



LEARNING DEEP ON CYBERBULLYING IS ALWAYS BETTER THAN BRUTE FORCE

MICHAL PTASZYNSKI¹

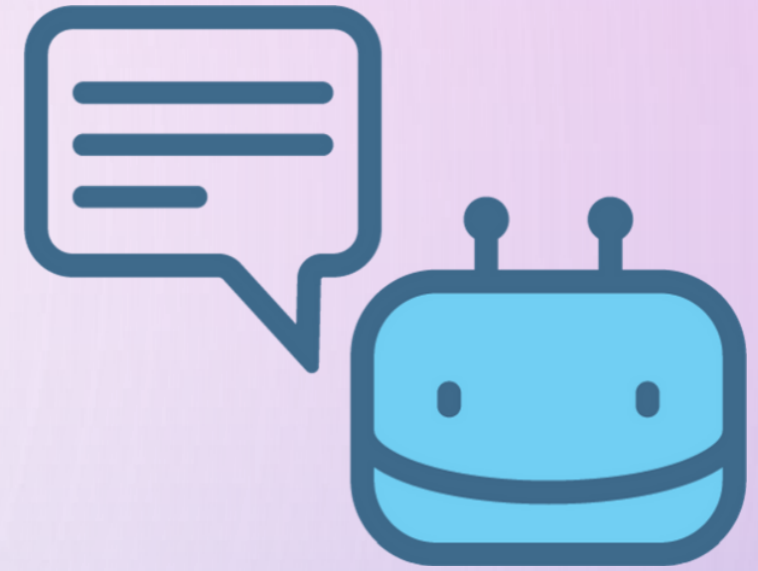
JUUSO KALEVI KRISTIAN ERONEN²

FUMITO MASUI¹

1. KITAMI INSTITUTE OF TECHNOLOGY

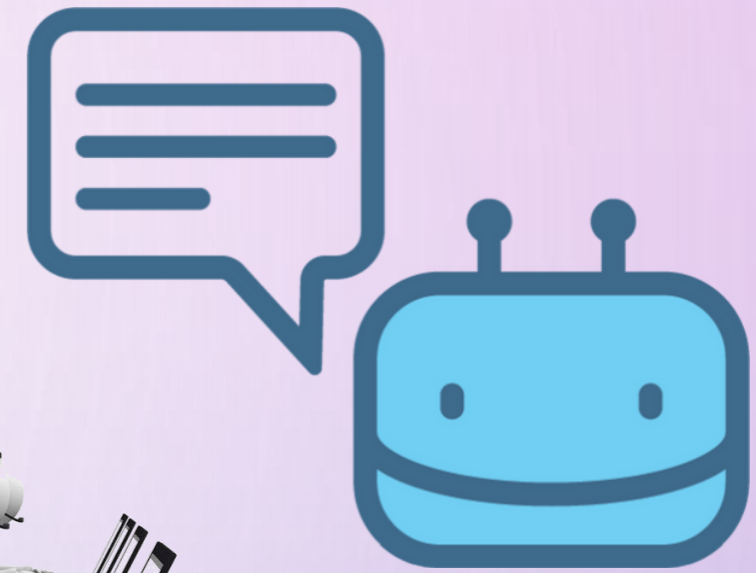
2. TAMPERE UNIVERSITY OF TECHNOLOGY

BACKGROUND



- DIALOG AGENTS APPLICATIONS

BACKGROUND



- DIALOG AGENTS APPLICATIONS
 - CALL CENTERS
 - CUSTOMER SUPPORT
 - CASUAL CONVERSATIONS
 - LANGUAGE EDUCATION
 - ADVERTISEMENT / PROPAGANDA
 - PR / ANTI-PR
 - **FORUM MODERATION**



Trump-AI



Chat bot responds with real quotes by "The Donald"



BACKGROUND

- FORUM MODERATION
 - ONE OF THE PROBLEMS ON FORA:

CYBERBULLYING



BACKGROUND

- CYBERBULLYING

“ USING **TECHNOLOGY** TO RIDICULE OR **HUMILIATE** OTHERS ”

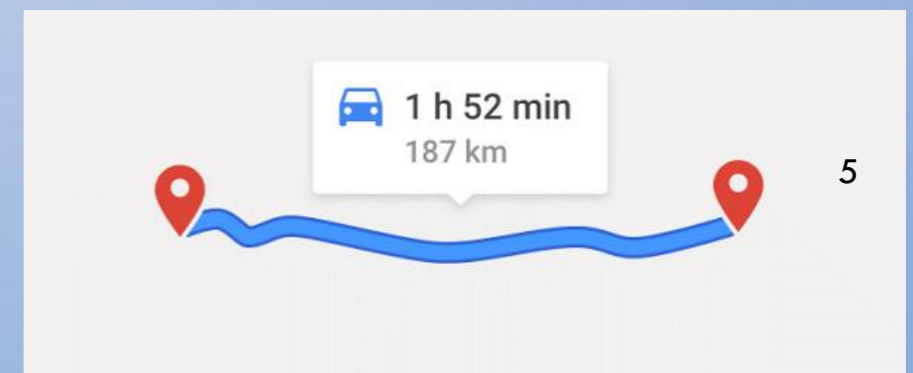
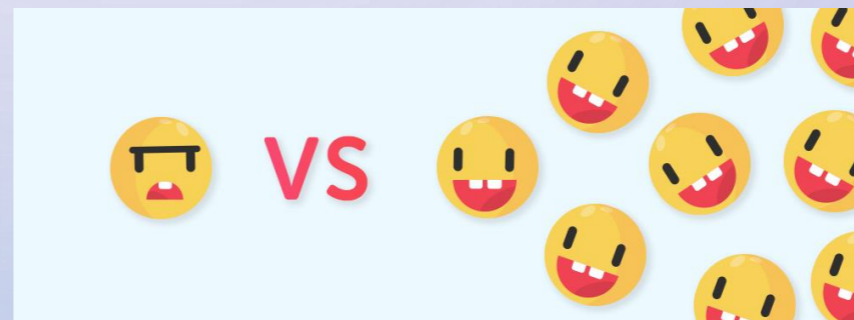
REPETITIVE

IMBALANCE OF POWER

(ONE VS MANY / WEAKER VS STRONGER)

EASY TO DO ON INETERNET

(LONGER PSYCHOLOGICAL DISTANCE)



BACKGROUND

- CYBERBULLYING

CAUSES :

AGGRESSION

ALIENATION

DEPRESSION

SELF-MUTILATION

SUICIDE



BACKGROUND

- CYBERBULLYING

CAUSES :

AGGRESSION

ALIENATION

DEPRESSION

SELF-MUTILATION

SUICIDE

REAL WORLD PROBLEM

**MORE LIFE ON INTERNET =
= MORE CYBERBULLYING**

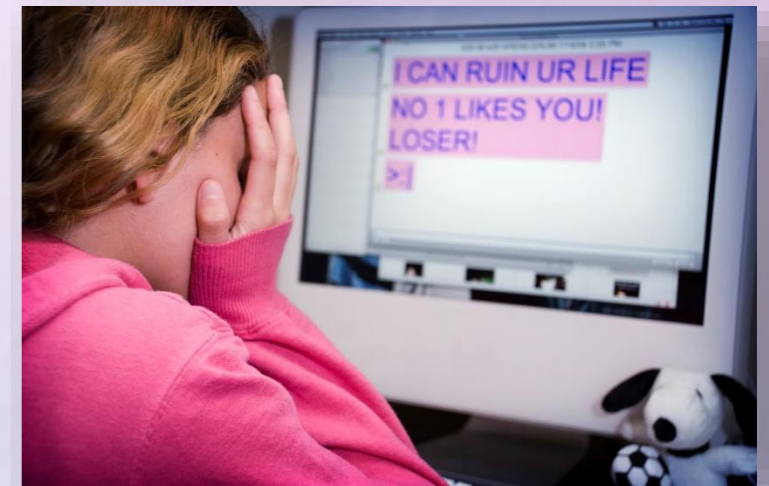


**8% TO EVEN
20% OF USERS
(OFTEN KIDS)**

BACKGROUND

- MANY INTERNET FORA (OVER 1 MIL.¹)
IMPOSSIBLE TO MODERATE EVERYTHING MANUALLY
- ONLINE PATROL (TEACHERS, PARENTS VOLUNTEERS)
 - READ EVERYTHING TO FIND CYBERBULLYING
 - NOT ENOUGH TIME
 - PSYCHOLOGICAL BURDEN

NEED TECHNOLOGY SUPPORT



1) <https://www.quora.com/How-many-online-forums-are-in-existence>

PREVIOUS RESEARCH

PREVIOUS RESEARCH

T. Matsuba, F. Masui, A. Kawai, N. Isu. 2011. **Study on the polarity classification model for the purpose of detecting harmful information on informal school sites** (in Japanese), In *Proceedings of NLP2011*, pp. 388-391.

SO-PMI-IR / phrases

Patent No. 2013-245813. Inventors: Fumito Masui, Michal Ptaszynski, Nitta Taisei. Patent name: *An Apparatus and Method for Detection of Harmful Entries on Internet*

2013
PATENT

M. Ptaszynski, F. Masui, Y. Kimura, R. Rzepka, K. Araki. 2015. **Extracting Patterns of Harmful Expressions for Cyberbullying Detection**, *7th Language & Technology Conference (LTC'15)*, 2015.11.27-29.

