# *New Services in the Cortana Analytics Suite*

## *Apache Spark in HDInsight*

Asad Khan
Principal Program Manager
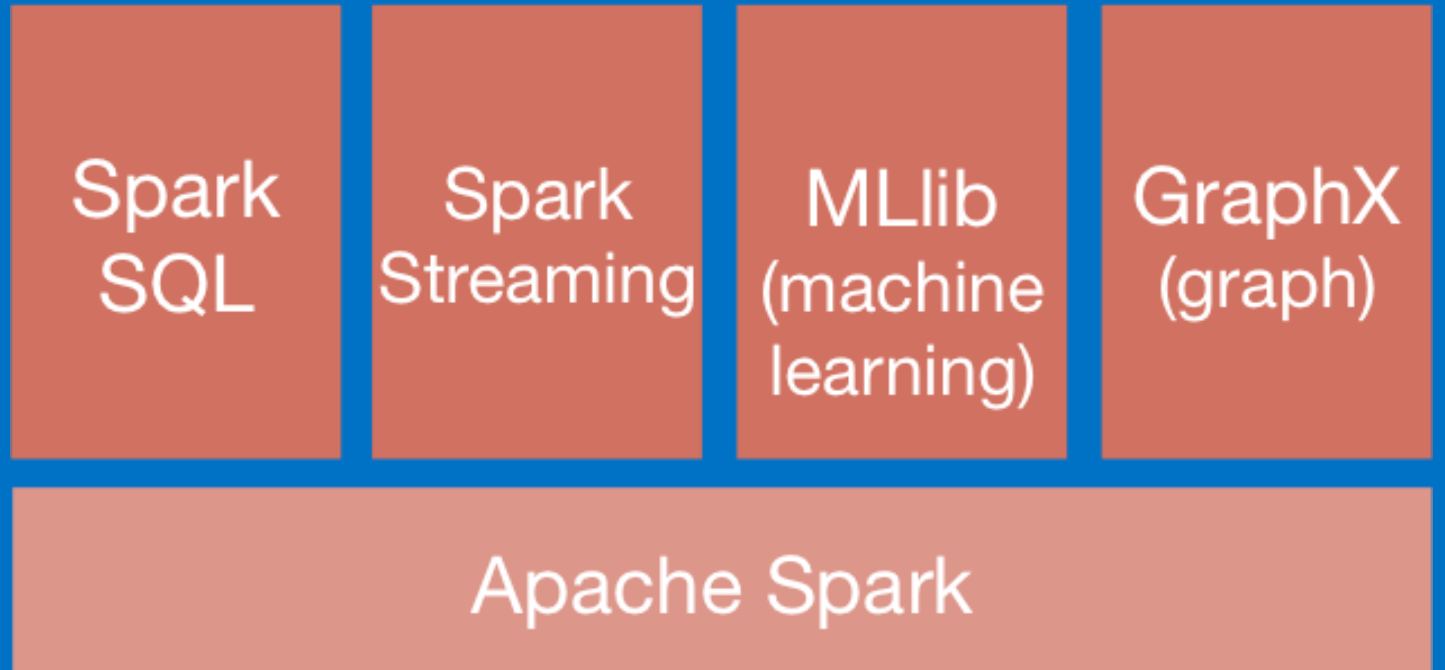Big Data Group

# Agenda

Introduce Apache Spark

Spark in Action

# Apache Spark – An Unified Framework

An unified, open source, parallel, data processing framework for Big Data Analytics

Spark Unifies:
- ☆ Batch Processing
- ☆ Real-time processing
- ☆ Stream Analytics
- ☆ Machine Learning
- ☆ Interactive SQL

| Spark SQL | Spark Streaming | MLlib (machine learning) | GraphX (graph) |
| --- | --- | --- | --- |

Apache Spark

https://spark.apache.org

# Spark – Benefits

## Performance

Using in-memory computing, Spark is considerably faster than Hadoop (100x in some tests).
Can be used for batch and real-time data processing.

