

지능형생물정보학

한양대학교 컴퓨터 소프트웨어학부
DATA MINING & BIOINFORMATICS LAB.

문희상

Index

1. Environment setup
2. Preparing input data
3. Construct & Training a model
4. Visualize the result with TensorBoard
5. Evaluate the result
6. Let's start the practice
7. References

Overview: environment setup

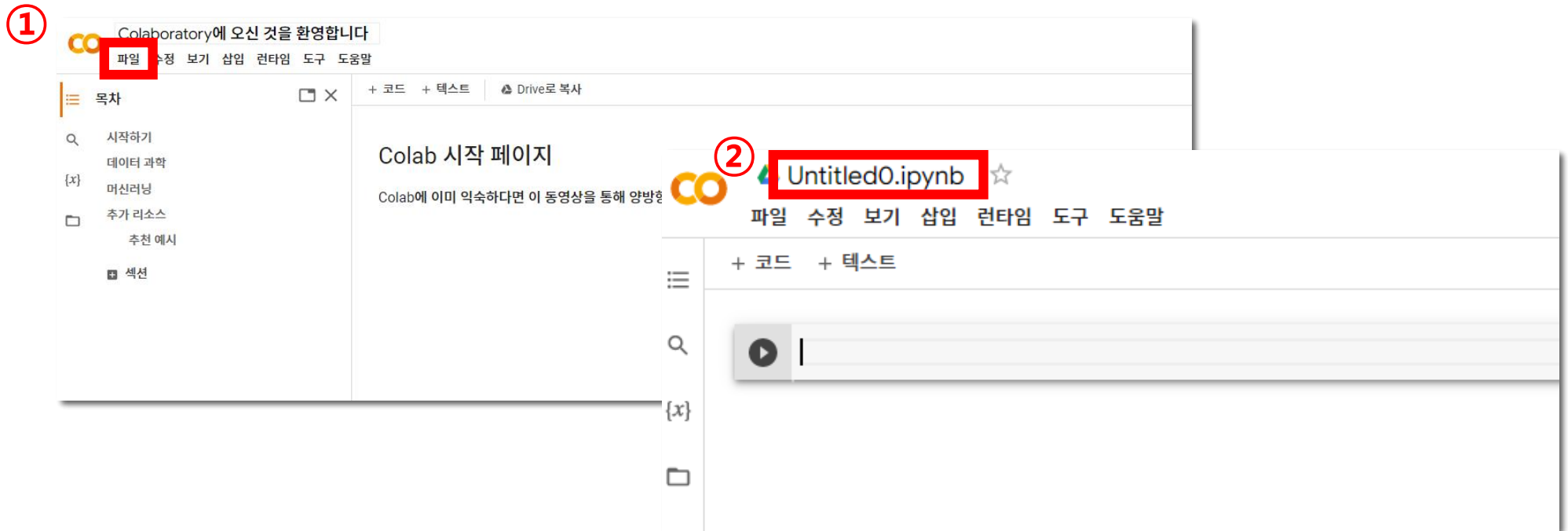
- Colab
 - Python 3.7.13
 - CUDA 11.2
 - Numpy 1.21.6
 - Scikit-learn 1.0.2
- Pytorch 1.12.0
- Pytorch-lightning 1.6.5
- Tensorboard 2.9.1
- Pytorch-geometric 2.0.4 (*graph modeling*)

Environment setup – Colab

<https://colab.research.google.com/?hl=ko>

① [파일]-[새 노트]

② 노트 명 변경



Environment setup – Colab

③ [런타임]-[런타임 유형 변경]

④ 하드웨어 가속기 → GPU

③ Colaboratory에 오신 것을 환영합니다

파일 수정 보기 삽입 런타임 도구 도움말

목차

- 시작하기
- 데이터 과학
- 머신러닝
- 추가 리소스
- 추천 예시
- 섹션

모두 실행 Ctrl+F9

이전 셀 실행 Ctrl+F8

초점이 맞춰진 셀 실행 Ctrl+Enter

선택항목 실행 Ctrl+Shift+Enter

이후 셀 실행 Ctrl+F10

실행 중단 Ctrl+M |

런타임 다시 시작 Ctrl+M .

다시 시작 및 모두 실행

런타임 연결 해제 및 삭제

런타임 유형 변경

세션 관리

런타임 로그 보기

이동

이동 이 동영상 을 통해 양방향 테이블,

3 Cool Colab

4 노트 설정

하드웨어 가속기

GPU ?

Colab을 최대한 활용하려면 필요하지 않은 경우 GPU를 사용하지 않는 것이 좋습니다. [자세히 알아보기](#)

☐ 백그라운드 실행

브라우저를 닫은 후에도 노트북을 계속 실행하고 싶으신가요? [Colab Pro+ 버전으로 업그레이드](#)

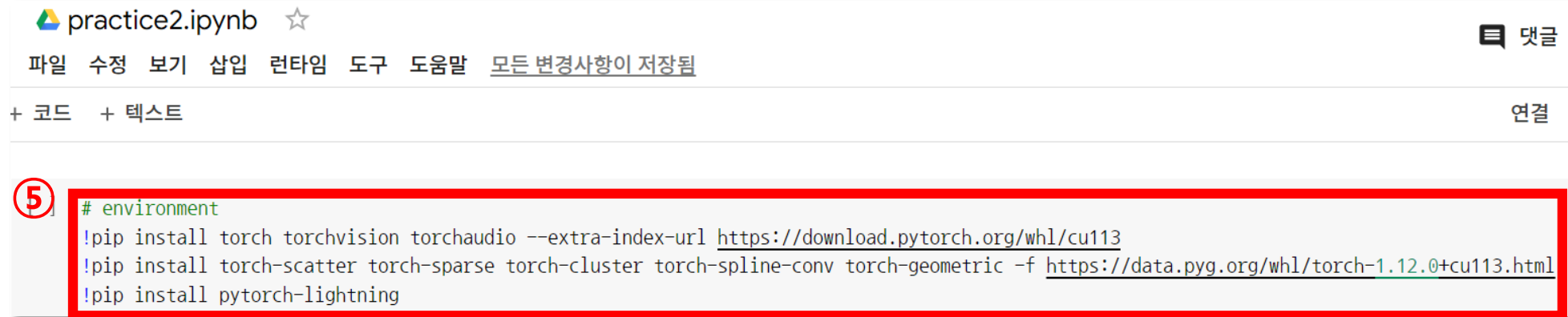
☐ 이 노트를 저장할 때 코드 셀 출력 생략

취소 저장

Environment setup – Colab

⑤ 코드 셀에 설치 명령어 작성

※ Shell 명령어 → !로 시작 (ex. !pip)



The screenshot shows a Google Colab notebook titled "practice2.ipynb". The interface includes a top bar with file management options (파일, 수정, 보기, 삽입, 런타임, 도구, 도움말) and a status message "모든 변경사항이 저장됨". Below the top bar, there are tabs for "+ 코드" and "+ 텍스트". The main content area displays a code cell with the following text:
⑤ `# environment`
`!pip install torch torchvision torchaudio --extra-index-url https://download.pytorch.org/whl/cu113`
`!pip install torch-scatter torch-sparse torch-cluster torch-spline-conv torch-geometric -f https://data.pyg.org/whl/torch-1.12.0+cu113.html`
`!pip install pytorch-lightning`

Environment setup – PyTorch

- <https://pytorch.org/>
- 1.12.0
- Linux
- Pip
- Python
- CUDA 11.3

INSTALL PYTORCH

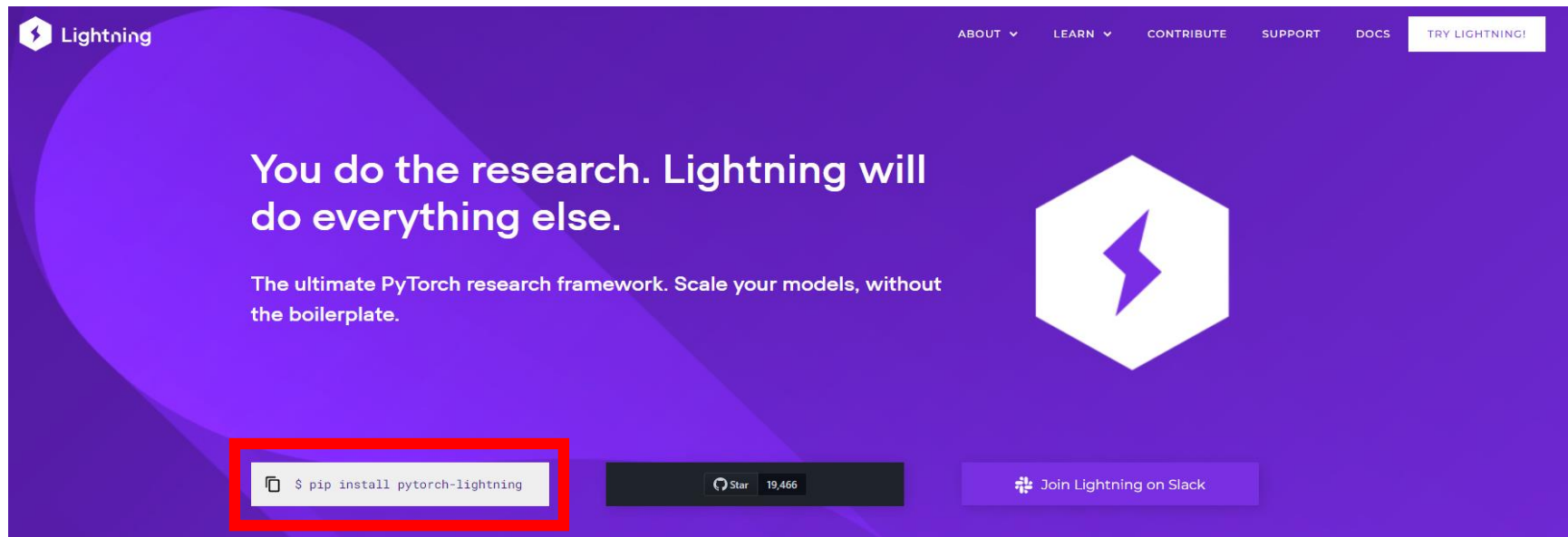
Select your preferences and run the install command. Stable represents the most currently tested and supported version of PyTorch. This should be suitable for many users. Preview is available if you want the latest, not fully tested and supported, 1.12 builds that are generated nightly. Please ensure that you have **met the prerequisites below (e.g., numpy)**, depending on your package manager. Anaconda is our recommended package manager since it installs all dependencies. You can also [install previous versions of PyTorch](#). Note that LibTorch is only available for C++.

Additional support or warranty for some PyTorch Stable and LTS binaries are available through the [PyTorch Enterprise Support Program](#).

| | | | | | |
|-------------------|--|-----------|-------------------|------------|-------------|
| PyTorch Build | Stable (1.12.0) | | Preview (Nightly) | | LTS (1.8.2) |
| Your OS | Linux | | Mac | | Windows |
| Package | Conda | Pip | | LibTorch | Source |
| Language | Python | | | C++ / Java | |
| Compute Platform | CUDA 10.2 | CUDA 11.3 | CUDA 11.6 | ROCm 5.1.1 | CPU |
| Run this Command: | pip3 install torch torchvision torchaudio --extra-index-url https://download.pytorch.org/whl/cu113 | | | | |

Environment setup – PyTorch-Lightning

- <https://www.pytorchlightning.ai/>
- Latest version



Environment setup – PyTorch Geometric

- <https://pytorch-geometric.readthedocs.io/en/latest/notes/installation.html>
- Pytorch 1.12.0
- Linux
- Pip
- CUDA 11.3

INSTALLATION

PyG is available for Python 3.7 to Python 3.10.

Note

We do not recommend installation as a root user on your system Python. Please setup a [Anaconda](#) or [Miniconda](#) environment or create a [Docker image](#).

Quick Start

| | | | | | |
|---------|----------------|------|----------------|---------|-----|
| PyTorch | PyTorch 1.12.* | | PyTorch 1.11.* | | |
| Your OS | Linux | Mac | | Windows | |
| Package | Conda | | Pip | | |
| CUDA | 10.2 | 11.3 | 11.5 | 11.6 | CPU |

Run:

```
pip install torch-scatter torch-sparse torch-cluster torch-spline-conv torch-geometric -f https://data.pyg.org/whl/torch-1.12.0+cu113.html
```

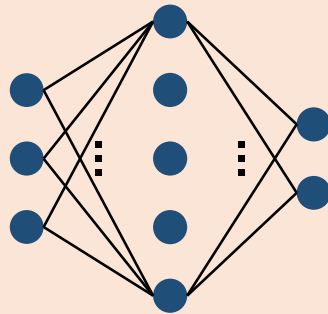
Deep learning process

Step1



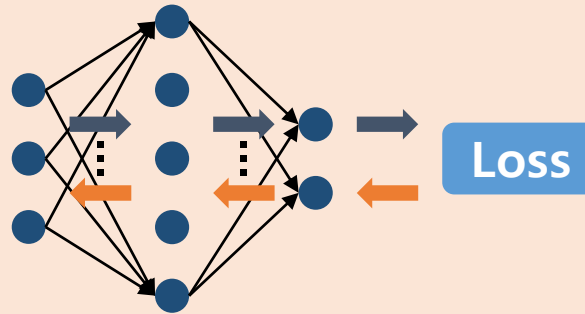
- Collecting data
- Preprocessing
- Dataset
- Data-loader

Step2



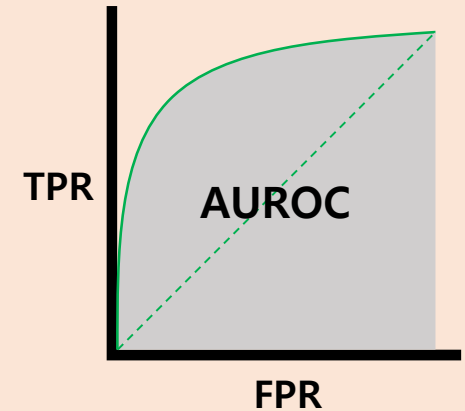
- Selecting method
- Building model

Step3



- Hyper-parameters
- Loss function
- Optimizer
- Training

Step4



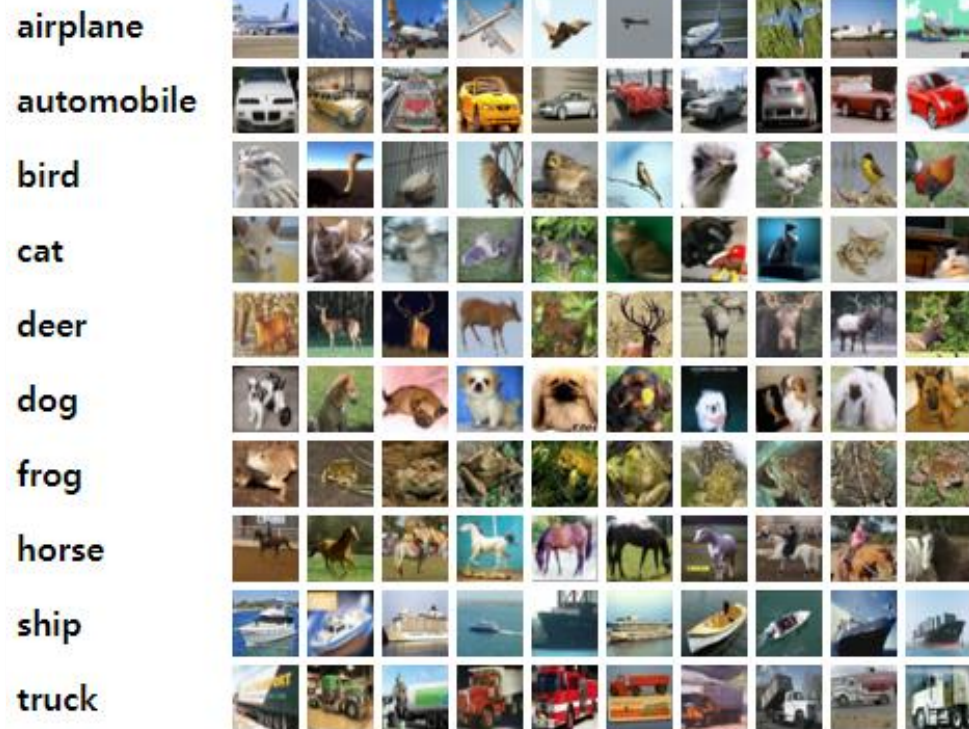
- Predicting
- Evaluating

Preparing input data

CIFAR 10 Dataset

- <https://www.cs.toronto.edu/~kriz/cifar.html>
- The CIFAR-10 dataset consists of 60000 32x32 colour images in 10 classes, with 6000 images per class.

Here are the classes in the dataset, as well as 10 random images from each:



Preparing input data

1. Load and normalize CIFAR10 datasets using torchvision.

```
# Get Dataset
# https://tutorials.pytorch.kr/beginner/blitz/cifar10\_tutorial.html
import torch
import torchvision
import torchvision.transforms as transforms

transform = transforms.Compose([transforms.ToTensor(), transforms.Normalize((0.5, 0.5, 0.5), (0.5, 0.5, 0.5))])

train_set = torchvision.datasets.CIFAR10(root='./data', train=True, download=True, transform=transform)
test_set = torchvision.datasets.CIFAR10(root='./data', train=False, download=True, transform=transform)
classes = ('plane', 'car', 'bird', 'cat', 'deer', 'dog', 'frog', 'horse', 'ship', 'truck')

print('*=====Train Set=====')
print(train_set)
print('*===== Test Set =====')
print(test_set)
```


Preparing input data

- Watch the image in datasets

```
import matplotlib.pyplot as plt
import numpy as np
```

```
img = train_set[1][0]
label = train_set[1][1]
img = img.numpy()
```

① `print(img.shape)`

② `print(f'{classes[label]}-{label}')`

```
origin_img = img / 2 + 0.5 # unnormalized
plt.imshow(np.transpose(origin_img, (1, 2, 0)))
```

③ `plt.show()`

```
plt.imshow(np.transpose(img, (1, 2, 0)))
```

④ `plt.show()`

Preparing input data

- Watch the image in datasets

(3, 32, 32) ①
truck-9 ②

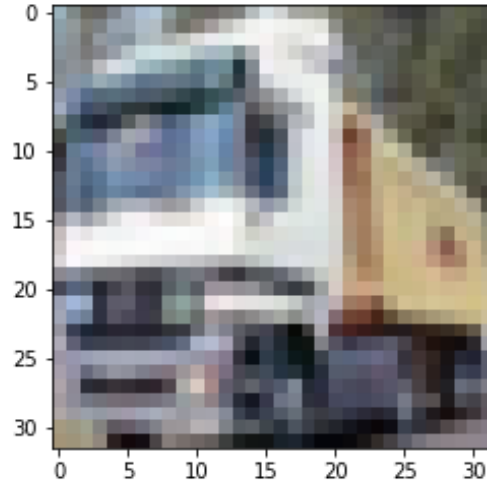
```
import matplotlib.pyplot as plt
import numpy as np

img = train_set[1][0]
label = train_set[1][1]
img = img.numpy()
```

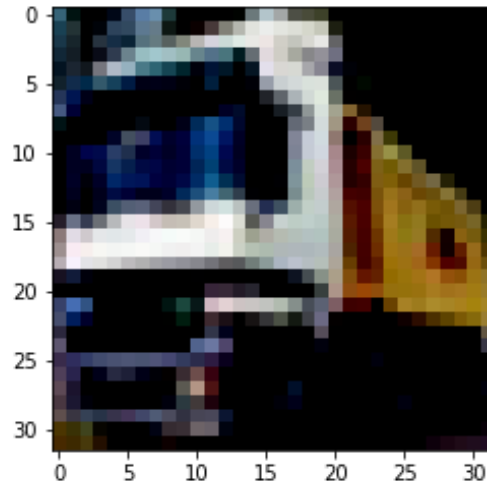
① `print(img.shape)`
② `print(f'{classes[label]}')`

`origin_img = img / 2 + 0.1`
③ `plt.imshow(np.transpose(origin_img, (2, 1, 0)))`
`plt.show()`

④ `plt.imshow(np.transpose(img, (2, 1, 0)))`
`plt.show()`



Clipping input data to the valid range for imshow with RGB data ([0..1] for floats or [0..255] for integers).



← Normalize 적용

Preparing input data

2. Define a data module

```
from torch.utils.data import DataLoader
from pytorch_lightning import LightningDataModule

class CustomDataModule(LightningDataModule):
    def __init__(self, training_set, test_set, batch_size=128, num_workers=2):
        super().__init__()
        self.training_set = training_set
        self.test_set = test_set

        # set parameters
        self.batch_size = batch_size
        self.num_workers = num_workers

    #how to make a batch
    def collate_function(self, batch):
        new_x = [x.flatten() for x, _ in batch]
        new_y = [torch.tensor(y, dtype=torch.long) for _, y in batch]
        return [torch.stack(new_x), torch.stack(new_y)]

    def train_dataloader(self):
        return DataLoader(self.training_set, batch_size=self.batch_size,
                          num_workers=self.num_workers, shuffle=True,
                          collate_fn=self.collate_function)

    def test_dataloader(self):
        return DataLoader(self.test_set, batch_size=self.batch_size,
                          num_workers=self.num_workers,
                          collate_fn=self.collate_function)

# data_module = CustomDataModule(train_set, test_set)
# trainer.fit(model, datamodule = data_module)
```


Construct a model

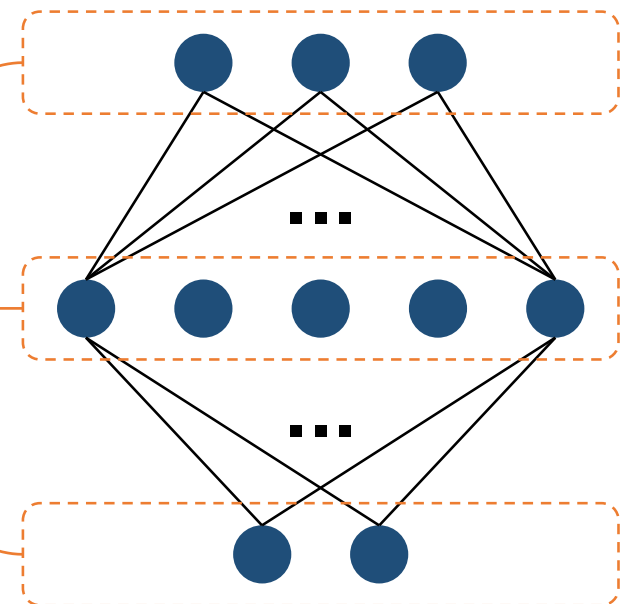
```
import torch.nn as nn
import torch.nn.functional as F

class DeepNeuralNet(nn.Module):
    def __init__(self, input_size, label_size, layer_size=64):
        super().__init__()
        self.input_size = input_size
        self.layer_size = layer_size
        self.label_size = label_size
        self.set_layer()

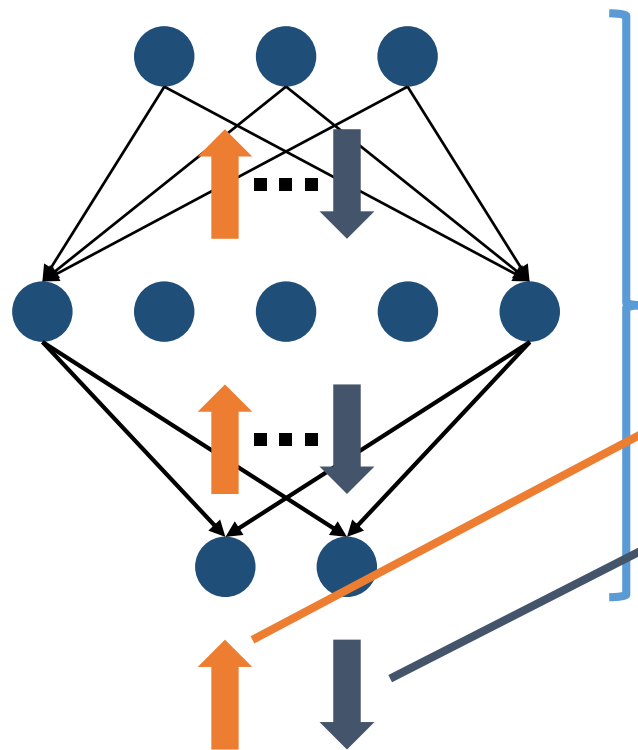
    def set_layer(self):
        # Build layer
        self.first_layer = nn.Linear(self.input_size, self.layer_size)
        self.second_layer = nn.Linear(self.layer_size, self.layer_size)
        self.last_layer = nn.Linear(self.layer_size, self.label_size)

    def forward(self, x):
        # Define operation
        x = self.first_layer(x)
        x = F.relu(x)
        x = self.second_layer(x)
        x = F.relu(x)
        x = self.last_layer(x)

    return x
```



Construct a model



Loss

```
from pytorch_lightning import LightningModule

from torch.nn import CrossEntropyLoss, Softmax
from torch.optim import Adam

class LitCIFAR10Model(LightningModule):
    def __init__(self, model, learning_rate=1e-3):
        super().__init__()
        self.model = model
        self.lr = learning_rate

    def forward(self, x):
        out = self.model(x)
        return Softmax(dim=1)(out)

    def configure_optimizers(self):
        optimizer=Adam(self.parameters(), lr=self.lr)
        return optimizer

    def training_step(self, batch, batch_idx):
        x, y = batch
        logits = self(x)
        loss = self.loss_function(logits, y)
        predict = torch.argmax(logits, dim=1)
        return {'loss': loss, 'predict': predict, 'answer': y}

    def training_epoch_end(self, outputs):
        loss = [output['loss'] for output in outputs]
        avg_loss = sum(loss) / len(outputs)
        self.logger.experiment.add_scalar("Loss/Epoch",
                                           avg_loss,
                                           self.current_epoch)

    def predict_step(self, batch, batch_idx):
        x, y = batch
        logits = self.model(x)
        predict = torch.argmax(logits, dim=1)
        return predict

    def loss_function(self, output, target):
        return CrossEntropyLoss()(output, target)
```

Training a model

```
from pytorch_lightning import Trainer
from pytorch_lightning.loggers import TensorBoardLogger

logger = TensorBoardLogger('./log', name="LitCIFAR10-DeepNN-Adam-128,1e-3")
data_module = CustomDataModule(train_set, test_set)

architecture = DeepNeuralNet(32*32*3, 10) #32*32(pixels)*3 channels, labels
model = LitCIFAR10Model(architecture)

trainer = Trainer(max_epochs=10, accelerator='gpu', devices=[0], logger=logger)
#trainer = Trainer(max_epochs=10, accelerator='cpu', logger=logger)

trainer.fit(model, datamodule=data_module)
```

| | Name | Type | Params |
|---|-------|---------------|--------|
| 0 | model | DeepNeuralNet | 201 K |

| | |
|-------|--|
| 201 K | Trainable params |
| 0 | Non-trainable params |
| 201 K | Total params |
| 0.806 | Total estimated model params size (MB) |

/usr/local/lib/python3.7/dist-packages/pytorch_lightning/loggers/tensorboard.py:251: UserWarning: Could not log computational graph since the `model.example`
"Could not log computational graph since the"

Epoch 9: 100%  391/391 [00:11<00:00, 35.33it/s, loss=1.14, v_num=0, train_loss=1.150]

INFO:pytorch_lightning.utilities.rank_zero:`Trainer.fit` stopped: `max_epochs=10` reached.

Visualize the result with TensorBoard

```
logger = TensorBoardLogger('./log', name="LitCIFAR10-DeepNN-Adam-128,1e-3")
```

```
trainer = Trainer(max_epochs=10, accelerator='gpu', devices=[0], logger=logger)
```

```
def training_epoch_end(self, outputs):  
    loss = [output['loss'] for output in outputs]  
    avg_loss = sum(loss) / len(outputs)  
    self.logger.experiment.add_scalar("Loss/Epoch",  
                                      avg_loss,  
                                      self.current_epoch)
```

```
%load_ext tensorboard  
%tensorboard --logdir "./log"
```

Visualize the result with TensorBoard

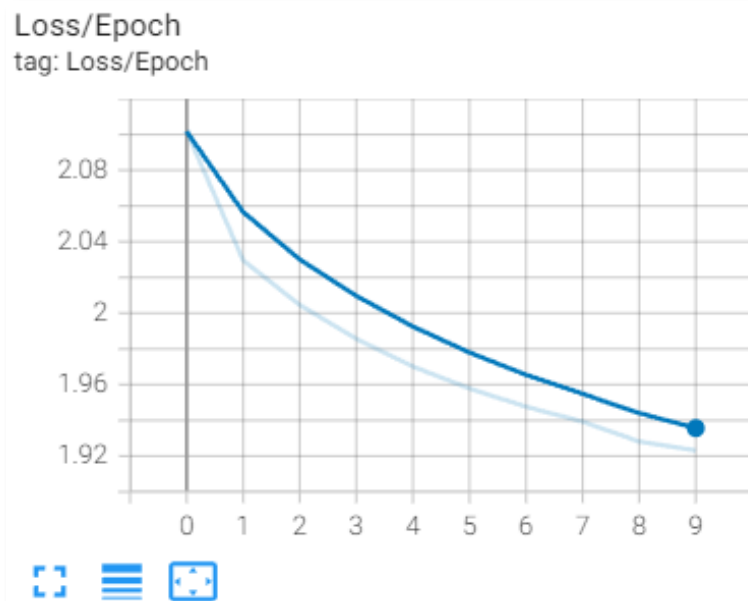
```
logger = TensorBoardLogger('./log', name="LitCIFAR10-DeepNN-Adam-128,1e-3")
```

```
trainer = Trainer(max_epochs=10, accelerator='gpu', devices=[0], logger=logger)
```

```
def training_epoch_end(self, outputs):  
    loss = [output['loss'] for output in outputs]  
    avg_loss = sum(loss) / len(outputs)  
    self.logger.experiment.add_scalar("Loss/Epoch",  
                                      avg_loss,  
                                      self.current_epoch)
```

```
%load_ext tensorboard  
%tensorboard --logdir "./log"
```

403 error -> 쿠키 허용 -> 재시작



Evaluate the result

- $Accuracy = \frac{\text{correct predictions}}{\text{all predictions}}$

```
predict = trainer.predict(model, dataloaders=data_module.test_dataloader())
predict = torch.cat(predict)

answer = [batch[1] for batch in data_module.test_dataloader()]
answer = torch.cat(answer)

correct = (predict == answer).sum()
total = answer.shape[0]
accuracy = correct / total

print(f'correct: {correct}, total: {total}, accuracy: {accuracy}')
```

- torchmetrics API can be used for this instead

```
# You can use torchmetric
import torchmetrics

acc = torchmetrics.functional.accuracy(predict, answer)
print(f"accuracy-torchmetrics: {acc}")
```

Evaluate the result

- $Accuracy = \frac{\text{correct predictions}}{\text{all predictions}}$

```
predict = trainer.predict(model, dataloaders=data_module.test_dataloader())
predict = torch.cat(predict)

answer = [batch[1] for batch in data_module.test_dataloader()]
answer = torch.cat(answer)

correct = (predict == answer).sum()
total = answer.shape[0]
accuracy = correct / total

print(f'correct: {correct}, total: {total}, accuracy: {accuracy}')
```

- torchmetrics API can be used for this instead

```
# You can use torchmetric
import torchmetrics
```

```
correct: 4879, total: 10000, accuracy: 0.48789998888896942
accuracy-torchmetrics: 0.48789998888896942
```

```
functional.accuracy(predict, answer)
torchmetrics: {acc}"
```

Practice – GNN

1. TUDataset

- <https://chrsmrrs.github.io/datasets/>
- 120+ datasets for graph classification and regression
- Tox21 AhR (example)

| | |
|---------------------------|--|
| Tox21_AhR_training(8169) | Data(edge_index=[2, 52], x=[25, 50], edge_attr=[52, 4], y=[1]) |
| Tox21_AhR_testing(272) | Data(edge_index=[2, 44], x=[20, 51], edge_attr=[44, 4], y=[1]) |
| Tox21_AhR_evaluation(607) | Data(edge_index=[2, 118], x=[53, 53], edge_attr=[118, 4], y=[1]) |

- x: Node features
- y: Answer
- Edge_index: Directed edge (Beginning node, End node)
- Edge_attr: Edge features

```
# Dataset
```

```
from torch_geometric.datasets import TUDataset
```

```
training_dataset = TUDataset('./dataset', 'Tox21_AhR_training')  
validation_dataset = TUDataset('./dataset', 'Tox21_AhR_testing')  
test_dataset = TUDataset('./dataset', 'Tox21_AhR_evaluation')
```


Practice – GNN

2. Construct model

```
# Pytorch-lightning datamodule
from torch.utils.data import DataLoader
from pytorch_lightning import LightningDataModule
from torch_geometric.data import Data, Batch
from torch.nn.functional import pad

class CustomData(LightningDataModule):
    def __init__(self, training_set, validation_set, test_set, batch_size=128, num_workers=1):
        super().__init__()
        self.training_set = training_set
        self.validation_set = validation_set
        self.test_set = test_set
        self.batch_size = batch_size
        self.num_workers = num_workers

    def collate_function(self, batch):
        return Batch.from_data_list([Data(edge_index=data.edge_index,
                                          x=pad(data.x, (0,3), 'constant', 0.)[:, :53],
                                          edge_attr=data.edge_attr,
                                          y=data.y.unsqueeze(0).float()) for data in batch])

    def train_dataloader(self):
        return DataLoader(self.training_set, batch_size=self.batch_size, num_workers=self.num_workers, shuffle=True, collate_fn=self.collate_function)

    def val_dataloader(self):
        return DataLoader(self.validation_set, batch_size=self.batch_size, num_workers=self.num_workers, collate_fn=self.collate_function)

    def test_dataloader(self):
        return DataLoader(self.test_set, batch_size=self.batch_size, num_workers=self.num_workers, collate_fn=self.collate_function)
```

Practice – GNN

2. Construct model

• GATConv

- In_channels
- out_channels
- Dropout
- Heads
 - Multi-head attention
- Concat
- ...

```
# Pytorch and pytorch-geometric module
from torch.nn import Module, Linear
from torch_geometric.nn import GATConv, global_mean_pool

class CustomGAT(Module):
    def __init__(self, input_size, label_size, layer_size=64, dropout=0.1, heads=2):
        super().__init__()
        self.input_size = input_size
        self.layer_size = layer_size
        self.dropout = dropout
        self.heads = heads
        self.label_size = label_size
        self.setup()

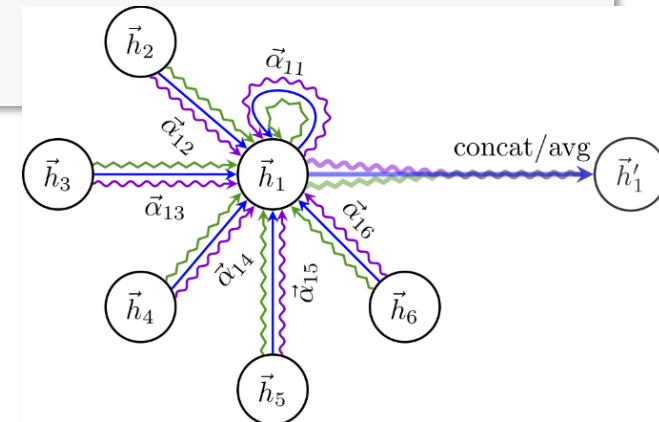
    def setup(self):
        self.first_layer = GATConv(self.input_size, self.layer_size, dropout=self.dropout, heads=self.heads, concat=False)
        self.last_layer = GATConv(self.layer_size, self.layer_size, dropout=self.dropout, heads=self.heads, concat=False)
        self.ffnn = Linear(self.layer_size, self.label_size)

    def convert_graph_into_single_vector(self, graph_hidden, batch_index):
        return global_mean_pool(graph_hidden, batch_index)

    def forward(self, batch):
        z = self.first_layer(batch.x, batch.edge_index)
        z = self.last_layer(z, batch.edge_index)
        z = self.convert_graph_into_single_vector(z, batch.batch)
        z = self.ffnn(z)
        return z
```

• Global_mean_pool

- Returns batch-wise graph-level-outputs by averaging node features across the node dimension



Practice – GNN

2. Construct model

- GATConv

- In_channels
- out_channels
- Dropout
- Heads
- Multi-head attention
- Concat
- ...

- Global_mean_pool

- Returns batch-wise graph-level-outputs across the node dimension

```
# Pytorch and pytorch-geometric module
from torch.nn import Module, Linear
from torch_geometric.nn import GATConv, global_mean_pool
```

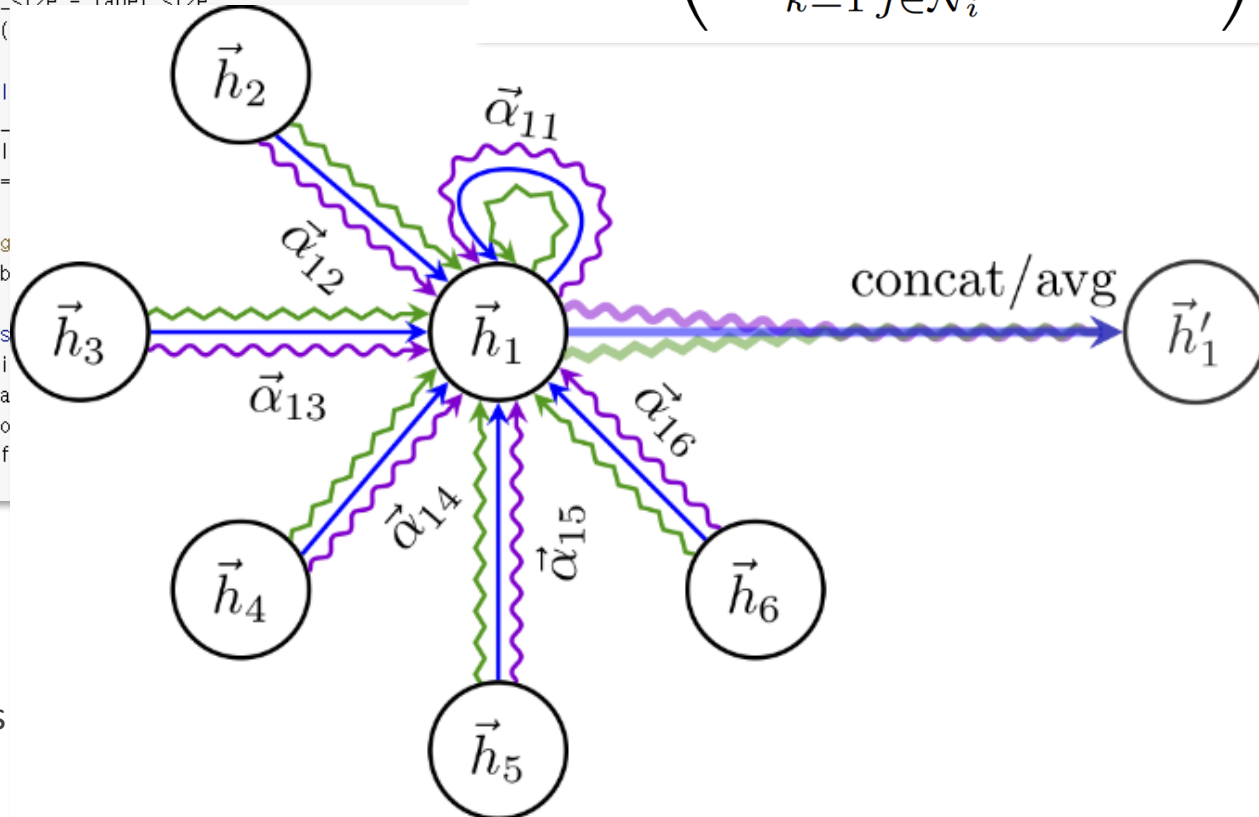
```
class CustomGAT(Module):
    def __init__(self, input_size, label_size, layer_size):
        super().__init__()
        self.input_size = input_size
        self.layer_size = layer_size
        self.dropout = dropout
        self.heads = heads
        self.label_size = label_size
        self.setup()

    def setup(self, self, first_layer, last_layer, ffn):
        self.first_layer = first_layer
        self.last_layer = last_layer
        self.ffn = ffn

    def convert_graph(self, graph):
        return graph

    def forward(self, x, edge_index, edge_weight):
        z = self.first_layer(x)
        z = self.last_layer(z)
        z = self.concat(z)
        z = self.ffn(z)
        return z
```

$$\vec{h}'_i = \sigma \left(\frac{1}{K} \sum_{k=1}^K \sum_{j \in \mathcal{N}_i} \alpha_{ij}^k \mathbf{W}^k \vec{h}_j \right)$$



GAT image from <https://github.com/PetarV-/GAT>

Practice – GNN

2. Construct model

```
# Pytorch-lightning module
from pytorch_lightning import LightningModule
from torch.nn import BCEWithLogitsLoss
from torch.optim import Adam

class CustomModel(LightningModule):
    def __init__(self, model, learning_rate=1e-3):
        super().__init__()
        self.model = model
        self.lr = learning_rate

    def forward(self, batch, mode):
        z = self.model(batch)
        loss = self.loss_function(z, batch.y)
        self.log(f"{mode}_loss", loss, batch_size=batch.y.size(0), prog_bar=True, on_step=False, on_epoch=True)
        return loss, z, batch.y

    def training_step(self, batch, batch_idx):
        loss, predict, answer = self(batch, 'train')
        return {'loss': loss, 'predict': predict, 'answer': answer}

    def validation_step(self, batch, batch_idx):
        loss, predict, answer = self(batch, 'val')
        return {'loss': loss, 'predict': predict, 'answer': answer}

    def test_step(self, batch, batch_idx):
        loss, predict, answer = self(batch, 'test')
        return {'loss': loss, 'predict': predict, 'answer': answer}

    def predict_step(self, batch, batch_idx):
        predict = self.model(batch)
        return predict

    def loss_function(self, output, target):
        return BCEWithLogitsLoss(reduction='mean')(output, target)

    def configure_optimizers(self):
        optimizer = Adam(self.parameters(), lr=self.lr)
        return optimizer
```

Practice – GNN

3. Training the model

```
# Training
from pytorch_lightning import Trainer

import warnings
warnings.filterwarnings(action='ignore')
# warnings.filterwarnings(action='default')

data_module = CustomData(training_dataset, validation_dataset, test_dataset)

gat = CustomGAT(53, 1)
model = CustomModel(gat)

trainer = Trainer(max_epochs=1, accelerator='gpu', devices=[0])
trainer.fit(model, datamodule=data_module)
```

| | Name | Type | Params |
|---|-------|-----------|--------|
| 0 | model | CustomGAT | 15.7 K |

| | |
|--------|--|
| 15.7 K | Trainable params |
| 0 | Non-trainable params |
| 15.7 K | Total params |
| 0.063 | Total estimated model params size (MB) |

Epoch 0: 100%  67/67 [00:12<00:00, 5.16it/s, loss=0.384, v_num=0, val_loss=0.361, train_loss=0.463]

Trainer.fit` stopped: `max_epochs=1` reached.

