

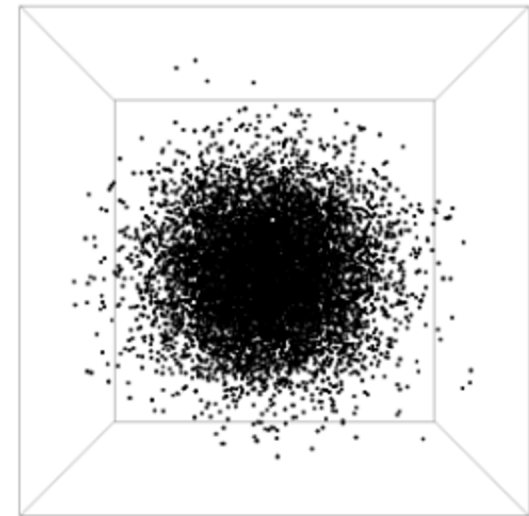
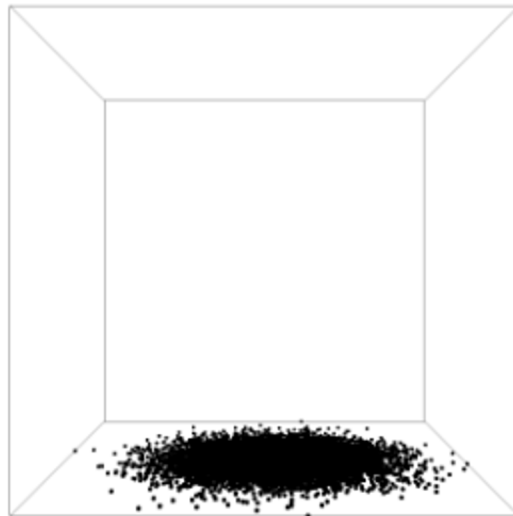
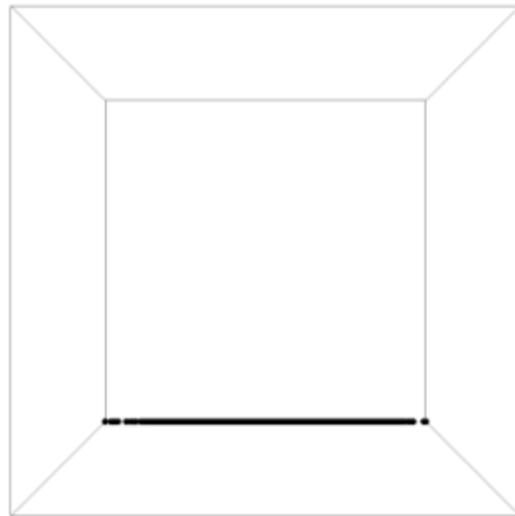


