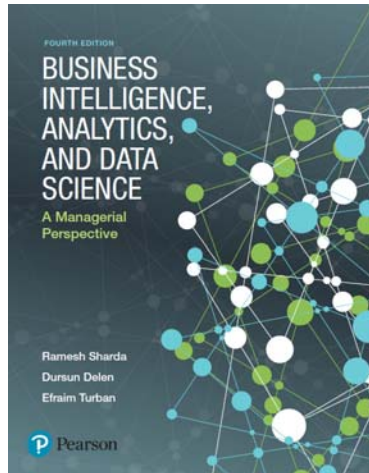


Business Intelligence, Analytics, and Data Science: A Managerial Perspective

Fourth Edition



Chapter 2

Descriptive Analytics I:
Nature of Data, Statistical
Modeling, and Visualization



Copyright © 2016, 2014, 2011 Pearson Education, Inc. All Rights Reserved

Learning Objectives (1 of 2)

- 2.1** Understand the nature of data as it relates to business intelligence (BI) and analytics
- 2.2** Learn the methods used to make real-world data analytics ready
- 2.3** Describe statistical modeling and its relationship to business analytics
- 2.4** Learn about descriptive and inferential statistics
- 2.5** Define business reporting, and understand its historical evolution



Slide 2-2

Copyright © 2016, 2014, 2011 Pearson Education, Inc. All Rights Reserved

Learning Objectives (2 of 2)

- 2.6** Understand the importance of data/information visualization
- 2.7** Learn different types of visualization techniques
- 2.8** Appreciate the value that visual analytics brings to business analytics
- 2.9** Know the capabilities and limitations of dashboards

OPENING VIGNETTE

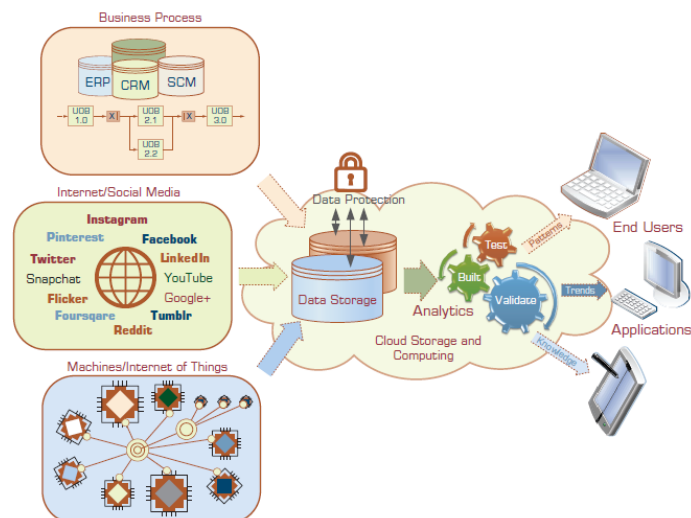
Attracts and Engages a New Generation of Radio Consumers with Data-Driven Marketing

1. What does SiriusXM do? In what type of market does it conduct its business?
2. What were the challenges? Comment on both technology and data-related challenges.
3. What were the proposed solutions?
4. How did they implement the proposed solutions? Did they face any implementation challenges?
5. What were the results and benefits? Were they worth the effort/investment?

The Nature of Data

- Data: a collection of facts
 - usually obtained as the result of experiences, observations, or experiments
- Data may consist of numbers, words, images, ...
- Data is the lowest level of abstraction (from which information and knowledge are derived)
- Data is the source for information and knowledge
- Data quality and data integrity → critical to analytics

The Nature of Data



Metrics for Analytics Ready Data

- Data source reliability
- Data content accuracy
- Data accessibility
- Data security and data privacy
- Data richness
- Data consistency
- Data currency/data timeliness
- Data granularity
- Data validity and data relevancy



Slide 2-7

Copyright © 2016, 2014, 2011 Pearson Education, Inc. All Rights Reserved

A Simple Taxonomy of Data

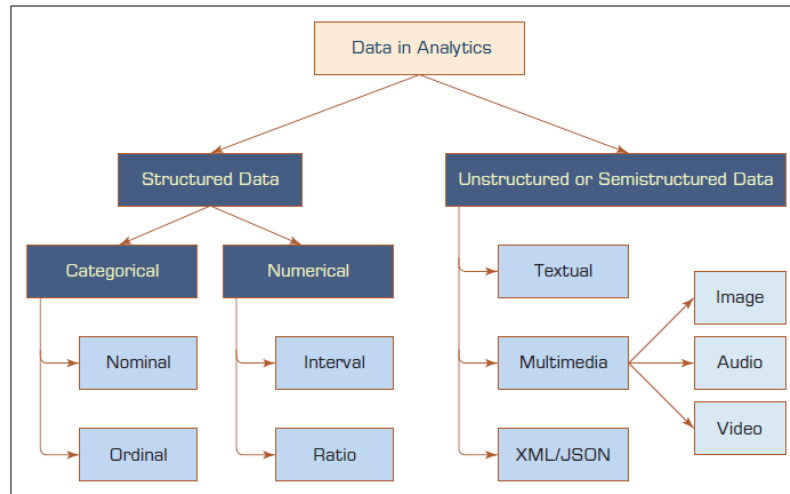
- Data (datum—singular form of data): facts
- Structured data
 - Targeted for computers to process
 - Numeric versus nominal
- Unstructured/textual data
 - Targeted for humans to process/digest
- Semi-structured data?
 - XML, HTML, Log files, etc.
- Data taxonomy...



