



Everything, Everywhere, All at Once: How learning order shapes the embedding space.

SLLD exam - Mod 1 & 2

Lucia Domenichelli - National PhD in Artificial Intelligence



TABLE OF CONTENTS

01

Some context

03

Our work

02

Study objectives

04

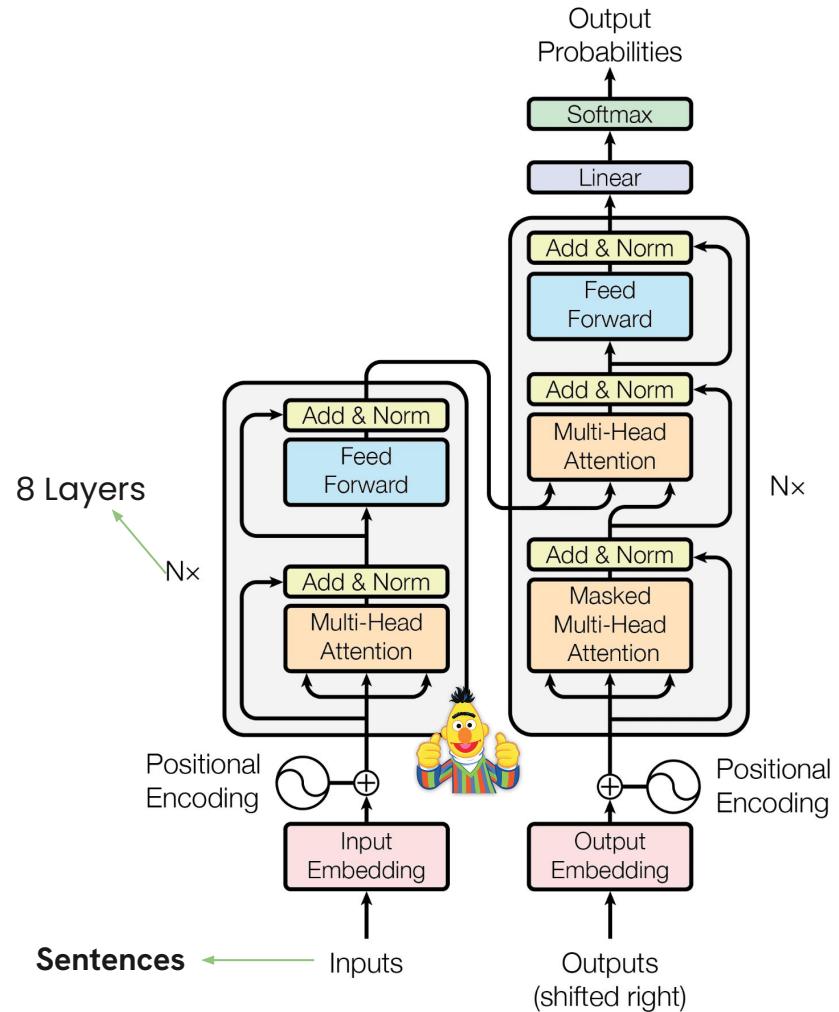
Conclusions

HOW IT WORKS?

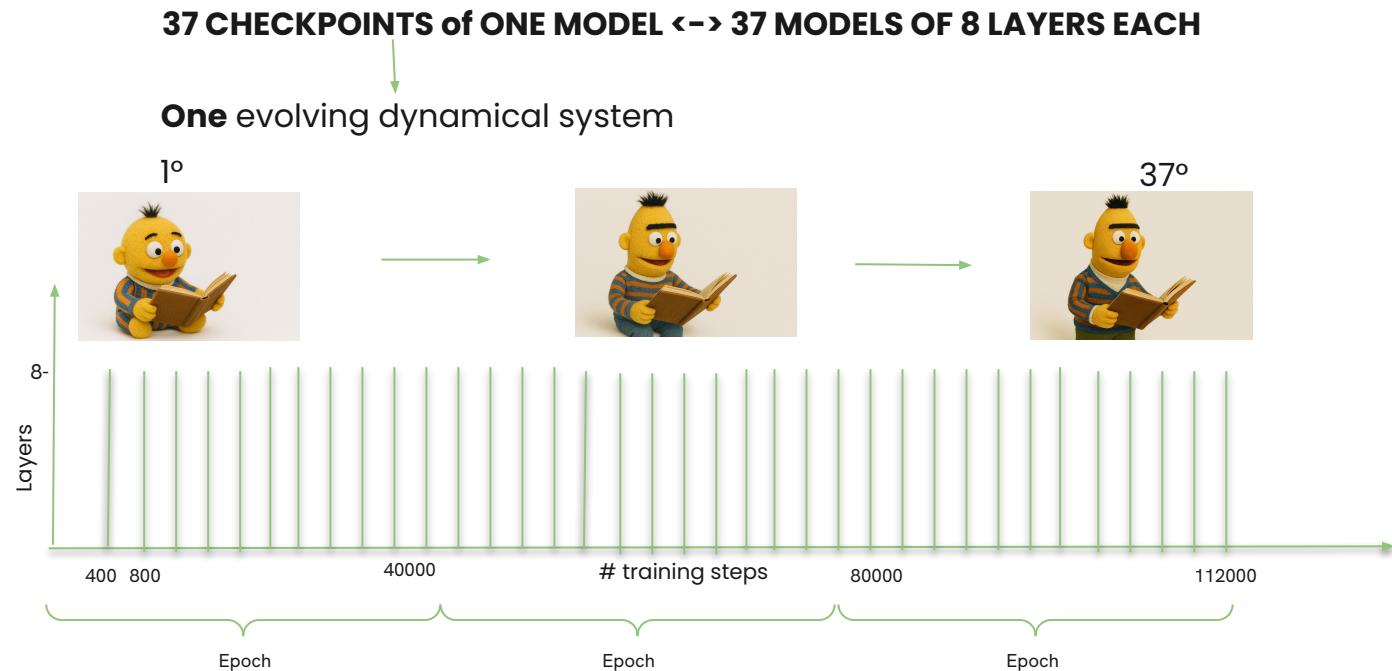
- ❑ Transformers Models have become ubiquitous in Natural Language Processing
- ❑ Pretraining + finetuning
- ❑ The task to solve in pretraining is called Language Modelling/Masked Language Modelling:

[CLS]Lucia likes to [MASK] climbing [SEP]

To predict



WE PRETRAINED BERT-MEDIUM (8 LAYERS - 512 FEATURES) ON AN ITALIAN DATASET!



