

Representation Constrastive Learning

Learning by comparison

Pablo Miralles González

October 23, 2023

Representation learning

Contrastive learning

Loss functions for contrastive learning

Generating data

Discussion on negative examples

Applications

Conclusions

2023-10-23

Representation Contrastive Learning

Representation learning

Contrastive learning

Loss functions for contrastive learning

Generating data

Discussion on negative examples

Applications

Conclusions

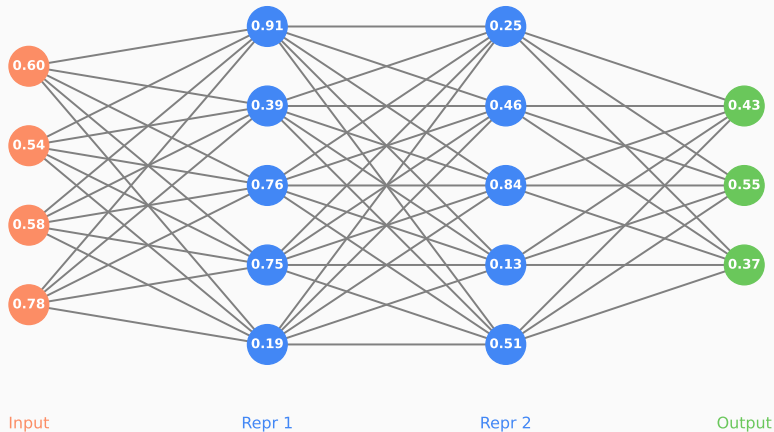
2023-10-23

Representation Constrastive Learning
└ Representation learning

Representation learning

Representation learning

What do we mean by representation?

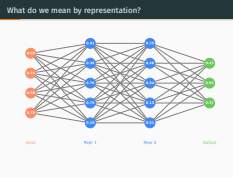


2023-10-23

Representation Contrastive Learning

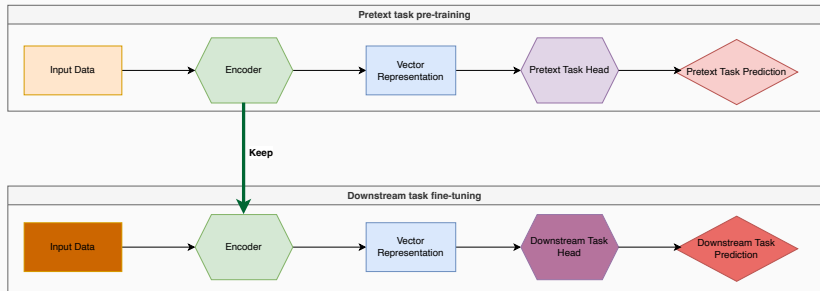
└ Representation learning

└ What do we mean by representation?



- Internal layers outputs can be viewed as different views or representations of the input data.
- They contain meaningful features for the task.

Why should we learn representations?



2023-10-23

Representation Contrastive Learning

└ Representation learning

└ Why should we learn representations?

Why should we learn representations?



- The possibility of transfer learning.
- Train for a complex pretext task to obtain a very general representation of the input data.
- Keep encoder, change head and fine-tune and use for downstream tasks.

Why should we learn representations?

Representation Contrastive Learning

└ Representation learning

└ Why should we learn representations?

1. We might not have enough data for the downstream task. We select pretext tasks for which we can generate data easily.

- Overcome data bottlenecks

2023-10-23

Why should we learn representations?

- Overcome data bottlenecks
- Outsource compute resources for training

Representation Contrastive Learning

└ Representation learning

└ Why should we learn representations?

1. People with greater resources can pre-train and upload weights. We can download and fine-tune for any task we want. If the encoder weights are frozen, it is much cheaper: we only update a smaller head. This allows people with fewer resources to use bigger models than they could otherwise.

2023-10-23

- Overcome data bottlenecks
- Outsource compute resources for training
- Better generalization

Why should we learn representations?

- Overcome data bottlenecks
- Outsource compute resources for training
- Better generalization

Representation Contrastive Learning

└ Representation learning

└ Why should we learn representations?

1. If the data for the downstream task is not representative, we might learn spurious correlations. By pre-training for a complex task with rich data, we make sure the model understand the latent distribution correctly. Still, fine-tuning might lead to representational collapse.

2023-10-23

- Overcome data bottlenecks
- Outsource compute resources for training
- Better generalization
- Zero-shot capabilities

Why should we learn representations?

- Overcome data bottlenecks
- Outsource compute resources for training
- Better generalization
- Zero-shot capabilities

Representation Contrastive Learning

└ Representation learning

└ Why should we learn representations?

1. For example, ChatGPT (see e.g. text classification, text transformations, code generation, code analysis...) or CLIP (zero-shot image classification with arbitrary classes).

2023-10-23

Why should we not learn representations?

Representation Contrastive Learning

└ Representation learning

└ Why should we not learn representations?

2023-10-23

• Sometimes overkill

- Sometimes overkill

1. If downstream task is very simple and the data is decent, we just don't need to.

Why should we not learn representations?

Representation Contrastive Learning

└ Representation learning

└ Why should we not learn representations?

1. It is very expensive and not everyone can pre-train massive deep learning models.

- Sometimes overkill
- Massive compute requirements

2023-10-23

- Sometimes overkill
- Massive compute requirements
- Generally poor zero-shot performance

Why should we not learn representations?

