Modeling data (with decision trees) 1. EDA - exploratory data analysis 2. Build a model 2.1 Define 2.2 Fit 2.3 Predict 2.4 Evalaute 3. Othe models In [20]: import pandas as pd import numpy as np import matplotlib.pyplot as plt import seaborn as sns import sklearn as sk from sklearn import tree from sklearn.tree import DecisionTreeClassifier In [21]: url = 'https://raw.githubusercontent.com/nlihin/data-analytics/main/datasets/iris.csv' iris_df = pd.read_csv(url) iris_df.head() Out[21]: sepal_length sepal_width petal_length petal_width class 0 5.1 3.5 1.4 0.2 Iris-setosa 1 4.9 3.0 1.4 0.2 Iris-setosa 2 0.2 Iris-setosa 4.7 3.2 1.3 0.2 Iris-setosa 3 4.6 3.1 1.5 5.0 0.2 Iris-setosa 4 3.6 1.4 1. EDA Get to know the data How many instances (rows) with each label (flower species)? In [22]: pd.pivot_table(iris_df,index=['class'],values=['sepal_length','sepal_width','petal_length','petal_width'], aggf Out[22]: petal_length petal_width sepal_length sepal_width class Iris-setosa 50.0 50.0 50.0 50.0 Iris-versicolor 50.0 50.0 50.0 50.0 50.0 50.0 Iris-virginica 50.0 50.0 Another way to look at it - species percentage In [23]: plt.figure(1, figsize=(5,5)) plt.title("Distribution of Species") iris df['class'].value counts().plot.pie(autopct="%1.1f%%") #try removing this - what happens? plt.ylabel("") Distribution of Species Iris-versicolor 33.3% Iris-virginica 33.3% 33.3% Iris-setosa In [48]: fig, axes = plt.subplots(2,2,figsize=(12,10)) sns.boxplot(x = 'class', y = 'sepal_length', data = iris_df, palette="gist_ncar_r", ax = axes[0,0]) sns.boxplot(x = 'class', y = 'sepal_width', data = iris_df, palette="gist_ncar_r", ax = axes[0,1]) sns.boxplot(x = 'class', y = 'petal_length', data = iris_df, palette="gist_ncar_r", ax = axes[1,0]) sns.boxplot(x = 'class', y = 'petal_width', data = iris_df, palette="gist_ncar_r", ax = axes[1,1]) 8.0 7.5 4.0 7.0 3.5 3.0 6.5 sepal length 6.0 5.5 2.5 5.0 4.5 2.0 Iris-setosa Iris-versicolor Iris-virginica Iris-setosa Iris-versicolor Iris-virginica dass 2.5 6 2.0 5 petal width petal length 2 0.5 1 0.0 Iris-versicolor Iris-versicolor Iris-setosa Iris-virginica Iris-setosa Iris-virginica dass In [25]: sns.pairplot(iris_df, hue='class') plt.show() sepal_length 5 4.5 4.0 sepal width 3.5 3.0 2.5 2.0 dass Iris-setosa 7 Iris-versicolor 6 Iris-virginica petal length 2.5 2.0 petal width 1.5 1.0 0.5 sepal_length sepal_width petal_length petal_width In [26]: plt.figure(figsize=(7,4)) sns.heatmap(iris_df.corr(),annot=True,cmap='Greens') #draws heatmap with input as the correlation matrix calcu 1.0 -0.11 0.87 0.82 sepal_length - 0.6 -0.11 -0.42 -0.36 sepal_width - 0.4 - 0.2 -0.42 0.87 petal_length - 0.0 -0.20.82 -0.36 0.96 petal_width - -0.4 sepal_length sepal_width petal_length petal_width **Observations** The Sepal Width and Length are not correlated The Petal Width and Length are highly correlated We will use all the features for training the algorithm and check the accuracy. Then we will use 1 Petal Feature and 1 Sepal Feature to check the accuracy of the algorithm as we are using only 2 features that are not correlated. Thus we can have a variance in the dataset which may help in better accuracy. We will check it later. 