# Self-training with Noisy Student improves ImageNet classification - Paper, Code (TF)

By Google Brain & CMU: Nov, 2019

| | |
|---|---|
| Contribution | • Noisy Student Training - a semi-supervised learning approach that works well even when labeled data is abundant. <br> • Authors show that noise helps in Student model to generalize better than Teacher model. <br> • Noisy Student Training boosts robustness in computer vision models. |
| Idea | • Paper extends the idea of self-training and distillation with equal-or-large student models and noise added to student during training. <br> • Idea is to learn more capacity Student model than Teacher as opposed to distillation process (smaller and without noise student model is used in distillation). <br> • A semi-supervised learning approach leveraging labeled as well as unlabeled data. |
| Training Process | 1. Train an EfficientNet on labeled data and use this as teacher to generate pseudo labels (soft or one-hot) on large unlabeled data. <br> 2. Train a larger EfficientNet as Student model (from scratch) on combination of labeled and pseudo labeled images with combined cross entropy loss. <br> 3. Make Student as Teacher and go back to first step. <br><br> This process is repeated few times. Authors achieves best results with 3 iterations. Resultant model is called **EfficientNet-L2** (scaled EfficientNet-B7). <br><br> • Two types of noise: 1) input noise - data augmentation with RandAugment. 2) model noise - dropout and stochastic depth are used. <br> • Specifically, the teacher produces high-quality pseudo labels by reading in clean images, while the student is required to reproduce those labels with augmented images as input. <br><br>  <br> Figure 1: Illustration of the Noisy Student Training. (All shown images are from ImageNet.) |
| Results | • Achieves 88.4% top-1 accuracy on ImageNet, which is 2.0% better than the SOTA model that requires 3.5B weakly labeled Instagram images. Proposed approach uses only 300M unlabeled images. <br> • Improves classification robustness on much harder test sets by large margins: ImageNet-A top-1 accuracy from 61.0% to 83.7%, ImageNet-C mean corruption error (mCE) from 45.7 to 28.3 and ImageNet-P mean flip rate (mFR) from 27.8 to 12.2. <br> • <br><br> <br> Table 1: Summary of key results compared to previous state-of-the-art models [86, 55]. Lower is better for mean corruption error (mCE) and mean flip rate (mFR). |
| Insights | 1. Among noisy student training and larger model, former contributes to majority of performance gains. <br> 2. Vision model can generally benefit from Noisy student training even without iterative training. (Ref: Table-16) <br><br> Better Performance with: <br> 1. Using soft probabilities than one-hot labels as pseudo labels works better. Soft pseudo labels are particularly helpful for out-of-domain unlabeled data. <br> 2. A large amount of unlabeled data. <br> 3. Joint training on labeled data and unlabeled data rather than two-stage process of pretraining with unlabeled data and then finetunes on labeled data (Ref: Table-13) <br> 4. Training Student from scratch performs better than initializing with teacher (Ref: Table-15) <br> 5. large Teacher model with better performance <br> 6. Class balancing is crucial from better performance during Student training. Techniques: Undersampling & Duplication |
| Few other details | • ImageNet data is used as labeled data and JFT 300M images data is used as unlabeled set. <br> • Data Filtering: only images with high confidence (>0.3 prob) are used during Student model training. <br> • Training time of EfficientNet-L2 is five times of Efficient-B7. |
| Utility | • Pre-trained checkpoints can be used a feature extractor for various classification and other vision tasks (code) |

Results table:

| | ImageNet top-1 acc. | ImageNet-A top-1 acc. | ImageNet-C mCE | ImageNet-P mFR |
|---|---|---|---|---|
| Prev. SOTA | 86.4% | 61.0% | 45.7 | 27.8 |
| Ours | **88.4%** | **83.7%** | **28.3** | **12.2** |

| Method | Top-1 Acc. | Top-5 Acc. |
|---|---|---|
| ResNet-50 | 77.6% | 93.8% |
| **Noisy Student Training (ResNet-50)** | **78.9%** | **94.3%** |

Table 16: Experiments on ResNet-50.

| Warm-start Epoch | Initializing student with teacher | | | | No Init |
|---|---|---|---|---|---|
| | 35 | 70 | 140 | 280 | 350 |
| Top-1 Acc. | 77.4% | 77.5% | 77.7% | 77.8% | **77.9%** |

Table 15: A student initialized with the teacher still requires at least 140 epochs to perform well. The baseline model, trained with labeled data only, has an accuracy of 77.3%.

| Model | B0 | B1 | B2 | B3 |
|---|---|---|---|---|
| Supervised Learning | 77.3% | 79.2% | 80.0% | 81.7% |
| Pretraining | 72.6% | 75.1% | 75.9% | 76.5% |
| Pretraining + Finetuning | 77.5% | 79.4% | 80.3% | 81.7% |
| Joint Training | **77.9%** | **79.9%** | **80.7%** | **82.1%** |

Table 13: Joint training works better than pretraining and finetuning. We vary the finetuning steps and report the best results. Models are trained for 350 epochs instead of 700 epochs without iterative training.

| Teacher | Teacher Acc. | Student | Student Acc. |
|---|---|---|---|
| B0 | 77.3% | B0 | 77.9% |
| | | B1 | **79.5%** |
| B2 | 80.0% | B2 | 80.7% |
| | | B3 | **82.0%** |
| B4 | 83.2% | B4 | 84.0% |
| | | B5 | **84.7%** |
| B7 | 86.9% | B7 | 86.9% |
| | | L2 | **87.2%** |