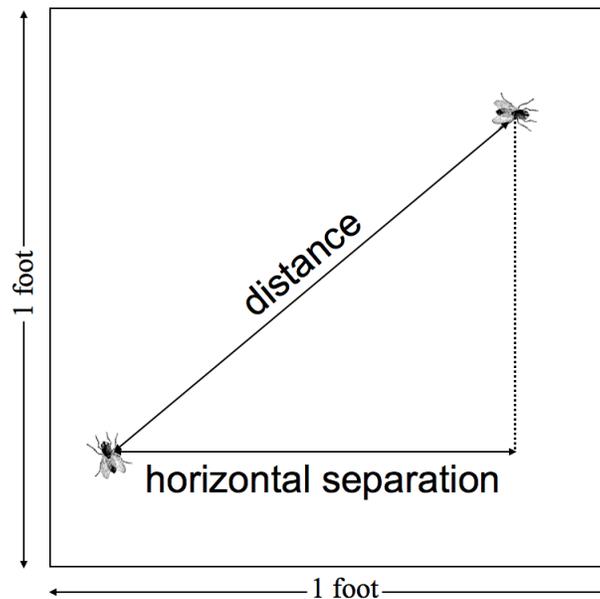# Python Demonstration
## MATH1900: Machine Learning

Location: http://people.sc.fsu.edu/~jburkardt/classes/ml_2019/python_demo/python_demo.pdf



*"Two flies land at random on a square piece of paper. On average, how far apart are they?"*

---

**Python Demo**

*How can we use a programming language like Python to:*

- *create a model of a mathematical problem;*
- *carry out a procedure to solve the problem;*
- *graphically illustrate the solution.*

---

Python is widely used for machine learning;

- it can be used interactively;
- a user can create new Python commands that extend the language;
- libraries are available for numerics, scientific computing, graphics;
- many machine learning packages have been built using Python.

The problem we will consider involves estimating the average distance of two flies who randomly settle on a square piece of paper. This particular problem can be solved analytically, but we are interested in how Python can simulate the problem, estimate the answer, and produce convincing tables and plots. We are studying a random process, and trying to say something (the average) about a huge number of cases.

Let's start by trying to simulate a single case.

# 1 One simulation

Let's assume the paper is one foot on each side. Then the position of any point on the square is simply a pair $(x, y)$ of numbers between 0 and 1. We need a way to:

- pick a pair of numbers `x` and `y` at random, and pack them into a point `p`:
- compute the distance between two points `p1` and `p2`;

The Python library **numpy** contains many useful functions, which are organized into collections called modules.

The module **numpy.random** includes the useful function `rand`:

- `rand()` returns a single random number;
- `rand( n )` returns n random numbers;
- `p=rand(2)` returns 2 random numbers as an object named "p";

To simulate two random flies:

```
1   import numpy as np
2   p1 = np.random.rand(2)
3   print ( p1 )
4   p2 = np.random.rand(2)
5   print ( p2 )
```
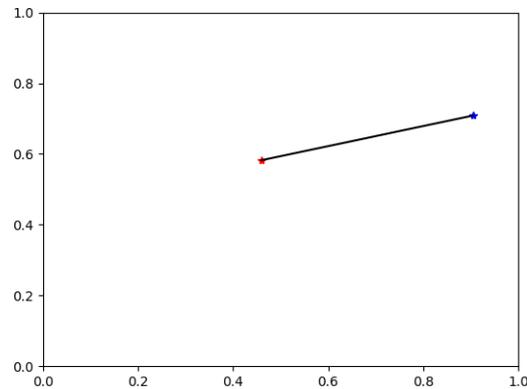
The module **numpy.linalg** includes the function `norm`:

- `norm(p)` returns the Euclidean length of the vector `p`;
- `norm(p1-p2)` returns the Euclidean distance between points `p1` and `p2`;

```
1   import numpy as np
2   p1 = np.random.rand(2)
3   p2 = np.random.rand(2)
4   d = np.linalg.norm ( p1 - p2 )
5   print ( d )
```

We can view the pair of flies using the **matplotlib** library:

```
1   import numpy as np
2   p1 = np.random.rand(2)
3   p2 = mp.random.rand(2)
4   import matplotlib.pyplot as plt
5   plt.plot ( p1[0], p1[1], 'r*' )
6   plt.plot ( p2[0], p2[1], 'b*' )
7   plt.plot ( [p1[0],p2[0]], [p1[1],p2[1]], 'k-' )
8   plt.axis ( [0,1,0,1] )
9   plt.show ( )
```

## 2 Many Simulations, Averaged

The average distance between random flies can be estimated by simulating many examples. In Python, the *for* loop lets us repeat an operation many times. We can initialize a variable `average` to zero, loop over `n` simulations, adding each distance to `average`, and then dividing by `n`:

```
1   import numpy as np
2   average = 0.0
3   n = 10
4   for test in range ( 0, n ):
5     p1 = np.random.rand(2)
6     p2 = np.random.rand(2)
7     d = np.linalg.norm ( p1 - p2 )
8     average = average + d
9
10  average = average / n
11  print ( average )
```

If we all run this program, using `n=10`, how much variation do we see in our answers? Should that make us worry? Let's try `n = 100` and compare answers. What about `n = 1000`?

