

Don't Forget To

ODTUG  
Kscope22  
grapevine, tx    june 19 - 23

Fill Out Your Evals



SOFTWARE SOLUTIONS

# Lessons from Implementing Machine Learning Using OAC and ADW

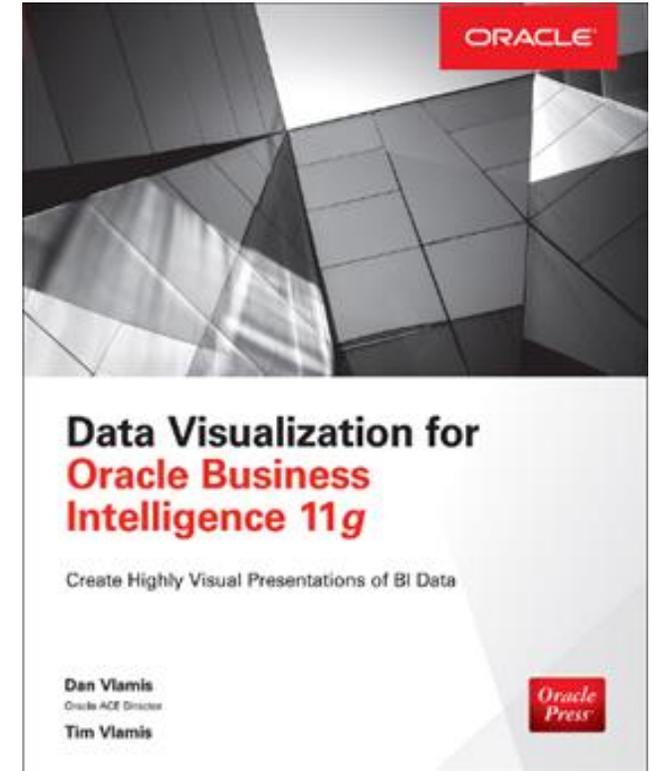
Tim Vlamis

June 20, 2022

[www.vlamis.com](http://www.vlamis.com)

# Vlamis Software Solutions

- Founded in 1992 in Kansas City, Missouri
- 400+ Enterprise Clients
- Consults in :
  - Enterprise Business Intelligence & Analytics
  - Analytic Warehousing
  - Machine Learning and Predictive Analytics
  - Data Visualization
  - ETL and data integration
- Vlamis consultants average 15+ years
- [www.vlamis.com](http://www.vlamis.com) (blog, papers, newsletters, services)
- Co-authors of book "Data Visualization for OBI 11g"



# Instructor's Background

## Tim Vlamis – Vice President & Analytics Strategist

- 30+ years in business modeling and valuation, forecasting, and scenario analyses
- Oracle ACE Director
- Co-author of “Data Visualization for Oracle Business Intelligence”
- Named contributor to and instructor for several of Oracle University’s Machine Learning and Predictive Analytics Courses
- Professional Certified Marketer (PCM) from AMA
- MBA Kellogg School of Management (Northwestern University)
- BA Economics Yale University
- tvlamis@vlamis.com

# Many Words Used for Similar Concepts

Predictive Analytics

Regression Data Mining

SQL Analytics Anomaly Detection

Python Adaptive Intelligence

Data Science

Diagnostic Analytics

Classification

AI

Advanced Analytics

Algorithm Descriptive Analytics

SQL R

Clustering

Artificial Intelligence

Prescriptive Analytics

Machine Learning

THIS IS YOUR MACHINE LEARNING SYSTEM?

YUP! YOU POUR THE DATA INTO THIS BIG PILE OF LINEAR ALGEBRA, THEN COLLECT THE ANSWERS ON THE OTHER SIDE.

