

Understanding the Multimodal World Results and Challenges

Sergey Nikolenko

Artificial Intelligence Lab, PDMI RAS

March 23, 2021





R

Multimodal Data Analysis



Plan

- Embeddings in modern ML
- The world is multimodal: how do we cope?
- Multimodal embeddings, different approaches
- A review of three state of the art 2020 papers on multimodal deep learning
- Our own example: recent work on symbolism prediction







Artificial Intelligence Lab, PDMI RAS







Learning Latent Representations

- Latent representations, or embeddings, is one of the most important concepts in modern machine learning
- Mathematically it's just a parametrization of a manifold: we assume that a set in a high-dimensional space actually has a lower dimension and try to find the parametrization





Russian Academy of Sciences

Learning Latent Representations

- The high-dimensional points can be images, sentences, sounds...
- We would like the latent space to have good properties; how can we find a good parametrization?









Learning Latent Representations

- Autoencoders: let's train a model to reconstruct a point from its latent space representation; datasets are not a problem now
- The important question here is how to ensure a good structure for the distribution of latent representations: VAE, AAE etc.; let's not go there



R

Multimodal Data Analysis

Russian Academy of Sciences

How does multimodal analysis work?

- Most multimodal analysis attempts to learn joint embeddings in some unified latent space
- The objective is to
 construct a vector space
 that brings semantically
 similar elements from
 different modalities close
 together





Russian Academy of Sciences

How can we learn joint embeddings?



> Learn a **joint** representation

- Learn coordinated representations

Learn an encoder-decoder transformation from one modality to another





Russian Academy of Sciences

Constraints via ranking losses

- The "constraint", bringing together vectors from different modalities, is usually implemented via some kind of ranking loss: contrastive, triplet loss etc.
- DeViSE (Frome et al., 2013); VSE++ (Faghri et al., 2018)





Children and adults

forms of martial arts

A reporter is talking about a movie scene

from the wolverines

A man playing guitar



A classical example

Example with video: \succ retrieval based on text queries (Mithun et al., 2018)





- Two different joint >spaces: Object-Text and Activity-Text space, pairwise ranking loss
- What has changed since 2018?





Transformers!

context2Vec Pre-trained seq2seq >By now, many state of the are ULMFiT ---- ELMo Transformer models are based on Multi-lingual MultiFiT Transformer-like self-attention BERT Cross-lingual Multi-task Multimodal too: ViLBERT, XLM UDify + Generation \succ +Knowledge Graph MT-DNN MASS Permutation LM LXMERT, VL-BERT, Transformer-XL Knowledge distillation UniLM More data Span prediction Remove NSP VisualBERT, UNITER, MT-DNN_{KD} Longer time Remove NSP ERNIE More data (Tsinghua) InterBERT, PixelBERT, XLNet SpanBERT Neural entity linker RoBERTa ERNIE-VIL, DeVLBert... KnowBert man wearing white shirt is walking

> on sidewalk along side other

> > Layer 3

Man

Laver 4

Shirt

Layer 5

Sidewalk

Layer 6

pedestrians



Semi-supervised Sequence Learning



Adversarial Training for Vision-and-Language

- Gan et al. (NeurIPS \succ 2020) develop VILLA (Vision-and-Language Large-scale Adversarial training): task-agnostic adversarial pretraining (representation learning), then task-specific adversarial fine-tuning
- V+L Transformers





A man sits on a colorful man-drawn carriage, while another man stands beside it.

















NeurIPS 2020

rd Conference on Neural Information Processing Systems Online • December 6-12, 2020





Multimodal Information in Large Transformers

- ➤ How do we model face-to-face communication?
- Rahman et al. (ACL 2020) propose a new gate for Transformers to accept multimodal nonverbal data during fine-tuning

#	Spoken words + acoustic and visual behaviors	Ground Truth	MAG- XLNet	XLNet
1	"And it really just lacked what made the other movies more enjoyable." + Frustrated and disappointed tone	-1.4	-1.41	-0.9
2	"But umm I liked it." + Emphasis on tone + positive shock through sudden eyebrow raise + +	1.8	1.9	1.2
3	"Except their eyes are kind of like this welcome to the polar express." + tense voice + frown expression	-0.6	-0.6	0.8
4	"Straight away miley cyrus acting miley cyrus, or lack of, she had this same expression throughout the entire film" + sarcastic voice + frustrated facial expression	-1.0	-1.2	0.2

58th Annual Conference of the Association for Computational Linguistics Seattle, Washington July 5–10, 2020

ACL 2020







12-in-1: Multi-Task Vision and Language

Lu et al. (CVPR 2020): \succ representation learning by solving 12 different problems with a ViLBERT-derived model

and Pattern Recognition



	Visual Question Answering What color is the child's outfit? Orange				
	Referring Expressions child sheep basket people sitting on chair				
-	Multi-modal Verification The child is petting a dog. false				
4	Caption-based Image Retrieval A child in orange clothes plays with sheep.				