Stable Anisotropic Regularization



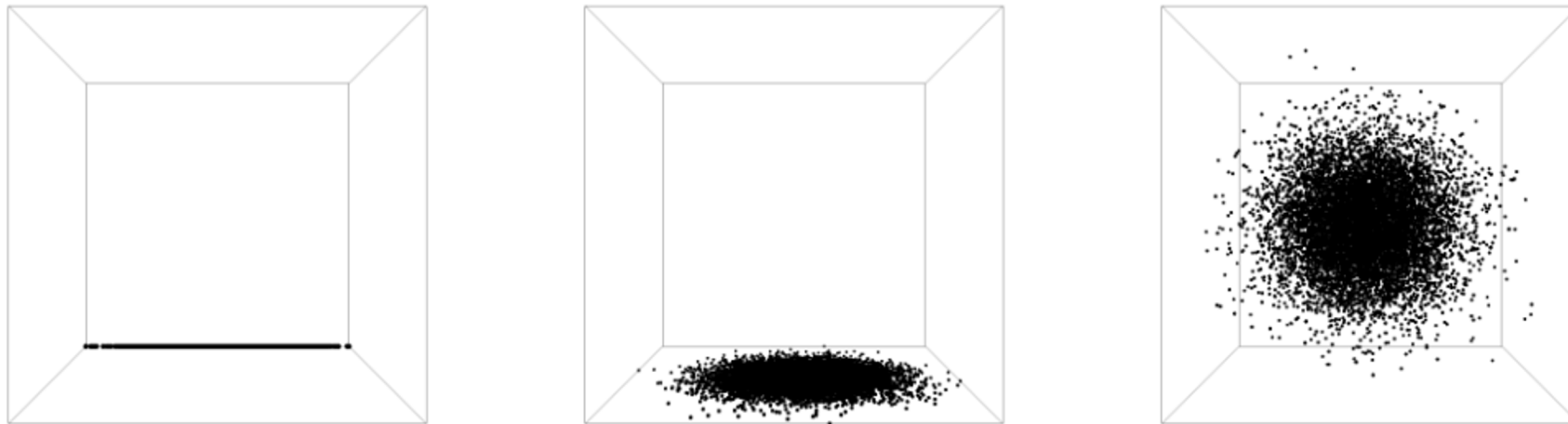
William Rudman¹ & Carsten Eickhoff²

Department of Computer Science, Brown University¹
School of Medicine, University of Tübingen²



Defining Isotropy

Isotropy: a distribution is *isotropic* if the variance of the data is uniformly distributed (i.e. the covariance matrix is proportional to the identity matrix).



Left: line embedded in 3D space. **Middle:** circle embedded in 3D space. **Right:** Sphere in 3D space.

Geometry of LLM Embeddings

- **Narrow Cone Hypothesis:**
 - Average random cosine similarity approaches 1 (Ethayarajh 2019).
 - Limits expressiveness of textual representations (Zhang et al. 2020; Timkey & Van Schijndel et al. 2021).
- **Limits Downstream Performance:**
 - Making embeddings more uniformly distributed improves performance (Rajaei et al. 2021; Zhou et al. 2021).

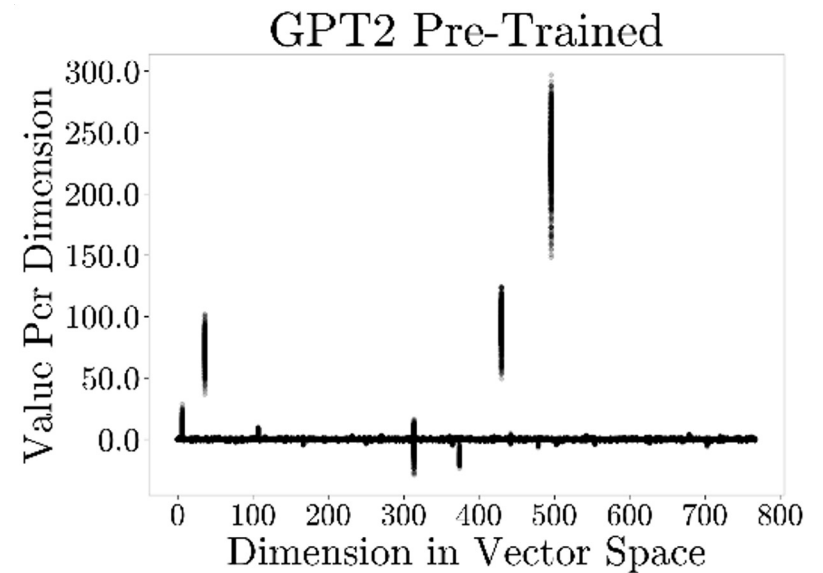


Figure 1: Average activation diagrams of sentence embeddings on the SST-2 validation dataset. The x-axis represents the index of the dimension, and the y-axis is the magnitude in that given dimension.

Methods for Improving Isotropy

- **All-But-The-Top:**
 - Post-processing algorithm that removes the *top* principal components from embeddings.
- **Cosine Similarity Regularization:**
 - CosReg adds a cosine similarity penalty to the standard CrossEntropy loss.
- **Can we regularize using IsoScore?:**
 - IsoScore is a more recent tool design for measuring isotropy.
 - *Not stable on mini-batch computations.*

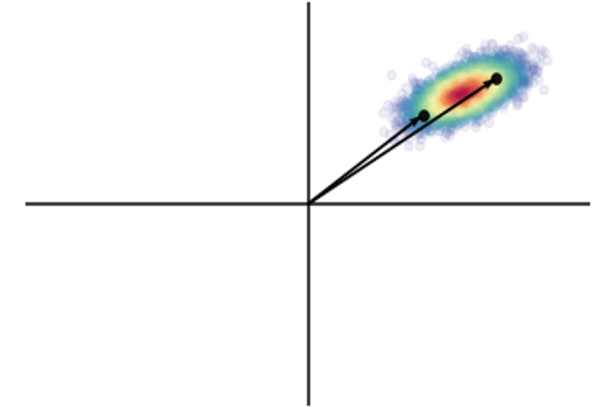


Figure 2: Average cosine similarity of **0.978**.

