

Quantifying and Reducing Gender Bias in Russian Word Embeddings

Team 2

Anna Koval

Ekaterina Kovalenko

Anastasiia Ryzhova

NLA Project
Skoltech 2018

Word embeddings and gender bias



Gender specific words:

$$\vec{\text{man}} - \vec{\text{woman}} \approx \vec{\text{king}} - \vec{\text{queen}}$$



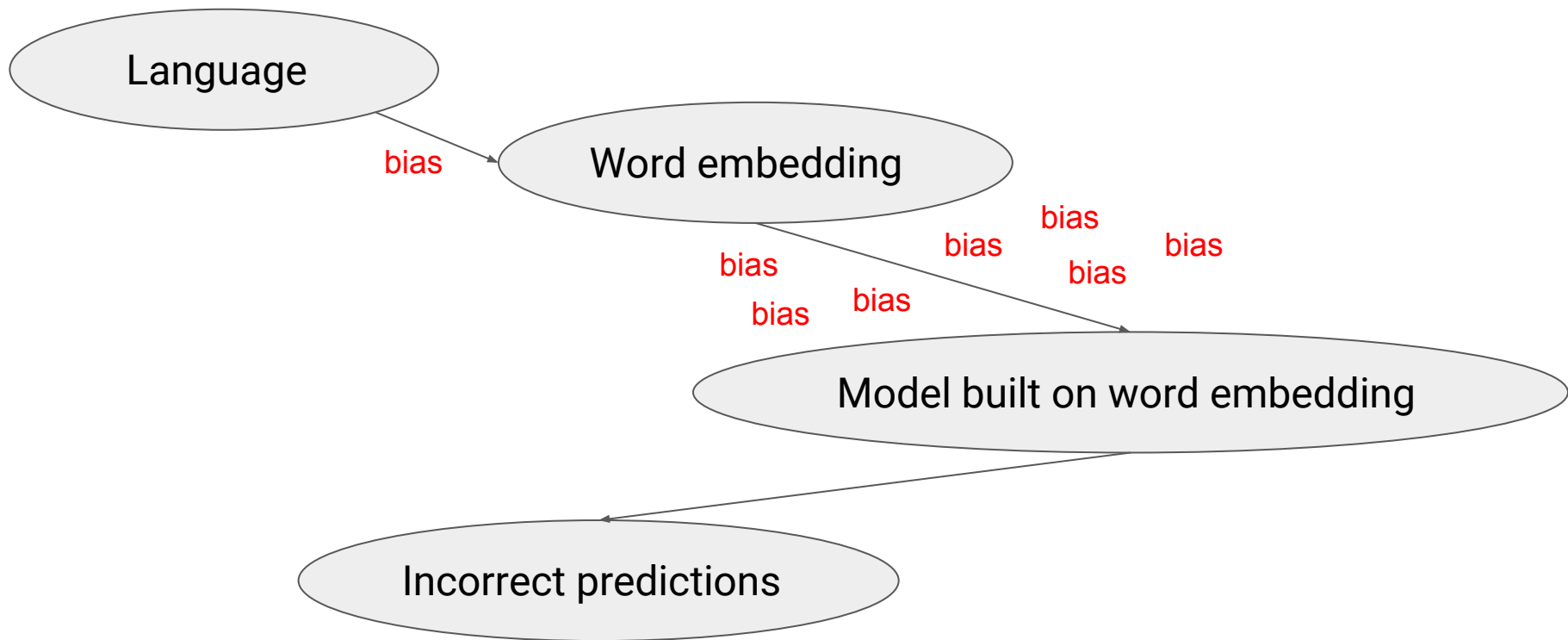
Gender neutral words with bias:

$$\vec{\text{man}} - \vec{\text{woman}} \approx \vec{\text{computer programmer}} - \vec{\text{homemaker}}$$

Problem

- **Gender bias** in a word embedding tends to be **amplified** in a model trained on this word embedding
- Training ML model on unbalanced sample is a well-known problem
- Web search example

