

# What can a corpus of YouTube videos tell us about grammatical variation in North American English?

---

**Joseph A. Stanley**

*Brigham Young University*

**Brett Hashimoto**

*Brigham Young University*

**Jack Grieve**

*University of Birmingham*

---

American Dialect Society 2026 Annual Meeting

January 11, 2026

New Orleans, Louisiana

# Background and Motivation

---

- In the 1950s, dialectologists relied on a mix of lexical, grammatical, phonological, and phonetic features.
  - Nowadays, dialect mapping primarily based on phonetic features (e.g. Labov, Ash & Boberg 2006)
- Some (but not much) corpus-based dialectology (Grieve 2016)
  - Some grammatical/lexical research, but the minority and mostly elicited rather than naturalistic (e.g. Kurath 1939, Carver 1987, Leemann et al 2018, Leemann et al 2020)
- Little that maps large areas or focuses on multiple features simultaneously (though see Kim et al 2019 and Stanley 2022)

Jack Grieve

# Regional Variation in Written American English

Table 3.1 *The complete set of 135 grammatical alternations and their 295 variants*

<b>Pronouns</b> (7, 14)	<i>As a result of/Resulting from Regardless of/No matter With the exception of/Except for/Except Get In/Into Depend+etc. On/Upon Compare+etc. To/With Concerned+etc. About/With/By Hear+etc. Of/About Complaint+etc. Of/About Concern+Concerns About/With Article+etc. About/On/Regarding/ Concerning Different from/than About/Around Number More than/Over Number Preposition Stranding/Fronting</i>
<b>Relative Pronouns</b> (4, 11)	
<i>Who/That Human Subject Zero/That/Whom/Who Human Object That/Which Non-human Subject Zero/That/Which Non-human Object</i>	
<b>Determiners</b> (8, 21)	
<i>Fewer/Less Most/Almost all With no/Without All/All of Half/Half of One/A Number Many/A lot of/Lots of/Plenty of Much/A lot of/Lots of/Plenty of/A great deal of</i>	
<b>Adjectives</b> (8, 16)	
<i>Previous/Prior Past/Last Numeral Only/Sole Un-/Not Negation In-/Not Negation Adjective -est/Most Adjective At all/Whatsoever Attributive/Predicative Adjective</i>	
<b>Nouns</b> (5, 11)	
<i>Kind of/Sort of/Type of Common Noun Genitive Of's Proper Noun Genitive Of's Pre-/Post-nominal Modification Nouns/Pronoun</i>	
<b>Prepositions</b> (22, 48)	
<i>Among/Amongst Toward/Towards Outside/Outside of Off/Off of In spite of/Despite Because of/Due to Instead of/Rather than</i>	
<b>Particles</b> (6, 12)	
	<i>Verb Away NP/Verb NP Away Verb Back NP/Verb NP Back Verb Down NP/Verb NP Down Verb Off NP/Verb NP Off Verb Out NP/Verb NP Out Verb Up NP/Verb NP Up</i>
<b>Subordinators</b> (7, 15)	
	<i>Until/Till Because/Since As long as/So long as As if/As though Although/Though/Even Though If/Whether If... then/Zero</i>
<b>Coordinators</b> (5, 10)	
	<i>But/Yet As well as/In addition to Neither... nor/or Not only... but/but also Sentence Initial And/No And</i>
<b>Verbs</b> (7, 14)	
	<i>Be Contraction/Full Pronoun Have Contraction/Full Modal Have Contraction/Full Strong/Weak Past Tense Strong/Weak Perfect Aspect Conditional Were/Was By-Passive/Non-By-Passive</i>

Table 3.1 (cont.)

<b>Modals</b> (7, 16)	<i>Almost/Nearly Numeral</i>
<i>Should/Ought</i>	<i>Only/Just Numeral</i>
<i>Will/Shall/Be going to</i>	<i>Much/Far Comparative Adjective</i>
<i>May/Might</i>	<i>Much more/Far more Adjective</i>
<i>Can/Could Question</i>	<b>Adverbials</b> (26, 60)
<i>Must/Need to/Have to</i>	<i>Previously/Formerly</i>
<i>Will Contraction/Full</i>	<i>Frequently/Often</i>
<i>Would Contraction/Full</i>	<i>Occasionally/Sometimes</i>
<b>Infinitives</b> (4, 8)	<i>Rarely/Seldom</i>
<i>To Contraction/Full</i>	<i>Repeatedly/Again and again</i>
<i>In order to/So as to</i>	<i>Immediately/Right away</i>
<i>Infinitive/Present Participle</i>	<i>Suddenly/All of a sudden</i>
<i>Seem/Seem to be</i>	<i>Meanwhile/In the meantime</i>
<b>Not</b> (6, 12)	<i>Simultaneously/At the same time</i>
<i>Be Not Contraction/Full</i>	<i>Currently/Presently/Right now</i>
<i>Do Not Contraction/Full</i>	<i>Usually/Normally/Most of the time</i>
<i>Have Not Contraction/Full</i>	<i>Clearly/Obviously</i>
<i>Modal Not Contraction/Full</i>	<i>Maybe/Perhaps</i>
<i>Ain't/Standard Negation</i>	<i>Probably/Likely</i>
<i>Not/No Negation</i>	<i>Not Only/Not Just</i>
<b>Adverbs</b> (13, 27)	<i>However/Nevertheless/Nonetheless</i>
<i>-ward/-wards</i>	<i>Therefore/Thus</i>
<i>-where/-place</i>	<i>For example/For instance</i>
<i>Very/Really</i>	<i>Actually/In fact</i>
<i>Especially/Particularly Adj.+Adv.</i>	<i>Additionally/In addition</i>
<i>Especially/Particularly Prep.+Subord.</i>	<i>Furthermore/Further/Moreover</i>
<i>Especially/Particularly NP</i>	<i>Ordinal-ol-ly</i>
<i>Totally/Completely/Entirely</i>	<i>Finally/Lastly/Last</i>
<i>Almost/Nearly Adjective</i>	Stance Adv.: Initial/Internal/Final
<i>Almost/Nearly NP</i>	Temporal Adv.: Initial/Internal/Final
	Linking Adv.: Initial/Internal/Final



