Uncompressed (virtual) Data Replicated by Hour of Day

Mark Chamness

Data Scientist

EMC  Corporation

Summer 2014

EMC²

# EMC Corporation

- Fortune "100" Corporation

- 61,000 employees in 86 countries

- Products:
  - Hardware: Data Storage systems (Petabyte scale)
  - Software/Security "data protection"

- $24 Billion revenue in 2013
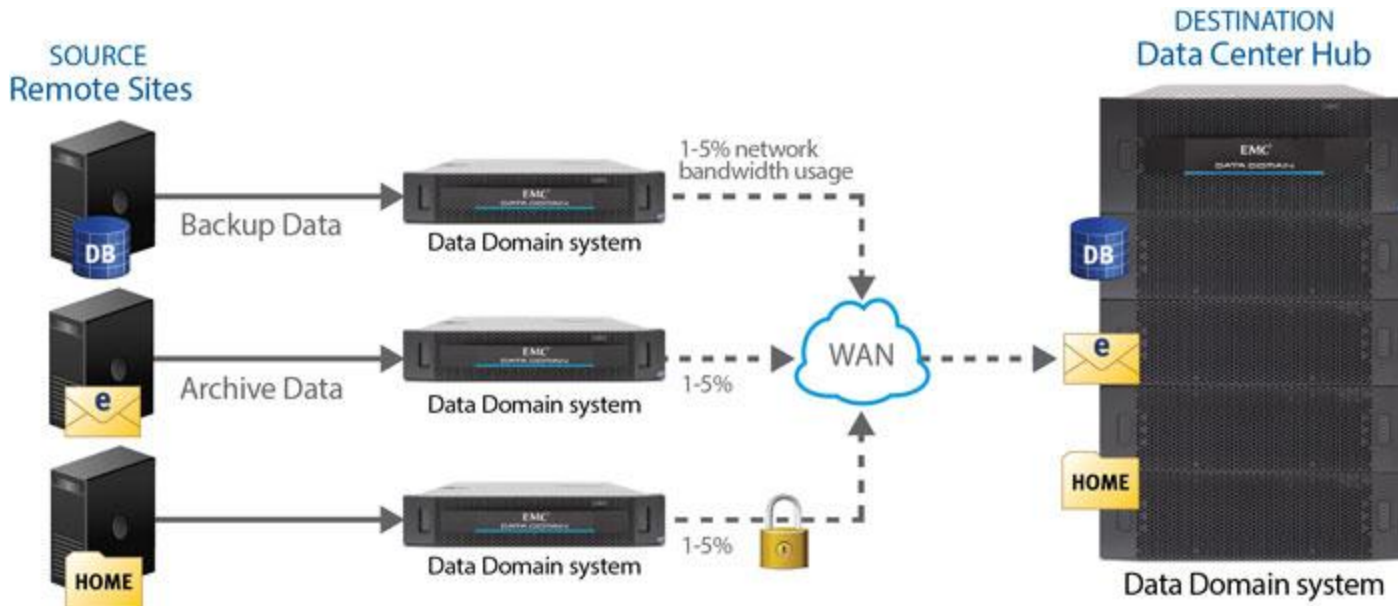
EMC²

# My Background

Academic:

- B.S. Physics, California Institute of Technology

- Ph.D. Candidate, Physics, Brown University

- M.S. Statistics, San Jose State.
  - (Expected 2015)