02

Study

■ Objectives



CURRICULUM LEARNING

$$\text{Gulpease} = 89 + \frac{300 * (\text{numero delle frasi}) - 10 * (\text{numero delle lettere})}{\text{numero delle parole}}$$

ReadIT → makes use of linguistics annotations as well (lexical, morpho-syntactic and syntactic information.)

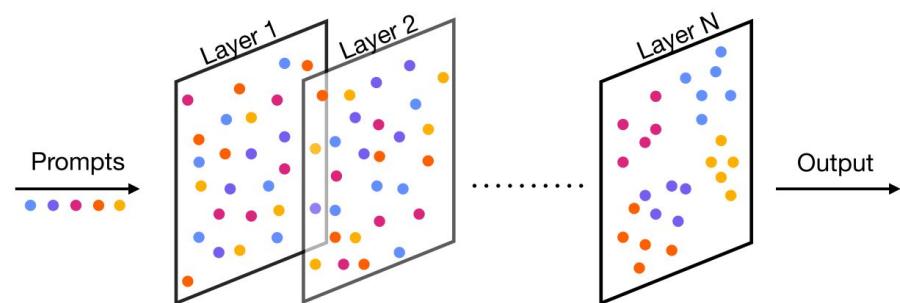
- 1) Does using a curriculum improve performance in **downstream tasks**?
- 2) Does using a curriculum **enhance linguistic** knowledge of the model?
- 3) How does the geometry of the **representation space** change when pretraining using a curriculum?

The strategies:

- Sentences' length
- Gulpease
- ReadIT
- 3 Random orderings



Linguistically Motivated



THE DATASET

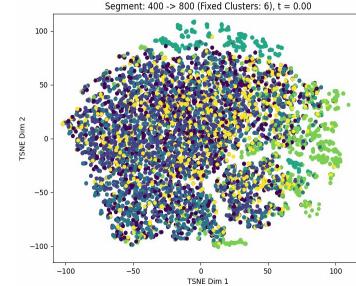
TRAINING



- 1000000 sentences longer than 6 tokens and up to 60 from half italian wikipedia.
- Avg length of a sentence 21.40 tokens
- 100000 training +20000 test samples randomly selected

PROBING

- The UD-Treebank in Italian
- 8000 train + 2000 test sentences annotated with 136 linguistic features.
- Similar **data domain** but different **topic** and distribution of tokens from the training data.



"Point cloud" of representations.

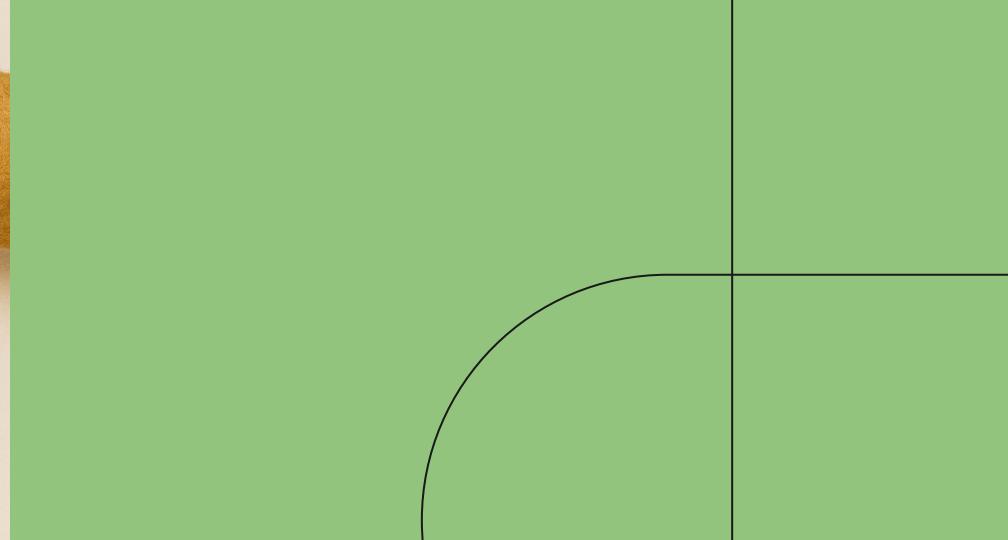
EMBEDDINGS SPACE

- We extract representations of the 20000 test sentences
- We keep the **mean pooling of tokens** of a sentence as representation for each sentence.

SIMILAR RESULTS ON THE CLS TOKEN!

