

## LEC.4.



# Introduction to Data Science



# 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

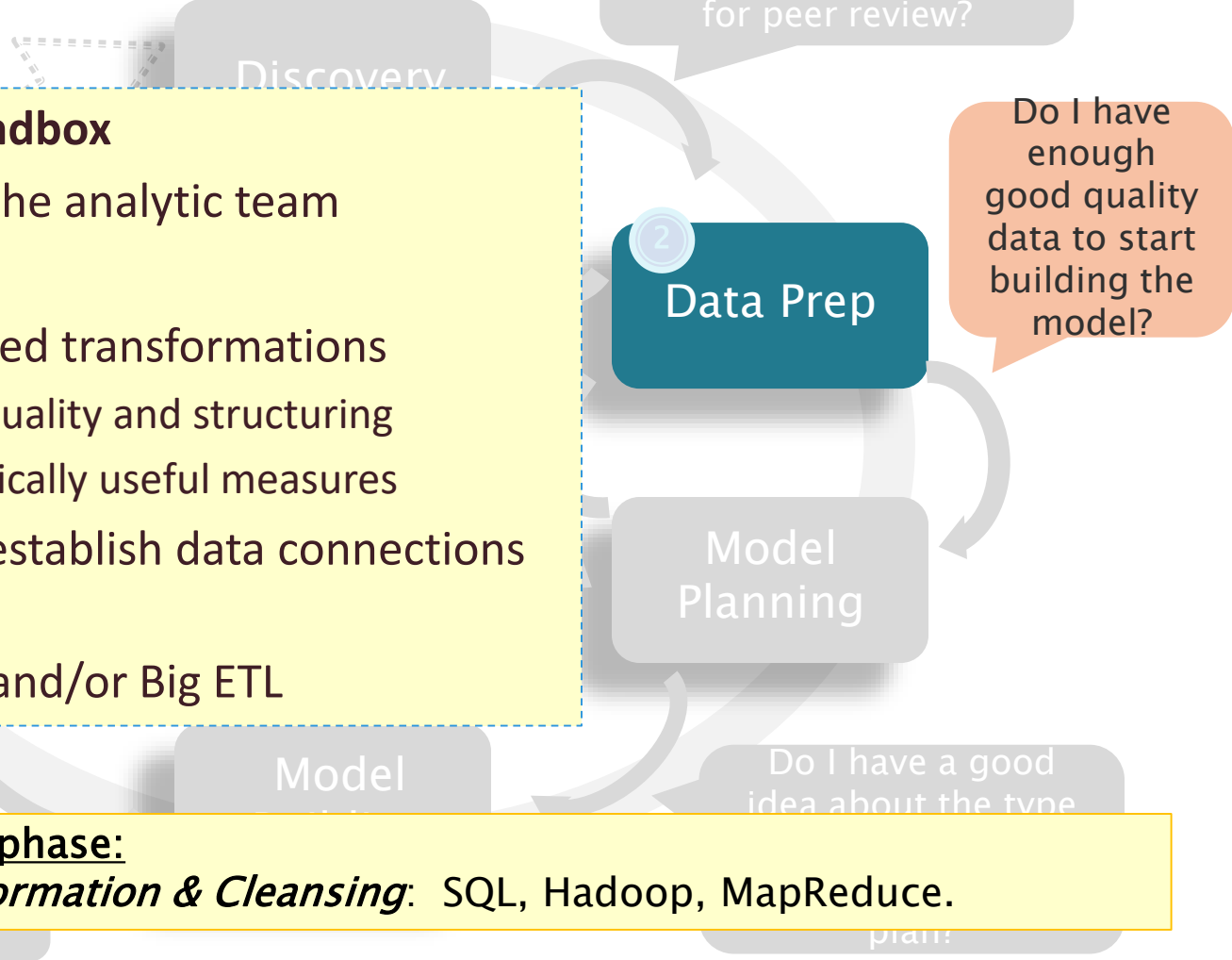
Mini Case  
Study



Sample <i>Business</i> Problems	Qualifiers	Analytical Approach
<ul style="list-style-type: none"><li>• How can we improve on x?</li><li>• What's happening real-time? Trends?</li><li>• How can we use analytics differentiate ourselves</li><li>• How can we use analytics to innovate?</li><li>• How can we stay ahead of our biggest competitor?</li></ul>	<p>Will the focus and scope of the problem change if the following dimensions change:</p> <ul style="list-style-type: none"><li>• Time</li><li>• People – how would x change this?</li><li>• Risk – conservative/aggressive</li><li>• Resources – none/unlimited</li><li>• Size of Data?</li></ul>	<p>Define an analytical approach, including key terms, metrics, and data needed.</p>
<p><i>Churn Prediction for Yoyodyne Bank</i></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"><li>• <b>Time:</b> Trailing 5 months</li><li>• <b>People:</b> Working team and business users from the Bank</li><li>• <b>Risk:</b> the project will fail if we cannot determine valid predictors of churn</li><li>• <b>Resources:</b> analytic sandbox, OLTP system</li><li>• <b>Data:</b> Use 24 months for the training set, then analyze 5 months of historical data for those customers who churned</li></ul>	<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
- **Perform ELT**
  - ▶ Determine needed transformations
    - ▶ Assess data quality and structuring
    - ▶ Derive statistically useful measures
  - ▶ Determine and establish data connections for raw data
  - ▶ Execute Big ELT and/or Big ETL

