Learning Representations with Deep InfoMax



Karan Grewal University of Toronto & Princeton University



learning unsupervised representations

- supervised annotation is very costly

- uses:
 - extract useful information for downstream tasks

- good representations:
 - high correlation between data & representation
 - captures signal, ignores noise



mutual information

a measure of how informative one variable is of the other

$$I(\mathcal{X};\mathcal{Z}) = \mathbb{E}_{\mathbb{P}_{\mathcal{XZ}}}\left[\log rac{\mathbb{P}_{\mathcal{XZ}}}{\mathbb{P}_{\mathcal{X}} \otimes \mathbb{P}_{\mathcal{Z}}}
ight]$$

notice that if $\mathbb{P}_{\mathcal{XZ}} = \mathbb{P}_{\mathcal{X}} \otimes \mathbb{P}_{\mathcal{Z}}$ then $I(\mathcal{X};\mathcal{Z}) = 0$



mutual information

1. contrastive loss lower bound (infoNCE)

$$I(\mathcal{X}; \mathcal{Z}) \geq \log(N) + \mathbb{E}_{S} \left[\log \frac{h(x^{*}, z^{*})}{\sum_{j} h(x_{j}, z_{j})} - \mathcal{L}_{contrast} \right]$$
where $h(\cdot, \cdot)$ estimates the ratio $\frac{\mathbb{P}_{\mathcal{XZ}}}{\mathbb{P}_{\mathcal{X}} \otimes \mathbb{P}_{\mathcal{Z}}}$,
$$S = \{(x^{*}, z^{*})\} \cup \{(x_{j}, z_{j})\}_{j=1}^{N}$$
from $\mathbb{P}_{\mathcal{XZ}}$ from $\mathbb{P}_{\mathcal{X}} \otimes \mathbb{P}_{\mathcal{Z}}$ van der

van den Oord et al., 2018

mutual information

2. Donsker-Varadhan representation of KL divergence

$$\mathcal{T}(\mathcal{X};\mathcal{Z}) = \mathbb{E}_{\mathbb{P}_{\mathcal{X}\mathcal{Z}}}\left[\log rac{\mathbb{P}_{\mathcal{X}\mathcal{Z}}}{\mathbb{P}_{\mathcal{X}}\otimes\mathbb{P}_{\mathcal{Z}}}
ight]$$

$$\mathcal{D} = D_{\mathrm{KL}}(\mathbb{P}_{\mathcal{XZ}} || \mathbb{P}_{\mathcal{X}} \otimes \mathbb{P}_{\mathcal{Z}})$$

$$e_{T_\omega:\,\Omega o\mathbb{R}} \,\,\, \mathbb{E}_{\mathbb{P}_{\mathcal{XZ}}}\left[T_\omega
ight] - \log\mathbb{E}_{\mathbb{P}_{\mathcal{X}}\otimes\mathbb{P}_{\mathcal{Z}}}\left[e^{T_\omega}
ight]$$

Donsker & Varadhan, 1983

framework



algorithm

- 1. sample (+) examples $x_+^{(1)},\ldots,x_+^{(n)}\sim \mathbb{P}_X$
- 2. compute representations $z^{(i)} = \mathsf{enc}_\psi(x^{(i)}_+) \quad orall i$
- 3. let $\{(x_+^{(i)}, z^{(i)})\}_i$ be the (+) pairs
- 4. sample (-) examples $x_{-}^{(1)},\ldots,x_{-}^{(n)}\sim\mathbb{P}_{X}$
- 5. let $\{(x_{-}^{(i)}, z^{(i)})\}_i$ be the (-) pairs

6. maximize
$$\frac{1}{n} \sum_{i=1}^{n} T_{\omega}(\underbrace{x_{+}^{(i)}, z^{(i)}}_{(+) \text{ pairs}}) - \log \frac{1}{n} \sum_{i=1}^{n} e^{T_{\omega}(x_{-}^{(i)}, z^{(i)})}$$

engineering footnotes

- encoder is a composition of convolutions followed by fully-connected layers

$$\mathsf{enc}_\psi = f_\psi \circ g_\psi$$

- first apply g_ψ to obtain an M imes M feature map, then apply f_ψ to get representation



- maximize MI between feature map and representation

local InfoMax

- hypothesis: encoder is more likely to capture information shared across all patches
 - pro \rightarrow global structure will be present
 - con → model has no incentive to focus on relevant information, pixel-level noise will be encoded



relevant information

local InfoMax





alternate MI objectives

objective	pros	cons
Donsker-Varadhan (DV)	tightest available bound on KL divergence	requires many (-) samples
Jensen-Shannon divergence (JSD)	stable; few (-) samples needed	not the tightest bound
noise contrastive estimation (NCE)	strongest results	requires many (-) samples



putting it all together

final objective for Deep InfoMax:

e.g. estimate KL lower bound

$$egin{aligned} &\max_{\omega,\psi}lpha\widehat{I}\left(\mathcal{X},\mathcal{Z}
ight)\ &+\max_{\omega,\psi}rac{eta}{M^2}\sum_{j=1}^{M^2}\widehat{I}\left(\mathcal{X}_j,\mathcal{Z}
ight) \end{aligned}$$

