

## Vertica Machine Learning V12.0.1 Cheat Sheet

Vertica Machine Learning supports the entire machine learning workflow via an SQL interface. For more information about the capabilities of Vertica ML, see the [Vertica ML documentation](#) or the [Vertica ML examples repository](#) on GitHub. To download the datasets used in this cheat sheet, you can follow the guidelines in the Vertica [documentation](#) or download the data directly from the [Vertica ML GitHub](#) repository. For a scikit-like machine learning library that integrates directly with the data in your Vertica database, see [VerticaPy](#).

## Preprocessing data

**Summarize data**

```
=> SELECT summarize_numcol(hits, salary) OVER() FROM baseball WHERE dob > '1975-7-1'::DATE; -- for numeric columns, displays 'COUNT', 'MEAN', 'STDDEV', 'MIN', 'PERC25', 'MEDIAN', 'PERC75', and 'MAX'
=> SELECT summarize_catcol(team USING PARAMETERS topk = 10) OVER() FROM baseball; --for a categorical column, displays 'CATEGORY', 'COUNT', and 'PERCENT'
```

**Detect outliers**

```
=> SELECT detect_outliers('baseball_outliers', 'baseball', 'hr, hits, avg, salary', 'robust_zscore' USING PARAMETERS outlier_threshold=3.0); -- outputs outliers to the 'baseball_outliers' table
=> SELECT iforest('baseball_outliers', 'baseball', 'team, hr, hits, avg, salary' USING PARAMETERS ntree=75, sampling_size=0.7, max_depth=15); --trains an iforest model on specified columns of the 'baseball' dataset
=> SELECT * FROM (SELECT first_name, last_name, apply_iforest(hr, hits, salary USING PARAMETERS model_name='baseball_outliers', contamination=0.1) AS predictions FROM baseball)
  AS outliers WHERE predictions.is_anomaly IS true; --applies 'baseball_outliers' to 'baseball' and list the outliers
```

**Measure correlations**

```
=> SELECT corr_matrix(hr, hits, avg, salary) OVER() FROM baseball; --calculates the Pearson Correlation Coefficient between each pair of input columns
```

**Normalize data**

```
=> SELECT normalize('baseball_normz', 'baseball', 'hr, hits', 'zscore'); --outputs a normalized result to view 'baseball_normz'
=> SELECT normalize_fit('baseball_normfitz', 'baseball', 'hr,hits', 'robust_zscore'); --computes and stores normalization parameters in model 'baseball_normfitz'
=> SELECT apply_normalize(* USING PARAMETERS model_name = 'baseball_normfitz') FROM baseball; --applies the normalization parameters from 'baseball_normfitz' to 'baseball'
=> SELECT reverse_normalize(* USING PARAMETERS model_name = 'baseball_normfitz') FROM baseball; --reverses the normalization transformation
```

**Dimensionality reduction**

```
=> SELECT pca('world_pca', 'world', '*' USING PARAMETERS exclude_columns='HDI, country'); --computes the principal components of 'world' and saves them in model 'world_pca'
=> CREATE TABLE worldPCA AS SELECT apply_pca(* USING PARAMETERS model_name='world_pca', exclude_columns='HDI, country', key_columns='HDI, country', cutoff=0.99) OVER() FROM world;
--uses the 'world_pca' model to transform the 'world' data such that the principal components account for 99 percent of the variance
=> SELECT apply_inverse_pca(HDI, country, col1, col2 USING PARAMETERS model_name='world_pca', exclude_columns='HDI, Country', key_columns='HDI, country') OVER() FROM worldPCA;
--reverts the apply_pca data transform
=> SELECT svd('world_svd', 'world', '*' USING PARAMETERS exclude_columns='HDI, country'); --computes the SVD decomposition of 'world' and saves the result in the model 'world_svd'
=> CREATE TABLE worldSVD AS SELECT apply_svd(* USING PARAMETERS model_name='world_svd', exclude_columns='HDI, country', key_columns='HDI, country', cutoff=0.99) OVER() FROM world;
-- computes the U matrix of the SVD decomposition of 'world' and stores the result in table 'worldSVD'
=> CREATE TABLE inverse_worldSVD AS SELECT apply_inverse_svd(* USING PARAMETERS model_name='world_svd', exclude_columns='HDI, country', key_columns='HDI, country') OVER() FROM worldSVD;
--transforms the table 'worldSVD' back to the original 'world' data
```