Problem



Existing solution

- Data: English *Word2vec* word embeddings
- Detect gender bias in names of occupations (*receptionist, computer programmer*)
- Reduce the bias by solving an optimization problem:

$$\min_{X \leq 0} \|AXA^T - AA^T\|_F^2 + \lambda \|PXb^T\|_F^2 \quad X = TT^T$$

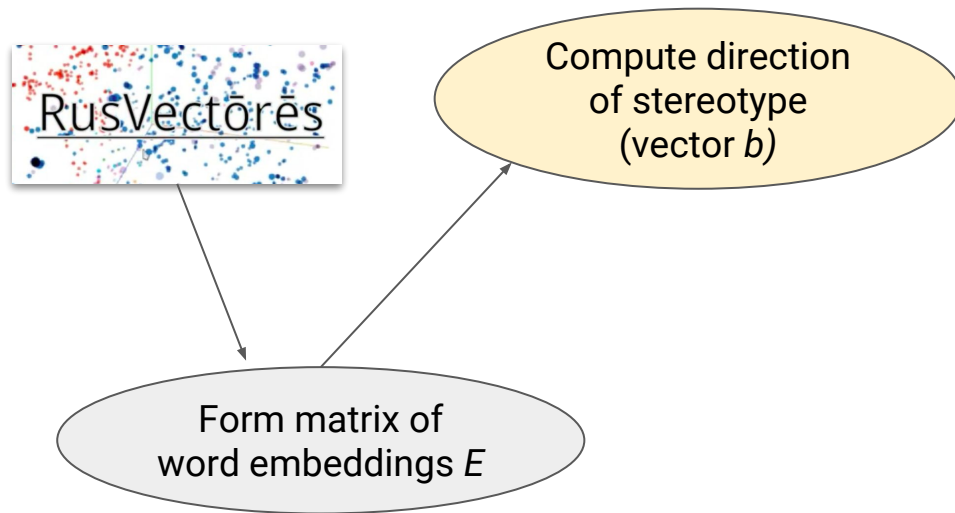
- SVD on A gives $\min_{X \leq 0} \|\Sigma V^T(X - I)V\Sigma\|_F^2 + \lambda \|PXb^T\|_F^2$

