





Malicious Behavior on the Web: Characterization and Detection

Srijan Kumar (@srijankr) Justin Cheng (@jcccf) Jure Leskovec (@jure)

Slides are available at http://snap.stanford.edu/www2017tutorial/

Tutorial Outline

Malicious users



Misinformation

Fake reviews



http://snap.stanford.edu/www2017tutorial

Vandalism

Vandalism is "an action involving deliberate destruction of or damage to public or private property."

Vandalism is common on Wikipedia

- Freely accessible
- Large reach
- Major source of information for many ____



Easy to add content

WikipediA

The free encyclopedia that anyone can edit

Vandalism: An edit that is:

- Non-value adding
- Offensive
- Destructive in removal

Vandalism



 \sim 7% edits are vandalism \sim 3-4 % editors are vandals



Emma Stone

Emily Jean "Emma" Stone is a hot American actress with a beautiful smile. In 1987, she fell out of the sky as an angel. Wikipedia

Bowe: November 6, 1988 (age 24), Scottsdale, Arizona, United Stat

Height: 1.68 m

Actres

Siblings: Spencer Stone Parents: Krista Stone, Jeff Stone Upcoming movie: The Amazing Spider-Man 2

Movies



Spider-Man

The Amazing

2012



The Croods

2013



Gangster

Squad

2013



Easy A

2010



2009

Zombieland

Tools to detect vandalism on Wikipedia

STiki: Metadata

EDITOR

registered?, account-age, geographical location, edit quantity, revert history, block history, is bot?, quantity of warnings on talk page

ARTICLE

age, popularity, length, size change, revert history

REVISION COMMENT length, section-edit?

TIMESTAMP time-of-day, day-of-week

West et al. (EuroSec, 2010)

ClueBot NG: Textual



- Vocabularies differ between vandalism and innocent edits
- Automatically assess individual word "goodness" probability

WikiTrust: Content driven



- Content that survives is good content
- Good content builds reputation for its author

Detection of vandals

Vandalism detection



Vandal detection

Using STiki to detect vandals



Stiki rule: Editor is a vandal if any edit's suspicion score exceeds threshold



Using ClueBot NG to detect vandals



ClueBot rule: Editor is a vandal if it reverts at least N edits

Using ClueBot NG to detect vandals



ClueBot rule: Editor is a vandal if it reverts at least N edits

Objective: Detect vandals in as few edits as possible

Data: Wikipedia Vandals

34,000 Editors Half are vandals

770,000 Edits 160,000 edits by vandals

Time: Jan 2013 - July 2014

Characteristics of vandals



WIKIPEDIA The Free Encyclopedia

Main page Contents Featured content Current events Random article Donate to Wikipedia Wikipedia store

Interaction

Not logged in Talk Contributions Create account Log in
Article Talk
Read Edit View history Search Wikipedia
Perth
From Wikipedia, the free encyclopedia

This article is about the capital of Western Australia. For the city in Scotland, see Perth, Scotland. For other uses, see Perth (disambiguation).

Perth (/ˈpɜːrθ/) is the capital and largest city of the Australian state of Western Australia. It is the fourth-most populous city in Australia, with an estimated population of 2.06 million (as of 30 June 2016) living in



Editors can edit article pages and talk pages

Vandals make visible edits



Vandals are quicker



Vandals do not discuss



Vandals make reversion driven edits



Detecting vandals

Pairwise Edit Features



Time x Type of page x First edit x Distance x Similarity x Reverted or not

Meta-Features: Transitions





Vandal Detection



VEWS identifies 87% vandals on or before first reversion. 44% vandal are identified before first reversion.

Early Warning System



Does reversion information help?



Combining Multiple Systems



Summary: Vandals

- Vandals: Users that make non-constructive contribution
- Vandals are aggressive: they make visible edits without discussing and edit war
- Vandals can be detected early by using temporal features and relation between edited pages
- Combination of metadata, text and human feedback is the best in detecting vandals

References

S. Kumar, F. Spezzano and V.S. Subrahmanian. VEWS: A Wikipedia Vandal Early Warning System. SIGKDD 2015.

B. T. Adler, L. de Alfaro, and I. Pye. Detecting wikipedia vandalism using wikitrust - lab report for PAN at CLEF 2010. CLEF, 2010

B. T. Adler, L. de Alfaro, S. M. Mola-Velasco, P. Rosso, and A. G. West. Wikipedia vandalism detection: Combining natural language, metadata, and reputation features. CICLing, 2011.

A. G. West, S. Kannan, and I. Lee. Detecting wikipedia vandalism via spatio-temporal analysis of revision metadata? in EUROSEC, 2010.

M. Potthast, B. Stein, and R. Gerling. Automatic vandalism detection in Wikipedia. Advances in Information Retrieval, ser. Lecture Notes in Computer Science, C. Macdonald, I. Ounis, V. Plachouras, I. Ruthven, and R. White, Eds. Springer Berlin Heidelberg, 2008.

S. Mola-Valesco. Wikipedia vandalism detection. WWW 2011.