

Self-supervised learning (SSL)

Self-supervision

$$\{(x_i, y_i)\}_{i=1}^L$$

$$\underbrace{\frac{1}{L} \sum_{i=1}^L L(f(x_i), y_i)}_{\text{supervised}} \rightarrow \min_f$$

Pre-train with out manual labels

- Embeddings
- Language models

GPT (Generative Pre-training of Transformers)

Architecture - transformer decoder
Task - language modeling

BERT (Bidirectional Encoder Representations from Transformers)

Architecture - transformer encoder

Task - Masked Language Modeling

2 | бозбун | 710 | урегнотекне | как | рунер

↓ mask

2 | [MASK] | 710 | урегнотекне | как | [MASK]

distribution
over tokens

Transformer
encoder

2 | отберу | 710 | урегнотекне | как | берга

- loss only over masked positions

All tokens

15%
masked

85%

unchanged

80%

[MASK]

10%

random
token

10%

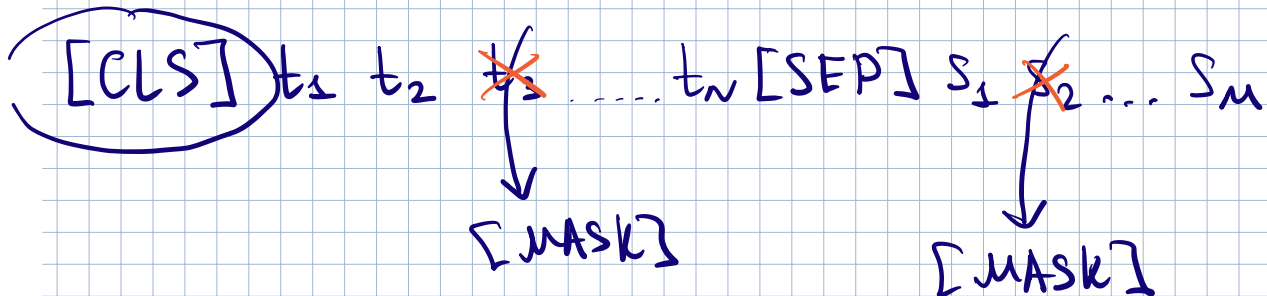
unchanged

calculate
loss

$$\text{BERT loss} = L_{\text{MLM}} + L_{\text{NSP}}$$

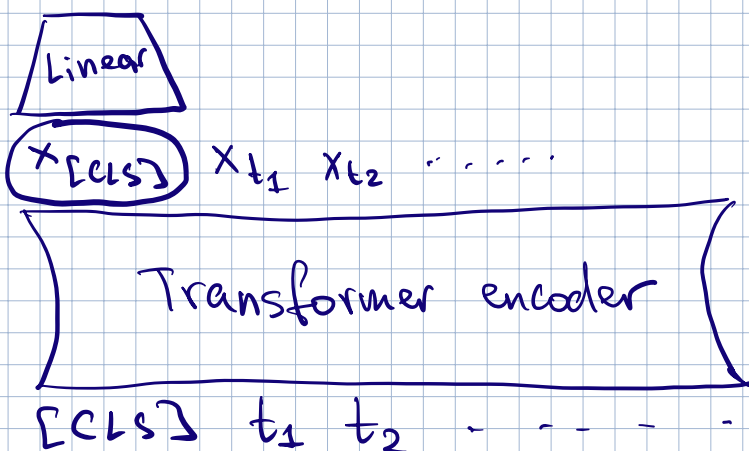
NSP - Next Sentence Prediction

Model input



$$\begin{array}{ccccccccccc} E_{[CLS]} & E_{t_1} & E_{t_2} & \dots & E_{t_n} & E_{[SEP]} & E_{s_1} & \dots & E_{s_m} \\ E_0 & E_1 & E_2 & \dots & E_N & E_{N+1} & E_{N+2} & \dots & E_{N+m+1} \\ E_A & E_A & E_A & \dots & E_A & E_A & E_B & \dots & E_B \end{array}$$

Linear probing



GLUE (General Language Understanding Evaluation)

- 9 NLP tasks, classification
 - CoLA - Corpus of Language Acceptability
 - MNLI - Multi-genre Natural Language Inference

• NLP

ROBERTA
ALBERT

• Images

BEiT

• Audio

HUBERT

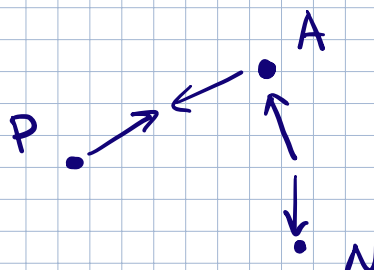
SSL for images

Pre-text tasks:

- Rotation angle prediction
- Jigsaw puzzle
- Image colorization

Triplet loss

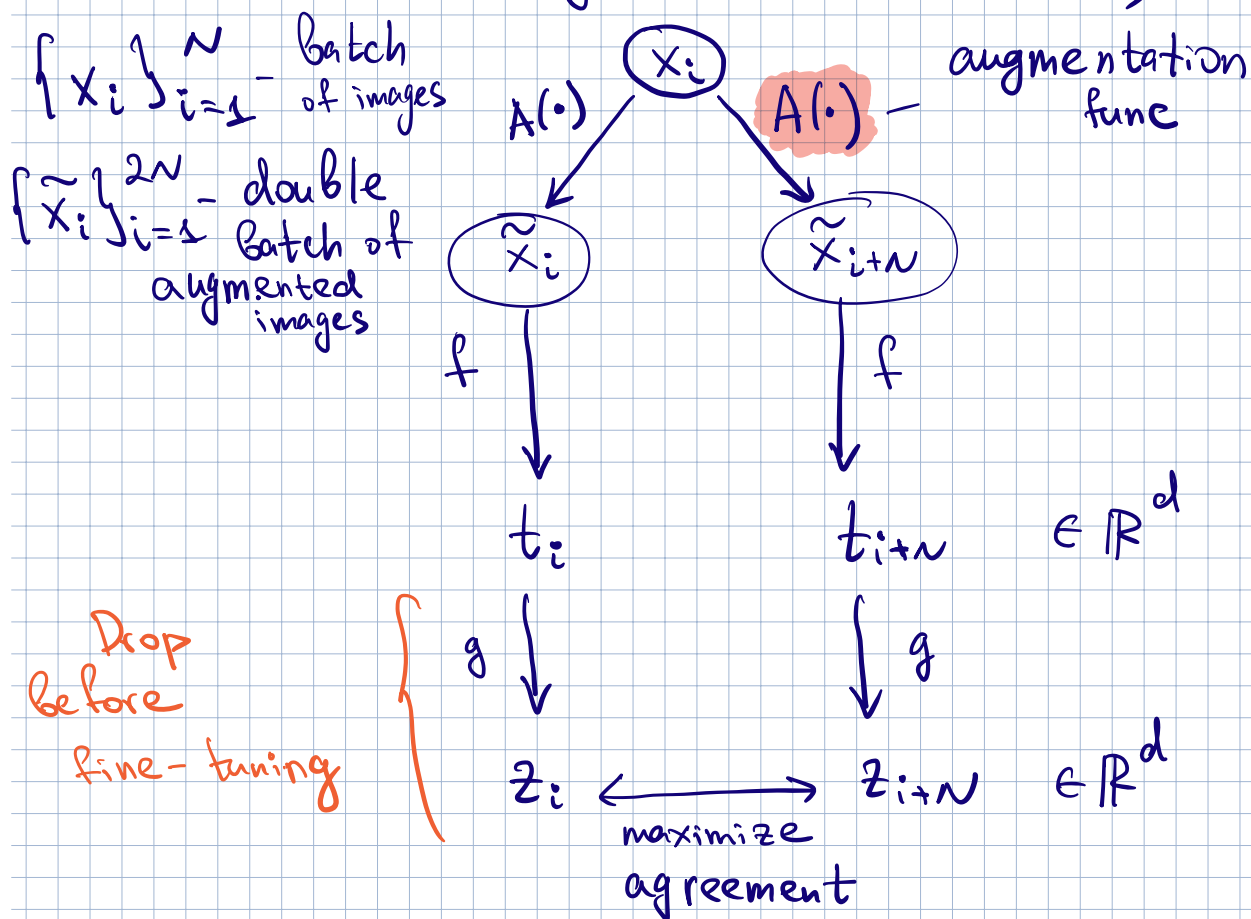
$\|f(x) - f(y)\|_2^2$ - semantic distance



$$\max(\|f(A) - f(P)\|_2^2 - \|f(A) - f(N)\|_2^2 + m, 0) \Rightarrow \min_f$$

 $m > 0$

SimCLR (A Simple framework for Contrastive Learning of Representations)



Contrastive loss

$$l(z_i, z_{i+N}) = -\log \frac{e^{\text{sim}(z_i, z_{i+N})}}{\sum_{k=1}^{2N} \mathbb{1}_{\{k \neq i\}} e^{\text{sim}(z_i, z_k)}}$$

$$\text{sim}(x, y) = \frac{\langle x, y \rangle}{\|x\|_2 \cdot \|y\|_2}$$

$$L(z_i, z_{i+N}) = \frac{1}{2} (l(z_i, z_{i+N}) + l(z_{i+N}, z_i))$$

- Random Crop
- Color Jitter
- Convert to Grayscale
- Gaussian Blur