

News Article Classification

(Multi-Label Learning: ML-KNN & BP-MLL)

Lauren Contard, Archit Datar, Yue Li, Bobby Lumpkin, Haihang Wu

The Ohio State University

STAT 6500



1 Introduction and Problem Statement

- Introduction and Description of Datasets
- Problem Statement

2 KNN Based Approaches

- Binary Relevance
- ML-KNN Algorithm
- Results

3 Neural Network Based Approaches

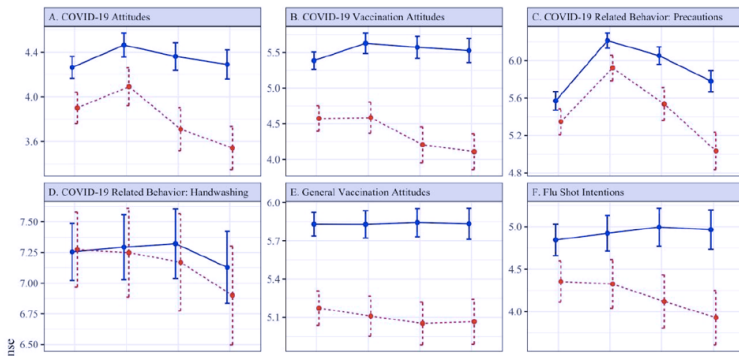
- Architectures: Feed Forward & Recurrent Networks
- Loss Functions: Cross Entropy vs BPMLL
- Results

4 Discussion and Conclusions

Introduction and Problem Statement

Introduction

- People in the U.S. have shown polarized attitudes and behaviors in response to COVID-19.
- What could be some of the drivers? We aim to explore whether and how mainstream news media in the U.S. play a role in the polarization of the public's attitudes and behaviors.



Description of Datasets

- We analyzed news content about elite communication from 10 news media, including newspapers (e.g., New York Times) and cable network news (e.g., ABC news).
- We got the list of news articles from the GDELT database and used the Python package "Scrapy" to get the full text.
- We have 8588 news articles in total and randomly selected 290 paragraphs of them to label.
- We manually labeled our paragraphs into 15 categories. One paragraph usually has multiple labels.

Two Examples from the Dataset

"Meanwhile, Bottoms has repeatedly urged residents to stay home. On Friday, she tweeted coronavirus fatality statistics for the state: 'The numbers speak for themselves. PLEASE STAY HOME.'" (From USAtoday, 04/25/2020)

- Threats/impacts
- Responses/actions
- Severity
- Self-efficacy
- Public Health
- Negative

"You know, I wonder, I think there's been a lot of self-congratulations every day that we see in those briefings, frankly, about the testing in the United States, and we're doing so well, we're doing now more than South Korea did." (From abcnews, 03/29/2020)

- Responses/actions
- External-efficacy
- Public health
- Political evaluation
- Negative

Problem Statement

- We have a multi-label classification problem. We adapted two learning algorithms: k-NN and Neural Networks to handle multi-label data directly.

RQ1: Which learning algorithm performs better in the multi-label classification problem: k-NN or neural networks??

- Our data has high dimensions.

RQ2: How can we effectively implement linear and non-linear dimensions reduction in the multi-label classification problem?

- The sequence of features need to be considered in the classification.

RQ3: Which type of the features perform better in the multi-label classification problem: with or without considering the order of words?

KNN Based Approaches

Binary Relevance (A Naive Approach)

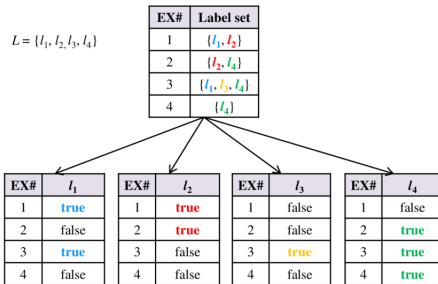
- An intuitive approach to deal with the multilabel paradigm.

Binary Relevance (A Naive Approach)

- An intuitive approach to deal with the multilabel paradigm.
- Works by decomposing the multi-label learning task into a number of independent binary learning tasks (one per class label).

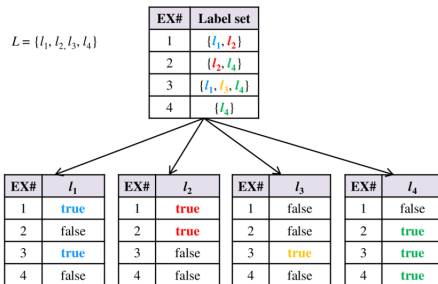
Binary Relevance (A Naive Approach)

- An intuitive approach to deal with the multilabel paradigm.
- Works by decomposing the multi-label learning task into a number of independent binary learning tasks (one per class label).



