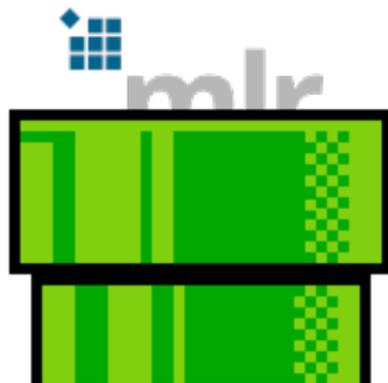


MLR3PIPELINES: MACHINE LEARNING PIPELINES AS GRAPHS

<https://bit.ly/32jDQOG>

Martin Binder, Florian Pfisterer, Michel Lang, **Bernd Bischl**

Computational Statistics, LMU



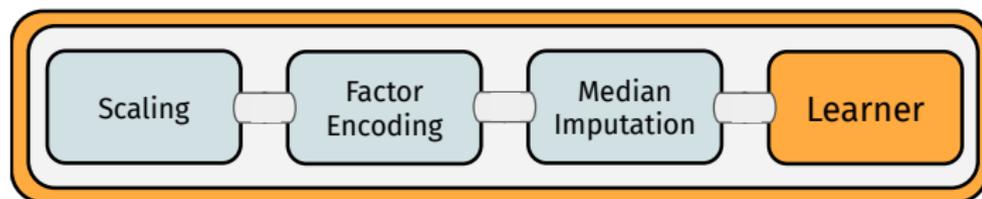
MLR3 PIPELINES

MACHINE LEARNING WORKFLOWS:

- **Preprocessing:** Feature extraction, feature selection, missing data imputation, . . .
- **Ensemble methods:** Model averaging, model stacking
- **mlr3:** modular model fitting

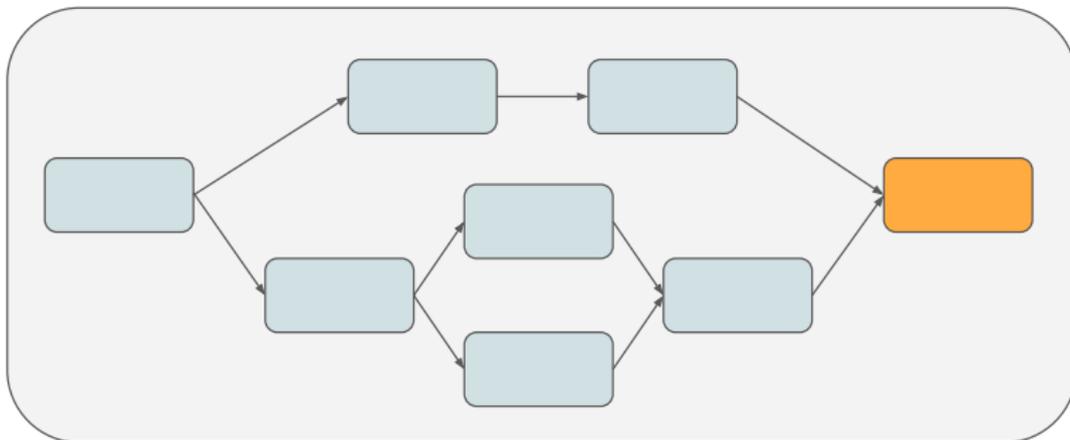
⇒ **mlr3pipelines:** modular ML workflows

(replaces mlr2's mlrCPO and most “wrappers”)



MACHINE LEARNING WORKFLOWS

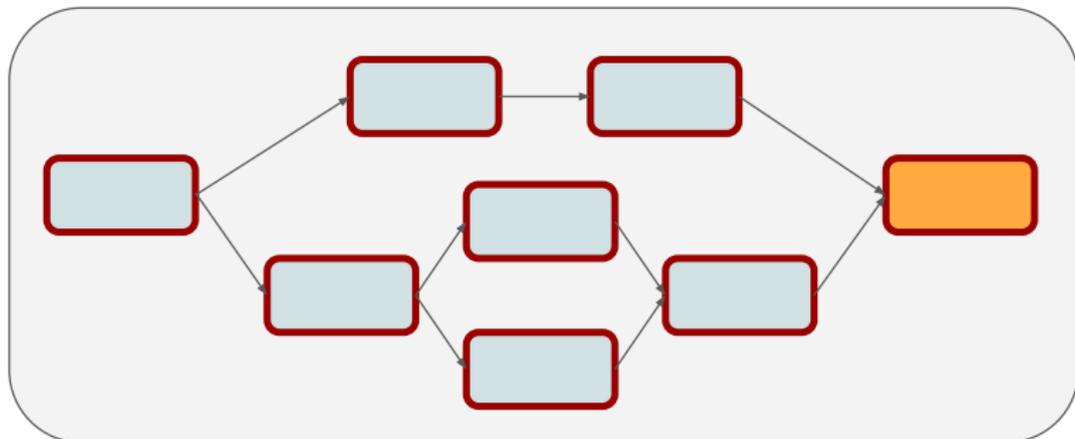
– what do they look like?



MACHINE LEARNING WORKFLOWS

– what do they look like?

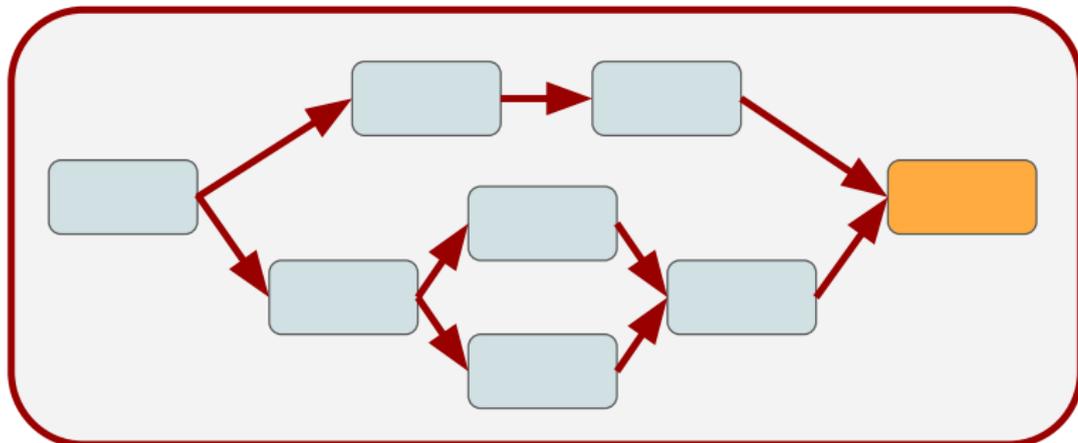
- **Building blocks:** *what is happening?* → PipeOp



MACHINE LEARNING WORKFLOWS

– what do they look like?

