



Kaggle text mining

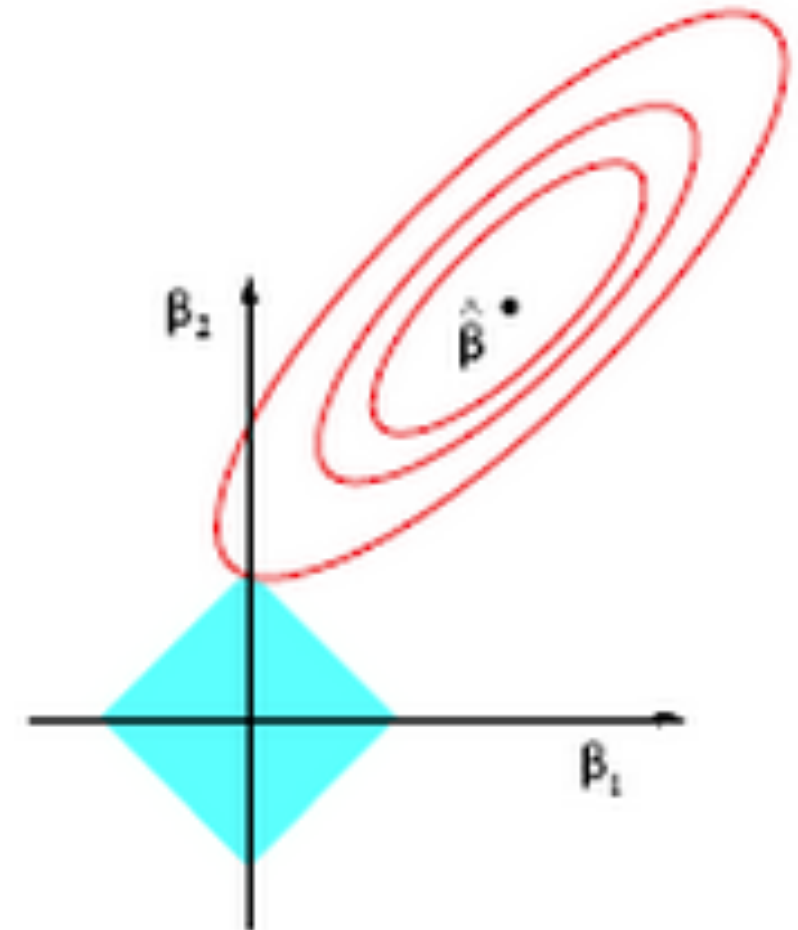