Binary Relevance (A Naive Approach)

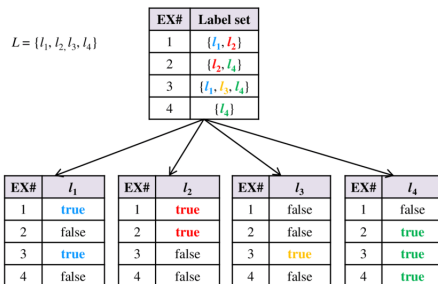
- An intuitive approach to deal with the multilabel paradigm.
- Works by decomposing the multi-label learning task into a number of independent binary learning tasks (one per class label).



→ Often criticized in the literature because of its label independence assumption.

Binary Relevance (A Naive Approach)

- An intuitive approach to deal with the multilabel paradigm.
- Works by decomposing the multi-label learning task into a number of independent binary learning tasks (one per class label).



- Often criticized in the literature because of its label independence assumption.
- We implement a KNN based binary relevance model and compare with a more novel adaptation: the ML-KNN model.

ML-KNN Algorithm: Overall Approach

Overall Approach: This ML-KNN algorithm takes a parametric, Bayesian approach towards estimating the Bayes Optimal Classifier. As with the single-label algorithm, it does this using the K nearest neighbors of an instance. Namely...

ML-KNN Algorithm: Overall Approach

Overall Approach: This ML-KNN algorithm takes a parametric, Bayesian approach towards estimating the Bayes Optimal Classifier. As with the single-label algorithm, it does this using the K nearest neighbors of an instance. Namely...

- Given a test instance, t , \vec{Y}_t is determined using the MAP estimate:

ML-KNN Algorithm: Overall Approach

Overall Approach: This ML-KNN algorithm takes a parametric, Bayesian approach towards estimating the Bayes Optimal Classifier. As with the single-label algorithm, it does this using the K nearest neighbors of an instance. Namely...

- Given a test instance, t , \vec{Y}_t is determined using the MAP estimate:

$$\begin{aligned}\vec{y}_t(\ell) &= \operatorname{argmax}_{b \in \{0,1\}} \mathbb{P} \left(H_b^\ell | E_{\vec{C}_t(\ell)}^\ell \right), \quad \ell \in \mathcal{Y} \\ &= \operatorname{argmax}_{b \in \{0,1\}} \frac{\mathbb{P} \left(H_b^\ell \right) \cdot \mathbb{P} \left(E_{\vec{C}_t(\ell)}^\ell | H_b^\ell \right)}{\mathbb{P} \left(E_{\vec{C}_t(\ell)}^\ell \right)} \\ &= \operatorname{argmax}_{b \in \{0,1\}} \mathbb{P} \left(H_b^\ell \right) \cdot \mathbb{P} \left(E_{\vec{C}_t(\ell)}^\ell | H_b^\ell \right)\end{aligned}$$

ML-KNN Algorithm: Overall Approach

Overall Approach: This ML-KNN algorithm takes a parametric, Bayesian approach towards estimating the Bayes Optimal Classifier. As with the single-label algorithm, it does this using the K nearest neighbors of an instance. Namely...

- Given a test instance, t , \vec{Y}_t is determined using the MAP estimate:

$$\begin{aligned}\vec{y}_t(\ell) &= \operatorname{argmax}_{b \in \{0,1\}} \mathbb{P} \left(H_b^\ell | E_{\vec{C}_t(\ell)}^\ell \right), \quad \ell \in \mathcal{Y} \\ &= \operatorname{argmax}_{b \in \{0,1\}} \frac{\mathbb{P} \left(H_b^\ell \right) \cdot \mathbb{P} \left(E_{\vec{C}_t(\ell)}^\ell | H_b^\ell \right)}{\mathbb{P} \left(E_{\vec{C}_t(\ell)}^\ell \right)} \\ &= \operatorname{argmax}_{b \in \{0,1\}} \mathbb{P} \left(H_b^\ell \right) \cdot \mathbb{P} \left(E_{\vec{C}_t(\ell)}^\ell | H_b^\ell \right)\end{aligned}$$

- Where we take a Bayesian approach towards estimating the prior probabilities, $\mathbb{P} \left(H_b^\ell \right)$, and conditional probabilities, $\mathbb{P} \left(E_{\vec{C}_t(\ell)}^\ell | H_b^\ell \right)$.

