



A CLOUD GURU

Text Feature Engineering



Brock Tubre

INSTRUCTOR



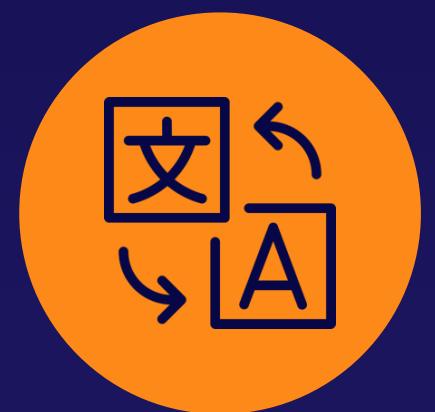
Feature Engineering

Transforming attributes within our data to make them more useful within our model for the problem at hand. Feature engineering is often compared to an art.



Text Feature Engineering

Transforming text within our data so Machine Learning algorithms can better analyze it.



Splitting text into bite size pieces.

Text Feature Engineering

AAPS PharmSciTech, Vol. 17, No. 1, February 2016 (© 2016)
DOI: 10.1208/s12249-016-0486-2



Editorial

Theme: Pharmaceutical Thermal Processing
Guest Editors: Feng Zhang and Michael A. Repka

Pharmaceutical Thermal Processing

Feng Zhang^{1,3} and Michael A. Repka²

published online 9 February 2016

Thermal processing has flourished as a novel technology to prepare a wide range of drug delivery systems over the past decade. Thermal processing is built on the principle of efficient conversion of the mechanical energy into the thermal energy of the formulation matrix. Single and twin screw extruders are the most commonly employed thermal processing equipment. Kinetisol® Dispersing technology, an innovative thermal process, was invented recently. This theme issue of AAPS PharmSciTech provides a broad overview of (a) the history and recent advancement of thermal processing, (b) scientific approaches to develop drug delivery systems using thermal processing, and (c) examples of a wide range of drug delivery systems prepared using thermal processing.

This issue of AAPS PharmSciTech offers three review articles. Martin provided a historical and technical review of a twin screw extruder as a continuous mixer for thermal processing in both the plastic and pharmaceutical industry from an equipment engineer's perspective (1). From his personal experience, Martin believed that the evolution of twin screw extrusion-based thermal processing in the pharmaceutical industry resembles what occurred for the plastics industry in 1980s. Patil et al. reviewed the application of melt extrusion in the production of diverse pharmaceutical delivery systems (2). The authors emphasized the critical role of the quality-by-design approach and process analytical technology in pharmaceutical melt extrusion monitoring and control. LaFountaine et al. discussed a number of formulation and process-driven strategies to enable thermal processing of challenging compositions (3). These included the use of traditional plasticizers and surfactants, temporary plasticizers utilizing sub- or supercritical carbon dioxide, designer polymers, and KinetiSol® Dispersing technology.

This theme issue also contains 15 original research articles. Bhagurkar et al. proposes the use of melt extrusion processing for the preparation of a polyethylene glycol base ointment (4). A

¹College of Pharmacy, University of Texas at Austin, 2409 University Ave, Austin, Texas 78712, USA.

²Department of Pharmaceutics and Drug Delivery, School of Pharmacy, University of Mississippi, University, Oxford, Mississippi 38677, USA.

³ To whom correspondence should be addressed. (e-mail: feng.zhang@austin.utexas.edu)

```
ogue: d92-1  
er of utterances files: 54  
th of dialogue: 189.448858  
nated number of turns: 32
```

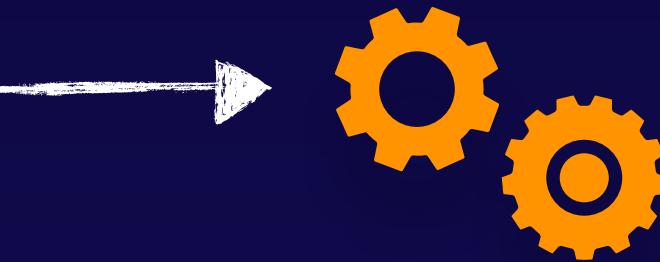
: s: hello <sil> can I help you
: u: yeah I want t- I want to determine the maximum number of boxcars of
oranges <sil> that I can get to Bath <sil> by seven a.m. <sil>
tomorrow morning
: so <brth> hm <sil>
so I guess all the boxcars will have to go through oran- <sil>
through Corning because that's where the orange juice <brth>
orange factory is
: so from Corning to Bath how far is that
: s: two hours
: u: and it's gonna take us also an hour to load <sil> boxcars right
: s: right + +
: u: + okay + so <sil> hm so <sil> every trip will take at least <sil>
three hours <sil> then
:
um
0 : s: right we can unload any amount of cargo onto a train in one hour
1 : so we can + <sil> do a maximum of three + boxcars in an hour
2 : u: + right <sil> okay +
3 : okay <sil> so I guess one thing we can do oh <brth> so <brth>
I guess one thing is that we should see how many boxcars we can actually
get to Corning in four hours
4 : um how far is it from Avon to Bath <sil> to Corning
5 : s: <click> <brth> that's six hours it's + shorter + through Dansville
6 : u: + okay +
7 : from Avon oh no but the thing is is that I was thinking if s- <sil>
I was wondering if we could actually pick up those two boxcars which are
at <sil> + Bath +
8 : s: + mm +
9 : u: but I think that we can't even put them in <sil> to <brth>
that they won't even come in to play
0 : unless how far is it from Elmira <brth> to Bath
1 : s: two hours
2 : u: oh really <sil> so then w- we could actually take <brth>
like engine E two
3 : have it go to bath <sil> pick up those two box cars there
4 : take it to Corn- how long from Bath to Corning an hour I guess
5 : s: two hours
6 : u: oh wait a second I thought <sil> well from from like Elmira <brth>
through to Corning to Bath is how many hours
7 : s: four hours + +
8 : u: + oh + that's four hours so wa- we're like screwed as far as tho- those
two box cars at Bath go
9 : like I don- I don't think we s- so to get the maximum number of box cars
of oranges <brth>
to bath by seven am I don't think we can even use those two at bath
0 : because there <sil> we because we can't get an engine <brth>
to Bath in time
1 : s: that's right
2 : u: okay <brth> so we have the three boxcars at <sil> Dansville
3 : so how far is it from Avon to Danville
4 : s: three hours

Example:

Raw Text: { “123 Main Street, Seattle, WA 98101” }

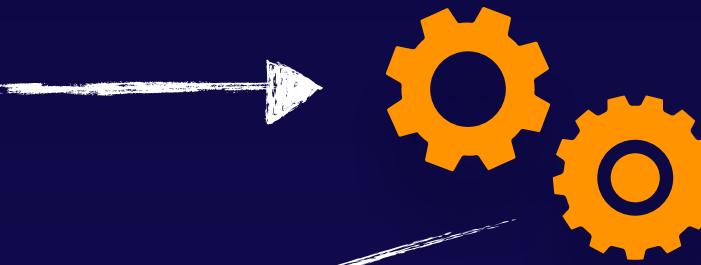
Example:

Raw Text: { “123 Main Street, Seattle, WA 98101” }



Example:

Raw Text: { “123 Main Street, Seattle, WA 98101” }

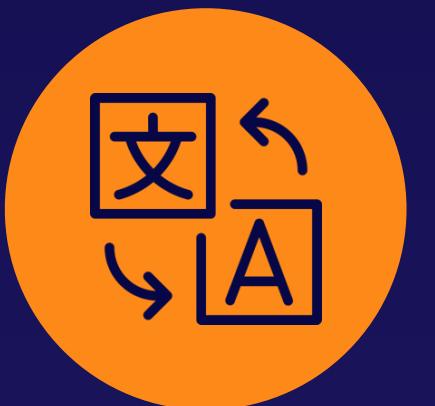


Address	City	State	Zip
123 Main Street	Seattle	WA	98101



Bag-of-Words

Tokenizes raw text and creates a statistical representation of the text.



