



UNIVERSIDAD DE BUENOS AIRES
FACULTAD DE CIENCIAS EXACTAS Y NATURALES
DEPARTAMENTO DE COMPUTACIÓN

Implementación y evaluación de un sistema de síntesis de habla con acento extranjero variable

Tesis de Licenciatura en Ciencias de la Computación

Franco Negri

Director: Agustin Gravano

Codirector: Master Yoda

Buenos Aires, 2018

IMPLEMENTACIÓN Y EVALUACIÓN DE UN SISTEMA DE SINTESIS DE HABLA

Un sistema de Text To Speech (TTS) es aquel que genera habla artificial a partir de un texto de entrada. En la actualidad estos sistemas se encuentran incluidos en muchas aplicaciones domesticas, desde navegación por GPS, asistentes personales inteligentes (como es el caso de SIRI), ayuda para personas no videntes, traducción automática, etc.

En las últimas décadas se han visto grandes progresos en este campo, siendo capaces de modelar con cierto grado de efectividad cuestiones tales como la prosodia del hablante, emociones, etc. Si bien actualmente se considera que el estado del arte para la síntesis es el entrenamiento con redes neuronales profundas (DNN), una técnica todavía utilizada es la que utiliza modelos ocultos de markov mas modelos de mezcla de gaussianas (HMM+GMM) que, a partir de un corpus de datos de entrenamiento, extrae información acústica y genera un modelo probabilístico que permita sintetizar habla.

Consideramos que este metodo si bien no nos permitirá obtener la mejor voz posible, nos permitirá entender mejor los features con los que trabajamos y por lo tanto modificarlos segun consideremos conveniente, algo que puede volverse mas complicado cuando se esta trabajando con redes neuronales profundas.

En este trabajo de tesis se estudia una manera posible de generar un TTS basado en HMMs capaz de sintetizar habla en español con acento extranjero. Las razones por las que podría querer diseñarse un sistema con estas características varían desde un punto de vista puramente técnico, ya que un sistema así permitiría la utilización de corpus de entrenamiento de hablantes no nativos para la generación de una nueva voz, pasando por cuestiones lingüísticas, como es poder vislumbrar el limite en que un acento deja de parecernos local para pasar a ser extranjero y cuestiones psicologicas: lograr distintos efectos sobre el usuario, quien podría reaccionar distinto ante diferentes acentos.

En el transcurso de este trabajo se espera además evaluar la prosodia y la fonética del modelo generado con estas características, como así también evaluar su inteligibilidad. Además pretendemos evaluar la efectividad de técnicas de speaker adaptation cuando se utilizan corpus de distintas nacionalidades con repertorios fonéticos disimiles (para este caso de estudio: castellano e ingles)

Para este trabajo nos basaremos fuertemente en la síntesis/análisis mel-cepstral, speech parameter modeling usando HMMs y speech parameter generation usando HMMs, como es descrito en la disertación doctoral *Simultaneous Modeling of phonetic and prosodic parameters, and characteristic conversion for hmm-based text-to-speech systems* del Profesor Tadashi Kitamura, Nagoya Institute of Technology[7].

Keywords: Síntesis de habla, HMM, HTS, GMM, acento extranjero, aprendizaje automatico.

AGRADECIMIENTOS

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Fusce sapien ipsum, aliquet eget convallis at, adipiscing non odio. Donec porttitor tincidunt cursus. In tellus dui, varius sed scelerisque faucibus, sagittis non magna. Vestibulum ante ipsum primis in faucibus orci luctus et ultrices posuere cubilia Curae; Mauris et luctus justo. Class aptent taciti sociosqu ad litora torquent per conubia nostra, per inceptos himenaeos. Mauris sit amet purus massa, sed sodales justo. Mauris id mi sed orci porttitor dictum. Donec vitae mi non leo consectetur tempus vel et sapien. Curabitur enim quam, sollicitudin id iaculis id, congue euismod diam. Sed in eros nec urna lacinia porttitor ut vitae nulla. Ut mattis, erat et laoreet feugiat, lacus urna hendrerit nisi, at tincidunt dui justo at felis. Class aptent taciti sociosqu ad litora torquent per conubia nostra, per inceptos himenaeos. Ut iaculis euismod magna et consequat. Mauris eu augue in ipsum elementum dictum. Sed accumsan, velit vel vehicula dignissim, nibh tellus consequat metus, vel fringilla neque dolor in dolor. Aliquam ac justo ut lectus iaculis pharetra vitae sed turpis. Aliquam pulvinar lorem vel ipsum auctor et hendrerit nisl molestie. Donec id felis nec ante placerat vehicula. Sed lacus risus, aliquet vel facilisis eu, placerat vitae augue.

A mi persona favorita: la transformada de Fourier

Índice general

0.1.	Preparación De Corpus	1
0.2.	Repertorio Fonetico y Mapeo De Fonemas	3
0.3.	Interpolación Entre Modelos	4
0.4.	Speaker-Adaptative Training	4
0.5.	Herramientas	4
0.5.1.	Festival y Festvox	5
0.5.2.	HTS	5
0.5.3.	HTS_engine	7
0.6.	Experimentación	9
0.6.1.	Interface	9
0.7.	Resultados	12
0.7.1.	Datos demográficos	12
0.7.2.	Inteligibilidad	13
0.7.3.	Datos normalizados	16
0.7.4.	Análisis de Nacionalidad	20
0.7.5.	Resultados Generales de la experimentación	22
0.8.	Trabajo Futuro	25
0.9.	Apendice	26
0.9.1.	Lista de Fonemas	26
0.9.2.	Mapeo Fonemas del Ingles-Castellano	26
0.9.3.	Parametros utilizados para el entrenamiento	26

0.1. Preparación De Corpus

Introducción

En esta introducción presentamos a modo de resumen, la metodología utilizada para el entrenamiento del sistema y otros detalles teóricos como el mapeo de fonemas necesario para adaptar el repertorio fonémico del inglés al castellano, etc.

Los pasos a seguir serán:

1. A partir de los corpus de datos disponibles, realizar un etiquetado fonémico de los corpus para su posterior utilización en el entrenamiento de los HMMs.
2. Realizar el entrenamiento de los sistemas (Uno por cada corpus disponible). Para esto contaremos con el framework de modelado de HMMs HTS.
3. Realizar un mapeo entre los fonos del castellano y los del inglés. Estos serán útiles en el paso siguiente donde será requerimiento indispensable tener modelados los mismos fonos para todos los modelos.
4. Utilizar las herramientas provistas por HTS para interpolar entre modelos y poder sintetizar habla con distintos grados de fonética y prosodia inglesa.

Más adelante continuaremos explicando algunos detalles implementativos de las herramientas utilizadas.

Preparación De los datos

Como fase inicial de este trabajo, es un paso indispensable conseguir el etiquetado fonético para las grabaciones de habla, que consistirá principalmente en una lista de fonos donde se indica donde comienza y donde termina cada uno. De esto dependerá fuertemente la calidad de las oraciones que logremos sintetizar, por lo que es necesario prestar especial cuidado en que estas estén lo mejor alineadas posible con los audios. Como ya adelantamos en la introducción, estas son necesarias para entrenar los HMM+GMM, extrayendo de aquí tanto la información acústica para cada fono (cosas tales como la frecuencia principal, la duración, etc) como así también información contextual (como por ejemplo: como suena un fonema cuando está seguido de algún otro, al principio de una oración, si se encuentra en un diptongo), por lo que una mala transcripción se traducirá indefectiblemente en un mal modelo y malas oraciones sintetizadas.

Al comienzo de la investigación comenzamos con tan solo un corpus de datos, *secyt-mujer*[1]. Este está compuesto por 741 oraciones de Español Rioplatense, equivalentes a 48 minutos de habla. Para el mismo también contábamos con sus transcripciones fonéticas y grafémicas anotadas de manera manual.

Realizamos varias pruebas de concepto utilizando HTS y este corpus, experimentando con diversos métodos para obtener el etiquetado fonético.

La primera estrategia consistió en utilizar alineamiento automático utilizando EHMM alignment [8] utilizando Festival y Festvox. Estos programas presentan como ventaja que tanto en la anotación del corpus como la generación de trazas fonéticas (utilizadas para la síntesis) presentan el mismo repertorio fonético. Esto es útil ya que podemos garantizar que el entrenamiento y la síntesis están utilizando el mismo método de generación de trazas/oraciones.

Aún así, los resultados preliminares con este metodo fueron bastante adversos: los audios generados resultaban poco inteligibles notandose claros defectos acústicos, El más notable siendo el fono /r/ que se asemejaba mas a /r/.

Utilizando Praat para visualizar el alineamiento entre utternaces y audios, descubrimos que la alineación estaba desfasada algunas milesimas de segundo. Sospechamos que esto se debió a algún problema con la normalización de los audios.

Dado que para este corpus contábamos con las transcripciones fonéticas anotadas de manera manual procedemos a implementar un híbrido con EHMM. En este híbrido tomamos las anotaciones hechas a mano para cada fonema y la información contextual y el repertorio fonético generado a partir del proceso de EHMM. De esta manera buscamos mejorar la alineación pero manteniendo el mismo repertorio fonetico y la misma meta-información brindada por el alineamiento automático.

El modelo generado con estos utternaces mixtos resulto ser superior a los generados solo con alineamiento automático. Aún así los audios sintetizados todavía no alcanzan una calidad aceptable, realizando pruebas internas todavía notamos que el sonido resultaba metalico y las frases poco inteligibles. Además se pudo percibir de manera informal otros detalles tales como que la voz original tenía un pitch mayor que la producida por los modelos, alrededor de un 10 %.

En este momento del trabajo obtenemos otro corpus de datos *loc1_pal*[2] con 1593 oraciones en Español Rioplatense con aproximadamente 2 horas y 26 minutos de habla. Con este nuevo conjunto de datos esperamos conseguir mejores resultados.