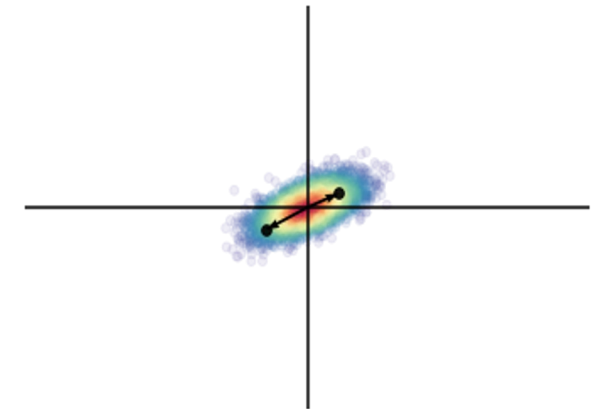
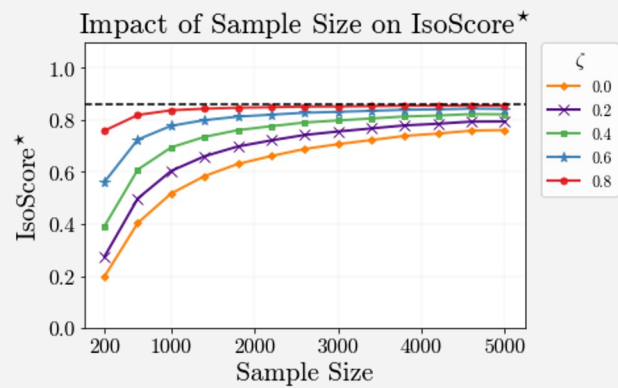
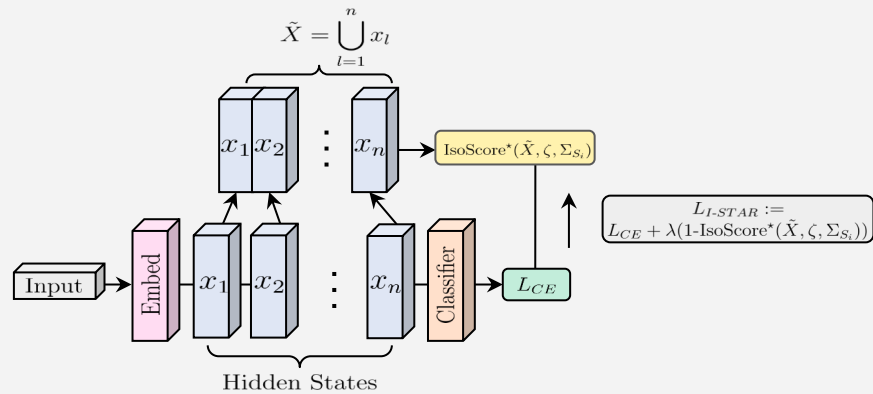


Figure 3: Average cosine similarity of **0.005**.

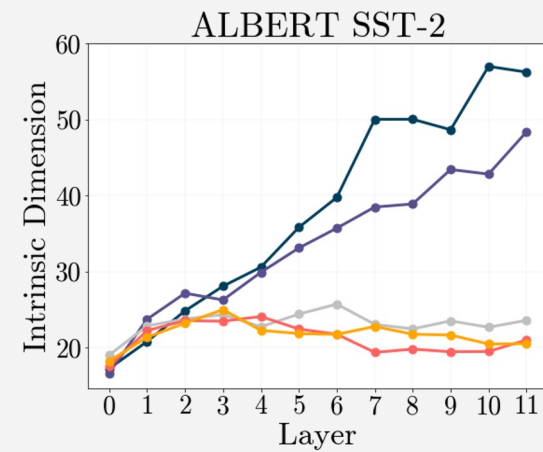
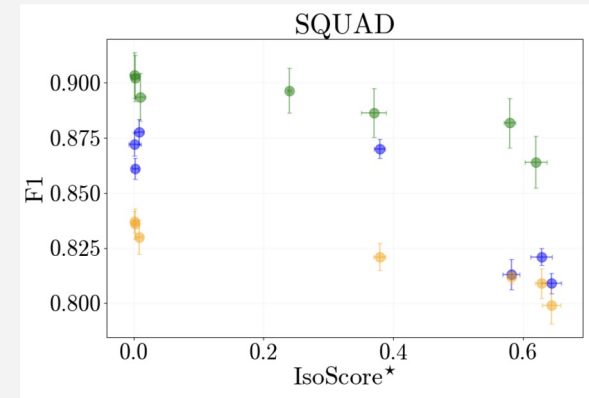
AIM 1

Develop a stable method of measuring isotropy.