**Encode categorical features**

```
=> SELECT one_hot_encoder_fit('bTeamEncoder', 'baseball', 'team' USING PARAMETERS extra_levels={'team' : ['Red Sox']}); --generates the encoding for column 'team' and stores it in model 'bTeamEncoder'
=> CREATE VIEW baseballEncoded AS SELECT apply_one_hot_encoder(* USING PARAMETERS model_name='bTeamEncoder', drop_first=True, ignore_null=False) FROM baseball;
--generates the encoded columns in view 'baseballEncoded'
```

**Impute missing values**

```
=> SELECT impute ('baseballImputed', 'baseball', 'hits', 'mean' USING PARAMETERS partition_columns='team'); --imputes missing values for 'hits' based on the relevant team's mean 'hits' value
```

**Process imbalanced data**

```
=> SELECT balance ('baseballBalanced', 'baseball', 'team', 'hybrid_sampling'); --returns a view where each team is equally represented
```

**Sample data**

```
=> CREATE TABLE baseball_sample AS SELECT * FROM baseball TABLESAMPLE(25); --creates a new table containing a 25% sample of 'baseball'
```

## Training and predicting

## Regression

## Linear Regression

=> SELECT linear\_reg('linear\_reg\_faithful', 'faithful\_training', 'eruptions', 'waiting' USING PARAMETERS optimizer='bfgs'); --trains a linear regression model that uses 'waiting' time to predict the duration of eruptions  
 => SELECT id, predict\_linear\_reg(waiting USING PARAMETERS model\_name='linear\_reg\_faithful') FROM faithful\_testing; --applies model 'linear\_reg\_faithful' to the 'faithful\_testing' returns predicted eruption durations

## Support Vector Machines (SVM)

=> SELECT svm\_regressor('svm\_reg\_faithful', 'faithful\_training', 'eruptions', 'waiting' USING PARAMETERS error\_tolerance=0.1); --trains an SVM regression model to predict duration of eruptions based on 'waiting' time  
 => SELECT id, predict\_svm\_regressor(waiting USING PARAMETERS model\_name='svm\_reg\_faithful') FROM faithful\_testing; --applies 'svm\_reg\_faithful' to the 'faithful\_testing' and returns predicted eruption durations

## Random Forest

=> SELECT rf\_regressor('rf\_reg\_cars', 'mtcars\_train', 'mpg', 'carb', cyl, hp, drat, wt' USING PARAMETERS ntree=100); --trains a random forest model to predict a car's MPG based on a number of the car's attributes  
 => SELECT mpg, predict\_rf\_regressor(carb, cyl, hp, drat, wt USING PARAMETERS model\_name='rf\_reg\_cars') FROM mtcars\_test; --applies 'rf\_reg\_cars' to 'mtcars\_test' and returns predicted MPGs

## XGBoost

=> SELECT xgb\_regressor ('xgb\_cars', 'mtcars\_train', 'mpg', 'carb', cyl, hp, drat, wt' USING PARAMETERS learning\_rate=0.5); --trains an XGBoost model to predict a car's MPG based on several of the car's attributes  
 => SELECT mpg, predict\_xgb\_regressor(carb,cyl,hp,drat,wt USING PARAMETERS model\_name='xgb\_cars') FROM mtcars\_test; --applies 'xgb\_cars' to 'mtcars\_test' and returns predicted MPGs

## Classification

## Logistic Regression

=> SELECT logistic\_reg('logistic\_cars', 'mtcars\_train', 'am', 'mpg, cyl, disp, hp, drat, wt, qsec, gear, carb' USING PARAMETERS optimizer='BFGS', regularization='L2');  
 --trains a logistic regression model to predict whether a car has an automatic or manual transmission  
 => SELECT car\_model, predict\_logistic\_reg(mpg, cyl, disp, hp, drat, wt, qsec, gear, carb USING PARAMETERS model\_name='logistic\_cars') FROM mtcars\_test; --returns 'logistic\_cars' predictions on the 'mtcars\_test' dataset

