DSC291: Machine Learning with Few Labels

Contrastive Learning Data Manipulation

Zhiting Hu Lecture 9, April 29, 2025



HALICIOĞLU DATA SCIENCE INSTITUTE

Logistics

- Homework 1 released
 - Due: May 23, Friday

Outline

• Contrastive learning

- Paper presentation:
 - Andrii Dovhaniuk, Wendi Tan: "Machine learning center-specific models"

Recap: 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.

Future \rightarrow

Time

 \leftarrow Past

Present



• Take a data example x, sample a "positive" sample x_{pos} and "negative" samples x_{neg} in some way

"positive" sample:

Data of the same labels Data of the same pseudo-labels

- Augmented/distorted version of x
- Data that captures the same target from different views



"negative" sample:

- Randomly sampled data
- Hard negative sample mining



Common contrastive learning functions

- Contrastive loss (Chopra et al. 2005)
- Triplet loss (Schroff et al. 2015; FaceNet)
- Lifted structured loss (Song et al. 2015)
- Multi-class n-pair loss (Sohn 2016)



- Noise contrastive estimation ("NCE"; Gutmann & Hyvarinen 2010)
- InfoNCE (van den Oord, et al. 2018)
- Soft-nearest neighbors loss (Salakhutdinov & Hinton 2007, Frosst et al. 2019)



• SimCSE ("Simple Contrastive learning of Sentence Embeddings"; Gao et al. 2021)

- SimCSE ("Simple Contrastive learning of Sentence Embeddings"; Gao et al. 2
 - Predict a sentence from itself with only dropout noise
 - One sentence gets two different versions of dropout augmentations



Figure 1: (a) Unsupervised SimCSE predicts the input sentence itself from in-batch negatives, with different hidden dropout masks applied.

- SimCSE ("Simple Contrastive learning of Sentence Embeddings"; Gao et al. 2021)
 - Predict a sentence from itself with only dropout noise
 - One sentence gets two different versions of dropout augmentations



(b) Supervised SimCSE

Question

Figure 1: (a) Unsupervised SimCSE predicts the input sentence itself from in-batch negatives, with different hidden dropout masks applied.

Sim

SimCSE ("Simple Contrastive learning of Sentence Embeddings"; Gao et al. 2021

- Predict a sentence from itself with only dropout noise
- One sentence gets two different versions of dropout augmentations



Figure 1: (a) Unsupervised SimCSE predicts the input sentence itself from in-batch negatives, with different hidden dropout masks applied. (b) Supervised SimCSE leverages the NLI datasets and takes the entailment (premise-hypothesis) pairs as positives, and contradiction pairs as well as other in-batch instances as negatives.

Contrastive learning: Ex 3 – InfoNCE (Noise-Contrastive Estimation)

- The CPC model
 - C_t: context representation from history \bigcirc
 - x_{t+k} (or z_{t+k}): future target \bigcirc



InfoNCE loss

- Define scoring function $f_k > 0$
- The InfoNCE loss:
 - Given $X = \{$ one positive sample from $p(x_{t+k} | c_t), N 1$ negative samples from the negative sampling distribution $p(x_{t+k}) \}$

$$\mathcal{L}_{N} = -\mathbb{E}_{X} \left[\log \frac{f_{k}(x_{t+k}, c_{t})}{\sum_{x_{j} \in X} f_{k}(x_{j}) c_{t}} \right]$$

Г

• InfoNCE is interesting because it's effectively maximizing the more than x_{t+k} and x_{t+k} next the product of the pr

InfoNCE loss

- Define scoring function $f_k > 0$
- The InfoNCE loss:
 - Given $X = \{$ one positive sample from $p(x_{t+k} | c_t), N 1$ negative samples from the negative sampling distribution $p(x_{t+k}) \}$

$$\mathcal{L}_{\mathrm{N}} = - \mathop{\mathbb{E}}_{X} \left[\log \frac{f_k(x_{t+k}, c_t)}{\sum_{x_j \in X} f_k(x_j, c_t)}
ight]$$

• InfoNCE is interesting because it's effectively maximizing the mutual information between c_t and x_{t+k}

16

Mutual Information (MI)

• How much is our uncertainty about x reduced by knowing c?

$$I(x;c) = \sum_{x,c} p(x,c) \log \frac{p(x,c)}{p(x)p(c)} = \sum_{x,c} p(x,c) \log \frac{p(x|c)}{p(x)}$$

