# Pivota

#### **BUILT FOR THE SPEED OF BUSINESS**

## The Emerging Role of Data Science in Pharma: How to Harness this Transformative Practice

Sarah Aerni, PhD Senior Data Scientist

saerni@gopivotal.com PRISME May 15, 2014



## Agenda

- Driving Transformation through Data Science
- Pivotal Technology Enabling Innovation and Driving Insights
  - The right technology for the right job
  - Reinventing data processing through paradigm shifts
- Use Case Discussion
  - Predicting pharmaceutical potency







## **Enabling Disruption with Data Science**

• To successfully build a data-driven enterprise, skilled individuals need to have access to data, tools and channels for operationalization







#### What Is Pivotal Data Labs?

٠





#### **Data Engineering**

**Data Science** 











## The Quantified Patient



## Leveraging healthcare data to drive predictive and personalized care





## Data driven drugs: From discovery to delivery



#### **RICH DATA SOURCES**

- Molecular data
  - Cellular drug screens
  - Animal models
- Clinical data including notes, images, markers (e.g. genomics, lab results)
- Sensor and assay data
- Internal and partner/purchased external data
- Contact center data
- Patient registries, public and federal data, clinical partnerships



### The landscape of technology for big data

Sample Applications

Challenges

Use Cases



Batch processing of large volumes of data	Not optimal for highly iterative methods, functions over windows	Word count on tweets



Analytics on large- scale structured data	Requires restructuring of data to manipulate very large files	Predicting mortality on clinical data from diverse sources
--	---	--

HAMSTER/MPI Operations on very GraphLab large matrices

Requires knowledge of OpenMP, mis-used for embarrassingly parallel problems

Protein docking, molecular dynamics



## Choosing the right environment for different analytics challenges

	Imaging	Clinical E Narratives	Genetics
HD	Good for processing many images rapidly	Many documents with no shared processing	Read mapping
HAWO	In-database processing of very large images stored as a table	Information retrieval	BAM file manipulations, counts
HAMSTER/MPI GraphLab	Processing very large images		Multiple sequence alignment

Pivotal

### A new architecture for improved pipeline



## In-database genome-wide association study

#### **COVARIATES**

Indiv	Covariates		iates SNP				
	1	2		10	1	2	М
1	F	23		18	AA	СС	TT
2	Μ	39		41	AT	CG	TT
3	М	50		23	AA	GG	TC
	:	:					
Ν	F	19		24	TT	CG	ТС





### In-database genome-wide association study

#### **COVARIATES**

#### **GENOTYPES**

Indiv	Covariates			
	1	2		10
1	F	23		18
2	М	39		41
3	М	50		23
Ν	F	19		24

Indiv	SNP	Geno			
1	1	AA			
2	1	AT			
3	1	AA			
1	2	CC			
2	2	CG			
3	2	GG			
Ν	М	ТС			





## In-database genome-wide association study



#### Visualize and analyze genomics data without movement



Generate relevant plots using tools like Tableau immediately after parallel statistical analysis in-database on Pivotal technology



#### Visualize and analyze genomics data without movement



#### Visualize and analyze genomics data without movement



## Data driven drugs: From discovery to delivery



#### **RICH DATA SOURCES**

- Molecular data
  - Cellular drug screens
  - Animal models
- Clinical data including notes, images, markers (e.g. genomics, lab results)
- Sensor and assay data
- Internal and partner/purchased external data
- Contact center data
- Patient registries, public and federal data, clinical partnerships



## Vaccine Potency Prediction

#### **Business Problem**

Predict potency and antigen levels of live virus vaccines based on manufacturing sensor data and manual data collected throughout the process.



#### Simplified Vaccine Manufacturing Process



#### Enabling predictive models through new architectures

#### Challenges

- Accessibility
  - Some data had never been used in predictive modeling due to poor data models
- Data Integrity
  - Manually entered data is prone to errors. There is no immediate feedback to examine the validity of the values entered
- Data Completeness
  - Manual data entry is time consuming. There is no feedback on what data is most useful in improving the efficiency and quality and hence no prioritization of what data should be collected



Purpose-built data models for rapid data querying and exploration



Automated data cleansing techniques



Opportunities to eliminate collection of incomplete or non-predictive data



## Model generation and evaluation

Predicting vaccine potency using manufacturing data



Feature engineering and transformation

 Enabled by rapid in-database processing

- Experimentation with model forms
  - Partial least squares
  - Random forest
  - Regularized regression
- Interpretation of model results for insight generation
  - Use cross-validation framework to assess variable importance



## Sample model insights

Interpreting the utility of a measure obtained during manufacturing based on model outcomes



- Some features may reveal tunable parameters to alter potency, others may simply be markers
- Features consistently absent from models *may* be uninformative for predicting potency
- Opportunities to provide realtime feedback on data entry errors and predicted potency outcomes



## Data driven drugs: From discovery to delivery





## **Enabling Disruption with Data Science**

• To successfully build a data-driven enterprise, skilled individuals need to have access to data, tools and channels for operationalization





# Pivota

#### **BUILT FOR THE SPEED OF BUSINESS**