# Unsupervised Learning

Computer Vision – Lecture 18

### **Further Reading**

- Slides from <u>S Savarese, A Zamir</u>
- Slides from <u>F Li</u>
- Slides from <u>Y Asano</u>

### **Basics: Supervised Learning**

Dataset  $D = \{(x_i, y_i) | 1 \le i \le N\}$ Inputs  $x_i$ Outputs  $y_i$ Training/Validation/Testing  $D = D_T \cup D_V \cup D_*$ 

Learn f(x) = y by minimizing  $\sum_{(x_i, y_i) \in D_T} \mathcal{L}(f(x_i), y_i)$ 

Hope to generalise:  $\sum_{(x_i,y_i)\in D_*} \mathcal{L}(f(x_i), y_i)$ 

### Basics: Unsupervised Learning

Dataset  $D = \{x_i | 1 \le i \le N\}$ Inputs  $x_i$ Outputs  $y_i$ Training/Validation/Testing  $D = D_T \cup D_V$ Testing on downstream task  $T = \{(\chi_i, v_i) | 1 \le i \le M\}$ 

Learn f(x) = ? by minimizing ??

Hope to generalise to another task  $\sum_{(\chi_i, v_i) \in T} \mathcal{L}(f(\chi_i), v_i)$ 

### What we need to do

- Trick the model to learn the downstream task without labels
- Build inductive biases into the model
- Find a learning signal for the training scheme
- Often: prevent trivial solutions and cheating

- Learning signals are general ideas to incorporate priors
- Priors model general assumptions about the world/task

This lecture contains a wide set of tools for learning from priors

- Can be used in all settings un/weakly/fully-supervised
- Not specific to tasks, apply to many areas
- For free\*! (no annotations needed)

- Recovery:  $f(M(x)) \leq x$
- Bottleneck:  $f(g(x)) \leq x$  with restriction on g(x)
- Dataset:  $f(x_1) <-> f(x_2)$
- Invariance:  $f(\pi(x)) \leq f(x)$
- Equivariance:  $f(\pi(x)) \leq \pi'(f(x))$  often  $\pi \equiv \pi'$  but not always
- Transformation estimation:  $f(\pi(x,\theta)) \leq \theta$
- Generative: *f*(*z*) <-> *D*
- Task-specific: f(x) <-> priors
- Uncertainty:  $f(x) \leq 0$  own error
- Many more!

### Reconstruction



### Learning Signal Toolbox Recovery



(Context autoencoder [Pathak et al.; CVPR '17], denoising autoencoder [Vincent et al.; ICML '08], Diffusion models [Sohl-Dickstein et al., ICML '15], etc.)

### **Rcovery:** Inpainting



D. Pathak, P. Krahenbuhl, J. Donahue, T. Darrell, A. Efros. <u>Context Encoders: Feature Learning by Inpainting</u>. CVPR 2016

### **Recovery: Colorization**



### **Colorization:** Architecture



### Failure Cases



<sup>13</sup>Source: A. Efros, R. Zhang

### Inherent Ambiguity



### Grayscale

Source: A. Efros, R. Zhang

### Inherent Ambiguity



Prediction



### Ground Truth

<sup>15</sup>Source: A. Efros, R. Zhang

### Biases



Source: A. Efros, R. Zhang

### Biases



### Learning Signal Toolbox Bottleneck



...

### Learning Signal Toolbox Equivariance



\* Often  $\pi \equiv \pi'$ 

 $\pi'^{-1}(f(\pi(x))) = f(x)$ is often more difficult because inverse is hard

 $f(\pi(x)) = \pi'(f(x))$ 

### **Transformation Estimation**



#### Rotation [Gidaris et al.; ICLR '18]



Jigsaw puzzle [Noroozi & Favaro; ECCV '16] 20

### Context prediction

- Pretext task: randomly sample a patch and one of 8 neighbors
- Guess the spatial relationship between the patches

Question 1:



A: Bottom right







Question 2:



A: Top center

### **Context prediction: Details**

Prevent "cheating": sample patches with gaps, pre-process to overcome chromatic aberration

AlexNet-like architecture

softmax



# Jigsaw puzzle solving

#### Crop out tiles



Claim: jigsaw solving is easier than context prediction, trains faster, transfers better

M. Noroozi and P. Favaro. Unsupervised Learning of Visual Representations by Solving Jigsaw Puzzles. ECCV 2016

### Jigsaw puzzle solving: Details



Predetermined set of 1000 permutations (out of 362,880 possible)

### Generative



#### Autoregressive



[van den Oord et al.; ICML '16, NeurIPS '16]

#### Variational Autoencoder



[Kingma et al.; ICLR '14]

**Diffusion Model** 



[Sohl-Dickstein et al.; ICML '15]



\* Different definitions of "related" and "unrelated" samples exist "Contrastive"

### SimCLR



- Introduce nonlinear projection (g) between representation (h) and feature used for computing contrastive loss (z).
- Use large mini-batch size.

