

A Unified Contextual Bandit Framework for Long- and Short-Term Recommendations

Maryam Tavakol and **Ulf Brefeld**
{tavakol,brefeld}@leuphana.de

Skopje - Sep 21, 2017

Recommendation



-20 %

POLO RALPH LAUREN
JULIE - Poloshirt - soft flannel heather
★ ★ ★ ★ ★

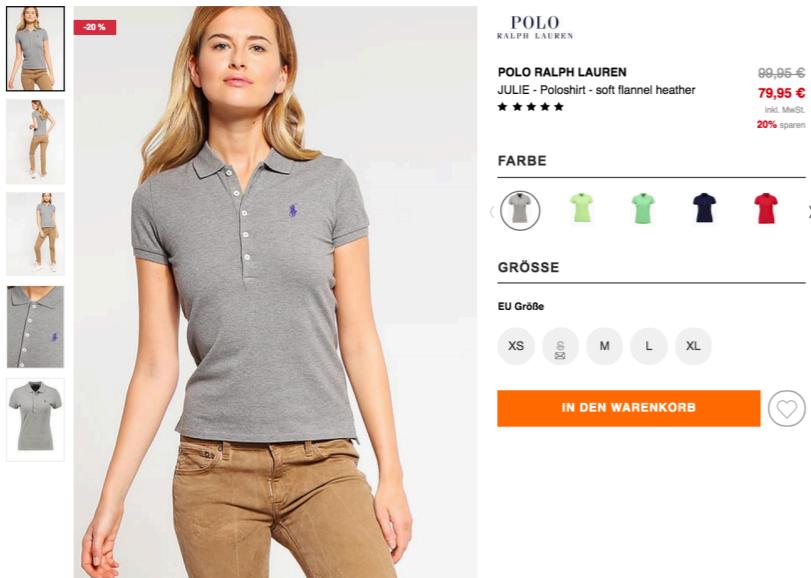
99,95 €
79,95 €
inkl. MwSt.
20% sparen

FARBE

XS S M L XL

IN DEN WARENKORB

Recommendation



DAS KÖNNTE DIR AUCH GEFALLEN



Polo Ralph Lauren
119,95 €



Polo Ralph Lauren
89,95 €



Tommy Hilfiger
59,95 €



Polo Ralph Lauren
129,95 € 103,95 €



Polo Ralph Lauren
109,95 € 87,95 €

Personalization

The Washington Post

SPORTS

THURSDAY, APRIL 19, 2012

SOCER
United's missed shot
A 1-1 draw against the expansion Montreal Impact disappoints a squad that could've used victory. D3



ON WASHINGTONPOST.COM/SPORTS
Tracee Hamilton Today, 11 a.m. Columnist discusses world of Washington-area sports in a Q&A.
Capitals Insider Tonight, 7:30 p.m. Can't get to Game 4? Talk Caps-Bruins on our open thread.
Schedule predictor Cast a vote on which games the Redskins win, and what their record will be.

PRO BASKETBALL
Crawford is Buck wild
Jordan Crawford scores 32 to lead the Wizards past the play-off-chasing Milwaukee Bucks, 123-112. D3



After Chelsea Ties Arsenal,
a Feeling Something Has
Been Lost

On the diamond, another Nats gem

Lucho elimina mejor...
...pero vuelve el verdugo
Aduriz, la gran amenaza • Arda, Busquets e Iniesta, dudas



Cholo y Berizzo se toman la
Copa como asunto personal
Simeone oculta sus cartas como nunca • Los celestes, sin Nolito

Miércoles 27 de enero de 2016 • 1€

104,7 MILLONES DE EUROS
CRISTIANO COSTÓ MÁS QUE BALE

El pago en tres plazos de 72 millones del fichaje del luso elevó el coste final: de los 96,9 acordados hasta 104,7. En el caso de Bale el traspaso se cerró en 91,5 millones y el abono en cuatro plazos aumentó la cantidad final hasta 99,7

Bentancur El Madrid frena su fichaje

Short-Term Zeitgeist



STATE OF THE ART

The iPhone 8: A Worthy Refinement Before the Next Generation

Apple's new iPhone 8 and 8 Plus have a ring of familiarity. But the phones may feel like a solid upgrade from older models because of their new processor.

