



A Taste of Applied Machine Learning

Carolyn Penstein Rosé

*Language Technologies Institute/
Human-Computer Interaction
Institute*

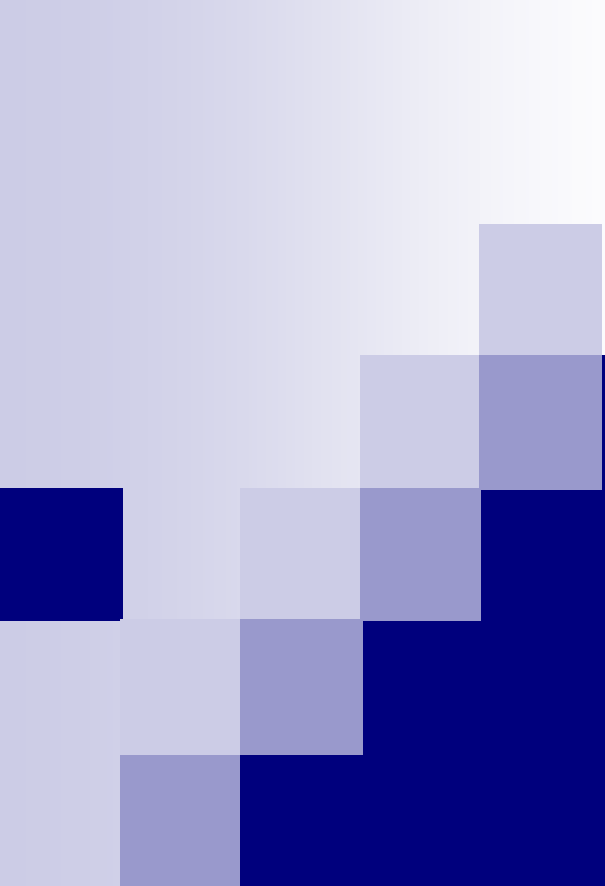
Carolyn Rosé



PhD in Language and
Information Technologies, 1998

*Enjoys Israeli folk dancing,
playing piano, long walks in the
woods, and knitting, crocheting,
and spinning yarn*

- Joint appointment between
Language Technologies and
Human-Computer Interaction
- President of the International
Society of the Learning Sciences
- Co-Chair of the CSCL community
committee
- Associate Editor of the
International Journal of Computer
Supported Collaborative Learning



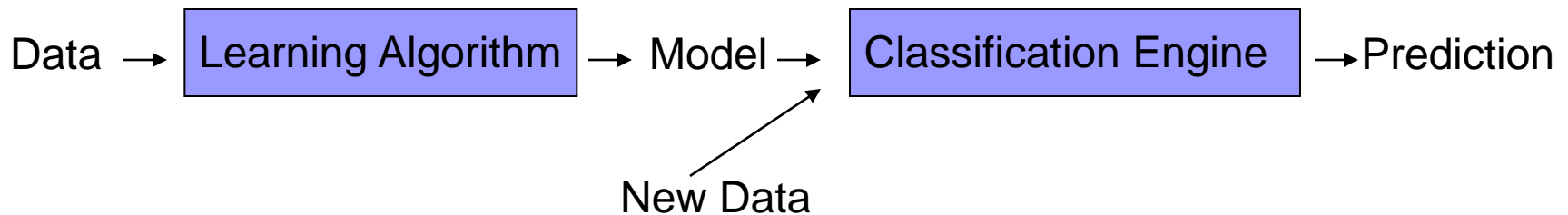
Machine Learning: Conceptual Overview

How does machine learning work?

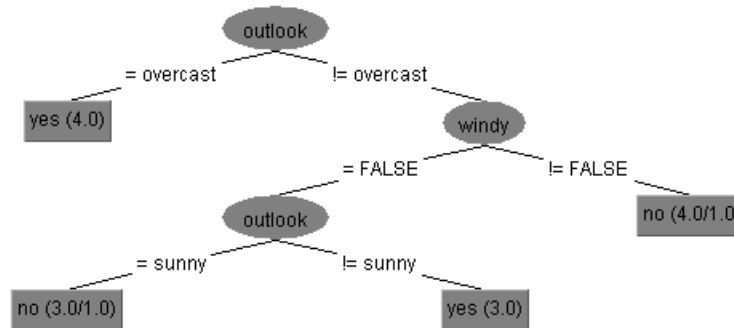
<u>outlook</u>	<u>temperature</u>	<u>humidity</u>	<u>windy</u>	<u>play</u>
sunny	hot	high	FALSE	no
sunny	hot	high	TRUE	no
<div> <div> A slightly different combination of features will find the most informative split. What do you think?</div> <div> <div>Outlook:</div> <div>Sunny -> No</div> <div>Overcast -> Yes</div> <div>Rainy-> Yes</div> </div> <div> <div>Class</div> <div>yes</div> <div>yes</div> <div>yes</div> <div>yes</div> <div>yes</div> <div>no</div> </div> <div> <div>What value?</div> <div>no</div> <div>yes</div> <div>yes</div> <div>yes</div> <div>yes</div> <div>no</div> </div> </div>				
sunny	moderate	normal	FALSE	yes
rainy	cool	normal	FALSE	yes
sunny	hot	normal	FALSE	yes
overcast	any	any	FALSE	yes
rainy	any	any	TRUE	no

What is machine learning?

- Automatically or *semi-automatically*
 - Inducing concepts (i.e., rules) from data
 - Finding patterns in data
 - Explaining data
 - Making predictions



More Complex Algorithm...



* Only makes 2 mistakes!



Why is it better?

- Not because it is more complex
 - Sometimes more complexity makes performance worse
- What is different in what the three rule representations assume about your data?
 - 0R
 - 1R
 - Trees
- The best algorithm for your data will give you exactly the power you need

Why is it better?

- Not b

- ☐ Som
- perf

- What
- repres

- ☐ 0R
- ☐ 1R
- ☐ Tree

- The b
- you e

Lets say you don't know the shape, what shape would you guess?



Why is it better?

- Not b

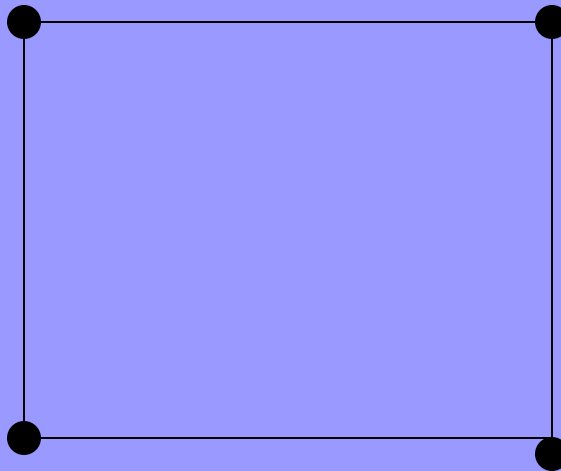
- ☐ Som
- perf

- What
- repres

- ☐ 0R
- ☐ 1R
- ☐ Tree

- The b
- you e

Lets say you don't know the shape, what shape would you guess?



Why is it better?

- Not b

- Som
 - perf

- What
- repres

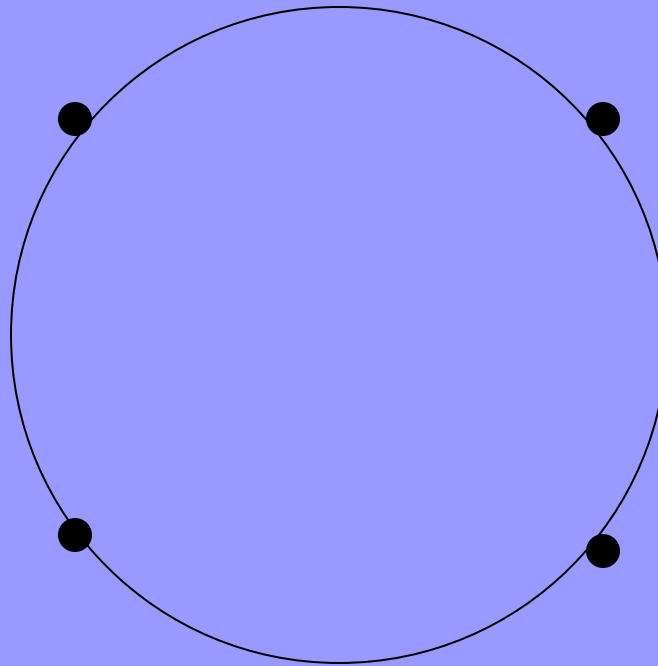
- 0R

- 1R

- Tree

- The b
- you e

But what if I told you it's really a circle?



Why is it better?

- Not b

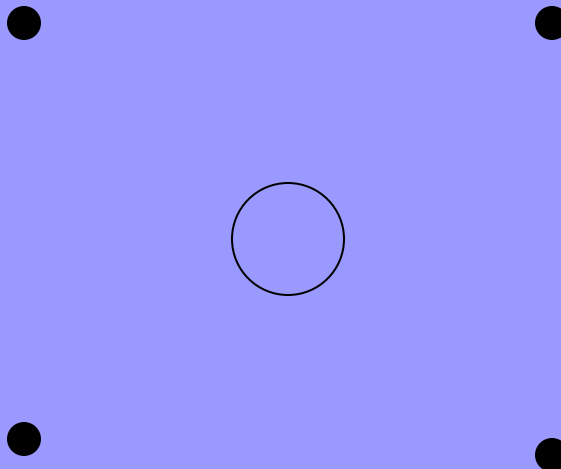
- Som
 - perf

- What
- repres

- 0R
 - 1R
 - Tree

- The b
- you e

If you know the shape, you have fewer degrees of freedom – less room to make a mistake.



Why is it better?

- Not b

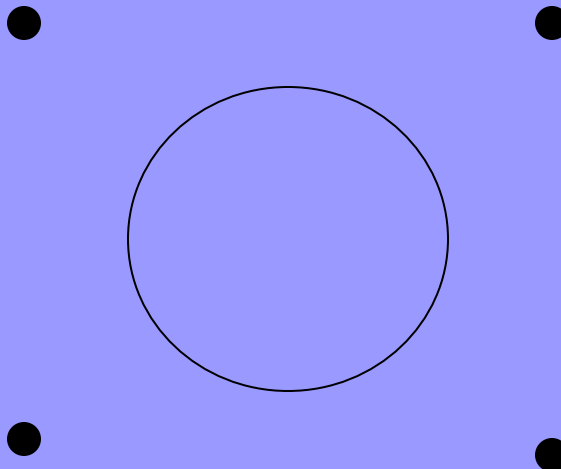
- Som
 - perf

- What
- repres

- 0R
 - 1R
 - Tree

- The b
- you e

If you know the shape, you have fewer degrees of freedom – less room to make a mistake.



Why is it better?

- Not b

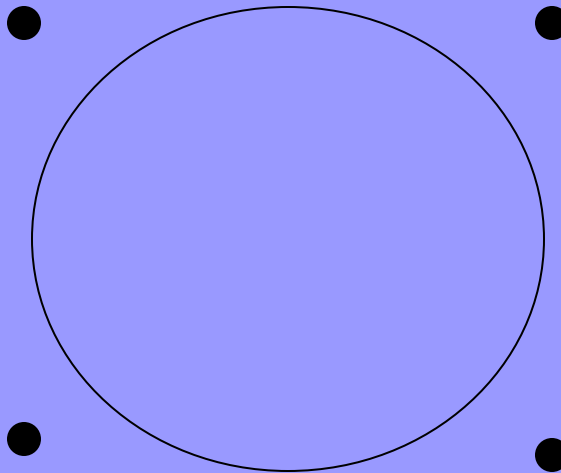
- Som
 - perf

- What
- repres

- 0R
 - 1R
 - Tree

- The b
- you e

If you know the shape, you have fewer degrees of freedom – less room to make a mistake.



Why is it better?

- Not b

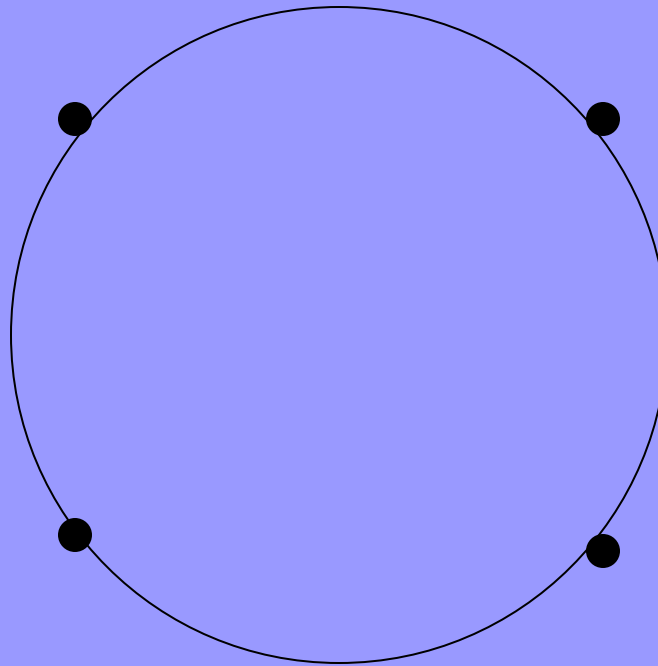
- Som
 - perf

- What
- repres

- 0R
 - 1R
 - Tree

- The b
- you e

If you know the shape, you have fewer degrees of freedom – less room to make a mistake.





Why is it better?

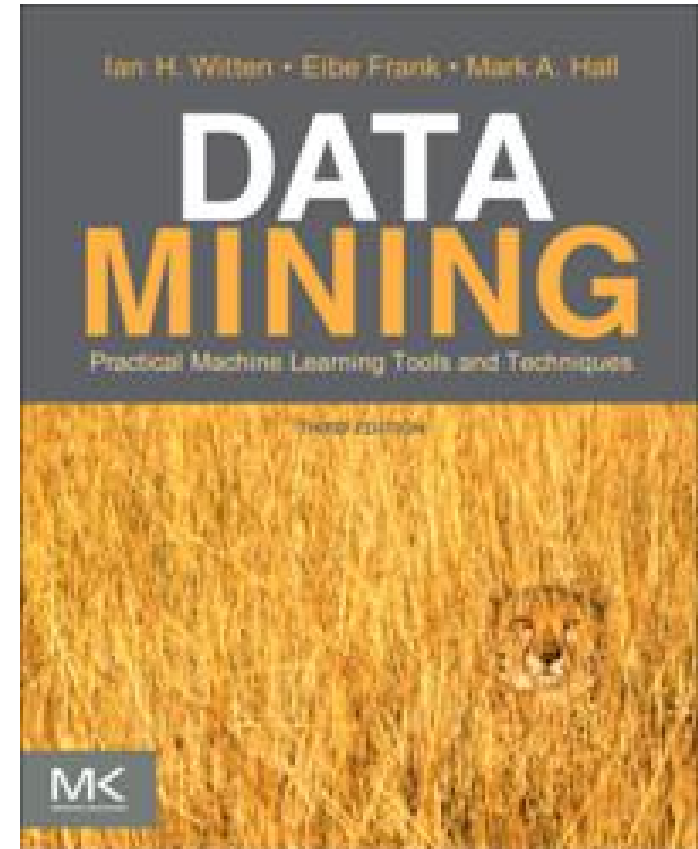
- Not because it is more complex
 - Sometimes more complexity makes performance worse
- What is different in what the three rule representations assume about your data?
 - 0R
 - 1R
 - Trees
- The best algorithm for your data will give you exactly the power you need



Tools and Resources

Essential Reading

- Witten, I. H., Frank, E., Hall, M. (2011). *Data Mining: Practical Machine Learning Tools and Techniques*, third edition, Elsevier: San Francisco





Other Suggested Readings

- Richard Cotton (2013). *Learning R*, O'Reilly and Associates
- Allen Downey (2013). *Think Bayes*, O'Reilly and Associates
- Mark Lutz (2013). *Learning Python*, O'Reilly and Associates
- Drew Conway & John White (2012). *Machine Learning for Hackers: Case Studies and Algorithms to Get You Started*, O'Reilly Media

Software Tools

- Data manipulation tools
 - Whatever you are comfortable with
 - Scripting language like R, Python, Perl
 - Excel
- Weka (<http://www.cs.waikato.ac.nz/ml/weka/>)
 - Open source machine learning toolkit
 - Includes Java API
- LightSIDE (<http://lightsidelabs.com/research>)
 - Weka add-on for text processing
 - Developed at CMU!

lightsidelabs.com/research/

LightSIDE Labs | Mac x

lightsidelabs.com/research/

LightSIDE Labs

Home Machine Learning Assessment Tutorials Hiring Contact Us

Search

Download LightSIDE Now

Latest Update 5.10.2013

We're happy that you're interested in using LightSIDE as part of your research! Here's what the freely available and completely GPLv3 open source version of LightSIDE does for you:

Easy, Fast Feature Extraction

CSV Files:
sentiment_sentences
DOCUMENT_LIST
Documents: sentiment_sentences
Class: <class>
Text Fields:
text

Feature Extractor Plugins:
Basic Features
Column Features
Previous Label Features

Configure Basic Features
Unigrams
Bigrams
Trigrams
POS Bigrams
Line Length
Contains Non-Stopwords
Binary N-grams?
Include Punctuation?
Remove Stopwords?
Stem N-grams?
Differentiate text columns?

Extract
Name: Features Rare Threshold: 5

Feature Table
Features: 2
Documents: sentiment_sentences
Feature Plugins
Feature Table Features

Evaluations to Display:
Target: pos
Basic Table Statistics
Target Hits
Precision
Total Hits
Correlation
F-Score
AUC

Features in Table:
Search:
Features: 2
Ranks: 2
Target Hits: 2

Recent Posts

[Survey: LightSIDE satisfaction](#)
[Tutorial: Quick Start to Error Analysis](#)
[Tutorial: Regular Expression and Stretchy Pattern Features](#)
[Tutorial: Saving and Exporting Results](#)
[Tutorial: Model Evaluation Settings](#)

Recent Comments

[Rick Weinberg on Contact Us](#)
[Rick Weinberg on Machine Learning](#)
[Amar Nath on Tutorial: Preparing My Data](#)
[elijah on Contact Us](#)
[Verb Washington on Contact Us](#)

Archives

[April 2013](#)
[February 2013](#)

Categories

[Uncategorized](#)

Meta

[Log in](#)
[Entries RSS](#)

Windows Explorer window showing the contents of the **Documents** library for the **LightSide** user.