AIM 2

Understand how isotropy correlates to model performance.



IsoScore*: A Stable and Differentiable Metric of Isotropy

- **IsoScore***:
 - Stable on mini-batch computations thanks to Regularized Discriminant Analysis (Friedman 1989).
 - $\text{IsoScore}^*(X) = 1$ implies that all principal components are equal.
 - $\text{IsoScore}^*(X) = 0$ implies that exactly principal component is non-zero.
- **IsoScore vs. IsoScore***:
 - Approaches IsoScore when the number of samples is large.
 - Equivalent IsoScore when no RDA-shrinkage is performed.

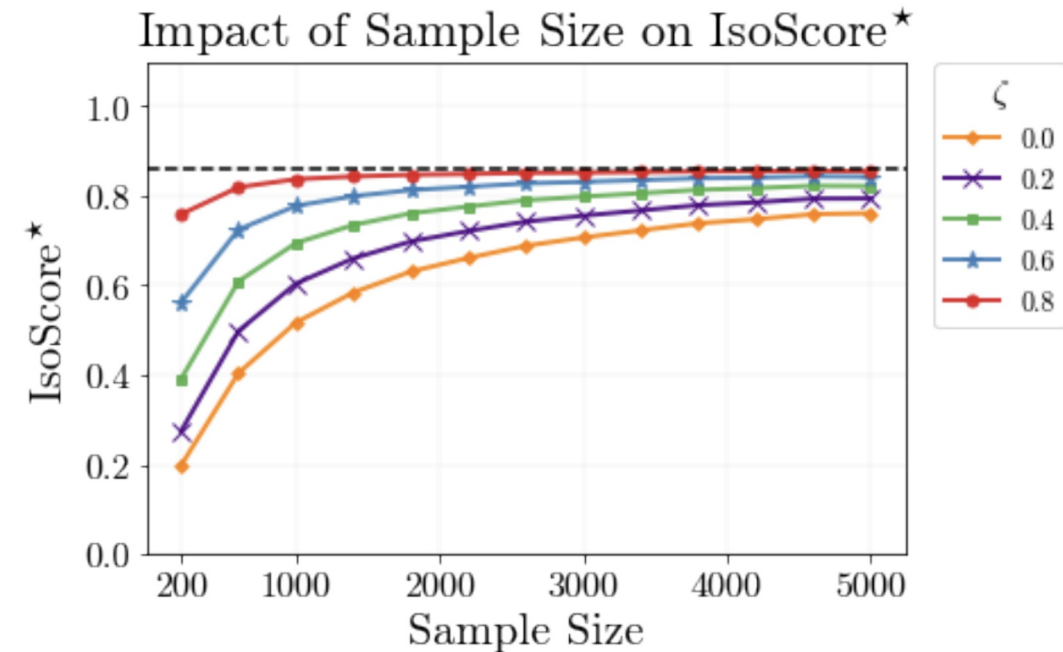


Figure 4: Impact of varying the RDA-shrinkage parameter on IsoScore*. Note that when $\zeta = 0$, IsoScore = IsoScore*. Dashed line is the true isotropy score of the point cloud.

