

# index

April 19, 2022

## 1 The Curse of Dimensionality - Lab

### 1.1 Introduction

In this lab, you'll conduct some mathematical simulations to further investigate the consequences of the curse of dimensionality.

### 1.2 Objectives

In this lab you will:

- Create and interpret a visual demonstrating how sparsity changes with n for n-dimensional spaces
- Demonstrate how training time increases exponentially as the number of features increases

### 1.3 Sparseness in N-Dimensional Space

As discussed, points in n-dimensional space become increasingly sparse as the number of dimensions increases. To demonstrate this, you'll write a function to calculate the Euclidean distance between two points. From there, you'll then generate random points in n-dimensional space, calculate their average distance from the origin, and plot the relationship between this average distance and n.

### 1.4 Euclidean Distance

To start, write a function which takes two points, p1 and p2, and returns the Euclidean distance between them. Recall that the Euclidean distance between two points is given by:

$$d(a, b) = \sqrt{(a_1 - b_1)^2 + (a_2 - b_2)^2 + \dots + (a_n - b_n)^2}$$

```
[ ]: import numpy as np
```

```
[1]: def euclidean_distance(p1, p2):  
    p1 = np.array(p1)  
    p2 = np.array(p2)  
    diff = p1 - p2  
    dis = np.sqrt(np.sum(np.power(diff, 2)))  
    return dis  
    # Your code here
```

## 1.5 Average Distance From the Origin

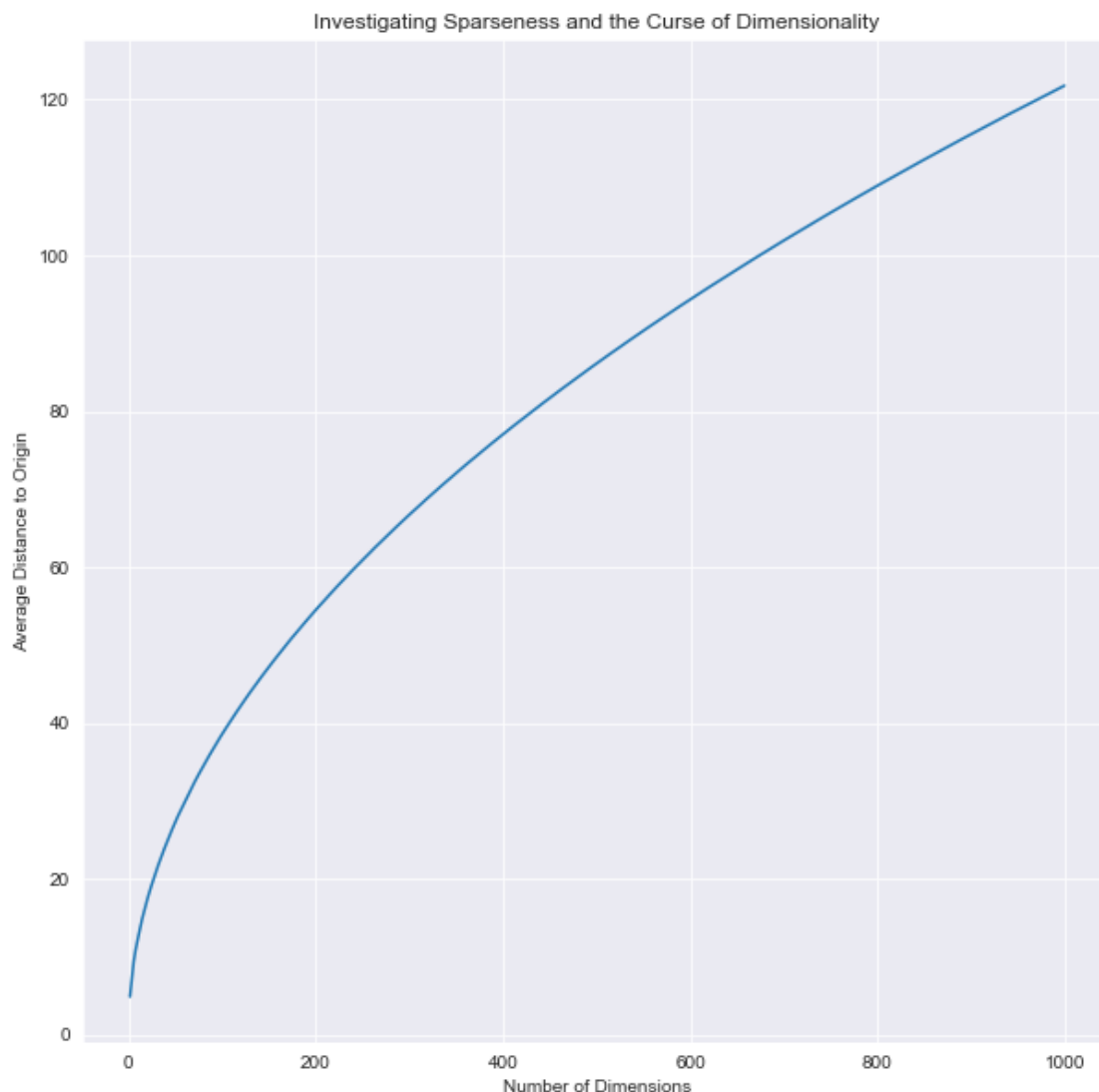
To examine the curse of dimensionality, you'll investigate the average distance to the center of n-dimensional space. As you'll see, this average distance increases as the number of dimensions increases. To investigate this, generate 100 random points for various n-dimensional spaces. Investigate n-dimensional spaces from n=1 to n=1000. In each of these, construct the 100 random points using a random number between -10 and 10 for each dimension of the point. From there, calculate the average distance from each of these points to the origin. Finally, plot this relationship on a graph; the x-axis will be n, the number of dimensions, and the y-axis will be the average distance from the origin.

```
[8]: import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
sns.set_style('darkgrid')
```

```
[16]: # Your code here
ave_dis = []
ave = []

for i in range(1, 1001):
    for j in range(101):
        p1 = np.random.uniform(low = -10, high = 10, size = i)
        p2 = np.zeros(p1.shape)
        ave_dis.append(euclidean_distance(p1, p2))
    ave.append(np.mean(np.array(ave_dis)))
```

```
[18]: plt.figure(figsize=(10, 10))
plt.plot(range(1, 1001), ave)
plt.xlabel('Number of Dimensions')
plt.ylabel('Average Distance to Origin')
plt.title('Investigating Sparseness and the Curse of Dimensionality');
```

