# Integrating Social Network Structure into Online Feature Selection

ANTONELA TOMMASEL

ADVISOR: DANIELA GODOY

ISISTAN

*Instituto Superior de Ingenieria de Software Tandil*

# Motivation

- Short-texts **accentuate** the **challenges** posed by the **high feature space dimensionality** of text learning tasks.

- The **linked** nature of **social data** causes **new dimensions** to be added to the feature space, which, also becomes **sparser**.

**Efficient and scalable online feature selection becomes a crucial requirement of numerous large-scale social applications.**
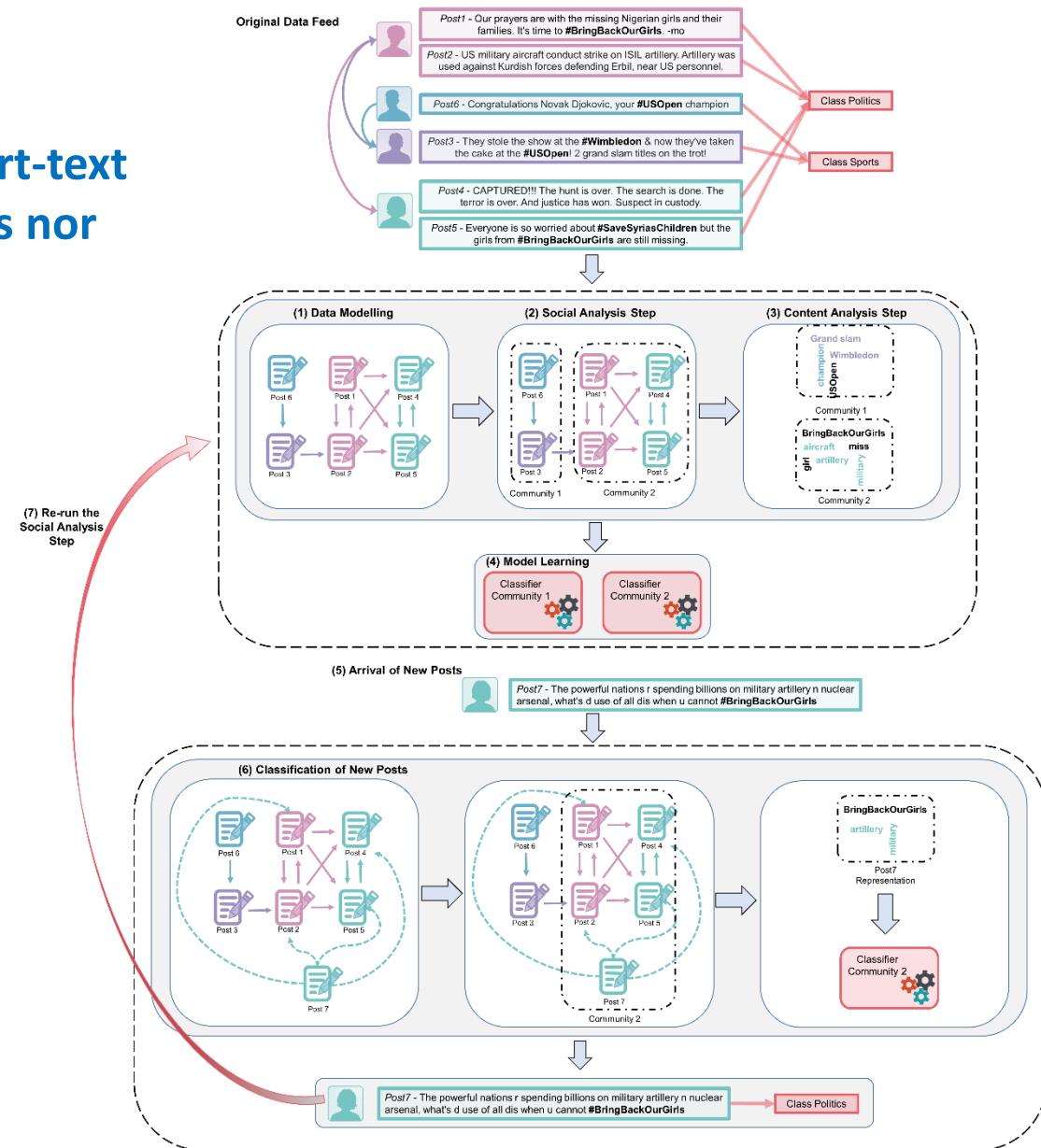
# Proposal

- An **Online Feature Selection** technique for **high-dimensional** data based on both social and **content-based information** for the **real-time** classification of **short-text** streams coming from social media.

