

```
cells: [
  {
    "cell_type": "markdown",
    "metadata": {},
    "source": [
      "# Module 5: Regression\n",
      "\n",
      "The following tutorial contains Python examples for solving regression problems. You should refer to the Appendix chapter on regression of the \"Introduction to Data Mining\" book to understand some of the concepts introduced in this tutorial. The notebook can be downloaded from http://www.cse.msu.edu/~ptan/dmbook/tutorials/tutorial5/tutorial5.ipynb.\n",
      "\n",
      "Regression is a modeling technique for predicting quantitative-valued target attributes. The goals for this tutorial are as follows:\n",
      "\n",
      "1. To provide examples of using different regression methods from the scikit-learn library package.\n",
      "2. To demonstrate the problem of model overfitting due to correlated attributes in the data.\n",
      "3. To illustrate how regularization can be used to avoid 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."
    ],
    "execution_count": null,
    "metadata": {},
    "outputs": [],
    "source": [
      "# Importing libraries\n",
      "import numpy as np\n",
      "import matplotlib.pyplot as plt\n",
      "\n",
      "# seed = 1 # seed for random number generation\n",
      "numInstances = 200 # number of data instances\n",
      "np.random.seed(seed)\n",
      "X = np.random.rand(numInstances,1).reshape(-1,1)\n",
      "y_true = -3*X + 1\n",
      "y = y_true + np.random.normal(size=numInstances).reshape(-1,1)\n",
      "\n",
      "plt.scatter(X, y, color='black')\n",
      "plt.plot(X, y_true, color='blue', linewidth=3)\n",
      "plt.xlabel('True function: y = -3X + 1')\n",
      "plt.ylabel('Y')\n",
      "plt.title('Y')\n",
    ],
    "cell_type": "markdown",
    "metadata": {},
    "source": [
      "# 5.1 Synthetic Data Generation\n",
      "\n",
      "To illustrate how linear regression works, we first generate a random 1-dimensional vector of predictor variables, X, from a uniform distribution. The response variable y has a linear relationship with X according to the following equation:  $y = -3x + 1 + \epsilon$ , where  $\epsilon$  corresponds to random noise sampled from a Gaussian distribution with mean 0 and standard deviation of 1."
    ],
    "execution_count": null,
    "metadata": {},
    "outputs": [],
    "source": [
      "# Regression\n",
      "1. Split the input data into their respective training and test sets.\n",
      "2. Fit multiple linear regression to the training data.\n",
      "3. Apply the model to the test data.\n",
      "4. Evaluate the performance of the model.\n",
      "5. Postprocessing: Visualizing the fitted model."
    ],
    "cell_type": "markdown",
    "metadata": {},
    "source": [
      "### Step 1: Split Input Data into Training and Test Sets"
    ],
    "cell_type": "code",
    "execution_count": null,
    "metadata": {},
    "outputs": [],
    "source": [
      "numTrain = 20 # number of training instances\n",
      "numTest = numInstances - numTrain\n",
      "\n",
      "X_train = X[1-numTest]\n",
      "X_test = X[numTest]\n",
      "y_train = y[1-numTest]\n",
      "y_test = y[numTest]"
    ],
    "cell_type": "markdown",
    "metadata": {},
    "source": [
      "### Step 2: Fit Regression Model to Training Set"
    ],
    "cell_type": "code",
    "execution_count": null,
    "metadata": {},
    "outputs": [],
    "source": [
      "from sklearn import linear_model\n",
      "from sklearn.metrics import mean_squared_error, r2_score\n",
      "\n",
      "# Create linear regression object\n",
      "regr = linear_model.LinearRegression()\n",
      "\n",
      "# Fit regression model to the training set\n",
      "regr.fit(X_train, y_train)"
    ],
    "cell_type": "markdown",
    "metadata": {},
    "source": [
      "### Step 3: Apply Model to the Test Set"
    ],
    "cell_type": "code",
    "execution_count": null,
    "metadata": {},
    "outputs": [],
    "source": [
      "# Apply model to the test set\n",
      "y_pred_test = regr.predict(X_test)"
    ],
    "cell_type": "markdown",
    "metadata": {},
    "source": [
      "### Step 4: Evaluate Model Performance on Test Set"
    ],
    "cell_type": "code",
    "execution_count": null,