Address Bar: Libraries > Documents > LightSide >

Menu Bar: File Edit View Tools Help

Toolbar: Organize Share with Burn New folder

Left Pane (Navigation):

- ★ Favorites
 - Desktop
 - _MACOSX
 - Recent Places
 - Dropbox
 - Google Drive
 - Documents
- Libraries
 - Documents** (selected)
 - Music
 - Pictures
 - Subversion
 - Videos
- Computer
 - OSDisk (C:)
- Network

Right Pane (Details):

Documents library
LightSide

Arrange by: Folder

Name	Date modified
LightSide_20130823	9/13/2013 8:44 AM
wekafiles	9/13/2013 8:44 AM
wekarefs	9/13/2013 8:44 AM
toolkits	9/13/2013 8:43 AM
tests	9/13/2013 8:43 AM
LightSide.app	9/13/2013 8:43 AM
plugins	9/13/2013 8:43 AM
src	9/13/2013 8:43 AM
bin	9/13/2013 8:43 AM
copyright	9/13/2013 8:43 AM
data	9/13/2013 8:43 AM
lib	9/13/2013 8:43 AM
saved	8/20/2013 2:35 PM
lightside_log	9/13/2013 4:02 PM
LightSide	9/13/2013 8:45 AM
run	8/20/2013 2:42 PM
LightSide_Researchers_Manual	8/4/2013 2:26 AM

What is machine learning

■ Algorithms?



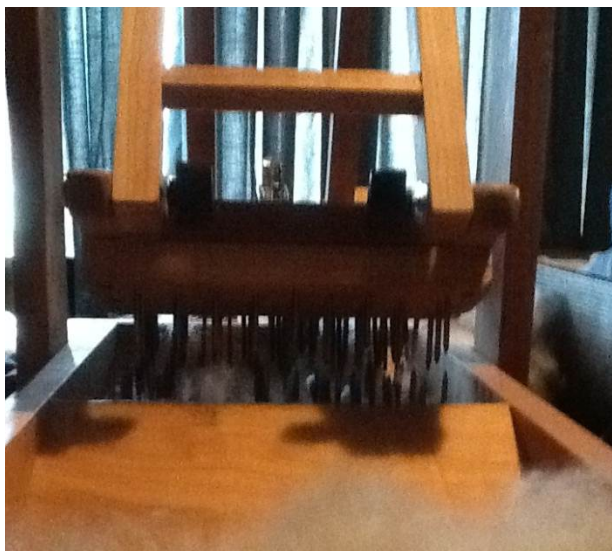
Ch. 48 (multiple of 7 plus 6)

1. Right Side Row. Hdc in third ch from hook and in each remaining ch. Ch.1, turn.
2. And ALL Wrong Side Rows. Sc in each stitch across. Ch.2, turn.
3. *Work hdc in 4 sc, work post dc as follows: yo, insert hook right to left under post of hdc below next sc, draw up loop, (yo and draw through 2 loops) twice to complete dc. Skp sc behind post dc, work hdc in next sc, work post dc in hdc below next sc*, skip sc behind post dc; repeat from * to * across, ending with hdc in last 4 sc. Ch 1, turn.
4. Wrong side row, repeat row 2.
5. *Hdc in 4 sc, skip first post dc, and work post dc on second post dc. Ch. 1, now work post dc around the skipped post dc (crossover made, skip the 3 sc behind crossover *; repeat from *to* across, ending with hdc in last 4 sc.
6. Wrong side row, repeat row 2.
7. * Hdc in 4 sc, work post dc in post dc below next sc (always skip sc behind each post dc), hdc in next sc, post dc in post dc below next sc*. Repeat from * to * ending with hdc in last 4 sc. Ch 1, turn. Repeat Rows 2 through 7 for pattern.





Sampling, Cleaning, Reformatting...





The rectangular batts
(like csv files with raw texts)
that come out of
the drum carder
are still not ready
to make something out of

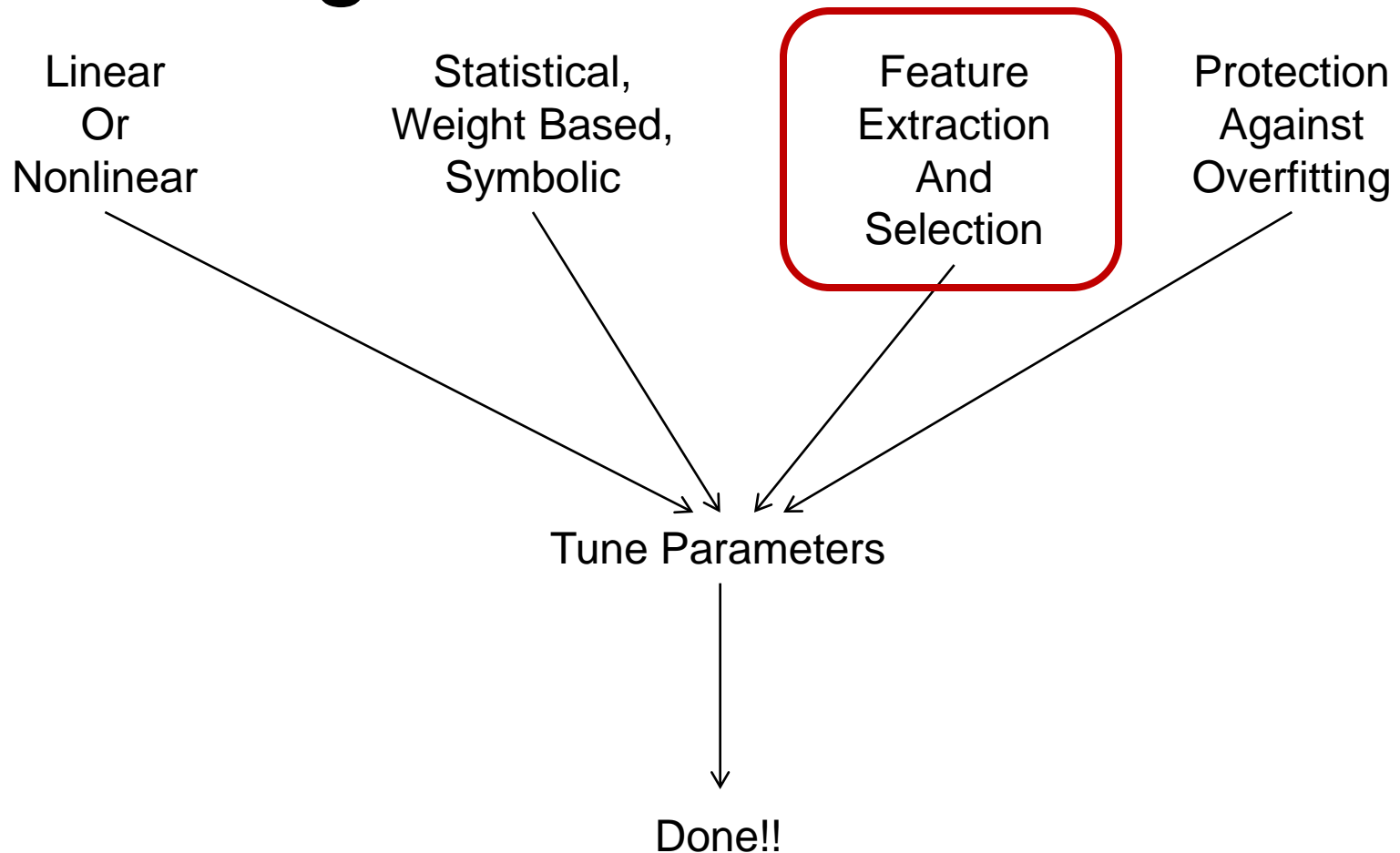
Here's where you come in!





Text Teaser

Decisions about Machine Learning Methods



Consider this simple example...

SimpleExample.xls		
	A	B
1	Code	text
2	Question	Tell me what your favorite color is.
3	Statement	I tell you my favorite color is blue.
4	Question	Where do you live?
5	Statement	I live where my family lives.
6	Question	Which kinds of baked goods do you prefer
7	Statement	I prefer to eat wheat bread.
8	Question	Which courses should I take?
9	Statement	You should take my applied machine learning course.
10	Question	Tell me when you get up in the morning.
11	Statement	I get up early.

Look for what distinguishes Questions and Statements in this dataset.

What clues do you see?

What are good features for text categorization?

SimpleExample.xls		
	A	B
1	Code	text
2	Question	Tell me what your favorite color is.
3	Statement	I tell you my favorite color is blue.
4	Question	Where do you live?
5	Statement	I live where my family lives.
6	Question	Which kinds of baked goods do you prefer?
7	Statement	I prefer to eat wheat bread.
8	Question	Which courses should I take?
9	Statement	You should take my applied machine learning course.
10	Question	Tell me when you get up in the morning.
11	Statement	I get up early.

What distinguishes Questions and Statements?

Not all questions end in a question mark.

What are good features for text categorization?

SimpleExample.xls		
	A	B
1	Code	text
2	Question	Tell me what your favorite color is.
3	Statement	I tell you my favorite color is blue.
4	Question	Where do you live?
5	Statement	I live where my family lives.
6	Question	Which kinds of baked goods do you prefer
7	Statement	I prefer to eat wheat bread.
8	Question	Which courses should I take?
9	Statement	You should take my applied machine learning course.
10	Question	Tell me when you get up in the morning.
11	Statement	I get up early.

What distinguishes Questions and Statements?

I versus you is not a reliable predictor

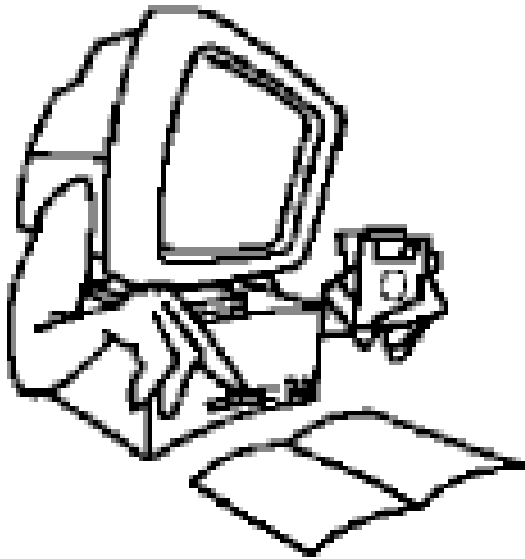
What are good features for text categorization?

SimpleExample.xls		
	A	B
1	Code	text
2	Question	Tell me <u>what</u> your favorite color is.
3	Statement	I tell you my favorite color is blue.
4	Question	<u>Where</u> do you live?
5	Statement	I live <u>where</u> my family lives.
6	Question	<u>Which</u> kinds of baked goods do you prefer
7	Statement	I prefer to eat wheat bread.
8	Question	<u>Which</u> courses should I take?
9	Statement	You should take my applied machine learning course.
10	Question	Tell me <u>when</u> you get up in the morning.
11	Statement	I get up early.

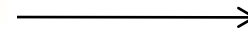
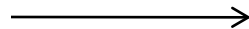
What distinguishes Questions and Statements?

Not all WH words occur in questions

Effective data representations make problems learnable...



- Machine learning isn't magic
- But it ***can*** be useful for identifying meaningful patterns in your data when used properly
- Proper use requires insight into your data





LightSIDE: A quick tour

CSV Files:

Class:

Type:

NOMINAL

Text Fields:

☐ Differentiate Text Fields

Feature Extractor Plugins:

- ☒ Basic Features
- ☐ Character N-Grams
- ☐ Column Features
- ☐ Parse Features
- ☐ Regular Expressions
- ☐ Stretchy Patterns

Configure Basic Features

- ☒ Unigrams
- ☐ Bigrams
- ☐ Trigrams
- ☐ POS Bigrams
- ☐ Word/POS Pairs
- ☐ Line Length
- ☐ Contains Non-Stopwords
- ☒ Binary N-grams?
- ☒ Include Punctuation?
- ☐ Remove Stopwords?
- ☐ Stem N-grams?

☐ Extract

Name:

features

Rare Threshold:

5

Feature Table:

Evaluations to Display:

Target:

Basic Table Statistics

- ☐ Recall
- ☐ Target Hits
- ☐ Precision
- ☐ Total Hits
- ☐ Correlation
- ☐ F-Score
- ☐ Kappa

Features in Table:

Search:

CSV Files:

Feature Extractor Plugins:

- ☒ Basic Features
- ☐ Character N-Grams
- ☐ Column Features
- ☐ Parse Features
- ☐ Regular Expressions
- ☐ Stretchy Patterns

Configure Basic Features

- ☒ Unigrams
- ☐ Bigrams
- ☐ Trigrams

Class:

Type: NOMINAL

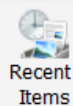
Text Fields:

☐ Differentiate

Feature Table:

Open

Look in: data



Recent Items



Desktop



My Documents



Computer



Network

- essay
- Gallup
- MovieReviews
- NewsgroupTopic
- sentiment_documents
- sentiment_sentences

File name: sentiment_sentences.csv

Files of type:

CSV

Open

Cancel

CSV Files:

sentiment_sentences

DOCUMENT_LIST

Documents: sentiment_sentences

Class: class

Type: NOMINAL

Text Fields:

☒ text

☐ Differentiate Text Fields

Feature Extractor Plugins:

- ☒ Basic Features
- ☐ Character N-Grams
- ☐ Column Features
- ☐ Parse Features
- ☐ Regular Expressions
- ☐ Stretchy Patterns

Configure Basic Features

- ☒ Unigrams
- ☐ Bigrams
- ☐ Trigrams
- ☐ POS Bigrams
- ☐ Word/POS Pairs
- ☐ Line Length
- ☐ Contains Non-Stopwords
- ☒ Binary N-grams?
- ☒ Include Punctuation?
- ☐ Remove Stopwords?
- ☐ Stem N-grams?

Extract

Name: features

Rare Threshold: 5

Feature Table:

Evaluations to Display:

Target:

Basic Table Statistics

- ☐ Recall
- ☐ Target Hits
- ☐ Precision
- ☐ Total Hits
- ☐ Correlation
- ☐ F-Score
- ☐ Kappa

Features in Table:

Search:

LightSIDE

Extract Features

Restructure Data

Build Models

Explore Results

Compare Models

Predict Labels

CSV Files:

sentiment_sentences

+

DOCUMENT_LIST

Documents: sentiment_sentences

Class: class

Type: NOMINAL

Text Fields:

text

☐ Differentiate Text Fields

Feature Extractor Plugins:

☒ Basic Features

☐ Character N-Grams

☐ Column Features

☐ Parse Features

☐ Regular Expressions

☐ Stretchy Patterns

Configure Basic Features

☒ Unigrams

☐ Bigrams

☐ Trigrams

☐ POS Bigrams

☐ Word/POS Pairs

☐ Line Length

☐ Contains Non-Stopwords

☒ Binary N-grams?

☒ Include Punctuation?

☐ Remove Stopwords?

☐ Stem N-grams?

Extract

Name: features1

Rare Threshold: 5

Feature Table:

features

+

FEATURE_TABLE

Documents: sentiment_sentences

Feature Plugins:

Feature Table: features

Evaluations to Display:

Target: neg

☐ Recall

☐ Target Hits

☐ Precision

☐ Total Hits

☐ Correlation

☐ F-Score

☐ Kappa

Features in Table:

Search:

Feature

'60s

'70s

'd

'em

'll

'm

're

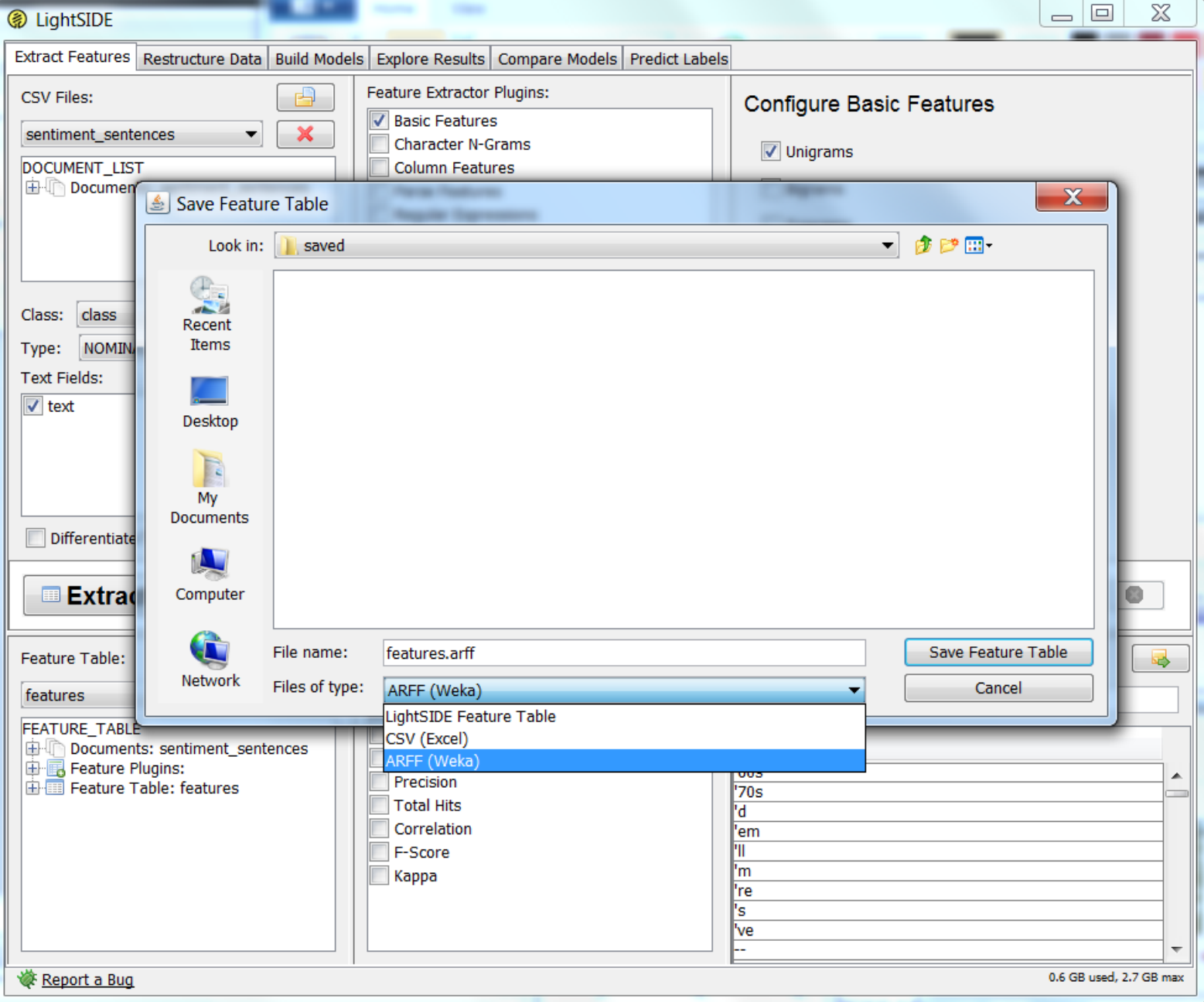
's

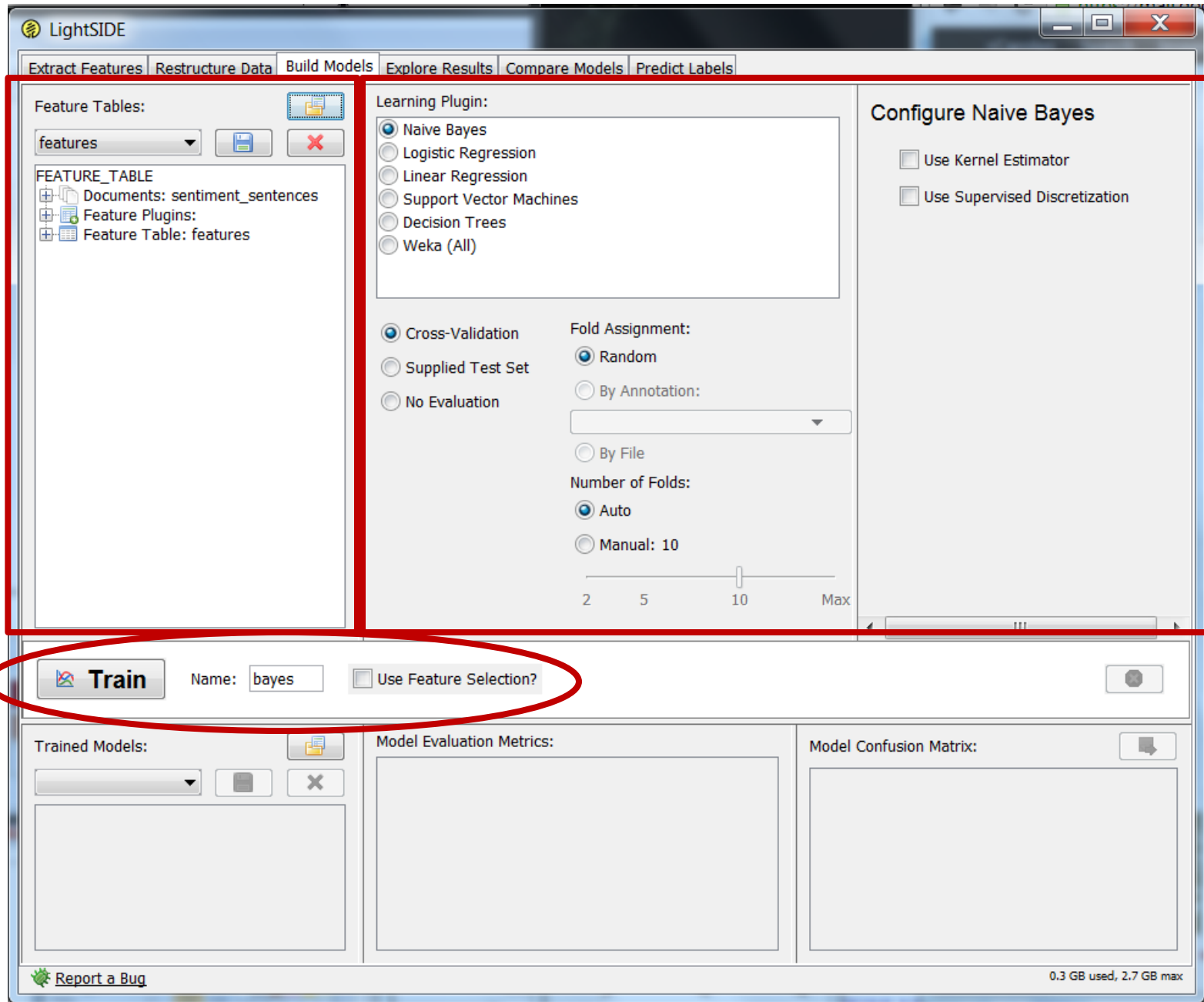
've

--

Report a Bug

0.3 GB used, 2.7 GB max





LightSIDE

Extract Features

Restructure Data

Build Models

Explore Results

Compare Models

Predict Labels

Feature Tables:

features

FEATURE_TABLE

Documents: sentiment_sentences

Feature Plugins:

Feature Table: features

Learning Plugin:

☒ Naive Bayes

☐ Logistic Regression

☐ Linear Regression

☐ Support Vector Machines

☐ Decision Trees

☐ Weka (All)

☒ Cross-Validation

☐ Supplied Test Set

☐ No Evaluation

Fold Assignment:

☒ Random

☐ By Annotation:

☐ By File

Number of Folds:

☒ Auto

☐ Manual: 10

2510Max

Configure Naive Bayes

☐ Use Kernel Estimator

☐ Use Supervised Discretization

Train

Name: bayes1

☐ Use Feature Selection?

Trained Models:

bayes

TRAINED_MODEL

Documents: sentiment_sentences

Feature Plugins:

Feature Table: features

Learning Plugin: Naive Bayes

Trained Model: bayes

Model Evaluation Metrics:

Metric	Value
Accuracy	0.774
Kappa	0.548

Model Confusion Matrix:

Act \ Pred	neg	pos
neg	4213	1118
pos	1292	4039

Report a Bug

0.6 GB used, 2.7 GB max

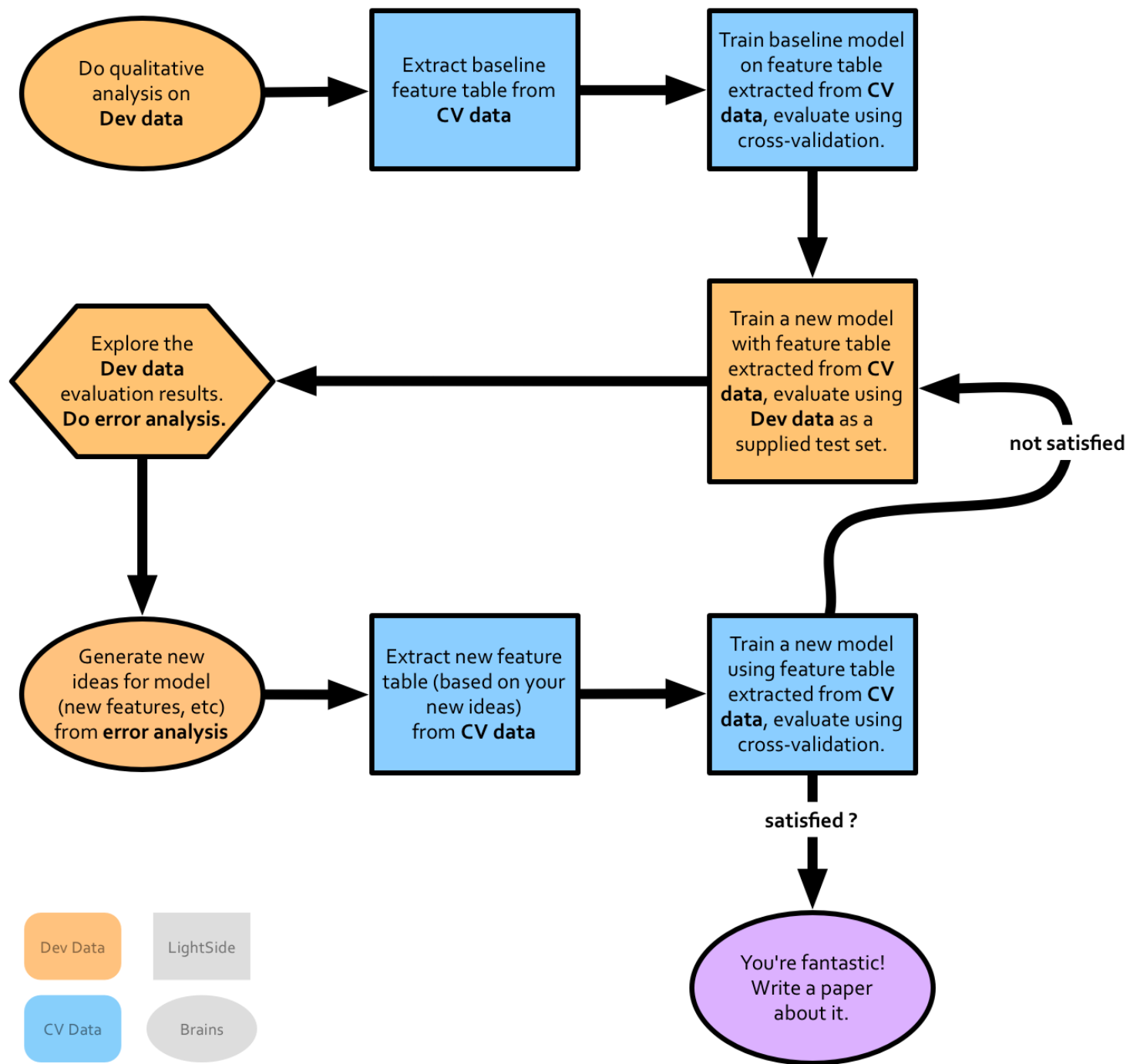


Effective Development and Evaluation Process in LightSIDE



Avoiding Overfitting!

- Separate data for evaluation from data for exploration
- We will refer to the exploration set as the Dev Set
- We will refer to the evaluation set as the cross-validation set
- You should also have a final test set you never look at until you think you are done!





Remember!!!!

- Use your development data for:
 - Qualitative analysis before ML
 - Error analysis
 - Ideas for design of new features
- Use your cross validation data for:
 - Evaluating your performance
- ***Never*** include the data you are testing on in the data you do feature selection with!!!



Basic Text Feature Extraction

Represent text as a vector where each position corresponds to a term

This is called the “bag of words” approach

Cheese
Cows
Eat
Hamsters
Make
Seeds

- Cows make cheese.
- 110010
- Hamsters eat seeds.
- 001101



Represent text as a vector where each position corresponds to a term

This is called the “bag of words” approach

*But same representation
for “Cheese makes cows.”!*

Cheese
Cows
Eat
Hamsters
Make
Seeds

■ Cows make cheese.

■ 110010

■ Hamsters eat seeds.

■ 001101



CSV Files:



Class:

Type:

Text Fields:

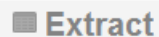
☐ Differentiate Text Fields

Feature Extractor Plugins:

- ☒ Basic Features
- ☐ Character N-Grams
- ☐ Column Features
- ☐ Parse Features
- ☐ Regular Expressions
- ☐ Stretchy Patterns

Configure Basic Features

- ☒ Unigrams
- ☐ Bigrams
- ☐ Trigrams
- ☐ POS Bigrams
- ☐ Word/POS Pairs
- ☐ Line Length
- ☐ Contains Non-Stopwords
- ☒ Binary N-grams?
- ☒ Include Punctuation?
- ☐ Remove Stopwords?
- ☐ Stem N-grams?



Extract

Name:

features

Rare Threshold:

5



Feature Table:



Evaluations to Display:

Target:

Basic Table Statistics

- ☐ Recall
- ☐ Target Hits
- ☐ Precision
- ☐ Total Hits
- ☐ Correlation
- ☐ F-Score
- ☐ Kappa

Features in Table:



Search:

U.S. government shuts down as Congress can't agree on spending bill

By Tom Cohen, CNN

updated 12:43 AM EDT, Tue October 1, 2013



STORY HIGHLIGHTS

- **NEW:** The House is expected to vote again overnight, including on appointing House negotiators
- "We will not go to conference with a gun to our head," says Sen. Harry Reid
- Obama says troops will get paid on time, but civilians may get more furloughs
- Conservatives wanted to undermine Obamacare before its private exchanges take effect Tuesday

Washington (CNN) -- The U.S. government shut down at 12:01 a.m. ET Tuesday after lawmakers in the House and the Senate could not agree on a spending bill to fund the government.

The two sides bickered and blamed each other for more than a week over Obamacare, the president's signature health care law. House Republicans insisted the spending bill include anti-Obamacare amendments. Senate Democrats were just as insistent that it didn't.

Federal employees who are considered essential will continue working. But employees deemed non-essential -- close to 800,000 -- will be furloughed.



Examples from Gallup Poll Data

- Male from Virginia, age 30, negative: “I think it’ll increase costs for everyone.”
- Female from Illinois, unknown age, positive: “Because the cost of healthcare is just outta sight crazy”
- Male from Michigan, age 70, positive: “the cost”

The Gallup Poll Dataset

The screenshot displays the LightSide software interface, which is used for extracting features from text data. The interface is divided into several sections:

- CSV Files:** A dropdown menu shows "Gallup.csv". Below it, a tree view shows "DOCUMENT_LIST" and "Documents: Gallup.csv".
- Class:** A dropdown menu shows "Vote".
- Type:** A dropdown menu shows "NOMINAL".
- Text Fields:** A list of fields with checkboxes: "Age", "Gender", "State", and "text". The "text" field is checked, and the entire section is circled in red.
- Feature Extractor Plugins:** A list of plugins with checkboxes: "Basic Features", "Character N-Grams", "Column Features", "False Features", "Regular Expressions", and "Stretchy Patterns". The "Column Features" plugin is checked and circled in red.
- Configure Basic Features:** A list of features with checkboxes: "Unigrams", "Bigrams", "Trigrams", "POS Bigrams", "Word/POS Pairs", "Line Length", "Contains Non-Stopwords", "Count Occurrences", "Include Punctuation", "Remove Stopwords", and "Stem N-Grams".
- Extract:** A button labeled "Extract".
- Name:** A text field showing "features1".
- Rare Threshold:** A text field showing "5".
- Feature Table:** A dropdown menu shows "features". Below it, a tree view shows "FEATURE_TABLE", "Documents: Gallup.csv", "Feature Plugins:", and "Feature Table: features".
- Evaluations to Display:** A dropdown menu shows "Target: Negative".
- Basic Table Statistics:** A list of statistics with checkboxes: "Correlation", "F-Score", "Kappa", "Precision", "Recall", "Target Hits", and "Total Hits".
- Features in Table:** A search bar and a list of features: "'ll", "'m", "'re", "'s", "'ve", "<COMMA>", "<ELLIPSE>", "<PERIOD>", "<QUESTIONMARK>", and "<SEMICOLON>".

At the bottom of the interface, there is a "Report a Bug" link and a status bar showing "0.2 GB used, 2.7 GB max".

Basic Types of Features

“Because the cost of healthcare is just outta sight crazy”

Configure Basic Features

☒ Unigrams

☐ Bigrams

☐ Trigrams

☐ POS Bigrams

☐ Word/POS Pairs

☐ Line Length

☐ Contains Non-Stopwords

☐ Count Occurences

☒ Include Punctuation

☐ Remove Stopwords

☐ Stem N-Grams

Basic Types of Features

“Because the cost of healthcare is just outta sight crazy”

Configure Basic Features

- ☒ Unigrams
- ☐ Bigrams
- ☐ Trigrams
- ☐ POS Bigrams
- ☐ Word/POS Pairs
- ☐ Line Length
- ☐ Contains Non-Stopwords
- ☐ Count Occurences
- ☒ Include Punctuation
- ☐ Remove Stopwords
- ☐ Stem N-Grams

Basic Types of Features

“Because the cost of healthcare is just outta sight crazy”

Configure Basic Features

- ☒ Unigrams
- ☐ Bigrams
- ☐ Trigrams
- ☐ POS Bigrams
- ☐ Word/POS Pairs
- ☐ Line Length
- ☐ Contains Non-Stopwords
- ☐ Count Occurences
- ☒ Include Punctuation
- ☐ Remove Stopwords
- ☐ Stem N-Grams

Basic Types of Features

“the cost of healthcare”

DT NN PRP NN

Configure Basic Features

- ☒ Unigrams
- ☐ Bigrams
- ☐ Trigrams
- ☐ POS Bigrams
- ☐ Word/POS Pairs
- ☐ Line Length
- ☐ Contains Non-Stopwords
- ☐ Count Occurences
- ☒ Include Punctuation
- ☐ Remove Stopwords
- ☐ Stem N-Grams



Part of Speech Tagging

<http://www.comp.leeds.ac.uk/ccalas/tagsets/upenn.html>

- | | |
|--------------------------------|------------------------------|
| 1. CC Coordinating conjunction | 12.NN Noun, singular or mass |
| 2. CD Cardinal number | 13.NNS Noun, plural |
| 3. DT Determiner | 14.NNP Proper noun, singular |
| 4. EX Existential there | 15.NNPS Proper noun, plural |
| 5. FW Foreign word | 16.PDT Predeterminer |
| 6. IN Preposition/subord | 17.POS Possessive ending |
| 7. JJ Adjective | 18.PRP Personal pronoun |
| 8. JJR Adjective, comparative | 19.PP Possessive pronoun |
| 9. JJS Adjective, superlative | 20.RB Adverb |
| 10.LS List item marker | 21.RBR Adverb, comparative |
| 11.MD Modal | 22.RBS Adverb, superlative |



Part of Speech Tagging

<http://www.comp.leeds.ac.uk/ccalas/tagsets/upenn.html>

23.RP Particle

24.SYM Symbol

25.TO to

26.UH Interjection

27.VB Verb, base form

28.VBD Verb, past tense

29.VBG Verb,
gerund/present participle

30.VBN Verb, past participle

31.VBP Verb, non-3rd ps.
sing. present

32.VBZ Verb, 3rd ps. sing.
present

33.WDT wh-determiner

34.WP wh-pronoun

35.WP Possessive wh-
pronoun

36.WRB wh-adverb

Basic Types of Features

“the cost of healthcare”

DT

NN

PRP

NN

Configure Basic Features

- ☒ Unigrams
- ☐ Bigrams
- ☐ Trigrams
- ☐ POS Bigrams
- ☐ Word/POS Pairs
- ☐ Line Length
- ☐ Contains Non-Stopwords
- ☐ Count Occurrences
- ☒ Include Punctuation
- ☐ Remove Stopwords
- ☐ Stem N-Grams

Basic Types of Features

“the cost of healthcare”

4

Configure Basic Features

- ☒ Unigrams
- ☐ Bigrams
- ☐ Trigrams
- ☐ POS Bigrams
- ☐ Word/POS Pairs
- ☐ Line Length
- ☐ Contains Non-Stopwords
- ☐ Count Occurences
- ☒ Include Punctuation
- ☐ Remove Stopwords
- ☐ Stem N-Grams

Basic Types of Features

“the cost of healthcare”

YES

Configure Basic Features

- ☒ Unigrams
- ☐ Bigrams
- ☐ Trigrams
- ☐ POS Bigrams
- ☐ Word/POS Pairs
- ☐ Line Length
- ☐ Contains Non-Stopwords
- ☐ Count Occurrences
- ☒ Include Punctuation
- ☐ Remove Stopwords
- ☐ Stem N-Grams

Basic Types of Features

“the cost is too great. The cost is immense!”

