

Model Fitting

Matthew McDonald

Fitting Models

How do you fit a linear model in R?

How many different ways can you think of?

- `lm` for linear model
- `glm` for generalized linear model (e.g. logistic regression)
- `glmnet` for regularized regression
- `keras` for regression using TensorFlow
- `stan` for Bayesian regression
- `spark` for large data sets

To specify a model using parsnip

- Choose a model
- Specify an engine
- Set the mode
- **fit** the model

To specify a model using parsnip

```
1 library(tidymodels)
2 linear_reg()
```

```
#> Linear Regression Model Specification (regression)
```

```
#>
```

```
#> Computational engine: lm
```

To specify a model using parsnip

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

```
#> Linear Regression Model Specification (regression)
```

```
#>
```

```
#> Computational engine: glmnet
```

To specify a model using parsnip

```
1 linear_reg() %>%  
2   set_engine("stan")
```

```
#> Linear Regression Model Specification (regression)
```

```
#>
```

```
#> Computational engine: stan
```

To specify a model using parsnip

```
1 decision_tree()
```

```
#> Decision Tree Model Specification (unknown mode)
```

```
#>
```

```
#> Computational engine: rpart
```

To specify a model using parsnip

```
1 decision_tree() %>%  
2   set_mode("regression")
```

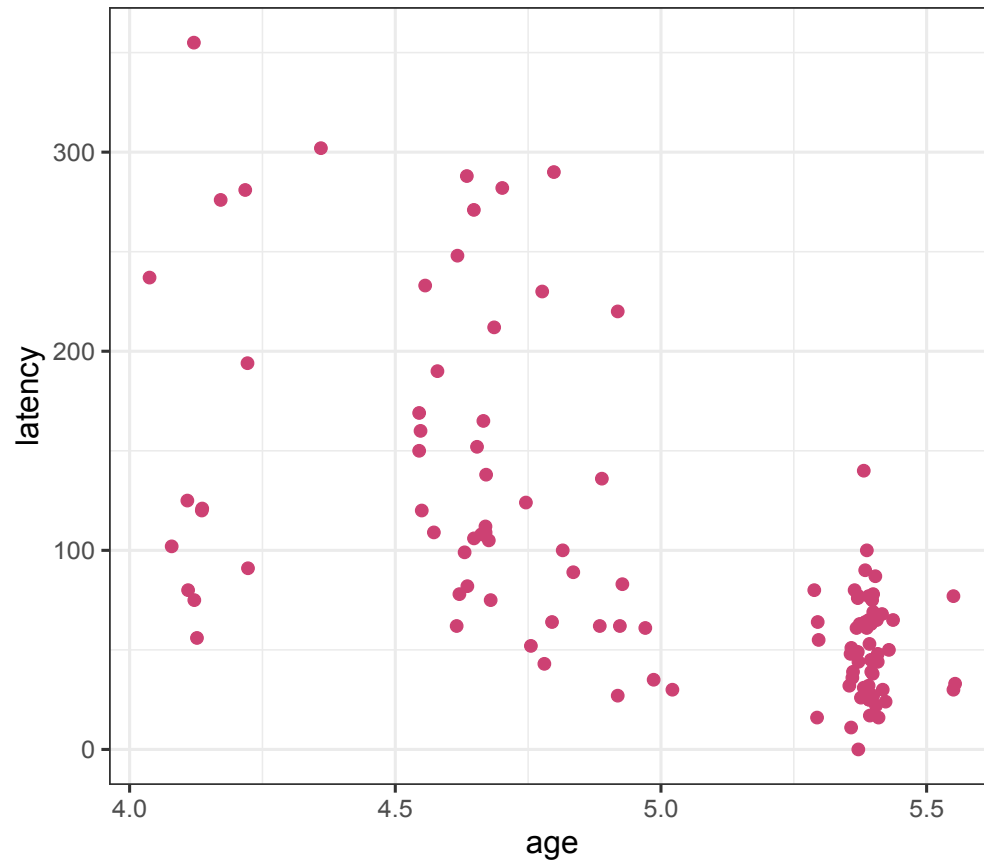
```
#> Decision Tree Model Specification (regression)
```

```
#>
```

