# CERN SUMMER STUDENT FINAL PRESENTATION

# Cross Validation Improvements in TMVA

Mohammad Uzair

**Supervisors:**

Dr. Lorenzo MONETA

Mr. Hans Kim ALBERTSSON BRANN

# About Me:

- **Name** : Mohammad Uzair
- **Country** : Pakistan 🇵🇰
- **Studies** : Final Year of bachelor's degree in Computer Science
- **University** : National University of Sciences and Technology (NUST)
- **Interested in** : Artificial Intelligence and Machine Learning in Computer Vision
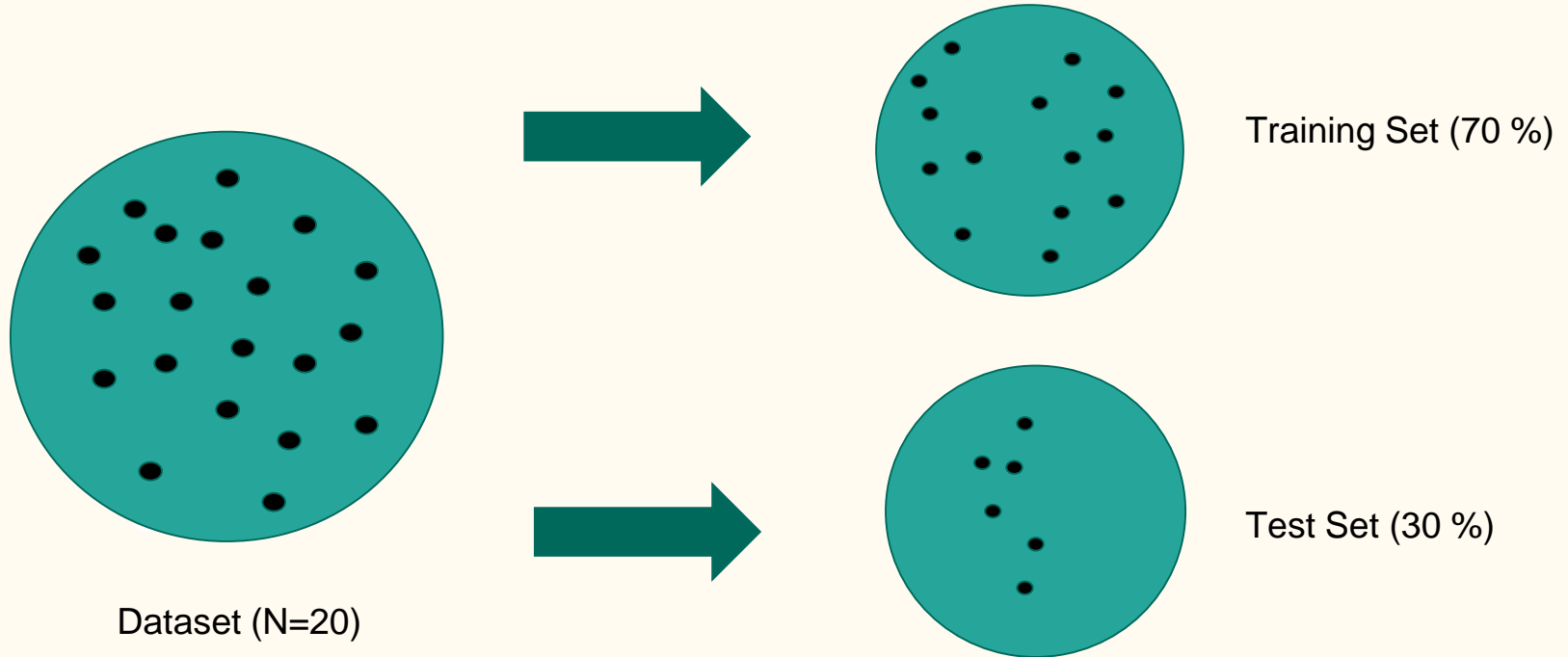
# Tasks Assigned:

Improve Cross Validation in TMVA

- New tutorial introducing Cross Validation in TMVA and make it available through SWAN.

