



基于贩毒案情图谱的辅助量刑研究

王治政¹, 王雷², 李帅驰¹, 孙媛媛¹, 陈彦光¹, 许策¹, 王刚¹, 林鸿飞¹

¹大连理工大学计算机科学与技术学院

²辽宁省锦州市人民检察院



大连理工大学

信息检索研究室

Information Retrieval Laboratory of DUT

搜人搜物搜信息 重情重义重认知

1

背景介绍

2

基于知识图谱的量刑预测

3

多视角知识图谱嵌入

4

实验分析与讨论

5

总结与展望

1

背景介绍

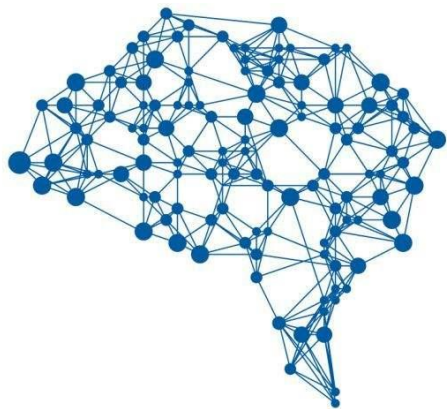
□ 智慧司法

- 人工智能技术与司法量刑的结合日益紧密
- 量刑结果是多种因素综合作用的结果
- 学习从法律文书到刑罚裁量的判决模式



□ 深度学习

- 端到端的方式学习案件事实的多源特征
- 减少人工成本但需要海量案例数据训练模型
- 判决模式对办案人员不可见，可解释性较弱



□ 司法知识图谱

- 描述案件要素、要素属性及案件要素间的关系
- 提供先验规则，支撑量刑模式的可解释性
- 保证量刑规则的动态扩展



□ 知识图谱嵌入

- 为量刑预测提供技术支持
- 兼容文本信息，符合司法业务的需求

基于知识图谱的量刑预测

2

基于知识图谱的量刑预测

基于知识图谱的量刑预测



□ 量刑预测重定义

✓ 基于分类模型的量刑预测

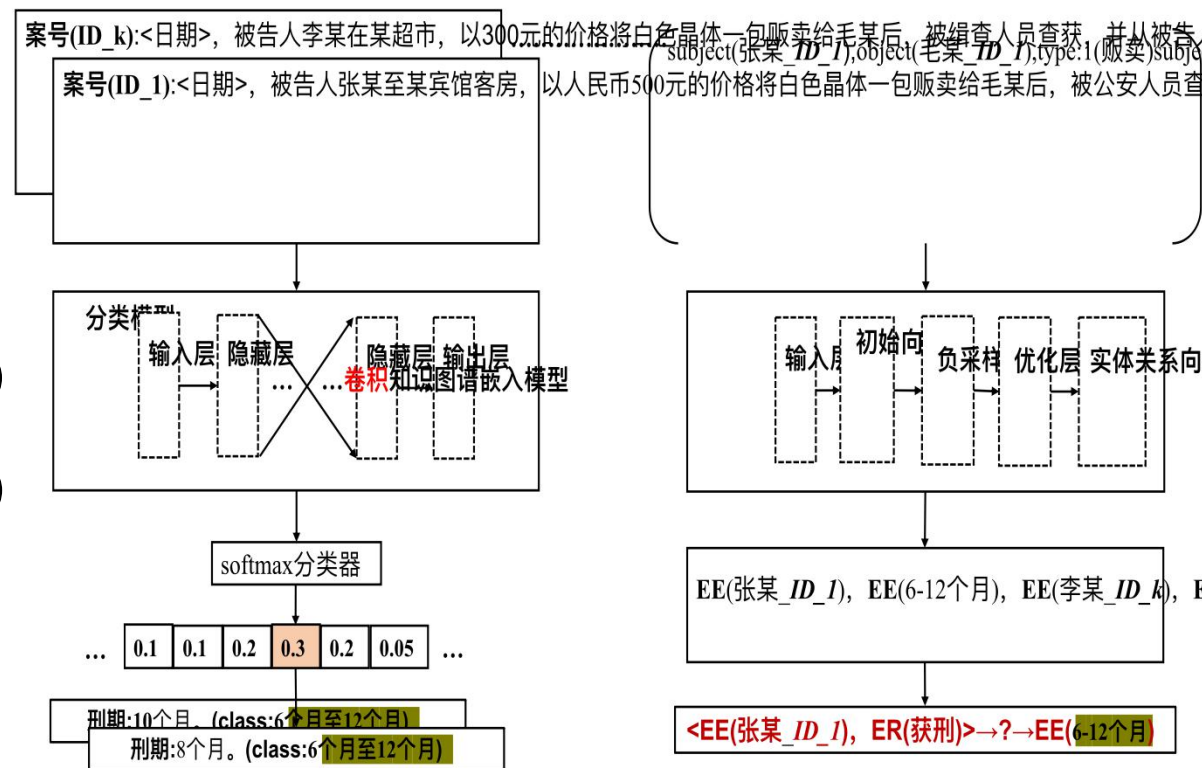
- 输入端: 案件事实的**自然文本**集合 ($T = \{t_1, t_2, \dots, t_n\}$)
- 输出端: 案件事实对应的**量刑区间** ($Y = \{y_1, y_2, \dots, y_k\}$)
- 中间层: 通过卷积操作学习**文本特征** $f(T)$ 。该特征经