Slide 2-8

Copyright © 2016, 2014, 2011 Pearson Education, Inc. All Rights Reserved

A Simple Taxonomy of Data



Application Case 2.1

Medical Device Company Ensures Product Quality While Saving Money

Questions for Discussion

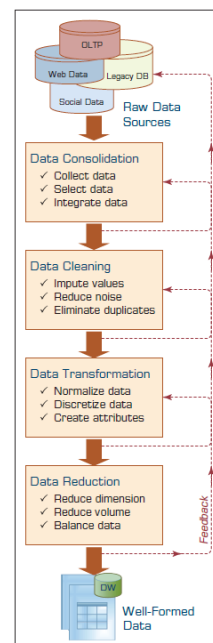
1. What were the main challenges for the medical device company? Were they market or technology driven?
2. What was the proposed solution?
3. What were the results? What do you think was the real return on investment (ROI)?

The Art and Science of Data Preprocessing

- The real-world data is dirty, misaligned, overly complex, and inaccurate
 - Not ready for analytics!
- Readying the data for analytics is needed
 - Data preprocessing
 - Data consolidation
 - Data cleaning
 - Data transformation
 - Data reduction
- Art – it develops and improves with experience

The Art and Science of Data Preprocessing

- Data reduction
 1. Variables
 - Dimensional reduction
 - Variable selection
 2. Cases/samples
 - Sampling
 - Balancing / stratification



Data Preprocessing Tasks and Methods

TABLE 2.1 A Summary of Data Preprocessing Tasks and Potential Methods

Main Task	Subtasks	Popular Methods
Data consolidation	Access and collect the data	SQL queries, software agents, Web services.
	Select and filter the data	Domain expertise, SQL queries, statistical tests.
	Integrate and unify the data	SQL queries, domain expertise, ontology-driven data mapping.
Data cleaning	Handle missing values in the data	Fill in missing values (imputations) with most appropriate values (mean, median, min/max, mode, etc.); recode the missing values with a constant such as "ML"; remove the record of the missing value; do nothing.
	Identify and reduce noise in the data	Identify the outliers in data with simple statistical techniques (such as averages and standard deviations) or with cluster analysis; once identified, either remove the outliers or smooth them by using binning, regression, or simple averages.
	Find and eliminate erroneous data	Identify the erroneous values in data (other than outliers), such as odd values, inconsistent class labels, odd distributions; once identified, use domain expertise to correct the values or remove the records holding the erroneous values.
Data transformation	Normalize the data	Reduce the range of values in each numerically valued variable to a standard range (e.g., 0 to 1 or -1 to +1) by using a variety of normalization or scaling techniques.
	Discretize or aggregate the data	If needed, convert the numeric variables into discrete representations using range- or frequency-based binning techniques; for categorical variables, reduce the number of values by applying proper concept hierarchies.
	Construct new attributes	Derive new and more informative variables from the existing ones using a wide range of mathematical functions (as simple as addition and multiplication or as complex as a hybrid combination of log transformations).
Data reduction	Reduce number of attributes	Principal component analysis, independent component analysis, chi-square testing, correlation analysis, and decision tree induction.
	Reduce number of records	Random sampling, stratified sampling, expert-knowledge-driven purposeful sampling.
	Balance skewed data	Oversample the less represented or undersample the more represented classes.

Application Case 2.2 (1 of 4)

Improving Student Retention with Data-Driven Analytics

Questions for Discussion

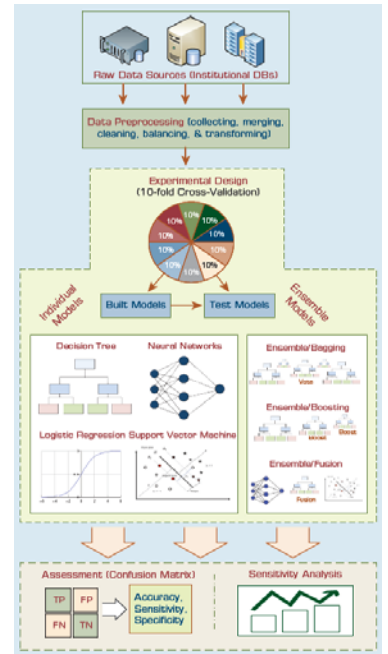
1. What is student attrition, and why is it an important problem in higher education?
2. What were the traditional methods to deal with the attrition problem?
3. List and discuss the data-related challenges within context of this case study.
4. What was the proposed solution? And, what were the results?

Application Case 2.2

Improving Student Retention with Data-Driven Analytics

(2 of 4)

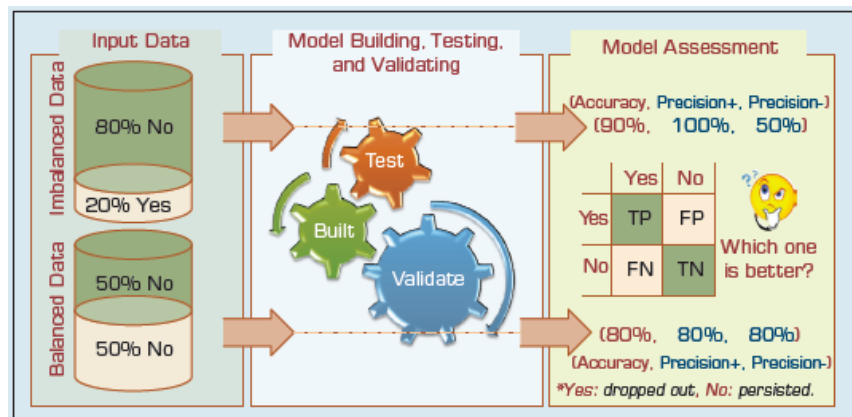
- Student retention
 - Freshmen class
- Why it is important?
- What are the common techniques to deal with student attrition?
- Analytics versus theoretical approaches to student retention problem



Application Case 2.2 (3 of 4)

