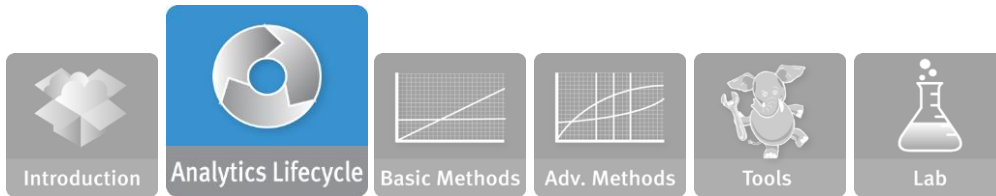




Module 2 – Data Analytics Lifecycle

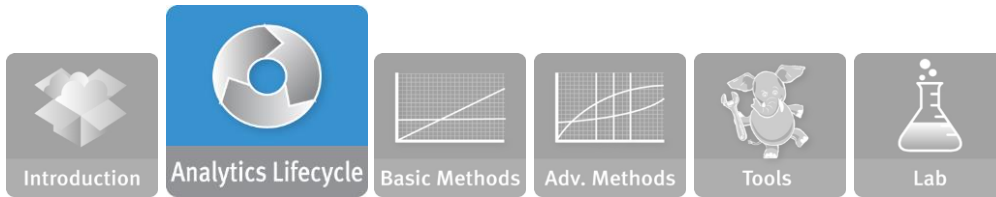
EMC² PROVEN PROFESSIONAL



Module 2: Data Analytics Lifecycle

Upon completion of this module, you should be able to:

- Apply the Data Analytics Lifecycle to a case study scenario
- Frame a business problem as an analytics problem
- Identify the four main deliverables in an analytics project



Module 2: Data Analytics Lifecycle

During this module the following topics are covered:

- Data Analytics Lifecycle
- Roles for a Successful Analytics Project
- Case Study to apply the data analytics lifecycle

How to Approach Your Analytics Problems



Your Thoughts?

- How do you currently approach your analytics problems?
- Do you follow a methodology or some kind of framework?
- How do you plan for an analytic project?



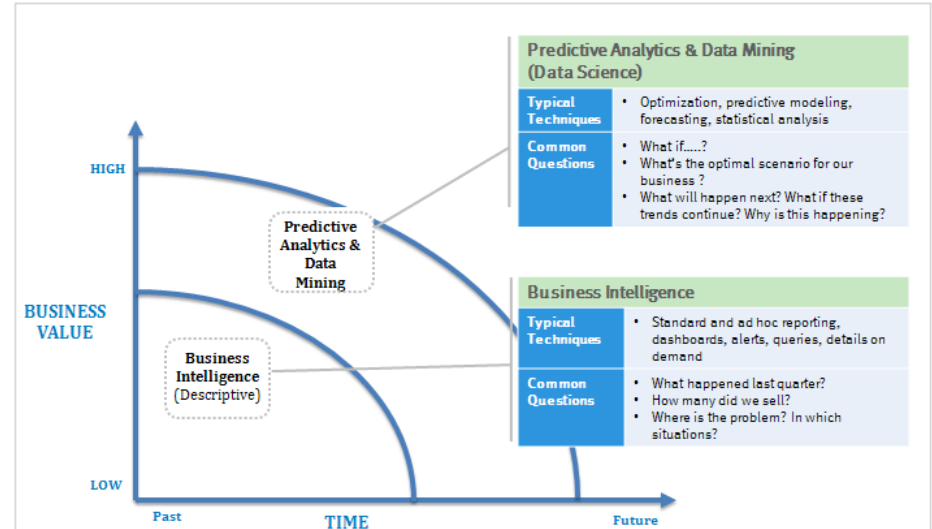
Value of Using the Data Analytics Lifecycle

- Focus your time
- Ensure rigor and completeness
- Enable better transition to members of the cross-functional analytic teams
 - ▶ Repeatable
 - ▶ Scale to additional analysts
 - ▶ Support validity of findings

“A journey of a thousand miles begins with a single step” (Lao Tzu)

Need For a Process to Guide Data Science Projects

1. Well-defined processes can help guide any analytic project
2. Focus of Data Analytics Lifecycle is on Data Science projects, not business intelligence



