

# **Semi-supervised Learning**

#### Temporal Ensembling for Semi-Supervised Learning - Samuli Laine, Timo Aila - ICLR 2017

Meta Pseudo Labels - Hieu Pham, Zihang Dai, Qizhe Xie, Minh-Thang Luong, Quoc V. Le - CVPR 2021





# **Semi-supervised learning**

Learning from both labeled **and** unlabeled data

#### Scenario 1:

You have ImageNet, but want to use even more data

Scenario 2: You have few labels, but loads of unlabeled data



Source: xview2 dataset



# How to make use of the unlabeled data?

- 1. Pre-train model on unlabeled data, then fine-tune on labeled dataset  $\rightarrow$  self-supervised learning
- 2. Train on labeled and unlabeled data at the same time.
  - a. Assign labels to unlabeled data
     → Pseudo labels
  - b. Train unsupervised objective in parallel to supervised → Consistency regularisation / Label propagation



### **Temporal Ensembling for Semi-Supervised Learning**

Samuli Laine, Timo Aila - ICLR 2017



### **Consistency Regularisation**

#### П-model



- Loss = Supervised + Unsupervised
- Unsupervised: **Consistency regularisation**. Different stochastic influences should lead to same results
- Unsupervised loss also used for labeled images



### Average predictions over time

#### П-model



#### Temporal ensembling





### **Self-ensembling over time**



#### **Temporal ensembling**



## **Self-ensembling over time**



#### Advantages:

- Compare prediction to moving average of past predictions
  - → Less noisy than current prediction!
- Only one pass per epoch  $\rightarrow$  Faster!

#### Disadvantage:

- Memory needed to store prediction for whole dataset
- Prediction targets are too old when training on large datasets



	Error rate (%) with # labels		
	4000	All (50000)	
Supervised-only	$35.56 \pm 1.59$	$7.33\pm0.04$	
with augmentation	$34.85 \pm 1.65$	$6.05\pm0.15$	
Conv-Large, $\Gamma$ -model (Rasmus et al., 2015)	$20.40\pm0.47$		
CatGAN (Springenberg, 2016)	$19.58\pm0.58$		
GAN of Salimans et al. (2016)	$18.63 \pm 2.32$		
Π-model	$16.55\pm0.29$	$6.90\pm0.07$	
$\Pi$ -model with augmentation	$12.36\pm0.31$	$\textbf{5.56} \pm \textbf{0.10}$	
Temporal ensembling with augmentation	$\textbf{12.16} \pm \textbf{0.24}$	$5.60\pm0.10$	



	Error rate (%) with # labels		
	4000	All (50000)	
Supervised-only	$35.56 \pm 1.59$	$7.33\pm0.04$	
with augmentation	$34.85 \pm 1.65$	$6.05\pm0.15$	
Conv-Large, $\Gamma$ -model (Rasmus et al., 2015)	$20.40\pm0.47$		
CatGAN (Springenberg, 2016)	$19.58\pm0.58$		
GAN of Salimans et al. (2016)	$18.63 \pm 2.32$		
Π-model	$16.55\pm0.29$	$6.90\pm0.07$	
$\Pi$ -model with augmentation	$12.36\pm0.31$	$\textbf{5.56} \pm \textbf{0.10}$	
Temporal ensembling with augmentation	$\textbf{12.16} \pm \textbf{0.24}$	$5.60\pm0.10$	

Why better than supervised on the full dataset?



	Error rate (%) with # labels		
	4000	All (50000)	
Supervised-only	$35.56 \pm 1.59$	$7.33\pm0.04$	
with augmentation	$34.85 \pm 1.65$	$6.05\pm0.15$	
Conv-Large, $\Gamma$ -model (Rasmus et al., 2015)	$20.40\pm0.47$		
CatGAN (Springenberg, 2016)	$19.58\pm0.58$		
GAN of Salimans et al. (2016)	$18.63 \pm 2.32$		
Π-model	$16.55\pm0.29$	$6.90\pm0.07$	
$\Pi$ -model with augmentation	$12.36\pm0.31$	$\textbf{5.56} \pm \textbf{0.10}$	
Temporal ensembling with augmentation	$\textbf{12.16} \pm \textbf{0.24}$	$5.60\pm0.10$	

Why better than supervised on the full dataset?  $\rightarrow$  Consistency regularisation



Table 3: CIFAR-100 results with 10000 labels, averages of 10 runs (4 runs for all labels).

