

```

{"cells": [{"metadata": {}, "cell_type": "markdown", "source": "**This notebook is an exercise in the [Intermediate Machine Learning] (https://www.kaggle.com/learn/intermediate-machine-learning) course. You can reference the tutorial at [this link] (https://www.kaggle.com/alexisbcook/pipelines).**"}, {"metadata": {}, "cell_type": "code", "source": "In this exercise, you will use pipelines to improve the efficiency of your machine learning code.\n\n# Setup\n\nThe questions below will give you feedback on your work. Run the following cell to set up the feedback system."}, {"metadata": {}, "cell_type": "code", "source": "# Set up code checking\nimport os\nif not os.path.exists(\"../input/train.csv\"):\n    os.symlink(\"../input/home-data-for-ml-course/train.csv\", \"../input/train.csv\")\n    os.symlink(\"../input/home-data-for-ml-course/test.csv\", \"../input/test.csv\")\n\nfrom learntools.core import binder\nbinder.bind(globals())\n\nfrom learntools.ml_intermediate.ex4 import *\n\nprint(\"Setup Complete\")", "execution_count": null, "outputs": []}, {"metadata": {}, "cell_type": "markdown", "source": "You will work with data from the [Housing Prices Competition for Kaggle Learn Users] (https://www.kaggle.com/c/home-data-for-ml-course). \n\n! [Ames Housing dataset image] (https://i.imgur.com/LTJV4e.png)\n\nRun the next code cell without changes to load the training and validation sets in `X_train`, `X_valid`, `y_train`, and `y_valid`. The test set is loaded in `X_test`."}, {"metadata": {}, "cell_type": "code", "source": "import pandas as pd\n\nfrom sklearn.model_selection import train_test_split\n\n# Read the data\nX_full = pd.read_csv(\"../input/train.csv\", index_col='Id')\nX_test_full = pd.read_csv(\"../input/test.csv\", index_col='Id')\n\n# Remove rows with missing target, separate target from predictors\nX_full.dropna(axis=0, subset=['SalePrice'], inplace=True)\ny = X_full.SalePrice\nX_full.drop(['SalePrice'], axis=1, inplace=True)\n\n# Break off validation set from training data\nX_train_full, X_valid_full, y_train, y_valid = train_test_split(X_full, y, random_state=0)\n\n# \"Cardinality\" means the number of unique values in a column\n# Select categorical columns with relatively low cardinality (convenient but arbitrary)\ncategorical_cols = [cname for cname in X_train_full.columns if\nX_train_full[cname].nunique() < 10 and\nX_train_full[cname].dtype == \"object\"]\n\n# Select numerical columns\nnumerical_cols = [cname for cname in X_train_full.columns if\nX_train_full[cname].dtype in ['int64', 'float64']]\n\n# Keep selected columns only\nmy_cols = categorical_cols + numerical_cols\nX_train = X_train_full[my_cols].copy()\nX_valid = X_valid_full[my_cols].copy()\nX_test = X_test_full[my_cols].copy()", "execution_count": null, "outputs": []}, {"metadata": {}, "cell_type": "code", "source": "X_train.head()", "execution_count": null, "outputs": []}, {"metadata": {}, "cell_type": "markdown", "source": "The next code cell uses code from the tutorial to preprocess the data and train a model. Run this code without changes."}, {"metadata": {}, "cell_type": "code", "source": "from sklearn.compose import ColumnTransformer\n\nfrom sklearn.pipeline import Pipeline\n\nfrom sklearn.impute import SimpleImputer\n\nfrom sklearn.preprocessing import OneHotEncoder\n\nfrom sklearn.ensemble import RandomForestRegressor\n\nfrom sklearn.metrics import mean_absolute_error\n\n# Preprocessing for numerical data\nnumerical_transformer = SimpleImputer(strategy='constant')\n\n# Preprocessing for categorical data\ncategorical_transformer = Pipeline(steps=[\n    ('imputer', SimpleImputer(strategy='most_frequent')), \n    ('onehot', OneHotEncoder(handle_unknown='ignore'))\n])\n\n# Bundle preprocessing for numerical and categorical data\npreprocessor = ColumnTransformer([\n    ('num', numerical_transformer, numerical_cols), \n    ('cat', categorical_transformer, categorical_cols)\n])\n\n# Define model\nmodel = RandomForestRegressor(n_estimators=100, random_state=0)\n\n# Bundle preprocessing and modeling code in a pipeline\nclf = Pipeline(steps=[('preprocessor', preprocessor),\n    ('model', model)])\n\n# Preprocessing of training data, fit model\nclf.fit(X_train, y_train)\n\n# Preprocessing of validation data, get predictions\npreds = clf.predict(X_valid)\n\nprint('MAE:', mean_absolute_error(y_valid, preds))", "execution_count": null, "outputs": []}, {"metadata": {}, "cell_type": "markdown", "source": "The code yields a value around 17862 for the mean absolute error (MAE). In the next step, you will amend the code to do better.\n\n# Step 1: Improve the performance\n\n### Part A\n\nNow, it's your turn! In the code cell below, define your own preprocessing steps and random forest model. Fill in values for the following variables:\n- `numerical_transformer`\n- `categorical_transformer`\n- `model`\n\nTo pass this part of the exercise, you need only define valid preprocessing steps and a random forest model."}, {"metadata": {}, "cell_type": "code", "source": "# Preprocessing for numerical data\nnumerical_transformer ="}

```

```

SimpleImputer(strategy='constant') # Your code here\n\n# Preprocessing for categorical data\ncategorical_transformer =
Pipeline(steps=[\n    ('imputer', SimpleImputer(strategy='most_frequent')),\n    ('onehot',
OneHotEncoder(handle_unknown='ignore'))\n]) # Your code here\n\n# Bundle preprocessing for numerical and categorical
