

https://introml.mit.edu/

6.390 Intro to Machine Learning

Lecture 8: Representation Learning

Shen Shen April 4, 2025 11am, Room 10-250

Outline

- Neural networks are *representation* learners
- Auto-encoders
- Unsupervised and self-supervised learning
- Word embeddings
- (Some recent representation learning ideas)

Outline

- Neural networks are *representation* learners
- Auto-encoders
- Unsupervised and self-supervised learning
- Word embeddings
- (Some recent representation learning ideas)



compositions of ReLU(s) can be quite expressive



in fact, asymptotically, can approximate any function!



[image credit: Phillip Isola]

https://playground.tensorflow.org/

Two different ways to visualize a function



Two different ways to visualize a function





Representation transformations for a variety of neural net operations



and stack of neural net operations











Neural networks are representation learners

Deep nets transform datapoints, layer by layer

Each layer gives a different *representation* (aka *embedding*) of the data



https://distill.pub/2017/feature-visualization/

Outline

- Neural networks are *representation* learners
- Auto-encoders
- Unsupervised and self-supervised learning
- Word embeddings
- (Some recent representation learning ideas)

Observed image

Drawn from memory



[Bartlett, 1932] [Intraub & Richardson, 1989]





LEONARDO 19 ANNI STUDENTE







[https://www.behance.net/gallery/35437979/Velocipedia]



"I stand at the window and see a house, trees, sky. Theoretically I might say there were 327 brightnesses and nuances of colour. Do I have "327"? No. I have sky, house, and trees."

— Max Wertheimer, 1923

humans also learn representations





Good representations are:

- Compact (*minimal*)
- Explanatory (*roughly sufficient*)



[See "Representation Learning", Bengio 2013, for more commentary]

Auto-encoder



"What I cannot create, I do not understand." Feynman

Auto-encoder

$$\min_W ||x- ilde{x}||^2$$





Image



$$ilde{x} = \mathrm{NN}(x;W)$$



Reconstructed image



 $egin{array}{l} {
m input} \ x\in \mathbb{R}^d \end{array}$

25





[Often, encoders can be kept to get "good representations" whereas decoders can serve as "generative models"]

Outline

- Neural networks are *representation* learners
- Auto-encoders
- Unsupervised and self-supervised learning
- Word embeddings
- (Some recent representation learning ideas)

Supervised Learning



Undergrads were time-consuming, algorithms were flawed, and the team didn't have money—Li said the project failed to win any of the federal grants she applied for, receiving comments on proposals that it was shameful Princeton would research this topic, and that the only strength of proposal was that Li was a woman.

A solution finally surface a graduate student who a Amazon Mechanical Turk sitting at computers arou online tasks for pennies.

"He showed me the webs



The Amazon Mechanical Turk backend for classifying images. Image: ImageNet

knew the ImageNet project was going to happen," she said. "Suddenly we found a tool that could scale, that we could not possibly dream of by hiring Princeton undergrads."

Unsupervised Learning



Label prediction (supervised learning)



Features

Feature reconstruction (unsupervised learning)



Self-supervised learning



Partial features

Other partial features

Masked Auto-encoder




predict color from gray-scale





[Zhang, Isola, Efros, ECCV 2016]

Masked Auto-encoder



Devlin et al., BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding, 2019

Self-supervised learning



Common trick:

- Convert "unsupervised" problem into "supervised" setup
- Do so by cooking up "labels" (prediction targets) from the raw data itself — called *pretext* task

Outline

- Neural networks are *representation* learners
- Auto-encoders
- Unsupervised and self-supervised learning
- Word embeddings
- (Some recent representation learning ideas)

Large Language Models (LLMs) are trained in a self-supervised way



- Scrape the internet for unlabeled plain texts.
- Cook up "labels" (prediction targets) from the unlabeled texts.
- Convert "unsupervised" problem into
 - "supervised" setup.

"To date, the cleverest thinker of all time was Issac."



[video edited from 3b1b]



input embedding (e.g. via a fixed encoder)

[image edited from 3b1b]

To	date		the	cle	ve	rest	think	er of	all	time	WBE	
1	Ť	Ĩ	1	Ĩ	Î	Ē	Ĩ	Î	i i	Ĩ	Ť	???
4	4			ţ.	Ļ		4	÷			Ļ	
5.4	7.8	[9.7]	2.6	[3.6]	5.6	1.6	[9.7]		3.2	[6.7]	4.4	
7.1	5.2	7.9	7.7	4.3	4.3	1.1	4.6	•	6.6	2.7	8.4	
6.0	5.6	4.6	4.5	6.9	9.8	6.5	9.7		1.3	7.3	6.9	
5.4	9.2	7.7	5.6	0.6	1.0	1.4	6.0		7.1	9.5	2.9	
4.2	0.7	1.2	0.2	6.6	2.1	1.9	7.3		2.9	2.5	8.1	
6.4	0.9	6.3	6.1	6.6	1.6	3,7	0.4		1.8	5.7	3.9	
4.3	0.2	1.4	6.1	2.1	6.5	8.1	2.8		5.8	5.9	8.7	
8.8	8.2	9.4	6.1	1.3	2,5	1.0	1.2		0.2	5.7	5.8	
3.8	8.6	4.1	6.8	3.6	2.4	1.0	1.2		0.0	9.4	6.9	

