



# COVID-19 Open Research Dataset Challenge

NATURAL LANGUAGE PROCESSING AND  
TOPIC MODELING: FINDING RELATED ARTICLES

*Team members:*

Tahani Alqahtani – Amani Albalawi.



# Outline:

---

- Objectives
- Workflow
- Dataset
- Tools
- Data preprocessing
- Topic Modelling
- Clustering Data

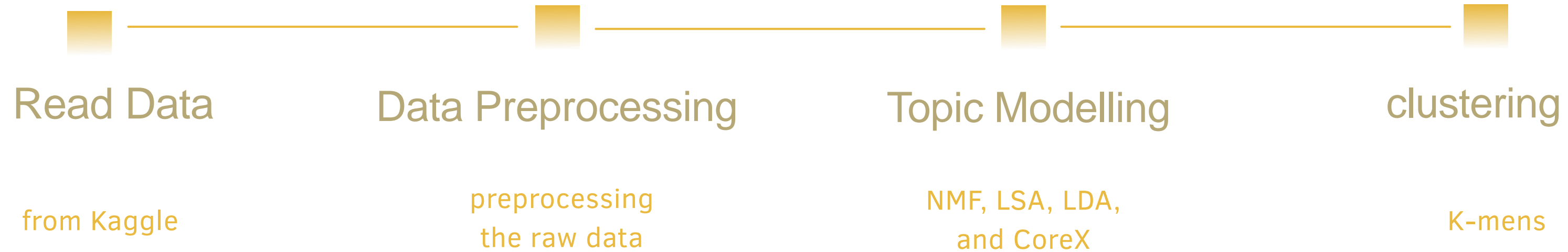


## ➤ Objective:

In this project we will find related articles by using Topic modelling. Here I am using, LSA , NMF , CorEx, and LDA. All this topic models are used to classify text in a document to a particular topic. It builds a topic per document model and words per topic model.

# WORKFLOW

FROM START TO FINISH





# Dataset

---

- We have made this dataset available on Kaggle and we use abstract column to make my model .