Para este corpus no cotábamos con transcripciones fonéticas manuales por lo que nos vimos forzados a utilizar EHMM nuevamente. Aún así, los resultados fueron superiores a los conseguidos con *secyt-mujer*. Los audios sintetizados resultaban inteligibles y con un marcado acento rioplatense. Tras algunas pruebas de concepto donde se experimentó con varios valores de GAMMA, rango de las frecuencias principales, y otros parámetros que consideramos podían afectar la calidad de la voz, logramos obtener resultados que superaban de manera significativa aquellos obtenidos previamente con el otro corpus. Por consiguiente concideramos que los audios generados habían alcanzado una buena calidad que serultaba ininteligible y aceptable para el objetivo de la investigación, por lo que decidimos utilizar uno de estos modelos para el resto del trabajo.

Especulamos que la disparidad en la calidad de los resultados es causada principalmente por la cantidad de audios y horas de habla de cada corpus[5]. Consideramos que esto juega un papel predominante en la calidad de los TTS generados, aún cuando se utiliza un método de etiquetado puramente automático y propenso a errores sistemáticos en el alineamiento.

Finalmente necesitamos generar otra voz con un idioma diferente que nos permita interpolar con el modelo previamente detallado. Para esto utilizamos el corpus *CMU-ARCTIC-SLT*[3] con 1132 oraciones en ingles y 56 minutos de habla, disponible en la pagina de hts [9]. Ya que este corpus venía a modo de demo con hts, asumimos que los parametros de entrenamiento y las transcripciones foneticas ya habían sido seleccionadas de manera apropiada, por lo que no intentamos mejorar la calidad de las transcripciones foneticas mas allá de lo que la demo ofrecía.

Para este trabajo todos los audios usarán sampling rate de 48kHz, precisión de 16 bits, mono.

El rango de extracción de frecuencia principal para utilizado fue de 100 hz a 350hz.

Una lista extensiva de los parámetros utilizados para el entrenamiento se puede ver en

el apéndice 3.

0.2. Repertorio Fonético y Mapeo De Fonemas

Para la generación de utterances, tanto de los audios en inglés como en castellano, utilizamos los repertorios fonéticos brindados por festvox (ver apéndice 1).

El primer desafío que se presenta es que estos repertorios fonéticos no tienen un mapeo directo con el Alfabeto Fonético Internacional: por ejemplo con este repertorio fonético, en el castellano existen tres fonos distintos para la /i/. Esta decisión por parte de festvox proviene de la necesidad de poder diferenciar la /i/ acentuada de la no acentuada y de aquella presente en los diptongos: /ia/, /ie/, /io/, /iu/.

Castellano	Ingles	Castellano	Ingles
a	aa	p	jh
a1	ae	r	k
b	ah	rr	l
ch	ao	s	m
d	aw	t	n
e	ax	u	ng
e1	ay	u0	ow
f	b	u1	oy
g	ch	x	p
i	d		r
i0	s		sh
i1	dh		t
k	eh		th
l	er		uh
ll	ey		uw
m	f		v
n	g		w
ny	hh		y
o	ih		z
o1	iy		zh

Tab. 0.1: Repertorios Fonéticos

De esto, surge el siguiente problema: como sintetizar oraciones en castellano utilizando un repertorio fonético en inglés, donde incluso la cantidad de fonos es diferente. La manera que encontramos de abordar esto fue realizando de manera perceptual una grilla donde cada fonema del castellano estuviera mapeado a uno del inglés. En otras palabras, antes del entrenamiento del modelo, tomamos los fonos generados por festival para el corpus en inglés y los reemplazamos con fonos del castellano, a partir de la siguiente tabla:

Castellano	Ingles	Castellano	Ingles
aa	a1	k	k
ae	a	l	l
ao	o	m	m
ou	o1	n	n
b	b	nx	n
ch	ch	ng	ng
d	d	p	p
dh	d	r	r
dx	dx	s	s
eh	e	t	t
el	e1	uh	u1
em	em	uw	u
en	en	w	u0
er	er	th	notUsed
ei	ei	v	v
f	f	jh	ll
g	notUsed3	y	y
hh	h	sh	sh
hv	hv	zh	zh
ih	i1	z	notUsed2
iy	i		

Tab. 0.2: Mapeo Fonetico

Aún así, para varios fonos tuvimos que hacer reglas especiales ya que no contabamos con ninguno del ingles lo suficientemente similar. Así, para el fono ɲ colapsamos las apariciones del fono /n/ seguido de /j/ . Si bien esta solución no es optima, ya que estamos tomando dos fonos para simular otro, esta solución se aproxima a la manera real en la que un hablante no nativo aprende un idioma con una carga fonética diferente al suyo. Citando un extracto del trabajo *Transcription of Spanish and Spanish-Influenced English*, Brian Goldstein, Temple University:

Consonants

As indicated in Table 5, there are many ways in which the features of Spanish influence the production of consonants in English. These influences cut across all sound classes, although the majority of influences will be in the fricative sound class. Several factors influence the extent to which one phonological system influences another. First, the influence may be due to the absence of phonemes or allophones in a language (Iglesias & Goldstein, 1998). For example, $\text{[p}^{\text{h}}\text{]}$, $\text{[t}^{\text{h}}\text{]}$, and $\text{[k}^{\text{h}}\text{]}$ do not occur in Spanish, and [f] , [v] , and $\text{[d}_3\text{]}$ do not

occur in most dialects of Spanish. In attempting to produce sounds in English that do not exist in Spanish, a native Spanish speaker might substitute a close relation. Thus, /ʃ/ might be produced as [f] ; /ʃo/ *show* \rightarrow [f]o . Second, there are differences in the phonotactic constraints of the two languages. In Spanish, word-initial clusters cannot begin with /s/ . Thus, Spanish speakers attempting to produce English clusters of that type might exhibit either *cluster reduction* (e.g., /stɑəz/ *stars* \rightarrow [tɑəz]) or *epenthesis* (or *prothesis*) (e.g., /stɑəz/ *stars* \rightarrow [estɑəz]) (Perez, 1994). Third, there are differences in the distribution of sounds. In Spanish, for example, the only word-final consonants are /s/ , /n/ , /r/ , /l/ , and /d/ .

Visto desde nuestra perspectiva, una persona que aprende una nueva lengua, realiza una aproximación entre los fonos conocidos y los fonos ‘objetivo’ de la nueva lengua.

De manera similar el ingles carece del fono trill alveolar /r/ y dado que el fono /r/ ya estaba siendo utilizado, no podíamos realizar un mapeo tan directo. Como solución tomamos la mitad de los fonos /r/ y los remplazamos con /r/ . Un problema similar surge

del fono /hh/, por lo que decidimos tomar la mitad de los fonos etiquetados como /hh/ y remplazarlos con /g/.

Aquellos fonos que consideramos suficientemente disimiles del castellano, como es el caso de la /sh/, /z/, etc los mapeamos a caracteres que no interfirieran para el entrenamiento (notUsed, notUsed2) ya que no los utilizaremos para la síntesis.

Utilizando estos fonos a la hora de entrenar nos permite sintetizar oraciones en castellano, aunque como es de esperar, dado que el corpus de entrenamiento es tan disimilar de las oraciones que queremos sintetizar, donde tanto las combinaciones de fonemas, las reglas prosódicas y las acentuaciones vocálicas son distintas, los audios sintetizados resultan incomprensibles y de muy baja calidad

0.3. Interpolación Entre Modelos

Una vez generados ambos modelos con el mismo repertorio fonético procedemos a experimentar de manera informal la efectividad del método. Tras algunos ajustes en los pesos fue posible generar oraciones donde la carga fonética podía reconocerse como un extranjero estadounidense o angloparlante hablando castellano. Concluimos esto en base a los detalles distintivos como el fono /r/ mas suavizada, o pronunciación de vocales mas abiertas, que comúnmente se atribuyen a extranjeros angloparlantes. Además para ciertos pesos en la interpolación, también pudimos apreciar cierto cambio en la prosodia.

0.4. Speaker-Adaptative Training

Se realizaron varias pruebas de concepto donde habiendo entrenado un HMM con *CMU-ARCTIC-SLT* se le realiza speaker adaptation con el corpus de *loc1_pal*, pero resultaron no concluyentes. El HMM final perdía completamente las características y el acento de *CMU-ARCTIC-SLT* por lo que se decidió no proseguir por este camino.

0.5. Herramientas

En esta sección presentamos las herramientas utilizadas para este trabajo.

0.5.1. Festival y Festvox

Festival es un framework que permite sintetizar habla. Además posee una gran variedad de APIs, para el procesamiento de audios y generación de nuevos TTS.

Festvox a su vez expande sobre Festival, agregando todavía mas herramientas relacionadas a la síntesis y generación de modelos, que van desde la generación de modelos prosódicos, hasta etiquetado automático de corpus.

Para este trabajo utilizaremos Festival y festvox para generar los Utternaces requeridos tanto para el entrenamiento como para la síntesis de audios. Estos consisten básicamente en una transcripción fonética de los audios dividida en segmentos temporales y datos contextuales tales como la cantidad de sílabas en la palabra siendo transcrita, fonemas que preceden y proceden al actual, etc.

En particular, Festival también cuenta con herramientas de etiquetado automático. Para este trabajo utilizaremos EHMM alignment, que a partir de un corpus y sus transcripciones, permite generar utternaces cuyos segmentos coinciden con aquellos de los audios.

Como veremos mas adelante, estos utternaces serán utilizados en el entrenamiento con HTS para modelar cada uno de los fonemas.

