# Mineração de Dados
## Aula 5 – parte 2

Especialização em Ciência de Dados e suas Aplicações

UTFPR

# Estudo de diferenças culturais

You are What you Eat (and Drink): Identifying Cultural Boundaries by Analyzing Food & Drink Habits in Foursquare.
ICWS'15

Thiago H. Silva

**Grande desafio:** encontrar dados apropriados para uso

- **Métodos tradicionais:** Questionários
  – Não escalam
  – Difícil de detectar mudanças dinâmicas

**Grande desafio:** encontrar dados apropriados para uso

- **Métodos tradicionais:** Questionários
  - Não escalam
  - Difícil de detectar mudanças dinâmicas

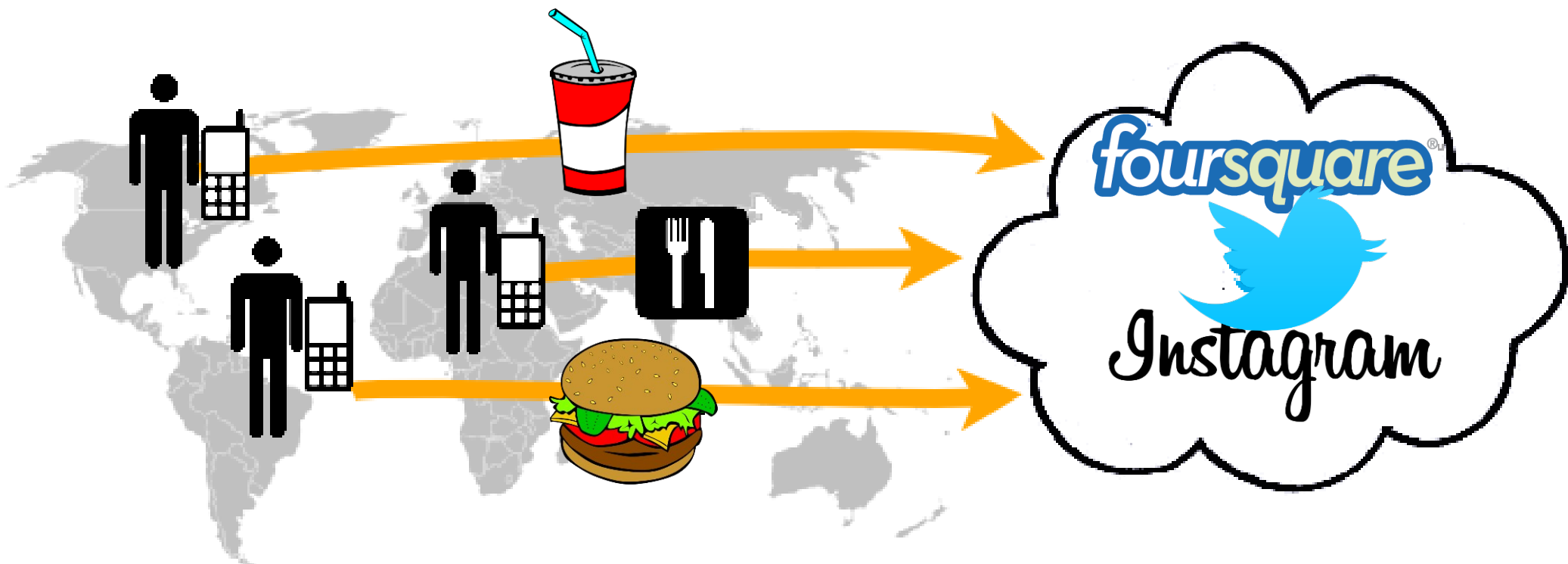É possível propor algum método alternativo?

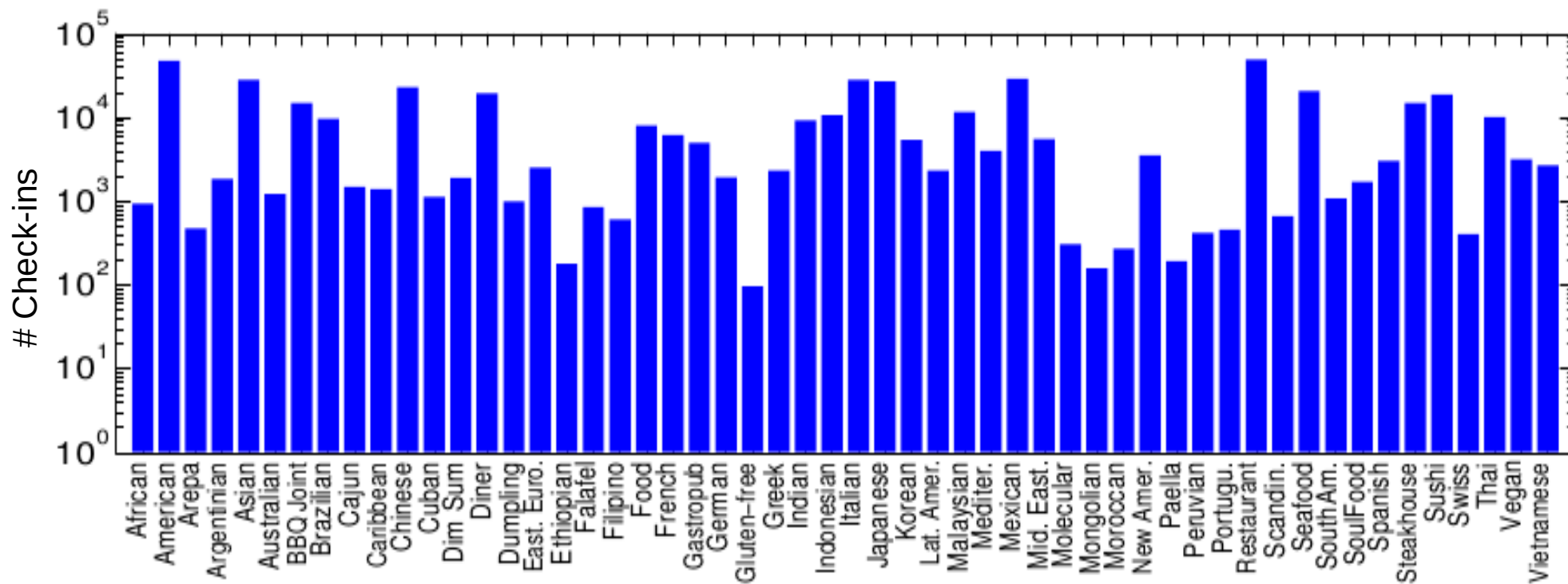Hábitos alimentares e de bebida são elementos fundamentais em uma cultura

Sensoriamento de atividades humanas em larga escala!



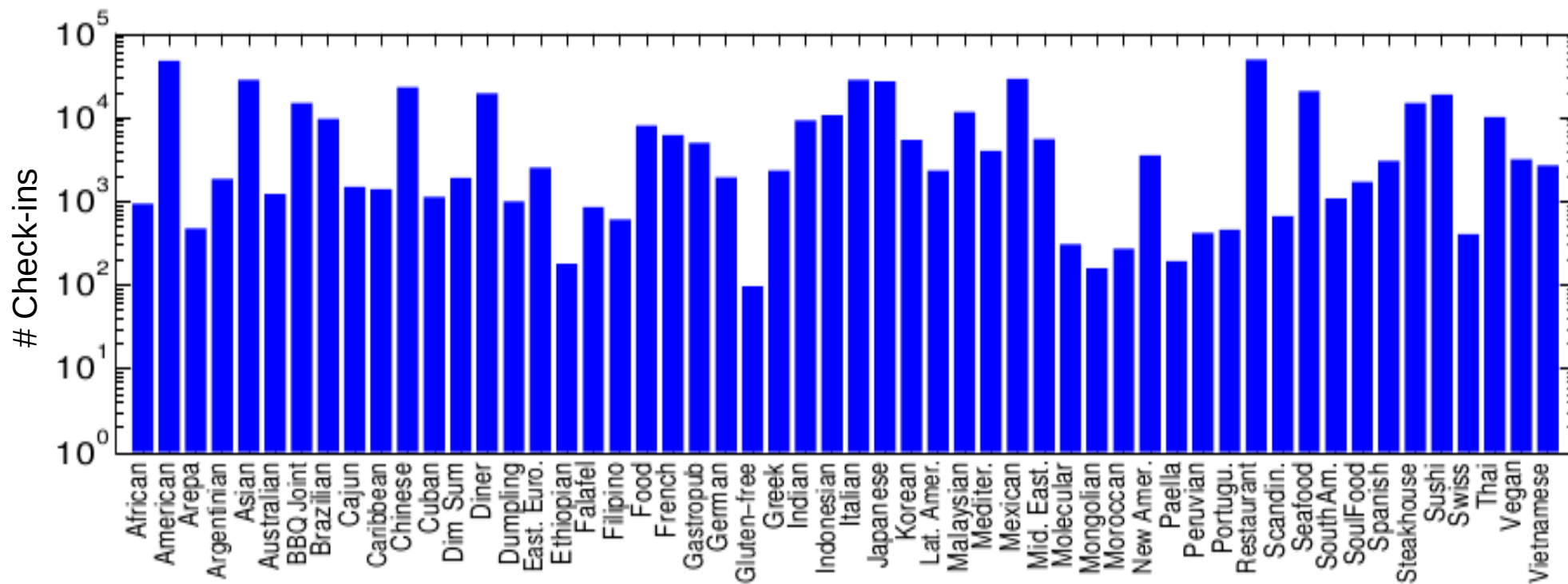Oportunidade sem precedentes para estudar diferenças culturais em escala global e baixo custo
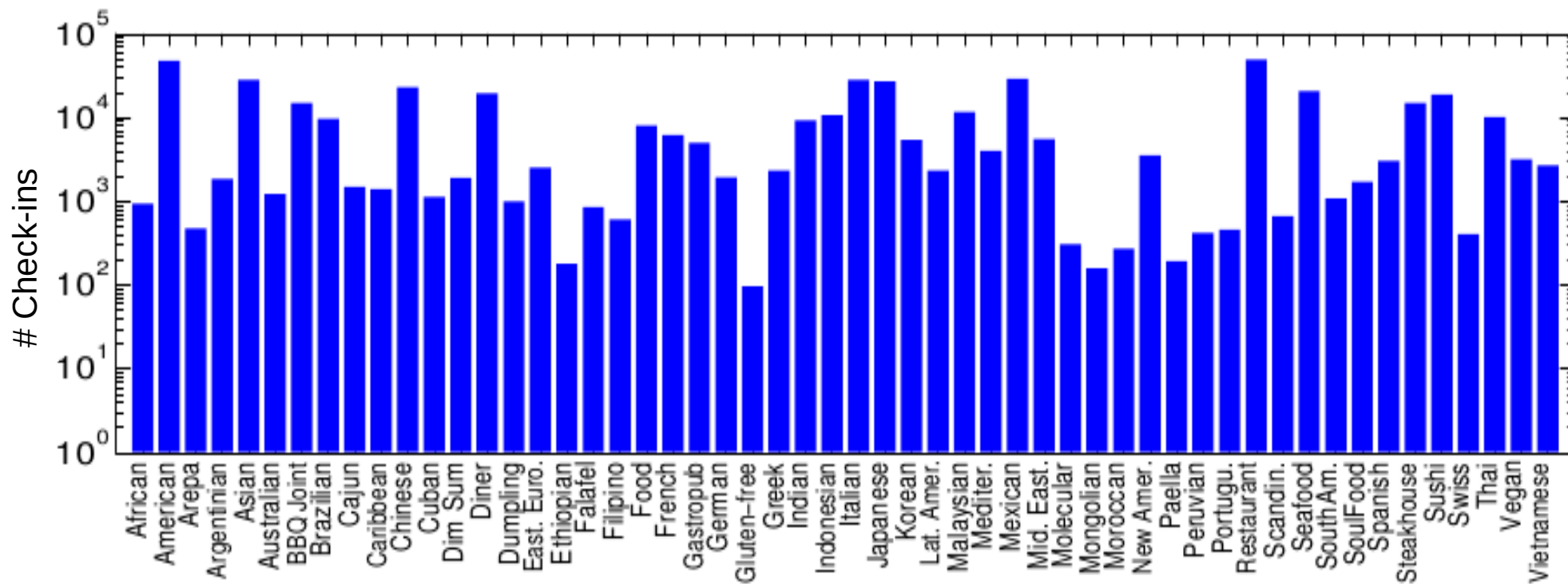
# Categorias Slow Food



Thiago H. Silva

Mapeia cada usuário $n_i$ n $\quad F_i = f_{1^i}, f_{2^i}, \ldots, f_{m^i}$
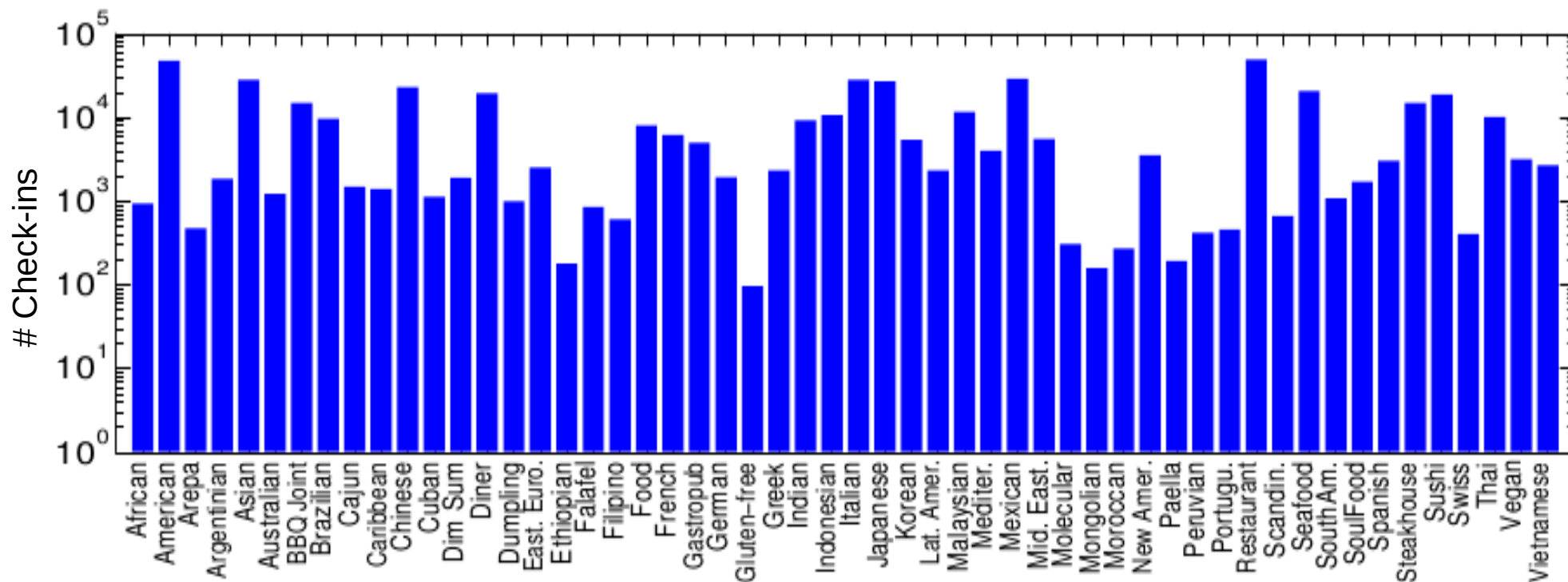
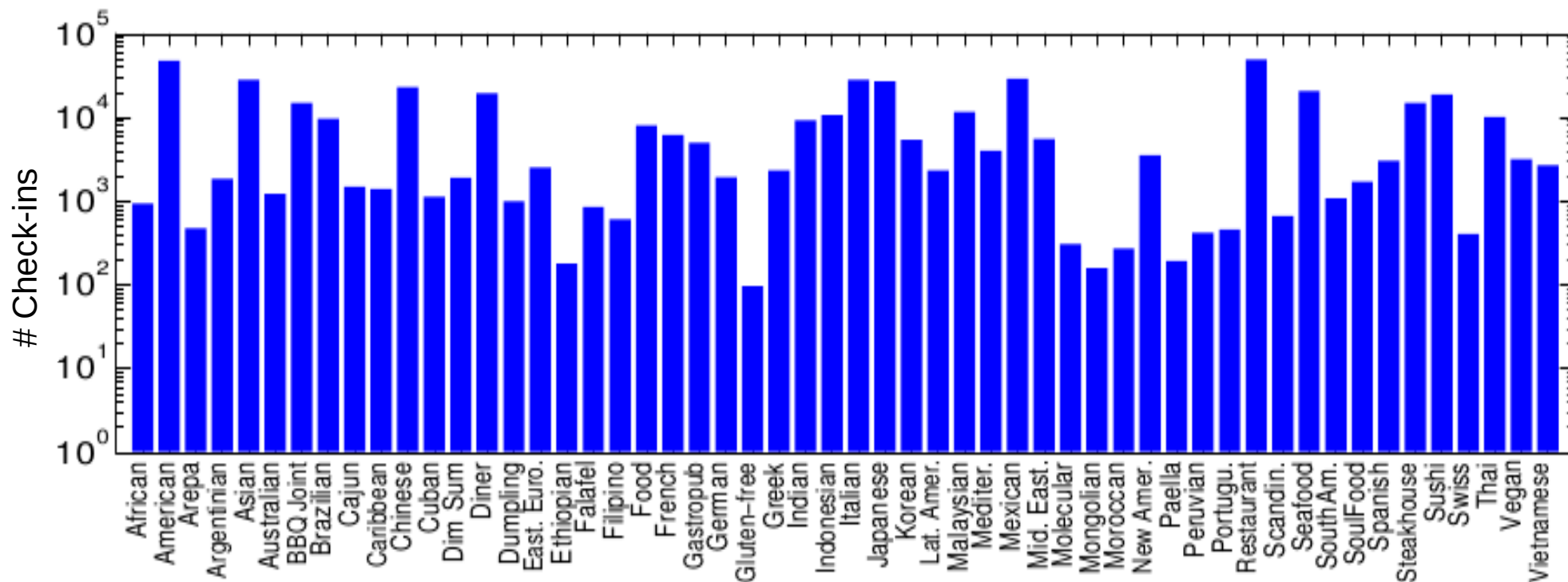Mapeia cada usuário $n_i$ $\mathfrak{n}$ $F_i = f_{1i}, f_{2i}, \ldots, f_{mi}$

