

The Goal

Our goal:

Propose ranking strategies for a stream of news information and a set of news sources.

We do not know other algorithms for comparing our results

... moreover ...

... ranking news articles is a different task than ranking web pages



Clustering Technique

We adopt a continuous measure of the lexical similarity between news posting

In our current implementation

- Similarity between news abstracts is represented using the canonical bag of words paradigm
- The abstracts are filtered out against a list of stop words
- The abstracts are ranked according to their similarity



Decay Rule

$R(n, t)$ is the rank of news n at time t .

Decay Rule:

$$R(n, t + \tau) = e^{-\alpha\tau} R(n, t), \quad t > t_i,$$

t_i is the time n_i was posted



Algorithms TA1

This class of algorithms assigns to every source the sum of the ranks assigned to the articles emitted by that source in the past.

Let $R(s, t)$ be the rank of source s at time t

With $S(n_i) = s_k$ we denote that n_i has been posted by s_k .

$$R(s_k, t) = \sum_{S(n_i)=s_k} R(n_i, t),$$

Possible definitions for the rank of pieces of news:

- $R(n_i, t_i) = 1$



Algorithms TA1

This class of algorithms assigns to every source the sum of the ranks assigned to the articles emitted by that source in the past.

Let $R(s, t)$ be the rank of source s at time t

With $S(n_i) = s_k$ we denote that n_i has been posted by s_k .

$$R(s_k, t) = \sum_{S(n_i)=s_k} R(n_i, t),$$

Possible definitions for the rank of pieces of news:

- $R(n_i, t_i) = 1$

However...

... the rank of n_i does not depend on the rank of the generating source



Algorithms TA1

This class of algorithms assigns to every source the sum of the ranks assigned to the news emitted by that source in the past. Let

$R(s, t)$ be the rank of source s at time t

With $S(n_i) = s_k$ we denote that n_i has been posted by s_k .

$$R(s_k, t) = \sum_{S(n_i)=s_k} R(n_i, t),$$

Possible definitions for the rank of pieces of news:

- $R(n_i, t_i) = c \lim_{\tau \rightarrow 0^+} R(S(n_i), t_i - \tau), \quad 0 < c < 1,$

It doesn't work for the limit case LC1



Algorithms TA1

This class of algorithms assigns to every source the sum of the ranks assigned to the news emitted by that source in the past. Let

$R(s, t)$ be the rank of source s at time t

With $S(n_i) = s_k$ we denote that n_i has been posted by s_k .

$$R(s_k, t) = \sum_{S(n_i)=s_k} R(n_i, t),$$

Possible definitions for the rank of pieces of news:

- $R(n_i, t_i) = [\lim_{\tau \rightarrow 0^+} R(S(n_i), t_i - \tau)]^\beta$, $0 < \beta < 1$.

β is similar to the magic ϵ in PageRank.



Algorithms TA1

This class of algorithms assigns to every source the sum of the ranks assigned to the news emitted by that source in the past. Let

$R(s, t)$ be the rank of source s at time t

With $S(n_i) = s_k$ we denote that n_i has been posted by s_k .

$$R(s_k, t) = \sum_{S(n_i)=s_k} R(n_i, t),$$

Possible definitions for the rank of pieces of news:

- $R(n_i, t_i) = [\lim_{\tau \rightarrow 0^+} R(S(n_i), t_i - \tau)]^\beta$, $0 < \beta < 1$.

β is similar to the magic ε in PageRank.



Algorithms TA1

Summarizing

Algorithm TA1

$$R(s_k, t) = \sum_{S(n_i)=s_k} R(n_i, t),$$

$$R(n_i, t_j) = \left[\lim_{\tau \rightarrow 0^+} R(S(n_i), t_j - \tau) \right]^\beta$$



Desiderata

- Property P1: Ranking for News posting and News sources
- Property P2: Important News articles are Clustered
- Property P3: Mutual Reinforcement between News Articles and News Sources
- Property P4: Time awareness
- Property P5: Online processing



Algorithms TA1

A possible algorithm is

Algorithm TA1

$$R(s_k, t) = \sum_{S(n_i)=s_k} R(n_i, t),$$

$$R(n_i, t_j) = \left[\lim_{\tau \rightarrow 0^+} R(S(n_i), t_j - \tau) \right]^\beta$$

It doesn't take into account the clustering process of news!



