

# 基于知识图谱嵌入的阿尔茨海默病药物重定位研究

## 摘要

阿尔茨海默病是一种起病隐匿、多因素、进行性神经退行性疾病，痴呆表现为主要特征，给社会带来巨大医疗负担，但目前还没有特效药物。然而，传统的药物开发存在成本高周期长等问题，且药物安全性需要大量的时间验证，而药物重定位能够极大的缓解上面问题。本文采用知识图谱嵌入研究阿尔茨海默病的药物重定位。首先，利用 4 种知识图谱嵌入模型对知识图谱进行表示学习；然后，使用多种评估指标评估了知识图谱嵌入模型的性能和学习到的嵌入向量的质量；最后，利用 RotatE 模型进行链接预测得出 16 种治疗阿尔茨海默病的候选药物。实验结果表明本文的研究方法能够有效的预测阿尔茨海默病的治疗药物，为研究人员提供了新的研究方法。本文的源代码可以从 <https://github.com/LuYF-Lemon-love/AD-KGE> 获得。

**关键词：**药物重定位；阿尔茨海默病；知识图谱；知识图谱嵌入；知识图谱补全

**中图分类号** TP399

## Study on drug repurposing for Alzheimer's disease based on knowledge graph embedding

Alzheimer's disease is a latent, multifactorial, progressive neurodegenerative disease characterized by dementia, which poses a huge medical burden to society, but there are currently no specific drugs available. However, traditional drug development has problems such as high costs and long cycles, and drug safety requires a lot of time to verify. Drug repurposing can greatly alleviate the above problems. This paper used knowledge graph embedding to study drug repurposing in Alzheimer's disease. Firstly, four knowledge graph embedding models are used to represent and learn knowledge graph; Then, the performance of the knowledge graph embedding model and the quality of the learned embedding vectors are evaluated using various evaluation metrics; Finally, sixteen candidate drugs for treating Alzheimer's disease were obtained through link prediction using the RotatE model. The experimental results show that the research method can effectively predict the therapeutic drugs for Alzheimer's disease, providing researchers with a new research method. The source code of this paper can be obtained from <https://github.com/LuYF-Lemon-love/AD-KGE>.

**Key words :** drug repurposing ; Alzheimer's disease ; knowledge graph ; knowledge graph embedding ; knowledge graph completion

# 1. 引言

阿尔茨海默病 (Alzheimer's disease, AD) 是一种常见的神经退行性疾病, 无法治愈且不可逆转<sup>[1]</sup>。其特征是伴有神经精神症状的渐进性严重痴呆<sup>[2]</sup>。据报道, 2021 年我国 60 岁以上人群中 有 983 万例 AD 患者<sup>[3]</sup>; 且另一份研究报告称, 到 2050 年, 我国 AD 患者的治疗费用将高达 18871.8 亿美元<sup>[4]</sup>。这些数据充分说明了 AD 对社会经济造成了巨大的负担, 开发 AD 的治疗药物势在必行。

然而, 研发一款新药至少花费 26 亿美元<sup>[5]</sup>和 10 年时间<sup>[6]</sup>, 这需要海量的金钱和时间成本。药物重定位, 又可以称为“老药新用”, 指的是从获批准的临床药物中发现新适用的病症或新用途的方法。该方法具有低成本、高效率的特点, 在突发性疾病和罕见病方面优势更为突出。近年来, 药物重定位得到了迅速发展, 领域内已经出现了很多用于探索药物和疾病之间关系的方法。其中, 知识图谱 (Knowledge Graph, KG) 就是实现药物重定位的一个重要举措<sup>[6]</sup>。

KG 是一种基于拓扑结构图存储知识的数据库。知识中的具体事物和抽象概念在 KG 中被表示为实体, 实体之间的联系被表示为关系, 进而知识被表示成格式为 (头实体, 关系, 尾实体) 的三元组。KG 是一个由大量的三元组组成的有向图结构, 图中的节点表示实体, 边表示实体间的关系。

然而, 许多 KG 都非常巨大, 如药物再利用知识图谱 (Drug Repurposing Knowledge Graph, DRKG) <sup>[7]</sup> 包含 97238 个实体和 5874261 个三元组。因此, 常采用知识图谱嵌入 (Knowledge Graph Embedding, KGE) 技术将实体和关系表示成低维稠密向量, 进而将 KG 建模成低维向量空间。在过去几年中, 研究人员提出了很多 KGE 模型, 如 TransE<sup>[8]</sup>、DistMult<sup>[9]</sup>、ComplEx<sup>[10]</sup> 和 RotatE<sup>[11]</sup> 等, 来学习实体和关系嵌入向量。KGE 模型能够利用各自对应的模型假设进行链接预测进而推测三元组中缺失的实体。因此使用 KG 进行药物重定位研究, 本质上就是使用 KGE 模型进行“疾病”实体和“药物”实体之间缺失关系的预测。

近年来, 研究人员提出了很多利用 KG 进行药物重定位的方法。Zeng 等<sup>[12]</sup> 建立了一个 1500 万个三元组的综合 KG, 包括药物、基因、疾病、药物副作用 4 种实体以及它们之间的 39 种关系, 然后利用 RotatE 学习实体和关系的表示, 进而确定了 41 种针对 COVID-19 的治疗药物。Zhang 等<sup>[13]</sup> 提出了一种基于神经网络和文献发现的方法, 首先利用 PubMed 和其他专注 COVID-19 的研究文献构建了一个生物学 KG, 然后利用多种 KGE 模型预测 COVID-19 的候选治疗药物, 并利用发现模式解释了 KGE 预测的合理性。目前也有研究人员利用 KGE 模型研究帕金森病的药物重定位, 并取得了不错的效果<sup>[14]</sup>。

Wang 等<sup>[6]</sup> 提出了一种基于 KG 的深度学习方法进行 AD 药物重定位。首先, 利用 DistMult 学习了预先构建的药物靶点对 KG 的实体和关系的嵌入表示, 然后利用一个 Conv-Conv 模块来提取药物-靶点对的特征, 提取到的特征被传入到一个全连接网络进行二分类, 最终通过载脂蛋白 E 作为靶点寻找治疗 AD 的药物。Nian 等<sup>[1]</sup> 从文献中构建一个 KG, 利用 TransE、DistMult 和 ComplEx 预测有助于 AD 治疗或预防的候选物质, 以研究 AD 与化学物质、药物和膳食补充剂之间的关系, 进而确定预防或延缓神经退行性进展的机会。

虽然上述 AD 药物重定位的研究工作都取得了很好的效果, 但依旧有一些可改进之处。在 Wang 等<sup>[6]</sup> 进行 AD 药物重定位的工作中, 构建了一个只有 4 种实体 (药物、靶标、酶、载体) 和 1 种关系类型 (药物-靶点) 的药物靶点 KG (包含 29607 个三元组), 该 KG 只能

