



HAL
open science

Recovering from accidental syncretisms with word2vec analogies

Camille Lavigne, Gilles Boyé, Gauvain Schalchli

► **To cite this version:**

Camille Lavigne, Gilles Boyé, Gauvain Schalchli. Recovering from accidental syncretisms with word2vec analogies. 20th International Morphology Meeting, Sep 2022, Budapest, Hungary. halshs-04872862

HAL Id: halshs-04872862

<https://shs.hal.science/halshs-04872862v1>

Submitted on 8 Jan 2025

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Recovering from accidental syncretisms with word2vec analogies

Recovering from accidental syncretisms with word2vec analogies

Camille Lavigne, Gilles Boyé & Gauvain Schalchli



Université Bordeaux-Montaigne



COGNITION, LANGUES, LANGAGE, ERGONOMIE

CLLE-Montaigne

Systematic idiosyncratic syncretism

- It has received a lot of attention
 - Aronoff (1994) introduced morphemes to capture it
 - Maiden (1992) questioned its stability
 - Baerman et al. (2005) discuss its possible orientations
 - Bonami & Boyé (2003) studied its structure
- And with the emergence of the Paradigm Cell Filling Problem, systematic syncretism and the associated morphemes have become a central part of the discussion

The Paradigm Cell Filling Problem

(1) Some formulations:

- a. *What licenses reliable inferences about inflected (and derived) surface forms of a lexical item?* (Ackerman et al. 2009)
- b. *Given exposure to an inflected word form of a novel lexeme, what licenses reliable inferences about the other wordforms in its inflectional family?* (Ackerman & Malouf 2013)

Some overlooked problems

- The Paradigm Cell **Finding** Problem (Boyé & Schalchli 2019):
 - (2) *Given exposure to a sample of inflected forms with enough lexemes, what licences reliable inferences about the paradigm shape of these lexemes?*
- The scope of the Paradigm Cell **Filling** Problem:
 - should be based on the knowledge available to the speaker
 - should not be limited to novel lexemes
 - (1b) *Given exposure to inflected wordforms of a lexeme, what licenses reliable inferences about the other wordforms in its inflectional family?*

Accidental syncretisms

The Paradigm Cell Association Problem:

- *Assuming that speakers have identified a morphomic paradigm shape based on contrasting forms, how can they place (syncretic) forms in the morphomic paradigm?*
- Hypothesis:
 - High frequency lexemes are the key to finding the paradigm cells
 - they have the most contrasting forms => larger set of forms
 - they appear in the most varied contexts => well-defined distributions
- ➔ Speakers can use them to classify other forms and recover from the accidental syncretisms

A small experiment

English preterit & past participle

- A large proportion of English verbs have the same form for their preterit and their past participle
 - all regular verbs: walked, opened, enjoyed, ...
 - many irregulars: thought, made, kept, ...
- But among the most frequent verbs, it is common to have two different forms
 - ate-eaten, blew-blown, did-done, ...
- We can use word-vectors to identify the morphomic cells occupied by a form in a corpus by comparing its word vector to our base of frequent contrastive items

Word2Vec

Relevance

- contextual embeddings
- no differentiation between homographic forms
- simple, efficient, available

Parameters

- skipgram negative sampling
- 152 dimensions
- 7-word window
- 20 tokens minimum
- 10 epochs
- negative sampling 15

Two corpora

- Brown
 - small corpus (1M words)
 - tagged
- BNC
 - large corpus (100M words)
 - divided in manageable parts
 - 100k to 15M words
 - tagged and untagged versions

Looking for a single shift vector

using the average shift from one set of verbs

Simple hypothesis for the tagged corpora: one shift vector fits all

1. we calculated the average $\vec{V}_{\text{PPART}} - \vec{V}_{\text{PRET}}$ for a set of verbs (\vec{V}_{SHIFT})
2. we added this average vector \vec{V}_{SHIFT} to all the \vec{V}_{PRET} available
3. cosine similarity results never converged uniformly
 - on the contrary, with larger sets of verbs, the cosine similarity was really bad

