```
In [1]:  #!/usr/bin/python

         import random
         import numpy
         import matplotlib.pyplot as plt
         import pickle

         from outlier_cleaner import outlierCleaner


         ### load up some practice data with outliers in it
         ages = pickle.load( open("practice_outliers_ages.pkl", "r") )
         net_worths = pickle.load( open("practice_outliers_net_worths.pkl", "r") )



         ### ages and net_worths need to be reshaped into 2D numpy arrays
         ### second argument of reshape command is a tuple of integers: (n_rows, n_columns)
         ### by convention, n_rows is the number of data points
         ### and n_columns is the number of features
         ages = numpy.reshape( numpy.array(ages), (len(ages), 1))
         net_worths = numpy.reshape( numpy.array(net_worths), (len(net_worths), 1))
         from sklearn.cross_validation import train_test_split
         ages_train, ages_test, net_worths_train, net_worths_test = train_test_split(ages,
```

```
C:\Users\Andrew\Anaconda3\envs\conda2\lib\site-packages\sklearn\cross_validatio
n.py:41: DeprecationWarning: This module was deprecated in version 0.18 in favo
r of the model_selection module into which all the refactored classes and funct
ions are moved. Also note that the interface of the new CV iterators are differ
ent from that of this module. This module will be removed in 0.20.
  "This module will be removed in 0.20.", DeprecationWarning)
```

Sebastian described to us an algorithm for improving a regression, which you will implement in this project. You will work through it in the next few quizzes. To summarize, what you'll do is fit the regression on all training points discard the 10% of points that have the largest errors between the actual y values, and the regression-predicted y values refit on the remaining points.

Start by running the starter code (outliers/outlier_removal_regression.py) and visualizing the points. A few outliers should clearly pop out. Deploy a linear regression, where net worth is the target and the feature being used to predict it is a person's age (remember to train on the training data!).
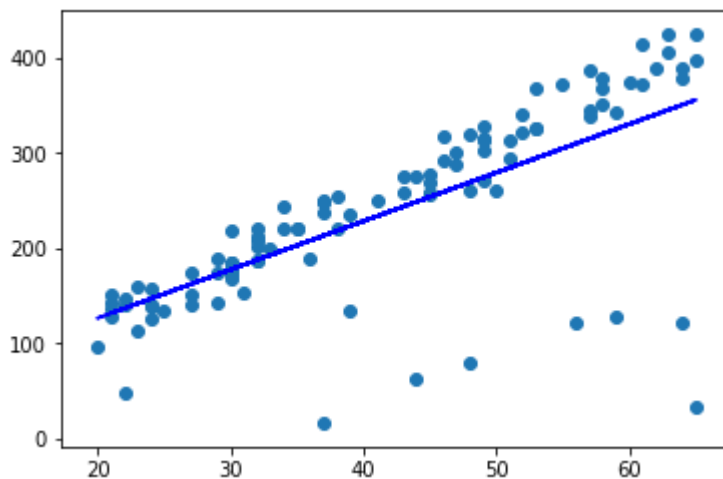
The "correct" slope for the main body of data points is 6.25 (we know this because we used this value to generate the data); what slope does your regression have?

```
In [3]:   ### fill in a regression here!  Name the regression object reg so that
          ### the plotting code below works, and you can see what your regression looks like

          from sklearn.linear_model import LinearRegression as lr
          reg = lr()
          reg.fit(ages_train,net_worths_train)

          try:
              plt.plot(ages, reg.predict(ages), color="blue")
          except NameError:
              pass
          plt.scatter(ages, net_worths)
          plt.show()

          print "coef ",reg.coef_
```



```
coef  [[5.07793064]]
```

What is the score you get when using your regression to make predictions with the test data?

```
In [4]:   reg.score(ages_test,net_worths_test)
```

```
Out[4]:   0.8782624703664671
```

In outliers/outlier_cleaner.py, you will find the skeleton for a function called outlierCleaner() that you will fill in with a cleaning algorithm. It takes three arguments: predictions is a list of predicted targets that come from your regression, ages is the list of ages in the training set, and net_worths is the actual value of the net worths in the training set. There should be 90 elements in each of these lists (because the training set has 90 points in it). Your job is to return a list called cleaned_data that has only 81 elements in it, which are the 81 training points where the predictions and the actual values (net_worths) have the smallest errors (90 * 0.9 = 81). The format of cleaned_data should be a list of tuples, where each tuple has the form (age, net_worth, error).

