---

This exercise will test your ability to read a data file and understand statistics about the data.

In later exercises, you will apply techniques to filter the data, build a machine learning model, and iteratively improve your model.

The course examples use data from Melbourne. To ensure you can apply these techniques on your own, you will have to apply them to a new dataset (with house prices from Iowa).

The exercises use a "notebook" coding environment. In case you are unfamiliar with notebooks, we have a [90-second intro video](#).

# Exercises

Run the following cell to set up code-checking, which will verify your work as you go.

In [ ]:

```
# Set up code checking
from learntools.core import binder
binder.bind(globals())
from learntools.machine_learning.ex2 import *
print("Setup Complete")
```

## Step 1: Loading Data

Read the Iowa data file into a Pandas DataFrame called `home_data`.

In [ ]:

```
import pandas as pd

# Path of the file to read
iowa_file_path = '../input/home-data-for-ml-course/train.csv'

# Fill in the line below to read the file into a variable home_data
home_data = pd.read_csv(iowa_file_path)

# Call line below with no argument to check that you've loaded the data correctly
step_1.check()
```

In [ ]:

```
# Lines below will give you a hint or solution code
#step_1.hint()
#step_1.solution()
```

## Step 2: Review The Data

Use the command you learned to view summary statistics of the data. Then fill in variables to answer the following questions

In [ ]:

## Print summary statistics in next line

In [ ]:

```
home_data.describe()
```

In [ ]:

```
import datetime
# What is the average lot size (rounded to nearest integer)?
avg_lot_size = int(round(home_data['LotArea'].mean()))

# As of today, how old is the newest home (current year - the date in which it was built)

#tmp = pd.datetime.now().year-home_data["YearBuilt"] deprecated !!
tmp = datetime.datetime.now().year-home_data["YearBuilt"]

newest_home_age = tmp.min()
# Checks your answers
step_2.check()
```

In [ ]:

```
#step_2.hint()
#step_2.solution()
```

## Think About Your Data

The newest house in your data isn't that new. A few potential explanations for this:

1.  They haven't built new houses where this data was collected.
2.  The data was collected a long time ago. Houses built after the data publication wouldn't show up.

If the reason is explanation #1 above, does that affect your trust in the model you build with this data? What about if it is reason #2?

How could you dig into the data to see which explanation is more plausible?

Check out this **discussion thread** to see what others think or to add your ideas.

## Keep Going

You are ready for **Your First Machine Learning Model**.

---

**Machine Learning Course Home Page**