# Predicting Murder Rates With Gradient Descent

### Goal

The goal of this project is to extend the techniques we learned in class, for one-dimensional gradient descent, to a two-dimensional domain space, solving a *linear regression problem*. This problem is also known as *curve fitting*. As part of this lab, you will learn how to compute with vectors instead of scalars. Please use file *userid*-regression/minimize.py for function *userid*-regression/plot\_regression.py for the code that draws the trace of the minimization in action (i.e., this is the file that has all the matplotlib junk). You'll see starter (and test) files in the regression starterkit. You will be doing your work in directory *userid*-regression.

# Discussion

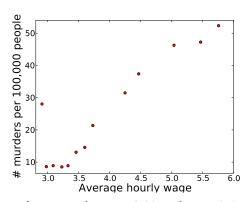
#### Problem statement

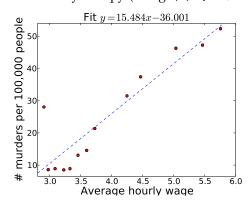
Given training data  $(x_i, y_i)$  for i = 1..n samples with dependent variable  $y_i$ , we would like to predict y for some x's not in our training set.  $x_i$  is generally a vector of independent variables but we'll use a scalar. If we assume there is a linear relationship between x and y, then we can draw a line through the data and predict future values with that line function. To do that, we need to compute the two parameters of our model: a slope and a y intercept. (We will see the model below.)

For example, if we compare the number of murders per 100,000 people in Detroit to the average hourly wage, our eyeballs easily detect a correlation. Here is data suitable to copy and paste into Python:

```
HOURLY_WAGE = [2.98, 3.09, 3.23, 3.33, 3.46, 3.6, 3.73, 2.91, 4.25, 4.47, 5.04, 5.47, 5.76]
MURDERS = [8.6, 8.9, 8.52, 8.89, 13.07, 14.57, 21.36, 28.03, 31.49, 37.39, 46.26, 47.24, 52.33]
```

and here is a scatter plot and best fit line as determined by numpy (using np.polyfit(HOURLY\_WAGE, MURDERS,1)).





Here, for example,  $x_0 = 2.98$  and  $y_0 = 8.6$ .

This might be a good point to remind everyone that correlation does not equal causation. I hardly think that paying people more makes them murderous, although I could see the opposite. ;) Correlation is a necessary

but not *sufficient* condition for causation. When you find a correlation, that gives you a candidate to check for cause-and-effect.

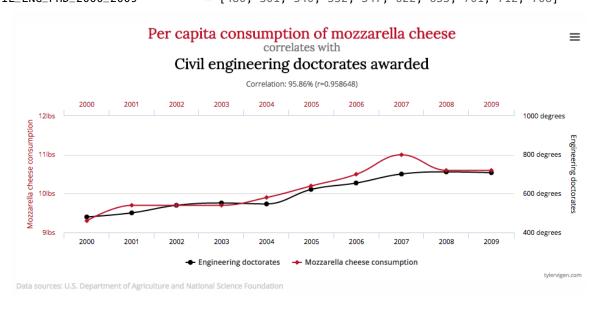
Here's another fun data set from http://www.tylervigen.com/spurious-correlations:

# http://www.census.gov/compendia/statab/2012/tables/12s0217.xls

MOZZ\_CHEESE\_PER\_CAPITA\_2000\_2009 = [9.3, 9.7, 9.7, 9.7, 9.9, 10.2, 10.5, 11.0, 10.6, 10.6]

# http://www.nsf.gov/statistics/infbrief/nsf11305/

CIVIL\_ENG\_PHD\_2000\_2009 = [480, 501, 540, 552, 547, 622, 655, 701, 712, 708]



#### Best fit line that minimizes squared error

Recall the formula for a line from high school: y = mx + b. We normally rewrite that using elements of a vector coefficients,  $\vec{B}$ , in preparation for describing it with vector notation from linear algebra. For simplicity, though, we'll stick with scalar coefficients for now:

$$\hat{y} = b_2 x + b_1$$

We use notation  $\hat{y}$  to indicate that it approximates y from our data set but is not necessarily equal to any of the  $y_i$ .

