

## VARITY FRAMWORK USER GUIDE

<b>1. INTRODUCTION</b>	<b>1</b>
<b>1.1 STEPS TO APPLY VARITY FRAMEWORK</b>	<b>1</b>
<b>1.2 TO BUILD VARITY MODELS</b>	<b>1</b>
<b>1.3 SYSTEM REQUIREMENT</b>	<b>2</b>
<b>1.4 TECHNICAL SUPPORT</b>	<b>2</b>
<b>2. VARITY DATA PREPARATION</b>	<b>2</b>
<b>3. VARITY SESSION CONFIGURATION</b>	<b>2</b>
<b>3.1 DATA OBJECT CONFIGURATION</b>	<b>2</b>
<b>3.2 ESTIMATOR OBJECT CONFIGURATION</b>	<b>3</b>
<b>3.3 PREDICTOR OBJECT CONFIGURATION</b>	<b>4</b>
<b>3.4 QUALITY INFORMATIVE PROPERTY CONFIGURATION</b>	<b>5</b>
<b>3.5 HYPERPARAMETER CONFIGURATION</b>	<b>6</b>
<b>4. VARITY COMMANDS</b>	<b>6</b>
<b>5. TABLE 1: REQUIRED PYTHON PACKAGES</b>	<b>10</b>

### 1. Introduction

VARITY is a supervised machine learning approach to build specialized predictive models using training examples with optimized differential weights. Training examples are assembled into different training sets: 1) one core set of training examples which are known to have high quality. 2) one or more add-on sets of training examples with uncertain quality. For each training set, the weights of examples are determined using one or more logistic functions each takes one quality-informative property as input. The parameters of each logistic function are treated as hyper-parameters and optimized for performance on the core set of examples using cross-validation.

#### 1.1 Steps to apply VARITY framework

1. Assemble training examples (high-quality core set and add-on sets with uncertain quality).
2. Identify quality-informative properties for each add-on set using domain knowledge (Optional: verify using moving window analysis).
3. Config hyper-parameters based on quality-informative properties (parameters of corresponding logistic functions).
4. Optional: Run nested cross-validation to evaluate performance.
5. Run hyper-parameter optimization to determine weight of each training example.
6. Train the final VARITY model with core set and weighted add-on sets.

## 1.2 System Requirement

Python version 3.7.2 and a few python packages listed in Table 1.

## 1.3 Technical Support

Please contact [joe.wu.ca@gmail.com](mailto:joe.wu.ca@gmail.com) for technical support.

## 2. VARITY project folder

A project folder that has sub-folders with pre-defined names is needed to start a VARITY project, please refer to <https://github.com/joewuca/verity/tree/master/project/> folder.

## 3. VARITY data preparation

All core and add-on sets of training examples need to be assembled into one single CSV file as an input, and data columns include:

- Feature columns (currently only supports Real or Integer type features)
- A “label” column: the dependent variable (0 or 1 for binary classification problem)
- A “extra\_data” column: 0 for core examples, and 1 for add-on examples.
- A “set\_name” column: name for different core and add-on sets.
- Candidate quality-informative property columns.

VARITY data files location: `/project_path/data`

The data used in manuscript “*Improved pathogenicity prediction for rare human missense variants*” is located in ([http://verity.varianteffect.org/downloads/VARITY\\_training.tar.gz](http://verity.varianteffect.org/downloads/VARITY_training.tar.gz))

## 4. VARITY session configuration

VARITY session uses a configuration file (`/project_path/config/YOUR_SESSION_ID.vsc`) to setup everything needed for training a VARITY model. The configuration file contains settings for four different type of objects used in VARITY framework: 1) Data 2) Estimator 3) Predictor 4) Quality informative property 5) Hyperparameter. Each object has a list of attributes, and their names can be found in a “object definition” line (starts with “\*” symbol) delimited by “|” sign. A new object instance can be created below the “object definition” line, starting with the object instance name and then values of each attribute. A valid session config file is needed before running any VARITY commands. The config file used in manuscript “*Improved pathogenicity prediction for rare human missense variants*” is located in <https://github.com/joewuca/verity/tree/master/project/config/> folder.

### 4.1 Data object configuration

The data object takes the input training data (consists of core and add-on training examples) and splits it in nested cross validation fashion. The core training examples were first splitted into  $k$  outer-loop folds. For each outer-loop, the core training examples are splitted into an outer-loop

training set and an outer-loop test set, and subsequently the outer-loop training examples are further splitted into  $k_2$  inner-loop folds. The purpose of outer-loop cross-validation is to fairly evaluate model performance on core set of examples, and inner-loop cross-validation is to optimize hyperparameter for each outer-loop. The attributes for data object:

- **test\_split\_method**

The way of splitting for outer-loop cross-validation.

