# GraphSAGE Tutorial

Feb 18, 2021

# Table of Contents

# Installation

# Steps

1. git clone [https://github.com/williamleif/GraphSAGE](https://github.com/williamleif/GraphSAGE)
2. (optional) create a virtual environment
   a. pip install --user virtualenv
   b. python -m venv [env_name]
3. pip install -r requirements.txt
   a. If using python3:
   b. Open requirements.txt
   c. Replace futures==3.2.0 with futures==3.1.1

```
8 futures==3.2.0  ⟶  8 futures==3.1.1
```

# Guides

# Guides

1. https://github.com/williamleif/GraphSAGE (README)
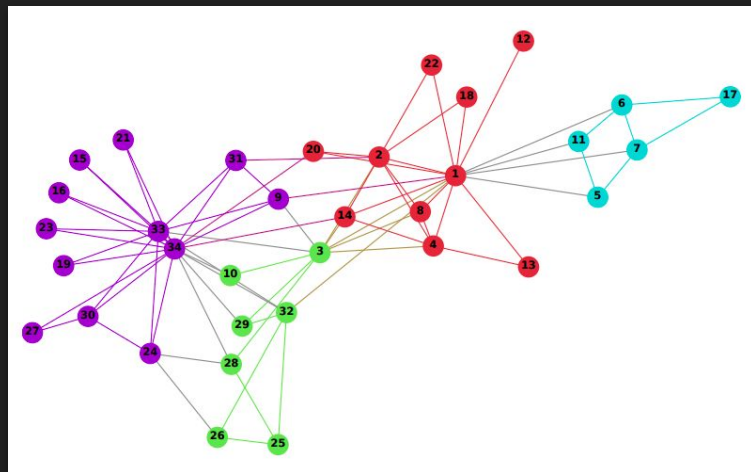   a. How to install
   b. How to use
   c. Models description
2. Graphsage/example_data
   a. Example of the required data format
3. Graphsage/example_supervised.sh & Graphsage/example_unsupervised.sh
   a. Example command to run the training and evaluation code in grarphsage/supervised_train.py

# Data

# Karate Toy Data

- 34 nodes
- Node classification (4 classes)
- Node attributes (toy-data only):
  - Last name
  - Age
- Neighboring nodes:
  - High chance of same last name (e.g. siblings)
  - High change of similar age (±5)

# Preparing Data

Files: (data/prepare_data.py)

```
 6 # Required:
 7 # 1. Load the graph
 8 # 2. Add validation and testing attributes to the graph
 9 # 3. Create the features
10 # 4. Save the data in the proper format
```

```
35 # Prepare graph
36 nx.set_node_attributes(G, 'val', False)
37 nx.set_node_attributes(G, 'test', False)
38
39 val = ['32', '17', '6']
40 test = ['19', '10', '10', '6']
41
42 for node in val:
43     G.node[node]['val'] = True
44
45 for node in test:
46     G.node[node]['test'] = True
47
```

# Generating Features

In this case, using:

1. Name
   a. One-hot vector
2. Age
   a. Integer
3. Node degree
   a. Integer

```python
48 # Generate features
49 feats = []
50
51 names_dict = {name: idx for idx, name in enumerate(set(names.values()))}
52 # print(names_dict)
53
54 for node in G:
55     name = names[node]
56     name_arr = [0]*len(names_dict)
57     name_arr[names_dict[name]] = 1
58     # 7-dim
59     # print(name_arr)
60
61     age = ages[node]
62     degree = G.degree()[node]
63
64     feat = name_arr + [age] + [degree]
65     # 9-dim
66     # print(feat)
67
68     feats.append(feat)
```

# Writing Data

Write the processed data in the required format (json or npy).

All files should include the same prefix.

```
71 # Writing:
72 # Required files:
73 # 1. -G.json
74 # 2. -id_map.json
75 # 3. -class_map.json
76 # 4. -feats.npy
```

# Writing Data

1. Graph written as json
2. ID map as json
   a. It is required because nodes can be indexed by strings instead of integers.
3. Class map as json
4. Features as numpy array

```python
81 # 1. Write graph
82 data = json_graph.node_link_data(G)
83 with open('processed/karate-G.json', 'w') as f:
84     json.dump(data, f)
85
86 # 2. Write id_map
87 id2idx = {n: idx for idx,n in enumerate(G.nodes())}
88 with open('processed/karate-id_map.json' ,'w') as f:
89     json.dump(id2idx, f)
90
91 # 3. Write class_map
92 with open('processed/karate-class_map.json', 'w') as f:
93     json.dump(labels, f)
94
95 # 4. Write features
96 feats = np.array(feats)
97 # shape: [number of nodes, number of features]
98 np.save('processed/karate-feats.npy', feats)
```