- **Building blocks:** *what is happening?* → PipeOp
- **Structure:** In what *sequence* is it happening? → Graph



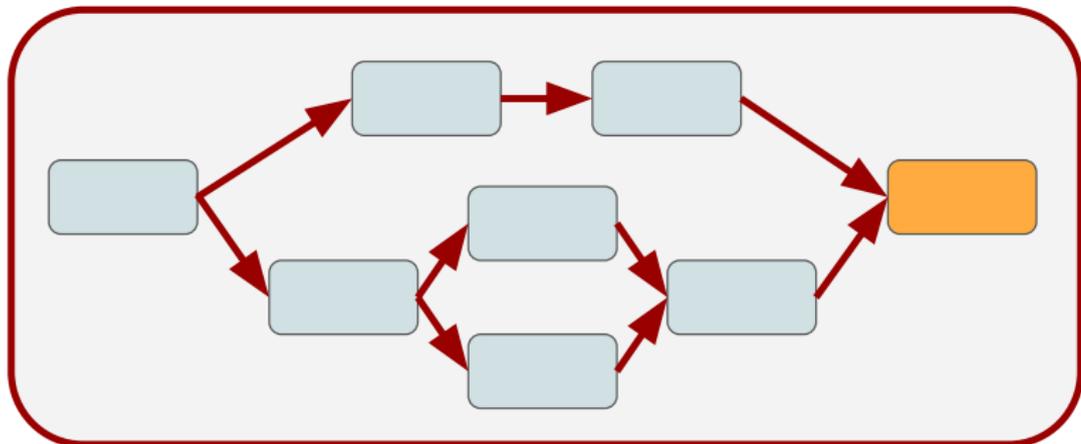
MACHINE LEARNING WORKFLOWS

– what do they look like?

- **Building blocks:** *what is happening?* → PipeOp

- **Structure:** In what *sequence* is it happening? → Graph

⇒ Graph: PipeOps as **nodes** with **edges** (data flow) between them

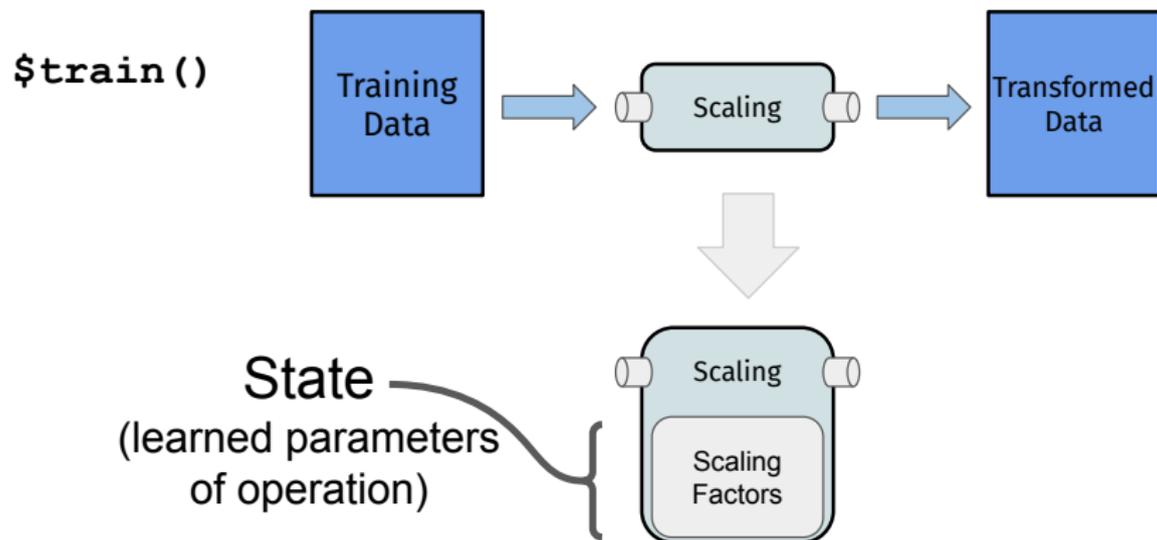


PIPEOPS

THE BUILDING BLOCKS

PIPEOP: SINGLE UNIT OF DATA OPERATION

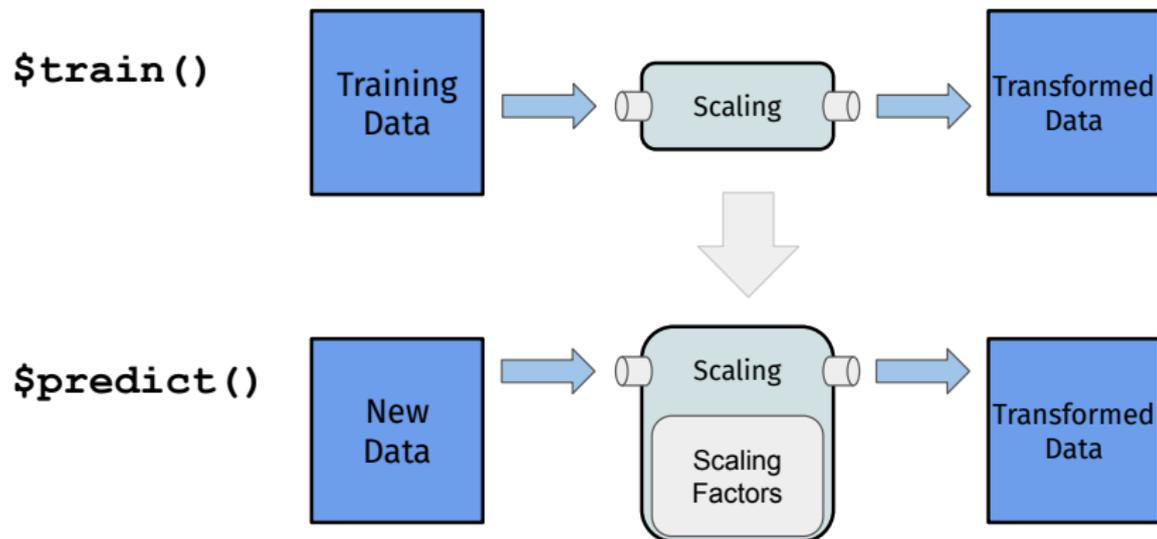
- `$train()`: process data and create `$state`



THE BUILDING BLOCKS

PIPEOP: SINGLE UNIT OF DATA OPERATION

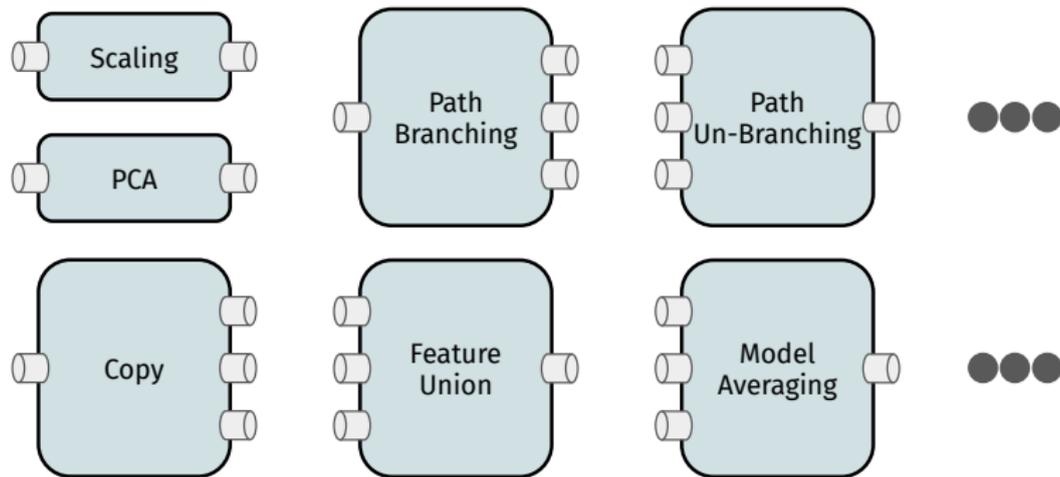
- `$train()`: process data and create `$state`
- `$predict()`: process data depending on the `$state`



