Minimum Spanning Tree Regularization for Self-Supervised Learning

Julie Mordacq Joint Work with David Loiseaux, Vicky Kalogeiton, Steve Oudot

Journées de Géométrie Algorithmique





Table of content

- 1. Introduction
- 2. Background: Minimum Spanning Trees & Dimension estimation
- 3. T-REG: Minimum Spanning Tree based Regularization
- 4. T-REGS: T-REG for Self-Supervised Learning

Deep Learning

Given a training set $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$, with $x_i \in \mathbb{R}^d$ an observation and $y_i \in \mathcal{Y} = \{0, ..., J\}$, we want to learn:

Deep Learning

Given a training set $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$, with $x_i \in \mathbb{R}^d$ an observation and $y_i \in \mathcal{Y} = \{0, ..., J\}$, we want to learn:

Neural network:

A parametrized function (neural network) $f_{\theta}: \mathbb{R}^d \to \mathcal{Y}$, indexed by some parameters $\theta \in \Theta \subset \mathbb{R}^b$

Deep Learning

Given a training set $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$, with $x_i \in \mathbb{R}^d$ an observation and $y_i \in \mathcal{Y} = \{0, ..., J\}$, we want to learn:

Neural network:

A parametrized function (neural network) $f_{\theta}: \mathbb{R}^d \to \mathcal{Y}$, indexed by some parameters $\theta \in \Theta \subset \mathbb{R}^b$

Optimization:

We want to find the **optimal parameters** θ , which minimise the empirical risk:

$$\hat{\theta} \in \arg\min_{\theta \in \Theta} M_N(\theta) \quad \text{with} \quad M_N := \frac{1}{N} \mathcal{L}(y_i, f_{\theta}(x_i))$$

Where \mathcal{L} is a cost function such as a mean squared error or cross-entropy.

Deep Learning

Given a training set $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$, with $x_i \in \mathbb{R}^d$ an observation and $y_i \in \mathcal{Y} = \{0, ..., J\}$, we want to learn:

Neural network:

A parametrized function (neural network) $f_{\theta}: \mathbb{R}^d \to \mathcal{Y}$, indexed by some parameters $\theta \in \Theta \subset \mathbb{R}^b$

Optimization:

We want to find the **optimal parameters** θ , which minimise the empirical risk:

$$\hat{\theta} \in \arg\min_{\theta \in \Theta} M_N(\theta) \quad \text{with} \quad M_N := \frac{1}{N} \mathcal{L}(y_i, f_{\theta}(x_i))$$

Where \mathcal{L} is a cost function such as a mean squared error or cross-entropy.

Gradient-based methods: back-propagation

Deep Learning

Given a training set $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$, with $x_i \in \mathbb{R}^d$ an observation and $y_i \in \mathcal{Y} = \{0, ..., J\}$, we want to learn:

Neural network:

A parametrized function (neural network) $f_{\theta}: \mathbb{R}^d \to \mathcal{Y}$, indexed by some parameters $\theta \in \Theta \subset \mathbb{R}^b$

Optimization:

We want to find the **optimal parameters** θ , which minimise the empirical risk:

$$\hat{\theta} \in \operatorname{arg\ min}_{\theta \in \Theta} M_N(\theta) \quad \text{with} \quad M_N := \frac{1}{N} \mathcal{L}(y_i, f_{\theta}(x_i))$$

Where \mathcal{L} is a cost function such as a mean squared error or cross-entropy.

Universal approximation theorem

Neural networks with a given structure can, in principle, approximate any continuous function to any desired degree of accuracy

Deep Unsupervised Learning

Given a training set $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$, with $x_i \in \mathbb{R}^d$ an observation and $y_i \in \mathcal{Y} = \{0, ..., J\}$ we want to learn:

Neural network: a parametrized function $f_{\theta}: \mathbb{R}^d \to \mathcal{Y}$ indexed by some parameters $\theta \in \Theta \subset \mathbb{R}^b$.

Optimization:

We want to find the **optimal parameters** θ , which minimise the empirical risk:

$$\hat{\theta} \in \operatorname{arg\ min}_{\theta \in \Theta} M_N(\theta) \quad \text{with} \quad M_N := \frac{1}{N} \mathcal{L}(y_i, f_{\theta}(x_i))$$

Deep Unsupervised Learning

Given a training set $\mathcal{D} = \{x_i\}_{i=1}^N$, with $x_i \in \mathbb{R}^d$ an observation

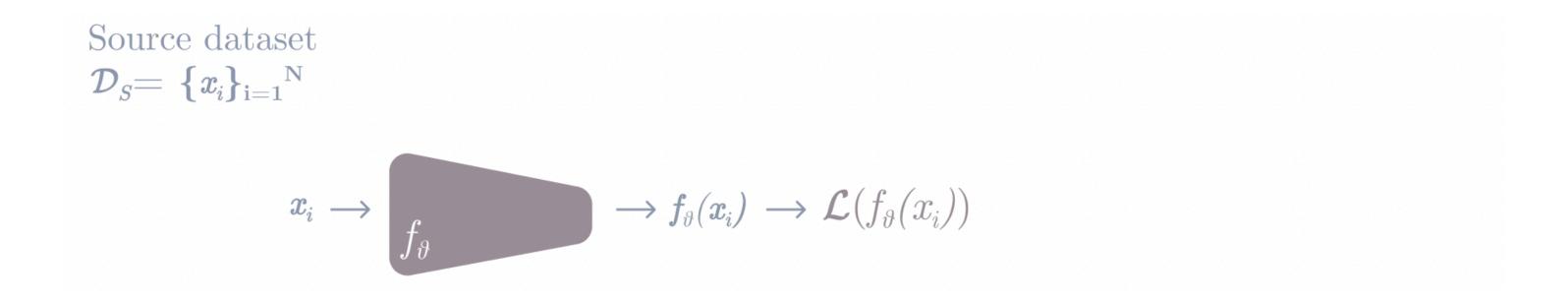
Neural network: a parametrized function $f_{\theta}: \mathbb{R}^d \to \mathbb{R}^p$ typically with d > p, indexed by some parameters $\theta \in \Theta \subset \mathbb{R}^b$.

Optimization:

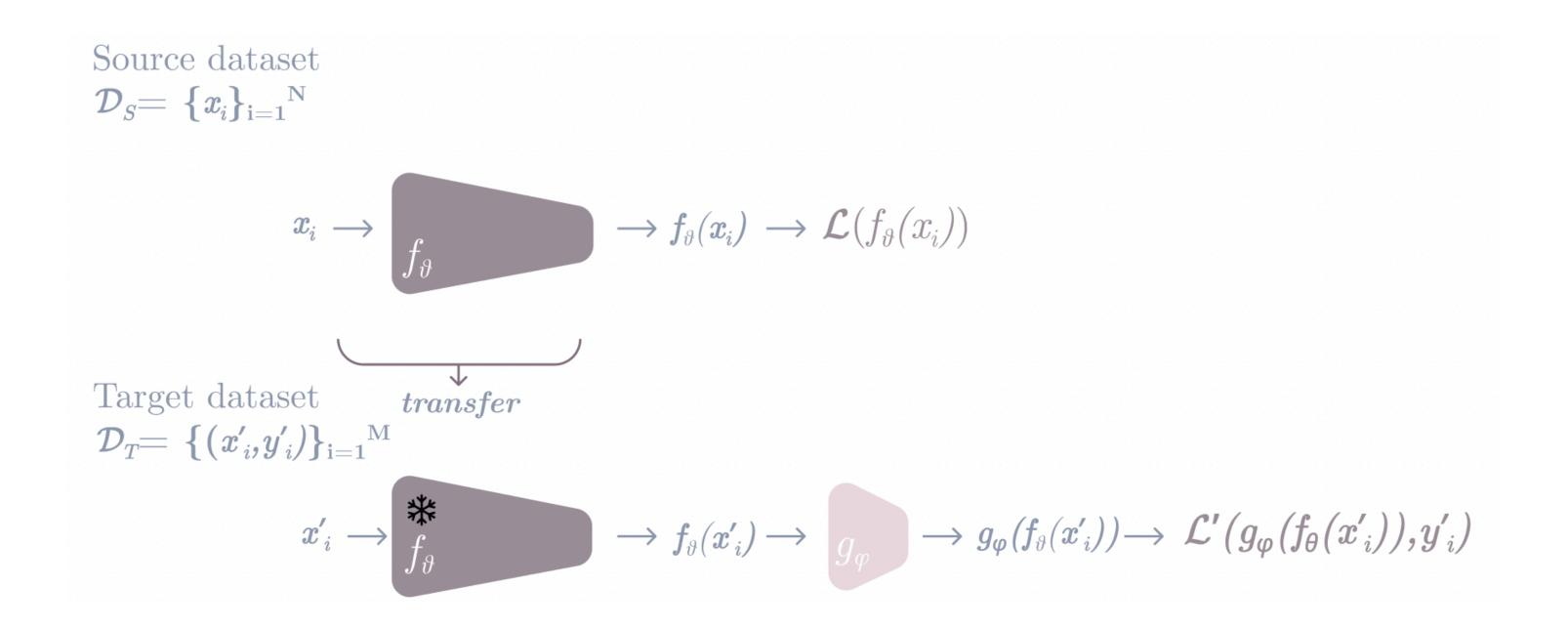
We want to find the **optimal parameters** θ , which minimise the empirical risk:

$$\hat{\theta} \in \arg\min_{\theta \in \Theta} M_N(\theta) \quad \text{with} \quad M_N := \frac{1}{N} \mathcal{L}(f_{\theta}(x_i))$$

