

1 Datasets

Escolhemos para esse trabalho 3 conjuntos de dados sendo 2 prontamente acessíveis pelo grupo FalaBrasil e o outro extraído diretamente de vídeos do YouTube. Todos os arquivos tem 22.050Hz com 16 bits.

Sigla	Nome	Tópico	Descrição Técnica
CONST	Constituição	Leitura da constituição nacional cujo texto e arquivos de áudio original foram processados pela equipe FalaBrasil de modo a adequar-se às necessidades do estudo de fala	Segmentados em aproximadamente 30 segundos, falante masculino, em ambiente de gravação de rádio. 9000 arquivos; aproximadamente 9h de gravações
LAPS	LAPS BenchMark	Corpus de voz utilizado para avaliação de desempenho de sistemas LVCSR.	700 frases, o corpus possui 35 locutores com 20 frases cada, sendo 25 homens e 10 mulheres, o que corresponde a aproximadamente 54 minutos de áudio
YT	Youtube (Cellbit)	Áudios extraídos de vídeos do YouTube de um <i>youtuber</i> no estilo <i>vlog</i> . Legendas manualmente inseridas nos vídeos pelo próprio autor dos vídeos.	23 vídeos com legenda publicados entre 04/07/2017 e 27/03/2018 totalizando aproximadamente 23h de áudio. O locutor é um youtuber do sexo masculino com 21 anos. Os arquivos foram seccionados baseado nos tempos de cada linha inteira do arquivo de legenda. 4203 arquivos ($2.9 \pm 4.2s$).

Table 1: Descrição dos Datasets

2 Questões

Queremos levantar como os atuais algoritmos de síntese neural se comportam perceptualmente na língua portuguesa com restrições de dados e de tempo na geração do modelo. Queremos levantar comparativamente se um modelo treinado com dados anteriores e que tenha apresentado bons resultados perceptuais pode ser benéfico para obter a síntese na língua portuguesa de maneira mais rápida, ou seja, com menor esforço computacional aproveitando o conhecimento já obtido do outro idioma.

A síntese de fala já é um problema com vários produtos disponíveis para uso pelo público [Mic, Goo] presente no cotidiano de muitos através das mais diversas aplicações mas esses modelos costumam permitir a síntese de apenas um falante e, muitas vezes, são modelos pagos e proprietários. Podemos obter a síntese de fala com os atuais algoritmos do estado da arte no português para um falante qualquer? Essa mesma síntese pode ser obtida com resultados satisfatórios com restrições de tempo e de dados em seu treinamento?

Podemos nos utilizar de *transfer learning* com um modelo pré-treinado no inglês para obter uma convergência mais rápida dos modelos em português com naturalidade similar percebida pelos ouvintes comparativamente a um modelo treinado exclusivamente com o português?

Conseguimos utilizar uma fração dos dados no treinamento com o modelo pré-treinado para obter um menor tempo de treinamento e uma síntese com naturalidade consistente com o mesmo modelo treinado com todos os dados em um período maior de tempo?

3 Métodos

Para responder aos questionamentos acima propomos os seguintes passos:

Tomar a implementação¹ do modelo Tacotron [SPW⁺17] reconhecido como estado da arte. O primeiro

¹<https://github.com/carpdm20/multi-speaker-tacotron-tensorflow>

passo para estabelecer um paralelo comparativo treinando o modelo completamente a partir de dados em português usando os mesmos parâmetros que a publicação original (arquitetura de rede, otimizador, taxa de aprendizagem, entre outros) e encerrando o treino no mesmo número de épocas que o modelo pré-treinado.

Tendo esse modelo como comparativo podemos então treinar modelos com restrições de tempo e de dados. Para as restrições de tempo determinamos fixar uma quantidade de épocas de treino em 1k, 5k, 33k e 50k épocas. Como nosso intuito é comparar a performance do *transfer learning* no auxílio do treino inicial do modelo não nos é interessante propor épocas muito maiores que nosso modelo de referência. Como nosso modelo de referência é treinado em 100k épocas estabelecemos 50k como faixa de corte. Para os parâmetros de dimensão dos dados propomos dividir os dados de entrada em porcentagens de suas totalidades, mais especificamente 5, 15, 33 e 50% dos dados, a serem escolhidos aleatoriamente dentro de cada conjunto de dados selecionado. Com esses parâmetros queremos gerar um total de 16 modelos para cada conjunto de dados. Com essas 16 hiperparametrizações fixadas para cada um dos 3 datasets descritos na Tabela 1 temos um total de 48 modelos.

Quantitativamente os modelos são analisados a partir da variação do valor do erro no tempo e o tempo proporcional para atingir aquela quantidade de épocas desejada. Qualitativamente desejamos estabelecer um teste com ouvintes humanos através de um teste de MOS (*Median Opinion Score*) conforme a literatura. Como a quantidade de modelos a ser testada é grande utilizamos uma métrica proveniente da fonoaudiologia para filtrar os modelos mais naturais de modo a oferecer apenas os melhores modelos para os ouvintes finais. O intuito é estabelecer um teste que não cause exaustão ou desistência dos ouvintes de modo a obter resultados consistentes. Além disso podemos estabelecer um indicativo da utilização dessa métrica de naturalidade estabelecida pela fonoaudiologia como futuro parâmetro quantitativo para estudos similares.