Algorithm 1 IsoScore* Forward Pass

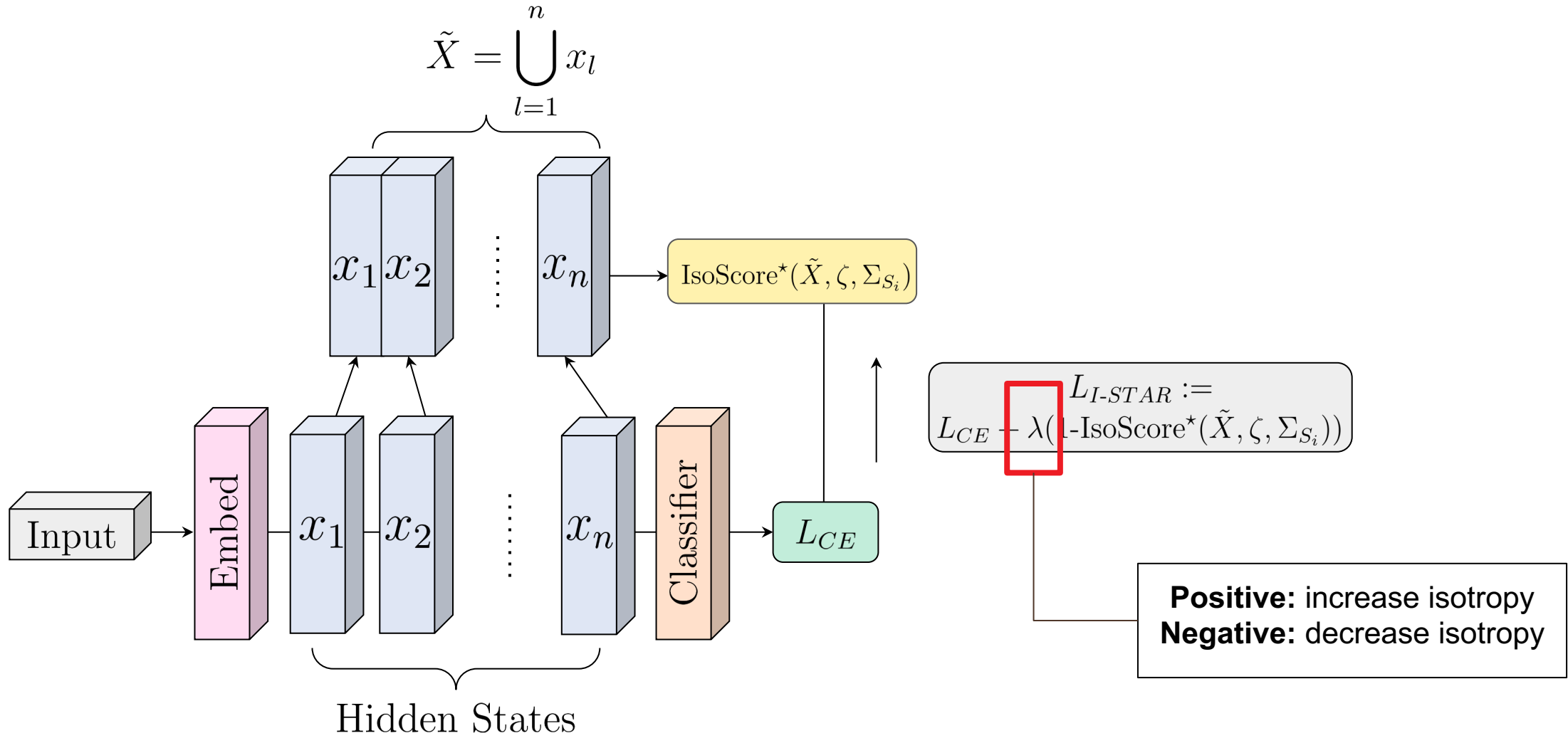
- 1: **Input:** $X \subset \mathbb{R}^d$ point cloud, $\Sigma_S \in \mathbb{R}^{d \times d}$ shrinkage covariance matrix, $\zeta \in (0, 1)$.
- 2: **Outputs:** I-STAR penalty of X .
- 3: calculate covariance matrix: Σ_X of X
- 4: calculate shrinkage matrix: $\Sigma_\zeta := (1 - \zeta) \cdot \Sigma_X + \zeta \cdot \Sigma_S$
- 5: calculate eigenvalues: $\Lambda := \{\lambda_1, \dots, \lambda_d\}$ of Σ_ζ
- 6: normalize eigenvalues: $\hat{\Lambda} := \sqrt{d} \cdot \Lambda / \|\Lambda\|_2$ such that $\|\hat{\Lambda}\| = \sqrt{d}$
- 7: calculate the isotropy defect:

$$\delta(\hat{\Lambda}) := \|\hat{\Lambda} - \mathbf{1}\| / \sqrt{2(d - \sqrt{d})}$$

where $\mathbf{1} = (1, \dots, 1)^\top \in \mathbb{R}^d$

- 8: calculate: $\phi(\hat{\Lambda}) := (d - \delta(\hat{\Lambda}))^2 (d - \sqrt{d})^2 / d^2$
 - 9: calculate: $\iota(\hat{\Lambda}) := (d \cdot \phi(\hat{\Lambda}) - 1) / (d - 1)$.
-

