

CSC 5741

Lecture 3: Data Mining and Data Processing

Lighton Phiri <lighton.phiri@unza.zm>
Department of Library and Information Science
University of Zambia

Lecture Series Outline

- Part I: Data Mining
- Part II: Data Processing and Transformation
- Part III: Paper Reading Discussion
- Part IV: Academic Talk

Lecture Series Outline

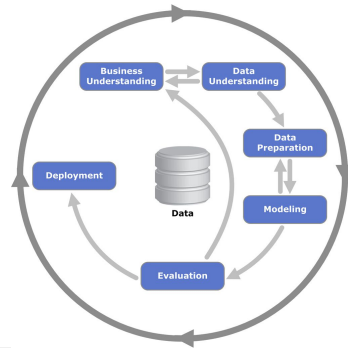
- Part I: Data Mining
 - Introduction
 - Data Mining Process
 - Data Mining Process Example
- Part II: Data Processing and Transformation
- Part III: Paper Reading
- Part IV: Academic Talk

Introduction

- Similar to the motivation for using computer systems, data mining involves processing of input data to yield information
 - Increasingly, massive amounts of data are being frequently produced
 - Raw data does not provide useful insight in comparison to information

CRISP-DM Open Standard (1/4)

- The Cross-industry standard process for data mining (CRISP-DM) is a model commonly used to highlight approaches in data mining
 - CRISP-DM segments a data mining project into six phases with no strict order of execution
 - Surveys conducted suggest CRISP-DM is the most widely used methodology

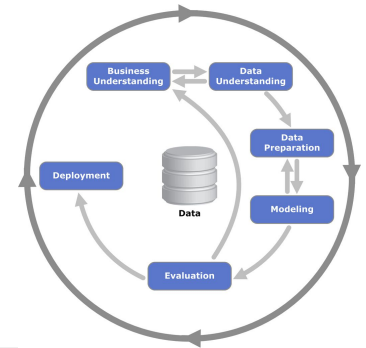


April 9 2019

CSC 5741 L03 - 5

CRISP-DM Open Standard (2/4)

- Business Understanding**
 - Situational analysis; problem definition, general and specific objectives objectives; **research question(s)** and general requirements analysis
- Data Understanding**
 - Identification of data sources; familiarisation of data sources and initial data collection

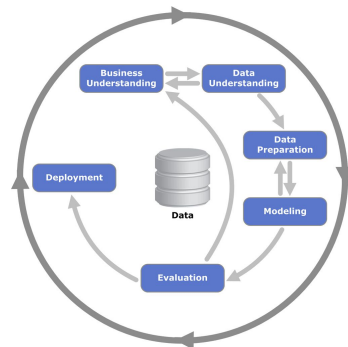


April 9 2019

CSC 5741 L03 - 6

CRISP-DM Open Standard (3/4)

- Data Preparation**
 - Data preprocessing; data cleaning and feature selection
- Modeling**
 - Creation of model—probably machine learning model—using data mining tools
- Evaluation**
 - Evaluation results against goals
- Deployment**
 - Deployment of models

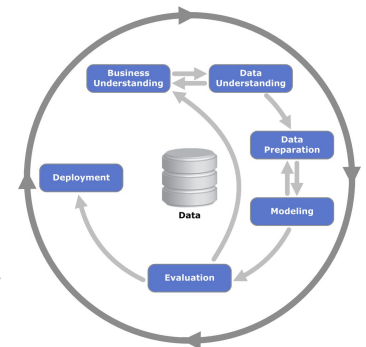


April 9 2019

CSC 5741 L03 - 7

CRISP-DM Open Standard (4/4)

- While CRISP-DM is the most widely used model [2], other data mining process models that exist include:
 - Analytics Solutions Unified Method for Data Mining/Predictive Analytics (ASUM-DM)
 - Sample, Explore, Modify, Model, and Assess (SEMMA)

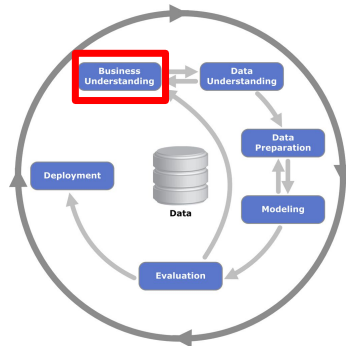


April 9 2019

CSC 5741 L03 - 8

CRISP-DM—Business Understanding (1/)