softmax后可以得到案件 t_i 的量刑 y_i 的概率

✓ 基于知识图谱的量刑预测

- 输入端: 从案件事实中提取的**三元组**集合 $Tr = \{tr_1, tr_2, \dots, tr_m\}$
- 输出端: 三元组 tr_i 的**量刑尾实体** (也是一种量刑区间)
- 中间层: 通过对**头实体向量和关系向量**进行二元操作, 计算当前的量刑

实体属于目标头实体的概率



基于知识图谱的量刑预测



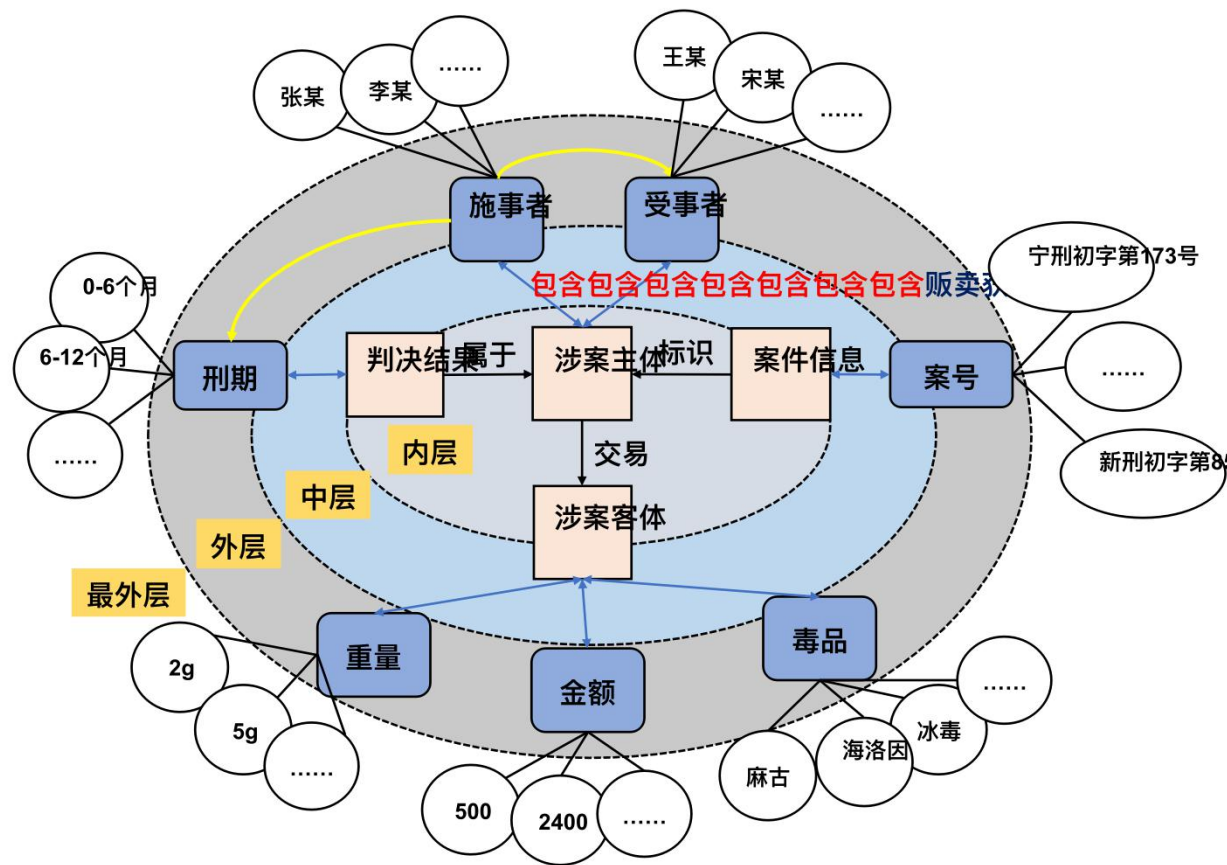
