

---

# Zero-Shot Knowledge Distillation in Deep Networks

Gaurav Kumar Nayak<sup>\*1</sup> Konda Reddy Mopuri<sup>\*2</sup> Vaisakh Shaj<sup>\*3</sup> R. Venkatesh Babu<sup>1</sup>  
Anirban Chakraborty<sup>1</sup>

ICML 2019

Cite: 77

2021.06.16

임진혁

1. Introduction
2. Proposed Approach
3. Experiments

## Limitation of KD(Knowledge Distillation)

KD is very usefully used in many fields

There is a limitation that there is no original data.

In such cases only trained model is available without training data

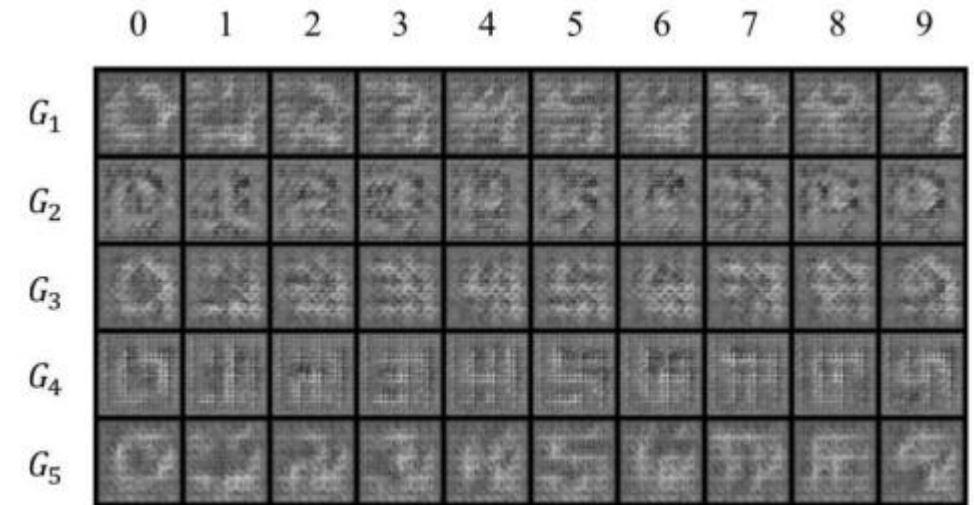
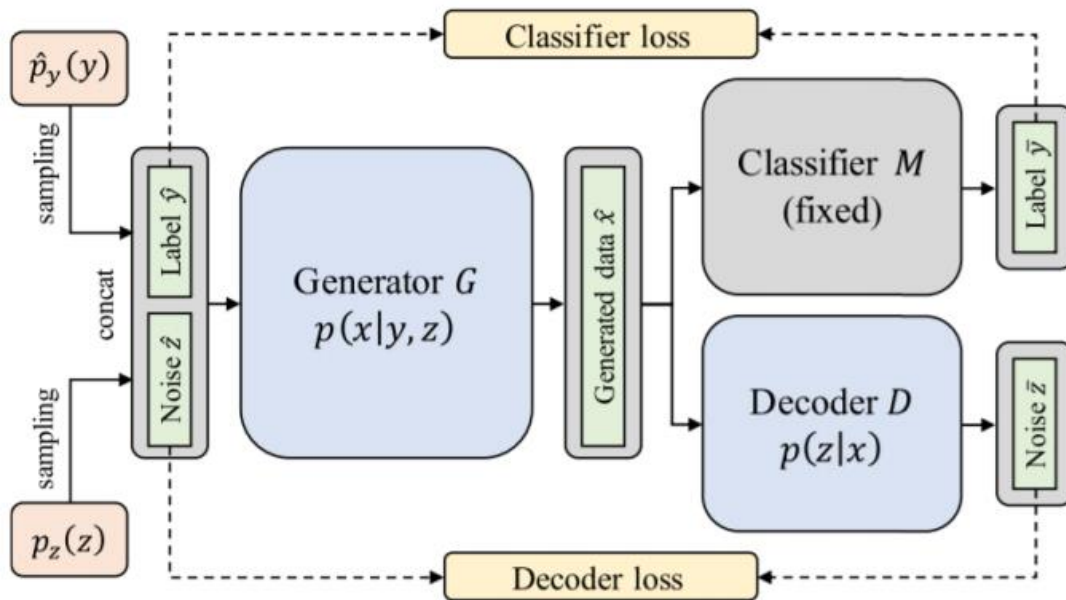
- Medical diagnosis, that patients' privacy prohibits distribution
- Commercial products with deployed models
- Cost from annotating data
- Proprietary data (JFT-300M, SFC. )

## **Data-free Knowledge Distillation (Zero Shot)**

with no data samples and no extracted prior information

So there should be “Transfer Set” instead of “original data”

## Knowledge Extraction with No Observable Data(Nips 2019)

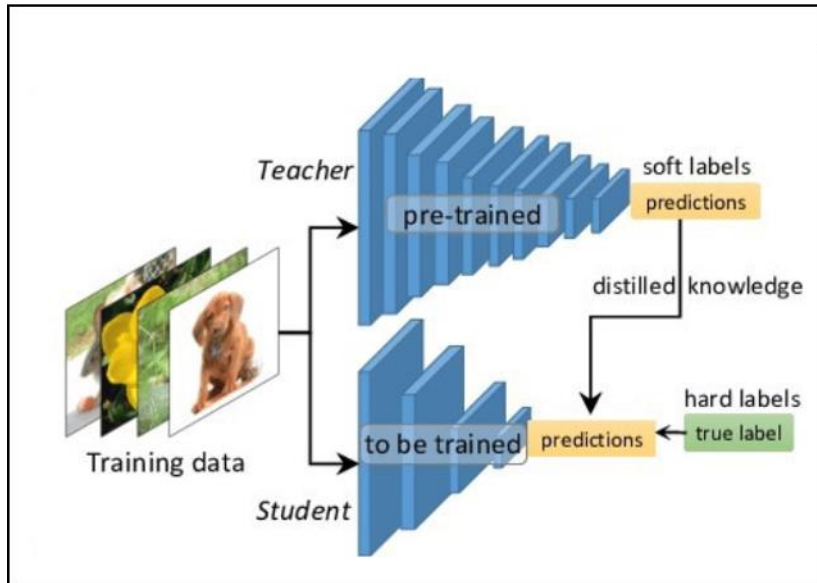


(a) MNIST ( $z = 0$ ).

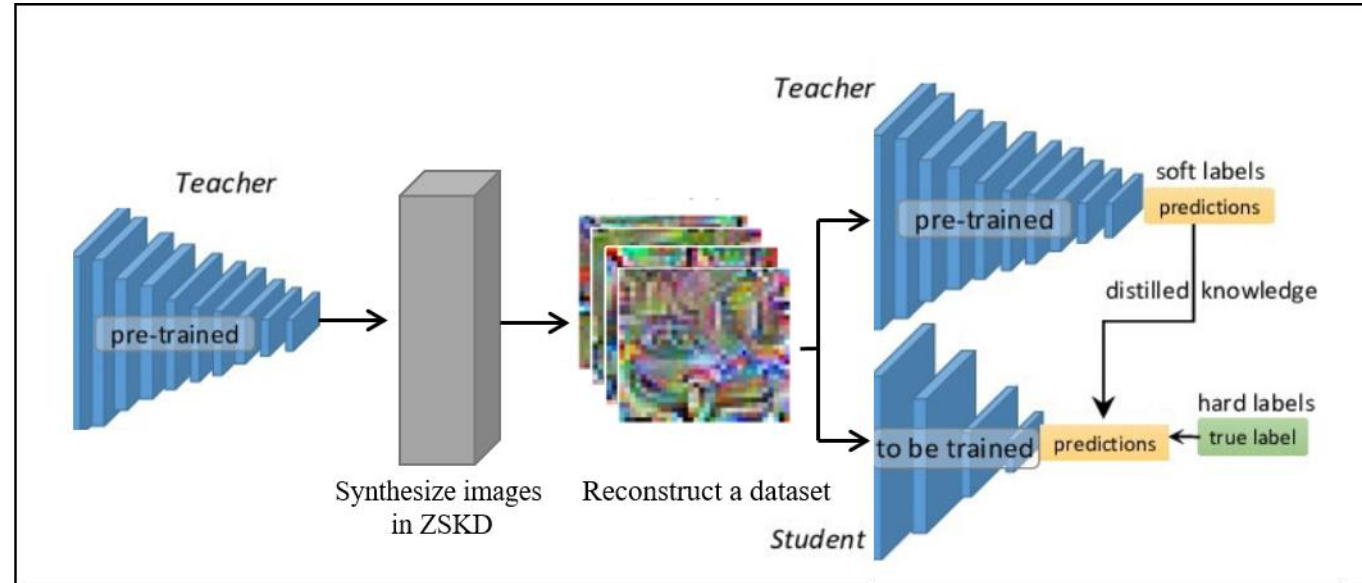
$$\mathcal{D} = \{\operatorname{argmax}_{\hat{x}} p(\hat{x}|\hat{y}, \hat{z}) \mid \hat{y} \sim \hat{p}_y(y) \text{ and } \hat{z} \sim p_z(z)\}$$

## Zero-Shot Knowledge Distillation (ICML 2019)

Make Train Data Distribution from Teacher's parameter



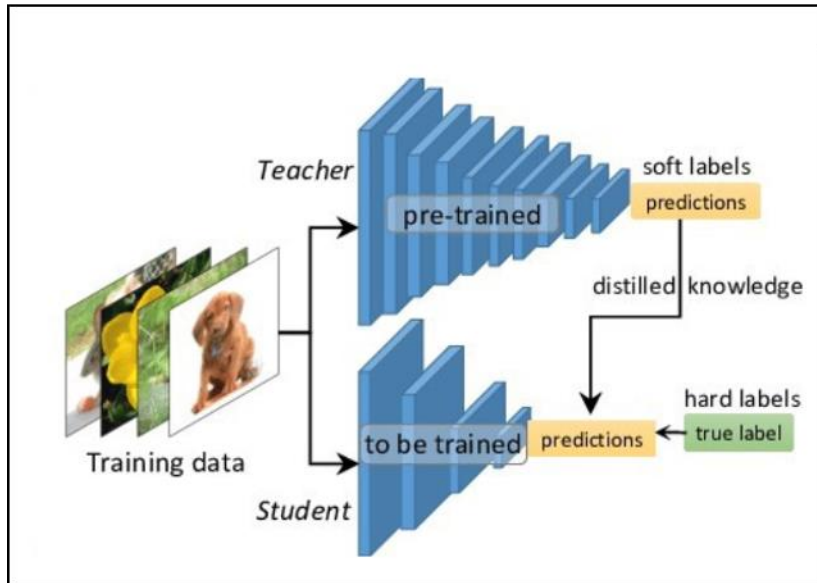
Conventional knowledge distillation



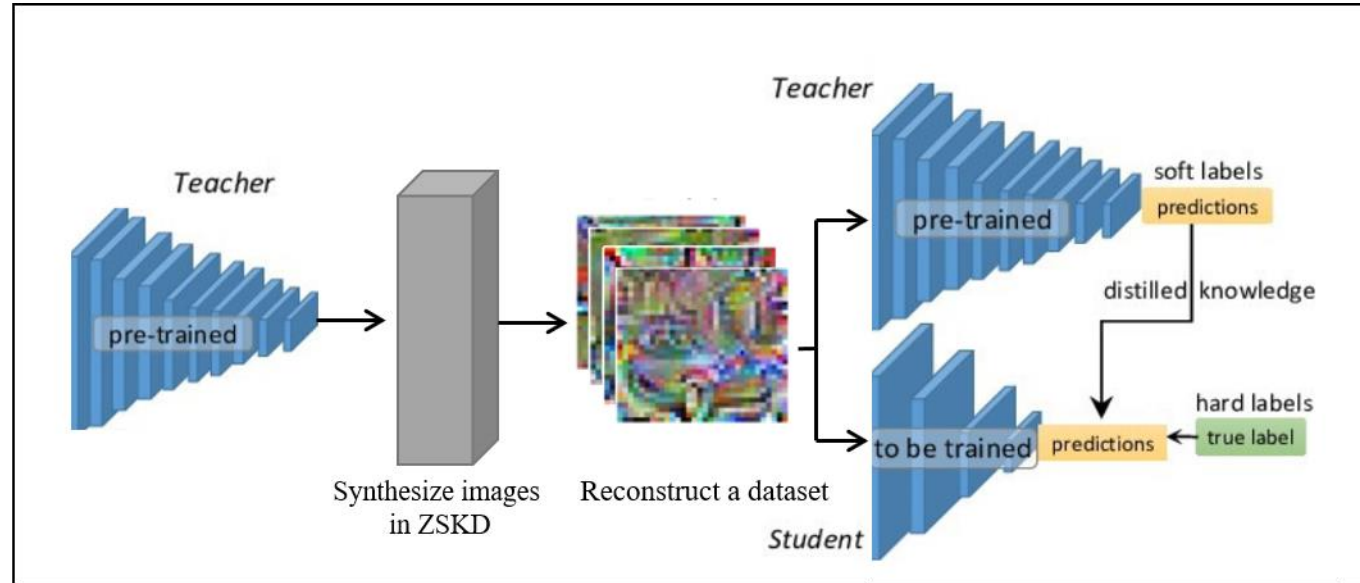
Data-free knowledge distillation  
(Zero-shot knowledge distillation)

## Zero-Shot Knowledge Distillation (ICML 2019)

Make Train Data Distribution from Teacher's parameter



Conventional knowledge distillation



Data-free knowledge distillation  
(Zero-shot knowledge distillation)



## Zero-Shot Knowledge Distillation (ICML 2019)

### Using Class Similarity Matrix

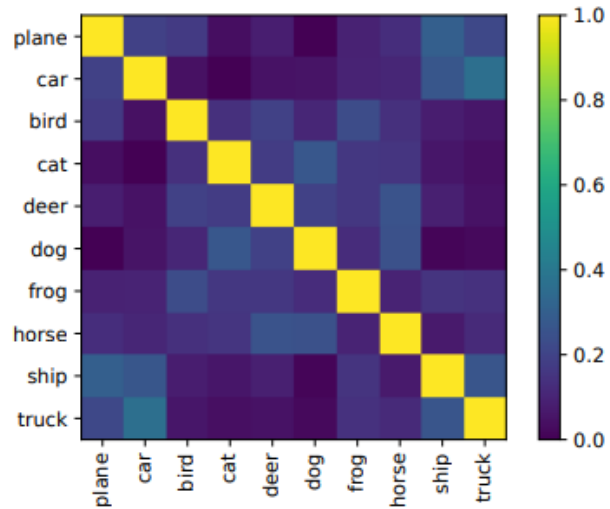
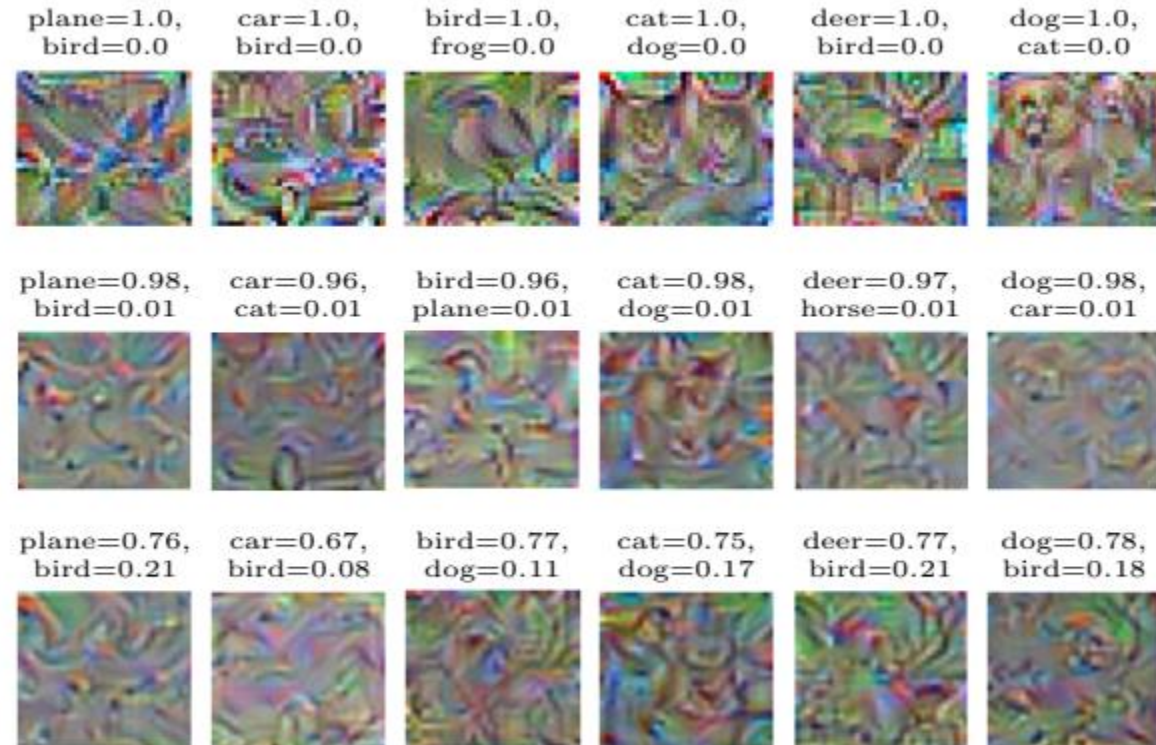
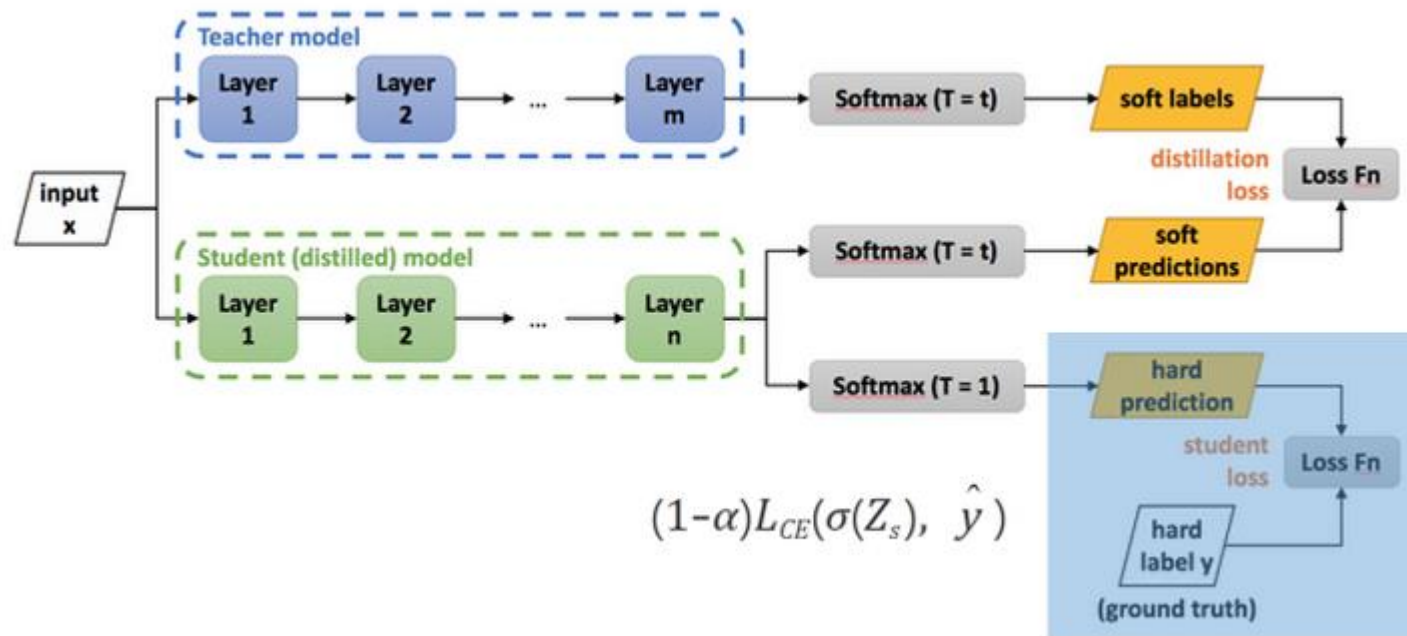


Figure 1. Class similarity matrix computed for the *Teacher* model trained over CIFAR-10 dataset. Note that the class labels are mentioned and the learned similarities are meaningful.



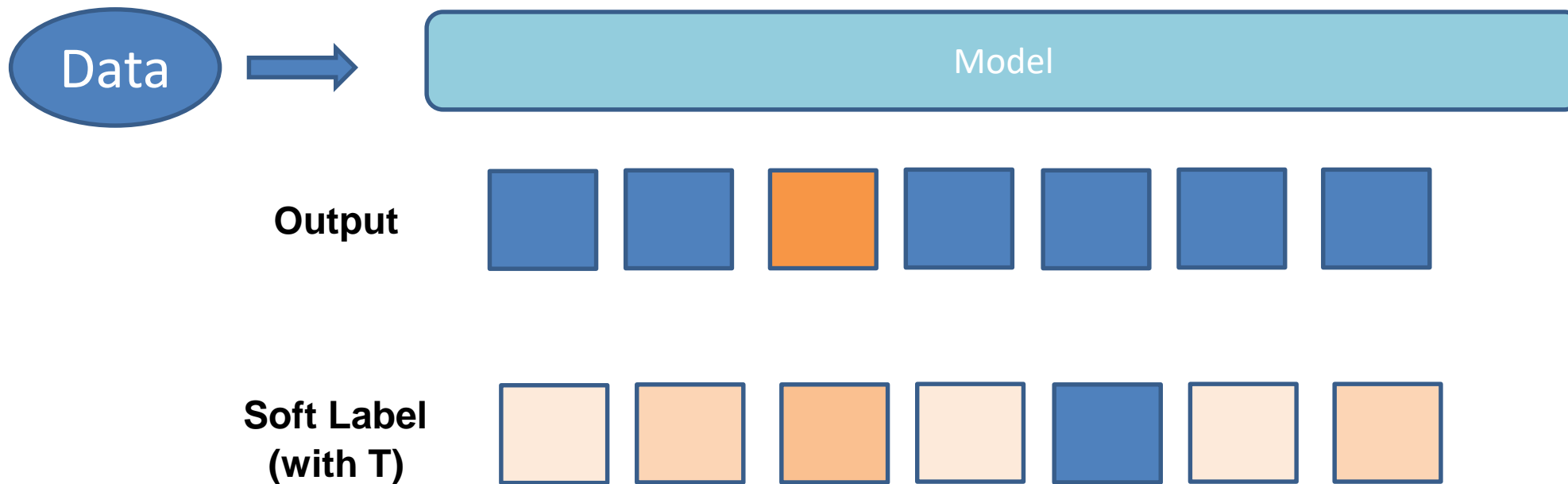


## Knowledge Distillation (Baseline)



$$L = \sum_{(x,y) \in \mathbb{D}} L_{KD}(S(x, \theta_S, \tau), T(x, \theta_T, \tau)) + \lambda L_{CE}(\hat{y}_S, y)$$

### Knowledge Distillation (Baseline)



### Knowledge Distillation (Baseline)



### ZSKD

Soft Label  
(with T)



model output space of the *Teacher* model.

- **Model(동사) output space of the Teacher model**
- **$\mathbf{s} \sim \mathbf{p}(\mathbf{s})$  : Random Vector** that represents the **softmax outputs** of the *Teacher* using **Dirichlet distribution**

## 2<sup>MIL</sup> Proposed Method

### ZSKD

**P(s)**  
**Dirichlet distribution**



x 없이 y\_pred를 만든다고 생각하면 된다. ( 이 경우에는 softmax output)  
y\_pred는 **Dirichlet distribution**에서 Sampling한 값이 될 것이다.

$$\bar{x}_i^k = \operatorname{argmin}_x L_{CE}(\mathbf{y}_i^k, T(x, \theta_T, \tau))$$

## Dirichlet Distribution(디리클레 분포)

model output space of the *Teacher* model.

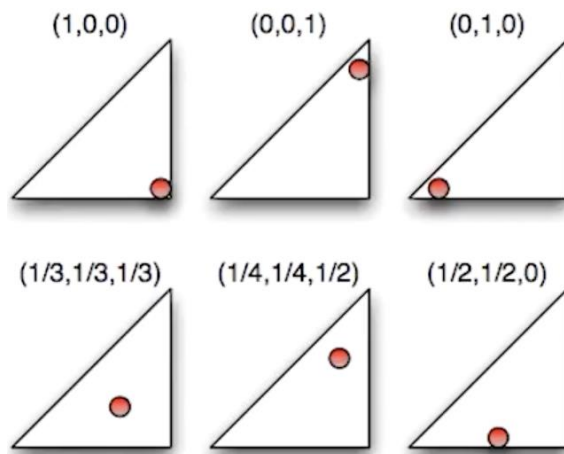
- 분포를 추정하는 방법은 크게 2가지 (Parametric , Nonparametric)
- Parametric
  - 모수를 가정( 평균,분산)
  - 데이터가 적어도 모수 분포를 잘 가정하면 좋은 추정이 된다
  - 예시: 가우시안 분포 /이항 분포/ 베타 분포/ 다변량 분포/ 디리클레 분포 (다변량 , 연속형)
- Nonparametric
  - 모수를 가정하지 않음
  - 데이터가 많을 수록 좋은 추정이 가능

### Dirichlet Distribution(디리클레 분포)

model output space of the *Teacher* model.

- $\alpha_k$ 에 대하여  $k$ 개의 연속형 확률변수에 대응되는  $k$ 개의 continuous values를 사용하여 분포 표현
- 확률 정의에 따라 해당 continuous random variables은 0보다는 크고 합하면 1이 된다.
- LDA에 쓰이는 바로 그 분포
- $k=3$ 일 때, 2차원으로 시각화

값이 1보다 클수록  
다양한 차원으로 퍼진다.





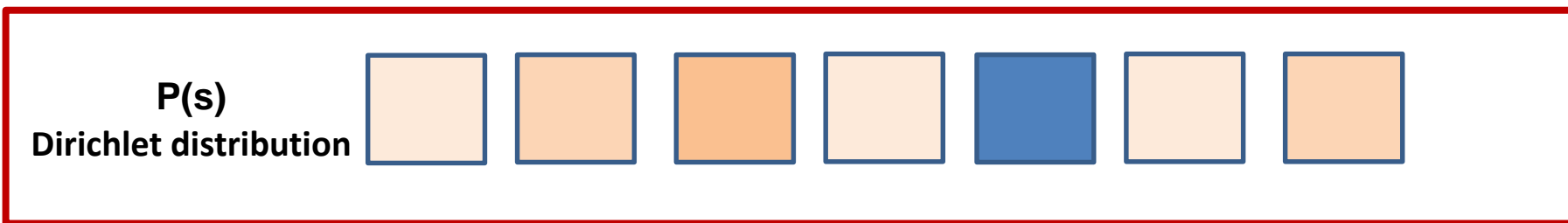
### Dirichlet Distribution(디리클레 분포)

Dirichlet distributions

$$f(x_1, \dots, x_k; \alpha_1, \dots, \alpha_k) = \frac{1}{B(\alpha)} \prod_{i=1}^k x_i^{\alpha_i-1}$$
$$B(\alpha) = \frac{\prod_{i=1}^k \Gamma(\alpha_i)}{\Gamma(\sum_{i=1}^k \alpha_i)} \quad (\Gamma \text{ is gamma function})$$

영제야 이게  
베타 분포의 관점에서 볼 때, 다항분포를 Control 하는 효과를 가진  
분포라는데 무슨 뜻인지 모르겠어 시바

### ZSKD



Softmax output  $s^k$  of class  $k$

$$Dir(K, \alpha^k),$$

$k \in 1 \dots K$  is the class index,

$K$  is the dimension of the output probability vector

$\alpha^k$  is the concentration parameter  $\alpha^k = [\alpha_1^k, \alpha_2^k, \dots, \alpha_K^k]$

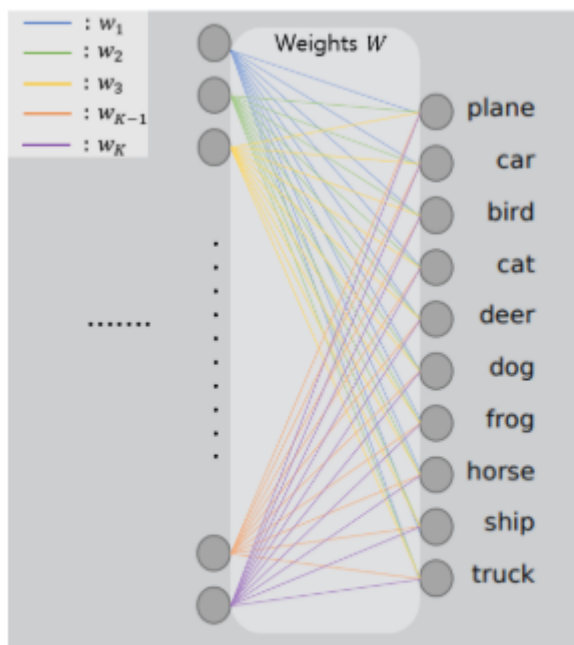
of the distribution modelling class  $k$  It should reflect the similarities across softmax vector

## ZSKD- Class similarity Matrix

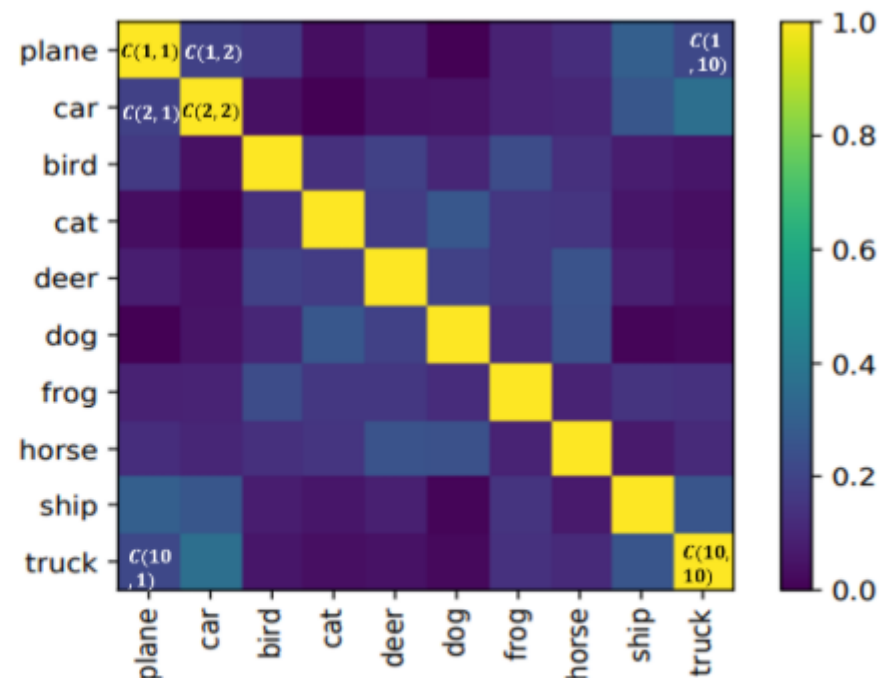
Matrix is consist of the weights connecting the final (softmax)  $W$

$C(i,j)$  denotes  
the visual similarity between  
the categories  $i$  and  $j$  in  $[0,1]$ .

$$C(i,j) = \frac{w_i^T w_j}{\|w_i\| \|w_j\|}$$



(a) Representing weights  
in the final and pre-final layers



(b) Class similarity matrix computed for the  
Teacher model trained over CIFAR-10 dataset

## ZSKD Process

$$\bar{x}_i^k = \underset{x}{\operatorname{argmin}} L_{CE}(\mathbf{y}_i^k, T(x, \theta_T, \tau))$$

N softmax vectors corresponding to class k sampled from **Dir(K,  $\alpha^k$ )** distribution.

