

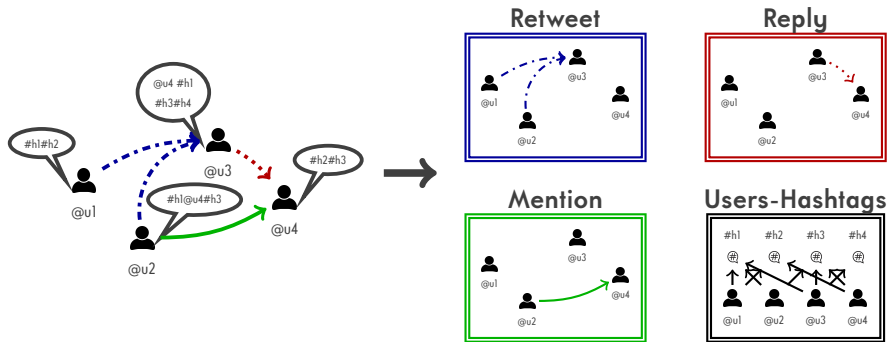
# Link weights recovery in heterogeneous information networks

Hông-Lan Botterman<sup>1</sup>, Robin Lamarche-Perrin<sup>2</sup>

<sup>1</sup>Sorbonne Université (LIP6)

<sup>2</sup>Institut des Systèmes Complexes Paris Île de France

# Link weights recovery: an illustrative example



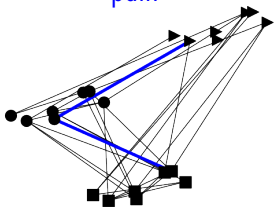
Expressing the UH link weights using other links/link weights

For some users, recovering their UH link weights, knowing some other links/link weights.

## heterogeneous information network

$$H := (V, E, w, \mu_s, \mu_t, \mathcal{V}, \mathcal{E}, \phi, \psi)$$

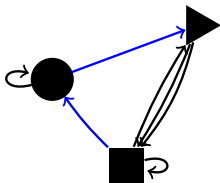
path



## network schema

$$T_H := (\mathcal{V}, \mathcal{E}, \nu_s, \nu_t)$$

metapath: ■ → ● → ►



## path-constrained random walk

$$\mathcal{V} = \{V_1, \dots, V_m\} \text{ and } \mathcal{E} = \{E_1, \dots, E_r\}$$

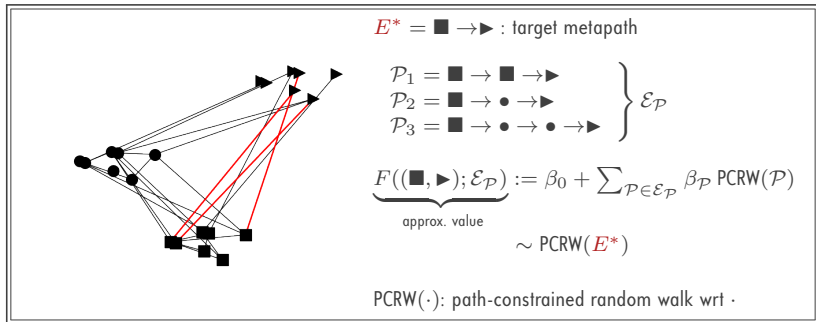
$$\mathcal{P} = V_1 \xrightarrow{E_{j_1}} V_{i_2} \dots V_{i_{n-1}} \xrightarrow{E_{j_{n-1}}} V_n, \quad i_2, \dots, i_{n-1} \in \{1, \dots, m\}, \\ j_1, \dots, j_{n-1} \in \{1, \dots, r\}.$$

$$\mathbb{P}((v_n|v_1) | \mathcal{P}) = \sum_{v_{n-1} \in V_{i_{n-1}}} \frac{w_{E_{j_{n-1}}}(v_{n-1}, v_n) f_{v_n}^\alpha}{\sum_k w_{E_{j_{n-1}}}(v_{n-1}, v_k) f_{v_k}^\alpha} \mathbb{P}\left((v_{n-1}|v_1) | \mathcal{P}^{1, i_{n-1}}\right)$$

# Linear combination of PCRW

Dependent variable: result w.r.t. a target link type -  $\text{PCRW}(E^*)$

Independent variables: results of random walks w.r.t. particular link types -  
 $\text{PCRW}(\mathcal{P}), \mathcal{P} \in \mathcal{E}_{\mathcal{P}}$



