

# Final Project

## Detection of Negative Reviews in Online Stores

---

### **Team #8**

Mikhail Sidorenko

Egor Baryshnikov

Roman Teplykh

### **Mentor**

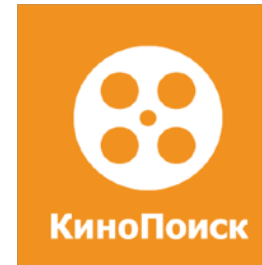
Evgeny Frolov

# Possible applications of semantic analysis

Яндекс.Маркет



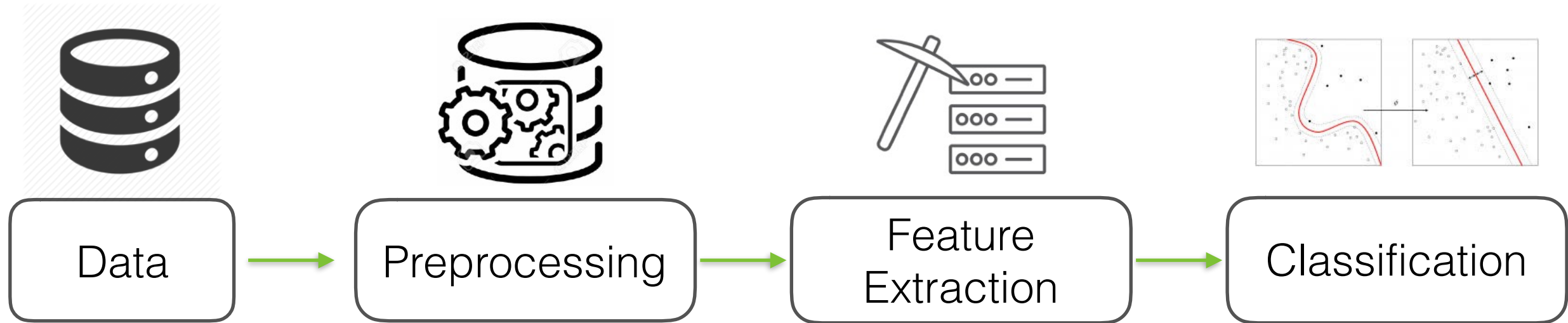
OZON.ru



amazon



# Possible workflow may be look like this



# Data we used

---

## Amazon Reviews dataset\*

\*<http://jmcauley.ucsd.edu/data/amazon/>



- 24 different categories of items
- Include ratings, reviews and other information
- «Cell Phones and Accessories» 194k reviews

