

Recommender Systems, McNamara, and the Illusion of Progress

Dietmar Jannach, University of Klagenfurt, Austria

dietmar.jannach@aau.at

Presented NeuRec Workhsop @ ICML 2021

December 2021

Recommender systems (RS)

- Automated recommendations everywhere

Recommended For You [Learn more](#)

BO BURNHAM
Bo Burnham: Words, Words, Words (TV Special 2010)
 TVMA Documentary | Comedy | Music
 ★★★★★★ 8.2/10
 The internet (and soon to be movie, TV, radio, etc.) phenomenon, Bo Burnham, brings you an hour stand-up comedy special, Words, Words, Words, of Blues in B...

WORDS WORDS WORDS
 Add to Watchlist
 Next »

Director: Sh...
 Stars: Bo Bu...

◀ Prev 6 Next 6 ▶

Who to follow

Customers Who Bought This Item Also Bought

Quantifying the User Experience: Practical...
 > Jeff Sauro
 ★★★★★ 8
 Paperback
 \$40.63

Recommended

<p>6 ABANDONED Disney Attractions That Are Still... Offhand Disney 1.9M views • 1 year ago</p>	<p>Kung Fu Monk vs Other Masters Don't Mess With... Fight Light 42M views • 8 months ago</p>	<p>Experiment: Coca Cola VS Mentos MrGear 21M views • 6 days ago</p>	<p>The Greatest Showman - Never Enough (Video con... Warner Music Spain 38M views • 1 year ago</p>	<p>German Shepherd Protects Babies and Kids Compilatio... MAI PM 23M views • 2 years ago</p>
--	--	--	--	--

Recommender systems (RS)

- “Everything is a recommendation”

The image displays three overlapping screenshots of popular digital platforms, each showcasing a recommender system:

- Netflix:** The top section, "Top Picks for Joshua," features personalized movie and TV show recommendations like "Breaking Bad" and "Sing." Below, "Trending Now" highlights "shameless," and "Because you watched Narcos" suggests "SURVIVING ESCOBAR." "New Releases" includes "BEYOND STRANGER THINGS."
- Facebook:** The "News Feed" shows a post from Matt Schlicht asking, "anyone else seeing this new facebook?" with options to "Update Status," "Events," "Photos," "Groups," and "Notes." Below, a post from Jessica Mah promotes a TEDxBerkeley registration, and another from Jason L. Baptiste and Mark Tongshuai Bao promotes a shredder.
- Twitter:** The "Home" feed shows tweets from Brie (@Skitch_ComedyFan) and Harold (@h_wang88). The right sidebar features "Trends for you" with categories like "#BreakingNews," "#WorldNews," "#BreakingNews," and "#GreatestOfAllTime."

There's substantial business value¹

- **Amazon in 2006:** About 35 % of revenue attributed to cross-selling, e.g., through “customers who bought”
- **YouTube in 2010:** About 60 % of the clicks on the home screen are on the recommendations
- **Netflix in 2012:** 75 % of what people watch is from some sort of recommendation

¹Jannach, D. and Jugovac, M.: "Measuring the Business Value of Recommender Systems". ACM Transactions on Management Information Systems, Vol. 10(4). 2019

Business value

- Direct profitability
 - Sales, revenue, profit
 - Effects on sales distribution
- Click-through-rate (CTR)
- Adoption of recommendation (“Long-CTR”)
- Engagement with service
 - Time spent on the platform, discovery effects
 - Number of interactions per session or user
- Retention rates
- ...

Direct,
short-term

Indirect,
long-term



Business value

- Key Performance Indicators
 - General ones
 - Revenue, Clicks/Page Impressions, etc.
 - Service-specific, business-model specific ones
 - eBay: “purchase-through-rate”, “bid-through-rate”
 - LinkedIn: Contact with employer made
 - Paper recommendation: “link-through”, “cite-through”
 - E-Commerce marketplace: “click-outs”
 - Online dating: “open communications”, “positive contacts per user”

Exciting times

- Huge economic value of recommendations
- Wide adoption in industry
- Modern technology used in organizations
 - Variety of machine learning techniques in place
 - Deep Learning adopted in industry
 - YouTube, Netflix, Facebook, Twitter ...
 - Shortage of experts (like you)
- General expectations in AI / machine learning

