



Erdős Goes Neural



Research Objective





Solve combinatorial problems with an unsupervised learning framework





Maximum Clique

Traveling Salesman

Minimum Cut

Example: Maximum Clique

Objective:

find the largest clique S (fully connected subgraph) within a given graph

min -|S| subject to $S \in \Omega_{clique}$ $S \subseteq V$



Why learning?

Many key combinatorial problems are NP-hard.

Learning Based

- Learns heuristics automatically from data
- Generalizes across problem sizes and distributions
- Fast inference after training

Classic algorithms — heuristics

- Hardcoded heuristics no adaptation to unseen instance distributions
- No learning from previous problem instances
- Require human-engineered problem-specific knowledge

Learning Paradigms

Supervised (Pointer Networks, Neurocore)



• Good performance

- High computational cost for labeling hard instances
- Limited generalization

Reinforcement Learning (S2V-DQN, REINFORCE)

- Flexible
- Straightforward for simple instances
- Combinatorial decision making problems

- Training stability
- Convergence issues
- Resource intensive

Learning Paradigms

Unsupervised

- Stable training
- Computationally efficient for labeling
- Better generalization

Solution Integrality









• Hard to obtain an integral solution



Background



Probabilistic Method

A non-constructive technique that proves the existence of a mathematical object by showing that the probability of randomly selecting such an object is greater than zero.

Approach to prove the existence of a set of nodes S with a desired property:

1. Define a probability distribution over the space of all possible set S on the graph

2. Show that P(S has property) > 0

Probabilistic Method - Max Cut

Objective: bipartition graph such that cut is maximized

 $\max_{S \subset V} cut(S, V - S)$





Probabilistic Method - Max Cut

Prove Max-Cut(G) $\geq \frac{|E|}{2}$

1. Define a probability distribution: sample every node (into S) with probability 1/2. $P(u \in S) = \frac{1}{2}$

$$\mathbb{E}[\operatorname{cut}(S, V - S)] = \sum_{\{u,v\} \in E} P(\{u, v\} \text{ is a cut edge})$$

$$= \sum_{\{u,v\} \in E} P[u \in S, v \notin S] + P[v \in S, u \notin S]$$

$$= \sum_{\{u,v\} \in E} \frac{1}{4} + \frac{1}{4}$$

$$= \frac{|E|}{2}$$
2. We can see that $P(\operatorname{cut}(S, V - S) \ge \frac{|E|}{2}) > 0$



Recap

- Each vertex is placed randomly in S or V S. $\mathbb{E}[\operatorname{cut}(S, V S)] = \frac{|E|}{2}$ Deterministically choose where to place each vertex to guarantee a cut of size $\geq \frac{|E|}{2}$

Derandomization via Conditional Expectation

- At each step, fix the value of one random choice.
- Choose it so that the conditional expectation of the total cut size does not decrease.
- After all choices are fixed, the total cut size is still \geq the original expected value.

Algorithm

For t = 1, ..., n do If $\mathbb{E}[\operatorname{cut} | x_1, ..., x_{t-1}, X_t = 1] > \mathbb{E}[\operatorname{cut} | x_1, ..., x_{t-1}]$ $x_t \leftarrow 1$ Else $x_t \leftarrow 0$ Return $\{x_t = 1 | v_t \in V\}$





Algorithm

For t = 1, ..., n do If $\mathbb{E}[\operatorname{cut} | x_1, ..., x_{t-1}, X_t = 1] > \mathbb{E}[\operatorname{cut} | x_1, ..., x_{t-1}]$ $x_t \leftarrow 1$ Else $x_t \leftarrow 0$ Return $\{x_t = 1 | v_t \in V\}$

$$\begin{array}{c} \bigcirc \quad X_t \\ \bullet \quad x_t = 1 \\ \bigcirc \quad x_t = 0 \end{array}$$

