# Week 1 Unit 1:
# Introduction to Data Science

SAP

openSAP
open.sap.com

# What to expect in the next 6 weeks?

# Introduction to Data Science
## Curriculum flow (weeks 1-3)

**1** Business & Data Understanding

- Introduction to Data Science
- Introduction to Project Methodologies
- Business Understanding Phase – Overview
- Defining Project Success Criteria
- Data Understanding Phase – Overview
- Initial Data Analysis & Exploratory Data Analysis

**Weekly Assignment**

**2** Data Preparation

- Data Preparation Phase – Overview
- Predictive Modeling Methodology – Overview
- Data Manipulation
- Selecting Data – Variable and Feature Selection
- Data Encoding

**Weekly Assignment**

**3** Modeling (1)

- Modeling Phase – Overview
- Detecting Anomalies
- Association Analysis
- Cluster Analysis
- Classification Analysis with Regression

**Weekly Assignment**

# Introduction to Data Science
## Curriculum flow (weeks 4-6)

**4** Modeling (2)

- Classification Analysis with Decision Trees
- Classification Analysis with KNN, NN, and SVM
- Time Series Analysis
- Ensemble Methods
- Simulation & Optimization
- Automated Modeling

**Weekly Assignment**

**5** Evaluation

- Evaluation Phase – Overview
- Model Performance Metrics
- Model Testing
- Improving Model Performance

**Weekly Assignment**

**6** Deployment & Maintenance

- Deployment Phase – Overview
- Deployment Options
- Monitoring & Maintenance
- Automating Deployment & Maintenance
- Myths & Challenges
- Data Science Applications and References

**Weekly Assignment**

**Final Exam**

# Introduction to Data Science
## Cumulative points lead to record of achievement

**Participate in Weekly Assignment** (Weeks1-6)

6 assignments
6 x 30 = 180 points

**Final Exam**
(Week 7)

180 points

**Record of Achievement**

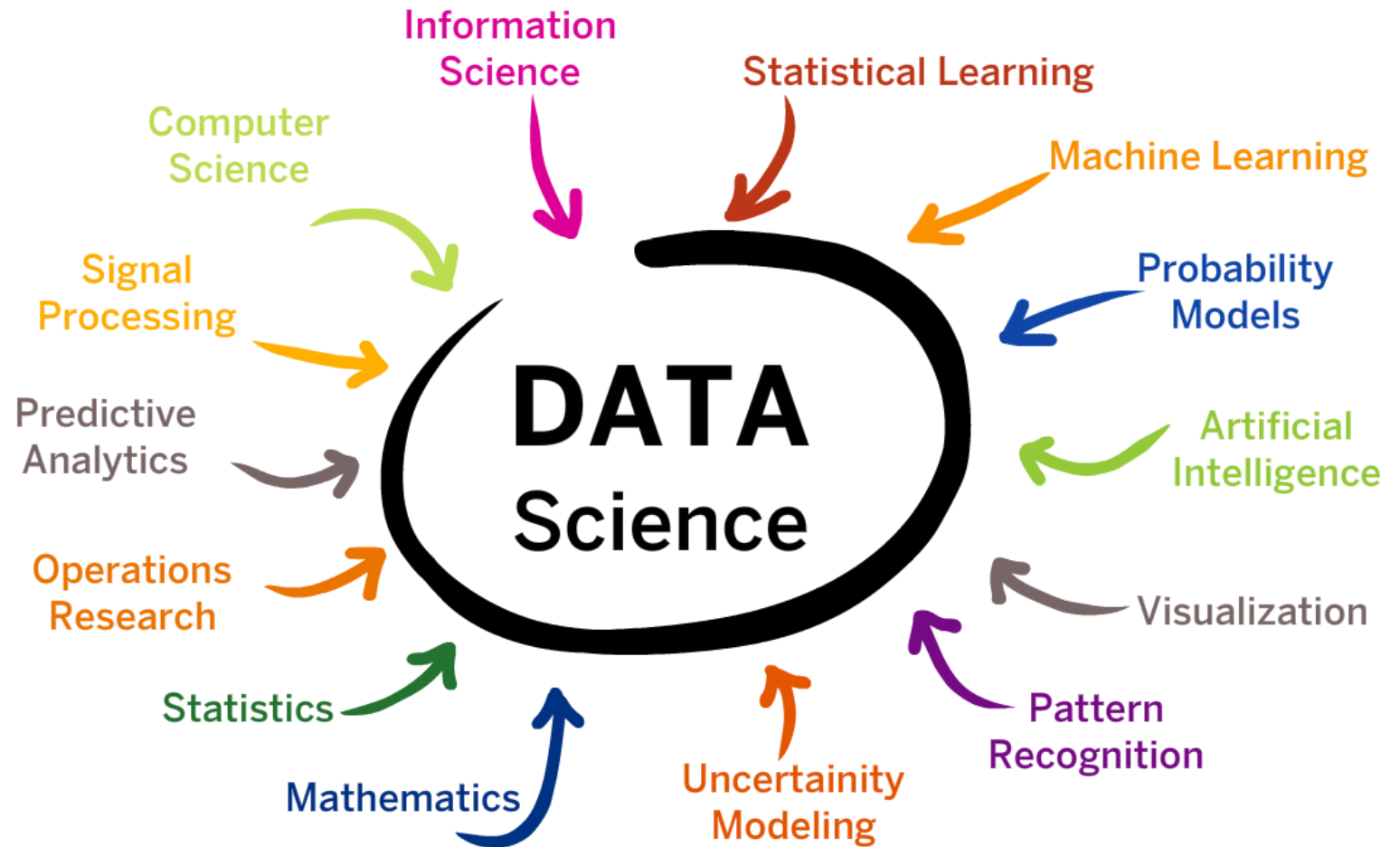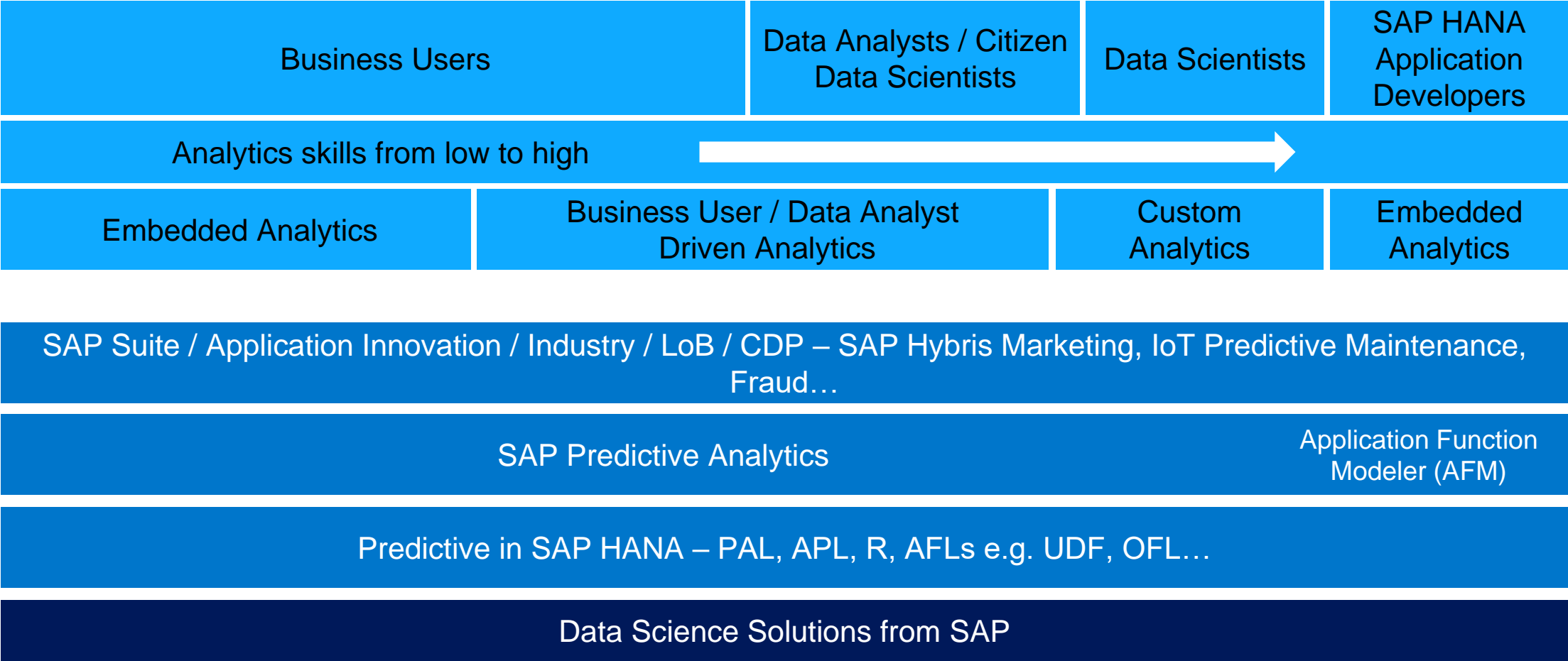When results above
**180 points**

