

Data Reliability Scoring

Yiling Chen¹ Shi Feng¹ Paul Kattuman² Fang-Yi Yu³

Harvard University¹ University of Cambridge² George Mason University³



Key Contributions

Can we assess the reliability of a dataset without access to the ground truth?

- Formalize the problem of reliability scoring based on observations from unknown experiments
- Propose dataset reliability orderings to benchmark reliability
- Design the Gram Determinant Score which can preserve dataset reliability orderings
- Conduct experiment on synthetic noise models and CIFAR-10 embeddings

Problem Formulation

Basic Setup

Reliability score $S(\hat{\mathbf{x}}, \mathbf{y})$ should assign a higher expected score to datasets that are closer to the true data.

Reports $\hat{\mathbf{x}}$	\hat{x}_1	\hat{x}_2	...	\hat{x}_N
			↑ Misreport matrix Q	
True data \mathbf{x}	x_1	x_2	...	x_N
	↓	↓	Experiment P	↓
Observations \mathbf{y}	y_1	y_2	...	y_N

- Three datasets of length N
 - (Hidden) True data $\mathbf{x} = (x_1, \dots, x_N) \in [d]^N$
 - Reported data $\hat{\mathbf{x}} = (\hat{x}_1, \dots, \hat{x}_N) \in [d]^N$
 - Observable data $\mathbf{y} = (y_1, \dots, y_N) \in \mathcal{Y}^N$
- Experiment $P \in \mathbb{R}^{\mathcal{Y} \times d}$: $y_n \sim P_{x_n}$ independently for all n .
- Misreport matrix $Q \in \mathbb{R}^{d \times d}$: $Q(i, j) = \frac{1}{N} \sum_{n=1}^N \mathbf{1}[x_n = i, \hat{x}_n = j]$. Conditional frequency matrices $Q_{\hat{\mathbf{x}}|\mathbf{x}}$ and $Q_{\mathbf{x}|\hat{\mathbf{x}}}$ are column stochastic.

Examples

- Insurance Reimbursement: Let \mathbf{x} represent the true disease states of patients, $\hat{\mathbf{x}}$ the diagnoses reported by the hospital, and \mathbf{y} auxiliary signals such as blood test results or imaging biomarkers.
- Image Labeling: If the dataset includes image labels, let \mathbf{x} be the true labels, $\hat{\mathbf{x}}$ the reported labels, and \mathbf{y} encoder-based representations or other observations derived from the images.

Reliability Orderings of Datasets

To compare the reliability of different reported datasets, $\hat{\mathbf{x}}$ and $\hat{\mathbf{x}}'$, relative to the true data \mathbf{x} , we propose strict partial orderings of reported datasets.

- Exact Match Ordering: $\hat{\mathbf{x}} \succ_{EXACT}^{\mathbf{x}} \hat{\mathbf{x}}'$ if $\hat{\mathbf{x}} = \mathbf{x}$ but $\hat{\mathbf{x}}' \neq \mathbf{x}$.
- Blackwell dominant ordering: $\hat{\mathbf{x}} \succ_{Blackwell}^{\mathbf{x}} \hat{\mathbf{x}}'$ if $\hat{\mathbf{x}}'$ is a garbling of $\hat{\mathbf{x}}$.
- Hamming ordering: $\hat{\mathbf{x}} \succ_{Hamming}^{\mathbf{x}} \hat{\mathbf{x}}'$ if $\hat{\mathbf{x}}$ is closer to \mathbf{x} than $\hat{\mathbf{x}}'$ in Hamming distance.

Reliability Scoring Preserving Reliability Orderings

Given a reliability ordering \succ , a reliability score S preserves the partial ordering \succ under experiment P , if

$$\mathbb{E}_{\mathbf{y} \sim P(\mathbf{x})}[S(\hat{\mathbf{x}}, \mathbf{y})] > \mathbb{E}_{\mathbf{y} \sim P(\mathbf{x})}[S(\hat{\mathbf{x}}', \mathbf{y})]. \quad (1)$$

for all $\mathbf{x}, \hat{\mathbf{x}}, \hat{\mathbf{x}}' \in \mathcal{X}^N$ with $\hat{\mathbf{x}} \succ^{\mathbf{x}} \hat{\mathbf{x}}'$.

Additionally, S asymptotically preserves \succ under a set of experiments \mathcal{P} and a set of misreport matrices \mathcal{Q} , if for all $P \in \mathcal{P}$ and $Q, Q' \in \mathcal{Q}$ and large enough N , S preserves \succ under P for all $\mathbf{x}, \hat{\mathbf{x}}, \hat{\mathbf{x}}'$ with $\hat{\mathbf{x}} \succ^{\mathbf{x}} \hat{\mathbf{x}}'$ and misreport matrices Q, Q' .