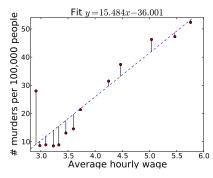
The "best line" is one that minimizes some cost function that compares the known each  $y_i$  value at  $x_i$  to the predicted  $\hat{y}$  of the linear model that we conjure up using parameters  $b_1$ ,  $b_2$ . A good measure is the *sum* of squared errors. The cost function adds up all of these squared errors to tell us how good of a fit our linear model is:

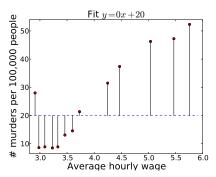
$$Cost(B) = \sum_{i=1}^{n} (\underbrace{\hat{y}}_{\text{linear model}} - \underbrace{y_i}_{\text{true value}})^2$$

Inlining the expression for our model, we get:

$$Cost(B) = \sum_{i=1}^{n} (b_2 x_i + b_1 - y_i)^2$$

As we wiggle the linear model parameters, the value of the cost function will change. The following graphs shows the errors/residuals that are squared and summed to get the overall cost for two different "curve fits."





The costs are 533.82 for the left and 3563.50 for the right.

The good news is that we know the cost function is a quadratic (by construction), which is convex and has an exact solution. All we have to do is figure out where the cost function flattens out. Mathematically, that is when the partial derivatives of the cost function are both zero. The partial derivative for a particular dimension is just the rate of change of that dimension at a particular spot. It's like a skier examining the slope of her skis against the mountain one at a time.

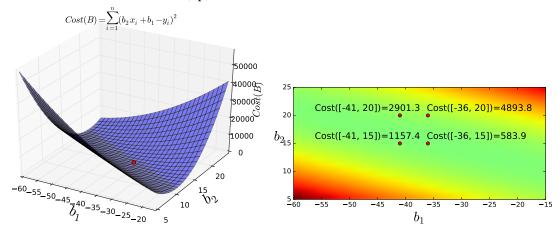
$$\nabla Cost(B) = 0$$
 (Analytic solution to optimization)

For our purposes, though, we'll use gradient descent to minimize the cost function. There are lots and lots of cost functions that are not simple little quadratics with symbolic solutions from calculus. In those cases, we have to use computer simulation.

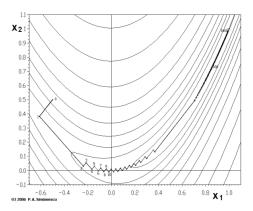
To show our prediction model in action, we can ask how many murders there would be in Detroit if the average wage were \$4.7? (Obviously, these wages are from 30 years ago.) To make a prediction, all we have to do is plug x = 4.7 into  $\hat{y} = 14.484x - 36.001$ , which gives us 32.074 murders.

## Gradient descent in 3D

Before trying to minimize the cost function, it's helpful to study what the surface looks like in three dimensions, as shown in the following two graphs. The X and Y dimensions are the coefficients of our linear model and the Z coordinate (up) is the cost function.



What surprised me is that changes to the slope of the linear model's coefficient  $b_2$ , away from the optimal  $b_2 = 15.484$ , cost much more than tweaks to the y intercept,  $b_1$ . Regardless, the surface is convex. Unfortunately, based upon the deep trough that grows slowly along the diagonal of  $(b_1, b_2)$ , gradient descent takes a while to converge to the minimum. We will examine the path of gradient descent for a few initial starting point. Wikipedia says that the Rosenbrock function is a pathological case for traditional gradient descent and it looks pretty similar to our surface with its shallow valley in this topographic view from above:



The recurrence relation for updating our estimate of  $\vec{B} = [b_1, b_2]$  that minimizes  $Cost(\vec{B})$  is the same as the single-variable gradient-descent we did in class but with a vector instead of a scalar:

$$\vec{B}_{i+1} = \vec{B}_i - \eta \nabla Cost(\vec{B}_i)$$

where we will approximate vectors of partial derivatives with partial finite differences:

$$\nabla F(\vec{X}) = \begin{bmatrix} \frac{\partial}{x_1} F(\vec{X}) \\ \frac{\partial}{x_2} F(\vec{X}) \end{bmatrix} \approx \begin{bmatrix} \frac{F(\begin{bmatrix} x_1 + h \\ x_2 \end{bmatrix}) - F(\vec{X})}{h} \\ \frac{F(\begin{bmatrix} x_1 \\ x_2 + h \end{bmatrix}) - F(\vec{X})}{h} \end{bmatrix}$$

The h variable is a very small delta step away from a specific coordinate in each dimension in turn. In our case, we will compute the components of a finite difference vector C' ignoring the division by the step h because it will get folded into our learning rate  $\eta$ .