Figure 4.3 *Anyone/Anybody* local spatial autocorrelation map



Figure 4.4 *Anyone* and *Anybody* simplified local spatial autocorrelation map



Figure 5.12 Northeastern dialect region cluster

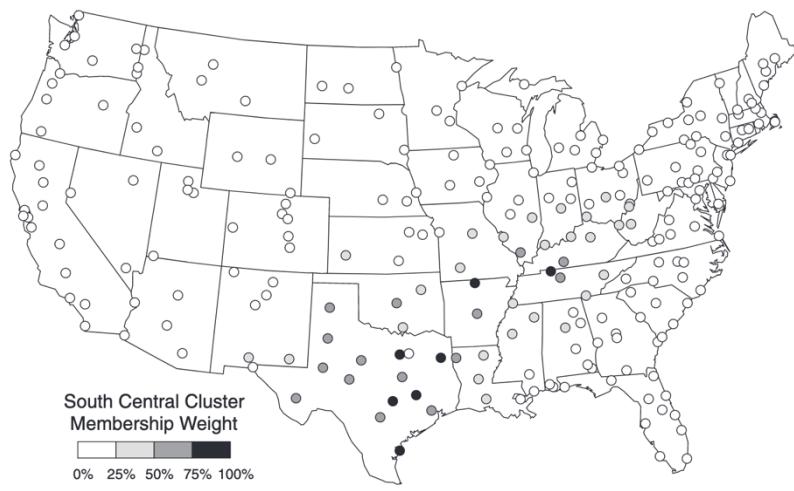


Figure 5.15 South Central dialect region cluster



Figure 5.16 Western dialect region fuzzy cluster

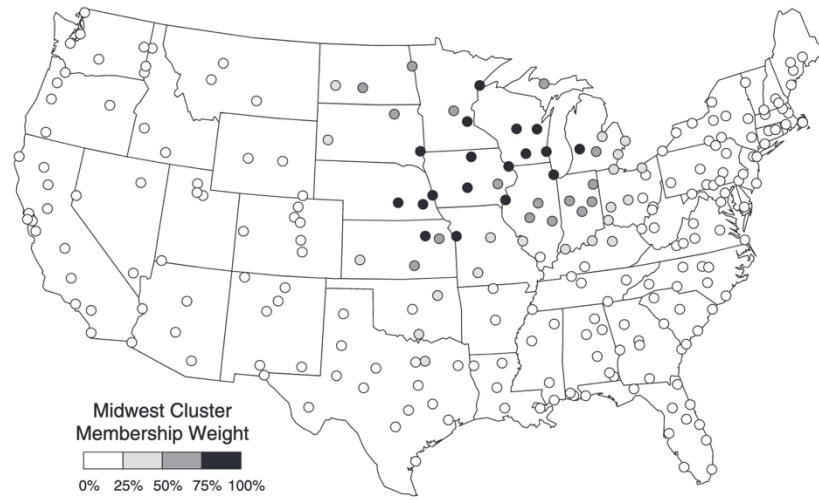


Figure 5.13 Midwestern dialect region cluster



Figure 5.14 Southeastern dialect region cluster

# Our Goal

---

Generate maps of the distributions of those  
100+ lexico-grammatical feature alternations  
in spoken North American English.

# Corpus

---

- Corpus of North American Spoken English (CoNASE: Coats, 2019; 2023)
  - YouTube channels of mainly regional and local government entities or other governmental/civic organizations.
  - Stratified sampling from counties across the US and Canada
  - 301,847 texts; 154,041 hours of spoken language; 1,252,066,371 words
  - Autotranscribed and geotagged.
  - Stanza lemmatized; Part-of-speech tagged
- Same 135 grammatical alternation variables as Grieve (2016)
  - But, algorithms for feature identification were altered from Grieve (2016) to be more suitable.
  - Accuracy checking of features

