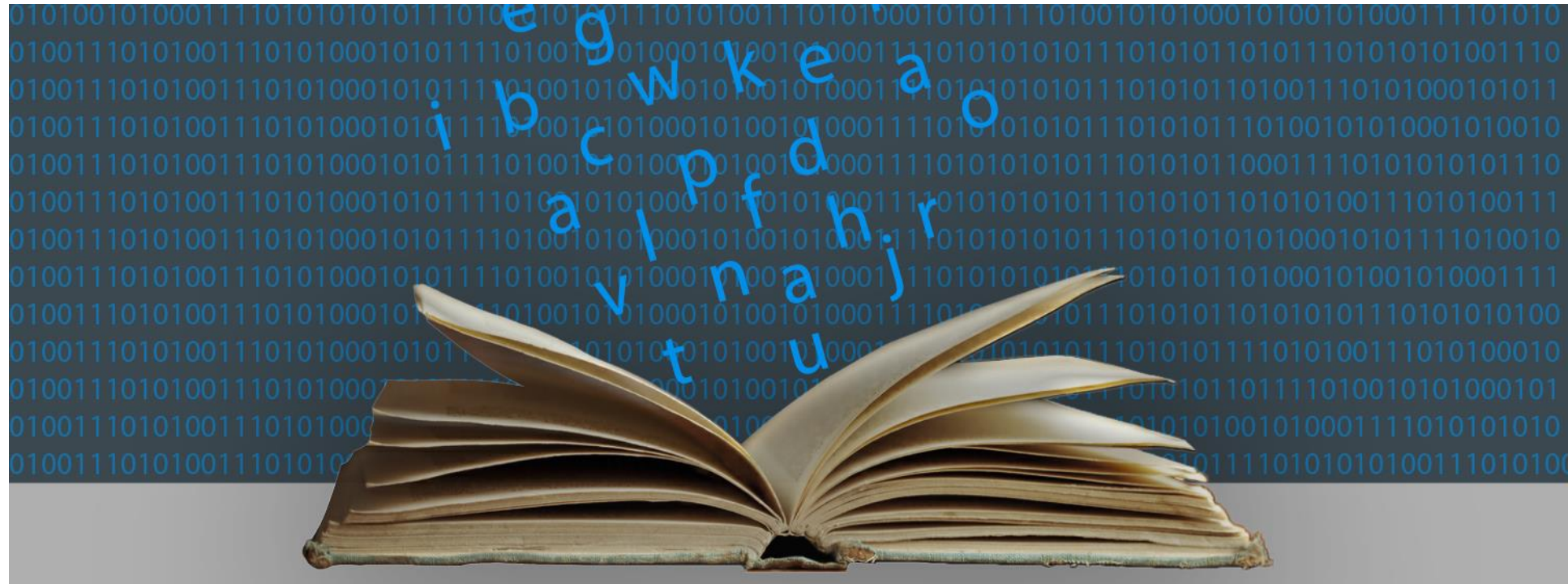


Examining a multi-layered approach for classification of OCR quality without Ground Truth



Mirjam Cuper, 17 March 2021

Introduction

- Mirjam Cuper
- Data scientist
- KB - National Library of the Netherlands
- Digital heritage

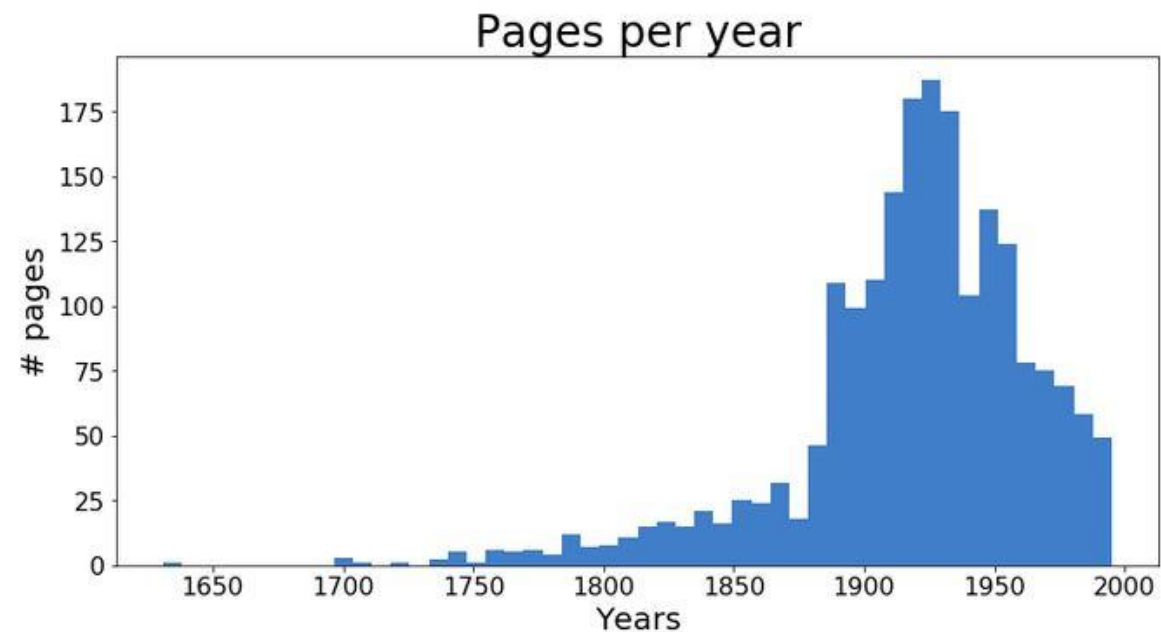
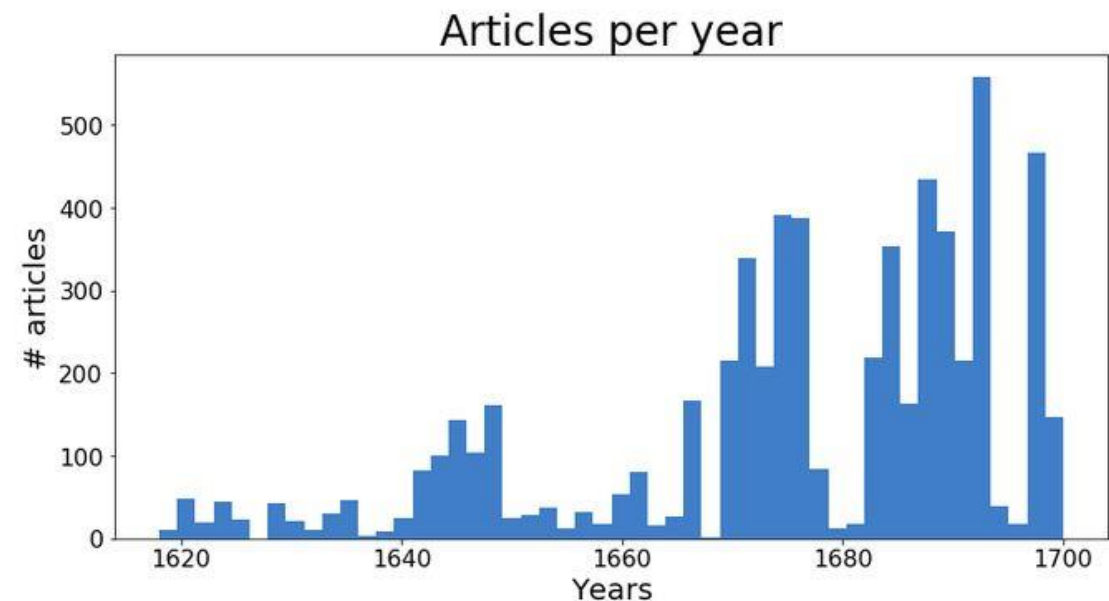


Overview

- Dataset
- Problem description
- Multi-layered approach
- First results
- Future work

Sample data

- Based on:
- +/- 30.000 news articles
- 2.000 newspaper pages
- Ground truth + OCR



WHAT IS THE PROBLEM?

OCR errors

WEST-INDIEN.
Suriname den 28 December 1682. Het Schip de Harderin, Schipper Dirck Barentsz. Kock, is alhier behouden aengekomen. Een dag a twee nae dese arriveerde hier oock Schipper Tobias Adriaenz., dewelcke aen St. Jago 21 Paerden ingenomen, hier levendigh overgevoert, en tot een hooge Prijs verkocht heeft. Noch een ander, zijnde een Engels Schipper, heeft mede, van Nieuw Engelandt komende, 20 Paerden hier gebracht. Met de Naturellen is het noch als vooren. Men heeft groot gebreck van Volck en Slaven. Anders staet hier alles wel.