**How large a multiple of 10 do we need before everyone's answer agrees to three decimal places?**

## 3 Using Arrays

Rather than working with one case at a time, we could instead generate all the cases at once. In that case, we will be working with arrays. For instance, `p1` will be a $2 \times n$ array, where each column represents a separate case.

The `norm()` function must be told to separately compute the norms of each column of the array. Columns count as "axis 0" in Python arrays.

```
1   n = 10
2   p1 = np.random.rand(2,n)
3   p2 = np.random.rand(2,n)
4   d = np.linalg.norm ( p1 - p2, axis = 0 )
5   average = np.mean ( d )
6   print ( average )
```

3

The advantages of writing a calculation in the vector version are that the code is somewhat more compact, it is easier to see that the computation could be done in parallel, and in particular, Python is able to compute many vector operations very efficiently.

# 4    Timing Comparison

Let's compare the scalar and vector codes when **n** is large.

First, let's rewrite our scalar calculation as a file, which we will call *simulate_many.py*, with the value of **n** an input parameter.

```
 1  import numpy as np
 2  def simulate_many ( n ):
 3     average = 0.0
 4     for test in range ( 0, n ):
 5        p1 = np.random.rand(2)
 6        p2 = np.random.rand(2)
 7        d = np.linalg.norm ( p1 - p2 )
 8        average = average + d
 9
10     average = average / n
11     return average
```

Doing this means we have created our own Python command, which can be executed with a single line:

```
 1     average = simulate_many ( n )
```

We can similarly create *simulate_vector.py*. Now we could compare the speed of the two methods by simply waiting for each to complete, or using a stopwatch, but instead, we will use the Python command **time**.

In order to increase **n** from 10 to 100 to ... to 1,000,000, we use the exponentiation operator **\*\***.

In order to print three things and have them all show up on one line, we can add to the **print()** statement the argument **end=""**.

```
 1     from time import time
 2     for logn in range ( 1 : 7 ):
 3        n = 10 ** logn
 4        print ( n, end = "" )
 5        t = time ( )
 6        simulate_many ( n )
 7        t = time - t
 8        print ( t, end = "" )
 9        t = time ( )
10        simulate_vector ( n )
11        t = time ( ) - t
12        print ( t )
```

We know that **simulate_many()** and **simulate_vector()** are computing the same things. When working with large sets of data, however, it can be a big time advantage to process all the data together, rather than doing them one at a time!

# 5    A convergence table

The average is an average over an infinite number of cases. We are estimating that average using **n** cases. How do we convince ourselves that the estimate is a good approximation? One way is to repeat the calculation for larger values of **n** and make a table. If the estimates seem to converge, then this suggests that our many choices of **n** seem to point towards a common answer *although we can't be sure this answer is right!*.

```
1    average_old = 0.0
2    for logn in range ( 1 : 7 ):
3      n = 10 ** logn
4      average = simulate_vector ( n )
5      print ( n, ' ', average, ' ', average − average_old )
6      average_old = average
```

Our results suggest that even a value of `n=1,000,000` only gives us confidence in the first two decimals of our estimate.

# 6 A histogram of the results

In the problem we are thinking about, possible distances $d$ range between 0 and $\sqrt{(2)}$. Some distances are more likely than others, as described by a probability density function. If we choose `n` large enough, we can make a histogram of our results that will suggest what the true PDF looks like.
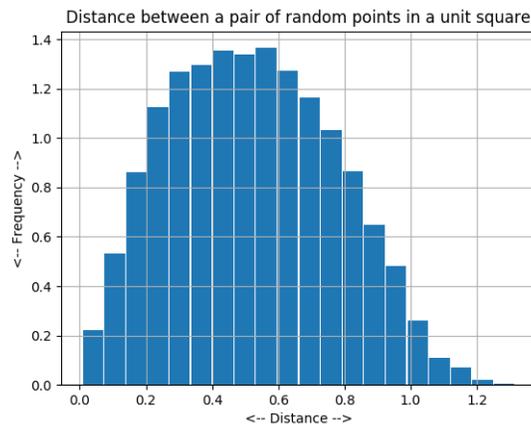The basic command is simply

```
1    plt.hist ( d )
```

but I want to specify the number of bins to be 20, I want the bars to be slightly separated, and I want the vertical $y$ axis to be scaled like a PDF, rather than a frequency count.

```
1    import matplotlib.pyplot as plt
2    import numpy as np
3    n = 10000
4    p1 = np.random.rand ( 2, n )
5    p2 = np.random.rand ( 2, n )
6    d = np.linalg.norm ( p1 − p2, axis = 0 )
7
8    plt.hist ( d, bins = 20, rwidth = 0.95, density = True )
9    plt.show ( )
```



Distance between a pair of random points in a unit square

# 7 Histogram versus PDF

For this particular problem, there happens to be an exact formula for the PDF. We can plot it together with the histogram, to see how we did. The formula for the exact PDF is a little messy, so I will skip the

programming details:

$$\text{pdf}(d) = \begin{cases} 2d(d^2 - 4d + \pi) & \text{if } d <= 1.0 \\ 2d(4\sqrt{d^2 - 1.0} - (d^2 + 2 - \pi) - 4\arctan(\sqrt{d^2 - 1})) & \text{if } 1.0 < d \end{cases}$$

Compare exact and observed PDF