

Gradient Descent

Overview

In this notebook, we will solve a simple linear regression problem by gradient descent.

We will see the effect of the learning rate on the trajectory in parameter space. We will show how Stochastic Gradient Descent (SGD) differs from the standard version, and the effect of "shuffling" your data during SGD.

```
In [3]: # Preliminaries - packages to load
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
%matplotlib inline
```

Generate Data from a known distribution

Below we will generate data a known distribution.
Specifically, the true model is:

$$Y = b + \theta_1 X_1 + \theta_2 X_2 + \epsilon$$

X_1 and X_2 have a uniform distribution on the interval $[0, 10]$, while const is a vector of ones (representing the intercept term).

We set actual values for b , θ_1 , and θ_2

Here $b = 1.5$, $\theta_1 = 2$, and $\theta_2 = 5$

We then generate a vector of y-values according to the model and put the predictors together in a "feature matrix" X_mat :

```
In [8]: np.random.seed(1234) ## This ensures we get the same data if all of the other parameters remain fixed

num_obs = 100
x1 = np.random.uniform(0,10,num_obs) # values btwn 0 and 10 and we want 100 obs
x2 = np.random.uniform(0,10,num_obs) # values btwn 0 and 10 and we want 100 obs
const = np.ones(num_obs) # const creates a array of 1's
eps = np.random.normal(0,.5,num_obs) # epsilon(error term) with mean of 0 and a standard deviation of 0.5,

b = 1.5
theta_1 = 2
theta_2 = 5

#y=b+theta_1*x1+theta_2*x2+eps

y = b*const + theta_1*x1 + theta_2*x2 + eps

x_mat = np.array([const,x1,x2]).T

Out[8]: array([[1., 1.9151945, 7.67116628],
 [1., 6.22198771, 7.08115362],
 [1., 4.13727739, 7.06871341],
 [1., 7.85358584, 5.57768828],
 [1., 7.79758088, 9.05836532],
 [1., 7.72592695, 1.471589 ],
 [1., 2.76464255, 0.29647881],
 [1., 0.01872378, 9.93093493],
 [1., 9.58139354, 1.14065699],
 [1., 8.75826335, 9.5089895 ],
 [1., 3.5781727 , 3.25787414],
 [1., 5.00995126, 1.9361869 ],
 [1., 6.83469355, 4.57611649],
 [1., 7.12762827, 9.26402571],
 [1., 3.76256755, 8.79689162],
 [1., 6.61186186, 5.52615755],
 [1., 5.03883165, 3.48088793],
 [1., 3.178845 , 5.2588732],
 [1., 7.72826622, 9.01798051],
 [1., 8.26411391, 7.06528163],
 [1., 6.64889884, 7.26038462],
 [1., 6.15396178, 9.0087837],
 [1., 9.75382424, 7.91638991],
 [1., 3.69824866, 5.00154781],
 [1., 9.33140102, 2.91125245],
 [1., 6.51378143, 1.51389264],
 [1., 3.97202578, 3.35174659],
 [1., 7.86736143, 6.97351777],
 [1., 1.0836122 , 0.73424444],
 [1., 6.68989653, 0.55066395],
 [1., 6.6912739 , 3.23194514],
 [1., 4.36173424, 5.99481884],
 [1., 8.02147642, 8.53898587],
 [1., 4.37869325, 2.87082425],
 [1., 7.04260971, 1.73867227],
 [1., 9.04891308, 1.34021206],
 [1., 2.18792106, 9.94653829],
 [1., 9.24867629, 1.79497889],
 [1., 4.42140755, 3.17648923],
 [1., 9.09315959, 5.68291405],
 [1., 0.59899233, 0.09348575],
 [1., 1.84287084, 9.06648621],
 [1., 6.4735275 , 7.7241431],
 [1., 6.74889844, 5.56934679],
 [1., 5.9482478 , 0.84773843],
 [1., 5.33316517, 8.53921466],