Representation Contrastive Learning

└ Representation learning

└ Why should we not learn representations?

1. Without fine-tuning, models that can perform zero-shot predictions are unlikely to perform very well.

- Sometimes overkill
- Massive compute requirements
- Generally poor zero-shot performance

2023-10-23

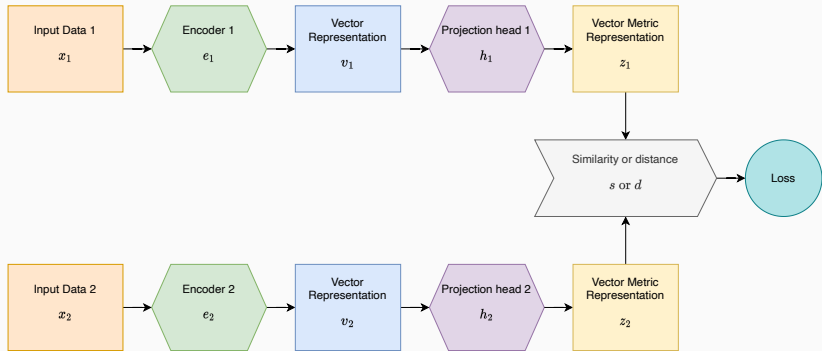
2023-10-23

Representation Contrastive Learning
└ Contrastive learning

Contrastive learning

Contrastive learning

Contrastive learning: learn by comparison



2023-10-23

Representation Contrastive Learning

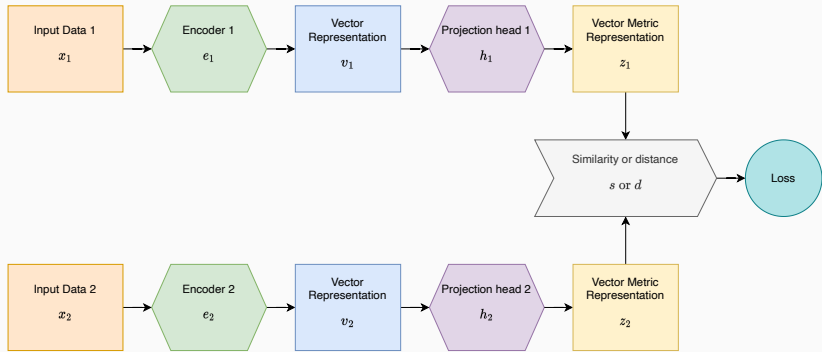
└ Contrastive learning

└ Contrastive learning: learn by comparison

Contrastive learning: learn by comparison



Contrastive learning: learn by comparison



Similar instances \implies close together

Dissimilar instances \implies far apart

2023-10-23

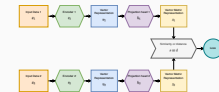
Representation Contrastive Learning

└ Contrastive learning

└ Contrastive learning: learn by comparison

- Go over full diagram.
- Input data might be of different modalities.
- For data of the same modality, we can use the same encoder and head with tied weights.

Contrastive learning: learn by comparison



Similar instances \implies close together
Dissimilar instances \implies far apart

Examples of distance and similarity

Representation Contrastive Learning

└ Contrastive learning

└ Examples of distance and similarity

2023-10-23

- Euclidean distance $d(z_1, z_2) = \|z_1 - z_2\|$
- Cosine similarity $s(z_1, z_2) = \frac{\langle z_1, z_2 \rangle}{\|z_1\| \|z_2\|}$

- Euclidean distance $d(z_1, z_2) = \|z_1 - z_2\|$
- Cosine similarity $s(z_1, z_2) = \frac{\langle z_1, z_2 \rangle}{\|z_1\| \|z_2\|}$

2023-10-23

Representation Contrastive Learning

└ Loss functions for contrastive learning

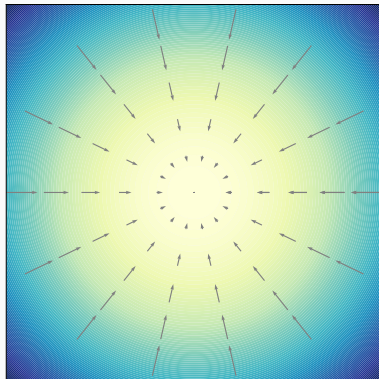
Loss functions for contrastive learning

Loss functions for contrastive learning

Pair loss¹

$$\begin{cases} \mathcal{L}(x, x^+) &= D(z, z^+)^2 \\ \mathcal{L}(x, x^-) &= \max(0, \varepsilon - D(z, z^-))^2, \end{cases}$$

Loss function for positive examples



¹Chopra, Hadsell, and LeCun, "Learning a Similarity Metric Discriminatively, with Application to Face Verification".