JIM WILSON/THE NEW YORK TIMES

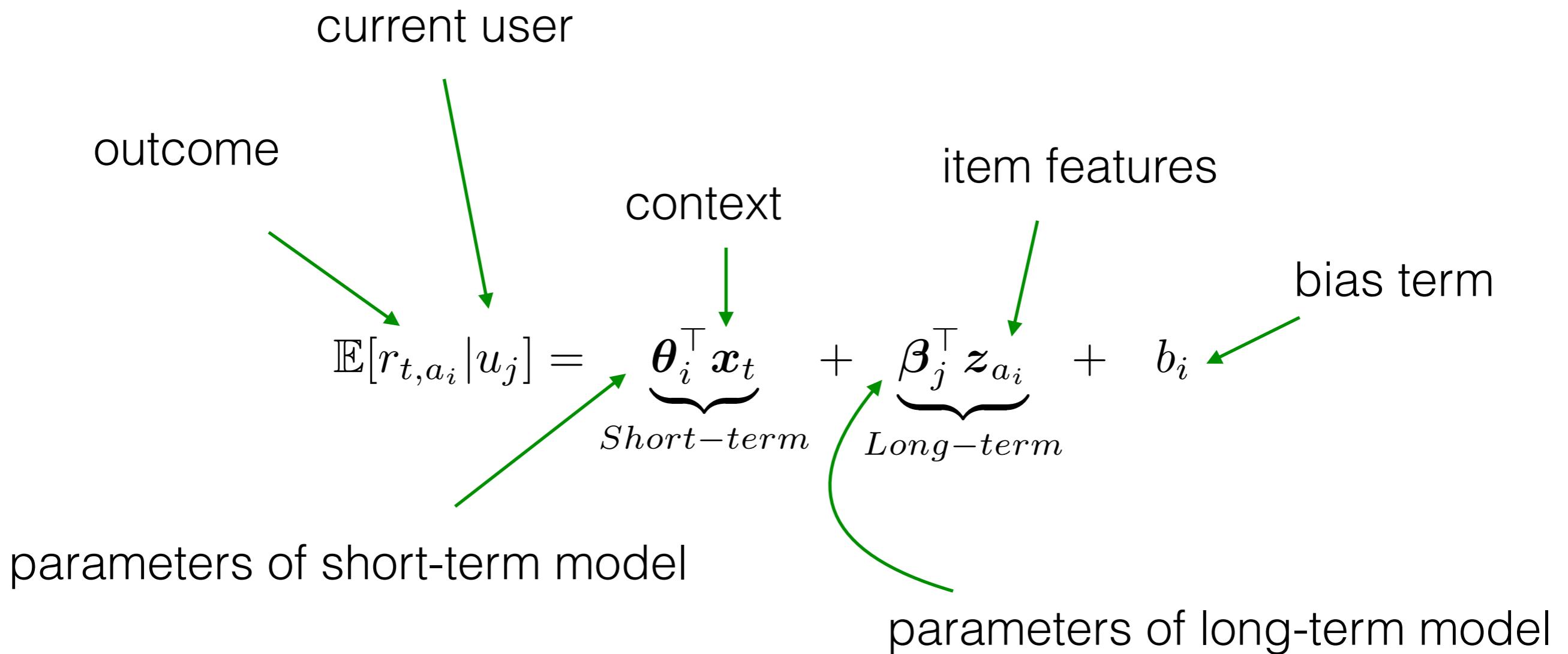


Flying to US to get an iPhone X is cheaper than buying in Europe. It's also illegal

Proposed Approach

- **Goal:** Combination of long- and short-term interests of users in one unified model
 - **Long-term** part + **Short-term** component
- **s.t.:** Generality in terms of optimization
- **Framework:** Contextual Multi-Armed Bandit (MAB)
 - e.g., LinUCB

Unified Model

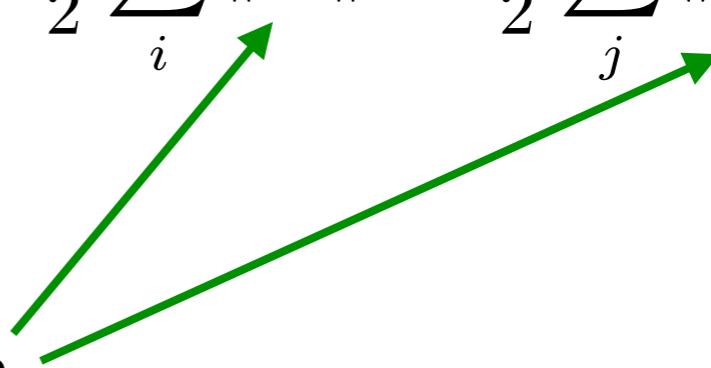


General Optimization

- Objective function with arbitrary loss, $V(\cdot, r_t)$

$$\inf_{\substack{\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_n \\ \boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_m \\ \boldsymbol{b}}} \frac{1}{T} \sum_{t=1}^T V(\boldsymbol{\theta}_t^\top \boldsymbol{x}_t + \boldsymbol{\beta}_t^\top \boldsymbol{z}_t + b_t, r_t) + \frac{\lambda}{2} \sum_i \|\boldsymbol{\theta}_i\|^2 + \frac{\hat{\mu}}{2} \sum_j \|\boldsymbol{\beta}_j\|^2$$