I-STAR Loss



Isotropy Negatively Correlates with Performance

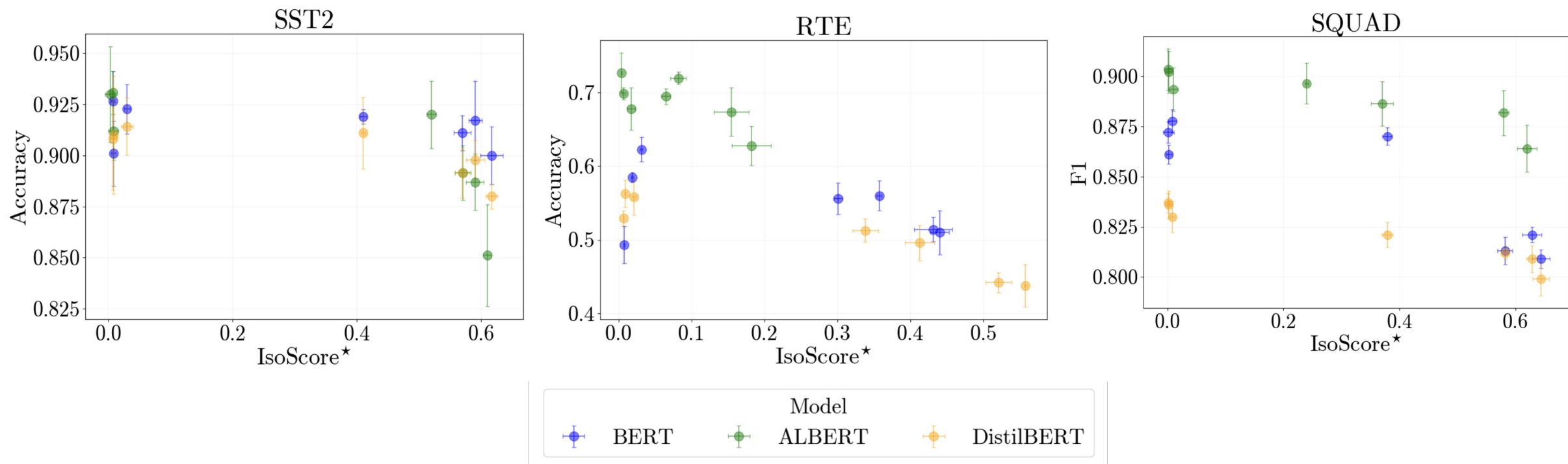


Figure 5: Relationship between IsoScore* (x-axis) and model performance (y-axis). We fine-tune each model with I-STAR using the tuning parameters, λ , in $\{-5, -3, -1, 0.50, 1, 3, 5\}$. We train each model over five random seeds and report the standard deviation of both performance and IsoScore*

