

215 **A Appendix: SoftStep Regression algorithm**

The SoftStep regression algorithm is built on top of neighborhood component analysis [Goldberger et al. [2004]]. Given a similarity measure sim an embedded sample z^* and embedded neighbors Z_N we predict

$$\hat{y} = \text{SoftMax}(sim(z^*, Z_N) + \ln(\text{SoftStepPred}(z^*, Z_N)))$$

216 where sim is a vector of similarities between z^* and the members of Z_N and SoftStepPred is the module described in Algorithm 1

Algorithm 1 SoftStep for prediction module

```

1: procedure SOFTSTEP( $Z, Z_N, \text{SoftStep\_fn}$ )
2:   Initialization (run once at module construction):
3:    $\text{params} \leftarrow \text{MLP or Linear layer with sigmoid activation}$ 
4:   Store  $\text{SoftStep\_fn}$ 
5:   Forward Pass:
6:    $(a_0, b_0, t) \leftarrow \text{params}(Z)$ 
7:    $sim \leftarrow \text{SIM}(Z, Z_N)$ 
8:    $sim\_norm \leftarrow \text{NORM}(sim)$ 
9:   if  $\text{SoftStep\_fn}$  requires  $a$  then
10:    if  $Z_N == Z$  (training) then
11:       $top\_sim \leftarrow \text{row-wise max of } sim\_norm \text{ excluding diagonal}$ 
12:    else
13:       $top\_sim \leftarrow \text{row-wise max of } sim\_norm$ 
14:    end if
15:     $a \leftarrow \min(a_0, top\_sim) - \epsilon$  ▷  $\epsilon > 0$  small
16:     $b \leftarrow a + b_0 \cdot (1 - a)$ 
17:   else
18:      $a \leftarrow a_0$ 
19:      $b \leftarrow b_0$ 
20:   end if
21:    $shift \leftarrow \text{SoftStep\_fn}(sim\_norm, a, b, t)$ 
22:   return  $sim + \log(shift)$ 
23: end procedure

```

217

218 **B SoftStep prediction experiments**

219 See Table 2 for results.

220 Exact model versions and pretrained weights are specified in the included GitHub repository. We
221 ensured that at least two distinct feature extractors were chosen per unstructured modality (text, audio,
222 and image) to demonstrate the generalizability of our proposed algorithm. The following is a list of
223 datasets and feature extractors used to conduct our experiments:

224 **RSNA Bone Age Prediction** The Radiological Society of North America (RSNA) Pediatric Bone
225 Age Machine Learning Challenge collected pediatric hand radiographs labeled with the age of the
226 subject in months [Halabi et al. [2019]]. We collected 14,036 images from this dataset. Images were
227 resized to 224x224 pixels, normalized with mean and standard deviation of 0.5 across the single
228 gray-scale channel and input to ResNet-18 pretrained on ImageNet [He et al. [2016]], [He and Jiang
229 [2021]].¹

230 **ADReSSo** The Alzheimer’s Dementia Recognition through Spontaneous Speech only (ADReSSo)
231 diagnosis dataset has 237 audio recordings of participants undergoing the Cookie Thief cognitive
232 assessment labeled with their score on the Mini Mental State Exam [Luz et al. [2021]]. Transcripts of
233 these recordings were tokenized and input to DistilBERT-base-uncased [Sanh [2019]], [Zolnoori et al.
234 [2023]].²