Once this cleaning function is working, you should see the regression result changes. What is the new slope? Is it closer to the "correct" result of 6.25?

In [ ]:
```python
# %load outlier_cleaner.py
#!/usr/bin/python


def outlierCleaner(predictions, ages, net_worths):
    """
        Clean away the 10% of points that have the largest
        residual errors (difference between the prediction
        and the actual net worth).

        Return a list of tuples named cleaned_data where
        each tuple is of the form (age, net_worth, error).
    """

    cleaned_data = []

    ### your code goes here
    import numpy as np

    err = (net_worths - predictions) ** 2
    cleaned_data = zip(ages, net_worths, err)
    cleaned_data = sorted(cleaned_data,key=lambda x: x[2][0], reverse=True)
    lim = int(len(net_worths)*0.1)

    return cleaned_data[lim:]
```

In [11]:
```python
### identify and remove the most outlier-y points
cleaned_data = []
try:
    predictions = reg.predict(ages_train)
    cleaned_data = outlierCleaner( predictions, ages_train, net_worths_train )
except NameError:
    print "your regression object doesn't exist, or isn't name reg"
    print "can't make predictions to use in identifying outliers"
### only run this code if cleaned_data is returning data
if len(cleaned_data) > 0:
    ages, net_worths, errors = zip(*cleaned_data)
    ages = numpy.reshape( numpy.array(ages), (len(ages), 1))
    net_worths = numpy.reshape( numpy.array(net_worths), (len(net_worths), 1))

    ### refit your cleaned data!
    try:
        reg.fit(ages, net_worths)
        plt.plot(ages, reg.predict(ages), color="blue")
    except NameError:
        print "you don't seem to have regression imported/created,"
        print "   or else your regression object isn't named reg"
        print "   either way, only draw the scatter plot of the cleaned data"
    plt.scatter(ages, net_worths)
    plt.xlabel("ages")
    plt.ylabel("net worths")
    plt.show()


else:
    print "outlierCleaner() is returning an empty list, no refitting to be done"
```
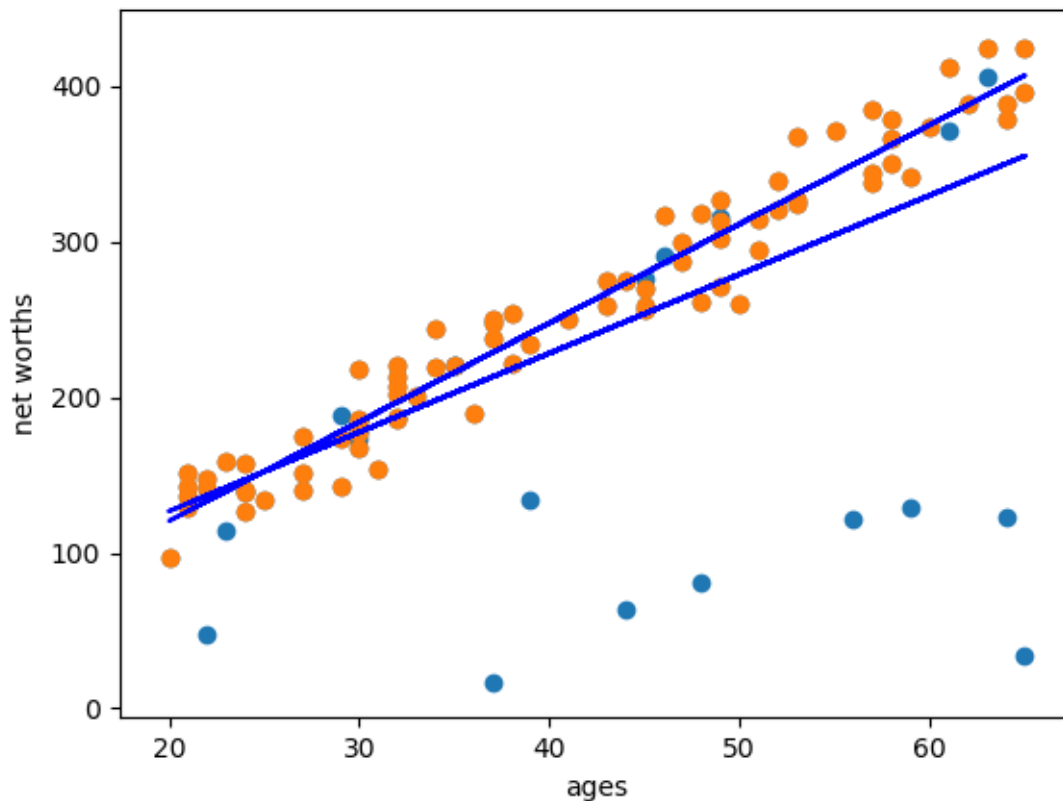
outlierCleaner() is returning an empty list, no refitting to be done


