

---

# Semi-supervised classification by reaching consensus among modalities

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We extend the Consensus Network [1] framework to Transductive Consensus  
2 Network (TCN), a semi-supervised multi-modal classification framework, and  
3 identify its two mechanisms: consensus and classification. By putting forward  
4 three variants as ablation studies, we show both mechanisms should be functioning  
5 together. Overall, TCNs outperform or align with the best benchmark algorithms  
6 when only 20 to 200 labeled data points are available.

## 7 1 Introduction

8 Traditionally, semi-supervised learning and multi-view learning are applied to increase data usage  
9 efficiency. On one hand, semi-supervised learning have shown good applications. TSVM [2]  
10 regularizes the decision boundaries using unlabeled data, Ladder [3] utilizes cascading autoencoder  
11 structures, and Categorical GAN [4] incorporates information theoretic optimization goals. On the  
12 other hand, multi-view learning distills information contained in multiple modalities. Co-training [5]  
13 and tri-training [6] directly sets up classifiers to supervise each other. PVAE [7] and SemiMVAE [8]  
14 set up variational autoencoder losses between modalities. Specifically, Consensus Networks [1] use  
15 adversarial training [9] that learns modality-invariant representations, and outperformed traditional  
16 algorithms on detecting cognitive impairments.

17 However, Consensus Networks are supervised learning algorithms, hence are limited by the avail-  
18 ability of labeled data. This motivates us to push it forward to semi-supervised regime, resulting  
19 in Transductive Consensus Networks (TCN). TCNs function in two mechanisms, which we call  
20 the consensus mechanism and the classification mechanism. We put forward several variants in  
21 ablation study manner to study the roles of these two mechanisms, and show that the existence of both  
22 mechanisms are crucial to good performance of TCNs. Overall, TCN accuracies are better than or  
23 align with those of benchmark algorithms (semi-supervised or supervised, multi-modal or uni-modal)  
24 on Bank Marketing and DementiaBank datasets, when 20-200 labeled data points are available.

## 25 2 Models

### 26 2.1 Consensus Networks

27 We first briefly review the CN framework [1] for supervised, multi-view classification. Consider a  
28 dataset,  $\{\mathbf{x}^{(i)}, y^{(i)}\} (\mathbf{x}^{(i)} \in \mathcal{X})$ , where each data point  $\mathbf{x}$  is composed of feature values from multiple  
29 modalities (i.e., 'views'). If  $M$  is the total number of modalities, then  $\mathbf{x} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_M]$ . For  
30  $m = 1, \dots, M$ ,  $\mathbf{x}_m$  could have different dimensions, but the dimension of  $\mathbf{x}_m^{(i)}$  is consistent throughout  
31 the dataset. E.g., there may be 200 acoustic features and 100 semantic features for a data point, but  
32 all data points are constrained to those dimensions.

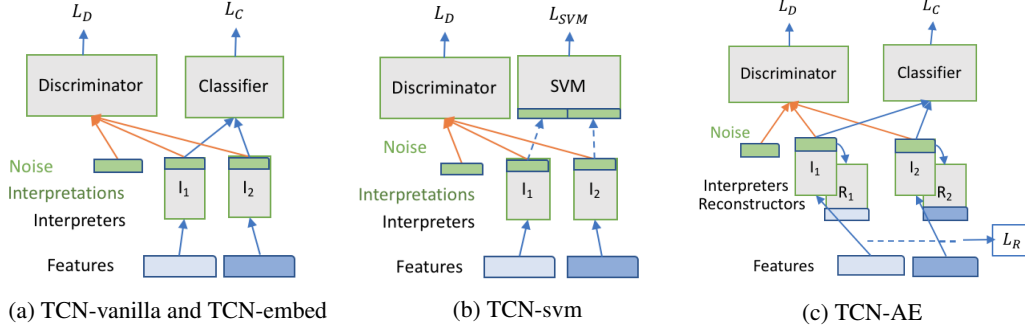


Figure 1: Network structures for TCN and TCN-embed (left), TCN-svm (middle) and TCN-AE (right), taking the example when there are two modalities in data ( $M=2$ ).

