535 confidence that Spanish learners' perception of Dutch /ɑ/~/aː/ is affected by the number of peaks
536 in a training distribution.
537
538 **3.3. Bayes factors**
539
540 From having found a *p*-value above 0.05 we cannot draw any conclusions about whether the null
541 hypothesis is true or false. Because we wanted to be able to quantify evidence in favor of both
542 the alternative *and* the null hypothesis, we computed Bayes factors (henceforth "BFs") (e.g.,
543 Kass and Raftery, 1995; Rouder et al., 2009; Gallistel, 2009; Kruschke, 2010). A BF denotes the
544 likelihood ratio of the data occurring under the null hypothesis ($H_0$) versus the data occurring
545 under the alternative hypothesis ($H_1$):
546

547 $$\mathrm{BF}_{01} = \frac{p(data|H_0)}{p(data|H_1)}$$

548
549 The "01" in this equation refers to $H_0$ and $H_1$ respectively. Thus, if $\mathrm{BF}_{01} = 10$, the observed data
550 are 10 times more likely to occur if $H_0$ is true than if $H_1$ is true; if BF01 = 0.1, the observed data
551 are 10 times more likely to occur if $H_1$ is true than if $H_0$ is true. If we assume that $H_0$ and $H_1$ are
552 equally likely a priori (as is common and as we do henceforth), the Bayes factor $\mathrm{BF}_{01}$ can be said
553 to quantify the evidence in support of $H_0$ over $H_1$. Thus, if $\mathrm{BF}_{01} = 10$, $H_0$ is 10 times more likely
554 to be true than $H_1$ (i.e., the odds are 10 to 1 in favor of $H_0$); if $\mathrm{BF}_{01} = 0.1$, $H_1$ is 10 times more
555 likely to be true than $H_0$; (i.e., the odds are 10 to 1 in favor of $H_1$). Whether a clear choice
556 between the two hypotheses is possible, depends on the magnitude of the Bayes factor. If $\mathrm{BF}_{01} >$
557 20, there is said to be strong support for $H_0$, and if $\mathrm{BF}_{01} < 1/20$, there is said to be strong support
558 for $H_1$; if, however, $\mathrm{BF}_{01}$ lies between 3 and 20, the data are said to moderately favor $H_0$, and if
559 $\mathrm{BF}_{01}$ lies between 1 and 3, the data are said to only trivially favor $H_0$ (Kass and Raftery, 1995).
560
561        In the current paper, the null and alternative hypotheses are defined in terms of the
562 standardized effect size of the difference in the improvement score (= the post-test minus the pre-
563 test accuracy percentage) between the Unimodal and Bimodal groups, i.e., in terms of how much
564 the two groups differ in their improvement of categorization accuracy after as compared to
565 before training. An observed effect size *d* can be calculated as the number of standard deviations
566 difference between two improvement scores:
567
568        *d* = (improvement score of group 1 – improvement score of group 2) / standard deviation
569
570 where the standard deviation is the pooled standard deviation.[12] In our case group 1 is the
571 Bimodal group and group 2 the Unimodal group.
572
573        The null hypothesis (Figure 5, top) is always the same, namely that there is no difference
574 in the improvement score between the Unimodal and Bimodal groups, and that accordingly the
575 effect size *d* is exactly zero:
576
577        $H_0$:    $d = 0$
578

---

[12] The pooled standard deviation is calculated as the within-sums-of-squares / (N1+N2-2).

