

Deceiving Authorship Detection

Tools to Write Anonymously & Current Trends in Adversarial Stylometry.

Michael Brennan, Sadia Afroz and Rachel Greenstadt. Drexel University.

Privacy, Security and Automation Lab

- Faculty
 - Dr. Rachel Greenstadt
- Graduate Students
 - Sadia Afroz (Deception Detection Lead)
 - Diamond Bishop
 - Michael Brennan
 - Aylin Caliskan
 - Ariel Stolerman (JStylo Lead Developer)
- Undergraduate Students
 - Pavan Kantharaju
 - Andrew McDonald (Anonymouth Lead Developer)

26C3/28C3 Diff

- Review & Updated Analysis of 26C3 Material
 - New Corpus (45 authors)
 - New Method (Writeprints)
 - Much more robust results.
- The tools we discussed are now built!
 - JStylo
 - Anonymouth
- Detecting Deception in Adversarial Writing

An Overview

- What is “Authorship Recognition” and “Adversarial Stylometry”?
- What is the anonymity threat?
- Analyzing & Deceiving Authorship Recognition
- Two Tools
 - JStylo
 - Anonymouth
- Detecting Deception

What is Authorship Recognition?

- The basic question: “who wrote this document?”
- **Stylometry:** The study of attributing authorship to documents based only on the linguistic style they exhibit.
 - “Linguistic Style” Features: sentence length, word choices, syntactic structure, etc.
 - Handwriting, content-based features, and contextual features are not considered.
 - Individuals have unique writing styles because language is learned on an individual basis.
- In this presentation, stylometry and authorship recognition are used interchangeably.

What is Adversarial Stylometry?

- Adversarial Stylometry: Applying deception to writing style in order to affect the outcome of stylometric analysis.
 - But, is writing style modifiable? (Yes!)
 - Is it possible to deceive stylometry through altered writing style? (Yes!)
 - What are the implications of looking at stylometry in an adversarial context?

How Can Stylometry be a Threat?

- Supervised Stylometry
 - Given a set of documents of known authorship, classify a document of unknown authorship.
 - Hypothetical Scenario: Alice the Anonymous Blogger vs. Bob the Abusive Employer.
- Unsupervised Stylometry
 - Given a set of documents of unknown authorship, cluster them into author groups.
 - Hypothetical Scenario: Anonymous Forum vs. Oppressive Government.

Purely Hypothetical?

- Previous examples are purely hypothetical. What about a real example?
- From “Inside WikiLeaks” by Daniel Domscheit-Berg:
 - “I nudged Julian with my foot. We exchanged glances and started giggling. If someone had run WikiLeaks documents through such a program, he would have discovered that the same two people were behind all the various press releases, document summaries, and correspondence issued by the project. The official number of volunteers we had was also, to put it mildly, grotesquely exaggerated.”



Adversarial Stylometry: A Review

- Understand the threat model
- Build a corpus.
- Evaluate current methods of stylometry against adversarial text.
- Analyze results and develop tools.

Threat Model

- Threat: Authorship recognition can identify you if there are sufficient writing samples and a set of suspects.
 - 6500+ words of training data per author
 - 500+ words of testing data
 - 50 or less suspects
 - These may be different:
 - Tweets (short messages)
 - Large numbers of authors (Writeprints)
- Old assumption: Writing style is invariant.
 - It's like a fingerprint, you can't really change it.

Circumvention Methods

- Challenge: conceive methods of circumventing writing style analysis.
 - **Obfuscation**
 - An author attempts to write a document in such a way that their personal writing style will not be recognized.
 - **Imitation**
 - An author attempts to write a document such that the writing style will be recognized as that of another specific author.
 - **Translation*:**
 - Machine translation is used to translate a document to one or more languages and then back to the original language.

Building a Corpus

- Corpus = Dataset of documents.
- Data sets for adversarial stylometry do not exist. Participants are required to craft intentionally adversarial passages.
- Participation had three parts:
 - Submit 6500 words of pre-existing writing from a formal source.
 - Write a new 500 word obfuscation passage.
 - Task: Describe your neighborhood.
 - Write a new 500 word imitation passage.
 - Task: Imitate Cormac McCarthy, describe your day.
- Authors had no formal training or knowledge in linguistics or stylometry.

Brennan-Greenstadt Corpus

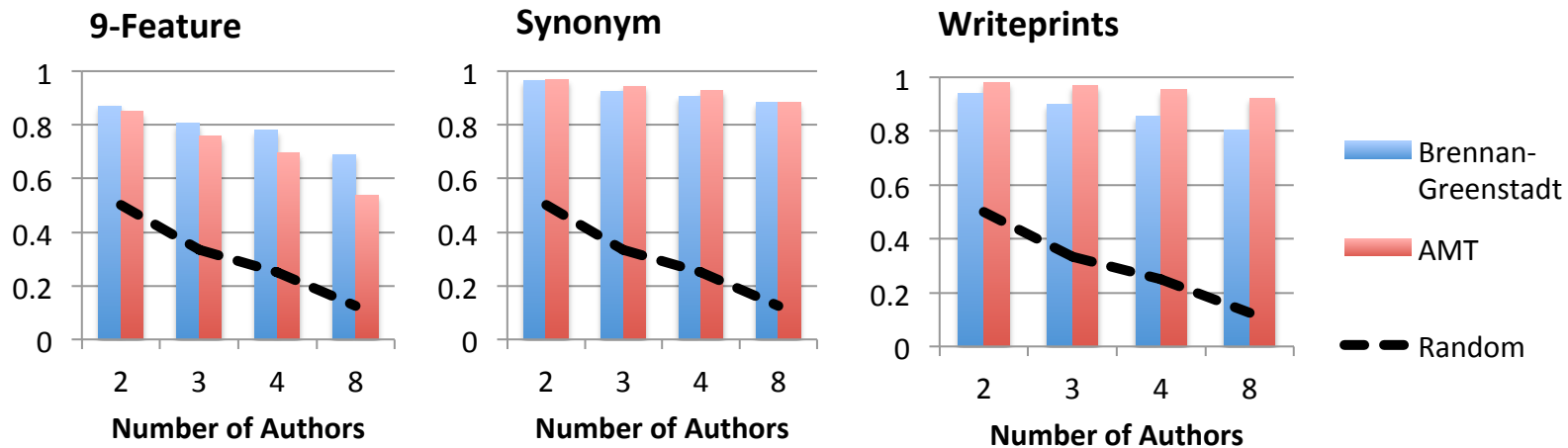
- 12 Individual Authors.
- Participants contacted through classes, colleagues, friends at Drexel University.
 - Motive for proper participation.
- One-on-one interaction with participants.
- Corpus is publicly available at <https://psal.cs.drexel.edu>
- Good for preliminary results, but we need something better.
 - Too small.
 - Too homogenous.

Building a Better Corpus with Amazon Mechanical Turk

- Drexel AMT Corpus
 - AMT = Amazon Mechanical Turk
- Same tasks as previous corpus.
- Only 45 of 101 of submissions are usable!
 - 45 Accepted Submissions.
 - Guidelines without spoiling data set. Must follow directions and:
 - Pre-existing writing must be formal in nature
 - Remove non-content
 - Minimal dialogue/quotations
 - Refrain from submitting: small samples, lab reports, Q&As, etc.
- Released today. Publicly available at <https://psal.cs.drexel.edu>
- This corpus is large, diverse, and unique.

Original vs. AMT Corpus

- AMT Corpus evaluated just as strongly as Drexel.
 - 9-Features does worse, Synonym does the same, Writeprints does better.



Evaluate Stylometry Methods Against the Corpus

- Three methods of Stylometry
 - 9-Feature / Neural Network
 - Synonym-Based Approach
 - Writeprints / SVM

Method 1:

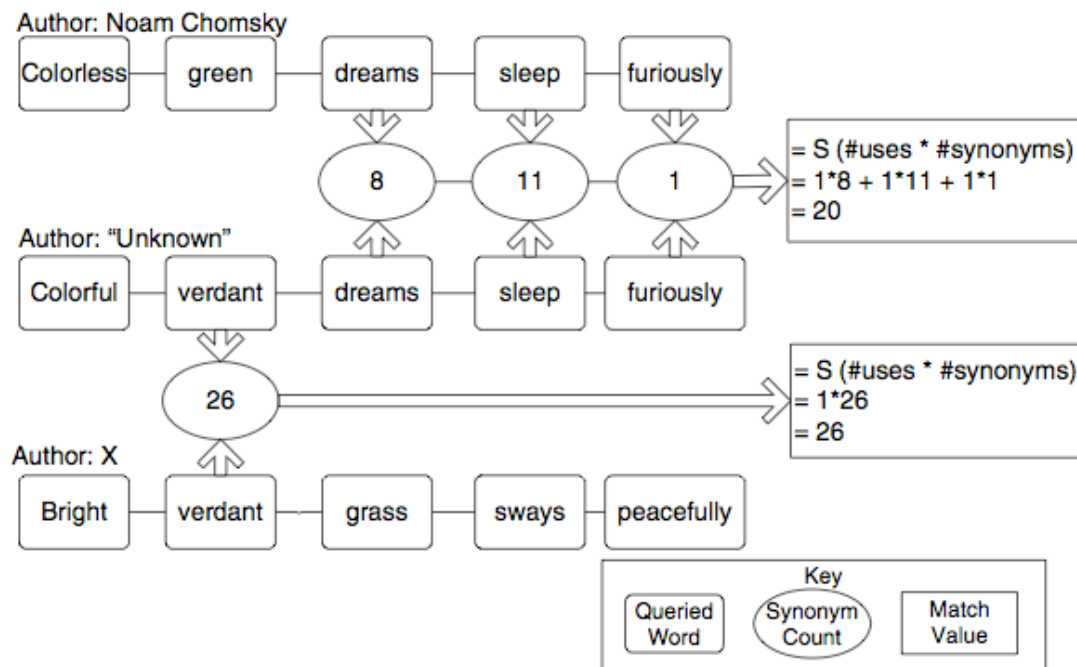
9-Feature Set Neural Network

- Simple stylometric approach. Demonstrates potential effectiveness with a small number of obscure metrics.
- 9-Feature Set
 - Unique words, Complexity, Sentence Count, Average Sentence Length, Average Syllable Count, Character Count, Letter Count, Gunning-Fog Readability Index, Flesch Reading Ease Score.
- Neural Network Classifier.

Method 2:

Synonym-Based Approach

- Examines word choices when compared to available synonyms and frequency of use.
- Clark & Hannon, 2007.
- Good demonstration of single feature type stylometry.



Method 3:

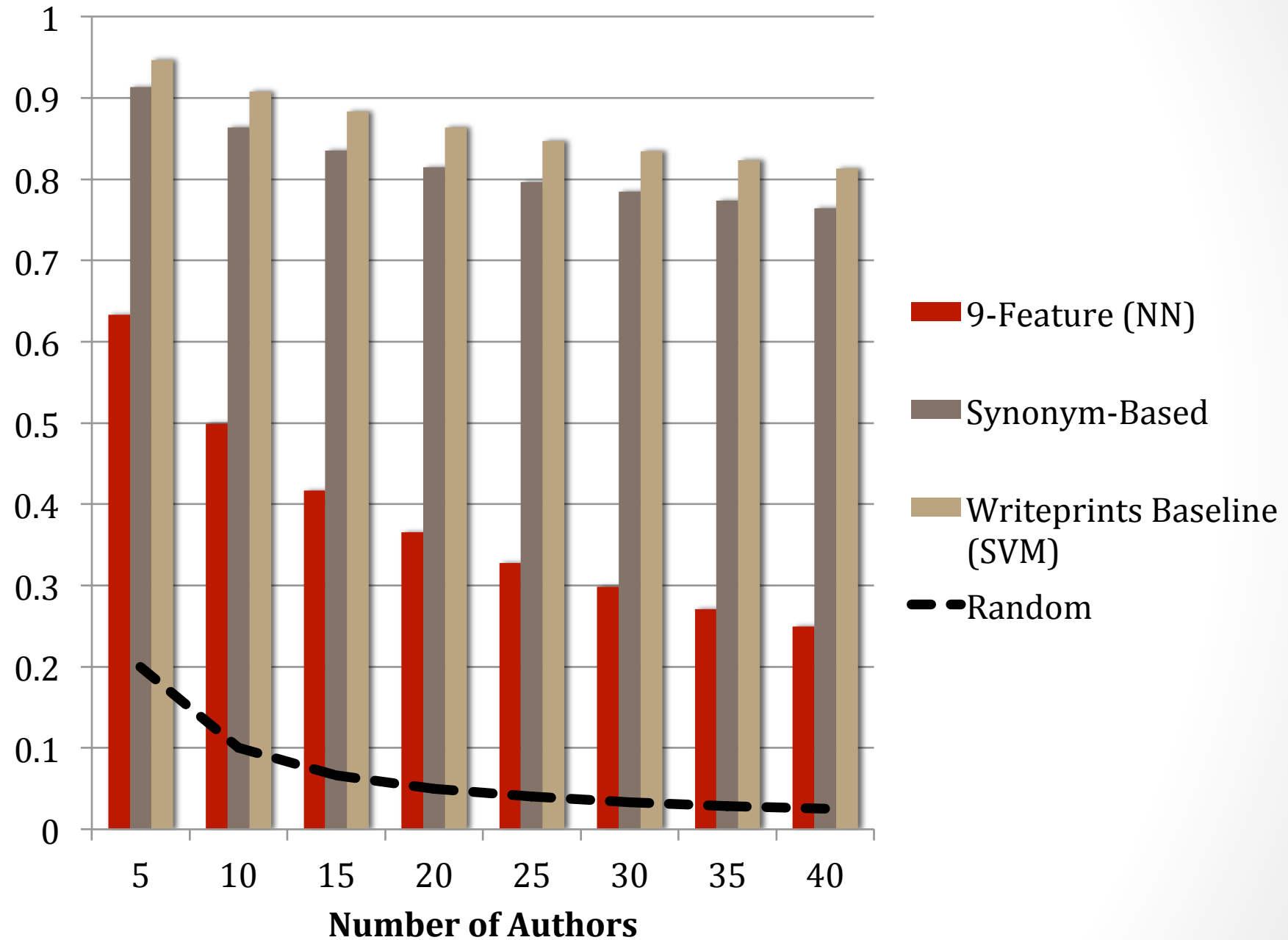
Writeprints (SVM)

- Based on the Writeprints approach by Abbasi & Chen, 2008.
 - Writeprints Baseline Feature Set.
 - Contains hundreds of features including character and word n-grams, function words, parts-of-speech tags, punctuation, and character level metrics.
- Support Vector Machine Classifier
 - Standard for multi-class classification in stylometry.
- Implementation of the full Writeprints approach uses a more extensive feature set and unique classification approach.

