# Machine Learning Demo / Tutorial

Big Data Management
Big Data Analytics

Ken Cottrell
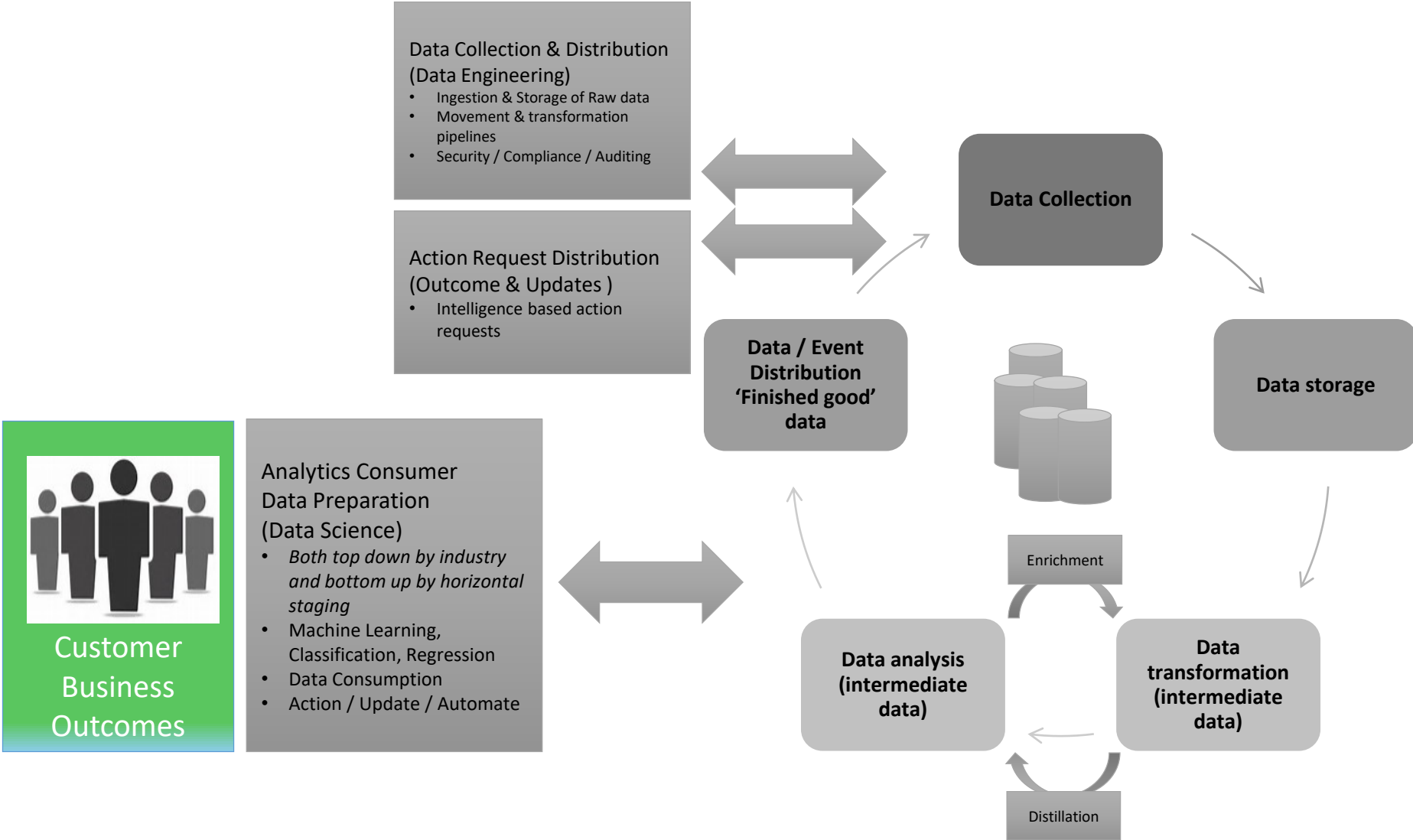
Cottrell.ken@gmail.com

