

Machine learning for algorithm design: Theoretical guarantees and applied frontiers

Ellen Vitercik

Stanford University

How to integrate **machine learning** into **algorithm design**?



Algorithm configuration

How to tune an algorithm's parameters?



Algorithm selection

Given a variety of algorithms, which to use?



Algorithm design

Can machine learning guide algorithm discovery?

How to integrate **machine learning** into **algorithm design**?



Algorithm configuration

How to tune an algorithm's parameters?



Algorithm selection

Given a variety of algorithms, which to use?



Algorithm design

Can machine learning guide algorithm discovery?

Algorithm configuration

Example: **Integer programming solvers**

Most popular tool for solving combinatorial (& nonconvex) problems



Routing



Manufacturing



Scheduling



Planning



Finance

Algorithm configuration

IP solvers (CPLEX, Gurobi) have a **ton** parameters

- CPLEX has **170-page** manual describing **172** parameters
- Tuning by hand is notoriously **slow, tedious,** and **error-prone**

CPX_PARAM_NODEFILEIND 100	CPX_PARAM_TRELIM 160	CPX_PARAM_RANDOMSEED 130	CPXPARAM_MIP_Pool_RelGap 148	CPX_PARAM_FLOWCOVERS 70	CPX_PARAM_BRDIR 39
CPX_PARAM_NODELIM 101	CPX_PARAM_TUNINGDETILIM 160	CPX_PARAM_REDUCE 131	CPXPARAM_MIP_Pool_Replace 151	CPX_PARAM_FLOWPATHS 71	CPX_PARAM_BTTL 40
CPX_PARAM_NODESEL 102	CPX_PARAM_TUNINGDISPLAY 162	CPX_PARAM_REINV 131	CPXPARAM_MIP_Strategy_Branch 39	CPX_PARAM_FPHEUR 72	CPX_PARAM_CALCQCPCDUALS 41
CPX_PARAM_NUMERICALEMPHASIS 102	CPX_PARAM_TUNINGMEASURE 163	CPX_PARAM_RELAXPREIND 132	CPXPARAM_MIP_Strategy_MIQCPStrat 93	CPX_PARAM_FRACCAND 73	CPX_PARAM_CLIQUES 42
CPX_PARAM_NZREADLIM 103	CPX_PARAM_TUNINGREPEAT 164	CPX_PARAM_RELOBJDIF 133	CPXPARAM_MIP_Strategy_StartAlgorithm 139	CPX_PARAM_FRACCUTS 73	CPX_PARAM_CLOCKTYPE 43
CPX_PARAM_OBJDIF 104	CPX_PARAM_TUNINGTILIM 165	CPX_PARAM_REPAIRTRIES 133	CPXPARAM_MIP_Strategy_VariableSelect 166	CPX_PARAM_FRACPASS 74	CPX_PARAM_CLONELOG 43
CPX_PARAM_OBJLLIM 105	CPX_PARAM_VARSEL 166	CPX_PARAM_REPEATPRESOLVE 134	CPXPARAM_MIP_SubMIP_NodeLimit 155	CPX_PARAM_GUBCOVERS 75	CPX_PARAM_COEREDIND 44
CPX_PARAM_OBJULIM 105	CPX_PARAM_WORKDIR 167	CPX_PARAM_RINSHEUR 135	CPXPARAM_OptimalityTarget 106	CPX_PARAM_HEURFREQ 76	CPX_PARAM_COLREADLIM 45
CPX_PARAM_PARALLELMODE 108	CPX_PARAM_WORKMEM 168	CPX_PARAM_RLT 136	CPXPARAM_Output_WriteLevel 169	CPX_PARAM_IMPLBD 76	CPX_PARAM_CONFLICTDISPLAY 46
CPX_PARAM_PERIND 110	CPX_PARAM_WRITELEVEL 169	CPX_PARAM_ROWREADLIM 141	CPXPARAM_Preprocessing_Aggregator 19	CPX_PARAM_INTSOLFILEPREFIX 78	CPX_PARAM_COVERS 47
CPX_PARAM_PERLIM 111	CPX_PARAM_ZEROHALFCUTS 170	CPX_PARAM_SCAIND 142	CPXPARAM_Preprocessing_Fill 19	CPX_PARAM_INTSOLLIM 79	CPX_PARAM_CPUMASK 48
CPX_PARAM_POLISHAFTERDETTIME 111	CPXPARAM_Benders_Strategy 30	CPX_PARAM_SCRIND 143	CPXPARAM_Preprocessing_Linear 120	CPX_PARAM_ITLIM 80	CPX_PARAM_CRAIN 50
CPX_PARAM_POLISHAFTEREPAGAP 112	CPXPARAM_Benders_Tolerances_feasibilitycut 35	CPX_PARAM_SIFTALG 143	CPXPARAM_Preprocessing_Reduce 131	CPX_PARAM_LANDPCUTS 82	CPX_PARAM_CUTLO 51
CPX_PARAM_POLISHAFTEREPGAP 113	CPXPARAM_Benders_Tolerances_optimalitycut 36	CPX_PARAM_SIFTDISPLAY 144	CPXPARAM_Preprocessing_Symmetry 156	CPX_PARAM_LBHEUR 81	CPX_PARAM_CUTPASS 52
CPX_PARAM_POLISHAFTERINTSOL 114	CPXPARAM_Conflict_Algorithm 46	CPX_PARAM_SIFTITLIM 145	CPXPARAM_Read_DataCheck 54	CPX_PARAM_LPMETHOD 136	CPX_PARAM_CUTSFACTOR 52
CPX_PARAM_POLISHAFTERNODE 115	CPXPARAM_CPUmask 48	CPX_PARAM_SIMDISPLAY 145	CPXPARAM_Read_Scale 142	CPX_PARAM_MFCUTS 82	CPX_PARAM_CUTUP 53
CPX_PARAM_POLISHAFTERTIME 116	CPXPARAM_DistMIP_Rampup_Duration 128	CPX_PARAM_SINGLIM 146	CPXPARAM_ScreenOutput 143	CPX_PARAM_MEMORYEMPHASIS 83	CPXPARAM_DATACHECK 54
CPX_PARAM_POLISHTIME (deprecated) 116	CPXPARAM_LPMethod 136	CPX_PARAM_SOLNPOOLGAP 146	CPXPARAM_Sifting_Algorithm 143	CPX_PARAM_MIPCBREDLP 84	CPX_PARAM_DEPIND 55
CPX_PARAM_POPULATELIM 117	CPXPARAM_MIP_Cuts_BQP 38	CPX_PARAM_SOLNPOOLCAPACITY 147	CPXPARAM_Sifting_Display 144	CPX_PARAM_MIPDISPLAY 85	CPX_PARAM_DETTILIM 56
CPX_PARAM_PPRIND 118	CPXPARAM_MIP_Cuts_LocallyImplied 77	CPX_PARAM_SOLNPOOLREPLACE 151	CPXPARAM_Sifting_Iterations 145	CPX_PARAM_MIPEMPHASIS 87	CPX_PARAM_DISJCUTS 57
CPX_PARAM_PREDUAL 119	CPXPARAM_MIP_Cuts_RLT 136	CPX_PARAM_SOLNPOOLINTENSITY 149	CPXPARAM_Simplex_Display 145	CPX_PARAM_MIPINTERVAL 88	CPX_PARAM_DIVETYPE 58
CPX_PARAM_PREIND 120	CPXPARAM_MIP_Cuts_ZeroHalfCut 170	CPX_PARAM_SOLNPOOLREPLACE 151	CPXPARAM_Simplex_Limits_Singularity 146	CPX_PARAM_MIPKAPPASTATS 89	CPX_PARAM_DPRIIND 59
CPX_PARAM_PRLINEAR 120	CPXPARAM_MIP_Limits_CutsFactor 52	CPX_PARAM_SOLUTIONTARGET (deprecated: see CPXPARAM_OptimalityTarget 106)	CPXPARAM_SolutionType 152	CPX_PARAM_MIPORDIND 90	CPX_PARAM_EACHCUTLIM 60
CPX_PARAM_PREPASS 121	CPXPARAM_MIP_Limits_RampupDetTimeLimit 127	CPX_PARAM_STARTALG 139	CPXPARAM_Threads 157	CPX_PARAM_MIPORDTYPE 91	CPX_PARAM_EPAGAP 61
CPX_PARAM_PRESLVND 122	CPXPARAM_MIP_Limits_RampupTimeLimit 128	CPX_PARAM_STRONGCANDLIM 154	CPXPARAM_TimeLimit 159	CPX_PARAM_MIPSEARCH 92	CPX_PARAM_EPGAP 61
CPX_PARAM_PRICELIM 123	CPXPARAM_MIP_Limits_Solutions 79	CPX_PARAM_STRONGCANDLIM 154	CPXPARAM_Tune_DefTimeLimit 160	CPX_PARAM_MIQCPSTRAT 93	CPX_PARAM_EPINT 62
CPX_PARAM_PROBE 123	CPXPARAM_MIP_Limits_StrongCand 154	CPX_PARAM_STRONGITLIM 154	CPXPARAM_Tune_Display 162	CPX_PARAM_MIRCUTS 94	CPX_PARAM_EPMRK 64
CPX_PARAM_PROBEDETTIME 124	CPXPARAM_MIP_Limits_StrongIt 154	CPX_PARAM_SUBALG 99	CPXPARAM_Tune_Measure 163	CPX_PARAM_MPSLONGNUM 94	CPX_PARAM_EPOPT 65
CPX_PARAM_PROBETIME 124	CPXPARAM_MIP_Limits_TreeMemory 160	CPX_PARAM_SUBMIPNODELIMIT 155	CPXPARAM_Tune_Repeat 164	CPX_PARAM_NETDISPLAY 95	CPX_PARAM_EPPER 65
CPX_PARAM_QPMAKEPSDIND 125	CPXPARAM_MIP_OrderType 91	CPX_PARAM_SYMMETRY 156	CPXPARAM_Tune_TimeLimit 165	CPX_PARAM_NETEPOPT 96	CPX_PARAM_EPRELAX 66
CPX_PARAM_QPMETHOD 138	CPXPARAM_MIP_Pool_AbsGap 146	CPX_PARAM_THREADS 157	CPXPARAM_WorkDir 167	CPX_PARAM_NETEPRHS 96	CPX_PARAM_EPRHS 67
CPX_PARAM_QPNZREADLIM 126	CPXPARAM_MIP_Pool_Capacity 147	CPX_PARAM_TILIM 159	CPXPARAM_WorkMem 168	CPX_PARAM_NETFIND 97	CPX_PARAM_FEASOPTMODE 68
	CPXPARAM_MIP_Pool_Intensity 149		CraInd 50	CPX_PARAM_NETITLIM 98	CPX_PARAM_FILEENCODING 69
				CPX_PARAM_NETPPRIIND 98	

Algorithm configuration

IP solvers (CPLEX, Gurobi) have a **ton** parameters

- CPLEX has **170-page** manual describing **172** parameters
- Tuning by hand is notoriously **slow, tedious**, and **error-prone**

What's the best **configuration** for the application at hand?



Best configuration for **routing** problems
likely not suited for **scheduling**



How to integrate **machine learning** into **algorithm design**?



Algorithm configuration

How to tune an algorithm's parameters?



Algorithm selection

Given a variety of algorithms, which to use?



Algorithm design

Can machine learning guide algorithm discovery?

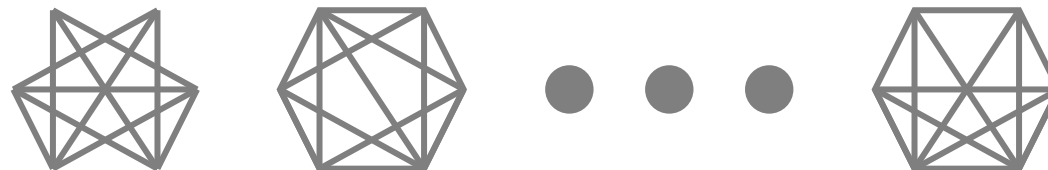
Algorithm selection in theory

Worst-case analysis has been the main framework for decades
Has led to beautiful, practical algorithms

Worst-case instances **rarely occur in practice**

In practice:

Instances solved in **past** are similar to **future** instances...




A stack of several cardboard boxes, each secured with red and white striped string. The boxes are labeled 'FRAGILE' with icons of a glass, a wine glass, and an umbrella. A person's hands are visible at the bottom, holding the stack. The background is blurred, showing a person in a white shirt.

**In practice, we have data about
the application domain**

Routing problems a shipping company solves

**In practice, we have data about
the application domain**



Clustering problems a biology lab solves

**In practice, we have data about
the application domain**



Scheduling problems an airline solves

Existing research



Constraint satisfaction

[Horvitz, Ruan, Gomes, Krautz, Selman, Chickering, UAI'01; ...]



Integer programming

[Hutter, Hoos, Leyton-Brown, CPAIOR '10; ...]



Economics (mechanism design)

[Likhodedov, Sandholm, AAI '04, '05; ...]



Computational biology

[Majoros, Salzberg, Bioinformatics'04; ...]

**Applied
research**

2000

2023

Existing research

Automated algorithm configuration and selection

[Gupta, Roughgarden, ITCS'16; Balcan, Nagarajan, **Vitercik**, White, COLT'17; ...]

Learning-augmented algorithms

[Lykouris, Vassilvitskii, ICML'18; Mitzenmacher, NeurIPS'18; ...]

Sample complexity of revenue maximization

[Balcan, Blum, Hartline, Mansour, FOCS'05; Elkind, SODA'07; ...]

**Applied
research**

**Theory
research**

2000

2023

ML + algorithm design: Potential impact

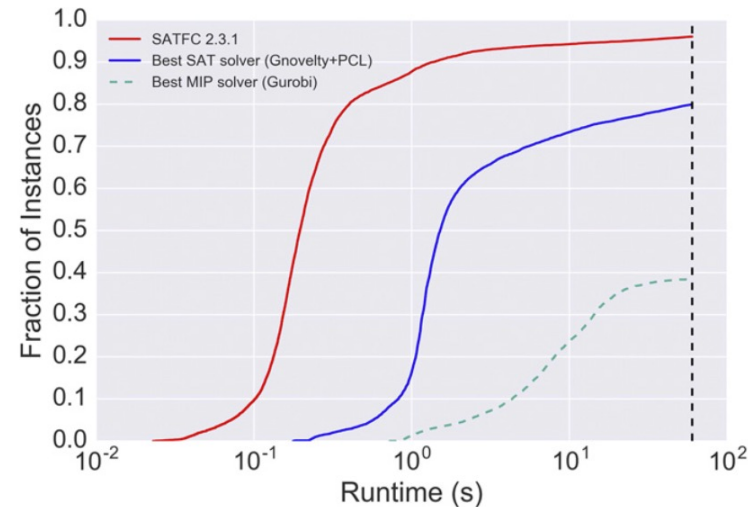
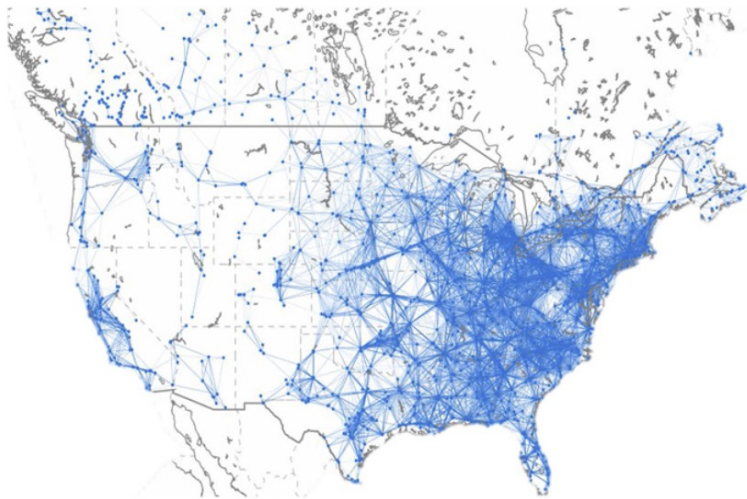
Example: integer programming

- Used heavily throughout industry and science
- **Many** different ways to incorporate **learning** into solving
- Solving is very difficult, so ML can make a huge difference



Example: Spectrum auctions

- In '16-'17, FCC held a \$19.8 billion radio spectrum auction
 - Involves solving huge graph-coloring problems



- SATFC uses algorithm configuration + selection
- Simulations indicate SATFC saved the government billions

Plan for tutorial

① **Theoretical guarantees**

- a. Statistical guarantees for algorithm configuration
- b. Online algorithm configuration

② **Applied techniques**

- a. Graph neural networks

Plan for tutorial

1 Theoretical guarantees

- a. **Statistical guarantees for algorithm configuration**
- b. Online algorithm configuration

2 Applied techniques

- a. Graph neural networks

Gupta, Roughgarden, ITCS'16

Balcan, DeBlasio, Dick, Kingsford, Sandholm, **Vitercik**, STOC'21

Balcan, Prasad, Sandholm, **Vitercik**, NeurIPS'21

Balcan, Prasad, Sandholm, **Vitercik**, NeurIPS'22

Running example: Sequence alignment

Goal: Line up pairs of strings

Applications: Biology, natural language processing, etc.



Did you mean: [vitercik](#)

Sequence alignment algorithms

Input: Two sequences S and S'

Output: Alignment of S and S'

$S = A C T G$
 $S' = G T C A$

Gap
↓
A - - C T G
- G T C A -
↑ ↑ ↑
Insertion/deletion (*indel*) Match Mismatch

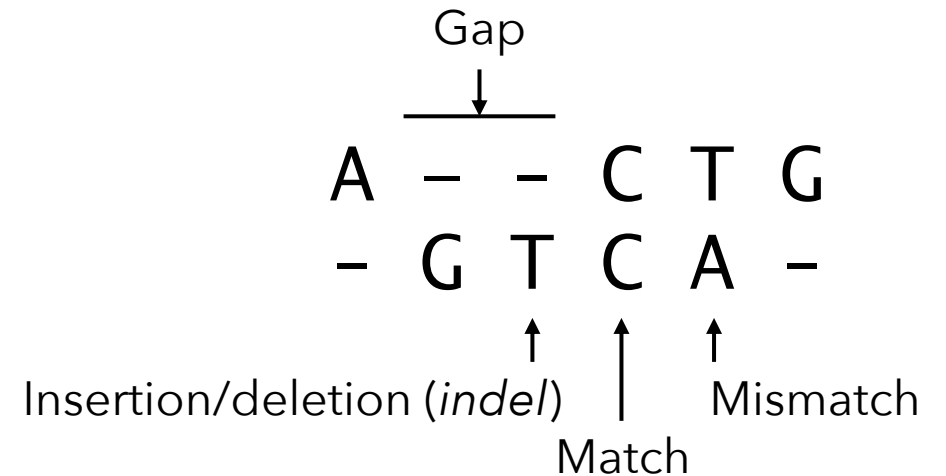
Sequence alignment algorithms

Standard algorithm with parameters $\rho_1, \rho_2, \rho_3 \geq 0$:

Return alignment maximizing:

$$(\# \text{ matches}) - \rho_1 \cdot (\# \text{ mismatches}) - \rho_2 \cdot (\# \text{ indels}) - \rho_3 \cdot (\# \text{ gaps})$$

$S = A C T G$
 $S' = G T C A$



Sequence alignment algorithms


Can sometimes access **ground-truth, reference** alignment

E.g., in computational biology: Bahr et al., Nucleic Acids Res.'01; Raghava et al., BMC Bioinformatics '03; Edgar, Nucleic Acids Res.'04; Walle et al., Bioinformatics'04

Requires extensive manual alignments
...rather just run parameterized algorithm

How to tune algorithm's parameters?

*"There is **considerable disagreement** among molecular biologists about the **correct choice**" [Gusfield et al. '94]*



A	-	-	C	T	G
-	G	T	C	A	-

Sequence alignment algorithms

-GRTCPKPDDLPFSTVVP-LKTFYEPGEEITYSCKPGYVSRGGMRKFICPLTGLWPINTLKCTP
E-VKCPFPSRPDNGFVNYPKPTLYYKDKATFGCHDGYSLDGP-EEIECTKLGNEWSAMPSC-KA

Ground-truth alignment of protein sequences

Sequence alignment algorithms

-GRTCPKPDDLPFSTVVP-LKTFYEPEEITYSCKPGYVSRGGMRKFICPLTGLWPINTLKCTP
E-VKCPFPSRPDNGFVNYPKPTLYYKDKATFGCHDGYSLDGP-EEIECTKLGNSAMPSC-KA

Ground-truth alignment of protein sequences

GRTCP---KPDDLPFSTVVPLKTFYEPEEITYSCKPGYVSRGGMRKFICPLTGLWPINTLKCTP
EVKCPFPSRPDN-GFVNYPKPTLYYK-DKATFGCHDGY-SLDGPEEIECTKLGNS-AMPSCKA

Alignment by algorithm with **poorly-tuned** parameters

Sequence alignment algorithms

-GRTCPKPDDLPFSTVVP-LKTFYEPGEEITYSCKPGYVSRGGMRKFICPLTGLWPINTLKCTP
E-VKCPFPSRPDNGFVNYPAPKPTLYYKDKATFGCHDGYSLDGP-EEIECTKLGNSAMPSC-KA

Ground-truth alignment of protein sequences

GRTCP---KPDDLPFSTVVPLKTFYEPGEEITYSCKPGYVSRGGMRKFICPLTGLWPINTLKCTP
EVKCPFPSRPDN-GFVNYPAPKPTLYYK-DKATFGCHDGY-SLDGPEEIECTKLGNS-AMPSCKA

Alignment by algorithm with **poorly-tuned** parameters

GRTCPKPDDLPFSTV-VPLKTFYEPGEEITYSCKPGYVSRGGMRKFICPLTGLWPINTLKCTP
EVKCPFPSRPDNGFVNYPAPKPTLYYKDKATFGCHDGY-SLDGPEEIECTKLGNSA-MPSCKA

Alignment by algorithm with **well-tuned** parameters

Automated parameter tuning procedure

1. Fix parameterized algorithm
2. Receive training set T of "typical" inputs



3. Find parameter setting w/ good avg performance over T

Runtime, solution quality, etc.

Automated parameter tuning procedure

1. Fix parameterized algorithm
2. Receive training set T of "typical" inputs



3. Find parameter setting w/ good avg performance over T

On average, output alignment is close to reference alignment

Automated parameter tuning procedure

1. Fix parameterized algorithm
2. Receive training set T of "typical" inputs



3. Find parameter setting w/ good avg performance over T

Key question:

How to find parameter setting with good avg performance?

Automated parameter tuning procedure

Key question:

How to find parameter setting with good avg performance?



E.g., for sequence alignment:
algorithm by Gusfield et al. ['94]

Many other generic search strategies

E.g., Hutter et al. [JAIR'09, LION'11], Ansótegui et al. [CP'09], ...

Automated parameter tuning procedure

1. Fix parameterized algorithm
2. Receive training set T of "typical" inputs

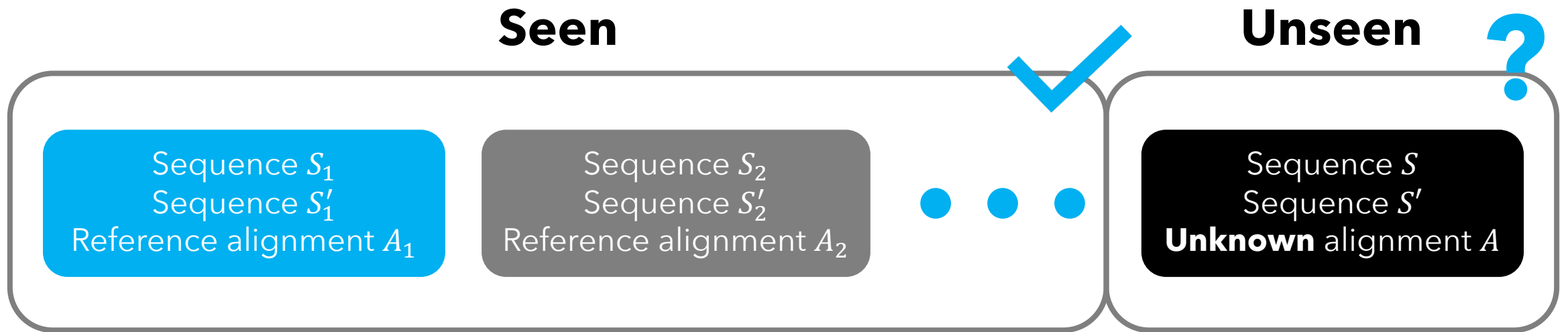


3. Find parameter setting w/ good avg performance over T

Key question (focus of this section):

Will that parameter setting have good **future** performance?

Automated parameter tuning procedure



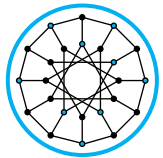
Key question (focus of this section):

Will that parameter setting have good **future** performance?

Generalization

Key question (focus of this section):

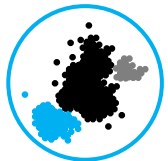
Good performance on **average** over **training set** implies good **future** performance?



Greedy algorithms

Gupta, Roughgarden, ITCS'16

First to ask question for algorithm configuration



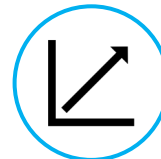
Clustering

Balcan, Nagarajan, **V**, White, COLT'17
Garg, Kalai, NeurIPS'18
Balcan, Dick, White, NeurIPS'18
Balcan, Dick, Lang, ICLR'20



Search

Sakaue, Oki, NeurIPS'22



Numerical linear algebra

Bartlett et al., COLT'22

And many other areas...

This section: Main result

Key question (focus of this section):

Good performance on **average** over **training set** implies good **future** performance?

Answer this question for any parameterized algorithm where:

Performance is **piecewise-structured** function of parameters

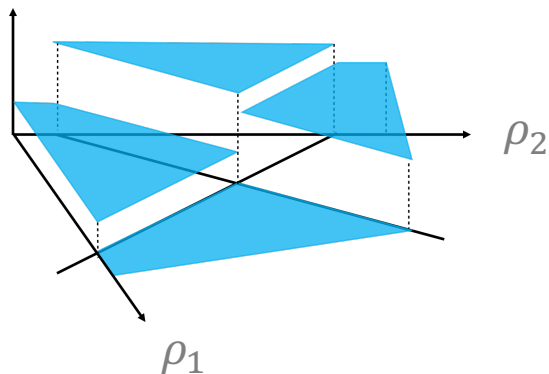
Piecewise constant, linear, quadratic, ...

This section: Main result

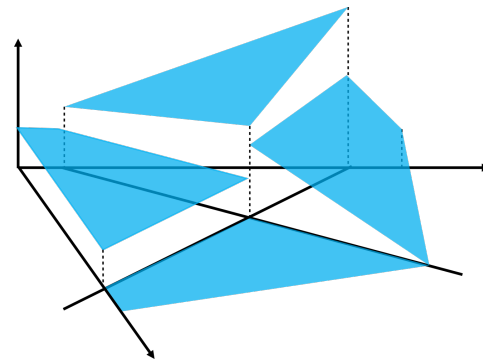
Performance is **piecewise-structured** function of parameters

Piecewise constant, linear, quadratic, ...

Algorithmic
performance
on fixed input



Piecewise constant



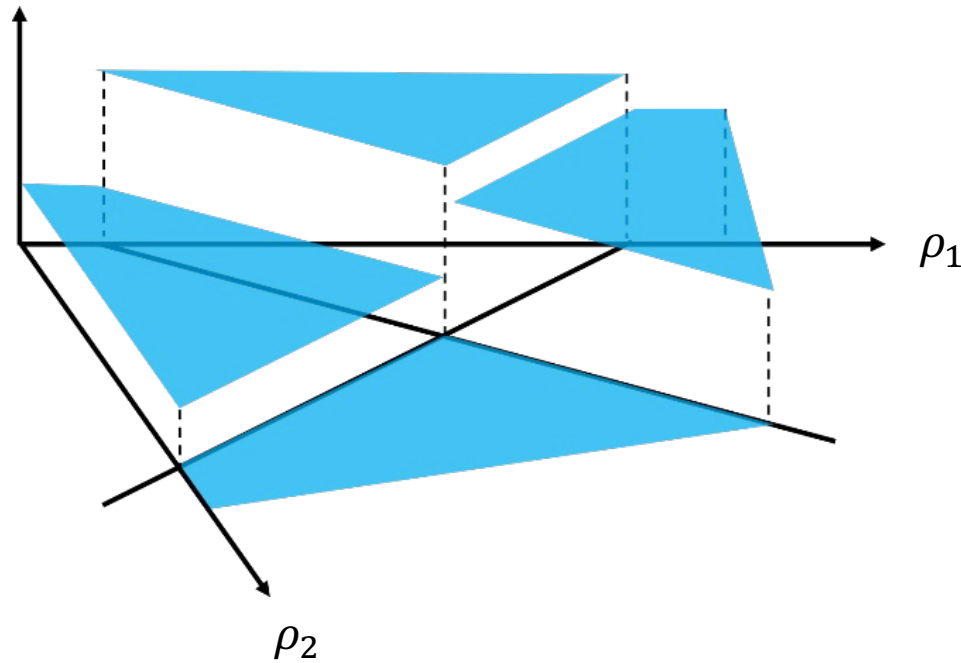
Piecewise linear



Piecewise ...

