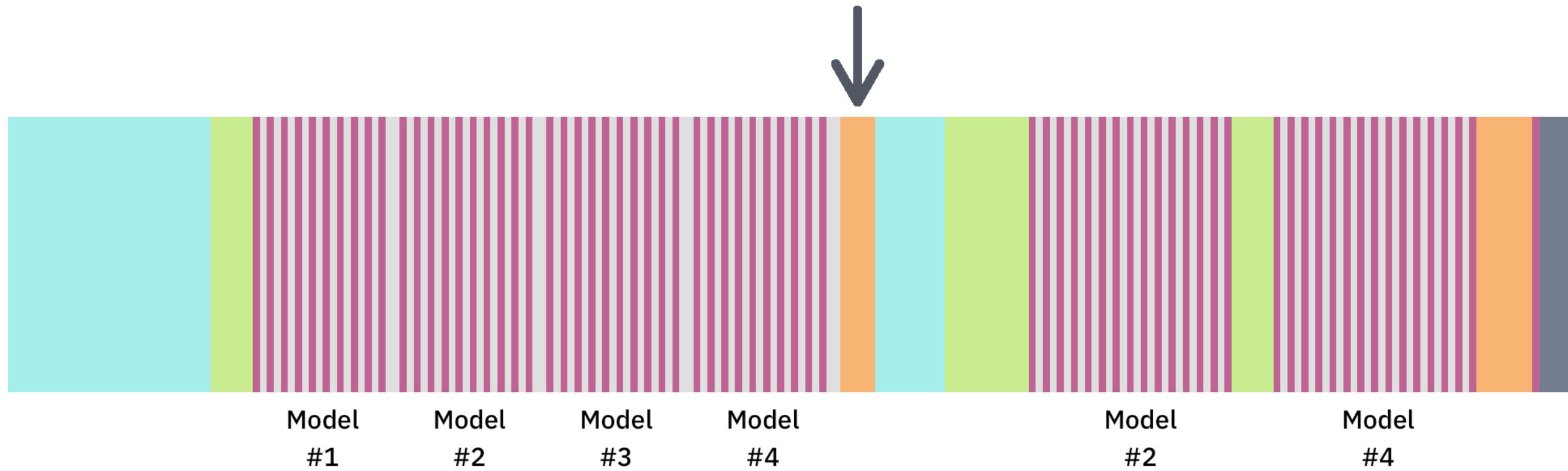


EDA  
Feature Engineering  
Model Fit  
Model Tuning  
Model Evaluation  
Communication, deployment, etc.

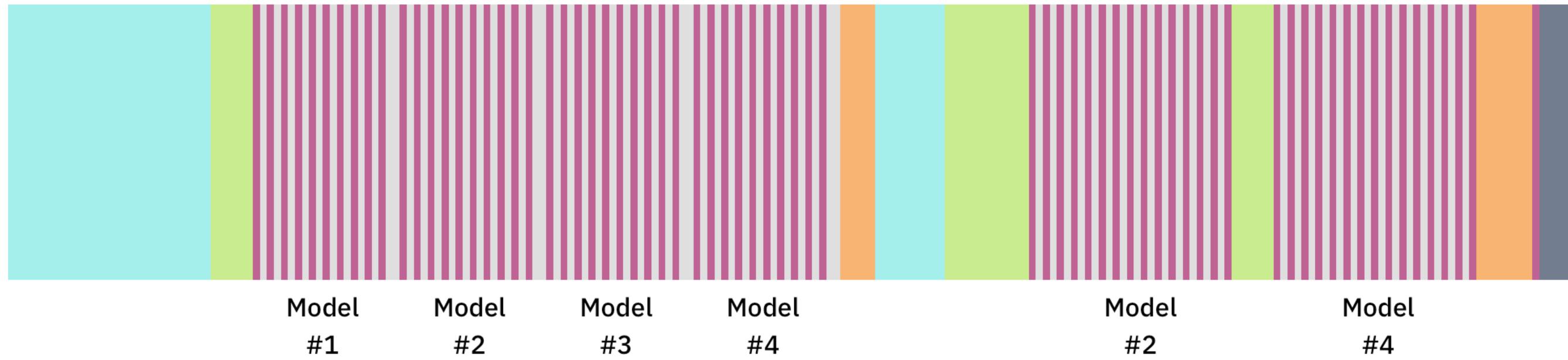


	EDA		Model Fit		Model Evaluation
	Feature Engineering		Model Tuning		Communication, deployment, etc.





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



<span style="background-color: cyan; border: 1px solid black; display: inline-block; width: 10px; height: 10px;"></span>	EDA	<span style="background-color: maroon; border: 1px solid black; display: inline-block; width: 10px; height: 10px;"></span>	Model Fit	<span style="background-color: orange; border: 1px solid black; display: inline-block; width: 10px; height: 10px;"></span>	Model Evaluation
<span style="background-color: lightgreen; border: 1px solid black; display: inline-block; width: 10px; height: 10px;"></span>	Feature Engineering	<span style="background-color: grey; border: 1px solid black; display: inline-block; width: 10px; height: 10px;"></span>	Model Tuning	<span style="background-color: darkgrey; border: 1px solid black; display: inline-block; width: 10px; height: 10px;"></span>	Communication, deployment, etc.