Geometric Intuition of Gram Determinant

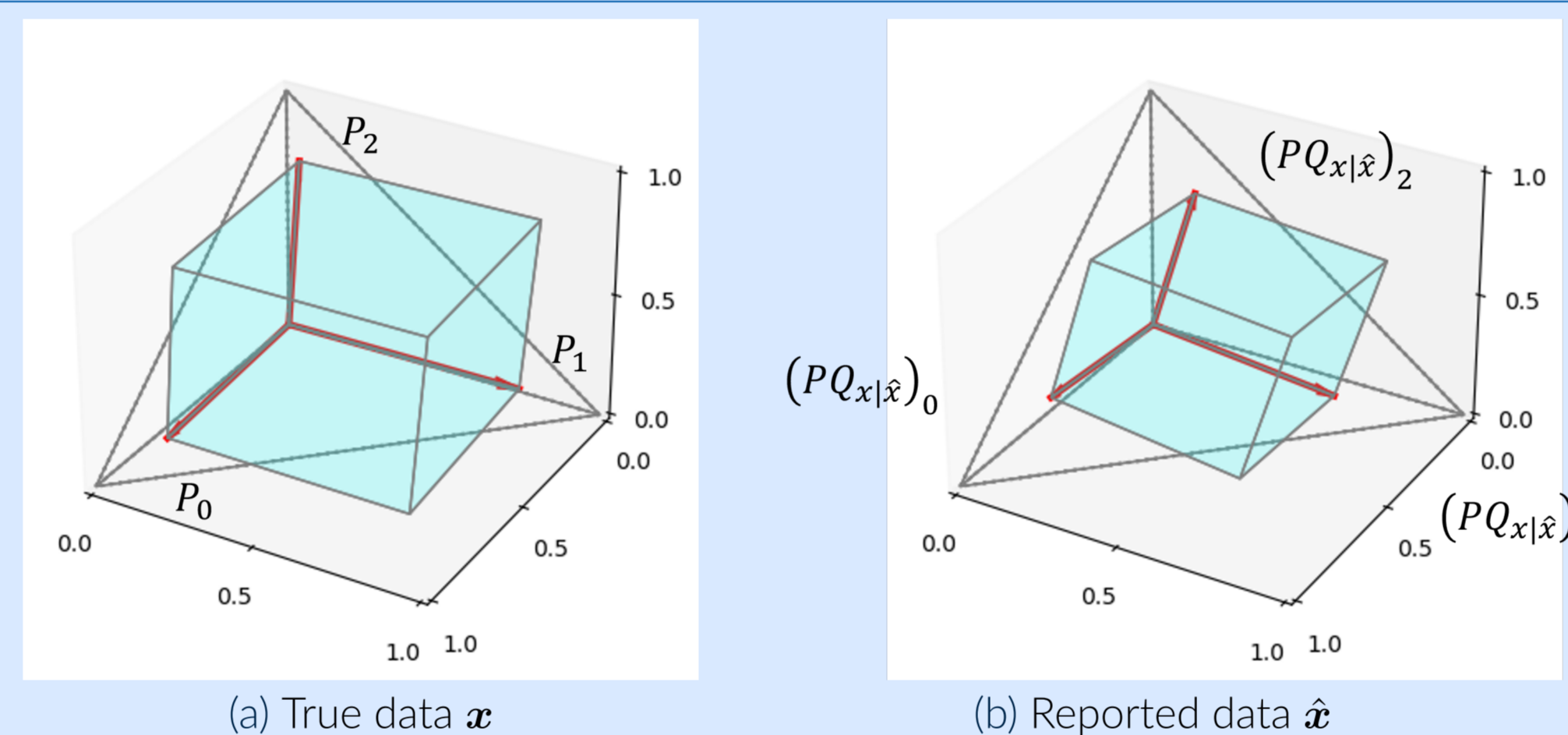


Figure 1. Parallelepipeds spanned by (normalized) observations with $d = 3$ and $|\mathcal{Y}| = 3$

Given P , the average distribution of observation under each true data label $x \in [d]$ is P_x , and the volume of the spanned parallelepiped is the Gram determinant of P

$$\text{vol}(P_0, \dots, P_{d-1}) = \sqrt{\det(P^\top P)}.$$

Given \mathbf{x} and $\hat{\mathbf{x}}$, the distribution under each reported label x becomes $(PQ_{\mathbf{x}|\hat{\mathbf{x}}})_x$, with volume $\sqrt{\det((PQ_{\mathbf{x}|\hat{\mathbf{x}}})^\top (PQ_{\mathbf{x}|\hat{\mathbf{x}}}))} \leq \sqrt{\det(P^\top P)}$

Gram Determinant Reliability Scores

Gram determinant scores measure the volume spanned by the vectors of observed experiment outcomes.

