This manuscript was submitted to the *Journal of the Acoustical Society of America* on February 28, 2022 and is currently in a second round of revisions. The published version will have some small changes.

This version reflects changes through October 13, 2022

**Sample Size Matters in Calculating Pillai Scores**

Joseph A. Stanley,[1] and Betsy Sneller[2]

[1] *Linguistics Department, Brigham Young University, Provo, Utah 84602, United States*

[2] *Department of Linguistics, Languages, and Cultures, Michigan State University, East Lansing, Michigan 48824, United States*

**Abstract**

Since their introduction to sociolinguistics by Hay et al. (2006), Pillai scores have become a standard metric for quantifying vowel overlap. However, there is no established threshold value for determining whether two vowels are merged, leading to conflicting ad hoc measures. Furthermore, as a parametric measure, Pillai scores are sensitive to sample size. In this paper, we use generated data from a simulated pair of underlyingly merged vowels to demonstrate (1) larger sample sizes yield reliably more accurate Pillai scores, (2) unequal group sizes across the two vowel classes is irrelevant in the calculation of Pillai scores, and (3) it takes many more data than many sociolinguistic studies typically analyze to return a reliably low Pillai score for underlyingly merged data. We provide some recommendations for maximizing reliability in the use of Pillai scores, and provide a formula to assist researchers in determining a reasonable threshold to use as an indicator of merged status given their sample size. We demonstrate these recommendations in action with a case study.

## I.    INTRODUCTION

Quantifying vowel overlap is an important component of many linguistic studies, ranging from sociolinguistics to laboratory phonology. While various methods of quantifying vowel overlap have been proposed (see, e.g., Nycz & Hall-Lew 2013), Pillai scores have emerged in recent years as the most commonly used method, particularly within sociolinguistics (Hay, Warren & Drager 2006; Nycz & Hall-Lew 2013), because of its ability to measure a distinction in multivariate space while also accounting for fixed effects like phonological context. However, there is no standard value of Pillai score that is broadly accepted to be a threshold for "merged" or "distinct"; indeed, individual studies make this determination primarily by comparison between individuals in a single data set, to show that some speakers are "more merged" while others are "less merged". In this paper, we provide a critical look into how Pillai scores are calculated and reported.

We focus on demonstrating how sample size plays a major role in the resulting Pillai score. We also show that the common approach of reporting Pillai scores alone is incomplete without also reporting sample sizes and *p*-values. Using simulation data drawn from an underlyingly merged data set, we show how larger samples produce lower Pillai scores (in other words, a more "merged" score). We further demonstrate that within Pillai, it is the total *n* across both samples that matters, and provide a formula that researchers can use to determine a threshold for "merged" given their own sample size. We highlight some important takeaways about using Pillai scores to measure vowel overlap and the potential risks for across- and within-study comparisons of speakers. We end with a case study which demonstrates how researchers can implement the recommendations in this paper when analyzing real data from sociolinguistic interviews.

## A. Pillai as a measure of vowel overlap

### 1. Pillai score overview

A multivariate analysis of variance (MANOVA) is an extension of the (univariate) analysis of variance (ANOVA). The difference is that while an ANOVA evaluates whether the difference between two or more groups in a single numeric variable can be predicted by some number of categorical independent variables (such as the F2 of /u/ across older and younger participants), a MANOVA can evaluate two or more dependent numeric variables simultaneously (such as F1, F2, F3, and duration of /u/ by older and younger participants). So, a researcher analyzing American English vowels (which typically are differentiated primarily on vowel quality in F1-F2 space) could use a MANOVA to see whether a speaker pronounces two historically distinct vowel classes differently, while including the effects of duration and place of articulation as independent variables. In this case, F1 and F2 would be the dependent variables, and vowel class, duration, and place of articulation would be the independent variables.

Simplified somewhat, the null hypothesis of a MANOVA is that category membership (for instance, two historically distinct vowel classes) offers no explanatory power for any of the dependent variables. In the case of a MANOVA that is fit to vowel data, the null hypothesis is that a merger is present. In other words, there would be no way to guess which historic vowel class a particular token came from by its acoustic measurements alone. Typically, the researcher's aim is to find evidence to reject that hypothesis. We note that high *p*-values associated with MANOVAs only indicate a lack of evidence to reject the null hypothesis that the two vowel classes are the same, rather than evidence *for* the null hypothesis. MANOVA cannot prove that two vowels are merged, only that there is little evidence to suggest that they're distinct.

There are four main test statistics associated with a MANOVA to compare what the data shows to the null hypothesis: Wilk's lambda, the Lawley-Hotelling trace, Roy's largest root, and the

47  Pillai-Bartlett trace. Of these four, the lattermost is the most robust for non-normally distributed data

48  and other violations of the assumptions of a MANOVA for tests that compare more than two groups

49  (Olson 1976). In sociophonetic data, where it is more common to compare only two groups, the Pillai-

50  Bartlett trace has less advantage over the other three in statistical validity, but does benefit from being

51  both easy to run in commonly used statistical environments and relatively easy to interpret. For

52  introductory overviews of these four test statistics, including their mathematical definitions and

53  conceptual explanations, see Bray & Maxwell (1985: 27–29), John & Wichern (2012: 336), Rencher &

54  Christensen (2012: 169–188), and Upton & Cook (2014, "multivariate analysis of variance

55  (MANOVA)").

56      The Pillai-Bartlett trace, often called the *Pillai score* or occasionally just *Pillai* in linguistics

57  studies (a convention that we adopt here), ultimately comes from Pillai (1955) and Bartlett (1939). It

58  returns a value that ranges between 0 and 1, with smaller values occurring when there is greater overlap

59  between the two groups in multivariate space, and larger numbers for less overlap. In other words,

60  small Pillai scores suggest a vowel merger. In reality, determining whether a merger is present is not

61  quite as simple as merely observing overlap. Two vowels may occupy the same F1-F2 space but the

62  distinction between phonemes may be maintained though some other cue like voice quality (Di Paolo

63  & Faber 1990), duration (Labov & Baranowski 2006), or vowel trajectory (Stanley 2020). We adopt a

64  simplified approach here since our focus is to explore the effects of sample size, but we acknowledge

65  that there is more to merger than overlapping midpoints in the F1-F2 space.

66      *2.  Meta-analyses of Pillai scores vs. other metrics*

67      Prior to the introduction of Pillai scores to sociophonetics, perhaps the most common way to

68  assess merger was through auditory coding. In its most basic form, the phonetically-trained researcher

69  would listen and evaluate whether there was a difference between the two sounds. But, as shown

70  below, Pillai scores are most commonly used on low vowels, which are notoriously difficult for

71  fieldworkers to accurately transcribe (Johnson 2010: 28–29). In fact, Moulton (1968: 464) rather

72  strongly states that early fieldworkers for the Linguistic Atlas Projects were "hopelessly and humanly

73  incompetent at transcribing phonetically the low and low back vowels that they heard from their

74  informants." Fortunately, formulaically quantifying overlap provides less subjectivity that earlier

75  researchers lacked.

76      While Pillai scores are currently the most common metric for quantifying mergers, especially

77  within sociolinguistic work, there are also other approaches available. In addition to the auditory

78  coding mentioned above, much early work used the Euclidean Distance between the point

79  representing the mean F1 and mean F2 of one vowel class and the point representing the mean F1

80  and mean F2 of a second vowel class as the primary metric for vowel merger (Hay, Warren & Drager

81  2006; Nycz & Hall-Lew 2013; Han & Kang 2013; Hall-Lew 2013). This early approach is not

82  particularly satisfying, as it fails to take into account the distributional properties of the data, including

83  the degree of overlap and the distribution of tokens within a vowel class (Kelley & Tucker 2020: 137).

84  The Spectral Overlap Assessment Metric (SOAM; Wassink 2006), which calculates overlap between

85  ellipses or ellipsoids fitted to the vowel distribution (see Wassink 2006 for details on the fitting) and

86  calculates the area or volume of their overlap, is one method adopted and recommended in other

87  sociophonetic work (Di Paolo, Yaeger-Dror & Wassink 2011: 103; Kendall & Fridland 2021: 56). We

88  refer interested readers to Nycz and Hall-Lew (2013) and Kelley and Tucker (2020) for in-depth

89  assessments of these measures and several others, compared with Pillai scores. Kelley and Tucker

90  assess four different metrics of vowel overlap, using a Monte Carlo simulation to test accuracy, and

91  find that Pillai scores produced the most accurate and precise values when compared to ground truth

92  values in their simulated data. They also recommend Pillai scores when sample sizes are "small," which

93  they define as 30 observations per group.

For sociophonetic data, and especially for naturalistic sociophonetic data, obtaining more than 30 observations per group is not always feasible, making Pillai an especially valuable tool for sociophonetics. As perhaps foregrounded by the title of this current paper, sample size plays an important role in the resulting Pillai score, which we highlight in detail in Sections 2–3. We note, however, that this is generally true of all measures of vowel overlap compared in the overview papers mentioned above (Nycz and Hall Lew 2013; Kelley and Tucker 2020). Indeed, any measure that uses means, standard deviations, and/or variance as inputs will be impacted by sample size. Likewise, larger sample sizes impact statistical significance, with larger sample sizes leading to smaller $p$-values. Our aim here is not to simply recommend that researchers obtain larger sample sizes, which in many cases is either not possible or may be at odds with other important data collection considerations, but rather to elucidate *how* sample size impacts resulting scores, so that researchers can be best informed when using Pillai as a measure of vowel overlap.

Perhaps an even bigger concern than total sample size with naturalistic data is the near impossibility of obtaining balanced token counts across two categories from naturalistic data. In wordlist data, researchers can carefully craft their wordlist to obtain balanced data: this typically means obtaining the same number of observations in each vowel class and ensuring that the wordlist is balanced for additional factors, such as phonological context. In naturalistic (i.e., conversational) speech, there is no way to ensure that a speaker produces balanced token counts across vowel category and across phonological contexts. As a result, researchers investigating something like vowel overlap need to understand precisely *how* Pillai scores may be affected by unbalanced data.

Finally, while Pillai has emerged as the best option so far for most sociophonetic data, there is one additional concern to address, which is that MANOVAs assume the data within each group follows a multivariate normal distribution (Bray & Maxwell 1985). Formant measurements for a given vowel from a given speaker in naturalistic vowel data, being nonnormally distributed, does not fit this

118  assumption (though see Whalen & Chen 2019 for some evidence that vowel formant data, even with

119  coarticulation, can be normally distributed). Recent work has suggested Bhattacharyya's Affinity[i] as a

120  better measure for non-normally distributed data (Fieberg & Kochanny 2005; Johnson 2015), and has

121  been taken up by some subsequent work (Strelluf 2018; Warren 2018; Jensen & Braber 2021). While

122  the benefits of Bhattacharyya's Affinity being a nonparametric measure make it particularly appealing

123  for the kind of distributions of data found in naturalistic speech, there is not yet a mechanism for

124  easily incorporating fixed effects like phonological context or speaker age.[ii] Furthermore,

125  Bhattacharyya's Affinity also works best on a relatively large data set of over 30 observations per

126  category (Seaman et al. 1999). While future work may make nonparametric methods like

127  Bhattacharyya's Affinity more easily integratable with the kind of statistical models linguists typically

128  use, for now we focus on Pillai score and the considerations necessary to make its use maximally

129  standardized.

