

# 8 Ways We Fail With Predictive Analytics in Business

**Stephen Smith**

Research Director, Data Science  
Eckerson Group

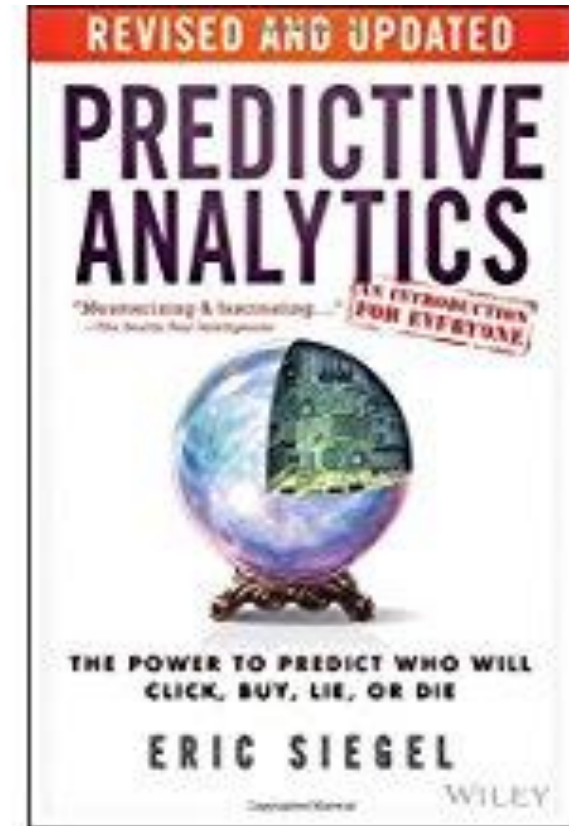
**Sponsored by:**



Sponsored by:  
**MARINER**  
INNOVATIONS



- **Best web ad selection**
  - *3.6% more revenue*
  - *\$1 million every 19 months*
  
- **Decreased loss ratio in insurance**
  - *0.5% reduction in loss ratio*
  - *\$50 million annually*



Great book!!!

Reference: "Predictive Analytics" p.25

Sponsored by:  
**MARINER**  
INNOVATIONS

1. Every business **should be using** predictive analytics
2. Predictive analytics is **not being used** often enough
3. Because it is a bit **complex and dangerous**

# We Interviewed Industry Experts



cloudera

alteryx



DOMINO



MONSANTO



data  
iku

TERADATA



nu  
tonian



Informatica



Microsoft



1. Every business **should be using** predictive analytics
2. Predictive analytics is **not being used** nearly enough
3. Because it is a bit **complex and dangerous**
4. It needs to be **automated and operationalized**



**Where's  
My Data?**



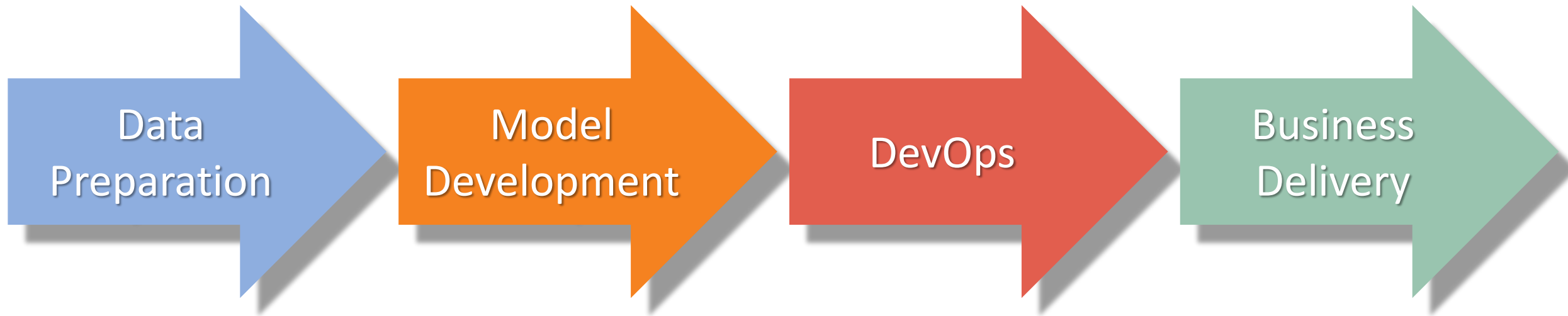
Predictive  
Analytics

**What do I  
do with  
this model?**



**Business**

# Full Data Science Lifecycle





Promise

- Expectations of reduced churn, increased cross sell etc.
- Hires 'data scientists'
- First models created

Fear

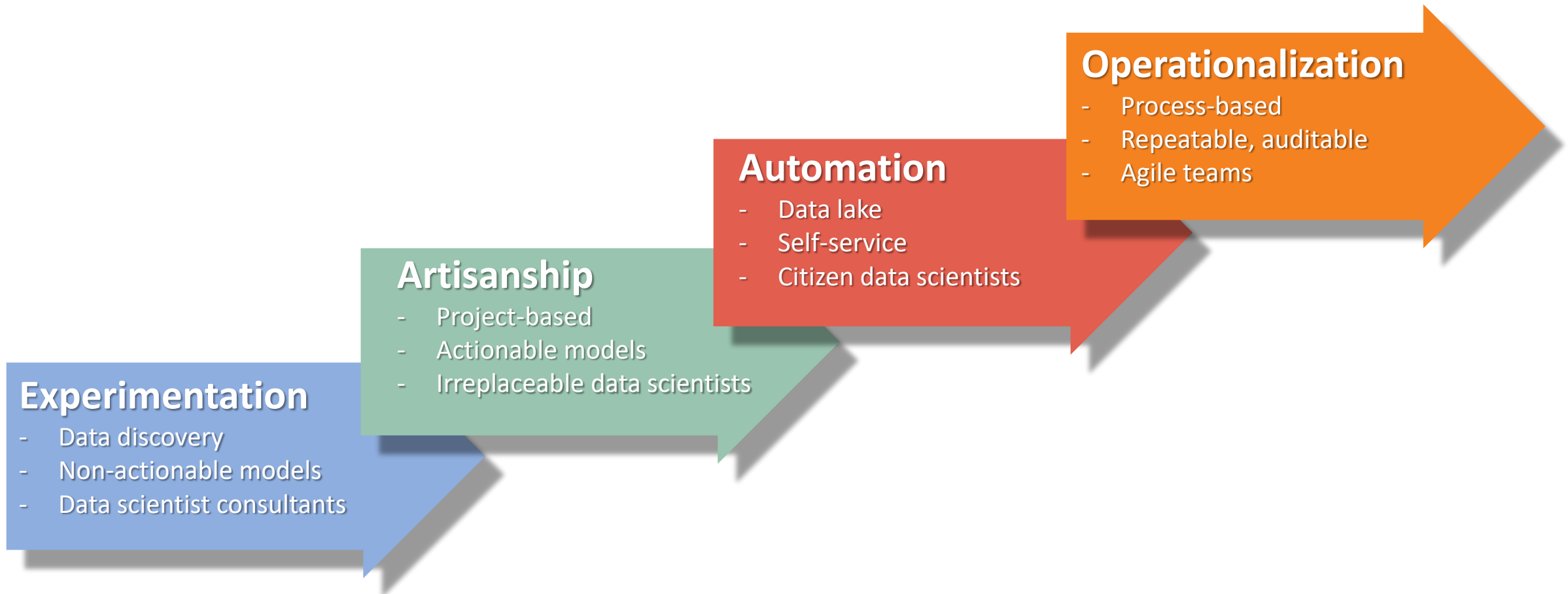
- Model doesn't work in production
- Models are late
- Data scientist quits