Improving Student Retention with Data-Driven Analytics

- Data imbalance problem



Application Case 2.2 (4 of 4) • Results... Improving Student Retention with Data-Driven Analytics

TABLE 2.2 Prediction Results for the Original/Unbalanced Dataset

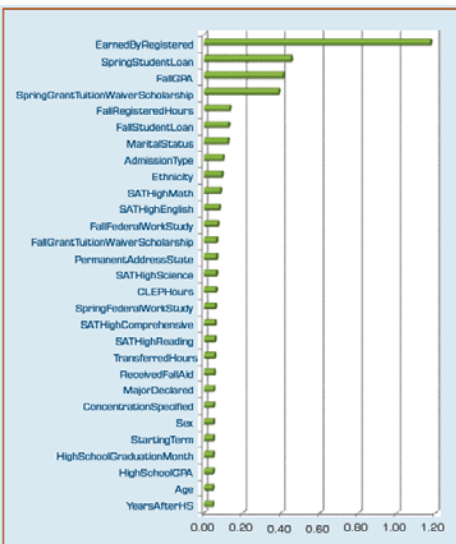
	ANN(MLP)		DT(C5)		SVM		LR	
	No	Yes	No	Yes	No	Yes	No	Yes
No	1494	384	1518	304	1478	255	1438	376
Yes	1596	11142	1572	11222	1612	11271	1652	11150
SUM	3090	11526	3090	11526	3090	11526	3090	11526
Per-Class Accuracy	48.35%	96.67%	49.13%	97.36%	47.83%	97.79%	46.54%	96.74%
Overall Accuracy	86.45%		87.16%		87.23%		86.12%	

TABLE 2.3 Prediction Results for the Balanced Data Set

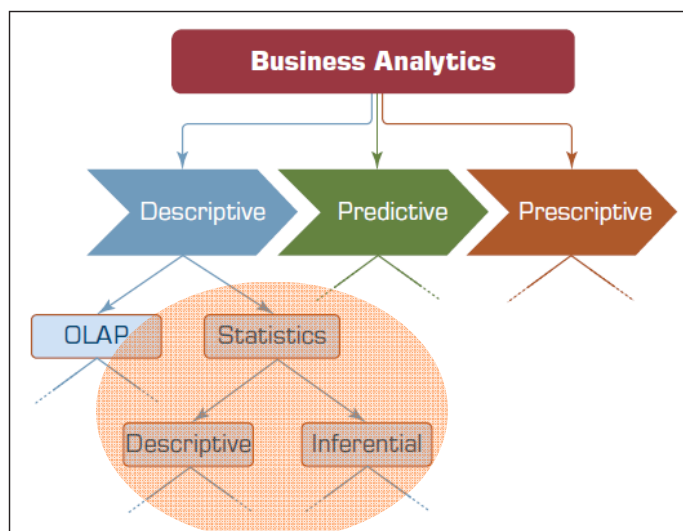
Confusion Matrix	ANN(MLP)		DT(C5)		SVM		LR	
	No	Yes	No	Yes	No	Yes	No	Yes
No	2309	464	2311	417	2313	386	2125	626
Yes	781	2626	779	2673	777	2704	965	2464
SUM	3090	3090	3090	3090	3090	3090	3090	3090
Per-Class Accuracy	74.72%	84.98%	74.79%	86.50%	74.85%	87.51%	68.77%	79.74%
Overall Accuracy	79.85%		80.65%		81.18%		74.26%	

TABLE 2.4 Prediction Results for the Three Ensemble Models

	Boosting		Bagging		Information Fusion	
	(Boosted Trees)		(Random Forest)		(Weighted Average)	
No	No	Yes	No	Yes	No	Yes
No	2242	375	2327	362	2335	351
Yes	948	2715	763	2738	755	2729
SUM	3090	3090	3090	3090	3090	3090
Per-Class Accuracy	72.56%	87.86%	75.31%	88.28%	75.57%	88.64%
Overall Accuracy	80.21%		81.80%		82.10%	



Statistical Modeling for Business Analytics



Statistical Modeling for Business Analytics

- **Statistics**
 - A collection of mathematical techniques to characterize and interpret data
- **Descriptive Statistics**
 - Describing the data (as it is)
- **Inferential statistics**
 - Drawing inferences about the population based on sample data
- Descriptive statistics for descriptive analytics

Descriptive Statistics Measures of Centrality Tendency

- **Arithmetic mean**

$$\bar{x} = \frac{x_1 + x_2 + \cdots + x_n}{n} \qquad \bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

- **Median**
 - The number in the middle
- **Mode**
 - The most frequent observation

Descriptive Statistics Measures of Dispersion

- **Dispersion**
 - Degree of variation in a given variable
- **Range**
 - Max - Min

- **Variance**

$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}$$

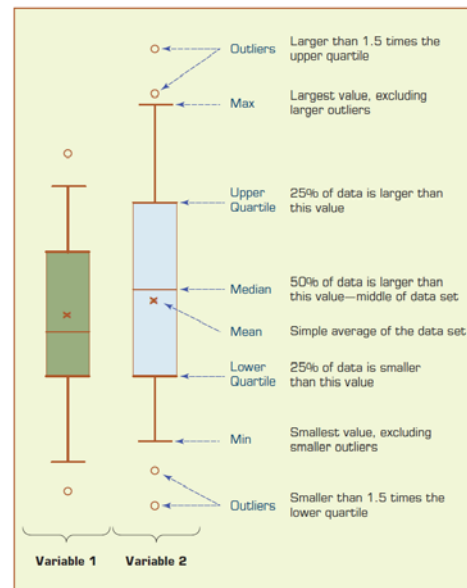
- **Standard Deviation**

$$s = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n - 1}}$$

- **Mean Absolute Deviation (MAD)**
 - Average absolute deviation from the mean

Descriptive Statistics Measures of Dispersion

- **Quartiles**
- **Box-and-Whiskers Plot**
 - a.k.a. box-plot
 - Versatile / informative