130

131  ***3.  How Pillai scores are used to measure overlap in sociophonetics***

132  Pillai scores have been used to analyze a variety of phenomena in several languages. Perhaps

133  the most common application of Pillai scores is to measure overlap between /ɑ/ and /ɔ/ in North

134  American English (Hall-Lew 2013; Kendall & Fridland 2017; Havenhill 2015; Stanford et al. 2019

135  inter alia). In fact, using Pillai scores to measure this merger was explicitly recommended in Becker

136  (2019), a volume of different studies all analyzing the spread of the /ɑ-ɔ/ merger in North American

137  English. Many conditioned mergers in English have been quantified with Pillai scores as well (Schmidt,

138  Diskin-Holdaway & Loakes 2021; Austin 2020; Freeman 2021; Newbert 2021 inter alia). To a lesser

139  extent, Pillai scores have been used in other languages when analyzing vowels that are marginally

140  contrastive (Galician: Amengual & Chamorro 2015; Italian: Nadeu & Renwick 2016; Bangla: Islam &

141  Ahmed 2020; Hawai'ian: Kettig 2021; Swiss German: Joo, Schwarz & Page 2018; Austrian German:

142    Sloos 2013). Some more innovative uses of Pillai scores include using them to analyze tones in varieties

143    of Cantonese (Fung & Lee 2019; Tse 2018) and in Spanish fricative mergers (Regan 2020).

144    Pillai scores have also been used to quantify splits, chiefly among phonological low vowels.

145    Fisher et al. (2015) assess the degree to which the Philadelphia short-*a* split is found in their sample

146    of Philadelphians. Relatedly, Hall-Lew et al. (2017) look at the *bath-trap* split in Scottish Parliament

147    data. Brozovsky (2020) uses Pillai scores to measure the raising (and separation) of prenasal /æ/, using

148    Pillai scores to measure the overlap between prenasal /æ/ and preobstruent /æ/ in Taiwanese Texans.

149    In a study looking at the possible effect of salience on a given lexical item compared to the rest of its

150    canonical vowel class, Bray (2021) analyzes the lowered realization of the vowel in *hockey* compared to

151    other relatively more raised /ɑ/ tokens in professional American hockey players.

152    As highlighted by the selection of citations above, since being introduced to the field, Pillai

153    scores have become widespread in sociophonetic studies. With the support of meta-analyses that

154    compare other competing measures, Pillai has become a useful go-to for measuring both the overlap

155    and the distinction of speakers' pronunciation of two phonological categories, especially in

156    sociolinguistic studies that need to compare across individual speakers. While Pillai scores are clearly

157    a valuable tool in measuring vowel overlap, there remain some outstanding issues with using it. In the

158    following sections, we highlight some of these issues.

159

160    **B. Issues with Pillai scores**

161    *4. What is considered merged?*

162    As useful as Pillai scores are for quantifying the degree of overlap, they do not necessarily

163    answer researchers' underlying question of whether two vowels are merged. Pillai scores range from

164    0 to 1, but there is no agreed-upon cutoff value or threshold for determining whether the two groups

165    are merged or not. As a result, many studies rely on an ad-hoc threshold to interpret the merged status

166  of their speakers. Some work has suggested specific thresholds for mergers. For instance, Jibson

167  (2021) suggested a Pillai threshold of 0.3 as an indicator of "merged" status, after a shuffling procedure

168  identified 0.3 as the 95[th] percentile of "merged" between 20 tokens of two vowel classes from his

169  speakers. Wassink (2006) likewise suggests some provisional thresholds for SOAM, where 0-20%

170  overlap represents "distinct", 20-40% represents "partially merged", and >40% represents "merged".

171  Relying on provisional or ad-hoc thresholds, however, is risky because sample size matters and it is

172  likely not comparable across studies or even between speakers.

173        One solution for determining whether a given Pillai score should be interpreted as an

174  indication of "merger" is to examine the $p$-values that are associated with the MANOVA model from

175  which the Pillai scores are generated. The model assumes that the vowel variable contributes no

176  information to differentiating between two groups. In other words, it assumes the two vowels are

177  underlyingly merged. A small $p$-value associated with the vowel variable would provide evidence

178  against that null hypothesis, allowing the researcher to conclude that the difference between the two

179  groups is likely true (i.e., that the speaker is not merged). We note that Pillai scores and $p$-values are

180  inversely correlated: lower Pillai scores typically accompany higher $p$-values. In fact, since Pillai scores

181  are just test statistics, they and $p$-values are functions of each other. Nevertheless, $p$-values are not

182  typically reported in sociophonetic studies that use Pillai (some exceptions include Wong & Hall-Lew

183  2014; Nadeu & Renwick 2016; Amengual & Chamorro 2015; Berry 2018; Sloos 2013).

184        There are a few points of caution to make about using and interpreting $p$-values, as we discuss

185  throughout this paper. One of these concerns the potential distinction between a statistically

186  significant difference (as defined by the model) and a ground truth difference for speakers. For a

187  speech community that has a ground truth merger in two vowel categories, there will be no difference

188  in their perception of these two vowel categories. However, as sample size increases, so does the

189  likelihood of a model returning a $p$-value below a given significance threshold (typically 0.05); it is

190      possible that even for pairs of sounds that are truly merged, a sufficiently large dataset can interpret

191      random variance in the data as a meaningful difference, as shown in the experiments in Section 3. An

192      additional caution to make regarding interpreting $p$-values alone as an indicator of merger is that a

193      statistically significant difference in formant values may not map onto a perceptible difference for the

194      human auditory system (see, e.g., Kewley-Port and Watson, 1994). *P*-values alone can likewise be

195      misleading in the opposite direction: previous work has shown that speakers and listeners can produce

196      and perceive reliable but small differences, including sub-phonemic differences, in what may otherwise

197      appear to be merged sounds (for instance, with cases of incomplete neutralization; Warner et al. 2004;

198      Pfiffner 2021). Vowel distinctions that are maintained by small effect sizes, or by sub-phonemic

199      distinctions not captured by the measurements in the model, may appear artificially to be merged

200      according to a $p$-value because it takes more data for smaller differences to be detected by the statistical

201      model. For these reasons, additional information such as Pillai scores can aid in the interpretation of

202      $p$-values, and vice versa.

203      *5. Sample size*

204      As suggested by the parametric nature of Pillai, and the discussion of ad-hoc thresholds above,

205      a major component of deciding which threshold should be used to determine merger status is the

206      number of tokens being analyzed. Previous work on Pillai scores and sample size have expressed

207      concern over too-small sample sizes (Gorman & Johnson 2013), and over unbalanced sample sizes

208      across the two vowel classes being analyzed (Nycz & Hall-Lew 2013; Johnson 2015).

209      Despite sample size having a major impact on Pillai scores, most studies in sociophonetics

210      that use Pillai as a measure of vowel overlap do not also clearly report sample size (exceptions include

211      Wong & Hall-Lew 2014; Holland & Brandenburg 2017; Berry 2018; Berry & Ernestus 2018). This in

212      turn makes it difficult both to assess the findings in an individual paper and to compare speakers

213      within and across studies. In the simulation experiments below, we show just how important sample

214     size is to the resulting Pillai score, and provide a formula to calculate a recommended Pillai score

215     threshold for merger status, given a particular sample size.

216     **II.     METHODS**

217     In this section, we present the results of Monte Carlo simulation experiments designed to test

218     the effect of different sample sizes on resulting Pillai score. In these simulations, we create two vowel

219     classes that are perfectly merged underlyingly, and alter (1) the sample sizes between the two vowel

220     classes to test the effect of unbalanced samples across vowel classes and (2) the overall sample size,

221     considering both vowel classes together, to test the effect of unbalanced total sample size across

222     speakers.

223     **A. Data generation**

224     We generated data using a Monte Carlo simulation (Metropolis & Ulam 1949). This is a

225     procedure where random draws are taken from an underlying probability distribution or existing

226     dataset and analyzed. This process is repeated independently many times and the information about

227     each iteration is aggregated. To begin the simulation, a bivariate normal distribution was generated in

228     R to simulate a single theoretical underlying vowel in the F1-F2 dimension for a single theoretical

229     speaker. For the sake of simplicity, the mean for F1 and F2 were both set to zero and the standard

230     deviation was 1. The correlation coefficient between the formants was zero, producing a circular

231     (rather than elliptical) distribution. Because the generated data was uncorrelated, the distribution is

232     most easily generated by combining two independent univariate normal distributions because the

233     product of their probability densities is equal to their joint probability density (Johnson & Wichern

234     2012 chapter 4 page 151). Specifically, in R, we generated F1 by running `rnorm(x, mean = 0, sd`

235     `= 1)`, where `x` is the number of tokens generated. F2 was generated using the same code, and the two

236     sets of numbers were combined to create the bivariate normal distribution.

237        We acknowledge that this generated data deviates from real vowel data in several ways. First,

238    the units here were standardized rather than simulating formant frequencies, though this approach

239    can also be thought of as simulating a $z$-scored (Lobanov) normalized vowel space. Vowel data is

240    often not normally distributed (though, see Whalen & Chen 2019), so this distribution is not

241    necessarily representative of actual acoustic data. Even for vowel data that does appear be multivariate

242    normally distributed, there is often some degree of correlation in F1 and F2 (seen as elliptical

243    distributions in the F1-F2 space). However, since the focus of this study is on how sample size affects

244    Pillai, rather than another variable like skewness, we wanted to keep the distribution as simple as

245    possible. In the future we hope to replicate this study with an underlying distribution that is based on

246    formant measurements that reflect naturalistic speech.

247        To simulate vowel data, random draws were taken from that single bivariate normal

248    distribution, and assigned to one of two "vowel class" labels. We note that it is somewhat nonsensical

249    to refer to these generated numbers as "vowels", especially since Pillai scores can be calculated on

250    non-vowel data. However, since the majority of Pillai scores in sociophonetic analyses are based on

251    vowels, for clarity, we will refer to this simulated data points as "vowels" and their arbitrary groups as

252    "vowel classes." These random draws represent a linguist sampling data from our theoretical speaker.

253    For the sake of illustration, we will say 30 such observations were generated. These 30 observations

254    were treated as tokens from a single underlying vowel class.[iii] Another 30 random draws were then

255    taken from the same bivariate normal distribution. These 30 observations were treated as tokens from

256    a different underlying vowel class. Generating two groups from the same underlying distribution

257    therefore creates a simulated pair of merged vowels. In theory, the two simulated vowel classes should

258    not be statistically different from each other in any way because they were drawn from the same

259    underlying distribution.

260

261      **B. Two experiments**

262      For this study, we ran two experiments. In Experiment 1, the two simulated vowel classes for

263      each "speaker" were of equal size. We began with a sample size of 5 observations per vowel class. We

264      then moved on to 6 observations per vowel class, and so on, until we reached 100 observations per

265      vowel class. For each of these 95 sample sizes, we repeated the simulation 1000 times, each

266      representing a different instance of a linguist sampling data from that one underlyingly merged

267      speaker. This produced 96,000 pairs of simulated vowel data, where each pair consisted of equal-sized

268      vowel classes, enabling us to test the effect of overall sample size on resulting Pillai score.

269      In Experiment 2, we varied the sample size between the two vowel classes for each "speaker".

270      We began with 5 tokens from one vowel and 6 from another. We then took 5 tokens of one and 7

271      from the other. We increased the size of the second group by steps of 1 until it contained 100 tokens.

272      We then repeated this process with the first group having 6 tokens, and increased the second group

273      from 5 to 100 in steps of 1. We iterated over these steps, increasing the sample size of the first group

274      up to 100, thereby generating pairs of vowel data where every combination of sample sizes from 5 to

275      100 was represented. We repeated this simulation 100 times per combination. This produced 921,600

276      pairs of simulated vowel data, where each pair consisted of different-sized vowel classes, enabling us

277      to test the effect of unbalanced vowel class size on resulting Pillai score.