 $\mathbb{E}[\operatorname{cut} | x_1 = 1] > \mathbb{E}[\operatorname{cut}]$



Algorithm

For t = 1, ..., n do If $\mathbb{E}[\operatorname{cut} | x_1, ..., x_{t-1}, X_t = 1] > \mathbb{E}[\operatorname{cut} | x_1, ..., x_{t-1}]$ $x_t \leftarrow 1$ Else $x_t \leftarrow 0$ Return $\{x_t = 1 | v_t \in V\}$





Algorithm

For t = 1, ..., n do If $\mathbb{E}[\operatorname{cut} | x_1, ..., x_{t-1}, X_t = 1] > \mathbb{E}[\operatorname{cut} | x_1, ..., x_{t-1}]$ $x_t \leftarrow 1$ Else $x_t \leftarrow 0$ Return $\{x_t = 1 | v_t \in V\}$ $\begin{array}{c} \bigcirc \quad X_t \\ \bullet \quad x_t = 1 \\ \bigcirc \quad x_t = 0 \end{array}$

 $\mathbb{E}[cut | x_1 = 1, x_2 = 1] > \mathbb{E}[cut | x_1 = 1]$



Algorithm

For t = 1, ..., n do If $\mathbb{E}[\operatorname{cut} | x_1, ..., x_{t-1}, X_t = 1] > \mathbb{E}[\operatorname{cut} | x_1, ..., x_{t-1}]$ $x_t \leftarrow 1$ Else $x_t \leftarrow 0$ Return $\{x_t = 1 | v_t \in V\}$





Algorithm

For t = 1, ..., n do If $\mathbb{E}[\operatorname{cut} | x_1, ..., x_{t-1}, X_t = 1] > \mathbb{E}[\operatorname{cut} | x_1, ..., x_{t-1}]$ $x_t \leftarrow 1$ Else $x_t \leftarrow 0$ Return $\{x_t = 1 | v_t \in V\}$ $\begin{array}{c} \bigcirc \quad X_t \\ \hline \bullet \quad x_t = 1 \\ \bigcirc \quad x_t = 0 \end{array}$

$$\mathbb{E}[cut | x_1 = 1, x_2 = 1, x_3 = 1] \le \mathbb{E}[cut | x_1 = 1, x_2 = 2]$$



Algorithm

For t = 1, ..., n do If $\mathbb{E}[\operatorname{cut} | x_1, ..., x_{t-1}, X_t = 1] > \mathbb{E}[\operatorname{cut} | x_1, ..., x_{t-1}]$ $x_t \leftarrow 1$ Else $x_t \leftarrow 0$ Return $\{x_t = 1 | v_t \in V\}$





Algorithm

For t = 1, ..., n do If $\mathbb{E}[\operatorname{cut} | x_1, ..., x_{t-1}, X_t = 1] > \mathbb{E}[\operatorname{cut} | x_1, ..., x_{t-1}]$ $x_t \leftarrow 1$ Else $x_t \leftarrow 0$ Return $\{x_t = 1 | v_t \in V\}$ $\begin{array}{c} \bigcirc \quad X_t \\ \bullet \quad x_t = 1 \\ \bigcirc \quad x_t = 0 \end{array}$

$$\mathbb{E}[cut | x_1 = 1, x_2 = 1, x_3 = 0, x_4 = 1] \le \mathbb{E}[cut | x_1 = 1, x_2 = 2, x_3 = 0]$$



Algorithm

For t = 1, ..., n do If $\mathbb{E}[\operatorname{cut} | x_1, ..., x_{t-1}, X_t = 1] > \mathbb{E}[\operatorname{cut} | x_1, ..., x_{t-1}]$ $x_t \leftarrow 1$ Else $x_t \leftarrow 0$ Return $\{x_t = 1 | v_t \in V\}$ $\begin{array}{c} \bigcirc \quad X_t \\ \bullet \quad x_t = 1 \\ \bigcirc \quad x_t = 0 \end{array}$



Graph Neural Networks (GNN)

Input: G = (V, E)