# exploratory data analysis

model

evaluation

# Why are these plots built?

# Visualization for decision making

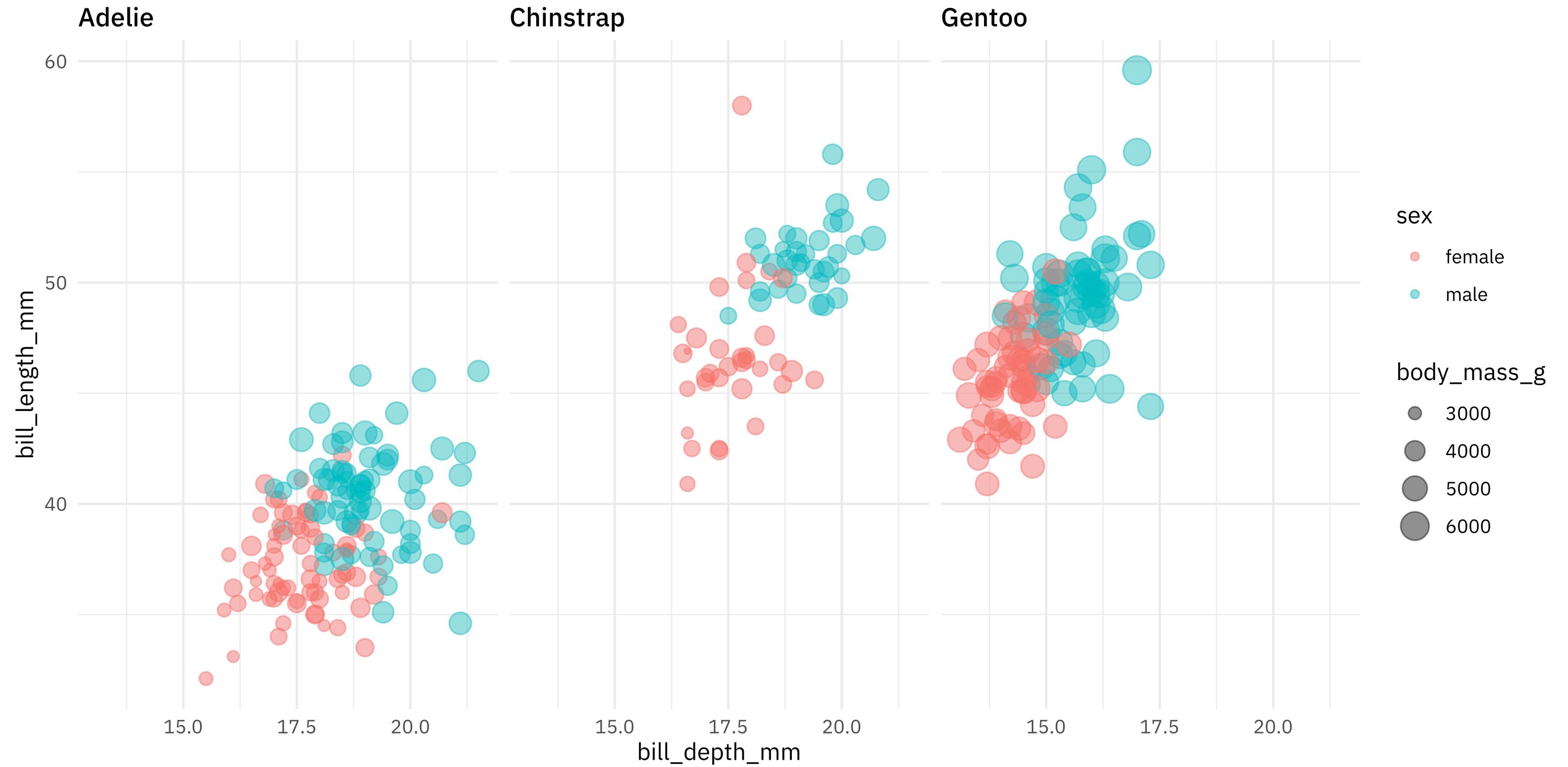
- What kind of model is appropriate?

# Visualization for decision making

- What kind of model is appropriate?
- What kind of preprocessing is needed?

# Visualization for decision making

- What kind of model is appropriate?
- What kind of preprocessing is needed?
- Is anything unusual going on with this data?



Palmer penguins



Trees in San Francisco

prioritize efficient iteration  
for EDA

# What kind of model is appropriate?

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

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

# What kind of model is appropriate?

```
linear_reg() %>%  
  set_engine("glmnet")
```

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

# What kind of model is appropriate?

```
rand_forest() %>%  
  set_engine("ranger") %>%  
  set_mode("regression")  
  
## Random Forest Model Specification (regression)  
##  
## Computational engine: ranger
```

# What kind of model is appropriate?

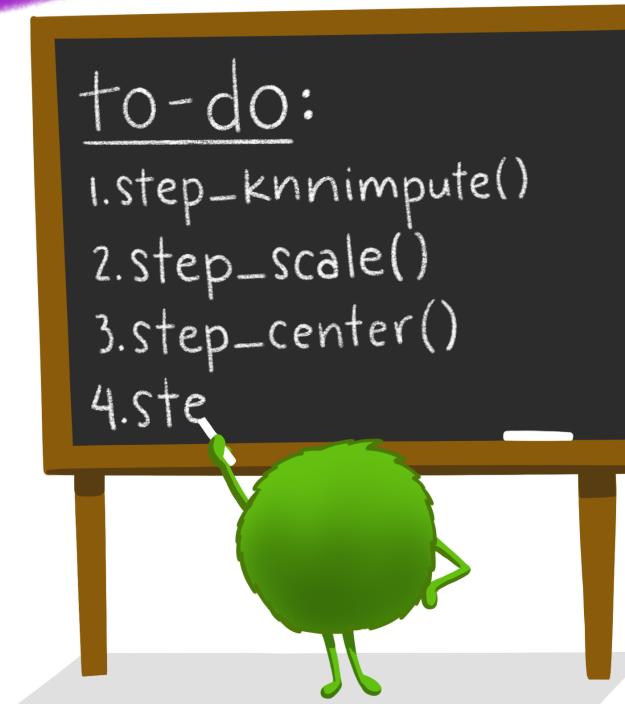
```
rand_forest() %>%  
  set_engine("ranger") %>%  
  set_mode("classification")
```

```
## Random Forest Model Specification (classification)  
##  
## Computational engine: ranger
```



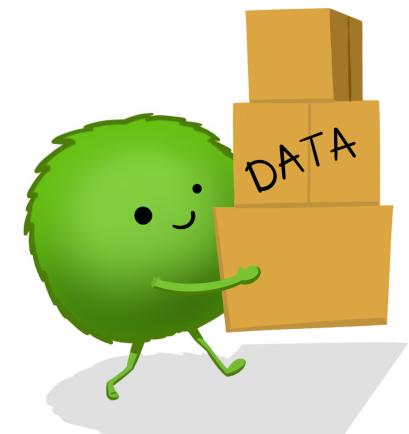
I. SPECIFY VARIABLES  
`recipe(y~a+b+..., data=pantry)`