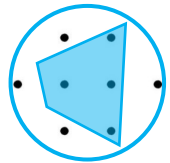
Example: Sequence alignment

Distance between **algorithm's output** given S, S'
and **ground-truth** alignment is p-wise constant



Piecewise structure

Piecewise structure unifies **seemingly disparate** problems:



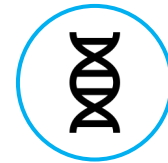
Integer programming

Balcan, Dick, Sandholm, **V**, ICML'18
Balcan, Prasad, Sandholm, **V**, NeurIPS'21
Balcan, Prasad, Sandholm, **V**, NeurIPS'22



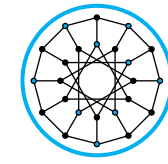
Clustering

Balcan, Nagarajan, **V**, White, COLT'17
Balcan, Dick, White, NeurIPS'18
Balcan, Dick, Lang, ICLR'20



Computational biology

Balcan, DeBlasio, Dick, Kingsford,
Sandholm, **V**, STOC'21



Greedy algorithms

Gupta, Roughgarden, ITCS'16



Mechanism configuration

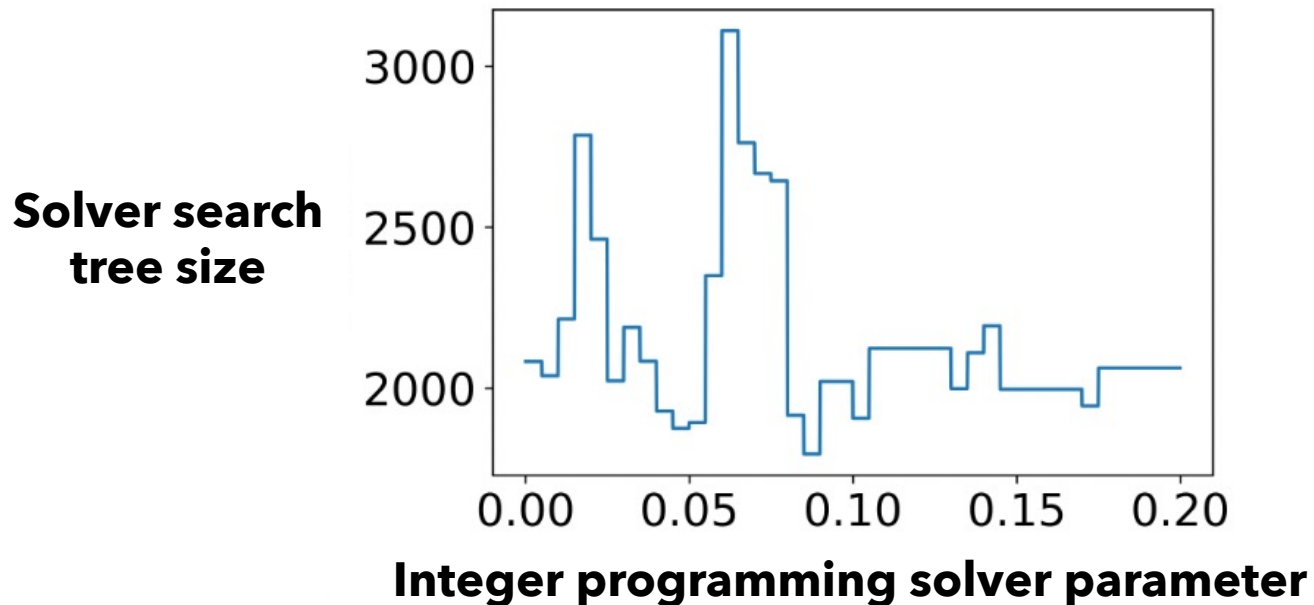
Balcan, Sandholm, **V**, EC'18

Ties to a long line of research on machine learning for **revenue maximization**

Likhodedov, Sandholm, AAAI'04, '05; Balcan, Blum, Hartline, Mansour, FOCS'05; Elkind, SODA'07;
Cole, Roughgarden, STOC'14; Mohri, Medina, ICML'14; Devanur, Huang, Psomas, STOC'16; ...

Primary challenge

Algorithmic performance is a **volatile** function of parameters
Complex connection between parameters and performance



Outline (theoretical guarantees)

1. Statistical guarantees for algorithm configuration
 - i. **Model**
 - ii. Piecewise-structured algorithmic performance
 - iii. Main result
 - iv. Applications
2. Online algorithm configuration

Model

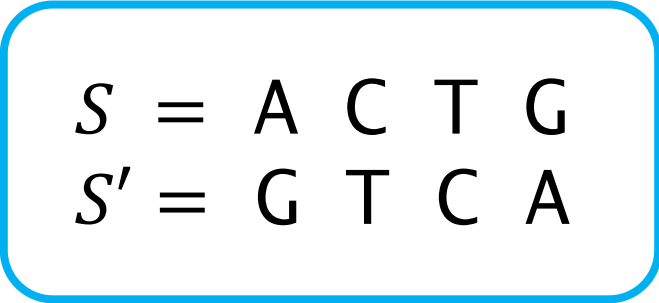
\mathbb{R}^d : Set of all parameters

\mathcal{X} : Set of all inputs

Example: Sequence alignment

\mathbb{R}^3 : Set of alignment algorithm parameters

\mathcal{X} : Set of sequence pairs



$S = A C T G$
 $S' = G T C A$

One sequence pair $x = (S, S') \in \mathcal{X}$

Algorithmic performance

$u_{\rho}(x)$ = utility of algorithm parameterized by $\rho \in \mathbb{R}^d$ on input x
E.g., runtime, solution quality, distance to ground truth, ...

Assume $u_{\rho}(x) \in [-1, 1]$

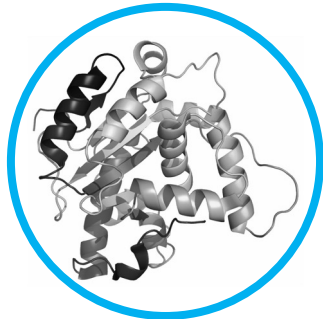
Can be generalized to $u_{\rho}(x) \in [-H, H]$

Model

Standard assumption: Unknown distribution \mathcal{D} over inputs
Distribution models specific application domain at hand



E.g., distribution over pairs of DNA strands



E.g., distribution over pairs of protein sequences

Generalization bounds

Key question: For any parameter setting ρ ,
is **average** utility on training set close to **expected** utility?

Formally: Given samples $x_1, \dots, x_N \sim \mathcal{D}$, for any ρ ,

$$\left| \underbrace{\frac{1}{N} \sum_{i=1}^N u_{\rho}(x_i)}_{\text{Empirical average utility}} - \underbrace{\mathbb{E}_{x \sim \mathcal{D}}[u_{\rho}(x)]}_{\text{Expected utility}} \right| \leq ?$$

Good **average empirical** utility \rightarrow Good **expected** utility

Outline (theoretical guarantees)

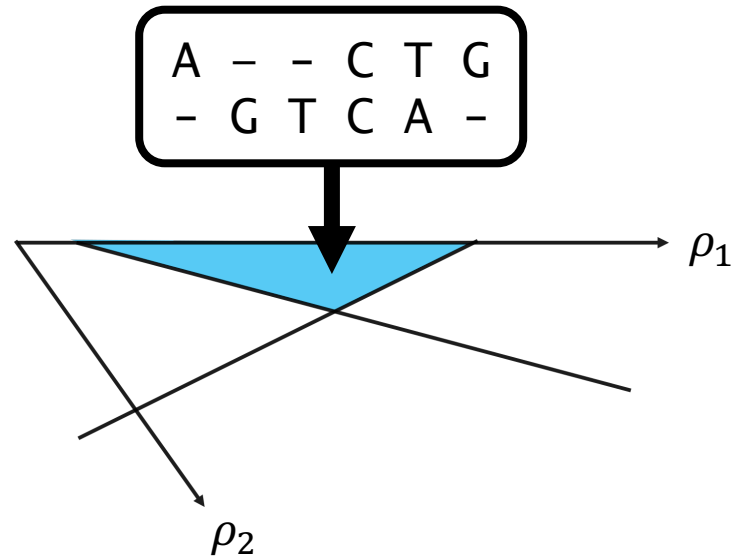
1. Statistical guarantees for algorithm configuration
 - i. Model
 - ii. Piecewise-structured algorithmic performance**
 - a. Example: Sequence alignment**
 - b. Dual function definition
 - iii. Main result
 - iv. Applications
2. Online algorithm configuration

Sequence alignment algorithms

Lemma:

For any pair S, S' , there's a partition of \mathbb{R}^3 s.t. in any region, algorithm's output is fixed across all parameters in region

$S = A C T G$
 $S' = G T C A$

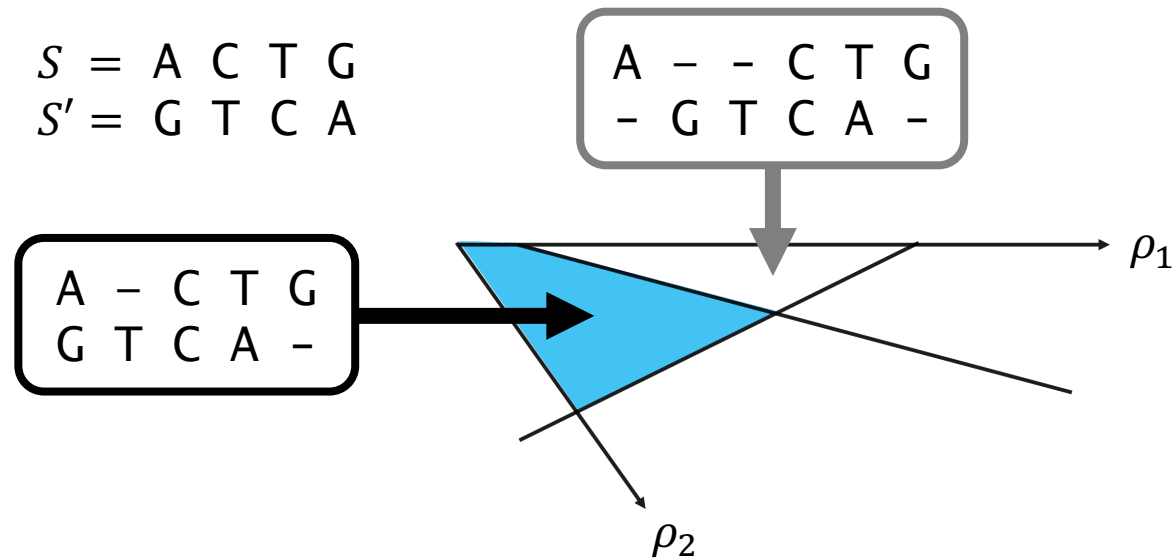


Sequence alignment algorithms

Lemma:

Defined by $(\max\{|S|, |S'|\})^3$ hyperplanes

For any pair S, S' , there's a partition of \mathbb{R}^3 s.t. in any region, algorithm's output is fixed across all parameters in region

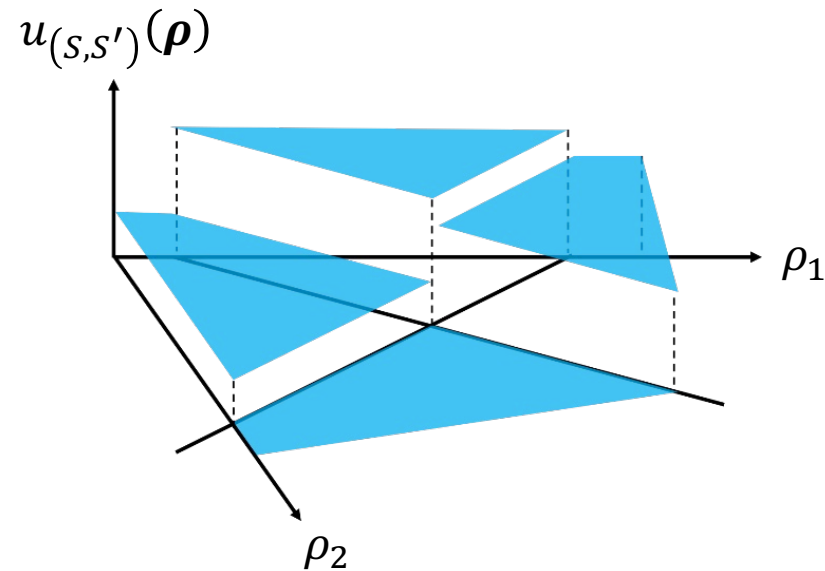


Piecewise-constant utility function

Corollary:

Utility is piecewise constant function of parameters

Distance between algorithm's output and ground-truth alignment



Outline (theoretical guarantees)

1. Statistical guarantees for algorithm configuration
 - i. Model
 - ii. Piecewise-structured algorithmic performance
 - a. Example: Sequence alignment
 - b. Dual function definition**
 - iii. Main result
 - iv. Applications
2. Online algorithm configuration

Primal & dual classes

$u_{\rho}(x)$ = utility of algorithm parameterized by $\rho \in \mathbb{R}^d$ on input x
 $\mathcal{U} = \{u_{\rho}: \mathcal{X} \rightarrow \mathbb{R} \mid \rho \in \mathbb{R}^d\}$ “Primal” function class

Typically, prove guarantees by bounding **complexity** of \mathcal{U}

Challenge: \mathcal{U} is gnarly

E.g., in sequence alignment:

- Each domain element is a pair of sequences
- Unclear how to plot or visualize functions u_{ρ}
- No obvious notions of Lipschitz continuity or smoothness to rely on

Primal & dual classes

$u_{\rho}(x)$ = utility of algorithm parameterized by $\rho \in \mathbb{R}^d$ on input x
 $\mathcal{U} = \{u_{\rho}: \mathcal{X} \rightarrow \mathbb{R} \mid \rho \in \mathbb{R}^d\}$ "Primal" function class

$u_x^*(\rho)$ = utility as function of parameters

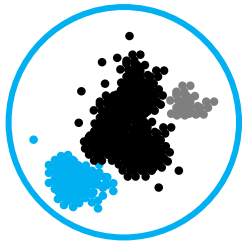
$$u_x^*(\rho) = u_{\rho}(x)$$

$\mathcal{U}^* = \{u_x^*: \mathbb{R}^d \rightarrow \mathbb{R} \mid x \in \mathcal{X}\}$ "Dual" function class

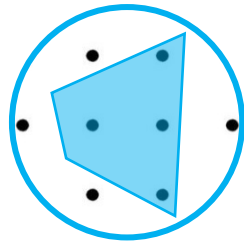
- Dual functions have simple, Euclidean domain
- Often have ample structure can use to bound complexity of \mathcal{U}

Piecewise-structured functions

Dual functions $u_x^*: \mathbb{R}^d \rightarrow \mathbb{R}$ are **piecewise-structured**



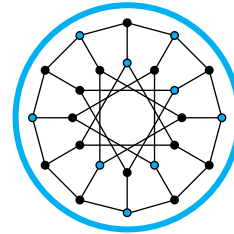
Clustering
algorithm
configuration



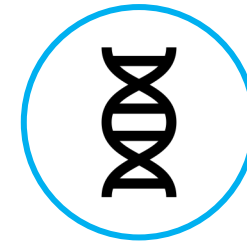
Integer programming
algorithm
configuration



Selling mechanism
configuration



Greedy
algorithm
configuration



Computational biology
algorithm
configuration



Voting mechanism
configuration

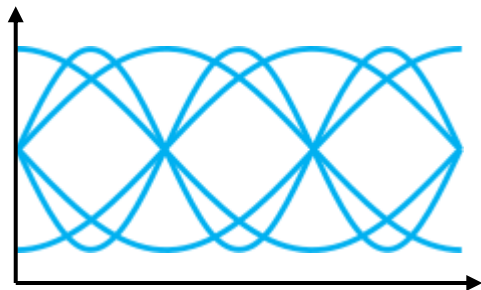
Outline (theoretical guarantees)

1. Statistical guarantees for algorithm configuration
 - i. Model
 - ii. Piecewise-structured algorithmic performance
 - iii. Main result**
 - iv. Applications
2. Online algorithm configuration

Intrinsic complexity

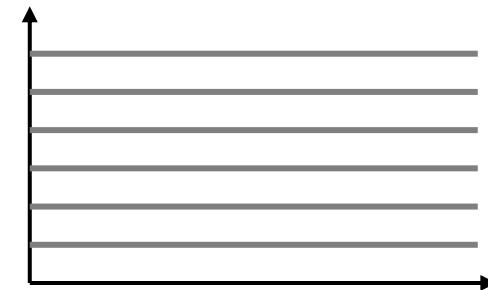
“Intrinsic complexity” of function class \mathcal{G}

- Measures how well functions in \mathcal{G} fit complex patterns
- Specific ways to quantify “intrinsic complexity”:
 - VC dimension
 - Pseudo-dimension



More complex

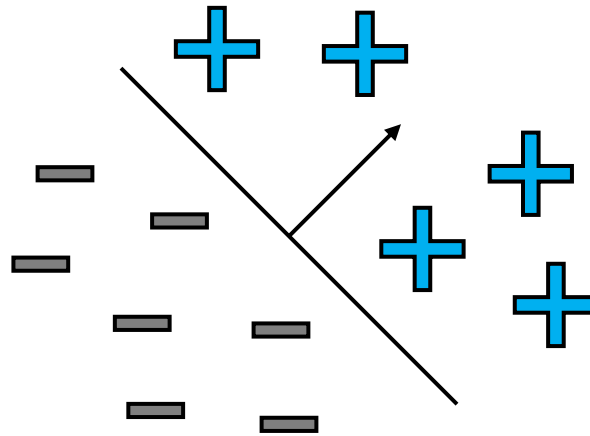
Less complex



VC dimension

Complexity measure for binary-valued function classes \mathcal{F}
(Classes of functions $f: \mathcal{Y} \rightarrow \{-1, 1\}$)

E.g., linear separators



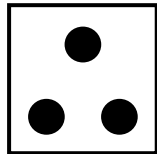
VC dimension

Size of the largest set $\mathcal{S} \subseteq \mathcal{Y}$

that can be labeled in all $2^{|\mathcal{S}|}$ ways by functions in \mathcal{F}

Example: \mathcal{F} = Linear separators in \mathbb{R}^2

$\text{VCdim}(\mathcal{F}) \geq 3$



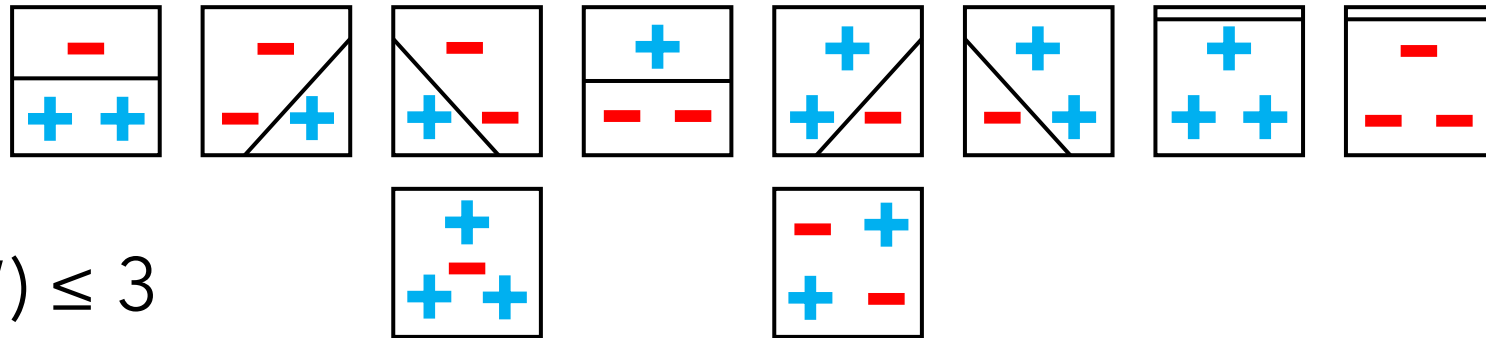
VC dimension

Size of the largest set $\mathcal{S} \subseteq \mathcal{Y}$

that can be labeled in all $2^{|\mathcal{S}|}$ ways by functions in \mathcal{F}

Example: \mathcal{F} = Linear separators in \mathbb{R}^2

$\text{VCdim}(\mathcal{F}) \geq 3$



$\text{VCdim}(\mathcal{F}) \leq 3$

$\text{VCdim}(\{\text{Linear separators in } \mathbb{R}^d\}) = d + 1$

VC dimension

Size of the largest set $\mathcal{S} \subseteq \mathcal{Y}$

that can be labeled in all $2^{|\mathcal{S}|}$ ways by functions in \mathcal{F}

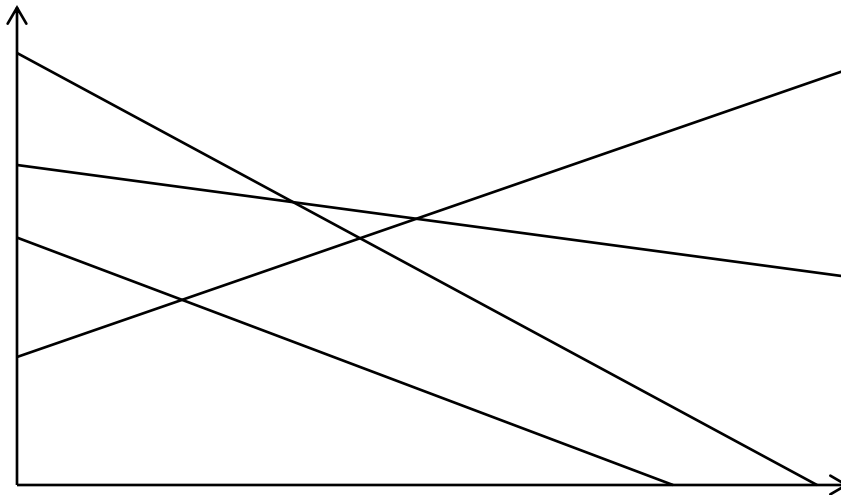
Mathematically, for $\mathcal{S} = \{y_1, \dots, y_N\}$,

$$\left| \left\{ \begin{pmatrix} f(y_1) \\ \vdots \\ f(y_N) \end{pmatrix} : f \in \mathcal{F} \right\} \right| = 2^N$$

Pseudo-dimension

Complexity measure for real-valued function classes \mathcal{G}
(Classes of functions $g: \mathcal{Y} \rightarrow [-1,1]$)

E.g., affine functions



Pseudo-dimension of \mathcal{G}

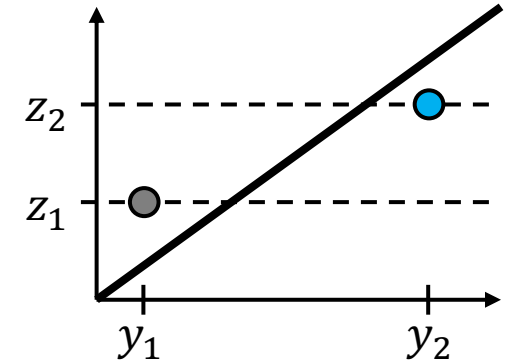
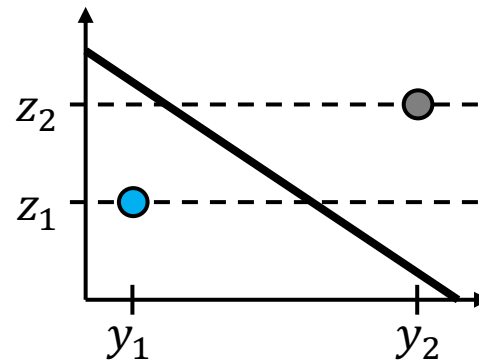
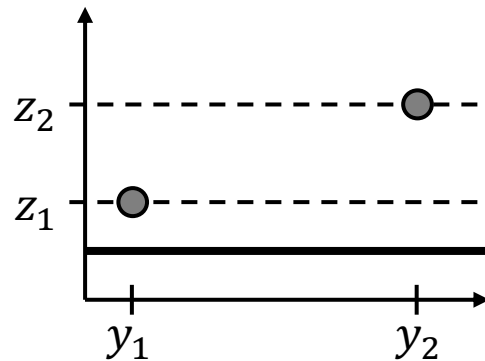
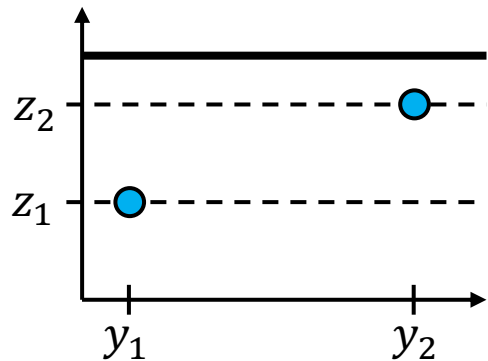
Size of the largest set $\{y_1, \dots, y_N\} \subseteq \mathcal{Y}$ s.t.:

for some *targets* $z_1, \dots, z_N \in \mathbb{R}$,

all 2^N above/below patterns achieved by functions in \mathcal{G}

Example: \mathcal{G} = Affine functions in \mathbb{R}

$\text{Pdim}(\mathcal{G}) \geq 2$



Can also show that $\text{Pdim}(\mathcal{G}) \leq 2$

Pseudo-dimension of \mathcal{G}

Size of the largest set $\{y_1, \dots, y_N\} \subseteq \mathcal{Y}$ s.t.:

for some *targets* $z_1, \dots, z_N \in \mathbb{R}$,

all 2^N above/below patterns achieved by functions in \mathcal{G}

Mathematically,

$$\left| \left\{ \begin{pmatrix} \mathbf{1}_{\{g(y_1) \geq z_1\}} \\ \vdots \\ \mathbf{1}_{\{g(y_N) \geq z_N\}} \end{pmatrix} : g \in \mathcal{G} \right\} \right| = 2^N$$

Sample complexity using pseudo-dim

In the context of **algorithm configuration**:

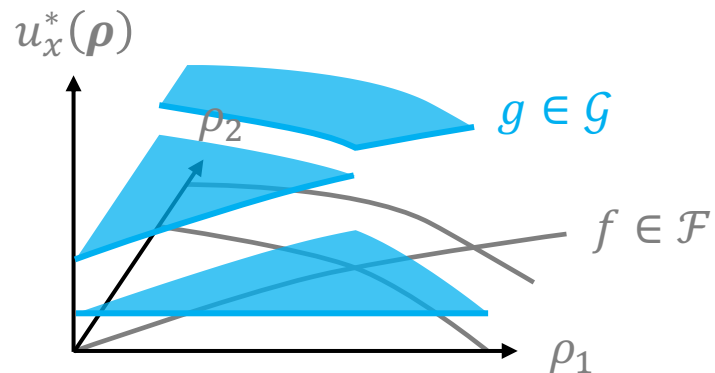
- $\mathcal{U} = \{u_{\rho} : \rho \in \mathbb{R}^d\}$ measure algorithm **performance**
- For $\epsilon, \delta \in (0,1)$, let $N = O\left(\frac{\text{Pdim}(\mathcal{U})}{\epsilon^2} \log \frac{1}{\delta}\right)$
- With probability at least $1 - \delta$ over $x_1, \dots, x_N \sim \mathcal{D}, \forall \rho \in \mathbb{R}^d$,

$$\left| \underbrace{\frac{1}{N} \sum_{i=1}^N u_{\rho}(x_i)}_{\text{Empirical average utility}} - \underbrace{\mathbb{E}_{x \sim \mathcal{D}}[u_{\rho}(x)]}_{\text{Expected utility}} \right| \leq \epsilon$$

Main result (informal)

Boundary functions $f_1, \dots, f_k \in \mathcal{F}$ partition \mathbb{R}^d s.t. in each region, $u_x^*(\boldsymbol{\rho}) = g(\boldsymbol{\rho})$ for some $g \in \mathcal{G}$.

Training set of size $\tilde{O}\left(\frac{\text{Pdim}(\mathcal{G}^*) + \text{VCdim}(\mathcal{F}^*) \log k}{\epsilon^2}\right)$ implies
WHP $\forall \boldsymbol{\rho}, |\text{avg utility over training set} - \text{exp utility}| \leq \epsilon$



Main result (informal)

Boundary functions $f_1, \dots, f_k \in \mathcal{F}$ partition \mathbb{R}^d s.t. in each region,
 $u_x^*(\boldsymbol{\rho}) = g(\boldsymbol{\rho})$ for some $g \in \mathcal{G}$.

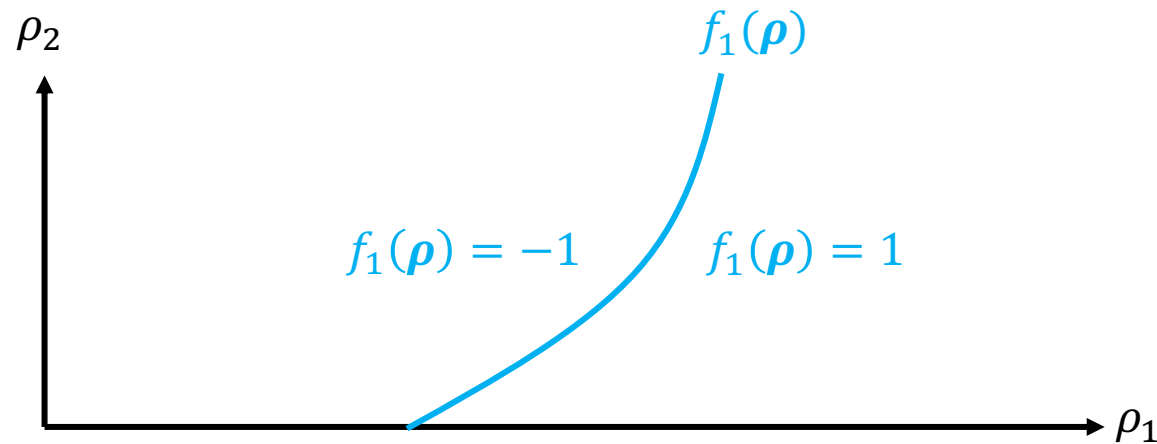
Theorem:

$$\text{Pdim}(\mathcal{U}) = \tilde{O}((\text{VCdim}(\mathcal{F}^*) + \text{Pdim}(\mathcal{G}^*)) \log k)$$

↑
Primal function class $\mathcal{U} = \{u_\rho \mid \rho \in \mathbb{R}^d\}$

Key lemma

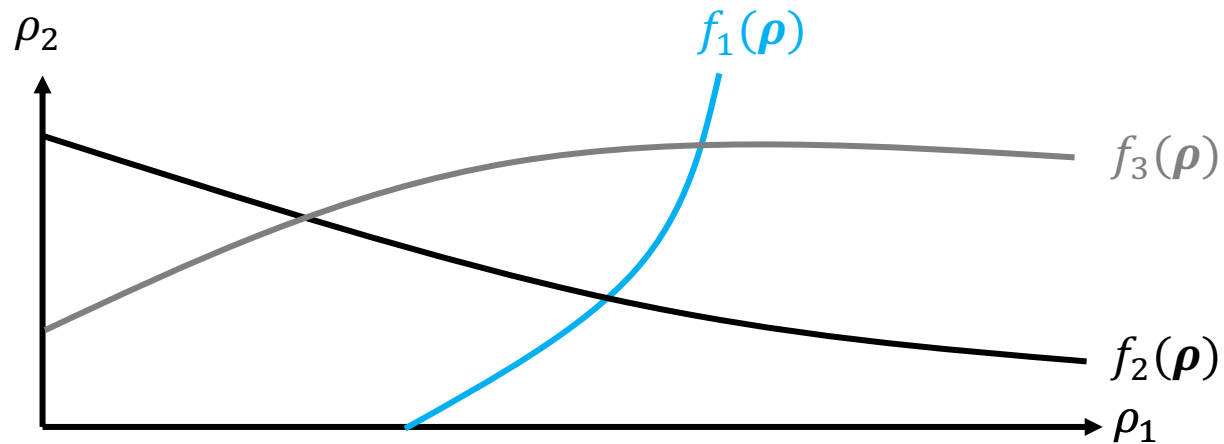
Each boundary function $f: \mathbb{R}^d \rightarrow \{-1, 1\}$ splits \mathbb{R}^d into 2 regions



Key lemma

Given D boundaries, how many sign patterns do they make?

$$\left| \left\{ \begin{pmatrix} f_1(\boldsymbol{\rho}) \\ \vdots \\ f_D(\boldsymbol{\rho}) \end{pmatrix} : \boldsymbol{\rho} \in \mathbb{R}^d \right\} \right| \leq ?$$



Key lemma

Given D boundaries, how many sign patterns do they make?

$$\left| \left\{ \begin{pmatrix} f_1(\boldsymbol{\rho}) \\ \vdots \\ f_D(\boldsymbol{\rho}) \end{pmatrix} : \boldsymbol{\rho} \in \mathbb{R}^d \right\} \right| \leq ?$$

Note: Sauer's lemma tells us that for any D points $\boldsymbol{\rho}_1, \dots, \boldsymbol{\rho}_D \in \mathbb{R}^d$

$$\left| \left\{ \begin{pmatrix} f(\boldsymbol{\rho}_1) \\ \vdots \\ f(\boldsymbol{\rho}_D) \end{pmatrix} : f \in \mathcal{F} \right\} \right| \leq (eD)^{\text{VCdim}(\mathcal{F})}$$

This is where transitioning to the dual comes in handy!