(214) 546-100
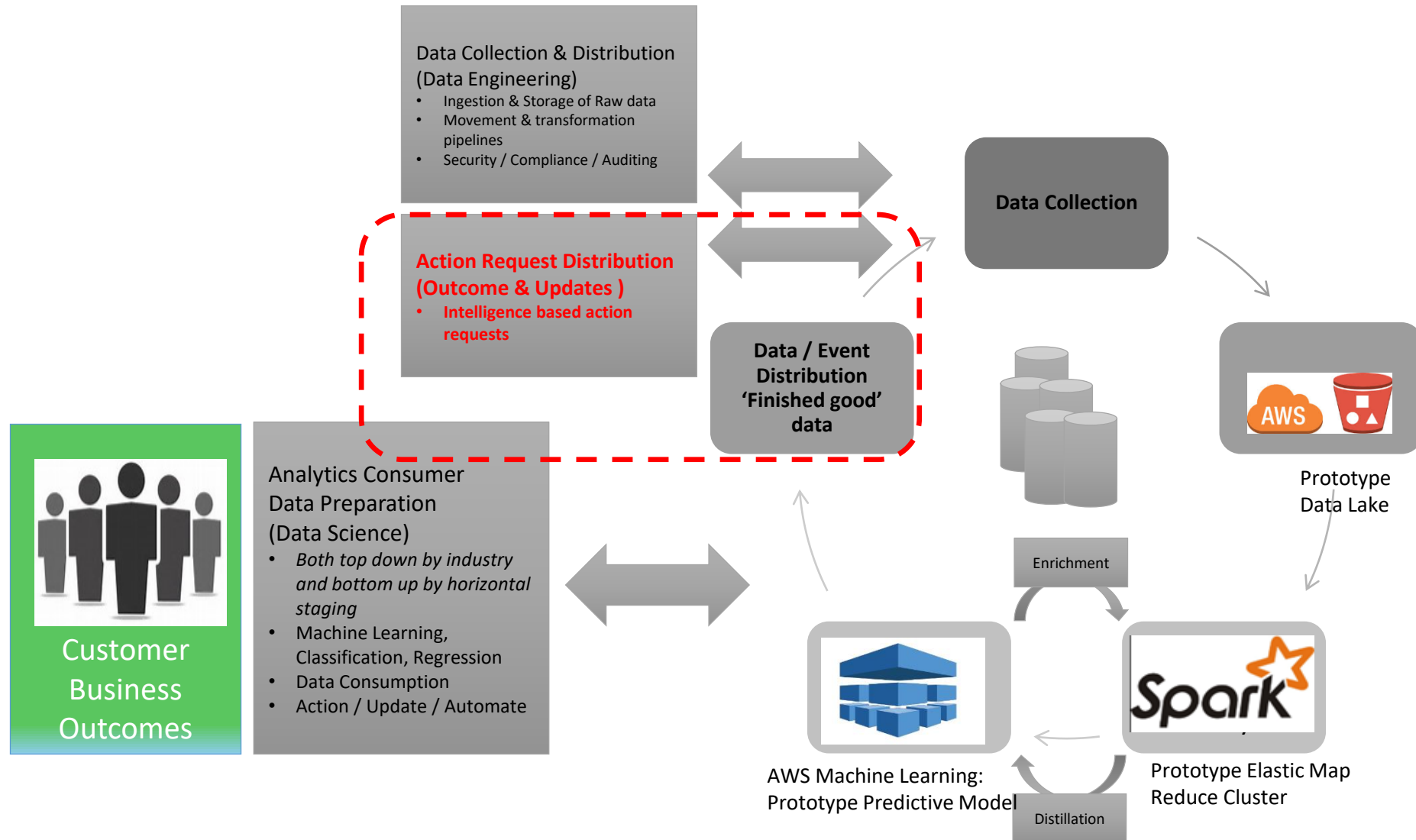
# agenda

- Problem Domain / Intro to the data set (in this case Public Healthcare Data)