Descriptive Statistics

Shape of a Distribution

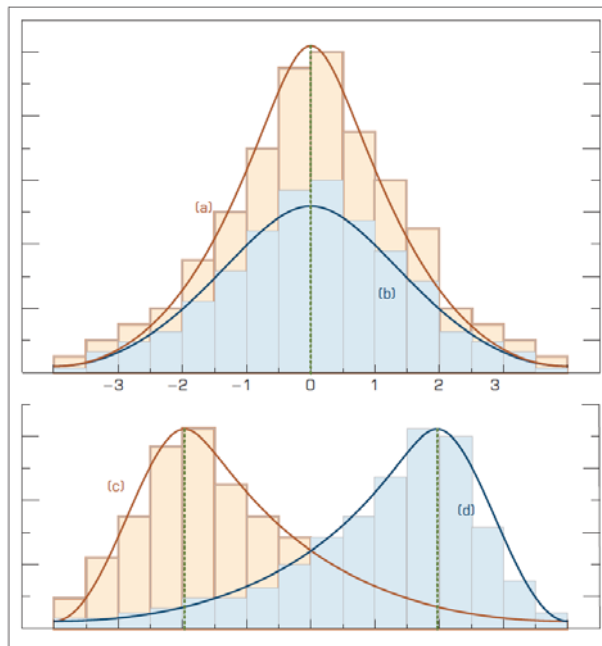
- **Histogram** – frequency chart
- **Skewness**
 - Measure of asymmetry

$$\text{Skewness} = S = \frac{\sum_{i=1}^n (x_i - \bar{x})^3}{(n - 1)s^3}$$

- **Kurtosis**
 - Peak/tall/skinny nature of the distribution

$$\text{Kurtosis} = K = \frac{\sum_{i=1}^n (x_i - \bar{x})^4}{ns^4} - 3$$

Relationship Between Dispersion and Shape Properties



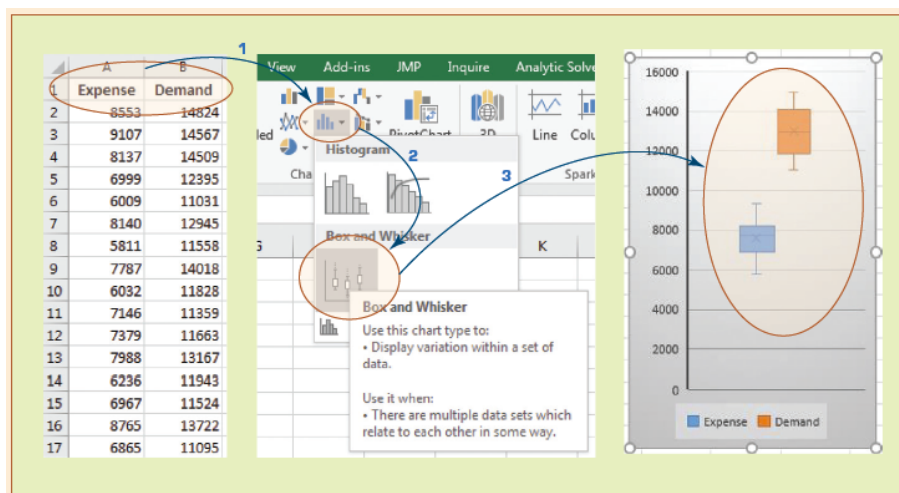
Technology Insights 2.1 – Descriptive Statistics in Excel

The screenshots illustrate the process of enabling the Data Analysis ToolPak in Microsoft Excel. The first three steps show navigating through the File menu to Options, then the Add-ins section, and finally selecting the Analysis ToolPak from the list of available add-ins. The fourth screenshot shows the Descriptive Statistics dialog box, where the input range is set to \$A\$1:\$B\$17 and the output range is set to \$D\$1:\$F\$17. The dialog box also includes options for output range, new worksheet, and various statistical outputs like summary statistics, confidence intervals, and kurtosis.

Expense	Demand
8553	14824
9107	14567
8137	14509
6999	12395
6009	11031
8140	12945
5811	11558
7787	14018
6032	11828
7146	11359
7988	13167
6236	11943
6967	11524
8765	13722
6865	11095

Technology Insights 2.1 – Descriptive Statistics in Excel

Creating box-plot in Microsoft Excel



Application Case 2.3

Town of Cary Uses Analytics to Analyze Data from Sensors, Assess Demand, and Detect Problems

Questions for Discussion

1. What were the challenges the Town of Cary was facing?
2. What was the proposed solution?
3. What were the results?
4. What other problems and data analytics solutions do you foresee for towns like Cary?



Slide 2-27

Copyright © 2015, 2014, 2011 Pearson Education, Inc. All Rights Reserved

Regression Modeling for Inferential Statistics

- **Regression**
 - A part of inferential statistics
 - The most widely known and used analytics technique in statistics
 - Used to characterize relationship between explanatory (input) and response (output) variable
- It can be used for
 - Hypothesis testing (explanation)
 - Forecasting (prediction)