THE BUILDING BLOCKS

PIPEOP: SINGLE UNIT OF DATA OPERATION

- `$train()`: process data and create `$state`
- `$predict()`: process data depending on the `$state`
- Multiple inputs or multiple outputs



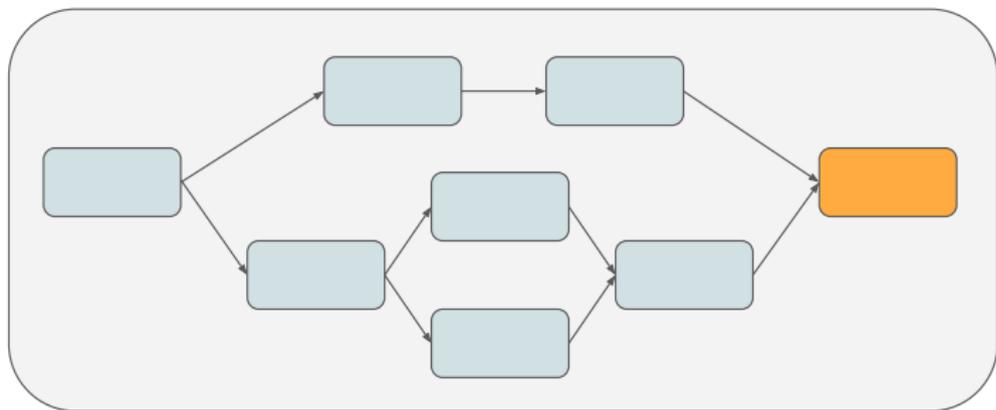
PIPEOPS SO FAR AND PLANNED

- Simple preprocessing operations: `scale`, `pca`, `apply`, `mutate`
- Missing value imputation: `impute`
- Feature selection and filtering: `select`, `filter`
- Categorical data encoding: `encode`
- Undersampling / subsampling: `balancesample`, `subsample`, `chunk`
- Learners: `learner`, `learner_cv`
- Ensemble methods on Predictions: `majorityvote`, `modelavg`
- Simultaneous and alternative branching: `copy`, `branch`, `unbranch`
- Combination of data: `featureunion`
- Backup prediction: `backuplearner`
- Text processing (*planned*)
- Time series and spatio-temporal data (*planned*)
- Multi-output and ordinal targets (*planned*)
- Outlier detection (*planned*)
- Hurdle models (*planned*)
- ...

GRAPH OPERATIONS

THE STRUCTURE

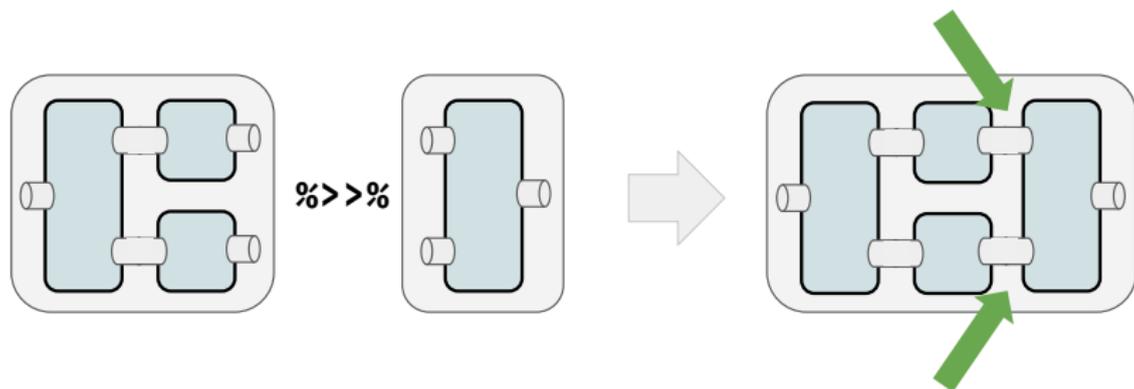
GRAPH OPERATIONS



THE STRUCTURE

GRAPH OPERATIONS

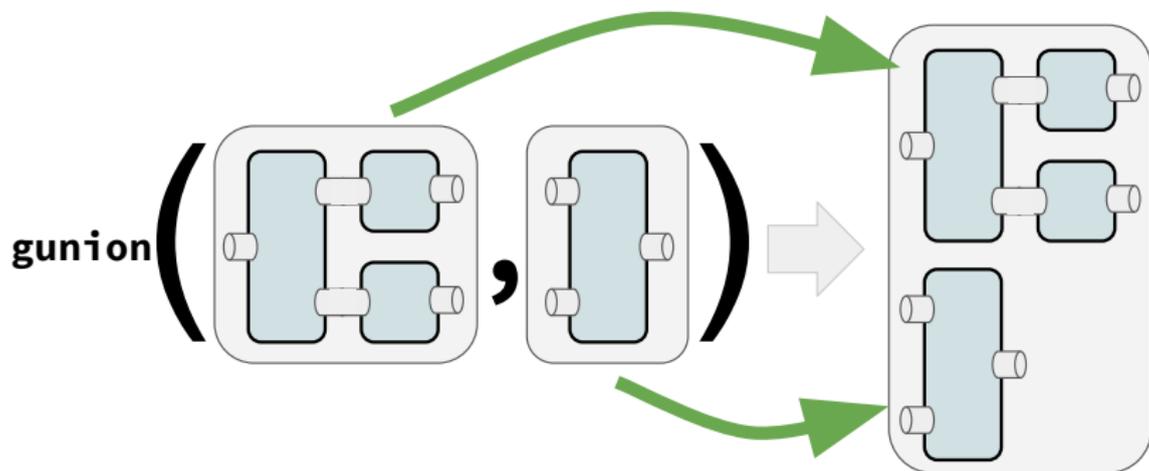
- The `%>>%`-operator concatenates Graphs and PipeOps