Looking for shift vectors

using multiple verbs as potential models

Second hypothesis for the tagged corpora: use a set of shift vectors

1. we calculated the $\vec{V}_{\text{PPART}} - \vec{V}_{\text{PRET}}$ for a set of verbs (\vec{V}_{SHIFT_V})
2. we added all the shift vectors \vec{V}_{SHIFT} to all the \vec{V}_{PRET} available
3. every addition gives a candidate vector for the expected form
 - ➔ in the resulting candidates, there were almost always one in close proximity to the target forms
 - not always the same using the same shift vector
 - many times several shift vectors yielded good results

Shift vectors and their cosine similarities

using multiple verbs as potential models

Cosine similarities between generated vector and true vector

```
[('usedpp', 0.4871392846107483), ('providedpp', 0.4284799098968506), ('recommendedpp', 0.4055340588092804),  
[('usedpp', 0.4552362561225891), ('soughtpp', 0.4343701899051666), ('foundpp', 0.41937941312789917),  
[('introducedpp', 0.43687403202056885), ('succeededpp', 0.41987723112106323), ('strippedpp', 0.4146502017974853),  
[('discoveredpp', 0.437302827835083), ('feltpp', 0.431774765253067), ('recommendedpp', 0.4293583035469055),  
[('usedpp', 0.5839736461639404), ('consideredpp', 0.5338331460952759), ('indicatedpp', 0.5320671200752258)  
[('usedpp', 0.5534017086029053), ('employedpp', 0.5186872482299805), ('indicatedpp', 0.5084676146507263),  
[('evokedpp', 0.4569481909275055), ('producedpp', 0.42817384004592896), ('seenpp', 0.4021519124507904),  
  
[('usedp', 0.46688348054885864), ('uses', 0.4503510594367981), ('builtp', 0.4356897473335266)  
[('foundp', 0.42088770866394043), ('usedp', 0.4088103771209717), ('uses', 0.3986565172672272)  
[('developedp', 0.42056742310523987), ('uses', 0.3877840042114258), ('refers', 0.3803362250328064)  
[('usedp', 0.4961053729057312), ('calledp', 0.3936977982521057), ('neededp', 0.3904203772544861)  
[('usedp', 0.47560277581214905), ('foundp', 0.42083126306533813), ('soldp', 0.4192705750465393)  
[('uses', 0.4365893006324768), ('usedp', 0.4298781156539917), ('adoptedp', 0.4227108359336853)  
[('usedp', 0.5334903001785278), ('playedp', 0.4617885947227478), ('enteredp', 0.4178486466407776)  
  
[('lookedp', 0.5686172246932983), ('staredp', 0.5274225473403931), ('glancedp', 0.514441728591919),  
[('glancedp', 0.5065970420837402), ('lookedp', 0.5036084651947021), ('staredp', 0.49536192417144775)  
[('lookedp', 0.5261794328689575), ('glancedp', 0.5252094864845276), ('staredp', 0.4943307936191559),  
[('lookedp', 0.5112943649291992), ('threwp', 0.5091012120246887), ('glancedp', 0.47928571701049805),  
[('lookedp', 0.5323918461799622), ('staredp', 0.49519649147987366), ('stoodp', 0.46771761775016785)  
[('prayedp', 0.523573100566864), ('lookedp', 0.5207690596580505), ('stoodp', 0.5074241161346436), (  
[('lookedp', 0.6486713886260986), ('glancedp', 0.6086246967315674), ('staredp', 0.5800936222076416),  
  
[('lookedpp', 0.5084061026573181), ('arrivedpp', 0.43247586488723755), ('seenpp', 0.41731923818588257)  
[('lookedpp', 0.5113502740859985), ('fallenpp', 0.4844769835472107), ('comepp', 0.44388478994369507)  
[('lookedpp', 0.49028122425079346), ('turnedpp', 0.4881010353565216), ('happenedpp', 0.4839658737182617)  
[('arrivedpp', 0.40553534030914307), ('dessert', 0.39110836386680603), ('feltpp', 0.374855637550354),  
[('lookedpp', 0.5882412195205688), ('turnedpp', 0.47441810369491577), ('pushedpp', 0.45221590995788574),  
[('arrivedpp', 0.43452924489974976), ('observedpp', 0.4045757055282593), ('mentionedpp', 0.3911161422729492),  
[('lookedpp', 0.4929914176464081), ('gownpp', 0.45018041133880615), ('becomepp', 0.4078204929828644),
```