Language Combinatorics / Preprocessing

Michal Ptaszynski, F. Masui, Y. Kimura, R. Rzepka, K. Araki. 2015. **Brute Force Works Best Against Bullying**, *IJCAI 2015 Workshop on Intelligent Personalization (IP 2015)*, Buenos Aires, 2015.07.25-31

Language Combinatorics = Brute Force Search

Michal Ptaszynski, P. Dybala, T. Matsuba, F. Masui, R. Rzepka, K. Araki, and Y. Momouchi. 2010. **In the Service of Online Order: Tackling Cyber-Bullying with Machine Learning and Affect Analysis**. *International Journal of Computational Linguistics Research*, Vol. 1, Issue 3, pp. 135-154, 2010.

SVM / optimization

T. Nitta, F. Masui, M. Ptaszynski, Y. Kimura, R. Rzepka, K. Araki. 2013. **Detecting Cyberbullying Entries on Informal School Websites Based on Category Relevance Maximization**. In *Proceedings of IJCNLP 2013*, pp. 579-586.

Category Relevance Optimization

S. Hatakeyama, F. Masui, M. Ptaszynski, K. Yamamoto. 2015. **Improving Performance of Cyberbullying Detection Method with Double Filtered Point-wise Mutual Information**. *ACM Symposium on Cloud Computing 2015 (SoCC'15)*, August 2015.

Automatic acquisition of harmful words

Michal Ptaszynski, P. Dybala, T. Matsuba, F. Masui, R. Rzepka and K. Araki. 2010. **Machine Learning and Affect Analysis Against Cyber-Bullying**. In *Proceedings of AISB'10*, 29th March – 1st April 2010.

Affect analysis of cyberbullying data

2009

2010

2011

2012

2013

2014

2015

PREVIOUS RESEARCH

T. Matsuba, F. Masui, A. Kawai, N. Isu. 2011. **Study on the polarity classification model for the purpose of detecting harmful information on informal school sites** (in Japanese), In *Proceedings of NLP2011*, pp. 388-391.

SO-PMI-IR / phrases

Patent No. 2013-245813. Inventors: Fumito Masui, Michal Ptaszynski, Nitta Taisei. Patent name: *An Apparatus and Method for Detection of Harmful Entries on Internet*

2013
PATENT

M. Ptaszynski, F. Masui, Y. Kimura, R. Rzepka, K. Araki. 2015. **Extracting Patterns of Harmful Expressions for Cyberbullying Detection**, *7th Language & Technology Conference (LTC'15)*, 2015.11.27-29.

Language Combinatorics / Preprocessing

Michal Ptaszynski, F. Masui, Y. Kimura, R. Rzepka, K. Araki. 2015. **Brute Force Works Best Against Bullying**, *IJCAI 2015 Workshop on Intelligent Personalization (IP 2015)*, Buenos Aires, 2015.07.25-31

Language Combinatorics = Brute Force Search

Michal Ptaszynski, P. Dybala, T. Matsuba, F. Masui, R. Rzepka, K. Araki, and Y. Momouchi. 2010. **In the Service of Online Order: Tackling Cyber-Bullying with Machine Learning and Affect Analysis**. *International Journal of Computational Linguistics Research*, Vol. 1, Issue 3, pp. 135-154, 2010.

SVM / optimization

T. Nitta, F. Masui, M. Ptaszynski, Y. Kimura, R. Rzepka, K. Araki. 2013. **Detecting Cyberbullying Entries on Informal School Websites Based on Category Relevance Maximization**. In *Proceedings of IJCNLP 2013*, pp. 579-586.

Category Relevance Optimization

S. Hatakeyama, F. Masui, M. Ptaszynski, K. Yamamoto. 2015. **Improving Performance of Cyberbullying Detection Method with Double Filtered Point-wise Mutual Information**. *ACM Symposium on Cloud Computing 2015 (SoCC'15)*, August 2015.

Automatic acquisition of harmful words

Michal Ptaszynski, P. Dybala, T. Matsuba, F. Masui, R. Rzepka and K. Araki. 2010. **Machine Learning and Affect Analysis Against Cyber-Bullying**. In *Proceedings of AISB'10*, 29th March – 1st April 2010.

Affect analysis of cyberbullying data

2009

2010

2011

2012

2013

2014

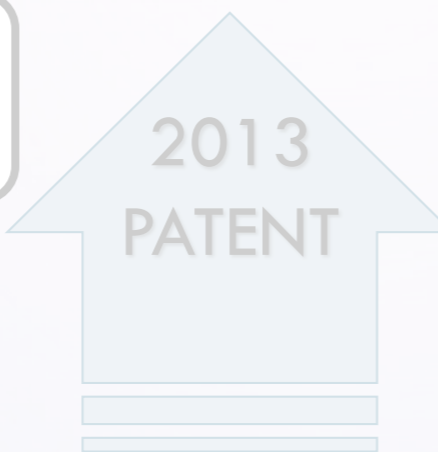
2015

PREVIOUS RESEARCH

T. Matsuba, F. Masui, A. Kawai, N. Isu. 2011. **Study on the polarity classification model for the purpose of detecting harmful information on informal school sites** (in Japanese), In *Proceedings of NLP2011*, pp. 388-391.

SO-PMI-IR / phrases

Patent No. 2013-245813. Inventors: Fumito Masui, Michal Ptaszynski, Nitta Taisei. Patent name: *An Apparatus and Method for Detection of Harmful Entries on Internet*



M. Ptaszynski, F. Masui, Y. Kimura, R. Rzepka, K. Araki. 2015. **Extracting Patterns of Harmful Expressions for Cyberbullying Detection**, *7th Language & Technology Conference (LTC'15)*, 2015.11.27-29.

Language Combinatorics / Preprocessing

Michal Ptaszynski, F. Masui, Y. Kimura, R. Rzepka, K. Araki. 2015. **Brute Force Works Best Against Bullying**, *IJCAI 2015 Workshop on Intelligent Personalization (IP 2015)*, Buenos Aires, 2015.07.25-31

Language Combinatorics = Brute Force Search

Michal Ptaszynski, P. Dybala, T. Matsuba, F. Masui, R. Rzepka, K. Araki, and Y. Momouchi. 2010. **In the Service of Online Order: Tackling Cyber-Bullying with Machine Learning and Affect Analysis**. *International Journal of Computational Linguistics Research*, Vol. 1, Issue 3, pp. 135-154, 2010.

SVM / optimization

T. Nitta, F. Masui, M. Ptaszynski, Y. Kimura, R. Rzepka, K. Araki. 2013. **Detecting Cyberbullying Entries on Informal School Websites Based on Category Relevance Maximization**. In *Proceedings of IJCNLP 2013*, pp. 579-586.

S. Hatakeyama, F. Masui, M. Ptaszynski, K. Yamamoto. 2015. **Improving Performance of Cyberbullying Detection Method with Double Filtered Point-wise Mutual Information**. *ACM Symposium on Cloud Computing 2015 (SoCC'15)*, August 2015.

Michal Ptaszynski, P. Dybala, T. Matsuba, F. Masui, R. Rzepka and K. Araki. 2010. **Machine Learning and Affect Analysis Against Cyber-Bullying**. In *Proceedings of AISB'10*, 29th March – 1st April 2010.

Affect analysis of cyberbullying data

Category Relevance Optimization

Automatic acquisition of harmful words

2009

2010

2011

2012

2013

2014

2015

PREVIOUS RESEARCH

T. Matsuba, F. Masui, A. Kawai, N. Isu. 2011. **Study on the polarity classification model for the purpose of detecting harmful information on informal school sites** (in Japanese), In *Proceedings of NLP2011*, pp. 388-391.

SO-PMI-IR / phrases

Patent No. 2013-245813. Inventors: Fumito Masui, Michal Ptaszynski, Nitta Taisei. Patent name: *An Apparatus and Method for Detection of Harmful Entries on Internet*

M. Ptaszynski, F. Masui, Y. Kimura, R. Rzepka, K. Araki. 2015. **Extracting Patterns of Harmful Expressions for Cyberbullying Detection**, *7th Language & Technology Conference (LTC'15)*, 2015.11.27-29.

Language Combinatorics / Preprocessing

2013 PATENT

Michal Ptaszynski, F. Masui, Y. Kimura, R. Rzepka, K. Araki. 2015. **Brute Force Works Best Against Bullying**, *IJCAI 2015 Workshop on Intelligent Personalization (IP 2015)*, Buenos Aires, 2015.07.25-31

Language Combinatorics = Brute Force Search

Michal Ptaszynski, P. Dybala, T. Matsuba, F. Masui, R. Rzepka, K. Araki, and Y. Momouchi. 2010. **In the Service of Online Safety: Tackling Cyber-Bullying with Machine Learning and Sentiment Analysis**. *International Journal of Computational Linguistics Research*, Vol. 1, Issue 3, pp. 135-154, 2010.

SVM / optimization

M. Ptaszynski, F. Masui, M. Ptaszynski, Y. Kimura, R. Rzepka, K. Araki. 2013. **Detecting Cyberbullying Entries on Informal Websites Based on Category Relevance Maximization**. In *Proceedings of IJCNLP 2013*, pp. 579-586.

S. Hatakeyama, F. Masui, M. Ptaszynski, K. Yamamoto. 2015. **Improving Performance of Cyberbullying Detection Method with Double Filtered Point-wise Mutual Information**. *ACM Symposium on Cloud Computing 2015 (SoCC'15)*, August 2015.

Michal Ptaszynski, P. Dybala, T. Matsuba, F. Masui, R. Rzepka and K. Araki. 2010. **Machine Learning and Affect Analysis Against Cyber-Bullying**. In *Proceedings of AISB'10*, 29th March – 1st April 2010.