□ 本体设计

✓ 概念设计

- 涉案主体：施事者和受事者
- 涉案客体：重量、金额和毒品
- 判决结果：刑期
- 案件信息：案号信息

✓ 关系设计

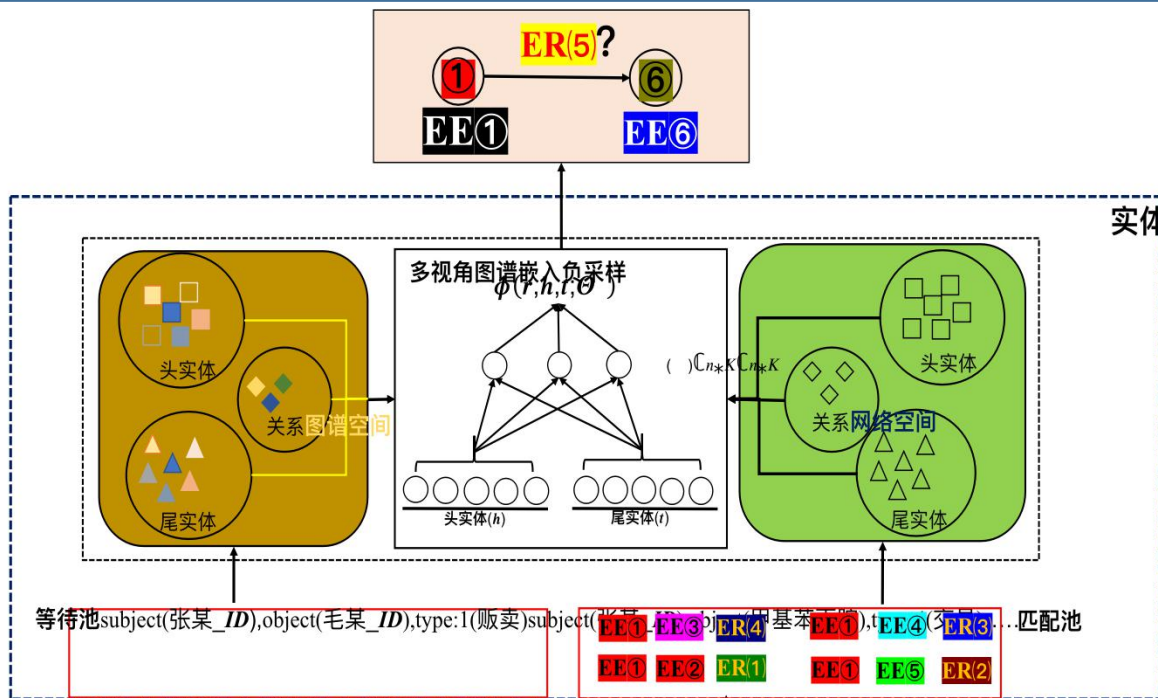
- 一阶关系（内层）：单向元关系类型，定义顶层概念间的联系
- 二阶关系（中层）：双向过渡关系，定义底层概念和顶层概念之间的扩展关系
- 三阶关系（外层）：单向实体关系，定义底层概念之间的具体关系