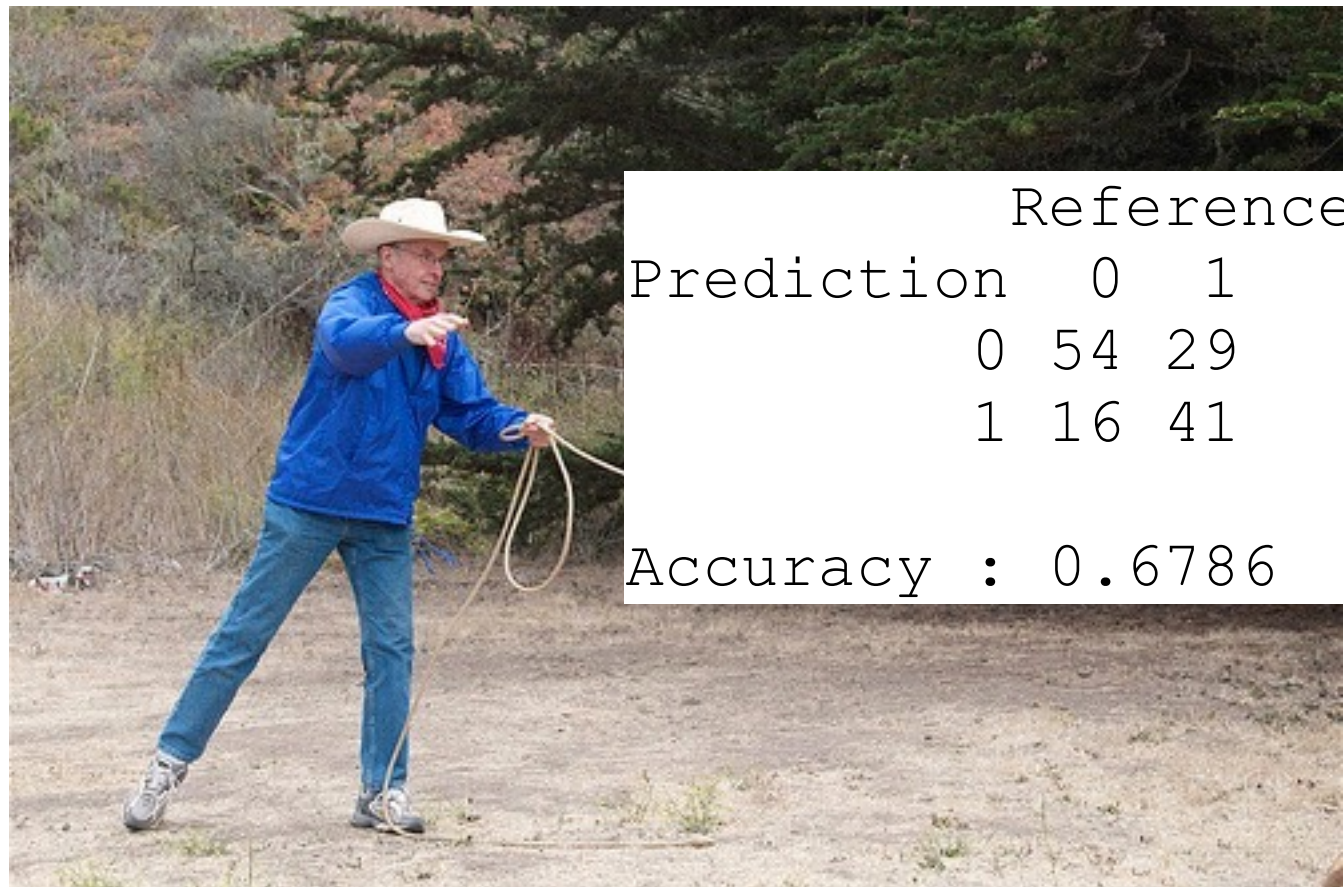
Winning solution: Anton Prokopyev



