# Graph Analysis using Spark, GraphX and Pregel

A graph is a mathematical representation for a Social Network. The graph consists of Vertices (which represent entities or people) and Edges connecting vertices (which represents the relationship between nodes).

A graph is useful as it provides a high level of information about certain features of the graph like:

1. **Connectedness** – The number of edges from each node. This will represent the 'popularity' of a node.

2. **Average Path Length** – The 'distance' between any two nodes. This can be used to represent the level of 'friendship' between two people, whether they know each other through common friends.

3. **Number of triangles** – Any 3 nodes connected in a network is called a triangle. A network with more triangles is considered more stable than a network with less triangles and more open nodes.

In this project, we will be using the GraphX library in Spark which provides new datatypes like VertexRDD, EdgeRDD and Edge to replicate a graph. We will also be using the Pregel API which will help us traverse the graph and find out properties like degrees of separation.

We will be working with two datasets. The first dataset contains the Superhero names and their unique ID, which we will use as a Vertex ID. The second dataset we will be using is the Superhero Co-appearance dataset, which contains all the superheroes in the Marvel Universe. Each row of the dataset contains the Unique ID of a Superhero followed by the IDs of the Superheroes he/she has made a co-appearance with in a comic.

**Creating the Superhero Vertices:**

The first step would be to parse the Superhero names dataset and get their names along with their IDs. We will be using the function parseNames to perform this task. This function will work in the following way:

1. For each row in the dataset, split based on the '\' delimiter.

2. Check if the row has length greater than 1. If yes, then trim the first field and convert it to Long and store it a heroID.

3. Check if the heroID is less than 6487 (The ones greater than that are arbitrary numbers). If yes, then return the first field (containing Hero ID) trimmed and converted to long and the second field (containing the Hero name).

This function will be called during the formation of the Vertices of our network. The HeroID will work as the VertexID.

**Creating the Superhero Network Edges:**

Now we will create the Superhero network. For this we will create a function makeEdges. The function will input each row of the Superhero Co-appearance file and return an Edge datatype list using the number of connections of the first ID in a row with the following IDs.

1. The function will first create a mutable List Buffer 'edges'.

2. It will split the input line based on the delimiter ' '.

3. The first ID in the line will be stored in the variable 'origin'.

4. For all the other IDs, an Edge variable will be created with the origin and current ID, both converted to Long. This Edge variable will be added to the edges list.

5. The edges list will be returned.

The Edge datatype also has an 'attribute' variable which can help keep track of physical distances and other parameters. But we have not used it.
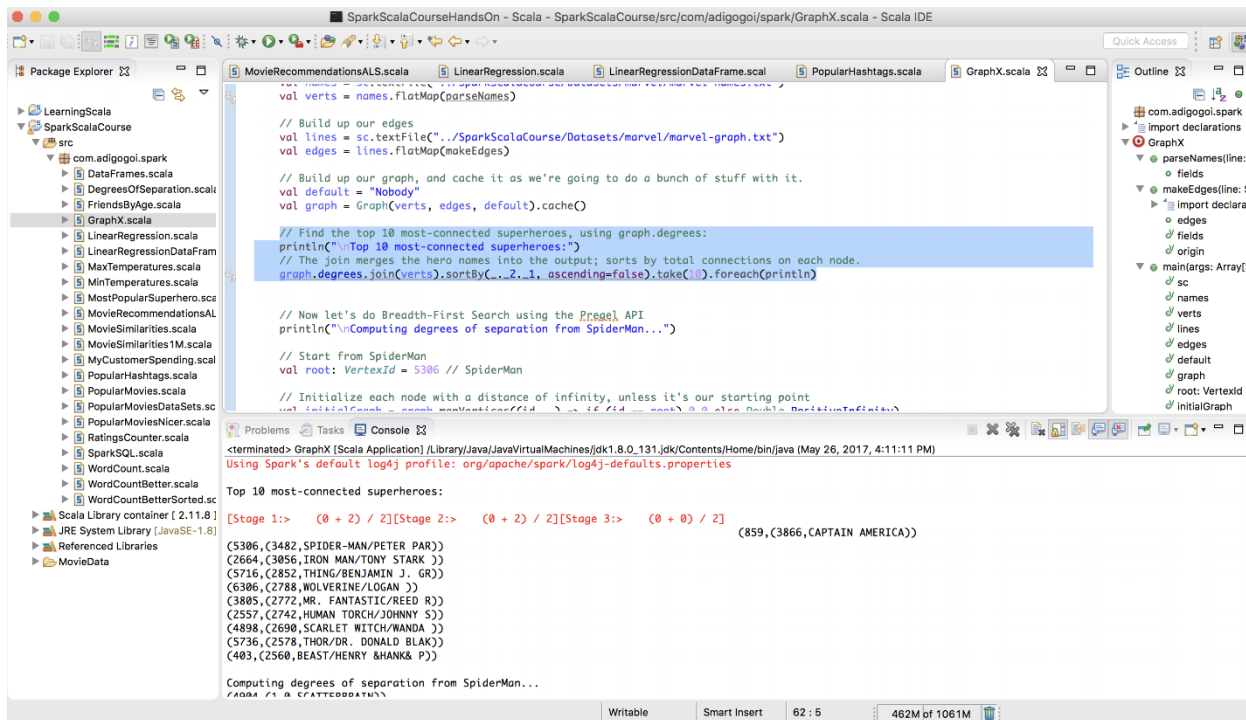
**Performing the Network Analytics:**

Our main function will first read the names and network files and store them in datasets in the right format. Then we will call the functions mentioned above to create the network.

The vertices and edges dataset will be merged into the 'graph' dataset, which will have the final network. We will cache it as we may need to use it more than once.

We will use Pregel as our network traversing API for this network. Pregel provides the general mechanism to perform general algorithms on our graph like Breadth First Algorithm that we will be using for this project. The explanation of this algorithm is out of scope for this project but in general it involves keeping a node as root node the sending of messages between nodes and recording the distance between them with respect to the root node.

Some of the analytics that we have performed on the network are:

1. **Getting the 'most-connected' Superheroes** – The nodes with the most edges connected to them represent the entities that have more connections, which can be friends or co-workers. In this operation, we are finding out the most famous hero based on the number of co-appearances they have made.



We can see that Captain America is the most famous hero with 3866 co-appearances with other heroes, followed by Spider-Man, with 3482 co-appearances. This is interesting, as I always thought that Spider-Man would be the most famous superhero of all times. But this is just based on the co-appearances made by each hero, so more data might be required to make a conclusion.

2. **Finding connections between Spider-Man and other heroes** – Having my faith that Spider-Man would be the most famous hero of all-times, I decided to find out the degrees of separation between him and the rest of the Marvel Superheroes. This I did using the afore-mentioned Breadth First Search algorithm of the Pregel API.

I noticed that even the lesser-known heroes have just 2 to 3 degrees of separation from Spider-Man. This maybe because many of the heroes Spider-Man has co-appeared with have also co-appeared in many comics with these heroes.

3. **Finding distance between Spider-Man and Adam** – I tried to focus on just one of the heroes who is not so popular in the Marvel Universe, Adam, and decided to find out his distance from Spider-Man. Adam's ID is 14 so that made it easy to find their Degrees of Separation.



Even Adam is just 2 degrees away from Spider-Man. These superheroes are a tight-knit bunch.

But this project did allow me to implement a network based graph data using Spark's GraphX and Pregel, so it is a success after all.