

Interpretable Deep Canonical Correlation Analysis Based Neural Networks

Yunlyu Wang
Matri.-Nr.:6898496

Master Thesis
Electrical System Engineering
University of Paderborn

First Examiner: Prof. Dr. Peter Schreier
Second Examiner: Prof. Dr.-Ing. Reinhold Häb-Umbach
Supervisors: Dr.-Ing. Tanuj Hasija, Maurice Kuschel



Contents

- 1 Introduction
- 2 Background
- 3 Implementation
- 4 Visualization Techniques
- 5 Experiments and Results
- 6 Conclusion and Outlook

Contents

- 1 Introduction
- 2 Background
- 3 Implementation
- 4 Visualization Techniques
- 5 Experiments and Results
- 6 Conclusion and Outlook

Introduction

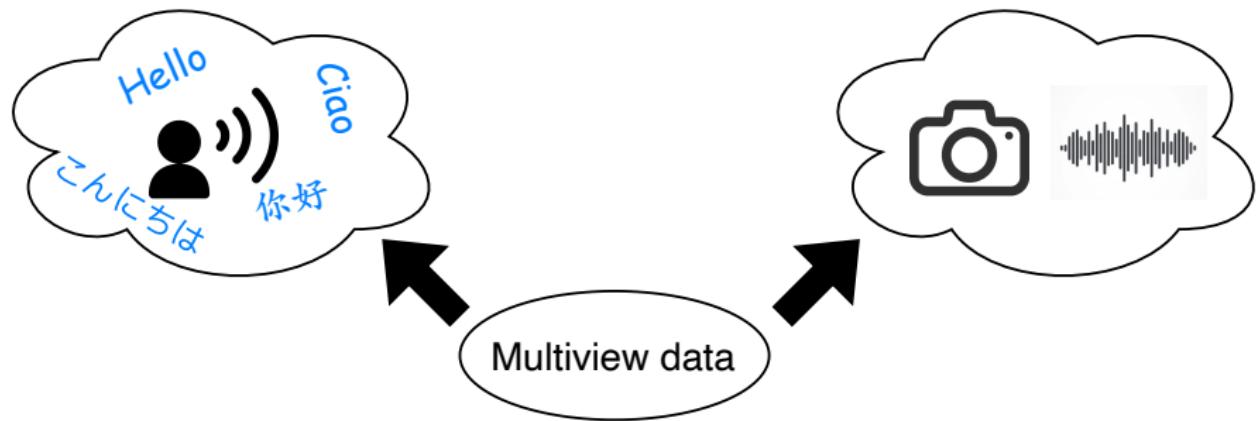
- Multi-view analysis is popular in machine learning

Introduction

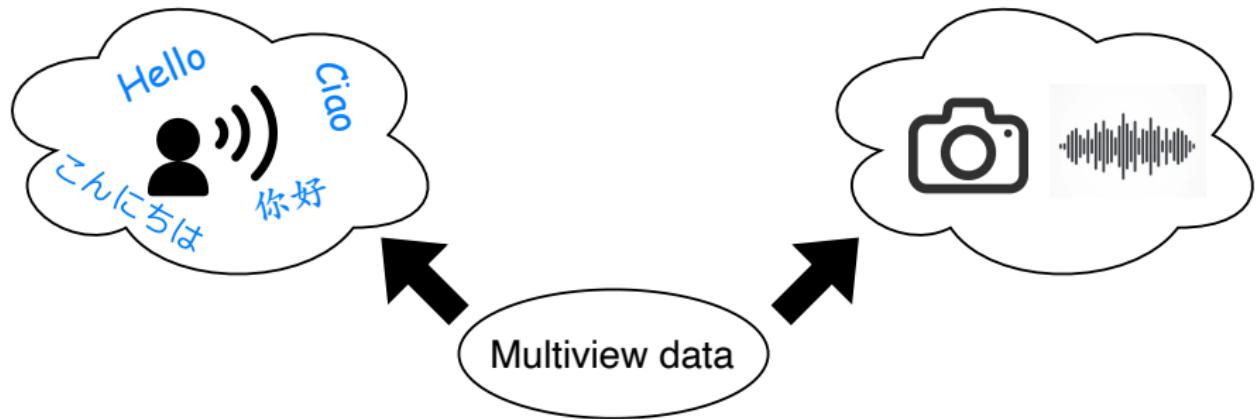
- Multi-view analysis is popular in machine learning
- Multi-view data acquires the data from different views

Introduction

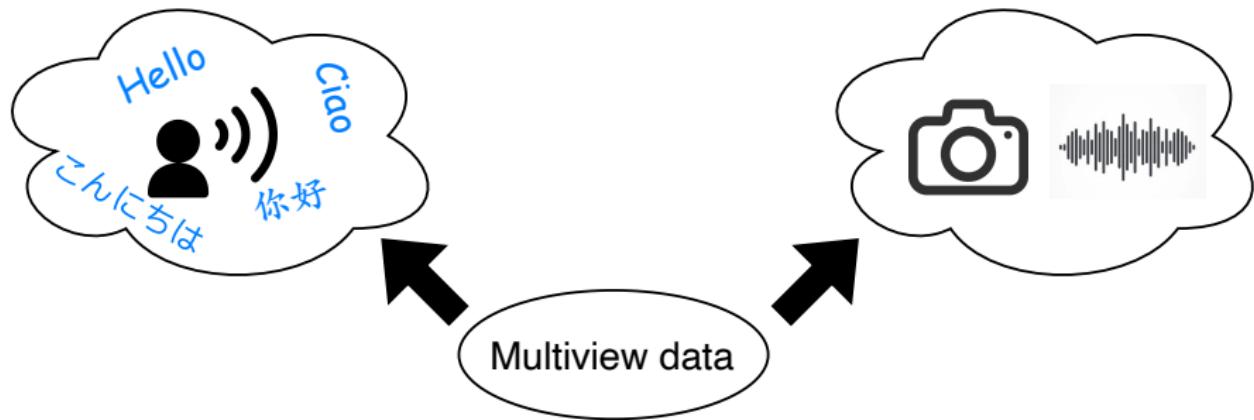
- Multi-view analysis is popular in machine learning
- Multi-view data acquires the data from different views



Introduction



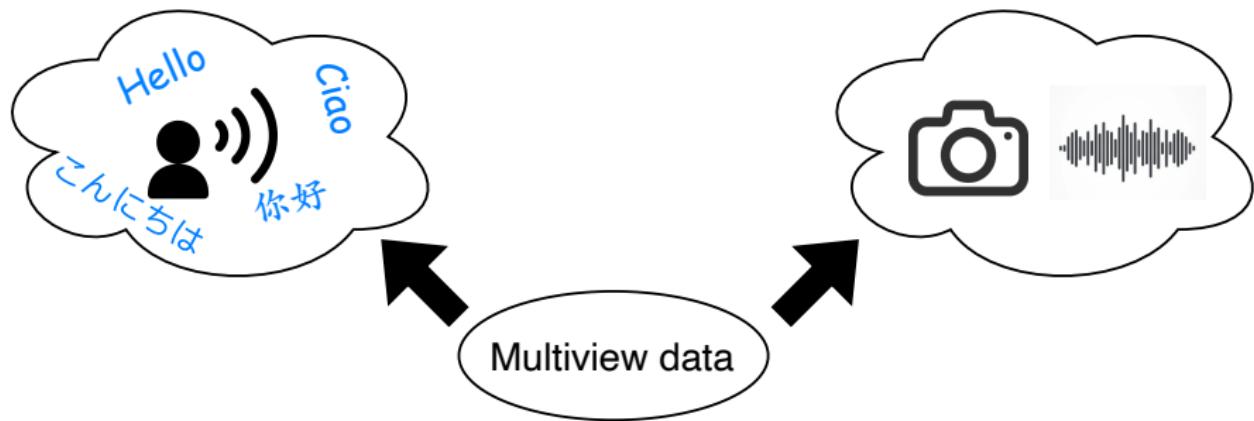
Introduction



Characteristics

- more information
- unsupervised learning
- to learn better representations

Introduction



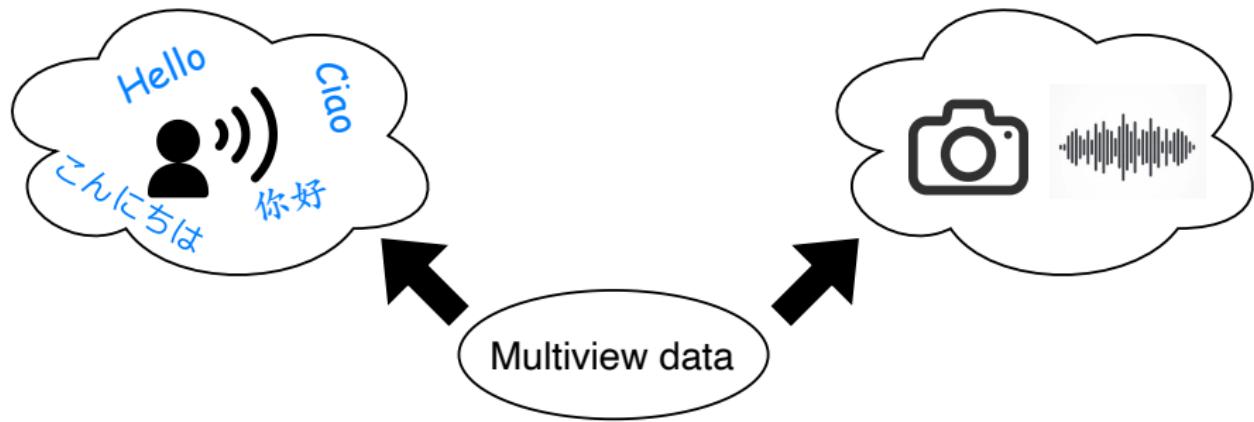
Characteristics

- more information
- unsupervised learning
- to learn better representations

Applications

- dimensionality reduction
- multi-label classification
- image denoising

Introduction



Multiview analysis

- correlation-maximization approaches

Contents

- 1 Introduction
- 2 Background
- 3 Implementation
- 4 Visualization Techniques
- 5 Experiments and Results
- 6 Conclusion and Outlook

CCA

Correlation-maximization approaches:

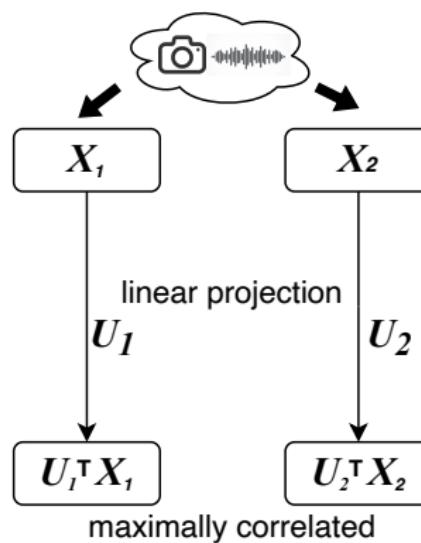
- Canonical Correlation Analysis(CCA) ¹

¹ Harold Hotelling. Relations between two sets of variates. In Breakthroughs in statistics, pages 162–190. Springer, 1992.

CCA

Correlation-maximization approaches:

- Canonical Correlation Analysis(CCA)¹



¹Harold Hotelling. Relations between two sets of variates. In Breakthroughs in statistics, pages 162–190. Springer, 1992.

CCA

The objective of CCA

- $\mathbf{X}_1 \in \mathbb{R}^{n_1 \times N}, \mathbf{X}_2 \in \mathbb{R}^{n_2 \times N}$ with covariances Σ_{11}, Σ_{22} and cross-covariance Σ_{12} . CCA finds the linear projections that:

$$\begin{aligned} & \max_{\mathbf{U}_1, \mathbf{U}_2} \text{Tr}(\mathbf{U}_1^\top \Sigma_{12} \mathbf{U}_2) \\ \text{s.t. } & \mathbf{U}_1^\top \Sigma_{11} \mathbf{U}_1 = \mathbf{I}, \\ & \mathbf{U}_2^\top \Sigma_{22} \mathbf{U}_2 = \mathbf{I}, \\ & \mathbf{u}_1^{(i)\top} \mathbf{X}_1 \mathbf{X}_2^\top \mathbf{u}_2^{(j)} = 0, \text{ for } i \neq j. \end{aligned}$$

- When covariances matrices are ill-conditioned, they can be replaced with their regularized version, e.g.

$$\Sigma_{11} = \frac{1}{N} \mathbf{X}_1 \mathbf{X}_1^\top + r_1 \mathbf{I}$$

where r_1 is regularization parameter

CCA

The solution of CCA

Define:

$$\mathbf{T} \triangleq \boldsymbol{\Sigma}_{11}^{-1/2} \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1/2}$$

Then take singular value decomposition of \mathbf{T} :

$$\mathbf{T} = \mathbf{F} \mathbf{K} \mathbf{G}^\top$$

Then the optimal objective value is:

$$(\mathbf{U}_1^*, \mathbf{U}_2^*) = (\boldsymbol{\Sigma}_{11}^{-\frac{1}{2}} \mathbf{F}, \boldsymbol{\Sigma}_{22}^{-\frac{1}{2}} \mathbf{G})$$

DCCA

Correlation-maximization approaches:

- Canonical Correlation Analysis(CCA)

DCCA

Correlation-maximization approaches:

- Canonical Correlation Analysis(CCA) → linear transformation

DCCA

Correlation-maximization approaches:

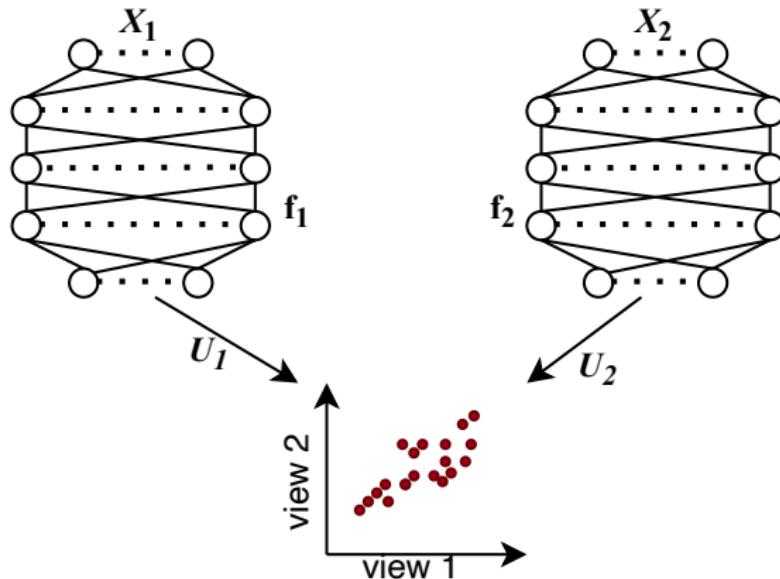
- Canonical Correlation Analysis(CCA) → linear transformation
- Deep Canonical Correlation Analysis(DCCA)²

²Andrew, Galen, et al. "Deep canonical correlation analysis." International conference on machine learning. PMLR, 2013.

DCCA

Correlation-maximization approaches:

- Canonical Correlation Analysis(CCA) → linear transformation
- Deep Canonical Correlation Analysis(DCCA)²



²Andrew, Galen, et al. "Deep canonical correlation analysis." International conference on machine learning. PMLR, 2013.