Regularization

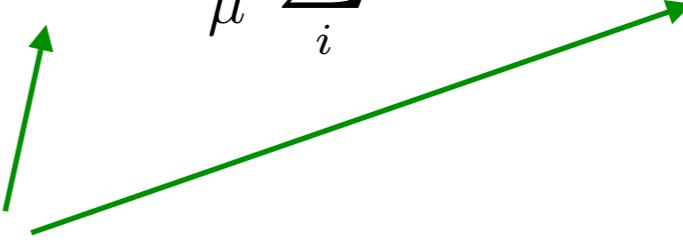


General Optimization

- Using the Fenchel-Legendre conjugate of loss function in the *dual space*:

$$\sup_{\boldsymbol{\alpha}, \mathbf{1}^\top \boldsymbol{\alpha} = 0} - C \sum_{t=1}^T V^*(-\frac{\alpha_t}{C}, r_t) - \frac{1}{2} \boldsymbol{\alpha}^\top [(\sum_i \boldsymbol{\delta}_i \otimes \boldsymbol{\delta}_i^\top) \circ XX^\top + \frac{1}{\mu} (\sum_i \boldsymbol{\phi}_i \otimes \boldsymbol{\phi}_i^\top) \circ ZZ^\top] \boldsymbol{\alpha}$$

Kernel trick



Optimization

- Gradient-based approaches (in dual or primal)
 - Calculating the gradient depends on the loss function
- Model parameters, (θ_i, β_j) , are obtained from α
- Kernel functions applicable

Algorithm

```
for  $t = 1, 2, \dots, T$  do
    Observe user  $u_t$  and context  $\mathbf{x}_t$ 
    for all  $a \in A_t$  do
        Observe arm features  $\mathbf{z}_a$ 
         $v_{t,a} = \text{mean reward} + \text{confidence bound}$ 
    end for
    Choose arm  $a_t = \arg \max_a v_{t,a}$ , and observe payoff  $r_t$ 
    Obtain  $\boldsymbol{\alpha}$  by optimizing the objective
    Compute  $\boldsymbol{\theta}_t$  and  $\boldsymbol{\beta}_t$  from  $\boldsymbol{\alpha}$ 
end for
```

Algorithm

```
for  $t = 1, 2, \dots, T$  do
    Observe user  $u_t$  and context  $\mathbf{x}_t$ 
    for all  $a \in A_t$  do
        Observe arm features  $\mathbf{z}_a$ 
         $v_{t,a} = \text{mean reward} + \text{confidence bound}$ 
    end for
    Choose arm  $a_t = \arg \max_a v_{t,a}$ , and observe payoff  $r_t$ 
    Obtain  $\boldsymbol{\alpha}$  by optimizing the objective
    Compute  $\boldsymbol{\theta}_t$  and  $\boldsymbol{\beta}_t$  from  $\boldsymbol{\alpha}$ 
end for
```

Algorithm

```
for  $t = 1, 2, \dots, T$  do
    Observe user  $u_t$  and context  $\mathbf{x}_t$ 
    for all  $a \in A_t$  do
        Observe arm features  $\mathbf{z}_a$ 
         $v_{t,a} = \text{mean reward} + \text{confidence bound}$ 
    end for
    Choose arm  $a_t = \arg \max_a v_{t,a}$ , and observe payoff  $r_t$ 
    Obtain  $\boldsymbol{\alpha}$  by optimizing the objective
    Compute  $\boldsymbol{\theta}_t$  and  $\boldsymbol{\beta}_t$  from  $\boldsymbol{\alpha}$ 
end for
```

Algorithm

```
for  $t = 1, 2, \dots, T$  do
    Observe user  $u_t$  and context  $\mathbf{x}_t$ 
    for all  $a \in A_t$  do
        Observe arm features  $\mathbf{z}_a$ 
         $v_{t,a} = \text{mean reward} + \text{confidence bound}$ 
    end for
    Choose arm  $a_t = \arg \max_a v_{t,a}$ , and observe payoff  $r_t$ 
    Obtain  $\boldsymbol{\alpha}$  by optimizing the objective
    Compute  $\boldsymbol{\theta}_t$  and  $\boldsymbol{\beta}_t$  from  $\boldsymbol{\alpha}$ 
end for
```

Algorithm

```
for  $t = 1, 2, \dots, T$  do
    Observe user  $u_t$  and context  $\mathbf{x}_t$ 
    for all  $a \in A_t$  do
        Observe arm features  $\mathbf{z}_a$ 
         $v_{t,a} = \text{mean reward} + \text{confidence bound}$ 
    end for
    Choose arm  $a_t = \arg \max_a v_{t,a}$ , and observe payoff  $r_t$ 
    Obtain  $\boldsymbol{\alpha}$  by optimizing the objective
    Compute  $\boldsymbol{\theta}_t$  and  $\boldsymbol{\beta}_t$  from  $\boldsymbol{\alpha}$ 
end for
```

Algorithm

```
for  $t = 1, 2, \dots, T$  do
    Observe user  $u_t$  and context  $\mathbf{x}_t$ 
    for all  $a \in A_t$  do
        Observe arm features  $\mathbf{z}_a$ 
         $v_{t,a} = \text{mean reward} + \text{confidence bound}$ 
    end for
    Choose arm  $a_t = \arg \max_a v_{t,a}$ , and observe payoff  $r_t$ 
    Obtain  $\boldsymbol{\alpha}$  by optimizing the objective
    Compute  $\boldsymbol{\theta}_t$  and  $\boldsymbol{\beta}_t$  from  $\boldsymbol{\alpha}$ 
end for
```

Algorithm

```
for  $t = 1, 2, \dots, T$  do
    Observe user  $u_t$  and context  $\mathbf{x}_t$ 
    for all  $a \in A_t$  do
        Observe arm features  $\mathbf{z}_a$ 
         $v_{t,a} = \text{mean reward} + \text{confidence bound}$ 
    end for
    Choose arm  $a_t = \arg \max_a v_{t,a}$ , and observe payoff  $r_t$ 
    Obtain  $\boldsymbol{\alpha}$  by optimizing the objective
    Compute  $\boldsymbol{\theta}_t$  and  $\boldsymbol{\beta}_t$  from  $\boldsymbol{\alpha}$ 
end for
```

Algorithm

```
for  $t = 1, 2, \dots, T$  do
    Observe user  $u_t$  and context  $\mathbf{x}_t$ 
    for all  $a \in A_t$  do
        Observe arm features  $\mathbf{z}_a$ 
         $v_{t,a} = \text{mean reward} + \text{confidence bound}$ 
    end for
    Choose arm  $a_t = \arg \max_a v_{t,a}$ , and observe payoff  $r_t$ 
    Obtain  $\boldsymbol{\alpha}$  by optimizing the objective
    Compute  $\boldsymbol{\theta}_t$  and  $\boldsymbol{\beta}_t$  from  $\boldsymbol{\alpha}$ 
end for
```

Instantiation: Squared Loss

- Conjugate of the loss function:

$$V^*(-\frac{\alpha_t}{C}, r_t) = \frac{1}{2C^2}\alpha_t^2 - \frac{1}{C}\alpha_t r_t$$

- Becomes a standard quadratic optimization with a constraint
- Confidence bound: $c\sqrt{\mathbf{x}_t^\top (X^\top X)^{-1} \mathbf{x}_t + \mathbf{z}_t^\top (Z^\top Z)^{-1} \mathbf{z}_t}$

Instantiation: Logistic Loss

- Conjugate of the loss function:

$$V^*(-\frac{\alpha_t}{r_t}, r_t) = (1 - \frac{\alpha_t}{Cr_t}) \log(1 - \frac{\alpha_t}{Cr_t}) + \frac{\alpha_t}{Cr_t} \log(\frac{\alpha_t}{Cr_t})$$

- Confidence bound:

Diagonal matrix of sigmoid model

$$c\sqrt{\mathbf{x}_t^\top (X^\top V_a X)^{-1} \mathbf{x}_t + \mathbf{z}_t^\top (Z^\top V_u Z)^{-1} \mathbf{z}_t}$$