 [1., 0.43324963, 7.28428676],
 [1., 5.6143398 , 1.42435373],
 [1., 3.29686446, 5.52468939],
 [1., 5.02968833, 2.7304326 ],
 [1., 1.11894319, 9.74495138],
 [1., 6.07193706, 6.67786906],
 [1., 6.65944643, 2.55653286],
 [1., 0.00764062, 1.00311484],
 [1., 6.17441709, 7.76180723],
 [1., 9.12122886, 7.82477993],
 [1., 7.8624133 , 6.01603914],
 [1., 9.92881466, 9.14403113],
 [1., 0.58991762, 6.58622782],
 [1., 7.91964135, 5.68367582],
 [1., 2.8525096 , 2.01755892],
 [1., 4.2493705 , 6.0826376 ],
 [1., 4.78993796, 9.5219541 ],
 [1., 1.98675179, 0.89063287],
 [1., 3.82317452, 9.93667363],
 [1., 6.53873085, 8.1870351 ],
 [1., 4.51644808, 5.45122166],
 [1., 9.82804742, 4.51254955],
 [1., 2.129427 , 0.90557188],
 [1., 1.19388898, 9.73264791],
 [1., 7.38523956, 5.9341133 ],
 [1., 8.87383633, 3.6674498 ],
 [1., 4.71632534, 3.23094693],
 [1., 1.07126817, 8.71423255],
 [1., 2.29218565, 2.10634063],
 [1., 9.9996195 , 7.34945189],
 [1., 4.16753538, 3.69319887],
 [1., 5.38851663, 8.01602599],
 [1., 0.06185875, 1.82735992],
 [1., 3.06641706, 7.01355379],
 [1., 4.36893172, 6.22776587],
 [1., 6.12148091, 4.93883446],
 [1., 9.18196075, 8.405377 ],
 [1., 6.2573667 , 7.12098987],
 [1., 6.9597565 , 4.43989881],
 [1., 1.49833716, 6.31834861],
 [1., 7.46985406, 3.6323976 ],
 [1., 8.31086992, 7.30721791],
 [1., 6.33725769, 4.75866372],
 [1., 4.38389881, 5.4441687 ],
 [1., 1.52572775, 6.40880435],
 [1., 5.68408635, 1.26285322],
 [1., 5.28224278, 1.71465261],
 [1., 5.51426764, 7.37086494],
 [1., 4.80389179, 1.27023384],
 [1., 5.62559563, 3.69649875],
 [1., 5.36878193, 0.04334065],
 [1., 8.19202067, 1.03104439],
 [1., 6.57115638, 8.02374182],
 [1., 6.65412742, 9.45032286]])
```

```
In [9]: y

Out[9]: array([44.00960834, 49.44119069, 50.57415314, 45.58928189, 65.35593861,
 14.93493564, 60.04285882, 46.6123024 , 26.42593365, 66.44763573,
 26.00356286, 21.26197374, 37.35515258, 62.48566207, 51.78454565,
 24.68094522, 29.14388523, 30.89742951, 62.68253254, 53.75444867,
 44.52288889, 58.51638392, 41.78856245, 38.12132178, 34.82375165,
 21.80988388, 25.46622631, 49.78398141, 12.05682541, 15.39651887,
 31.96172805, 40.19232745, 60.38296962, 19.20268827, 24.07886969,
 21.98268992, 55.68703273, 28.68851837, 26.74897285, 47.70514593,
 9.90129951, 58.25411775, 52.26455884, 43.23633517, 17.89728884,
 25.64311941, 39.30888874, 20.98432084, 36.44677963, 25.71305238,
 52.83723555, 46.89545889, 25.82169038, 7.39534188, 52.51954727,
 59.82824402, 55.99628893, 67.50716147, 53.72035595, 44.71835284,
 17.09886971, 48.8695994 , 59.23483932, 56.83522379, 58.88538389,
 43.66214118, 37.71952691, 43.33256307, 47.88288694, 52.67558512,
 45.21155426, 32.93371396, 27.28704805, 47.54794318, 18.72818759,
 56.49688477, 28.54755755, 61.77134925, 49.06492289, 43.1575642 ,
 41.96786483, 38.62279785, 61.40832297, 49.78122273, 38.7479899 ,
 5.19228941, 33.87971392, 54.71962001, 37.95472396, 26.852549 ,
 36.17695578, 19.45541981, 29.38572724, 56.88867987, 18.22335359,
 38.13641244, 42.61877866, 23.47252364, 42.91850714, 62.70436696]])
```