DCCA

The objective of DCCA

There are $\mathbf{X}_1 \in \mathbb{R}^{n \times N}$, $\mathbf{X}_2 \in \mathbb{R}^{n \times N}$, \mathbf{f}_1 and \mathbf{f}_2 denote mapping implemented by DNNs, with a corresponding set of learnable parameters, $\mathbf{W}_{\mathbf{f}_1}$ and $\mathbf{W}_{\mathbf{f}_2}$

$$\begin{aligned} & \max_{\mathbf{W}_{\mathbf{f}_1}, \mathbf{W}_{\mathbf{f}_2}, \mathbf{U}_1, \mathbf{U}_2} \frac{1}{N} \text{Tr}(\mathbf{U}_1^\top \mathbf{f}_1(\mathbf{X}_1) \mathbf{f}_2(\mathbf{X}_2)^\top \mathbf{U}_2) \\ & \text{s.t. } \mathbf{U}_1^\top \left(\frac{1}{N} \mathbf{f}_1(\mathbf{X}_1) \mathbf{f}_1(\mathbf{X}_1)^\top + r_1 \mathbf{I} \right) \mathbf{U}_1 = \mathbf{I} , \\ & \quad \mathbf{U}_2^\top \left(\frac{1}{N} \mathbf{f}_2(\mathbf{X}_2) \mathbf{f}_2(\mathbf{X}_2)^\top + r_2 \mathbf{I} \right) \mathbf{U}_2 = \mathbf{I} , \\ & \quad \mathbf{u}_1^{(i)\top} \mathbf{f}_1(\mathbf{X}_1) \mathbf{f}_2(\mathbf{X}_2)^\top \mathbf{u}_2^{(j)} = 0, \text{for } i \neq j . \end{aligned}$$

where \mathbf{U}_1 and \mathbf{U}_2 are the CCA projections that project the DNN outputs, r_1 and r_2 are regularization parameters for sample covariance estimation.

DCCAE

Correlation-maximization approaches:

- Canonical Correlation Analysis(CCA)
- Deep Canonical Correlation Analysis(DCCA)

DCCAE

Correlation-maximization approaches:

- Canonical Correlation Analysis(CCA)
- Deep Canonical Correlation Analysis(DCCA)
 - ▶ shared representation →[?] reconstruct original views

DCCAE

Correlation-maximization approaches:

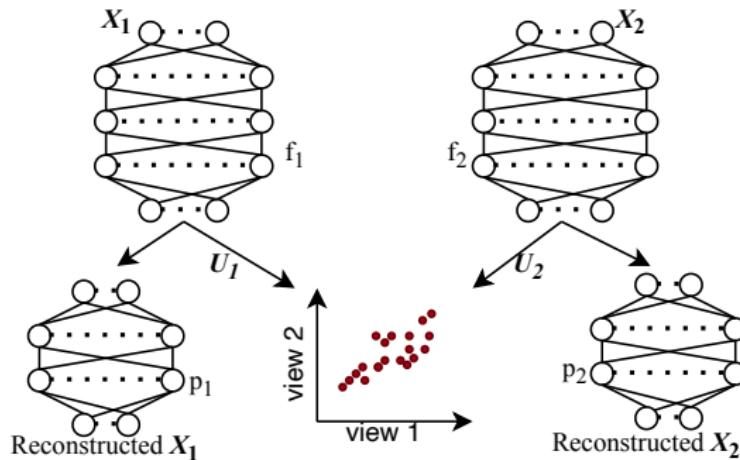
- Canonical Correlation Analysis(CCA)
- Deep Canonical Correlation Analysis(DCCA)
 - ▶ shared representation →[?] reconstruct original views
- Deep Canonically Correlated Autoencoders(DCCAE)³

³Wang, Weiran, et al. "On deep multi-view representation learning." International conference on machine learning. PMLR, 2015.

DCCAE

Correlation-maximization approaches:

- Canonical Correlation Analysis(CCA)
- Deep Canonical Correlation Analysis(DCCA)
 - ▶ shared representation →[?] reconstruct original views
- Deep Canonically Correlated Autoencoders(DCCAE)³



³Wang, Weiran, et al. "On deep multi-view representation learning." International conference on machine learning. PMLR, 2015.

DCCAE

The objective of DCCAE

There are two view data $\mathbf{X}_1 \in \mathbb{R}^{n \times N}$, $\mathbf{X}_2 \in \mathbb{R}^{n \times N}$. $\mathbf{f}_1, \mathbf{f}_2$ are the two DNNs of finding the maximal correlated outputs, and $\mathbf{p}_1, \mathbf{p}_2$ denote the two DNNs of reconstruction, with a corresponding set of learnable parameters, $\mathbf{W}_{\mathbf{f}_1}, \mathbf{W}_{\mathbf{f}_2}, \mathbf{W}_{\mathbf{p}_1}$ and $\mathbf{W}_{\mathbf{p}_2}$.

$$\begin{aligned} & \min_{\mathbf{W}_{\mathbf{f}_1}, \mathbf{W}_{\mathbf{f}_2}, \mathbf{W}_{\mathbf{p}_1}, \mathbf{W}_{\mathbf{p}_2}, \mathbf{U}_1, \mathbf{U}_2} - \frac{1}{N} \text{Tr}(\mathbf{U}_1^\top \mathbf{f}_1(\mathbf{X}_1) \mathbf{f}_2(\mathbf{X}_2)^\top \mathbf{U}_2) \\ & + \frac{\lambda}{N} \sum_{i=1}^N (\|\mathbf{x}_{1i} - \mathbf{p}_1(\mathbf{f}_1(\mathbf{x}_{1i}))\|^2 + \|\mathbf{x}_{2i} - \mathbf{p}_2(\mathbf{f}_2(\mathbf{x}_{2i}))\|^2) \end{aligned}$$

s.t. the same constraints in DCCA ,

where $\lambda > 0$ is a trade-off parameter.

Motivation

Disadvantages

Motivation

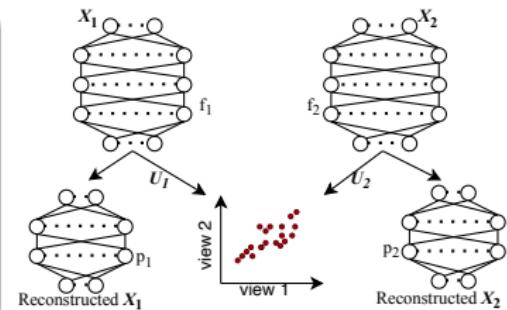
Disadvantages

- DCCAE based techniques

Motivation

Disadvantages

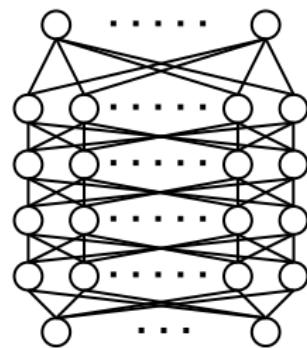
- DCCAE based techniques
 - ▶ feedforward neural network



Motivation

Disadvantages

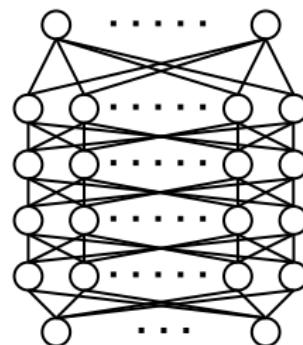
- DCCAE based techniques
 - ▶ feedforward neural network
 - ▶ not suitable for images



Motivation

Disadvantages

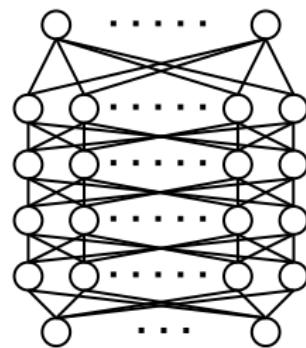
- DCCAE based techniques
 - ▶ feedforward neural network
 - ▶ not suitable for images
- Correlation $\rightarrow 1 \Rightarrow$ Overfitting



Motivation

Disadvantages

- DCCAE based techniques
 - ▶ feedforward neural network
 - ▶ not suitable for images
- Correlation $\rightarrow 1 \Rightarrow$ Overfitting
- Not interpretable



Motivation

Disadvantages

- DCCAE based techniques
 - ▶ feedforward neural network
 - ▶ not suitable for images
- Correlation $\rightarrow 1 \Rightarrow$ Overfitting
- Not interpretable

This thesis proposes

- CNN-based DCCAE model

Motivation

Disadvantages

- DCCAE based techniques
 - ▶ feedforward neural network
 - ▶ not suitable for images
- Correlation $\rightarrow 1 \Rightarrow$ Overfitting
- Not interpretable

This thesis proposes

- CNN-based DCCAE model
- Analysing overfitting

Motivation

Disadvantages

- DCCAE based techniques
 - ▶ feedforward neural network
 - ▶ not suitable for images
- Correlation $\rightarrow 1 \Rightarrow$ Overfitting
- Not interpretable

This thesis proposes

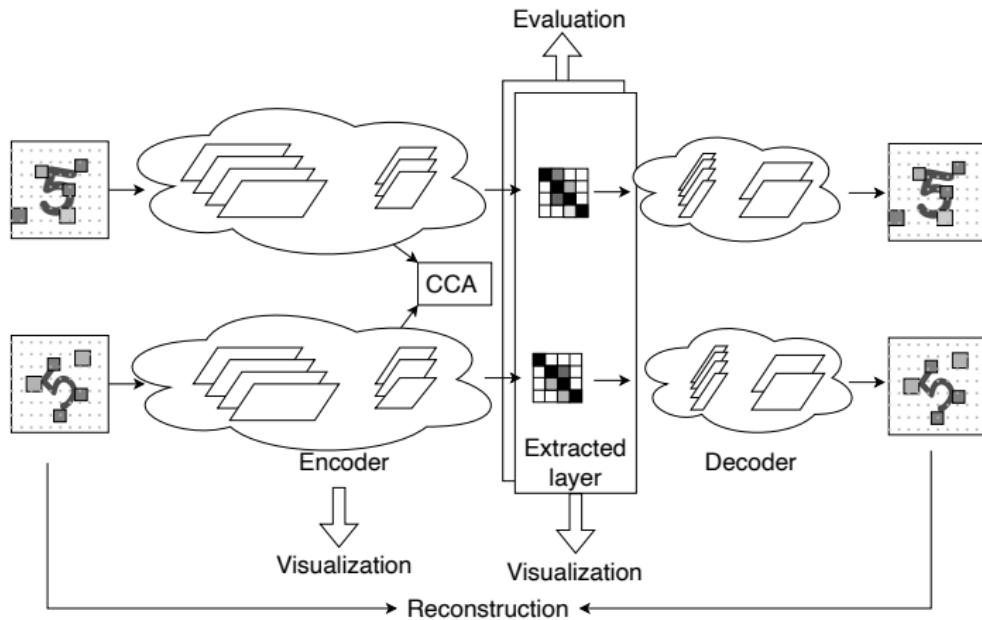
- CNN-based DCCAE model
- Analysing overfitting
- Interpretability

Contents

- 1 Introduction
- 2 Background
- 3 Implementation
- 4 Visualization Techniques
- 5 Experiments and Results
- 6 Conclusion and Outlook

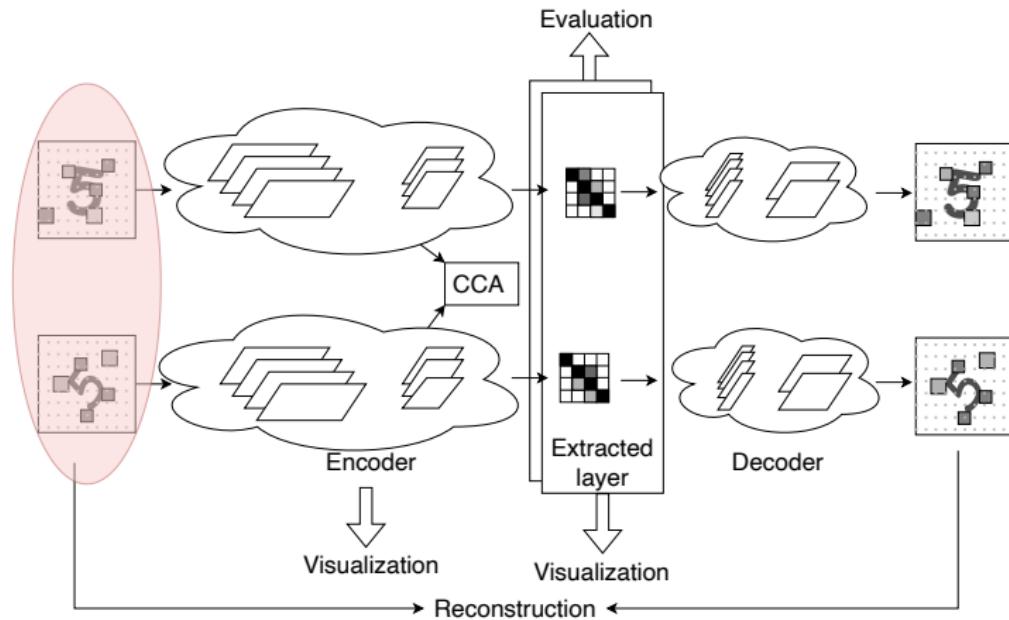
Overview

The CNN-based DCCAE network model



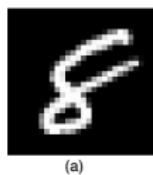
Proposed Pipeline

The CNN-based DCCAE network model



Data Generation

- Inputs: MNIST dataset



(a)



(b)



(c)



(d)

Figure: MNIST

Data Generation

- Inputs: MNIST dataset with some box noise



(a)



(b)



(c)



(d)

Figure: MNIST



(a)



(b)



(c)



(d)

Figure: MNIST with Box Nosie

Data Generation

- Inputs: MNIST dataset with some box noise



(a)



(b)



(a)



(b)



(c)



(d)



(c)



(d)

Figure: MNIST

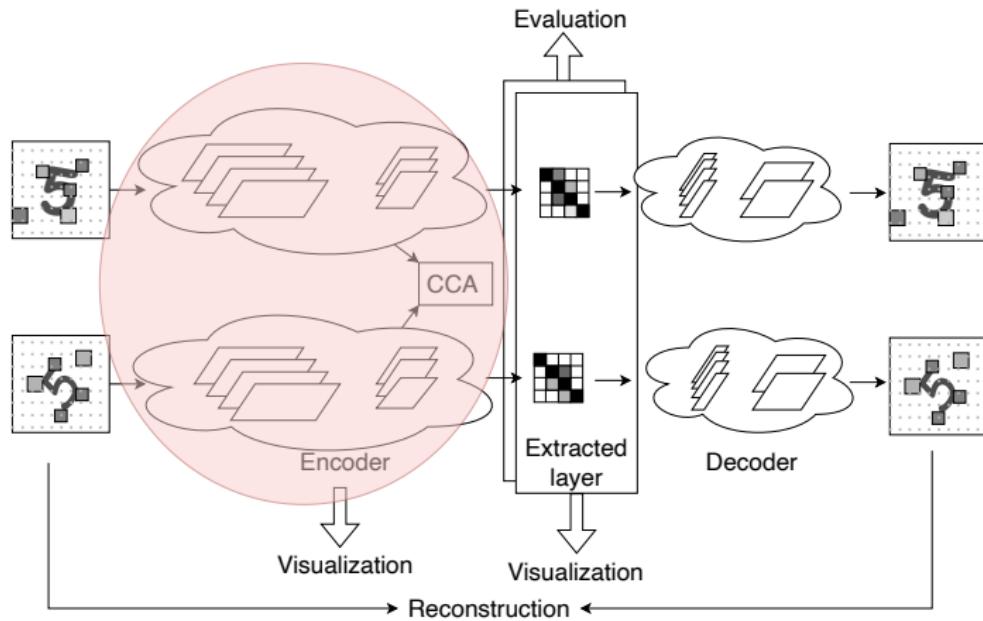
Figure: MNIST with Box Nosie

The character of the data

- The number of boxes is customizable
- The size and the position of boxes are random in each examples of both views