0.5.2. HTS

HTS es un framework de entrenamiento y síntesis de sistemas TTS basado en HMMs que modela simultáneamente la duración, el espectro (mel-cepstrum) y la frecuencia principal (f_0) de utilizando una combinación de HMMs:

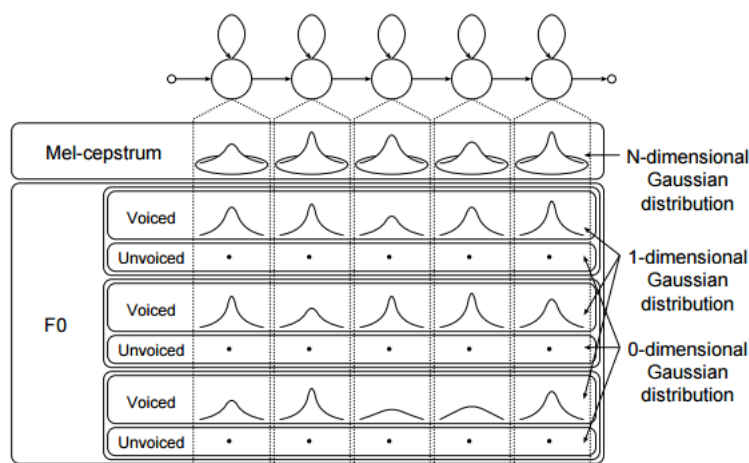


Figure 5.3: structure of HMM.

Fig. 0.1: Estructura de un hmm (simultaneous modeling of phonetic and prosodic parameters, and characteristic conversion for hmm-based text-to-speech systems, takayoshi yoshimura, january 2002, pag. 41)

Mel-cepstum (espectro) y los tres vectores del f_0 (frecuencia principal) son modelados en paralelo.

Por otro lado HTS toma la decisión de modelar la información prosódica dentro de este mismo framework. Para esto, las distribuciones para el espectro, la frecuencia principal y las duraciones son clusterizadas independientemente utilizando la información contextual extraída de los audios de entrenamiento. A continuación se presenta una vista esquemática de la estructura del HMM generado clusterizado con arboles de decisión:

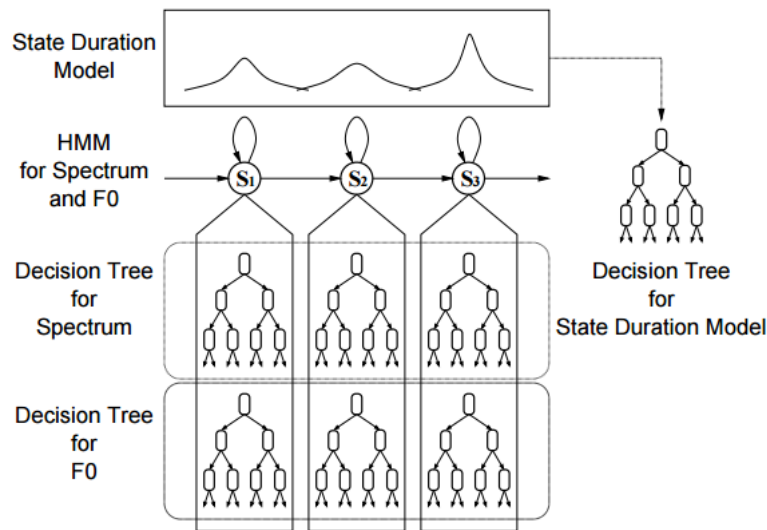


Figure 5.4: Decision trees.

Fig. 0.2: Esquema HMM generado utilizando arboles de decisión (simultaneous modeling of phonetic and prosodic parameters, and characteristic conversion for hmm-based text-to-speech systems, takayoshi yoshimura, january 2002, pag. 45)

En particular para este trabajo la clusterización de datos se realizó generando arboles de decisión y para cada fonema se tomó los dos fonemas precedentes y los dos fonemas procedentes y se extrajo la siguiente información contextual.

- Modo de articulación del fonema.
- Punto de articulación del fonema.
- La perspectiva articulatoria (anterior, central o posterior).
- Si el fonema es una vocal o una consonante.
- En caso de ser una vocal, a que categoría pertenece: por ejemplo para el fonema $/i/$: i (no acentuada), $i0$ (diptongo) , $i1$ (acentuada).
- En caso de ser una vocal, su redondeamiento vocálico.
- En caso de ser una consonante, si es lenis o fortis.

En la siguiente imagen se muestra un fragmento de un árbol de decisión generado para modelar la duración de los fonemas en un hmm:

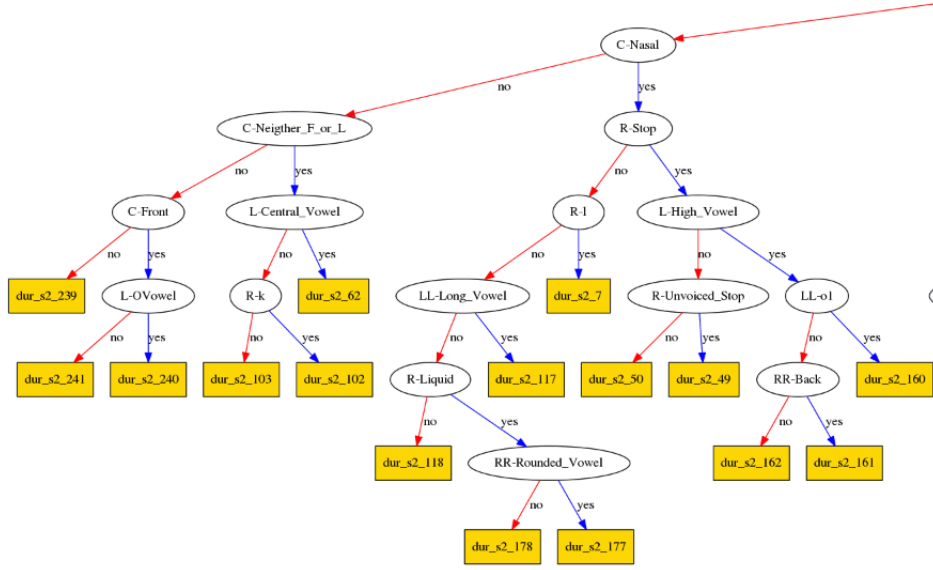


Fig. 0.3: Árbol de decisión generado para la duración de un HMM

Con este modelo, el sistema podrá inferir por ejemplo cosas como: si el fonema actual no es nasal (C-Nasal) seguido de un stop (R-Stop), que no es el fonema *l* estará modelado por función de probabilidad gaussiana definida en *dur_s2_7*.

En las primeras iteraciones del desarrollo no contábamos con la información acústica por lo que se generaron modelos carentes de información contextual. En estos primeros modelos se pudo apreciar una calidad mucho peor en los audios generados, sonando estos sumamente metálicos y carentes de prosodia. Tras agregar los factores contextuales y realizar algunas pruebas de concepto con ellas pudimos comprobar que las voces sonaban mucho mas humanas.

HTS también brinda la posibilidad de realizar speaker-adaptive training. Esta técnica permite tomar un modelo ya entrenado y adaptarlo para asimilar características de un nuevo hablante. Esta técnica nace de la idea que construir un corpus de datos es costoso tanto en espacio de almacenamiento, tiempo de grabación y etiquetado, por lo que resulta mas económico generar una nueva voz sintética a partir de un modelo bien generado y adaptándolo luego con características del nuevo corpus.

Dentro del adaptative training existen varias técnicas, en este trabajo utilizaremos offline supervised adaptation, que tiene como requisito adicional conocer los utterances del segundo corpus.

0.5.3. HTS_engine

Finalmente para generar voces con acento extranjero se utilizó *hts_engine*. Esta herramienta permite interpolar con pesos arbitrarios entre varios modelos para producir un nuevo modelo con una mezcla de la carga fonética y prosódica de ambos hablantes y sintetizar audios. Esto nos brinda un gran rango exploratorio y nos permitirá ajustar la carga fonética de los modelos originales para acercarnos al modelo deseado.

0.6. Experimentación

En la siguiente aparatado intentaremos validar dos hipótesis: que el modelo generado realmente puede ser identificado como un hablante extranjero de habla inglesa y al mismo tiempo que este posee un grado de inteligibilidad aceptable.

Para eso se condujo una encuesta perceptual donde a cada participante, se le presentó un audio con una oración semánticamente impredecible, fonéticamente balanceada y con distintos grados de mezcla de español e ingles y se le pidió que la transcribiera y que intentara identificar la nacionalidad del hablante.

La encuesta se realizó a través de internet, con el mismo set de instrucciones para todos los participantes y pidiendo como requerimiento la utilización de auriculares. Cada participante podía contestar un máximo 5 veces (otorgándoles siempre audios distintos).

La misma se llevó a cavo desde el 18 de octubre de 2017 hasta el primero de diciembre del mismo año, tiempo durante el cual fue publicada en distintas redes sociales y listas de emails de la facultad.

Los utternaces utilizados fueron generados tomando palabras de manera aleatoria de una lista de sustantivos, adjetivos y verbos y realizando correcciones donde fuera necesario para mantener el balanceo fonético.

Los utternaces utilizados fueron:

- Utternace 1: Mi montaña aguileña recorrió la esquina
- Utternace 2: Aquel fuerte vidrio prefirió aquel botón
- Utternace 3: Este enjoyado juez comprará nuestro corchete
- Utternace 4: Tu estrecho posavasos gritó la fechoría
- Utternace 5: Nuestro nublado tigre concluyó a este chupetín
- Utternace 6: Su profundo riñón apoyó a Julio
- Utternace 7: El frío churrasco oyó lo de Polonia
- Utternace 8: Las acongojadas cotorras sonrieron a mi círculo
- Utternace 9: Ese gruñón perro prometió a esos cuñados
- Utternace 10: El nudillo Argentino perdió su vaso

Para cada uno de estos diez utternaces se varió el nivel de mezcla entre 30 % de ingles - 70 % castellano hasta 70 % de ingles - 30 % de castellano, 10 % a la vez.

0.6.1. Interface

A continuación se presenta la interfase utilizada para realizar la encuesta junto con las decisiones de diseño mas relevantes tomadas a lo largo de su creación.

Al entrar en la pagina de la encuesta todos los participantes fueron presentados con la siguiente pantalla principal:

Estudio de Percepción

¡Gracias por participar!

Con este estudio queremos evaluar la calidad de distintas voces artificiales.

Es fundamental que lo hagas con auriculares y en un ambiente silencioso.

