

Digital Skills  
Training  
Programs at  
Knowledge  
Technology

INTRODUCTION  
TO NATURAL  
LANGUAGE  
PROCESSING

- Participants will be introduced to Natural Language Processing (NLP) and its applications through examples and exercises.
- Build intelligent applications that can interpret the human language to deliver impactful results



**KNOWLEDGE TECHNOLOGY**  
R E S E A R C H U N I T

<b>Course Title</b>	Introduction to Natural Language Processing
<b>Duration</b>	10 Days
<b>Trainer</b>	Assoc. Prof. Dr. Rayner Alfred
<b>Cost</b>	Email <a href="mailto:ralfred121@gmail.com">ralfred121@gmail.com</a> or call 013-881-9966 for quotations
<b>Max Participants</b>	25

**SYNOPSIS**

Participants will be introduced to Natural Language Processing (NLP) and its applications through examples and exercises. This will be followed by an introduction to the initial stages of solving a problem, which includes problem definition, getting text data, and preparing text data for modeling. With exposure to concepts such as advanced Natural Language Processing algorithms and visualization techniques, participants will learn how to create applications that can extract information from unstructured data and present it as impactful visuals. Although participants will continue to learn NLP-based techniques, the focus will gradually shift to developing useful applications. In those sections, participants will gain an understanding of how to apply Natural Language Processing techniques to answer questions, as can be used for chatbots.

**LEARNING OUTCOMES**

What participants will gain at the end of the course

- Obtain, verify, and clean data before transforming it into a correct format for use
- Perform data analysis and machine learning tasks using Python
- Gain an understanding of the basics of computational linguistics
- Build models for general NLP tasks
- Evaluate the performance of a model with the right metrics
- Visualize, quantify, and perform exploratory analysis from any text data

**JUSTIFICATION TO LEARN NATURAL LANGUAGE PROCESSING**

In the past years, the tech world has seen a surge of Natural Language Processing (NLP) applications in various areas, including adtech, publishing, customer service and market

intelligence. According to Gartner’s hype cycle, NLP has reached the peak of inflated expectations in 2018. Many businesses see it as a “go-to” solution to generate value from the 80% of business-relevant data that comes in unstructured form. To put it simply—NLP is wildly adopted with wildly variable success, it is a unique NLP component critical for the core business of our company and most important is that there is a large demand to have an internal competence to develop IP-relevant NLP technology.

This Natural Language Processing course is designed for novice and mid-level data scientists and machine learning developers who want to gather and analyze text data to build an NLP-powered product. It’ll help participants to have prior experience of coding in Python using data types, writing functions, and importing libraries. Some experience with linguistics and probability is useful but not necessary.

### TOPICS LIST

- [1] Introduction to Natural Language Processing
- [2] Basic Feature Extraction Methods
- [3] Developing a Text classifier
- [4] Collecting Text Data from the Web
- [5] Topic Modelling
- [6] Text Summarization and Text Generation
- [7] Vector Representation
- [8] Sentiment Analysis

### COURSE SYLLABUS (10 DAYS)

DAY	TOPICS COVERED	TIME
One	<b>MODULE 1: INTRODUCTION TO NATURAL LANGUAGE PROCESSING</b> <ul style="list-style-type: none"> <li>➤ History of NLP, Text Analytics and NLP</li> <li>➤ Various Steps in NLP</li> <li>➤ Tokenization, Part of Speech Tagging, Stop Word Removal, Text Normalization, Spelling Correction, Stemming, Lemmatization, Named-Entity Recognition (NER), Word Sense Disambiguation, Sentence Boundary Detection</li> <li>➤ Kick Starting an NLP Project</li> </ul>	8am – 5pm
Two	<b>MODULE 2: BASIC FEATURE EXTRACTION METHODS</b> <ul style="list-style-type: none"> <li>➤ Types of Data</li> <li>➤ Categorizing Data Based on Structure</li> <li>➤ Categorization of Data Based on Content</li> <li>➤ Cleaning Text Data</li> <li>➤ Tokenization and Types of Tokenizers, Issues with Tokenization</li> <li>➤ Stemming, RegexpStemmer, The Porter Stemmer</li> <li>➤ Lemmatization</li> <li>➤ Language Translation</li> <li>➤ Stop-Word Removal</li> <li>➤ Feature Extraction from Texts, Extracting General Features from Raw Text</li> <li>➤ Bag of Words, Zipf's Law</li> <li>➤ TF-IDF, Feature Engineering</li> <li>➤ Word Clouds and Other Visualizations</li> </ul>	8am – 5pm
Three and Four	<b>MODULE 3: DEVELOPING A TEXT CLASSIFIER</b> <ul style="list-style-type: none"> <li>➤ Machine Learning</li> <li>➤ Unsupervised Learning</li> </ul>	8am – 5pm

