

Humans outperform machines at unsupervised learning

- Humans are very good at unsupervised learning, e.g. a 2 year old knows intuitive physics
- Babies construct an approximate but sufficiently reliable model of physics, how do they manage that? Note that they interact with the world, not just observe it.



Invariance and Disentangling

- Invariant features
- Which invariances?
- Alternative: learning to disentangle factor.



 Good disentangling → avoid the curse of dimensionality:

Dependencies are "simple" when the data is projected in the right abstract space

Disentangling from denoising objective (Glorot, Bordes & Bengio ICML 2011)

- Early deep learning research already is looking for possible disentangling arising from unsupervised learning of representations
- Experiments on stacked denoising auto-encoders with ReLUs, on BoW text classification
- Features tend to specialize to either sentiment or domain



Unsupervised Learning of Representations: Simple Auto-Encoders



- Code = new coordinate system
- Encoder and decoder can have more layers
- Reconstruction can be probability distribution

Denoising Auto-Encoder (Vincent et al 2008)



- Corrupt the input during training only
- Train to reconstruct the uncorrupted input



- Encoder & decoder: any parametrization
- As good or better than RBMs for unsupervised pre-training



Auto-Encoders Learn Salient Variations, like a non-linear PCA

- Minimizing reconstruction error forces to keep variations along manifold.
- Regularizer wants to throw away all variations.
- With both: keep ONLY sensitivity to variations ON the manifold.

Manifold Learning = Representation Learning





Interpolating in Latent Space

If the model is good (unfolds the manifold), interpolating between latent values yields plausible images.



woman with glasses

woman without glasses

without glasses

man with glasses

11

Deep Unsupervised Generative Models

Texture



Shakespeare

Why, Salisbury must find his flesh and thought That which I am not aps, not a man and in fire, To show the reining of the raven and the wars To grace my hand reproach within, and not a fair are hand, That Caesar and my goodly father's world; When I was heaven of presence and our fleets, We spare with hours, but cut thy council I am great, Murdered and by thy master's ready there My power to give thee but so much as hell: Some service in the noble bondman here, Would show him to her wine.

Chinese characters

	Let		車鳥	月、思	51	
,	木ケ目	搅	坎女	鴙	拍 打	
	挌	职	松	誷	壑十	
	兵	坷	拹	虎	晋	
	摀	樢	髡	地受	Ē	
Bedrooms						

Hand-writing

More of national temperament More of national temperament More of national temperament more of national temperament more of national temperament



Latent Variables and Abstract Representations

- Encoder/decoder view: maps between low & high-levels
- Encoder does inference: interpret the data at the abstract level
- Decoder can generate new configurations
- Encoder flattens and disentangles the data manifold



Extracting Structure By Gradual Disentangling and Manifold Unfolding (Bengio 2014, arXiv 1407.7906)

Each level transforms the data into a representation in which it is easier to model, unfolding it more, contracting the noise dimensions and mapping the signal dimensions to a factorized (uniform-like) distribution.

$$\min KL(Q(x,h)||P(x,h))$$

for each intermediate level h



Helmholtz Machines (Hinton et al 1995) and Variational Auto-Encoders (VAEs)

 $P(h_3)$

 $P(h_2|h_3|)$

 $P(h_1|h_2)$

 $P(x|h_1)$

Decoder = generator

(Kingma & Welling 2013, ICLR 2014) (Gregor et al ICML 2014; Rezende et al ICML 2014) (Mnih & Gregor ICML 2014; Kingma et al, NIPS 2014) h_3

- Parametric approximate inference
- Successors of Helmholtz machine (*Hinton et al '95*)
- Maximize variational lower bound on log-likelihood: $\min KL(Q(x,h)||P(x,h))$ where Q(x) = data distr.

or equivalently

$$\sum_{x,h} Q(x)Q(h|x) \log \frac{P(x,h)}{Q(h|x)} = \sum_{x,h} Q(x)Q(h|x) \log P(x|h) + KL(Q(h|x)||P(h))$$

inference

Encoder =

 $Q(h_3|h_2)$

 $Q(h_2|h_1)$

 $Q(h_1|x)$

 h_2

 h_1

 \mathcal{X}

Q(x)

GAN: Generative Adversarial Networks A radical alternative to max, likelihood

Goodfellow et al NIPS 2014



Early Days of GAN Samples







CIFAR-10 (fully connected)



CIFAR-10 (convolutional)

Convolutional GANs

(Radford et al, arXiv 1511.06343)

Strided convolutions, batch normalization, only convolutional layers, ReLU and leaky ReLU



Convolutional Networks

- Scale up neural networks to process very large images / video sequences
 - Sparse connections
 - Parameter sharing
- Automatically generalize across spatial translations of inputs
- Applicable to any input that is laid out on a grid (1-D, 2-D, 3-D, ...)

Convnets: Key Idea

- Replace matrix multiplication in ordinary neural nets with convolution
- Everything else stays the same
 - Maximum likelihood
 - Back-propagation
 - etc.

Convolutional Neural Networks

- A special kind of deep learning tailored for images
- Exploits the invariance to translations
- Exploits multi-scale hierarchy

Convolutional neural network for imaging data



2D Convolution



Figure 9.1, Deep Learning book, Goodfellow et al 2016

Sparse Connectivity

Sparse connections due to small convolution kernel



Dense connections





Sparse Connectivity

Sparse connections due to small convolution kernel



Dense connections





Growing Receptive Fields



Parameter Sharing



Cross-Channel Pooling and Invariance to Learned Transformations



Pooling with Downsampling



Convolution with Stride



Major ConvNet Architectures

- Spatial Transducer Net: input size scales with output size, all layers are convolutional
- All Convolutional Net: no pooling layers, just use strided convolution to shrink representation size
- Inception: complicated architecture designed to achieve high accuracy with low computational cost
- ResNet: blocks of layers with same spatial size, with each layer's output added to the same buffer that is repeatedly updated. Very many updates = very deep net, but without vanishing gradient.

ResNets: Skip Connections

• Identity paths make it possible for gradients to flow through deeper networks (He et al 2015), SOTA on object recognition



Deep Data Fusion

- Deep nets are very good at combining multiple sources of data, multiple sensors or modalities
- Can have separate pre-processing stages for each modality, then CONCATENATE the representations before continuing processing



Need to map to the same spatial scale, or 'copy' a non-spatial modality at all positions.

Generating Text from Images

- (Kiros *et al.*, 2014; Mao *et al.*, 2014; Donahue *et al.*, 2014;
 Vinyals *et al.*, 2014; Fang *et al.*, 2014; Chen and Zitnick, 2014; Karpathy and Li, 2014; Venugopalan *et al.*, 2014).
- Convolutional net → generative RNN





A group of people shopping at an outdoor market.

There are many vegetables at the fruit stand.



A close up of a child holding a stuffed animal



Two pizzas sitting on top of a stove top oven.

(GT: Three different types of pizza on top of a stove.)

U-Net Architecture for CNNs with Pixel-Wise Outputs



Generative Adversarial Networks











Image 2 Image



<u>Isola et al. 2016</u>

Text 2 Image, B&W 2 Color

This bird is red and brown in color, with a stubby beak

The bird is short and stubby with yellow on its body

A bird with a medium orange bill white body gray wings and webbed feet

This small black bird has a short, slightly curved bill and long legs

A small bird with varying shades of brown with white under the eyes

A small yellow bird with a black crown and a short black pointed beak

This small bird has a white breast, light grey head, and black wings and tail













Lucy Li

Horse 2 Zebra: matching 2 domains by analogy of their distribution structure

Input video

Output video





The Future of Deep AI



- Scientific progress is slow and continuous, but social and economic impact can be disruptive
- Many fundamental research questions are in front of us, with much uncertainty about when we will crack them, but we will
- Importance of continued investment in basic & exploratory AI research, for both practical (recruitment) short-term and long-term reasons
- Let us continue to keep the field open and fluid, be mindful of social impacts, and make sure AI will bloom for the benefit of all

Montreal Institute for Learning Algorithms

Universit

de Mon