

数据

进展

将四张关于药物的英文描述转化为向量， 向量维度为 384。

dataset1-embed

#	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W
drug_ID	name	text	remark	embeddings																			
D801296	Glucosam	Glucosamine is a col		[-0.020988]35656249523, -0.008754121139645576, 0.0420364402245453, 0.00309414928779006, 0.020215753465890884, 0.04261597990989685, 0.05394284427165985, 0.05138896405696869, -0.02740561																			
D809230	Azelinidip	Azelinidipine is a dthi		[-0.039345961064100266, 0.012494326569139957, 0.0172696802765131, 0.0014667290199920535, -0.00694510450810194, -0.0014920956455171108, 0.05455440282821655, 0.06688935309648514, -0.0018																			
D805012	Abiratorer	Abiraterone is an an		[-0.007257240624401379, -0.014835622161628616, -0.01233113700151443, 0.003621977539372444, 0.018908174386660057, -0.006891153621481228, 0.06691475449800491, 0.020731637254357338, -0																			
D801195	Flecainide	Flecainide is a class I		[-0.016198281198740005, -0.024711614474654198, 0.015311497263610363, -0.00821769330650568, 0.03201591596007347, -0.017098283635616302, 0.04555404558777809, 0.058438532054424296, 0.0072																			
D800201	Caffeine	Caffeine is a stimula		[0.0013237999519333243, -0.019451841711997986, 0.02181737683713436, 0.01887153834104538, 0.033295243978500366, 0.03812200203537941, 0.06601182371377945, 0.026519063860177994, -0.015121																			
D801020	Isonoribide	Isonoribide monitri		[-0.038880717009305954, -0.006452600471675396, -0.01203270722180605, -0.005468334071338177, 0.034212470054626465, 0.032815784215927124, 0.12399627268314362, 0.04473685473203659, -0.012																			
D800278	Argatroba	Argatroban is a synti		[-0.11045777052640915, 0.00297142518684268, -0.027401477098464966, -0.03348169103264809, 0.0032956646755337715, -0.006009288132190704, 0.07751692831516266, 0.0296223964521505, -0.0355																			
D801403	Methotrin	Methotrimprazine i		[-0.036387164145708084, -0.015518652275204659, 0.010442624102938175, -0.0017884275875985622, -0.008049173280596733, 0.01622184179723268, 0.0776083841919899, 0.016423484310507774, -0.0																			

dataset2-embed

#	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S
1	id	name	text	embeddings															
2	D801296	Glucosam	9101014108	[-0.035623881965875626, -0.03385728597640991, -0.011325006373226643, -0.02503310516476631, 0.017765864729881287, 0.01183082815259695, 0.0134650785475961															
3	D809230	Azelinidip	9101011112	[-0.030684709548950195, -0.015542080625891685, -0.008758020587265491, -0.009046725928783417, 0.03146114945411682, 0.005918900482356548, 0.0142298741266															
4	D805812	Abiratorer	9101011112	[-0.031172320246696472, -0.02019496262073517, -0.00864656176418066, -0.008801450952887535, 0.0346277616918087, 0.008282174374927094, 0.0154255889356136															
5	D801195	Flecainide	9101011112	[-0.03214859217405319, -0.016464434564113617, -0.00648346827098451, -0.010862707160413265, 0.03264220058917999, 0.0032184196170419455, 0.0188205726444															
6	D800201	Caffeine	9101011114	[-0.03287119418382645, -0.026214882731437683, -0.01654323935508728, -0.01743941754102707, 0.02061687596142292, 0.01269600261002779, 0.01687655970454216															
7	D801020	Isonoribide	9101014108	[-0.0379858821630478, -0.04304181418623734, -0.013344445265829563, -0.031228944659233093, 0.006323614157736301, 0.01980564557015896, 0.01826548576354															
8	D800278	Argatroba	9101011112	[-0.0315649492710865, -0.01667488974716187, -0.006394927168891907, -0.010945974849164486, 0.0192759473178067, 0.007473705408476544, 0.014811124653674															