THE STRUCTURE

GRAPH OPERATIONS

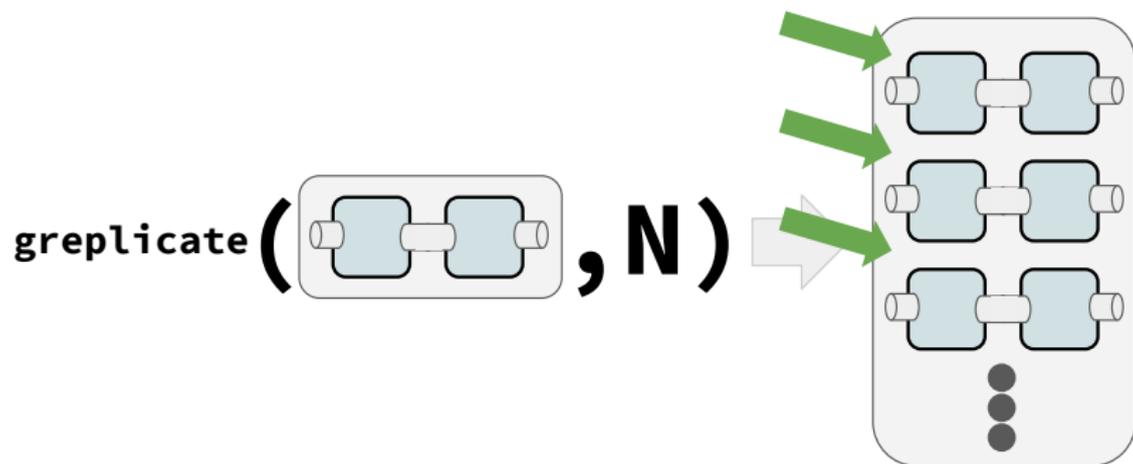
- The `%>>%`-operator concatenates Graphs and PipeOps
- The `gunion()`-function unites Graphs and PipeOps



THE STRUCTURE

GRAPH OPERATIONS

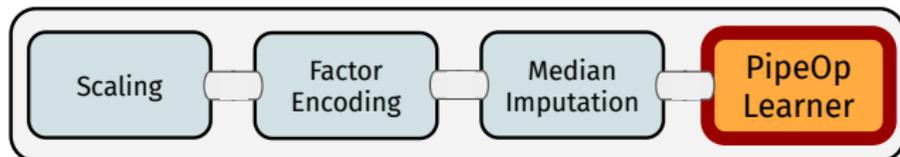
- The `%>>%`-operator concatenates Graphs and PipeOps
- The `gunion()`-function unites Graphs and PipeOps
- The `grePLICATE()`-function unites copies of Graphs and PipeOps



LEARNERS AND GRAPHS

PIPEOPLEARNER

- **Learner** as a PipeOp
- Fits a model, output is Prediction



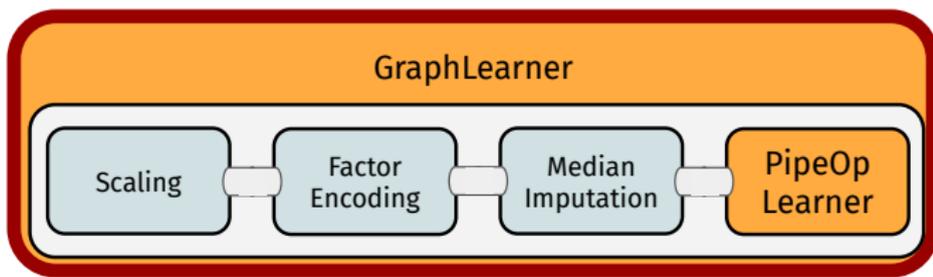
LEARNERS AND GRAPHS

PIPEOPLEARNER

- **Learner as a PipeOp**
- Fits a model, output is Prediction

GRAPHLEARNER

- **Graph as a Learner**
- All benefits of `mlr3`: **resampling, tuning, nested resampling, ...**

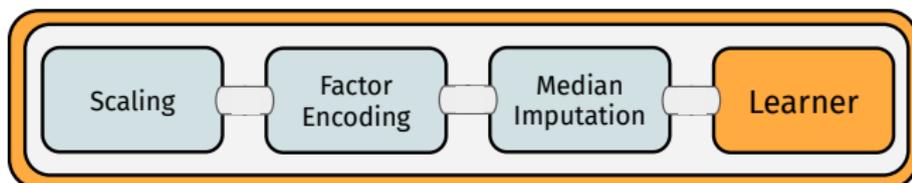


LINEAR PIPELINES

MLR3 PIPELINES IN ACTION

LINEAR PREPROCESSING PIPELINE

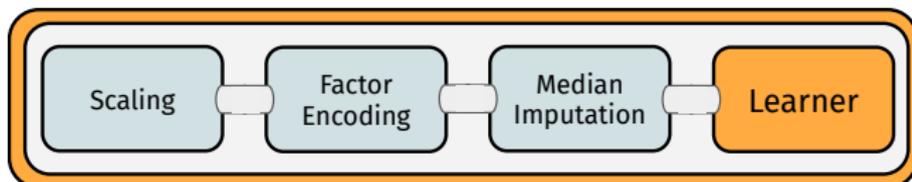
```
graph_pp = mlr_pipeops$get("scale") %>>%  
  mlr_pipeops$get("encode") %>>%  
  mlr_pipeops$get("impute") %>>%  
  mlr_pipeops$get("learner",  
    learner = mlr_learners$get("classif.rpart"))
```



MLR3 PIPELINES IN ACTION

LINEAR PREPROCESSING PIPELINE

```
graph_pp = "scale" %>>% "encode" %>>% "impute" %>>%  
  mlr_pipeops$get("learner",  
    learner = mlr_learners$get("classif.rpart"))
```

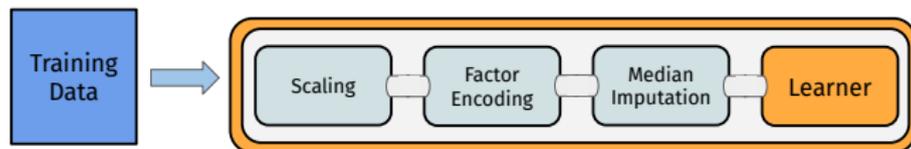


MLR3 PIPELINES IN ACTION

LINEAR PREPROCESSING PIPELINE

- `train()`ing: Data propagates and creates `$states`

```
glrn = GraphLearner$new(graph_pp)  
glrn$train(task)
```

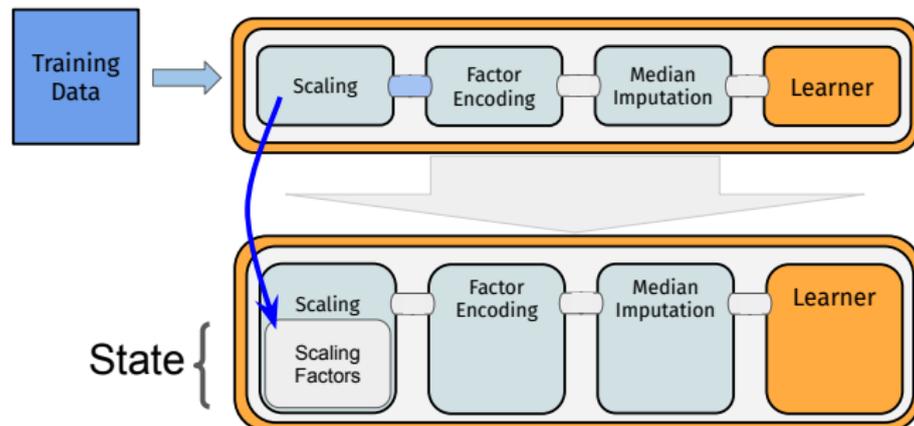


MLR3 PIPELINES IN ACTION

LINEAR PREPROCESSING PIPELINE

- `train()`ing: Data propagates and creates `$states`

```
glrn = GraphLearner$new(graph_pp)  
glrn$train(task)
```

