

1 **Interpreting the Order of Operations in a Sociophonetic Analysis**

2
3 Joseph A. Stanley

4 Brigham Young University

5 joey.stanley@byu.edu

6 ORCID: 0000-0002-9185-0048

7 8 Abstract

9 Sociophonetic data analysis involves a pipeline of processing steps to convert a raw
10 spreadsheet of acoustic measurements to interpretable results. While most studies report
11 the steps used in their pipeline, very few explicitly report their order in which those steps
12 were applied. This study analyzes a dataset containing vowel formant data from 53 speakers
13 by processing it 5,040 unique ways, each representing a different permutation of seven
14 processing steps. To analyze the effect that an order has on the overall results, pairs of
15 pipelines that differed only by swapping two adjacent steps, were compared. The most
16 important steps in the pipeline were when normalization happened, how outliers were
17 detected, and when good data was excluded. This study illustrates the what happens when
18 these steps are rearranged relative to each other in order to justify and recommend the
19 following order of operations: classifying allophones, removing outliers, normalizing, and
20 then subsetting.

21
22 Keywords: sociophonetics, data processing, quantitative methods

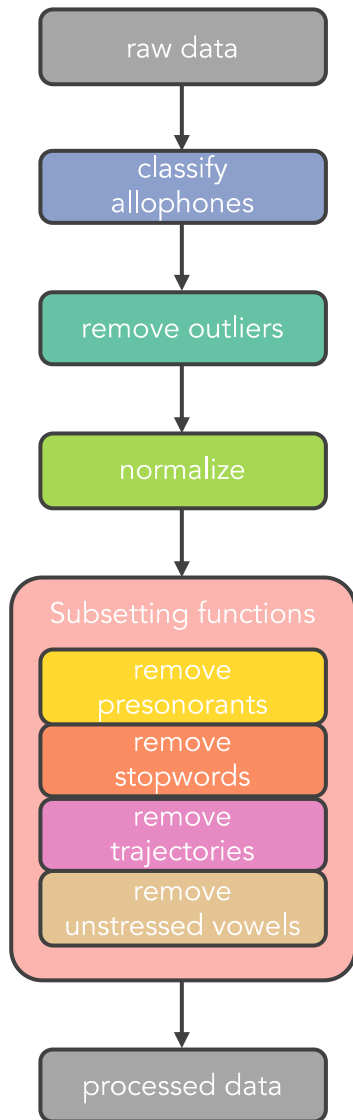
23

24 1 Order of operations

25 As an increasingly quantitative field, sociophoneticians and other linguists who work
26 with numeric data are constantly bombarded by new methods, techniques, tools, and
27 statistics. Many of these tools are meant to facilitate tasks that were once labor-intensive,
28 giving the current generation of researchers access to larger datasets and opening doors to
29 new research questions. For example, automatic formant extraction has made it possible to
30 process and analyze an order of magnitude more data than previous techniques (Labov,
31 Rosenfelder & Fruehwald 2013: 35) and it is not uncommon to hear of datasets containing
32 hundreds of thousands, if not millions, of vowel tokens (e.g. Olsen et al. 2017; Brand et al.
33 2021; Kendall & Farrington 2021).

34 Regardless of whether a project uses fully automatic or completely manual methods, the
35 typical goal of data processing is to convert audio into a spreadsheet of acoustic
36 measurements to allow for statistical analysis. To accomplish this task, one must send the
37 data through a pipeline of processing steps that includes aligning text to speech, extracting
38 acoustic measurements like vowel formants, and filtering the full dataset to include only the
39 tokens of interest. Since modern quantitative research values transparency and
40 reproducibility in methodology sections, most sociophonetics papers today fortunately
41 provide detail on what processing procedures were done to the data.

42 However, Stanley (2022) finds that it is not enough to simply report the processing steps
43 in a methods section because the *order* in which those steps applies matters. Stanley
44 observed that the overall interpretation of a dataset can be different if, for example,
45 normalization occurs *before* removing outliers as opposed to occurring *after* removing
46 outliers. To test this more rigorously, he processed the same spreadsheet 5,040 times, each
47 with a different pipeline of data analysis steps. The component parts of those analyses were
48 identical from one pipeline to the next, but it was the order in which those steps were applied
49 that varied. He found that what the results show and what the research may ultimately
50 conclude about a particular dataset can change somewhat dramatically depending on the
51 order of operations. He ends the study with a recommended order that sociophonetic
52 researchers can adopt (Figure 1), as well as a call to action, urging researchers to report their
53 analyses with greater detail.



54

55 *Figure 1: The order of operations recommended by Stanley (2022) and justified in this paper.*

56

57 Reporting the order of operations in methods sections is unusual (see Brand et al. 2021
 58 fig. 2, for a rare example) but very good for transparency. However it does little good if
 59 researchers are not familiar with how the order affects the overall results. A detailed
 60 methods section may explain that normalization happened before outliers were removed,
 61 but it is not currently clear what effect that order had on the results. How should a reader
 62 evaluate the results of one study that normalized the data before removing outliers against
 63 another study that transposed those two steps?

64 This paper addresses this gap and explores in more detail the effect that some orderings
65 are likely to have on the results of a study. It is obviously not feasible to describe every
66 possible permutation, but instead, some of the most important patterns will be illuminated
67 and explained. The goal for this paper is to arm readers with some working knowledge of
68 how to interpret others' order of operations¹ and to justify the order presented in Stanley
69 (2022).