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

# 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

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

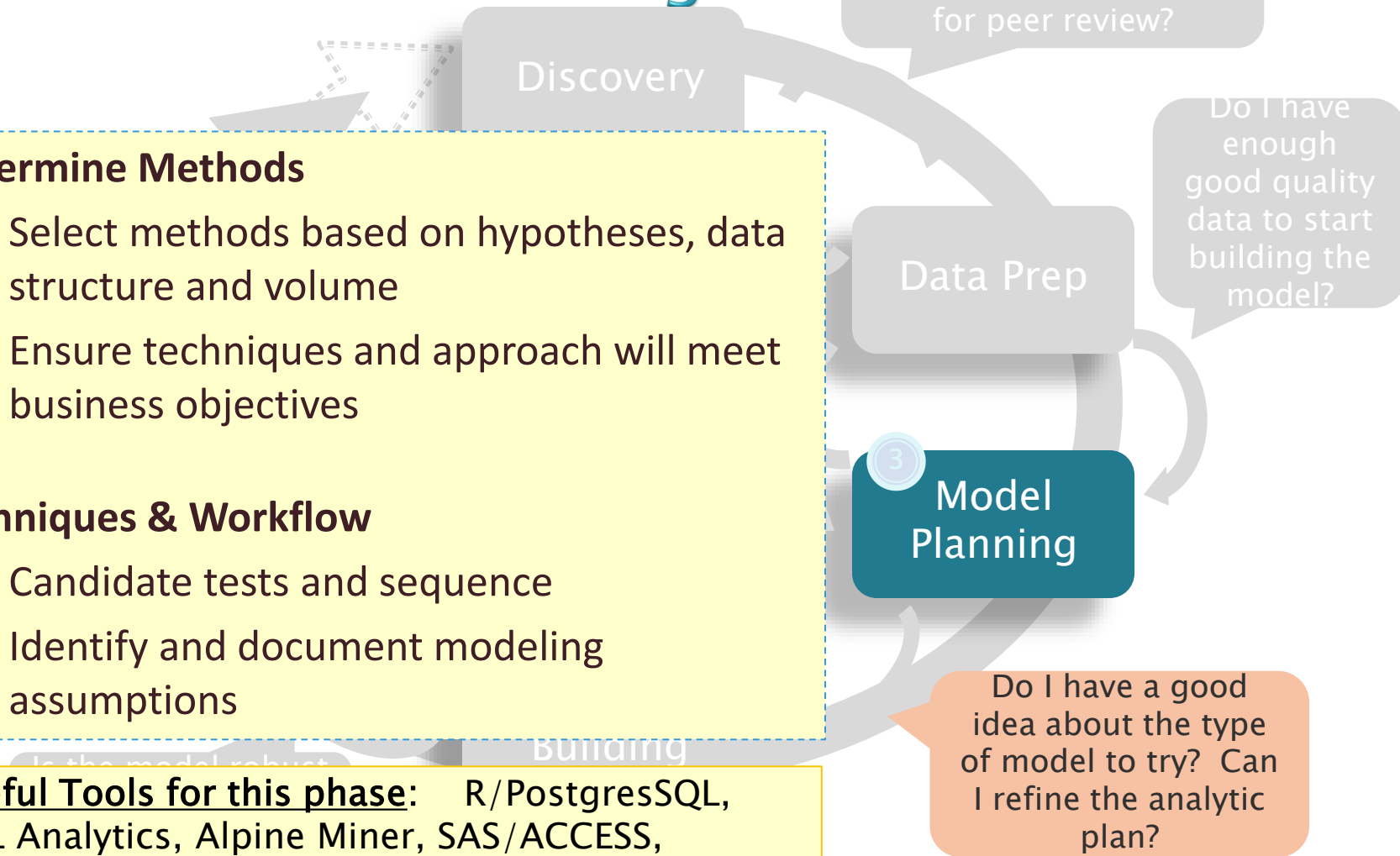
# Data Analytics Lifecycle

## Phase 3: Model Planning



- **Determine Methods**
  - ▶ Select methods based on hypotheses, data structure and volume
  - ▶ Ensure techniques and approach will meet business objectives
- **Techniques & Workflow**
  - ▶ Candidate tests and sequence
  - ▶ Identify and document modeling assumptions

- **Useful Tools for this phase:** R/PostgreSQL, SQL Analytics, Alpine Miner, SAS/ACCESS, SPSS/OBDC



# 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

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

Data Prep

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

3  
Model Planning

Do I have a good idea about the type of model to try? Can I refine the analytic plan?

# 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.
- At this point, a Data Scientist would assess the methods and select the best model for the situation

Market Sector	Analytic Techniques /Methods Used
Retail Business	<b>decision tree</b>
Daily Grocery	<b>MLR (multiple linear regression), and decision tree</b>
Wireless Telecom	Neural network, data mining , <b>decision tree</b> , hierarchical neurofuzzy systems, rule evolver
Retail Banking	<b>Multiple regression</b>
Wireless Telecom	<b>Logistic regression, neural network, decision tree</b>