278      Across both experiments, Pillai scores were calculated for each pair of simulated vowel data.

279      Pillai scores were calculated by fitting a MANOVA model to the data using the `manova()` function

280      in R. The simulated F1-F2 measurements were the dependent variables and the vowel class was the

281      only independent variable. While it would be possible to incorporate additional simulated independent

282      variables such as place of articulation and duration into the MANOVA, we consider this to be beyond

283      the scope of the current paper, which focuses on the effect of sample size on resulting Pillai scores.

284      We therefore include only historical vowel class in our models, and leave more complex MANOVA

13

285    models to future work. The Pillai scores and *p*-values associated with the vowel class variable were

286    then extracted from that MANOVA model. To reiterate, the Pillai scores for all of these distribution

287    should be very close to zero (indicating complete overlap) because every data point was generated

288    from the same underlying bivariate normal distribution. Because the data is randomly generated, some

289    Pillai scores will be higher than others, but by rerunning the simulation many times per sample size,

290    we can begin to see patterns that may emerge at a given sample size.

291         The coding for this study was done in the R programming language (R Core Team 2021) with

292    the help of the `tidyverse` suite of packages (Wickham et al. 2019) and `joeyr` (Stanley 2021).

293    Visualizations were generated using `ggplot2` (Wickham 2015) and `see` (Lüdecke, Patil, et al. 2021)

294    with color palettes from `ggthemes` (Arnold 2018) and `scico` (Pedersen & Crameri 2020).


295    **III.    RESULTS**

296    **A. Experiment 1: Equal sample sizes**

297         To address how sample size affects Pillai scores, we first present the results from Experiment

298    1, where the two simulated vowel classes were the same size. Before inspecting the results of all sample

299    sizes though, it is important to understand how the 1000 Pillai scores were distributed within a given

300    sample size. Figure 1 shows two different views of the distribution of Pillai scores when the sample

301    size for both groups was 10. We see that the distribution of points representing resulting Pillai scores

302    is rather wide, a consequence of using such a small sample size for inferential statistics, ranging from

303    less than 0.001 to 0.568. Much of the data is clustered near the bottom of the plot but there is a long

304    "tail" extending upwards. This is not a haphazard pattern, but rather follows a distribution that can

305    be easily transformed into an *F* distribution and reflects the underlying mathematical properties of

306    how Pillai scores are calculated (cf Rencher & Christensen 2012: 182). For this particular sample size,

307    the mean Pillai score was 0.104, the median was 0.077, and the $95^{th}$ percentile was 0.294. As seen

308   below, these numbers change depending on the sample size, but the underlying distribution of the

309   Pillai scores is consistent across sample sizes. We show this distribution to dispel any misconceptions

310   that Pillai scores are uniformly distributed within a particular range, and to highlight that generally they
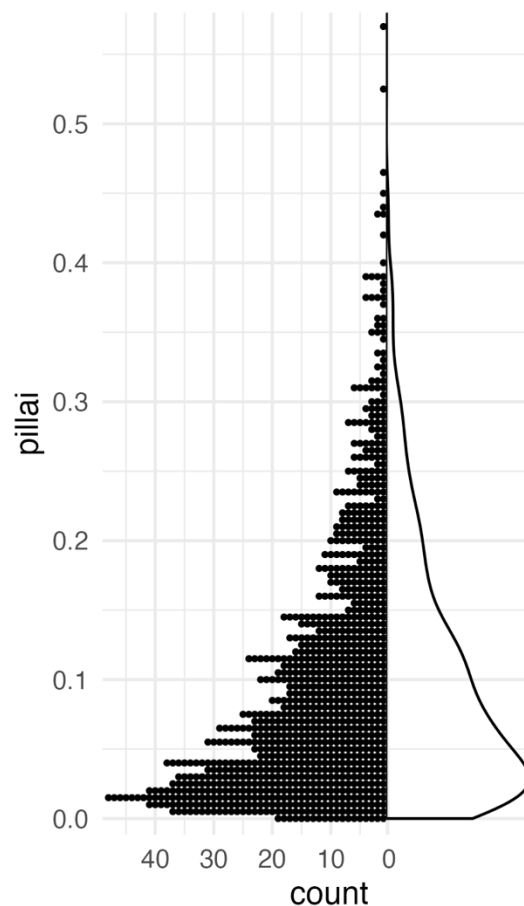
311   fall near the lower end of the distribution.

312



313

314        Figure 1: Distribution of Pillai scores on 1000 pairs of simulated groups, each with a size of

315   10.

316        With that distribution in mind, we can now zoom out to view all samples at once. Figure 2

317   shows all 96,000 Pillai scores by their sample size. Though present at all sample sizes, the "bottom-

318   heavy" distribution shown in in Figure 1 is not displayed in Figure 2, in order to make the general

319   trend across samples easier to see. It is immediately apparent that larger groups more consistently

320   produced lower Pillai scores. With very small sample sizes (groups of fewer than 10 observations per

321   vowel class), Pillai scores were quite high. For these small sample sizes, Pillai scores were sometimes

322   closer to 1 than 0, even for these underlyingly merged vowel classes which should, in principle, return

323   a Pillai score of 0. In other words, these small sample sizes sometimes resulted in very misleading Pillai

324   scores that may cause a researcher to interpret two vowel classes as distinct even if the true underlying

325   distribution was perfectly merged. As the sample size increases, Pillai scores were more consistently

326   low, as we would expect for underlyingly merged vowels.

327       The black line overlayed on Figure 2 indicates the 95$^{th}$ percentile for each sample size. This

328   line also very closely corresponds to the threshold for vowel class being statistically significant in the

329   MANOVA models: almost all points above that line had $p$-values less than 0.05 while almost all points

330   below it had greater $p$-values. Because we are modeling the null hypothesis (*i.e.* underlyingly merged

331   vowel classes), the distribution of $p$-values is uniform. It therefore is unsurprising, and in fact expected,

332   that the highest 5% of Pillai scores within a given sample size also return $p$-values less than 0.05. This

333   line shows that if there are just 10 observations per vowel class, 95% of the Pillai scores were under

334   0.3. However, as is evident in Figure 2, this threshold is only applicable to two groups of 10 since Pillai

335   scores decrease with larger samples. For example, the 95$^{th}$ percentile of returned Pillai scores does not

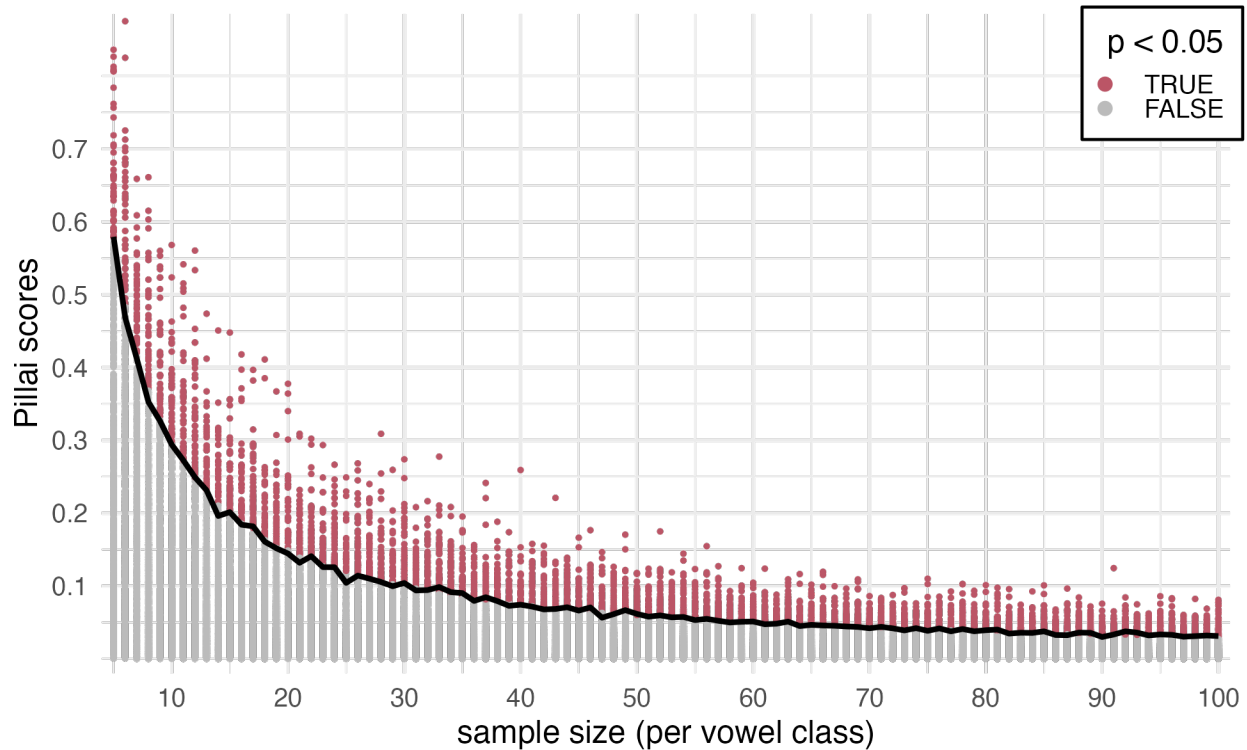336   drop to 0.1 until there are 30 observations per group.

Figure 2 (color online): Pillai scores for all simulations of equally-sized groups simulations, by sample size.

**B. Experiment 2: Unequal sample sizes**

While the previous section found that sample size affects Pillai scores in a predictable way, with larger samples returning more reliable Pillai scores, in this section we conduct further simulations to explore what effect, if any, an unbalanced sample has on Pillai scores. Unless data collection is carefully controlled to include a fixed number of tokens per vowel class, Pillai scores are run on vowel classes that are not comprised of the same number of tokens. Here we ask: what effect do grossly unbalanced groups have on Pillai scores?
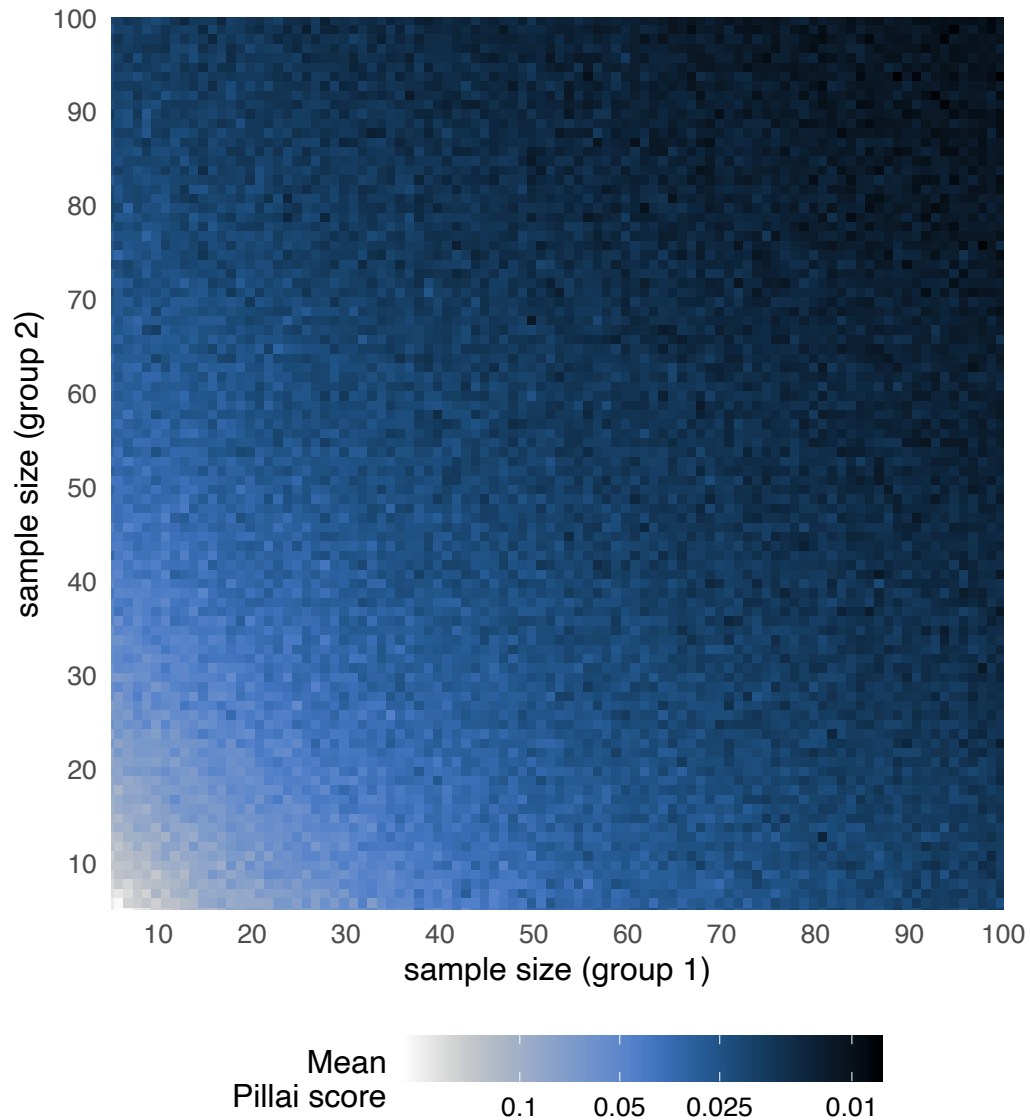
Figure 3 (color online): Mean Pillai scores for all simulations.

