Deep Model Adaptation using Domain Adversarial Training



Skolkovo Institute of Science and Technology (Skoltech) Moscow region, Russia

Victor Lempitsky, *joint work with* Yaroslav Ganin

Deep supervised neural networks



are a "big thing" in computer vision and beyond are hungry for labeled data





Where to get the data?

- Biomedical
- Unusual cameras / image types
- Videos
- Data with expert-level annotation (not mTurkable)
- Surrogate training data often available: Borrow from adjacent modality • Generate synthetic imagery (computer graphics) Use data augmentation to amplify data (imagebased rendering, morphing, re-synthesis,...)

Resulting training data are shifted. Domain adaptation needed.

Deep Model Adaptation using Domain Adversarial Training



Lots of modalities do not have large labeled data sets:

Example: Internet images -> Webcam sensor



Example: (semi-)synthetic to real





Assumptions and goals

(e.g. synthetic images) (e.g. real images) Goal: train a deep neural net that does well on the target domain

Large-scale deep unsupervised domain adaptation

Deep Model Adaptation using Domain Adversarial Training



• Lots of *labeled* data in the source domain • Lots of unlabeled data in the target domain



Domain shift in a deep architecture

feature extractor

When trained on source only, feature distributions do not match:

X



Idea 1: domain-invariant features wanted

Feature distribution without adaptation:



Our goal (after adaptation):





Idea 2: measuring domain shift





Domain loss low

Deep Model Adaptation using Domain Adversarial Training

Domain loss high

Learning with adaptation



1. Build this network 2. Train feature extractor + class predictor on source data 3. Train feature extractor + domain classifier on source+target data 4. Use feature extractor + class predictor at test time

Idea 3: minimizing domain shift



Emerging features: Discriminative (good for predicting y)

Deep Model Adaptation using Domain Adversarial Training

Domain-discriminative (good for predicting d)



LOSS

Idea 3: minimizing domain shift



Gradient reversal layer: • Multiplies the gradient by $-\lambda$ at backprop

Deep Model Adaptation using Domain Adversarial Training

Copies data without change at forwardprop



loss

Idea 3: minimizing domain shift



Emerging features: Discriminative (good for predicting y)

Deep Model Adaptation using Domain Adversarial Training

Domain-invariant (not good for predicting d)



LOSS

Gradient reversal layer

class GradReversalLayer : Layer {

- float lambda;
- blob forward (blob x) return x

Deep Model Adaptation using Domain Adversarial Training

blob backward(blob dzdy) { return multiply(*dzdy*, -*lambda*)

Saddle point interpretation

Our objective (small label prediction loss + large domain classification loss wanted)



The backprop converges to a saddle point:

 $\hat{\theta}_d = \arg\max_{\theta_d} E(\hat{\theta}_f, \hat{\theta}_y, \theta_d).$

Similar idea for generative networks: In NIPS, 2014]

Deep Model Adaptation using Domain Adversarial Training

$E(\theta_f, \theta_y, \theta_d) = \sum L_y^i(\theta_f, \theta_y) - \lambda \sum L_d^i(\theta_f, \theta_d)$ i = 1..N

$(\hat{\theta}_f, \hat{\theta}_y) = \arg\min_{\theta_f, \theta_y} E(\theta_f, \theta_y, \hat{\theta}_d)$

[Goodfellow et al. Generative adversarial nets.