There are  $M$  interpreter networks  $I_m$  ( $m = 1, \dots, M$ ), each compressing one modality of features into a representation, which we call **consensus interpretation vector**. A discriminator  $D$  tries to distinguish the origin of each latent representation. A classifier  $C$  makes predictions based on all representations.

$$\mathbf{v}_m = I_m(\mathbf{x}_m) \quad P(\hat{m}) = D(\mathbf{v}_m) \quad P(\hat{y}) = C([\mathbf{v}_1, \dots, \mathbf{v}_m])$$

33 The training is done by iteratively optimizing two targets:

$$\min_{C, I} \mathcal{L}_C \text{ and } \min_D \max_I \mathcal{L}_D, \text{ where} \quad (1)$$

$$\mathcal{L}_C = \mathbb{E}_{\mathbf{x}}[-\log P(y|\mathbf{x})] \quad \mathcal{L}_D = \mathbb{E}_{\mathbf{x}} \mathbb{E}_m[-\log P(\hat{m} = m|\mathbf{v}_m)]$$

34 Note that empirically, an additional noise modality  $\mathbf{v}_0 \sim \mathcal{N}(\mu_{1..M}, \sigma_{1..M}^2)$  is injected to enhance the  
35 ability of the discriminator.

## 36 2.2 Transductive Consensus Networks

37 In this paper we extend CN to TCN. Formally, the input data include those labeled,  $\{\mathbf{x}^{(i)}, y^{(i)}\}$   
38 ( $\mathbf{x}^{(i)} \in \mathcal{X}_L$ ), and unlabeled,  $\{\mathbf{x}^{(i)}\}$  ( $\mathbf{x}^{(i)} \in \mathcal{X}_U$ ). In the semi-supervised learning setting, there  
39 can be a lot more unlabeled data points than labeled:  $|\mathcal{X}_U| \gg |\mathcal{X}_L|$ , where the whole dataset is  
40  $\mathcal{X} = \mathcal{X}_L \cup \mathcal{X}_U$ .

41 Here each data point  $\mathbf{x}$  contains feature values from multiple modalities (i.e., ‘views’), and the  
42 interpreter networks  $I_m$  ( $m = 1, \dots, M$ ), discriminator  $D$  and classifier  $C$  are set up identical to CN  
43 as well. Different from CN, the classification loss is defined on only those labeled data, while the  
44 discriminator loss is defined across both labeled and unlabeled data:

$$\mathcal{L}_C = \mathbb{E}_{\mathbf{x} \in \mathcal{X}_L} \text{ and } \mathcal{L}_D = \mathbb{E}_{\mathbf{x} \in \mathcal{X}} \mathbb{E}_m[-\log P(\hat{m} = m|\mathbf{v}_m)] \quad (2)$$

## 45 2.3 TCN variants

46 TCNs function in two mechanisms: The consensus mechanism compresses each data sample into  
47 "consensus interpretations", and the classifier mechanism tries to make these interpretations meaning-  
48 ful. To perform ablation studies on these mechanisms, we test the following three variants.

49 **TCN-embed** consists of the same networks as TCN but  $N_p = 30$  iterations of  $\min_D \max_I \mathcal{L}_D$  are  
50 carried out before the iterative optimizations (1). TCN-embed enhances consensus mechanism, yet  
51 allowing both mechanisms to cooperate. TCN-embed results align with TCN.

52 **TCN-svm** removes the classifier network from TCN-embed. After the pretraining phase across the  
53 whole dataset, we extract the consensus interpretations of those labeled data samples to train an SVM.  
54 TCN-svm lets the consensus mechanism to function alone, resulting in almost trivial classifiers.