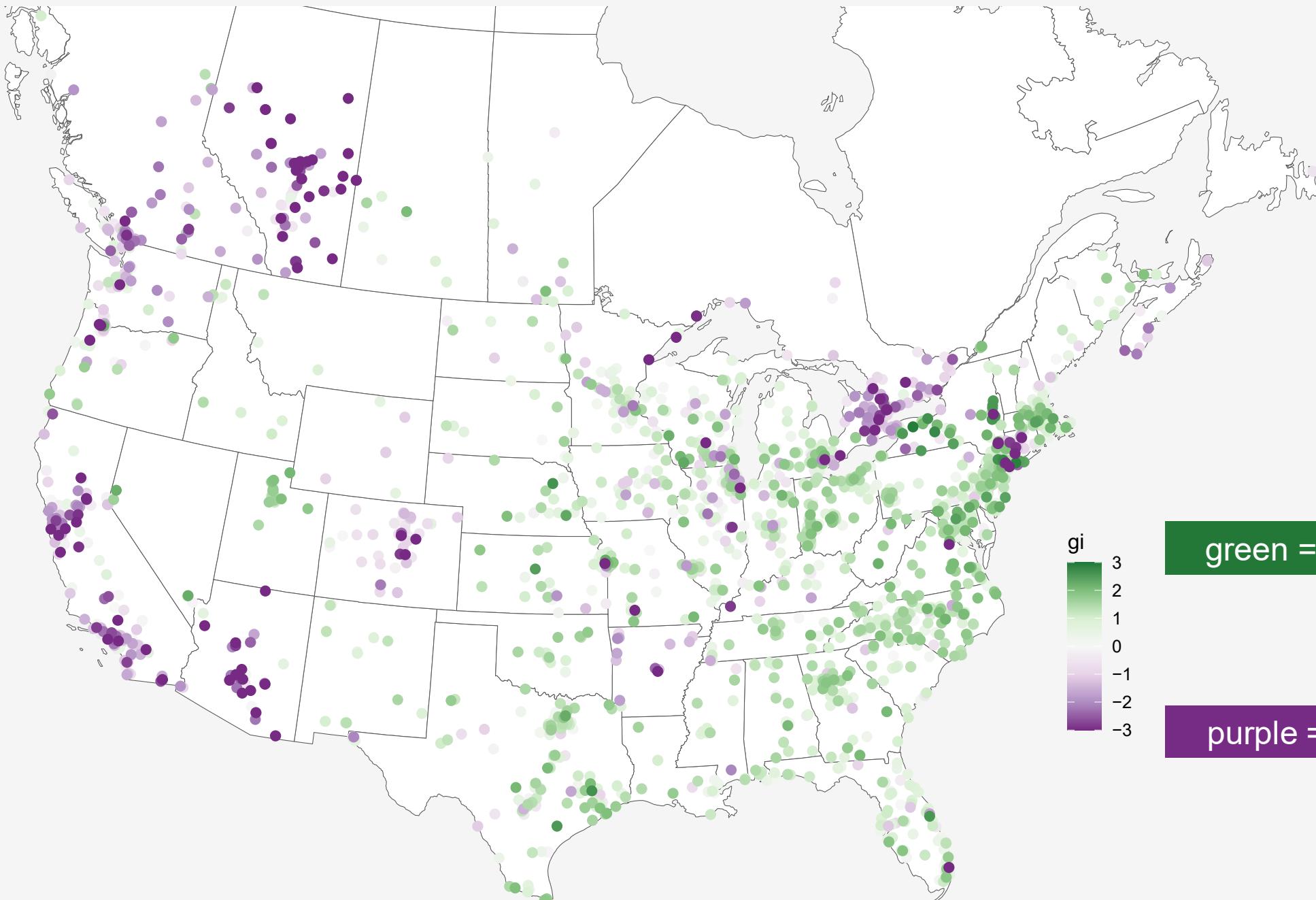
# Quantitative analysis

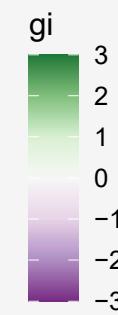
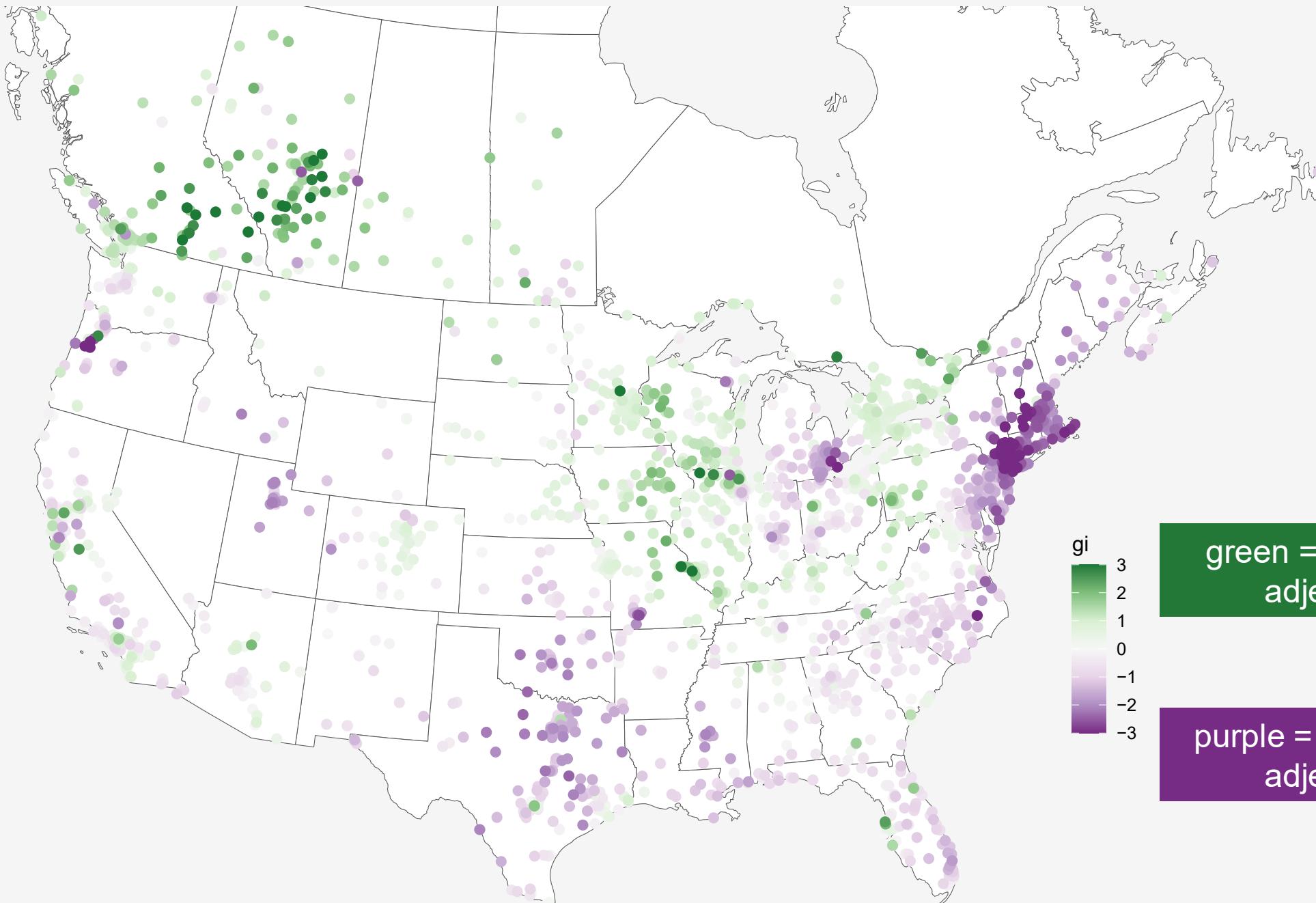
---

- Proportions by location
  - We calculated the proportion of each variant for each feature, i.e.  $A/(A+B)$ .
  - Weighted average per location (301K texts → 2,537 locations)
- Spatial stats following Grieve (2016)
  - Getis Ord-Gi statistic: For each location, indicates whether there is high/low clustering at that location (without regard to political boundaries)
  - Interpret this like a z-score, so high absolute values = statistically significant.
- Maps
  - Plot points (if there's enough data).
  - One variant is green; the other is purple.

