

Correlation

November 30, 2017

```
In [1]: from pyspark.sql.types import StructType, StructField, FloatType, LongType, StringType
        from pyspark.shell import spark

        feats = []
        f = open('features.txt')
        for line_num, line in enumerate(f):
            if line_num == 0:
                # Timestamp
                feats.append(StructField(line.strip(), LongType(), True))
            elif line_num == 1:
                # Geohash
                feats.append(StructField(line.strip(), StringType(), True))
            else:
                # Other features
                feats.append(StructField(line.strip(), FloatType(), True))

        schema = StructType(feats)
```

Welcome to

```
  _--_
 /  _/  _--_  _--_  _--_  _--_
\_  \/_  \/_  \/_  \/_  \/_  \/_
/_  /  .--/\_  \/_  \/_  \/_  \/_  version 2.2.0
/_  /
/_  /
```

Using Python version 3.6.3 (default, Oct 6 2017 12:04:38)
SparkSession available as 'spark'.

```
In [3]: df = spark.read.format('csv').option('sep', '\t').schema(schema).load('inputs/nam_mini

In [17]: col_names = []
        for i in range(2, len(df.columns)):
            col_names.append(df.columns[i])
        df_features = df.select(*col_names)
        rdd_df = df_features.rdd

In [18]: from pyspark.mllib.stat import Statistics
        coeff = Statistics.corr(rdd_df.map(list), method="pearson")
```

```
In [21]: import numpy as np
         np.savetxt('./heatmap-generation/correlation_matrix.txt', coeff)
```