Trust

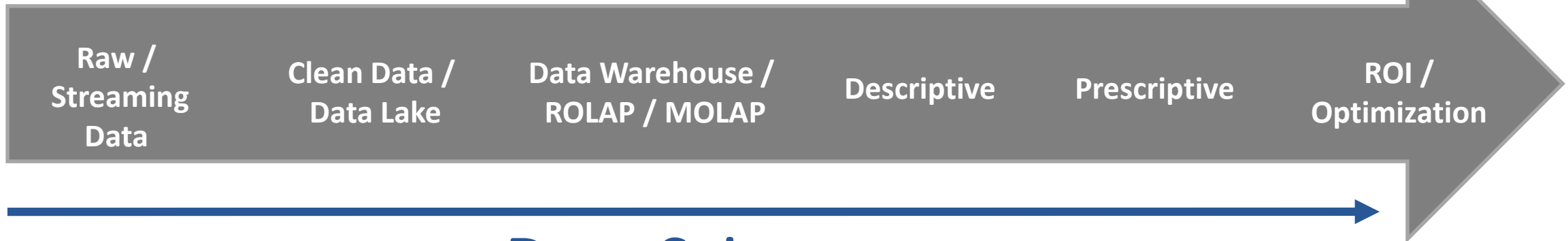
- Focused projects begin to succeed
- Business champions emerge requesting models
- Best practices evolve and are enforced

Acceleration

- Hundreds to thousands of active models
- Doubling data scientists quadruples number of models
- Business users expect to use models



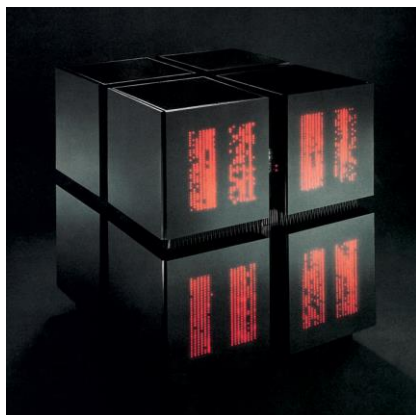
## Business Intelligence



## Data Science

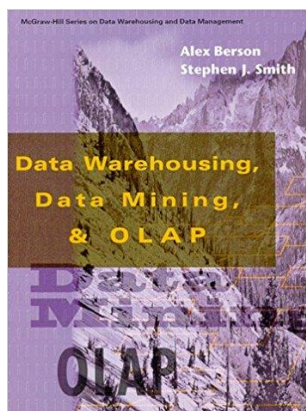


MPP

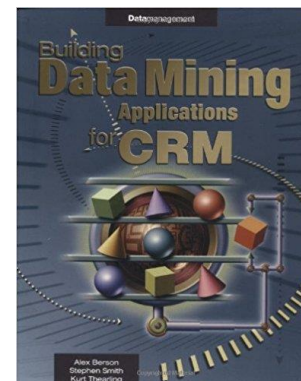


1990s

DW + OLAP



Data Mining



2000

Pharmaceuticals



2010

Education



“Data Science”

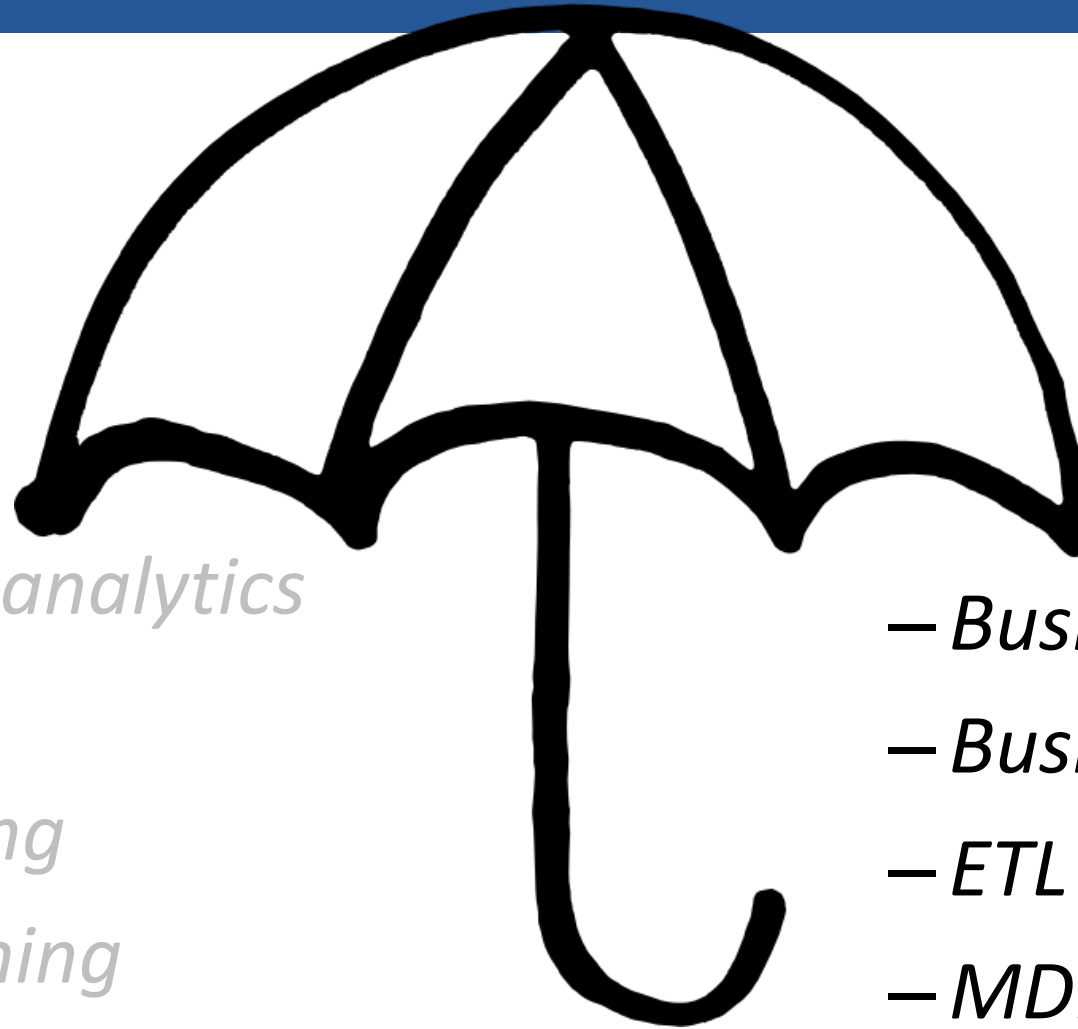


- Massively Parallel Processing => Hadoop
- OLAP => Business Intelligence
- Data Mining => Predictive Analytics
- Neural Networks => Deep Learning
- Data Warehousing => Data Lake (+MDM,DG,ETL...)
- - => Data Science

- ‘Computer Science’ was Born in 1959
- A need to describe the practice of using computers
- Interestingly:
  - *The alternative name proposed for “Computer Science” was “Data Science”*



- *Predictive analytics*
- *Statistics*
- *Data mining*
- *Deep learning*
- *Artificial Intelligence*
- *Machine Learning*



- *Predictive analytics*
- *Statistics*
- *Data mining*
- *Deep learning*
- *Artificial Intelligence*
- *Machine Learning*

- *Business intelligence*
- *Business analytics*
- *ETL*
- *MDM*
- *Data engineering*



# The Elephant in the Room

A close-up photograph of an elephant's eye, showing the intricate texture of its wrinkled, brownish-grey skin. The eye itself is a deep, reddish-brown color with a dark pupil and a prominent, dark, hairy eyelid. The lighting is soft, highlighting the texture of the skin and the details of the eye.

Why isn't data science used more?

“I have data?”



“Where’s my data?”



“What will happen next?”

“What should I do next?”

“How do I start using predictive analytics?”



**WARNING**

“Today we can process Exabytes of data at lightning speed, and this gives us the potential to make bad decisions far more quickly, efficiently and with greater impact than we did in the past.”