Our workflow

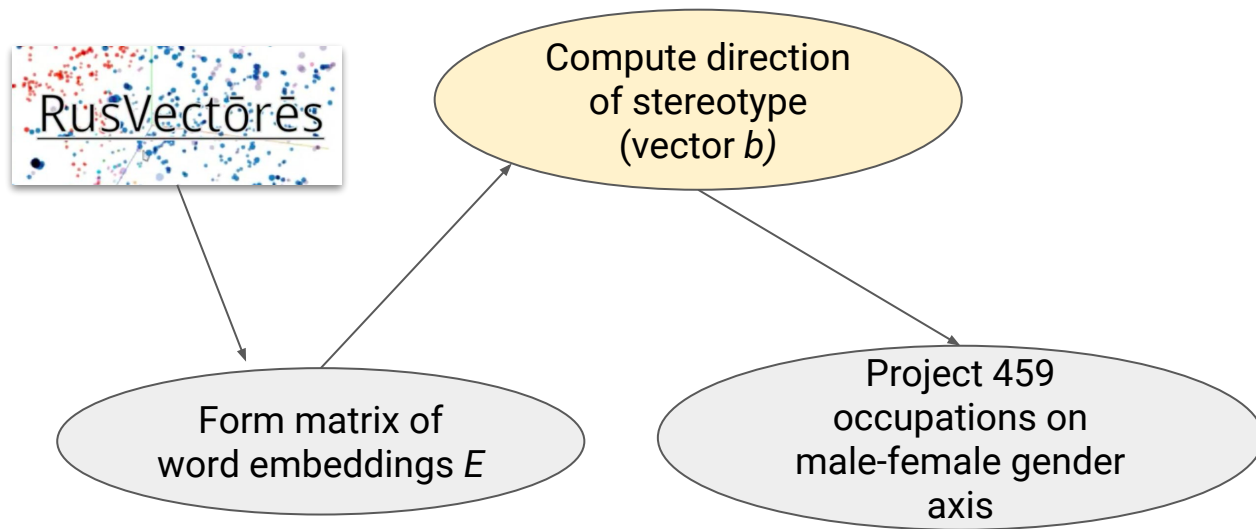


Form matrix of
word embeddings E

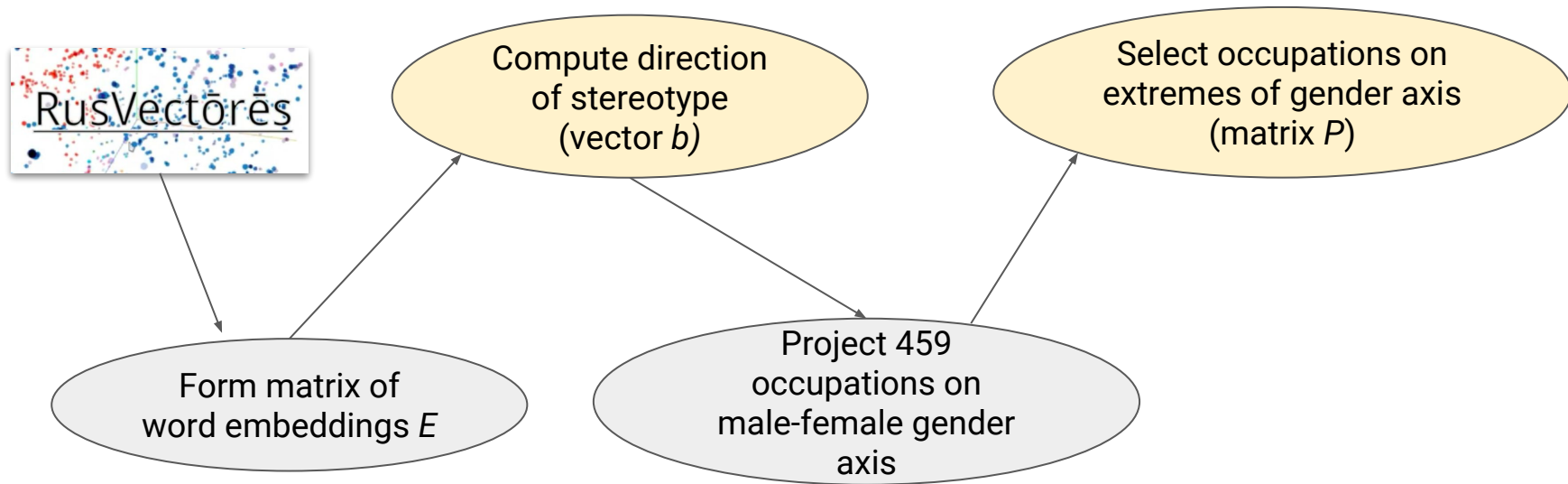
Our workflow



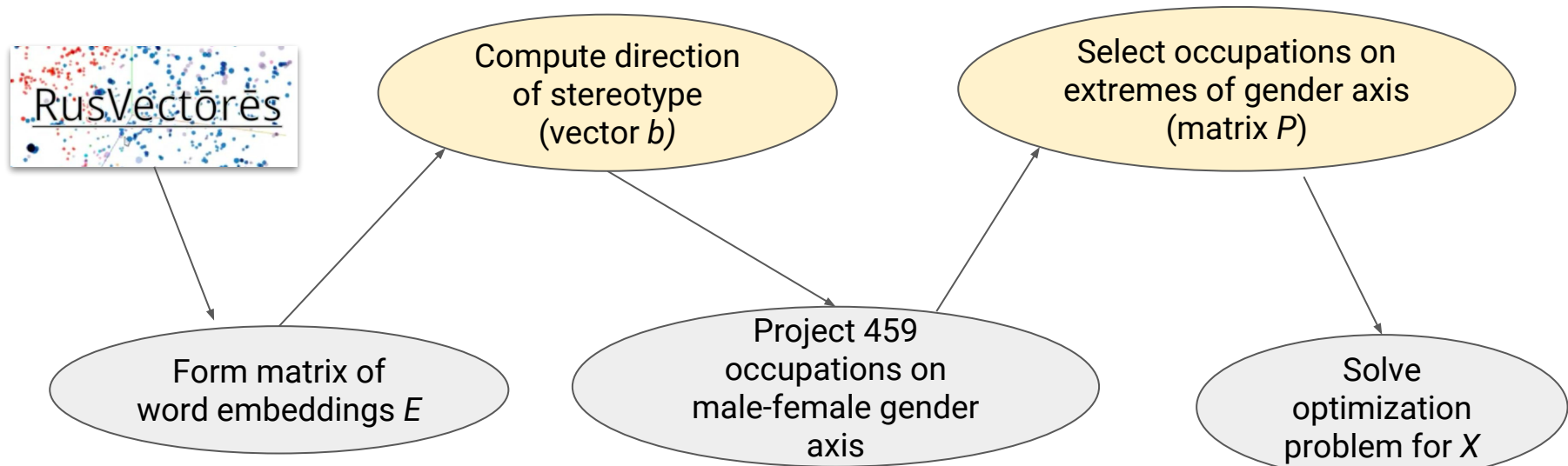
Our workflow



Our workflow

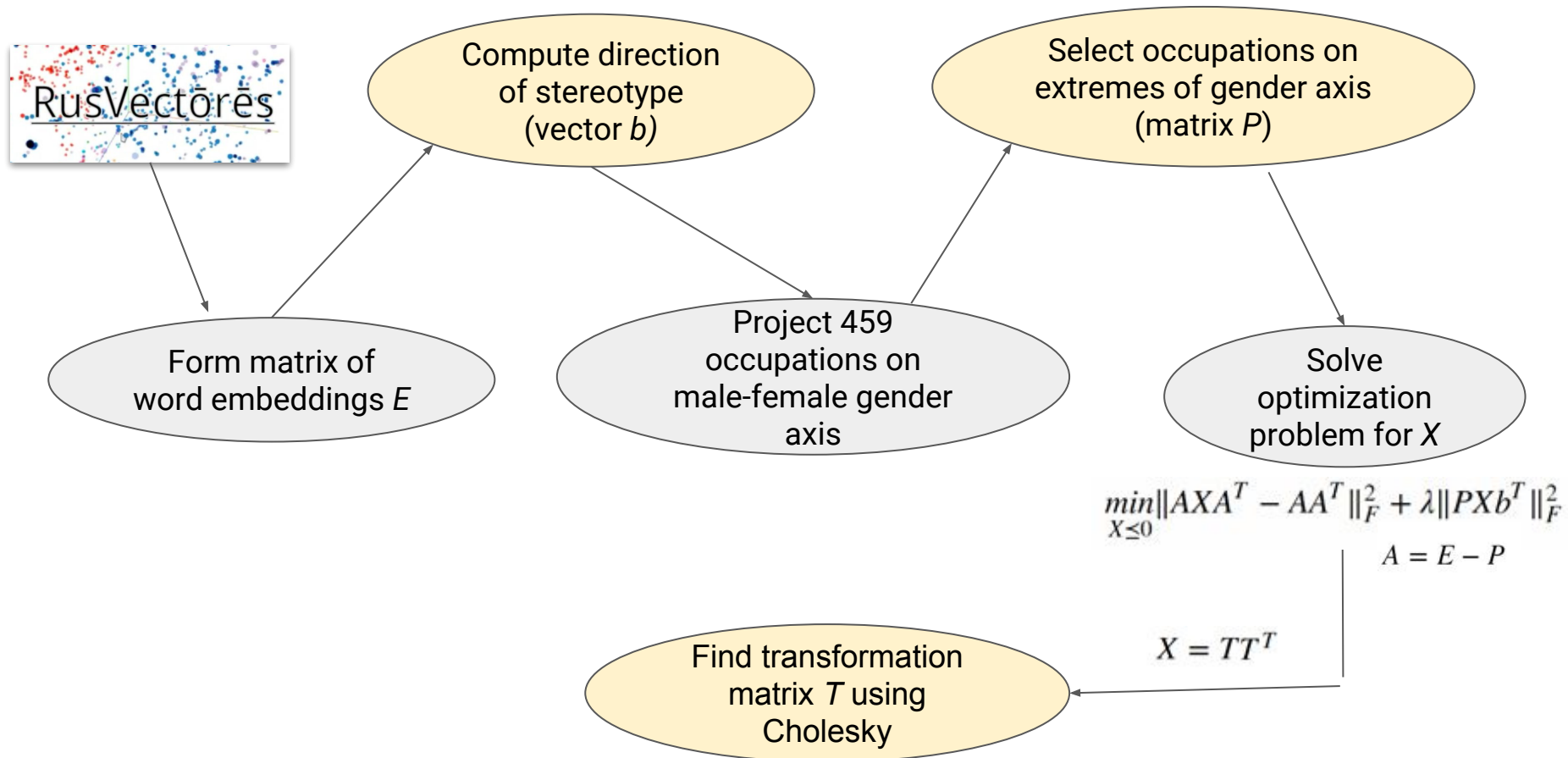


Our workflow

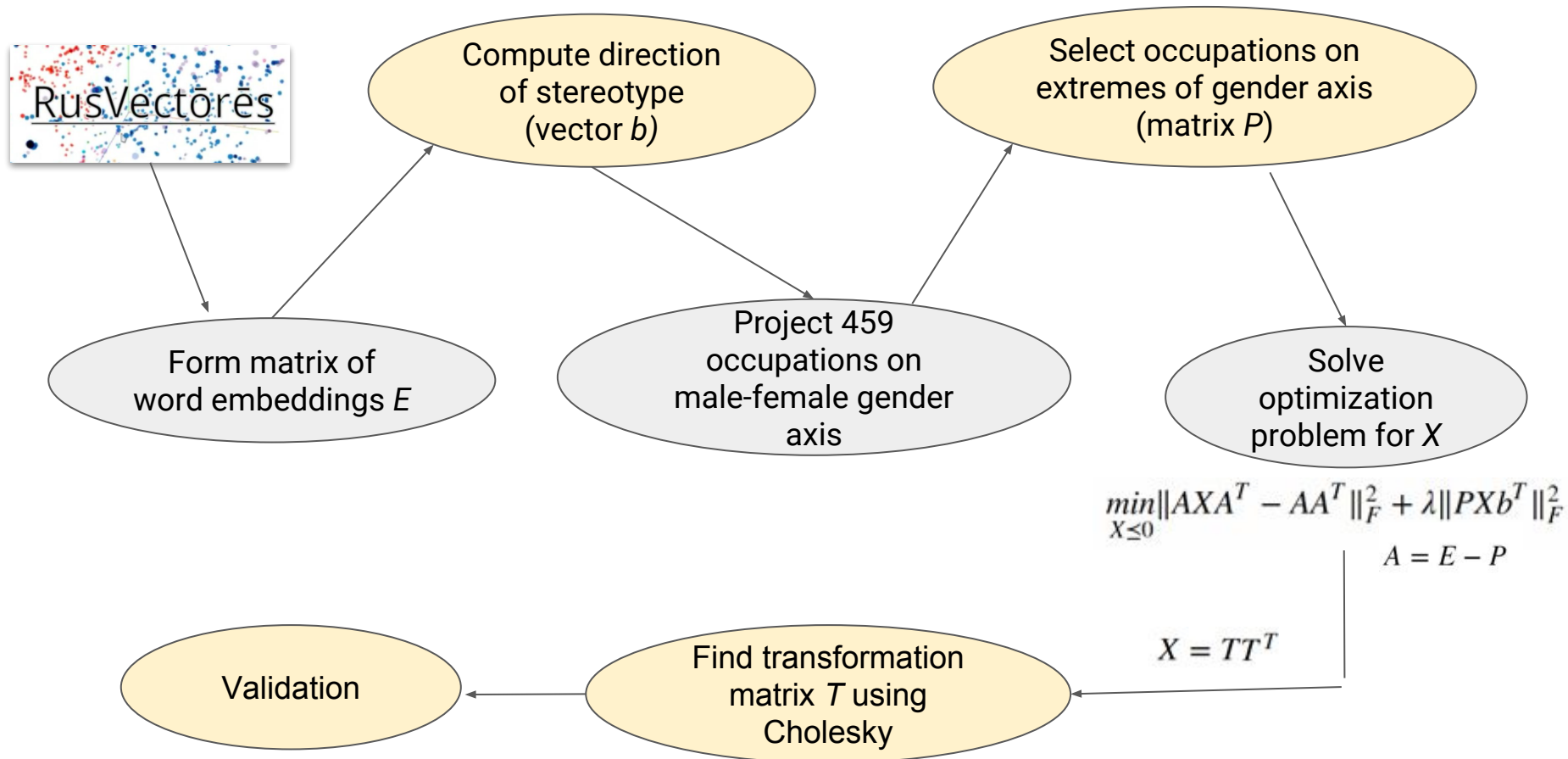


$$\min_{X \leq 0} \|AXA^T - AA^T\|_F^2 + \lambda \|PXb^T\|_F^2$$
$$A = E - P$$

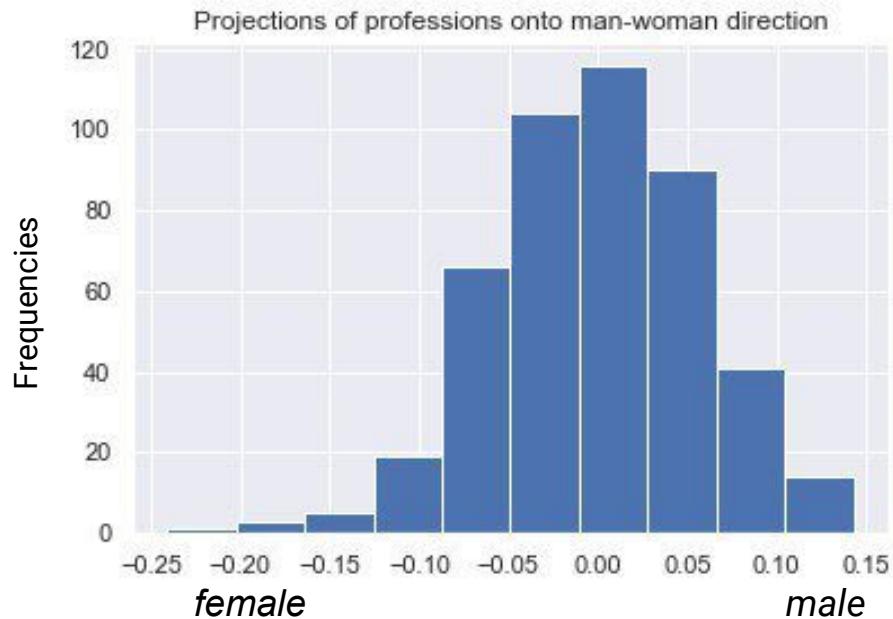
Our workflow



Our workflow



Results: bias detection



Projections on gender axis

Extreme *man* occupations (1 stdev)

менеджер
шеф-повар
губернатор
директор
президент
тренер
программист

Extreme *woman* occupations (1.5 stdev)

учитель
повар
библиотекарь
корректор
художник
врач