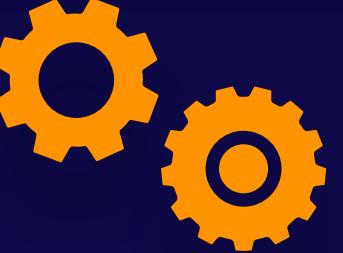
Breaks up text by white space into single words.



Example:

Raw Text: { “he is a jedi and he will save us” }

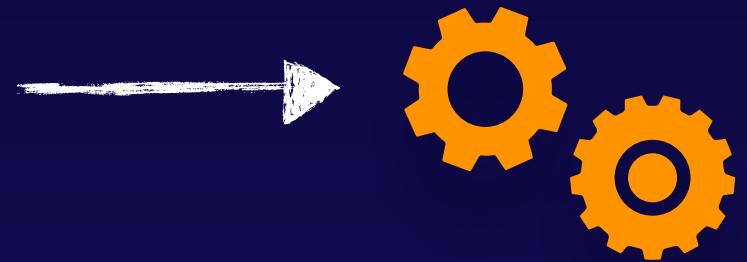
Example:

Raw Text: { “he is a jedi and he will save us” } → 

Bag-of-Words

Example:

Raw Text: { “he is a jedi and he will save us” }



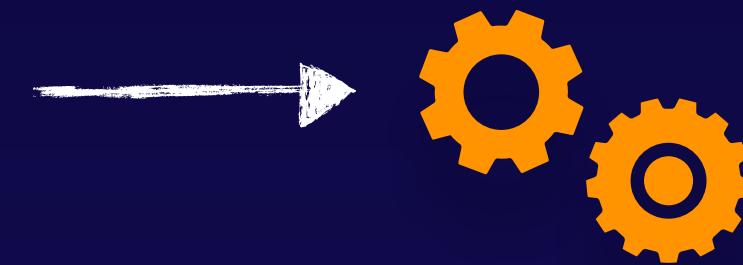
Bag-of-Words



{ “he”, “is”, “a”, “jedi”, “and”, “he”, “will”, “save”, “us” }

Example:

Raw Text: { “he is a jedi and he will save us” }



Bag-of-Words

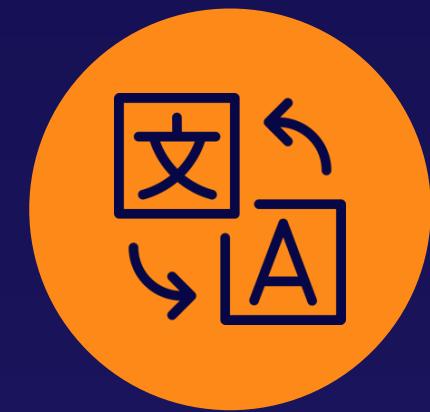
ID	word	count
1	he	2
2	is	1
3	a	1
4	jedi	1
5	and	1
6	will	1
7	save	1
8	us	1

{ “he”, “is”, “a”, “jedi”, “and”, “he”, “will”, “save”, “us” }



N-Gram

An extension of Bag-of-Words which produces groups of words of n size.



Breaks up text by white space into groups of words.

N-Gram Example (1-gram)

Example:

Raw Text: { “he is a jedi and he will save us” }

N-Gram Example (1-gram)

Example:

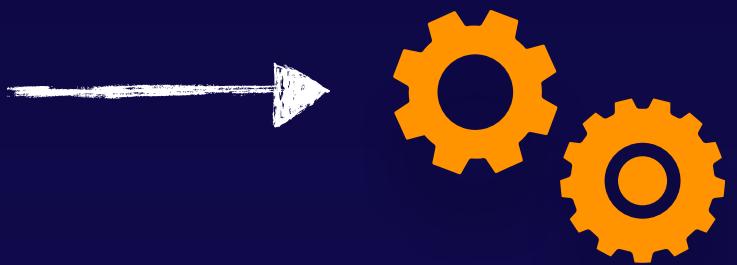
Raw Text: { “he is a jedi and he will save us” } → 

N-gram, size = 1

N-Gram Example (1-gram)

Example:

Raw Text: { “he is a jedi and he will save us” }



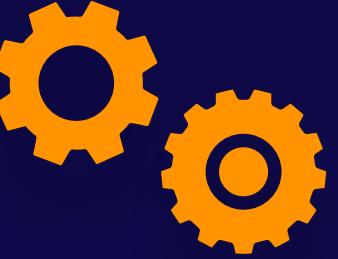
N-gram, size = 1



{ “he”, “is”, “a”, “jedi”, “and”, “he”, “will”, “save”, “us” }

N-Gram Example (2-gram)

Example:

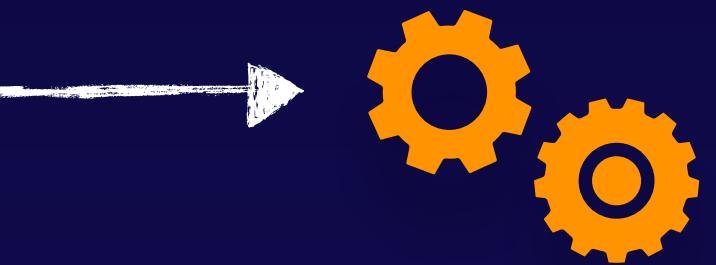
Raw Text: { “he is a jedi and he will save us” } → 

N-gram, size = 2

N-Gram Example (2-gram)

Example:

Raw Text: { “he is a jedi and he will save us” }



N-gram, size = 2



{ “he is”, “is a”, “a jedi”, “jedi and”, “and he”,
“he will”, “will save”, “save us”, “he”, “is”, “a”,
“jedi”, “and”, “he”, “will”, “save”, “us” }

N-Gram Example (2-gram)

Example:

Raw Text: { “he is a jedi and he will save us” } → 

“sliding window”

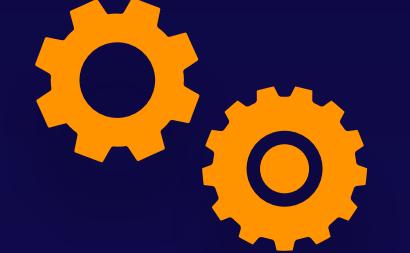
$n = 2$

N-gram, size = 2

{ “he is”, “is a”, “a jedi”, “jedi and”, “and he”,
“he will”, “will save”, “save us”, “he”, “is”, “a”,
“jedi”, “and”, “he”, “will”, “save”, “us” }

N-Gram Example (2-gram)

Example:

Raw Text: { “he is a jedi and he will save us” } → 

N-gram, size = 2



{ “he is”, “is a”, “a jedi”, “jedi and”, “and he”,
“he will”, “will save”, “save us”, “he”, “is”, “a”,
“jedi”, “and”, “he”, “will”, “save”, “us” }

N-Gram Example (2-gram)

Example:

Raw Text: { “he is a jedi and he will save us” } → 

N-gram, size = 2



{ “he is”, “is a”, “a jedi”, “jedi and”, “and he”,
“he will”, “will save”, “save us”, “he”, “is”, “a”,
“jedi”, “and”, “he”, “will”, “save”, “us” }

N-Gram Example (2-gram)

Example:

Raw Text: { “he is a jedi and he will save us” } → 

N-gram, size = 2

{ “he is”, “is a”, “a jedi”, “jedi and”, “and he”,
“he will”, “will save”, “save us”, “he”, “is”, “a”,
“jedi”, “and”, “he”, “will”, “save”, “us” }

N-Gram Example (2-gram)

Example:

Raw Text: { “he is a jedi **and he** will save us” } → 

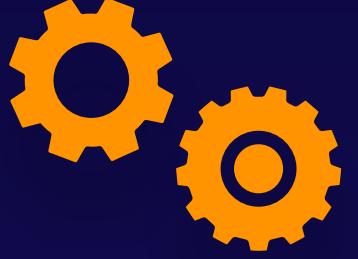
N-gram, size = 2



{ “he is”, “is a”, “a jedi”, “jedi **and**”, “**and he**”,
“he will”, “will save”, “save us”, “he”, “is”, “a”,
“jedi”, “and”, “he”, “will”, “save”, “us” }