    "metadata": {},
    "outputs": [],
    "source": [
      "# Comparing true versus predicted values\n",
      "plt.scatter(y_test, y_pred_test, color='black')\n",
      "plt.title('Comparing true and predicted values for test set')\n",
      "plt.xlabel('True values for y')\n",
      "plt.ylabel('Predicted values for y')\n",
      "\n",
      "# Model evaluation\n",
      "print('Root mean squared error = %.4f' % np.sqrt(mean_squared_error(y_test, y_pred_test)))\n",
      "print('R-squared = %.4f' % r2_score(y_test, y_pred_test))"
    ],
    "cell_type": "markdown",
    "metadata": {},
    "source": [
      "### Step 5: Postprocessing"
    ],
    "cell_type": "code",
    "execution_count": null,
    "metadata": {},
    "outputs": [],
    "source": [
      "# Display model parameters\n",
      "print('Slope = ', regr.coef_[0][0])\n",
      "print('Intercept = ', regr.intercept_[0])\n",
      "\n",
      "# Plot outputs\n",
      "plt.scatter(X_test, y_test, color='black')\n",
      "plt.plot(X_test, y_pred_test, color='blue', linewidth=3)\n",
      "plt.title('Predicted Function: y = %.2fX + %.2f' % (regr.coef_[0], regr.intercept_[0]))\n",
      "plt.xlabel('X')\n",
      "plt.ylabel('Y')\n",
    ],
    "cell_type": "markdown",
    "metadata": {},
    "source": [
      "### 5.3 Effect of Correlated Attributes\n",
      "\n",
      "In this example, we illustrate how the presence of correlated attributes can affect the performance of regression models. Specifically, we create 4 additional variables, X2, X3, X4, and X5 that are strongly correlated with the previous variable X created in Section 5.1. The relationship between X and y remains the same as before. We then fit y against the predictor variables and compare their training and test set errors."
    ],
    "cell_type": "code",
    "execution_count": null,
    "metadata": {},
    "outputs": [],
    "source": [
      "# First, we create the correlated attributes below."
    ],
    "cell_type": "code",
    "execution_count": null,
    "metadata": {},
    "outputs": [],
    "source": [
      "seed = 1\n",
      "np.random.seed(seed)\n",
      "X2 = 0.5*X + np.random.normal(0, 0.04, size=numInstances).reshape(-1,1)\n",
      "X3 = 0.5*X2 + np.random.normal(0, 0.01, size=numInstances).reshape(-1,1)\n",
      "X4 = 0.5*X3 + np.random.normal(0, 0.01, size=numInstances).reshape(-1,1)\n",
      "X5 = 0.5*X4 + np.random.normal(0, 0.01, size=numInstances).reshape(-1,1)\n",
      "\n",
      "fig, (ax1, ax2, ax3, ax4) = plt.subplots(2, 2, figsize=(12,9))\n",
      "ax1.scatter(X, X2, color='black')\n",
      "ax1.set_xlabel('X')\n",
      "ax1.set_ylabel('X2')\n",
      "c = np.corrcoef(np.column_stack([X[1-numTest], X2[1-numTest]]).T)\n",
      "titlestr = 'Correlation between X and X2 = %.4f' % (c[0,1])\n",
      "ax1.set_title(titlestr)\n",
      "\n",
      "ax2.scatter(X2, X3, color='black')\n",
      "ax2.set_xlabel('X2')\n",
      "ax2.set_ylabel('X3')\n",
      "c = np.corrcoef(np.column_stack([X2[1-numTest], X3[1-numTest]]).T)\n",
      "titlestr = 'Correlation between X2 and X3 = %.4f' % (c[0,1])\n",
      "ax2.set_title(titlestr)\n",
      "\n",
      "ax3.scatter(X3, X4, color='black')\n",
      "ax3.set_xlabel('X3')\n",
      "ax3.set_ylabel('X4')\n",
      "c = np.corrcoef(np.column_stack([X3[1-numTest], X4[1-numTest]]).T)\n",
      "titlestr = 'Correlation between X3 and X4 = %.4f' % (c[0,1])\n",
      "ax3.set_title(titlestr)\n",
      "\n",
      "ax4.scatter(X4, X5, color='black')\n",
      "ax4.set_xlabel('X4')\n",
      "ax4.set_ylabel('X5')\n",
      "c = np.corrcoef(np.column_stack([X4[1-numTest], X5[1-numTest]]).T)\n",
      "titlestr = 'Correlation between X4 and X5 = %.4f' % (c[0,1])\n",
      "ax4.set_title(titlestr)"
    ],
    "cell_type": "markdown",
    "metadata": {},
    "source": [
      "Next, we create 4 additional versions of the training and test sets. The first version, X_train2 and X_test2 have 2 correlated predictor variables, X and X2. The second version, X_train3 and X_test3 have 3 correlated predictor variables, X, X2, and X3. The third version have 4 correlated variables, X, X2, X3, and X4 whereas the last version have 5 correlated variables, X, X2, X3, X4, and X5."