- *Objectives*?
    - Enhancing the process of knowledge discovery in social-media.

    - Helping in the development of new and more effective models for personalisation and recommendation of content in social environments.

# Social-based OFS

**Addresses the massive-scale OFS task for high-dimensional short-text data arriving in a continuous stream, in which neither features nor data instances are fully known in advance.**

# Social-based OFS

**Addresses the massive-scale OFS task for high-dimensional short-text data arriving in a continuous stream, in which neither features nor data instances are fully known in advance.**

# Social-based OFS

**Addresses the massive-scale OFS task for high-dimensional short-text data arriving in a continuous stream, in which neither features nor data instances are fully known in advance.**

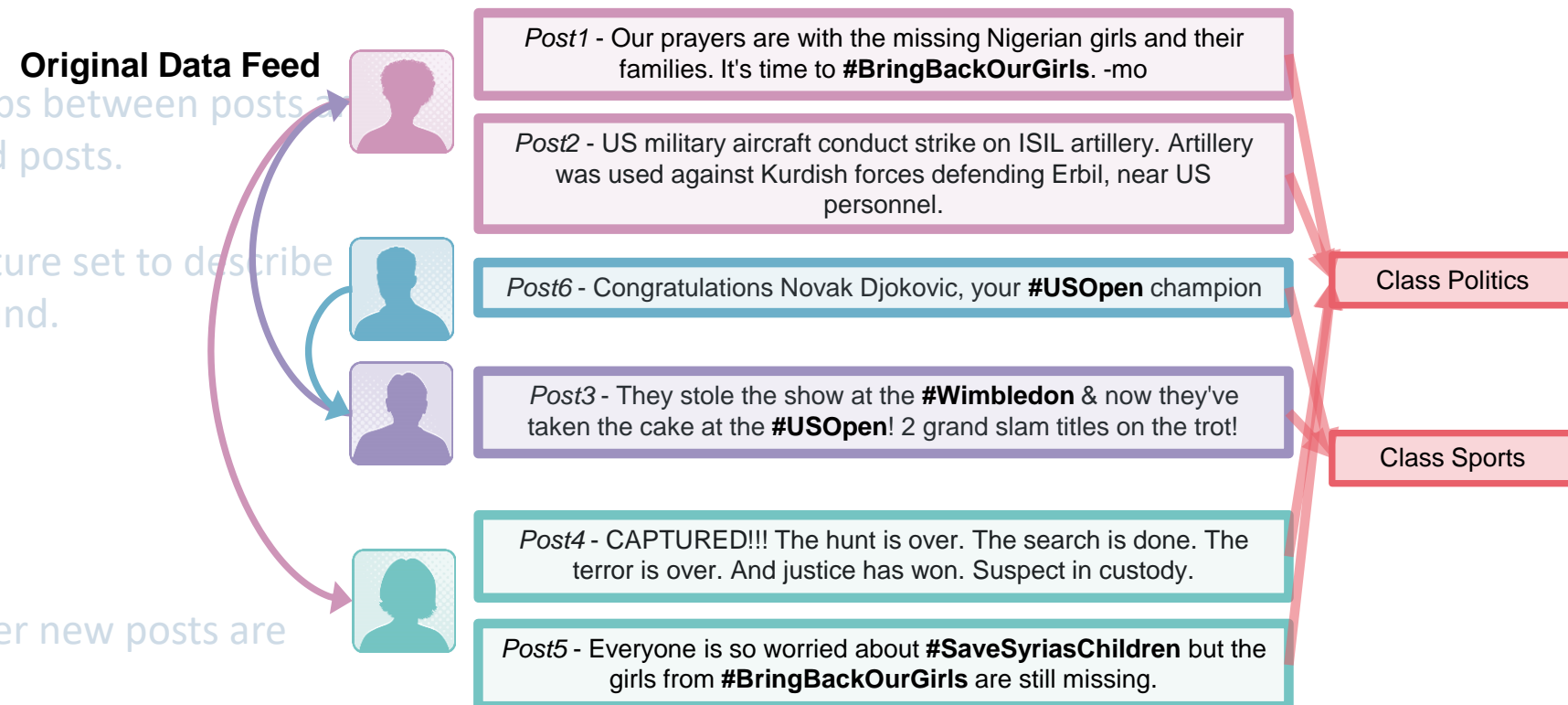1. *Data Modelling* as a graph representing the social posts and their relations.