03

■ Our Work



FINETUNING

Sentiment analysis
(GLUE - sst2, 2 classes)

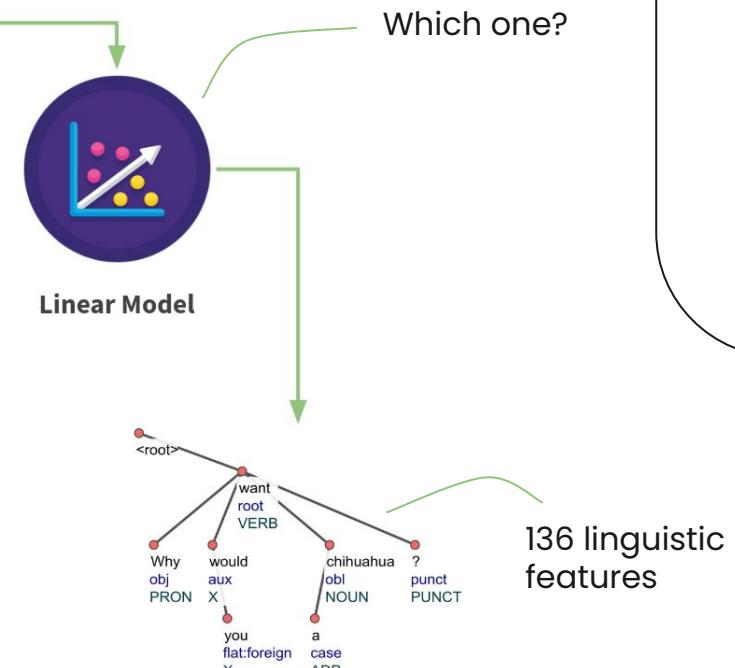
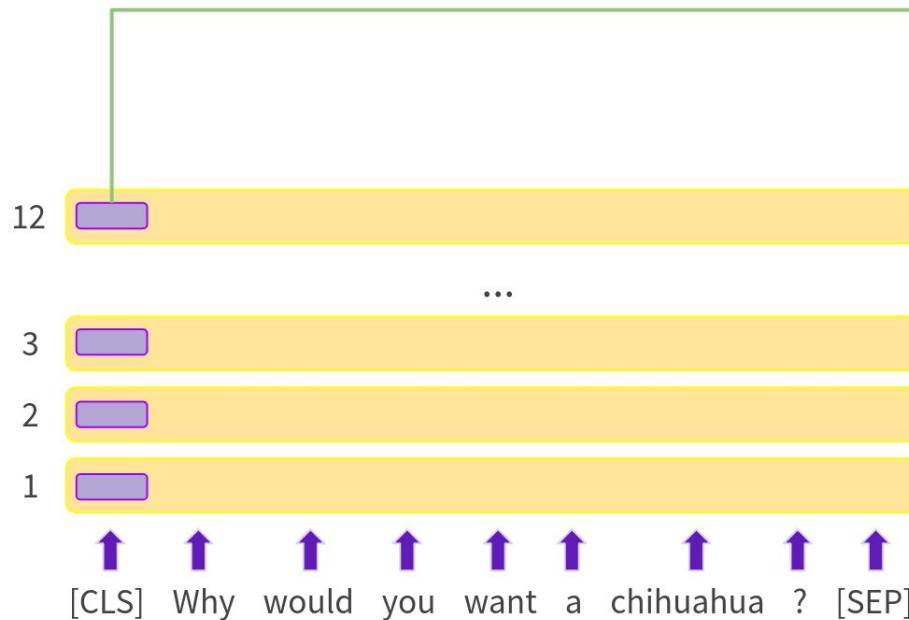
Complexity (*Human Complexity Judgment*, 7 classes)

POS tagging (235 classes)



Not so promising!

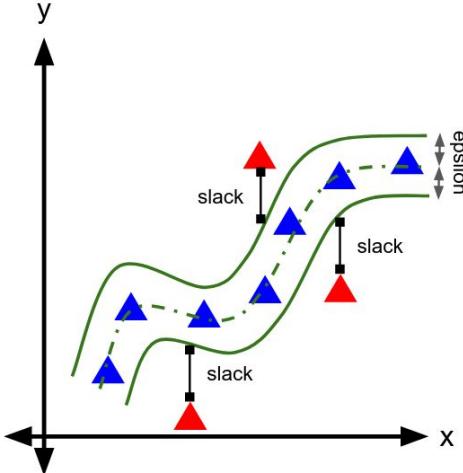
PROBING TASKS



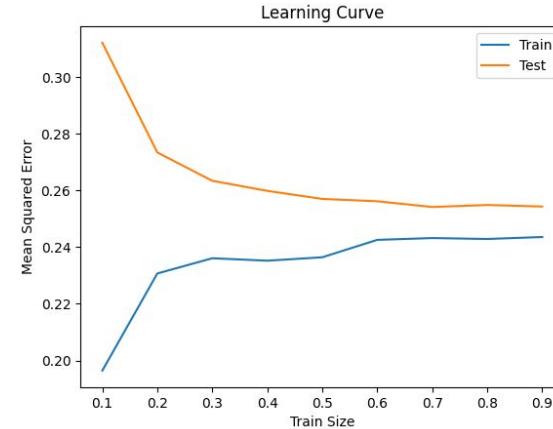
SUPPORT VECTOR MACHINE



- ❑ A simple Support Vector Regressor (LinearSVR class from scikit-learn) with Linear Kernel was employed!
- ❑ Using “simple” linear models is usually the way to go in this context
- ❑ Prior to data processing, we normalize the data using a MinMax scaler to ensure features fall within the range [0, 1] for comparability.
- ❑ 5-Fold Cross Validation has been used to mitigate estimated variance



Source: <https://medium.com/it-paragon/support-vector-machine-regression-cf65348b6345>