Na literatura de fonoaudiologia recomenda-se utilizar as frases estabelecidas no CAPE-V para avaliar a naturalidade da fala. O CAPE-V é composto de 6 frases a serem pronunciadas. Com essas frases e o auxílio do software Pratt é possível obter métricas quanto ao formato da onda, amplitude

Com essas proporções para os três conjuntos de dados listados pretendemos concluir quanto a variabilidade mínima necessária de um mesmo falante e também como o comprimento do conjunto em média afeta a naturalidade final sintetizada.

As métricas quantitativas do treino serão:

- O tempo de treino de cada modelo sob as mesmas condições
- O valor da função de perda na última iteração e no menor valor alcançado durante o treino avaliado simultaneamente no conjunto de treino no conjunto de validação

3.1 Proposta de Sentenças 1

Finalizados os treinos dos modelos e podendo-se finalmente estabelecer a síntese selecionamos 3 sentenças completas de cada conjunto de dados presentes nos arquivos originais para permitir um conjunto de controle. Selecionamos ainda 3 trava línguas pela dificuldade usual que eles representariam para um falante humano e pela presença de termos inexistentes nos conjuntos de treino². Conforme o teste nas publicações originais o ouvinte será orientado a ouvir o áudio pelo menos 2x antes de determinar em qual dos 5 níveis de naturalidade ele se encaixa.

Totalizando assim 12 sentenças a serem sintetizadas por 5 modelos, totalizando 60 sentenças sintetizadas e 9 sentenças originais. A partir desses dados a proposta é estabelecer um comparativo da evolução do MOS médio à medida que adicionamos frações dos datasets originais. A partir de um teste de significância estatística desejo determinar o nível de significância com a qual pode-se afirmar que houve alguma melhoria nas médias entre os modelos.

Com o teste podemos afirmar então se as médias do MOS score de fato crescem e se crescem qual é esse nível de significância. A partir dessa informação também somos capazes de apontar os limites de variação mínimos do conjunto de dados para a evolução com significância da média perceptual dos ouvintes. Podemos por fim apontar também na dimensão do tempo o impacto do corte no conjunto de dados e seu respectivo impacto no resultado do MOS de naturalidade do modelo final.

²<http://www.fonoaudiologia.med.br/voz/7-teste-sua-diccao>

3.2 Proposta de Sentenças 2

As referências originais se usam das frases de Harvard extraídas do apêndice do: IEEE Subcommittee on Subjective Measurements IEEE Recommended Practices for Speech Quality Measurements de 1969³. Entretanto essas sentenças foram pensadas de maneira anglocêntrica e podem não ser adequadas ao português. Seguindo a literatura de fonoaudiologia recomenda-se usar as frases do CAPE-V para detecção de problemas de fala^{4 5 6} totalizando assim 6 frases. Podemos então avaliar a síntese através da análise AVQI 03.01 proposta na literatura de fonoaudiologia para detecção de desvios e vícios de fala⁷. Esse método se utiliza do Software Praat que é vastamente utilizado no campo de fala para detecção de parâmetros técnicos, como os necessários para equação conforme podemos ver abaixo

$$AVQI_{03.01} = (4.152 - 0.177 * CPPs - (0.006 * HNR) - (0.037 * Shim) + (0.941 * ShdB) + 0.01 * Slope + (0.093 * Tilt)) * 2.8902 \quad (1)$$

Para consolidar então os resultados obtidos do AVQI podemos executar uma pesquisa buscando obter o MOS dessas 30 sínteses (6 frases * 5 modelos). Com os resultados desse questionário podemos validar a correlação do MOS com o AVQI e correlacionar o volume de dados disponível a nota do AVQI e ao valor atribuído de média de MOS de Naturalidade com seus respectivos níveis de significância.

References

- [Goo] Google. Cloud text-to-speech. <https://cloud.google.com/text-to-speech/>. acessado em 03/2019.
- [Mic] Microsoft. Text to speech. <https://azure.microsoft.com/en-us/services/cognitive-services/text-to-speech/>. acessado em 03/2019.
- [SPW⁺17] Jonathan Shen, Ruoming Pang, Ron J. Weiss, Mike Schuster, Navdeep Jaitly, Zongheng Yang, Zhifeng Chen, Yu Zhang, Yuxuan Wang, R. J. Skerry-Ryan, Rif A. Saurous, Yannis Agiomyriannakis e Yonghui Wu. Natural TTS synthesis by conditioning wavenet on mel spectrogram predictions. *CoRR*, abs/1712.05884, 2017.

³<http://www.cs.columbia.edu/hgs/audio/harvard.html>

⁴https://www.pucsp.br/laborvox/dicas_pesquisa/downloads/CAPEV.pdf

⁵<https://www.asha.org/uploadedFiles/members/divs/D3CAPEVprocedures.pdf>

⁶http://www.scielo.br/scielo.php?script=sci_arttext&pid=S2317-17822019000100303&lng=en&nrm=iso&tlng=en

⁷<https://journals.sagepub.com/doi/abs/10.1177/0003489416636131?journalCode=aora>