- Susan Etlinger

*Ted Talk: What do we do with all this big data?*





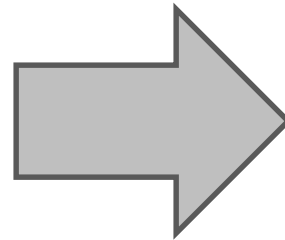


Uranium Ore = Perfectly Safe

# You Make Plutonium from Uranium

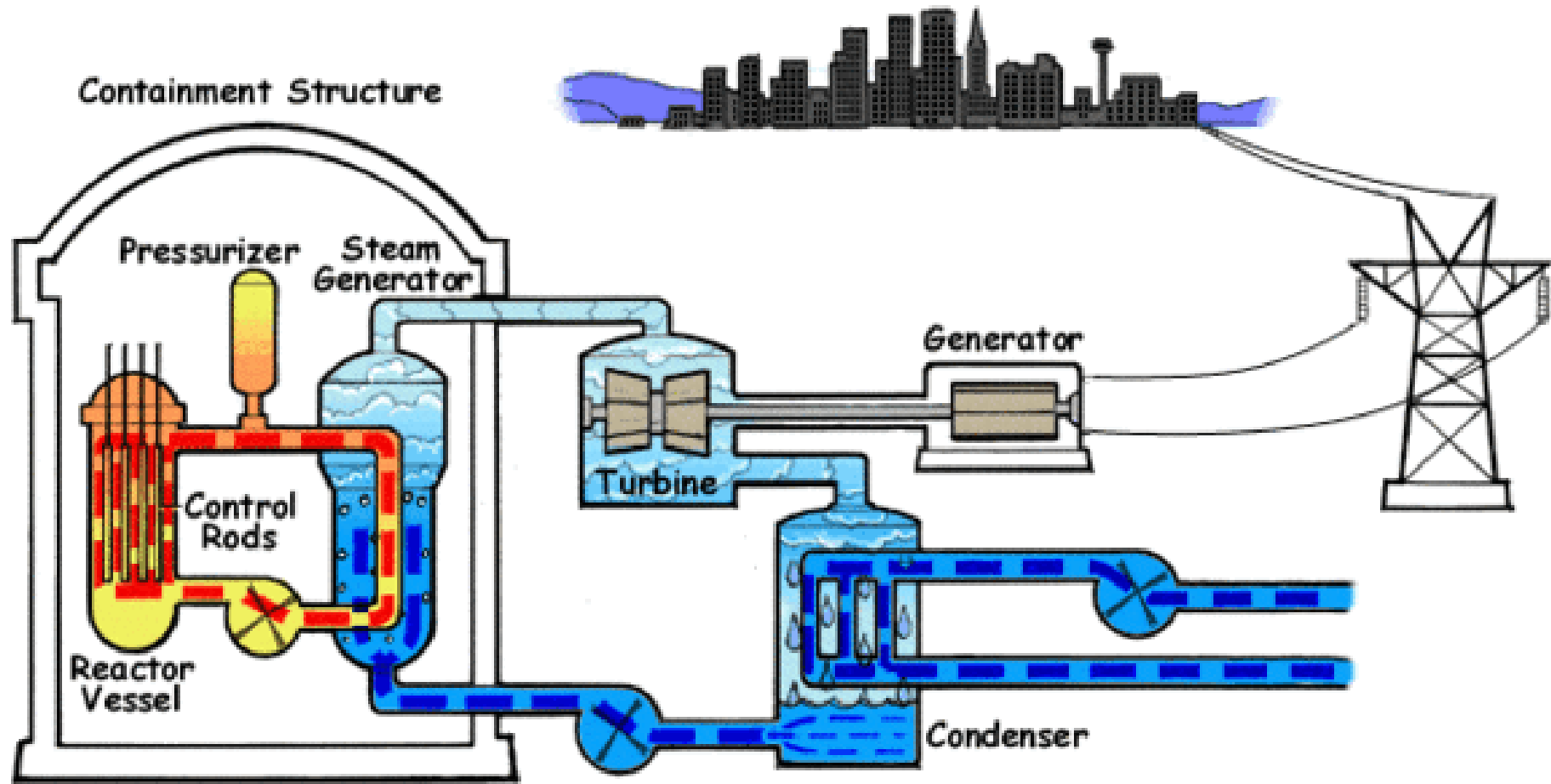


Data = Uranium Ore



Data Science = Plutonium

# Careful Process Control Required



	<b>Plutonium</b>	<b>Predictive Analytics</b>
<b>Powerful</b>	Generates electricity	Drives new revenue with little investment
<b>Dangerous</b>	Explosions	Mistakes can cost hundreds of millions of \$\$
<b>Handle with care</b>	Requires operationalized processes and tools	Requires operationalized processes and tools



# WHERE DATA SCIENCE FAILS

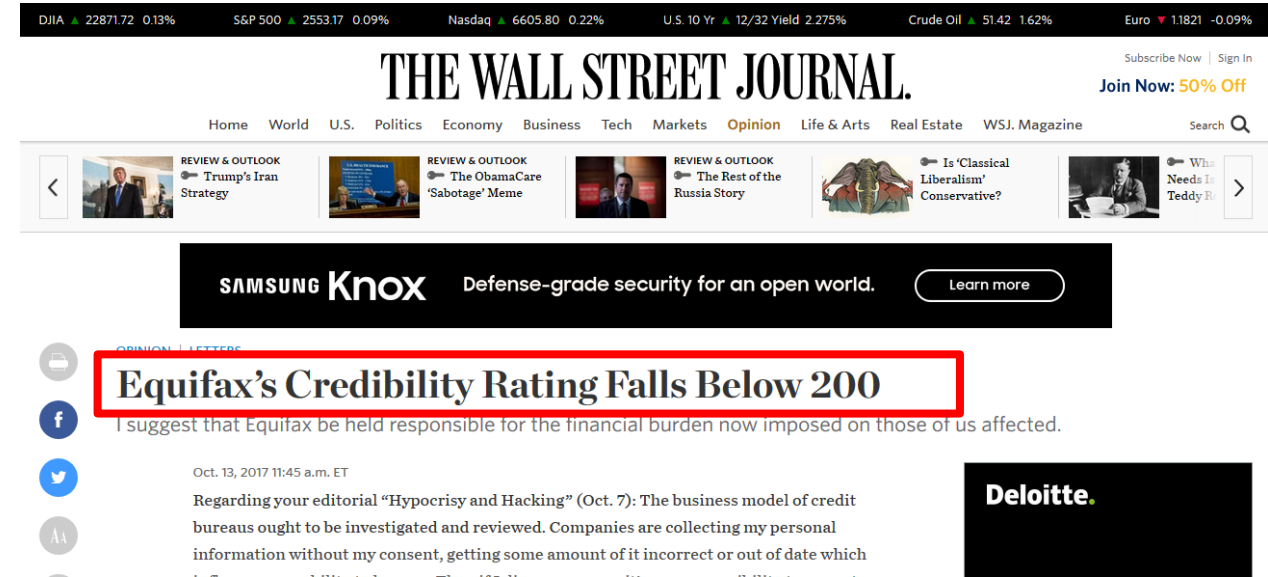
# How costly are data science fails?