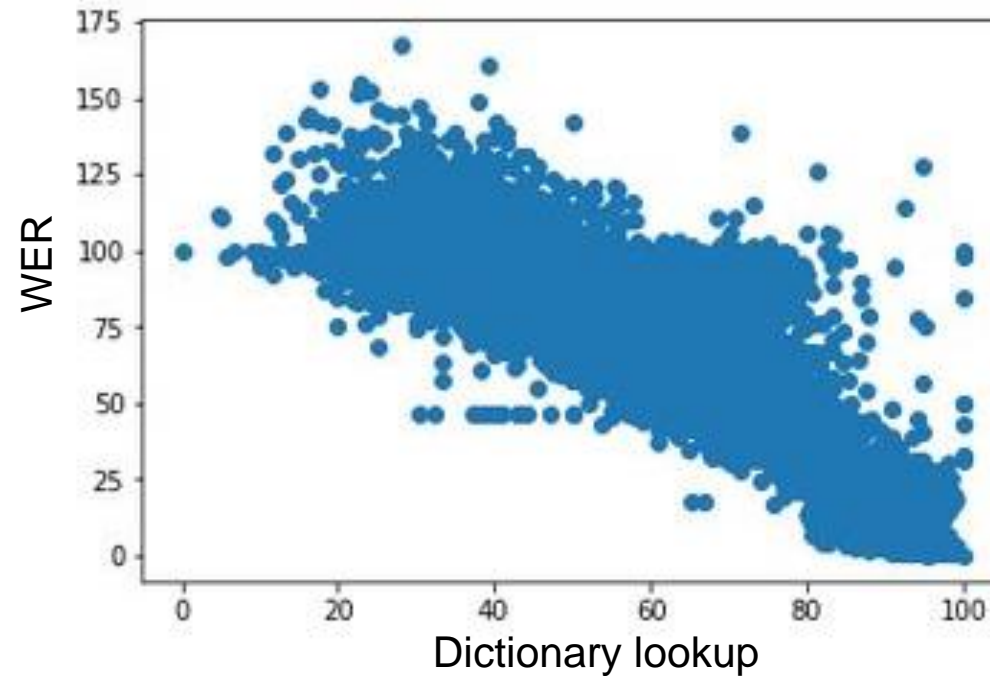
WEST-INDIEN.

SUrinimÃ den 28 December 1682. Het Schip de Harderin, Schipper Dirck Barentsz. Koek, is alhier behouden aengekomen. Een dag a twee nae defe arriveerde hier oock Schipper Tobias Adriaenz., dewelcke aen St. Jago2i Paerden ineenomenhier levendigh overgevoert, en tot eenhooge Prijs verkocht heeft. Noch een ander, zijnde een Engels Schipper, heeft mede, van Nieuw Engelandt komendeÂ» 20 Paerden hier gebracht. Met de Naturellen is het noch als vooren. Men heeft groot gebreck van Volck en Slaven. Anders ftaet hier alles wel.