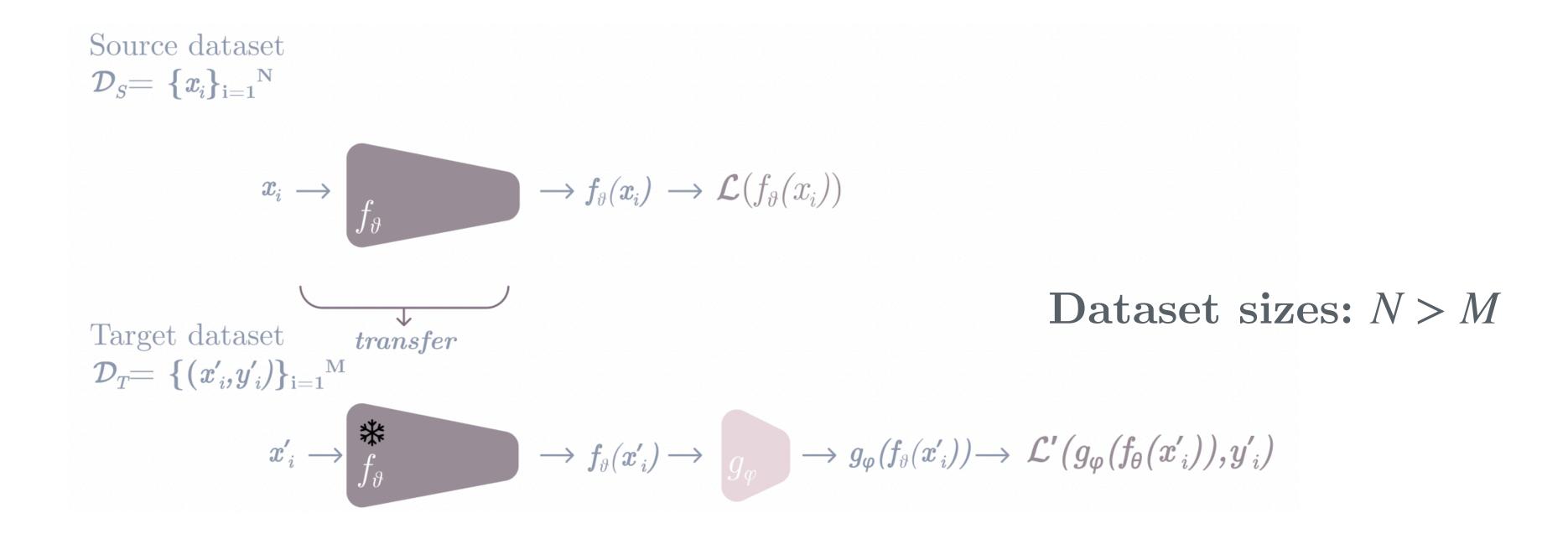
Self-supervised learning



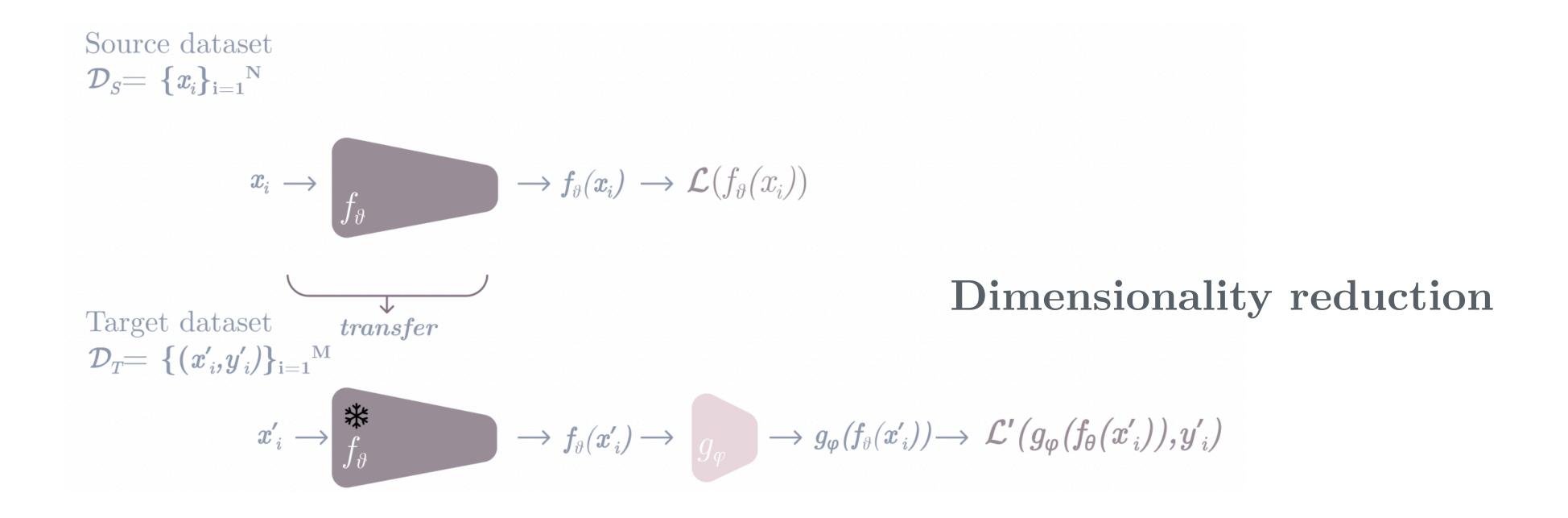
Self-supervised learning



Self-supervised learning



Self-supervised learning



Self-supervised learning

Transfer Learning: transferring knowledge (e.g., learned weights: f_{θ}) from a source domain to improve performance in a target domain and task

How to get 'good' and 'general' representations without knowing the downstream task(s)?

Self-supervised learning

Self-supervised learning: paradigm to learn meaningful data representations without relying on annotations, by leveraging supervision from the raw data itself

Self-supervised learning

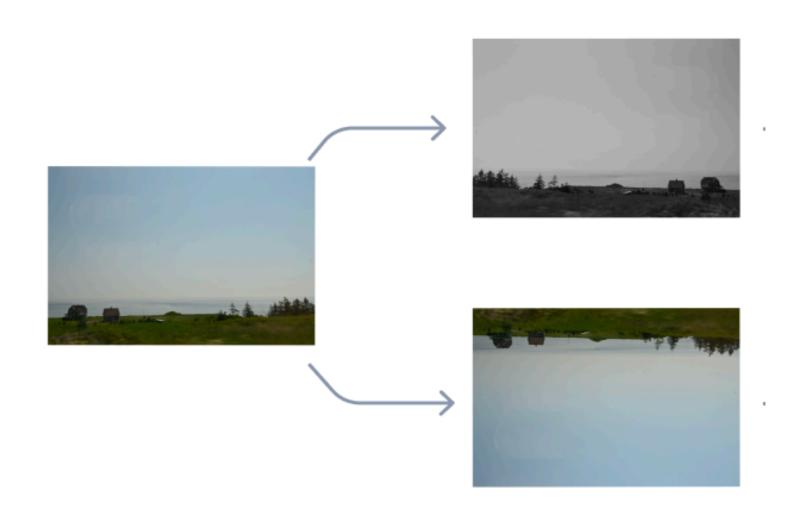
Self-supervised learning: paradigm to learn meaningful data representations without relying on annotations, by leveraging supervision from the raw data itself.