Past Participle

Preterit

Vector set

Cosine similarities between generated vector and true vector

```
[('usedpp', 0.4871392846107483), ('providedpp', 0.4284799098968506), ('recommendedpp', 0.4055340588092804),  
[('usedpp', 0.4552362561225891), ('soughtpp', 0.4343701899051666), ('foundpp', 0.41937941312789917),  
[('introducedpp', 0.43687403202056885), ('succeededpp', 0.41987723112106323), ('strippedpp', 0.4146502017974853),  
[('discoveredpp', 0.437302827835083), ('feltpp', 0.431774765253067), ('recommendedpp', 0.4293583035469055),  
[('usedpp', 0.5839736461639404), ('consideredpp', 0.5338331460952759), ('indicatedpp', 0.5320671200752258)  
[('usedpp', 0.5534017086029053), ('employedpp', 0.5186872482299805), ('indicatedpp', 0.5084676146507263),  
[('evokedpp', 0.4569481909275055), ('producedpp', 0.42817384004592896), ('seenpp', 0.4021519124507904),
```

```
[('usedp', 0.46688348054885864), ('uses', 0.4503510594367981), ('builtp', 0.4356897473335266)  
[('foundp', 0.42088770866394043), ('usedp', 0.4088103771209717), ('uses', 0.3986565172672272)  
[('developedp', 0.42056742310523987), ('uses', 0.3877840042114258), ('refers', 0.3803362250328064)  
[('usedp', 0.4961053729057312), ('calledp', 0.3936977982521057), ('neededp', 0.3904203772544861)  
[('usedp', 0.47560277581214905), ('foundp', 0.42083126306533813), ('soldp', 0.4192705750465393)  
[('uses', 0.4365893006324768), ('usedp', 0.4298781156539917), ('adoptedp', 0.4227108359336853)  
[('usedp', 0.5334903001785278), ('playedp', 0.4617885947227478), ('enteredp', 0.4178486466407776)
```

```
[('lookedp', 0.5686172246932983), ('staredp', 0.5274225473403931), ('glancedp', 0.514441728591919),  
[('glancedp', 0.5065970420837402), ('lookedp', 0.5036084651947021), ('staredp', 0.49536192417144775)  
[('lookedp', 0.5261794328689575), ('glancedp', 0.5252094864845276), ('staredp', 0.4943307936191559),  
[('lookedp', 0.5112943649291992), ('threwp', 0.5091012120246887), ('glancedp', 0.47928571701049805),  
[('lookedp', 0.5323918461799622), ('staredp', 0.49519649147987366), ('stoodp', 0.46771761775016785)  
[('prayedp', 0.523573100566864), ('lookedp', 0.5207690596580505), ('stoodp', 0.5074241161346436), (  
[('lookedp', 0.6486713886260986), ('glancedp', 0.6086246967315674), ('staredp', 0.5800936222076416),
```

```
[('lookedpp', 0.5084061026573181), ('arrivedpp', 0.43247586488723755), ('seenpp', 0.41731923818588257)  
[('lookedpp', 0.5113502740859985), ('fallenpp', 0.4844769835472107), ('comepp', 0.44388478994369507)  
[('lookedpp', 0.49028122425079346), ('turnedpp', 0.4881010353565216), ('happenedpp', 0.4839658737182617)  
[('arrivedpp', 0.40553534030914307), ('dessert', 0.39110836386680603), ('feltpp', 0.374855637550354),  
[('lookedpp', 0.5882412195205688), ('turnedpp', 0.47441810369491577), ('pushedpp', 0.45221590995788574),  
[('arrivedpp', 0.43452924489974976), ('observedpp', 0.4045757055282593), ('mentionedpp', 0.3911161422729492)  
[('lookedpp', 0.4929914176464081), ('grownpp', 0.45018041133880615), ('becomepp', 0.4078204929828644),
```