Affect analysis of cyberbullying data

Category Relevance Maximization

Automatic acquisition of harmful words

2009

2010

2011

2012

2013

2014

2015

PREVIOUS RESEARCH

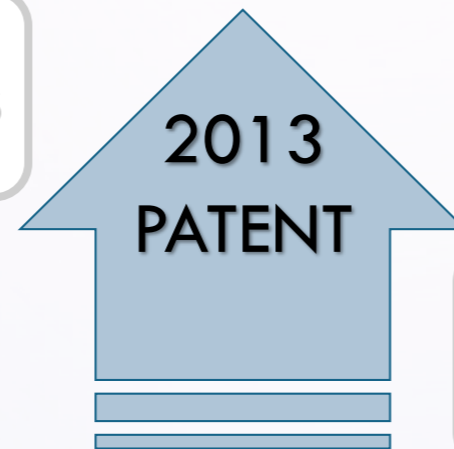
T. Matsuba, F. Masui, A. Kawai, N. Isu. 2011. **Study on the polarity classification model for the purpose of detecting harmful information on informal school sites (in Japanese)**, In *Proceedings of NLP2011*, pp. 388-391.

SO-PMI-IR / phrases

Patent No. 2013-245813. Inventors: Fumito Masui, Michal Ptaszynski, Nitta Taisei. Patent name: *An Apparatus and Method for Detection of Harmful Entries on Internet*

M. Ptaszynski, F. Masui, Y. Kimura, R. Rzepka, K. Araki. 2015. **Extracting Patterns of Harmful Expressions for Cyberbullying Detection**, *7th Language & Technology Conference (LTC'15)*, 2015.11.27-29.

Language Combinatorics / Preprocessing



Michal Ptaszynski, F. Masui, Y. Kimura, R. Rzepka, K. Araki. 2015. **Brute Force Works Best Against Bullying**, *IJCAI 2015 Workshop on Intelligent Personalization (IP 2015)*, Buenos Aires, 2015.07.25-31

Language Combinatorics = Brute Force Search

Michal Ptaszynski, P. Dybala, T. Matsuba, F. Masui, R. Rzepka, K. Araki, and Y. Momouchi. 2010. **In the Service of Online Order: Tackling Cyber-Bullying with Machine Learning and Affect Analysis**. *International Journal of Computational Linguistics Research*, Vol. 1, Issue 3, pp. 135-154, 2010.

SVM / optimization

T. Nitta, F. Masui, M. Ptaszynski, Y. Kimura, R. Rzepka, K. Araki. 2013. **Detecting Cyberbullying Entries on Informal School Websites Based on Category Relevance Maximization**. In *Proceedings of IJCNLP 2013*, pp. 579-586.

S. Hatakeyama, F. Masui, M. Ptaszynski, K. Yamamoto. 2015. **Improving Performance of Cyberbullying Detection Method with Double Filtered Point-wise Mutual Information**. *ACM Symposium on Cloud Computing 2015 (SoCC'15)*, August 2015.

Category Relevance Optimization

Automatic acquisition of harmful words

Michal Ptaszynski, P. Dybala, T. Matsuba, F. Masui, R. Rzepka and K. Araki. 2010. **Machine Learning and Affect Analysis Against Cyber-Bullying**. In *Proceedings of AISB'10*, 29th March – 1st April 2010.

Affect analysis of cyberbullying data

2009

2010

2011

2012

2013

2014

2015

PREVIOUS RESEARCH

T. Matsuba, F. Masui, A. Kawai, N. Isu. 2011. **Study on the polarity classification model for the purpose of detecting harmful information on informal school sites** (in Japanese), In *Proceedings of NLP2011*, pp. 388-391.

SO-PMI-IR / phrases

Patent No. 2013-245813. Inventors: Fumito Masui, Michal Ptaszynski, Nitta Taisei. Patent name: *An Apparatus and Method for Detection of Harmful Entries on Internet*

2013
PATENT

M. Ptaszynski, F. Masui, Y. Kimura, R. Rzepka, K. Araki. 2015. **Extracting Patterns of Harmful Expressions for Cyberbullying Detection**, *7th Language & Technology Conference (LTC'15)*, 2015.11.27-29.

Language Combinatorics / Preprocessing

Michal Ptaszynski, F. Masui, Y. Kimura, R. Rzepka, K. Araki. 2015. **Brute Force Works Best Against Bullying**, *IJCAI 2015 Workshop on Intelligent Personalization (IP 2015)*, Buenos Aires, 2015.07.25-31

Language Combinatorics = Brute Force Search

Michal Ptaszynski, P. Dybala, T. Matsuba, F. Masui, R. Rzepka, K. Araki, and Y. Momouchi. 2010. **In the Service of Online Order: Tackling Cyber-Bullying with Machine Learning and Affect Analysis**. *International Journal of Computational Linguistics Research*, Vol. 1, Issue 3, pp. 135-154, 2010.

SVM / optimization

T. Nitta, F. Masui, M. Ptaszynski, Y. Kimura, R. Rzepka, K. Araki. 2013. **Detecting Cyberbullying Entries on Informal School Websites Based on Category Relevance Maximization**. In *Proceedings of IJCNLP 2013*, pp. 579-586.

Category Relevance Optimization

S. Hatakeyama, F. Masui, Y. Kimura, R. Rzepka, K. Araki. 2015. **Improving Performance of Cyberbullying Detection with Double Filtered Word Embeddings**. *ACM Symposium on Cloud Computing (SoCC'15)*, August 2015.

Automatic Acquisition of harmful words

Michal Ptaszynski, P. Dybala, T. Matsuba, F. Masui, R. Rzepka and K. Araki. 2010. **Machine Learning and Affect Analysis Against Cyber-Bullying**. In *Proceedings of AISB'10*, 29th March – 1st April 2010.

Affect analysis of cyberbullying data

Michal Ptaszynski, K. Yamamoto. 2015. **Improving Cyberbullying Detection Method with Mutual Information**. *ACM Symposium on Cloud Computing (SoCC'15)*, August 2015.

2009

2010

2011

2012

2013

2014

2015

PREVIOUS RESEARCH

T. Matsuba, F. Masui, A. Kawai, N. Isu. 2011. **Study on the polarity classification model for the purpose of detecting harmful information on informal school sites** (in Japanese), In *Proceedings of NLP2011*, pp. 388-391.

SO-PMI-IR / phrases

Patent No. 2013-245813. Inventors: Fumito Masui, Michal Ptaszynski, Nitta Taisei. Patent name: *An Apparatus and Method for Detection of Harmful Entries on Internet*

2013
PATENT

M. Ptaszynski, F. Masui, Y. Kimura, R. Rzepka, K. Araki. 2015. **Extracting Patterns of Harmful Expressions for Cyberbullying Detection**, *7th Language & Technology Conference (LTC'15)*, 2015.11.27-29.

Language Combinatorics / Preprocessing

Michal Ptaszynski, F. Masui, Y. Kimura, R. Rzepka, K. Araki. 2015. **Brute Force Works Best Against Bullying**, *IJCAI 2015 Workshop on Intelligent Personalization (IP 2015)*, Buenos Aires, 2015.07.25-31

Language Combinatorics = Brute Force Search

Michal Ptaszynski, P. Dybala, T. Matsuba, F. Masui, R. Rzepka, K. Araki, and Y. Momouchi. 2010. **In the Service of Online Order: Tackling Cyber-Bullying with Machine Learning and Affect Analysis**. *International Journal of Computational Linguistics Research*, Vol. 1, Issue 3, pp. 135-154, 2010.

SVM / optimization

T. Nitta, F. Masui, M. Ptaszynski, Y. Kimura, R. Rzepka, K. Araki. 2013. **Detecting Cyberbullying Entries on Informal School Websites Based on Category Relevance Maximization**. In *Proceedings of IJCNLP 2013*, pp. 579-586.

S. Hatakeyama, F. Masui, M. Ptaszynski, K. Yamamoto. 2015. **Improving Performance of Cyberbullying Detection Method with Double Filtered Point-wise Mutual Information**. *ACM Symposium on Cloud Computing 2015 (SoCC'15)*, August 2015.

Michal Ptaszynski, P. Dybala, T. Matsuba, F. Masui, R. Rzepka and K. Araki. 2010. **Machine Learning and Affect Analysis Against Cyber-Bullying**. In *Proceedings of AISB'10*, 29th March – 1st April 2010.

Affect analysis of cyberbullying data

Category Relevance Optimization

Automatic acquisition of harmful words

2009

2010

2011

2012

2013

2014

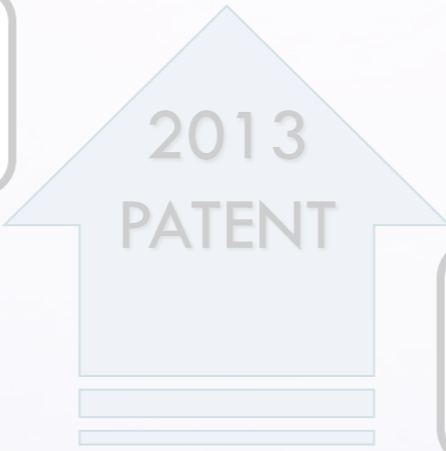
2015

PREVIOUS RESEARCH

T. Matsuba, F. Masui, A. Kawai, N. Isu. 2011. **Study on the polarity classification model for the purpose of detecting harmful information on informal school sites** (in Japanese), In *Proceedings of NLP2011*, pp. 388-391.

SO-PMI-IR / phrases

Patent No. 2013-245813. Inventors: Fumito Masui, Michal Ptaszynski, Nitta Taisei. Patent name: *An Apparatus and Method for Detection of Harmful Entries on Internet*



M. Ptaszynski, F. Masui, Y. Kimura, R. Rzepka, K. Araki. 2015. **Extracting Patterns of Harmful Expressions for Cyberbullying Detection**, *7th Language & Technology Conference (LTC'15)*, 2015.11.27-29.