## What is data science?

**Data science** is an interdisciplinary field about processes and systems that enable the extraction of knowledge or insights from data.
Data science employs techniques and theories drawn from a wide range of disciplines.

# Introduction to Data Science
## Data science personas

| Business Users | Data Analysts / Citizen Data Scientists | Data Scientists | SAP HANA Application Developers |
|---|---|---|---|

**Analytics skills from low to high** →

| Embedded Analytics | Business User / Data Analyst Driven Analytics | Custom Analytics | Embedded Analytics |
|---|---|---|---|

SAP Suite / Application Innovation / Industry / LoB / CDP – SAP Hybris Marketing, IoT Predictive Maintenance, Fraud…

SAP Predictive Analytics · Application Function Modeler (AFM)

Predictive in SAP HANA – PAL, APL, R, AFLs e.g. UDF, OFL…

Data Science Solutions from SAP

# Introduction to Data Science
## Data science solutions from SAP

| SAP Predictive Analytics | SAP Lumira | SAP HANA Studio / AFM | SAP RDS Analytics Solutions | SAP Industry & LoB Solutions | Partner Analytical BI & Tools |
|---|---|---|---|---|---|

**SAP HANA**

| Predictive Analysis Library (PAL) | Business Function Library | Automated Predictive Library | Simulation | Optimization |
|---|---|---|---|---|
| Text Search | Text Analysis and Mining | Spatial Analysis | Graph Engine | Rules Engine |

R

| SAP IQ | HADOOP | SAP ESP | 3rd Party Data Source | SAP Data Services |
|---|---|---|---|---|

Data Connectors

**Data types**
Connect to SAP HANA directly or via Sybase IQ / Hadoop / ESP / Data Services

| Transaction Data | Unstructured Data | Real-Time Data | Location Data | Machine Data | Others |
|---|---|---|---|---|---|

# Introduction to Data Science
## SAP HANA Predictive Analysis Library (PAL)

**Build High-Performance Predictive Apps**

- The SAP HANA Predictive Analysis Library (PAL) is a built-in C++ library for performing in-memory data mining and statistical calculations.

- PAL is designed to provide high performance on large datasets for real-time analytics.

Hadoop / Sybase IQ, Sybase ASE, Teradata

Spatial, Machine, Real-Time Data

**SAP HANA**

**Main Memory**

**Virtual Tables**

SQLScript Optimized Query Plan

PAL

R Scripts

Unstructured

**Text Analysis**

**Spatial Data**

KNN classification

K-means

ABC classification

Weighted score tables

Regression

C4.5 decision tree

Association analysis: market basket

**R Engine**

**SAP HANA Studio/AFM, Apps & Tools**

# Introduction to Data Science
## SAP HANA Predictive Analysis Library (PAL) algorithms

- SAP HANA Predictive Analysis Library (PAL) contains a wide range of algorithms that can be deployed for in-HANA and standalone data science applications.