Algorithms TA2

A good news ranking algorithm working on a stream of information should exploit some data stream **clustering** technique.

$$R(n_i, t_i) = \left[\lim_{\tau \rightarrow 0^+} R(S(n_i), t_i - \tau) \right]^\beta + \sum_{t_j < t_i} e^{-\alpha(t_i - t_j)} \sigma_{ij} R(n_j, t_j)^\beta,$$



Algorithms TA2

A good news ranking algorithm working on a stream of information should exploit some data stream **clustering** technique.

$$R(n_i, t_i) = \left[\lim_{\tau \rightarrow 0^+} R(S(n_i), t_i - \tau) \right]^\beta + \sum_{t_j < t_i} e^{-\alpha(t_i - t_j)} \sigma_{ij} R(n_j, t_j)^\beta,$$



Algorithms TA2

A good news ranking algorithm working on a stream of information should exploit some data stream **clustering** technique.

$$R(n_i, t_i) = \left[\lim_{\tau \rightarrow 0^+} R(S(n_i), t_i - \tau) \right]^\beta + \sum_{t_j < t_i} e^{-\alpha(t_i - t_j)} \sigma_{ij} R(n_j, t_j)^\beta,$$

The rank of a piece of news depends on

- the rank of the source
- the rank of “similar” pieces of news



Algorithms TA2

A good news ranking algorithm working on a stream of information should exploit some data stream **clustering** technique.

$$R(n_i, t_i) = \left[\lim_{\tau \rightarrow 0^+} R(S(n_i), t_i - \tau) \right]^\beta + \sum_{t_j < t_i} e^{-\alpha(t_i - t_j)} \sigma_{ij} R(n_j, t_j)^\beta,$$

The rank of a piece of news depends on

- the rank of the **source**
- the rank of “similar” pieces of news



Algorithms TA2

A good news ranking algorithm working on a stream of information should exploit some data stream **clustering** technique.

$$R(n_i, t_i) = \left[\lim_{\tau \rightarrow 0^+} R(S(n_i), t_i - \tau) \right]^\beta + \sum_{t_j < t_i} e^{-\alpha(t_i - t_j)} \sigma_{ij} R(n_j, t_j)^\beta,$$

The rank of a piece of news depends on

- the rank of the source
- the rank of **“similar” pieces of news**



Algorithms TA2

Too bad!

LC2 case

A news source mirroring another, gets a finite rank significantly greater than the rank of the mirrored one!



The Final TA algorithm: TA3

Idea

Modify **a posteriori** the rank of a source

A source which has emitted in the past news stories highly mirrored in the future, will receive a “bonus” acknowledging the importance



The Final TA algorithm: TA3

Idea

Modify **a posteriori** the rank of a source

A source which has emitted in the **past** news stories highly mirrored in the **future**, will receive a “bonus” acknowledging the importance



Algorithms TA3

The algorithm is

$$R(s_k, t) = \sum_{S(n_i)=s_k} R(n_i, t) + \sum_{S(n_i)=s_k} e^{-\alpha(t-t_i)} \sum_{\substack{t_j \in [t_i, t] \\ S(n_i) \neq s_k}} \sigma_{ij} R(n_j, t_j)^\beta,$$

$$R(n_i, t_i) = \left[\lim_{\tau \rightarrow 0^+} R(S(n_i), t_i - \tau) \right]^\beta + \sum_{t_j < t_i} e^{-\alpha(t_i-t_j)} \sigma_{ij} R(n_j, t_j)^\beta.$$



