**Step1. Open the data/house.arff Dataset**

Click the "*Open file…*" button to open a data set and double click on the "*data*" directory.

Weka provides a number of small common machine learning datasets that you can use to practice on.

Select the "**house.***arff*" file to load the house dataset.



**Step2. Creating the regression model with WEKA**

To create the model, click on the **Classify** tab. The first step is to select the model we want to build, so WEKA knows how to work with the data, and how to create the appropriate model:

1. Click the **Choose** button, then expand the **functions** branch
2. Select the **LinearRegression** leaf with using **training set** test option.
3. Click Start to create a model.

## Step3: Interpreting the regression model

Regression output

```
sellingPrice=(-26.6882*houseSize)+(7.0551*lotSize)+(43166.0767*bedrooms)+
(42292.0901*bathroom)-21661.1208
```

House value using regression model

```
SellingPrice=(-26.6882*2983)+(7.0551*9365)+(43166.0767*5)+(42292.0901*1)       -
21661.1208
```

**sellingPrice = 222,921**

**Step4: Interpret the patterns and conclusions that our model tells us**

- **Granite doesn't matter :** It will throw out and ignore columns that don't help in creating a good model. So this regression model is telling us that granite in your kitchen doesn't affect the house's value.

- **Bathrooms do matter:** Since we use a simple 0 or 1 value for an upgraded bathroom, we can use the coefficient from the regression model to determine the value of an upgraded bathroom on the house value. The model tells us it adds $42,292 to the house value.

- **Bigger houses reduce the value:** Model is telling us that the bigger our house is, the lower the selling price? This can be seen by the negative coefficient in front of the houseSize variable. The model is telling us that every additional square foot of the house reduces its

price by $26? That doesn't make any sense at all. How should we interpret this?  The house size, unfortunately, isn't an independent variable because it's related to the bedrooms variable, which makes sense, since bigger houses tend to have more bedrooms. So our model isn't perfect. But we can fix this.  On the **Preprocess** tab, you can remove columns from the data set. Remove the **houseSize** column and create another model. How does it affect the price of my house? How does this new model make more sense? (My amended house value: $215,554).