Motivation

#	Δ1d	Team Name
1		
2		
3		
4		
5		
6	↓2	King's New Clothes 

[1] LASSO Regression



Least Absolute Shrinkage and Selection Operator



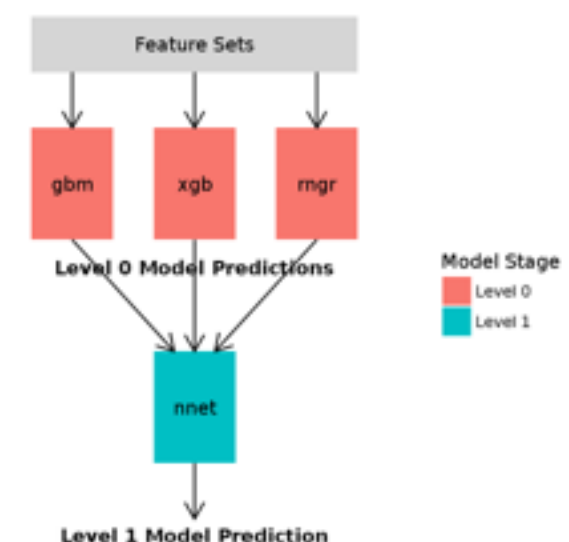
Competing on Kaggle with one model

$Cost(W) = RSS(W) + \lambda * (\text{sum of absolute value of weights})$

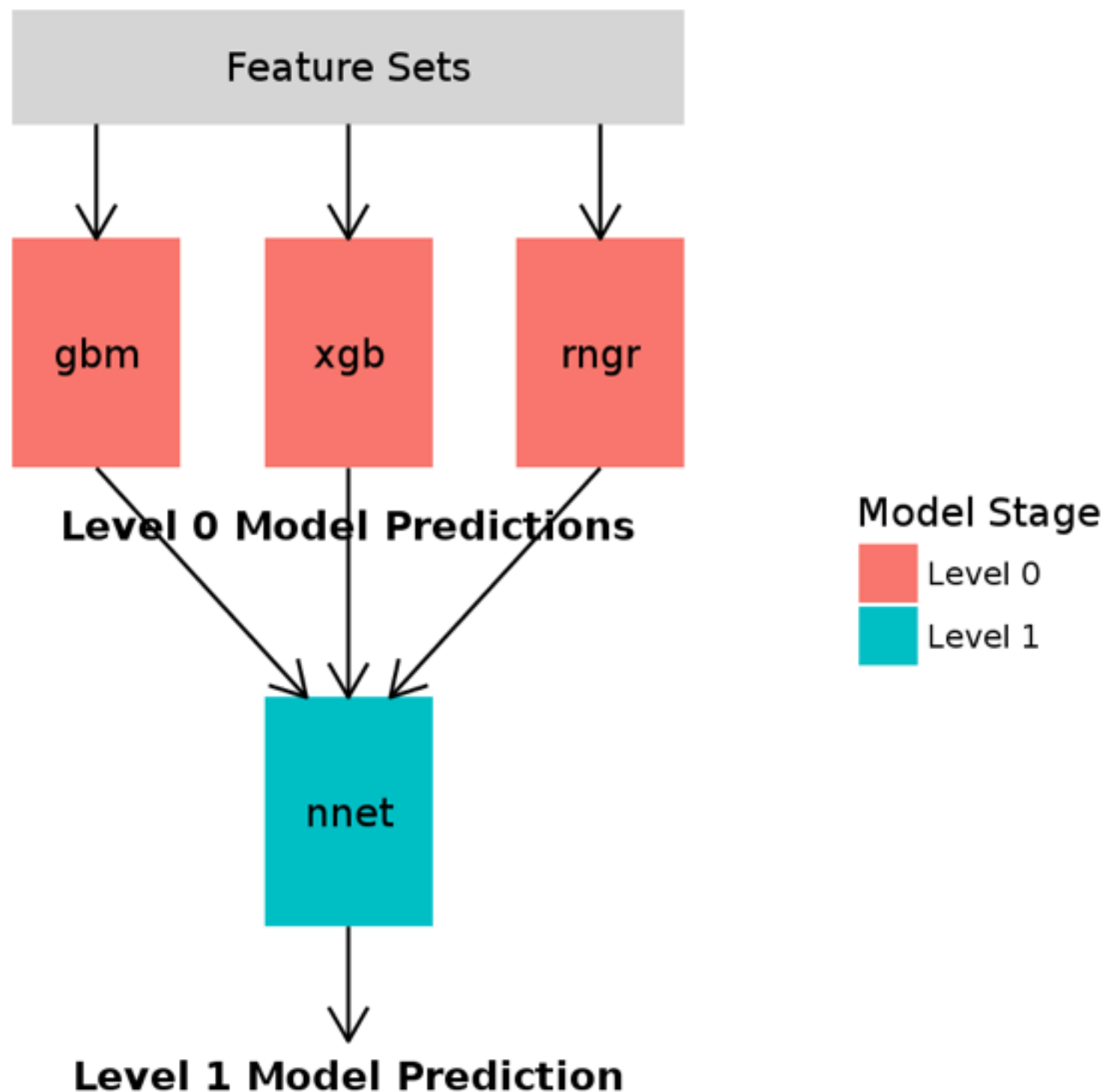
$$= \sum_{i=1}^N \left\{ y_i - \sum_{j=0}^M w_j x_{ij} \right\}^2 + \lambda \sum_{j=0}^M |w_j|$$