	<ul style="list-style-type: none"> <li>➤ Hierarchical Clustering</li> <li>➤ K-Means Clustering</li> <li>➤ Supervised Learning</li> <li>➤ Classification</li> <li>➤ Logistic Regression</li> <li>➤ Naive Bayes Classifiers</li> <li>➤ K-Nearest Neighbors</li> <li>➤ Regression</li> <li>➤ Linear Regression</li> <li>➤ Tree Methods</li> <li>➤ Random Forest</li> <li>➤ GBM and XGBoost</li> <li>➤ Sampling</li> <li>➤ Developing a Text Classifier</li> <li>➤ Feature Extraction</li> <li>➤ Feature Engineering</li> <li>➤ Removing Correlated Features</li> <li>➤ Dimensionality Reduction</li> <li>➤ Deciding on a Model Type</li> <li>➤ Evaluating the Performance of a Model</li> <li>➤ Building Pipelines for NLP Projects</li> <li>➤ Saving and Loading Models</li> </ul>	
Five	<p><b>MODULE 4: COLLECTING TEXT DATA FROM THE WEB</b></p> <ul style="list-style-type: none"> <li>➤ Collecting Data by Scraping Web Pages</li> <li>➤ Requesting Content from Web Pages</li> <li>➤ Dealing with Semi-Structured Data</li> <li>➤ JSON</li> <li>➤ XML</li> <li>➤ Using APIs to Retrieve Real-Time Data</li> <li>➤ API Creation</li> <li>➤ Extracting Data from Local Files</li> </ul>	8am – 5pm
Six	<p><b>MODULE 5: TOPIC MODELING</b></p> <ul style="list-style-type: none"> <li>➤ Topic Discovery</li> <li>➤ Discovering Themes</li> <li>➤ Exploratory Data Analysis</li> <li>➤ Document Clustering</li> <li>➤ Dimensionality Reduction</li> <li>➤ Historical Analysis</li> <li>➤ Bag of Words</li> <li>➤ Topic Modeling Algorithms</li> <li>➤ Latent Semantic Analysis</li> <li>➤ LSA – How It Works</li> <li>➤ Latent Dirichlet Allocation</li> <li>➤ LDA – How It Works</li> <li>➤ Topic Fingerprinting</li> </ul>	8am – 5pm
Seven	<p><b>MODULE 6: TEXT SUMMARIZATION AND TEXT GENERATION</b></p> <ul style="list-style-type: none"> <li>➤ What is Automated Text Summarization?</li> <li>➤ Benefits of Automated Text Summarization</li> <li>➤ High-Level View of Text Summarization</li> <li>➤ Extractive Text Summarization</li> <li>➤ Abstractive Text Summarization</li> </ul>	8am – 5pm

	<ul style="list-style-type: none"> <li>➤ Sequence to Sequence</li> <li>➤ Encoder Decoder</li> <li>➤ TextRank</li> <li>➤ Summarizing Text Using Gensim</li> <li>➤ Summarizing Text Using Word Frequency</li> <li>➤ Generating Text with Markov Chains</li> <li>➤ Markov Chains</li> </ul>	
Eight	<p><b>MODULE 7: VECTOR REPRESENTATION</b></p> <ul style="list-style-type: none"> <li>➤ Vector Definition</li> <li>➤ Why Vector Representations?</li> <li>➤ Encoding</li> <li>➤ Character-Level Encoding</li> <li>➤ Positional Character-Level Encoding</li> <li>➤ One-Hot Encoding</li> <li>➤ Key Steps in One-Hot Encoding</li> <li>➤ Word-Level One-Hot Encoding</li> <li>➤ Word Embeddings</li> <li>➤ Word2Vec</li> <li>➤ Using Pre-Trained Word Vectors</li> <li>➤ Document Vectors</li> <li>➤ Uses of Document Vectors</li> </ul>	8am – 5pm
Nine and Ten	<p><b>MODULE 8: SENTIMENT ANALYSIS</b></p> <ul style="list-style-type: none"> <li>➤ Why is Sentiment Analysis Required?</li> <li>➤ Growth of Sentiment Analysis 0</li> <li>➤ Monetization of Emotion</li> <li>➤ Types of Sentiments</li> <li>➤ Key Ideas and Terms</li> <li>➤ Applications of Sentiment Analysis</li> <li>➤ Tools Used for Sentiment Analysis</li> <li>➤ NLP Services from Major Cloud Providers</li> <li>➤ Online Marketplaces</li> <li>➤ Python NLP Libraries</li> <li>➤ Deep Learning Libraries</li> <li>➤ TextBlob</li> <li>➤ Understanding Data for Sentiment Analysis</li> <li>➤ Training Sentiment Models</li> </ul>	8am – 5pm