# recipes:



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

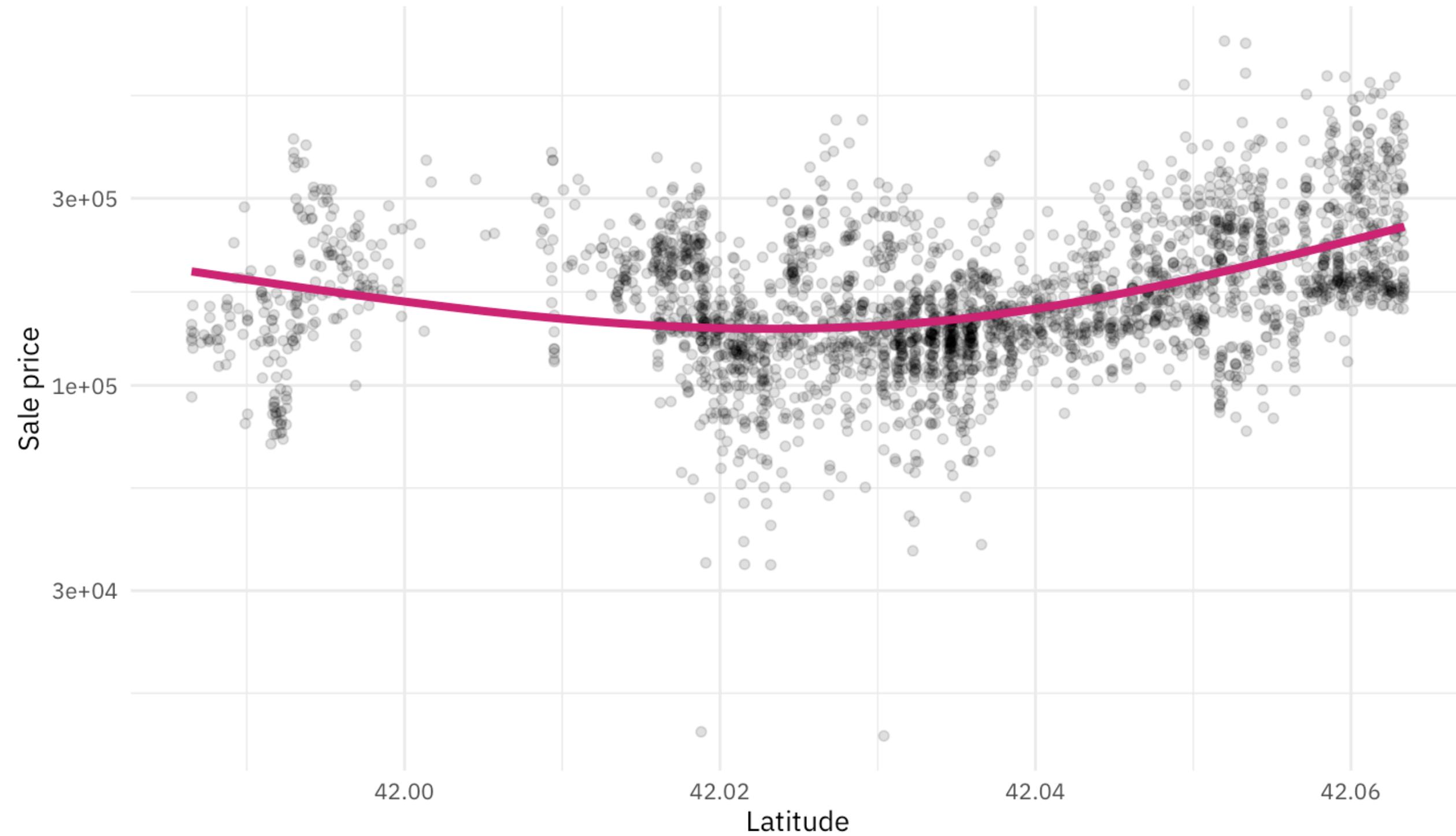
STREAMLINED DATA PRE-PROCESSING FOR  
STATISTICAL + MACHINE LEARNING MODELS



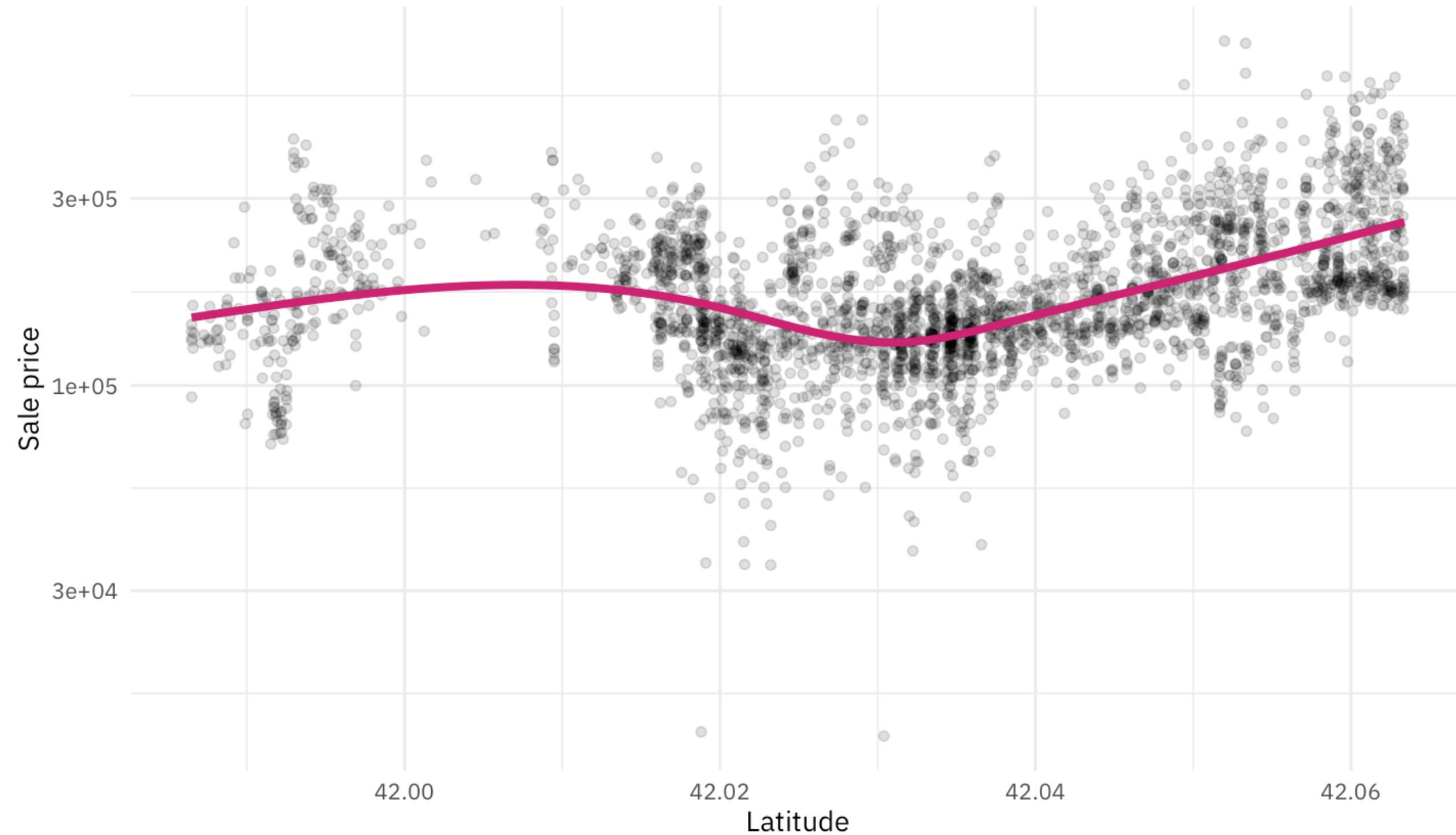
3. PROVIDE  
DATASET(S) FOR  
RECIPE STEPS  
`prep()`