Proposed Pipeline

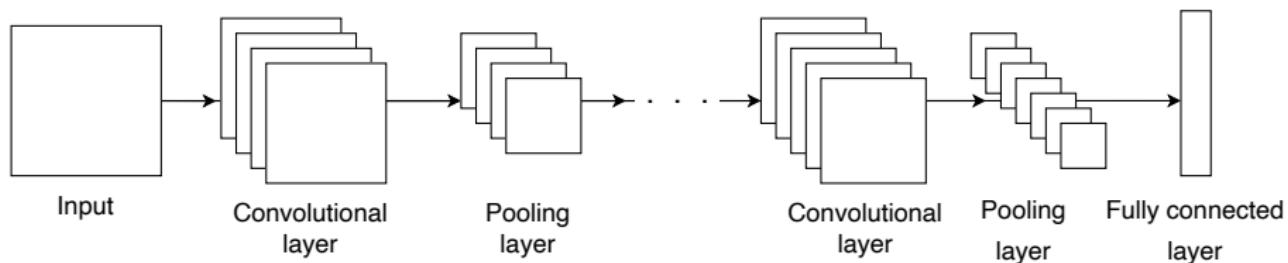
The CNN-based DCCAE network model



Convolution Neural Network

Convolution neural network ⁴ is mainly used for applications in image and speech recognition

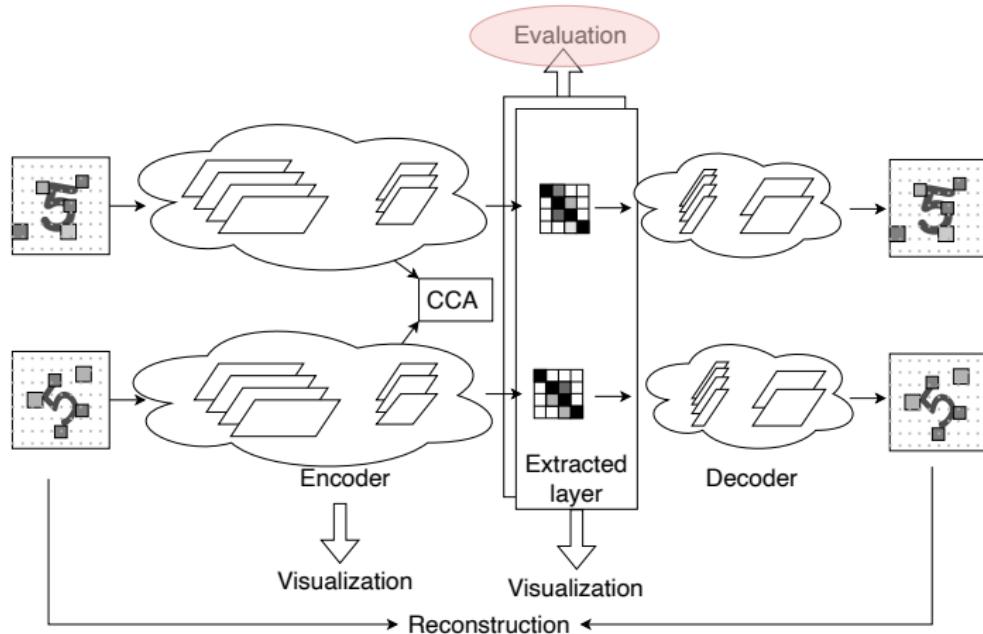
- small number of weights (shared)
- reduce high dimensionality
- without losing image information



⁴Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient- based learning applied to document recognition. Proceedings of the IEEE, 86(11):2278–2324, 1998.

Proposed Pipeline

The CNN-based DCCAE network model



Evaluation

Accuracy of the classification \implies Performance of the model.

Evaluation

Accuracy of the classification \implies Performance of the model.

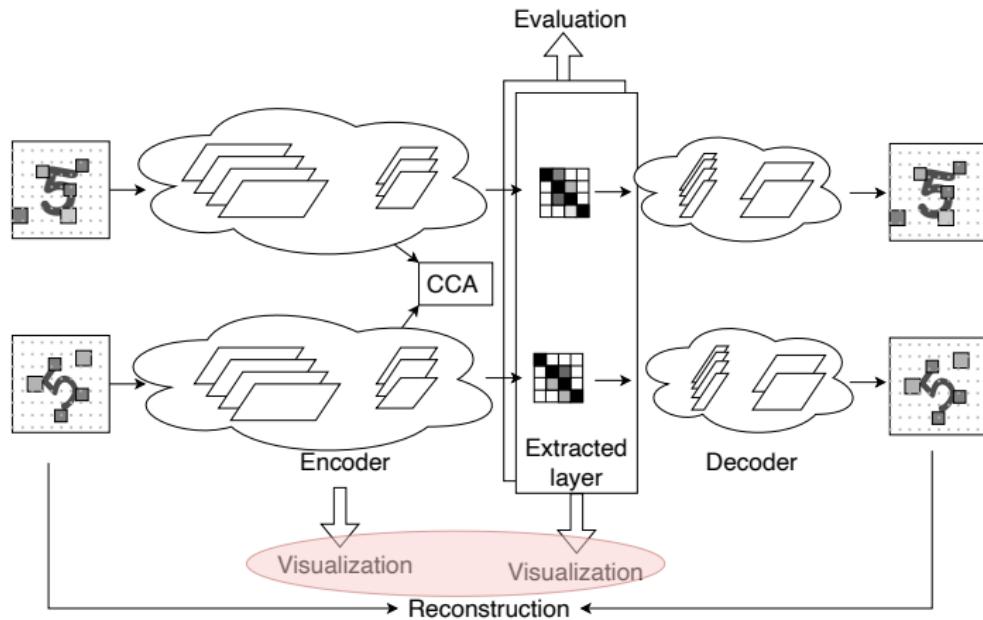
- Support vector machines (SVM)⁵
- K-means clustering⁶

⁵Vladimir N Vapnik. Support vector networks. *Machine learning*, 20(3):273–297, 1995.

⁶Anil K Jain and Richard C Dubes. *Algorithms for clustering data*. Prentice-Hall, Inc., 1988.

Proposed Pipeline

The CNN-based DCCAE network model



Visualization

Visualization of neural network \implies which parts of input the network learns

Visualization

Visualization of neural network \implies which parts of input the network learns

- Saliency Map⁷
- SmoothGrad⁸
- Gradient-weighted Class Activation Mapping(GradCAM)⁹

⁷ Niels JS Mørch, Ulrik Kjems, Lars Kai Hansen, Claus Svarer, Ian Law, Benny Lautrup, Steve Strother, and Kelly Rehm. Visualization of neural networks using saliency maps. In Proceedings of ICNN' 95-International Conference on Neural Networks, volume 4, pages 2085–2090. IEEE, 1995.

⁸ Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. Smoothgrad: removing noise by adding noise. arXiv preprint arXiv:1706.03825, 2017.

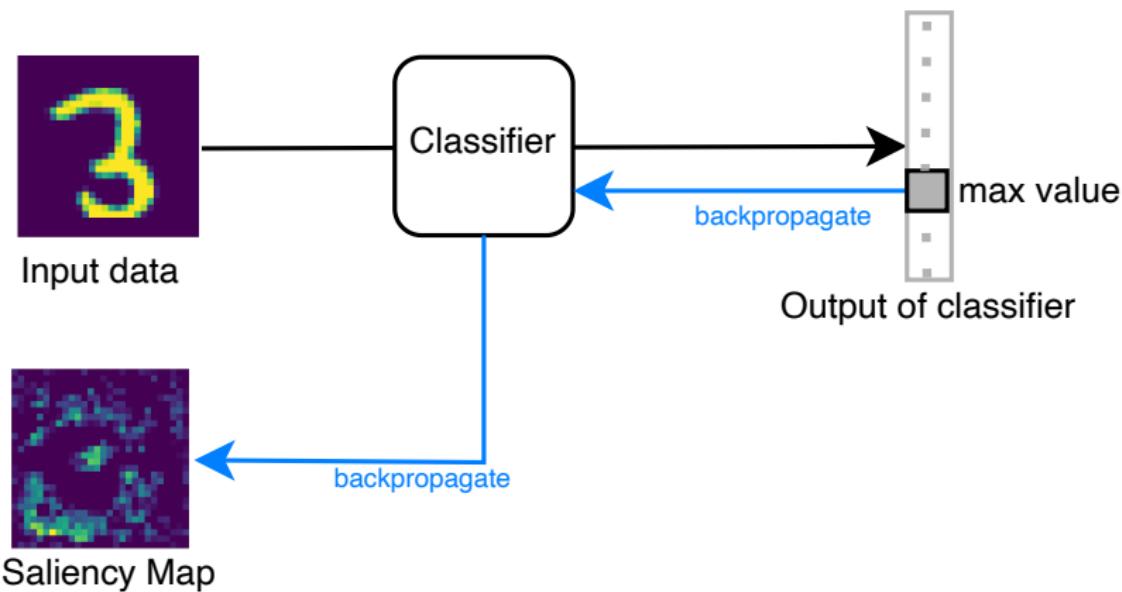
⁹ Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In Proceedings of the IEEE international conference on computer vision, pages 618–626, 2017.

Contents

- 1 Introduction
- 2 Background
- 3 Implementation
- 4 Visualization Techniques
- 5 Experiments and Results
- 6 Conclusion and Outlook

Saliency Map⁷

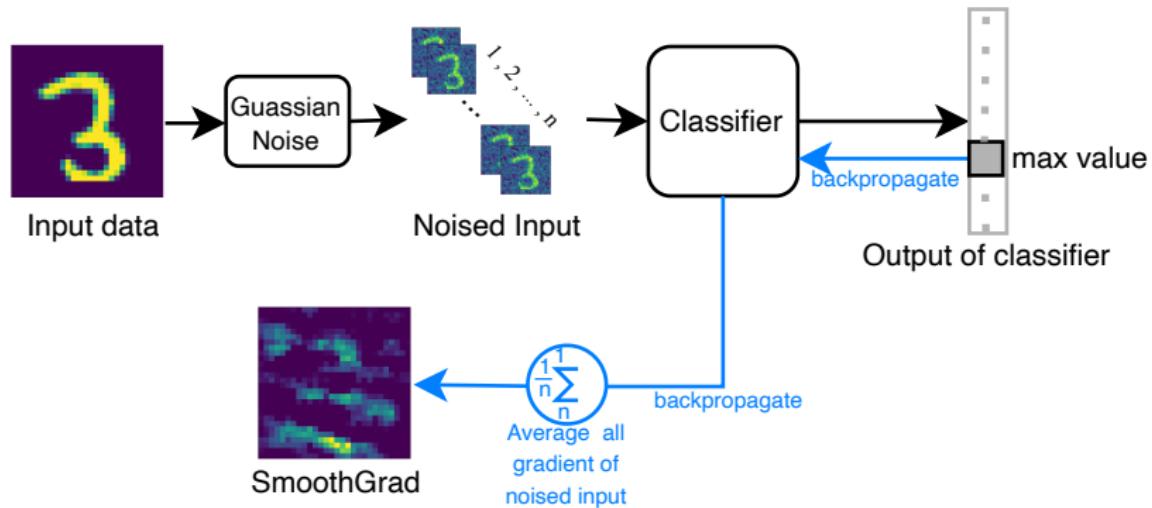
Workflow



⁷ Niels JS Mørch, Ulrik Kjems, Lars Kai Hansen, Claus Svarer, Ian Law, Benny Lautrup, Steve Strother, and Kelly Rehm. Visualization of neural networks using saliency maps. In Proceedings of ICNN'95-International Conference on Neural Networks, volume 4, pages 2085–2090. IEEE, 1995.

SmoothGrad⁸

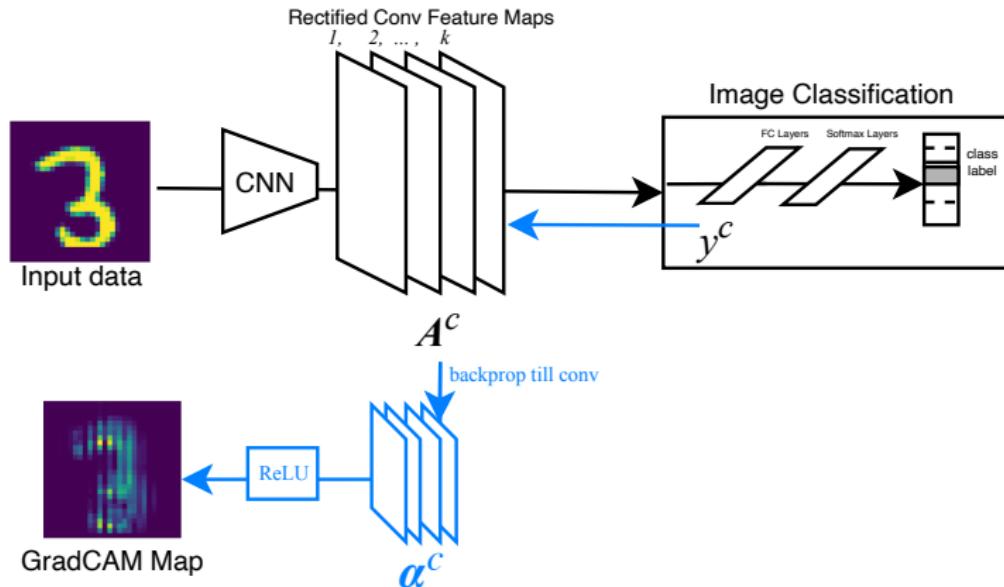
Workflow



⁸ Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. Smoothgrad: removing noise by adding noise. arXiv preprint arXiv:1706.03825, 2017.

GradCAM⁹

Workflow



⁹ Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In Proceedings of the IEEE international conference on computer vision, pages 618–626, 2017.