Knowledge Distill with S

$$\theta_S = \underset{\theta_S}{\operatorname{argmin}} \sum_{\bar{x} \in \bar{X}} L_{KD}(T(\bar{x}, \theta_T, \tau), S(\bar{x}, \theta_S, \tau))$$

---

### Algorithm 1 Zero-Shot Knowledge Distillation

---

**Input :** *Teacher* model  $T$

$N$ : number of DIs crafted per category,

$[\beta_1, \beta_2, \dots, \beta_B]$ :  $B$  scaling factors,

$\tau$ : Temperature for distillation

**Output :** Learned *Student* model  $S(\theta_S)$ ,

$\bar{X}$ : *Data Impressions*

```

1 Obtain  $K$ : number of categories from  $T$ 
2 Compute the class similarity matrix
    $C = [\mathbf{c}_1^T, \mathbf{c}_2^T, \dots, \mathbf{c}_K^T]$  as in eq. (2)
3  $\bar{X} \leftarrow \emptyset$ 
4 for  $k=1:K$  do
5   Set the concentration parameter  $\alpha^k = \mathbf{c}_k$ 
6   for  $b=1:B$  do
7     for  $n=1:\lfloor N/B \rfloor$  do
8       Sample  $\mathbf{y}_n^k \sim \text{Dir}(K, \beta_b \times \alpha^k)$ 
9       Initialize  $\bar{x}_n^k$  to random noise and craft  $\bar{x}_n^k =$ 
          $\underset{x}{\operatorname{argmin}} L_{CE}(\mathbf{y}_n^k, T(x, \theta_T, \tau))$ 
10       $\bar{X} \leftarrow \bar{X} \cup \bar{x}_n^k$ 
11    end
12  end
13 end
14 Transfer the Teacher's knowledge to Student using the DIs
    via  $\theta_S = \underset{\theta_S}{\operatorname{argmin}} \sum_{\bar{x} \in \bar{X}} L_{KD}(T(\bar{x}, \theta_T, \tau), S(\bar{x}, \theta_S, \tau))$ 

```

## ZSKD Process

- Step 1 : Train the Teacher network with cifar 10
- Step 2 : Extract final layer weights from the Pretrained Teacher Network
- Step 3 : Compute and save the Class Similarity for scales of 1.0 and 0.1
- Step 4 : Generate the Data Impressions (DI's)
- Step 5 : Train the Student network with generated DI's

---

### Algorithm 1 Zero-Shot Knowledge Distillation

---

**Input :** *Teacher* model  $T$

$N$ : number of DIs crafted per category,

$[\beta_1, \beta_2, \dots, \beta_B]$ :  $B$  scaling factors,

$\tau$ : Temperature for distillation

**Output :** Learned *Student* model  $S(\theta_S)$ ,

$\bar{X}$ : *Data Impressions*

```

1 Obtain  $K$ : number of categories from  $T$ 
2 Compute the class similarity matrix
    $C = [\mathbf{c}_1^T, \mathbf{c}_2^T, \dots, \mathbf{c}_K^T]$  as in eq. (2)
3  $\bar{X} \leftarrow \emptyset$ 
4 for  $k=1:K$  do
5   Set the concentration parameter  $\alpha^k = \mathbf{c}_k$ 
6   for  $b=1:B$  do
7     for  $n=1:\lfloor N/B \rfloor$  do
8       Sample  $\mathbf{y}_n^k \sim \text{Dir}(K, \beta_b \times \alpha^k)$ 
9       Initialize  $\bar{x}_n^k$  to random noise and craft  $\bar{x}_n^k =$ 
          $\underset{x}{\operatorname{argmin}} L_{CE}(\mathbf{y}_n^k, T(x, \theta_T, \tau))$ 
10       $\bar{X} \leftarrow \bar{X} \cup \bar{x}_n^k$ 
11    end
12  end
13 end
14 Transfer the Teacher's knowledge to Student using the DIs
   via  $\theta_S = \underset{\theta_S}{\operatorname{argmin}} \sum_{\bar{x} \in \bar{X}} L_{KD}(T(\bar{x}, \theta_T, \tau), S(\bar{x}, \theta_S, \tau))$ 
```

# 3 MIL Experiments

---



Experiement.pdf

---

감사합니다  
죄송합니다.