Como perguntas em um survey!

Mapeia cada usuário $n_i$ ⋔ $F_i = f_{1^i}, f_{2^i}, \ldots, f_{m^i}$

$f_{k^i} = 0|1$ representa se o usuário $n_i$ gosta de $f_k$

Mapeia cada usuário $n_i$ ŋ $F_i = f_{1^i}, f_{2^i}, \ldots, f_{m^i}$

$f_{k^i} = 0|1$ representa se o usuário $n_i$ gosta de $f_k$

Respostas dos usuários

Data from LBSNs can be used if and only if:

# Data requirements

Data from LBSNs can be used if and only if:

**1** - Associate a user to its location;

Brazilian

Data from LBSNs can be used if and only if:

**1** - Associate a user to its location;
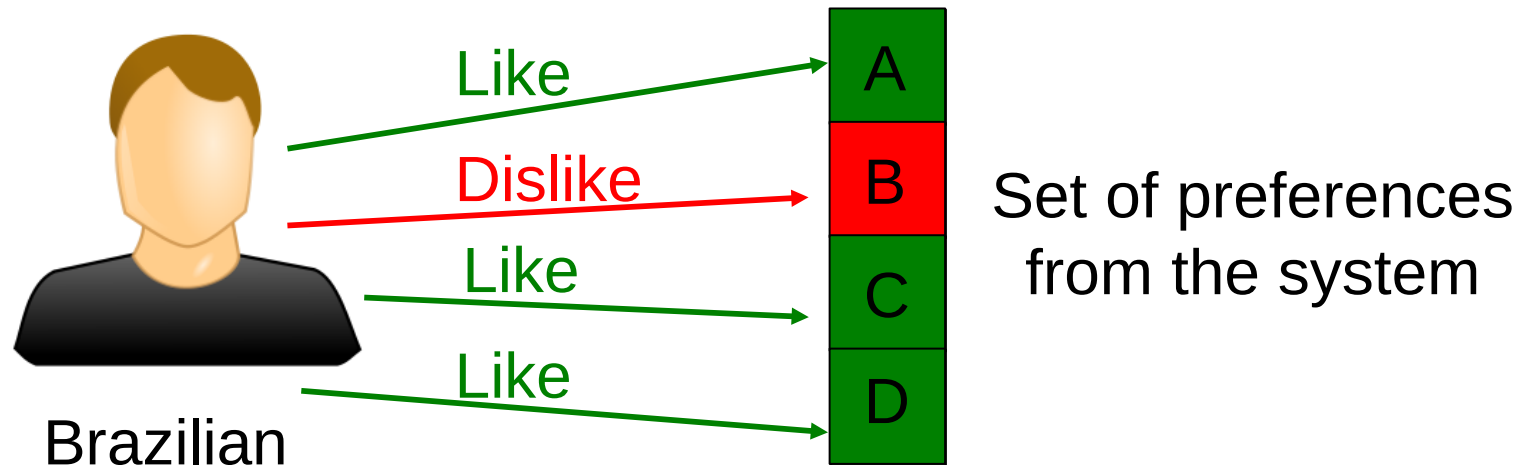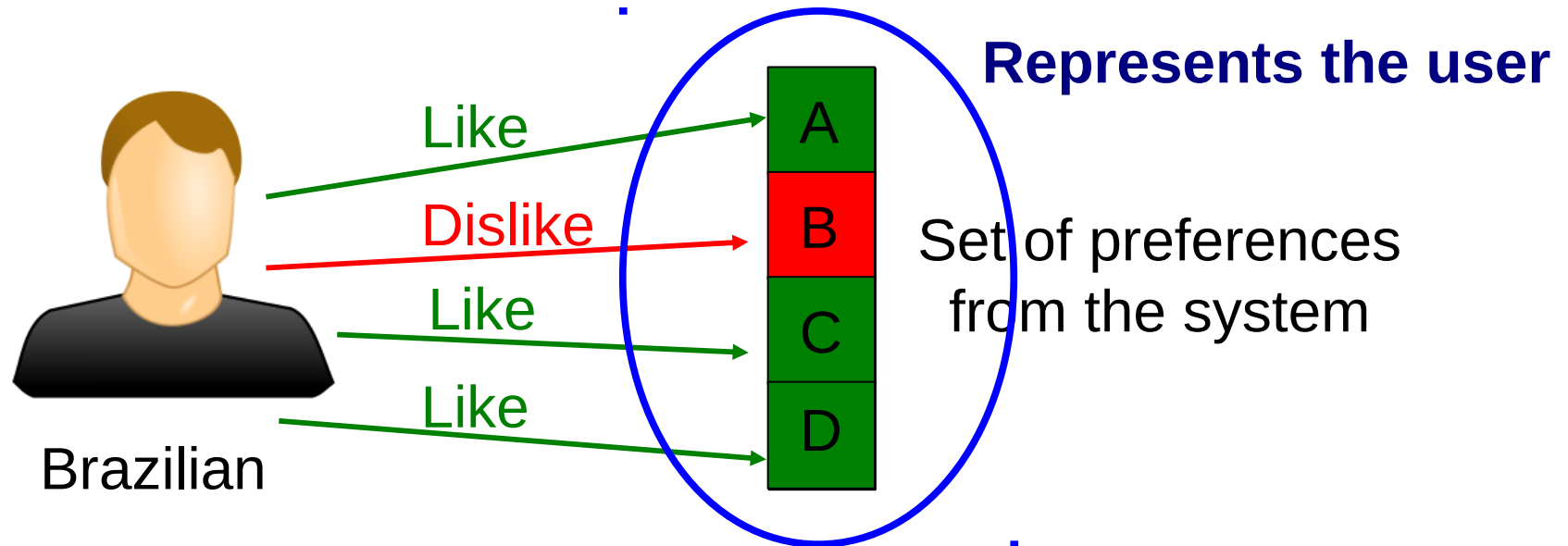**2** - Extract finite set of preferences from the data;

Brazilian

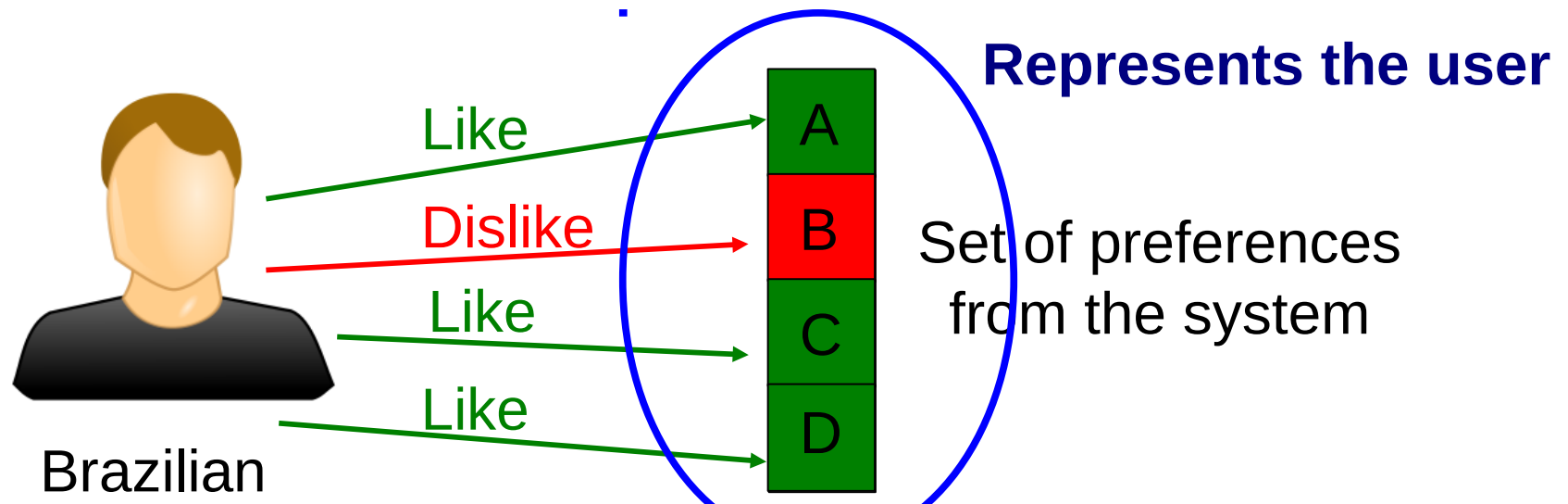| |
|---|
| A |
| B |
| C |
| D |

Set of preferences
from the system

Data from LBSNs can be used if and only if:

**1** - Associate a user to its location;
**2** - Extract finite set of preferences from the data;
**3** - Map users' actions into the preferences.



Brazilian

Like → A

Dislike → B

Like → C

Like → D

Set of preferences from the system

# Data requirements

Data from LBSNs can be used if and only if:

**1** - Associate a user to its location;
**2** - Extract finite set of preferences from the data;
**3** - Map users' actions into the preferences.



**Represents the user**

Set of preferences from the system

Like

Dislike

Like

Like

Brazilian

A B C D

Data from LBSNs can be used if and only if:

**1** - Associate a user to its location;
**2** - Extract finite set of preferences from the data;
**3** - Map users' actions into the preferences.



**Represents the user**

Like

Dislike

Like

Like

A

B

C

D

Set of preferences from the system

Brazilian

We demonstrate with Foursquare data

Thiago H. Silva

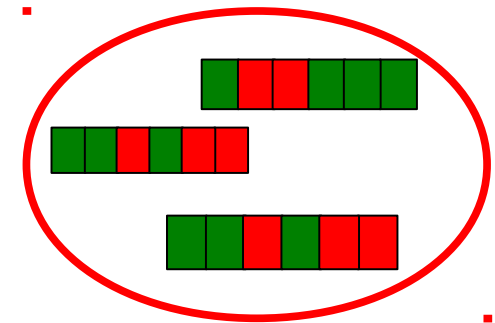**Can we define cultural signatures of different areas around the world?**

Spatial evaluation

For a given geographical area:



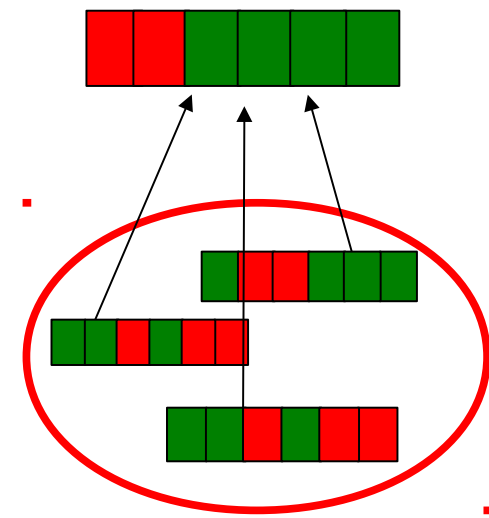Area a (Belo Horizonte)

Area b (New York City)

Spatial evaluation

For a given geographical area:

- Aggregate all users' preference in normalized vectors
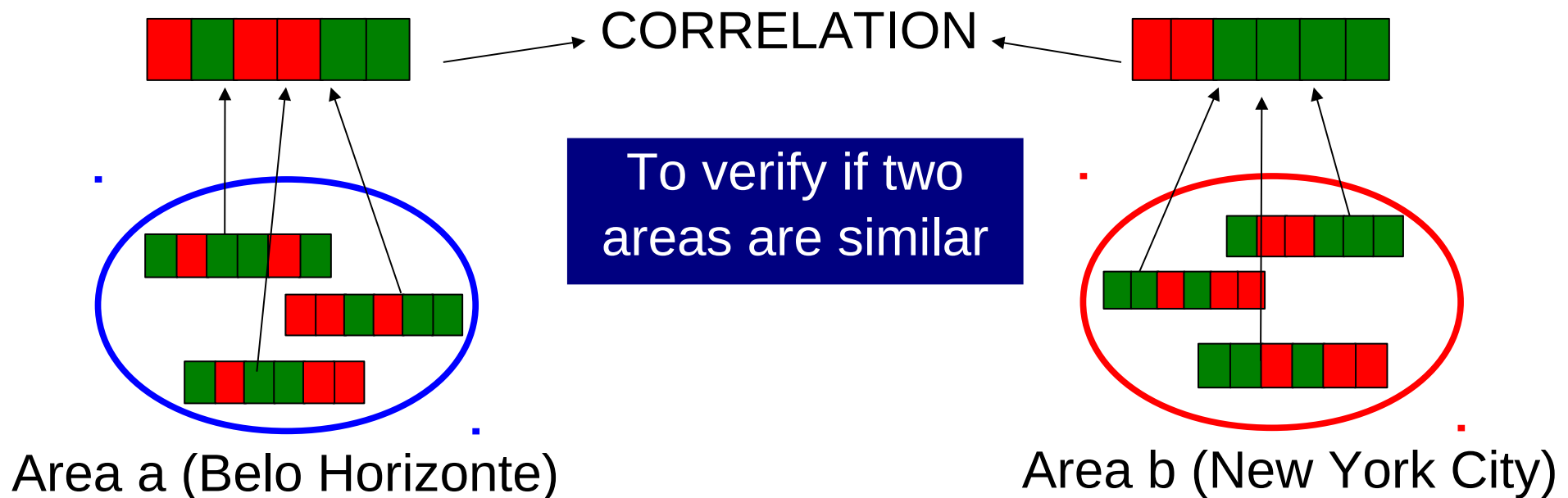