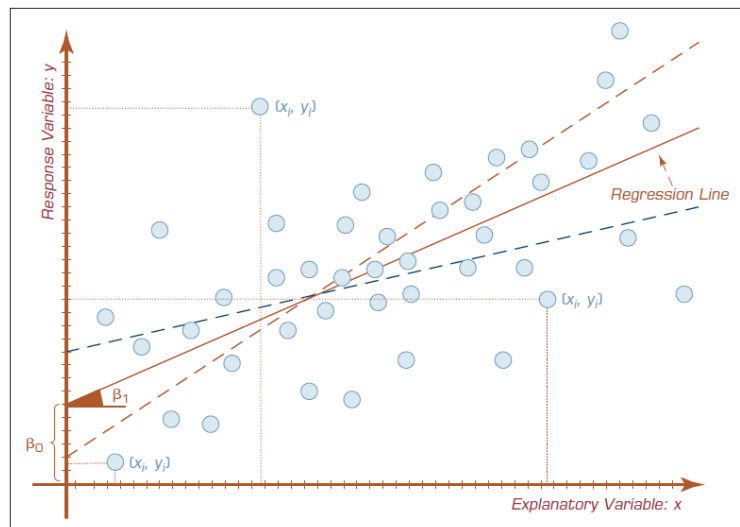
Slide 2-28

Copyright © 2015, 2014, 2011 Pearson Education, Inc. All Rights Reserved

Regression Modeling

- Correlation versus Regression
 - What is the difference (or relationship)?
- Simple Regression versus Multiple Regression
 - Base on number of input variables
- How do we develop linear regression models?
 - Scatter plots (visualization—for simple regression)
 - Ordinary least squares method
 - A line that minimizes squared of the errors

Regression Modeling



Regression Modeling

- x : input, y : output
- Simple Linear Regression

$$y = \beta_0 + \beta_1 x$$

- Multiple Linear Regression

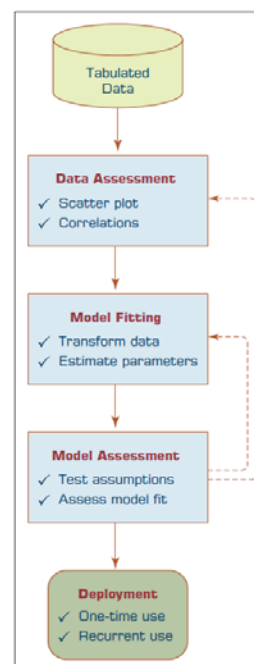
$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots + \beta_n x_n$$

- The meaning of Beta (β) coefficients
 - Sign (+ or -) and magnitude

Process of Developing a Regression Model

How do we know if the model is good enough?

- R^2 (R-Square)
- p Values
- Error measures (for prediction problems)
 - MSE, MAD, RMSE



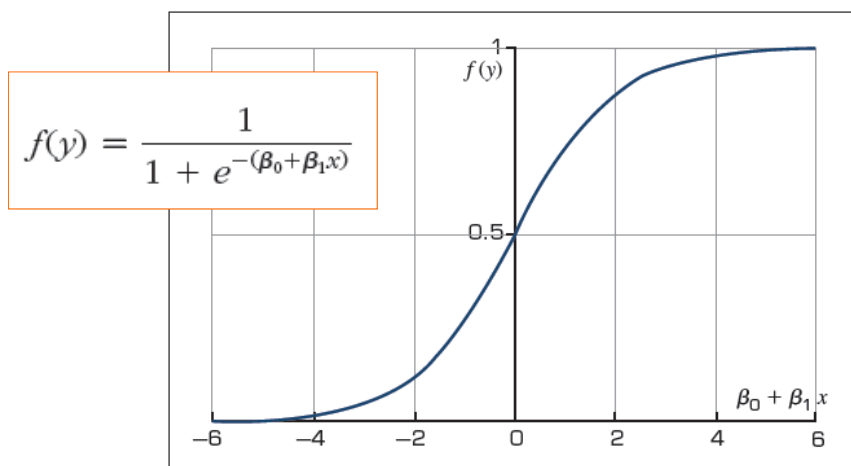
Regression Modeling Assumptions

- Linearity
- Independence
- Normality (Normal Distribution)
- Constant Variance
- Multicollinearity
- What happens if the assumptions do NOT hold?
 - What do we do then?

Logistic Regression Modeling

- A very popular statistics-based classification algorithm
- Employs supervised learning
- Developed in 1940s
- The difference between Linear Regression and Logistic Regression
 - In Logistic Regression Output/Target variable is a binomial (binary classification) variable (as opposed to numeric variable)

Logistic Regression Modeling



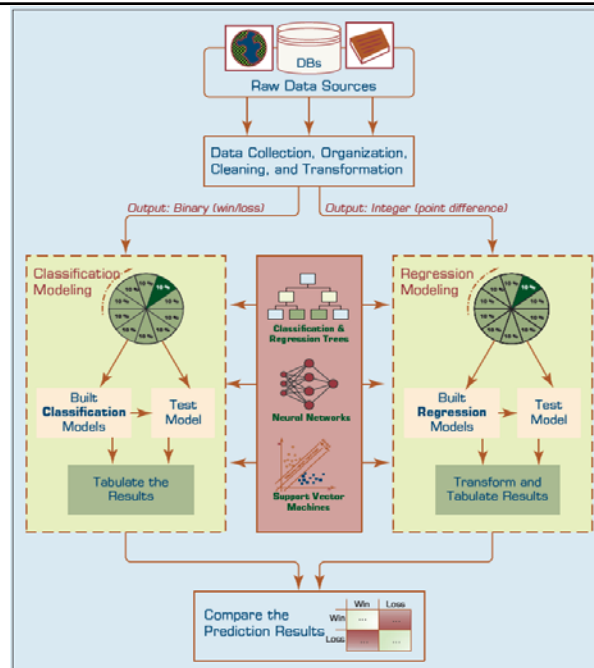
Application Case 2.4 (1 of 4)

Predicting NCAA Bowl Game Outcomes



Application Case 2.4 (2 of 4) Predicting NCAA Bowl Game Outcomes

- The analytics process to develop prediction models (both regression and classification type) for NCAA Bowl Game outcomes



Application Case 2.4 (3 of 4) Predicting NCAA Bowl Game Outcomes

Prediction Results

- Classification
- Regression

Prediction Method (Classification)		Confusion Matrix		Accuracy** (in %)	Sensitivity (in %)	Specificity (in %)
		Win	Loss			
ANN (MLP)	Win	92	42	75.00	68.66	82.73
	Loss	19	91			
SVM (RBF)	Win	105	29	79.51	78.36	80.91
	Loss	21	89			
DT (C&RT)	Win	113	21	86.48	84.33	89.09
	Loss	12	98			

*The output variable is a binary categorical variable (Win or Loss); differences were sig (**p < 0.01).