Dictionary lookup - overview

- Correlation WER
- Relatively good result

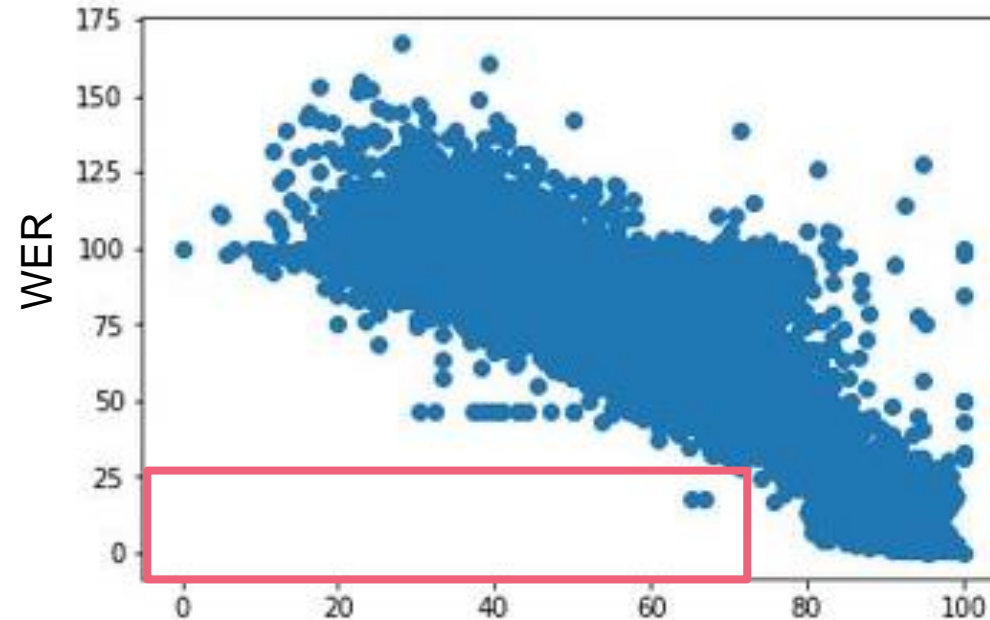
Correlation Dictionary lookup and WER



Dictionary lookup - overview

- Correlation WER
- Relatively good result

Correlation Dictionary lookup and WER



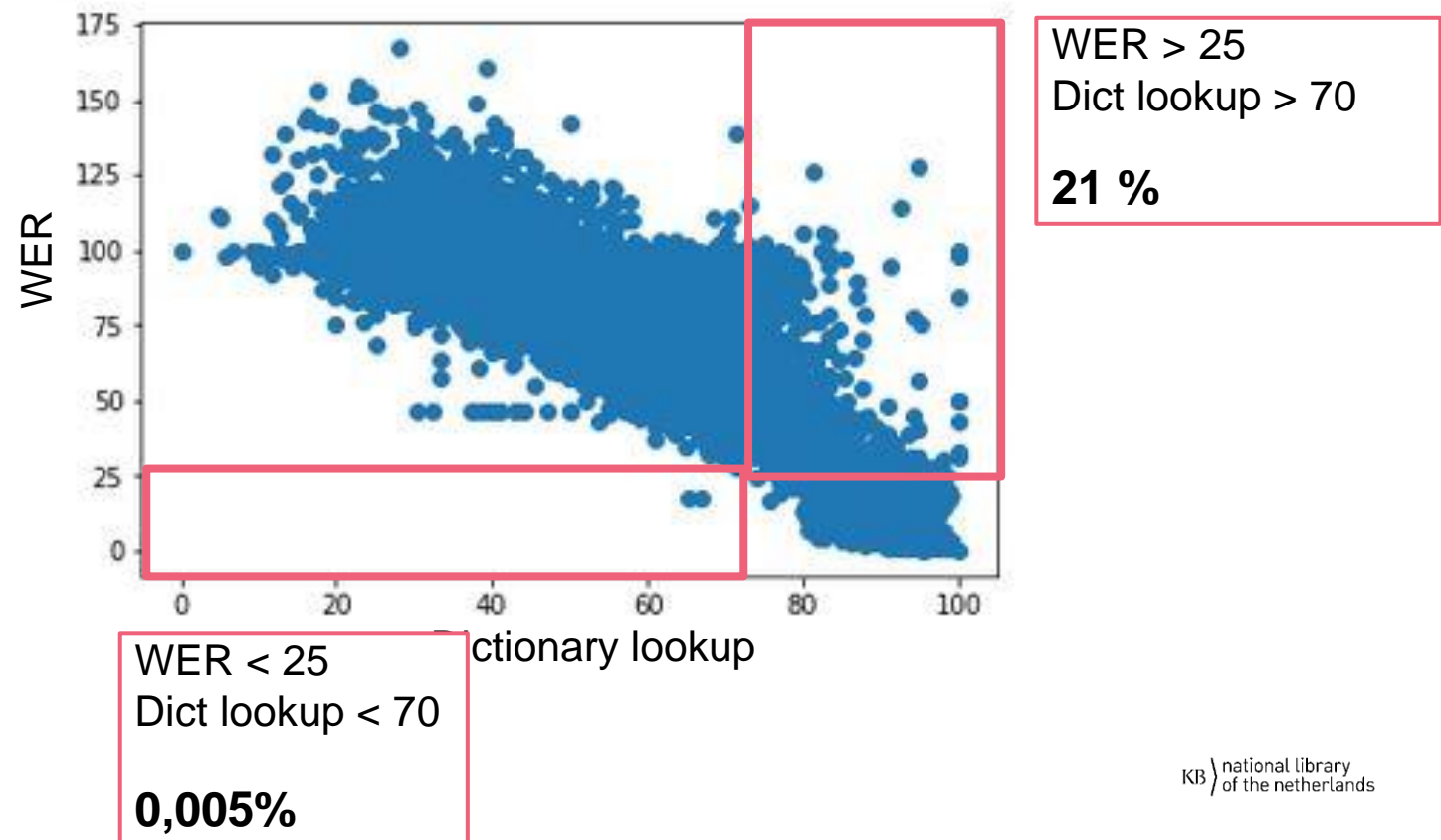
WER < 25
Dict lookup < 70

0,005%

Dictionary lookup - overview

- Correlation WER
- Relatively good result

Correlation Dictionary lookup and WER



Missing text

Wt Maestricht den 13 dito.
De Stadt Chiene / dooz dien sy in 't parlementeren
't Beschut onder de Fransche Troupen lossen / is bechtens
der-handt ingenomen / ende de goede Stadt aen bier plaets
sen / sonder te weten hoe aen vrandt gheraecht / ende een
Clooster / daer in veel Volcks geblycht was / dooz 't vper
hergaen / de blamme heeft men eenige upzen wijt ende vrees
ghesien.



Wt Maestricht den 13 dito.
g_<Â«^taÃ¶tiâ, -i)ieiiien/boo? bien
fptn'tpartementer

Missing text

Wt Maestricht den 13 dito.
De Stadt Chiene / dooz dien sy in 't parlementeren
't Beschut onder de Fransche Troupen lossen / is bechtens
der-handt ingenomen / ende de goede Stadt aen bier plaets
sen / sonder te weten hoe aen brandt
Clooster / daer in veel Dolcks geblu
hergaen / de blamme heefmen eenige
ghesien.

Wt Maestricht den 13 dito.
g_<Â«^taÃ¶tiâ,¬i)ieiiien/boo? bien
fptn'tpartementer

dictionary:
75% of words found

Word importance

dictionary:
76% of words found

Een handleiding of een gebriksaanwijzing is een schriftelijke instrvctie die met een product meegeleverd wordt.

De inhovd van een handlejding is:

- hoe een produot te assernbleren;
- hoe een product te reporeren;
- hoe een prodvct te instaileren;
- hoe een product gebruikt dient te worden;
- hoe een poduct te onderhovden;
- hoe de instellngen van een product aan te passen;
- hoe een storimg op te lossen;
- hoe een produci niet te misbruiken;
- hoe comtact op te nemen voor servioe.

Word importance

dictionary:
76% of words found

Een **handleiding** of een **gebruiksaanwijzing** is een **schriftelijke** instructie die met een **product** meegeleverd wordt.

De **inhoud** van een **handleiding** is:

hoe een **product** te **assembleren**;

hoe een product te **repareren**;

hoe een **product** te **installeren**;

hoe een product gebruikt dient te worden;

hoe een **product** te **onderhouden**;

hoe de **instellingen** van een product aan te passen;

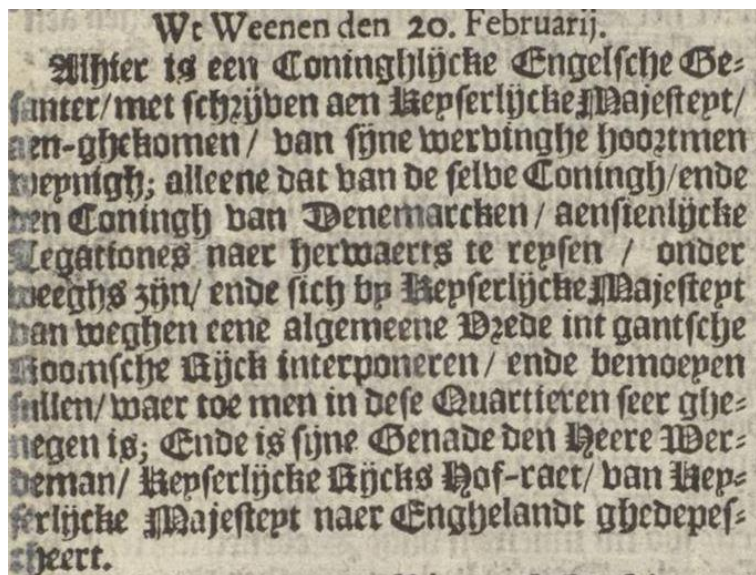
hoe een **storing** op te lossen;

hoe een **product** niet te misbruiken;

hoe **contact** op te nemen voor **service**.

Time period

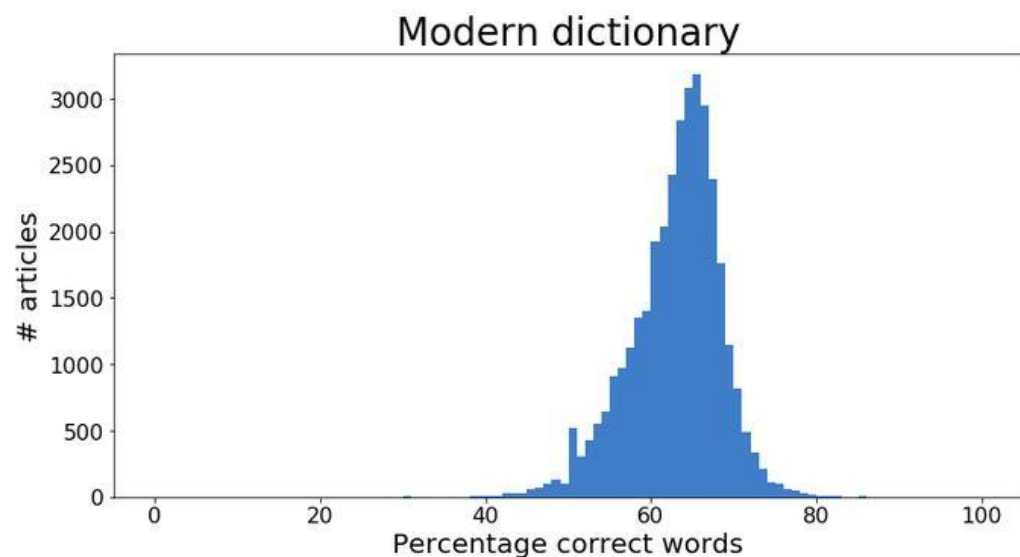
Article, 1636 – Ground Truth text



Wt Weenen den 20. Februarij. Alhier is een Coninghlijcke Engelsche Gesanter, met schrijven aen Keyserlijcke Majestejt, aen-ghekomen, van sijne werwinghe hoortmen weynigh; alleene dat van de selve Coningh, ende den Coningh van Denemarcken, aensienlijcke Legationes naer herwaerts te reysen, onder weeghs zijn, ende sich by Keyserlijcke Majestejt van weghen eene algemeene Vrede int gantsche Roomsche Rijk interponeren, ende bemoeyen sullen, waer toe men in dese Quartieren seer ghenegen is; Ende is sijne Genade den Heere Werdeman, Keyserlijcke Rijcks Hof-raet, van Keyserlijcke Majestejt naer Engelandt ghedepescheert. De werwingen gaen alhier noch sterck voort, ende wert veel Meel ende Coorn, tot provianderinghe vande Keyserlijcke Armaden, in 't Rijk ghesonden.