利用单一的药物靶点信息，忽略了大量对 AD 药物重定位有用的其他信息，如基因、药物副作用和症状等。在 Nian 等<sup>[1]</sup>利用文献挖掘进行 AD 药物重定位的工作中，利用规则和 BERT 分类器仅仅保留了与 AD 相关的 791827 个三元组，虽然取得了很好的效果，但也忽略了很多其他疾病的信息，如 AD 病人有很多精神症状，其他精神疾病实体也许对 AD 的药物重定位很有帮助。

针对以上问题，本文采用 KGE 模型在大型 KG (DRKG) 上研究了 AD 药物重定位。首先，利用多种 KGE 模型 (TransE、DistMult、ComplEx 和 RotatE) 在 DRKG 上学习实体和关系的嵌入向量，通过 3 种经典的 KG 评估指标评估了 4 种 KGE 模型；然后，在整个 KG 上重新训练 KGE 模型，并利用多种嵌入向量分析手段评估了模型学习到的嵌入向量的质量；最终，根据 KGE 模型的评估结果选择 RotatE 作为最终的药物重定位模型，找到了 16 种治疗 AD 的候选药物。

## 2. 方法

### 2.1. 数据

DRKG<sup>[7]</sup>是一个涉及基因、药物、疾病、生物过程、副作用和症状的综合生物 KG，包括来自 DrugBank、Hetionet、GNBR、String、IntAct 和 DGIdb 等六个现有数据库的信息，以及从最近发表的 Covid19 出版物（截止到 2020 年 3 月 22 日）中收集的数据（后文标记为 bioarx 数据库）。它有属于 13 种实体类型的 97238 个实体；以及属于 107 种关系类型的 5874261 个三元组。DRKG 使用“实体类型::ID”的格式表示一个实体，如“Disease::MESH:D000544”，其中“Disease”是实体类型，“MESH:D000544”是 ID；使用“数据源名::关系名::头实体类型:尾实体类型”的格式表示关系，如“DRUGBANK::treats::Compound:Disease”，其中“DRUGBANK”是数据源名，“treats”是关系名，“Compound”是头实体类型，“Disease”是尾实体类型。

### 2.2. KGE 模型基本原理

为了实现在 DRKG 上学习实体和关系的嵌入向量，考虑到算力限制，本文仅研究和对比了 4 种经典且具有线性时间复杂度的 KGE 模型，即 TransE<sup>[8]</sup>、DistMult<sup>[9]</sup>、ComplEx<sup>[10]</sup>、RotatE<sup>[11]</sup>。在利用 KGE 模型来推断现有 KG 的缺失关系，从而达到补全 KG 的任务中，KG 通常被标记为  $T$ ，是一组格式为  $(h, r, t)$  三元组的集合，其中  $h, t \in E$ ,  $r \in R$ ,  $E$  是 KG 的实体集合， $R$  是 KG 的关系集合。KGE 模型一般都具有一个度量  $(h, r, t)$  成立概率的评分函数，该评分函数是特定 KGE 模型对 KG 的建模假设<sup>[11]</sup>。

#### 2.2.1. TransE 模型基本原理

TransE<sup>[8]</sup>是一个代表性的平移模型，它假设实体和关系属于同一向量空间  $\mathbb{R}^d$ ,  $d$  是向量空间的维度。关系  $r$  被建模为实体向量的平移，如果三元组  $(h, r, t)$  成立，那么  $h + r \approx t$ ，即  $t$  应该是  $h + r$  最近的实体向量；如果不成立， $h + r$  应该远离  $t$ 。TransE 只能建模 1 对 1 的关系类型；但是从另一种关系分类角度，它能捕获反对称、反转和组成三种关系但不能

捕获对称关系<sup>[11]</sup>。TransE 的评分函数如公式(1)所示。

$$f(h, r, t) = -\|h + r - t\|_{L_1/L_2} \quad (1)$$

如公式(1)所示，TransE 依据距离函数（ $L_1$  范数和  $L_2$  范数）选择的不同有两个变体分别为 TransE<sub>l1</sub> 和 TransE<sub>l2</sub>。

### 2.2.2. DistMult 模型的基本原理

DistMult<sup>[9]</sup>是一个双线性模型，它为每一种关系提供了一个对角矩阵来建模实体之间的交互进而捕获 KG 的潜在语义。DistMult 也假设实体和关系属于同一向量空间  $\mathbb{R}^d$ ，其评分函数如公式(2)所示。

$$f(h, r, t) = h^T \text{diag}(r) t \quad (2)$$

其中， $\text{diag}(r)$ 是关系  $r$  的对角矩阵。

### 2.2.3. ComplEx 模型的基本原理

由于 DistMult<sup>[9]</sup>使用的是对角矩阵，因此仅仅能捕获对称关系。为了捕获反对称和反转关系，ComplEx<sup>[10]</sup>将向量空间从实数域扩展到复数域，极大的提升了模型的表现能力。ComplEx 假设实体和关系属于同一复数向量空间  $\mathbb{C}^d$ ，其评分函数如公式(3)所示。

$$f(h, r, t) = \text{Real}(h^T \text{diag}(r) \bar{t}) \quad (3)$$

其中， $\text{Real}(\cdot)$ 表示复数的实部， $\bar{t}$ 表示  $t$  的共轭。

### 2.2.4. RotatE 模型的基本原理

受到 TransE 和欧拉恒等式的启发，RotatE<sup>[11]</sup>将头实体和尾实体映射到复数向量空间，即当  $h, t \in \mathbb{C}^d$ ， $r \in \mathbb{C}^d$ ， $|r_i| = 1$ ，将关系  $r$  建模为从头实体  $h$  到尾实体  $t$  的逐元素旋转。RotatE 模型能够捕获对称、反对称、反转和组成四种类型关系，其评分函数如公式(4)所示。

$$f(h, r, t) = -\|h \circ r - t\|^2 \quad (4)$$

其中， $\circ$ 表示哈达玛积。

### 2.2.5. 优化

本文使用最大间隔方法训练模型，以最小化正确三元组的排名<sup>[8]</sup>，其损失函数如公式(5)所示。

$$\mathcal{L} = \sum_{(h,r,t) \in T} \sum_{(h',r,t') \in T^-} \max(0, \gamma - f(h,r,t) + f(h',r,t')) \quad (5)$$

