

```

{"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/introduction).**\n\n---\n\n"}, {"metadata":{"cell_type":"markdown","source":"As a warm-up, you'll review some machine learning fundamentals and submit your initial results to a Kaggle competition.\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":{"trusted":false,"cell_type":"code","source":"# Set up code checking\nimport os\nif not os.path.exists(\"../input/train.csv\"):\nos.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\") \nfrom learntools.core import binder\nbinder.bind(globals())\nfrom learntools.ml_intermediate.ex1 import *\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) to predict home prices in Iowa using 79 explanatory variables describing (almost) every aspect of the homes. \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 features in `X_train` and `X_valid`, along with the prediction targets in `y_train` and `y_valid`. The test features are loaded in `X_test`. (_If you need to review **features** and **prediction targets**, please check out [this short tutorial](https://www.kaggle.com/dansbecker/your-first-machine-learning-model). To read about model **validation**, look [here](https://www.kaggle.com/dansbecker/model-validation). Alternatively, if you'd prefer to look through a full course to review all of these topics, start [here](https://www.kaggle.com/learn/machine-learning).)_"}, {"metadata":{"trusted":false,"cell_type":"code","source":"import pandas as pd\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# Obtain target and predictors\ny = X_full.SalePrice\nfeatures = ['LotArea', 'YearBuilt', '1stFlrSF', '2ndFlrSF', 'FullBath', 'BedroomAbvGr', 'TotRmsAbvGrd']\nX = X_full[features].copy()\nX_test = X_test_full[features].copy()\n\n# Break off validation set from training data\nX_train, X_valid, y_train, y_valid = train_test_split(X, y, train_size=0.8, test_size=0.2, random_state=0)","execution_count":null,"outputs":[]}, {"metadata":{"cell_type":"markdown","source":"Use the next cell to print the first several rows of the data. It's a nice way to get an overview of the data you will use in your price prediction model."}, {"metadata":{"trusted":false,"cell_type":"code","source":"X_train.head()","execution_count":null,"outputs":[]}, {"metadata":{"cell_type":"markdown","source":"The next code cell defines five different random forest models. Run this code cell without changes. (_To review **random forests**, look [here](https://www.kaggle.com/dansbecker/random-forests)._)"}, {"metadata":{"trusted":false,"cell_type":"code","source":"from sklearn.ensemble import RandomForestRegressor\n\n# Define the models\nmodel_1 = RandomForestRegressor(n_estimators=50, random_state=0)\nmodel_2 = RandomForestRegressor(n_estimators=100, random_state=0)\nmodel_3 = RandomForestRegressor(n_estimators=100, criterion='mae', random_state=0)\nmodel_4 = RandomForestRegressor(n_estimators=200, min_samples_split=20, random_state=0)\nmodel_5 = RandomForestRegressor(n_estimators=100, max_depth=7, random_state=0)\n\nmodels = [model_1, model_2, model_3, model_4, model_5]","execution_count":null,"outputs":[]}, {"metadata":{"cell_type":"markdown","source":"To select the best model out of the five, we define a function `score_model()` below. This function returns the mean absolute error (MAE) from the validation set. Recall that the best model will obtain the lowest MAE. (_To review **mean absolute error**, look [here](https://www.kaggle.com/dansbecker/model-validation)._) \n\nRun the code cell without changes."}, {"metadata":{"trusted":false,"cell_type":"code","source":"from sklearn.metrics import mean_absolute_error\n\n# Function for comparing different models\ndef score_model(model, X_t=X_train, X_v=X_valid, y_t=y_train, y_v=y_valid):\n    model.fit(X_t, y_t)\n    preds = model.predict(X_v)\n    return mean_absolute_error(y_v, preds)\n\nfor i in range(0, len(models)):\n    mae = score_model(models[i])\n    print(\"Model %d MAE: %d\" % (i+1, mae))","execution_count":null,"outputs":[]}, {"metadata":{"cell_type":"markdown","source":"# Step 1: Evaluate several models\n\nUse the above results to fill in the line below. Which model is the best model? Your answer should be one of `model_1`, `model_2`, `model_3`, `model_4`, or `model_5`."}, {"metadata":{"trusted":true,"cell_type":"code","source":"from sklearn.ensemble import RandomForestRegressor\n\n# Define the models\nmodel_1 = RandomForestRegressor(n_estimators=50, random_state=0)\nmodel_2 = RandomForestRegressor(n_estimators=100,"}]}

```

```

random_state=0)\nmodel_3 = RandomForestRegressor(n_estimators=100, criterion='mae', random_state=0)\nmodel_4 =
RandomForestRegressor(n_estimators=200, min_samples_split=20, random_state=0)\nmodel_5 = RandomForestRegressor(n_estimators=100,
max_depth=7, random_state=0)\n\nmodels = [model_1, model_2, model_3, model_4, model_5]", "execution_count": null, "outputs": [],
{"metadata": {"trusted": true}, "cell_type": "code", "source": "from sklearn.metrics import mean_absolute_error\n\n# Function for
comparing different models\ndef score_model(model, X_t=X_train, X_v=X_valid, y_t=y_train, y_v=y_valid):\n    model.fit(X_t, y_t)\n
preds = model.predict(X_v)\n    return mean_absolute_error(y_v, preds)\n\nfor i in range(0, len(models)):\n    mae =
score_model(models[i])\n    print(\"Model %d MAE: %d\" % (i+1, mae))", "execution_count": null, "outputs": [], {"metadata":
{"trusted": false}, "cell_type": "code", "source": "# Fill in the best model\nbest_model = model_3\n\n# Check your
answer\nstep_1.check()", "execution_count": null, "outputs": [], {"metadata": {"trusted": false}, "cell_type": "code", "source": "# Lines
below will give you a hint or solution code\nstep_1.hint()\nstep_1.solution()", "execution_count": null, "outputs": [], {"metadata":
{"trusted": false}, "cell_type": "code", "source": "# Define a model\nmy_model =
RandomForestRegressor(random_state=0) # Your code here\n\n# Check your answer\nstep_2.check()", "execution_count": null, "outputs":
[], {"metadata": {"trusted": false}, "cell_type": "code", "source": "# Lines below will give you a hint or solution
code\nstep_2.hint()\nstep_2.solution()", "execution_count": null, "outputs": [], {"metadata": {"trusted": false}, "cell_type": "code", "source": "# Fit the model to the training data\nmy_model.fit(X, y)\n\n# Generate test
predictions\npreds_test = my_model.predict(X_test)\n\n# Save predictions in format used for competition scoring\noutput =
pd.DataFrame({'Id': X_test.index,\n              'SalePrice': preds_test})\noutput.to_csv('submission.csv',
index=False)", "execution_count": null, "outputs": [], {"metadata": {"trusted": false}, "cell_type": "code", "source": "# Submit your results\n\nOnce
you have successfully completed Step 2, you're ready to submit your results to the leaderboard! First, you'll need to join the
competition if you haven't already. So open a new window by clicking on [this link](https://www.kaggle.com/c/home-data-for-ml-
course). Then click on the **Join Competition** button.\n\n!join competition image](https://i.imgur.com/wLmFtH3.png)\n\nNext,
follow the instructions below:\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"}, {"metadata": {"trusted": false}, "cell_type": "code", "source": "# Keep
going\n\nYou've made your first model. But how can you quickly make it better?\n\nLearn how to improve your competition results by
incorporating columns with **[missing values](https://www.kaggle.com/alexisbcook/missing-values)**.", {"metadata":
{"trusted": false}, "cell_type": "code", "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}

```