IR COCO/Fli Three zebras	ickr-like are grazing in a grass field.	GQA-like Is the baby zebra stand to the zebra on the GuessWhat Guesser-like Q: Which entity is it? Q: Is it of A: Zebra A: No	thing next Yes VQA-like How main there on the left? Q: Is it eating grass? A: Yes VQA-like How main there on RefCOCOg-li baby zel	hy zebras are h the right? Two ke bra RefCOCO-like trees
IR COCO/Flid Elephants are	ckr-like e bathing in the river water.	Visual7w-like Which is the baby elephant RefCOCO+-like swimming elephant	Yisual Genome QA-like Where are the elephants? In water SNLI-VE-like No elephants in the image are swimming. contradiction	NLVR2-like At least one of the animals in either image is swimming. True
}	CVPR 2020		> And now for	a detailed
山山の	Conference on Computer Vision	Seattle, Washington	example	

June 14-19, 2020



Russian Heademy of Sciences

Example: symbolism in advertising

- > Ads are all about symbols
- Often they are not easy to parse even for a human
- Sometimes an object resembles another



Sometimes it's not even an object...





Example: symbolism in advertising

But text really helps! What is this ad about?







Example: symbolism in advertising

But text really helps! What is this ad about? The text explains it all.







Example: symbolism in advertising

In 2017, \succ researchers from the University of Pittsburgh published a dataset on understanding ads; thousands of labeled images











Labeling: topics, sentiment, and symbols





Russian Heademy of Sciences

Labeling: visual understanding strategy







Labeling: Q&A



What should I do? I should be careful what words I use on my kid.Why?Because words can hurt as much as fists.

Q: Why should I be careful what words I use on my kid?A: Because words can hurt as much as fists.



I should stop smoking because my lungs are extremely sensitive and could go up in smoke.

I should buy this candy because it is unique and rises above the rest, like the Swiss Alps.





Russian Academy of Sciences

So what did people do?

- (Hussein et al., 2018): in the original work, symbols are predicted with object detection
- But these results were obviously treated as a weak baseline





Speed (fast, quick, speed)





Symbolism detection by joint embeddings

- (Ye, Kovashka, 2018): later, joint multimodal representations were used to map symbols into a joint latent space with statements describing the ad's message
- Siamese networks with triplet loss and a knowledge base for symbols





Symbolism detection by joint embeddings



So this work was "very multimodal" from the start





Russian Academy of Sciences

Personalized symbolism detection

(Murrugarra-Llerena, >Kovashka, 2019): let's add the viewer's personality! They collect caption-gaze samples from users



I should buy this furniture because it is sustainable even in water

I should leave my work behind because a beautiful vacation awaits



I should buy this bottle because it is chilled and refreshing



I should drink Miller lite because its new vortex bottle gives you the smoothest pour.



I should buy this car because it's good for my family.

I should buy this car because it's elegant.

I should buy this car because it's safe for my children. (b)



... fami

Style



Russian Heademy of Sciences

Personalized symbolism detection

- The dataset contains text annotations and the annotator's personality
- The approach is again to learn joint embeddings, this time with three modalities











I should take more time to put myself together because it will allow me to find more **appropriate companionship.**



I should buy Gucci Guilty perfume because it will make me a **sexier person**

R

Multimodal Data Analysis

Russian Academy of Sciences

But what about the text?

- All of these works treat the ad as an image; they use NLP only to process Q&A statements
- But often an ad is completely incomprehensible without reading the text!



A new era of entertainment is about to begin.





IRRESISTIBLE SWISS MILK FILLED CHOCOLATE





But what about the text?

- This was noted already in the original competition on this dataset, in a work by Otani et al. (2018)
- They applied it to a different problem: matching text statements to ads; their pipeline is actually rather straightforward IR

Statement text Visual and OCR clues from ad





What we did

- Our approach was to combine text-based models that read the text with OCR and image processing:
 - image-based classifier: we already exceeded state of the art just with an EfficientNet-based classifier!
 - object detection: Faster R-CNN with Inception-Resnet backbone
 - text-based models:

Transformer-based models (BERT, RoBERTa) and simpler classifiers



' Sciences



Russian Heademy of Sciences

An aside: OCR is hard!

