Self-supervised Learning Learning without labels

Antonia Marcu & Jonathon Hare

What and why?

- How might we train large models to do something useful with very little labelled data (but lots of unlabelled data)?
- Can we somehow learn embeddings that are useable for downstream tasks?

• Self-supervised Learning might be the answer!

referred to as learning the "pretext task")



Basic idea: train model on the unlabelled data in some way (this is often

encoder) on a small labelled dataset



Basic idea: then train a small task-specific network (and possibly fine-tune

Types of SSL

- Auto-regressive SSL
- Auto-encoders
 - Noisy/masked auto-encoders
- Contrastive SSL
- Non-contrastive SSL

Types of SSL

- Auto-regressive SSL
- Auto-encoders
 - Noisy/masked auto-encoders
- Contrastive SSL
- Non-contrastive SSL

Encoder-decoder

Siamese



Auto-regressive SSL Learning what comes next



Auto-regressive Image Modelling PixelRNN

x_1	



$$p(\mathbf{x}) = \prod_{i=1}^{n^2} p(x_i | x_1, \dots, x_n)$$

Pixel Recurrent Neural Networks



Autoencoders Learning to compress / Learning to reconstruct



Autoencoders: Image Modelling Masked Autoencoders



Masked Autoencoders Are Scalable Vision Learners

Autoencoders: Text Modelling BERT



Contextual Word Embeddings

[CLS]	
and	
,	
50	
far	
as	
they	
knew	
,	
they	
were	
quite	
right	
•	
[SEP]	

Autoencoders: Text Modelling BERT



Masked prediction

Next sentence prediction

Autoencoders: Text Modelling BART



train on noisy text

BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension.



predict clean text BCDE Α Autoregressive Decoder $\langle s \rangle A B C D$



Exploiting Similarity in the Embedding Space

Quick Nomenclature

Embedding space = Latent space = Representation Space*

*ish. Some of the papers we will discuss today use **Representation Space** to mean something very specific. Other use **Embedding space** to mean something specific

Siamese Nets



Signature Verification using a "Siamese" Time Delay Neural Network

Can we turn this on its head?

• similar inputs in similar ways



• similar inputs in similar ways



similar inputs in similar ways



How can we get similar inputs in an unsupervised way?

similar inputs in similar ways







similar inputs in similar ways





 \bullet similar inputs in similar ways





 \bullet similar inputs in similar ways





similar inputs in similar ways





What can go wrong?

Preventing Collapse Pushing dissimilar embeddings away from each other

- Contrastive Learning
- Information Maximisation

distance between dissimilar samples



distance between dissimilar samples





distance between dissimilar samples







"A Simple Framework for Contrastive Learning of Visual Representations"

- In an N-dimensional batch:

 - Negative samples (different): all other samples in the batch

"A Simple Framework for Contrastive Learning of Visual Representations"

• Positive samples (similar): augmented versions of each image (2N total)

- In an N-dimensional batch: lacksquare

 - Negative samples (different): all other samples in the batch

"A Simple Framework for Contrastive Learning of Visual Representations"

• Positive samples (similar): augmented versions of each image (2N total)



"A Simple Framework for Contrastive Learning of Visual Representations"



What can go wrong?

distance between dissimilar samples





distance between dissimilar samples





BYOL No negative samples needed



the representation* of the augmented sample will look like?

Intuition: Given the representation* of a sample, can we learn to predict what

* used in the generic sense here

Reminder



"A Simple Framework for Contrastive Learning of Visual Representations"

Reminder



"A Simple Framework for Contrastive Learning of Visual Representations"

BYOL No negative samples needed



BYOL No negative samples needed



Information Maximisation

BarlowTwins



"Barlow Twins: Self-Supervised Learning via Redundancy Reduction"

BarlowTwins



"Barlow Twins: Self-Supervised Learning via Redundancy Reduction"

BarlowTwins



"Barlow Twins: Self-Supervised Learning via Redundancy Reduction"

Distortions

 Defining rich enough augmentations is crucial for the success of such methods

What makes an augmentation "good"?