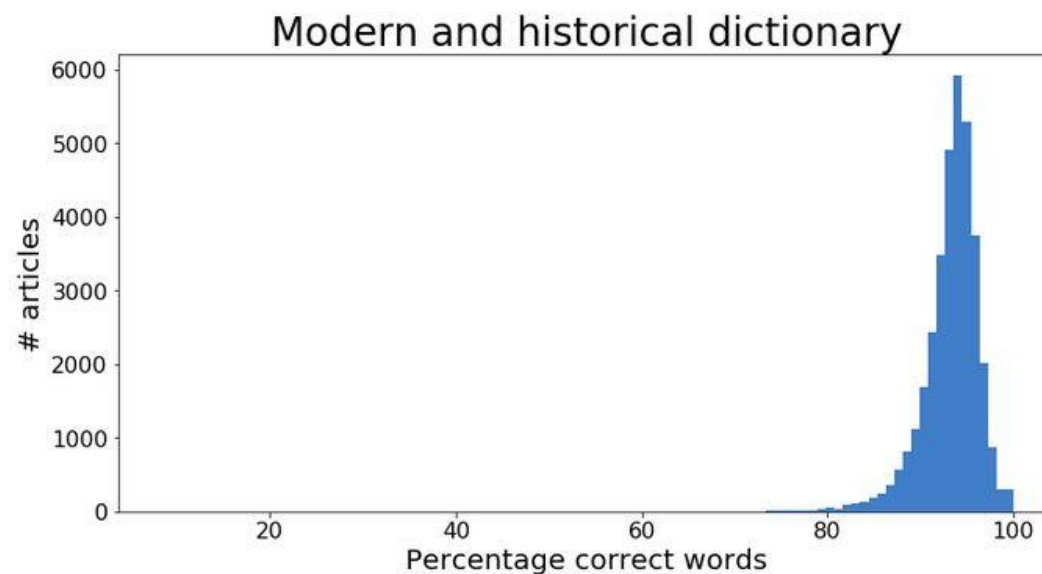
Time period

17th century articles: ground truth data



Mean: 63%

Median: 64%



Mean: 93%

Median: 94%

Language

JOURNAL DE LA PROVINCE DE LIMBOURG. (N.° 5.)

Samedi, 6 Janvier 1821.

PORTUGAL Lisbonne, le 13 décembre. Les élections des députés se poursuivent avec activité dans cette capitale et dans les provinces; elles se font en tout conformément aux ordres qui ont été

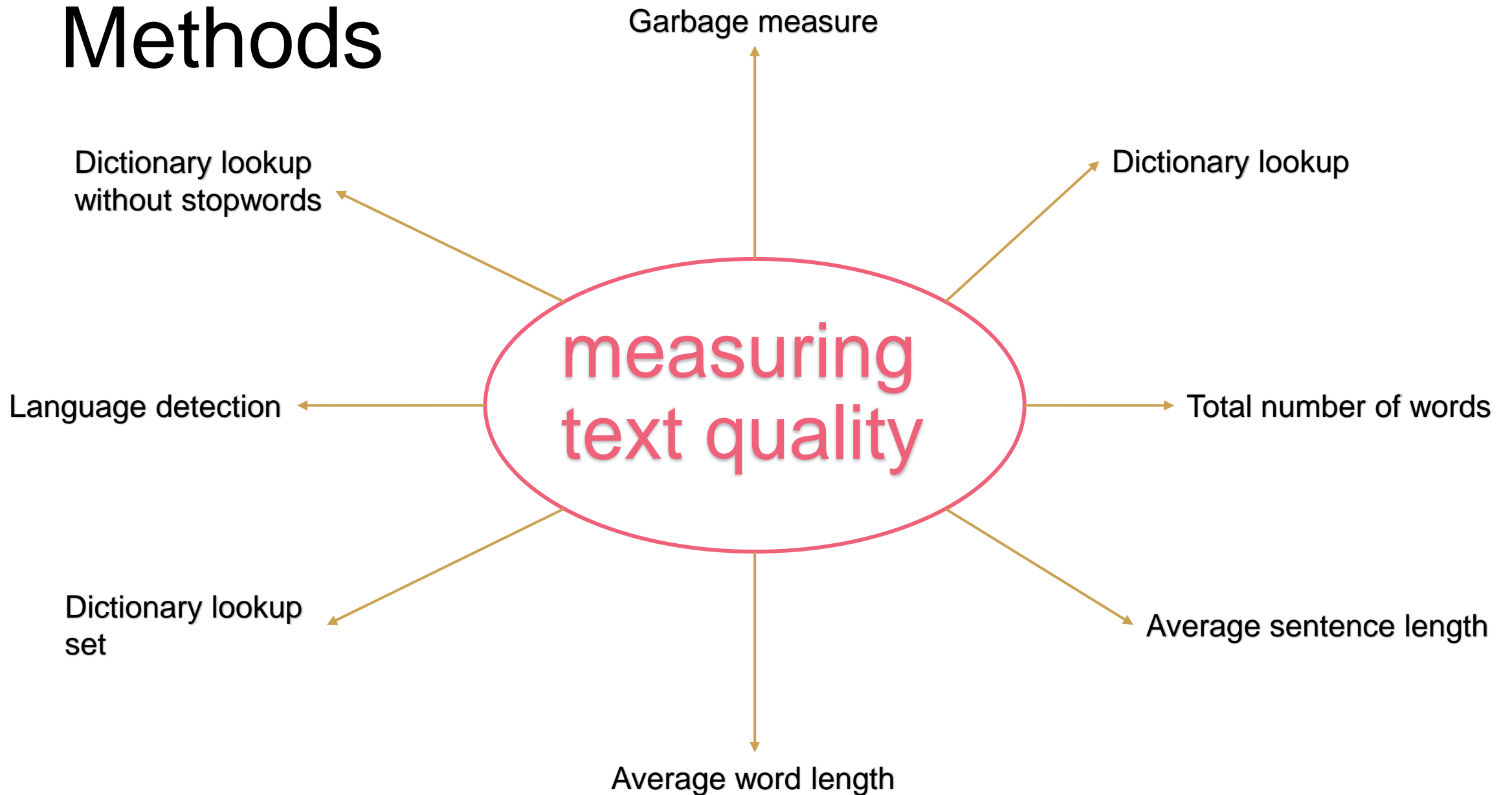
Combined dictionary:

86% of words found

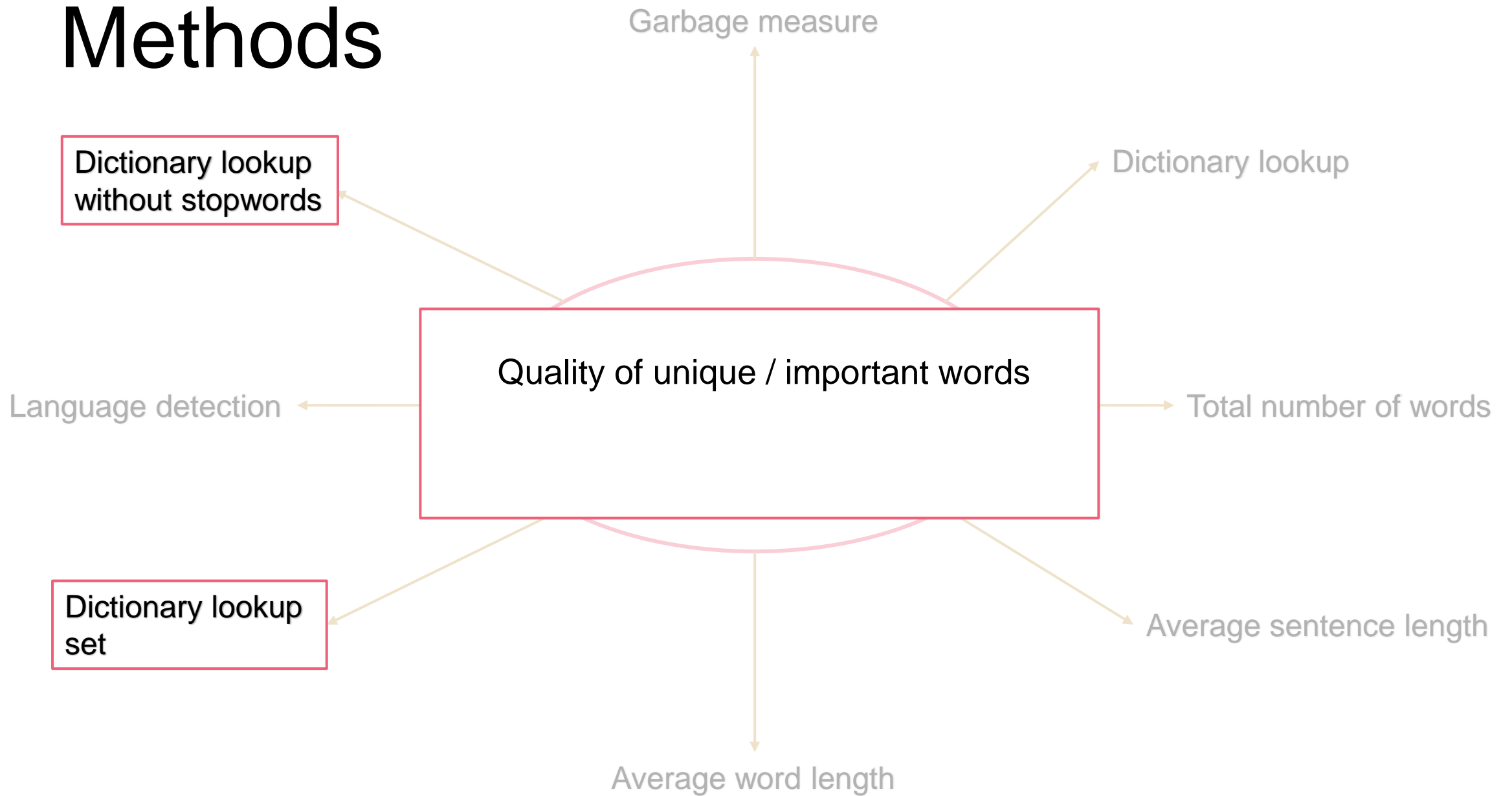
és à ce sujet. Cependant, dans l'Alentejo on a découvert voudrait réunie des cortès à Lamego, mais afin d'étouffer incipie cette division, on a envoyé de cette capitale dans la vince quelques régimens d'infanterie et de cavalerie. Nos ec le cabinet de Madrid deviennent chaque jour plus et plus amicales. Avant hier, nous avons vu arriver une estarette espagnole venant de Madrid, avec des dépêches pour notre gouvernement. Depuis quelques jours, un grand nombre d'Espagnols affluent à Lisbonne. Le bruit court que des mouvemens sérieux ont eu lieu à Rio-Janeiro aussitôt qu'on y a appris la nouvelle des événemens du Portugal. On dit que les troupes ont fait dans cette capitale la même chose qu'à Oporto; mais jusqu'à présent, on ne connaît rien d'officiel à cet égard. ESPAGNE. Madrid, le 21 décembre. On dit que le nombre des Espagnols qui ont ordre de quitter la capitale est de 54, mais on ne connaît pas encore leurs noms. On

