Contrastive Learning: A Mathematical Perspective

Allowing Image And Text Data To Communicate

Ricardo Baptista

Statistical Sciences University of Toronto



IMSI: Data Assimilation and Inverse Problems for Digital Twins Workshop
October 7, 2025

Collaborators

A Mathematical Perspective on Contrastive Learning

Ricardo Baptista¹, Andrew Stuart^{2,3}, Son Tran³

¹Statistical Sciences University of Toronto

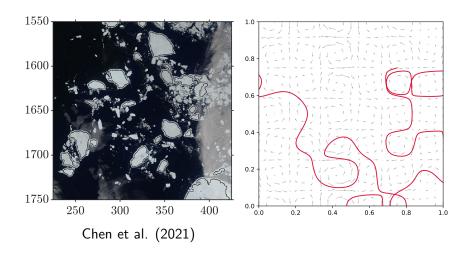
²Computing+Mathematical Sciences California Institute of Technology

> ³Stores Foundational AI Amazon

> > arXiv:2505.24134

Sampling Velocity Given Tracer Positions

- ▶ Infer ocean currents conditioned on data from passive tracers (e.g., sea ice floes)
- ▶ Relies on a model relating velocity to observations



Sampling Images Given Text

- Generate a distribution of images conditioned on a text prompt
- For example, using Stable Diffusion XL: Podell et al. [3] (2023)
- Relies on a model to link text and images

A cat and a frog







An elephant in the jungle





A skateboarder in California





Aligning Text and Images

- ► Key contrastive learning methodology CLIP: Radford et al. [4] (2021)
- 44135 Google Scholar citations as of October 6th 2025
- ▶ Cosine similarity between $a, b \in \mathbb{S}^{\ell-1}$: $\langle a, b \rangle$
- ▶ CLIP represents text as $a \in \mathbb{S}^{\ell-1}$ and image as $b \in \mathbb{S}^{\ell-1}$: then calculate $\langle a, b \rangle$.

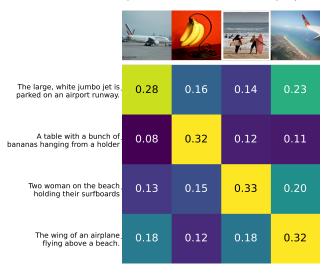


Table of Contents

CLIP Learning Problem

Generalizations of CLIP

Analytical Solutions for CLIP

Numerical Results

Conclusions

Table of Contents

CLIP Learning Problem

Generalizations of CLIF

Analytical Solutions for CLIP

Numerical Results

Conclusions

Setup for Contrastive Learning

Data is generated from a common reality

Text: $u \in \mathcal{U}$

Images: $v \in \mathcal{V}$

Distribution: $\mu(du, dv)$

Data pairs: $\{(u^i, v^i)\}_{i=1}^N \sim \mu \text{ i.i.d.}$

Embed data into a common low-dimensional space

$$g_u \colon \mathcal{U} \times \Theta \to \mathbb{S}^{\ell-1}$$

 $g_v \colon \mathcal{V} \times \Theta \to \mathbb{S}^{\ell-1}$

$$g_{\mathsf{v}} \colon \mathcal{V} \times \Theta \to \mathbb{S}^{\ell-1}$$

Encoders are represented using data-dependent architectures:

- ► Text: Byte-pair encoding and transformers
- ► Images: Convolution layers and transformers
- $ightharpoonup L^2$ normalization to map to the sphere

Contrastive Learning Problem

CLIP Objective Function Radford et al [4] (2021)

Find $\theta = (\theta_u, \theta_v, \tau)$ using only samples from μ

$$\begin{split} \mathsf{L}^{N}_{\mathsf{clip}}(\theta) &:= \frac{1}{N} \sum_{i=1}^{N} \langle g_{u}(u^{i}; \theta_{u}), g_{v}(v^{i}; \theta_{v}) \rangle / \tau \\ &- \frac{1}{2N} \sum_{i=1}^{N} \log \left(\sum_{j=1}^{N} \exp(\langle g_{u}(u^{i}; \theta_{u}), g_{v}(v^{j}; \theta_{v}) \rangle / \tau) \right) \\ &- \frac{1}{2N} \sum_{j=1}^{N} \log \left(\sum_{i=1}^{N} \exp(\langle g_{u}(u^{i}; \theta_{u}), g_{v}(v^{j}; \theta_{v}) \rangle / \tau) \right) \\ \theta^{\star} &= \arg \max_{\theta} \ \mathsf{L}^{N}_{\mathsf{clip}}(\theta). \end{split}$$

Objective aligned pair samples and penalizes unaligned pairs

Population-level Picture | B, Stuart and Tran (2025) [1]

Data Notation

Data Measure: $\mu(du, dv)$

Data Marginals: $\mu_u(du), \mu_v(dv)$

Data Conditionals: $\mu_{u|v}(du|v), \mu_{v|u}(dv|u)$

Target Notation

Target Measure: $\nu(du, dv; \theta) = \rho(u, v; \theta)\mu_u(du)\mu_v(dv)$

$$\rho(u, v; \theta) \propto \exp(\langle g_u(u; \theta_u), g_v(v; \theta_v) \rangle / \tau)$$

Target Conditionals: $\nu_{u|v}(du|v;\theta) = \rho(u|v;\theta)\mu_u(du)$

$$\nu_{v|u}(dv|u;\theta) = \rho(v|u;\theta)\mu_v(dv)$$

Conditionals are weighted by the marginal data measures.

Population-level Picture II B, Stuart and Tran (2025) [1]

Population Objective Function

$$\begin{split} \mathsf{L}_{\mathsf{cond}}(\theta) &= \frac{1}{2} \mathbb{E}_{(u,v) \sim \mu} \Big[\log \rho(u|v;\theta) + \log \rho(v|u;\theta) \Big], \\ &= \mathbb{E}_{(u,v) \sim \mu} \langle g_u(u;\theta_u), g_v(v;\theta_v) \rangle / \tau \\ &\quad - \frac{1}{2} \mathbb{E}_{v \sim \mu_v} \log \mathbb{E}_{u' \sim \mu_u} \exp \big(\langle g_u(u';\theta_u), g_v(v;\theta_v) \rangle / \tau \big) \\ &\quad - \frac{1}{2} \mathbb{E}_{u \sim \mu_u} \log \mathbb{E}_{v' \sim \mu_v} \exp \big(\langle g_u(u;\theta_u), g_v(v';\theta_v) \rangle / \tau \big). \\ \theta_{\mathsf{cond}} &= \arg \max_{\theta} \ \mathsf{L}_{\mathsf{cond}}(\theta). \end{split}$$

In practice, we minimize an empirical objective $L_{cond}^{N} \approx L_{cond}$

Theorem: Related to empirical CLIP Objective

$$\mathsf{L}^{N}_{\mathsf{clip}}(\theta) = \mathsf{L}^{N}_{\mathsf{cond}}(\theta) - \mathsf{log}(N)$$

For any N, the minimizers are the same, but only L_{cond}^N is well-defined in the population loss limit.

Population-level Picture III B, Stuart and Tran (2025) [1]

Recall: KL divergence $D_{kl}(\mu_1||\mu_2) = \int \log \frac{\mu_1(u)}{\mu_2(u)} \mu_1(du)$

Theorem: CLIP minimizes KL between conditionals

$$\mathsf{J}_{\mathsf{cond}}(\theta) = \frac{1}{2} \mathbb{E}_{\nu \sim \mu_{\nu}} \big[\mathsf{D}_{\mathsf{kl}} (\mu_{u|\nu} || \nu_{u|\nu} (\cdot; \theta)) \big] + \frac{1}{2} \mathbb{E}_{u \sim \mu_{u}} \big[\mathsf{D}_{\mathsf{kl}} (\mu_{\nu|u} || \nu_{\nu|u} (\cdot; \theta)) \big]$$

Assuming $\mu_{u|v} \ll \mu_u$ and $\mu_{v|u} \ll \mu_v$, then:

$$\mathop{\arg\min}_{\theta} \mathsf{J}_{\mathsf{cond}}(\theta) = \mathop{\arg\min}_{\theta} \mathsf{L}_{\mathsf{cond}}(\theta)$$

The non-parametric minimizer is $\nu_{u|v}(\cdot;\theta)=\mu_{u|v}$ and $\nu_{v|u}(\cdot;\theta)=\nu_{v|u}$.

Table of Contents

CLIP Learning Problem

Generalizations of CLIP

Analytical Solutions for CLIP

Numerical Results

Conclusions

Generalizations of CLIP I

Alignment Metric

Joint $\nu(u, v; \theta) \propto \rho(u, v; \theta) \mu_u(du) \mu_v(dv)$ is well-defined for any ρ such that

$$\int_{\mathcal{U}\times\mathcal{V}}\rho(u,v;\theta)\mu_u(du)\mu_v(dv)<\infty$$

Examples:

Normalized encoders: $|g_u(u;\theta_u)|_2 = 1$ and $|g_v(v;\theta_v)|_2 = 1$ with

$$\rho(u, v; \theta) = \exp\left(\langle g_u(u; \theta_u), g_v(v; \theta_v) \rangle\right) \propto \exp\left(-\frac{1}{2}|g_u(u; \theta_u) - g_v(v; \theta_v)|^2\right)$$

Un-normalized encoders with:

$$\rho(u, v; \theta) \propto \exp\left(-\frac{1}{2}|g_u(u; \theta_u) - g_v(v; \theta_v)|^2\right).$$

Generalizations of CLIP II

Loss Function

Given a divergence D and weights $\lambda_u, \lambda_v \geq 0$:

$$\begin{split} \mathsf{J}_{\mathsf{cond},\mathsf{D}}(\theta;\lambda_u,\lambda_v) &:= \frac{\lambda_u}{2} \mathbb{E}_{\nu \sim \mu_v} \big[\mathsf{D}(\mu_{u|\nu}||\nu_{u|\nu}(\cdot;\theta)) \big] + \frac{\lambda_v}{2} \mathbb{E}_{u \sim \mu_u} \big[\mathsf{D}(\mu_{\nu|u}||\nu_{\nu|u}(\cdot;\theta)) \big] \\ \theta_{\mathsf{cond}}^* &= \arg\min_{\theta} \mathsf{J}_{\mathsf{cond},\mathsf{D}}(\theta;\lambda_u,\lambda_v) \end{split}$$

Examples:

- ▶ Choose $(\lambda_u, \lambda_v) = (1, 0)$ to match the u|v conditional
- ▶ Choose divergence D computable from samples, e.g. maximum mean discrepancy:

$$\mathsf{D}(\mu_1,\mu_2) = \mathbb{E}_{(u,u')\sim\mu_1\otimes\mu_1}k(u,u') - 2\mathbb{E}_{(u,v)\sim\mu_1\otimes\mu_2}k(u,v) + \mathbb{E}_{(v,v')\sim\mu_2\otimes\mu_2}k(v,v')$$

Generalizations of CLIP III

Joint loss function

$$\begin{aligned} \mathsf{J}_{\mathsf{joint}}(\theta) &= \mathsf{D}_{\mathsf{kl}}(\mu||\nu(\cdot;\theta)) \\ \theta_{\mathsf{joint}}^* &= \operatorname*{arg\,min}_{\theta} \mathsf{J}_{\mathsf{joint}}(\theta) \end{aligned}$$

Implementable loss function

$$\begin{split} \mathsf{L}_{\mathsf{joint}}(\theta) &= \mathbb{E}_{(u,v) \sim \mu} \langle g_u(u;\theta_u), g_v(v;\theta_v) \rangle / \tau \\ &- \log \mathbb{E}_{(u,v) \sim \mu_u \otimes \mu_v} \exp \big(\langle g_u(u;\theta_u), g_v(v;\theta_v) \rangle / \tau \big) \\ \theta_{\mathsf{joint}}^* &= \arg \max_{\theta} \mathsf{L}_{\mathsf{joint}}(\theta) \end{split}$$

Objective maximizes alignment of paired data and minimizes un-alignment of un-paired data.

Empirical approximation L_{joint}^{N} is more efficient to evaluate than L_{clip}^{N} .

Table of Contents

CLIP Learning Problem

Generalizations of CLIP

Analytical Solutions for CLIP

Numerical Results

Conclusions

Low-Rank Matrix Approximations I

Goal: Analyze closed-form solutions for CLIP optimization

Setting: Gaussian data distribution $\mu = \mathcal{N}(0,\mathcal{C})$ with block covariance matrix

$$\mathcal{C} = \begin{bmatrix} \mathcal{C}_{uu} & & \mathcal{C}_{uv} \\ \mathcal{C}_{vu} & & \mathcal{C}_{vv} \end{bmatrix}$$

Model: Linear encoders project $u \in \mathbb{R}^{n_u}$ and $v \in \mathbb{R}^{n_v}$ into \mathbb{R}^{ℓ} with $\ell \leq \min(n_u, n_v)$

$$g_u(u) = Gu, \quad G \in \mathbb{R}^{n_u \times \ell},$$

 $g_v(v) = Hv \quad H \in \mathbb{R}^{n_v \times \ell}.$

CLIP Probabilistic Model:

$$\nu(u, v; \theta) \propto \exp(\langle Gu, Hv \rangle) \mu_u(u) \mu_v(v)$$

= $\exp(\langle u, Av \rangle) \mu_u(u) \mu_v(v), \quad A = G^\top H.$

Low-Rank Matrix Approximations II

Theorem: Minimizing Two-Sided Loss

Minimizing L_{cond} over A with $\ell = min(n_u, n_v)$ has solution

$$A = \mathcal{C}_{uu}^{-1} \mathcal{C}_{uv} \mathcal{C}_{vv}^{-1}$$

resulting in the conditional distributions

$$\nu_{u|v}(u|v;\theta^*) = \mathcal{N}(\mathcal{C}_{uv}\mathcal{C}_{vv}^{-1}v;\mathcal{C}_{uu})
\nu_{v|u}(v|u;\theta^*) = \mathcal{N}(\mathcal{C}_{vu}\mathcal{C}_{uu}^{-1}u;\mathcal{C}_{vv})$$