Let $G = P^\top P$ where $G(x, x') = \langle P_x, P_{x'} \rangle$, the Gram matrix of reports $\hat{\mathbf{x}}$ be $\hat{G} = Q^\top G Q$ where $\hat{G}(x, x') := \frac{1}{N^2} \sum_{n, n': \hat{x}_n = x, \hat{x}_{n'} = x'} \langle P_{x_n}, P_{x_{n'}} \rangle$, and the Gram determinant score is

$$\Gamma := \det(\hat{G}). \quad (2)$$

The Gram determinant score in preserves

- exact match ordering under \mathcal{P}_{indep} and $\mathcal{Q}_{nonperm}$,
- Blackwell ordering under \mathcal{P}_{indep} and \mathcal{Q}_{reg} , and
- $\frac{1}{4L}$ -Hamming ordering under \mathcal{P}_{indep} and $\mathcal{Q}_{L, 1/64L^2 d^2}$.

Properties of Gram determinant reliability scores

Estimating Gram Determinant Scores

- \hat{G} admits estimator $\bar{G}(x, x') = \frac{1}{N^2} \sum_{n, n' \in [N]: \hat{x}_n = x, \hat{x}_{n'} = x'} \mathbf{1}[y_n = y_{n'}]$ which is consistent as $N \rightarrow \infty$.
- We can decrease the bias of the above plug-in estimator \bar{G} through a carefully designed sampling scheme.

These estimators for the Gram determinant scores asymptotically preserve the above reliability orderings.

Kernelized Gram Determinant Scores

If \mathcal{Y} has a kernel K , we can define G_K so that $G_K(x, x') = \langle P_x, P_{x'} \rangle_K := \mathbb{E}_{y \sim P_x, y' \sim P_{x'}}[K(y, y')]$, $\hat{G}_K = Q^\top G_K Q$ and Gram determinant score with kernel K as $\Gamma_K := \det(\hat{G}_K)$.

- Linear feature: $K(y, y') = \langle \phi(y), \phi(y') \rangle$, $\phi: \mathcal{Y} \rightarrow \mathbb{R}^k$ e.g. one-hot encoder.
- Gaussian RBF: $K(y, y') = \exp\left(-\frac{\|y - y'\|_2^2}{2\sigma^2}\right)$.
- Dot product: $K(y, y') = y^\top y'$ e.g., observation are forecast vectors.