# Results

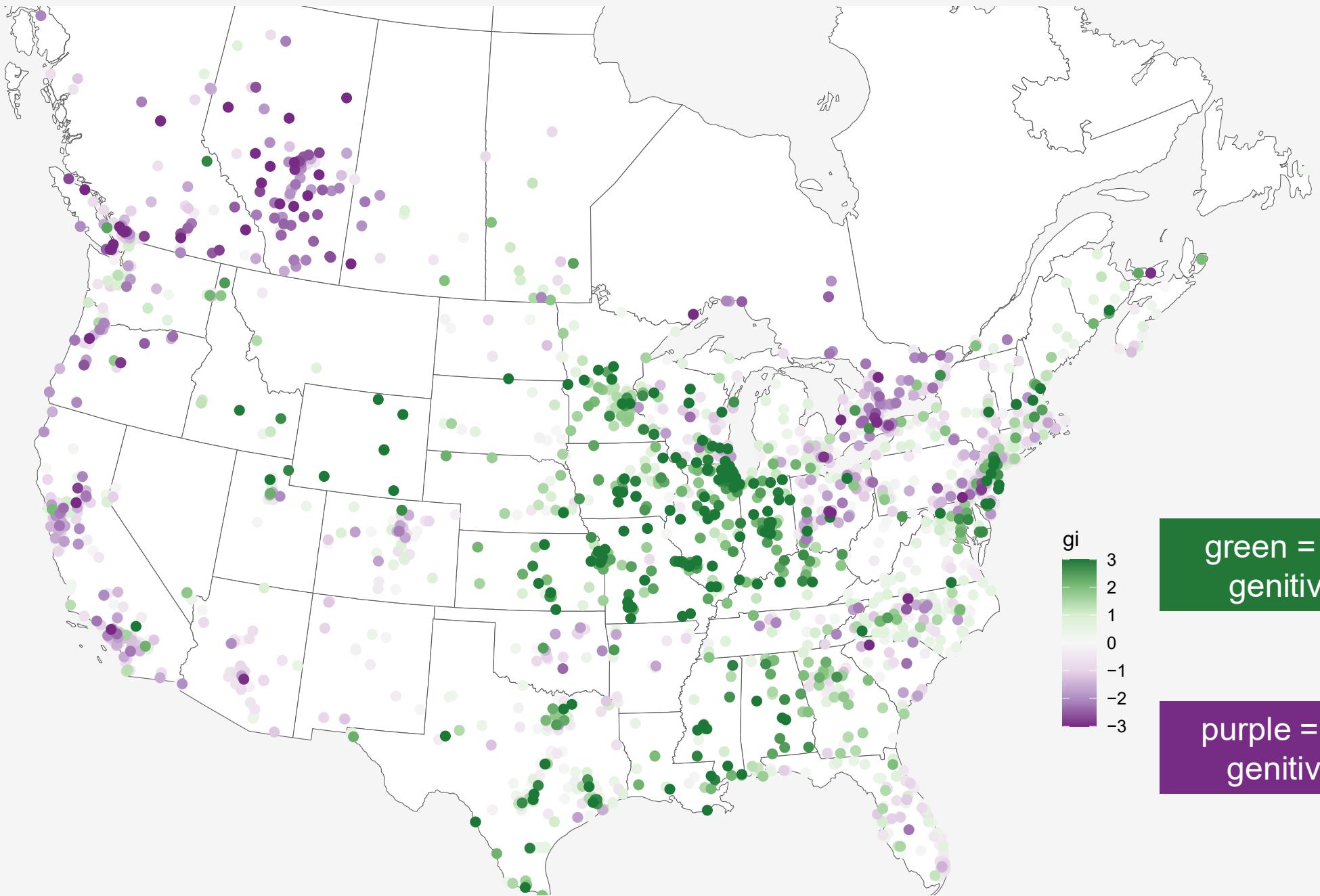
---





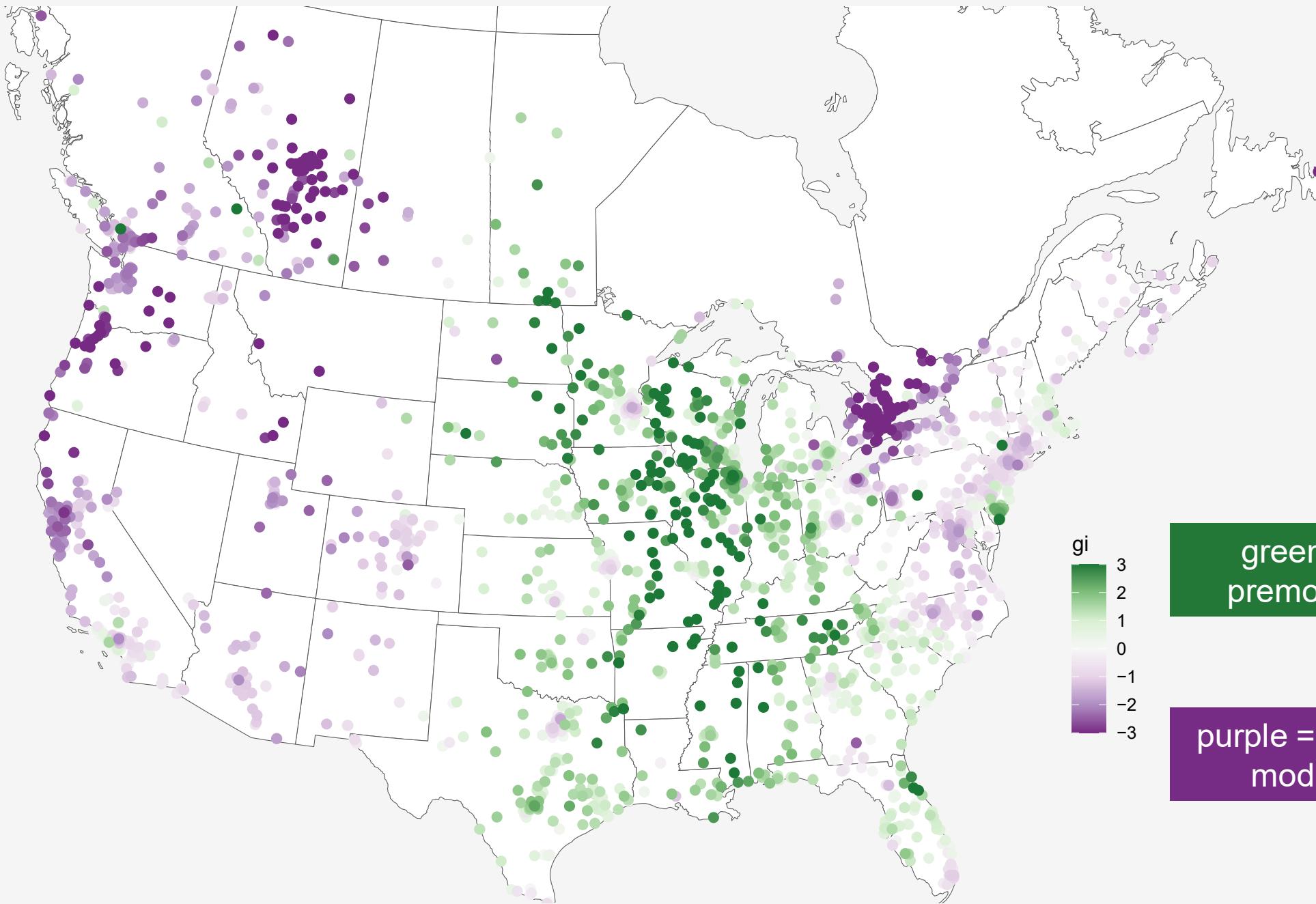
green = attributive  
adjectives

purple = predicative  
adjectives



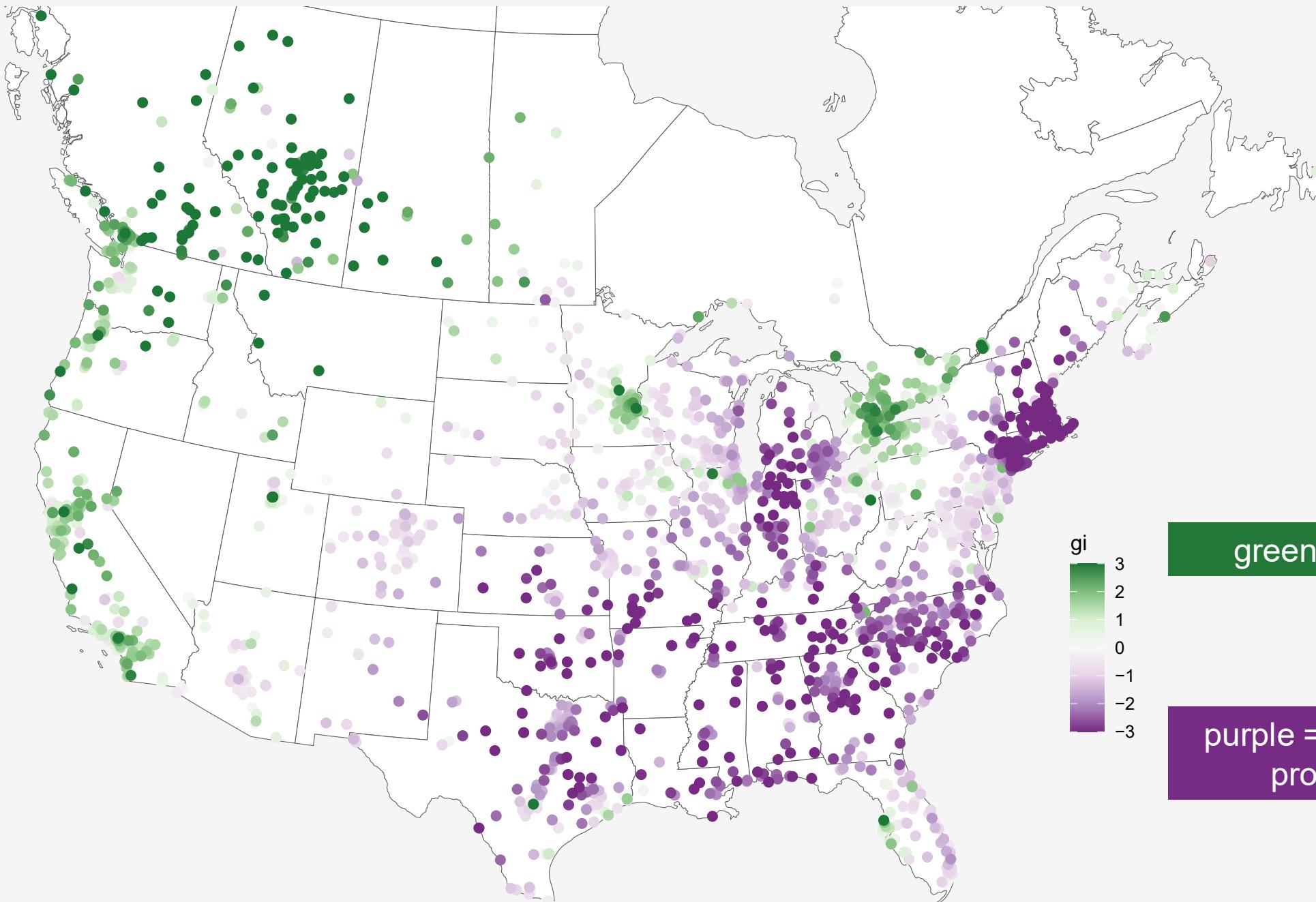
green = noun  
genitive 's'