Past Participle

Preterit

Vector set

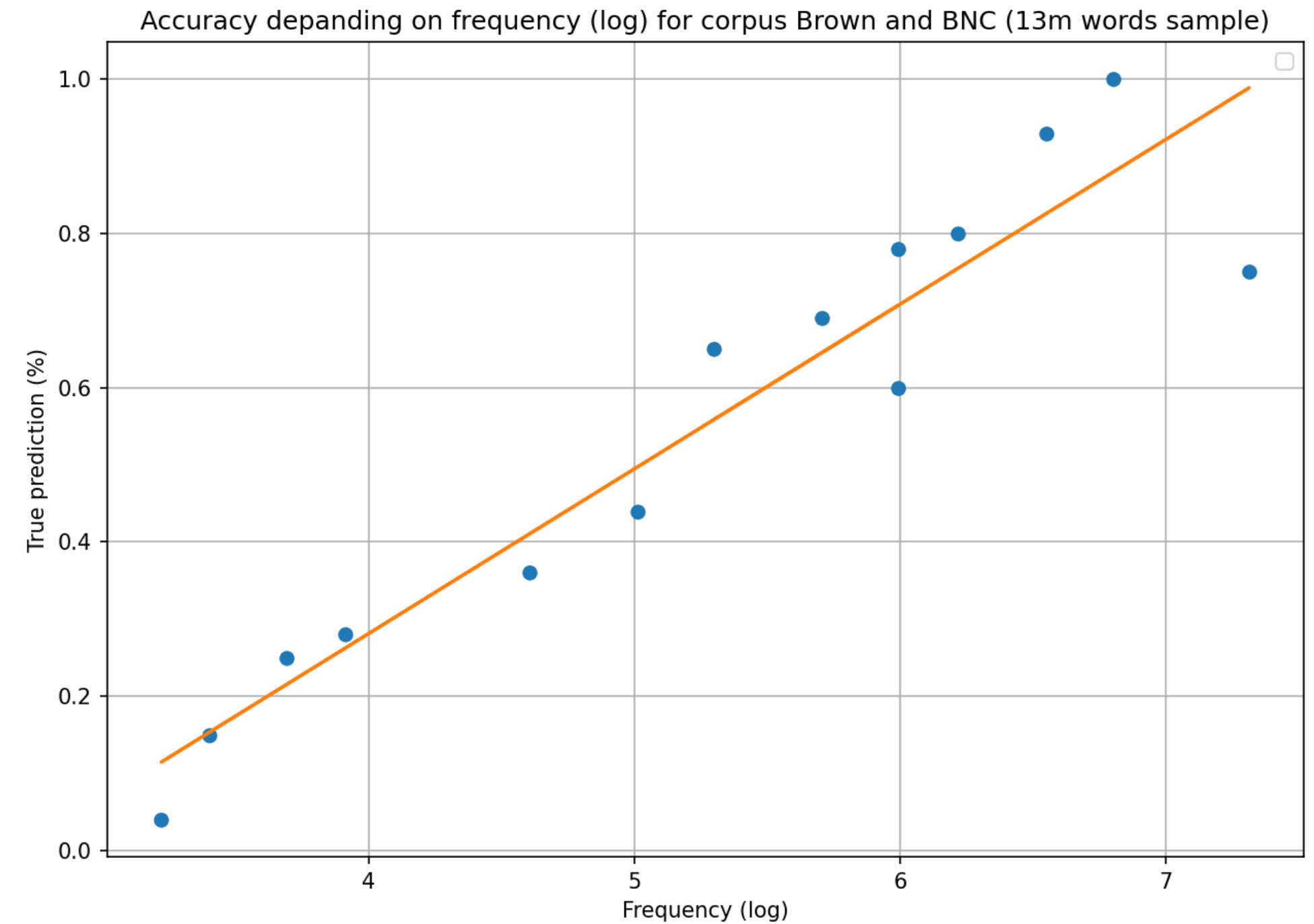
Predictions from the shift vector set

Problem: how to evaluate results?