## Support Vector Machines (SVM)

=> SELECT svm\_classifier('svm\_cars', 'mtcars\_train', 'am', 'mpg, cyl, disp, hp, drat, wt, qsec, vs, gear, carb' USING PARAMETERS exclude\_columns='hp,drat');  
 --trains an SVM model to predict whether a car has an automatic or manual transmission  
 => SELECT car\_model, am, predict\_svm\_classifier(mpg, cyl, disp, wt, qsec, vs, gear, carb USING PARAMETERS model\_name='svm\_cars') FROM mtcars\_test; --returns 'svm\_cars' predictions on the 'mtcars\_test' dataset

## Naive Bayes

=> SELECT naive\_bayes('naive\_congress', 'house84\_train', 'party', '\*' USING PARAMETERS exclude\_columns='party, id'); --trains a Naive Bayes model to predict a Congress member's party association  
 => SELECT party, predict\_naive\_bayes(vote1, vote2, vote3 USING PARAMETERS model\_name='naive\_congress', type='response') AS predicted\_party FROM house84\_test; --returns the actual party and the predicted party  
 => SELECT predict\_naive\_bayes\_classes(id, vote1, vote2, vote3 USING PARAMETERS model\_name='naive\_congress', key\_columns='id', exclude\_columns='id', classes='democrat, republican', match\_by\_pos='false') OVER() FROM house84\_test; --returns the predicted party and the probability for both 'democrat' and 'republican'

## Random Forest

=> SELECT rf\_classifier('rf\_iris', 'iris1', 'species', 'sepal\_length, sepal\_width, petal\_length, petal\_width' USING PARAMETERS ntree=100, sampling\_size=0.3); --trains a Random Forest model to predict 'species' based on feature measurements  
 => SELECT id, predict\_rf\_classifier(sepal\_length, sepal\_width, petal\_length, petal\_width USING PARAMETERS model\_name='rf\_iris') FROM iris2; --returns the predicted 'species' classifications  
 => SELECT predict\_rf\_classifier\_classes(id, sepal\_length, sepal\_width, petal\_length, petal\_width USING PARAMETERS model\_name='rf\_iris', key\_columns='id', exclude\_columns='id') OVER () FROM iris2;  
 --returns the predicted 'species' class and its probability

## XGBoost

=> SELECT xgb\_classifier('xgb\_iris', 'iris1', 'species', 'Sepal\_Length, Sepal\_Width, Petal\_Length, Petal\_Width' USING PARAMETERS max\_depth=5); --trains an XGBoost model to predict 'species' based on feature measurements  
 => SELECT predict\_xgb\_classifier(Sepal\_Length, Sepal\_Width, Petal\_Length, Petal\_Width USING PARAMETERS model\_name='xgb\_iris') FROM iris2; --returns the predicted 'species' classifications  
 => SELECT predict\_xgb\_classifier\_classes(Sepal\_Length, Sepal\_Width, Petal\_Length, Petal\_Width USING PARAMETERS model\_name='xgb\_iris', classes='virginica, versicolor, setosa', probability\_normalization='softmax') OVER() FROM iris2; --returns the predicted 'species' and the probability for each 'species' class

## Clustering

## K-means

=> SELECT kmeans('agar\_kmeans', 'agar\_dish\_1', '\*' , 5 USING PARAMETERS exclude\_columns='id', output\_view='agar\_view', key\_columns='id'); --trains a K-means model to group microbes into 5 clusters based on x-y coordinates  
 => SELECT id, apply\_kmeans(x, y USING PARAMETERS model\_name='agar\_kmeans') FROM agar\_dish\_2; --applies 'agar\_kmeans' to 'agar\_dish\_2' and assigns each input row to a cluster

**Bisecting k-means**

```
=> SELECT bisecting_kmeans('agar_bkmeans', 'agar_dish_1', '*', 5 USING PARAMETERS exclude_columns='id', key_columns='id', output_view='agar_bk_view');
```

--trains a bisecting k-means model to group microbes into five or fewer clusters based on x-y coordinates

```
=> SELECT id, apply_bisecting_kmeans(x,y USING PARAMETERS model_name='agar_bkmeans', number_clusters=3) FROM agar_dish_2;--applies 'agar_bkmeans' to 'agar_dish_2' and assigns each row to one of three clusters
```