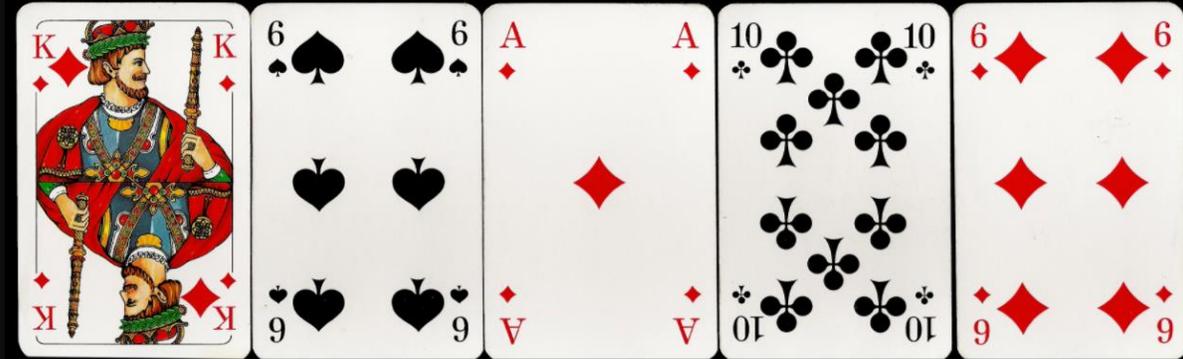
WHAT IF THE ANSWERS ARE WRONG?

JUST STIR THE PILE UNTIL THEY START LOOKING RIGHT.



# Probabilities, not Outcomes

- The likelihood of an event defines what is a good bet, not whether it wins or loses.



# Four Realms of Analytics

Probability Based	<b>Diagnostic Analytics</b>	<b>Predictive Analytics</b>
Rules Based	<b>Descriptive Analytics</b>	<b>Prescriptive Analytics</b>
	Past	Future

# ML Decision Making Matrix

<p><b>Longevity of Decision</b></p> <p>Years</p>	<p><b>Exception Analysis and Anomaly Detection</b></p>	<p><b>Scenario Analysis and Statistical Modeling</b></p>
	<p><b>Prescriptive Analytics and Real Time Decisions</b></p> <p>Days</p>	<p><b>Predictive Analytics and Workflow Optimization</b></p>
<p>Reaction</p>		<p>Planned/ Guided</p>
<p><b>Latency of Decision</b></p>		