其中,  $\gamma > 0$  是正负例三元组得分的间隔距离。 $T$ 是正例三元组集合,  $T^-$ 是负例三元组的集合, 如公式(6)所示, 它是通过破坏原有三元组中的实体和关系得到的<sup>[15]</sup>。

$$T^- = E \times R \times E - T \quad (6)$$

## 2.3. KGE 模型的评估

### 2.3.1. 经典评估

KGE 模型可以通过链接预测技术预测 KG 中缺失的三元组, 即给定 $(h, r, ?)$ 预测缺失的尾实体  $t$ , 或者给定 $(?, r, t)$ 预测缺失的头实体  $h$ 。可以通过链接预测给出正确实体的排名。常使用三种经典指标来评估链接预测的性能: 正确实体评分函数的平均排名 (Mean Rank, MR)<sup>[8]</sup>, 正确实体评分函数的平均倒数排名 (Mean Reciprocal Rank, MRR)<sup>[11]</sup>和正确实体评分函数的前  $N$  的比例即前  $N$  命中率 Hits@N ( $N = 1, 3, 10$ )<sup>[8]</sup>。

如果用 $\text{rank}_h$ 和 $\text{rank}_t$ 分别表示预测正确头实体和尾实体的排名,  $T$ 表示需要评估的三元组集合, 那么 MR、MRR 和 Hits@N 的具体的计算方法分别如公式(7)、(8)和(9)所示。

$$MR = \frac{1}{2|T|} \sum_{(h,r,t) \in T} \text{rank}_h + \text{rank}_t \quad (7)$$

$$MRR = \frac{1}{2|T|} \sum_{(h,r,t) \in T} \frac{1}{\text{rank}_h} + \frac{1}{\text{rank}_t} \quad (8)$$

$$\text{Hits}@N = \frac{1}{2|T|} \sum_{(h,r,t) \in T} I[\text{rank}_h \leq N] + I[\text{rank}_t \leq N] \quad (9)$$

在公式(9)中, 如果条件为真,  $I[*]$ 等于 1, 否则等于 0。从式公式(7)、(8)和(9)可知, 对于相同的  $T$ , MR 值越小, 代表正确实体的排名越靠前, 说明链接预测越精确; MRR 和 Hits@N 值越大, 代表正确实体的排名越靠前, 说明链接预测越精确。

### 2.3.2. 嵌入评估

由于 DRKG 结合了来自不同数据源的信息, 本文通过嵌入评估来定性验证 KGE 模型是否生成了有意义的实体和关系嵌入。理想的情况是, KGE 模型能够学习到不同关系嵌入向量的差异之处和相同类型实体的相似之处。

本文首先采用  $t$  分布随机近邻嵌入 (T-distributed Stochastic Neighbor Embedding, t-SNE)<sup>[16]</sup>将关系嵌入向量进行降维并可视化。DRKG 共有来源于 7 个数据库的 107 种关系类型, 如果相同数据来源的关系向量在可视化图中越分散, 就说明 KGE 模型越能学习到不同关系嵌入向量的差异之处, 即使它们来源于同一数据库。进一步地, 如公式(10)所示, 本文还使

用了余弦相似度来计算 DRKG 的关系嵌入向量对之间的相似性，并通过对比相似度分布的直方图来评估各种 KGE 模型。不同关系嵌入向量的相似度越低，表示 KGE 模型越能捕捉到不同关系嵌入向量的差异。使用这样的 KGE 模型进行链接预测的效果也就越好。

$$similarity = \cos(\theta) = \frac{a \cdot b}{\|a\| \|b\|} = \frac{\sum_{i=1}^n a_i \times b_i}{\sqrt{\sum_{i=1}^n (a_i)^2} \times \sqrt{\sum_{i=1}^n (b_i)^2}} \quad (10)$$

在公式(10)中,  $a_i$  和  $b_i$  分别表示向量  $a$  和  $b$  的第  $i$  个分量, 余弦相似度的取值范围为  $[-1, 1]$ ,  $-1$  表示两个向量方向相反,  $1$  表示方向相同,  $0$  表示相互独立。

接下来本文使用主成分分析将实体嵌入向量降到 30 维<sup>[16]</sup>, 并利用 t-SNE 将其降到 2 维空间进行可视化。使用主成分分析的原因在于本文的研究对象中共有 97238 个实体, 数量众多, 若直接利用 t-SNE 降维和可视化, 可能会引入大量噪声。DRKG 有共有 13 种实体类型, 相同类型的实体在可视化图中越聚集, KGE 模型对实体嵌入的效果就越好。

## 2.4. AD 药物重定位

使用 KGE 模型做药物重定位时, 将 Drugbank 中被 FDA 批准的药物作为候选药物 (相对分子质量  $\geq 250$ , 共 8104 个), 它们构成了头实体集合。选择 DRKG 中所有治疗关系作为链接预测的关系, 共有 DRUGBANK::treats::Compound:Disease, GNBR::T::Compound:Disease, Hetionet::CtD::Compound:Disease 三种, 其中 treats、T、CtD 分别是 DrugBank 数据库、GNBR 数据库、Hetionet 数据库中的治疗关系。选择 DRKG 中全部 AD 实体作为尾实体集合, 共有 Disease::DOID:10652, Disease::MESH:C536599, Disease::MESH:D000544 三种, 其中 Disease::DOID:10652 是来自 Hetionet 数据源的 AD 实体, Disease::MESH:C536599 和 Disease::MESH:D000544 是被映射到 MESH ID 的 AD 实体 (其中 Disease::MESH:C536599 是无神经纤维缠结 AD 的实体)。将上面实体和关系集合进行格式为 (h, r, t) 排列组合 (总共  $8104 \times 3 \times 3 = 72936$  种可能), 然后计算所有组合评分函数的得分, 最后选择得分前  $N$  的药物作为 AD 的治疗药物, 其中  $N$  的值取决于不同 KGE 模型在测试集上的 MR 指标结果。

## 2.5. 实验设置

将 DRKG 的三元组按照 90%、5%、5% 的比例划分为训练集、验证集和测试集, 分别为 5286834 个、293713 个和 293714 个。

综合 5 个经典的 KGE 评估指标 (即 MR、MRR、Hits@1、Hits@3、Hits@10) 的综合表现, 在验证集上利用网格搜索所有模型的超参数 (TransE\_l1、TransE\_l2、DistMult、Complex 和 RotatE)。所有模型的训练批处理大小和每个正例三元组使用的负例三元组的数量分别固定为 4096 和 256, 学习率 (learning rate, lr) 则都从 {0.01, 0.05, 0.1} 中选择。由于 RotatE 模型实体维度是超参数嵌入维度 (the embedding dimension, hidden\_dim) 的 2 倍, 本文选择将 RotatE 模型的 hidden\_dim 固定为 200, 其他模型的 hidden\_dim 则从 {200, 400} 中选择。对于超参数  $\gamma$ , TransE\_l1、TransE\_l2 和 RotatE 从 {6, 12, 18} 中选择, 而 DistMult、Complex 模型

则从{50,125,200}中进行选择。

本文的实验是利用 Zheng 等<sup>[17]</sup>开发 DGL-KE 工具包实现的。

### 3. 结果

#### 3.1. KGE 模型的经典评估

表 1 列出了在 KG 补全任务中, 4 种 KGE 模型在测试集上的结果。如表 1 所示, 对于 MR 指标, TransE 两种变体分别取得了最优结果 60.83 和次优结果 62.64; 对于 MRR 指标, ComplEx 取得了最优结果 0.621, RotatE 次之为 0.614; 对于 Hits@1 指标, ComplEx 取得了最优结果为 0.537, RotatE 次之为 0.515; 对于 Hits@3 和 Hits@10, RotatE 取得了最优结果分别为 0.681 和 0.780, ComplEx 取得了次优结果分别为 0.673 和 0.768。而 DistMult 在 3 种指标上都没有取得最优和次优结果。

**Table 1** The traditional evaluation results of the KGE model. The best results are in **bold** and the second best results are in underline.

Model	MRR	MR	Hits@1	Hits@3	Hits@10
TransE_l1	0.530	<u>62.64</u>	0.412	0.606	0.740
TransE_l2	0.437	<b>60.83</b>	0.302	0.515	0.693
DistMult	0.484	105.55	0.401	0.515	0.643
ComplEx	<b>0.621</b>	112.74	<b>0.537</b>	<u>0.673</u>	<u>0.768</u>
RotatE	<u>0.614</u>	63.51	<u>0.515</u>	<b>0.681</b>	<b>0.780</b>

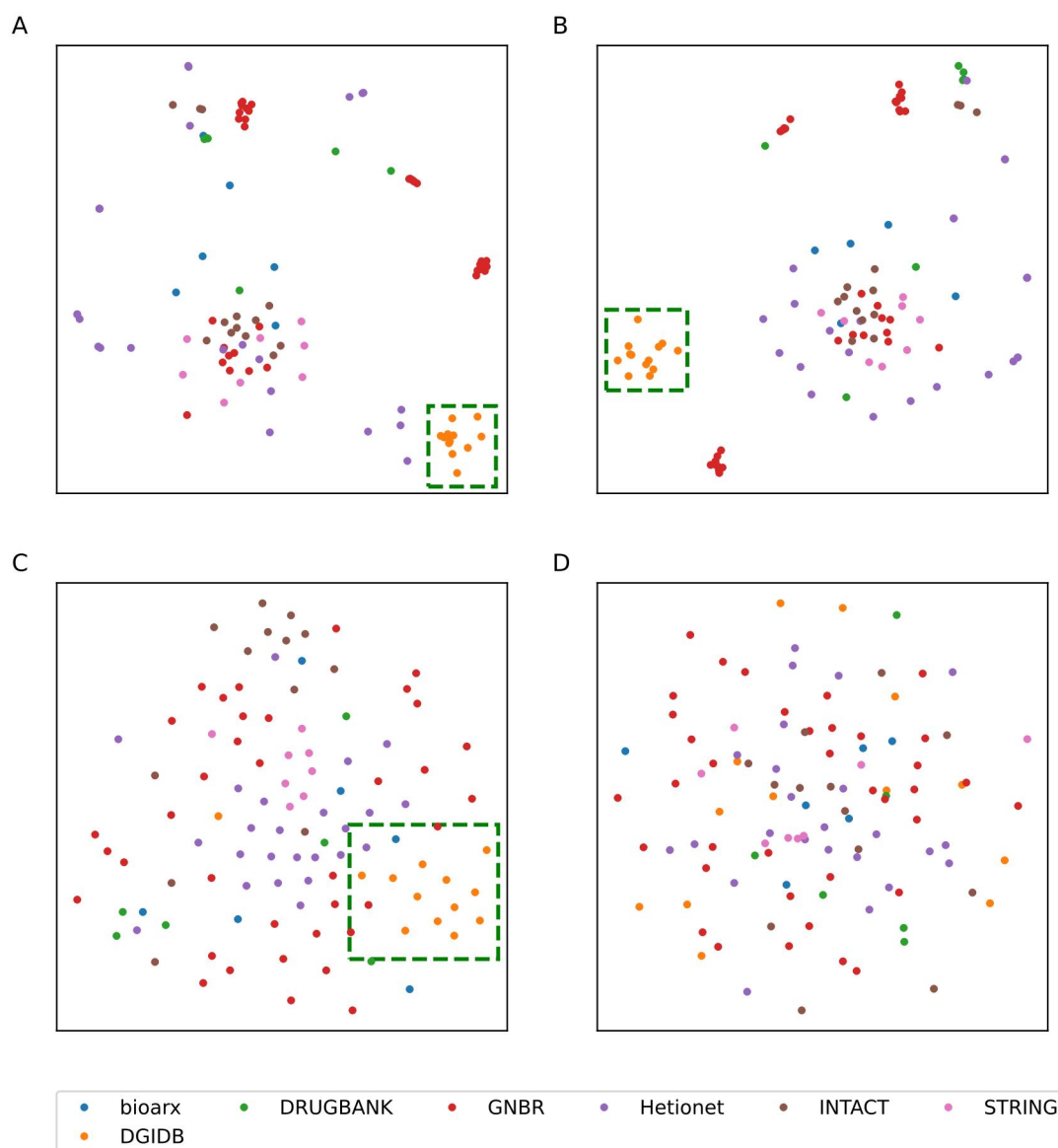
各个模型超参数的最佳配置是: 对于 TransE\_l1, hidden\_dim=400,  $\gamma=18$ , lr=0.05; 对于 TransE\_l2, hidden\_dim=400,  $\gamma=12$ , lr=0.1; 对于 DistMult, hidden\_dim=400,  $\gamma=50$ , lr=0.1; 对于 ComplEx, hidden\_dim=400,  $\gamma=50$ , lr=0.1; 对于 RotatE, hidden\_dim=200,  $\gamma=18$ , lr=0.05。

鉴于 DistMult 模型在经典评估中并不出色的表现, 本文仅选择 TransE\_l1、TransE\_l2、ComplEx 和 RotatE 模型, 利用最佳超参数, 重新在整个 DRKG 上进行训练, 并进一步进行模型的嵌入评估和 AD 药物重定位。

#### 3.2. KGE 模型的嵌入评估

图 1A、1B、1C、1D 分别展示了 TransE\_l1、TransE\_l2、ComplEx 和 RotatE 的关系嵌入向量在 2D 空间的可视化图。图中每一个圆点代表 DRKG 中一种关系类型, 因此共有 107 个圆点; 相同颜色的圆点代表关系来自相同的 DRKG 中相同的数据库。从图 1A、1B 和 1C 中可以看出, TransE\_l1、TransE\_l2 和 ComplEx 的关系嵌入向量出现不同程度的同数据源聚集现象, 如代表虚线框中标注出来的、代表 DGIdb 数据源的橙色点; 而 RotatE 的关系嵌入向量广泛的分布在 2D 的空间中, 即便来自相同源数据集的关系都没有出现聚集的现象,