Invariance to a given class of transforms

Joint-Embedding Self-Supervised Learning

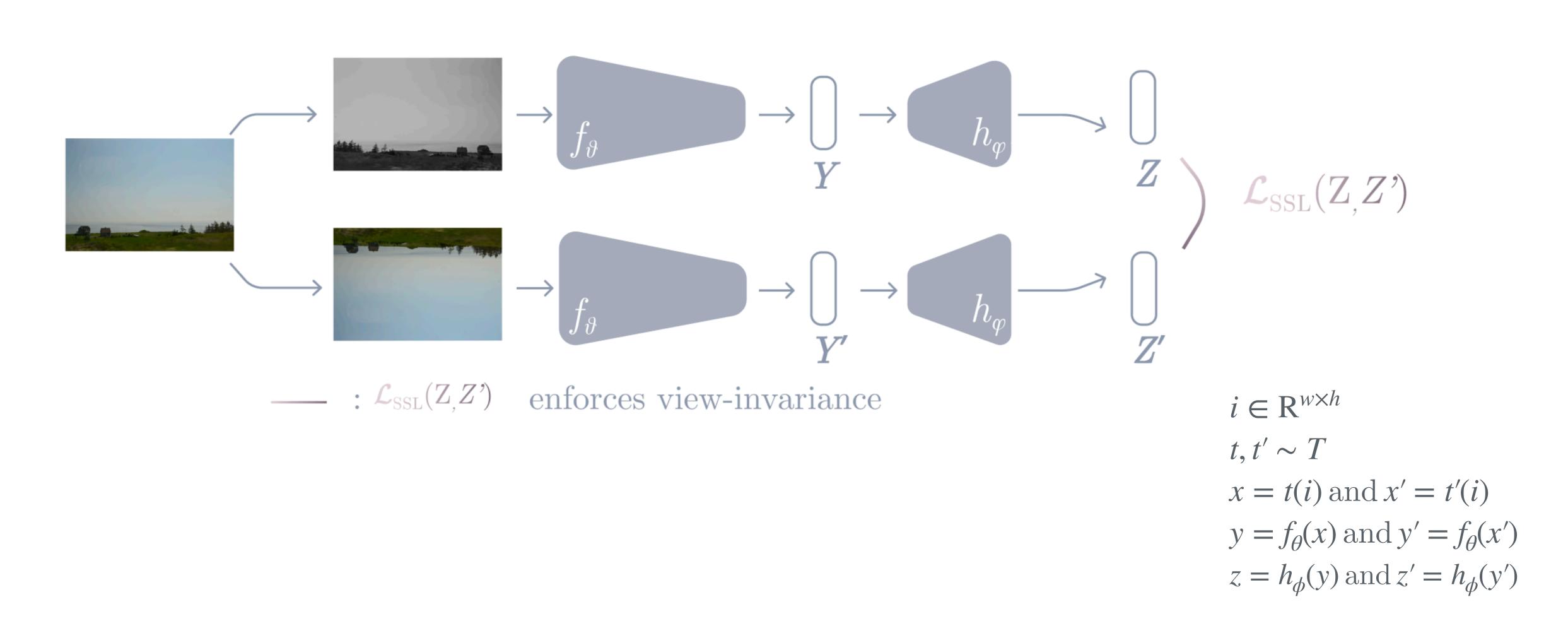


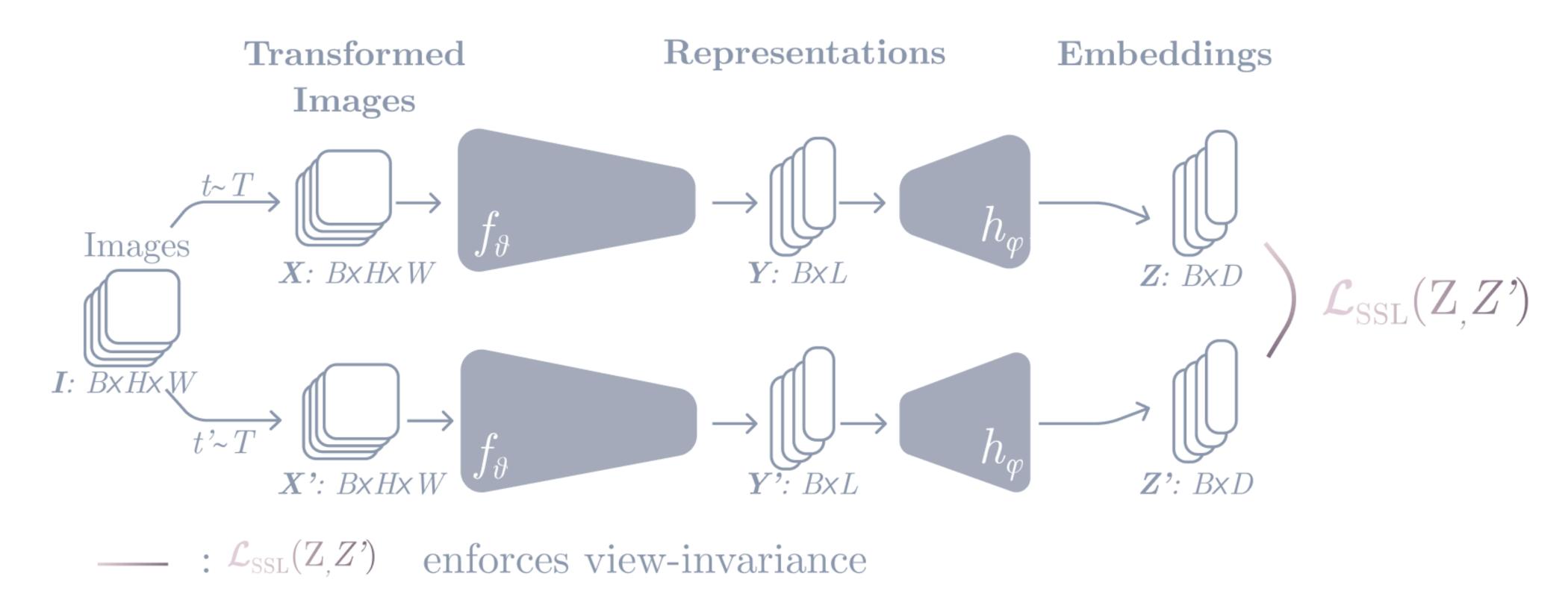
 $i \in \mathbb{R}^{w \times h}$



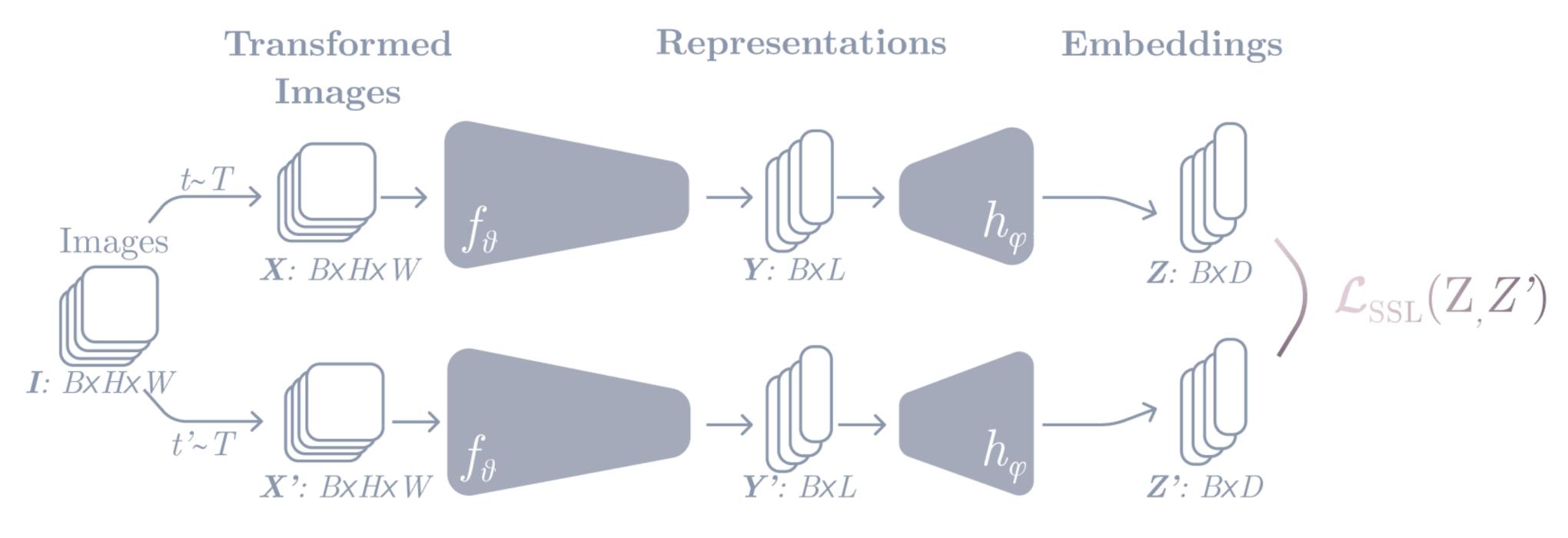
$$i \in \mathbb{R}^{w \times h}$$

 $t, t' \sim T$
 $x = t(i) \text{ and } x' = t'(i)$





Joint-Embedding Self-Supervised Learning



 \ldots : $\mathcal{L}_{\text{SSL}}(\mathbf{Z},\mathbf{Z}')$ enforces view-invariance

In practise:

- we may consider f_{θ} and h_{ϕ} as one network
- can be viewed as a dimensionality reduction problem as $D < H \times W$

Joint-Embedding Self-Supervised Learning

The fundamental challenge in JE-SSL: preventing **collapse** either representational or dimensional

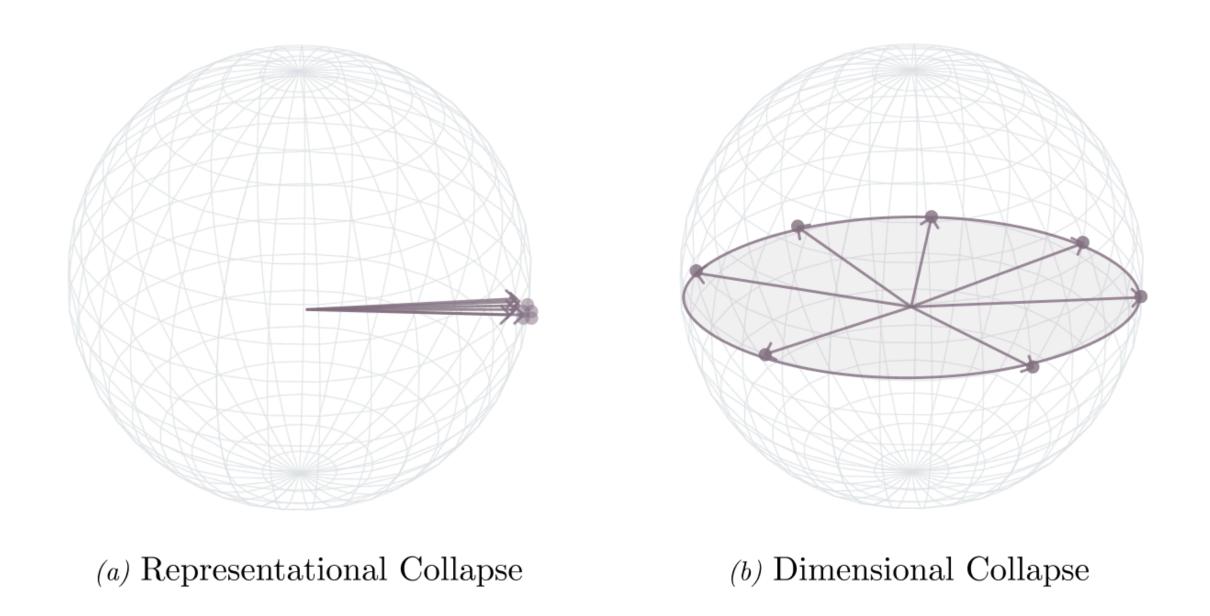


Figure: Illustration of collapse.

- i. Contrastive approaches, e.g, SimCLR [1] encourage embeddings of different views of the same image to be similar while pushing away embeddings of different images.
- ii. Asymmetric architecture, e.g., BYOL [2]
- iii. Covariance-based approaches, e.g., VICReg [3] enforce embedding decorrelation by promoting an identity covariance matrix

^[1] A Simple Framework for Contrastive Learning of Visual Representations, ICML 2020

^[2] Bootstrap your own latent: A new approach to self-supervised Learning, NeurIPS 2020

^[3] VICReg: Variance-Invariance-Covariance Regularization For Self-Supervised Learning, ICLR 2022

Joint-Embedding Self-Supervised Learning

Fang et al. [1] formalized desirable properties of JE-SSL:

- 1. Mitigating dimensional collapse: encourage the embeddings to span the entire space
- 2. Promoting sample uniformity: ensure the embeddings are evenly distributed across the representation space

25

^[1] Rethinking uniformity in self-supervised learning, Fang et al, ICLR 2024

Joint-Embedding Self-Supervised Learning

Fang et al. [1] formalized desirable properties of JE-SSL:

1. Mitigating dimensional collapse: encourage the embeddings to span the entire space

2. Promoting sample uniformity: ensure the embeddings are evenly distributed across the

representation space



^[1] Rethinking uniformity in self-supervised learning, Fang et al, ICLR 2024

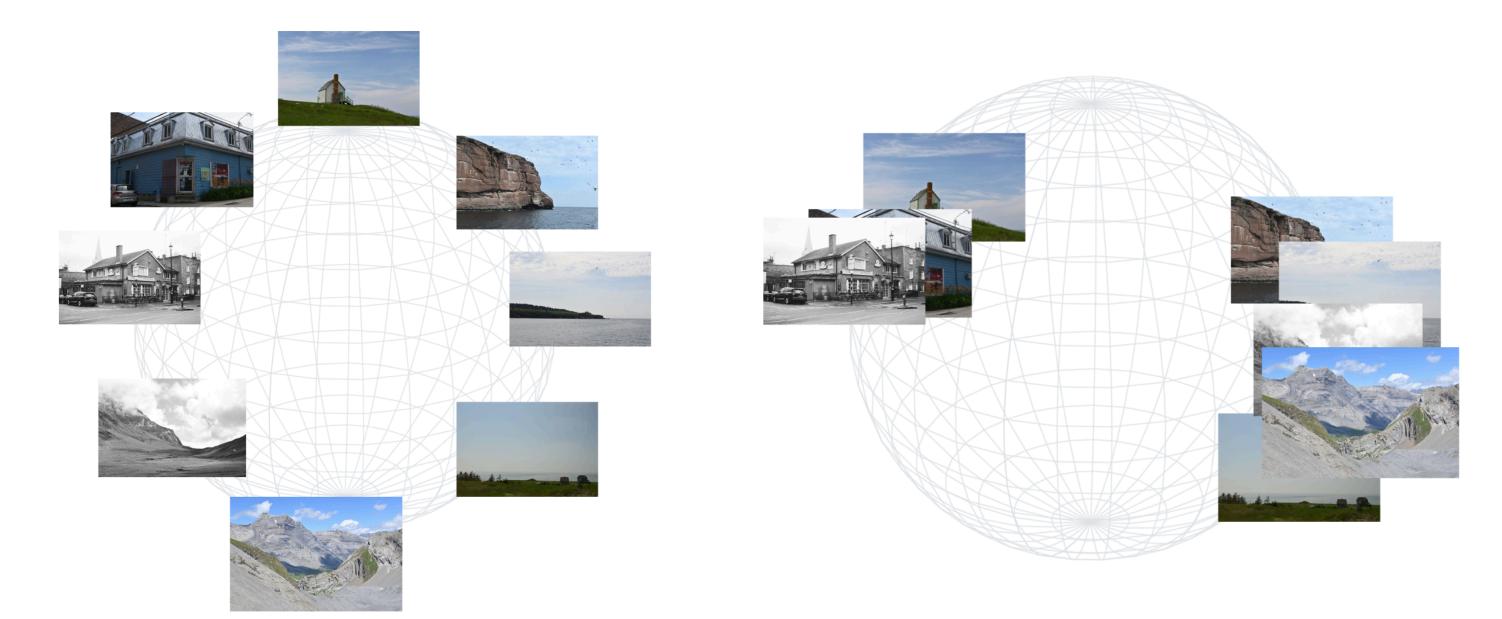
Joint-Embedding Self-Supervised Learning

Fang et al. [1] formalized desirable properties of JE-SSL:

1. Mitigating dimensional collapse: encourage the embeddings to span the entire space

2. Promoting sample uniformity: ensure the embeddings are evenly distributed across the

representation space



^[1] Rethinking uniformity in self-supervised learning, Fang et al, ICLR 2024

Joint-Embedding Self-Supervised Learning

Fang et al. [1] formalized desirable properties of JE-SSL:

- 1. Mitigating dimensional collapse: encourage the embeddings to span the entire space
- 2. Promoting sample uniformity: ensure the embeddings are evenly distributed across the representation space

How? Maximizing the length of the Minimum Spanning Tree

28

^[1] Rethinking uniformity in self-supervised learning, Fang et al, ICLR 2024

Minimum spanning tree

Definition:

Given a point cloud Z in a Euclidean space, a **spanning tree** (ST) of Z is an undirected graph G = (V, E) with the vertex set V = Z and edge set $E \subset V \times V$ such that G is connected without cycles.

Minimum spanning tree

Definition:

Given a point cloud Z in a Euclidean space, a **spanning tree** (ST) of Z is an undirected graph G = (V, E) with the vertex set V = Z and edge set $E \subset V \times V$ such that G is connected without cycles.

The length of G is:

$$E(G) := \sum_{(z,z')\in E} ||z-z'||_2$$

A MST(Z), is an ST of Z that minimises length E.

Minimum spanning tree

Definition:

Given a point cloud Z in a Euclidean space, a **spanning tree** (ST) of Z is an undirected graph G = (V, E) with the vertex set V = Z and edge set $E \subset V \times V$ such that G is connected without cycles.

The length of G is:

$$E(G) := \sum_{(z,z')\in E} ||z-z'||_2$$

A MST(Z), is an ST of Z that minimises length E.

The MST also relates to the persistence in degree 0 of the Rips filtration.

Minimum spanning tree

Definition:

Given a point cloud Z in a Euclidean space, a **spanning tree** (ST) of Z is an undirected graph G = (V, E) with the vertex set V = Z and edge set $E \subset V \times V$ such that G is connected without cycles.

The length of G is:

$$E(G) := \sum_{(z,z')\in E} ||z-z'||_2$$

A MST(Z), is an ST of Z that minimises length E.

The MST also relates to the persistence in degree 0 of the Rips filtration.

Specifically, there is a bijection between the edges of the minimal spanning tree of a finite metric space $x = \{x_1, ..., x_n\}$ and the points in the persistence diagram $PH_0(x)$ obtained from the Rips filtration.

Minimum spanning tree

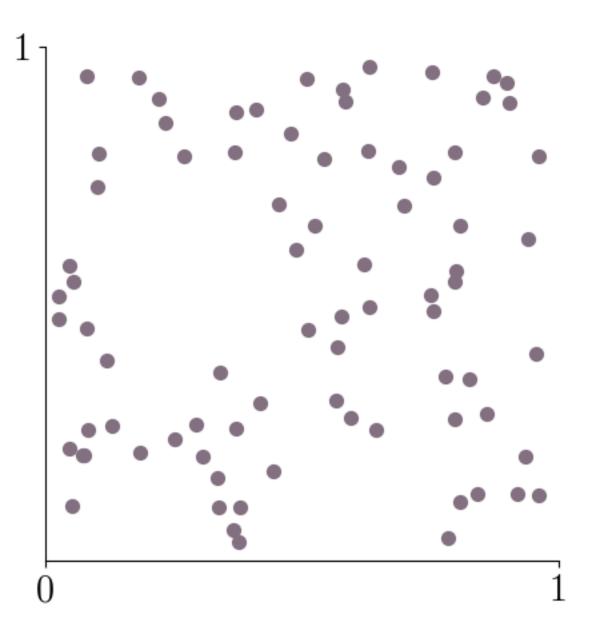


Figure: Example of MST in 2-d

Minimum spanning tree

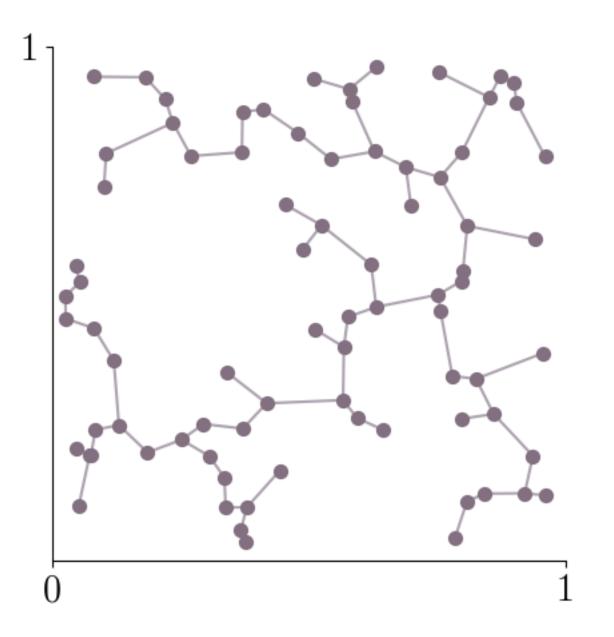


Figure: Example of MST in 2-d

2. Background

Background

Minimum spanning tree and dimension estimation

Steele [1] studied the length of the minimum spanning tree of random Euclidean spaces. Let E_n be an i.i.d n-sample drawn from a probability measure P_X with compact support on \mathbb{R}^d . For $d \geq 2$, Theorem 1 [1] controls the growth rate of the length of the MST(X_n) as follows:

$$E(MST(X_n) \sim Cn^{(d-1)/d} \text{ almost surely, as } n \to \infty$$

The asymptotic rate allows to derive several estimators of intrinsic dimension

^[1] Growth rates of euclidean minimal spanning trees with power weighted edges, Steele, 1988

Background

Minimum spanning tree and persistence

Persistence provides a mathematical framework to optimize the length of the MST [1]:

$$\forall x \in X, \quad \nabla_x E(MST(X)) = \sum_{(x,z) \text{ edge}} \nabla_x \parallel x - z \parallel_2 = \sum_{(x,z) \text{ edge}} \parallel x - z \parallel_2^{-1} (x - z).$$
of $MST(X)$ of $MST(X)$

Each pair of points forming an edge in the MST exerts a repulsive force on the other during optimization.

[1] Optimising persistence homology-based functions, Carrière et al, ICML 2021

3. T-REG: Minimum Spanning Tree based regularization

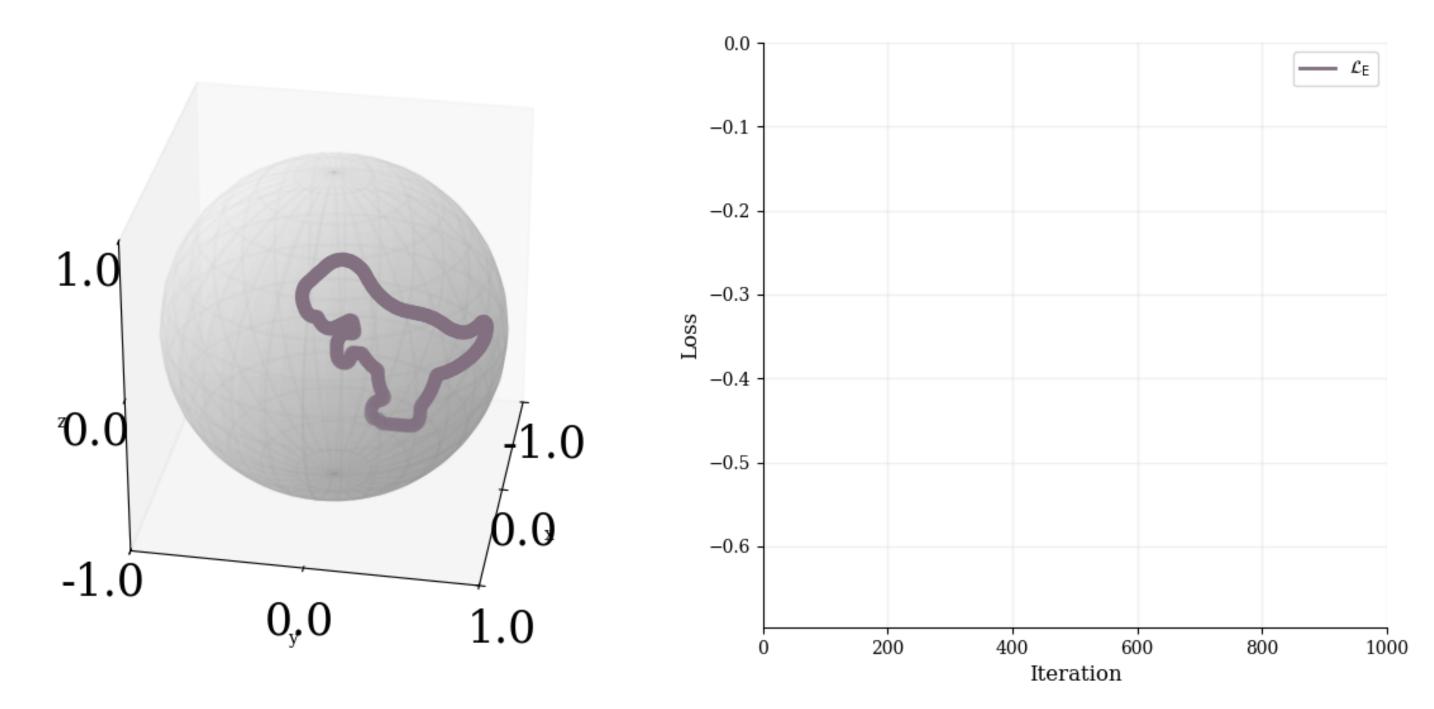
Given, $Z = \{z_0, ..., z_n\} \in \mathbb{R}^d$, the MST length is defined as:

$$\mathcal{L}_{\mathbf{E}} = -\frac{1}{n} E(\mathbf{MST}(Z))$$

Given, $Z = \{z_0, ..., z_n\} \in \mathbb{R}^d$, the MST length is defined as:

$$\mathcal{L}_{\mathbf{E}} = -\frac{1}{n} E(\mathbf{MST}(Z))$$

Point Cloud optimization with \mathcal{L}_{E}



Given, $Z = \{z_0, ..., z_n\} \in \mathbb{R}^d$, the MST length is defined as:

$$\mathcal{L}_{\mathbf{E}} = -\frac{1}{n} E(\mathbf{MST}(Z))$$

The soft-sphere constraint is given by:

$$\mathcal{L}_{\mathbf{S}} = \frac{1}{n} \sum_{i} (||z_{i}||_{2} - 1)^{2}$$

Given, $Z = \{z_0, ..., z_n\} \in \mathbb{R}^d$, the MST length is defined as:

$$\mathcal{L}_{\mathbf{E}} = -\frac{1}{n} E(\mathbf{MST}(Z))$$

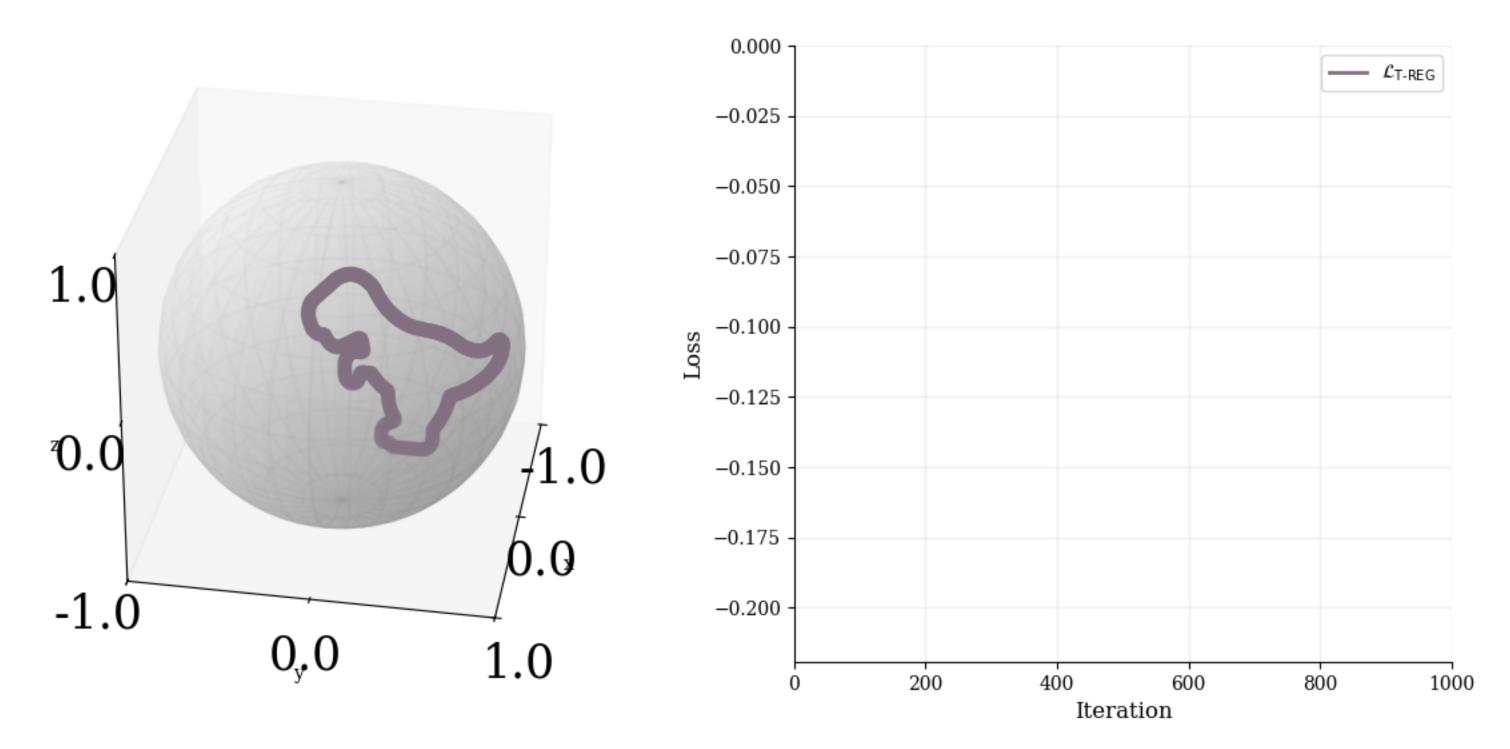
The soft-sphere constraint is given by:

$$\mathcal{Z}_{\mathbf{S}} = \frac{1}{n} \sum_{i} (||z_i||_2 - 1)^2$$

The overall loss is defined as: $\mathcal{L}_{T-REG} = \gamma \mathcal{L}_E + \lambda \mathcal{L}_S$ (with γ , λ are hyper-parameters)

The overall loss is defined as: $\mathcal{L}_{\text{T-REG}} = \gamma \mathcal{L}_{\text{E}} + \lambda \mathcal{L}_{\text{S}}$

Point Cloud optimization with $\mathcal{L}_{T\text{-REG}}$



3. T-REG: Minimum Spanning Tree based regularization

Theoretical analysis

Theoretical analysis

Asymptotic behavior on large samples (n > d + 1):

Derived from Steele [1], let E_n be an i.i.d n-sample drawn from a probability measure P_X with compact support on \mathbb{R}^d .

For $d \ge 2$, Theorem 1 [1] controls the growth rate of the length of the MST (X_n) as follows:

$$E(MST(X_n) \sim Cn^{(d-1)/d} \text{ almost surely, as } n \to \infty$$

^[1] Growth rates of euclidean minimal spanning trees with power weighted edges, Steele, 1988

Theoretical analysis

Asymptotic behavior on large samples (n > d + 1):

We fix a compact Riemannian d-manifold, \mathcal{M} , equipped with the d-dimensional Hausdorff measure μ .

Theorem 4.4 [1]: Let X_n be an iid n-sample of a probability measure on \mathcal{M} with density f_X w.r.t. μ . Then, there exists a constant independent of f_X and of \mathcal{M} such that:

$$n^{(d-1)/d} \cdot E(MST(X_n) \xrightarrow[n \to \infty]{} C' \int f^{\frac{d-1}{d}} d\mu$$

- [1] Determining Intrinsic Dimension and Entropy of High-Dimensional Shape Spaces. Costa et al, 2006
- [2] Shannon entropy, Renyi entropy, and information. Bromiley et al. 2004
- [3] Maximum entropy autoregressive conditional heteroskedasticity model. Park et al, 2009

Theoretical analysis

Asymptotic behavior on large samples (n > d + 1):

We fix a compact Riemannian d-manifold, \mathcal{M} , equipped with the d-dimensional Hausdorff measure μ .

Theorem 4.4 [1]: Let X_n be an iid n-sample of a probability measure on \mathcal{M} with density f_X w.r.t. μ . Then, there exists a constant independent of f_X and of \mathcal{M} such that:

$$n^{(d-1)/d} \cdot E(MST(X_n) \xrightarrow[n \to \infty]{} C' \int f^{\frac{d-1}{d}} d\mu$$

The limit is related to the intrinsic Rényi $\frac{d}{d-1}$ -entropy which is known to converge to the Shannon entropy as $\frac{d-1}{d} \to 1[2]$. The Shannon entropy, in turn, achieves its **maximum at the uniform** distribution on compact sets [3].

- [1] Determining Intrinsic Dimension and Entropy of High-Dimensional Shape Spaces. Costa et Hero, 2006
- [2] Shannon entropy, Renyi entropy, and information. Bromiley et al. 2004
- [3] Maximum entropy autoregressive conditional heteroskedasticity model. Park et al, 2009

Theoretical analysis

Behaviour on small samples ($n \le d + 1$, e.g. batch sizes are often smaller than or comparable to the ambient dimension):

Theoretical analysis

Behaviour on small samples ($n \le d + 1$, e.g. batch sizes are often smaller than or comparable to the ambient dimension): To account for the effect of the soft sphere constraint, we assume the points of X lie inside some fixed closed Euclidean d-ball B of radius r centered at the origin

Theorem 4.1:

Under the above conditions, the maximum of the E(MST(Z)), over the point sets $X \subset B$ of fixed cardinality n is attained when the points of X lie on the sphere $S = \partial B$, at the vertices of a regular (n-1)-simplex that has S as its smallest circumscribing sphere.

Theoretical analysis

Behaviour on small samples ($n \le d + 1$, e.g. batch sizes are often smaller than or comparable to the ambient dimension): In order to account for the effect of the soft sphere constraint, we assume the points of X lie inside some fixed closed Euclidean d-ball B of radius r centered at the origin

Theorem 4.1:

Under the above conditions, the maximum of the E(MST(Z)), over the point sets $X \subset B$ of fixed cardinality n is attained when the points of X lie on the sphere $S = \partial B$, at the vertices of a regular (n-1)-simplex that has S as its smallest circumscribing sphere.

Behavior of T-REG:

- First, \mathscr{L}_E in $\mathscr{L}_{T\text{-REG}}$ expands the point cloud until the sphere constraint term \mathscr{L}_S becomes the dominating term
- Then, the points stop expanding and start spreading themselves out uniformly along the sphere of directions.

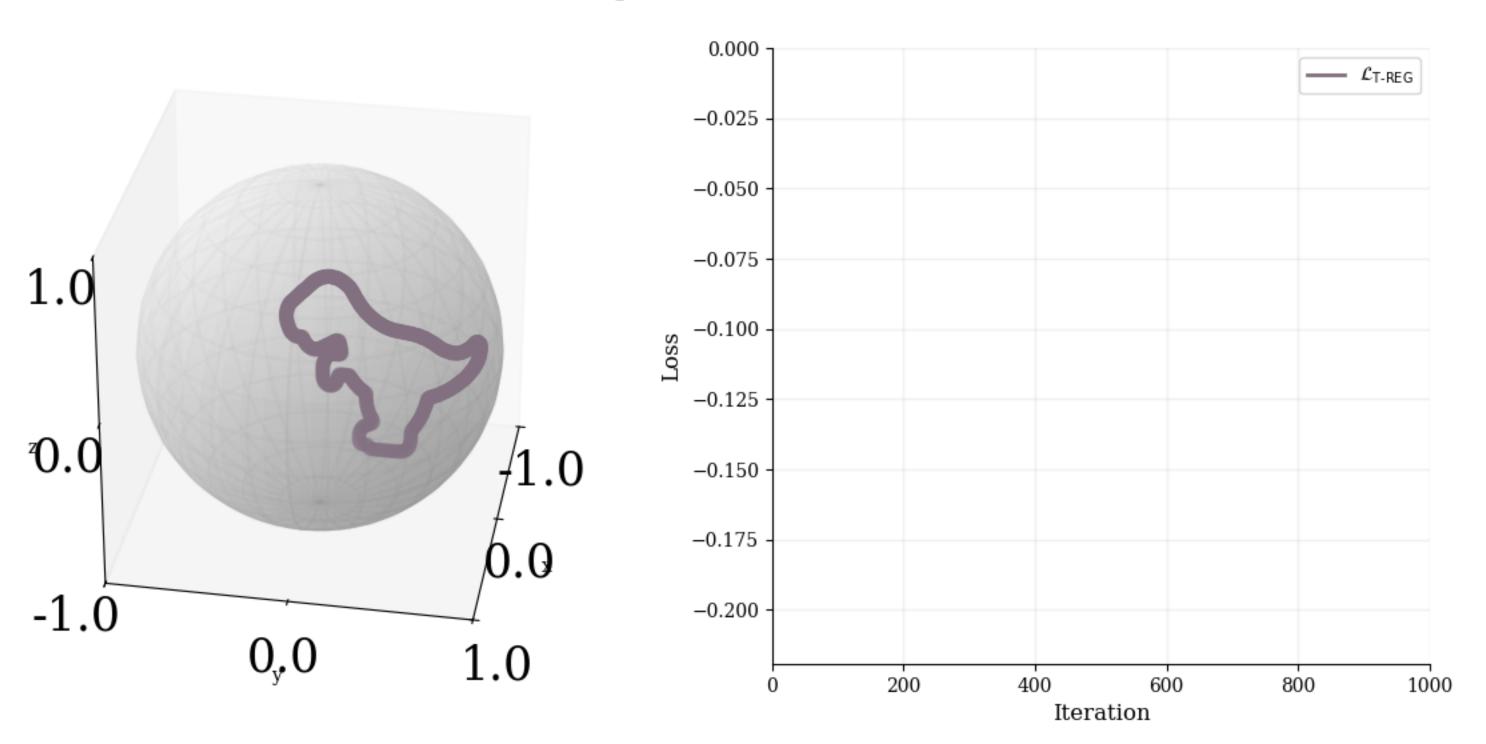
3. T-REG: Minimum Spanning Tree based regularization