- Proposition
 1. take the k best cosine similarities of a given set
 2. count the corresponding forms of in the results
 3. count as true prediction when the most counted form is the expected one, otherwise count as fail

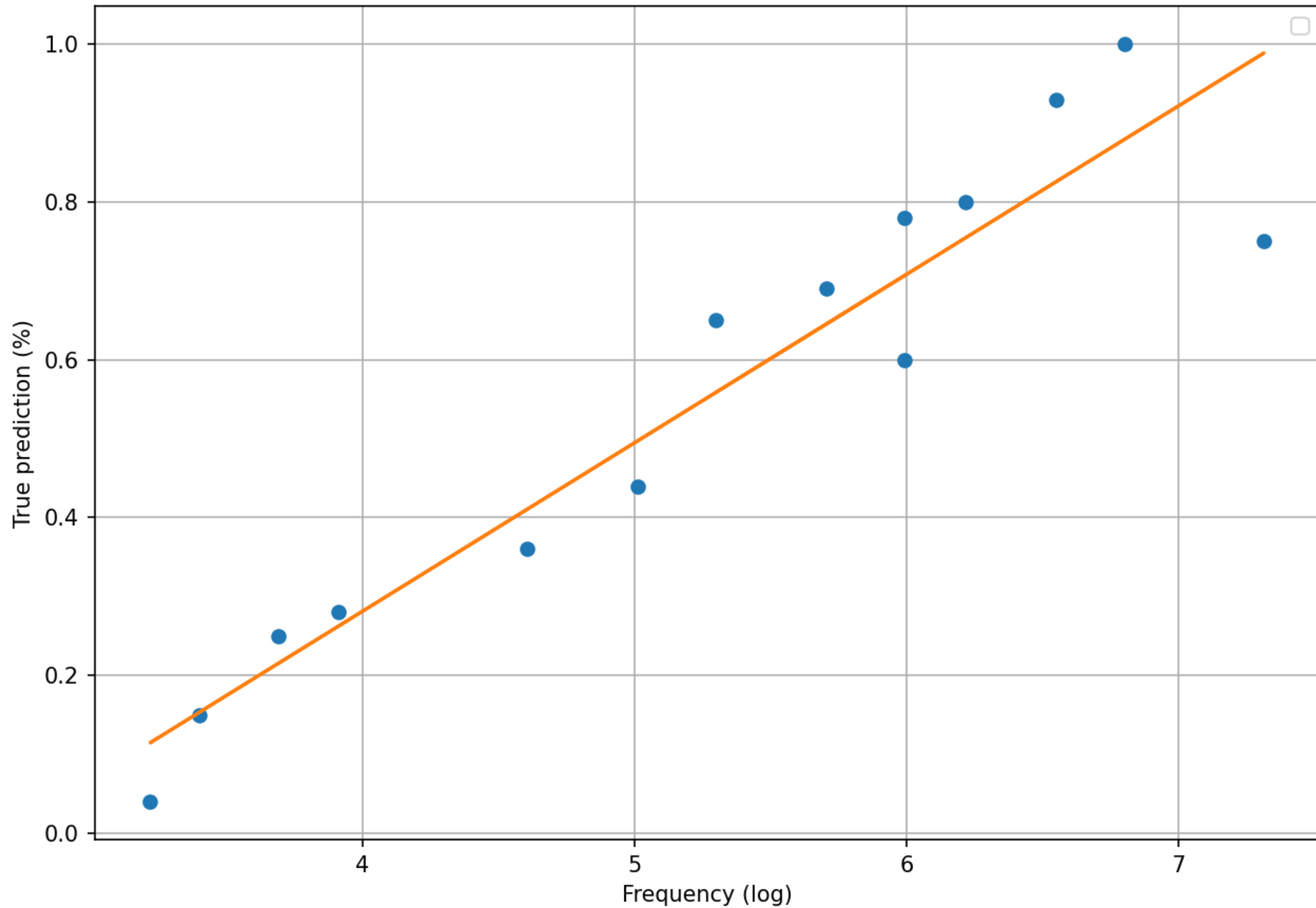
Results from the set of shift vectors

- accuracy is log-linearly correlated with the frequency of a given verb ($r=0.94$, $p\text{-value}<10^{-6}$)
- overall accuracy is 25% only on Brown (probably too small) but reaches 81% on the BNC



LinregressResult(slope=0.21350889635596002, intercept=-0.5725907469078314, rvalue=0.945112792670326, pvalue=3.504946844732338e-07, stderr=0.021308342717363406, intercept_stderr=0)

Accuracy depending on frequency (log) for corpus Brown and BNC (13m words sample)