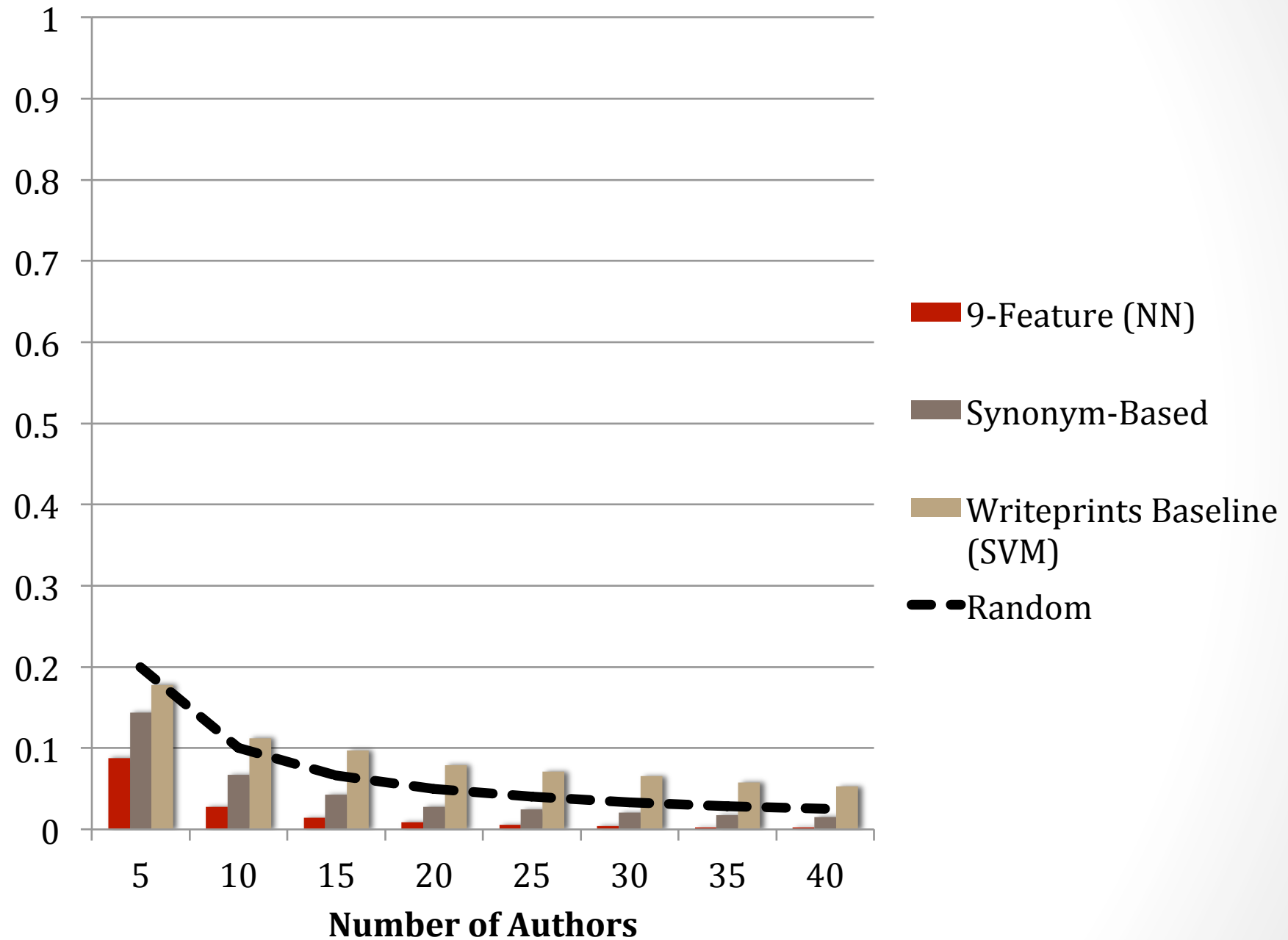
Detecting Authorship

- Four experiments.
 - **Baseline:** Evaluate the non-adversarial passages to establish a baseline accuracy for each method and feature set.
 - **Obfuscation:** Attempt classification of the obfuscation passages when classifiers are trained on the baseline data.
 - **Imitation:** Attempt classification of the imitation passages.
 - **Imitation Success:** Analyze the number of successful imitations.

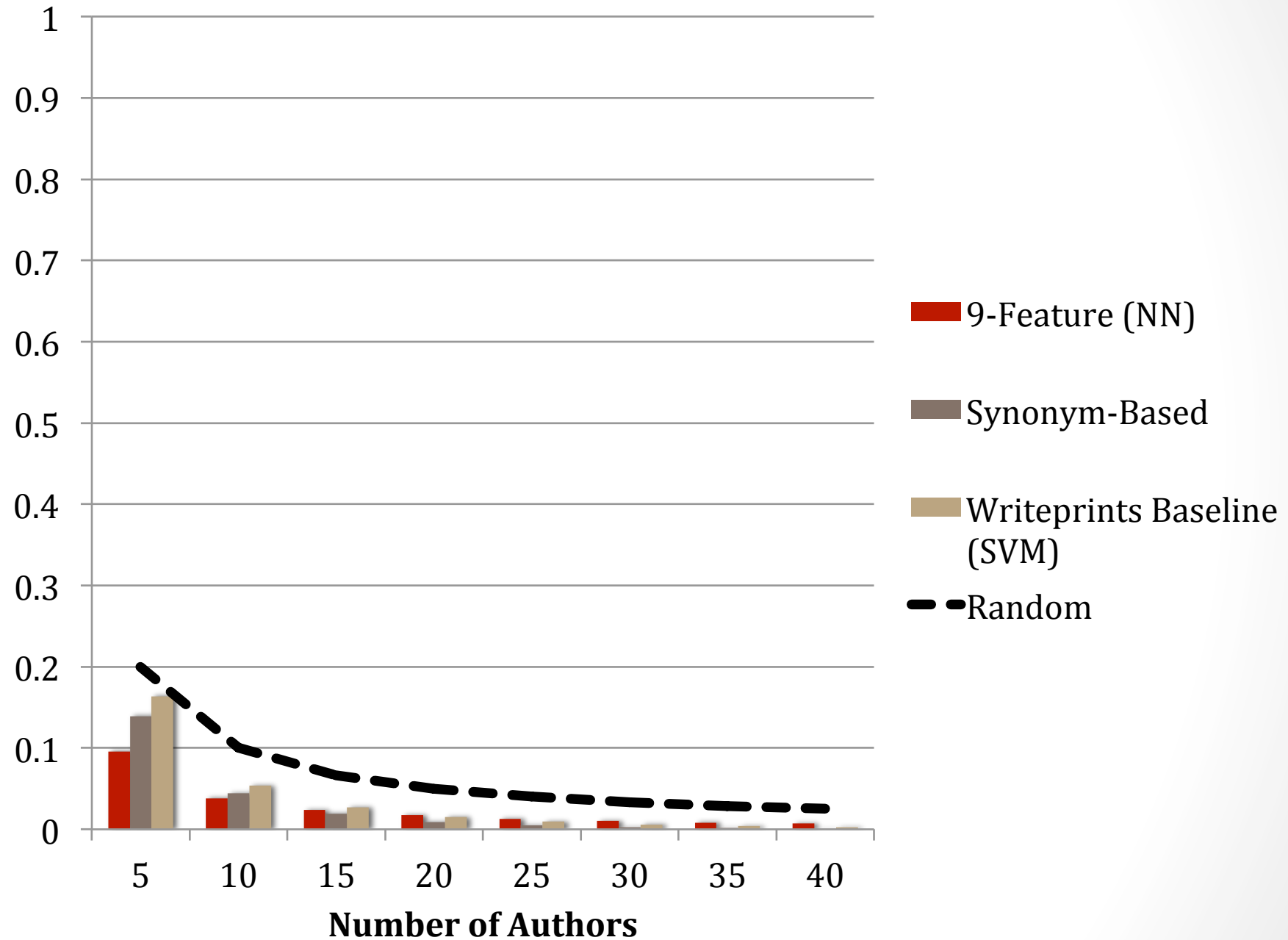
Baseline Precision



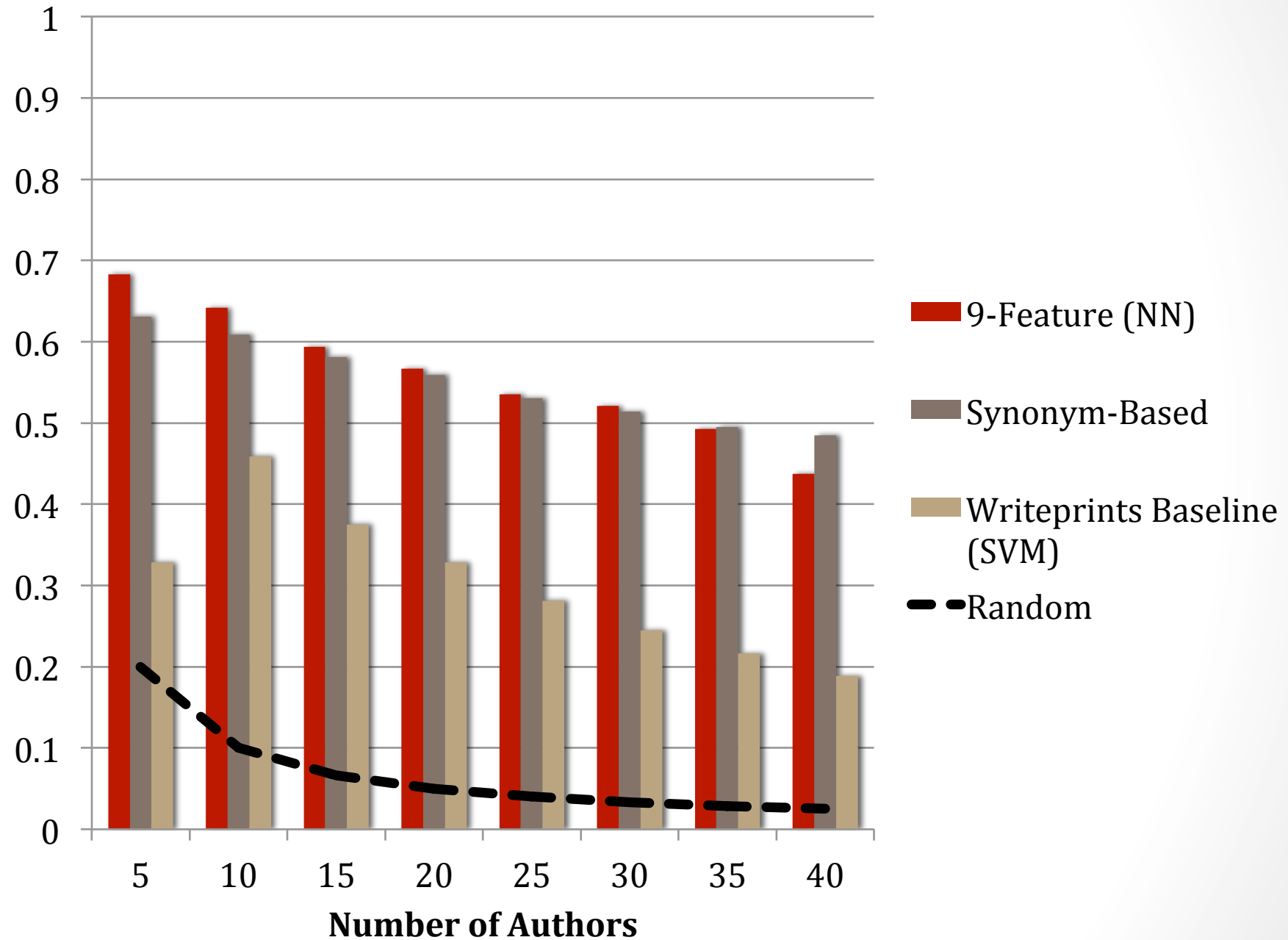
Obfuscation Precision



Imitation Precision



Imitation Success (Framing Cormac McCarthy)



Two Tools

- JStylo: Authorship Recognition Analysis Tool.
- Anonymouth: Authorship Recognition Evasion Tool.
- Free, Open Source. (GNU GPL)
- Alpha releases available today at <https://psal.cs.drexel.edu>
 - Migrating to GitHub soon.