2. Build a classification model 2.1 Define the model: we will use a decision tree based on the gini index (gini index is the default) In [27]: my model = sk.tree.DecisionTreeClassifier(criterion="gini") 2.2 Fit the model on the data: Define what are the features (x) and what is the target (y) In [28]: features = ['sepal_length', 'sepal_width', 'petal_length', 'petal_width'] X = iris df[features] y = iris_df['class'] Use the model we chose on the features we chose (X and y) In [29]: my_model.fit(X, y) DecisionTreeClassifier() Out[29]: now model1 is a tree, and it has tree features: In [30]: type(my_model) sklearn.tree._classes.DecisionTreeClassifier Out[30]: Let's see what we got: In [31]: text_representation = tree.export_text(my_model) print(text_representation) |---| feature 2 <= 2.45 |--- class: Iris-setosa |--- feature_2 > 2.45 |--- feature 3 <= 1.75 |--- feature 2 <= 4.95 |--- feature 3 <= 1.65 | |--- class: Iris-versicolor |--- feature_3 > 1.65 | |--- class: Iris-virginica |---| feature 2 > 4.95 |--- feature_3 <= 1.55 | |--- class: Iris-virginica |--- feature_3 > 1.55 |--- feature_0 <= 6.95 | |--- class: Iris-versicolor |--- feature_0 > 6.95 |--- class: Iris-virginica -- feature 3 > 1.75 |--- feature 2 <= 4.85 |--- feature_1 <= 3.10 | |--- class: Iris-virginica |--- feature_1 > 3.10 |--- class: Iris-versicolor |---| feature 2 > 4.85 Vizualize: In [32]: fig = plt.figure(figsize=(15,10)) tree.plot_tree(my_model, feature names = features, class names = 'model1.classes ', filled=True, rounded = True) plt.show() petal_length <= 2.45 gini = 0.667 samples = 150 value = [50, 50, 50]class = m $petal_width <= 1.75$ gini = 0.0gini = 0.5 samples = 100 value = [50, 0, 0] class = m value = [0, 50, 50]class = o $petal_length <= 4.85$ $petal_length <= 4.95$ gini = 0.168 samples = 54 gini = 0.043 samples = 46 value = [0, 49, 5] class = 0 value = [0, 1, 45] class = d sepal_width <= 3.1 gini = 0.444 petal_width <= 1.65 gini = 0.041 samples = 48 petal_width <= 1.55 gini = 0.0gini = 0.444 samples = 3 value = [0, 1, 2] samples = 6 value = [0, 0, 43]value = [0, 2, 4] class = d value = [0, 47, 1] class = 0 class = dclass = d $sepal_length <= 6.95$ gini = 0.0qini = 0.0gini = 0.0gini = 0.444 samples = 2samples = 1samples = 47samples = 1samples = 3samples = 3value = [0, 47, 0] class = 0 value = [0, 0, 1] class = d value = [0, 0, 3] class = d value = [0, 0, 2] class = d value = [0, 1, 0] class = 0 value = [0, 2, 1]class = o gini = 0.0samples = 2samples = 1 value = [0, 2, 0]value = [0, 0, 1]That's nice. But how well will it predict on new data? 2.3 Predict Split the data to 70% train set and 30% test In [33]: X_train, X_test, y_train, y_test = sk.model_selection.train_test_split(X, y, test_size=0.3, random_state=1) In [34]: X train.head(2) Out[34]: sepal_length sepal_width petal_length 118 7.7 6.9 2.3 2.6 5.7 1.7 0.3 In [35]: X test.head() Out[35]: sepal_length sepal_width petal_length petal_width 14 5.8 4.0 1.2 0.2 98 5.1 2.5 3.0 1.1 6.6 3.0 4.4 1.4 5.4 1.3 0.4 16 131 7.9 3.8 6.4 2.0 In [36]: y test.head() Iris-setosa Out[36]: Iris-versicolor 75 Iris-versicolor Iris-setosa 131 Iris-virginica Name: class, dtype: object Use the training data to train the model (a new model - model2) In [37]: model2 = sk.tree.DecisionTreeClassifier(criterion="gini") model2.fit(X_train,y_train) DecisionTreeClassifier() Out[37]: Finally - predict: In [38]: prediction2 = model2.predict(X_test) prediction2 array(['Iris-setosa', 'Iris-versicolor', 'Iris-versicolor', 'Iris-setosa', Out[38]: 'Iris-virginica', 'Iris-versicolor', 'Iris-virginica', 'Iris-setosa', 'Iris-setosa', 'Iris-virginica', 'Iris-versicolor', 'Iris-setosa', 'Iris-virginica', 'Iris-versicolor', 'Iris-versicolor', 'Iris-setosa', 'Iris-versicolor', 'Iris-versicolor', 'Iris-setosa', 'Iris-setosa', 'Iris-versicolor', 'Iris-versicolor', 