# Data preprocessing

---

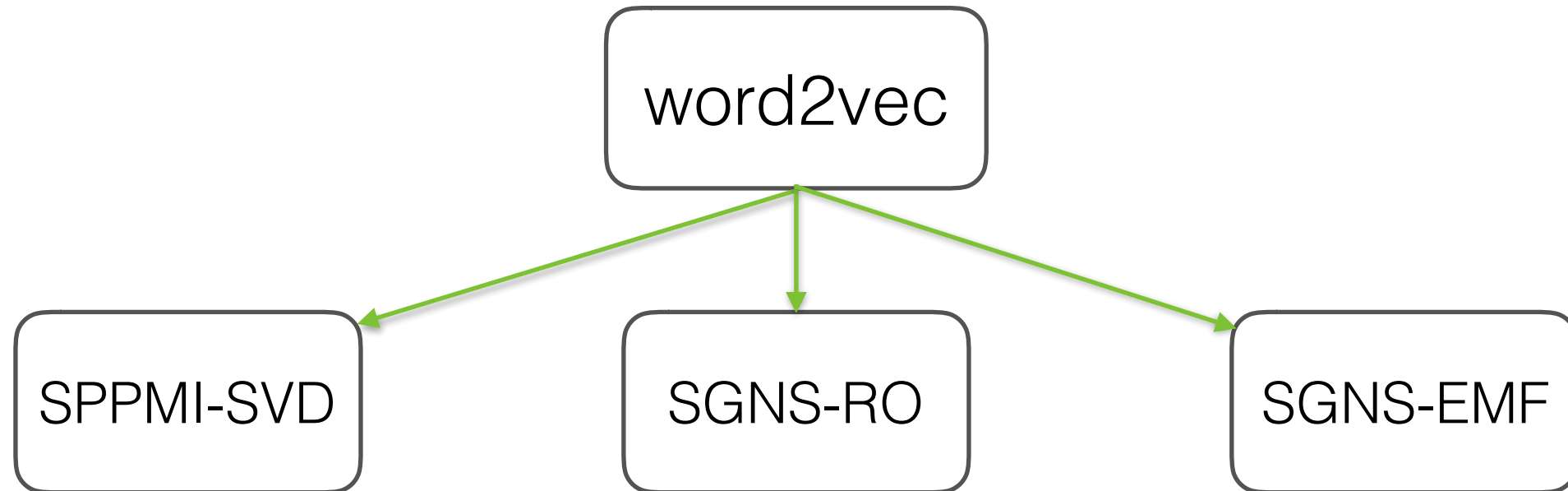
## **Summary**

- vocabulary size: 3723 words
- sliding window size: 2
- two labels: positive/negative review

# Feature extraction

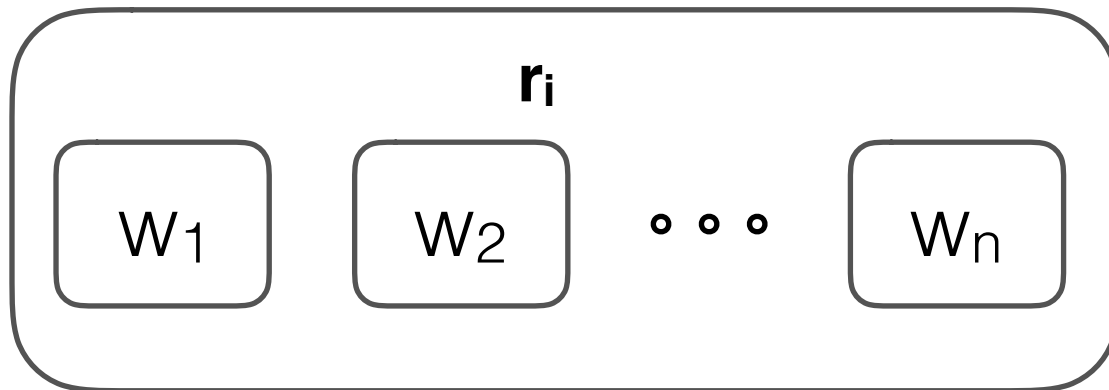
---

## Word Embeddings



# Feature extraction

**i-th review**



$$r_i = \frac{\sum w_j}{n}$$

# word2vec algorithms

---

## SPPMI-SVD

**Idea:** find  $W$  and  $C$  using SVD decomposition of SPPMI matrix

**Disadvantage:** such approach doesn't lead to minimization of SGNS objective

## SGNS-RO

**Idea:** optimize SGNS objective directly on the low-rank matrices space

**Disadvantage:** works in assumption of independence of  $wc$  values



# word2vec algorithms

## SGNS-EMF

**Idea:** explicitly factorize co-occurrence matrix

**Disadvantage:** too many cycles

---

**Algorithm 1:** Alternating minimization for explicit matrix factorization

---

**Input:** Co-occurrence matrix  $\mathbf{D}$ , step-size of gradient descent  $\eta$ , maximum number of iterations  $K$