# GraphSAGE Code

# Train/ Val/ Test Nodes and Edges

- Graph should have '**val**' and '**test**' attribute on every node.
- Edges connected to 'val' or 'test' nodes are **removed.**

```
45    broken_count = 0
46    for node in G.nodes():
47        if not 'val' in G.node[node] or not 'test' in G.node[node]:
48            G.remove_node(node)
49            broken_count += 1
```
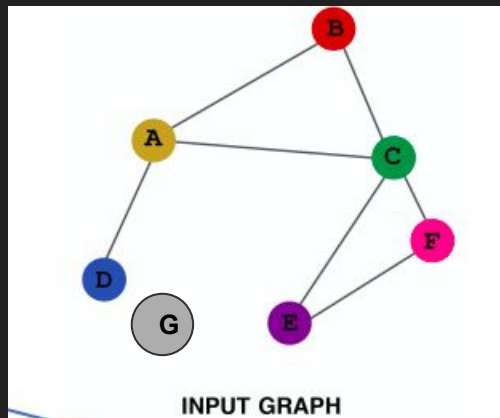
```
55    for edge in G.edges():
56        if (G.node[edge[0]]['val'] or G.node[edge[1]]['val'] or
57            G.node[edge[0]]['test'] or G.node[edge[1]]['test']):
58            G[edge[0]][edge[1]]['train_removed'] = True
59        else:
60            G[edge[0]][edge[1]]['train_removed'] = False
```

# Train/ Val/ Test Nodes and Edges

- Val and test nodes determined by the node attribute.
- The same node can be used for both val and test.
- Train nodes are all nodes except val and test.
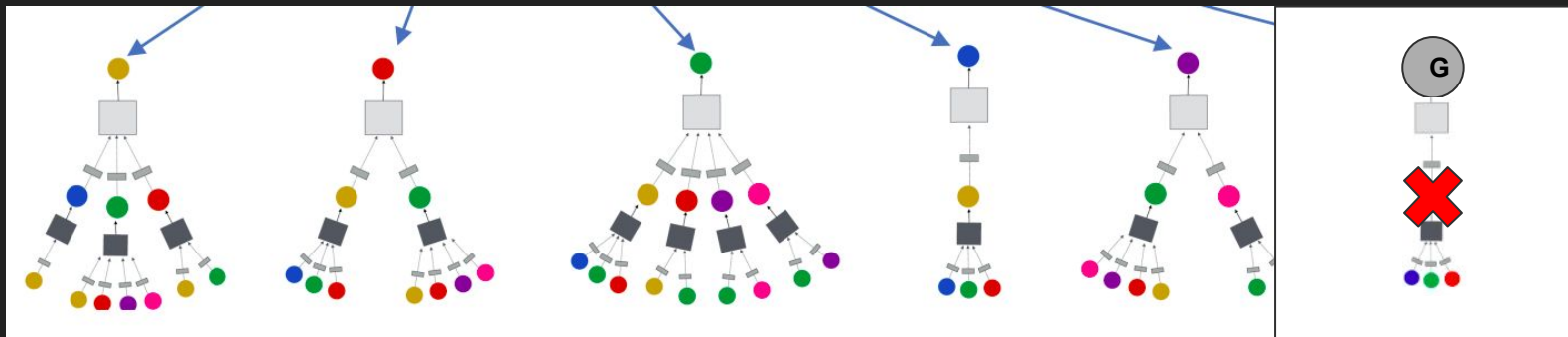
```
209        self.val_nodes = [n for n in self.G.nodes() if self.G.node[n]['val']]
210        self.test_nodes = [n for n in self.G.nodes() if self.G.node[n]['test']]
211
212        self.no_train_nodes_set = set(self.val_nodes + self.test_nodes)
213        self.train_nodes = set(G.nodes()).difference(self.no_train_nodes_set)
214        # don't train on nodes that only have edges to test set
215        self.train_nodes = [n for n in self.train_nodes if self.deg[id2idx[n]] > 0]
```

# Implementation Challenge



INPUT GRAPH

Problem:
Isolated nodes have no neighbors for aggregation

Solution:
Create a dummy vector to be the neighbor for isolated nodes.