Prediction Method (Regression-Based*)		Confusion Matrix		Accuracy**	Sensitivity	Specificity
		Win	Loss			
ANN (MLP)	Win	94	40	72.54	70.15	75.45
	Loss	27	83			
SVM (RBF)	Win	100	34	74.59	74.63	74.55
	Loss	28	82			
DT (C&RT)	Win	106	28	77.87	76.36	79.10
	Loss	26	84			

*The output variable is a numerical/integer variable (point-diff); differences were sig (**p < 0.01).

Application Case 2.4 (4 of 4)

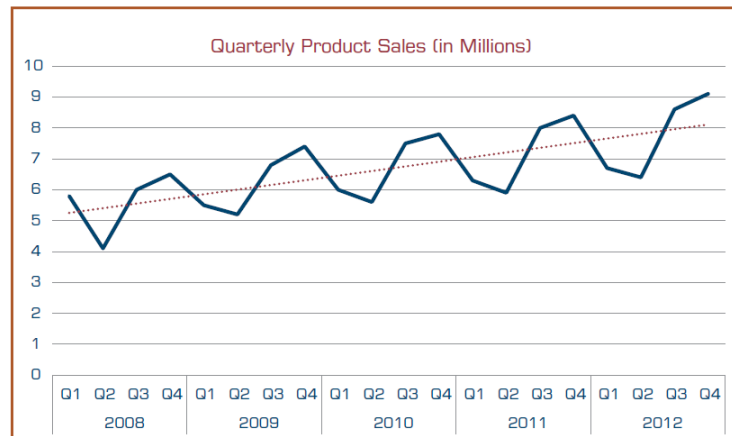
Predicting NCAA Bowl Game Outcomes

Questions for Discussion

1. What are the foreseeable challenges in predicting sporting event outcomes (e.g., college bowl games)?
2. How did the researchers formulate/design the prediction problem (i.e., what were the inputs and output, and what was the representation of a single sample—row of data)?
3. How successful were the prediction results? What else can they do to improve the accuracy?

Time Series Forecasting

- Is it different than Simple Linear Regression? How?



Business Reporting Definitions and Concepts

- Report = Information → Decision
- Report?
 - Any communication artifact prepared to convey specific information
- A report can fulfill many functions
 - To ensure proper departmental functioning
 - To provide information
 - To provide the results of an analysis
 - To persuade others to act
 - To create an organizational memory...



Slide 2-41

Copyright © 2016, 2014, 2011 Pearson Education, Inc. All Rights Reserved

What is a Business Report?

- A written document that contains information regarding business matters.
- **Purpose:** to improve managerial decisions
- **Source:** data from inside and outside the organization (via the use of ETL)
- **Format:** text + tables + graphs/charts
- **Distribution:** in-print, email, portal/intranet

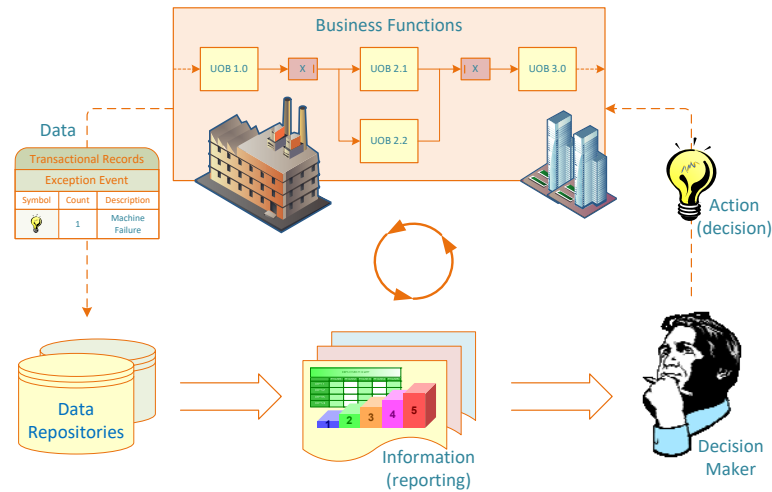
Data acquisition → Information generation → Decision making → Process management



Slide 2-42

Copyright © 2016, 2014, 2011 Pearson Education, Inc. All Rights Reserved

Business Reporting



Types of Business Reports

- **Metric Management Reports**
 - Help manage business performance through metrics (SLAs for externals; KPIs for internals)
 - Can be used as part of Six Sigma and/or TQM
- **Dashboard-Type Reports**
 - Graphical presentation of several performance indicators in a single page using dials/gauges
- **Balanced Scorecard–Type Reports**
 - Include financial, customer, business process, and learning & growth indicators

Application Case 2.5

Flood of Paper Ends at FEMA

Questions for Discussion

1. What is FEMA, and what does it do?
2. What are the main challenges that FEMA faces?
3. How did FEMA improve its inefficient reporting practices?



Slide 2-45

Copyright © 2015, 2014, 2011 Pearson Education, Inc. All Rights Reserved

Data Visualization

“The use of visual representations to explore, make sense of, and communicate data.”

