

# *What next? Modeling human behavior using smartphone usage data and (deep) recommender systems*

## Master Thesis Presentation

Simon Wiegrebbe

*Supervisor: Dr. David Rügamer*



*What next?  
Modeling human behavior using  
smartphone usage data and (deep)  
recommender systems*

Simon Wiegrebbe

1 Motivation and Data

2 Modeling

3 Evaluation and Tuning

4 App-level Results

5 Sequence-level Results

6 Discussion

7 References

October 01, 2021

## 1 Motivation and Data

## 2 Modeling

## 3 Evaluation and Tuning

## 4 App-level Results

## 5 Sequence-level Results

## 6 Discussion

## 7 References

Simon Wiegrefe

1 Motivation and Data

2 Modeling

3 Evaluation and Tuning

4 App-level Results

5 Sequence-level Results

6 Discussion

7 References

# 1 Motivation and Data

1 Motivation and  
Data

2 Modeling

3 Evaluation and  
Tuning

4 App-level  
Results

5 Sequence-level  
Results

6 Discussion

7 References

# Motivation

- ▶ Smartphone usage data becoming a valuable data source
- ▶ Most behavioral research: association between smartphone usage patterns and pre-established personality traits

Research Idea:

- ▶ Model behavioral sequences by means of next-event prediction
- ▶ Smartphone usage data from the PhoneStudy project (Stachl et al. 2019)
- ▶ Use natural sequential order (app sessions) in the data:
  - ▶ App sessions start by switching on the screen (ON) and end by switching it off (OFF)
  - ▶ e.g.: *ON, WhatsApp, Calendar, Safari, OFF*
- ▶ Many possible events + sequential data → Use sequence-aware recommender system (RS) algorithms

1 Motivation and Data

2 Modeling

3 Evaluation and Tuning

4 App-level Results

5 Sequence-level Results

6 Discussion

7 References

- ▶ 310 users, study period from October 29, 2017 through January 22, 2018
- ▶ Each app usage assigned exact opening date and time

userID	timestamp	sessionID	appID
1	1.511423e+09	1	1392
1	1.511423e+09	1	1389
1	1.511424e+09	2	1392
1	1.511424e+09	2	1389
1	1.511424e+09	3	1392
...	...	...	...
310	1.515952e+09	844295	1389
310	1.515953e+09	844296	1392
310	1.515953e+09	844296	1389

Table 1: Excerpt of anonymized app-level data. The timestamp column contains Unix timestamps, i.e., seconds passed since January 01, 1970 (UTC).

1 Motivation and  
Data

2 Modeling

3 Evaluation and  
Tuning

4 App-level  
Results

5 Sequence-level  
Results

6 Discussion

7 References

# App-level Representation

- ▶ In language modeling:
  - ▶ Tokens  $\hat{=}$  words
  - ▶ Sentence  $\hat{=}$  sequence of tokens ending with a period
- ▶ Here:
  - ▶ Tokens  $\hat{=}$  apps
  - ▶ Sentences  $\hat{=}$  sessions
- ▶ Objective: next-app prediction
  - ▶ Predicting the next app a user is going to use in a given session
- ▶ Mostly very short sessions

Simon Wiegrefe

1 Motivation and Data

2 Modeling

3 Evaluation and Tuning

4 App-level Results

5 Sequence-level Results

6 Discussion

7 References

# Sequence-level Representation

- ▶ How to address the issue of short session length?
- ▶ Focus on behavior, not individual apps
  - ▶ App-level sessions := concatenations of app-level *categories*
  - ▶ These categories were pre-established by Stachl et al. (2020)
    - ▶ E.g.: "WhatsApp" → "Messaging"
- ▶ Now:
  - ▶ Tokens  $\hat{=}$  categorized app-level sessions
  - ▶ Sentences  $\hat{=}$  daily concatenations of a user's sessions
- ▶ For the sake of unambiguousness: use the terms "event" and "sequence"

Simon Wiegrefe

1 Motivation and Data

2 Modeling

3 Evaluation and Tuning

4 App-level Results

5 Sequence-level Results

6 Discussion

7 References

# Summary Statistics

What next?  
Modeling human behavior using smartphone usage data and (deep) recommender systems

Simon Wiegrefe

1 Motivation and Data

2 Modeling

3 Evaluation and Tuning

4 App-level Results

5 Sequence-level Results

6 Discussion

7 References

	App-level	Sequence-level
Number of events	4,314,830	720,379
Number of unique events	2,488	2,454
Number of sequences	844,296	9,377
Number of users	310	310
Sequences per user	2,723.54	30.25
Mean number of events per sequence	5.11	76.82
1st quartile of number of events per sequence	2.0	34.0
Median number of events per sequence	4.0	66.0
3rd quartile of number of events per sequence	6.0	106.0

Table 2: Summary statistics of app-level and sequence-level data.

- ▶ Drawback of sequence-level analysis: data size
  - ▶ events in sequence-level  $\approx$  sequences in app-level

## 2 Modeling

1 Motivation and Data

2 Modeling

3 Evaluation and Tuning

4 App-level Results

5 Sequence-level Results

6 Discussion

7 References

# Definitions and Terminology

- ▶ Baseline model := non-NN-based model
- ▶ Session-based model: no incorporation of user-level information (user ID)
- ▶ Session-aware model: incorporation of user-level information
- ▶  $s_{complete} := (s_1, s_2, \dots, s_m)$ : sequence of chronologically ordered events
  - ▶  $s_s$ : last “known” event in the sequence
  - ▶  $s_{s+1}$ : event we seek to predict
  - ▶  $s := (s_1, s_2, \dots, s_s)$ : current (known) sequence
- ▶  $i$ : candidate event for  $s_{s+1}$

# Overview of Algorithms

Session-based baseline models:

- ▶ Co-occurrence frequency-based models:
  - ▶ *AR, SR*
- ▶ Neighborhood-based models:
  - ▶ *SKNN, STAN, VSTAN*
- ▶ Tree-based models:
  - ▶ *CT*

Session-based neural models:

- ▶ *GRU4Rec*

Session-aware neural models:

- ▶ *HGRU4Rec*

Extension-equipped models:

- ▶ *SR\_BR, VSTAN\_EBR, GRU4Rec\_R*

Simon Wiegrefe

1 Motivation and Data

2 Modeling

3 Evaluation and Tuning

4 App-level Results

5 Sequence-level Results

6 Discussion

7 References

# Session-based Baseline Models (I)

- ▶ *AR* and *SR* (Ludewig and Jannach 2018)
  - ▶ are based on co-occurrence frequencies
  - ▶ only take into account  $s_s$  when making a prediction
- ▶ *AR*
  - ▶ simply counts co-occurrences of  $s_s$  with *any*  $i$
  - ▶ normalizes this count by the number of all co-occurrences
- ▶ *SR*
  - ▶ accounts for sequential event order
  - ▶ only counts co-occurrences where  $s_s$  precedes *any*  $i$
  - ▶ decreases the weight if other events occurred in-between

Simon Wiegrefe

1 Motivation and Data

2 Modeling

3 Evaluation and Tuning

4 App-level Results

5 Sequence-level Results

6 Discussion

7 References

# Session-based Baseline Models (II)

- ▶ The neighborhood-based *SKNN* (Jannach and Ludewig 2017)
  - ▶ defines a neighborhood of most similar past sequences
  - ▶ determines similarity between  $s$  and neighbor sequences
  - ▶ computes score as sum of similarity scores across all sequences which contain  $i$
- ▶ *STAN* (Garg et al. 2019) and *VSTAN* (Ludewig et al. 2021) extend *SKNN*, for instance by
  - ▶ accounting for event recency in  $s$  using decay functions
  - ▶ accounting for sequence recency of neighbor sequences

Simon Wiegrefe

1 Motivation and Data

2 Modeling

3 Evaluation and Tuning

4 App-level Results

5 Sequence-level Results

6 Discussion

7 References

# Session-based Baseline Models (III)

- ▶ *CT* by Mi and Faltings (2018)
  - ▶ uses a Markov model with context-dependent depth of the Markov chain
  - ▶ builds a tree where
    - ▶ each node is a set of sequences with a specific suffix (context)
    - ▶ each edge combines a parent and child node
  - ▶ defines experts as conditional probability of  $i$ , given a context
  - ▶ generates predictions through combination of all experts activated by  $s$

Simon Wiegrefe

1 Motivation and Data

2 Modeling

3 Evaluation and Tuning

4 App-level Results

5 Sequence-level Results

6 Discussion

7 References

# Session-based Neural Models

What next?  
Modeling human  
behavior using  
smartphone usage  
data and (deep)  
recommender  
systems

Simon Wiegrefe

## ► GRU4Rec (Hidasi et al. 2015)

- ▶ initially one-hot encodes single input events
- ▶ feeds input vectors into a Gated Recurrent Unit (GRU) layer
- ▶ uses pairwise ranking losses for training
- ▶ outputs, for each event, the likelihood of being next in the sequence

1 Motivation and Data

2 Modeling

3 Evaluation and Tuning

4 App-level Results

5 Sequence-level Results

6 Discussion

7 References

# Session-aware Neural Models

What next?  
Modeling human behavior using smartphone usage data and (deep) recommender systems

Simon Wiegrefe

- ▶ *HGRU4Rec* (Quadrana et al. 2017)
  - ▶ is a user-aware extension of *GRU4Rec*
  - ▶ contains a short- and long-term memory GRU layer
  - ▶ generates recommendations for each event in a sequence through a session-level GRU (like *GRU4Rec*)
  - ▶ updates an additional user-level GRU at the end of each sequence
  - ▶ employs its hidden state for initialization of the session-level GRU at the beginning of the next sequence

1 Motivation and Data

2 Modeling

3 Evaluation and Tuning

4 App-level Results

5 Sequence-level Results

6 Discussion

7 References

# Extensions

- ▶ We use (a combination of) 3 different heuristics for some session-based algorithms
  - ▶ Extensions provide user-level information from past sequences → session-awareness
- 1) *Extend* prepends events from the user's preceding sequence if  $s$  is short
  - 2) *Boost* increases the score of  $i$  if  $i$  has occurred in the user's past sequences
  - 3) *Remind* adds a reminder score to the original model score

# Implementation

What next?  
Modeling human  
behavior using  
smartphone usage  
data and (deep)  
recommender  
systems

Simon Wiegrefe

- ▶ Implementation of models and extensions based on Latifi, Mauro, and Jannach (2021)<sup>1</sup>
- ▶ We perform all modeling, evaluation, and analysis tasks in Python

1 Motivation and Data

2 Modeling

3 Evaluation and Tuning

4 App-level Results

5 Sequence-level Results

6 Discussion

7 References

---

<sup>1</sup><https://github.com/rn51/session-rec/>

## 3 Evaluation and Tuning

1 Motivation and Data

2 Modeling

3 Evaluation and Tuning

4 App-level Results

5 Sequence-level Results

6 Discussion

7 References

# Train-Validation-Test Split

- ▶ Time-ordered and user-clustered data
- ▶ Standard time-agnostic cross-validation not applicable
- ▶ Last-event split method, applied twice:
  - 1) Clip off each user's last sequence → test set
  - 2) Clip off each user's last sequence from the remaining data → validation set
- ▶ Additionally: split study period into 5 equally long sub-periods (windows)
  - ▶ Apply train-validation-test split to all 5 windows
  - ▶ Average performance results across all 5 test sets

Simon Wiegrefe

1 Motivation and Data

2 Modeling

3 Evaluation and Tuning

4 App-level Results

5 Sequence-level Results

6 Discussion

7 References

# Evaluation Protocol and Metrics (I)

What next?  
Modeling human behavior using smartphone usage data and (deep) recommender systems

Simon Wiegrefe

- ▶ Evaluate predictions for all test sequence events, except the first one
- ▶ Ground truth: (only) the actually observed event
- ▶ How to choose a performance metric?
  - ▶ Quantify the goodness of recommendation lists