2. *Social Analysis Step.* Social relationships between posts are analysed to find groups of socially related posts.

3. *Content Analysis Step.* An optimal feature set to describe each group of socially related posts is found.

4. *Model Learning.*

5. *Arrival* and classification of new posts.

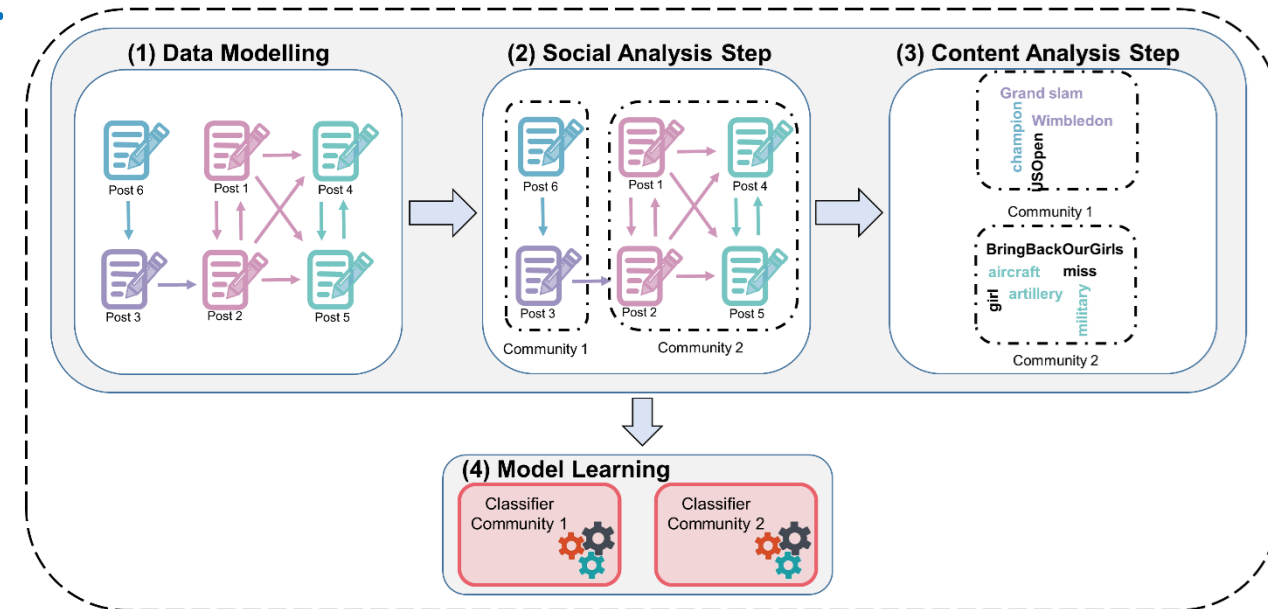6. *Re-run of the Social Analysis Step.* After new posts are classified, the feature space is updated.

**Original Data Feed**

*Post1* - Our prayers are with the missing Nigerian girls and their families. It's time to **#BringBackOurGirls**. -mo

*Post2* - US military aircraft conduct strike on ISIL artillery. Artillery was used against Kurdish forces defending Erbil, near US personnel.

*Post6* - Congratulations Novak Djokovic, your **#USOpen** champion

*Post3* - They stole the show at the **#Wimbledon** & now they've taken the cake at the **#USOpen**! 2 grand slam titles on the trot!

*Post4* - CAPTURED!!! The hunt is over. The search is done. The terror is over. And justice has won. Suspect in custody.

