Continuous Learning of Deep Neural Networks

Eric Koepke, Sebastian Freytag, Martin König, Sabrina Richter

TUM Data Innovation Lab

In Cooperation with



05.08.2019

Roadmap

1. Problem Definition

- 2. Methods Supervised
- 3. Results Supervised
- 4. Methods Semi-Supervised
- 5. Results Semi-Supervised

What is PreciBake?

PreciBake is an AI company, that among other things, works on automatic baking program selection.



What is behind this technology? Or: The life cycle of a data scientist at PreciBake

- 1. The data scientist gets an almost infinite stream of incoming data and lets someone label it.
- 2. He trains a model, applies it, watches it excitedly.
- 3. Then: Performance drops!
- 4. Repeat.

What happened?

Q: Why does the performance drop after some time?

A: The input distribution is slightly different to the distribution of the training data, e.g.

- the class distribution changed (it is carnival and no one wants pretzels, everyone wants 'Krapfen')
- the image data changed (the camera got dirty or the lighting changed)
- maybe even new products are being baked

Solution(?): Iteratively feed new data to the model for training.

The problem: Catastrophic Forgetting

Solution(?): Iteratively feed new data to the model for training.

 \rightarrow Model will overfit on recent data and loose performance on old data!

This phenomenon is called *Catastrophic Forgetting*.

Our Benchmark Dataset: CORe50



- 50 objects grouped into 10 classes
- each in 11 different settings
- images per object are frames of 15s films, delivering 300 images each

How can our model adapt to new conditions without forgetting previously learned knowledge.

And can we even improve our model by feeding in more and more data without training it from scratch?

The research field that deals with this kind of issues is called *Continual Learning*.

Roadmap

- 1. Problem Definition
- 2. Methods Supervised
- 3. Results Supervised
- 4. Methods Semi-Supervised
- 5. Results Semi-Supervised

State of the Art

Catastrophic forgetting has mainly been addressed with three types of different approaches:

- ► Ensembles: Accumulate different classifiers for different tasks → Learn++
- Regularization: Protect parameters which are important for previous tasks
 - \rightarrow Synaptic Intelligence
- Memory: Keep fractions of old data and feed in gradually
 - \rightarrow Gradient Episodic Memory
 - $\rightarrow \text{Memory Replay}$



Task 2











Failure of Regularization

















Memory Replay



Memory Replay

Joined Training



Roadmap

- 1. Problem Definition
- 2. Methods Supervised
- 3. Results Supervised
- 4. Methods Semi-Supervised
- 5. Results Semi-Supervised

Results on CORe50 dataset



The PreciBake Dataset



32k pictures of 12 classes over a period of 9 months

Results on PreciBake Data



Class Distribution of the PreciBake Dataset



Roadmap

- 1. Problem Definition
- 2. Methods Supervised
- 3. Results Supervised
- 4. Methods Semi-Supervised
- 5. Results Semi-Supervised

Motivation

Reduce annotation effort by using unlabeled data from the oven's camera.



Motivation

Reduce annotation effort by using unlabeled data from the oven's camera.



Network and Loss Design



Feature Space Regularization: \mathcal{L}_{RW} (unsupervised)

Random walks based on similarity graphs

Intuition: "Points forming tight structures over the feature space should hold similar labels."



Feature Space Regularization: \mathcal{L}_{RW} (unsupervised)

Random walks based on similarity graphs

Realization: Similarity matrix $\Gamma \in [0, 1]^{N_c \times N_c}$, where N_c denotes the number of classes.



Feature Space Regularization: \mathcal{L}_{RW} (unsupervised)

Random walks based on similarity graphs

Realization: Similarity matrix $\Gamma \in [0, 1]^{N_c \times N_c}$, where N_c denotes the number of classes.



Feature Space Regularization: \mathcal{L}_{VAT} (unsupervised)

Virtual Adversarial Training

Intuition: "Points close in the input space should be close in the feature space."

$$\mathcal{L}_{\mathsf{VAT}} = \sum_{i=1}^{B_u} D(f_{ heta}(x_i), f_{ heta}(x_i + \epsilon_{\mathsf{adv}})),$$

where f_{θ} denotes the feature space embedding, D denotes the Kullback-Leibler divergence, and B_u denotes the batch size of unlabeled data.