70 **2 Method**

71 *2.1 Data*

72 The data for this study is the same as what was used in Stanley (2022). To begin, vowel
73 formant measurements were automatically extracted from interviews with 53 Western
74 American English speakers. Such raw data is typically not analyzed directly and usually goes
75 through a pipeline of processing steps that ultimately trims it down to the subset that will be
76 used for analysis. For example, this subset may contain only midpoints of stressed vowels in
77 certain phonetic environments from content words, with outliers removed and the data
78 normalized to facilitate comparison across speakers.

79 For this study, seven steps were selected, based on their apparent ubiquity in
80 sociophonetic studies on North American English, to ultimately accomplish this processing:
81 defining allophones, removing outliers, normalization, removing presonorant tokens,
82 removing stopwords,² isolating midpoints from vowel trajectories, and removing unstressed
83 vowels. These seven steps can be rearranged into 5,040 permutations or orders of
84 operations, so the original raw data was therefore processed 5,040 times, each representing

¹ The information presented in this study can, unfortunately, be misused. If a particular result is desired, such as a higher LBMS Index or lower Pillai score, a simple “fix” would be to process the data using a different pipeline. What is presented in the next sections are essentially “cheat codes” to do just that. Hopefully, as more researchers, reviewers, and editors become familiar with order of operations, they will be in a better position to identify poorly conducted and reported methods.

² Stopwords are words that are very frequent and are often subject to vowel reduction. While there is no established list of what is or is not a stopword, such lists typically contain pronouns, prepositions, other grammatical particles, and perhaps some very frequent content words.

85 a different hypothetical analysis of the same data. These permutations collectively generated
86 5,040 spreadsheets, each representing a hypothetical analysis conducted on the original
87 data.

88 For each of the resulting spreadsheets, the following three metrics were then extracted
89 for each speaker:

- 90 • Pillai scores (Pillai 1955; Nycz & Hall-Lew 2013) were calculated between /ɑ/ and
91 /ɔ/ to quantify the Low Back Merger (*i.e.* the merger of *cot* and *caught*). Pillai scores
92 were also calculated between preobstruent and prenasal allophones of /æ/ to assess
93 the degree of prenasal raising in words like *ban*, *sand*, and *hand*.
- 94 • The normalized F1 and F2 measurements of /ɛ/ and /æ/ were compared to the
95 “benchmarks” derived from the *Atlas of North American English* (ANAE; Labov, Ash
96 & Boberg 2006) as a way to determine whether a particular speaker has shifted those
97 vowels.
- 98 • The Low-Back-Merger Shift (LBMS; Becker 2019a) was calculated using the LBMS
99 Index (Becker 2019b). The shift involves the lowering and centralizing of /ɪ/, /ɛ/,
100 and /æ/ and is a more gradient measure than the binary comparison to ANAE
101 benchmarks. The LBMS Index calculates the average Euclidean distance between
102 those three vowels and /i/ in the Lobonov-normalized space.

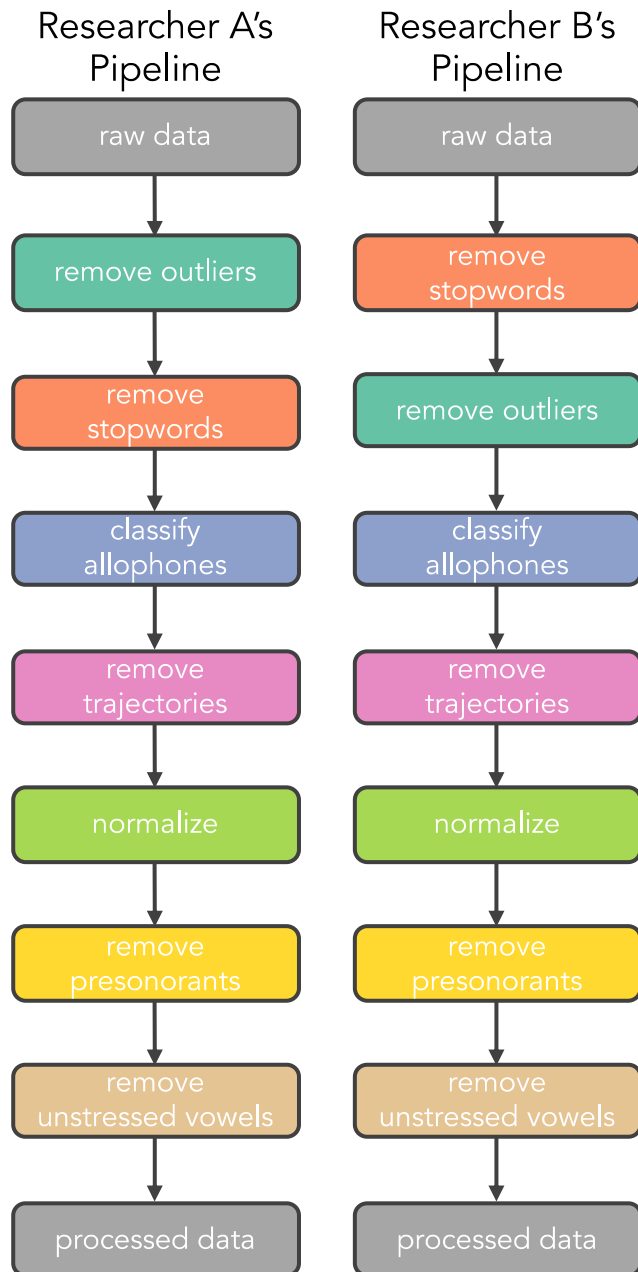
103 These three measures were chosen simply because they are commonly used in current North
104 American sociophonetic research, particularly among those who study the spread of Low-
105 Back-Merger shift across time and space.

