IPSW 2023 Radio-Canada Final Presentation Determining the right moment for suggesting the creation of an account

Mario Canche (CIMAT), Andrea Ek (CIMAT), Michael R Lindstrom (UTRGV), Corentin Lonjarret (Radio-Canada), Patrick Mesana (HEC Montreal), Carlos Montes (UTRGV), Nicolas Schönau (Radio-Canada), Omar Sharif (UTRGV), Marziyeh Talebian (Concordia), Louis Willems (Radio-Canada)

August 25, 2023











### What is the problem?

- CBC/Radio-Canada supports Canadian culture and democratic life with a diversity of content
- A strategic aim of the website is to encourage visitors to create accounts to enable personalized experiences
- Users interact with the platform in manifold ways and very few users create accounts
- Radio-Canada wants to determine the optimal moment to suggest account creation to users while also mitigating the annoyance of such a prompt

### Where is the math?

- We seek to identify sequences of events that are closely tied with account creation
- We seek to estimate the probability a user on their journey will create an account
- We seek to identify similarities (clustered features) between the users who create accounts; likewise for those who do not

#### DATA STRUCTURE



### How data are handled?

- data are a collection of webpage interactions with timestamps from Apr-May 2023, including the type of activity, and unique ID per visitor
- Some filtering was needed, e.g., filtering out events after account creation or users with presence before sample
- some users may delete their cookies and appear under different IDs, bots may be in the dataset, etc.

Radio-Canada Working Group IPSW

IPSW 2023 Radio-Canada Final Presentation

### Summary of our methods

- Visualizations and Clustering: clustering visitor activity to understand user groups and identifying features that distinguish account-creators
- Interaction Driven Creation Model: modelling the probability a visitor will create an account given their current webpage interaction status
- Sevents Embedding Analysis: representing events as vectors to find correlated events with account creation events

#### Temporal variations in account creation



### Most recent visit types by different user groups



### Final visit types recorded before account creation



### Kmeans analysis



### Assumptions and derivation

Let there be U users with data

- Given the current user interaction x ∈ ℝ<sup>k</sup>, let Y<sub>x</sub> ∈ {0,1} be a r.v. for a subsequent account creation (0=no, 1=yes)
- Model  $\Pr(Y_x = 0|x) = \sigma(\theta^T \hat{x})$  where  $\theta \in \mathbb{R}^{k+1}$  is a parameter,  $\hat{x} = (1, x^{(1)}, ..., x^{(k)})^T$ , and  $\sigma : z \mapsto \frac{1}{1 + \exp(-z)}$  is the sigmoid function
- To user *i*, let N<sub>i</sub> be either the number of interactions leading up to but not including account creation (if they do) or their number of interactions otherwise
- So Let  $x_{ij} \in \mathbb{R}^k$  be the *j*th interaction of user *i*
- The account dataset is  $\mathcal{D}|_x = \{\{y_{ij}\}_{j=1}^{N_i}\}_{i=1}^U$ , a collection of creation statuses with  $y_{ij}$  the realization of  $Y_{x_{ij}}$

### Assumptions and derivation

- Assume (approximate) decoupling to generate separate experiments effectively a Logistic Regression
- 2 The log-likelihood is then

$$egin{split} \mathcal{L}(\mathcal{D}|_xig| heta) &= \sum_{i=1}^U \mathbbm{1}_{y_{i,N_i}=1} \left(\sum_{j=1}^{N_i-1}\log\sigma( heta^ op\hat{x}_{ij}) + \log(1-\sigma( heta^ op\hat{x}_{i,N_i}))
ight) \ &+ \sum_{i=1}^U \mathbbm{1}_{y_{i,N_i}=0} \left(\sum_{j=1}^{N_i}\log\sigma( heta^ op\hat{x}_{ij})
ight) \end{split}$$

where  $\mathbbm{1}$  is the indicator function

This can be maximized with gradient ascent

# Interaction Driven Creation Model



#### Results

• Looking at  $\theta$  gives insight into user account creation patterns and probabilities of account creation

heta component	value
intercept	-0.022
duration	0.003
day_of_visit_4	0.030
day_of_visit_1	-0.048
cookies₋N	-0.044
detailed_event_Écoute_OHdio	-0.039

Table: Positive coefficients make a user less likely to create an account, negative coefficients make them more likely to create an account.

## What is the immediate context of an account creation event?



# **Events Embedding Analysis**



16 / 22

## Skip-Gram with Negative Sampling (Word2Vec)

#### Estimate the prob of being in the same context

$$p(w_O|w_I) = \frac{\exp\left(v'_{w_O}^{\top} v_{w_I}\right)}{\sum_{w=1}^{W} \exp\left(v'_w^{\top} v_{w_I}\right)}$$

#### Negative Sampling trick - predicting if pairs are in the same context



## "Vocabulary" has 553 different events

post_prop5	visit_type	detailed_event	post_page_event	
abitibitemiscamingue	Info	Page_Info	0	0
acadie	Info	Page_Info	0	
alberta	Info	Page_Info	0	2
alimentation	Info	Page_Info	0	3
ar	Info	Page_Info	0	4
no_post_prop5_info	OHdio	Action_creation_compte	100	548
no_post_prop5_info	no_visit_type_info	Action_creation_compte	100	549
no_post_prop5_info	Info	Action_autre	101	550
no_post_prop5_info	Info	Action_autre	102	551
recettes	Mordu	Action_autre	102	552

### DEMO

### Main Findings

- Sunday, Monday and Tuesday are the days where more accounts have been created
- Visitors who experienced multiple events show a higher rate of account creation
- Probability of account creation can be modelled, suggesting: Monday is a better day for account creation, Écoute\_OHdio is a good event for account creation, and that users who have longer events may be less likely to create an account
- Correlation between creation action from a section with events from the same section

### More to explore...

- Use LSTM recurrent networks to model the sequential problem with the numerical variables.
- Interaction Driven Creation Model could perhaps be improved in performance with a better model/handling of features
- Add more features, e.g., navigation levels, to increase number of events
- Incode sequences as context to have sequence embedding
- Strengthening and identifying the common interpretations between the complementary models

### We would like to thank...

- Radio Canada and their representatives for contributing this problem and their availability throughout the workshop in resolving our confusions
- Odile Marcotte and Nancy Laramée for their organization
- ORM and IVADO for hosting the event