0: random split

1: stratified split (keep same prior for the training and test in each outer-loop fold)

- **test\_split\_folds**

Number of folds ( $k_1$ ) for outer-loop cross-validation

- **test\_split\_ratio**

Only used when the “test\_split\_folds” attribute is set to 1. When there is only one outer-loop fold, this value determines the fraction of the core set of examples as test set. If the value is set to 0, the corresponding inner-loop cross-validation will be based on the whole core set of examples, which is the appropriate configuration for the data object used for building a final VARIETY model.

- **cv\_split\_method**

The way of splitting for inner-loop cross-validation.

0: random split

1: stratified split (keep same prior for the training and validation in each inner-loop fold)

- **cv\_split\_folds**

Number of folds for inner-loop cross-validation

- **cv\_split\_ratio**

Only used when the “cv\_split\_folds” attribute is set to 1. When there is only one inner-loop fold, this value determines the fraction of the outer-loop training examples as inner-loop validation set.

- **data\_file**

The absolute path of the input training data.

#### 4.2 Estimator object configuration

VARIETY framework is designed to support a list of different machine learning algorithms. However, currently it only support the gradient boosted trees algorithm. The attributes for data object:

- **algo\_name**

The name of the learning algorithm. The name for gradient boosted trees for binary classification is “xgb\_c”

- **round\_digits**

The significant digits of the output metrics (e.g., AUROC)

#### 4.3 Predictor object configuration

The predictor object is associated with one data object instance and one estimator instance, and has following attributes:

- **type**

0: VARITY model predictor

1: Other existing predictor (predictions can be found as one column in the data file)

- **ml\_type**

VARITY framework is designed to support different type of machine learning tasks. However, currently only the following type is supported:

“classification\_binary”: binary classification

- **data**

the name of the associated data object instance

- **estimator**

the name of the associated estimator object instance

- **tune\_obj**

VARITY framework is designed to support a list of metrics as the objective function for hyperparameter optimization and performance evaluation using cross-validation. Currently only the following metrics are supported:

“macro\_cv\_aubprc”: The Area Under Balanced Precision Recall Curve via cross-validation

“macro\_cv\_auroc”: The Area Under ROC curve via cross-validation

- **hyperopt\_trials**

The number of hyperparameter tuning trials (HyperOpt trials)

- **trials\_mv\_size**

The number of trials in one moving window (used for select the best trial, see Section 4 save\_best\_hp command).

- **features**

A list of features used for the predictor. For other predictors with existing predictions (type = 1), just use the corresponding column name in the data file as a single feature. All features need to be put in square brackets ([]) and separated by comma.

#### 4.4 Quality informative property (qip) configuration

- **predictors**

A list of predictors (see 3.3) associated with the current quality informative property.

- **weight\_function**

The weight function that takes the current quality informative property as input. The output of the weight function is used for weighting training examples. Currently only 'logistic' is supported.

- **hyperparameters**

The name of the parameters for the weight function that are treated as hyperparameters.

- **set\_list**

A list of training sets. The current weight function (quality informative property) applies on the training examples in the all the sets listed here.

- **set\_type**

The type of the listed sets, either 'core' or 'addon'

- **qip\_col**

The column name of quality informative property in the data files.

- **direction**

0: order the quality-informative property from low to high in moving windows analysis.

1: order the quality-informative property from high to low in moving windows analysis.

- **mv\_size\_precent**

For moving window analysis, this parameter determines the number of examples in each moving window, which equals to product of [number of examples in an add-on set] and [mv\_size\_percent] (rounded).

- **mv\_data\_points**

Number of moving windows for moving window analysis

- **enable**

1: The quality informative property is enabled for hyper-parameter tuning

0: The quality informative property is disabled for hyper-parameter tuning

## 4.5 Hyperparameter configuration

There are two types of hyperparameters for a VARITY predictor; 1) algorithm level parameters with respect to the estimator. 2) weighting parameters for training sets. Each hyperparameter has following attributes:

- **qip**

The quality informative property defined in 3.4

- **hp\_type**

1: logistic parameter (growth rate  $k$  and maximum value  $L$ )

2: logistic parameter (mid-point  $x_0$ )

3: algorithm parameter

- **from**

The lower bound of the hyperparameter value.

- **to**

The higher bound of the hyperparameter value. For filtering parameters, the higher bound is determined automatically based on number of examples in the associated add-on set.

- **step**

The difference between each hyperparameter value from the lower bound to higher bound. For filtering parameters, this parameter indicates the number of examples in a “filtering block”. The add-on set is filtered out in blocks instead of individual example.

- **default**

The default value of the hyperparameter.

- **data\_type**

The data type of the hyperparameter (int or real)

- **data\_interval**

number of possible mid-points.

- **significant\_digits**

number of significant digits for the value of parameters, use ‘None’ for no restriction.