2023-10-23

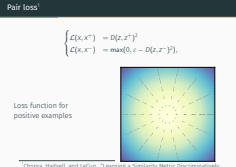
Representation Contrastive Learning

└ Loss functions for contrastive learning

└ Pair loss^a

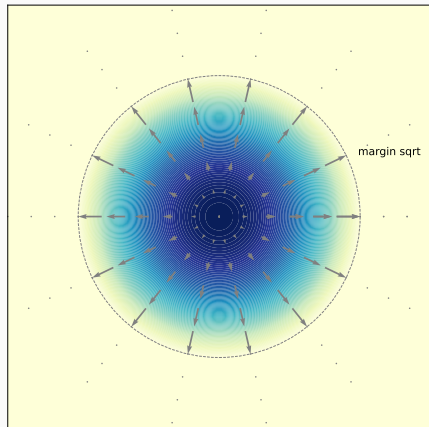
^aChopra, Hadsell, and LeCun, "Learning a Similarity Metric Discriminatively, with Application to Face Verification".

- Different for positive and negative examples.
- Explain plot.



$$\begin{cases} \mathcal{L}(x, x^+) = D(z, z^+)^2 \\ \mathcal{L}(x, x^-) = \max(0, \varepsilon - D(z, z^-))^2, \end{cases}$$

Loss function for negative examples



Representation Contrastive Learning

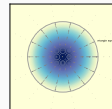
- Loss functions for contrastive learning

- Pair loss

- Dissimilar instances separated by margin ε .

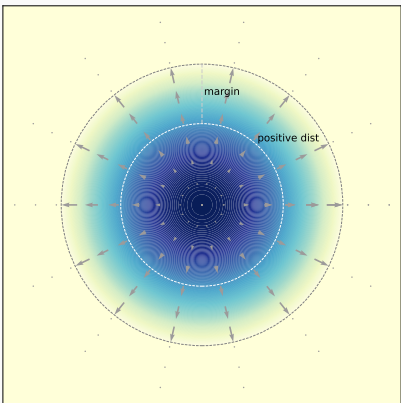
$$\begin{cases} \mathcal{L}(x, x^+) = D(z, z^+)^2 \\ \mathcal{L}(x, x^-) = \max(0, \varepsilon - D(z, z^-))^2, \end{cases}$$

Loss function for negative examples



Triplet loss⁴

$$\mathcal{L}(x, x^+, x^-) = \max(0, D(z, z^+)^2 - D(z, z^-)^2 + \epsilon)$$



⁴Schroff, Kalenichenko, and Philbin, "FaceNet: A Unified Embedding for Face Recognition and Clustering".

2023-10-23

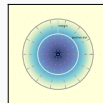
Representation Constrastive Learning

└ Loss functions for contrastive learning

└ Triplet loss^a

Triplet loss^a

$$\mathcal{L}(x, x^+, x^-) = \max(0, D(z, z^+)^2 - D(z, z^-)^2 + \epsilon)$$



^aSchroff, Kalenichenko, and Philbin, "FaceNet: A Unified Embedding for Face Recognition and Clustering".

- Distance between similar and dissimilar instances separated by margin.

$\{x_i\}_{i=1}^n = \text{set of examples}$

$P = \{\text{pairs of similar instances}\} \quad N = \{\text{pairs of dissimilar instances}\}$

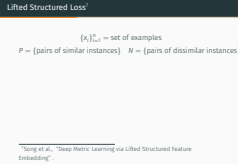
⁷Song et al., “Deep Metric Learning via Lifted Structured Feature Embedding” .

2023-10-23

Representation Contrastive Learning

└ Loss functions for contrastive learning

└ Lifted Structured Loss



$\{x_i\}_{i=1}^n$ = set of examples

$P = \{\text{pairs of similar instances}\}$ $N = \{\text{pairs of dissimilar instances}\}$

$$\mathcal{L}(N, P) = \frac{1}{2|P|} \sum_{(i,j) \in P} L_{i,j}^2$$

$$L_{i,j} = D_{i,j} + \log \left(\sum_{(i,k) \in N} e^{\varepsilon - D_{i,k}} + \sum_{(j,l) \in N} e^{\varepsilon - D_{j,l}} \right)$$

$$D_{i,j} = D(z_i, z_j)$$

└ Loss functions for contrastive learning

└ Lifted Structured Loss

Lifted Structured Loss

$\{x_i\}_{i=1}^n$ = set of examples
 $P = \{\text{pairs of similar instances}\}$ $N = \{\text{pairs of dissimilar instances}\}$

$$\mathcal{L}(N, P) = \frac{1}{2|P|} \sum_{(i,j) \in P} L_{i,j}^2$$
$$L_{i,j} = D_{i,j} + \log \left(\sum_{(i,k) \in N} e^{\varepsilon - D_{i,k}} + \sum_{(j,l) \in N} e^{\varepsilon - D_{j,l}} \right)$$
$$D_{i,j} = D(z_i, z_j)$$

$\{x_i\}_{i=1}^n$ = set of examples

$P = \{\text{pairs of similar instances}\}$ $N = \{\text{pairs of dissimilar instances}\}$

$$\mathcal{L}(N, P) = \frac{1}{2|P|} \sum_{(i,j) \in P} L_{i,j}^2$$

$$\underbrace{L_{i,j}}_{\text{smooth}} \geq \hat{L}_{i,j} = D_{i,j} + \max \left(\max_{(i,k) \in N} \varepsilon - D_{i,k}, \max_{(j,l) \in N} \varepsilon - D_{j,l} \right)$$

$$D_{i,j} = D(z_i, z_j)$$

└ Loss functions for contrastive learning

└ Lifted Structured Loss

- We are actually penalizing the small differences between the distance with a positive example and the hardest negative example, up to some margin, similar to the triplet loss.