Empirical analysis

Empirical analysis, promoting sample uniformity

Set-up: We apply T-REG alone to optimize the positions of a given point cloud, and we analyze its behavior when n > d + 1.

Point Cloud optimization with \mathcal{L}_{T-REG}

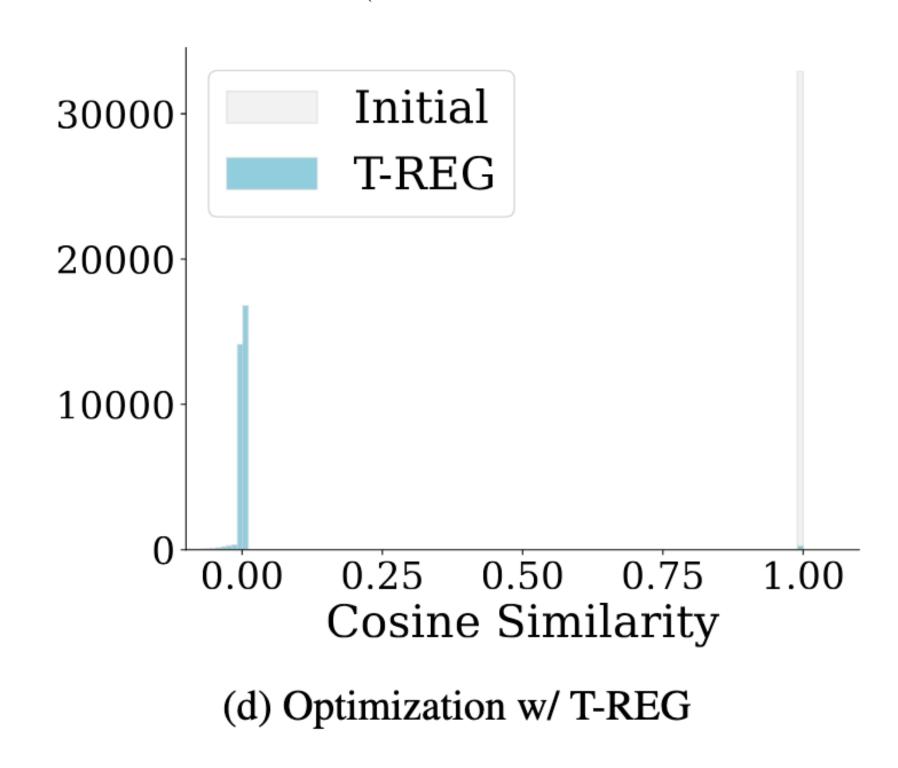


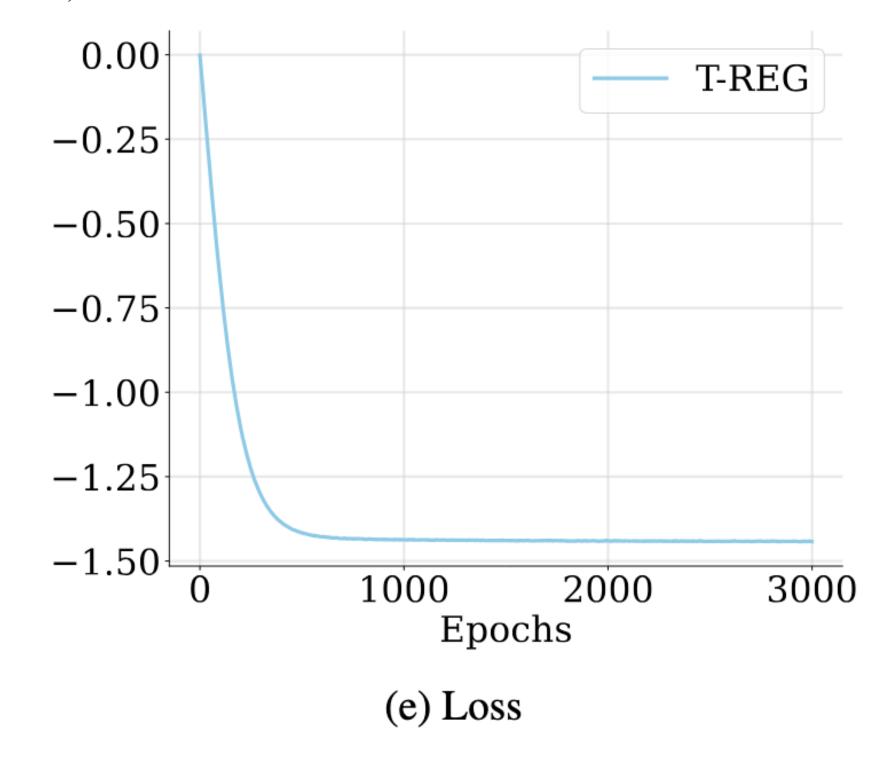
Empirical analysis, promoting sample uniformity

Set-up: We apply T-REG alone to optimize the positions of a given point cloud, and we analyze its behavior when $n \le d+1$ (specifically n=1024, d=1024).

Empirical analysis, promoting sample uniformity

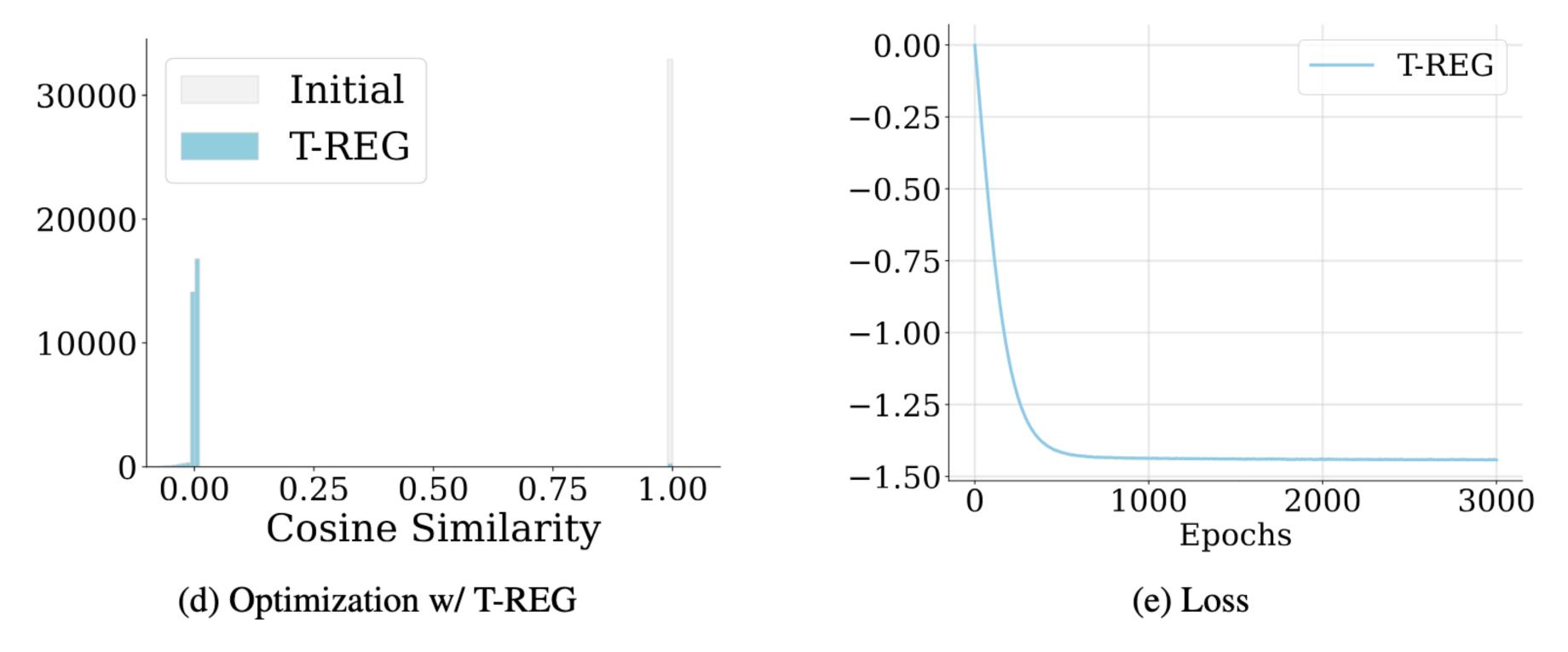
Set-up: We apply T-REG alone to optimize the positions of a given point cloud, and we analyze its behavior when $n \le d+1$ (specifically n=1024, d=1024).