Configure Basic Features

- ☒ Unigrams
- ☐ Bigrams
- ☐ Trigrams
- ☐ POS Bigrams
- ☐ Word/POS Pairs
- ☐ Line Length
- ☐ Contains Non-Stopwords
- ☐ Count Occurrences
- ☒ Include Punctuation
- ☐ Remove Stopwords
- ☐ Stem N-Grams

The value of the feature is the number of times it occurs, rather than 1 if it occurs or 0 otherwise, which is the default.

Basic Types of Features

“the cost is too great. The cost is immense!”

Configure Basic Features

- ☒ Unigrams
- ☐ Bigrams
- ☐ Trigrams
- ☐ POS Bigrams
- ☐ Word/POS Pairs
- ☐ Line Length
- ☐ Contains Non-Stopwords
- ☐ Count Occurences
- ☒ Include Punctuation
- ☐ Remove Stopwords
- ☐ Stem N-Grams

If you uncheck this, punctuation will be ignored and stripped out of the representation.

Basic Types of Features

~~“the cost of healthcare”~~

Configure Basic Features

- ☒ Unigrams
- ☐ Bigrams
- ☐ Trigrams
- ☐ POS Bigrams
- ☐ Word/POS Pairs
- ☐ Line Length
- ☐ Contains Non-Stopwords
- ☐ Count Occurrences
- ☒ Include Punctuation
- ☐ Remove Stopwords
- ☐ Stem N-Grams

Basic Types of Features

“healthcare costss” → “healthcare cost”

Configure Basic Features

- ☒ Unigrams
- ☐ Bigrams
- ☐ Trigrams
- ☐ POS Bigrams
- ☐ Word/POS Pairs
- ☐ Line Length
- ☐ Contains Non-Stopwords
- ☐ Count Occurences
- ☒ Include Punctuation
- ☐ Remove Stopwords
- ☐ Stem N-Grams



Clarification on Basic text feature extractor

- POS tagging happens before stemming or stopword removal
- POS bigrams are not affected by stopword removal – POS tags for stopwords will still be included
- On word n-grams, the only n-grams that will be dropped in the case of stopword removal are ones that consist only of stopwords

Feature Space Customizations

■ Feature Space Design

- Think like a computer!
- Machine learning algorithms look for features that are good predictors, not features that are necessarily meaningful
- Look for approximations
 - If you want to find questions, you don't need to do a complete syntactic analysis
 - Look for question marks
 - Look for wh-terms that occur immediately before an auxiliary verb



Error Analysis

Error Analysis Process

High Level Overview

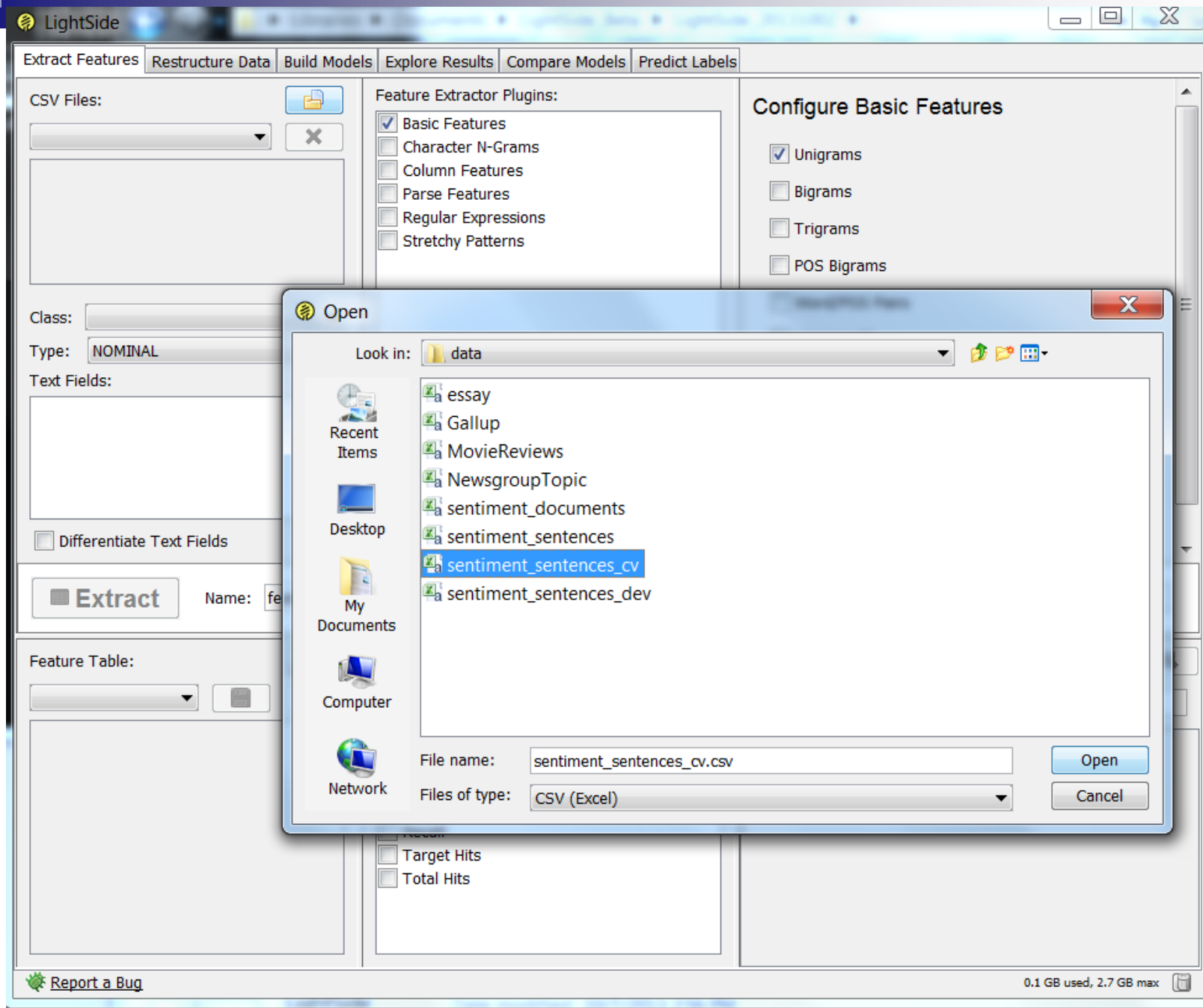
=== Confusion Matrix ===

a	b	
44	25	
10	149	

<-- classified as

Goal: We want to discover how to re-represent the data so that instances with the same class value look more similar to one another and instances with different class values look more different

- Identify large error cells
- Make comparisons
 - Ask yourself how it is similar to the instances that were correctly classified with the same class (**vertical comparison**)
 - How it is different from those it was incorrectly not classified as (**horizontal comparison**)



CSV Files:

sentiment_sentences_cv.csv

DOCUMENT_LIST

Documents: sentiment_sentences_cv.

Class: class

Type: NOMINAL

Text Fields:

☒ text

☐ Differentiate Text Fields

Feature Extractor Plugins:

- ☒ Basic Features
- ☐ Character N-Grams
- ☐ Column Features
- ☐ Parse Features
- ☐ Regular Expressions
- ☐ Stretchy Patterns

Configure Basic Features

- ☒ Unigrams
- ☐ Bigrams
- ☐ Trigrams
- ☐ POS Bigrams
- ☐ Word/POS Pairs
- ☐ Line Length
- ☐ Contains Non-Stopwords
- ☐ Count Occurrences
- ☒ Include Punctuation
- ☐ Remove Stopwords
- ☐ Stem N-Grams

Extract

Name: features1

Rare Threshold: 5

Feature Table:

features

FEATURE_TABLE

Documents: sentiment_sentences_cv.

Feature Plugins:

Feature Table: features

Evaluations to Display:

Target: neg

Basic Table Statistics

- ☐ Correlation
- ☐ F-Score
- ☐ Kappa
- ☐ Precision
- ☐ Recall
- ☐ Target Hits
- ☐ Total Hits

Features in Table:

Search:

Feature

'60s
'70s
'd
'em
'll
'm
're
's
've
--

Feature Tables:

features

FEATURE_TABLE

- Documents: sentiment_sentences_cv.
- Feature Plugins:
- Feature Table: features

Learning Plugin:

- ☐ Naive Bayes
- ☐ Logistic Regression
- ☐ Linear Regression
- ☐ Support Vector Machines
- ☐ Decision Trees
- ☒ Weka (All)

- ☒ Cross-Validation
- ☐ Supplied Test Set
- ☐ No Evaluation

Fold Assignment:

- ☒ Random
- ☐ By Annotation:
- ☐ By File

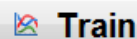
Number of Folds:

- ☒ Auto
- ☐ Manual: 10

2 5 10 Max

Configure Weka (All)

- weka
 - classifiers
 - bayes
 - functions
 - GaussianProcesses
 - LibLINEAR
 - LibSVM
 - LinearRegression
 - Logistic
 - MultilayerPerceptron
 - SGD
 - SGDText
 - SimpleLinearRegression
 - SimpleLogistic
 - SMO
 - SMOreg
 - VotedPerceptron
 - lazy
 - meta
 - misc
 - rules
 - trees



Train

Name: weka

☐ Feature Selection

Trained Models:

Model Evaluation Metrics:

Model Confusion Matrix:

Feature Tables:

features

FEATURE_TABLE

- Documents: sentiment_sentences_cv.
- Feature Plugins:
- Feature Table: features

Learning Plugin:

- ☐ Naive Bayes
- ☐ Logistic Regression
- ☐ Linear Regression
- ☐ Support Vector Machines
- ☐ Decision Trees
- ☒ Weka (All)

- ☒ Cross-Validation
- ☐ Supplied Test Set
- ☐ No Evaluation

Fold Assignment:

- ☒ Random
- ☐ By Annotation:

- ☐ By File

Number of Folds:

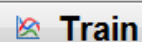
- ☒ Auto
- ☐ Manual: 10

2 5 10 Max

Configure Weka (All)

Choose

SMD -C 1.0 -L 0.001 -P 1.0



Train

Name: weka1

☒ Feature Selection #: 1000

Trained Models:

weka

TRAINED_MODEL

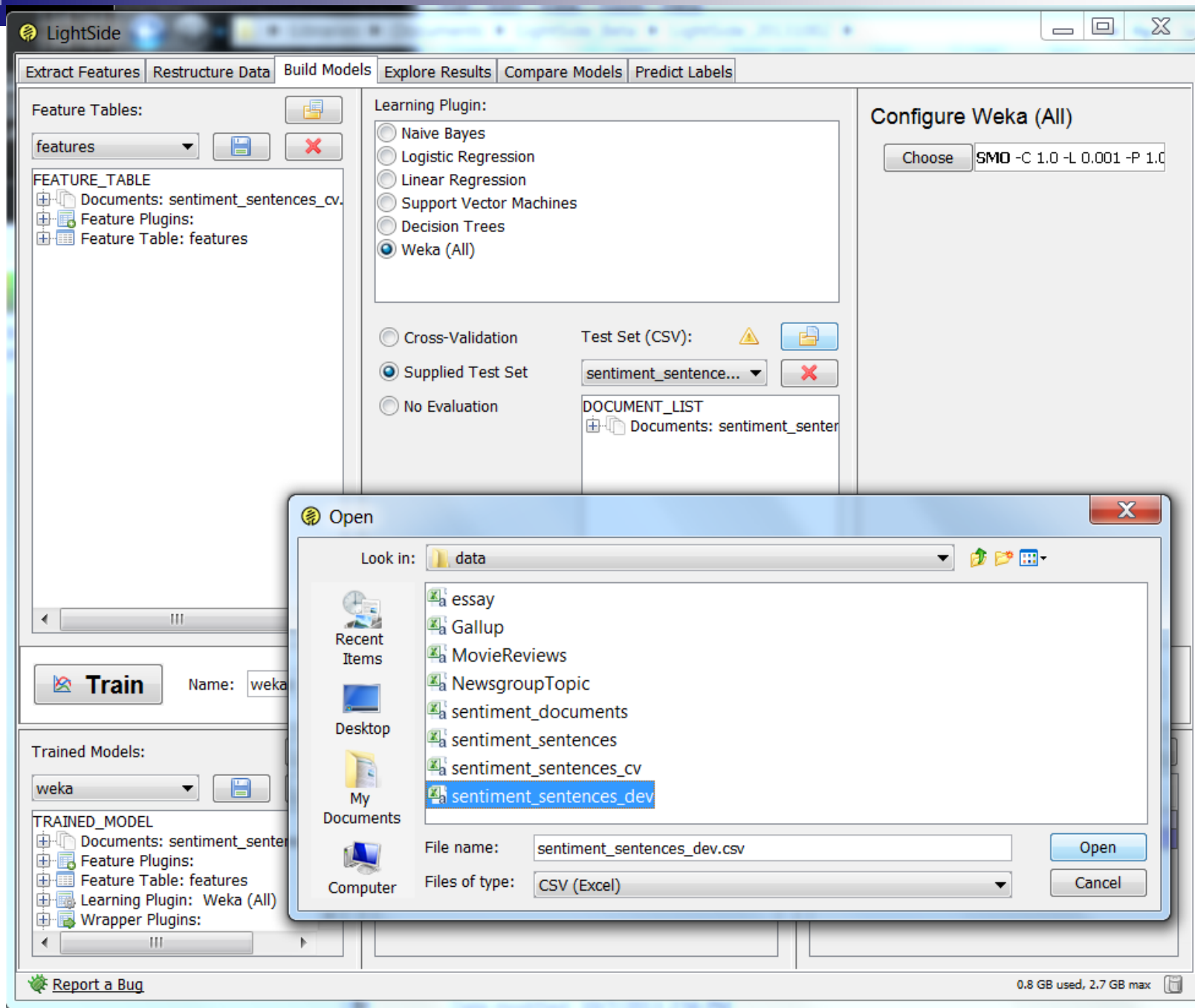
- Documents: sentiment_sentences
- Feature Plugins:
- Feature Table: features
- Learning Plugin: Weka (All)
- Wrapper Plugins:

Model Evaluation Metrics:

Metric	Value
Accuracy	0.7343
Kappa	0.4685

Model Confusion Matrix:

Act \ Pred	neg	pos
neg	2905	941
pos	1095	2721



Extract Features | Restructure Data | Build Models | Explore Results | Compare Models | Predict Labels

Feature Tables:

features

FEATURE_TABLE

- Documents: sentiment_sentences_cv
- Feature Plugins:
- Feature Table: features

Learning Plugin:

- ☐ Naive Bayes
- ☐ Logistic Regression
- ☐ Linear Regression
- ☐ Support Vector Machines
- ☐ Decision Trees
- ☒ Weka (All)

☐ Cross-Validation☒ Supplied Test Set☐ No Evaluation

Test Set (CSV):

sentiment_sentence...

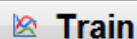
DOCUMENT_LIST

Documents: sentiment_sen...

Configure Weka (All)

Choose

SMD -C 1.0 -L 0.001 -P 1.0



Train

Name: weka2

☒ Feature Selection #: 1000

Trained Models:

weka1

TRAINED_MODEL

- Documents: sentiment_sentences
- Feature Plugins:
- Feature Table: features
- Learning Plugin: Weka (All)
- Wrapper Plugins:

Model Evaluation Metrics:

Metric	Value
Accuracy	0.7507
Kappa	0.5015

Model Confusion Matrix:

Act \ Pred	neg	pos
neg	1149	336
pos	412	1103

LightSide

Extract Features | Restructure Data | Build Models | Explore Results | Compare Models | Predict Labels

Highlight: weka1

Cell Highlight:

Act \ Pred	neg	pos
neg	1149	336
pos	412	1103

Evaluations to Display:

- ☐ Feature Confusion Ranking
- ☐ Average Value
- ☐ Frequency
- ☐ Horizontal Absolute Difference
- ☐ Horizontal Difference
- ☐ Vertical Absolute Difference

Features in Table:

Search:

Feature

- ☐ '60s
- ☐ '70s
- ☐ 'd
- ☐ 'em
- ☐ 'il
- ☐ 'm
- ☐ 're
- ☐ 's
- ☐ 've
- ☐ --
- ☐ -lsb-
- ☐ -rsb-
- ☐ 1
- ☐ 10
- ☐ 100
- ☐ 101

Exploration Plugin: Model Output (Text)

+ -0.5744 * (normalized) witless
+ 1 * (normalized) witty
+ 1.1405 * (normalized) wonderful
+ 1.0185 * (normalized) wonderfully
+ -0.516 * (normalized) wooden
+ -1 * (normalized) woody
+ 0.1701 * (normalized) work
+ 1.1733 * (normalized) works
+ 0.6076 * (normalized) world
+ -0.3756 * (normalized) worse
+ -1.3245 * (normalized) worst
+ 1.06 * (normalized) worth
+ -0.2368 * (normalized) would
+ -1.0217 * (normalized) writers
+ -0.9047 * (normalized) wrong
+ 1 * (normalized) wry
+ 1.1511 * (normalized) y
+ 0.6517 * (normalized) year
+ 0.4901 * (normalized) years
+ 0.4529 * (normalized) yet
- 0.2855

Number of kernel evaluations: 144356432 (53.759% cached)

Report a Bug

0.8 GB used, 2.7 GB max

LightSide

Extract Features | Restructure Data | Build Models | Explore Results | Compare Models | Predict Labels

Highlight: **weka1**

Cell Highlight:

Act \ Pred	neg	pos
neg	1149	336
pos	412	1103

Evaluations to Display:

- ☒ Frequency
- ☐ Horizontal Absolute Difference
- ☐ Horizontal Difference
- ☒ Vertical Absolute Difference
- ☐ Vertical Difference

Model Analysis

Features in Table:

Search:

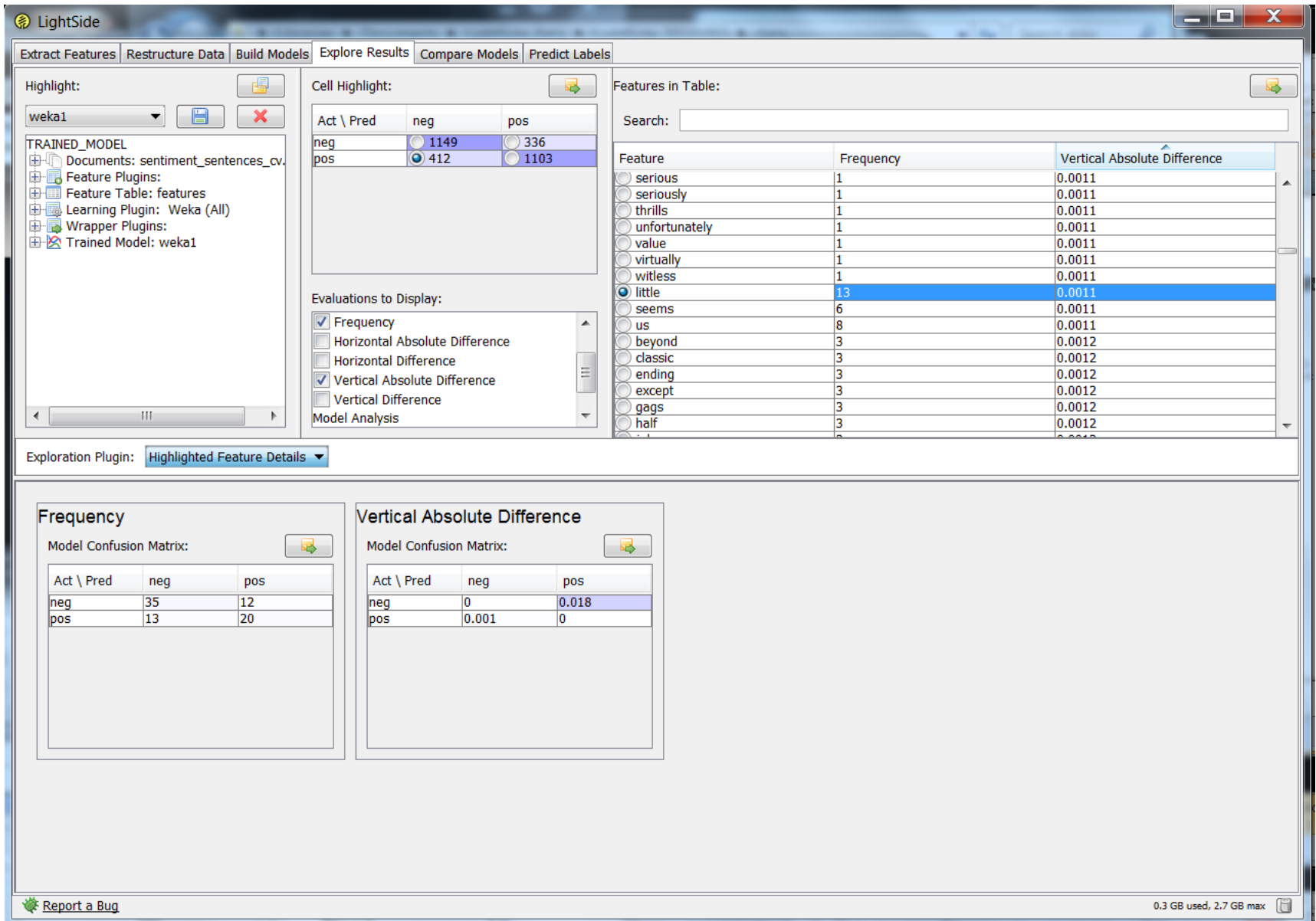
Feature	Frequency	Vertical Absolute Difference
serious	1	0.0011
seriously	1	0.0011
thrills	1	0.0011
unfortunately	1	0.0011
value	1	0.0011
virtually	1	0.0011
witless	1	0.0011
little	13	0.0011
seems	6	0.0011
us	8	0.0011
beyond	3	0.0012
classic	3	0.0012
ending	3	0.0012
except	3	0.0012
gags	3	0.0012
half	3	0.0012

Exploration Plugin: **Model Output (Text)**

- + 0.3548 * (normalized) lee
- + -0.2567 * (normalized) left
- + 0.4641 * (normalized) legacy
- + 0.6121 * (normalized) leigh
- + -0.4566 * (normalized) less
- + -0.5215 * (normalized) let
- + 0.4781 * (normalized) life
- + -1.1239 * (normalized) lifeless
- + -0.2275 * (normalized) like
- + 0.7136 * (normalized) lips
- + -0.5328 * (normalized) list
- + -1 * (normalized) listless
- + 0.5816 * (normalized) literary
- + 0.939 * (normalized) **iterate**
- + -0.4418 * (normalized) little
- + 1.1512 * (normalized) lively
- + 0.0912 * (normalized) lives
- + -0.0151 * (normalized) long
- + 0.4274 * (normalized) longing
- + 0.2859 * (normalized) look
- + -0.5328 * (normalized) looking
- + -0.9967 * (normalized) loses
- + -0.9503 * (normalized) loud
- + -0.6698 * (normalized) lousy
- + 0.3711 * (normalized) love

Report a Bug

0.3 GB used, 2.7 GB max



LightSide

Extract Features | Restructure Data | Build Models | Explore Results | Compare Models | Predict Labels

Highlight:
weka1

TRAINED_MODEL
Documents: sentiment_sentences_cv.
Feature Plugins:
Feature Table: features
Learning Plugin: Weka (All)
Wrapper Plugins:
Trained Model: weka1

Cell Highlight:

Act \ Pred	neg	pos
neg	1149	336
pos	412	1103

Evaluations to Display:
☒ Frequency
☐ Horizontal Absolute Difference
☐ Horizontal Difference
☒ Vertical Absolute Difference
☐ Vertical Difference
Model Analysis

Features in Table:

Search:

Feature	Frequency	Vertical Absolute Difference
<input type="radio"/> serious	1	0.0011
<input type="radio"/> seriously	1	0.0011
<input type="radio"/> thrills	1	0.0011
<input type="radio"/> unfortunately	1	0.0011
<input type="radio"/> value	1	0.0011
<input type="radio"/> virtually	1	0.0011
<input type="radio"/> witless	1	0.0011
<input checked="" type="radio"/> little	13	0.0011
<input type="radio"/> seems	6	0.0011
<input type="radio"/> us	8	0.0011
<input type="radio"/> beyond	3	0.0012
<input type="radio"/> classic	3	0.0012
<input type="radio"/> ending	3	0.0012
<input type="radio"/> except	3	0.0012
<input type="radio"/> gags	3	0.0012
<input type="radio"/> half	3	0.0012

Exploration Plugin: Documents Display

☒ Filter documents by selected feature
☐ Reverse document filter
☒ Documents from selected cell only

Instance	Predicted	Actual	Text
<input checked="" type="checkbox"/> 124	neg	pos	we know the plot...
<input checked="" type="checkbox"/> 189	neg	pos	pretty good little ...
<input checked="" type="checkbox"/> 229	neg	pos	shyamalan offers...
<input checked="" type="checkbox"/> 563	neg	pos	daughter from da...
<input type="checkbox"/> 697	neg	pos	there are times ...
<input type="checkbox"/> 1414	neg	pos	reign of fire may ...
<input type="checkbox"/> 1690	neg	pos	it's excessively q...
<input type="checkbox"/> 1700	neg	pos	an original little fi...
<input type="checkbox"/> 2629	neg	pos	there's very little ...
<input type="checkbox"/> 2733	neg	pos	grown-up quibble...
<input type="checkbox"/> 2809	neg	pos	a sensitive and e...
<input type="checkbox"/> 2815	neg	pos	affectionately re...
<input type="checkbox"/> 2979	neg	pos	all the pieces fall ...

Instance 189 (Predicted neg, Actual pos)

Highlighting little feature hits

pretty good little movie .

Instance 229 (Predicted neg, Actual pos)

Highlighting little feature hits

shyamalan offers copious hints along the way -- myriad signs , if you will -- that beneath the familiar , funny surfa
ce is a far bigger , far more meaningful story than one in which little green men come to earth for harvesting pur
poses .

Report a Bug

0.4 GB used, 2.7 GB max

LightSide

FileHomeInsertDesignTransitionsAnimationsSlide Show

Extract FeaturesRestructure DataBuild ModelsExplore ResultsCompare ModelsPredict Labels

Highlight:
weka1

TRAINED_MODEL

- Documents: sentiment_sentences_cv.
- Feature Plugins:
- Feature Table: features
- Learning Plugin: Weka (All)
- Wrapper Plugins:
- Trained Model: weka1

Cell Highlight:
Act \ Prednegpos
neg1149336
pos4121103
Evaluations to Display:

- ☒ Frequency
- ☐ Horizontal Absolute Difference
- ☐ Horizontal Difference
- ☒ Vertical Absolute Difference
- ☐ Vertical Difference
- Model Analysis

Features in Table:
Search:

Feature	Frequency	Vertical Absolute Difference
sandler	0	0
sandra	0	0
santa	0	0
sara	0	0
satisfy	0	0
saturday	0	0
save	0	0
saves	0	0
saving	0	0
saying	0	0
scares	0	0
scattered	0	0
scenery	0	0
schaeffer	0	0
scherfig	0	0
schmidt	0	0

Exploration Plugin: Documents Display

☒ Filter documents by selected feature
☐ Reverse document filter
☒ Documents from selected cell only

Instance	Predicted	Actual	Text
<input checked="" type="checkbox"/> 20	pos	pos	who is the audie...
<input checked="" type="checkbox"/> 28	pos	pos	first-time writer...
<input checked="" type="checkbox"/> 60	pos	pos	a rather brilliant...
<input checked="" type="checkbox"/> 209	pos	pos	'possession , ' b...
<input checked="" type="checkbox"/> 343	pos	pos	a weird , arresti...
<input checked="" type="checkbox"/> 453	pos	pos	amari has dress...
<input checked="" type="checkbox"/> 521	pos	pos	it's a great deal ...
<input type="checkbox"/> 536	pos	pos	witty and often s...
<input type="checkbox"/> 712	pos	pos	disney has alwa...
<input type="checkbox"/> 733	pos	pos	director juan jos...
<input type="checkbox"/> 1003	pos	pos	the tasteful little...
<input type="checkbox"/> 1307	pos	pos	there is a refres...
<input type="checkbox"/> 1556	pos	pos	the stunning , dr...
<input type="checkbox"/> 1674	pos	pos	a solidly enterta...
<input type="checkbox"/> 1692	pos	pos	a smart little ind...

Instance 453 (Predicted pos, Actual pos)
Highlighting little feature hits
amari has dressed up this little parable in a fairly irresistible package full of privileged moments and memorable p
erformances .

Instance 521 (Predicted pos, Actual pos)
Highlighting little feature hits
it's a great deal of sizzle and very little steak . but what spectacular sizzle it is ! . . . in this incarnation its fizz is inf
ectious .

Report a Bug

0.4 GB used, 2.7 GB max

* Testing bigrams as an alternative....

LightSide

Extract Features | Restructure Data | Build Models | Explore Results | Compare Models | Predict Labels

CSV Files:

sentiment_sentences_cv.csv

DOCUMENT_LIST

Documents: sentiment_sentences_cv.

Class: class

Type: NOMINAL

Text Fields:

☒ text

☐ Differentiate Text Fields

Feature Extractor Plugins:

☒ Basic Features
☐ Character N-Grams
☐ Column Features
☐ Parse Features
☐ Regular Expressions
☐ Stretchy Patterns

☒ Unigrams
☒ Bigrams
☐ Trigrams
☐ POS Bigrams
☐ Word/POS Pairs
☐ Line Length
☐ Contains Non-Stopwords
☐ Count Occurrences
☒ Include Punctuation
☐ Remove Stopwords
☐ Stem N-Grams
☒ Track Feature Hit Location

Extract Name: features3 Rare Threshold: 5

Feature Table:

features2

FEATURE_TABLE

Documents: sentiment_sentences_cv.
Feature Plugins:
Feature Table: features2

Evaluations to Display:

Target: neg

Basic Table Statistics

☐ Correlation
☐ F-Score
☐ Kappa
☐ Precision
☐ Recall
☐ Target Hits
☐ Total Hits

Features in Table:

Search:

Feature

indeed
indeed_COMMA
independent
indian
indie
individual
industry
inept
inevitable
infectious
infomercial
infuses
ingenious
ingredients
inhabit
inner

[Report a Bug](#)

0.4 GB used, 2.7 GB max

LightSide

Extract Features | Restructure Data | Build Models | Explore Results | Compare Models | Predict Labels

Feature Tables:

features2

FEATURE_TABLE

Documents: sentiment_sentences_cv.

Feature Plugins:

Feature Table: features2

Learning Plugin:

☐ Naive Bayes

☐ Logistic Regression

☐ Linear Regression

☐ Support Vector Machines

☐ Decision Trees

☒ Weka (All)

☒ Cross-Validation

☐ Supplied Test Set

☐ No Evaluation

Fold Assignment:

☒ Random

☐ By Annotation:

☐ By File

Number of Folds:

☒ Auto

☐ Manual: 10

2

5

10

Max

Configure Weka (All)

Choose SMO -C 1.0 -L 0.001 -P 1.0E-12 -N 0 -V -1 -W 1 -K "weka.classifiers.functions.s

Train

Name: weka4

☒ Feature Selection # : 1000

Trained Models:

weka3

TRAINED_MODEL

Documents: sentiment_sentences_cv.

Feature Plugins:

Feature Table: features2

Learning Plugin: Weka (All)

Wrapper Plugins:

Trained Model: weka3

Model Evaluation Metrics:

Metric	Value
Accuracy	0.7219
Kappa	0.4436

Model Confusion Matrix:

Act \ Pred	neg	pos
neg	2898	948
pos	1183	2633

Report a Bug

0.4 GB used, 2.7 GB max

LightSide

Extract Features | Restructure Data | Build Models | Explore Results | Compare Models | Predict Labels

Baseline Model:
weka

TRAINING_MODEL
Documents: sentiment_sentences_cv.csv
Feature Plugins:
Feature Table: features
Learning Plugin: Weka (All)
Wrapper Plugins:
Trained Model: weka

Competing Model:
weka3

TRAINING_MODEL
Documents: sentiment_sentences_cv.csv
Feature Plugins:
Feature Table: features2
Learning Plugin: Weka (All)
Wrapper Plugins:
Trained Model: weka3

Comparison Plugin: Basic Model Comparison

Baseline Model Metrics:

Metric	Value
Accuracy	0.7343
Kappa	0.4685

Competing Model Metrics:

Metric	Value
Accuracy	0.7219
Kappa	0.4436

Baseline Confusion Matrix:

Act \ Pred	neg	pos
neg	2905	941
pos	1095	2721

Competing Confusion Matrix:

Act \ Pred	neg	pos
neg	2898	948
pos	1183	2633

Highly significant improvement (p=0.008**, t=2.635)

Report a Bug

0.4 GB used, 2.7 GB max

LightSide

Extract Features | Restructure Data | Build Models | Explore Results | Compare Models | Predict Labels

Highlight:
weka1

TRAINED_MODEL

- Documents: sentiment_sentences_cv.
- Feature Plugins:
- Feature Table: features
- Learning Plugin: Weka (All)
- Wrapper Plugins:
- Trained Model: weka1

Cell Highlight:

Act \ Pred	neg	pos
neg	1149	336
pos	412	1103

Evaluations to Display:

- ☒ Frequency
- ☒ Horizontal Absolute Difference
- ☐ Horizontal Difference
- ☐ Vertical Absolute Difference
- ☐ Vertical Difference

Model Analysis

Features in Table:

Search:

Feature	Frequency	Horizontal Absolute Difference
's	101	0.0568
with	39	0.0577
<COMMA>	233	0.0591
like	40	0.0617
does	35	0.0632
movie	59	0.0734
n't	51	0.0839
film	38	0.0855
but	88	0.0885
a	188	0.1076
of	155	0.1088
and	148	0.2038
an	50	0.0001
still	10	0.0002
cast	7	0.0002
ever	7	0.0002

Exploration Plugin: Model Output (Text)

+ -0.8604 * (normalized) lame

+ -1 * (normalized) laughable

+ -0.8779 * (normalized) lawrence

+ -0.5057 * (normalized) lazy

+ -1 * (normalized) leaden

+ -1.0732 * (normalized) leaves

+ 0.3548 * (normalized) lee

+ -0.2567 * (normalized) left

+ 0.4641 * (normalized) legacy

+ 0.6121 * (normalized) leigh

+ -0.4566 * (normalized) less

+ -0.5215 * (normalized) let

+ 0.4781 * (normalized) life

+ -1.1239 * (normalized) lifeless

+ -0.2275 * (normalized) like

+ 0.7136 * (normalized) lips

+ -0.5328 * (normalized) list

+ -1 * (normalized) listless

+ 0.5816 * (normalized) literary

+ 0.939 * (normalized) literate

+ -0.4418 * (normalized) little

+ 1.1512 * (normalized) lively

+ 0.0912 * (normalized) lives

+ -0.0151 * (normalized) long

+ 0.4274 * (normalized) longing

Report a Bug

0.3 GB used, 2.7 GB max

Highlight:

logit

TRAINED_MODEL

- Documents: sentiment_sentences_cv.
- Feature Plugins:
- Feature Table: features
- Learning Plugin: Logistic Regression
- Trained Model: logit

Cell Highlight:

Act \ Pred	neg	pos
neg	1137	348
pos	379	1136

Evaluations to Display:

- ☐ Horizontal Difference
- ☐ Vertical Absolute Difference
- ☐ Vertical Difference

Model Analysis

- ☐ Feature Influence
- ☒ Feature Weight

Features in Table:

Search:

Feature	Frequency	Horizontal Absolute Differ...	Feature Weight
fun	2	0.022	-1.0294
funeral	0	0	-0.1445
funnier	2	0.0053	0.5206
funniest	0	0	-0.4535
funny	12	0.0035	-0.8932
further	1	0.0018	0.0474
future	2	0.0018	0.8045
fuzzy	0	0	0.139
g	0	0	0.7943
gag	0	0.0009	0.2412
gaghan	0	0	0.3002
gags	3	0.0062	0.7382
game	2	0.0026	0.3971
games	0	0.0009	-0.1543
gang	1	0.0009	0.111
gangs	0	0.0009	-0.6682

Exploration Plugin: Highlighted Feature Details

Frequency

Model Confusion Matrix:

Act \ Pred	neg	pos
neg	12	12
pos	12	40

Horizontal Absolute Difference

Model Confusion Matrix:

Act \ Pred	neg	pos
neg	0	0.024
pos	0.004	0

Feature Weight

Model Confusion Matrix:

Act \ Pred	neg	pos
neg	-0.893	0.893
pos	-0.893	0.893

Highlight:

logit

TRAINED_MODEL
 Documents: sentiment_sentences_cv.
 Feature Plugins:
 Feature Table: features
 Learning Plugin: Logistic Regression
 Trained Model: logit

Cell Highlight:

Act \ Pred	neg	pos
neg	1137	348
pos	379	1136

Evaluations to Display:

- ☐ Horizontal Difference
☐ Vertical Absolute Difference
☐ Vertical Difference
 Model Analysis
☐ Feature Influence
☒ Feature Weight

Features in Table:

Search:

Feature	Frequency	Horizontal Absolute Differ...	Feature Weight
the	229	0.0532	0.0508
a	187	0.055	-0.1195
do	26	0.0554	0.3847
movie	51	0.0597	0.3268
does	32	0.0607	0.3848
of	154	0.0655	-0.0984
with	31	0.0731	-0.3772
n't	47	0.0818	0.5158
film	33	0.0899	-0.3272
but	89	0.1142	0.2248
and	135	0.2019	-0.3715
'll	5	0	-0.3673
two	5	0	-0.3729
screen	6	0	0.0045
sense	6	0	0.0931

Exploration Plugin: Highlighted Feature Details

Frequency

Model Confusion Matrix:

Act \ Pred	neg	pos
neg	112	35
pos	47	48

Horizontal Absolute Difference

Model Confusion Matrix:

Act \ Pred	neg	pos
neg	0	0.002
pos	0.082	0

Feature Weight

Model Confusion Matrix:

Act \ Pred	neg	pos
neg	0.516	-0.516
pos	0.516	-0.516

Extract Features

Restructure Data

Build Models

Explore Results

Compare Models

Predict Labels

CSV Files:

sentiment_sentences_cv.csv

DOCUMENT_LIST

Documents: sentiment_sentences_cv.

