# **Lab #4**: Logistic Regression_kNN

This lab is to continous dealing with Logistic Regression, kNN, and Decision Tree algorithms applied to classification tasks.

=================================================================

## LogisticRegression

**Build a model:**

```
from sklearn.linear_model import LogisticRegression
classifier = LogisticRegression(random_state = 0)
classifier.fit(Xtrain, ytrain)
```

**Test the model:**

```
y_pred = classifier.predict(Xtest)
```

**Confustion matrix:**

```
from sklearn.metrics import confusion_matrix
cm = confusion_matrix(ytest, y_pred)
```

**Accuracy**:

```
from sklearn.metrics import accuracy_score
print ("Accuracy : ", accuracy_score(ytest, y_pred))
```

**Some additional values on metrics:**

- metrics.confusion_matrix(y_test, y_pred[, ...]): Compute confusion matrix to evaluate the accuracy of a classification
- metrics.precision_score(y_test, y_pred[, ...]) Compute the precision
- metrics.recall_score(y_test, y_pred[, ...]): Compute the recall
- metrics.f1_score(y_test, y_pred[, labels, ...]): Compute the F1 score, also known as balanced F-score or F-measure
- metrics.accuracy_score(y_test, y_pred[, ...]): Accuracy classification score

=================================================================

**kNN algorithm** (https://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html):

**Syntax:**

*class* sklearn.neighbors.KNeighborsClassifier(*n_neighbors=5*, *, *weights='uniform'*, *algorithm='auto'*, *leaf_size=30*, *p=2*, *metric='minkowski'*, *metric_params=None*, *n_jobs=None*)

**where,**

**n_neighbors:** int, default=5. Number of neighbors to use by default for `kneighbors` queries.
**weights:**{'uniform', 'distance'} or callable, default='uniform'. Weight function used in prediction. Possible values:

- '**uniform**' : uniform weights. All points in each neighborhood are weighted equally.
- '**distance**' : weight points by the inverse of their distance. in this case, closer neighbors of a query point will have a greater influence than neighbors which are further away.
- [**callable**] : a user-defined function which accepts an array of distances, and returns an array of the same shape containing the weights.

**algorithm:** {'auto', 'ball_tree', 'kd_tree', 'brute'}, default='auto'. Algorithm used to compute the nearest neighbors:

- 'ball_tree' will use `BallTree`
- 'kd_tree' will use `KDTree`
- 'brute' will use a brute-force search.
- 'auto' will attempt to decide the most appropriate algorithm based on the values passed to `fit` method.

Note: fitting on sparse input will override the setting of this parameter, using brute force.

**leaf_size:** int, default=30

Leaf size passed to BallTree or KDTree. This can affect the speed of the construction and query, as well as the memory required to store the tree. The optimal value depends on the nature of the problem.

**p:** int, default=2

Power parameter for the Minkowski metric. When $p = 1$, this is equivalent to using manhattan_distance (l1), and euclidean_distance (l2) for $p = 2$. For arbitrary p, minkowski_distance (l_p) is used.

**metric:** str or callable, default='minkowski'

The distance metric to use for the tree. The default metric is minkowski, and with $p=2$ is equivalent to the standard Euclidean metric. For a list of available metrics, see the documentation of `DistanceMetric`. If metric is "precomputed", X is assumed to be a distance matrix and must be square during fit. X may be a sparse graph, in which case only "nonzero" elements may be considered neighbors.

**metric_params:** dict, default=None

>Additional keyword arguments for the metric function.

**n_jobs:** int, default=None

>The number of parallel jobs to run for neighbors search. `None` means 1 unless in a `joblib.parallel_backend` context. `-1` means using all processors. See Glossary for more details. Doesn't affect `fit` method.

**Usage:**

```
from sklearn.neighbors import KNeighborsClassifier
model = KNeighborsClassifier(n_neighbors=3)
# Train the model using the training sets
model.fit(X_train,y_train)
#Predict Output
y_pred = model.predict(X_test)
```

======================================================================

**Decision Tree** (https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html)**:**

**Syntax***:*

*class* sklearn.tree.DecisionTreeClassifier(*, *criterion='gini'*, *splitter='best'*, *max_depth=None*, *min_samples_split=2*, *min_samples_leaf=1*, *min_weight_fraction_leaf=0.0*, *max_features=None*, *random_state=None*, *max_leaf_nodes=None*, *min_impurity_decrease=0.0*, *class_weight=None*, *ccp_alpha=0.0*)

**where**,

**criterion:**{"gini", "entropy"}, default="gini"

>The function to measure the quality of a split. Supported criteria are "gini" for the Gini impurity and "entropy" for the information gain.