# Adjacency List

Let max_degree = 128

Number of neighbors:

- If (>128):
  - Randomly choose 128 neighbors
- If (<128):
  - Randomly repeat some neighbors
- If (0):
  - Only points to dummy node

```python
227    def construct_adj(self):
228        # Nodes: [0...n-1]
229        # n = len(id2idx)
230        adj = len(self.id2idx)*np.ones((len(self.id2idx)+1, self.max_degree))
231        deg = np.zeros((len(self.id2idx),))
232
233        for nodeid in self.G.nodes():
234            if self.G.node[nodeid]['test'] or self.G.node[nodeid]['val']:
235                continue
236            neighbors = np.array([self.id2idx[neighbor]
237                for neighbor in self.G.neighbors(nodeid)
238                if (not self.G[nodeid][neighbor]['train_removed'])])
239            deg[self.id2idx[nodeid]] = len(neighbors)
240            if len(neighbors) == 0:
241                continue
242            if len(neighbors) > self.max_degree:
243                neighbors = np.random.choice(neighbors, self.max_degree, replace=False)
244            elif len(neighbors) < self.max_degree:
245                neighbors = np.random.choice(neighbors, self.max_degree, replace=True)
246            adj[self.id2idx[nodeid], :] = neighbors
247        return adj, deg
```

```python
133    if not features is None:
134        # pad with dummy zero vector
135        features = np.vstack([features, np.zeros((features.shape[1],))])
```

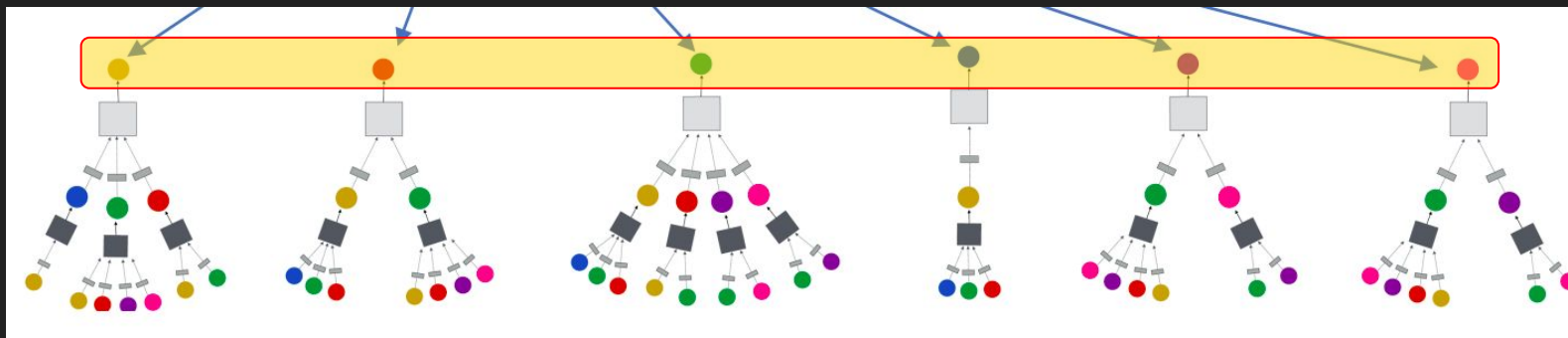- Nodes #0…#(n-1)
  - Actual nodes
- Node #n
  - Dummy node

# Identity Features

- Identity features: randomly initialized features for each node.
- If "features" and "identity features" are both provided, they will be concatenated together.
  - E.g. 9-dimensional features + 50-dimensional identity features → 59-dimensional total features.
  - E.g. shape of "features" is (501, 9), shape of "embeds" is (501, 50) → shape of total features is (501, 59).

```
51    if identity_dim > 0:
52        self.embeds = tf.get_variable("node_embeddings", [adj.get_shape().as_list()[0], identity_dim])
53    else:
54        self.embeds = None
55    if features is None:
56        if identity_dim == 0:
57            raise Exception("Must have a positive value for identity feature dimension if no input features given.")
58        self.features = self.embeds
59    else:
60        self.features = tf.Variable(tf.constant(features, dtype=tf.float32), trainable=False)
61        if not self.embeds is None:
62            self.features = tf.concat([self.embeds, self.features], axis=1)
```
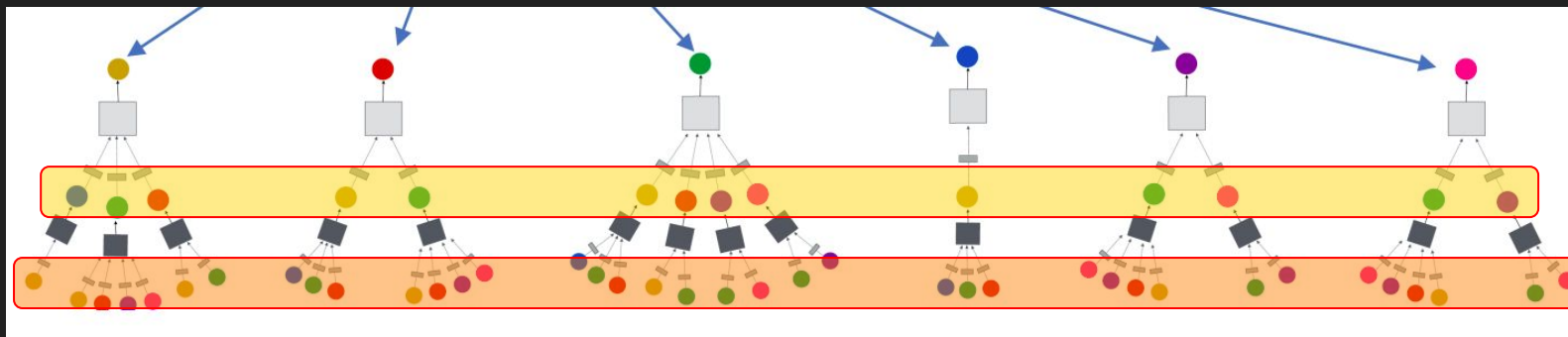
# Node Sampling