= $H(x) + H(c) - H(x,c)$
= $H(x) - H(x|c)$
= $KL(p(x,c) || p(x)p(c))$

Minimizing InfoNCE <=> Maximzing MI

• InfoNCE loss

$$\mathcal{L}_{\mathrm{N}} = - \mathop{\mathbb{E}}_{X} \left[\log \frac{f_k(x_{t+k}, c_t)}{\sum_{x_j \in X} f_k(x_j, c_t)}
ight]$$

$$I(x_{t+k}, c_t) \ge \log(N) - \mathcal{L}_N$$

Minimizing InfoNCE <=> Maximzing MI

• InfoNCE loss

$$\mathcal{L}_{\mathrm{N}} = - \mathop{\mathbb{E}}_{X} \left[\log \frac{f_k(x_{t+k}, c_t)}{\sum_{x_j \in X} f_k(x_j, c_t)}
ight]$$

• The loss is optimized when

$$f_k(x_{t+k}, c_t) \propto rac{p(x_{t+k}|c_t)}{p(x_{t+k})}$$

• Proof:

$$p(sample \ i \ is \ positive | X, c_t) = \frac{p(x_i | c_t) \prod_{l \neq i} p(x_l)}{\sum_{j=1}^N p(x_j | c_t) \prod_{l \neq j} p(x_l)}$$
$$= \frac{\frac{p(x_i | c_t)}{p(x_i)}}{\sum_{j=1}^N \frac{p(x_j | c_t)}{p(x_j)}}.$$

$$\mathcal{L}_{\mathrm{N}} = - \mathop{\mathbb{E}}_{X} \left[\log rac{f_k(x_{t+k}, c_t)}{\sum_{x_j \in X} f_k(x_j, c_t)}
ight]$$



$$\mathcal{L}_{\mathrm{N}} = - \mathop{\mathbb{E}}_{X} \left[\log rac{f_k(x_{t+k}, c_t)}{\sum_{x_j \in X} f_k(x_j, c_t)}
ight]$$

$$\begin{split} \mathcal{L}_{\mathrm{N}}^{\mathrm{opt}} &= -\mathop{\mathbb{E}}_{X} \log \left[\frac{\frac{p(x_{t+k}|c_{t})}{p(x_{t+k})}}{\frac{p(x_{t+k}|c_{t})}{p(x_{t+k})} + \sum_{x_{j} \in X_{\mathrm{neg}}} \frac{p(x_{j}|c_{t})}{p(x_{j})}} \right] \\ &= \mathop{\mathbb{E}}_{X} \log \left[1 + \frac{p(x_{t+k})}{p(x_{t+k}|c_{t})} \sum_{x_{j} \in X_{\mathrm{neg}}} \frac{p(x_{j}|c_{t})}{p(x_{j})} \right] \\ \end{split}$$

$$\mathcal{L}_{\mathrm{N}} = - \mathop{\mathbb{E}}_{X} \left[\log \frac{f_k(x_{t+k}, c_t)}{\sum_{x_j \in X} f_k(x_j, c_t)}
ight]$$

$$\begin{split} \mathcal{L}_{\mathrm{N}}^{\mathrm{opt}} &= -\mathop{\mathbb{E}}_{X} \log \left[\frac{\frac{p(x_{t+k}|c_t)}{p(x_{t+k})}}{\frac{p(x_{t+k}|c_t)}{p(x_{t+k})} + \sum_{x_j \in X_{\mathrm{neg}}} \frac{p(x_j|c_t)}{p(x_j)}}{p(x_j)} \right] \\ &= \mathop{\mathbb{E}}_{X} \log \left[1 + \frac{p(x_{t+k})}{p(x_{t+k}|c_t)} \sum_{x_j \in X_{\mathrm{neg}}} \frac{p(x_j|c_t)}{p(x_j)} \right] \\ &\approx \mathop{\mathbb{E}}_{X} \log \left[1 + \frac{p(x_{t+k})}{p(x_{t+k}|c_t)} (N-1) \mathop{\mathbb{E}}_{x_j} \frac{p(x_j|c_t)}{p(x_j)} \right] \end{split}$$