# Business Use Cases of Machine Learning

- Production
- Sales
- Marketing
- Finance
- HR
- IT
- Logistics

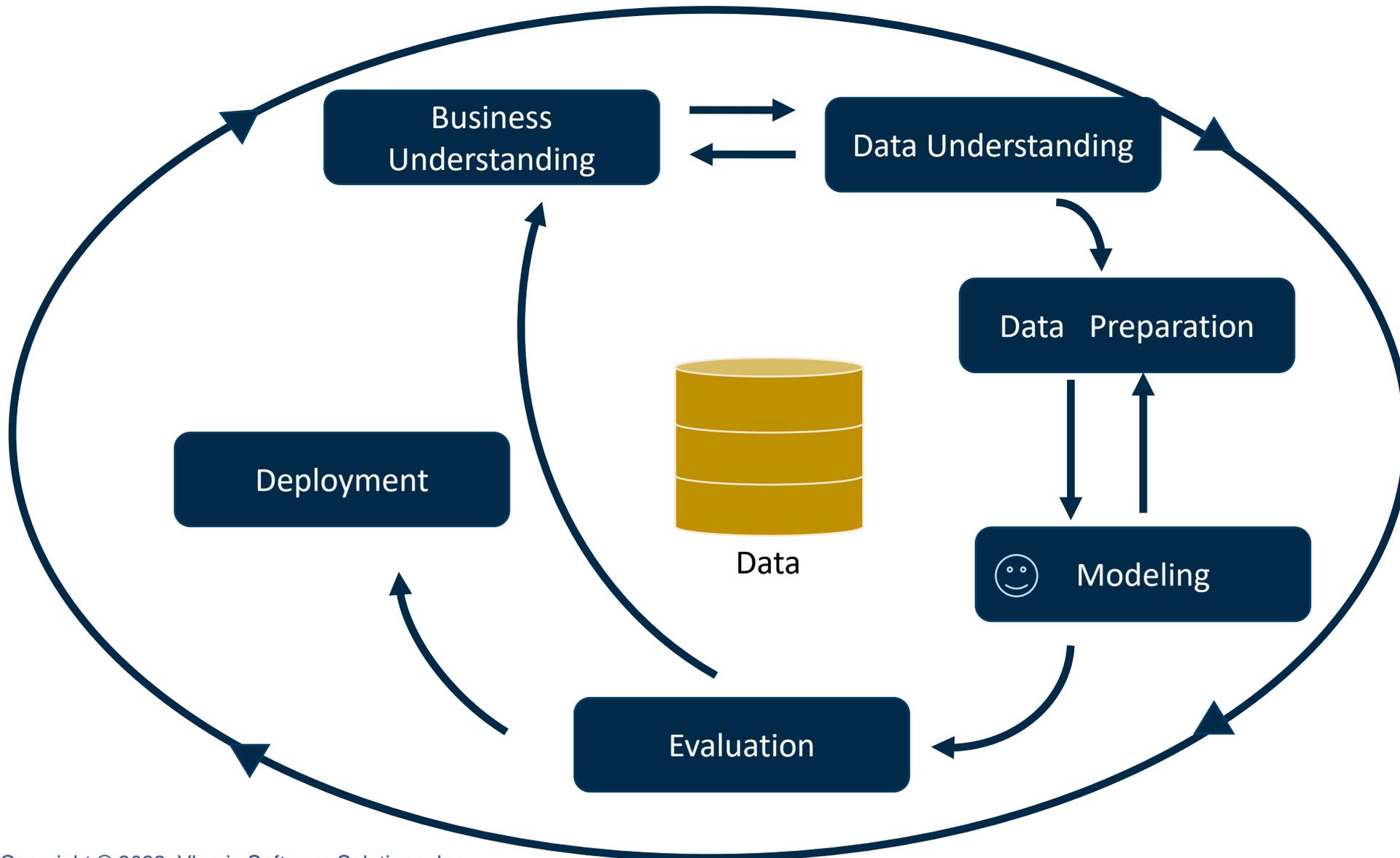
# Classification of Potential Customers

- Predict which people were most likely to be customers of health clinics.
- Several hundred variables
- Significant issues with data provided by syndicator
- Significant issues with data provided by clients
- Lack of clinic attributes
- Lack of timing attributes
- How many models to build

# Behavioral Segmentation

- Clustered approx 1000 restaurant locations into natural groups
- Rely on behavior, not traditional attributes for segmentation
- Many expected and unexpected patterns revealed
- Segmentation informed the development of marketing programs

# CRISP-DM Phases



# Questions for Data Scientists

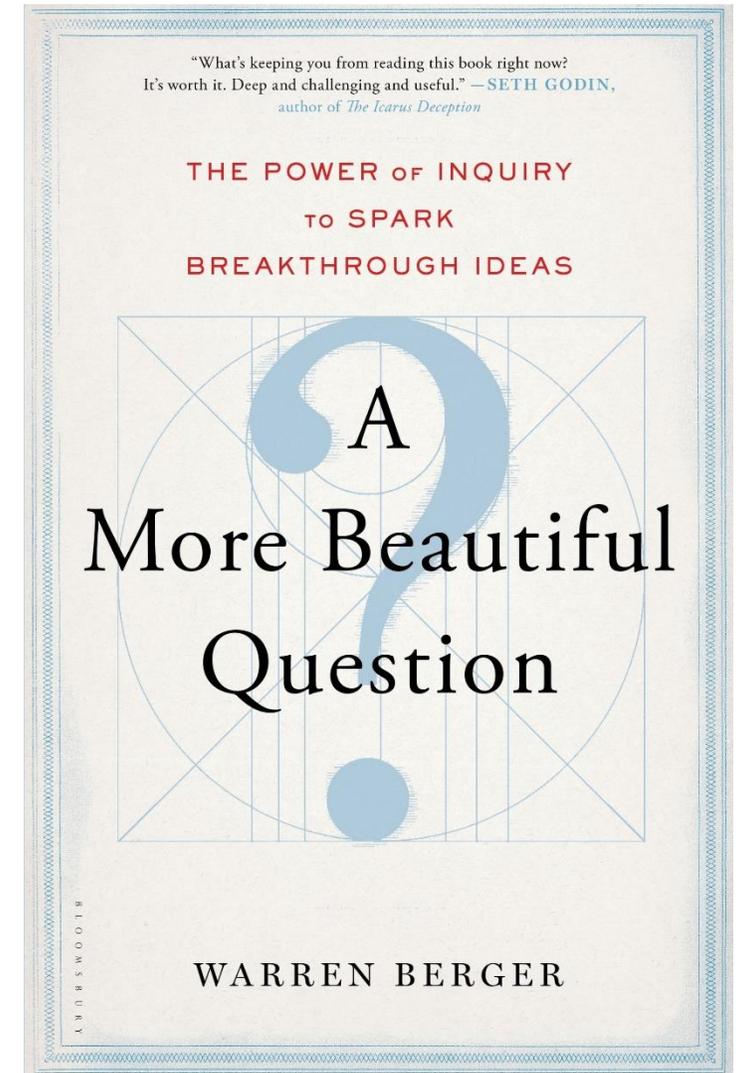
- What problems are you trying to solve?
- What algorithmic strategies provide the most value?
- Who is your audience?
  - Executive, professional analyst, IT, functional business management?
  - What background in consu
- How clean is the data?
  - Data created by transactions or analysis?
- How consistent is the data?
- Data used for reporting/analysis or in automated business process?
- When was the data gathered and when refreshed?
- Where does the data reside? Which platforms, applications, interfaces, and environments?