Key lemma

Given D boundaries, how many sign patterns do they make?

$$\left| \left\{ \begin{pmatrix} f_1(\boldsymbol{\rho}) \\ \vdots \\ f_D(\boldsymbol{\rho}) \end{pmatrix} : \boldsymbol{\rho} \in \mathbb{R}^d \right\} \right| \leq (eD)^{\text{VCdim}(\mathcal{F}^*)}$$

Note: Sauer's lemma tells us that for any D points $\boldsymbol{\rho}_1, \dots, \boldsymbol{\rho}_D \in \mathbb{R}^d$

$$\left| \left\{ \begin{pmatrix} f(\boldsymbol{\rho}_1) \\ \vdots \\ f(\boldsymbol{\rho}_D) \end{pmatrix} : f \in \mathcal{F} \right\} \right| \leq (eD)^{\text{VCdim}(\mathcal{F})}$$

This is where transitioning to the dual comes in handy!

Proof ideas

For any problem instances x_1, \dots, x_N and targets $z_1, \dots, z_N \in \mathbb{R}$,

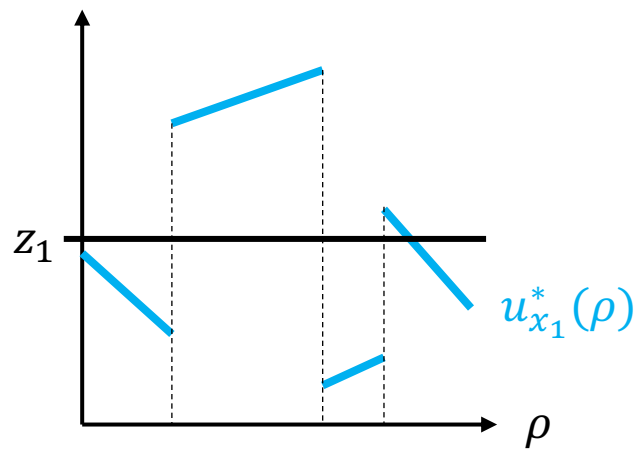
$$\left| \left\{ \begin{pmatrix} \operatorname{sgn}(u_{\boldsymbol{\rho}}(x_1) - z_1) \\ \vdots \\ \operatorname{sgn}(u_{\boldsymbol{\rho}}(x_N) - z_N) \end{pmatrix} : \boldsymbol{\rho} \in \mathbb{R}^d \right\} \right| \leq ?$$

Switching to the dual functions,

$$\left| \left\{ \begin{pmatrix} \operatorname{sgn}(u_{x_1}^*(\boldsymbol{\rho}) - z_1) \\ \vdots \\ \operatorname{sgn}(u_{x_N}^*(\boldsymbol{\rho}) - z_N) \end{pmatrix} : \boldsymbol{\rho} \in \mathbb{R}^d \right\} \right| \leq ?$$

Proof ideas

$$\left| \left\{ \begin{pmatrix} \operatorname{sgn}(u_{x_1}^*(\boldsymbol{\rho}) - z_1) \\ \vdots \\ \operatorname{sgn}(u_{x_N}^*(\boldsymbol{\rho}) - z_N) \end{pmatrix} : \boldsymbol{\rho} \in \mathbb{R}^d \right\} \right| \leq ?$$

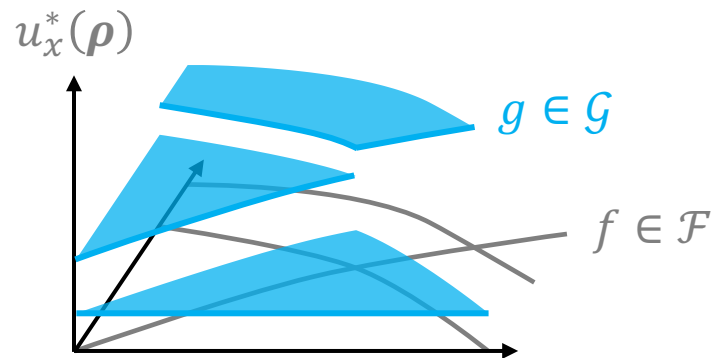


Proof ideas

$$\left| \left\{ \begin{pmatrix} \text{sgn}(u_{x_1}^*(\boldsymbol{\rho}) - z_1) \\ \vdots \\ \text{sgn}(u_{x_N}^*(\boldsymbol{\rho}) - z_N) \end{pmatrix} : \boldsymbol{\rho} \in \mathbb{R}^d \right\} \right| \leq ?$$

The duals $u_{x_1}^*, \dots, u_{x_N}^*$ correspond to Nk boundary functions in \mathcal{F}

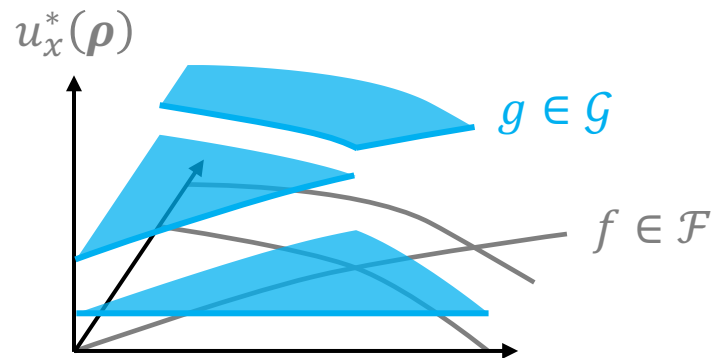
How many regions R_1, \dots, R_M in \mathbb{R}^d ? $M \leq (eNk)^{\text{VCdim}(\mathcal{F}^*)}$



Proof ideas

$$\left| \left\{ \begin{pmatrix} \text{sgn}(u_{x_1}^*(\boldsymbol{\rho}) - z_1) \\ \vdots \\ \text{sgn}(u_{x_N}^*(\boldsymbol{\rho}) - z_N) \end{pmatrix} : \boldsymbol{\rho} \in R_j \right\} \right| \leq ?$$

$\forall \boldsymbol{\rho} \in R_j$, duals are simultaneously structured: $u_{x_i}^*(\boldsymbol{\rho}) = g_i(\boldsymbol{\rho}), \forall i$



Proof ideas

$$\left| \left\{ \begin{pmatrix} \operatorname{sgn}(u_{x_1}^*(\boldsymbol{\rho}) - z_1) \\ \vdots \\ \operatorname{sgn}(u_{x_N}^*(\boldsymbol{\rho}) - z_N) \end{pmatrix} : \boldsymbol{\rho} \in R_j \right\} \right| \leq ?$$

$\forall \boldsymbol{\rho} \in R_j$, duals are simultaneously structured: $u_{x_i}^*(\boldsymbol{\rho}) = g_i(\boldsymbol{\rho}), \forall i$

$$\left| \left\{ \begin{pmatrix} \operatorname{sgn}(g_1(\boldsymbol{\rho}) - z_1) \\ \vdots \\ \operatorname{sgn}(g_N(\boldsymbol{\rho}) - z_N) \end{pmatrix} : \boldsymbol{\rho} \in R_j \right\} \right| \leq ?$$

Proof ideas

$$\left| \left\{ \begin{pmatrix} \text{sgn}(u_{x_1}^*(\boldsymbol{\rho}) - z_1) \\ \vdots \\ \text{sgn}(u_{x_N}^*(\boldsymbol{\rho}) - z_N) \end{pmatrix} : \boldsymbol{\rho} \in R_j \right\} \right| \leq ?$$

$\forall \boldsymbol{\rho} \in R_j$, duals are simultaneously structured: $u_{x_i}^*(\boldsymbol{\rho}) = g_i(\boldsymbol{\rho}), \forall i$

$$\left| \left\{ \begin{pmatrix} \text{sgn}(g_1(\boldsymbol{\rho}) - z_1) \\ \vdots \\ \text{sgn}(g_N(\boldsymbol{\rho}) - z_N) \end{pmatrix} : \boldsymbol{\rho} \in R_j \right\} \right| \leq \underbrace{(eN)^{\text{Pdim}(g^*)}}$$

Follows from key lemma

Proof ideas

$$\left| \left\{ \begin{pmatrix} \text{sgn}(u_{x_1}^*(\boldsymbol{\rho}) - z_1) \\ \vdots \\ \text{sgn}(u_{x_N}^*(\boldsymbol{\rho}) - z_N) \end{pmatrix} : \boldsymbol{\rho} \in \mathbb{R}^d \right\} \right|$$

$$\leq \underbrace{(eNk)^{\text{VCdim}(\mathcal{F}^*)}}_{\text{Number of regions}} \underbrace{(eN)^{\text{Pdim}(\mathcal{G}^*)}}_{\text{Number of sign patterns within each region}}$$

Number of regions

Number of sign patterns within each region

$\text{Pdim}(\mathcal{U})$ equals largest N s.t. $2^N \leq (eNk)^{\text{VCdim}(\mathcal{F}^*)} (eN)^{\text{Pdim}(\mathcal{G}^*)}$,
so $\text{Pdim}(\mathcal{U}) = \tilde{O}((\text{VCdim}(\mathcal{F}^*) + \text{Pdim}(\mathcal{G}^*)) \log k)$

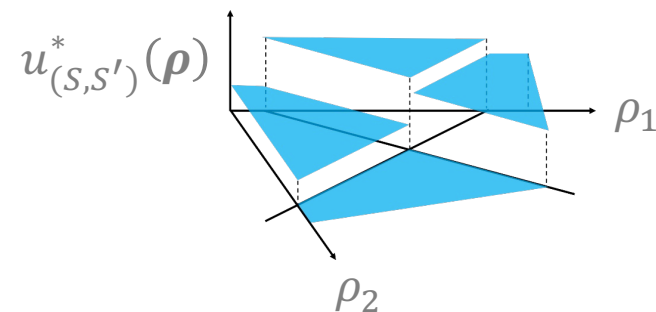
Outline (theoretical guarantees)

1. Statistical guarantees for algorithm configuration
 - i. Model
 - ii. Piecewise-structured algorithmic performance
 - iii. Main result
 - iv. Applications**
 - a. Sequence alignment**
 - b. Greedy algorithms
 - c. Cutting planes
2. Online algorithm configuration

Piecewise constant dual functions

Lemma:

Utility is piecewise constant function of parameters

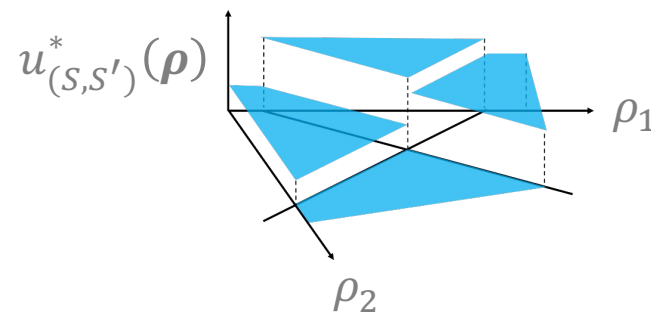


Sequence alignment guarantees

Theorem: Training set of size

$$\tilde{O}\left(\frac{\text{Pdim}(\mathcal{G}^*) + \text{VCdim}(\mathcal{F}^*) \log k}{\epsilon^2}\right) = \tilde{O}\left(\frac{\log(\text{max seq. length})}{\epsilon^2}\right)$$

implies WHP $\forall \rho$, |**avg** utility over training set - **exp** utility| $\leq \epsilon$



Sequence alignment guarantees

Theorem: Training set of size

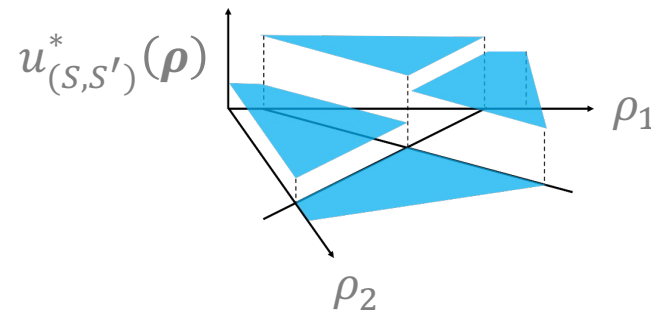
$$\tilde{O}\left(\frac{\text{Pdim}(\mathcal{G}^*) + \text{VCdim}(\mathcal{F}^*) \log k}{\epsilon^2}\right) = \tilde{O}\left(\frac{\log(\text{max seq. length})}{\epsilon^2}\right)$$

\mathcal{G} = constant
functions in \mathbb{R}^3
 $\text{Pdim}(\mathcal{G}^*) = O(1)$

\mathcal{F} = hyperplanes in \mathbb{R}^3
 $\text{VCdim}(\mathcal{F}^*) = O(1)$

$(\text{max sequence length})^3$

implies WHP $\forall \rho, |\mathbf{avg}$ utility over training set - \mathbf{exp} utility $| \leq \epsilon$



Outline (theoretical guarantees)

1. Statistical guarantees for algorithm configuration
 - i. Model
 - ii. Piecewise-structured algorithmic performance
 - iii. Main result
 - iv. Applications
 - a. Sequence alignment
 - b. Greedy algorithms**
 - c. Cutting planes
2. Online algorithm configuration

Example: MWIS

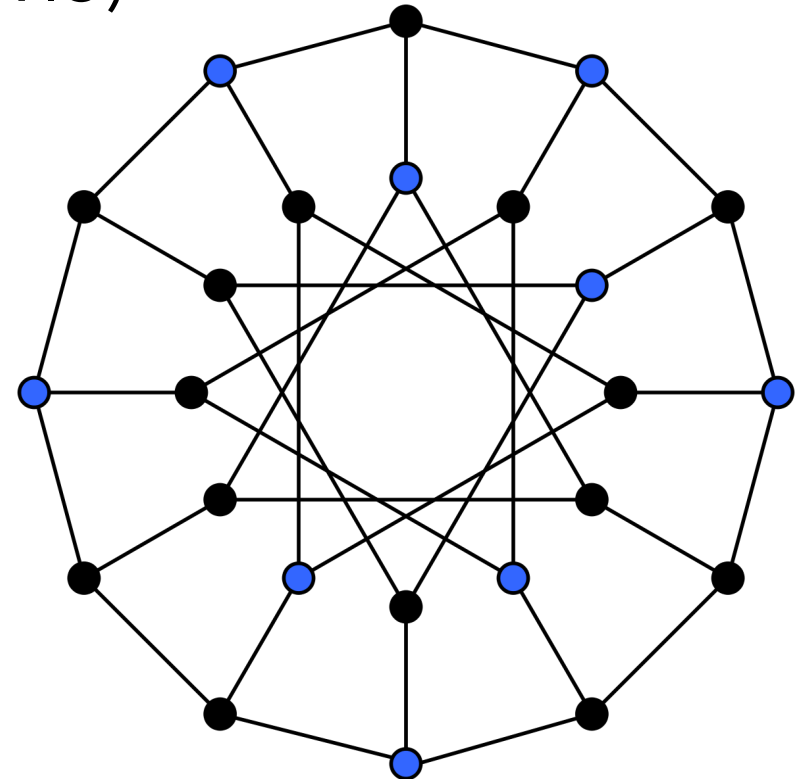
Maximum weight independent set (MWIS)

Problem instance:

- Graph $G = (V, E)$
- n vertices with weights $w_1, \dots, w_n \geq 0$

Goal: find subset $S \subseteq [n]$

- Maximizing $\sum_{i \in S} w_i$
- No nodes $i, j \in S$ are connected: $(i, j) \notin E$



Example: MWIS

Greedy heuristic:

Greedily add vertices v in decreasing order of $\frac{w_v}{(1+\deg(v))}$

Maintaining independence

Parameterized heuristic [Gupta, Roughgarden, ITCS'16]:

Greedily add nodes in decreasing order of $\frac{w_v}{(1+\deg(v))^\rho}$, $\rho \geq 0$

[Inspired by knapsack heuristic by Lehmann et al., JACM'02]

Example: MWIS

Given a MWIS instance x , $u_x^*(\rho) =$ weight of IS algorithm returns

Theorem [Gupta, Roughgarden, ITCS'16]:

$u_x^*(\rho)$ is piecewise-constant with at most n^2 pieces

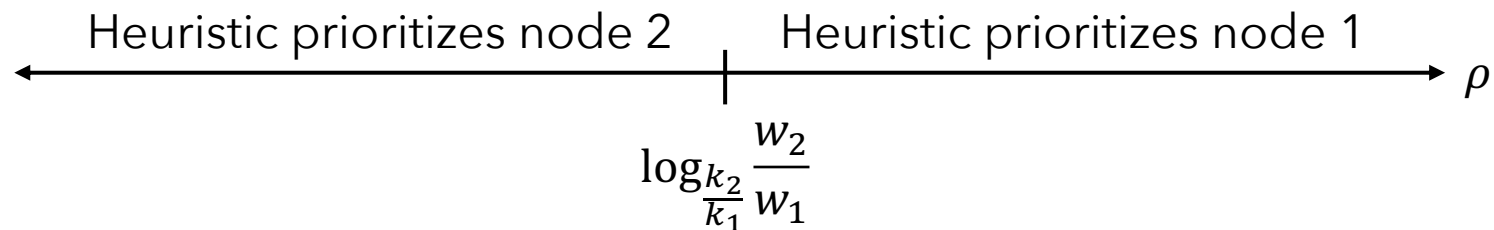
Example: MWIS

Given a MWIS instance x , $u_x^*(\rho) =$ weight of IS algorithm returns

- Weights $w_1, \dots, w_n \geq 0$
- $\deg(i) + 1 = k_i$

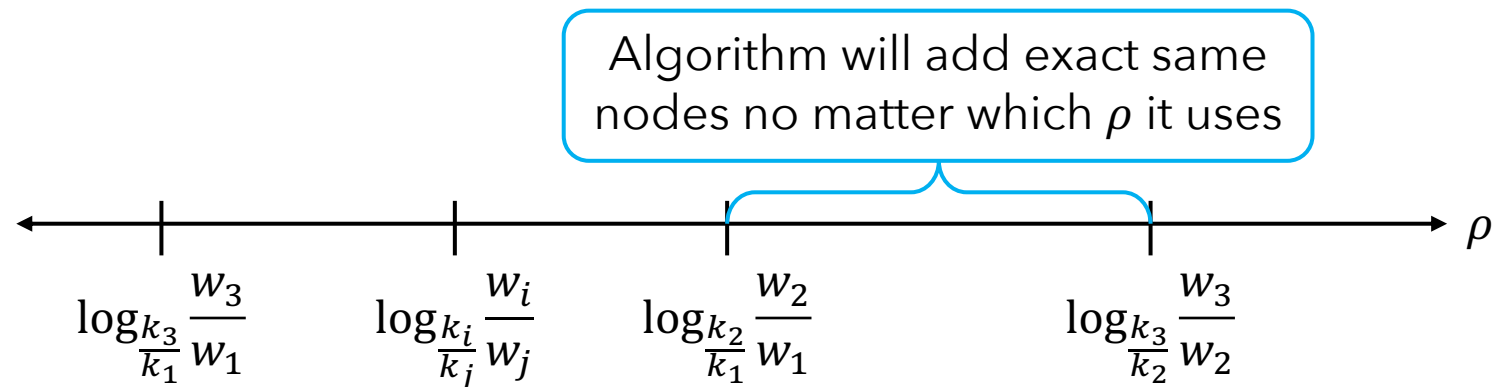
Algorithm parameterized by ρ would add **node 1** before **2** if:

$$\frac{w_1}{k_1^\rho} \geq \frac{w_2}{k_2^\rho} \quad \Leftrightarrow \quad \rho \geq \log_{\frac{k_2}{k_1}} \frac{w_2}{w_1}$$



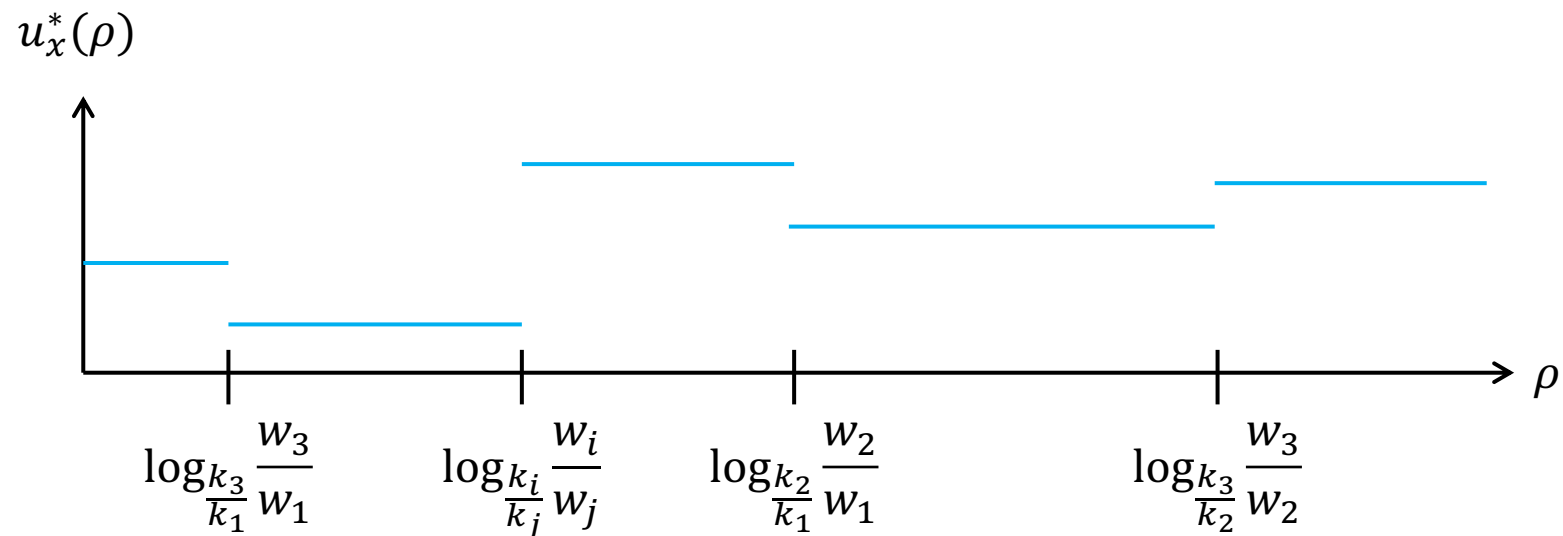
Example: MWIS

- $\binom{n}{2}$ thresholds per instance
- Partition \mathbb{R} into regions where algorithm's output is fixed



Example: MWIS

- $\binom{n}{2}$ thresholds per instance
- Partition \mathbb{R} into regions where algorithm's output is fixed
 $\Rightarrow u_x^*(\rho)$ is constant



MWIS guarantees

Theorem: Training set of size

$$\tilde{O}\left(\frac{\text{Pdim}(\mathcal{G}^*) + \text{VCdim}(\mathcal{F}^*) \log k}{\epsilon^2}\right) = \tilde{O}\left(\frac{\log n}{\epsilon^2}\right)$$

implies WHP $\forall \rho$, **avg** utility over training set - **exp** utility $\leq \epsilon$

MWIS guarantees

Theorem: Training set of size

$$\tilde{O}\left(\frac{\text{Pdim}(\mathcal{G}^*) + \text{VCdim}(\mathcal{F}^*) \log k}{\epsilon^2}\right) = \tilde{O}\left(\frac{\log n}{\epsilon^2}\right)$$

\mathcal{G} = constant functions
 $\text{Pdim}(\mathcal{G}^*) = O(1)$

\mathcal{F} = thresholds
 $\text{VCdim}(\mathcal{F}^*) = O(1)$

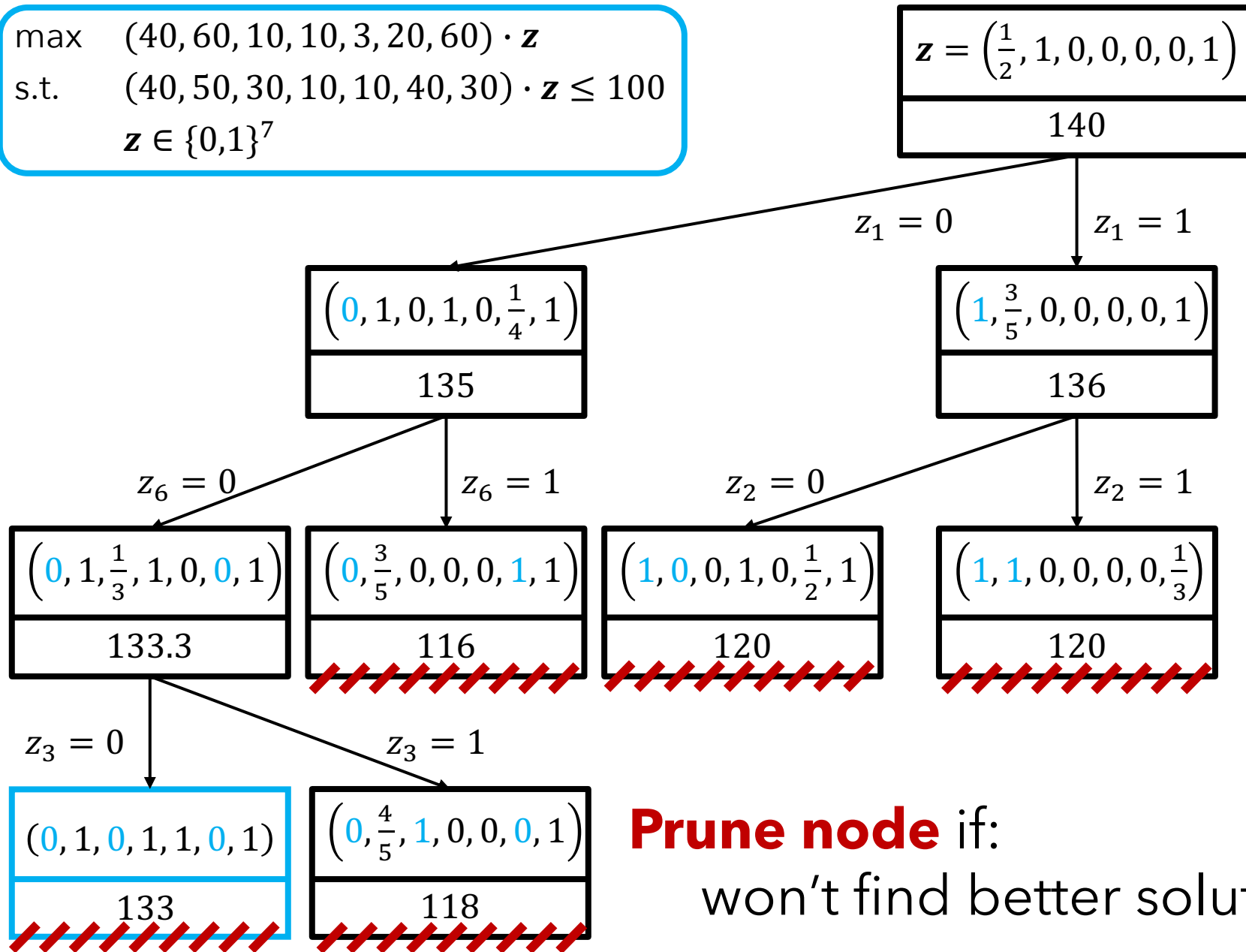
n^2

implies WHP $\forall \rho$, **avg** utility over training set - **exp** utility $\leq \epsilon$

Outline (theoretical guarantees)

1. Statistical guarantees for algorithm configuration
 - i. Model
 - ii. Piecewise-structured algorithmic performance
 - iii. Main result
 - iv. Applications
 - a. Sequence alignment
 - b. Greedy algorithms
 - c. Cutting planes**
2. Online algorithm configuration

$$\begin{aligned} \max & (40, 60, 10, 10, 3, 20, 60) \cdot \mathbf{z} \\ \text{s.t.} & (40, 50, 30, 10, 10, 40, 30) \cdot \mathbf{z} \leq 100 \\ & \mathbf{z} \in \{0,1\}^7 \end{aligned}$$



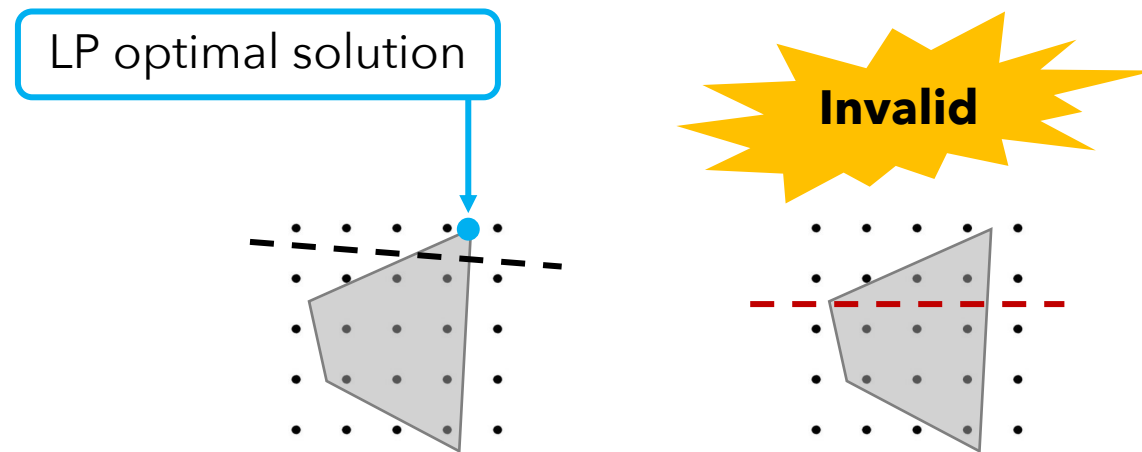
Branch
and
bound
(B&B)

Prune node if:
won't find better solution along branch

Cutting planes

Additional constraints that:

- Separate the LP optimal solution
 - Tightens LP relaxation to prune nodes sooner
- Don't separate any integer point



Cutting planes

Modern IP solvers add cutting planes through the B&B tree
"Branch-and-cut"

Responsible for breakthrough speedups of IP solvers
Cornuéjols, *Annals of OR* '07

Challenges:

- Many different types of cutting planes
 - Chvátal-Gomory cuts, cover cuts, clique cuts, ...
- How to choose which cuts to apply?