N-Gram Example (2-gram)

Example:

Raw Text: { “he is a jedi and he will save us” } → 

N-gram, size = 2



{ “he is”, “is a”, “a jedi”, “jedi and”, “and he”,
“he will”, “will save”, “save us”, “he”, “is”, “a”,
“jedi”, “and”, “he”, “will”, “save”, “us” }

N-Gram Example (2-gram)

Example:

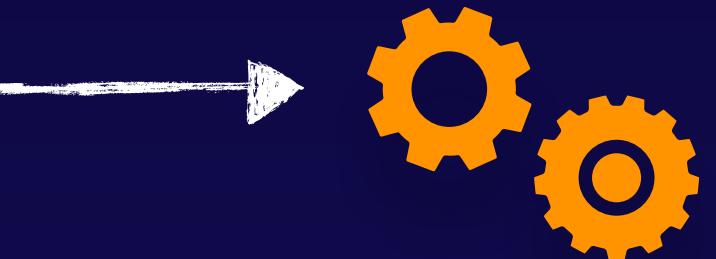


{ “he is”, “is a”, “a jedi”, “jedi and”, “and he”,
“he will”, “will save”, “save us”, “he”, “is”, “a”,
“jedi”, “and”, “he”, “will”, “save”, “us” }

N-Gram Example (2-gram)

Example:

Raw Text: { “he is a jedi and he will **save us**” }



N-gram, size = 2

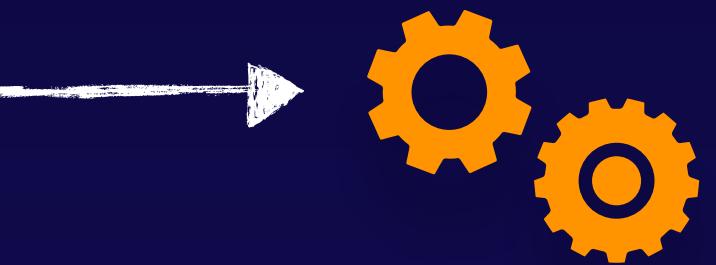


{ “he is”, “is a”, “a jedi”, “jedi and”, “and he”,
“he will”, “will save”, “**save us**”, “he”, “is”, “a”,
“jedi”, “and”, “he”, “will”, “**save**”, “us” }

N-Gram Example (2-gram)

Example:

Raw Text: { “he is a jedi and he will save us” }



N-gram, size = 2



{ “he is”, “is a”, “a jedi”, “jedi and”, “and he”,
“he will”, “will save”, “save us”, “he”, “is”, “a”,
“jedi”, “and”, “he”, “will”, “save”, “us” }

unigram = 1 word tokens

bigram = 2 word tokens

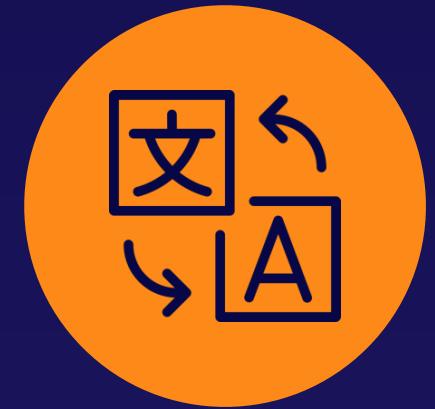
trigram = 3 word tokens

...



Orthogonal Sparse Bigram (OSB)

Creates groups of words of size n and outputs every pair of words that includes the first word.



Creates groups of words that always include the first word.

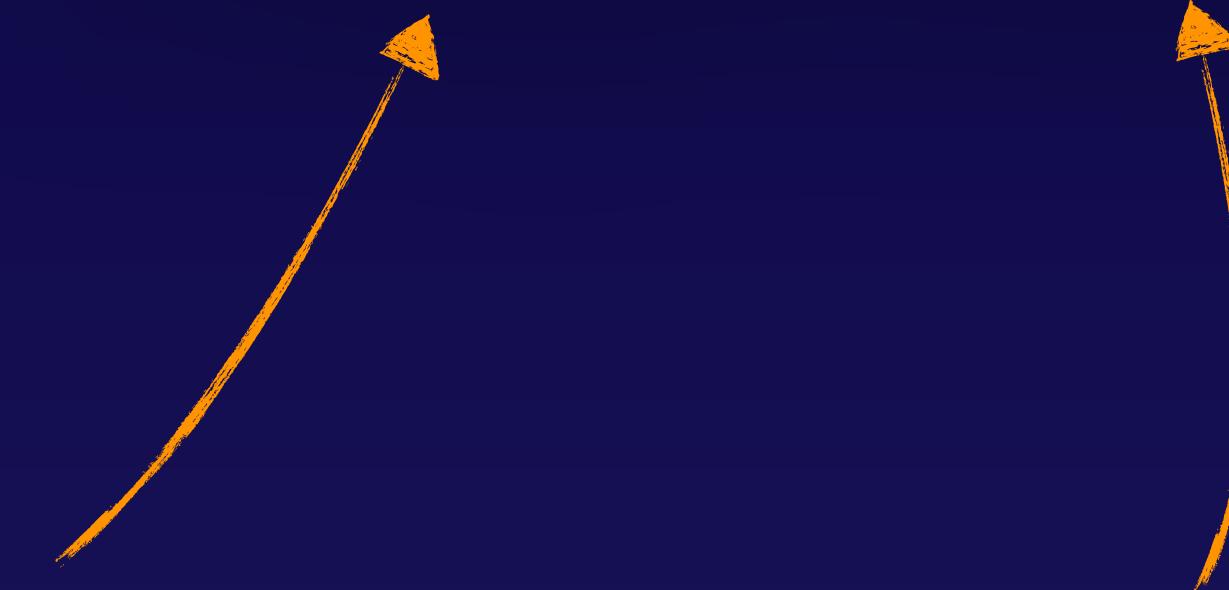
Orthogonal Sparse Bigram

Orthogonal Sparse Bigram



When two things are
independent of each other

Orthogonal Sparse Bigram



When two things are
independent of each other

scattered or thinly
distributed

Orthogonal Sparse Bigram



When two things are
independent of each other

scattered or thinly
distributed

2-gram or two words

Orthogonal Sparse Bigram

When two things are
independent of each other

scattered or thinly
distributed

2-gram or two words

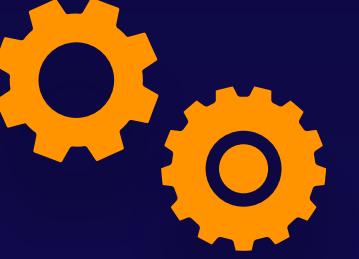
OSB is not “better” or “worse” than N-Gram,
it’s just another technique to use...

Example:

Raw Text: { “he is a jedi and he will save us” }

Orthogonal Sparse Bigram (OSB)

Example:

Raw Text: { “he is a jedi and he will save us” } → 
OSB, size = 4

Orthogonal Sparse Bigram (OSB)

Example:

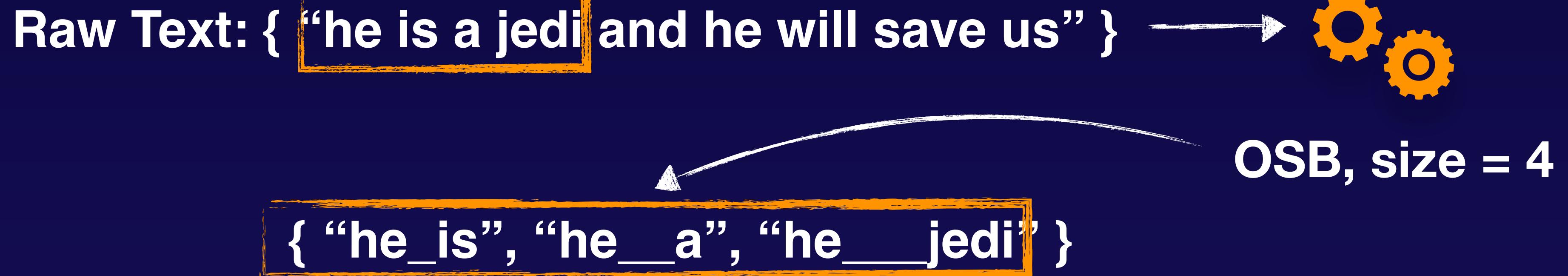
Raw Text: { “he is a jedi and he will save us” } → 