## Developer Productivity

Easy-to-use APIs for processing large datasets.
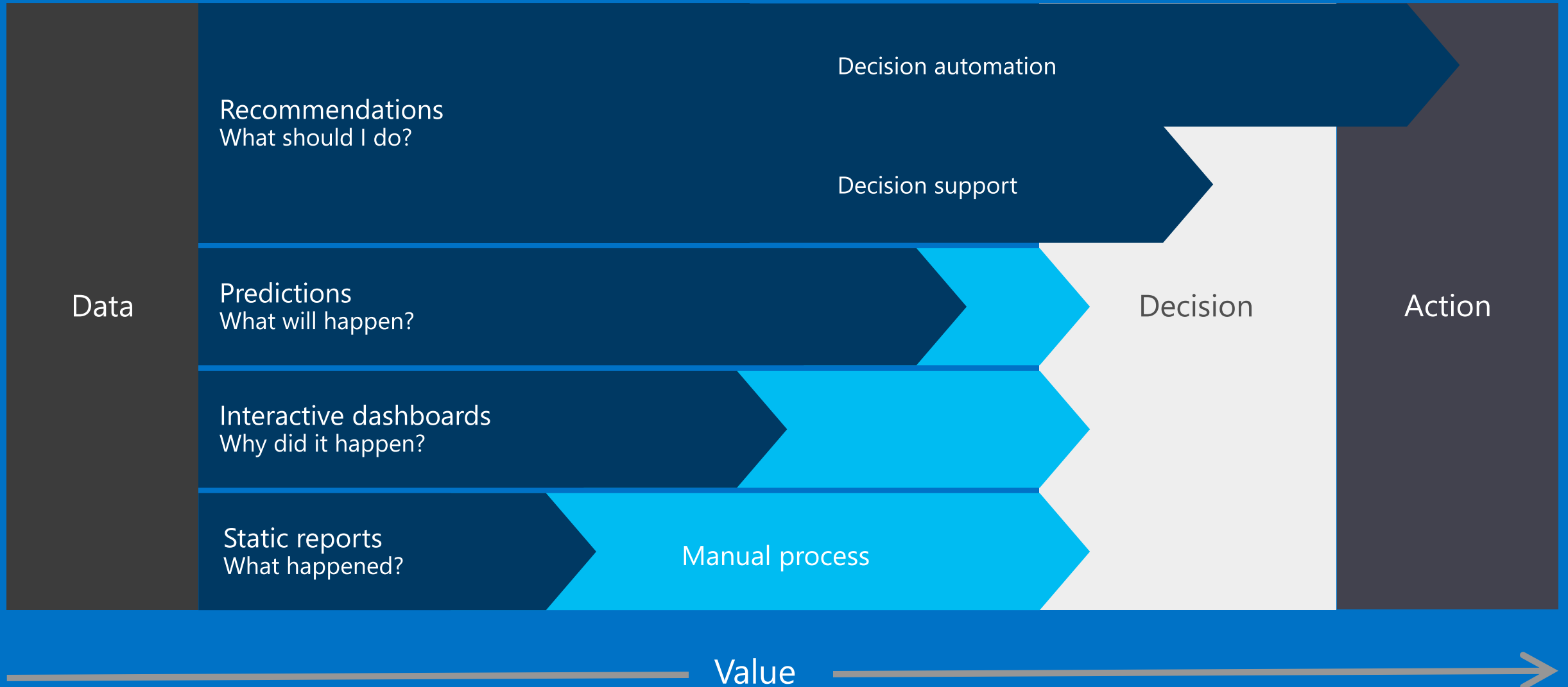Includes 100+ operators for transforming.

## Unified Engine

Integrated framework includes higher-level libraries for interactive SQL queries, processing streaming data, machine learning and graph processing.
A single application can combine all types of processing

## Ecosystem

Spark has built-in support for many data sources such as HDFS, RDBMS, S3, Apache Hive, Cassandra and MongoDB.

Runs on top the Apache YARN resource manager.

