Spark
ooo

The Data
ooo

Analysis
ooooooooo

# Master the Mainframe World Championship 2016
## Apache Spark on IBM z/OS

Stephen Solis-Reyes

The University of Western Ontario
London, Ontario, Canada

September 17, 2016

Spark
○○○

The Data
○○○

Analysis
○○○○○○○○○

## Outline

1. About Spark (what I learned, what I liked)

2. About the data

3. My analysis/insights

Spark
●○○

The Data
○○○

Analysis
○○○○○○○○○

## About Spark



Spark is a platform for large-scale data analytics.

Spark
●○○

The Data
○○○

Analysis
○○○○○○○○○

## About Spark



Spark is a platform for large-scale data analytics.

- Multiple languages: Scala, Python, Java
- Large variety of input data sources
- Several libraries:
    - MLlib (machine learning)
    - GraphX (graph processing)
    - Spark Streaming (near-real-time streaming analytics)

## What I liked about Spark

- Extremely scalable, automatically
  (10000s of nodes, PBs of data)

## What I liked about Spark

- Extremely scalable, automatically
  (10000s of nodes, PBs of data)
- Very fault-tolerant
  - on failure, work is moved to another worker automatically – no data loss

Spark
○●○

The Data
○○○

Analysis
○○○○○○○○○

## What I liked about Spark

- Extremely scalable, automatically
  (10000s of nodes, PBs of data)
- Very fault-tolerant
  - on failure, work is moved to another worker automatically – no data loss
- Cross-platform, runs anywhere Java runs *without code changes*

Spark
○●○

The Data
○○○

Analysis
○○○○○○○○○

## What I liked about Spark

- Extremely scalable, automatically
  (10000s of nodes, PBs of data)
- Very fault-tolerant
    - on failure, work is moved to another worker automatically – no data loss
- Cross-platform, runs anywhere Java runs *without code changes*
- Interactive shell
    - can easily explore data 'live', no need to compile and wait

Spark
○○●

The Data
○○○

Analysis
○○○○○○○○○

## Spark on z/OS



- Run analytics right where the data is
  - less network use, lower latency, security, etc.

Spark
○○●

The Data
○○○

Analysis
○○○○○○○○○

## Spark on z/OS



- Run analytics right where the data is
  - less network use, lower latency, security, etc.
- Can use many z/OS-native data sources
  - DB2, VSAM, IMS, PDSE, etc.

Spark
○○●

The Data
○○○

Analysis
○○○○○○○○○

## Spark on z/OS



- Run analytics right where the data is
  - less network use, lower latency, security, etc.
- Can use many z/OS-native data sources
  - DB2, VSAM, IMS, PDSE, etc.

### Contest system

z13 running z/OS, data stored in DB2

Spark
ooo

The Data
●oo

Analysis
ooooooooo

# Dataset



~6000 clients

~1.5M transactions
(year 2013, ~$75.2M total)

Spark
○○○

The Data
○●○

Analysis
○○○○○○○○○

# Clients



Gender (A/B)



Age



Annual income



Education level



Service discontinued

Spark
ooo

The Data
oo●

Analysis
ooooooooo

# Transactions



Client



Date and Time



Amount



Merchant name/category



Card type/brand

Spark
ooo

The Data
ooo

Analysis
●oooooooo

## About the Analysis

**Predictive analytics drives smart business decisions**

- What actionable business insights can we draw from the data, using Spark?

Spark
○○○
The Data
○○○
Analysis
○●○○○○○○○○

# Transaction Behaviour Over Time

**Transactions per Month**



### Business Insights

Transaction volume is fairly steady per month.

- Probably not much benefit in monthly card promotions

# Transaction Behaviour Over Time
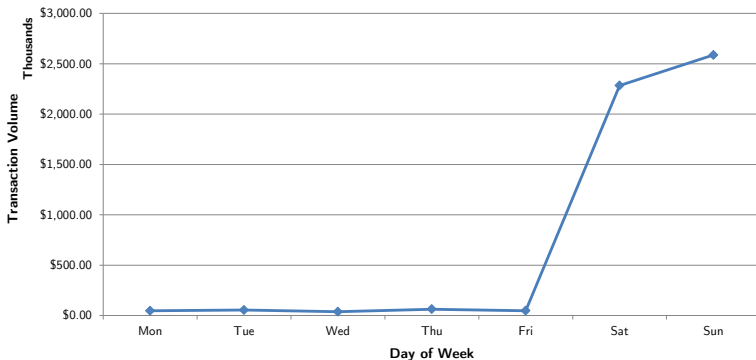
**Transactions per Day of Month**



### Business Insights

Huge spike on the 2nd of every month.

- Businesses should be prepared for the volume and use promotions to take advantage

# Merchant Category Insights
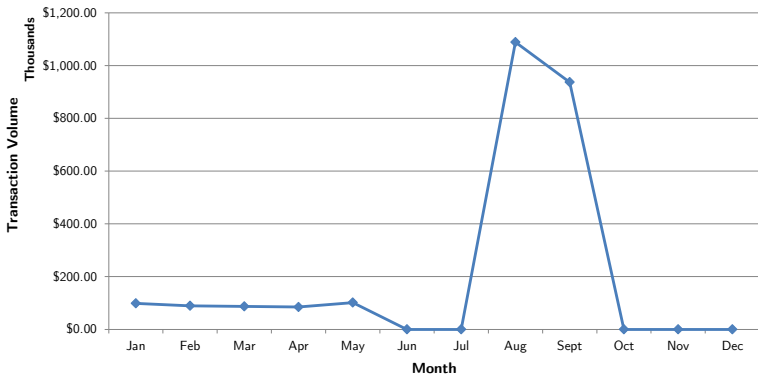


**Movie Theater Transactions per Day of Week**

### Business Insights

- Friday nights are surprisingly unpopular
- Weekday discounts (not just Tuesdays) could be effective