ML-KNN Algorithm: More Notation

Notation:

ML-KNN Algorithm: More Notation

Notation:

- Let $N(x)$ denote the set of K nearest neighbors of x , identified in the training set.

ML-KNN Algorithm: More Notation

Notation:

- Let $N(x)$ denote the set of K nearest neighbors of x , identified in the training set.
- Let $\vec{C}_x(\ell) = \sum_{a \in N(x)} \vec{y}_a(\ell)$ ($\ell \in \mathcal{Y}$) define a membership counting vector.

ML-KNN Algorithm: More Notation

Notation:

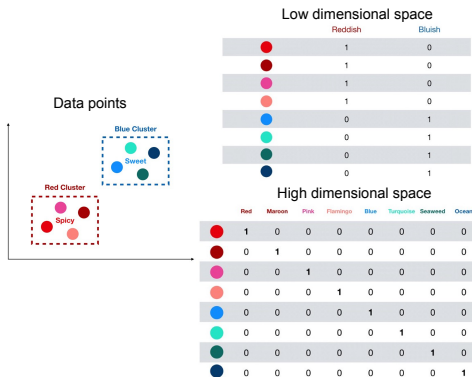
- Let $N(x)$ denote the set of K nearest neighbors of x , identified in the training set.
- Let $\vec{C}_x(\ell) = \sum_{a \in N(x)} \vec{y}_a(\ell)$ ($\ell \in \mathcal{Y}$) define a membership counting vector.
- Let H_0^ℓ denote the event that test instance t does not have a label ℓ and let H_1^ℓ denote the event that it does have label ℓ .

ML-KNN Algorithm: More Notation

Notation:

- Let $N(x)$ denote the set of K nearest neighbors of x , identified in the training set.
- Let $\vec{C}_x(\ell) = \sum_{a \in N(x)} \vec{y}_a(\ell)$ ($\ell \in \mathcal{Y}$) define a membership counting vector.
- Let H_0^ℓ denote the event that test instance t does not have a label ℓ and let H_1^ℓ denote the event that it does have label ℓ .
- Let E_j^ℓ ($j \in \{1, \dots, K\}$) denote the event that, among the K nearest neighbors of t , there are exactly j instances which have label ℓ .

ML-KNN Algorithm: Reasons for dimension reduction



Having a high dimensional feature space causes Euclidian distances between points to be fairly similar as the distance vector components are partitioned across many dimensions.

Nonlinear Dimension Reduction (ANN Autoencoder)

Nonlinear Dimension Reduction (ANN Autoencoder)

- PCA allows only for reductions in dimension due to linear representations.

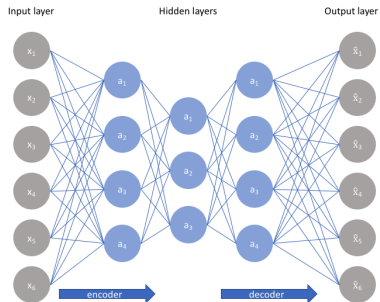
Nonlinear Dimension Reduction (ANN Autoencoder)

- PCA allows only for reductions in dimension due to linear representations.
- Autoencoders can learn data projections with suitable dimensionality and sparsity limitations that allow for nonlinear representations, as well.

Nonlinear Dimension Reduction (ANN Autoencoder)

- PCA allows only for reductions in dimension due to linear representations.
- Autoencoders can learn data projections with suitable dimensionality and sparsity limitations that allow for nonlinear representations, as well.

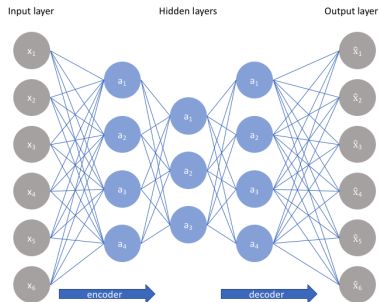
→ Train a feed forward neural network with a "bottleneck layer" to perform the identity mapping.



Nonlinear Dimension Reduction (ANN Autoencoder)

- PCA allows only for reductions in dimension due to linear representations.
- Autoencoders can learn data projections with suitable dimensionality and sparsity limitations that allow for nonlinear representations, as well.

- Train a feed forward neural network with a "bottleneck layer" to perform the identity mapping.
- Use bottleneck layer values as encodings of the data.



Threshold Function Learning

Perfect classification, using a constant threshold requires two conditions:

- 1 Predicted logit values for "true" labels be separated from predicted logit values for "false" labels.
- 2 This separation be around some constant value (usually either 0.5 or 0)

Learning a threshold function aims to relax the second condition.