EMC²

# Finding the Position

- A previous manager contacted me

- After ~1 year, you create your own role

- Be proactive, not reactive in career goals

EMC²

# Why Study Replication?

- Most common customer query: "Replication"
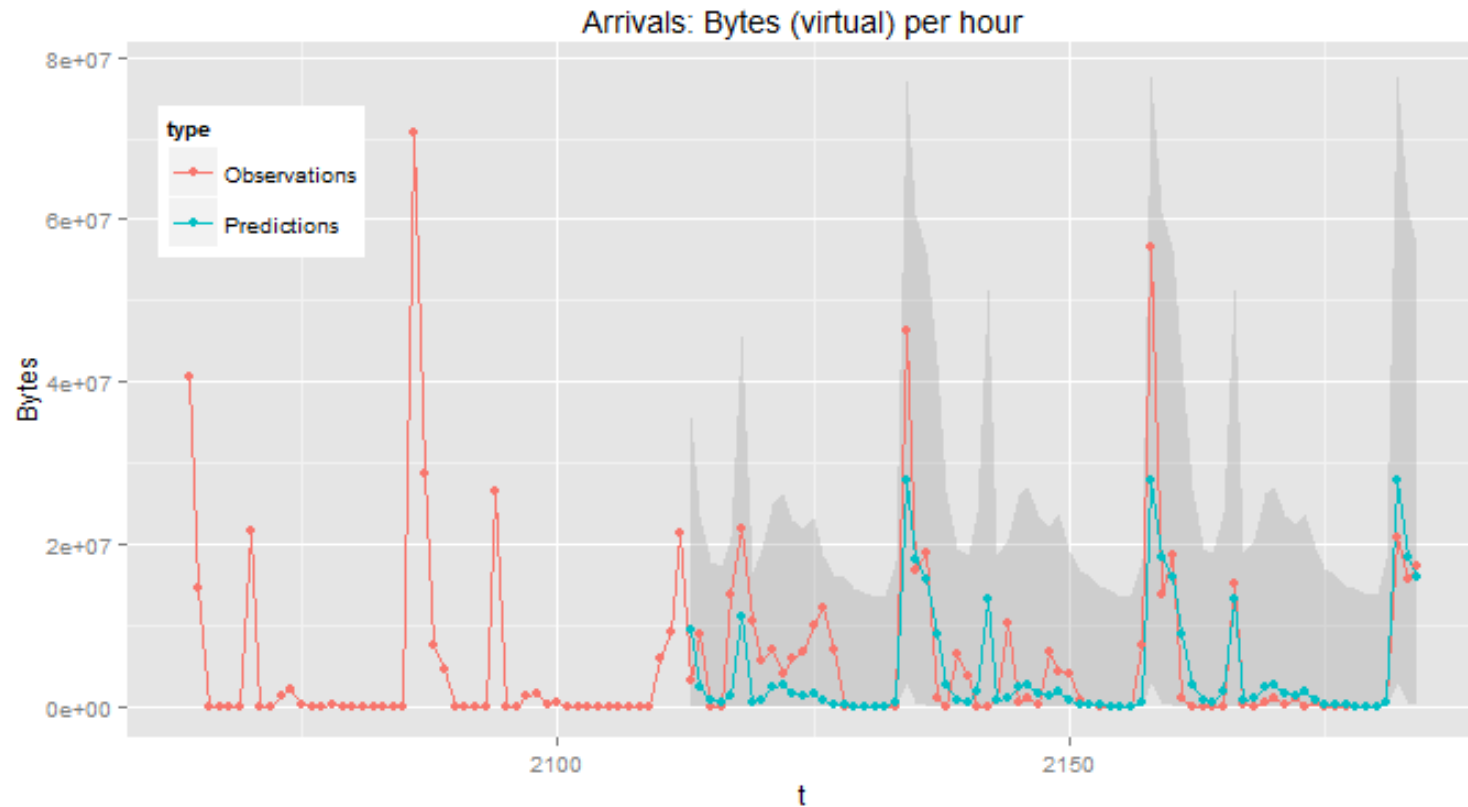
- Long time to resolve

