

# Domain Generalization

CS 330

# Logistics

Project milestone on **Wednesday, November 16**

Homework 4 (optional) due **Monday, November 14**

# Plan for Today

## Domain Generalization

- Problem formulation
- Algorithms
  - Adding explicit regularizers
  - Data augmentation

## Goals for this lecture:

- Understand domain generalization: intuition, problem formulation
- Familiarize mainstream DG approaches: regularization-based, augmentation-based

# Recap: Domain Adaptation

Perform well on target domain  $p_T(x, y)$ ,  
using training data from source domain(s)  $p_S(x, y)$

A form of **transfer learning**, with access to target domain data during training  
("transductive" learning)

Unsupervised domain adaptation: access to unlabeled target domain data

Common assumptions:

- Source and target domain only differ in domain of the function, i.e.  $p_S(y | x) = p_T(y | x)$
- There exists a single hypothesis with low error on both source and target domains.

Revisiting: A "domain" is a special case of a "task" 

→ 한글자 한글인 vs. 다글자 한글인 example

A task:  $\mathcal{T}_i \triangleq \{p_i(\mathbf{x}), p_i(\mathbf{y} | \mathbf{x}), \mathcal{L}_i\}$

A domain:  $d_i \triangleq \{p_i(\mathbf{x}), p(\mathbf{y} | \mathbf{x}), \mathcal{L}\}$

# Recap: Domain Adaptation

Perform well on target domain  $p_T(x, y)$ ,  
using training data from source domain(s)  $p_S(x, y)$

A form of **transfer learning**

Unsupervised domain

Can we always access unlabeled data  
from the target domain?



Not always.

data during training  
(e.g. "learning")

get domain data

Common assumptions:

- Source and target domain only differ in domain of the function, i.e.  $p_S(y|x) = p_T(y|x)$
- There exists a single hypothesis with low error on both source and target domains.

Revisiting: A “domain” is a special case of a “task”

$$\text{A task: } \mathcal{T}_i \triangleq \{p_i(\mathbf{x}), p_i(\mathbf{y}|\mathbf{x}), \mathcal{L}_i\} \quad \text{A domain: } d_i \triangleq \{p_i(\mathbf{x}), p(\mathbf{y}|\mathbf{x}), \mathcal{L}\}$$

# Recap: Domain Adaptation

Perform well on target domain  $p_T(x, y)$ ,  
using training data from source domain(s)  $p_S(x, y)$

A form of transfer learning with access to target domain data during training

Cases where we don't have access to unlabeled data in the target domain

- Real-time deployment and don't have time to collect target domain data
- Obtaining target data may be restricted by privacy policy

Common:

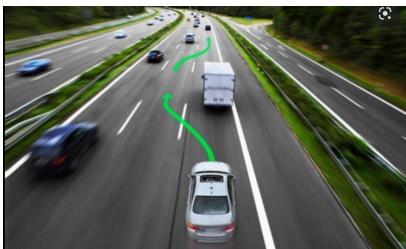
- Source and target domain only differ in domain of the function, i.e.  $p_S(y|x) = p_T(y|x)$
- There exists a single hypothesis with low error on both source and target domains.

Revisiting: A "domain" is a special case of a "task"

$$\text{A task: } \mathcal{T}_i \triangleq \{p_i(\mathbf{x}), p_i(\mathbf{y}|\mathbf{x}), \mathcal{L}_i\} \quad \text{A domain: } d_i \triangleq \{p_i(\mathbf{x}), p(\mathbf{y}|\mathbf{x}), \mathcal{L}\}$$

# Real-Time Deployment

Real-time deployment and don't have time to collect data



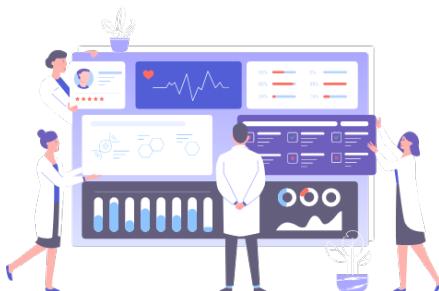
Trained on three types of roads

Can imagine many similar scenarios  
↳ smart factory images, different image collectors, etc.

Deploy to a new road

↳ Night time driving conditions

# Privacy Concerns



Trained on 3 hospitals

Deploy to a new hospital

Can't access training data

↳ Privacy concerns



# Plan for Today

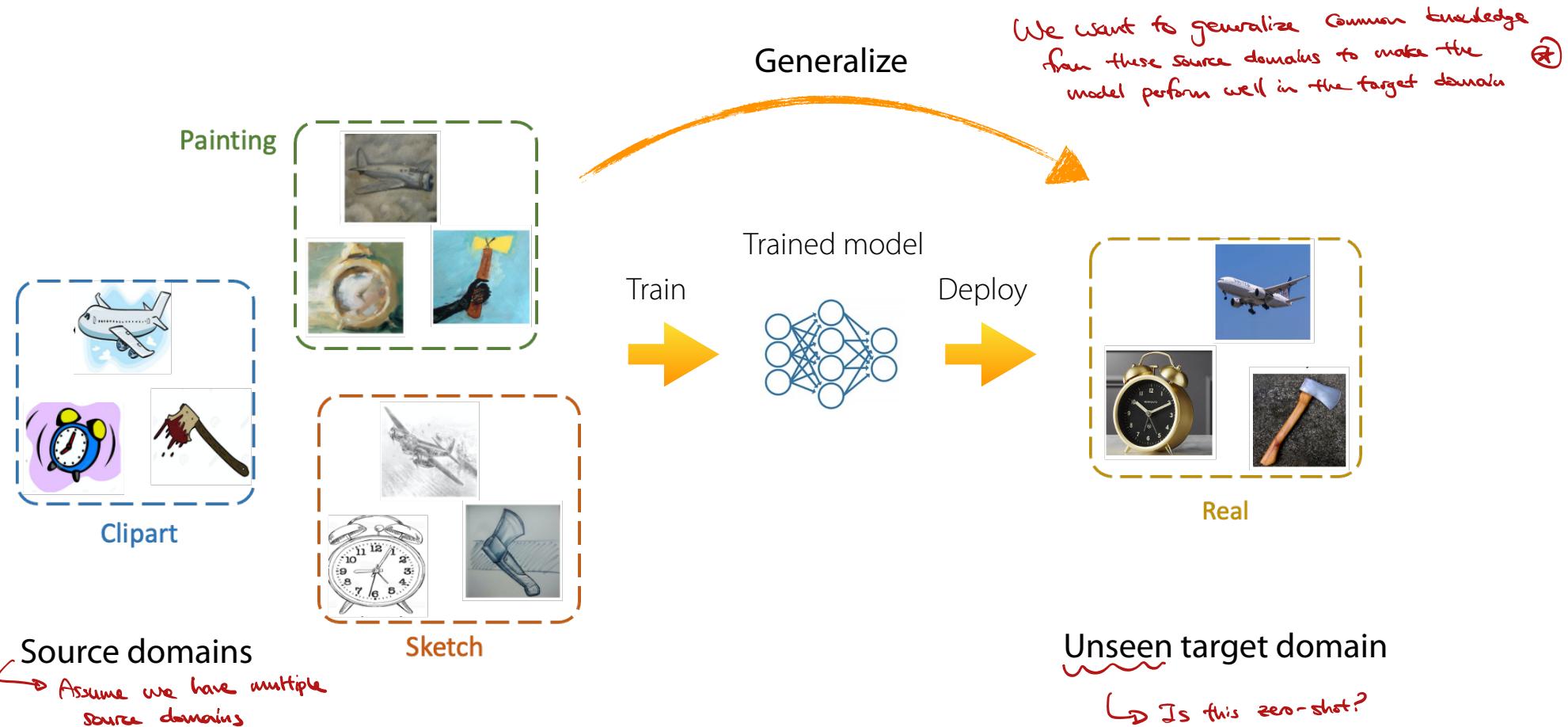
## Domain Generalization

- **Problem formulation**
- Algorithms
  - Adding explicit regularizers
  - Data augmentation

## Goals for this lecture:

- Understand [domain generalization](#): intuition, problem formulation
- Familiarize mainstream DG approaches: [regularization](#)-based, [augmentation](#)-based

# Domain Generalization



# Domain Generalization Problem

Given source domains  $p_1(x, y), \dots, p_n(x, y)$ , solve unseen target domain  $p_T(x, y)$   
without accessing the data from it.

!!!