3

多视角知识图谱嵌入

多视角知识图谱嵌入

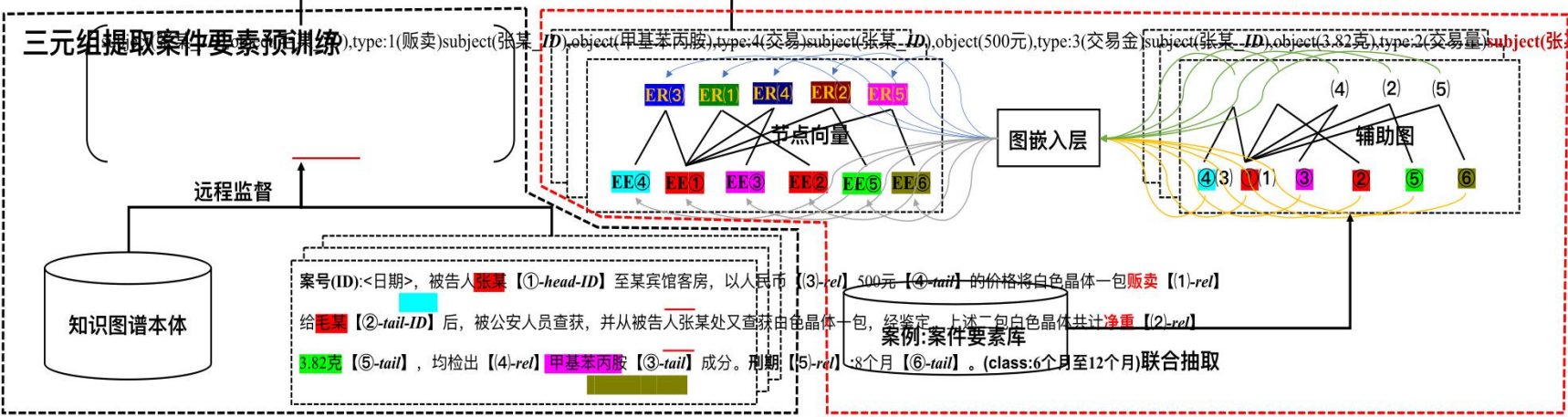


➤ 结合远程监督和知识图谱本体抽取

案情三元组

➤ 构建实体和关系的辅助图学习案件要素的初始表示

➤ 根据Complex算法设计多视角的案件要素再表示方法

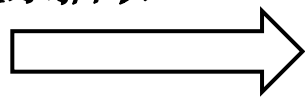


□ 三元组提取

对于知识图谱中的某个三元组而言，如果外部语料有包含该三元组中的实体对的句子，那么这些句子都可以反映该三元组中的关系

远程监督

关系抽取



三元组提取

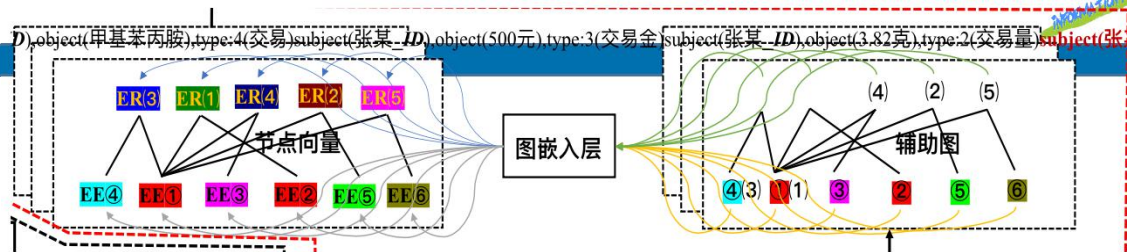
在外部语料中，只有那些包含本体概念的句子才可能反映概念间的具体关系类型。从语料中筛选出句子后，使用联合抽取的方式提取案件要素，以此构建模型的输入三元组

头实体	尾实体	关系	所属元关系
施事者	受事者	贩卖	无
施事者	毒品	交易	交易
施事者	金额	交易金	交易
施事者	重量	交易量	交易
施事者	刑期	获刑	属于

多视角知识图谱嵌入



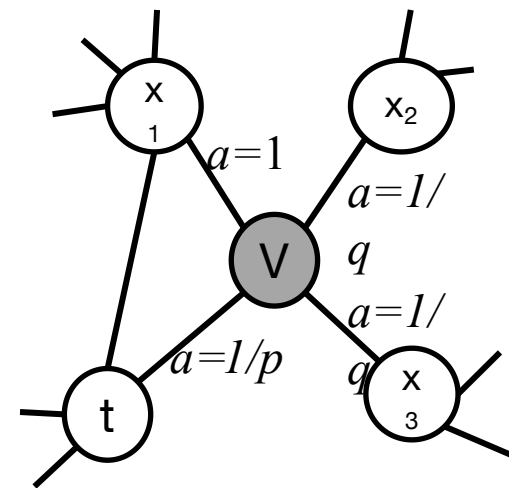
□ 案件要素预训练



- 三元组是案件要素的一种关系型结构表示，反映了案件要素在知识图谱中的局部信息
- 提供案件要素的全局依存关系，将案件要素组织成辅助二部图 $G = (V, E)$

$$\max_f \sum_{u \in V} I_g P_r(\mathcal{N}_s(u) | f(u))$$

- 目标是学习映射函数 f ，将案件要素从网络空间映射到向量空间
- 通过案件要素采样来构建当前案件要素 u 的邻域信息 $\mathcal{N}_s(u)$
- 最大化 $\mathcal{N}_s(u)$ 与 u 产生连边的条件概率实现对 u 的向量表示 ($\mathbb{R}^{V \times d}$)