- A wide range of algorithms are available for the following types of analysis:

| SAP HANA Predictive Analysis Library |
|---|
| **Association Analysis** |
| **Classification Analysis** |
| **Regression** |
| **Cluster Analysis** |
| **Time Series Analysis** |
| **Probability Distribution** |
| **Outlier Detection** |
| **Link Prediction** |
| **Data Preparation** |
| **Statistic Functions (Univariate)** |
| **Statistic Functions (Multivariate)** |

# Introduction to Data Science
## SAP HANA Automated Predictive Library (APL) algorithms

- SAP HANA APL is an application function library (AFL) that lets you use the data mining capabilities of the SAP Predictive Analytics automated analytics engine on your customer datasets stored in SAP HANA.

- You can create a wide range of models to answer your business questions.



Classification Models

Clustering Models

Regression Models

APL

Time Series Analysis

Social Network Analysis

Recommendation

# Introduction to Data Science
## R integration for SAP HANA and standalone

# Introduction to Data Science
## SAP Predictive Analytics

- SAP Predictive Analytics is built for both data scientists and business / data analysts, making predictive analytics accessible to a broad spectrum of users.

- Automated and expert modes

- Used to automate data preparation, predictive modeling, and deployment tasks

- Rich pre-built modelling functionality

- PAL, APL, and R  language support

- Advanced visualization

- Native integration with SAP HANA

# Introduction to Data Science
## Application function modeler (AFM)

- Graphical tool to build advanced applications in SAP HANA

- Web-based flow-graph editor
  - Support for AFL, R, SDI, & SDQ
  - Used to create procedures or task runtime operations
  - Interoperability with SAP HANA studio AFM

- SAP HANA studio-based AFM
  - PAL function support including time series, clustering, classification, and statistics
  - General usability enhancements for an easier, simpler, and more functional experience

# Thank you

**Contact information:**

**open@sap.com**

# © 2016 SAP SE or an SAP affiliate company. All rights reserved.

# Week 1 Unit 2: Introduction to Project Methodologies

openSAP
open.sap.com

# Introduction to Project Methodologies
## Why should there be a project methodology?

- The data science process must be reliable and repeatable by people with little data science background.

- A project methodology:
  - Provides a framework for recording experience
  - Allows projects to be replicated
  - Provides an aid to project planning and management
  - Is a "comfort factor" for new adopters
  - Reduces dependency on "stars"

**TIME**

Task 1

Task 2

Task 3

Task 4

# Introduction to Project Methodologies
Cross-industry standard process for data mining (CRISP-DM)

# Introduction to Project Methodologies
## CRISP-DM – Phase 1: Business Understanding

# Introduction to Project Methodologies
## CRISP-DM – Phase 2: Data Understanding



Business Understanding → **Data Understanding** → Data Preparation → Modeling → Evaluation → Deployment

| Tasks | Outputs |
|---|---|
| Collect Initial Data | Initial Data Collection Report |
| Describe Data | Data Description Report |
| Explore Data | Data Exploration Report |
| Verify Data Quality | Data Quality Report |