Booming academic research

- The Beginnings
 - Tapestry news filtering (1992)
 - GroupLens (1994)
 - Matrix completion, nearest neighbors
- Success reports in e-commerce (early 2000s)
- Netflix Prize (2006-2009)
 - Matrix Factorization
- AI/ML Boom (- present)
 - Deep Learning

Booming academic research

- Countless papers published each year
 - KDD, TheWebConf, IJCAI, SIGIR, NeurIPS ...
 - 117 papers alone in AAAI and IJCAI in 2018/2019¹
 - ACM TOIS, IEEE TKDE, UMUAI, ...
- ACM Conference on Recommender Systems
 - Constantly growing since 2007
- ACM Transactions on Recommender Systems
 - Established in 2021



¹Jannach, D. and Bauer, C.: "Escaping the McNamara Fallacy: Towards more Impactful Recommender Systems Research". AI Magazine, Vol. 41(4). 2020, pp. 79-95

Continued progress

interactions. Extensive experiments on a real-world dataset show that DAM outperforms the state-of-the-art solution, verifying the effectiveness of our attention design and multi-

level and individual-level. Extensive experiments

enables us to directly learn a ranking function. Extensive online and offline experiments deployed on a commercial platform demonstrate that our models significantly increase diversity while preserving accuracy compared to the state-of-the-art sequential recommendation model, and consequently our models improve user satisfaction.

different domains
significant perfor-
state-of-the-art

tion
two
pro
tain
sup
ommendation models

In addition, experiments on two real-world datasets demonstrate the superiority of our approach against other state-of-the-art approaches in terms of ranking accuracy and efficiency.

But wait ...

ables us to directly learn a ranking function. Extensive online and offline experiments deployed on a commercial platform demonstrate that our models significantly increase diversity while preserving accuracy compared to the state-of-the-art sequential recommendation model, and consequently our models improve user satisfaction.

But wait ...

- Higher prediction accuracy **consequently** leads to higher user satisfaction?



Notorious B.I.G.: Mo money mo problems

Measurements, revisited

- Real-world deployment:
 - Sales, Revenue, CTR, Engagement, Retention rate, Sales distributions, Conversion Rate, Time on site, ...
- Academic papers, predominantly:
 - RMSE, MAE, NDCG, Precision/Recall, ...

Use of **proxy measures / offline metrics**:

- **Academia**: Because we mostly don't have access to real systems
- **Industry**: Because A/B tests are expensive/risky

An important question

arXiv.org > cs > arXiv:2011.07931

Search...

Help | Adv

Computer Science > Information Retrieval

[Submitted on 7 Nov 2020]

Do Offline Metrics Predict Online Performance in Recommender Systems?

[Karl Krauth](#), [Sarah Dean](#), [Alex Zhao](#), [Wenshuo Guo](#), [Mihaela Curmei](#), [Benjamin Recht](#), [Michael I. Jordan](#)

Recommender systems operate in an inherently dynamical setting. Past recommendations influence future behavior, including which data points are observed and how user preferences change. However, experimenting in production systems with real user dynamics is often infeasible, and existing simulation-based approaches have limited scale. As a result, many state-of-the-art algorithms are designed to solve supervised learning problems, and progress is judged only by offline metrics. In this work we investigate the extent to which offline metrics predict online performance by evaluating eleven recommenders across six controlled simulated environments. We observe that

A plausible assumption

- Ranking performance **is** important
 - More relevant items appear high up in the list ➡
 - Making it easier for consumers to find them ➡
 - Leading to more consumption/sales
 - Leading to more satisfaction
 - e.g., because something useful was found
- Nothing wrong with this

Consider this case

- Imagine the user liked this



- We recommend



A good recommendation?



- Highly accurate
- But **useful?**
 - Probably already seen
 - Monotone and obvious, satisfying?
 - Maybe not leading to additional engagement
 - Maybe a missed sales opportunity

Maybe still useful as a reminder or recommendation shortcut¹

¹Lerche, L., Jannach, D. and Ludewig, M.: "On the Value of Reminders within E-Commerce Recommendations". In: Proceedings UMAP 2016

Purpose and usefulness

- What is a good (or: **useful**) recommendation?
- Recommender systems can serve a **variety of purposes**¹ for **different stakeholders**²
 - Consumer, e.g.:
 - Reduction of choice overload, discovery, understanding the space of options, entertainment, ...
 - Provider, e.g.:
 - Sales, revenue, customer retention, ...

