Fine-tuning Embedding Models

Where and how to fine-tune models for dense retrieval



James Briggs

Staff Developer Advocate Pinecone (we're hiring)

Educator 🔽 🌒 Û

Director Aurelio Consulting **QLO**

Freelance ML Engineering and Advisory

Other

Code: https://github.com/jamescalam/train-sentence-transformers

Ebook: https://www.pinecone.io/learn/nlp/







Haystack EU 2022

There are 1000s of *pretrained* embedding models



Why fine-tune anything?



uniqueness of domain

3 days ago

if I wanted to do TSDAE and GPL together is there a specific order? Like pretrained using TSDAE first then fine tune with GPL

· 3 weeks ago

Hi James! You mentioned a project you did involving the Dhivehi and using TSDAE. Given it's hard to find a doc2query and cross-encoder model in Dhivehi, could you have used GPL as well?

1 month ago

Hi James, I'm looking to finetune a transformer model on my custom dataset for which I'm annotating documents by myself, My doubt is how many documents do I need to annotate for finetuning?

· 3 months ago

I like that approach with the 3-4 sentences. I made something similar but I grouped one minute worth of audio together. Overlapping sentences sounds like it can solve some of the issues I have seen. I actually have many questions: Hope that's okay

- How does this way differ from using Haystack QA and pipeline with Pinecone? Is this more accurate?

- If I wanted to train the QA model to better answer questions regarding a specific topic, say only include psychology videos TTS and then train the model to better understand psychological terms, how would I do that?

4 months ago

how to train sbert with a specific domain?

6 months ago

Awesome video James, how can we fine tune a model on our custom data, is there any way to train the model with out custom data. One more question which is a better approach this one or facebook/DPR model?

🖯 • 10 months ago

Hello, i have one question : is it possible to fine tune sbert with unsupervised data ? I have 2 job catalogs that i need to match by textual similarity but i dont have labels (for each job I only have it's title and description)

• 1 year ago

Great series! This is exactly what i have been looking for the project im working on. Thanks for making this. It would really great if you could share more on how to finetune the three models for new languages or share some codes we can explore. Thanks again.

09/21/2022

I hope to achieve good results, it is a challenge to train a model in Spanish, I don't know any dataset specialised in legal-QA, so I will make a fine tuning to a legal dataset in my language.

08/31/2022

On the github, there is literally the code to do semantic similarity for his top performaing models

https://github.com/Muennighoff/sgpt#asymmetricsemantic-search

There are 4 types tho and I don't understand what the difference is.

Cross-Encoder vs BiEncoder with Symmetric vs

Asymmetric semantic search

I don't know which one to use for the best performance.



09/15/2022

I've trained a model for sentence pair classification, but it seems unable to manage misspelt words, even if they are similar. Still, there is no such problem in sentence transformer its cosine similarity; it can recognise it with a higher score. Any ideas?



Hello 🙂

@jamesbriggs do you have any examples of pretrain using mlm or tsdae on allmpnet models



08/31/2022

Has any one performed any form of evaluation of haystack? (How else can we improve it right?) I mean to check score and improve. I am trying to implement their evaluation system, but cant really figure it out.



07/14/2022

Hello and thank you for all the help you are providing! I am working on semantic search for my internship/thesis and for the evaluation of different models, I have gathered around 100 similar sentencepairs from the company data and I have been getting the cosine similarity for each pair with different models (sbert/use/elmo/simCSE etc. and non-contextual algorithms as well). Does it make sense to evaluate them this way, looking for example at the mean of the similarity scores or should I be looking into something else? I also want to try and finetune models on this list. The domain is healthtech (a mixture of medical/biomedical and digital health) and it's usergenerated text (so the sentences are not purely technical). I would really appreciate any ideas you may have!! Thank you!!

07/21/2022



I have a question: When trying to fine-tune a pretrained BERT model to a new domain (Domain adaptation) should I perform the MLM task on few epochs or should I also apply the NSP task?