1. Lose all the investment in a hedge fund.
2. Send expensive offers to the wrong customers.
3. Ruin a brand.

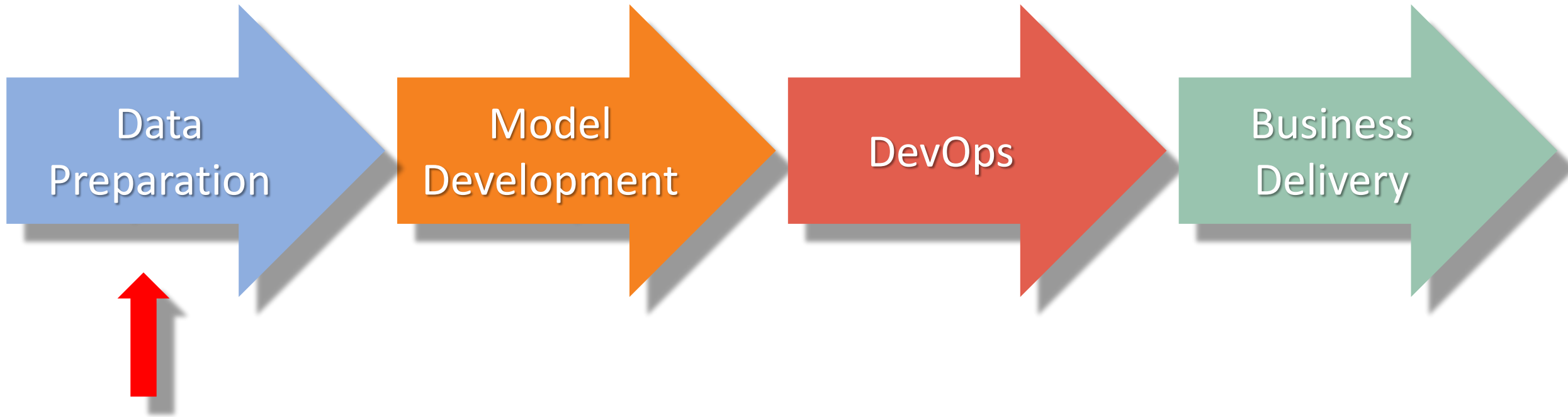
\$184-\$334 million



*“Reputation Impact of a Data Breach”*







# #1 Data Drift, Anomalies, Errors

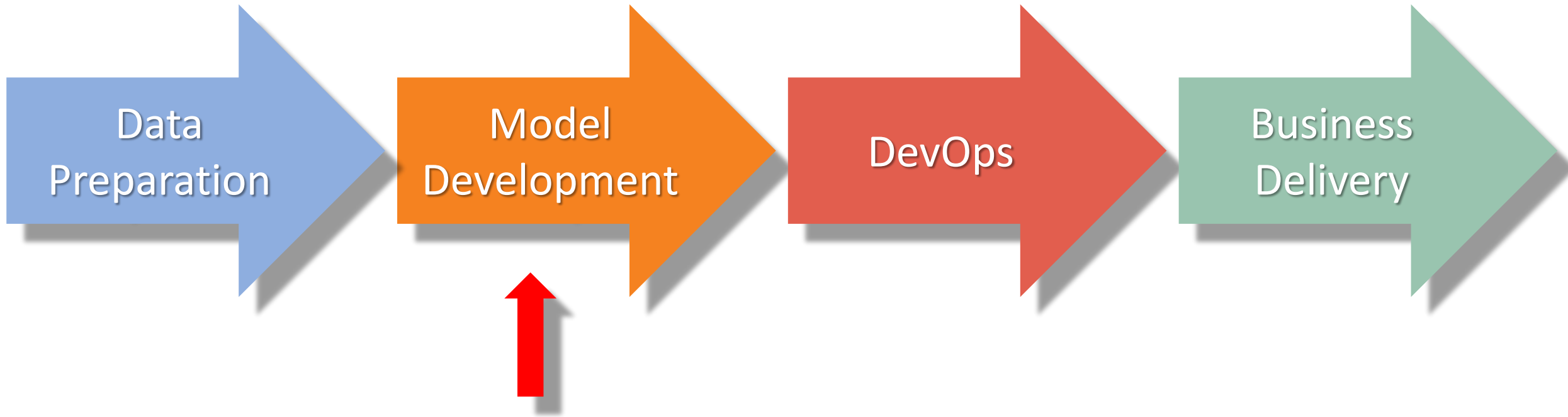
- **Solution**
  - *Proprietary algorithms: radial basis function, Black-Scholes, ANN*
- **Failure**
  - *Model prediction didn't match actual*
- **Why?**
  - *Errors in industry-wide data sets used by all hedge funds*

- **Challenge**

- *Input data values can change over time*
- *Errors can creep into data*

- **Best Practices**

- *Tools that have automated alerts for data changes*
- *Algorithms that manage missing or corrupt values*



## #2 Bad Model Validation

- **Solution**

- *Proprietary time-series algorithms*

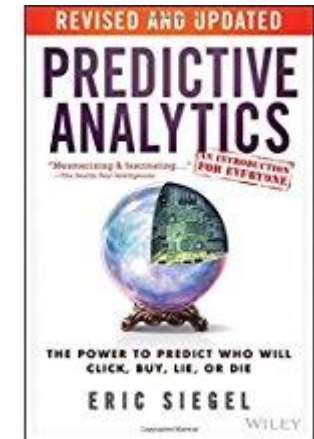
- **Failure**

- *None – error detected before deployment*

- **Why?**

- *Temporal leakage*

- *Time window mistakenly shifted by one day*

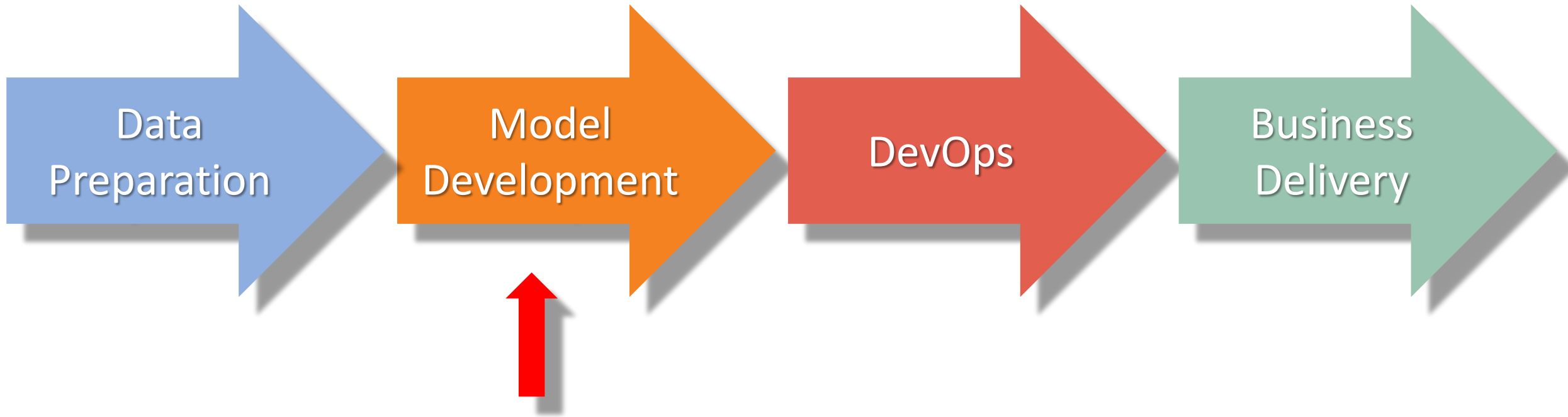


- **Challenge**

- *Concepts like target leakage are too complex for citizen data scientists*

- **Best Practices**

- *Trial model in non-actionable live environment*
  - *Have 3<sup>rd</sup> party reverse engineer model*
  - *Be careful!*



## #3 Regulatory Compliance

- **Solution**
  - *Target looks at market basket analysis*
- **Failure**
  - *Father notices teenage daughter is receiving coupons for baby items from Target*
  - *Uncomfortable conversation with daughter ensues ...*
- **Why?**
  - *Dealing with dangerous data*