### Invariance



\* Often with strong augmentations, but without changing the identity of the image

### SimCLR: Augmentations



Figure 4. Illustrations of the studied data augmentation operators. Each augmentation can transform data stochastically with some internal parameters (e.g. rotation degree, noise level). Note that we *only* test these operators in ablation, the *augmentation policy used to train our models* only includes *random crop (with flip and resize), color distortion*, and *Gaussian blur*. (Original image cc-by: Von.grzanka)

Crop	33.1	33.9	56.3	46.0	39.9	35.0	30.2	39.2		- 50
Cutout	32.2	25.6	33.9	40.0	26.5	25.2	22.4	29.4		
Color	55.8	35.5	18.8	21.0	11.4	16.5	20.8	25.7		-40
Sopel	46.2	40.6	20.9	4.0	9.3	6.2	4.2	18.8		- 30
st trai	38.8	25.8	7.5	7.6	9.8	9.8	9.6	15.5		- 20
Blur	35.1	25.2	16.6	5.8	9.7	2.6	6.7	14.5		
Rotate	30.0	22.5	20.7	4.3	9.7	6.5	2.6	13.8		-10
	CLOB	Cutout	Color	Sobel	Noise	Blur	Rotate	NVerage		
2nd transformation										

*Figure 5.* Linear evaluation (ImageNet top-1 accuracy) under individual or composition of data augmentations, applied only to one branch. For all columns but the last, diagonal entries correspond to single transformation, and off-diagonals correspond to composition of two transformations (applied sequentially). The last column reflects the average over the row.

### DINO: Self-Distillation with No Labels

- Student-Teacher training.
- Teacher's weights are exponential moving average (EMA) of the student.
- Teacher sees global view.
- Student sees local view.
- Student tries to predict teacher's distribution.



## DINO – sharpening and centring

• Collapse: same prediction for all samples.

• Centring: 
$$f_T(x_i) - \frac{1}{N} \sum_{j=1}^N f_T(x_j)$$

- Sharpening: low temperature for teacher softmax.
- Loss: Entropy between student and teacher distributions.

### DINO features are popular



(a) Input images

Appearance



(b) Co-segmented objects and parts



(c) Input image pair



(d) Correspondences

[Amir et al.; CVPR '22]



(a) DINO [6] (b) LOST [45].

#### (c) TokenCut (ours)



(d) Attention maps associated to different patches [Wang et al.; CVPR '22]

32











Structure



Output

[Tumanyan et al.; CVPR '22]

### Evaluation

- Occasionally simple: when training aligns fully with the task
- Often: some processing is needed
- Bridging the final gap between model and task
- In a practical setting: unsupervised learning is just the beginning
- In research: how far do we get with as little supervision as possible?

### Unsupervised vs. Self-Supervised?

...



Andrew Davison @AjdDavison

They are the same, 1

Can anyone explain if there is a difference between unsupervised and self-supervised learning? To me they seem the same and I find myself using both terms interchangeably (I prefer unsupervised), but I feel like I'm confusing people who understand them to mean different things.

7:16 PM · Jun 14, 2022 · Twitter for Android

Unlike SSL, it's difficult to abbreviate unsupervised learning

Self supervised uses a different loss function

To be a bit cute, unsupervised is the term for older techniques that don't work. Self supervised is the term for newer methods that do work.

labels are derived (perhaps implicitly) using problem-specific principles

"self-supervised" when the code looks exactly like supervised

There is no common definition (and someone will always complain)!

### Unsupervised vs. Weakly Supervised

Weak supervision means supervision but for a different task

- Segmentation from bounding boxes/captions/classes
- (Dense) depth from stereo
- Captioning from object labels
- Human in the loop annotations
- Objects from sound

### **Unsupervised Image Classification**





### **Unsupervised Image Classification**



### Hungarian Matching

- Find the lowest cost 1-to-1 assignments between *N* clusters and *N* labels.
- Cost: a *N* × *N* matrix *C* that contains the errors we induce when we match cluster *i* to label *j*.
- Find a row permutation matrix *P* that minimizes the diagonal.

$$\min_{P} \operatorname{Tr}(PC)$$

• Complexity:  $\mathcal{O}(N^3)$ 

# Classifying Images without Labels

- Learn a self-supervised representation
- Loss: Neighbouring images same class + all classes have equal size
- Even better with self-training



Method	Backbone	Labels	Top-1	Top-5
Supervised Baseline	ResNet-50	$\checkmark$	25.4	48.4
Pseudo-Label	ResNet-50	$\checkmark$	-	51.6
VAT + Entropy Min. 56	ResNet-50	$\checkmark$	-	47.0
InstDisc 51	ResNet-50	$\checkmark$	-	39.2
BigBiGAN 15	ResNet-50(4x)	$\checkmark$	-	55.2
PIRL 32	ResNet-50	$\checkmark$	-	57.2
CPC v2 20	ResNet-161	$\checkmark$	52.7	77.9
SimCLR 7	ResNet-50	$\checkmark$	48.3	75.5
SCAN (Ours)	ResNet-50	X	39.9	60.0

# Self-labelling by Clustering



[Asano et al.; ICLR '20]

### Clusters are interpretable



cluster 393, purity: 0.668

cluster 503, purity: 0.930



### Clusters are interpretable



cluster 406, purity: 0.455 clus

cluster 0, purity: 0.558

cluster 2568, purity: 0.377

### **GAN-based Segmentation**

• ReDo: Layer-wise generative models for unsupervised object discovery



- Advantage: You get segmentation for free!
- Drawbacks: fragile training, difficulty scaling due to custom architecture

### 3D from a single image



3D ground truth or shape models



multi-view





depth maps



silhouettes



keypoints



camera viewpoint

### 3D from a single image



3D ground truth or shape models



multi-view



depth maps



silhouettes



keypoints



camera viewpoint

### Unsupervised Learning of 3D Objects

**Training Data** 

Output



#### single-view images of a category

NO other supervision!

instance-specific 3D shapes

Unsupervised Learning of Probably Symmetric Deformable 3D Objects from Images in the Wild, S Wu et al., CVPR 2020

### Observation I

• Symmetry is a strong constraint!



input

flipped

### **Observation II**

• Shading is a strong constraint! (Shape from shading)



Photometric method for determining surface orientation from multiple images Robert J Woodham, Optical engineering, 1980













#### **Photo-Geometric Autgencoding Q3:** Non-symmetric albedo, deformation, etc?































































reconstruction