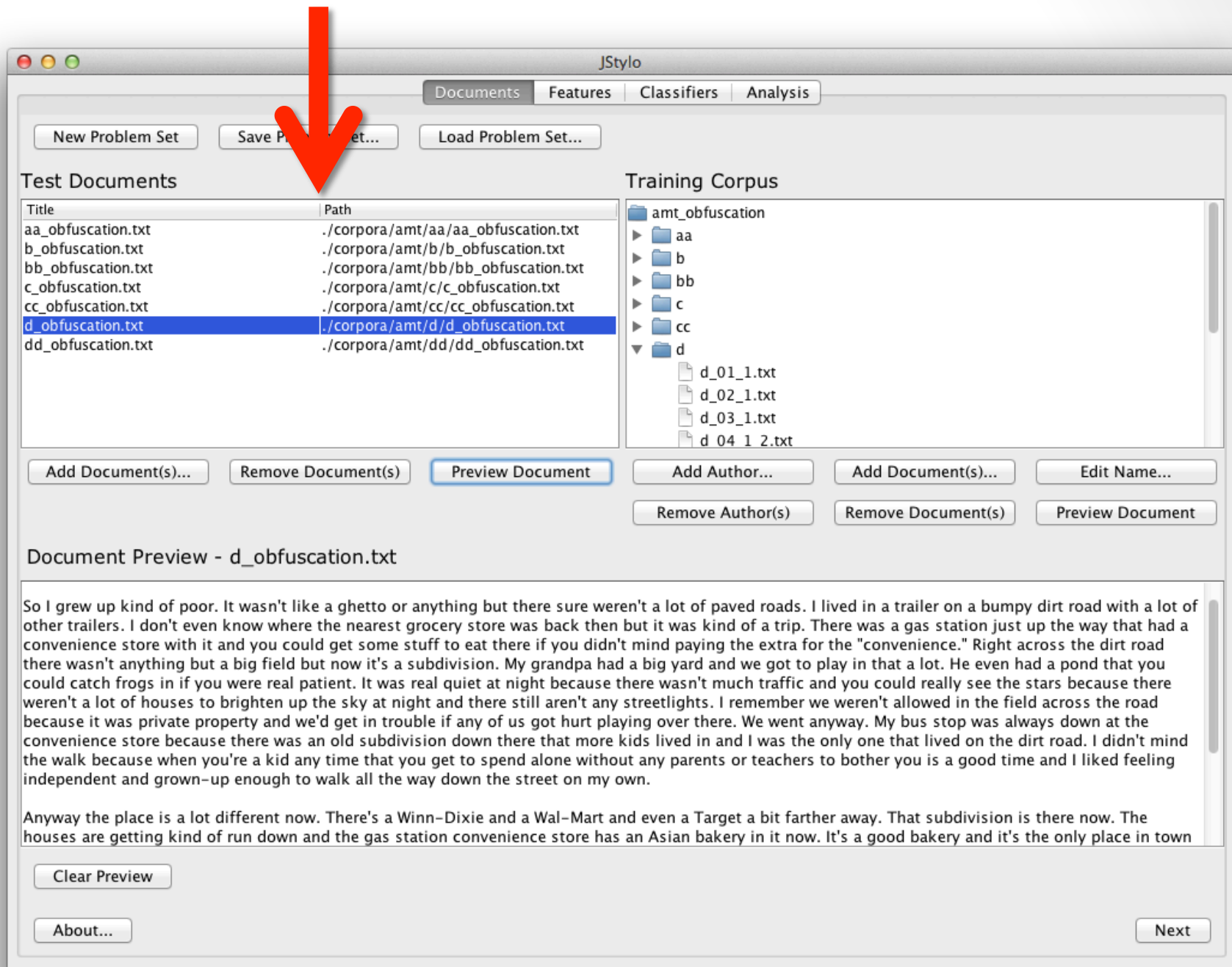
JStylo: The Problem

- Stylometry-based research is difficult.
- Existing tools are good but limited.
 - Weka provides a suite of machine-learning classification tools.
 - Not tailored for text analysis – no feature extraction ability.
 - Functions better as an API for software development.
 - JGAAP has a strong basic toolset for stylometry.
 - Limited in running multiple feature sets.
 - Strong API.
 - Extendable. Intended to be used this way.
- Nuances of stylometry are not easy to grasp.
- Many open research questions related to authorship. We need an easy-to-use tool that both researchers and non-technical users can understand.

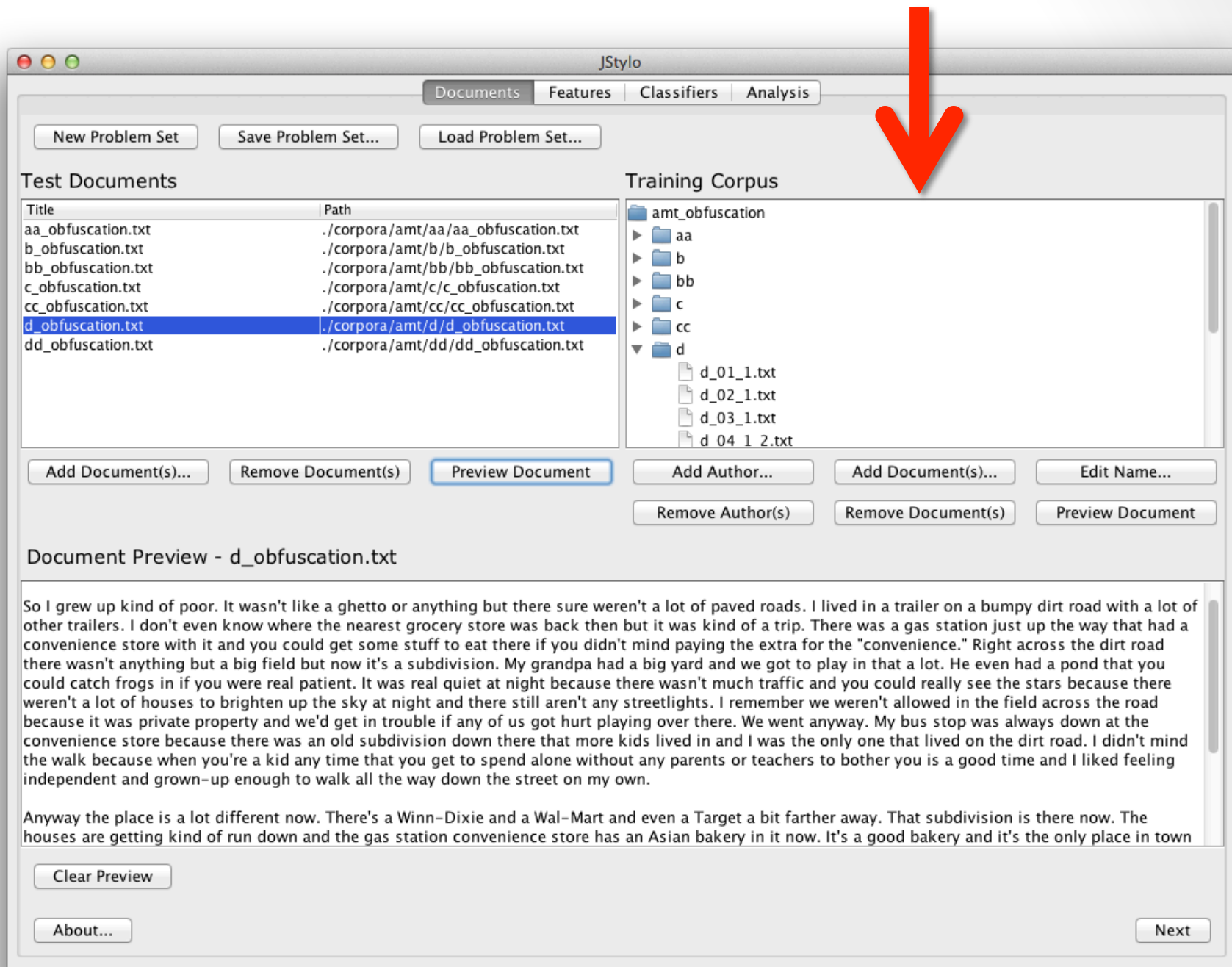
JStylo

- JStylo is an authorship recognition analysis tool. It is built upon a framework of:
 - JGAAP (Java Graphical Authorship Attribution Project)
 - Weka 3 Data Mining Software
- Features
 - Two existing adversarial corpora, featured in this presentation, and new corpus building functionality.
 - Wide selection of writing feature extractors and ability to add new extractors.
 - Wide selection of machine learning based classifiers.
 - Intuitive GUI.
- Alpha Release Available Now: <https://psal.cs.drexel.edu>

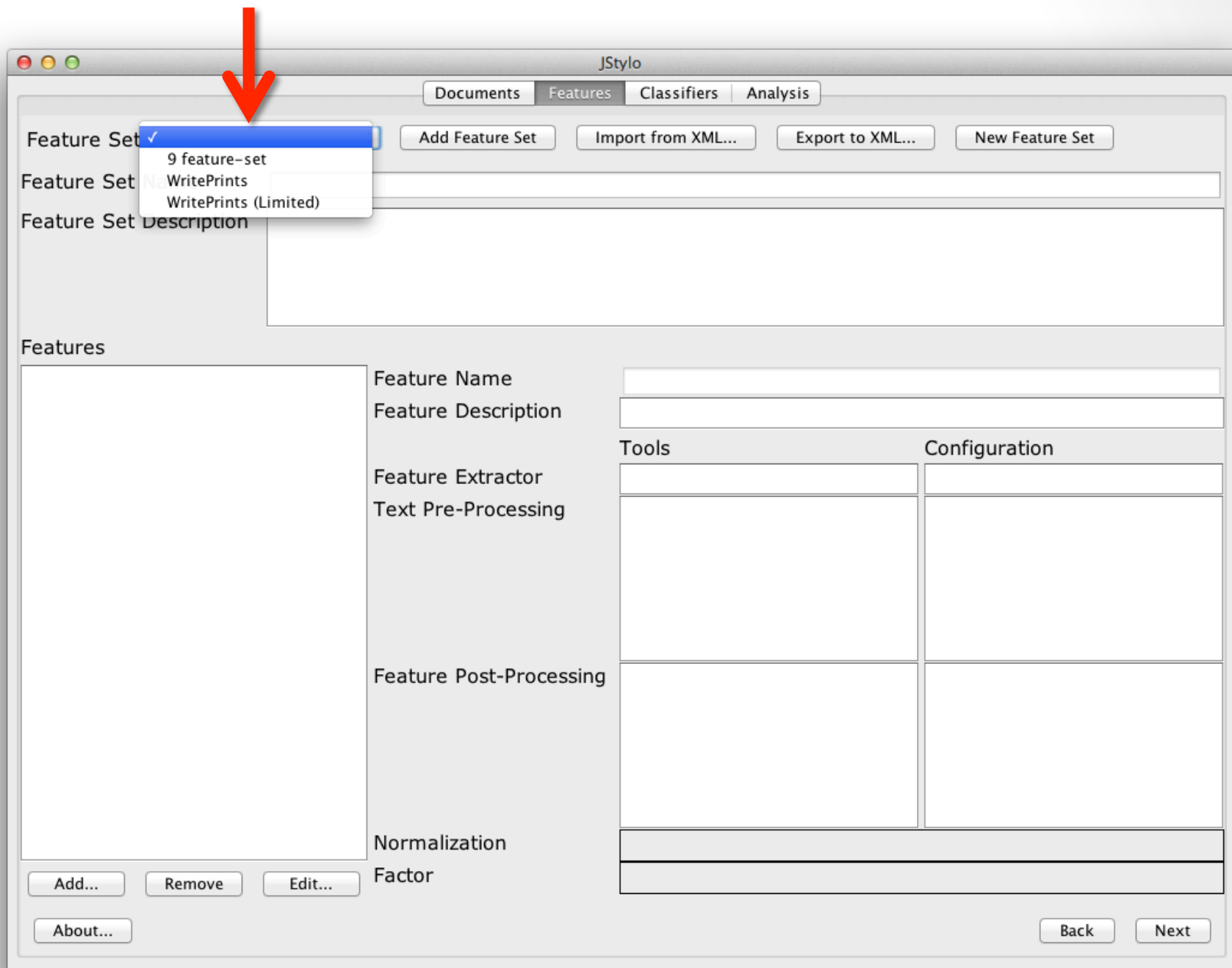
JStylo Demo



(send bug reports, suggestions, questions to ariels@drexel.edu)



(send bug reports, suggestions, questions to ariels@drexel.edu)



(send bug reports, suggestions, questions to ariels@drexel.edu)

JStylo

Documents Features Classifiers Analysis

Feature Set 9 feature-set Add Feature Set Import from XML... Export to XML... New Feature Set

Feature Set Name 9 feature-set

Feature Set Description 9 features used by Brennan and Greenstadt.

