

事实证明，在文本上预训练语言模型（LM）有助于各种下游 NLP 任务。最近的工作表明，知识图（KG）可以补充文本数据，提供结构化的背景知识，为推理提供有用的支架。然而，这些作品并未经过预先训练来大规模学习两种模式的深度融合，从而限制了获得文本和 KG 完全联合表示的潜力。在这里，我们提出了 DRAGON（深度双向语言知识图预训练），这是一种自监督方法，用于大规模地从文本和知识图谱中预训练深度联合的语言知识基础模型。具体来说，我们的模型将文本片段对和相关 KG 子图作为输入，并双向融合来自两种模式的信息。我们通过统一两个自监督推理任务（掩码语言建模和知识图谱链接预测）来预训练该模型。DRAGON 在各种下游任务（包括一般和生物医学领域的问答）上优于现有 LM 和 LM+KG 模型，平均绝对增益 +5%。特别是，DRAGON 在有关语言和知识的复杂推理（涉及长上下文或多步骤推理的问题上+10%）和低资源 QA（在 OBQA 和 RiddleSense 上+8%）以及新的状态方面取得了出色的表现在各种 BioNLP 任务上取得了最先进的结果。我们的代码和经过训练的模型可 <https://github.com/michiyasunaga/dragon> 上获取。