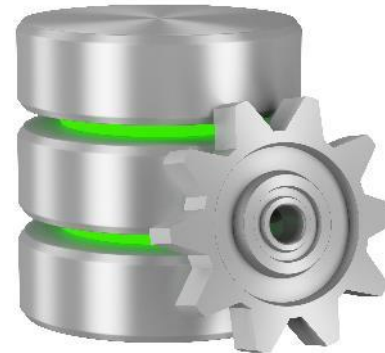
- **Challenge**
  - *Privacy violations are complex*
- **Best Practices**
  - *Agile teams with privacy advocate*
  - *Evaluation of risks vs. benefits*
  - *Build processes for CDO/CIO/CPO approval for things that could damage the corporate brand*
  - *Ability to rollback model creation environment*

## Data Lake / Data Warehouse



- Rollback of data
- Random sampling
- Flagging dangerous data
- Data lineage
- Data governance

## Predictive Analytics Engine

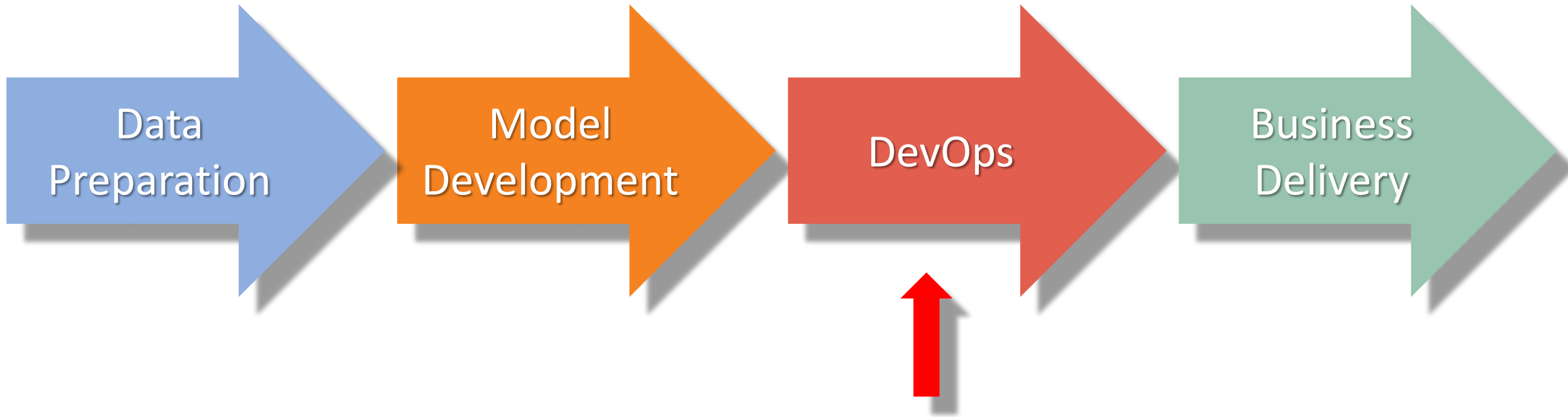


- Test and control creation
- Detect highly correlated variables
- Alerts on variable variation
- Model lineage
- Model governance
- Rollback to models
- Model audit trails
- Temporal leakage detection
- Alert on model aging

## Business User

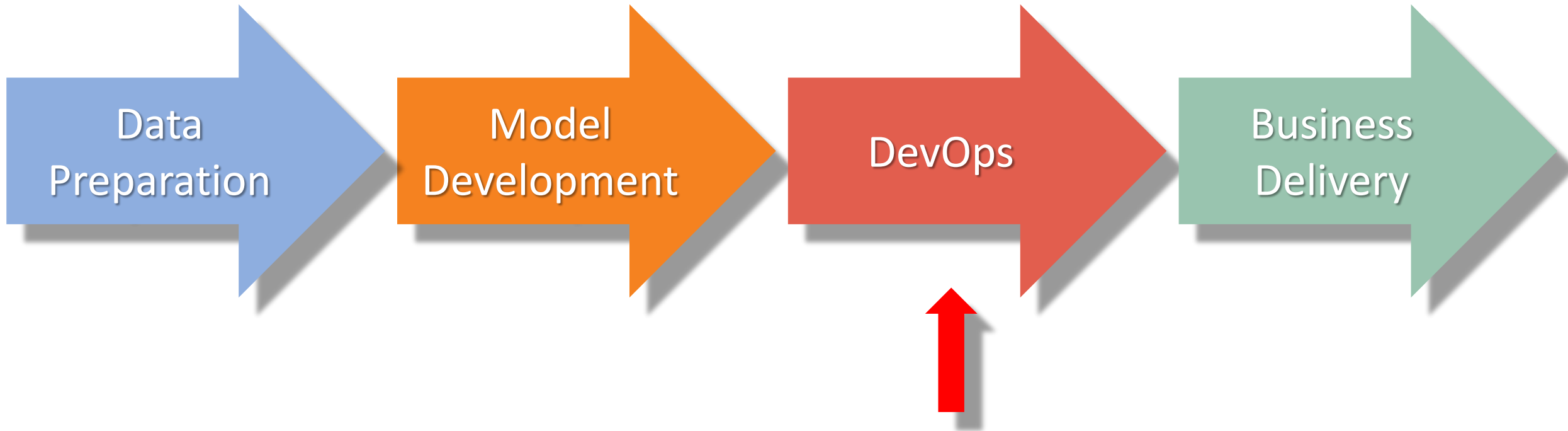


- Business modeling / ROI
- Workflow
- Collaboration
- Model marketplace
- Model ranking



## #4 Model Degradation

- **Challenge**
  - *Models decay in performance*
  - *Data scientists are overwhelmed*
- **Best Practices**
  - *Automate retraining or refreshing a model*
  - *Utilize tools that have alerting systems*
  - *Empower operations to offload data scientists*



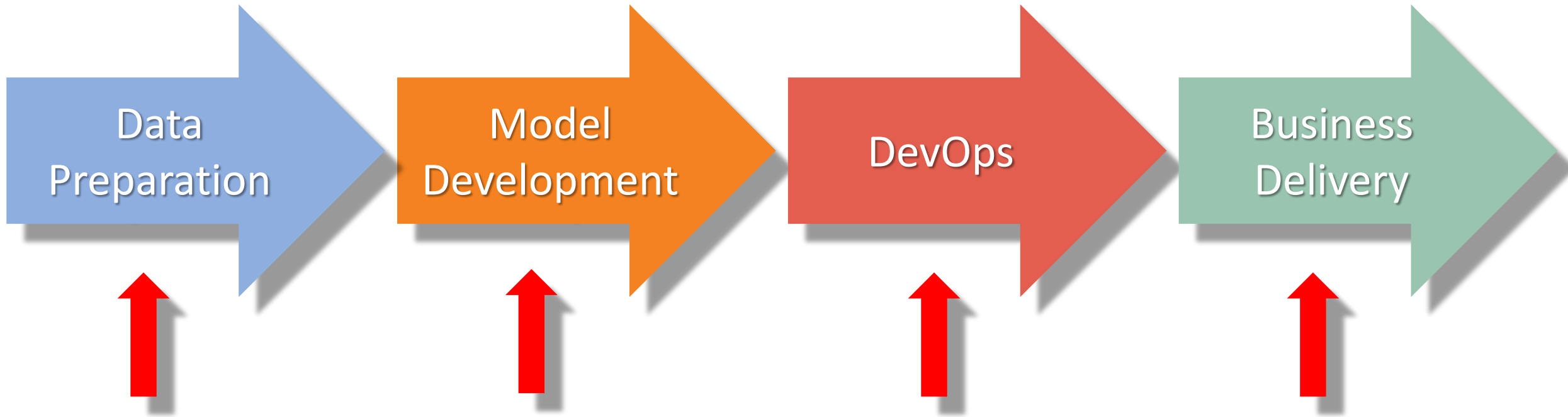
## #5 Deployment Disconnect

- **Challenge**

- *Model is created in one database and deployed elsewhere*
- *Must be recoded and tested*
- *Can delay deployment by 3-6 months*