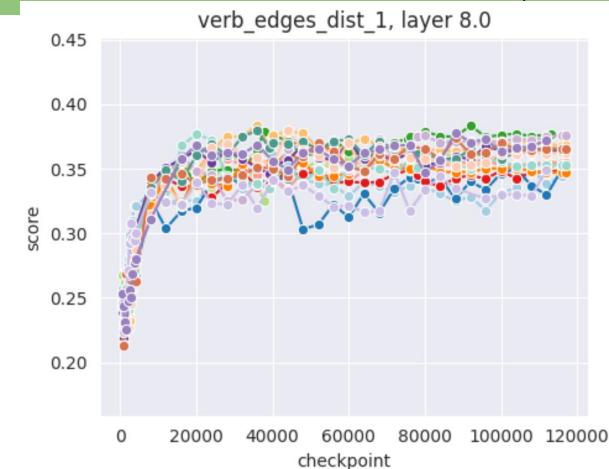
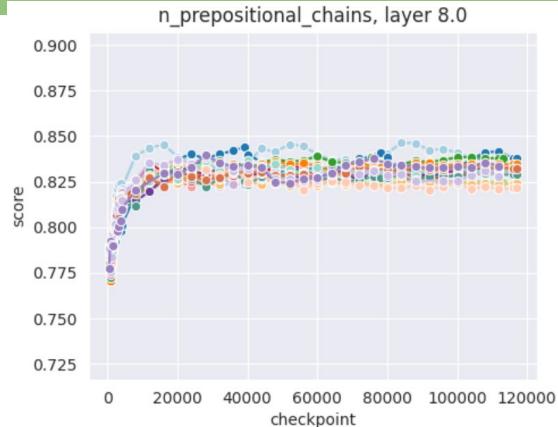


SOME RESULTS

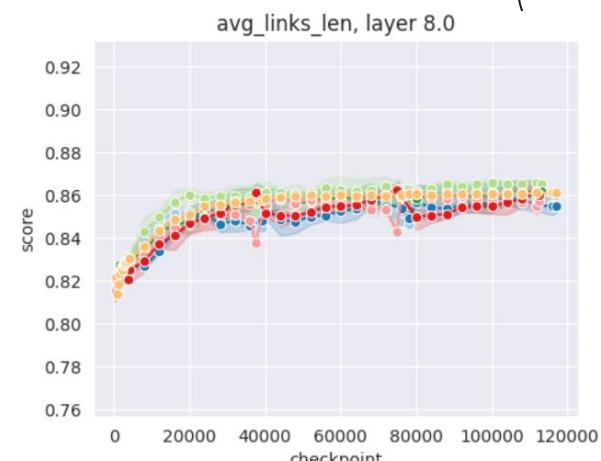
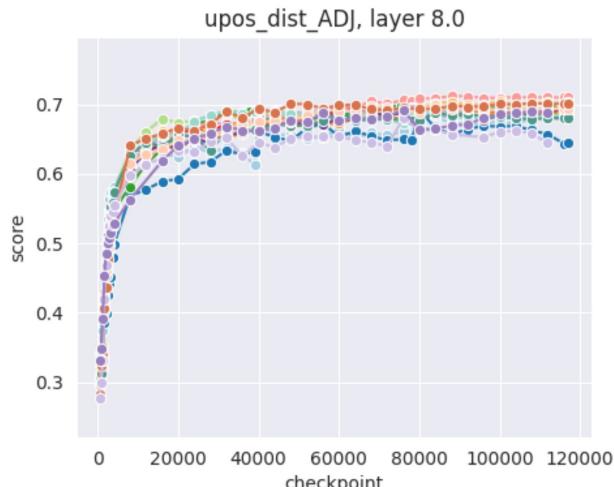
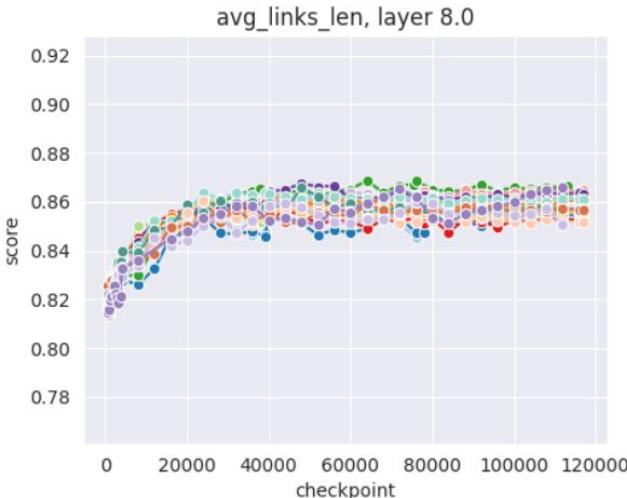


Not so much here...

- sentence_length_inverted
- sentence_length
- readit_global_inverted
- readit_global
- gulpease_inverted
- gulpease
- random

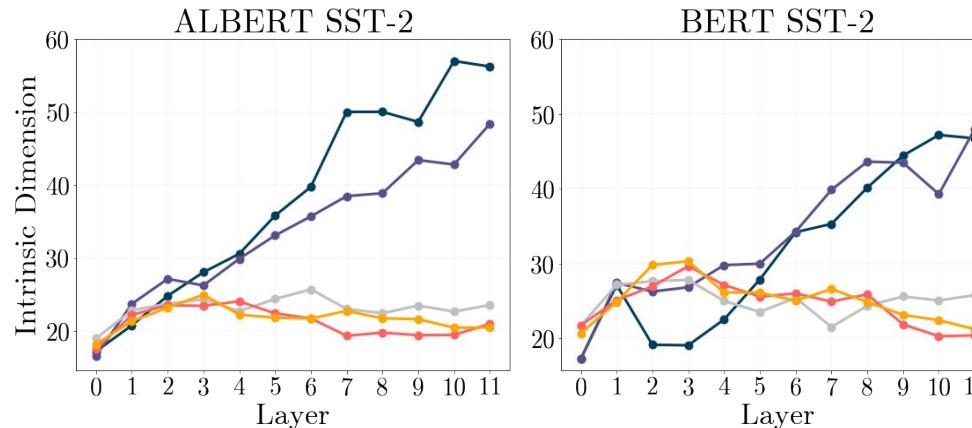


And many more!