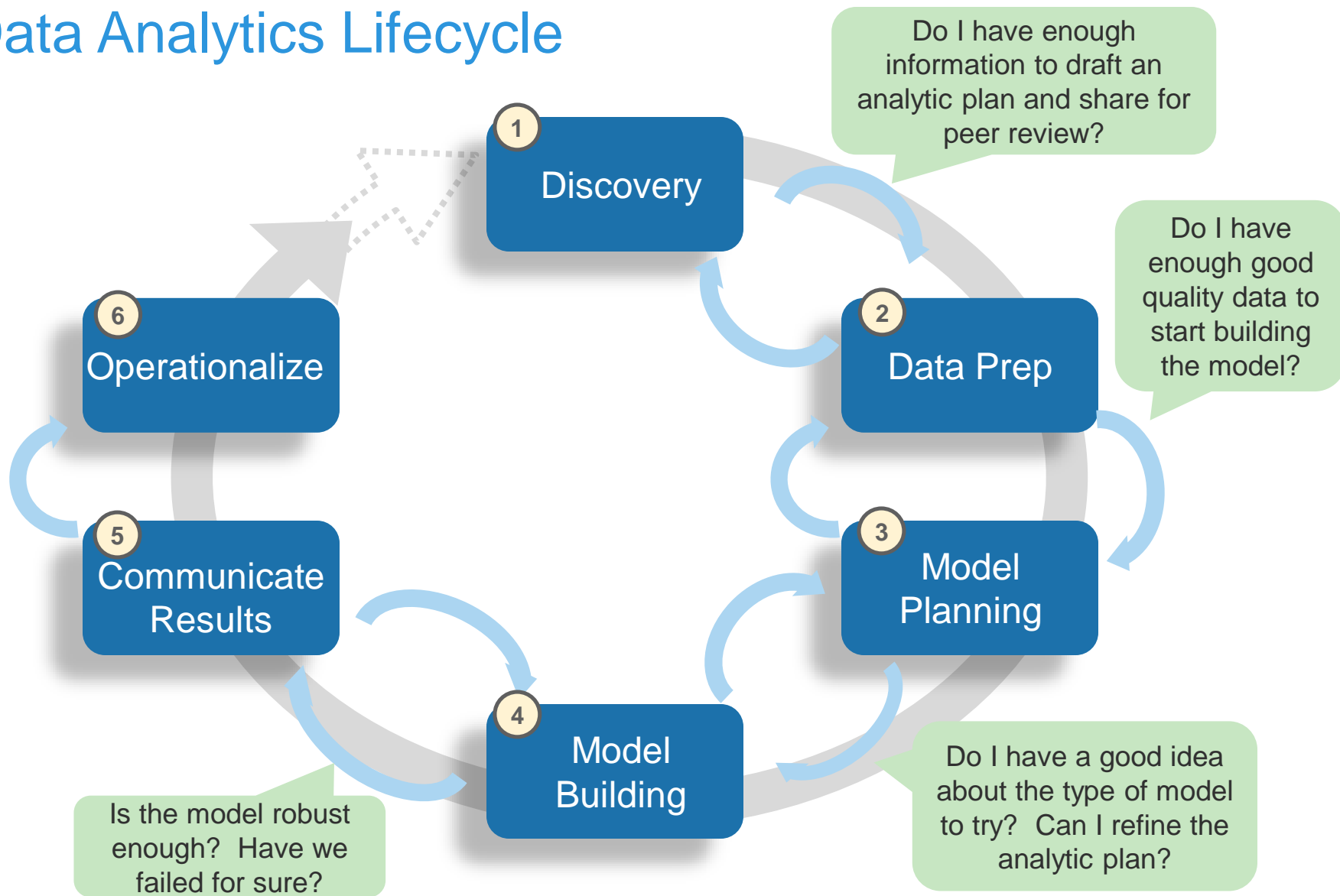
3. Data Science projects tend to require a more consultative approach, and differ in a few ways
 - ▶ More due diligence in Discovery phase
 - ▶ More projects which lack shape or structure
 - ▶ Less predictable data

Key Roles for a Successful Analytic Project



Role	Description
Business User	Someone who benefits from the end results and can consult and advise project team on value of end results and how these will be operationalized
Project Sponsor	Person responsible for the genesis of the project, providing the impetus for the project and core business problem, generally provides the funding and will gauge the degree of value from the final outputs of the working team
Project Manager	Ensure key milestones and objectives are met on time and at expected quality.
Business Intelligence Analyst	Business domain expertise with deep understanding of the data, KPIs, key metrics and business intelligence from a reporting perspective
Data Engineer	Deep technical skills to assist with tuning SQL queries for data management, extraction and support data ingest to analytic sandbox
Database Administrator (DBA)	Database Administrator who provisions and configures database environment to support the analytical needs of the working team
Data Scientist	Provide subject matter expertise for analytical techniques, data modeling, applying valid analytical techniques to given business problems and ensuring overall analytical objectives are met

Data Analytics Lifecycle



Data Analytics Lifecycle

Phase 1: Discovery



1

Discovery

Do I have enough information to draft an analytic plan and share for peer review?

Do I have enough good

- **Learn the Business Domain**

- ▶ Determine amount of domain knowledge needed to orient you to the data and interpret results downstream
- ▶ Determine the general analytic problem type (such as clustering, classification)
- ▶ If you don't know, then conduct initial research to learn about the domain area you'll be analyzing

- **Learn from the past**

- ▶ Have there been previous attempts in the organization to solve this problem?
- ▶ If so, why did they fail? Why are we trying again? How have things changed?