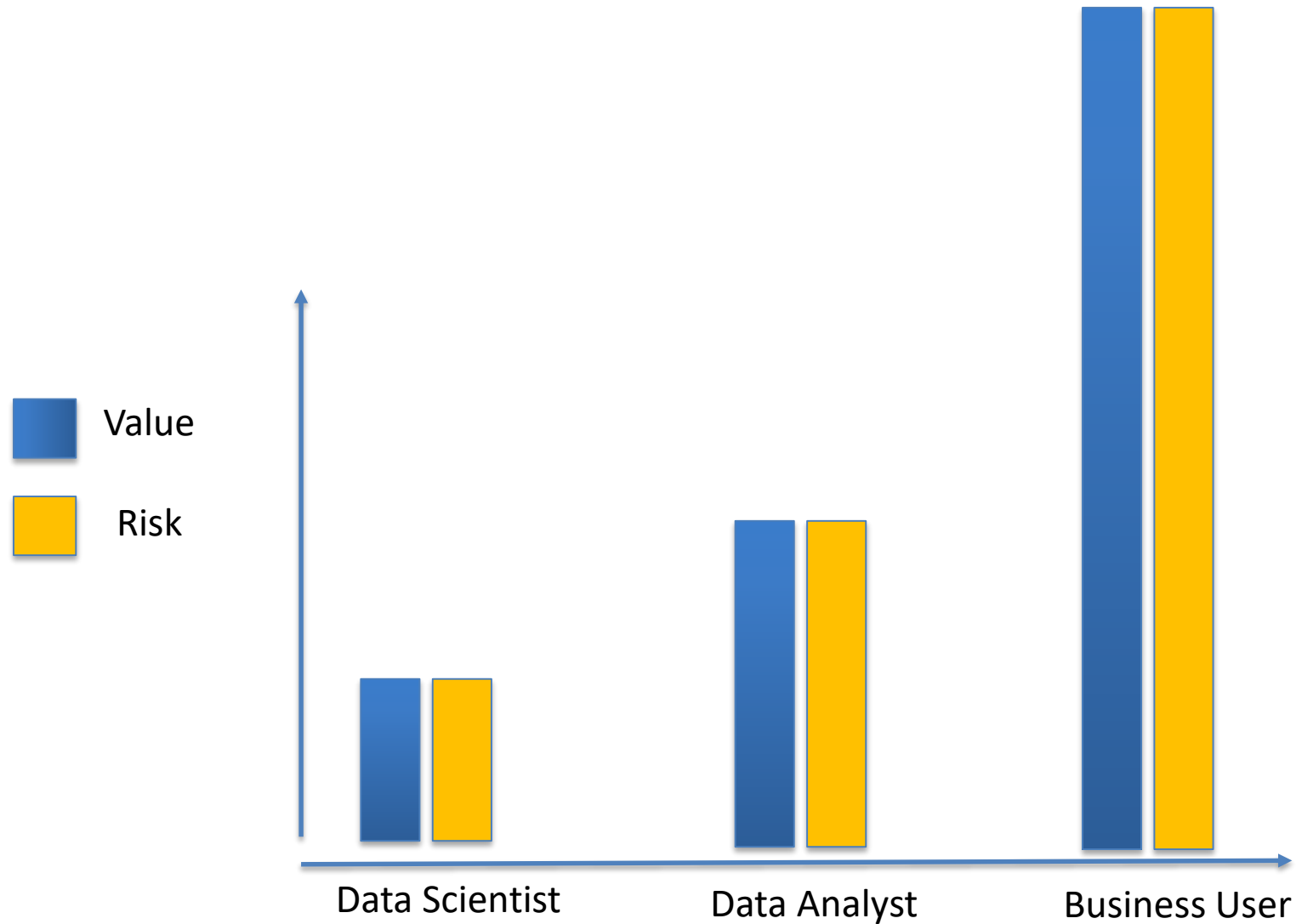
- **Best Practices**

- *Use tools that automatically generate code*
- *Build and deploy in the same database*