**Key**
- TASKS
- OUTPUTS

# Introduction to Project Methodologies
## CRISP-DM – Phase 3: Data Preparation

# Introduction to Project Methodologies
CRISP-DM – Phase 4: Modeling

# Introduction to Project Methodologies
CRISP-DM – Phase 5: Evaluation

# Introduction to Project Methodologies
## CRISP-DM – Phase 6: Deployment

# Introduction to Project Methodologies
CRISP-DM – Update

# Thank you

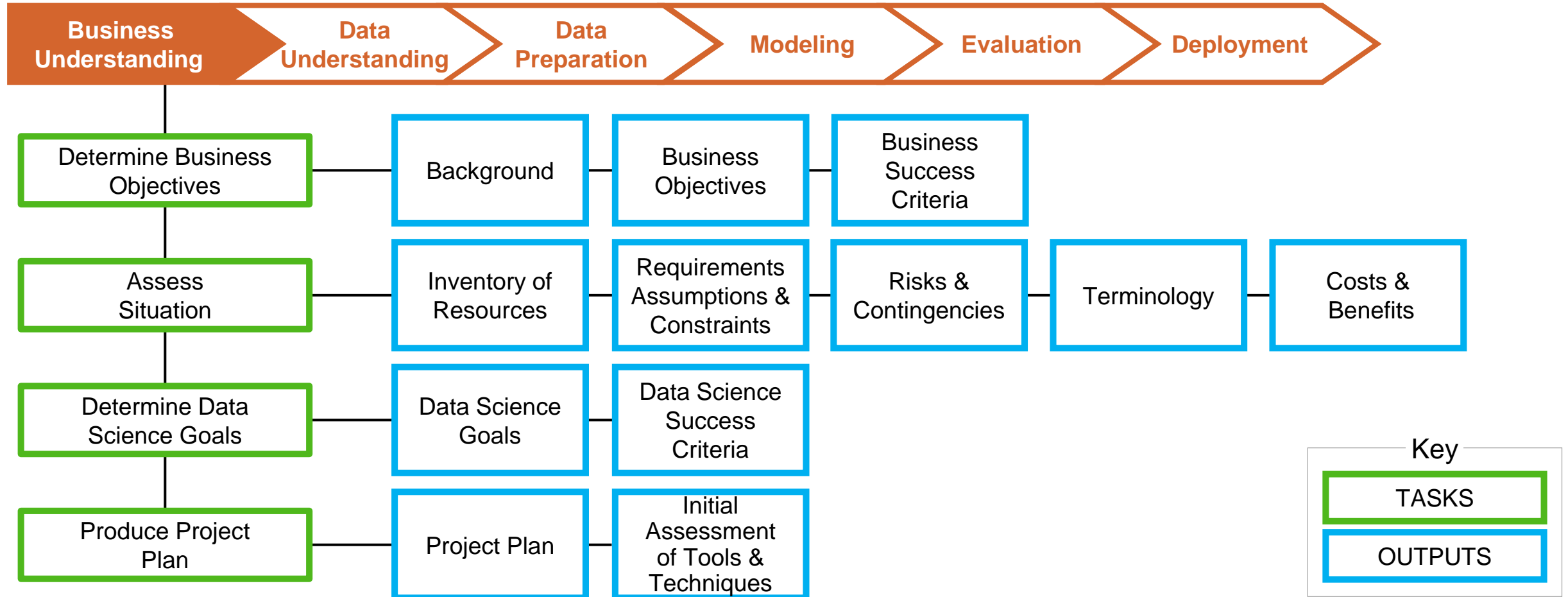**Contact information:**

**open@sap.com**

Week 1 Unit 3: **Business Understanding Phase – Overview**

SAP

# Business Understanding Phase – Overview
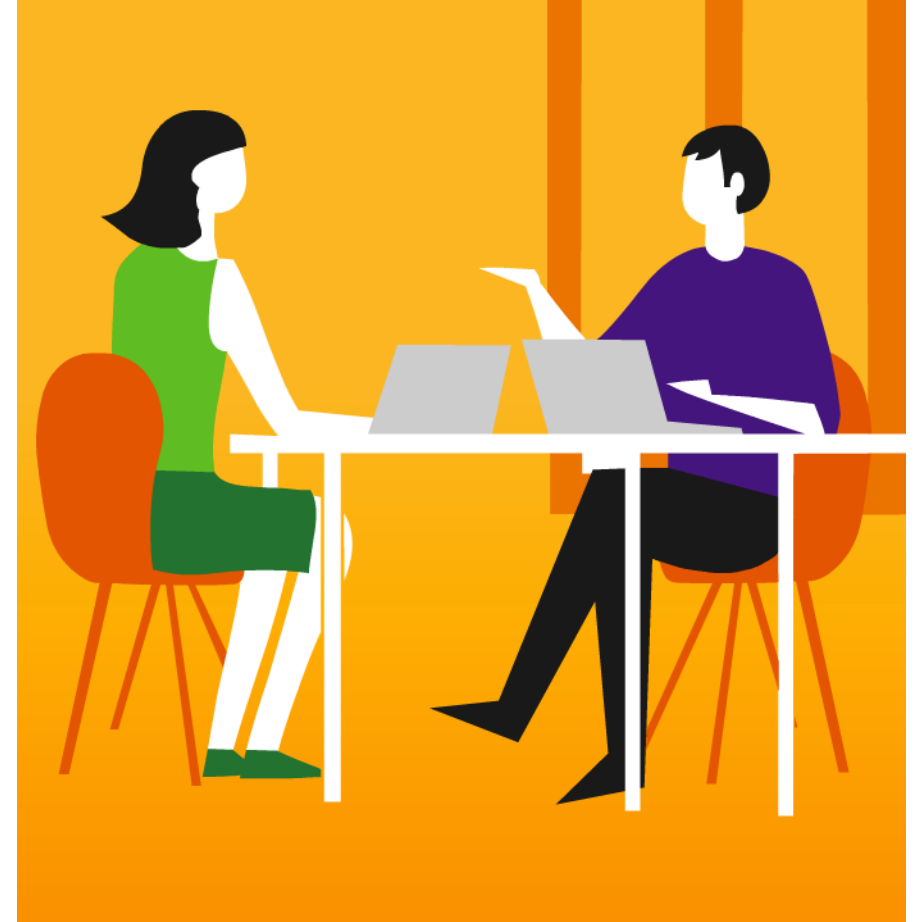## CRISP-DM – Phase 1: Business Understanding

# Business Understanding Phase – Overview

## Phase 1.1: Determine Business Objectives

- **Task**
  - The first objective of the data analyst is to thoroughly understand, from a business perspective, what the client really wants to accomplish.

- **Outputs**
  - Background
  - Business Objectives
  - Business Success Criteria

# Business Understanding Phase – Overview
## Phase 1.2: Assess Situation

- **Task**
  - In the previous task, your objective is to quickly get to the crux of the situation. Here, you want to flesh out the details.

- **Outputs**
  - Inventory of Resources
  - Requirements, Assumptions, & Constraints
  - Risks & Contingencies
  - Terminology
  - Costs & Benefits

# Business Understanding Phase – Overview
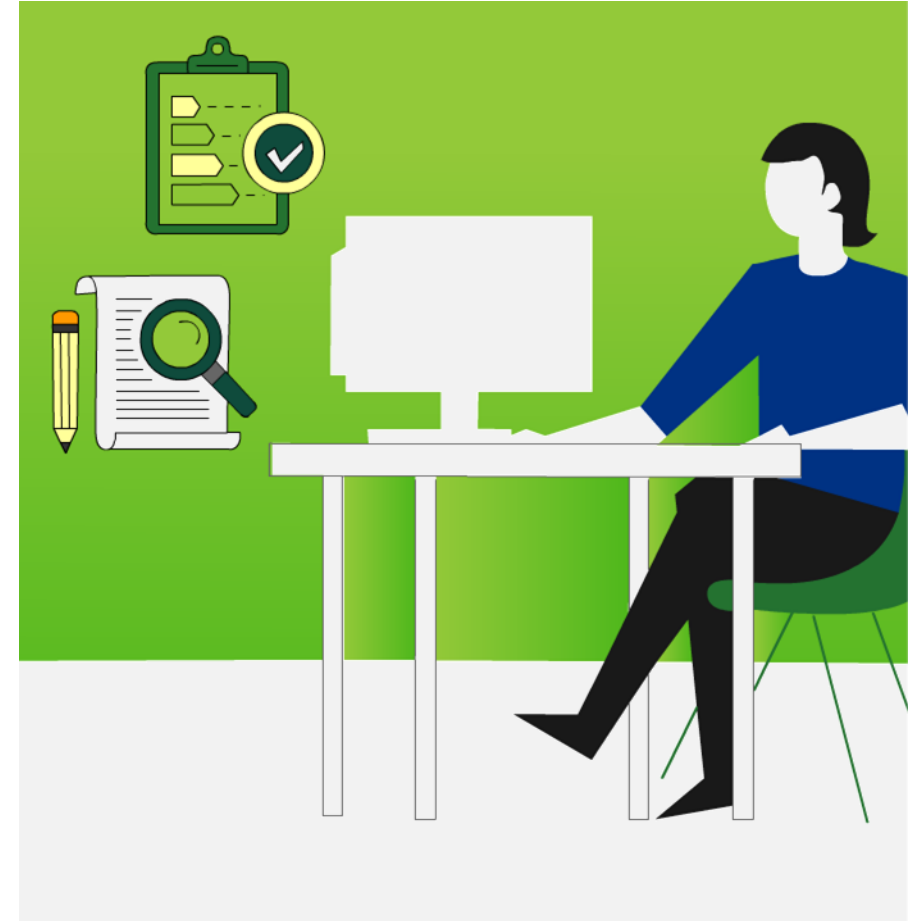## Phase 1.3: Determine Data Science Goals

- **Task**
  - A ***business goal*** states objectives in business terminology.
  - A ***data science goal*** states project objectives in technical terms.