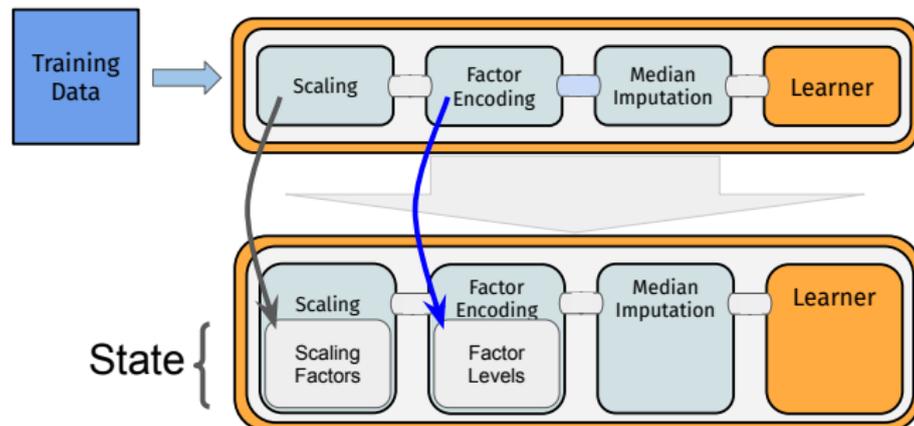


MLR3 PIPELINES IN ACTION

LINEAR PREPROCESSING PIPELINE

- `train()`ing: Data propagates and creates `$states`

```
glrn = GraphLearner$new(graph_pp)  
glnr$train(task)
```

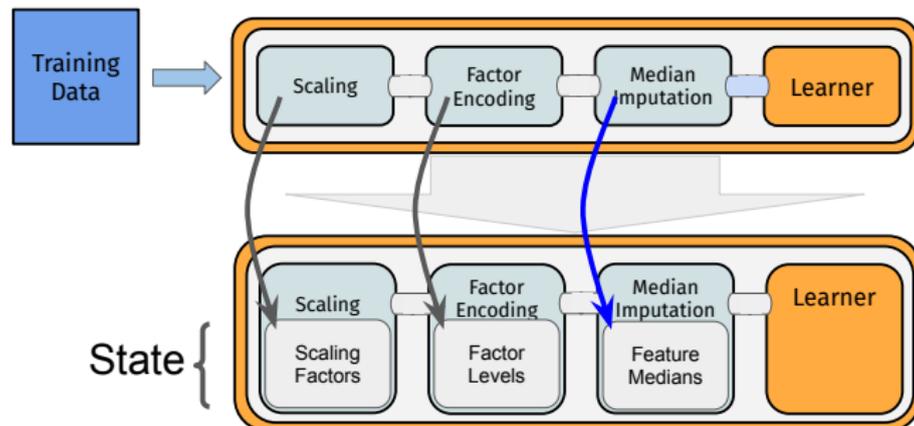


MLR3 PIPELINES IN ACTION

LINEAR PREPROCESSING PIPELINE

- `train()`ing: Data propagates and creates `$states`

```
glrn = GraphLearner$new(graph_pp)  
glrn$train(task)
```

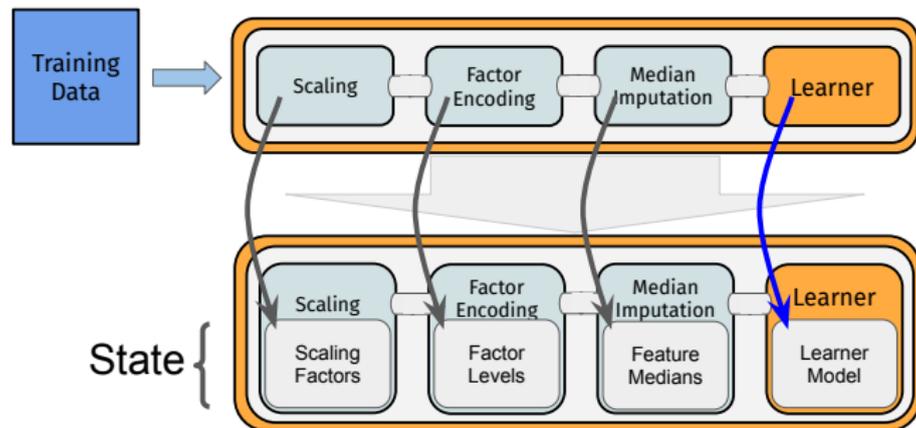


MLR3 PIPELINES IN ACTION

LINEAR PREPROCESSING PIPELINE

- `train()`ing: Data propagates and creates `$states`

```
glrn = GraphLearner$new(graph_pp)  
glrn$train(task)
```

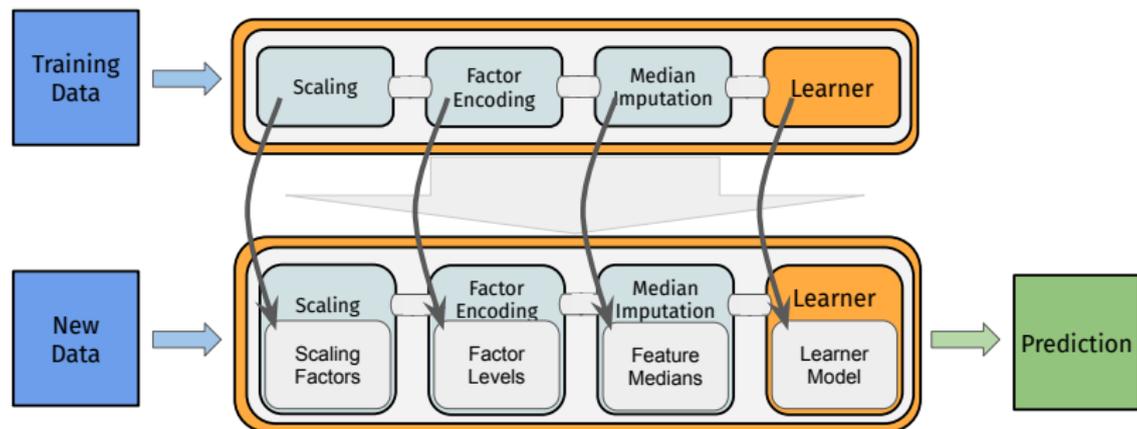


MLR3 PIPELINES IN ACTION

LINEAR PREPROCESSING PIPELINE

- `train()`ing: Data propagates and creates `$states`
- `predict()`ion: Data propagates, uses `$states`

```
glrn$predict(task)
```



MLR3 PIPELINES IN ACTION

LINEAR PREPROCESSING PIPELINE

- Setting / retrieving parameters: `$param_set`

```
graph_pp$pipeops$impute$param_set$values$method_num = "mean"
```