node2vec

node2vec所得的案件要素向量并不区分实体和关系，为统一的 **节点** 表示

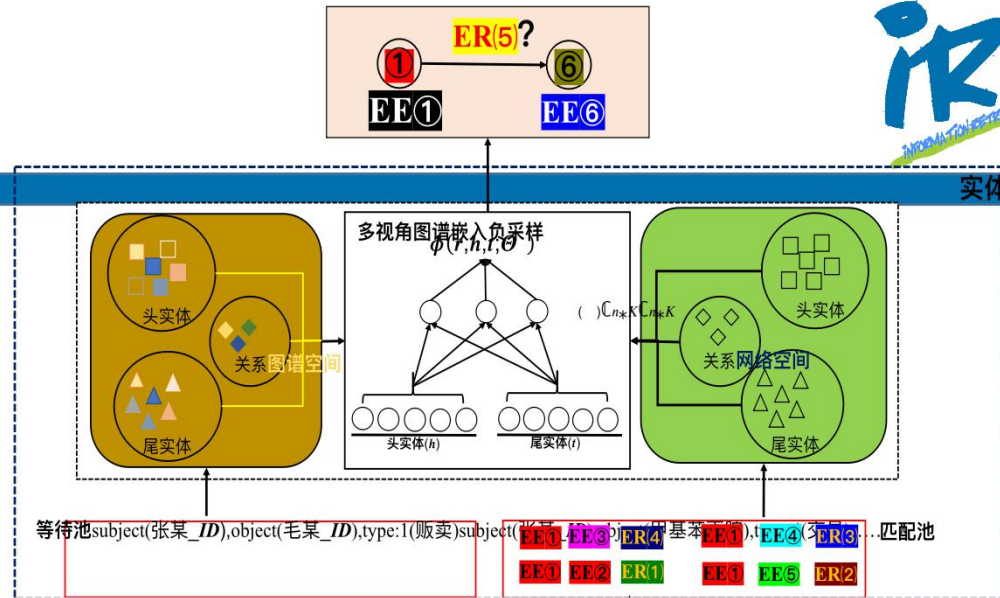


多视角知识图谱嵌入



□ 实体和关系向量再表示

✓ 初始向量



- 匹配池 ($V \times d$) : 包含 V 个 d 维向量的案件要素, 其中 V 是由实体和关系组成的
- 等待池 ($k \times b$) : 包含 $k \subset Tr$ 个正例三元组和为每个三元组采样的 b 个负例三元组

$$h = EE(V_h^{D_{en}}); t = EE(V_t^{D_{en}}); r = ER(V_r^{D_{re}})$$

在“匹配池”中, 实体编码从零开始, 关系编码从实体编码的结束位置开始。以查表的方式将“匹配池”中的实体和关系向量映射到“等待池”中的三元组, 得到三元组中实体和关系的初始化向量

多视角知识图谱嵌入



□ 实体和关系向量再表示

✓ 向量再表示

关系类型均为非对称关系



Complex知识图谱嵌入方法

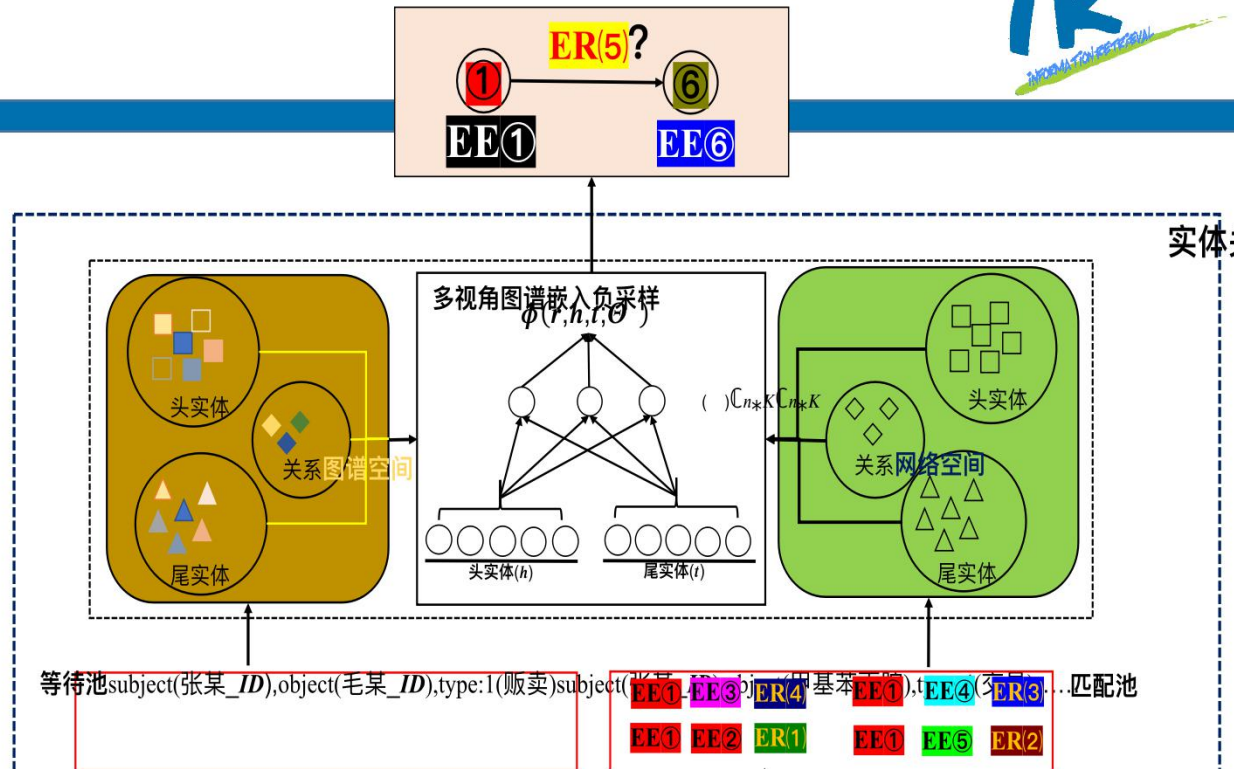


1、实体和关系向量映射到复数空间

$$\begin{matrix} \square & \square & \square & \square & \square & \square & \square & \square \\ h = h + ih; & t = t + it; & r = r + ir \end{matrix}$$