Backprop updates Label prediction loss for the *i*th example



Domain classification loss for the *i*th example $\theta_f - \mu \left(rac{\partial L_y^i}{\partial \theta_f} \right)$ $\mathcal{P}^{\mathcal{O}}$ $\mu \frac{\partial L_d^i}{\partial \theta_d}$

Initial experiments: baselines



Shallow adaptation baseline: [Fernando et al.,

Lower bound: training on source domain only

Deep Model Adaptation using Domain Adversarial Training

Upper bound: training on target domain with labels

Unsupervised visual domain adaptation using subspace alignment. ICCV, 2013] applied to the last-but-one layer



Example: from synthetic to real



"Windows digits"

0,93	
0,92	-
0,91	
0,9	-
0,89	-
0,88	4
0,87	-
0,86	-
0,85	-
0,84	_
0,83	-





"House numbers"

			-	
		11011		/
			-	
			-	
			-	
	_		-	

Baseline Deep adapt Upper bound

Example: large gap



"House numbers	, " 1	
	0,9	
	0,8	
	0,7	
Reverse	0,6	
direction does	0,5	
not work ③	0,4	
		No adar

Deep Model Adaptation using Domain Adversarial Training



No adaptBaselineDeepUpperadaptbound





Traffic signs: semi-supervised adaptation



• 43 classes 430 Real examples Testing on Real only

erro]

Validation

Sample architectures for image classification



(a) MNIST architecture; inspired by the classical LeNet-5 (LeCun et al., 1998).



SVHN architecture; adopted from Srivastava et al. (2014). (b)



(c) GTSRB architecture; we used the single-CNN baseline from Cireşan et al. (2012) as our starting point.

Office dataset



Results on Office dataset

Method

GFK(PLS, PCA) (Gong et al., SA* (Fernando et al., 2013) DLID (Chopra et al., 2013) DDC (Tzeng et al., 2014) DAN (Long and Wang, 2015) SOURCE ONLY DANN

Most similar approach (matches means of distributions): [Tzeng et al. Deep domain confusion: Maximizing for domain invariance. CoRR, abs/1412.3474, 2014]

Л

OURCE	Amazon	DSLR	WEBG
ARGET	Webcam	Webcam	DSI
2012)	.197	.497	.665
	.450	.648	.69
	.519	.782	.89
	.618	.950	.98
	.685	.960	.99
	.642	.961	.97
	.730	.964	.99



Beyond image classification Domain-Adversarial Training of Neural Networks Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, Victor Lempitsky, JMLR 2016

		Original data		mSDA representation			
Source	Target	DANN	NN	SVM	DANN	NN	SVM
BOOKS	DVD	.784	.790	.799	.829	.824	.830
BOOKS	ELECTRONICS	.733	.747	.748	.804	.770	.766
BOOKS	KITCHEN	.779	.778	.769	.843	.842	.821
DVD	BOOKS	.723	.720	.743	.825	.823	.826
DVD	ELECTRONICS	.754	.732	.748	.809	.768	.739
DVD	KITCHEN	.783	.778	.746	.849	.853	.842
ELECTRONICS	BOOKS	.713	.709	.705	.774	.770	.762
ELECTRONICS	DVD	.738	.733	.726	.781	.759	.770
ELECTRONICS	KITCHEN	.854	.854	.847	.881	.863	.847
KITCHEN	BOOKS	.709	.708	.707	.718	.721	.769
KITCHEN	DVD	.740	.739	.736	.789	.789	.788
KITCHEN	ELECTRONICS	.843	.841	.842	.856	.850	.861
	(a) Classification accuracy on the Amazon reviews data set						

Adaptation for Person Re-identification



VIPER

VIPER to CUHK

Deep Model Adaptation using Domain Adversarial Training











PRID

CUHK



Adaptation for Person Re-identification







PRID

CUHK



Caveats

- Domains should not be too far apart
- Early on, the gradient from the domain classification loss should not be too strong
- The trick used to obtain the results: gradually increase λ from 0 to 1



Deep Model Adaptation using Domain Adversarial Training

too far apart rom the domain d not be too strong the results: gradually

Conclusion

- Scalable method for deep unsupervised domain adaptation
- Based on simple idea. Takes few lines of code (+ defining a specific network architecture). Caffe implementation available.
- State-of-the-art results
- Unsupervised parameter tuning is easy (look at the domain classifier error)
- Main challenge: initialization and stepsize

http://sites.skoltech.ru/compvision/projects/qrl/