MLR3 PIPELINES IN ACTION

LINEAR PREPROCESSING PIPELINE

- Setting / retrieving parameters: `$param_set`

```
graph_pp$pipeops$impute$param_set$values$method_num = "mean"
```

- Retrieving state: `$state` of individual PipeOps (*after* `$train()`)

```
graph_pp$pipeops$scale$state %>% head(1)
## $center
## Petal.Length  Petal.Width  Sepal.Length  Sepal.Width
##      3.758000      1.199333      5.843333      3.057333
```

MLR3 PIPELINES IN ACTION

LINEAR PREPROCESSING PIPELINE

- Setting / retrieving parameters: `$param_set`

```
graph_pp$pipeops$impute$param_set$values$method_num = "mean"
```

- Retrieving state: `$state` of individual PipeOps (*after* `$train()`)

```
graph_pp$pipeops$scale$state %>% head(1)

## $center
## Petal.Length  Petal.Width Sepal.Length  Sepal.Width
##      3.758000      1.199333      5.843333      3.057333
```

- Retrieving intermediate results: `$.result` (set debug option before)

```
graph_pp$pipeops$scale$.result[[1]]$data() %>% head(3)

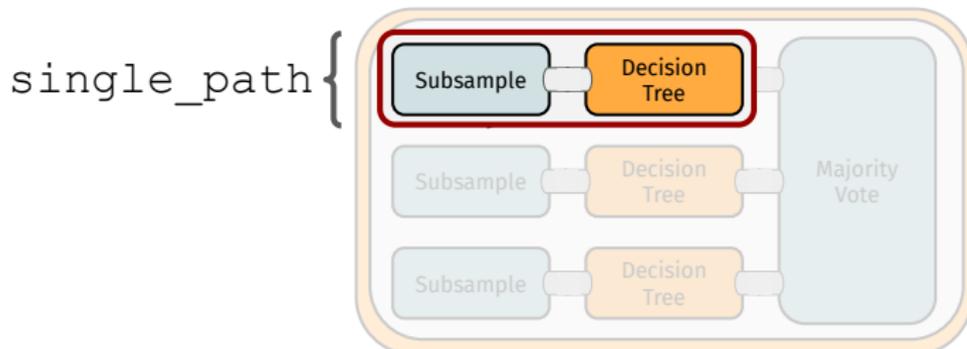
##   Species Petal.Length Petal.Width Sepal.Length Sepal.Width
## 1:  setosa   -1.335752   -1.311052   -0.8976739    1.0156020
## 2:  setosa   -1.335752   -1.311052   -1.1392005   -0.1315388
## 3:  setosa   -1.392399   -1.311052   -1.3807271    0.3273175
```

NONLINEAR PIPELINES

MLR3 PIPELINES IN ACTION

ENSEMBLE METHOD: BAGGING

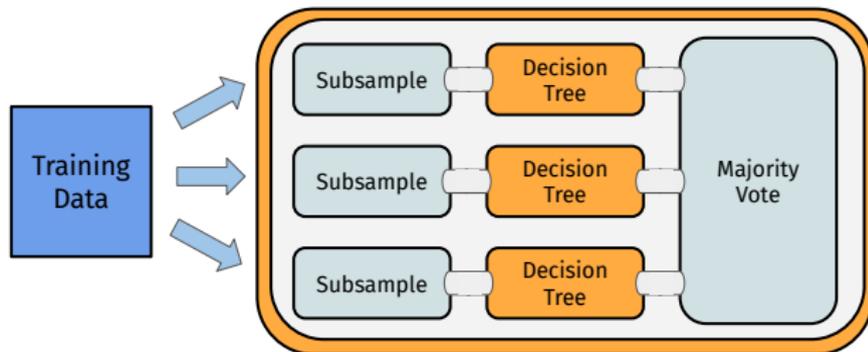
```
single_path = "subsample" %>>%  
  mlr_pipeops$get("learner",  
    learner = mlr_learners$get("classif.rpart"))
```



MLR3 PIPELINES IN ACTION

ENSEMBLE METHOD: BAGGING

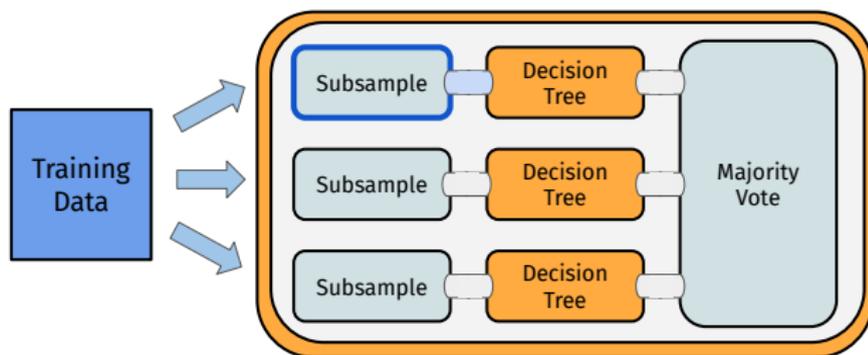
```
single_path = "subsample" %>>%  
  mlr_pipeops$get("learner",  
    learner = mlr_learners$get("classif.rpart"))  
  
graph_bag = greplicate(single_path, n = 3) %>>%  
  mlr_pipeops$get("majorityvote", innum = 3)
```



MLR3 PIPELINES IN ACTION

ENSEMBLE METHOD: BAGGING

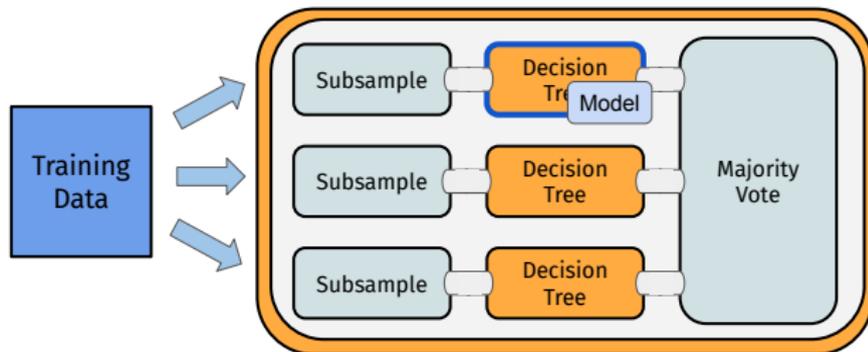
```
single_path = "subsample" %>>%  
  mlr_pipeops$get("learner",  
    learner = mlr_learners$get("classif.rpart"))  
  
graph_bag = grePLICATE(single_path, n = 3) %>>%  
  mlr_pipeops$get("majorityvote", innum = 3)
```



