

## GOALS

- ① COMPARE NB, LSTM, SVM ON EMAIL DATASET
- ② EXPLORE DIFFERENT TOKENIZATION OPTIONS TO IMPROVE RESULTS
- ③ DETERMINE SAMPLE SIZE FOR GOOD RESULTS

## EXPERIMENTS

BASELINE NB/LSTM/SVM

EXP A: CHANGING SAMPLE SIZE OF TRAIN SET (PRETTY SIMILAR > 5000)

EXP B: Change "IMPORTANT WORDS" (15-50 SIMILAR)

EXP C: USE POS/NER TAGS (70% ACC)

## WHAT'S NOVEL

DATASET

TOKEN DEF

## CURRENT TOKEN DEFINITION

CHANGE 1-9 → <DIGIT> TAG

REMOVE NON @/#/\$/%/!/&/& PUNCTUATION

REMOVE CASE + MANDATORY D/C DATASET

# OUTLINE

## ① INTRO

NEED: SPAM STILL A PROBLEM → LEVCHENKO ✓

SPAM CLASS AN EVEN EVOLVING FIELD → HAS TO KEEP UP  
"ADVERSARIAL"

NOVELTY OF RESULTS - WHY INTERESTING TO OTHERS

WHY INTERESTING TO ME?

## ② RELATED WORK

PANTEL → SUMMARIZE + WHAT WAS NEW → NEW = NB ✓ SAHAMI ✓

PAUL GRAHAM → NEW = BETTER TOKENIZATION FOR NB

DRUGLER  
~~ANDREWS & POWERS~~ → DEMONSTRATED SVM > NB

BASYAR → ~~1ST~~ 1<sup>ST</sup> LSTM FOR SPAM CLASS

## ③ GOALS

PRJ A → "IMPROVE CLASS TECHNIQUES"

✓ VARY SIZE OF TRAIN CORPUS

✓ DEPLOY NLP TECHNIQUES (POS/NER)

PRJ B → MINIMIZE TRAIN SIZE  
→ UNIMPORTANT WORDS  
→ NLP TECHNIQUES  
→ LSTM NN

MODELS → SVM + NB

DATASET → DESCRIBE SOURCE, DESCRIBE PREPROCESSING

EXPERIMENTS → STATE PROBLEM/HYP, DESCRIBE EXP,

DESCRIBE RESULTS, INTERPRET RESULTS

FUTURE WORK → LSTM, LARGER DATASET / PERSONALIZED

BIBLIOGRAPHY

## Project C Requirements

- ① 2 MEETINGS ✓
- ② WRITTEN NOTES ✓
- ③ CODE/DATA/ANALYSIS/README -> IN GITHUBS
- ④ FINAL REPORT (4-8 PGS)

↳ ACM Conference Template

↳ TITLE AUTHOR + "MENTOR:" + URL

SEC #1 → INTRO : NEED, WHAT HAS BEEN DONE, NOVELTY +

NEW RESULTS, WHY WORK = INTERESTING

SEC #2 → RELATED WORK SUMMARY → SUMMARIZE + "WHY IT'S DIFFERENT"

↳ INCLUDE CITES

SEC #3 → GOALS FROM PRJ A/B, IF THEY ✓, WHY DIDN'T ✓

SEC #4 → DISCUSSION OF WORK -> GOAL / METHODS / DATA / RESULTS

SEC #5 → Future work

SEC #6 → BIBLIOGRAPHY

## Related Work

PANTIL ET AL. (NB)

~~Graham et al. 2002~~

Graham 2002 (NB)

Graham 2003 (NB)

SINGH ET AL (SVMs)

RAHMAN ET AL (LSTM)

KAMAMA ET AL (NB/SVM)

Exp A/B/c  
STATE

HYPOTHESIS : There is a min train

EXPERIMENT

RESULTS

INTERPRET

- 38 PTS  
TOTAL

## PROJECT B FEEDBACK

- 10 PTS [- IMPROVE NOTES + MEET w/ PROF

- 9 PTS [- INTRO - INCLUDE WHY INTERESTING TO FIELD / STUDENT + WHY NOVEL

- 3 PTS [- ALIGN CITATIONS w/ BIBLIOGRAPHY

- 8 PTS [- RELATED WORK → ORDER BY PAPER NOT FEATURE  
↳ ADD more THROUGH RW SECTION  
↳ maybe CITE / OTHER NN INFO for STATE OF ART  
↳ FIND ADDITIONAL PAPERS from ERANUM / PANTEL TL

- DATASET ✓

### WRITEN UP SUGGESTIONS

- 5 PTS

- (1) # / HIERARCHY TO SECTIONS
- (2) OVERVIEW TO SUMMARIZE METHODS → BREAKDOWN FROM THESE
- (3) FIGURE OUT HOW TO EXPAND BEYOND TOKENIZATION
- (4) BREAK UP METHODOLOGY / DATASET / RESULTS
- (5) TITLE EXP. APPROPRIATELY
- (6) USE ACM / IEEE template