# Good Questions/Hypotheses are Needed

What behaviors in the past year are most significant in terms of segmenting our customers?

What's the Life Time Value of each customer?  
What's a potential new customer worth?

Which products are purchased together most often? Which products are purchased with our most profitable products?



# Bounded vs. Unbounded Domains

- Bounded games like poker, baseball, elections, web A/B testing, etc.
  - Defined rules, time, and results
  - Can use “classic” statistics for prediction
  - Scale space is predetermined
  - Neural nets are excellent for classification exercises in bounded domains
- Unbounded games like economic growth, forests, profitability, etc.
  - Significant challenges exist using “classic” statistics
  - Assumptions are both necessary and more important than anything else
  - Scale space is undetermined
  - AI does not do well with unbounded predictions

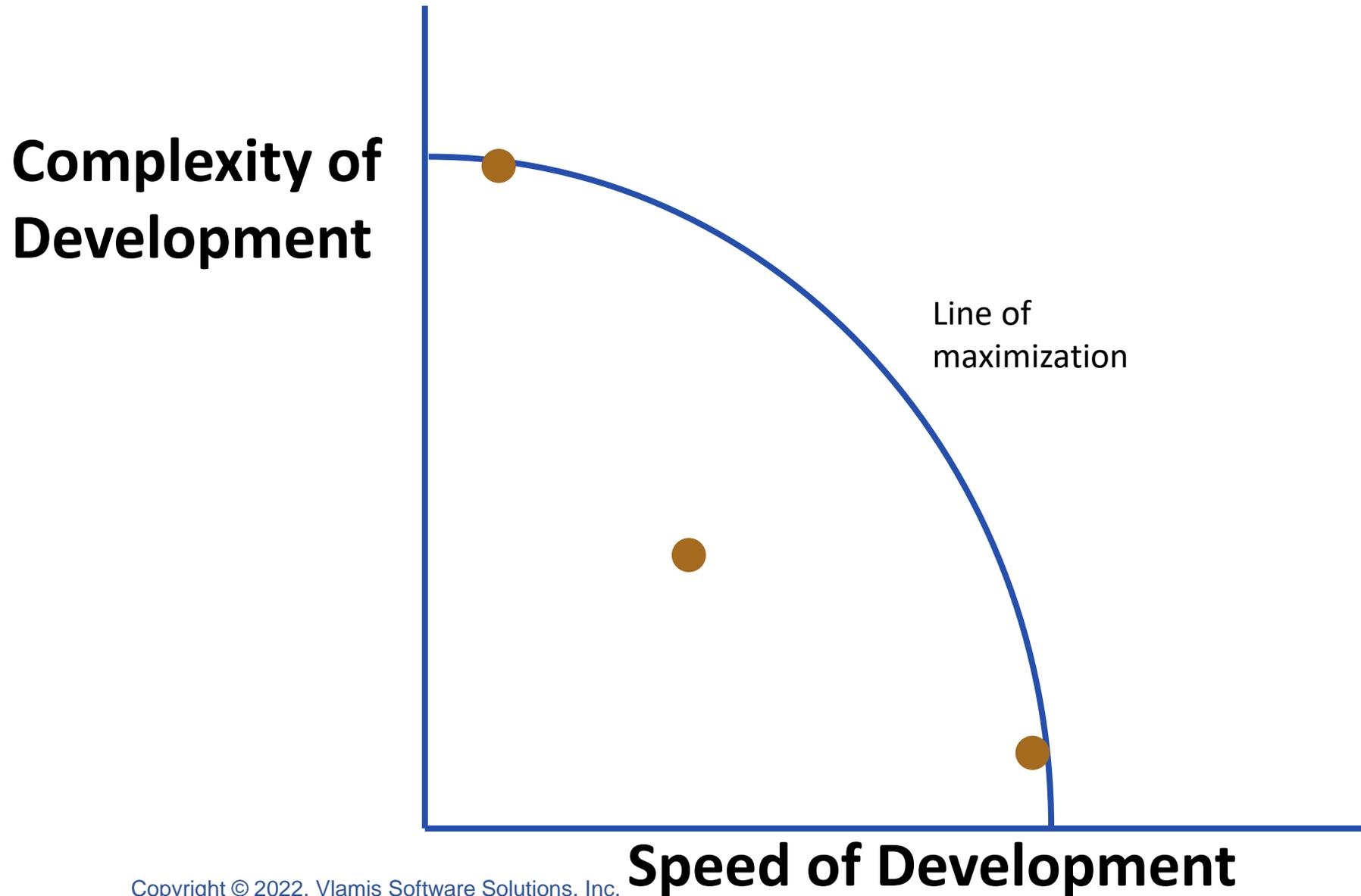
# Move the Algorithms to the Data

- Moving data is expensive
- Replicating data makes validation hard
- Aggregation requires consistent data

# Ethical/Risk Frameworks for AI & ML

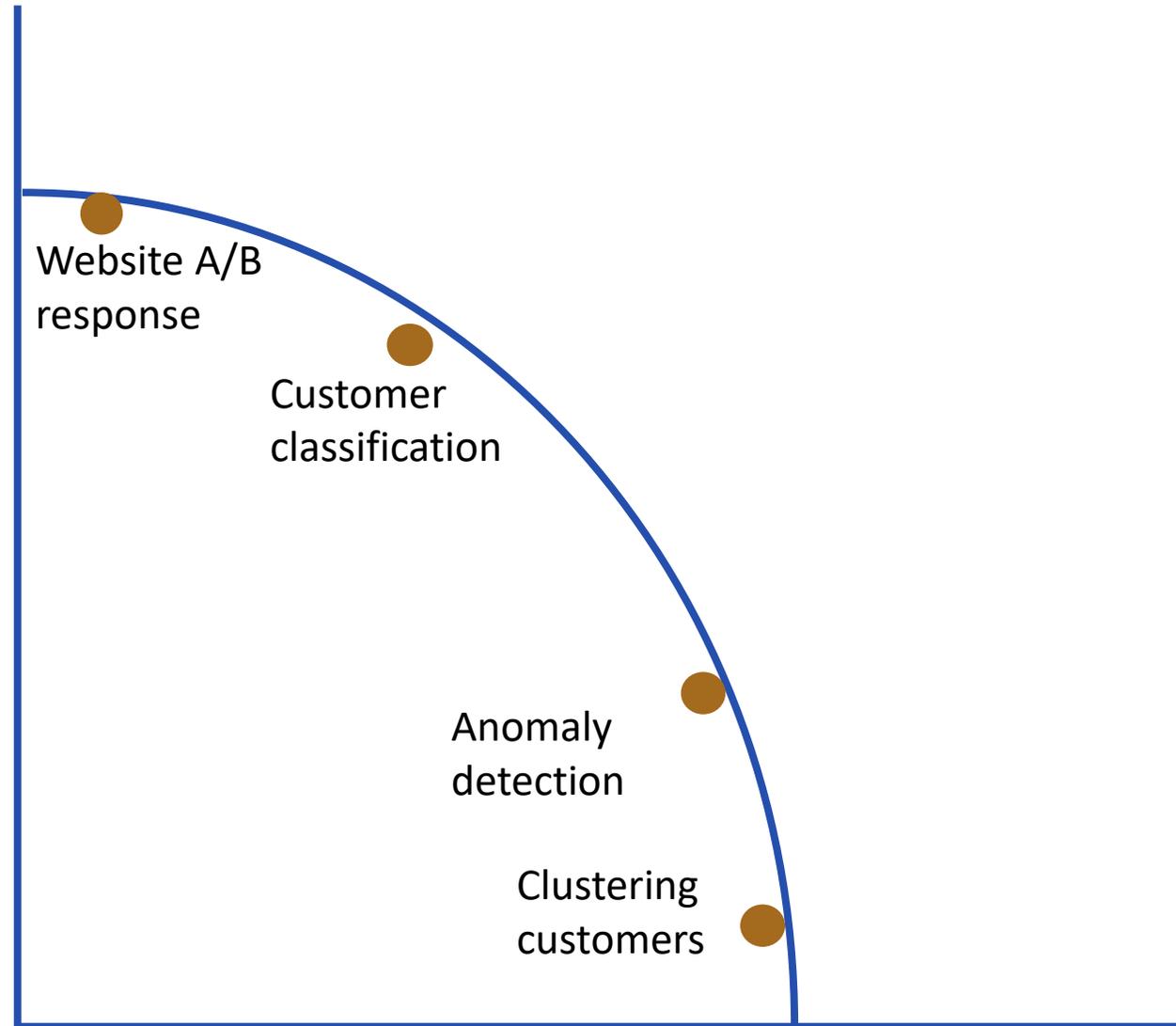
- Stakeholder analysis
- Negotiations/shared interests
- Fiduciary responsibility
- Risk management
- Security
- Data governance and Master Data Management
- Distributive Justice, Ethics, and Moral Philosophy
- Legal framework (HIPAA, FCRA, EU GDPR, etc.)
- Data Mining Frameworks (KDD, CRISP-DM, etc.)
- Complex Adaptive Systems, Systems Dynamics

