

CITIC Research Center, Universidade da Coruña

### Can data placement be effective for Neural Networks classification tasks? Introducing the Orthogonal Loss

B. Cancela, V. Bolón-Canedo and A. Alonso-Betanzos

< □ > < □ > < □ > < □ > < □ > < □ >

#### Introduction

Classic classification losses aim to separate the deep features between classes as much as possible:

• It does not matter where the deep features are placed, as long as they are clearly separated.

Furthermore, confidence scores do not provide any information about how confident is in the decision.

• For instance, an image is classified as a car with a 90% confidence, but it is a plane, which is not part of the training data.

< □ > < □ > < □ > < □ > < □ > < □ >

## Introduction

Imagine we have two features, two classes and a linear classifier which is the identity matrix.

• Where the training data is placed for a loss value close to zero?



## Research goal

Having an orthogonal classifier, is it possible to place the data over the classifier's autovectors, one per each class?

• Positive values for the direction of the correct class and zeros for the rest.



## Research goal

Confidence scores do not provide any information about how confident is in the decision.

- For instance, an image is classified as a car with a 90% confidence, but it is a plane, which is not part of the training data.
- SOLUTION: A variant to the softmax function, aiming to provide a more accurate label confidence.

< □ > < □ > < □ > < □ > < □ > < □ >

# Requisites

**Requisite 1:** It should be a Neural Network with one hidden layer at list.

**Requisite 2:** The last layer will be fixed during the training procedure.

**Requisite 3:** The last layer will be initialized as an orthogonal matrix.

イロト イボト イヨト イヨト

#### Exponential Orthogonal Loss

Forces the correct class to have the highest positive value.

- $\tilde{y}^{(i)}$ : model prediction
- y<sup>(i)</sup>: input label (one-hot encoding)

$$\mathcal{L}^{EOL}(\tilde{\mathbf{y}}^{(i)}, \mathbf{y}^{(i)}) = \sum \underbrace{\mathbf{y}^{(i)}}_{\mathbf{y}^{(i)}} e^{R - \tilde{\mathbf{y}}^{(i)}} + \lambda \underbrace{(1 - \mathbf{y}^{(i)})}_{\mathcal{L}^{EOL}_{neel}} (\tilde{\mathbf{y}}^{(i)})^2, \quad R > 0,$$

イロト 不良 トイヨト イヨト

#### Trench Orthogonal Loss

Forces the correct class to have the fixed positive value R.

- $\tilde{y}^{(i)}$ : model prediction
- y<sup>(i)</sup>: input label (one-hot encoding)

$$\mathcal{L}^{TOL}(\tilde{\mathbf{y}}^{(i)}, \mathbf{y}^{(i)}) = \sum \underbrace{\mathbf{y}^{(i)}}_{\mathbf{y}^{(i)}} \underbrace{(R - \tilde{\mathbf{y}}^{(i)})^{2}}_{\mathcal{L}_{regL}} + \lambda \underbrace{(1 - \mathbf{y}^{(i)})}_{\mathcal{L}_{regL}^{TOL}} (\tilde{\mathbf{y}}^{(i)})^{2}, \quad R > 0,$$

### Softmax variation

- $\tilde{y}^{(i)}$ : model prediction
- Checks how good the projection is in each class.
- $\sum p \leq 1$

$$p(\tilde{\mathrm{y}}^{(i)}) = rac{\max(0, \tilde{\mathrm{y}}^{(i)})}{\max(P, \sum \max(0, \tilde{\mathrm{y}}^{(i)}))}, \quad P > 0.$$

イロト イポト イヨト イヨト

э.

#### **Experimental Results**

- Tested in 5 different datasets using two networks.
- The results are better (or comparable) than using the cross-entropy loss.

VGG	CE (fixed classifier)	CE	EOL	TOL
MNIST	$99.57\pm0.01$	$99.57\pm0.04$	$99.57\pm0.06$	$\textbf{99.64} \pm \textbf{0.01}$
Fashion-MNIST	$94.26\pm0.13$	$94.27\pm0.07$	$\textbf{94.35} \pm \textbf{0.05}$	$\textbf{94.45} \pm \textbf{0.14}$
CIFAR-10	$91.70\pm0.13$	$91.55\pm0.17$	$\textbf{92.03} \pm \textbf{0.13}$	$\textbf{91.82} \pm \textbf{0.12}$
CIFAR-100	$66.38\pm0.67$	$63.13\pm0.80$	$\textbf{69.95} \pm \textbf{0.19}$	$\textbf{69.90} \pm \textbf{0.34}$
STL-10	$79.78\pm0.56$	$78.27\pm0.23$	$\textbf{81.39} \pm \textbf{1.08}$	$\textbf{80.44} \pm \textbf{0.79}$

< ロ > < 同 > < 回 > < 回 > < 回 > <

э

#### **Experimental Results**

- Tested in 5 different datasets using two networks.
- The results are better (or comparable) than using the cross-entropy loss.

WRN-16-4	CE (fixed classifier)	CE	EOL	TOL
MNIST	$99.67\pm0.02$	$99.64\pm0.02$	$99.64\pm0.02$	$99.66\pm0.01$
Fashion-MNIST	$95.31\pm0.10$	$95.26\pm0.14$	$\textbf{95.47} \pm \textbf{0.13}$	$95.21\pm0.18$
CIFAR-10	$94.59\pm0.13$	$94.59\pm0.17$	$\textbf{94.88} \pm \textbf{0.13}$	$94.75\pm0.10$
CIFAR-100	$74.95\pm0.09$	$74.70\pm0.24$	$74.92\pm0.09$	$\textbf{75.39} \pm \textbf{0.18}$
STL-10	$86.19\pm0.35$	$85.70\pm0.21$	$\textbf{86.65} \pm \textbf{0.34}$	$\textbf{86.71} \pm \textbf{0.28}$

< 日 > < 同 > < 回 > < 回 > .

э

## Using a reduced training data version

Both EOL and TOL obtain a better accuracy when using a subset of the training data.



< ロ > < 同 > < 回 > < 回 > < 回 >

#### Effect of the softmax variant

The number of errors is reduced as our variant output increases its confidence.



(日) (四) (日) (日) (日)

# Conclusions

- Deep features can be effectively placed without losing performance (sometimes even better).
- It obtains better results when the number of training samples is low.
- The softmax variant provides a better degree of certainty in the probability outputs.
- Deep features placed in orthogonal projections have interesting properties.

< ロ > < 同 > < 回 > < 回 > < 回 > <



CITIC Research Center, Universidade da Coruña

## Can data placement be effective for Neural Networks classification tasks? Introducing the Orthogonal Loss

B. Cancela, V. Bolón-Canedo and A. Alonso-Betanzos

< ロ > < 同 > < 回 > < 回 > < 回 > <