

UST seminar

Towards domain-agnostic Video action recognition

Hyungmin Kim

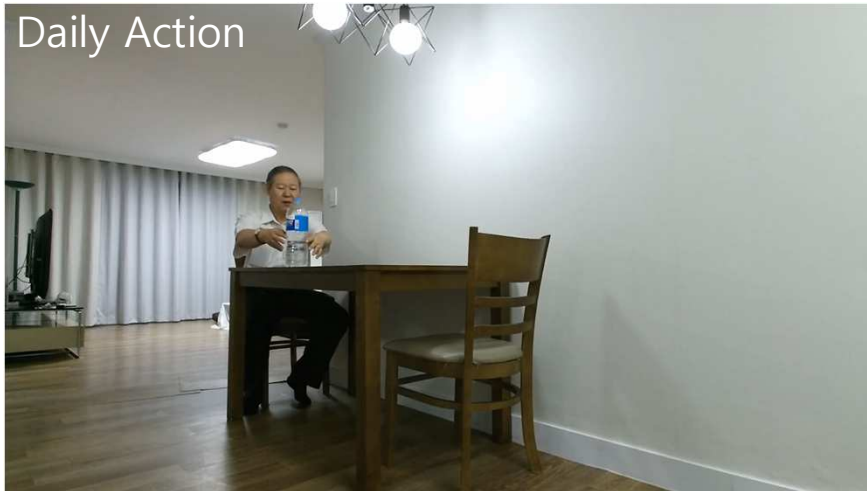
UST-ETRI

Ph.D. student

khm159@etri.re.kr/ust.ac.kr

21th Sep, 2023

Human action recognition in Human-Robot Interaction



Video from [1]

- Problem definition :

Given input video v_i

Model predicts corresponding action label y_i

Video
acquisition



"Drinking something"

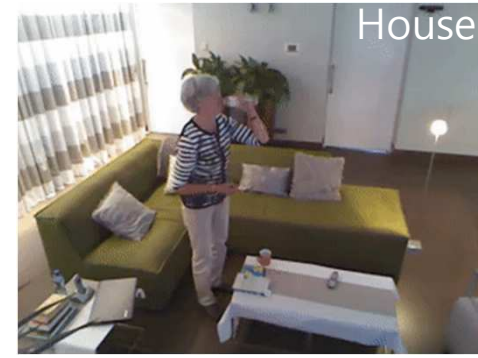
Service robot encounters numerous domains



Video from [1]



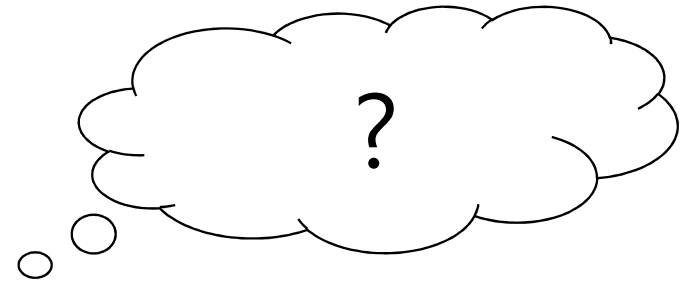
Video from [2]



Video from [3]

- Care-robot encounters numerous domains

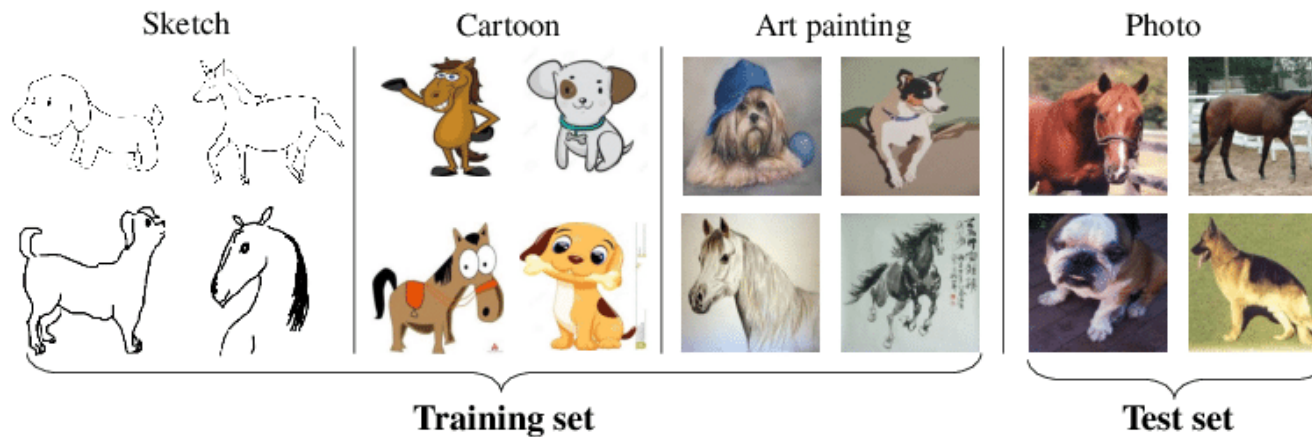
Apartment, House, Building, Office



- [1] Jang, Jinhyeok, et al. "ETRI-activity3D: A large-scale RGB-D dataset for robots to recognize daily activities of the elderly." *IROS*, 2020.
[2] Shahroury, Amir, et al. "Ntu rgb+ d: A large scale dataset for 3d human activity analysis." *CVPR*. 2016.
[3] Das, Srijan, et al. "Toyota smarthome: Real-world activities of daily living." *ICCV*. 2019.

What is Domain Generalization?

Domain Generalization?



Train a generalized model with datasets from different domains that share the same semantic to improve performance on Un-seen domain.

Motivation

- Why DG is required?

Most of the existing deep learning/machine learning tasks operated under **i.i.d. assumption**. However, this is an impossible in real-world.

ICML21

Consider a classification task where the learning algorithm has access to i.i.d. data from m domains, $\{(d_i, \mathbf{x}_i, y_i)\}_{i=1}^n \sim (D_m, \mathcal{X}, \mathcal{Y})^n$ where $d_i \in D_m$ and $D_m \subset \mathcal{D}$ is a set of m domains. Each training in-

AAAI-22

In most statistical machine learning algorithms, a fundamental assumption is that the training data and test data are independently and identically distributed (i.i.d.). However, the data we have in many real-world applications are not i.i.d. Distributional shifts are ubiquitous. Under such

arXiv20

Convolutional Neural Networks (CNNs) show impressive performance in the standard classification setting where training and testing data are drawn i.i.d. from a given domain. However, CNNs do not readily generalize to new domains with different statistics, a setting that is simple for humans. In this work, we address

NeuralPS21

Independent and identically distributed (i.i.d.) condition is the underlying assumption of machine learning experiments. However, this assumption may not hold in real-world scenarios, i.e., the training and the test data distribution may differ significantly by distribution shifts. For example, a

ICLR21

success of CNNs heavily relies on the i.i.d. assumption, i.e. training and test data should be drawn from the same distribution; when such an assumption is violated even just slightly, as in most real-world application scenarios, severe performance degradation is expected (Hendrycks & Dietterich,

ICML13

Since in general $\mathbb{P}_{XY}^i \neq \mathbb{P}_{XY}^j$, the samples in \mathcal{S} are not i.i.d. Let $\hat{\mathbb{P}}^i$ denote empirical distribution associated

NeurIPS18

discrete set $\{1, 2, \dots, N_c\}$, where N_c denotes the number of classes. Let $\{\mathcal{D}_i\}_{i=1}^{p+q}$ represent the $p+q$ distributions, each of which exists on the joint space $\mathcal{X} \times \mathcal{Y}$. Let $D_i = \{(\mathbf{x}_j^{(i)}, y_j^{(i)})\}_{j=1}^{N_i}$ represent the dataset sampled from the i^{th} distribution, i.e. each $(x_j^i, y_j^i) \stackrel{i.i.d.}{\sim} \mathcal{D}_i$. In the rest of the paper,

NeurIPS21

Problem 3.1 (Domain generalization). Let $\mathcal{E}_{\text{train}} \subsetneq \mathcal{E}_{\text{all}}$ be a finite subset of training domains, and assume that for each $e \in \mathcal{E}_{\text{train}}$, we have access to a dataset $\mathcal{D}^e := \{(x_j^e, y_j^e)\}_{j=1}^{n_e}$ sampled i.i.d. from $\mathbb{P}(X^e, Y^e)$. Given a function class \mathcal{F} and a loss function $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_{>0}$, our goal is to learn

ECCV20

We expand our notation set for the theoretical analysis. As we study the domain-agnostic cross-domain setting, we no longer work with i.i.d data. There-

IEEE Trans. On Knowledge and Data Engineering
to new (test) data. Traditional ML models are trained based on the i.i.d. assumption that training and testing data are identically and independently distributed. However, this assumption does not always hold in reality. When the prob-


What is the i.i.d. assumption?

Independently and Identically Distributed Assumption

(1) **Independent** (Observations that acquire each data point are independent of each other)

Random samples from an unknown joint probability distribution (A.K.A. training dataset)

Each event are
independent
each other



Sample num	Feature 1	Feature 2	...	Feature N
1	1.24	True
2
3
...
M-1
M	0.11	False

What is the i.i.d. assumption?

Independently and Identically Distributed Assumption

(1) **Independent** (Observations that acquire each data point are independent of each other)

Random samples from an unknown joint probability distribution (A.K.A. training dataset)

Each event are
independent
each other

	Sample num	Feature 1	Feature 2	...	Feature N
	1	1.24	True
	2
	3

	M-1
	M	0.11	False

Why independent each other?



What is the i.i.d. assumption?

Independently and Identically Distributed Assumption

(1) **Independent** (Observations that acquire each data point are independent of each other)

Purpose :

→ **To calculate Likelihood easily**

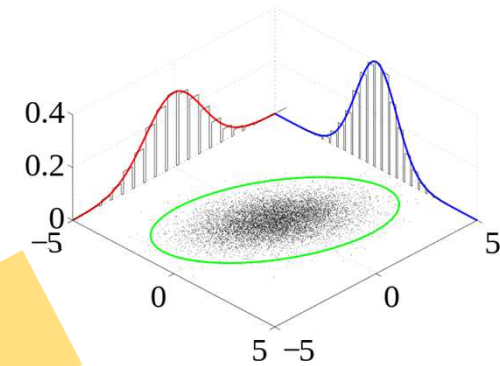


Likelihood?

Preliminary : Likelihood?

Likelihood :

The **probability** that the data is **derived from a particular probability distribution**



Unknown joint
probability distribution

Sample num	Feature 1	Feature 2	...	Feature N
1	1.24	True
2
3
...
M-1
M	0.11	False

sampling

Preliminary : Likelihood?

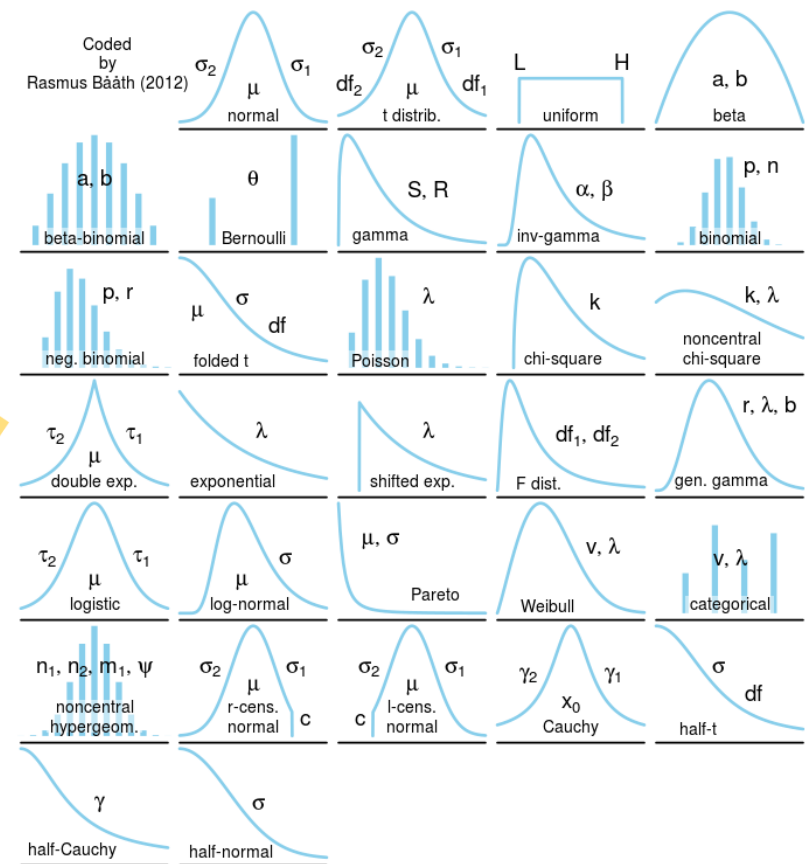
Likelihood :

Let there are the data : $[1, 1, 1, 1]$

To **which probability distribution** are the **given data points**
most likely to belong?

$[1, 1, 1, 1]$

compare



Preliminary : Likelihood?

Likelihood :

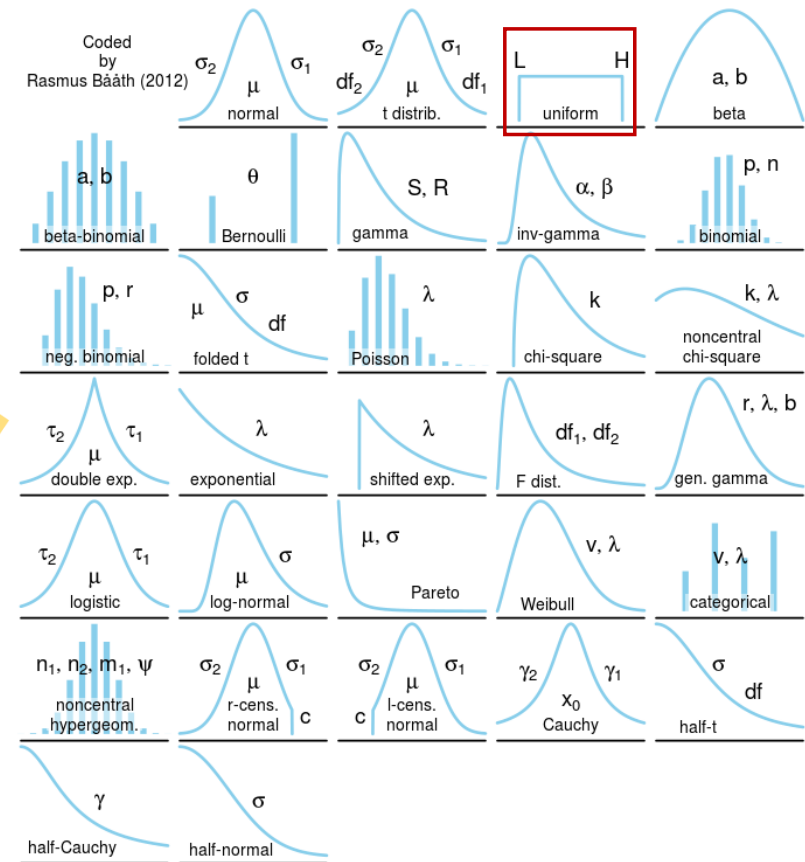
Let there are the data : $[1, 1, 1, 1]$

To **which probability distribution** are the **given data points**
most likely to belong?

$[1, 1, 1, 1]$

compare

Uniform distribution



Well defined distributions

Preliminary : Likelihood?

Likelihood :

We basically **parameterize the unknown distribution** with parameters $\boldsymbol{\theta} = [\theta_1, \theta_2, \dots, \theta_k]^T$

Preliminary : Likelihood?

Likelihood :

We basically **parameterize the unknown distribution** with parameters $\boldsymbol{\theta} = [\theta_1, \theta_2, \dots, \theta_k]^T$

For example gaussian ...

$$g(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{(-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2})} \approx \exp(\alpha x^2 + \beta x + \gamma)$$

Preliminary : Likelihood?

Likelihood :

We basically **parameterize the unknown distribution** with parameters $\theta = [\theta_1, \theta_2, \dots, \theta_k]^T$

For example gaussian ...

$$g(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{(-\frac{1(x-\mu)^2}{2\sigma^2})} \approx \exp(\alpha x^2 + \beta x + \gamma)$$

In this case.... The **Maximum Likelihood Estimation(MLE)** is to find optimal α, β, γ

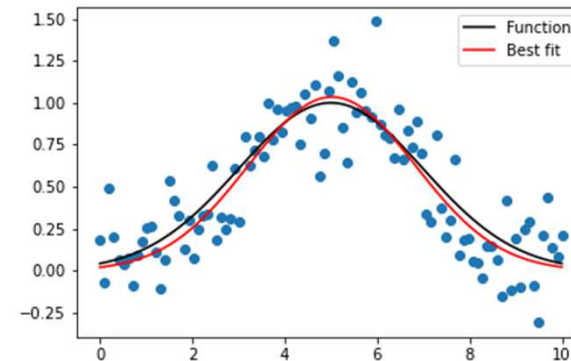
Preliminary : Likelihood?

Likelihood :

We basically **parameterize the unknown distribution** with parameters $\theta = [\theta_1, \theta_2, \dots, \theta_k]^T$

For example gaussian ...

$$g(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{(-\frac{1(x-\mu)^2}{2\sigma^2})} \approx \exp(\alpha x^2 + \beta x + \gamma)$$



In this case.... The **Maximum Likelihood Estimation(MLE)** is to find optimal α, β, γ

To find **mostly likely belonging probability distribution** when given data points...

Preliminary : Likelihood?

Likelihood :

We basically **parameterize the unknown distribution** with parameters $\boldsymbol{\theta} = [\theta_1, \theta_2, \dots, \theta_k]^T$

Now we can define the likelihood with the given parameter $\boldsymbol{\theta}$ and dataset X like this...

$$\text{Likelihood}(\boldsymbol{\theta}) = p(X|\boldsymbol{\theta})$$

Preliminary : Likelihood?

Likelihood :

We basically **parameterize the unknown distribution** with parameters $\theta = [\theta_1, \theta_2, \dots, \theta_k]^T$

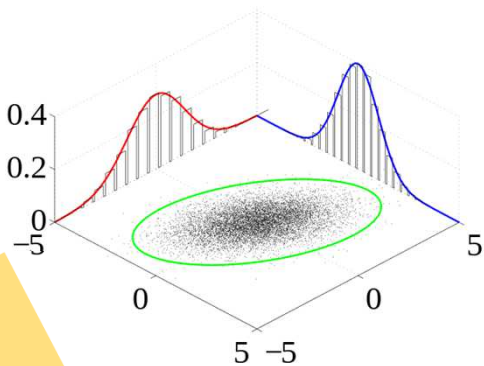
Now we can define the likelihood with the given parameter θ and dataset X like this...

$$\text{Likelihood}(\theta) = p(X|\theta)$$

If $p(X|\theta)$ is high, then we say that the estimated θ **effectively explain Unknown (original) distribution of X**

So.. Now we know the likelihood

We can model transfer function with some parameters θ
The probability that the data is **derived from a particular probability distribution**



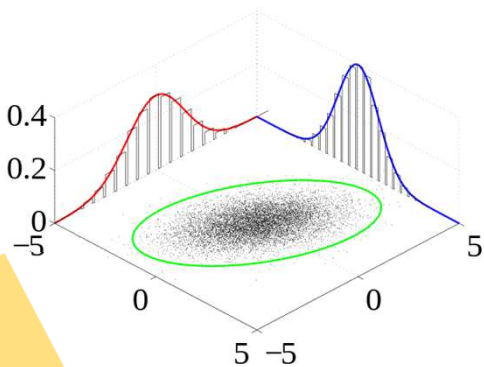
Unknown joint probability distribution

Sample num	Feature 1	Feature 2	...	Feature N
1	1.24	True
2
3
...
M-1
M	0.11	False

sampling

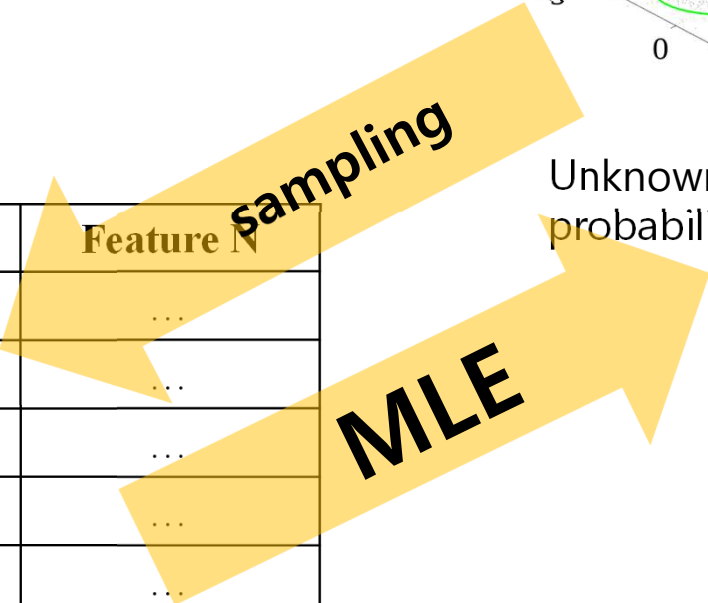
So.. Now we know the likelihood

We can model transfer function with some parameters θ
The **probability** that the data is **derived from a particular probability distribution**
We **estimate** θ when **likelihood is maximized**.



Unknown joint probability distribution

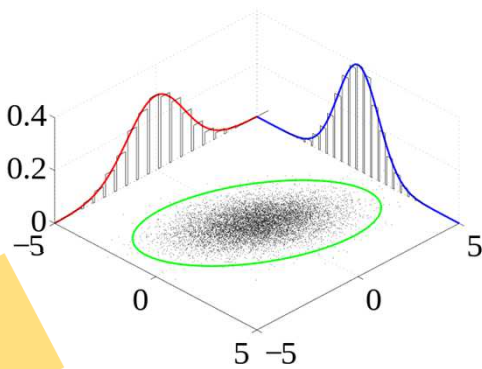
Sample num	Feature 1	Feature 2	...	Feature N
1	1.24	True
2
3
...
M-1
M	0.11	False



Remember, we need to find the model from dataset

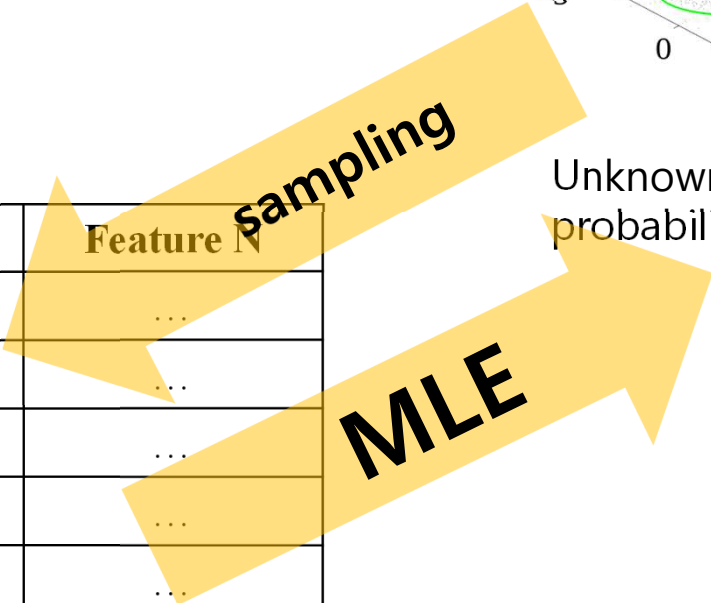
So.. Now we know the likelihood

We can model transfer function with some parameters θ
The probability that the data is derived from a particular probability distribution
We estimate θ when likelihood is maximized.



Unknown joint probability distribution

Sample num	Feature 1	Feature 2	...	Feature N
1	1.24	True
2
3
...
M-1
M	0.11	False



$Likelihood(\theta) = p(X|\theta) \uparrow$

What is the i.i.d. assumption?

Independently and Identically Distributed Assumption

(1) **Independent** (Observations that acquire each data point are independent of each other)

Purpose :

→ **To calculate Likelihood easily**

it can be expressed **in the form of a product**, so the complexity of the operation can be minimized by adding log.



I don't get it...

What is the i.i.d. assumption?

Independently and Identically Distributed Assumption

(1) **Independent** (Observations that acquire each data point are independent of each other)

Purpose :

→ **To calculate Likelihood easily**

it can be expressed **in the form of a product**, so the complexity of the operation can be minimized by adding log.

Likelihood

$$l(\theta) = \underbrace{P(x_1, x_2, x_3, \dots, x_n | \theta)}_{\text{INDEPENDENT!}} = \prod_{n=1}^N P(x_n | \theta)$$



What is the i.i.d. assumption?

Independently and Identically Distributed Assumption

(1) **Independent** (Observations that acquire each data point are independent of each other)


Purpose :

→ **To calculate Likelihood easily**

it can be expressed **in the form of a product**, so the complexity of the operation can be minimized by adding log.

Likelihood

$$\begin{aligned} l(\theta) &= P(x_1, x_2, x_3, \dots, x_n | \theta) = \overbrace{P(x_1 | \theta) P(x_2 | \theta) P(x_3 | \theta) \dots P(x_n | \theta)}^{\text{INDEPENDENT!}} = \prod_{n=1}^N P(x_n | \theta) \\ &= \log \prod_{n=1}^N P(x_n | \theta) \\ &= \sum_{n=1}^N \log P(x_n | \theta) \end{aligned}$$

SUMMATION! 

What is the i.i.d. assumption?

Independently and Identically Distributed Assumption

If you are interested in more **relationship between MLE and cross-entropy or MSE loss...**

Then I recommend following materials.. (because of the presentation time limit I'm not explain this)

I was able to get a good insight from the materials below.

- (1) [Miranda and Lester James, "Understanding softmax and the negative log-likelihood", 2017.](#)
- (2) [Hwalseok Lee, "Everything about autoencoder", 2018 \(Korean\)](#)
- (3) [Doersch, Carl. "Tutorial on variational autoencoders." *arXiv preprint*, 2016\).](#)
- (4) [Scaling Up Deep Learning - Yoshua Bengio, *ICML*, 2014.](#)

What is the i.i.d. assumption?

Independently and Identically Distributed Assumption

If not independent?

$$l(\theta) = P(x_1, x_2, x_3, \dots, x_n | \theta)$$

The joint probability are **extremely hard** to define and calculate....



Remember.. No dependent....

What is the i.i.d. assumption?

Independently and Identically Distributed Assumption

(2) **Identically distributed**

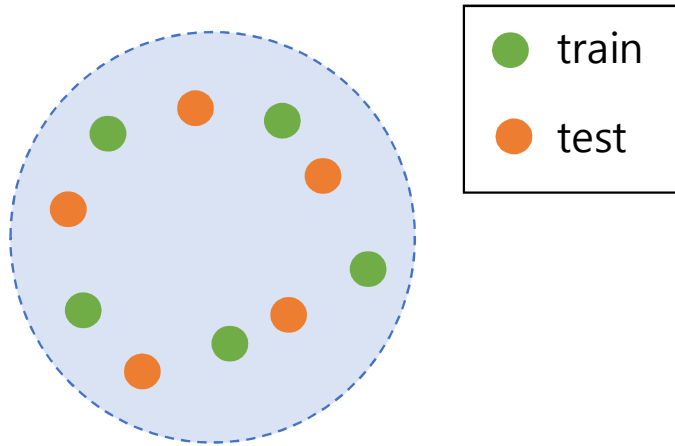
There are no overall trends the distribution **doesn't fluctuate** and all items in the sample are **taken from the same probability distribution.**

What is the i.i.d. assumption?

Independently and Identically Distributed Assumption

(2) **Identically distributed**

There are no overall trends the distribution **doesn't fluctuate** and all items in the sample are **taken from the same probability distribution**.



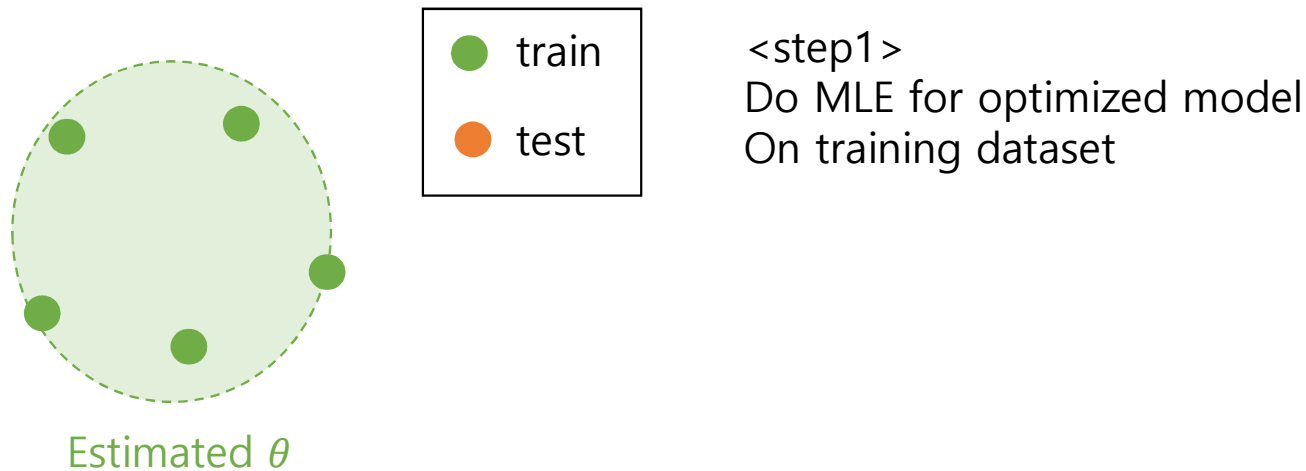
Ex. gaussian

What is the i.i.d. assumption?

Independently and Identically Distributed Assumption

(2) **Identically distributed**

There are no overall trends the distribution **doesn't fluctuate** and all items in the sample are **taken from the same probability distribution**.

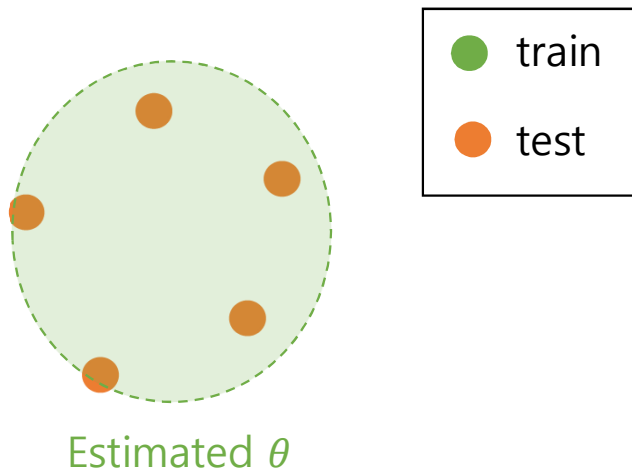


What is the i.i.d. assumption?

Independently and Identically Distributed Assumption

(2) **Identically distributed**

There are no overall trends the distribution **doesn't fluctuate** and all items in the sample are **taken from the same probability distribution**.



<step2>

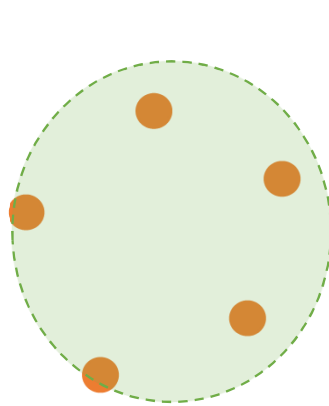
The optimized model can effectively explain Testing dataset **if train and test dataset are identically distributed**

What is the i.i.d. assumption?

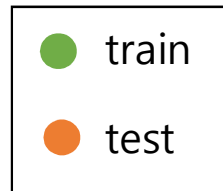
Independently and Identically Distributed Assumption

(2) **Identically distributed**

There are no overall trends the distribution **doesn't fluctuate** and all items in the sample are **taken from the same probability distribution**.



Estimated θ



<step2>

The optimized model can effectively explain Testing dataset **if train and test dataset are identically distributed**



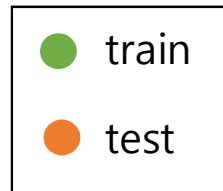
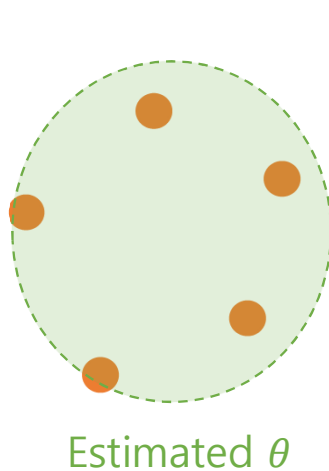
Is it possible?

What is the i.i.d. assumption?

Independently and Identically Distributed Assumption

(2) **Identically distributed**

There are no overall trends the distribution **doesn't fluctuate** and all items in the sample are **taken from the same probability distribution**.



<step2>

The optimized model can effectively explain Testing dataset **if train and test dataset are identically distributed**



Is it possible?

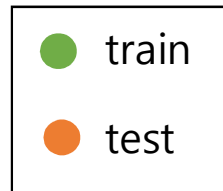
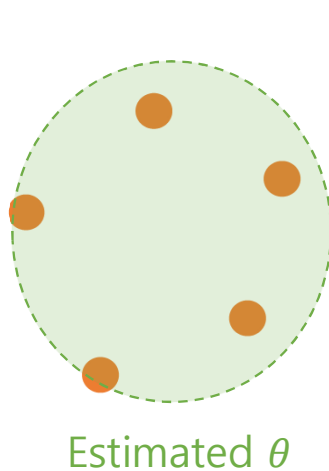
Yes!

What is the i.i.d. assumption?

Independently and Identically Distributed Assumption

(2) **Identically distributed**

There are no overall trends the distribution **doesn't fluctuate** and all items in the sample are **taken from the same probability distribution**.



<step2>

The optimized model can effectively explain Testing dataset **if train and test dataset are identically distributed**



Is it possible?

Yes!

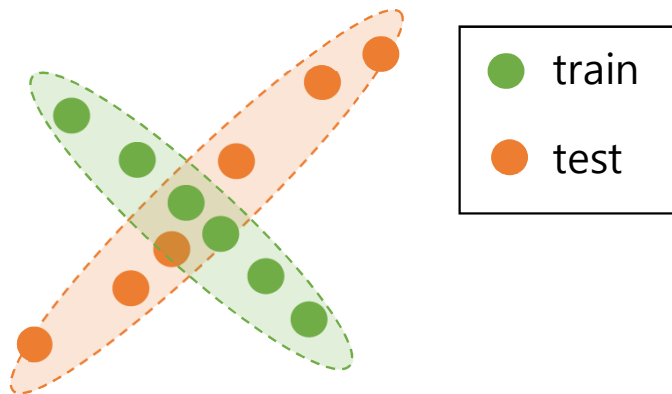
If the data obtained is representative of same population

What is the i.i.d. assumption?

Independently and Identically Distributed Assumption

(2) **Identically distributed**

If not? **The training is useless**



The training set and the test set represent **different distributions.**

Summary

Independently and Identically Distributed Assumption

(1) Independent

To calculate likelihood easily

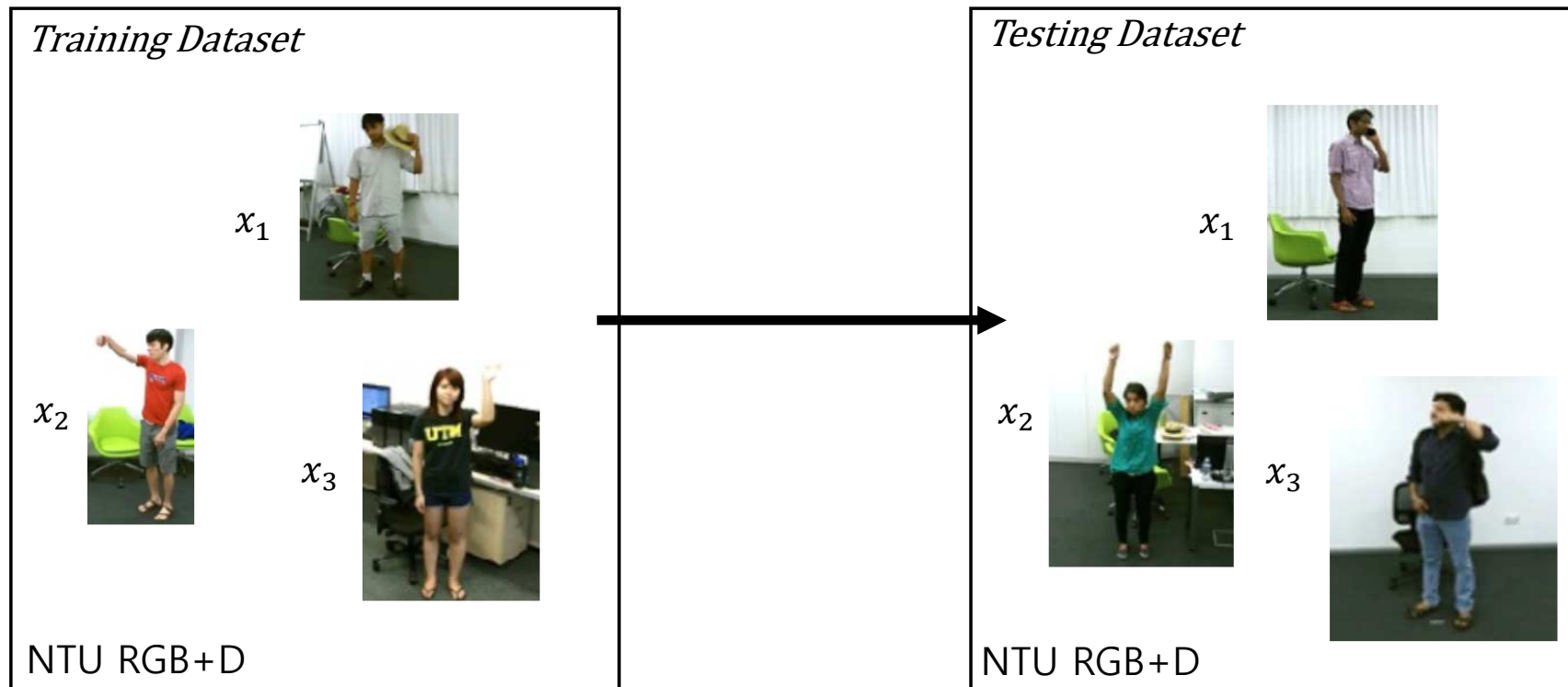
(2) Identically distributed

The population represented by the learning set and the test set is the same.

In i.i.d. assumption...

- Usually, the training set and test set are obtained in the same or similar way. Each sample is Independent Identically distributed.

Similar back ground, similar camera view point, etc...



i.i.d. assumption is easily violated in **real world**

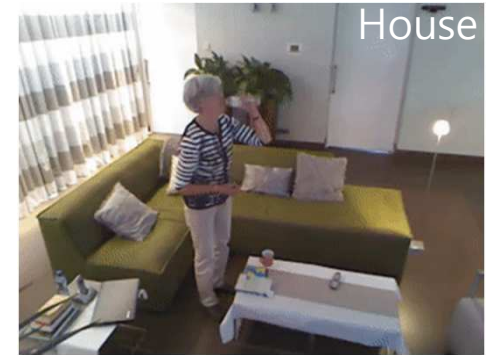
- Care-robot encounters numerous domains

Apartment, House, Building, Office

Inferred domains



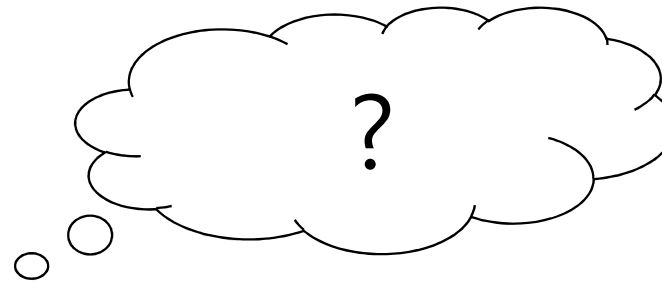
Video from [2]



Video from [3]



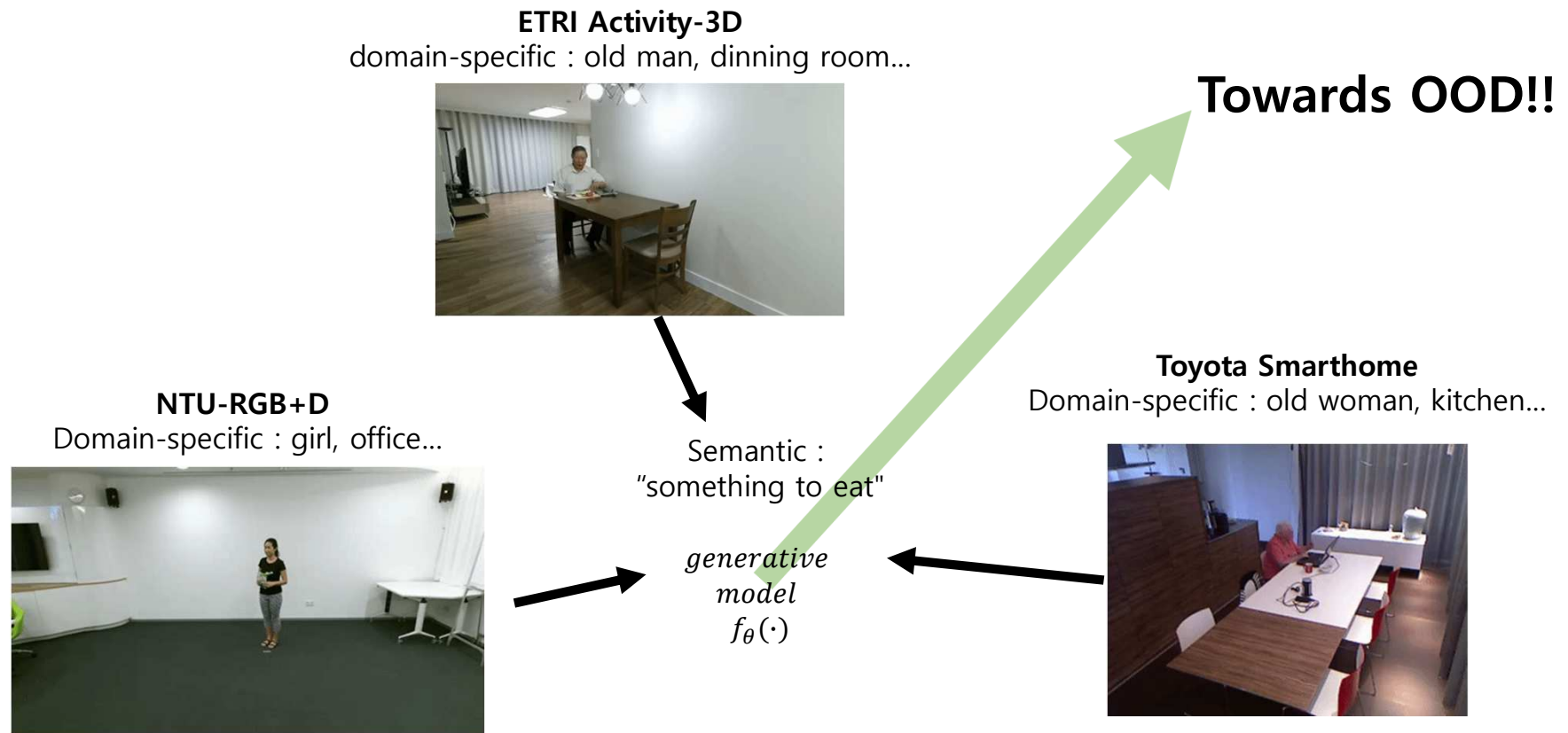
Video from [1]



- [1] Jang, Jinhyeok, et al. "ETRI-activity3D: A large-scale RGB-D dataset for robots to recognize daily activities of the elderly." *IROS*, 2020.
[2] Shahrudy, Amir, et al. "Ntu rgb+ d: A large scale dataset for 3d human activity analysis." *CVPR*. 2016.
[3] Das, Srijan, et al. "Toyota smarthome: Real-world activities of daily living." *ICCV*. 2019.

The goal of DG

Using many training domains well
to train models that are **robust to domain differences**.



Evaluation protocol

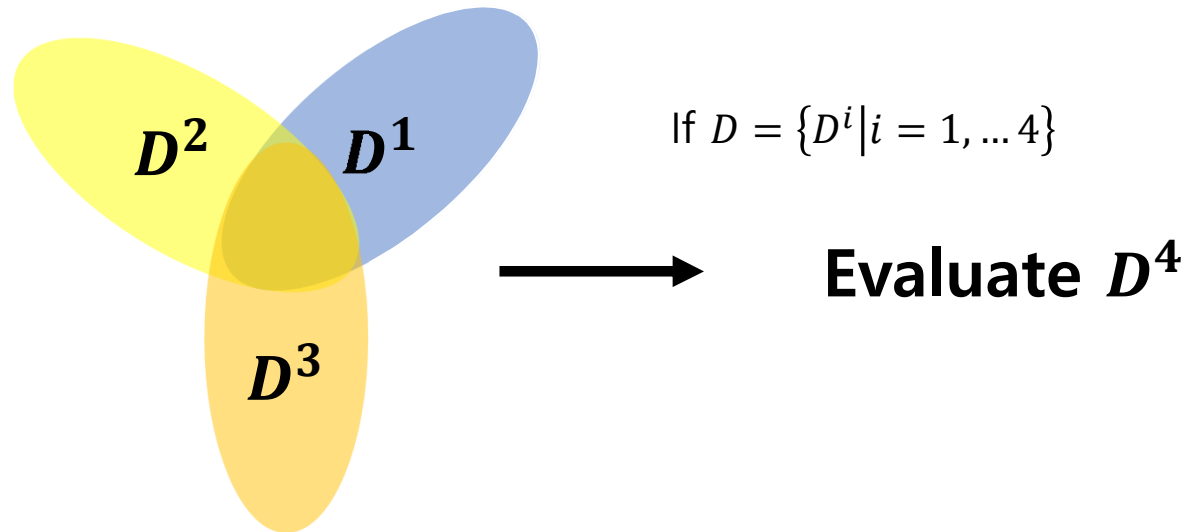
- **Leave-one-domain-out protocol**

Set of domains $D = \{D^i | i = 1, \dots, M\}$

Each domain $D^i = \{(x^{i,j}, y^{i,j})\}_{j=1}^{n_i}$

one domain from D is excluded or left out during training,
while the remaining domains of D are used for training the model.

- n_j : number of samples of D^i
- M : number of domains



DG in video domain

Unfortunately, there are few papers and little research.

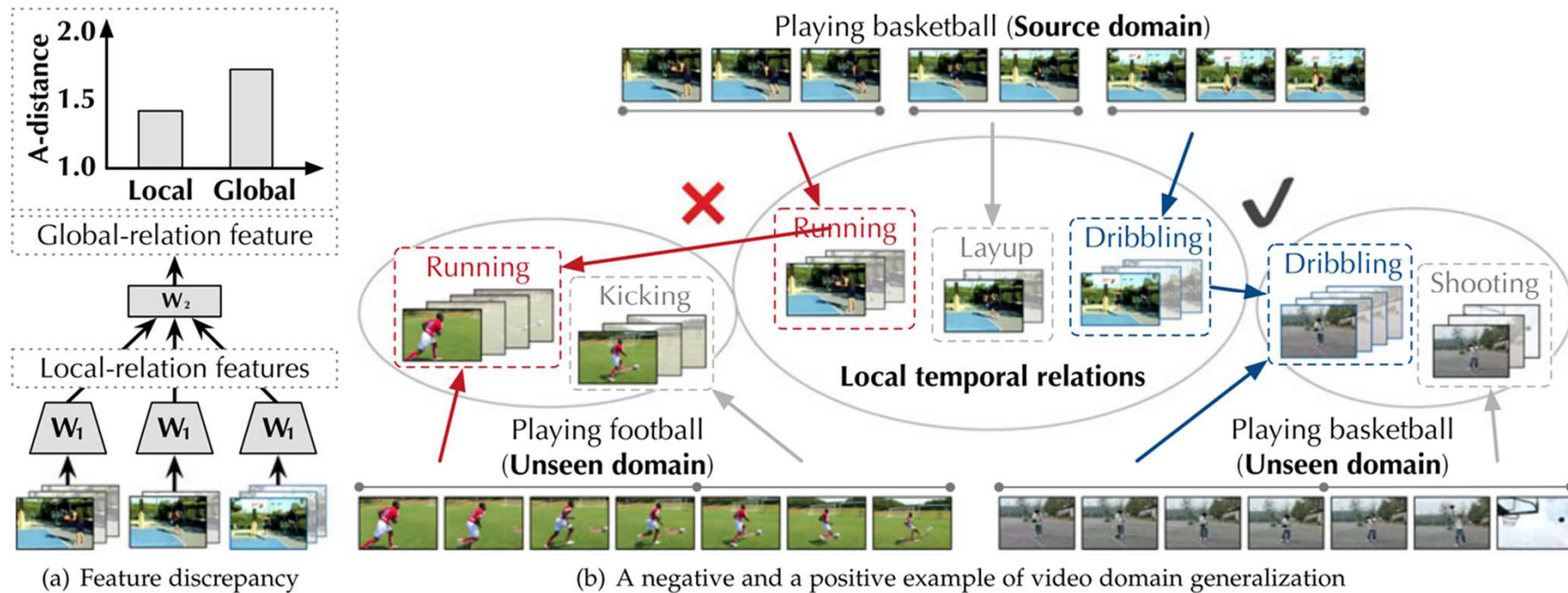
So, I will introduce a pioneer paper.

(best of my knowledge, **in real domain, only two**(and one is my work...)).



VideoDG review

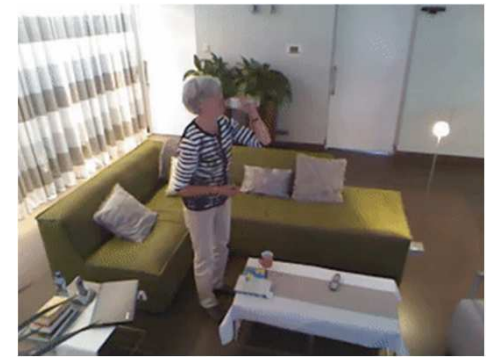
Motivation : **The “temporal” discrepancy** is **co-existed** with “spatial” discrepancy in video domain.



VideoDG review

(1) Spatial domain difference

Caused by the variations of the appearance of video frames
Partially solve by image-based domain generalization methods



Example of the spatial domain difference in video (ENT dataset*)

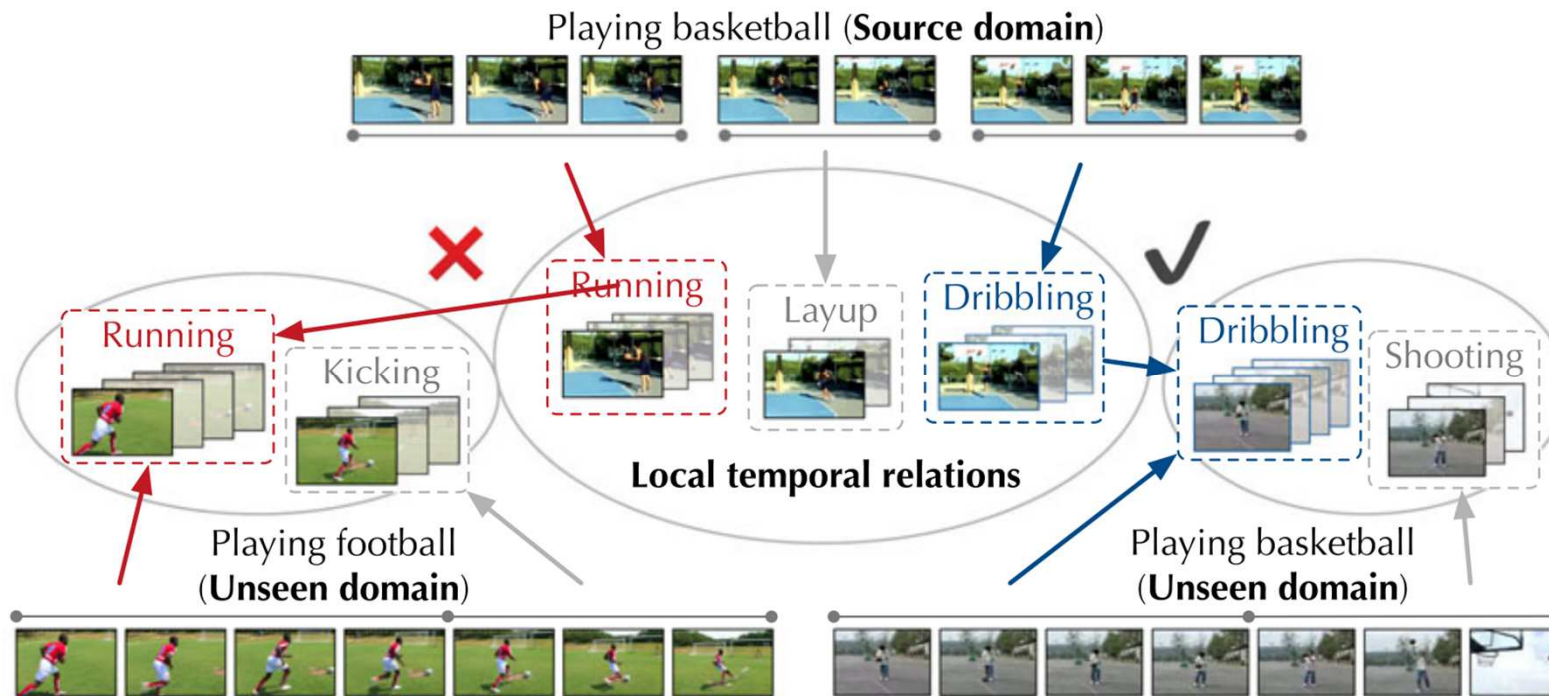
Yao, Zhiyu, et al. "Videodg: Generalizing temporal relations in videos to novel domains." *IEEE TPAMI*, 2021.

*Hyungmin, Kim et al. "A simple baseline for domain generalization of action recognition and a realistic out-of-domain scenario." *International Conference on Ubiquitous Robots (UR)*. IEEE, 2023.

VideoDG review

(2) Temporal domain difference

The unexpected absence or misalignment of short-term video events



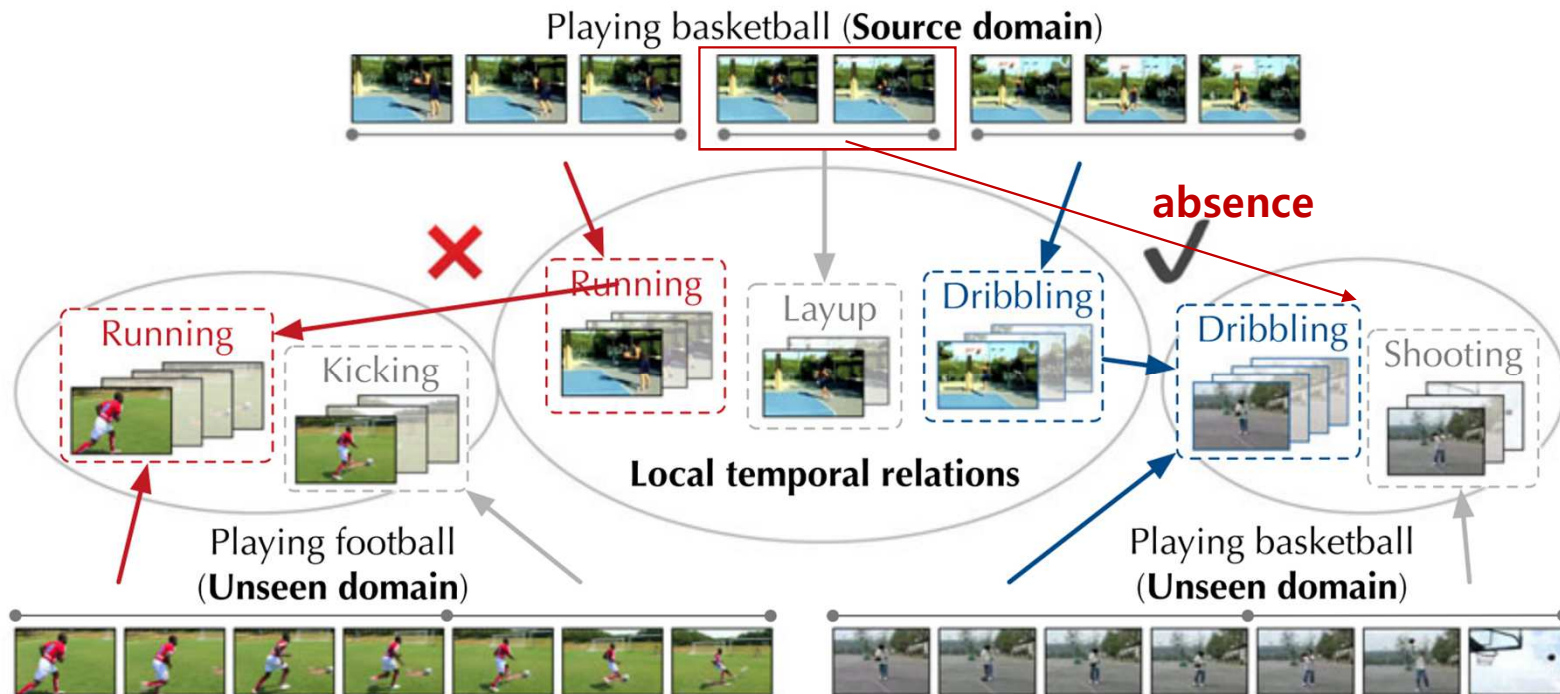
(b) A negative and a positive example of video domain generalization

VideoDG review

(2) Temporal domain difference

Short-term events can be existed or confused
When domain-shift is occurred.

The unexpected **absence** or misalignment of short-term video events



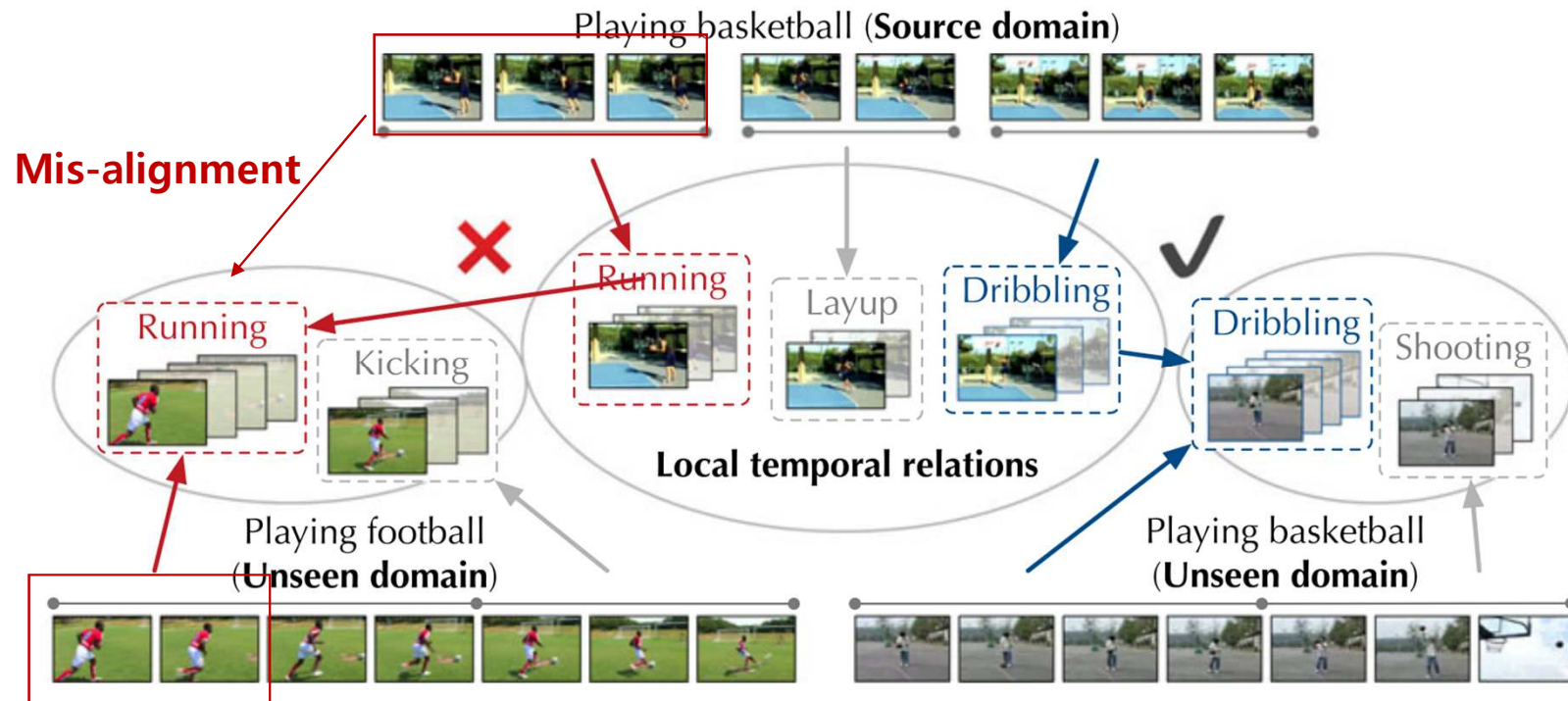
(b) A negative and a positive example of video domain generalization

VideoDG review

(2) Temporal domain difference

Short-term events can be existed or confused
When domain-shift is occurred.

The unexpected absence or **misalignment** of short-term video events

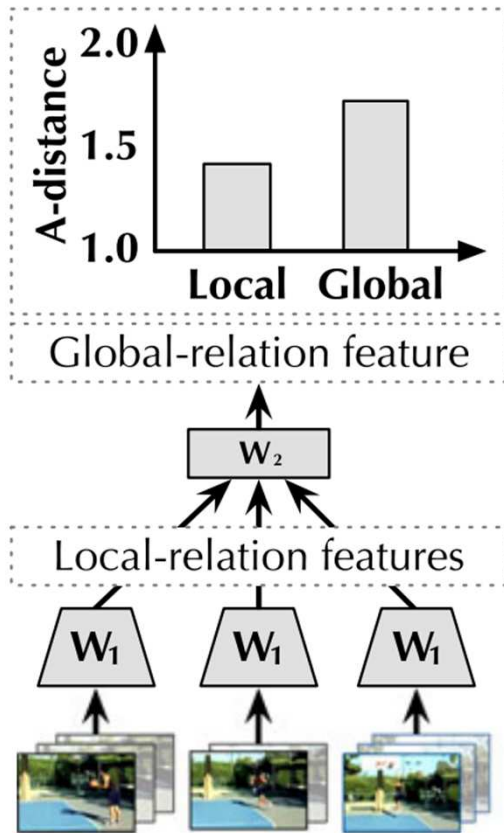


(b) A negative and a positive example of video domain generalization

VideoDG review

In case of UCF→HMDB

Motivation : **The "temporal" discrepancy** is **co-existed** with "spatial" discrepancy in video domain.



A-distance*

If A-distance large : two distribution far apart

If A-distance small : two distribution close to each other

(a) Feature discrepancy

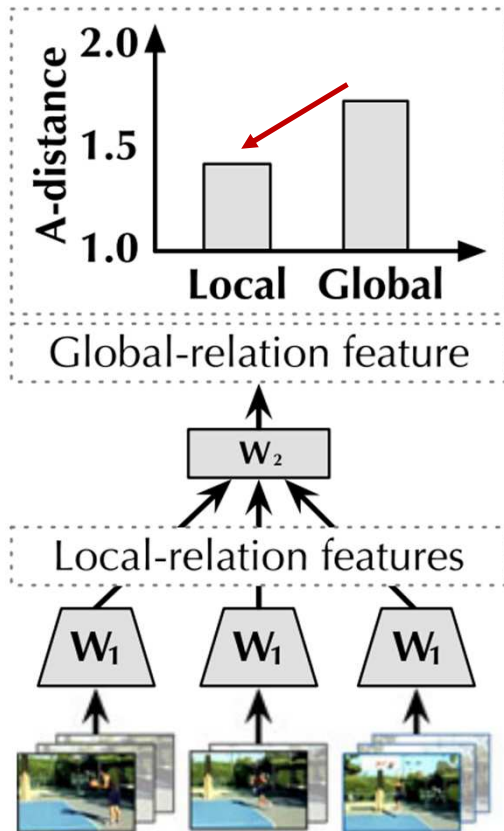
Yao, Zhiyu, et al. "Videodg: Generalizing temporal relations in videos to novel domains." *IEEE TPAMI*, 2021.

*Ben-David, Shai, et al. "Analysis of representations for domain adaptation." *Advances in neural information processing systems* 19 (2006).

VideoDG review

In case of UCF→HMDB

Motivation : **The "temporal" discrepancy** is **co-existed** with "spatial" discrepancy in video domain.



(a) Feature discrepancy

A-distance*

If A-distance large : two distribution far apart

If A-distance small : two distribution close to each other

The local feature are more generalizable then global feature

Yao, Zhiyu, et al. "Videodg: Generalizing temporal relations in videos to novel domains." *IEEE TPAMI*, 2021.

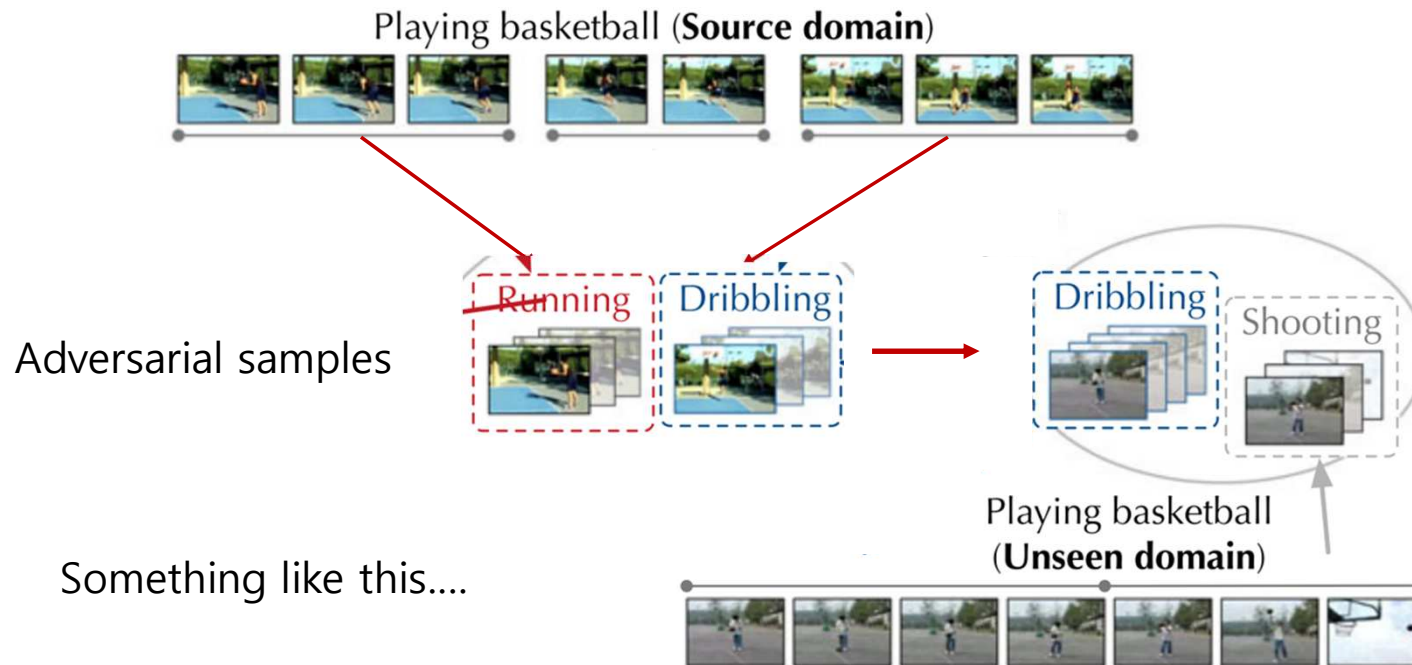
*Ben-David, Shai, et al. "Analysis of representations for domain adaptation." *Advances in neural information processing systems* 19 (2006).

VideoDG review

Main Idea

The local feature are more generalizable then global feature

Using local features **to generate adversarial samples** to augment training dataset



Yao, Zhiyu, et al. "Videodg: Generalizing temporal relations in videos to novel domains." *IEEE TPAMI*, 2021.

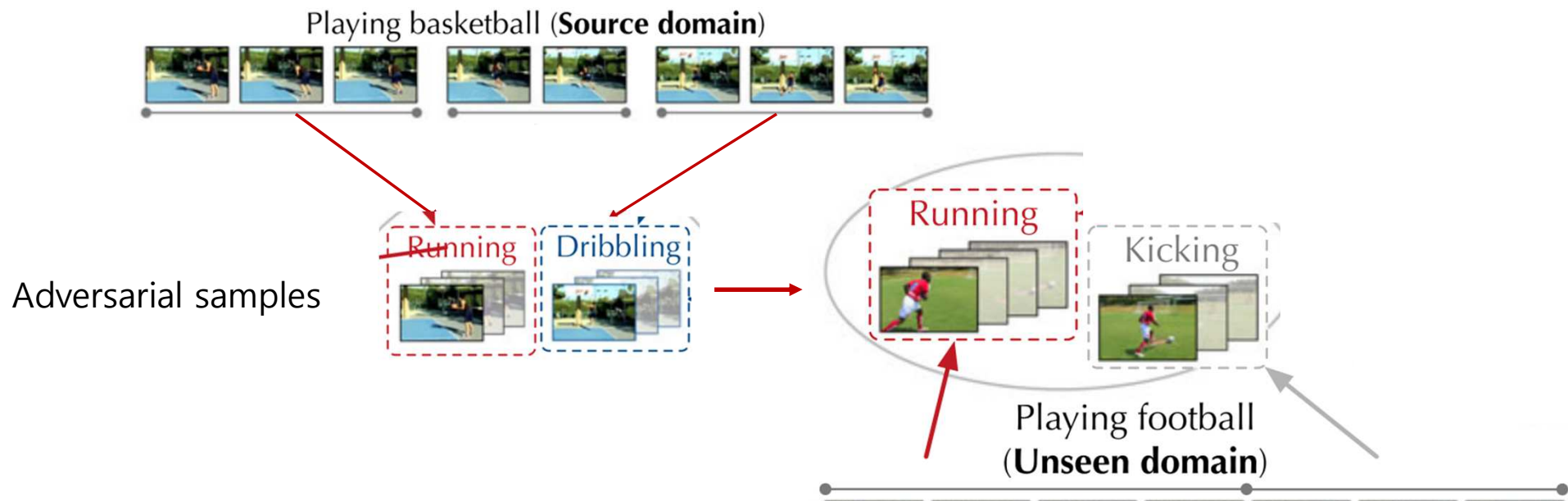
*Ben-David, Shai, et al. "Analysis of representations for domain adaptation." *Advances in neural information processing systems* 19 (2006).

VideoDG review

Main Idea

The local feature are more generalizable then global feature

Using local features **to generate adversarial samples** to augment training dataset



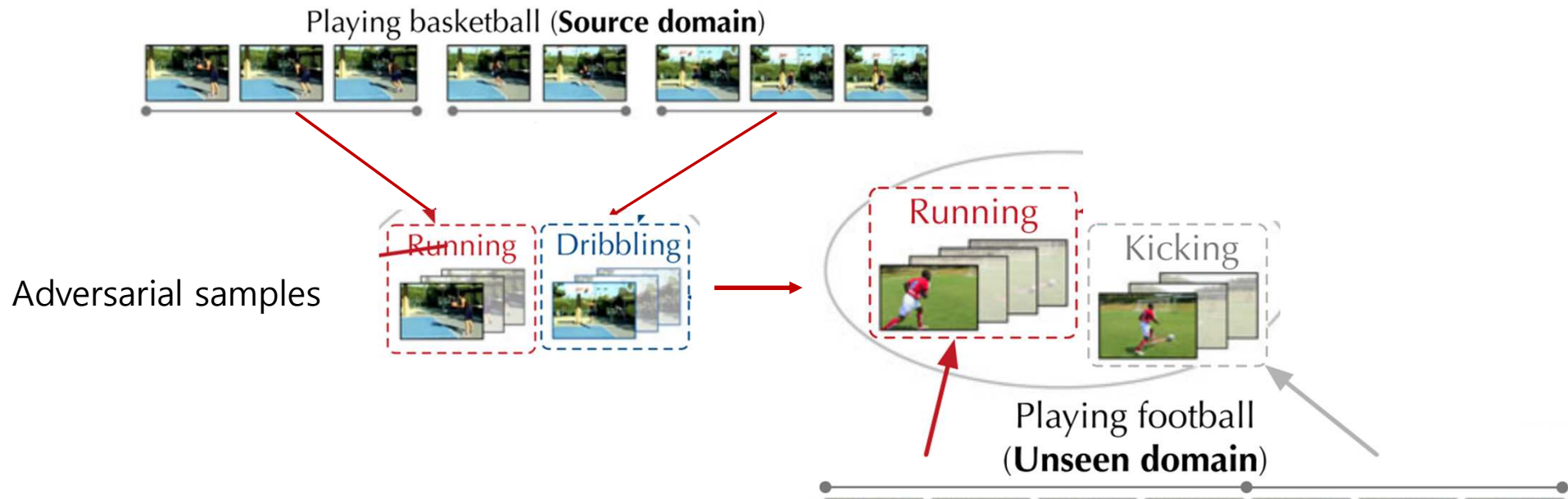
But also **degrade the discriminability** and can miss classify into similar categories

VideoDG review

Main Idea

The local feature are more generalizable then global feature

Using local features **to generate adversarial samples** to augment training dataset



But also degrade the discriminability and can miss classify into similar categories

The global feature are also important

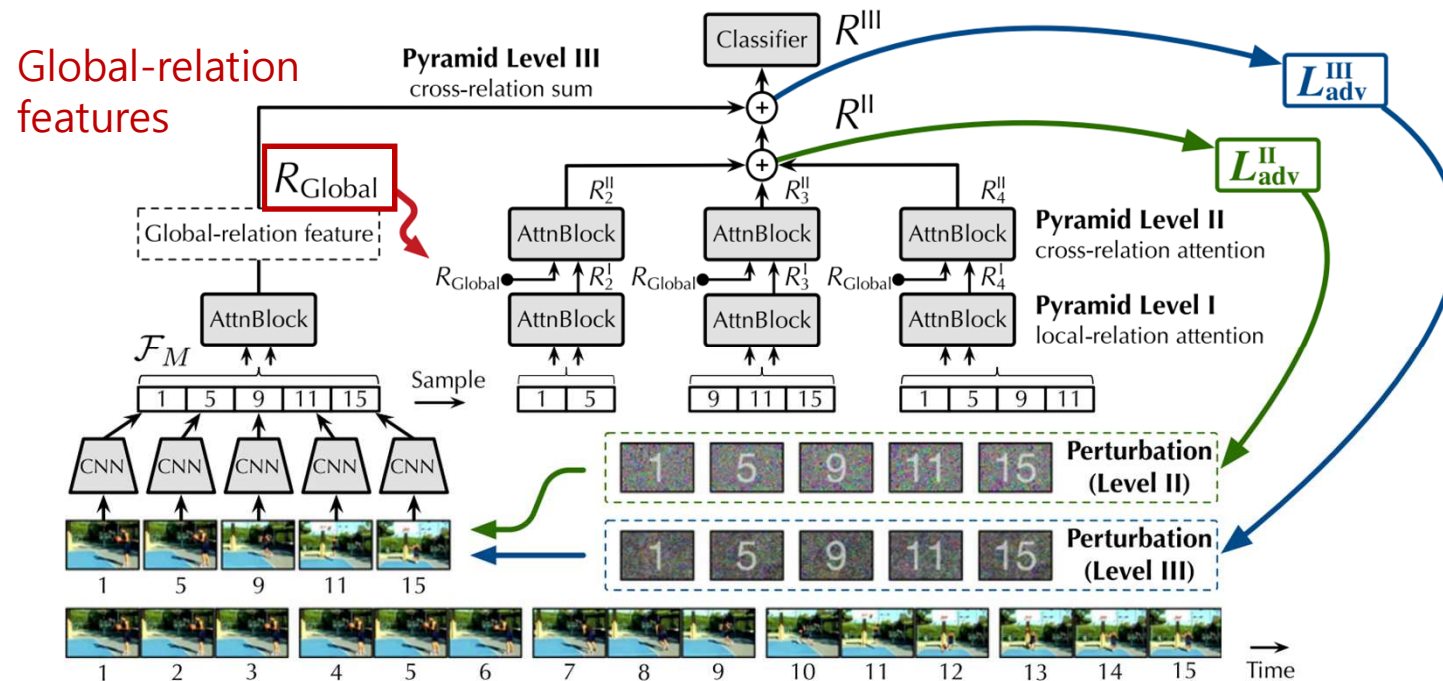
Yao, Zhiyu, et al. "Videodg: Generalizing temporal relations in videos to novel domains." *IEEE TPAMI*, 2021.

*Ben-David, Shai, et al. "Analysis of representations for domain adaptation." *Advances in neural information processing systems* 19 (2006).

VideoDG review : APN

APN(Adversarial Pyramid Network)

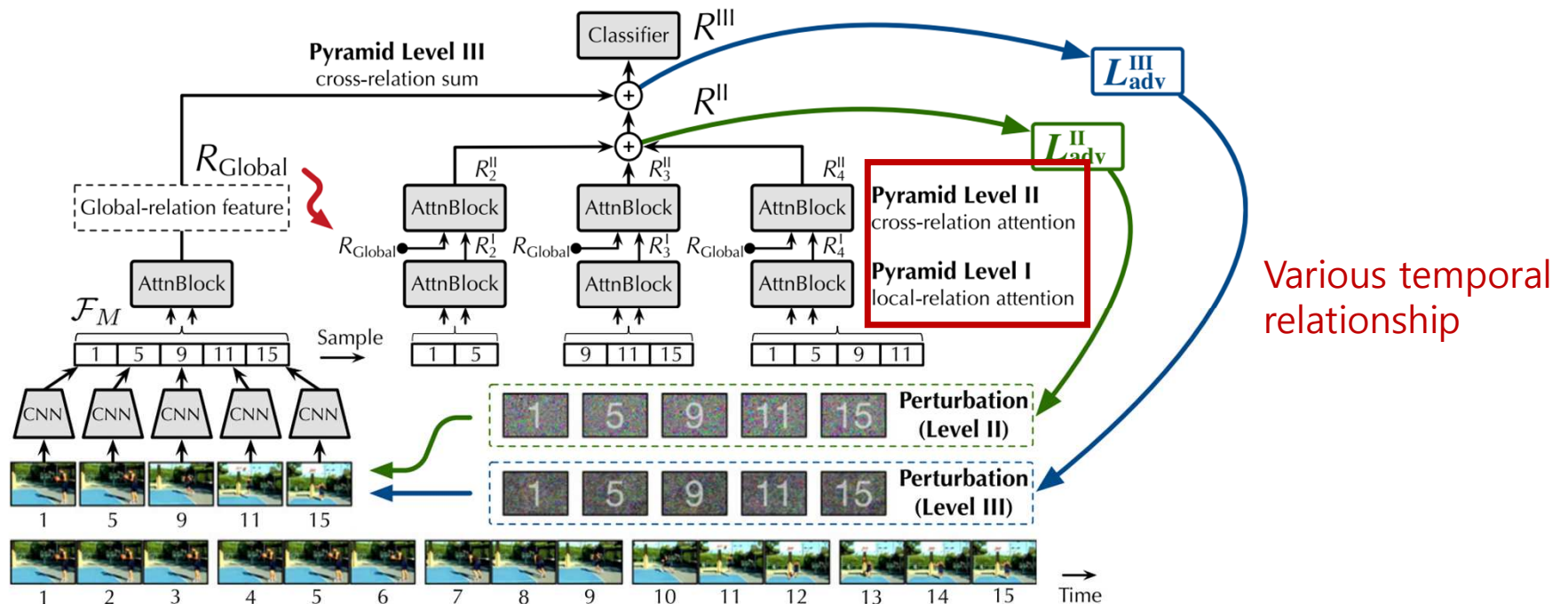
To use the **global-relation features** to guide the generalization of local events



VideoDG review : APN

APN(Adversarial Pyramid Network)

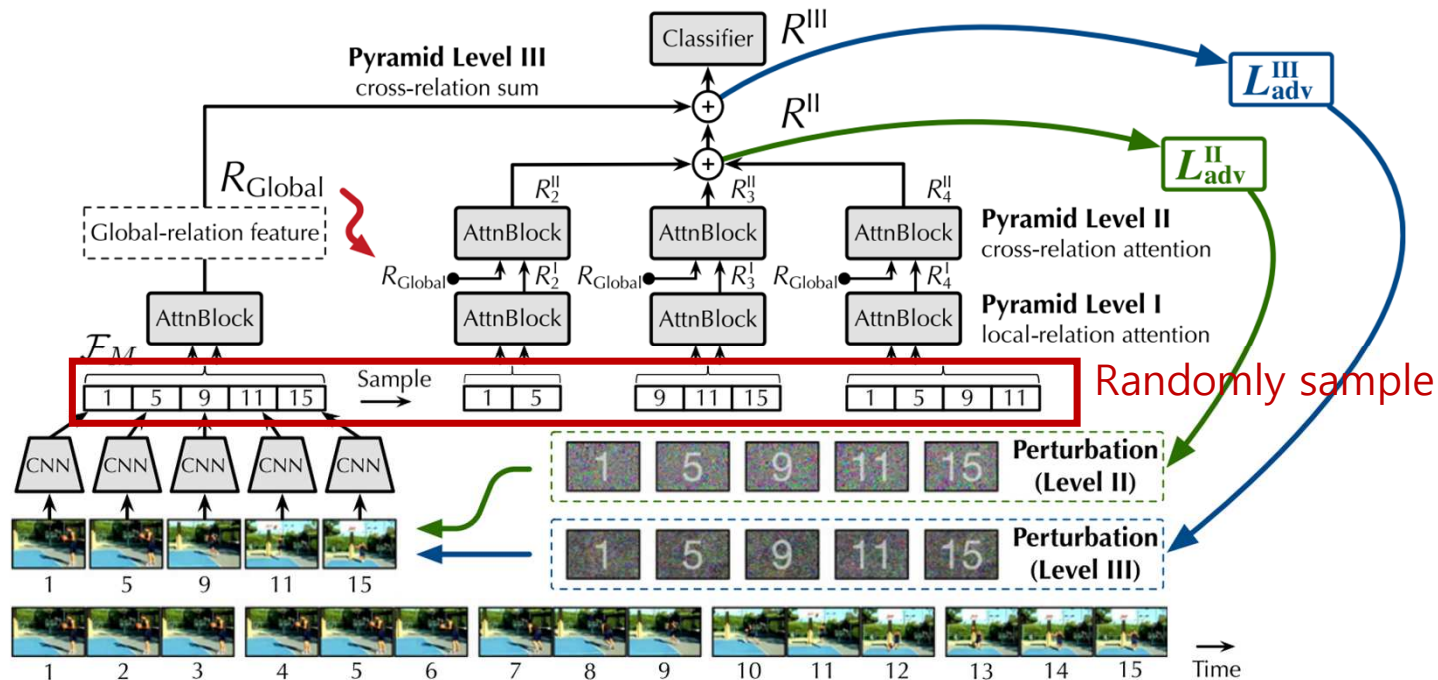
To use the global-relation features to guide the generalization of local events
And dynamically **find the events that are highly relevant to the overall video representation**



VideoDG review : APN

Pyramid I : Local-Relation Attention

Obtain R_1^I, \dots, R_4^I : **attention in various local snippets**

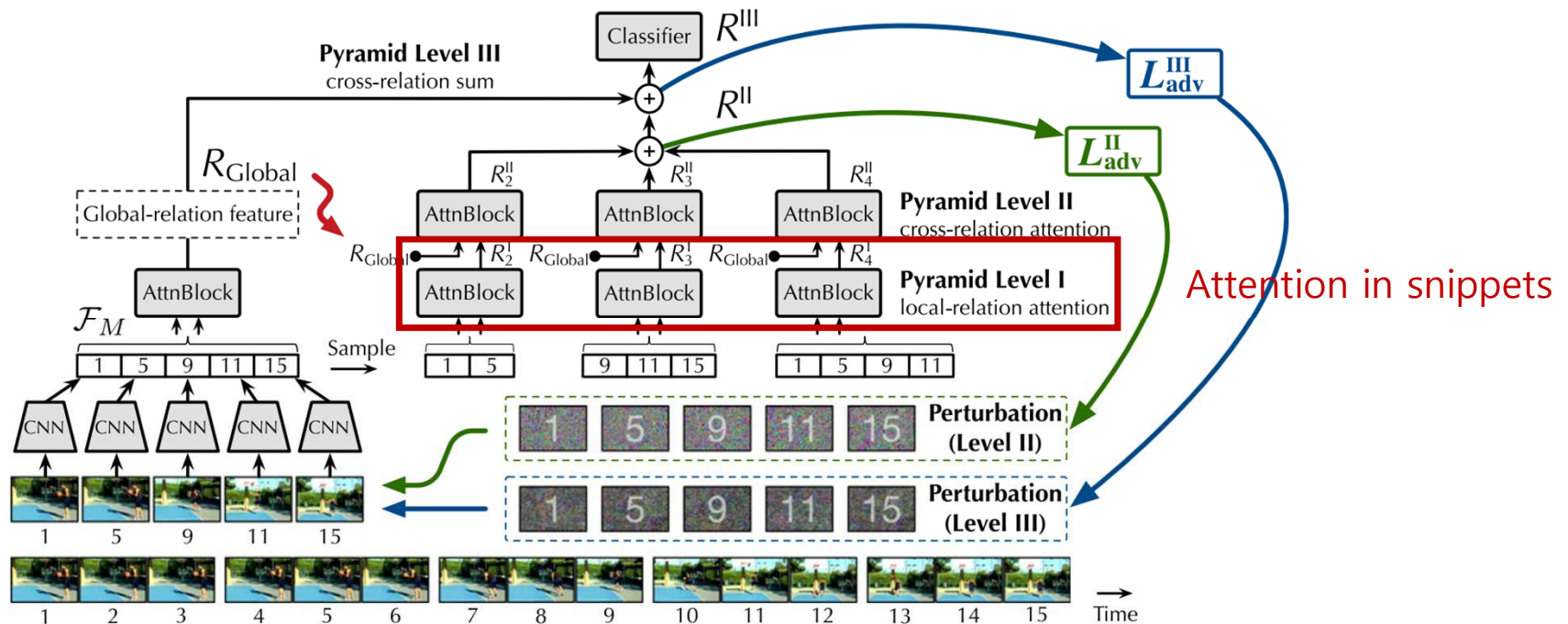


VideoDG review : APN

Pyramid I : Local-Relation Attention

Obtain R_1^I, \dots, R_4^I : attention in various local snippets

Mining local information (more generative)

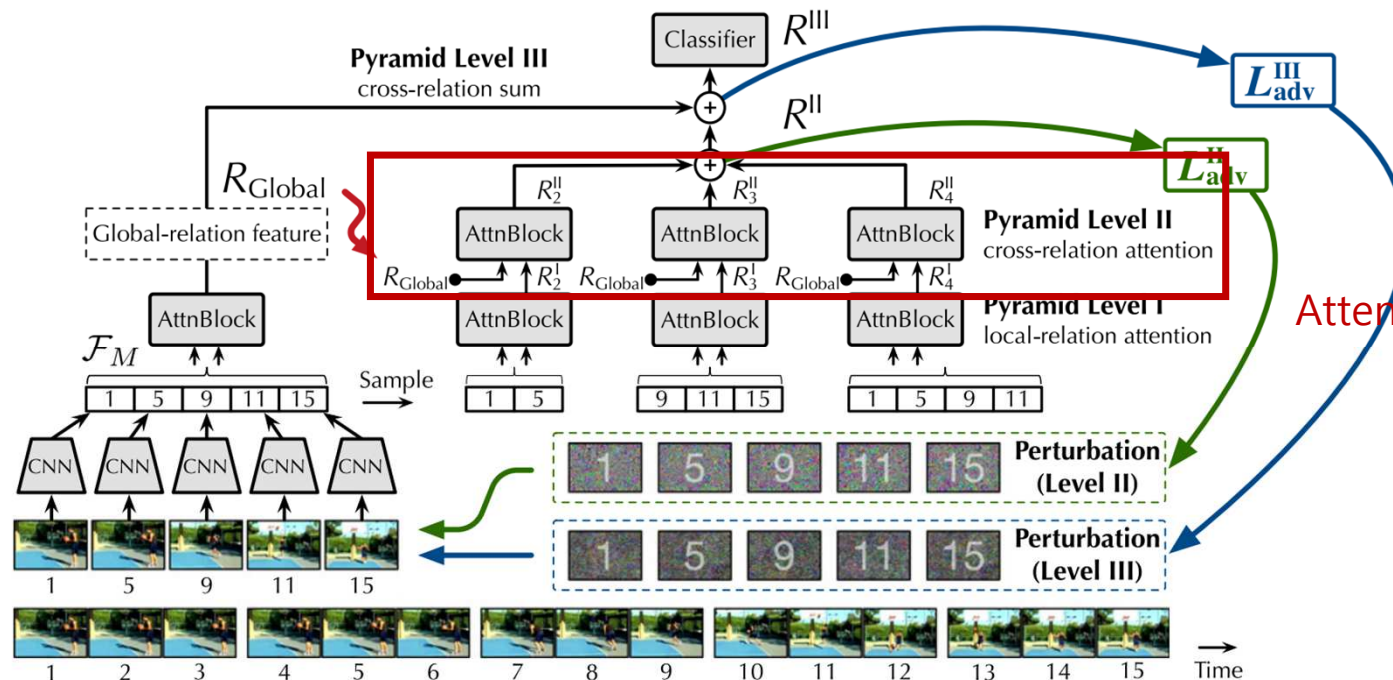


VideoDG review : APN

Pyramid II : Cross-Relation Attention

Obtain $R_1^{II}, \dots, R_4^{II}$: attention in global and local

**Constraint local feature with categorical information
To more discriminative recognition**



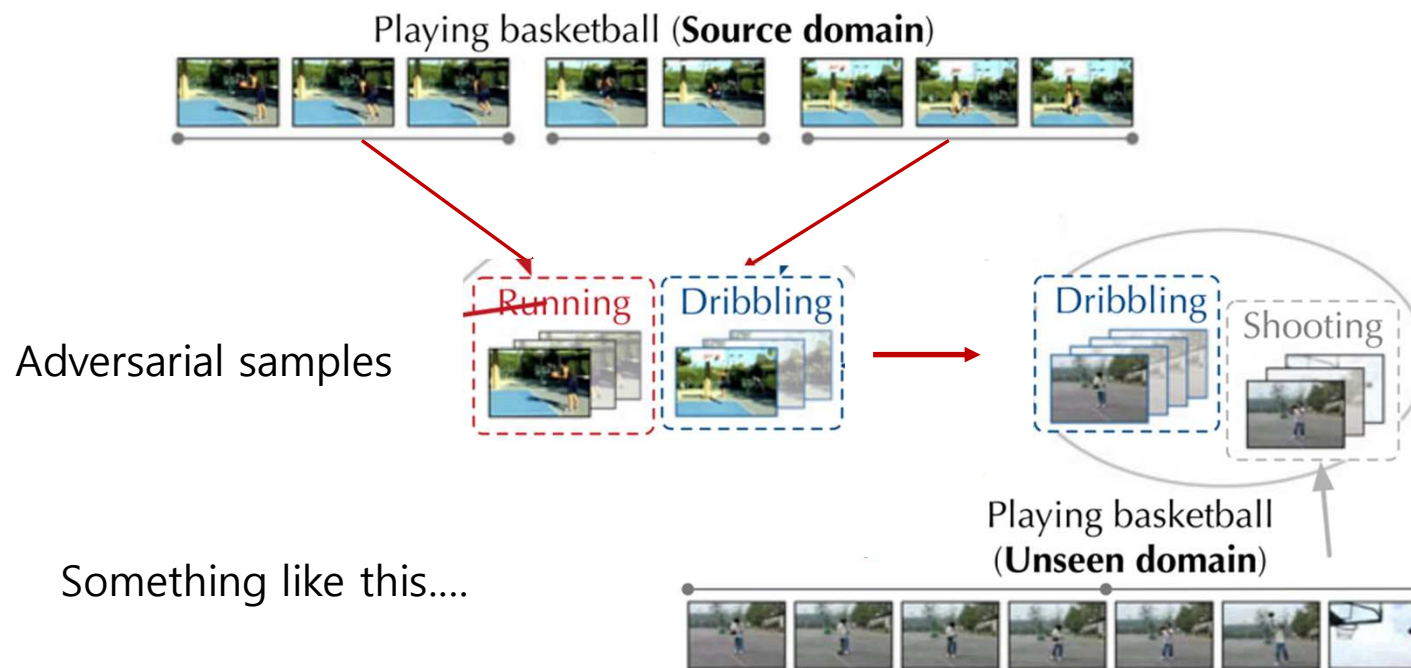
Attention in global and local

VideoDG review : RADA

RADA algorithm

Purpose :

- Generate the adversarial samples using training domain samples **close as possible as to that of the invisible target domain.**



VideoDG review : RADA

RADA algorithm

Purpose :

- Generate the adversarial samples using training domain samples close as possible as to that of the invisible target domain.

VideoDG review : RADA

RADA algorithm

Purpose :

- Generate the adversarial samples using training domain samples close as possible as to that of the invisible target domain.

Phases :

(1) T_max maximization phases :

Generate adversarial examples from the multi-level relational features

VideoDG review : RADA

RADA algorithm

Purpose :

- Generate the adversarial samples using training domain samples close as possible as to that of the invisible target domain.

Phases :

(1) T_max maximization phases :

Generate adversarial examples from the multi-level relational features

(2) **Minimization** phase of **classification error** with a robustness **regularization** :

Train generalizable features to be unaffected by overly divergent new data points

VideoDG review : RADA

RADA algorithm

Purpose :

- Generate the adversarial samples using training domain samples close as possible as to that of the invisible target domain.

Phases :

(1) T_max maximization phases :

Generate adversarial examples from the multi-level relational features

(2) **Minimization** phase of **classification error** with a robustness **regularization** :

Train generalizable features to be unaffected by overly divergent new data points

VideoDG review : RADA

RADA algorithm (1) minimize source classification loss

1.1 randomly sample data from batch

Algorithm 2. The RADA Framework for Training APN

Input: A source video dataset $\mathcal{S} = \{(X_i, Y_i)\}_{i=1}^n$, the penalty parameter of the transportation cost γ , the robustness regularization parameter λ , the number of maximization phases T_{\max} , and the learning rate α

Output: Learned APN weights θ

```
1: Initialize  $\theta \leftarrow \theta_0$ 
2: repeat
3:    $(X, Y) \sim \mathcal{S}$   $\triangleright$  Randomly sample a batch of data
4:    $(\mathcal{R}_0^{\text{II}}, \mathcal{R}_0^{\text{III}}) = \text{APN}(X)$ 
5:    $\theta \leftarrow \theta - \alpha \nabla_{\theta} \ell(h(\mathcal{R}_0^{\text{III}}), Y)$   $\triangleright$  Minimize the classification loss of the source data
6:    $X^{\text{II}} = X; X^{\text{III}} = X$ 
7:   for  $t = 1, \dots, T_{\max}$  do  $\triangleright$  For each maximization phase
8:      $(\mathcal{R}^{\text{II}}, \_) = \text{APN}(X^{\text{II}})$ 
9:      $X^{\text{II}} \leftarrow X^{\text{II}} + \nabla_X L_{\text{adv}}^{\text{II}}(\theta; (X, Y))$   $\triangleright$  Generate new data according to Eq. (9)
10:     $(\_, \mathcal{R}^{\text{III}}) = \text{APN}(X^{\text{III}})$ 
11:     $X^{\text{III}} \leftarrow X^{\text{III}} + \nabla_X L_{\text{adv}}^{\text{III}}(\theta; (X, Y))$ 
12:  end for
13:   $(\mathcal{R}^{\text{II}}, \_) = \text{APN}(X^{\text{II}})$   $\triangleright$  For minimization phase with robust training
14:   $(\_, \mathcal{R}^{\text{III}}) = \text{APN}(X^{\text{III}})$ 
15:   $\theta \leftarrow \theta - \alpha \nabla_{\theta} (L_{\text{cls}}^{\text{II}}(\theta; (X^{\text{II}}, Y)) + \lambda L_{\text{robust}}^{\text{II}}(\theta; (X^{\text{II}}, Y)))$   $\triangleright$  According to Eq. (12)
16:   $\theta \leftarrow \theta - \alpha \nabla_{\theta} (L_{\text{cls}}^{\text{III}}(\theta; (X^{\text{III}}, Y)) + \lambda L_{\text{robust}}^{\text{III}}(\theta; (X^{\text{III}}, Y)))$ 
17: until Convergence
```

VideoDG review : RADA

RADA algorithm (1) minimize source classification loss

1.2 APN inference

Algorithm 2. The RADA Framework for Training APN

Input: A source video dataset $\mathcal{S} = \{(X_i, Y_i)\}_{i=1}^n$, the penalty parameter of the transportation cost γ , the robustness regularization parameter λ , the number of maximization phases T_{\max} , and the learning rate α

Output: Learned APN weights θ

```
1: Initialize  $\theta \leftarrow \theta_0$ 
2: repeat
3:    $(X, Y) \sim \mathcal{S}$   $\triangleright$  Randomly sample a batch of data
4:    $(\mathcal{R}_0^{\text{II}}, \mathcal{R}_0^{\text{III}}) = \text{APN}(X)$ 
5:    $\theta \leftarrow \theta - \alpha \nabla_{\theta} \ell(h(\mathcal{R}_0^{\text{III}}), Y)$   $\triangleright$  Minimize the classification loss of the source data
6:    $X^{\text{II}} = X; X^{\text{III}} = X$ 
7:   for  $t = 1, \dots, T_{\max}$  do  $\triangleright$  For each maximization phase
8:      $(\mathcal{R}^{\text{II}}, \_) = \text{APN}(X^{\text{II}})$ 
9:      $X^{\text{II}} \leftarrow X^{\text{II}} + \nabla_X L_{\text{adv}}^{\text{II}}(\theta; (X, Y))$   $\triangleright$  Generate new data according to Eq. (9)
10:     $(\_, \mathcal{R}^{\text{III}}) = \text{APN}(X^{\text{III}})$ 
11:     $X^{\text{III}} \leftarrow X^{\text{III}} + \nabla_X L_{\text{adv}}^{\text{III}}(\theta; (X, Y))$ 
12:  end for
13:   $(\mathcal{R}^{\text{II}}, \_) = \text{APN}(X^{\text{II}})$   $\triangleright$  For minimization phase with robust training
14:   $(\_, \mathcal{R}^{\text{III}}) = \text{APN}(X^{\text{III}})$ 
15:   $\theta \leftarrow \theta - \alpha \nabla_{\theta} (L_{\text{cls}}^{\text{II}}(\theta; (X^{\text{II}}, Y)) + \lambda L_{\text{robust}}^{\text{II}}(\theta; (X^{\text{II}}, Y)))$   $\triangleright$  According to Eq. (12)
16:   $\theta \leftarrow \theta - \alpha \nabla_{\theta} (L_{\text{cls}}^{\text{III}}(\theta; (X^{\text{III}}, Y)) + \lambda L_{\text{robust}}^{\text{III}}(\theta; (X^{\text{III}}, Y)))$ 
17: until Convergence
```

VideoDG review : RADA

RADA algorithm (1) minimize source classification loss

1.3 classification loss backward (on source domain)

Algorithm 2. The RADA Framework for Training APN

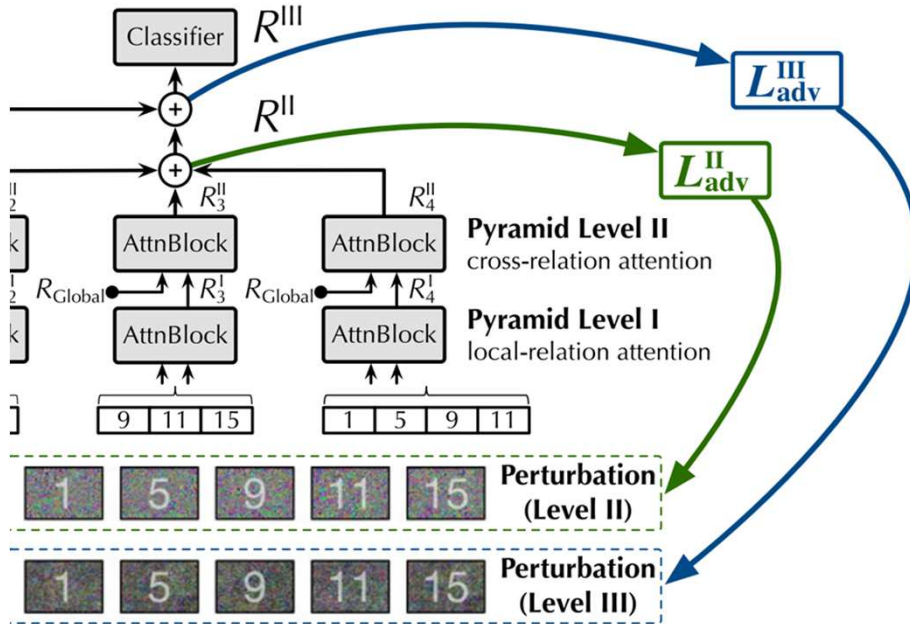
Input: A source video dataset $\mathcal{S} = \{(X_i, Y_i)\}_{i=1}^n$, the penalty parameter of the transportation cost γ , the robustness regularization parameter λ , the number of maximization phases T_{\max} , and the learning rate α

Output: Learned APN weights θ

```
1: Initialize  $\theta \leftarrow \theta_0$ 
2: repeat
3:    $(X, Y) \sim \mathcal{S}$   $\triangleright$  Randomly sample a batch of data
4:    $(\mathcal{R}_0^{\text{II}}, \mathcal{R}_0^{\text{III}}) = \text{APN}(X)$ 
5:    $\theta \leftarrow \theta - \alpha \nabla_{\theta} \ell(h(\mathcal{R}_0^{\text{III}}), Y)$   $\triangleright$  Minimize the classification loss of the source data
6:    $X^{\text{II}} = X; X^{\text{III}} = X$ 
7:   for  $t = 1, \dots, T_{\max}$  do  $\triangleright$  For each maximization phase
8:      $(\mathcal{R}^{\text{II}}, \_) = \text{APN}(X^{\text{II}})$ 
9:      $X^{\text{II}} \leftarrow X^{\text{II}} + \nabla_X L_{\text{adv}}^{\text{II}}(\theta; (X, Y))$   $\triangleright$  Generate new data according to Eq. (9)
10:     $(\_, \mathcal{R}^{\text{III}}) = \text{APN}(X^{\text{III}})$ 
11:     $X^{\text{III}} \leftarrow X^{\text{III}} + \nabla_X L_{\text{adv}}^{\text{III}}(\theta; (X, Y))$ 
12:  end for
13:   $(\mathcal{R}^{\text{II}}, \_) = \text{APN}(X^{\text{II}})$   $\triangleright$  For minimization phase with robust training
14:   $(\_, \mathcal{R}^{\text{III}}) = \text{APN}(X^{\text{III}})$ 
15:   $\theta \leftarrow \theta - \alpha \nabla_{\theta} (L_{\text{cls}}^{\text{II}}(\theta; (X^{\text{II}}, Y)) + \lambda L_{\text{robust}}^{\text{II}}(\theta; (X^{\text{II}}, Y)))$   $\triangleright$  According to Eq. (12)
16:   $\theta \leftarrow \theta - \alpha \nabla_{\theta} (L_{\text{cls}}^{\text{III}}(\theta; (X^{\text{III}}, Y)) + \lambda L_{\text{robust}}^{\text{III}}(\theta; (X^{\text{III}}, Y)))$ 
17: until Convergence
```

VideoDG review : RADA

RADA algorithm (2) maximization phases



2.2 Adversarial sample generation!

Algorithm 2. The RADA Framework for Training APN

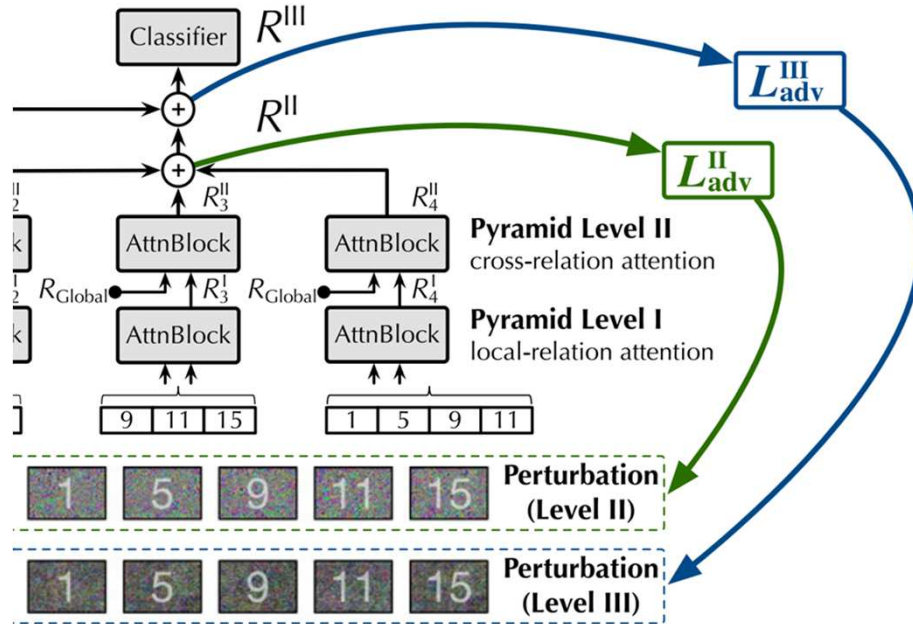
Input: A source video dataset $\mathcal{S} = \{(X_i, Y_i)\}_{i=1}^n$, the penalty parameter of the transportation cost γ , the robustness regularization parameter λ , the number of maximization phases T_{max} , and the learning rate α

Output: Learned APN weights θ

- 1: **Initialize** $\theta \leftarrow \theta_0$
- 2: **repeat**
- 3: $(X, Y) \sim \mathcal{S}$ \triangleright Randomly sample a batch of data
- 4: $(\mathcal{R}_0^{II}, \mathcal{R}_0^{III}) = \text{APN}(X)$
- 5: $\theta \leftarrow \theta - \alpha \nabla_{\theta} \ell(h(\mathcal{R}_0^{III}), Y)$ \triangleright Minimize the classification loss of the source data
- 6: $X^{II} = X; X^{III} = X$
- 7: **for** $t = 1, \dots, T_{max}$ **do** \triangleright For each maximization phase
- 8: $(\mathcal{R}^{II}, _) = \text{APN}(X^{II})$
- 9: $X^{II} \leftarrow X^{II} + \nabla_X L^{II}_{adv}(\theta; (X, Y))$ \triangleright Generate new data according to Eq. (9)
- 10: $(_, \mathcal{R}^{III}) = \text{APN}(X^{III})$
- 11: $X^{III} \leftarrow X^{III} + \nabla_X L^{III}_{adv}(\theta; (X, Y))$
- 12: **end for**
- 13: $(\mathcal{R}^{II}, _) = \text{APN}(X^{II})$ \triangleright For minimization phase with robust training
- 14: $(_, \mathcal{R}^{III}) = \text{APN}(X^{III})$
- 15: $\theta \leftarrow \theta - \alpha \nabla_{\theta} (L^{II}_{cls}(\theta; (X^{II}, Y)) + \lambda L^{II}_{robust}(\theta; (X^{II}, Y)))$ \triangleright According to Eq. (12)
- 16: $\theta \leftarrow \theta - \alpha \nabla_{\theta} (L^{III}_{cls}(\theta; (X^{III}, Y)) + \lambda L^{III}_{robust}(\theta; (X^{III}, Y)))$
- 17: **until** Convergence

VideoDG review : RADA

RADA algorithm (2) maximization phases



2.2 Adversarial sample generation!



Algorithm 2. The RADA Framework for Training APN

Input: A source video dataset $\mathcal{S} = \{(X_i, Y_i)\}_{i=1}^n$, the penalty parameter of the transportation cost γ , the robustness regularization parameter λ , the number of maximization phases T_{max} , and the learning rate α