Class: class

Type: NOMINAL

Text Fields:

☒ text

☐ Differentiate Text Fields

Feature Extractor Plugins:

☒ Basic Features

☐ Character N-Grams

☐ Column Features

☐ Parse Features

☐ Regular Expressions

☒ Stretchy Patterns

0 1 2 3 4 5 6 7 8 Gap Length

About Stretchy Patterns

☒ Include surface words in patterns

☒ Include POS tags in patterns

Categories: Add... Clear

nt: [n't]

☒ Require at least one category per pattern

☐ Don't include surface/POS form when a category matches

☒ Categories match against surface words

☐ Categories match against POS tags

☐ Count pattern hits

Prune Rare Features after N documents:

0 100 200 500 1000

Extract

Name: features2

Rare Threshold: 5

Feature Table:

features1

FEATURE_TABLE

Documents: sentiment_sentences_cv.

Feature Plugins:

Feature Table: features1

Evaluations to Display:

Target: neg

☐ Correlation

☐ F-Score

☐ Kappa

☐ Precision

☐ Recall

☐ Target Hits

☐ Total Hits

Features in Table:

Search:

Feature

does nt [GAP] JJ

does nt [GAP] JJR

does nt [GAP] NN

does nt [GAP] PRP

does nt [GAP] RB

does nt [GAP] VB

does nt [GAP] a

does nt [GAP] it

does nt [GAP] much

does nt [GAP] the

film VBZ nt

film [GAP] nt

film [GAP] nt VB

has nt

have nt

have nt VBN

have nt [GAP] IN

d, 2.7 GB max

0.8 GB used, 2.7 GB max

Feature Tables:

features

FEATURE_TABLE

- Documents: sentiment_sentences_cv.
- Feature Plugins:
- Feature Table: features

Learning Plugin:

- ☐ Naive Bayes
- ☒ Logistic Regression
- ☐ Linear Regression
- ☐ Support Vector Machines
- ☐ Decision Trees
- ☐ Weka (All)

- ☒ Cross-Validation
- ☐ Supplied Test Set
- ☐ No Evaluation

Fold Assignment:

☒ Random☐ By Annotation:☐ By File

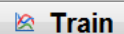
Number of Folds:

☒ Auto☐ Manual: 10

2 5 10 Max

Configure Logistic Regression

- ☐ L2 Regularization
- ☒ L1 Regularization
- ☐ L2 Regularization (Dual)



Train

Name: logit6

☐ Feature Selection

Trained Models:

logit5

TRAINED_MODEL

- Documents: sentiment_sentences_cv.
- Feature Plugins:
- Feature Table: features
- Learning Plugin: Logistic Regression
- Trained Model: logit5

Model Evaluation Metrics:

Metric	Value
Accuracy	0.7254
Kappa	0.4508

Model Confusion Matrix:

Act \ Pred	neg	pos
neg	2820	1026
pos	1078	2738

Baseline Model:

logit4

TRAINED_MODEL

- Documents: sentiment_sentences_cv.csv
- Feature Plugins:
- Feature Table: features1
- Learning Plugin: Logistic Regression
- Trained Model: logit4

Competing Model:

logit5

TRAINED_MODEL

- Documents: sentiment_sentences_cv.csv
- Feature Plugins:
- Feature Table: features
- Learning Plugin: Logistic Regression
- Trained Model: logit5

Comparison Plugin: Basic Model Comparison

Baseline Model Metrics:

Metric	Value
Accuracy	0.7317
Kappa	0.4633

Competing Model Metrics:

Metric	Value
Accuracy	0.7254
Kappa	0.4508

Baseline Confusion Matrix:

Act \ Pred	neg	pos
neg	2845	1001
pos	1055	2761

Competing Confusion Matrix:

Act \ Pred	neg	pos
neg	2820	1026
pos	1078	2738

Highly significant improvement ($p=0.006^{**}$, $t=2.745$)



Special Text Features

Stretchy Patterns in LightSIDE

Looking at sentiment_sentences.csv

The screenshot displays the LightSIDE application window with the 'Extract Features' tab selected. The interface is organized into several panels:

- CSV Files:** Shows 'sentiment_sentences.csv' loaded. Below it, a 'DOCUMENT_LIST' section lists 'Documents: sentiment_sentences.csv'. The 'Class' is set to 'class' and the 'Type' is 'NOMINAL'. Under 'Text Fields', 'text' is checked. There is a 'Differentiate Text Fields' checkbox at the bottom.
- Feature Extractor Plugins:** A list of plugins with checkboxes. 'Basic Features' and 'Stretchy Patterns' are checked. Other options include 'Character N-Grams', 'Column Features', 'Parse Features', and 'Regular Expressions'.
- Configure Stretchy Patterns:** Contains two sliders for 'Pattern Length' (set to 4) and 'Gap Length' (set to 2). There is an 'About Stretchy Patterns' button. Checkboxes for 'Include surface words in patterns' (checked), 'Include POS tags in patterns' (unchecked), 'Require at least one category per pattern' (unchecked), and 'Don't include surface/POS form when a category matches' (checked) are present. An 'Add...' button for categories is also shown.
- Extract:** A button to start the extraction process. The 'Name' is 'features' and the 'Rare Threshold' is '5'.
- Feature Table:** A panel for viewing the extracted features, currently empty.
- Evaluations to Display:** A section for selecting metrics. Under 'Basic Table Statistics', 'Correlation', 'F-Score', 'Kappa', 'Precision', 'Recall', 'Target Hits', and 'Total Hits' are listed with checkboxes.
- Features in Table:** A panel for searching through the features, currently empty.

At the bottom of the window, there is a 'Report a Bug' link and a memory usage indicator showing '0.0 GB used, 2.7 GB max'.

Configuring Stretchy Patterns

Configure Stretchy Patterns

Pattern Length **2**

Gap Length **3**

[About Stretchy Patterns](#)

☒ Include surface words in patterns **4**

☐ Include POS tags in patterns

Categories:

5

☐ Require at least one category per pattern

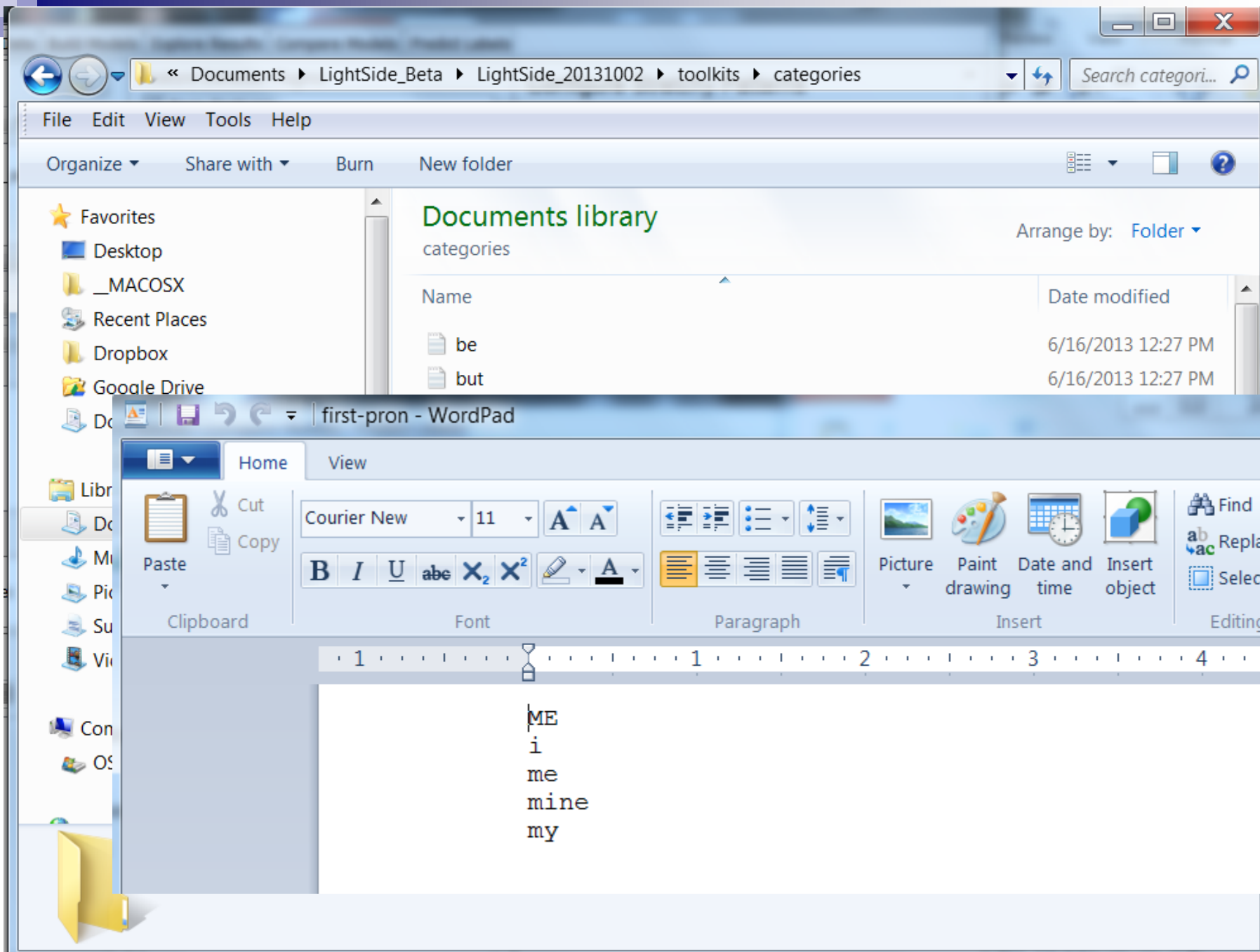
☒ Don't include surface/POS form when a category matches **6**

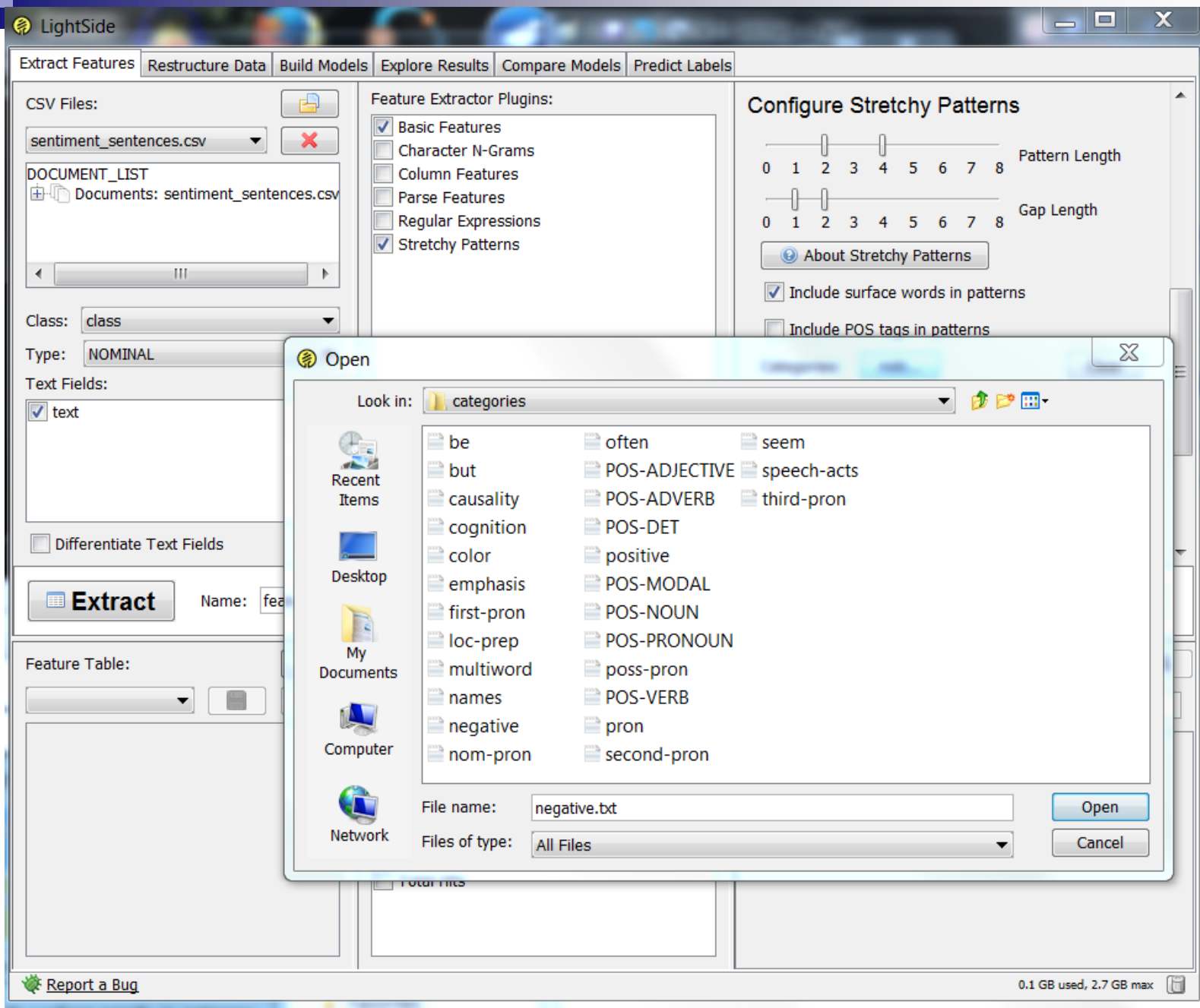
☒ Categories match against surface words

☐ Categories match against POS tags

☐ Count pattern hits **7**

- Longer patterns and longer gaps lead to larger numbers of features
- Categories are useful both for abstraction and for anchoring the patterns





CSV Files:

sentiment_sentences.csv

DOCUMENT_LIST

Documents: sentiment_sentences.csv

Class: class

Type: NOMINAL

Text Fields:

☒ text

☐ Differentiate Text Fields

Feature Extractor Plugins:

- ☒ Basic Features
- ☐ Character N-Grams
- ☐ Column Features
- ☐ Parse Features
- ☐ Regular Expressions
- ☒ Stretchy Patterns

☐ Include POS tags in patterns

Categories:

Add...

Clear

STRONG-NEG: [awful, bad, badly, disgusting, horrible, poorly, terrible, worse, worst]

☒ Require at least one category per pattern

☐ Don't include surface/POS form when a category matches

☒ Categories match against surface words

☐ Categories match against POS tags

☐ Count pattern hits

Prune Rare Features after N documents:



Extract

Name: features

Rare Threshold: 5

Feature Table:

Evaluations to Display:

Target:

Basic Table Statistics

- ☐ Correlation
- ☐ F-Score
- ☐ Kappa
- ☐ Precision
- ☐ Recall
- ☐ Target Hits
- ☐ Total Hits

Features in Table:

Search:

CSV Files:

sentiment_sentences.csv

DOCUMENT_LIST

Documents: sentiment_sentences.csv

Class: class

Type: NOMINAL

Text Fields:

☒ text☐ Differentiate Text Fields

Feature Extractor Plugins:

- ☒ Basic Features
- ☐ Character N-Grams
- ☐ Column Features
- ☐ Parse Features
- ☐ Regular Expressions
- ☒ Stretchy Patterns

☐ Include POS tags in patterns

Categories:

Add...

Clear

STRONG-NEG: [awful, bad, badly, disgusting, horrible, poorly, terrible, worse, worst]

☒ Require at least one category per pattern☐ Don't include surface/POS form when a category matches☒ Categories match against surface words☐ Categories match against POS tags☐ Count pattern hits

Prune Rare Features after N documents:

0 100 200 500 1000

Extract

Name: features1

Rare Threshold: 5

Feature Table:

features

FEATURE_TABLE

- Documents: sentiment_sentences.csv
- Feature Plugins:
- Feature Table: features

Evaluations to Display:

Target: neg

Basic Table Statistics

- ☐ Correlation
- ☐ F-Score
- ☐ Kappa
- ☐ Precision
- ☐ Recall
- ☐ Target Hits
- ☐ Total Hits

Features in Table:

Search:

Feature

be [GAP] STRONG-NEG
is [GAP] STRONG-NEG
it 's [GAP] STRONG-NEG
it [GAP] STRONG-NEG
not [GAP] STRONG-NEG
of [GAP] STRONG-NEG
so STRONG-NEG
the STRONG-NEG
the STRONG-NEG [GAP] of
the STRONG-NEG [GAP] the

Baseline Model:

logit

TRAINED_MODEL

- + Documents: sentiment_sentences.csv
- + Feature Plugins:
- + Feature Table: unigrams
- + Learning Plugin: Logistic Regression
- + Trained Model: logit

Competing Model:

logit1

TRAINED_MODEL

- + Documents: sentiment_sentences.csv
- + Feature Plugins:
- + Feature Table: features
- + Learning Plugin: Logistic Regression
- + Trained Model: logit1

Comparison Plugin: Basic Model Comparison

Baseline Model Metrics:

Metric	Value
Accuracy	0.7605
Kappa	0.5209

Competing Model Metrics:

Metric	Value
Accuracy	0.7609
Kappa	0.5219

Baseline Confusion Matrix:

Act \ Pred	neg	pos
neg	4089	1242
pos	1312	4019

Competing Confusion Matrix:

Act \ Pred	neg	pos
neg	4087	1244
pos	1305	4026

Insignificant improvement (p=0.579, t=-0.556)

Regular Expressions

- ◆ * allows the previous part of the regex to repeat, but it is not necessary.
- ◆ + is the same, but requires the previous part to match at least once.
- ◆ ? allows the previous part to happen either once or not at all, but does not match further.
- ◆ . is a wildcard, matching any one character.
- ◆ Certain character classes are predefined, like \w (any character A-Z), \d (any digit 0-9), and \s (any type of space character).