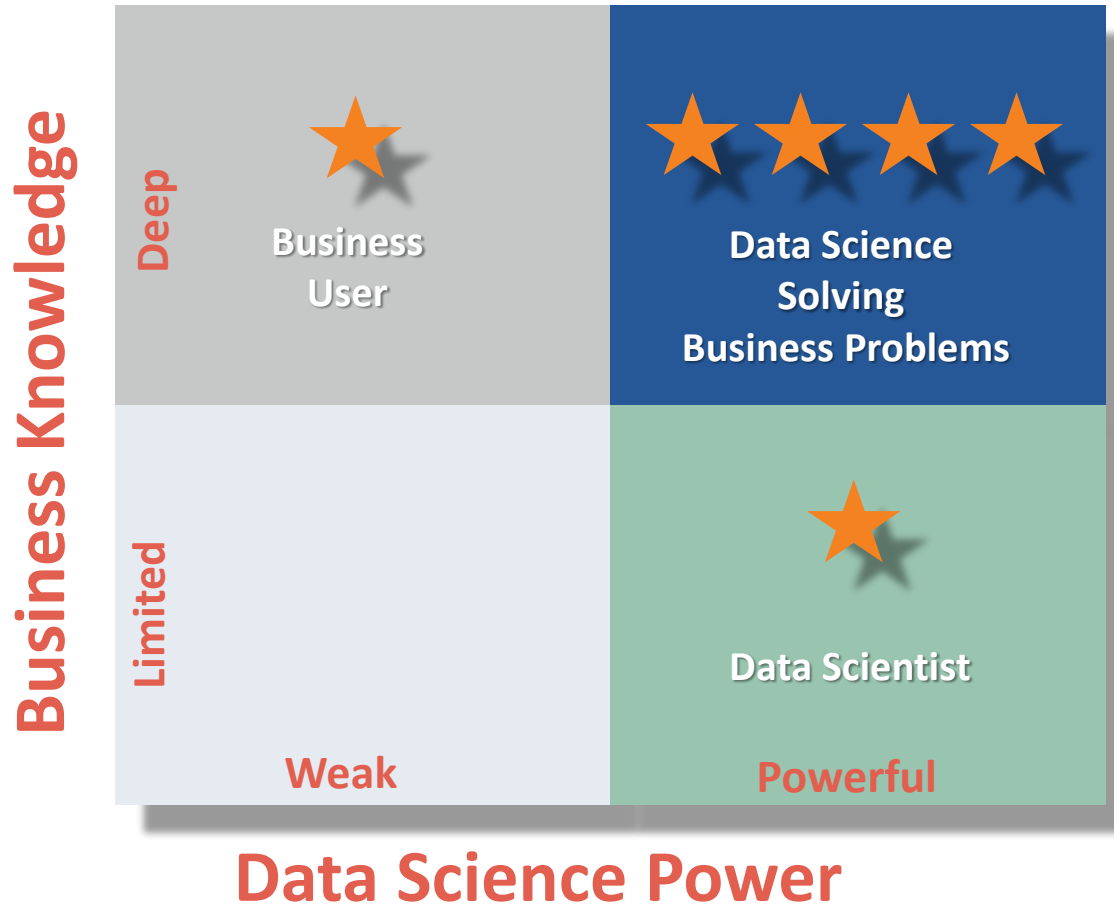
## #6 Finding Data Scientists

# Value and Risk Increase Together





# Acceleration = Business + DS

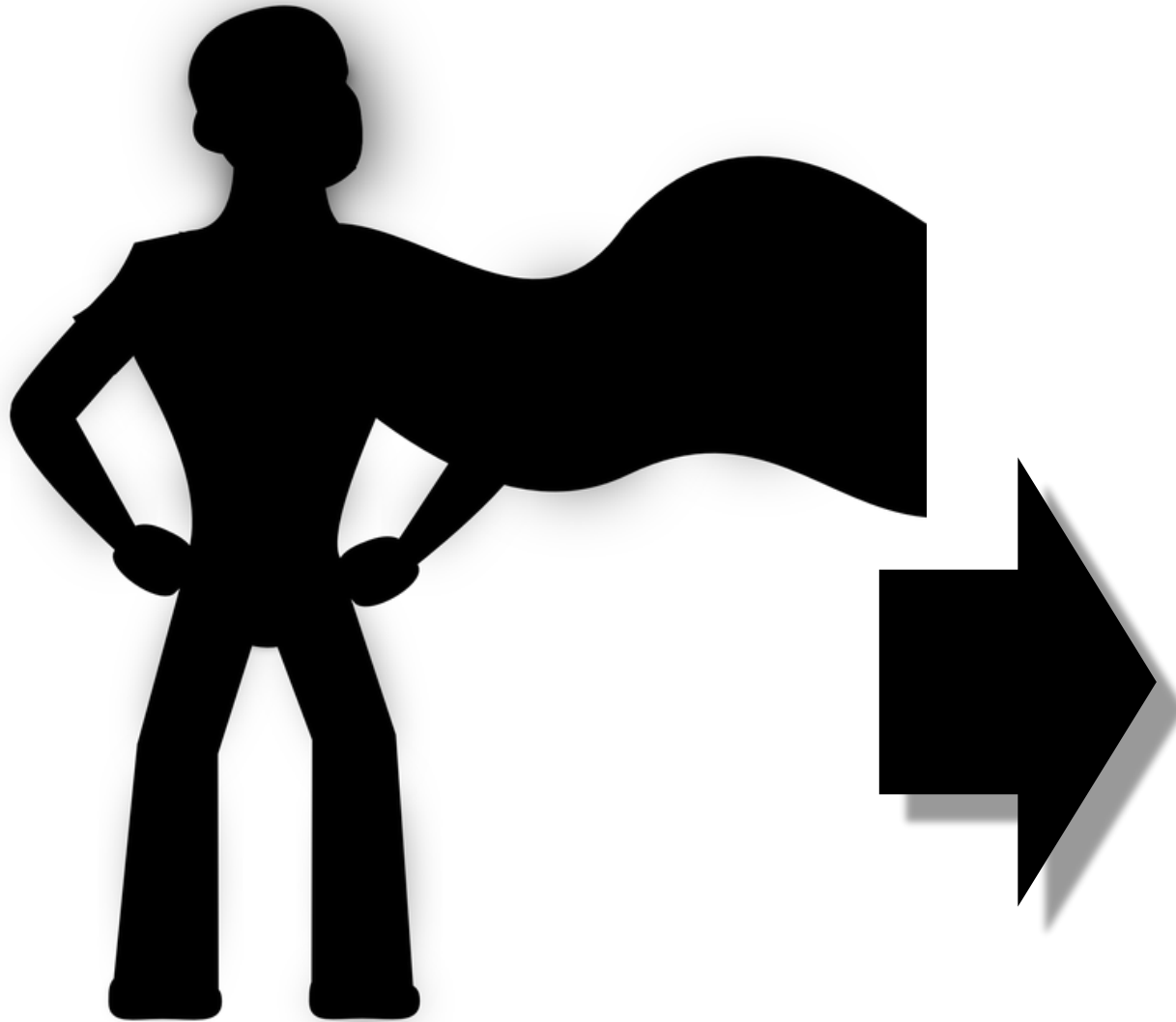




## Mythical Data Scientist

- Data expert
- Data engineer
- Deep learning expert
- Statistician
- Coder
- Cloud expert
- Operations expert
- Parallel processing developer
- Deep thinker
- Good communicator
- Understands business problem
- Business rules expert

## Required Skills



**Mythical Superman Data Scientist**



**Multidisciplinary Agile Team**

## 1. Data Steward

- Knows what data is available, lineage, and governance

## 2. Data Engineer

- Can access and transform data and optimize database performance

## 3. Privacy Advocate

## 4. Data Scientist

- Understands predictive analytics, statistics, machine learning
- Provides high accuracy models and scores

## 5. Operational Engineer

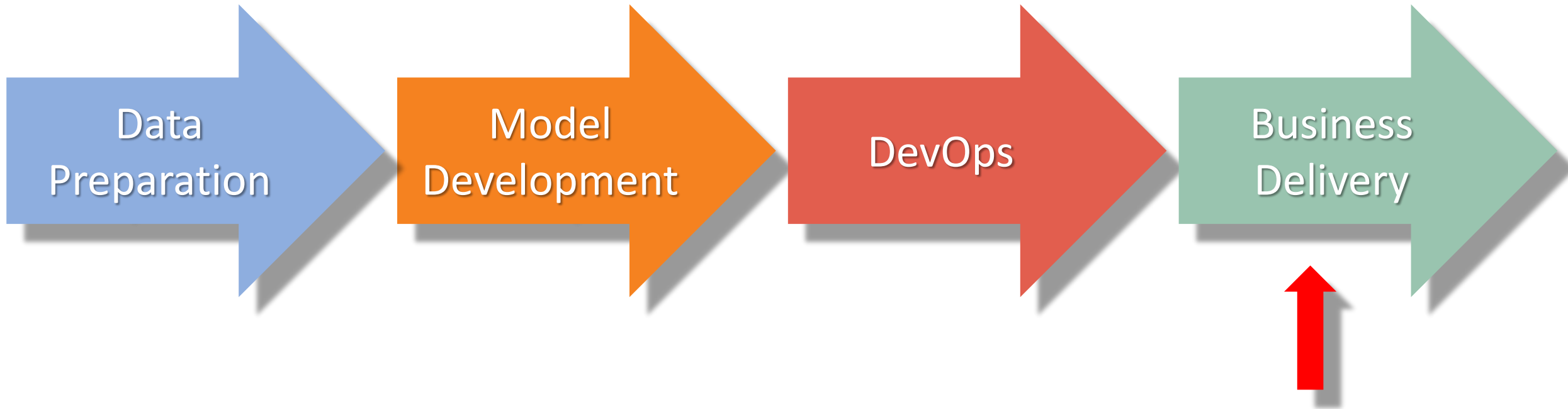
## 6. Data Analyst

- Can translate between business requests and data realities
- Provides ad hoc query and report support

## 7. Business User

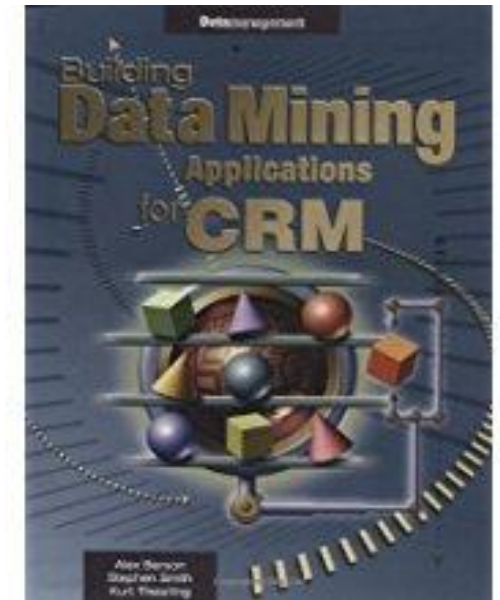
- Understands the business problem and responsible for ROI
- Provides business case and logical validation of causality

- **Challenge**
  - *Heroic project-based model creation*
  - *High error rates*
  - *Can't scale*
  
- **Best Practices**
  - *Agile teams*
  - *Embed DS representative into business*