Lifted Structured Loss

$\{x_i\}_{i=1}^n$ = set of examples
 $P = \{\text{pairs of similar instances}\}$ $N = \{\text{pairs of dissimilar instances}\}$

$$\mathcal{L}(N, P) = \frac{1}{2|P|} \sum_{(i,j) \in P} L_{i,j}^2$$

$$\underbrace{L_{i,j}}_{\text{smooth}} \geq \hat{L}_{i,j} = D_{i,j} + \max \left(\max_{(i,k) \in N} \varepsilon - D_{i,k}, \max_{(j,l) \in N} \varepsilon - D_{j,l} \right)$$

$$D_{i,j} = D(z_i, z_j)$$

$$X(x_1, x_2) = \begin{cases} 1 & \text{if } x_1 \text{ and } x_2 \text{ are similar} \\ 0 & \text{if } x_1 \text{ and } x_2 \text{ are dissimilar} \end{cases}$$

$$X(x_1, x_2) \sim P(\cdot | x_1, x_2) \quad \rightarrow \quad P(1 | x_1, x_2) = \sigma(s(z_1, z_2))$$

⁷Gutmann and Hyvärinen, “Noise-Contrastive Estimation: A New Estimation Principle for Unnormalized Statistical Models” .

$$X(x_1, x_2) = \begin{cases} 1 & \text{if } x_1 \text{ and } x_2 \text{ are similar} \\ 0 & \text{if } x_1 \text{ and } x_2 \text{ are dissimilar} \end{cases}$$

$$X(x_1, x_2) \sim P(\cdot | x_1, x_2) \quad \rightarrow \quad P(1 | x_1, x_2) = \sigma(s(z_1, z_2))$$

$$X(x_1, x_2) = \begin{cases} 1 & \text{if } x_1 \text{ and } x_2 \text{ are similar} \\ 0 & \text{if } x_1 \text{ and } x_2 \text{ are dissimilar} \end{cases}$$

$$X(x_1, x_2) \sim P(\cdot | x_1, x_2) \rightarrow P(1 | x_1, x_2) = \sigma(s(z_1, z_2))$$

$$\begin{aligned} \mathcal{L}_{\text{Bin-NCE}} &= -\mathbb{E}_{p^+} \log P(1 | x_1, x_2) - \mathbb{E}_{p^-} \log(1 - P(1 | x_1, x_2)) \approx \\ &= -\frac{1}{|P|} \sum_{(i,j) \in P} \log \sigma(s(z_i, z_j)) - \frac{1}{|N|} \sum_{(i,j) \in N} \log(1 - \sigma(s(z_i, z_j))) \end{aligned}$$

Loss functions for contrastive learning

Binary Noise-Contrastive Estimation^a

- Explain probabilistic approach and population sample with batch.

2023-10-23

$$X(x_1, x_2) = \begin{cases} 1 & \text{if } x_1 \text{ and } x_2 \text{ are similar} \\ 0 & \text{if } x_1 \text{ and } x_2 \text{ are dissimilar} \end{cases}$$

$$X(x_1, x_2) \sim P(\cdot | x_1, x_2) \rightarrow P(1 | x_1, x_2) = \sigma(s(z_1, z_2))$$

$$\begin{aligned} \mathcal{L}_{\text{Bin-NCE}} &= -\mathbb{E}_{p^+} \log P(1 | x_1, x_2) - \mathbb{E}_{p^-} \log(1 - P(1 | x_1, x_2)) \approx \\ &= -\frac{1}{|P|} \sum_{(i,j) \in P} \log \sigma(s(z_i, z_j)) - \frac{1}{|N|} \sum_{(i,j) \in N} \log(1 - \sigma(s(z_i, z_j))) \end{aligned}$$

$x; S = \{x_0^+, x_1^-, \dots, x_n^-\} \rightarrow$ rank the positive one!

$$P(i|x, S) = \frac{\exp(s(x, x_i))}{\sum_{j=0}^n \exp(s(x, x_j))}$$

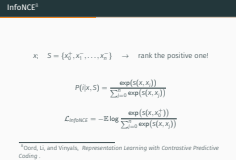
$$\mathcal{L}_{\text{InfoNCE}} = -\mathbb{E} \log \frac{\exp(s(x, x_0^+))}{\sum_{j=0}^n \exp(s(x, x_j))}$$

⁸Oord, Li, and Vinyals, *Representation Learning with Contrastive Predictive Coding*.

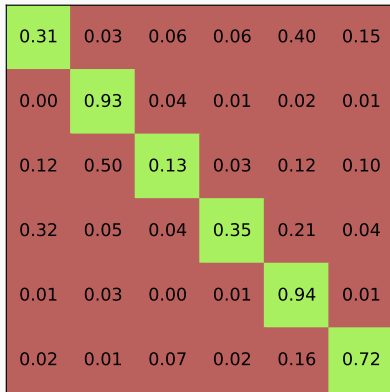
└ Loss functions for contrastive learning

└ InfoNCE^a

- Instance x , examples S , only x_0^+ positive.



$$B = \{(x_0, x'_0), (x_1, x'_1), \dots, (x_n, x'_n)\} \rightarrow \text{softmax}(S(z_i, z'_j))$$



