2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS) October 25-29, 2020, Las Vegas, NV, USA (Virtual)

Serket: A Framework for Construction of Multimodal Learning Models

The University of Electro-Communications Tomoaki Nakamura

978-1-7281-6211-9/20/\$31.00 ©2020 IEEE

Background

 Our group has developed many multimodal learning models based on probabilistic models



Learning object and language [Nakamura+ 15]





Place concept formation [Taniguchi+ 17]

Formation of integrated concepts with object and ⁸²⁷ motion [Attamimi+ 14]

Learning object concepts and language model

Multimodal information is classified into categories (MLDA)

- Robot obtain multimodal (visual, auditory and haptic) information by observing, grasping and shaking objects
- User teaches object features by speech
- We assumed robot does not have predefined language knowledge
 Parameters of speech recognition (SR) are learned simultaneously



Spatial Concept Acquisition [Taniguchi+ 2017]

- Taniguchi et al. proposed more complicated model
- Robot builds a map using SLAM and simultaneously learns:
 - space name and appearance (MLDA), space region (GMM), parameters of speech recognition (LM)



Multimodal Learning=Complementary Learning

This model is constructed by connecting two models:

- Speech recognition and clustering (latent Dirichlet allocation)
- Shared variable w^w is determined with mutual influence
 - Possibility that speech o is recognized as w^w
 - Possibility that w^w is co-occurred with category z

complementary learning



Background

These models are also constructed by connecting smallscale models



Problem in Constructing Models

- These models have complicated structure
- To realize human-like learning models, much more complicated models are required



832

This part becomes more difficult
 Framework to easily construct models is required

SERKET

To easily construct multimodal learning models, we have proposed the framework SERKET



833

Shared latent variables are optimized with mutual influence on each model

Modules in SERKET

- Each module has observations and latent variables
- Models are constructed by connecting shared latent variables of modules hierarchically



 Parameters are estimated by communication between modules while programmatic independence maintains Estimation of shared latent variables

Parameters are estimated by exchanging messages



- I. Receive messages from other modules and latent variable $z_{m,n}$ and parameters are updated
- 2. New messages are sent to other modules based on updated parameters

- By executing this procedure sequentially in each module, the parameters are optimized mutually
- Currently, two methods are implemented for message exchange

Parameter optimization 1

Message Passing (MP) Approach



- Assume that z is determined with mutual influence : z ~ P(z|Θ₁, Θ₂, o)
 ∝ P(z|Θ₁)P(z|Θ₂, o)
- Module $I \rightarrow$ Module 2:
 - Module I sends $P(z|\Theta_1)$ to Module 2
 - Module 2 updates Θ_2 using $P(z|\Theta_1)$
- Module $2 \rightarrow$ Module I:
 - Module 2 sends $P(z|o, \Theta_2)$ to Module 1
 - Module I updates Θ_1 using $P(z|o, \Theta_2)$
- Θ_1 : Parameters of Module I
- Θ_2 : Parameters of Module 2
- z : Shared latent variables
- o: Observation

Parameter optimization 2

Sampling Importance Resampling (SIR) approach



- In the case that w has a large number of possibilities such as speech recognition results
 Monte Carlo approximation
- Module 2 generates samples w_n and sends them to Module I

$$w_n \sim P(w|\Theta_2, o)$$

- Module I resamples based on P(w|Θ₁) and sends selected samples w^{*} to Module 2
- Parameters are updated using selected w^*

Implementation Examples

- We confirm that integrated models by SERKET improve their performance
- Used modules:
 - Variational Auto-Encoder (VAE)
 - Gaussian Mixture Model (GMM)
 - Multimodal Latent Dirichlet Allocation (MLDA)
 - Markov Model (MM)
 - Speech Recognition (SR)
- Multimodal Dataset:
 - Image : MNIST
 - Speech : Spoken Arabic Digit Dataset

Implementation Examples

Five examples

- I. VAE+GMM
- 2. VAE+GMM+MLDA
- 3. VAE+GMM+MLDA+MM

We share Jupyter notebook that you can execute on the Google Collaboratory

```
http://iros20.naka-lab.org/
```

- 4. Model for language acquisition by robots
- 5. Model for learning object feature extractor by robots

Example1 : VAE+GMM

Unsupervised image classification

D

Image data: MNIST (784 dims)



- Dimensional compression by VAE (18 dims)
- Classification by GMM



Example1 : VAE+GMM



```
import serket as srk
 1
   import vae
 \mathbf{2}
   import gmm
 3
   import numpy as np
 4
5
   # Load observations
6
   obs1 = srk.Observation(np.loadtxt("data.txt"))
 \overline{7}
   category = np.loadtxt("category.txt")
8
9
   # Define modules
10
   vae1 = vae.VAE(18, itr=200, batch_size=500)
11
   gmm1 = gmm.GMM(10, category=category)
12
13
   # Define connection between modules
14
   vae1.connect(obs1)
15
   gmm1.connect(vae1)
16
17
   # Optimize the parameters
18
   for i in range(5):
19
       vae1.update()
20
       gmm1.update()
21
```

