# Query Embedding Strategy

Connor Henderson

February 2021

# 1 Sub-graph Generation

---

**Algorithm 1:** $get\_subgraphs(\mathcal{G}, n)$

---

**input** : $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ : Directed acyclic graph from which sub-graphs have
        to be generated
        $n$ : Node which acts as the root of the sub-graph

**output** : $joins$ : A list of tuples representing the join node and the
        tables involved with that join
        $terminals$ : A list representing the tables present in a graph

$n := \text{root}$
$joins := [\,]$
$terminals := [\,]$

**do**
$\quad\big|\quad \mathcal{N}_n := \{\ n' \mid (n, n')\ \in\ \mathcal{E}\ \}$
$\quad\big|\quad n := n'$
**while** $len(\mathcal{N}_n) == 1$;

$\mathcal{N}_n := \{\ n' \mid\ (n, n')\ \in\ \mathcal{E}\ \}$

**if** $len(\mathcal{N}_n) > 1$ **then**
$\quad\big|\quad$ **foreach** $n' \in \mathcal{N}_n$ **do**
$\quad\big|\quad\big|\quad joins',\ terminals' = get\_subgraphs(\mathcal{G},\ n')$
$\quad\big|\quad\big|\quad joins := joins \cup \text{joins'}$
$\quad\big|\quad\big|\quad terminals := terminals \cup \text{terminals'}$
$\quad\big|\quad$ **end**
$\quad\big|\quad joins := joins \cup [(\text{n, terminals})]$
**end**

**if** $len(\mathcal{N}_n) == 0$ **then**
$\quad\big|\quad terminals := terminals \cup [\text{n}]$
**end**

**return** $joins, terminals$

---

2

# 2 Embedding

---

**Algorithm 2:** *embedding_process*

---

**input** : $\mathbb{G} = \{G_1, G_2, ..., G_n\}$ : Set of graphs for which embeddings are desired

$\quad\quad\quad D$: The maximum number of tables to consider for a sub-graph

$\quad\quad\quad \delta$: number of dimensions in the latent space

$\quad\quad\quad \epsilon$: number of training epochs

$\quad\quad\quad \alpha$ : learning rate

**output** : $\Phi \in \mathbb{R}^{\|\mathbb{G}\| \times \delta}$ : Matrix of graph representation vectors

Initialize $\Phi$ with random values

**for** $e \leftarrow 0$ **to** $\epsilon$ **do**

$\quad$ **foreach** $G_i \in \mathbb{G}$ **do**

$\quad\quad$ sg$_n$ = get_subgraphs(G$_i$, 0)

$\quad\quad$ $J(\Phi) := -log Pr(sg_n \mid \Phi(G_n))$

$\quad\quad$ $\Phi = \Phi - \alpha \frac{\partial J}{\partial \Phi}$

$\quad$ **end**

**end**

**return** $\Phi$

---