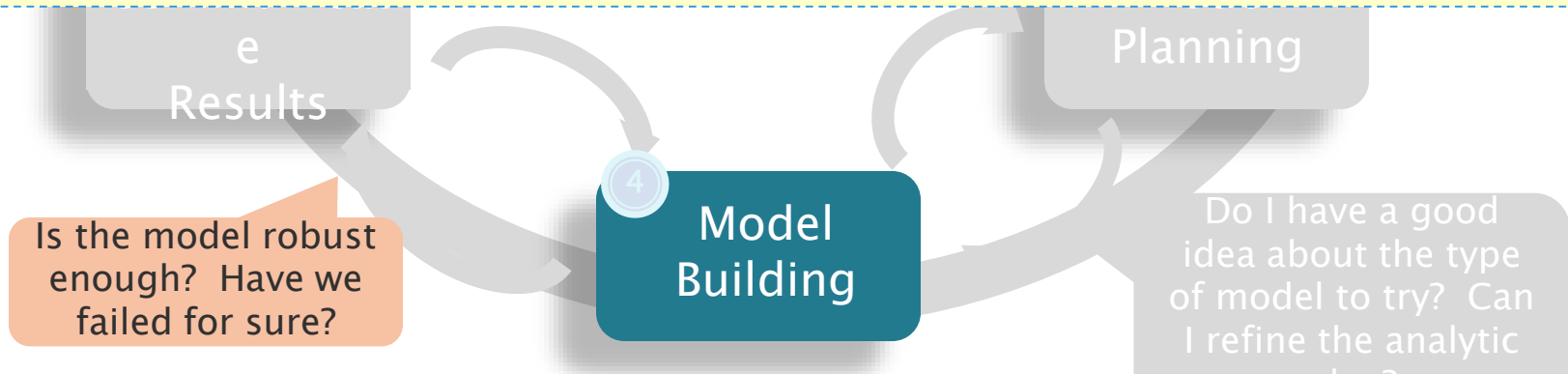
# Data Analytics Lifecycle

## 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

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

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

***Did we succeed? Did we fail?***

- Interpret the results
- Compare to IH's from Phase 1
- Identify key findings
- Quantify business value
- Summarizing findings, depending on audience

Operationali  
ze

5  
Communicate  
Results

Discovery

Building

*Mini Case Study:  
Churn Prediction for  
Yoyodyne Bank*

For the YoyoDyne Case Study,  
what would be some possible results and key  
findings?

# Data Analytics Lifecycle

## Phase 6: Operationalize



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

Discovery

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

6  
Operationalize

Communicate  
Results

- Run
- Assess the benefits
- Provide final deliverables
- Implement the model in the production environment
- Define process to update, retrain, and retire the model, as needed

Model  
Building

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

Do I have a good idea about the type of model to try? Can I refine the analytic plan?



# Analytic Plan

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"> <li>• <b>Sponsor Presentation</b> addressing:                             <ul style="list-style-type: none"> <li>• Are the results good for me?</li> <li>• What are the benefits of the findings?</li> <li>• What are the implications of this for me?</li> </ul> </li> </ul>
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"> <li>• <b>Sponsor Presentation</b> addressing:                             <ul style="list-style-type: none"> <li>• What's the business impact of doing this?</li> <li>• What are the risks? ROI?</li> <li>• How can this be evangelized within the organization (and beyond)?</li> </ul> </li> </ul>
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"> <li>• Show the <b>analyst presentation</b></li> <li>• Determine if the reports will change</li> </ul>
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"> <li>• Share the <b>code</b> from the analytical project</li> <li>• Create <b>technical document</b> on how to implement it.</li> </ul>
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"> <li>• Share the <b>code</b> from the analytical project</li> <li>• Create <b>technical document</b> on how to implement it.</li> </ul>
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"> <li>• Show the <b>analyst presentation</b></li> <li>• Share the <b>code</b></li> </ul>

# 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

# Applying the Data Analytics Lifecycle



- In a typical Data Analytics Problem - you would have gone through:
  - Phase 1 – Discovery - have the problem framed
  - Phase 2 – Data Preparation - have the data prepared
- Now you need to plan the model and determine the method to be used.

# Phase 3 – Model Planning

Model planning is the process of determining the appropriate analytic method based on the problem. It also depends on the type of data and the computational resources available.

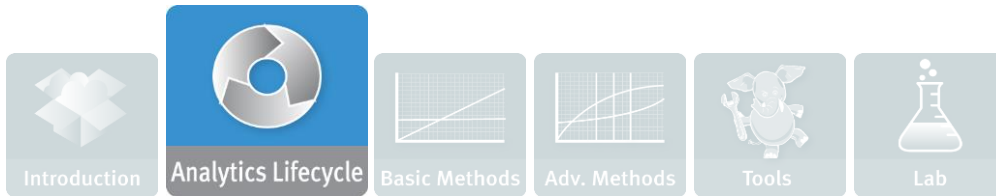
Discovery

Data Prep

Model Planning

Do I have a good idea about the type of model to try? Can I refine the analytic plan?





# 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