

```
cells: [
  {
    "cell_type": "markdown",
    "metadata": {},
    "source": [
      "# Module 6: Classification\n",
      "\n",
      "The following tutorial contains Python examples for solving classification problems. You should refer to the Chapters 3 and 4 of the 'Introduction to Data Mining' book for more details. For example, mammals can be identified as warm-blooded vertebrates that give birth to their young. The goals for this tutorial are as follows:\n",
      "1. To provide examples of using different classification techniques from the scikit-learn library package.\n",
      "2. To demonstrate the problem of model overfitting.\n",
      "\n",
      "Read the step-by-step instructions below carefully. To execute the code, click on the corresponding cell and press the SHIFT-ENTER keys simultaneously.\n"
    ]
  },
  {
    "cell_type": "markdown",
    "metadata": {},
    "source": [
      "\n",
      "## 6.1 Vertebrate Dataset\n",
      "\n",
      "We use a variation of the vertebrate data described in Example 3.1 of Chapter 3. Each vertebrate is classified into one of 5 categories: mammals, reptiles, birds, fishes, and amphibians, based on a set of explanatory attributes (predictor variables). Except for 'name', the rest of the attributes have been converted into a 'one hot encoding' binary representation. To illustrate this, we will first load the data into a Pandas DataFrame object and display its content."
    ]
  },
  {
    "cell_type": "code",
    "execution_count": null,
    "metadata": {},
    "outputs": [],
    "source": [
      "import pandas as pd\n",
      "\n",
      "pd.read_csv('vertebrate.csv', header='infer')\n",
      "data"
    ]
  },
  {
    "cell_type": "markdown",
    "metadata": {},
    "source": [
      "\n",
      "Given the limited number of training examples, suppose we convert the problem into a binary classification task (mammals versus non-mammals). We can do so by replacing the class labels of the instances to 'non-mammals' except for those that belong to the 'mammals' class."
    ]
  },
  {
    "cell_type": "code",
    "execution_count": null,
    "metadata": {},
    "outputs": [],
    "source": [
      "data['Class'] = data['Class'].replace(['fishes', 'birds', 'amphibians', 'reptiles'], 'non-mammals')\n",
      "data"
    ]
  },
  {
    "cell_type": "markdown",
    "metadata": {},
    "source": [
      "\n",
      "We can apply Pandas cross-tabulation to examine the relationship between the Warm-blooded and Gives Birth attributes with respect to the class."
    ]
  },
  {
    "cell_type": "code",
    "execution_count": null,
    "metadata": {},
    "outputs": [],
    "source": [
      "pd.crosstab([data['Warm-blooded'], data['Gives Birth']], data['Class'])"
    ]
  },
  {
    "cell_type": "markdown",
    "metadata": {},
    "source": [
      "\n",
      "The results above show that it is possible to distinguish mammals from non-mammals using these two attributes alone since each combination of their attribute values corresponds to a unique class. For example, mammals can be identified as warm-blooded vertebrates that give birth to their young. Such a relationship can also be derived using a decision tree classifier, as shown by the example given in the next subsection."
    ]
  },
  {
    "cell_type": "markdown",
    "metadata": {},
    "source": [
      "\n",
      "## 3.2 Decision Tree Classifier\n",
      "\n",
      "In this section, we apply a decision tree classifier to the vertebrate dataset described in the previous subsection."
    ]
  },
  {
    "cell_type": "code",
    "execution_count": null,
    "metadata": {},
    "outputs": [],
    "source": [
      "from sklearn import tree\n",
      "X = data[['Name', 'Class']]\n",
      "y = data['Class']\n",
      "clf = tree.DecisionTreeClassifier(criterion='entropy', max_depth=3)\n",
      "clf = clf.fit(X, y)"
    ]
  },
  {
    "cell_type": "markdown",
    "metadata": {},
    "source": [
      "\n",
      "The preceding commands will extract the predictor (X) and target class (Y) attributes from the vertebrate dataset and create a decision tree classifier object using entropy as its impurity measure for splitting criterion. The decision tree class in Python sklearn library also supports using 'gini' as impurity measure. The classifier above is also constrained to generate trees with a maximum depth equals to 3. Next, the classifier is trained on the labeled data using the fit() function."
    ]
  },
  {
    "cell_type": "code",
    "execution_count": null,
    "metadata": {},
    "outputs": [],
    "source": [
      "from sklearn import tree\n",
      "X = data[['Name', 'Class']]\n",
      "y = data['Class']\n",
      "clf = tree.DecisionTreeClassifier(criterion='entropy', max_depth=3)\n",
      "clf = clf.fit(X, y)"
    ]
  },
  {
    "cell_type": "code",
    "execution_count": null,
    "metadata": {},
    "outputs": [],
    "source": [
      "dot_data = tree.export_graphviz(clf, feature_names=X.columns, class_names=['mammals', 'non-mammals'], filled=True, \n",
      "                                out_file=None)\n",
      "graph = pydotplus.graph_from_dot_data(dot_data)\n",
      "image = graph.create_png()"
    ]
  },
  {
    "cell_type": "markdown",
    "metadata": {},
    "source": [
      "\n",
      "Next, suppose we apply the decision tree to classify the following test examples."