    ],
    "cell_type": "code",
    "execution_count": null,
    "metadata": {},
    "outputs": [],
    "source": [
      "X_train2 = np.column_stack([X[1-numTest], X2[1-numTest]])\n",
      "X_test2 = np.column_stack([X[1-numTest], X2[1-numTest]])\n",
      "X_train3 = np.column_stack([X[1-numTest], X2[1-numTest], X3[1-numTest]])\n",
      "X_test3 = np.column_stack([X[1-numTest], X2[1-numTest], X3[1-numTest]])\n",
      "X_train4 = np.column_stack([X[1-numTest], X2[1-numTest], X3[1-numTest], X4[1-numTest]])\n",
      "X_test4 = np.column_stack([X[1-numTest], X2[1-numTest], X3[1-numTest], X4[1-numTest]])\n",
      "X_train5 = np.column_stack([X[1-numTest], X2[1-numTest], X3[1-numTest], X4[1-numTest], X5[1-numTest]])\n",
      "X_test5 = np.column_stack([X[1-numTest], X2[1-numTest], X3[1-numTest], X4[1-numTest], X5[1-numTest]])"
    ],
    "cell_type": "markdown",
    "metadata": {},
    "source": [
      "Below, we train 4 new regression models based on the 4 versions of training and test data created in the previous step."
    ],
    "cell_type": "code",
    "execution_count": null,
    "metadata": {},
    "outputs": [],
    "source": [
      "regr2 = linear_model.LinearRegression()\n",
      "regr2.fit(X_train2, y_train)\n",
      "\n",
      "regr3 = linear_model.LinearRegression()\n",
      "regr3.fit(X_train3, y_train)\n",
      "\n",
      "regr4 = linear_model.LinearRegression()\n",
      "regr4.fit(X_train4, y_train)\n",
      "\n",
      "regr5 = linear_model.LinearRegression()\n",
      "regr5.fit(X_train5, y_train)"
    ],
    "cell_type": "markdown",
    "metadata": {},
    "source": [
      "All 4 versions of the regression models are then applied to the training and test sets."
    ],
    "cell_type": "code",
    "execution_count": null,
    "metadata": {},
    "outputs": [],
    "source": [
      "y_pred_train = regr.predict(X_train)\n",
      "y_pred_test = regr.predict(X_test)\n",
      "y_pred_train2 = regr2.predict(X_train2)\n",
      "y_pred_test2 = regr2.predict(X_test2)\n",
      "y_pred_train3 = regr3.predict(X_train3)\n",
      "y_pred_test3 = regr3.predict(X_test3)\n",
      "y_pred_train4 = regr4.predict(X_train4)\n",
      "y_pred_test4 = regr4.predict(X_test4)\n",
      "y_pred_train5 = regr5.predict(X_train5)\n",
      "y_pred_test5 = regr5.predict(X_test5)"
    ],
    "cell_type": "markdown",
    "metadata": {},
    "source": [
      "For postprocessing, we compute both the training and test errors of the models. We can also show the resulting model and the sum of the absolute weights of the regression coefficients, i.e.,  $\sum_{j=0}^d |w_j|$ , where  $d$  is the number of predictor attributes."