- Improve presentation of tutorial by extending the plotting functionality (Introducing the feature to draw an average ROC Curve).

- Improve CV fold generation targeting unbalanced datasets (Introducing  the feature of Stratified Splitting)

# Validation in Machine Learning

- Divide the dataset into Train and Test sets.

- Tries to estimate the expected error of the model.

- Allows us to get an honest assessment of the trained model.

# Validation in Machine Learning



Training Set (70 %)

Test Set (30 %)

Dataset (N=20)

# Drawbacks of Validation

- Test errors can be highly variable depending on how much data we use for each set.
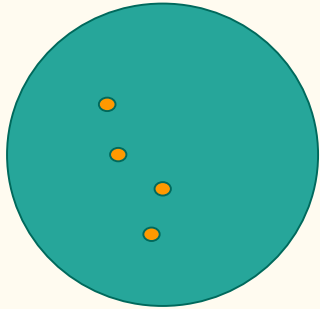- Model Developed on only a subset of data. Ideally we want 100% data for training and 100% for validation.

# K-Fold Cross Validation



Training Set

Validation Fold 1

Dataset (N=20)

# K-Fold Cross Validation



Training Set

Validation Fold 2

Dataset (N=20)

# K-Fold Cross Validation



Dataset (N=20)

Training Set

Validation Fold 3

**And so on….**

# Tutorial Introducing TMVA Cross Validation

—

# Tutorial introducing TMVA Cross Validation

- The tutorial was outdated.

- Was not easy to understand for new users.

- A new and updated tutorial in python jupyter notebook with proper explanation so that it is easy to use and understand and make it available through SWAN.

# Python Jupyter Notebook for basic Tutorial on TMVA Cross Validation

# Python Jupyter Notebook for basic Tutorial on TMVA Cross Validation

# Extending the plotting functionality

# Improve the presentation of tutorial

- Currently, we can only show ROC's of individual folds.

- Often, we can be interested in the average behaviour.

- An added feature of visualising the average ROC Curve.

- Addition of this feature in the tutorial to improve the presentation.

# Average ROC Curve

## With drawFolds = False



## With drawFolds = True

# Improve CV fold generation

# Random Splitting VS Stratified Splitting

- Determines the distribution of input data.

- Random Splitting just randomly splits the data equally.

- Data of a fold can be distributed differently than the whole.

- Problems, in particular, arise with unbalanced classes
  (can more easily occur in multi-class classification).

- Stratified Splitting ensures that each fold follows the same
  distribution as the whole.

# Random Splitting VS Stratified Splitting

- Random Splitting just randomly splits the data equally



30 samples
03 Classes
- 10 samples
- 15 samples
- 05 samples

6 random samples in each fold



| Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5 |

# Random Splitting VS Stratified Splitting

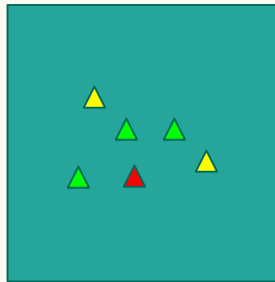- Stratified Splitting randomly splits data ensuring that each fold is good representative of the whole.
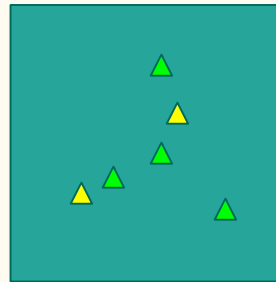