55 **TCN-AE** contains an additional reconstructor network per modality  $R_{1..M}(\cdot)$ , each recovering  
 56 the input modality from latent interpretations:  $\hat{\mathbf{x}}_{\mathbf{m}} = R_{\mathbf{m}}(\mathbf{v}_{\mathbf{m}} + \epsilon)$  Defining reconstruction loss as  
 57  $\mathcal{L}_{\mathcal{R}} = \mathbb{E}_{x \in \mathcal{X}} \mathbb{E}_{\mathbf{m}} |\hat{\mathbf{x}}_{\mathbf{m}} - \mathbf{x}_{\mathbf{m}}|^2$ , the optimization target in TCN-AE can be expressed as:

$$\min_{C, I_{1..M}, R_{1..M}} \mathcal{L}_{\mathcal{C}}, \text{ and } \max_{I_{1..M}} \min_D \mathcal{L}_{\mathcal{D}}, \text{ and } \min_{I_{1..M}, R_{1..M}} \mathcal{L}_{\mathcal{R}} \quad (3)$$

58 As shown in Figure 3, TCN-AE has inferior performances than TCN. Reconstruction in an autoen-  
 59 coder style counteracts the consensus mechanisms, and should *not* be used with CN models.

### 60 3 Experiments and Results

61 We run experiments on two classification datasets, DementiaBank [10] and Bank Marketing [11].

Dataset	N. of samples	%pos / %neg	N. features per modality
Bank Marketing ('BM')	9640	48.13 / 51.87	10 / 22 / 12
DementiaBank ('DB')	473	50.76 / 49.26	185 / 117 / 110

Table 1: In BM, the three modalities correspond to basic information, statistical data, and employment. In DB, the three modalities correspond to acoustic, syntactic-semantic, and lexical.

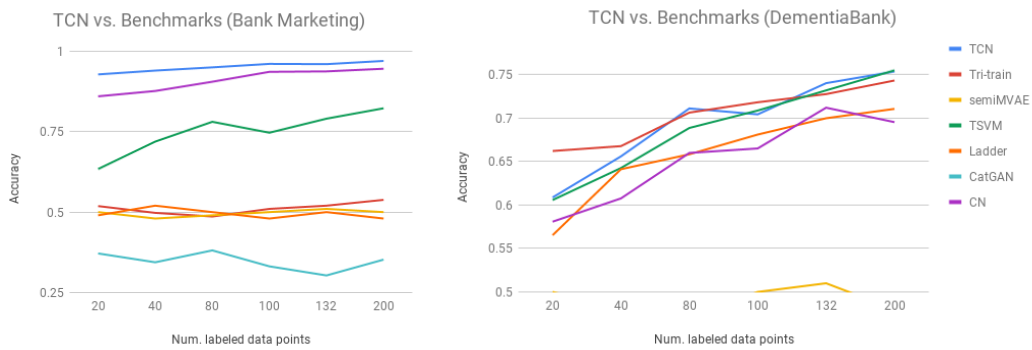


Figure 2: TCN (top blue line) outperforms or aligns with benchmark algorithms, including multi-modal semi-supervised (tri-train [6]), uni-modal semi-supervised (TSVM [2], Ladder [3], CatGAN [4]), and multi-modal supervised (CN [1]).

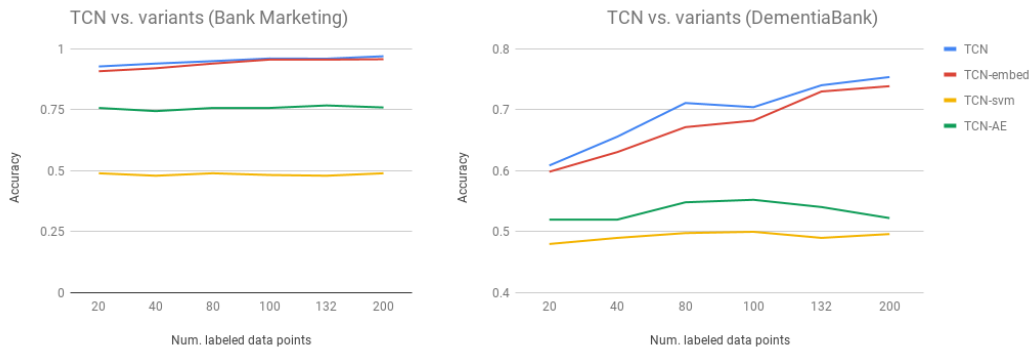


Figure 3: Accuracy plots for TCN vs its variants. Best viewed in colors. TCN-embed accuracies aligns with TCN, both significantly outperforming TCN-AE, which is better than TCN-svm. The consensus and classification mechanisms should both be present.