Area a (Belo Horizonte)

Area b (New York City)

# Extraction of cultural signatures

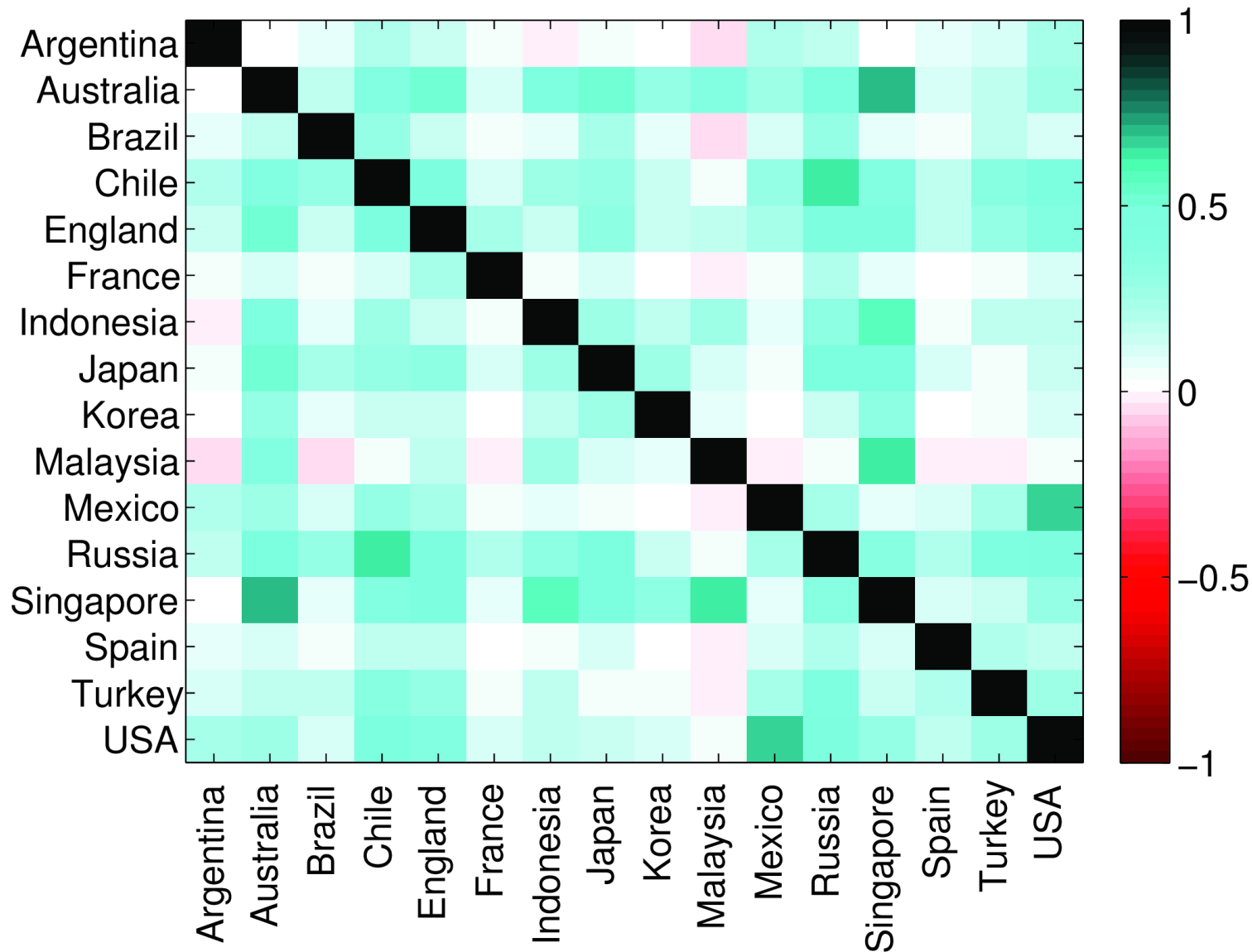Spatial evaluation

For a given geographical area:

- Aggregate all users' preference in normalized vectors



CORRELATION

To verify if two areas are similar

Area a (Belo Horizonte)

Area b (New York City)

Thiago H. Silva – DCC/UFMG

# Extraction of cultural signatures

Results for **countries**
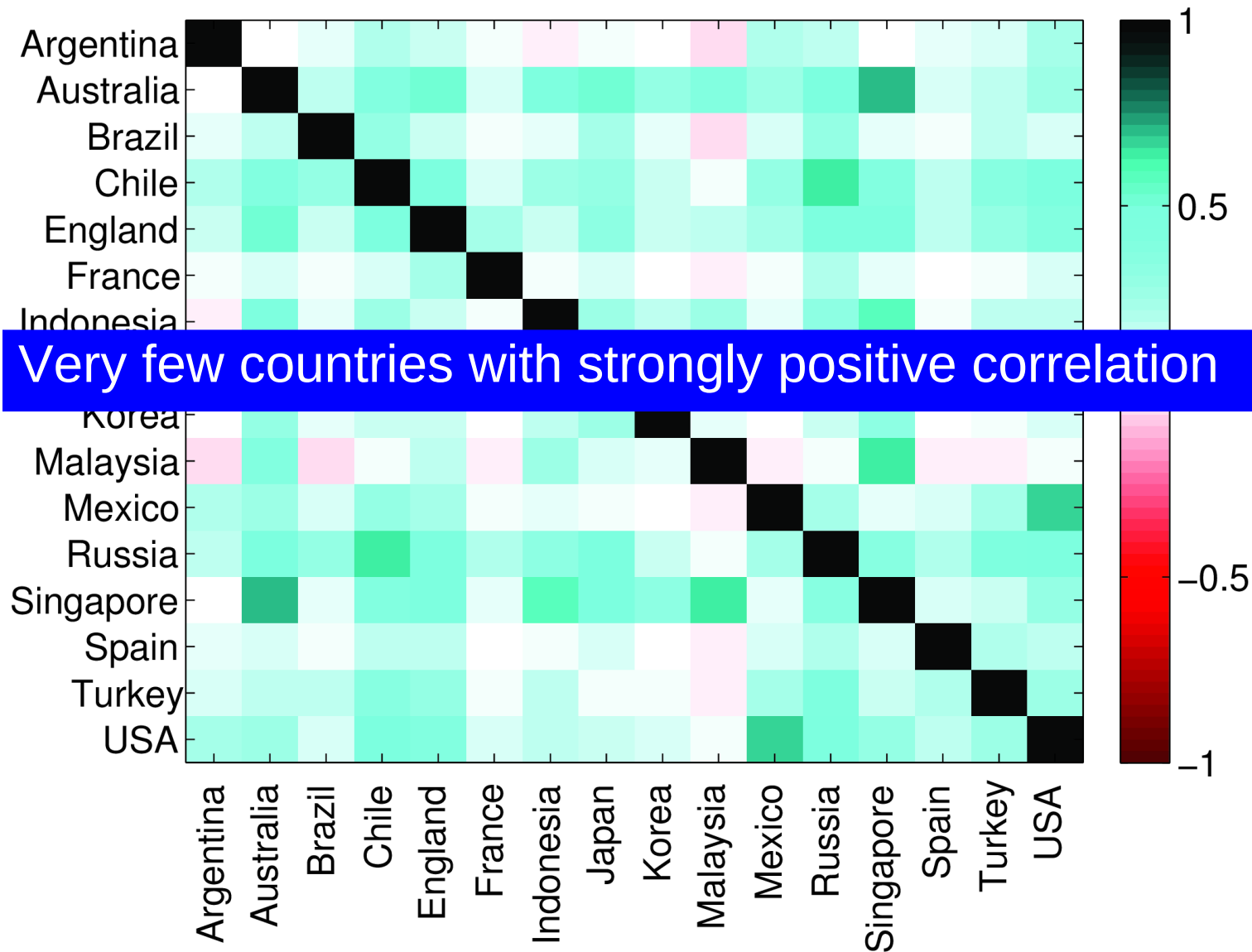
# Extraction of cultural signatures



Spatial evaluation

Results for **countries**
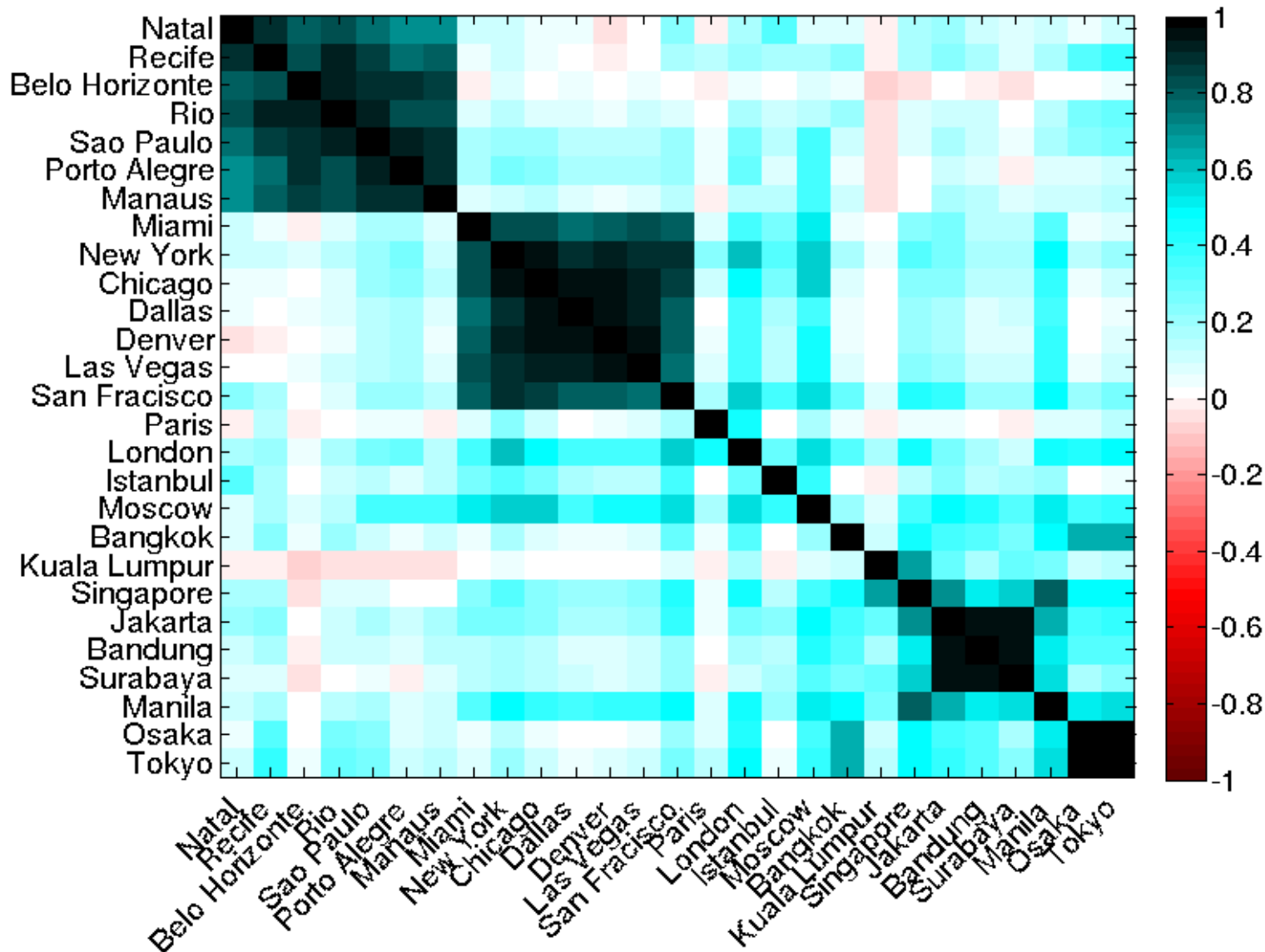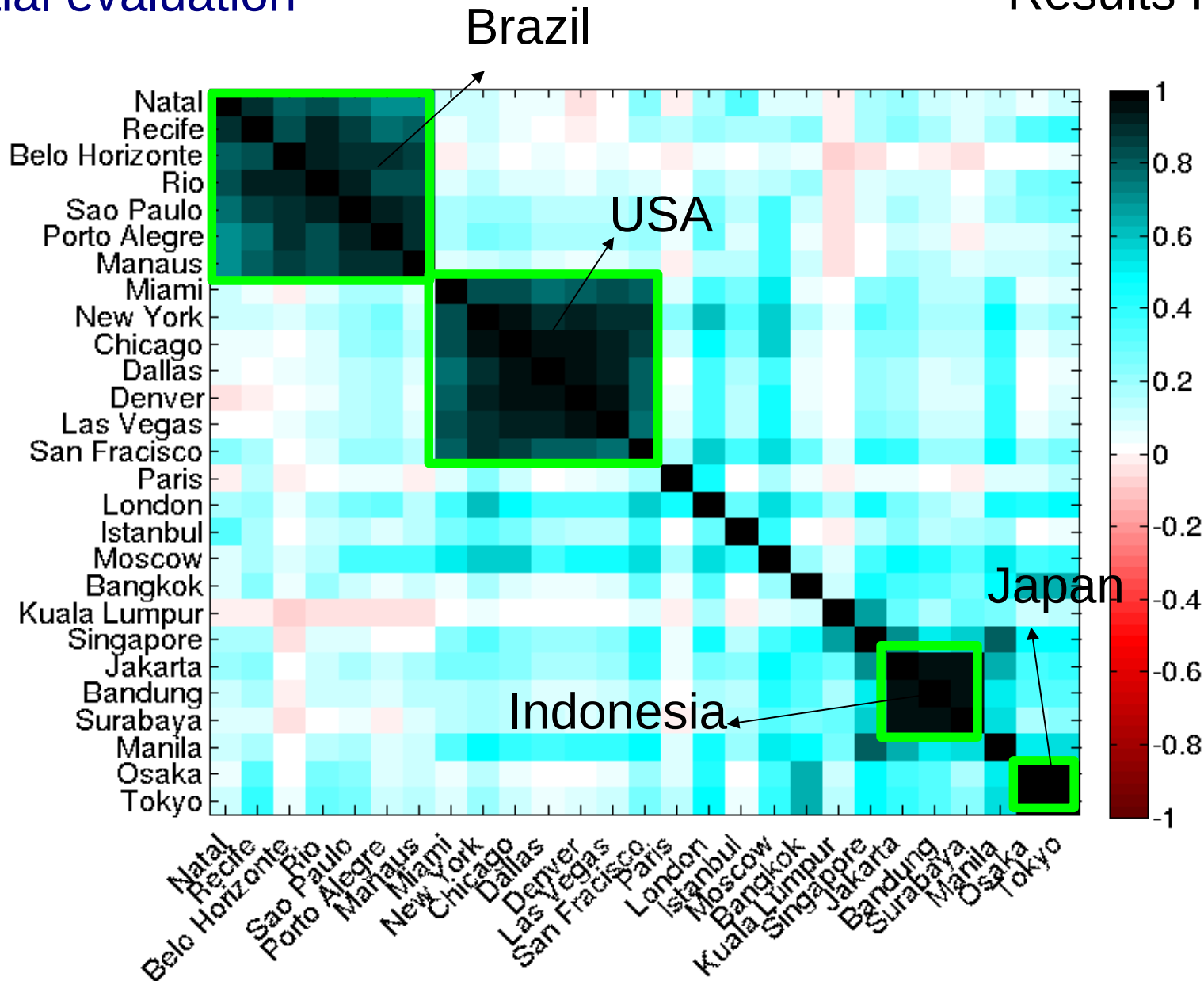
Positive correlation

