

THE UNIVERSITY OF BRITISH COLUMBIA

Topics in AI (CPSC 532S): **Multimodal Learning with Vision, Language and Sound**

Lecture 11: Unsupervised Learning, Autoencoders



Autoencoders

Self (i.e. self-encoding)

Feed forward network intended to reproduce the input

- Encoder/Decoder architecture Encoder: $f = \sigma(\mathbf{W}\mathbf{x})$ Decoder: $g = \sigma(\mathbf{W}'\mathbf{h})$



*slide from Louis-Philippe Morency



er

Autoencoders

Self (i.e. self-encoding)

Feed forward network intended to reproduce the input

- Encoder/Decoder architecture Encoder: $f = \sigma(\mathbf{W}\mathbf{x})$ Decoder: $g = \sigma(\mathbf{W}'\mathbf{h})$
- Score function

$$\mathbf{x}' = f(g(\mathbf{x}))$$

 $\mathcal{L}(\mathbf{x}',\mathbf{x})$



*slide from Louis-Philippe Morency



er

Autoencoders

A standard neural network architecture (linear layer followed by non-linearity)

- Activation depends on type of data (e.g., sigmoid for binary; linear for real valued)
- Often use tied weights

 $\mathbf{W}' = \mathbf{W}$



*slide from Louis-Philippe Morency



De-noising Autoencoder

Idea: add noise to input but learn to reconstruct the original

- Leads to better representations
- Prevents copying

Note: different noise is added during each epoch



*slide from Louis-Philippe Morency



er

Stacked (deep) Autoencoders and Denoising Autoencoders

What **can we do** with them?

- Good for compression (better than PCA)
- Disregard the decoder and use the middle layer as a representation
- Fine-tune the autoencoder for a task



*slide from Louis-Philippe Morency









[Pathak et al., 2016]



(a) Central region





(b) Random block



(c) Random region

















Pretraining Method	Supervision	Pretraining time	Classification	Detection	Segmentation
ImageNet [26]	1000 class labels	3 days	78.2%	56.8%	48.0%
Random Gaussian	initialization	< 1 minute	53.3%	43.4%	19.8%
Autoencoder	_	14 hours	53.8%	41.9%	25.2%
Agrawal <i>et al</i> . [1]	egomotion	10 hours	52.9%	41.8%	-
Doersch et al. [7]	context	4 weeks	55.3%	46.6%	-
Wang <i>et al</i> . [39]	motion	1 week	58.4%	44.0%	-
Ours	context	14 hours	56.5%	44.5%	29.7%



Spatial Context Networks



[Wu, Sigal, Davis, 2017]



Spatial Context Networks



	Initialization	Supervision	Pretraining time	Classification	Detection
Random Gaussian	random	N/A	< 1 minute	53.3	43.4
Wang <i>et al</i> . [32]	random	motion	1 week	58.4	44.0
Doersch et al. [3]	random	context	4 weeks	55.3	46.6
*Doersch et al. [3]	1000 class labels	context	—	65.4	50.4
Pathak <i>et al</i> . [21]	random	context inpainting	14 hours	56.5	44.5
Zhang <i>et al</i> . [36]	random	color	—	65.6	46.9
ImageNet [21]	random	1000 class labels	3 days	78.2	56.8
*ImageNet	random	1000 class labels	3 days	76.9	58.7
SCN-EdgeBox	1000 class labels	context	10 hours	79.0	59.4

[Wu, Sigal, Davis, 2017]





Every layer could be treated as a random variable, then entire network is a Markov Chain

Data processing theorem: if the only connection between X and Z is through Y, the information that Z gives about X cannot be bigger than the information

that Y gives about X.



50 networks of same topology being optimized



Observation: In the information plane layers first increase the mutual information between themselves and the output and then reduce information between themselves and the input (which leads to "forgetting" of irrelevant inputs and ultimately generalization)



Limitation: Does not seem to work for non-Tanh activations (e.g., ReLU)