## 5. VARITY commands

To run VARITY commands, please take the following steps first:

- 1) Download code in this git repository (<https://github.com/joewuca/varity/tree/master/python>) to a local folder as your VARITY script folder.
- 2) Make sure you have installed python 3, and all associated packages needed for VARITY framework (See Table 1)
- 3) Make sure your python PATH has included the VARITY script folder.

VARITY framework currently supports the following commands:

- **init\_session**

```
python3 varity_run.py actions=init_session session_id=YOUR_SESSION_ID
project_path=PATH_OF_YOUR_PROJECT_FOLDER
```

This command initiates a session (creating all necessary VARITY framework objects) based on the current configuration file. Unless you reinitiate the session again, changes to the configuration file will not affect the initialized VARITY objects. You can reinitiate the session by adding argument *reinitiate=1* to the command line, but you might need to generate all your results again after session re-initiation.

**output:**

- */project\_path/output/npv/[session\_id]\_[EACH DATA\_INSTANCE\_NAME]\_savedata.npy*
- */project\_path/output/npv/[session\_id]\_[EACH PREDICTOR\_INSTANCE\_NAME]\_hp\_config\_dict.npy*

- **mv\_analysis**

```
python3 varity_run.py actions=mv_analysis session_id=YOUR_SESSION_ID
predictor=PREDICTOR_INSTANCE_NAME mv_qip=QIP_INSTANCE_NAME
project_path=PATH_OF_YOUR_PROJECT_FOLDER
```

Moving analysis first order the add-on sets (“set\_list” attribute of *mv\_qip*) examples by the informative property (“qip\_col” attribute of *mv\_qip*) in ascending or descending order (“direction” attribute of *mv\_qip*), then create [“mv\_data\_points” attribute of *mv\_qip*] number of moving windows (size of each window equals to [number of examples in add-on sets] \* [“mv\_size\_precent” attribute of *mv\_qip*]). The predictive utility of each window is estimated using 10-fold cross validation on the core training set, where the training examples in each fold were supplemented by all of the add-on examples in that moving window.

**output:**

- */project\_path/output/csv/[session\_id]\_mv\_analysis\_[predictor]\_[mv\_qip].csv*

- **plot\_mv\_result**

```
python3 varity_run.py action=plot_mv_result session_id=YOUR_SESSION_ID predictor=  
PREDICTOR_INSTANCE_NAME mv_qip =QIP_INSTANCE_NAME  
project_path=PATH_OF_YOUR_PROJECT_FOLDER
```

Plot the moving analysis result (the predictive utility of each moving window)

**output:**

- /project\_path/output/img/[session\_id]\_plot\_mv\_result\_[predictor]\_[mv\_qip].png

- **hp\_tuning**

```
python3 varity_run.py action=hp_tuning session_id=YOUR SESSION ID  
predictor=PREDICTOR_INSTANCE_NAME cur_test_fold= THE OUTER-LOOP FOLD  
project_path=PATH_OF_YOUR_PROJECT_FOLDER
```

Hyperparameter tuning for the input **predictor** on the specified outer-loop fold. For the predictors for nested cross validation, the possible outer-loop fold value is from 0 to “test\_split\_folds” attribute of the corresponding data object minus one. For the final VARIETY model predictor (compare to predictors used for nested cross-validation), there is only one dummy outer loop therefore the outer-loop fold should be set to 0.

**output:**

- /project\_path/output/npy/[session\_id]\_[predictor]\_[filtering\_hp]\_tf[cur\_test\_fold]\_trials.pkl
- /project\_path/output/csv/[session\_id]\_[predictor]\_[filtering\_hp]\_tf[cur\_test\_fold]\_trial\_results.txt

- **save\_best\_hp**

```
python3 varity_run.py action=save_best_hp session_id=YOUR_SESSION_ID  
predictor=PREDICTOR_INSTANCE_NAME cur_test_fold=THE OUTER-LOOP FOLD  
project_path=PATH_OF_YOUR_PROJECT_FOLDER
```

Select and save the optimum hyperparameter setting from all hyperparameter optimization (HyperOpt) trials using the following procedure: 1) Re-order all trials by mean metric (“tune\_obj” attribute of **predictor**) on training sets (averaged over 10 training sets) from low to high; 2) calculate a moving window (we used window size 100) average of mean metric on validation sets; 3) define an “early stopping” point at the first moving window (the “fittest” region) for which mean metric on validation sets begins to descend; 4) Select as final the hyperparameters from the trial within this “fittest” region that achieved the highest mean metric on validation sets.