Negative correlation

# Extraction of cultural signatures

Spatial evaluation

Results for **countries**



Very few countries with strongly positive correlation

**Spatial evaluation**

Results for **cities**

# Extraction of cultural signatures

**Spatial evaluation**

Results for **cities**

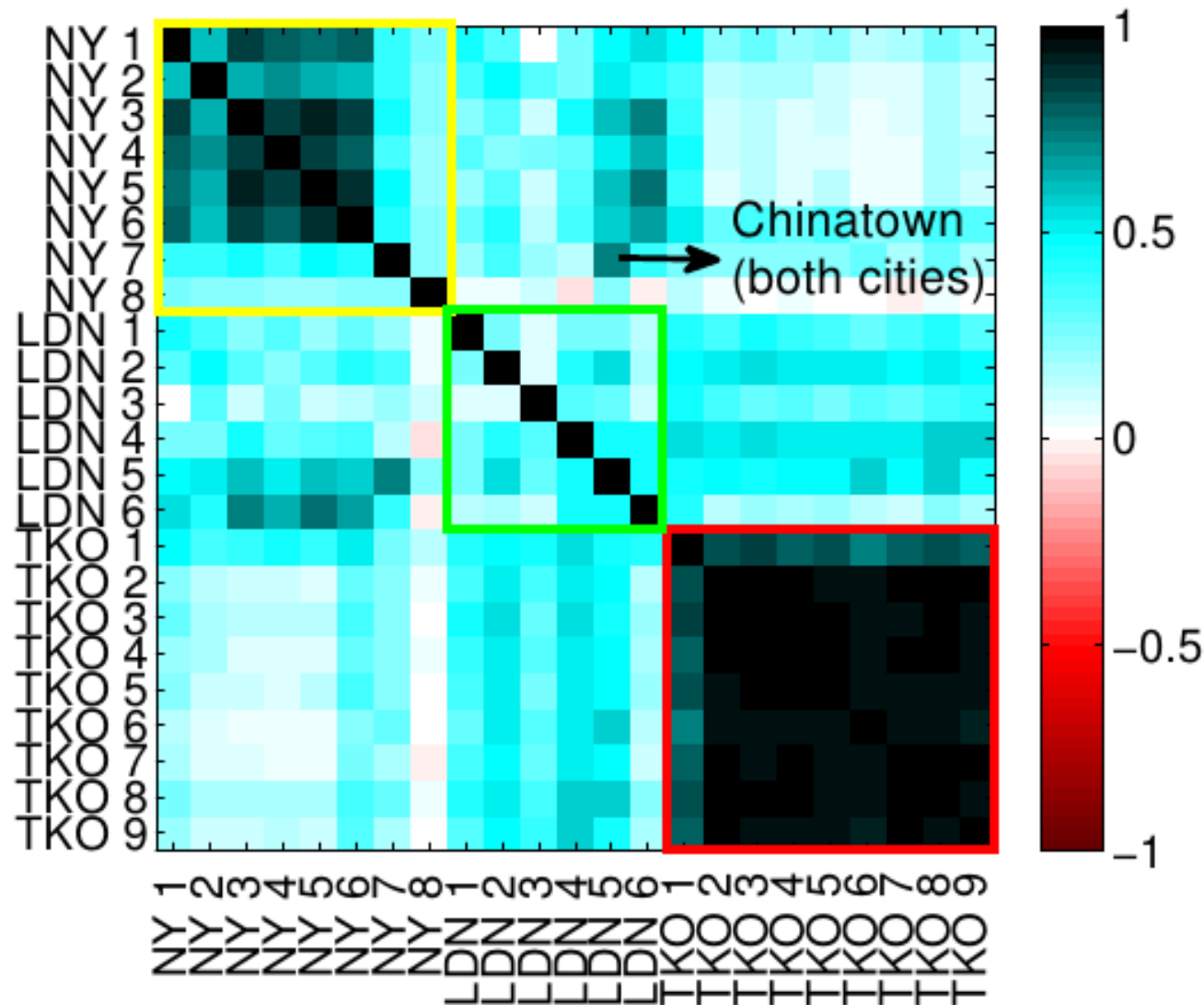Spatial evaluation



New York (NY)     Tokyo (TKO)

Popular areas

London (LND)

Spatial evaluation
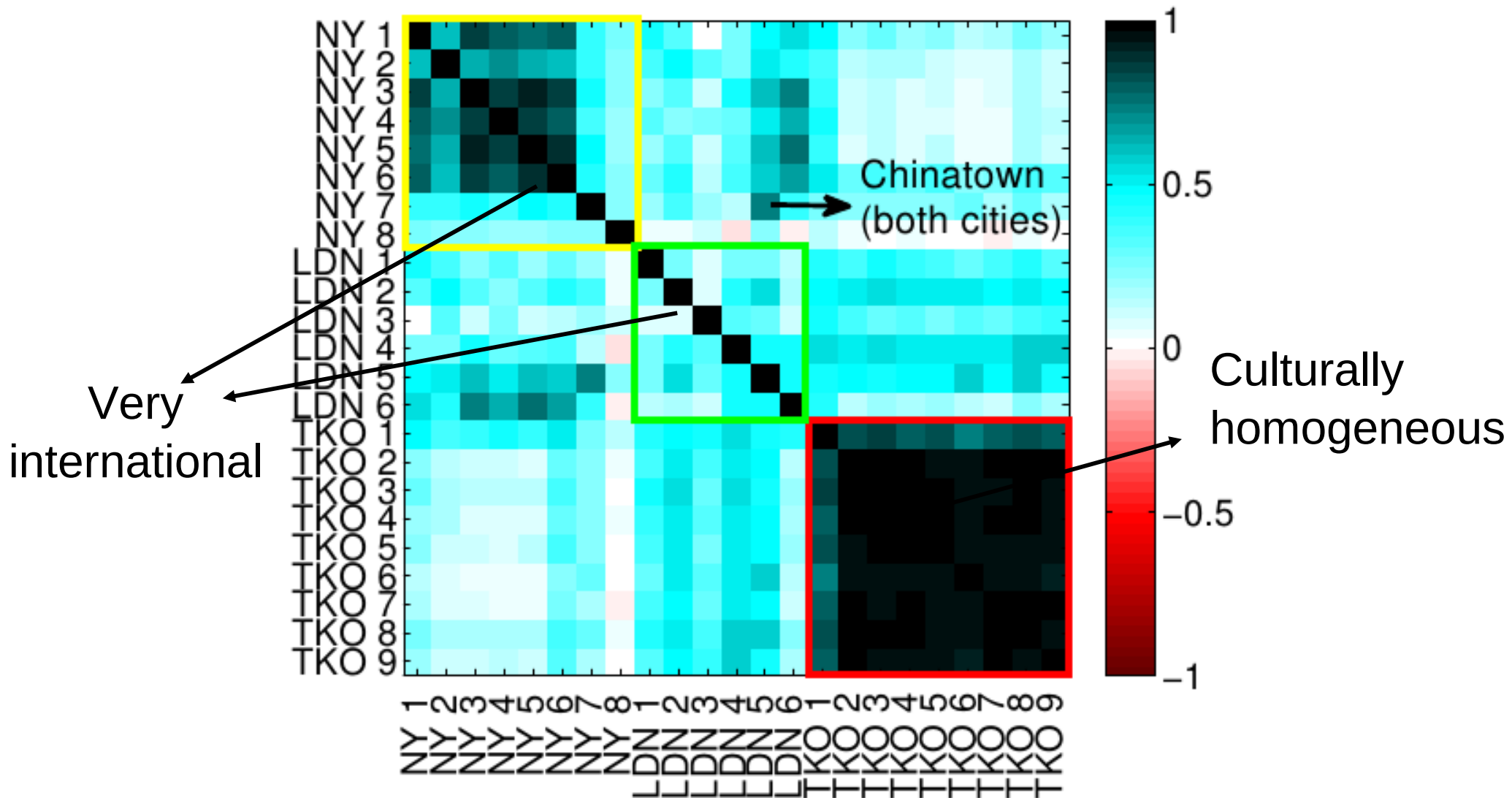
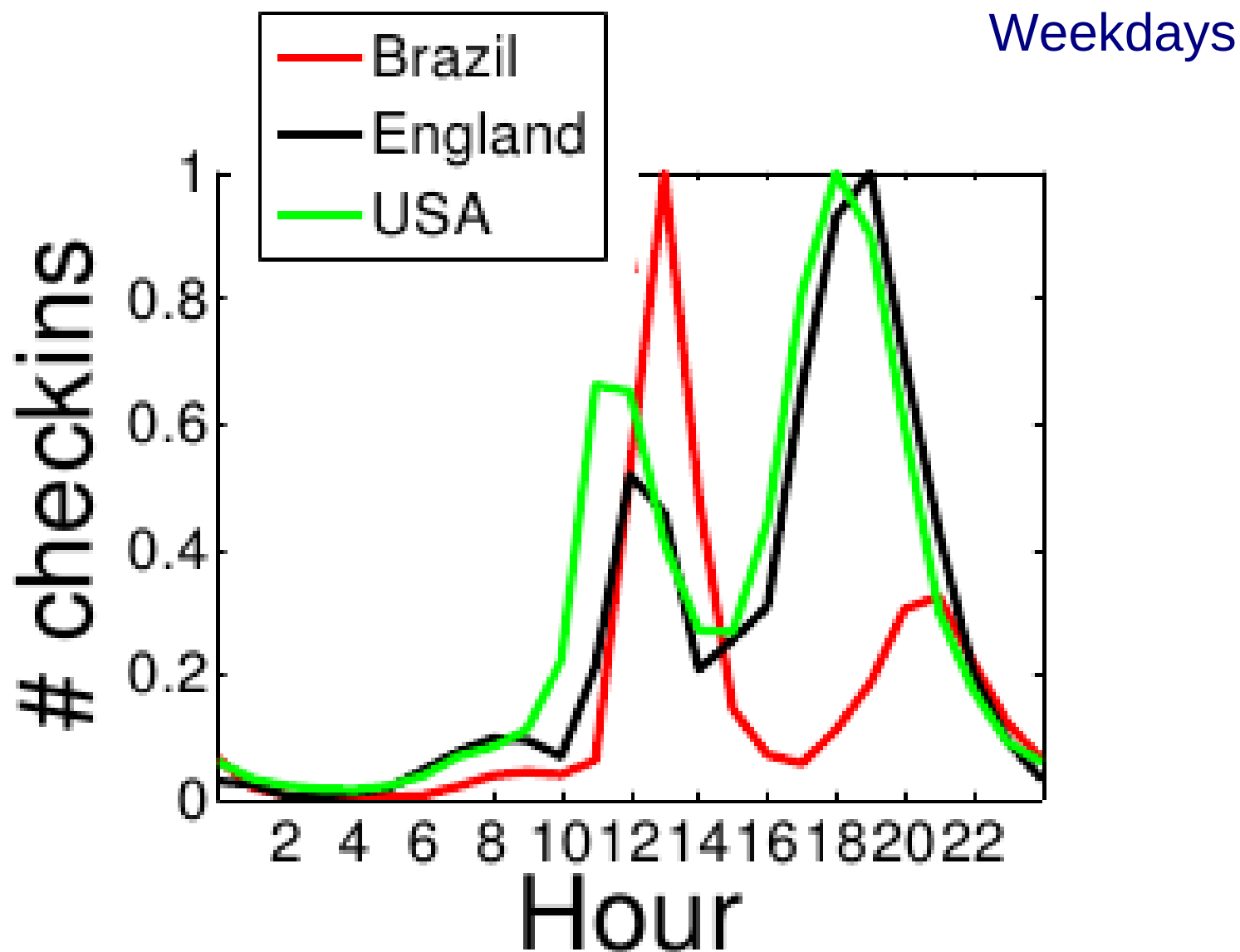Results for **areas inside cities**

# Extraction of cultural signatures

Spatial evaluation

Results for **areas inside cities**



Chinatown
(both cities)

Very
international

Culturally
homogeneous

Thiago H. Silva – DCC/UFMG
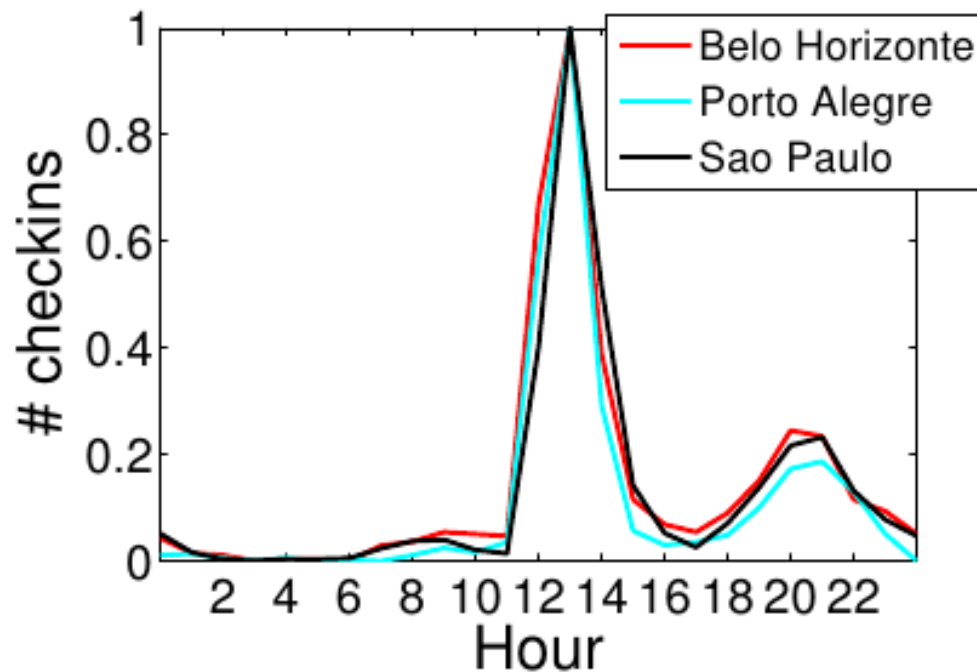
# Extraction of cultural signatures