The screenshot shows a software interface with several panels. The top navigation bar includes tabs: 'Extract Features', 'Restructure Data', 'Build Models', 'Explore Results', 'Compare Models', and 'Predict Labels'. The 'Extract Features' tab is active.

CSV Files: A dropdown menu shows 'sentiment_sentences'. Below it, a 'DOCUMENT_LIST' section shows a tree view with 'Documents: sentiment_sentences'.

Class: A dropdown menu shows 'class'.

Type: A dropdown menu shows 'NOMINAL'.

Text Fields: A checkbox labeled 'text' is checked.

Feature Extractor Plugins: A list of plugins is shown, with 'Regular Expressions' selected and marked with a yellow circle with the number 1.

Configure Regular Expressions: A text input field contains 'but.*bad', marked with a yellow circle with the number 2. Below it, a list of matches is shown: 'but.*bad' (highlighted in blue and marked with a yellow circle with the number 3) and 'good|great|awesome' (marked with a yellow circle with the number 4). A 'Regex Cheat Sheet' button is visible, marked with a yellow circle with the number 5. A checkbox labeled 'Count Occurrences' is shown, marked with a yellow circle with the number 6.

American Street Gangs

Predict gang affiliation from posts

- **Crips, Bloods, Hoovers**
 - crips started in South Central LA
 - Pirus, Bloods, Hoovers from crips
- Chicago based
 - People Nation
 - **vice lords, latin kings, stones**
 - Folk nation
 - **gangster disciples**
- **Trinitarios**
 - hispanic gang based in NYC



Graffiti Based Style Features



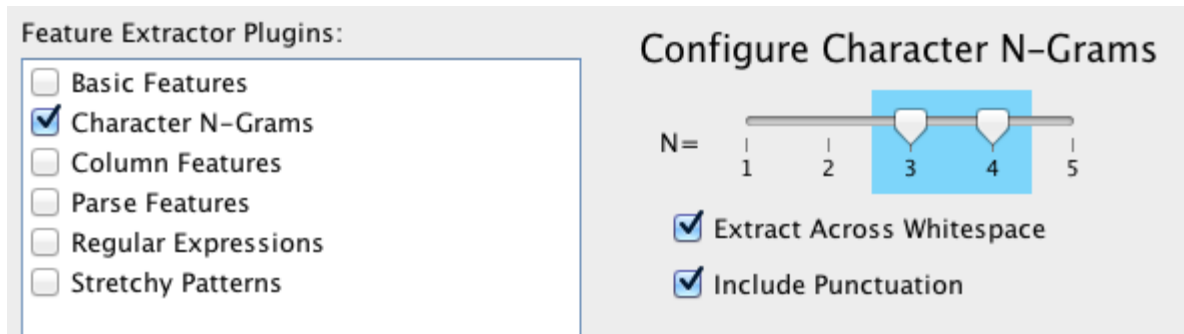
Graffiti

Social messages
Stylistic writing
crossing out other gangs

On the board

c	ck ckcrab, ckcome
ck	cc fucc, blocc
p	pk pkut, ...
h	hk whky, hkappens
b	bk bk1, bkang
e	3 3ast
s	5 5hit
c	c^ c^rime, c^uh

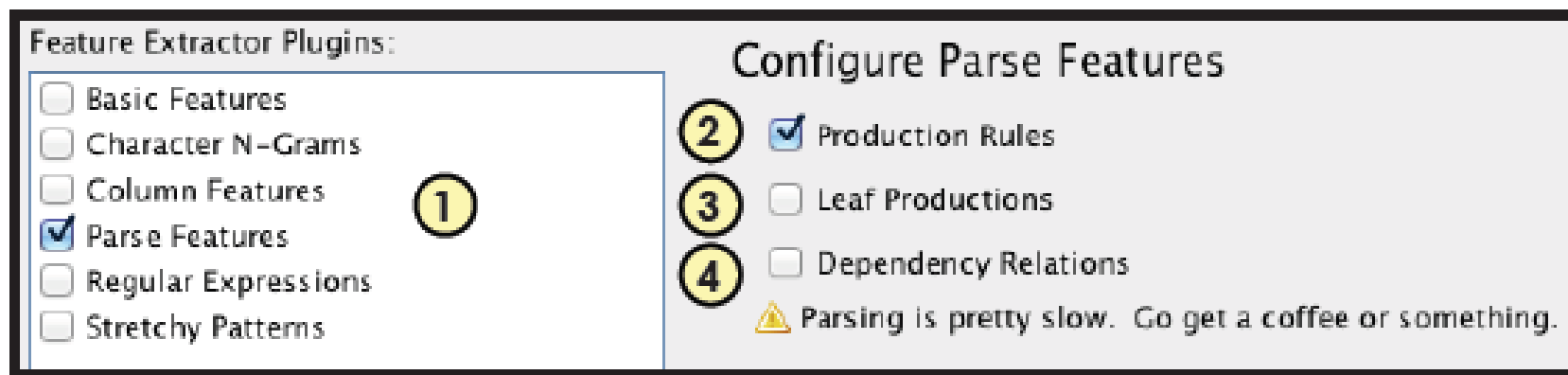
Character N-grams



- Character bigrams can detect graffiti style features
- Could also be used to identify consistent endings on words (i.e., that indicate formality or gender)

Parse Features

- Word based features lose all structure and order within sentences
- Parse features can capture that
- But they are SLOW!!





Leveraging Subpopulations through Multi-Level Modeling

Evaluation

The screenshot shows the LightSide software interface with the 'Explore Results' tab selected. The 'Feature Tables' panel on the left shows 'features2' as the selected table. The 'Learning Plugin' panel in the center shows 'Naive Bayes' as the selected plugin. The 'Configure Naive Bayes' panel on the right shows 'Use Kernel Estimator' and 'Use Supervised Discretization' as options. The 'Cross-Validation' option is selected for evaluation, with 'Fold Assignment' set to 'By Annotation' and 'Gender' as the annotation. The 'Number of Folds' is set to 'Auto'. The 'Train' button is visible, and the 'Name' field is set to 'bayes1'. The 'Trained Models' panel at the bottom left shows 'bayes' as the trained model. The 'Model Evaluation Metrics' panel at the bottom center shows the following metrics:

Metric	Value
Accuracy	0.7558
Kappa	0.5109

The 'Model Confusion Matrix' panel at the bottom right shows the following matrix:

Act \ Pred	Negative	Positive
Negative	379	129
Positive	101	333

The status bar at the bottom indicates '0.2 GB used, 2.7 GB max'.

Evaluation

The screenshot shows the LightSide software interface with the 'Explore Results' tab selected. The 'Feature Tables' panel on the left shows 'features2' as the selected table. The 'Learning Plugin' panel in the center shows 'Naive Bayes' as the selected plugin. The 'Configure Naive Bayes' panel on the right shows 'Use Kernel Estimator' and 'Use Supervised Discretization' as options. The 'Cross-Validation' option is selected for evaluation, with 'Fold Assignment' set to 'By Annotation' and 'Gender' as the annotation. The 'Number of Folds' is set to 'Auto'. The 'Train' button is visible, and the 'Name' field is set to 'bayes1'. The 'Trained Models' panel at the bottom left shows 'bayes' as the trained model. The 'Model Evaluation Metrics' panel at the bottom center shows the following metrics:

Metric	Value
Accuracy	0.7558
Kappa	0.5109

The 'Model Confusion Matrix' panel at the bottom right shows the following matrix:

Act \ Pred	Negative	Positive
Negative	379	129
Positive	101	333

The status bar at the bottom indicates '0.2 GB used, 2.7 GB max'.

Evaluation

The screenshot shows the LightSide software interface with the 'Explore Results' tab selected. The 'Learning Plugin' section is configured with 'Naive Bayes' and 'Cross-Validation'. The 'Fold Assignment' is set to 'Random', and the 'Number of Folds' is set to 'Manual: 2'. A red dashed box highlights the 'Cross-Validation' section. The 'Configure Naive Bayes' section on the right has 'Use Kernel Estimator' and 'Use Supervised Discretization' unchecked. The 'Train' button is visible, and the 'Trained Models' section shows 'bayes1' as the trained model. The 'Model Evaluation Metrics' table shows Accuracy of 0.7654 and Kappa of 0.5292. The 'Model Confusion Matrix' table shows the following data:

Act \ Pred	Negative	Positive
Negative	389	119
Positive	102	332

The status bar at the bottom indicates '0.3 GB used, 2.7 GB max'.

Evaluation

The screenshot displays the LightSide software interface, which is used for machine learning model evaluation. The interface is divided into several sections:

- Navigation Tabs:** Extract Features, Restructure Data, Build Models, Explore Results, Compare Models, Predict Labels.
- Baseline Model:** Set to 'bayes'. The TRAINED_MODEL list includes: Documents: Gallup.csv, Feature Plugins, Feature Table: features2, Learning Plugin: Naive Bayes, and Trained Model: bayes.
- Competing Model:** Set to 'bayes1'. The TRAINED_MODEL list includes: Documents: Gallup.csv, Feature Plugins, Feature Table: features2, Learning Plugin: Naive Bayes, and Trained Model: bayes1.
- Comparison Plugin:** Basic Model Comparison.
- Baseline Model Metrics:**

Metric	Value
Accuracy	0.7558
Kappa	0.5109
- Competing Model Metrics:**

Metric	Value
Accuracy	0.7654
Kappa	0.5292
- Baseline Confusion Matrix:**

Act \ Pred	Negative	Positive
Negative	379	129
Positive	101	333
- Competing Confusion Matrix:**

Act \ Pred	Negative	Positive
Negative	389	119
Positive	102	332
- Statistical Significance:** Insignificant improvement ($p=0.455$, $t=-0.747$)
- Footer:** Report a Bug, 0.3 GB used, 2.7 GB max.



Why is performance different?

- Men and women used language differently
- Different focus
 - Women had a more personal focus
 - Men had a more national/objective focus

What is different in how men and women talk?

- Word-based features capture more content than style, and are thus vulnerable to domain specificity.

male	female
linux	shopping
microsoft	mom
gaming	cried
server	freaked
software	pink

(Schler 2006)



What is different in how men and women talk?

- Women's language as "deviant" - *Lakoff (1975)* or "more varied" - *Chambers (1992)*
- Extrathematic details in conversational storytelling - time and location (male), people and speech acts (female).
Johnstone (1993)
"...after a full three years..."
"...he would sit and talk to my mother..."

What is different in how men and women talk?

- Hedging, qualifiers, and intensifiers -
“I think I might have said ...”
“So he brought to me...”
“I’m sometimes so jealous of people”
- “like” particle - gender variations in placement and usage
Iyeiri, Yaguchi, Okabe (2005)
*“...and then, we asked **like** four and one...”*
*“**Like**, instead of advanced, basic, proficient, and whatever...”*



Confounded with other variables

- Men sound older and women sound younger (Argamon et al., 2007)
- Men sound more like non-fiction and women sound more like fiction (Argamon et al., 2003)



Why do low level features overfit?

- In a linear model, positive weights push the decision towards one class while negative weights push the decision towards the other class
- The magnitude of the weight indicates how much of a push that feature gives

Why do low level features overfit?

- What happens if the same feature predicts age, gender, and social class?
 - If you are predicting gender, then the average value for each feature assumes the mix of age and social class in the data set you trained for
 - The weights normalize for this mix
 - If the mix changes, then the normalization will be wrong
 - So the weights won't predict gender correctly anymore on datasets where the mix of those other factors is different



Train

A word cloud of three-letter codes arranged in a circular pattern. The codes are: FYH, MYH, MOL, MYL, and FOH. The frequency of each code is as follows:

Code	Frequency
FYH	10
MYH	10
MOL	10
MYL	10
FOH	10

Test

A word cloud of chemical formulas arranged in a circular pattern. The formulas include FYL, MOH, FOL, MOL, and MYH. The words are scattered throughout the circle, with some appearing more frequently than others. For example, FYL and MOH appear multiple times, while MYH appears fewer times. The arrangement is non-uniform, with some words clustered together and others isolated.

Evaluation of Domain Generality

Occupation	Unigram	Unigram + Bigram	POS Bigram	Stretchy Patterns
Engineering	49.5 (-.01)	53 (.06)	49 (-.02)	50.5 (.01)
Education	49 (-.02)	52 (.04)	54.5 (.09)	51 (.02)
Internet	55.5 (.11)	47.5 (-.05)	55.5 (.11)	56.5 (.13)
Law	51.5 (.03)	46.5 (-.07)	46.5 (-.07)	50.5 (.01)
Non-Profit	50 (0)	54 (.08)	49 (-.02)	51. (.02)
Technology	50 (0)	53.5 (.07)	50 (0)	51.5 (.03)
Arts	48 (-.04)	46.5 (-.07)	51 (.02)	55.4 (.11)
Media	53 (.06)	50 (0)	45 (-.1)	51.5 (.02)
Science	52 (.04)	48 (-.04)	40.5 (-.19)	59.5 (.19)
Student	51 (.02)	46 (-.09)	55 (.10)	62 (.24)
Average	50.95 (.002)	49.7 (-.007)	49.6 (.01)	53.94 (.08)
Random CV	61.05 (.22)	59.65 (.19)	57.95 (.16)	62.8 (.26)

- Contrast random CV and leave-one-occupation-out CV
- All feature space representations show significant drop between random CV and leave-one-occupation-out CV
- Only stretchy patterns remain significantly above random performance

Feature Splitting (Daumé III, 2007)

General



General

Domain A

Domain B

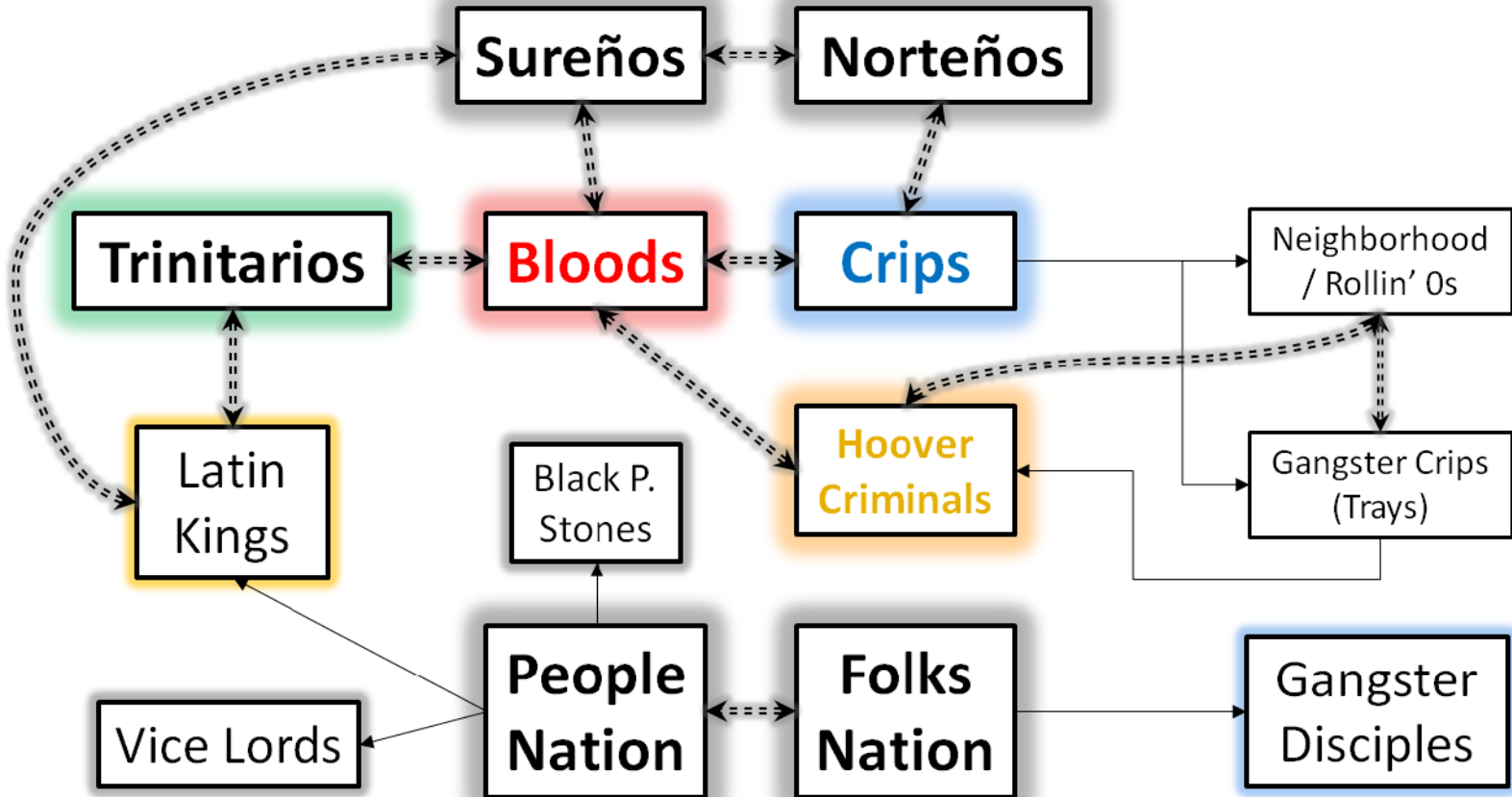
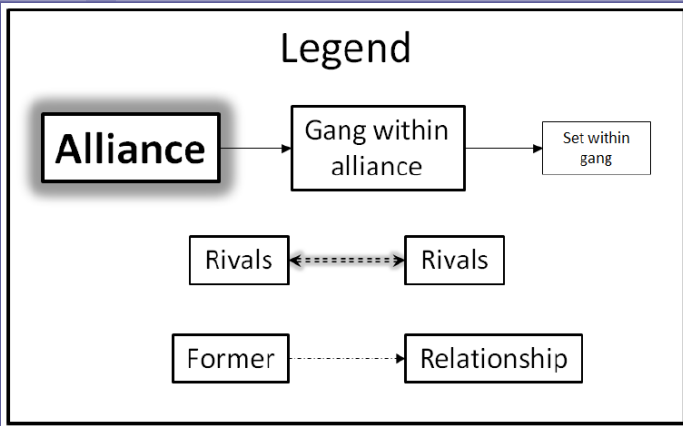


Why is this nonlinear?