“sliding window”

$n = 4$

Orthogonal Sparse Bigram (OSB)

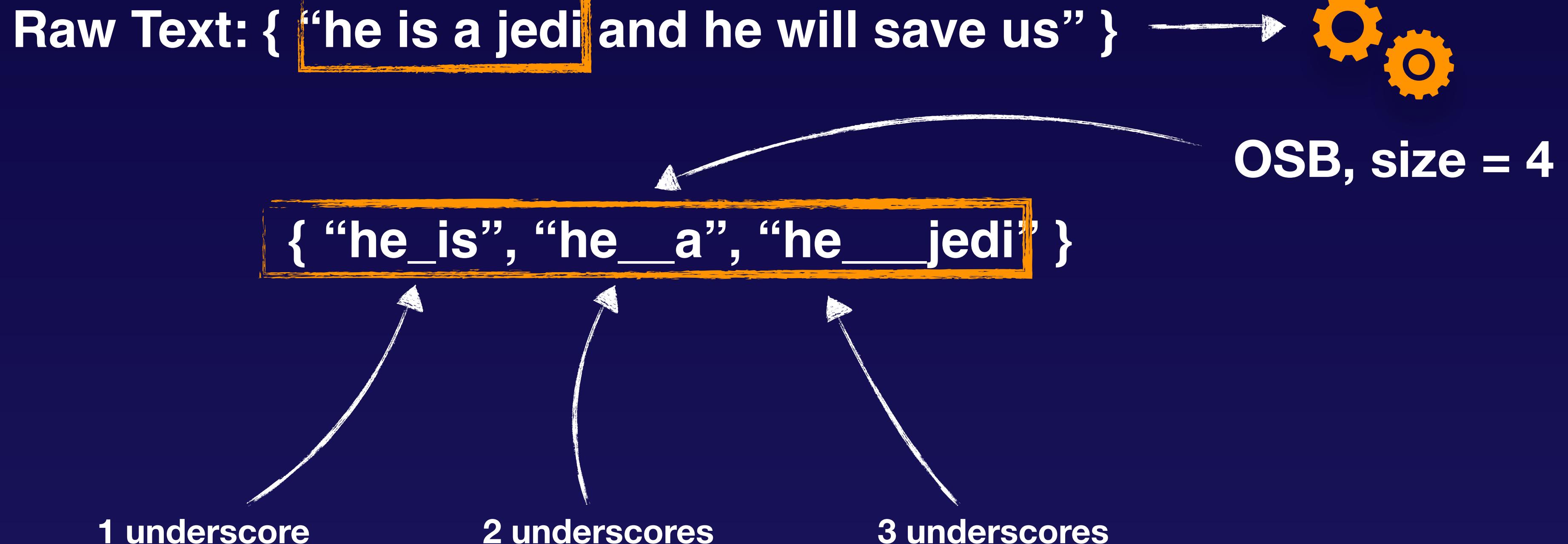
Example:



underscore = white space

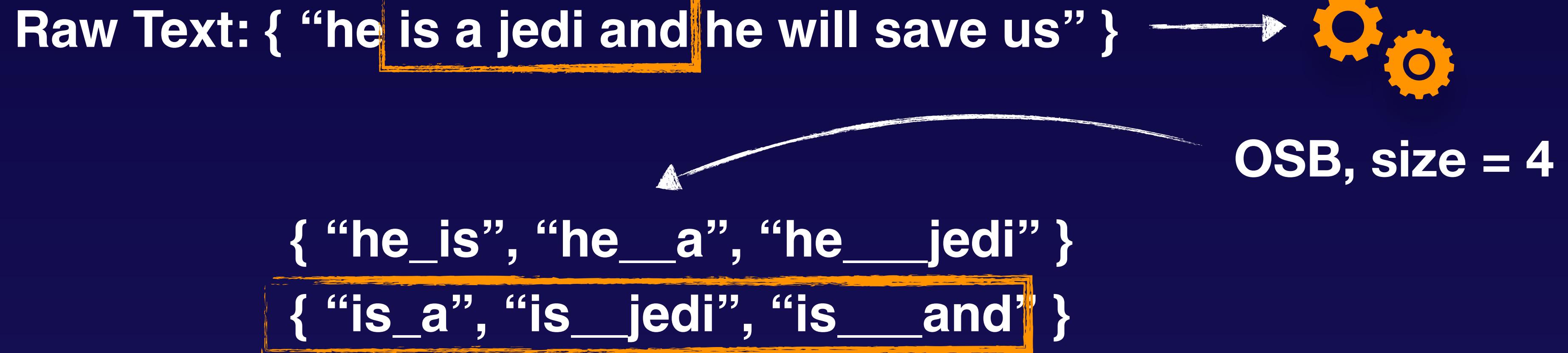
Orthogonal Sparse Bigram (OSB)

Example:



Orthogonal Sparse Bigram (OSB)

Example:



Orthogonal Sparse Bigram (OSB)

Example:

Raw Text: { “he is a jedi and he will save us” }



OSB, size = 4

{ “he_is”, “he_a”, “he_jedi” }

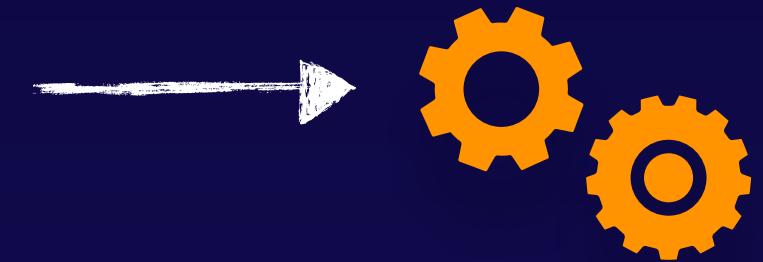
{ “is_a”, “is_jedi”, “is_and” }

{ “a_jedi”, “a_and”, “a_he” }

Orthogonal Sparse Bigram (OSB)

Example:

Raw Text: { “he is a jedi and he will save us” }



OSB, size = 4

{ “he_is”, “he_a”, “he_jedi” }

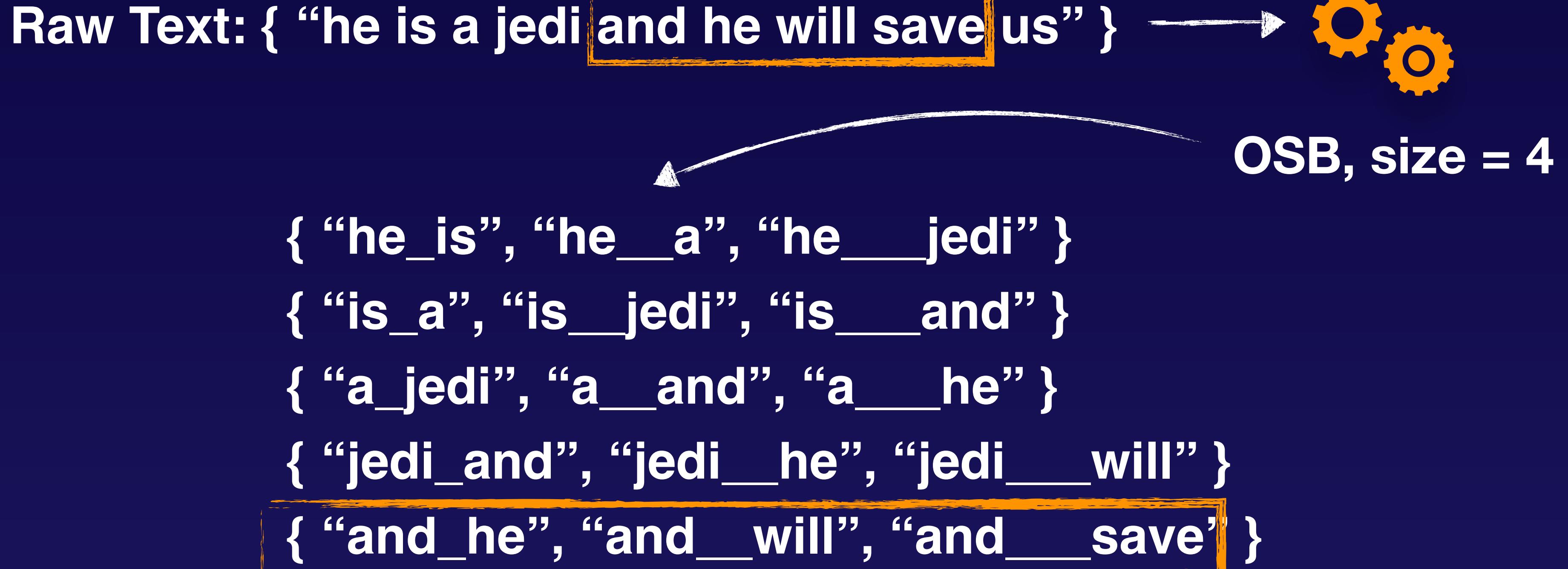
{ “is_a”, “is_jedi”, “is_and” }

{ “a_jedi”, “a_and”, “a_he” }

{ “jedi_and”, “jedi_he”, “jedi_will” }

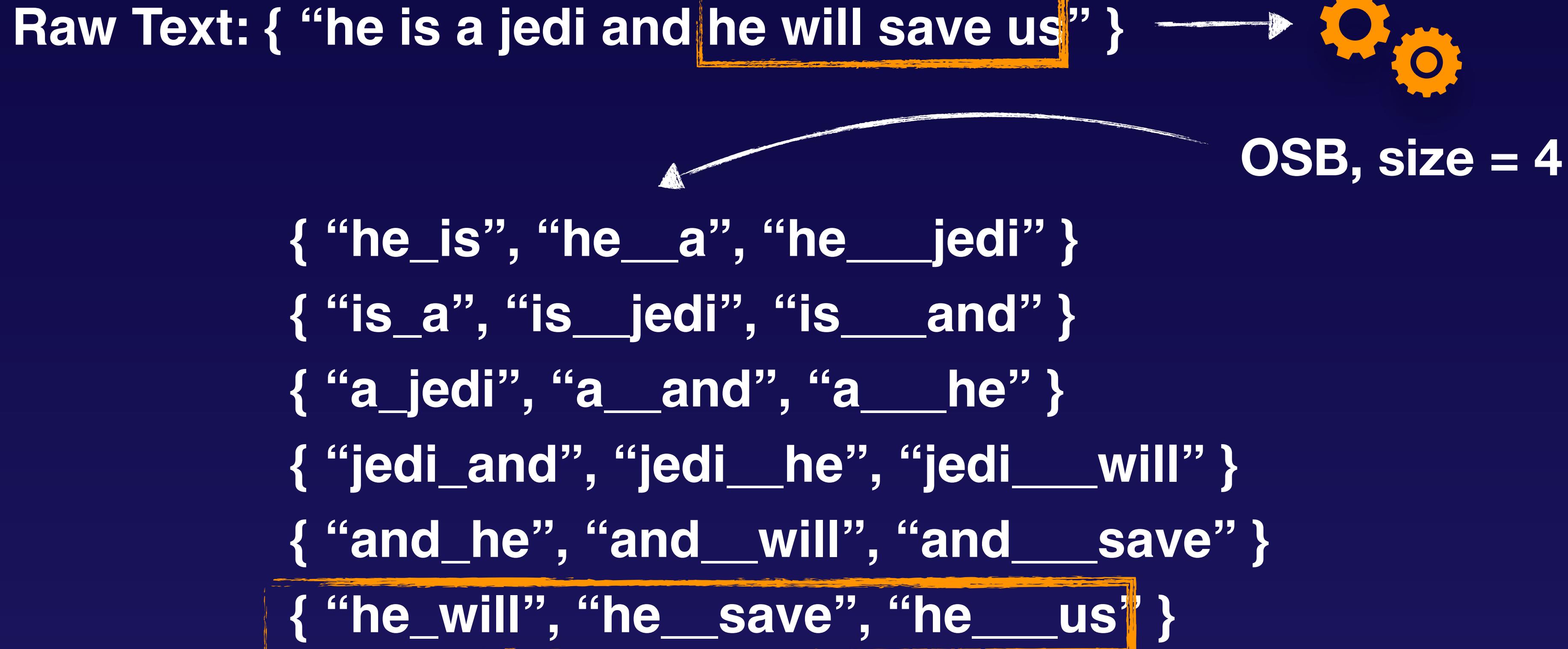
Orthogonal Sparse Bigram (OSB)

Example:



Orthogonal Sparse Bigram (OSB)

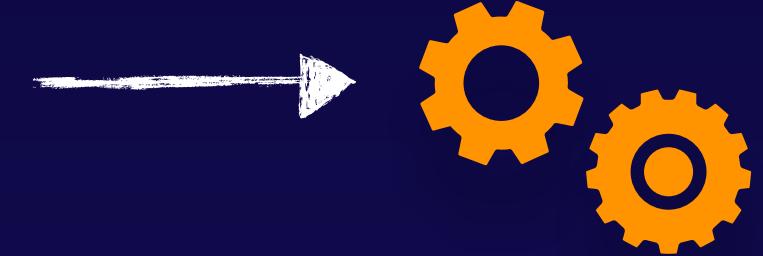
Example:



Orthogonal Sparse Bigram (OSB)

Example:

Raw Text: { “he is a jedi and he will save us” }



OSB, size = 4

{ “he_is”, “he_a”, “he_jedi” }

{ “is_a”, “is_jedi”, “is_and” }

{ “a_jedi”, “a_and”, “a_he” }

{ “jedi_and”, “jedi_he”, “jedi_will” }

{ “and_he”, “and_will”, “and_save” }

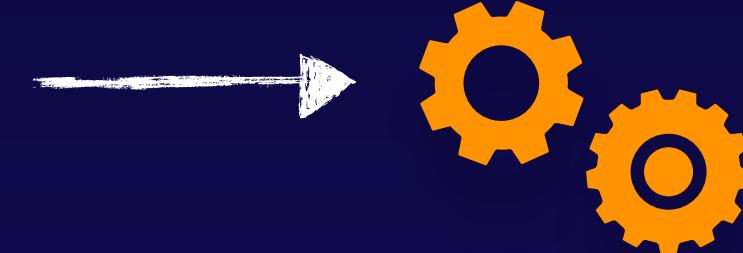
{ “he_will”, “he_save”, “he_us” }

{ “will_save”, “will_us” }

Orthogonal Sparse Bigram (OSB)

Example:

Raw Text: { “he is a jedi and he will **save us**” }



OSB, size = 4

{ “he_is”, “he_a”, “he_jedi” }

{ “is_a”, “is_jedi”, “is_and” }

{ “a_jedi”, “a_and”, “a_he” }

{ “jedi_and”, “jedi_he”, “jedi_will” }

{ “and_he”, “and_will”, “and_save” }

{ “he_will”, “he_save”, “he_us” }

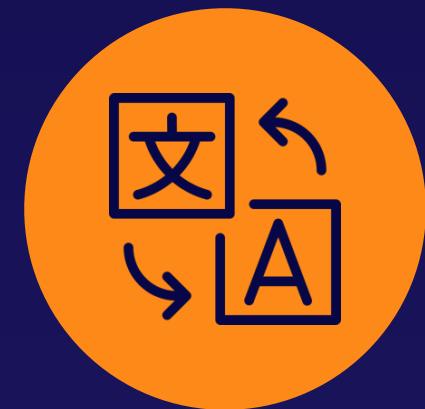
{ “will_save”, “will_us” }

{ “**save_us**” }



Term Frequency - Inverse Document Frequency (tf-idf)

Represents how important a word or words are to a given set of text by providing appropriate weights to terms that are common and less common in the text.