# Horizon Functions Show Tradeoffs



# Machine Learning Tradeoffs

**Specific hypothesis test in bounded realm**



**Exploration – Hypotheses development**

# Dos and Don'ts for ML & Predictive Analytics

## ■ Do

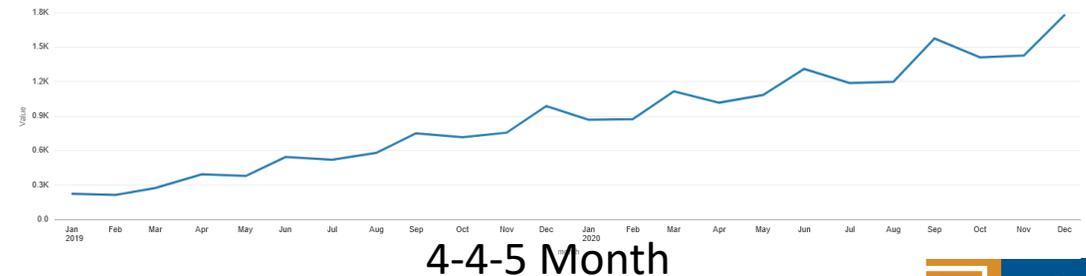
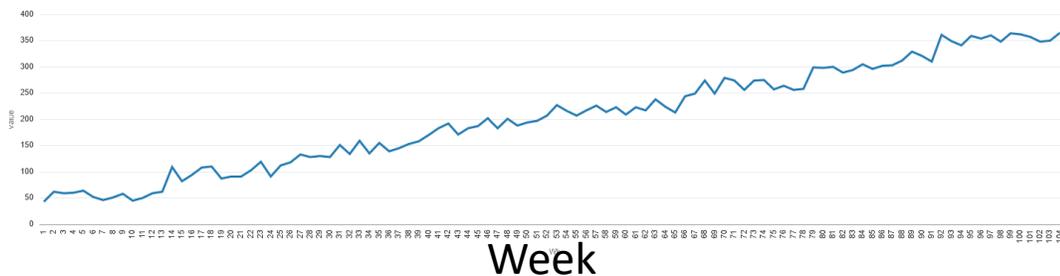
- See ML analytics as a continuous process
- Invest in an overall strategy, not in isolated tools and projects
- Put enough resources in place
- Grow your own talent and systems and involve IT
- Seek guidance and expertise early on

## ■ Don't

- Export your internal data and outsource to "experts"
- Try to develop predictive analytics "on the cheap"
- Anoint then isolate internal experts
- Build multiple, parallel infrastructure systems without IT
- Struggle with basics and then reinvent the wheel

# Forecasting for 30,000 Products

- A medical equipment company had complex forecasting needs
- International demand and distribution
- International production facilities
- High degree of seasonality (sports injuries)
- 445 pattern of months was creating problems for exponential smoothing algorithm
- Automated algorithms are powerful, but need experienced implementers.