106 In the end, a final spreadsheet was compiled and contained these five metrics (Pillai
107 scores for the Low Back Merger and prenasal raising, whether /ɛ/ or /æ/ were shifted past
108 the “benchmarks,” and the LBMS Index) for each of the 53 speakers, for each of the 5,040
109 permutations.

110 2.2 Analysis and processing

111 While Stanley (2022) describes the results from this spreadsheet in general terms, the
112 current study dives deeper by uncovering patterns in the results. Consider, for example, a
113 hypothetical speaker with the pseudonym “Justin” whose raw formant data is freely
114 available in an online repository. Researcher A processes Justin’s data using these steps

115 presented on the left side of Figure 2. Researcher B, who doing an independent analysis, uses
116 a similar pipeline, only the first two steps, removing stopwords and removing outliers, were
117 swapped, as seen in the right side of Figure 2. Researcher A concludes that Justin's low back
118 vowels had a Pillai score of 0.051. Researcher B finds it to be 0.055. This is admittedly a small
119 difference, but it is a difference nonetheless. Stanley (2022) shows that the magnitude of this
120 difference is typical across the 53 speakers in this dataset. However, the question that
121 remains is whether it is always the case that the first order of operations produces a lower
122 Pillai score than the second.



123

124 *Figure 2: Two similar pipelines used by hypothetical researchers. Note that these do not reflect the order of*
 125 *operations recommended in this paper.*

126

127 After attempting different approaches, I found that the best way to systematically
 128 explore the permutations was to compare all pairs of permutations where the only difference
 129 between them was that two adjacent steps were switched. I focus on steps that are adjacent
 130 in the recommended order of operations presented in Figure 1. By focusing on small changes

131 like swapping adjacent steps, general trends vis-à-vis those steps can still be observed. These
132 comparisons were done using the `near()` function in the `dplyr` package (Wickham et al. 2018)
133 in the R programming language (R Core Team 2021), which allows for a very small amount
134 of tolerance (on the order of eight decimal places) so that numbers that were not truly
135 identical, but for all intents and purposes are the same, could be treated as identical.

136 With these comparisons in hand, I asked two broad questions: what permutations
137 resulted in identical output and for permutations that did change the output, what reliable
138 patterns were there? The following sections answer these questions. Section 3 identifies
139 functions that, when swapped, do not change the overall result. Section 4 addresses the first
140 step in the recommended pipeline, allophone classification, and considers its relationship to
141 the second step, outlier removal. Section 5 briefly describes outlier removal and why it
142 should happen before normalization. Section 6 then focuses on normalization and its
143 relationship to the final step, the subsetting functions, described in section 3.

144 **3 Permutations that make no difference**

145 The first category of results includes identifying which permutations resulted in
146 identical output. For some of the seven functions used here, it makes no difference if they get
147 applied in different orders so long as they are adjacent to each other in the pipeline. For
148 example, removing presonorant tokens and removing stopwords have no effect on each
149 other because there are no tokens that are defined as presonorant only after stopwords have
150 been removed, and vice versa. The set of observations that are excluded in these steps is
151 fixed: regardless of the normalization procedure, whether outliers have been removed, or
152 how the vowels are classified, the exact same set of observations will be excluded each time.
153 What other pairs like these are there?

154 Upon examination of the 5,040 permutations, it appears that most of the procedures that
155 involved subsetting the data are independent of each other and can be swapped with no
156 effect on the overall results. Specifically, these functions are isolating midpoints from their
157 trajectories, removing presonorant tokens, removing tokens from stopwords, and removing
158 tokens with lexical stress. Regardless of whether these functions all happen at once or in
159 adjacent steps in the pipeline, there is no effect on the overall results.

160 For the remainder of this paper, these four functions will be collapsed down into a single
161 “subsetting” function. This is justified theoretically because they accomplish the same
162 purpose of removing data that is good but of no interest to the researcher’s current project.
163 It is also justified from a programming standpoint since they can easily be accomplished
164 using the same function. As shown in Figure 1 and described in Section 6, it is recommended
165 that these subsetting steps all happen at the end of the data analysis pipeline, after defining
166 vowel classes, removing outliers, and normalization have occurred.

167 We now turn to the permutations that do have an effect on the overall result and the
168 more complicated matter of identifying what effect they have.