Get the "Right" answer directly

In the below cells we solve for the optimal set of coefficients. Note that even though the true model is given by:

$$b = 1.5, \theta_1 = 2, \text{ and } \theta_2 = 5$$

The maximum likelihood (least-squares) estimate from a finite data set may be slightly different.

We can solve the problem two ways:

1. By using the scikit-learn `LinearRegression` model
2. Using matrix algebra directly via the formula $\theta = (X^T X)^{-1} X^T y$

Note: The scikit-learn solver may give a warning message, this can be ignored.

```
In [11]: ### Solve directly using sklearn
from sklearn.linear_model import LinearRegression

lr_model = LinearRegression(fit_intercept=False) #intercepts are included in our feature values in our x_mat that we defined earlier.
lr_model.fit(x_mat, y)

lr_model.coef_

Out[11]: array([1.49904618, 1.99675416, 5.01156315])

b=1, theta1 = 2 and theta2 = 5
```

```
In [12]: ## Solve by matrix calculation

#theta=(X^T X)^-1 X^T y

#So that's going to be the inverse of the dot product of the transpose and the value itself,
#and then the dot product of that with x transpose and then the dot product of that with y.

np.linalg.inv(np.dot(x_mat.T,x_mat)).dot(x_mat.T).dot(y)

Out[12]: array([1.49904618, 1.99675416, 5.01156315])
```

Same values as with scikit-learn

Solving by Gradient Descent

For most numerical problems, we don't / can't know the underlying analytical solution. This is because we only arrive at analytical solutions by solving the equations mathematically, with pen and paper. That is more often than not just impossible. Fortunately, we have a way of converging to an approximate solution, by using **Gradient Descent**.

lets explore this very useful method because Neural Networks, along with many other complicated algorithms, are trained using Gradient Descent. Seeing how gradient descent works on a simple example will build intuition and help us understand some of the nuances around setting the learning rate and other parameters. lets also explore Stochastic Gradient Descent and compare its behavior to the standard approach.

The next several cells have code to perform (full-batch) gradient descent.