By the way, OCR is hard! (later we had some more luck with CharNet)



Tesseract (Smith, 2007)	Bite the boredom, TCT am 5 bite-sized treats with a crunchy		
	outside PUR RCo e eB Pree ee Rar ed oe 2 ; RN good		
EAST + Tesseract (Kopeykina and	Bite dale boredom, unleash the fun! bite sized treats rh Pe		
Savchenko, 2019)	tah outside and delicious aay filled Sita bo ii te STi) me KFc		
	macaroni and cheese. Sl PPPS m)/7-) UC So good		
Google Android Vision API	FC ma heese A Bite the boredom, unleash the fun! 5 bite-		
	sized treats with a cr good chy and a delicious center filled		
	with ni and creamy cheese. soft		
CloudVision (Otani et al., 2018)	Bite the boredom, unleash the fun! 5 bite-sized treats with a		
	crunchy outside and a delicious center filled with soft maca-		
	roni and creamy cheese. KFC ma and heese good KFC		
AR-Net Caption	a poster for the film		





Our results

By reading the text and >linear blending, we improve over state of the art results: even OCR quality is not as crucial as it might seem

		Dataset: 221 labels			Dataset: 53 labels			
OCR	Model	Image-only results: 0.1967		Image-only results: 0.2814				
		Text-based	Ensemble	Ensemble +	Text-based	Ensemble	Ensemble +	
		(w/text)	(w/text)	Backoff (all)	(w/text)	(w/text)	Backoff (all)	
	BoN	0.0881	0.1996	0.2002	0.1345	0.2889	0.2865	
Tesseract	BERT	0.0220	0.1932	0.1967	0.1794	0.2935	0.2892	
	RoBERTa	0.0220	0.1934	0.1967	0.1765	0.2918	0.2882	
EAST	BoN	0.1198	0.2087	0.2122	0.1457	0.2955	0.2957	
EASI +	BERT	0.0689	0.1949	0.1989	0.1933	0.3050	0.3046	
resseract	RoBERTa	0.0225	0.1975	0.2013	0.1994	0.2914	0.2919	
AD Not	BoN	0.1014	0.1918 0.1975 0.2012		0.1586	0.2862		
AK-INCL Contions	BERT	0.1089			0.1615	0.2834		
Captions	RoBERTa	0.0226			0.1733	0.2871		
Google	BoN	0.1407	0.2102	0.2106	0.2076	0.3026	0.2994	
Android	BERT	0.0971	0.1975	0.2000	0.2278	0.3054	0.3017	
Vision API	RoBERTa	0.1201	0.1992	0.2014	0.2409	0.3189	0.3128	
CloudVision	BoN	0.1830	0.2310	0.2292	0.2420	0.3186	0.3177	
(Otani et al.,	BERT	0.1520	0.2014	0.2023	0.2653	0.3048	0.3051	
2018)	RoBERTa	0.1580	0.2006	0.2016	0.2898	0.3074	0.3075	



OCR

Text-based Image-based Blend **Ground truth** There's no lace like home sex:love;fashion nature;fun;love;sports None fun;safety;comfort;humor

CONVERSE



LIZER BIKE HELMETS

danger;safety;protection

danger;violence;death;safety

injury;safety



Put your shirt and join our team WWF nature;environment danger;violence;death;love;injury nature;danger;death;environment nature environment





Takeaway point from this example

- Key takeaway: a simple straightforward approach can be more >useful than complex models (in this case, multimodal joint embeddings)
- This paper was published at COLING 2020 >

Conference on Computational

CSSLING The 28th International Computation Ad Lingua: Text Classification Improves Symbolism Prediction in Image Advertisements

Andrey Savchenko^{1,3}, Anton Alekseev¹, Sejeong Kwon², Elena Tutubalina¹, Evgeniy Miasnikov¹ and Sergey Nikolenko^{1,4}

Interestingly, at the same conference we saw... >





Empirically Multimodally Additive Projection

- (Hessel, Lee, 2020):
 "Does my multimodal model learn cross-modal interactions? It's harder to tell than you might think!"
- They try to approximate multimodal models with
 linear image/text ensembles, removing cross-modal interactions
- And they find very little or no loss in performance!







So what's next? Let's find out together!



THANK YOU FOR YOUR ATTENTION!