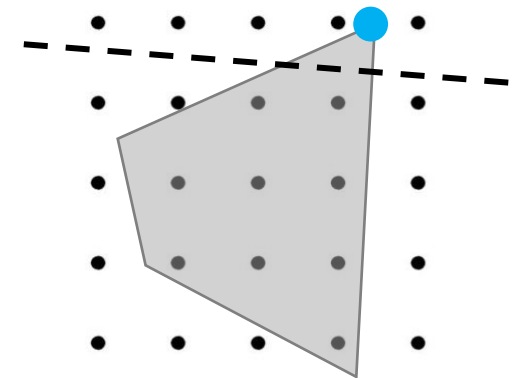
Chvátal-Gomory cuts

We study the canonical family of *Chvátal-Gomory* (CG) cuts

CG cut parameterized by $\boldsymbol{\rho} \in [0,1)^m$ is $\lfloor \boldsymbol{\rho}^T A \rfloor \mathbf{z} \leq \lfloor \boldsymbol{\rho}^T \mathbf{b} \rfloor$

Important properties:

- CG cuts are valid
- Can be chosen so it separates the LP opt



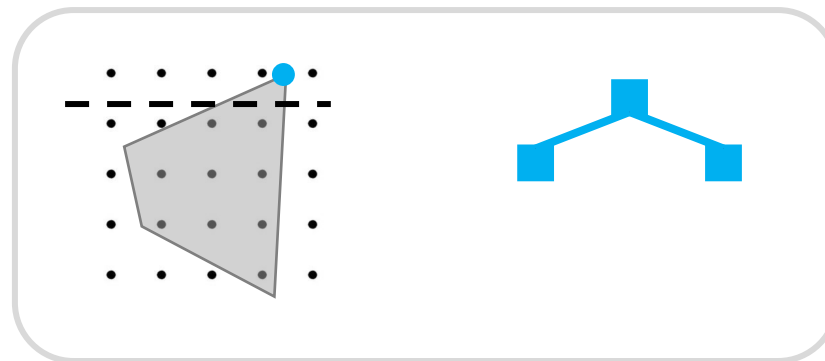
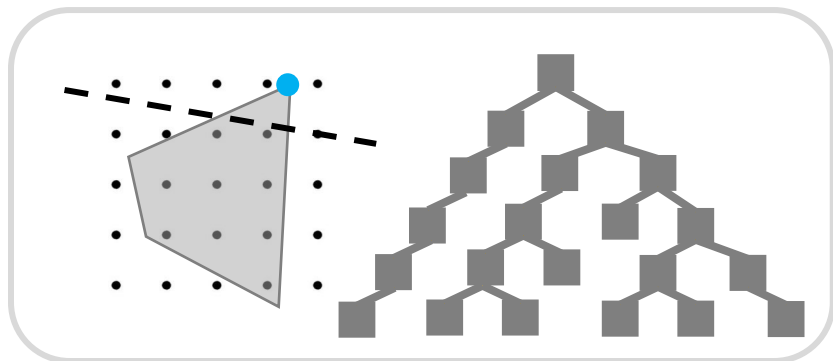
Key challenge

Cut (typically) remains in LPs throughout **entire** tree search

Every aspect of tree search depends on LP guidance

Node selection, variable selection, pruning, ...

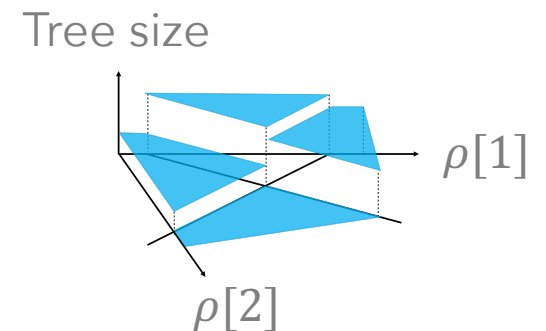
Tiny change in cut can cause **major changes to tree**



Key lemma

Lemma: $O(\|A\|_{1,1} + \|\mathbf{b}\|_1 + n)$ hyperplanes partition $[0,1)^m$ into regions
s.t. in any one region, B&C tree is fixed

Tree size is a piecewise-constant function of $\boldsymbol{\rho} \in [0,1)^m$



Key lemma

Lemma: $O(\|A\|_{1,1} + \|\mathbf{b}\|_1 + n)$ hyperplanes partition $[0,1)^m$ into regions
s.t. in any one region, B&C tree is fixed

Proof idea:

- CG cut parameterized by $\boldsymbol{\rho} \in [0,1)^m$ is $\lfloor \boldsymbol{\rho}^T A \rfloor \mathbf{z} \leq \lfloor \boldsymbol{\rho}^T \mathbf{b} \rfloor$
- For any $\boldsymbol{\rho}$ and column \mathbf{a}_i , $\lfloor \boldsymbol{\rho}^T \mathbf{a}_i \rfloor \in [-\|\mathbf{a}_i\|_1, \|\mathbf{a}_i\|_1]$
- For each integer $k_i \in [-\|\mathbf{a}_i\|_1, \|\mathbf{a}_i\|_1]$:

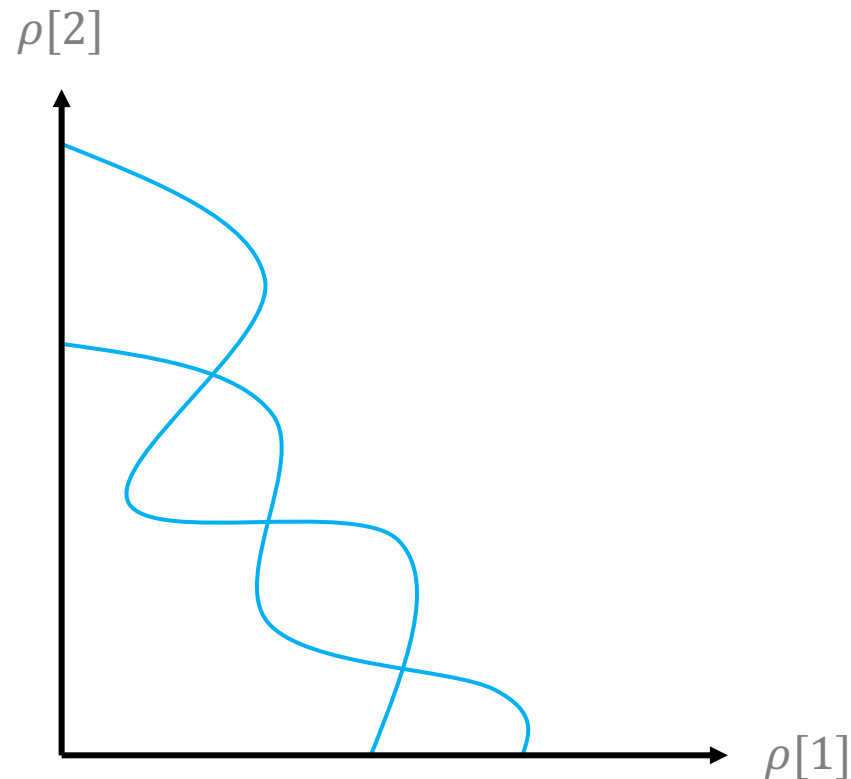
$$\lfloor \boldsymbol{\rho}^T \mathbf{a}_i \rfloor = k_i \text{ iff } k_i \leq \boldsymbol{\rho}^T \mathbf{a}_i < k_i + 1$$

$O(\|A\|_{1,1} + n)$
halfspaces

- In any region defined by intersection of halfspaces:
 $(\lfloor \boldsymbol{\rho}^T \mathbf{a}_1 \rfloor, \dots, \lfloor \boldsymbol{\rho}^T \mathbf{a}_m \rfloor)$ is constant

Beyond Chvátal-Gomory cuts

For more complex families, boundaries can be more complex



Cutting plane guarantees

Theorem: Training set of size

$$\tilde{O}\left(\frac{\text{Pdim}(\mathcal{G}^*) + \text{VCdim}(\mathcal{F}^*) \log k}{\epsilon^2}\right) = \tilde{O}\left(\frac{m \log(\|A\|_{1,1} + \|\mathbf{b}\|_1 + n)}{\epsilon^2}\right)$$

implies WHP $\forall \rho$, **avg** utility over training set - **exp** utility $\leq \epsilon$

Cutting plane guarantees

Theorem: Training set of size

$$\tilde{O}\left(\frac{\text{Pdim}(\mathcal{G}^*) + \text{VCdim}(\mathcal{F}^*) \log k}{\epsilon^2}\right) = \tilde{O}\left(\frac{m \log(\|A\|_{1,1} + \|\mathbf{b}\|_1 + n)}{\epsilon^2}\right)$$

$$\|A\|_{1,1} + \|\mathbf{b}\|_1 + n$$

\mathcal{F} = hyperplanes in \mathbb{R}^m
 $\text{VCdim}(\mathcal{F}^*) = O(m)$

\mathcal{G} = constant functions in \mathbb{R}^m
 $\text{Pdim}(\mathcal{G}^*) = O(m)$

implies WHP $\forall \rho, |\mathbf{avg}$ utility over training set - \mathbf{exp} utility $| \leq \epsilon$

Outline (theoretical guarantees)

1. Statistical guarantees for algorithm configuration
- 2. Online algorithm configuration**

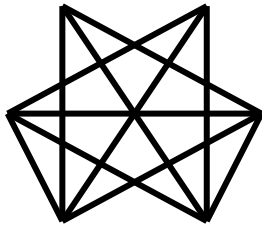
Gupta, Roughgarden, ITCS'16
Balcan, Dick, **Vitercik**, FOCS'18
Balcan, Dick, Pegden, UAI'20

Online algorithm configuration

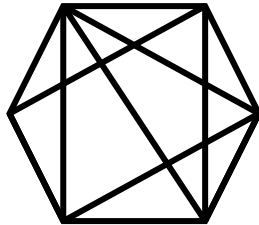
What if inputs are not i.i.d., but even adversarial?

E.g., MWIS:

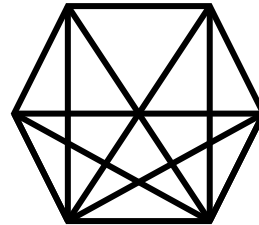
Day 1: ρ_1



Day 2: ρ_2



Day 3: ρ_3




Goal: Compete with best parameter setting in hindsight

- Impossible in the worst case
- Under what conditions is online configuration possible?


Online model

Over T timesteps $t = 1, \dots, T$:

1. Learner chooses **parameter setting** ρ_t
2. Nature (or adversary ) chooses **problem instance** x_t
3. Learner obtains **reward** $u_{\rho_t}(x_t) = u_{x_t}^*(\rho_t)$
4. Learner **observes function** $u_{x_t}^*$ (full information feedback)
 - Simplest setting so we'll start here
 - Will look at other feedback models later (e.g., bandit)

Online model

Over T timesteps $t = 1, \dots, T$:

1. Learner chooses **parameter setting** ρ_t
2. Nature (or adversary ) chooses **problem instance** x_t
3. Learner obtains **reward** $u_{\rho_t}(x_t) = u_{x_t}^*(\rho_t)$
4. Learner **observes function** $u_{x_t}^*$ (full information feedback)

Goal: Minimize **regret** $\max_{\rho} \sum_{t=1}^T u_{\rho}(x_t) - \sum_{t=1}^T u_{\rho_t}(x_t)$

Ideally, $\frac{1}{T} \cdot (\text{Regret}) \rightarrow 0$ as $T \rightarrow \infty$

On average, competing with best algorithm in hindsight

Outline (theoretical guarantees)

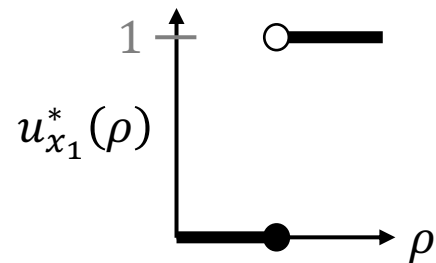
1. Statistical guarantees for algorithm configuration
2. Online algorithm configuration
 - i. **Worst-case instance**
 - ii. Dispersion
 - iii. Semi-bandit model

Worst-case MWIS instance

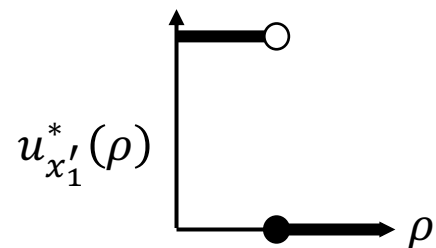
Exists adversary choosing MWIS instances s.t.:

Every full information online algorithm has **linear regret**

Round 1:



Dual function: Utility on instance x_1 as function of ρ



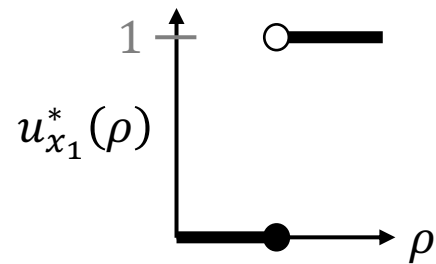
Dual function: Utility on instance x'_1 as function of ρ

Worst-case MWIS instance

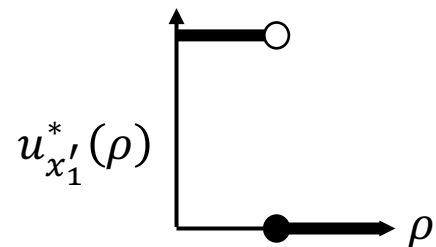
Exists adversary choosing MWIS instances s.t.:

Every full information online algorithm has **linear regret**

Round 1:



Adversary chooses x_1 or x'_1 with equal probability

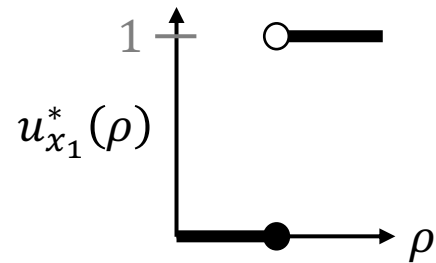


Worst-case MWIS instance

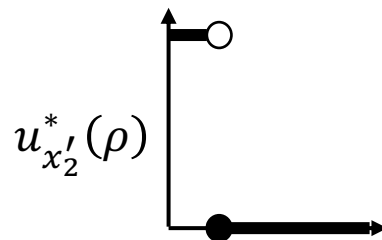
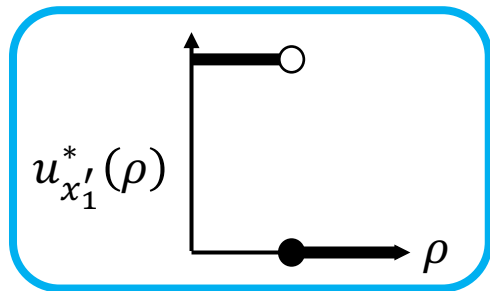
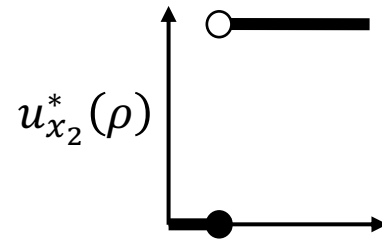
Exists adversary choosing MWIS instances s.t.:

Every full information online algorithm has **linear regret**

Round 1:



Round 2:

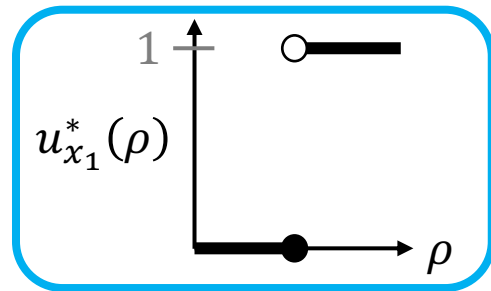


Worst-case MWIS instance

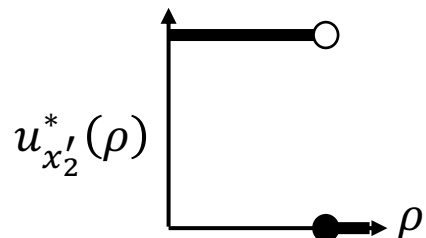
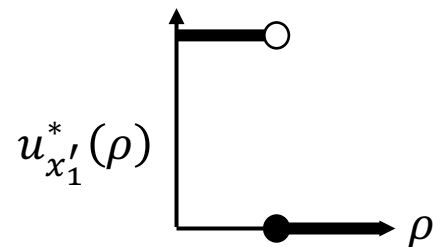
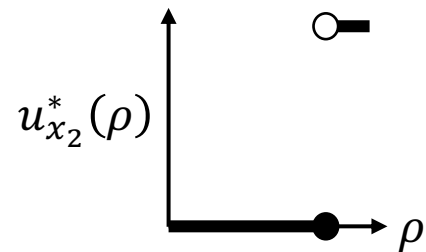
Exists adversary choosing MWIS instances s.t.:

Every full information online algorithm has **linear regret**

Round 1:



Round 2:



Repeatedly halves optimal region

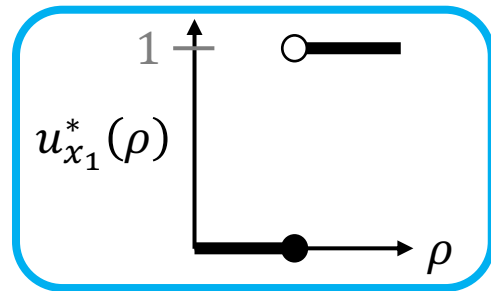


Worst-case MWIS instance

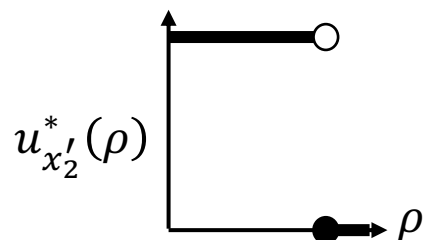
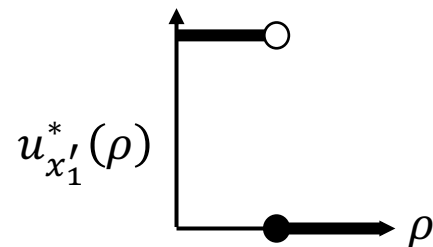
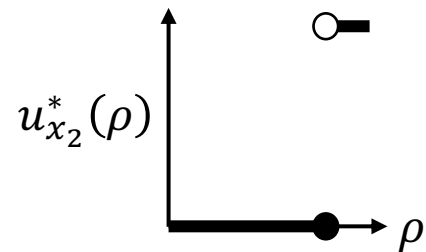
Exists adversary choosing MWIS instances s.t.:

Every full information online algorithm has **linear regret**

Round 1:



Round 2:



Repeatedly halves optimal region

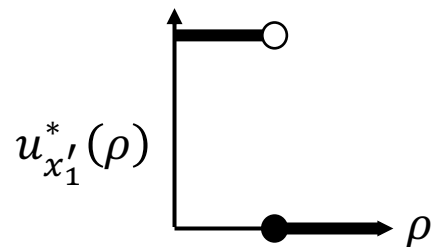
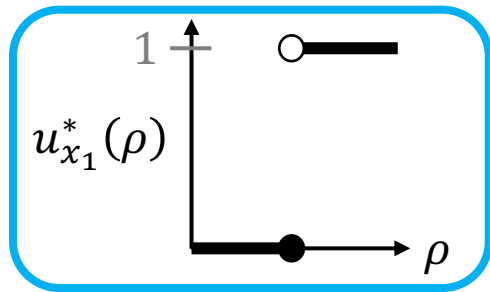


Worst-case MWIS instance

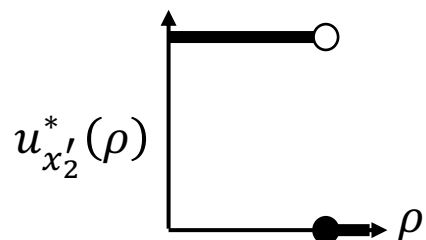
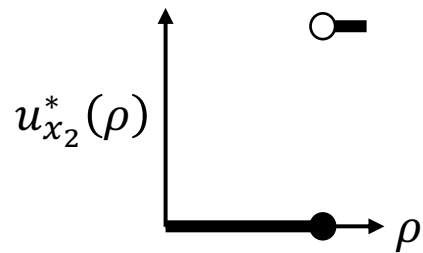
Exists adversary choosing MWIS instances s.t.:

Every full information online algorithm has **linear regret**

Round 1:



Round 2:



Repeatedly halves optimal region



Learner's expected reward: $\frac{T}{2}$

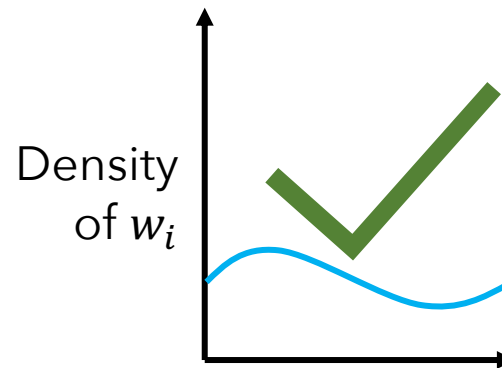
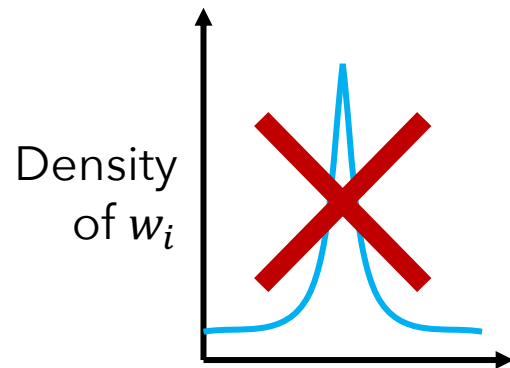
Reward of best ρ in hindsight: T

Expected regret = $\frac{T}{2}$

Smoothed adversary: MWIS

Sub-linear regret is possible if adversary has a “shaky hand”:

- Node weights w_1, \dots, w_n and degrees k_1, \dots, k_n are stochastic
- Joint density of (w_i, w_j, k_i, k_j) is bounded



Later generalized by Cohen-Addad, Kanade [AISTATS, '17]; Balcan, Dick, Vitercik [FOCS'18]; Balcan et al. [UAI'20]; ...

Outline (theoretical guarantees)

1. Statistical guarantees for algorithm configuration
2. Online algorithm configuration
 - i. Worst-case instance
 - ii. Dispersion**
 - iii. Semi-bandit model

Dispersion

Mean adversary concentrates discontinuities near maximizer ρ^*
Even points very close to ρ^* have low utility!

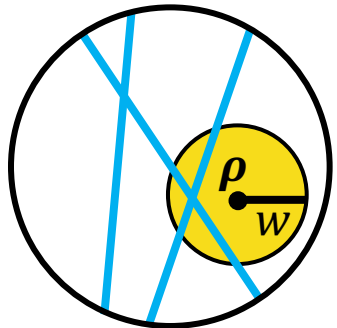
$u_{x_1}^*, \dots, u_{x_T}^* : \underline{B(\mathbf{0}, 1)} \rightarrow [-1, 1]$ are **(w, k) -dispersed at point ρ** if:

Can be generalized to any bounded subset

Dispersion

Mean adversary concentrates discontinuities near maximizer ρ^*
Even points very close to ρ^* have low utility!

$u_{x_1}^*, \dots, u_{x_T}^*: B(\mathbf{0}, 1) \rightarrow [-1, 1]$ are **(w, k) -dispersed at point ρ** if:
 ℓ_2 -ball $B(\rho, w)$ contains discontinuities for $\leq k$ of $u_{x_1}^*, \dots, u_{x_T}^*$



Ball of radius w about ρ contains 2 discontinuities
 $\Rightarrow (w, 2)$ -dispersed at ρ

Outline (theoretical guarantees)

1. Statistical guarantees for algorithm configuration
2. Online algorithm configuration
 - i. Worst-case instance
 - ii. Dispersion
 - a. Algorithm**
 - b. Regret bound
 - c. Bandit feedback
 - d. Proving dispersion holds
 - iii. Semi-bandit model

Exponentially weighted forecaster

[Freund, Schapire, JCSS'97, Cesa-Bianchi & Lugosi '06, ...]

input: Learning rate $\eta > 0$

initialization: $U_0(\boldsymbol{\rho}) = 0$ is the constant function

for $t = 1, \dots, T$:

choose distribution \mathbf{q}_t over \mathbb{R}^d such that $\mathbf{q}_t(\boldsymbol{\rho}) \propto \exp(\eta U_{t-1}(\boldsymbol{\rho}))$

Exponentially upweight high-performance parameter settings

choose parameter setting $\boldsymbol{\rho}_t \sim \mathbf{q}_t$, receive reward $u_{x_t}^*(\boldsymbol{\rho}_t)$

observe utility function $u_{x_t}^*: B(\mathbf{0}, 1) \rightarrow [0, 1]$

update $U_t = U_{t-1} + u_{x_t}^*$

Outline (theoretical guarantees)

1. Statistical guarantees for algorithm configuration
2. Online algorithm configuration
 - i. Worst-case instance
 - ii. Dispersion
 - a. Algorithm
 - b. Regret bound**
 - c. Bandit feedback
 - d. Proving dispersion holds
 - iii. Semi-bandit model

Regret

$$\text{Regret} = \sum_{t=1}^T u_{x_t}^*(\boldsymbol{\rho}^*) - \sum_{t=1}^T u_{x_t}^*(\boldsymbol{\rho}_t)$$

Theorem: Suppose $u_{x_1}^*, \dots, u_{x_T}^*: B(\mathbf{0}, 1) \rightarrow [0, 1]$ are:

1. Piecewise L -Lipschitz
2. (w, k) -dispersed at $\boldsymbol{\rho}^*$

EWf has regret $O\left(\sqrt{Td \log \frac{1}{w}} + TLw + k\right)$

When is this a good bound?