1. Picking a learning rate, and a number of iterations, run the code, and then plot the trajectory of your gradient descent.
2. Find examples where the learning rate is too high, too low, and "just right".
3. Look at plots of loss function under these conditions.

```
In [13]: ## Parameters to play with
learning_rate = 1e-3
num_iter = 10000
theta_initial = np.array([3,3,3])
```

```
In [16]: def gradient_descent(learning_rate, num_iter, theta_initial):

    ## Initialization steps
    theta = theta_initial
    theta_path = np.zeros((num_iter+1,3))
    theta_path[0,:] = theta_initial

    loss_vec = np.zeros(num_iter)

    ## Main Gradient Descent loop (for a fixed number of iterations)
    for i in range(num_iter):
        y_pred = np.dot(theta.T,x_mat.T)
        loss_vec[i] = np.sum((y-y_pred)**2)
        grad_vec = (y-y_pred).dot(x_mat)/num_obs #sum up the gradients across all observations and divide by num_obs
        grad_vec = grad_vec
        theta = theta + learning_rate*grad_vec
        theta_path[i+1,:]=theta
    return theta_path, loss_vec

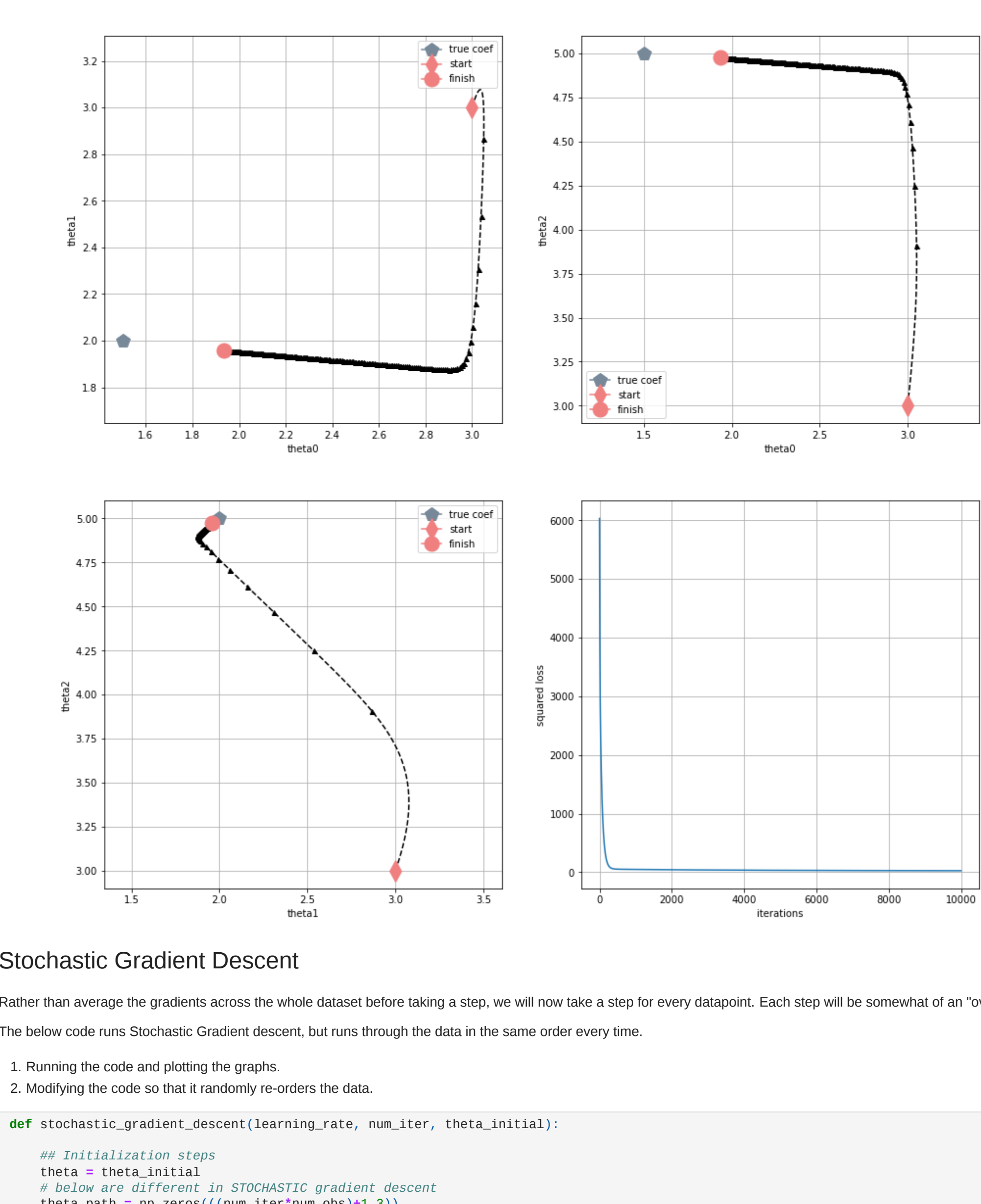
In [20]: true_coef = [b, theta_1, theta_2]

def plot_i(theta_path, i, j, ax):
    ax.plot(true_coef[i], true_coef[j],
            marker='p', markersize=15, label='true coef',
            color='r778899')