This approximation becomes more accurate as N increases, so it is preferable to use large negative samples

$$\mathcal{L}_{ ext{N}} = - \mathop{\mathbb{E}}\limits_{X} \left[\log rac{f_k(x_{t+k}, c_t)}{\sum_{x_j \in X} f_k(x_j, c_t)}
ight]$$

$$\begin{split} \mathcal{L}_{N}^{\text{opt}} &= -\mathop{\mathbb{E}}_{X} \log \left[\frac{\frac{p(x_{t+k}|c_t)}{p(x_{t+k})}}{\frac{p(x_{t+k}|c_t)}{p(x_{t+k})} + \sum_{x_j \in X_{\text{neg}}} \frac{p(x_j|c_t)}{p(x_j)}}{p(x_j)} \right] \\ &= \mathop{\mathbb{E}}_{X} \log \left[1 + \frac{p(x_{t+k})}{p(x_{t+k}|c_t)} \sum_{x_j \in X_{\text{neg}}} \frac{p(x_j|c_t)}{p(x_j)} \right] \\ &\approx \mathop{\mathbb{E}}_{X} \log \left[1 + \frac{p(x_{t+k})}{p(x_{t+k}|c_t)} (N - 1) \mathop{\mathbb{E}}_{x_j} \frac{p(x_j|c_t)}{p(x_j)} \right] = 1 \\ &= \mathop{\mathbb{E}}_{X} \log \left[1 + \frac{p(x_{t+k})}{p(x_{t+k}|c_t)} (N - 1) \right] \end{split}$$

$$\mathcal{L}_{ ext{N}} = - \mathop{\mathbb{E}}\limits_{X} \left[\log rac{f_k(x_{t+k}, c_t)}{\sum_{x_j \in X} f_k(x_j, c_t)}
ight]$$

$$\begin{split} \mathcal{L}_{\mathrm{N}}^{\mathrm{opt}} &= - \mathop{\mathbb{E}}_{X} \log \left[\frac{\frac{p(x_{t+k}|c_t)}{p(x_{t+k})}}{\frac{p(x_{t+k}|c_t)}{p(x_{t+k})} + \sum_{x_j \in X_{\mathrm{neg}}} \frac{p(x_j|c_t)}{p(x_j)}}{p(x_j)} \right] \\ &= \mathop{\mathbb{E}}_{X} \log \left[1 + \frac{p(x_{t+k})}{p(x_{t+k}|c_t)} \sum_{x_j \in X_{\mathrm{neg}}} \frac{p(x_j|c_t)}{p(x_j)} \right] \\ &\approx \mathop{\mathbb{E}}_{X} \log \left[1 + \frac{p(x_{t+k})}{p(x_{t+k}|c_t)} (N-1) \mathop{\mathbb{E}}_{x_j} \frac{p(x_j|c_t)}{p(x_j)} \right] \\ &= \mathop{\mathbb{E}}_{X} \log \left[1 + \frac{p(x_{t+k})}{p(x_{t+k}|c_t)} (N-1) \right] \\ &\geq \mathop{\mathbb{E}}_{X} \log \left[\frac{p(x_{t+k})}{p(x_{t+k}|c_t)} N \right] \\ &= -I(x_{t+k}, c_t) + \log(N), \end{split}$$

$$\mathcal{L}_{ ext{N}} = - \mathop{\mathbb{E}}\limits_{X} \left[\log rac{f_k(x_{t+k}, c_t)}{\sum_{x_j \in X} f_k(x_j, c_t)}
ight]$$

$$I(x_{t+k}, c_t) \ge \log(N) - \mathcal{L}_{N_t}$$

Key Takeaways: Contrastive learning

- Contrastive learning is a way of doing self-supervised learning
- Positive samples, negative samples
- Mutual information

$$I(x;c) = \sum_{x,c} p(x,c) \log \frac{p(x,c)}{p(x)p(c)} = \sum_{x,c} p(x,c) \log \frac{p(x|c)}{p(x)}$$
$$= H(x) + H(c) - H(x,c)$$
$$= H(x) + H(x|c)$$
$$= KL(p(x,c) || p(x)p(c))$$

Data Manipulation

Data manipulation

- Data augmentation
- Om Pue Applies label-preserving transformations on original data points to expand the data \bigcirc size
- Data reweighting
 - Assigns an importance weight to each instance to adapt its effect on learning \bigcirc

Shasp

Data synthesis

Generates entire artificial examples

- Curriculum learning
 - Makes use of data instances in an order based on "difficulty" \bigcirc

Onenial leanny

made (

Questions?