15

579   *<Insert Figure 5 around here>*

581   The value of the BF depends on the definition of the alternative hypothesis. To accommodate
582   different *a priori* beliefs about the effect size, we computed the BF in four different ways, i.e.,
583   with four different alternative hypotheses, which are increasingly less specific about the expected
584   value of the effect size. The first and second alternative hypotheses ($H_1$ and $H_2$) include
585   information about the effect size obtained from EBW2011, WER2013 and WB2013; the third
586   and fourth alternative hypotheses ($H_3$ and $H_4$) do not. Table 4 provides an overview of the four
587   alternative hypotheses and the resultant BFs, which we will now discuss in detail.[13]

589   **Table 4**: The four alternative hypotheses (H) and the resulting Bayes factors (BF).

| H | | BF |
|---|---|---|
| $H_1$: | $d = + 0.50$ | $BF_{01} = 137.86$ |
| $H_2$: | $d$ is a random value drawn from a uniform distribution between 0 and 1. | $BF_{02} = 5.97$ |
| $H_3$: | $d$ is a random value drawn from a Gaussian distribution with mean 0 and standard deviation 1. | $BF_{03} = 5.32$ |
| $H_4$: | $d$ is a random value drawn from a Cauchy distribution | $BF_{04} = 4.73$ |

593         Alternative hypothesis 1 (Figure 5, second from top) stipulates that the effect size *d* is a
594   specific value:

596         *$H_1$*:     $d = + 0.50$

598   This value of +0.50 is based on effect sizes derived from the improvement scores observed in
599   EBW2011, WER2013 and WB2013, as follows. In EBW2011 and WER2013, one group of
600   listeners was exposed to a non-enhanced bimodal distribution (the Bimodal group), a second
601   group to an enhanced bimodal distribution (the Enhanced group), and a third group to music (the
602   Music group). In WB2013, improvement in categorization was compared between a Music group
603   and two Enhanced groups, one presented with a discontinuous distribution and the other to a
604   continuous distribution. As mentioned in the Introduction (section 1.4), in all three studies the

---

13 The four Bayes factors can be computed in R (R Core Team, 2013) with the equation **dt** (*t*, *df*) /
(**mean** (*weight* * **dt** (*t*, *df*, **ncp** = *d* * **sqrt**(*n*))) / **mean** (*weight*)). In this equation, **dt** is the R function that computes
the *t* probability density, and **ncp** is the non-centrality parameter of this density; *t* is the between-groups *t* value of
our experiment, i.e. -0.43; *df* is the number of degrees of freedom for a *t* test, i.e. 60+60-2 = 118; *n* is half the
geometric mean of the two group sizes (Rouder et al. 2009, p.234), i.e. 60*60/(60+60) = 30; *d* is the hypothesized
range of possible effect sizes, and *weight* is the shape of the distribution for all these *d* values. For $H_1$, *d* is 0.5
and *weight* is 1. For $H_2$, *d* is (-0.5+1:1e5)/1e5 and *weight* is 1. For $H_3$, *d* is ((-10e5**width*+0.5):(10e5**width*-
0.5))/1e5 and weight is exp(-0.5*(*d*/*width*)^2), where *width* is 1. For $H_4$, *d* is ((-
1000*1e4**width*+0.5):(1000*1e4**width*-0.5))/1e4 and *weight* is 1/(1+(*d*/*width*)^2)), where *width* is sqrt(2)/2 (our
equations for $H_3$ and $H_4$ are formulated in such a way that they will also work for other values of *width*). At the time
of writing the computations for $H_3$ and $H_4$ are also available on Rouder's website
(http://pcl.missouri.edu/bayesfactor).