purple = noun  
genitive *of*



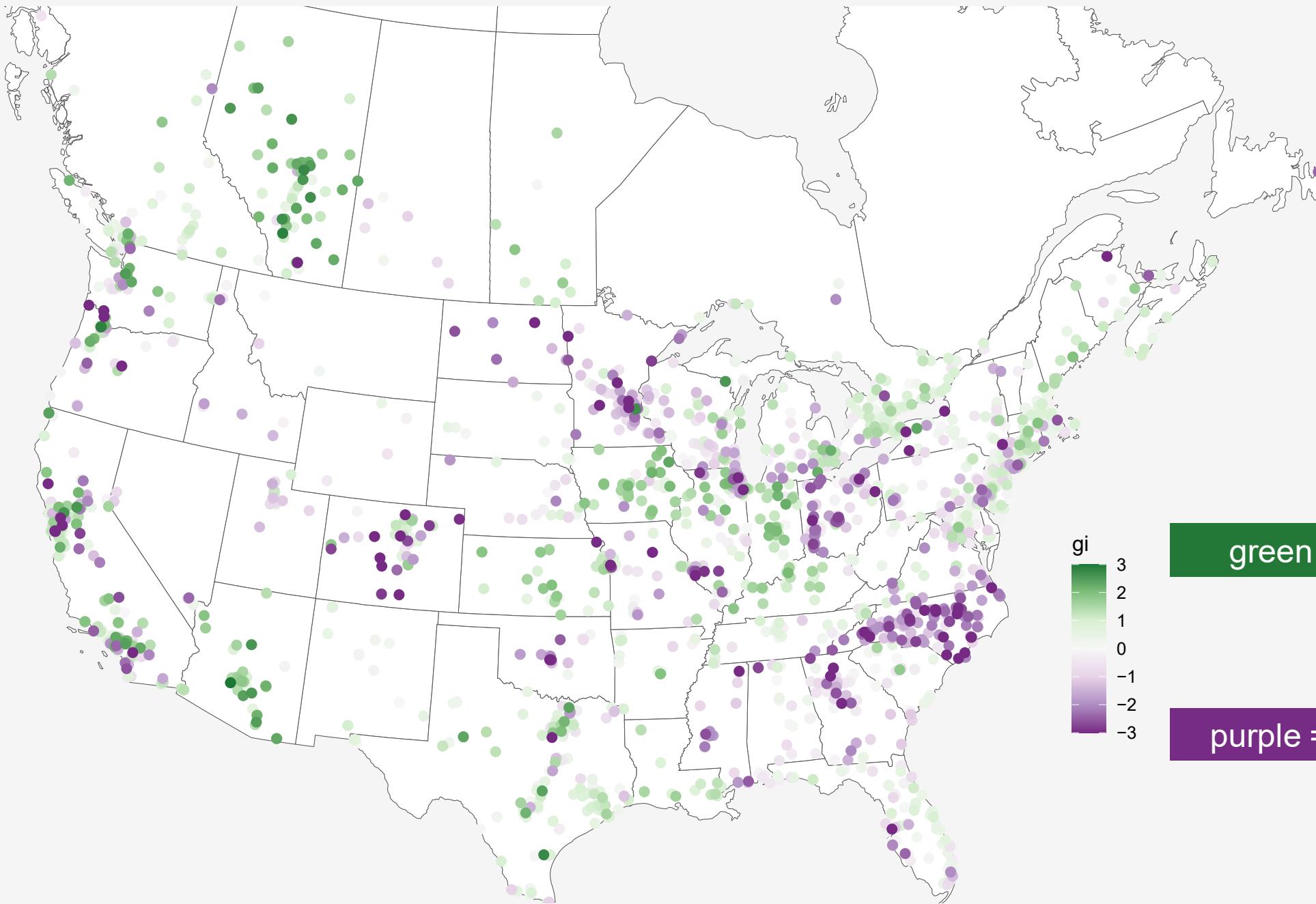
green = noun  
premodification

purple = noun post-  
modification



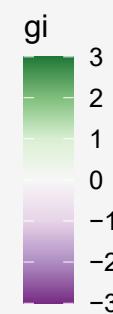
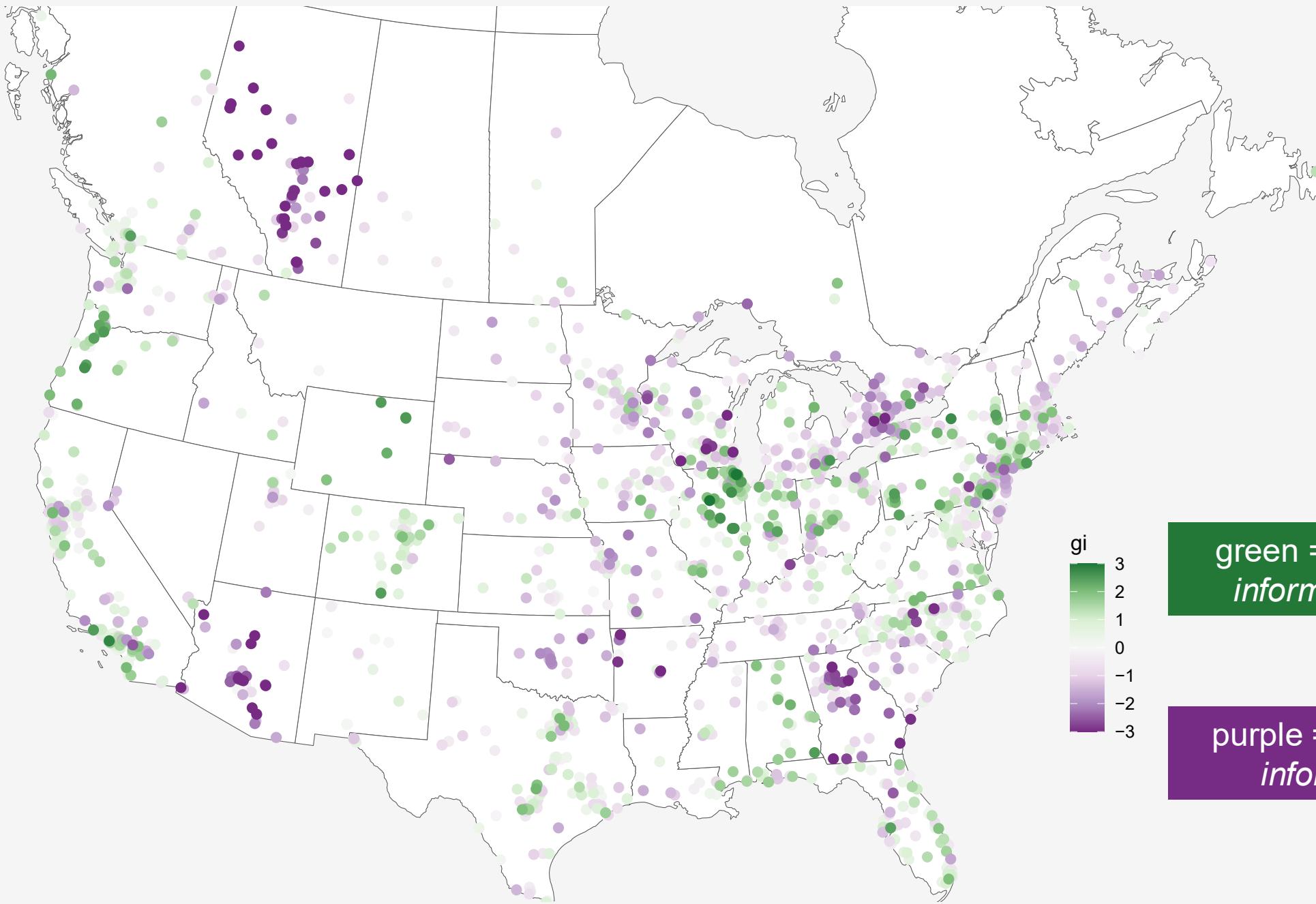
green = noun

