

# Package ‘SAMGEP’

May 7, 2026

**Type** Package

**Title** A Semi-Supervised Method for Prediction of Phenotype Event Times

**Version** 0.1.0-1

**Description** A novel semi-supervised machine learning algorithm to predict phenotype event times using Electronic Health Record (EHR) data.

**URL** <https://github.com/celehs/SAMGEP>

**BugReports** <https://github.com/celehs/SAMGEP/issues>

**License** GPL-3

**Encoding** UTF-8

**RoxygenNote** 7.1.1

**Depends** R (>= 3.5.0)

**Imports** stats, mvtnorm, nlme, pROC, abind, nloptr, foreach,  
doParallel, parallel, Rcpp

**LinkingTo** Rcpp, RcppArmadillo

**Suggests** knitr, rmarkdown

**VignetteBuilder** knitr

**LazyData** true

**NeedsCompilation** yes

**Author** Yuri Ahuja [aut, cre],  
Tianxi Cai [aut],  
PARSE LTD [aut]

**Maintainer** Yuri Ahuja <Yuri\_Ahuja@hms.harvard.edu>

**Repository** CRAN

**Date/Publication** 2021-01-06 10:00:02 UTC

## Contents

SAMGEP-package	2
samgep	2
simdata	4

**Index****5**


---

SAMGEP-package	<i>SAMGEP: A Semi-supervised Method for Prediction of Phenotype Event Times Using the Electronic Health Record</i>
----------------	--

---

**Description**

Semi-supervised Adaptive Markov Gaussian Embedding Process (SAMGEP) is a novel semi-supervised machine learning algorithm to predict phenotype event times using Electronic Health Record (EHR) data.

---

samgep	<i>Semi-supervised Adaptive Markov Gaussian Process (SAMGEP)</i>
--------	--

---

**Description**

Semi-supervised Adaptive Markov Gaussian Process (SAMGEP)

**Usage**

```

samgep(
  dat_train = NULL,
  dat_test = NULL,
  Cindices = NULL,
  w = NULL,
  w0 = NULL,
  V = NULL,
  observed = NULL,
  nX = 10,
  covs = NULL,
  survival = FALSE,
  Estep = Estep_partial,
  Xtrain = NULL,
  Xtest = NULL,
  alpha = NULL,
  r = NULL,
  lambda = NULL,
  surrIndex = NULL,
  nCores = 1
)

```

**Arguments**

<code>dat_train</code>	(optional if <code>Xtrain</code> is supplied) Raw training data set, including patient IDs (ID), healthcare utilization feature (H) and censoring time (C)
<code>dat_test</code>	(optional) Raw testing data set, including patient IDs (ID), a healthcare utilization feature (H) and censoring time (C)
<code>Cindices</code>	(optional if <code>Xtrain</code> is supplied) Column indices of EHR feature counts in <code>dat_train/dat_test</code>
<code>w</code>	(optional if <code>Xtrain</code> is supplied) Pre-optimized EHR feature weights
<code>w0</code>	(optional if <code>Xtrain</code> is supplied) Initial (i.e. partially optimized) EHR feature weights
<code>V</code>	(optional if <code>Xtrain</code> is supplied) <code>nFeatures</code> x <code>nEmbeddings</code> embeddings matrix
<code>observed</code>	(optional if <code>Xtrain</code> is supplied) IDs of patients with observed outcome labels
<code>nX</code>	Number of embedding features (defaults to 10)
<code>covs</code>	(optional) Baseline covariates to include in model; not yet operational
<code>survival</code>	Binary indicator of whether target phenotype is of type survival (i.e. stays positive after incident event) or relapsing-remitting (defaults to FALSE)
<code>Estep</code>	E-step function to use ( <code>Estep_partial</code> or <code>Estep_full</code> ; defaults to <code>Estep_partial</code> )
<code>Xtrain</code>	(optional) Embedded training data set, including patient IDs (ID), healthcare utilization feature (H) and censoring time (C)
<code>Xtest</code>	(optional) Embedded testing data set, including patient IDs (ID), healthcare utilization feature (H) and censoring time (C)
<code>alpha</code>	(optional) Relative weight of semi-supervised to supervised MGP predictors in SAMGEP ensemble
<code>r</code>	(optional) Scaling factor of inter-temporal correlation
<code>lambda</code>	(optional) L1 regularization hyperparameter for feature weight ( <code>w</code> ) optimization
<code>surrIndex</code>	(optional) Index (within <code>Cindices</code> ) of primary surrogate index for outcome event
<code>nCores</code>	Number of cores to use for parallelization (defaults to 1)

**Value**

`w_opt` Optimized feature weights (`w`)

`r_opt` Optimized inter-temporal correlation scaling factor (`r`)

`alpha_opt` Optimized semi-supervised:supervised relative weight (`alpha`)

`lambda_opt` Optiized L1 regularization hyperparameter (`lambda`)

`margSup` Posterior probability predictions of supervised model (MGP Supervised)

`margSemisup` Posterior probability predictions of semi-supervised model (MGP Semi-supervised)

`margMix` Posterior probability predictions of SAMGEP

`cumSup` Cumulative probability predictions of supervised model (MGP Supervised)

`cumSemisup` Cumulative probability predictions of semi-supervised model (MGP Semi-supervised)

`cumMix` Cumulative probability predictions of SAMGEP

---

`simdata`*Simulated Dataset*

---

**Description**

Click [HERE](#) to view details.

**Usage**

```
simdata
```

**Format**

An object of class `list` of length 3.

**Examples**

```
str(simdata)
```

# Index

\* **datasets**

simdata, [4](#)

\* **package**

SAMGEP-package, [2](#)

samgep, [2](#)

SAMGEP-package, [2](#)

simdata, [4](#)