Feature Space Regularization: \mathcal{L}_{Center} (supervised)

Centering of clusters in feature space

Intuition: "Penalize points that are far from their class center."

$$\mathcal{L}_{ ext{Center}} = \sum_{i=1}^{B_l} ||f_ heta(x_i) - c_{y_i}||_2^2$$

where f_{θ} denotes the feature space embedding, c_{y_i} denotes the y_i th class center in feature space, and B_i denotes the batch size of labeled data.



Roadmap

- 1. Problem Definition
- 2. Methods Supervised
- 3. Results Supervised
- 4. Methods Semi-Supervised
- 5. Results Semi-Supervised

Validation Results



Validation Results



Validation Results

Reduced prediction error from 0.0175 to 0.015 (ca. 15%).



Center Loss Performance

Explanation: Joint supervision of cross-entropy and center loss **increases inter-class distance** and **smoothens intra-class variation** respectively.



Thank you for your attention.



Appendix 1: Formulas Synaptic Intelligence

$$ilde{\mathcal{L}}_i = \mathcal{L}_i + c \sum_k \Omega^i_k (riangle^i_k)^2$$

$$\begin{split} & \bigtriangleup_k^i = \theta_k(t_i) - \theta_k(t_{i-1}) \text{ denotes how far } \theta_k \text{ moved} \\ & \Omega_k^i = \sum_{j=0}^{i-1} \frac{\omega_k^j}{(\bigtriangleup_k^i)^2 + \epsilon} \text{ measures how much } \theta_k \text{ contributed to a drop} \\ & \text{ in the loss with} \end{split}$$

$$egin{aligned} \mathcal{L}(heta(t_i)) &- \mathcal{L}(heta(t_{i-1})) = \int_{ heta(t_{i-1})}^{ heta(t_i)} rac{\partial L}{\partial heta}(heta) d heta &= \int_{t_{i-1}}^{t_i} rac{\partial L}{\partial heta}(heta(t)) \cdot heta'(t) dt \ &= \sum_k \int_{t_{i-1}}^{t_i} rac{\partial L}{\partial heta_k}(heta(t)) heta'_k(t) dt =: -\sum_k \omega_k^i \end{aligned}$$

Appendix 2: Formulas \mathcal{L}_{RW} and \mathcal{L}_{VAT}

$$\begin{split} & \text{Minimize } L_{\text{CE}}(f_{\theta}) + L_{\text{SSL}}(f_{\theta}) \\ & \text{where } f_{\theta} \text{ denotes the feature space layer,} \\ & L_{\text{SSL}} = L_{\text{VAT}} + L_{\text{RW}}, \\ & L_{\text{VAT}} = \sum_{\mathbf{x} \in \mathbb{B}_{U}} \mathbb{KL}(f_{\theta}(\mathbf{x}) - f_{\theta}(\mathbf{x} + \epsilon)), \quad (\text{local consistency}) \\ & L_{RW} = \sum_{i=0}^{\tau} \alpha_{i} H(I, \Gamma^{(i)}) \qquad (\text{global consistency}) \\ & \Gamma^{(\tau)} = \Gamma^{(p \to x)} \cdot (\Gamma^{(x \to x)})^{\tau} \cdot \Gamma^{(x \to p)} \\ & \text{where } \Gamma^{(p \to x)} \in [0, 1]^{N_{e} \times M}, N_{c} \text{ number of classes, } M \text{ batch size,} \\ & \Gamma^{(p \to x)}_{i,j} \text{ denotes the transition probability from } p_{i} \text{ to } x_{i}, \end{split}$$

 $H(I,\Gamma)$ is the average cross-entropy between the rows of I and Γ .

$$\begin{split} H(I,\Gamma) &= -\frac{1}{N_c} \sum_{i=0}^{N_c} \log \Gamma^{(i)} \\ \Gamma^{(p \rightarrow x)} &= \operatorname{Softmax}(\mathbf{A}^{\mathrm{T}}) \\ \Gamma^{(x \rightarrow p)} &= \operatorname{Softmax}(\mathbf{A}) \\ \Gamma^{(x \rightarrow x)} &= \operatorname{Softmax}(\mathbf{B}) \\ A_{i,j} &= -||f_{\theta}(x_i) - p_j||^2 \\ B_{i,j} &= -||f_{\theta}(x_i) - f_{\theta}(x_j)||^2 \\ p_j &= \frac{1}{N_c} \sum_{x_i \in class(j)} f_{\theta}(x_i) \end{split}$$

Appendix 2: \mathcal{L}_{RW} and Absence of Classes

Problem: \mathcal{L}_{RW} unstable wrt. absence of classes.