	Error rate (%) with # labels		
	10000	All (50000)	
Supervised-only	$51.21 \pm 0.33$	$29.14 \pm 0.25$	
with augmentation	$44.56 \pm 0.30$	$26.42 \pm 0.17$	
Π-model	$43.43 \pm 0.54$	$29.06 \pm 0.21$	
Π-model with augmentation	$39.19\pm0.36$	$26.32\pm0.04$	
Temporal ensembling with augmentation	$\textbf{38.65} \pm \textbf{0.51}$	$\textbf{26.30} \pm \textbf{0.15}$	



### Semi-supervised learning on CIFAR-100 + TinyImages

Table 3: CIFAR-100 results with 10000 labels, averages of 10 runs (4 runs for all labels).

	Error rate (%) with # labels		
	10000	All (50000)	
Supervised-only	$51.21 \pm 0.33$	$29.14 \pm 0.25$	
with augmentation	$44.56 \pm 0.30$	$26.42\pm0.17$	
Π-model	$43.43 \pm 0.54$	$29.06 \pm 0.21$	
$\Pi$ -model with augmentation	$39.19\pm0.36$	$26.32\pm0.04$	
Temporal ensembling with augmentation	$\textbf{38.65} \pm \textbf{0.51}$	$\textbf{26.30} \pm \textbf{0.15}$	

Table 4: CIFAR-100 + Tiny Images results, averages of 10 runs.

	Error rate (%) with # unlabeled auxiliary inputs from Tiny Images		
	Random 500k	Restricted 237k	
Π-model with augmentation	$25.79 \pm 0.17$	$25.43 \pm 0.32$	
Temporal ensembling with augmentation	$\textbf{23.62}\pm\textbf{0.23}$	$\textbf{23.79} \pm \textbf{0.24}$	



## **Robustness to noisy labels on SVHN dataset**

Standard supervised

Temporal ensembling





## **Robustness to noisy labels on SVHN dataset**

Standard supervised

Temporal ensembling



Temporal ensembling gets >90% accuracy with 80% false labels. **Discussion**: How can that be?



## **Connections to other approaches**

- Ensembling over past predictions takes a lot of memory, changes slowly
   → Ensemble over past weights of the model: Mean Teacher
- Making outputs for differently augmented versions match: Consistency Regularisation, FixMatch, Contrastive Learning
- Soft knowledge of old model (temporal ensemble prediction) gets transferred to new model: ~ knowledge distillation, teacher-student approaches
- Model's predictions are used as soft labels for the next time step. Hard labels → Pseudo Label / Self-Training, a form of Entropy Regularization
- Encouraging local smoothness by adding noise: Virtual Adversarial Training (VAT), Noisy Student



### **Meta Pseudo Labels**

Hieu Pham, Zihang Dai, Qizhe Xie, Minh-Thang Luong, Quoc V. Le - CVPR 2021



### **Meta Pseudo Labels**



Problem in regular pseudo-labels: Confirmation bias

Solution: Use labeled data as validation of student's learning success, adjust the teacher to construct better pseudo-labels



# Meta Pseudo Labels: Algorithm

- 1. Pre-train the teacher model on the labeled dataset
- 2. Create hard pseudo labels from teacher's prediction on unlabeled data
- 3. Train student on this batch of pseudo-labeled data
- 4. Check student's learning progress on batch *L* of labeled data
- 5. Compute teacher's gradient based on:
  - a. Student's performance on *L*
  - b. Teacher's performance on the labeled batch *L*
  - c. Teacher's performance on consistency regulariser on L
- 6. Update teacher's parameters
- 7. Go to 2.