Corollary: Conditional means of $\mu_{u|v}, \mu_{v|u}$, but not variances, are recovered.

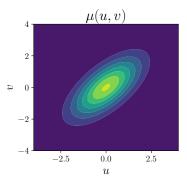
Theorem: With Dimension Reduction

Minimizing L_{cond} over A with $\ell < \min(n_u, n_v)$ has solution

$$A = C_{uu}^{-\frac{1}{2}} (C_{uu}^{-\frac{1}{2}} C_{uv} C_{vv}^{-\frac{1}{2}})_{\ell} C_{vv}^{-\frac{1}{2}}$$

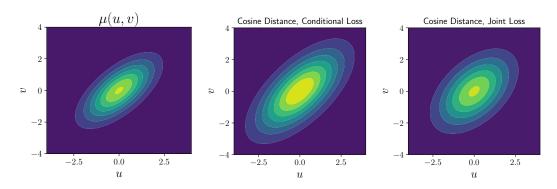
Visualizations of CLIP Generalizations I

- lacktriangle Two-dimensional Gaussian data distribution $\mu=\mathcal{N}(0,\mathcal{C})$
- ▶ Used un-normalized linear encoders to preserve Gaussian structure



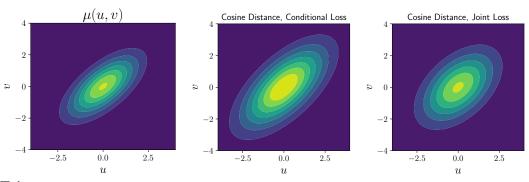
Visualizations of CLIP Generalizations I

- lacktriangle Two-dimensional Gaussian data distribution $\mu = \mathcal{N}(0, \mathcal{C})$
- Used un-normalized linear encoders to preserve Gaussian structure



Visualizations of CLIP Generalizations I

- lacktriangle Two-dimensional Gaussian data distribution $\mu = \mathcal{N}(0, \mathcal{C})$
- ▶ Used un-normalized linear encoders to preserve Gaussian structure



Takeaway:

► Conditional means, but not variance are matched with two-sided conditional loss

Visualizations of CLIP Generalizations II

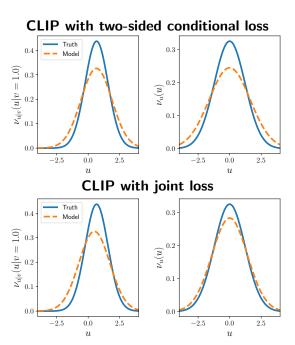


Table of Contents

CLIP Learning Problem

Generalizations of CLIF

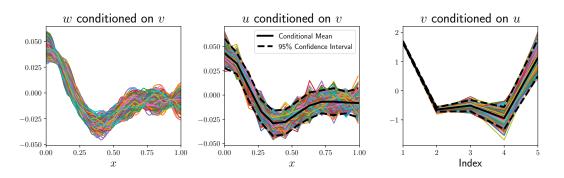
Analytical Solutions for CLIP

Numerical Results

Conclusions

Gaussian Example I

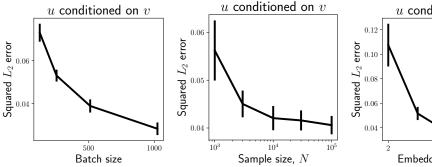
- ▶ Reality: $w \sim \mathcal{N}(0, C)$ is an un-observed Gaussian process in $L^2((0, 1); \mathbb{R})$
- $u \in \mathbb{R}^{12}$ is a noisy pointwise evaluation of w at uniformly-spaced grid locations
- ▶ $v \in \mathbb{R}^5$ is five leading Fourier coefficients of w

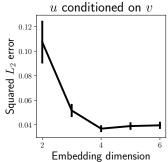


Note: $\mu_{u|v}$ and $\mu_{v|u}$ are Gaussian with closed-form conditional means and covariances