30 samples
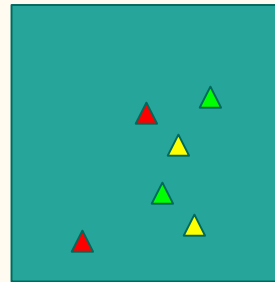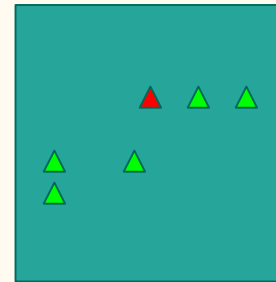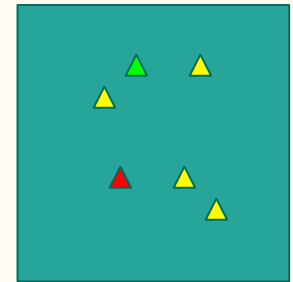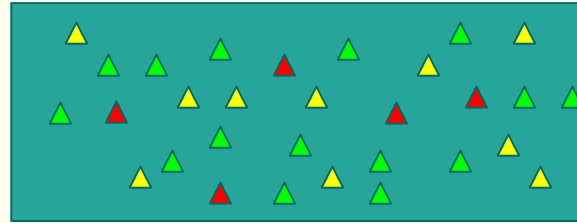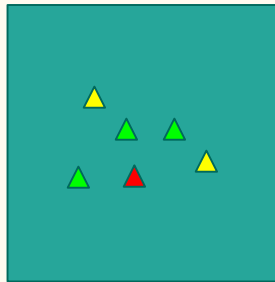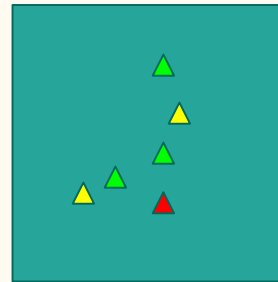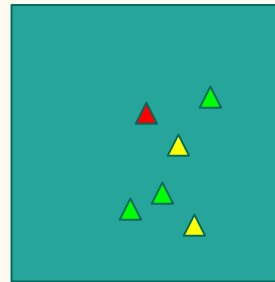03 Classes
△ 10 samples
▲ 15 samples
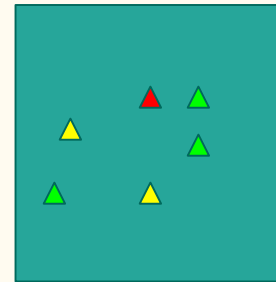▲ 05 samples

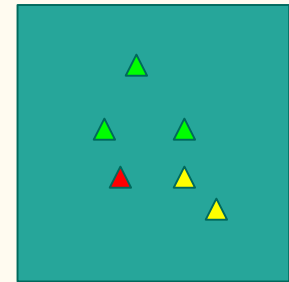6 random samples in each fold

| Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5 |

# Stratified Splitting in TMVA

```cpp
TMVA::DataLoader *d = new TMVA::DataLoader("dataset");

d->AddSignalTree(std::get<1>(data_class0));
d->AddBackgroundTree(std::get<1>(data_class1));

d->AddVariable("x", 'D');
d->AddSpectator("id", "id", "");
d->PrepareTrainingAndTestTree(
    "", Form("SplitMode=Block:nTrain_Signal=%i:nTrain_Background=%i:!V", nPointsSig, nPointsBkg));

d->GetDataSetInfo().GetDataSet(); // Force creation of dataset.
```

```cpp
//For Random Splitting
TMVA::CvSplitKFolds split{NUM_FOLDS, "", kFALSE, 0};
   d->MakeKFoldDataSet(split);
```

```cpp
//For Stratified Splitting
TMVA::CvSplitKFolds split1{NUM_FOLDS, "", kTRUE, 0};
   d->MakeKFoldDataSet(split1);
```

# Future Improvements:

- Weighted kFolds Splitting

- Investigate feasibility of integrating TMVA GUI with tutorials

- Investigate feasibility of adding validation set (this would mainly be beneficial for methods with a large number of parameters, e.g. DNN and BDT's)

**Special Thanks to my Supervisors:**

Dr. Lorenzo MONETA

Mr. Hans Kim ALBERTSSON BRANN