## **Comparison on TwoMoon dataset**

Supervised

Pseudo Labels

Meta Pseudo Labels



- MLP with two hidden layers, eight hidden nodes each
- 1000 points per cluster, three points labeled per cluster
- Pseudo labels based on supervised model



## **Comparison on TwoMoon dataset**

Supervised

Pseudo Labels

Meta Pseudo Labels



- MLP with two hidden layers, eight hidden nodes each
- 1000 points per cluster, three points labeled per cluster
- Pseudo labels based on supervised model

**Discussion:** Why does the pseudo labels model not learn the supervised decision boundary?



## **Experimental settings**

#### Training details

- Architectures that are commonly used in previous work
- Hyperparameters the same as in previous work
- Only dropped some augmentations in RandAugment that don't make sense
- Fine-tune on labeled data after teacher-student-training
- Last model checkpoint is evaluated, not chosen via validation set, since labeled data is limited



## **Experimental settings**

#### Training details

- Architectures that are commonly used in previous work
- Hyperparameters the same as in previous work
- Only dropped some augmentations in RandAugment that don't make sense
- Fine-tune on labeled data after teacher-student-training
- Last model checkpoint is evaluated, not chosen via validation set, since labeled data is limited

#### Comparisons

- Only compare to methods using the same architecture
- No comparison to methods using self-distillation or distillation from bigger teacher

"[...] since it is known that larger architectures and distillation can improve any method, possibly including Meta Pseudo Labels."



### **Experiments on small datasets**

	Mathad	CIFAR-10-4K	SVHN-1K	ImageNet-10%		
	Memou	$(\text{mean} \pm \text{std})$	$(\text{mean} \pm \text{std})$	Top-1	Top-5	
	Temporal Ensemble [35]	$83.63 \pm 0.63$	$92.81 \pm 0.27$		_	
	Mean Teacher [64]	$84.13 \pm 0.28$	$94.35\pm0.47$			
	VAT + EntMin [44]	$86.87 \pm 0.39$	$94.65\pm0.19$	—	83.39	
	LGA + VAT [30]	$87.94 \pm 0.19$	$93.42\pm0.36$		_	
Label Propagation Methods	ICT [71]	$92.71 \pm 0.02$	$96.11\pm0.04$		_	
Laber Propagation Methods	MixMatch [5]	$93.76\pm0.06$	$96.73 \pm 0.31$		_	
	ReMixMatch [4]	$94.86 \pm 0.04$	$97.17 \pm 0.30$		_	
	EnAET [72]	94.65	97.08		_	
	FixMatch [58]	$95.74 \pm 0.05$	$97.72 \pm 0.38$	71.5	89.1	
	UDA* [76]	$94.53 \pm 0.18$	$97.11 \pm 0.17$	68.07	88.19	
	SimCLR [8, 9]	_	_	71.7	90.4	
	MOCOv2 [10]	—	—	71.1	—	
Self-Supervised Methods	PCL [38]	—	—	—	85.6	
	PIRL [43]	—	—	—	84.9	
	BYOL [21]	—	—	68.8	89.0	
	Meta Pseudo Labels	$\textbf{96.11} \pm \textbf{0.07}$	$\textbf{98.01} \pm \textbf{0.07}$	73.89	91.38	
	Supervised Learning with full dataset*	$94.92 \pm 0.17$	$97.41 \pm 0.16$	76.89	93.27	