- Outline business and data mining goals and objectives
- Conduct a situational analysis to identify how problem is current resolved
- Prepare an overall project plan

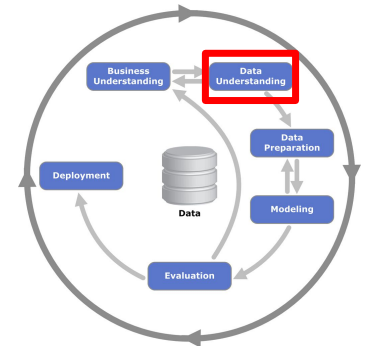


April 9 2019

CSC 5741 L03 - 9

CRISP-DM—Data Understanding (1/)

- Identify data sources
- Extract/collect required data
- Described and explore the data collected to gain some sense of what insights to derive
- Ascertain quality of data collected



April 9 2019

CSC 5741 L03 - 10

CRISP-DM—Data Preparation (1/)

- Select data required for modeling process/phase
- Clean the data
- Reconstruct the data and derive necessary attributes
- Merge different data sources
- Reformat the data
 - e.g. Naming conventions

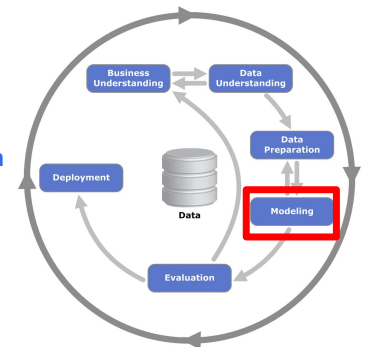


April 9 2019

CSC 5741 L03 - 11

CRISP-DM—Modeling (1/)

- Define the model components, features, how it behaves and how to interpret it
- Evaluate the various alternative techniques that can be integrated with the model
 - e.g. Evaluate different classification algorithms

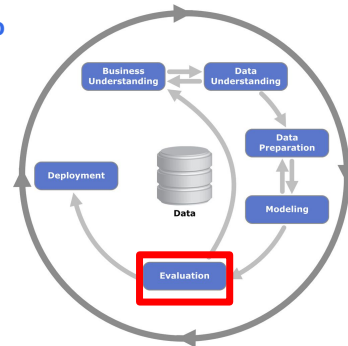


April 9 2019

CSC 5741 L03 - 12

CRISP-DM—Evaluation (1/)

- Devise evaluation techniques to be used
 - Efficiency vs effectiveness/efficacy
- Interpret model results to ascertain if model should be deployed
- Review the process if necessary

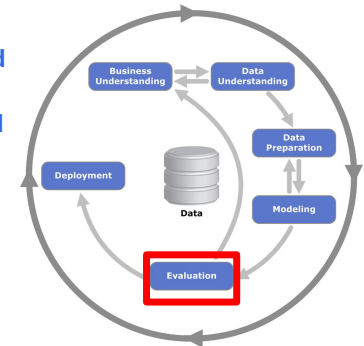


April 9 2019

CSC 5741 L03 - 13

CRISP-DM—Deployment (1/)

- Determine how the model results will be presented to end users
- Identify end user that will need to use the model results



April 9 2019

CSC 5741 L03 - 14

CRISP-DM—Random Example (1/)

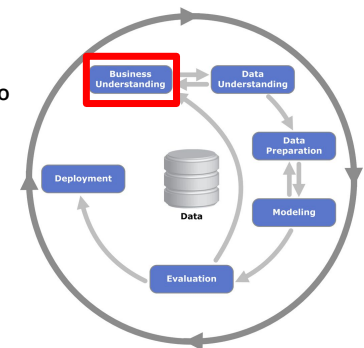
- ICT 1110 performance is bad. The poor performance transcends all the various assessments written by students: quizzes, tests and practical programming questions.