## 62 References

- 63 [1] Zining Zhu, Jekaterina Novikova, and Frank Rudzicz. Detecting cognitive impairments by  
64 agreeing on interpretations on linguistic features. *arxiv 1808.06570*, 2018.
- 65 [2] Thorsten Joachims. Transductive inference for text classification using support vector ma-  
66 chines. In *Proceedings of the 16th International Conference of Machine Learning (ICML-99)*,  
67 volume 99, pages 200–209, 1999.
- 68 [3] Antti Rasmus, Mathias Berglund, Mikko Honkala, Harri Valpola, and Tapani Raiko. Semi-  
69 supervised learning with ladder networks. In *Proceedings of Advances in Neural Information  
70 Processing Systems*, pages 3546–3554, 2015.
- 71 [4] Jost Tobias Springenberg. Unsupervised and semi-supervised learning with categorical genera-  
72 tive adversarial networks. *Proceedings of International Conference on Learning Representations  
73 (ICLR)*, 2016.
- 74 [5] Avrim Blum and Tom Mitchell. Combining labeled and unlabeled data with co-training. In  
75 *Proceedings of the eleventh annual conference on Computational learning theory*, pages 92–100.  
76 ACM, 1998.
- 77 [6] Zhi-Hua Zhou and Ming Li. Tri-training: Exploiting unlabeled data using three classifiers.  
78 *IEEE Transactions on knowledge and Data Engineering*, 17(11):1529–1541, 2005.
- 79 [7] Wei-Ning Hsu and James Glass. Disentangling by Partitioning: A Representation Learning  
80 Framework for Multimodal Sensory Data. 2018.
- 81 [8] Changde Du, Changying Du, Jinpeng Li, Wei-long Zheng, Bao-liang Lu, and Huiguang He.  
82 Semi-supervised Bayesian Deep Multi-modal Emotion Recognition. 2017.
- 83 [9] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil  
84 Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Proceedings of  
85 Advances in neural information processing systems*, pages 2672–2680, 2014.
- 86 [10] James T Becker, François Boiler, Oscar L Lopez, Judith Saxton, and Karen L McGonigle. The  
87 natural history of Alzheimer’s disease: description of study cohort and accuracy of diagnosis.  
88 *Archives of Neurology*, 51(6):585–594, 1994.
- 89 [11] Dua Dheeru and Efi Karra Taniskidou. {UCI} Machine Learning Repository, 2017.
- 90 [12] Sérgio Moro, Paulo Cortez, and Paulo Rita. A data-driven approach to predict the success of  
91 bank telemarketing. *Decision Support Systems*, 62:22–31, 2014.
- 92 [13] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training  
93 by reducing internal covariate shift. In *Proceedings of the 32nd International Conference on  
94 Machine Learning (ICML-15)*, pages 448–456, 2015.
- 95 [14] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito,  
96 Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in  
97 PyTorch. 2017.
- 98 [15] F Pedregosa, G Varoquaux, A Gramfort, V Michel, B Thirion, O Grisel, M Blondel, P Pret-  
99 tenhofer, R Weiss, V Dubourg, J Vanderplas, A Passos, D Cournapeau, M Brucher, M Perrot,  
100 and E Duchesnay. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning  
101 Research*, 12:2825–2830, 2011.
- 102 [16] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *Proceesings of  
103 International Conference on Learning Representations (ICLR)*, 2014.
- 104 [17] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine  
105 learning research*, 9(Nov):2579–2605, 2008.