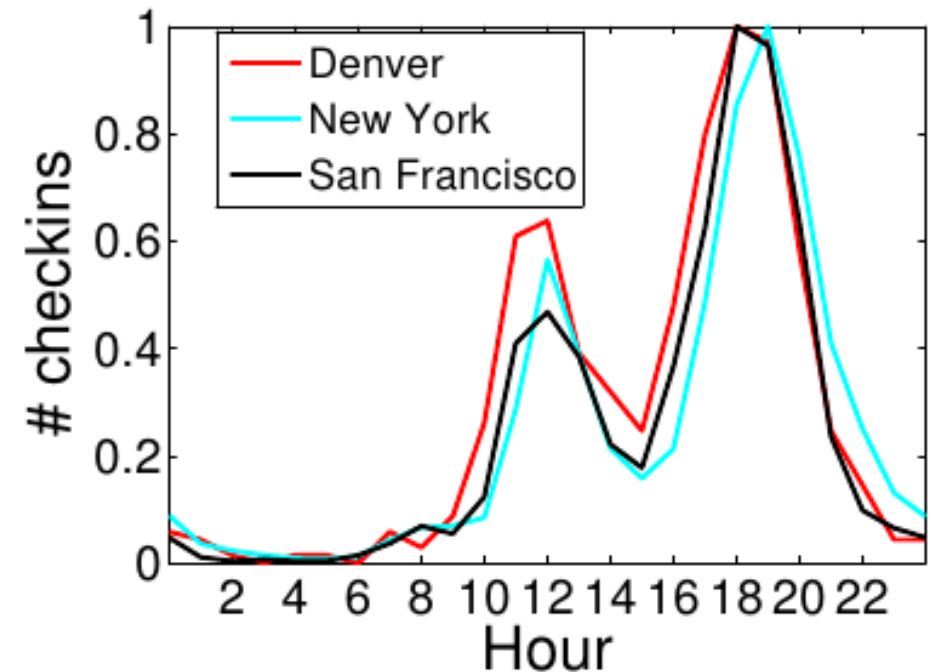Temporal evaluation

Weekdays

# Extraction of cultural signatures
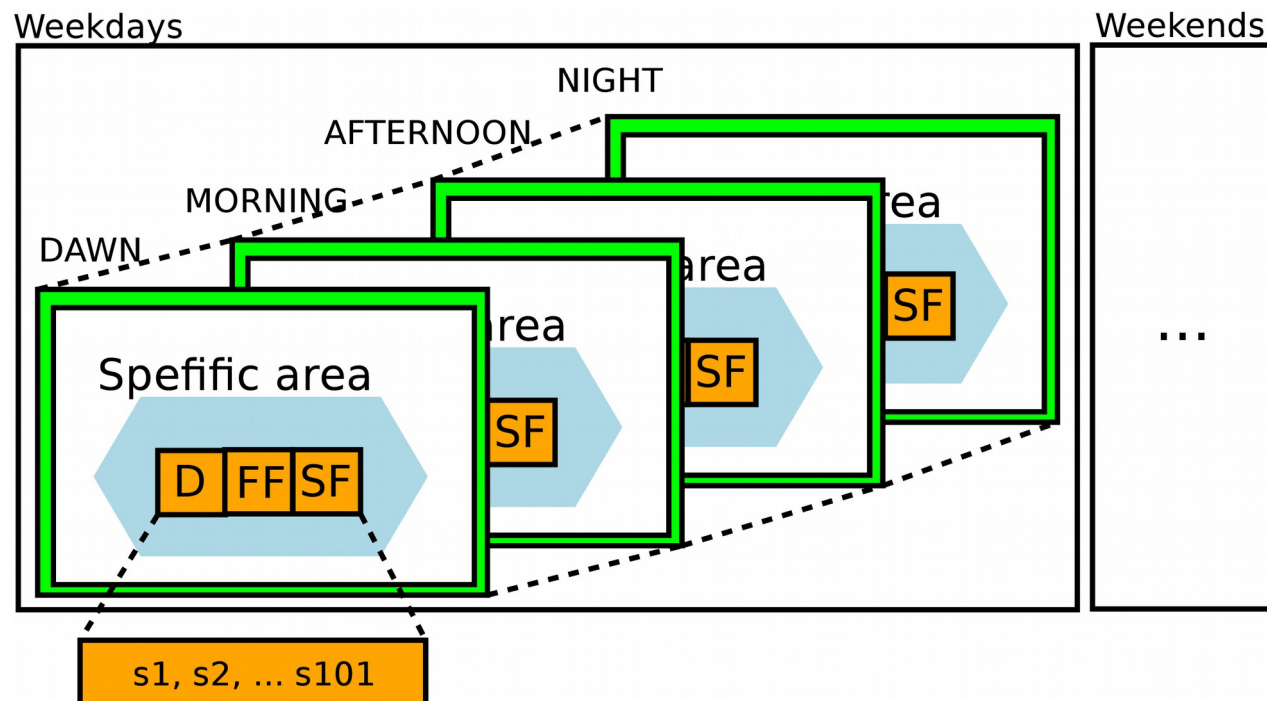
Temporal evaluation

Weekdays



Brazilian cities

American cities

Most of the cities follow the general pattern of the country

Thiago H. Silva – DCC/UFMG

**Considered features: spatial / temporal**

- Each area *a* has a normalized preference vector in **4 disjoint periods of the day** and on **weekdays** and **weekends**



Weekdays — Weekends

NIGHT
AFTERNOON
MORNING
DAWN

Spefific area

D | FF | SF

SF

SF

SF

...

s1, s2, ... s101

General preference vector

**D**=drink / **FF**=fast food / **SF**=slow food

34

# Identifying cultural boundaries

Preference vector for area
(time and space)

Preference vector for area
(time and space)

↓

Principal Component Analysis (PCA)

# Identifying cultural boundaries
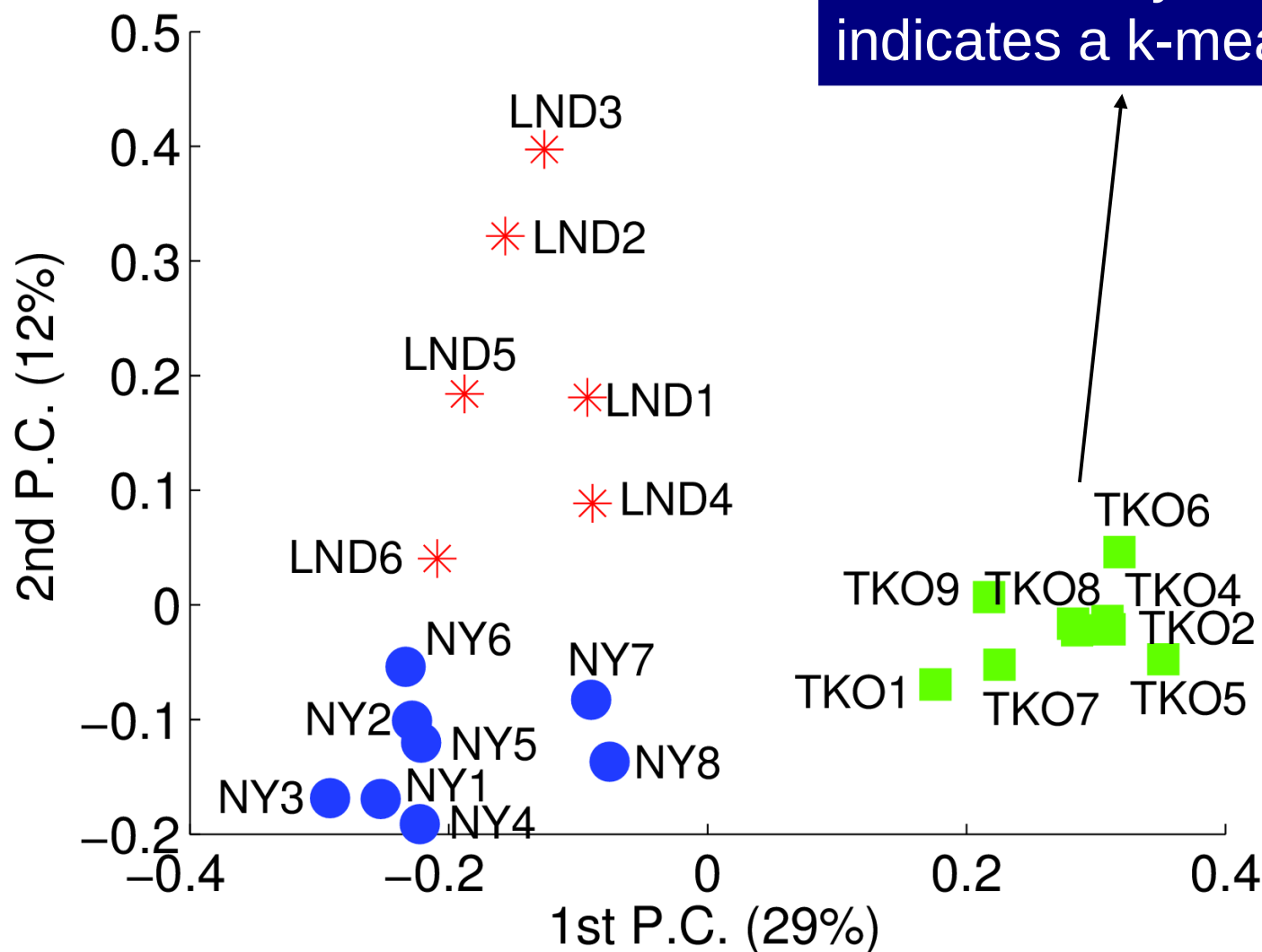
Preference vector for area
(time and space)

Principal Component Analysis (PCA)

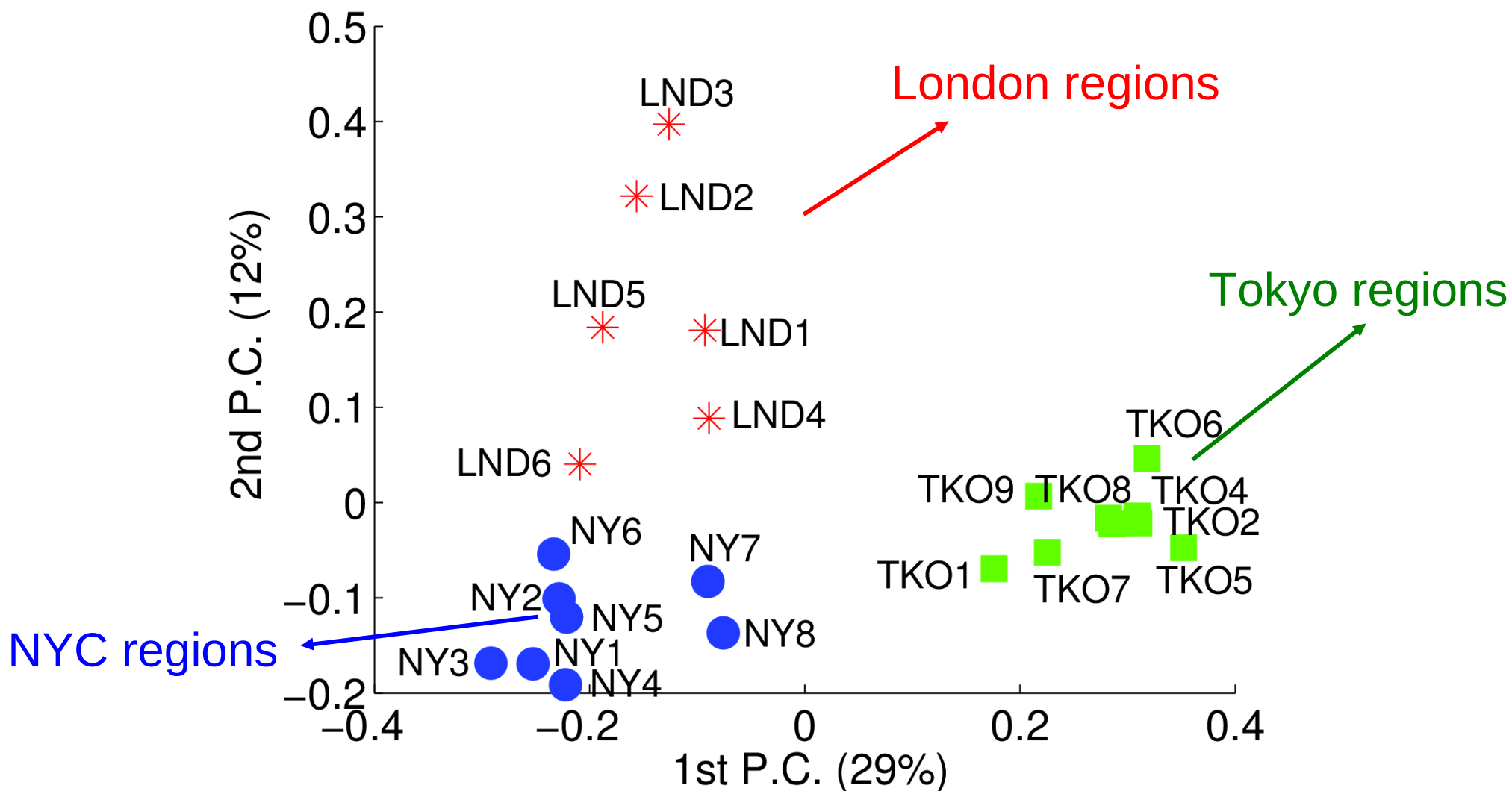k-means to group areas in the space defined by the PCs

# Identifying cultural boundaries

Clustering areas inside cities

Each color/symbol indicates a k-means cluster
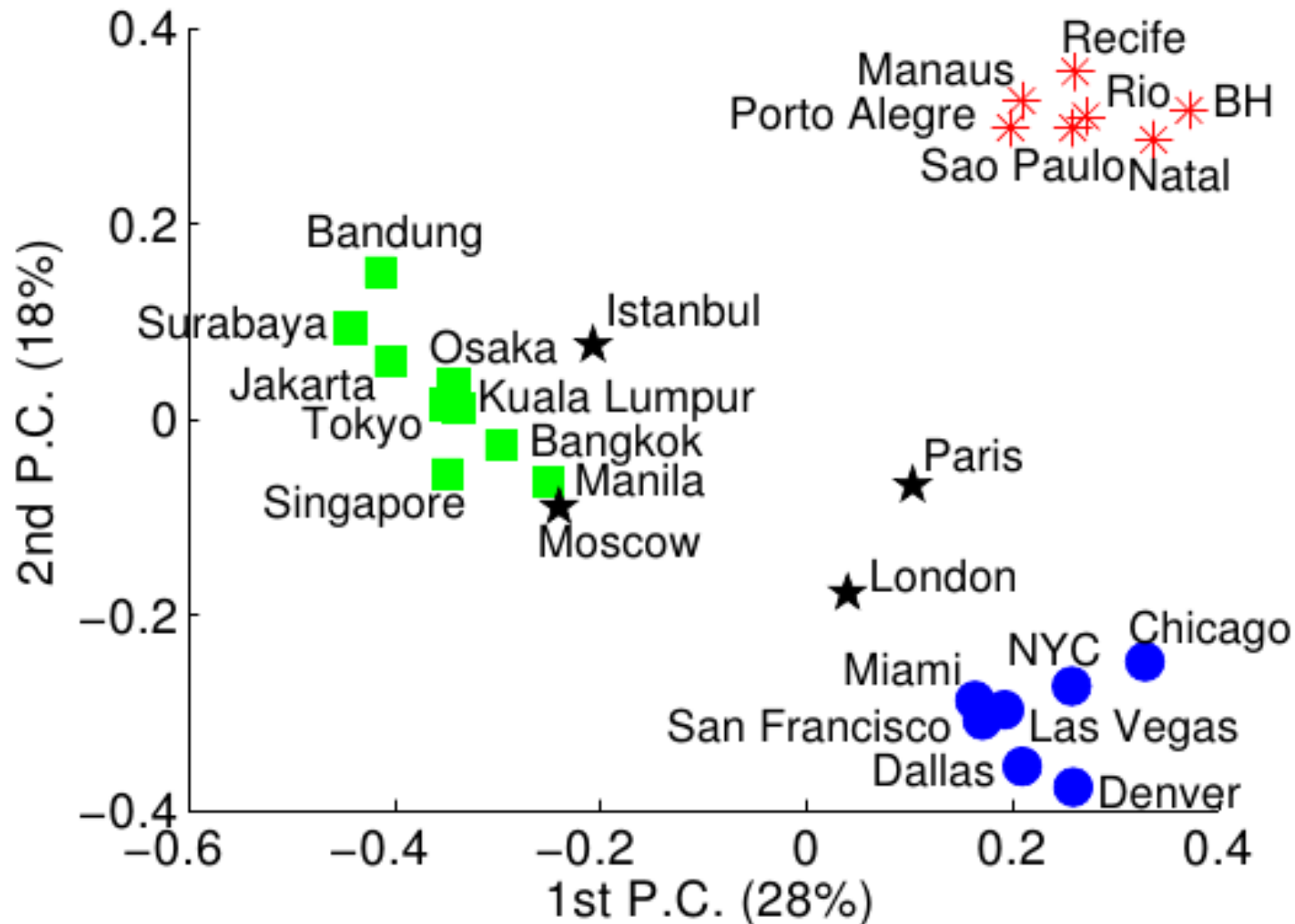
# Identifying cultural boundaries