# Machine Learning Advice

- Use the same technology for development as production
- Focus on using clean data that is important to the organization
- Remember it's about value delivered, not technology used
- Don't start with the hardest problems
- Do start with the obvious use cases
- Work in the analytic warehouse

# Dirty Data is a Pollutant



Copyright © 2022, Vlamis Software Solutions, Inc.

# Clean Data is Essential



Copyright © 2022, Viamis Software Solutions, Inc.

# Unsuspected Dirty Data

- Undefined nulls or zeros (mean zero or no data?)
- Numeric scores that are not normally distributed
  - Grades
    - Data elements scored 1 – 10 with strong clusters
- Numeric scores with inconsistent ranking
  - low-high or high-low?
- Hidden hierarchical attributes that are inconsistently applied
  - Manufacturing error/reason codes example
  - State, city, county, metro codes

# Keeping Data Clean is Easier than Cleaning It

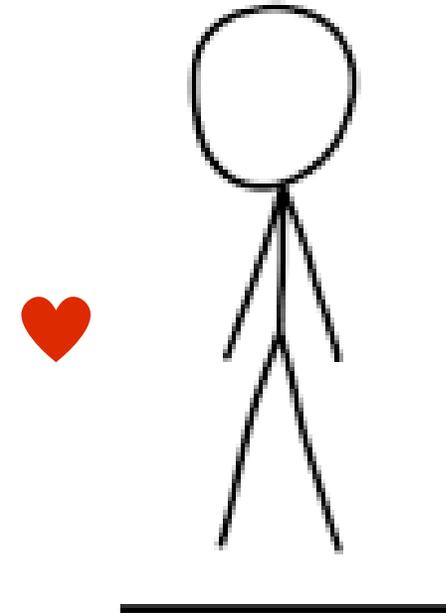
- Start with clean data
- Keep it clean
- Don't move dirty data
- Isolate dirty data

# Clean Data is

- Accurate
- Current
- Consistent
- Complete
- Conformed
- Congruent

# Head Hands Heart

- Listen carefully, speak truthfully and humbly
- Handle data skillfully
- Be responsible at all times



# Abstraction

- Abstraction can reduce replication and increase dynamic integration
- Too many layers of abstraction can create “black box” systems that are difficult to understand
- Be careful “embedding” abstractions in code that are not documented. Alias of an alias of an alias of an alias from different subsystems with no consistency or pattern or documentation or organization.

# What tools to use?

- For those who think through coding, use Zeppelin notebooks in Oracle Machine Learning
- For those who think visually, use Oracle Data Miner in SQL Developer
- For who use Python, learn to leverage the transparency layer and execute inside ADW or Oracle Database

# Tradeoffs

- Don't duplicate data  $\leftrightarrow$  Have data where you need it
- Conform and rationalize data  $\leftrightarrow$  leave data in its original form
- Assume that all data is the same  $\leftrightarrow$  assume that all data is unique
- Names are for systems  $\leftrightarrow$  names are for people
- Save everything  $\leftrightarrow$  save nothing

# Bill Inmon on Analytic Warehousing

“The problem is that after you have created the database, no one really wants a database. Managers of the world want visualisation, they don’t want a database. What they really want are things presented to them rapidly and beautifully.”

Analytics India Magazine 04/06/2020 “Data Scientists are Actually Data Garbage Collectors”