# Research Goals

- Define, validate, & document the data set
  - Mountain of undocumented data

- Algorithm to indicate the probability and/or severity of replication lag
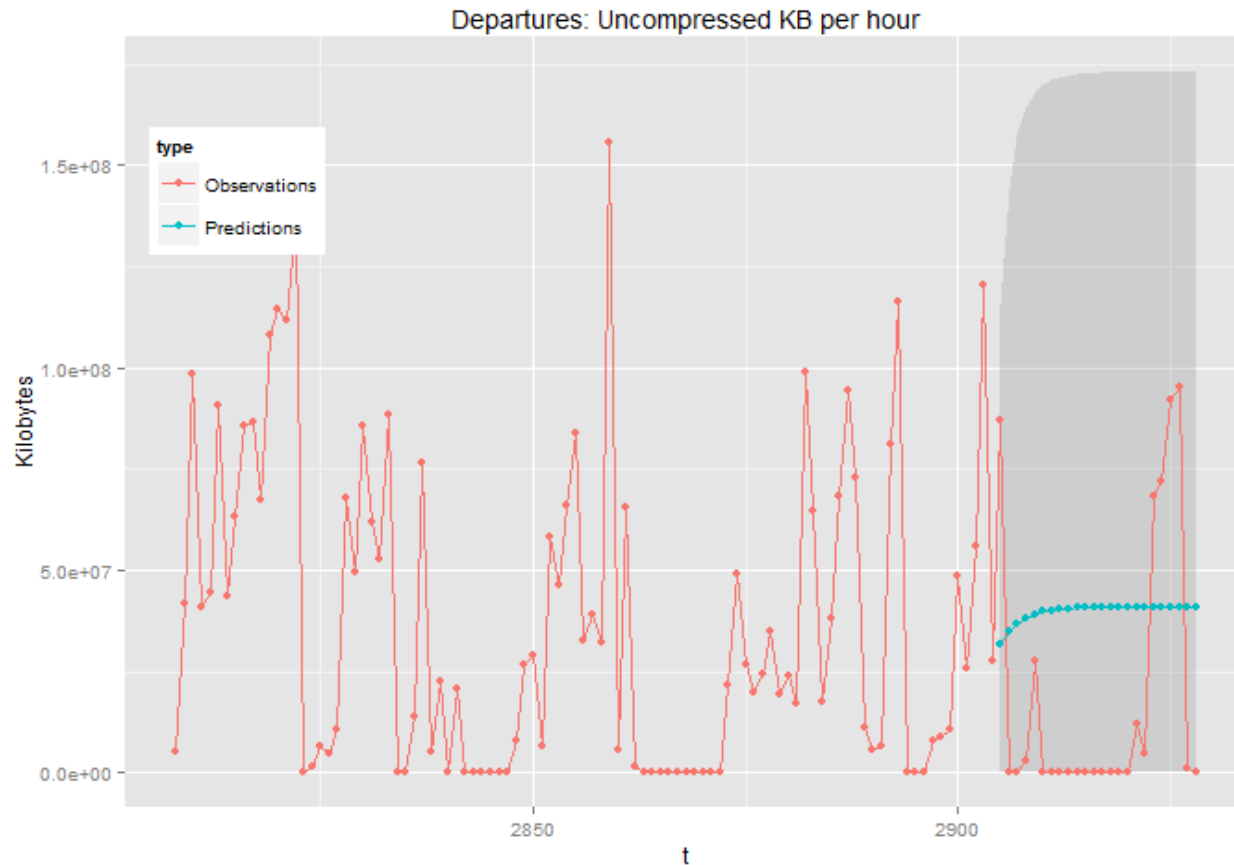
EMC²