4. APPLY  
PRE-PROCESSING!  
`bake()`

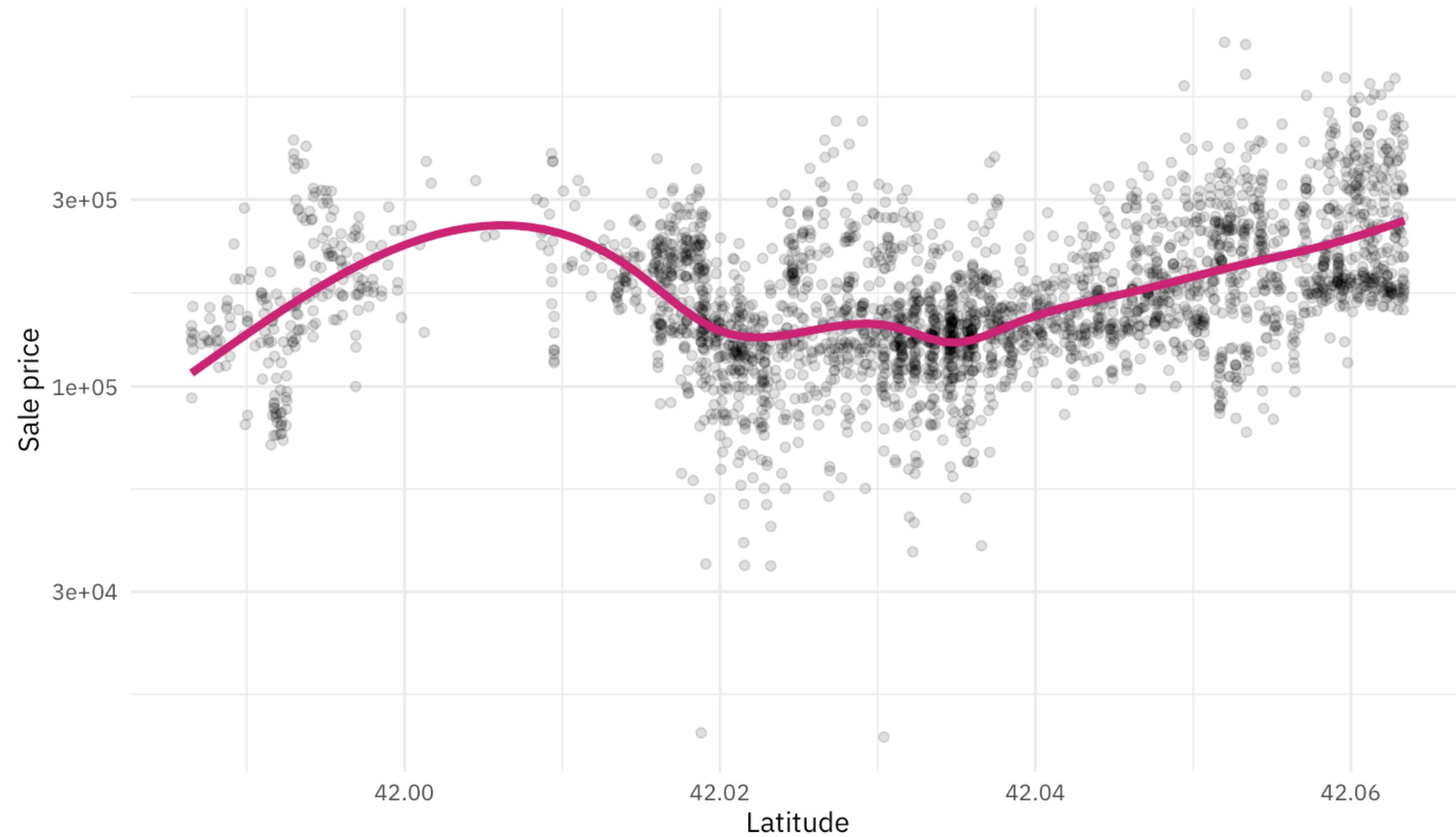
## 2 spline terms



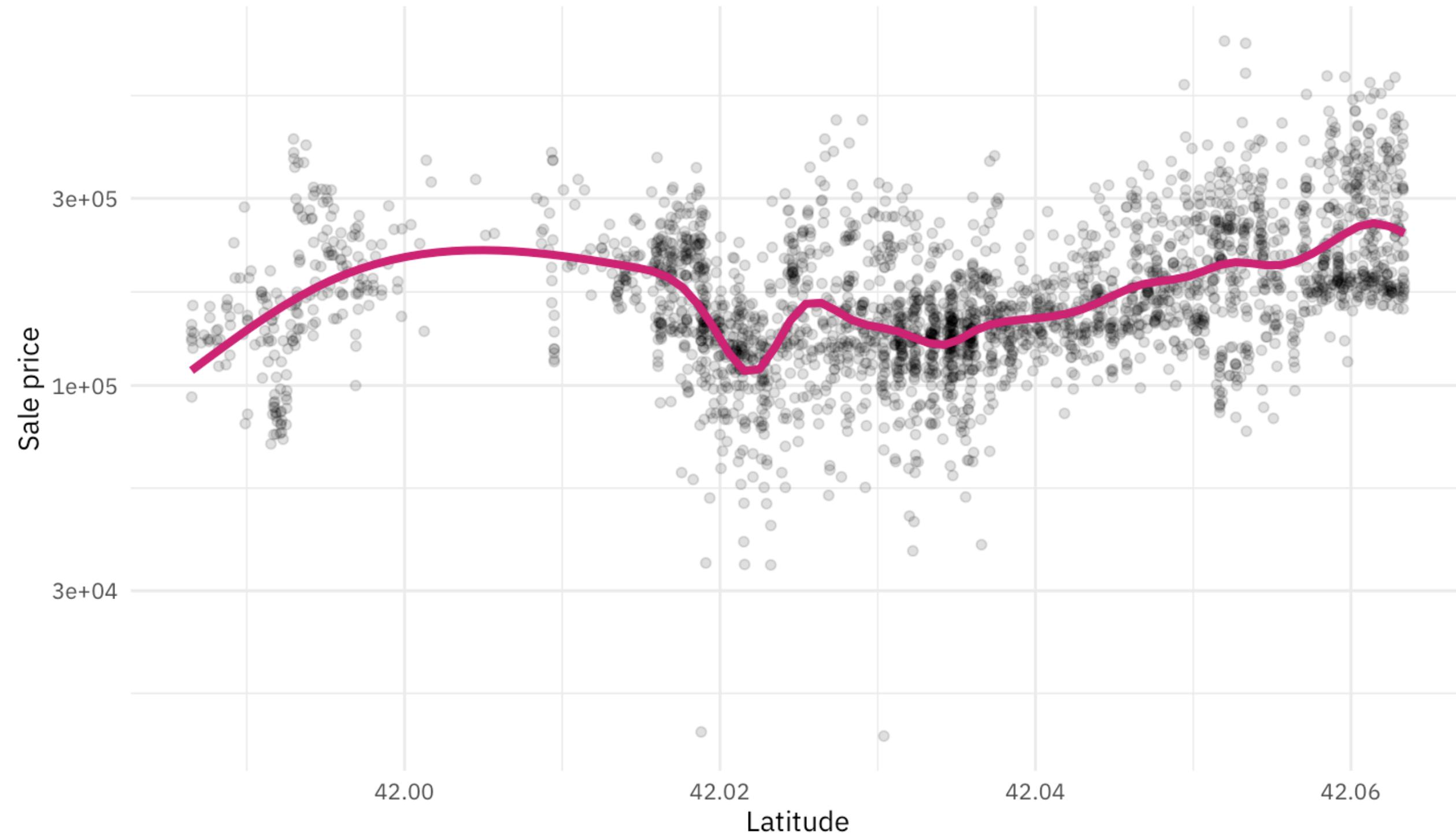
## 5 spline terms



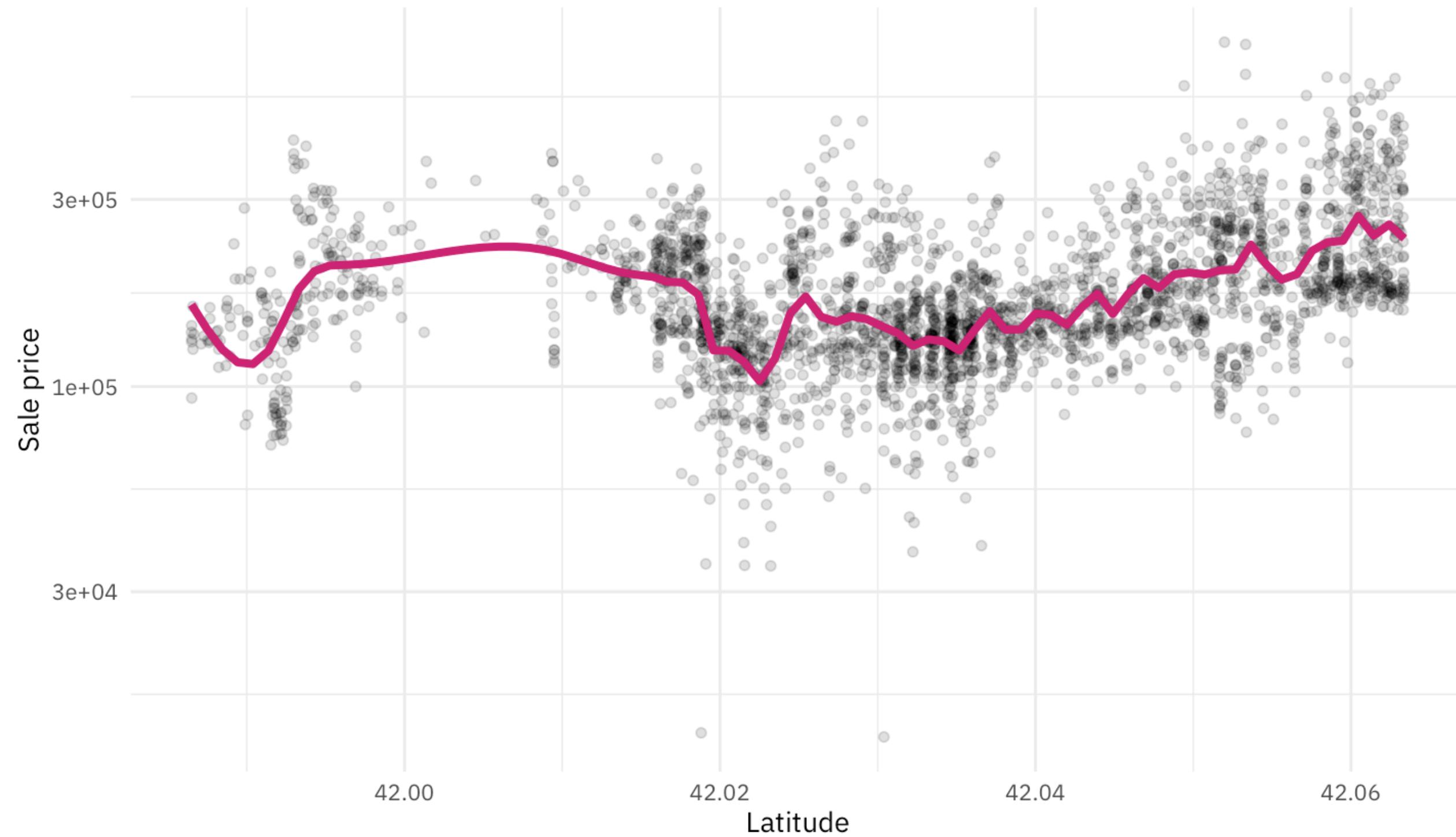
# 10 spline terms



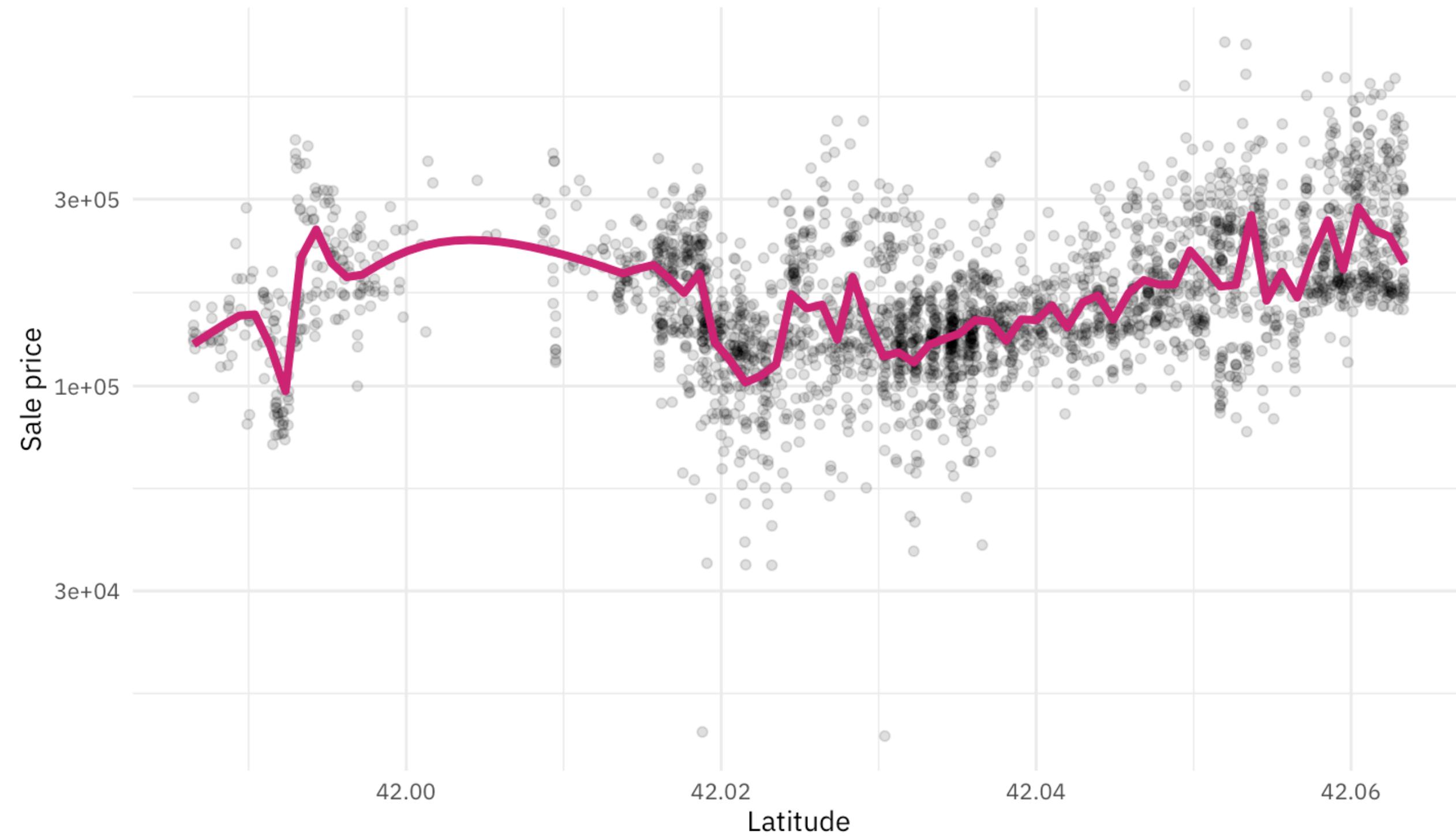
## 20 spline terms



## 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

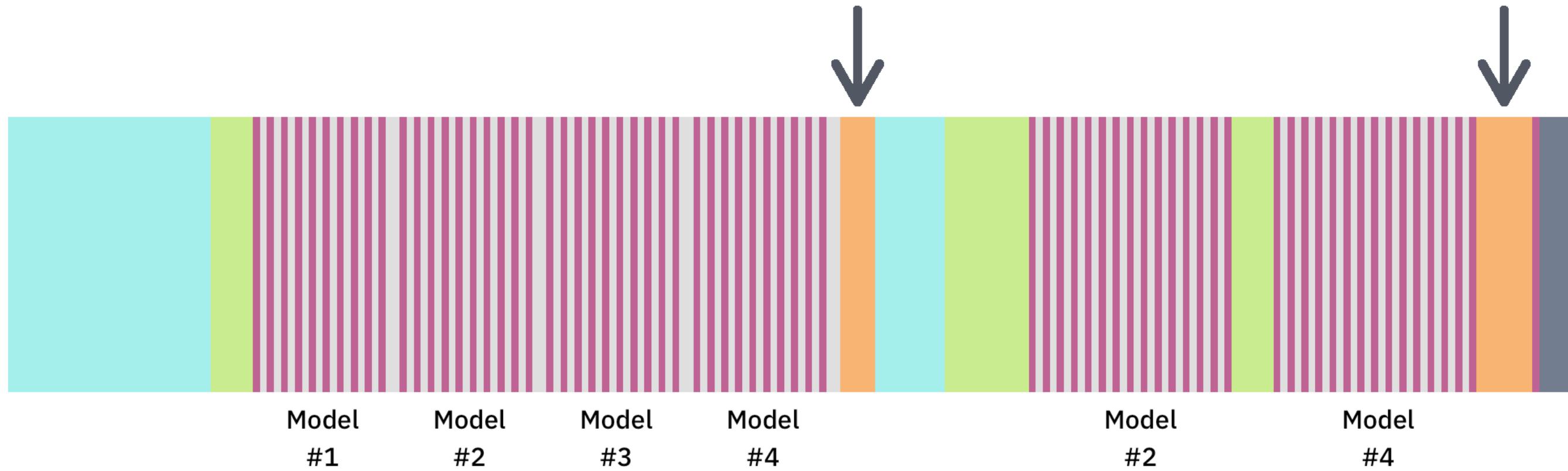
10:16 PM · Jul 6, 2016 · [YoruFukurou](#)

---

**321** Retweets and comments    **501** Likes

---