???? works in pycharm

In [16]:
```python
from IPython.display import Image
Image("Figure_1.png")
```

Out[16]:



In [17]:
```python
Image("Capture.png")
```

Out[17]:



```
In[2]: reg.coef_
Out[2]: array([[6.36859481]])
```

What's the new score when you use the regression to make predictions on the test set?

In [18]:
```python
reg.score(ages_test,net_worths_test)
```

Out[18]:  0.8782624703664671

In the mini-project for the regressions lesson, you used a regression to predict the bonuses for Enron employees. As you saw, even a single outlier can make a big difference on the regression result. There was something we didn't tell you, though, which was that the dataset we had you use in that project had already been cleaned of some significant outliers. Identifying and cleaning away outliers is something you should always think about when looking at a dataset for the first time, and now you'll get some hands-on experience with the Enron data.

You can find the starter code in outliers/enron_outliers.py, which reads in the data (in dictionary form) and converts it into a sklearn-ready numpy array. Since there are two features being extracted from the dictionary ("salary" and "bonus"), the resulting numpy array will be of dimension N x 2, where N is the number of data points and 2 is the number of features. This is perfect input

for a scatterplot; we'll use the matplotlib.pyplot module to make that plot. (We've been using pyplot for all the visualizations in this course.) Add these lines to the bottom of the script to make your scatterplot:

```
In [20]: # %load enron_outliers.py
#!/usr/bin/python

import pickle
import sys
import matplotlib.pyplot
sys.path.append("../tools/")
from feature_format import featureFormat, targetFeatureSplit


### read in data dictionary, convert to numpy array
data_dict = pickle.load( open("../final_project/final_project_dataset.pkl", "r")
features = ["salary", "bonus"]
data = featureFormat(data_dict, features)


### your code below
print data.max()
for point in data:
    salary = point[0]
    bonus = point[1]
    matplotlib.pyplot.scatter( salary, bonus )

matplotlib.pyplot.xlabel("salary")
matplotlib.pyplot.ylabel("bonus")
matplotlib.pyplot.show()
```
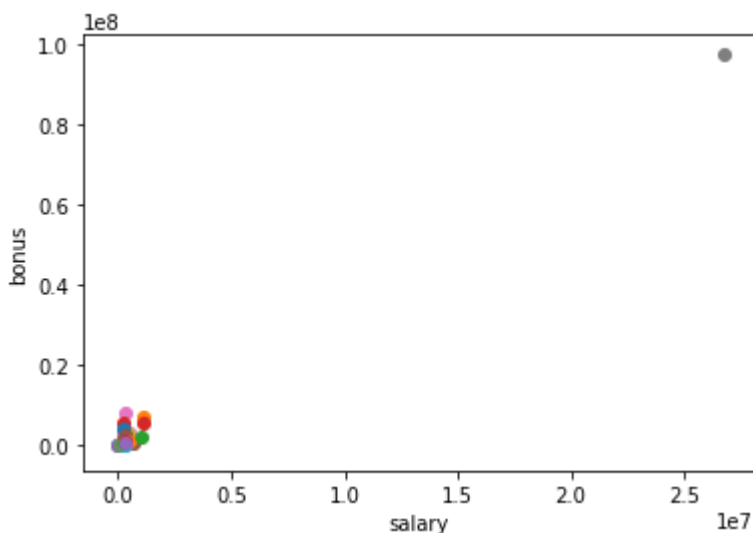
97343619.0