For $w = \frac{1}{L\sqrt{T}}$ and $k = \tilde{O}(\sqrt{T})$, regret is $\tilde{O}(\sqrt{Td})$

Regret upper bound: Proof sketch

$$W_t = \int_{B(\mathbf{0},1)} \exp(\eta U_t(\boldsymbol{\rho})) d\boldsymbol{\rho} \quad \left(U_t(\boldsymbol{\rho}) = \sum_{\tau=1}^t u_{\tau}^*(\boldsymbol{\rho}) \right)$$

Goal:

Something in terms
of OPT = $\sum_{t=1}^T u_t^*(\boldsymbol{\rho}^*)$

$$\leq \frac{W_T}{W_0} \leq$$

Something in terms
of ALG = $\sum_{t=1}^T u_t^*(\boldsymbol{\rho}_t)$

Learner's performance (ALG) is sufficiently large compared to OPT

Regret upper bound: Proof sketch

$$W_t = \int_{B(\mathbf{0},1)} \exp(\eta U_t(\boldsymbol{\rho})) d\boldsymbol{\rho} \quad \left(U_t(\boldsymbol{\rho}) = \sum_{\tau=1}^t u_{\tau}^*(\boldsymbol{\rho}) \right)$$

Goal:

Something in terms
of $\text{OPT} = \sum_{t=1}^T u_t^*(\boldsymbol{\rho}^*)$

$$\leq \frac{W_T}{W_0} \leq \exp(\text{ALG}(e^{\eta} - 1))$$

Standard
EWF analysis

Regret upper bound: Proof sketch

$$W_t = \int_{B(\mathbf{0},1)} \exp(\eta U_t(\boldsymbol{\rho})) d\boldsymbol{\rho} \quad \left(U_t(\boldsymbol{\rho}) = \sum_{\tau=1}^t u_{\tau}^*(\boldsymbol{\rho}) \right)$$

Goal: Something in terms of $\text{OPT} = \sum_{t=1}^T u_t^*(\boldsymbol{\rho}^*)$ $\leq \frac{W_T}{W_0} \leq \exp(\text{ALG}(e^\eta - 1))$

$$W_T = \int_{B(\mathbf{0},1)} \exp\left(\eta \sum_{t=1}^T u_t^*(\boldsymbol{\rho})\right) d\boldsymbol{\rho} \geq \int_{B(\boldsymbol{\rho}^*,w)} \exp\left(\eta \sum_{t=1}^T u_t^*(\boldsymbol{\rho})\right) d\boldsymbol{\rho}$$

Regret upper bound: Proof sketch

Goal: $\text{Something in terms of } \text{OPT} = \sum_{t=1}^T u_t^*(\boldsymbol{\rho}^*) \leq \frac{W_T}{W_0} \leq \exp(\text{ALG}(e^\eta - 1))$

$$\begin{aligned} W_T &= \int_{B(\mathbf{0},1)} \exp\left(\eta \sum_{t=1}^T u_t^*(\boldsymbol{\rho})\right) d\boldsymbol{\rho} \geq \int_{B(\boldsymbol{\rho}^*,w)} \exp\left(\eta \sum_{t=1}^T u_t^*(\boldsymbol{\rho})\right) d\boldsymbol{\rho} \\ &\geq \int_{B(\boldsymbol{\rho}^*,w)} \exp(\eta(\text{OPT} - k - TLw)) d\boldsymbol{\rho} \\ &= \text{Vol}(B(\boldsymbol{\rho}^*,w)) \exp(\eta(\text{OPT} - k - TLw)) \end{aligned}$$

Regret upper bound: Proof sketch

$$\frac{\text{Vol}(B(\boldsymbol{\rho}^*, w)) \exp(\eta(\text{OPT} - k - TLw))}{\text{Vol}(B(\mathbf{0}, 1))} \leq \frac{W_T}{W_0} \leq \exp(\text{ALG}(e^\eta - 1))$$

Rearranging and setting $\eta = \sqrt{\frac{d}{T} \log \frac{1}{w}}$:

$$\text{Regret} = \text{OPT} - \text{ALG} = O\left(\sqrt{Td \log \frac{1}{w}} + TLw + k\right)$$

Matching lower bound

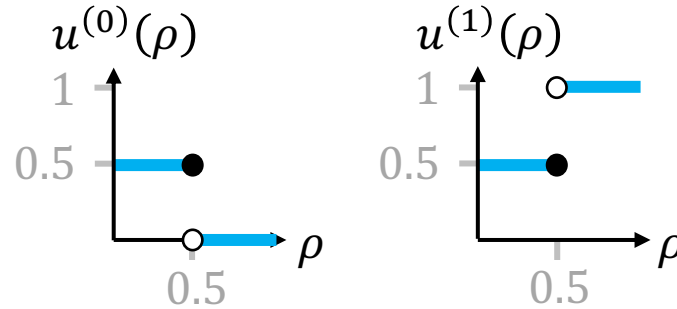
Theorem: For any algorithm, exist PW-constant u_1^*, \dots, u_T^* s.t.:

$$\text{Algorithm's regret is } \Omega \left(\inf_{(w,k)} \sqrt{Td \log \frac{1}{w}} + k \right)$$

Inf over all (w, k) -dispersion parameters that u_1^*, \dots, u_T^* satisfy at ρ^*

$$\text{Upper bound} = O \left(\inf_{(w,k)} \sqrt{Td \log \frac{1}{w}} + k \right)$$

Regret lower bound: Proof sketch



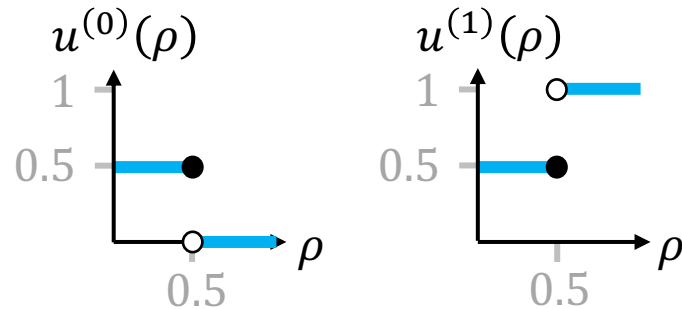
Lemma [Weed et al., COLT'16]:

Exist distributions μ_U, μ_L over $\{u^{(0)}, u^{(1)}\}$ s.t. for any algorithm,

$$\max_{\mu_U, \mu_L} \max_{\rho \in [0, 1]} \mathbb{E} \left[\sum_{t=1}^T u_t^*(\rho) - \sum_{t=1}^T u_t^*(\rho_t) \right] \geq \frac{\sqrt{T}}{32}$$

u_1^*, \dots, u_T^* drawn from worse of μ_U, μ_L

Regret lower bound: Proof sketch



Lemma [Weed et al., COLT'16]:

Exist distributions μ_U, μ_L over $\{u^{(0)}, u^{(1)}\}$ s.t. for any algorithm,

$$\max_{\mu_U, \mu_L} \max_{\rho \in [0, 1]} \mathbb{E} \left[\sum_{t=1}^T u_t^*(\rho) - \sum_{t=1}^T u_t^*(\rho_t) \right] \geq \frac{\sqrt{T}}{32}$$

Any $\rho > 0.5$ is optimal under μ_U , any $\rho \leq 0.5$ is optimal under μ_L

Regret lower bound: Proof sketch

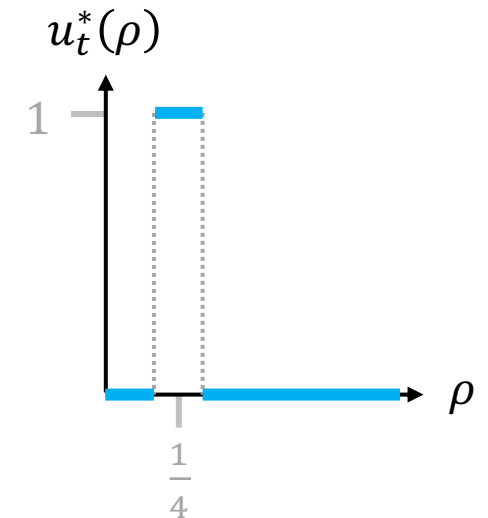
Worst case instance:

1. Draw $u_1^*, \dots, u_{T-\sqrt{T}}^*$ from worse of μ_U, μ_L and define:

$$\rho^* = \operatorname{argmax}_{\rho \in \{\frac{1}{4}, \frac{3}{4}\}} \sum_{t=1}^{T-\sqrt{T}} u_t^*(\rho)$$

2. Define $u_t^*(\rho) = \mathbf{1}_{\{|\rho - \rho^*| \leq \frac{1}{10}\}}$ for $t > T - \sqrt{T}$

Note: $\rho^* \in \operatorname{argmax} \sum_{t=1}^T u_t^*(\rho)$



Regret lower bound: Proof sketch

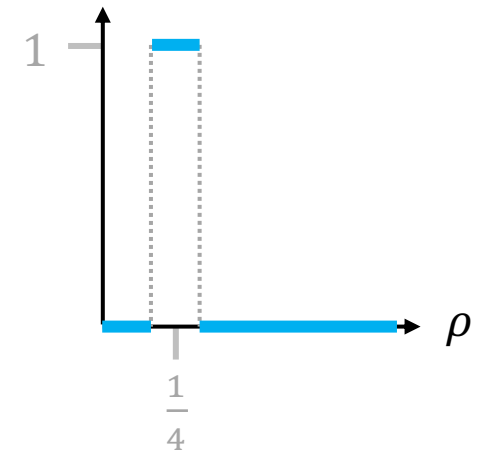
Analysis:

- Regret $\geq \frac{\sqrt{T}}{64}$ (follows from lemma by Weed et al., [COLT'16])
- Lower bound follows from fact that $\frac{\sqrt{T}}{64} = \Omega\left(\inf_{(w,k)} \sqrt{T \log \frac{1}{w} + k}\right)$

Only last $k = \sqrt{T}$ functions have discontinuities in

$$\left[\rho^* - \frac{1}{8}, \rho^* + \frac{1}{8}\right]$$

$\Rightarrow u_1^*, \dots, u_T^*$ are $(w = \frac{1}{8}, k = \sqrt{T})$ -dispersed around ρ^*




Outline (theoretical guarantees)

1. Statistical guarantees for algorithm configuration
2. Online algorithm configuration
 - i. Worst-case instance
 - ii. Dispersion
 - a. Algorithm
 - b. Regret bound
 - c. Bandit feedback**
 - d. Proving dispersion holds
 - iii. Semi-bandit model

Bandit feedback

Over T timesteps $t = 1, \dots, T$:

1. Learner chooses **parameter setting** $\boldsymbol{\rho}_t$
2. Nature (or adversary ) chooses **problem instance** x_t
3. Learner obtains **reward** $u_{\boldsymbol{\rho}_t}(x_t) = u_{x_t}^*(\boldsymbol{\rho}_t)$
4. Learner **only** observes $u_{x_t}^*(\boldsymbol{\rho}_t)$ (not entire function)

Bandit feedback

Theorem: If $u_1^*, \dots, u_T^*: B(\mathbf{0}, 1) \rightarrow [0,1]$ are:

1. Piecewise L -Lipschitz
2. (w, k) -dispersed at ρ^*

The UCB algorithm has regret $\tilde{O} \left(\sqrt{Td \left(\frac{1}{w}\right)^d} + TLw + k \right)$

- If $d = 1$, $w = \frac{1}{\sqrt[3]{T}}$, and $k = \tilde{O}(T^{2/3})$, regret is $\tilde{O}(LT^{2/3})$
- If $w = T^{\frac{d+1}{d+2}-1}$, $k = \tilde{O}(T^{\frac{d+1}{d+2}})$, then regret is $\tilde{O} \left(T^{\frac{d+1}{d+2}} (\sqrt{d3^d} + L) \right)$

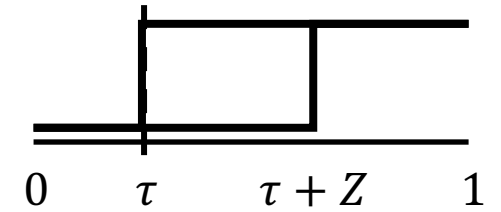
Outline (theoretical guarantees)

1. Statistical guarantees for algorithm configuration
2. Online algorithm configuration
 - i. Worst-case instance
 - ii. Dispersion
 - a. Algorithm
 - b. Regret bound
 - c. Bandit feedback
 - d. Proving dispersion holds**
 - iii. Semi-bandit model

Smooth adversaries and dispersion

Adversary chooses thresholds $u_t^*: [0,1] \rightarrow \{0,1\}$

Discontinuity τ "smoothed" by adding $Z \sim N(0, \sigma^2)$



Lemma: WHP, $\forall w, u_1^*, \dots, u_T^*$ are $\left(w, \tilde{O}\left(\frac{Tw}{\sigma} + \sqrt{T}\right)\right)$ -dispersed

Corollary: $w = \frac{\sigma}{\sqrt{T}} \Rightarrow$ **Full information regret** = $O\left(\sqrt{T \log \frac{T}{\sigma}}\right)$

Simple example: knapsack

Problem instance:

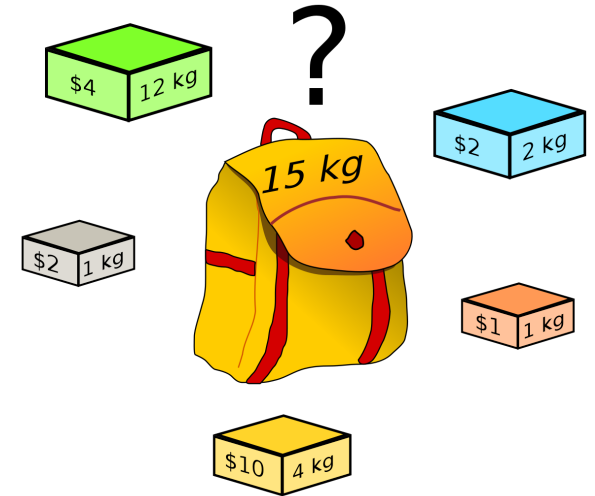
- n items, item i has value v_i and size s_i
- Knapsack with capacity K

Goal: find most valuable items that fit

Algorithm (parameterized by $\rho \geq 0$):

Add items in decreasing order of $\frac{v_i}{s_i^\rho}$

[Gupta and Roughgarden, ITCS'16]



Dispersion for knapsack

Theorem: If instances randomly distributed s.t. on each round:

1. Each v_i independent from s_i
2. All (v_i, v_j) have joint density functions with range $\subseteq [0, \kappa]$,

W.h.p., for any $\alpha \geq \frac{1}{2}$, u_1^*, \dots, u_T^* are

$\left(\tilde{O}\left(\frac{T^{1-\alpha}}{\kappa}\right), \tilde{O}\left((\# \text{ items})^2 T^\alpha\right) \right)$ -dispersed

Corollary: Full information regret = $\tilde{O}\left((\# \text{ items})^2 \sqrt{T}\right)$

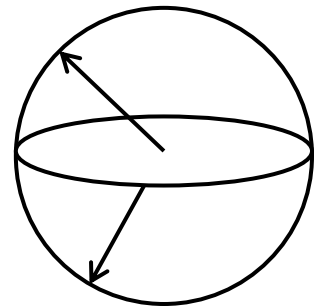
More results for algorithm configuration

Under **no assumptions**, we show dispersion for
Integer quadratic programming approximation algs

Based on semi-definite programming relaxations

- s -linear rounding [Feige & Langberg '06]
- Outward rotations [Zwick '99]
 - Both generalizations of Goemans-Williamson max-cut alg ['95]

Leverage algorithm's randomness to prove dispersion



Outline (theoretical guarantees)

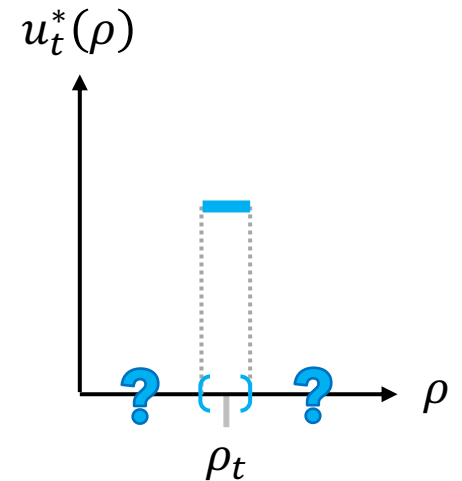
1. Statistical guarantees for algorithm configuration
2. Online algorithm configuration
 - i. Worst-case instance
 - ii. Dispersion
 - iii. Semi-bandit model**

Semi-bandit model

- Computing the entire function $u_t^*(\rho)$ can be challenging
- Often, it's easy to compute interval in which $u_t^*(\rho_t)$ is constant
 - E.g., in IP, simple bookkeeping with CPLEX callbacks
- **Semi-bandit model:** learner learns $u_t^*(\rho_t)$ and interval

Balcan, Dick, Pegden [UAI'20]:

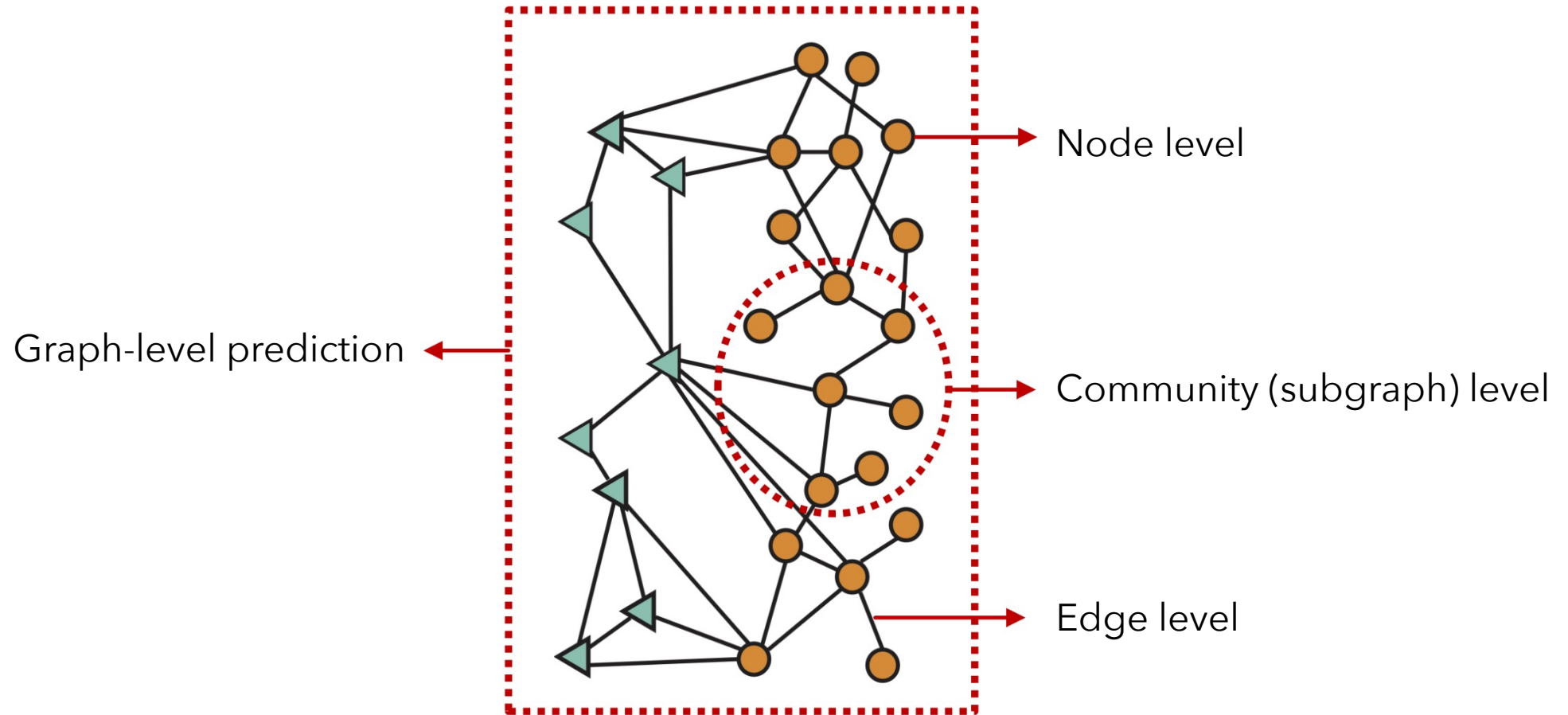
- Regret bounds that are nearly as good as full info
- Introduce a more general definition of dispersion



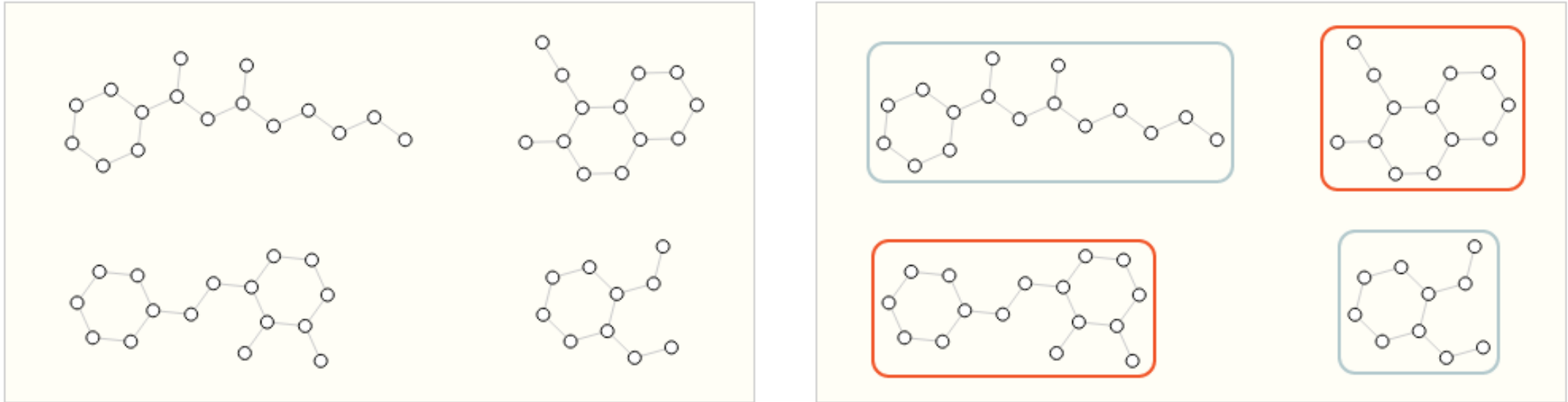
Outline (applied techniques)

- 1. GNNs overview**
2. Neural algorithmic alignment
3. Integer programming with GNNs

Different types of tasks



Prediction with graphs: Examples

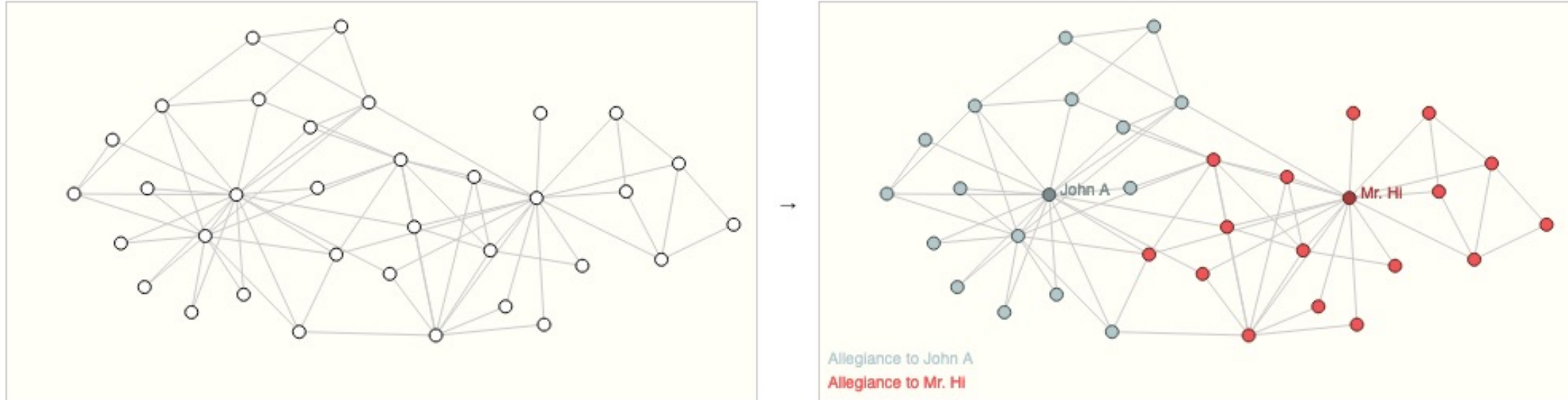


Graph-level tasks:

E.g., for a molecule represented as a graph, could predict:

- What the molecule smells like
- Whether it will bind to a receptor implicated in a disease

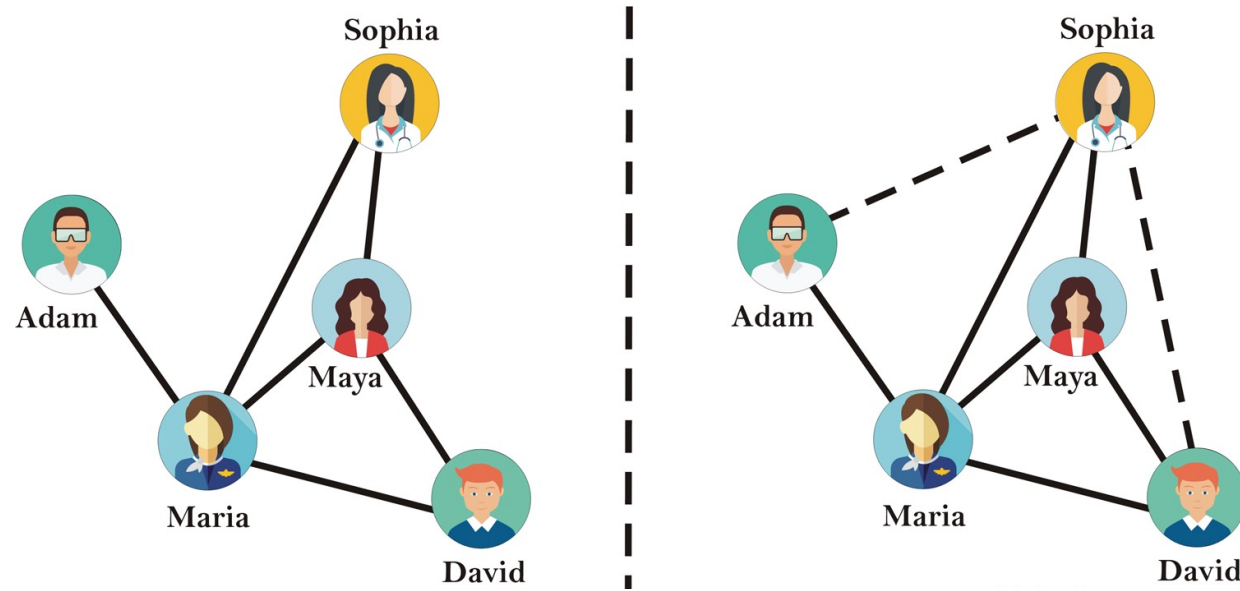
Prediction with graphs: Examples



Node-level tasks:

E.g., political affiliations of users in a social network

Prediction with graphs: Examples



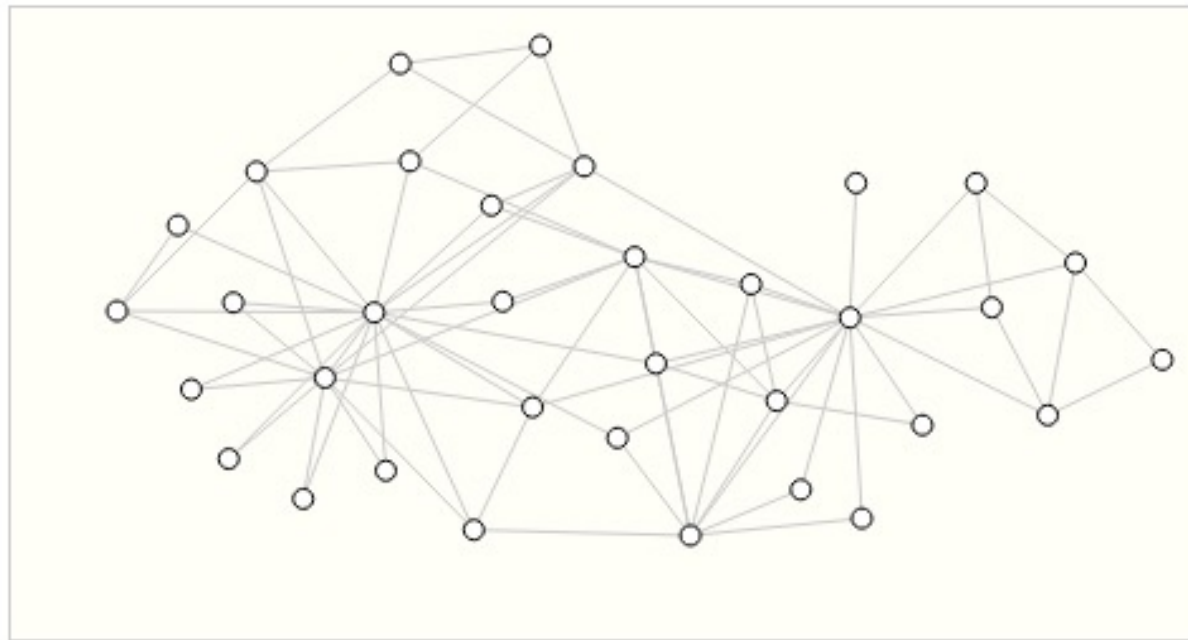
Edge-level tasks: E.g.:

- Suggesting new friends
- Recommendations on Amazon, Netflix, ...

GNN motivation

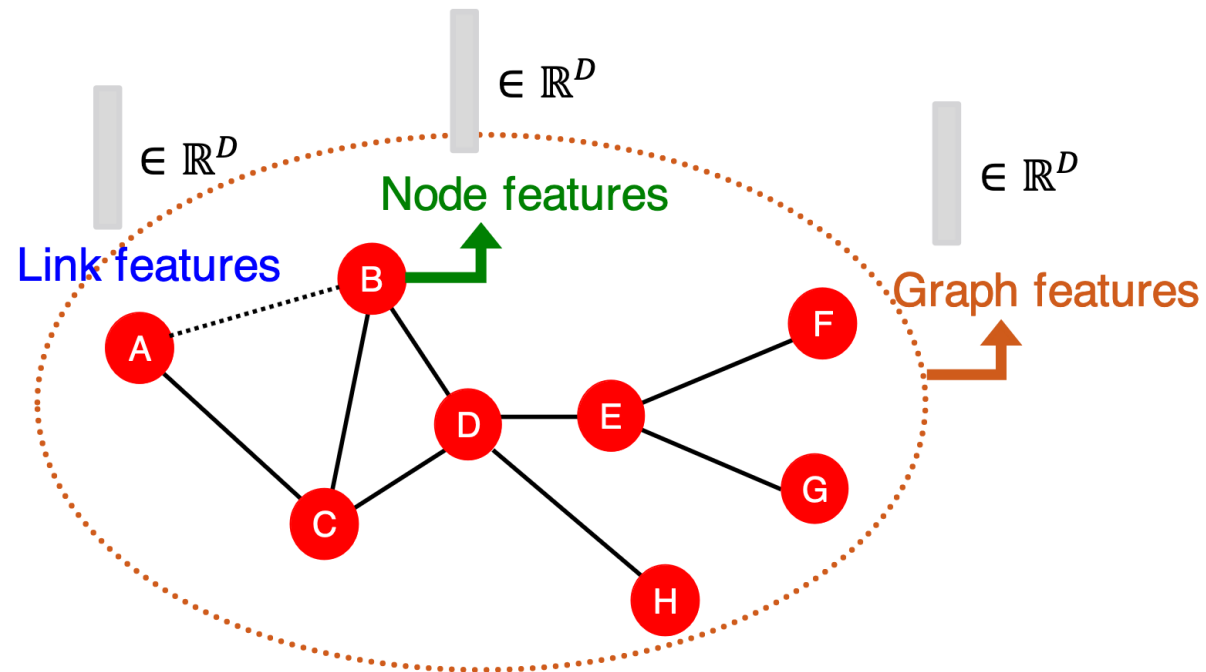
Main question:

How to utilize relational structure for better prediction?