### Large-scale experiment: Adding unlabeled data

Mathad	# Danama	aroma Extra Data		geNet	ImageNet-ReaL [6]
Method	# Parains	Extra Data	Top-1	Top-5	Precision@1
ResNet-50 [24]	26M	_	76.0	93.0	82.94
ResNet-152 [24]	60M	—	77.8	93.8	84.79
DenseNet-264 [28]	34M	—	77.9	93.9	—
Inception-v3 [62]	24M		78.8	94.4	83.58
•••					
EfficientNet-B7 [63]	66M	_	85.0	97.2	_
EfficientNet-B7 + FixRes [70]	66M	_	85.3	97.4	_
EfficientNet-L2 [63]	480M	_	85.5	97.5	_
ResNet-50 Billion-scale SSL [79]	26M	3.5B labeled Instagram	81.2	96.0	_
ResNeXt-101 Billion-scale SSL [79]	193M	3.5B labeled Instagram	84.8	_	—
ResNeXt-101 WSL [42]	829M	3.5B labeled Instagram	85.4	97.6	88.19
FixRes ResNeXt-101 WSL [69]	829M	3.5B labeled Instagram	86.4	98.0	89.73
Big Transfer (BiT-L) [33]	928M	300M labeled JFT	87.5	98.5	90.54
Noisy Student (EfficientNet-L2) [77]	480M	300M unlabeled JFT	88.4	98.7	90.55
Noisy Student + FixRes [70]	480M	300M unlabeled JFT	88.5	98.7	—
Vision Transformer (ViT-H) [14]	632M	300M labeled JFT	88.55	_	90.72
EfficientNet-L2-NoisyStudent + SAM [16]	480M	300M unlabeled JFT	88.6	98.6	_
Meta Pseudo Labels (EfficientNet-B6-Wide)	390M	300M unlabeled JFT	90.0	98.7	91.12
Meta Pseudo Labels (EfficientNet-L2)	480M	300M unlabeled JFT	90.2	<b>98.8</b>	91.02



### **Meta Pseudo Labels - lite version**

Reduced Meta Pseudo Labels

- 1. Train large teacher model T until convergence
- 2. Pre-compute all *soft* pseudo-labels
- Small MLP as teacher T' that is trained together with the student: Input: pre-computed distributions Output: adjusted distributions
- 4. Train T' and student with Meta Pseudo Labels

Performs slightly better than Noisy Student (+1% on ImageNet, less on smaller datasets)

86.9% on ImageNet instead of 90.2%, but smaller architecture and smaller, different unlabeled dataset.



## Conclusion

#### Pros:

- Semi-supervised method that works on rather small and very large datasets
- Can beat pure supervised learning while using much fewer labels (!)
- Very interesting way to use labeled data as validation set
- Lite version for 'regular' users

#### Cons:

- Unclear how good the reduced model is, since it uses a different experimental setup
- Possibly unfair comparisons on small datasets, since they have the advantage of distillation
- Unclear why they don't use soft pseudo labels and gradient descent



#### References

**Temporal Ensembling for Semi-Supervised Learning** - Samuli Laine, Timo Aila - ICLR 2017 - <u>https://arxiv.org/abs/1610.02242</u>

Meta Pseudo Labels - Hieu Pham, Zihang Dai, Qizhe Xie, Minh-Thang Luong, Quoc V. Le - CVPR 2021 https://arxiv.org/abs/2003.10580v4

Mean teachers are better role models: Weight-averaged consistency targets improve semisupervised deep learning results - Antti Tarvainen, Harri Valpola - NIPS 2017 https://arxiv.org/abs/1703.01780

**FixMatch: Simplifying Semi-Supervised Learning with Consistency and Confidence** - Kihyuk Sohn, David Berthelot, Chun-Liang Li, Zizhao Zhang, Nicholas Carlini, Ekin D. Cubuk, Alex Kurakin, Han Zhang, Colin Raffel - NeurIPS 2020 - <u>https://arxiv.org/abs/2001.07685</u>

Virtual Adversarial Training: A Regularization Method for Supervised and Semi-Supervised Learning -Takeru Miyato, Shin-Ichi Maeda, Masanori Koyama, Shin Ishii - IEEE Transactions on Pattern Analysis and Machine Intelligence (Volume: 41, Issue: 8, Aug. 1 2019) https://ieeexplore.ieee.org/abstract/document/8417973

Self-training with Noisy Student improves ImageNet classification - Qizhe Xie, Minh-Thang Luong, Eduard Hovy, Quoc V. Le - CVPR 2020 - <u>https://arxiv.org/abs/1911.04252</u>

Semi-supervised Learning by Entropy Minimization - Yves Grandvalet, Yoshua Bengio - NIPS 2014 https://www.researchgate.net/profile/Y\_Bengio/publication/221618545\_Semi-supervised\_Learning\_by\_Entro py\_Minimization/links/546b702a0cf2f5eb18091df0.pdf