Contents

- 1 Introduction
- 2 Background
- 3 Implementation
- 4 Visualization Techniques
- 5 Experiments and Results
- 6 Conclusion and Outlook

Experiments and Results

Experiments

- Comparison between two network structures
- Visualizations evaluation
- Overfitting analysis

Two Network Structures

DNN model structure

| Layers | Output Shape | Activation function |
|---------------------|--------------|---------------------|
| Input | (None, 784) | |
| Layer 1(Dense) | (None, 1024) | ReLU |
| Layer 2(Dense) | (None, 1024) | ReLU |
| Layer 3(Dense) | (None, 1024) | ReLU |
| Output layer(Dense) | (None, 10) | None |

Two Network Structures

DNN model structure

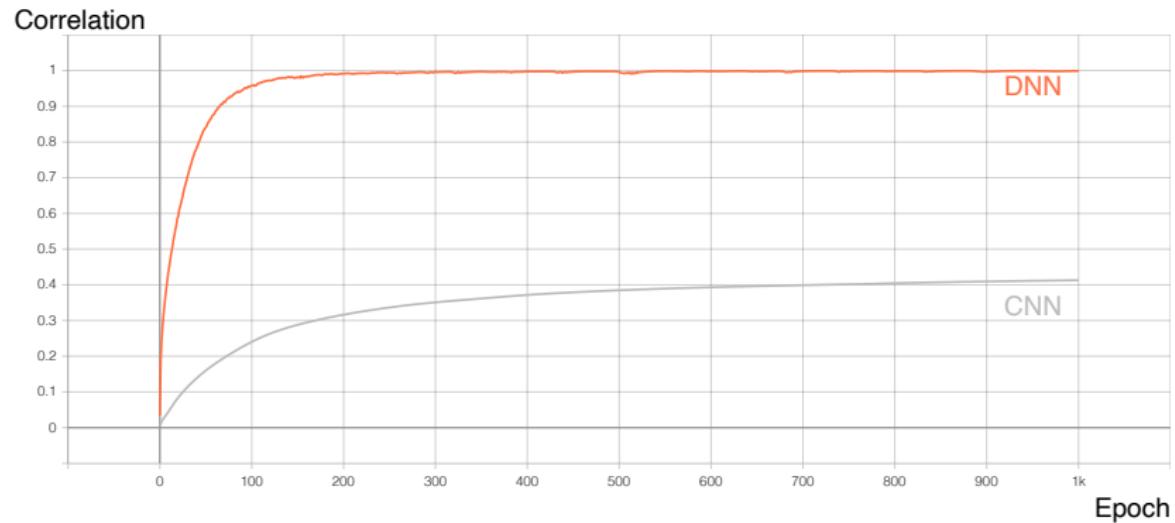
| Layers | Output Shape | Activation function |
|---------------------|--------------|---------------------|
| Input | (None, 784) | |
| Layer 1(Dense) | (None, 1024) | ReLU |
| Layer 2(Dense) | (None, 1024) | ReLU |
| Layer 3(Dense) | (None, 1024) | ReLU |
| Output layer(Dense) | (None, 10) | None |

CNN model structure

| Layers | Output Shape | Filters | Kernel / Pooling | Stride | Padding | Activation function |
|-------------------------|-----------------|---------|------------------|--------|------------|---------------------|
| Input | (None, 28,28,1) | | | | | |
| Layer 1(Conv2D) | (None, 28,28,8) | 8 | (5,5) | (1,1) | same | ReLU |
| Layer 2(MaxPooling2D) | (None, 14,14,8) | | (2,2) | (2,2) | no padding | |
| Layer 3(Conv2D) | (None, 14,14,6) | 6 | (5,5) | (1,1) | same | ReLU |
| Layer 4(MaxPooling2D) | (None, 7,7,6) | | (2,2) | (2,2) | no padding | |
| Layer 5(Conv2D) | (None, 7,7,4) | 4 | (5,5) | (1,1) | same | ReLU |
| Layer 6(MaxPooling2D) | (None, 3,3,4) | | (2,2) | (2,2) | no padding | |
| Layer 7(Conv2D) | (None, 3,3,4) | 2 | (5,5) | (1,1) | same | ReLU |
| Output layer((Flatten)) | (None, 18) | | | | | |