*Post5* - Everyone is so worried about **#SaveSyriasChildren** but the girls from **#BringBackOurGirls** are still missing.

Class Politics

Class Sports

# Social-based OFS

**Addresses the massive-scale OFS task for high-dimensional short-text data arriving in a continuous stream, in which neither features nor data instances are fully known in advance.**

1. *Data Modelling* as a graph representing the social posts and their relations.

2. *Social Analysis Step.* Social relationships between posts are analysed to find groups of socially related posts.

3. *Content Analysis Step.* An optimal feature set to describe each group of socially related posts is found.

4. *Model Learning.*

5. *Arrival* and classification of new posts.

6. *Re-run of the Social Analysis Step.* After new posts are classified, the feature space is updated.

# Social-based OFS

**Addresses the massive-scale OFS task for high-dimensional short-text data arriving in a continuous stream, in which neither features nor data instances are fully known in advance.**

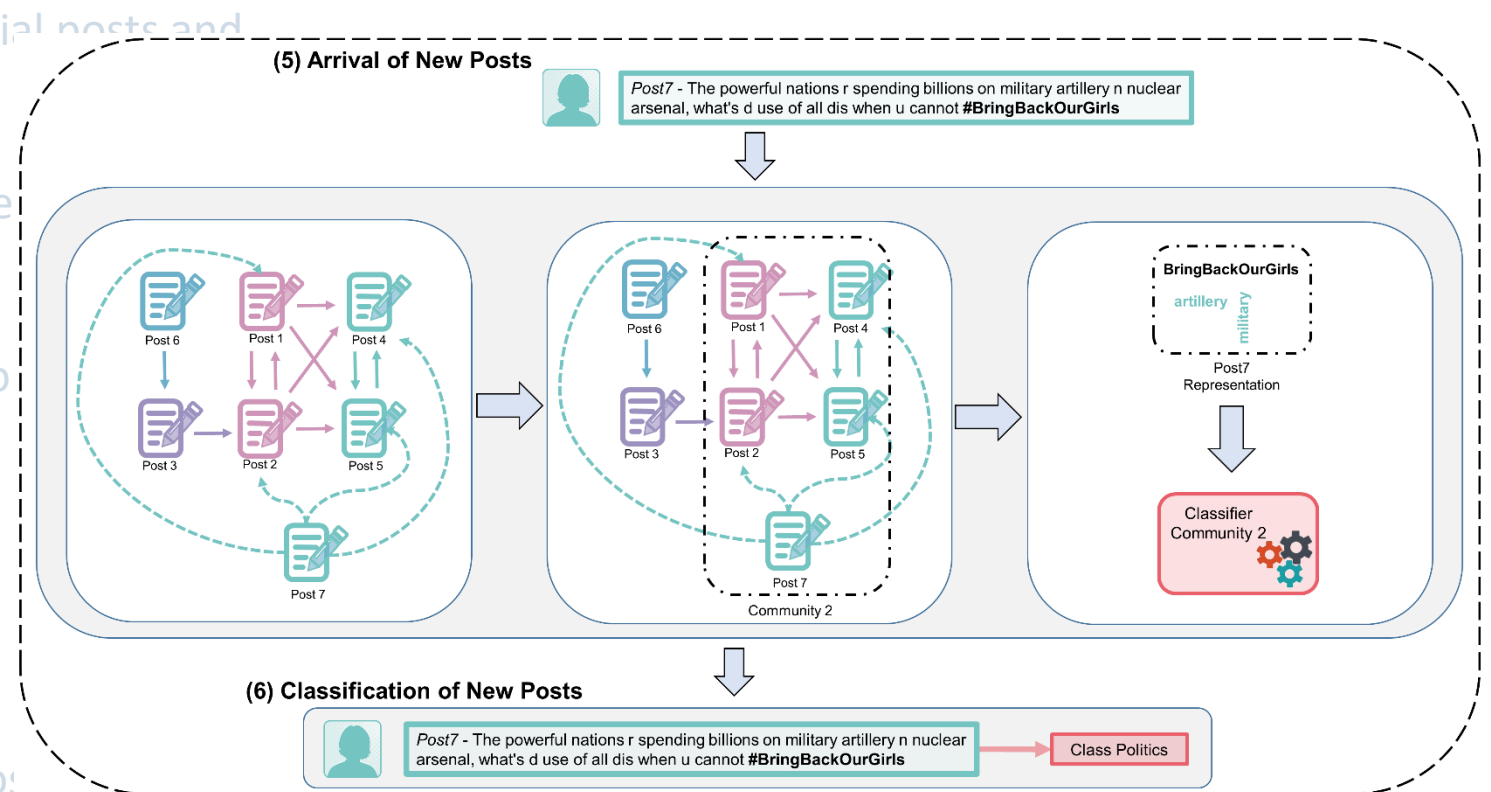1. *Data Modelling* as a graph representing the social posts and their relations.