Competing on Kaggle with ensemble of models

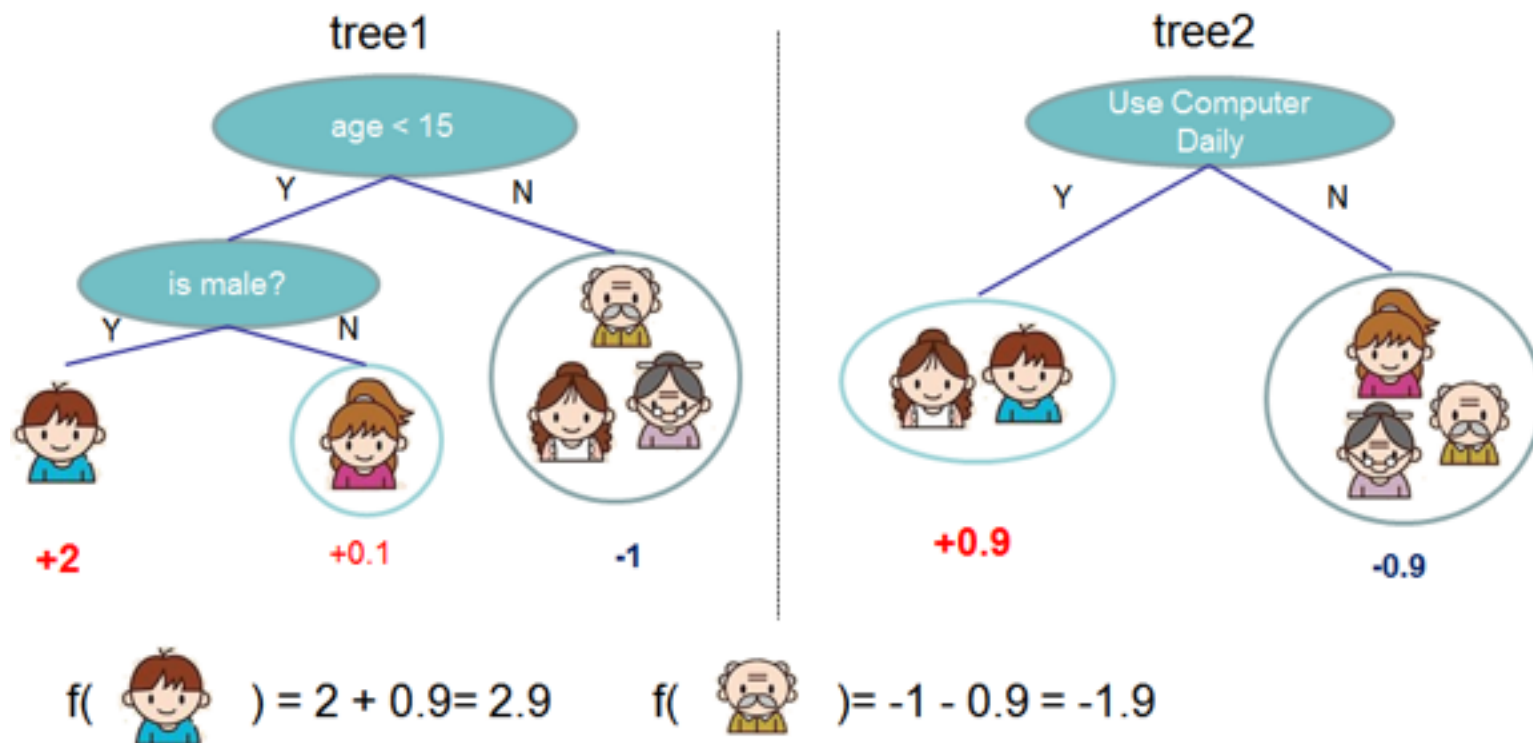