\*\*\*

Previously in domain adaptation,  
they had access to both domain  
data + (unlabeled) target data  
during training → for unsupervised  
domain adaptation

Common assumptions

v. Similar  
set of  
assumptions  
as before

- All domains only differ in domain of the function, i.e.,  $p_1(y|x) = \dots = p_n(y|x) = p_T(y|x)$ .
- Only  $p(x)$  can change
- There exists a single hypothesis with low error in all domains.

(↳ eg) is the object an airplane, axe or a clock?

Guarantees that we can learn a model  
that can perform well across different domains.

Revisiting: A “domain” is a special case of a “task”

$$\text{A task: } \mathcal{T}_i \triangleq \{p_i(\mathbf{x}), p_i(\mathbf{y}|\mathbf{x}), \mathcal{L}_i\} \quad \text{A domain: } d_i \triangleq \{p_i(\mathbf{x}), p(\mathbf{y}|\mathbf{x}), \mathcal{L}\}$$



Key distinction!

# Meta-Learning v.s. Domain Generalization

Revisiting: A “domain” is a special case of a “task”

A task:  $\mathcal{T}_i \triangleq \{p_i(\mathbf{x}), p_i(\mathbf{y} | \mathbf{x}), \mathcal{L}_i\}$     A domain:  $d_i \triangleq \{p_i(\mathbf{x}), p(\mathbf{y} | \mathbf{x}), \mathcal{L}\}$

## Meta-Learning Problem

Transfer learning with many source tasks

Given data from  $\mathcal{T}_1, \dots, \mathcal{T}_n$ , solve new task  $\mathcal{T}_t$  more quickly / proficiently / stably

## Domain Generalization

A special case of meta-learning

Given data from domains  $d_1, \dots, d_n$ , perform well on new domain  $d_t$

- Only  $p_i(x)$  changes across tasks
- direct generalization/no adaptation

# Ⓐ Domain Adaptation v.s. Domain Generalization

## Domain Adaptation

"transductive" setting

Given labeled data from source domain  $p_S(x, y)$  and unlabeled data from target domain  $p_T(x, y)$ , perform well on this target domain

Ⓐ Target data access during training ✓ (unlabeled data)  *Unsupervised domain adaptation*

Only one source domain

The model is specialized for the target domain

## Domain Generalization

"inductive" setting

Given labeled data from a set of source domains  $p_1(x, y), \dots, p_n(x, y)$ , perform well on target domain  $p_T(x, y)$

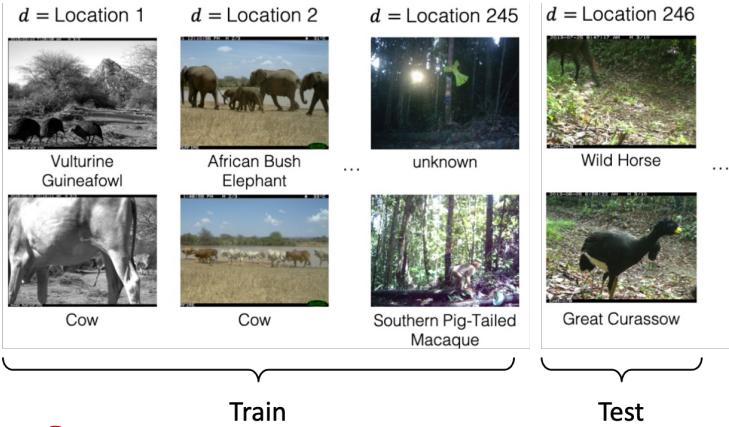
Ⓐ Test data access during training 

Need more than one source domain

The model can be applied to all domains

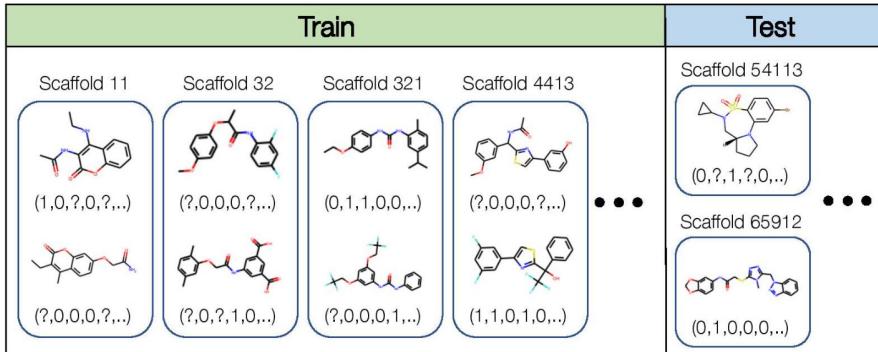
# Domain Generalization: Applications

## Wildlife recognition

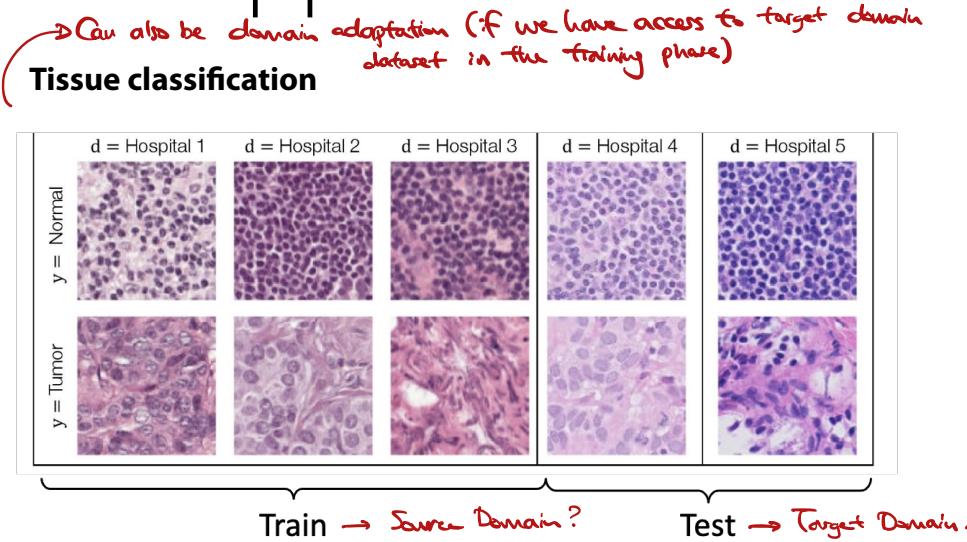


→ Dry Discovery

## Molecule property prediction



## Tissue classification



Train → Source Domain?

Test → Target Domain?

## Code completion

	Repository ID ( $d$ )	Source code context ( $x$ )	Next tokens ( $y$ )
Train	Repository 1	... from easyrec.gateway import EasyRec <EOL> gateway = EasyRec('tenant','key') <EOL> item_type = gateway. <span style="background-color: pink;">_____</span>	get_item_type
		... response = gateway.get_other_users() <EOL> get_params = HTTPretty. <span style="background-color: pink;">_____</span>	last_request
	Repository 2	import numpy as np ... <EOL> if np.linalg.norm(target - prev_target) > far_threshold: <EOL> norm = np. <span style="background-color: pink;">_____</span>	linalg
		... new_trans = np.zeros((n_beats + max_beats, n_beats)) <EOL> new_trans[:n_beats,:n_beats] = np. <span style="background-color: pink;">_____</span>	max
		...	
Test	Repository 6,001	... if e.errno == errno.ENOENT: <EOL> continue <EOL> p = subprocess.Popen () <EOL> stdout = p. <span style="background-color: pink;">_____</span>	communicate
		... command = shlex.split(command) <EOL> command = map(str, command) <EOL> env = os. <span style="background-color: pink;">_____</span>	environ
		...	

# Plan for Today

## Domain Generalization

- Problem formulation
- Algorithms
  - **Adding explicit regularizers**
- Data augmentation

## Goals for this lecture:

- Understand [domain generalization](#): intuition, problem formulation
- Familiarize mainstream DG approaches: [regularization](#)-based, [augmentation](#)-based

# How to Learn Generalizable Representations?

Why do machine learning models fail to generalize?