data\npreprocessor = ColumnTransformer(\n    transformers=[\n        ('num', numerical_transformer, numerical_cols),\n        ('cat', categorical_transformer, categorical_cols)\n    ])\n\n# Define model\nmodel =
RandomForestRegressor(random_state=0)\nprint(model.get_params().keys())\n\nn_estimators = [100, 80, 120, 140,160,180,
200]\nparam_grid = dict(n_estimators = n_estimators)#转化为字典格式, 网络搜索要求\n\n# Check your
answer\nstep_1.a.check()", "execution_count": null, "outputs": [], {"metadata": {"trusted": true}, "cell_type": "code", "source": "# Lines
below will give you a hint or solution code\n#step_1.a.hint()\n#step_1.a.solution()", "execution_count": null, "outputs": [],
{"metadata": {}, "cell_type": "markdown", "source": "### Part B\n\nRun the code cell below without changes.\n\nTo pass this step, you
need to have defined a pipeline in Part A that achieves lower MAE than the code above. You're encouraged to take your time
here and try out many different approaches, to see how low you can get the MAE! (_If your code does not pass, please amend the
preprocessing steps and model in Part A._)"}], {"metadata": {"trusted": true}, "cell_type": "code", "source": "# Bundle preprocessing and
modeling code in a pipeline\n# Bundle preprocessing and modeling code in a pipeline\nmy_pipeline = Pipeline(steps=[('preprocessor',
preprocessor),\n    ('model', GridSearchCV(model,param_grid,cv = 5))\n])\n\n# print(\n"Best: %f using %s" % (grid_result.best_score_,grid_search.best_params_))\n# params =
grid_result.cv_results_['params']\n\n# Preprocessing of training data, fit model \nmy_pipeline.fit(X_train, y_train)\n\n#
Preprocessing of validation data, get predictions\npreds = my_pipeline.predict(X_valid)\n\n# Evaluate the model\nscore =
mean_absolute_error(y_valid, preds)\nprint('MAE:', score)\n\n# Check your
answer\nstep_1.b.check()", "execution_count": null, "outputs": [], {"metadata": {"trusted": true}, "cell_type": "code", "source": "# Line
below will give you a hint\n#step_1.b.hint()", "execution_count": null, "outputs": [], {"metadata":
{"cell_type": "markdown", "source": "# Step 2: Generate test predictions\n\nNow, you'll use your trained model to generate
predictions with the test data."}], {"metadata": {"trusted": true}, "cell_type": "code", "source": "# Preprocessing of test data, fit
model\npreds_test = my_pipeline.predict(X_test) # Your code here\n\n# Check your
answer\nstep_2.check()", "execution_count": null, "outputs": [], {"metadata": {"trusted": true}, "cell_type": "code", "source": "# Lines
below will give you a hint or solution code\n#step_2.hint()\n#step_2.solution()", "execution_count": null, "outputs": [], {"metadata":
{"cell_type": "markdown", "source": "Run the next code cell without changes to save your results to a CSV file that can be submitted
directly to the competition."}], {"metadata": {"trusted": true}, "cell_type": "code", "source": "# Save test predictions to file\noutput =
pd.DataFrame({'Id': X_test.index,\n    'SalePrice': preds_test})\noutput.to_csv('submission.csv',
index=False)", "execution_count": null, "outputs": [], {"metadata": {}, "cell_type": "markdown", "source": "# Submit your results\n\nOnce
you have successfully completed Step 2, you're ready to submit your results to the leaderboard! If you choose to do so, make sure
that you have already joined the competition by clicking on the Join Competition button at [this link]
(https://www.kaggle.com/c/home-data-for-ml-course). \n1. Begin by clicking on the blue Save Version button in the top right
corner of the window. This will generate a pop-up window. \n2. Ensure that the Save and Run All option is selected, and then
click on the blue Save button.\n3. This generates a window in the bottom left corner of the notebook. After it has finished
running, click on the number to the right of the Save Version button. This pulls up a list of versions on the right of the
screen. Click on the ellipsis (...) to the right of the most recent version, and select Open in Viewer. This brings you
into view mode of the same page. You will need to scroll down to get back to these instructions.\n4. Click on the Output tab on
the right of the screen. Then, click on the file you would like to submit, and click on the blue Submit button to submit your
results to the leaderboard.\n\nYou have now successfully submitted to the competition!\n\nIf you want to keep working to improve
your performance, select the blue Edit button in the top right of the screen. Then you can change your code and repeat the
process. There's a lot of room to improve, and you will climb up the leaderboard as you work.\n\n# Keep going\n\nMove on to learn
about cross-validation (https://www.kaggle.com/alexisbcook/cross-validation), a technique you can use to obtain more accurate
estimates of model performance!"), {"metadata": {}, "cell_type": "markdown", "source": "Have questions or comments? Visit
the [Learn Discussion forum](https://www.kaggle.com/learn-forum/161289) to chat with other Learners.*"}], "metadata": {"kernelspec":

```

```
{"language":"python","display_name":"Python 3","name":"python3"},"language_info":  
{"pygments_lexer":"ipython3","nbconvert_exporter":"python","version":"3.6.4","file_extension":".py","codemirror_mode":  
{"name":"ipython","version":3},"name":"python","mimetype":"text/x-python"}},"nbformat":4,"nbformat_minor":4}
```