Datos Personales:

Antes de empezar, por favor completá estos datos, que usaremos sólo para generar métricas de los participantes. Tu participación es totalmente anónima y confidencial.

Edad:

Género:

Dónde pasaste la mayor parte de
tus primeros 10 años de vida:

Guardar

Con el objetivo de no influir en las respuestas de los participantes participantes, se procuró darles a los participantes la información indispensable para completar la encuesta. En ningún momento de la encuesta se especifica el objetivo del estudio.

Con la intención de estandarizar los resultados, fue requerimiento obligatorio utilizar auriculares para la encuesta. También se le pidió a cada participante que la realizara en un lugar silencioso.

Además, le pedimos a cada participante que indique su género, su edad y la provincia en la que transcurrido la mayor parte de su infancia.

Una vez que completados estos datos, se les presentaba otra vista con las instrucciones específicas para completar la encuesta:

Estudio de Percepción

Instrucciones

- Se te presentará un breve audio con alguien hablando, que dura aproximadamente 3 segundos.
- ¡Recordá utilizar auriculares!
- Tu tarea es transcribir las palabras del audio, e indicar el origen/nacionalidad del hablante.
- Esta tarea puede resultar difícil. Hacela lo mejor que puedas. Si solo lográs entender palabras sueltas, transcribirlas en el orden en que las escuchás.
- Podés escuchar el audio solamente dos veces.

Entendido!

Una vez presionado el botón de “entendido!” se les presentaba un audio, que podían escuchar un máximo de 2 veces, una caja de texto libre donde plasmar la transcripción del mismo y una caja de texto libre donde podían escribir la nacionalidad correspondiente a la voz.

Reproducir el audio

Te quedan 2 reproducciones

Transcripción:

Origen/Nacionalidad del hablante:

Guardar!

Una vez que la respuesta era guardada, se le preguntaba si quería continuar transcribiendo otro audio, o de caso de haber contestado cinco veces, se le presentaba un mensaje donde se le indicaba que ya podía cerrar la encuesta.

0.7. Resultados

A modo de introducción, comenzaremos esta sección mostrando los datos demográficos obtenidos. Mas adelante, continuaremos con un análisis mas exhaustivo de inteligibilidad y por separado se realizará otro análisis respecto a la nacionalidad atribuida a los utternaces. Por ultimo para evaluar la hipótesis original, compondremos estos dos ejes para dilucidar su grado de validez.

0.7.1. Datos demográficos

Se encuestaron 109 participantes de los cuales se obtuvieron 352 resultados.

Del total de participantes, 49 pertenecían al rango comprendido entre 18 y 25 años, 43 estaban en el rango 26-35. 17 de los participantes eran mayores a 35 años:

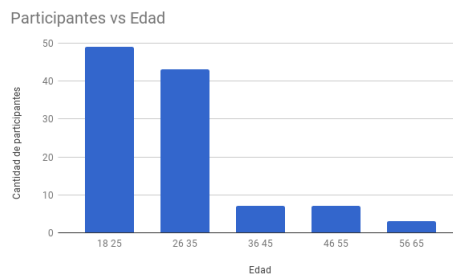


Fig. 0.4: Edad de los participantes

Con respecto al genero de los participantes, 187 respuestas fueron brindadas por participantes del genero femenino mientras que 163 respuestas fueron brindadas por participante del genero masculino.

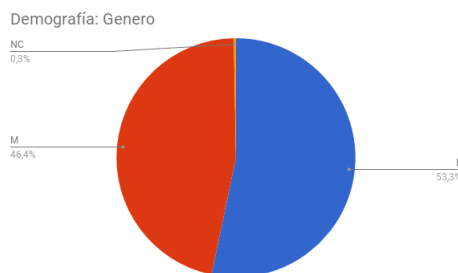


Fig. 0.5: Genero

En los datos referentes a la región en que cada participante pasó su infancia puede verse una predominancia de personas del Gran Buenos Aires con 45 %, seguido por un 30 % que pasaron su infancia en la Capital Federal. Menos del 25 % pertenece al resto de las provincias Argentinas. Además, 10 personas contestaron que se criaron fuera del país.

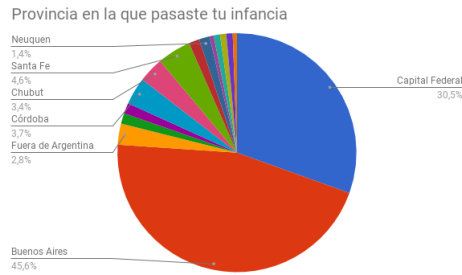


Fig. 0.6: Distribución Territorial

0.7.2. Inteligibilidad

Para el análisis de resultados utilizaremos la distancia de Levenshtein con inserciones, remociones y reemplazos. Respetando los acentos pero sin tener en cuenta mayúsculas o minúsculas.

Presentamos aquí los resultados obtenidos sin ningún tipo de modificación:

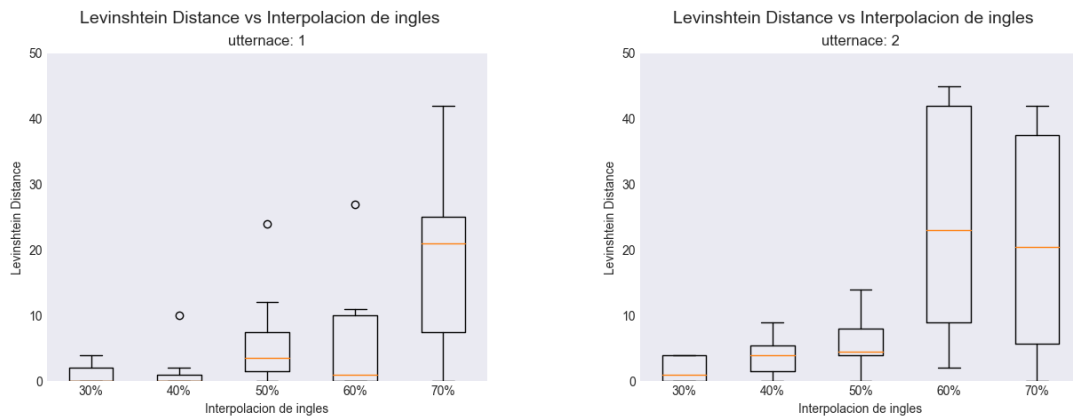


Fig. 0.7: Utternace 1 y 2

Como puede verse en la mayoría de los utternaces se puede observar que hasta el 50 % de mezcla castellano-ingles, se conserva el un buen grado de inteligibilidad, rondando una distancia que varía entre 10 y 20 caracteres. Pasados el 60 % de ingles, en la mayoría de los utternaces se observa una disminución brusca en la inteligibilidad, llegando a una distancia de 45 caracteres.

Analizando detenidamente las transcripciones obtenidas pudimos observar algunas fallas sistemáticas que podrían generar ruido en el análisis, tales como:

- Los participantes escribieron de manera diferente las secciones del utternace que no comprendieron.
 - Muchos de ellos escribieron: "...", "..." o simplemente omitieron la palabra.
 - En casos menos comunes: "****", "???", "blablabla".

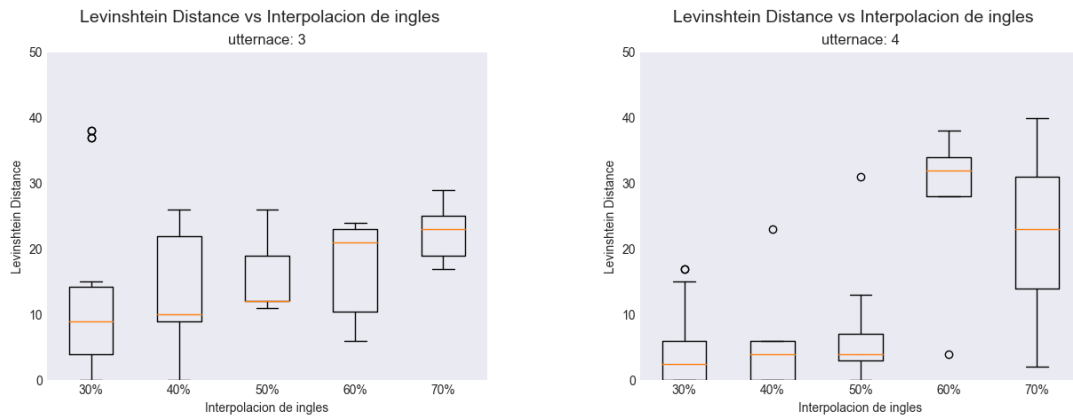


Fig. 0.8: Utternace 3 y 4

- En casos donde no comprendieron ningún segmento del utternace, es común observar expresiones como “no entendí nada”, “nada”, “”, etc.
- Utilización de signos de puntuación:
 - Puntos finales para expresar el final de la oración o expresiones como “(?)”.
 - En un caso extremo, un participante el participante transcribió “ “tu estrecho posavasos”, grito la fechoría”, cuando el Utternace original solo decía “tu estrecho portavasos gritó la fechoría”.
- Omisión de acentos y faltas ortográficas en palabras que no presentan ambigüedades. Como por ejemplo: “grunion” en vez de “gruñón”

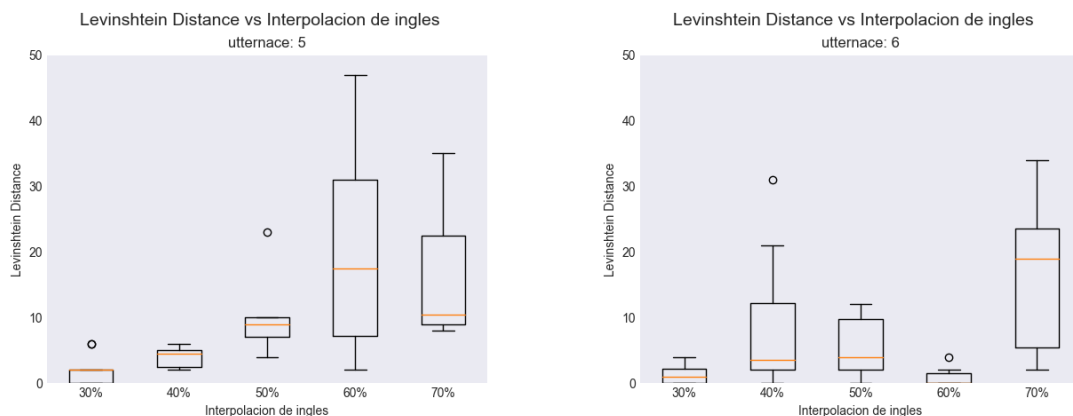


Fig. 0.9: Utternace 5 y 6

Para disminuir ruido de la muestra, se decidió realizar una limpieza de los datos donde consideramos que no era disruptiva.