# Node Sampling

Files: minibatch.py

- (line 308) nodes are sampled from train_nodes in order.
- Every epoch, train_nodes is shuffled.

```python
304    def next_minibatch_feed_dict(self):
305        start_idx = self.batch_num * self.batch_size
306        self.batch_num += 1
307        end_idx = min(start_idx + self.batch_size, len(self.train_nodes))
308        batch_nodes = self.train_nodes[start_idx : end_idx]
309        return self.batch_feed_dict(batch_nodes)

317    def shuffle(self):
318        """ Re-shuffle the training set.
319            Also reset the batch number.
320        """
321        self.train_nodes = np.random.permutation(self.train_nodes)
322        self.batch_num = 0
```

# Neighbor Sampling

# Neighbor Sampling

Files:
supervised_train.py
models.py

```python
207            sampler = UniformNeighborSampler(adj_info)
208            layer_infos = [SAGEInfo("node", sampler, FLAGS.samples_1, FLAGS.dim_1),
209                                    SAGEInfo("node", sampler, FLAGS.samples_2, FLAGS.dim_2)]
```

```python
254        def sample(self, inputs, layer_infos, batch_size=None):
255            """ Sample neighbors to be the supportive fields for multi-layer convolutions.
256
257            Args:
258                inputs: batch inputs
259                batch_size: the number of inputs (different for batch inputs and negative samples).
260            """
261
262            if batch_size is None:
263                batch_size = self.batch_size
264            samples = [inputs]
265            # size of convolution support at each layer per node
266            support_size = 1
267            support_sizes = [support_size]
268            for k in range(len(layer_infos)):
269                t = len(layer_infos) - k - 1
270                support_size *= layer_infos[t].num_samples
271                sampler = layer_infos[t].neigh_sampler
272                node = sampler((samples[k], layer_infos[t].num_samples))
273                samples.append(tf.reshape(node, [support_size * batch_size,]))
274                support_sizes.append(support_size)
275            return samples, support_sizes
```

# Neighbor Sampling

Files: neigh_samplers.py

```python
15  class UniformNeighborSampler(Layer):
16      """
17      Uniformly samples neighbors.
18      Assumes that adj lists are padded with random re-sampling
19      """
20      def __init__(self, adj_info, **kwargs):
21          super(UniformNeighborSampler, self).__init__(**kwargs)
22          self.adj_info = adj_info
23
24      def _call(self, inputs):
25          ids, num_samples = inputs
26          adj_lists = tf.nn.embedding_lookup(self.adj_info, ids)
27          adj_lists = tf.transpose(tf.random_shuffle(tf.transpose(adj_lists)))
28          adj_lists = tf.slice(adj_lists, [0,0], [-1, num_samples])
29          return adj_lists
```

# Test Stats

```
326    print("Writing test set stats to file (don't peak!)")
327    val_cost, val_f1_mic, val_f1_mac, duration = incremental_evaluate(sess, model, minibatch, FLAGS.batch_size, test=True)
328    with open(log_dir() + "test_stats.txt", "w") as fp:
329        fp.write("loss={:.5f} f1_micro={:.5f} f1_macro={:.5f}".
330                format(val_cost, val_f1_mic, val_f1_mac))
331
```

- Log dir is:
  - ./sup-[data folder name]
- Test results are saved there.

```
81 def log_dir():
82     log_dir = FLAGS.base_log_dir + "/sup-" + FLAGS.train_prefix.split("/")[-2]
83     log_dir += "/{model:s}_{model_size:s}_{lr:0.4f}/".format(
84             model=FLAGS.model,
85             model_size=FLAGS.model_size,
86             lr=FLAGS.learning_rate)
87     if not os.path.exists(log_dir):
88         os.makedirs(log_dir)
89     return log_dir
```

# Execution Script

# Execution Script

In "GraphSAGE" folder:

- python -m graphsage.supervised_train --train_prefix data/processed/karate


Because the training data is small, adjust batch_size, print_every, etc.

- python -m graphsage.supervised_train --train_prefix data/processed/karate --batch_size 4 --print_every 2 --validate_iter 5

Done