purple = personal  
pronoun



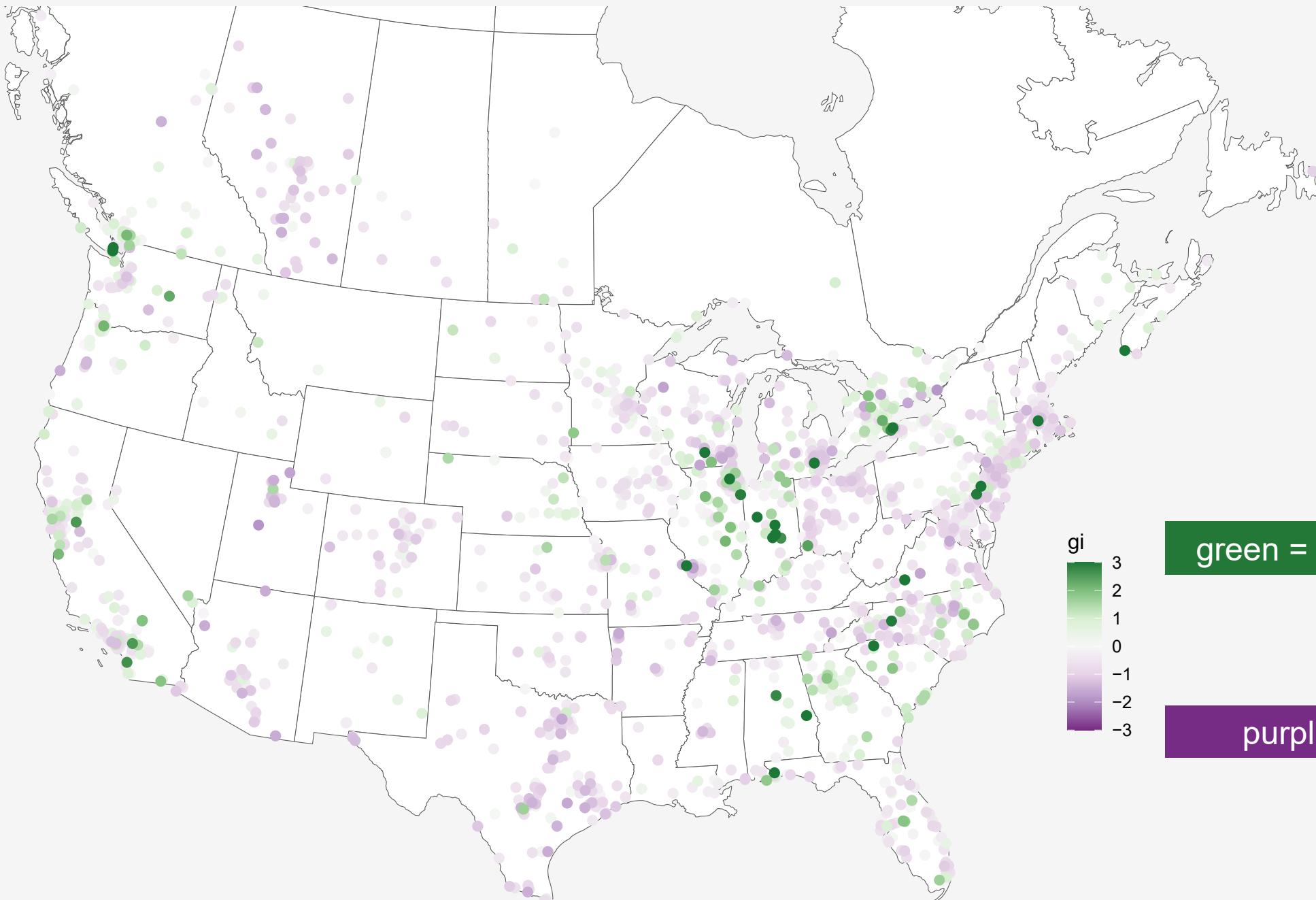
green = compare to

purple = compare with



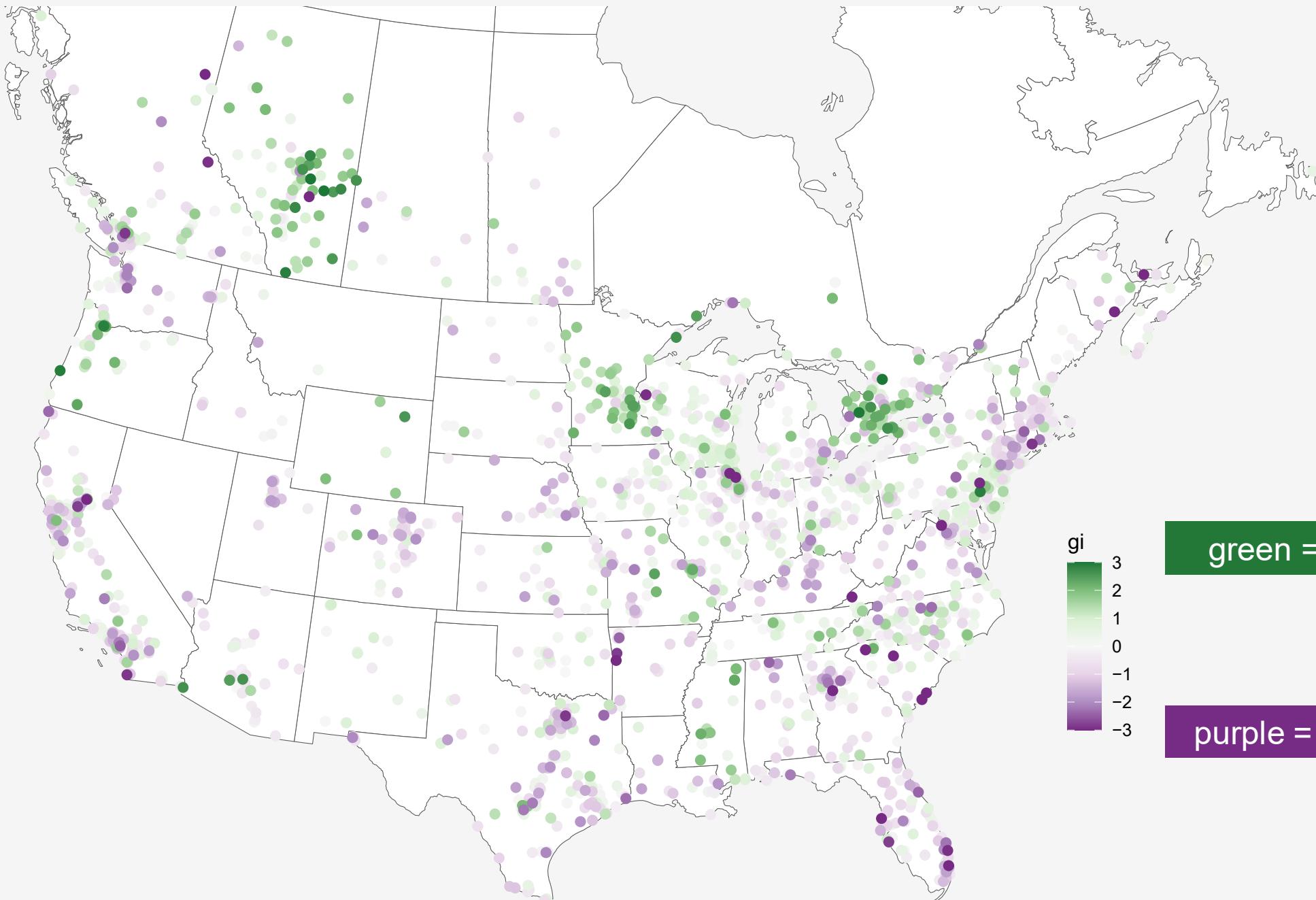
green = *article/book/  
information about*