Figure 3 shows the mean Pillai scores from the 100 simulations for each combination of sample sizes between the two vowel classes. Going from the bottom left corner (two small sample sizes from each vowel class) to the top right corner (two large sample sizes from each vowel class), we see a general trend of decreasing Pillai scores as sample sizes increase. Reflecting Figure 2, we see these decreasing Pillai scores dropping more sharply as sample sizes are small (under around 30 tokens per vowel class).

357       A surprising pattern emerges when we look closely at the resulting Pillai scores from unequal

358       sample sizes: namely, that unequal samples across the two vowel classes do not impact Pillai score—

359       it is only the *total* sample size taking both vowel classes together that matters. For example, the mean

360       Pillai score for a pair of 10 tokens and 50 tokens drawn from this underlyingly merged distribution is

361       0.0358 and the mean Pillai score for a pair of 20 and 40 tokens is 0.0355. A pair of 30 tokens and 30

362       tokens drawn from an underlyingly merged distribution is nearly identical: 0.0381. We see this pattern

363       visually reflected along the diagonal between the top left corner and the bottom right corner of Figure

364       3, which shows a symmetrical resulting Pillai score for all unequal pairs of samples that sum to the

365       same total. In other words, we find that neither the existence of an unequal sample size between vowel

366       classes nor the degree of unequalness impact Pillai score.

367       These findings bring us to recommend simply using as many tokens as researchers have

368       available for an individual speaker, to bring the *total* sample size as high as possible. When comparing

369       across speakers, it is important to recall that total sample size impacts the resulting Pillai score, and

370       we recommend normalizing and reporting total sample sizes across speakers in a study to make

371       resulting Pillai scores maximally comparable.

372       **IV.     IMPLICATIONS**

373       **C. Choosing a threshold based on sample size**

374       As a result of these simulations, we are in a position to recommend a formula that researchers

375       can use as a guide to determine whether a given resulting Pillai score indicates a merger. We sought a

376       formula that was a function of sample size and could provide the 95th percentile for Pillai scores given

377       that two groups come from the same underlying distribution. In other words, we wanted to find the

378       function for the black line in Figure 2. We chose the 95th percentile for this distribution to align with

379       the common benchmark of $p < 0.05$.

380    The formula is based on the observation that by taking the natural log of the Pillai score and

381    the natural log of half of the total sample size of both groups (essentially log-transforming the axes in

382    Figure 2), the 95[th] percentile per sample size followed a straight line with an intercept of 1 and a slope

383    of –1. The formula therefore begins as *pillai* = 1 − *n*. But because this straight line only makes sense

384    when both *x* and *y* are log-transformed, it must be modified to be log(*pillai*) = 1 − log(*n*). At this point,

385    we solve for *pillai*, using e[x] as the antilog function, yielding $p_{95} = e^{1 - log\left(\frac{n}{2}\right)}$. Reducing this equation[iv]

386    produces the formula that we recommend for determining a Pillai score cut off,

387

388                                    $$p_{95} = \frac{2e}{n} \qquad\qquad (1)$$

389

390    where $p_{95}$ is the 95[th] percentile of Pillai scores given a total sample size *n*. The curve that is

391    generated by this formula very nearly follows the black line in Figure 2, which represents the 95[th]

392    percentile of Pillai scores given the *total* sample size (summing the samples from each group together).

393    To ensure the validity of the formula, we reran the entire simulation 100 more times and found that

394    in every case the formula is a very good approximation of the 95[th] percentile for Pillai scores. While

395    we do not have the mathematical expertise to prove whether Equation 1 represents the true underlying

396    distribution from which the 95[th] percentiles are drawn, we believe it is close enough to be fruitfully

397    used to determine a reasonable cutoff for "merged" status for a given total sample size.

398    Using Equation 1, we see that, for example, a sample of 20 total observations would produce

399    a Pillai score less than or equal to 0.2718 95% of the time if the two vowel classes were underlyingly

400    merged (we note that this is a close approximation to the cutoff of 0.3 that Jibson (2021) chose for

401    his total sample size of 20). For illustration, Table 1 presents threshold suggestions drawn from this

402    formula for different total sample sizes, highlighting that it takes a great deal of data to reliably return

403 Pillai scores that are close to zero (recall that Pillai scores closer to zero reflect a more "merged"

404 production). We recommend that researchers use Equation 1 in conjunction with *p*-values to make a

405 more informed decision about the merged status of two vowel classes rather than an ad hoc or

406 arbitrary cutoff.

407

408 Table 1: Pillai thresholds at various sample sizes based on Equation 1.

| Total sample size | Pillai threshold for "merged" |
|---|---|
| 20 | 0.2718 |
| 40 | 0.1359 |
| 60 | 0.0906 |
| 80 | 0.0680 |
| 100 | 0.0544 |
| 120 | 0.0453 |

409

410 As pointed out to us by a reviewer, one potential caveat for Equation 1 is that it is based on

411 uncorrelated data, as we mention in Section II.A above. Recall that the formants were generated

412 independently of each other and the correlation coefficient between them was zero, producing a

413 circular rather than elliptical distribution. Since vowel formant data is typically correlated, the utility

414 of this function is limited if it is not applicable to real data. To evaluate whether the correlation of the

415 vowel formants affects Pillai scores, we generated yet another set of simulated data. We followed the

416 methods described in Experiment 2 sample sizes ranging from 5 to 100 tokens per vowel, though for

417 the sake of creating a manageable amount of data, each vowel's sample size was limited to multiple of

418 5. This time though, we used the `mvrnorm()` function in the MASS package (Venables, Ripley &

419  Venables 2002) to generate correlated bivariate normal data, with correlation coefficients ranging from

420  0 to 0.9, in intervals of 0.1. We generated 100 vowel pairs per combination of sample sizes and

421  correlation coefficients, resulting in 361,000 new sets of simulated vowel data. We again calculated the

422  Pillai score for each vowel pair and then ran a linear regression model on these Pillai scores, with the

423  log-transformed Pillai score as the dependent variable and the log-transformed total sample size and

424  the correlation as predictors. We find that correlation was not a significant predictor in these 361,000

425  simulations, meaning there was no significant difference between the Pillai scores of uncorrelated data

426  and correlated data, given a particular total sample size. We also find that unbalanced data did not

427  affect Pillai scores, even in these correlated datasets. Thus we can be confident that Equation 1 is

428  applicable even to correlated, unbalanced data like what is typically found in real vowel formant

429  measurements.

430  **D. Sample size matters across speakers but not across vowels**

431  One important takeaway from these simulations is that it takes a relatively large amount of

432  data to reliably (meaning 95% of the time) return a low Pillai score such as 0.1 from two underlyingly

433  merged vowels. In an analysis of English /ɑ/ and /ɔ/, for example, conversational data can typically

434  provide a sufficient number of observations for a robust analysis of overlap. However, few studies

435  that analyze wordlist data contain many more than 10 tokens of these two vowels. And even within

436  long-form conversational data, if the research question focuses on an infrequent phonological variable,

437  the total sample size quickly drops. Because total sample size has a major impact on the resulting Pillai

438  score, we recommend researchers choose a relevant threshold for "merger" status, based on their total

439  sample size, and use it in conjunction with the resulting *p*-value to make a determination about merger

440  status for their data.

441  Perhaps the most surprising takeaway from these simulations, for both authors, was that

442  although Pillai is not a nonparametric test, it does not actually matter if the token counts across the

443     two categories being investigated are unequal. Instead, the most critical consideration is the total

444     number of tokens, summed across both categories. This is particularly important for naturalistic

445     sociolinguistic work, which relies on casual conversation rather than carefully constructed word lists

446     for data, meaning that is it often not possible to obtain balanced token counts across categories.

447     Following the results of Experiment 2, we can reassure researchers that unbalanced tokens across

448     vowel classes will not impact the resulting Pillai score. Instead, and following the results of Experiment

449     1, we recommend using as many total tokens possible for an analysis of a single speaker, regardless of

450     unbalanced samples across vowel classes for that speaker.

451         At the same time, any study aiming to compare the "merger" status of multiple individual

452     speakers should take into account the total sample sizes for each speaker, and especially consider the

453     fact that sample sizes (and therefore, the interpretation of resulting Pillai scores) may be different

454     across speakers. One of the primary goals for a robust measure of vowel overlap in sociolinguistics

455     has been to track the development of vowel mergers and splits across speakers in a given corpus (Nycz

456     & Hall-Lew 2013; Johnson 2010; Strelluf 2018; Labov et al. 2016) to analyze the trajectory of large-

457     scale language change over time. Because distributions with lower token counts produce inflated Pillai

458     scores, this means speakers who are less talkative will artificially appear to have more distinct vowels

459     than speakers who are very talkative. We recommend using one of two ways to account for unequal

460     sample sizes across speakers.

461         One option is for researchers to conduct an analysis of individual speakers, incorporating all

462     the relevant pieces of evidence (the recommended Pillai threshold given an individual speaker's sample

463     size, Pillai scores, $p$-values, and visualizations), to make a determination about the merged status of

464     each speaker. Section V.A presents an example of this method in action, where we demonstrate how

465     we leveraged all these pieces of evidence together to try to understand the merged status of each

466     speaker. This method allows researchers to obtain a fairly robust understanding of individuals as they

467 compare to each other and across styles. However, since this method requires researchers to synthesize

468 a number of gradient measures into discrete categories (such as "merged", "ummerged", and perhaps

469 "partially merged"), it makes it more difficult to track fine grained changes in the development of a

470 merger across a large speech community over time .