Los cambios fueron:

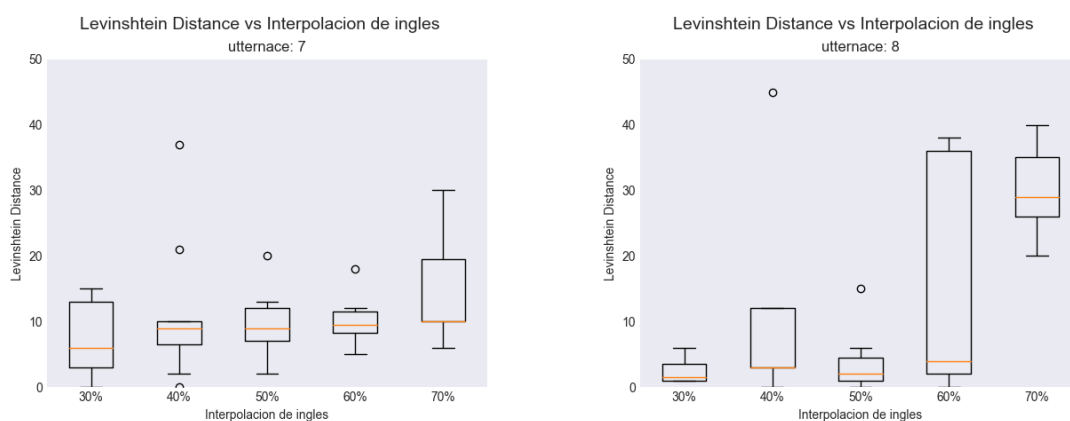


Fig. 0.10: Utternace 7 y 8

- corregir “ni” por “ñ ” en la palabra grunion.
- Remoción de todos los signos de puntuación.
- Remoción de oraciones como “blabla”, “no entendí” o cualquier otra expresión que indique ininteligibilidad de una palabra u oración.
- Corrección de acentos en palabras no ambiguas: “botón”, “prefirió”, “recorrió”, “chupetín”, “riñón”, “gruñón”.

Aquellas palabras que presentan ambivalencia, como : “concluyó” no fueron modificadas ya que concluyó/concluyo son validas. El participante podría haber interpretado la palabra con cualquiera de las dos connotaciones cambiando el significado de la interpretación.

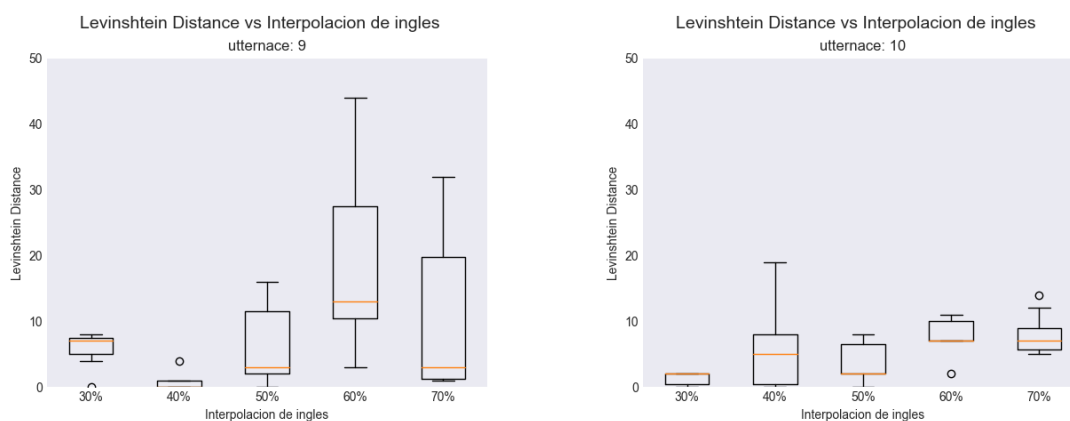


Fig. 0.11: Utternace 9 y 10

Esperamos que esta normalización no solamente ayude a disminuir la propagación de los resultados, sino que además, nos permitirá interpretar de manera mas intuitiva el significado la distancia de Levenshtein en cada caso.

0.7.3. Datos normalizados

Con esta estandarización de los datos, trataremos de darles un peso intuitivo que nos permitan sistematizar el análisis.

Por ejemplo, tomando el Utternace 8 de las frases utilizadas en la experimentación:

- “Las acongojadas cotorras sonrieron a mi círculo”

Podemos observar las siguientes transcripciones extraídas de los resultados:

- Distancia 0: “Las acongojadas cotorras sonrieron a mi círculo”
- Distancia 10: “Las acongojadas culturas sonrieron en semicírculo”
- Distancia 20: “Plaza sombreada con sombrero sonrieron en mi círculo”
- Distancia 30: “sonrieron en mi círculo”
- Distancia 40: “círculo”
- Distancia 48: “”

A partir de esto, para este apartado vamos a tomar que una distancia de Levenshtein

- Distancia 0-10: Buena Inteligibilidad
- Distancia 10-20: Mediana Inteligibilidad
- Distancia 20-40: Baja Inteligibilidad
- Distancia 40-: Inteligibilidad nula

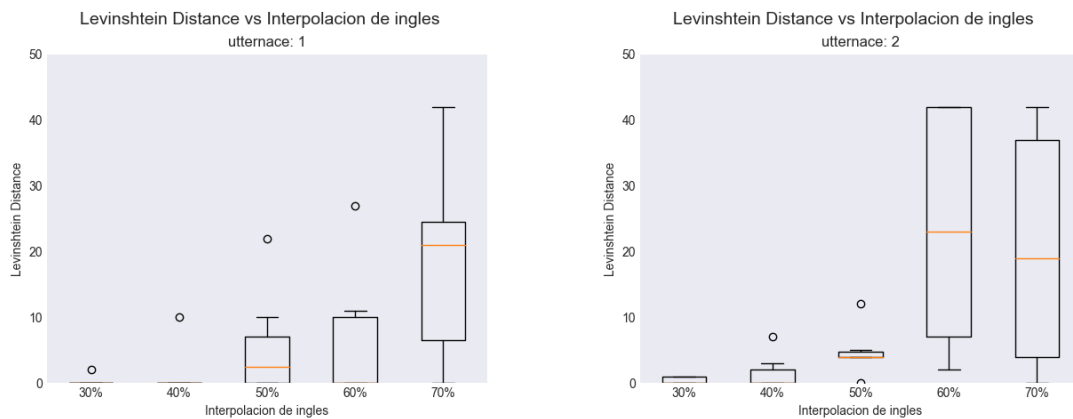


Fig. 0.12: Utternace 1 y 2 Normalizados

Con este baseline, podemos ver que para una interpolación de ingles de 30 %, 96 de los 106 participantes comprendieron de manera adecuada el texto con una inteligibilidad alta, mientras que los 10 restantes obtuvieron una inteligibilidad media.

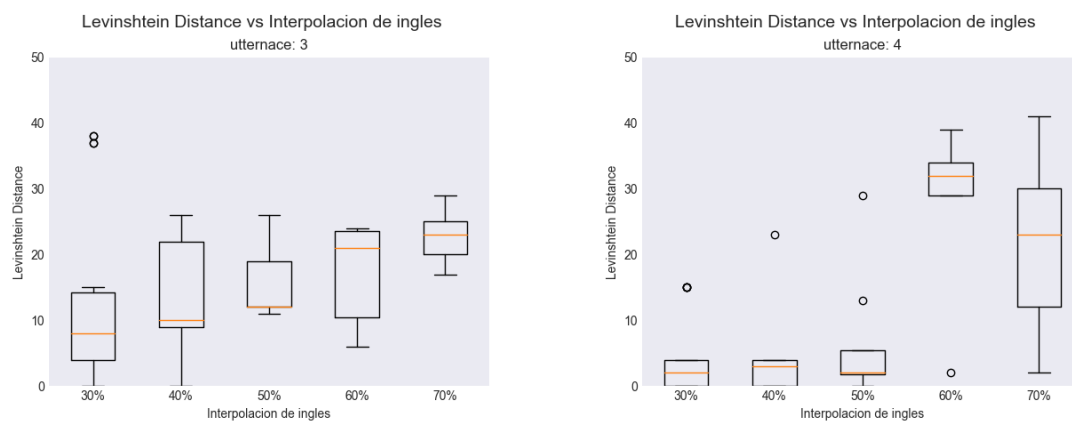


Fig. 0.13: Utternace 3 y 4 Normalizados

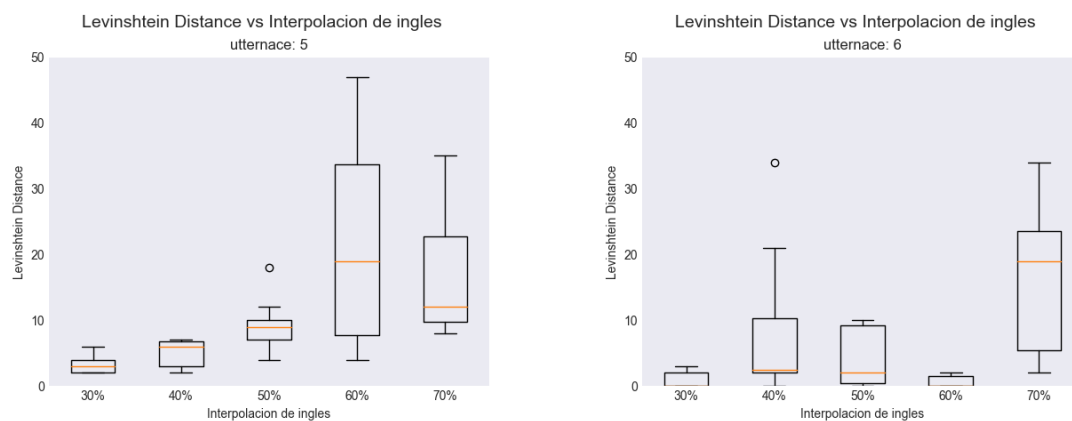


Fig. 0.14: Utternace 5 y 6 Normalizados

Para la mezcla 40 % inglés - 60 % castellano, de un total de 67 participantes, 57 obtuvieron una inteligibilidad alta, 2 una inteligibilidad media, 5 una baja y 3 una inteligibilidad nula.