Goal: classify dog vs. cat

Domain 1:  
water



45% of train data



5% of train data

Domain 2:  
grass

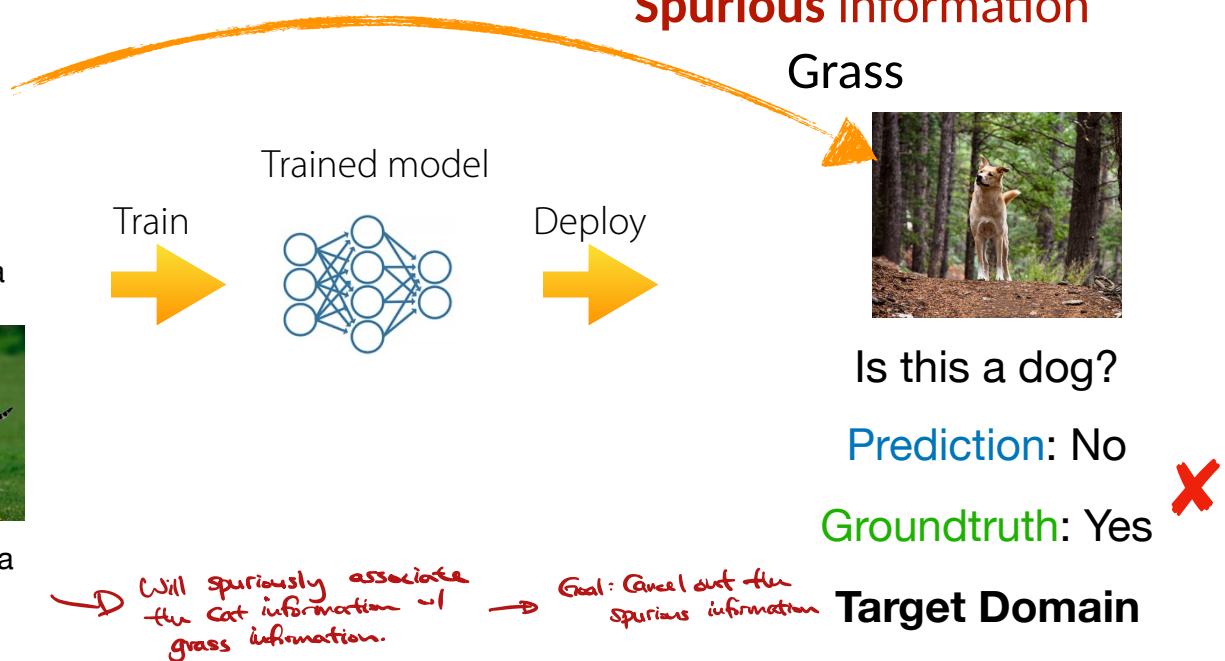


5% of train data



45% of train data

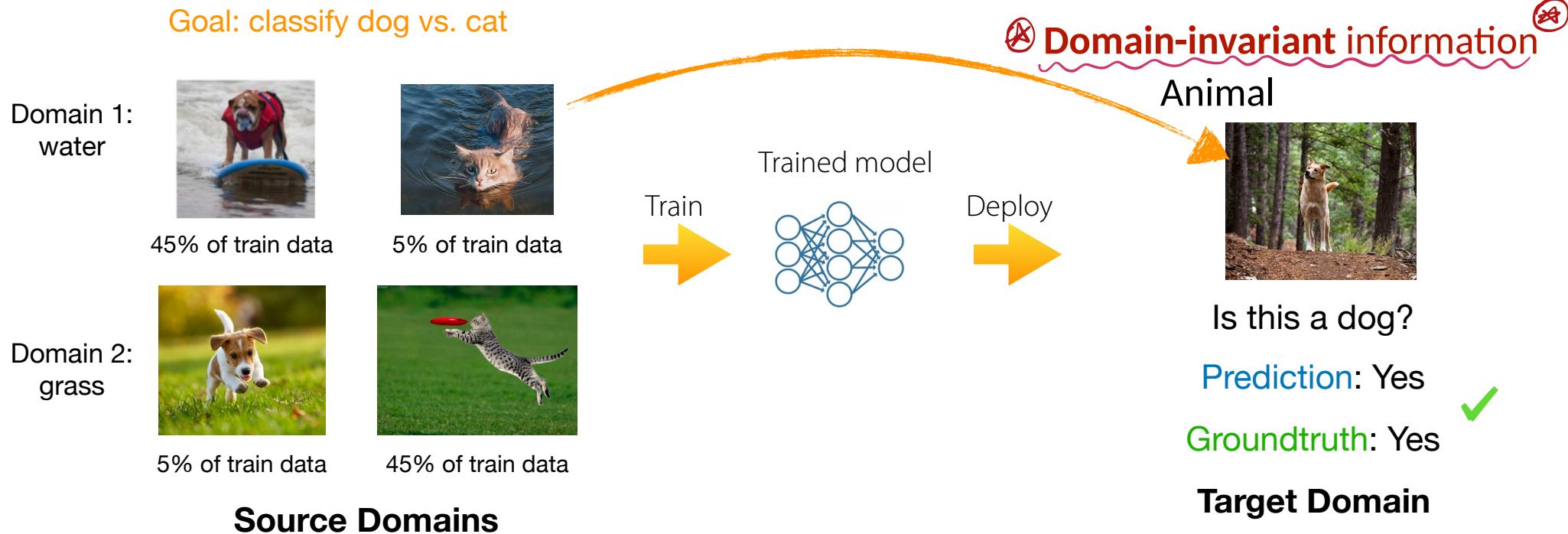
Source Domains



# How to Learn Generalizable Representations?

To overcome spurious correlation —> train a neural network to learn **domain invariance**

Domain invariance: we want to learn features that don't change across domains



# Regularization-based Method

~~Key idea: Use a regularizer to align representations across domains~~

→ get domain-invariant representation

Domain 1:  
water



45% of train data



5% of train data

Domain 2:  
grass

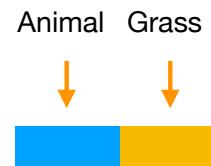
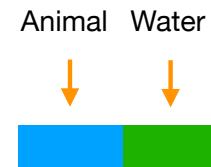


5% of train data



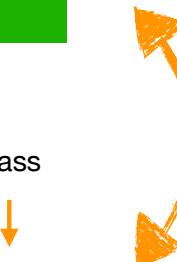
45% of train data

**Source Domains**



*More granular/  
specific representations*

**Representations**



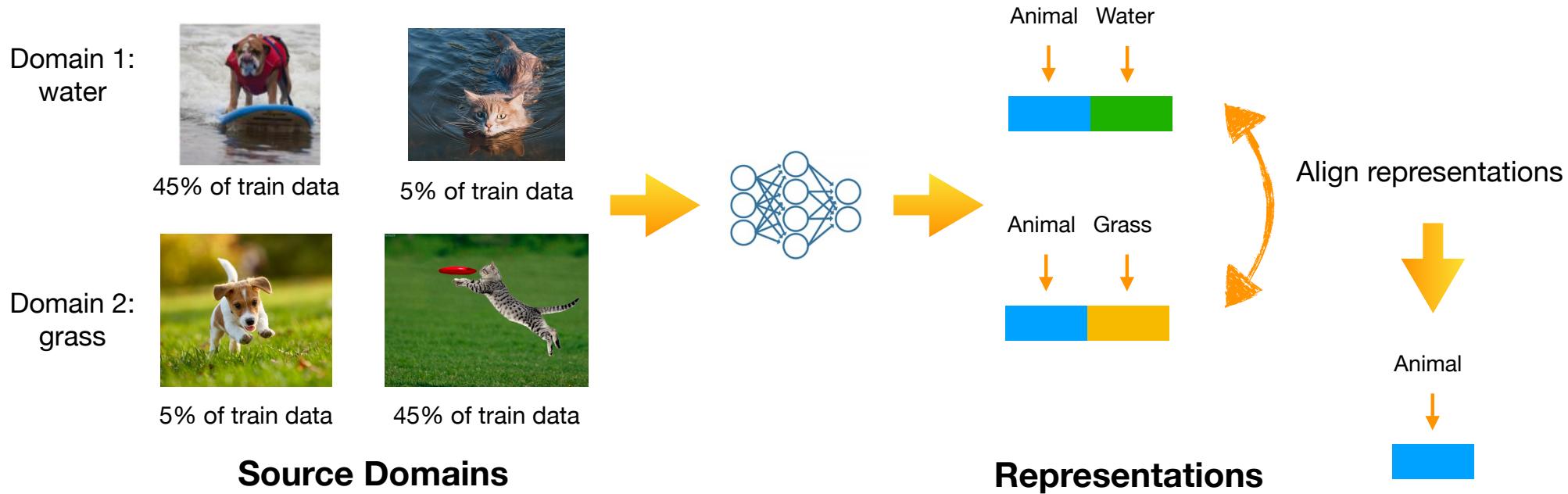
Align representations



Animal



# Regularization-based Method



## Source Domains

## Representations

Label classification loss

$$\min_{\theta} \mathbb{E}_{(x,y)}[\ell(f_{\theta}(x), y)] + \lambda \mathcal{L}_{reg}$$