<span style="background-color: lightblue; border: 1px solid black; display: inline-block; width: 10px; height: 10px;"></span>	EDA	<span style="background-color: darkred; border: 1px solid black; display: inline-block; width: 10px; height: 10px;"></span>	Model Fit	<span style="background-color: orange; border: 1px solid black; display: inline-block; width: 10px; height: 10px;"></span>	Model Evaluation
<span style="background-color: lightgreen; border: 1px solid black; display: inline-block; width: 10px; height: 10px;"></span>	Feature Engineering	<span style="background-color: lightgray; border: 1px solid black; display: inline-block; width: 10px; height: 10px;"></span>	Model Tuning	<span style="background-color: darkgray; border: 1px solid black; display: inline-block; width: 10px; height: 10px;"></span>	Communication, deployment, etc.

# Why are these plots built?

# Visualization as validation tool

- Did this model "work"?

# Visualization as validation tool

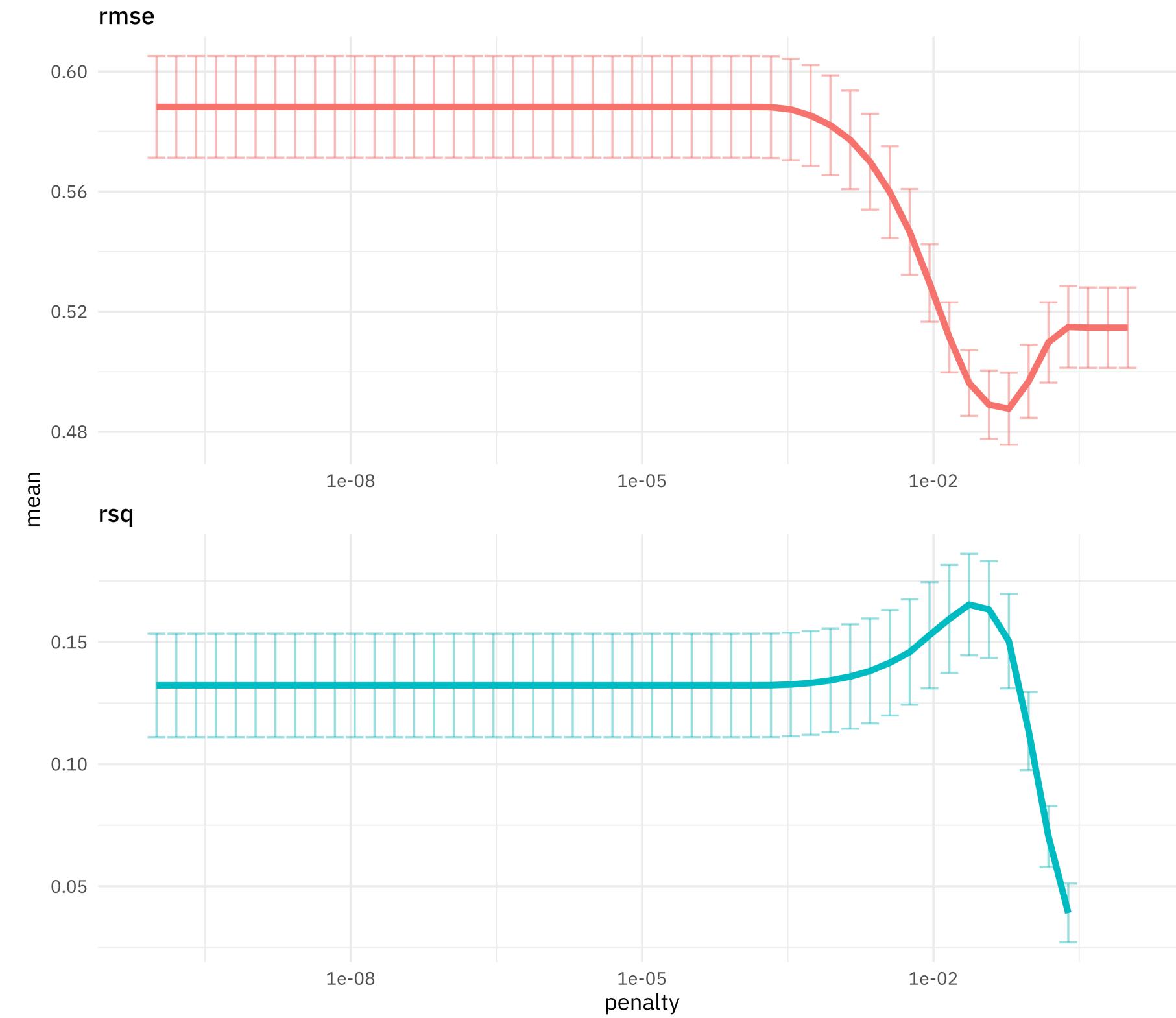
- Did this model "work"?
- Which model performed best?