Comparison between Two Network Structures

Correlation



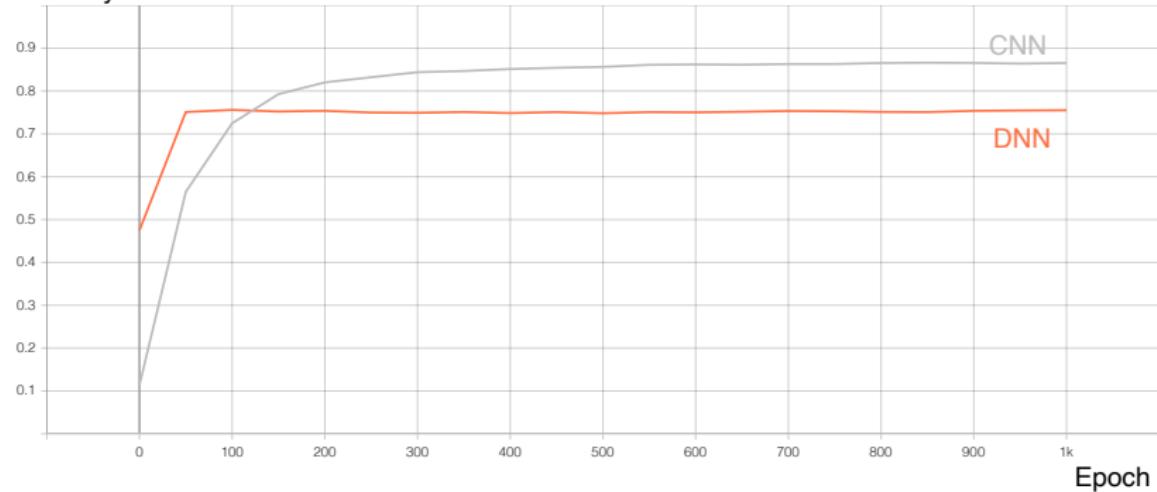
Setting

- Input: 4 boxes MNIST

Comparison between Two Network Structures

Accuracy

Accuracy of SVM



Setting

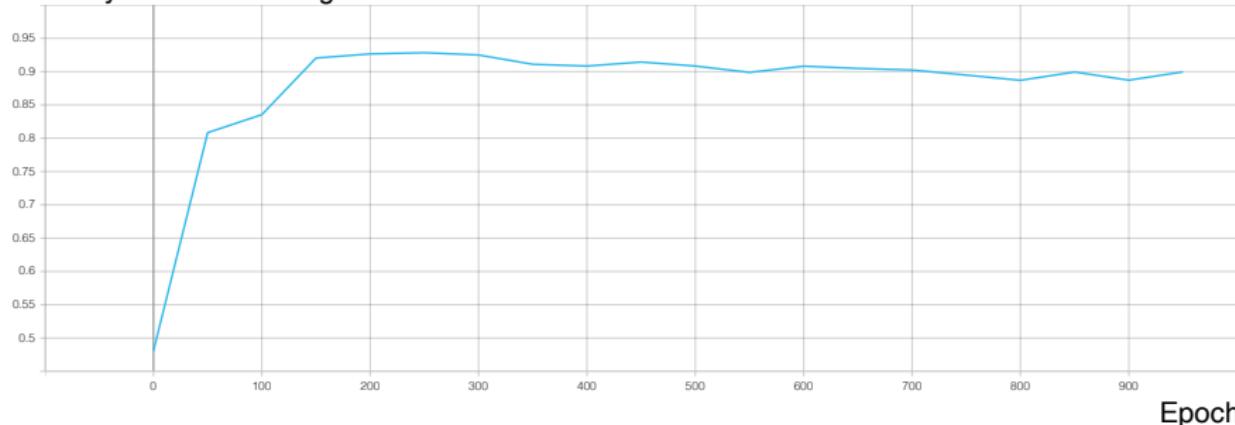
- Input: 4 boxes MNIST

Visualization Evaluation

Visualization Evaluation

DNN model

Accuracy of the Clustering



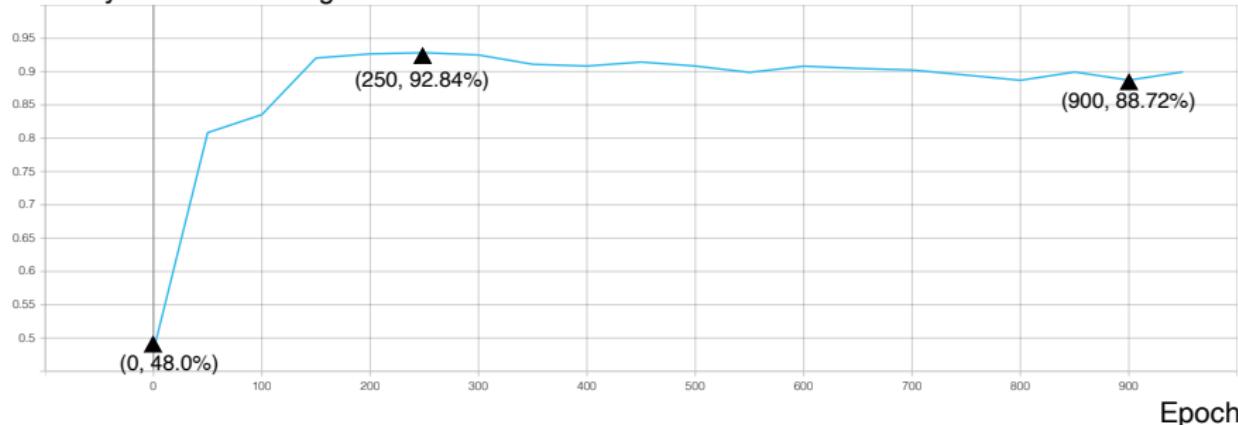
Setting

- Input: 1 box MNIST

Visualization Evaluation

DNN model

Accuracy of the Clustering

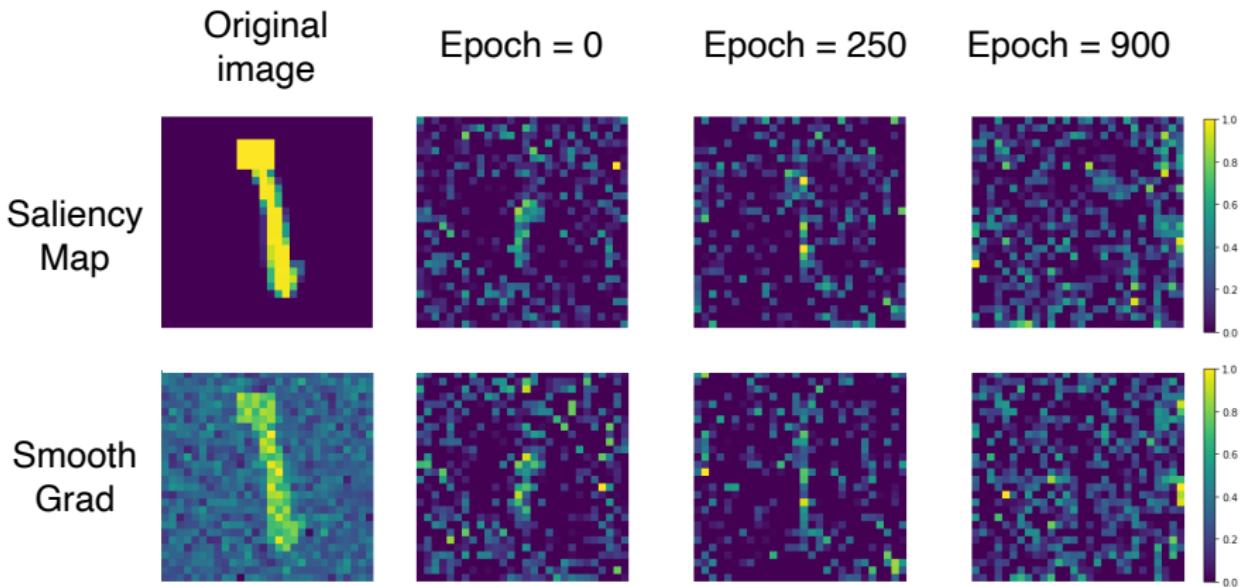


Setting

- Input: 1 box MNIST

Visualization Evaluation

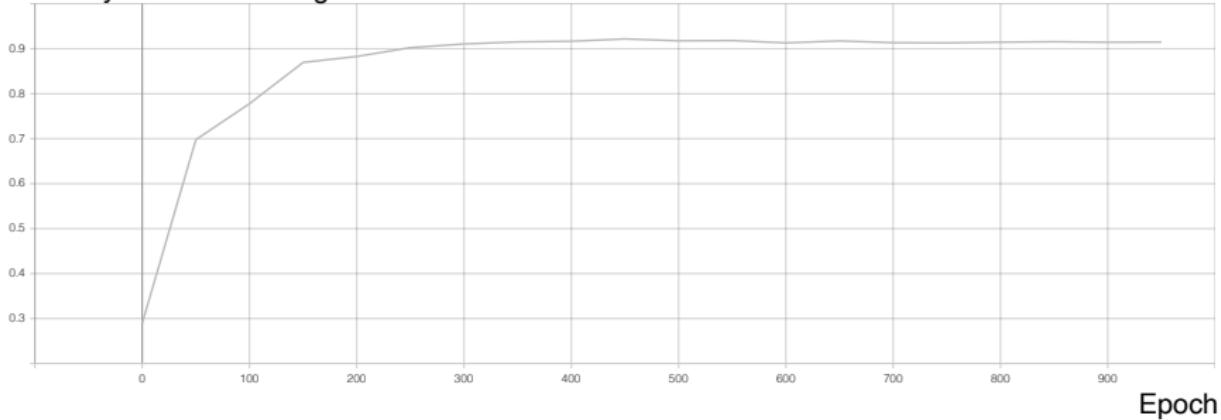
DNN model



Visualization Evaluation

CNN model

Accuracy of the Clustering



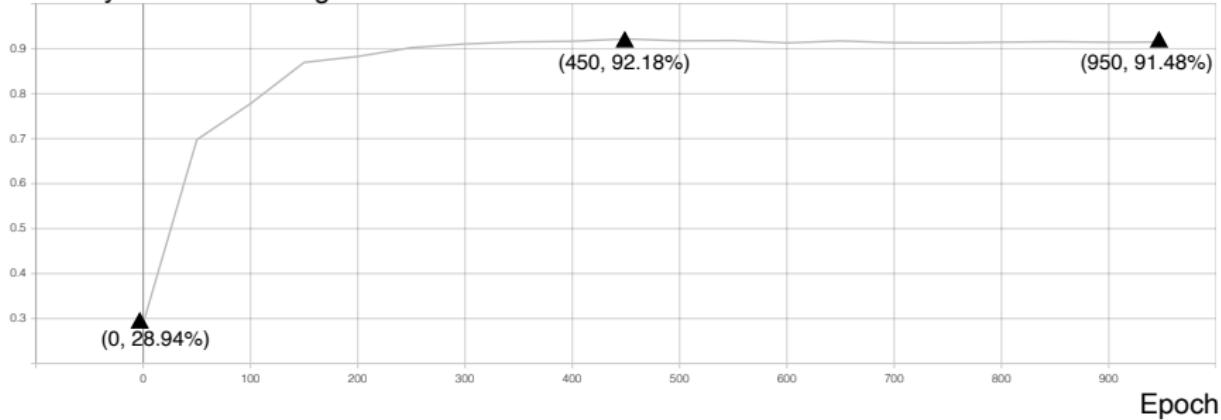
Setting

- Input: 1 box MNIST

Visualization Evaluation

CNN model

Accuracy of the Clustering

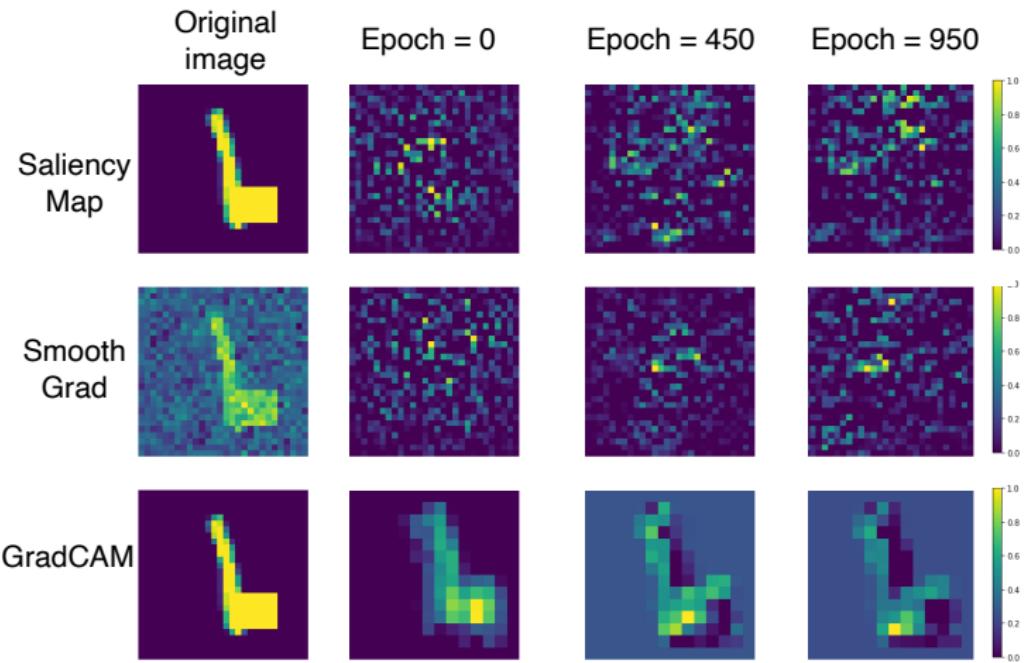


Setting

- Input: 1 box MNIST

Visualization Evaluation

CNN model



Overfitting Evaluation

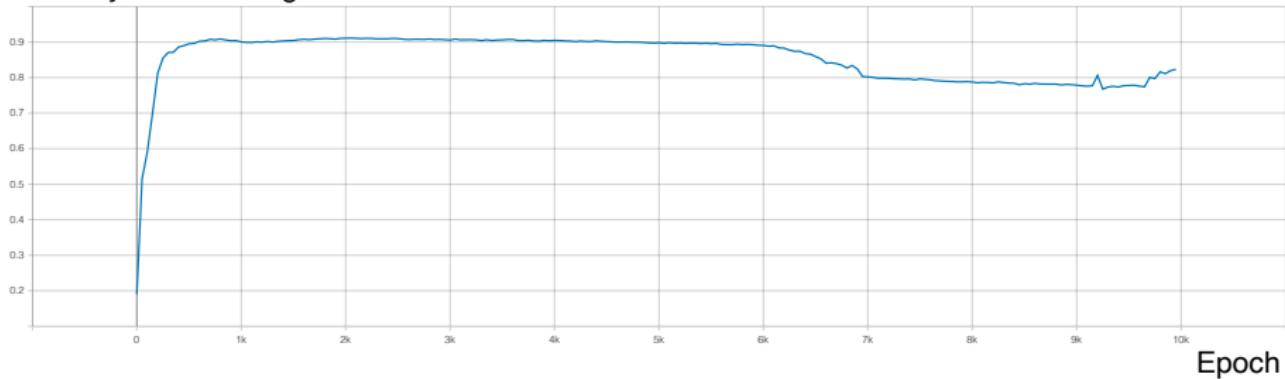
Setting

- Input: 2 boxes MNIST

Overfitting Evaluation

CNN model

Accuracy of Clustering



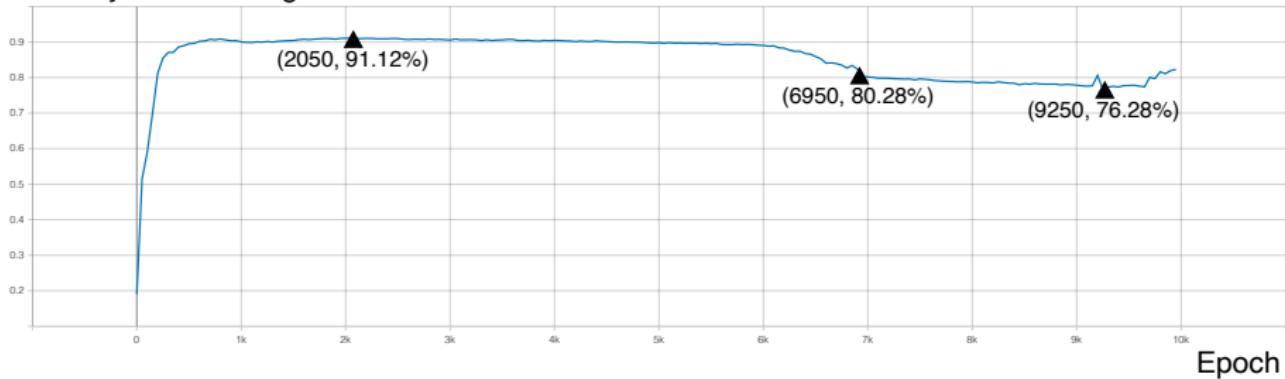
Setting

- Input: 2 boxes MNIST

Overfitting Evaluation

CNN model

Accuracy of Clustering

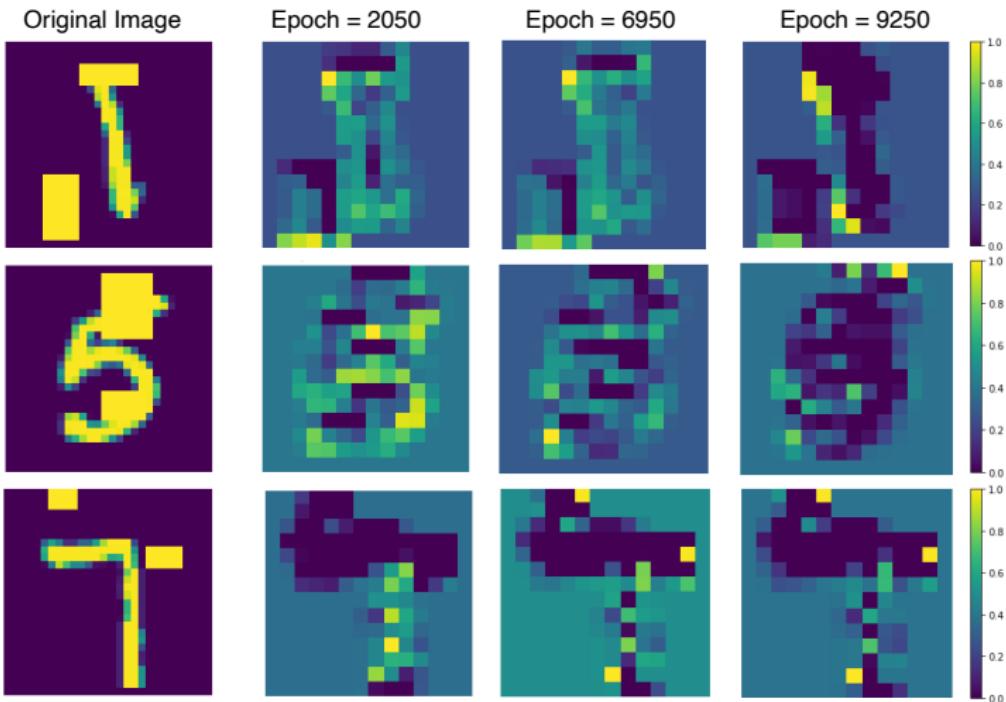


Setting

- Input: 2 boxes MNIST

Overfitting Evaluation

CNN model



Overfitting Evaluation

- Quantitative metric \Rightarrow Digit attention

Overfitting Evaluation

- Quantitative metric \Rightarrow Digit attention

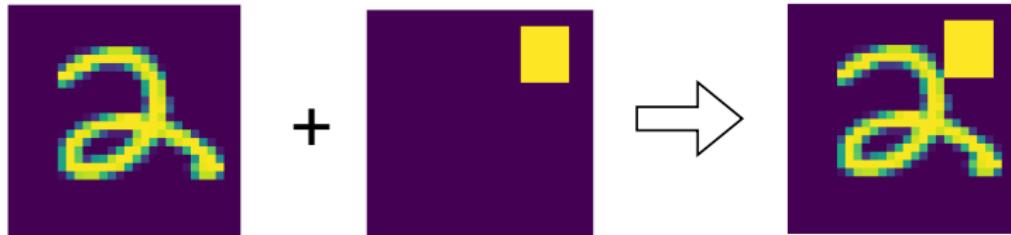
Digit Attention



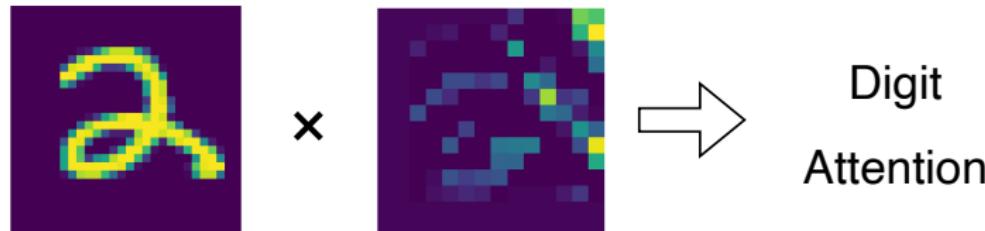
Overfitting Evaluation

- Quantitative metric \Rightarrow Digit attention

Digit Attention



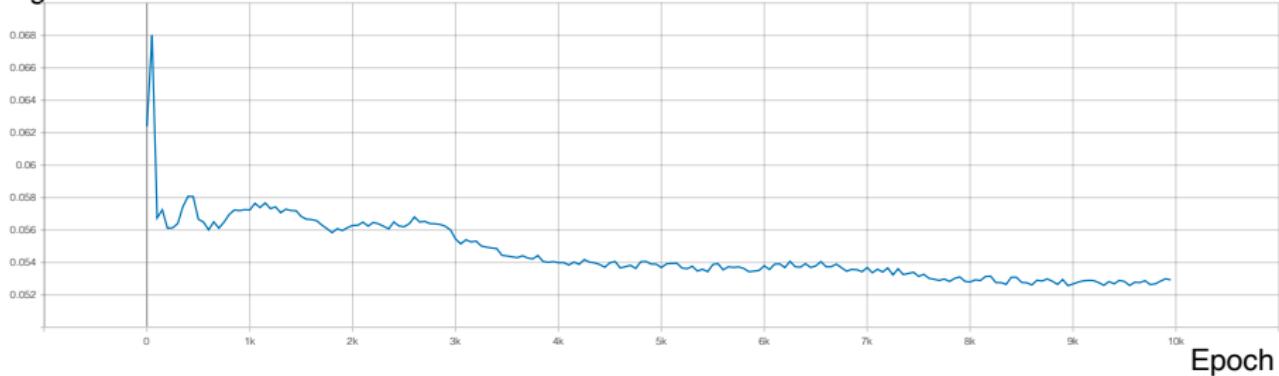
Digit Attention



Overfitting Evaluation

CNN model

Digit Attention

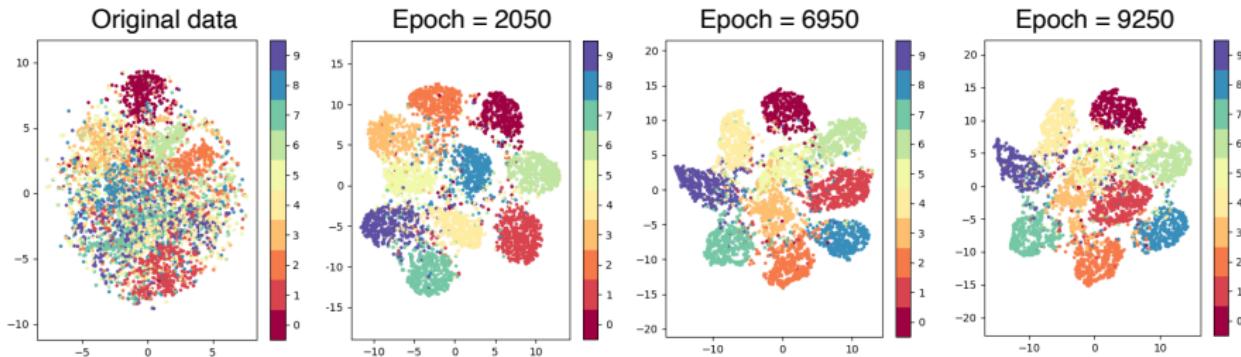


Setting

- Input: 2 boxes MNIST

Overfitting Evaluation

CNN model



Setting

- Input: 2 boxes MNIST

Contents

- 1 Introduction
- 2 Background
- 3 Implementation
- 4 Visualization Techniques
- 5 Experiments and Results
- 6 Conclusion and Outlook

Conclusion

Conclusion

Contribution

- Traditional DCCAE model

Conclusion

Contribution

- Traditional DCCAE model
- ✓ CNN-based DCCAE model
 - ▶ outperformed traditional DCCAE

Conclusion

Contribution

- Traditional DCCAE model
- ✓ CNN-based DCCAE model
▶ outperformed traditional DCCAE

Contribution

- Correlation $\rightarrow 1 \Rightarrow$ Overfitting

Conclusion

Contribution

- Traditional DCCAE model
- ✓ CNN-based DCCAE model
▶ outperformed traditional DCCAE

Contribution

- Correlation $\rightarrow 1 \Rightarrow$ Overfitting
- ✓ Analysing overfitting
▶ more accurate representations
▶ less likely to overfitting

Conclusion

Contribution

- Traditional DCCAE model
- ✓ CNN-based DCCAE model
▶ outperformed traditional DCCAE

Contribution

- Correlation $\rightarrow 1 \Rightarrow$ Overfitting
- ✓ Analysing overfitting
▶ more accurate representations
▶ less likely to overfitting

Contribution

- Not interpretable

Conclusion

Contribution

- Traditional DCCAE model
- ✓ CNN-based DCCAE model
- ▶ outperformed traditional DCCAE

Contribution

- Correlation $\rightarrow 1 \Rightarrow$ Overfitting
- ✓ Analysing overfitting
- ▶ more accurate representations
 - ▶ less likely to overfitting

Contribution

- Not interpretable
- ✓ Interpretability
- ▶ Saliency Map
 - ▶ SmoothGrad
 - ▶ GradCAM

Outlook

Outlook

Limitation

- Visualization results → Not ideal

Outlook

Limitation

- Visualization results → Not ideal
- Saliency Map, SmoothGrad → Extra classifier

Outlook

Limitation

- Visualization results → Not ideal
- Saliency Map, SmoothGrad → Extra classifier

Outlook

- More visualization techniques

Outlook

Limitation

- Visualization results → Not ideal
- Saliency Map, SmoothGrad → Extra classifier

Outlook

- More visualization techniques
 - ▶ Guided Backpropagation

Outlook

Limitation

- Visualization results → Not ideal
- Saliency Map, SmoothGrad → Extra classifier

Outlook

- More visualization techniques
 - ▶ Guided Backpropagation
- More quantitative metrics
 - ▶ Attention corporate more techniques
 - ▶ Relate with the correlation

This endeavour would not have been possible without the help of my supervisors:

Dr.-Ing. Tanuj Hasija and Mr. Maurice Kuschel

I am truly grateful for their support and guidance.

Thank you!

CCA

Moreover, notice that the solution is independent with the scalar of the projections \mathbf{u}_1 and \mathbf{u}_2 . For example, if $u_1 = \alpha\mathbf{u}_1$ where α is a scalar,

$$\begin{aligned}\frac{\alpha\mathbf{u}_1^\top \Sigma_{12} \mathbf{u}_2}{\sqrt{\alpha\mathbf{u}_1^\top \Sigma_{11} \alpha\mathbf{u}_1 \mathbf{u}_2^\top \Sigma_{22} \mathbf{u}_2}} &= \frac{\alpha\mathbf{u}_1^\top \Sigma_{12} \mathbf{u}_2}{\sqrt{\alpha^2 \mathbf{u}_1^\top \Sigma_{11} \mathbf{u}_1 \mathbf{u}_2^\top \Sigma_{22} \mathbf{u}_2}} \\ &= \frac{\mathbf{u}_1^\top \Sigma_{12} \mathbf{u}_2}{\sqrt{\mathbf{u}_1^\top \Sigma_{11} \mathbf{u}_1 \mathbf{u}_2^\top \Sigma_{22} \mathbf{u}_2}}.\end{aligned}$$

CCA

For solving the (8), there are two ways to solve it.