1 Motivation and Data

2 Modeling

3 Evaluation and Tuning

4 App-level Results

5 Sequence-level Results

6 Discussion

7 References

## Evaluation Protocol and Metrics (II)

- 1) Hit Rate (HR):  $HR@k$  is simply the fraction of all  $n$  events for which the corresponding recommendation list of length  $k$ ,  $rl(k)_i$ , includes the ground truth,  $y_i$ :

$$HR@k = \frac{\sum_{i=1}^n \mathbb{1}_{rl(k)_i}(y_i)}{n}$$

- 2) Mean Reciprocal Rank (MRR):  $MRR@k$  additionally accounts for the ranking within the recommendation list.  $MRR@k$  computes the reciprocal rank of the ground truth within the recommendation list,  $rr_i$ , then averages this reciprocal rank across all  $n$  events:

$$MRR@k = \frac{\sum_{i=1}^n rr_i}{n}$$

- We consider  $HR@k$  and  $MRR@k$  for  $k \in \{1, 5, 10, 20\}$

- ▶ Simple random search with budget 100 for each algorithm
- ▶ Hyperparameter search spaces as in Latifi, Mauro, and Jannach (2021)
- ▶ Tuning on five-window data, then averaging performance to determine optimal hyperparameter configuration
- ▶ Tuning metric:  $HR@1$

1 Motivation and Data

2 Modeling

3 Evaluation and Tuning

4 App-level Results

5 Sequence-level Results

6 Discussion

7 References

## 4 App-level Results

1 Motivation and Data

2 Modeling

3 Evaluation and Tuning

4 App-level Results

5 Sequence-level Results

6 Discussion

7 References

# Overall Performance

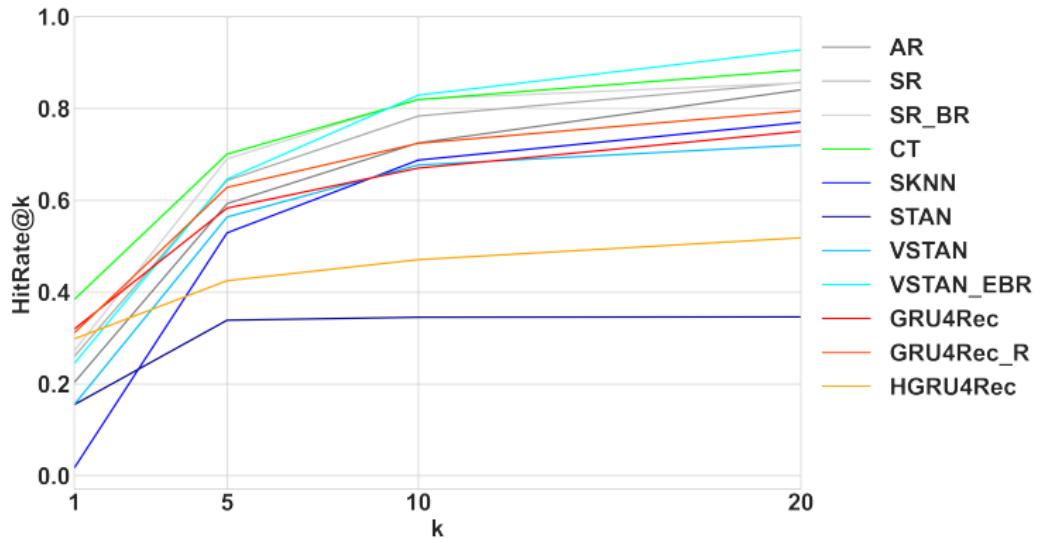


Figure 1:  $HR@k$  performance for  $k = 1, 5, 10$ , and 20 on five-window app-level data.

- ▶ Best performer i.t.o.  $HR@1$  and  $HR@5$ : *CT*
- ▶ Best performer i.t.o.  $HR@10$  and  $HR@20$ : *VSTAN\_EBR*
- ▶ Strong  $HR@1$  performance of NN-based models

# Minimum Sequence Length (I)

- ▶ Background:
  - ▶ *GRU4Rec*, *GRU4Rec\_R*, and *HGRU4Rec* employ RNNs, which learn from the present sequence
  - ▶ App-level sequences are typically short → RNN-based methods do not have “much to learn from”
  - ▶ Co-occurrence frequency-based models simply “look up” co-occurrences with the last known event
- ▶ Hypotheses:
  - ▶ Better performance of NN-based models on longer sequences
  - ▶ No impact of sequence length on performance of *AR*, *SR*, and *SR\_BR*

→ Train and evaluate our models on a subset containing only sequences with at least 20 events

# Minimum Sequence Length (II)

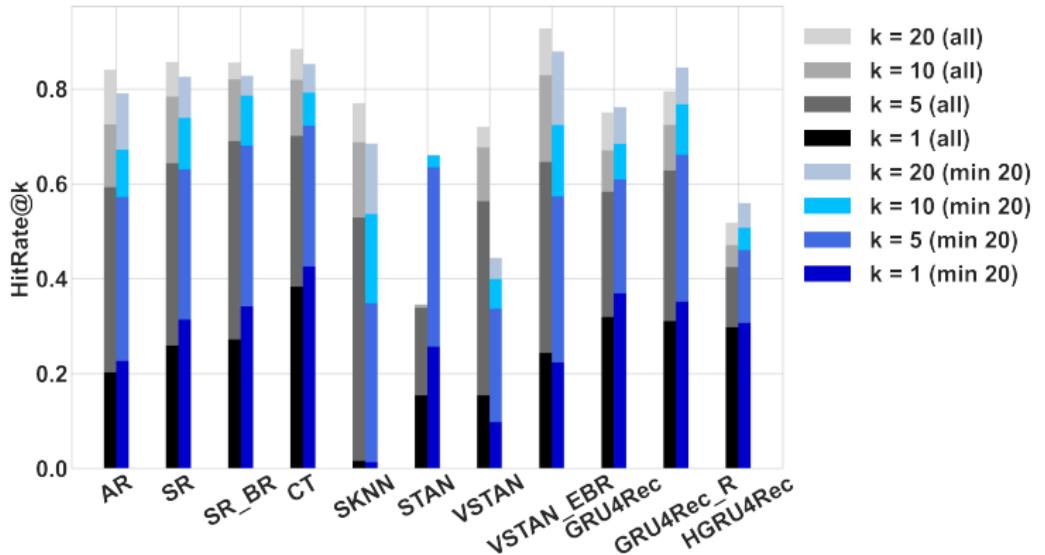


Figure 2:  $HR@k$  comparison between performance on full five-window app-level data (left bars) and performance on five-window app-level data when only training and evaluating on sequences with a minimum length of 20 (right bars), for  $k \in \{1, 5, 10, 20\}$ .