169 **4 Steps 1 and 2: Allophone classification and outlier removal**

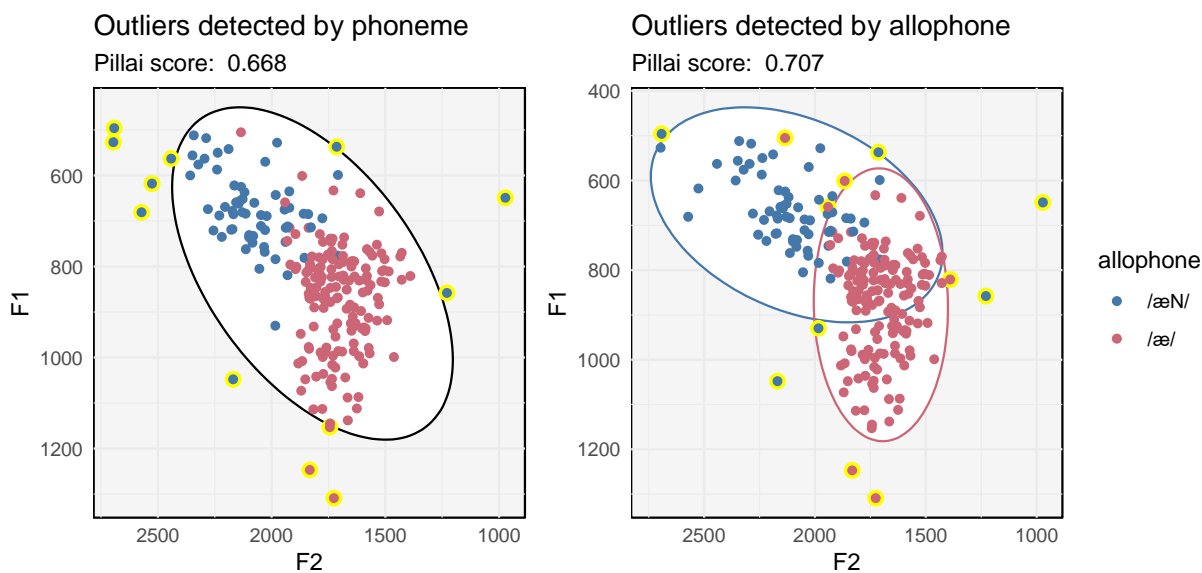
170 While the previous section described subsetting functions, which remove good but
171 uninteresting data, this section focuses on the effect of outlier removal, which excludes data
172 considered to be bad, and its relationship to allophone classification. This step is different
173 from the other subsetting functions discussed above because what is considered an outlier
174 changes depending on the steps that have already occurred.

175 Identifying and (potentially) excluding outliers is an important step in any data analysis
176 project. Ideally, all extreme values and outliers should be hand-checked, though this is
177 infeasible when analyzing large datasets. One compromise is to do a blanket exclusion of all
178 extreme observations under the assumption that they are true outliers. This is sometimes
179 done by taking the *z*-scores for each formant, for each vowel, for each speaker and remove
180 observations that are more than, say, two standard deviations from the mean. An alternative
181 is to use Mahalanobis distances, which can be thought of as finding the centroid for each
182 vowel for each speaker, and fitting an ellipse centered around that point, rotated and
183 stretched to best match the distribution of the data. How much data is encompassed by the
184 ellipse is determined by the researcher, and observations that fall outside of it are excluded.

185 To examine the effect of outlier exclusion and its relationship to classifying allophones,
186 consider prenasal raising, which is sometimes measured using Pillai scores. The vowel in
187 words like *bat* and *ban* are underlyingly /æ/, and what the Pillai scores quantify is the
188 amount of overlap between the two allophones. A higher Pillai score suggests less overlap
189 and, consequently, a more raised /æN/. As shown in the left panel of Figure 3, if outlier

190 detection happens by phoneme, that vowel's centroid will be somewhere between the two
 191 allophones. For speakers with little overlap in their allophones (or in smaller, cleaner
 192 datasets), this centroid may be in an area where no actual tokens occur. Based on
 193 Mahalanobis distances from that single center point, only the most raised tokens of /æN/
 194 and the most lowered tokens of preobstruent /æ/ will be considered outliers. Because only
 195 these tokens are removed, and low observations of /æN/ that are in the "territory" of
 196 preobstruent /æ/ and high tokens of preobstruent /æ/ that are in the "territory" of /æN/
 197 are retained, the two clusters are artificially drawn together, resulting in greater overlap and
 198 a lower Pillai score.

199 However, if allophones have already been defined, outliers can be determined and
 200 excluded by allophone. As shown in the right panel of Figure 3, this creates two centroids,
 201 appropriately centered among the cluster of tokens for the speaker's allophones. Now, some
 202 of the extreme observations that were considered outliers in the other pipeline are retained
 203 in the analysis and some of the observations that were in the territory of the other allophone
 204 have been removed. The result is that the amount of overlap between the two is smaller since
 205 they are more distinct from one another, yielding a higher Pillai score.



206
 207 *Figure 3: F1 and F2 measurements of ban (blue) and bat (red) from "Sabrina," a female speaker born in Idaho*
 208 *in 1977. Ellipses, which were determined by Mahalanobis distances, show the inclusion region(s); any*
 209 *observations outside of the region are highlighted in yellow and would presumably be excluded from further*

210 *analysis. Note that, for the purposes of providing a clean visual, the subsetting functions have already been*
211 *applied, though that is not the order of operations recommended in this paper.*

212

213 So, a simple swap of processing steps – defining allophones and outlier removal – can
214 influence the results of the study. A researcher who removes outliers by *phoneme* will find
215 less prenasal-raising (or any allophonic split) than a researcher who removes outliers by
216 *allophone*. For a single speaker this difference may be inconsequential, but when such
217 differences are consistent across many speakers, it can influence how an allophonic split is
218 interpreted in a sample, perhaps leading the researcher to determine whether that split is
219 present in the speech community. Figure 3 makes it apparent that outliers should be
220 detected at the allophonic level, so defining allophones should come before outliers are
221 detected.