## 106 4 Appendix

### 107 4.1 Detailed description of datasets

108 **The Bank Marketing dataset** is from the UCI machine learning repository[11]. used for predicting  
109 whether the customer will subscribe a term deposit in a bank marketing campaign via telephone[12].  
110 There are originally 4,640 positive samples (subscribe) and 36,548 negative ones (did not subscribe).  
111 Since consensus network models do not work well on imbalanced datasets, we randomly sample  
112 5,000 negative samples to create an (almost) balanced dataset. We also convert the categorical raw  
113 features<sup>1</sup> into one-hot representations. We then divide the features into three somewhat arbitrary  
114 modalities: basic information, statistical data, and employment-related features.

115 Three modalities are determined as following. The division are somewhat arbitrary, except that we  
116 try to make the binary features resulting from one categorical feature be in the same modality.

- 117 1. Basic information: age, marital, education, housing, loan, contact, duration, pdays, previous,  
118 management
- 119 2. Statistical information: campaign, poutcome, emp.var.rate, unknown, cons.conf.idx, eu-  
120 ribor3m, day\_in\_week (converted to 7 binary features), month (converted to 12 binary  
121 features)
- 122 3. Employment-related information: consumer price index, never employed, retired, self-  
123 employed, technician, services, student, housemaid, entrepreneur, blue-collar

124 **DementiaBank**<sup>2</sup> contains 473 narrative picture descriptions of the clinical “cookie-theft  
125 picture”[10], containing 240 positive samples (the Dementia class) and 233 negative samples (the  
126 Control class). We extract 413 linguistic features from each speech sample and their transcriptions,  
127 including acoustic (e.g., pause durations), semantic-syntactic (e.g., complexity of the syntactic parse  
128 structures), and lexical modality(e.g., average word length).

- 129 1. Acoustic-related features: phonation rate, mean pause duration, pause word ratio, total  
130 speech duration, short/medium/long pause count, speech rate, word/audio/(filled or un-  
131 filled) pauses durations, the mean/variance/kurtosis/skewness of the first 42 Mel Frequency  
132 Cepstral Coefficients
- 133 2. Syntactic-semantic features: probabilistic context-free grammar parsing tree heights (average  
134 / max / etc.), and the occurrences of 104 production rules (e.g: NP → DT).
- 135 3. Lexical and POS-derived features: the occurrences of part-of-speech tags, Brunet’s index,  
136 Honore’s statistics, word length, cosine distances between words in sentences, etc.

### 137 4.2 Implementation

138 For simplicity, we use fully connected networks for all of  $I_{1..M}$ ,  $D$ ,  $C$ , and  $R_{1..M}$  in this paper. To  
139 enable faster convergence, all fully connected networks have a batch normalization[13] layer. For  
140 training neural networks, the batch size is set to 10. The neural network models are implemented  
141 using PyTorch[14], and supervised learning benchmark algorithms (SVM, MLP) in scikit-learn[15].

142 We use the Adam optimizer[16] with an initial learning rate of 0.001. In training TCN, TCN-embed,  
143 and TCN-AE, optimizations are stopped when the classification loss does not change by more than  
144  $10^{-5}$  in comparison to the previous step, or when the step count reaches 100. In the pre-training  
145 phase of TCN-embed and TCN-svm, training is stopped when the discrimination loss changes by  
146 less than  $10^{-5}$ , or when pretraining step count reaches 20.

147 Sometimes the iterative optimization (i.e., the I-D-CI cycle for TCN / TCN-embed, and the I-D-RI-CI  
148 cycle for the TCN-AE variant) is trapped in local saddle points – the training classification loss does  
149 not change while the training classification loss is higher than  $\log_2 \approx 0.693^3$ . We check once more  
150 when training stops. If the training classification loss is higher than  $\log_2$ , the model is re-initialized

---

<sup>1</sup><https://archive.ics.uci.edu/ml/datasets/bank+marketing>

<sup>2</sup><https://dementia.talkbank.org>

<sup>3</sup>expected loss of a binary classifier with zero knowledge

151 with a new random seed and the training is restarted. Empirically this re-initialization only happen no  
 152 more than once per ten runs, but the underlying cause need to be examined further.