From data to decisions and actions

# Spark is fast

Spark is the current (2014) Sort Benchmark winner.
3x faster than 2013 winner (Hadoop).

|  | 2013 Record (Hadoop) | Spark 100 TB | Spark 1 PB |
|---|---|---|---|
| Data Size | 102.5 TB | 100 TB | 1000 TB |
| Time | 72 min | 23 min | 234 min |
| Nodes | 2100 | 206 | 190 |
| Cores | 50400 | 6592 | 6080 |
| Rate/Node | 0.67 GB/min | 20.7 GB/min | 22.5 GB/min |

Spark is fast not just for In-Memory but On-Disk computation as well

tinyurl.com/spark-sort

# Interactive Data analysis through Zeppelin
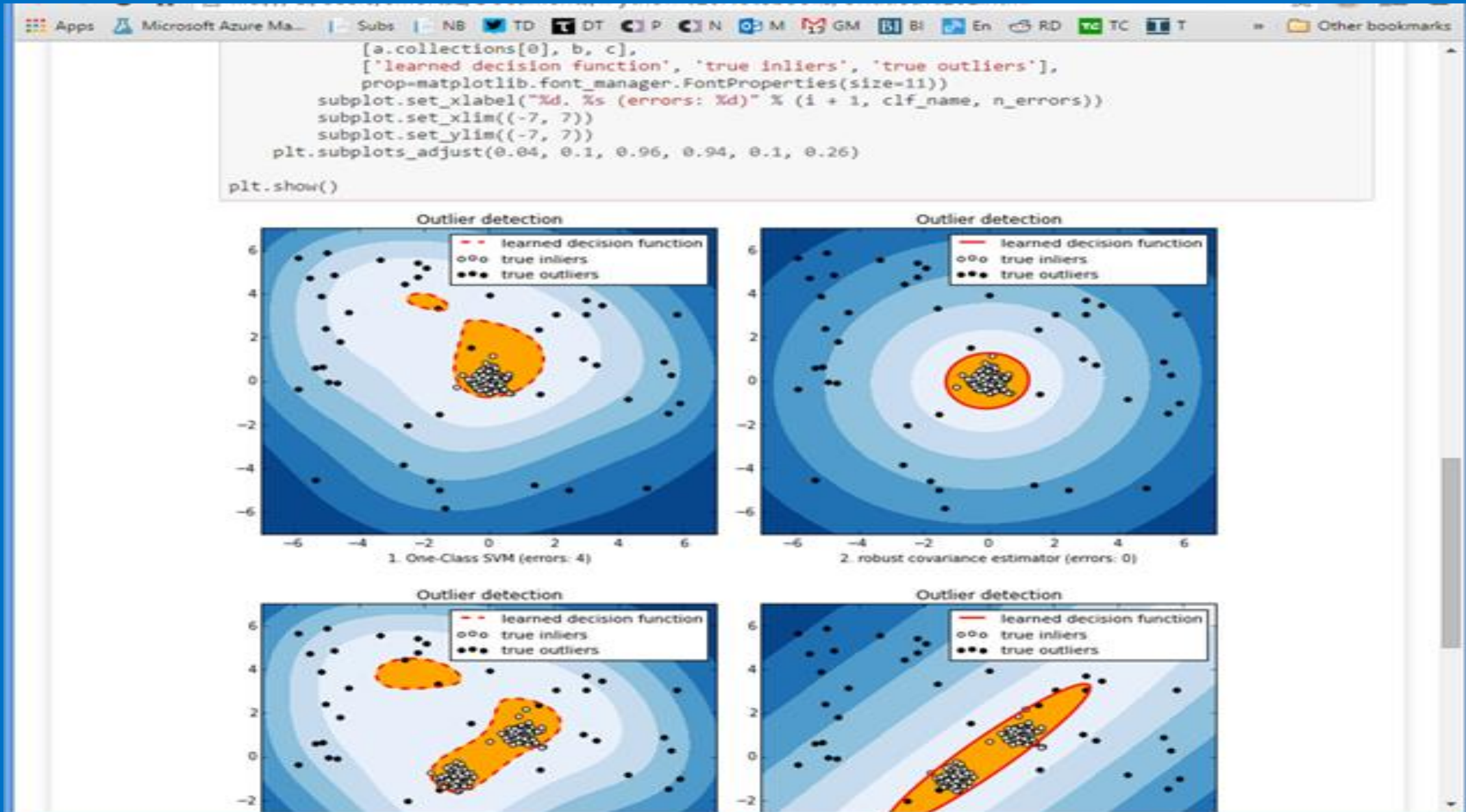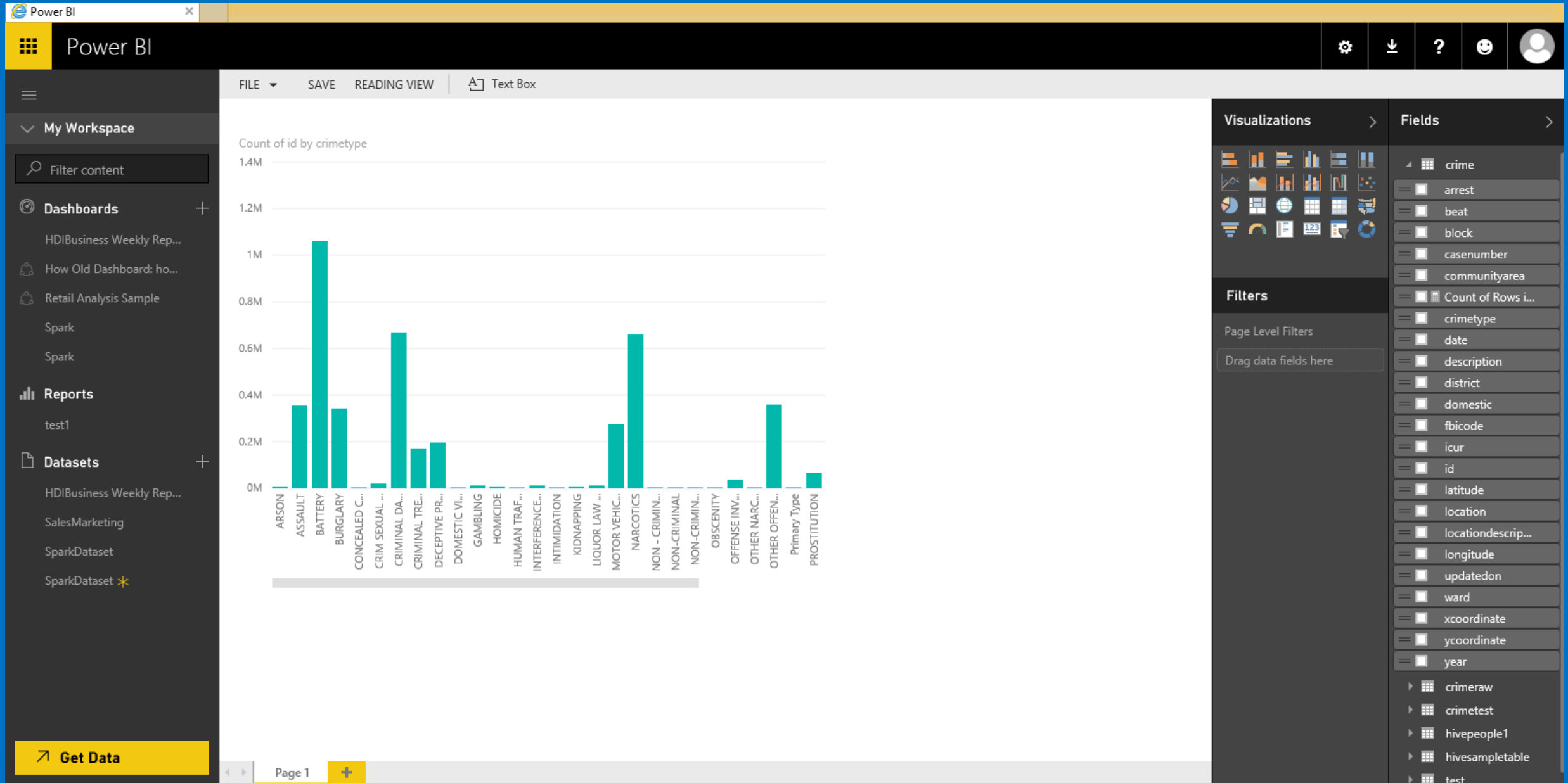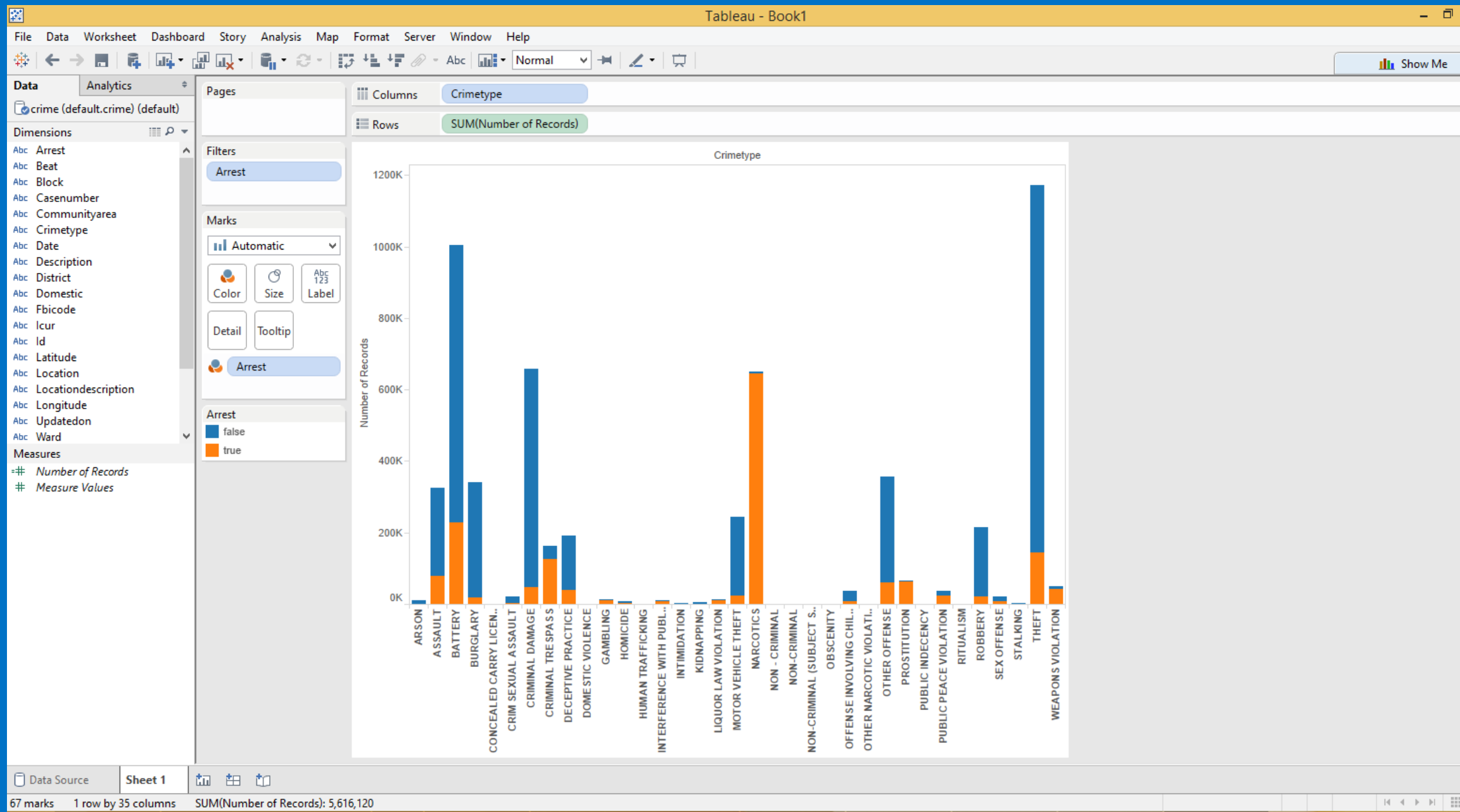
# Machine learning though Juypter

# Visual exploration through BI

# Visual exploration through BI

# Agenda

Introduce Apache Spark

Spark in Action

**Contact:**

**asadk@microsoft.com**