Gradient-descent algorithm

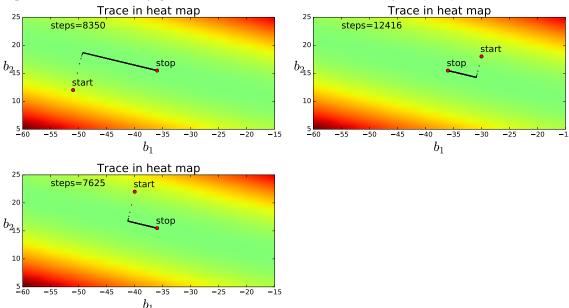
The minimization algorithm looks like:

- 1. Pick an initial  $B_0$
- 2. Let  $B = B_0$

3. 
$$C' = \begin{bmatrix} Cost(\begin{bmatrix} b_1+h \\ b_2 \end{bmatrix}) - Cost(B_i) \\ Cost(\begin{bmatrix} b_1 \\ b_2+h \end{bmatrix}) - Cost(B_i) \end{bmatrix}$$
, the finite difference in each direction computed separately

- 4. Let  $B_{i+1} = B_i \eta C'$
- 5. Goto step 3 until  $abs(Cost(B_{i+1}) Cost(B_i)) < precision$  and when the Cost starts going back up

Using a small learning rate, my solution takes 78604 steps starting from coordinate (-45,10) using a very small step size, which gives me a fairly decent approximation of the minimum: [-36.00063291 15.48437211] compared to the analytic solution [ -36.000625 15.484375 ]. Starting from about the same distance away in the shallow valley at (-45,25), my solution takes 81551 steps. Cutting my step size by 5x, takes 431803 steps from (-45,25). Surprisingly, when I start very close to the minimum at (-36,15), my solution takes 71963 steps and does not exactly give a more accurate result. "Your mileage may vary."



When I crank up the learning rate and use a very small step size, my solution converges must faster and with the same accuracy. For example, with 10 times the learning rate as before, (-45,25) converges in 8174 steps instead of 81551 steps.

# Getting fancy

For the most part, the algorithm described above converges on the analytic solution with a high degree of precision, assuming that your h is small enough. In some cases however it only gets 3 digits of precision versus about six, depending on where the starting point is. To improve that accuracy, we can improve our termination condition from simply comparing previous and new value of the cost function to use the Euclidean distance between  $B_{i+1}$  and  $B_i$  or the magnitude of the finite difference.

To increase convergence rate, it's a good idea to have different  $\eta$  learning rates per dimension, particularly in this case because we want to speed along the valley of  $b_1$ , but not jump too far on  $b_2$  as it is already steep; the finite difference will cause it to move along quickly. *Note:* Because numpy will overload the \* multiple operator to work with vectors as well as scalars,  $\eta$  can be a 2-vector and  $\eta C'$  will "just work."

For extremely large data sets or very high dimension data sets, you can use *stochastic gradient descent* which takes a single random or a small random subset of the samples when computing the cost function. Surprisingly, this still converges and takes much less time to compute. I believe that the math says that the algorithm will converge no matter how many samples (don't quote me on that). Further, picking a random subset adds noise to the process, which has been shown to pull gradient descent out of local minima and get it moving again towards a lower minimum. In our case, there is an exact solution and therefore a single minimum so we would only get a speed benefit from stochastic gradient descent.