¹Jannach, D. and Adomavicius, G.: "Recommendations with a Purpose". In: Proceedings RecSys 2016

²Abdollahpouri, H., Adomavicius, G., Burke, R., Guy, I., Jannach, D., Kamishima, T., Krasnodebski, J. and Pizzato, L.: "Multistakeholder Recommendation: Survey and Research Directions". UMUAI, Vol. 30. 2020, pp. 127–158

Purpose and usefulness

- Recommender systems can serve a **variety of purposes** for **different stakeholders**
 - Moreover: usefulness may depend on a variety of factors like user context, revenue models, etc.
- But what do we do?
 - Abstract from domain/application specifics
 - Use **ranking accuracy as only and universal proxy** for everything
 - (Yes, there is research on diversity, serendipity etc., but suffering from similar problems)

The McNamara Fallacy

- Robert McNamara
 - US Secretary of Defense ('61 – '68)
- Vietnam War:
 - Assessing progress (also) with the enemy **body count**
- The Fallacy ...
 - “... involves making a decision based solely on quantitative observations (or metrics) and ignoring all others.” (Wikipedia)



In recommender systems research

- We focus on what we can easily measure
 - Prediction accuracy in offline experiments
- But do our assumptions actually hold?
 - Will (sometimes tiny) increases in accuracy on historical data lead to better recommender systems?
 - Where is the evidence?

Some studies

A Comparison of Offline Evaluations, Online Evaluations, and User Studies in the Context of Research-Paper Recommender Systems

Joeran Beel¹ and Stefan Langer¹

¹Otto-von-Guericke University

Abstract. The evaluation of recommender systems in practice. The attention in the recommendation of different evaluation methods, the context of research-paper recommender filtering approaches in the research paper.

Investigating the Persuasion Potential of Recommender Systems from a Quality Perspective: An Empirical Study

PAOLO CREMONESI and FRANCA GARZOTTO, Politecnico di Milano
ROBERTO TURRIN, Microsoft

User Perception of Differences in Recommender Algorithms

strand^{1,2}, F. Maxwell Harper², Martijn C. Willemsen³, and Joseph A. Konstan²
Computer Science Department, University of Texas at Austin, TX, USA
²GroupLens Research, Department of Computer Science and Engineering, University of Minnesota, Minneapolis, MN, USA
³Human-Technology Interaction Group, School of Innovation Sciences, Eindhoven University of Technology, Eindhoven, The Netherlands
{harper, konstan}@cs.umn.edu, M.C.Willemsen@tue.nl

user evaluation of recommender systems. We present new tools for understanding what makes a user study aimed at understanding the subjective differences users perceive between different collaborative filtering algorithms and how

Offline and Online Evaluation of News Recommender Systems at swissinfo.ch

Florent Garcin
Artificial Intelligence Lab
Ecole Polytechnique Fédérale
de Lausanne
Switzerland
firstname.lastname@epfl.ch

Boi Faltings
Artificial Intelligence Lab
Ecole Polytechnique Fédérale
de Lausanne
Switzerland
firstname.lastname@epfl.ch

Olivier Donatsch
SWI swissinfo.ch
Swiss Broadcasting Corp.
Switzerland

Ayar Alazzawi
SWI swissinfo.ch
Swiss Broadcasting Corp.
Switzerland

Christophe Bruttin
SWI swissinfo.ch
Swiss Broadcasting Corp.
Switzerland

Amr Huber
SWI swissinfo.ch
Swiss Broadcasting Corp.
Switzerland

ABSTRACT

We report on the live evaluation of various news recommender systems conducted on the website *swissinfo.ch*. We demonstrate that there is a major difference between offline and online accuracy evaluations. In an offline setting, we

on Context Trees (CT) [7], systems adapts its model to differences. The model evolves it is always up to date and in real-time.

User-Centric Evaluation of Session-Based Recommendations for an Automated Radio Station

Malte Ludewig
TU Dortmund, Germany
malte.ludewig@tu-dortmund.de

Dietmar Jannach
University of Klagenfurt, Austria
dietmar.jannach@aau.at

What does the literature tell us?