There's one outlier that should pop out to you immediately. Now the question is to identify the source. We found the original data source to be very helpful for this identification; you can find that PDF in final_project/enron61702insiderpay.pdf What's the name of the dictionary key of this data point? (e.g. if this is Ken Lay, the answer would be "LAY KENNETH L").

the Total amt

Does this outlier seem like a data point that we should include when running machine learning on this dataset? Or should we remove it?

remove

A quick way to remove a key-value pair from a dictionary is the following line: dictionary.pop( key, 0 ) Write a line like this (you'll have to modify the dictionary and key names, of course) and remove the outlier before calling featureFormat(). Now rerun the code, so your scatterplot doesn't have this outlier anymore. Are all the outliers gone?

In [22]:
```python
# %load enron_outliers.py
#!/usr/bin/python

import pickle
import sys
import matplotlib.pyplot
sys.path.append("../tools/")
from feature_format import featureFormat, targetFeatureSplit


### read in data dictionary, convert to numpy array
data_dict = pickle.load( open("../final_project/final_project_dataset.pkl", "r")
features = ["salary", "bonus"]

data_dict.pop('TOTAL',0)

data = featureFormat(data_dict, features)


### your code below
print data.max()
for point in data:
    salary = point[0]
    bonus = point[1]
    matplotlib.pyplot.scatter( salary, bonus )

matplotlib.pyplot.xlabel("salary")
matplotlib.pyplot.ylabel("bonus")
matplotlib.pyplot.show()
```
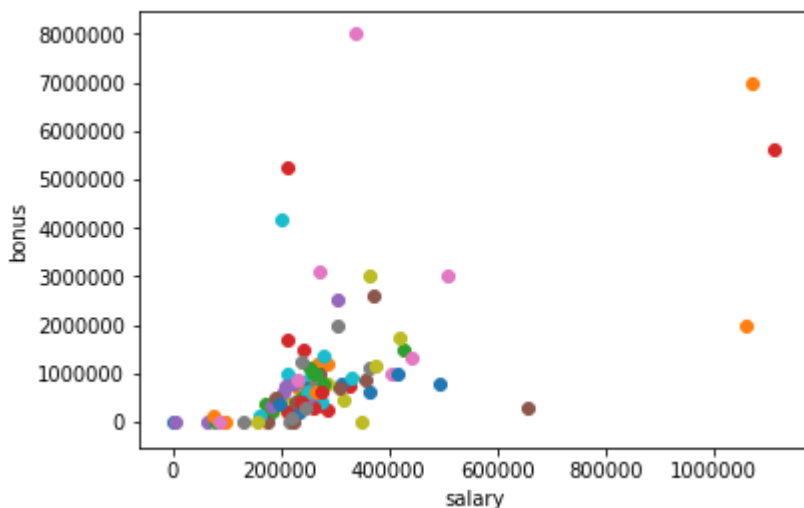
8000000.0



A quick way to remove a key-value pair from a dictionary is the following line: dictionary.pop( key, 0 ) Write a line like this (you'll have to modify the dictionary and key names, of course) and remove the outlier before calling featureFormat(). Now rerun the code, so your scatterplot doesn't have this outlier anymore. Are all the outliers gone?

no

We would argue that there's 4 more outliers to investigate; let's look at a couple of them. Two people made bonuses of at least 5 million dollars, and a salary of over 1 million dollars; in other words, they made out like bandits. What are the names associated with those points?

jeff skilling

ken lay

Would you guess that these are typos or weird spreadsheet lines that we should remove, or that there's a meaningful reason why these points are different? (In other words, should they be removed before we, say, try to build a POI identifier?)

leave in

In [ ]: