
Homework 14

Regularization-based Lifelong Learning

ML TAs

ntu-ml-2022spring-ta@googlegroups.com

Machine Learning Exit Evaluation

Friendly reminder:

This questionnaire is only for educational analysis. Will never affect your grades, please fill in with peace of mind.

<https://forms.gle/jAQPSoFQaoXcA22U6>



Questionnaire of Education Course Series for AI Technologies and Applications

Friendly reminder:

This questionnaire is only for educational analysis. Will not be sent to school or teacher and will never affect your grades, all information is absolutely confidential, please fill in with peace of mind.

<https://forms.gle/boFFRPwvyMeaExxB6>

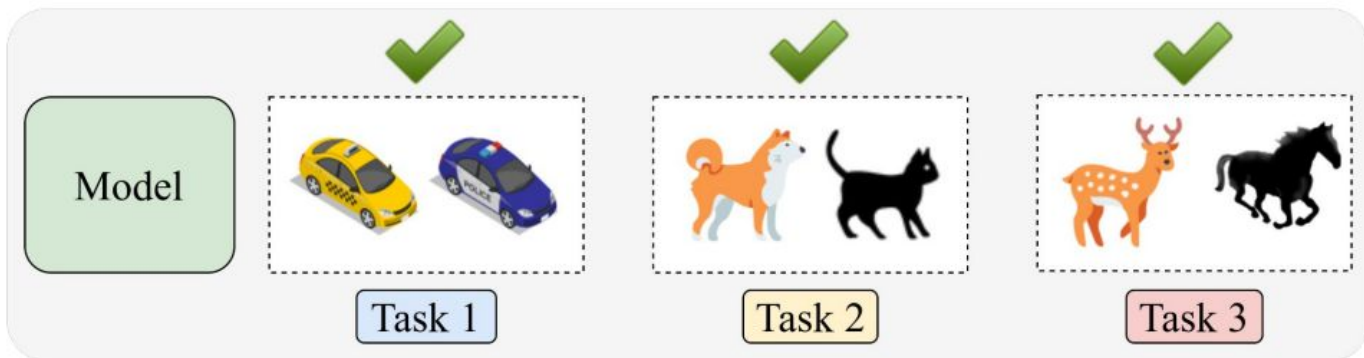


Outline

- Introduction
- Dataset
- Sample Code
- Grading
- Submission

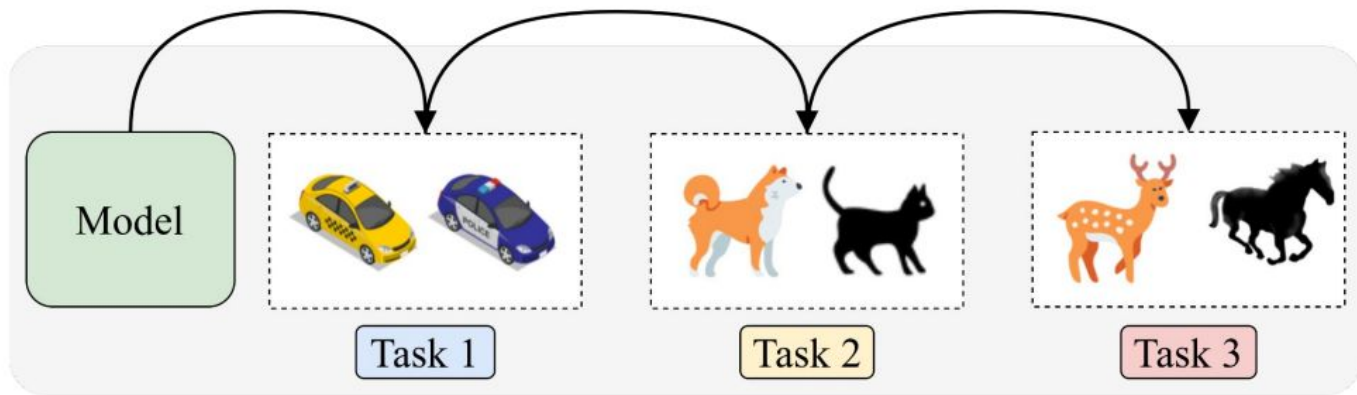
Introduction - LifeLong Learning

Goal: A model can beat all task!



Introduction - LifeLong Learning

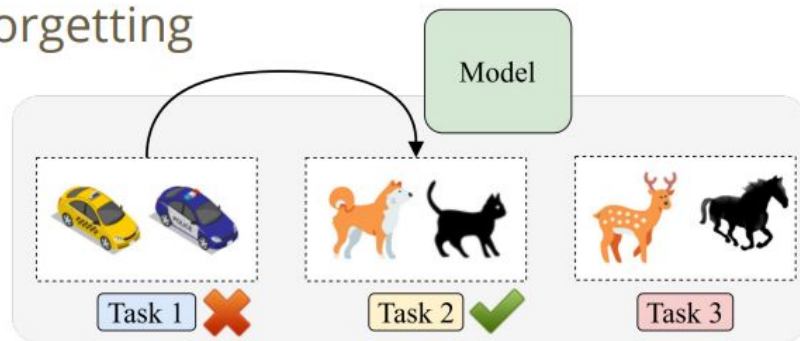
Condition: Model Sequentially Learn Different Task! (In Training Time)



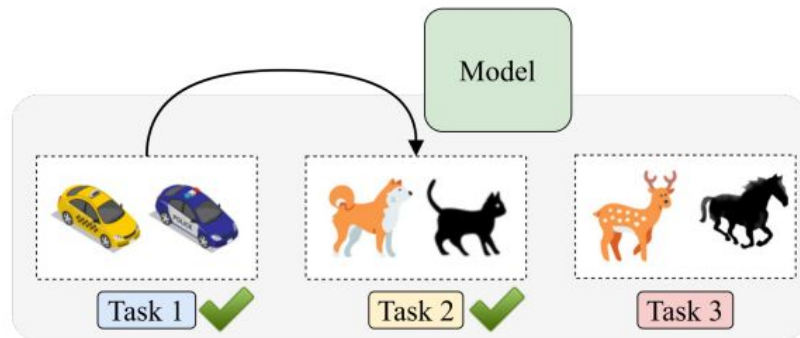
Introduction - LifeLong Learning

KeyPoint: Avoid Catastrophic Forgetting

Catastrophic Forgetting



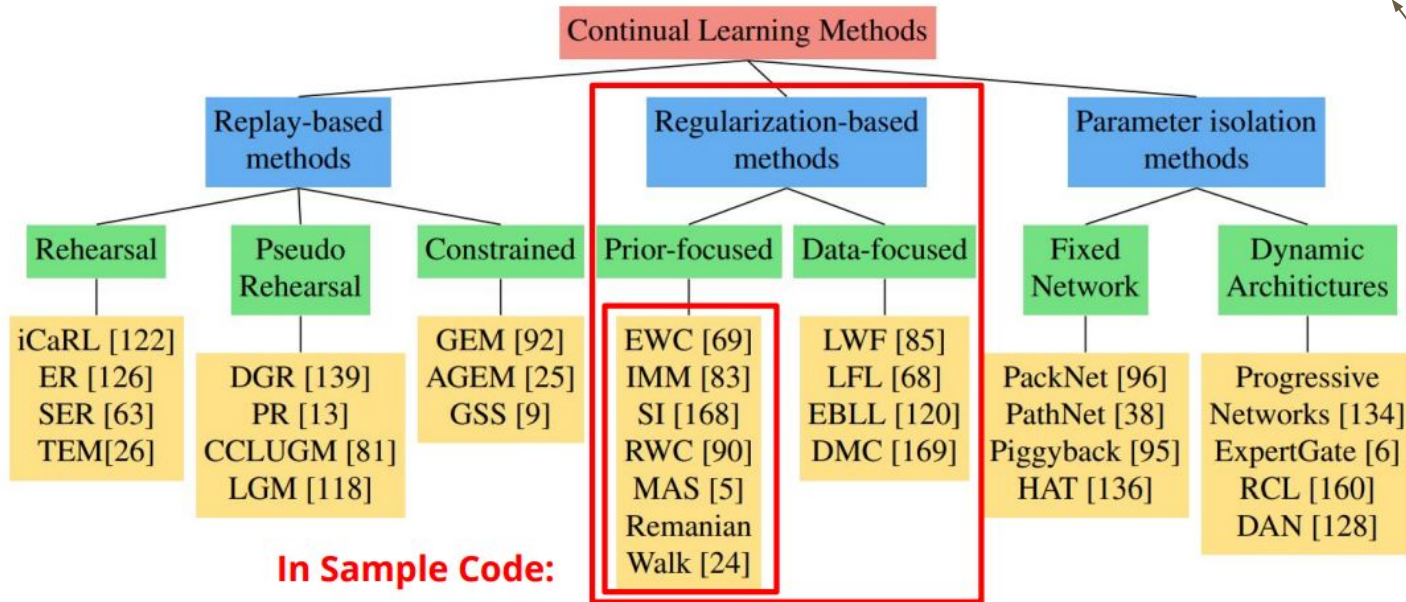
Avoid Catastrophic Forgetting



Introduction

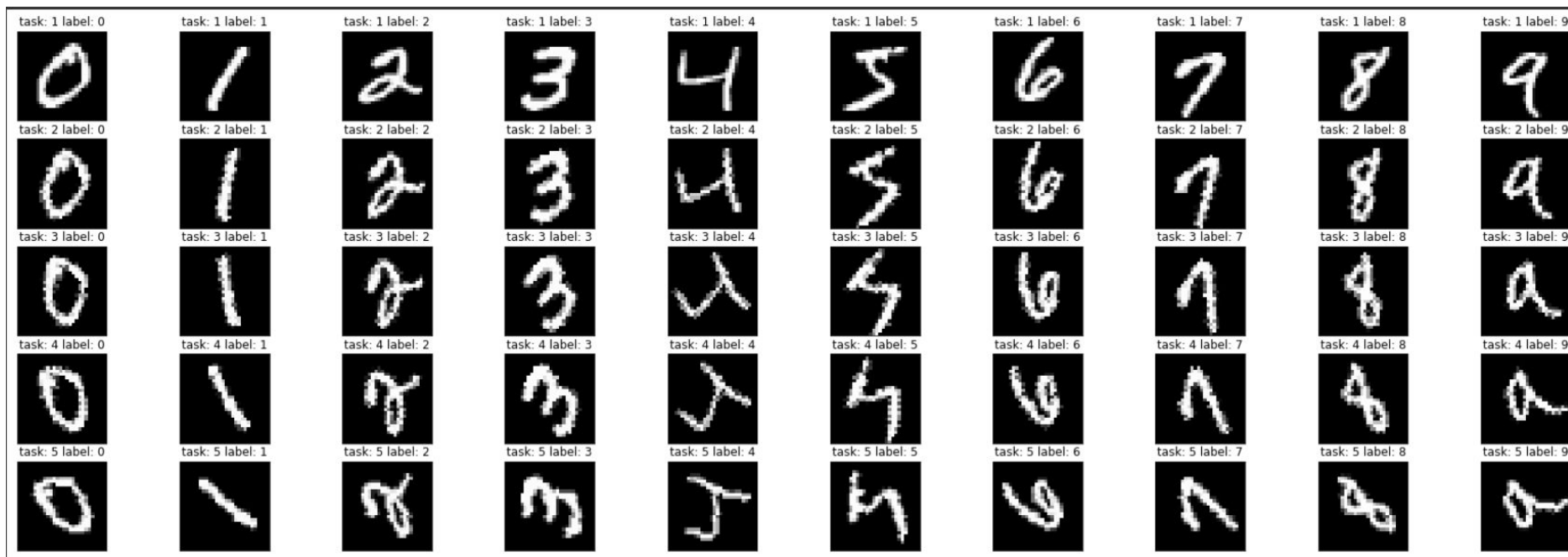
$$L'(\theta) = L(\theta) + \lambda \sum_i b_i (\theta_i - \theta_i^b)^2$$

Regularization term



Dataset

Rotated MNIST (Generated by TAs)



Sample Code - Training Details

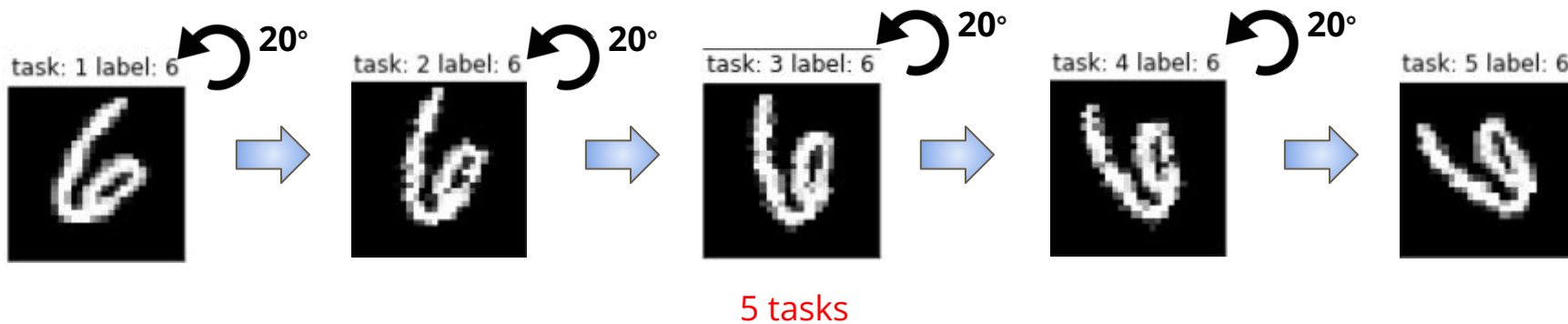
- 5 task / Each task has 10 epoches for training.
- Each method cost ~20 minutes for training model. (Tesla T4)
- Each method cost ~60 minutes for training model. (Tesla K80)
- [Colab link](#)

Sample Code - Guideline

- Utilities
- Prepare Data
- Prepare Model
- Train and Evaluate
- Methods
- Plot function

Sample Code - Prepare Data

- Prepare Data
 - Rotation and Transformation
 - Dataloaders and Arguments
 - Visualization



Sample Code - Prepare Model

- Prepare Model
 - Model Architecture

Fixed model size!

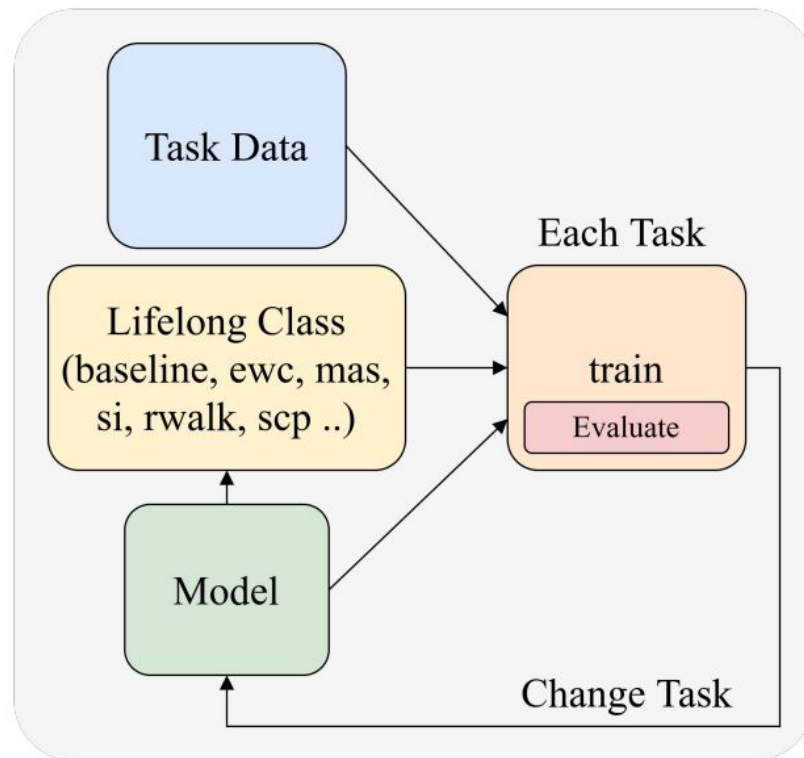
```
Model(  
  (fc1): Linear(in_features=784, out_features=1024, bias=True)  
  (fc2): Linear(in_features=1024, out_features=512, bias=True)  
  (fc3): Linear(in_features=512, out_features=256, bias=True)  
  (fc4): Linear(in_features=256, out_features=10, bias=True)  
  (relu): ReLU()  
)
```

Sample Code - Train and Evaluate

- Train:
 - Sequentially train.
 - Add regularization term and update it.
- Evaluate:
 - Evaluate by using a special metric.
 - (Please read sample code and describe it in your report.)

Sample Code - Before Going to Methods

Training Pipeline:



Sample Code - Before Going to Methods

Lifelong Class
(baseline, ewc, mas,
si, rwalk, scp ..)

```
6 class baseline(object):
7     """
8     baseline technique: do nothing in regularization term [initialize and all weight is zero]
9     """
10    def __init__(self, model, dataloaders, device):
11
12        self.model = model
13        self.dataloaders = dataloaders
14        self.device = device
15
16        self.params = {n: p for n, p in self.model.named_parameters() if p.requires_grad} #extract all parameters in models
17        self.p_old = {} # store current parameters
18        self._precision_matrices = self._calculate_importance() # generate weight matrix
19
20        for n, p in self.params.items():
21            self.p_old[n] = p.clone().detach() # keep the old parameter in self.p_old
22
23    def _calculate_importance(self):
24        precision_matrices = {}
25        for n, p in self.params.items(): # initialize weight matrix (fill zero)
26            precision_matrices[n] = p.clone().detach().fill_(0)
27
28        return precision_matrices
29
30    def penalty(self, model: nn.Module):
31        loss = 0
32        for n, p in model.named_parameters():
33            _loss = self._precision_matrices[n] * (p - self.p_old[n]) ** 2
34            loss += _loss.sum()
35        return loss
36
37    def update(self, model):
38        # do nothing
39        return
```


Sample Code - Before Going to Methods

Lifelong Class
(baseline, ewc, mas,
si, rwalk, scp ..)

```
6 class baseline(object):
7     """
8     baseline technique: do nothing in regularization term [initialize and all weight is zero]
9     """
10    def __init__(self, model, dataloaders, device):
11
12        self.model = model
13        self.dataloaders = dataloaders
14        self.device = device
15
16        self.params = {n: p for n, p in self.model.named_parameters() if p.requires_grad} #extract all parameters in models
17        self.p_old = {} # store current parameters
18        self._precision_matrices = self._calculate_importance() # generate weight matrix
19
20        for n, p in self.params.items():
21            self.p_old[n] = p.clone().detach() # keep the old parameter in self.p_old
22
23    def _calculate_importance(self):
24        precision_matrices = {}
25        for n, p in self.params.items(): # initialize weight matrix (fill zero)
26            precision_matrices[n] = p.clone().detach().fill_(0)
27
28        return precision_matrices
29
30    def penalty(self, model: nn.Module):
31        loss = 0
32        for n, p in model.named_parameters():
33            _loss = self._precision_matrices[n] * (p - self.p_old[n]) ** 2
34            loss += _loss.sum()
35        return loss
36
37    def update(self, model):
38        # do nothing
39        return
```

train

Sample Code - Methods

- Baseline (Do nothing in regularization term)

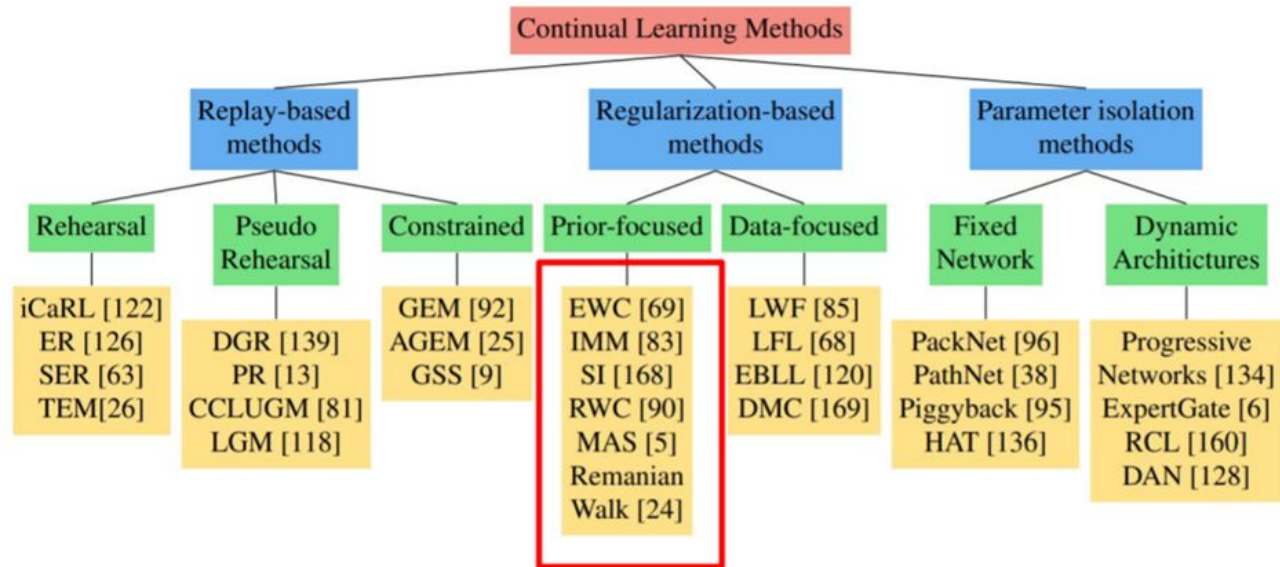
- [EWC](#)

- [MAS](#)

- [SI](#)

- [RWalk](#)

- [SCP](#)



EWC - Elastic Weight Consolidation

1. You need to know how to generate Guardiance weight from EWC!
2. Do this method need to use label ?
3. Hint: (Trace the class ewc and its calculate_importance function)

Paper Link: <https://arxiv.org/abs/1612.00796>

MAS - Memory Aware Synapse

1. You need to know how to generate Guardiance weight from MAS!
2. Do this method need to use label ?
3. We want you to implement Omega Matrix for MAS! Please read page 21 carefully and paste your code (only TODO block) in report.

Please do not modify any part of the sample code except the TODO block.

Paper Link: <https://arxiv.org/abs/1711.09601>

MAS - Memory Aware Synapse

- The method proposed in the paper is the local version by taking squared L2-norm outputs from each layer of the model.
- Here we only want you to implement the global version by taking outputs from the last layer of the model.
- Hint: (It is similar to the way you generate the Fisher matrix for EWC, the only difference is the calculation of the important weight.)

$$\mathcal{L}_B = \mathcal{L}(\theta) + \sum_i \frac{\lambda}{2} \Omega_i (\theta_i - \theta_{A,i}^*)^2$$

Paper Link: <https://arxiv.org/abs/1711.09601>

$$\Omega_i = \left\| \frac{\partial \ell_2^2(M(x_k; \theta))}{\partial \theta_i} \right\|$$

SI - Synaptic Intelligence

1. You need to know how to generate Guardiance weight from SI!
2. Do this method need to use label ?
3. Hint: (Accumulated loss change in each update step)

Paper Link: <https://arxiv.org/abs/1703.04200>, [Talk Slide](#)

SI - Main Idea

$$L(\theta) = L_2(\theta) + c \sum_i \Omega_i (\theta_i - \theta_{1,i}^*)^2$$

From learning trajectory

Parameter importance on-line from learning trajectory!

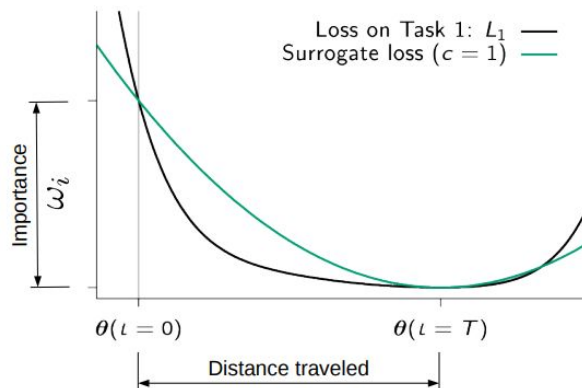
Picture comes from: [Talk Slide](#)

SI - Main Idea

Leveraging per-parameter importance for continual learning

$$L(\theta) = L_2(\theta) + c \sum_i \Omega_i (\theta_i - \theta_{1,i}^*)^2$$

$$\Omega_i \equiv \frac{\omega_i}{(\Delta_i)^2 + \epsilon}$$



Picture comes from: [Talk Slide](#)

SI - Main Idea

Total change in loss is given by the path integral over the gradient field

$$\int_C \mathbf{g}(\boldsymbol{\theta}(t)) d\boldsymbol{\theta} = \int_{t_0}^{t_1} \mathbf{g}(\boldsymbol{\theta}(t)) \cdot \boldsymbol{\theta}'(t) dt = L(t_1) - L(t_0)$$

$$= \sum_k \underbrace{\int_{t_0}^{t_1} g_k(t) \theta'_k(t) dt}_{\omega_k} \equiv - \sum_k \omega_k$$

- Is a parameter-specific quantity
- Can be computed on-line during training (running sum)

\mathbf{g} : Gradient

$\boldsymbol{\theta}$: Parameters

$\boldsymbol{\theta}'$: Updates

Natural way of assigning credit for a global change to local parameters

$$L(t_1) - L(t_0) = - \sum_k \omega_k^\mu$$

Picture comes from: [Talk Slide](#)

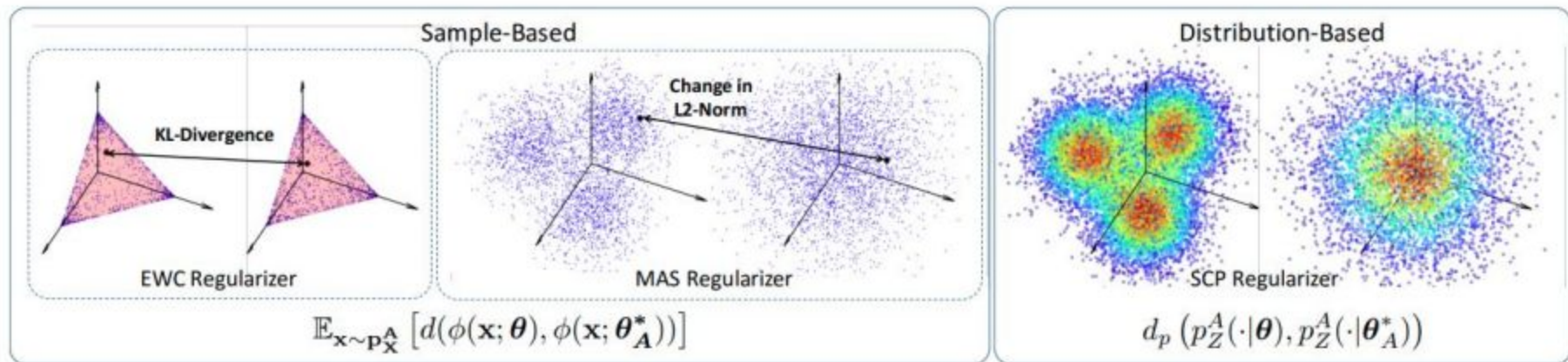
RWalk - Remanian Walk

1. Trace Rwalk class and its update function!
2. Do this method need to use label ?
3. Hint: (The code is similar to two method which mentioned in sample code)

Paper Link: <https://arxiv.org/abs/1801.10112>

SCP - Sliced Cramer Preservation

1. Paper Link: <https://openreview.net/pdf?id=Blge3TNKwH>
2. Do this method need to use label ?



SCP - Main Idea

- Propose Distributed-based Distance to **prevent fast intransigence** and **avoid overestimating the importance of parameters**.

intransigence



Model do not want to learn new task,
and it just keep old task performance

Other Methods and Scenarios

- Only in Multiple Choice Questions
- iCaRL (<https://arxiv.org/abs/1611.07725>)
- LwF (<https://arxiv.org/abs/1606.09282>)
- GEM (<https://arxiv.org/abs/1706.08840>)
- DGR (<https://arxiv.org/abs/1705.08690>)
- Three scenarios for continual learning (<https://arxiv.org/abs/1904.07734>)

Grading

- 20 multiple choice questions (**8pts, 0.4pt each**)
- Report (**2pts**)
- You have to choose ALL the correct answers for each question
- No leaderboards are needed!!

Grading - Multiple Choice Questions

- 20 multiple choice questions (**8pts, 0.4pt each**)
 - Basic Concept: 3 Questions
 - EWC: 2 Questions
 - MAS: 2 Questions
 - SI: 2 Questions
 - RWalk: 2 Questions
 - SCP: 3 Questions
 - Other Methods & scenarios: 6 Questions
 - ICaRL, LwF, GEM, DGR
 - Three Scenarios

Grading - Report

- Plot the learning curve of the metric with every method. (The Plotting function is provided in the sample code.) **(0.5pt)**
- Describe the metric. **(0.5pt)**
- Paste the code that you implement Omega Matrix for MAS. **(1pt)**

Please do not modify any part of the sample code except the TODO block.

Please just paste the TODO block.

If you plot the right learning curve, you will still get the point of the first part no matter whether you implement Omega Matrix for MAS or not.

Submission

- The questions are on gradescope
- Submit your report to gradescope
- Running the code may need some time!
- You can answer the questions unlimited times
- The length of answering time of the assignment is unlimited
- We will consider the latest submission as the final score
- You will see the scores after the deadline only!
- **No late submission!**
- **Remember to save the answer when answering the questions!**
- Deadline: **2022/06/24 23:59**

Link

- Code: [Colab](#)

If any questions, you can ask us via...

- NTU COOL (Recommended)
- Email
 - mlta-2022-spring@googlegroups.com
 - The title should begin with “[hw14]”
- TA hour
 - Tuesday, 20:00 ~ 21:00
 - Friday, 15:00 ~ 17:00