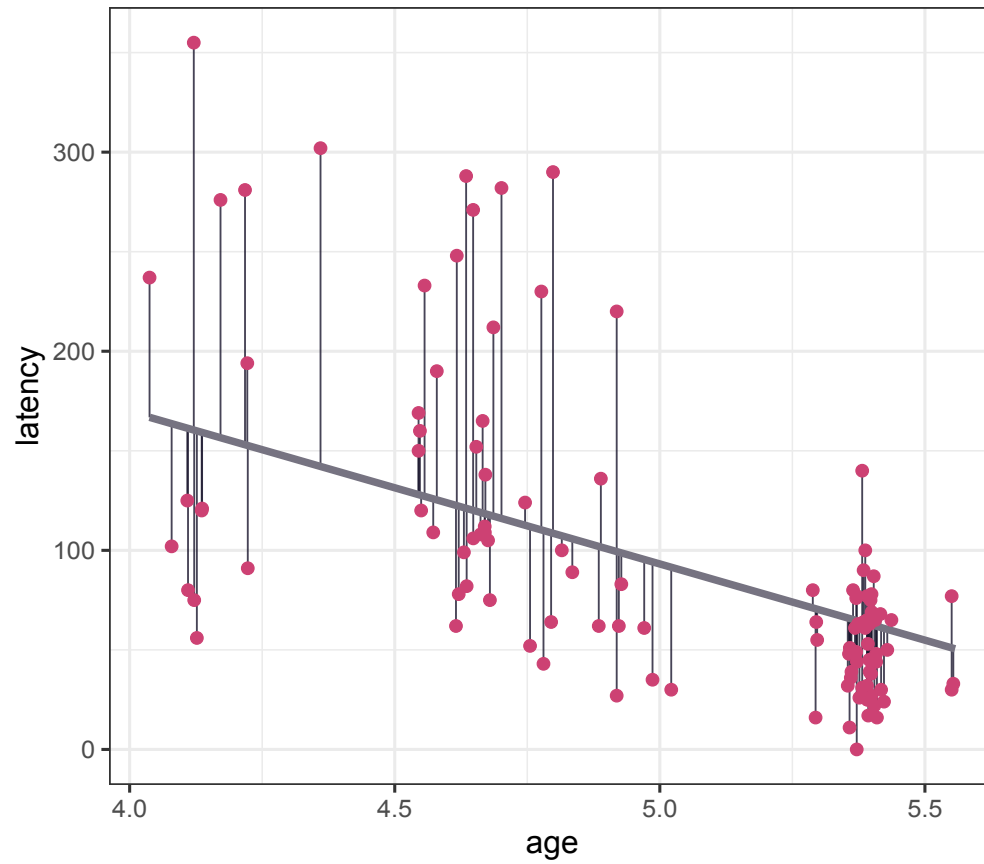
```
#> Computational engine: rpart
```

All available models are listed at <https://www.tidymodels.org/find/p>

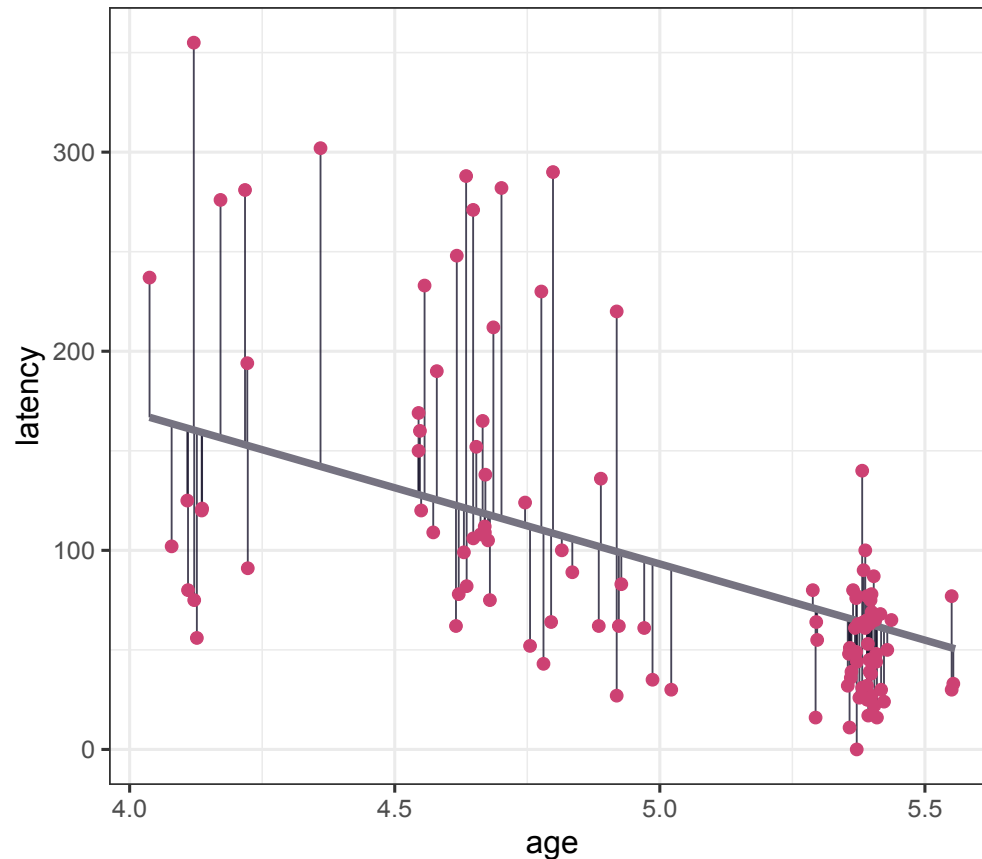
Linear regression



Linear regression



Linear regression



- Outcome modeled as linear combination of predictors:

$$\text{latency} = \beta_0 + \beta_1 \cdot \text{age} +$$

- Find a line that minimizes the mean squared error (MSE)

Comparing base R vs tidymodels

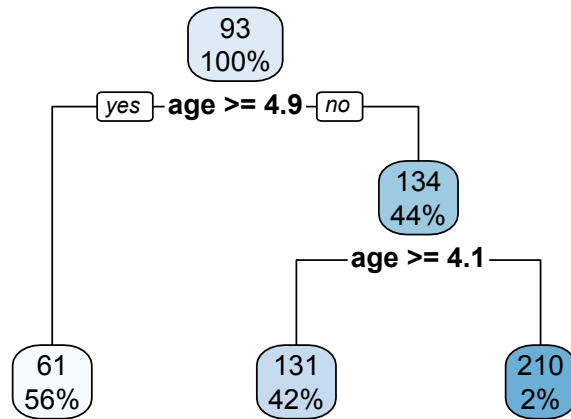
```
1 lm(latency ~ age, data = frog_train) %>%  
2 broom::tidy()
```

```
#> # A tibble: 2 × 5  
#>   term          estimate std.error statistic  p.value  