Gaussian Example II

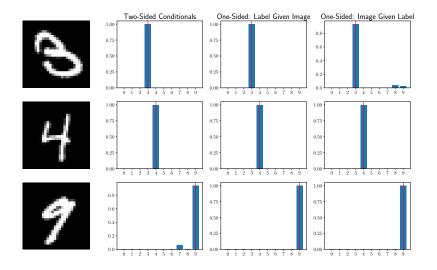
- Used un-normalized linear encoders to represent Gaussian conditionals
- Minimized L_{cond}^{N} to learn encoders for both modalities
- Compared the approximation to the true conditional expectations for u|v and v|u
- Varied the batch size, total training samples N and embedding dimension ℓ





MNIST Example

- lacksquare Data modalities: images $u\in\mathcal{U}=[0,1]^{28 imes28}$ and digits $v\in\mathcal{V}=\{0,\dots,9\}$
- ▶ Performed classifications using models learned with different loss functions

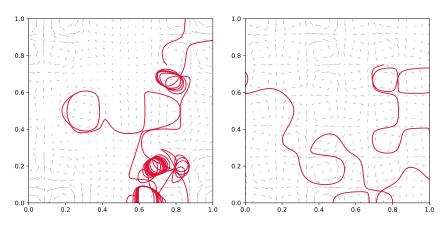


Lagrangian Data Assimilation

Original problem formulation: Kuznetsov, Ide and Jones (2003)[2]

- ▶ Velocity Field (Image); Lagrangian Trajectory (Text).
- ► Retrieval: identify the most likely velocity from a database given a trajectory

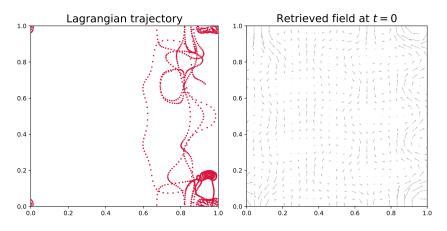
Paired potentials (background) and trajectories (red)



Lagrangian Data Assimilation

- ► CLIP model can align trajectory data (e.g., time-series or text) with time-dependent potential functions (e.g., images)
- ▶ Model predicts Fourier representation of the potential directly from trajectories

Paired potentials (background) and trajectories (red)



Lagrangian Data Assimilation

Measured accuracy of retrieval between both modalities

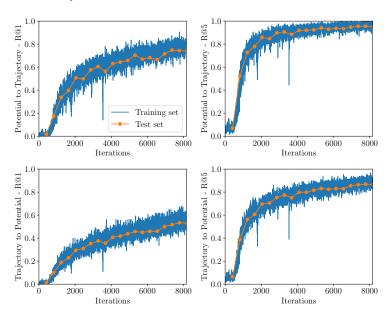


Table of Contents

CLIP Learning Problem

Generalizations of CLIF

Analytical Solutions for CLIF

Numerical Results

Conclusions

Conclusions

Main Messages

- ► Contrastive learning relates two data modalities by tilting a product distribution
- ► Generalization yields new probabilistic loss functions and alignment metrics
- ► Encoders have closed-form low-rank matrix solutions in linear-Gaussian settings
- Application to Lagrangian data assimilation

Outlook

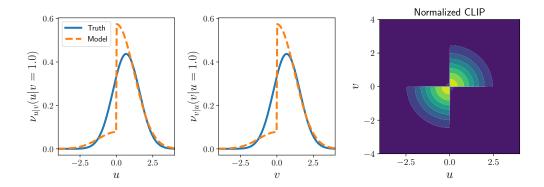
- Extending framework to more than two modalities
- Theoretical study of compositional generalization
- Other applications in science and engineering

References I

- R. Baptista, A. M. Stuart, and S. Tran. A mathematical perspective on contrastive learning. arXiv:2505.24134, 2025.
- [2] L. Kuznetsov, K. Ide, and C. K. Jones. A method for assimilation of Lagrangian data. Monthly Weather Review, 131(10):2247–2260, 2003.
- [3] D. Podell, Z. English, K. Lacey, A. Blattmann, T. Dockhorn, J. Müller, J. Penna, and R. Rombach. SDXL: Improving latent diffusion models for high-resolution image synthesis. arXiv:2307.01952, 2023.
- [4] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, et al. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning*, pages 8748–8763. PMLR, 2021.

Visualizations of CLIP Generalizations IV

- Evaluated approximation with commonly-used normalized encoders
- \triangleright Encoders capture the sign of the correlation between u and v



MNIST Example II

- Generated images conditioned on a digit using models learned with different losses
- lacktriangleright 16 images are sampled from the conditional distribution $u_{u|v}$



Takeaway: One-sided loss for classification shows mode collapse onto a single image