Shows us the popularity of a word or words in text data by making common words like “the” or “and” less important.

Term Frequency - Inverse Document Frequency

Term Frequency - Inverse Document Frequency

How frequent
does a word
appear.

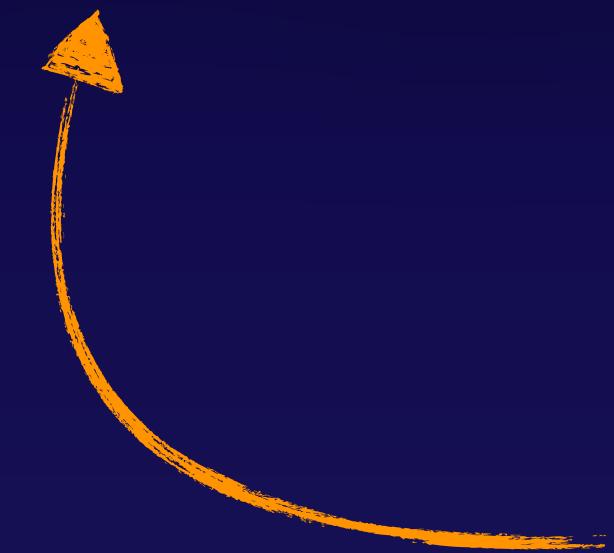


Term Frequency - Inverse Document Frequency

How frequent does a word appear.



Number of documents in which terms occur.



Term Frequency - Inverse Document Frequency

How frequent does a word appear.



Makes common words less meaningful.

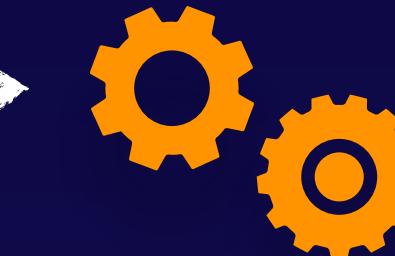


Number of documents in which terms occur.

Example:



Example:



tf-idf

Document 1

word	count
the	3
force	1
Luke	1
Skywalker	1
a	2

Document 2

word	count
the	2
jedi	1
a	1
empire	1



Example:



tf-idf

Document 1

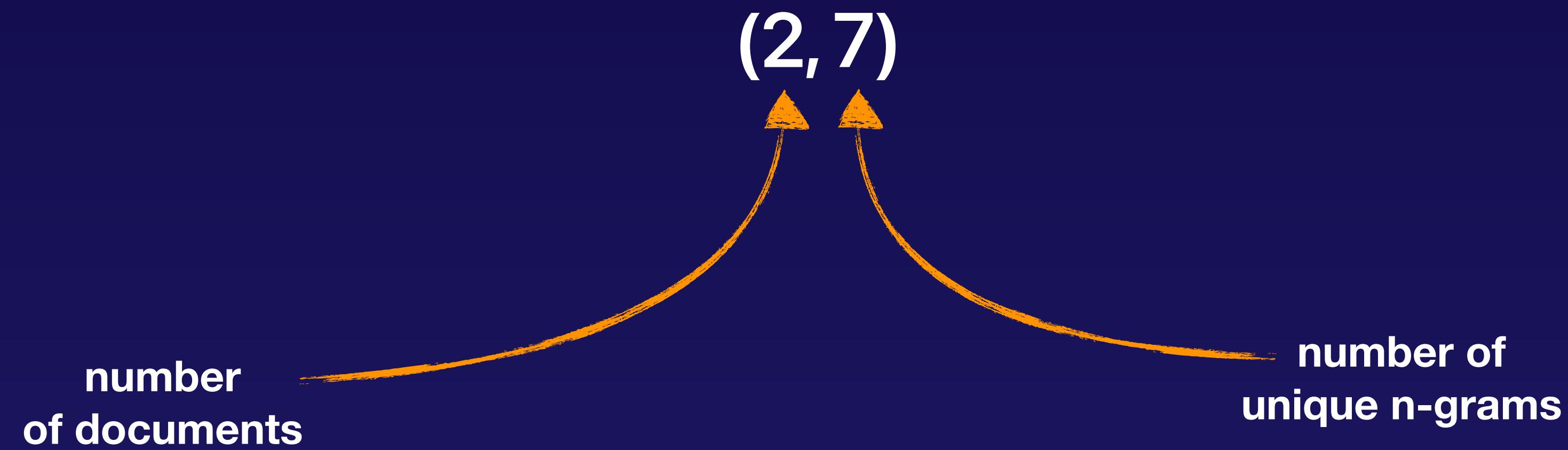
word	count
the	3
force	1
Luke	1
Skywalker	1
a	2

Document 2

word	count
the	2
jedi	1
a	1
empire	1

Since the tokens “the” and “a” showed up in **BOTH** documents many times, these are deemed as **less important**.

(number of documents, number of unique n-grams)



Problem

Matching phrases in spam emails

Transformation

N-Gram

Easier to compare whole phrases like “click here now” or “you’re a winner”.

Problem

Determining the subject matter of multiple PDFs

Transformation

Tf-idf
Orthogonal Sparse Bigram

Filter less important words in PDFs.
Find common word combinations repeated throughout PDFs.

- **Remove Punctuation**

Sometimes removing punctuation is a good idea if we do not need them.

- **Lowercase Transformation**

Using lowercase transformation can help standardize raw text.

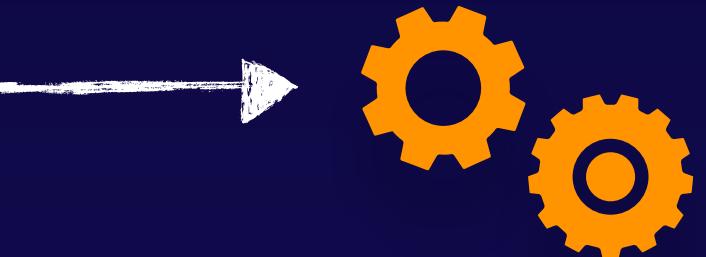
Example:

Raw Text: { “Help me, Obi-Wan Kenobi. You’re my only hope.” }

Removing Punctuation

Example:

Raw Text: { “Help me, Obi-Wan Kenobi. You’re my only hope.” }

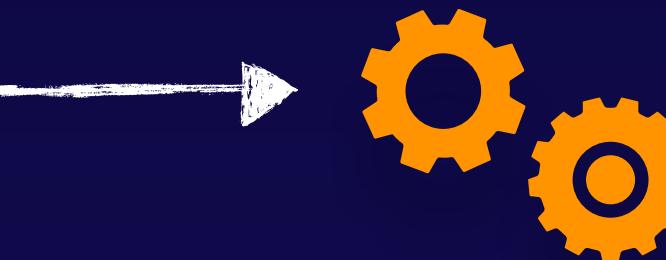


remove
punctuation

Removing Punctuation

Example:

Raw Text: { “Help me, Obi-Wan Kenobi. You’re my only hope.” }



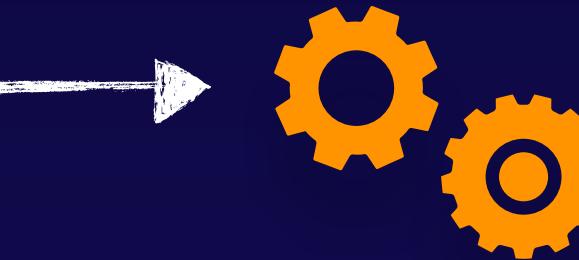
remove
punctuation

{ “Help”, “me”, “Obi-Wan”, “Kenobi”, “You’re”, “my”, “only”, “hope” }

Removing Punctuation

Example:

Raw Text: { “Help me, Obi-Wan Kenobi. You’re my only hope.” }



remove
punctuation

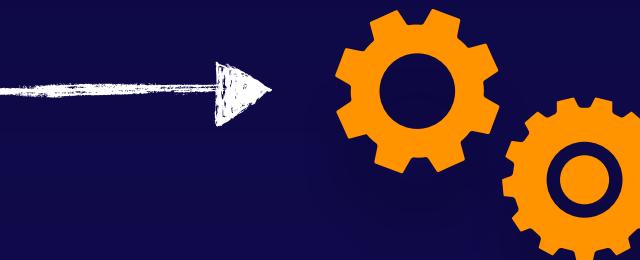
{ “Help”, “me”, “Obi-Wan”, “Kenobi”, “You’re”, “my”, “only”, “hope” }