#>   <chr>          <dbl>    <dbl>    <dbl>    <dbl>  
#> 1 (Intercept)    476.      31.9     14.9 2.76e-41  
#> 2 age           -76.5      6.34    -12.1 2.98e-29
```

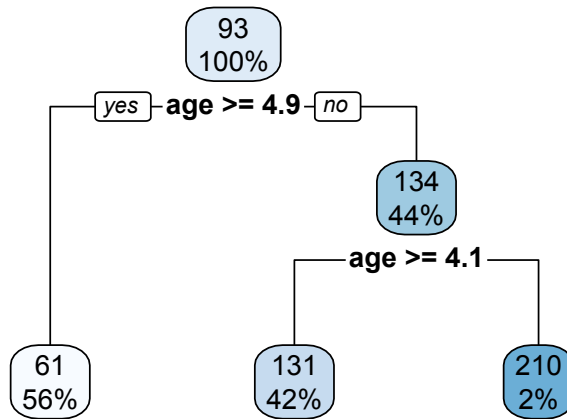
```
1 linear_reg() %>%  
2   set_engine('lm') %>%  
3   set_mode('regression') %>%  
4   fit(latency ~ age, data = frog_train) %>%  
5   broom::tidy()
```

```
#> # A tibble: 2 × 5  
#>   term          estimate std.error statistic  p.value  
#>   <chr>          <dbl>    <dbl>    <dbl>    <dbl>  
#> 1 (Intercept)    476.      31.9     14.9 2.76e-41  
#> 2 age           -76.5      6.34    -12.1 2.98e-29
```