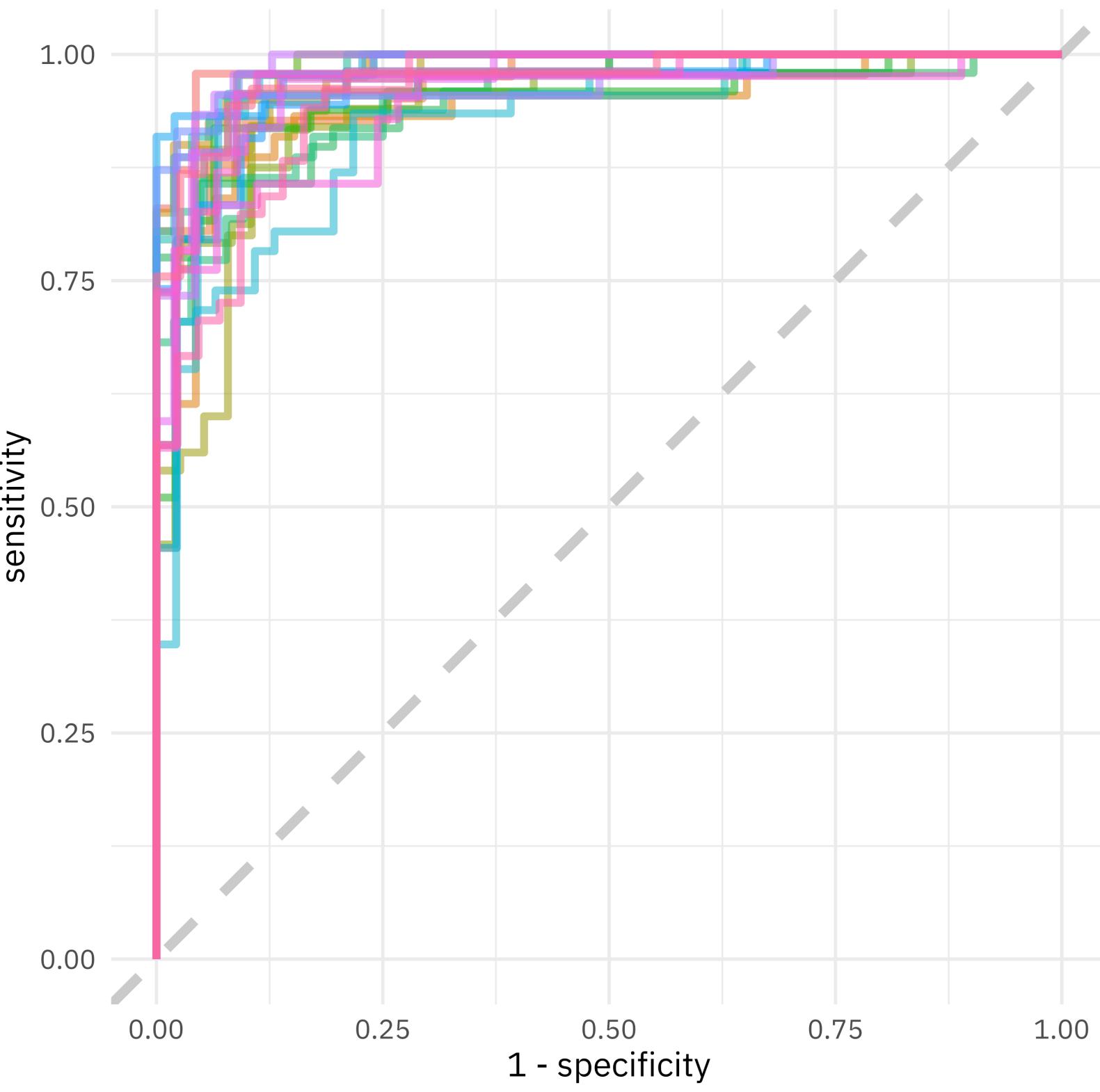
# Visualization as validation tool

- Did this model "work"?
- Which model performed best?
- For which kind of observations?

# Visualization as validation tool

- Did this model "work"?
- Which model performed best?
- For which kind of observations?
- Which model features are important?