**Time Series****Autoregression**

```
=> SELECT autoregressor('AR_temp', 'temp_data', 'Temperature', 'time' USING PARAMETERS p=3); --trains an autoregressive model to predict 'Temperature' based on the previous three values
```

```
=> SELECT predict_autoregressor(Temperature USING PARAMETERS model_name='AR_temp', npredictions=20) OVER(ORDER BY time) FROM temp_data; --returns 20 predictions beginning at the end of 'temp_data'
```

**Moving-average**

```
=> SELECT moving_average('MA_temp', 'temp_data', 'temperature', 'time' USING PARAMETERS q=3); --trains a moving-average model to predict 'Temperature' based on the errors of the previous three predictions
```

```
=> SELECT predict_moving_average(Temperature USING PARAMETERS model_name='MA_temp', npredictions=20) OVER(ORDER BY time) FROM temp_data; --returns 20 predictions beginning at the end of 'temp_data'
```

**Evaluating model performance****Regression metrics****Mean Squared Error**

```
=> SELECT mse(obs, pred) OVER() FROM (SELECT eruptions AS obs, predict_linear_reg(waiting USING PARAMETERS model_name='linear_reg_faithful') AS pred FROM faithful_testing) AS prediction_output;
```

--returns the mean squared error between the predictions from 'linear\_reg\_faithful' and the ground truth values in 'faithful\_testing'

**R-squared**

```
=> SELECT rsquared(obs, pred) OVER() FROM (SELECT eruptions AS obs, predict_linear_reg(waiting USING PARAMETERS model_name='linear_reg_faithful') AS pred FROM faithful_testing) AS prediction_output;
```

--returns the R-squared value between the predictions from 'linear\_reg\_faithful' and the ground truth values in 'faithful\_testing'

**Classification metrics****Confusion Matrix**

```
=> SELECT confusion_matrix(obs::int, pred::int USING PARAMETERS num_classes=2) OVER() FROM (SELECT am AS obs, predict_logistic_reg(mpg, cyl, disp, hp, drat, wt, qsec, gear, carb USING PARAMETERS model_name='logistic_cars') AS pred FROM mtcars_test) AS prediction_output; --computes a confusion matrix for the 'logistic_cars' model using the predictions and the observed values in 'mtcars_test'
```

**Error Rate**

```
=> SELECT error_rate(obs::int, pred::int USING PARAMETERS num_classes=2) OVER() FROM (SELECT am AS obs, predict_logistic_reg(mpg, cyl, disp, hp, drat, wt, qsec, gear, carb USING PARAMETERS model_name='logistic_cars', num_classes=2) AS pred FROM mtcars_test) AS prediction_output;
```

--calculates 'logistic\_cars' rate of incorrect classifications for each class as well as the total error rate across all classes

**Lift Table**

```
=> SELECT lift_table(obs::int, prob USING PARAMETERS num_bins=50) OVER() FROM (SELECT am AS obs, predict_logistic_reg(mpg, cyl, disp, hp, drat, wt, qsec, gear, carb USING PARAMETERS model_name='logistic_cars', type='probability') AS prob FROM mtcars_test) AS prediction_output; --returns a lift chart for the 'logistic_cars' model on the 'mtcars_test' dataset
```

**ROC**

```
=> SELECT roc(obs::int, prob USING PARAMETERS num_bins=50) OVER() FROM (SELECT am AS obs, predict_logistic_reg(mpg, cyl, disp, hp, drat, wt, qsec, gear, carb USING PARAMETERS model_name='logistic_cars', type='probability') AS prob FROM mtcars_test) AS prediction_output; --returns a table containing points on an ROC curve for the 'logistic_cars' model and 'mtcars_test' dataset
```

**Cross Validation**

```
=> SELECT cross_validate('svm_cars', 'mtcars_train', 'am', 'mpg,cyl,disp,hp,drat,wt,qsec,vs,gear,carb' USING PARAMETERS cv_fold_count= 5, cv_hyperparams='{\"C\":[0.1,1,5]}', cv_model_name='svm_cv', cv_metrics='error_rate'); --performs k-fold cross validation for an SVM classifier on 'mtcars_train' data -- evaluates 15 models (5 folds and 3 values of hyperparameter 'C') and saves the result in 'svm_cv'
```