Average over training examples

Explicit regularizer to learn  
domain-invariant representation

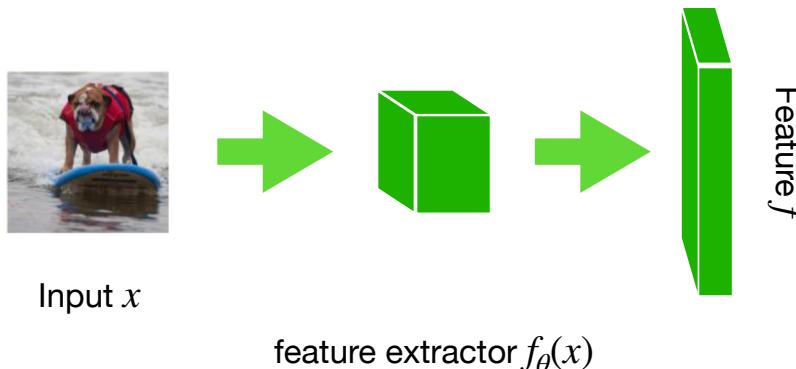
Q: How do we define  
such a regularizer?

# Recap: Domain Adversarial Training in DA

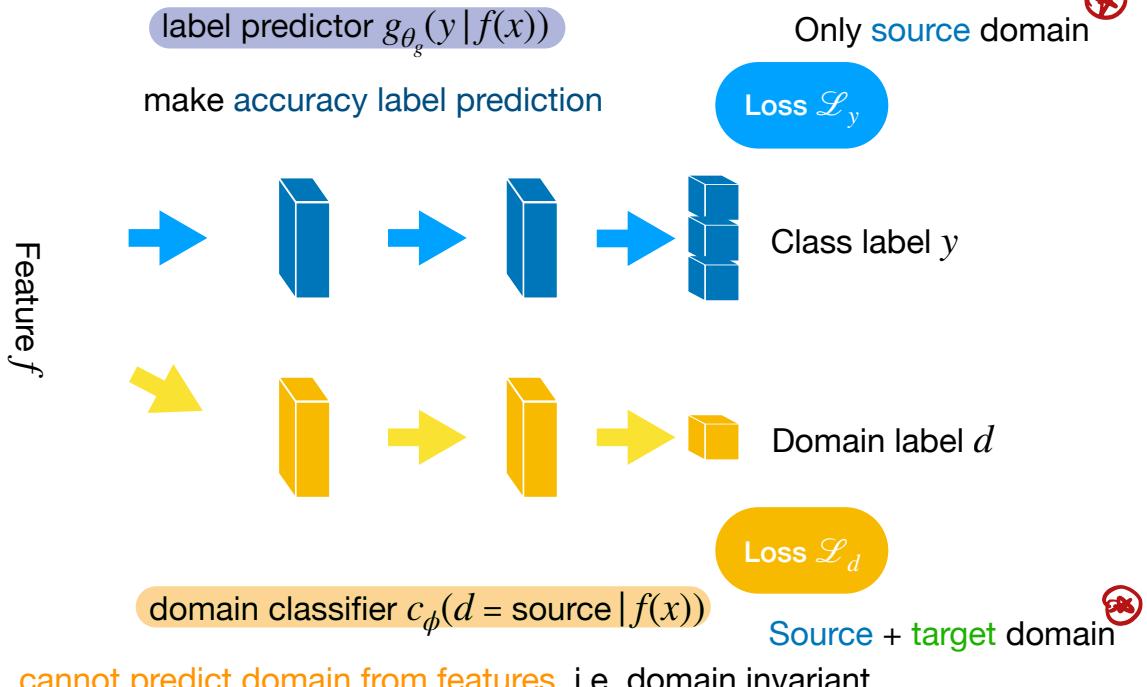
↳ Domain Adaptation  
(from last lecture)

**Key idea:** predictions must be made based on features that cannot be discriminated between the domains

Same architecture as p18  
of previous lecture's slides.



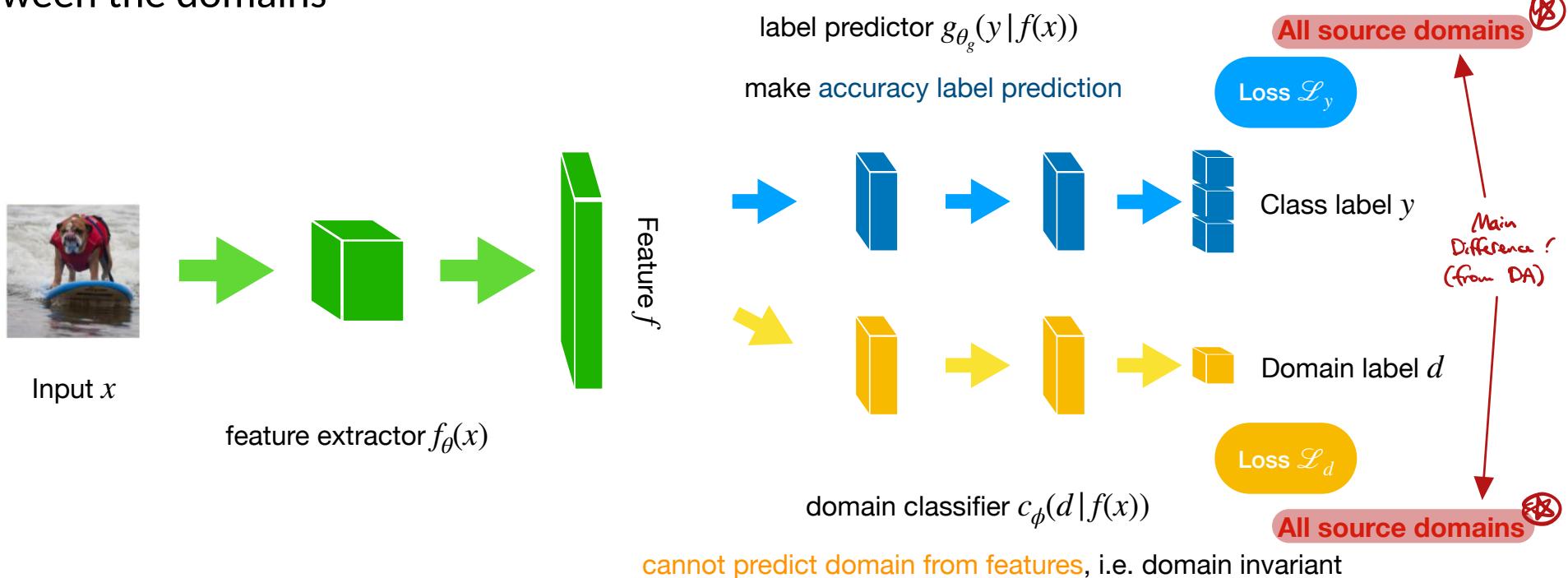
**Question:** Does anyone have ideas on how to use domain adversarial training in the domain generalization setting?