- **Outputs**
  - Describe data science goals.
  - Define data science success criteria.

# Business Understanding Phase – Overview
## Phase 1.4: Produce Project Plan

- **Task**

  - Describe the intended plan for achieving the data mining goals and thereby achieving the business goals.

- **Output**

  - Project plan with project stages, duration, resources, etc.
  - Initial assessment of tools & techniques.

# Thank you

**Contact information:**

**open@sap.com**

# Week 1 Unit 4: Defining Project Success Criteria

# Defining Project Success Criteria

Business and data science project success criteria: reminder

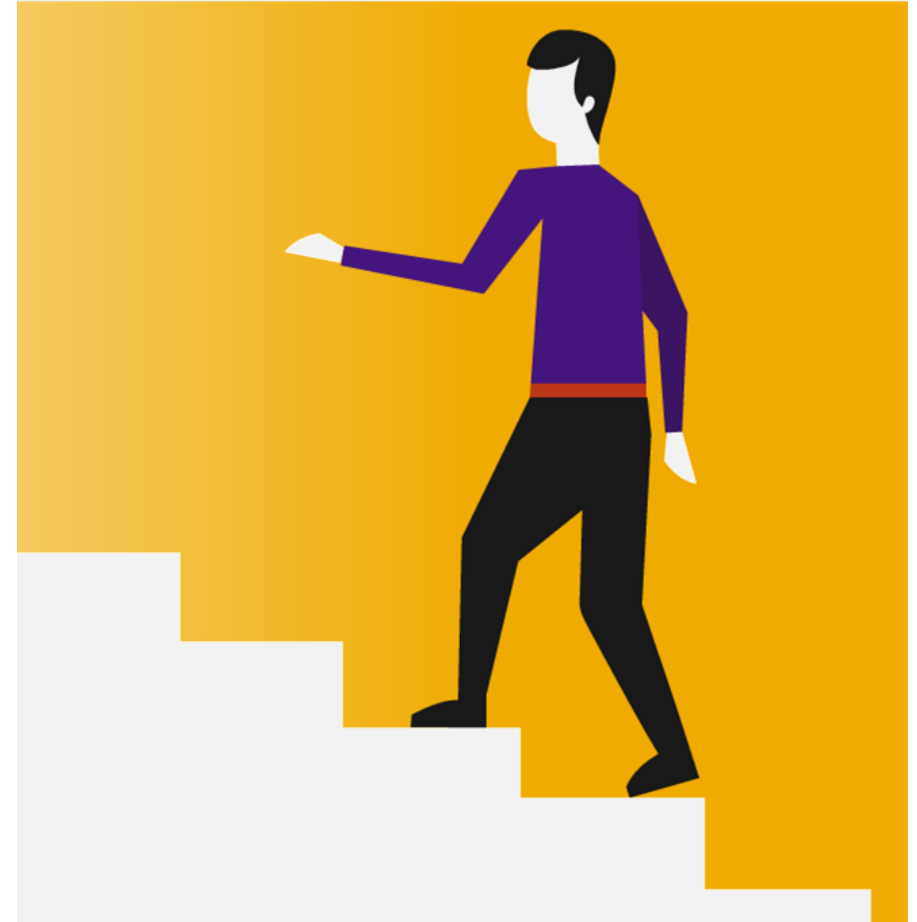**Business success criteria**

- Describe the criteria for a successful or useful outcome to the project from the business point of view.

**Data science success criteria**

- Define the criteria for a successful outcome to the project in technical terms.

# Defining Project Success Criteria
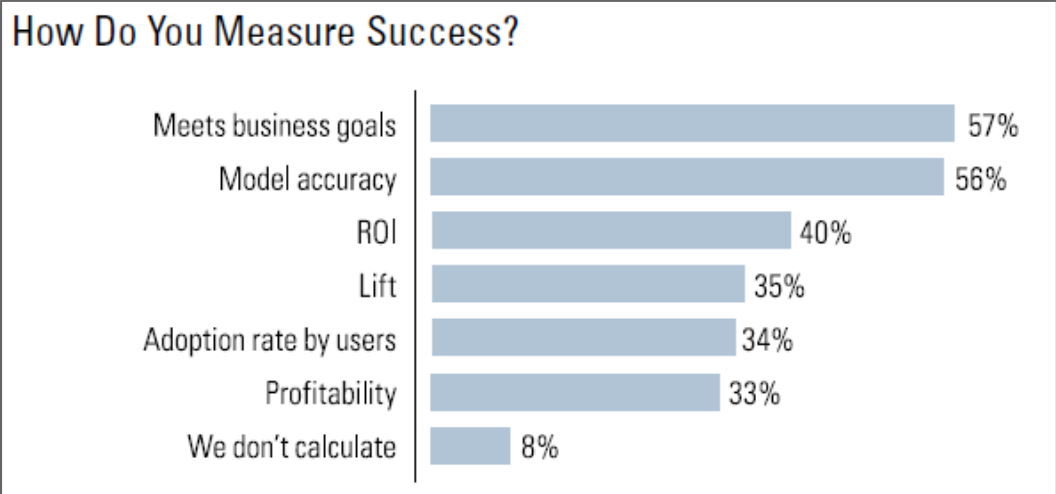## Recent industry surveys



### How Do You Measure Success?

| | |
|---|---|
| Meets business goals | 57% |
| Model accuracy | 56% |
| ROI | 40% |
| Lift | 35% |
| Adoption rate by users | 34% |
| Profitability | 33% |
| We don't calculate | 8% |

*Figure 5. Based on 110 users who have implemented predictive analytics initiatives that offer "very high" or "high" value. Respondents could select multiple choices.*