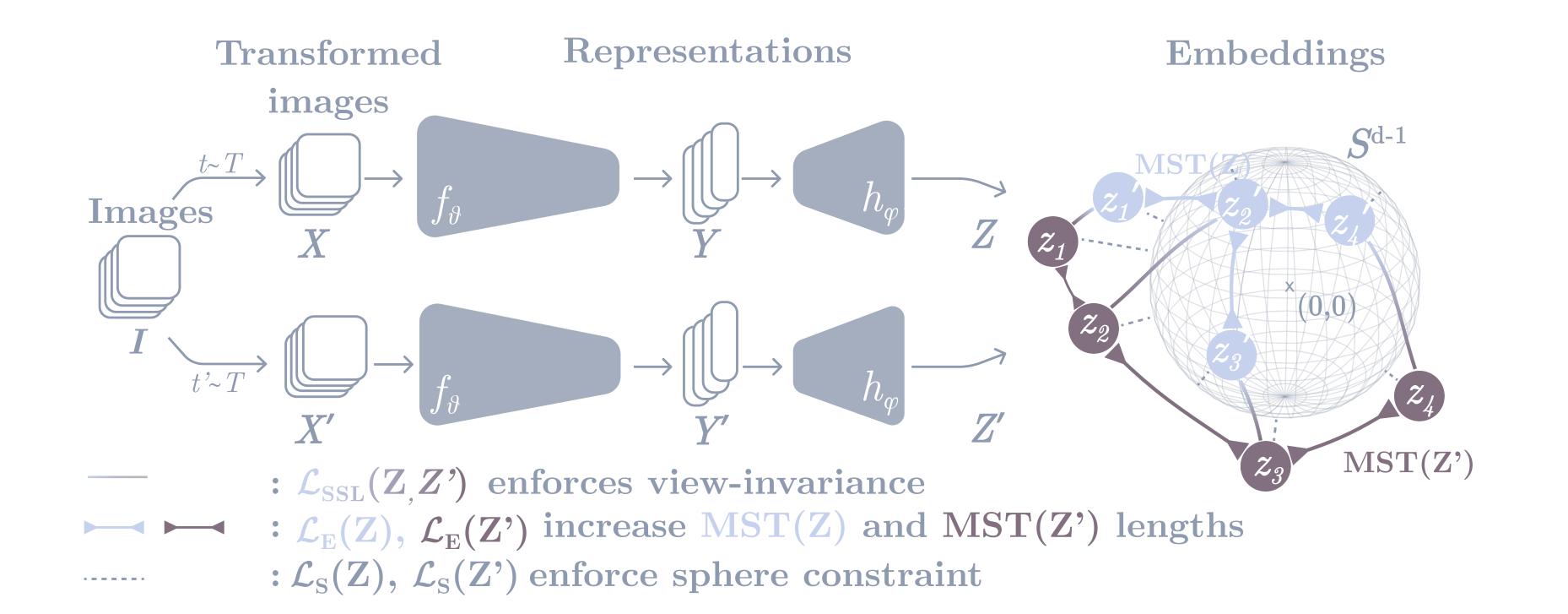
Empirical analysis, promoting sample uniformity

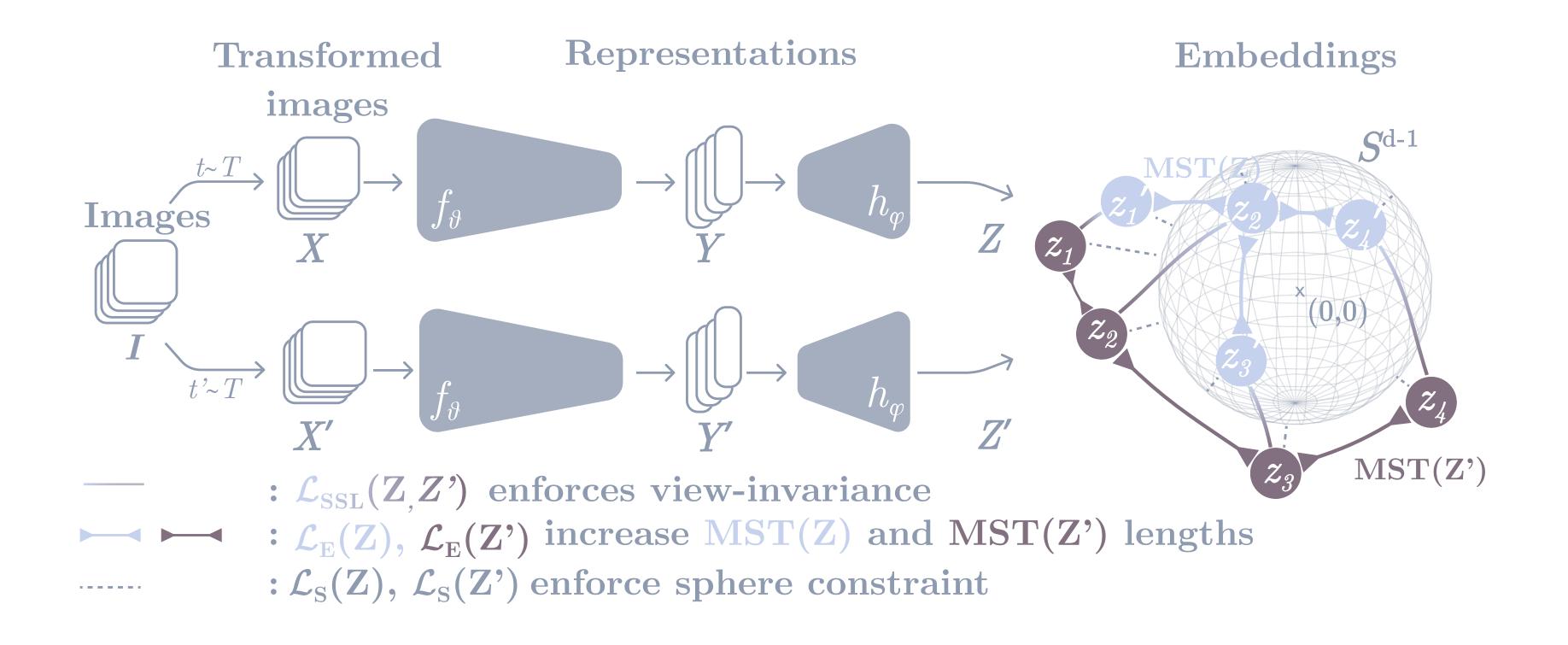
Set-up: We apply T-REG alone to optimize the positions of a given point cloud, and we analyze its behavior when $n \le d+1$ (specifically n=1024, d=1024).



After optimization with T-REG, the distribution becomes almost a Dirac slightly below 0, which indicates that the configuration of the points is close to that of the vertices of the regular simplex.

4. T-REGS: T-REG for Self-Supervised Learning





 $\mathcal{L}(Z,Z') = \beta \mathcal{L}_{\mathrm{SSL}}(Z,Z') + \gamma \mathcal{L}_{\mathrm{E}}(Z) + \lambda \mathcal{L}_{\mathrm{S}}(Z) + \gamma \mathcal{L}_{\mathrm{E}}(Z') + \lambda \mathcal{L}_{\mathrm{S}}(Z')$

T-REGS: T-REG for Self-supervised learning Evaluation on Standard SSL benchmark

Evaluation on standard SSL Benchmark

Method		CIFAR-10 [37]	CIFAR-100 [37]	
		91.3	68.5	
Zero-CL [67]	+ \mathcal{L}_u	91.3	68.4	
	+ \mathcal{W}_2	91.4	68.5	
		90.7	60.3	
MoCo v2 [14]	+ \mathcal{L}_u	91.0	61.2	
MOCO V2 [14]	+ \mathcal{W}_2	91.4	63.7	
		89.5	63.7	
BYOL [28]	$^+\mathcal{L}_u^{} \ ^+\mathcal{W}_2^{}$	90.1	62.7	
D I OL [26]	+ \mathcal{W}_2	90.1	65.2	
	+ T-REGS	90.4	65.7	
		91.2	68.2	
Barlow Twins [64]	+ \mathcal{L}_u	91.4	68.4	
Dariow Twills [04]	$^+\mathcal{L}_u \ ^+\mathcal{W}_2$	91.4	68.5	
	+ T-REGS	91.8	68.5	
\mathcal{L}_{MSE}	+ T-REGS	91.3	67.4	

Table 1: Comparison with W_2 regularization [23] on CIFAR-10/100. The table is inherited from Fang et al. [23], and we follow the same protocol: ResNet-18 models are pre-trained for 500 epochs on CIFAR-10/100, with a batch size of 256, followed by linear probing. We report Top-1 accuracy (%), **bold** indicates best performance.

On CIFAR-10/100.

T-REGS demonstrate strong standalone performances.

Evaluation on standard SSL Benchmark

Method		CIFAR-10 [37]	CIFAR-100 [37]	
		91.3	68.5	
Zero-CL [67]	+ \mathcal{L}_u	91.3	68.4	
	+ \mathcal{W}_2	91.4	68.5	
		90.7	60.3	
MoCo v2 [14]	+ \mathcal{L}_u	91.0	61.2	
WIOCO V2 [14]	+ \mathcal{W}_2	91.4	63.7	
		89.5	63.7	
BYOL [28]	$^+\mathcal{L}_u^{} \ ^+\mathcal{W}_2^{}$	90.1	62.7	
D I OL [26]	+ \mathcal{W}_2	90.1	65.2	
	+ T-REGS	90.4	65.7	
		91.2	68.2	
Barlow Twins [64]	+ \mathcal{L}_u	91.4	68.4	
Dariow Twills [04]	+ \mathcal{W}_2	91.4	68.5	
	+ T-REGS	91.8	68.5	
\mathcal{L}_{MSE}	+ T-REGS	91.3	67.4	

Table 1: Comparison with W_2 regularization [23] on CIFAR-10/100. The table is inherited from Fang et al. [23], and we follow the same protocol: ResNet-18 models are pre-trained for 500 epochs on CIFAR-10/100, with a batch size of 256, followed by linear probing. We report Top-1 accuracy (%), bold indicates best performance.

On CIFAR-10/100.

Using T-REGS as an auxiliary loss consistently improves performance over the respective baselines, and over variants that use \mathcal{L}_{U} or \mathcal{W}_{2} as additional regularization terms.

