



# Personalized Machine Learning

## Ethics in PML

Rodrigo Alves

December 04, 2025

# Filter Bubbles

- A filter bubble refers to a state in which a user is **only exposed** to information and content that **aligns** with their preferences and beliefs.
- Created by **recommendation algorithms** tailoring content based on user behavior, preferences, and past interactions.
- Filter bubbles can lead to a **limited perspective**, reinforcing existing beliefs, and hindering exposure to diverse opinions and information.
- Users may be exposed to a narrow range of content, limiting their exposure to diverse perspectives.
- Filter bubbles contribute to the creation of **echo chambers**, reinforcing and amplifying existing opinions without exposure to dissenting views.
- Filter bubbles can have implications for **democratic discourse** by limiting the exposure of users to a balanced range of information.

# Filter Bubbles: Mitigation Strategies

**Increased transparency:** recommendation systems should be transparent about their algorithms and provide users with options to understand and control personalized content.

**Diverse recommendations:** algorithms can be designed to intentionally introduce diversity in recommendations, ensuring users are exposed to a broader range of content.

**User Empowerment:** Users should have the ability to customize and adjust their preferences, enabling them to break out of narrow filter bubbles.

# Fairness in Recommendation Systems

- Fairness in recommendation systems refers to the **equitable treatment of users (or items)**, ensuring that recommendations are not biased or discriminatory.
- It is crucial to avoid reinforcing or exacerbating existing **social biases** and to **provide equal opportunities** for all users (items) to discover diverse content.
- Fairness in recommendation systems involves addressing issues of bias, discrimination, and ensuring equitable representation for all user groups.
- Recommendation algorithms can unintentionally reflect and perpetuate biases present in training data, leading to unfair treatment of certain user groups.
- Ensuring fairness across **diverse user groups** (often anonymous) with different preferences, backgrounds, and cultural contexts poses a challenge for recommendation systems.
- Common recommendation evaluation metrics may not capture the nuances of fairness, requiring the development of specialized metrics to assess equity.

# Strategies for Fair Recommendations

**Diverse representation:** ensuring that recommendations represent a diverse set of items and perspectives, avoiding over-representation of certain content or viewpoints.

**Bias mitigation techniques:** implementing techniques to identify and mitigate biases in recommendation algorithms, such as re-ranking, diversity-aware training, and fairness-aware models.

**User empowerment:** incorporating user feedback mechanisms and allowing users to adjust recommendations can empower individuals and contribute to fairer outcomes.

# Explainability and Interpretability

- Explainability and interpretability in recommender systems refer to the **ability to understand** and communicate the reasons behind recommendations to users.
- Transparent and interpretable recommendations enhance **user trust, satisfaction**, and enable users to make **more informed decisions** based on the system's suggestions.
- Challenges:
  - Advanced algorithms are often accurate, especially those based on deep learning, but can be inherently complex to explain.
  - Simple algorithms, like kNN, are easier to explain but might not offer high accuracy.
  - Ensuring that explanations are understandable to users with varying levels of technical expertise poses a significant hurdle in achieving effective explainability.

# Explainability and Interpretability

**Simplifying model complexity:** developing models with a balance between accuracy and simplicity, making it easier to provide clear and understandable explanations.

**User feedback integration:** incorporating user feedback into the recommendation process and explaining how it influences subsequent recommendations, fostering a sense of user control.

**Feature Importance:** highlighting key features or factors that influenced a recommendation, helping users understand the basis of the system's decision-making process.

# GDPR in Personalization

- Explicit consent;
- Data minimization;
- Transparency and Explainability.

Ensure users provide clear and explicit consent for personalized recommendations. Collect only necessary data and be transparent in how it's used. Algorithms should be explainable for transparency.



# GDPR Compliance in Personalization

- Right to access and portability;
- Data security;
- Automated decision-making.

Individuals have the right to access and move their data. Implement robust security measures for data protection. Users should have the option to opt out of automated decision-making, including personalized recommendations.

# GDPR: Implementation and Documentation

- Data protection impact assessments (DPIA);
- Data processing records;
- Data retention and DPO.

In certain cases, conduct DPIA to assess privacy impact. Maintain clear records of data processing. Establish data retention policies and consider appointing a Data Protection Officer (DPO).



Obrigado :) - Faculty of Information Technology