Language Combinatorics / Preprocessing

Michal Ptaszynski, F. Masui, Y. Kimura, R. Rzepka, K. Araki. 2015. **Brute Force Works Best Against Bullying**, *IJCAI 2015 Intelligent Personalization (IP 2015)*, Buenos Aires, 2015.

Language Combinatorics = Brute Force Search

Michal Ptaszynski, P. Dybala, T. Matsuba, F. Masui, R. Rzepka, K. Araki, and Y. Momouchi. 2010. **In the Service of Online Order: Tackling Cyber-Bullying with Machine Learning and Affect Analysis**. *International Journal of Computational Linguistics Research*, Vol. 1, Issue 3, pp. 135-154, 2010.

SVM / optimization

T. Nitta, F. Masui, M. Ptaszynski, Y. Kimura, R. Rzepka, K. Araki. 2013. **Detecting Cyberbullying Entries on Informal School Websites Based on Category Relevance Maximization**. In *Proceedings of IJCNLP 2013*, pp. 579-586.

S. Hatakeyama, F. Masui, M. Ptaszynski, K. Araki, Y. Kimura. 2015. **Improving Performance of Cyberbullying Detection with Double Filtered Point-wise Mutual Information**. *ACM Symposium on Cloud Computing 2015 (SoCC)*, August 2015.

Michal Ptaszynski, P. Dybala, T. Matsuba, F. Masui, R. Rzepka and K. Araki. 2010. **Machine Learning and Affect Analysis Against Cyber-Bullying**. In *Proceedings of AISB'10*, 29th March – 1st April 2010.

Affect analysis of cyberbullying data

Category Relevance Optimization

Automatic acquisition of harmful words



PREVIOUS RESEARCH

T. Matsuba, F. Masui, A. Kawai, N. Isu. 2011. **Study on the polarity classification model for the purpose of detecting harmful information on informal school sites** (in Japanese), In *Proceedings of NLP2011*, pp. 388-391.

SO-PMI-IR / phrases

Patent No. 2013-245813. Inventors: Fumito Masui, Michal Ptaszynski, Nitta Taisei. Patent name: *An Apparatus and Method for Detection of Harmful Entries on Internet*

2013
PATENT

M. Ptaszynski, F. Masui, Y. Kimura, R. Rzepka, K. Araki. 2015. **Extracting Patterns of Harmful Expressions for Cyberbullying Detection**, *7th Language & Technology Conference (LTC'15)*, 2015.11.27-29.

Language Combinatorics / Preprocessing

Michal Ptaszynski, F. Masui, Y. Kimura, R. Rzepka, K. Araki. 2015. **Brute Force Works Best Against Bullying**, *IJCAI 2015 Workshop on Intelligent Personalization (IP 2015)*, Buenos Aires, 2015.07.25-31

Language Combinatorics = Brute Force Search

Michal Ptaszynski, P. Dybala, T. Matsuba, F. Masui, R. Rzepka, K. Araki, and Y. Momouchi. 2010. **In the Service of Online Order: Tackling Cyber-Bullying with Machine Learning and Affect Analysis**. *International Journal of Computational Linguistics Research*, Vol. 1, Issue 3, pp. 135-154, 2010.

SVM / optimization

T. Nitta, F. Masui, M. Ptaszynski, Y. Kimura, R. Rzepka, K. Araki. 2013. **Detecting Cyberbullying Entries on Informal School Websites Based on Category Relevance Maximization**. In *Proceedings of IJCNLP 2013*, pp. 579-586.

S. Hatakeyama, F. Masui, M. Ptaszynski, K. Yamamoto. 2015. **Improving Performance of Cyberbullying Detection Method with Double Filtered Point-wise Mutual Information**. *ACM Symposium on Cloud Computing 2015 (SoCC'15)*, August 2015.

Michal Ptaszynski, P. Dybala, T. Matsuba, F. Masui, R. Rzepka and K. Araki. 2010. **Machine Learning and Affect Analysis Against Cyber-Bullying**. In *Proceedings of AISB'10*, 29th March – 1st April 2010.

Affect analysis of cyberbullying data

Category Relevance Optimization

Automatic acquisition of harmful words

2009 2010 2011 2012 2013 2014 2015

PREVIOUS RESEARCH

T. Matsuba, F. Masui, A. Kawai, N. Isu. 2011. Study on the polarity classification model for the purpose of detecting harmful information on informal school sites (in Japanese), In *Proceedings of NLP2011*, pp. 383-390.

Patent No. 2013-245813. Inventors: Fumito Masui, Michal Ptaszynski, Nitta Taisei.

M. Ptaszynski, F. Masui, Y. Kimura, R. Rzepka, K. Araki. 2015. Extracting Patterns of Harmful Expressions for Cyberbullying Detection, *7th Language & Technology Conference (LTC'15)*, 2015.11.27-29.

Feature sophistication

simple →

→ sophisticated

DEEP LEARNING

syntactic pat.

word patterns

phrases

bag-of-words

words

2009

2010

2011

2012

2013

2014

2015

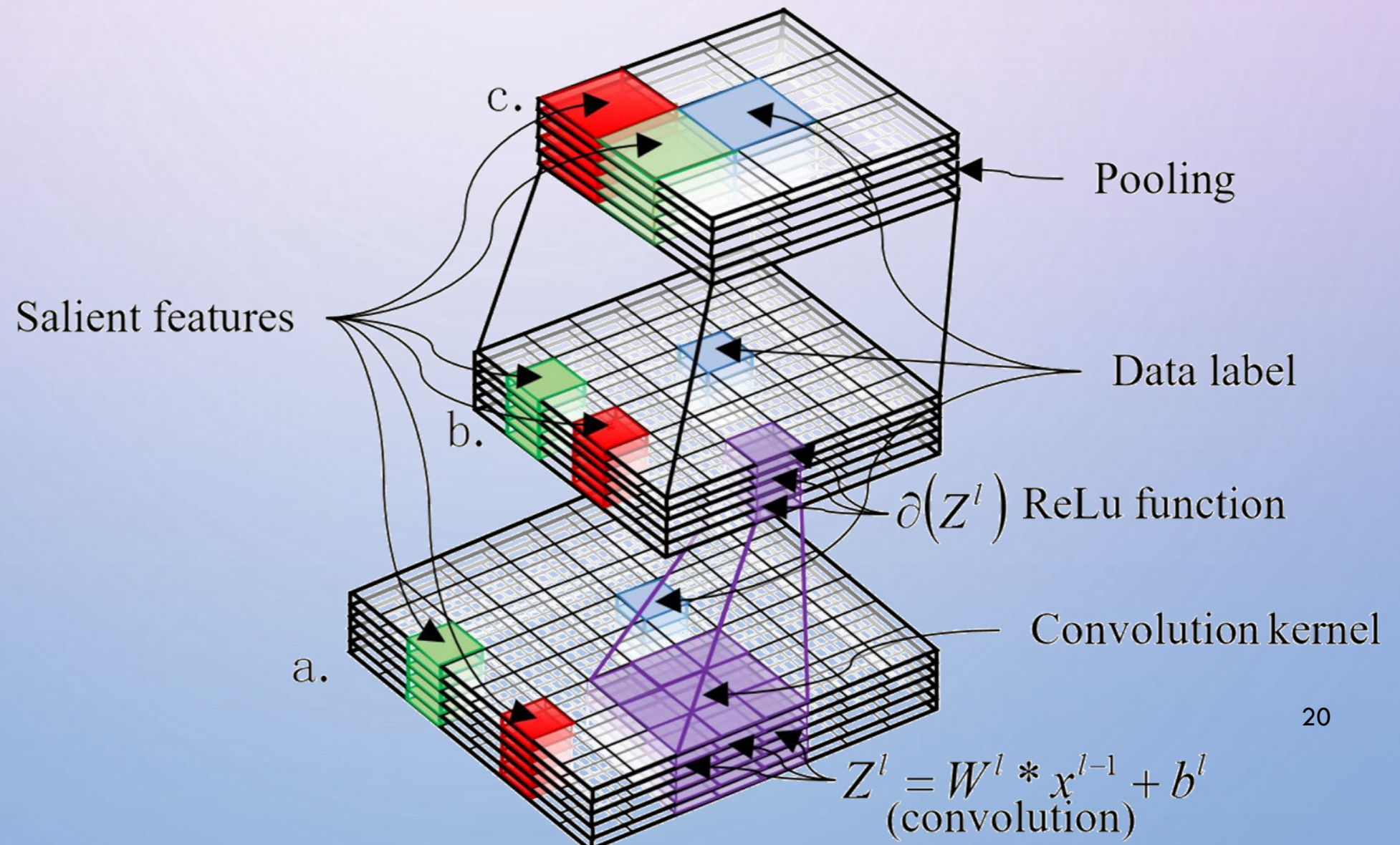
2017

Now!