dataset3-embed

#	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S
1	name	text	embeddings																
2	Abemacic Drug Abem		[-0.07470092922449112, -0.01814381219446659, -0.0062376041896641254, -0.008868244476616383, -0.03253835439682007, -0.02478274144232273, 0.0022831058595329523, 0																
3	Dabrafenil Drug Dab		[-0.06942513585090637, 0.005917823873460293, -0.020922495052218437, -0.010210166685283184, -6.933228632988175e-06, -0.01423491071909666, 0.0348786786198616, 0.0																
4	Dasatinib Drug Das		[-0.05765390396118164, -0.0016053157160058618, -0.02879844605922699, -0.01372817438095808, -0.0043209572322666645, 0.0023034221958369017, 0.03605161979794502, 0																
5	Dihydroerol Drug Dih		[-0.06774567067623138, 0.002911383518949151, -0.015215851366519928, -0.0010666117304936051, 0.003946242853999138, 0.00030461576534435153, 0.04329787194728851, 0																
6	Diltiazem Drug Dilti		[-0.06650730967521667, -0.001232348964549601, 0.009667130187153816, 0.039643172174692154, -0.0076605756767094135, -0.01059085877459588, 0.03255016356706619, 0.0																
7	Doxorubic Drug Doxi		[-0.06357309222221375, -0.00442147161815902, -0.0027675717137753963, -0.002237686887383461, -0.0036964244209229946, -0.0151503945508599281, 0.02871409989893436																
8	Doxorubic Drug Doxi		[-0.06357309222221375, -0.00442147161815902, -0.0027675717137753963, -0.002237686887383461, -0.0036964244209229946, -0.0151503945508599281, 0.02871409989893436																

dataset4-embed

#	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V
name	text	embeddings																				
Glucosam	Glucosam	[-0.040664125233888626, -0.01442786306142807, -0.032501090317964554, -0.02311810292303562, -0.052117906510829926, -0.005631760694086552, 0.021722840145230293, 0.00711920578032732, -0.03343564																				
Azelinidip	Azelinidip	[-0.06278902292251587, -0.011920550838112831, -0.021952062845230103, -0.022612914443016052, -0.03965534642338753, -0.0013634719653055072, 0.0314444899559021, 0.003527092980220914, -0.02201487																				
Abiratorer	Abiratorer	[-0.04875269904732704, -0.006710604764521122, -0.036351755261421204, -0.011774021200835705, -0.03725462406873703, -0.004845349118113518, 0.016354622319340706, -0.005023746285587549, -0.014667																				
Caffeine	Caffeine	[-0.04520631954077906, -0.02227089063955307, -0.019181381911039352, -0.008428539149463177, -0.021478816966874695, -0.017831066623330116, 0.022584386363529472, 0.01128732692450285, 0.00907742																				
Isonoribide	Isonoribide	[-0.045787370611429214, -0.017566561698913574, -0.020353883504867554, -0.02425614558160305, -0.03885207325220108, 0.0002876250073313713, 0.03084552101790905, 0.01408976037055254, 0.010410354																				
Argatroba	Argatroba	[-0.0595758892595768, -0.031037576496601105, -0.03870357573032379, -0.031071174889802933, -0.031674452126026154, -0.011135323904454708, 0.03469545315317056, 0.01272706874054009, -0.0357176251																				
		[-0.06970905512571335, 0.0046306937001645565, -0.045629214495420456, -0.04207830876111984, -0.045688338577747345, -0.01238231174647808, 0.02444538101553917, 0.00032224846654571593, -0.0285039																				

问题

1、dataset3-embed (药物 ddi 信息) 和 dataset4-embed (药物 MACCS 键和物理特性信息) 还没找到合适的英文描述， 还需要找找， 或由 LLM 生成。暂时先用着。

代码

进度

```

class GNN3(nn.Module): 3用法
    def __init__(self, dataset, tail_len, relation_len, args, dict1, drug_name, **kwargs):
        super(GNN3, self).__init__(**kwargs)
        self.text_embed_3 = self.load_pretrain_embeds(dict1, embed_path: '../data-text-embed/dataset3-embed.csv')
        self.kg, self.dict1 = dataset["dataset3"], dict1
        self.drug_name, self.args = drug_name, args
        self.drug_embed = nn.Embedding(num_embeddings=572, embedding_dim=args.embedding_num)
        self.rela_embed = nn.Embedding(num_embeddings=67, embedding_dim=args.embedding_num)
        self.ent_embed = nn.Embedding(num_embeddings=572, embedding_dim=args.embedding_num)
        self.W1 = nn.Parameter(torch.randn(size=(572, args.embedding_num, args.embedding_num)))
        self.b1 = nn.Parameter(torch.randn(size=(args.neighbor_sample_size, args.embedding_num)))
        self.W2 = nn.Parameter(torch.randn(size=(572, args.embedding_num, args.embedding_num)))
        self.b2 = nn.Parameter(torch.randn(size=(args.neighbor_sample_size, args.embedding_num)))
        self.Linear1 = nn.Sequential(nn.Linear(args.embedding_num * 2, args.embedding_num),
                                     nn.ReLU(),
                                     nn.BatchNorm1d(args.embedding_num))

        self.relu = nn.ReLU()
        self.soft = nn.Softmax(dim=1)
        for m in self.modules():
            if isinstance(m, nn.Linear):
                nn.init.xavier_uniform_(m.weight)
                nn.init.zeros_(m.bias)
            elif isinstance(m, nn.BatchNorm1d):
                nn.init.constant_(m.weight, 1)