Features

Unique Words Count
Complexity
 Sentence Count
 Average Sentence Length
 Average Syllables in Word
 Gunning-Fog Readability Index
 Character Space
 Letter Space
 Flesch Reading Ease Score

Feature Name Complexity

Feature Description Ratio of unique words to total number of words in the document.

Tools	Configuration
Unique words count	
Word-edges Punctuation Stripper Unify Case	
Number of words in the document	
1.0	

Add... Remove Edit... About...

Back Next

(send bug reports, suggestions, questions to ariels@drexel.edu)

JStylo

Documents Features Classifiers Analysis

Feature Set 9 feature-set Add Feature Set Import from XML... Export to XML... New Feature Set

Feature Set Name 9 feature-set

Feature Set Description 9 features used by Brennan and Greenstadt.

Features

- Unique Words Count
- Complexity
- Sentence Count
- Average Sentence Length
- Average Syllables in Word
- Gunning-Fog Readability Index
- Character Space
- Letter Space
- Flesch Reading Ease Score

Feature Name Complexity

Feature Description Ratio of unique words to total number of words in the document.

Feature Extractor

Text Pre-Processing

Feature Post-Processing

Normalization

Factor

Tools Configuration

Unique words count

Word-edges Punctuation Stripper

Unify Case

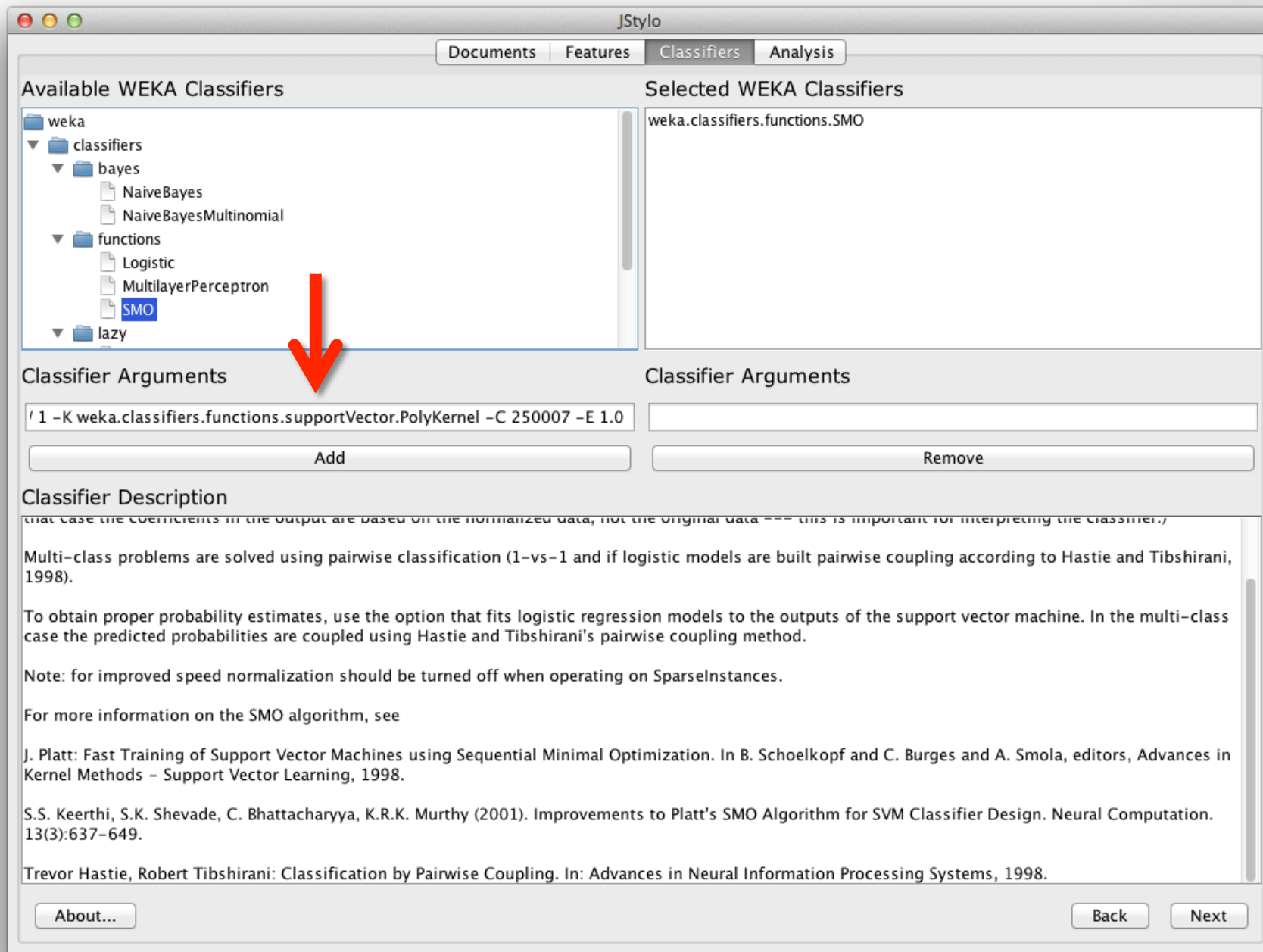
Number of words in the document

1.0

Add... Remove Edit... About...

Back Next

(send bug reports, suggestions, questions to ariels@drexel.edu)



(send bug reports, suggestions, questions to ariels@drexel.edu)

DocumentsFeaturesClassifiersAnalysis

Analysis Type

☒ Train on training corpus and classify test doc...
 ☐ Run 10-folds cross validation on training cor...

Configuration

☒ Output feature vectors (ARFF format)
 ☒ Use sparse representation for feature vectors

Post Analysis

Training to ARFF...

Training to CSV...

Test to ARFF...

Test to CSV...

Run Analysis

Stop

Results

2011-12-28, 13:59:312011-12-28, 14:04:392011-12-28, 14:12:272011-12-28, 14:23:35

```
{0 233,1 0.47166,2 26,3 19,4 1.550607,5 12.215385,6 2663,7 2084,8 56.368623,9 _dummy_}
2011-12-28, 14:23:40 Starting training and testing phase...

=====

Running analysis with classifier 1 out of 1:
> Classifier: weka.classifiers.functions.SMO
> Options: -C 1.0 -L 0.0010 -P 1.0E-12 -N 0 -V -1 -W 1 -K weka.classifiers.functions.supportVector.PolyKernel -C 250007 -E 1.0

2011-12-28, 14:23:40 Starting classification...
2011-12-28, 14:23:40 done!

Results:
=====
doc \ author    aa          b          bb          c          cc          d
-----
aa_obfuscation.txt 0.222222    0.185185    0.148148    0.259259 +    0.074074    0.111111
b_obfuscation.txt  0.185185    0.148148    0.259259 +    0.222222    0.074074    0.111111
bb_obfuscation.txt 0.222222    0.185185    0.259259 +    0.148148    0.074074    0.111111
c_obfuscation.txt  0.222222    0.148148    0.259259 +    0.185185    0.074074    0.111111
cc_obfuscation.txt 0.185185    0.148148    0.259259 +    0.222222    0.074074    0.111111
d_obfuscation.txt  0.222222    0.148148    0.259259 +    0.185185    0.111111    0.074074
dd_obfuscation.txt 0.222222    0.148148    0.259259 +    0.185185    0.111111    0.074074
```

Save Results...

About...

Back

(send bug reports, suggestions, questions to ariels@drexel.edu)

Results:

=====

doc \ author	aa	b	bb	c
aa_obfuscation.txt	0.222222	0.185185	0.148148	0.259259 +
b_obfuscation.txt	0.185185	0.148148	0.259259 +	0.222222
bb_obfuscation.txt	0.222222	0.185185	0.259259 +	0.148148
c_obfuscation.txt	0.222222	0.148148	0.259259 +	0.185185
cc_obfuscation.txt	0.185185	0.148148	0.259259 +	0.222222
d_obfuscation.txt	0.222222	0.148148	0.259259 +	0.185185
dd_obfuscation.txt	0.222222	0.148148	0.259259 +	0.185185

(send bug reports, suggestions, questions to ariels@drexel.edu)

Documents

Features

Classifiers

Analysis

Analysis Type

Configuration

Post Analysis

☐ Train on training corpus and classify test doc...
 ☒ Output feature vectors (ARFF format)

☒ Run 10-folds cross validation on training cor...
 ☒ Use sparse representation for feature vectors

Training to ARFF...

Training to CSV...

Test to ARFF...

Test to CSV...

Run Analysis

Stop

Results

2011-12-28, 13:59:31

2011-12-28, 14:04:39

2011-12-28, 14:12:27

2011-12-28, 14:23:35

2011-12-28, 14:23:54

Correctly Classified Instances

61

61.6162 %

Incorrectly Classified Instances

38

38.3838 %

Kappa statistic

0.5338

Mean absolute error

0.2127

Root mean squared error

0.315

Relative absolute error

89.4837 %

Root relative squared error

91.4539 %

Total Number of Instances

99

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	ROC Area	Class
	0	0	0	0	0	?	
	0.357	0.071	0.455	0.357	0.4	0.687	aa
	0.647	0.159	0.458	0.647	0.537	0.841	b
	0.842	0.075	0.727	0.842	0.78	0.95	bb
	0.611	0.062	0.688	0.611	0.647	0.88	c
	0.944	0.099	0.68	0.944	0.791	0.947	cc
	0.077	0	1	0.077	0.143	0.817	d
Weighted Avg.	0.616	0.081	0.663	0.616	0.579	0.863	

=== Confusion Matrix ===

a

b

c

d

e

f

g

<-- classified as

Save Results...

About...

Back

=== Summary ===

Correctly Classified Instances	61	61.6162 %
Incorrectly Classified Instances	38	38.3838 %
Kappa statistic	0.5338	
Mean absolute error	0.2127	
Root mean squared error	0.315	
Relative absolute error	89.4837 %	
Root relative squared error	91.4539 %	
Total Number of Instances	99	

=== Detailed Accuracy By Class ===

	TP Rate	FP Rate	Precision	Recall	F-Measure	ROC Area	Class
	0	0	0	0	0	?	
	0.357	0.071	0.455	0.357	0.4	0.687	aa
	0.647	0.159	0.458	0.647	0.537	0.841	b
	0.842	0.075	0.727	0.842	0.78	0.95	bb
	0.611	0.062	0.688	0.611	0.647	0.88	c
	0.944	0.099	0.68	0.944	0.791	0.947	cc
	0.077	0	1	0.077	0.143	0.817	d
Weighted Avg.	0.616	0.081	0.663	0.616	0.579	0.863	

(send bug reports, suggestions, questions to ariels@drexel.edu)

JStylo Dev Goals

- Wider selection of classification methods and features.
 - Writeprints, Synonym-based, more Weka methods.
 - Ensemble classifiers, weighted averaging.
 - Greater pre and post-processing options.
- Easier to use and understand for non-technical users.
 - Adding an online tutorial.
 - GUI installs of new feature extractors and classifiers.
- Logging and graphing results over multiple experiments.
 - Visualization of documents, authors, and classifications.

Anonymouth: The Problem

- Authorship recognition can be a legitimate threat to privacy and anonymity.
- Intuition in changing writing style goes a long way, but may not be enough and may not be sustainable over multiple documents.
 - We already see methods that offer some resistance to adversarial passages.
- Fully automated text anonymization is an intractable problem.
 - We need a solution that explains authorship recognition nuances as needed and assists the authoring making the most useful changes towards anonymity.