- Demo & Tools discussion

- Summary, Lessons Learned

# Goal of this demo / Tutorial: Describe a phased Approach to a Data-Driven Architecture

**Data Collection & Distribution (Data Engineering)**
- Ingestion & Storage of Raw data
- Movement & transformation pipelines
- Security / Compliance / Auditing

**Action Request Distribution (Outcome & Updates )**
- Intelligence based action requests

**Analytics Consumer Data Preparation (Data Science)**
- *Both top down by industry and bottom up by horizontal staging*
- Machine Learning, Classification, Regression
- Data Consumption
- Action / Update / Automate

**Customer Business Outcomes**

**Data Collection**

**Data / Event Distribution 'Finished good' data**

**Data storage**

Enrichment

**Data analysis (intermediate data)**

**Data transformation (intermediate data)**

Distillation

# Goal of this demo / Tutorial: Describe a phased Approach to a Data-Driven Architecture

**Data Collection & Distribution (Data Engineering)**
- Ingestion & Storage of Raw data
- Movement & transformation pipelines
- Security / Compliance / Auditing

**Action Request Distribution (Outcome & Updates )**
- Intelligence based action requests

**Data / Event Distribution 'Finished good' data**

**Data Collection**

Prototype Data Lake

Analytics Consumer
Data Preparation
**(Data Science)**
- *Both top down by industry and bottom up by horizontal staging*
- Machine Learning, Classification, Regression
- Data Consumption
- Action / Update / Automate

Customer Business Outcomes

Enrichment

AWS Machine Learning: Prototype Predictive Model

Prototype Elastic Map Reduce Cluster

Distillation

# Problem domain (Healthcare)

**Functional**: *Providers who increasingly need to move to a Value-based care model, away from per-encounter payment model.*

- Providers want to predict & prevent Readmissions within 30 days, that treat the same condition within a 30 day window. Otherwise face non-reimbursements.

- Need to make use of multiple data sources: Electronic health records, Lab systems, Claims systems, Population Health evidence-based data (for baseline criteria), etc.

- Massive number of Attributes from many data sources need to be analyzed: difficult to manually create the business rules needed to predict

## Technical

- Massive capacity required for Data ingest, preparation, staging, analysis, refinement

- Time consuming to Find out (a) which measures ("features") have the most predictive power

- Have to iteratively experiment with and measure different algorithms to find the best predictive model



UCI Machine Learning Repository
Center for Machine Learning and Intelligent Systems

**Diabetes 130-US hospitals for years 1999-2008 Data Set**
Download: Data Folder, Data Set Description

Abstract: This data has been prepared to analyze factors related to readmission as well as other outcomes pertaining to patients with diabetes.

| Data Set Characteristics: | Multivariate | Number of Instances: | 100000 | Area: | Life |
|---|---|---|---|---|---|
| Attribute Characteristics: | Integer | Number of Attributes: | 55 | Date Donated | 2014-05-03 |
| Associated Tasks: | Classification, Clustering | Missing Values? | Yes | Number of Web Hits: | 133553 |

**Source:**

# Problem domain (All verticals)

**"Featurization" of data (Features in raw data may become Columns / Attributes in BI / DW)**

- Features often need to be extracted from lots of input data examples

- What is best tool to transform raw data into types that lend themselves to regression / classification? Target use cases (binary Classification , multi-class Classification, Regression) work best with different tools

- Which columns are just noise (depends on other columns or have useless data)  and can be removed to improve throughput?