Model Simplification

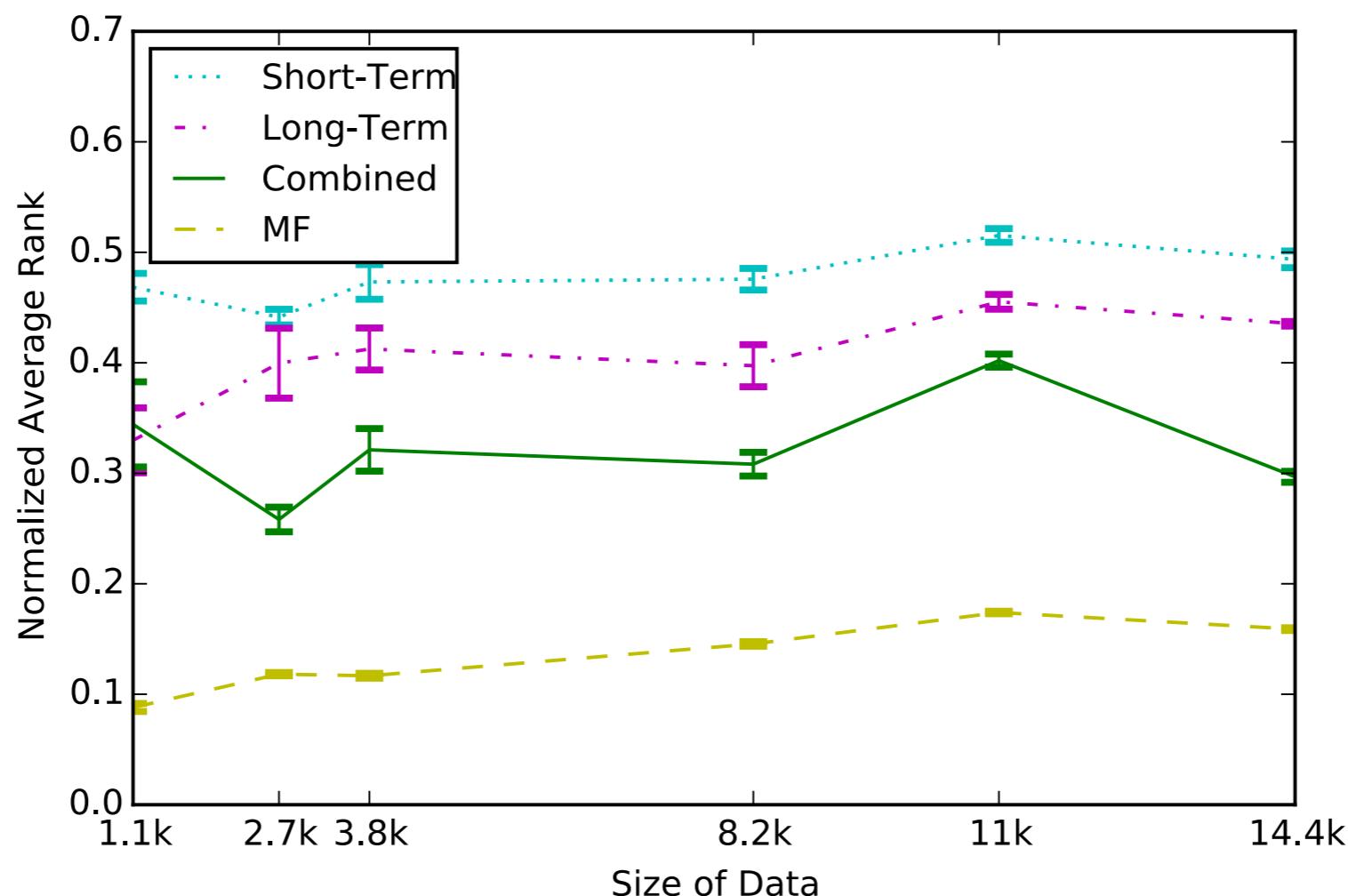
- Focus on the item model
 - **Short-Term:** $\mathbb{E}[r_{t,a_i}] = \theta_i^\top x_t$
 - **Short-Term+Average:** $\mathbb{E}[r_{t,a_i}] = \theta_i^\top x_t + \beta^\top z_{a_i}$
- Focus on the user model
 - **Long-Term:** $\mathbb{E}[r_{t,a_i}|u_j] = \beta_j^\top z_{a_i}$
 - **Long-Term+Average:** $\mathbb{E}[r_{t,a_i}|u_j] = \beta_j^\top z_{a_i} + \theta^\top z_{a_i}$

Empirical Study

- Using **squared** loss function
- Dataset: User transactions from Zalando*
- Baseline: Matrix Factorization (MF)
- Performance measure: normalized average rank