LinregressResult(slope=0.21350889635596002, intercept=-0.5725907469078314, rvalue=0.945112792670326, pvalue=3.504946844732338e-07, stderr=0.021308342717363406, intercept_stderr=0)

Summing up the exploration of the tagged corpora

- No single shift vector solution seems viable but finding separate shift vectors for different groups of syncretic forms seems possible
 - but some clustering methods should help determine which shift vector is relevant for which syncretic form

Next hypothesis for the **untagged** corpus using the **tagged** corpus for control:

- *if a syncretic form is most similar to the sum of the vectors of the preterit and the past participle of a discriminating verb in the **untagged corpus**, the shift vector between the discriminating forms in the **tagged corpus** applied to the tagged preterit should be in close proximity to the tagged past participle*

The vectors in the untagged corpus

- a base of 10 discriminating verbs with:
 - 1 vector \vec{V}_{R1} for the preterit (e.g. ate)
 - 1 vector \vec{V}_{R2} for the past participle (e.g. eaten)
 - 1 vector \vec{V}_{R3} for the sum of the two $\vec{V}_{R1} + \vec{V}_{R2}$ (e.g. ate+eaten)
- ➔ 30 reference vectors
- 190 syncretic verbs
 - 1 vector \vec{V}_S for the common form of the preterit and the past participle
- ➔ 190 syncretic vectors

Finding the best analogy

- For each syncretic vector \vec{V}_S
 - we look for the closest reference vector \vec{V}_R
 - if it is a combination vector, $\vec{V}_{R3} (\vec{V}_{R1} + \vec{V}_{R2})$
 - in the tagged corpus, we calculate the corresponding difference, $\vec{V}'_{S1} + \vec{V}'_{R2} - \vec{V}'_{R1}$, and look for the closest vectors in the tagged corpus
 - if the results are good enough ($\vec{V}'_{S1} + \vec{V}'_{R2} - \vec{V}'_{R1} \cong \vec{V}'_{S2}$), we can introduce pairs of separate vectors in the place of the original syncretic ones:
 - $\vec{V}_{S1} = (\vec{V}_S + \vec{V}_{R1} - \vec{V}_{R2})/2$
 - $\vec{V}_{S2} = (\vec{V}_S + \vec{V}_{R2} - \vec{V}_{R1})/2$