Practice – GNN

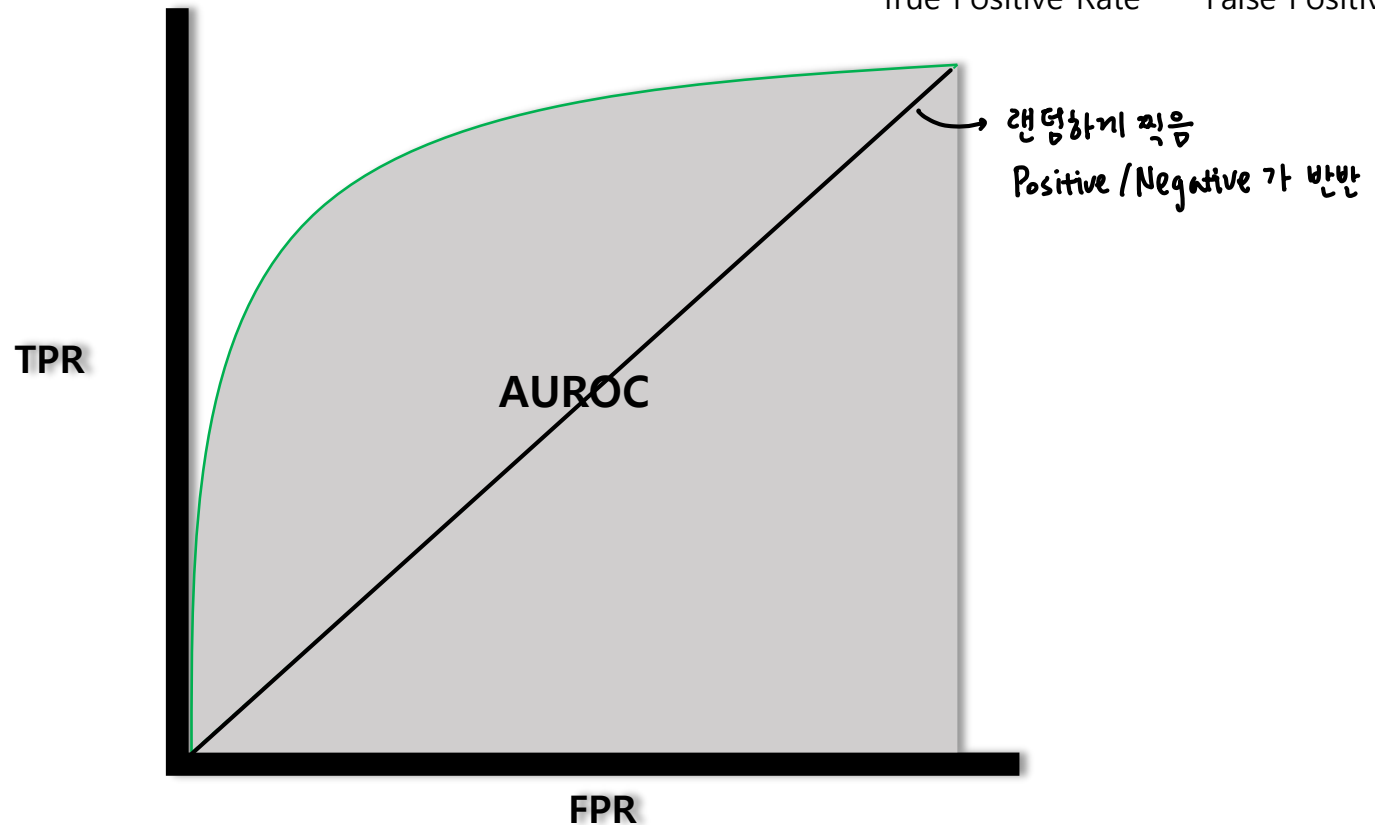
4. Evaluate the result

- Metric – `torchmetrics.functional.auroc`
 - $0.5 == \text{random}$

| Confusion Matrix | | Real | |
|------------------|----------|----------------|----------------|
| | | Positive | Negative |
| Predict | Positive | True Positive | False Positive |
| | Negative | False Negative | True Negative |

True Positive Rate

False Positive Rate



Practice – GNN

4. Evaluate the result

- Metric – `torchmetrics.functional.auroc`

```
torchmetrics.functional.auroc ( preds , target , num_classes = None , pos_label = None ,
```

```
# Evaluation
import torch
from torchmetrics.functional import auroc

outputs = trainer.predict(model, dataloaders=data_module.test_data_loader())
y = torch.concat(outputs)
x = torch.concat([batch.y for batch in data_module.test_data_loader()]).int()

evaluation = auroc(y, x)
print(f"auc-roc: {evaluation}")
```

Predicting DataLoader 0: 100%  5/5 [00:00<00:00, 20.68it/s]

auc-roc: 0.4389864206314087

References

- <https://www.cs.toronto.edu/~kriz/cifar.html>
- https://tutorials.pytorch.kr/beginner/blitz/cifar10_tutorial.html
- https://pytorch-lightning.readthedocs.io/en/stable/notebooks/lightning_examples/cifar10-baseline.html
- Golovko, Vladimir, et al. "A shallow convolutional neural network for accurate handwritten digits classification." *International Conference on Pattern Recognition and Information Processing*. Springer, Cham, 2016.
- Veličković, Petar, et al. "Graph attention networks." *arXiv preprint arXiv:1710.10903* (2017).

Thank you