- ▶  $CT$  still best performer for  $HR@1$  and  $HR@5$
- ▶ No large changes for  $AR$ ,  $SR$ , and  $SR\_BR$
- ▶ Performance of NN-based models improves

# Minimum Sequence Length (III)

What next?  
Modeling human behavior using smartphone usage data and (deep) recommender systems

Simon Wiegrefe

- ▶ What if instead we train on all sequences and only evaluate on long sequences?
  - ▶  $CT$  still best performer
  - ▶ All neural models perform considerably worse
  - ▶ Surprising because the full training dataset is considerably larger
- ▶ Conclusion: performance on long sequences benefits from training on long sequences only

1 Motivation and Data

2 Modeling

3 Evaluation and Tuning

4 App-level Results

5 Sequence-level Results

6 Discussion

7 References

# Position in Test Sequence (I)

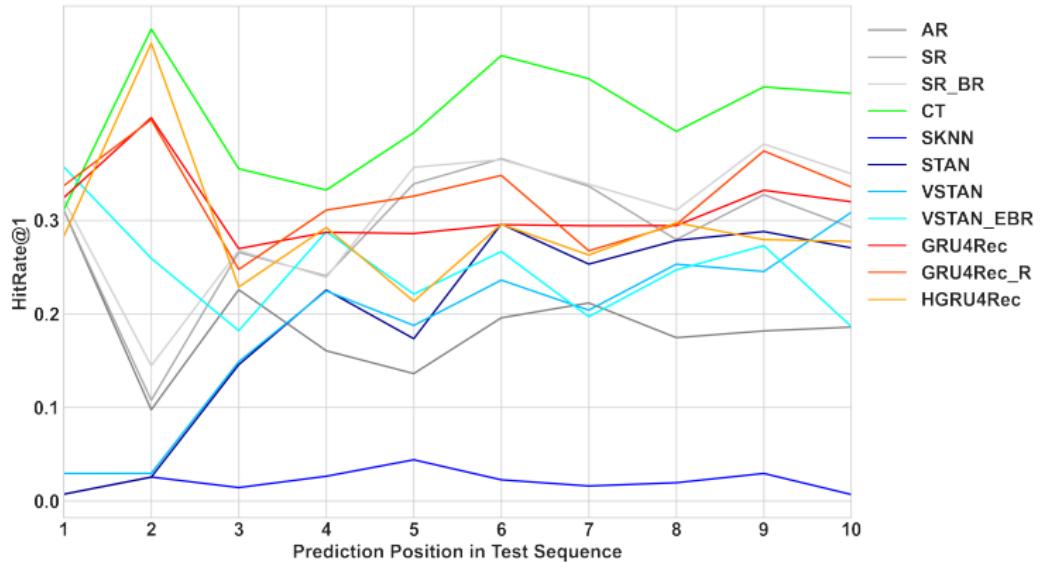


Figure 3: HR@1 performance across the first ten prediction positions on five-window app-level data.

- ▶ Initial performance boost for VSTAN\_EBR
- ▶ No clear trend for all other models

# Position in Test Sequence (II)

Algorithm	position <= 2	position > 2	position <= 5	position > 5	position <= 10	position > 10
AR	0.2269	0.1892	0.2082	0.1934	0.2045	0.1898
SR	0.2307	0.2768	0.2514	0.2798	0.2676	0.2516
SR_BR	0.2520	0.2854	0.2660	0.2904	0.2838	0.2468
CT	0.3878	0.3818	0.3766	0.4012	0.3911	0.3710
SKNN	0.0142	0.0167	0.0193	0.0111	0.0193	0.0061
STAN	0.0145	0.2298	0.0843	0.2602	0.1268	0.2385
VSTAN	0.0295	0.2230	0.0950	0.2469	0.1270	0.2577
VSTAN_EBR	0.3180	0.2058	0.2807	0.1903	0.2709	0.1405
GRU4Rec	0.3581	0.2984	0.3264	0.3098	0.3208	0.3173
GRU4Rec_R	0.3659	0.2827	0.3342	0.2816	0.3304	0.2311
HGRU4Rec	0.3639	0.2593	0.3132	0.2665	0.3073	0.2542

Table 3:  $HR@1$  performance results on five-window app-level data, by positional cutoff within test sequence.

- ▶ Worse performance for NN-based models on later positions
- ▶ If training is not tailored towards them: NN-based models struggle with later positions in the prediction sequences and, consequently, with long prediction sequences

- ▶ Key issue and potential performance bottleneck: short sequence length
- ▶ ON and OFF events are hardly informative
- ▶ ON-OFF sequences make up 38.91% of all app-level sequences
- ▶ Effect of dropping all ON and OFF events from the app-level data?