**output:**

- /project\_path/output/npy/[session\_id]\_[predictor]\_tf[cur\_test\_fold]\_hp\_dict.npy
- /project\_path/output/csv/[session\_id]\_[predictor]\_tf[cur\_test\_fold]\_best\_hps.csv



- */project\_path/output/img/[session\_id]\_[predictor]\_tf[cur\_test\_fold]\_hp\_selection.png*
- **plot\_hp\_weight**

*python3 varity\_run.py action=plot\_hp\_weight session\_id=YOUR\_SESSION\_ID  
predictor=PREDICTOR\_INSTANCE\_NAME cur\_test\_fold=THE OUTER-LOOP FOLD  
filtering\_hp=HYPERPARAMETER\_INSTANCE\_NAME*

This command plots the weight of each example in an add-on set or combined add-on sets (“source” attribute of **filtering\_hp**) ordered by the informative property (“orderby” of **filtering\_hp**). The filtering threshold and weight are based on the optimized hyperparameters.

**output:**

  - */project\_path/output/img/[session\_id]\_[predictor]\_[filtering\_hp].png*
- **test\_cv\_prediction**

*python3 varity\_run.py action=test\_cv\_prediction session\_id=YOUR\_SESSION\_ID  
predictor=PREDICTOR\_INSTANCE\_NAME*

This command runs nested cross-validation. For each outer-loop, It makes predictions the test set using model trained with the optimized hyper-parameters via inner-loop cross-validation.

**output:**

  - */project\_path/output/npy/[session\_id]\_[predictor]\_test\_cv\_results.npy*
  - */project\_path/output/csv/[session\_id]\_[predictor]\_hp\_test\_cv\_results.csv*
- **plot\_test\_result**

*python3 varity\_run.py action=plot\_test\_result session\_id=YOUR\_SESSION\_ID predictor=PREDICTOR\_INSTANCE\_NAME compare\_predictors=[PREDICTORS\_FOR\_COMPARISON]*

Plot the balanced precision recall curve and ROC curve using the results from nested cross-validation. The statistical test is carried out between each predictor specified in **compare\_predictors** (usually non-VARITY predictors) and the predictor specified in **predictor** (usually VARITY predictor). For each outer-loop fold, the test set is filtered if there is a missing prediction from any of the predictor specified in **compare\_predictors**.

**output:**

  - */project\_path/output/npy/[session\_id]\_[predictor]\_filter\_1\_auroc\_interp.png*
  - */project\_path/output/csv/[session\_id]\_[predictor]\_filter\_1\_aubprc\_interp.png*
- **target\_prediction**

*python3 varity\_run.py action=target\_prediction session\_id=YOUR\_SESSION\_ID predictor=PREDICTOR\_INSTANCE\_NAME cur\_test\_fold=THE OUTER-LOOP FOLD  
target\_file=TARGET\_FILE\_NAME loo= [0 OR 1]*

Predict the examples in **target\_file**. If the **loo** is set to 1, then only the target examples that have been used in training will be predicted using leave-one-example-out strategy.

**output:**

- `/project_path/output/npv/[session_id]_predictor_tf[cur_test_fold]_target_predicted.csv`
- `/project_path/output/npv/[session_id]_predictor_tf[cur_test_fold]_target_loo_predicted.csv`

**6. Table 1: Required python packages**

Package Name	Version	Description (Link to document)
Cython	0.29.14	C extension for python ( <a href="https://cython.org/">https://cython.org/</a> )
graphviz	0.13.2	Open source graph visualization software ( <a href="https://graphviz.org/">https://graphviz.org/</a> )
hyperopt	0.2.2	Bayesian hyperparameter optimization ( <a href="https://github.com/hyperopt/hyperopt">https://github.com/hyperopt/hyperopt</a> )
matplotlib	3.1.0	Python visualization ( <a href="https://matplotlib.org/">https://matplotlib.org/</a> )
numpy	1.16.0	Scientific computing with Python ( <a href="https://numpy.org/">https://numpy.org/</a> )
pandas	0.24.0	Data analysis and manipulation tool ( <a href="https://pandas.pydata.org">https://pandas.pydata.org</a> )
scikit-learn	0.20.2	Machine learning in Python ( <a href="https://scikit-learn.org/stable/">https://scikit-learn.org/stable/</a> )
scipy	1.2.0	Python-based ecosystem of open-source software for mathematics, science, and engineering ( <a href="https://www.scipy.org/">https://www.scipy.org/</a> )
seaborn	0.9.0	Statistical data visualization ( <a href="https://seaborn.pydata.org/">https://seaborn.pydata.org/</a> )
shap	0.34.0	A game theoretic approach to explain the output of any machine learning model ( <a href="https://github.com/slundberg/shap">https://github.com/slundberg/shap</a> )
xgboost	0.90	An optimized distributed gradient boosting library ( <a href="https://xgboost.readthedocs.io/en/latest/">https://xgboost.readthedocs.io/en/latest/</a> )