    ],
    "cell_type": "code",
    "execution_count": null,
    "metadata": {},
    "outputs": [],
    "source": [
      "import pandas as pd\n",
      "import matplotlib.pyplot as plt\n",
      "\n",
      "columns = ['Model', 'Train error', 'Test error', 'Sum of Absolute Weights']\n",
      "model1 = '%.2f X + %.2f' % (regr.coef_[0][0], regr.intercept_[0])\n",
      "value1 = [model1, np.sqrt(mean_squared_error(y_train, y_pred_train)),\n",
      "          np.sqrt(mean_squared_error(y_test, y_pred_test)),\n",
      "          np.absolute(regr.coef_[0]).sum() + np.absolute(regr.intercept_[0])]\n",
      "\n",
      "model2 = '%.2f X + %.2f X2 + %.2f' % (regr2.coef_[0][0], regr2.coef_[0][1], regr2.intercept_[0])\n",
      "value2 = [model2, np.sqrt(mean_squared_error(y_train, y_pred_train2)),\n",
      "          np.sqrt(mean_squared_error(y_test, y_pred_test2)),\n",
      "          np.absolute(regr2.coef_[0]).sum() + np.absolute(regr2.intercept_[0])]\n",
      "\n",
      "model3 = '%.2f X + %.2f X2 + %.2f X3 + %.2f' % (regr3.coef_[0][0], regr3.coef_[0][1], regr3.coef_[0][2], regr3.intercept_[0])\n",
      "value3 = [model3, np.sqrt(mean_squared_error(y_train, y_pred_train3)),\n",
      "          np.sqrt(mean_squared_error(y_test, y_pred_test3)),\n",
      "          np.absolute(regr3.coef_[0]).sum() + np.absolute(regr3.intercept_[0])]\n",
      "\n",
      "model4 = '%.2f X + %.2f X2 + %.2f X3 + %.2f X4 + %.2f' % (regr4.coef_[0][0], regr4.coef_[0][1], regr4.coef_[0][2], regr4.coef_[0][3], regr4.intercept_[0])\n",
      "value4 = [model4, np.sqrt(mean_squared_error(y_train, y_pred_train4)),\n",
      "          np.sqrt(mean_squared_error(y_test, y_pred_test4)),\n",
      "          np.absolute(regr4.coef_[0]).sum() + np.absolute(regr4.intercept_[0])]\n",
      "\n",
      "model5 = '%.2f X + %.2f X2 + %.2f X3 + %.2f X4 + %.2f X5 + %.2f' % (regr5.coef_[0][0], regr5.coef_[0][1], regr5.coef_[0][2], regr5.coef_[0][3], regr5.coef_[0][4], regr5.intercept_[0])\n",
      "value5 = [model5, np.sqrt(mean_squared_error(y_train, y_pred_train5)),\n",
      "          np.sqrt(mean_squared_error(y_test, y_pred_test5)),\n",
      "          np.absolute(regr5.coef_[0]).sum() + np.absolute(regr5.intercept_[0])]\n",
      "\n",
      "results = pd.DataFrame([value1, value2, value3, value4, value5], columns=columns)\n",
      "\n",
      "plt.plot(results['Sum of Absolute Weights'], results['Train error'], 'r--')\n",
      "plt.plot(results['Sum of Absolute Weights'], results['Test error'], 'k--')\n",
      "plt.legend(['Train error', 'Test error'])\n",
      "plt.xlabel('Sum of Absolute Weights')\n",
      "plt.ylabel('Error rate')\n",
      "plt.title('Error rate')\n",
      "plt.show()"
    ],
    "cell_type": "markdown",
    "metadata": {},
    "source": [
      "The results above show that the first model, which fits y against X only, has the largest training error, but smallest test error, whereas the fifth model, which fits y against X and other correlated attributes, has the smallest training error but largest test error. This is due to a phenomenon known as model overfitting, in which the model is too closely tailored to the training data and thus performs poorly on new, unseen test instances. From the plot, we can see that the disparity between the training and test errors becomes wider as the sum of absolute weights of the model (which represents the model complexity) increases. Thus, one should control the complexity of the regression model to avoid the model overfitting problem."