471       For researchers aiming to analyze fine grained differences over real or apparent time across

472 many speakers, however, it may be more beneficial to keep Pillai as a gradient measure. To analyze

473 Pillai across many speakers at once, we recommend a second approach to account for sample size,

474 which we present two options for. The basic goal for both options is to make Pillai scores comparable

475 across speakers. This can be done either by accounting for sample size in the model (see our example

476 in Section V.B.3), or by analyzing the same number of tokens per speaker (see our example in Section

477 V.B.4). Our recommendations for how best to account for sample size in the model and how best to

478 analyze the same number of tokens per speaker are discussed in detail in Section V.B.3-4, and an

479 example of the R code used to apply these recommendations is provided in the supplementary file.

480       A similar issue arises when comparing merger status from a single individual speaker but across

481 multiple speech styles. In particular, sociophonetic work often compares casual speech styles like

482 conversation to more formal speech styles, such as a minimal pair list or a reading task. Critically for

483 Pillai scores, since lower token counts will artificially appear more distinct, speech styles such as word

484 lists and reading tasks with lower token counts will likewise artificially appear more distinct than speech

485 styles with higher token counts (such as casual speech). In sum, there is a strong risk that a Pillai score

486 difference between speech styles may be interpreted as a stylistic difference (whereby speakers appear

487 to be maintaining a distinction in wordlists that they do not maintain in casual speech), when the effect

488 is entirely driven by differences in total sample sizes across the two styles. It's not uncommon, for

489 instance, to hear of speakers apparently undoing mergers in read speech in comparison to their casual

490 speech (cf. Labov 1994: 80; Berry 2018; Berry & Ernestus 2018); while Labov (1994) correctly points

491 out that apparent unmergers in careful speech may be speakers hypercorrecting in response to

492 orthographic differences across historically distinct vowel classes, we can add that distinctions

493 measured by Pillai score will additionally be impacted by sample size. We urge researchers who use

494 Pillai as a metric of overlap to pay close attention to differences in total sample sizes across speech

495 style,[v] and either control for total sample size across styles or adjust their threshold of "merger" in

496 each style accordingly, following the same recommendations provided above and elaborated on in

497 Section V.

498

499 **E. What to report when using Pillai scores**

500 Finally, we end with some general recommendations for what researchers should include in

501 their results when using Pillai scores as a measure of overlap. First, because total sample size impacts

502 Pillai score so strongly, we recommend always reporting total sample size per speaker (or per style, for

503 studies that compare merger across speech styles). This practice will have the added benefit of enabling

504 better comparisons across sociolinguistic studies, in turn enabling researchers to gain a clearer picture

505 of the large-scale spread of some ongoing mergers such as the /ɑ/-/ɔ/ merger.

506 Second, in addition to reporting and controlling for total sample sizes, we recommend

507 reporting $p$-values where possible. In other words, where studies report individual speakers' (or styles')

508 Pillai scores, those should always include the total sample size and the $p$-value as well. We are not

509 necessarily advocating for an increased reliance upon $p$-values as a binary meaningful threshold,

510 particularly since some statisticians are urging quantitative researchers to abandon declarations of

511 "statistical significance" and to instead understand $p$-values as one gradient measure in concert with

512 additional evidence (Wasserstein, Schirm & Lazar 2019). However, reporting a Pillai score without its

513 accompanying $p$-value is akin to reporting a regression estimate without a $p$-value, a $p$-value without

514 an effect size, or a mean without a standard deviation. A Pillai score is a test statistic and should be

515 reported as such. These two numbers in conjunction paint a better picture of the merged status of a

516 pair of vowels than either one in isolation.

517

518 **V.    A CASE STUDY**

519     The following is a case study to illustrate how one might conduct an analysis of vowel merger

520 using the recommendations in this paper. As illustrative data, we draw from sociolinguistic interviews

521 conducted and analyzed by the first author with residents of southwest Washington State, a region

522 where the low back vowels (/ɑ-ɔ/) are often merged, though with some indication of separation in

523 some speakers (see Stanley 2020 for more details). On average across the 52 participants analyzed

524 here, interviews lasted 46 minutes and yielded 143 tokens containing either /ɑ/ or /ɔ/ in preobstruent

525 position. After the interviews, 30 of those participants then read a wordlist containing another 20

526 tokens in a more careful style.

527

528 **A.  Analysis of individual speakers**

529     We performed a MANOVA on each speaker's F1 and F2 measurements, separately for each

530 style, with historic vowel class as the only predictor. Given the number of observations produced by

531 each speaker, Equation 1 was implemented to establish a potential cutoff value for each style. The *p*-

532 values and the Pillai scores were extracted from the MANOVA model and the latter were compared

533 to the cutoff values.

534     In this sample, there are some cases where the evidence overwhelmingly points towards a

535 merger. For example, 48-year-old Donna produced 179 low back vowels in the conversational portion

536 of her interview. With that many tokens, Equation 1 suggests that if her vowels were underlyingly

537 merged, her Pillai score should be less than 0.0304 about 95% of the time. In other words, anything

538    less than 0.0304 would be evidence for a merger. As it turns out, the MANOVA model performed on

539    her data yielded a Pillai score of 0.0289 with a $p$-value of 0.0756. The fact that her Pillai score is less

540    than the threshold for her token count and that the $p$-value above 0.05 means that historical vowel

541    class category does not predict Donna's acoustic realization, suggesting that Donna's two vowel

542    classes are likely underlyingly merged. In the wordlist data, there were only 20 tokens, so the suggested

543    cutoff value determined by Equation 1 is much higher at 0.2718 because of the smaller sample size.

544    The MANOVA performed on Donna's wordlist data produced a Pillai score of 0.1648 ($p$ = 0.216).

545    We reiterate that sample sizes also impact $p$-values, with smaller sample sizes producing higher $p$-

546    values. Taken all together, this data makes a strong case that /ɑ/ and /ɔ/ are underlyingly merged in

547    Donna's speech in both speech styles.
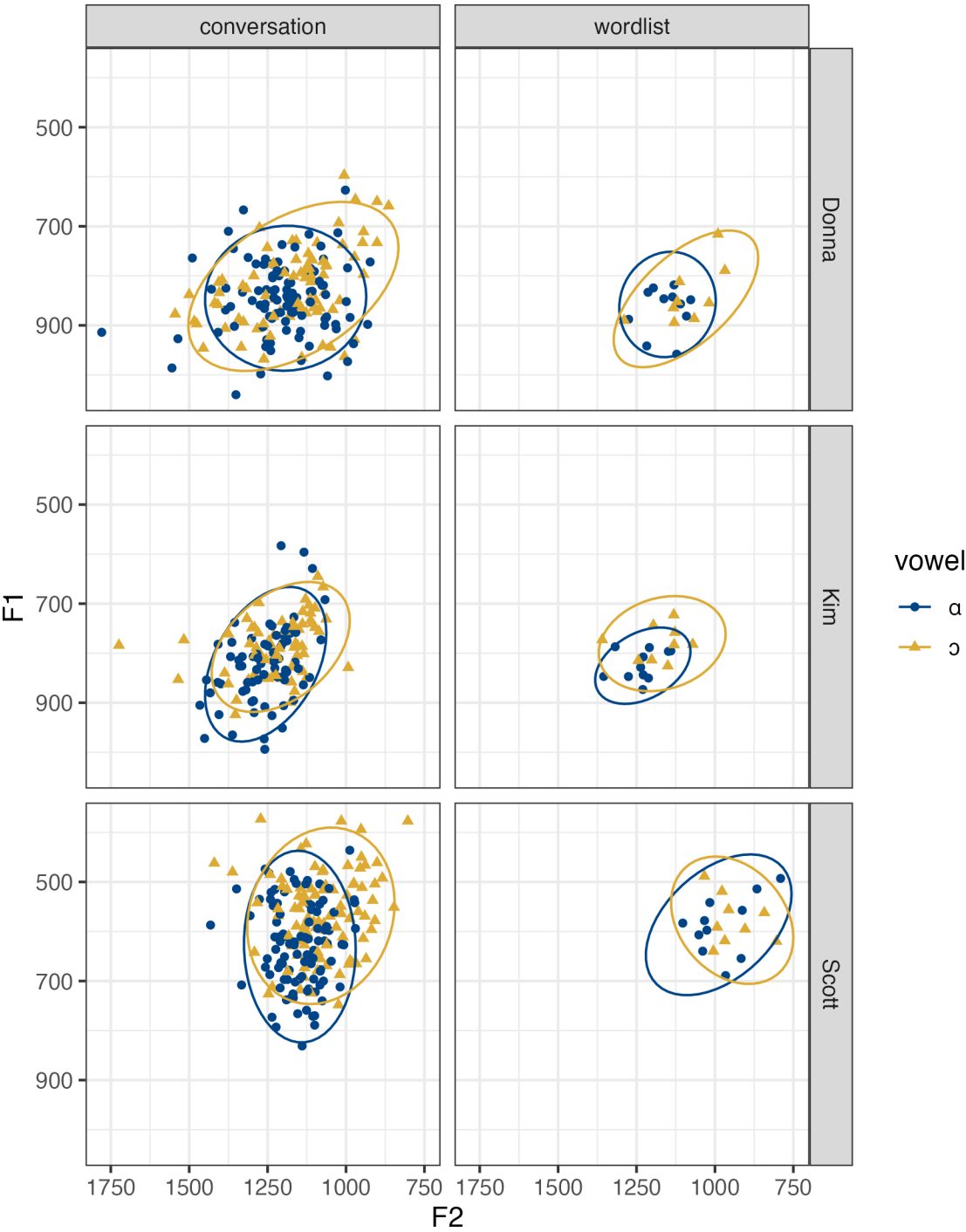
548         On the other hand, some speakers' data are indicative of a distinction. Kim, another 48-year-

549    old woman, produced 137 low back tokens in the conversational portion of her interview. Equation

550    1 suggests that a Pillai score less than 0.0397 would be evidence for a merger, but the MANOVA

551    performed on her data returned a Pillai score of 0.0658 ($p$ = 0.010). Though the Pillai score is relatively

552    close to zero, we do not interpret her data as underlying merged, since with that many observations

553    from a truly merged distribution we would expect an even lower Pillai score (less than 0.0397). Based

554    on the 20 tokens from her wordlist, the cutoff would be 0.2718 but the MANOVA on those 20 tokens

555    yielded a Pillai score of 0.3700 ($p$ = 0.020), further suggesting a distinction. Because Kim's Pillai scores

556    were higher than the thresholds and were accompanied by low $p$-values for both speech styles, we

557    conclude that Kim's low back vowels, while close in acoustic space, are not fully merged.

558         However, even when considering $p$-values alongside recommended thresholds, not all cases

559    are as straightforwardly interpretable as Donna's and Kim's. Scott is a 28-year-old man whose

560    interview contained 195 low back tokens. The Pillai score based on his data was 0.1975, far higher

561    than the threshold (0.0279) produced by Equation 1, and was accompanied by a low $p$-value ($p$ <

562  0.001), suggesting two distinct underlying vowel distributions. However, the Pillai score based on the

563  20 tokens he produced in his wordlist (0.0397) was much lower than the threshold (again, 0.2718 ),

564  and had a high *p*-value (*p* = 0.7090). When the Pillai score is lower than the threshold and is

565  accompanied by a high *p*-value his wordlist data, it is tempting to interpret Scott as being a rare case

566  of producing a merger in the wordlist that he does not produce in the conversational portion of the

567  interview. However, because the sample size is so small in the wordlist, the Pillai score (and indeed,

568  any result of an inferential statistical test) should be taken with a large grain of salt. We include Scott's

569  data here in part to demonstrate that even when leveraging Pillai scores alongside a recommended

570  threshold and a *p*-value, data with low token counts can still be difficult to interpret.

