# Multioutput Classification in Machine Learning

| ⊙ Type | @datasciencebrain |
| --- | --- |

## Introduction

Multioutput classification is a machine learning problem where the goal is to predict multiple target variables from a single input. In this scenario, instead of predicting a single target variable (as in traditional classification), we aim to predict a set of targets for each input data instance. Each output variable can either be treated as an independent classification task or as a correlated task.

This topic extends the concepts of both **multiclass** and **multi-label** classification to the situation where more than one output needs to be predicted.

## Key Differences

1. **Multiclass Classification**: A single target variable has multiple possible classes, e.g., categorizing images of animals into classes like `cat`, `dog`, or `rabbit`.

2. **Multi-label Classification**: A single target variable can have multiple labels simultaneously, e.g., categorizing an email into "spam", "promotion", or "urgent" — it can belong to multiple categories at once.

3. **Multioutput Classification**: Multiple independent or dependent target variables are predicted, where each target may have different classes. For example, predicting both the age and income of a person based on their features like age, gender, and profession.

## Types of Multioutput Classification Problems

1. **Independent Multioutput Classification**: Each target is treated as a separate binary or multiclass classification task. For example, predicting whether a person has a specific disease and their income range might be two independent tasks.

2. **Multivariate Multioutput Classification**: Targets are related or dependent on each other. For example, predicting multiple medical conditions where some conditions may be dependent on others.

# Techniques for Multioutput Classification

There are two main approaches to solving multioutput classification problems:

1. **Problem Transformation Methods**

2. **Algorithm Adaptation Methods**

## 1. Problem Transformation Methods

These methods convert the multioutput classification problem into multiple independent classification problems that can be handled using any standard machine learning algorithm.

## Binary Relevance (BR)

In binary relevance, each output variable is treated as an independent binary classification task. If the output variables are independent of each other, this approach works well. For each target, a separate binary classifier is trained, and each classifier predicts whether its corresponding label is true or false.

## Example Code for Binary Relevance

Here, we will use the `scikit-learn` library to demonstrate binary relevance with the Random Forest classifier.

```
import numpy as np
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import train_test_split
from sklearn.multioutput import MultiOutputClassifier
from sklearn.datasets import make_multilabel_classification
```

```python
from sklearn.metrics import accuracy_score

# Generating a synthetic multi-output classification dataset
X, y = make_multilabel_classification(n_samples=1000, n_features=20, n_classes=3, n_labels=2, random_state=42)

# Splitting data into training and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)

# Training RandomForest for multi-output classification using Binary Relevance
rf = RandomForestClassifier(n_estimators=100, random_state=42)
multi_target_rf = MultiOutputClassifier(rf, n_jobs=-1)
multi_target_rf.fit(X_train, y_train)

# Making predictions
y_pred = multi_target_rf.predict(X_test)

# Evaluating the model
print("Accuracy:", accuracy_score(y_test, y_pred))
```

**Explanation**:

- `make_multilabel_classification` : Generates synthetic data with multiple output labels.

- `MultiOutputClassifier` : Applies binary relevance, creating one classifier per output label.

- The model is trained on the training data, and predictions are made on the test set.

## Classifier Chains

Classifier chains are an extension of binary relevance. Instead of training independent classifiers, classifier chains use the predictions from previous classifiers as additional features for subsequent classifiers. This approach accounts for potential dependencies between output variables.

## Example Code for Classifier Chains

```python
from sklearn.multioutput import MultiOutputClassifier
from sklearn.svm import SVC
from sklearn.model_selection import train_test_split
from sklearn.datasets import make_multilabel_classification
from sklearn.metrics import accuracy_score

# Generating synthetic data for multi-output classification
X, y = make_multilabel_classification(n_samples=1000, n_features=20, n_classes=3, n_labels=2, random_state=42)

# Splitting data into training and test sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)

# Training a Classifier Chain using Support Vector Classifiers (SVC)
chain_svc = MultiOutputClassifier(SVC(), n_jobs=-1)
chain_svc.fit(X_train, y_train)

# Making predictions
y_pred = chain_svc.predict(X_test)

# Evaluating the model
print("Accuracy:", accuracy_score(y_test, y_pred))
```

**Explanation**:

- Classifier chains are trained sequentially, where each model uses the predictions from the previous model as input.

- The method can handle dependencies between the output labels more effectively than binary relevance.