Spark
○○○

The Data
○○○

Analysis
○○○○●○○○○

# Merchant Category Insights

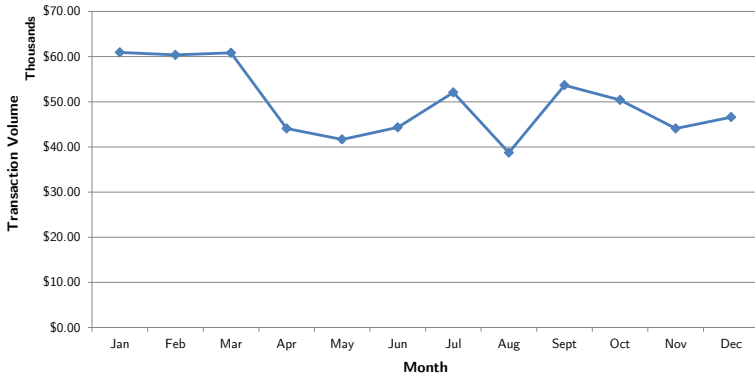**Travel (Air/Hotel) Transactions per Month**



### Business Insights

- Travel-related businesses should anticipate the spike in Aug/Sept
- Could offer discounts in Jul/July, Oct/Nov/Dec to encourage spending

Spark
○○○

The Data
○○○

Analysis
○○○○○●○○○

# Merchant Category Insights

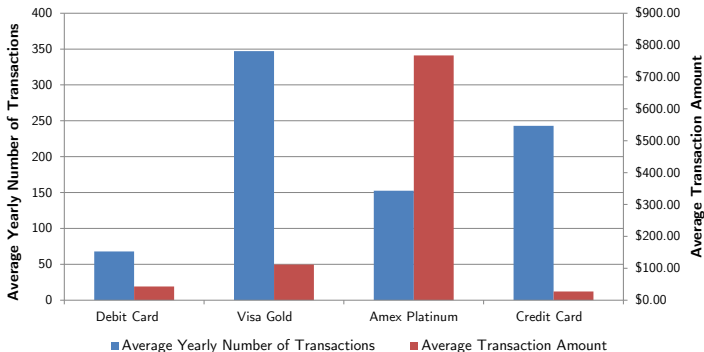**Heating and Plumbing Transactions per Month**



### Business Insights

Definite monthly trends, perhaps caused by weather/climate.

- Could consider discounts in off periods, hire seasonal employees during high demand

Spark
○○○

The Data
○○○

Analysis
○○○○○○●○○

# Card Brands
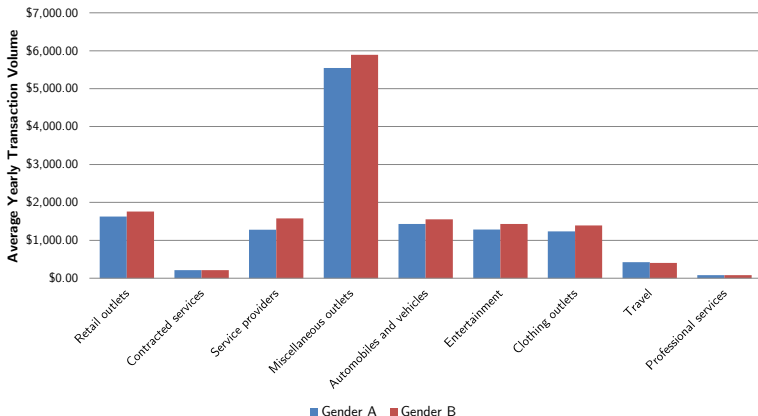
**Transactions per Card Brand**



**Business Insights**

- Banks should push credit cards, since people make more transactions with them
- Platinum card customers make much larger transactions, could push more expensive promos

Spark
○○○

The Data
○○○

Analysis
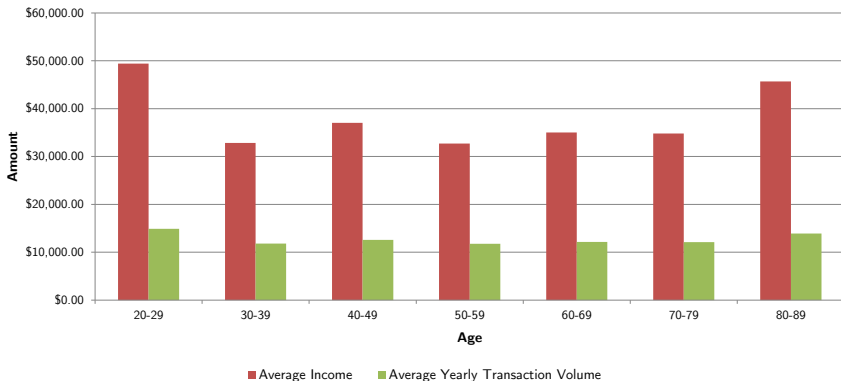○○○○○○○○●○

# Gender

## Merchant Categories per Gender



### Business Insights

Different genders do not seem to have different spending patterns.
- Gender-specific marketing may not be so effective

Spark
○○○

The Data
○○○

Analysis
○○○○○○○○○●

# Age



**Client Income and Transaction Volume per Age**

## Business Insights

Young and old clients have more income, but do not spend more money.

- Age-specific marketing may not be so effective

Spark
000

The Data
000

Analysis
000000000

"Information is the oil of the 21st century, and analytics is the combustion engine."

— Peter Sondergaard, Gartner Research

**Thanks! Questions?**



Me: Stephen Solis-Reyes