 $+\min_{\psi}\max_{ heta}\gamma D_{ extsf{JSD}}(\mathbb{Q}_{ extsf{prior}}||\mathbb{Q}_{\psi})$

global InfoMax

local InfoMax

prior matching

just like adversarial autoencoders

 $lpha,eta,\gamma\;$ are tunable hyperparameters

results - linear classification

train SVM on learned representations

model	conv	fc	Z
VAE	53.8	42.1	39.6
adversarial autoencoder	55.2	43.3	37.8
Bigan	56.4	38.4	44.9
NAT	48.6	42.6	39.6
DIM	57.6	45.6	18.6
DIM - global only	46.8	28.8	29.1
DIM - local only	63.3	54.1	49.6

CIFAR-10

 $conv \rightarrow last conv layer, fc \rightarrow 2nd last fc layer, z \rightarrow representation$

results - nonlinear classification (1/2)

train shallow neural network on learned representations

model	conv	fc	Z	conv	fc	Z
fully supervised		75.4			42.3	
VAE β -VAE adversarial autoencoder BiGAN DIM - global DIM - local (DV) DIM - local (JSD) DIM - local (NCE)	60.7 62.4 59.4 62.6 52.2 72.7 73.3 75.2	60.5 57.9 57.2 62.7 52.8 70.6 73.6 75.6	54.6 55.4 52.8 52.5 43.2 64.7 67.0 69.1	37.2 32.3 36.2 37.6 27.7 48.5 48.1 49.7	34.1 26.9 33.4 33.3 24.4 44.4 45.9 47.7	24.2 29.0 23.3 21.5 20.0 39.3 39.6 41.6

CIFAR-10

 $conv \rightarrow last conv layer$, $fc \rightarrow 2nd last fc layer$, $z \rightarrow representation$

CIFAR-100

results - nonlinear classification (2/2)

train shallow neural network on learned representations

	Tiny ImageNet			STL-10			
model	conv	fc	Y		conv	fc	Z
fully supervised		75.4				42.3	
VAE β -VAE adversarial autoencoder BiGAN DIM - global DIM - local (DV) DIM - local (JSD) DIM - local (NCE)	18.6 19.3 18.0 24.4 11.3 30.4 33.5 34.2	16.9 16.8 17.3 20.2 6.3 29.5 36.9 38.1	11.9 12.4 11.5 13.1 5.0 28.2 31.7 33.3		58.3 57.2 59.5 71.5 42.0 69.2 72.9 72.6	56.7 55.1 54.5 67.2 30.8 63.8 70.9 70.0	46.5 46.9 43.9 58.5 28.1 61.9 65.9 67.1

 $conv \rightarrow last conv layer$, $fc \rightarrow 2nd last fc layer$, $z \rightarrow representation$

results - MI neural estimate

train a neural network to estimate MI between input and representations

model	MINE estimate
VAE adversarial autoencoder BiGAN NAT DIM DIM - global only DIM - local only	93.0 87.5 37.7 6.0 101.7 49.6

CIFAR-10

results - neural dependency measure

shuffle representations along batch axis and train discriminator to tell real from fake



CIFAR-10

model	discriminator loss
VAE	1.6
adversarial autoencoder	0.1
BiGAN	24.5
NAT	0.1
DIM	22.9
DIM - global only	10.0
DIM - local only	9.2

Brakel & Bengio, 2017

results - occlusion



Contrastive Predictive Coding





Deep InfoMax + occlusion

results - occlusion

non-linear classification with occlusion

model	CIFAR-10	STL-10
CPC	77.5	77.8
DIM - original	81.0	77.0
DIM - multiple representations	77.5	78.2

extensions

- (a) maximize MI between different views of the same object/scene
 - minimize contrastive loss between representations obtained from 2 separate views / encoders

Tian et al., 2019; Bachman et al., 2019

- (b) maximize MI between audio representations from the same speaker
 - similar to (a), but the two "views" correspond to audio waveforms taken from the same speaker

conclusions

pros

- Deep InfoMax doesn't require a decoder
- MI-based objectives can be extended to other tasks
- local information can be encoded advantageous for downstream tasks

cons

- still requires 2 discriminators
- may be hard to scale encoding local information can be harmful

collaborators



Devon Hjelm Microsoft Research, Montréal & Mila



Alex Fedorov Georgia Tech



Samuel Lavoie-Marchildon Mila



Phil Bachman Microsoft Research, Montréal



Adam Trischler Microsoft Research, Montréal



Yoshua Bengio Mila