1 Motivation and Data

2 Modeling

3 Evaluation and Tuning

4 App-level Results

5 Sequence-level Results

6 Discussion

7 References

# Removing ON and OFF Events (II)

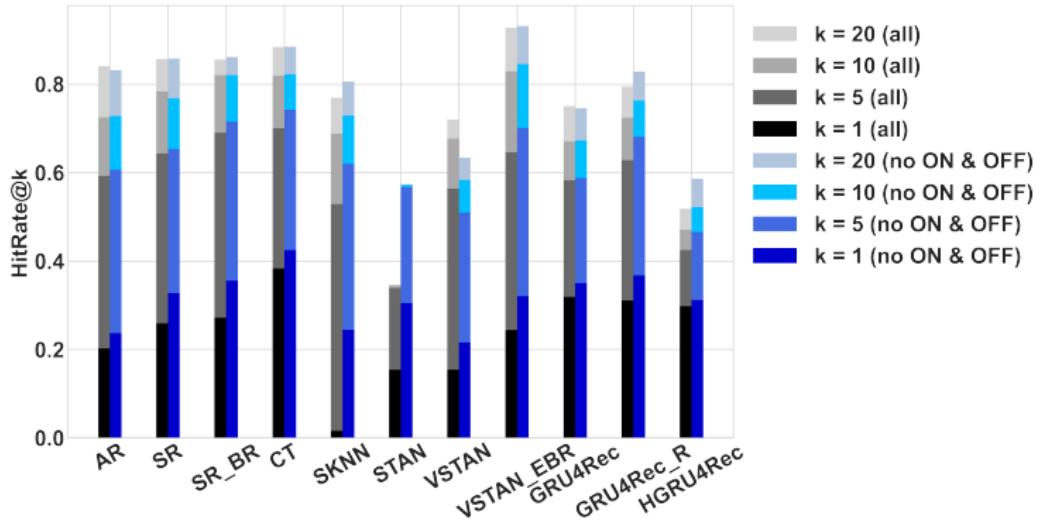


Figure 4:  $HR@k$  performance comparison between full five-window app-level data (left bars) and five-window app-level data after dropping all ON and OFF events (right bars), for  $k \in \{1, 5, 10, 20\}\}.$

- ▶ Improvements i.t.o.  $HR@1$  across the board
- ▶ Substantial improvements for neighborhood-based models
- ▶ Drawback: limited representativeness of results

- ▶ Ultimate goal: predict human behavioral sequences → consider next-category prediction instead of next-app prediction
- ▶ For evaluation, simply consider app category: e.g., “Messaging” instead of “WhatsApp”
- ▶ If performance improves considerably: models learn more about behavioral sequences than previously thought

1 Motivation and Data

2 Modeling

3 Evaluation and Tuning

4 App-level Results

5 Sequence-level Results

6 Discussion

7 References

# Category-level Prediction (II)

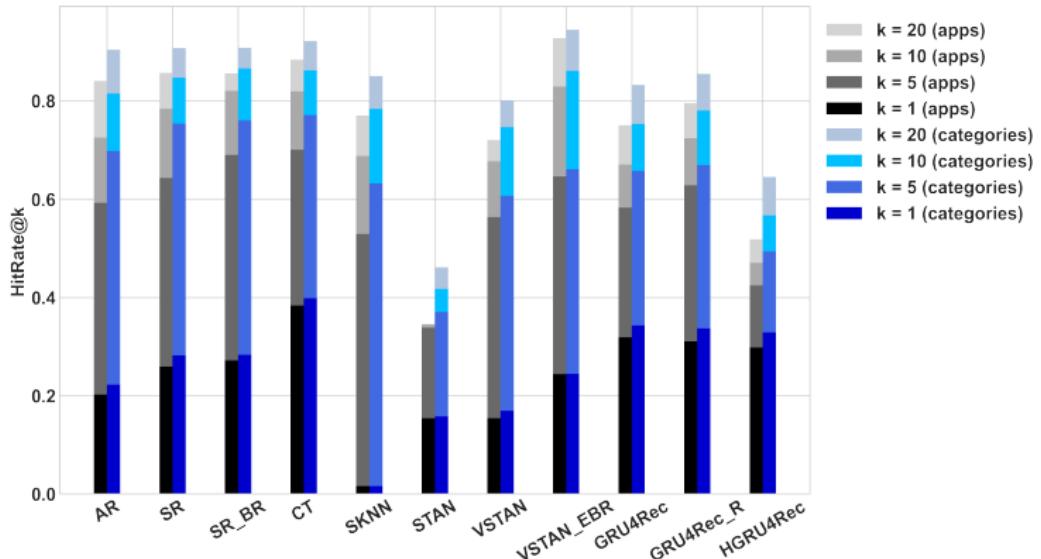


Figure 5:  $HR@k$  performance increases on five-window app-level data when only considering app categories for evaluation (left bars), instead of considering the individual apps as well (right bars), for  $k \in \{1, 5, 10, 20\}$ .

- ▶ Performance increases especially for larger  $k$ , more pronounced for NN-based methods, and proportional to app-level performance

# Embedding Analysis (I)

What next?  
Modeling human behavior using smartphone usage data and (deep) recommender systems

Simon Wiegrefe

- ▶ Can deep learning models learn smartphone app semantics?
- ▶ Do apps from a common app category form clusters in the embedding space? → Add an embedding layer ( $d = 128$ ) to *GRU4Rec*
- ▶ Apply TSNE (Hinton and Roweis 2002) to obtain two-dimensional app embeddings

1 Motivation and Data

2 Modeling

3 Evaluation and Tuning

4 App-level Results

5 Sequence-level Results

6 Discussion

7 References

# Embedding Analysis (II)

What next?  
Modeling human behavior using smartphone usage data and (deep) recommender systems

Simon Wiegrefe

1 Motivation and Data

2 Modeling

3 Evaluation and Tuning

4 App-level Results

5 Sequence-level Results

6 Discussion

7 References

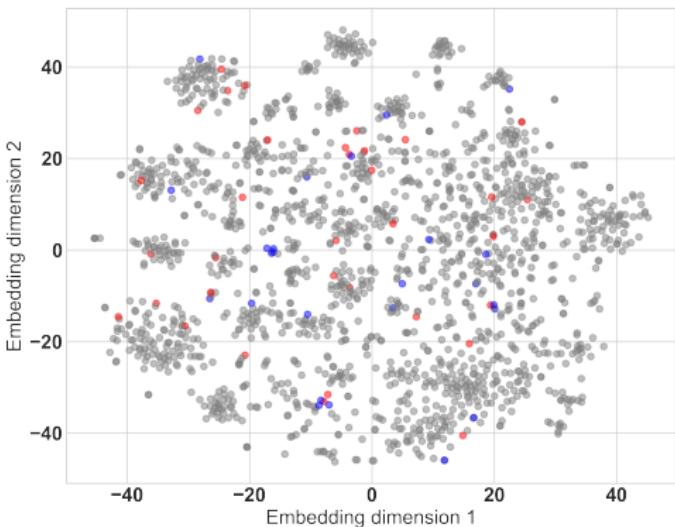


Figure 6: App category-based clustering of app-level embeddings. Blue dots represent apps categorized as *Messaging*, red dots represent apps categorized as *Social Networks*. For illustration, app embeddings are reduced to a dimensionality of two.