222 **5 Steps 2 and 3: Outlier removal and normalization**

223 Figure 1 suggests that outliers be removed before normalization. The goal of
224 normalization is to map acoustic data to the perceptual vowel space, so that an [æ] produced
225 by a smaller person and a perceptually identical [æ] produced by a larger person line up even
226 though they are acoustically different (Barreda & Nearey 2018); essentially it aims to mimic
227 the human ear by eliminating physiological differences while maintaining sociolinguistic
228 differences. One method is the Lobonov (1971) technique which transforms formant
229 measurements into *z*-scores independently for each speaker and each formant. Another
230 technique is the one described in the ANAE (Labov, Ash & Boberg 2006: 39–40) which, for
231 each speaker, finds the mean log-transformed formant frequency, compares it to the mean
232 log-transformed formant frequency for all speakers in the sample, and uses that single
233 parameter to adjust all formants for that speaker.

234 Rather than intensely scrutinize what happens when these two steps are swapped, I
235 instead offer a logical explanation. Both normalization procedures used here rely on an
236 average vowel formant frequency, whether that be per formant (as in Lobanov) or with all
237 formants pooled together (as in ANAE). Mathematically, the mean is sensitive to outliers and
238 the presence of a single extreme observation can shift the mean away from its “true” value
239 in the direction of the outlier. It makes sense then that outliers be removed before

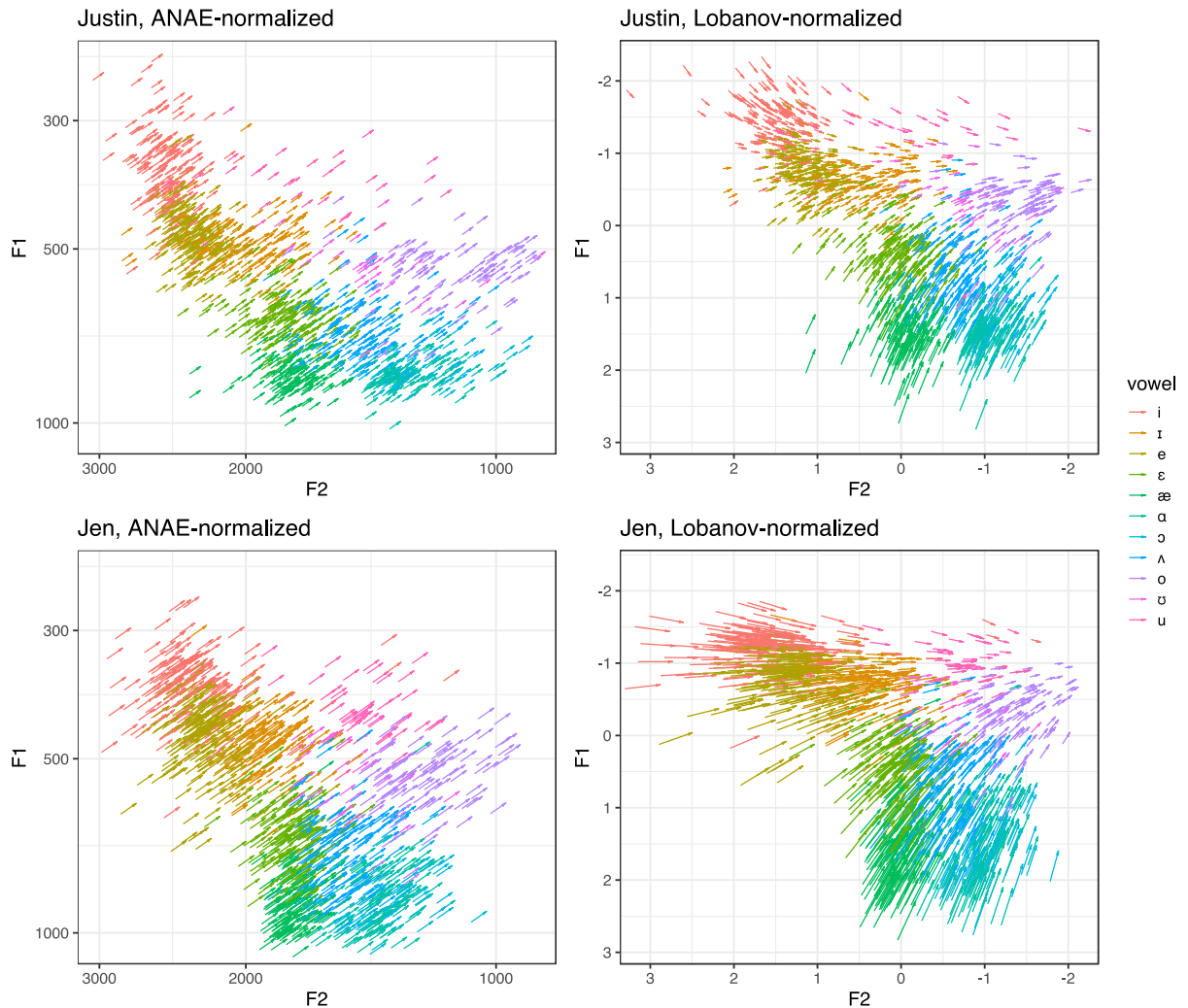
240 normalizing so that their presence does not interfere with the resulting normalized values
241 of the remaining “good” data.