    ],
    "cell_type": "code",
    "execution_count": null,
    "metadata": {},
    "outputs": [],
    "source": [
      "### 5.4 Ridge Regression\n",
      "\n",
      "Ridge regression is a variant of MLR designed to fit a linear model to the dataset by minimizing the following regularized least-square loss function:\n",
      "
$$S(\mathbf{w}) = \sum_{i=1}^N \|\mathbf{y}_i - \mathbf{X}_i \mathbf{w} - w_0\|^2 + \lambda \|\mathbf{w}\|^2 + w_0^2$$
\n",
      "where  $\lambda$  is the hyperparameter for ridge regression. Note that the ridge regression model reduces to MLR when  $\lambda = 0$ . By increasing the value of  $\lambda$ , we can control the complexity of the model as will be shown in the example below."
    ],
    "cell_type": "code",
    "execution_count": null,
    "metadata": {},
    "outputs": [],
    "source": [
      "In the example shown below, we fit a ridge regression model to the previously created training set with correlated attributes. We compare the results of the ridge regression model against those obtained using MLR."
    ],
    "cell_type": "code",
    "execution_count": null,
    "metadata": {},
    "outputs": [],
    "source": [
      "from sklearn import linear_model\n",
      "\n",
      "ridge = linear_model.Ridge(alpha=0.4)\n",
      "ridge.fit(X_train5, y_train)\n",
      "y_pred_train_ridge = ridge.predict(X_train5)\n",
      "y_pred_test_ridge = ridge.predict(X_test5)\n",
      "\n",
      "model6 = '%.2f X + %.2f X2 + %.2f X3 + %.2f X4 + %.2f X5 + %.2f' % (ridge.coef_[0][0],\n",
      "          ridge.coef_[0][1],\n",
      "          ridge.coef_[0][2],\n",
      "          ridge.coef_[0][3],\n",
      "          ridge.coef_[0][4],\n",
      "          ridge.intercept_[0])\n",
      "value6 = [model6, np.sqrt(mean_squared_error(y_train, y_pred_train_ridge)),\n",
      "          np.sqrt(mean_squared_error(y_test, y_pred_test_ridge)),\n",
      "          np.absolute(ridge.coef_[0]).sum() + np.absolute(ridge.intercept_[0])]\n",
      "\n",
      "ridge_results = pd.DataFrame([value6], columns=columns, index=['Ridge'])\n",
      "pd.concat([results, ridge_results])"
    ],
    "cell_type": "markdown",
    "metadata": {},
    "source": [
      "By setting an appropriate value for the hyperparameter,  $\lambda$ , we can control the sum of absolute weights, thus producing a test error that is quite comparable to that of MLR without the correlated attributes."
    ],
    "cell_type": "code",
    "execution_count": null,
    "metadata": {},
    "outputs": [],
    "source": [
      "### 5.5 Lasso Regression\n",
      "\n",
      "One of the limitations of ridge regression is that, although it was able to reduce the regression coefficients associated with the correlated attributes and reduce the effect of model overfitting, the resulting model is still not sparse. Another variation of MLR, called Lasso regression, is designed to produce sparser models by imposing an  $\ell_1$  regularization on the regression coefficients as shown below:\n",
      "
$$S(\mathbf{w}) = \sum_{i=1}^N \|\mathbf{y}_i - \mathbf{X}_i \mathbf{w} - w_0\|^2 + \lambda \|\mathbf{w}\|_1 + |w_0|$$
\n",
      "The example code below shows the results of applying Lasso regression to the previously used correlated dataset."