April 9 2019

CSC 5741 L03 - 15

CRISP-DM—Random Example (2/)

- Outline business and data mining goals and objectives
 - Monitor student performance to prevent poor performance
 - Identify at risk students and devise corrective measures
- Conduct a situational analysis to identify how problem is current resolved
 - How are at-risk students currently identified

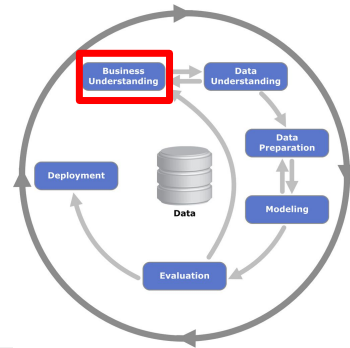


April 9 2019

CSC 5741 L03 - 16

CRISP-DM—Random Example (3/)

- Prepare an overall project plan
 - Timeline of project execution

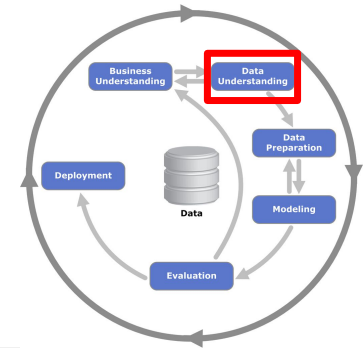


April 9 2019

CSC 5741 L03 - 17

CRISP-DM—Random Example (4/)

- Identify data sources
 - (i) assessment results (ii) student demographics (iii) student past experience (iv) Moodle interaction logs (v) tutorial attendance (vi) lecture attendance (vii) tutor feedback
- Extract/collect required data
 - (i) assessment results (ii) SIS extraction (iii) questionnaire??

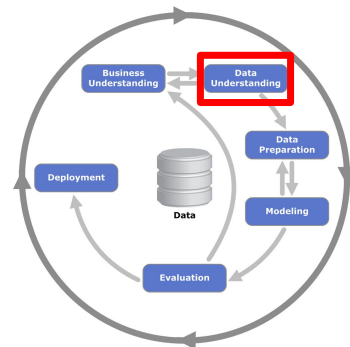


April 9 2019

CSC 5741 L03 - 18

CRISP-DM—Random Example (5/)

- Described and explore the data collected to gain some sense of what insights to derive
- Ascertain quality of data collected



April 9 2019

CSC 5741 L03 - 19

CRISP-DM—Data Preparation (1/)

- Some information for this basic example is not yet available for 2017/18
 - Lecture attendance
 - Tutorial attendance



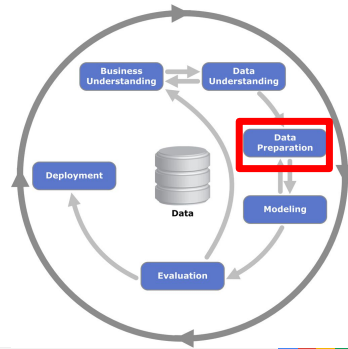
April 9 2019

CSC 5741 L03 - 20

CSC 5741 L03 - 24

CRISP-DM—Random Example (6/)

- **Select data required for modeling process/phase**
 - Will all the data sources be used?
- **Clean the data**
 - Normalise student names
 - Normalise their demographic details (e.g. Home Towns)

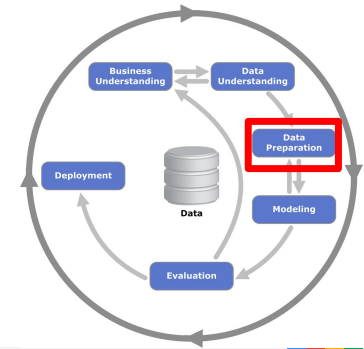


April 9 2019

CSC 5741 L03 - 25

CRISP-DM—Random Example (7/)

- **Reconstruct the data and derive necessary attributes**
 - Student ages perhaps?
- **Merge different data sources**
- **Reformat the data**
 - e.g. Naming conventions

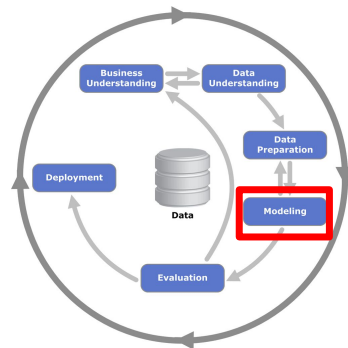


April 9 2019

CSC 5741 L03 - 26

CRISP-DM—Random Example (7/)

- **Define the model components, features, how it behaves and how to interpret it**
- **Evaluate the various alternative techniques that can be integrated with the model**
 - Evaluate potentially useful classification algorithms