### 153 4.3 Monitoring the similarity of interpretations

154 **The similarity of interpretations** It is important to evaluate whether the adversarial and classifier  
 155 mechanisms make the interpretations more similar. To evaluate the similarity, we treat the hidden  
 156 dimensions of each interpretation vector  $\mathbf{v}_m = [v_{m,1}, v_{m,2}, \dots, v_{m,j}, \dots]$  (after normalization by their  
 157 sum) as discrete values of a probability mass function<sup>4</sup>, which we write as  $p_m$ . The  $M$  modalities for  
 158 each data point are therefore approximated by  $M$  probability distributions. Now, we can measure the  
 159 relative JS divergences between each pair of interpretation vectors  $\mathbf{v}_m$  and  $\mathbf{v}_n$  derived from the same  
 160 data sample ( $\hat{D}(p_m||p_n)$ ). To acquire the relative value, we normalize the JS divergence by the total  
 161 entropy in  $p_m$  and  $p_n$ :

$$\hat{D}(p_m||p_n) = \frac{1}{2(\mathbb{H}_{p_m} + \mathbb{H}_{p_n})} (D_{KL}(p_m||p_n) + D_{KL}(p_n||p_m))$$

$$\text{where } D_{KL}(p_m||p_n) = \sum_j v_{m,j} \log \frac{v_{n,j}}{v_{m,j}}$$

162 where  $v_{m,j}$  and  $v_{n,j}$  are the  $j^{\text{th}}$  component of  $\mathbf{v}_m$  and  $\mathbf{v}_n$  respectively. In total, for each data sample,  
 163  $\frac{M(M-1)}{2}$  pairs of relative divergences are calculated. We average the negative of these divergences to  
 164 get the similarity for the interpretations:

$$\text{Similarity} = \mathbb{E}_i \mathbb{E}_{m,n \in \{1..M\} \text{ and } m \neq n} \{ -\hat{D}(p_m^{(i)}||p_n^{(i)}) \}$$

165 Note that the ‘‘similarity’’ is defined such that its maximum possible value is 0 (where there is no JS  
 166 divergence between any pair of the interpretation vectors), and it has no theoretical lower bound.

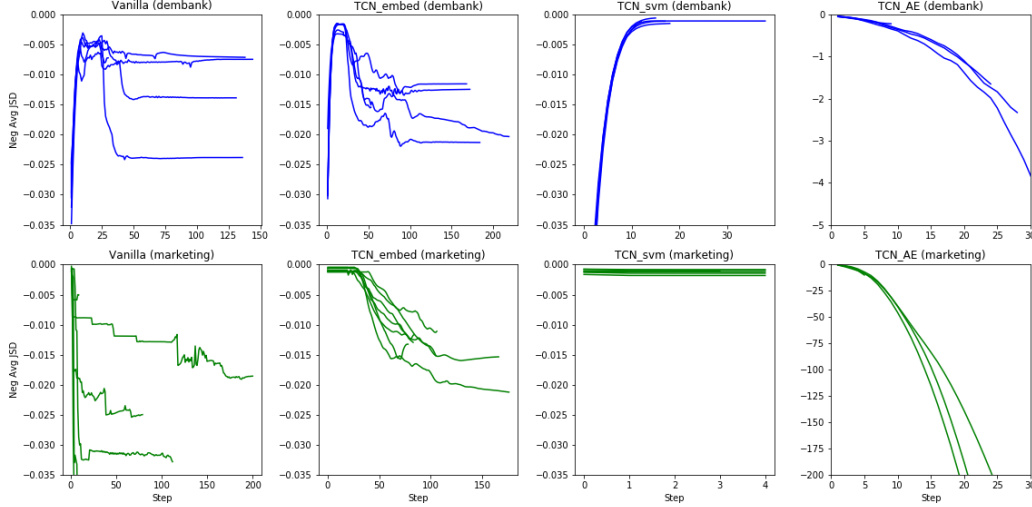


Figure 4: Examples of similarity plots against the number of steps taken, for DementiaBank using 80 labeled samples (‘‘DB80’’, blue) and Bank Marketing using 20 labeled samples (‘‘BM20’’, green). The y axis are scaled to (-0.035, 0) except TCN-AE, where the relative JS divergences ‘‘explode’’. Note that training stops when losses converge (as detailed in §4.2), so the trials may stop at different steps.

