

SMate: Synthetic Minority Adversarial Technique

Pablo Rodriguez Bertorello, Liang Ping Koh

1. Problem Statement

- In important classification problems, like cancer detection, datasets can be imbalanced
- Imbalanced can lead to poor classification on minority classes
- Current methods to address class imbalance are more suited for non-image problems
- Contribution: algorithm based on GANs and Transfer Learning solve the class imbalance

2. Dataset and Pre-processing

- CIFAR10 consists of 60,000 images, in 10 categories
- Train-test split of 50,000 - 10,000
- To induce imbalance, Truck images are intentionally undersampled: only 500 images of trucks for training



3. Description of Alternative Methods

- The following prior art baseline methods are evaluated:
 - Oversampling: sampling with replacement from minority class to rebalance the dataset
 - Synthetic Minority Oversampling Technique (SMOTE): creates new data points as random linear combinations of minority training examples
 - Adaptive synthesis (ADASYN): identifies minority class examples that are difficult to separate. SMOTE is applied to these examples

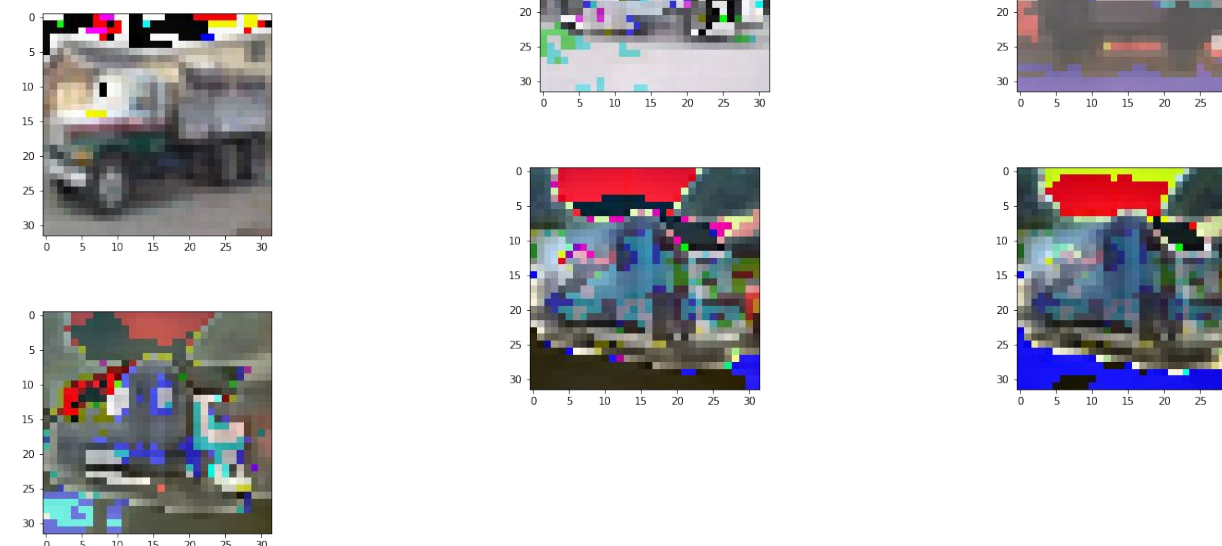
4. Baseline Method Evaluation

- Train a classifier for 100 epochs on the following training sets, and evaluate the test accuracy and confusion matrix:
 - With all 50,000 training samples
 - Imbalanced with 45,500 training examples: 4,500 removed for trucks
 - Oversampling: imbalanced with 45,500 training examples, plus 4,500 generated by sampled from the minority class
 - SMOTE: imbalanced with 45,500 training examples, plus 4,500 synthesized examples
 - ADASYN: imbalanced with 45,500 training examples, plus 4,500 synthesized by SMOTE

5. Prior Art

Model	Test Accuracy	Truck Class Recall
Full Data	77.5%	85%
Undersampled Data	76.1%	56%
Oversampled Data	66.6%	62%
SMOTE Data	75.7%	51%
ADASYN Data	76.6%	56%

- Performance worsens for minority class
- Oversampling improves recall, but worsens accuracy. Likely overfitting to minority
- SMOTE and ADASYN do not improve accuracy, but do not improve minority class image examples



- Linear combinations of image examples produce poor output with a lot of noise. It is an unnatural operation for image data and does not improve the error rate.