STUDYING THE GEOMETRY of LLMs

- ❑ **Narrow Cone Hypothesis (Ethayarajh 2019)** -> Models are not isotropic, i.e., they do not uniformly utilize the embedding space.
- ❑ **Manifold Hypothesis (Clayton, 2015)** -> High Dimensional data lies on a manifold of much lower dimensionality than the number of features

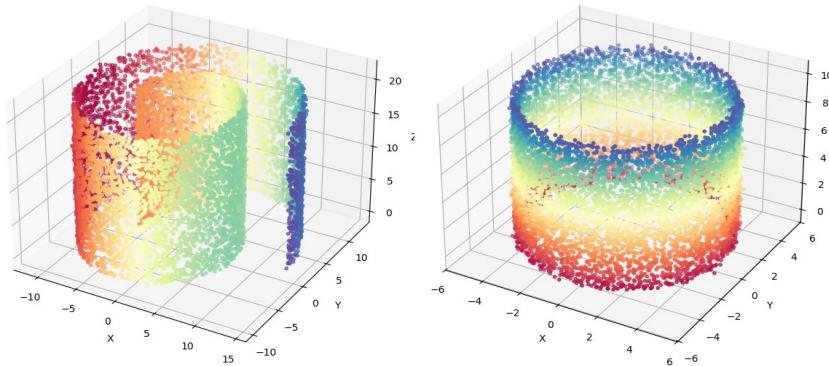


Intrinsic dimensionality estimate of ALBERT and BERT sentence embeddings from SST-2. (Rudman, 2023)

STUDYING THE GEOMETRY of LLMs

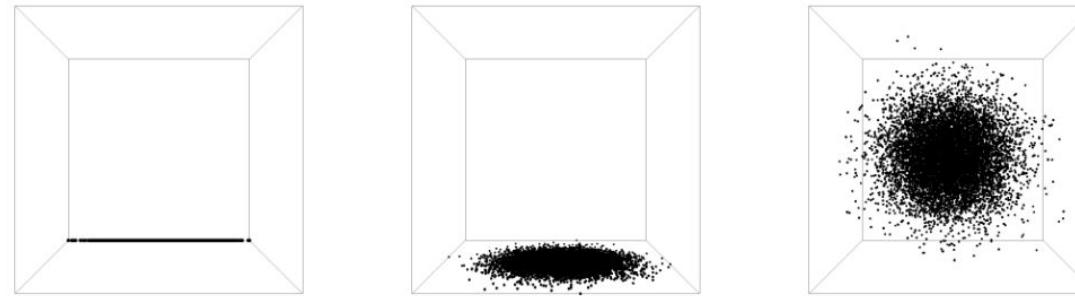
However...

- ❑ Curved and twisted hypersurfaces → Can we “iron” them?
- ❑ Complex topologies (no hyperplanes)



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).

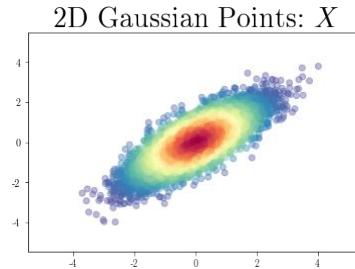


Left: line embedded in 3D space. **Middle:** circle embedded in 3D space. **Right:** Sphere in 3D space.

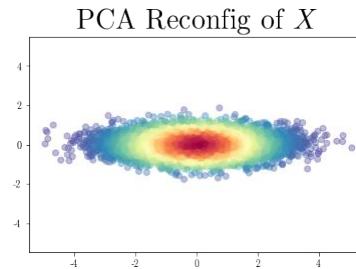
Why is it useful?

1. Low intrinsic dimensionality in later layers correlates to better model performance. (Recanatesi et al., 2019; Ansuini et al., 2019).

IsoScore



1) Point cloud X in \mathbb{R}^2 .



2) Project X using PCA to get X^{PCA} .

$$\begin{pmatrix} 1.80 & 0.00 \\ 0.00 & 0.20 \end{pmatrix}$$

The covariance matrix of the PCA-reconfigured points X^{PCA} . The bottom-right element (0.20) is circled in red.

3) Compute covariance of X^{PCA} .

$$\frac{\sqrt{2}}{\| (1.80 \quad 0.20) \|} \cdot (1.80 \quad 0.20)$$

$$\frac{\| (1.41 \quad 0.16) - (1 \quad 1) \|}{\sqrt{2(2 - \sqrt{2})}}$$

$$0.22$$

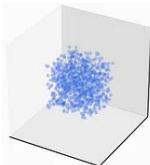
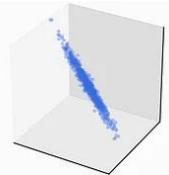
4) Normalize the diagonal of X^{PCA} to have the same norm as $(1,1)$ to get V^{PCA} .

5) Calculate the Euclidean distance between V^{PCA} and $(1,1)$ then normalize.

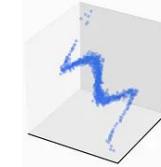
6) Linearly rescale to be in the interval $[0,1]$.