Loss functions for contrastive learning

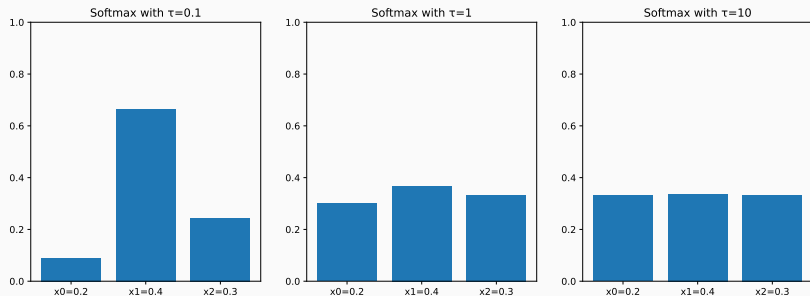
InfoNCE: example setting

2023-10-23



- Batches of pairs of similar instances $\{(x_1, x'_1), \dots, (x_n, x'_n)\}$.
- Instances across pairs are considered to be dissimilar.
- We can compute a similarity matrix $S = (s(z_i, z'_j))_{i,j}$, where the main diagonal values should be high and the rest should be low.
- We can calculate the InfoNCE across rows or columns. It is also possible to average both options, yielding a *symmetric InfoNCE* loss.

$$\mathcal{L}_{NT-Xent} = -\mathbb{E} \log \frac{\exp(s(x, x_0^+)/\tau)}{\sum_{j=0}^n \exp(s(x, x_j)/\tau)}$$



¹¹Chen et al., “A Simple Framework for Contrastive Learning of Visual Representations”.

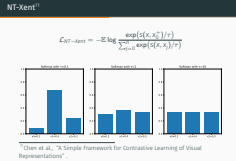
Representation Contrastive Learning

Loss functions for contrastive learning

NT-Xent^a

^aChen et al., “A Simple Framework for Contrastive Learning of Visual

- A small value of τ makes the softmax sharper, and small differences between the similarity of positive and negative examples already produces a high likelihood.
- A large value of τ forces the difference in similarity to be large.
- This parameter can be viewed as the margin parameter in previous functions.



2023-10-23

Generating data

Generating data for contrastive learning

Representation Contrastive Learning

2023-10-23

└ Generating data

└ Generating data for contrastive learning

Data = pairs of positive and negative examples.

Data = pairs of positive and negative examples.

Generating data for contrastive learning

Data = pairs of positive and negative examples.

- Human supervision.
- Data augmentation.
- Multi-sensor input.
- Local-global relationship.
- Sequential coherence/consistency.

2023-10-23

Representation Contrastive Learning

└ Generating data

└ Generating data for contrastive learning

Data = pairs of positive and negative examples.

- Human supervision.
- Data augmentation.
- Multi-sensor input.
- Local-global relationship.
- Sequential coherence/consistency.



2023-10-23

Representation Constrastive Learning

└ Generating data

└ Human supervision

Human supervision





Costly and painful!

2023-10-23

Representation Contrastive Learning

└ Generating data

└ Human supervision

Human supervision

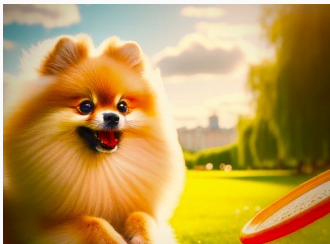


Costly and painful!

Specially for massive NN models

Data augmentation

Small modifications that don't alter anything meaningful.



2023-10-23

Representation Constrastive Learning

└ Generating data

└ Data augmentation

Data augmentation

Small modifications that don't alter anything meaningful.



- Augmented versions of similar instances are similar.
- Augmented versions of dissimilar instances are dissimilar.

- **Images.** Rotations, translations, cutouts, cropping, resizing...
- **Text.** More complex. E.g. back-translation, masking words, adding noise...

└ Generating data

└ Data augmentation

2023-10-23

- **Images.** Rotations, translations, cutouts, cropping, resizing...
- **Text.** More complex. E.g. back-translation, masking words, adding noise...

- **Images.** Rotations, translations, cutouts, cropping, resizing...
- **Text.** More complex. E.g. back-translation, masking words, adding noise...

Remark: performance can be highly sensitive to data augmentation strategies.¹⁴

¹⁴Chen et al., “A Simple Framework for Contrastive Learning of Visual Representations” .

└ Generating data

└ Data augmentation

- **Images.** Rotations, translations, cutouts, cropping, resizing...
- **Text.** More complex. E.g. back-translation, masking words, adding noise...

Remark: performance can be highly sensitive to data augmentation strategies.¹⁴

¹⁴Chen et al., “A Simple Framework for Contrastive Learning of Visual Representations” .

Multi-sensor input

Video (audio and image), multiple cameras, cameras and other sensors...



Representation Contrastive Learning

└ Generating data

└ Multi-sensor input

2023-10-23

Multi-sensor input

Video (audio and image), multiple cameras, cameras and other sensors...



Local-global relationship

Local and global features should be similarly represented.