### PREDICTIVE ANALYTICS
#### Extending the Value of Your Data Warehousing Investment
By Wayne W. Eckerson

In their Third Annual Data Miner Survey, Rexer Analytics, an analytics and renowned CRM consulting firm based in Winchester, Massachusetts asked the BI community "How do you evaluate project success in data mining?" Out of 14 different criteria, a massive 58% ranked "Model performance" (lift, R2, etc) as the primary factor.

| | |
|---|---|
| Model performance (lift, R2, etc) | 58% |
| Improve efficiency | 49% |
| Produced new business insights | 48% |
| Revenue growth | 44% |
| Increased sales | 42% |
| Return on Investment (ROI) | 42% |
| Increased profit | 38% |
| Improved quality of product/service | 35% |
| Increased customer satisfaction | 34% |
| Reduce costs of producing products/.. | 27% |
| Customer service improvements | 21% |
| Results were published | 13% |
| Produced new scientific insights | 12% |

# Defining Project Success Criteria
## Model success criteria: descriptive or predictive models
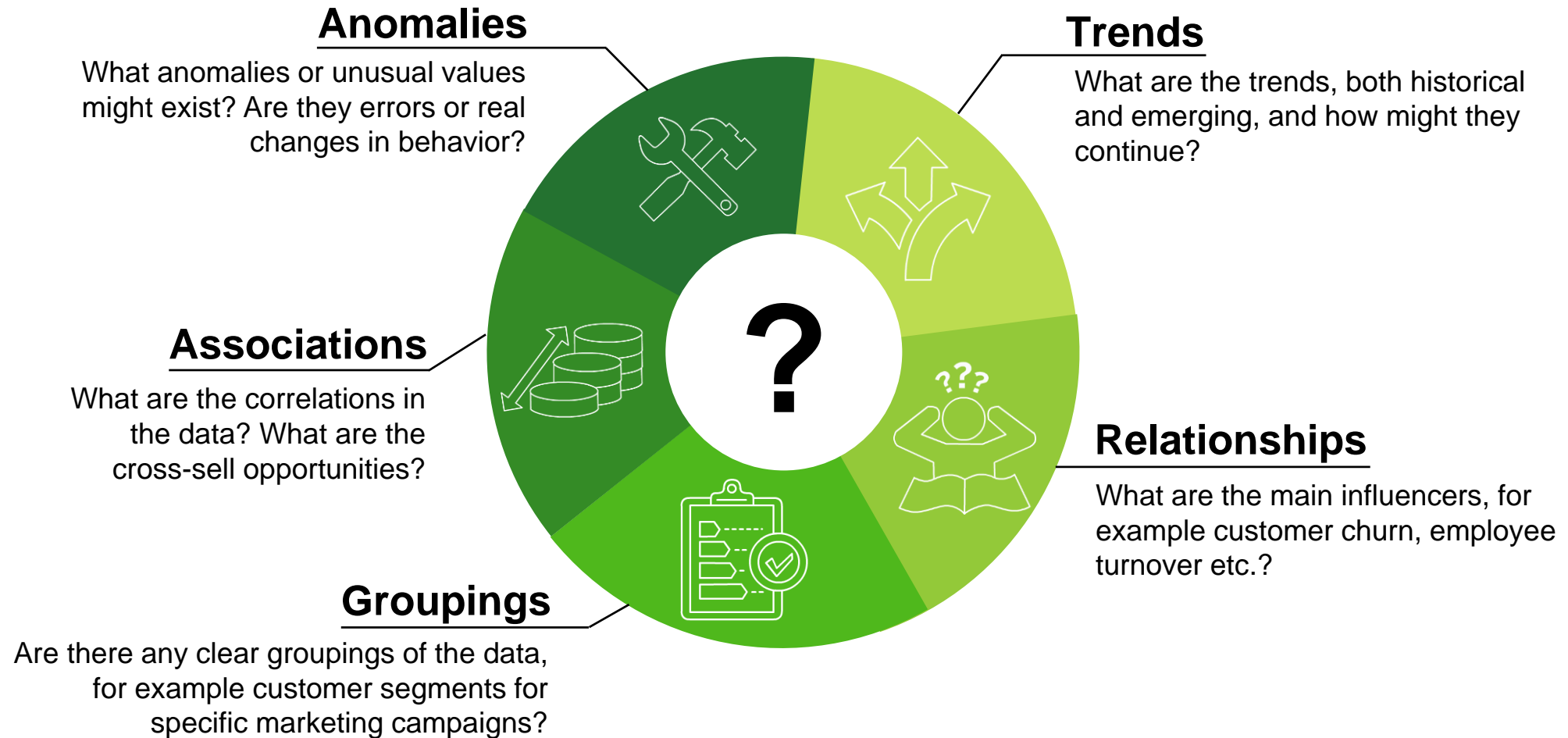
### Descriptive Models

- Descriptive analysis describes or summarizes raw data and makes it more interpretable. It describes the past – i.e. any point of time that an event occurred, whether it was one minute ago or one year ago.

- Descriptive analytics are useful because they allow us to learn from past behaviors and understand how these might influence future outcomes.

- Common examples of descriptive analytics are reports that provide historical insights regarding a company's production, financials, operations, sales, finance, inventory and customers.

- Descriptive analytical models include cluster models, association rules, and network analysis.

### Predictive Models

- Predictive analysis predicts what might happen in the future – providing estimates about the likelihood of a future outcome.

- One common application is the use of predictive analytics to produce a credit score. These scores are used by financial services to determine the probability of customers making future credit payments on time.

- Typical business uses include: understanding how sales might close at the end of the year, predicting what items customers will purchase together, or forecasting inventory levels based upon a myriad of variables.