Decision trees

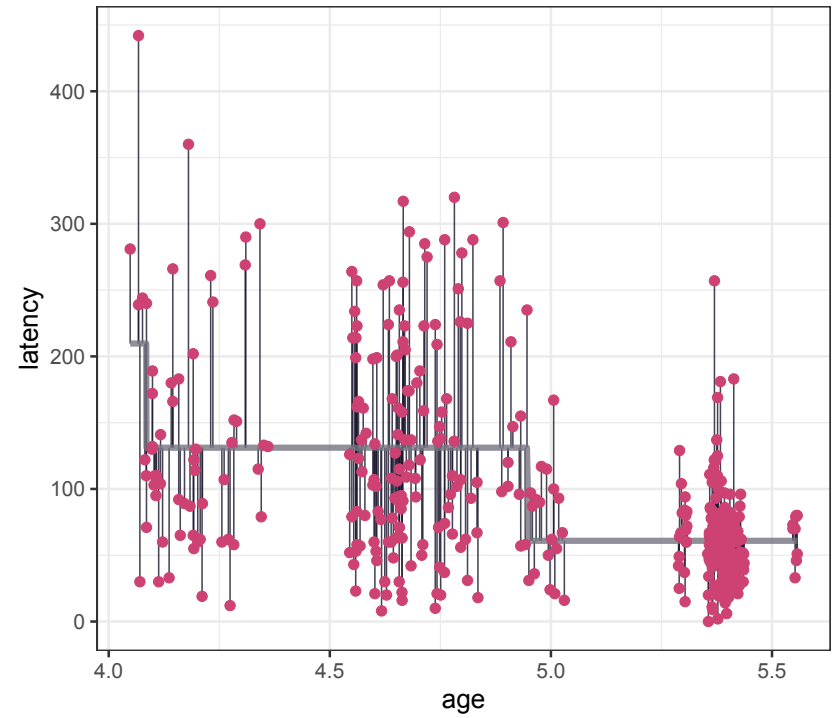
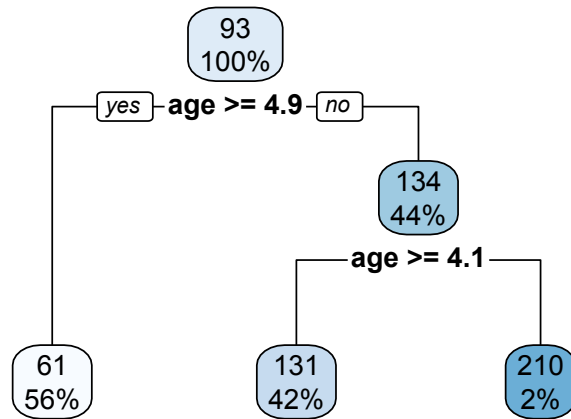


Decision trees



- Series of splits or if/then statements based on predictors
- First the tree *grows* until some condition is met (maximum depth, no more data)
- Then the tree is *pruned* to reduce its complexity

Decision trees



Comparing base R vs tidymodels

```
1 rpart(latency ~ age,  
2       data=frog_train,  
3       control=rpart.control(cp=0.015))
```