Empirical Results

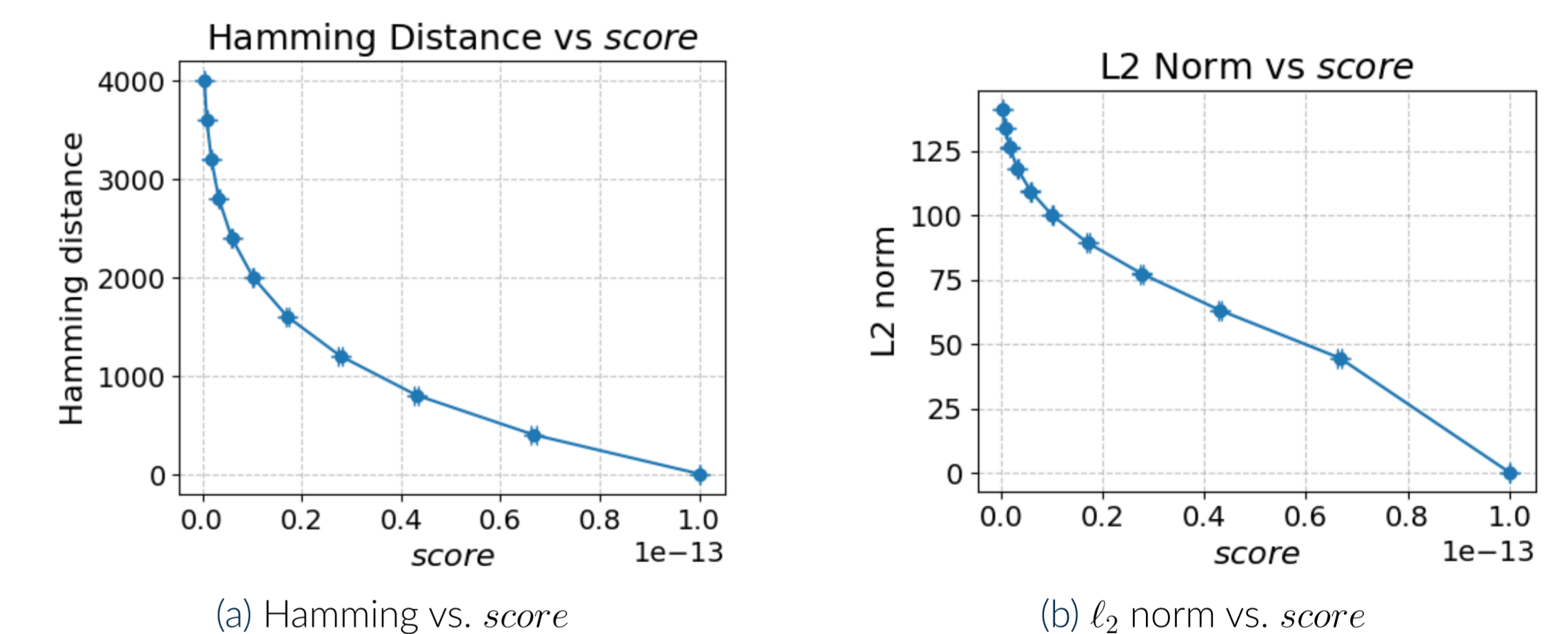


Figure 2. Gram Determinant Score on Synthetic Data

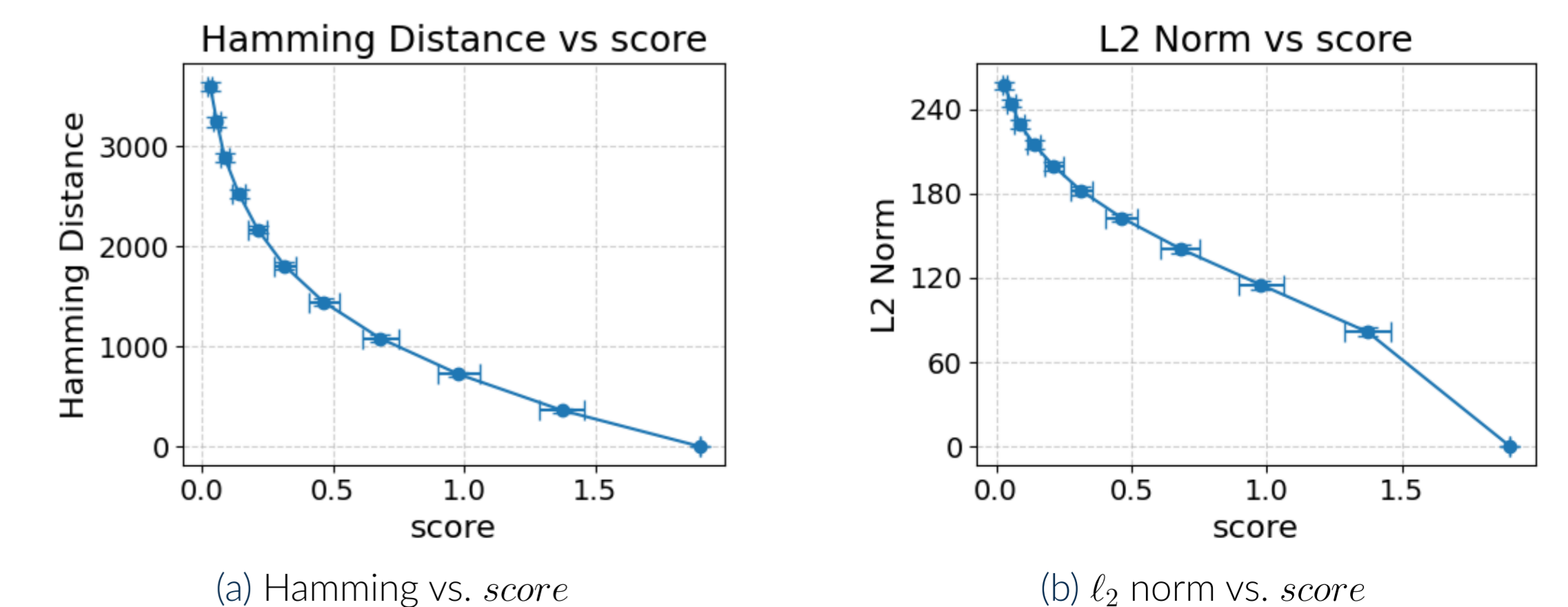


Figure 3. Gram Determinant Score with Kernels on Image Data