DSC291: Machine Learning with Few Labels

Self-supervised Learning

Zhiting Hu Lecture 8, April 24, 2025



HALICIOĞLU DATA SCIENCE INSTITUTE

Outline

- Self-Supervised Learning (SSL)
 - Contrastive learning

- Paper presentation:
 - Kaiming Tao, Wenqi Li: "Transformers without Normalization"

"X"-supervised learning

- Supervised learning
- Unsupervised learning
- Self-supervised learning ~ 2018
- Weakly-/distantly-supervised learning

X

- Semi-supervised learning
- • •

Construction,

Self-Supervised Learning

- Given an observed data instance t
- One could derive various supervision signals based on the structure of the data
- By applying a "split" function that artificially partition t into two parts
 - $\circ (\mathbf{x}, \mathbf{y}) = split(\mathbf{t})$
 - sometimes split in a stochastic way
- Treat x as the input and y as the output
- Train a model $p_{\theta}(\mathbf{y}|\mathbf{x})$







Self-Supervised Learning: Examples

- Predict any part of the input from any other part.
- Predict the future from the past.

Predict the future from the recent past.

- Predict the past from the present.
- Predict the top from the bottom.



[Courtesy: Lecun "Self-supervised Learning"]

Self-Supervised Learning: Examples

- Predict any part of the input from any other part.
- Predict the future from the past.
- Predict the future from the recent past.
- Predict the past from the present.
- Predict the top from the bottom.
- Predict the occluded from the visible
- Pretend there is a part of the input you don't know and predict that.



[Courtesy: Lecun "Self-supervised Learning"]

Self-Supervised Learning: Motivation (I)

Our brains do this all the time

- Filling in the visual field at the retinal blind spot
- Filling in occluded images, missing segments in speech
- Predicting the state of the world from partial (textual) descriptions
- Predicting the consequences of our actions
- Predicting the sequence of actions leading to a result
- Predicting any part of the past, present or future percepts from whatever information is available.



[Courtesy: Lecun "Self-supervised Learning"]

Self-Supervised Learning: Motivation (I)

- Successfully learning to predict everything from everything else would result in the accumulation of lots of background knowledge about how the world works
- The model is forced to learn what we really care about, e.g. a semantic representation, in order to solve the prediction problem

[Courtesy: Lecun "Self-supervised Learning"] [Courtesy: Zisserman "Self-supervised Learning"]

Self-Supervised Learning: Motivation (II)

- The machine predicts any part of its input from any observed part
 - A lot of supervision signals in each data instance
- Untapped/availability of vast numbers of unlabeled text/images/videos..
 - Facebook: one billion images uploaded per day
 - 300 hours of video are uploaded to YouTube every minute

Self-Supervised Learning (SSL): Examples

- SSL from text
- SSL from images
- SSL from videos

Self-Supervised Learning from Text

Examples:

- Language models
- Learning contextual text representations

Language Models

- Calculates the probability of a sentence:
 - Sentence:



Language Models

- Calculates the probability of a sentence:
 - Sentence: Ο

 $\mathbf{y} = (y_1, y_2, \dots, y_T)$

Example:

(*I*, *like*, *this*, ...)





Language Models: Training

- Given data example y^*
- Minimizes negative log-likelihood of the data

$$\min_{\theta} \mathcal{L}(\theta) = -\log p_{\theta}(\boldsymbol{y}^*) = -\prod_{t=1}^{T} p_{\theta}(\boldsymbol{y}^*_t \mid \boldsymbol{y}^*_{1:t-1})$$

• Next word prediction

Self-Supervised Learning from Text

Examples:

- Language models
- Learning contextual text representations







BERT: Pre-training with Self-supervised Learning

0.1% Aardvark Masked LM Use the output of the Possible classes: . . . masked word's position All English words 10% Improvisation to predict the masked word . . . 0% Zyzzyva FFNN + Softmax 512 2 3 8 4 5 6 BERT Randomly mask . . . 512 2 8 3 15% of tokens Let's stick this skit [CLS] to in Input skit stick to improvisation in this [CLS] Let's

BERT: Downstream Fine-tuning

• Use BERT for sentence classification



BERT Results

• Huge improvements over SOTA on 12 NLP task

System	MNLI-(m/mm)	QQP	QNLI	SST-2	CoLA	STS-B	MRPC	RTE	Average
	392k	363k	108k	67k	8.5k	5.7k	3.5k	2.5k	-
Pre-OpenAI SOTA	80.6/80.1	66.1	82.3	93.2	35.0	81.0	86.0	61.7	74.0
BiLSTM+ELMo+Attn	76.4/76.1	64.8	79.9	90.4	36.0	73.3	84.9	56.8	71.0
OpenAI GPT	82.1/81.4	70.3	88.1	91.3	45.4	80.0	82.3	56.0	75.2
BERT _{BASE}	84.6/83.4	71.2	90.1	93.5	52.1	85.8	88.9	66.4	79.6
BERTLARGE	86.7/85.9	72.1	91.1	94.9	60.5	86.5	89.3	70.1	81.9

Table 1: GLUE Test results, scored by the GLUE evaluation server. The number below each task denotes the number of training examples. The "Average" column is slightly different than the official GLUE score, since we exclude the problematic WNLI set. OpenAI GPT = (L=12, H=768, A=12); BERT_{BASE} = (L=12, H=768, A=12); BERT_{LARGE} = (L=24, H=1024, A=16). BERT and OpenAI GPT are single-model, single task. All results obtained from https://gluebenchmark.com/leaderboard and https://blog.openai.com/language-unsupervised/.