- ▶ No category-level clustering recognizable
- ▶ Only for 11.67% of apps their most similar app (i.t.o. cosine similarity) is from the same category

# Embedding Analysis (III)

What next?  
Modeling human  
behavior using  
smartphone usage  
data and (deep)  
recommender  
systems

Simon Wiegrefe

- ▶ Alternatively: start off with data-driven clustering approach k-means
- ▶ Look at potential accumulations of app categories within each cluster

1 Motivation and Data

2 Modeling

3 Evaluation and Tuning

4 App-level Results

5 Sequence-level Results

6 Discussion

7 References

# Embedding Analysis (IV)

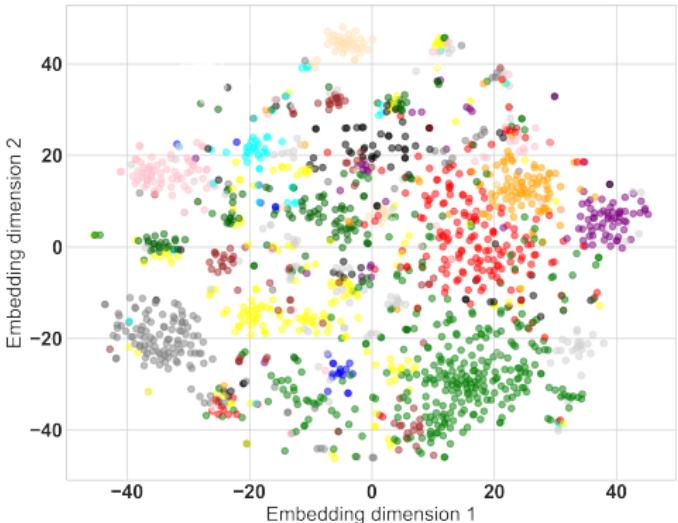


Figure 7: k-means clustering of app-level embeddings ( $k = 15$ ). For illustration, app embeddings are reduced to a dimensionality of two.

- ▶ Moccasin-colored cluster: 32 out of 52 apps (>60%) are camera or image editing apps
- ▶ However: vast majority of clusters dispersed across app space, with little intra-cluster app category clustering.

- ▶ Experimentally construct app analogies such as "Messaging 1 - Social Network 1 + Social Networks 2 = ???"
- ▶ We find no meaningful app analogies in our embeddings:
  - ▶ App analogies conceptually much less intuitive than word analogies
  - ▶ Low overall quality of *GRU4Rec* embeddings

1 Motivation and Data

2 Modeling

3 Evaluation and Tuning

4 App-level Results

5 Sequence-level Results

6 Discussion

7 References

## 5 Sequence-level Results

1 Motivation and Data

2 Modeling

3 Evaluation and Tuning

4 App-level Results

5 Sequence-level Results

6 Discussion

7 References

# Overall Performance

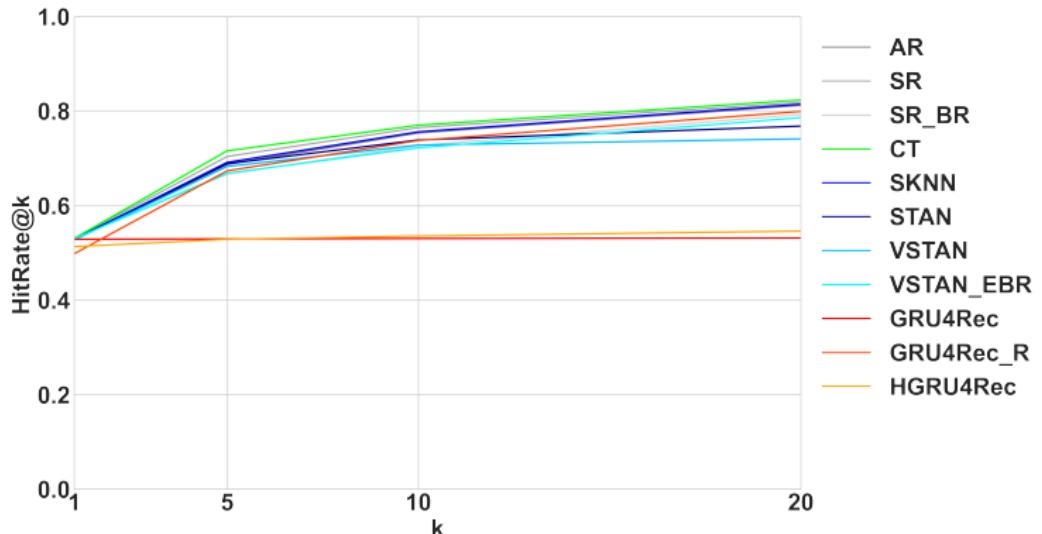


Figure 8:  $HR@k$  performance for  $k = 1, 5, 10$ , and 20 on five-window sequence-level data.

- ▶ Strong  $HR@1$  performance by all algorithms
- ▶ Low performance increases with increasing  $k$
- ▶  $GRU4Rec$  and  $HGRU4Rec$  weakest performers for  $k > 1$

# Removing ON-OFF Tokens (I)

What next?  
Modeling human  
behavior using  
smartphone usage  
data and (deep)  
recommender  
systems

Simon Wiegrefe

- ▶ Suspiciously high  $HR@1$  performance across all algorithms
- ▶ High prevalence of ON-OFF tokens (51.06%)
- ▶ All algorithms predict ON-OFF tokens (almost) everywhere
  - ▶ Predictive performance on other tokens ~0%
- ▶ Effect of removing ON and OFF events from underlying app-level data?