19

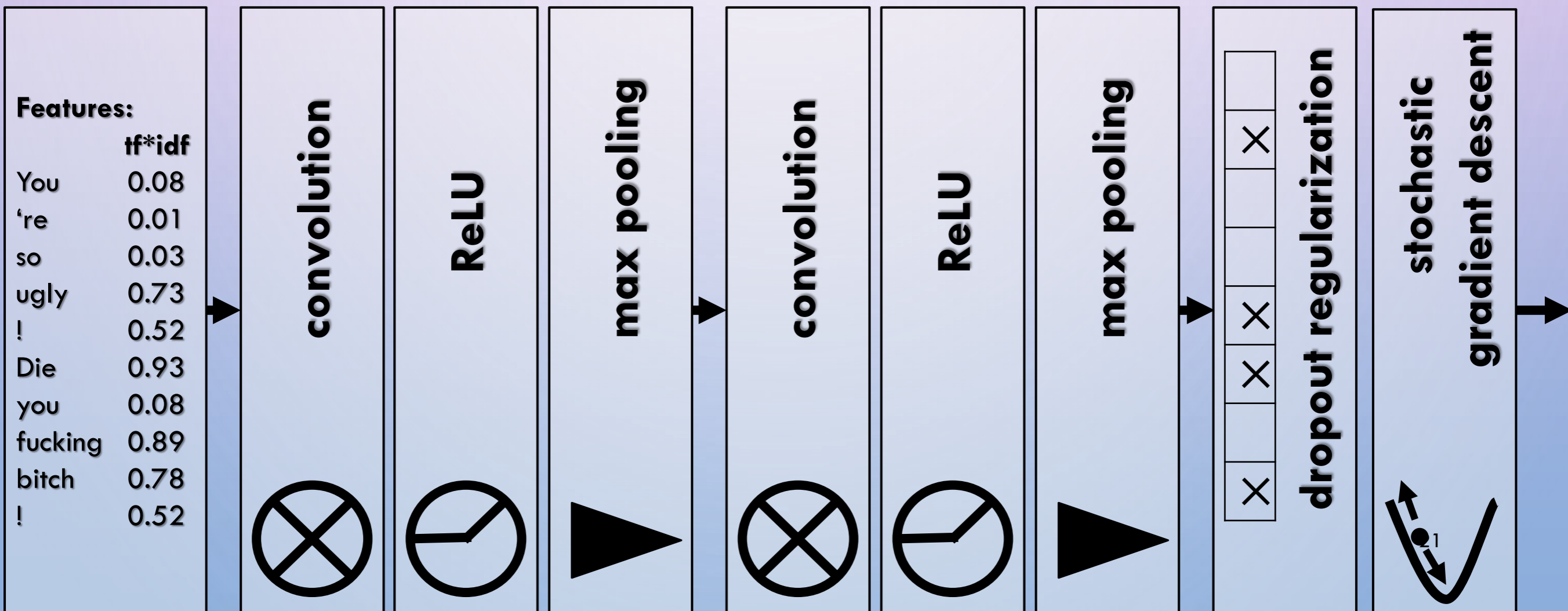
PROPOSED METHOD

- DEEP CONVOLUTIONAL NEURAL NETWORK



PROPOSED METHOD

- DEEP CONVOLUTIONAL NEURAL NETWORK



*) dummy weights only for explanation.

PROPOSED METHOD

- CONVOLUTION**

Average of weights in batch for each feature

Features:	tf*idf
You	0.08
're	0.01
so	0.03
ugly	0.73
!	0.52
Die	0.93
you	0.08
fucking	0.89
bitch	0.78
!	0.52

Batch = 5x5

You	're	so	ugly	!
Die	you	f*ng	b*ch	!
...				

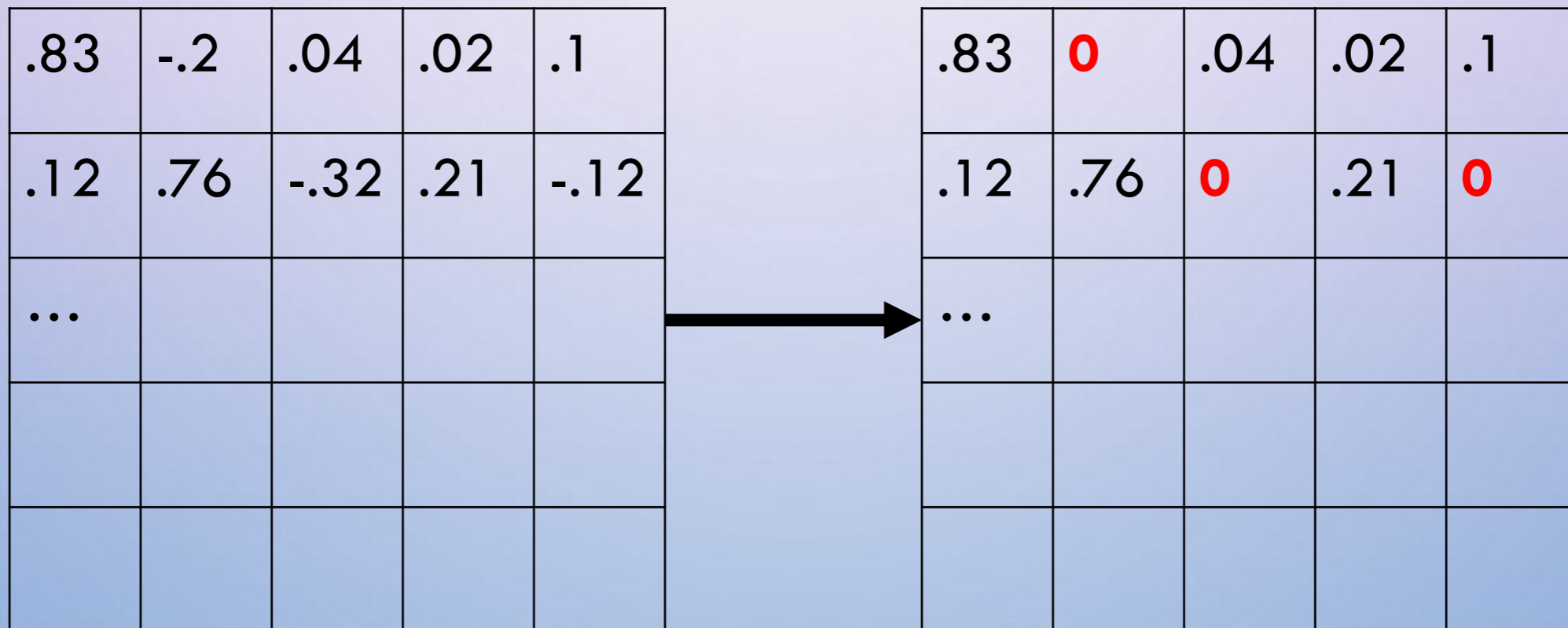
Slide batch window



.08	.01	.03	.73	.52
.93	.08	.89	.78	.52
...				