Goal: learn node or graph representations

GNN Core Idea

- Each node updates its feature by combining information from its neighbors
- Repeat over multiple layers

Limitation

- All neighbors treated equally
- Sometimes cannot distinguish different graphs



Message Passing

Key Idea

- Each node updates its feature by receiving messages from neighbors
- This is repeated over multiple layers to capture larger neighborhoods

Message passing steps

- 1. **Message**: Each node collects features from its neighbors.
- 2. **Aggregation**: Combines the incoming messages (e.g., mean, sum)
- 3. **Update**: Applies a transformation (like a neural network) to update its own feature

ETH zürich

Graph Attention Networks (GAT)

Input: G = (V, E)

Goal: learn node or graph representations

GAT Core Idea

- Learn to weigh neighbors differently
- Important neighbors contribute more

Key Features

- Focuses on relevant parts of the graph
- Adapts to node importance dynamically



Erdős goes Neural



Erdos goes Neural

To solve the problem

 $\min_{S \subseteq V} f(S;G) \text{ subject to } S \in \Omega$

1. Probability function learned by a GNN

2. P(S has property) > 0 is optimized trough the derived loss function

Pipeline

Construct a GNN that outputs a probability p_i on each node v_i .



Training

Minimizing the probabilistic penalty loss ensures the learned distribution contains a small cost and feasible set with sufficient probability.



Decoding

Recover the *discrete* solution by sampling or by the method of conditional expectation.



Erdos goes neural: Loss function



Theorem (informal): let f be non-negative. There exists a set $S^* \sim D$ that satisfies:

```
f(S^*;G) < \mathcal{L}(D;G), \quad S^* \in \Omega,
```

with positive probability.





Experimental Setup

Problem: Maximum clique

Baseline: neural networks (Toenshoff 2019), classic heuristics, and solvers (CBC, Gurobi 9.0)

Datasets:

- Real world: TWITTER, IMDB, COLLAB
- Hard instances for maximum clique

Architecture

- GIN and GAT layers
- Residual connection and batch norm
- One hot encoding of random node as input
- Train with Adam and early stopping

Results

Test set approximation ratios for all methods on real-world datasets

	IMDB	COLLAB	TWITTER
Erdős' GNN (fast)	$1.000 \ (0.08 \ s/g)$	$0.982 \pm 0.063 \; (0.10 \; { m s/g})$	$0.924 \pm 0.133 \; (0.17 \; { m s/g})$
Erdős' GNN (accurate)	$1.000 \ (0.10 \ s/g)$	$0.990 \pm 0.042 \; (0.15 \; \mathrm{s/g})$	$0.942 \pm 0.111 \; (0.42 \; { m s/g})$
RUN-CSP (fast)	$0.823 \pm 0.191 \; (0.11 \; \mathrm{s/g})$	$0.912 \pm 0.188 \; (0.14 \; \mathrm{s/g})$	$0.909 \pm 0.145 \; (0.21 \; \mathrm{s/g})$
RUN-CSP (accurate)	$0.957 \pm 0.089 \; (0.12 \; {\rm s/g})$	$0.987 \pm 0.074 \; (0.19 \; \mathrm{s/g})$	$0.987 \pm 0.063 \; (0.39 \; {\rm s/g})$
Bomze GNN	$0.996 \pm 0.016 \; (0.02 \; \mathrm{s/g})$	$0.984 \pm 0.053 \ (0.03 \ s/g)$	$0.785 \pm 0.163 \ (0.07 \ s/g)$
MS GNN	$0.995\pm0.068~(0.03~{\rm s/g})$	$0.938 \pm 0.171 \ (0.03 \ s/g)$	$0.805 \pm 0.108 \ (0.07 \ s/g)$
NX MIS approx.	$0.950\pm0.071~(0.01~{\rm s/g})$	$0.946\pm0.078~(1.22~{\rm s/g})$	$0.849 \pm 0.097 \; (0.44 \; {\rm s/g})$
Greedy MIS Heur.	$0.878 \pm 0.174 \; (1e-3 \; s/g)$	$0.771 \pm 0.291 \; (0.04 \; {\rm s/g})$	$0.500\pm0.258~(0.05~{\rm s/g})$
Toenshoff-Greedy	$0.987 \pm 0.050 \; (1e-3 \; s/g)$	$0.969\pm0.087~(0.06~{\rm s/g})$	$0.917\pm0.126(0.08{\rm s/g})$
CBC (1s)	$0.985 \pm 0.121 \; (0.03 \; \rm s/g)$	$0.658 \pm 0.474 \; (0.49 \; {\rm s/g})$	$0.107 \pm 0.309 ~(1.48 ~\rm s/g)$
CBC (5s)	$1.000 \ (0.03 \ s/g)$	$0.841 \pm 0.365 \; (1.11 \; \mathrm{s/g})$	$0.198\pm0.399(4.77{\rm s/g})$
Gurobi $9.0 (0.1s)$	$1.000 \ (1e-3 \ s/g)$	$0.982 \pm 0.101 \; (0.05 \; \mathrm{s/g})$	$0.803 \pm 0.258 \; (0.21 \; { m s/g})$
Gurobi $9.0 (0.5s)$	$1.000 \ (1e-3 \ s/g)$	$0.997 \pm 0.035 \; (0.06 \; \rm s/g)$	$0.996 \pm 0.019 \; (0.34 \; { m s/g})$
Gurobi $9.0 (1s)$	$1.000 \ (1e-3 \ s/g)$	$0.999 \pm 0.015 \; (0.06 \; {\rm s/g})$	$1.000 (0.34 { m s/g})$
Gurobi 9.0 (5s)	1.000 (1e-3 s/g)	$1.000 \ (0.06 \ s/g)$	$1.000 \ (0.35 \ s/g)$

Results

Approximation ratios for best methods on hard instances

	Training set	Test set	Large Instances
Erdős' GNN (fast)	$0.899\pm0.064~(0.27~{\rm s/g})$	$0.788\pm0.065(0.23~{\rm s/g})$	$0.708 \pm 0.027 \; (1.58 \; { m s/g})$
Erdős' GNN (accurate)	$0.915 \pm 0.060 \; (0.53 \; { m s/g})$	$0.799 \pm 0.067 \; (0.46 \; { m s/g})$	$0.735 \pm 0.021 \; (6.68 \; \mathrm{s/g})$
$\overline{\text{RUN-CSP}}$ (fast)	$0.833 \pm 0.079 \; (0.27 \; { m s/g})$	$0.738 \pm 0.067 \; (0.23 \; { m s/g})$	$0.771 \pm 0.032 \; (1.84 \; \mathrm{s/g})$
RUN-CSP (accurate)	$0.892\pm0.064~(0.51~{\rm s/g})$	$0.789\pm0.053~(0.47~{\rm s/g})$	$0.804\pm0.024~(5.46~{\rm s/g})$
Toenshoff-Greedy	$0.924 \pm 0.060 ~(0.02 ~{ m s/g})$	$0.816\pm0.064~(0.02~{\rm s/g})$	$0.829\pm0.027(0.35{\rm s/g})$
Gurobi $9.0 (0.1s)$	$0.889 \pm 0.121 \; (0.18 \; { m s/g})$	$0.795 \pm 0.118 \; (0.16 \; { m s/g})$	$0.697 \pm 0.033 \; (1.17 \; {\rm s/g})$
Gurobi $9.0 \ (0.5s)$	$0.962 \pm 0.076 \; (0.34 \; { m s/g})$	$0.855 \pm 0.083 \; (0.31 \; { m s/g})$	$0.697 \pm 0.033 \; (1.54 \; { m s/g})$
Gurobi $9.0 (1.0s)$	$0.980\pm0.054(0.45{\rm s/g})$	$0.872\pm0.070(0.40{ m s/g})$	$0.705 \pm 0.039 \; (2.05 \; { m s/g})$
Gurobi $9.0 (5.0s)$	$0.998 \pm 0.010 (0.76 { m s/g})$	$0.884 \pm 0.062 \; (0.68 \; { m s/g})$	$0.790 \pm 0.285 \; (6.01 \; \mathrm{s/g})$
Gurobi $9.0 (20.0s)$	$0.999\pm0.003~(1.04~{ m s/g})$	$0.885 \pm 0.063 \; (0.96 \; { m s/g})$	$0.807 \pm 0.134 \; (21.24 \; {\rm s/g})$