## Your task

You will use gradient descent to solve the linear regression problem above, using the same data. As part of your final submission, you must provide heat maps with traces that indicate the steps taken by your gradient descent as I have shown above. Have your program choose two random starting  $B_0$  vectors (you can use random.randrange()) to produce your heat maps. In *userid*-regression/minimize.py, define a function called minimize that takes the indicated parameters and returns the minimum B parameters of your linear model, the number of steps, and the trace array of intermediate  $B_i$  values (B is a list of 2-vectors).

```
def minimize(f, B0, eta, h, precision):
    trace = []
    B = B0
    steps = 0
    while True:
        steps += 1
        if steps % 10 == 0: # only capture every 10th value
            trace.append(B)
        ...
    return (B, steps, trace)
```

As an example, in my userid-regression/plot\_regression.py script, I call that function like this:

```
(m,steps,trace) = minimize(Cost, B0, LEARNING_RATE, h, PRECISION)
```

Use pylab.imshow() to draw the heat map whose  $b_1$  are the y intercepts,  $b_2$  coordinates are the slopes, and heat value is the cost of row, col. The row is the  $b_2$  slope values and the col is the  $b_1$  y-intercept values. Hide all of your heat map construction in function in userid-regression/plot\_regression.py:

```
def heatmap(f, trace=None): # trace is a list of [b1, b2] pairs
    ...
so we can call it like this:
heatmap(Cost, trace)
plt.show() # finally show the heatmap
```

Your heatmap function will create two ranges for the  $b_1$  and  $b_2$  coordinates and then simply compute the cost, storing it in a matrix such as C[i][j] = Cost([b1,b2]) where i is the index of  $b_2$  range values, .... The trace is plotted on top of the map once you create it. The coordinate ranges of the heat map are not

related to the trace at all; just show the ranges in the heatmaps displayed in this document. It took me a while to figure out all of the crazy methods to draw the heat maps so I'll toss you a bone and just provide it for you (it goes in your heatmap() function):

Please show the information as I have shown in the graphs to make it easier to compare results and for me to grade.

You will have to pick an appropriately small step value h to get a decent approximation of the derivative through finite differences. You want that number to be small enough that your algorithm does not oscillate around the minimum. If the number is too big it will compute a finite difference that leads to  $B_{i+1}$  leaping across the minimum to the other wall of the function. You must pick a learning rate  $\eta$  that allows you to go as fast as you can but not so fast that it overruns the minimum back and forth. When I crank up my learning rate too far, I also see the algorithm go off into the weeds and stops with a minimum of [ -4.86000929e+10 -2.85744570e+12].

# **Testing**

You can run *userid*-regression/test\_regression.py to check your miminize function. Here's my sample run:

```
$ python test_regression.py
starting at [-23, 12]
gradient descent gives [-36.00034331 15.48430239] Cost 533.815093809
exact is [-36.000625000000007, 15.484375]
starting at [-50, 10]
gradient descent gives [-36.00062579 15.48437039] Cost 533.815093755
exact is [-36.000625000000007, 15.484375]
starting at [-36, 17]
gradient descent gives [-36.00062582 15.4843704 ] Cost 533.815093755
exact is [-36.000625000000007, 15.484375]
starting at [-36, 15]
gradient descent gives [-36.00034331 15.48430239] Cost 533.815093809
exact is [-36.000625000000007, 15.484375]
starting at [100, 17]
gradient descent gives [-36.00034316 15.48430235] Cost 533.815093809
exact is [-36.000625000000007, 15.484375]
starting at [-15, 0]
gradient descent gives [-36.00034325 15.48430237] Cost 533.815093809
exact is [-36.000625000000007, 15.484375]
```

#### Resources

There is a lot of material out there on the web that can be helpful.

- Finite difference at Wikipedia
- Stochastic Gradient Descent Tricks
- Numerical recipes (See Chap 10 on minimization of functions)
- Single verbal minimization in line searches
- Andrew Ng's CS229 Lecture notes
- Data analysis with Python

You must tweak the step size and other parameters so that your results agree with the first three decimal points of the analytic solution [-36.00062500000007, 15.484375]. (My solution is much better than that, except for a couple of weird starting positions where it only gets three decimal places.)

- Your script *userid*-regression/minimize.py with the minimize() function.
- Your script
   userid-regression/plot\_regression.py, which
   must have your heatmap function. This file
   imports the previous.



• A PDF or PDFs of your graph with two visible traces on two heat maps or the same heat map if the traces are clear. The graph(s) should include the start, stop location and the number of steps as I have done on mine. As part of your PDF(s), please indicate the *B* parameters you compute with your minimize function. Please name the file plot.pdf, or if two, plot1.pdf and plot2.pdf.