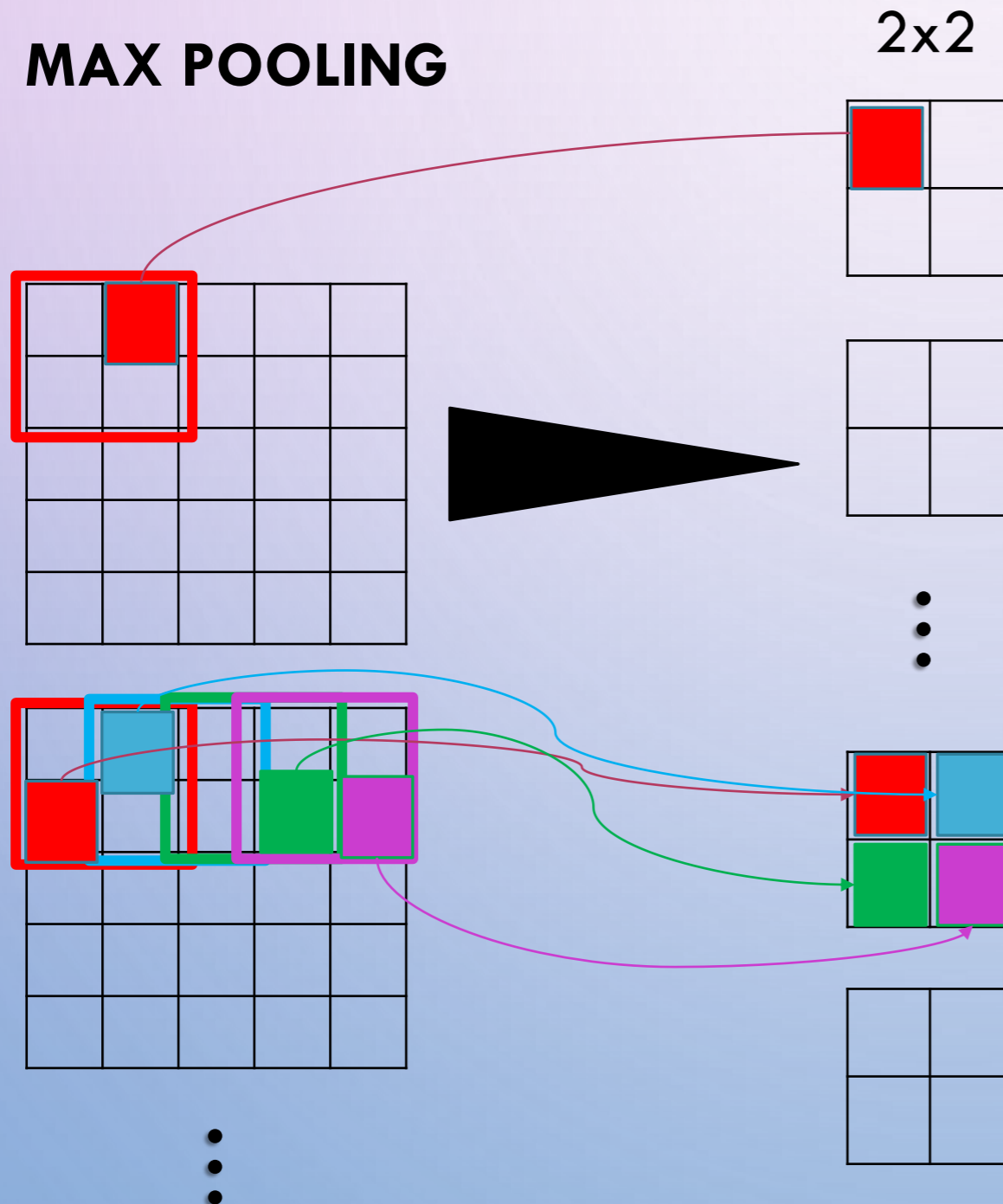
PROPOSED METHOD

- **RECTIFIED LINEAR UNITS (RELU)**



PROPOSED METHOD

- **MAX POOLING**

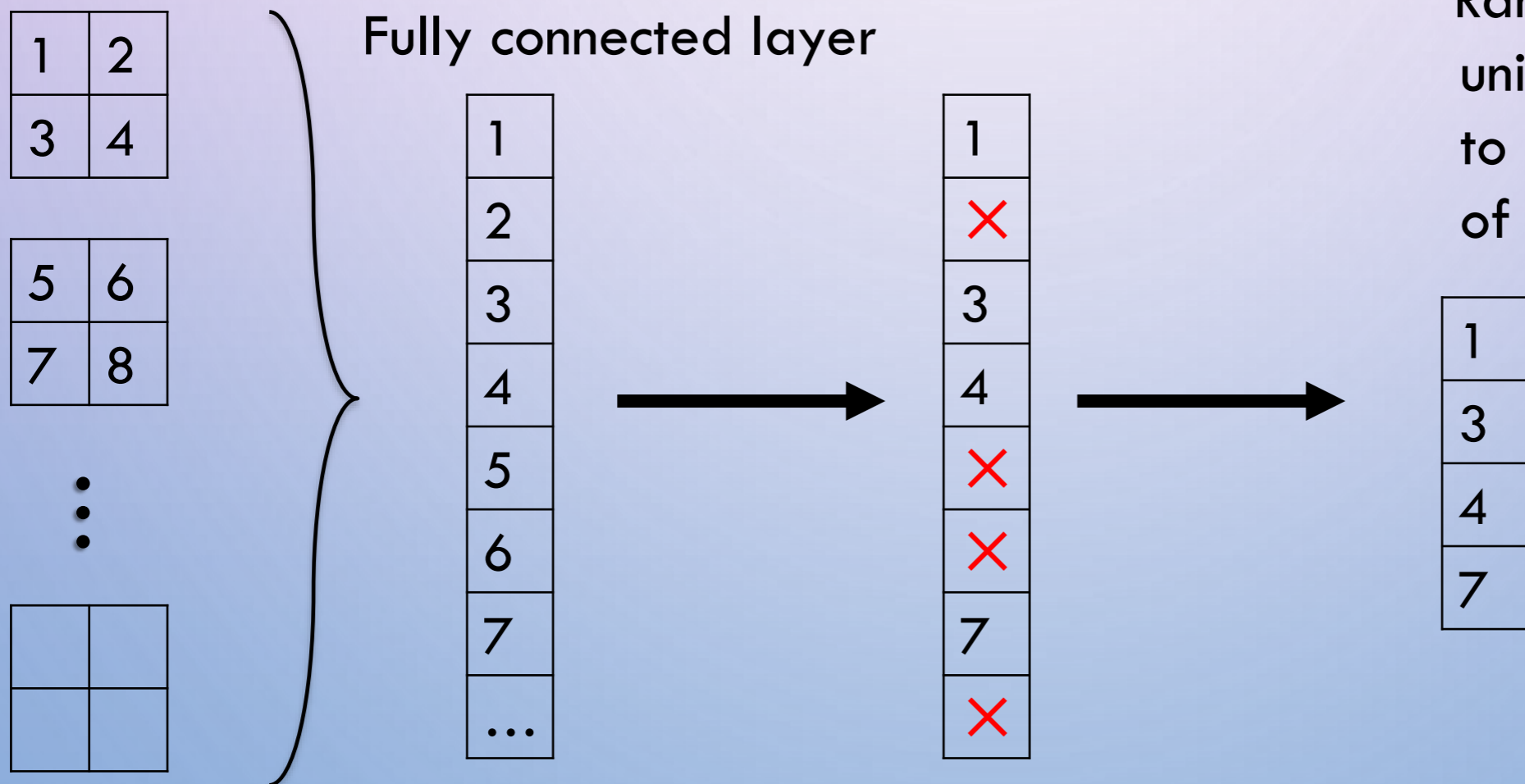


2x2

* Reduce dimensionality and correct over-fitting

PROPOSED METHOD

- **DROPOUT REGULARIZATION**

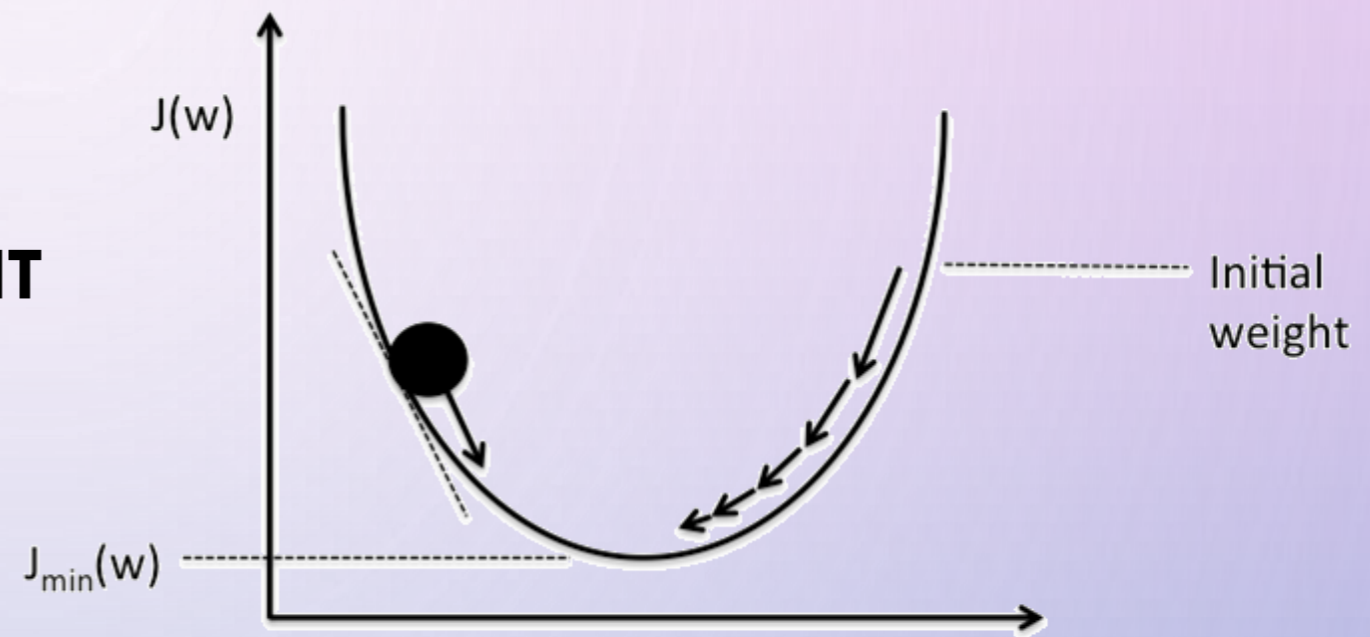
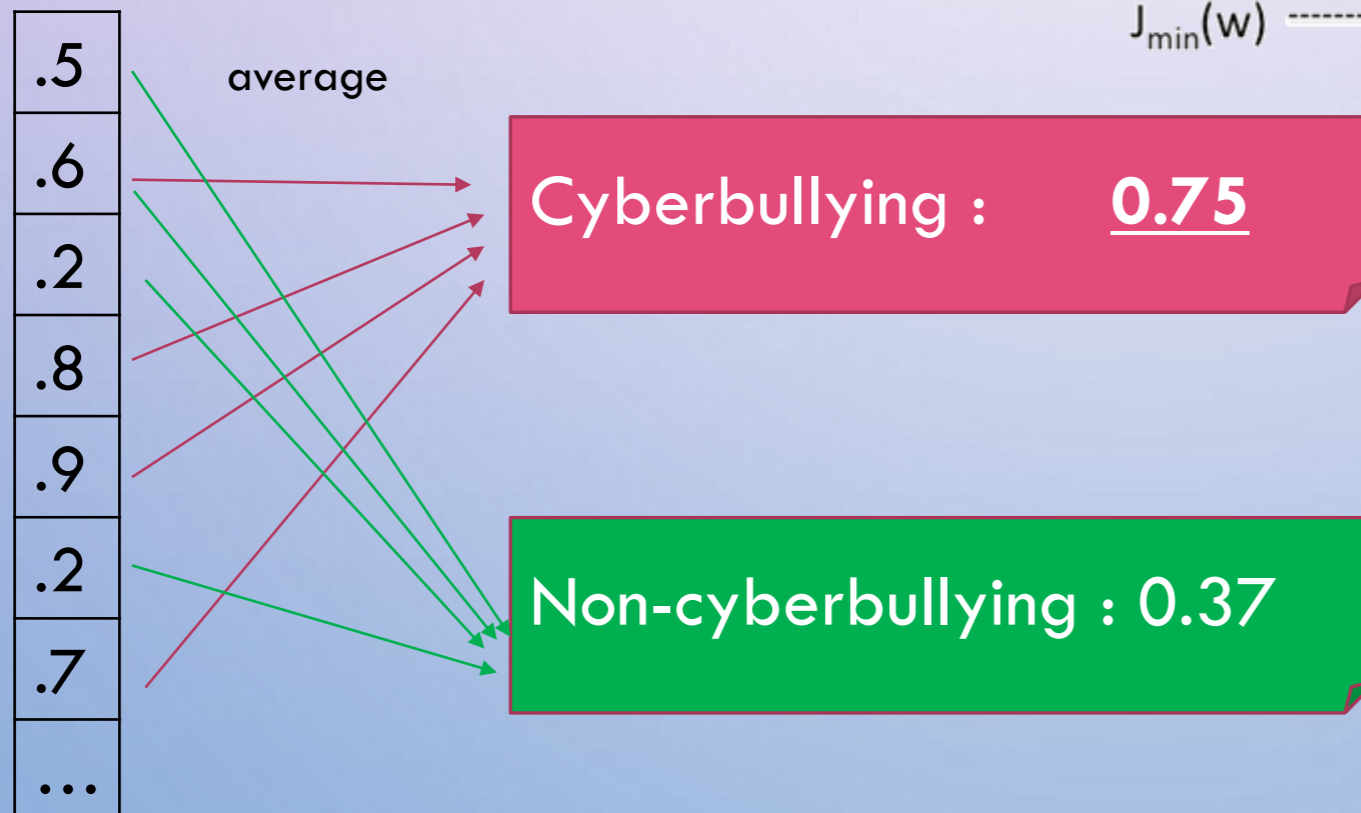