> @jamesbriggs mpnet is generally a good choice though for ... 08/16/2022

Yeah, using mpnet only... Don't think model has any issue... Might be fine-tuning would help here

Do you have any references for tuning mpnet?

09/11/2022

01/10/2022

hello. I have a QA model and want to do some reranking on answers. I produced multiple answers just like the qa pipeline in huggingface.

To re-rank them i'm thinking to use sentence bert with MSE loss on their F1 scores where the input is (question, answer_i), does this approach make sense?

We can all benefit from some fine-tuning



Comparing items

Fine-tuning embedding models revolves around **contrast**

A is similar to B

B is not similar to C

etc

Labeled text-pairs

Fine-tuning with labels

Labeled text-pairs

- + Good performance
- Hard to find data

sentence1	sentence2	label
"a plane is taking off"	"an air plane is taking off"	5
"a man is playing a large flute"	"a man is playing a flute"	3.8
"a man is spreading shreded"	"a man is spreading shredding"	3.8
"three men are playing chess"	"two men are playing chess"	2.6

GLUE Sentence Textual Similarity benchmark (STSb)

Fine-tuning with labels

It's hard to find labeled data



Unlabeled text-pairs

Unlabeled text pairs

Most datasets don't contain specific (or any) labels

premise	hypothesis	label
"choir sings to church audience"	"church with cracks in ceiling"	1 (neutral)
"choir sings to church audience"	"church is filled with song"	O (entailment)
"choir sings to church audience"	"choir sings at baseball game"	2 (contradiction)
"woman with a green scarf"	"the woman is young"	1 (neutral)

Stanford Natural Language Inference (SNLI)

All we need are (anchor, positive) pairs

premise	hypothesis	label
"choir sings to church audience"	"church is filled with song"	O (entailment)
"woman with a big grin"	"the woman is very happy"	O (entailment)
"old man poses in front of advert"	"man poses in front of ad"	O (entailment)

Stanford Natural Language Inference (SNLI)

M. Henderson, et. al., Efficient Natural Language Response Suggestion for Smart Reply (2017), Google

N. Reimers, Losses - MultipleNegativesRanking, Sentence Transformers Docs

We mix-and-match to get negatives

$$A_0P_0 = E$$
"choir sings to church audience", "church is filled with song"]
 $A_1P_1 = E$ "woman with a big grin", "the woman is very happy"]
 $A_2P_2 = E$ "old man poses in front of advert", "man poses in front of ad"]

$$\begin{array}{ccc}
\mathcal{A}_{0}P_{0} & \mathcal{A}_{1}P_{0} & \mathcal{A}_{2}P_{0} \\
\text{Mix different pairs o create negatives} & \mathcal{A}_{0}P_{1} & \mathcal{A}_{1}P_{1} & \mathcal{A}_{2}P_{1} \\
\mathcal{A}_{0}P_{2} & \mathcal{A}_{1}P_{2} & \mathcal{A}_{2}P_{2}
\end{array}$$

We mix-and-match to get negatives



Source: <u>pinecone.io/learn/geng/</u> pinecone.io/learn/fine-tune-sentence-transformers-mnr/

Then optimize



Then optimize



This works, but could be better

Randomly chosen pairs are often easy to distinguish

It's like asking the model to spot the difference



Ideally, we want to challenge the model



Haystack EU 2022

Hard-negatives are items that the model struggles to distinguish as negative



(Hard-negatives are not MNR-specific)

Finding hard-negatives

- 1. In the dataset, if we're lucky
- 2. Hard-negative mining





Haystack EU 2022



Haystack EU 2022

How many pairs are needed?