605  improvement score was significantly larger for the Enhanced group than for the Music group. In
606  EBW2011 and WER2013, the improvement score for the Bimodal group was not significantly
607  different from that of the Music group and also not from that of the Enhanced group. For the
608  current analysis, we considered the improvement scores of the previous Enhanced groups as
609  proxies for the expected improvement score of our Bimodal group (which was also exposed to an
610  enhanced bimodal distribution, just as the Enhanced groups in the previous studies; section 1.6).
611  Because it was not clear whether our Unimodal group would behave more similarly to the
612  previous Music groups or to the previous Bimodal groups, we considered the improvement
613  scores of the previous Music and Bimodal groups as proxies for the expected improvement score
614  of our Unimodal group. When calculating the effect sizes observed in the three studies, we used
615  the above-mentioned formula for the effect size $d$, and took a previous Enhanced group as group
616  1, and either a previous Bimodal group or a previous Music group as group 2. The improvement
617  scores for the Enhanced, Bimodal and Music groups were 6.04% (CI = +2.76 ~ +9.31%), 0.80%
618  (CI = –2.22 ~ +3.83%) and –0.15% (CI = –3.50 ~ +3.21%) respectively in EBW2011, and 6.63%
619  (CI = +4.05 ~ +9.20%), 3.83% (CI = +0.97 ~ 6.68%) and 2.00% (CI = –0.50 ~ +4.50%)
620  respectively in WER2013. The improvement scores for the Enhanced and Music groups in
621  WB2013 were 9.68% (CI=+6.80%~+12.55) and 2.00% (CI= –0.50~+4.50) respectively.[14] The
622  pooled standard deviation for the Enhanced and Bimodal groups was 12.00% in EBW2011 and
623  9.57% in WER2013. The pooled standard deviation for the Enhanced and Music groups was
624  12.09% in EBW2011, 8.94% in WER2013 and 9.50% in WB2013. Table 5 shows the resulting
625  effect sizes $d$.
626
627  **Table 5**: Effect size $d$ in previous studies (see text).
628

| Previous study | Enhanced–Bimodal | Enhanced–Music |
|---|---|---|
| EBW (2011) | +0.44 | +0.51 |
| WER (2013) | +0.29 | +0.52 |
| WB (2013) | | +0.81 |

629
630
631        The average of the five listed effect sizes is +0.51, which we rounded to +0.50 in
632  hypothesis 1. Notice that this value is explicitly positive, i.e., it reflects the belief that our
633  Bimodal group will have a *higher* improvement score, and thus improve *more* after distributional
634  training than the Unimodal group. The BF calculated on the basis of the null hypothesis versus
635  this first alternative hypothesis expresses strong support for the null:
636
637        $BF_{01} = 137.86$
638
639  Specifically, $BF_{01}$ indicates that the observed data are 137.86 times more likely to have occurred
640  under $H_0$ (that $d$ is exactly 0), than under $H_1$ (that $d$ is exactly 0.5).
641

---

14 The Enhanced group referred to here is the group presented with a continuous enhanced distribution in WB2013
(the Continuous Enhanced group). In WB2013 the group presented with a discontinuous enhanced distribution (the
Discontinuous Enhanced group) and the Music group were taken from WER2013.

642        In alternative hypotheses 2 through 4, the effect size is no longer defined as a specific
643    value, but as a probability density function (Figure 5, as explained below): $d$ is expected not to
644    be one specific value, but a random value drawn from a distribution whose form defines the
645    likelihood of that value. In alternative hypothesis 2, the effect size is any value between 0 and 1
646    with equal probability (Figure 5, middle):
647
648        $H_2$:    $d$ is a random value drawn from a uniform distribution between 0 and 1.
649
650    The hypothesis still includes the information mentioned in Table 5 about previously obtained
651    effect sizes (i.e., all effect sizes in Table 5 fall within the range of the distribution), but it is
652    vaguer about the precise value of the expected effect size than hypothesis 1. Since $d$ is defined as
653    0 or positive, hypothesis 2 expresses the belief that the Bimodal group will improve *at least as*
654    *much* as the Unimodal group. The BF calculated on the basis of the null hypothesis versus this
655    second alternative hypothesis also expresses support for the null:
656
657        $BF_{02} = 5.97$
658
659    That is, $BF_{02}$ implies that the observed data are 5.97 times more likely to have occurred under $H_0$
660    (that $d$ is exactly 0) than under $H_2$ (that $d$ is somewhere between 0 and 1).
661
662        Hypotheses 1 and 2 show that previous observations can be incorporated in the
663    alternative hypothesis to different extents, depending on the researcher's belief in the truth value
664    of these observations. Previous observations can also be deemed inappropriate for incorporation
665    in the alternative hypothesis, for example if concerns (such as mentioned in the section 1.2)
666    about the earlier observations create uncertainty about the applicability of the information to the
667    experiment to be performed. In this case, the alternative hypothesis should reflect the assumption
668    that we do not have a clear expectation about the effect size. This is done in alternative
669    hypotheses 3 and 4. In alternative hypothesis 3, the effect size is any value around 0, with values
670    closer to the mean being more likely than values further away from the mean as defined by a
671    Gaussian distribution (Figure 5, fourth from top):
672
673        $H_3$:    $d$ is a random value drawn from a Gaussian distribution with a mean of 0 and a
674                    standard deviation of 1.
675
676    Since $d$ can be positive, zero or negative, the belief that the Bimodal group will improve at least
677    as much as the Unimodal group, which was inherent in alternative hypotheses 1 and 2, is now
678    dropped. The BF calculated on the basis of the null hypothesis versus the third alternative
679    hypothesis still expresses support for the null:
680
681        $BF_{03} = 5.32$
682
683    In other words, $BF_{03}$ indicates that the observed data are 5.32 times more likely to have occurred
684    under $H_0$ (that $d$ is exactly 0) than under $H_3$, (that $d$ is a value around zero, whose probability is
685    defined by a Gaussian distribution).
686