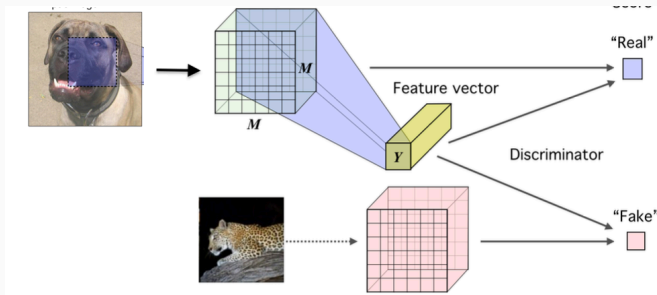


Figure 1: Local-global relationship in images¹⁵

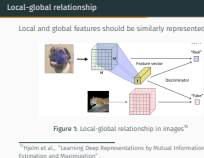
¹⁵Hjelm et al., "Learning Deep Representations by Mutual Information Estimation and Maximization".

Representation Constrastive Learning

└ Generating data

└ Local-global relationship

2023-10-23



The green box comes from an image of a dog, and the red box from an image of a leopard. The local features of the dog are aggregated into a global representation. The local and global features of the dog are similar, and the global features of the dog and the local features of the leopard are dissimilar.

Sequential coherence/consistency

Exploit continuity in smaller sub-units.

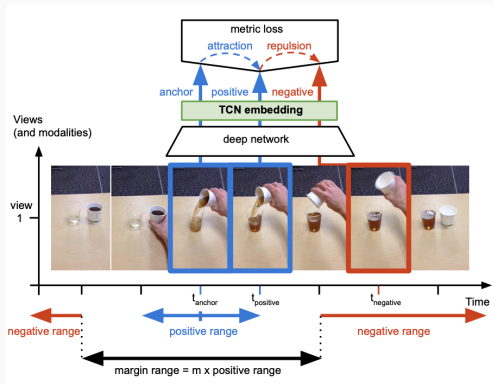


Figure 2: Example of consistency in videos¹⁶

¹⁶Sermanet et al., "Time-Contrastive Networks: Self-Supervised Learning from Video".

2023-10-23

Representation Contrastive Learning

└ Generating data

└ Sequential coherence/consistency

Sequential coherence/consistency

Exploit continuity in smaller sub-units.

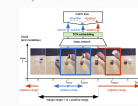


Figure 2: Example of consistency in videos¹⁶
¹⁶Sermanet et al., "Time-Contrastive Networks: Self-Supervised Learning from Video".

For example, in videos, images that are very close in time are likely to be similar, to contain the same concept. Images far apart in time are likely to be different.

2023-10-23

Representation Contrastive Learning
└ Discussion on negative examples

Discussion on negative examples

Discussion on negative examples

The importance of negative samples

- Negative examples prevent representational collapse $f \equiv \text{constant}$.
- Empirical evidence of better performance.¹⁷

Several discussions around the use of negative examples.

¹⁷Chen et al., “A Simple Framework for Contrastive Learning of Visual Representations” .

Representation Contrastive Learning

└ Discussion on negative examples

└ The importance of negative samples

2023-10-23

- Negative examples prevent representational collapse $f \equiv \text{constant}$.
- Empirical evidence of better performance.¹⁷

Several discussions around the use of negative examples.

¹⁷Chen et al., “A Simple Framework for Contrastive Learning of Visual Representations” .



Negative examples?

2023-10-23

Representation Contrastive Learning

└ Discussion on negative examples

└ False negatives in self-supervision



Negative examples?

False negatives in self-supervision

How to **remain self-supervised** and **mitigate bias** from **false negatives**?

Representation Contrastive Learning

└ Discussion on negative examples

└ False negatives in self-supervision

2023-10-23

False negatives in self-supervision

Representation Contrastive Learning

└ Discussion on negative examples

└ False negatives in self-supervision

2023-10-23

How to **remain self-supervised** and **mitigate bias** from **false negatives**?

Work like Debiased Contrastive Learning.¹⁸

Improved performance by increasing the number of positive examples for a given instances, and adding a lot of complexity.

¹⁸Chuang et al., "Debiased Contrastive Learning".

└ Discussion on negative examples

└ Hardware bottlenecks

2023-10-23

Number of negatives \approx Batch size \implies $(\# \text{ negatives})^2$ scaling

1. It is common to use other examples in the batch as negatives.
2. If we do this, we tie the batch size and the number of negatives.
3. Quadratic scaling of complexity with $\#$ of negatives.

Solution:

- Encoding of all samples stored offline.
- Negatives sampled from offline storage.
- $e_{\text{online}} \neq e_{\text{offline}}$

Number of negatives \approx Batch size \implies $(\# \text{ negatives})^2$ scaling

Solution:

- Encoding of all samples stored offline.
- Negatives sampled from offline storage.
- $e_{\text{online}} \neq e_{\text{offline}}$

1. It is common to use other examples in the batch as negatives.
2. If we do this, we tie the batch size and the number of negatives.
3. Quadratic scaling of complexity with $\#$ of negatives.

When and how to update offline encoder/encodings?

- All the samples after each checkpoint.¹⁹
- Queue of mini-batches and moving average.²⁰

¹⁹Wu et al., “Unsupervised Feature Learning via Non-parametric Instance Discrimination” .

²⁰He et al., “Momentum Contrast for Unsupervised Visual Representation Learning” .

└ Discussion on negative examples

└ Hardware bottlenecks

- @wu2018UnsupervisedFeature sampled negative representations randomly from a memory bank with the full dataset. At the end of each epoch, all the representations in the memory bank are updated with the new checkpoint of the model.
- @he2020MomentumContrast use a queue with a fixed number of mini-batches. After each mini-batch, the new examples are added, and the oldest mini-batch is removed. The queue is used to sample negative examples for the current mini-batch. They separated the online encoder that is being trained from an offline encoder that produces the representations for the queue for empirical reasons. The parameters of the offline encoder are updated through a momentum update rule with the parameters of the online one.

- All the samples after each checkpoint.¹⁹
- Queue of mini-batches and moving average.²⁰

Hard negative mining

Representation Contrastive Learning

└ Discussion on negative examples

└ Hard negative mining

2023-10-23

↑ # negatives \implies harder negatives \implies performance?

²¹Kalantidis et al., "Hard Negative Mining for Contrastive Learning"

↑ # negatives \implies harder negatives \implies performance?

²¹Kalantidis et al., "Hard Negative Mining for Contrastive Learning" .

Hard negative mining

Representation Contrastive Learning

└ Discussion on negative examples

└ Hard negative mining

2023-10-23

$\uparrow \# \text{ negatives} \implies \text{harder negatives} \implies \text{performance?}$

Some work on the topic.²¹

²¹Kalantidis et al., “Hard Negative Mixing for Contrastive Learning” .

$\uparrow \# \text{ negatives} \implies \text{harder negatives} \implies \text{performance?}$

Some work on the topic.²¹

²¹Kalantidis et al., “Hard Negative Mixing for Contrastive Learning” .

↑ # negatives \implies harder negatives \implies performance?

Some work on the topic.²¹

Increased false negatives?

- What if increasing negatives only improved performance by having a better chance of drawing hard negatives (negatives that are close to the instance)?
- Some work on mining hard negatives, that is, trying to select the hardest negatives to train.
- This increases the chance of false negatives when the encoder improves, in the self-supervised setting.

²¹Kalantidis et al., "Hard Negative Mining for Contrastive Learning" .

Are negative examples really necessary?

Do they **only** avoid collapse?

2023-10-23

Representation Contrastive Learning

└ Discussion on negative examples

└ Are negative examples really necessary?

Are negative examples really necessary?

Do they **only** avoid collapse?

Are negative examples really necessary?

Do they **only** avoid collapse? Can this be done in a different way?

2023-10-23

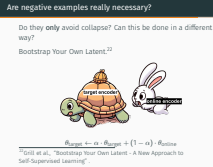
Representation Contrastive Learning

└ Discussion on negative examples

└ Are negative examples really necessary?

Are negative examples really necessary?

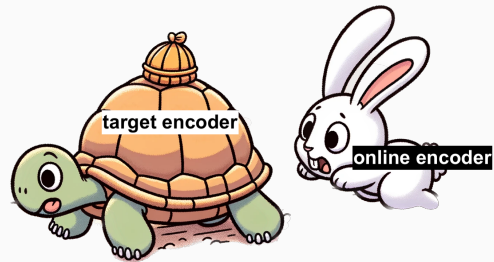
Do they **only** avoid collapse? Can this be done in a different way?



Are negative examples really necessary?

Do they **only** avoid collapse? Can this be done in a different way?

Bootstrap Your Own Latent.²²



$$\theta_{\text{target}} \leftarrow \alpha \cdot \theta_{\text{target}} + (1 - \alpha) \cdot \theta_{\text{online}}$$

²²Grill et al., "Bootstrap Your Own Latent - A New Approach to Self-Supervised Learning".

2023-10-23

Representation Contrastive Learning

└ Discussion on negative examples

└ Are negative examples really necessary?

- 2 neural networks, an online (predictive) network and a target network. They use only positive examples, and for those the online network tries to predict the metric representation from the target network.
- The parameters of the target network are updated after every iteration with an exponential moving average of the online parameters.
- The authors argue that since the update to the target parameters is not exactly according to the gradient of the loss with respect to θ_{target} , there is no a priori reason to believe that the target network would converge to a collapsed representation.
- Informal discussion on whether batch normalization plays a role.

2023-10-23

Representation Constrastive Learning
└ Applications

Applications

Applications

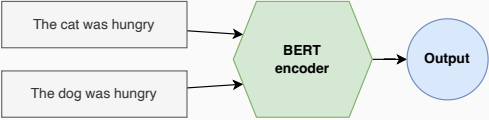


Figure 3: BERT diagram

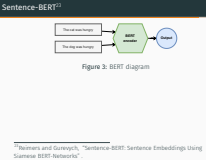
²³Reimers and Gurevych, "Sentence-BERT: Sentence Embeddings Using Siamese BERT-Networks" .

2023-10-23

Representation Constrastive Learning

└ Applications

└ Sentence-BERT



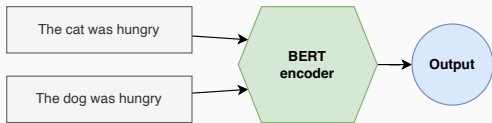


Figure 3: BERT diagram

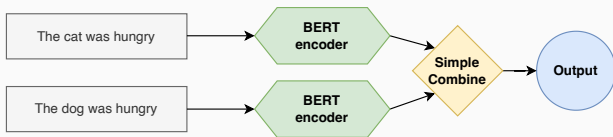


Figure 4: Sentence-BERT diagram

2023-10-23

Representation Constrastive Learning

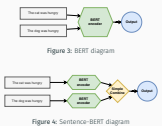
└ Applications

└ Sentence-BERT^a

^aReimers and Gurevych, "Sentence-BERT: Sentence Embeddings Using

Explain the difference in approach using the diagrams

Sentence-BERT



+ 10000 sentences, find most similar pair, **65 hour** → **5 seconds**.