571       We can add one additional tool to the suite of evidence we consider when diagnosing

572  individual speakers: a visual inspection of a plot. Figure 4 shows the distributions of /ɑ/ and /ɔ/

573  tokens in F1-F2 space for conversational data (left) and wordlist data (right) for both Donna (top),

574  Kim (middle), and Scott (bottom). Ellipses represent one standard deviation for each vowel class. A

575  visual inspection of these plots shows that both vowel classes exhibit a fair amount of overlap in both

576  styles, with what appears to be more overlap in Donna's, as her /ɔ/ vowel class actually encompasses

577  /ɑ/ in both speech styles (a distributional property that indicates merger). Adding the measures of

578  Pillai scores, using the recommended thresholds for "merged" given the specific token counts, along

579  with *p*-values for these speakers in these two styles, enables us to more confidently state that Donna's

580  two vowels are underlyingly merged, while Kim's are distinct in both styles and Scott's are distinct at

581  least in the conversational style. For both Donna and Kim, the Pillai scores were higher in the wordlist

582  style compared to the conversation. Seeing these differences in Pillai scores alone, a researcher may

583  be tempted to conclude that there is style shifting occurring for both speakers, such that the merger

584  undoes itself in more careful speech. Leveraging all of the evidence together – Pillai scores alongside

585  the recommended threshold for a given sample size, as well as *p*-values and a visual inspection of the

586     data – allows us to reject this interpretation and instead see important differences between speakers

587     in the sample.



588

589     Figure 4 (color online): Low back vowels from three speakers across two styles.

590

591

592    We note, in fact, that interpreting Pillai score alone without the additional evidence of

593    threshold and $p$-value provides a misleading interpretation of the entire dataset. Across the 30 speakers

594    in the sample who completed both tasks, a one-sided paired $t$-test comparing Pillai scores in

595    conversational data to Pillai scores in wordlist data suggests that there is a statistically significant trend

596    towards higher Pillai scores (i.e., a "less merged" pronunciation) in the wordlist data ($t$ = -3.3296, df

597    = 29, $p$ = 0.001). This pattern obtains across most speakers, suggesting on the surface that almost

598    every speaker in the sample "unmerges" their vowels in wordlist style data or that they only merge the

599    vowels in casual conversation. Given that the low-back merger is a change in progress (cf. Labov,

600    Yaeger & Steiner 1972 for other examples), this pattern obtaining across this many speakers of all ages

601    is suspiciously regular – much more regular than we would expect given patterns of variation and

602    change in sociophonetic work. In fact, it was this apparently regular unmerging found in these

603    participants' wordlists that led us to investigate the effect of sample size in the first place.

604    Incorporating all of the relevant pieces of evidence (recommended threshold given sample size, along

605    with Pillai scores and $p$-values) allows us to understand the suspiciously regular finding as an artefact

606    of wordlists having far smaller sample sizes than conversational speech. Likewise, leveraging all of the

607    evidence, including sample size differences across styles, allows a clearer understanding of the

608    individual speakers in the sample and whether they are likely to be truly merged or not.

609

610    **B. Analysis of many speakers**

611    While the level of scrutiny in the previous section is appropriate and encouraged for analyzing

612    individual speakers, we acknowledge that researchers may have a different goal in mind. For instance,

613    researchers tracking the development of a merger as it becomes closer together in phonetic space (and

614    before a categorical merging) will want to understand how the two vowel classes become less distinct

615    across generations – even before a categorical merger has taken place.  Analyzing data speaker-by-

616    speaker requires us to take gradient measures (Pillai scores, threshold, and *p*-value) and interpret them

617    into discrete categories for each speaker and style ("merged", "not merged" or "partially merged",

618    depending on the researcher's interpretation); this discrete interpretation in turn makes it difficult to

619    track how large-scale change proceeds across generations. In this section we explore how Pillai scores

620    have changed over time in this sample by fitting linear models to the data. The first and more

621    straightforward approach simply accounts for sample size in the modeling; the second and more

622    robust method implements a bootstrapping procedure to remove difference between sample sizes

623    across speakers.

624

625        ***6.   Incorporating sample size in a regression model***

626            We next wish to identify whether there were patterns across genders or time, so we would like

627    to fit the Pillai scores to a linear model. However, since Pillai scores are typically non-normally

628    distributed, the results from a regression model fit to them may be unreliable. However, as mentioned

629    above, Pillai scores and *p*-values are functions of each other, and as seen in Figure 2, if the Pillai score

630    for a given sample size is higher than the 95[th] percentile threshold, the *p*-value will be less than 0.05.

631    Pillai scores and *p*-values essentially say the same thing. We can therefore fit a model to the *p*-values

632    instead of Pillai scores without loss of interpretation. As mentioned above, *p*-values are uniformly

633    distributed under the null hypothesis of underlyingly merged vowel classes, which is what we model

634    here. And a uniform distribution can be easily transformed into a normal distribution using the

635    `qnorm()` function in R. We can therefore can fit a linear regression model to the transformed *p*-values

636    to study change over time. In this case, modeling on normalized *p*-values resulted in a better model fit

637    than modeling raw *p*-values or on raw or transformed Pillai scores.[vi] As independent variables, we

638    included gender, birth year (scaled and centered around zero using `scale()`), and their interaction.

639    Since each speaker's Pillai scores were calculated on different amounts of data, we incorporated

640    speaker sample size into the model as well, after transforming it using Equation 1. In R, the formula