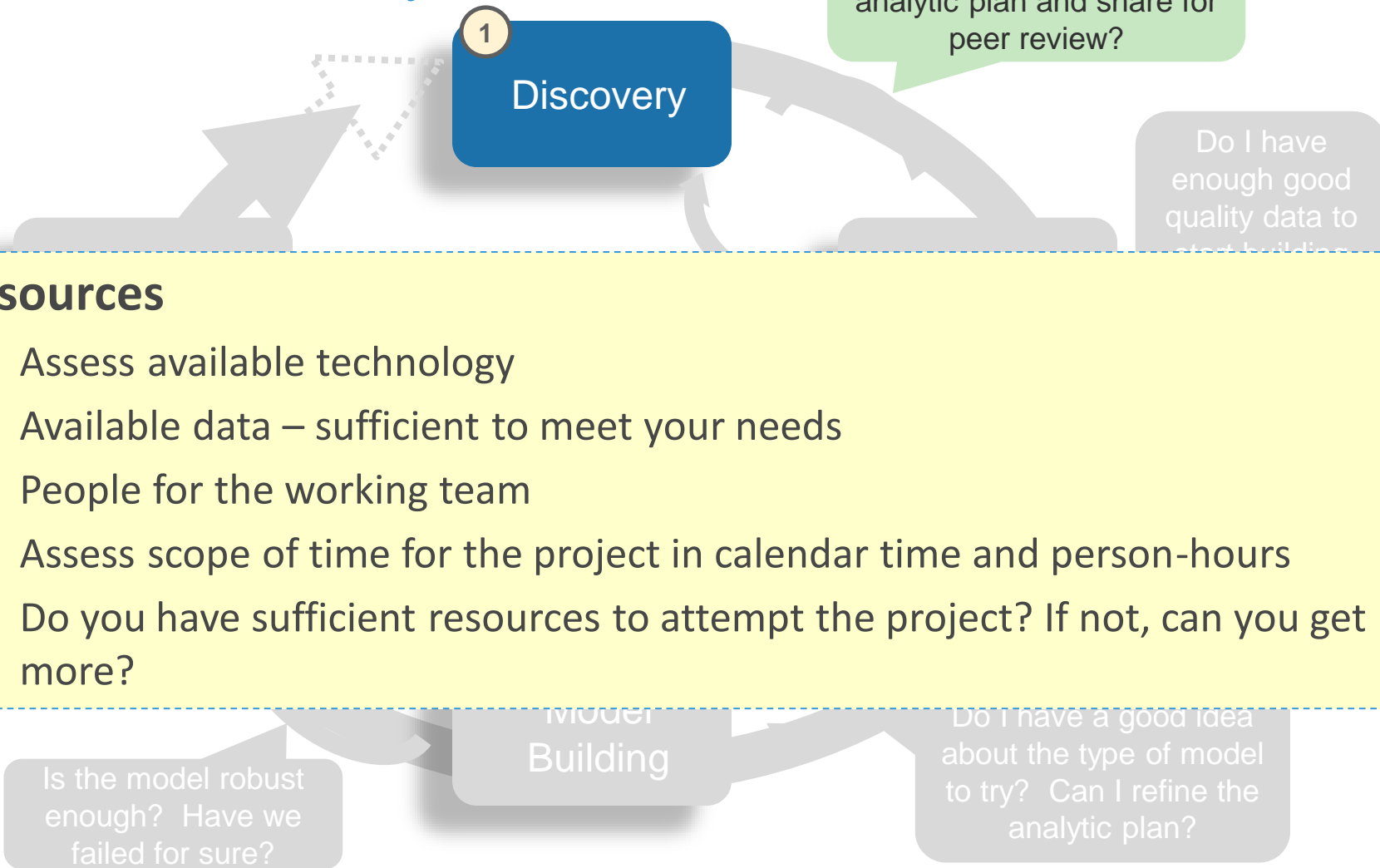
Building

Is the model robust enough? Have we failed for sure?

about the type of model to try? Can I refine the analytic plan?

Data Analytics Lifecycle

Phase 1: Discovery



Data Analytics Lifecycle

Phase 1: Discovery



- **Frame the problem.....***Framing is the process of stating the analytics problem to be solved*
 - ▶ *State the analytics problem*, why it is important, and to whom
 - ▶ Identify key stakeholders and their interests in the project
 - ▶ Clearly articulate the current situation and ***pain points***
 - ▶ Objectives – identify what needs to be achieved in business terms and what needs to be done to meet the needs
 - ▶ What is the goal? What are the criteria for success? What’s “good enough”?
 - ▶ What is the failure criterion (when do we just stop trying or settle for what we have)?
 - ▶ Identify the success criteria, key risks, and stakeholders (such as RACI)

failed for sure?

...



Tips for Interviewing the Analytics Sponsor

- Even if you are “given” an analytic problem you should work with clients to clarify and frame the problem
 - ▶ You’re typically handed solutions, you need to identify the problem and their desired outcome



Sponsor Interview Tips

- Prepare for the interview – draft your questions, review with colleague, team
- Use open-ended questions, don’t ask leading questions
- Probe for details, follow-up
- Don’t fill every silence – give them time to think
- Let them express their ideas, don’t put words in their mouth, let them share their feelings
- Ask clarifying questions, ask why – is that correct? Am I on target? Is there anything else?
- Use active listening – repeat it back to make sure you heard it correctly
- Don’t express your opinions
- Be mindful of your body language and theirs – use eye contact, be attentive
- Minimize distractions
- Document what you heard and review it back with the sponsor



Tips for Interviewing the Analytics Sponsor

Interview Questions

- What is the business problem you're trying to solve?
- What is your desired outcome?
- Will the focus and scope of the problem change if the following dimensions change:
 - Time – analyzing 1 year or 10 years worth of data?
 - People – how would this project change this?
 - Risk – conservative to aggressive
 - Resources – none to unlimited (tools, tech,)
 - Size and attributes of Data
- What data sources do you have?
- What industry issues may impact the analysis?
- What timelines are you up against?
- Who could provide insight into the project? Consulted?
- Who has final say on the project?



Data Analytics Lifecycle

Phase 1: Discovery



1

Discovery

Do I have enough information to draft an analytic plan and share for peer review?

- **Formulate Initial Hypotheses**

- ▶ $IH, H_1, H_2, H_3, \dots H_n$
- ▶ Gather and assess hypotheses from stakeholders and domain experts
- ▶ Preliminary data exploration to inform discussions with stakeholders during the hypothesis forming stage

- **Identify Data Sources – Begin Learning the Data**

- ▶ Aggregate sources for previewing the data and provide high-level understanding
- ▶ Review the raw data
- ▶ Determine the structures and tools needed
- ▶ Scope the kind of data needed for this kind of problem

Do I have enough good quality data to start building the model?



Good idea
of model
refine the
plan?

Using a Sample Case Study to Track the Phases in the Data Analytics Lifecycle



Mini Case Study: Churn Prediction for Yoyodyne Bank

Situation Synopsis

- Retail Bank, Yoyodyne Bank wants to improve the Net Present Value (NPV) and retention rate of customers
- They want to establish an effective marketing campaign targeting customers to reduce the churn rate by at least five percent
- The bank wants to determine whether those customers are worth retaining. In addition, the bank also wants to analyze reasons for customer attrition and what they can do to keep them
- The bank wants to build a data warehouse to support Marketing and other related customer care groups

How to Frame an Analytics Problem

Mini Case Study



Sample <i>Business</i> Problems	Qualifiers	Analytical Approach
<ul style="list-style-type: none"> How can we improve on x? What's happening real-time? Trends? How can we use analytics differentiate ourselves How can we use analytics to innovate? How can we stay ahead of our biggest competitor? 	<p>Will the focus and scope of the problem change if the following dimensions change:</p> <ul style="list-style-type: none"> Time People – how would x change this? Risk – conservative/aggressive Resources – none/unlimited Size of Data? 	<p>Define an analytical approach, including key terms, metrics, and data needed.</p>
<p>Mini Case Study: Churn Prediction for Yoyodyne Bank</p> <p><u>Yoyodyne Bank</u> How can we improve Net Present Value (NPV) and retention rate of the customers?</p>	<ul style="list-style-type: none"> Time: Trailing 5 months People: Working team and business users from the Bank Risk: the project will fail if we cannot determine valid predictors of churn Resources: EDW, analytic sandbox, OLTP system Data: Use 24 months for the training set, then analyze 5 months of historical data for those customers who churned 	<p>How do we identify churn/no churn for a customer?</p> <p>Pilot study followed full scale analytical model</p>