可以说，RotatE 更好地学习到了各个关系本身的差异，受数据源的影响较小。

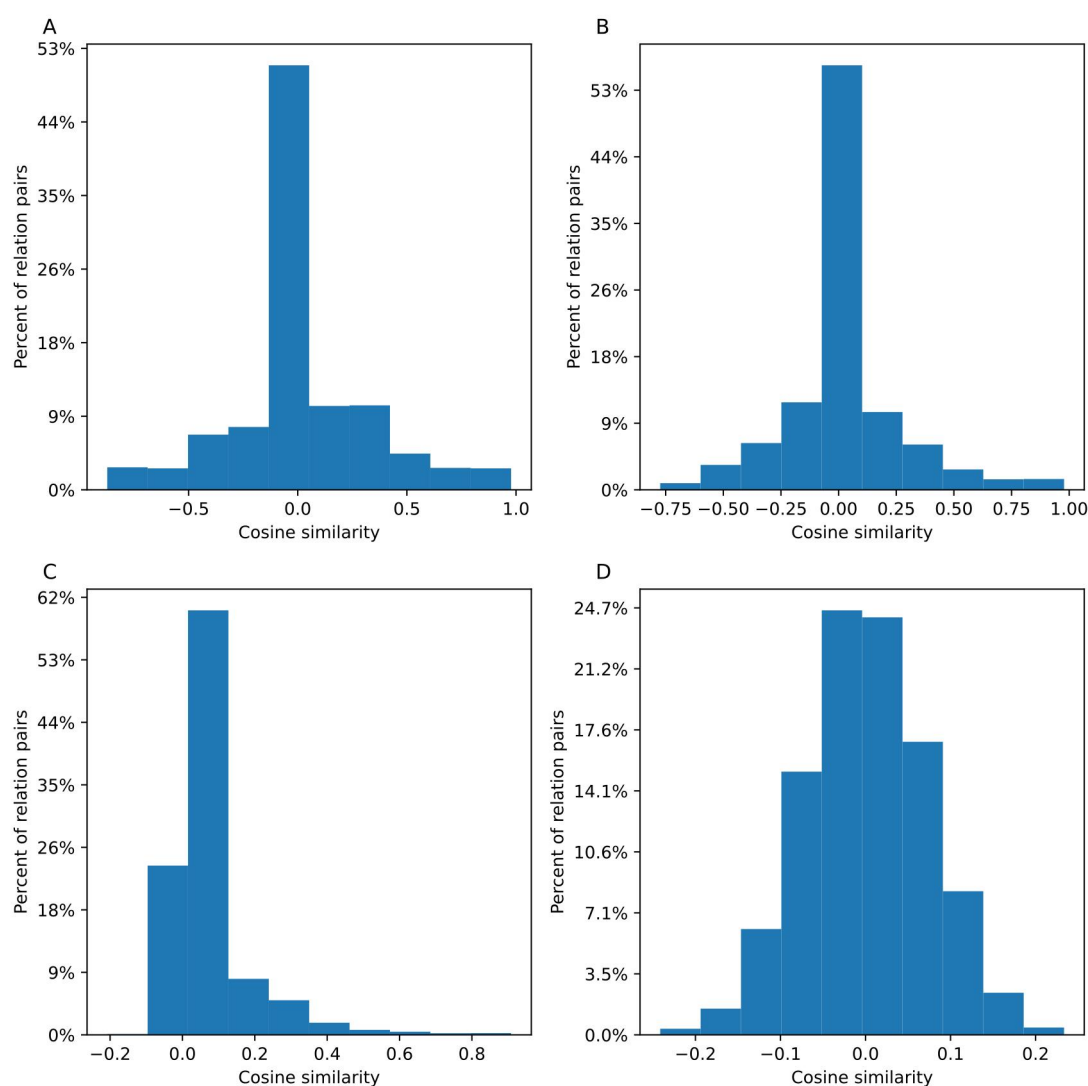


**Figure 1** Distribution of relation embeddings in 2D euclidean space for 4 models. Subgraphs A, B, C and D are the results of TransE\_l1, TransE\_l2, ComplEx and RotatE respectively.

图 2A、2B、2C、2D 显示了 TransE\_l1、TransE\_l2、ComplEx 和 RotatE 的不同关系嵌入向量对之间的余弦相似度分布直方图。对于 TransE\_l1，相似度值分布在 $[-0.873, 0.977]$ 范围内，其中约有 7%相似度值大于 0.50 的关系对；TransE\_l2 与 TransE\_l1 类似，也存在着 5%相似度大于 0.50 的关系对。ComplEx 模型的相似度值分布在 $[-0.208, 0.908]$ 范围内，存在 1%相似度大于 0.50 的关系对。相比而言，RotatE 模型的相似度整体都较小，分布在 $[-0.241, 0.233]$ 的范围内。进一步地，本文考察了包含并且只包含一种治疗关系的嵌入向量对之间余弦相似度的最大值，TransE\_l1 为 0.917，TransE\_l2 为 0.841，ComplEx 为 0.225，RotatE 为 0.180。这就说明对于 TransE\_l1 和 TransE\_l2，存在着与治疗关系非常相似的其他类型的关系向量，这很可能会干扰链接预测的结果。而对于 RotatE 模型，治疗关系向量与其他类型的关系向量之间的相似度值最高也仅为 0.18，说明治疗关系与其他类型的关系有着极小的相

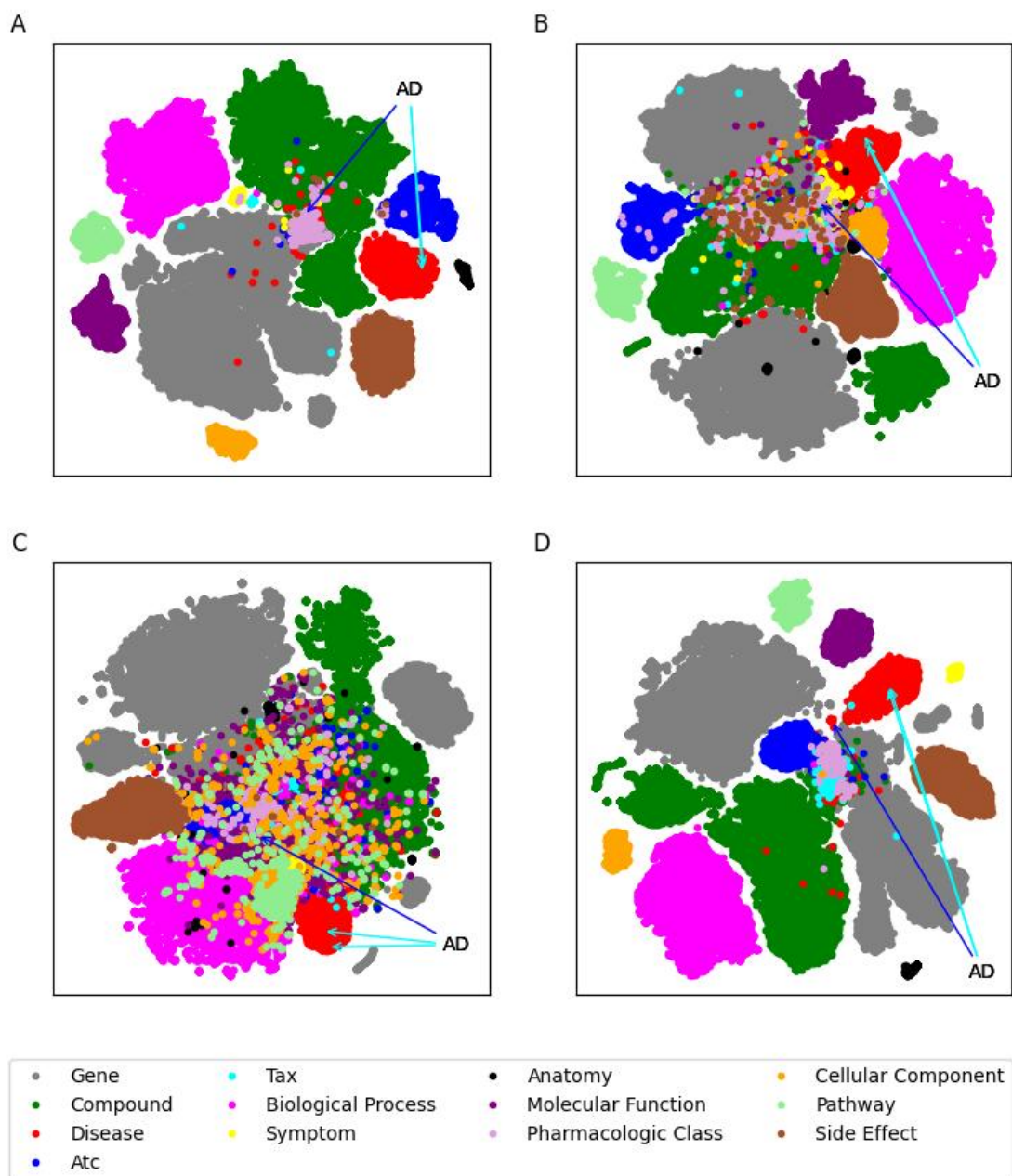


似性，在链接预测时，不易受到其他关系类型的影响。



**Figure 2** Histogram of cosine similarity between relations for 4 models. Subgraphs A, B, C and D are the results of TransE\_l1, TransE\_l2, ComplEx and RotatE respectively.

图 3A、3B、3C、3D 是 TransE\_l1、TransE\_l2、ComplEx 和 RotatE 的实体嵌入 2D 空间的可视化图，每一个圆点代表了一个实体，不同的颜色代表不同的实体类型。用蓝色和蓝绿色箭头指出了药物重定位 3 个 AD 实体，蓝色箭头指向的是 Disease::DOID:10652 实体，它是来自 Hetionet 数据源的 AD 实体。从图 3 中可以看到，在所有模型中，相同类别的实体正如期望的那样聚集到了一起，其中 TransE\_l1 和 RotatE 的结果要优于另外 2 个模型。4 个模型都将来自 Hetionet 数据源的 AD 实体和来自 MESH ID 空间中两种 AD 实体区分开了。2 种 MESH ID 空间的 AD 实体在 TransE\_l1、TransE\_l2 和 RotatE 的 2D 空间中距离很近，但在 ComplEx 的 2D 空间中这两种实体还有较大距离。



**Figure 3** Distribution of entity embeddings in 2D euclidean space for 4 models. Subgraphs A, B, C and D are the results of TransE\_l1, TransE\_l2, ComplEx and RotatE respectively.

### 3.3. AD 药物重定位

综合 KGE 的经典评估和嵌入评估结果，我们将使用 RotatE 模型作为 AD 药物重定位的最终模型。由于 RotatE 的 MR 指标结果是 63.51，因此我们将得分前 50 的药物作为候选药物。

在得分前 10 的候选药物列表中，只有第 9 名的药物没有被 DRKG 标注为对 AD 有治疗作用，说明该方法能够正确表达 KG 中原有的三元组，证明了该方法的有效性。因此本文将得分前 50 的候选药物列表中，没有被 DRKG 标注为对 AD 有治疗作用的药物作为重定位得

到的 AD 候选药物。考虑到其中得分排名在第 23 的西布曲明已退市<sup>[26]</sup>，因此最终剩余 16 种候选药物。表 2 列出了这些药物在 RotatE 模型中排名、在 DRKG 中的名称以及文献中提及的该药物与 AD 的关系（其中 None 代表我们暂时并未找到文献报道它们之间的关系）。从表 2 中可以看到，仅有 4 种排名较靠后（31、41、44、46）的药物，我们暂未从文献中找到它们与 AD 之间的关系外，其余药物都被文献证实可能是 AD 的潜在药物。

**Table 2** Candidate drugs obtained from drug repurposing.

Rank	Drug name	Literature support
9	Glutathione	The beneficial effect of many nutrients on the course of AD has been demonstrated. These include: glutathione, polyphenols, curcumin, coenzyme Q10, vitamins B6, B12, folic acid, unsaturated fatty acids, lecithin, UA, caffeine and some probiotic bacteria <sup>[18]</sup> .
11	Haloperidol	Haloperidol inactivates AMPK and reduces tau phosphorylation in a tau mouse model of Alzheimer's disease <sup>[19]</sup> .
13	Capsaicin	In Alzheimer's disease, capsaicin reduces neurodegeneration and memory impairment <sup>[20]</sup> .
16	Quercetin	Quercetin has demonstrated antioxidant, anti-inflammatory, hypoglycemic, and hypolipidemic activities, suggesting therapeutic potential against type 2 diabetes mellitus (T2DM) and Alzheimer's disease (AD) <sup>[21]</sup> .
17	Estradiol	Mounting evidence indicates that the neurosteroid estradiol (17 $\beta$ -estradiol) plays a supporting role in neurogenesis, neuronal activity, and synaptic plasticity of AD. This effect may provide preventive and/or therapeutic approaches for AD <sup>[22]</sup> .
18	Glucose	Specifically, decreased O-GlcNAcylation levels by glucose deficiency alter mitochondrial functions and together contribute to Alzheimer's disease pathogenesis <sup>[23]</sup> .
20	Disulfiram	Identification of disulfiram as a secretase-modulating compound with beneficial effects on Alzheimer's disease hallmarks <sup>[24]</sup> .
21	Adenosine	Emerging evidence suggests adenosine G protein-coupled receptors (GPCRs) are promising therapeutic targets for Alzheimer's disease <sup>[25]</sup> .
23	Sibutramine	In October 2010, Sibutramine was withdrawn from U.S. <sup>[26]</sup> .
29	Paroxetine	Paroxetine ameliorates prodromal emotional dysfunction and late-onset memory deficit in Alzheimer's disease mice <sup>[27]</sup> .
31	Cocaine	None.
39	Paclitaxel	In addition to NSAIDs, an anticancer drug, paclitaxel, has considerable potential as an AD treatment <sup>[28]</sup> .
41	Cholesterol	None.

43	Glyburide	Our findings suggest that a pharmacologic approach to inhibit galanin in the brain, either by glibenclamide or pioglitazone might dramatically improve symptoms in Alzheimer's disease <sup>[29]</sup> .
44	Staurosporine	None.
46	Cortisone	None.
48	Amitriptyline	These results indicate that amitriptyline has significant beneficial actions in aged and damaged AD brains and that it shows promise as a tolerable novel therapeutic for the treatment of AD <sup>[30]</sup> .

## 4. 讨论

在这项研究中，利用 KGE 模型研究了 AD 的药物重定位。具体来说，采用 RotatE 学习 DRKG 的实体和关系的嵌入向量表示，然后利用链接预测技术发现 AD 的候选治疗药物。实验结果表明，RotatE 能够整合 DRKG 的多源信息，进而完成 AD 的药物重定位任务。相比于仅仅是利用 AD 相关三元组和单一药物靶点相互作用构建 KG，本文发现利用大型多实体类型和多关系类型的 KG 对 AD 药物重定位是有益的。

本文使用的数据集是涉及 13 种实体和 107 种关系，包含 5874261 个三元组的 DRKG。DRKG 包含 5103 个疾病实体，能够使 KGE 模型很好利用与 AD 相关的精神疾病实体信息；13 种实体类型和 107 种关系包含了各种各样的生物信息，KGE 模型能够将上述多种生物信息嵌入到一个统一的向量空间中，进而能够利用整个 KG 的信息进行 AD 的药物重定位。

通过嵌入评估结果，发现 RotatE 能够整合 DRKG 中来自多个数据源的信息，避免了不同数据源的实体和关系相互独立，进而影响 AD 药物重定位的效果。在 RotatE 得分前 10 的候选药物列表中，只有第 9 名的药物没有被 DRKG 标注为对 AD 有治疗作用，因此可以认为 RotatE 拟合了 DRKG，证明了该方法的有效性。16 种候选药物中仅有 4 种排名较靠后（31、41、44、46）的药物，我们暂未从文献中找到它们与 AD 之间的关系外，其余药物都被文献证实可能是 AD 的潜在药物。上述 AD 药物重定位实验结果，证明了利用大型多实体类型和多关系类型的 KG 对 AD 药物重定位是有益的。

由于 Wang 等<sup>[6]</sup>仅仅针对载脂蛋白 E 靶点进行 AD 药物重定位，得到的候选药物是一些金属和金属化合物，如锌、铜、银、氯化锌、醋酸锌、硫酸锌，其中锌和铜已经被 DRKG 标注为对 AD 有治疗作用；而本文使用 DRKG 学习到了各种各样的生物信息，得到的 AD 重定位结果是全方位的，因此，它们的结果与本文的结果没有重合也是合理的，这也说明了使用综合 KG 进行药物重定位的必要性。本文预测结果排名第 46 位可的松也出现在了 Nian 等<sup>[1]</sup>通过预防关系进行 AD 药物重定位的结果中，并且可的松在他们的结果中取得了第 1 的位置，但是并没有发现可的松直接治疗 AD 的证据；在本文的预测结果中有很多治疗其他精神疾病的药物，如排名第 11 位的氟哌啶醇（治疗精神分裂症）、第 29 位的帕罗西汀（治疗重度抑郁症、恐慌症、强迫症）和排名第 48 位的阿米替林（抗抑郁药），证明了 RotatE 确实利用其他疾病和 AD 的关系进行药物重定位，进一步证明了利用大型多实体类型和多关系类型的 KG 对 AD 药物重定位是有益的。

由于 DRKG 没有将所有的疾病都映射到统一的 ID 空间，如 Disease::DOID:10652，这

对药物重定位的效果产生了一定的影响。在构建 KG 时，有必要将同类型的实体映射到一个统一的 ID 空间，这对 KGE 模型学习嵌入向量将有很大的帮助。未来，我们将研究更多种类的 KGE 模型在药物重定位中的应用；也将研究实体对齐技术，来将多种数据源的实体映射到统一的命名空间中，进而使得 KGE 模型学习到更好的嵌入向量。

## References

- [1] Nian Y, Hu XY, Zhang R, et al. Mining on Alzheimer's diseases related knowledge graph to identify potential AD-related semantic triples for drug repurposing[J]. BMC Bioinformatics, 2022, 23(Suppl 6): 407. <https://doi.org/10.1186/s12859-022-04934-1>.
- [2] Moya-Alvarado G, Gershoni-Emek N, Perlson E, et al. Neurodegeneration and Alzheimer's disease (AD).What can proteomics tell us about the Alzheimer's brain?[J]. Mol Cell Proteomics, 2016, 15(2): 409-25. <https://doi.org/10.1074/mcp.R115.053330>.
- [3] Ren RJ, Yin P, Wang ZH, et al. China Alzheimer disease report 2021[J]. Journal of Diagnostics Concepts & Practice(诊断学理论与实践), 2021, 20(04): 317-337. <https://doi.org/10.16150/j.1671-2870.2021.04.001>.
- [4] Jia JP, Wei CB, Chen SQ, et al. The cost of Alzheimer's disease in China and re-estimation of costs worldwide[J]. Alzheimers Dement, 2018, 14(4): 483-491. <https://doi.org/10.1016/j.jalz.2017.12.006>.
- [5] Avorn J. The \$2.6 billion pill—methodologic and policy considerations[J]. N Engl J Med, 2015, 372(20): 1877-1879. <https://doi.org/10.1056/NEJMp1500848>.
- [6] Wang SD, Du ZZ, Ding M, et al. KG-DTI: a knowledge graph based deep learning method for drug-target interaction predictions and Alzheimer's disease drug repositions[J]. Applied Intelligence, 2022, 52(1): 846–857. <https://doi.org/10.1007/s10489-021-02454-8>.
- [7] Ioannidis VN, Song X, Manchanda S, et al. DRKG - drug repurposing knowledge graph for Covid-19[J]. <https://github.com/gnn4dr/DRKG/>, 2020.
- [8] Bordes A, Usunier N, Garcia-Duran A, et al. Translating embeddings for modeling multi-relational data[C]//Advances in Neural Information Processing Systems. Curran Associates, Inc., 2013. [https://proceedings.neurips.cc/paper\\_files/paper/2013/hash/1cecc7a77928ca8133fa24680a88d2f9-Abstract.html](https://proceedings.neurips.cc/paper_files/paper/2013/hash/1cecc7a77928ca8133fa24680a88d2f9-Abstract.html).
- [9] Yang BS, Yih S, He XD, et al. Embedding entities and relations for learning and inference in knowledge bases[C]//Proceedings of ICLR. 2015. <http://arxiv.org/abs/1412.6575>.
- [10] Trouillon T, Welbl J, Riedel S, et al. Complex embeddings for simple link prediction[C]//Proceedings of the 33rd International Conference on International Conference on Machine Learning. Journal of Machine Learning Research, 2016, 48: 2071-2080. <https://arxiv.org/abs/1606.06357>.
- [11] Sun ZQ, Deng ZH, Nie JY, et al. RotatE: knowledge graph embedding by relational rotation in complex space[C]//Proceedings of