- The choice of algorithm matters **in practice**
 - But often different families of approaches compared
- Offline/Online comparisons with **user studies**
 - Accuracy-based ranking of algorithms in most cases **not aligned** with user-related quality factors
 - Often only small differences in offline rankings
- Industry/Netflix:
 - “.. we do not find [offline experiments] to be as highly predictive of A/B test outcomes as we would like.”¹

¹Carlos A. Gomez-Urbe and Neil Hunt. 2016. The Netflix Recommender System: Algorithms, Business Value, and Innovation. ACM Trans. Manage. Inf. Syst. 6, 4

McNamara Fallacy in RS

- We are obsessed with accuracy improvements
 - Even though algorithms with the same accuracy might recommend very different things¹
 - Even though we do not know if small improvements matter in practice
 - Even though we do not know if offline accuracy is a good proxy for any business KPI or user-perceived qualities

¹Jannach, D., Lerche, L., Kamehkhosh, I. and Jugovac, M.: "What recommenders recommend: an analysis of recommendation biases and possible countermeasures". User Modeling and User-Adapted Interaction, Vol. 25(5).

The question

- Why aren't we looking at **relevant** or **interesting**) questions?
 - Is it easier to publish such research?
 - Algorithms paper needs no explicit research question
 - Paper needs no underlying theory (e.g., about human behavior), even though recommending is an HCI problem
 - No need to involve users in the research
 - No need to craft an experimental design and defend it against reviewers
 - No need for complex statistical analyses
 - No need to discuss research limitations

The question

- Why aren't we looking at **relevant** or **interesting**) questions?
 - Is it easier to publish such research?
 - Getting a paper through at a good (ML) conference or journal remains difficult, of course
 - A lot of competition
 - Mixed review quality
 - Beefing it up with math might help¹

¹Zachary C. Lipton and Jacob Steinhardt. Troubling trends in machine learning scholarship. Queue, 17(1), February 2019.

A general “applied ML” problem?

- Known hyper-focus on accuracy in applied ML¹
- Critical voices: NLP²



🤖 (((yoav' ()J)()J)))

@yoavgo

...

my two cents on why NLP as a field is focusing on the ML-ish / algorithmic / leaderboard-ish aspects (incl., now, LLMs) and not on the underlying language phenomena: it is just so much easier, on so many levels.


¹Kiri L. Wagstaff. Machine learning that matters. In Proceedings of the Twenty-Ninth International Conference on Machine Learning, ICML '12, pages 529–534, 2012.

² <https://twitter.com/yoavgo/status/1431284873151528960>

A general “applied ML” problem?

- Critical voices: NLP



 (((yoav' (J)(J))) @yoavgo · 27. Aug. ...

ML/algo/leaderboarding is conceptually very simple. esp. for people from CS background. there is a clear, well defined, clean problem, with an easy to measure metrics of success, now improve some metric.



4




10



111



 (((yoav' (J)(J))) @yoavgo · 27. Aug. ...

otoh, the "real language problems" are freakishly inherently messy. there are no clear metrics, there are no clear goals. everything is vague, and due to the zipfian nature of language there are corner cases everywhere.

A general “applied ML” problem?

- Critical voices

Buckman's Homepage



Please Commit More Blatant Academic Fraud

Posted on May 29, 2021

¹ <https://jacobbuckman.com/2021-05-29-please-commit-more-blatant-academic-fraud/>

A general “applied ML” problem?

- Critical voices
 - Day-to-day fraud
 - “Trying that shiny new algorithm out on a couple dozen seeds, and then only reporting the best few. “
 - “Running a big hyperparameter sweep on your proposed approach but using the defaults for the baseline.”
 - “Cherry-picking examples where your model looks good, or cherry-picking whole datasets to test on, where you’ve confirmed your model’s advantage”
 - As a result
 - “... even at top conferences, the median published paper contains no truth or insight.”

It gets even worse

- But we sometimes don't even make true progress on these problematic measures
- We reproduced more than 2 dozen algorithms
 - Traditional matrix completion setups
 - Session-based recommendation
 - Session-aware recommendation

Ferrari Dacrema, M., Boglio, S., Cremonesi, P. and Jannach, D.: "A Troubling Analysis of Reproducibility and Progress in Recommender Systems Research". *ACM Transactions on Information Systems*, Vol. 39(2). 2021

Ferrari Dacrema, M., Cremonesi, P. and Jannach, D.: "Are We Really Making Much Progress? A Worrying Analysis of Recent Neural Recommendation Approaches". In: *Proceedings of the 2019 ACM Conference on Recommender Systems (RecSys 2019)*. Copenhagen, 2019

Ludewig, M., Latifi, S., Mauro, N. and Jannach, D.: "Empirical Analysis of Session-Based Recommendation Algorithms". *User Modeling and User-Adapted Interaction*, Vol. 31(1). 2021, pp. 149–181

Latifi, S., Mauro, N. and Jannach, D.: "Session-aware Recommendation: A Surprising Quest for the State-of-the-art". *Information Sciences*, Vol. 573. 2021, pp. 291-315

It gets even worse

- We reproduced more than 2 dozen algorithms
 - Using the original code
 - Using the original datasets
 - Using the original evaluation
- We found
 - Only one single case, where the proposed neural method was **not** outperformed by existing techniques (e.g., nearest neighbors)
 - In one case even recommending popular items was better
 - Severe methodological issues
 - Low reproducibility in general

What to do?