MEASURING METHODS

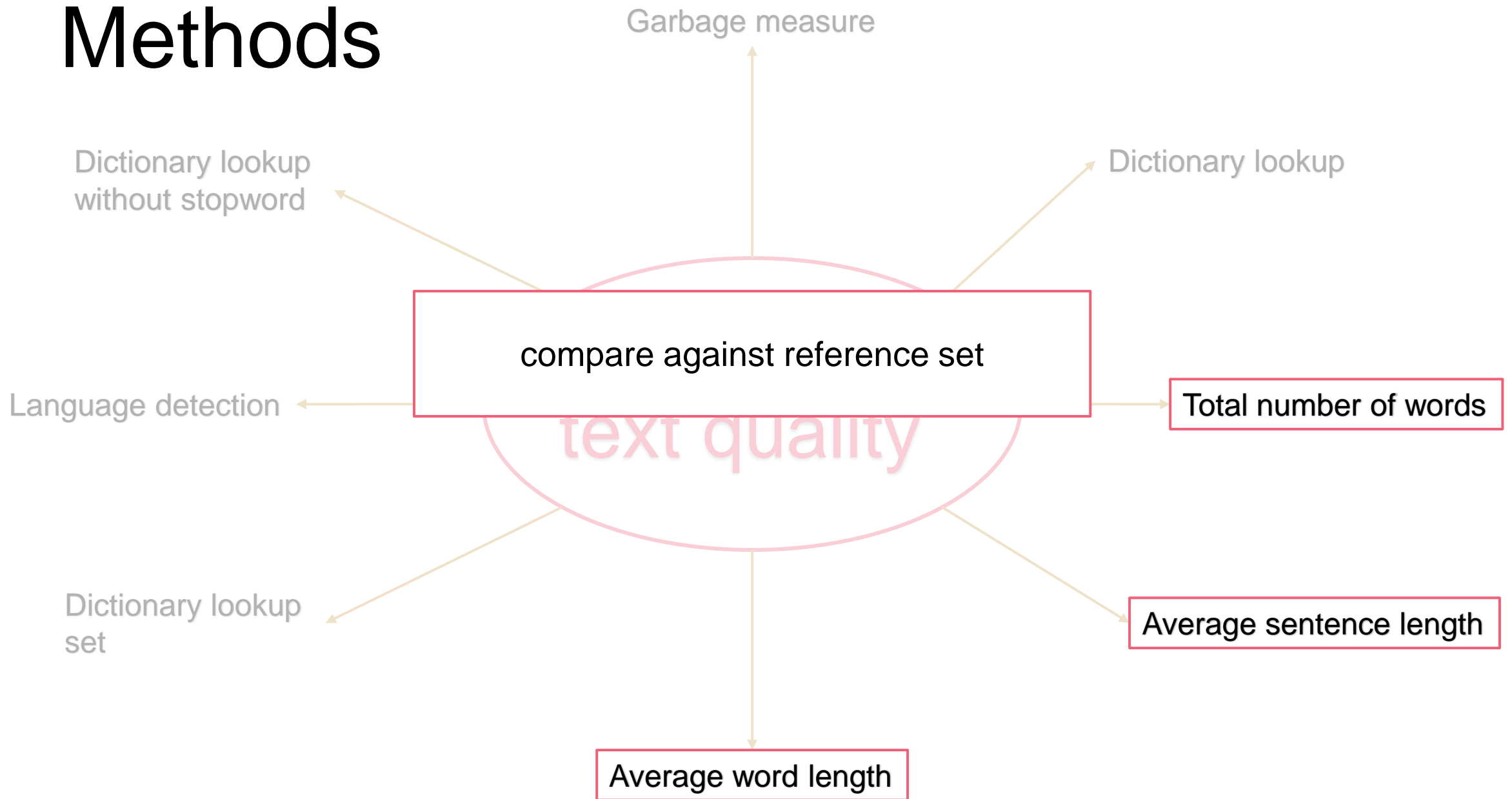
Methods



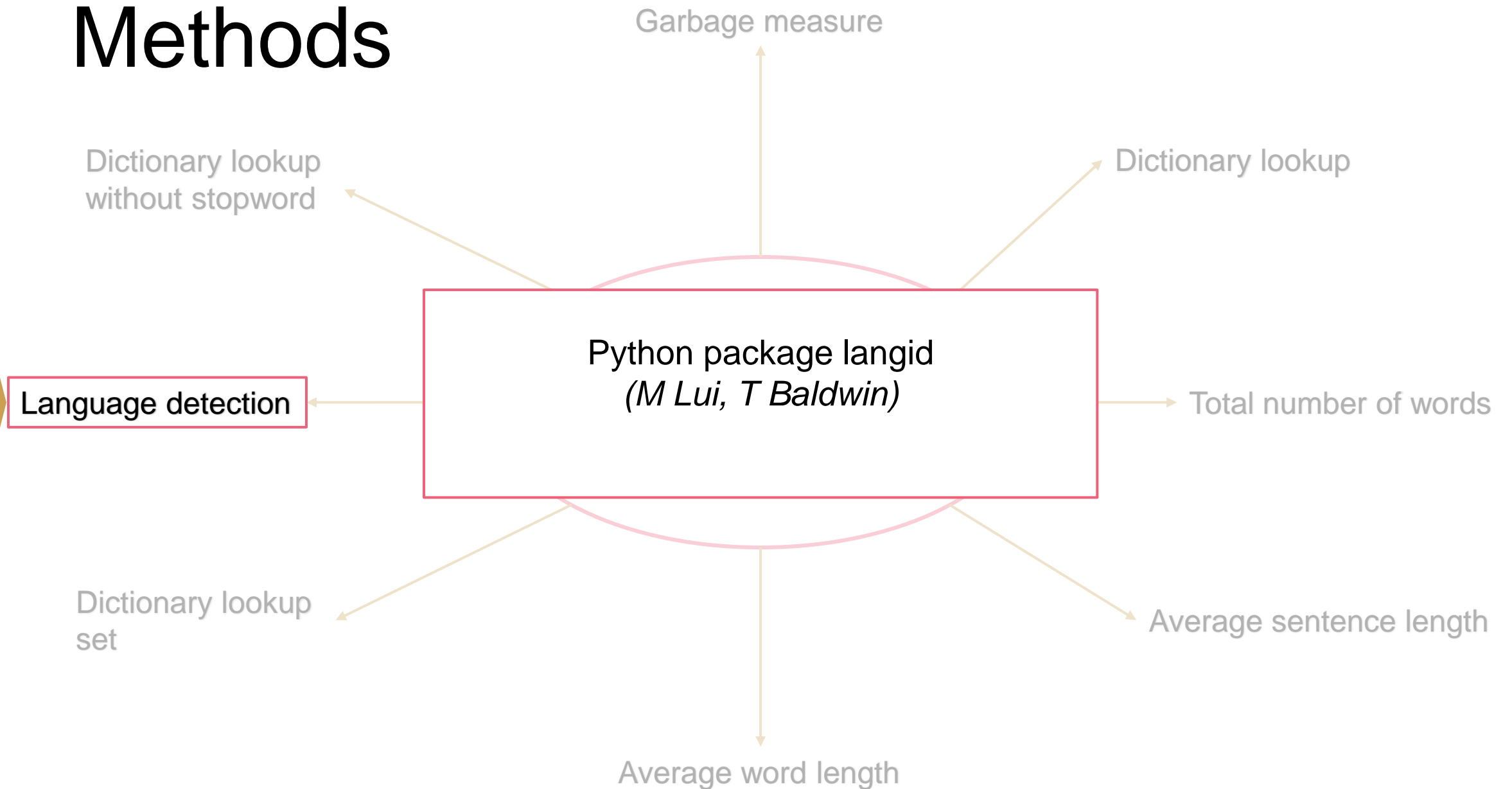
Methods



Methods



Methods



Methods

Garbage measure

Dictionary lookup
without stopword

Dictionary lookup

Detect 'garbage' strings
(*Taghva, K. et al., 2001*)

Language detection

Total number of words

Dictionary lookup
set

Average sentence length

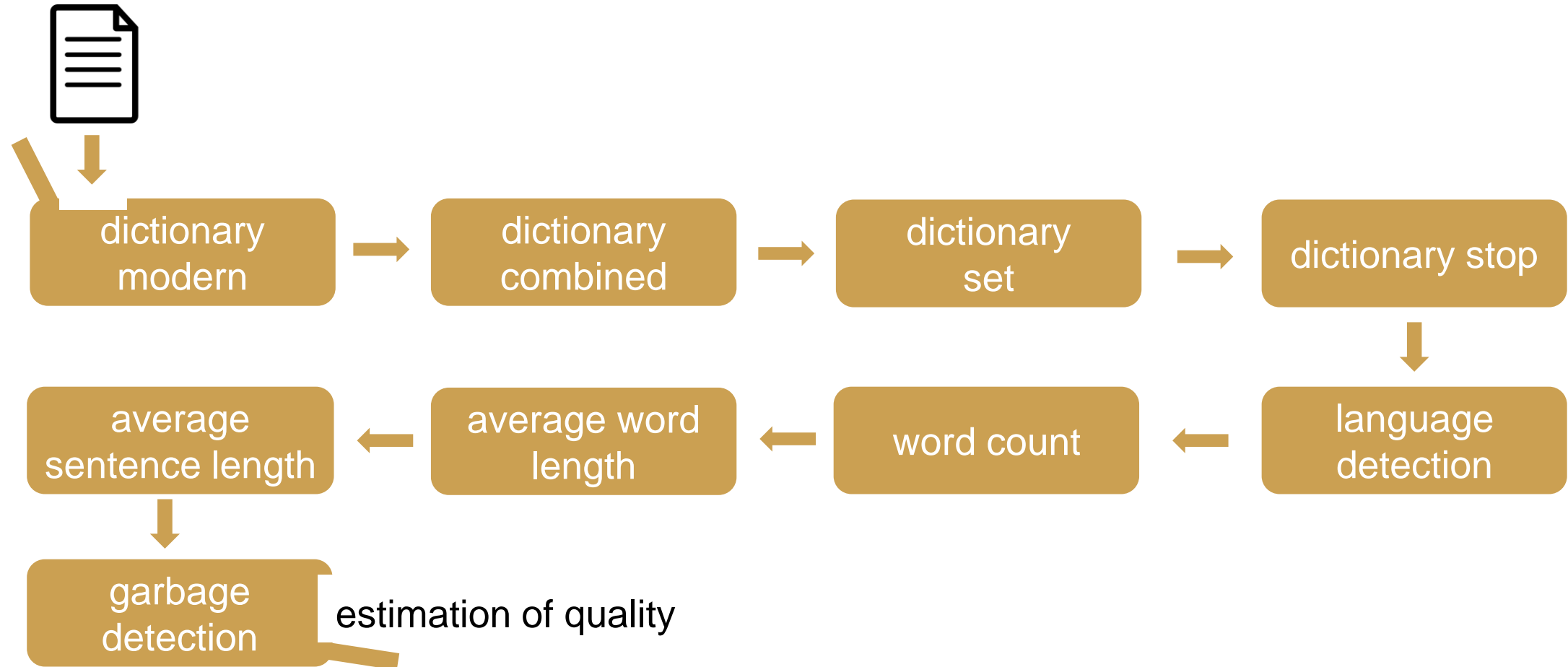
Average word length