Anonymouth

- Anonymouth is an authorship recognition circumvention tool. It is built upon a framework of:
 - JStylo (JGAAP & Weka)
 - WordNet
- Features
 - Corpora, feature extractor, and classifier functionality from JStylo.
 - Suggestion system for modifying documents to evade authorship detection. Ideal value for each feature is calculated, existence of the features is highlighted, user is assisted in changing them.
 - Iterative approach to anonymizing writing style.
 - Dictionary / Synonyms / Interactive Editing Console
- Alpha Release Available Now: <https://psal.cs.drexel.edu>

Anonymouth Demo

Anonymouth

Documents | Features | Classifiers | Editor

New Problem Set | Save Problem Set... | Load Problem Set...

Your Document to Anonymize

Title	Path
-------	------

Add Docu... Remove D... Preview Do...

Your Sample Documents

Title	Path
-------	------

Add Docu... Remove D... Preview Do...

Other Sample Documents

Authors

Add Autho... Add Docu... Edit Name...
Remove Au... Remove D... Preview Do...

Document Preview

Clear Preview

Next

Anonymouth

Documents Features Classifiers Editor

New Problem Set Save Problem Set... Load Problem Set...

Your Document to Anonymize

Title	Path
a_09_3.txt	/Users/robotcaptain/Dow...

Add Docu... Remove D... Preview Do...

Your Sample Documents

Title	Path
a_01_1.txt	/Users/robotcaptain/Dow...
a_02_1_2.txt	/Users/robotcaptain/Dow...
a_03_2.txt	/Users/robotcaptain/Dow...
a_04_2_3.txt	/Users/robotcaptain/Dow...
a_05_3.txt	/Users/robotcaptain/Dow...
a_07_3.txt	/Users/robotcaptain/Dow...
a_08_3.txt	/Users/robotcaptain/Dow...
a_10_3.txt	/Users/robotcaptain/Dow...
a_11_3.txt	/Users/robotcaptain/Dow...
a_12_3.txt	/Users/robotcaptain/Dow...
a_06_3.txt	/Users/robotcaptain/Dow...

Add Docu... Remove D... Preview Do...

Other Sample Documents

- Authors
 - d
 - b
 - c
 - c_01_1.txt
 - c_02_1.txt
 - c_03_1.txt
 - c_04_1.txt
 - c_05_1.txt
 - c_06_1.txt
 - c_07_1.txt

Add Autho... Add Docu... Edit Name... Remove Au... Remove D... Preview Do...

Document Preview - a_09_3.txt

life often contains.

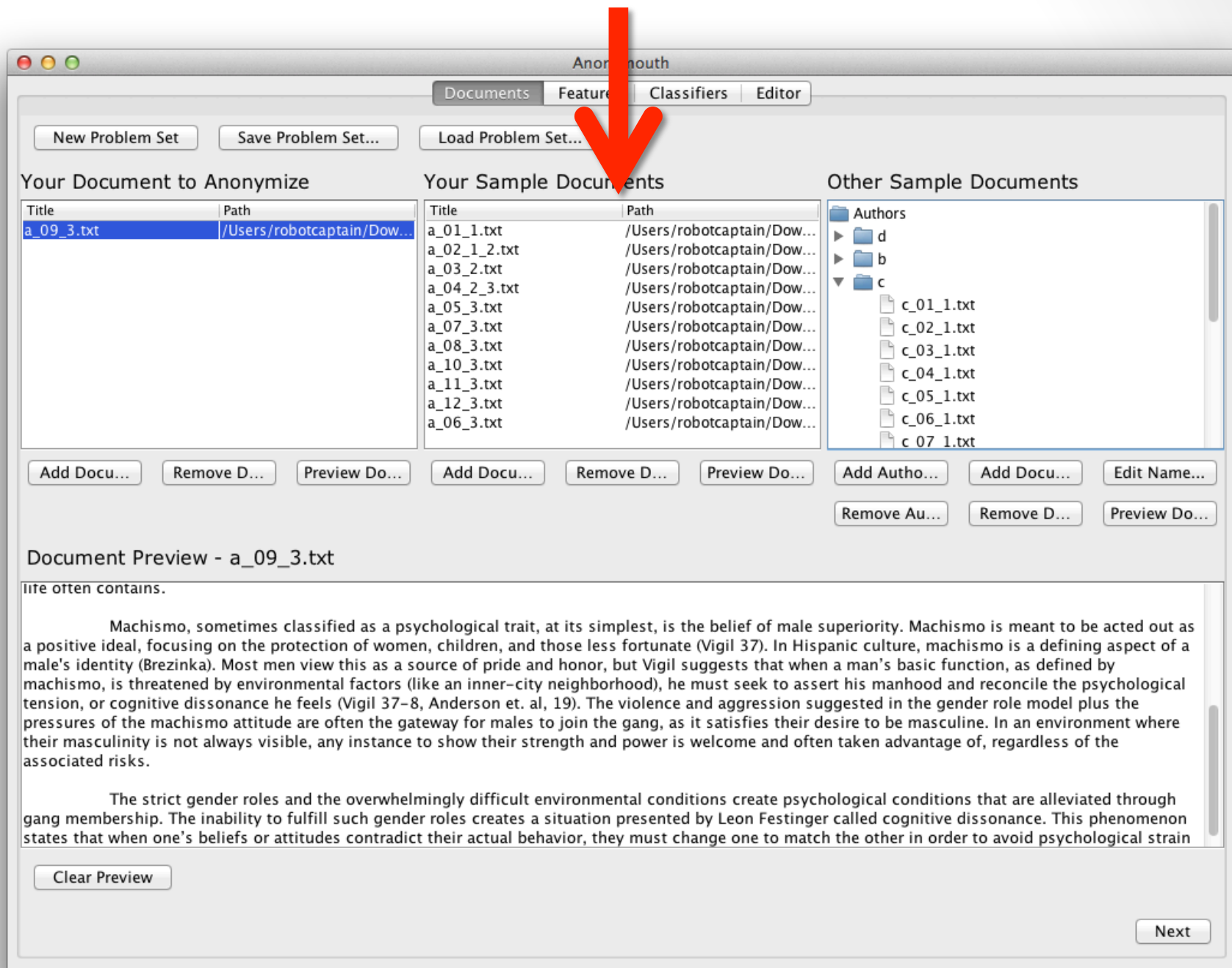
Machismo, sometimes classified as a psychological trait, at its simplest, is the belief of male superiority. Machismo is meant to be acted out as a positive ideal, focusing on the protection of women, children, and those less fortunate (Vigil 37). In Hispanic culture, machismo is a defining aspect of a male's identity (Brezinka). Most men view this as a source of pride and honor, but Vigil suggests that when a man's basic function, as defined by machismo, is threatened by environmental factors (like an inner-city neighborhood), he must seek to assert his manhood and reconcile the psychological tension, or cognitive dissonance he feels (Vigil 37-8, Anderson et. al, 19). The violence and aggression suggested in the gender role model plus the pressures of the machismo attitude are often the gateway for males to join the gang, as it satisfies their desire to be masculine. In an environment where their masculinity is not always visible, any instance to show their strength and power is welcome and often taken advantage of, regardless of the associated risks.

The strict gender roles and the overwhelmingly difficult environmental conditions create psychological conditions that are alleviated through gang membership. The inability to fulfill such gender roles creates a situation presented by Leon Festinger called cognitive dissonance. This phenomenon states that when one's beliefs or attitudes contradict their actual behavior, they must change one to match the other in order to avoid psychological strain

Clear Preview

Next

(send bug reports, suggestions, questions to awm32@drexel.edu)



(send bug reports, suggestions, questions to awm32@drexel.edu)

Anonymous

Documents Features Classifiers Editor

New Problem Set Save Problem Set... Load Problem Set...

Your Document to Anonymize

Title	Path
a_09_3.txt	/Users/robotcaptain/Dow...

Add Docu... Remove D... Preview Do...

Your Sample Documents

Title	Path
a_01_1.txt	/Users/robotcaptain/Dow...
a_02_1_2.txt	/Users/robotcaptain/Dow...
a_03_2.txt	/Users/robotcaptain/Dow...
a_04_2_3.txt	/Users/robotcaptain/Dow...
a_05_3.txt	/Users/robotcaptain/Dow...
a_07_3.txt	/Users/robotcaptain/Dow...
a_08_3.txt	/Users/robotcaptain/Dow...
a_10_3.txt	/Users/robotcaptain/Dow...
a_11_3.txt	/Users/robotcaptain/Dow...
a_12_3.txt	/Users/robotcaptain/Dow...
a_06_3.txt	/Users/robotcaptain/Dow...

Add Docu... Remove D... Preview Do...

Other Sample Documents

- Authors
 - d
 - b
 - c
 - c_01_1.txt
 - c_02_1.txt
 - c_03_1.txt
 - c_04_1.txt
 - c_05_1.txt
 - c_06_1.txt
 - c_07_1.txt

Add Autho... Add Docu... Edit Name... Remove Au... Remove D... Preview Do...

Document Preview - a_09_3.txt


life often contains.

Machismo, sometimes classified as a psychological trait, at its simplest, is the belief of male superiority. Machismo is meant to be acted out as a positive ideal, focusing on the protection of women, children, and those less fortunate (Vigil 37). In Hispanic culture, machismo is a defining aspect of a male's identity (Brezinka). Most men view this as a source of pride and honor, but Vigil suggests that when a man's basic function, as defined by machismo, is threatened by environmental factors (like an inner-city neighborhood), he must seek to assert his manhood and reconcile the psychological tension, or cognitive dissonance he feels (Vigil 37-8, Anderson et. al, 19). The violence and aggression suggested in the gender role model plus the pressures of the machismo attitude are often the gateway for males to join the gang, as it satisfies their desire to be masculine. In an environment where their masculinity is not always visible, any instance to show their strength and power is welcome and often taken advantage of, regardless of the associated risks.

The strict gender roles and the overwhelmingly difficult environmental conditions create psychological conditions that are alleviated through gang membership. The inability to fulfill such gender roles creates a situation presented by Leon Festinger called cognitive dissonance. This phenomenon states that when one's beliefs or attitudes contradict their actual behavior, they must change one to match the other in order to avoid psychological strain

Clear Preview

Next



(send bug reports, suggestions, questions to awm32@drexel.edu)

Documents

Features

Classifiers

Editor

Feature Set

9 feature-set

Add Feature Set
Import from XML...
Export to XML...
New Feature Set

Feature Set Name

9 feature-set

Feature Set Description

9 features used by Brennan and Greenstadt.

Unique Words Count

Complexity

Sentence Count

Average Sentence Length

Average Syllables in Word

Gunning-Fog Readability Index

Character Space

Letter Space

Flesch Reading Ease Score

Feature Name

Feature Description

Feature Extractor

Text Pre-Processing

Feature Post-Processing

Normalization

Factor

Tools

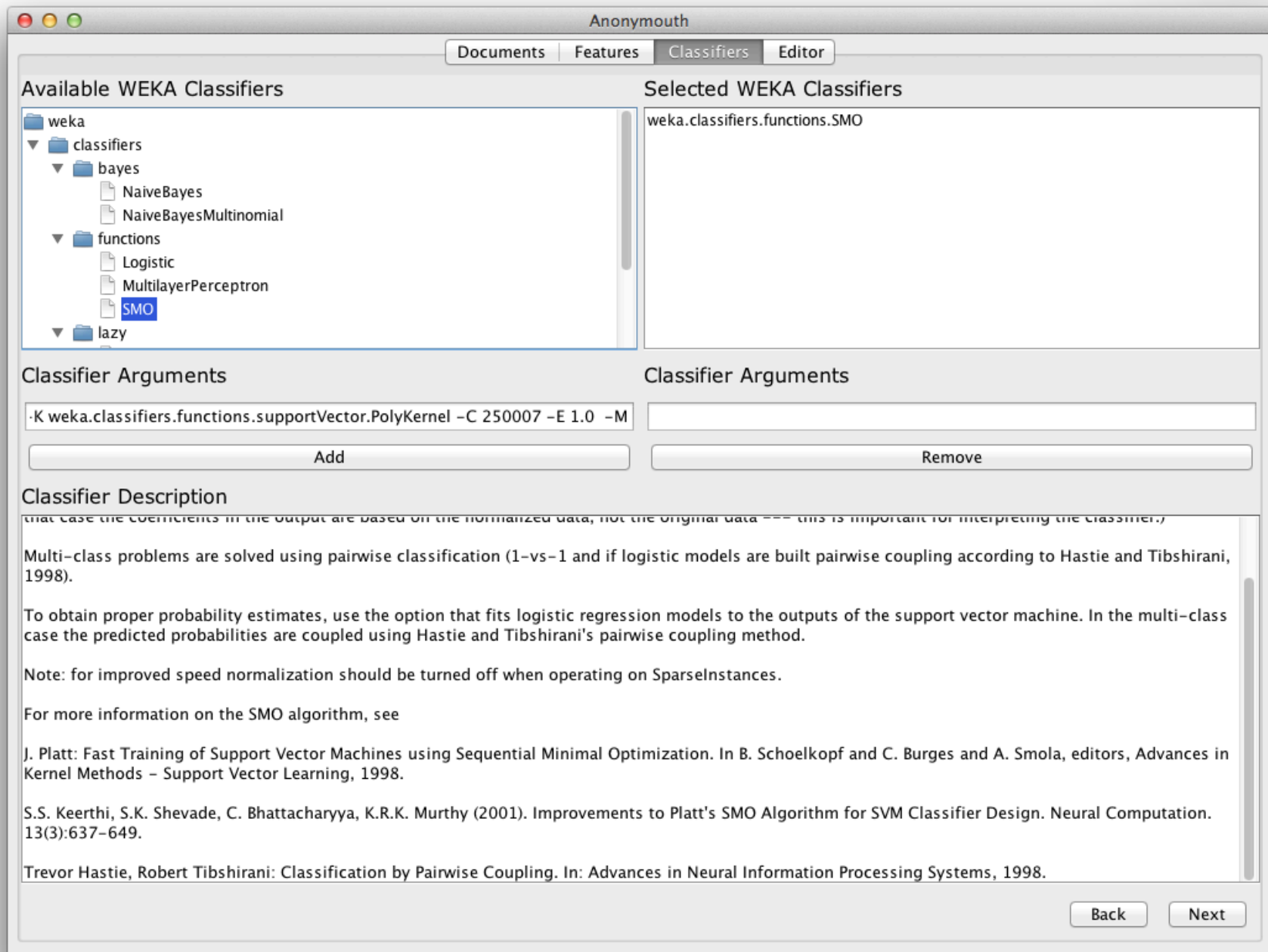
Configuration

Remove

Back

Next

(send bug reports, suggestions, questions to awm32@drexel.edu)