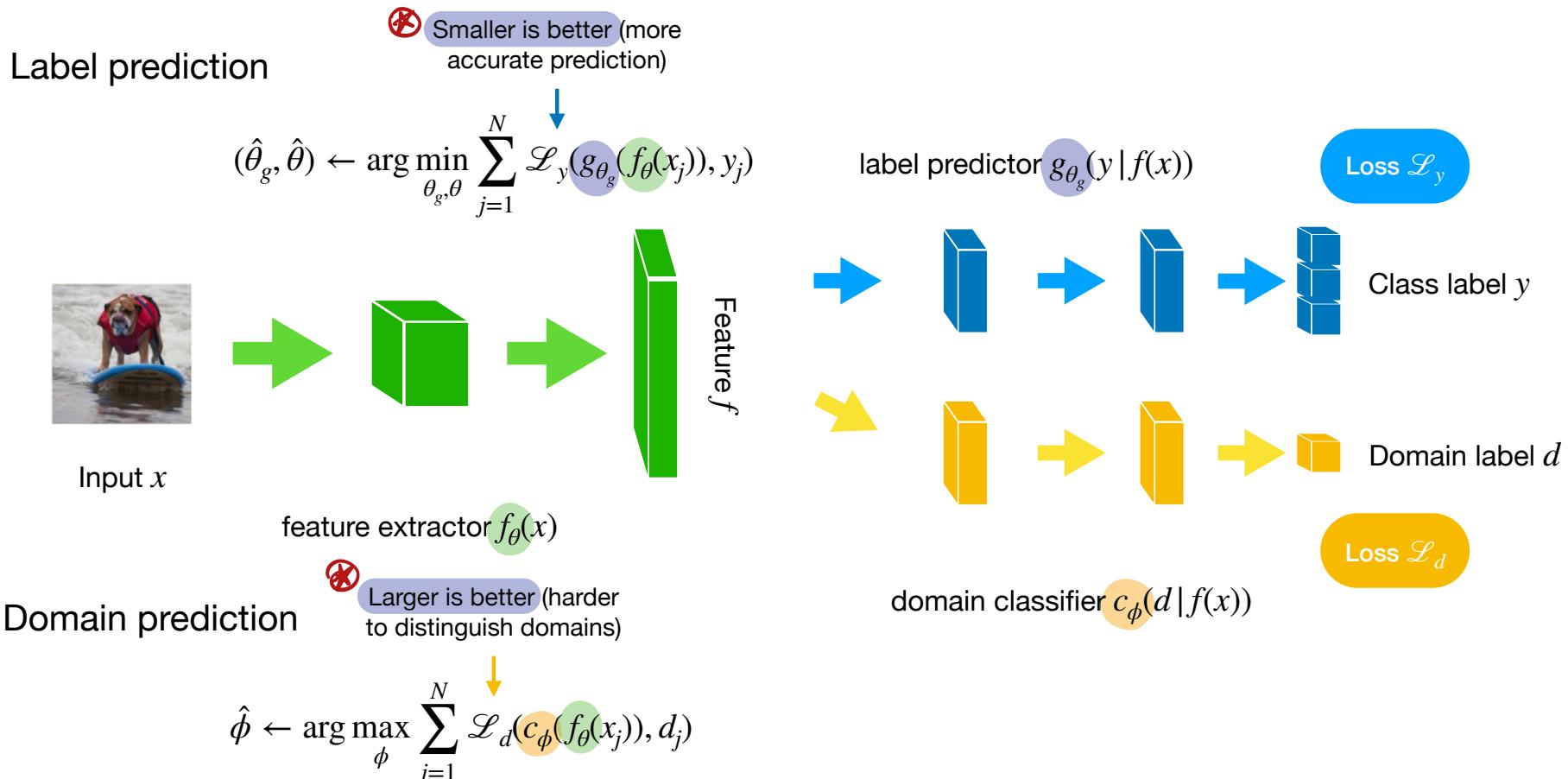
# Domain Adversarial Training in DG

↳ Domain Generalization

**Key idea:** predictions must be made based on features that cannot be discriminated between the domains



# Domain Adversarial Training in DG



# Domain Adversarial Training in DG

DANN loss in DG

Label classification loss

$$\min_{\theta} \mathbb{E}_{(x,y)}[\ell(f_{\theta}(x), y)] + \lambda \mathcal{L}_{reg}$$

Explicit regularizer to learn domain-invariant representation

Full algorithm

$$\mathcal{L} = \sum_{j=1}^N \mathcal{L}_y(g_{\theta_g}(f_{\theta}(x_j)), y_j) - \lambda \mathcal{L}_d(c_{\phi}(f_{\theta}(x_j)), d_j)$$

1. Randomly initialize encoder  $f_{\theta}$ , label classifier  $g_{\theta_g}$ , domain classifier  $c_{\phi}$
2. Update domain classifier:  $\min_{\phi} \mathcal{L} = \sum_{i=1}^n -\mathcal{L}_d(c_{\phi}(f_{\theta}(x_i)), d_i)$ 

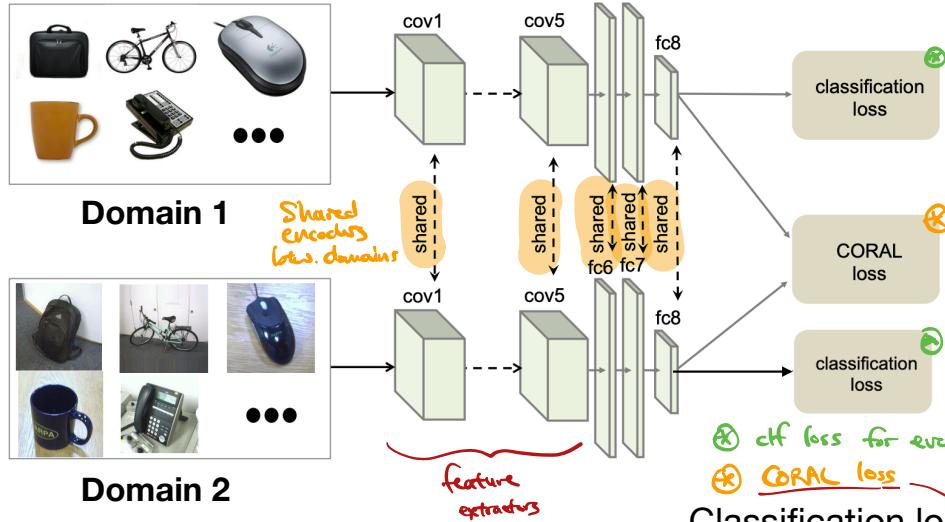
*minimize the (-)ve  $\leftrightarrow$  maximize the (+)ve*
3. Update label classifier & encoder:  $\min_{\theta, \theta_g} \mathcal{L} = \sum_{i=1}^n \mathcal{L}_y(g_{\theta_g}(f_{\theta}(x_i)), y_i) - \lambda \mathcal{L}_d(c_{\phi}(f_{\theta}(x_i)), d_i)$
4. Repeat steps 2 & 3  $\rightarrow$  Until convergence

Are there any other ways to **learn domain-invariant features without adversarial optim**?

# Alternative Approach — CORAL

**Key idea:** directly aligning representations between different domains with some similarity metrics

**CORAL:** Correlation Alignment for Domain Adaptation (usually also used in DG)



You can change the model architecture to improve expressive power, but typically applying CORAL will give better results in this setting.

Notations

$k$ : num of features

Calculate covariance matrices



Domain 2

Classification loss

$\text{CORAL loss}$

Directly align between the representations

$$\mathcal{L} = \sum_{j=1}^{n_1+n_2} \mathcal{L}_c(f_\theta(x_i), y_i) + \lambda \mathcal{L}_{coral}$$

$$\mathbf{X}_1 \in \mathbb{R}^{n_1 \times k}$$

$$\mu_1 = \frac{1}{n_1} \mathbf{1}^T \mathbf{X}_1 \in \mathbb{R}^{1 \times k}$$

$$C_1 = \frac{1}{n_1 - 1} \sum_{i=1}^{n_1} (\mathbf{X}_1 - \mu_1)^T (\mathbf{X}_1 - \mu_1)$$

$$C_2 = \frac{1}{n_2 - 1} \sum_{i=1}^{n_2} (\mathbf{X}_2 - \mu_2)^T (\mathbf{X}_2 - \mu_2)$$

$$\mathcal{L}_{coral} = \frac{1}{4k^2} \|C_1 - C_2\|_F^2$$

Try to match the covariance matrices  
⇒ Make the cov. matrix match for every domain

Explicit regularizer to learn domain-invariant representation  
⇒ Extensible to more than 2 domains

Use pairwise CORAL loss

# Results

↗ Empirical Risk Minimization  
 ↗ General form of Max. Likelihood Est. (MLE)

ERM

CORAL

DANN

OfficeHome

↗ Domains  
 ↗ 3 → Train  
 ↗ 1 → Test



66.5%

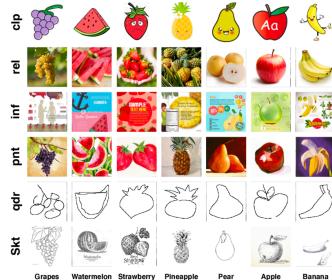
<

**68.7%**

>

65.9%

DomainNet



40.9%

<

**41.5%**

>

38.3%

iWildCam



30.8%

<

**32.7%**

>

n/a

↗ Designing a good regularizer  
 is v. important.

# Pros and Cons of Regularization-based Methods

+ General to all kinds of data and networks

Domain 1:  
water



+ Some theoretical guarantee

45% of train data



- The regularizer being too harsh / too constraining on the representation

Domain 2:  
grass



5% of train data



45% of train data

$$\min_{\theta} \mathbb{E}_{(x,y)}[\ell(f_{\theta}(x), y)] + \lambda \mathcal{L}_{reg}$$

Explicit regularizer encourages internal representation to contain no info about the background

But in some cases we may need some info to do this clf, or sth similar. Or, if we don't include the bg info, it may hurt us to learn better domain representations/expressive power of NN.