- Often Requires many cycles to find the best models, and  therefore benefits from parallelization and elastic compute and storage services

Schema is not always clear from reading the values

Some columns may be redundant, or have unusable data

| | A | | C | D | E | F | G | H | I | J | K | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | discharge_disposition_id | admission_source_id | time_in_hospital | medical_specialty | num_lab_procedures | nu |
| 1 | encounter_id | patient_nbr | race | gender | age | admission_type_id | | | | | | ce |
| 2 | 2278392 | 8222157 | Caucasian | Female | [0-10) | 6 | 25 | 1 | 1 | Pediatrics | 41 | |
| 3 | 149190 | 55629189 | Caucasian | Female | [10-20) | 1 | 1 | 7 | 3 | ? | 59 | |
| 4 | 64410 | 86047875 | AfricanAmerican | Female | [20-30) | 1 | 1 | 7 | 2 | ? | 11 | |
| 5 | 500364 | 82442376 | Caucasian | Male | [30-40) | 1 | 1 | 7 | 2 | ? | 44 | |
| 6 | 16680 | 42519267 | Caucasian | Male | [40-50) | 1 | 1 | 7 | 1 | ? | 51 | |

# Problem domain: Data Analysis today is constrained by (a) inbound Data pipelines and (b) speed of analytics workflow

**Predefined** analytics processes and information

Capture data from **carefully defined** Enterprise sources

Stores only carefully formatted data

*Traditional ETL processes into BI/DW repository*

*Reporting processes*

- *New sources constantly added*
- *Constantly growing in volume, variety, velocity*

- *Existing Data transform & repository is Server-based, Static, Schema-first designs*
- *Expensive to reformat new data at scale: volume, velocity, variety requirements*
- *Don't really support self-service pipelines*

*Existing analytics tools can't model, evaluate, retest, deploy new insights fast enough for new business objectives*

- ***Current architecture:*** *historical Analysis of structured data using* **hard-coded** *business rules*

- ***New architecture : Predictive*** *analysis that* **learns** *business rules with massive amounts of training data*

# Removing the Constraints: Let Data Science / Data Analysts better utilize the upstream data pipeline

IOT / Streaming
Free Form text
Images
Voice
Video

*High volume capture of any / all kinds of data, schema format applied based on need*

*High volume transform, cleaning, enrichment, merging, forking*

*High throughput model, evaluate, compare, deploy*

Data Lake

**Parallel, Distributed In-memory Elastic Map Reduce**

**Parallel Distributed ML , AI**

amazon web services

Spark

hadoop

Machine Learning  AWS

H₂O.ai   IBM Watson

# Demo tools discussion : Elastic Map Reduce (EMR) & Machine Learning (ML)

## What is EMR and how does it augment "traditional" ETL, BI, DW?

- ETL, BI, DW still fit into the overall data pipeline but EMR provides the elastic compute and storage capacity to handle massive intermediate data processing via parallelization over large clusters

- For example, AWS provides ETL and DW systems for source and target stages in the EMR pipeline

- EMR used to mean Hadoop framework for disk-based batch, but now has evolved to include in-memory Apache Spark and other in-memory distributed process frameworks

## Why is it called Machine Learning?

- Automates many of the feature extraction, model creation, model evaluation, and model refinement steps by using parallel processing

- The business rules for prediction are inferred automatically, without explicit coding, based on training examples that contain the correct answer.

- It is able to infer the "rules" from data that is either too voluminous, too variable, or has too many attributes to specify in a rules table. Some examples of data include free text, voice, image, video, or tables with a large number of attributes such as health records.

- Algorithms include Classification such as Binary , Multi-class, Regression, Neural nets, etc

# Demo Data Workflow

# Demo analytics workflow

# Demo

# Summary / Lessons Learned

# Demo Summary / lessons learned (cont.)

The ML Predictive model created by my demo is considered poor (not much better than random guessing)

- *AWS ML suggests that these iterations may improve Model.*

- *Collect more data: Increase the number of training examples*

- *Feature processing: Add more variables and better feature processing*

- *Model parameter tuning: Consider alternate values for the training parameters used by your learning algorithm*

Ken's ideas about provided data set

- *for HIPAA compliance, had to strip out needed diagnostic, procedure codes and the like that fall under PHI*

- *Additional demo graphic data (such as Zip codes) or Census Block data, would likely add predictive value to a model. Again, too specific for public data.*

For Healthcare and other complex needs, a more specialized ML not provided by AWS may be required.