It represents the interaction between each feature and the Domain variable

Now that the feature space represents the nonlinearity, the algorithm to train the weights can be linear.

Gang Alliances



Gangs Data

LightSide

Extract Features | Restructure Data | Build Models | Explore Results | Compare Models | Predict Labels

CSV Files:

ThreadCompStyles.csv

DOCUMENT_LIST

Documents: ThreadCompStyles.csv

Class: composition

Type: NOMINAL

Text Fields:

☐ 3e

☐ 5s

☐ 6b

☐ bCaret

☐ bk

☐ Differentiate Text Fields

Feature Extractor Plugins:

☐ Basic Features

☐ Character N-Grams

☒ Column Features

☐ Parse Features

☐ Regular Expressions

☐ Stretchy Patterns

Configure Column Features

Column Name

☒ 3e

☒ 5s

☒ 6b

☒ bCaret

☒ bk

☒ cCaret

☒ cc

☒ ck

☐ dominant

☒ hCaret

☒ hk

☒ length

☒ numUsers

☒ pCaret

☒ pk

☐ starter

☐ subfeature

All None

Extract Name: features1 Rare Threshold: 5

Feature Table:

features

FEATURE_TABLE

Documents: ThreadCompStyles.csv

Feature Plugins:

Feature Table: features

Evaluations to Display:

Target: allied

Basic Table Statistics

☐ Correlation

☐ F-Score

☐ Kappa

☐ Precision

☐ Recall

☐ Target Hits

☐ Total Hits

Features in Table:

Search:

Feature

3e_column

5s_column

6b_column

bCaret_column

bk_column

cCaret_column

cc_column

ck_column

hk_column

length_column

Report a Bug

0.1 GB used, 2.7 GB max

Extract Features

Restructure Data

Build Models

Explore Results

Compare Models

Predict Labels

Feature Tables:

features

FEATURE_TABLE

Documents: ThreadCompStyles.csv

Feature Plugins:

Feature Table: features

Filters Available:

☐ Combine Features

☐ Filter Feature Values

☒ Multilevel Modeling

☐ Regroup Instances

Configure Multilevel Modeling

Select Levels:

Domain	A	B
<input checked="" type="checkbox"/> dominant	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/> length	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/> 6b	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/> bCaret	<input type="checkbox"/>	<input type="checkbox"/>
<input type="checkbox"/> starter	<input type="checkbox"/>	<input type="checkbox"/>

A*B

A[B]

Select Features in Level:

Feature Source

☐ All Column Features

☒ 3e__column

☒ 5s__column

☒ 6b__column

☒ bCaret__column

All

None

v

Add Domain

v

Select Levels:

Restructure

Name: :structure

Rare Threshold: 5

Restructured Tables:

Evaluations to Display:

Target:

☐ Correlation

☐ F-Score

☐ Kappa

☐ Precision

☐ Recall

☐ Target Hits

☐ Total Hits

Features in Table:

Search:

Report a Bug

0.0 GB used, 2.7 GB max

Extract Features

Restructure Data

Build Models

Explore Results

Compare Models

Predict Labels

Feature Tables:

features

FEATURE_TABLE

Documents: ThreadCompStyles.csv

Feature Plugins:

Feature Table: features

Filters Available:

☐ Combine Features
☐ Filter Feature Values
☒ Multilevel Modeling
☐ Regroup Instances

Configure Multilevel Modeling

Select Levels:

Domain	A	B
dominant	<input type="checkbox"/>	<input type="checkbox"/>
length	<input type="checkbox"/>	<input type="checkbox"/>
6b	<input type="checkbox"/>	<input type="checkbox"/>
bCaret	<input type="checkbox"/>	<input type="checkbox"/>
starter	<input type="checkbox"/>	<input type="checkbox"/>

A*B

A[B]

Add Domain

Select Features in Level:

Feature Source

All Column Features

3e_column

5s_column

6b_column

bCaret_column

All

None

Select Levels:

dominant

Intercepts

Slopes

Features

Restructure

Name: structure1

Rare Threshold: 5

Restructured Tables:

restructure

MODIFIED_TABLE

Documents: ThreadCompStyles.csv

Feature Plugins:

Feature Table: features

Restructure Plugins:

Restructured Table: restructure

Evaluations to Display:

Target: allied

Basic Table Statistics

☐ Correlation
☐ F-Score
☐ Kappa
☐ Precision
☐ Recall
☐ Target Hits
☐ Total Hits

Features in Table:

Search:

Feature

dominant::bloods_cc_column

dominant::bloods_ck_column

dominant::bloods_length_column

dominant::bloods_numUsers_column

dominant::bloods_xo_column

dominant::crisps (Intercept)

dominant::crisps_bk_column

dominant::crisps_cCaret_column

dominant::crisps_cc_column

dominant::crisps_ck_column

Report a Bug

0.0 GB used, 2.7 GB max

LightSide

Extract Features | Restructure Data | Build Models | Explore Results | Compare Models | Predict Labels

Feature Tables:

features

FEATURE_TABLE

- Documents: ThreadCompStyles.csv
- Feature Plugins:
- Feature Table: features

Learning Plugin:

- ☐ Naive Bayes
- ☒ Logistic Regression
- ☐ Linear Regression
- ☐ Support Vector Machines
- ☐ Decision Trees
- ☐ Weka (All)

☒ Cross-Validation
☐ Supplied Test Set
☐ No Evaluation

Fold Assignment:

- ☒ Random
- ☐ By Annotation:
- 3e
- ☐ By File

Number of Folds:

- ☒ Auto
- ☐ Manual: 10

2 5 10 Max

Configure Logistic Regression

- ☒ L2 Regularization
- ☐ L1 Regularization
- ☐ L2 Regularization (Dual)

Train

Name: logit1

☐ Feature Selection

Trained Models:

logit

TRAINED_MODEL

- Documents: ThreadCompStyles.csv
- Feature Plugins:
- Feature Table: features
- Learning Plugin: Logistic Regression
- Trained Model: logit

Model Evaluation Metrics:

Metric	Value
Accuracy	0.4817
Kappa	0.188

Model Confusion Matrix:

Act \ Pred	allied	homogeneous	mixed	opposing
allied	20	39	2	179
homogeneous	9	107	3	54
mixed	4	20	3	43
opposing	12	84	3	290

Report a Bug

0.1 GB used, 2.7 GB max

LightSide

Extract FeaturesRestructure DataBuild ModelsExplore ResultsCompare ModelsPredict Labels

Feature Tables:

restructure

MODIFIED_TABLE

Documents: ThreadCompStyles.csv

Feature Plugins:

Feature Table: features

Restructure Plugins:

Restructured Table: restructure

Learning Plugin:

☐ Naive Bayes

☒ Logistic Regression

☐ Linear Regression

☐ Support Vector Machines

☐ Decision Trees

☐ Weka (All)

☒ Cross-Validation

☐ Supplied Test Set

☐ No Evaluation

Fold Assignment:

☒ Random

☐ By Annotation:

3e

☐ By File

Number of Folds:

☒ Auto

☐ Manual: 10

2

5

10

Max

Configure Logistic Regression

☒ L2 Regularization

☐ L1 Regularization

☐ L2 Regularization (Dual)

Train

Name: logit2

☐ Feature Selection

Trained Models:

logit1

TRAINED_MODEL

Documents: ThreadCompStyles.c

Feature Plugins:

Feature Table: features

Restructure Plugins:

Restructured Table: restructure

Model Evaluation Metrics:

Metric	Value
Accuracy	0.6124
Kappa	0.4108

Model Confusion Matrix:

Act \ Pred	allied	homogeneous	mixed	opposing
allied	87	24	3	126
homogeneous	21	132	4	16
mixed	8	5	21	36
opposing	48	40	7	294

Report a Bug

0.1 GB used, 2.7 GB max

LightSide

Extract FeaturesRestructure DataBuild ModelsExplore ResultsCompare ModelsPredict Labels

Baseline Model:

logit

TRAINED_MODEL

- Documents: ThreadCompStyles.csv
- Feature Plugins:
- Feature Table: features
- Learning Plugin: Logistic Regression
- Trained Model: logit

Competing Model:

logit1

TRAINED_MODEL

- Documents: ThreadCompStyles.csv
- Feature Plugins:
- Feature Table: features
- Restructure Plugins:
- Restructured Table: restructure
- Learning Plugin: Logistic Regression
- Trained Model: logit1

Comparison Plugin: Basic Model Comparison

Baseline Model Metrics:

Metric	Value
Accuracy	0.4817
Kappa	0.188

Competing Model Metrics:

Metric	Value
Accuracy	0.6124
Kappa	0.4108

Baseline Confusion Matrix:

Act \ Pred	allied	homogeneous	mixed	opposing
allied	20	39	2	179
homogeneous	9	107	3	54
mixed	4	20	3	43
opposing	12	84	3	290

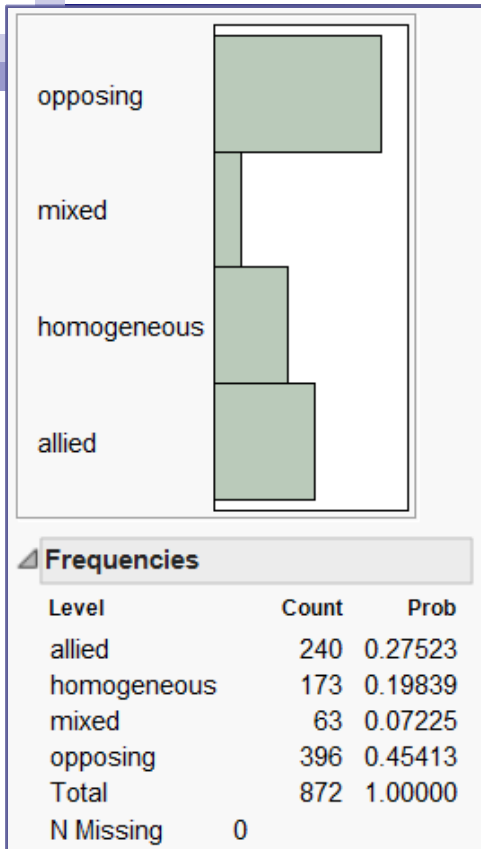
Competing Confusion Matrix:

Act \ Pred	allied	homogeneous	mixed	opposing
allied	87	24	3	126
homogeneous	21	132	4	16
mixed	8	5	21	36
opposing	48	40	7	294

Highly significant improvement (p=0**, t=-7.024)

Report a Bug

0.1 GB used, 2.7 GB max



Feature Analysis

- Style features that distinguish Allied from Opposing differ by dominant gang

- **Crips:**

- Allied: **bCaret**
 - Opposing: CC, PK, cCaret

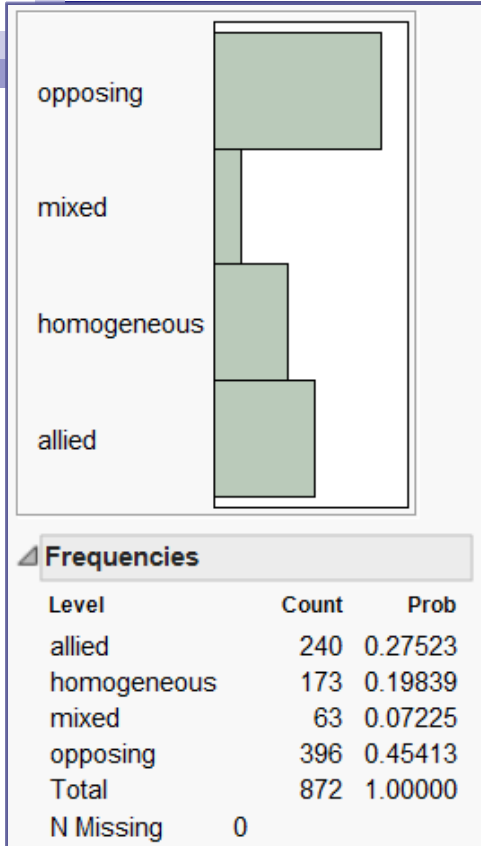
When the dominant gang is in an allied thread, we see style features that unite them against opposing gangs.

- **Bloods:**

- Allied: **XO, CC**
 - Opposing: hCaret, BK

- **Latin Kings:**

- Allied: CC, **XO**
 - Opposing: 5S



Feature Analysis

- Style features that distinguish Allied from Opposing differ by dominant gang

- **Crips:**

- ☐ Allied: bCaret
- ☐ Opposing: **CC**, **PK**, **cCaret**

When the dominant gang is in an opposing thread, we also see features that unite the opposing gangs against them.

- **Bloods:**

- ☐ Allied: XO, CC
- ☐ Opposing: hCaret, **BK**

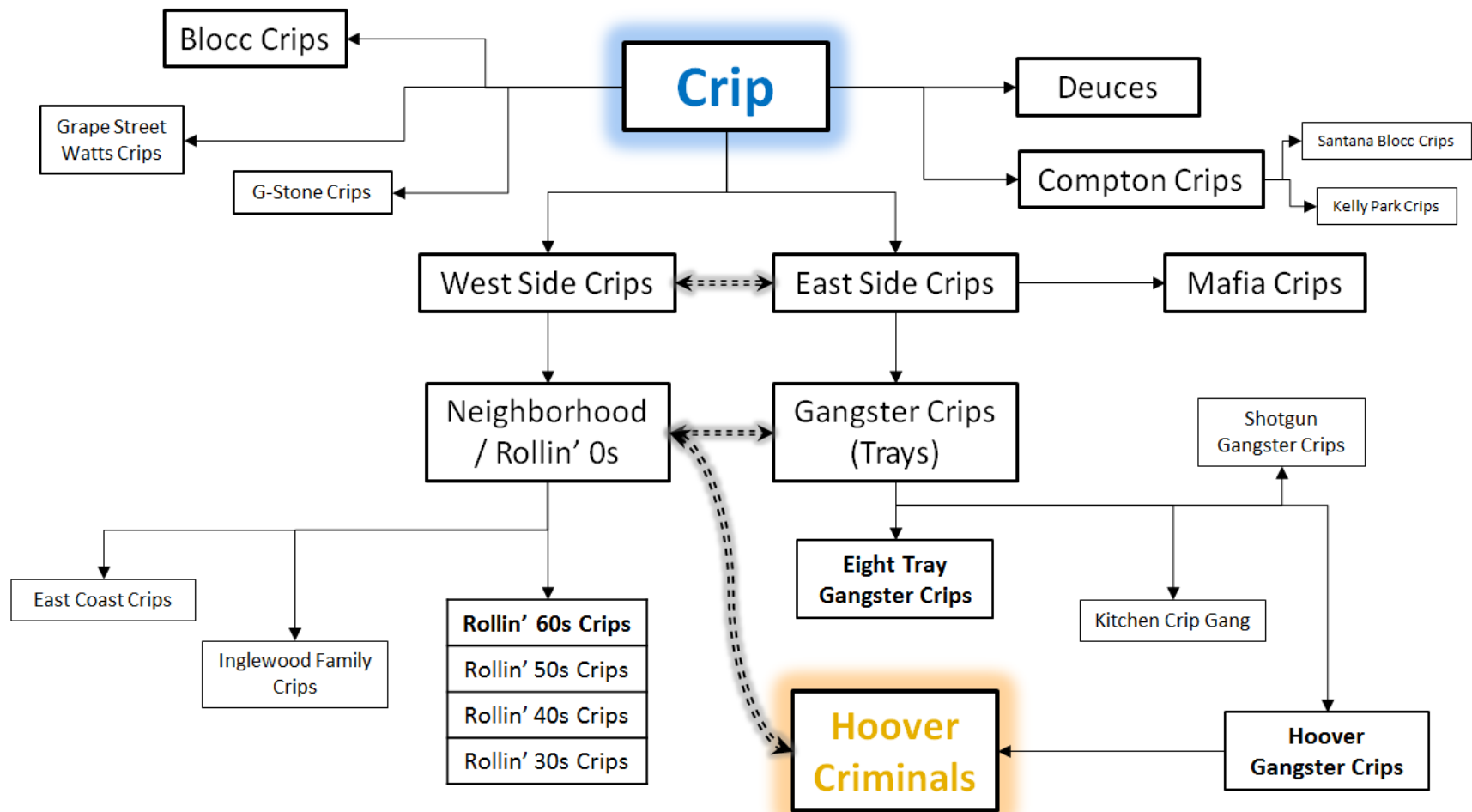
- **Latin Kings:**

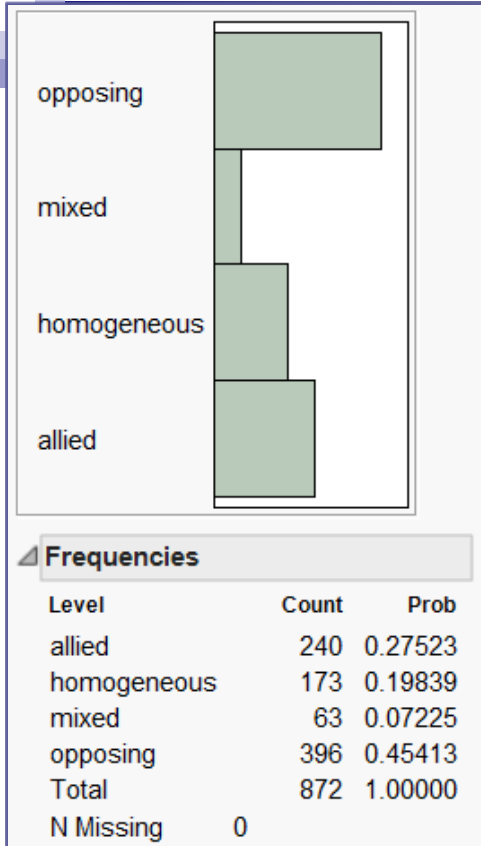
- ☐ Allied: CC, XO
- ☐ Opposing: **5S**

opposing

mixed

Feature Analysis





Feature Analysis

- Unigram features that distinguish Allied from Opposing don't differ by dominant gang as much as style features

- **Universal:**

- ☐ Allied: Imao, you, crew
- ☐ Opposing: forever, wtf, where

We see relationship words, but not gang identity words.

- **Crips:**

- ☐ Allied: lol
- ☐ Opposing: know, about

- **Bloods:**

- ☐ Allied: niggas, the
- ☐ Opposing: at