Data Analytics Lifecycle

Phase 2: Data Preparation



- **Prepare Analytic Sandbox**
 - ▶ Work space for the analytic team
 - ▶ 10x+ vs. EDW
- **Perform ELT**
 - ▶ Determine needed transformations
 - ▶ Assess data quality and structuring
 - ▶ Derive statistically useful measures
 - ▶ Extract data and determine data connections for raw data, OLTP transactions, OLAP cubes or data feeds
 - ▶ Big ELT and Big ETL

Do I have enough information to draft an analytic plan and share for peer review?

2
Data Prep

Do I have enough good quality data to start building the model?

Model Planning

Do I have a good idea about the type of model

- **Useful Tools for this phase:**
 - ***For Data Transformation & Cleansing:*** SQL, Hadoop, MapReduce, Alpine Miner

Data Analytics Lifecycle

Phase 2: Data Preparation



- **Familiarize yourself with the data thoroughly**

- ▶ List your data sources
- ▶ What's needed vs. what's available

- **Data Conditioning**

- ▶ Clean and normalize data
- ▶ Discern what you keep vs. what you discard

- **Survey & Visualize**

- ▶ Overview, zoom & filter, details-on-demand
- ▶ Descriptive Statistics
- ▶ Data Quality

- **Useful Tools for this phase:**

- Descriptive Statistics on candidate variables for diagnostics & quality
- **Visualization:** R (base package, ggplot and lattice), GnuPlot, Ggobi/Rggobi, Spotfire, Tableau

Do I have enough information to draft an analytic plan and share for peer review?

2

Data Prep

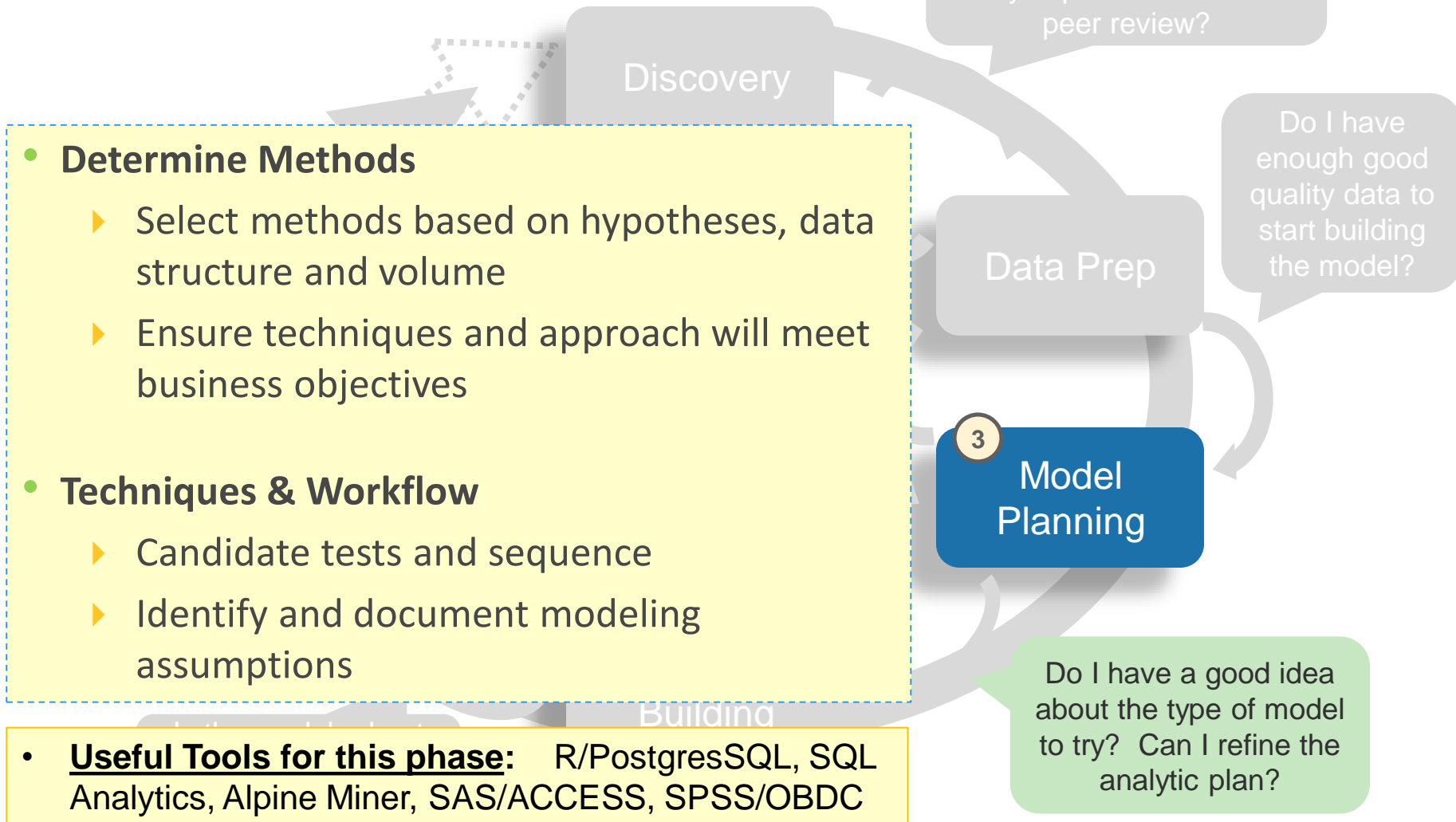
Do I have enough good quality data to start building the model?