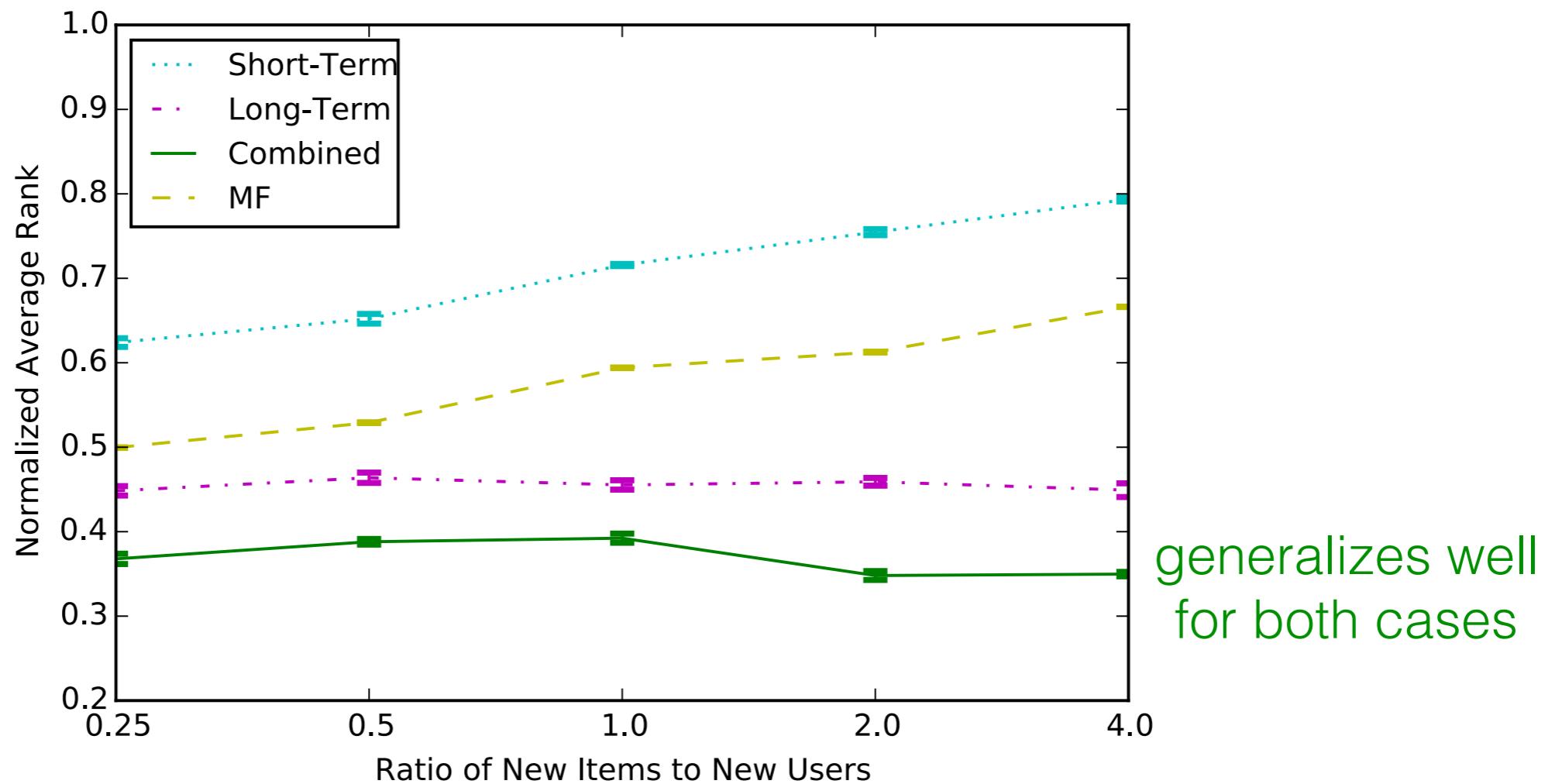
No New User/Item

- The combined approach outperforms either short- or long-term models —but not the baseline!