- Eigenvalue Decomposition (EVD)

$$\begin{aligned}\Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \mathbf{U}_1 &= \lambda_1^2 \mathbf{U}_1 \\ \Sigma_{22}^{-1} \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} \mathbf{U}_2 &= \lambda_2^2 \mathbf{U}_2\end{aligned}\tag{1}$$

- Singular Value Decomposition (SVD)

$$\begin{aligned}\Sigma_{11}^{-1/2} \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \Sigma_{11}^{-1/2} \mathbf{U}_1 &= \eta_1 \mathbf{U}_1 \\ \Sigma_{22}^{-1/2} \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22}^{-1/2} \mathbf{U}_2 &= \eta_2 \mathbf{U}_2\end{aligned}\tag{2}$$

DeepCCA

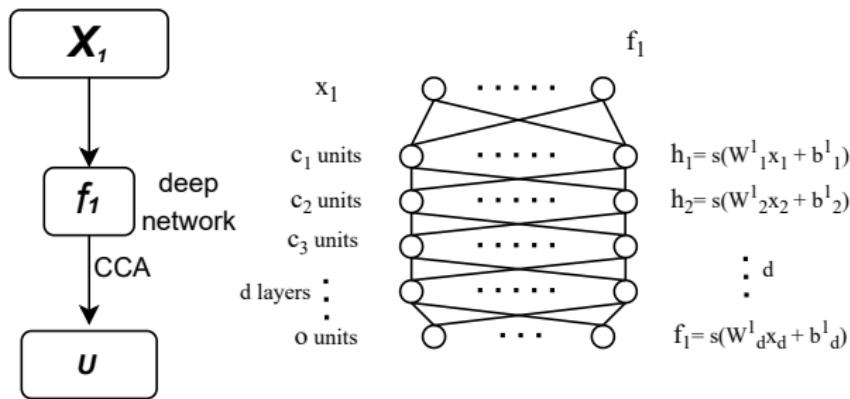


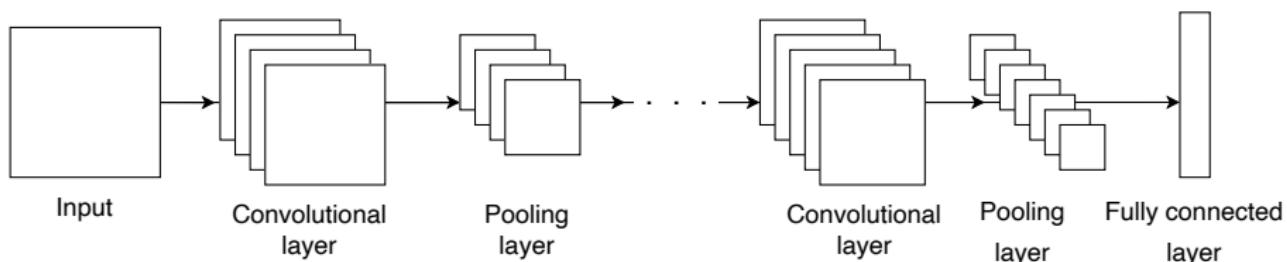
Figure: A schematic of deep CCA

Given an instance x_2 of the second view, the representation $f_2(x_2)$ is computed the same way, with different parameters W_l^2 and b_l^2 .

Convolution Neural Network

A simple convolutional neural network structure

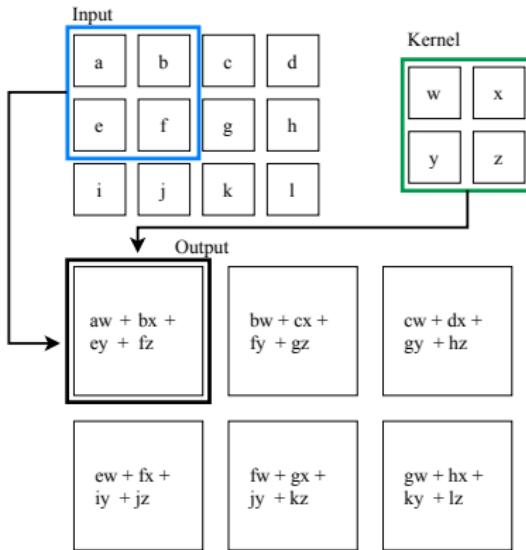
- Convolutional layers
 - ▶ Convolution operation
 - ▶ Padding and strides
- Pooling layers
- A fully connected layer



Convolution Neural Network

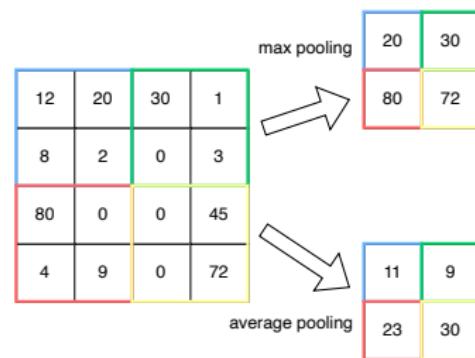
Convolutional layers

- Convolution operation
- Padding and strides



Pooling layers

- Average pooling
- Max pooling



Saliency Map

Saliency Map

Given a gray scalar image $\mathbf{I}_1 \in \mathbb{R}^{w \times h}$ and the corresponding class c_1 , the output can be calculated by the classifier function S_c and the predicted label $class(\mathbf{I}_1)$ is the index of the maximal value in the output.

$$class(\mathbf{I}_1) = \operatorname{argmax}_c S_c(\mathbf{I}_1) \quad (3)$$

To get the Saliency Map $\mathbf{M} \in \mathbb{R}^{w \times h}$, first it has to get the derivative ω of $class(\mathbf{I}_1)$ with respect to the image \mathbf{I}_1 by backpropagation

$$\begin{aligned} \omega &= \frac{\partial class(\mathbf{I}_1)}{\partial \mathbf{I}_1} \\ &= \frac{\partial \operatorname{argmax}_c S_c(\mathbf{I}_1)}{\partial \mathbf{I}_1} . \end{aligned} \quad (4)$$

Saliency Map

Saliency Map

The number of elements in ω is equal to the number of pixels in I_1 . Therefore the Saliency Map M can be taken directly from the absolute value of the ω with corresponding pixels.

$$M(i, j) = |\omega_{p(i,j)}| , \quad (5)$$

where the $p(i, j)$ denotes the index of the element of ω corresponding to the image I_1 in the i -th row and j -th column.

If the input I_1 is a multi-channel image, e.g. **RGB** picture. The Saliency Map M can take the maximal value of the channels h to be obtained.

$$M(i, j) = \max_h |\omega_{p(i,j,h)}| , \quad (6)$$

where the $p(i, j, h)$ represents the index of the pixel in the (i, j) position and the c color channel of image I_1 .

SmoothGrad

SmoothGrad

Adding noise has a de-noise effect

$$\tilde{\mathbf{I}} = \frac{1}{n} \sum_1^n (\hat{\mathbf{I}}_1 + \cdots + \hat{\mathbf{I}}_n) , \quad (7)$$

where $\hat{\mathbf{I}}_n$ denotes the noised \mathbf{I} for $n = 1, 2, \dots, n$.

$$\begin{aligned}\hat{\mathbf{I}}_1 &= \mathbf{I}_1 + \mathcal{N}(0, \sigma^2) , \\ \hat{\mathbf{I}}_2 &= \mathbf{I}_1 + \mathcal{N}(0, \sigma^2) , \\ &\vdots \\ \hat{\mathbf{I}}_n &= \mathbf{I}_1 + \mathcal{N}(0, \sigma^2) ,\end{aligned} \quad (8)$$

where $\mathcal{N}(0, \sigma^2)$ represents Gaussian noise with standard deviation σ .

SmoothGrad

SmoothGrad

Then using each noised input can calculate the corresponding noised gradient $\hat{\omega}_1, \hat{\omega}_2, \dots, \hat{\omega}_n$ by backpropagation:

$$\hat{\omega}_n = \frac{\partial \text{class}(\hat{\mathbf{I}}_n)}{\partial \hat{\mathbf{I}}_n} = \frac{\partial \text{argmax}_c S_c(\hat{\mathbf{I}}_n)}{\partial \hat{\mathbf{I}}_n} . \quad (9)$$

Taking the de-noised thinking in (7) to calculate the smooth gradient of input \mathbf{I}

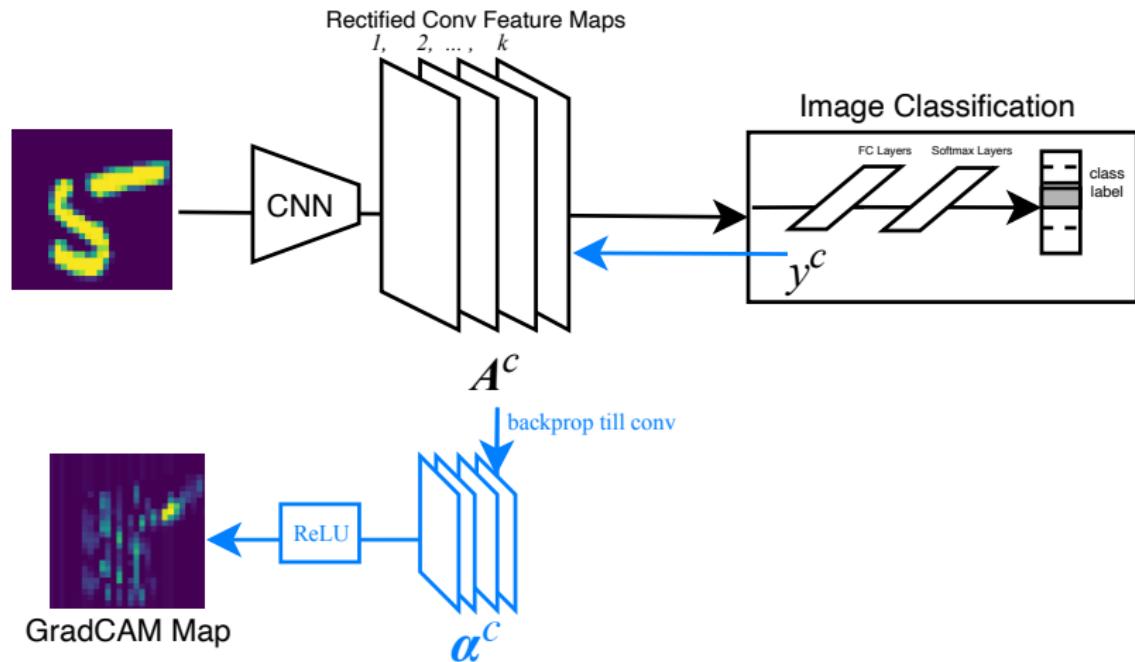
$$\tilde{\omega} = \frac{1}{n} \sum_1^n (\hat{\omega}_1 + \hat{\omega}_2 + \dots + \hat{\omega}_n) . \quad (10)$$

SmoothGrad map can be computed:

$$\begin{aligned} \tilde{\mathbf{M}}(i, j) &= |\hat{\omega}_{p(i,j)}| , \\ \tilde{\mathbf{M}}(i, j) &= \max_h |\hat{\omega}_{p(i,j,h)}| , \text{ if the input is multi-channel image} \end{aligned} \quad (11)$$

GradCAM

GradCAM



Vsiualization

GradCAM

It first has to compute the gradient of the score for class c , \mathbf{y}^c (before the softmax), with respect to feature map activation \mathbf{A}^k of a convolutional layer,

$$\frac{\partial \mathbf{y}^c}{\partial \mathbf{A}^k}$$

These gradients then are global-average-pooled over the width and height dimensions to obtain the neuron importance weights α_k^c .

$$\alpha_k^c = \underbrace{\frac{1}{Z} \sum_i \sum_j}_{\text{global average pooling}} \underbrace{\frac{\partial \mathbf{y}^c}{\partial A_{ij}^k}}_{\text{gradient via backprop}}, \quad (12)$$

where A_{ij}^k represents the element value in \mathbf{A}^k with corresponding pixel index (i, j) and Z is the sum of all pixels.

Vsiualization

GradCAM

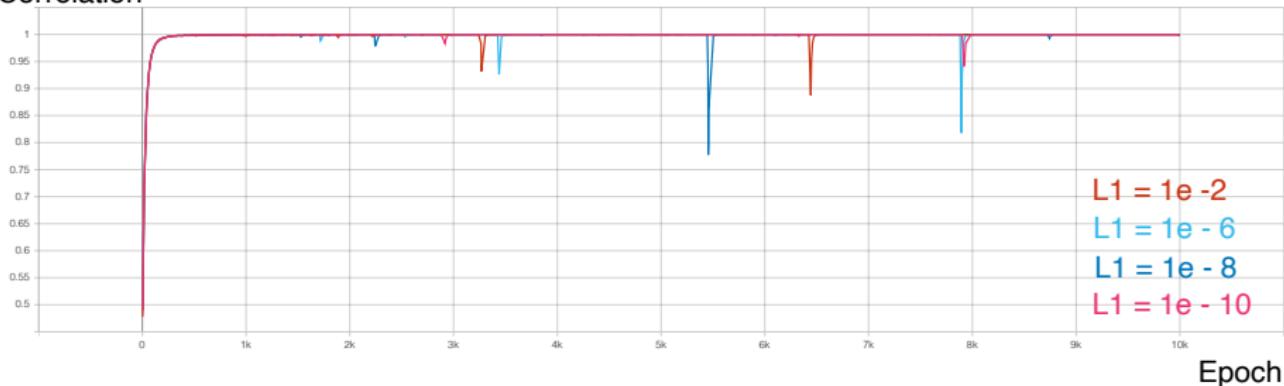
Then GradCAM map L_{GradCAM}^c can be a weighted combination of forward activation maps then be the answers after through a ReLU.

$$L_{\text{GradCAM}}^c = \text{ReLU} \underbrace{\left(\sum_k \alpha_k^c A^k \right)}_{\text{linear combination}} \quad (13)$$