2. *Social Analysis Step.* Social relationships between analysed to find groups of socially related posts.

3. *Content Analysis Step.* An optimal feature set to each group of socially related posts is found.

4. *Model Learning.*

5. *Arrival* and classification of new posts.

6. *Re-run of the Social Analysis Step.* After new post classified, the feature space is updated.

# Social-based OFS

**Addresses the massive-scale OFS task for high-dimensional short-text data arriving in a continuous stream, in which neither features nor data instances are fully known in advance.**

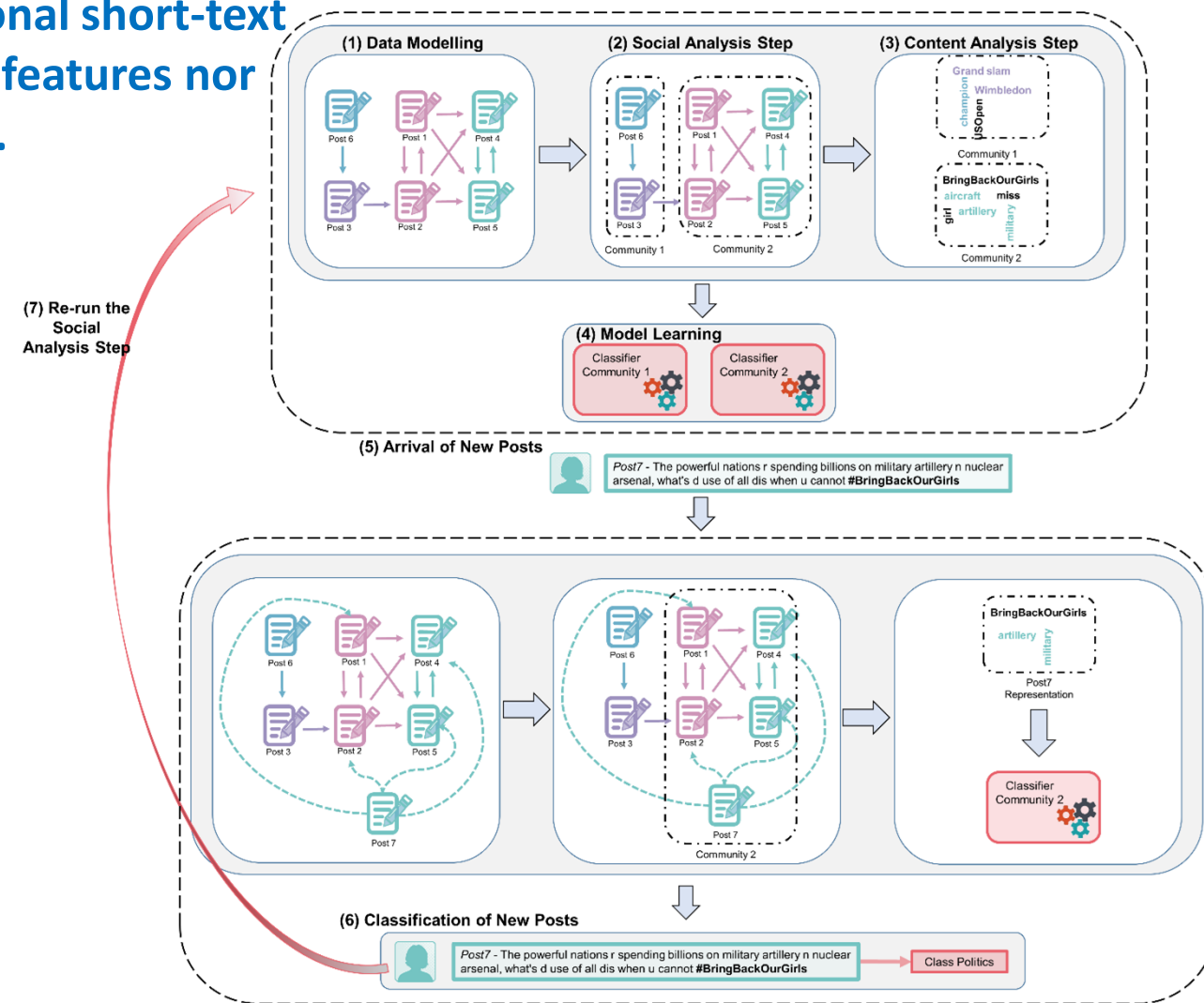1. *Data Modelling* as a graph representing the social posts and their relations.