Ensemble modelling

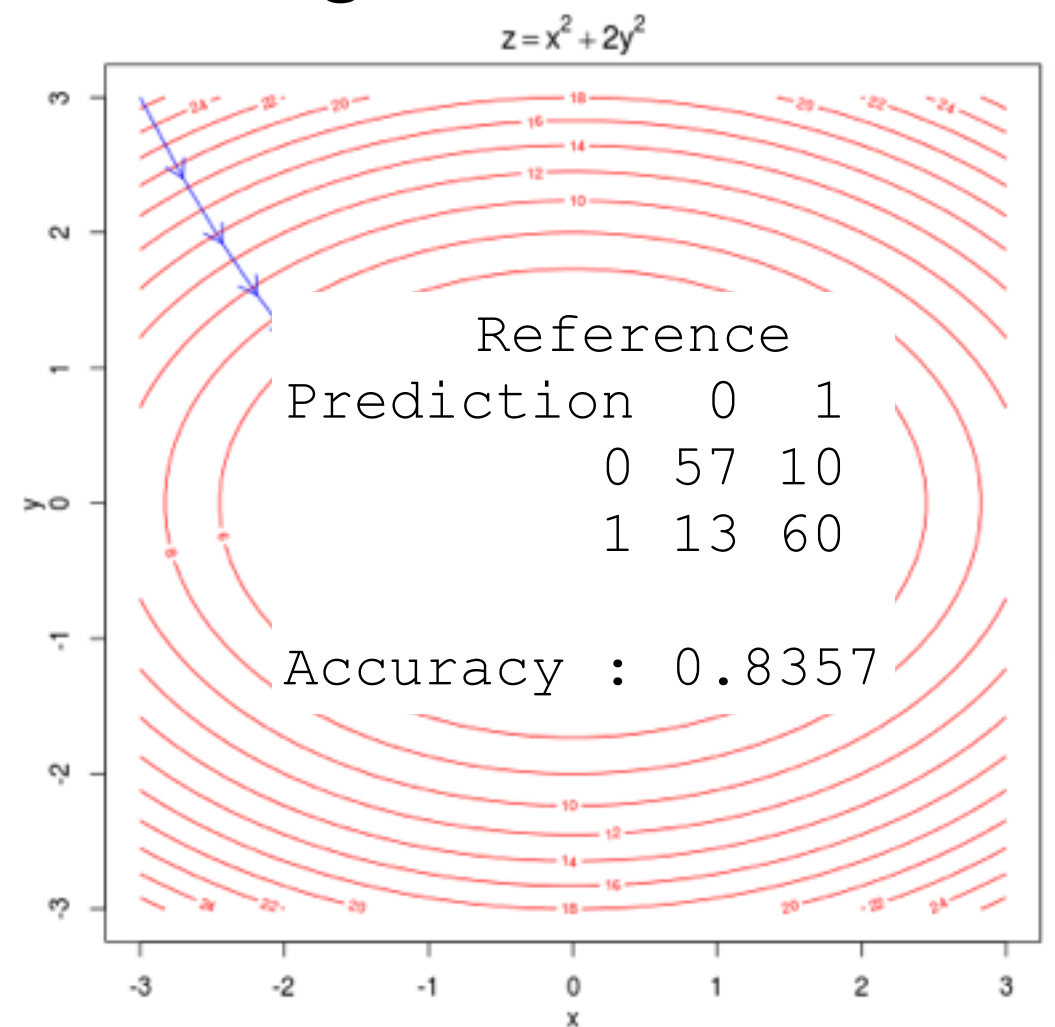


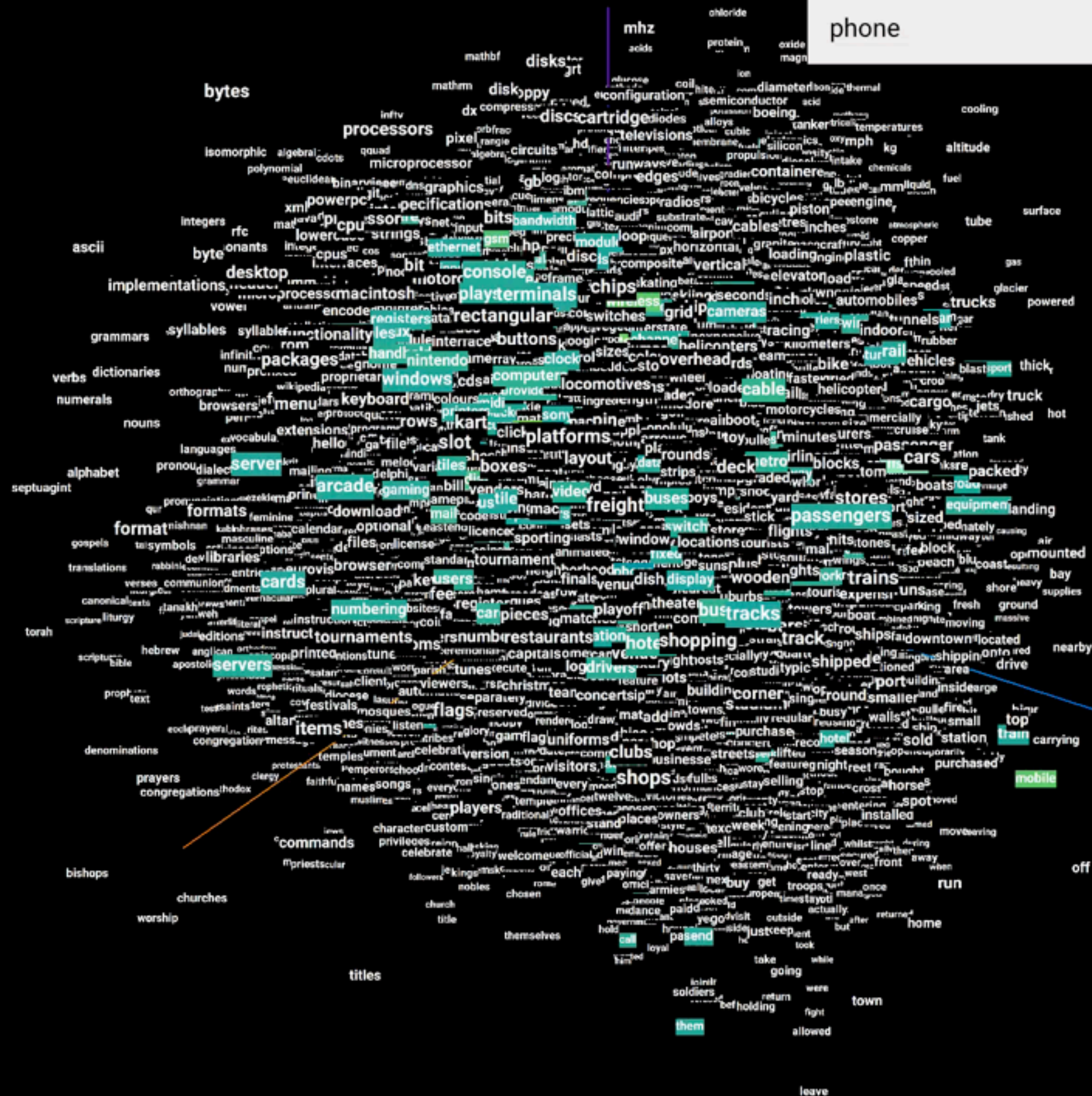
[2] XGBoost

- Combines many weak trees
Does the person like video games?