    ax.plot(theta_path[i, :], theta_path[i, j],
            color='k', linestyle='--', marker='o',
            markersize=5, markeredgecolor='r',
            markersize=5, markeredgecolor='r')
    ax.plot(theta_path[0, :], theta_path[0, j], markers='d',
            markersize=15, label='start', color='FF0000')
    ax.plot(theta_path[-1, :], theta_path[-1, j], markers='o',
            markersize=15, label='finish', color='FF0000')
    ax.set(
        xlabel='theta'+str(i),
        ylabel='theta'+str(j))
    ax.axis('equal')
    ax.grid(True)
    ax.legend(locs='best')
```

```
def plot_all(theta_path, loss_vec, learning_rate, num_iter, theta_initial, gdtype='Gradient Descent'):
    fig = plt.figure(figsize=(16, 16))
    title = '(gdtype) in the 3d parameter space - Learning rate is {lr} // {iters} // starting point {initial}'
    title = title.format(gdtype=gdtype, lr=learning_rate, iters=num_iter, initial=theta_initial)
    fig.suptitle(title, fontsize=20)
    ax = fig.add_subplot(2, 2, 1)
    plot_i(theta_path, 0, 1, ax)
    ax = fig.add_subplot(2, 2, 2)
    plot_i(theta_path, 0, 2, ax)
    ax = fig.add_subplot(2, 2, 3)
    plot_i(theta_path, 1, 1, ax)
    ax = fig.add_subplot(2, 2, 4)
    plot_i(theta_path, 1, 2, ax)
    ax.set_xlabel='iterations', ylabel='squared loss')
    ax.grid(True)
```