L1-norm regularization

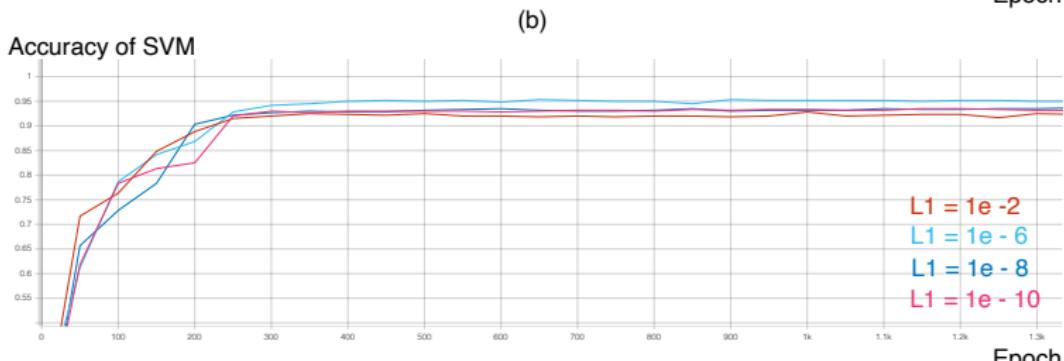
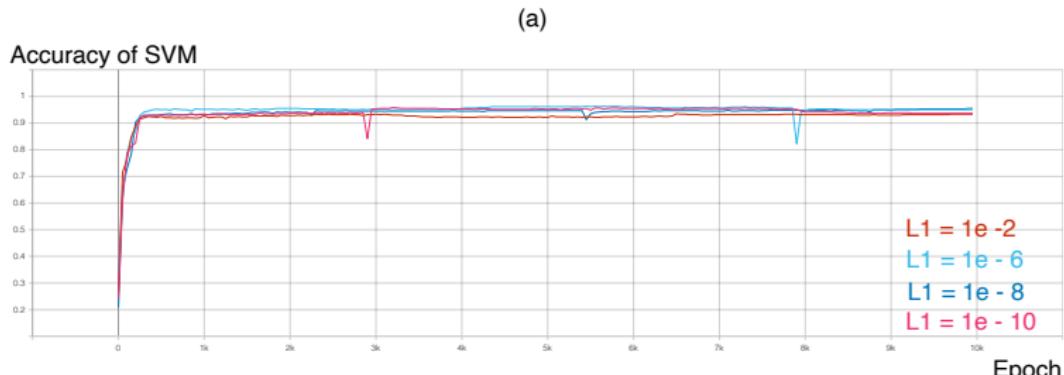
DNN model

Correlation



L1-norm regularization

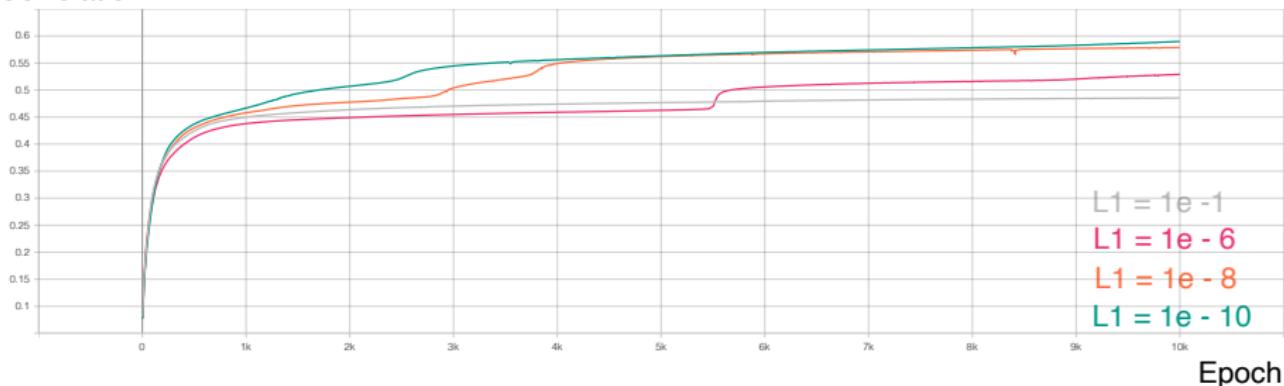
DNN model



L1-norm regularization

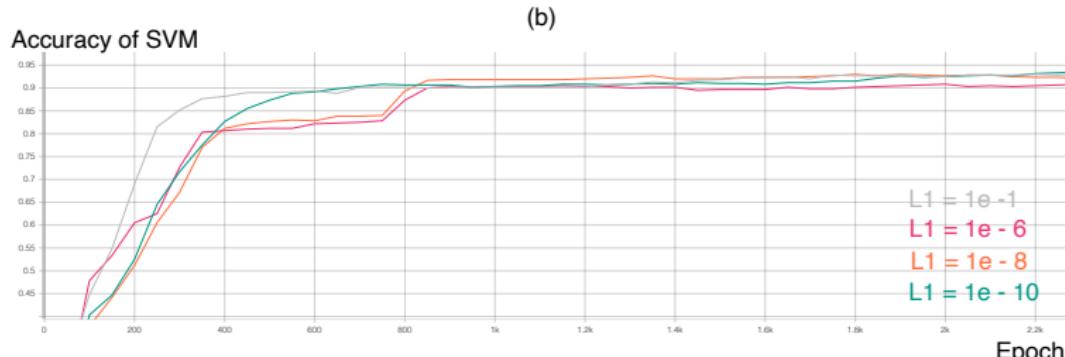
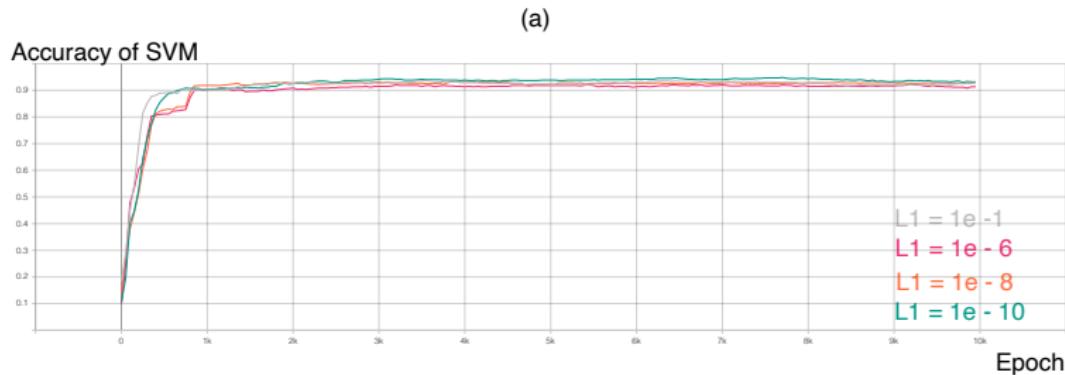
CNN model

Correlation



L1-norm regularization

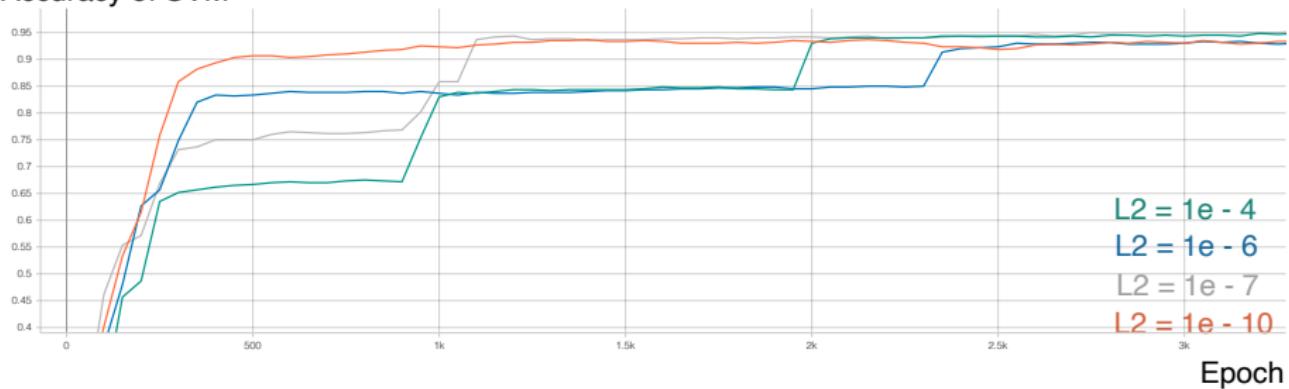
CNN model



L2-norm regularization

CNN model

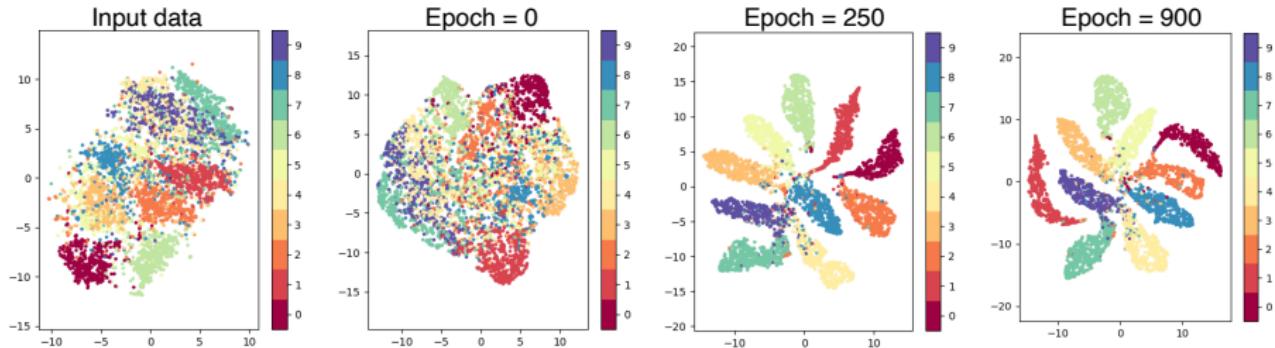
Accuracy of SVM



The comparable optimal hyper-parameters setting is $L2 = 1e - 6$

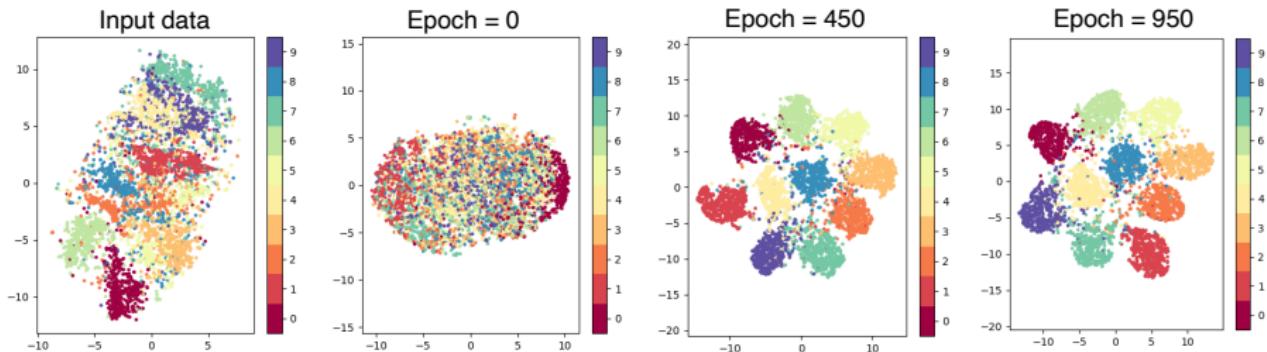
Visualization Evaluation

DNN model



Visualization Evaluation

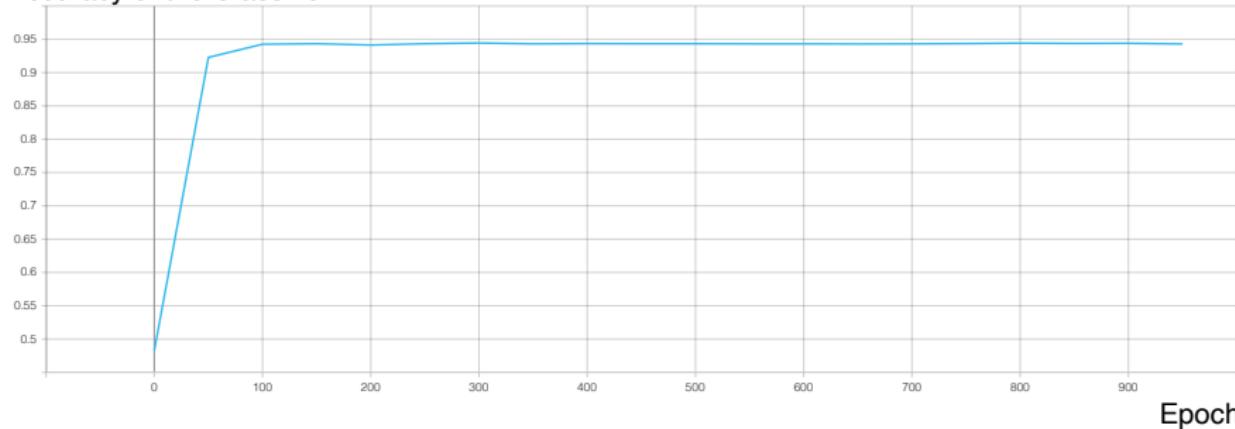
CNN model



Overfitting Evaluation

CNN model

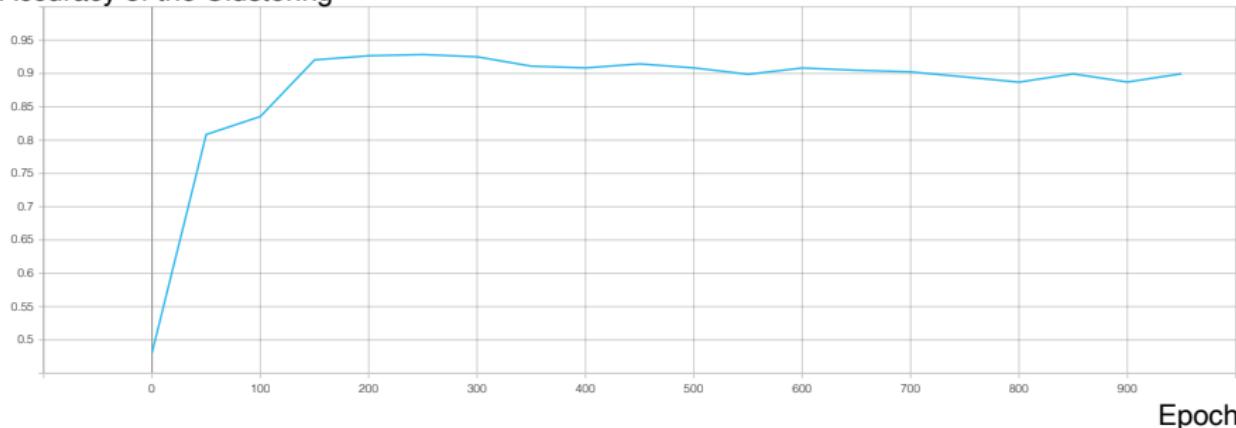
Accuracy of the Classifier



Overfitting Evaluation

CNN model

Accuracy of the Clustering



- k-means clustering¹⁰

Overfitting Evaluation

New performance metric

- box attention

box attention

An original image \mathbf{I} has the corresponding GradCAM map $\mathbf{G} \in \mathbb{R}^{m_g \times n_g}$ and the box filter is $\mathbf{B} \in \mathbb{R}^{m \times n}$. First scale the \mathbf{G} in the same size of \mathbf{I} by a linear transformation,