- Data visualization vs. Information visualization
- Information = aggregation, summarization, and contextualization of data
- Related to information graphics, scientific visualization, and statistical graphics
- Often includes charts, graphs, illustrations, ...



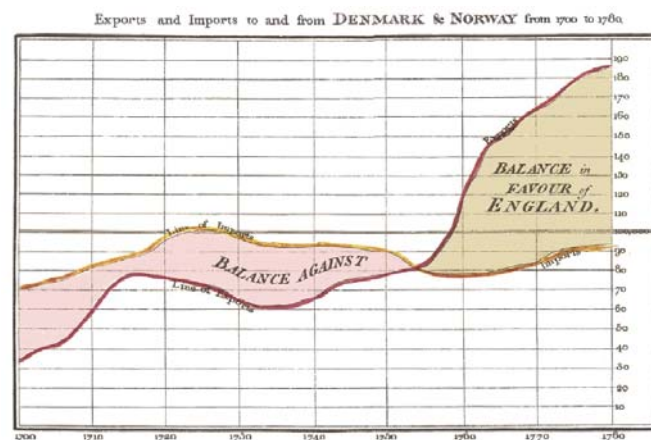
Slide 2-46

Copyright © 2015, 2014, 2011 Pearson Education, Inc. All Rights Reserved

A Brief History of Data Visualization

- Data visualization can date back to the second century AD
- Most developments have occurred in the last two and a half centuries
- Until recently it was not recognized as a discipline
- Today's most popular visual forms date back a few centuries

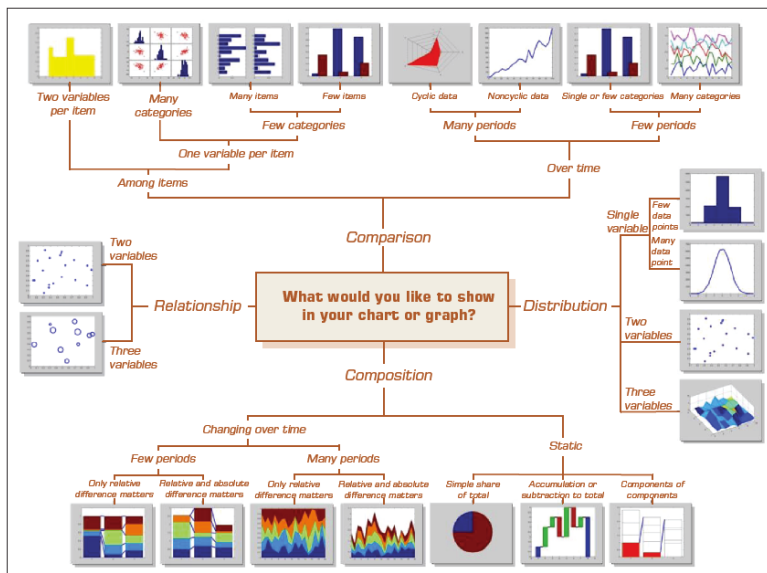
The First Pie Chart Created by William Playfair in 1801



William Playfair is widely credited as the inventor of the modern chart, having created the first line and pie charts.

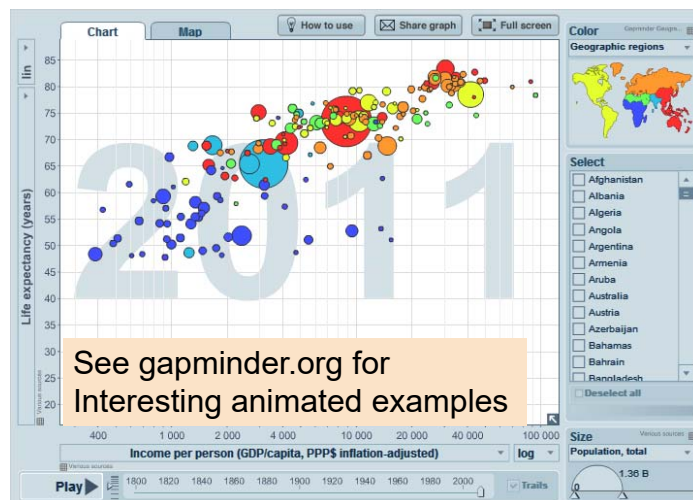
The Bottom line is divided into Years, the Right hand line into £10,000 each.
Published at the Art Office, 17 May 1801, by W^m Playfair.

Which Chart or Graph Should You Use?



An Example Gapminder Chart

Wealth and Health of Nations



The Emergence of Data Visualization and Visual Analytics

- Magic Quadrant for Business Intelligence and Analytics Platforms (Source: Gartner.com)
- Many data visualization companies are in the 4th quadrant
- There is a move towards visualization



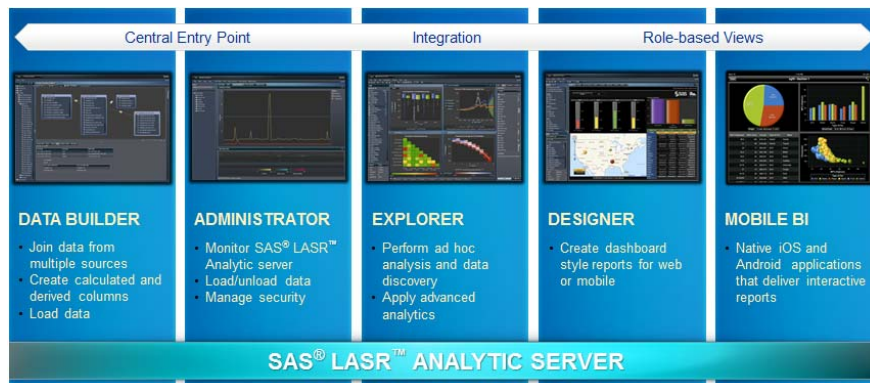
The Emergence of Data Visualization and Visual Analytics