```

在每个GNN中输入文本向量，并用dict1与GNN输出的向量对齐

```

182 class GNN3(nn.Module): 3用法
207
208     def forward(self, arguments):
209         kg, dict1, args, drug_name = self.kg, self.dict1, self.args, self.drug_name
210         gnn2_embedding, gnn1_embedding, idx = arguments
211         text_embed_3 = self.text_embed_3
212         adj_tail, adj_relation = self.arrge(kg, dict1, args.neighbor_sample_size)
213
214         drug_name = torch.LongTensor(drug_name)
215         adj_tail = torch.LongTensor(adj_tail)
216         adj_relation = torch.LongTensor(adj_relation)
217
218         drug_embedding = self.drug_embed(drug_name)
219         rela_embedding = self.rela_embed(adj_relation)
220         ent_embedding = self.ent_embed(adj_tail)
221         drug_rel = drug_embedding.reshape((572, 1, args.embedding_num)) * rela_embedding
222         drug_rel_weigh = drug_rel.matmul(self.W1) + self.b1
223         drug_rel_weigh = self.relu(drug_rel_weigh)
224         drug_rel_weigh = drug_rel_weigh.matmul(self.W2) + self.b2
225         drug_rel_score = torch.sum(drug_rel_weigh, axis=-1, keepdims=True)
226         drug_rel_score = self.soft(drug_rel_score)
227         weighted_ent = drug_rel_score.reshape((572, 1, args.neighbor_sample_size)).matmul(ent_embedding)
228         drug_e = torch.cat(tensors=[weighted_ent.reshape(572, args.embedding_num), drug_embedding.reshape((572, args.embedding_num))], dim=1)
229         drug_f = self.Linear1(drug_e)
230         drug_f_text_embed = torch.cat(tensors=[drug_f, text_embed_3], dim=1)
231         return drug_f_text_embed, gnn2_embedding, gnn1_embedding, idx

```

拼接，drug_f的形状为 (572, 128)
text_embed_3的形状为 (572, 384)

```

class FusionLayer(nn.Module): 2用法
    def __init__(self, args):
        super().__init__()
        self.fullConnectionLayer = nn.Sequential(
            # (args.embedding_num + 384)
            nn.Linear((args.embedding_num + 384) * 4 * 2, (args.embedding_num + 384) * 4),
            nn.ReLU(),
            nn.BatchNorm1d((args.embedding_num + 384) * 4),
            nn.Dropout(args.dropout),
            nn.Linear((args.embedding_num + 384) * 4, (args.embedding_num + 384) * 2),
            nn.ReLU(),
            nn.BatchNorm1d((args.embedding_num + 384) * 2),
            nn.Dropout(args.dropout),
            nn.Linear((args.embedding_num + 384) * 2, out_features=65))
        for m in self.modules():
            if isinstance(m, nn.Linear):
                nn.init.xavier_uniform_(m.weight)
                nn.init.zeros_(m.bias)
            elif isinstance(m, nn.BatchNorm1d):
                nn.init.constant_(m.weight, val=1)
                nn.init.constant_(m.bias, val=0)

```

修改融合层的维度

问题

1、向量拼接已完成，还缺少向量融合

```

class FusionLayer(nn.Module): 2用法
    def forward(self, arguments):
        idx = idx.numpy().tolist()
        drugA = []
        drugB = []
        for i in idx:
            drugA.append(i[0])
            drugB.append(i[1])

        Embedding = torch.cat(
            tensors=[gnn1_embedding[drugA], gnn2_embedding[drugA], gnn3_embedding[drugA], gnn4_embedding[drugA],
                     gnn1_embedding[drugB], gnn2_embedding[drugB], gnn3_embedding[drugB], gnn4_embedding[drugB]]
        )

        return self.fullConnectionLayer(Embedding)

```

在该部分添加LLM融合向量部分

打算使用 Llama-3.2-1B 来完成融合


[ModelScope](#)

[首页](#)
[模型库](#)
[数据集](#)
[创空间](#)
[AIGC专区](#)
[文档中心](#)



Llama-3.2-1B

[LLM-Research / Llama-3.2-1B](#)

[文本生成](#)
[PyTorch](#)
[Transformers](#)
[Safetensors](#)
[Llama](#)
[开源协议: llama3.2](#)
[英语](#)

[@LLM-Research 提供](#)
[50,527 下载](#)
[2024-10-25 更新](#)

流程

向量拼接部分已完成，后续再冻结 LLM 参数来融合向量。