1 Motivation and Data

2 Modeling

3 Evaluation and Tuning

4 App-level Results

5 Sequence-level Results

6 Discussion

7 References

# Removing ON-OFF Tokens (II)

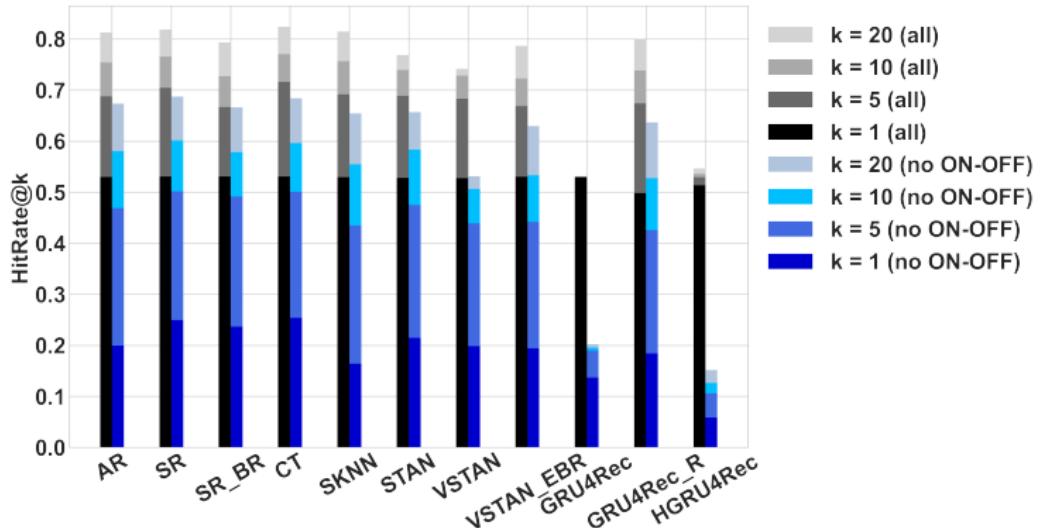


Figure 9:  $HR@k$  performance comparison for all selected algorithms between full five-window sequence-level data (left bars) and five-window sequence-level data after dropping all ON and OFF events (right bars), for  $k \in \{1, 5, 10, 20\}$ .

- ▶ Performance drops for all algorithms, especially i.t.o.  $HR@1$
- ▶ *CT* best, *GRU4Rec* and *HGRU4Rec* worst performers

# Position in Test Sequence (I)

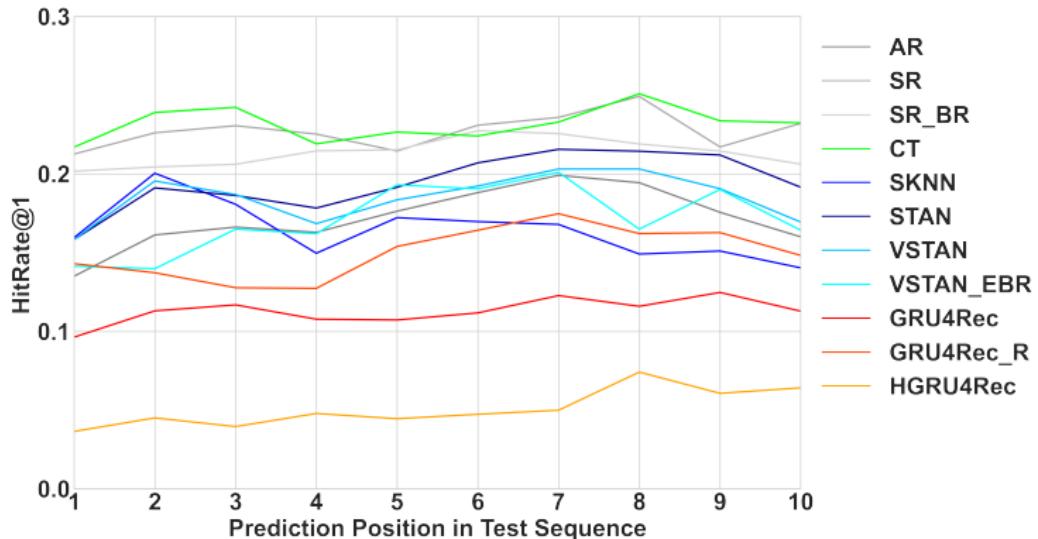


Figure 10: HR@1 performance across the first ten prediction positions on five-window sequence-level data for all selected algorithms.

- ▶ ON and OFF events removed from the underlying app-level data
- ▶ No clear trend for any of the models

# Position in Test Sequence (II)

Algorithm	position <= 2	position > 2	position <= 5	position > 5	position <= 10	position > 10
AR	0.1478	0.2024	0.1598	0.2057	0.1711	0.2096
SR	0.2193	0.2515	0.2218	0.2543	0.2270	0.2581
SR_BR	0.2030	0.2387	0.2083	0.2413	0.2132	0.2454
CT	0.2278	0.2556	0.2288	0.2583	0.2314	0.2622
SKNN	0.1798	0.1629	0.1725	0.1625	0.1645	0.1635
STAN	0.1746	0.2174	0.1808	0.2206	0.1936	0.2226
VSTAN	0.1768	0.1999	0.1783	0.2019	0.1846	0.2037
VSTAN_EBR	0.1406	0.1972	0.1595	0.1991	0.1704	0.2019
GRU4Rec	0.1045	0.1389	0.1081	0.1416	0.1127	0.1462
GRU4Rec_R	0.1402	0.1870	0.1377	0.1916	0.1494	0.1963
HGRU4Rec	0.0405	0.0597	0.0426	0.0614	0.0504	0.0621

Table 4:  $HR@1$  performance results on five-window sequence-level data, by positional cutoff within test sequence.

- ▶ All models except *SKNN* perform better on later positions of the test sequences
- ▶ The precise positioning of the cutoff not very relevant

# Position in Test Sequence (III)

What next?  
Modeling human behavior using smartphone usage data and (deep) recommender systems

Simon Wiegrefe

- ▶ For NN-based models: performance improvement for later events in line with expectations
- ▶ Comparison app- versus sequence-level data:
  - ▶ App-level setting: predominantly short sequences
  - ▶ Sequence-level setting: mostly long sequences
- ▶ Corroborates our previous conclusion: training on short and evaluating on long sequences negatively affects the performance of NN-based algorithms

