# UBC Capstone Project 2023

## Masters of Data Science
## Computational Linguistics

# About the UBC Capstone Project

**Objective**: students apply their knowledge to help partner organization solve a business problem.

**Duration**: 2 months (May - June)

**Application process:**

- November – online proposal.

- January – meeting with students.

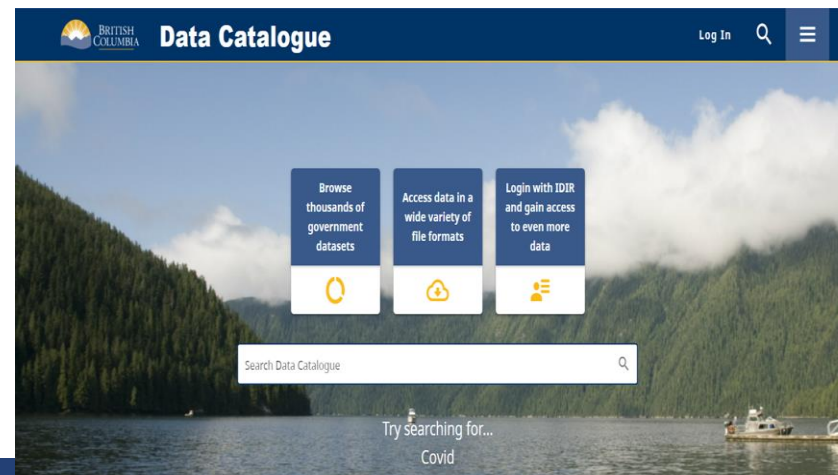- February – decisions announced.

# Business Problem

## Problem Statements

1. The **search function is limited** in its ability to display all the relevant datasets.

2. The **system is limited in using analytics to portray a big picture** understanding of the datasets in the catalogue.
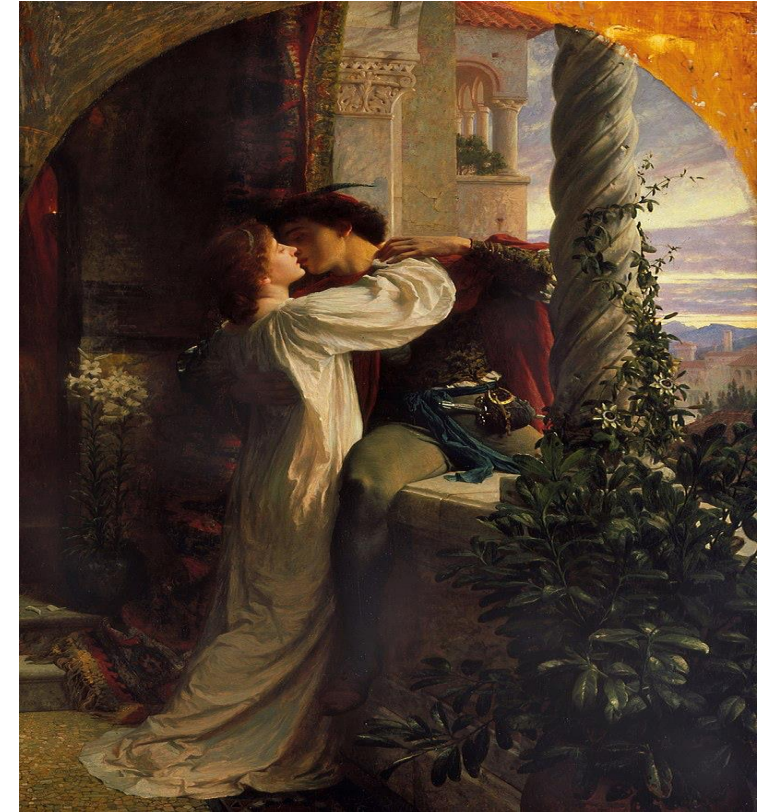
## Opportunities