Self-Supervised Learning (SSL): Examples

- SSL from text
- SSL from images
- SSL from videos

SSL from Images, EX (I): relative positioning

Train network to predict relative position of two regions in the same image



[Courtesy: Zisserman "Self-supervised Learning"]

SSL from Images, EX (I): relative positioning



22

SSL from Images, EX (I): relative positioning Evaluation: PASCAL VOC Detection

• 20 object classes (car, bicycle, person, horse ...)

• Predict the bounding boxes of all objects of a given class in an image (if any)



[Courtesy: Zisserman "Self-supervised Learning"]

SSL from Images, EX (I): relative positioning Evaluation: PASCAL VOC Detection

- Pre-train CNN using self-supervision (no labels)
- Train CNN for detection in R-CNN object category detection pipeline

R-CNN



[Girshick et al. 2014]

SSL from Images, EX (I): relative positioning Evaluation: PASCAL VOC Detection



[Courtesy: Zisserman "Self-supervised Learning"]

SSL from Images, EX (II): colorization

Train network to predict pixel colour from a monochrome input





[Courtesy: Zisserman "Self-supervised Learning"]

Colorful Image Colorization, Zhang et al., ECCV 2016

SSL from Images, EX (II): colorization

Train network to predict pixel colour from a monochrome input



[Courtesy: Zisserman "Self-supervised Learning"]

Colorful Image Colorization, Zhang et al., ECCV 2016

SSL from Images, EX (III): exemplar networks

- Exemplar Networks (Dosovitskiy et al., 2014)
- Perturb/distort image patches, e.g. by cropping and affine transformations
- Train to classify these exemplars as same class



[Courtesy: Zisserman "Self-supervised Learning"]

SSL from Images, EX (IV): masked autoencoder (MAE)



Figure 1. **Our MAE architecture**. During pre-training, a large random subset of image patches (*e.g.*, 75%) is masked out. The encoder is applied to the small subset of *visible patches*. Mask tokens are introduced *after* the encoder, and the full set of encoded patches and mask tokens is processed by a small decoder that reconstructs the original image in pixels. After pre-training, the decoder is discarded and the encoder is applied to uncorrupted images (full sets of patches) for recognition tasks.

SSL from Images, EX (IV): masked autoencoder (MAE)



random subset of image patches (e.g., 75%) is masked out. The encoder is applied to the small subset of visible patches. Mask tokens are introduced after the encoder, and the full set of encoded patches and mask tokens is processed by a small decoder that reconstructs the original image in pixels. After pre-training, the decoder is discarded and the encoder is applied to uncorrupted images (full sets of patches) for recognition tasks.

Question: Why is this (75%)

in BERT (15%)?

much larger than the mask rate

Question: What're your ideas of SSL from videos?



"Sequence" of data

Four example tasks:

- Video sequence order
 - Sequential Verification: Is this a valid sequence?









"Sequence" of data

[Courtesy: Zisserman "Self-supervised Learning"]

Wei et al., 2018 Arrow of Time 32

Four example tasks:

- Video sequence order
 - Sequential Verification: Is this a valid sequence?
- Video direction
 - Predict if video playing forwards or backwards

Four example tasks:

- Video sequence order
 - Sequential Verification: Is this a valid sequence?
- Video direction
 - Predict if video playing forwards or backwards
- Video tracking
 - Given a color video, colorize all frames of a gray scale version using a reference frame



[Courtesy: Zisserman "Self-supervised Learning"]



Vondrec et al., 2018

Four example tasks:

- Video sequence order
 - Sequential Verification: Is this a valid sequence?
- Video direction
 - Predict if video playing forwards or backwards
- Video tracking
 - Given a color video, colorize all frames of a gray scale version using a reference frame
- Video next frame prediction

[Courtesy: Zisserman "Self-supervised Learning"]

Vondric et al., 2018 35

Key Takeaways

- Self supervision learning
 - Predicting any part of the observations given any available information
 - The prediction task forces models to learn semantic representations
 - Massive/unlimited data supervisions
- SSL for text:
 - Language models: next word prediction
 - BERT text representations: masked language model (MLM)
 - SSL for images/videos:
 - Various ways of defining the prediction task

Questions?