'Iris-virginica', 'Iris-setosa', 'Iris-virginica', 'Iris-versicolor', 'Iris-setosa', 'Iris-setosa', 'Iris-versicolor', 'Iris-virginica', 'Iris-versicolor', 'Iris-virginica', 'Iris-versicolor', 'Iris-virginica', 'Iris-virginica', 'Iris-setosa', 'Iris-versicolor', 'Iris-setosa', 'Iris-versicolor', 'Iris-virginica', 'Iris-virginica', 'Iris-setosa', 'Iris-versicolor', 'Iris-virginica', 'Iris-versicolor'], dtype=object) present in a nicer format: In [39]: df = pd.DataFrame(columns = ['class', 'predictions']) In [40]: df['class'] = y_test df = df.reset_index(drop=True) In [41]: df['predictions'] = pd.Series(prediction2) Out[41]: predictions class 0 Iris-setosa Iris-setosa **1** Iris-versicolor Iris-versicolor **2** Iris-versicolor Iris-versicolor Iris-setosa Iris-setosa Iris-virginica Iris-virginica **5** Iris-versicolor Iris-versicolor 7 Iris-setosa 8 Iris-setosa Iris-setosa Iris-virginica Iris-virginica Iris-versicolor Iris-versicolor 11 Iris-setosa Iris-setosa Iris-virginica Iris-virginica Iris-versicolor Iris-versicolor Iris-versicolor Iris-versicolor 15 Iris-setosa Iris-setosa Iris-versicolor Iris-versicolor Iris-versicolor Iris-versicolor 18 Iris-setosa Iris-setosa 19 Iris-setosa Iris-setosa Iris-versicolor Iris-versicolor Iris-versicolor Iris-versicolor 22 Iris-versicolor Iris-virginica 23 Iris-setosa Iris-setosa Iris-virginica Iris-virginica 25 Iris-versicolor Iris-versicolor 26 Iris-setosa Iris-setosa 27 Iris-setosa Iris-setosa Iris-versicolor Iris-versicolor 29 Iris-virginica Iris-virginica 30 Iris-versicolor Iris-versicolor 31 Iris-virginica Iris-virginica 32 Iris-versicolor Iris-versicolor 33 Iris-virginica Iris-virginica 34 Iris-virginica Iris-virginica 35 Iris-setosa Iris-setosa Iris-versicolor Iris-versicolor 37 Iris-setosa Iris-setosa 38 Iris-versicolor Iris-versicolor 39 Iris-virginica Iris-virginica 40 Iris-virginica Iris-virginica 41 Iris-setosa Iris-setosa Iris-versicolor 42 Iris-virginica 43 Iris-virginica Iris-virginica Iris-versicolor Iris-versicolor 2.4 Evaluate There are various ways to measure accuracy: The most obvious - (number of correct)/(total) In [42]: sk.metrics.accuracy_score(prediction2, y_test) 0.95555555555556 Out[42]: A truth table In [43]: sk.metrics.confusion_matrix(y_test, prediction2) array([[14, 0, 0], Out[43]: [0, 17, 1], [0, 1, 12]], dtype=int64) In [44]: sk.metrics.plot_confusion_matrix(model2, X_test, y_test, cmap=plt.cm.Blues) <sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0x296166d9850> Out[44]: - 16 Iris-setosa 14 - 12 - 10 True label 1 Iris-versicolor 8 - 6 - 4 0 1 Iris-virginica - 2 lris-versicolor lris-virginica Iris-setosa Predicted label Precision, Recall and F1 score - most common. But out of scope for this course In [45]: print(sk.metrics.classification_report(y_test, prediction2, target_names=model2.classes_)) precision recall f1-score support 1.00 1.00 1.00 Iris-setosa 0.94 0.94 0.94 Iris-versicolor 18 Iris-virginica 0.92 0.92 0.92 13 0.96 45 accuracy 0.96 0.96 0.96 45 macro avg weighted avg 0.96 45 0.96 0.96 Precision, Recall and F1 score - most common. But out of scope for this course In [52]: print(sk.metrics.classification_report(y_test, prediction2, target_names=model2.classes_)) precision recall f1-score support 1.00 1.00 Iris-setosa 1.00 0.94 0.94 0.94 Iris-versicolor 18 0.92 0.92 0.92 Iris-virginica 13 0.96 45 accuracy 0.96 0.96 0.96 0.96 45 macro avg 0.96 0.96 45 weighted avg 3. Other models 3.1 Petal vs Sepal Which attributes were the most important in the training data? In [55]: dt = sk.tree.DecisionTreeClassifier(random_state=0) dt.fit(X_train, y_train) pred = dt.predict(X_test) acc = sk.metrics.accuracy_score(pred, y_test) In [56]: dt.feature_importances_ array([0.02146947, 0.02146947, 0.06316954, 0.89389153]) Out[56]: organize it in a nicer format: In [57]: feature_importances = pd.DataFrame(dt.feature_importances_, index = iris df.columns[:4], columns=['importance']).sort_values('importance', ascending=False) feature_importances.head() Out[57]: importance 0.893892 petal_width petal_length 0.063170 sepal_length 0.021469 sepal_width 0.021469 We see that petal is most important. What happens if we only use sepal? In [58]: X s = iris df[['sepal_length','sepal_width']] y_s = iris_df['class'] X_train_s, X_test_s, y_train_s, y_test_s = sk.model_selection.train_test_split(X_s, y_s, test_size=0.3, random_ model_s = sk.tree.DecisionTreeClassifier(criterion="gini") model_s.fit(X_train_s, y_train_s) prediction_s = model_s.predict(X_test_s) print('The accuracy of the Decision Tree is', sk.metrics.accuracy_score(prediction_s, y_test_s)) The accuracy of the Decision Tree is 0.62222222222222 In [59]: sk.metrics.plot_confusion_matrix(model_s, X_test_s, y_test_s, cmap=plt.cm.Blues) <sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0x296174155b0> 12 Iris-setosa 14 - 10 Frue label - 8 1 Iris-versicolor - 6 4 0 6 Iris-virginica 2 Iris-versicolor Iris-virginica Iris-setosa Predicted label In [60]: df2 = pd.DataFrame(columns = ['class', 'predictions']) df2['class'] = y_test_s df2 = df2.reset index(drop=True) df2['predictions'] = pd.Series(prediction s) Now let's try a model with two features: one sepal and one petal: In [61]: #X_s = iris_df[['sepal_length','sepal_width']] X_m = iris_df[['sepal_length','petal width']] y_m = iris_df['class'] X_train_m, X_test_m, y_train_m, y_test_m = sk.model_selection.train_test_split(X_m, y_m, test_size=0.3, random) model mix = sk.tree.DecisionTreeClassifier(criterion="gini", random state=1) model_mix.fit(X_train_m,y_train_m) prediction_mix = model_mix.predict(X_test_m) print('The accuracy of the Decision Tree is', sk.metrics.accuracy score(prediction mix, y test m)) The accuracy of the Decision Tree is 0.955555555555556 3.2 an SVM model In [62]: from sklearn.svm import SVC svm model = SVC(random state = 0)svm model.fit(X train, y train) svm_pred = svm_model.predict(X test) sk.metrics.accuracy score(svm pred, y test) 0.977777777777777 Out[62]: Terminology: features - attributes class/target - the feature we want to predict • A guide to trees - scikit-learn A summary of new functions: • plot.pie(autopct="%1.1f%%") - a pie plot. autopct adds percentages to each of the slices of the pie chart, 1.1 percentages with 1 number after the decimal dot. • DecisionTreeClassifier documentation • my_model = sk.tree.DecisionTreeClassifier(criterion="gini") - decision tree model with the gini index as the node splitting criterion $my_model.fit(X, y)$ - train the model. X are the features and y is the target/class. my_model.feature_importances_ - shows how important is each of the features in model my_model sk.model_selection.train_test_split - splits the data to train and test tree.export_text(my_model) - shows the tree of the model in text format. documentation tree.plot_tree - plots a nice decision tree. documentation Metrics documentation sk.metrics.accuracy_score - computes the model's jaccard accuracy score sk.metrics.confusion_matrix - a truth table sk.metrics.plot_confusion_matrix - a confusion matrix