# ➤ Tools & libraries:

## Libraries



## Technologie





# Data Preprocessing

---

1

**Remove  
Numbers,  
Capital letters  
and  
Punctuation**

2

**Remove stop  
word**

3

**Tokenizing  
and  
Lemmatizing**

# ➤ Topic Modelling

## NMF

Topics(10)  
CV & TFIDF

## LSA

Topics(10)  
CV & TFIDF

## CorEx

Topics(5)  
CV & TFIDF

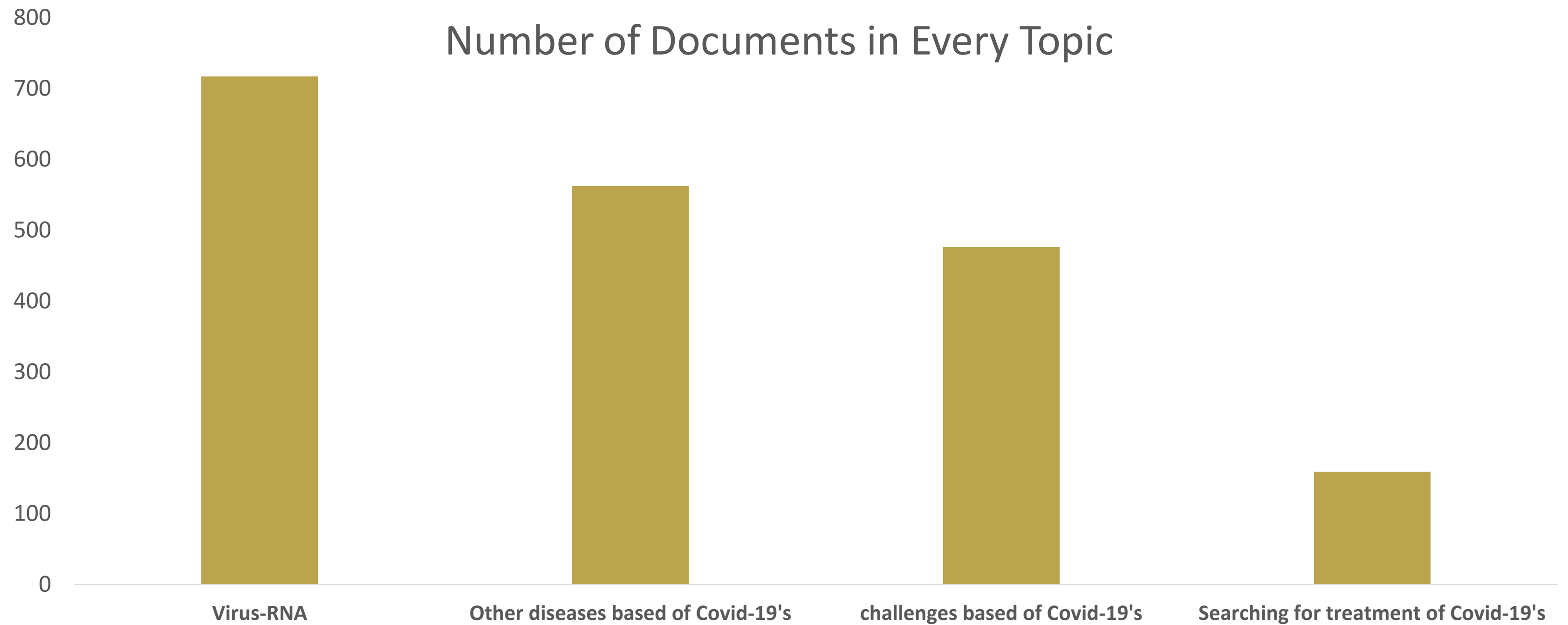
## LDA

Topics(2-10)  
CV & TFIDF



Best model was Count Vectorizer LDA with 4 topics



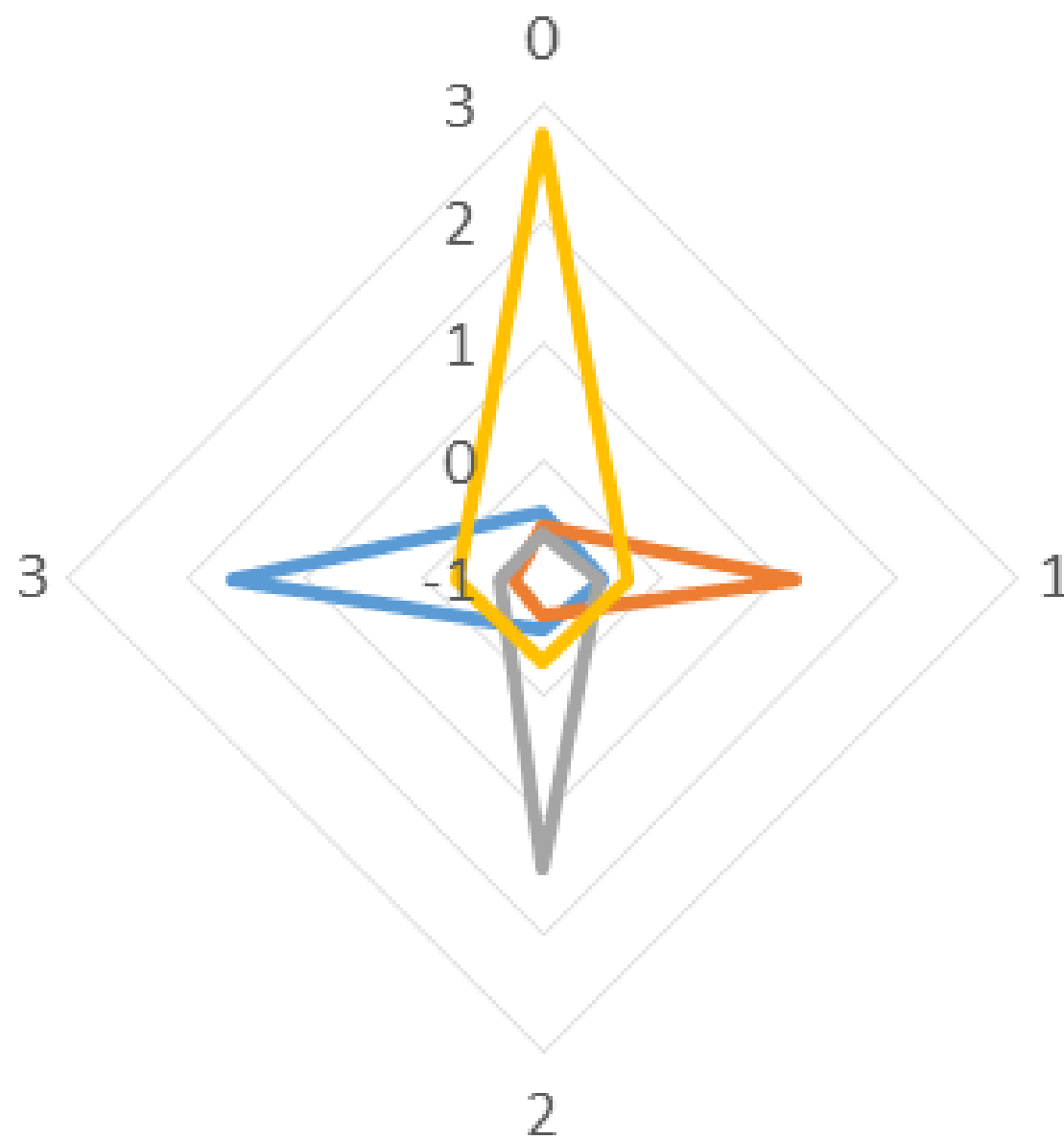




# Clustering Model

## K-men cluster

— Topic0 — Topic1 — Topic2 — Topic3





Any question?

Thanks