Example1: Mutual learning of VAE and GMM Modified evidence lower bound (ELBO) of VAE $\mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\phi}; \boldsymbol{o}) = -D_{KL}(q_{\boldsymbol{\phi}}(\boldsymbol{z}_1 | \boldsymbol{o}) \| \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{I})) + \mathbb{E}_{q_{\boldsymbol{\phi}}(\boldsymbol{z}_1 | \boldsymbol{o})}[\log p_{\boldsymbol{\theta}}(\boldsymbol{o} | \boldsymbol{z}_1)]$ Standard VAE: $\mu = 0$ GMM Latent space suitable z_2 for classification can be learned Compute μ of cluster with mutual influence z_1 \boldsymbol{z}_1 μ z_1 784 dimensional observations z_1 is estimated using modified ELBO are compressed into and received μ 18 dimensional z_1 VAE

Example1: Classification Results

Adjusted Rand Index (ARI)

Mutual learning	average	best
-	0.477	0.478
\checkmark	0.503	0.568
average, Average of 10 trials		

* best: The best result in 10 trials





Confusion Matrix

- Latent space suitable for classification was learned with mutual influence
- Classification accuracy increased

Example2: VAE+GMM+MLDA

D

Unsupervised classification of image and speech

Pairwise dataset of image and speech



Example 2: VAE+GMM+MLDA

Source code

Example 2: Classification Results

Adjusted Rand Index (ARI)

Mutual learning	average	best
-	0.604	0.638
\checkmark	0.637	0.735

Classification accuracy improved by
multimodal learning
using images and speech

* average: Average of 10 trials best:The best result in 10 trials

Confusion Matrix

Example 3: VAE+GMM+MLDA+MM

Category and their sequence rules are learned

Example 3: VAE+GMM+MLDA+MM


```
import vae
  import gmm
3
  import mlda
  import mm
5
  import numpy as np
6
   # Load observations
8
  obs1 = srk.Observation(np.loadtxt("data1.txt")) # 画像
9
  obs2 = srk.Observation(np.loadtxt("data2.txt")) # 音声
10
  category = np.loadtxt("category.txt")
11
12
   # Define modules
13
  vae1 = vae.VAE(18, itr=200, batch size=500)
14
  gmm1 = gmm.GMM(10, category=category)
15
  mlda1 = mlda.MLDA(10, category=category)
16
  mm1 = mm.MarkovModel()
17
18
   # Define connection between modules
19
  vae1.connect(obs1)
20
   gmm1.connect(vae1)
21
  mlda1.connect(obs2, gmm1)
22
  mm1.connect(mlda1)
23
\mathbf{24}
   # Optimize the parameters
25
  for i in range(5):
26
      vae1.update()
27
      gmm1.update()
28
      mlda1.update()
29
      mm1.update()
30
```

import serket as srk

Suorce code

Example 3: Classification Results

Adjusted Rand Index (ARI)

Mutual learning	average	best
-	0.575	0.524
\checkmark	0.834	0.980

Classification accuracy significantly improved by learning transition rules

* average: Average of 10 trials best:The best result in 10 trials

Confusion Matrix

Conclusion of Example 1,2 and 3

Adjusted Rand Index (ARI)

	VAE+GMM	VAE+GMM+MLDA	VAE+GMM+MLDA+MM
average	0.503	0.637	0.834

Confusion Matrix

D

Complementary learning is realized by using SERKET

Example 4: Model for language acquisition

Unsupervised classification of multimodal dataset

- Robot obtain multimodal (visual, auditory and haptic) information by observing, grasping and shaking objects
- User teaches object features by speech
- We assumed robot does not have language knowledge

Parameters of speech recognition (SR) are learned simultaneously

Example 4: Model for language acquisition

Learn not only categories but also language model (LM)

Language model and object categories are learned mutually

Example 4: Learning of LM

SIR is used for learning language model

- Learn parameters of MLDA - Compute $P(w^{(l)}|w^v, w^a, w^t)$ $(l=1,\cdots,L)$

$$P(w^{(l)}|w^v, w^a, w^t)$$

- w is determined by resampling based on $P(w^{(l)}|w^v, w^a, w^t)$
- Update parameters of LM

Example 4: Classification Results

Classify 50 multimodal data

- Classification accuracy
 - w/o mutual learning: 80%
 - w/ mutual learning: 94%
- Speech recognition accuracy
 - w/o mutual learning: 64%
 - w/ mutual learning: 74%

Classification and speech recognition accuracies improve by multimodal learning based on SERKET

Example 5: Model for learning object feature extractor by robots

- We used a dataset obtained through the robot observed objects and the human taught object feature by speech
 - # of the objects: 499
 - # of the categories: 81
- Words and images were used in this example
 - Words
 - The speech was recognized by phoneme recognizer
 - Recognized strings were segmented into words by unsupervised word segmentation (NPYLM)
 - Images

Example 5: MLDA+VAE

The model to learn object categories and image features from words and images

Example 5: Classification Results

- Compare with the model where pre-trained CNN is used for feature extractor
- Classification accuracies:
 - Integrated model of VAE and MLDA: 67.8%
 - Model with pre-trained CNN: 66.8%
- Learned latent space:

Suitable latent space for classification were learned

Example 5: Cross modal inference

Images were generated from words input

"Cup noodle"

"Plastic bottle"

"Snack" "Sponge"
 Images that represent the characteristics of the categories were generated

Conclusion

- We implemented VAE, GMM, LDA and MM modules, and integrated model with them
- We showed implementation examples and it is easy to construct the integrated models
- Accuracy improved by mutual influence of modules
- Moreover, we implemented speech recognition module and its integrated model

- Future work
 - Integrate deep neural network using Pixyz
 - Improve the efficiency of parameter inference