Model Planning

Do I have a good idea about the type of model?

Data Analytics Lifecycle

Phase 3: Model Planning

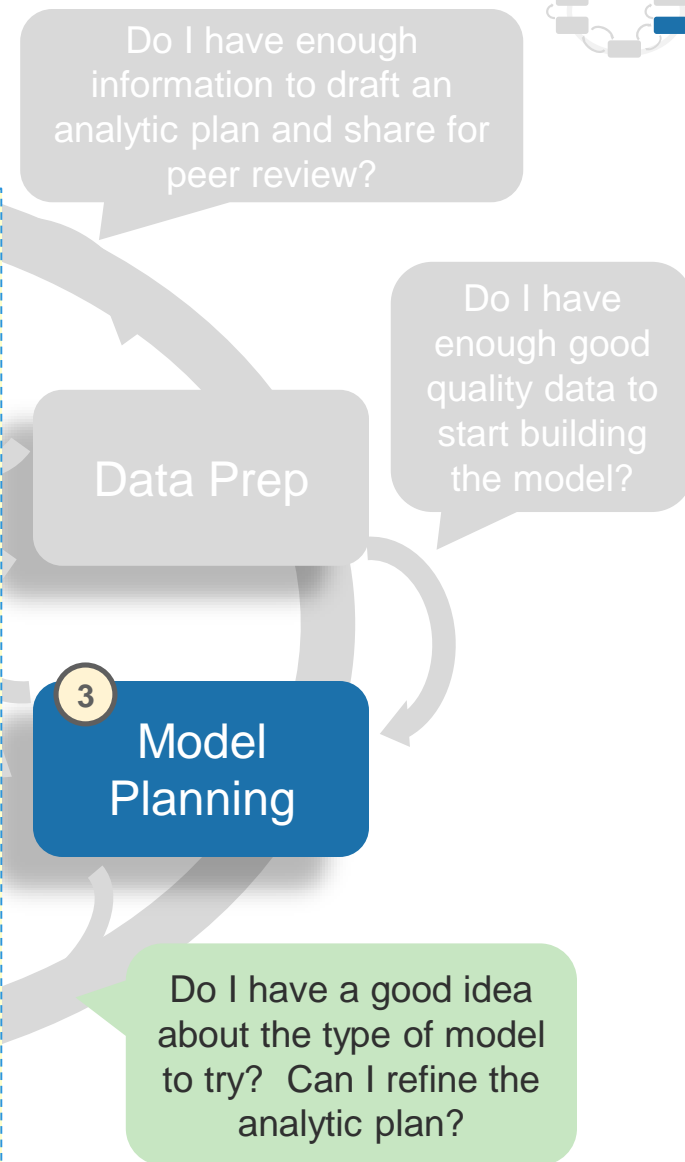


Data Analytics Lifecycle

Phase 3: Model Planning



- **Data Exploration**
- **Variable Selection**
 - ▶ Inputs from stakeholders and domain experts
 - ▶ Capture essence of the predictors, leverage a technique for dimensionality reduction
 - ▶ Iterative testing to confirm the most significant variables
- **Model Selection**
 - ▶ Conversion to SQL or database language for best performance
 - ▶ Choose technique based on the end goal





Sample Research: Churn Prediction in Other Verticals

*Mini Case Study:
Churn Prediction for
Yoyodyne Bank*

- After conducting research on churn prediction, you have identified many methods for analyzing customer churn across multiple verticals (those in **bold** are taught in this course)
- At this point, a Data Scientist would assess the methods and select the best model for the situation

Market Sector	Analytic Techniques/Methods Used
Wireless Telecom	DMEL method (data mining by evolutionary learning)
Retail Business	Logistic regression , ARD (automatic relevance determination), decision tree
Daily Grocery	MLR (multiple linear regression), ARD, and decision tree
Wireless Telecom	Neural network, decision tree , hierarchical neurofuzzy systems, rule evolver
Retail Banking	Multiple regression
Wireless Telecom	Logistic regression , neural network, decision tree

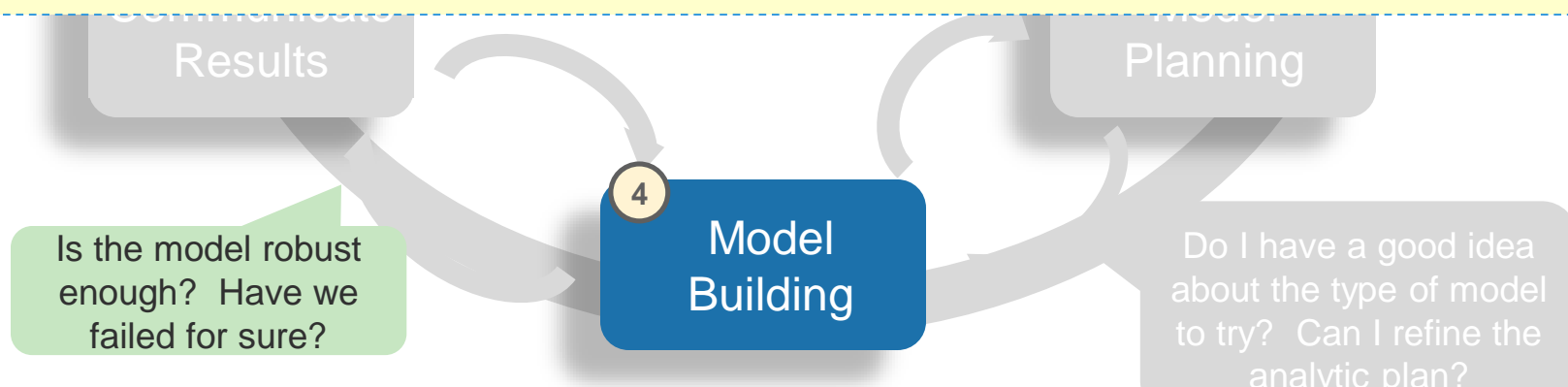
Data Analytics Lifecycle

Phase 4: Model Building



Do I have enough information to draft an analytic plan and share for peer review?

