

DATASETS: CSDML2010, ENRON, SPAMASSASSIN, NCI,
KAGGLE PHISING DATA

PROJECT B: CLEAN DATA, CREATE BASELINE NAIVE BAYES

EXP A: VARY TRAINING DATA SIZE BY 2,000 INTERVALS

EXP B: VARY "IMPORTANT WORDS"

EXP C: USE NLP TOOLKIT (POS TAGGING, NER)

EXP D: COMPARE NAIVE BAYES TO LSTM NEURAL NET

DOES IT
GENERALIZE?

TEST SET? → TOTALLY NEW(?) ON NEW PORTION OF
5 DATASETS

WOULD GUESS IT PERFORMS WORSE BUT
MIGHT SPEAK TO IMPORTANCE OF PERSONALIZATION?

Romain Garcke - Next to You

Blake.08 - SUPERCHONSTER

TINLICKER - REGISTRY

Eli + Fun - Come Back(%) Anorexia

Michael Casette - Bells of Konkusha

Kaki - Om

49

Next Steps:

- ① Fix ① Exp A to shorten/efficient code
- ② Exp B: Figure out how to token MNR on top x words
- ③ Exp C: Adjust input tokens w/ NER / pos tagging

FINISH LSTM + SVM

Experiments & Tests - 4/9

- ✓ #1 → Vary Train DATA SIZE
- ✓ #2 → Vary "IMPORTANT WORDS"
- ✓ #3 → USE POS TAGGING + NER
- #4 → PRODUCE LSTM / FIND PAPER + RESPONSE
- #5 → RECREATE PAMEL
- #6 → RECREATE GRAHAM

finished research to CHANGE ADD/SUB FUTURE TESTS

- SPAM CLASS STATE OF their SURVEY
- LSTM for SPAM CLASS

Top - λ^k Important Words

~~Model. log-probs~~

FIT MODEL ON ALL WORDS

(COMBINE GRAIN LOG PROBS → CONVERT TO 1 SCORE (THE GREATER (?))
+ RANK WORDS ASSOCIATED

~~REFIT w/ COMBINED FIT NEW MODEL w/ TOP λ^k WORDS
PROJECT~~

TOKENIZATION PROCESSES

GRAHAM 02

AZ19, DASHES, APOSTROPHES, \$

IGNORES "ALL DIGIT" TOKENS

LOWERCASE ONLY

TOP 15 INTERESTING TOKENS USED

GRAHAM 03

KEEPS CASE, ADDS !, ADDS . /, in b/wn #'s

ADDS SUBJ TAG