Removing Punctuation

Example:

Raw Text: { “Help me, Obi-Wan Kenobi. You’re my only hope.” }



remove
punctuation

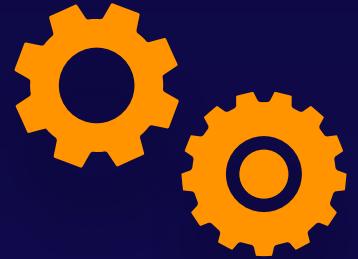
{ “Help”, “me”, “Obi-Wan”, “Kenobi”, “You’re”, “my”, “only”, “hope” }

Example:

Raw Text: { “I Find Your Lack of Faith Disturbing” }

Lowercase Transformation

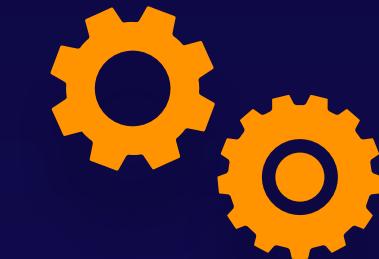
Example:

Raw Text: { “I Find Your Lack of Faith Disturbing” } → 
lowercase

Lowercase Transformation

Example:

Raw Text: { “I Find Your Lack of Faith Disturbing” }



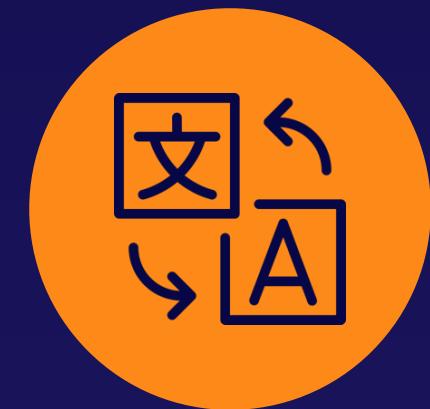
lowercase

{ “i find your lack of faith disturbing” }



Cartesian Product Transformation

Creates a new feature from the combination of two or more text or categorical values.



Combining sets of words together.

Cartesian Product



Cartesian Product



René Descartes

c.1596 – c.1650

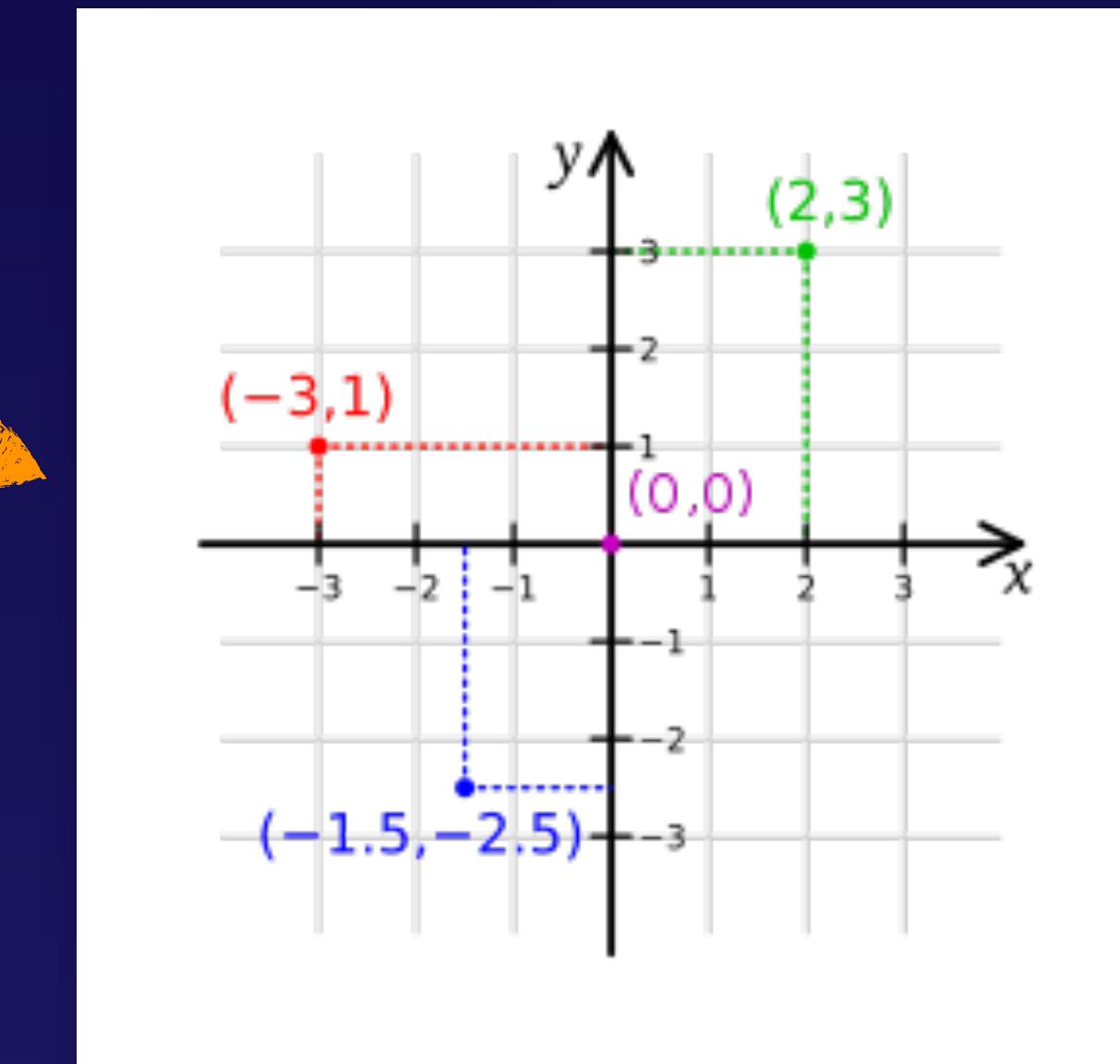
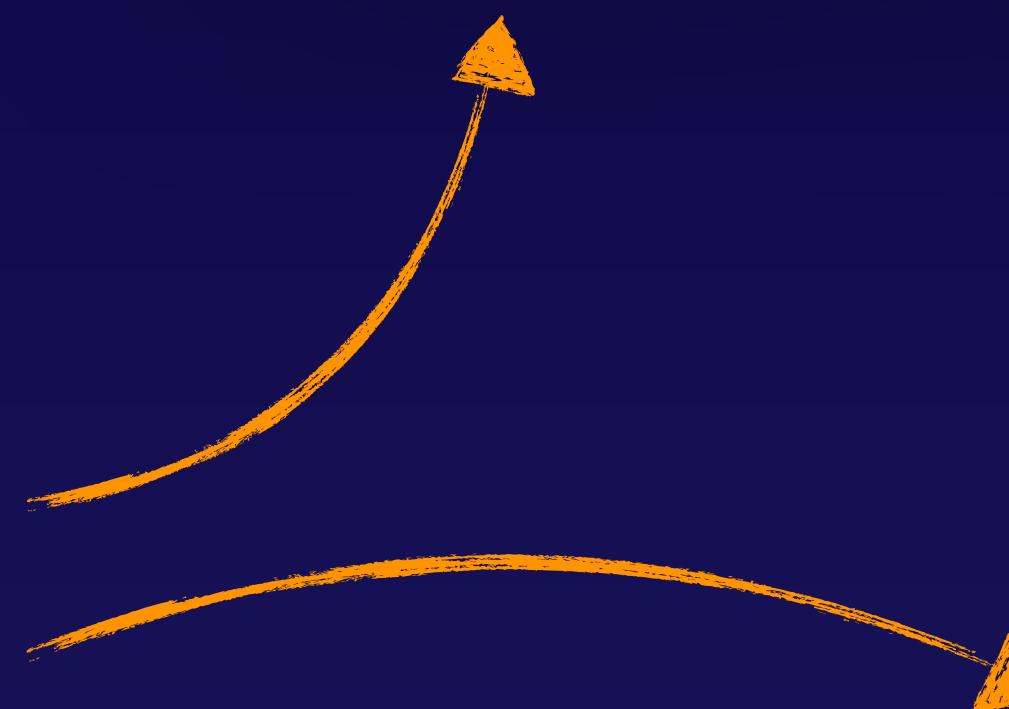
Latinized: Renatus Cartesius