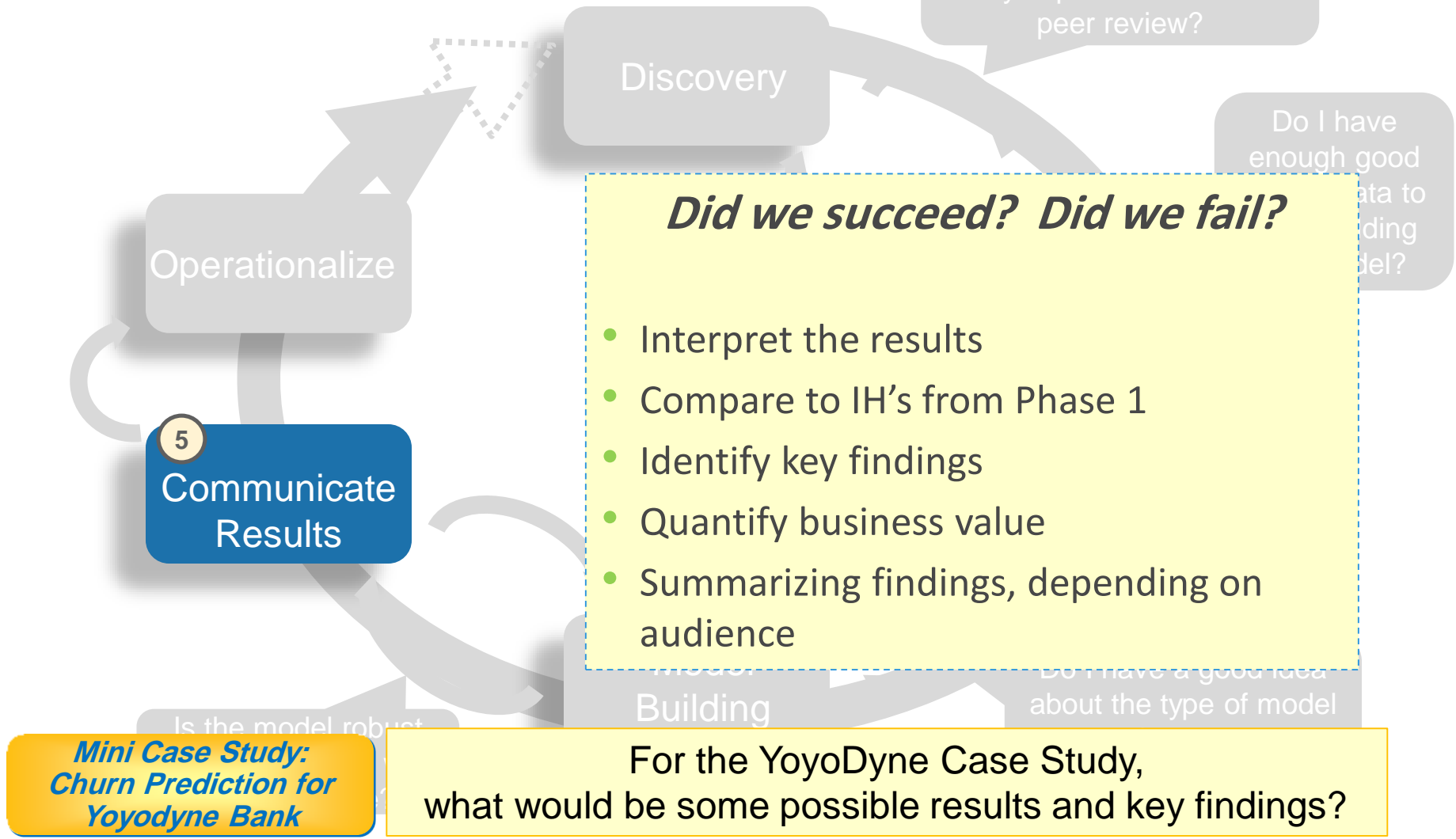
- **Develop data sets for testing, training, and production purposes**
 - ▶ Need to ensure that the model data is sufficiently robust for the model and analytical techniques
 - ▶ Smaller, test sets for validating approach, training set for initial experiments
- **Get the best environment you can for building models and workflows...fast hardware, parallel processing**



- **Useful Tools for this phase:** R, PL/R, SQL, Alpine Miner, SAS Enterprise Miner