Documents

Features

Classifiers

Editor

Original

The inner-city, ethnic neighborhoods that are breeding grounds for street gang formation are the perfect conditions to create the gender role strain-conflict in many males (Vigil 48). Many aspects of a gang member's life prior to joining the gang are not prime conditions to create a healthy self-image of one's self, especially when attempting to fill traditional gender roles. This model suggests that males will engage in negative masculine acts, like violence and aggression, if they are not able to assert their masculinity in traditional male roles, like the roles of provider or protector (Vigil 24). In the case of many male members of the Latin Kings, they are unable to assert their masculinity as their culture dictates. In places like an inner-city neighborhood where finances are often a troubling aspect of life, males may struggle to care for their families financially and protect them from the violence inner-city life often contains.

Machismo, sometimes classified as a psychological trait, at its simplest, is the belief of male superiority. Machismo is meant to be acted out as a positive ideal, focusing on the protection of women, children, and those less fortunate (Vigil 37). In Hispanic culture, machismo is a defining aspect of a male's identity (Brezinka). Most men view this as a source of pride and honor, but Vigil suggests that when a man's basic function, as defined by machismo, is threatened by environmental factors (like an inner-city neighborhood), he must seek to assert his manhood and reconcile the psychological tension or cognitive dissonance he feels (Vigil 37-8, Anderson et. al, 19). The violence and aggression suggested in the gender role model plus the pressures of the machismo attitude are often the gateway for males to join the gang, as it satisfies their desire to be masculine. In an environment where their masculinity is not always visible, by instance to show their strength and power is welcome and often taken advantage of, regardless of the associated risks.

The strict gender roles and the overwhelmingly difficult environmental conditions create psychological conditions that are alleviated through gang membership. The inability to fulfill such gender roles creates a situation presented by Leon Festinger's role conflict dissonance. This phenomenon states that when one's beliefs or attitudes contradict their actual behavior, they must change one to match the other in order to avoid psychological strain (Myers 147).

Results of this Document's Classification (% probability of authorship per author)

d	~* you ~*	b	c
11.66	53.84	11.42	23.08
Actual Author:	~* you ~*		

Unfortunately, your document seems to have been written by: ~* you ~*

Suggestion

Feature Name:

Present Value: null

Target Value: null

List of Suggestions

No.	Feature Name
1	COMPLEXITY
2	CHARACTER_SPACE
3	UNIQUE_WORDS_COUNT
4	SENTENCE_COUNT
5	AVERAGE_SENTENCE_LENGTH
6	AVERAGE_SYLLABLES_IN_WORD
7	GUNNING_FOG_READABILITY_...
8	LETTER_SPACE
9	FLESCH_READING_EASE_SCORE

Editing Progress...

Re-process

Clear Highlights

Dictionary

Verbose

Save...

Exit

(send bug reports, suggestions, questions to awm32@drexel.edu)

Documents

Features

Classifiers

Editor

Original

Original -> 1

All these feelings of inadequacy and neglect, as stated, push children to seek emotional and social needs elsewhere, as well as to find a place where they feel a sense of security in an environment where security is not always guaranteed day-to-day (Vigil 23). These negative feelings, that are a result of the social order, work in the same way in which traditional Hispanic gender roles push males into gang membership. In order to protect themselves, whether it be physically, emotionally, or psychologically, gang members are constantly playing a demented game of survival of the fittest, finding the gang to be their best chance for survival (Emmer 30).

The inner-city, ethnic neighborhoods that are breeding grounds for street gang formation are the perfect conditions to create the gender role strain-conflict in many males (Vigil 48). Many aspects of a gang member's life prior to joining the gang are not prime conditions to create a healthy self-image of one's self, especially when attempting to fill traditional gender roles. This model suggests that males will engage in negative masculine acts, like violence and aggression, if they are not able to assert their masculinity in traditional male roles, like the roles of provider or protector (Vigil 24). In the case of many male members of the Latin Kings, they are unable to assert their masculinity as their culture dictates. In places like an inner-city neighborhood where finances are often a troubling aspect of life, males may struggle to care for their families financially and protect them from the violence inner-city life often contains.

Machismo, sometimes classified as a psychological trait, at its simplest, is the belief of male superiority. Machismo is meant to be acted out as a positive ideal, focusing on the protection of women, children, and those less fortunate (Vigil 37). In Hispanic culture, machismo is a defining aspect of a male's identity (Brezinka). Most men view this as a source of pride and honor, but Vigil suggests that when a man's basic function, as defined by machismo, is threatened by environmental factors (like an inner-city neighborhood), he must seek to assert his manhood and reconcile the psychological tension, or cognitive dissonance he feels (Vigil 37-8, Anderson et. al, 19). The violence and aggression suggested in the gender role model plus the pressures of the machismo attitude are often the gateway for males to join the gang, as it satisfies their desire to be masculine. In an environment where their masculinity is not always visible, any instance to show their strength and power is welcome and often taken.

Results of **Last** Document's Classification (% probability of authorship per author)

d	~* you ~*	b	c
11.66	53.84	11.42	23.08
Actual Author:	~* you ~*		

Unfortunately, your document seems to have been written by: ~* you ~*

Suggestion

It Appears as if the present value of this feature in your document is within the acceptable range of potential target values. However, because of the interrelated nature of stylometric features, it may be useful to re-visit this feature after you have made changes to other features.

Note: because of the clustering algorithm and method to choose the optimal target value used, it is possible that that a feature will naturally fall into this category.

CHARACTER_SPACE :

Present Value: 3161.0

Target Value: 3161.0

List of Suggestions

No.	Feature Name
1	COMPLEXITY
2	CHARACTER_SPACE
3	UNIQUE_WORDS_COUNT
4	SENTENCE_COUNT
5	AVERAGE_SENTENCE_LENGTH
6	AVERAGE_SYLLABLES_IN_WORD
7	GUNNING_FOG_READABILITY_...
8	LETTER_SPACE
9	FLESCH_READING_EASE_SCORE

Editing Progress...

Re-process

Clear Highlights

Dictionary

Verbose

Save...

Exit

(send bug reports, suggestions, questions to awm32@drexel.edu)

Documents

Features

Classifiers

Editor

Original

Original -> 1

All these feelings of inadequacy and neglect, as stated, push children to seek emotional and social needs elsewhere, as well as to find a place where they feel a sense of security in an environment where security is not always guaranteed day-to-day (Vigil 23). These negative feelings, that are a result of the social order, work in the same way in which traditional Hispanic gender roles push males into gang membership. In order to protect themselves, whether it be physically, emotionally, or psychologically, gang members are constantly playing a demented game of survival of the fittest, finding the gang to be their best chance for survival (Emmer 30).

The inner-city, ethnic neighborhoods that are breeding grounds for street gang formation are the perfect conditions to create the gender role strain-conflict in many males (Vigil 48). Many aspects of a gang member's life prior to joining the gang are not prime conditions to create a healthy self-image of one's self, especially when attempting to fill traditional gender roles. This model suggests that males will engage in negative masculine acts, like violence and aggression, if they are not able to assert their masculinity in traditional male roles, like the roles of provider or protector (Vigil 24). In the case of many male members of the Latin Kings, they are unable to assert their masculinity as their culture dictates. In places like an inner-city neighborhood where finances are often a troubling aspect of life, males may struggle to care for their families financially and protect them from the violence inner-city life often contains.

Machismo, sometimes classified as a psychological trait, at its simplest, is the belief of male superiority. Machismo is meant to be acted out as a positive ideal, focusing on the protection of women, children, and those less fortunate (Vigil 37). In Hispanic culture, machismo is a defining aspect of a male's identity (Brezinka). Most men view this as a source of pride and honor, but Vigil suggests that when a man's basic function, as defined by machismo, is threatened by environmental factors (like an inner-city neighborhood), he must seek to assert his manhood and reconcile the psychological tension, or cognitive dissonance he feels (Vigil 37-8, Anderson et. al, 19). The violence and aggression suggested in the gender role model plus the pressures of the machismo attitude are often the gateway for males to join the gang, as it satisfies their desire to be masculine. In an environment where their masculinity is not always visible, any instance to show their strength and power is welcome and often taken.

Results of **Last** Document's Classification (% probability of authorship per author)

d	~* you ~*	b	c
11.66	53.84	11.42	23.08
Actual Author:	~* you ~*		

Unfortunately, your document seems to have been written by: ~* you ~*

Suggestion

It Appears as if the present value of this feature in your document is within the acceptable range of potential target values. However, because of the interrelated nature of stylometric features, it may be useful to re-visit this feature after you have made changes to other features.

Note: because of the clustering algorithm and method to choose the optimal target value used, it is possible that that a feature will naturally fall into this category.

UNIQUE_WORDS_COUNT :

Present Value: 249.0

Target Value: 249.0

List of Suggestions

No.	Feature Name
1	COMPLEXITY
2	CHARACTER_SPACE
3	UNIQUE_WORDS_COUNT
4	SENTENCE_COUNT
5	AVERAGE_SENTENCE_LENGTH
6	AVERAGE_SYLLABLES_IN_WORD
7	GUNNING_FOG_READABILITY_...
8	LETTER_SPACE
9	FLESCH_READING_EASE_SCORE

Editing Progress...

Re-process

Clear Highlights

Dictionary

Verbose

Save...

Exit

(send bug reports, suggestions, questions to awm32@drexel.edu)

Documents

Features

Classifiers

Editor

Original

Original -> 1

All these feelings of inadequacy and neglect, as stated, push children to seek emotional and social needs elsewhere, as well as to find a place where they feel a sense of security in an environment where security is not always guaranteed day-to-day (Vigil 23). These negative feelings, that are a result of the social order, work in the same way in which traditional Hispanic gender roles push males into gang membership. In order to protect themselves, whether it be physically, emotionally, or psychologically, gang members are constantly playing a demented game of survival of the fittest, finding the gang to be their best chance for survival (Emmer 30).

The inner-city, ethnic neighborhoods that are breeding grounds for street gang formation are the perfect conditions to create the gender role strain-conflict in many males (Vigil 48). Many aspects of a gang member's life prior to joining the gang are not prime conditions to create a healthy self-image of one's self, especially when attempting to fill traditional gender roles. This model suggests that males will engage in negative masculine acts, like violence and aggression, if they are not able to assert their masculinity in traditional male roles, like the roles of provider or protector (Vigil 24). In the case of many male members of the Latin Kings, they are unable to assert their masculinity as their culture dictates. In places like an inner-city neighborhood where finances are often a troubling aspect of life, males may struggle to care for their families financially and protect them from the violence inner-city life often contains.