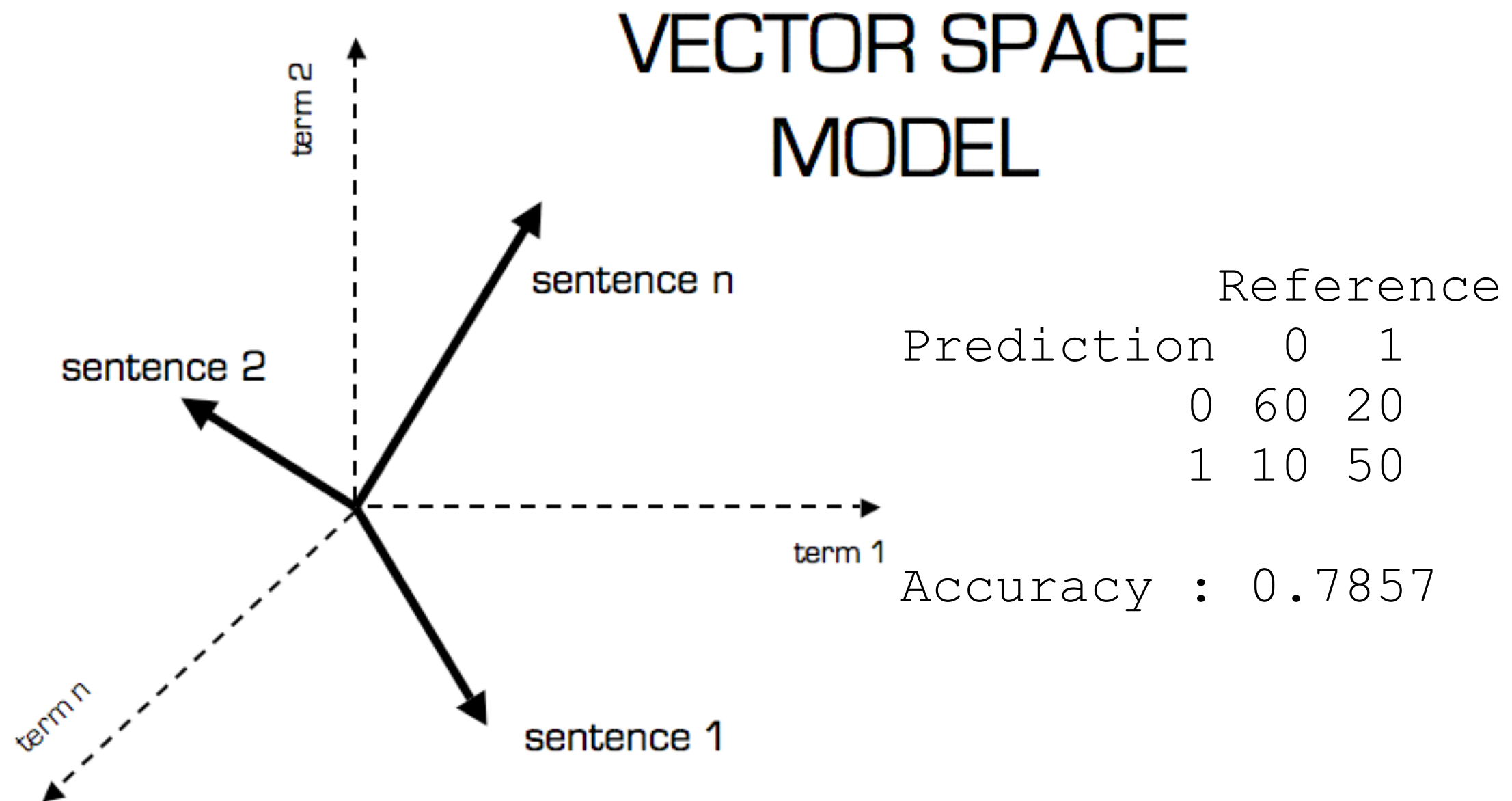


- Uses gradient descent





[3] text2vec



Combined performance

1st Anton Prokopyev

0.80556

Select all of Anton Prokopyev's submissions

Id	File	Description	Public	Private	Submission Date	Selected?
<input type="checkbox"/> 4660472	submit8_stack.csv	stack May 8	0.80556	0.81667	Tue, 09 May 2017 05:11:00	Yes
<input type="checkbox"/> 4657763	submit5.csv	text2vec only	0.80000	0.81667	Mon, 08 May 2017 18:35:19	Yes



I'd like to add you to my professional network on LinkedIn

Connect: [linkedin.com/in/prokopyev](https://www.linkedin.com/in/prokopyev)

Thanks!

UC San Diego. POLI274: Text as Data, M. Roberts

High-dimensional word space: <http://projector.tensorflow.org>

XGBoost for Sentiment Analysis: <https://github.com/wush978/FeatureHashing/blob/master/vignettes/SentimentAnalysis.Rmd>

Text2Vec for Sentiment Analysis: <https://www.r-bloggers.com/twitter-sentiment-analysis-with-machine-learning-in-r-using-doc2vec-approach>