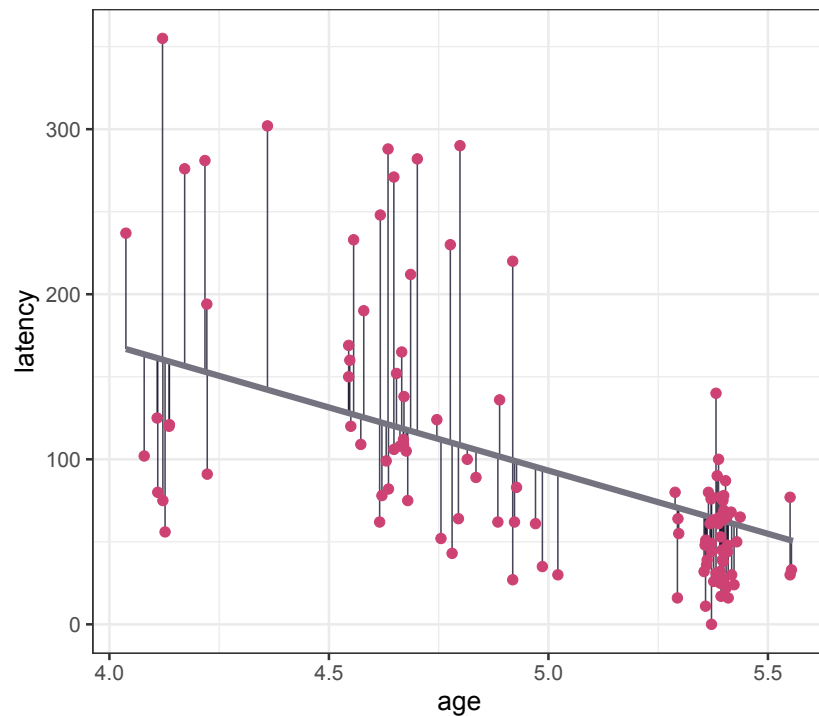
```
#> n= 456  
#>  
#> node), split, n, deviance, yval  
#>      * denotes terminal node  
#>  
#> 1) root 456 2197966.0  92.90351  
#>   2) age>=4.947975 256  252347.4  60.89844 *  
#>   3) age< 4.947975 200 1347741.0 133.87000  
#>     6) age>=4.085127 193 1195008.0 131.11920 *  
#>     7) age< 4.085127 7  111005.4 209.71430 *
```

```
1 decision_tree(cost_complexity = 0.015) %>%  
2   set_engine('rpart') %>%  
3   set_mode('regression') %>%  
4   fit(latency ~ age, data = frog_train)
```

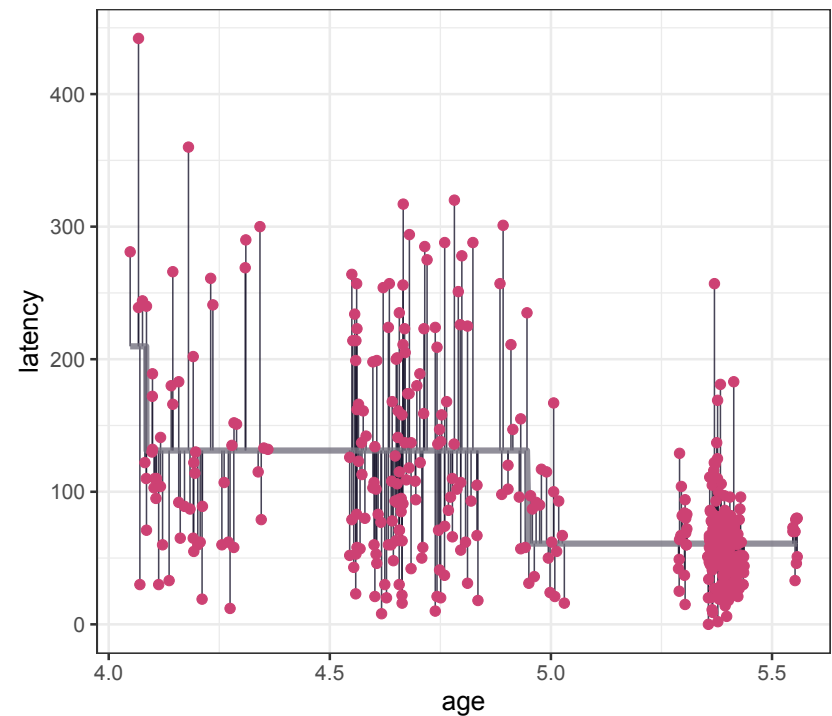
```
#> parsnip model object  
#>  
#> n= 456  
#>  
#> node), split, n, deviance, yval  
#>      * denotes terminal node  
#>  
#> 1) root 456 2197966.0  92.90351  
#>   2) age>=4.947975 256  252347.4  60.89844 *  
#>   3) age< 4.947975 200 1347741.0 133.87000  
#>     6) age>=4.085127 193 1195008.0 131.11920 *  
#>     7) age< 4.085127 7  111005.4 209.71430 *
```


All models are wrong, but some are useful!