Characteristics

1. Type of text (advertisement, article, page)
2. Published year
3. Language

MULTI-LAYERED APPROACH

Pipeline

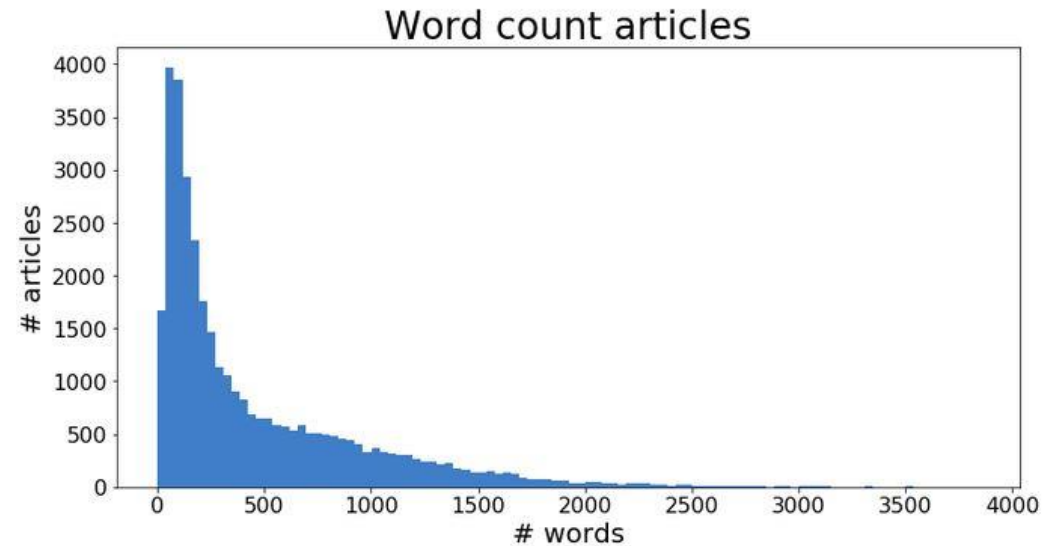


Comparison measures

17th century news articles Ground Truth

80% of articles:

- Word count:
 - 56 – 1182 words
- Word length:
 - Mean: 4 – 5 characters
 - Median: 4 – 5 characters
- Sentence length:
 - Mean: 12-36 words
 - Median: 5-26 words

