# Data Analytics

# Lesson 12.

## Recap and advanced topics

Dr. Hai Tran

hai.tran@sbsuni.edu.vn

Scholar: https://scholar.google.com/citations?user=kHZvlTkAAAAJ&hl=en&oi=ao

Co-Founder: XAI - https://xai.foo/

Saigon Business School

In partnership with

UCW UNIVERSITY CANADA WEST

MACQUARIE University SYDNEY·AUSTRALIA

UON University of Northampton

AMITY GLOBAL INSTITUTE

# Learning materials

● Textbook
  ● Evans, J. (2016) Business Analytics. 2nd edn. Pearson.
  ● Runkler, T. (2016) Data Analytics: Models and Algorithms for Intelligent Data Analysis. 2nd edn. Vieweg+Teubner Verlag.

● Online reference materials
  ● archive.ics.uci.edu/ml/
  ● powerbi.microsoft.com
  ● https://github.com/topics/data-analysis-python
  ● https://media.pearsoncmg.com/ph/esm/esm_evans_eba3e_20/tools/eba3e_analytic_solver.html
  ● https://data.imf.org/

# Agenda

- Lesson 1: Understanding Data Analytics Terminologies.
- Lesson 2: Foundation of Business Analytics
- Lesson 3: Visualizing and Exploring data
- Lesson 4: Applying Descriptive Analytic Techniques
- Lesson 5: Data Modeling
- Lesson 6: Predictive Analytics
- Lesson 7: Regression, Classification and Clustering
- Lesson 8: Forecasting Techniques
- Lesson 9: Investigating Predictive Analytic Techniques
- Lesson 10: Introduction to Data Mining
- Lesson 11: Demonstrating Prescriptive Analytic Methods
- Lesson 12: Recap and advanced topics

sbsedu.vn

# Recap and advanced topics

# Recap and advanced topics
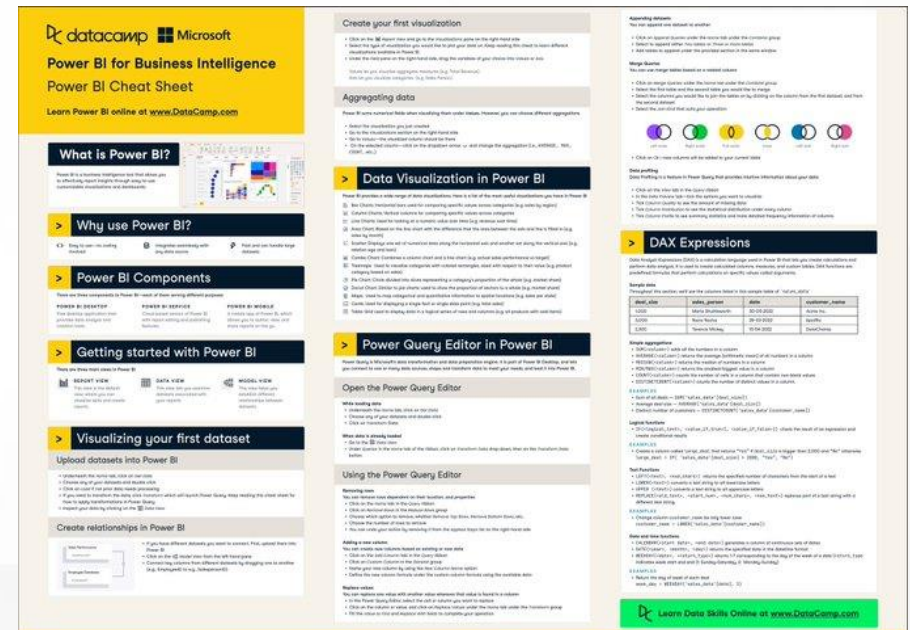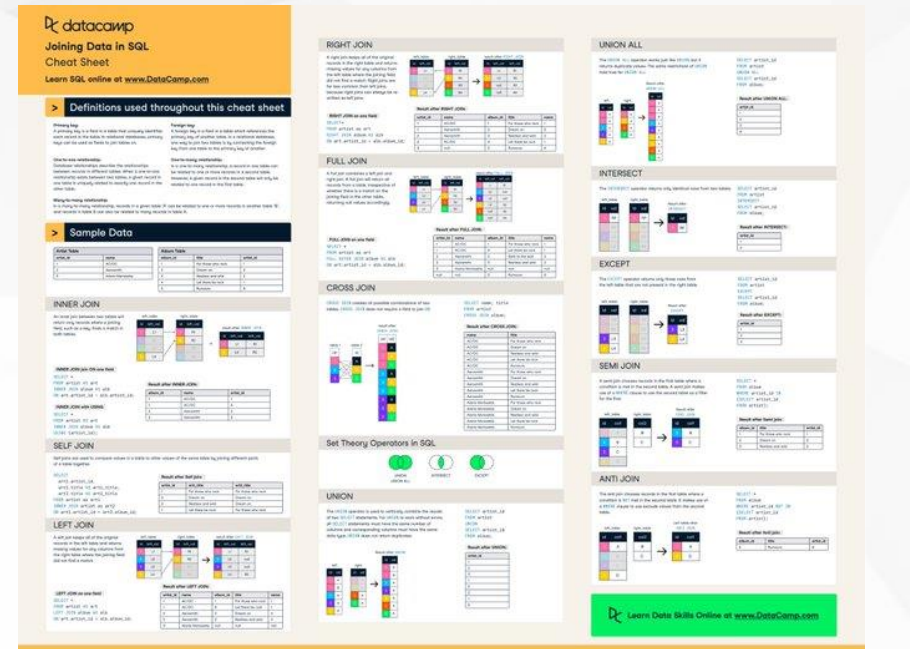
# Conclusion & Questions

Here are five advanced topics in data analytics:

1. **Deep Learning and Neural Networks:** Explore advanced techniques in deep learning, including convolutional neural networks (CNNs) and recurrent neural networks (RNNs), to enhance pattern recognition and prediction capabilities.

2. **Natural Language Processing (NLP) and Text Analytics:** Delve into the complexities of analyzing and understanding human language, including sentiment analysis, text summarization, and language generation, to derive insights from unstructured text data.

3. **Time Series Analysis and Forecasting:** Focus on analyzing data that varies over time, such as stock prices, weather patterns, or sales figures, using sophisticated methods like ARIMA (AutoRegressive Integrated Moving Average) and machine learning models for accurate forecasting.

4. **Geospatial Analytics:** Explore the integration of geographic information systems (GIS) with data analytics, allowing for the analysis and visualization of spatial patterns, location-based insights, and geospatial data to make informed decisions.

5. **Explainable AI (XAI):**Address the interpretability and transparency challenges associated with complex machine learning models, ensuring that the decisions made by these models can be understood and trusted by stakeholders, especially in sensitive domains like healthcare and finance.

# Thank you