Machismo, sometimes classified as a psychological trait, at its simplest, is the belief of male superiority. Machismo is meant to be acted out as a positive ideal, focusing on the protection of women, children, and those less fortunate (Vigil 37). In Hispanic culture, machismo is a defining aspect of a male's identity (Brezinka). Most men view this as a source of pride and honor, but Vigil suggests that when a man's basic function, as defined by machismo, is threatened by environmental factors (like an inner-city neighborhood), he must seek to assert his manhood and reconcile the psychological tension, or cognitive dissonance he feels (Vigil 37-8, Anderson et. al, 19). The violence and aggression suggested in the gender role model plus the pressures of the machismo attitude are often the gateway for males to join the gang, as it satisfies their desire to be masculine. In an environment where their masculinity is not always visible, any instance to show their strength and power is welcome and often taken

SENTENCE_COUNT :

Present Value: 19.0

Target Value: 27.2789

List of Suggestions

No.	Feature Name
1	COMPLEXITY
2	CHARACTER_SPACE
3	UNIQUE_WORDS_COUNT
4	SENTENCE_COUNT
5	AVERAGE_SENTENCE_LENGTH
6	AVERAGE_SYLLABLES_IN_WORD
7	GUNNING_FOG_READABILITY_...
8	LETTER_SPACE
9	FLESCH_READING_EASE_SCORE

Results of **Last** Document's Classification (% probability of authorship per author)

d	~* you ~*	b	c
11.66	53.84	11.42	23.08

Actual Author: ~* you ~*

Unfortunately, your document seems to have been written by: ~* you ~*

Editing Progress...

Re-process

Clear Highlights

Dictionary

Verbose

Save...

Exit

(send bug reports, suggestions, questions to awm32@drexel.edu)

Documents

Features

Classifiers

Editor

Original

Original->1

All these feelings of inadequacy and neglect, as stated, push children to seek emotional and social needs elsewhere, as well as to find a place where they feel a sense of security in an environment where security is not always guaranteed day-to-day (Vigil 23). These negative feelings, that are a result of the social order, work in the same way in which traditional Hispanic gender roles push males into gang membership. In order to protect themselves, whether it be physically, emotionally, or psychologically, gang members are constantly playing a demented game of survival of the fittest, finding the gang to be their best chance for survival (Emmer 30).

The inner-city, ethnic neighborhoods, that are breeding grounds for street gang formation are the perfect conditions to create the gender role strain-conflict in many males (Vigil 48). Many aspects of a gang member's life prior to joining the gang are not prime conditions to create a healthy self-image of one's self, especially when attempting to fill traditional gender roles. This model suggests that males will engage in negative masculine acts, like violence and aggression, if they are not able to assert their masculinity in traditional male roles, like the roles of provider or protector (Vigil 24). In the case of many male members of the Latin Kings, they are unable to assert their masculinity as their culture dictates. In places like an inner-city neighborhood where finances are often a troubling aspect of life, males may struggle to care for their families financially and protect them from the violence inner-city life often contains.

Machismo, sometimes classified as a psychological trait, at its simplest, is the belief of male superiority. Machismo is meant to be acted out as a positive ideal, focusing on the protection of women, children, and those less fortunate (Vigil 37). In Hispanic culture, machismo is a defining aspect of a male's identity (Brezinka). Most men view this as a source of pride and honor, but Vigil suggests that when a man's basic function, as defined by machismo, is threatened by environmental factors (like an inner-city neighborhood), he must seek to assert his manhood and reconcile the psychological tension, or cognitive dissonance, he feels (Vigil 37-8, Anderson et. al, 19). The violence and aggression suggested in the gender role model plus the pressures of the machismo attitude are often the gateway for males to join the gang, as it satisfies their desire for masculine. In an environment where their masculinity is not always visible, any instance to show their strength and power is welcome and often taken

Suggestion

The present value of this feature is: '19.0' and it should be increased to '26.58081363101645'. Helpful insights (and possibly highlighting) that will aid in accomplishing this task are (or will soon be) in the works.

Feature Name:

Present Value:

null

Target Value:

null

List of Suggestions

No.	Feature Name
1	COMPLEXITY
2	CHARACTER_SPACE
3	UNIQUE_WORDS_COUNT
4	SENTENCE_COUNT
5	AVERAGE_SENTENCE_LENGTH
6	AVERAGE_SYLLABLES_IN_WORD
7	GUNNING_FOG_READABILITY...
8	LETTER_SPACE
9	FLESCH_READING_EASE_SCORE

Results of this Document's Classification (% probability of authorship per author)

d	~* you ~*	b	c
6.59	15.91	41.33	36.16
Actual Author:	~* you ~*		

Your document appears as if 'b' wrote it!

Editing Progress...

Re-process

Clear Highlights

Dictionary

Verbose

Save...

Exit

(send bug reports, suggestions, questions to awm32@drexel.edu)

Documents

Features

Classifiers

Editor

Original

Original -> 1

Two years later, in November of 1964, Hoover held a rare press conference in which he called King "the most notorious liar in the country," creating a public relations debacle. In fact, he had also called King "one of the lowest characters in the country" and "gave an affirmative response" to a question regarding the rumored links between King and communism. This led to a meeting between Hoover and King in early December at King's request, where publicly they appeared to resolve their differences and emerge in a tenuous state of cooperation. While King emerged feeling as though he was simply talked at, Hoover's subordinates actively sought to meet with other African American leaders to discuss replacing King. The news article discussing the meeting notes that Johnson had waived the mandatory retirement age of 70 that Hoover was soon approaching. Hoover had been in power for decades with no plans of retiring, and even though Johnson had waived the rule Hoover saw his position as the weakest in years, perhaps helping to explain his willingness to make public amends with his foremost enemy at the time.

While the public's knowledge of the activities of Hoover and King in 1964 and 1965 was limited to this and other minor comments in the press, this was the year when Hoover switched tactics and tried to destroy King personally instead of only through linking him to communism. The wiretapping revealed little about King's ties with communists, but did record King's sexual promiscuity in detail. In a January 1964 memo, Hoover ordered surveillance in Milwaukee although his agents believed visible police presence would dissuade King from engaging in sexual activity, reasoning, "I don't share the conjecture. King is a 'tom cat' with obsessive degenerate sexual urges." Hoover believed that evidence of King's immoral behavior would discredit him, ending the threat. Threats and blackmail related to an individual's private behavior, particularly when it came to immorality or, worse yet, homosexuality, was both an often-used and effective tool for Hoover. One oft-cited rumor surrounding Hoover was that he himself might be homosexual, the evidence for this being his unusual relationship with his number two at the FBI, Clyde Tolson. However, none of his usual news sources were interested in using the information, and it did not have the desired effect on other civil rights leaders or public figures.

The culmination of the personal attacks, while most likely not directly involving Hoover though there was little in the FBI he was unaware of, was a particularly ominous package sent by the FBI, under Sullivan's orders, to King in November 1964. Opened on January 5th, 1965, by King's wife Coretta, the package contained a solicited-

Suggestion

The highlighted words represent the top two groups of words with a syllable count greater than, or equal to, '1' You should make an effort to break up as many of these words as possible (e.g. instead of 'intelligently' use 'in a smart way').

AVERAGE_SYLLABLES_IN_WORD :

Present Value: 1.7729

Target Value: 1.5736

List of Suggestions

No.	Feature Name
1	COMPLEXITY
2	AVERAGE_SENTENCE_LENGTH
3	CHARACTER_SPACE
4	SENTENCE_COUNT
5	GUNNING_FOG_READABILITY_...
6	UNIQUE_WORDS_COUNT
7	AVERAGE_SYLLABLES_IN_WORD
8	LETTER_SPACE
9	FLESCH_READING_EASE_SCORE

Results of **Last** Document's Classification (% probability of authorship per author)

d	~* you ~*	b	c
12.12	78.86	8.02	1.0

Actual Author: ~* you ~*

Unfortunately, your document seems to have been written by: ~* you ~*

Editing Progress...

Re-process

Clear Highlights

Dictionary

Verbose

Save...

Exit

(send bug reports, suggestions, questions to awm32@drexel.edu)

Anonymouth Challenges

- Features are often not independent.
 - Increasing the number of complex words will also increase average syllable count.
 - Reducing the number of times a specific word occurs will also affect the lexical density.
- How can we create an algorithm for anonymity that generates an obfuscated document with minimal effort and without circular feature modification?

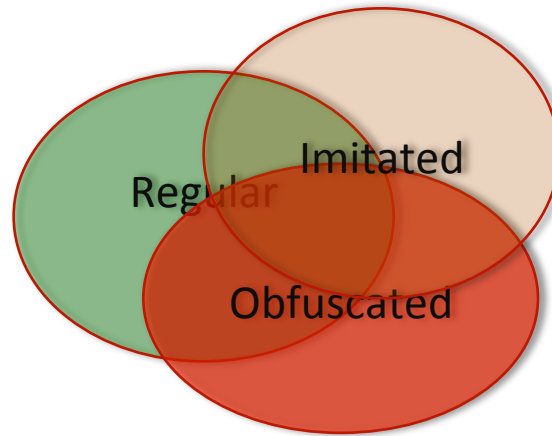
Anonymouth Dev Goals

- Streamlined suggestion system.
 - Improved automation on applicable features.
 - Improved clustering algorithm to provide optimal path to anonymity.
- Improved editing interface.
 - Increased phrase and word synonym set support.
 - Edit by blocks of text, not simply feature-by-feature.
- Wider set of features and classification methods.
 - Multi-method and feature collection analysis.
- Usability and anonymity user studies.

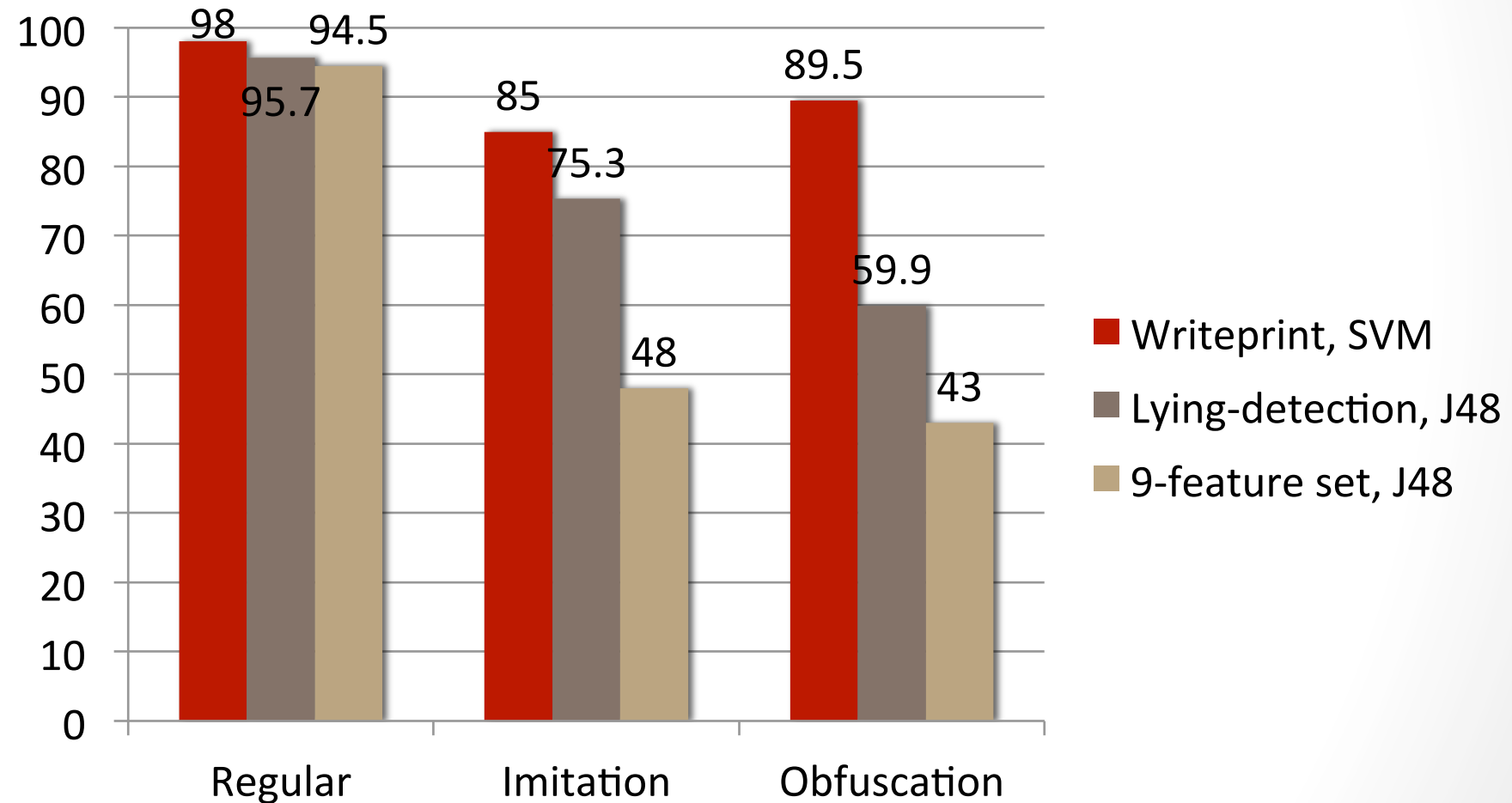
Opening Development

- Project will continue to be developed by PSAL at Drexel, but we welcome collaboration and participation.
- We are interested in...
 - Linguistic Experts
 - Security Advisors
 - UI Experts