# Selection by a forward linear regression

## 1. Least squares problem:

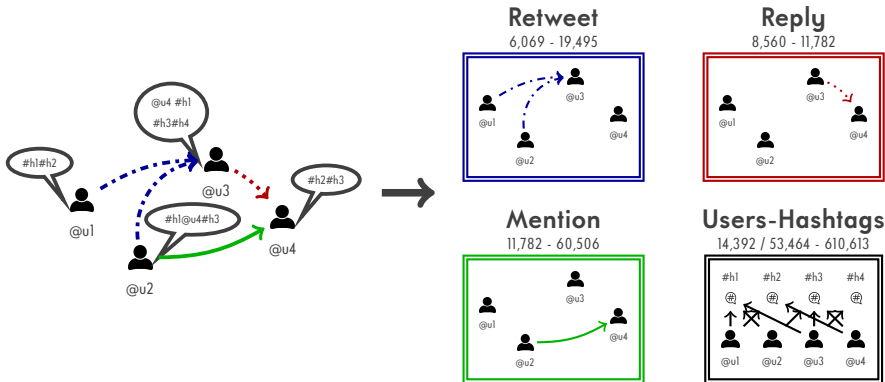
$$\text{Let } S = \sum r_i^2: \quad \frac{\partial S}{\partial \beta_j} = 2 \sum r_i \frac{\partial r_i}{\partial \beta_j} = 0$$

$$\rightarrow \hat{\beta} = (X^T X)^{-1} X^T y$$

## 2. Forward linear regression:

- starting with no variables in the model
- testing the addition of each variable using a chosen model fit criterion:  
*p*-value and *t*-test:  $t^* = (\bar{x} - \mu_0)/(s/\sqrt{n})$
- adding the variable (if any) whose inclusion gives the most statistically significant improvement of the fit:  
coef. of determination  $R^2 = 1 - \sum_{i=1}^n (y_i - \hat{y}_i)^2 / \sum_{i=1}^n (y_i - \bar{y}_i)^2$
- repeating this process until there is no further improvement

# Twitter data related to FIFA World Cup 2014



From June 12 to July 13, 2014  
32 teams - 64 matches played

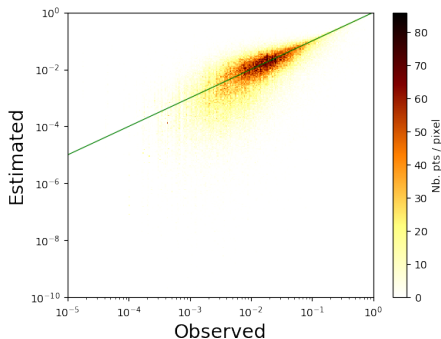
HIN with  $\mathcal{V} = \{\text{users, hashtags}\}$  and  $\mathcal{E} = \{\text{retweet, reply, mention, post}\}$ .

# Describing UH from other link types (locally)

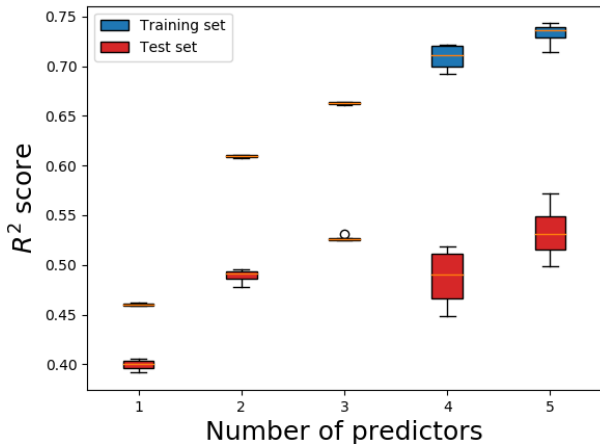
$\mathcal{E}_{\mathcal{P}} = \{ \text{Metapaths of length } \leq 4 \}$

UH = 0.1974 MT-UH (iter 1)  
+ 0.5556 RP-UH (iter 2)  
+ 0.0650 RT-RP-UH (iter 3)  
+ 0.1591 RP-MT-UH (iter 4)  
+ 0.0074 MT-RT-UH (iter 5)

with  $r^2 = 0.7129$



# Recovering UH from other link types and other users



Training set 80% - Test set: 20% of users

Overfitting



# Summary

## **Main idea**

Adequate linear combination of path-constrained random walks results to describe and, to some extent, retrieve the strength of the links between different entities.

## **Several improvements**

- Temporal aspects

- Searching for candidate metapaths

# Link weights recovery in heterogeneous information networks

Hông-Lan Botterman<sup>1</sup>, Robin Lamarche-Perrin<sup>2</sup>

<sup>1</sup>Sorbonne Université (LIP6)

<sup>2</sup>Institut des Systèmes Complexes Paris Île de France