INTRINSIC DIMENSIONALITY

Linear



Nonlinear



TwoNN score

$$\Delta v_l = \omega_d(r_l^d - r_{l-1}^d),$$

$$P(\Delta v_l \in [v, v + dv]) = \rho e^{-\rho v} dv.$$

$$P(R \in [\bar{R}, \bar{R} + d\bar{R}])$$

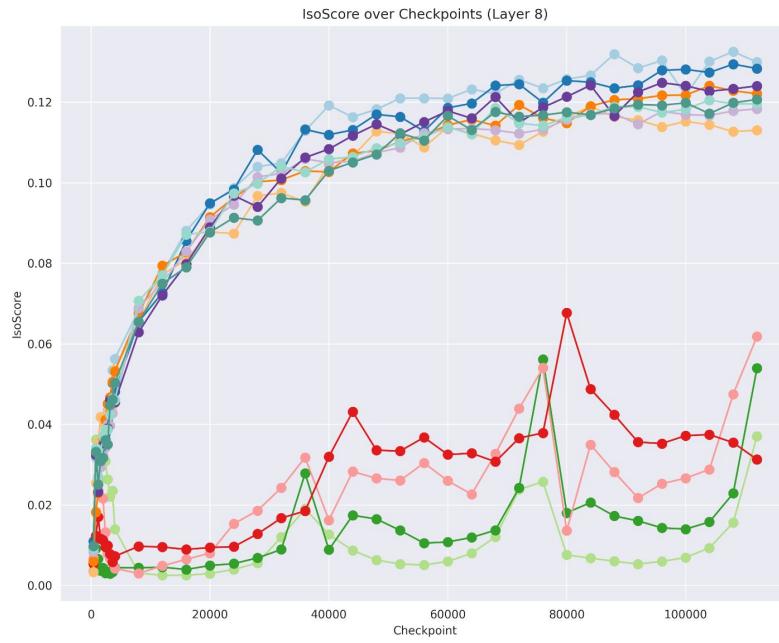
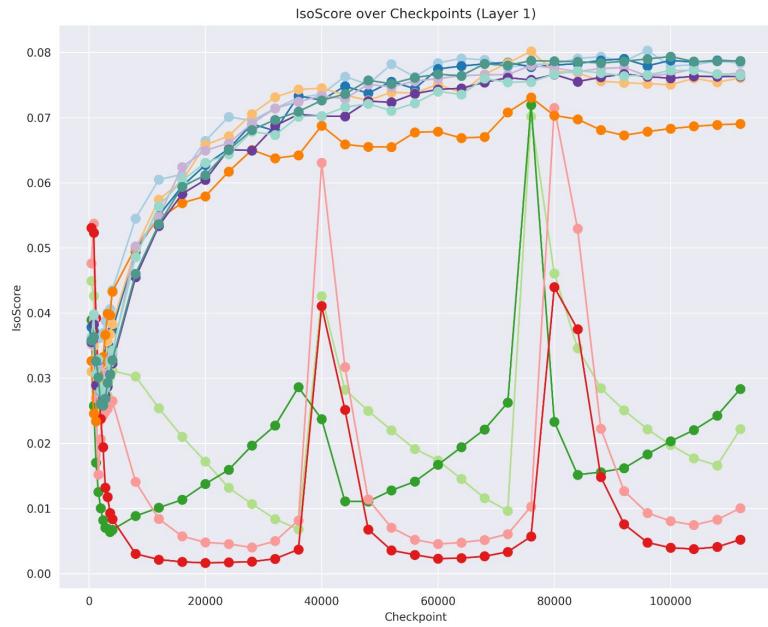
$$= \int_0^\infty dv_i \int_0^\infty dv_j \rho^2 e^{-\rho(v_i+v_j)} \mathbf{1}_{\left\{ \frac{v_j}{v_i} \in [\bar{R}, \bar{R} + d\bar{R}] \right\}}$$
$$= d\bar{R} \frac{1}{(1 + \bar{R})^2},$$

We employ a simple measure derived from PCA.

How many **principal components** are needed to retain **99% of the variance** in the data.

1. Compute the pairwise distances for each point in the dataset $i=1, \dots, N$.
2. For each point i find the two shortest distances r_1 and r_2 .
3. For each point i compute $\mu_i = \frac{r_2}{r_1}$.
4. Compute the empirical cumulative $F^{emp}(\mu)$ by sorting the values of μ in an ascending order through a permutation σ , then define $F^{emp}(\mu_{\sigma(i)}) \doteq \frac{i}{N}$.
5. Fit the points of the plane given by coordinates $\{(log(\mu_i), -log(1 - F^{emp}(\mu_i))) | i=1, \dots, N\}$ with a straight line passing through the origin.