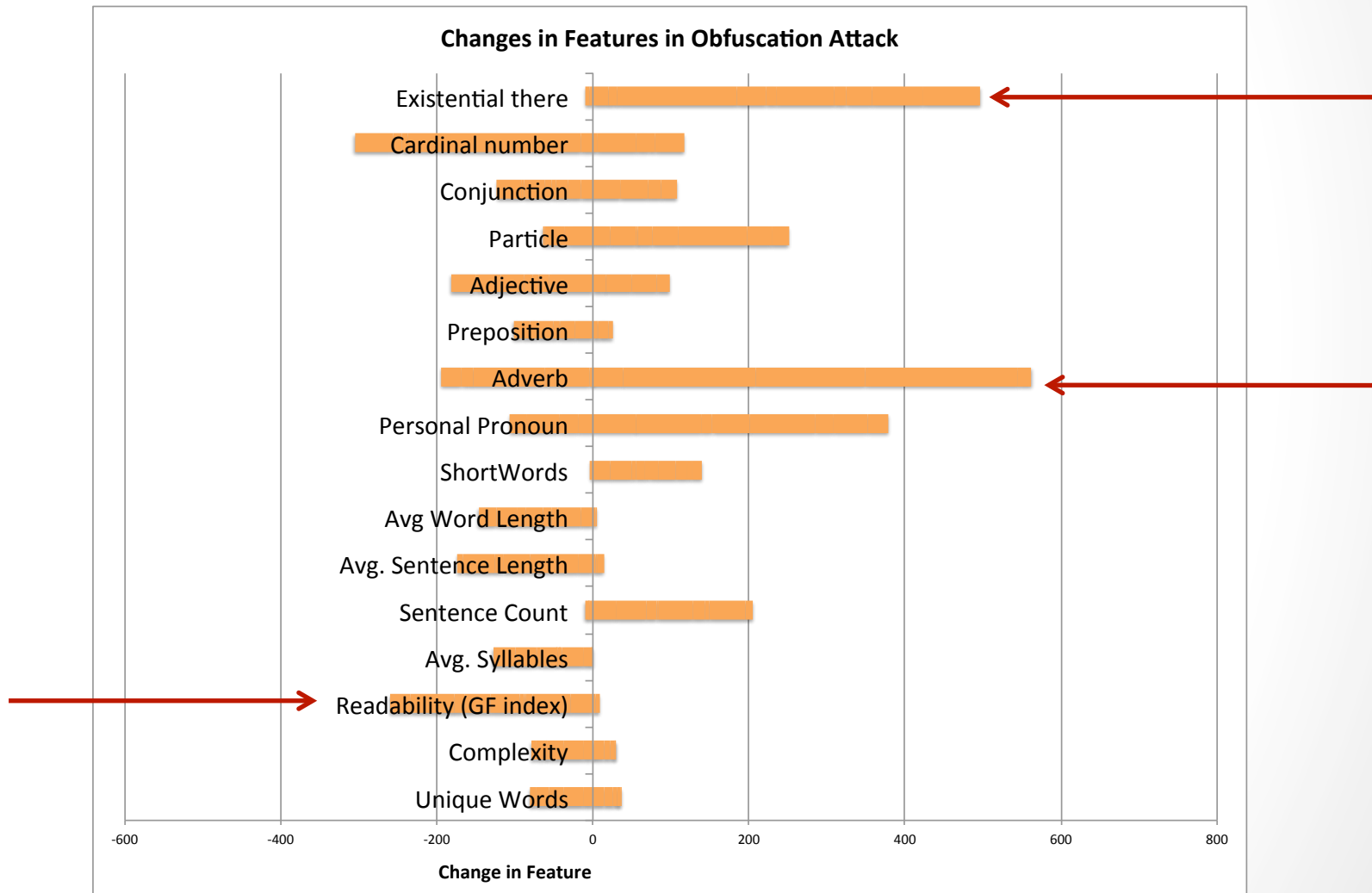
Can we detect stylistic deception?



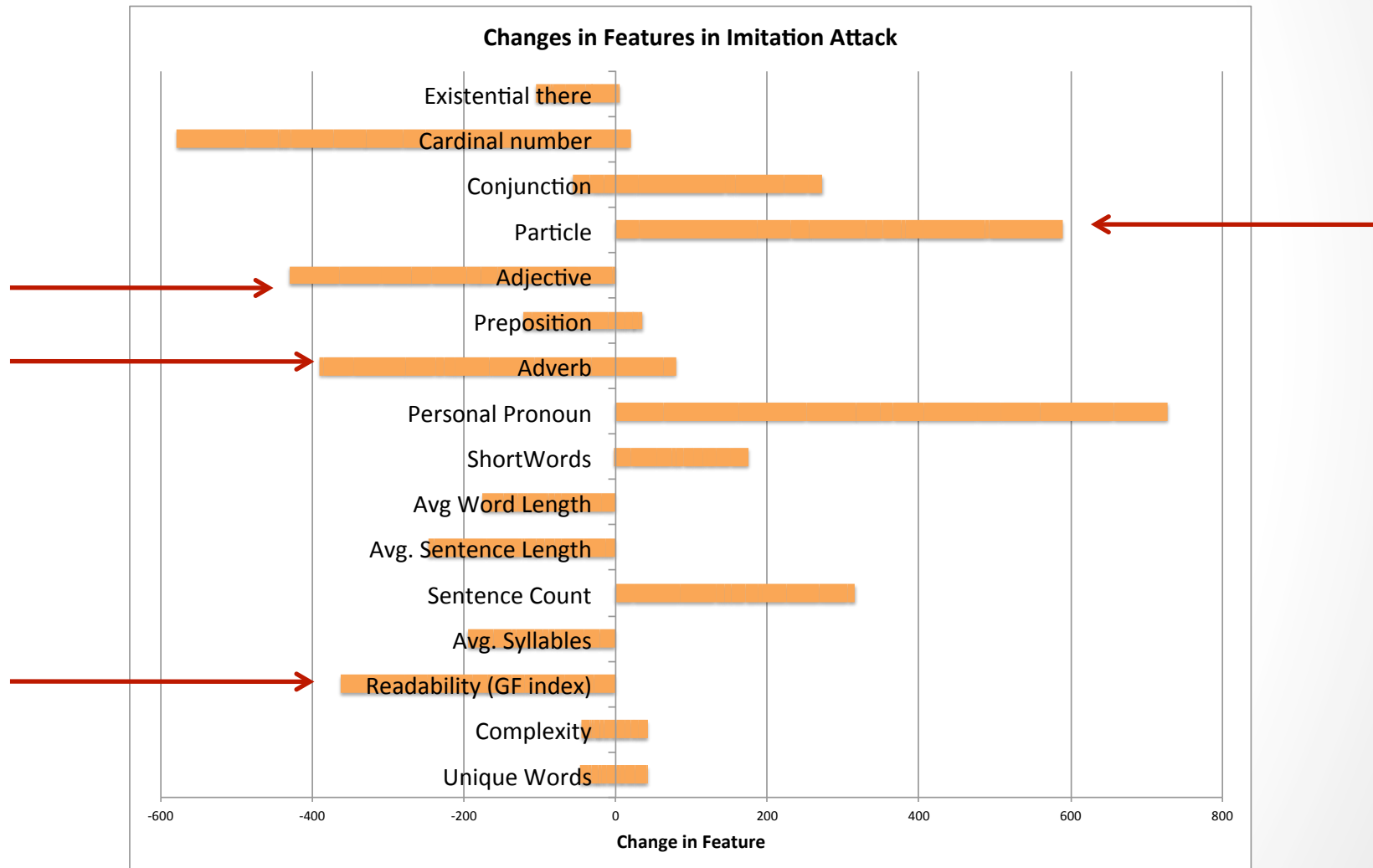
Detecting stylistic deception is possible



Feature Changes in Obfuscated Passages

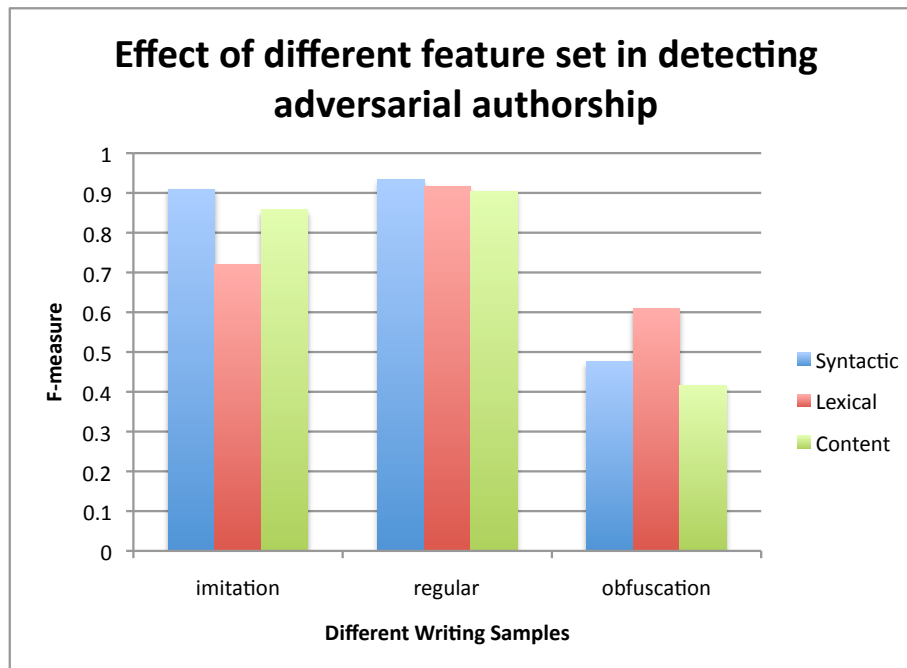


Feature Changes in Imitated Passages



Problem with the dataset: Topic Similarity

- All the deceptive documents were of same topic.
- Non-content-specific features have same effect as content-specific features.



Hemingway-Faulkner Imitation Corpus

- Articles from the International Imitation Hemingway Contest (2000-2005)
- Articles from the Faux Faulkner Contest (2001-2005)
- Original excerpts of Ernest Hemingway and William Faulkner

Deception detection is possible even when the topic is not similar

- 81.2% accurate in detecting imitated documents.

Long term deception

- A Gay Girl In Damascus blog:
 - Original author was a 40-year old American citizen, Thomas MacMaster.
 - Pretended to be a Syrian gay woman, Amina Arraf.
 - The author worked for at least 5 years to create a new style.

Long term deception is hard to detect

- None of the blog posts were found to be deceptive.
- But regular authorship recognition can help.
- We tried to attribute authorship of the blog posts using Thomas (as himself), Thomas (as Amina), Britta (Thomas's wife).
- 54.3% of the blog posts were attributed to Thomas (as himself)

Recap

- Available Now:
 - Brennan-Greenstadt Adversarial Stylometry Corpus (12 Authors)
 - Drexel AMT Adversarial Stylometry Corpus (45 Authors)
 - JStylo Alpha Release
 - Anonymouth Alpha Release
- Future Work:
 - Beta releases of JStylo and Anonymouth
 - Academic publication of new results
 - Continued analysis of deception detection and short message classification
 - Continued research on improving partially automated anonymization

Thanks.

- We want to hear from you.
 - Mike Brennan (mb553@drexel.edu)
 - Rachel Greenstadt (greenie@cs.drexel.edu)
 - Ariel Stolerman, JStylo Lead (ariels@drexel.edu)
 - Andrew McDonald, Anonymouth Lead (ams23@drexel.edu)
 - Sadia Afroz, Deception Detection Lead (sa499@drexel.edu)
 - Aylin Caliskan, Translation & Stylometry (ac993@drexel.edu)
- PSAL: <https://psal.cs.drexel.edu>
- **We are looking for interested grad students and post-docs!**

Addendum Slides

Research Questions, Practical Implications.

- Our upcoming research questions have substantial practical implications.
- How do you anonymize a document sufficiently in a reasonable period of time?
 - What is *sufficient*? What is *reasonable*?
- Can Anonymouth be used to successfully imitate other authors?
- Can Anonymouth maintain long-term deception? Can its usage be detected?
- JStylo vs. Anonymouth – who wins?
 - Based on JStylo, Anonymouth will have everything it needs to help evade detection by the methods it contains.

Two Tools?

- Aren't we creating a tool that enables surveillance and de-anonymization?
 - Anonymouth can't exist without JStylo. But it also shows that you can't necessarily depend on stylometry to assign authorship.
 - JStylo allows for easier use of authorship recognition tools, but is extensible and open-source. Implementing a method in JStylo will enable counter-attacks in Anonymouth.
- JStylo vs. Anonymouth – who wins?
 - Based on JStylo, Anonymouth will have everything it needs to help evade detection by the methods it contains.
 - Note that nothing prevents others from plugging in proprietary stylometric methods into their version of JStylo.