# Challenges: Big Bad Data

- Data origin: computers all over world

- Hundreds of variables

- Billions of observations

- Not controlled experiment

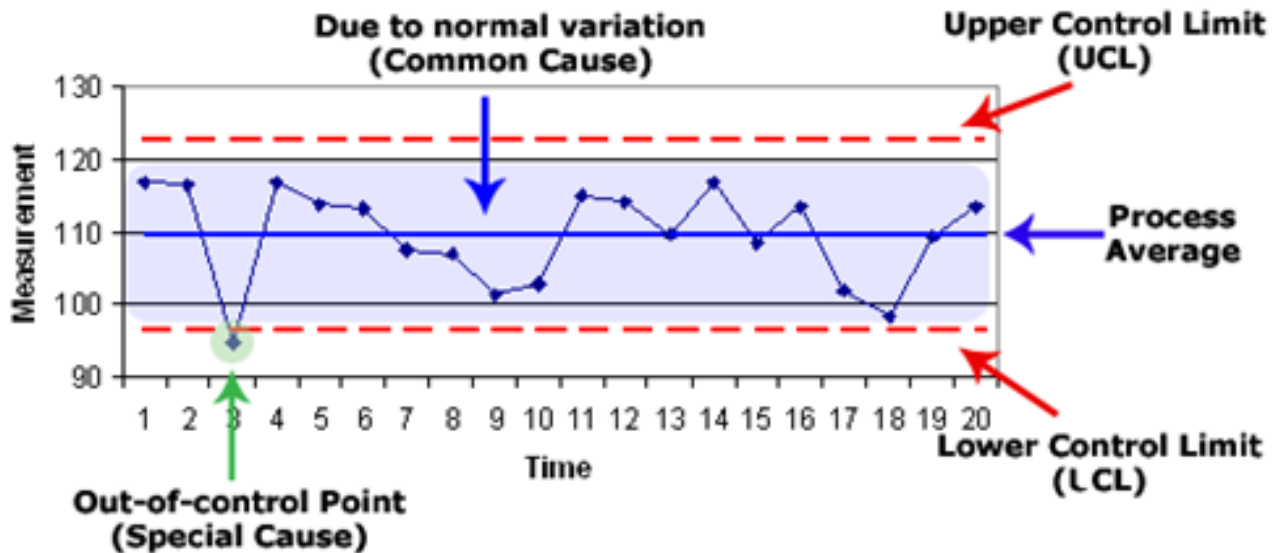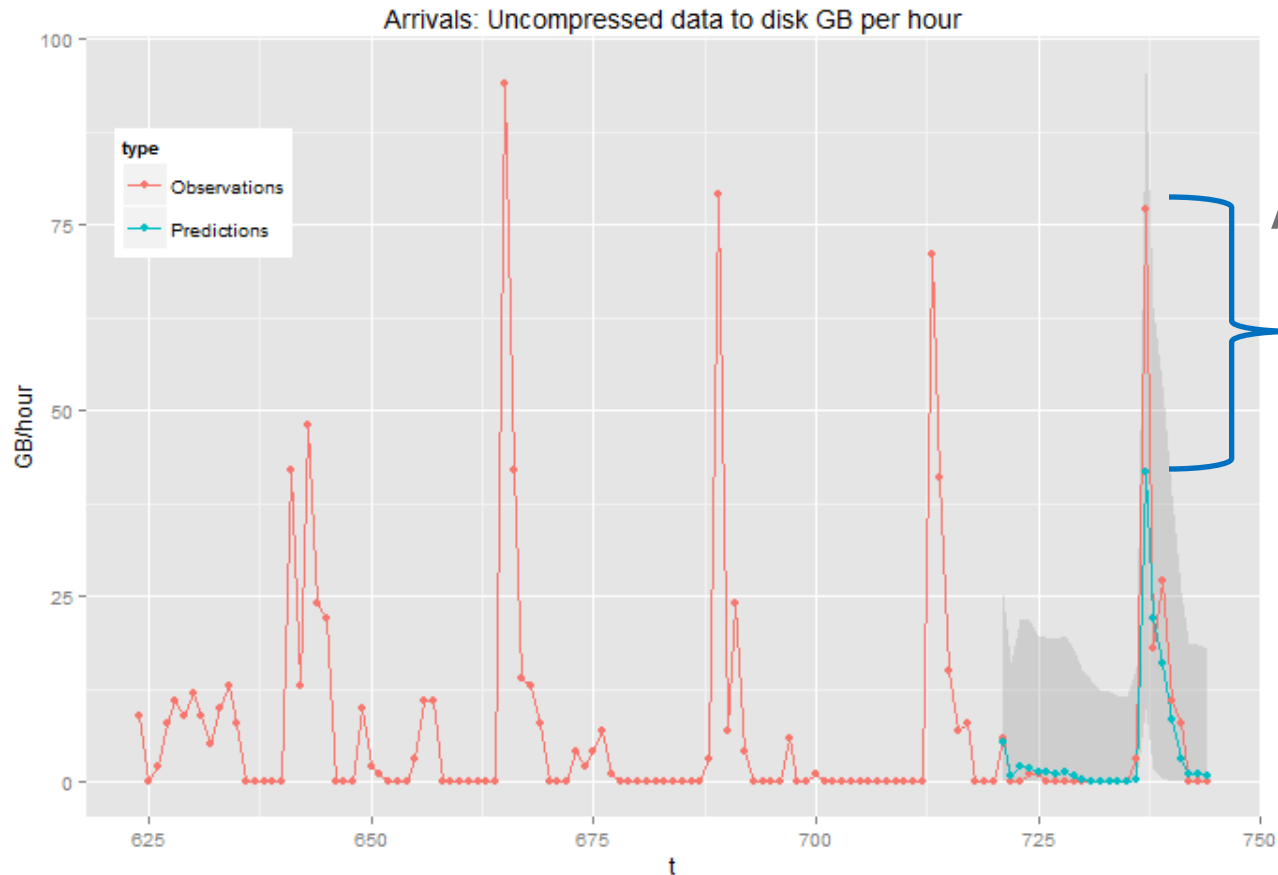- Significant Data Scrubbing

EMC²

# Time series analysis



Arrivals: Bytes (virtual) per hour

# Some forecasts less useful



Departures: Uncompressed KB per hour

# SPC - Statistical Process Control

- Define rules for "Out-of-Control" Conditions
- Test for rule match



**Due to normal variation (Common Cause)**

**Upper Control Limit (UCL)**

**Process Average**

**Out-of-control Point (Special Cause)**

**Lower Control Limit (LCL)**

Measurement — Time

EMC²

# p-values are useful



Arrivals: Uncompressed data to disk GB per hour

Almost exceeds 95% confidence interval

~2σ

EMC²

# Replication has hourly correlation



Uncompressed (virtual) Data Replicated by Hour of Day

Visualization Helps!

EMC²

# Software/Statistical Methods used

- R/RStudio
  - ggplot
  - Time series analysis

- Database
  - SQL (Postgresql database)
  - pgAdmin
  - Madlib: "Big Data SQL for Data Scientists"

- Documentation via Excel/PowerPoint/Word

EMC²

# Recommendations for Finding Internship

- LinkedIn
  - Educate yourself: find profiles in target role
  - Find relevancy of your background/experience
  - Keep up-to-date & accurate

- Salary.com
  - Gather data on your worth

- Job sites: Dice/BrassRing/etc

- Apply to companies you like

EMC²

# Q&A

EMC²