**splitter:** {"best", "random"}, default="best"

>The strategy used to choose the split at each node. Supported strategies are "best" to choose the best split and "random" to choose the best random split.

**max_depth:** int, default=None

>The maximum depth of the tree. If None, then nodes are expanded until all leaves are pure or until all leaves contain less than min_samples_split samples.

**min_samples_split:** int or float, default=2

>The minimum number of samples required to split an internal node:

- If int, then consider `min_samples_split` as the minimum number.
- If float, then `min_samples_split` is a fraction and `ceil(min_samples_split * n_samples)` are the minimum number of samples for each split.

**min_samples_leaf:** int or float, default=1

The minimum number of samples required to be at a leaf node. A split point at any depth will only be considered if it leaves at least `min_samples_leaf` training samples in each of the left and right branches. This may have the effect of smoothing the model, especially in regression.

- If int, then consider `min_samples_leaf` as the minimum number.
- If float, then `min_samples_leaf` is a fraction and `ceil(min_samples_leaf * n_samples)` are the minimum number of samples for each node.

**min_weight_fraction_leaf:** float, default=0.0

The minimum weighted fraction of the sum total of weights (of all the input samples) required to be at a leaf node. Samples have equal weight when sample_weight is not provided.

**max_features:** int, float or {"auto", "sqrt", "log2"}, default=None

The number of features to consider when looking for the best split:

- If int, then consider `max_features` features at each split.
- If float, then `max_features` is a fraction and `int(max_features * n_features)` features are considered at each split.
- If "auto", then `max_features=sqrt(n_features)`.
- If "sqrt", then `max_features=sqrt(n_features)`.
- If "log2", then `max_features=log2(n_features)`.
- If None, then `max_features=n_features`.

Note: the search for a split does not stop until at least one valid partition of the node samples is found, even if it requires to effectively inspect more than `max_features` features.

**random_state:** int, RandomState instance or None, default=None

Controls the randomness of the estimator. The features are always randomly permuted at each split, even if `splitter` is set to `"best"`. When `max_features < n_features`, the algorithm will select `max_features` at random at each split before finding the best split among them. But the best found split may vary across different runs, even if `max_features=n_features`. That is the case, if the improvement of the criterion is identical for several splits and one split has to be selected at random. To obtain a deterministic behaviour during fitting, `random_state` has to be fixed to an integer.

**max_leaf_nodes:** int, default=None

Grow a tree with `max_leaf_nodes` in best-first fashion. Best nodes are defined as relative reduction in impurity. If None then unlimited number of leaf nodes.

**min_impurity_decrease:** float, default=0.0

A node will be split if this split induces a decrease of the impurity greater than or equal to this value.

The weighted impurity decrease equation is the following:

```
N_t / N * (impurity - N_t_R / N_t * right_impurity
                    - N_t_L / N_t * left_impurity)
```

where `N` is the total number of samples, `N_t` is the number of samples at the current node, `N_t_L` is the number of samples in the left child, and `N_t_R` is the number of samples in the right child.

`N`, `N_t`, `N_t_R` and `N_t_L` all refer to the weighted sum, if `sample_weight` is passed.

**class_weight:** dict, list of dict or "balanced", default=None

Weights associated with classes in the form `{class_label: weight}`. If None, all classes are supposed to have weight one. For multi-output problems, a list of dicts can be provided in the same order as the columns of y.

Note that for multioutput (including multilabel) weights should be defined for each class of every column in its own dict. For example, for four-class multilabel classification weights should be [{0: 1, 1: 1}, {0: 1, 1: 5}, {0: 1, 1: 1}, {0: 1, 1: 1}] instead of [{1:1}, {2:5}, {3:1}, {4:1}].

The "balanced" mode uses the values of y to automatically adjust weights inversely proportional to class frequencies in the input data as `n_samples / (n_classes * np.bincount(y))`

For multi-output, the weights of each column of y will be multiplied.

Note that these weights will be multiplied with sample_weight (passed through the fit method) if sample_weight is specified.

**ccp_alpha:** non-negative float, default=0.0

Complexity parameter used for Minimal Cost-Complexity Pruning. The subtree with the largest cost complexity that is smaller than `ccp_alpha` will be chosen. By default, no pruning is performed.

**Usage:**

```
from sklearn.tree import DecisionTreeClassifier
clf_model = DecisionTreeClassifier(criterion="gini", random_state=42,
max_depth=3, min_samples_leaf=5)
clf_model.fit(X_train,y_train)
```

```
# Plot decision tree
tree.plot_tree(clf)
# Predict X_test
y_predict = clf_model.predict(X_test)
```