Algorithms TA3

The rank of s_k is

$$\begin{aligned}
 R(s_k, t) &= \sum_{S(n_i)=s_k} R(n_i, t) + \\
 &+ \sum_{S(n_i)=s_k} e^{-\alpha(t-t_i)} \sum_{\substack{t_j \in [t_i, t] \\ S(n_i) \neq s_k}} \sigma_{ij} R(n_j, t_j)^\beta
 \end{aligned}$$

- the ranks of the pieces of news generated in the past
- a factor of the rank of news articles similar and posted later on by other sources



Algorithms TA3

The rank of s_k is

$$\begin{aligned}
 R(s_k, t) = & \sum_{S(n_i)=s_k} R(n_i, t) + \\
 & + \sum_{S(n_i)=s_k} e^{-\alpha(t-t_i)} \sum_{\substack{t_j \in [t_i, t] \\ S(n_i) \neq s_k}} \sigma_{ij} R(n_j, t_j)^\beta
 \end{aligned}$$

- the ranks of the pieces of news generated in the **past**
- a factor of the rank of news articles similar and posted later on by other sources



Algorithms TA3

The rank of s_k is

$$\begin{aligned}
 R(s_k, t) &= \sum_{S(n_i)=s_k} R(n_i, t) + \\
 &+ \sum_{S(n_i)=s_k} e^{-\alpha(t-t_i)} \sum_{\substack{t_j \in [t_i, t] \\ S(n_i) \neq s_k}} \sigma_{ij} R(n_j, t_j)^\beta
 \end{aligned}$$

- the ranks of the pieces of news generated in the past
- a factor of the rank of news articles similar and posted **later on** by other sources



Algorithms TA3

The rank of a single piece of news n_i is still

$$R(n_i, t_i) = \left[\lim_{\tau \rightarrow 0^+} R(S(n_i), t_i - \tau) \right]^\beta + \sum_{t_j < t_i} e^{-\alpha(t_i - t_j)} \sigma_{ij} R(n_j, t_j)^\beta.$$

Note:

If n_i aggregates with a set stories posted in the future, we do not assign to n_i an extra bonus

The idea is that we want to privilege the **freshness** of a news article rather than its clustering importance



Algorithms TA3

The rank of a single piece of news n_i is still

$$R(n_i, t_i) = \left[\lim_{\tau \rightarrow 0^+} R(S(n_i), t_i - \tau) \right]^\beta + \sum_{t_j < t_i} e^{-\alpha(t_i - t_j)} \sigma_{ij} R(n_j, t_j)^\beta.$$

Note:

If n_i aggregates with a set stories posted in the future, we do not assign to n_i an extra bonus

The idea is that we want to privilege the **freshness** of a news article rather than its clustering importance



Algorithms TA3

- 1 TA3 is coherent with all the desirable properties P1–P5
- 2 unfortunately is more complicated than the others
- 3 it is not easy to write down a formula for the stationary mean value of the sources for the limit cases LC1 and LC2



Algorithms TA3

- 1 TA3 is coherent with all the desirable properties P1–P5
- 2 **unfortunately is more complicated than the others**
- 3 it is not easy to write down a formula for the stationary mean value of the sources for the limit cases LC1 and LC2



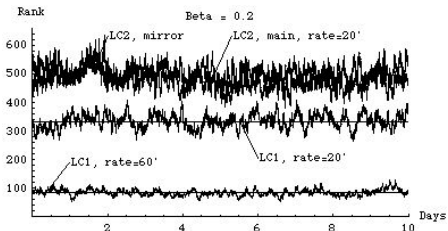
Algorithms TA3

- 1 TA3 is coherent with all the desirable properties P1–P5
- 2 unfortunately is more complicated than the others
- 3 it is not easy to write down a formula for the stationary mean value of the sources for the limit cases LC1 and LC2



TA3 behavior on the limit cases

Limit cases LC1 and LC2 are satisfied.