- The situation is difficult
 - Established research paths
 - Academic incentive system will not change
- Yoshua Bengio:
 - Time to rethink the publication process in machine learning (Blog post, 2020)
 - “[...] our current system incentivizes incremental work and creates a lot of pressure and stress on grad students (and researchers in general) to submit as many papers as possible at each deadline.”
 - (Bengio’s name appeared on > 80 papers in 2020)

What to do?

- Such a change in publication culture will take time
- Too late for current students
- No incentives or signs for such a change



The absolute minimum

- Make sure that our work is **reproducible**, share:
 - Code
 - Own model, baseline models, hyperparameter optimization procedure, evaluation, data pre-processing
 - Data
 - Original data, pre-processed data, data splits for evaluation
 - Documentation & “Infrastructure”
 - Readme, installation instructions
 - Scripts to run everything end to end
 - Maybe: Docker image

The absolute minimum

- Make sure that our work **progresses** the field
 - Systematically tune all baselines properly
 - Don't copy results from other papers without reproducing
 - Do not only include algorithms of one type (e.g., the latest deep learning models)
 - Sanity check: popular item recommendations, nearest neighbors
 - Justify the choice of
 - Datasets
 - Evaluation protocol, metrics and cut-offs based on the claims you make

What would be better?

- Investigate more **relevant, but challenging** questions, there are indeed many
 - (This is an advice for professors)
- For example:
 - Recommendation largely is a problem **of human-computer interaction (HCI)**
 - How to **explain** recommendations, how to **persuade** users, how to help them make **better decisions**?
 - How to increase their **trust**?
 - How to interact with them in a **conversation**?
 - How to **visually present** things?

The HCI Perspective

- Often more important than algorithms
 - A/B at swissinfo.ch
 - Around 30% increase in CTR with algorithm
 - **But:** “When the recommendations are at the top-right corner of the page, the CTR is more than double (2.25) the one at the bottom.”
 - Also, recently:

The ‘Unreasonable’ Effectiveness of Graphical User Interfaces for Recommender Systems

Joeran Beel, University of Siegen, Germany, joeran.beel@uni-siegen.de

Haley Dixon, Trinity College Dublin, Ireland, dixonh@tcd.ie

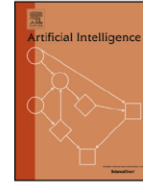
Social Sciences Perspective - Explanations



Contents lists available at [ScienceDirect](#)