Output: Learned APN weights θ

- 1: **Initialize** $\theta \leftarrow \theta_0$
- 2: **repeat**
- 3: $(X, Y) \sim \mathcal{S}$ \triangleright Randomly sample a batch of data
- 4: $(\mathcal{R}_0^{II}, \mathcal{R}_0^{III}) = \text{APN}(X)$
- 5: $\theta \leftarrow \theta - \alpha \nabla_{\theta} \ell(h(\mathcal{R}_0^{III}), Y)$ \triangleright Minimize the classification loss of the source data
- 6: $X^{II} = X; X^{III} = X$
- 7: **for** $t = 1, \dots, T_{max}$ **do** \triangleright For each maximization phase
- 8: $(\mathcal{R}^{II}, _) = \text{APN}(X^{II})$
- 9: $X^{II} \leftarrow X^{II} + \nabla_X L_{adv}^{II}(\theta; (X, Y))$ \triangleright Generate new data according to Eq. (9)
- 10: $(_, \mathcal{R}^{III}) = \text{APN}(X^{III})$
- 11: $X^{III} \leftarrow X^{III} + \nabla_X L_{adv}^{III}(\theta; (X, Y))$
- 12: **end for**
- 13: $(\mathcal{R}^{II}, _) = \text{APN}(X^{II})$ \triangleright For minimization phase with robust training
- 14: $(_, \mathcal{R}^{III}) = \text{APN}(X^{III})$
- 15: $\theta \leftarrow \theta - \alpha \nabla_{\theta} (L_{cls}^{II}(\theta; (X^{II}, Y)) + \lambda L_{robust}^{II}(\theta; (X^{II}, Y)))$ \triangleright According to Eq. (12)
- 16: $\theta \leftarrow \theta - \alpha \nabla_{\theta} (L_{cls}^{III}(\theta; (X^{III}, Y)) + \lambda L_{robust}^{III}(\theta; (X^{III}, Y)))$
- 17: **until** Convergence

VideoDG review : RADA

RADA algorithm (2) maximization phases

- Make perturbations

$$\underbrace{X^{\text{II}}}_{\text{input}} \leftarrow \underbrace{X^{\text{II}} + \nabla_X L^{\text{II}}_{\text{adv}}(\theta; (X, Y))}_{\text{perturbations}}$$

2.1 Add perturbations to adversarial samples



VideoDG review : RADA

RADA algorithm (2) maximization phases

- Make perturbations

$$\underbrace{X^{\text{II}}}_{\text{input}} \leftarrow \underbrace{X^{\text{II}} + \nabla_X L_{\text{adv}}^{\text{II}}(\theta; (X, Y))}_{\text{perturbations}}$$

$$L_{\text{adv}}^k(\theta; (X, Y)) := \sup_{X \in \mathcal{X}} \{ \ell(h(\mathcal{R}^k), Y) - \gamma c(\mathcal{R}^k, \mathcal{R}_0^k) \}, \quad (9)$$



VideoDG review : RADA

RADA algorithm (2) maximization phases

- Make perturbations

$$\underbrace{X^{\text{II}}}_{\text{input}} \leftarrow \underbrace{X^{\text{II}} + \nabla_X L_{\text{adv}}^{\text{II}}(\theta; (X, Y))}_{\text{perturbations}}$$

$$L_{\text{adv}}^k(\theta; (X, Y)) := \sup_{X \in \mathcal{X}} \{ \ell(h(\mathcal{R}^k), Y) - \gamma c(\mathcal{R}^k, \mathcal{R}_0^k) \}, \quad (9)$$

K=1,2

$$l(h(R^k), Y)$$

Just cross entropy loss With FC layer h

: The **added perturbation must not move away from semantic.**



VideoDG review : RADA

RADA algorithm (2) maximization phases

- Make perturbations

$$\underbrace{X^{\text{II}}}_{\text{input}} \leftarrow \underbrace{X^{\text{II}} + \nabla_X L_{\text{adv}}^{\text{II}}(\theta; (X, Y))}_{\text{perturbations}}$$

$$L_{\text{adv}}^k(\theta; (X, Y)) := \sup_{X \in \mathcal{X}} \{ \ell(h(\mathcal{R}^k), Y) - \gamma c(\mathcal{R}^k, \mathcal{R}_0^k) \}, \quad (9)$$

K=1,2

$\gamma c(R^k, R_0^k)$

Just weighted mse loss between R^k, R_0^k

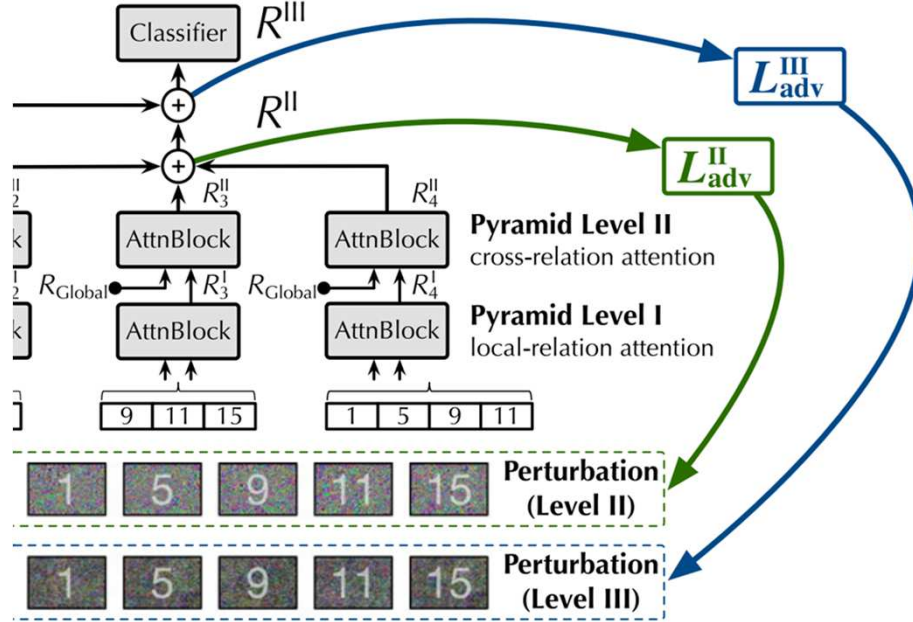
: Regulates features so that they are not too far from aggregated features at each level **to maintain consistency in advertising samples**



VideoDG review : RADA

RADA algorithm (2) maximization phases

- Make perturbations



2.2 Add perturbation few times to adversarial sample

Algorithm 2. The RADA Framework for Training APN

Input: A source video dataset $\mathcal{S} = \{(X_i, Y_i)\}_{i=1}^n$, the penalty parameter of the transportation cost γ , the robustness regularization parameter λ , the number of maximization phases T_{max} , and the learning rate α

Output: Learned APN weights θ

```

1: Initialize  $\theta \leftarrow \theta_0$ 
2: repeat
3:    $(X, Y) \sim \mathcal{S}$   $\triangleright$  Randomly sample a batch of data
4:    $(\mathcal{R}_0^{II}, \mathcal{R}_0^{III}) = \text{APN}(X)$ 
5:    $\theta \leftarrow \theta - \alpha \nabla_{\theta} \ell(h(\mathcal{R}_0^{III}), Y)$   $\triangleright$  Minimize the classification loss of the source data
6:    $X^{II} = X; X^{III} = X$ 
7:   for  $t = 1, \dots, T_{max}$  do  $\triangleright$  For each maximization phase
8:      $(\mathcal{R}^{II}, \_) = \text{APN}(X^{II})$ 
9:      $X^{II} \leftarrow X^{II} + \nabla_X L_{adv}^{II}(\theta; (X, Y))$   $\triangleright$  Generate new data according to Eq. (9)
10:     $(\_, \mathcal{R}^{III}) = \text{APN}(X^{III})$ 
11:     $X^{III} \leftarrow X^{III} + \nabla_X L_{adv}^{III}(\theta; (X, Y))$ 
12:  end for
13:   $(\mathcal{R}^{II}, \_) = \text{APN}(X^{II})$   $\triangleright$  For minimization phase with robust training
14:   $(\_, \mathcal{R}^{III}) = \text{APN}(X^{III})$ 
15:   $\theta \leftarrow \theta - \alpha \nabla_{\theta} (L_{cls}^{II}(\theta; (X^{II}, Y)) + \lambda L_{robust}^{II}(\theta; (X^{II}, Y)))$   $\triangleright$  According to Eq. (12)
16:   $\theta \leftarrow \theta - \alpha \nabla_{\theta} (L_{cls}^{III}(\theta; (X^{III}, Y)) + \lambda L_{robust}^{III}(\theta; (X^{III}, Y)))$ 
17: until Convergence
    
```

VideoDG review : RADA

RADA algorithm (3) minimization phases

- Compensate overly divergent generated samples

$$L_{\text{cls}}^k(\theta; (X^k, Y)) := \underbrace{\ell(h(\mathcal{R}^k), Y)}_{\text{classify adversarial examples}}, \quad (10)$$

$$L_{\text{robust}}^k(\theta; (X^k, Y)) := \underbrace{\ell(h(\mathcal{R}^k), h(\mathcal{R}_0^{\text{III}}))}_{\text{robustness regularization}}, \quad (11)$$

$$L_{\text{min}} := L_{\text{cls}}^k(\theta; (X^k, Y)) + \lambda L_{\text{robust}}^k(\theta; (X^k, Y)), \quad (12)$$

3.1 Cross entropy loss on adversarial example



Algorithm 2. The RADA Framework for Training APN

Input: A source video dataset $\mathcal{S} = \{(X_i, Y_i)\}_{i=1}^n$, the penalty parameter of the transportation cost γ , the robustness regularization parameter λ , the number of maximization phases T_{max} , and the learning rate α

Output: Learned APN weights θ

```

1: Initialize  $\theta \leftarrow \theta_0$ 
2: repeat
3:    $(X, Y) \sim \mathcal{S}$   $\triangleright$  Randomly sample a batch of data
4:    $(\mathcal{R}_0^{\text{II}}, \mathcal{R}_0^{\text{III}}) = \text{APN}(X)$ 
5:    $\theta \leftarrow \theta - \alpha \nabla_{\theta} \ell(h(\mathcal{R}_0^{\text{III}}), Y)$   $\triangleright$  Minimize the classification loss of the source data
6:    $X^{\text{II}} = X; X^{\text{III}} = X$ 
7:   for  $t = 1, \dots, T_{\text{max}}$  do  $\triangleright$  For each maximization phase
8:      $(\mathcal{R}^{\text{II}}, \_) = \text{APN}(X^{\text{II}})$ 
9:      $X^{\text{II}} \leftarrow X^{\text{II}} + \nabla_X L_{\text{adv}}^{\text{II}}(\theta; (X, Y))$   $\triangleright$  Generate new data according to Eq. (9)
10:     $(\_, \mathcal{R}^{\text{III}}) = \text{APN}(X^{\text{III}})$ 
11:     $X^{\text{III}} \leftarrow X^{\text{III}} + \nabla_X L_{\text{adv}}^{\text{III}}(\theta; (X, Y))$ 
12:  end for
13:   $(\mathcal{R}^{\text{II}}, \_) = \text{APN}(X^{\text{II}})$   $\triangleright$  For minimization phase with robust training
14:   $(\_, \mathcal{R}^{\text{III}}) = \text{APN}(X^{\text{III}})$ 
15:   $\theta \leftarrow \theta - \alpha \nabla_{\theta} (L_{\text{cls}}^{\text{II}}(\theta; (X^{\text{II}}, Y)) + \lambda L_{\text{robust}}^{\text{II}}(\theta; (X^{\text{II}}, Y)))$   $\triangleright$  According to Eq. (12)
16:   $\theta \leftarrow \theta - \alpha \nabla_{\theta} (L_{\text{cls}}^{\text{III}}(\theta; (X^{\text{III}}, Y)) + \lambda L_{\text{robust}}^{\text{III}}(\theta; (X^{\text{III}}, Y)))$ 
17: until Convergence
  
```

VideoDG review : RADA

RADA algorithm (3) minimization phases

- Compensate overly divergent generated samples

$$L_{\text{cls}}^k(\theta; (X^k, Y)) := \underbrace{\ell(h(\mathcal{R}^k), Y)}_{\text{classify adversarial examples}}, \quad (10)$$

$$L_{\text{robust}}^k(\theta; (X^k, Y)) := \underbrace{\ell(h(\mathcal{R}^k), h(\mathcal{R}_0^{\text{III}}))}_{\text{robustness regularization}}, \quad (11)$$

$$L_{\text{min}} := L_{\text{cls}}^k(\theta; (X^k, Y)) + \lambda L_{\text{robust}}^k(\theta; (X^k, Y)), \quad (12)$$

3.2 robustness regularization

Cross-entropy loss between
predicted distribution from Adversarial sample and
Predicted distribution from source(original) sample



Algorithm 2. The RADA Framework for Training APN

Input: A source video dataset $\mathcal{S} = \{(X_i, Y_i)\}_{i=1}^n$, the penalty parameter of the transportation cost γ , the robustness regularization parameter λ , the number of maximization phases T_{max} , and the learning rate α

Output: Learned APN weights θ

```

1: Initialize  $\theta \leftarrow \theta_0$ 
2: repeat
3:    $(X, Y) \sim \mathcal{S}$   $\triangleright$  Randomly sample a batch of data
4:    $(\mathcal{R}_0^{\text{II}}, \mathcal{R}_0^{\text{III}}) = \text{APN}(X)$ 
5:    $\theta \leftarrow \theta - \alpha \nabla_{\theta} \ell(h(\mathcal{R}_0^{\text{III}}), Y)$   $\triangleright$  Minimize the classification
      loss of the source data
6:    $X^{\text{II}} = X; X^{\text{III}} = X$ 
7:   for  $t = 1, \dots, T_{\text{max}}$  do  $\triangleright$  For each maximization phase
8:      $(\mathcal{R}^{\text{II}}, \_) = \text{APN}(X^{\text{II}})$ 
9:      $X^{\text{II}} \leftarrow X^{\text{II}} + \nabla_X L_{\text{adv}}^{\text{II}}(\theta; (X, Y))$   $\triangleright$  Generate new data
      according to Eq. (9)
10:     $(\_, \mathcal{R}^{\text{III}}) = \text{APN}(X^{\text{III}})$ 
11:     $X^{\text{III}} \leftarrow X^{\text{III}} + \nabla_X L_{\text{adv}}^{\text{III}}(\theta; (X, Y))$ 
12:  end for
13:   $(\mathcal{R}^{\text{II}}, \_) = \text{APN}(X^{\text{II}})$   $\triangleright$  For minimization phase with
      robust training
14:   $(\_, \mathcal{R}^{\text{III}}) = \text{APN}(X^{\text{III}})$ 
15:   $\theta \leftarrow \theta - \alpha \nabla_{\theta} (L_{\text{cls}}^{\text{II}}(\theta; (X^{\text{II}}, Y)) + \lambda L_{\text{robust}}^{\text{II}}(\theta; (X^{\text{II}}, Y)))$ 
       $\triangleright$  According to Eq. (12)
16:   $\theta \leftarrow \theta - \alpha \nabla_{\theta} (L_{\text{cls}}^{\text{III}}(\theta; (X^{\text{III}}, Y)) + \lambda L_{\text{robust}}^{\text{III}}(\theta; (X^{\text{III}}, Y)))$ 
17: until Convergence
  
```

Conclusion

- The i.i.d. assumption are the fundamental but crucial assumption when training machine learning model
- However the i.i.d. assumption are easily violated in real world.
- The Domain Generalization(DG) has been proposed for solving these problem.
- However, the Video Domain Generalization is an un-charted area
- In Video, the temporal and spatial domain shifts simultaneously occurs
- Therefore temporal domain shifts are main challenge in video domain generalization
- In previous work, the temporal attention on diverse relation in time axis is tried.
- The adversarial samples are one of the method to solve VDG.

UST seminar

Towards domain-agnostic Video action recognition

Hyungmin Kim

UST-ETRI

Ph.D. student

khm159@etri.re.kr/ust.ac.kr

21th Sep, 2023