641    for this model was

642

643        ```
lm(qnorm_p ~ gender * scale(yob) + equation1(n), data = df)
```

644

645        where `qnorm_p` represents the *p*-values normalized with the `qnorm()` function, `scale(yob)`

646    is the scaled birth years, `equation1(n)` is the application of Equation 1 on each speaker's sample

647    size (`n`), and `df` is the data frame of containing one row for each of the 52 speakers, and columns for

648    metadata, sample size, and *p*-values. We direct readers to the supplementary file for an example of the

649    R code used to implement this model.

650        As seen in Figure 5, the model suggests a significant interaction between gender and birth year.

651    Men's *p*-values were relatively low compared to women's, which correspond to higher Pillai scores

652    and greater separation between /ɑ/ and /ɔ/. We interpret this as suggesting that women are, overall,

653    more merged than men in this sample – in other words, women are leading the change towards a

654    merger in Washington. In fact, 19 of the 24 men had Pillai scores above their respective thresholds,

655    suggesting a distinction in their low vowels, while 18 of the 28 women had Pillai scores lower than

656    their thresholds, suggesting a merger. The significantly positive interaction term suggests difference

657    between genders shrinks among younger people.

658

659

Figure 5 (color online): Normalized *p*-values by gender and birth year with predicted regression
lines. The horizontal dotted line crosses the *y*-axis at -1.67, corresponding to *p*-values of 0.05; speakers
below that line have *p*-values less than 0.05 and Pillai scores higher than the threshold, given their
sample size.

We note that threshold calculated using Equation 1, which serves as a proxy for sample size,
was a significant predictor in this model ($p = 0.0485$). It's positive coefficient suggests that as the
threshold increases (corresponding to *smaller* sample sizes), the normalized *p*-value increases
(corresponding to an increase in raw *p*-values as well). We interpret this to mean that, even after
incorporating sample size in the model, people with less data had larger Pillai scores. For datasets that
include a wider range of sample sizes, particularly where each person contributed fewer tokens than
those in dataset analyzed here, we would expect to find a larger effect of sample size on *p*-values, and
consequently Pillai scores.

674    *7.  Downsampling and bootstrapping*

675        As mentioned in Section IV, the other method that we recommend for comparing Pillai scores

676    across a large number of speakers is to reduce the size of each speaker's dataset down to the sample

677    size of the speaker who contributed the least amount of data. Downsampling in this way will allow us

678    to obtain Pillai scores that are comparable across all speakers, and therefore allow us to track fine-

679    grained changes in merger status across apparent time.

680        To begin our downsampling example, we first restrict our analysis to the more prolific

681    conversational portion of the interview, to maximize the possible total sample size per speaker. The

682    least talkative speaker in this sample produced only 82 tokens in this style. So, even though some

683    speakers produced many more (as many as 327 in one case), we took a random sample (with

684    replacement[vii]) of just 82 tokens from each speaker. However, we found that Pillai scores from a

685    random sample of one speaker's data varied considerably from a different random sample from that

686    same speaker. So, we implemented a bootstrapping procedure and took 1000 random samples (with

687    replacement) of 82 tokens from each speaker's dataset, calculated the Pillai scores and other summary

688    statistics from each sample, and aggregated them by speaker. We find that these aggregated values

689    were very similar to values taken from the full dataset: the correlation between speakers' Pillai scores

690    on the full dataset and speakers' mean Pillai scores across the 1000 samples was 0.9984, suggesting

691    that the aggregated bootstrapped values are a very good approximation of the full dataset's values.

692    The difference is that they are based on equal sample sizes rather than different sample sizes, which

693    means that the resulting Pillai scores will be comparable across speakers.

694        To identify whether there were patterns across genders or time, we fit the normalized *p*-values

695    to a linear model with gender and birth year as predictors. Since the data was downsampled, sample

696    size does not need to be incorporated into the model; Pillai scores can be directly compared and the

697    influence of particularly talkative or reticent speakers is minimized. We compared two models, one

698 where the dependent variable was the raw $p$-value and the other where the dependent variable was

699 normalized (again using the `qnorm()` function in R). In this case, the model on the normalized $p$-

700 values was a better fit, based on the adjusted $R^2$ and model assumption checks (see also footnote 5).

701 Because the values from this bootstrapped method were similar to the values obtained from the full

702 dataset, it is unsurprising that the model suggests a similar trend: men's Pillai scores start higher and

703 decrease across time while women's Pillai scores are lower and increase very slightly over time, with

704 the two converging among the youngest speakers.

705 　　In this dataset, the two approaches to comparing Pillai scores across speakers (incorporating

706 sample size in the model and downsampling with bootstrapping) yielded very similar results. We again

707 attribute this to the fact that the sample sizes were all somewhat large and relatively homogeneous.

708 For datasets where the number of observations per speaker is roughly similar, the downsampling and

709 bootstrapping approach may be more appropriate since it will result in little loss of data. Similarly, for

710 datasets where the total sample size per speaker is large, the downsampling and bootstrapping

711 approach may also be more appropriate since downsampling to a large sample size from a very large

712 sample size will result in little change. The difference between 30 and 60 total tokens has a larger

713 impact on resulting Pillai score than the difference between 300 and 600. On the other hand, for

714 samples with more variation in how many tokens were extracted from each speaker, the approach that

715 incorporates sample size into the model may be more appropriate since the effect of sample size will

716 be more apparent. And for samples where the number of tokens is small, incorporating the sample

717 size into the model may also be more appropriate so that none of the already scant data is lost. In the

718 end, we recommend researchers use the approach that makes the most sense for their data.

719

## VI.    CONCLUSION

720

721      In this paper, we present a close view into the effect of sample size on resulting Pillai scores,

722  a common measure for quantifying vowel overlap. We use a series of simulation experiments drawing

723  from an underlyingly merged pair of vowels to demonstrate (1) that larger sample sizes yield reliably

724  lower Pillai scores, (2) that unequal group sizes across the two vowel classes is irrelevant in the

725  calculation of Pillai scores, and (3) that it takes more data than many sociolinguistic studies collect to

726  reliably return a low Pillai score (e.g., under 0.1) even for underlyingly merged data. These results have

727  implications for how Pillai scores are compared across studies and between speakers or speech styles

728  within the same study. We provide some recommendations for maximizing reliability in the use of

729  Pillai scores, and provide a formula to assist researchers in determining a reasonable Pillai score

730  threshold to use as an indicator of merged status given their sample size. We recommend the use of

731  Equation 1 a in conjunction with the Pillai scores' accompanying *p*-values to make informed decisions

732  about the merged status of two vowels in a given speaker. By properly using and reporting aspects of

733  Pillai score, researchers can come closer to accurately quantifying vowel overlap, identifying vowel

734  mergers, and ultimately understanding broad patterns of variation and change in vowel merger.

735

745

746       See supplementary material at [URL will be inserted by AIP] for a brief tutorial on how to

747 implement these recommendations in R.

748

749

750       **REFERENCES**

751 Amengual, Mark & Pilar Chamorro. 2015. The Effects of Language Dominance in the Perception
752      and Production of the Galician Mid Vowel Contrasts. *Phonetica* 72(4). 207–236.
753      https://doi.org/10.1159/000439406.
754 Arnold, Jeffrey B. 2018. ggthemes: Extra Themes, Scales and Geoms for "ggplot2."
755      https://CRAN.R-project.org/package=ggthemes.
756 Austin, Martha. 2020. Mismatches between Linguistic and Sociolinguistic Perception. Presented
757      at the The 94th Annual Meeting of the Linguistic Society of America, New Orleans, LA.
758 Bartlett, M. S. 1939. A note on tests of significance in multivariate analysis. *Mathematical*
759      *Proceedings of the Cambridge Philosophical Society* 35(2). 180–185.
760      https://doi.org/10.1017/S0305004100020880.
761 Becker, Kara. 2019. Introduction. In Kara Becker (ed.), *The Low-Back-Merger Shift: Uniting the*
762      *Canadian Vowel Shift, the California Vowel Shift, and short front vowel shifts across*
763      *North America* (Publication of the American Dialect Society 104). Durham, NC: Duke
764      University Press.
765 Berry, Grant M & Mirjam Ernestus. 2018. Phonetic alignment in English as a *lingua franca* :
766      Coming together while splitting apart. *Second Language Research* 34(3). 343–370.
767      https://doi.org/10.1177/0267658317737348.
768 Berry, Grant Michael. 2018. *Liminal voices, central constraints: Minority adoption of majority*
769      *sound change*. The Pennsylvania State University Ph.D. Dissertation.
770 Bray, Andrew. 2021. [ˈhɒki]: An Emerging Third-Order Index of a Hockey-Based Persona.
771      Presented at the New Ways of Analyzing Variation 49, Austin, Texas.
772 Bray, James H. & Scott E. Maxwell. 1985. *Multivariate Analysis of Variance* (Quantitative
773      Applications in the Social Sciences). Vol. 07–054. Newbury Park, CA: Sage Publications.
774 Brozovsky, Erica Sharon. 2020. *Taiwanese Texans: A Sociolinguistic Study of Language and*
775      *Cultural Identity*. Austin, TX: University of Texas at Austin Ph.D. Dissertation.
776 Di Paolo, Marianna & Alice Faber. 1990. Phonation differences and the phonetic content of the
777      tense-lax contrast in Utah English. *Language Variation and Change* 2(02). 155–204.
778      https://doi.org/10.1017/S0954394500000326.
779 Di Paolo, Marianna, Malcah Yaeger-Dror & Alicia Beckford Wassink. 2011. Analyzing Vowels. In
780      Marianna Di Paolo & Malcah Yaeger-Dror (eds.), *Sociophonetics: A Student's Guide*, 87–
781      106. 1st edn. London: Routledge.
782 Fieberg, John & Christopher O. Kochanny. 2005. Quantifying Home-Range Overlap: The
783      Importance of the Utilization Distribution. (Ed.) Lanham. *Journal of Wildlife*

784    *Management* 69(4). 1346–1359. https://doi.org/10.2193/0022-
785    541X(2005)69[1346:QHOTIO]2.0.CO;2.

786  Fisher, Sabriya, Hilary Prichard & Betsy Sneller. 2015. The Apple Doesn't Fall Far From the Tree:
787    Incremental Change in Philadelphia Families. *University of Pennsylvania Working Papers*
788    *in Linguistics* 21(2). http://repository.upenn.edu/pwpl/vol21/iss2/7.

789  Freeman, Valerie. 2021. Vague eggs and tags: Prevelar merger in Seattle. *Language Variation*
790    *and change* 1–24. https://doi.org/doi:10.1017/S0954394521000028.

791  Fung, Roxana S. Y. & Chris K. C. Lee. 2019. Tone mergers in Hong Kong Cantonese: An
792    asymmetry of production and perception. *The Journal of the Acoustical Society of*
793    *America* 146(5). EL424–EL430. https://doi.org/10.1121/1.5133661.

794  Gorman, Kyle & Daniel Ezra Johnson. 2013. Quantitative Analysis. In Robert Bayley, Richard
795    Cameron & Ceil Lucas (eds.), *The Oxford Handbook of Sociolinguistics*, vol. 1. Oxford
796    University Press. https://doi.org/10.1093/oxfordhb/9780199744084.013.0011.

797  Hall-Lew, Lauren. 2013. 'Flip-flop' and mergers-in-progress. *English Language and Linguistics*
798    17(02). 359–390.

799  Hall-Lew, Lauren, Ruth Friskney & James M. Scobbie. 2017. Accommodation or political identity:
800    Scottish members of the UK Parliament. *Language Variation and Change* 29(3). 341–
801    363. https://doi.org/10.1017/S0954394517000175.

802  Han, Jeong-Im & Hyunsook Kang. 2013. Cross-generational Change of /o/ and /u/ in Seoul
803    Korean I: Proximity in Vowel Space. *Phonetics and Speech Sciences* 5(2). 25–31.
804    https://doi.org/10.13064/KSSS.2013.5.2.025.

805  Havenhill, Jonathan. 2015. Maintenance of the COT-CAUGHT Contrast Among Metro Detroit
806    Speakers: A Multimodal Articulatory Analysis. *University of Pennsylvania Working Papers*
807    *in Linguistics* 21(2).

808  Hay, Jennifer, Paul Warren & Katie Drager. 2006. Factors influencing speech perception in the
809    context of a merger-in-progress. *Journal of Phonetics* (Modelling Sociophonetic
810    Variation) 34(4). 458–484. https://doi.org/10.1016/j.wocn.2005.10.001.

811  Holland, Cory & Tara Brandenburg. 2017. Beyond the Front Range: The Coloradan Vowel Space.