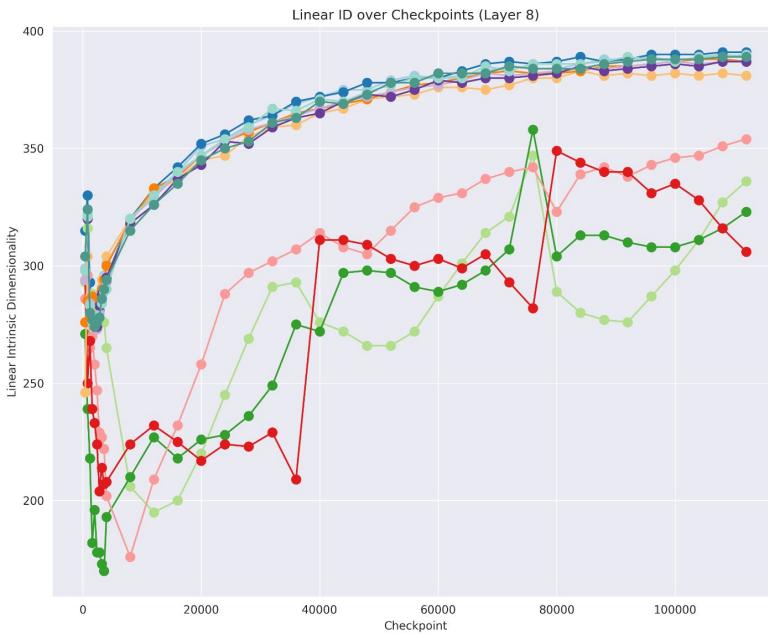
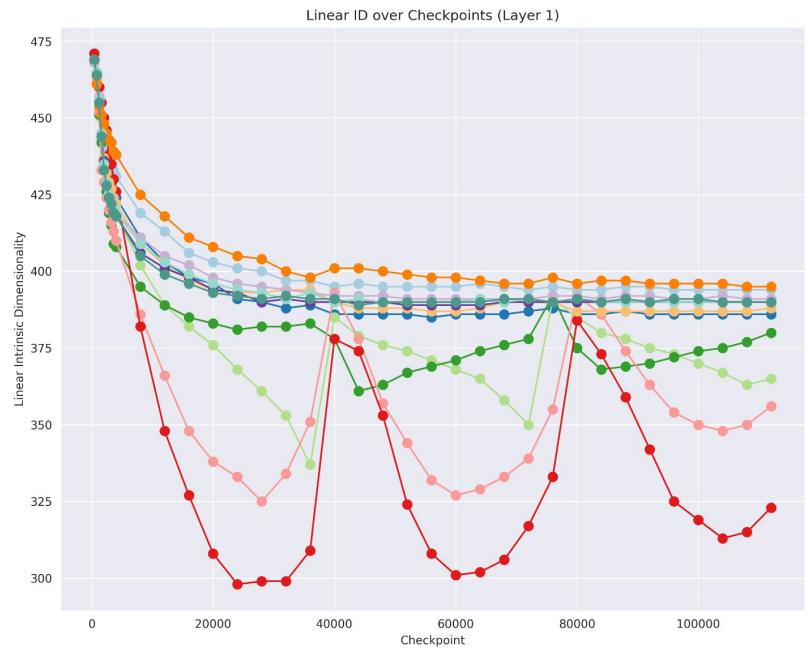
IsoScore



A legend listing the 12 model configurations used in the IsoScore charts, each associated with a colored circle and line segment:

- orig_inverted
- orig
- gulpease_inverted
- gulpease
- sentence_length_inverted
- sentence
- readit_inverted
- readit
- rand14142_inverted
- rand14142
- rand42_inverted
- rand42

Linear Intrinsic Dimensionality

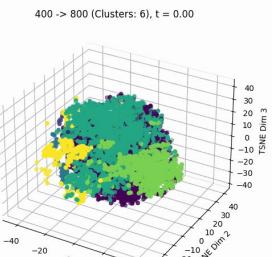


A legend listing 12 models, each associated with a colored circle and a line segment:

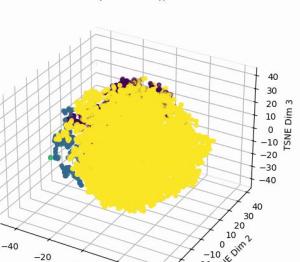
- orig_inverted
- orig
- gulpase_inverted
- gulpase
- sentence_length_inverted
- sentence
- readit_inverted
- readit
- rand14142_inverted
- rand14142
- rand42_inverted
- rand42

COHERENT CLUSTERS?

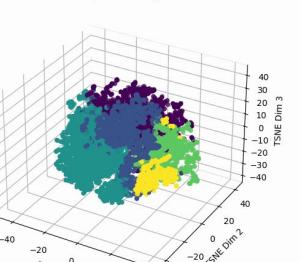
- Feature selection/reduction with PCA → **Cumulative Explained Variance** plot
- KMeans for Clustering → **Knee plot** for best k
- **Visualization** with T-SNE



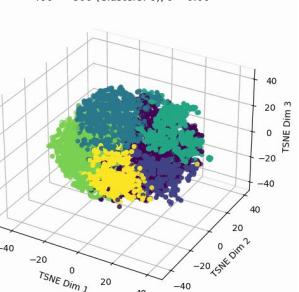
Sentence length



Gulpease



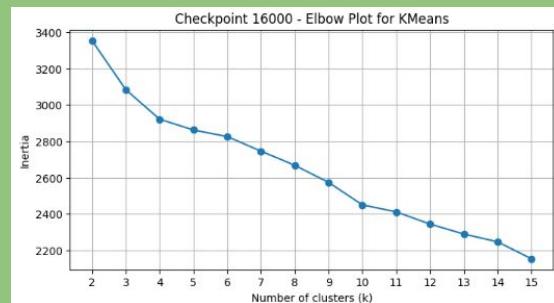
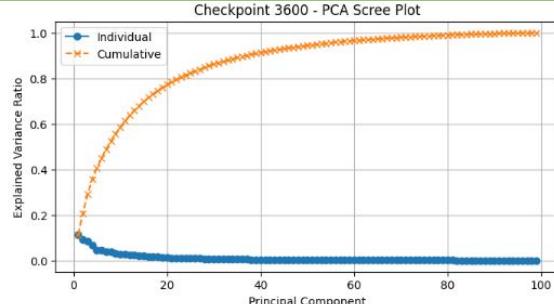
ReadIT