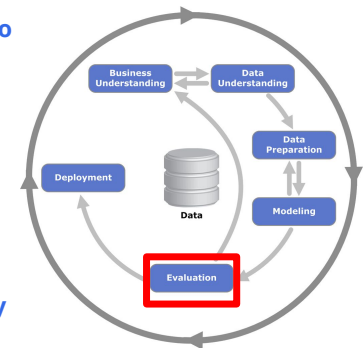


April 9 2019

CSC 5741 L03 - 27

CRISP-DM—Evaluation (8/)

- **Devise evaluation techniques to be used**
 - Evaluate classification model's effectiveness
 - Notice that efficiency is not important here
- **Interpret model results to ascertain if model should be deployed**
- **Review the process if necessary**

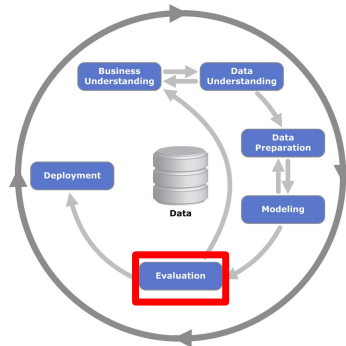


April 9 2019

CSC 5741 L03 - 28

CRISP-DM—Deployment (1/)

- Determine how the model results will be presented to end users
 - Implement a Web application that will be used to determine and present at-risk students
- Identify end user that will need to use the model results
 - Figure out if lecturers and tutors will have access to results. Perhaps HoD as well?



April 9 2019

CSC 5741 L03 - 29

Q & A Session

- Comments, concerns and complaints?

April 9 2019

CSC 5741 L03 - 30

Lecture Series Outline

- Part I: Data Mining
- Part II: Data Processing and Transformation
 - Data Collection and Cleaning
 - Transforming and merging data
- Part II: Paper Reading Discussion
- Part IV: Academic Talk

April 9 2019

CSC 5741 L03 - 31

Lecture Series Outline

- Part I: Data Mining
- Part II: Data Processing and Transformation
- Part II: Paper Reading Discussion
 - M. Mgala and A. Mbogho (2015) "Data-driven intervention-level prediction modeling for academic performance"
- Part IV: Academic Talk

April 9 2019

CSC 5741 L03 - 32

Paper Reading Session



University of Cape Town

My Author Page My Binders SIGN OUT:
Lighton Phiri

SEARCH

Check out a preview of the [next ACM DL](#)

Data-driven intervention-level prediction modeling for academic performance

Full Text: [PDF](#)

Authors: [Mvuruya Mgaala](#) [University of Cape Town, Cape Town](#)
[Audrey Mboogho](#) [University of Cape Town, Cape Town](#)



2015 Article

Published in:

- Proceeding
ICTD '15 Proceedings of the Seventh International Conference on
Information and Communication Technologies and Development
Article No. 2

Singapore, Singapore — May 15 - 18, 2015

ACM New York, NY, USA ©2015

[table of contents](#) ISBN: 978-1-4503-3163-0 doi> [10.1145/2737856.2738012](#)

[Bibliometrics](#)

Citation Count: 3
Downloads (cumulative): 152
Downloads (12 Months): 29
Downloads (6 Weeks): 3

Tools and Resources

[Request Permissions](#)

TOC Service:

[Email](#) [RSS](#)

[Save to Binder](#)

[View My Binders](#)

Export Formats:

[BibTeX](#) [EndNote](#) [ACM Ref](#)

Share:

[Facebook](#) [Twitter](#) [LinkedIn](#) [YouTube](#) [Google+](#)

Author Tags:

April 9 2019

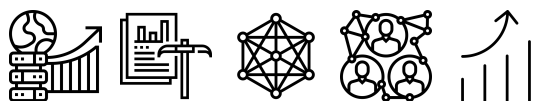
CSC 5741 L03 - 33

Bibliography

- [1] Witten, I. H., Frank, E., Hall, M. A., Pal, C. J. (2017) Data Mining: Practical Machine Learning Tools and Techniques. Chapter 1 <https://www.cs.waikato.ac.nz/ml/weka/book.html>
- [2] Kurgan, L. A. and Muselek, P. (2006) A Survey of Knowledge Discovery and Data Mining Process Models <https://doi.org/10.1017/S0269888906000737>

April 9 2019

CSC 5741 L03 - 34



CSC 5741 Lecture 3: Data Mining and Data Processing

Lighton Phiri <lighton.phiri@unza.zm>
Department of Library and Information Science
University of Zambia