242 **6 Steps 3 and 4: Normalization and subsetting**

243 Finally, we explore last pair of functions in Figure 1, the effect that vowel normalization
244 has on a dataset that includes good but irrelevant observations (such as the ones removed
245 in the subsetting step), and on one that does not. Examining the effect of normalization is
246 more complicated than other steps in this study because two different procedures are used.
247 Comparison to ANAE “benchmarks” of course necessitates the use of ANAE normalization so
248 that the values are comparable. When calculating the LBMS Index, Becker (2019b)
249 recommends using Lobonov (z-score) normalization. The Pillai score analysis here used
250 ANAE normalization, though Stanley (2021) shows that Pillai scores are unaffected by either
251 of these two normalization procedures.

252 In the recommended order of operations, normalization happens before subsetting has
253 occurred. How does this compare to normalizing after subsetting? Figure 4 illustrates this
254 comparison with both normalization procedures applied to two speakers’ data. Within each
255 of the four panels, the data is shown in two ways. For each observation, arrows start at the
256 point in the normalized F1-F2 space where that observation is located when normalization
257 occurs before subsetting (which is the recommended order). Arrows point to the location
258 where those same observations are located in the normalized F1-F2 space when
259 normalization occurs after subsetting the data. That is, each vowel token is represented by
260 an arrow and the arrows collectively show the effect on the full vowel space when
261 normalization occurs after subsetting.

262



263

264 *Figure 4: Vowel plots from two representative speakers showing how data is shifted when normalizing before*
 265 *subsetting the data to after subsetting, for two methods of normalization. Note that the panels on the left are*
 266 *on a logarithmic scale. Also note that data that was removed in the subsetting step was removed from the*
 267 *plot so that its effect on the remaining “good” data is clearer.*

268

269 In the left two panels, which show the effect of the ANAE normalization procedure, all
 270 the arrows point in the same direction towards the upper right, and within speakers the
 271 arrow length is equal for all observations. Since there is just one parameter per speaker and
 272 it applies equally to both F1 and F2, the effect of normalization is the same on all vowels for
 273 both formants. This diagonal shift has to do with finding the average formant value for all
 274 tokens for both formants – essentially locating the center of the vowel space – with higher

275 average formants resulting in a lower/fronter global position in the vowel space. A higher
276 average is compensated by the normalization procedure by lowering all formants more so
277 than a vowel space with a lower average.

278 In terms of order of operations, subsetting before normalizing causes the ANAE
279 procedure to lower both formants for all observations, resulting in a higher/backer position
280 in the vowel space. This means that the data that is removed in the subsetting step is, on
281 average, higher/backer than the remaining data. It is difficult to determine precisely why
282 this shift is consistently upward for all speakers since the tokens that are removed
283 (unstressed vowels, stopwords, irrelevant allophones) are not typically thought of as being
284 higher/backer than the data that is retained. But one possible reason is that because
285 trajectories, particularly those of elsewhere allophones, are often U-shaped (Nearey 2013
286 Fig. 11; Stanley 2020 Fig 4.4), midpoints are, on average, lower in the vowel space than their
287 trajectories. Therefore, only retaining midpoints (e.g. the data represented by the tail end of
288 the arrows) brings the average formant value up, meaning a lower position in the F1-F2
289 space. This lower position means the normalization procedure must compensate by pulling
290 all formants down, which results in a higher/backer position in the vowel space (e.g. the data
291 represented by the arrowheads).

292 The right two panels of Figure 4 are more complicated. This time, instead of adjusting all
293 points equally, the Lobonov normalization procedure has a varying types and amounts of
294 influence, depending on the speaker and the formant. For the two speakers shown there, we
295 see that higher formants lower and lower formants raise, resulting in a vertically
296 compressed vowel space. All their vowels back somewhat, though the degree of backing is
297 conditioned by F2, with front vowels backing more than back vowels. The pattern is similar
298 but not identical across speakers; in Figure 4 we see that the two speakers' data was affected
299 in slightly different ways. Indeed, each of the 53 speakers in this sample had a different
300 instantiation of this pattern, suggesting that Lobonov exerts an influence that is less
301 predictable than what the ANAE procedure does.

302 As far as the reason for the compression/backing pattern, this can be explained by
303 thinking about the kind of data that is removed in the subsetting step and how that excluded
304 data might affect the mean and dispersion of the vowel space. Because trajectories are U-
305 shaped, vowel data that includes trajectories (e.g., the version of data represented by the

306 arrowheads) will be, on average, higher in the F1-F2 space than vowel data that only includes
307 midpoints (e.g. the version of the data represented by the tail end of the arrows). So to
308 compensate for this difference in height, the normalization procedure brings the higher
309 portion of the vowel space down, more than it would for the lower part of the vowel space.
310 This is represented visually in Figure 4 by the lowered version of the vowel space in the right
311 panel. The compression in the vowel space is the result of the reduced vowels typical of
312 unstressed tokens and stopwords. The vowel space before subsetting is considered less
313 disperse and has a lower standard deviation because of the dense clustering of tokens in the
314 center of the vowel space. Again, the normalization procedure compensates for this by doing
315 the opposite; in this case, the denser vowel space is spread out more.