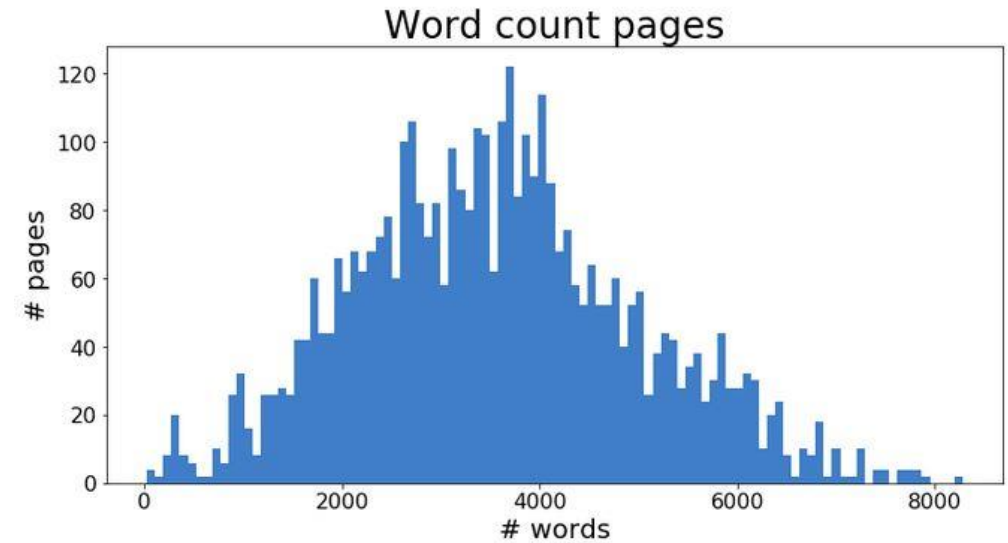


Comparison measures

Random newspaper pages Ground Truth

80% of articles:

- Word count:
 - 1760 – 5558 words
- Word length:
 - Mean: 4.9 – 5.4 characters
 - Median: 4 characters
- Sentence length:
 - Mean: 7-19 words
 - Median: 3-16 words



Pipeline – example outcomes

Dict modern	Dict combined	Dict set	Dict stop	Word count	Language	Median word length	Mean word length	Median Sentence length	Mean sentence length	Garbage
17 th century news article										
60%	80%	75%	65%	2000	nl (1.0)	4	4.5	9	15	[0.25, 0, 0, 0.25, 0.12, 0]
20 th century newspaper page										
60%	80%	75%	65%	2000	nl (1.0)	4	4.5	9	15	[0.25, 0, 0, 0.25, 0.12, 0]

 = good

 = bad

Pipeline – example outcomes

Dict modern	Dict combined	Dict set	Dict stop	Word count	Language	Median word length	Mean word length	Median Sentence length	Mean sentence length	Garbage
17 th century news article										
60%	80%	75%	65%	2000	nl (1.0)	4	4.5	9	15	[0.25, 0, 0, 0.25, 0.12, 0]
20 th century newspaper page										
60%	80%	75%	65%	2000	nl (1.0)	4	4.5	9	15	[0.25, 0, 0, 0.25, 0.12, 0]

Cut-off point: 70%

Pipeline – example outcomes

Dict modern	Dict combined	Dict set	Dict stop	Word count	Language	Median word length	Mean word length	Median Sentence length	Mean sentence length	Garbage
17 th century news article										
60%	80%	75%	65%	2000	nl (1.0)	4	4.5	9	15	[0.25, 0, 0, 0.25, 0.12, 0]
20 th century newspaper page										
60%	80%	75%	65%	2000	nl (1.0)	4	4.5	9	15	[0.25, 0, 0, 0.25, 0.12, 0]

Based on reference set

Pipeline – example outcomes

Dict modern	Dict combined	Dict set	Dict stop	Word count	Language	Median word length	Mean word length	Median Sentence length	Mean sentence length	Garbage
17 th century news article										
60%	80%	75%	65%	2000	nl (1.0)	4	4.5	9	15	[0.25, 0, 0, 0.25, 0.12, 0]
20 th century newspaper page										
60%	80%	75%	65%	2000	nl (1.0)	4	4.5	9	15	[0.25, 0, 0, 0.25, 0.12, 0]

Dutch

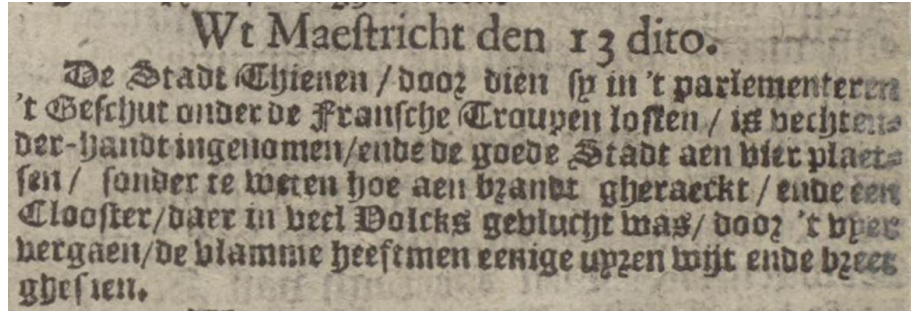
Pipeline – example outcomes

Dict modern	Dict combined	Dict set	Dict stop	Word count	Language	Median word length	Mean word length	Median Sentence length	Mean sentence length	Garbage
17 th century news article										
60%	80%	75%	65%	2000	nl (1.0)	4	4.5	9	15	[0.25, 0, 0, 0.25, 0.12, 0]
20 th century newspaper page										
60%	80%	75%	65%	2000	nl (1.0)	4	4.5	9	15	[0.25, 0, 0, 0.25, 0.12, 0]

Each value < 0.25

FIRST RESULTS

Missing text



Wt Maestricht den 13 dito.
 g_<^ta^tiâ, -i)ieiiien/boo? bien
 fptn'tpartementer

Dict modern	Dict combined	Dict set	Dict stop	Word count	Language	Median word length	Mean word length	Median Sentence length	Mean sentence length	Garbage
50%	75%	75%	66,67%	8	Lb (0.9)	4.5	8.88	4.5	4.5	[0.25, 0, 0, 0.25, 0.12, 0]

Remarkable:

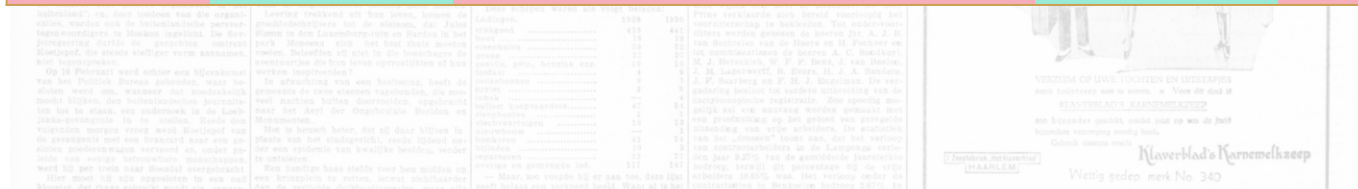
- Language Luxembourg's
- Very low word count
- Garbage detected

Dictionary lookup – missing text



KOETJEPOF UIT MOSKOU VERWIJDERD? Hij zou te Soesdal gevangen worden gehouden. Nieuwe lezing over zijn ontvoering. — Met de „Spartak” van Trouville via Antwerpen naar Rusland. - Thans in een oud klooster. I ■ kou i lse inke.!

Dict modern	Dict combined	Dict set	Dict stop	Word count	Language	Median word length	Mean word length	Median Sentence length	Mean sentence length	Garbage
62.4%	70.8%	69%	48.8%	867	nl (1.0)	3	3.9	7	8.4	[0.33, 0, 0, 0.33, 0.02, 0.37]



MEERDERHEID IN TRANSVAAL. Als resultaat der provinciale verkiezingen. JOH • zingen voor den provincialen raad van 'l • • kaaneche partij en Z l.ab- Twee uitslagen zijn nog niet l.- i zul

Remarkable:
 - Low word count for page – less than half of the minimum value
 - High garbage measure



WEINIG OPLEVING IN DE HAVEN. Ook hier het algemeene e malaise, maakt, b. v. met het ■ n ■ hepen ida -ttge ' alles ii" een n . ral nu •rs aar let np i •n de -treatjittl U den A ' . i de Poll 19:» itso tti 27 II 4 8 | i S is 61 1 13. — 1 14 3 al ls be; - r_an, de toestand In «een of « , e hotu en In aan i ;. het oogenblik het beeld

OOST-INDIË. ONTSLAGKWESTIE BIJ DE BATAAFSCHE. ■ - . H. i. M. raadt -chf C en, of i ZUID-SUMATRA LANDBOUW- EN NIJVER HEIDSVEREENIGING. ■ J. ! M. J. 1 i i ■ WEEROVERZICHT. ' A ' ■ ■ LiUkt'i. tn 7. 55 gr. Woedende

Word importance

Een **handleiding** of een **gebruiksaanwijzing** is een **schriftelijke instructie** die met een **product** meegeleverd wordt.