812    In Valerie Fridland, Alicia Beckford Wassink, Tyler Kendall & Besty E. Evans (eds.), *Speech*
813    *in the Western States, Volume 2: The Mountain West* (Publication of the American
814    Dialect Society 102), 9–30. Durham, NC: Duke University Press. DOI: 10.1215/00031283-
815    4295277.

816  Islam, Md Jahurul & Iftakhar Ahmed. 2020. Mid-front and back vowel mergers in Mymensingh
817    Bangla: An acoustic investigation. *Linguistics Journal* 14(1). 206–232.

818  Ismay, Chester & Albert Y. Kim. 2020. *Statistical Inference via Data Science: A ModernDive into*
819    *R and the Tidyverse* (The R Series). Boca Raton, FL: Taylor & Francis Group.
820    https://moderndive.com/index.html.

821  Jensen, Sandra & Natalie Braber. 2021. The BATH-TRAP split in the East Midlands. Poster
822    presented at the New Ways of Analyzing Variation 49, Austin, Texas.

823  Jibson, Jonathan. 2021. Merged status thresholds for Pillai scores. Poster presented at the New
824    Ways of Analyzing Variation 49, Austin, Texas.

825  Johnson, Daniel Ezra. 2010. *Stability and Change Along a Dialect Boundary: The Low Vowels of*
826    *Southeastern New England* (Publication of the American Dialect Society 95). Durham,
827    NC: Duke University Press.

828 Johnson, Daniel Ezra. 2015. Quantifying vowel overlap with Bhattacharyya's affinity. Presented
829     at the New Ways of Analyzing Variation (NWAV44), Toronto.
830 Johnson, Richard A. & Dean W. Wichern. 2012. *Applied Multivariate Statistical Analysis*. Phi
831     Learning Private Limited.
832 Joo, Hyoun-A, Lara Schwarz & B. Richard Page. 2018. Nonconvergence and Divergence in
833     Bilingual Phonological and Phonetic Systems: Low Back Vowels in Moundridge
834     Schweitzer German and English. *Journal of Language Contact* 11(2). 304–323.
835     https://doi.org/10.1163/19552629-01102006.
836 Kelley, Matthew C. & Benjamin V. Tucker. 2020. A comparison of four vowel overlap measures.
837     *The Journal of the Acoustical Society of America* 147(1). 137–145.
838     https://doi.org/10.1121/10.0000494.
839 Kendall, Tyler & Valerie Fridland. 2017. Regional relationships among the low vowels of U.S.
840     English: Evidence from production and perception. *Language Variation and Change*
841     29(2). 245–271. https://doi.org/10.1017/S0954394517000084.
842 Kendall, Tyler & Valerie Fridland. 2021. *Sociophonetics* (Key Topics in Sociolinguistics).
843     Cambridge: Cambridge University Press.
844 Kettig, Thomas T. 2021. *Ha'ina 'ia Mai Ana Ka Puana: The Vowels of 'Ōlelo Hawai'i*. Mānoa,
845     Hawai'i: University of Hawai'i at Mānoa Ph.D. Dissertation.
846 Labov, William. 1994. *Principles of linguistic change. Vol. 1: Internal features* (Language in
847     Society). Oxford: Wiley-Blackwell.
848 Labov, William & Maciej Baranowski. 2006. 50 msec. *Language Variation and Change* 18(03).
849     https://doi.org/10.1017/S095439450606011X.
850 Labov, William, Sabriya Fisher, Duna Gylfadottír, Anita Henderson & Betsy Sneller. 2016.
851     Competing systems in Philadelphia phonology. *Language Variation and Change* 28(3).
852     273–305. https://doi.org/10.1017/S0954394516000132.
853 Labov, William, Malcah Yaeger & Richard Steiner. 1972. *A quantitative study of sound change in
854     progress: Volume 1*. Philadelphia, PA: US Regional Survey.
855 Lüdecke, Daniel, Mattan Ben-Shachar, Indrajeet Patil, Philip Waggoner & Dominique Makowski.
856     2021. performance: An R Package for Assessment, Comparison and Testing of Statistical
857     Models. *Journal of Open Source Software* 6(60). 3139.
858     https://doi.org/10.21105/joss.03139.
859 Lüdecke, Daniel, Indrajeet Patil, Mattan Ben-Shachar, Brenton Wiernik, Philip Waggoner &
860     Dominique Makowski. 2021. see: An R Package for Visualizing Statistical Models. *Journal
861     of Open Source Software* 6(64). 3393. https://doi.org/10.21105/joss.03393.
862 Metropolis, Nicholas & S. Ulam. 1949. The Monte Carlo Method. *Journal of the American
863     Statistical Association* 44(247). 335–341.
864     https://doi.org/10.1080/01621459.1949.10483310.
865 Moulton, William G. 1968. Structural Dialectology. *Language* 44(3). 17.
866 Nadeu, Marianna & Margaret E.L. Renwick. 2016. Variation in the lexical distribution and
867     implementation of phonetically similar phonemes in Catalan. *Journal of Phonetics* 58.
868     22–47. https://doi.org/10.1016/j.wocn.2016.05.003.
869 Newbert, Cornelia. 2021. *Language variation in South Africa: A sociophonetic study of the vowel
870     system of Black South African English*. Regensburg, Germany: University of Regensburg
871     Ph.D. Dissertation.

872  Nycz, Jennifer & Lauren Hall-Lew. 2013. Best practices in measuring vowel merger. *Proceedings*
873      *of Meetings on Acoustics* 20(1). 060008. https://doi.org/10.1121/1.4894063.
874  Olson, Chester L. 1976. On choosing a test statistic in multivariate analysis of variance.
875      *Psychological Bulletin* 83(4). 579–586. https://doi.org/10.1037/0033-2909.83.4.579.
876  Pedersen, Thomas Lin & Fabio Crameri. 2020. scico: Colour Palettes Based on the Scientific
877      Colour-Maps. https://CRAN.R-project.org/package=scico.
878  Pfiffner, Alexandra M. 2021. *Cue-Based Features: Modeling change and varitaion in the voicing*
879      *contrasts of Minnesotan English, Afrikaans, and Dutch*. Washington, DC: Georgetown
880      University Ph.D. Dissertation.
881  Pillai, K. C. S. 1955. Some New Test Criteria in Multivariate Analysis. *The Annals of Mathematical*
882      *Statistics* 26(1). 117–121. https://doi.org/doi:10.1214/aoms/1177728599.
883  R Core Team. 2021. R: A Language and Environment for Statistical Computing. Vienna, Austria:
884      R Foundation for Statistical Computing. http://www.R-project.org.
885  Regan, Brendan. 2020. Extending Pillai Scores to Fricative Mergers: Advancing a Gradient
886      Analysis of a Split-in-Progress in Andalusian Spanish. *University of Pennsylvania Working*
887      *Papers in Linguistics* 26(2).
888  Rencher, Alvin C. & William F. Christensen. 2012. *Methods of multivariate analysis* (Wiley Series
889      in Probability and Statistics). Third Edition. Hoboken, New Jersey: Wiley.
890  Schmidt, Penelope, Chloé Diskin-Holdaway & Debbie Loakes. 2021. New insights into /el/-/æl/
891      merging in Australian English. *Australian Journal of Linguistics* 41(1). 66–95.
892      https://doi.org/10.1080/07268602.2021.1905607.
893  Seaman, D. Erran, Joshua J. Millspaugh, Brian J. Kernohan, Gary C. Brundige, Kenneth J.
894      Raedeke & Robert A. Gitzen. 1999. Effects of Sample Size on Kernel Home Range
895      Estimates. *The Journal of Wildlife Management* 63(2). 739.
896      https://doi.org/10.2307/3802664.
897  Sloos, Marjoleine. 2013. The reversal of the BÄREN-BEEREN merger in Austrian Standard
898      German. *The Mental Lexicon* 8(3). 353–371. https://doi.org/10.1075/ml.8.3.05slo.
899  Sneller, Betsy. 2018. *Mechanisms Of Phonological Change*. Philadelphia: University of
900      Pennsylvania Ph.D. Dissertation.
901  Stanford, James N., Monica Nesbitt, James King & Sebastian Turner. 2019. Pioneering a dialect
902      shift in the Pioneer Valley: Evidence for the Low-Back-Merger Shift in Western
903      Massachusetts. Presented at the New Ways of Analyzing Variation 48, Eugene, Oregon.
904  Stanley, Joseph A. 2020. *Vowel dynamics of the Elsewhere Shift: A sociophonetic analysis of*
905      *English in Cowlitz County, Washington*. Athens, Georgia: University of Georgia Ph.D.
906      Dissertation.
907  Stanley, Joseph A. 2021. joeyr: Functions for Vowel Data. https://joeystanley.github.io/joeyr/.
908  Strelluf, Christopher. 2018. *Speaking from the heartland: The midland vowel system of Kansas*
909      *City* (Publication of the American Dialect Society 103). Durham, NC: Duke University
910      Press.
911  Tse, Holman. 2018. *Beyond the Monolingual Core and out into the Wild: A Variationist Study of*
912      *Early Bilingualism and Sound Change in Toronto Heritage Cantonese*. Pittsburgh, PA:
913      University of Pittsburgh Ph.D. Dissertation.
914  Upton, Graham & Ian Cook. 2014. *A Dictionary of Statistics*. Oxford University Press.
915      https://doi.org/10.1093/acref/9780199679188.001.0001.

916  Venables, W. N., Brian D. Ripley & W. N. Venables. 2002. *Modern applied statistics with S*
917       (Statistics and Computing). 4th ed. New York: Springer.
918  Warner, Natasha, Allard Jongman, Joan Sereno & Rachèl Kemps. 2004. Incomplete
919       neutralization and other sub-phonemic durational differences in production and
920       perception: evidence from Dutch. *Journal of Phonetics* 32(2). 251–276.
921       https://doi.org/10.1016/S0095-4470(03)00032-9.
922  Warren, Paul. 2018. Quality and quantity in New Zealand English vowel contrasts. *Journal of the*
923       *International Phonetic Association* 48(3). 305–330.
924       https://doi.org/10.1017/S0025100317000329.
925  Wasserstein, Ronald L., Allen L. Schirm & Nicole A. Lazar. 2019. Moving to a World Beyond "p <
926       0.05." *The American Statistician* 73(sup1). 1–19.
927       https://doi.org/10.1080/00031305.2019.1583913.
928  Wassink, Alicia Beckford. 2006. A geometric representation of spectral and temporal vowel
929       features: Quantification of vowel overlap in three linguistic varieties. *The Journal of the*
930       *Acoustical Society of America* 119(4). 2334–2350.
931  Whalen, D. H. & Wei-Rong Chen. 2019. Variability and Central Tendencies in Speech Production.
932       *Frontiers in Communication* 4. 49. https://doi.org/10.3389/fcomm.2019.00049.
933  Wickham, Hadley. 2015. *ggplot2: Elegant Graphics for Data Analysis* (Use R!). 2nd edn. New
934       York: Springer.
935  Wickham, Hadley, Mara Averick, Jennifer Bryan, Winston Chang, Lucy McGowan, Romain
936       François, Garrett Grolemund, et al. 2019. Welcome to the Tidyverse. *Journal of Open*
937       *Source Software* 4(43). 1686. https://doi.org/10.21105/joss.01686.
938  Wong, Amy Wing-mei & Lauren Hall-Lew. 2014. Regional variability and ethnic identity: Chinese
939       Americans in New York City and San Francisco. *Language & Communication* 35. 27–42.
940       https://doi.org/10.1016/j.langcom.2013.11.003.
941

---

[i] See also Kelley & Tucker (2020) for details on the overlapping coefficient, another nonparametric method for calculating overlap.

[ii] Preliminary work for Sneller (2018) attempted, with mixed results, a modified model approach of Bhattaharyya's Affinity by extracting model values for fixed effects and manually adjusting the data. However, the onerousness of this approach makes it not widely implementable in comparison with Pillai, which has the advantage of being able to easily integrate commonly used mixed-effects models.

[iii] In this paper, we focus on the effect of sample size for underlyingly merged speakers. Future work may fruitfully investigate how sample size interacts with Pillai score for underlyingly unmerged speakers as well.

[iv] Let $p_{95} = e^{1-log\left(\frac{n}{2}\right)}$. Since $e^{x-y} = e^{x+(-y)} = e^x e^{-y} = e^x \frac{1}{e^y} = \frac{e^x}{e^y}$, and since $e^1 = e$, then $p_{95} = \frac{e}{e^{\ln\left(\frac{n}{2}\right)}}$. Since $e^{\ln(x)} = x$, then $p_{95} = \frac{e}{\frac{n}{2}}$. Multiplying the numerator and the denominator by 2 produces the final equation, $\frac{2e}{n}$. (We are grateful to the anonymous review who pointed out these simplification steps for this formula.) This is most easily implemented in R as `2*exp(1)/n`.

[v] One reviewer asked why we did not recommend accounting for the difference in sample sizes across speech styles by simply adding a term for style in the linear model. The reasoning is fairly straightforward: researchers cannot control how many tokens of interest are produced in the conversational portion of the interview, allowing the difference in token count across style to vary wildly by speaker. Since sample sizes influence Pillai score in a predictable way but style influences Pillai score in an unpredictable way (as it's dependent on the differences in sample sizes), it is inadvisable to use style as the predictor rather than a more straightforward way to control for or incorporate sample size.

[vi] We determined this using the `compare_performance()` and `check_model()` functions in the `performance` package in R (Lüdecke, Ben-Shachar, et al. 2021). The model fit to the raw $p$-values had a lower adjusted $R^2$ than the one fit to the transformed data (0.2169 vs. 0.4048). It also satisfied assumptions of linearity, homogeneity of variance, and normality of residuals, unlike the

model fit to raw *p*-values. Models fit to raw or log-transformed Pillai scores had lower $R^2$ values and failed to meet satisfy of the model assumptions.

[vii] Note that we sample with replacement (`replace=TRUE`), in order to make our resampling comparable across all speakers including our least talkative speaker.

If we were to resample *without* replacement, this introduces a confound related to sample size: the amount of error introduced per speaker is proportional to their sample size. Resampling without replacement 1000 times would produce 1000 identical distributions for the least talkative speaker, but 1000 different distributions for every other speaker. Resampling *with* replacement allows us to obtain a standard deviation for each speaker with similar confidence, and introduces a similar amount of uncertainty across speakers in the means of their distributions. For more information on bootstrapping with replacement, we refer the reader to Ismay & Kim (2020).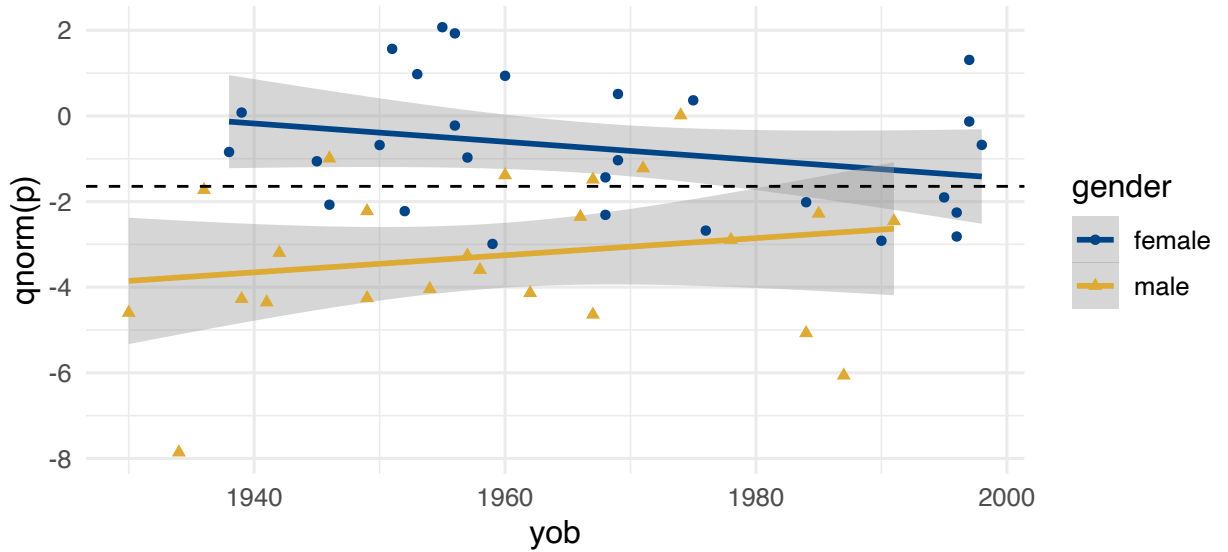