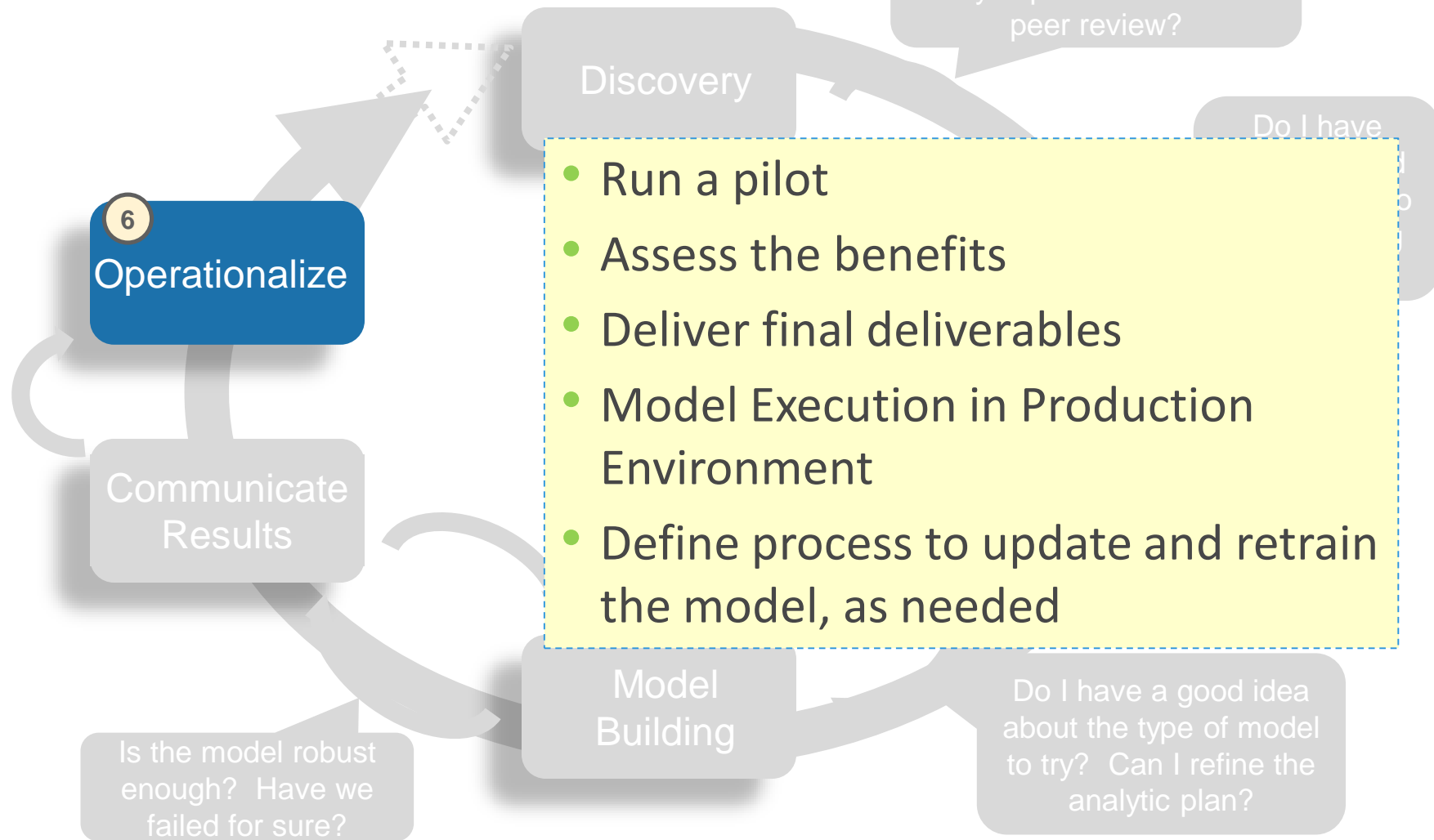
Data Analytics Lifecycle

Phase 5: Communicate Results



Data Analytics Lifecycle

Phase 6: Operationalize



Analytic Plan

Mini Case Study: Churn Prediction for Retail Banking



Components of Analytic Plan	Retail Banking: Yoyodyne Bank
Phase 1: Discovery Business Problem Framed	How do we identify churn/no churn for a customer?
Initial Hypotheses	Transaction volume and type are key predictors of churn rates.
Data	5 months of customer account history.
Phase 3: Model Planning - Analytic Technique	Logistic regression to identify most influential factors predicting churn.
Phase 5: Result & Key Findings	Once customers stop using their accounts for gas and groceries, they will soon erode their accounts and churn. If customers use their debit card fewer than 5 times per month, they will leave the bank within 60 days.
Business Impact	If we can target customers who are high-risk for churn, we can reduce customer attrition by 25%. This would save \$3 million in lost of customer revenue and avoid \$1.5 million in new customer acquisition costs each year.

Key Outputs from a Successful Analytic Project, by Role



Role	Description	What the Role Needs in the Final Deliverables
Business User	Someone who benefits from the end results and can consult and advise project team on value of end results and how these will be operationalized	<ul style="list-style-type: none"> • Sponsor Presentation addressing: <ul style="list-style-type: none"> • Are the results good for me? • What are the benefits of the findings? • What are the implications of this for me?
Project Sponsor	Person responsible for the genesis of the project, providing the impetus for the project and core business problem, generally provides the funding and will gauge the degree of value from the final outputs of the working team	<ul style="list-style-type: none"> • Sponsor Presentation addressing: <ul style="list-style-type: none"> • What's the business impact of doing this? • What are the risks? ROI? • How can this be evangelized within the organization (and beyond)?
Project Manager	Ensure key milestones and objectives are met on time and at expected quality.	
Business Intelligence Analyst	Business domain expertise with deep understanding of the data, KPIs, key metrics and business intelligence from a reporting perspective	<ul style="list-style-type: none"> • Show the analyst presentation • Determine if the reports will change
Data Engineer	Deep technical skills to assist with tuning SQL queries for data management, extraction and support data ingest to analytic sandbox	<ul style="list-style-type: none"> • Share the code from the analytical project • Create technical document on how to implement it.
Database Administrator (DBA)	Database Administrator who provisions and configures database environment to support the analytical needs of the working team	<ul style="list-style-type: none"> • Share the code from the analytical project • Create technical document on how to implement it.
Data Scientist	Provide subject matter expertise for analytical techniques, data modeling, applying valid analytical techniques to given business problems and ensuring overall analytical objectives are met	<ul style="list-style-type: none"> • Show the analyst presentation • Share the code

4 Core Deliverables to Meet Most Stakeholder Needs



1. Presentation for Project Sponsors

- “Big picture” takeaways for executive level stakeholders
- Determine key messages to aid their decision-making process
- Focus on clean, easy visuals for the presenter to explain and for the viewer to grasp

2. Presentation for Analysts

- Business process changes
- Reporting changes
- Fellow Data Scientists will want the details and are comfortable with technical graphs (such as ROC curves, density plots, histograms)

