
Shaping the representations space:

Geometric Diagnostics of Representations in Transformers

Lucia Domenichelli

Istituto di Linguistica Computazionale “Antonio Zampolli” (ILC-CNR)
ItaliaNLP Lab – www.italianlp.it

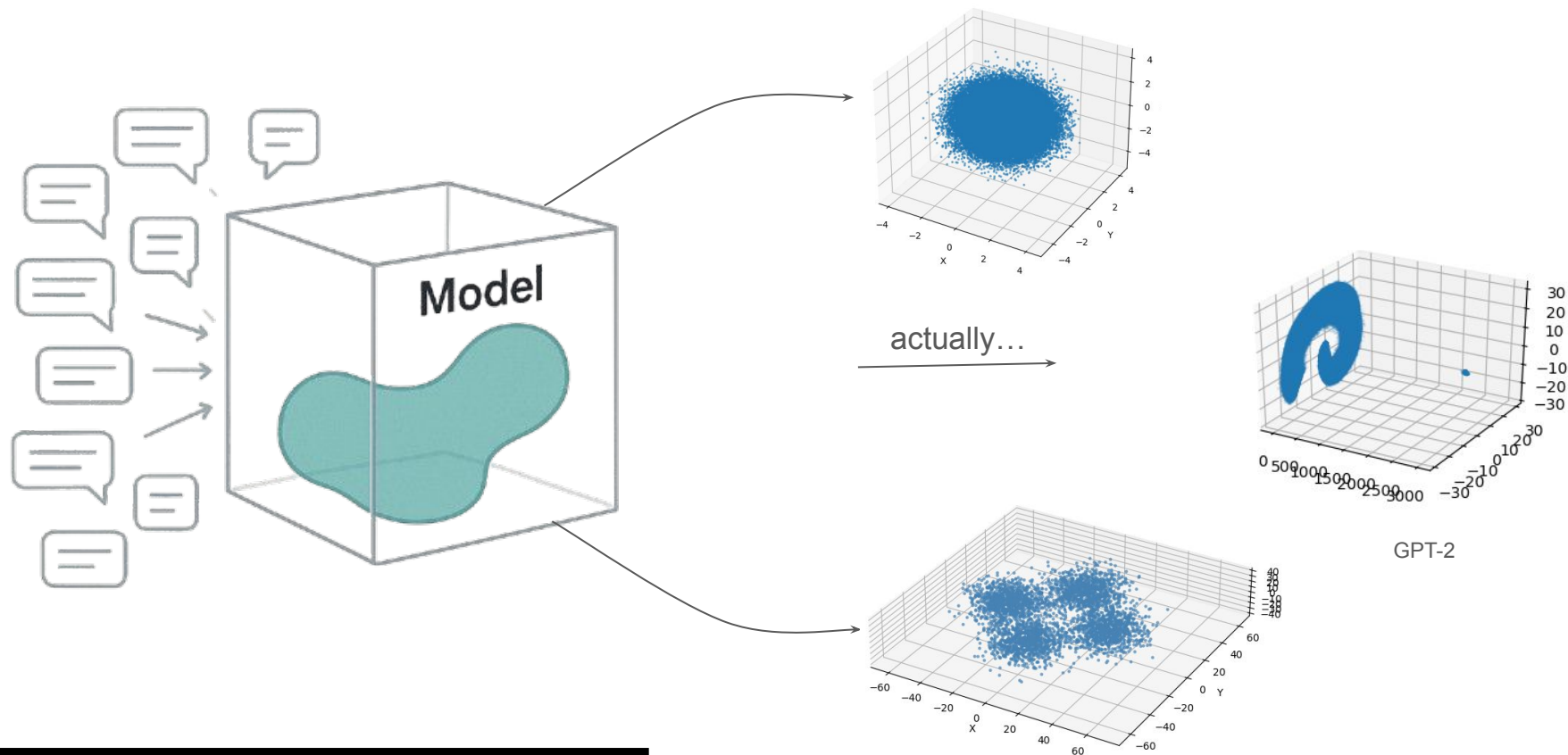


Istituto di Linguistica
Computazionale
“Antonio Zampolli”



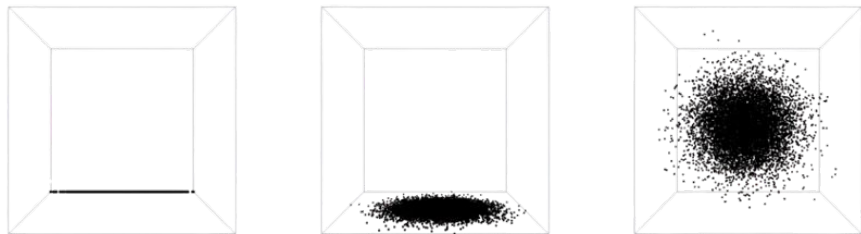
Consiglio Nazionale delle Ricerche

The idea



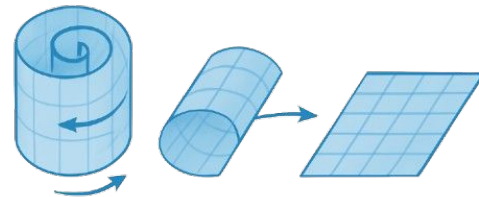
The metrics

Isotropy



A distribution is isotropic if its variance is uniformly distributed across all dimensions.

Intrinsic dimensionality



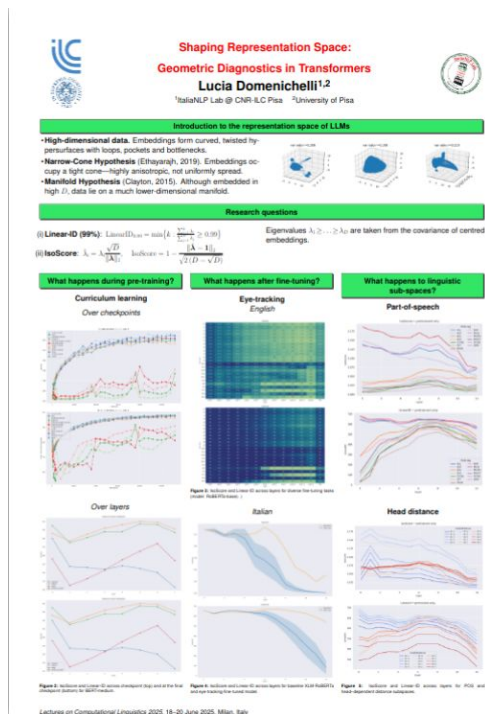
The dimension of a manifold, or *intrinsic dimensionality*, is the minimal number of degrees of freedom needed to describe it without information loss.

Research questions

Given some selected sentences and words representations:

- How do these spaces change between internal layers of Transformers?
- How do these spaces change during pretraining?
- How do these spaces change after fine-tuning?
- What about the individual subspaces of linguistic features of text?

Want to know more?



Want to share your opinion?

Please come to my poster!