Linear regression

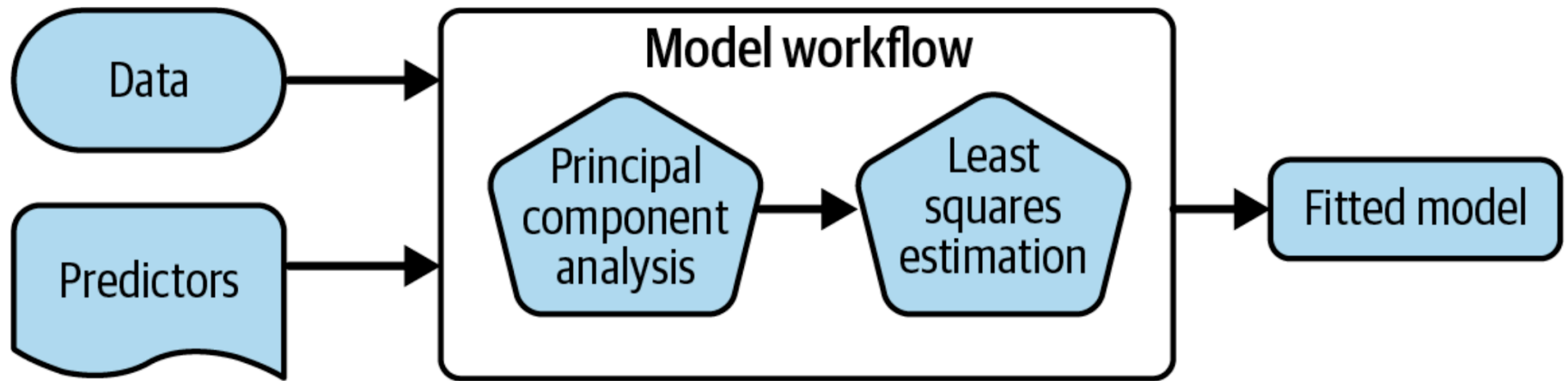


Decision trees

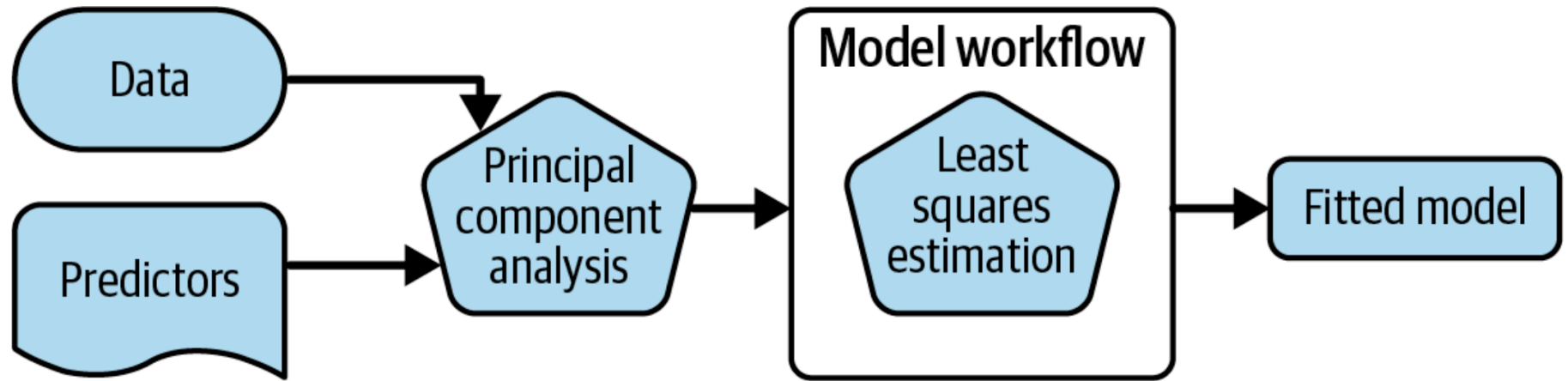


A model workflow

Workflows bind preprocessors and models



What is wrong with this?



Why a `workflow()`?

