



HAL
open science

Priors Based On Time-Sensitive Social Signals

Ismail Badache, Mohand Boughanem

► **To cite this version:**

Ismail Badache, Mohand Boughanem. Priors Based On Time-Sensitive Social Signals. 37th European Conference on Information Retrieval (ECIR 2015), Mar 2015, Vienna, Austria. . hal-03154382

HAL Id: hal-03154382

<https://hal.science/hal-03154382>

Submitted on 28 Feb 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Priors Based On Time-Sensitive Social Signals

Ismail Badache and Mohand Boughanem
IRIT - Paul Sabatier University, Toulouse, France
{Badache, Boughanem}@irit.fr

1. Introduction

- Majority of search engines include social signals (e.g. +1, like) as non-textual features to relevance. However, in the existing works signals are considered time-independent.
- Hypothesis 1:** signals are time-dependent, the date when the user action has happened is important to distinguish between recent and old signals. Therefore, the recency of signals may indicate some recent interests to the resource, which may improve the a priori relevance of document.
- Hypothesis 2:** the number of signals of a resource depends on the resource age, an old resource may have much more signals than a recent one.

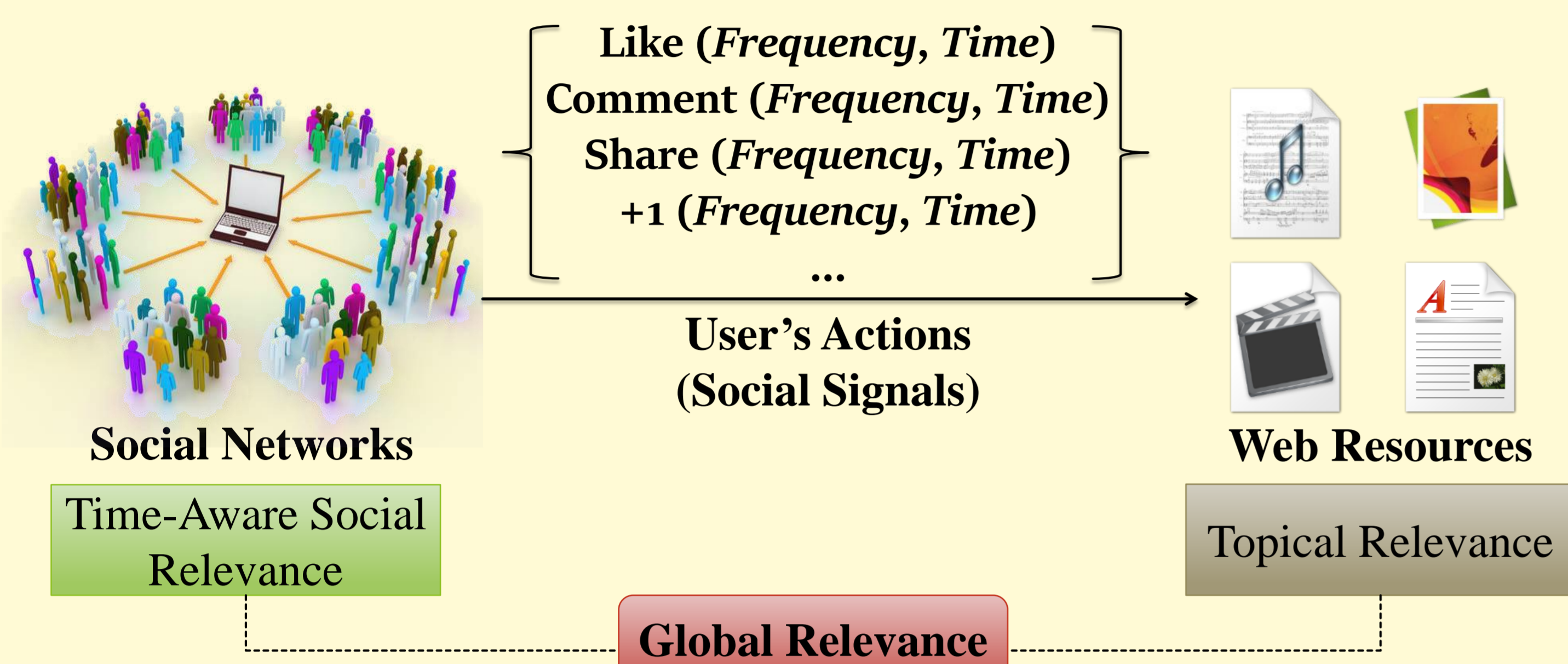


Fig. Global presentation of our approach

Research questions are the following:

- How to take into account signals and their date to estimate the priors?
- What is the impact of temporally-aware signals on IR system performance?