# Pros and Cons of Regularization-based Methods

+ General to all kinds of data and networks

+ Some theoretical guarantee

- The regularizer being too harsh / too constraining on the representation

These methods can help the performance,  
but do not always work



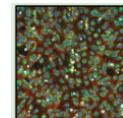
iWildCam

Empirical Risk  
Minimization

30.8%

32.7%

CORAL



RxRx1

29.9%

28.4%

CORAL

Are there any other approaches to relax the dependency of the regularizer?

→ Data  
Augmentation

# Plan for Today

## Domain Generalization

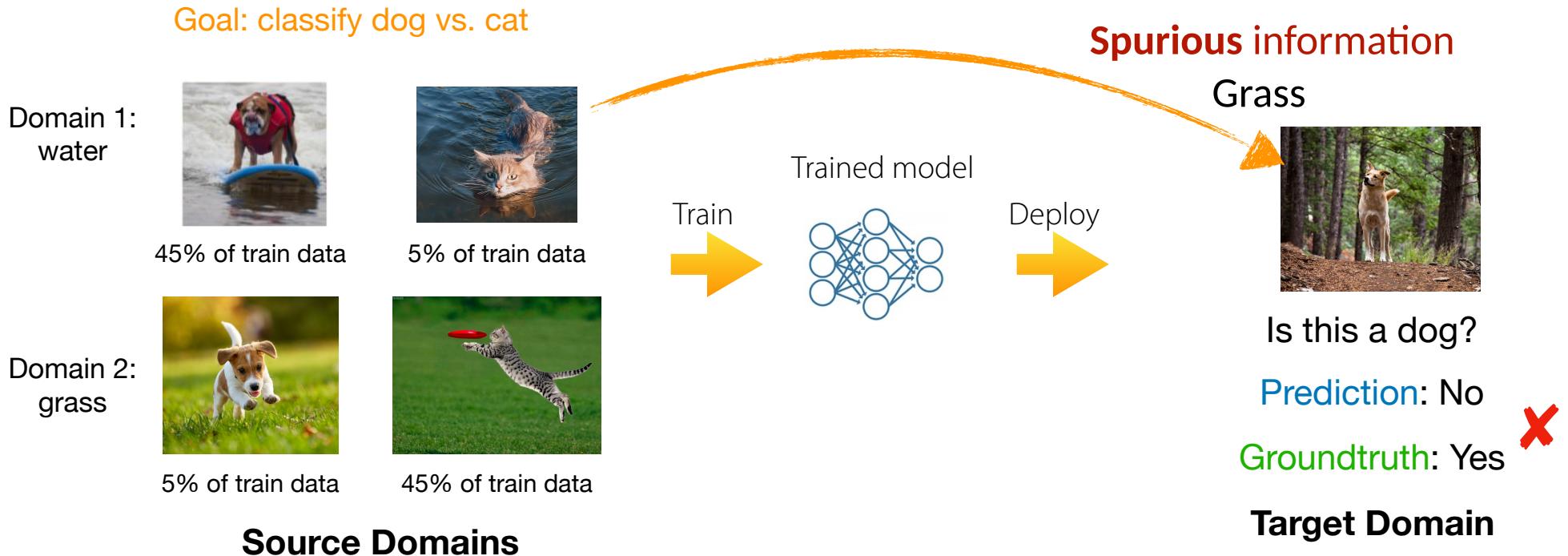
- Problem formulation
- Algorithms
  - Adding explicit regularizers
- **Data augmentation**

## Goals for this lecture:

- Understand [domain generalization](#): intuition, problem formulation
- Familiarize mainstream DG approaches: [regularization](#)-based, [augmentation](#)-based

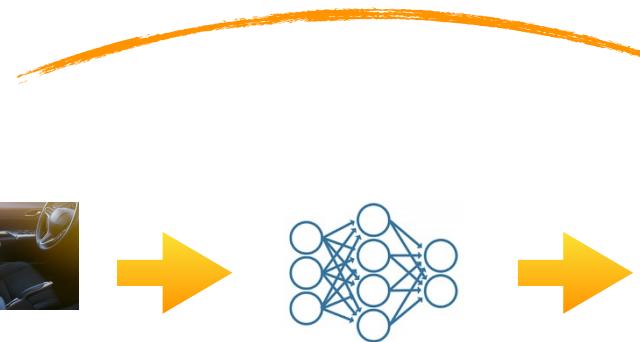
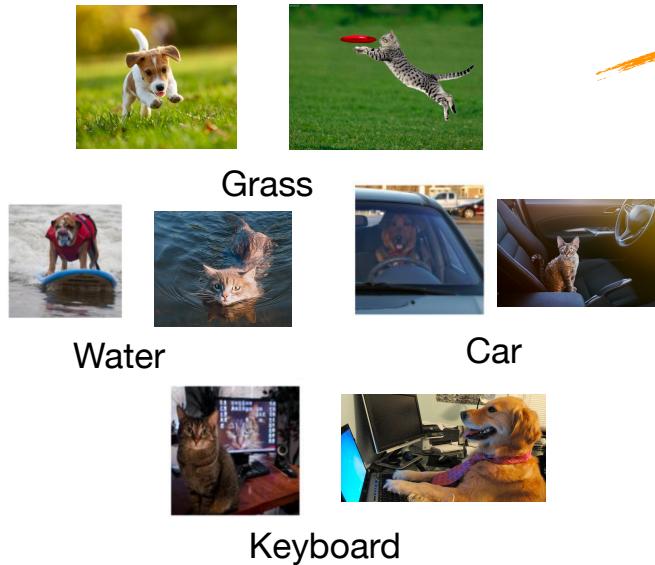
# Recap: Spurious Correlation

**Recap:** spurious correlation between domains and labels



# Data Augmentation

✗ If we can collect more data



Question: Will the network still associate dogs with water background in source domains?

✗ NO! There are many more backgrounds. We can't recognize dogs only with grass background.

→ The model will become more robust  
to different types of backgrounds.  
(w more, diverse data)



Is this a dog?

Prediction: Yes



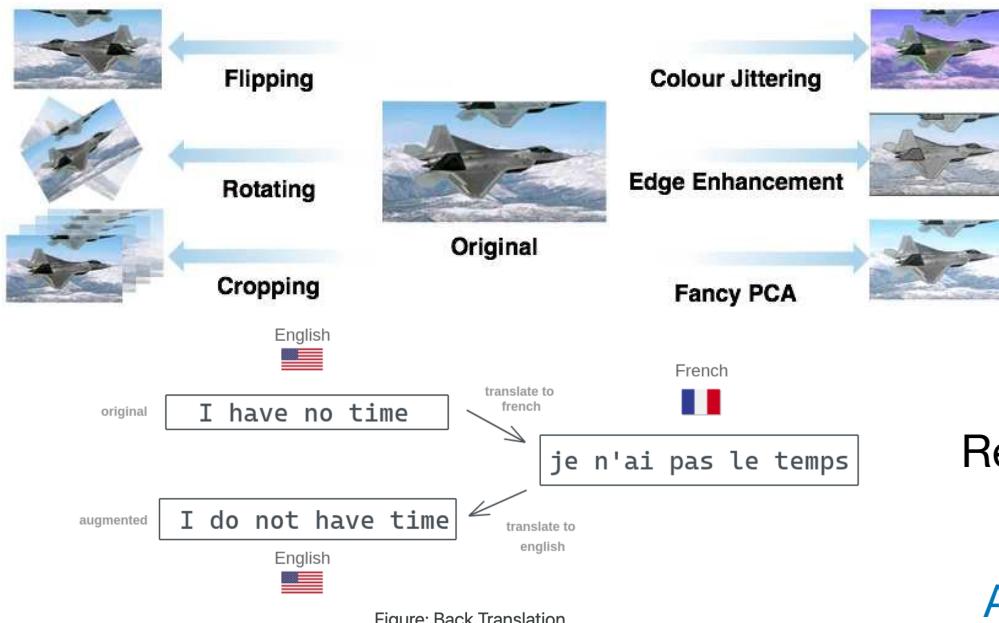
Groundtruth: Yes

Target Domain

✗ We can not collect more data → Let's generate data! → Augmentation

# Data Augmentation

Generating data with **simple operators**



Requires knowledge of the problem domain

Any general approaches?

