

we find that dropout acts as minimal "data augmentation" of hidden representations while removing it leads to a representation collapse

$$\ell_i = -\log \frac{e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_i^+)/\tau}}{\sum_{j=1}^N e^{\text{sim}(\mathbf{h}_i, \mathbf{h}_j^+)/\tau}}, \quad (1)$$

$$f_{\theta}(x_i) = h_i \quad f_{\theta}(x_i^+) = h_i^+$$

$f$ : pretreated language model  
BERT & ROBERTa  
 $h = f_{\theta}(x)$

Key 1  $\ell_{\text{align}} \triangleq \mathbb{E}_{(x, x^+) \sim p_{\text{pos}}} \|f(x) - f(x^+)\|^2. \quad (2)$

On the other hand, uniformity measures how well the embeddings are uniformly distributed:

$$\ell_{\text{uniform}} \triangleq \log \mathbb{E}_{x, y \sim \text{i.i.d. } p_{\text{data}}} e^{-2\|f(x) - f(y)\|^2} \quad (3)$$

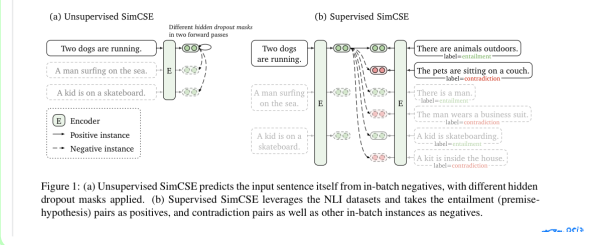


Figure 1: (a) Unsupervised SimCSE predicts the input sentence itself from in-batch negatives, with different hidden dropout masks applied. (b) Supervised SimCSE leverages the NLI datasets and takes the entailment (premise hypothesis) pairs as positives, and contradiction pairs as well as other in-batch instances as negatives.

To better understand the strong performance of SimCSE, we borrow the [analysis tool](#) from Wang and Isola (2020), which takes *alignment* between semantically-related positive pairs and *uniformity* of the whole representation space to measure the *quality of learned embeddings*. Through empirical analyses, we find that our unsupervised SimCSE consistently improves uniformity while avoiding degradation of alignment. In addition, by improving the expressiveness of the representations, the same analysis shows that the NLI training signal can further improve alignment between positive pairs and produce better sentence embeddings. We also draw a connection to the recent findings that pre-trained word embeddings suffer from [word co-occurrence bias](#) (Wang and Isola, 2020). From a structural perspective—the *contrastive learning objective* “flattens” the singular value distribution of the sentence embedding space, hence improving uniformity.

Training objective	$f_\theta$	$(f_{\theta_1}, f_{\theta_2})$
Next sentence	67.1	68.9
Next 3 sentences	67.4	68.8
Delete one word	75.9	73.1
Unsupervised SimCSE	<b>82.5</b>	80.7

↑  
Single encoder.  
 $f_0$

$$\ell_i = -\log \frac{e^{\text{sim}(\mathbf{h}_i^{z_i}, \mathbf{h}_i^{z'_i})/\tau}}{\sum_{j=1}^N e^{\text{sim}(\mathbf{h}_i^{z_i}, \mathbf{h}_j^{z'_i})/\tau}}, \quad (4)$$

○ negative pair & positive pair.  
term의 의미를 어떻게 구분하는지

$p$	0.0	0.01	0.05	0.1
STS-B	71.1	72.6	81.1	<b>82.5</b>

$p$	0.15	0.2	0.5	Fixed 0.1
STS-B	81.4	80.5	71.0	43.6

Table 3: Effects of different dropout probabilities  $p$  on the STS-B development set (Spearman’s correlation, BERT<sub>base</sub>). *Fixed 0.1*: default 0.1 dropout rate but apply the same dropout mask on both  $x_i$  and  $x_i^+$ .