→ Namely, we fit a linear regression model to learn threshold values from the logit outputs of our models.

Neural Network Based Approaches

A Brief History

A Brief History

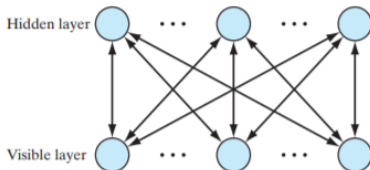
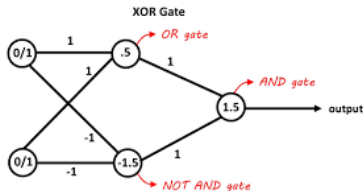
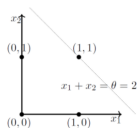
- Inspired by biological nervous systems, neural networks date back to the first half of the 20th century with works such as those by McCulloch and Pitts, which could model simple logical operations.

AND Function



AND function

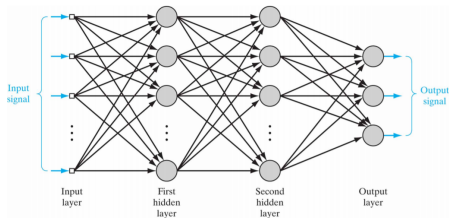
$$x_1 + x_2 = \sum_{i=1}^2 x_i \geq 2$$



Network Architectures (Feed Forward vs Recurrent)

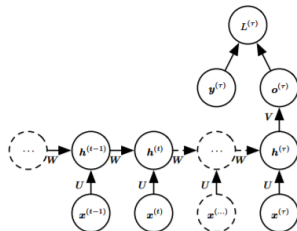
Feed Forward Networks

- Neurons in the first layer represent components of the input vectors.
- The output of the neuron in the next layer is determined by applying a non-linear “activation function” to a linear combination of the input components, plus a bias.



Recurrent Neural Networks (RNNs)

- RNNs are a popular adaptation for NLP problems.
- They utilize hidden unit connections with shared weights.
- Unfolding an RNN let's us visualize it like a feed forward network (see below).



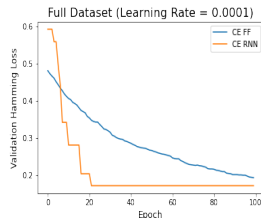
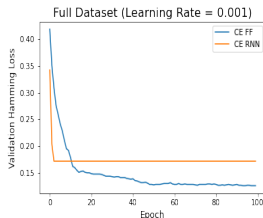
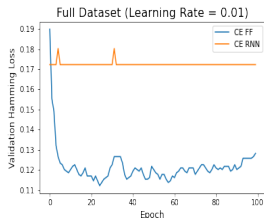
Naive vs Novel Approaches (Cross Entropy vs BPMML)

- By "naive" we refer to multilabel networks that utilize a cross entropy loss for training.
- By comparison, Zhang and Zhou bpmll proposed a novel loss, that emphasizes pairwise ranking accuracy.

$$E = \sum_{i=1}^m E_i = \sum_{i=1}^m \frac{1}{|Y_i| |\overline{Y}_i|} \sum_{(k,l) \in Y_i \times \overline{Y}_i} \exp(-(c_k^i - c_l^i))$$

so that the i^{th} error term is severely penalized if c_k^i is much smaller than c_l^i .

Artificial Neural Network Results: Full Dataset

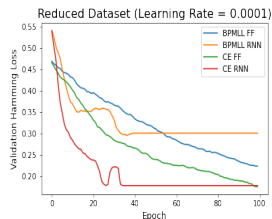
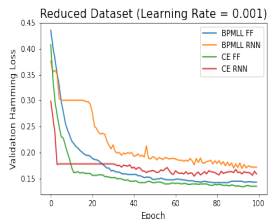
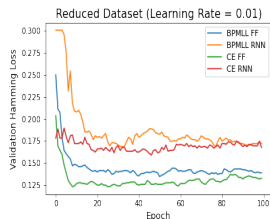


Learning Rate	CE FF	CE RNN
0.01	0.3020833333	0.2379807692
0.001	0.17868589740	0.1730769231
0.0001	0.2115384615	0.1770833333

Table 1: Hamming Loss with Threshold Function Learning

CE FF outperform CE RNN in constant threshold but underperform in learned threshold; The effect of learning rate.