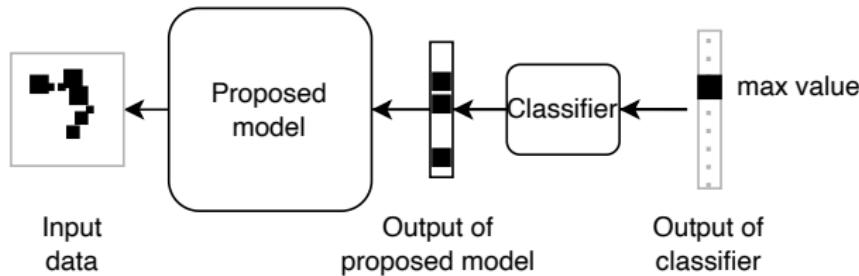
$$\mathbf{G}'(i, j) = \frac{m - n}{m_g - n_g} \mathbf{G}(i, j) ,$$

where \mathbf{G}' is the scaled GradCAM map and (i, j) represents the pixel value in the i -th row and j -th column. Then box attention b is defined as the average result of scaled GradCAM map times the box filter,

$$b = \mathbb{E}(\mathbf{G}' \cdot \mathbf{B}) . \quad (14)$$

Visualization Techniques

Saliency Map and SmoothGrad workflow in the proposed model



Setting

- Gaussian noise $\mathcal{N}(0, 0.4^2)$

Experiments and Results

Experiments

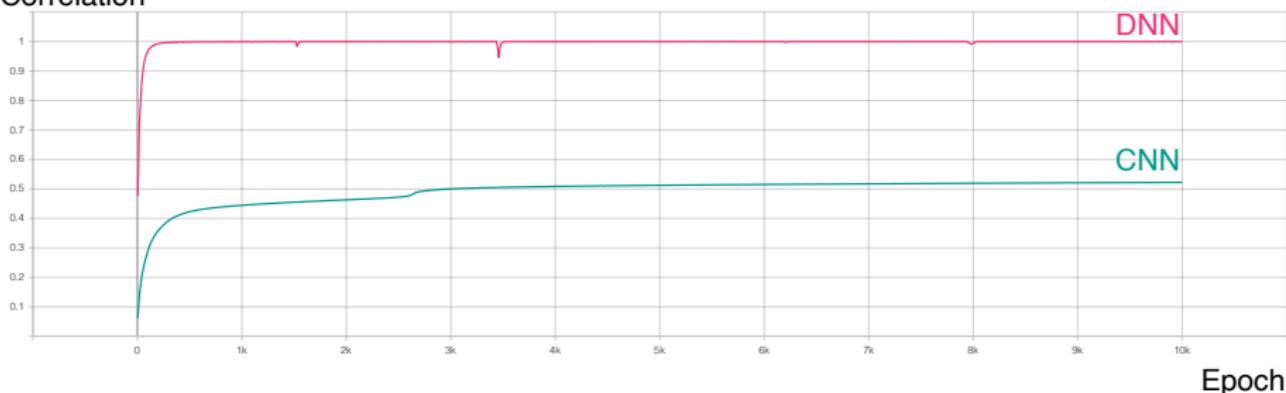
- Comparison between Two Network Structures
- Hyper-parameter settings
 - ▶ L1-norm
 - ▶ L2-norm
- Visualizations evaluation
- Overfitting analysis

Optimal setting

- The comparable optimal hyper-parameters setting is $L2 = 1e - 6$

L2-norm

Correlation



Setting

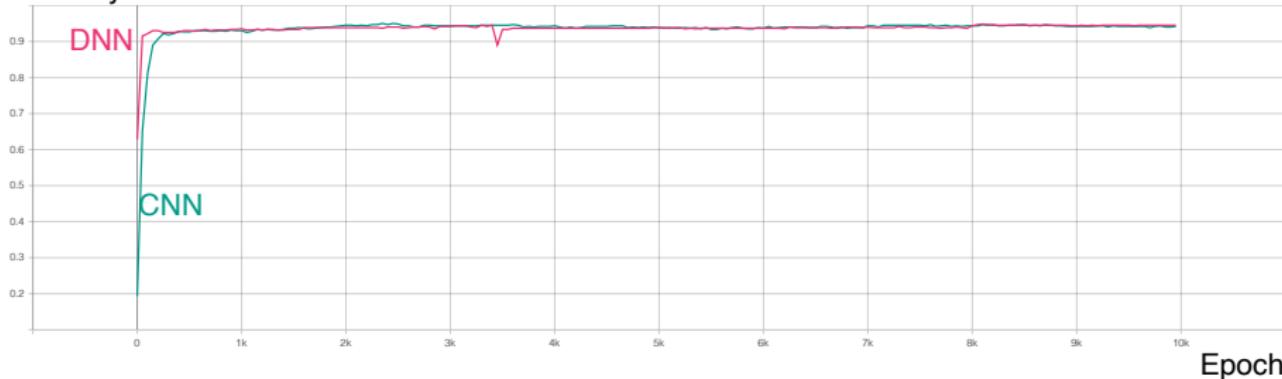
- Input: 1 box MNIST

Optimal setting

- The comparable optimal hyper-parameters setting is $L2 = 1e - 6$

L2-norm

Accuracy of SVM

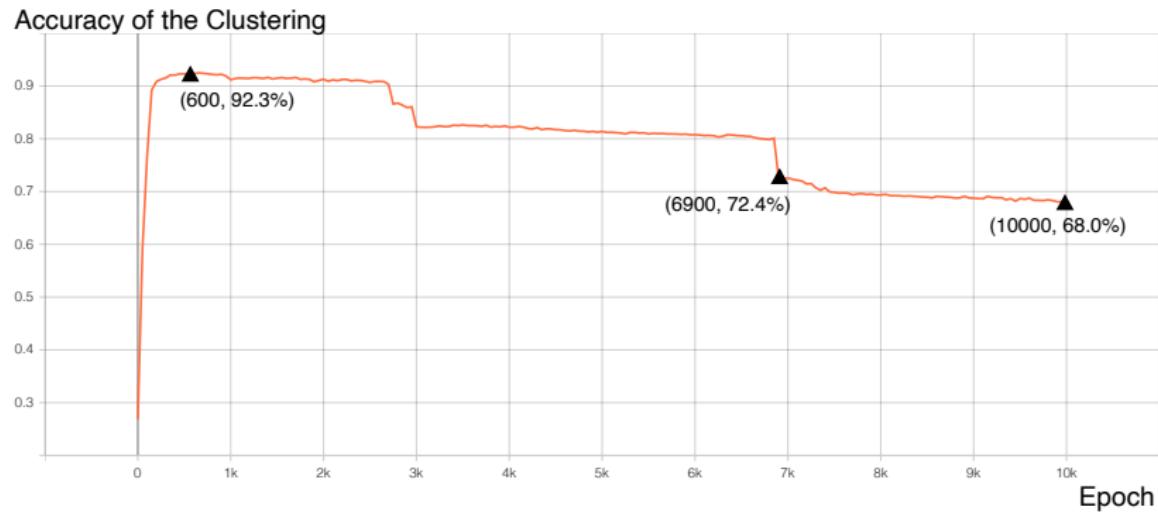


Setting

- Input: 1 box MNIST

Overfitting Evaluation

CNN model

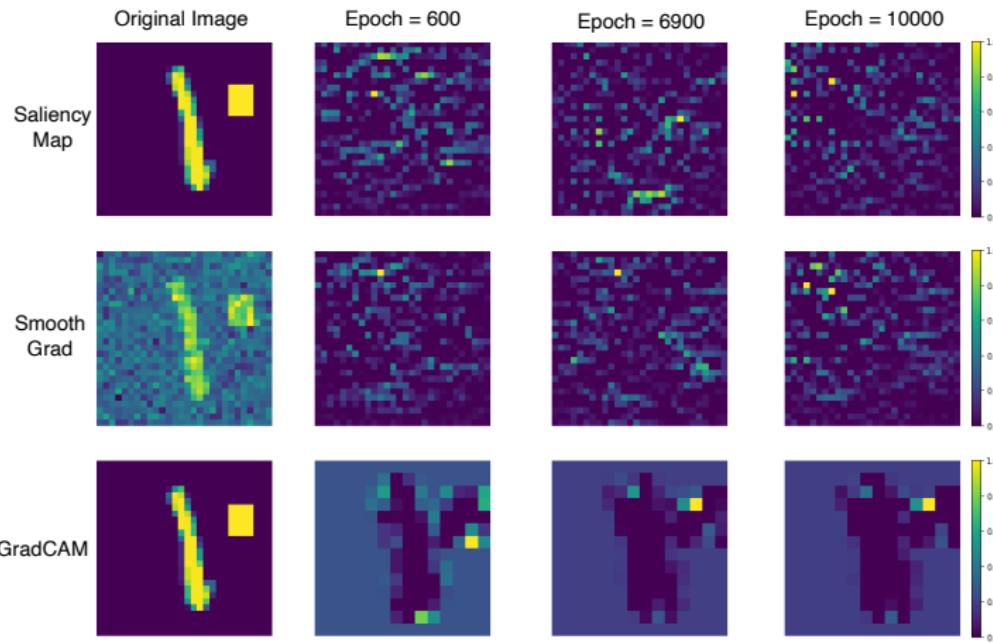


Setting

- Input: 1 box MNIST

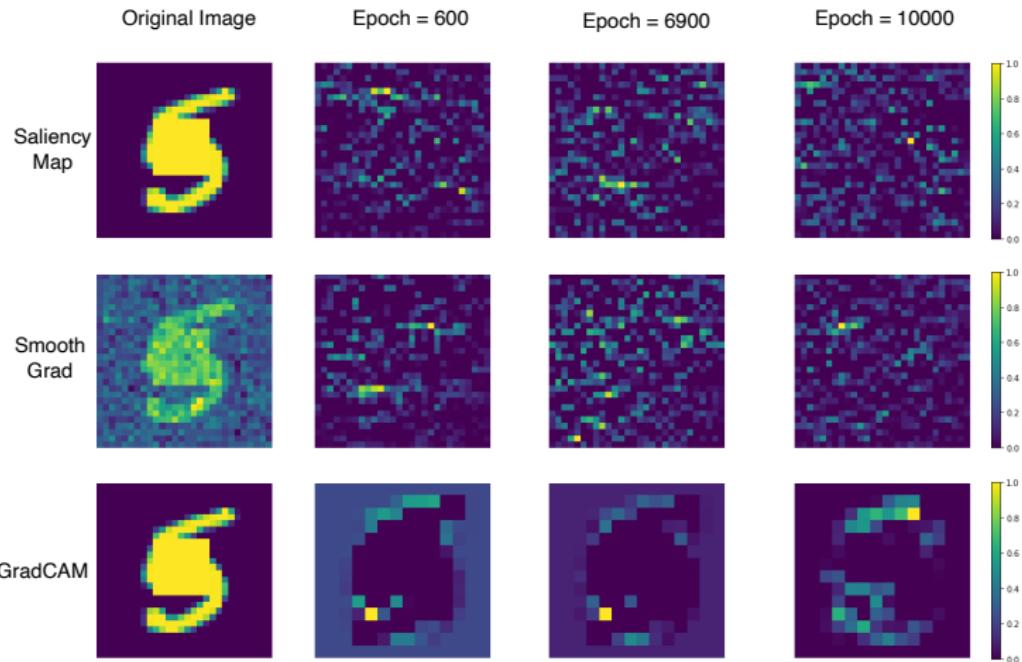
Overfitting Evaluation

CNN model



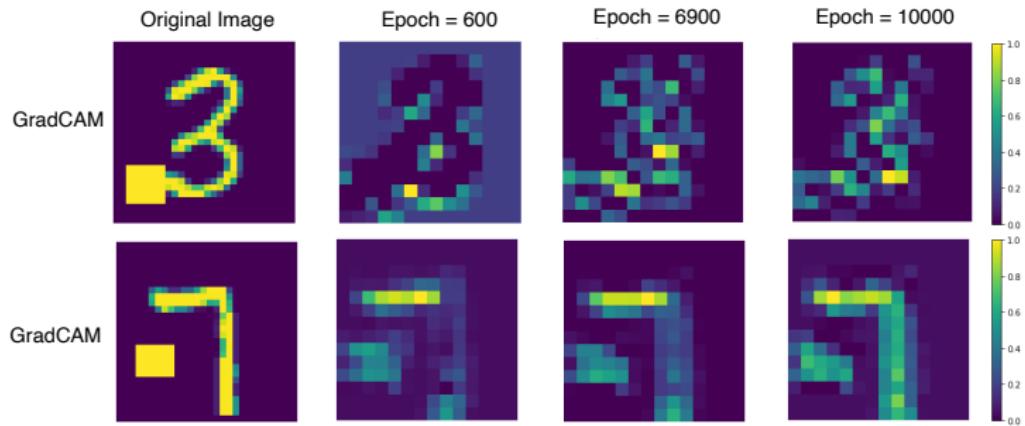
Overfitting Evaluation

CNN model



Overfitting Evaluation

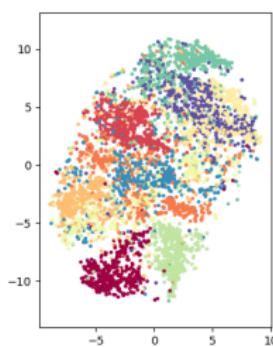
CNN model



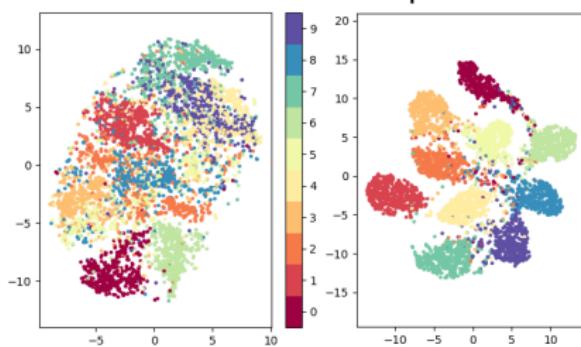
Overfitting Evaluation

CNN model

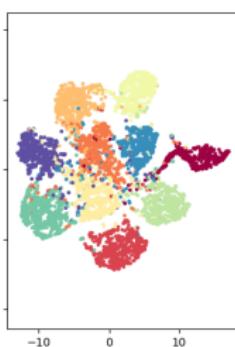
Original data



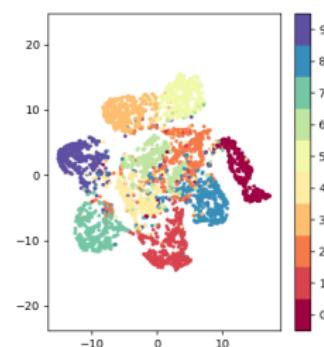
Epoch = 600



Epoch = 6900



Epoch = 10000



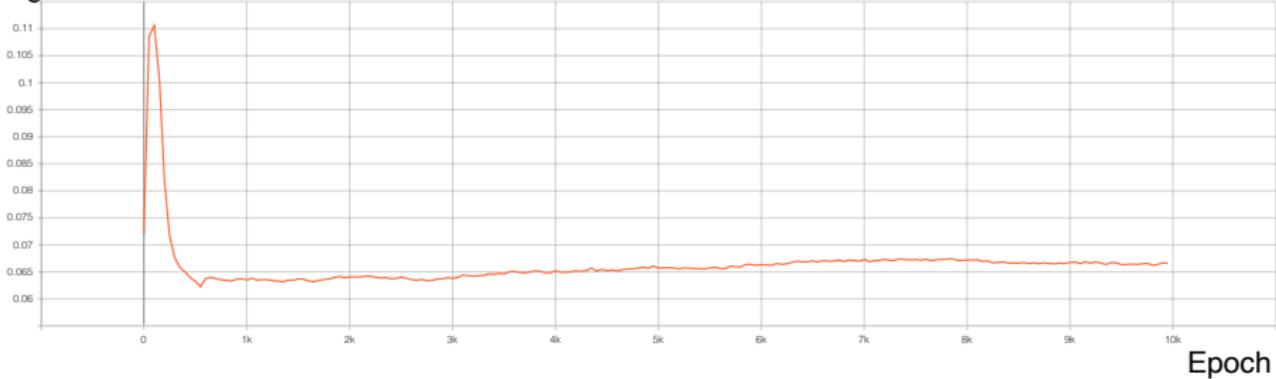
Setting

- Input: 1 box MNIST

Overfitting Evaluation

CNN model

Digit Attention



Setting

- Input: 1 box MNIST