Clustering areas inside cities



k = 3  (3 different cities)
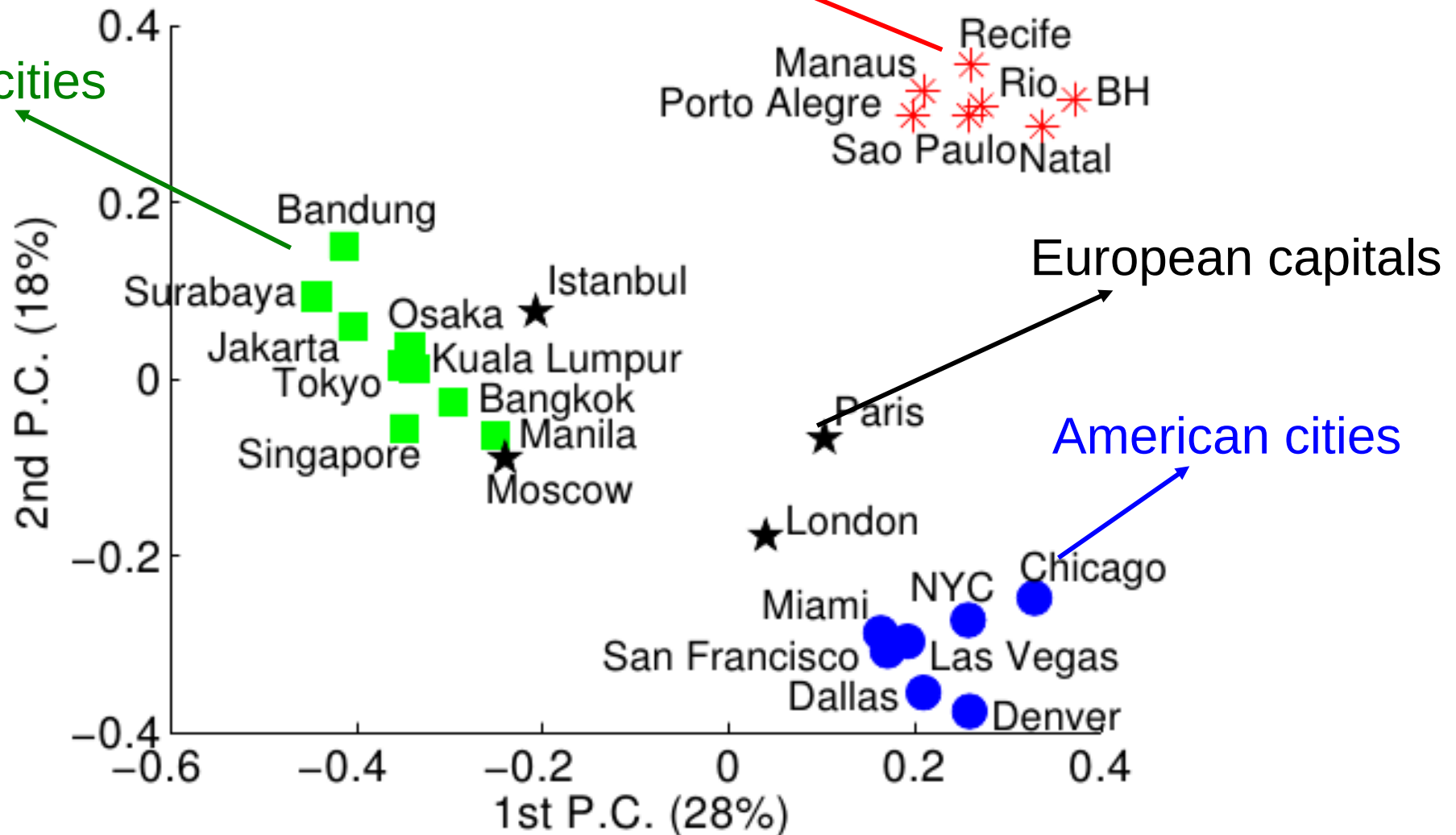
# Identifying cultural boundaries

Clustering cities



k = 4 (4 distinct regions)

Thiago H. Silva – DCC/UFMG

# Identifying cultural boundaries

We can use a **partial set of the features** and **specific time**:
- E.G. Drink at weekend



K = 3 (cities)

42

We can use a **partial set of the features** and **specific time**:
- E.G. Drink at weekend



Note the potential for area recommendation!

E outros traços culturais?

E outros traços culturais?

# Features

➔ Cada cidade c é representada por um vetor de preferências composto por classes de cervejas (features)
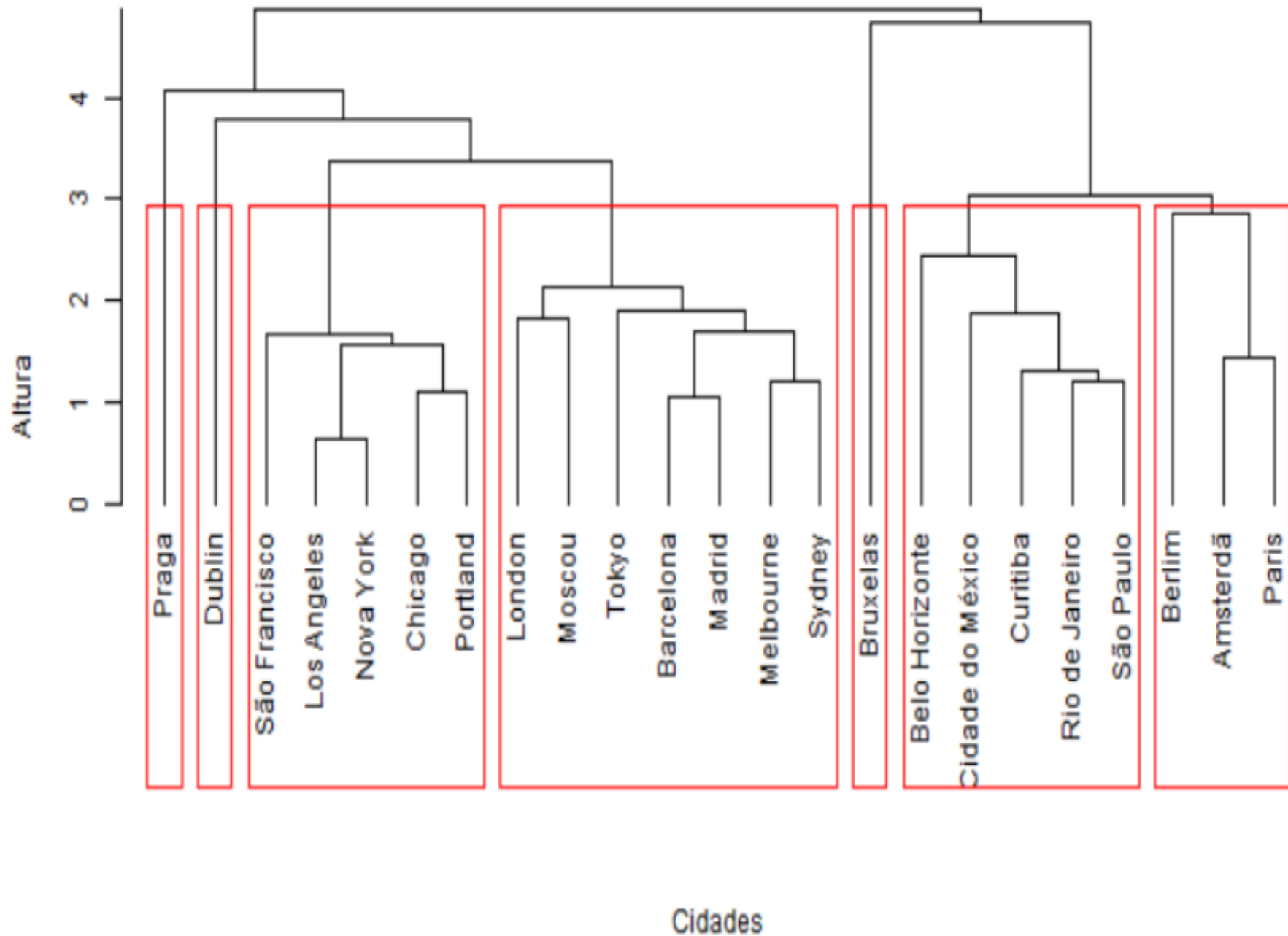


Cidade A

Cidade B

Thiago H. Silva

# Resultados

➔ **Dendograma para o agrupamento realizado com as cidades**

# Engagement of Polarized Groups



Jordan Kobellarz     Alexandre Graeml
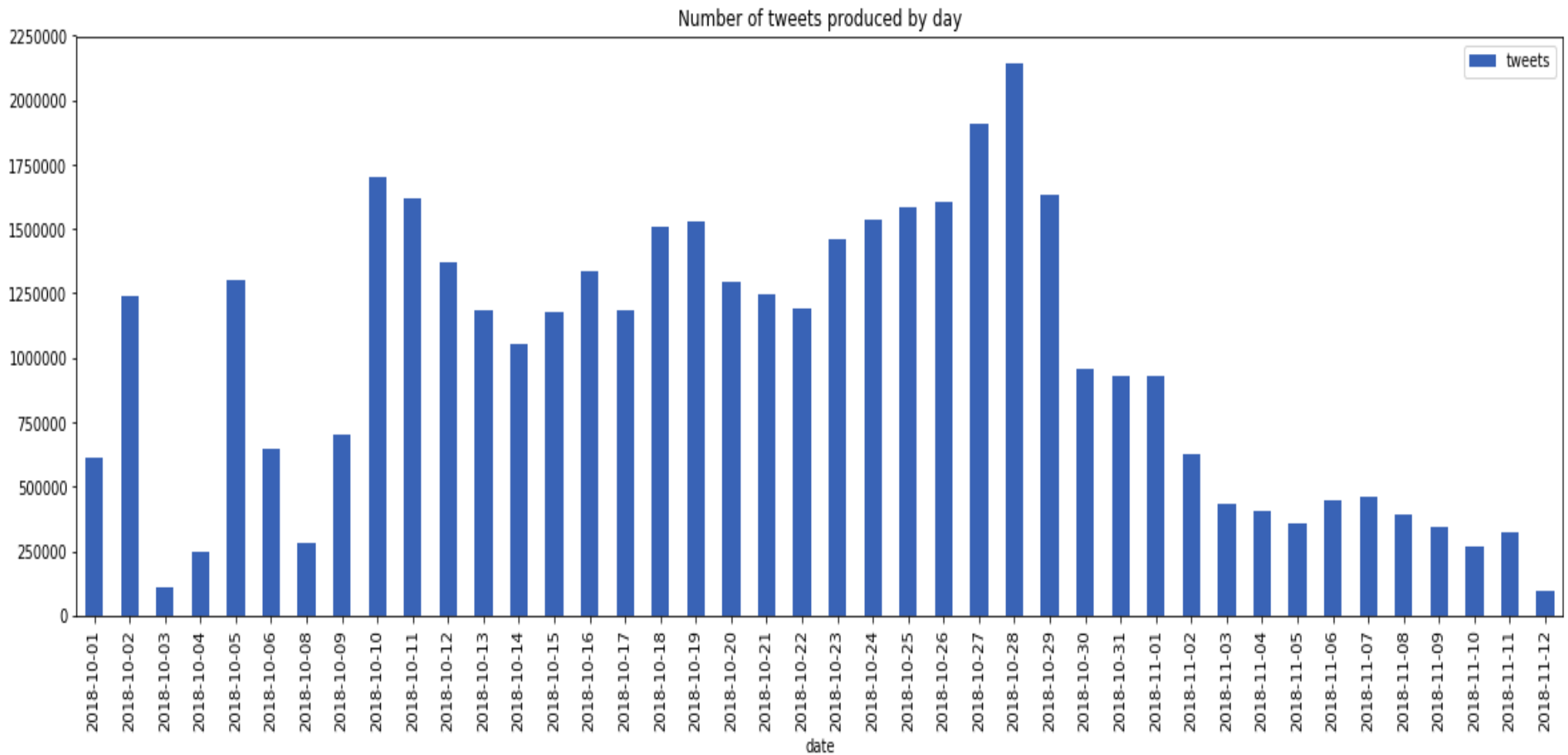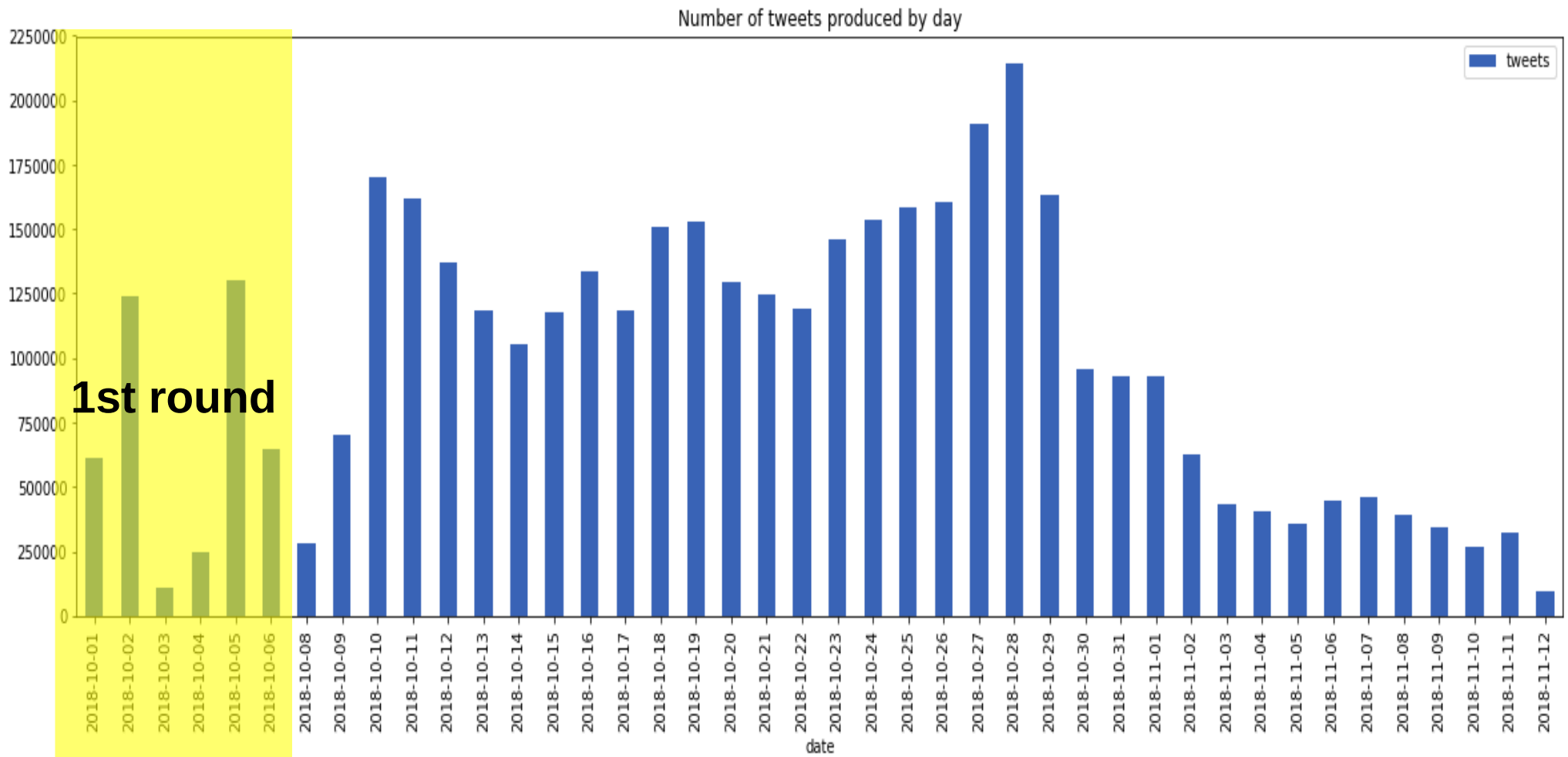
UTFPR