<https://amitness.com/2020/02/back-translation-in-google-sheets/>

# Data Augmentation — Mixup

✗ **Interpolating** training examples *→ Key Idea!*

A learning model

$$\mathcal{D}_{tr} = \{x_i, y_i\}_{i=1}^N \rightarrow \text{Classifier},$$

Mixup

$$\tilde{\mathcal{D}}_{tr} = \{\tilde{x}_i, \tilde{y}_i\}_{i=1}^N \rightarrow \text{Classifier},$$

where

$$\begin{aligned} \text{✗ } \tilde{x}_i &= \lambda \textcolor{blue}{x}_i + (1 - \lambda) \textcolor{green}{x}_j, \tilde{y}_i = \lambda \textcolor{blue}{y}_i + (1 - \lambda) \textcolor{green}{y}_j \text{ ✗ } \\ \lambda &\sim \text{Beta}(\alpha, \beta) \end{aligned}$$

Generating some virtual  
examples between two  
classes

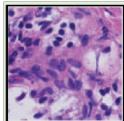


# Data Augmentation — Mixup

Mixup can improve the performance on domain generalization



Empirical Risk Minimization



70.3%

mixup

71.2%

Camelyon17



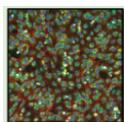
FMoW

32.8%

34.2%

→ CIF land type  
→ generalize for  
year ~ region

But it is not always good!



RxRx1

29.9%

26.5%

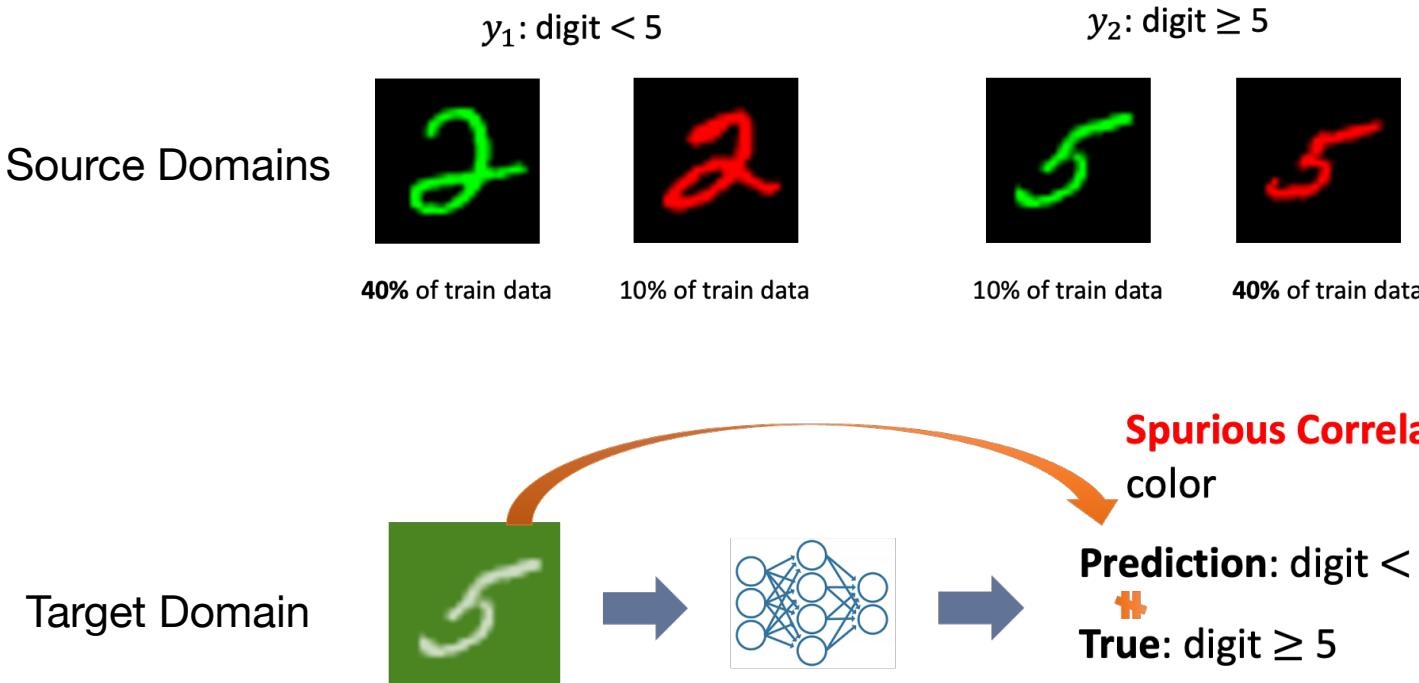
Original mixup only focuses  
on data augmentation instead  
of learning domain invariance.

How to Improve it?

# Data Augmentation — Mixup

→ Do these methods apply  
in the text domain as well?  
⇒ Yes, will discuss later.

A simpler example with spurious correlation

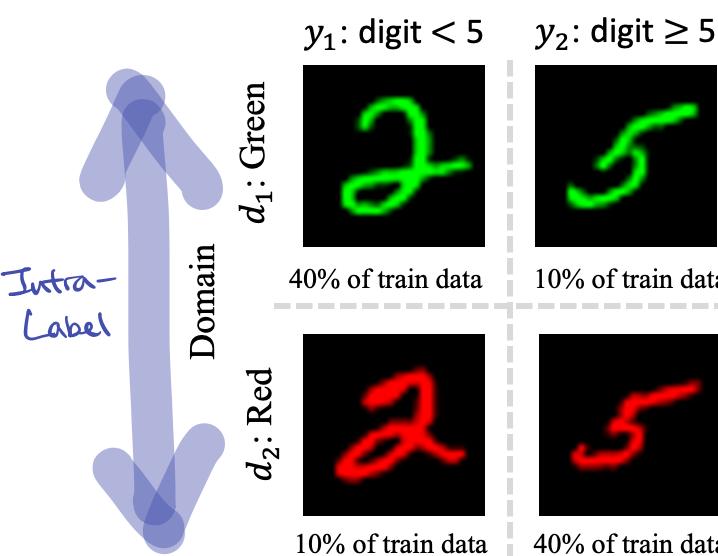


Hard to apply when  
we don't have the  
domain labels

# Can we Improve Mixup? — LISA

↳ Learning Invariant Predictors  
w/ Selective Augmentation.

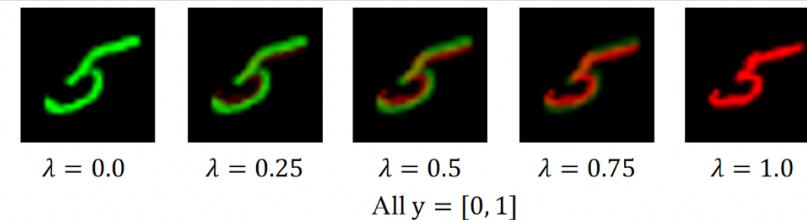
Key idea: selective interpolate examples to emphasize invariant information



Colored MNIST

Mixup:  $x_{mix} = \lambda \mathbf{x}_i + (1 - \lambda) \mathbf{x}_j, y_{mix} = \lambda \mathbf{y}_i + (1 - \lambda) \mathbf{y}_j$   
 $\lambda \sim \text{Beta}(\alpha, \beta)$

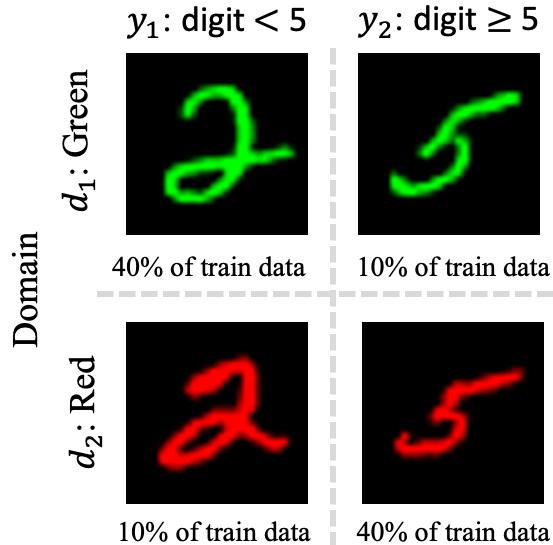
Intra-label LISA — Interpolates samples with the same ~~label~~ but different domains ( $d_i \neq d_j, y_i = y_j$ )  
Key Idea



