

## Node classification and link prediction

---

You are given six real networks with class information associated with each node.

- Zachary karate club network with 2 clubs ([karate\\_club.net](#))
- Davis southern women network with 3 groups ([southern\\_women.net](#))
- US college football network with 12 conferences ([american\\_football.net](#))
- Computer scientists collaboration network with 8 fields ([sicris\\_collaboration.net](#))
- Java class dependency network with 54 packages ([cdn\\_java.net](#))
- Norwegian directors collaboration network with 2 genders ([board\\_directors.net](#))

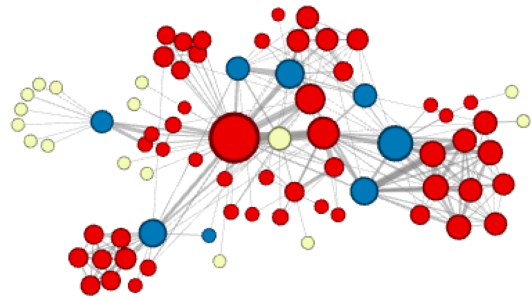
In the following exercises you will evaluate node classification and link prediction using either node or edge features or their embeddings. For this you can use a machine learning library such as [scikit-learn](#) that contains implementations of popular classifiers. Another way is to generate a tab-separated file for each exercise with nodes or edges in rows and features or embeddings in columns, and upload the file to the [Orange](#) data mining software.

### I. Node classification with features

1. **(code)** Consider different features of nodes you can compute from network structure. These include measures of node centrality, link analysis algorithms, graphlet orbit counts, node community affiliation, stochastic block models, ego-network analysis etc. Train a machine learning classifier on selected node features and evaluate its performance using some standard measure.
2. **(discuss)** How well can you classify the nodes by using their network-based features? Is the performance sufficient for practical applications? Which features are most informative for each network?

### II. Node classification with embeddings

1. **(code)** Consider the [node2vec](#) algorithm for computing low-dimensional node embeddings using biased random walks. Since the networks are relatively small, you might consider reducing the number of dimensions to  $\leq 32$  and also adjust some other default parameters (e.g., number and length of random walks, random walk biases). Train a machine learning classifier on generated node embeddings and evaluate its performance using some standard measure. (*Note that the implementation works correctly only if node identifiers are strings.*)



2. **(discuss)** How well can you classify the nodes by using their network-based embeddings? Is the performance sufficient for practical applications? How does the performance compare to using only network-based features?

### III. Link prediction with features

1. **(homework)** Consider different features of pairs of nodes (either connected or not) you can compute from network structure. These include link prediction indices, node distances, structural equivalence or topological overlap, node community co-affiliation, stochastic block models etc. Train a machine learning classifier on selected features and evaluate its performance using some standard measure. *(Make sure to compute node features from a network with test links removed.)*
2. **(discuss)** How well can you predict the links by using network-based features of pairs of nodes? Is the performance sufficient for practical applications? Which features are most informative for each network?

### IV. Link prediction with embeddings

1. **(homework)** Consider the [node2vec](#) package for computing low-dimensional embeddings of pairs of nodes (either connected or not) using biased random walks. These can be defined as a concatenation or some aggregation of embeddings of individual nodes (e.g., element-wise average or multiplication). Train a machine learning classifier on generated embeddings and evaluate its performance using some standard measure. *(Make sure to compute node embeddings from a network with test links removed.)*
2. **(discuss)** How well can you predict the links by using network-based embeddings of pairs of nodes? Is the performance sufficient for practical applications? How does the performance compare to using only network-based features?