





### Discovering and Overcoming Limitations of Noise-engineered Data-free Knowledge Distillation

<u>Piyush Raikwar</u><sup>1</sup>, Deepak Mishra<sup>2</sup>
 <sup>1</sup>ABV-IIITM, Gwalior, India <sup>2</sup>IIT Jodhpur, India

# Data-free knowledge distillation

Traditionally, we assume the availability of original data.

Data-free distillation; we do not have the original data.

**Prior works** 

• Optimize input.



<sup>•</sup> Train a generative model.



Figure 1: Prior works

# The most straightforward alternative

#### **Random noise**

- Easy to generate, almost no computational burden.
- Previous attempts at distillation using Gaussian noise.
- Should not work directly, as it is basically some gibberish to teacher.

But, technically..

- Different input distribution.
- Covariate shift in hidden layer activations.

### How to make it work?

Use current statistics instead of running statistics in teacher.



### A toy example



Figure 3: (Left) Circles data on which the MLP is trained. (Middle) Gaussian noise used as input to the trained MLP. (Right) Scatter plot for embeddings in different cases.

### Student's perspective

Make student accustomed to original data.

- 1. Use current statistics in student while evaluation.
- 2. Adjust student's running statistics by just feed forwarding some original data.

```
Algorithm 1 Training - KD
```

```
Requires: pretrained teacher T(.)

Initialize: student S(.;\theta) with parameters \theta

for B in 1, 2, ..., \mathcal{B}_1 do

G \sim \mathcal{N}(0, 1)

y_T \leftarrow T(G|\mu_B, \sigma_B)

y_S \leftarrow S(G|\theta, \mu_B, \sigma_B)

\theta \leftarrow \theta - \eta \frac{\partial L_{KD}}{\partial \theta}

end for
```

#### Algorithm 2 Evaluation

```
Requires: pretrained student S(.; \theta)
for B in 1, 2, ..., \mathcal{B}_2 do
X \sim D
y_S \leftarrow S(X|\theta, \mu_B, \sigma_B)
y_{label} \leftarrow argmax(y_S)
end for
```

### Experiments



Table 2: Results on other datasets

### Other observations

- 1. Larger the batch size during training the better.
- 2. Larger the batch size during inference the better.

But, they have to be just enough, e.g., 256 batch size is sufficient.



3. Handling partial BN layers helps partially.

${\cal P}$	Student accuracy
100 (Running statistics)	13.49
90	18.26
75	57.73
50	79.54
25	82.24
0 (Current statistics)	89.4

#### Table 3: Percent BN layers using running statistics

4. More the data for adjusting the running statistics of student the better.

### Conclusion

- We show how covariate shift interferes with data-free distillation.
- We propose an approach to mitigate it to a significant extent and show that KD is possible using just Gaussian noise.
- We might not necessarily need realistic data, at least for KD. Thus we lay the foundations for noise-engineered data-free distillation.

### Future work

- Noise of lower resolutions.
- Various other noises, such as fractals.
- Applying the proposed method to other domains like transfer learning and domain adaptation.
- Use proposed method to complement other data-free distillation approaches.

# Thank you for the attention!



Contact at: piyush.raikwar@cern.ch