purple = *article/book/  
information on*



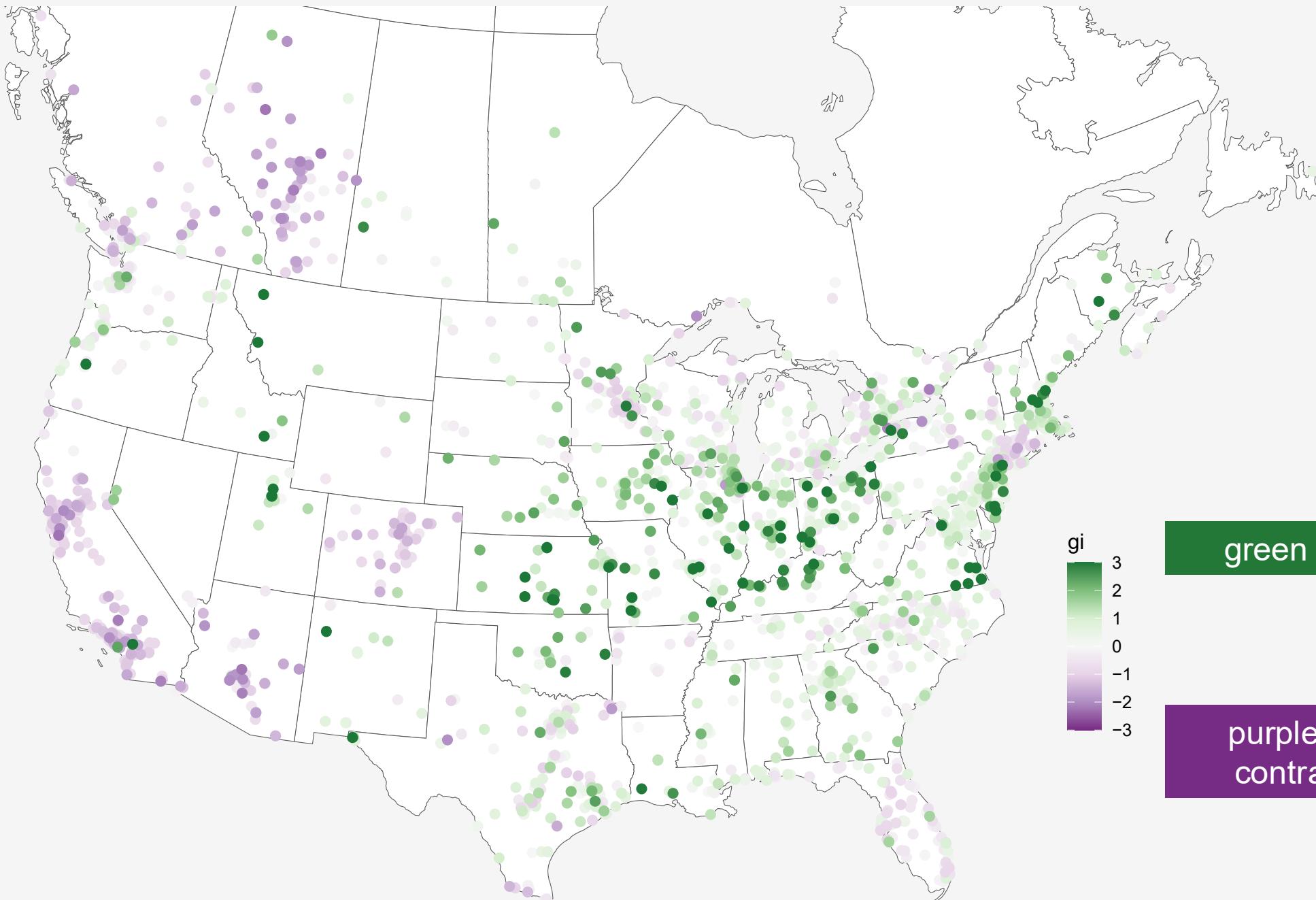
green = *whether*

purple = *if*



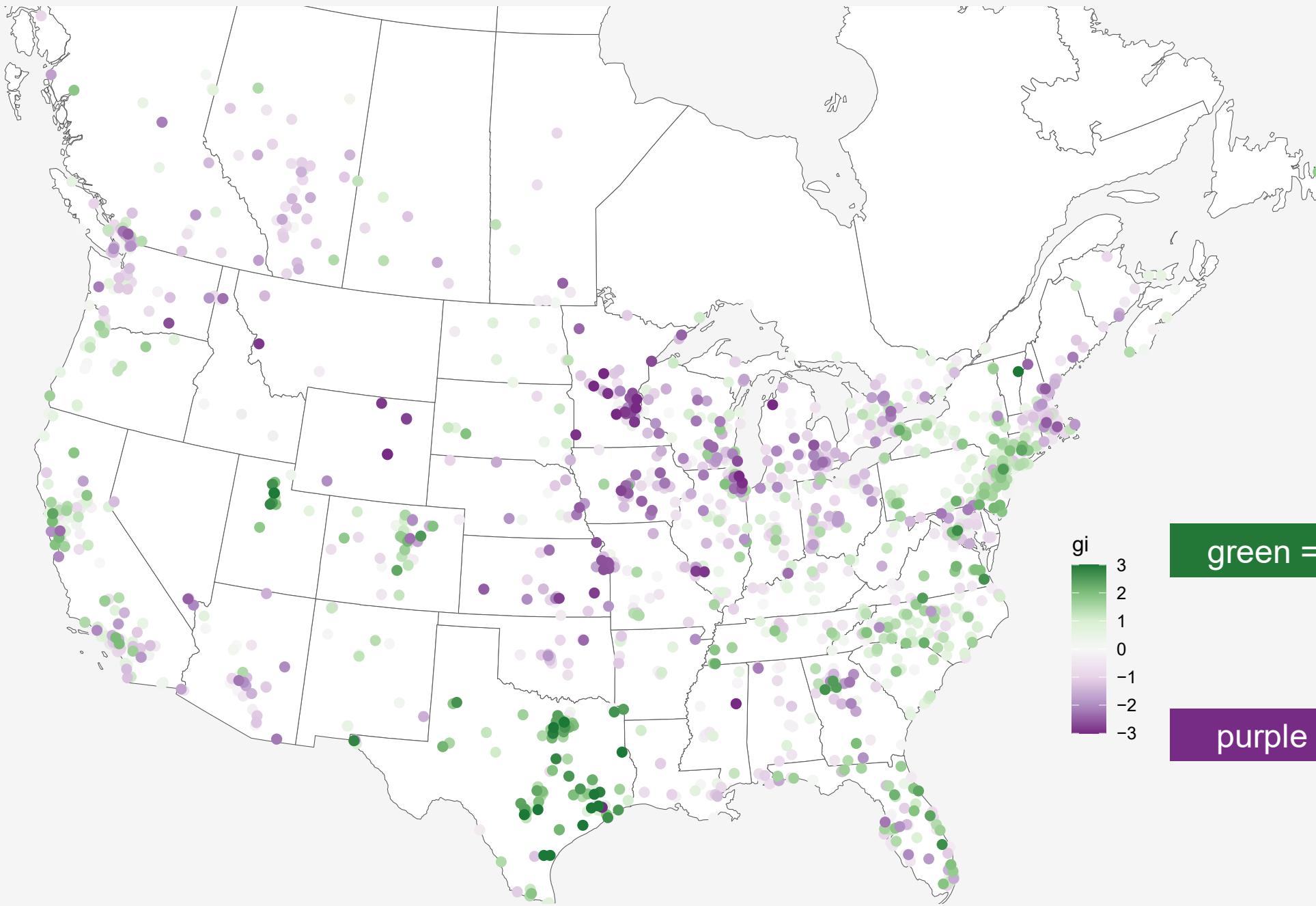
green = will

purple = shall



green = *ain't*

purple = *not*  
contraction



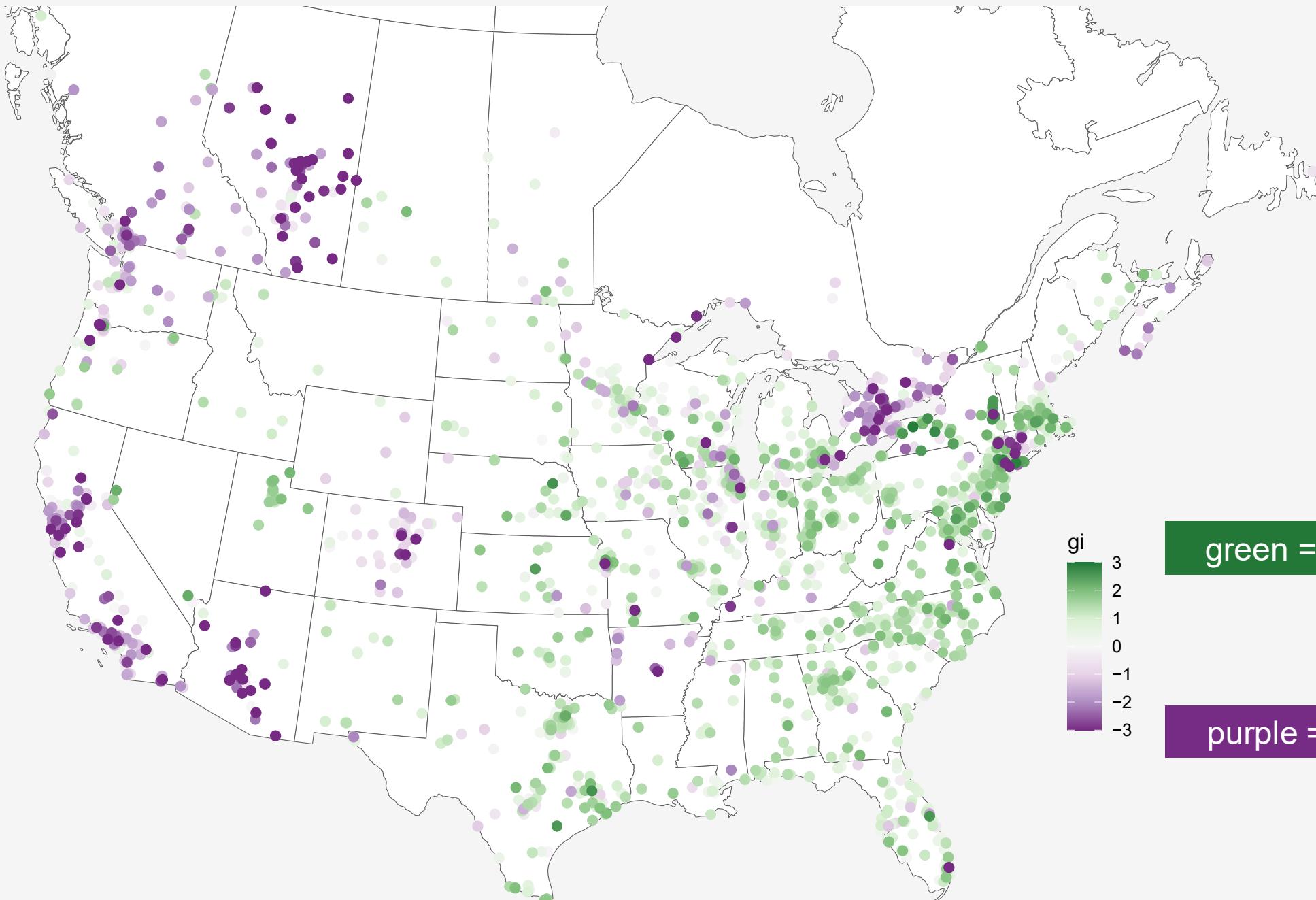
# Discussion

---

# Discussion

---

- Summary
  - Many features align nicely with known dialectal isoglosses
  - Many reveal interesting new geographic patterns.
- Limitations
  - Stylistic effects?
  - Speaker metadata?
  - Exploring the effect of spatial density of sampling locations for our spatial methods, especially when it diverges from underlying variation in population density



green = *a lot of*

purple = *lots of*

# Future Work

---

- Bigger project:
  - Compare our work with previous dialect mapping
  - Multivariate analyses (factor analysis, cluster analysis etc.)
- Future research
  - Improve accuracy of features
  - Sample underrepresented regions
  - More features!

# References

---

- Carver, Craig M. 1987. *American Regional Dialects: A Word Geography*. Ann Arbor: The University of Michigan Press.
- Coats, Steven. 2019. A corpus of regional American language from YouTube. In *Proceedings of the Digital Humanities in the Nordic Countries 4th Conference Copenhagen*, 79–91. Denmark. [http://ceur-ws.org/Vol-2364/7\\_paper.pdf](http://ceur-ws.org/Vol-2364/7_paper.pdf).
- Coats, Steven. 2024. Naturalistic Double Modals in North America. *American Speech: A Quarterly of Linguistic Usage* 99(1). 47–77. <https://doi.org/10.1215/00031283-9766889>.