167 **Experiments monitoring similarities** We monitor the similarities (defined as negative relative JS  
 168 divergence) between interpretation vectors. Several trends can be observed in Figure 4:

- 169 1. In vanilla TCN on DementiaBank, the similarity usually first rises, then drops to a stable  
 170 final value. On Bank Marketing, the similarity drops without first rising. This might be

<sup>4</sup>There is a ReLU layer at output of each interpreter, so the probability mass will be non-negative.

171 attributed to Bank Marketing modalities (containing only  $\approx 15$  features per modality) being  
 172 not as “sufficient and redundant” (borrowing from [6]) as DementiaBank (containing  $\approx 110$   
 173 per modality).

- 174 2. In the absence of the classifier mechanism, similarities converge to almost the highest  
 175 possible value under the consensus mechanism. This can be seen in the pretraining phase of  
 176 TCN-embed and TCN-svm. Note that there is no explicit steps of ‘converging to the high  
 177 similarity’ on the Bank Marketing dataset – they directly go to the max values – because each  
 178 Bank Marketing step contains more pretraining iterations than DementiaBank ( $\approx 10,000$   
 179 samples vs.  $\approx 500$  samples), resulting in much stronger consensus mechanisms.
- 180 3. In TCN-embed, the classifier mechanism later “pulls down” the similarity. Note that the  
 181 accuracy of TCN-svm is around 50%; we can infer that a meaningful consensus state need  
 182 not have perfect similarity.
- 183 4. The addition of reconstructors inhibit the consensus mechanism in terms of reaching a high  
 184 similarity between interpretations, as shown by the exploding JS divergences in TCN-AE  
 185 models. This further illustrates that TCN distills information in a different aspect from  
 186 denoising autoencoders.

#### 187 4.4 Visualizing the interpretations

188 Figure 5 shows several 2D visualizations of interpretation vectors drawn from an arbitrary run, as  
 189 an example of interpretations with low, medium, and high similarity. In §??, we illustrate how the  
 190 similarities between interpretations evolve during optimization in TCN models.

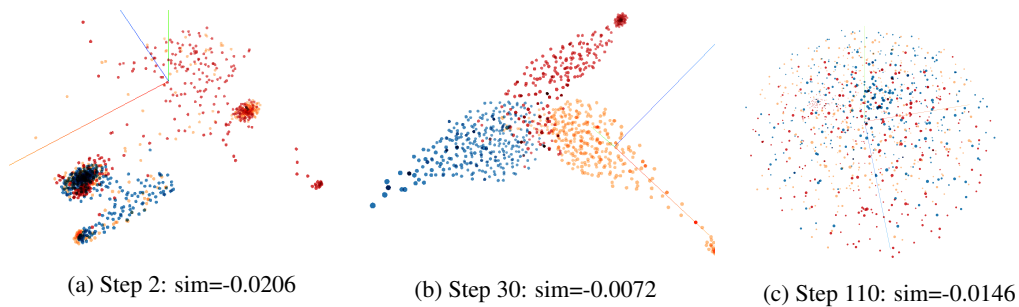


Figure 5: Three 2-D T-SNE[17] visualizations comparing interpretation vectors among modalities, taken from a run of the vanilla TCN (on DementiaBank dataset with 80 labeled data). The three colors represent three modalities. At step 2, the interpretations are distributed randomly. At step 110, they become mixed evenly. The most interesting embedding happens at step 30, when interpretations of the three modalities form three ‘drumstick’ shapes. With the highest symmetry visually, this configuration of interpretations also has the highest *similarity* among the three.