Further Decreasing Isotropy with I-STAR Improves Performance

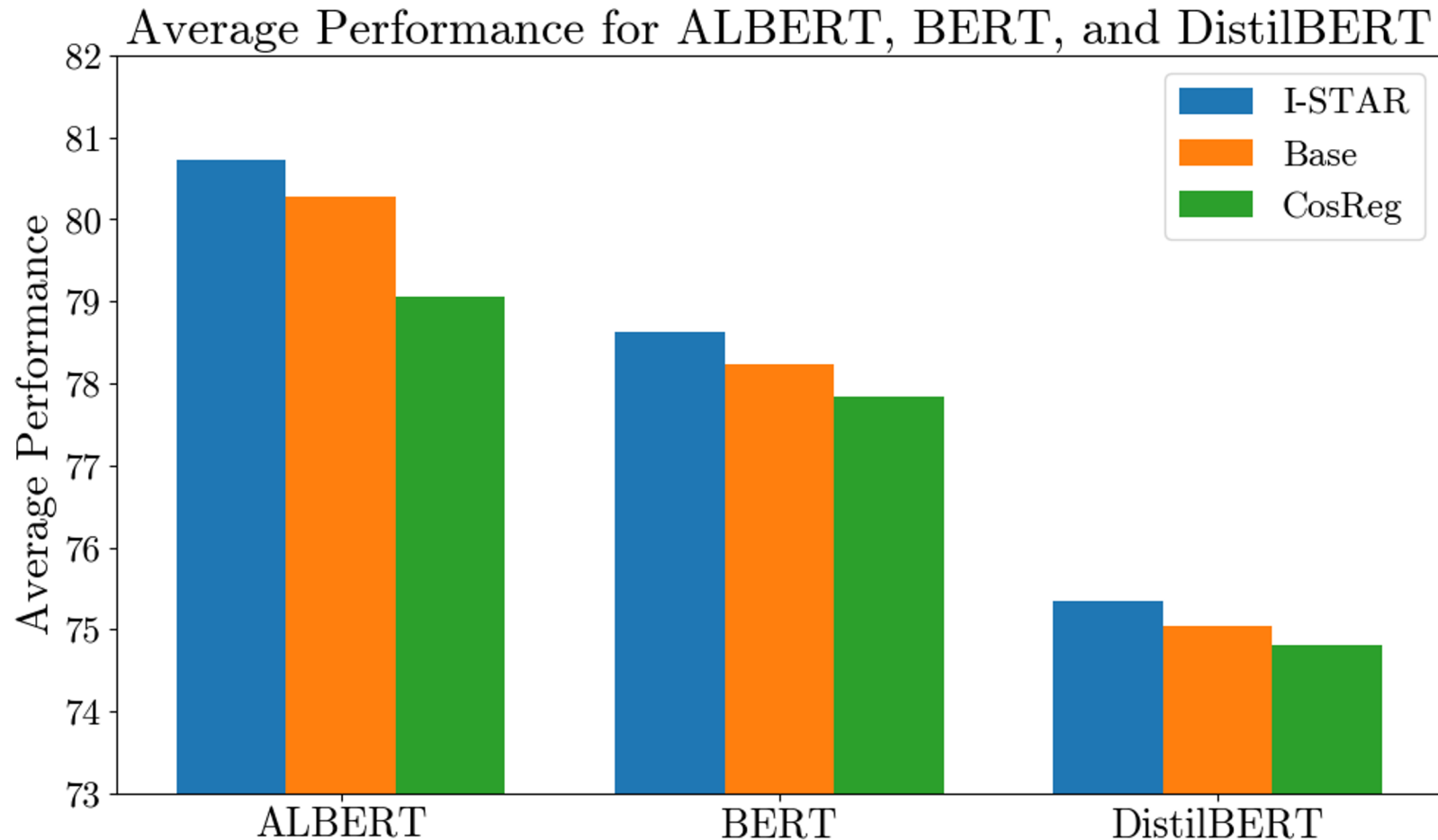


Figure 6: Average performance of ALBERT, BERT and DistilBERT fine-tuned using I-STAR (negative tuning params only), no regularization and Cosine Similarity regularization. Average is computed across 9 common NLP Benchmarks. GLUE (7 tasks), SST-5 and SQUAD.

Why do models prefer anisotropy?

- **Benefits of Anisotropic Noise:**
 - Helps models escape local minima in the loss landscape (Zhu et al. 2018).
- **I-STAR Alters Intrinsic Dimension:**
 - Low intrinsic dimension correlates with improved classification performance (Ansuini et al. 2019).
- **Minimizing IsoScore Maximizes Silhouette Scores:**
 - Isotropy objectives are incompatible with clustering objectives (Mickus et al. 2024).