2023-10-23

Representation Constrastive Learning

└ Applications

└ Sentence-BERT

Sentence-BERT

+ 10000 sentences, find most similar pair, **65 hour** → **5 seconds**.

- + 10000 sentences, find most similar pair, **65 hour** → **5 seconds**.
- + Allows fast k-NN approximations.

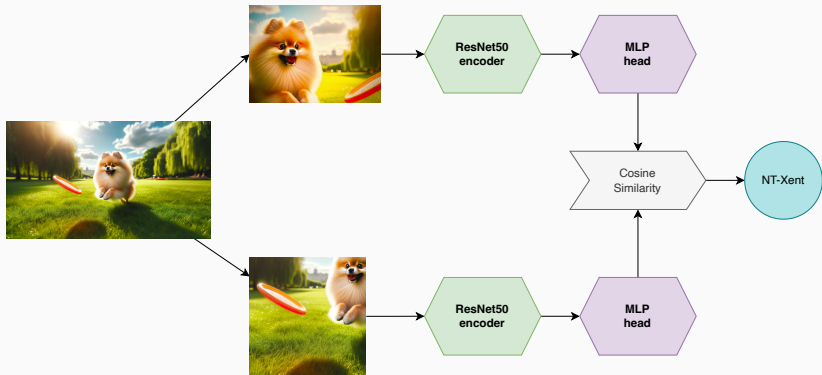
- 10000 sentences, find most similar pair, **65 hour** → **5 seconds**.
- Allows fast k-NN approximations.

- + 10000 sentences, find most similar pair, **65 hour** → **5 seconds**.
- + Allows fast k-NN approximations.
- + Great performance for STS.

- 10000 sentences, find most similar pair, **65 hour** → **5 seconds**.
- Allows fast k-NN approximations.
- Great performance for STS.

- + 10000 sentences, find most similar pair, **65 hour** → **5 seconds**.
- + Allows fast k-NN approximations.
- + Great performance for STS.
- Worse performance on more complex tasks.

- + 10000 sentences, find most similar pair, **65 hour** → **5 seconds**.
- + Allows fast k-NN approximations.
- + Great performance for STS.
- Worse performance on more complex tasks.



²⁴Chen et al., “A Simple Framework for Contrastive Learning of Visual Representations” .

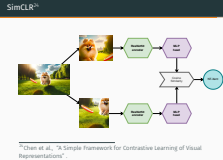
2023-10-23

Representation Contrastive Learning

└ Applications

└ SimCLR^a

^aChen et al. “A Simple Framework for Contrastive Learning of Visual



²⁴Chen et al., “A Simple Framework for Contrastive Learning of Visual Representations” .

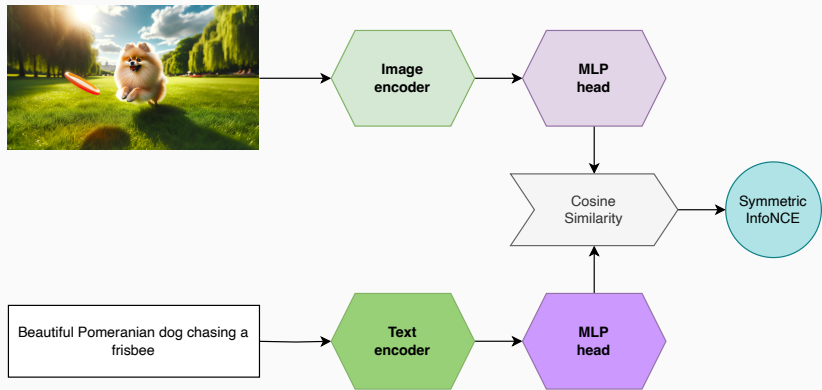
- Go over diagram.
- Other examples in batch assumed negatives.

- State-of-the-art performance in self-supervised, semi-supervised (1% and 10% of labels) and transfer learning for classification tasks.

- State-of-the-art performance in self-supervised, semi-supervised (1% and 10% of labels) and transfer learning for classification tasks.
- Competitive performance with fully-supervised models on ImageNet (with 4x the model size, but without labeled data).

- State-of-the-art performance in self-supervised, semi-supervised (1% and 10% of labels) and transfer learning for classification tasks.
- Competitive performance with fully-supervised models on ImageNet (with 4x the model size, but without labeled data).

Evaluations using a linear classifier on the learned representation



2023-10-23

Representation Contrastive Learning

└ Applications

└ CLIP

CLIP



- + Zero-shot classification (including geo-localization, OCR, facial emotion recognition, action recognition...)
- Well below SOTA performance.

2023-10-23

Conclusions

- Great potential for representation learning, specially self-supervised.
- Many challenges: resource intensity, zero-shot low performance, false negatives...

Conclusions

Representation Constrastive Learning

└─ Conclusions

└─ Conclusions

2023-10-23

- Great potential for representation learning, specially self-supervised.
- Many challenges: resource intensity, zero-shot low performance, false negatives...