Para la interpolación 50 % inglés - 50 % castellano, de un total de 75 participantes, 57 obtuvieron transcribir el audio demostrando una buena inteligibilidad, mientras que 14 tuvieron una inteligibilidad media y 4 una inteligibilidad baja.

Para todas las interpolaciones enunciadas previamente, los errores mas comunes varían desde falta de acentos en palabras como “concluyó” hasta faltas de inteligibilidad en palabras con cierta complejidad fonética como “aguileña” o “gruñón”.

Como caso particular el Utternace 3: “este enjoyado juez comprará nuestro corchete” podemos observar que la mayoría de los participantes cometieron errores al transcribir la palabra “juez” que confundieron de manera sistemática con palabras sonoramente similares como “fue”, y “enjoyado” que transcribieron como “enfollado”, “enrollado” y la conjugación exacta del verbo comprar.

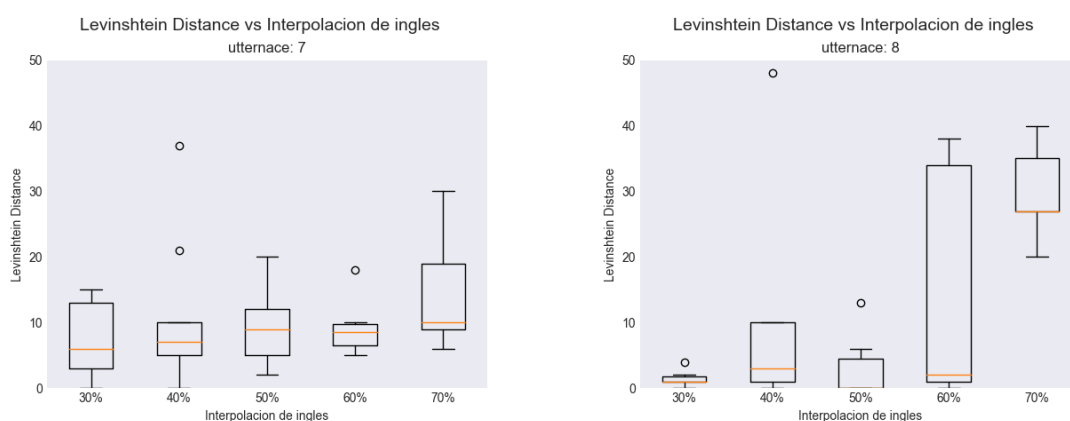


Fig. 0.15: Utternace 7 y 8 Normalizados

Para los grados de interpolación 60 % inglés - 40 % castellano y 70 % inglés - 30 % castellano, podemos observar un aumento notable de la variabilidad en las respuestas. Para el primero, de las 70 respuestas obtenidas, 40 participantes lograron transcribir con un buen grado de inteligibilidad los audios, 6 obtuvieron una inteligibilidad media, y 24 transcribieron el audio con una inteligibilidad baja o nula.

Para 70 % inglés - 30 % castellano, la diferencia es todavía mas marcada, de los 68 resultados obtenidos, 28 lograron transcribir el audio con un buen grado de inteligibilidad, 8 con un grado medio y 32 con un grado bajo o nulo de inteligibilidad.

Consideramos estos resultados tan dispares pueden deberse a dos motivos:

El primero es que existen características particulares de los participantes y sus capacidades para discernir palabras incluso cuando presentan defectos en la pronunciación del hablante. En particular, el Utternace 4 muestra como para el mismo grado de interpolación, con características similares 2 participantes de los 9 que realizaron la transcripción, obtuvieron distancias 2 y 6 en sus transcripciones.

El segundo motivo puede deberse a que existen características particulares de los utternaces o del modelo utilizado para generar la voz que afectan la comprensión del audio: el Utternace 10, donde el 6 de los 8 participantes obtuvieron una buena transcripción del

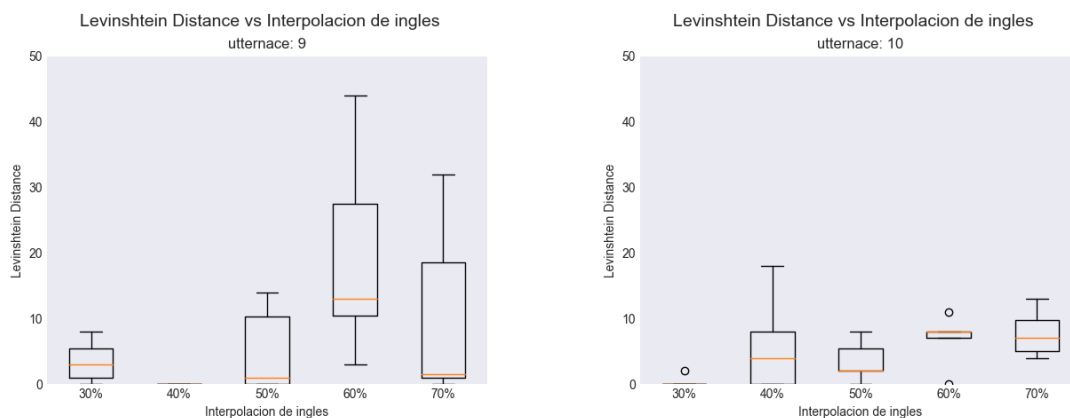


Fig. 0.16: Utternace 9 y 10 Normalizados

audio, y el Utternace 8, donde todos los participantes transcribieron el audio con inteligibilidad baja o nula, parecen demostrar esto. O bien la dificultad de los utternaces es variable o, lo que es todavía mas probable, llegado cierto punto en la interpolación, algunos fonemas empiezan a “romperse” o se alejan demasiado del fonema castellano correcto y terminan por disminuir la claridad de la voz.

En conclusión, en esta sección pudimos demostrar que fue posible generar una voz con cierto grado de inteligibilidad interpolando hasta un 50 % de ingles y 50 % de castellano. Pasado el 50 % de ingles, la variabilidad de las respuestas se vuelve mucho mas grande pudiendo haber participantes que comprenden el audio hasta algunos que no logran comprender segmentos aislados del mismo.

0.7.4. Análisis de Nacionalidad

En esta sección analizaremos los resultados de las nacionalidades que los participantes de la encuesta atribuyeron a la voz.

Dado que en esta instancia se le permitió a los participantes ingresar texto libre las respuestas resultaron bastante heterogéneas. Los participantes tomaron la consigna de manera diferente, pudiendo encontrarse respuestas que no pueden ser atribuidas a una nacionalidad. Como ejemplo de algunas respuestas pueden encontrarse cosas como: Latino, Anglo, Robot, España (sur).

Consideramos que las respuestas de la índole “robot”, “es una voz artificial”, no son validas ya que no aportan información para esta investigación.

Por esta razón, en esta instancia decidimos agrupar las respuestas en cuatro grupos lógicos:

- Hispanohablante: Latino, Argentino, Español, Uruguayo, Centroamericano, Boliviano, Mexicano, Colombiano
- Angloparlante: Estadounidense, Ingles, Irlandés, Canadiense, Anglo
- No sabe/No contesta: Robot, no se
- Otro: Ruso, Brasiltiño

Con estas agrupaciones, presentamos las nacionalidades atribuidas a la voz generada para cada punto de la interpolación.

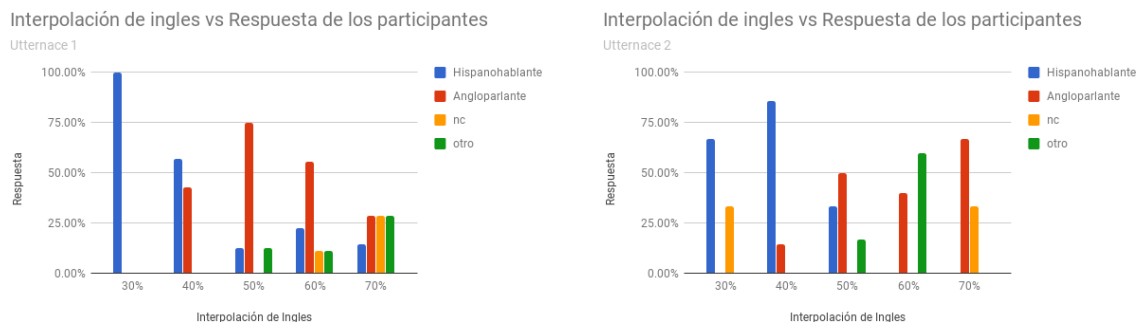


Fig. 0.17: Utternace 1 y 2

De estos resultados podemos observar que con 30 % de interpolación de ingles, los participantes coinciden en que la voz puede atribuirse a una persona de habla nativa española.

Para una mezcla de 40 % ingles, puede verse que no hay una decisión concluyente con respecto a la nacionalidad del hablante. Por ejemplo en el Utternace 3 el 80 % de los participantes coincide que la voz pertenece a un hablante de habla hispana, mientras que en el Utternace 9 el 60,00 % de los participantes considera que la voz pertenece a un hablante de habla anglosajona.