Michelle Reddy

Stanford University

# User engagement on Twitter regarding 2018 Brazil presidential election

Collection of

# ~ 36 Milion tweets

related to 2018 Brazilian Presidential Election
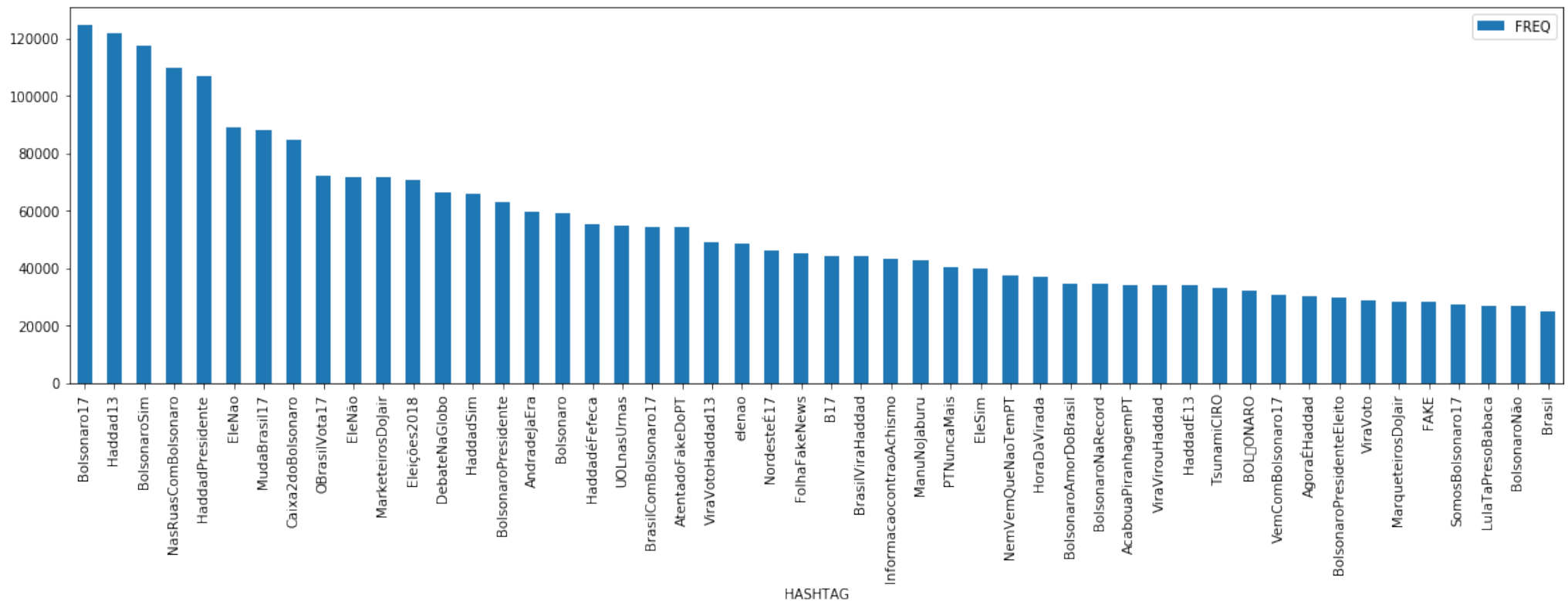
Number of tweets produced by day

# Data



Number of tweets produced by day

1st round

# Data



Number of tweets produced by day

# Data

# 114.512
Unique hashtags

UTFPR

**114.312**
Unique hashtags

**Hard do classify manually!**

# Classification of Hashtags

100 most popular

Help of volunteers

#AndradeJaEra, #Eleições2018, #FolhaFakeNews, #UOLnasUrnas, #Brasil, #InformacaocontraoAchismo, #HoraDaVirada, #VemProDebate, #ViraVoto, #VotoEmCedula, #RodaViva, #AFalhaéCafonérrima, #G1, #SeuVotoMePõeEmRisco, #FAKE, #NãoAceitaremosFraude, #SanatórioGeral, #EAgoraTSE, #Folha, #NaoAceitaremosFraude, #IbopeFake, #viravoto, #democracia, #SuasticaFake, #FakeNews, #Brazil, #SomosTodosReginaDuarte, #delegadofrancischini, #Eleicoes2018 e #BrasilDecide.

#NasRuasComBolsonaro, #BolsonaroSim, #MudaBrasil17, #Bolsonaro17, #Bolsonaro, #BrasilComBolsonaro17, #NordesteÉ17, #BolsonaroPresidente, #NemVemQueNaoTemPT, #PTNuncaMais, #BolsonaroPresidenteEleito, #LulaTaPresoBabaca, #HaddadNãoÉCristão, #B17, #Nordeste17, #OLulaTaPresoBabaca, #EleSim, #bolsonaro17, #PTNão, #PTnão, #FolhaP******DoPT (Conteúdo impróprio), #AcabouaPiranhagemPT, #BolsonaroPresidente17, #elesim, #bolsonaro e #PTnao.
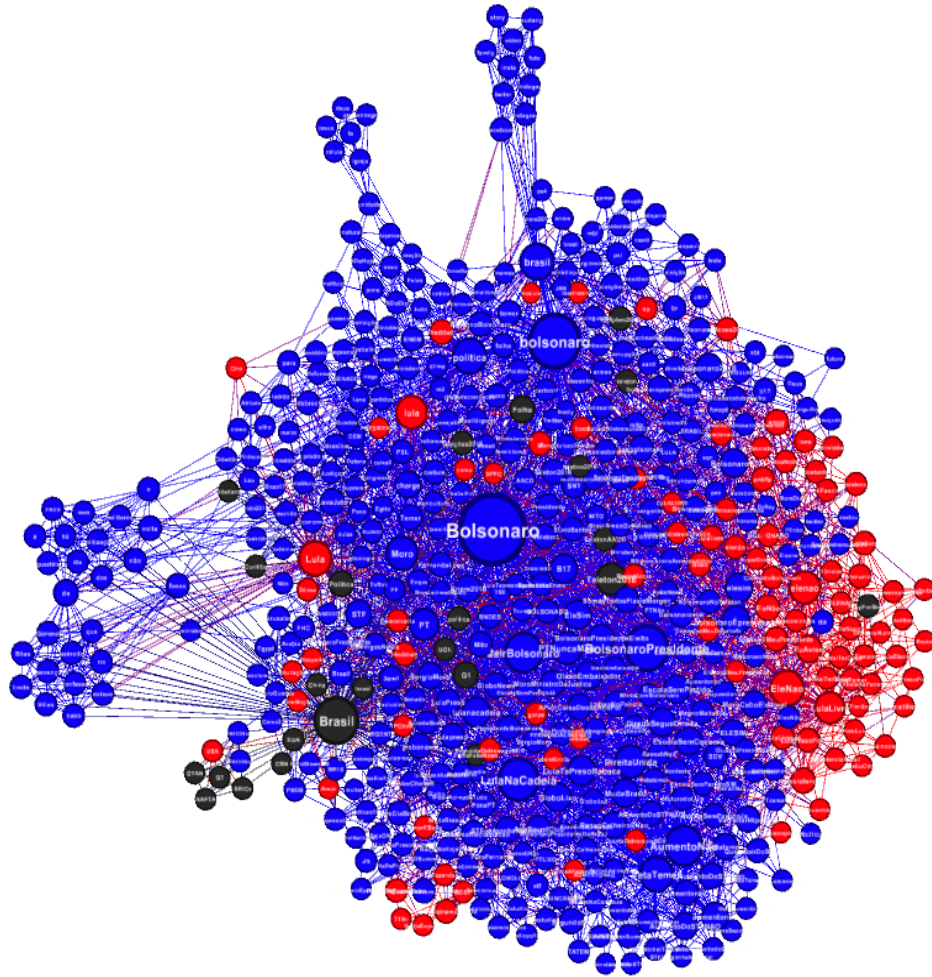
#Haddad13, #Caixa2doBolsonaro, #HaddadPresidente, #HaddadSim, #EleNao, #BrasilViraHaddad, #EleNão, #AgoraÉHaddad, #BolsonaroNão, #ViraVotoHaddad13, #CassaçãoDoBolsonaro, #LulaLivre, #bolsonaroCagao, #Caixa2DoBolsonaro, #Haddad, #elenao, #ViraVirouHaddad13, #ELENAO, #MaisLivrosMenosArmas e #haddadpresidente.

76 agreements

#AndradeJaEra, #Eleições2018, #FolhaFakeNews, #UOLnasUrnas, #Brasil, #InformacaocontraoAchismo, #HoraDaVirada, #VemProDebate, #ViraVoto, #VotoEmCedula, #RodaViva, #AFalhaéCafonérrima, #G1, #Seu... AceitaremosFraude, #SanatórioGeral, #EAgoraTSE, #Folha, #NaoAceitaremosFraude, #IbopeFake, #viravoto, #democracia, #SuasticaFake, #FakeNews, #Brazil, #SomosTodosReginaDuarte, #delegadofrancischini, #Eleicoes2018 e #BrasilDecide.

**Uncertain (?)**

#NasRuasComBolsonaro, #BolsonaroSim, #MudaBrasil17, #Bolsonaro17, #Bolsonaro, #BrasilComBolsonaro17, #NordesteÉ17, #BolsonaroPresidente, #NemVemQueNaoTemPT, #PTNuncaMais, #BolsonaroPresidenteEleito, #LulaTaPresoBabaca, #HaddadNãoÉCristão, #B17, #Nordeste17, #OLulaTaPresoBabaca, #EleSim, #bolsonaro17, #PTNão, #PTnão, #FolhaP***** ... prio), #AcabouaPiranhagemPT, #BolsonaroPresidente17, #elesim, #bolsonaro e #PTnao.

**Right (R)**

#Haddad13, #Caixa2doBolsonaro, #HaddadPresidente, #HaddadSim, #EleNao, #BrasilViraHaddad, #EleNão, #AgoraÉHaddad, #BolsonaroNão, #ViraVotoHaddad13, #CassaçãoDoBolsonaro, #LulaLivre, #bolsonar... Bolsonaro, #Haddad, #elenao, #ViraVirouHaddad13, #ELENAO, #MaisLivrosMenosArmas e #haddadpresidente.

**Left (R)**

Network of co-occurrences of hashtags (all tweets)



Semi-supervised learning using gaussian fields and harmonic functions (Zhu, 2003)
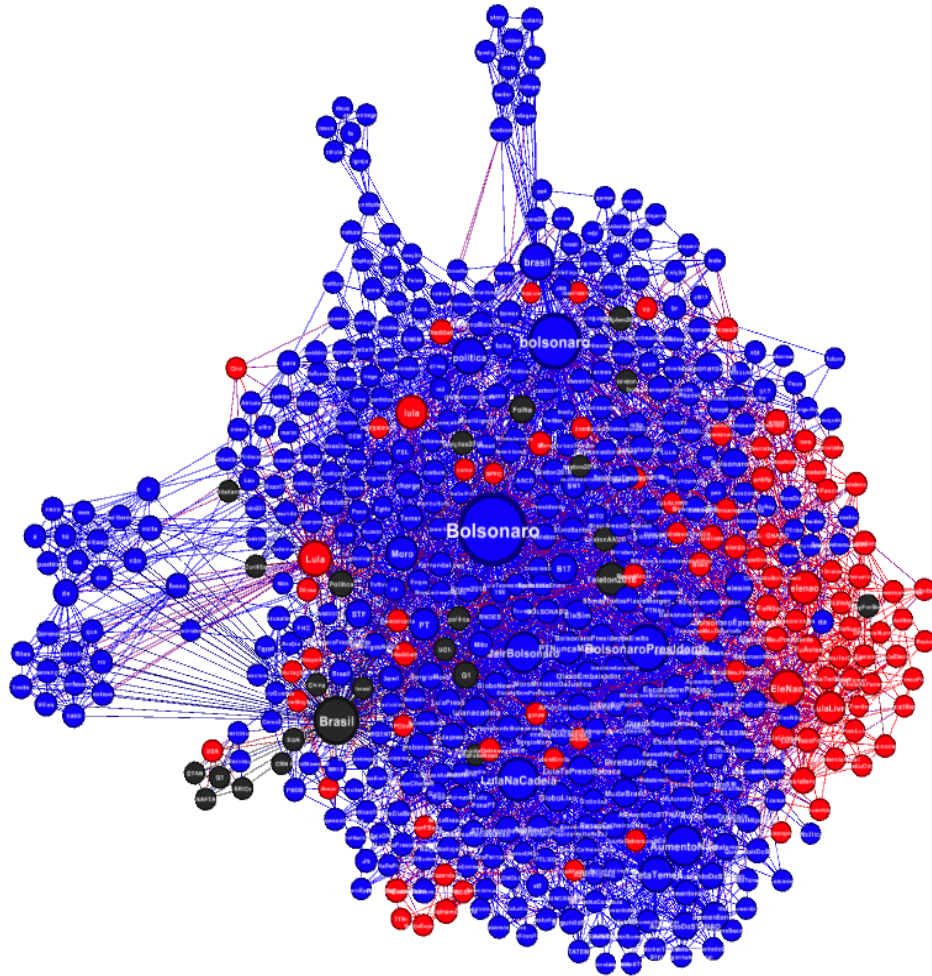
**UTFPR**

Network of co-occurrences of hashtags (all tweets)



Semi-supervised learning using gaussian fields and harmonic functions (Zhu, 2003)

Classification of **78,649** hashtags (68.7% of total)