Different background, same label

# Can we Improve Mixup? — LISA

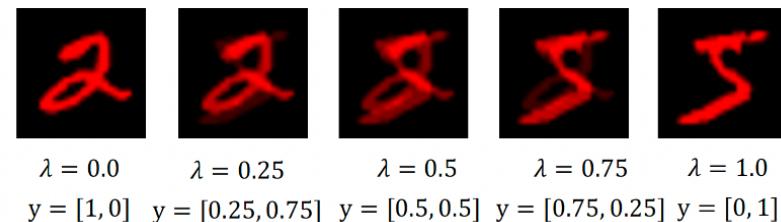
**Key idea:** selective interpolate examples to emphasize invariant information



Mixup:  $x_{mix} = \lambda \mathbf{x}_i + (1 - \lambda) \mathbf{x}_j, y_{mix} = \lambda \mathbf{y}_i + (1 - \lambda) \mathbf{y}_j$   
 $\lambda \sim \text{Beta}(\alpha, \beta)$

$\neq$  Intra-Label LISA (prev. slide)

✳️ Intra-domain LISA – Interpolates samples with the different label but same domains ( $d_i = d_j, y_i \neq y_j$ ) ✳️



✳️✳️ Domain information is not the reason for the label change ✳️✳️

↳ Basically why this works

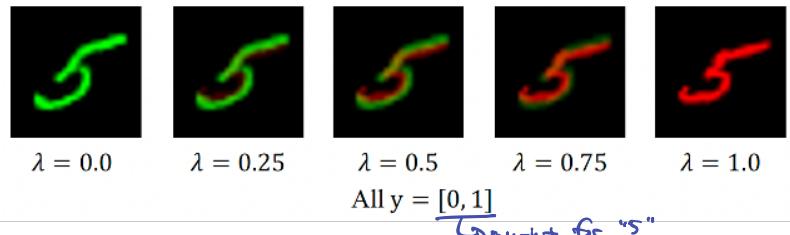
# Can we Improve Mixup? — LISA

In practice, LISA combines the two  
→ Try to get the best of both worlds.

Intra-label  
LISA

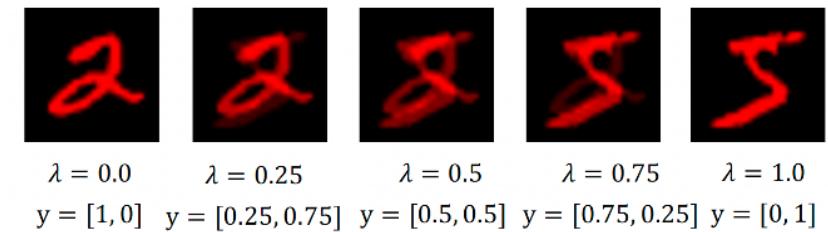
= Inter-domain

Which one  
Should you  
use?



Intra-domain  
LISA

= Inter-label



+ more domains

+ spurious correlations  
are not very strong

+ domain information is highly  
spuriously correlated with the  
label

$p_{sel}$ :

Determine intra-label LISA or intra-domain LISA at each iteration

# Full Algorithm of LISA

1. Randomly initialize the model parameter  $\theta$
2. Sample strategy  $s \sim Bernoulli(p_{sel})$
3. Sample a batch of examples  $\mathcal{B}$ 
  - (i) If  $s=0$ , for each example  $(x_i, y_i)$  in  $\mathcal{B}$ , sample  $(x_j, y_j)$  that satisfies  $(y_i = y_j)$  and  $(d_i \neq d_j)$   
Same label      diff. domain
  - (ii) If  $s=1$ , for each example  $(x_i, y_i)$  in  $\mathcal{B}$ , sample  $(x_j, y_j)$  that satisfies  $(y_i \neq y_j)$  and  $(d_i = d_j)$   
diff. label      Same domain
4. Use interpolated examples to update the model
5. Repeat steps 3 & 4

← Intra-label LISA

← Intra-domain LISA

→ Applicable to multiple domains + multiple classes.

# Results

		ERM	Regularization-based (CORAL)	Augmentation-based (LISA)
Camelyon17		70.3%	74.7%	<b>77.1%</b>
FMoW		32.3%	34.6%	<b>35.5%</b>
RxRx1		29.9%	28.4%	<b>31.9%</b>
Amazon		53.8%	53.8%	<b>54.7%</b>
iWildCAM		30.8%	<b>32.7%</b>	27.6%
OGB-MolPCBA		<b>28.3%</b>	17.9%	27.5%

LISA can also work on text data, how to apply mixup?

# Manifold Mixup

Original Mixup

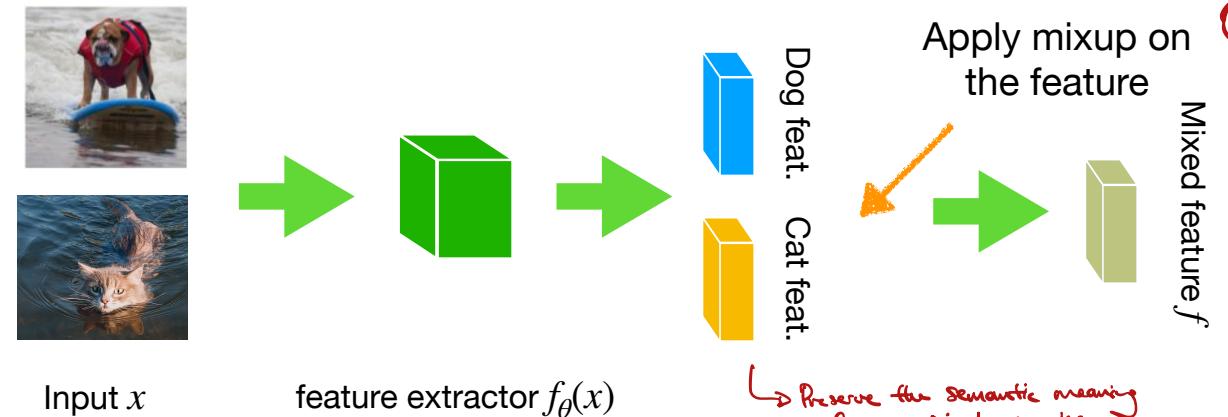
Apply mixup on  
the input



feature extractor  $f_{\theta}(x)$

Apply mixup on  
the feature

Manifold Mixup  
Term used in topology  
Apply mixup @ the  
feature level!



Preserve the semantic meaning  
of the original examples.  
Can make the NN more robust  
to adversarial examples.

Applicable to  
image clf, NLP,  
machine translation

# Invariance Analysis

Metrics:

Accuracy of domain prediction      Divergence of predictions among domains

Use this to measure whether we can find a better domain invariance or not

→ the less is this?

→ Proxy measure for domain invariance  $\textcircled{R}$  How does this work?

	IP <sub>adp</sub> ↓				IP <sub>kl</sub> ↓			
	CMNIST	Waterbirds	Camelyon17	MetaShift	CMNIST	Waterbirds	Camelyon17	MetaShift
ERM	82.85%	94.99%	49.43%	67.98%	6.286	1.888	1.536	1.205
Vanilla mixup	92.34%	94.49%	52.79%	69.36%	4.737	2.912	0.790	1.171
IRM	69.42%	95.12%	47.96%	67.59%	7.755	1.122	0.875	1.148
IB-IRM	74.72%	94.78%	48.37%	67.39%	1.004	3.563	0.756	1.115
V-REX	63.58%	93.32%	61.38%	68.38%	3.190	3.791	1.281	1.094
<b>LISA (ours)</b>	<b>58.42%</b>	<b>90.28%</b>	<b>45.15%</b>	<b>66.01%</b>	<b>0.567</b>	<b>0.134</b>	<b>0.723</b>	<b>1.001</b>

Smaller is better  
→ Less accurate domain prediction

Smaller is better

**LISA leads to greater domain invariance than prior methods with explicit regularizers**

→ Harder to tell which domain the example is coming from

# Regularization-based v.s. Augmentation-based Methods

## Regularization-based Method

- + General to all kinds of data and networks
- + Some theoretical guarantee
- Rely on the design of regularizers

## Augmentation-based Method

- + Easy to understand and simple to implement
- + No need to worry about how to design regularizers
- Largely limited to classification

# Plan for Today

## Domain Generalization

- Problem formulation
- Algorithms
  - Adding explicit regularizers
  - Data augmentation

## Goals for this lecture:

- Understand **domain generalization**: intuition, problem formulation
- Familiarize mainstream DG approaches: **regularization-based**, **augmentation-based**

# Reminders

Project milestone on **Wednesday, November 16**

Homework 4 (optional) due **Monday, November 14**

**Next time:** Lifelong learning