3. Code for technical people

4. Technical specs of implementing the code

Analyst Wish List for a Successful Analytics Project



Data & Workspaces

- Access to all the data, including aggregated OLAP data, BI tools, raw data, structured and various states of unstructured data as needed
- Up-to-date data dictionary to describe the data
- Area for staging and production data sets
- Ability to move data back and forth between workspaces and staging areas
- Analytic sandbox with strong compute power to experiment and play with the data

Tools

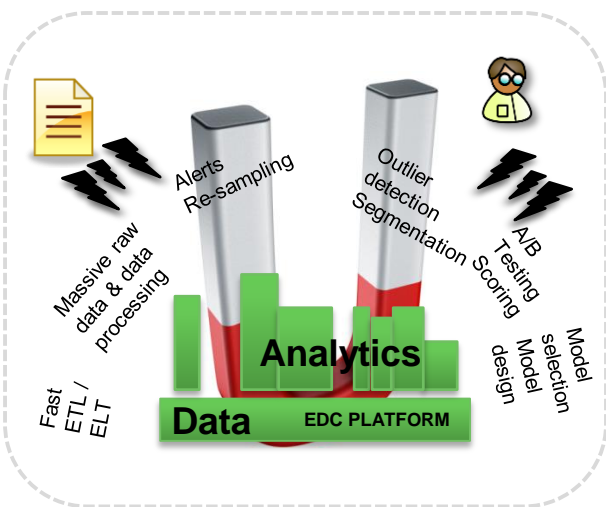
- Statistical/mathematical/visual software of choice for a given situation and problem set, such as SAS, Matlab, R, java tools, Tableau, Spotfire
- Collaboration: an online platform or environment for collaboration and communicating with team members
- Tool or place to log errors with systems, environments or data sets

Concepts in Practice

Greenplum's Approach to Analytics

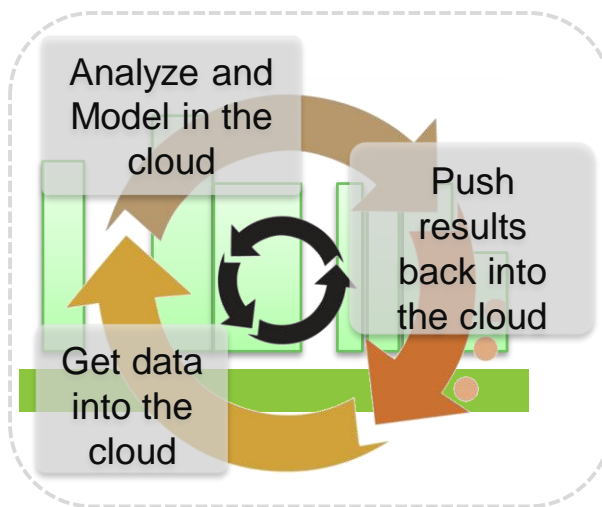
Magnetic

Attract all kinds of data



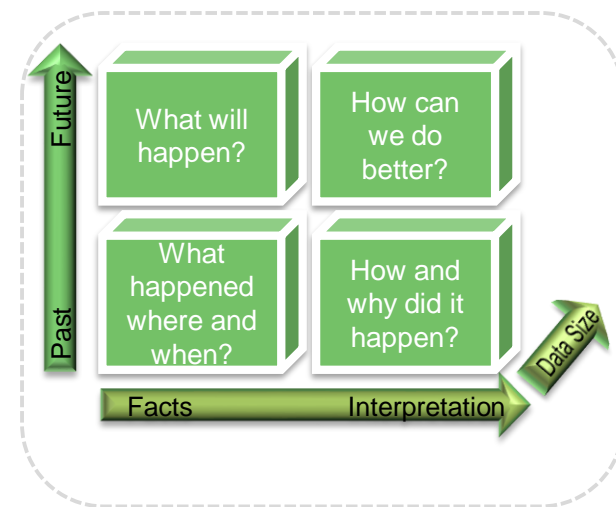
Agile

Flexible and elastic data structures



Deep

Rich data repository and algorithmic engine



Source: MAD Skills: New Analysis Practices for Big Data, March 2009



“The pessimist –
complains about the wind

The optimist –
expects it to change

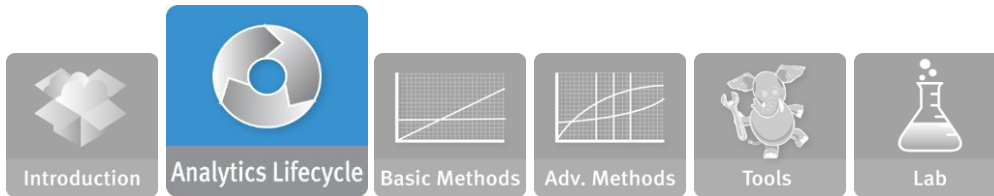
The leader –
adjusts the sails
John Maxwell
(Leadership Author)

Check Your Knowledge



Your Thoughts?

- In which phase would you expect to invest most of your project time and why? Where would expect to spend the least time?
- What are the benefits of doing a pilot program before a full scale rollout of a new analytical methodology? Discuss this in the context of the mini case study.
- What kinds of tools would be used in the following phases, and for which kinds of use scenarios?
 - ▶ Phase 2: Data Preparation
 - ▶ Phase 4: Model Execution
- Now that you have completed the analytical project at Yoyodyne, you have an opportunity to repurpose this approach for an online eCommerce company. What phases of the lifecycle do you need to focus on to identify ways to do this?



Module 2: Summary

Key points covered in this module:

- The Data Analytics Lifecycle was applied to a case study scenario
- A business problem was framed as an analytics problem
- The four main deliverables in an analytics project were identified

Lab Exercise 1: Introduction to Data Environment



This first lab introduces the Analytics Lab Environment you will be working on throughout the course.

After completing the tasks in this lab you should be able to:

- Authenticate and access the Virtual Machine (VM) assigned to you for all of your lab exercises
- Locate data sets you will be working with for the course's labs
- Use meta commands and PSQL to navigate through the data sets
- Create sub-sets of the big data, using table joins and filters to analyze subsequent lab exercises