ICLR. 2019. <https://openreview.net/forum?id=HkgEQnRqYQ>.

[12] Zeng XX, Song X, Ma TF, et al. Repurpose open data to discover therapeutics for COVID-19 using deep learning[J]. *J Proteome Res*, 2020, 19(11): 4624-4636. <https://doi.org/10.1021/acs.jproteome.0c00316>.

[13] Zhang R, Hristovski D, Schutte D, et al. Drug repurposing for COVID-19 via knowledge graph completion[J]. *J Biomed Inform*, 2021, 115(1): 103696. <https://doi.org/10.1016/j.jbi.2021.103696>.

[14] 李宗贤. 基于知识图谱的帕金森病药物重定位[J]. *信息技术与信息化*, 2022, No.268(07): 28-32. <https://doi.org/10.3969/j.issn.1672-9528.2022.07.006>.

[15] Han X, Cao SL, Lv X, et al. OpenKE: an open toolkit for knowledge embedding[C]//*Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*. Brussels, Belgium: Association for Computational Linguistics, 2018: 139-144. <https://aclanthology.org/D18-2024/>.

[16] Maaten LVD, Hinton G. Visualizing data using t-SNE[J]. *Journal of Machine Learning Research*, 2008, 9(86): 2579-2605. <http://jmlr.org/papers/v9/vandemaaten08a.html>.

[17] Zheng Da, Song X, Ma C, et al. DGL-KE: training knowledge graph embeddings at scale[C]//*Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*. New York, NY, USA: Association for Computing Machinery, 2020: 739-748. <https://arxiv.org/abs/2004.08532>.

[18] Sliwinska S, Jeziorek M. The role of nutrition in Alzheimer's disease[J]. *Rocz Panstw Zakl Hig*, 2021, 72(1): 29-39. <https://doi.org/10.32394/rpzh.2021.0154>.

[19] Koppel J, Jimenez H, Adrien L, et al. Haloperidol inactivates AMPK and reduces tau phosphorylation in a tau mouse model of Alzheimer's disease[J]. *Alzheimers Dement*, 2016, 2(2): 121-130. <https://doi.org/10.1016/j.trci.2016.05.003>.

[20] Pasiński M, Szulczyk B. Beneficial effects of capsaicin in disorders of the central nervous system[J]. *Molecules*, 2022, 27(8): 2484. <https://doi.org/10.3390/molecules27082484>.

[21] Zu GX, Sun KY, Li L, et al. Mechanism of quercetin therapeutic targets for Alzheimer disease and type 2 diabetes mellitus[J]. *Sci Rep*, 2021, 11(1): 22959. <https://doi.org/10.1038/s41598-021-02248-5>.

[22] Sahab-Negah S, Hajali V, Moradi HR, et al. The impact of estradiol on neurogenesis and cognitive functions in Alzheimer's disease[J]. *Cell Mol Neurobiol*, 2020, 40(3): 283-299. <https://doi.org/10.1007/s10571-019-00733-0>.

[23] Huang CW, Rust NC, Wu HF, et al. Altered O-GlcNAcylation and mitochondrial dysfunction, a molecular link between brain glucose dysregulation and sporadic Alzheimer's disease[J]. *Neural Regen Res*, 2023, 18(4): 779-783. <https://doi.org/10.4103/1673-5374.354515>.

[24] Reinhardt S, Stoye N, Luderer M, et al. Identification of disulfiram as a secretase-modulating compound with beneficial effects on Alzheimer's disease hallmarks[J]. *Sci Rep*, 2018, 8(1): 1329. <https://doi.org/10.1038/s41598-018-19577-7>.

- [25] Trinh PNH, Baltos JA, Hellyer SD, et al. Adenosine receptor signalling in Alzheimer's disease[J]. *Purinergic Signal*, 2022, 18(3): 359-381. <https://doi.org/10.1007/s11302-022-09883-1>.
- [26] U. S. Food and Drug Administration. FDA drug safety communication: FDA recommends against the continued use of Meridia (sibutramine)[EB/OL]. (2010-10-08)[2018-02-06]. <https://www.fda.gov/drugs/drug-safety-and-availability/fda-drug-safety-communication-fda-recommends-against-continued-use-meridia-sibutramine>.
- [27] Ai PH, Chen S, Liu XD, et al. Paroxetine ameliorates prodromal emotional dysfunction and late-onset memory deficit in Alzheimer's disease mice[J]. *Transl Neurodegener*, 2020, 9(1): 18. <https://doi.org/10.1186/s40035-020-00194-2>.
- [28] Lehrer S, Rheinstein PH. Transspinal delivery of drugs by transdermal patch back-of-neck for Alzheimer's disease: a new route of administration[J]. *Discov Med*, 2019, 27(146): 37-43. <https://pubmed.ncbi.nlm.nih.gov/30721650/>.
- [29] Baraka A, ElGhotny S. Study of the effect of inhibiting galanin in Alzheimer's disease induced in rats[J]. *Eur J Pharmacol*, 2010, 641(2): 123-127. <https://doi.org/10.1016/j.ejphar.2010.05.030>.
- [30] Chadwick W, Mitchell N, Carroll J, et al. Amitriptyline-mediated cognitive enhancement in aged 3×Tg Alzheimer's disease mice is associated with neurogenesis and neurotrophic activity[J]. *PLoS One*, 2011, 6(6): e21660. <https://doi.org/10.1371/journal.pone.0021660>.