Graph neural networks: First step

- Design features for nodes/links/graphs
- Obtain features for all training data



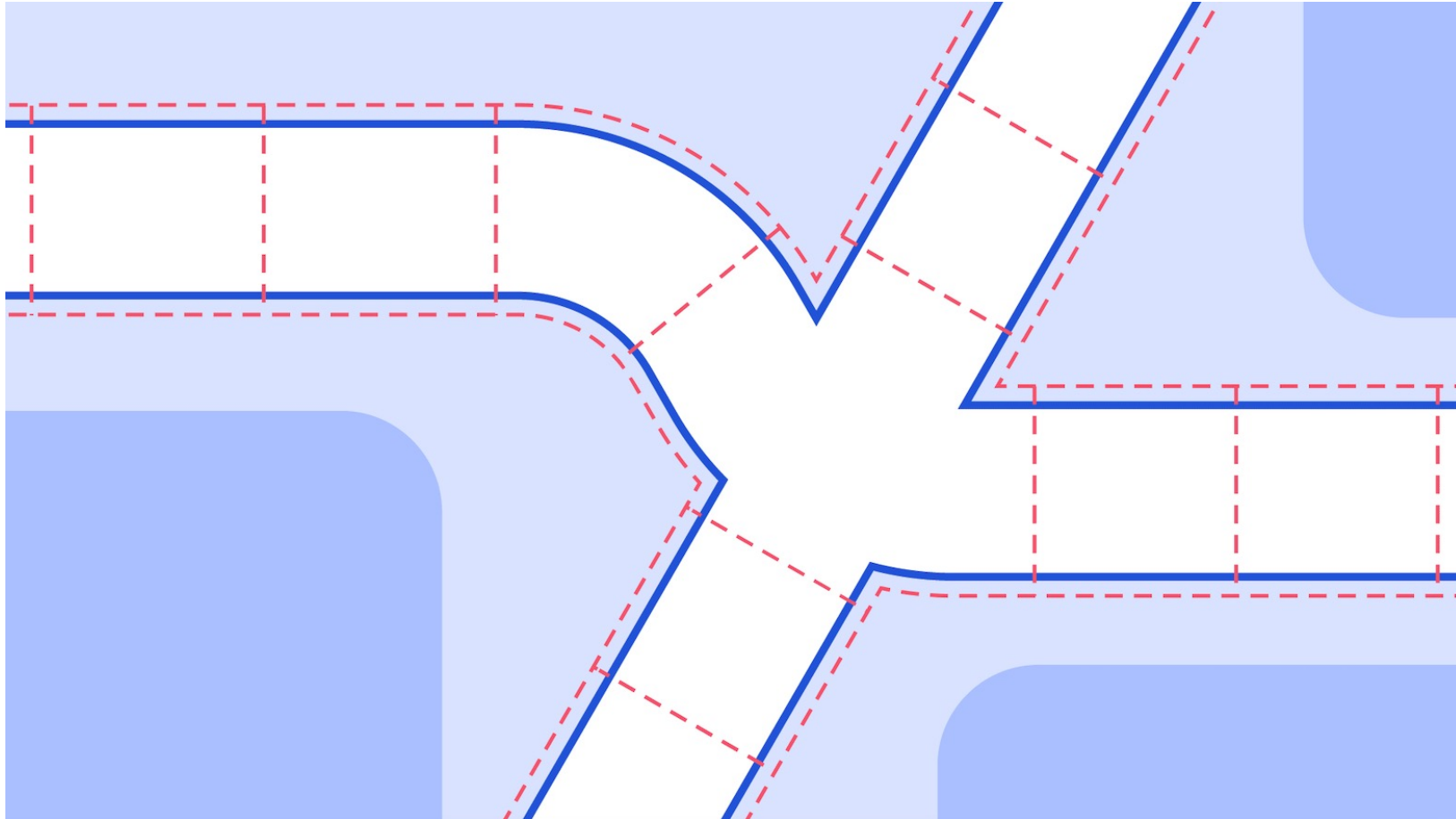
Graph neural networks: Objective

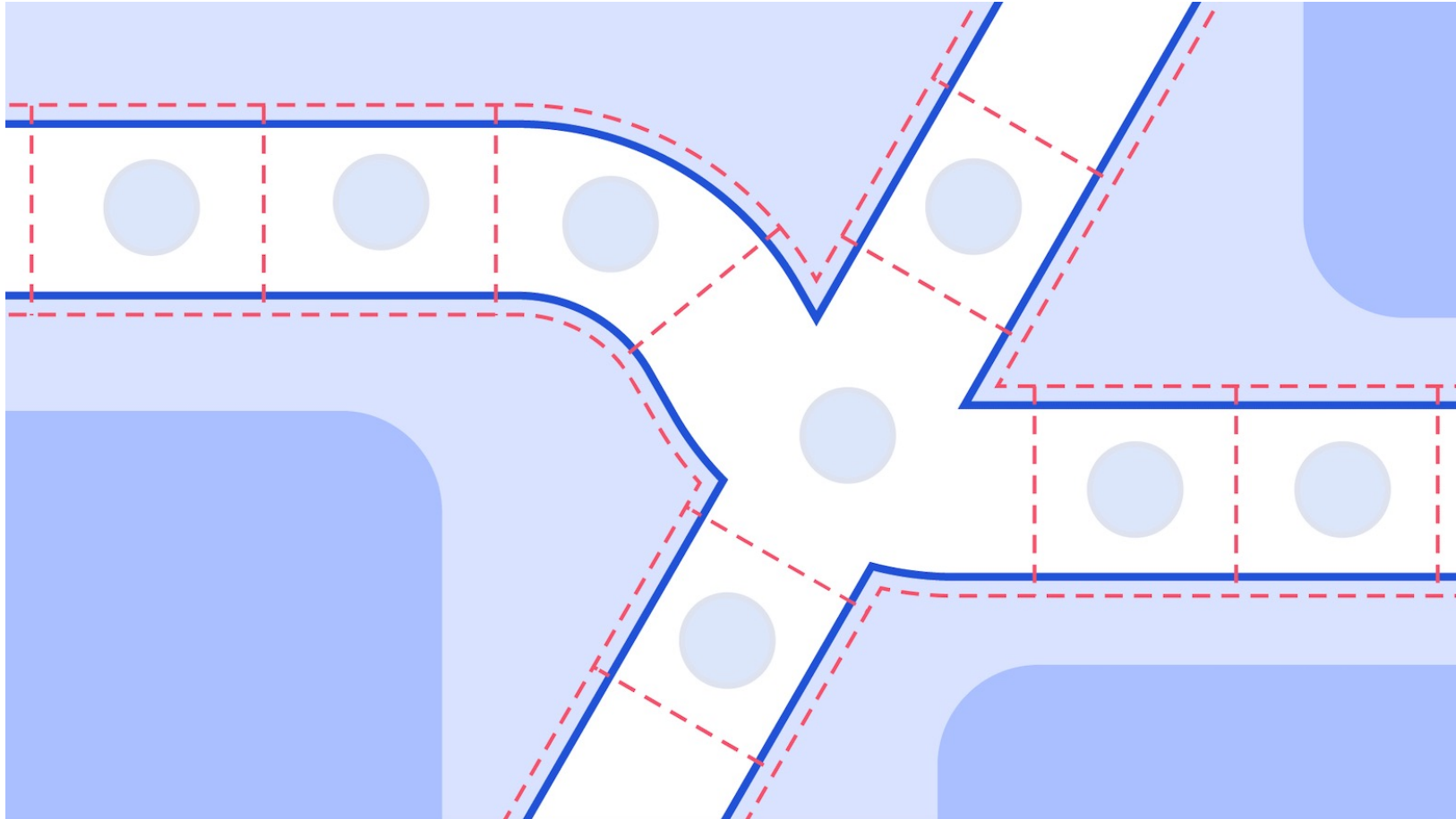
Idea:

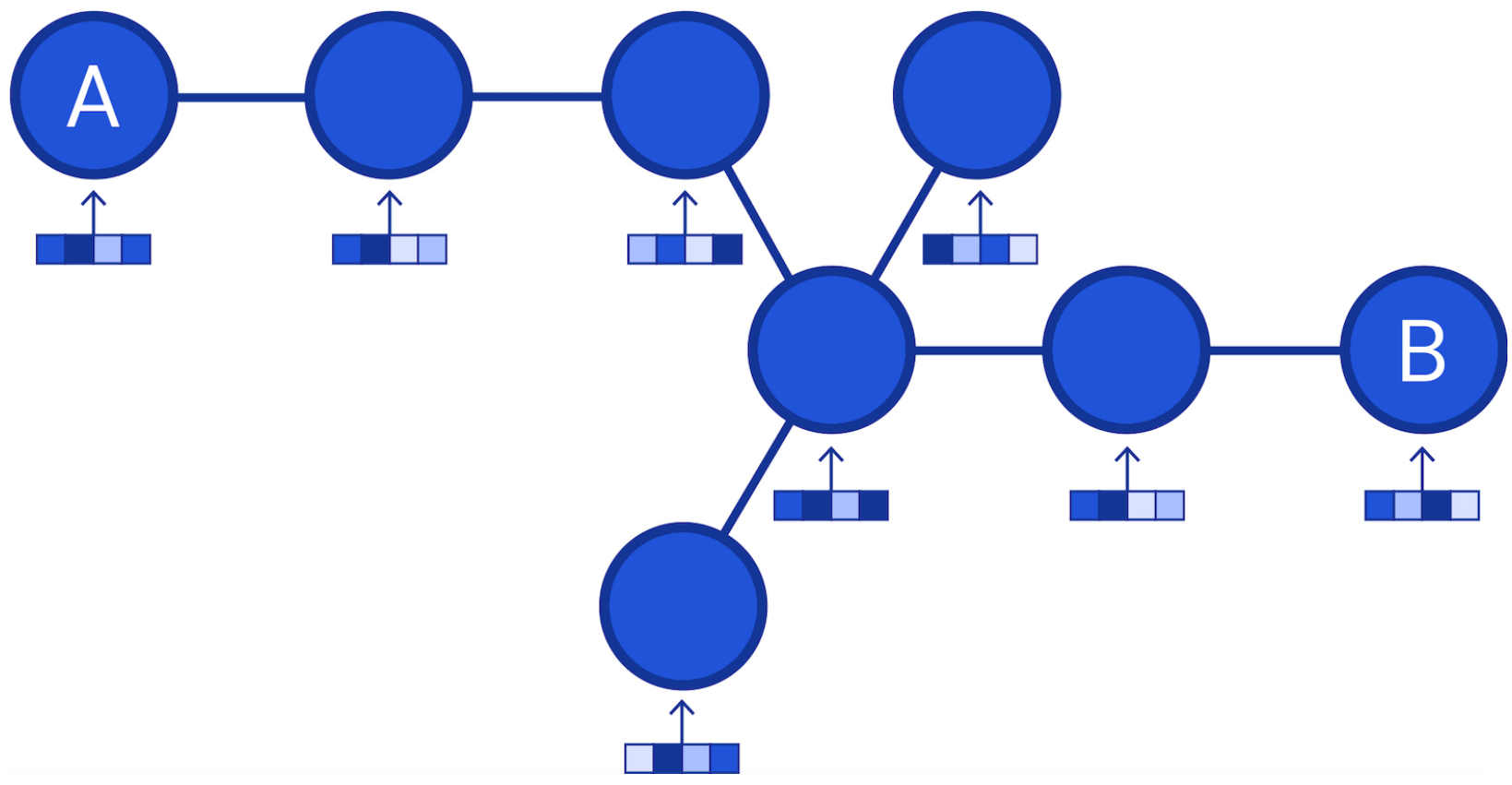
1. Encode each node and its neighborhood with embedding
2. Aggregate set of node embeddings into graph embedding
3. Use embeddings to make predictions

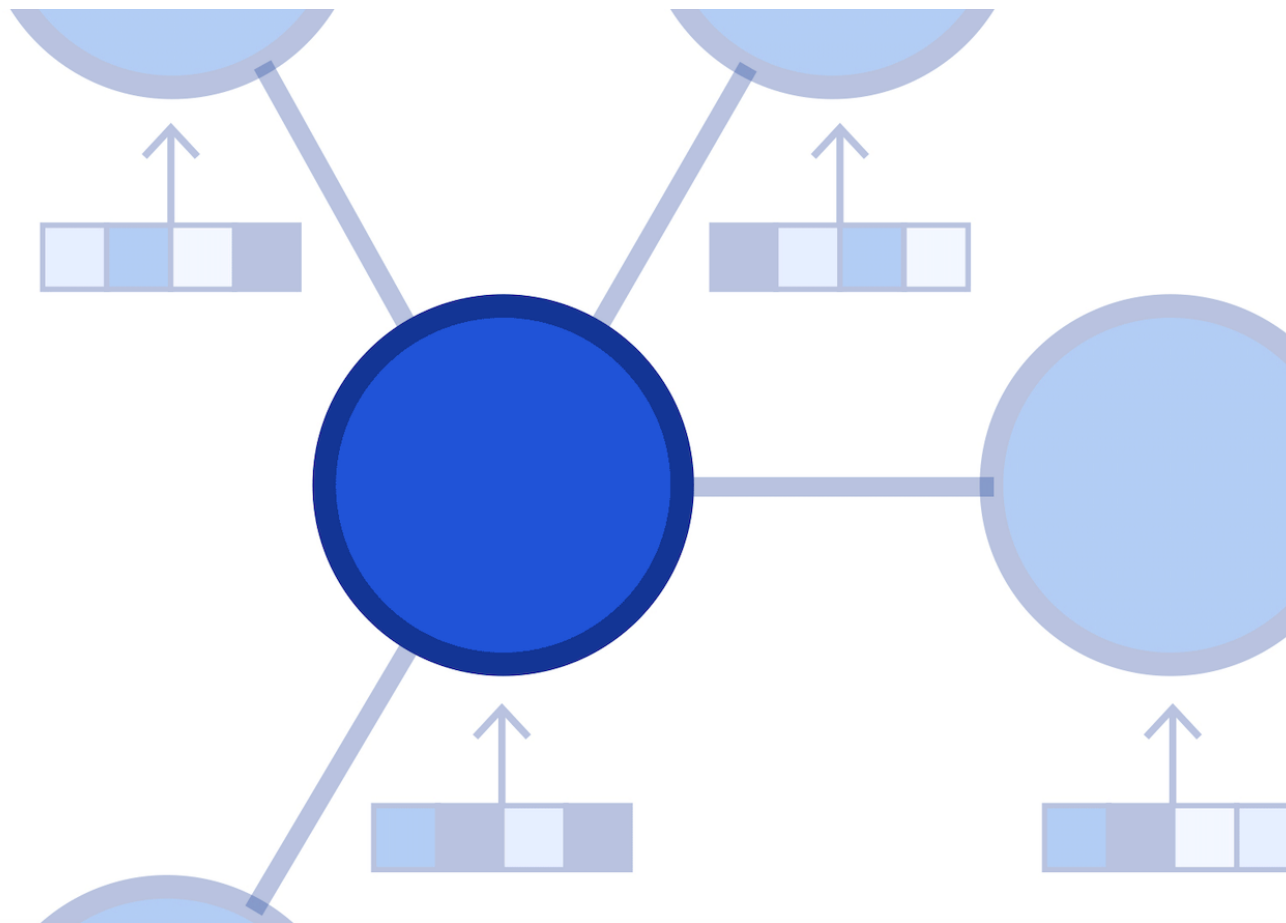


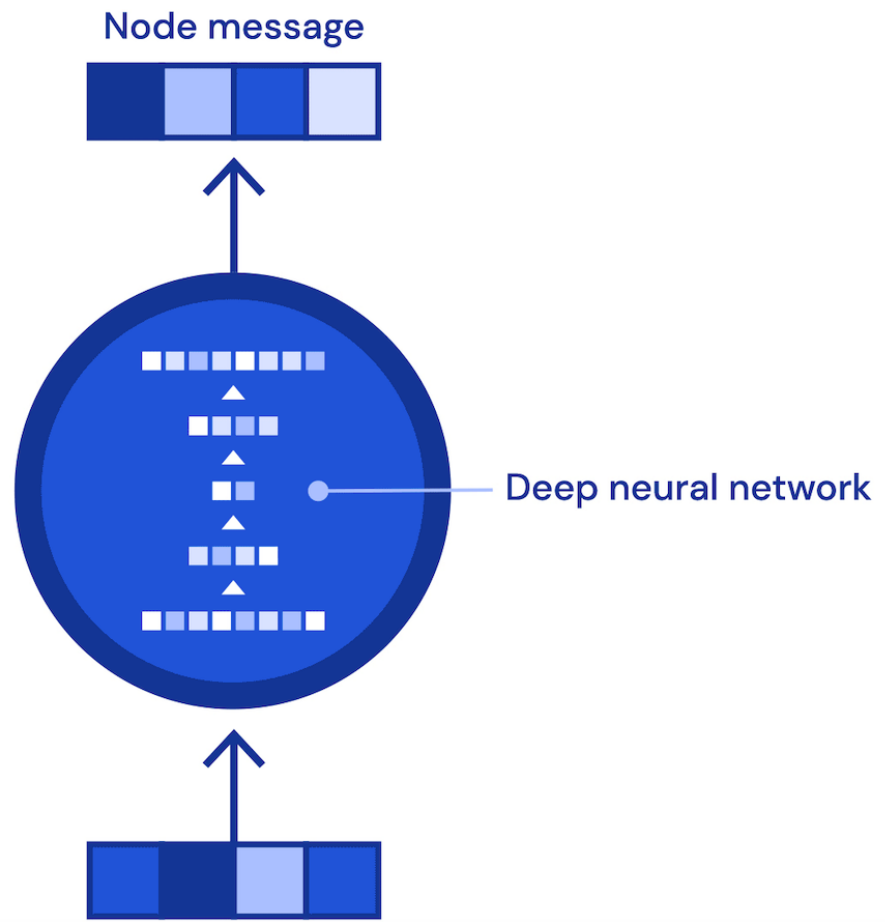


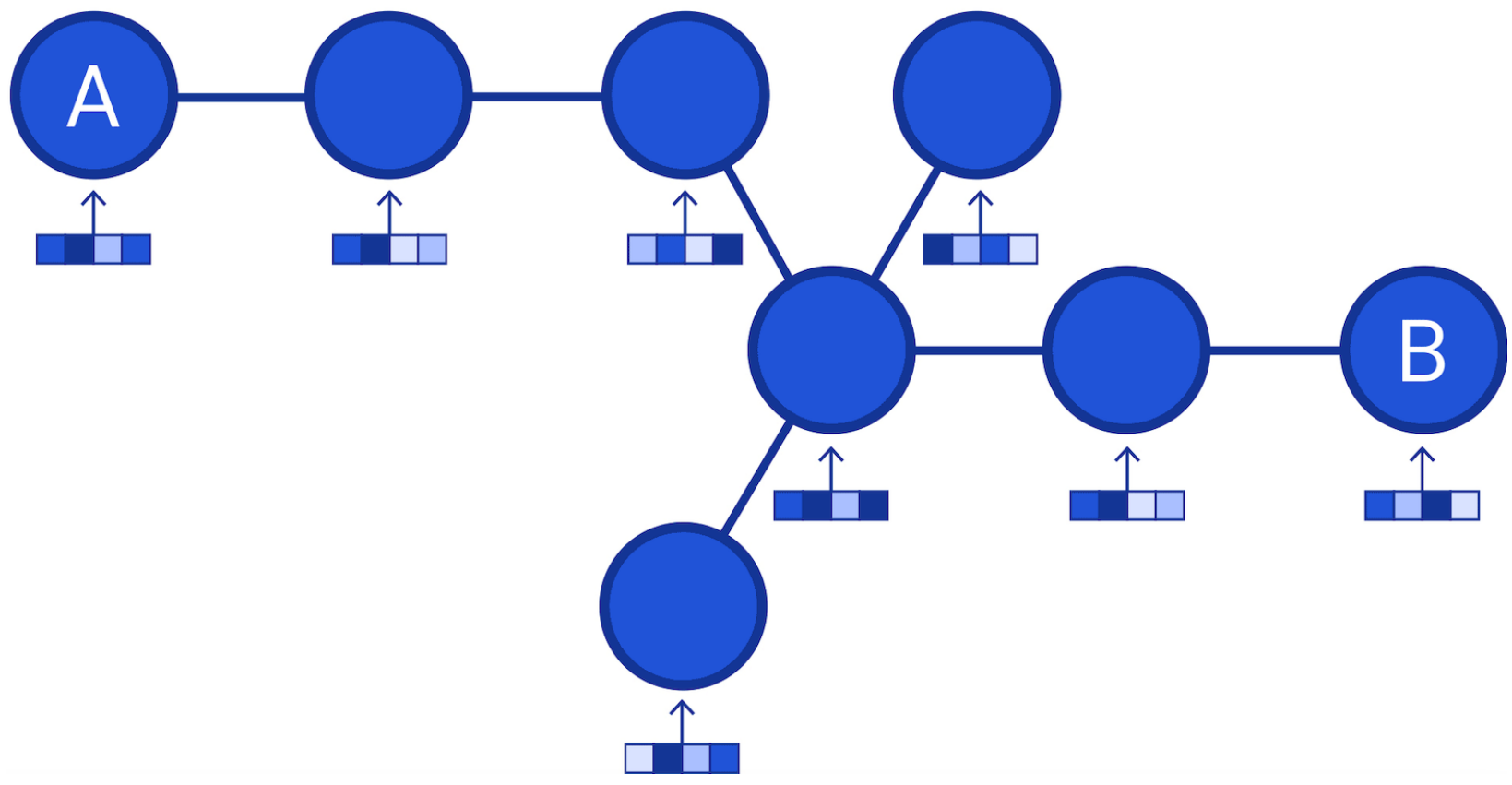


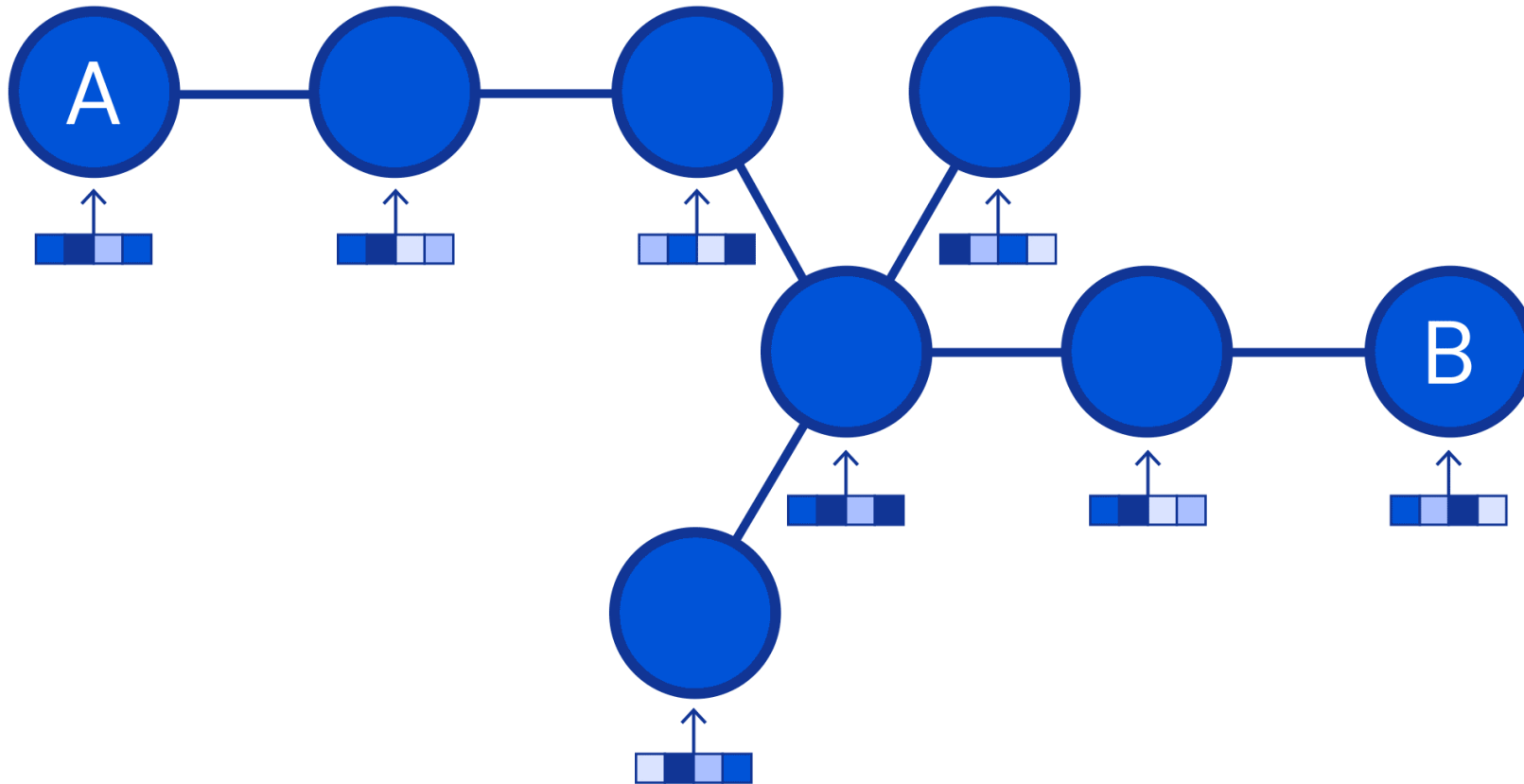












Encoding neighborhoods: General form

$$\mathbf{h}_u^{(0)} = \mathbf{x}_u \text{ (feature representation for node } u\text{)}$$

In each round $k \in [K]$, for each node v :

1. **Aggregate** over neighbors

$$\mathbf{m}_{N(v)}^{(k)} = \text{AGGREGATE}^{(k)} \left(\left\{ \mathbf{h}_u^{(k-1)} : u \in N(v) \right\} \right)$$

Neighborhood of v

Encoding neighborhoods: General form

$\mathbf{h}_u^{(0)} = \mathbf{x}_u$ (feature representation for node u)

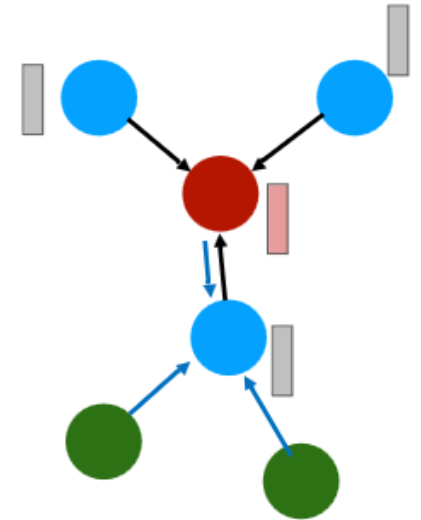
In each round $k \in [K]$, for each node v :

1. **Aggregate** over neighbors

$$\mathbf{m}_{N(v)}^{(k)} = \text{AGGREGATE}^{(k)} \left(\left\{ \mathbf{h}_u^{(k-1)} : u \in N(v) \right\} \right)$$

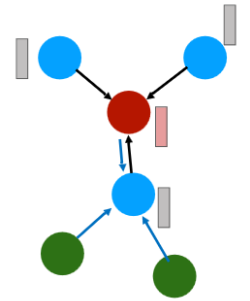
2. **Update** current node representation

$$\mathbf{h}_v^{(k)} = \text{COMBINE}^{(k)} \left(\mathbf{h}_v^{(k-1)}, \mathbf{m}_{N(v)}^{(k)} \right)$$



The basic GNN

[Merkwirth and Lengauer '05; Scarselli et al. '09]



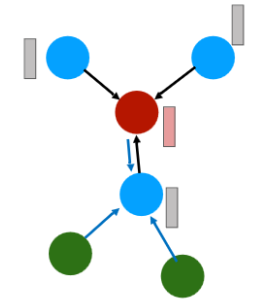
$$\mathbf{m}_{N(v)} = \text{AGGREGATE}(\{\mathbf{h}_u : u \in N(v)\}) = \sum_{u \in N(v)} \mathbf{h}_u$$

$$\text{COMBINE}(\mathbf{h}_v, \mathbf{m}_{N(v)}) = \sigma(W_{\text{self}}\mathbf{h}_v + W_{\text{neigh}}\mathbf{m}_{N(v)} + \mathbf{b})$$

Trainable parameters

Non-linearity (e.g.,
tanh or ReLU)

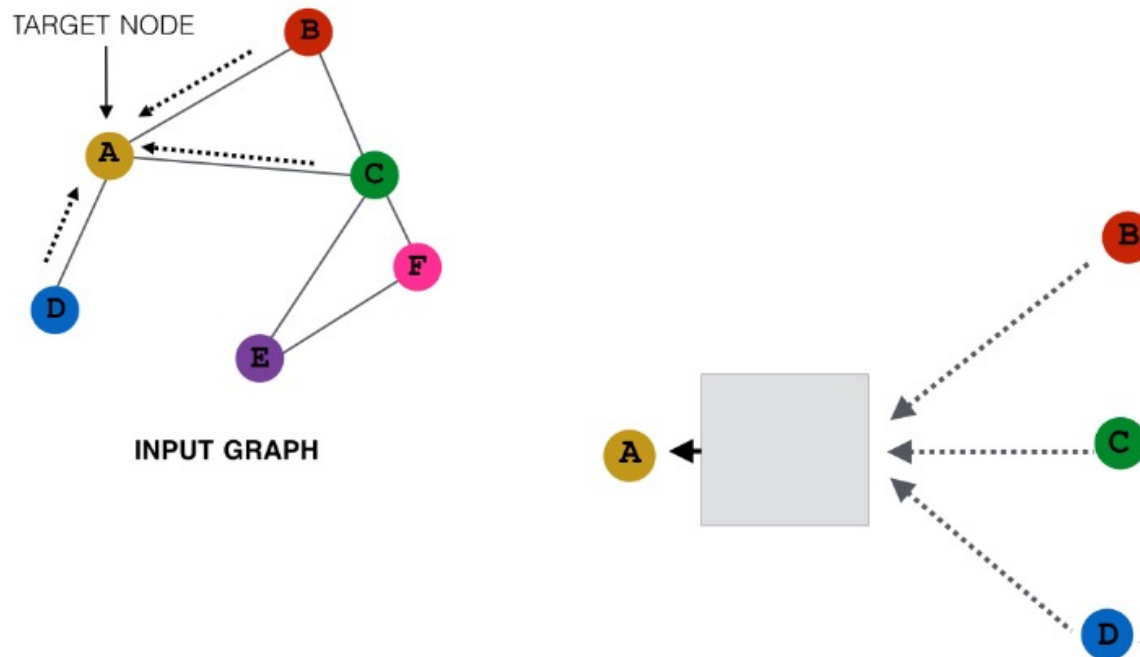
Aggregation functions



$$\mathbf{m}_{N(v)} = \text{AGGREGATE}(\{\mathbf{h}_u : u \in N(v)\}) = \bigoplus_{u \in N(v)} \mathbf{h}_u$$

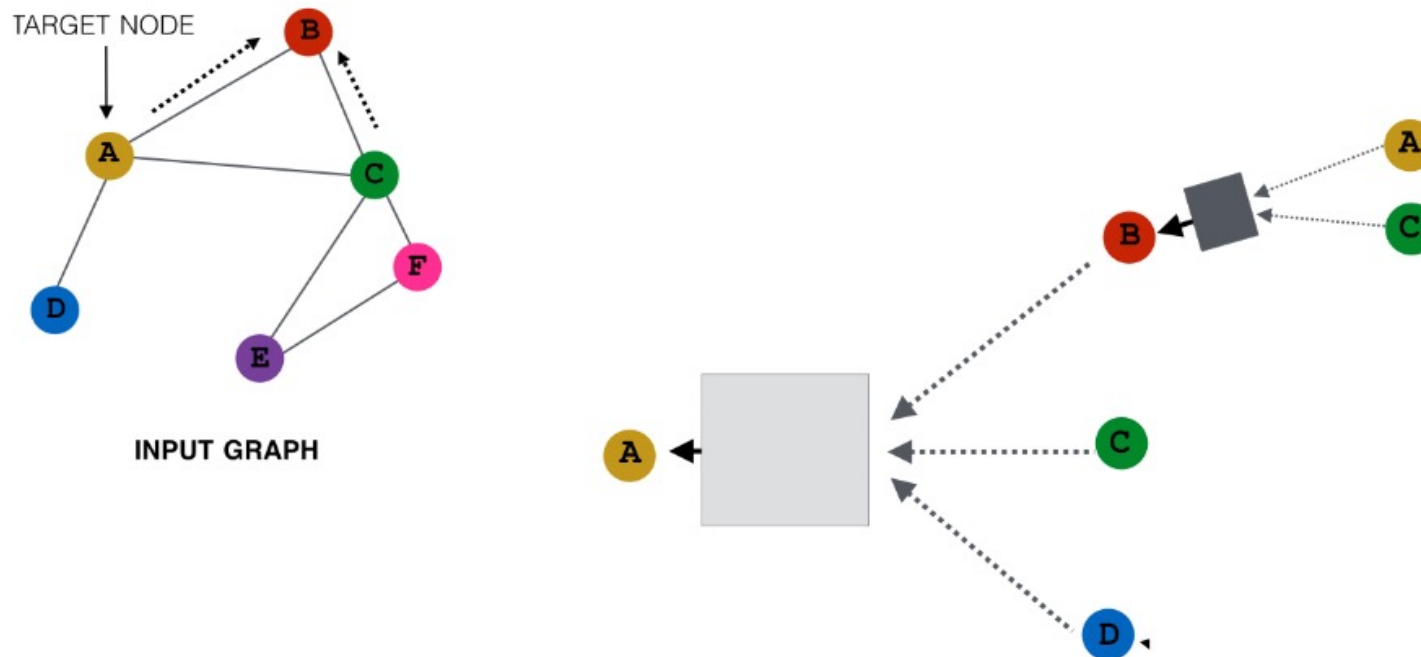
Other element-wise aggregators, e.g.:
Maximization, averaging

Node embeddings unrolled



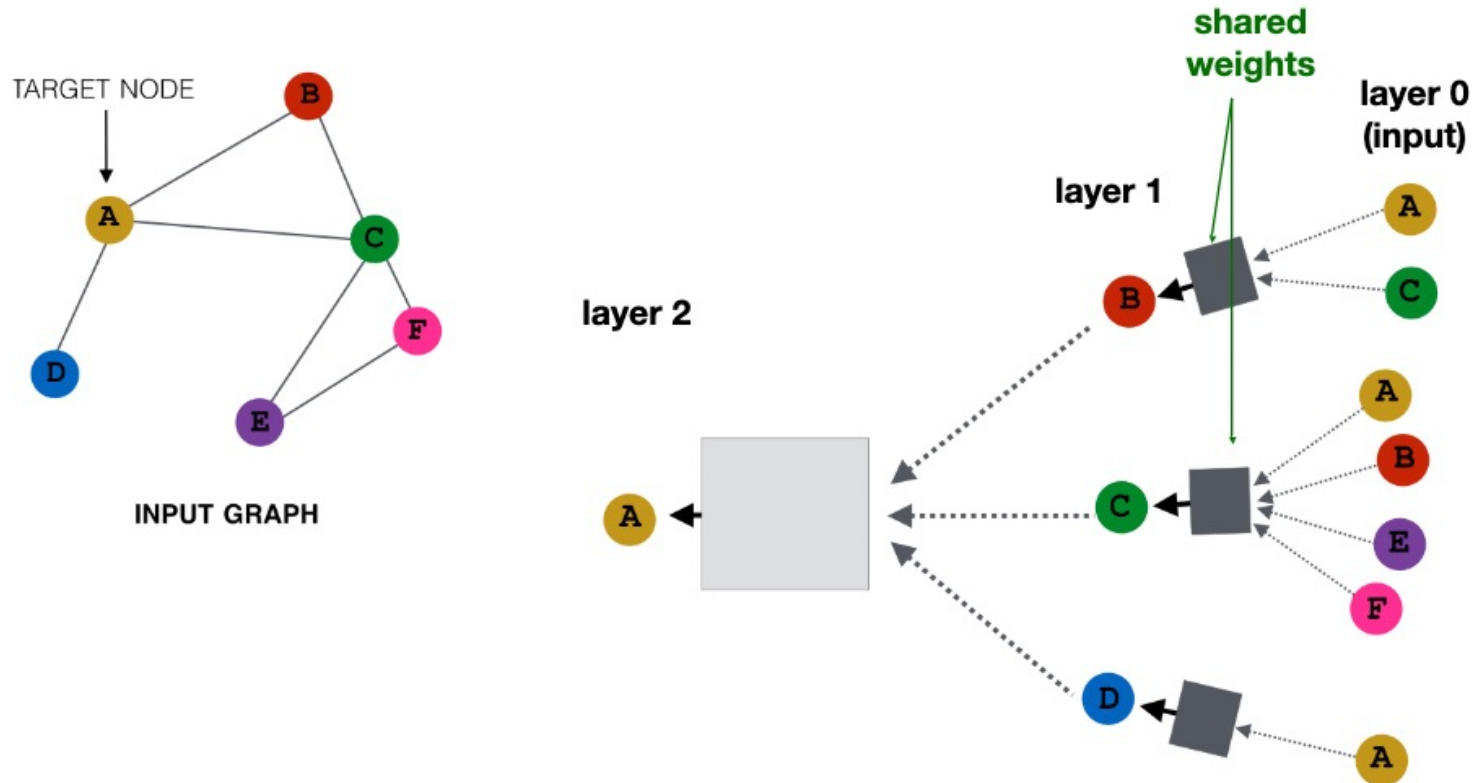
Grey boxes: aggregation functions that we learn

Node embeddings unrolled



Grey boxes: aggregation functions that we learn

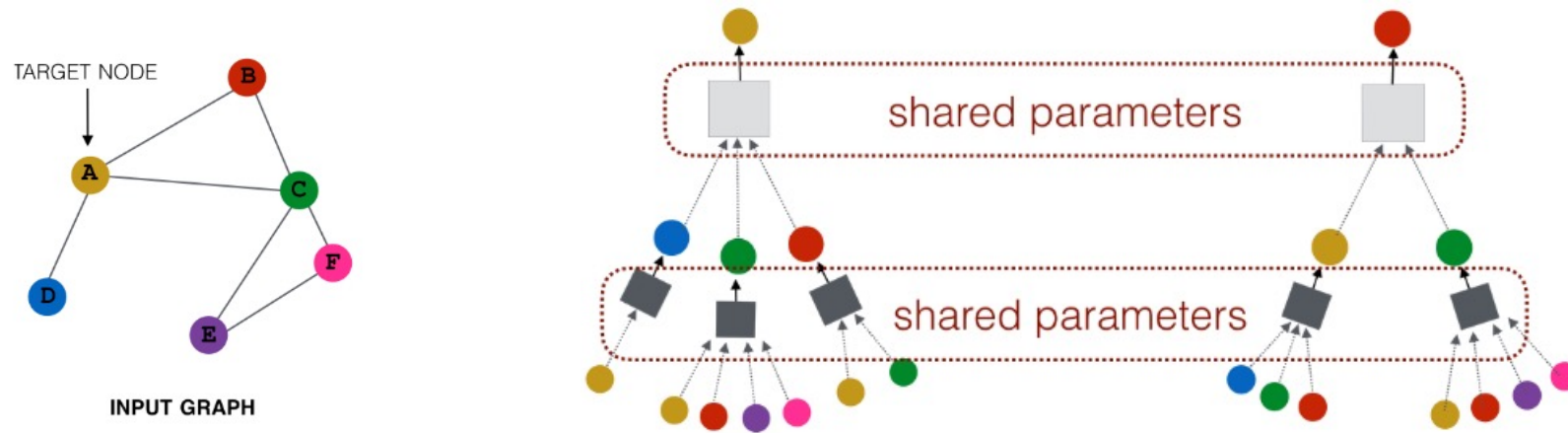
Node embeddings unrolled



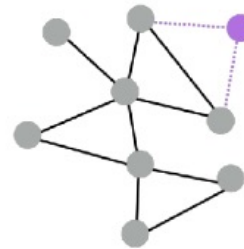
Grey boxes: aggregation functions that we learn

Weight sharing

Use the same aggregation functions for all nodes

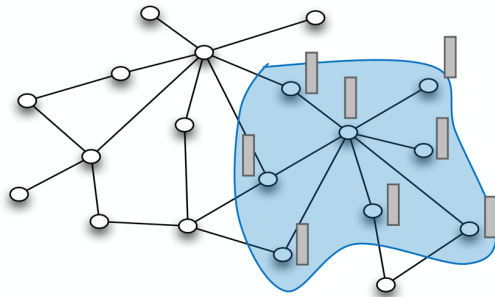


Can generate encodings for previously unseen nodes & graphs!

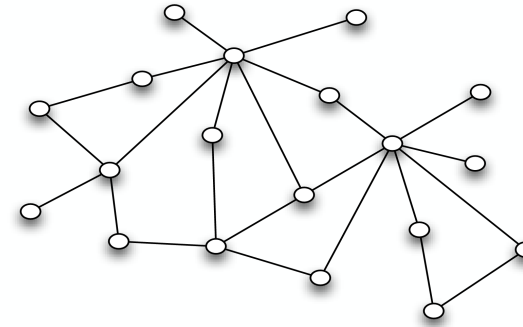


Training a GNN

- What is a data point?



Node and its neighborhood



Entire graph

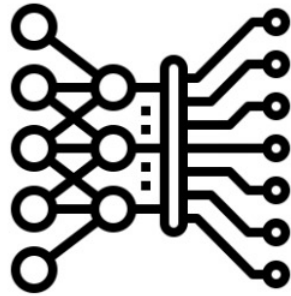
- What to specify?
 - Aggregate and combine functions
 - Readout function: combines node embeddings \rightarrow graph embedding
 - Loss function on prediction
- Train with SGD

Outline (applied techniques)

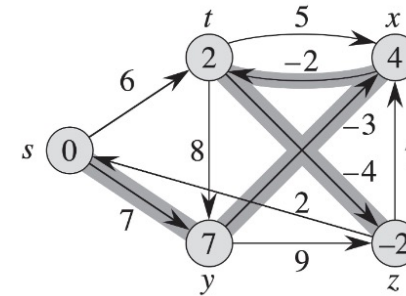
1. GNNs overview
- 2. Neural algorithmic alignment**
3. Integer programming with GNNs

Veličković, Ying, Padovano, Hadsell, Blundell, ICLR'20
Cappart, Chételat, Khalil, Lodi, Morris, Veličković, arXiv'21

Problem-solving approaches



- + Operate on raw inputs
- + Generalize on noisy conditions
- + Models reusable across tasks
- Require big data
- Unreliable when extrapolating
- Lack of interpretability



- + Trivially strong generalization
- + Compositional (subroutines)
- + Guaranteed correctness
- + Interpretable operations
- Input must match spec
- Not robust to task variations

Is it possible to get the best of both worlds?

Previous work

Previous work:

- Shortest path [Graves et al. '16; Xu et al., '19]
- Traveling salesman [Reed and De Freitas '15]
- Boolean satisfiability [Vinyals et al. '15; Bello et al., '16; ...]
- Probabilistic inference [Yoon et al., '18]

Ground-truth solutions used to drive learning

Model has **complete freedom** mapping raw inputs to solutions

Neural graph algorithm execution

Key observation: Many algorithms share related **subroutines**
E.g. Bellman-Ford, BFS enumerate sets of edges adjacent to a node

Neural graph algorithm execution

- Learn several algorithms **simultaneously**
- Provide intermediate supervision signals
Driven by how a known classical algorithm would process the input

Outline (applied techniques)

1. GNNs overview
2. Neural algorithmic alignment
 - i. Example algorithms**
 - ii. Experiments
 - iii. Additional motivation
 - iv. Additional research
3. Integer programming with GNNs

Breadth-first search

- Source node s

- Initial input $x_i^{(1)} = \begin{cases} 1 & \text{if } i = s \\ 0 & \text{if } i \neq s \end{cases}$

- Node is reachable from s if any of its neighbors are reachable:

$$x_i^{(t+1)} = \begin{cases} 1 & \text{if } x_i^{(t)} = 1 \\ 1 & \text{if } \exists j \text{ s.t. } (j, i) \in E \text{ and } x_j^{(t)} = 1 \\ 0 & \text{else} \end{cases}$$

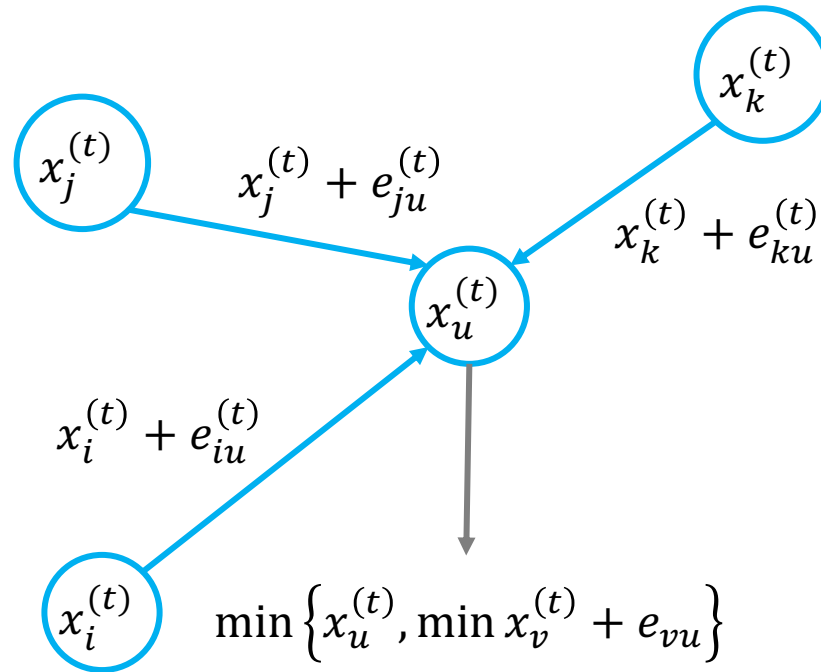
- Algorithm output at round t : $y_i^{(t)} = x_i^{(t+1)}$

Bellman-Ford (shortest path)

- Source node s
- Initial input $x_i^{(1)} = \begin{cases} 0 & \text{if } i = s \\ \infty & \text{if } i \neq s \end{cases}$
- Node is reachable from s if any of its neighbors are reachable
Update distance to node as minimal way to reach neighbors

$$x_i^{(t+1)} = \min \left\{ x_i^{(t)}, \min_{(j,i) \in E} x_j^{(t)} + e_{ji}^{(t)} \right\}$$

Bellman-Ford: Message passing

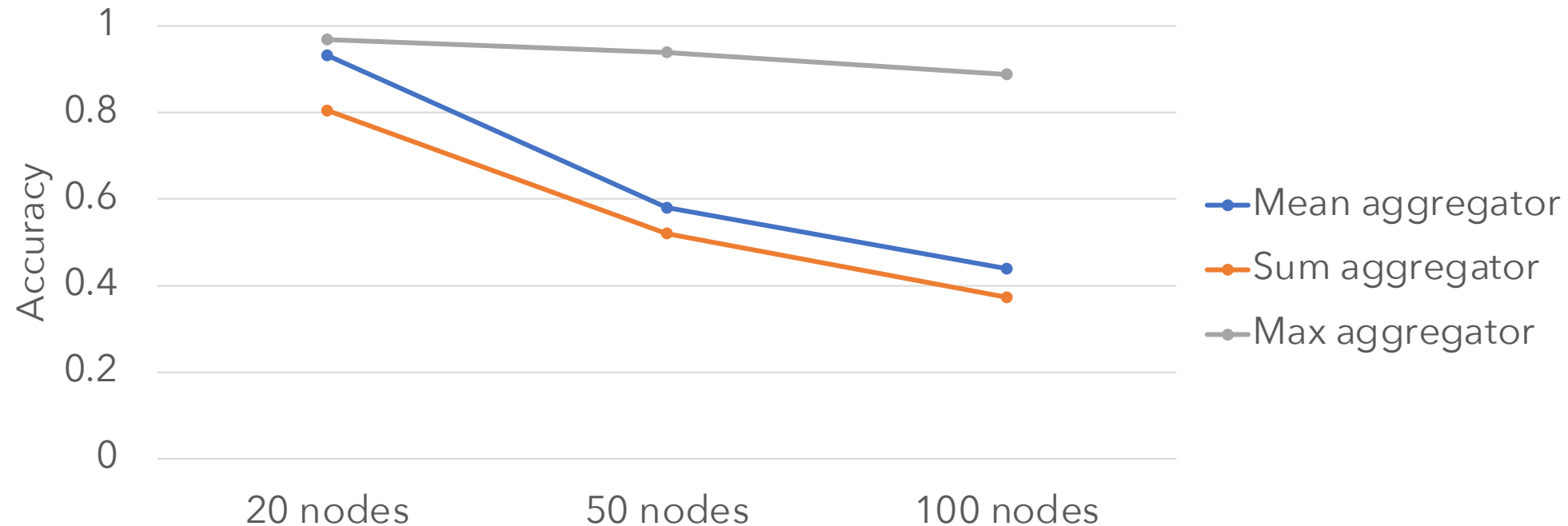


Key idea (roughly speaking): Train GNN so that $\mathbf{h}_u^{(t)} \approx x_u^{(t)}, \forall t$
(Really, so that a function of $\mathbf{h}_u^{(t)} \approx x_u^{(t)}$)

Outline (applied techniques)

1. GNNs overview
2. Neural algorithmic alignment
 - i. Example algorithms
 - ii. Experiments**
 - iii. Additional motivation
 - iv. Additional research
3. Integer programming with GNNs

Shortest-path predecessor prediction



Improvement of max-aggregator increases with size

It **aligns** better with underlying algorithm [Xu et al., ICLR'20]

Learning multiple algorithms

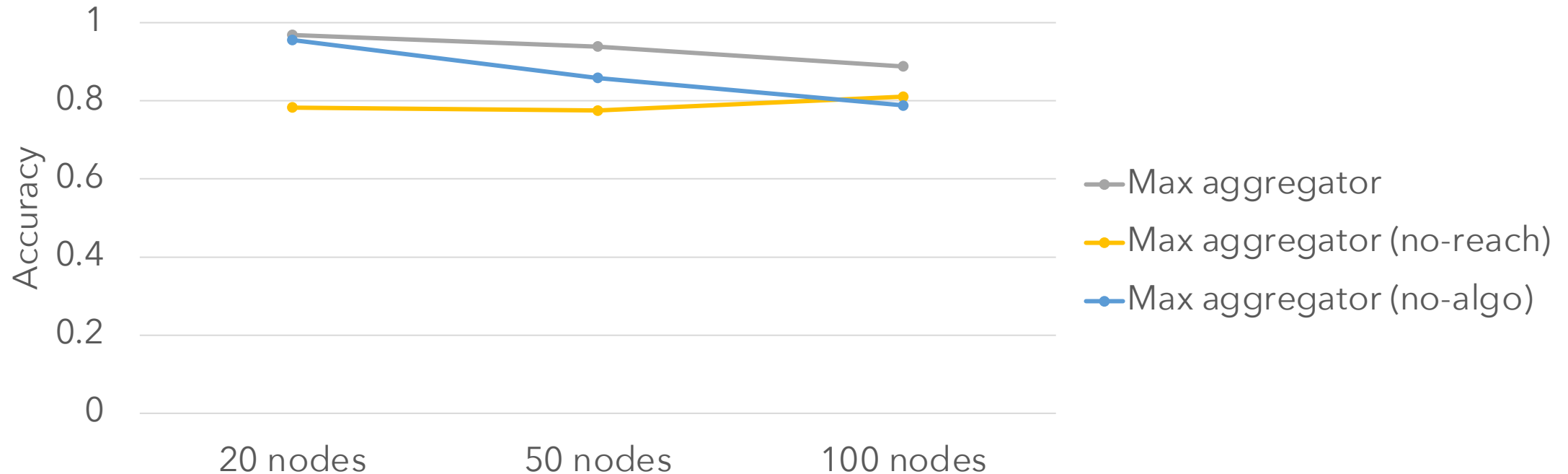
Learn to execute both BFS and Bellman-Ford **simultaneously**

- At each step t , concatenate relevant $x_i^{(t)}$ and $y_i^{(t)}$ values

Comparisons

- (*no-reach*): Learn Bellman-Ford alone
 - Doesn't simultaneously learn reachability
- (*no-algo*):
 - Don't supervise intermediate steps
 - Learn predecessors directly from input $x_i^{(1)}$

Shortest-path predecessor prediction



- **(no-reach) results:** positive knowledge transfer
- **(no-algo) results:** benefit of supervising intermediate steps

Outline (applied techniques)

1. GNNs overview
2. Neural algorithmic alignment
 - i. Example algorithms
 - ii. Experiments
 - iii. Additional motivation**
 - iv. Additional research
3. Integer programming with GNNs

Key question

Key question in neural algorithmic alignment:

If we're just teaching a NN to **imitate** a classical algorithm...

Why not just run that algorithm?

Why use GNNs for algorithm design?

Classical algorithms are designed with **abstraction** in mind
Enforce their inputs to conform to stringent preconditions

However, we design algorithms to solve **real-world** problems!



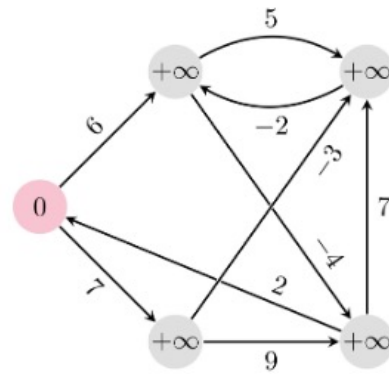
Natural inputs

Abstractifying the core problem

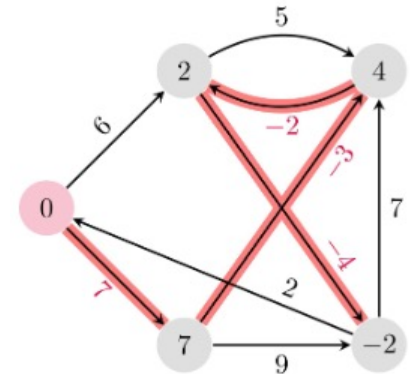
- Assume we have real-world inputs
 - ...but algorithm only admits abstract inputs
- Could try **manually** converting from one input to another



Natural inputs



Abstract inputs



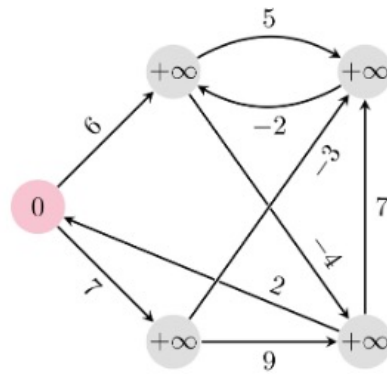
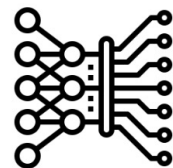
Abstract outputs

Attacking the core problem

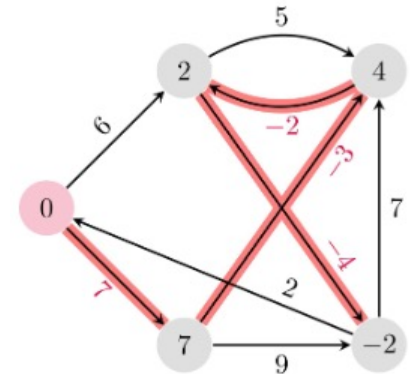
- Alternatively, **replace** human feature extractor with NN
 - Still apply same combinatorial algorithm
- Issue: algorithms typically perform **discrete optimization**
 - Doesn't play nicely with **gradient-based** optimization of NNs



Natural inputs



Abstract inputs



Abstract outputs

Algorithmic bottleneck

Second (more fundamental) issue: **data efficiency**

- Real-world data is often incredibly rich
- We still have to compress it down to scalar values

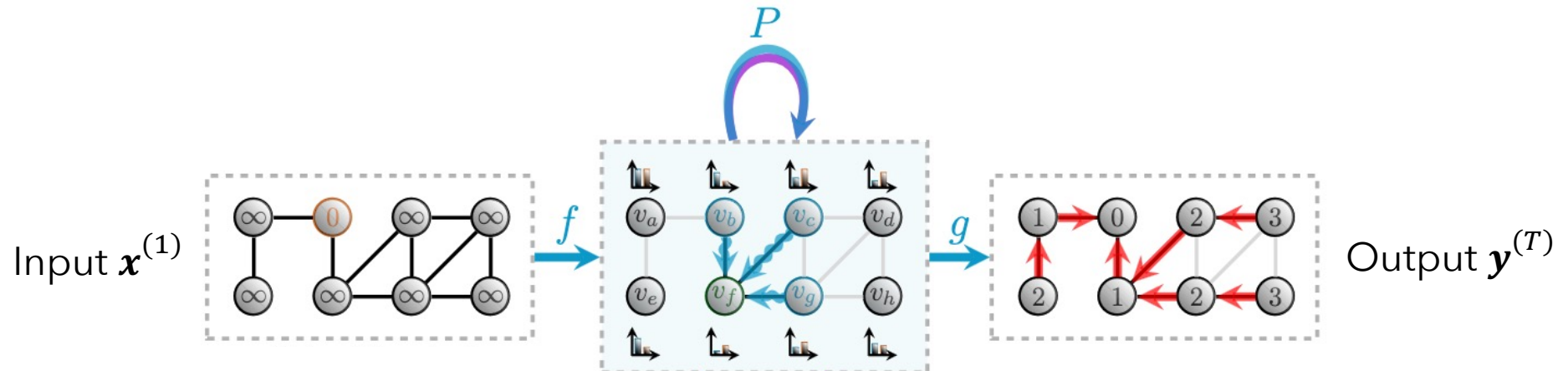
The algorithmic solver commits to using this scalar

Assumes it is perfect!

If there's insufficient training data to estimate the scalars:

- Alg will give a **perfect solution**
- ...but in a **suboptimal environment**

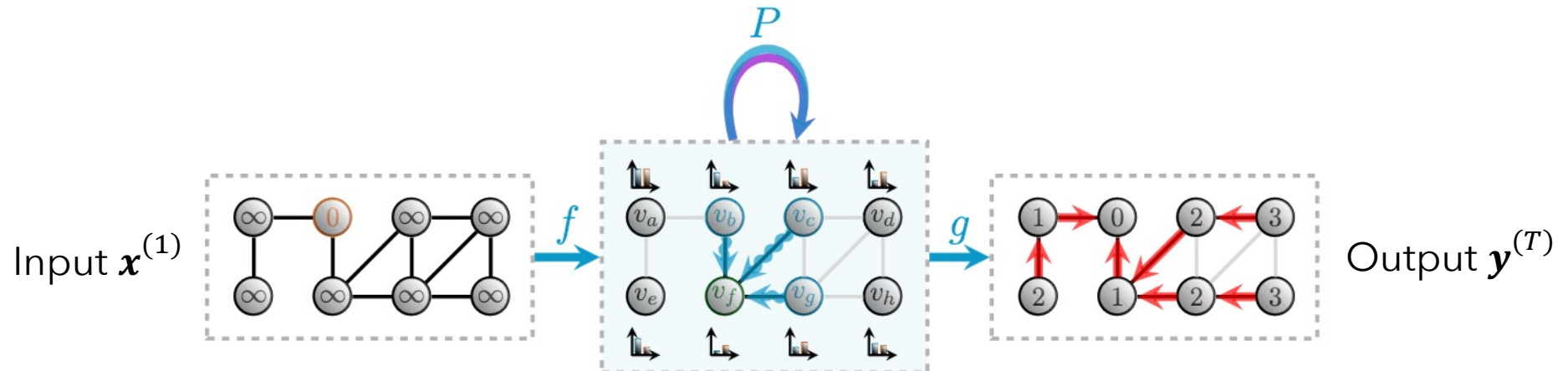
Neural algorithmic pipeline



Encoder network f

- E.g., makes sure input is in correct dimension for next step

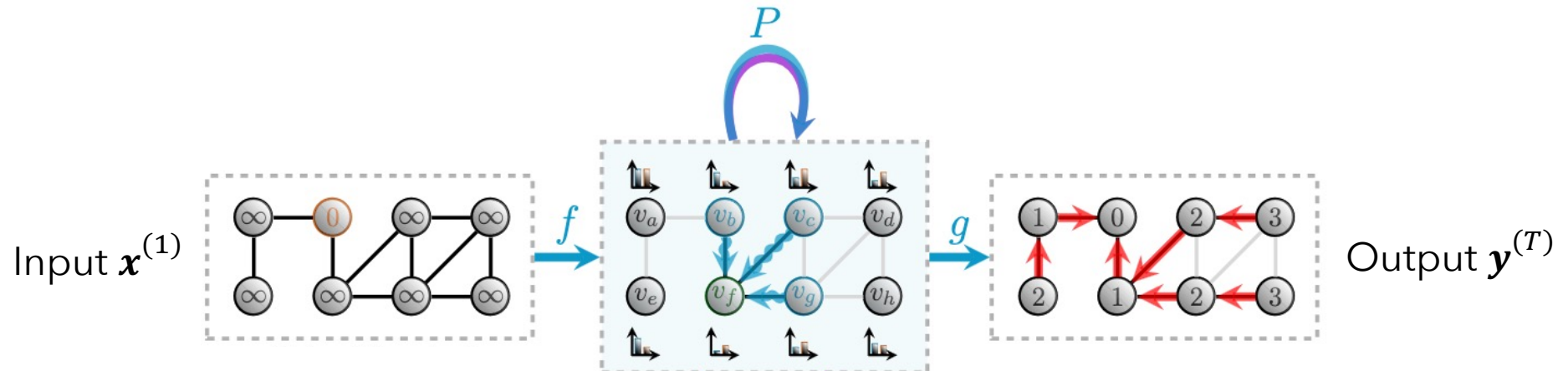
Neural algorithmic pipeline



Processor network P

- Graph neural network
- Run multiple times (termination determined by a NN)

