# Ensemble methods

## TP 5 – Visualizing several algorithms

The objective of this TP is to visualize the performance of several algorithms, using the same dataset as input. We will use iris database to perform classification using decision trees, AdaBoost, Random Forest and Extra trees, using scikit learn implementation.

1. Import libraries

```python
import matplotlib.pyplot as plt
import numpy as np
from matplotlib.colors import ListedColormap

from sklearn.datasets import load_iris
from sklearn.ensemble import (
    AdaBoostClassifier,
    ExtraTreesClassifier,
    RandomForestClassifier,
)
from sklearn.tree import DecisionTreeClassifier
```

2. Prepare dataset, hyperparameters and algorithms to be used.

```python
# Parameters
n_classes = 3 #Number of clases to detect
n_estimators = 30 #number of trees to train
max_depth=3 #maximum depth of each tree to train
RANDOM_SEED = 13  # fix the seed on each iteration

#parameter for plotting
cmap = plt.cm.RdYlBu
plot_step = 0.02  # fine step width for decision surface contours
plot_step_coarser = 0.5  # step widths for coarse classifier guesses
plot_idx = 1

# Load data
iris = load_iris()


#defining models to train
models = [
    DecisionTreeClassifier(max_depth=None),
    RandomForestClassifier(n_estimators=n_estimators, max_depth=max_depth),
    ExtraTreesClassifier(n_estimators=n_estimators, max_depth=max_depth),
    AdaBoostClassifier(DecisionTreeClassifier(max_depth=max_depth), n_estimators=n_estimators),
]
```

3. Run all algorithms and plot results

```python
for pair in ([0, 1], [0, 2], [2, 3]):
    for model in models:
        # We only take the two corresponding features
        X = iris.data[:, pair]
        y = iris.target

        # Shuffle
        idx = np.arange(X.shape[0])
        np.random.seed(RANDOM_SEED)
        np.random.shuffle(idx)
        X = X[idx]
        y = y[idx]

        # Standardize
        mean = X.mean(axis=0)
        std = X.std(axis=0)
        X = (X - mean) / std

        # Train
        model.fit(X, y)

        scores = model.score(X, y)
        # Create a title for each column and the console by using str() and
        # slicing away useless parts of the string
        model_title = str(type(model)).split(".")[-1][:-2][: -len("Classifier")]

        model_details = model_title
        if hasattr(model, "estimators_"):
            model_details += " with {} estimators".format(len(model.estimators_))
        print(model_details + " with features", pair, "has a score of", scores)

        plt.subplot(3, 4, plot_idx)
        if plot_idx <= len(models):
            # Add a title at the top of each column
            plt.title(model_title, fontsize=9)

        # Now plot the decision boundary using a fine mesh as input to a
        # filled contour plot
        x_min, x_max = X[:, 0].min() - 1, X[:, 0].max() + 1
        y_min, y_max = X[:, 1].min() - 1, X[:, 1].max() + 1
        xx, yy = np.meshgrid(
            np.arange(x_min, x_max, plot_step), np.arange(y_min, y_max, plot_step)
        )
```

```python
            # Plot either a single DecisionTreeClassifier or alpha blend the
            # decision surfaces of the ensemble of classifiers
            if isinstance(model, DecisionTreeClassifier):
                Z = model.predict(np.c_[xx.ravel(), yy.ravel()])
                Z = Z.reshape(xx.shape)
                cs = plt.contourf(xx, yy, Z, cmap=cmap)
            else:
                # Choose alpha blend level with respect to the number
                # of estimators
                # that are in use (noting that AdaBoost can use fewer estimators
                # than its maximum if it achieves a good enough fit early on)
                estimator_alpha = 1.0 / len(model.estimators_)
                for tree in model.estimators_:
                    Z = tree.predict(np.c_[xx.ravel(), yy.ravel()])
                    Z = Z.reshape(xx.shape)
                    cs = plt.contourf(xx, yy, Z, alpha=estimator_alpha, cmap=cmap)

            # Build a coarser grid to plot a set of ensemble classifications
            # to show how these are different to what we see in the decision
            # surfaces. These points are regularly space and do not have a
            # black outline
            xx_coarser, yy_coarser = np.meshgrid(
                np.arange(x_min, x_max, plot_step_coarser),
                np.arange(y_min, y_max, plot_step_coarser),
            )
            Z_points_coarser = model.predict(
                np.c_[xx_coarser.ravel(), yy_coarser.ravel()]
            ).reshape(xx_coarser.shape)
            cs_points = plt.scatter(
                xx_coarser,
                yy_coarser,
                s=15,
                c=Z_points_coarser,
                cmap=cmap,
                edgecolors="none",
            )

            # Plot the training points, these are clustered together and have a
            # black outline
            plt.scatter(
                X[:, 0],
                X[:, 1],
                c=y,
                cmap=ListedColormap(["r", "y", "b"]),
                edgecolor="k",
                s=20,
            )
        plot_idx += 1  # move on to the next plot in sequence
plt.suptitle("Classifiers on feature subsets of the Iris dataset", fontsize=12)
plt.axis("tight")
plt.tight_layout(h_pad=0.5, w_pad=0.5, pad=2.5)
plt.show()
```

- Why the random seed must be the same for all algorithms?
- Change max_depth from 3 to 1 and to 10, what happens? (in terms of performance and speedness)
- Are there any unnecessary hyperparameters?