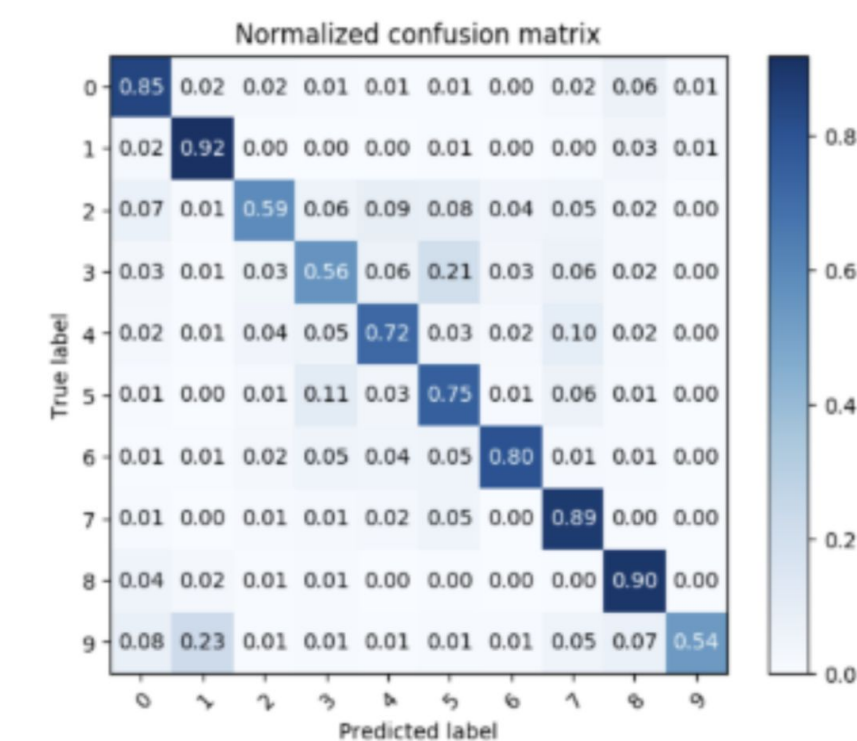
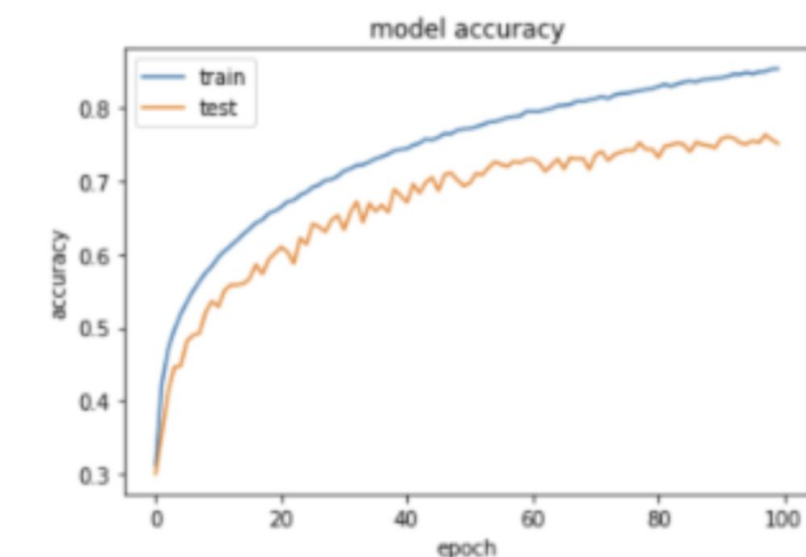
6. SMate Algorithm

- Augment the minority class data-set randomly: Flip, Crop, Gaus-sianBlur, ContrastNormalization, AdditiveGaussianNoise, Multiply, Affine
- Train a generator for all the majority classes
- Perform Transfer Learning, freezing the first six layers of the Generator
- Thereafter, train a minority class Generator. These are the images that our GAN produces:



- Relying on the minority class Generator, balance the minority class with 4500 generated images

7. SMate Performance



- SMate performs better than all the other methods
- Our images look somewhat like trucks, but tend to get confused with cars
- Intuition: starting from a Generator for majority classes, GAN over-fitting is prevented

Conclusion

- SMate outperforms under-sampling, over-sampling, SMOTE, and ADASYN
- Current synthetic methods cannot simulate true distribution for image data
- Future research in loss functions should seek to fit GAN generated images to a minority class, while penalizing for semblance of any of the majority classes
- SMate appears generalizable to every type of data, with high potential for time-series