Sensitivity to the parameters

A first group of experiments address the sensitivity at changes of the parameters ρ and β

Algorithm TA3 is not much sensitive to changes in the parameters involved



Sensitivity to the parameters

A first group of experiments address the sensitivity at changes of the parameters ρ and β

ρ is the half-life decay time, that is

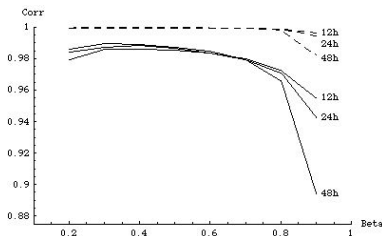
$$e^{-\alpha\rho} = \frac{1}{2}$$

Algorithm TA3 is not much sensitive to changes in the parameters involved



Sensitivity to the parameters

A first group of experiments address the sensitivity at changes of the parameters ρ and β

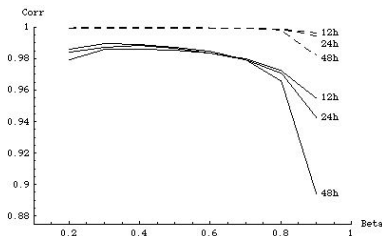


Algorithm TA3 is not much sensitive to changes in the parameters involved



Sensitivity to the parameters

A first group of experiments address the sensitivity at changes of the parameters ρ and β



Algorithm TA3 is not much sensitive to changes in the parameters involved



TA3 vs NTA1

It is necessary to have such a complicate algorithm?



TA3 vs NTA1

It is necessary to have such a complicate algorithm?

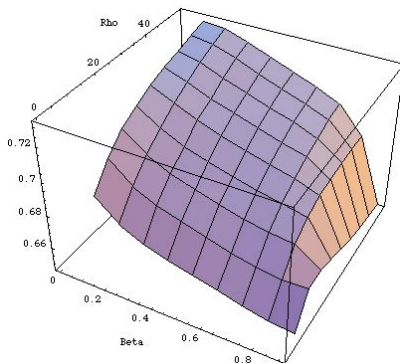
Yes



Sensitivity to the parameters

TA3 vs NTA1

It is necessary to have such a complicate algorithm?

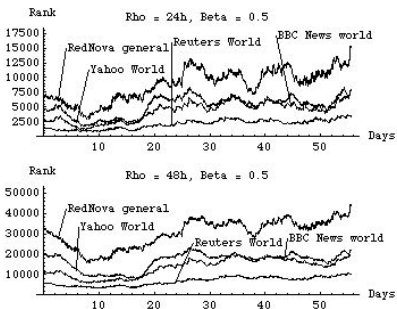


Lower correlation for large values of β and low values of ρ



Ranking news articles and news sources

Rank evolution over 55 days of the top 4 sources in the category World



Top News Sources

Source	# Postings
RedNova general	3154
Yahoo World	1924
Reuters World	1363
Yahoo Politics	900
BBC News world	1368
Reuters	555
Xinhua	339
New York Times world	549

Remark

Some news agency are considered more important than others even if they release a lower number of pieces of news



Top News Articles

Note

For top pieces of news it is common to recognize the same piece of information re-posted by other agencies

The rank of a singular news article is deeply dependent on the rank of the source posting it



Conclusions

- Algorithms for ranking News articles and news agencies
- Step-by-step construction
- Extensive testing on more than 300,000 pieces of news and 2,000 news sources
- The ranking is done online
- Same ideas for ranking publications, authors and scientific journals etc



Conclusions

- Algorithms for ranking News articles and news agencies
- **Step-by-step construction**
- Extensive testing on more than 300,000 pieces of news and 2,000 news sources
- The ranking is done online
- Same ideas for ranking publications, authors and scientific journals etc



Conclusions

- Algorithms for ranking News articles and news agencies
- Step-by-step construction
- Extensive testing on more than 300,000 pieces of news and 2,000 news sources
- The ranking is done online
- Same ideas for ranking publications, authors and scientific journals etc