Evaluation on standard SSL Benchmark

#			Imagenet-100 [58]	ImageNet-1k [18]	
views	Method	_	Top-1	Batch Size	Top-1
	SwAV [9]		74.3	4096	66.5
8	FroSSL [55]		79.8	-	-
	SSOLE [34]		82.5	256	73.9
	SimCLR [12]		77.0	4096	66.5
	MoCo v2 [14]		79.3	256	67.4
	SimSiam [13]		78.7	256	68.1
	W-MSE [20]		69.1	512	65.1
	Zero-CL [67]		79.3	1024	68.9
	VICReg [4]		79.4	1024	68.3
2 -	CW-RGP [61]		77.0	512	67.1
	INTL [62]		81.7	512	69.5
	BYOL [28]		80.3	1024	66.5
		+ T-REGS	80.8	1024	67.2
	Darlass Trying [64]		80.2	2048	67.7
	Barlow Twins [64]	+ T-REGS	80.9	2048	67.8
	\mathcal{L}_{MSE}	+ T-REGS	80.3	512	68.8

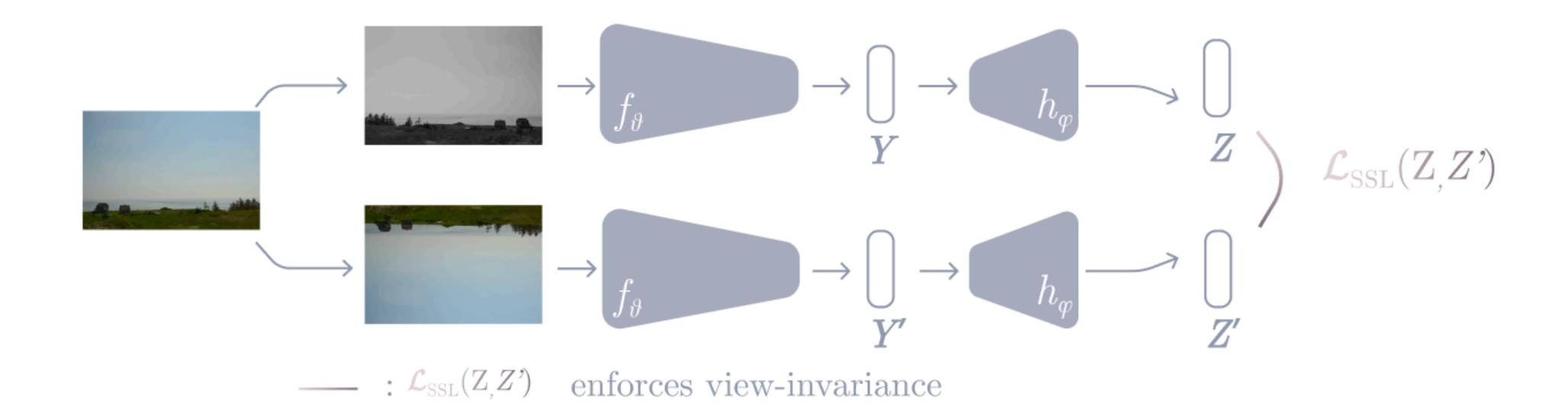
Table 2: Linear Evaluation on ImageNet-100/1k. We report Top-1 accuracy (%). Top-4 best methods are bolded. For ImageNet-100, ResNet-18 are pre-trained for 400 epochs using a batch size of 256; for ImageNet-1k, ResNet-50 are pre-trained for 100 epochs. The table is mostly inherited from Weng et al. [62].

On ImageNet-100/1k.

T-REGS is competitive with method that use the same number of views.

T-REGS: T-REG for Self-supervised learning
Further applications

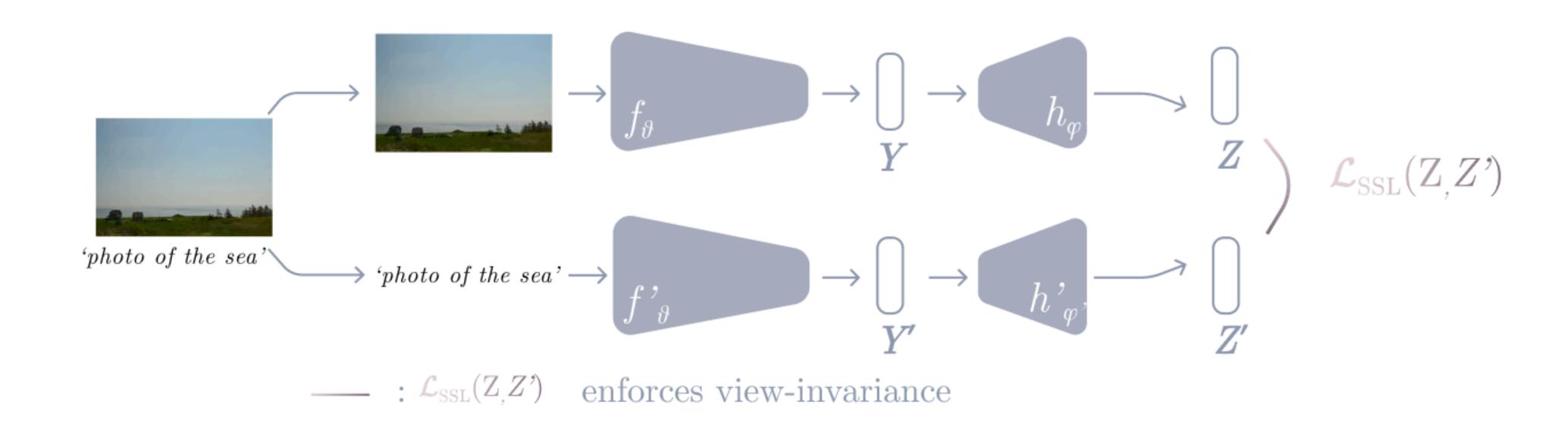
Further applications



Footer

Further applications

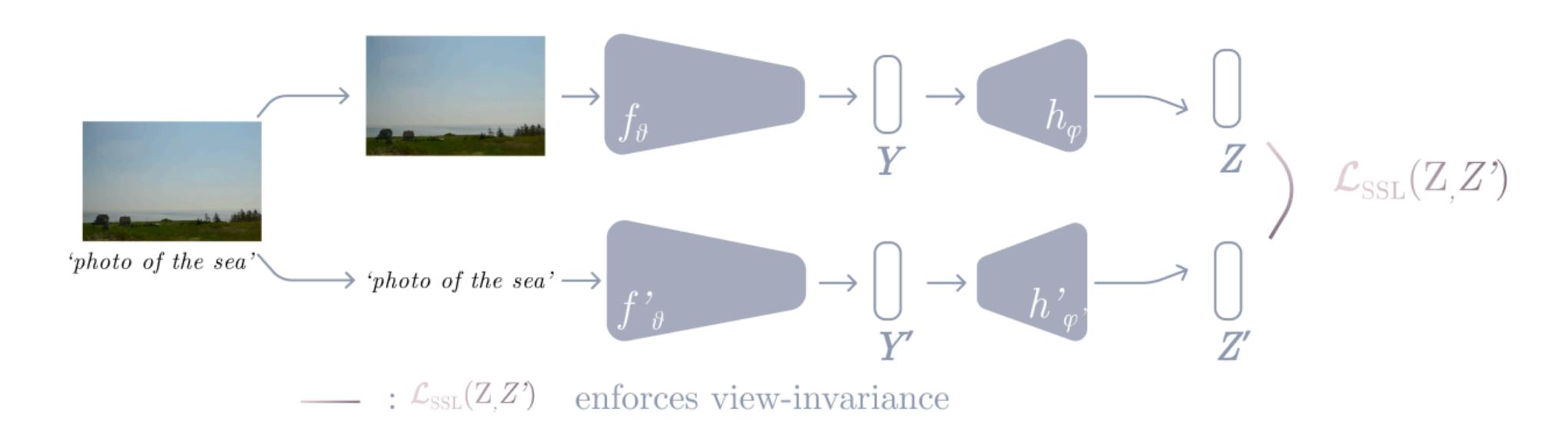
Joint-Embedding multimodal models [1]:



[1] Learning Transferable Visual Models From Natural Language Supervision, Radford et al, ICML 2021

Further applications

Joint-Embedding multimodal models [1]:



Such models preserves distinct subspaces for text and image: the modality gap [2].

- [1] Learning Transferable Visual Models From Natural Language Supervision, Radford et al, ICML 2021
- [2] Mind the Gap: Understanding the Modality Gap in Multi-modal Contrastive Representation Learning, Liang et al, NeurIPS 2022

Further applications

We fine-tune CLIP using T-REGS as regularization to the standard CLIP loss to encourage more robust and uniformly distributed representations.

	Flickr30k [48]				MS-COCO [42]			
	$i \rightarrow t$		$t \rightarrow i$		$t \rightarrow i$		$t \rightarrow i$	
Method	R@1	R@5	R@1	R@5	R@1	R@5	R@1	R@5
Zero-Shot	71.1	90.4	68.5	88.9	31.9	56.9	28.5	53.1
Finetune	81.2	95.5	80.7	95.8	36.7	63.6	36.9	63.9
ES [40]	71.8	90.0	68.5	88.9	31.9	56.9	28.7	53.0
i-Mix [39]	72.3	91.7	69.0	91.1	34.0	63.0	34.6	62.2
Un-Mix [53]	78.5	95.4	74.1	91.8	38.8	66.2	33.4	61.0
m^3 -Mix [44]	82.3	95.9	82.7	96.0	41.0	68.3	39.9	67.9
$\mathcal{L}_{\text{CLIP}}$ + T-REGS	83.2	96.0	80.8	96.4	41.6	68.7	41.5	68.7

Table 3: Cross-Modal Retrieval after finetuning CLIP. Image-to-text ($i \rightarrow t$) and text-to-image ($t \rightarrow i$) retrieval results (top 1/5 Recall: R@1, R@5). The table is mostly inherited from Oh et al. [44]. Bold indicates the best performance.

T-REGS: T-REG for Self-supervised learning Analysis

Complexity

The MSTs are computed with Kruskal's algorithm, whose worst-case time is $\mathcal{O}(B^2D)$, with B the batch size and D the embedding dimension.

Empirically T-REGS matches the per-step wall-clock of VICReg and simCLR.

Complexity

The MSTs are computed with Kruskal's algorithm, whose worst-case time is $\mathcal{O}(B^2D)$, with B the batch size and D the embedding dimension.

Empirically T-REGS matches the per-step wall-clock of VICReg and simCLR.

Method	Complexity	B range	D range	Wall-clock time
SimCLR [12] VICReg [4]	$\mathcal{O}(B^2 \cdot D) \ \mathcal{O}(B \cdot D^2)$	[2048-4096] [1024-4096]	[256-1024] [4096-8192]	$0.22 \pm 0.03 \\ 0.23 \pm 0.02$
\mathcal{L}_{MSE} + T-REGS	$\mathcal{O}(B^2(D \cdot \mathrm{log}B))$	[512-1024]	[512-2048]	0.20 ± 0.001

Table 5: Complexity and computational cost. Comparison between different methods is performed, with training on ImageNet-1k distributed across 4 Tesla H100 GPUs. The wall-clock time (sec/step) is averaged over 500 steps. B, D ranges are reported from Bardes et al. [4], Garrido et al. [25].

Analysis

$$\mathcal{L}(Z,Z') = \beta \mathcal{L}_{\mathrm{SSL}}(Z,Z') + \gamma \mathcal{L}_{\mathrm{E}}(Z) + \lambda \mathcal{L}_{\mathrm{S}}(Z) + \gamma \mathcal{L}_{\mathrm{E}}(Z') + \lambda \mathcal{L}_{\mathrm{S}}(Z')$$

T-REGS coefficients		Scal	ling		
β	γ	λ	$rac{eta}{\gamma}$	$\frac{\gamma}{\lambda}$	Top-1
1	-	-	-	-	collapse
1	1	-	1	-	collapse
10	1	1	10	1	collapse
10	0.5	5e-2	20	10	25.7
10	0.2	2e-2	50	10	45.4
10	0.5	2.5e-3	20	200	65.0
10	0.2	1e-3	50	200	65.3
10	0.5	2e-3	20	250	64.9
10	0.2	8e-4	50	250	66.1
10	0.02	8e-5	100	300	63.3

Table 4: Impact of coefficients. \mathcal{L}_{MSE} + T-REGS top-1 accuracy (%) on ImageNet-1k with online evaluation protocol over 50 epochs. Bold indicates best performance.

Thank you!