Results: bias reduction

Variance of projections before and after transformation

decreased by **1.5** for words with bias

	Projection before transformation	Projection after transformation
губернатор	0.09	0.04
менеджер	0.10	0.06
шеф-повар	0.09	0.04
повар	-0.10	-0.07
библиотекарь	-0.09	-0.05

Discussion

Results:

1. Bias detected but not very extreme compared to research on w2v Google news
2. Debiasing transformation worked to a limited extent only

Conclusions for future:

1. Apply algorithm on word embeddings trained on **news** (we expect more bias)
2. Take a larger matrix P (more occupations)
3. Tune in *lambda* parameter in the optimization task

Methodology can be applied to other languages and other types of bias (e.g. racial)

Thank you for your attention!

Our data

- Word embedding trained on the Russian National Corpus
- Corpus size: 250 million words
- Vocabulary size: 195 071
- Vector size: 300



Computing vector b with PCA

мужчина
↗
женщина

отчим
↗
мачеха

официант
↗
официантка

рубашка
↗
блуза

- Select 45 pairs of words that reflect gender opposites
- Compute 45 differences between pairs of vectors and stack them into a matrix
- Do PCA to find a principal vector - direction of stereotype

PCA returned singular values with very slow decay → we used $v_{\text{мужчина}} - v_{\text{женщина}}$

Optimization problem

$$\min_{X \succeq 0} \|AXA^T - AA^T\|_F^2 + \lambda \|PXB^T\|_F^2$$

$$X = TT^T$$

- Direction of stereotype b
- Matrix of biased words P
- Matrix of background words A
- Transformation matrix T

Validation

- Compute projections of biased words on the gender axis
- Compute the variance of these projections before and after transformation
- The variance should decrease and get close to zero

Contribution of team members

Top 3 for each where 1 is largest contribution:

Anna Koval

1. Presentation and report
2. Linguistics-related decisions
3. Programming (data preprocessing, optimization)

Ekaterina Kovalenko

1. Programming (data preprocessing, svd, optimization)
2. Math
3. Presentation and report

Anastasiia Ryzhova

1. Programming (data preprocessing, projections, optimization)
2. Math
3. Presentation and report