Esta gran disparidad de resultados entre distintos utternaces se puede atribuir a las características particulares de cada Utternace. En particular el Utternace 9: “Ese gruñón

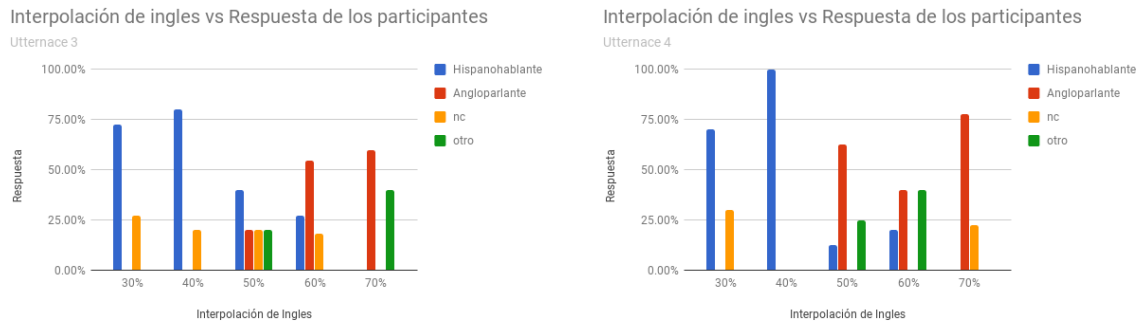


Fig. 0.18: Utternace 3 y 4

perro prometió a esos cuñados” contiene una /r/ que resulta muy notoria al pronunciarse con una intensidad menor a la esperada (mas similar a una /r/) y es atribuida, en general, a un hablante extranjero.

Bajo esta suposición observamos que los otros utternaces que presentan este fonema:

- Utternace 1: “Mi montaña aguileña recorrió la esquina”
- Utternace 6: “Su profundo riñón apoyó a Julio”
- Utternace 7: “El frío churrasco oyó lo de Polonia”
- Utternace 8: “Las acongojadas cotorras sonrieron a mi círculo”

También presentan un mayor porcentaje de atribuciones a nacionalidad anglosajona.

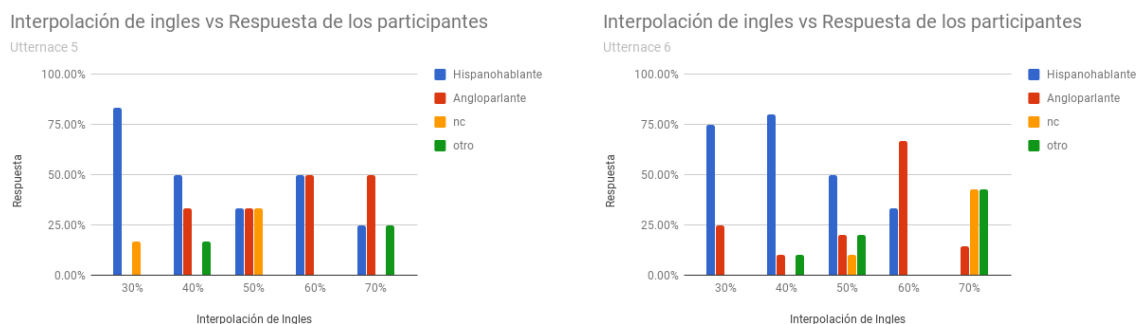


Fig. 0.19: Utternace 5 y 6

Con 50 % y 60 % de ingles los resultados son similares. obtenemos que aproximadamente en el 50 % de los utternaces, mas de la mitad de los participantes consideraron que la voz pertenecía a un anglosajón hablando castellano. Para estos grados de interpolación también podemos observar que en un 80 % de los utternaces al menos un 20 % de los participantes atribuyen la nacionalidad del hablante a un no nativo no anglosajón hablando castellano.

Con 70 % de interpolación, en el 80 % de los utternaces se puede apreciar que al menos 50 % de los participantes dijo que el hablante era de origen anglosajón. Mas aún, en el 40 %

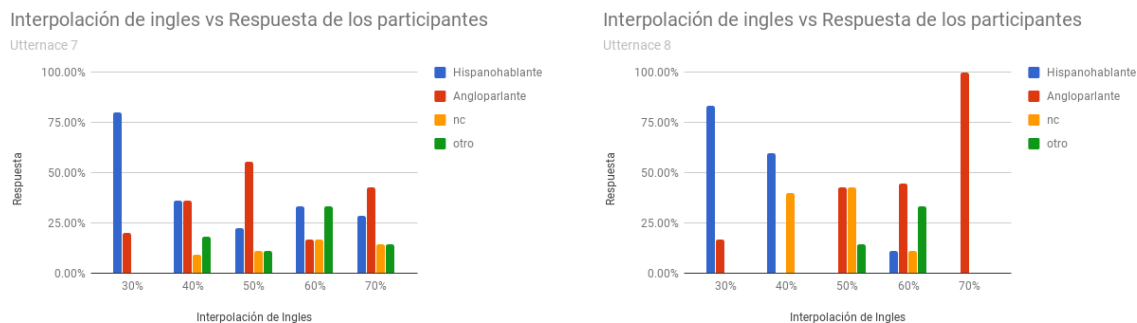


Fig. 0.20: Utternace 7 y 8

de los utternaces el 75 % de los participantes coincidió que la voz era de angloparlante. También podemos ver que para este grado de interpolación en el 70 % de los utternaces ningún participante considera que la voz sea de habla hispana. En el 30 % restante, 25 % de los participantes o menos consideran que la voz pertenezca a un hispanohablante.

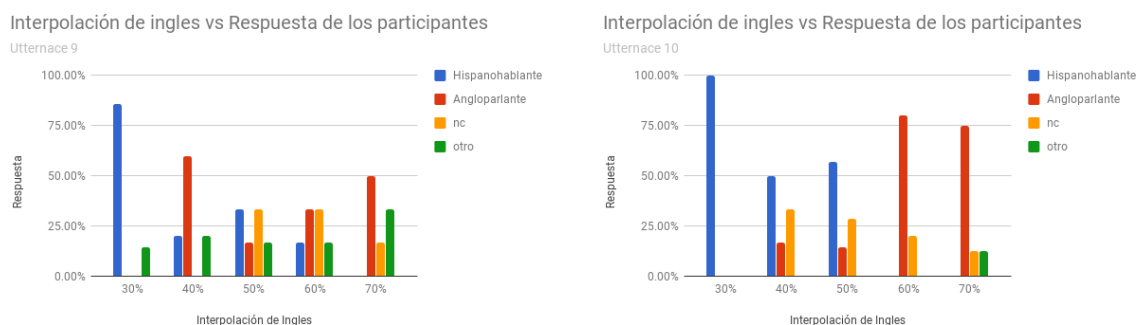


Fig. 0.21: Utternace 9 y 10

Hasta ahora analizamos los dos ejes de nuestra hipótesis por separado (por un lado, inteligibilidad, por otro, nacionalidad atribuida a la voz). En el ultimo apartado de la investigación buscaremos sacar conclusiones al componer ambos ejes en un mismo análisis.

0.7.5. Resultados Generales de la experimentación

Por ultimo visualizamos mostraremos la distancia de Levenshtein superpuesto con con la probabilidad de un participante de reconocer la voz como un hablante anglosajón.

Volviendo a la hipótesis original, podemos ver que esta técnica permite generar una voz que pueda ser identificada como un extranjero hablando ingles, con un grado de efectividad que varía desde el 60 % hasta el 100 % dependiendo del Utternace elegido y el grado de interpolación.

También es interesante observar casos como el que se presenta comparando el Utternace 10 y el Utternace 8, ambos con 70 % de mezcla de ingles, que en cuanto a inteligibilidad se encuentran en extremos opuestos, muestran que aproximadamente un 80 % y un 100 % de participantes identificaron como nativo anglosajón. Esto nos da a pensar que la inteli-