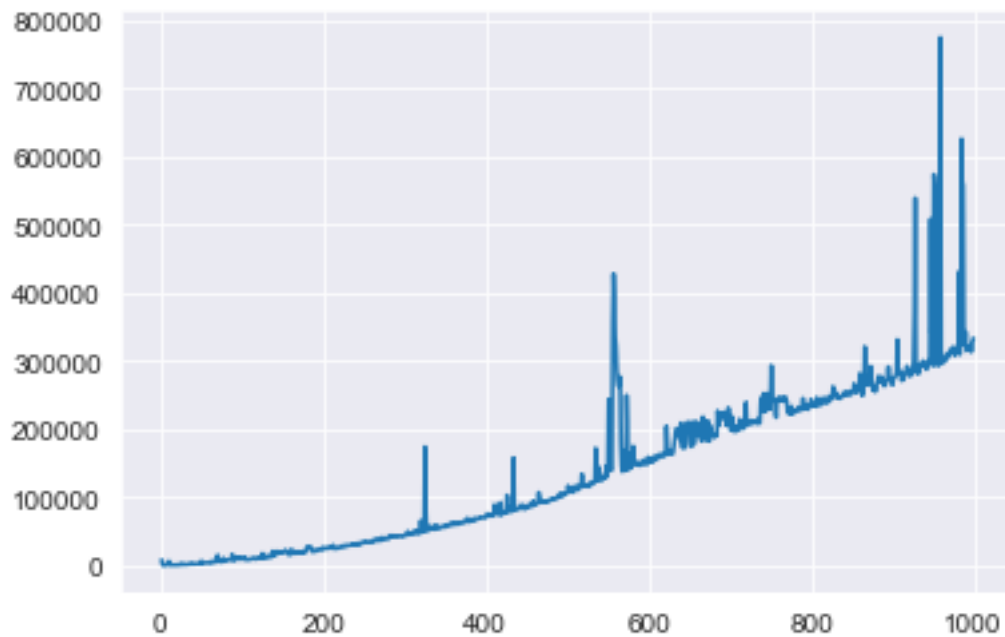


## 1.6 Convergence Time

As mentioned before, another issue with increasing the dimension of the feature space is the training time required to fit a machine learning model. While more data will generally lead to better predictive results, it will also substantially increase training time. To demonstrate this, generate lists of random numbers as you did above. Then, use this list of random numbers as a feature in a mock dataset; choose an arbitrary coefficient and multiply the feature vector by this coefficient. Then, sum these feature-coefficient products to get an output,  $y$ . To spice things up (and not have a completely deterministic relationship), add a normally distributed white noise parameter to your output values. Fit an ordinary least squares model to your generated mock data. Repeat this for a varying number of features, and record the time required to fit the model. (Be sure to only record the time to train the model, not the time to generate the data.) Finally, plot the number of features,  $n$ , versus the training time for the subsequent model.

```
[23]: import pandas as pd
import datetime
from sklearn.linear_model import LinearRegression, Lasso
```

```
[38]: # Your code may take some time to run
# LR = LinearRegression()
sample_size = 10**3
times = []
for i in range(1, 1001):
    x = np.array([np.random.uniform(low = -10, high = 10, size = i) for j in
    range(sample_size)])
    coeff = np.array(range(1, i+1))
    y = x.dot(coeff) + np.random.normal(loc=0, scale=.1, size=sample_size)
    LR = LinearRegression()
    start = datetime.datetime.now()
    LR.fit(x, y)
    end = datetime.datetime.now()
    elapsed = end - start
    times.append(elapsed)
plt.plot(range(1,1001), [t.microseconds for t in times]);
```