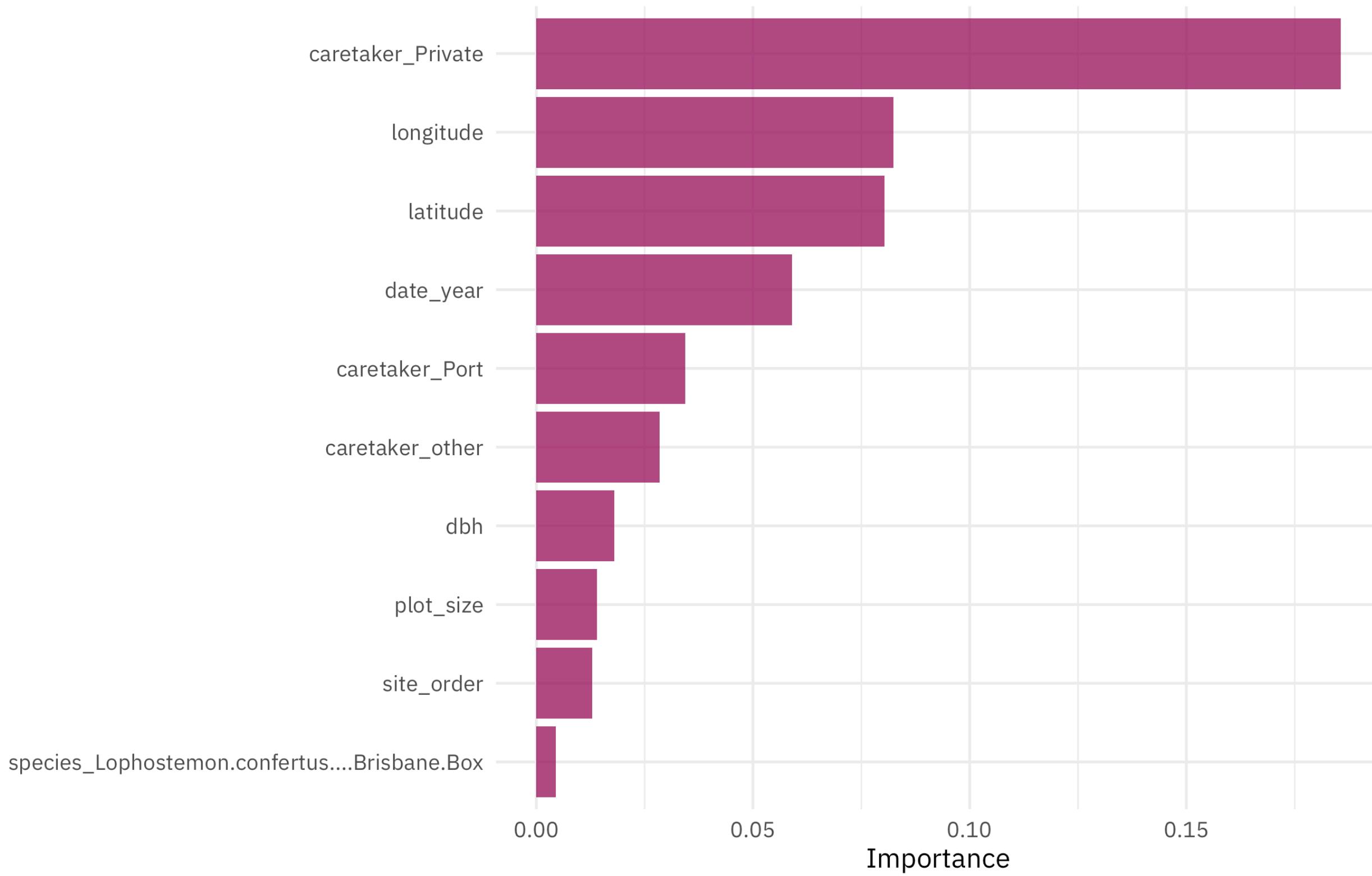


# automatic plotting methods

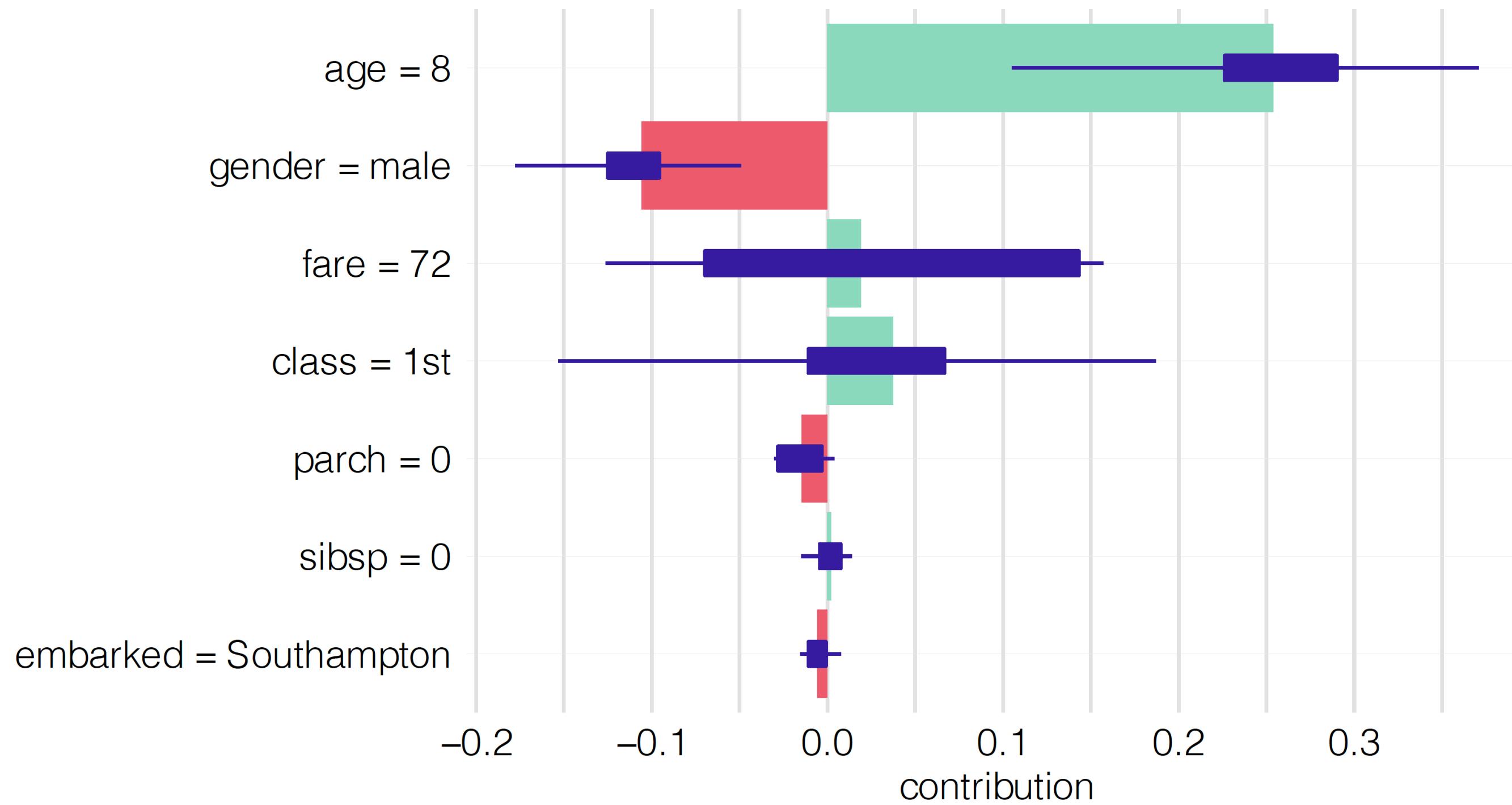


Volcano classification

# explaining models



## Average attributions for Johnny D



# Who are these plots for?

understand  
1  
gatua

understand  
models

# Thank you!

---

Julia Silge