- Workflows handle new data better than base R tools in terms of new factor levels
- You can use other preprocessors besides formulas (more on feature engineering tomorrow!)
- They can help organize your work when working with multiple models
- Most importantly, a workflow captures the entire modeling process: `fit()` and `predict()` apply to the preprocessing steps in addition to the actual model fit

A model workflow

```
1 tree_spec <-  
2   decision_tree() %>%  
3   set_mode("regression")  
4  
5 tree_spec %>%  
6   fit(latency ~ ., data = frog_train)
```

```
#> parsnip model object
```

```
#>
```

```
#> n= 456
```

```
#>
```

```
#> node), split, n, deviance, yval
```

```
#>   * denotes terminal node
```

```
#>
```

```
#> 1) root 456 2197966.00  92.90351
```

```
#>   2) age>=4.947975 256  252347.40  60.89844
```

```
#>     4) treatment=control 131  91424.06  48.42748 *
```

```
#>     5) treatment=gentamicin 125  119197.90  73.96800 *
```

```
#>   3) age< 4.947975 200 1347741.00 133.87000
```

```
#>     6) treatment=control 140  986790.70 118.25710
```

```
#>     12) reflex=mid,full 129  754363.70 111.56590 *
```

```
#>     13) reflex=low 11  158918.20 196.72730 *
```

```
#>   7) treatment=gentamicin 60  247194.60 170.30000
```

```
#>     14) age< 4.664439 30  102190.20 147.83330
```

```
#>     28) age>=4.566638 22  53953.86 129.77270 *
```

```
#>     29) age< 4.566638 8  21326.00 197.50000 *
```

```
#>   15) age>=4.664439 30  114719.40 192.76670 *
```

A model workflow

```
1 tree_spec <-  
2   decision_tree() %>%  
3   set_mode("regression")  
4  
5 workflow() %>%  
6   add_formula(latency ~ .) %>%  
7   add_model(tree_spec) %>%  
8   fit(data = frog_train)
```

```
#> == Workflow [trained] ==  
#> Preprocessor: Formula  
#> Model: decision_tree()  
#>  
#> — Preprocessor —  
#> latency ~ .  
#>  
#> — Model —  
#> n= 456  
#>  
#> node), split, n, deviance, yval  
#> * denotes terminal node  
#>  
#> 1) root 456 2197966.00 92.90351  
#> 2) age>=4.947975 256 252347.40 60.89844  
#> 4) treatment=control 131 91424.06 48.42748 *  
#> 5) treatment=gentamicin 125 119197.90 73.96800 *  
#> 3) age< 4.947975 200 1347741.00 133.87000  
#> 6) treatment=control 140 986790.70 118.25710  
#> 12) reflex=mid,full 129 754363.70 111.56590 *  
#> 13) reflex=low 11 158918.20 196.72730 *  
#> 7) treatment=gentamicin 60 247194.60 170.30000  
#> 14) age< 4.664439 30 102190.20 147.82320
```

A model workflow

```
1 tree_spec <-  
2   decision_tree() %>%  
3   set_mode("regression")  
4  
5 workflow(latency ~ ., tree_spec) %>%  
6   fit(data = frog_train)
```

```
#> == Workflow [trained] ==  
#> Preprocessor: Formula  
#> Model: decision_tree()  
#>  
#> — Preprocessor —  
#> latency ~ .  
#>  
#> — Model —  
#> n= 456  
#>  
#> node), split, n, deviance, yval  
#>   * denotes terminal node  
#>  
#> 1) root 456 2197966.00  92.90351  
#>   2) age>=4.947975 256  252347.40  60.89844  
#>     4) treatment=control 131  91424.06  48.42748 *  
#>     5) treatment=gentamicin 125  119197.90  73.96800 *  
#>   3) age< 4.947975 200 1347741.00 133.87000  
#>     6) treatment=control 140  986790.70 118.25710  
#>       12) reflex=mid,full 129  754363.70 111.56590 *  
#>       13) reflex=low 11  158918.20 196.72730 *  
#>     7) treatment=gentamicin 60  247194.60 170.30000  
#>     14) age< 4.664439 30  102190.20 147.82220
```


Predict with your model

How do you use your new `tree_fit` model?