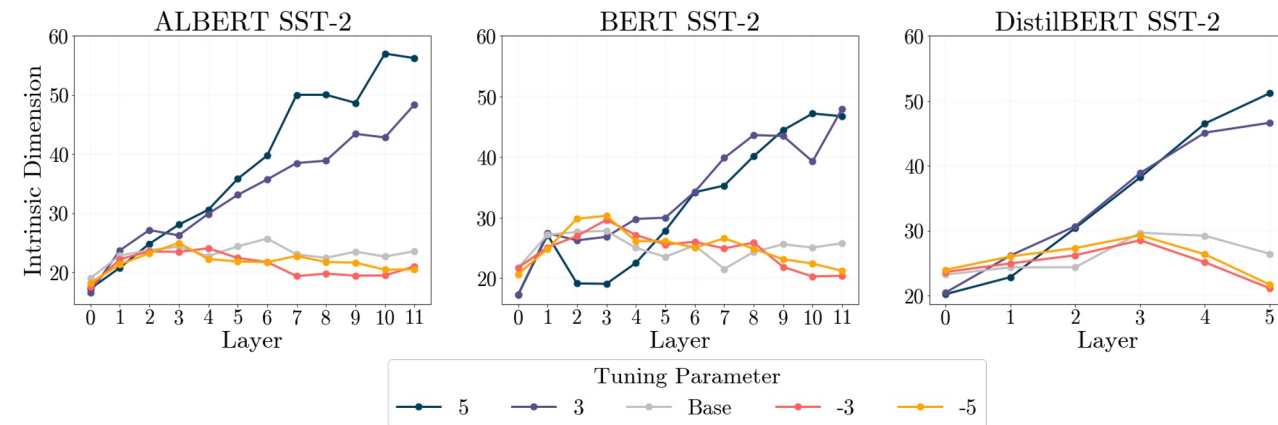
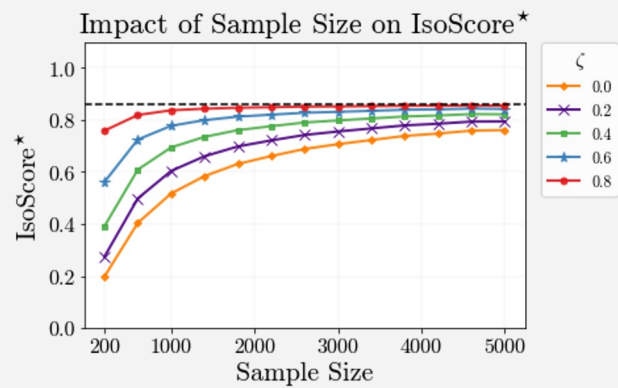
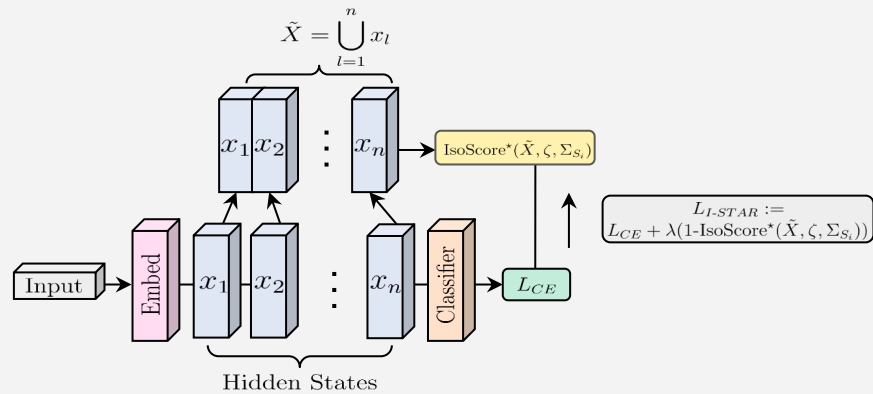


Figure 7: TwoNN Intrinsic Dimensionality estimate of ALBERT, BERT, and DistilBERT sentence embeddings obtained from the SST-2 validation data for models fine-tuned on the SST-2 using I-STAR with tuning-parameters. Base' represents the case where no regularization is used.

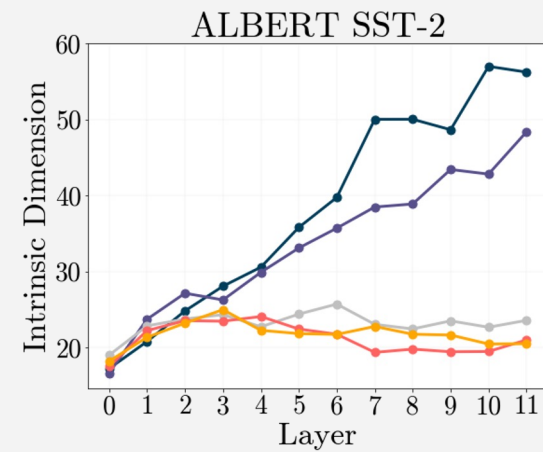
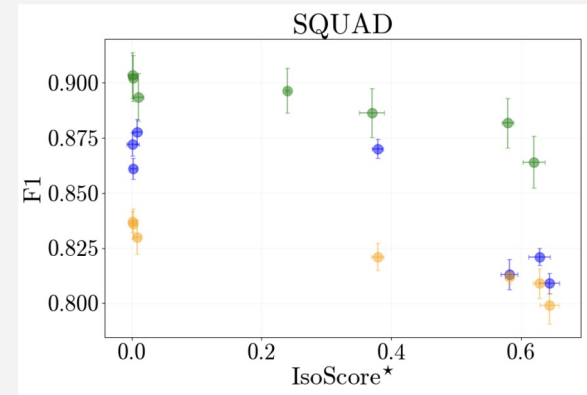
AIM 1

Develop a stable method of measuring isotropy.




AIM 2

Understand how isotropy correlates to model performance.



Thank you!

- **GitHub:** https://github.com/bcbi-edu/p_eickhoff_isoscore.git
- **Pip installation:**


```
IsoScore 2.0.1  
pip install IsoScore
```
- **Paper:** <https://arxiv.org/pdf/2305.19358.pdf>
- **Email:** william_rudman@brown.edu