Artificial Neural Network Results: Reduced Dataset



Learning Rate	CE FF	BPMLL FF	CE RNN	BPMLL RNN
0.01	0.1520979021	0.145979021	0.2447552448	0.2194055944
0.001	0.1853146853	0.2578671329	0.1791958042	0.20454545
0.0001	0.2071678322	0.1844405594	0.1896853147	0.2132867133

Table 2: Hamming Loss with Threshold Function Learning

Same conclusion for RNN and FF; BPMLL shows NO better performance in hamming loss than cross entropy.

Discussion and Conclusions

Discussion and Conclusions: ANN Approaches

Discussion and Conclusions: ANN Approaches

- 1 The FF networks tend to outperform the RNNs.

Discussion and Conclusions: ANN Approaches

- ① The FF networks tend to outperform the RNNs.
 - More parameters lead to overfitting in the bidirectional LSTM.

Discussion and Conclusions: ANN Approaches

- ① The FF networks tend to outperform the RNNs.
 - More parameters lead to overfitting in the bidirectional LSTM.
 - Single directional LSTM can't escape locally optimal solution.

Discussion and Conclusions: ANN Approaches

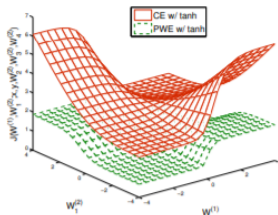
- ① The FF networks tend to outperform the RNNs.
 - More parameters lead to overfitting in the bidirectional LSTM.
 - Single directional LSTM can't escape locally optimal solution.
 - Training on more data may improve.

Discussion and Conclusions: ANN Approaches

- ① The FF networks tend to outperform the RNNs.
 - More parameters lead to overfitting in the bidirectional LSTM.
 - Single directional LSTM can't escape locally optimal solution.
 - Training on more data may improve.
- ② The cross entropy networks tend to outperform the BPMLL networks.

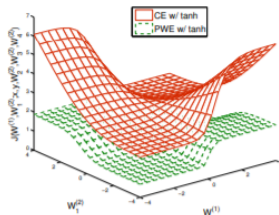
Discussion and Conclusions: ANN Approaches

- 1 The FF networks tend to outperform the RNNs.
 - More parameters lead to overfitting in the bidirectional LSTM.
 - Single directional LSTM can't escape locally optimal solution.
 - Training on more data may improve.
- 2 The cross entropy networks tend to outperform the BPMLL networks.
 - The surface for the BPMLL loss with tanh activations has plateaus in which gradient descent can be very slow in comparison with the cross-entropy (Nam et al. [2014]). The case may be similar for ReLU activations.



Discussion and Conclusions: ANN Approaches

- ① The FF networks tend to outperform the RNNs.
 - More parameters lead to overfitting in the bidirectional LSTM.
 - Single directional LSTM can't escape locally optimal solution.
 - Training on more data may improve.
- ② The cross entropy networks tend to outperform the BPMLL networks.
 - The surface for the BPMLL loss with tanh activations has plateaus in which gradient descent can be very slow in comparison with the cross-entropy (Nam et al. [2014]). The case may be similar for ReLU activations.



- While BPMLL is supposed to leverage correlations between labels, Nam et al. [2014] conjecture that these correlations also may cause overfitting.

Discussion and Conclusions: ANN Approaches (continued)

- ③ Models saw no improvement from using a learned threshold function.

Discussion and Conclusions: ANN Approaches (continued)

- ③ Models saw no improvement from using a learned threshold function.
 - Models could sufficiently separate predicted logits about 0.5.

Discussion and Conclusions: ANN Approaches (continued)

- ③ Models saw no improvement from using a learned threshold function.
 - Models could sufficiently separate predicted logits about 0.5.
 - This approach may be more useful for models that require extensive resources for sufficient training (improved less extensively trained models).

- Jinseok Nam, Jungi Kim, Eneldo Loza Mencía, Iryna Gurevych, and Johannes Fürnkranz. Large-scale multi-label text classification — revisiting neural networks. In Toon Calders, Floriana Esposito, Eyke Hüllermeier, and Rosa Meo, editors, *Machine Learning and Knowledge Discovery in Databases*, pages 437–452, Berlin, Heidelberg, 2014. Springer Berlin Heidelberg.
- Min-Ling Zhang and Zhi-Hua Zhou. Ml-knn: A lazy learning approach to multi-label learning. *Pattern Recognition*, 40(7):2038–2048, 2007. doi: 10.1016/j.patcog.2006.12.019.
- Min-Ling Zhang and Zhi-Hua Zhou. Multilabel neural networks with applications to functional genomics and text categorization. *IEEE Transactions on Knowledge and Data Engineering*, 18(10):1338–1351, 2006. doi: doi:10.1109/TKDE.2006.162.