1 Motivation and Data

2 Modeling

3 Evaluation and Tuning

4 App-level Results

5 Sequence-level Results

6 Discussion

7 References

## 6 Discussion

1 Motivation and Data

2 Modeling

3 Evaluation and Tuning

4 App-level Results

5 Sequence-level Results

6 Discussion

7 References

- ▶ By and large, strong predictive performance of most algorithms
- ▶ NN-based models mostly perform well i.t.o.  $HR@1$  and  $HR@5$ 
  - ▶ Amongst them, *HGRU4Rec* is often the weakest one
- ▶ NN-based model performance is prone to sequence length and data size
- ▶ NN-based models are very expensive i.t.o. runtime and computational effort
- ▶ Baseline models are the preferable modeling choice for our data

1 Motivation and Data

2 Modeling

3 Evaluation and Tuning

4 App-level Results

5 Sequence-level Results

6 Discussion

7 References

## Conclusion (II)

- ▶ *CT* recommendable i.t.o.  $HR@1$  and  $HR@5$ , no tuning
- ▶ *SR* exhibits strong performance i.t.o.  $HR@10$  and  $HR@20$ , fast
- ▶ No overarching user-level effects in our data
  - ▶ For predicting future behavioral sequences of a particular user, not overly helpful to know this particular person's past smartphone usage patterns
- ▶ User-level extensions mostly effective, especially for short sequences and early positions
  - ▶ not due to some profound user-level learning
  - ▶ instead, addressing technical weaknesses
  - ▶ e.g., *VSTAN\_EBR* alleviates poor early-position performance of other neighborhood-based models stemming from low informational content in short sequences

1 Motivation and Data

2 Modeling

3 Evaluation and Tuning

4 App-level Results

5 Sequence-level Results

6 Discussion

7 References

- ▶ Dataset size: potentially giving a relative advantage to non-neural methods
- ▶ Algorithm selection: not including some of the modern sophisticated approaches, e.g., *BERT4Rec* (Sun et al. 2019)
  - ▶ Attention-based models require even more training data
  - ▶ Their main advantage is the better handling of *long-term* dependencies while we mostly have *short* sequences

1 Motivation and Data

2 Modeling

3 Evaluation and Tuning

4 App-level Results

5 Sequence-level Results

6 Discussion

7 References

- ▶ Increased dataset size: new PhoneStudy dataset → Investigate impact of data size on (NN-based) model performance
- ▶ Information extraction: incorporation of duration, exact daytime, and geolocation of app usage
- ▶ Transfer learning: use of pre-trained transformers?

1 Motivation and Data

2 Modeling

3 Evaluation and Tuning

4 App-level Results

5 Sequence-level Results

6 Discussion

7 References

Simon Wiegrefe

1 Motivation and Data

2 Modeling

3 Evaluation and Tuning

4 App-level Results

5 Sequence-level Results

6 Discussion

7 References

**Thank you for your attention!**

## 7 References

1 Motivation and Data

2 Modeling

3 Evaluation and Tuning

4 App-level Results

5 Sequence-level Results

6 Discussion

7 References

- 1 Motivation and Data
- 2 Modeling
- 3 Evaluation and Tuning
- 4 App-level Results
- 5 Sequence-level Results
- 6 Discussion
- 7 References

# References I

- Garg, Diksha, Priyanka Gupta, Pankaj Malhotra, Lovekesh Vig, and Gautam Shroff. 2019. "Sequence and Time Aware Neighborhood for Session-Based Recommendations: Stan." In *Proceedings of the 42nd International Acm Sigir Conference on Research and Development in Information Retrieval*, 1069–72.
- Hidasi, Balázs, Alexandros Karatzoglou, Linas Baltrunas, and Domonkos Tikk. 2015. "Session-Based Recommendations with Recurrent Neural Networks." *arXiv Preprint arXiv:1511.06939*.
- Hinton, Geoffrey, and Sam T Roweis. 2002. "Stochastic Neighbor Embedding." In *NIPS*, 15:833–40. Citeseer.

## References II

Jannach, Dietmar, and Malte Ludewig. 2017. "When Recurrent Neural Networks Meet the Neighborhood for Session-Based Recommendation." In *Proceedings of the Eleventh Acm Conference on Recommender Systems*, 306–10.

Latifi, Sara, Noemi Mauro, and Dietmar Jannach. 2021. "Session-Aware Recommendation: A Surprising Quest for the State-of-the-Art." *Information Sciences* 573: 291–315.

Ludewig, Malte, and Dietmar Jannach. 2018. "Evaluation of Session-Based Recommendation Algorithms." *User Modeling and User-Adapted Interaction* 28 (4-5): 331–90.

Ludewig, Malte, Noemi Mauro, Sara Latifi, and Dietmar Jannach. 2021. "Empirical Analysis of Session-Based Recommendation Algorithms." *User Modeling and User-Adapted Interaction* 31 (1): 149–81.

## References III

Mi, Fei, and Boi Faltings. 2018. "Context Tree for Adaptive Session-Based Recommendation." *arXiv Preprint arXiv:1806.03733*.

Quadrana, Massimo, Alexandros Karatzoglou, Balázs Hidasi, and Paolo Cremonesi. 2017. "Personalizing Session-Based Recommendations with Hierarchical Recurrent Neural Networks." In *Proceedings of the Eleventh Acm Conference on Recommender Systems*, 130–37.

Stachl, Clemens, Quay Au, Ramona Schoedel, Samuel D Gosling, Gabriella M Harari, Daniel Buschek, Sarah Theres Völkel, et al. 2020. "Predicting Personality from Patterns of Behavior Collected with Smartphones." *Proceedings of the National Academy of Sciences* 117 (30): 17680–7.

1 Motivation and Data

2 Modeling

3 Evaluation and Tuning

4 App-level Results

5 Sequence-level Results

6 Discussion

7 References

## References IV

Stachl, Clemens, Ramona Schoedel, Quay Au, Sarah Völkel, Daniel Buschek, Heinrich Hussmann, Bernd Bischl, and Markus Bühner. 2019. "The Phonestudy Project." OSF. <https://doi.org/10.17605/OSF.IO/UT42Y>.

Sun, Fei, Jun Liu, Jian Wu, Changhua Pei, Xiao Lin, Wenwu Ou, and Peng Jiang. 2019. "BERT4Rec: Sequential Recommendation with Bidirectional Encoder Representations from Transformer." In *Proceedings of the 28th Acm International Conference on Information and Knowledge Management*, 1441–50.