**PRC**

```
=> SELECT prc(obs::int, prob::float USING PARAMETERS num_bins=50, f1_score=true) OVER() FROM (SELECT am AS obs, predict_logistic_reg(mpg, cyl, disp, hp, drat, wt, qsec, gear, carb USING PARAMETERS model_name='logistic_cars', type='probability') AS prob FROM mtcars_test) AS prediction_output; --returns a table containing points on a precision recall curve (PRC) for the 'logistic_cars' model and 'mtcars_test' dataset
```

**Decision tree metrics**

```
=> SELECT read_tree(USING PARAMETERS model_name='rf_iris', format='tabular'); --returns information for decision trees in the 'rf_iris' model, such as 'is_leaf' and 'node_depth'
```

```
=> SELECT rf_predictor_importance( USING PARAMETERS model_name = 'rf_iris'); --measures the importance of each predictor in the 'rf_iris' model using the Mean Decrease Impurity (MDI) method
```

## Managing models

### List models

=> SELECT \* FROM models; --returns information about all models in the database

### Delete a model

=> DROP MODEL xgb\_iris; --removes the model from the database

### Upgrade a model

=> SELECT upgrade\_model( USING PARAMETERS model\_name = 'xgb\_cars'); --upgrades 'xgb\_cars' to the Vertica version currently in use by the database

### Change model owner, schema, and name

=> ALTER MODEL xgb\_cars OWNER TO analyst;

=> ALTER MODEL xgb\_cars SET SCHEMA public;

=> ALTER MODEL xgb\_cars RENAME to xgb\_autos;

### Change model privileges

=> GRANT USAGE ON MODEL linear\_reg\_faithful TO analyst; --allows the user 'analyst' to run the 'linear\_reg\_faithful' model

=> REVOKE USAGE ON MODEL linear\_reg\_faithful FROM analyst; --removes the usage privileges for 'analyst' on the 'linear\_reg\_faithful' model

### Read model attributes

=> SELECT get\_model\_summary(USING PARAMETERS model\_name='linear\_reg\_faithful'); --displays a summary of the model 'linear\_reg\_faithful'

=> SELECT get\_model\_attribute(USING PARAMETERS model\_name='linear\_reg\_faithful'); --lists all attributes of the model (e.g. the type of regularization used, number of iterations, etc.)

=> SELECT get\_model\_attribute(USING PARAMETERS model\_name='linear\_reg\_faithful', attr\_name='data'); --returns the value of attribute 'data' in the model 'linear\_reg\_faithful'

### Import/export models to other Vertica clusters

=> SELECT export\_models('/home/dbadmin/myModels', 'AR\_temp'); --exports model 'AR\_temp' to directory 'myModels'

=> SELECT import\_models('/home/newDir/myModels/\*' USING PARAMETERS new\_schema='analyst'); --imports all models under 'myModels' to schema 'analyst'

## Using external models

--As the following examples are dependent on external models, they will not work out of the box. For more information on using external models in Vertica, see the [Vertica documentation](#) or [github TensorFlow example](#).

### PMML models

=> SELECT import\_models('/data/username/temp/spark\_logistic\_reg' USING PARAMETERS category='pmml'); --imports a PMML model named 'spark\_logistic\_reg'

=> SELECT SELECT predict\_pmml(\* USING PARAMETERS model\_name='spark\_logistic\_reg') AS prediction FROM test\_data; --makes predictions on the 'test\_data' using the 'spark\_logistic\_reg' model

=> SELECT export\_models('/path/to/export/to', 'spark\_logistic\_reg' USING PARAMETERS category='pmml'); --exports 'spark\_logistic\_reg' in PMML format

### TensorFlow models

=> SELECT import\_models('/path/tf\_models/tf\_mnist\_keras' USING PARAMETERS category='tensorflow'); --imports a TensorFlow model named 'tf\_mnist\_keras'

=> SELECT prdict\_tensorflow(\* USING PARAMETERS model\_name='tf\_mnist\_keras') OVER(PARTITION BEST) FROM tf\_mnist\_test\_images; --makes predictions on the 'tf\_mnist\_test\_images' data with the 'tf\_mnist\_keras' model

=> SELECT export\_models('/path/to/export/to', 'tf\_mnist\_keras' USING PARAMETERS category='tensorflow'); --exports 'tf\_mnist\_keras' as a frozen graph