# #7 Business Disconnect

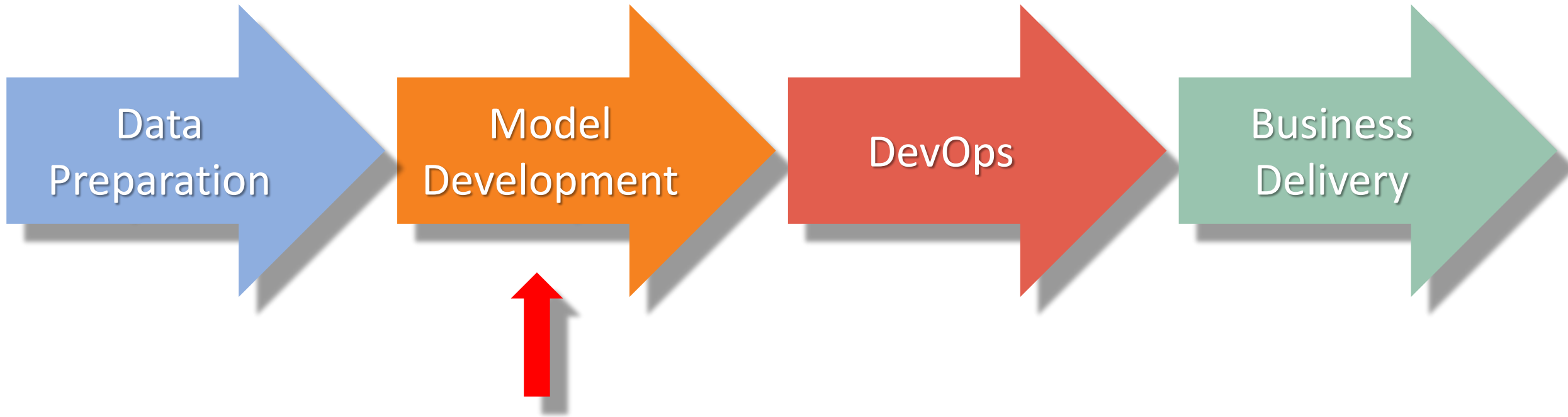
- **Solution**
  - *CART model in specialized PA tool*
- **Failure**
  - *Created higher churn than control group*
- **Why?**
  - *Predictive model excellent*
  - *BUT - Offer reminded customers of their anniversary*



Great book!!!  
(completely biased)

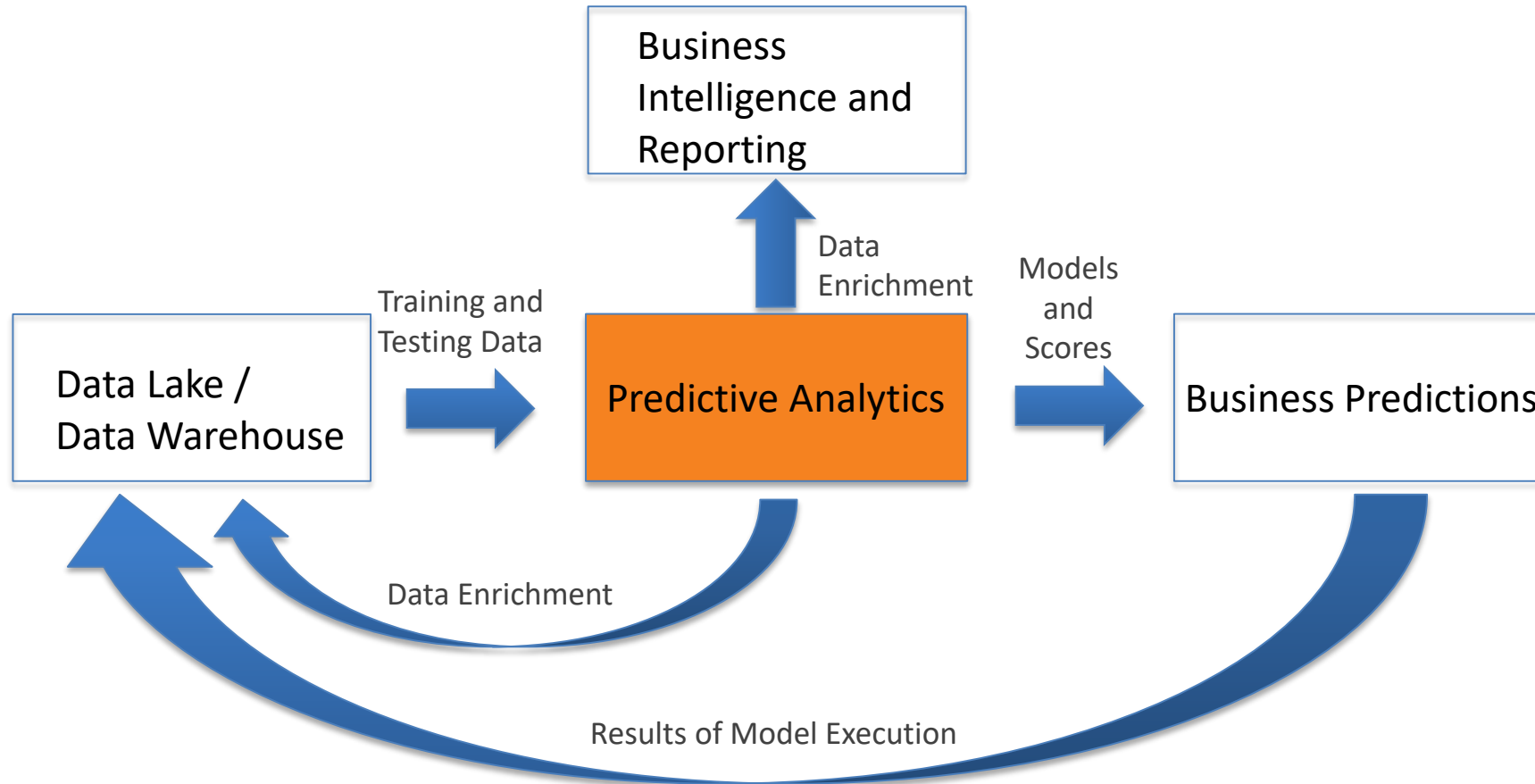
- **Challenge**
  - *Possibility of large errors*
  - *Gun shy business users*
  - *Frustrated data scientists*
- **Best Practices**
  - *Predict offer effect with small trials*
  - *Include business user in top down design*





## #8 Lack of a Data Science Platform

- **Challenge**
  - *High error rates*
  - *Long delivery cycles*
  - *Can't scale*
- **Best Practices**
  - *Centralize and focus on building a 'platform'*
  - *Deliver like an internal software company*
  - *Plan for scale but execute for success today*



1. Bad Data
2. Bad Validation
3. Regulatory Compliance
4. Model Degradation
5. Deployment Disconnect
6. Finding Data Scientists
7. Business Disconnect
8. No Platform



## Path to Digital Transformation

1. **Self-service data science** will be the dominant lifeform
2. **Citizen data scientists** will be effective
3. Data scientists will be **more effective**
4. Data science best practices will look a lot like **computer science best practices**
5. Operationalized data science will provide a **strategic competitive advantage**

## Questions?

### Stephen Smith

Research Director, Data Science

Eckerson Group

Email: [ssmith@Eckerson.com](mailto:ssmith@Eckerson.com)

Twitter: @steve4years



### Karen Fegarty

Regional VP, Mariner Innovations

T/ 902-499-4983

[Karen.fegarty@marinerinnovations.com](mailto:Karen.fegarty@marinerinnovations.com)

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

## Resources

### Papers on [www.Eckerson.com](http://www.Eckerson.com)

- **Best Practices in Data Science : Ten Keys**
- **The Demise of the Data Warehouse**
- **Eckerson Eight Innovations in Data Science**
- **Data Science is Plutonium**

### Books

- **“Building Data Mining Applications for CRM”**  
– Stephen Smith, Alex Berson, Kurt Thearling
- **“Data Warehousing, Data Mining, & OLAP”**  
– Stephen Smith, Alex Berson
- **“Predictive Analytics”**  
- Eric Siegel
- **“Data Science for Business”**  
– Foster Provost, Tom Fawcett

Sponsored by:

