

Model2Vec加速RAG：模型小15倍，速度快500倍：

原创 南七无名氏 PyTorch研习社 2025年01月28日 11:01 安徽

在机器学习的世界里，嵌入（Embedding）是一个基础且关键的技术，广泛应用于自然语言处理（NLP）、搜索引擎、推荐系统等多个领域。然而，尽管嵌入技术已经取得了显著进展，但传统的嵌入方法依然面临着模型庞大、计算资源消耗大、推理速度慢等问题。

那么，如何才能突破这些限制，提高嵌入技术的效率和性能呢？

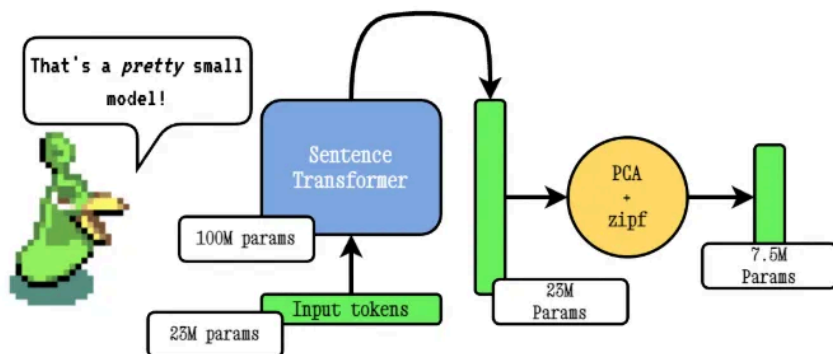
今天，我们要为大家介绍一个新的技术突破——Model2Vec，这款嵌入技术的“新宠”正通过其革命性的设计，使得嵌入模型的规模缩小 15 倍，速度提升 500 倍，同时还能保持优异的性能表现，堪称嵌入技术的“强力增强版”！



The Fastest State-of-the-Art Static Embeddings in the World

🤖 Models | 📖 Tutorials | 🌐 Website | 🏆 Results

pypi package v0.3.7 python 3.9 | 3.10 | 3.11 | 3.12 downloads 35k codecov 92% license MIT



什么是 Model2Vec?

Model2Vec 是一个全新的嵌入模型，提供了极为高效、轻便且快速的静态嵌入解决方案。与传统的动态嵌入模型（如 Sentence Transformers）不同，Model2Vec 通过对单词或短语进行预计算的方式，在不牺牲性能的情况下大大缩小了模型的体积和提高了速度。

具体来说，Model2Vec 的优势包括：

- 缩小模型体积 15 倍：在同等计算能力下，Model2Vec 的嵌入模型比传统模型小 15 倍，节省了存储空间和计算资源。

- **速度提升500倍**：得益于预计算的静态嵌入技术，Model2Vec 的推理速度比传统的动态嵌入模型快 500 倍，几乎可以实现即时响应。
- **零配置，无需预索引**：与其他需要预先构建索引的搜索方法不同，Model2Vec 支持直接对文档进行向量搜索，极大简化了使用过程。

为什么选择静态嵌入？

Sentence Transformers 这样的传统的动态嵌入模型，通常需要每次处理句子时都计算出新的嵌入，这意味着它们在运行时对资源的需求非常高。而 Model2Vec 采用的是**静态嵌入**，也就是**将每个单词或短语的嵌入提前计算好并保存，避免了每次计算的开销，进而提升了系统的整体效率。**

Model2Vec 的强大功能

1. **即时向量搜索**：无论是数百万文档还是大规模数据集，Model2Vec 都能提供秒级响应，轻松完成向量搜索。
2. **模型压缩与加速**：将模型压缩 15 倍，速度提升 500 倍，Model2Vec 在性能上几乎没有折扣，依然保持了高精度和高效能。
3. **简便易用的蒸馏**：只需几秒钟，就能将复杂的 Sentence Transformers 模型转化为静态嵌入模型，极大简化了开发者的操作流程。
4. **预训练模型**：在 HuggingFace 上，Model2Vec 提供了预训练的最先进的静态嵌入模型，让开发者无需从零开始训练，直接应用。

Model2Vec 与 RAG 的完美结合



Model2Vec 通过对静态嵌入的优化，彻底改变了传统嵌入技术的局限。与传统的动态嵌入模型不同，Model2Vec 的静态嵌入经过预计算，能够快速对大规模数据进行向量检索，**极大加速了 RAG 的“R (Retrieval, 检索)”部分。**

在 RAG 模型中，检索模块是决定生成质量和效率的关键，而 Model2Vec 的优势就在于它能够通过以下方式提升 RAG 的整体表现：

- **超高效的向量检索**：Model2Vec 支持对数百万篇文档进行即时向量检索，无需复杂的预索引过程。这种高效的检索速度直接提升了 RAG 模型中检索模块的响应速度。
- **大幅度压缩与加速**：Model2Vec 将嵌入模型的体积缩小 15 倍，同时速度提升 500 倍，使得 RAG 在执行时的计算开销大大降低，特别适合在需要快速响应的大规模应用场景中使用。
- **与现有 RAG 架构无缝集成**：Model2Vec 可以轻松与像 LangChain 等 RAG 工具集成，帮助开发者优化现有的 RAG workflow，提升数据检索和生成的速度。

开始使用 Model2Vec

开始使用 Model2Vec 的最简单方法是从 HuggingFace 中心加载其旗舰模型之一。这些模型经过预先训练并可立即使用。以下代码片段展示了如何加载模型并进行嵌入：

```
4 from model2vec import StaticModel
5
6 # Load a model from the HuggingFace hub (in this case the potion-base-8M model)
7 model = StaticModel.from_pretrained("minishlab/potion-base-8M")
8
9 # Make embeddings
10 embeddings = model.encode(["It's dangerous to go alone!",
11 "It's a secret to everybody."])
12
13 # Make sequences of token embeddings
14 token_embeddings = model.encode_as_sequence(["It's dangerous to go alone!",
15 "It's a secret to everybody."])
```

官方 GitHub 链接：

<https://github.com/MinishLab/model2vec>



PyTorch研习社

打破知识壁垒，做一名知识的传播者

655篇原创内容

公众号