    ]
  },
  {
    "cell_type": "code",
    "execution_count": null,
    "metadata": {},
    "outputs": [],
    "source": [
      "testData = [[['gila monster', 0, 0, 0, 1, 1, 'non-mammals']],\n",
      "            [['platypus', 1, 0, 0, 0, 1, 1, 'mammals']],\n",
      "            [['owl', 1, 0, 0, 1, 0, 0, 'non-mammals']],\n",
      "            [['dolphin', 1, 1, 1, 0, 0, 0, 'mammals']]]\n",
      "testData = pd.DataFrame(testData, columns=data.columns)\n",
      "testData"
    ]
  },
  {
    "cell_type": "markdown",
    "metadata": {},
    "source": [
      "\n",
      "We first extract the predictor and target class attributes from the test data and then apply the decision tree classifier to predict their classes."
    ]
  },
  {
    "cell_type": "code",
    "execution_count": null,
    "metadata": {},
    "outputs": [],
    "source": [
      "testX = testData[['Name', 'Class']]\n",
      "testY = testData['Class']\n",
      "preds = clf.predict(testX)\n",
      "predictions = pd.concat([testData['Name'], pd.Series(preds, name='Predicted Class')], axis=1)\n",
      "predictions"
    ]
  },
  {
    "cell_type": "markdown",
    "metadata": {},
    "source": [
      "\n",
      "Except for platypus, which is an egg-laying mammal, the classifier correctly predicts the class label of the test examples. We can calculate the accuracy of the classifier on the test data as shown by the example given below."
    ]
  },
  {
    "cell_type": "code",
    "execution_count": null,
    "metadata": {},
    "outputs": [],
    "source": [
      "from sklearn.metrics import accuracy_score\n",
      "print('Accuracy on test data is %.2f' % (accuracy_score(testY, preds)))"
    ]
  },
  {
    "cell_type": "markdown",
    "metadata": {},
    "source": [
      "\n",
      "## 3.3 Model Overfitting\n",
      "\n",
      "To illustrate the problem of model overfitting, we consider a two-dimensional dataset containing 1500 labeled instances, each of which is assigned to one of two classes, 0 or 1. Instances from each class are generated as follows:\n",
      "1. Instances from class 1 are generated from a mixture of 3 Gaussian distributions, centered at [6, 14], [10, 6], and [14, 14], respectively.\n",
      "2. Instances from class 0 are generated from a uniform distribution in a square region, whose sides have a length equals to 20.\n",
      "\n",
      "For simplicity, both classes have equal number of labeled instances. The code for generating and plotting the data is shown below. All instances from class 1 are shown in red while those from class 0 are shown in black."
    ]
  },
  {
    "cell_type": "code",
    "execution_count": null,
    "metadata": {},
    "outputs": [],
    "source": [
      "import numpy as np\n",
      "import matplotlib.pyplot as plt\n",
      "from numpy.random import random\n",
      "N = 1500\n",
      "mean1 = [6, 14]\n",
      "mean2 = [10, 6]\n",
      "mean3 = [14, 14]\n",
      "cov = [[3.5, 0], [0, 3.5]]\n",
      "np.random.seed(50)\n",
      "X = np.concatenate((X, np.random.multivariate_normal(mean1, cov, int(N/6))))\n",
      "X = np.concatenate((X, np.random.multivariate_normal(mean2, cov, int(N/6))))\n",
      "X = np.concatenate((X, np.random.multivariate_normal(mean3, cov, int(N/6))))\n",
      "Y = np.concatenate((np.ones(int(N/2)), np.zeros(int(N/2))))\n",
      "plt.plot(X[int(N/2), 0], X[int(N/2), 1], 'r', X[int(N/2), 0], X[int(N/2), 1], 'k', ms=4)"
    ]
  },
  {
    "cell_type": "markdown",
    "metadata": {},
    "source": [
      "\n",
      "In this example, we reserve 80% of the labeled data for training and the remaining 20% for testing. We then fit decision trees of different maximum depths (from 2 to 50) to the training set and plot their respective accuracies when applied to the training and test sets."
    ]
  },
  {
    "cell_type": "code",
    "execution_count": null,
    "metadata": {},
    "outputs": [],
    "source": [
      "from sklearn.model_selection import train_test_split\n",
      "X_train, X_test, Y_train, Y_test = train_test_split(X, Y, test_size=0.8, random_state=1)\n",
      "from sklearn.metrics import accuracy_score\n",
      "maxdepths = [2, 3, 4, 5, 6, 7, 8, 9, 10, 15, 20, 25, 30, 35, 40, 45, 50]\n",
      "trainAcc = np.zeros(len(maxdepths))\n",
      "testAcc = np.zeros(len(maxdepths))\n",
      "index = 0\n",
      "for depth in maxdepths:\n",
      "    clf = tree.DecisionTreeClassifier(max_depth=depth)\n",
      "    clf = clf.fit(X_train, Y_train)\n",
      "    Y_predTrain = clf.predict(X_train)\n",
      "    Y_predTest = clf.predict(X_test)\n",
      "    trainAcc[index] = accuracy_score(Y_train, Y_predTrain)\n",
      "    testAcc[index] = accuracy_score(Y_test, Y_predTest)\n",
      "    index += 1\n",
      "plt.plot(maxdepths, trainAcc, 'r-', maxdepths, testAcc, 'b--')\n",
      "plt.legend(['Training Accuracy', 'Test Accuracy'])\n",
      "plt.xlabel('Max depth')\n",
      "plt.ylabel('Accuracy')"
```