MLR3 PIPELINES IN ACTION

ENSEMBLE METHOD: BAGGING

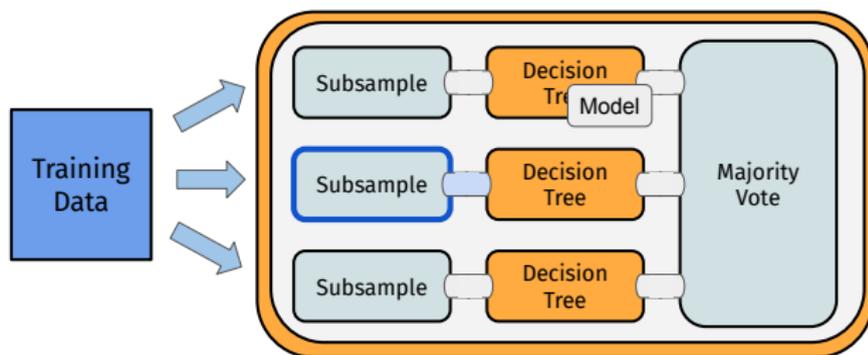
```
single_path = "subsample" %>>%  
  mlr_pipeops$get("learner",  
    learner = mlr_learners$get("classif.rpart"))  
  
graph_bag = grePLICATE(single_path, n = 3) %>>%  
  mlr_pipeops$get("majorityvote", innum = 3)
```



MLR3 PIPELINES IN ACTION

ENSEMBLE METHOD: BAGGING

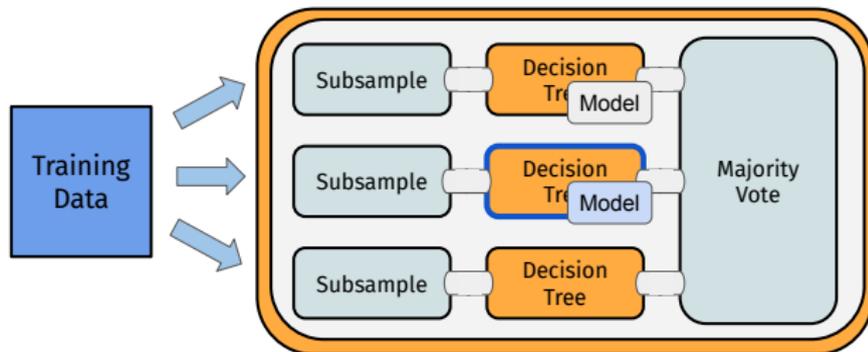
```
single_path = "subsample" %>>%  
  mlr_pipeops$get("learner",  
    learner = mlr_learners$get("classif.rpart"))  
  
graph_bag = greplicate(single_path, n = 3) %>>%  
  mlr_pipeops$get("majorityvote", innum = 3)
```



MLR3 PIPELINES IN ACTION

ENSEMBLE METHOD: BAGGING

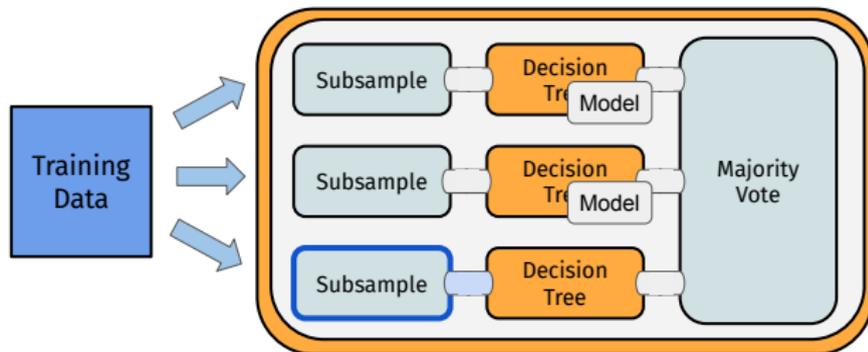
```
single_path = "subsample" %>>%  
  mlr_pipeops$get("learner",  
    learner = mlr_learners$get("classif.rpart"))  
  
graph_bag = grePLICATE(single_path, n = 3) %>>%  
  mlr_pipeops$get("majorityvote", innum = 3)
```



MLR3 PIPELINES IN ACTION

ENSEMBLE METHOD: BAGGING

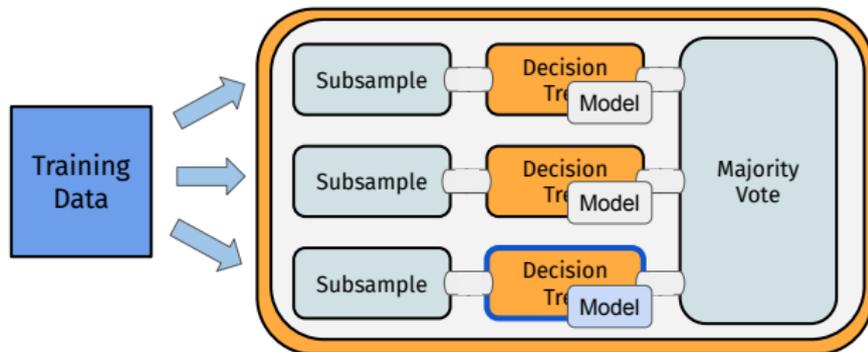
```
single_path = "subsample" %>>%  
  mlr_pipeops$get("learner",  
    learner = mlr_learners$get("classif.rpart"))  
  
graph_bag = greplicate(single_path, n = 3) %>>%  
  mlr_pipeops$get("majorityvote", innum = 3)
```



MLR3 PIPELINES IN ACTION

ENSEMBLE METHOD: BAGGING

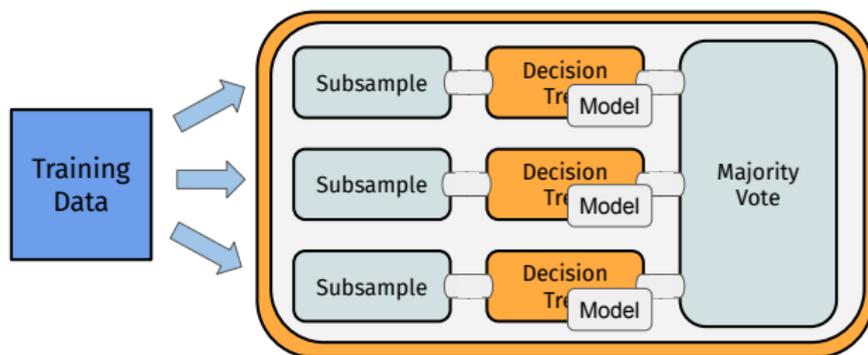
```
single_path = "subsample" %>>%  
  mlr_pipeops$get("learner",  
    learner = mlr_learners$get("classif.rpart"))  
  
graph_bag = greplicate(single_path, n = 3) %>>%  
  mlr_pipeops$get("majorityvote", innum = 3)
```