316 Disregarding the difference between the normalization procedures for the moment,³
317 Figure 4 should cause some concern when considering why these particular normalization
318 procedures are used. For the ANAE procedure, all vowels for a speaker appear higher and
319 backer if normalization happens after subsetting the data. Recall that the ANAE procedure is
320 often used to see whether certain vowels are, on average, past the “thresholds” that have
321 been derived from ANAE’s maps. Subsetting the data before normalization may result in a
322 speaker’s /ε/ to be higher than the threshold of 650Hz, while subsetting the data after
323 normalization could result in that speaker’s /ε/ lower than the threshold. In Stanley (2020),
324 I show that this is precisely what happened for 39 of the 53 speakers in this dataset. Since
325 the order of operations used in the ANAE is not known, it is impossible to make an accurate
326 comparison.

327 For the Lobonov procedure it is more complicated because normalizing after subsetting
328 the data resulted in a vertically compressed vowel space and backer vowels. Artificial vowel
329 space compression like this means vowels appear closer together in the normalized vowel
330 space, and since the LBMS Index (among many other measures not mentioned in this study)
331 quantifies distances between vowels in the Euclidean space, those distances will be shorter.
332 In other words, normalizing after subsetting the data will typically result in a lower LBMS

³ It is out of the scope of this paper to evaluate the pros and cons of these two normalization procedures, though it appears that the error introduced by the ANAE procedure is smaller and more predictable than that by the Lobonov method. I refer interested readers to Barreda & Nearey (2018) for a more in-depth treatment.

333 Index and those speakers will appear less shifted than if normalization happened before
334 subsetting the data.

335 So which pipeline is preferred? Unlike the subsetting, removing outliers, and allophone
336 classification steps, the reason for why normalization should occur where it does in Stanley's
337 (2022) recommended order is not as clear. While it is true that swapping the normalization
338 and subsetting steps produces different results, it is difficult to determine which one of those
339 results is better. For now, I can only recommend that sociophonetic pipelines normalize
340 before subsetting. The reason is that it may be important to include *all* data from any one
341 speaker – including stopwords, irrelevant allophones, unstressed vowels, and vowel
342 trajectories – so that the resulting normalized dataset represents that speaker's full vowel
343 space in the normalized, perceptual vowel space.

344 One may argue that subsetting should happen before normalization so that comparisons
345 across studies are more meaningful. It is true that a study that collects interview data will
346 exclude many more tokens than another study that only elicits wordlist data. An argument
347 can be made that by subsetting before normalization, what remains from the interview study
348 more closely matches what is even collected in the wordlist study, and therefore the effect
349 of normalization is more comparable between the two. However, a follow-up study of the
350 interview data that happens to focus on something that was previously excluded, like vowel
351 trajectories, will end up with a different input into the normalization step than the first study.
352 There will be a difference between the normalized data in the first study and the normalized
353 data in the second study. In other words, a single token of [æ] will have different normalized
354 F1-F2 measurements between the two studies. This makes no sense since the token has not
355 changed whatsoever. For this reason, I believe it is important to normalize before subsetting
356 so that any analysis of that particular token of /æ/ will be based on the same normalized F1-
357 F2 measurements.⁴

⁴ This order may cause issues with imbalanced across speakers (some speakers producing far fewer vowel tokens while others producing many more) or within speakers (an over- or underrepresentation of a particular vowel, which pulls the mean in its direction) (cf. Brand et al. 2021). Again, it is outside the scope of this paper to evaluate different normalization procedures, but these problems may be indicative of an imperfect

358 7 Conclusions

359 The previous sections have highlighted just a few of the consequences that may occur
360 when a dataset is processed with different orders of operations. Classifying vowels into
361 allophones, stopwords, and unstressed vowels should happen first so that outlier detection
362 can do a better job at finding truly deviant tokens. Outliers should be removed before
363 normalization so that bad data does not have an effect on where in the normalized F1-F2
364 space the good data is located. Subsetting functions like isolating midpoints from their
365 trajectories, removing uninteresting allophones, removing stopwords, and removing tokens
366 with lexical stress do not interact and should be done more or less simultaneously and
367 should happen after normalization occurs so that the vowels of interest end up in the same
368 position in the normalized vowel space regardless of whether vowel trajectories, unstressed
369 vowels, etc. are studied. The order of operations is important and should be included in
370 methods sections and preregistrations⁵ of studies (Nosek et al. 2018). This study has shown
371 why the order is important and how to interpret the order of operations that are presented
372 in future papers.

373 An important takeaway of this study is that the full range of vowel pronunciations for a
374 given speaker should be collected, if possible. If a script is set up to only extract acoustic
375 measurements from midpoints, or if a filter is set in place to skip over stopwords and
376 uninteresting allophones, that essentially forces part of the subsetting step to occur first in
377 the pipeline. The consequences of subsetting that data before normalizing is a shift in the
378 vowel space in the direction of the arrows in Figure 4. If those filters were removed from
379 that script and new acoustic measurements were extracted from the same audio, the
380 normalized data will be in the opposite direction of the arrows shown in Figure 4, even
381 though they represent the exact same vowel tokens.