**Output:**  $\mathbf{C}_K, \mathbf{W}_K$

```

1 initialize  $\mathbf{C}_i$  and  $\mathbf{W}_i$  randomly,  $i = 1$ ;
2 while  $i \leq K$  do
3    $\mathbf{W}_i = \mathbf{W}_{i-1}$ ;
4   //minimize over  $\mathbf{W}$ ;
5   repeat
6      $\mathbf{W}_i = \mathbf{W}_i - \eta \mathbf{C}_{i-1} (\mathbb{E}_{\mathbf{D}'|\mathbf{W}_i, \mathbf{C}_{i-1}} \mathbf{D}' - \mathbf{D})$ ;
7   until Convergence;
8    $\mathbf{C}_i = \mathbf{C}_{i-1}$ ;
9   //minimize over  $\mathbf{C}$ ;
10  repeat
11     $\mathbf{C}_i = \mathbf{C}_i - \eta (\mathbb{E}_{\mathbf{D}'|\mathbf{W}_i, \mathbf{C}} \mathbf{D}' - \mathbf{D}) \mathbf{W}_i^T$ ;
12  until Convergence;
13   $i = i + 1$ ;

```

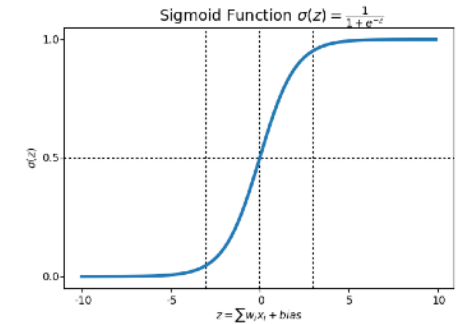
---

\*Li, et al., 2015, «Word Embedding Revisited: A New Representation Learning and Explicit Matrix Factorization Perspective»

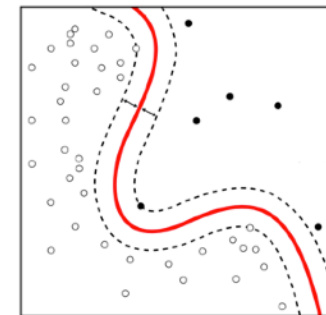
# Classification

## Summary

- two simple classifiers were used «out of the box»
- classification metric: f1-score



Logistic  
Regression



SVM

# Similarity test's results

**Spearman's correlation between predicted similarities and the manually assessed ones (k = 5, alpha=0.5), simlex999 dataset**

d=100	SVD-SPPMI	0.13284
	SGNS-RO	<b>0.13466</b>
	SGNS-EMF	0.03252
d=200	SVD-SPPMI	0.12277
	SGNS-RO	<b>0.13051</b>
	SGNS-EMF	0.06966
d=500	SVD-SPPMI	0.18781
	SGNS-RO	<b>0.18920</b>
	SGNS-EMF	0.06119

# SGNS's objective function values

**The values of SGNS objective function at the optimal point (all values are multiplied by  $10^{-9}$ ,  $k=5$ ,  $\alpha=0.5$ )**

	SVD-SPPMI	SGNS-RO	SGNS-EMF
d=100	-0.2383	<b>-0.2321</b>	-0.3841
d=200	-0.2381	<b>-0.2316</b>	-0.5406
d=500	-0.2357	<b>-0.2300</b>	-0.8484

# SGNS's objective function values

**The values of SGNS objective function at the optimal point (all values are multiplied by  $10^{-9}$ ,  $d=200$ ,  $\alpha=0.5$ )**

	SVD-SPPMI	SGNS-RO	SGNS-EMF
<b>k=1</b>	<b>-0.0758</b>	<b>-0.0742</b>	<b>-0.3467</b>
k=5	-0.2381	-0.2316	-0.5406
<b>k=15</b>	<b>-0.6354</b>	<b>-0.6157</b>	<b>-0.6779</b>

# Classification results

## F1-score values (k=5, alpha=0.5)

		LR	SVC
d=100	SVD-SPPMI	<b>0.87892</b>	<b>0.87901</b>
	SGNS-RO	0.87890	0.87888
	SGNS-EMF	0.86754	0.86849
d=200	SVD-SPPMI	0.88341	0.88338
	SGNS-RO	<b>0.88345</b>	<b>0.88341</b>
	SGNS-EMF	0.87446	0.87492
d=500	SVD-SPPMI	0.89012	0.89016
	SGNS-RO	<b>0.89019</b>	<b>0.89023</b>
	SGNS-EMF	0.88568	0.88580

Thank you for your attention!

---

Any questions?

<https://github.com/Bulldogger/NLA-Project>