Conclusions

Main Features

- Merges probabilistic method with deep learning
- Principled way to ensure feasibility of solutions for various CO problems
- Competitive performance

Conclusions

Limitations

- Theoretical guarantees No worst-case approximation bounds
- Lacks strong generalization
- MIP Solvers like Gurobi are still SoA and greedy algorithms can be extremely efficient in practice
- Can be hard to obtain analytically a differentiable expression for certain constraints (e.g., imposing path or tree structure on solutions).

Follow-up

Erdos goes neural laid the groundwork for unsupervised neural approaches to combinatorial optimization

- Meta-EGN (ICLR 2023) Learns to generalize across tasks by optimizing initial weights great generalization but higher training cost
- Objective Relaxation (NeurIPS 2022) Relax objectives to ensure quality of relaxed solution Theoretically grounded
- DIFUSCO (NeurIPS 2023) leverages graph-based diffusion models to iteratively refine solutions

Thank You



Erdos Goes Neural: architecture

- 6 GIN layers followed by a multi-head GAT layer
- Residual connection and batch normalization at every layer
- Graph size normalization, for optimization stability
- The output of GNN to a 2 layers MLP, giving as output one value per node
- Min-Max normalization to rescale the output values in the interval [0,1]

Case study: Maximum clique

$$G = (V, E) \quad |V| = N, \quad |E| = M$$

Objective $\min_{S \subseteq V} |\gamma - |S| \text{ subject to } S \in \Omega_{\text{clique}}$

Probabilistic penalty
$$\mathscr{L}_{clique}(D,G) = \gamma + E[|S|] + \beta E[|V-S|]$$

Loss $= \gamma - (\beta + 1) \sum_{(v_i,v_j)\in E} w_{ij}p_ip_j + \frac{\beta}{2} \sum_{v_i \neq v_j} p_ip_j$

With probability at least *t*, set $S^* \sim D$ satisfies

$$\gamma - \frac{\mathscr{L}(D,G)}{1-t} \le |S^*| \quad S^* \in \Omega_{\text{clique}}$$







Maximum Clique: Step 1

 $G = (V, E) \quad |V| = N, \quad |E| = M$

1. Use a Bernoulli random variable for each node

 $x_i = \begin{cases} 1, & \text{with probability } p_i \\ 0, & \text{with probability } 1 - p_i \end{cases}$

2. Write the cost function for the problem

3. Build a GNN

1. Take in input some node features

2. Output Nx1 vector of probabilities

 $\min_{S \subseteq V} \gamma - |S| \quad \text{subject to } S \text{ is a clique}$

Maximum Clique: Step 2

Derivation of the loss function





Maximum Clique: Step 3

- 1. Train the network
- 2. Retrieve the set S from the probabilities of the network using the method of conditional expectation:

Sort the nodes according to their probabilities. Starting from the high probability nodes, for each node v_i do:

1. Evaluate the loss for $p_i = 1$ and for $p_i = 0$

- 2. Set p_i to either 1 or 0 depending on what achieved the better loss
- 3. Move on to the next node and repeat from step 1

43

Graph Isomorphism Networks (GIN)

Input: G = (V, E)

Goal: learn node or graph representations

GIN Core Idea

- Stronger aggregation: sum all neighbor features
- Apply a neural network (MLP) after aggregation

Key Features

- Best at capturing subtle graph structures
- MLP can learn complex local patterns