- Repeat the same experiment for a Lasso penalized regression model

```
[40]: # Your code may take some time to run

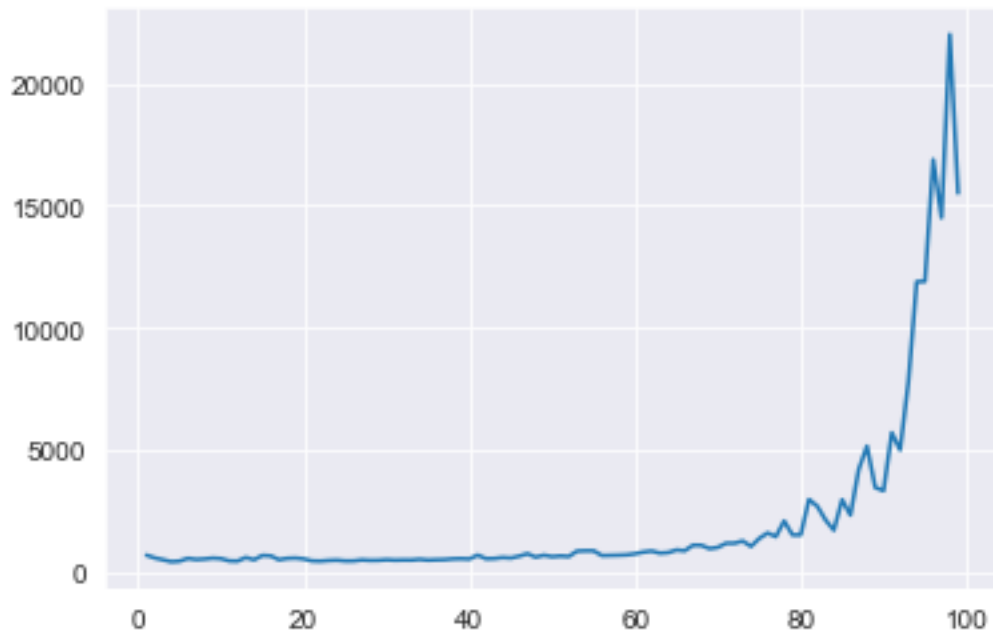
# Your code may take some time to run
```

```

# LR = LinearRegression()
sample_size = 10**2
times = []
r = 100
for i in range(1, r):
    x = np.array([np.random.uniform(low = -10, high = 10, size = i) for j in
    ↪range(sample_size)])
    coeff = np.array(range(1, i+1))
    y = x.dot(coeff) + np.random.normal(loc=0, scale=.1, size=sample_size)
    LS = Lasso()
    start = datetime.datetime.now()
    LS.fit(x, y)
    end = datetime.datetime.now()
    elapsed = end - start
    times.append(elapsed)
plt.plot(range(1,r), [t.microseconds for t in times]);

```

/opt/anaconda3/envs/learn-env/lib/python3.8/site-  
packages/sklearn/linear\_model/\_coordinate\_descent.py:647: ConvergenceWarning:  
Objective did not converge. You might want to increase the number of iterations,  
check the scale of the features or consider increasing regularisation. Duality  
gap: 1.027e+05, tolerance: 8.721e+04  
model = cd\_fast.enet\_coordinate\_descent(



## 1.7 Optional (Level Up)

### 1.7.1 Show Just How Slow it Can Go!

If you're up for putting your computer through the wringer and are very patient to allow the necessary computations, try increasing the maximum  $n$  from 1000 to 10,000 using Lasso regression. You should see an interesting pattern unveil. See if you can make any hypotheses as to why this might occur!

*Note: You can expect your code to take over an hour to run on a 2.7 GHz speed CPU!*

```
[ ]: # Your code may take some time to run
```

## 1.8 Summary

In this lab, you conducted various simulations to investigate the curse of dimensionality. This demonstrated some of the caveats of working with large datasets with an increasing number of features. With that, the next section will explore principal component analysis, a means of reducing the number of features in a dataset while preserving as much information as possible.