* Randomly delete some units during training to prevent co-adaptation of hidden units

PROPOSED METHOD

- STOCHASTIC GRADIENT DESCENT**

Fully connected layer



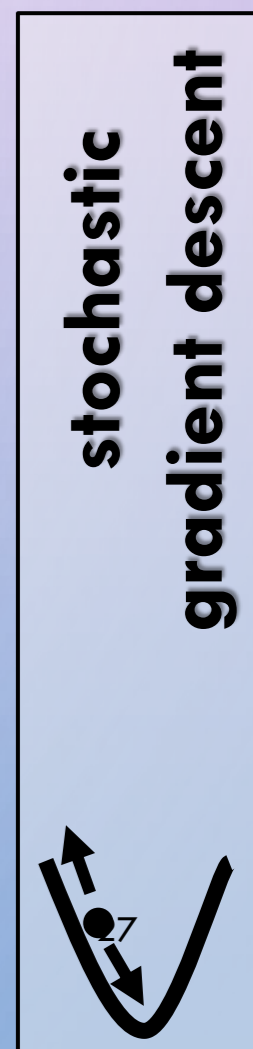
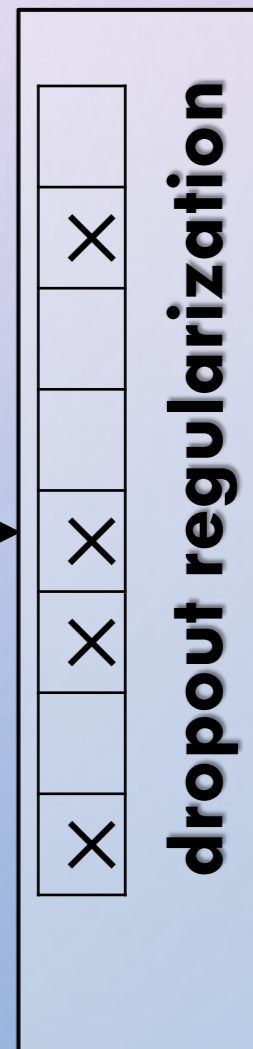
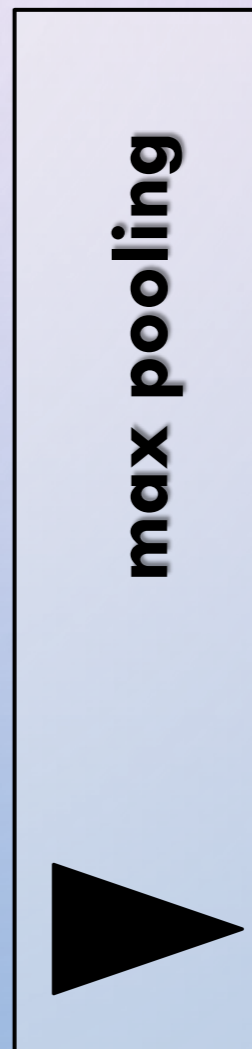
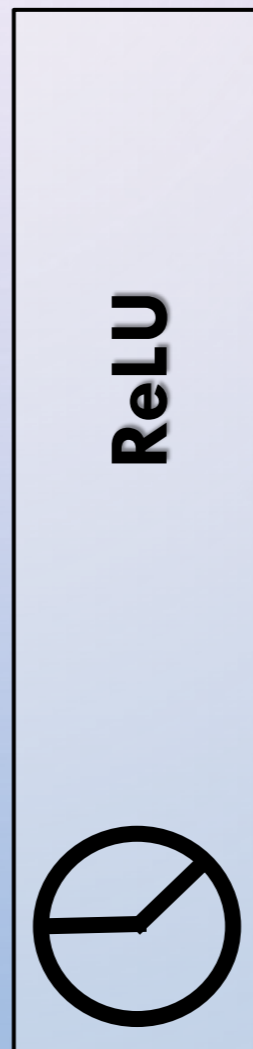
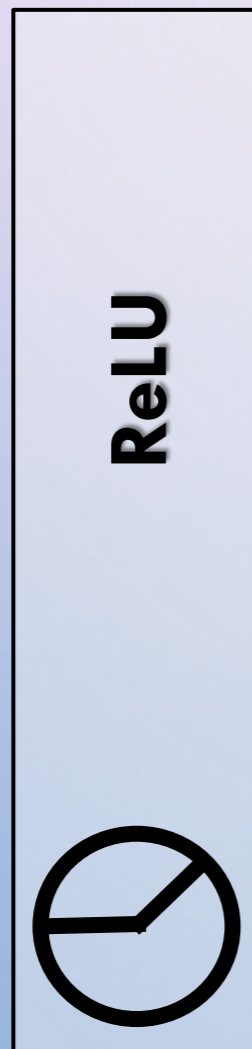
	correct	answer	Error
CB	1	.75	.25
N-CB	0	.37	.63
	Total error =		0.88

* Slide weight a little back and forth to²⁶ minimize error

PROPOSED METHOD

- DEEP CONVOLUTIONAL NEURAL NETWORK

Features:	tf*idf
You	0.08
're	0.01
so	0.03
ugly	0.73
!	0.52
Die	0.93
you	0.08
fucking	0.89
bitch	0.78
!	0.52



PROPOSED METHOD

- LANGAUGE COMBINATORICS
- DEEP CONVOLUTIONAL NEURAL NETWORK
- SHALLOW CNN
- SVM
 - LINEAR
 - POLYNOMIAL
 - RADIAL BASED FUNCTION
 - SIGMOID
- DECISION TREES (J48)
- RANDOM FORESTS
- KNN
- NAÏVE BAYES
- JRIP

[Ptaszynski et al., 2010]
[Sood et al., 2012]
[Dinakar et al., 2012]
[Sarna and Bhatia, 2015]
[Ptaszynski et al., 2015a,b]
...AND A FEW OTHERS

PROPOSED METHOD

- LANGAUGE COMBINATORICS → BRUTE FORCE
- DEEP CONVOLUTIONAL NEURAL NETWORK
- SHALLOW CNN → Neural Nets
- SVM
 - LINEAR
 - POLYNOMIAL → SVM
 - RADIAL BASED FUNCTION
 - SIGMOID
- DECISION TREES (J48) → trees
- RANDOM FORESTS
- KNN
- NAÏVE BAYES → lazy / rules
- JRIP

10-fold
X-validation

DATASET

- ACTUAL DATA COLLECTED BY INTERNET PATROL (ANNOTATED BY EXPERTS)
- FROM UNOFFICIAL SCHOOL FORUMS (BBS)
- PROVIDED BY HUMAN RIGHT CENTER IN JAPAN (MIE PREFECTURE)
- ACCORDING TO THE DEFINITION BY JAPANESE MINISTRY OF EDUCATION (MEXT)
- 1,490 HARMFUL AND 1,508 NON-HARMFUL ENTRIES

FEATURE SELECTION

Example: John McDonald killed Mary Poppins!

- **TOKENIZATION:** JOHN MCDONNARD KILLED MARY POPPINS !
- **LEMMATIZATION:** JOHN MCDONNARD KILL MARY POPPINS !
- **PARTS OF SPEECH:** NOUN NOUN VERB NOUN NOUN EXCL.
- **TOKENS WITH POS:** JOHN_NOUN MCDONNARD_NOUN KILLED_VERB MARY_NOUN POPPINS_NOUN !_EXCL.
- **LEMMAS WITH POS:** JOHN_NOUN MCDONNARD_NOUN KILL_VERB MARY_NOUN POPPINS_NOUN !_EXCL.
- **TOKENS WITH NAMED ENTITY RECOGNITION:** JOHN [COMPANY] KILLED MARY POPPINS !
- **LEMMAS WITH NER:** JOHN [COMPANY] KILL MARY POPPINS !
- **CHUNKING:** JOHN_MCDONNARD_KILLED MARY_POPPINS_!
- **DEPENDENCY STRUCTURE:** 1-JOHN 1-MCDONNARD 2-KILLED 3-MARY 3-POPPINS 3-!
- **CHUNKING WITH NER:** JOHN_[COMPANY]_KILLED MARY_POPPINS_!
- **DEPENDENCY STRUCTURE WITH NAMED ENTITIES:** 1-JOHN 1-[COMPANY] 2-KILLED 3-MARY 3-POPPINS³¹ 3-!