theta_path, loss_vec = gradient_descent(learning_rate, num_iter, theta_initial)
plot_all(theta_path, loss_vec, learning_rate, num_iter, theta_initial)

Gradient Descent in the 3d parameter space - Learning rate is 0.001 // 10000 iters // starting point [3 3 3]



Stochastic Gradient Descent

Rather than average the gradients across the whole dataset before taking a step, we will now take a step for every datapoint. Each step will be somewhat of an "overreaction" but they should average out.

The below code runs Stochastic Gradient descent, but runs through the data in the same order every time.

1. Modifying the code and plotting the graphs.
2. Running the code so that it randomly re-orders the data.

```
In [23]: def stochastic_gradient_descent(learning_rate, num_iter, theta_initial):

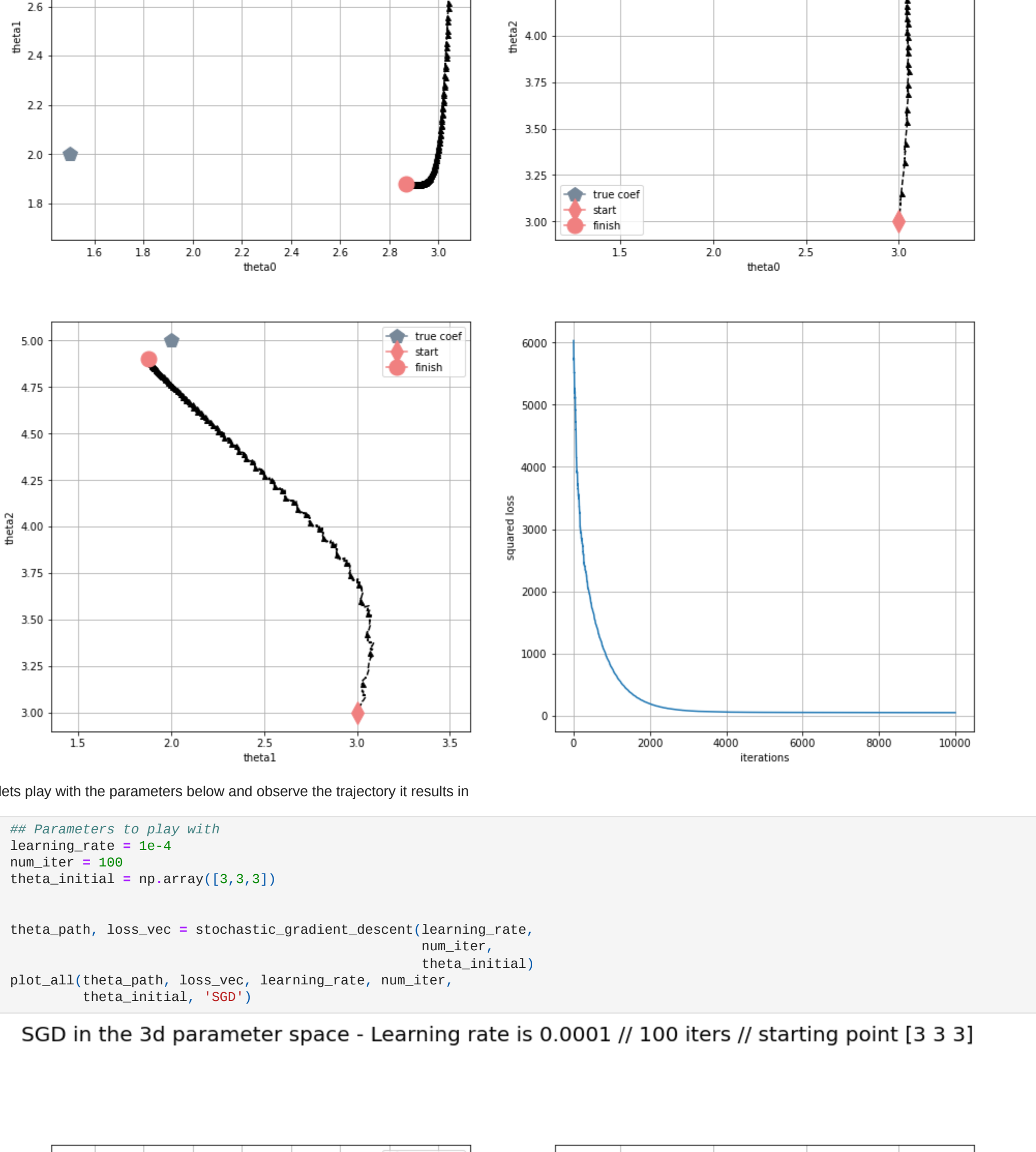
    ## Initialization steps
    theta = theta_initial
    # below are different in STOCHASTIC gradient descent
    theta_path = np.zeros((num_iter+1,num_obs+1,3))
    theta_path[0,:] = theta_initial
    loss_vec = np.zeros(num_iter+1,num_obs)

    ## Main SGD loop
    count = 0
    for i in range(num_iter):
        for j in range(num_obs):
            count+=1
            y_pred = np.dot(theta.T,x_mat.T)
            loss_vec[count-1] = np.sum((y-y_pred)**2)
            grad_vec = (y[j]-y_pred[j])*(x_mat[j,:])
            theta = theta + learning_rate*grad_vec
            theta_path[count,:]=theta
        return theta_path, loss_vec

    ## Parameters to play with
    learning_rate = 1e-4
    num_iter = 100
    theta_initial = np.array([3,3,3])

theta_path, loss_vec = stochastic_gradient_descent(learning_rate,
num_iter,
theta_initial)
plot_all(theta_path, loss_vec, learning_rate,
num_iter,
theta_initial, 'SGD')
```

SGD in the 3d parameter space - Learning rate is 0.0001 // 100 iters // starting point [3 3 3]



lets play with the parameters below and observe the trajectory it results in

```
In [26]: ## Parameters to play with
learning_rate = 1e-4
num_iter = 100
theta_initial = np.array([3,3,3])

theta_path, loss_vec = stochastic_gradient_descent(learning_rate,
num_iter,
theta_initial)
plot_all(theta_path, loss_vec, learning_rate,
num_iter,
theta_initial, 'SGD')
```

SGD in the 3d parameter space - Learning rate is 0.0001 // 100 iters // starting point [3 3 3]