- Grieve, Jack. 2016. *Regional Variation in Written American English*. 1st edn. Cambridge University Press. <https://doi.org/10.1017/CBO9781139506137>.
- Kim, Chaeyoon, Sravana Reddy, James Stanford, Ezra Wyschogrod & Jack Grieve. 2019. Bring on the crowd! Using online audio crowdsourcing for large-scale New England dialectology and acoustic sociophonetics. *American Speech* 94(2). 151–194. <https://doi.org/10.1215/00031283-7251252>.
- Kurath, Hans. 1939. *Linguistic Atlas of New England*. 6 vols. bound as 3. Providence: Brown University for the American Council of Learned Societies.
- Labov, William, Sharon Ash & Charles Boberg. 2006. *The Atlas of North American English: Phonetics, phonology and sound change*. Berlin: Walter de Gruyter.
- Leemann, Adrian, Péter Jeszenszky, Carina Steiner, Melanie Studerus & Jan Messerli. 2020. Linguistic fieldwork in a pandemic: Supervised data collection combining smartphone recordings and videoconferencing. *Linguistics Vanguard* 6(s3). <https://doi.org/10.1515/lingvan-2020-0061>.
- Stanley, Joseph A. 2022. Regional Patterns in Prevelar Raising. *American Speech* 97(3). 374–411. <https://doi.org/10.1215/00031283-9308384>.

# What can a corpus of YouTube videos tell us about grammatical variation in North American English?

---

**Joseph A. Stanley**

*Brigham Young University*

**Brett Hashimoto**

*Brigham Young University*

**Jack Grieve**

*University of Birmingham*

---

Download these slides at  
[joeystanley.com/ads2026](http://joeystanley.com/ads2026)