RESULTS AND DISCUSSION

RESULTS AND DISCUSSION

- LAZY CLASSIFIERS :
 - GENERALLY POOR PERFORMANCE (F1 ~ 50% - 60%)
 - BETTER WITH NER (F1 ~ 70%)

RESULTS AND DISCUSSION

- LAZY CLASSIFIERS :
 - GENERALLY POOR PERFORMANCE (F1 ~ 50% - 60%)
 - BETTER WITH NER (F1 ~ 70%)
- TREE-BASED :
 - J48 – LOW
 - RANDOM FOREST – PRETTY GOOD (F1 ~ 80%), BUT TIME INEFFICIENT

RESULTS AND DISCUSSION

- LAZY CLASSIFIERS :
 - GENERALLY POOR PERFORMANCE (F1 ~ 50% - 60%)
 - BETTER WITH NER (F1 ~ 70%)
- TREE-BASED :
 - J48 – LOW
 - RANDOM FORREST – PRETTY GOOD (F1 ~ 80%), BUT TIME INEFFICIENT
- SVM
 - LINEAR – NICE (UP TO F1=82.5%)
 - OTHER – POOR
 - SUPER FAST TO TRAIN (BEST TIME TO PERFORMANCE RATIO)

RESULTS AND DISCUSSION

- BEST FEATURE SETS:
 - TOKENS + NER
 - LEMMAS + NER

RESULTS AND DISCUSSION

- BEST FEATURE SETS:

- TOKENS + NER
- LEMMAS + NER

** CYBERBULLYING IS OFTEN ABOUT REVEALING PRIVATE INFORMATION, NOT ONLY ABOUT SLANDERING **

RESULTS AND DISCUSSION

- BEST FEATURE SETS:

- TOKENS + NER
- LEMMAS + NER

** CYBERBULLYING IS OFTEN ABOUT REVEALING PRIVATE INFORMATION, NOT ONLY ABOUT SLANDERING **

- 2ND BEST METHOD (EXCEPT PROPOSED)

- BRUTE FORCE (F1 = 80.3%)
- EXCEPT ONE SVM CASE
 - LINEAR SVM ON LEMMAS (F1 = 82.5%)



RESULTS AND DISCUSSION

- BEST METHOD
 - CNN WITH 2HIDDEN LAYERS (PROPOSED)
 - $F1 = 93.5\%$
 - NER ALWAYS HELPED



RESULTS AND DISCUSSION

- BEST METHOD
 - CNN WITH 2HIDDEN LAYERS (PROPOSED)
 - F1=93.5%
 - NER ALWAYS HELPED



...WHY?

RESULTS AND DISCUSSION

- GREATEST PROBLEM WITH NEURAL NETS
(AND ALSO MANY OTHER MACHINE LEARNING METHODS)

RESULTS AND DISCUSSION

- GREATEST PROBLEM WITH NEURAL NETS
(AND ALSO MANY OTHER MACHINE LEARNING METHODS)
- INTERPRETABILITY
 - WHY RESULTS WERE AS GOOD?
 - WHAT EXACTLY WAS SO GOOD ABOUT IT?
 - WHAT INFLUENCED THE RESULTS?

RESULTS AND DISCUSSION

- GREATEST PROBLEM WITH NEURAL NETS
(AND ALSO MANY OTHER MACHINE LEARNING METHODS)

- INTERPRETABILITY

- WHY RESULTS WERE AS GOOD?
- WHAT EXACTLY WAS SO GOOD ABOUT IT?
- WHAT INFLUENCED THE RESULTS?



**IF YOU DON'T KNOW
THE CAUSE YOU CAN
ALWAYS LOOK FOR
CORRELATION**

RESULTS AND DISCUSSION

1. NAMED ENTITY RECOGNITION USUALLY HELPED ESPECIALLY WITH CNN
2. ...?

→ GENERAL LOOK AT THE DATA

→ WHAT IS DIFFERENT ABOUT THE DATA?

RESULTS AND DISCUSSION

- LEXICAL DENSITY [URE, 1971]

ALL UNIQUE WORDS IN CORPUS / ALL WORDS IN CORPUS

RESULTS AND DISCUSSION

- LEXICAL DENSITY [URE, 1971]

ALL UNIQUE WORDS IN CORPUS / ALL WORDS IN CORPUS

NOT JUST WORDS:

TOKENS, LEMMAS, POS, TOKEN-POS, TOKEN-NER, LEMMA-POS, LEMMA-NER,
CHUNKS, CHUNKS-NER, DEPENDENCY, DEP-NER,

RESULTS AND DISCUSSION

- ~~LEXICAL DENSITY [URE, 1971]~~ → **FEATURE DENSITY**

ALL UNIQUE WORDS IN CORPUS / ALL WORDS IN CORPUS

NOT JUST WORDS:

TOKENS, LEMMAS, POS, TOKEN-POS, TOKEN-NER, LEMMA-POS, LEMMA-NER,
CHUNKS, CHUNKS-NER, DEPENDENCY, DEP-NER,

FEATURE DENSITY

- CALCULATE FD FOR ALL DATASETS
- CHECK CORRELATION BETWEEN FD AND CLASSIFIER RESULTS

FEATURE DENSITY

Classifier	ρ value	2-sided p-value
CNN-2L	0.638	*p=0.035
SVM-pol	-0.431	p=0.185
SVM-sig	-0.534	p=0.091
SPEC(BEP)	-0.550	p=0.133
RF	-0.560	p=0.073
SVM-lin	-0.564	p=0.076
SPEC(F1)	-0.636	p=0.066
SVM-rad	-0.639	*p=0.034
CNN-1L	-0.709	*p=0.019
JRip	-0.729	*p=0.011
NB	-0.736	*p=0.013
J48	-0.791	**p=0.006
kNN	-0.809	**p=0.004

FEATURE DENSITY

Classifier	ρ value	2-sided p-value
CNN-2L	0.638	*p=0.035
SVM-pol	-0.431	p=0.185
SVM-sig	-0.534	p=0.091
SPEC(BEP)	-0.550	p=0.133
RF	-0.560	p=0.073
SVM-lin	-0.564	p=0.076
SPEC(F1)	-0.636	p=0.066
SVM-rad	-0.639	*p=0.034
CNN-1L	-0.709	*p=0.019
JRip	-0.729	*p=0.011
NB	-0.736	*p=0.013
J48	-0.791	**p=0.006
kNN	-0.809	**p=0.004

**DENSIER
DATA KILLS
CLASSIFIER**

POORLY PERFORMING
CLASSIFIERS:
NEGATIVE STRONG₅₀
CORRELATION WITH FD

FEATURE DENSITY

Classifier	ρ value	2-sided p-value
CNN-2L	0.638	*p=0.035
SVM-pol	-0.431	p=0.185
SVM-sig	-0.534	p=0.091
SPEC(BEP)	-0.550	p=0.133
RF	-0.560	p=0.073
SVM-lin	-0.564	p=0.076
SPEC(F1)	-0.636	p=0.066
SVM-rad	-0.639	*p=0.034
CNN-1L	-0.709	*p=0.019
JRip	-0.729	*p=0.011
NB	-0.736	*p=0.013
J48	-0.791	**p=0.006
kNN	-0.809	**p=0.004

PROBABLY*
DENSIER DATA
HARMS
CLASSIFIER
 (*low significance)



MODERATELY PERFORMING
 CLASSIFIERS:
 NEGATIVE MEDIUM
 CORRELATION WITH FD

FEATURE DENSITY

Classifier	ρ value	2-sided p-value
CNN-2L	0.638	*p=0.035
SVM-pol	-0.431	p=0.185
SVM-sig	-0.534	p=0.091
SPEC(BEP)	-0.550	p=0.133
RF	-0.560	p=0.073
SVM-lin	-0.564	p=0.076
SPEC(F1)	-0.636	p=0.066
SVM-rad	-0.639	*p=0.034
CNN-1L	-0.709	*p=0.019
JRip	-0.729	*p=0.011
NB	-0.736	*p=0.013
J48	-0.791	**p=0.006
kNN	-0.809	**p=0.004



CNN – BEST PERFORMANCE
STRONG POSITIVE
CORRELATION WITH FD

**DENSIER
DATA =
BETTER
RESULTS**

FEATURE DENSITY

Classifier	ρ value	2-sided p-value
CNN-2L	0.638	*p=0.035
SVM-pol	-0.431	p=0.185
SVM-sig	-0.534	p=0.091
SPEC(BEP)	-0.550	p=0.133
RF	-0.560	p=0.073
SVM-lin	-0.564	p=0.076
SPEC(F1)	-0.636	p=0.066
SVM-rad	-0.639	*p=0.034
CNN-1L	-0.709	*p=0.019
JRip	-0.729	*p=0.011
NB	-0.736	*p=0.013
J48	-0.791	**p=0.006
kNN	-0.809	**p=0.004



CNN – BEST PERFORMANCE
STRONG POSITIVE
CORRELATION WITH FD



**FUTURE:
LET'S CHECK
EVEN
DENSIER DATA**

CONCLUSIONS

- PROBLEM: CYBERBULLYING DETECTION
- MANY FEATURE SETS
- MANY CLASSIFIERS
- PROPOSED DEEP CNN SOLUTION
- NAMED ENTITY ANNOTATION USUALLY HELPED DETECT CYBERBULLYING
- FEATURE DENSITY POSITIVELY CORELATED WITH CNN RESULTS
 - WILL CHECK EVEN DENSIER FEATURE SETS
 - WILL CHECK FOR OTHER TASKS : SENTIMENT, DECEPTION, SARCASM, ETC.



THANK YOU FOR YOUR KIND
ATTENTION!

MICHAL PTASZYNSKI

PTASZYNSKI@IEEE.ORG