- Emergence of new companies
 - Tableau, Spotfire, QlikView, ...
- Increased focus by the big players
 - MicroStrategy improved Visual Insight
 - SAP launched Visual Intelligence
 - SAS launched Visual Analytics
 - Microsoft bolstered PowerPivot with Power View
 - IBM launched Cognos Insight
 - Oracle acquired Endeca

Visual Analytics

- A recently coined term
 - Information visualization + predictive analytics
- Information visualization
 - Descriptive, backward focused
 - “what happened” “what is happening”
- Predictive analytics
 - Predictive, future focused
 - “what will happen” “why will it happen”
- There is a strong move toward **visual analytics**

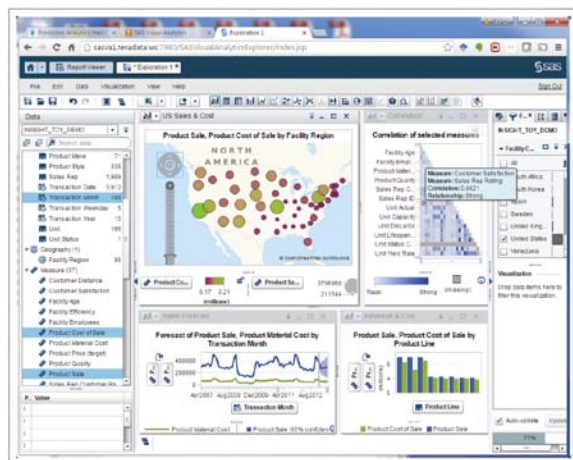
Visual Analytics by SAS Institute



- SAS Visual Analytics Architecture
 - Big data + In memory + Massively parallel processing + ..

Visual Analytics by SAS Institute

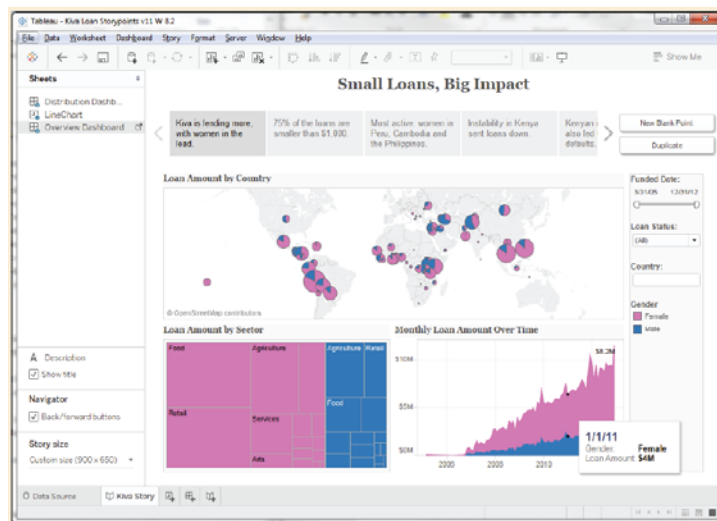
- At teradatauniversitynetwork.com, you can learn more about SAS VA, experiment with the tool



Slide 2-57

Copyright © 2016, 2014, 2011 Pearson Education, Inc. All Rights Reserved

Technology Insight 2.3 Telling Great Stories with Data and Visualization



Slide 2-58

Copyright © 2016, 2014, 2011 Pearson Education, Inc. All Rights Reserved

Performance Dashboards

- Performance dashboards are commonly used in BPM software suites and BI platforms
- Dashboards provide visual displays of important information that is consolidated and arranged on a single screen so that information can be digested at a single glance and easily drilled in and further explored

Performance Dashboards



Application Case 2.7

Dallas Cowboys Score Big with Tableau and Teknion

Questions for Discussion

1. How did the Dallas Cowboys use information visualization?
2. What were the challenge, the proposed solution, and the obtained results?



Slide 2-61

Copyright © 2016, 2014, 2011 Pearson Education, Inc. All Rights Reserved

Performance Dashboards

- Dashboard design
 - The fundamental challenge of dashboard design is to display all the required information on a single screen, clearly and without distraction, in a manner that can be assimilated quickly
- Three layer of information
 - Monitoring
 - Analysis
 - Management



Slide 2-62

Copyright © 2016, 2014, 2011 Pearson Education, Inc. All Rights Reserved

Performance Dashboards

- What to look for in a dashboard
 - Use of visual components to highlight data and exceptions that require action
 - Transparent to the user, meaning that they require minimal training and are extremely easy to use
 - Combine data from a variety of systems into a single, summarized, unified view of the business
 - Enable drill-down or drill-through to underlying data sources or reports
 - Present a dynamic, real-world view with timely data
 - Require little coding to implement, deploy, and maintain

Best Practices in Dashboard Design

- Benchmark KPIs with Industry Standards
- Wrap the Metrics with Contextual Metadata
- Validate the Design by a Usability Specialist
- Prioritize and Rank Alerts and Exceptions
- Enrich Dashboard with Business-User Comments
- Present Information in Three Different Levels
- Pick the Right Visual Constructs
- Provide for Guided Analytics

Application Case 2.8

Visual Analytics Helps Energy Supplier Make Better Connections

Questions for Discussion

1. Why do you think energy supply companies are among the prime users of information visualization tools?
2. How did Electrabel use information visualization for the single version of the truth?
3. What were their challenges, the proposed solution, and the obtained results?



Slide 2-65

Copyright © 2016, 2014, 2011 Pearson Education, Inc. All Rights Reserved

End of Chapter 2

- Questions / Comments



Slide 2-66

Copyright © 2016, 2014, 2011 Pearson Education, Inc. All Rights Reserved