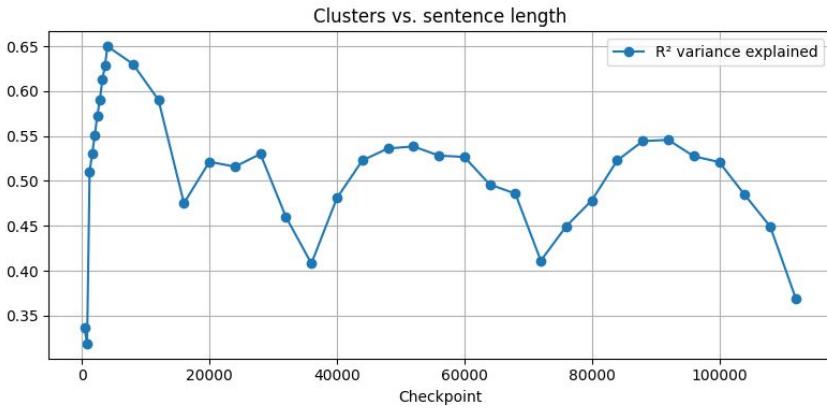
Random

Nonlinear dimensionality reduction technique, it models each high-dimensional object by a two- or three-dimensional point in such a way that similar objects are modeled by nearby points and dissimilar objects are modeled by distant points with high probability.

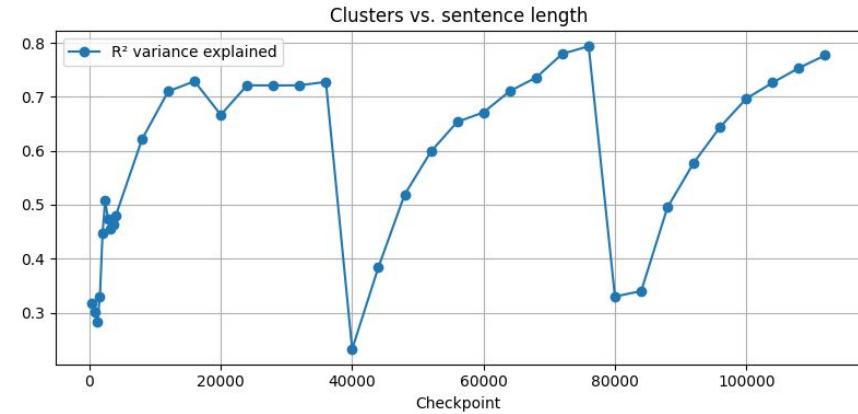
Some example plots.



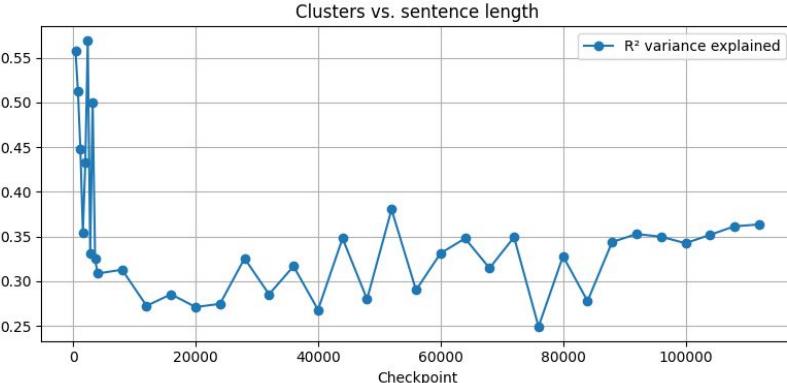
COHERENT CLUSTERS?



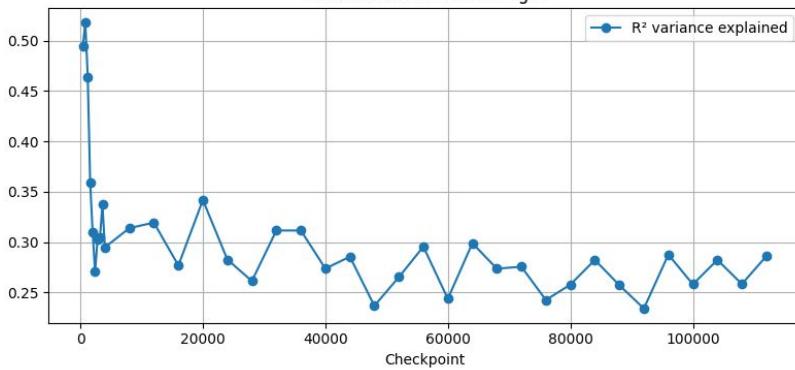
Gulpease



Sentence Length



Readit



Random



FUTURE WORK ...

EMPLOY NONLINEAR METRICS AND APPROACHES

PERFORM THE SAME ANALYSIS ON WORDS RATHER THAN SENTENCES

APPLY FUNCTIONAL ANALYSIS AND CONSIDER DEPENDENCE FROM “TIME”

DOES THIS MEAN REPRESENTATIONS CAN BE COMPRESSED ?



REFERENCES

- ❑ Curriculum Learning: A Survey – Petru Soviany, Radu Tudor Ionescu, Paolo Rota, Nicu Sebe
- ❑ Modeling Easiness for Training Transformers with Curriculum Learning – Leonardo Ranaldi, Giulia Pucci, Fabio Massimo Zanzotto
- ❑ Probing Linguistic Knowledge in Italian Neural Language Models across Language Varieties – Alessio Miaschi, Dominique Brunato, Felice Dell’Orletta, Giulia Venturi
- ❑ How Contextual are Contextualized Word Representations? Comparing the Geometry of BERT, ELMo, and GPT-2 Embeddings – Kawin Ethayarajh
- ❑ Mechanics and Geometry of Solids and Surfaces – J. D. Clayton
- ❑ IsoScore: Measuring the Uniformity of Embedding Space Utilization – William Rudman, Nate Gillman, Taylor Rayne, Carsten Eickhoff
- ❑ Stable Anisotropic Regularization – William Rudman, Carsten Eickhoff