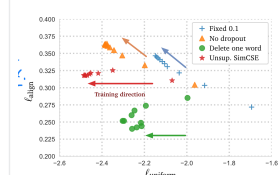


Figure 2:  $\ell_{\text{align}}$ - $\ell_{\text{uniform}}$  plot for unsupervised SimCSE, “no dropout”, “fixed 0.1”, and “delete one word”. We visualize checkpoints every 10 training steps and the arrows indicate the training direction. For both  $\ell_{\text{align}}$  and  $\ell_{\text{uniform}}$ , lower numbers are better.

## 5 Connection to Anisotropy

$$-\frac{1}{\tau} \mathbb{E}_{x, x^+ \sim p_{\text{pos}}} [f(x)^\top f(x^+)] + \mathbb{E}_{x \sim p_{\text{data}}} \left[ \log \mathbb{E}_{x^- \sim p_{\text{data}}} [e^{f(x)^\top f(x^-)/\tau}] \right], \quad (6)$$

$$\sum_{i=1}^n Q_i = - \sum_{i=1}^n \log e^{\sin(h_i, h_i^*)/2} + \sum_{i=1}^n \sum_{j=1}^N e^{\sin(h_i, h_j^*)/2}$$

$$= -\frac{1}{2} \sum_{(x_i, z_i) \sim p_{\text{data}}} \left[ f(x_i; z_i)^T f(x_i^*; z_i) \right]$$

$$+ E_{x \sim p_{data}} \left[ \log E_{x \sim p_{data}} \left[ e^{f(x)^T f(x_i)/2} \right] \right]$$

$p_{data}$  is uniform over finite samples  $\{x_i\}_{i=1}^m$ , with  $\mathbf{h}_i = f(x_i)$ , we can derive the following formula from the second term with Jensen's inequality:

non-convex  
omg.  
2nd 2nd

$$\begin{aligned} & \mathbb{E}_{x \sim p_{\text{data}}} \left[ \log \mathbb{E}_{x \sim p_{\text{data}}} \left[ e^{f(x)^T (f(x') - f)} \right] \right] \\ &= \frac{1}{m} \sum_{i=1}^m \log \left( \frac{1}{m} \sum_{j=1}^m e^{h_i^T h_j / r} \right) \\ &\geq \frac{1}{2m^2} \sum_i \sum_j h_i^T h_j. \end{aligned} \quad (7)$$

$W$ : sentence embedding matrix.

$i$ -th row of  $W$

$$\text{sum}(WW^T) = \sum_{i=1}^m \sum_{j=1}^m h_i^T h_j$$

$$\forall \lambda, |h_\lambda| = 1 \Rightarrow \text{diag}(W W^T) = \mathbb{1}$$

$WW^T$ : Symmetric matrix.

$$p(t) = \det(A - tI) = 0. \text{ if, } \underline{\text{gee eigenwerte}}$$

$$= c(t - \lambda_1) \cdots (t - \lambda_n).$$

$$= \left( C^n \dots \sum_{i=1}^n \lambda_i - \frac{n}{\lambda_1} \lambda_i \right)$$

$$\frac{\partial \det(A)}{\partial t} = \alpha \sum_{i=1}^n \lambda_i$$

$$A = P D P^{-1}$$
$$\text{tr}(A) = \text{tr}(D) = \text{eigenvalues}$$

Symmetric  $\Rightarrow$  1. real values eigenvalue.  
2.  $x_i, x_j = \delta_{ij}$   
3. diagonalizable.

$$[w^w]_{ij} > 0.$$

## F Distribution of Singular Values

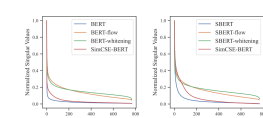


Figure F.1: Singular value distributions of sentence embedding matrix from sentences in STS-B. We normalize the singular values so that the largest one is 1.

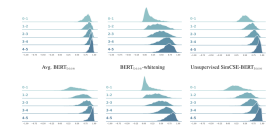


Figure G.1: Density plots of cosine similarities between sentence pairs in STS-B. Pairs are divided into 5 groups based on ground truth labels (higher means more similar along the x-axis, and  $\lambda$  axis is the cosine similarity).