3. Time-Aware Social Signals

- We propose to consider the date associated with a signal and the age of a resource. To estimate priors, we distinguish two ways to handle it:

Proposal 1: Time of Signal

Resource associated with fresh signals are more likely to interest user and should be promoted comparing to those associated with old signals. Therefore, instead of counting each occurrence of a given signal, we bias the counting, noted $Count_{t_a}$, by the date of the occurrence of the signal.

$$Count_{t_a}(t_{j,a_i}, D) = \sum_{j=1}^k f(t_{j,a_i}, D) = \sum_{j=1}^k \exp\left(-\frac{\|t_{current} - t_{j,a_i}\|^2}{2\sigma^2}\right) \quad (4)$$

Where $f(t_{j,a_i}, D)$ represents signal-time function, we use Gaussian Kernel to estimate a distance between current time $t_{current}$ and t_{j,a_i} with $\sigma \in \mathbb{R}^+$. k represents the number of moments (*datetime*) at which action a_i was produced.

- $P(D)$ is estimated using formula 3 but by replacing $Count()$ by $Count_{t_a}()$. Notice that if the signal time is not considered $f(t_{j,a_i}, D) = 1 \forall t_{j,a_i}$

Proposal 2: Age of Resource

The simple counting of signals may boost old resources compared to recent ones, because resources with long life in the Web has much more chance to get more signals than recent ones. So to cope with this issue we propose to normalize the distribution of signals associated with a resource through resource publication date. We divide the number of signals by the current lifespan of the resource.

$$Count_{t_D}(a_i, D) = \frac{Count(a_i, D)}{Age(D)} = \frac{Count(a_i, D)}{\exp\left(-\frac{\|t_{current} - t_D\|^2}{2\sigma^2}\right)} \quad (5)$$

- The prior $P(D)$ is estimated using formula 3 but by replacing $Count()$ by $Count_{t_D}()$ for document and $Count_{t_C}$ for collection.

2. Time-Independent Social Signals

- We exploit language model to estimate the relevance of document D to a query Q .

$$P(D|Q) = Rank P(D) \cdot P(Q|D) = P(D) \cdot \prod_{w_i \in Q} P(w_i|Q) \quad (1)$$

- $P(D)$ is a document prior. w_i represents words of query Q . Estimating $P(w_i|Q)$ can be performed using different models.
- The priors are estimated by a counting of actions a_i performed on D .

$$P(D) = \prod_{a_i \in A} P(a_i) \quad (2)$$

With $P(a_i)$ is estimated using maximum-likelihood: $P(a_i) = \frac{Count(a_i, D)}{Count(a_*, D)}$

- We smooth $P(a_i)$ by collection C using *Dirichlet*:

$$P(D) = \prod_{a_i \in A} \left(\frac{Count(a_i, D) + \mu \cdot P(a_i|C)}{Count(a_*, D) + \mu} \right) \quad (3)$$

With $P(a_i|C)$ is estimated using maximum-likelihood: $P(a_i|C) = \frac{Count(a_i, C)}{Count(a_*, C)}$

Where $P(D)$ represents the a priori probability of D . $Count(a_i, D)$ represents number of occurrence of action a_i on resource D . a_* is the total number of social signals in document D or in collection C .

4. Experimental Evaluation

Dataset

- INEX IMDb Dataset.
- 30 INEX IMDb Topics and their relevance judgments.
- Social data from 5 social networks: *Like*, *Share* and *Comment* (Facebook); *Tweet* (Twitter); *+1* (Google+); *Bookmark* (Delicious); *Share* (LinkedIn).

Results