MLR3 PIPELINES IN ACTION

ENSEMBLE METHOD: BAGGING

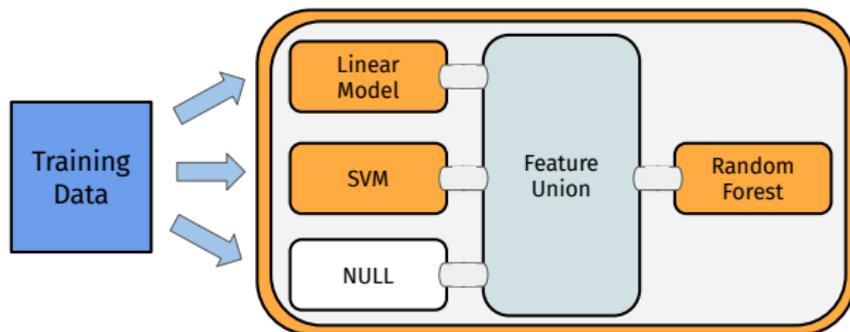
```
single_path = "subsample" %>>%  
  mlr_pipeops$get("learner",  
    learner = mlr_learners$get("classif.rpart"))  
  
graph_bag = greplicate(single_path, n = 3) %>>%  
  mlr_pipeops$get("majorityvote", innum = 3)
```



MLR3PIPELINES IN ACTION

ENSEMBLE METHOD: STACKING

```
graph_stack = gunion(list(  
  mlr_pipeops$get("learner_cv",  
    learner = mlr_learners$get("regr.lm")),  
  mlr_pipeops$get("learner_cv",  
    learner = mlr_learners$get("regr.svm")),  
  "null")) %>>%  
mlr_pipeops$get("featureunion", innum = 3) %>>%  
mlr_pipeops$get("learner",  
  learner = mlr_learners$get("regr.ranger"))
```

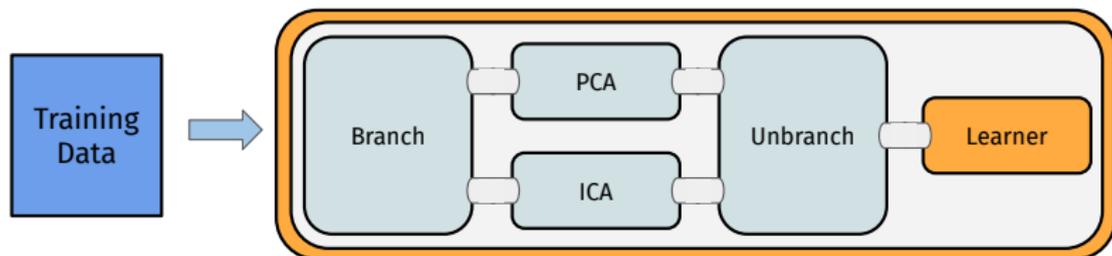


MLR3 PIPELINES IN ACTION

BRANCHING

```
graph_branch = mlr_pipeops$get("branch", c("pca", "ica")) %>>%  
  union(list("pca", "ica")) %>>%  
  mlr_pipeops$get("unbranch", c("pca", "ica")) %>>%  
  mlr_pipeops$get("learner",  
    learner = mlr_learners$get("classif.kknn"))
```

Execute only one of several alternative paths

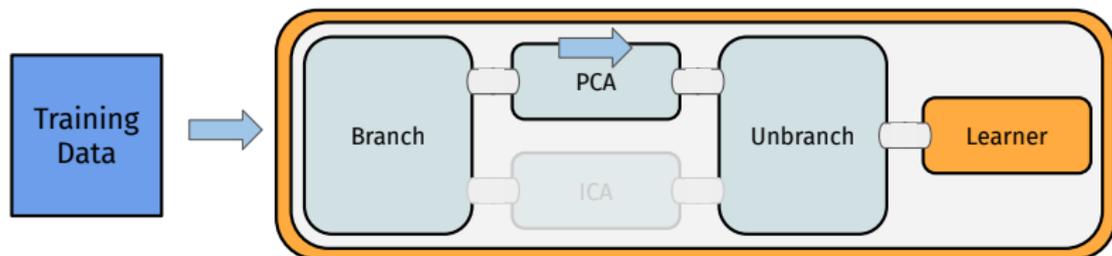


MLR3 PIPELINES IN ACTION

BRANCHING

```
graph_branch = mlr_pipeops$get("branch", c("pca", "ica")) %>>%  
  union(list("pca", "ica")) %>>%  
  mlr_pipeops$get("unbranch", c("pca", "ica")) %>>%  
  mlr_pipeops$get("learner",  
    learner = mlr_learners$get("classif.kknn"))
```

```
> graph_branch$pipeops$branch$  
  param_set$values$selection = "pca"
```

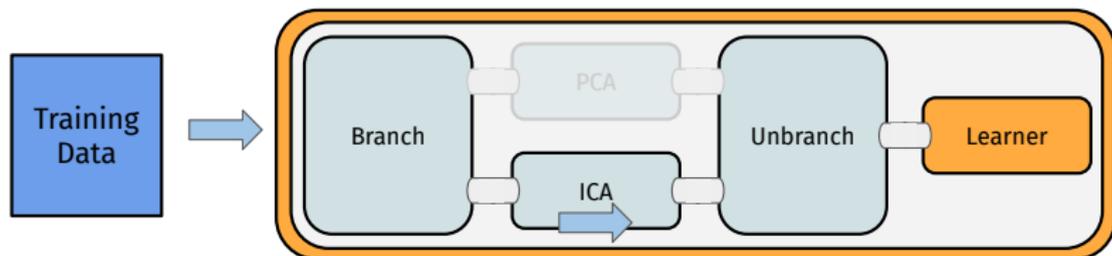


MLR3 PIPELINES IN ACTION

BRANCHING

```
graph_branch = mlr_pipeops$get("branch", c("pca", "ica")) %>>%  
  union(list("pca", "ica")) %>>%  
  mlr_pipeops$get("unbranch", c("pca", "ica")) %>>%  
  mlr_pipeops$get("learner",  
    learner = mlr_learners$get("classif.kknn"))
```

```
> graph_branch$pipeops$branch$  
  param_set$values$selection = "ica"
```



HYPERPARAMETERS AND TUNING

HYPERPARAMETERS AND TUNING

- PipeOps have *hyperparameters* (using `paradox` pkg)
- Graphs have hyperparameters of all components *combined*
- \Rightarrow simultaneous **Tuning** of Learner and preprocessing (mlr3tuning package)

```
library("paradox") ; library("mlr3tuning")
glrn = "scale" %>>% mlr_pipeops$get("learner",
  mlr_learners$get("classif.rpart"))
ps = ParamSet$new(list(
  ParamLgl$new("scale.scale"),
  ParamInt$new("classif.rpart.minsplit", 1, 20)
))
ff = PerformanceEvaluator$new(task, glrn, "cv", "classif.ce", ps)
tuner = TunerRandomSearch$new(ff, TerminatorEvaluations$new(10))

tuner$tune()

tuner$tune_result()
```

MLR3PIPELINES

NOT SHOWN HERE:

- Many more PipeOps: select, apply, encode, ...
- Automatic type-checking when constructing Graphs
- Interactive (html + javascript) plots
- Extensible by R6 inheritance of PipeOp base class

UPCOMING FEATURES

- More PipeOps
- Caching of expensive results
- Automatically parallel execution of concurrent operations

Thanks! Questions? Comments? Comment on Github?

mlr3: <https://github.com/mlr-org/mlr3>

mlr3pipelines: <https://github.com/mlr-org/mlr3pipelines>