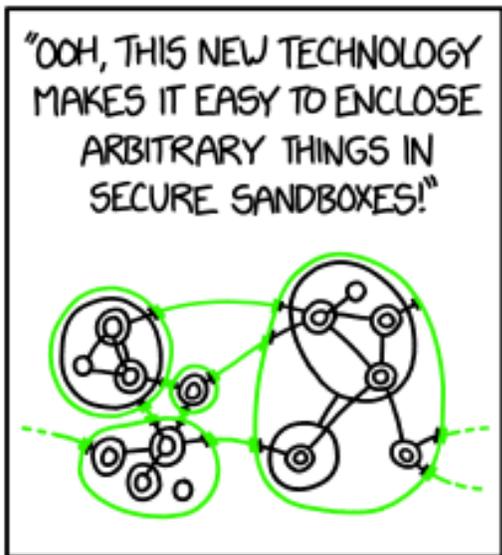
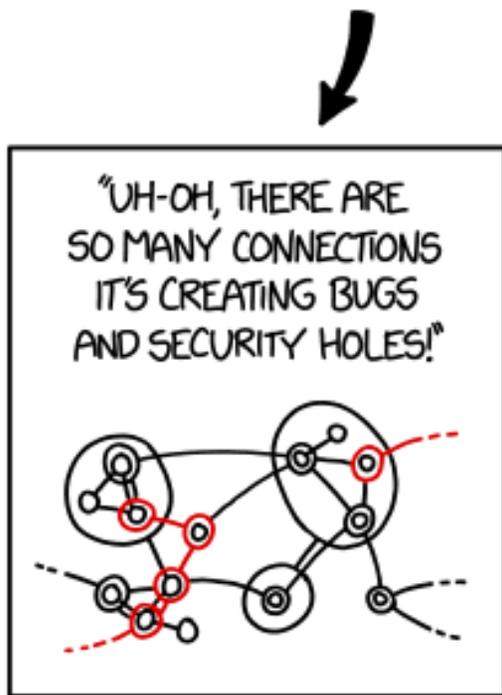
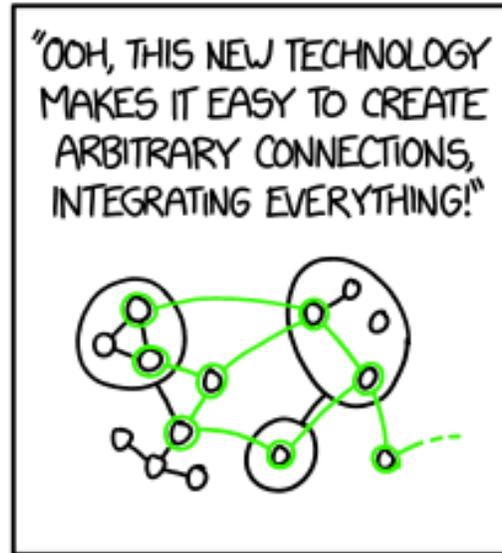
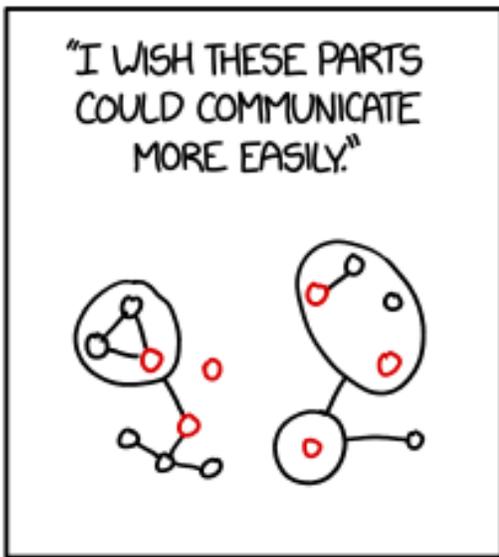
# Listen to Data to Discover Structures

- Relative importance
- Natural relationships
- Similarities/differences
- Predictions



# Discovered Structures in Data

- Unsupervised Machine Learning Algorithms
  - Hierarchical clustering (Distance, Density, Partitioning)
  - Dimensionality reduction (PCA, SVD)
  - Association Rules
- Property Graph Algorithms (network theory)
  - Casts data sets as nodes and edges (elements and relationships)
  - Customer A buys from Factory B (two elements one relationship)
  - Centrality, connectedness, degrees of separation, page-rank, dozens more



# Accessibility vs. Security



# Limitations and Challenges

1. An information retrieval system will tend not to be used whenever it is more painful and troublesome for a customer to have information than for him not to have it.
    - a) Where an information retrieval system tends not be used, a more capable information retrieval system may tend to be used even less.
- Calvin Mooers 1959
  - [https://en.wikipedia.org/wiki/Calvin\\_Mooers](https://en.wikipedia.org/wiki/Calvin_Mooers)



# Thank You!!

- Tim Vlamis

[tvlamis@vlamis.com](mailto:tvlamis@vlamis.com)

Don't Forget To

ODTUG  
Kscope22  
grapevine, tx    june 19 - 23

Fill Out Your Evals