¹https://pytorch.org/hub/pytorch_vision_resnet/

²https://huggingface.co/docs/transformers/en/model_doc/distilbert

Dataset	Linear	SoftStep
RSNA [Halabi et al. [2019]]	5.12 ± 0.608	4.13 ± 0.235
MedSegBench [Kuş and Aydin [2024]]	6.69 ± 1.29	4.02 ± 0.608
ADReSSo [Luz et al. [2021]]	96.2 ± 33.3	28.0 ± 4.13
CoughVid [Orlandic et al. [2021]]	39.7 ± 27.0	21.2 ± 0.195
NoseMic [Butkow et al. [2024]]	7.96 ± 0.947	6.01 ± 0.472
Udacity [Du et al. [2019]]	0.435 ± 0.170	0.358 ± 0.104
Pitchfork [Pinter et al. [2020]]	69.6 ± 167.0	10.2 ± 1.01
Houses [Ahmed and Moustafa [2016]]	303 ± 81.3	13.2 ± 1.77
Books (see below)	8.33 ± 0.680	7.48 ± 0.588
Austin (see below)	2.29 ± 0.216	2.05 ± 0.231

Table 2: Average mean squared error (MSE) \pm standard deviation across ten random splits of each dataset. The best mean results for each dataset are shown in bold. All values are scaled by 10^3 for readability. A complete description of each dataset, including preprocessing pipelines and feature extractors, is provided.

- 235 **MedSegBench** The MedSegBench BriFiSegMSBench dataset is comprised of single-channel
 236 microscopy images and corresponding segmentation masks [Kuş and Aydin [2024]]. We estimated the
 237 size of segmentation mask areas using EfficientNet trained on ImageNet [Rizk et al. [2014]], [Tan and
 238 Le [2019]]³.
 239 **CoughVid** The CoughVid dataset provides over 25,000 crowdsourced cough recordings, with 6,250
 240 recordings labeled with participant age in years [Orlandic et al. [2021]]. One second of cough audio
 241 was input to HuBERT pretrained on LibriSpeech and mean-pooled across time [Hsu et al. [2021]],
 242 [Feng et al. [2024]]⁴.
 243 **NoseMic** The NoseMic dataset collected 1,297 30-second audio recordings of heart rate-induced
 244 sounds in the ear canal using an in-ear microphone under several activities [Butkow et al. [2024]].
 245 Audio clips were denoised⁵, encoded with the Whisper tiny audio encoder [Radford et al. [2023]], and
 246 mean-pooled across time⁶.
 247 **Udacity** The Udacity self-driving car dataset is comprised of dashcam videos labeled with the angle
 248 of the car’s steering wheel [Du et al. [2019]]. Videos were downsampled to 4 frames per second for a
 249 total of 6,762 images and individual frames were input to EfficientNet trained on ImageNet [Tan and
 250 Le [2019]]³.
 251 **Pitchfork** 24,649 reviews from the website Pitchfork were collected [Pinter et al. [2020]], where
 252 albums are scored from 0 to 10 in 0.1 increments. 1,500 randomly selected reviews were tokenized
 253 and input to BERT base [Warner et al. [2024]]⁷.
 254 **Houses** The Houses dataset collects 535 curbside images of houses as well as their log-scaled
 255 list price [Ahmed and Moustafa [2016]]. Images were resized to 256x256 pixels, center cropped to
 256 224x224, ImageNet normalized, and input to ResNet-34 [He et al. [2016]]⁸.
 257 **Books** The MachineHack Book Price Prediction dataset collated 6237 synopses of books labeled
 258 with their log-normalized price⁹. Synopses were tokenized and input to DistilBERT-base-uncased.¹⁰
 259 **Austin** The Kaggle Austin Housing Prices dataset collects over 15000 descriptions of homes labeled
 260 with their log-scaled list price⁹. Descriptions were tokenized and input to DistilBERT-base-uncased.¹⁰

³<https://docs.pytorch.org/vision/main/models/efficientnet.html>

⁴<https://huggingface.co/facebook/hubert-base-1s960>

⁵<https://pypi.org/project/noisereduce>

⁶<https://huggingface.co/openai/whisper-tiny>

⁷<https://huggingface.co/google-bert/bert-base-uncased>

⁸https://machinehack.com/hackathons/predict_the_price_of_books/overview

⁹<https://www.kaggle.com/datasets/ericpierce/austinhousingprices>