1. Explore options to **enhance the search function** in the Data Catalogue

2. Seek an independent perspective on **how to visualize the data** and **improve the experience for users**
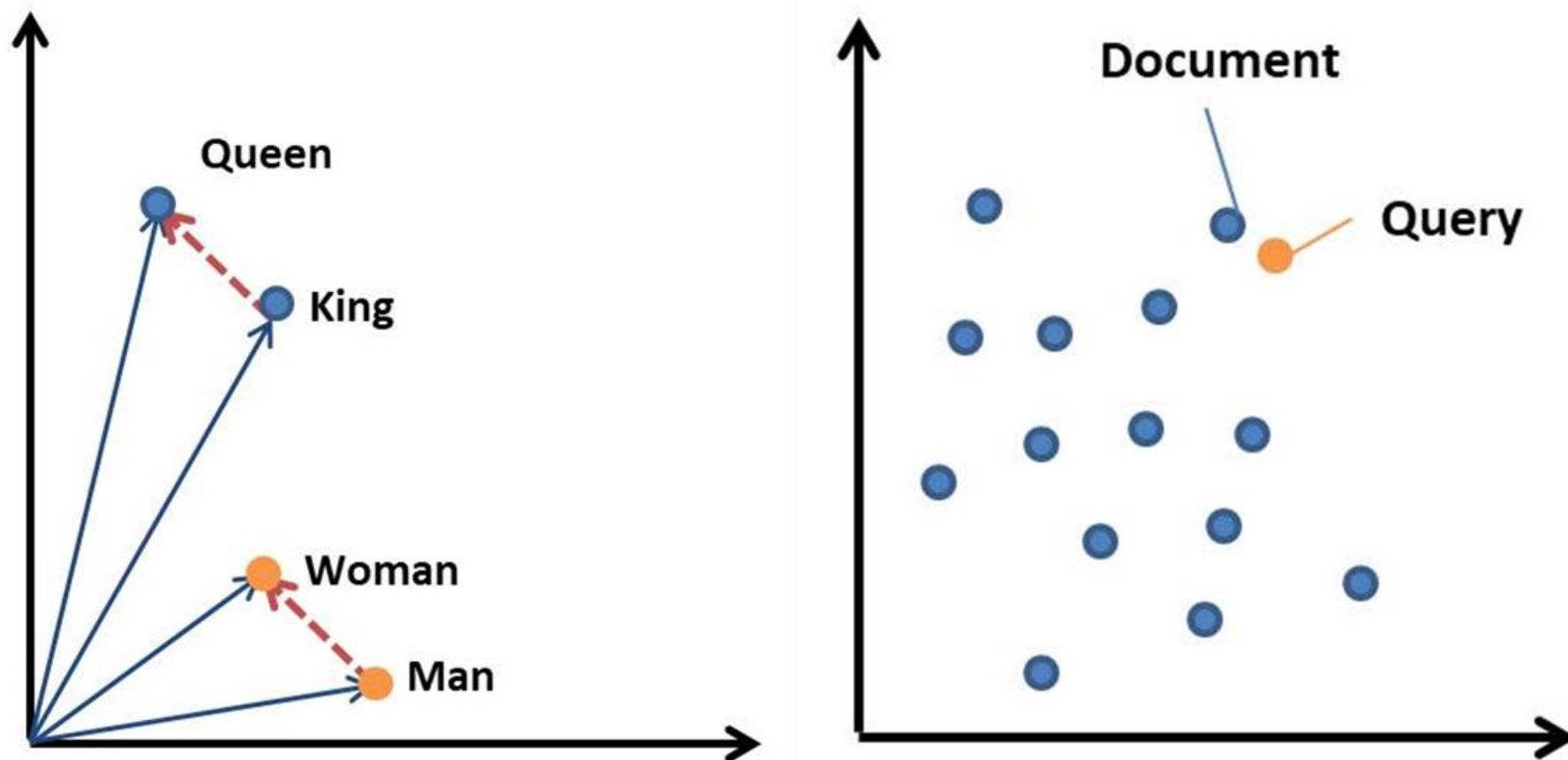
# Causes of the Problem

- Search engines (ad hoc IR systems)

- Vocabulary: Mismatch Problem
  - Example: A tragic love story

- Search engine as a binary test:
  - False positives
  - False negatives

# Solution

- Semantic Search (embedding learning)

- Example
  - A bottle of tezguino is present at the table.
  - Tezguino is highly popular among individuals.
  - Consumption of tezguino induces intoxication.
  - Corn is used in the production of tezguino.

# Solution

# Demo – Increases Accuracy

# Demo – Understands Synonyms

# Demo – Understands Phrases



BC DATA SERVICE

# Demo – Handles Spelling Errors

# Topic Modeling



## Hierarchical Clustering

14_financialeconomic_financ...
10_operating_tbls_quarterly
13_provincial_province_debt
5_forecast_forecasts_interest
7_revenue_revenues_budget
18_taxes_tax_assessed
15_municipality_taxes_provi...
19_census_demographics_bc
17_metadata_icbc_e02
24_maps_columbia_provinces
2_cariboo_caribou_habitat
9_wildfire_fire_bc
16_bcgnws_drivebc_bcts
1_ecosystem_ecosystems_poly...
4_forest_timber_trees
25_mta_coal_mining
0_schools_enrollment_tuition
3_hydrology_hydrologic_wate...
11_bcer_permits_geophysical
8_tantalis_crown_land
23_courts_court_reports
12_wildlife_species_habitat
6_habitat_fish_lakes
21_owl_habitat_spotted
20_coastal_shorebirds_sea
22_fishery_fisheries_fish

# Implementation Plan

- **Finish upgrade** to ckan and solr (search platform)

- **Integrate the enhanced search solution** with new version of solr

- **Explore topic modeling as an enhancement** to the BC Data Catalogue

# Lessons Learned Along the Way

## BCDS

- Personal information cannot be collected from the web.

- Need to prepare a "gold standard" dataset and evaluation metrics in advance.

## UBC Students

- How things work in the industry

- How we can apply our knowledge

- Learned BC Data Service work ethics and tools

# Benefits

## BCDS

- Collaboration opportunity for BC Stats and DSS.

- Established project documentation to facilitate future capstone projects

- Got a fresh perspective on enhancing the product

- Learned about practical applications of computational linguistics

## UBC Students

- Learned new skills

- Worked with real-world data

- Teamwork experience

# How Labour Intensive Was It?

**BCDS**

- 35 hours for proposal and project set-up.

- 6 hours for weekly meetings

- 15 hours for ad hoc student support

**UBC Students**

- 6 weeks full time work

# References

- Furnas, G. W., Landauer, T. K., Gomez, L. M., and Dumais, S. T. (1987). The vocabulary problem in human-system communication. Communications of the ACM 30(11), 964971.

- Jurafsky, D., & Martin, J. H. (2020). Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition. Pearson.

- Lin, J., Nogueira, R., and Yates, A. (2021). Pretrained transformers for text ranking: BERT and beyond. arXiv:2010.06467v3