De **inhoud** van een **handleiding** is:

hoe een **product** te **assembleren**;

hoe een product te **repareren**;

hoe een **product** te **installeren**;

hoe een product gebruikt dient te worden;

hoe een **product** te **onderhouden**;

hoe de **instellingen** van een product aan te passen;

hoe een **storing** op te lossen;

hoe een **product** niet te misbruiken;

hoe **contact** op te nemen voor **service**.

Word importance

Een handleiding of een gebruiksaanwijzing is een schriftelijke instructie die met een product meegeleverd wordt.

Dict modern	Dict combined	Dict set	Dict stop	Word count	Language	Median word length	Mean word length	Median Sentence length	Mean sentence length	Garbage
73.3%	76%	58.1%	40%	76	nl (1.0)	3	5.43	15	25.67	[0.01, 0, 0, 0, 0, 0.01]

hoe een product te installeren;
hoe een product gebruikt dient te worden;
hoe een product te onderhouden;
hoe de instellingen van een product aan te passen;
hoe een storing op te lossen;

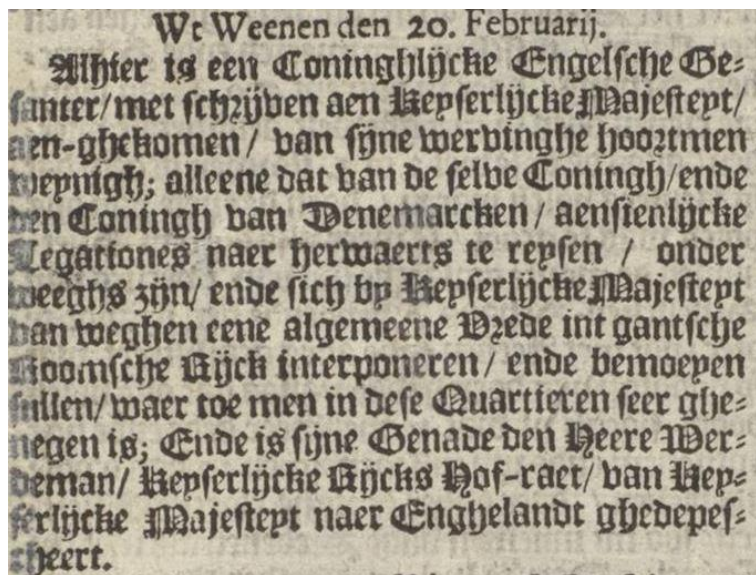
Not every measure can be evaluated

Remarkable:

- Dict set and dict stop considerably lower

Time period

Article, 1636



Wt Weenen den 20. Februarij. Alhier is een Coninghlijcke Engelsche Gesanter, met schrijven aen Keyserlijcke Majestejt, aen-ghekomen, van sijne werwinghe hoortmen weynigh; alleene dat van de selve Coningh, ende den Coningh van Denemarcken, aensienlijcke Legationes naer herwaerts te reysen, onder weeghs zijn, ende sich by Keyserlijcke Majestejt van weghen eene algemeene Vrede int gantsche Roomsche Rijk interponeren, ende bemoeyen sullen, waer toe men in dese Quartieren seer ghenegen is; Ende is sijne Genade den Heere Werdeman, Keyserlijcke Rijcks Hof-raet, van Keyserlijcke Majestejt naer Engelandt ghedepescheert. De wervingen gaen alhier noch sterck voort, ende wert veel Meel ende Coorn, tot provianderinghe vande Keyserlijcke Armaden, in 't Rijk ghesonden.

Dictionary lookup – timeperiod

Article, 1636

Dict modern	Dict combined	Dict set	Dict stop	Word count	Language	Median word length	Mean word length	Median Sentence length	Mean sentence length	Garbage
41%	93.4%	91.7%	89.3%	107	nl (1.0)	5	6.13	4	22.4	[0.04, 0, 0, 0.03, 0, 0.09]

Remarkable:
- Combined dict significantly higher!

Dictionary lookup – languages

JOURNAL DE LA PROVINCE DE LIMBOURG. (N.° 5.)

Samedi, 6 Janvier 1821.

PORTUGAL Lisbonne, le 13 décembre. Les élections des députés se poursuivent avec activité dans cette capitale et dans les provinces; elles se font en tout conformément aux ordres qui ont été

Combined dictionary:

86% of words found

és à ce sujet. Cependant, dans l'Alentejo on a découvert voudrait réunie des cortès à Lamego, mais afin d'étouffer incipie cette division, on a envoyé de cette capitale dans la vince quelques régimens d'infanterie et de cavalerie. Nos ec le cabinet de Madrid deviennent chaque jour plus et plus amicales. Avant hier, nous avons vu arriver une estarette espagnole venant de Madrid, avec des dépêches pour notre gouvernement. Depuis quelques jours, un grand nombre d'Espagnols affluent à Lisbonne. Le bruit court que des mouvemens sérieux ont eu lieu à Rio-Janeiro aussitôt qu'on y a appris la nouvelle des événemens du Portugal. On dit que les troupes ont fait dans cette capitale la même chose qu'à Oporto; mais jusqu'à présent, on ne connaît rien d'officiel à cet égard. ESPAGNE. Madrid, le 21 décembre. On dit que le nombre des Espagnols qui ont ordre de quitter la capitale est de 54, mais on ne connaît pas encore leurs noms. On

Dictionary lookup – languages

JOURNAL DE LA PROVINCE DE LIMBOURG. (N.° 5.)
 Samedi, 6 Janvier 1821.
 PORTUGAL Lisbonne, le 13 décembre. Les élections des députés se poursuivent avec activité dans cette capitale et dans les provinces; elles se font en tout conformément aux ordres qui ont été
 émis à ce sujet. Cependant, dans l'Alentejo on a découvert

Dict modern	Dict combined	Dict set	Dict stop	Word count	Language	Median word length	Mean word length	Median Sentence length	Mean sentence length	Garbage
42.7%	82.8%	71,8%	76.4%	803	fr (1.0)	4	4.74	16	18.06	[0.11, 0, 0, 0.1, 0, 0.09]

Remarkable:
 - Predicted language

gouvernement. Depuis quelques jours, un grand nombre d'Espagnols affluent à Lisbonne. Le bruit court que des mouvements l'on y a appris la nouvelle troupes ont fait dans cette capitale la même chose qu'à Oporto; mais jusqu'à présent, on ne connaît rien d'officiel à cet égard. ESPAGNE. Madrid, le 21 décembre. On dit que le nombre des Espagnols qui ont ordre de quitter la capitale est de 54, mais on ne connaît pas encore leurs noms. On

FUTURE WORK

More testing

- More tests and more data
- Prediction of correct texts
- Comparison measured quality with ground truth

Quantifying results

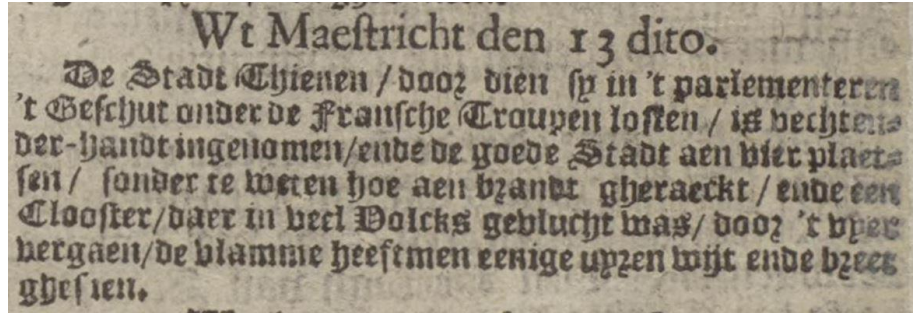
From table to one final score:

- 'good', 'average', 'bad'

Weights

- Importance
- Characteristics

Weights - example



Wt Maestricht den 13 dito.
 g_<^«^taÃ¶tiâ,¬i)ieiiien/boo? bien
 fptn'tpartementer

Dict modern	Dict combined	Dict set	Dict stop	Word count	Language	Median word length	Mean word length	Median Sentence length	Mean sentence length	Garbage
50%	75%	75%	66,67%	8	Lb (0.9)	4.5	8.88	4.5	4.5	[0.25, 0, 0, 0.25, 0.12, 0]

Word count and language more important?

Adding more measures

- N-grams
- Part of speech
- Letter ratio
- WER (without ground truth)


Ideas?

mirjam.cuper@kb.nl

@CuperMirjam

Questions?





KB } national library
of the netherlands