382 Quantitative sociophonetic analysis has come a long way in the past several decades. An
383 explosion of methods, techniques, functions, and procedures have been proposed and used.
384 The dust has settled somewhat and some of these procedures have become standard because

procedure. Additional work that simulates these imbalances (e.g. Barreda & Nearey 2018) may be fruitful in the quest for a normalization procedure that truly models the human ear.

⁵ I thank the anonymous reviewer who recommended that I add this suggestion.

385 they have been shown to be empirically superior in some way; we are becoming more
386 equipped to properly analyze sociophonetic data. The dust has continued to settle enough
387 for this study to point out the importance of order of operations. Future scholars should
388 continue to scrutinize our linguistic methods so that we can better analyze sociolinguistic
389 variation.

390

391 References

392 Barreda, Santiago & Terrance M. Nearey. 2018. A regression approach to vowel
393 normalization for missing and unbalanced data. *The Journal of the Acoustical Society*
394 *of America* 144(1). 500–520. <https://doi.org/10.1121/1.5047742>.

395 Becker, Kara (ed.). 2019a. *The Low-Back-Merger Shift: Uniting the Canadian Vowel Shift, the*
396 *California Vowel Shift, and short front vowel shifts across North America* (Publication
397 of the American Dialect Society 104). Durham, NC: Duke University Press.

398 Becker, Kara. 2019b. Introduction. In Kara Becker (ed.), *The Low-Back-Merger Shift: Uniting*
399 *the Canadian Vowel Shift, the California Vowel Shift, and short front vowel shifts across*
400 *North America* (Publication of the American Dialect Society 104). Durham, NC: Duke
401 University Press.

402 Brand, James, Jen Hay, Lynn Clark, Kevin Watson & Márton Sóskuthy. 2021. Systematic co-
403 variation of monophthongs across speakers of New Zealand English. *Journal of*
404 *Phonetics* 88. 101096. <https://doi.org/10.1016/j.wocn.2021.101096>.

405 Kendall, Tyler & Charlie Farrington. 2021. The Corpus of Regional African American
406 Language. Eugene, Oregon: The Online Resources for African American Language
407 Project. <http://oraal.uoregon.edu/coraal>.

408 Labov, William, Sharon Ash & Charles Boberg. 2006. *The atlas of North American English:*
409 *Phonetics, phonology and sound change*. Berlin: Walter de Gruyter.

410 Labov, William, Ingrid Rosenfelder & Josef Fruehwald. 2013. One hundred years of sound
411 change in Philadelphia: Linear incrementation, reversal, and reanalysis. *Language*
412 89(1). 30–65. <https://doi.org/10.1353/lan.2013.0015>.

413 Lobonov, B. M. 1971. Classification of Russian vowels spoken by different listeners. *The*
414 *Journal of the Acoustical Society of America* 49. 606–608.

415 Nearey, Terrance M. 2013. Vowel inherent spectral change in the vowels of North American
416 English. In Geoffrey Stewart Morrison & Peter F. Assmann (eds.), *Vowel inherent*
417 *spectral change*, 49–85. Berlin, Heidelberg: Springer. [https://doi.org/10.1007/978-](https://doi.org/10.1007/978-3-642-14209-3_4)
418 [3-642-14209-3_4](https://doi.org/10.1007/978-3-642-14209-3_4).

419 Nosek, Brian A., Charles R. Ebersole, Alexander C. DeHaven & David T. Mellor. 2018. The
420 preregistration revolution. *Proceedings of the National Academy of Sciences* 115(11).
421 2600–2606. <https://doi.org/10.1073/pnas.1708274114>.

422 Nycz, Jennifer & Lauren Hall-Lew. 2013. Best practices in measuring vowel merger.
423 *Proceedings of Meetings on Acoustics* 20(1). 060008.
424 <https://doi.org/10.1121/1.4894063>.

425 Olsen, Rachel M., Michael L. Olsen, Joseph A. Stanley, Margaret E. L. Renwick & William A.
426 Kretschmar Jr. 2017. Methods for transcription and forced alignment of a legacy
427 speech corpus. *Proceedings of Meetings on Acoustics* 30(1). 060001.
428 <https://doi.org/10.1121/2.0000559>.

429 Pillai, K. C. S. 1955. Some new test criteria in multivariate analysis. *The Annals of*
430 *Mathematical Statistics* 26(1). 117–121.
431 <https://doi.org/doi:10.1214/aoms/1177728599>.

432 R Core Team. 2021. R: A language and environment for statistical computing. Vienna, Austria:
433 R Foundation for Statistical Computing. <http://www.R-project.org>.

434 Stanley, Joseph A. 2020. *Vowel dynamics of the Elsewhere Shift: A sociophonetic analysis of*
435 *English in Cowlitz County, Washington*. Athens, Georgia: University of Georgia
436 Dissertation.

437 Stanley, Joseph A. 2021. Pillai scores don't change after normalization.
438 <https://joestanley.com/blog/pillai-scores-dont-change-after-normalization>. (2
439 November, 2021).

440 Stanley, Joseph A. 2022. Order of operations in sociophonetic analysis. In *University of*
441 *Pennsylvania Working Papers in Linguistics*. 28(2), Article 17. Available at:
442 <https://repository.upenn.edu/pwpl/vol28/iss2/17>.

443 Wickham, Hadley, Romain François, Lionel Henry & Kirill Müller. 2018. dplyr: A grammar of
444 data manipulation. <https://CRAN.R-project.org/package=dplyr>.

445