Only tweets that have one those

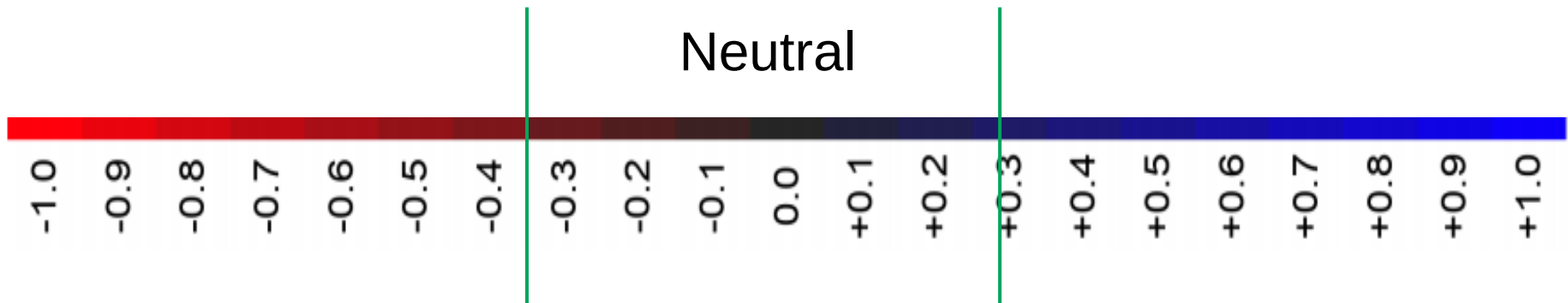$$P(H) = \frac{|H_R| - |H_L|}{|H|},$$
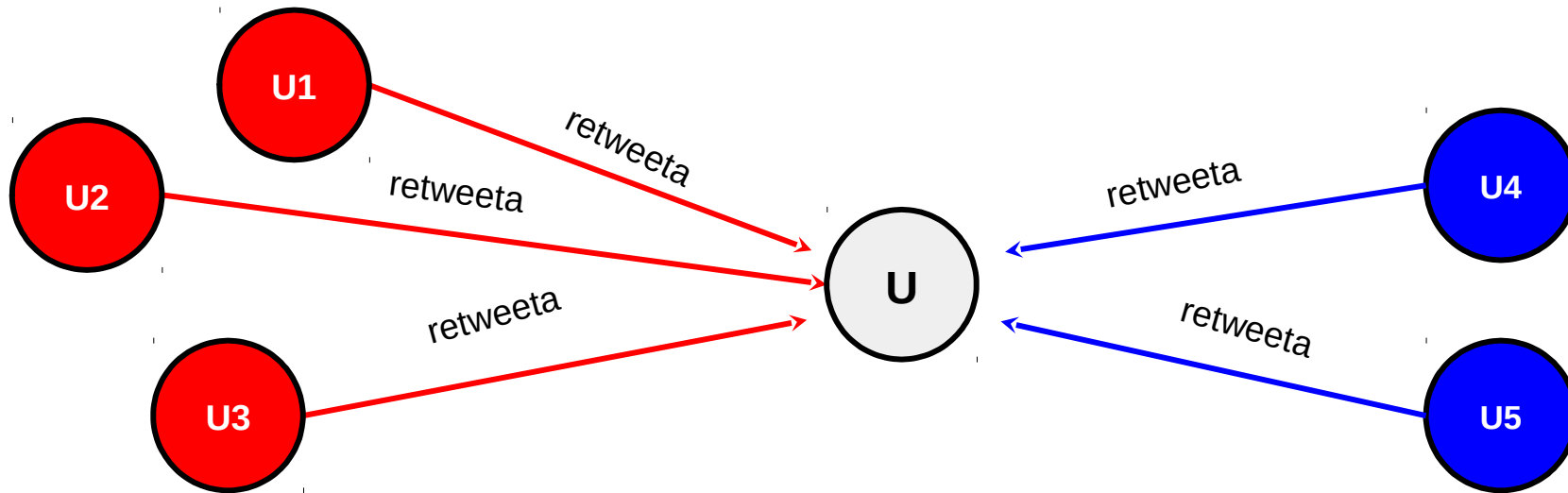
$$H = H_R + H_L + H_?$$

$$P(H) = \frac{|H_R| - |H_L|}{|H|},$$

$$H = H_R + H_L + H_?$$

-1.0  -0.9  -0.8  -0.7  -0.6  -0.5  -0.4  -0.3  -0.2  -0.1  0.0  +0.1  +0.2  +0.3  +0.4  +0.5  +0.6  +0.7  +0.8  +0.9  +1.0

$$P(H) = \frac{|H_R| - |H_L|}{|H|},$$

$$H = H_R + H_L + H_?$$

Neutral

-1.0  -0.9  -0.8  -0.7  -0.6  -0.5  -0.4  -0.3  -0.2  -0.1  0.0  +0.1  +0.2  +0.3  +0.4  +0.5  +0.6  +0.7  +0.8  +0.9  +1.0

# Engagement Graph



Weight represents the amount of retweets

Left

Right

"Echo Chambers"
[Barberá et al. 2015]

# Engagement Graph



Very few connections between bubbles!

Few central nodes regarding this aspect

# Polarity distribution



After the election

Binominal for polarization
[Fiorina e Abrams, 2008]

# Interaction between users



The truth that I want to hear

Bridges made by some
traditional media

G1

Bridge mechanism

**UT**FPR

Willy Muller
(visiting student)

Method to quantify gender preferences in different regions
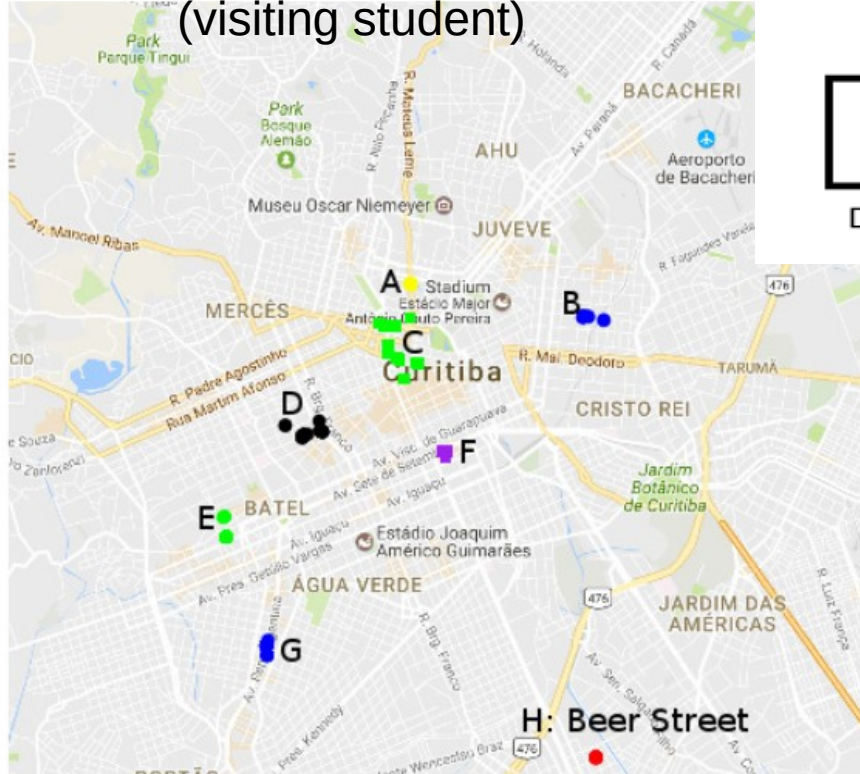
Identification of "anomalous" areas

Ladies section        Men section



**Gender Matters! Analyzing Global Cultural Gender Preferences for Venues Using Social Sensing**
EPJ Data Science, 2017

# Urban Planning and Place Branding

**UTFPR**

Giovani

Ville Santalla
(visiting student)

Luiz Celso

Adaptive and reactive planning (identifying hidden potential)
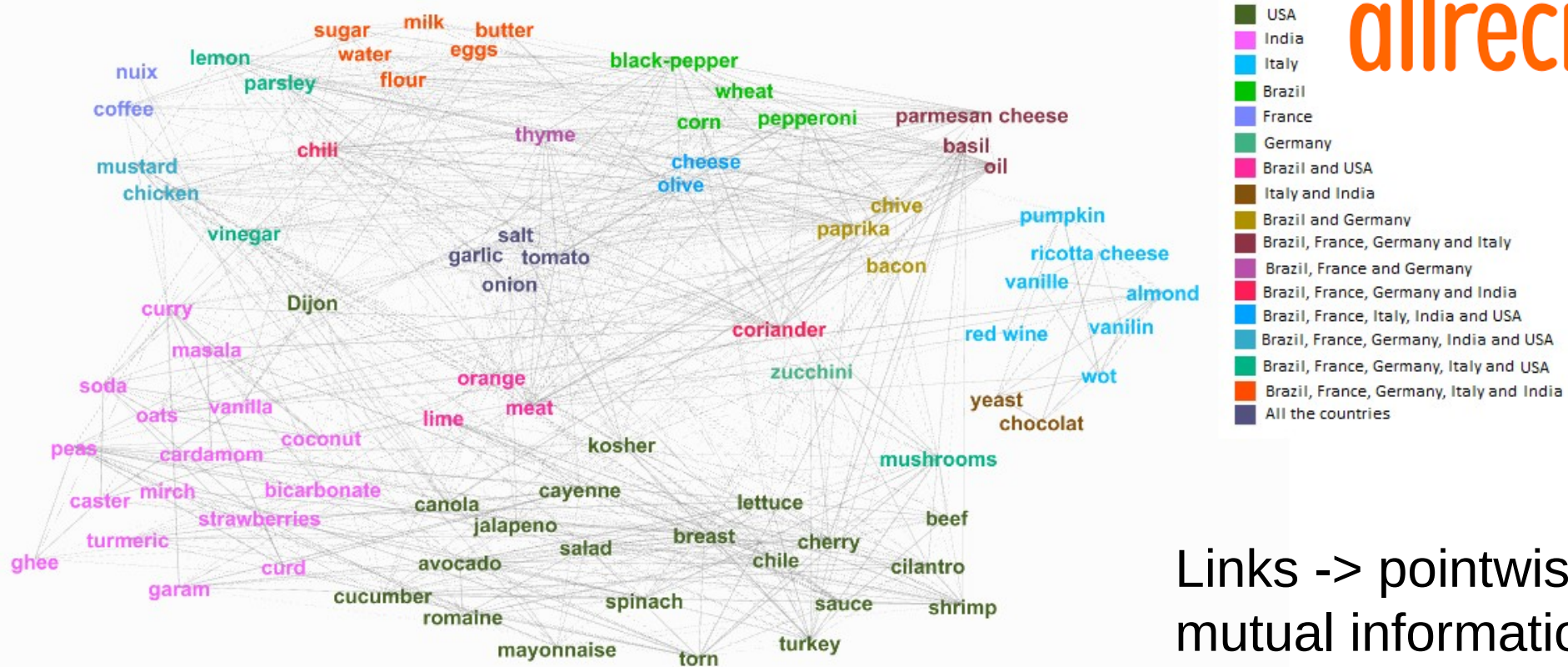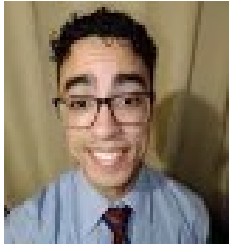Things that are already happening!

# Understanding Success

Better understand the success
recipes around the world
(New recommendation systems)

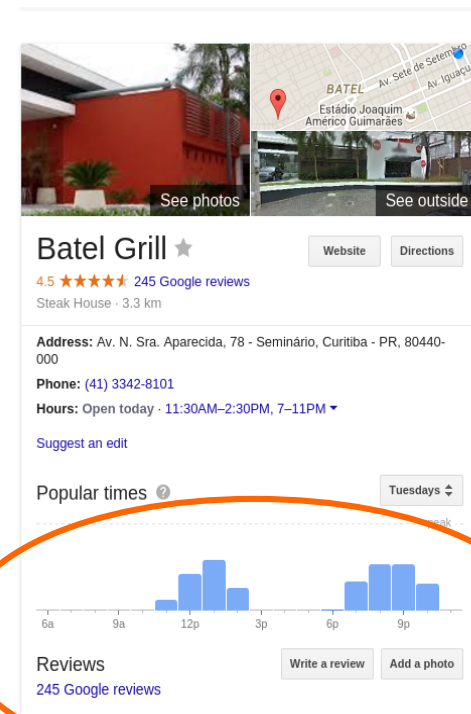Juliana Viscenheski

allrecipes



Links -> pointwise
mutual information

**From Pizza to Curry: Preferences for Recipes Around the World**
Webmedia 2019
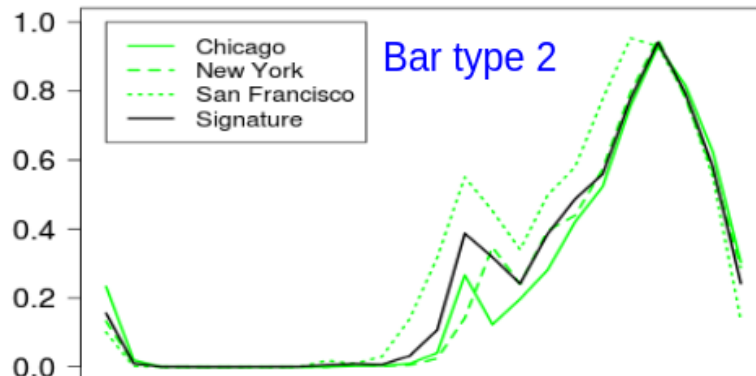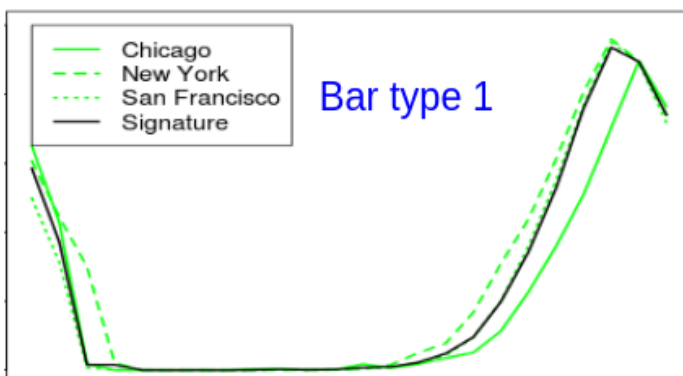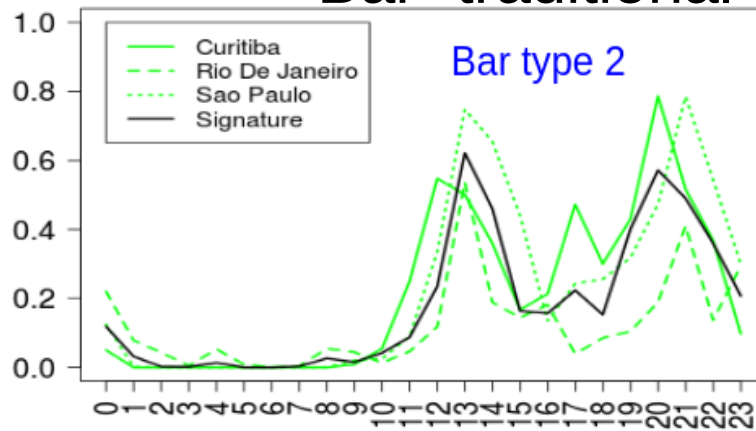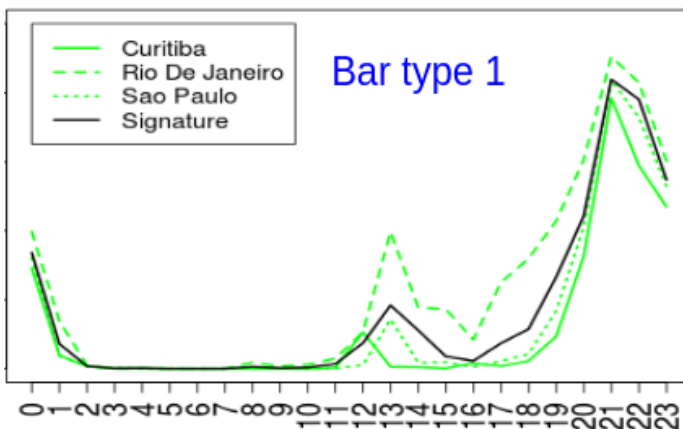
# Business Functioning Dynamics

Leonardo de Assis

Popularity time series of places is an important descriptor ("signature" of a place)

Bar "traditional"

Bar "nightclub"



**Extraction and Exploration of Business Categories Signature**
VLDB Workshops, 2018