## TRAINER'S BIOGRAPHIES



### **RAYNER ALFRED**

#### **ASSOCIATE PROFESSOR OF COMPUTER SCIENCE**

Certified IBM DB2 Academic Associate, Certified Tester Foundation Level (CTFL)

**AREAS OF SPECILIZATION:** Advanced Machine Intelligence, Data Analytics, Data Mining, Information Retrieval, Artificial Intelligence, Machine Learning, Knowledge Discovery

**ADDRESS:** Knowledge Technology Research Group, Faculty of Computing and Informatics, Universiti Malaysia Sabah, Jalan UMS, 88400 Kota Kinabalu, Sabah.

**CONTACT:** Mobile: 6013-881-9966, eMail: [ralfred@ums.edu.my](mailto:ralfred@ums.edu.my)

Rayner Alfred is an Associate Professor of Computer Science at the Faculty of Computing and Informatics, Universiti Malaysia Sabah in Malaysia that focuses on Data Science and Software Engineering programmes. He leads and defines projects around knowledge discovery, information retrieval and machine learning that focuses on building smarter mechanism that enables knowledge discovery in structured and unstructured data. His work addresses the challenges related to big data problem: How can we create and apply smarter collaborative knowledge discovery and machine learning technologies that bridge the structured and unstructured data mining and cope with the big data problem.

Rayner completed his PhD in 2008 looking at intelligent techniques using machine learning to model and optimize the dynamic and distributed processes of knowledge discovery for structured and unstructured data. He holds a PhD degree in Computer Science from York University (United Kingdom), a master's degree in computer science from Western Michigan University, Kalamazoo (USA) and a Computer Science degree from Polytechnic University of Brooklyn, New York (USA) where he was the recipient of the *Myron M. Rosenthal Academic Achievement Award* for the outstanding academic achievement in Computer Science in 1994. He has authored and co-authored more than 100 journals/book chapters and conference papers, editorials, and served on the program and organizing committees of numerous national and international conferences and workshops.

Rayner is currently a member of IEEE, a Certified Software Tester (CTFL) from the International Software Testing Qualifications Board (*ISTQB*), and a certified IBM DB2 Academic Associate (IBM DB2 AA). He leads the Advanced Machine Intelligence (AMI) research group in UMS and he has led several projects related to knowledge discovery and machine learning on Big Data. Rayner is also the recipient of the Research Fellow at Japan Advanced Institute of Science and Technology (JAIST), Japan. He is also the recipient of multiple GOLD awards at national and international research exhibitions in Data Mining and Machine Learning based solutions (Face Recognition and Knowledge Discovery), that include International Trade Fair Ideas in Nuremberg, Germany (iNEA2018) International Invention Innovation Competition in Toronto, Canada (iCAN 2018), Seoul International Invention Exhibition in Seoul, Korea (SIIF 2010). He has secured RM6,931.433.00 worth of project grants. Some of his project researches include biometric authentication using face recognition, building security based on plate number recognition using deep learning, sentiment analysis for Malay and English in measuring public opinion, news-news correlation trending, machine learning algorithm-based solution for predicting diseases in health care, smart monitoring using an ensemble based face recognition system and smart information management and retrieval to name a few. Some of the completed projects include Semantic Multi-Agent For Knowledge Sharing, developing an Evolutionary-Based Ensemble Classifier Framework for Learning Big Relational Data, developing a genetic-based hierarchical agglomerative clustering technique for parallel clustering of bilingual corpora based on reduced terms, enhancing document Clustering By Integrating Semantic Background Knowledge and Syntactic

Features Into the BOW Representation and the fundamental Study on an Evolutionary Based Features Construction Methods for Data Summarization Approach to Predict Survival Factors of Coral Reefs in Malaysia, to name a few and also infrared face recognition based on ensemble approach.