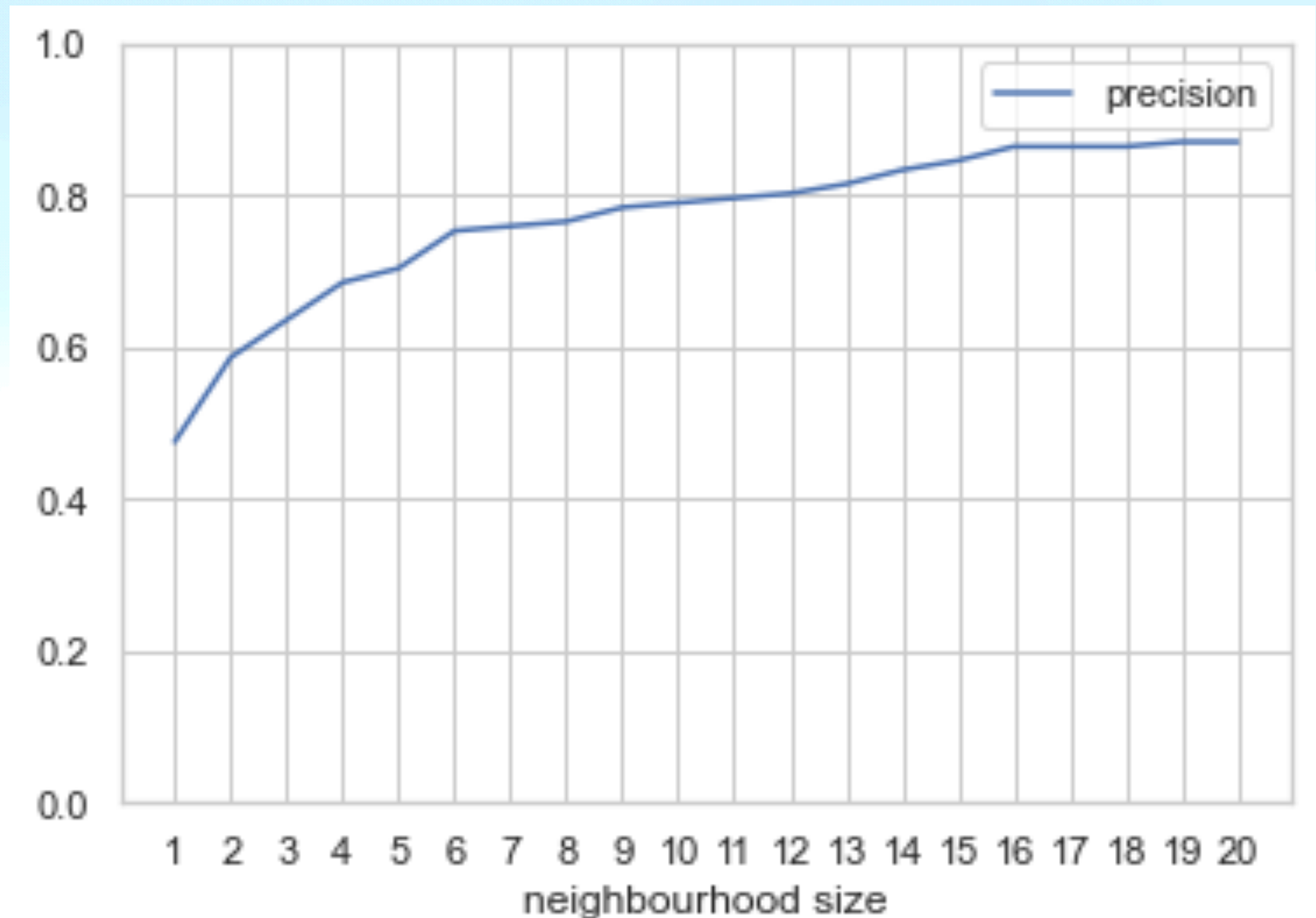
Classification results

reference vectors	syncretic vectors	proximity rank	nb
ate+eaten	found, used, provided, included, produced, bou...	51, 242, 682, 841, 287, 8, 156, 79, 98, 248, 2...	26
blew	opened	38	1
blew+blown	shot	2	1
did	wanted, let, refused	10, 68, 90	3
did+done	needed, helped, tried, meant, understood, paid...	60, 108, 20, 6, 65, 235, 53, 45, 191, 7, 282, ...	14
drew+drawn	made, brought, set, held, obtained, sought, in...	32, 1, 11, 66, 222, 229, 607, 104, 4, 88, 406,...	21
drove	started, stopped, stayed, returned, watched, hit	42, 15, 55, 35, 11, 7	6
drove+driven	moved, followed, offered, called, built, creat...	7, 16, 436, 24, 30, 529, 99, 616, 23, 378, 10,...	21
forgot	had, said, felt, talked	18, 43, 11, 2	4
forgot+forgotten	thought, told, liked, believed, remembered, he...	6, 25, 12, 20, 7, 0, 330, 56, 116, 271, 19, 31...	40
grew+grown	developed, remained, increased, tended, result...	353, 64, 775, 124, 130, 36	6
sank+sunk	formed, reached, settled, dropped, fought, dis...	176, 18, 22, 1, 75, 8, 6	7
spoke	seemed, lived, appeared, met, played, joined, ...	14, 36, 14, 20, 155, 47, 23, 218, 223, 55, 106	11
spoke+spoken	asked, read, described, shared, referred, cont...	48, 11, 731, 260, 145, 568, 267, 12, 23, 405, ...	14
threw	looked, stood	16, 14	2

Precision and neighbourhood size

Depending on the neighbourhood size chosen the evaluation of the precision varies:

- in 50% of cases, the closest vector to the computed one is the expected one
- in 90% of cases, the expected vector is in the 20 closest neighbours



Conclusion

- no single shift vector fitted all our syncretic verbs
 - but using a base of reference vectors and finding the closest reference allowed us to find sets of syncretic forms sharing the same shifting vector and recover from the accidental syncretisms
- This is exploratory work
 - the dataset was small, the syncretism was simple and we did not look at the error distribution
 - we would like to try this on French accidental syncretisms in conjugation where the patterns are more complex and the references are more ambiguous

Thank you