    ],
    "cell_type": "code",
    "execution_count": null,
    "metadata": {},
    "outputs": [],
    "source": [
      "from sklearn import linear_model\n",
      "\n",
      "lasso = linear_model.Lasso(alpha=0.02)\n",
      "lasso.fit(X_train5, y_train)\n",
      "y_pred_train_lasso = lasso.predict(X_train5)\n",
      "y_pred_test_lasso = lasso.predict(X_test5)\n",
      "\n",
      "model7 = '%.2f X + %.2f X2 + %.2f X3 + %.2f X4 + %.2f X5 + %.2f' % (lasso.coef_[0],\n",
      "          lasso.coef_[1],\n",
      "          lasso.coef_[2],\n",
      "          lasso.coef_[3],\n",
      "          lasso.coef_[4],\n",
      "          lasso.intercept_[0])\n",
      "value7 = [model7, np.sqrt(mean_squared_error(y_train, y_pred_train_lasso)),\n",
      "          np.sqrt(mean_squared_error(y_test, y_pred_test_lasso)),\n",
      "          np.absolute(lasso.coef_[0]).sum() + np.absolute(lasso.intercept_[0])]\n",
      "\n",
      "lasso_results = pd.DataFrame([value7], columns=columns, index=['Lasso'])\n",
      "pd.concat([results, lasso_results])"
    ],
    "cell_type": "code",
    "execution_count": null,
    "metadata": {},
    "outputs": [],
    "source": [
      "In this next example, we illustrate how to apply cross-validation to select the best hyperparameter value for fitting a Lasso regression model."
    ],
    "cell_type": "code",
    "execution_count": null,
    "metadata": {},
    "outputs": [],
    "source": [
      "from sklearn import linear_model\n",
      "\n",
      "lasso = linear_model.LassoCV(cv=5, alphas=[0.2, 0.4, 0.6, 0.8, 1.0])\n",
      "lasso.fit(X_train5, y_train)\n",
      "y_pred_train_lasso = lasso.predict(X_train5)\n",
      "y_pred_test_lasso = lasso.predict(X_test5)\n",
      "\n",
      "model8 = '%.2f X + %.2f X2 + %.2f X3 + %.2f X4 + %.2f X5 + %.2f' % (lasso.coef_[0],\n",
      "          lasso.coef_[1],\n",
      "          lasso.coef_[2],\n",
      "          lasso.coef_[3],\n",
      "          lasso.coef_[4],\n",
      "          lasso.intercept_[0])\n",
      "value8 = [model8, np.sqrt(mean_squared_error(y_train, y_pred_train_lasso)),\n",
      "          np.sqrt(mean_squared_error(y_test, y_pred_test_lasso)),\n",
      "          np.absolute(lasso.coef_[0]).sum() + np.absolute(lasso.intercept_[0])]\n",
      "\n",
      "lasso_results = pd.DataFrame([value8], columns=columns, index=['LassoCV'])\n",
      "pd.concat([results, lasso_results])"
    ],
    "cell_type": "code",
    "execution_count": null,
    "metadata": {},
    "outputs": [],
    "source": [
      "This section presents example Python code for fitting linear regression models to a dataset. We also illustrate the problem of model overfitting and shows two alternative methods, called ridge and Lasso regression, that can help alleviate such a problem. While the model overfitting problem shown here is illustrated in the context of correlated attributes, the problem is more general and may arise due to other factors such as noise and other exceptional values in the data."
    ],
    "cell_type": "code",
    "execution_count": null,
    "metadata": {},
    "outputs": [],
    "source": [
      "kernelspec: {\n",
      "  \"display_name\": \"Python 3\",\n",
      "  \"language\": \"python\",\n",
      "  \"name\": \"python3\"\n",
      },\n",
      "language_info: {\n",
      "  \"codemirror_mode\": {\n",
      "    \"name\": \"ipython\",\n",
      "    \"version\": 3\n",
      "  },\n",
      "  \"file_extension\": \".py\",\n",
      "  \"mimetype\": \"text/x-python\",\n",
      "  \"name\": \"python\",\n",
      "  \"nbconvert_exporter\": \"python\",\n",
      "  \"pygments_lexer\": \"ipython3\",\n",
      "  \"version\": \"3.6.4\"\n",
      },\n",
      "nbformat: 4,\n",
      "nbformat_minor: 2
    ]
```