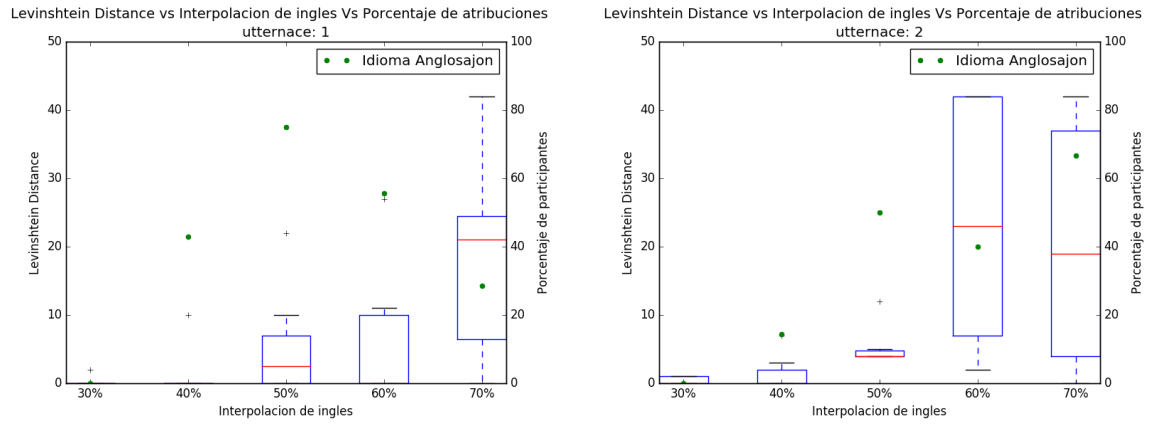


Fig. 0.22: Utternace 1 y 2

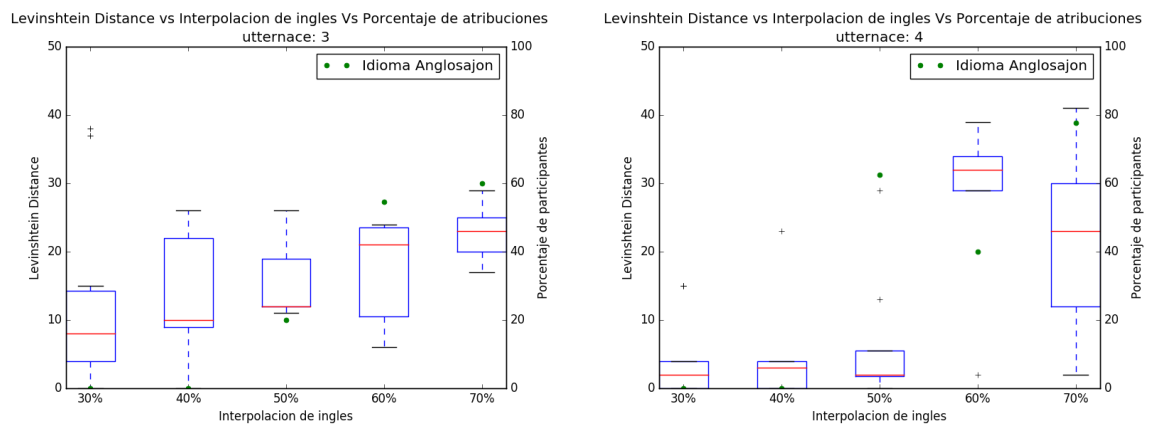


Fig. 0.23: Utternace 3 y 4

gibilidad de una oración y su probabilidad de ser identificado como un hablante ingles son variables independientes y que este ultimo factor este mas ligado a otros factores como la sonoridad de ciertos fonemas o la prosodia general de la voz.

Este no es un caso aislado, véase que lo mismo sucede con el Utternace 4 y 6 con 60 % de mezcla de ingles, si bien las inteligibilidades están en extremos opuestos, sus probabilidades de ser identificados como hablantes extranjeros difieren en menos del 20 %.

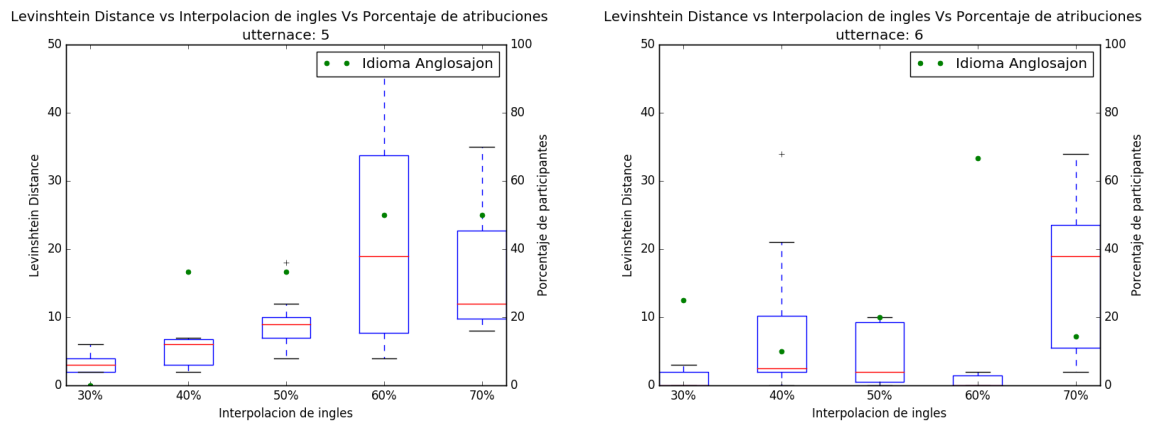


Fig. 0.24: Utternace 5 y 6

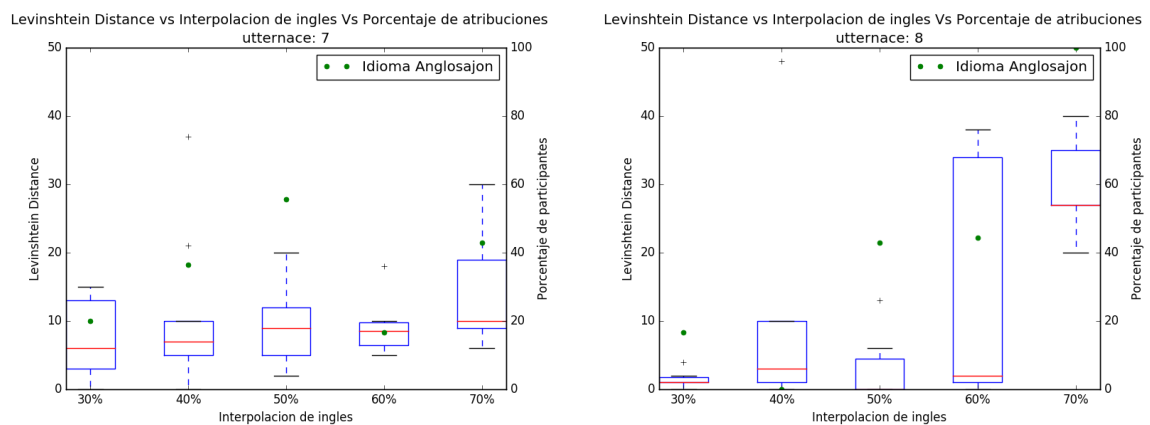


Fig. 0.25: Utternace 7 y 8

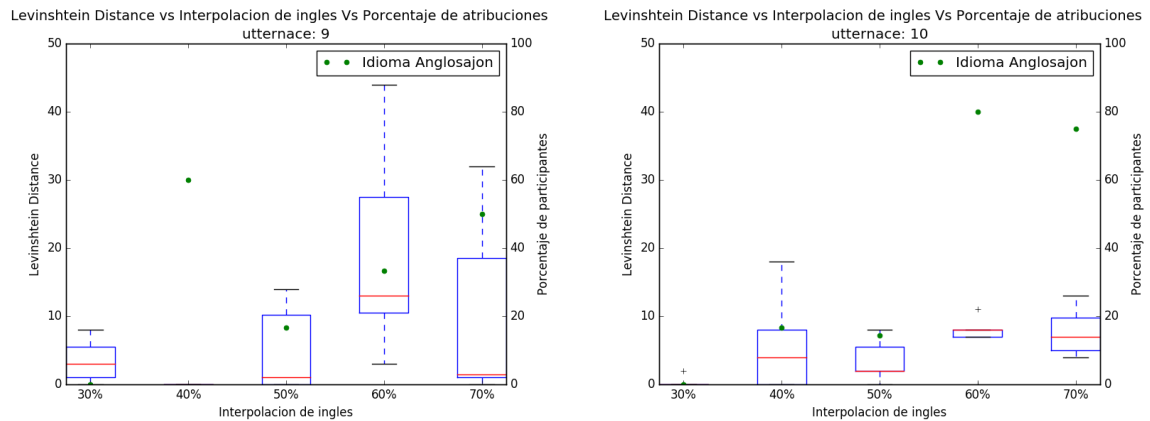


Fig. 0.26: Utternace 9 y 10

0.8. Trabajo Futuro

Como fue discutido en la sección de experimentación, una interpolación general puede producir que ciertos fonemas se alejen demasiado del fonema real del castellano, disminuyendo la inteligibilidad de la voz sintetizada. Un posible camino a seguir es realizar una interpolación controlada que permita regular cada fonema por separado. Para fonemas que puedan resultar problemáticos como el caso de la /r vibrante el grado de interpolación podría dejarse mas cercano al castellano, mientras que para fonemas con comportamientos mas similares el grado de interpolación podría llevarse mas cerca del modelo ingles.

0.9. Apendice

0.9.1. Lista de Fonemas

Castellano	Ingles	Castellano	Ingles
a	aa	p	jh
a1	ae	r	k
b	ah	rr	l
ch	ao	s	m
d	aw	t	n
e	ax	u	ng
e1	ay	u0	ow
f	b	u1	oy
g	ch	x	p
i	d	-	r
i0	s	-	sh
i1	dh	-	t
k	eh	-	th
l	er	-	uh
ll	ey	-	uw
m	f	-	v
n	g	-	w
ny	hh	-	y
o	ih	-	z
o1	iy	-	zh

Tab. 0.3: Mapeo Fonetico

0.9.2. Mapeo Fonemas del Ingles-Castellano

0.9.3. Parametros utilizados para el entrenamiento

Bibliografía

- [1] Automatic determination of phrase breaks for Argentine Spanish, Humberto M. Torres & Jorge A. Gurlekian, Laboratorio de Investigaciones Sensoriales CONICET, University of Buenos Aires, Argentina
- [2] Torres, Humberto & Gurlekian, Jorge & Cossio-Mercado, Christian. (2012). Aromo: Argentine Spanish TTS System.
- [3] Kominek, John & W Black, Alan. (2004). The CMU Arctic speech databases. SSW5-2004.
- [4] Statistical Parametric Speech Synthesis Based on Speaker and Language Factorization, Heiga Zen, Member, IEEE, Norbert Braunschweiler, Sabine Buchholz, Mark J. F. Gales, Fellow, IEEE, Kate Knill, Member, IEEE, Sacha Krstulovic, and Javier Latorre, Member, IEEE
- [5] Simultaneous Modeling of phonetic and prosodic parameters, and characteristic conversion for hmm-based text-to-speech systems, takayoshi yoshimura, pag. 28
- [6] Speaker Similarity Evaluation Of Foreign-accented Speech Synthesis Using Hmm-based Speaker Adaptation, Mirjam Wester, Reima Karhila
- [7] Simultaneous Modeling of phonetic and prosodic parameters, and characteristic conversion for hmm-based text-to-speech systems, takayoshi yoshimura
- [8] SUB-PHONETIC MODELING FOR CAPTURING PRONUNCIATION VARIATIONS FOR CONVERSATIONAL SPEECH SYNTHESES, Kishore Prahallad, Alan W Black and Ravishankhar Mosur
<https://www.cs.cmu.edu/~awb/papers/ICASSP2006/0100853.pdf>
- [9] <http://hts.sp.nitech.ac.jp/?Download>