- *Note that you can still use AWS as an infrastructure platform in these cases. For example, vendors like H20.ai or DataRobot may be a good fit.*

# Demo tools: AWS Machine Learning

**Provides a ready-to-use Diabetes example, with instructions and sample data**

1. AWS ML Low cost

2. Automatically divide the data into Train and Test sets (70%/30% split)

3. Ask analyst to refine the schema as needed from the inferred schema

4. Ask analyst to specify the output (Binary Classification) value if present

5. Create a Model & Evaluation of Model's predictive quality0

*I tailored portions of the example*

• *AWS Tutorial Example uses Redshift, but I just used S3 in my account*

• *Enabled remote read / write into S3 from an external EMR cluster, to simulate a more complex Data Pipeline*

• *In the external EMR cluster (actually runs on AWS itself, but not in my account) added some Python "data engineering" DataFrame code, very simple*

• *Omitted a couple of the Redshift SQL steps (omitted Joins of Admit, Discharge codes as Numerical Categories instead of String Categories)*

# Demo tools: Apache Spark

## Elastic Map Reduce (EMR)

- de-facto standard for in-memory elastic processing

- Platform, API, SDK standards

- Improves throughput via integration with Hadoop framework

- Rish Support for Parallel distributed processing

- Rich language and ML library support (Python seems to be most popular but Scala and R also have their following)



**AWS has core EMR / Spark services, but I used Databricks as an integrated Front-end to AWS for convenience**

- Free community edition (includes Notebooks to run interactive code, file-import features, but you need licensed version for Jobs management on a larger cluster)

- Fully supports Apache Spark standard as a major contributor

- Convenient UI with Notebook code execution on Clusters, for Data Engineering and Data Science

- Built-in support for AWS clusters and Data import from S3 and other repositories

appendix

# Demo Summary / lessons learned

- ## Created S3 bucket as a prototype Data Lake

  - *stores raw data of any format, scales to almost infinite size and provides tools to overlay schemas for extraction*

- ## Python

  - *Why Python versus Scala or R? Python seems more (Big) Data Scientist friendly than R. R has been used a lot for workstation-level statistics but not as common in EMR. Scala may have better performance for some intensive data transformation processes.*

- ## Importance of a Data Pipeline approach to feed Advanced Analytics

  - *Key components: EMR clusters, Micro-services & Containers (Kubernetes seems to be gaining traction)*

  - *Need to have flexibility for Ad-hoc, agile way to insert, merge, split, serialize, parallelize data streams*

  - *Tradeoff between batch and interactive needed for ML. we should be conversant in both Hadoop and Apache Spark architectures. (They're actually converging anyway)*

# How I setup ML and EMR accounts

- ## S3 account, Free tier

1. Create buckets / folders
2. Upload the Diabetes dataset
3. Setup Spark Users, permissions, access keys

- ## Databricks account (free community edition)

1. Run a Cluster, create a Notebook to run in cluster
2. In Notebook key in Python Code (Shift-Return for each command cell)
3. Insert my Access Key and Secret Key from a special Spark user in IAM
4. code to Setup Read access to S3
5. code to transform data
6. code to Write cleaned Dataframe to S3

**AWS ML**

1. Set ML DataSource to S3 bucket folder with Cleaned from external EMR
2. Change some of the inferred column schemas
3. Select default Train/test mix (70/30)
4. Select binary classification output label (Readmitted = yes / no)
5. Run the train / evaluation session
6. Refresh screen to see when completed
7. Look at evaluation of model quality (ability to predict Positives / negatives
8. (not tried – run a sample record for a sample prediction)

# MFA for your account



aws

Amazon Web Services Sign In With Authentication Device

**The page you are trying to access requires users with authentication devices to sign in using an authentication code.**

