

# Building Ensemble Models

*Xiaoyong Pan, Jenna Reys, Peter R. Rijnbeek*

*2018-09-09*

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Usage</b>	<b>2</b>
<b>3</b>	<b>Apply Ensemble model</b>	<b>3</b>
<b>4</b>	<b>Acknowledgments</b>	<b>3</b>

## 1 Introduction

Ensemble models combine several models to improve the overall performance. Traditionally, weak learners were combined to boost performance but recent results show that combining several strong approaches can result in a better performance. There are many examples in literature where ensemble models outperform individual models using stacking, i.e. a final logistic regression layer across the individual model outputs, but other approaches like weighing has also shown promising results.

This vignette describes how you can use the Observational Health Data Sciences and Informatics (OHDSI) `PatientLevelPrediction` package to build ensemble models. This vignette assumes you have read and are comfortable with building single patient level prediction models as described in the `BuildingPredictiveModels` vignette.

This will enable studying ensemble methods at scale in the OHDSI data network.

In `PatientLevelPrediction` package, four ensemble strategies have been implemented:

1. average ensemble: Calculate the average probability from individual models
2. product ensemble: Calculate the product of probabilities from individual models.
3. weighted ensemble: Calculate the weighted average probability from individual models using train AUC as weights.
4. stacked ensemble: Train a logistics regression on outputs from individual models

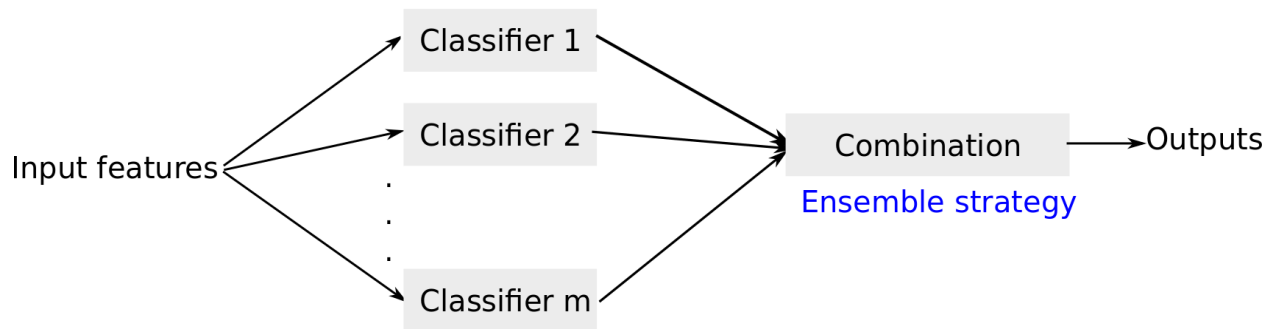


Figure 1: Ensemble model

## 2 Usage

Use the `PatientLevelPrediction` package to generate a `population` and `plpData` object. Alternatively, you can make use of the data simulator. The following code snippet creates a population of 12000 patients.

```
set.seed(1234)
data(plpDataSimulationProfile)
sampleSize <- 12000
plpData <- simulatePlpData(
  plpDataSimulationProfile,
  n = sampleSize
)

population <- createStudyPopulation(
  plpData,
  outcomeId = 2,
  binary = TRUE,
  firstExposureOnly = FALSE,
  washoutPeriod = 0,
  removeSubjectsWithPriorOutcome = FALSE,
  priorOutcomeLookback = 99999,
  requireTimeAtRisk = FALSE,
  minTimeAtRisk = 0,
  riskWindowStart = 0,
  addExposureDaysToStart = FALSE,
  riskWindowEnd = 365,
  addExposureDaysToEnd = FALSE,
  verbosity = futile.logger::INFO
)
```

Specify the prediction algorithms to be combined.

```
# Use LASSO logistic regression and Random Forest as base predictors
model1 <- setLassoLogisticRegression()
model2 <- setRandomForest()
```

Specify a test fraction and a sequence of training set fractions.

```
testFraction <- 0.2
```

Specify an `ensembleStrategy` to combine multiple predictors. The strategy used for ensembling the outputs from different models, it can be 'mean', 'product', 'weighted' and 'stacked': 'mean' the average probability from different models 'product' the product rule 'weighted' the weighted average probability from different models using train AUC as weights. 'stacked' the stacked ensemble trains a logistics regression on different models.

```
ensembleStrategy <- 'stacked'
```

Specify the test split to be used.

```
# Use a split by person, alternatively a time split is possible
testSplit <- 'person'
```

Run the ensemble learning to combine `model1` and `model2`. You can also use different `plpData` for different models.

```
results <- PatientLevelPrediction::runEnsembleModel(population,
  dataList = list(plpData, plpData),
```

```
modellist = list(model1, model2),
testSplit=testSplit,
testFraction=testFraction,
nfold=3, splitSeed=1000,
ensembleStrategy = ensembleStrategy)
```

### 3 Apply Ensemble model

```
plpData <- loadPlpData("<data file>")
modellist <- loadEnsemblePlpModel("<model>")
populationSettings <- plpModel$populationSettings
populationSettings$plpData <- plpData
population <- do.call(createStudyPopulation, populationSettings)
```

Get the prediction, please make sure the ensemble strategy for training and applying is the same:

```
prediction <- applyEnsembleModel(population,
                                dataList = list(plpData, plpData),
                                modellist = modellist,
                                analysisId = NULL,
                                save = NULL,
                                ensembleStrategy = "stacked")$prediction
```

### 4 Acknowledgments

Considerable work has been dedicated to provide the PatientLevelPrediction package.

```
citation("PatientLevelPrediction")
```

```
##
## Jenna Reys, Martijn J. Schuemie, Marc A. Suchard, Patrick B.
## Ryan and Peter R. Rijnbeek (2018). PatientLevelPrediction:
## Package for patient level prediction using data in the OMOP
## Common Data Model. R package version 2.0.5.
##
## A BibTeX entry for LaTeX users is
##
## @Manual{,
##   title = {PatientLevelPrediction: Package for patient level prediction using data in the OMOP
## Common Data Model},
##   author = {Jenna Reys and Martijn J. Schuemie and Marc A. Suchard and Patrick B. Ryan and Peter R.
##   year = {2018},
##   note = {R package version 2.0.5},
## }
```

**Please reference this paper if you use the PLP Package in your work:**

Reys JM, Schuemie MJ, Suchard MA, Ryan PB, Rijnbeek PR. Design and implementation of a standardized framework to generate and evaluate patient-level prediction models using observational healthcare data. J Am Med Inform Assoc. 2018;25(8):969-975.