2、Complex算法优化

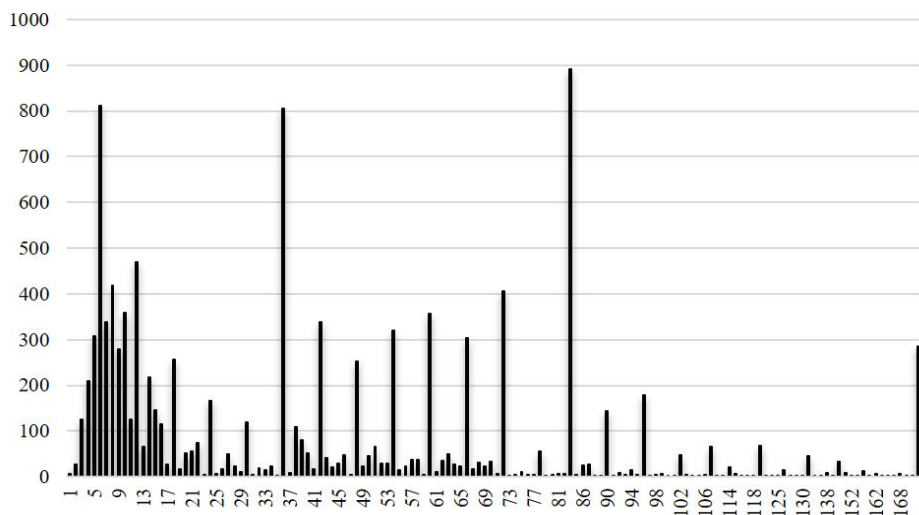
$$\begin{aligned} \phi(r, h, t; \Theta) &= \text{Re}(\langle w_r^{\square}, e_h^{\square}, e_t^{\square} \rangle) \\ &= \text{Re}(\sum_{k=1}^K w_{rk}^{\square}, e_{hk}^{\square}, e_{tk}^{\square}) \end{aligned}$$



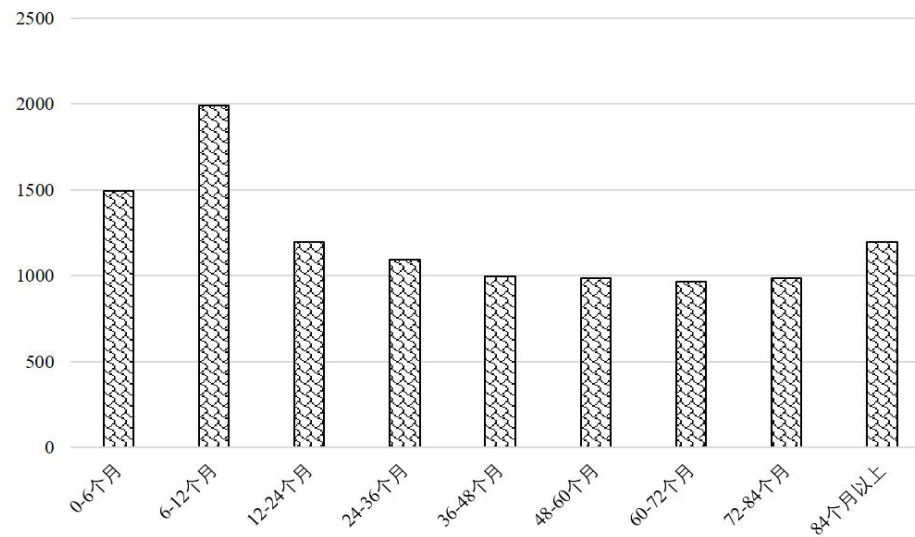
4 实验分析与讨论

□ 数据集

- 从**186,912**份涉毒类刑事案件判决书筛选获得**10,913**份符合要求的贩卖毒品类案件
- 抽取**50,684**个三元组。从所有指明量刑结果的三元组中随机抽取**2,000**个构成测试集，再抽取**2,000**个作为验证集，其余的**46,684**个三元组为训练集
- 统计量刑标签，并划分量刑结果区间，分为**9**个类别