**Provide your authentication code in the field below to complete sign in.**

Your Email Address:       kenneth.cottrell@verizon.com

Authentication Code:      771739

Sign In

Having problems with your authentication device? Click here

**About Amazon.com Sign In**

Amazon Web Services uses information from your Amazon.com account to identify you and allow access to Amazon Web Services. Your use of this site is governed by our Terms of Use and Privacy Policy linked below. Your use of Amazon Web Services products and services is governed by the AWS Customer Agreement linked below unless you purchase these products and services from an AWS Value Added Reseller. The AWS Customer Agreement was updated on June 28, 2017. For more information about these updates, see Recent Changes.

Terms of Use Privacy Policy AWS Customer Agreement © 1996-2017, Amazon.com, Inc. or its affiliates

An amazon.com company

# AWS IAM to create Spark users from external EMR

# IAM Spark user Access Key

# AWS Spark user (in IAM) Credentials to use in Python code

# DateBricks Community Edition

# Data Bricks landing page

# Starting a cluster for the Notebook

# Create a Notebook

# Data Bricks to S3 connection: many options, I used Boto framework

Python coding , you can cut and paste from my Textfile  into Notebook. Run one cell at a time or all together. At the end Upload cleaned data to S3

AWS ML screen shots – load cleaned data (sent from EMR) from S3

# AWS ML asks for permission to use S3, then asks you to proceed

# Some of the datatypes need to be reset from the inferred schema

# Set the last column (readmitted) as the binary classification output (i.e. the "Label"

# Review and complete datasource creation

To create just the ML model, using an existing datasource - Make sure you don't select one of the split train/test sets (unless that is your intent)

When you run the ML model, or ML Model and Evaluation (which does Train and test) – you may need to refresh screen to see if completed

View the evaluation. This demo's evaluation says our model is not much better than random (=50%) so considered a poor predictor

Model is considered poor – reason is AUC of 6.5 isn't much better than random (.5). Should be greater than .9 to be considered a good model

# Appendix

# Demo / tools discussion

- Why not use "traditional" ETL, BI, DW systems?

- Traditional centralized (server-based) tools can't handle capacity and unstructured data

- Need new parallel elastic scaling approaches to handle Volume, Velocity, Variety of data

- Enforce schema layout / data model too soon in the data factory process

- BI and DW are still good for "finished good" data for specific use cases, but not as suitable for "work in process" data engineering and data science stages

- Better approach: *New elastic cloud architectures* coming to market, coupled with a rich ecosystem of data management and data science tools that run in cloud

- Intermediate data:   Elastic Map Reduce, both batch-to-disk (Hadoop) and interactive in-memory (Apache Spark)

- Movement of traditional statistical languages like R and Python to a cloud parallel processing capability

- Machine Learning and AI are now able to leverage cloud capacity to automate many of the analytical workflows to calculate, compare, refine, and deploy the predictive models

- Pay only for temporary capacity needed to run huge jobs, then free it up when you get the insight needed

# Demo / tools discussion

"push" data model

- Why not use "traditional" ETL, BI, DW systems?

- Traditional tools can't handle capacity and unstructured data

- Need new parallel elastic scaling approaches to handle Volume, Velocity, Variety of data

"pull", self-service data model for ad-hoc constituencies

- layout / data model too soon in the data factory process

- good for "finished good" data for specific use cases, but not as suitable for "work engineering and data science stages

- etter approach: *New elastic cloud architectures* available, coupled with a rich ecosystem of data management and data science tools that run in cloud

- Intermediate data:   Elastic Map Reduce, both batch-to-disk (Hadoop) and interactive in-memory (Apache Spark)

- Movement of traditional statistical languages like R and Python to a cloud parallel processing capability

- Machine Learning and AI are now able to leverage cloud capacity to automate many of the analytical workflows to calculate, compare, refine, and deploy the predictive models

- Pay only for temporary capacity needed to run huge jobs, then free it up when you get the insight needed