Artificial Intelligence

www.elsevier.com/locate/artint

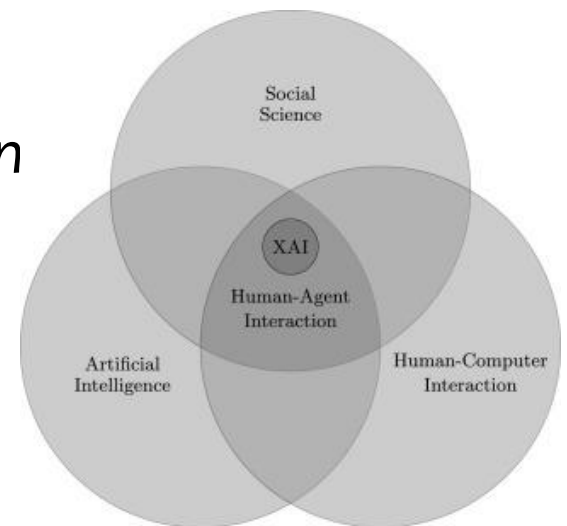


Explanation in artificial intelligence: Insights from the social sciences



Tim Miller

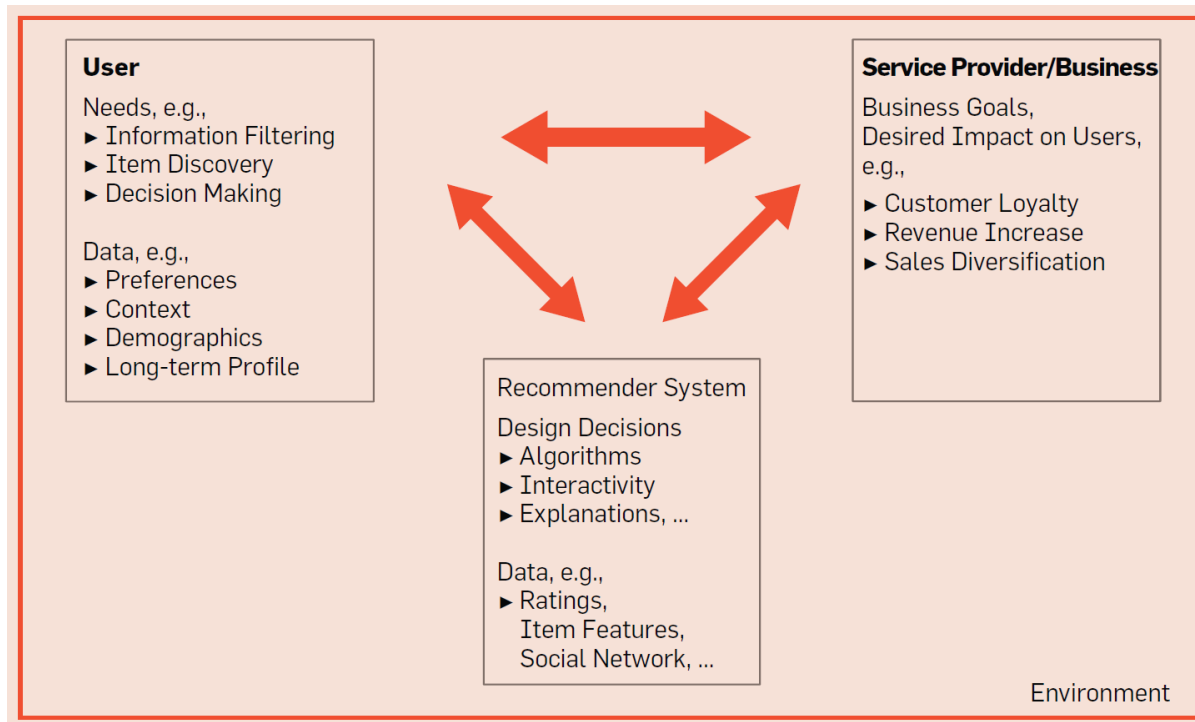
“[...] explainable AI is not just ‘more AI’. Ultimately, it is a human–agent interaction problem.”



The Information Systems Perspective

– Impact on Users and Providers

- Business value, multi-stakeholder perspectives



What about offline + algorithms?

- Offline analyses remain relevant as well
- But maybe a shift would be helpful
 - Understand the problem first
 - What is a good recommendation in a given application context?
 - Diversity/popularity may be good in one application (or a specific user situation) but not in another
 - More analytics-based approaches required first
 - Choose evaluation designs with goal and purpose in mind
 - Validate the proxies
 - Only then try to generalize

What about algorithms?

- Evaluate **with users!**
- Exploring Simulation-based approaches
 - Offline A/B tests, counterfactual evaluation
 - Reduce the offline/online gap
 - “Ground-truth” often remains simulated, various assumptions
 - Agent-based Modeling and Simulation
 - Model system to observe emergent behavior
 - e.g. effects on provider strategy on long-term profitability
 - Typically: various simplifications, mainly to detect trends

Towards a more comprehensive research approach

- Avoid to be a “one-trick pony”
- Let’s grow our methodological repertoire
 - Algorithms-based research
 - Studies with users (randomized controlled trials)
 - Qualitative research methods
 - Surveys, focus groups, towards a better understanding of the problem first
- Let’s take a more scientific approach
 - What are our research hypotheses? What are our underlying theories? Are our experiments suited to support our hypothesis? How can our claims be falsified?

Thank you for your attention

dietmar.jannach@aau.at

<https://tinyurl.com/rs-mc-namara>



Announcement



Association for
Computing Machinery

ACM Transactions on Recommender Systems

- A new journal entirely devoted to recommender systems
- Submissions open in January 2022
 - Technical papers, surveys, industry reports and case studies, opinion pieces