2. *Social Analysis Step.* Social relationships between posts are analysed to find groups of socially related posts.

3. *Content Analysis Step.* An optimal feature set to describe each group of socially related posts is found.
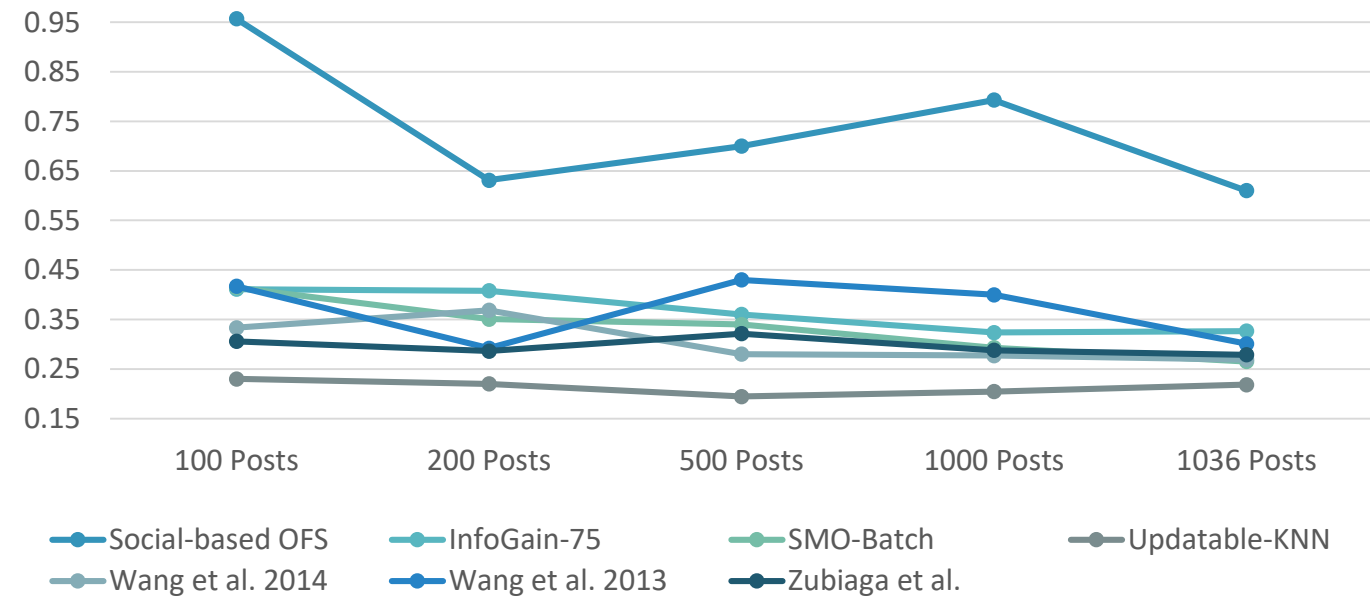
4. *Model Learning.*

5. *Arrival* and classification of new posts.

6. *Re-run of the Social Analysis Step.* After new posts are classified, the feature space is updated.

# Current State

- Preliminary evaluations conducted on two real-world short-texts datasets achieved **promising** results when compared to traditional and state-of-the-art in both batch and online settings!!