- Predictive analytical models include classification models, regression models, and neural network models.
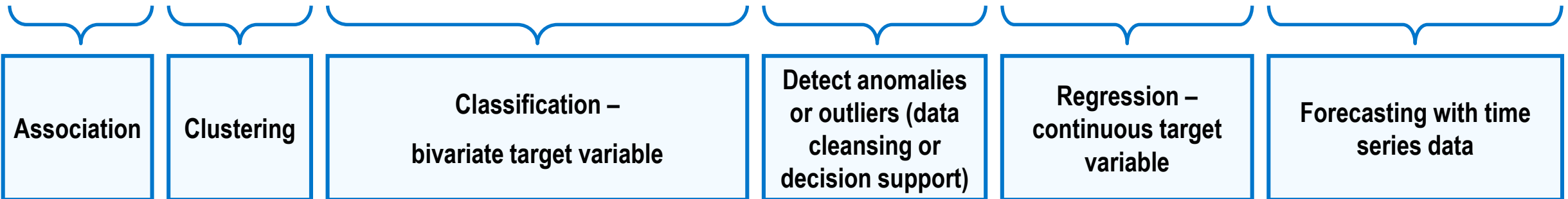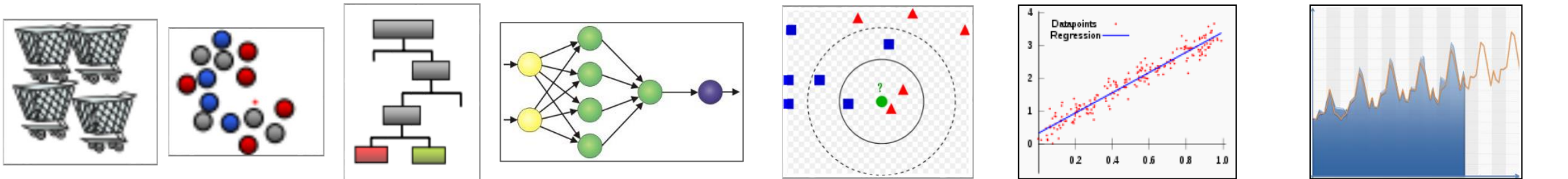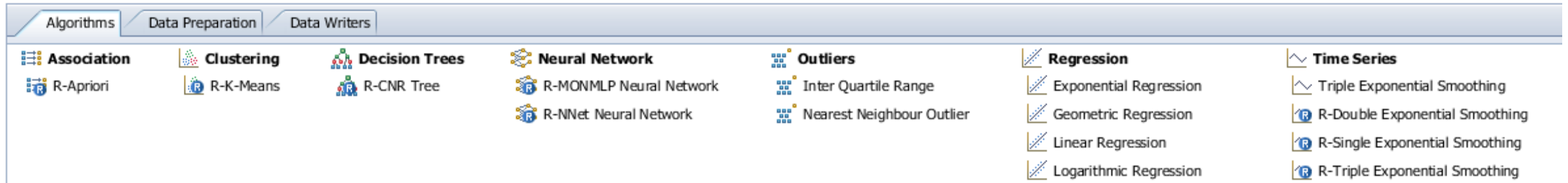
# Defining Project Success Criteria
## Model success criteria: choosing the algorithm



**Anomalies**
What anomalies or unusual values might exist? Are they errors or real changes in behavior?

**Trends**
What are the trends, both historical and emerging, and how might they continue?

**Associations**
What are the correlations in the data? What are the cross-sell opportunities?

**Relationships**
What are the main influencers, for example customer churn, employee turnover etc.?

**Groupings**
Are there any clear groupings of the data, for example customer segments for specific marketing campaigns?

# Defining Project Success Criteria
There is a wide range of algorithms to choose from…



| Algorithms | Data Preparation | Data Writers |
| --- | --- | --- |

**Association**
R-Apriori

**Clustering**
R-K-Means

**Decision Trees**
R-CNR Tree

**Neural Network**
R-MONMLP Neural Network
R-NNet Neural Network

**Outliers**
Inter Quartile Range
Nearest Neighbour Outlier

**Regression**
Exponential Regression
Geometric Regression
Linear Regression
Logarithmic Regression

**Time Series**
Triple Exponential Smoothing
R-Double Exponential Smoothing
R-Single Exponential Smoothing
R-Triple Exponential Smoothing

| Association | Clustering | Classification – bivariate target variable | Detect anomalies or outliers (data cleansing or decision support) | Regression – continuous target variable | Forecasting with time series data |
| --- | --- | --- | --- | --- | --- |

**Descriptive Models**          **Predictive Models**

# Defining Project Success Criteria
## Which business question do you need to answer?

### Classification

**Who** will (buy | fraud | churn …) next (week | month | year…)?

### Regression

**What** will the (revenue | # churners) be next (week | month…)?

### Segmentation or Clustering

**What are the groups** of customers with similar (behavior | profile …)?

### Forecasting (Time Series Analysis)

**What** will the (revenue | # churners…) be over next year on a monthly basis?

### Link Analysis

**Analyze interactions** to identify (communities | influencers…)

### Association or Recommendation Engines

**Provides** recommendations on web sites or to retailers – basket analysis

# Defining Project Success Criteria
## Model accuracy and robustness

- The accuracy and robustness of the model are two major factors to determine the quality of the prediction, which reflects how successful the model is.

- **Accuracy** is often the starting point for analyzing the quality of a predictive model, as well as an obvious criterion for prediction. Accuracy measures the ratio of correct predictions to the total number of cases evaluated.

- The **robustness** of a predictive model refers to how well a model works on alternative data. This might be hold-out data or new data that the model is to be applied onto.

# Defining Project Success Criteria
Training and testing: data cutting strategies

- Central to developing predictive models and assessing if they are successful is a **train-and-test regime**.

- Data is partitioned into **training** and **test** subsets. There are a variety of **cutting strategies** (e.g. random/sequential/periodic).

- We build our model on the training subset (called the **estimation** subset) and evaluate its performance on the test subset (a **hold-out** sample called the **validation** subset).

- Simple two and three-way data partitioning is shown in the diagram.