## Label Powerset (LP)

In label powerset, each unique combination of labels in the training set is treated as a single label in a multiclass classification problem. This can lead to a large

number of classes if the number of output labels is high.

## Example Code for Label Powerset

```python
from sklearn.multioutput import MultiOutputClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import train_test_split
from sklearn.datasets import make_multilabel_classification
from sklearn.metrics import accuracy_score

# Generating synthetic data
X, y = make_multilabel_classification(n_samples=1000, n_features=20, n_classes=3, n_labels=2, random_state=42)

# Splitting data
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)

# Using RandomForest with MultiOutputClassifier to perform Label Powerset
rf_model = MultiOutputClassifier(RandomForestClassifier(), n_jobs=-1)
rf_model.fit(X_train, y_train)

# Predictions and evaluation
y_pred = rf_model.predict(X_test)
print("Accuracy:", accuracy_score(y_test, y_pred))
```

**Explanation**:

- Label powerset method treats each combination of labels as a unique class. However, this could become computationally expensive for problems with many classes.

## 2. Algorithm Adaptation Methods

These methods modify existing algorithms to directly handle multioutput problems.

## Decision Trees

Decision trees can be extended for multioutput classification. In this case, a tree is built to predict multiple outputs simultaneously. Each node of the tree can predict multiple labels, and the decision tree will try to minimize the error for all outputs simultaneously.

## Example Code for Decision Trees

```python
from sklearn.tree import DecisionTreeClassifier
from sklearn.datasets import make_multilabel_classification
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score

# Generating synthetic multi-output classification data
X, y = make_multilabel_classification(n_samples=1000, n_features=20, n_classes=3, n_labels=2, random_state=42)

# Splitting data
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)

# Using Decision Tree Classifier for Multioutput Classification
tree_clf = DecisionTreeClassifier(random_state=42)
multi_output_tree = MultiOutputClassifier(tree_clf, n_jobs=-1)
multi_output_tree.fit(X_train, y_train)

# Making predictions
y_pred = multi_output_tree.predict(X_test)

# Evaluating performance
print("Accuracy:", accuracy_score(y_test, y_pred))
```

# Neural Networks

Neural networks can be adapted to multioutput classification by having multiple output units in the final layer, one for each target variable. Each output unit will represent a class for its respective output variable.

## Example Code for Neural Networks

```
import tensorflow as tf
from sklearn.datasets import make_multilabel_classification
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import LabelBinarizer
from sklearn.metrics import accuracy_score

# Generating synthetic data
X, y = make_multilabel_classification(n_samples=1000, n_features=20, n_classes=3, n_labels=2, random_state=42)

# Splitting data
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)

# Define a neural network model for multioutput classification
model = tf.keras.Sequential([
    tf.keras.layers.Dense(64, activation='relu', input_shape=(X_train.shape[1],)),
    tf.keras.layers.Dense(32, activation='relu'),
    tf.keras.layers.Dense(y_train.shape[1], activation='sigmoid')  # Output layer
with as many neurons as output variables
])

model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])

# Training the model
model.fit(X_train, y_train, epochs=10, batch_size=32)

# Making predictions
y_pred = (model.predict(X_test) > 0.5).astype(int)  # Threshold predictions
```

```
# Evaluating performance
print("Accuracy:", accuracy_score(y_test, y_pred))
```

**Explanation**:

- The model has multiple output neurons (one for each class of each output variable).

- The `sigmoid` activation function is used since we are dealing with binary classification problems for each output.

- Binary cross-entropy is used as the loss function for each output.

# Evaluation Metrics for Multioutput Classification

When evaluating a multioutput classification model, standard metrics like **accuracy** might not suffice. Here are some common evaluation metrics:

- **Hamming Loss**: This measures the fraction of incorrect labels over all predictions.

  $$\text{Hamming Loss} = \frac{1}{N \times L} \sum_{i=1}^{N} \sum_{j=1}^{L} 1(y_{ij} \neq \hat{y}_{ij})$$

  where N is the number of samples and LL is the number of output labels.

- **Subset Accuracy**: This metric measures the number of samples for which all the output labels are correctly predicted.

- **Macro and Micro Average**: These averages combine performance across all outputs:

  - **Macro Average**: Calculate the metric for each output and then take the average.

  - **Micro Average**: Aggregate the contributions of all outputs and then compute the metric.

# References

- Scikit-learn documentation on <u>multi-output classification</u>

- "Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow" by Aurélien Géron