[video edited from 3b1b]



Cross-entropy loss encourages the internal weights update so as to make this probability higher "A robot must obey the orders given it by human beings ..."







dot-product similarity





good embedding representation => sensible dot-product similarity => enables effective attention in transformers next week.

For now, let's think about dictionary look-up:







What if we run





What if we run



But we can probably see the rationale behind this:



If we are to generalize this idea, we need to:

















Outline

- Neural networks are *representation* learners
- Auto-encoders
- Unsupervised and self-supervised learning
- Word embeddings
- (Some recent representation learning ideas)

Good representations are:

Auto-encoders try to achieve these

these may just emerge as well

- Compact (*minimal*)
 Explanatory (*roughly sufficient*)
- Disentangled (independent factors)
- Interpretable
- Make subsequent problem solving easy



• pre-training

- contrastive
- multi-modality



Often, what we will be "tested" on is not what we were trained on.



Final-layer adaptation: freeze *f*, train a new final layer to new target data



Finetuning: initialize f as f, then continue training for f' as well, on new target data



A lot of data

A little data

[images credit: visionbook.mit.edu]

The allegory of the cave





Contrastive learning



[images credit: visionbook.mit.edu]

Contrastive learning





Figure 1. ImageNet top-1 accuracy of linear classifiers trained on representations learned with different self-supervised methods (pretrained on ImageNet). Our method, SimCLR, is shown in bold.

Multi-modality



[images credit: visionbook.mit.edu]

Video - audio



[Owens et al, Ambient Sound Provides Supervision for Visual Learning, ECCV 2016] [Slide credit: Andrew Owens]
What did the model learn?



[Owens et al, Ambient Sound Provides Supervision for Visual Learning, ECCV 2016] [Slide credit: Andrew Owens]





[Owens et al, Ambient Sound Provides Supervision for Visual Learning, ECCV 2016] [Slide credit: Andrew Owens]

Image classification (done in the contrastive way)



```
# extract feature representations of each modality
I_f = image_encoder(I) #[n, d_i]
T_f = text_encoder(T) #[n, d_t]
```

```
# joint multimodal embedding [n, d_e]
I_e = l2_normalize(np.dot(I_f, W_i), axis=1)
T_e = l2_normalize(np.dot(T_f, W_t), axis=1)
```

```
# scaled pairwise cosine similarities [n, n]
logits = np.dot(I_e, T_e.T) * np.exp(t)
```

```
# symmetric loss function
labels = np.arange(n)
loss_i = cross_entropy_loss(logits, labels, axis=0)
loss_t = cross_entropy_loss(logits, labels, axis=1)
loss = (loss_i + loss_t)/2
```

Image classification (done in the contrastive way)



1. Contrastive pre-training

2. Create dataset classifier from label text



Dall-E: text-image generation



Figure 2: A high-level overview of unCLIP. Above the dotted line, we depict the CLIP training process, through which we learn a joint representation space for text and images. Below the dotted line, we depict our text-to-image generation process: a CLIP text embedding is first fed to an autoregressive or diffusion prior to produce an image embedding, and then this embedding is used to condition a diffusion decoder which produces a final image. Note that the CLIP model is frozen during training of the prior and decoder.

How Much Information is the Machine Given during Learning?

"Pure" Reinforcement Learning (cherry)

- The machine predicts a scalar reward given once in a while.
- A few bits for some samples

Supervised Learning (icing)

- The machine predicts a category or a few numbers for each input
- Predicting human-supplied data
- ▶ 10→10,000 bits per sample

Self-Supervised Learning (cake génoise)

- The machine predicts any part of its input for any observed part.
- Predicts future frames in videos

Millions of bits per sample

© 2019 IEEE International Solid-State Circuits Conference

1.1: Deep Learning Hardware: Pa



Y. LeCun

Summary

- We looked at the mechanics of neural net. Today we see deep nets learn representations, just like our brains do.
- This is useful because representations transfer they act as prior knowledge that enables quick learning on new tasks.
- Representations can also be learned without labels, e.g. as we do in unsupervised, or selfsupervised learning. This is great since labels are expensive and limiting.
- Without labels there are many ways to learn representations. We saw today:
 - representations as compressed codes, auto-encoder with bottleneck
 - (representations that are predictive of their context)
 - (representations that are shared across sensory modalities)

We'd love to hear

your thoughts. Thanks!