**Two-Way Partition**

Complete Data Set

Training Set | Test Set

Develop models | Evaluate models

**Three-Way Partition**

Complete Data Set

Training Set | Validation Set | Test Set

Develop models | Evaluate models | Evaluate the selected models

Select the best model based upon validation set performance

# Defining Project Success Criteria
## Example success criteria for predictive models

**Business Criteria:**

- Models meet business goals – the model meets the business objectives specified in CRISP-DM phase 1.1 as defined by the customer

- The model's contributing variables and the variable categories make "business sense"

**Model Performance Criteria:**

- Depends on the algorithm. For a classification model for example:

  - Model accuracy compared to any previous, similar models. There are a variety of accuracy measures that will be discussed in this course

  - Model robustness – Models have acceptable robustness

**Software Usability Criteria:**

- Speed and ease of model <u>development</u> – so the customer can build new models and update existing models quickly

- Speed and ease of model <u>deployment</u> – so the customer can create Apply datasets easily and deploy models quickly with the required outputs (probabilities, deciles, etc.) – speed to market

- Ease of model <u>maintenance</u> – so the customer can easily define when models require refreshing/rebuilding and undertake this quickly and easily

- Integration capability with other systems

# Defining Project Success Criteria
Example success criteria for descriptive models

## Business Criteria:

- Models meet business goals – the model meets the business objectives specified in CRISP-DM phase 1.1
- Contributing variables and categories make business sense

## Model Performance Criteria:

- Depends on the algorithm. For a cluster model for example:
  - Determining the clustering tendency of a set of data (distinguishing whether non-random structure actually exists in the data)
  - Comparing the results to given class labels (comparing model results to existing cluster groups)
  - Evaluating how well the results of the analysis fit the data without reference to external information
  - Comparing the results of different cluster models to determine which is better
  - Determining the best number of clusters

## Software Usability Criteria:

- Speed and ease of model development – so the customer can build new models and update existing models quickly
- Speed and ease of model deployment – so the customer can create Apply datasets easily and deploy models quickly with the required outputs (probabilities, deciles, etc.) – speed to market
- Ease of model maintenance – so the customer can easily define when models require refreshing/rebuilding and undertake this quickly and easily
- Integration capability with other systems

# Thank you

**Contact information:**

**open@sap.com**

# © 2016 SAP SE or an SAP affiliate company. All rights reserved.

Week 1 Unit 5: **Data Understanding Phase – Overview**

# Data Understanding Phase – Overview
## CRISP-DM – Phase 2: Data Understanding

# Data Understanding Phase – Overview
## Phase 2.1: Collect Initial Data

- **Task**
  - Acquire the data (or access to the data) listed in the project resources.
  - This initial collection includes data loading into the data exploration tool and data integration if multiple data sources are acquired.

- **Output – Initial Data Collection Report**
  - List the following:
    - The dataset (or datasets) acquired
    - The dataset locations
    - The methods used to acquire the datasets
    - Any problems encountered
  - Record problems encountered and any solutions.

# Data Understanding Phase – Overview
## Phase 2.2: Describe Data

- **Task**
  - Examine the "gross" or "surface" properties of the acquired data and report on the results.

- **Output – Data Description Report**
  - Describe the data that has been acquired, including:
    - The format of the data.
    - The quantity of data, e.g. the number of records and fields in each table.
    - The identities of the fields.
    - Any other surface features of the data that have been discovered.

# Data Understanding Phase – Overview
## Phase 2.3: Explore Data

- **Task**
  - This task tackles the data mining questions, which can be addressed using querying, visualization, and reporting.

- **Output – Data Exploration Report**
  - Describe results of this task including:
    - First findings or initial hypothesis and their impact on the remainder of the project.
    - If appropriate, include graphs and plots.

# Data Understanding Phase – Overview
## Phase 2.4: Verify Data Quality

- **Task**
  - Examine the quality of the data, addressing questions such as:
    - Is the data complete?
    - Is it correct or does it contain errors?
    - Are there missing values in the data?

- **Output – Data Quality Report**
  - List the results of the data quality verification.
  - If quality problems exist, list possible solutions.

# Thank you

**Contact information:**

**open@sap.com**

# © 2016 SAP SE or an SAP affiliate company. All rights reserved.

# Week 1 Unit 6: Initial Data Analysis & Exploratory Data Analysis

# Initial Data Analysis & Exploratory Data Analysis
## Initial data analysis

- *"Initial data analysis (IDA) is an essential part of nearly every analysis"*
  Problem Solving, A Statisticians Guide
  Christopher Chatfield

- Chatfield defines the various steps in IDA. It includes analysis of:
  - The structure of the data
  - The quality of the data
    - errors, outliers, and missing observations
  - Descriptive statistics
  - Graphs
- The data are modified according to the analysis:
  - Adjust extreme observations, estimate missing observations, transform variables, bin data, form new variables.



PROBLEM SOLVING
A statistician's guide

Second edition

Chris Chatfield

Texts in Statistical Science

CHAPMAN & HALL/CRC

# Initial Data Analysis & Exploratory Data Analysis
## Exploratory data analysis

- Exploratory data analysis (EDA) is an approach to analyzing data for the purpose of formulating hypotheses that are worth testing, and complements the tools of conventional statistics for testing hypotheses.

- It was so named by John Tukey.
  Wikipedia

- *"It is important to understand what you CAN DO before you learn to measure how WELL you seem to have done it."*

- *"To learn about data analysis, it is right that each of us try many things that do not work – that we tackle more problems than we make expert analyses of. We often learn less from an expertly done analysis than from one where, by not trying something, we missed an opportunity to learn more."*
  John Tukey, Exploratory Data Analysis

# Initial Data Analysis & Exploratory Data Analysis
## Example – US Stores data demonstration

- Demonstration using the SAP Predictive Analytics expert system.

- We will use US Stores retail data to walk you through this topic.

- The dataset contains the following variables:

  - Store location

  - Turnover

  - Margin

  - Staff

  - Store size

# Initial Data Analysis & Exploratory Data Analysis
## Example – Data visualization

# Initial Data Analysis & Exploratory Data Analysis
## Example – Scatter plot matrix

# Initial Data Analysis & Exploratory Data Analysis
## Example – Bubble plot

# Initial Data Analysis & Exploratory Data Analysis
## Example – Parallel co-ordinate plot

# Initial Data Analysis & Exploratory Data Analysis
## Example – Box plot

# Initial Data Analysis & Exploratory Data Analysis
## Example – Statistical summary chart

# Thank you

**Contact information:**

**open@sap.com**

openSAP
open.sap.com

# © 2016 SAP SE or an SAP affiliate company. All rights reserved.