量刑标签数量分布



量刑区间数量分布

□ 实验设置

案件要素预训练

- 节点向量维度 (d_1) : $d_1 = 200$
- 采样长度 (L) : 构造案件要素的邻域信息 $\mathcal{N}_s(u)$ 的长度 $L = 100$, 即为案件要素 u 采样100个节点作为其上下文信息
- 训练轮数 ($epc\ h_1$) : node2vec算法使用语言模型对采样的案件要素序列进行训练, 其训练轮数 $epc\ h_1 = 10$

实体关系向量再表示

- 实体和关系向量维度 (d_2) : $d_2 = d_1 = 200$
- 计算批次 ($batch$) : 每个计算批次的大小与模型的“等待池”尺寸一致, 为 100×50 , 其中100为正例三元组的个数, 50为每个三元组的负采样个数
- 训练轮数 ($epc\ h_2$) : 模型训练阶段的迭代次数 $epc\ h_2 = 2000$

□ 实验结果

✓ 纵向对比 (知识图谱方法对比)

目的: 验证多视角图谱嵌入方法比其他知识图谱嵌入模型能学习到更丰富的实体和关系特征

Method	Hit@1	Hit@3	MRR
Random	0.112	0.340	0.244
TransE	0.172	0.431	0.376
TransH	0.164	0.419	0.365
TransD	0.159	0.431	0.366
TransR	<u>0.195</u>	0.434	0.387
HolE	0.120	0.415	0.347
DistMult	0.185	0.427	0.385
Simple	0.105	0.436	0.338
Complex	0.181	<u>0.442</u>	<u>0.396</u>
Ours	0.197	0.467	0.406
Ours-money	0.216	0.504	0.426

□ 实验结果

✓ 纵向对比 (结果分析)

Method	Hit@1	Hit@3	MRR
Complex	0.181	<u>0.442</u>	<u>0.396</u>
Ours	0.197	0.467	0.406
Ours-money	0.216	0.504	0.426

- 多视角知识图谱嵌入学习到了更丰富的**案件要素依赖特征** (全局特征)
- 从表中可以看出去掉“交易金”后的模型结果又分别提升了1.9%、3.7%和2%
 - 买卖双方的交易金额存在不确定性因素，它是由时间、地域和涉案人员主观动机等多种因素所决定的。同时在《中华人民共和国刑法（2017年修正）》第三百四十七条中明确规定了贩卖毒品的重量和种类，而对涉案金额并没有做出规定，因此交易金额在量刑中属于一种**弱特征**
 - 案件事实描述中对于交易金额的提及并不**全面**，因此抽取的金额可能只是涉案总金额的一部分

□ 实验结果

✓ 知识图谱 V.S 深度学习

目的: 验证基于知识图谱开展量刑预测任务的可行性, 同时为该任务的进一步研究提供可能的改进方案和思路

	Hit@1	Hit@3
多视角知识图谱嵌入	0.216	0.504

	Acc.	macro_F
TextRNN	0.202	0.136
TextRNN_Att	0.221	0.137
TextRCNN	0.223	0.139
DPCNN	0.234	0.164
TextCNN	0.282	0.220
FastText	0.271	0.210
Bert_Transformer	0.511	0.491

□ 实验结果

✓ 知识图谱 V.S 深度学习 (结果分析)

	Hit@1	Hit@3
多视角知识图谱嵌入	0.216	0.504
	Acc.	macro_F
Bert_Transformer	0.511	0.491

- 与使用词向量的文本分类方法**基本可比**，而与Bert相比差距较大
 - Bert利用大量的参数和更深的网络层数学习了文本语义的更**高阶表示**
 - Transformer在编码文本时使用注意力机制捕获了文本的**长距离依赖**
- 多视角图谱嵌入方法的Hit@3与Bert_Transformer模型是可比