Low resource scenarios



Haystack EU 2022

Low-resource == very little data



N. Reimers, I. Gurevych, <u>Making Monolingual Sentence Embeddings</u> <u>Multilingual using Knowledge Distillation</u> (2020), EMNLP



K. Wang, et al., <u>TSDAE: Using Transformer-based Sequential</u> Denoising Auto-Encoder for Unsupervised Sentence Embedding Learning (2021), EMNLP

Translation pairs

Translation pairs

"Parallel data" - text-pairs in a **source** language and a **target** language

- + It's easy to find parallel data
- Requires existing use-case model in source language (teacher)
- Requires pretrained model for source and target language (student)

With translation pairs we can use:

Multilingual knowledge distillation

Multilingual Knowledge Distillation



Training Overview

Teacher model performance on **source** set == perf ceiling



Student should approach perf ceiling set by teacher on **target** set

... for several epochs

Only have text?

No pairs, just text?

What if we don't have any text pairs?

no pairs, just text? TSDAE

K. Wang, et al., <u>TSDAE: Using Transformer-based Sequential</u> <u>Denoising Auto-Encoder for Unsupervised Sentence Embedding</u> <u>Learning</u> (2021), EMNLP

TSDAE Similar to MLM

Masked-Language Modeling (MLM) is the pretraining approach of models like BERT



But at the sentence-level



TSDAE

- + Easy to fine-tune
- + All we need is unstructured text data
- Only works for vanilla similarity (i.e. no QA)
- Performance cannot compare to previous methods

Data Augmentation

AugSBERT Synthetic dataset augmentation?

If we have some data, or are willing to annotate a small 1-5K dataset

Data augmentation with transformers

Can we generate more pairs, and label them?

More pairs == random mix-and-match

Synthetic similarity scores?

Cross-encoders

If we have a small 1-5K dataset of labeled pairs, we can fine-tune a cross-encoder

Process outline



pinecone.io/learn/data-augmentation/

Considerations

- + Can start with small 1-5K dataset
- + Performance is reasonable
- Initial dataset must be labeled with similarity scores
- Must also fine-tune cross-encoder
- More complicated process

Asymmetric Search

For asymmetric search

Search can be symmetric or asymmetric



Asymmetric search, source: pinecone.io/learn/genq

Symmetric: query texts and passages texts are expected to be the same (size, form, etc)

Asymmetric: query texts are not equal size/form to passages texts, for example QA



Question-answering



GenQ and GPL **Training methods**

MNR is good when we have many query-context pairs

But with low-resource domains?

- TSDAE cannot (only for vanilla similarity)
- AugSBERT format with scores not ideal

Query gen models

There are seq2seq transformer models

- + Given a "context" or passage, they can generate queries
- Needs pretrained seq2seq model or fine-tune on *(query, context)* pairs

Query generation

Passage

"Python is an interpreted, high-level and general-purpose programming language. Python's design philosophy emphasizes code readability with its notable use of significant whitespace. Its language constructs and object-oriented approach aim to help programmers write clear, logical code for small and large-scale projects."



source: /learn/geng





Find contexts/passages from target domain





3 Train on (query, passage) pairs with MNR loss



GenQ

- + Only needs plain text data
- + Good performance if quality data
- Synthetic query generation are prone to poor quality
- Needs pretrained query gen model
- Whole process relies on high quality queries

Can we avoid bad performance from poor queries?

Generative Pseudo Labeling (GPL)



Haystack EU 2022

GPL and negative mining



GPL and pseudo labeling



GPL and other techniques



Haystack EU 2022

source: pinecone.io/learn/gpl

How to fine-tune?

Pairs and labels? Cosine similarity loss

Just pairs? MNR

Not much of anything but translation pairs? Multilingual knowledge distillation

Small set of pairs and labels? AugSBERT

Asymmetric use-case, small dataset? GenQ

GenQ performing badly? GPL

Any Questions?

James

- YouTube.com/c/JamesBriggs
- discord.qq/c5QtDB9RAP
 - twitter.com/jamescalam

Pinecone



 $\left\{ \begin{array}{c} & & \\ & & \\ & & \end{array} \right\}$ Vector database for millions or billions of records, free and paid services available



Semantic Search ebook:

pinecone.io/learn/nlp



Careers: pinecone.io/careers



twitter.com/pinecone