René Descartes
c.1596 – c.1650

Latinized: Renatus Cartesius

Cartesian Product



Cartesian Product



Cartesian Product



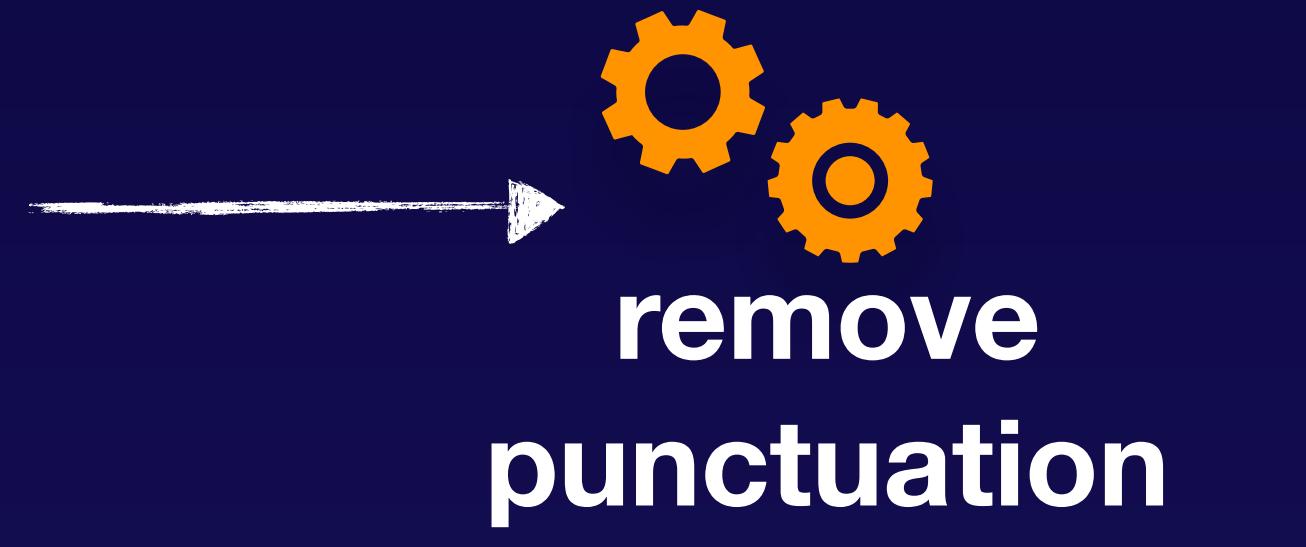
We are multiplying a set of words by another set of words to create a new feature...

Cartesian Product Example

ID	textbook	binding
1	Python Data Science Handbook	Softcover
2	Visualization Analysis & Design	Hardcover
3	Machine Learning Algorithms	Softcover

Cartesian Product Example

ID	textbook	binding
1	Python Data Science Handbook	Softcover
2	Visualization Analysis & Design	Hardcover
3	Machine Learning Algorithms	Softcover



ID	cartesian_product
1	{"Python_Softcover", "Data_Softcover", "Science_Softcover", "Handbook_Softcover"}
2	{"Visualization_Hardcover", "Analysis_Hardcover", "Design_Hardcover"}
3	{"Machine_Softcover", "Learning_Softcover", "Algorithms_Softcover"}



Feature Engineering Dates

Translating dates into useful information.



Feature Engineering Dates

Translating dates into useful information.



Learn from the best to
ensure success
Reasons we will be
successful

25 great
love to

DATE FORMATS

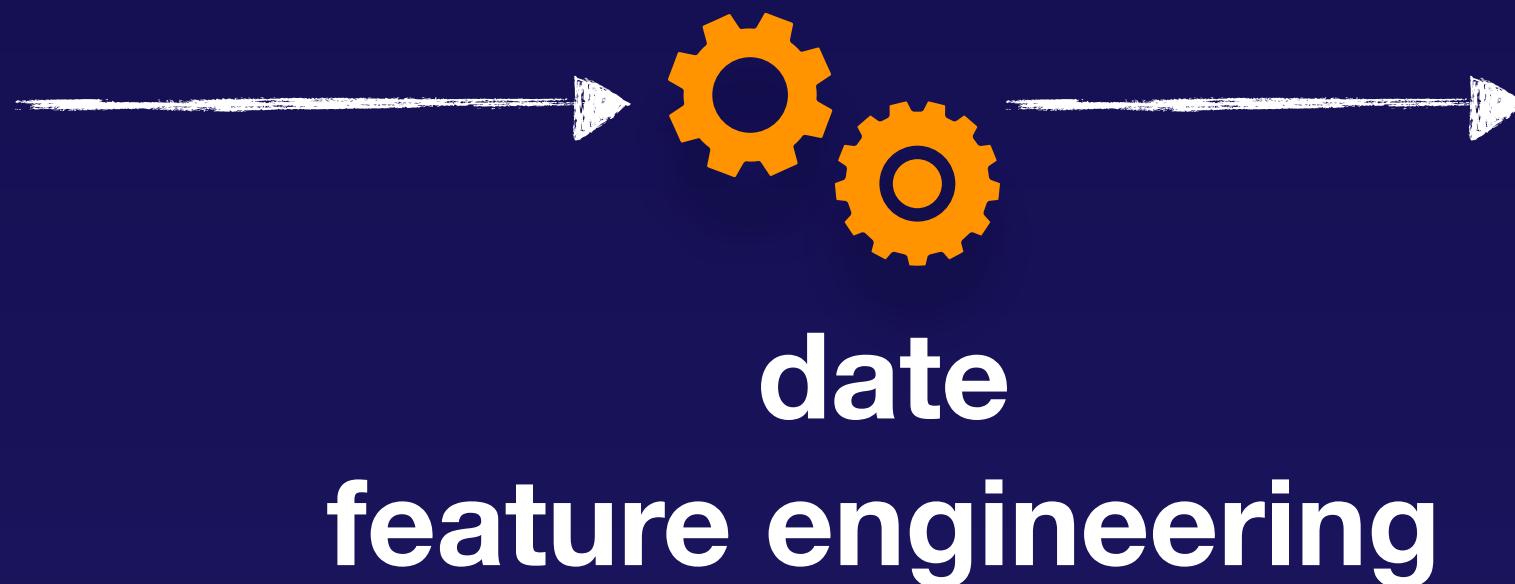
- 2013-02-08T09
- 6 Mar 17 21:22 UT
- Mon 06 Mar 2017 21:22:23 z
- 09/28/2019

EXTRACTED INFORMATION

- Was it a weekend? Was it a week day?
- Was the date the end of a quarter?
- What was the season?
- Was the day a holiday?
- Was it during business hours?
- Was the World Cup taking place on this date?

Feature Engineering Dates

Date
2015-06-17
2015-04-24
2015-02-12
2015-12-16
2015-03-14



is_weekend	day_of_week	month	year
0	2	6	2015
0	4	4	2015
0	3	2	2015
0	2	12	2015
1	5	3	2015

Text Feature Engineering Summary

Technique	Function
N-Gram	Splits text by whitespace with window size n
Orthogonal Sparse Bigram (OSB)	Splits text by keeping first word and uses delimiter with remaining whitespaces between second word with window size n
Term Frequency - Inverse Document Frequency (tf-idf)	Helps us determine how important words are within multiple documents
Removing Punctuation	Removes punctuation from text
Lowercase	Lowercase all the words in the text
Cartesian Product	Combines words together to create new feature
Feature Engineering Dates	Extracts information from dates and creates new features