- The obtained results exposed the limitations of pure content-based techniques for classifying social media short-texts.
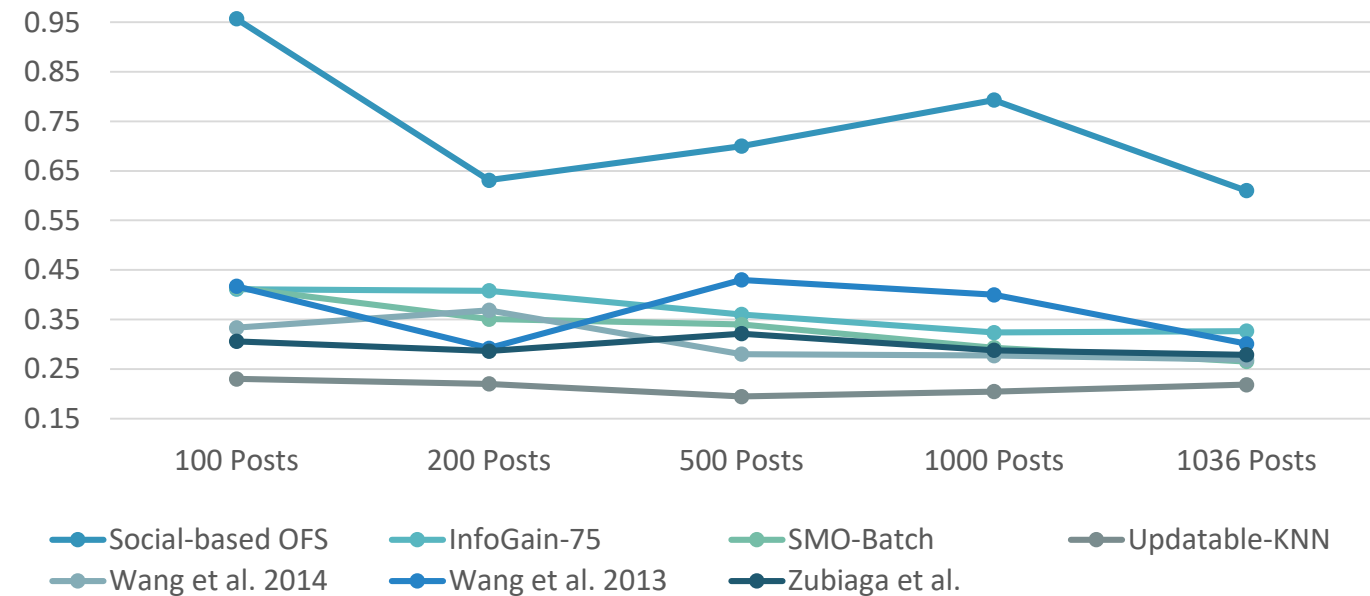
# Current State

- Preliminary evaluations conducted on two real-world short-texts datasets achieved **promising** results when compared to traditional and state-of-the-art in both batch and online settings!!

- The obtained results exposed the limitations of pure content-based techniques for classifying social media short-texts.

**Leveraging on social information becomes crucial for OFS.**



Legend: Social-based OFS, InfoGain-75, SMO-Batch, Updatable-KNN, Wang et al. 2014, Wang et al. 2013, Zubiaga et al.

X-axis: 100 Posts, 200 Posts, 500 Posts, 1000 Posts, 1036 Posts

Y-axis: 0.15, 0.25, 0.35, 0.45, 0.55, 0.65, 0.75, 0.85, 0.95

# Contributions

- This thesis tackles the **challenging** problem of **Online Feature Selection**.

- Addresses the problem of **how** to **exploit** the **linked nature** of social media data.

- Proposes a technique for **leveraging on social relations**.

- **Combines** **social information** with **content** for effectively and efficiently performing feature selection.

- **Scalability**. Appropriate for real-time environments in which neither features nor instances are known in advance.

# Questions?

# Integrating Social Network Structure into Online Feature Selection

ANTONELA TOMMASEL

ADVISOR: DANIELA GODOY

ISISTAN

*Instituto Superior de Ingenieria de Software Tandil*