5 总结与展望

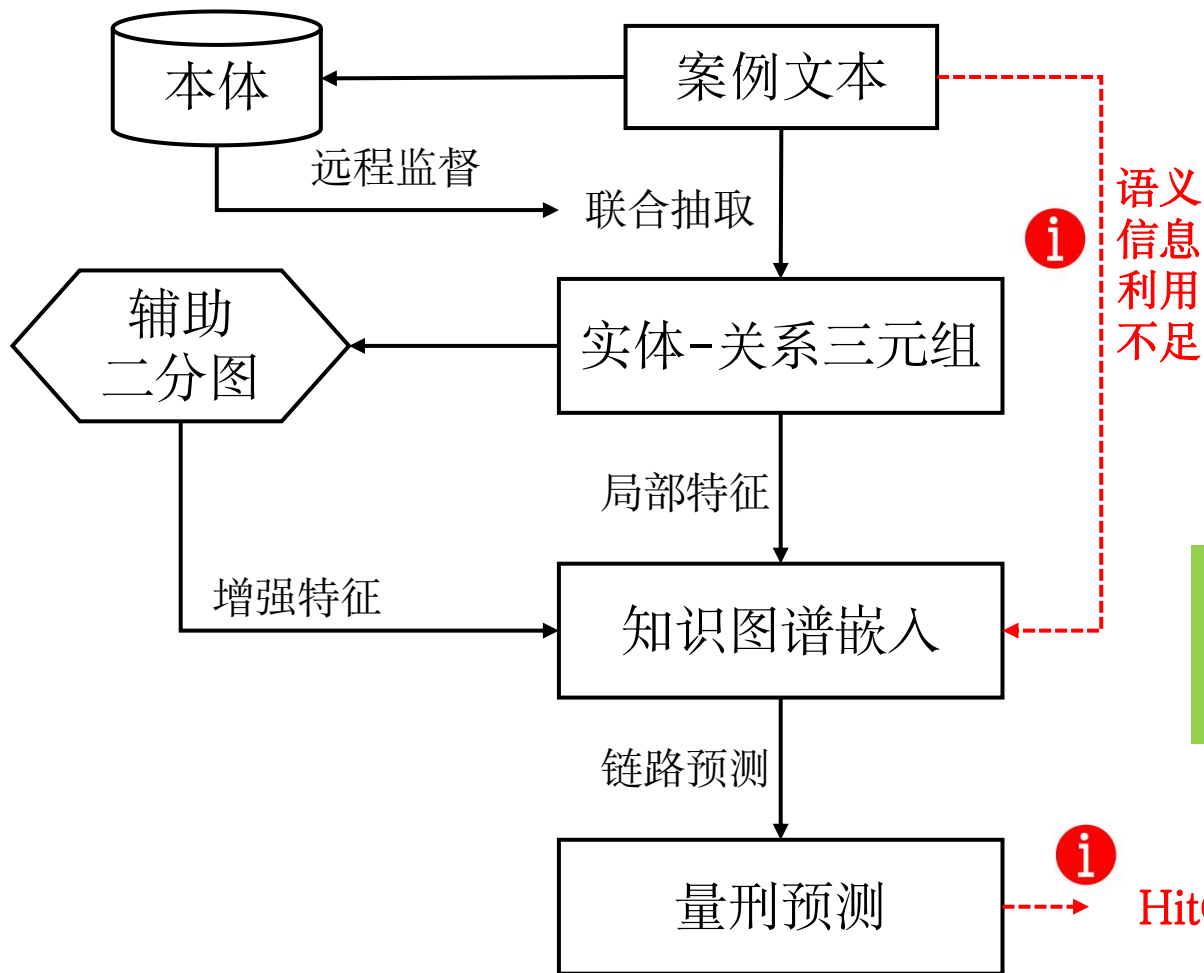
□ 主要贡献

- 将量刑预测与知识图谱的链路预测任务相结合，提出了基于知识图谱的量刑预测任务。使用知识图谱本体为辅助量刑提供先验规则，提升模型的可解释性，保证办案人员对量刑规则的掌控。
- 设计一种多视角的知识图谱嵌入方法，同时学习案件要素的网络特征（全局特征）和知识图谱的结构特征（局部特征），以提高量刑预测的准确性。
- 针对贩卖毒品案提出了一种知识图谱本体的设计规则。使用该规则并结合远程监督方法指导贩卖毒品罪的案情图谱构建，为基于知识图谱的量刑预测提供数据集

总结与展望



未来工作



1、引入**语言模型**为知识图谱中的实体和关系表示提供语义特征，增强案件文本之间的关联性，捕获案件的语义特征

2、引入**排序学习**的思想，实现模型对预测结果的偏重排序，将候选结果中的正确项前移

Hit@3与Bert_Transformer模型是可比



1. 谭红叶,张博文,张虎,李茹.面向法律文书的量刑预测方法研究[J].中文信息学报,2020,34(03):107-114.
2. 陈彦光, 刘海顺, 李春楠, 等. 基于刑事案例的知识图谱构建技术[J]. 郑州大学学报 (理学版), 2019, 51(3): 85-90.
3. Trouillon T, Welbl J, Riedel S, et al. Complex embeddings for simple link prediction[C]. International Conference on Machine Learning (ICML), 2016.
4. Mintz M, Bills S, Snow R, et al. Distant supervision for relation extraction without labeled data[C]// International Joint Conference on Acl. Association for Computational Linguistics, 2009.
5. Grover A, Leskovec J. node2vec: Scalable feature learning for networks[C]//Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining. 2016: 855-864.
6. Yao L, Mao C, Luo Y. KG-BERT: BERT for knowledge graph completion[J]. arXiv preprint arXiv:1909.03193, 2019
7. GitHub, <https://github.com/649453932/Chinese-Text-Classification-Pytorch>, last accessed 2020/7/22.

感谢聆听
请批评指正