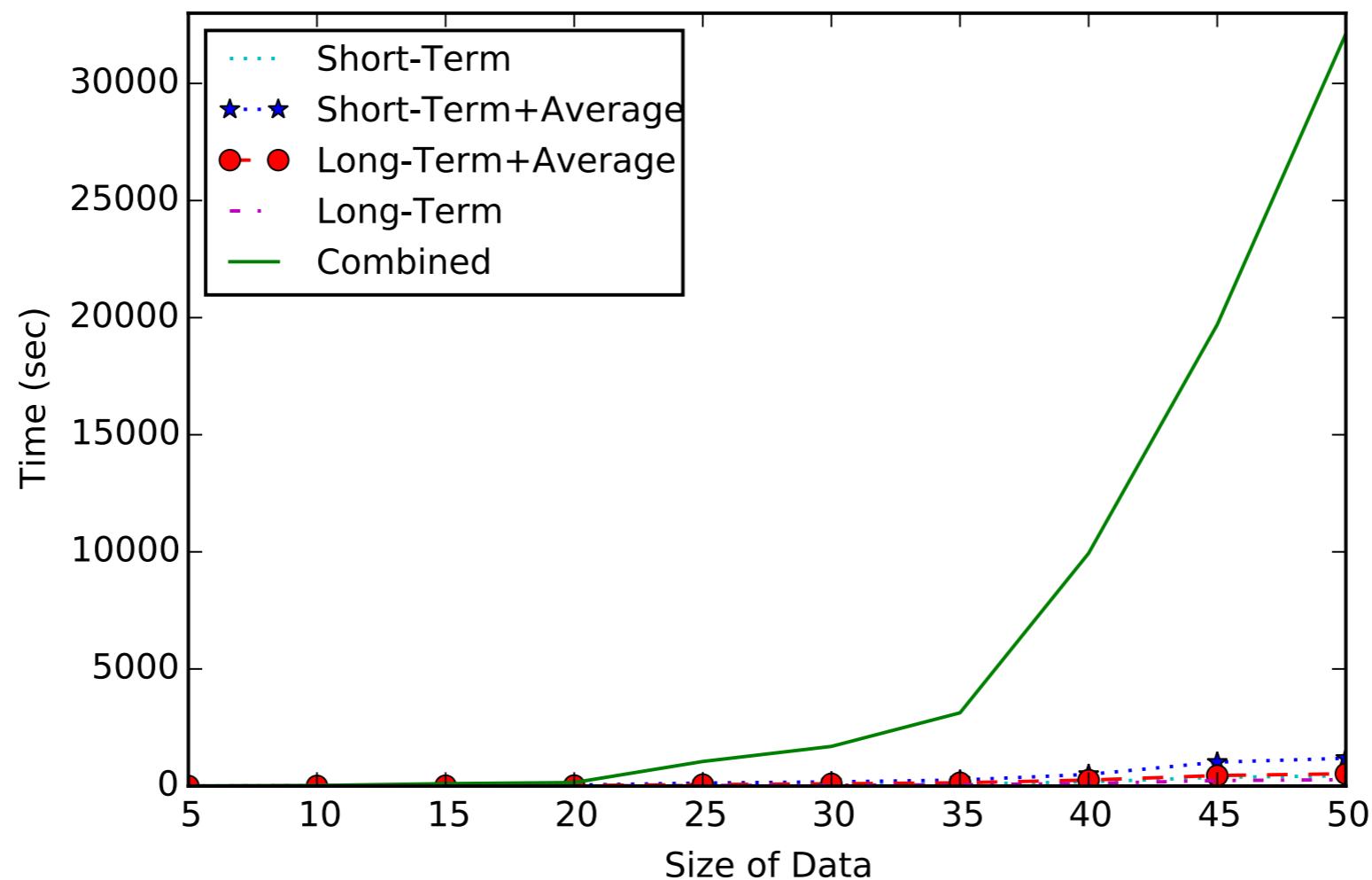
Cold Start Scenarios

- Robustness of combined model in case of new user/item



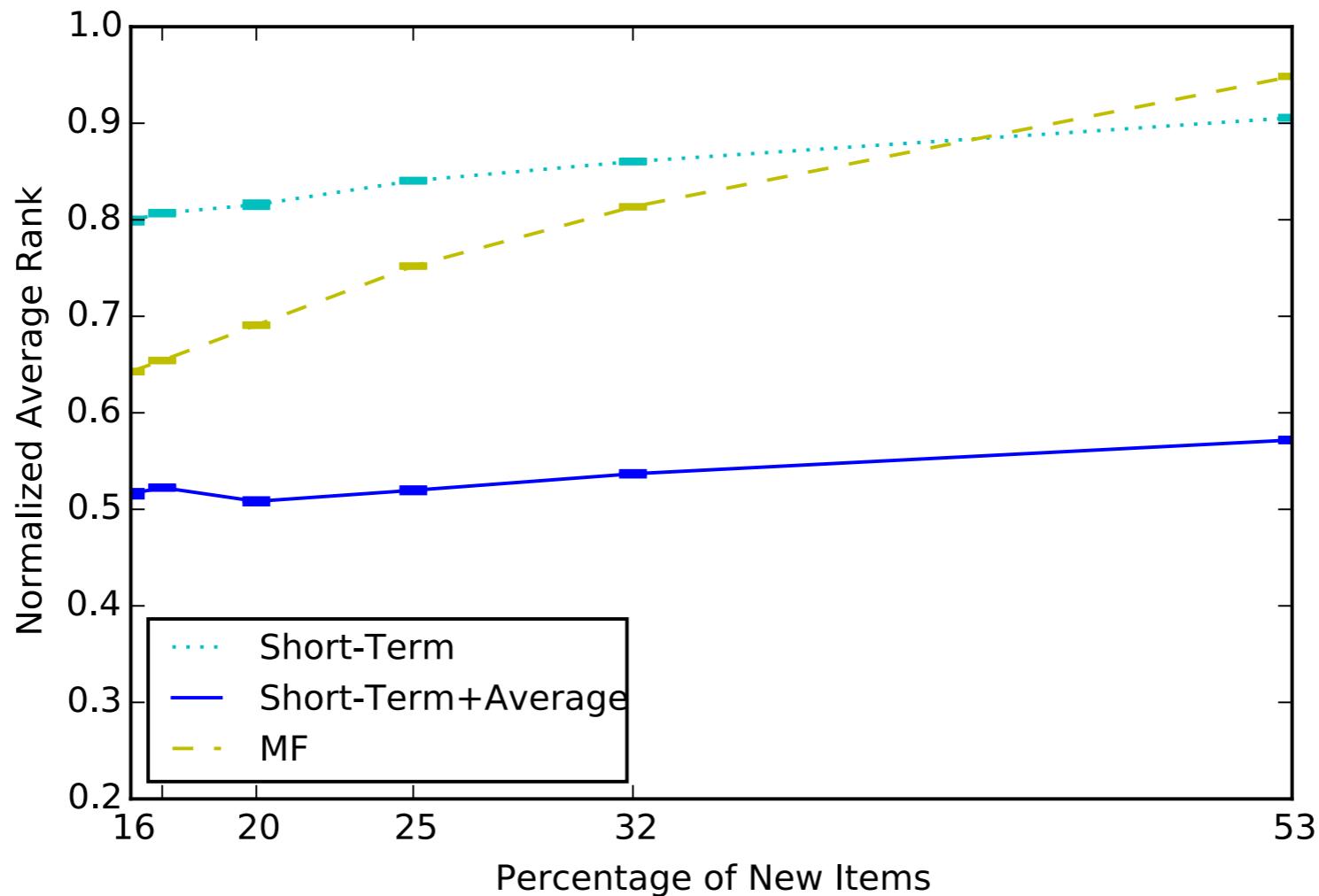
Time Complexity

- The optimization time in combined model is exponential



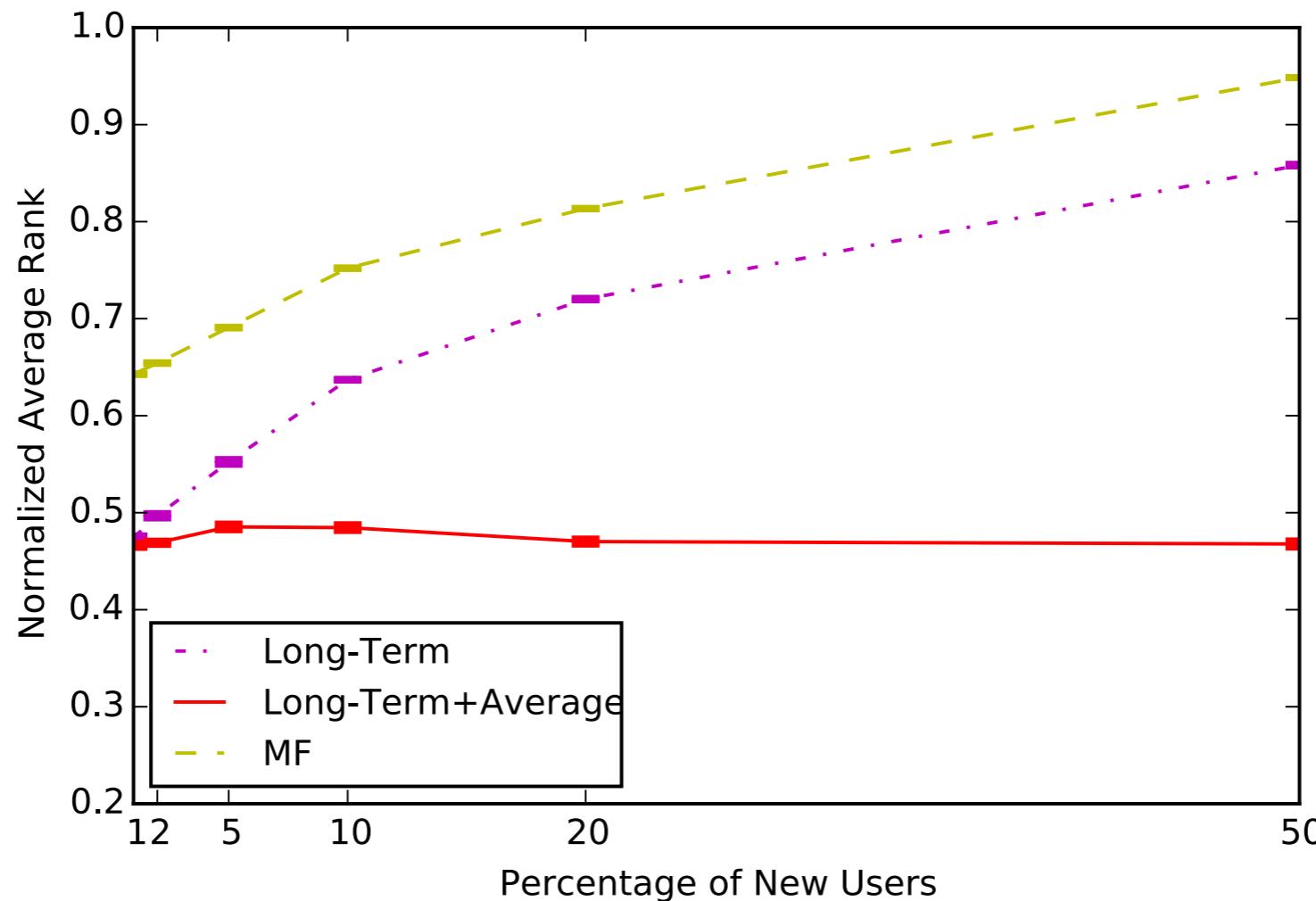
Short-Term Models

- The average term compensates for the new items



Long-Term Models

- The average term compensates for the new users



Conclusion

- The short- and long-term interests of users are combined in one model
- Free choice of *loss function* and *model complexity*
 - There is not one best model: the choice depends on the application

Questions?

Thanks for your attention

A Unified Contextual Bandit Framework for Long-
and Short-Term Recommendations

Maryam Tavakol & Ulf Brefeld
{tavakol,brefeld}@leuphana.de

Source code available at <https://github.com/marytavakol/Bandits>