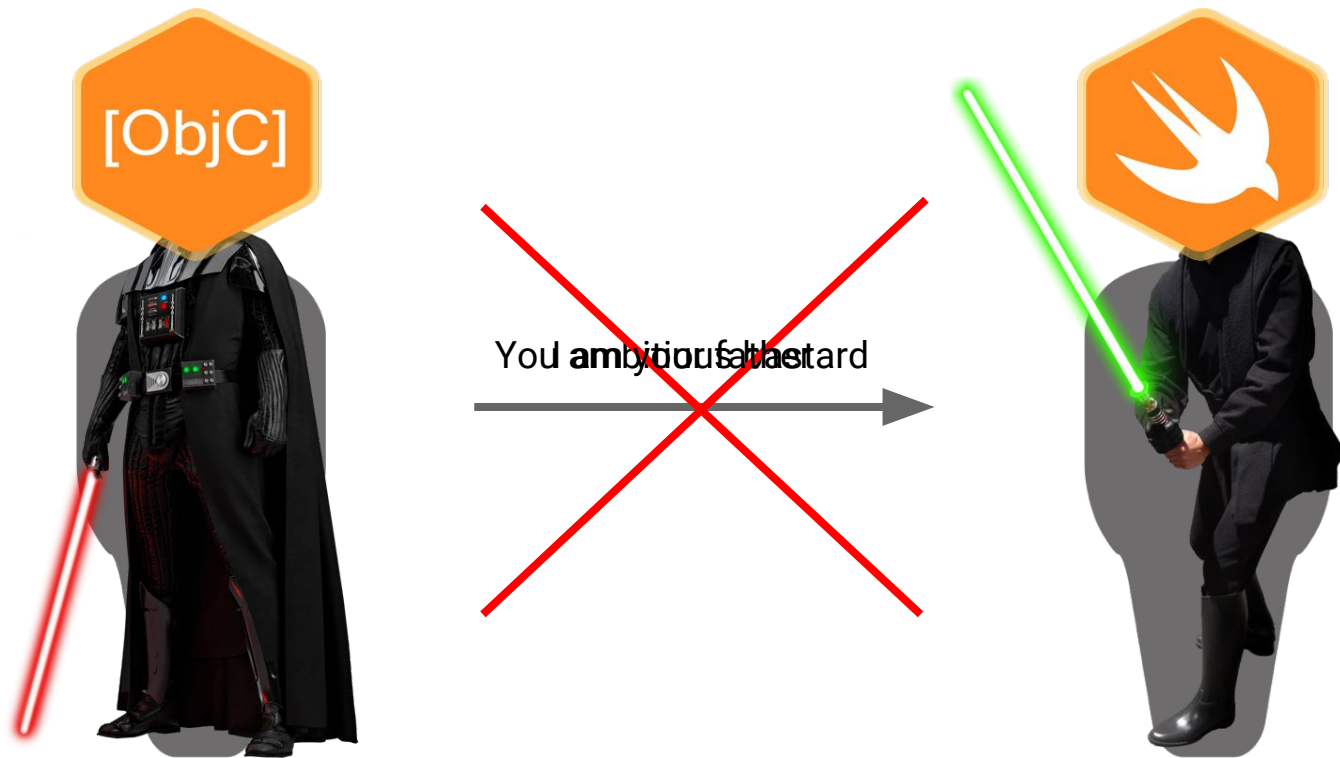


HATE BLOCKER WITH NATURAL LANGUAGE FRAMEWORK

Jan Schwarz

STRV

HATE BLOCKER



HOW TO FIX IT

- Prevent haters from using the app
- Prevent haters from connecting with good users
- Prevent haters from sending hateful messages

NATURAL LANGUAGE PROCESSING

STRV

WHAT IS NLP

- Semantics - meaning
- Sentiment - subjective information
- Entities - names, places, companies...
- Interaction between computers and human languages



NLP SUBTASKS

- Language identification
 - *Ninakuchukia*
- Tokenization - paragraphs, sentences, clauses, phrases, words...
 - *Mr. Jan Kaltoun joined STRV s.r.o. in 2016.*
 - *Mr || Jan Kaltoun joined STRV s || r || o || in 2016*
- Part of speech - noun, verb, preposition...
 - I want to fly like a fly.
 - Pronoun Verb Particle Verb Preposition Determiner Noun
- Lemmatization
 - I haven't eaten for hours
 - I have eat for hour

NLP IN IOS

- `NSLinguisticTagger`
 - Language identification
 - Tokenization
 - Part of speech
 - Lemmatization
 - Named entity recognition
- Natural Language Framework
 - `NSLinguisticTagger`
 - Language models

NATURAL LANGUAGE FRAMEWORK VS. HATERS

STRV

DETECT HATEFUL WORDS

- Create a list of hateful words
- Analyze a message
- Ban messages that contain words from the blacklist



DEMO

STRV

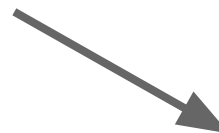
ISSUES

- Potentially long list
- Ambiguity
 - Git - moron vs. versioning system
- Negation
 - You are a moron vs. You are not a moron
- (In)direct speech
 - Sara said you were a moron
- Context
 - You have a memory like an elephant vs. You have a bottom like an elephant
- ...

CAN WE DO BETTER?



● 1.8 Ready for Sale



● 2.0 Metadata Rejected

TEXT CLASSIFIER

- Training dataset
- Build a model with CreateML
- Use CoreML to load a model and classify any text

TRAINING DATASET

- <https://github.com/t-davidson/hate-speech-and-offensive-language>
- 24783 labeled tweets

Hateful	Offensive	Neither	Class	Tweet
0	0	3	2	@mayasolovely: As a woman you shouldn't complain about cleaning up your house. & as a man you should always take the trash out...
0	3	0	1	All I wanna do is get money and fuck model bitches!
2	1	0	0	#JesusChrist was STRAIGHT> That's why the #faggots killed him. #PERIOD #SonOfGod>

**WARNING:
THE FOLLOWING DEMO
CONTAINS EXPLICIT LANGUAGE**

STRV

ISSUES

- Training dataset
 - Neither - 4163
 - Hateful - 1430
 - Offensive - 19190
- CreateML is a blackbox

RECAP

- Natural Language Framework
 - Replacement of NSLinguisticTagger
 - MLTextClassifier, MLWordTagger
- Your model is just as accurate as your data is

- Demo: <https://github.com/schwarja/hateblocker>
- jan.schwarz@strv.com

THANK YOU

Jan Schwarz

STRV

QUESTIONS

STRV