687   It is possible to be even less specific about the expected value of the effect size than in
688   alternative hypothesis 3, by loosening the belief that the effect size is more likely to occur close
689   to zero. This is done with a Cauchy distribution (for an explanation, see Rouder et al., 2009), as
690   used in alternative hypothesis 4 (Figure 5, bottom):
691
692   *H₄*:   *d* is a random value drawn from a Cauchy distribution, with a width of $(\sqrt{2})/2$.[15]
693
694   Notice in Figure 5 that the tails of the Cauchy distribution are much heavier than those of the
695   Gaussian distribution, thus reflecting a much smaller confidence that the effect size should be
696   relatively close to zero. Again, the BF calculated on the basis of the null hypothesis versus the
697   fourth alternative hypothesis expresses support for the null:
698
699   $BF_{04} = 4.73$
700
701   Thus, $BF_{04}$ indicates that the observed data are 4.73 times more likely to have occurred under $H_0$
702   (that *d* is exactly 0) than under $H_4$ (that *d* is a value around zero, whose probability is defined by
703   a Cauchy distribution, i.e., with more uncertainty as to the effect size than expressed in the
704   Gaussian distribution used for $H_3$).
705
706   In sum, four different calculations of the Bayes factor, which differ in the extent to which
707   they incorporate *a priori* beliefs about the expected effect size, unanimously support the null
708   hypothesis that there is no difference between bimodally and unimodally trained Spanish
709   participants in improvement of categorization of Dutch [ɑ]- and [a]-tokens. If we follow the
710   interpretation of Bayes factors by Kass and Raftery (1995; section 3.3), the support for the null
711   hypothesis ranges from moderate support (hypotheses 2 through 4, which represent less strong *a*
712   *priori* beliefs about the effect size than hypothesis 1) to strong support (hypothesis 1, which
713   incorporates the most explicit *a priori* beliefs).
714
715   **4. Discussion**
716
717   In the present study we trained Spanish adult participants on a bimodal or a unimodal
718   distribution encompassing the Dutch vowel contrast /ɑ/~/a/, and then tested their improvement in
719   categorization of Dutch [ɑ]- and [a]-tokens after training. For the first time in the research on
720   distributional learning of speech sounds, the bimodal and unimodal distributions had nearly
721   identical dispersions, as defined by the range, standard deviation and edge strength. The results
722   show that Spanish adult participants improve their categorization of Dutch [ɑ]- and [a]-tokens
723   irrespective of the training distribution, and that categorization accuracy does not improve
724   significantly more after exposure to one distribution than after exposure to the other distribution.
725   Additionally, four different Bayes factors (ranging from incorporating *a priori* beliefs about the
726   expected effect size as much as possible to not incorporating previous knowledge at all) provided
727   unanimous evidence for the null hypothesis that there is no difference between bimodally and

---

15 The equation used for the Cauchy distribution is: ((-1000*1e4**width*+0.5):(1000*1e4**width*-0.5))/1e4,
where *width* is sqrt(2)/2 (see also note 12).