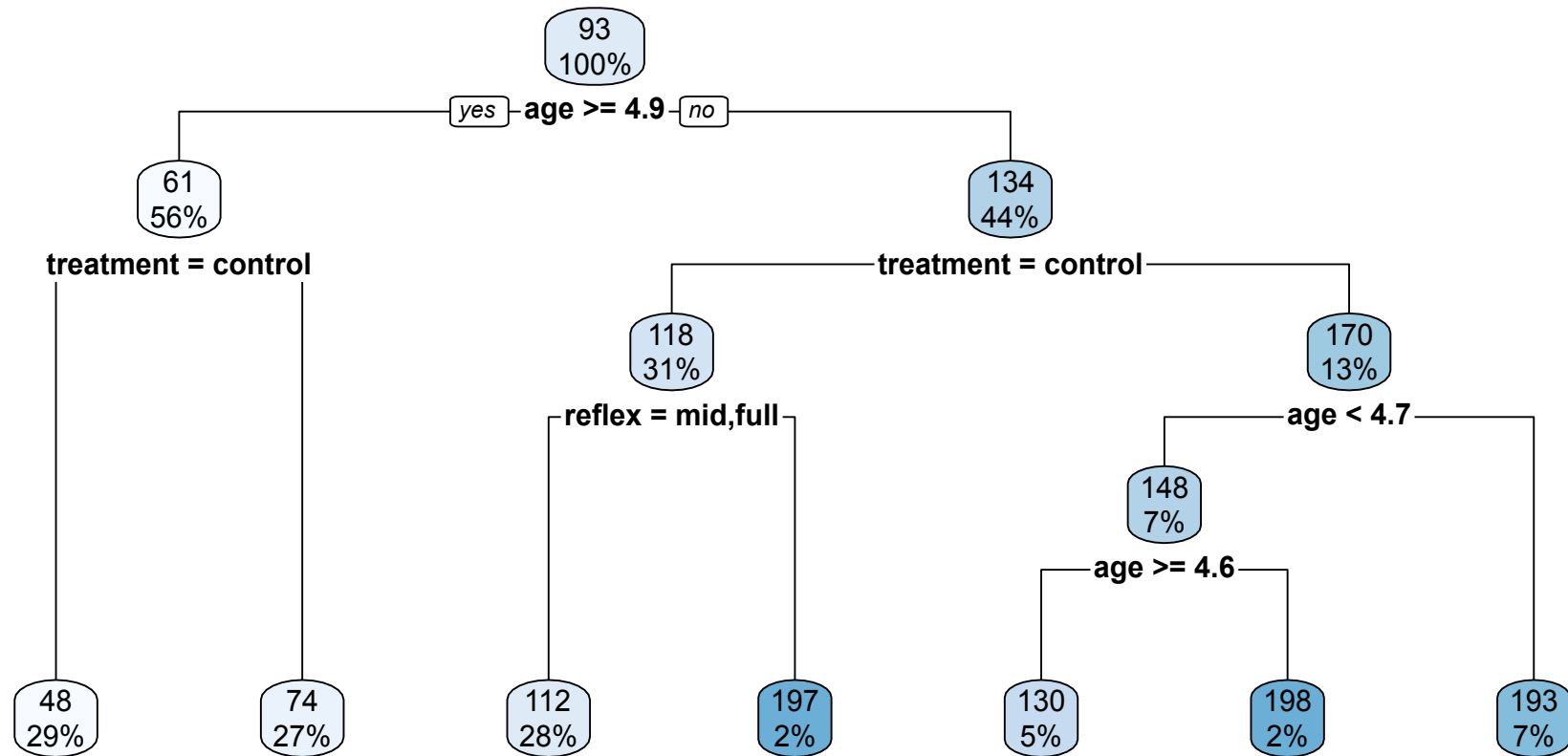
```
1 tree_spec <-  
2   decision_tree() %>%  
3   set_mode("regression")  
4  
5 tree_fit <-  
6   workflow(latency ~ ., tree_spec) %>%  
7   fit(data = frog_train)
```

The tidymodels prediction guarantee!

- The predictions will always be inside a **tibble**
- The column names and types are **unsurprising** and **predictable**
- The number of rows in **new_data** and the output **are the same**

Understand your model

How do you understand your new `tree_fit` model?



Understand your model

How do you understand your new `tree_fit` model?

```
1 library(rpart.plot)
2 tree_fit %>%
3   extract_fit_engine() %>%
4   rpart.plot(roundint = FALSE)
```

You can `extract_*`() several components of your fitted workflow.

Understand your model

How do you **understand** your new `tree_fit` model?

You can use your fitted workflow for model and/or prediction explanations:

- overall variable importance, such as with the `vip` package
- flexible model explainers, such as with the `DALEXtra` package

Learn more at <https://www.tmwr.org/explain.html>