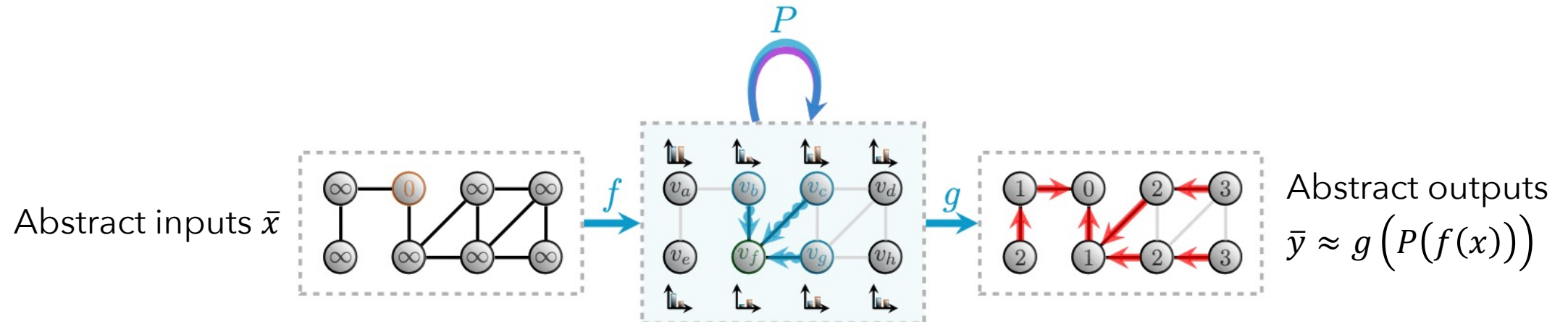
Neural algorithmic pipeline



Decoder network g

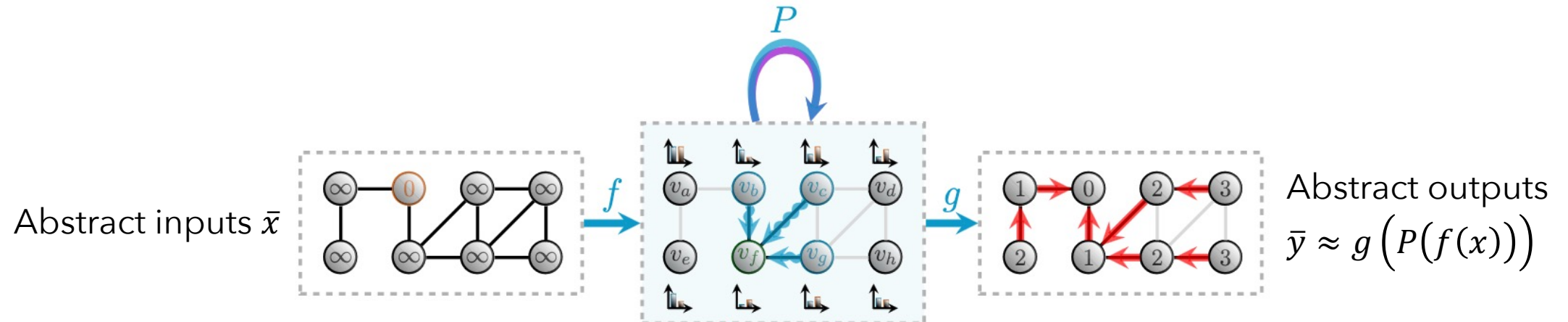
- Transform's GNNs output into algorithmic output

Neural algorithmic pipeline



1. On abstract inputs, learn encode-process-decode functions

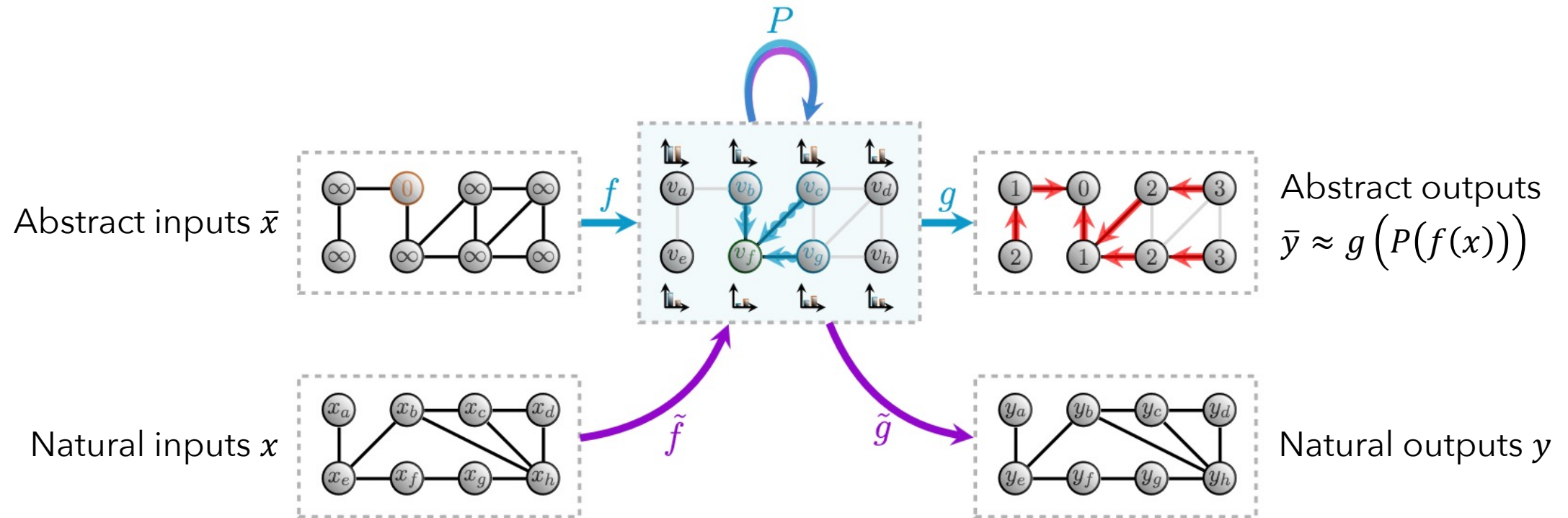
Neural algorithmic pipeline



After training on abstract inputs, processor P :

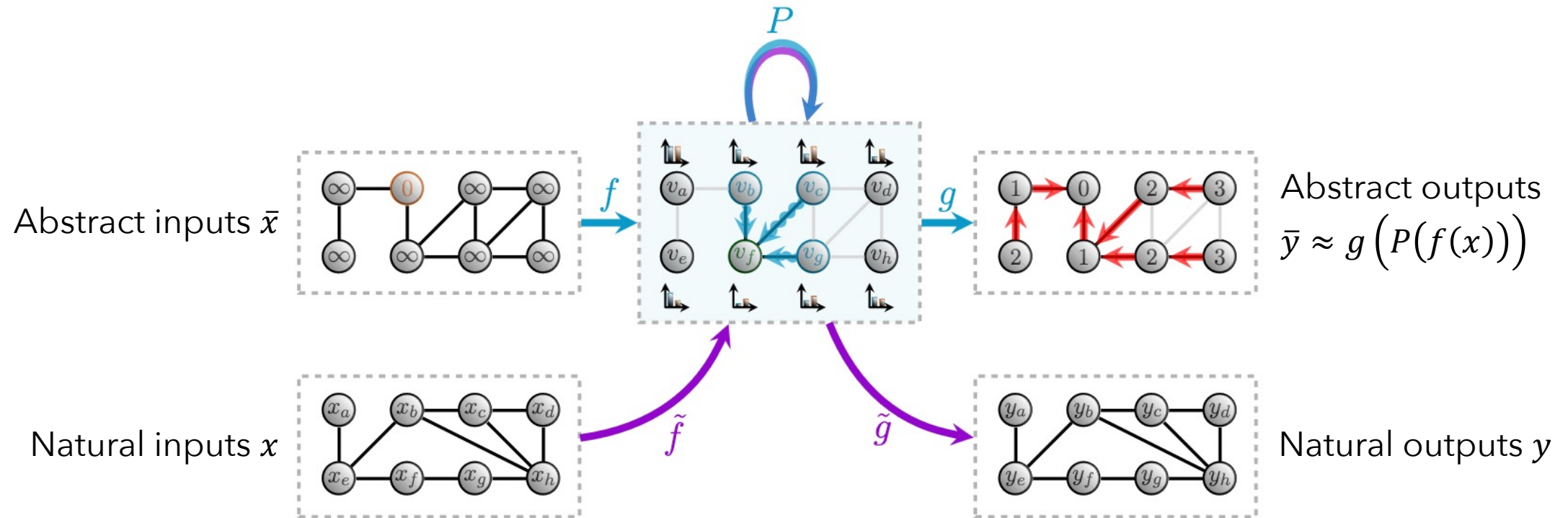
1. Is aligned with computations of target algorithm
2. Admits useful gradients
3. Operates over high-dim latent space (better use of data)

Neural algorithmic pipeline



2. Set up encode-decode functions for natural inputs/outputs

Neural algorithmic pipeline



3. Learn parameters using loss that compares $\tilde{g}\left(P\left(\tilde{f}(x)\right)\right)$ to y

Outline (applied techniques)

1. GNNs overview
2. Neural algorithmic alignment
 - i. Example algorithms
 - ii. Experiments
 - iii. Additional motivation
 - iv. Additional research**
3. Integer programming with GNNs

Additional research

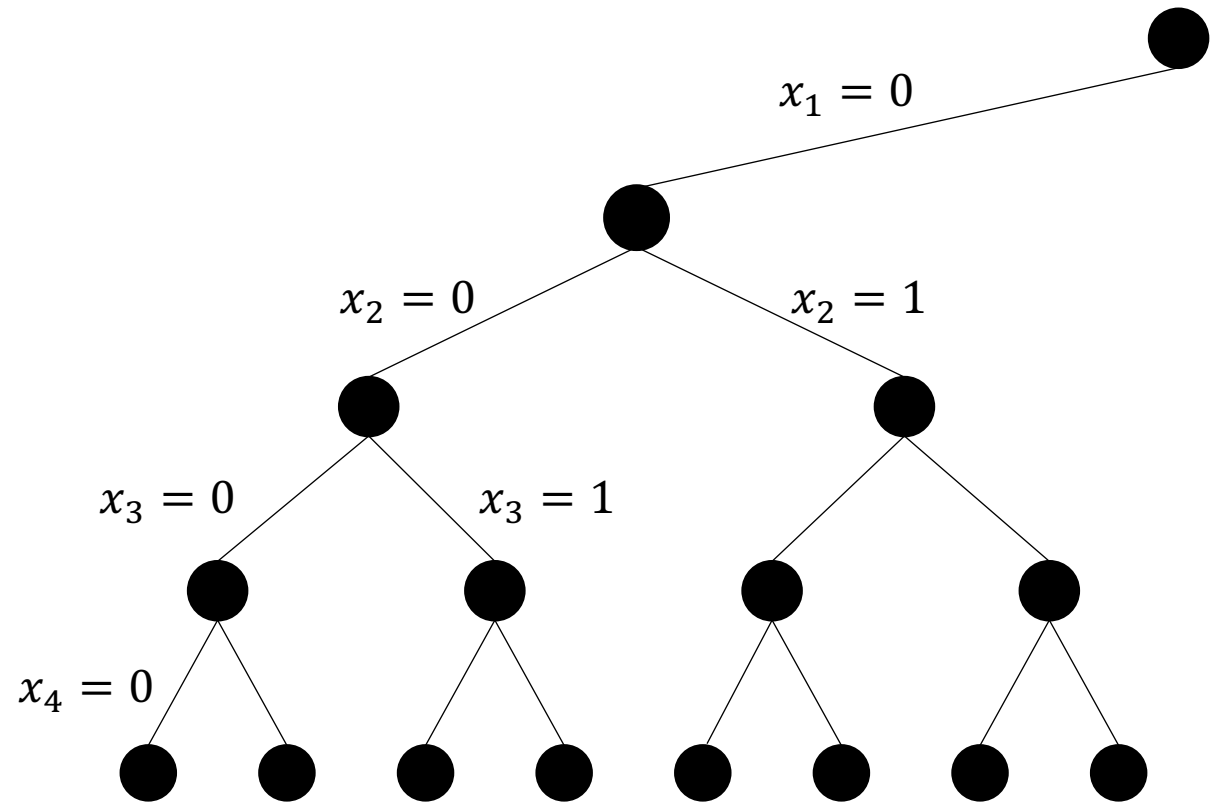
Lots of research in the past few years! E.g.:

- How to achieve **algorithmic alignment** & **theory guarantees**
 - Xu et al., ICLR'20; Dudzik, Veličković, NeurIPS'22
- **CLRS** benchmark
 - Sorting, searching, dynamic programming, graph algorithms, etc.
 - Veličković et al. ICML'22; Ibarz et al. LoG'22; Bevilacqua et al. ICML'23
- **Primal-dual** algorithms
 - Numeroso et al., ICLR'23

Outline (applied techniques)

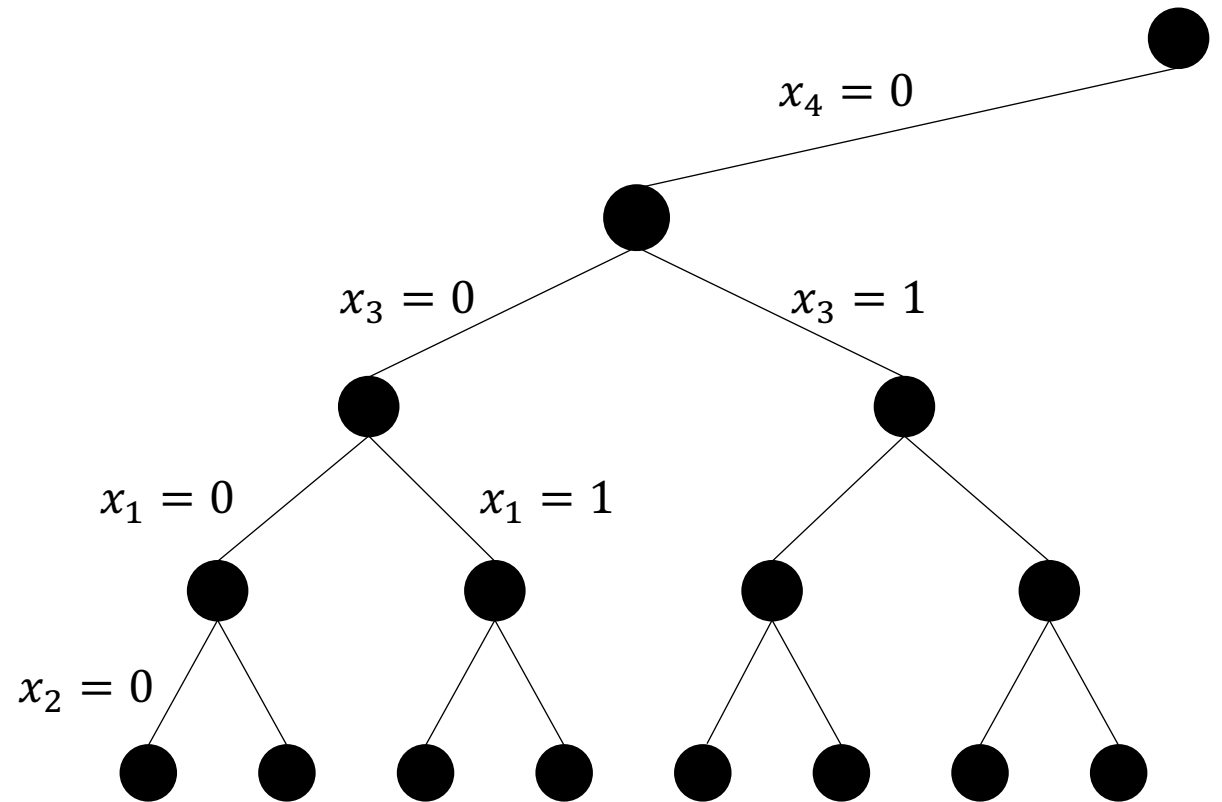
1. GNNs overview
2. Neural algorithmic alignment
- 3. Integer programming with GNNs**

Variable selection policy (VSP)



Better branching order than x_1, x_2, x_3, x_4 ?

Variable selection policy (VSP)



Better branching order than x_1, x_2, x_3, x_4 ? E.g., x_4, x_3, x_1, x_2

Variable selection policy (VSP)

At node j with LP objective value $z(j)$:

- Let $z_i^+(j)$ be the LP objective value after setting $x_i = 1$
- Let $z_i^-(j)$ be the LP objective value after setting $x_i = 0$

VSP example:

Branch on the variable x_i that maximizes

$$\max\{z(j) - z_i^+(j), 10^{-6}\} \cdot \max\{z(j) - z_i^-(j), 10^{-6}\}$$

If score was $(z(j) - z_i^+(j))(z(j) - z_i^-(j))$ and $z(j) - z_i^+(j) = 0$:
would lose information stored in $z(j) - z_i^-(j)$

Strong branching

Challenge: Computing $z_i^-(j), z_i^+(j)$ requires solving a lot of LPs

- Computing all LP relaxations referred to as **strong-branching**
- Very **time intensive**

Pro: Strong branching leads to small search trees

Idea: Train an ML model to imitate strong-branching

Khalil et al. [AAAI'16], Alvarez et al. [INFORMS JoC'17], Hansknecht et al. [arXiv'18]

This paper: using a GNN

Outline (applied techniques)

1. GNNs overview
2. Neural algorithmic alignment
3. Integer programming with GNNs
 - i. Machine learning formulation**
 - ii. Baselines
 - iii. Experiments
 - iv. Additional research

Problem formulation

Goal: learn a policy $\pi(a_t | s_t)$

Probability of branching on variable a_t when solver is in state s_t

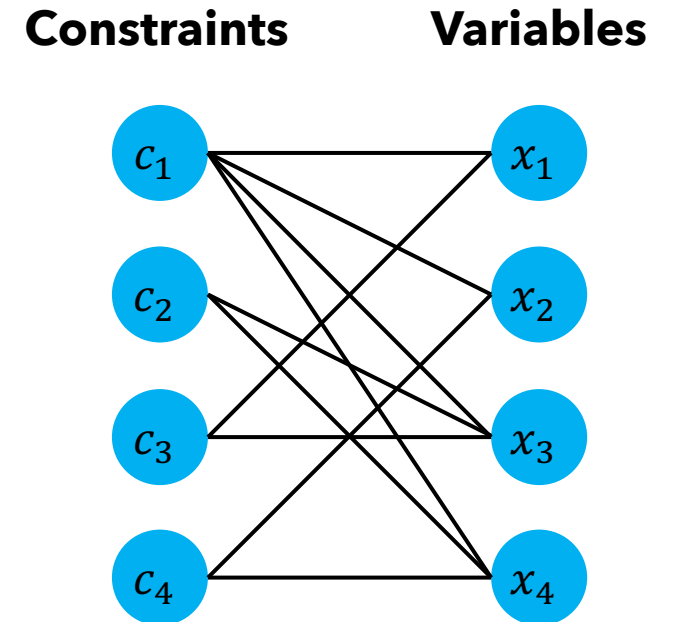
Approach (imitation learning):

- Run strong branching on training set of instances
- Collect dataset of (state, variable) pairs $S = \{(s_i, a_i^*)\}_{i=1}^N$
- Learn policy π_θ with training set S

State encoding

State s_t of B&B encoded as a **bipartite graph**
with **node** and **edge features**

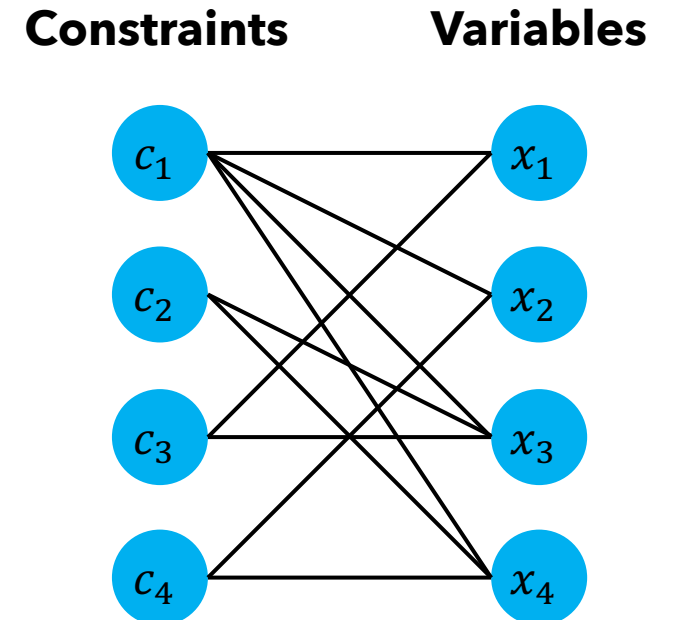
$$\begin{aligned} \max \quad & 9x_1 + 5x_2 + 6x_3 + 4x_4 \\ \text{s.t.} \quad & 6x_1 + 3x_2 + 5x_3 + 2x_4 \leq 10 \quad (c_1) \\ & x_3 + x_4 \leq 10 \quad (c_2) \\ & -x_1 + x_3 \leq 0 \quad (c_3) \\ & -x_2 + x_4 \leq 0 \quad (c_4) \\ & x_1, x_2, x_3, x_4 \in \{0,1\} \end{aligned}$$



State encoding

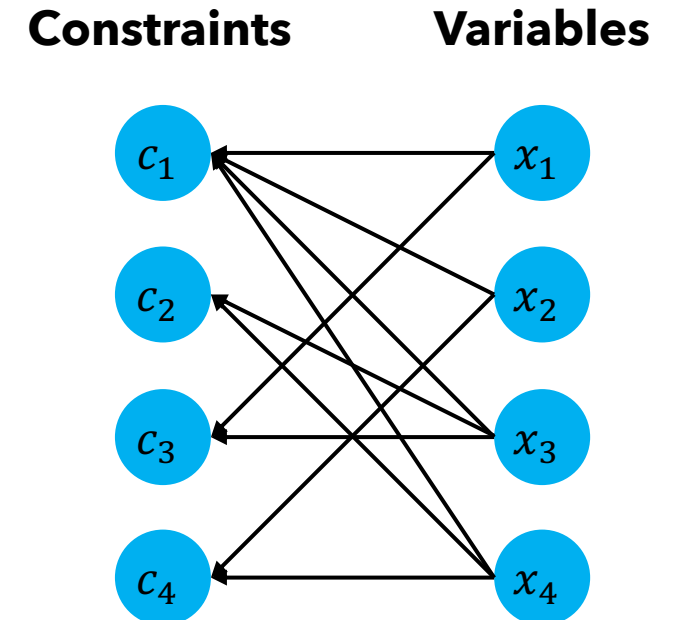
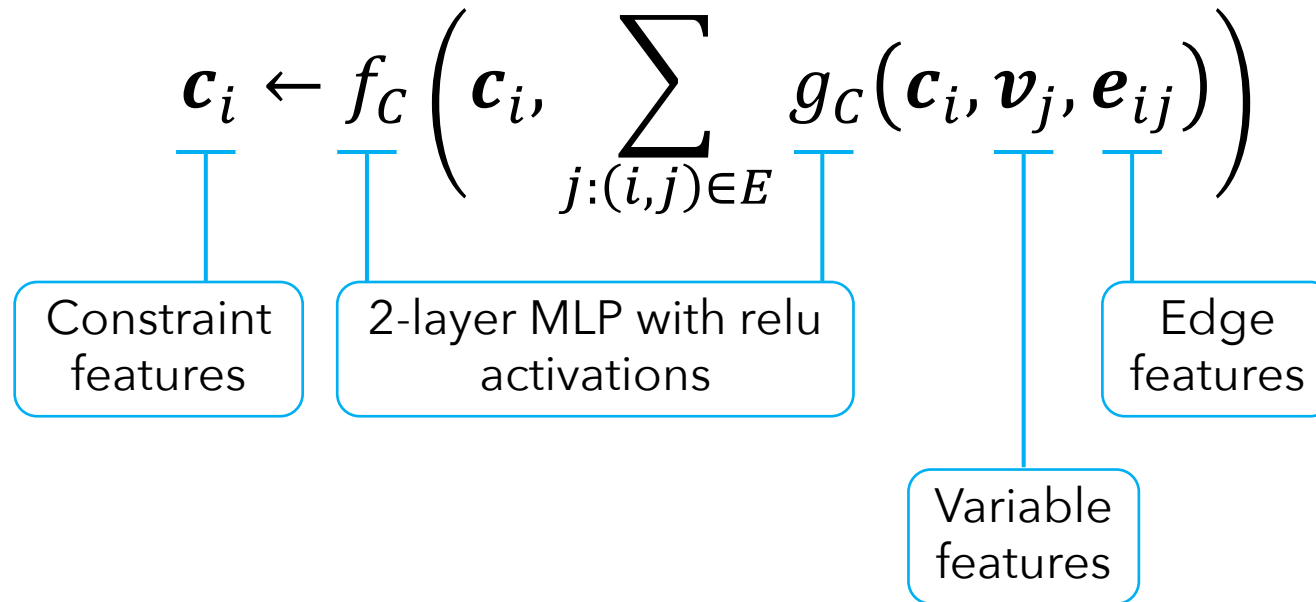
State s_t of B&B encoded as a **bipartite graph** with **node** and **edge features**

- **Edge feature:** constraint coefficient
- **Example node features:**
 - Constraints:
 - Cosine similarity with objective
 - Tight in LP solution?
 - Variables:
 - Objective coefficient
 - Solution value equals upper/lower bound?



GNN structure

1. Pass from variables \rightarrow constraints



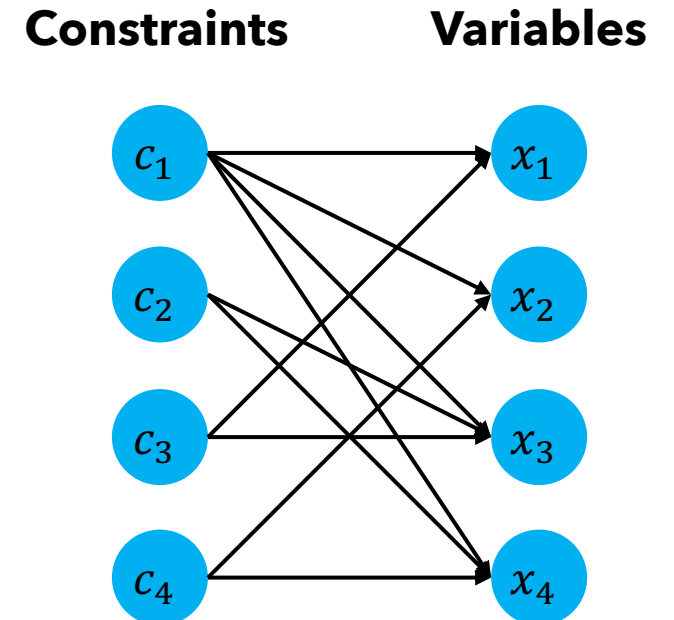
GNN structure

1. Pass from variables \rightarrow constraints

$$\mathbf{c}_i \leftarrow f_C \left(\mathbf{c}_i, \sum_{j:(i,j) \in E} g_C(\mathbf{c}_i, \mathbf{v}_j, \mathbf{e}_{ij}) \right)$$

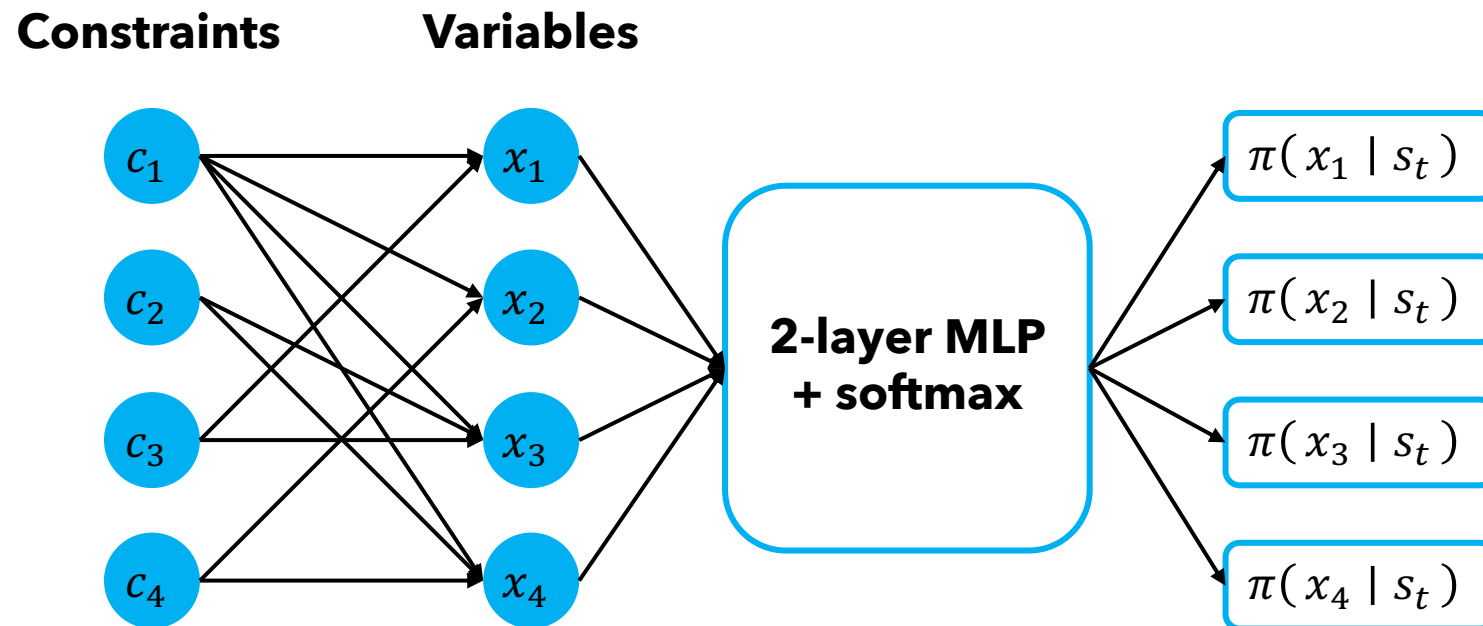
2. Pass from constraints \rightarrow variables

$$\mathbf{v}_j \leftarrow f_V \left(\mathbf{v}_j, \sum_{i:(i,j) \in E} g_V(\mathbf{c}_i, \mathbf{v}_j, \mathbf{e}_{ij}) \right)$$



GNN structure

3. Compute distribution over variables



Outline (applied techniques)

1. GNNs overview
2. Neural algorithmic alignment
3. Integer programming with GNNs
 - i. Machine learning formulation
 - ii. Baselines**
 - iii. Experiments
 - iv. Additional research

Reliability pseudo-cost branching (RPB)

Rough idea:

- Goal: estimate $z(j) - z_i^+(j)$ w/o solving the LP with $x_i = 1$
- Estimate = avg change after setting $x_i = 1$ elsewhere in tree
This is the "pseudo-cost"
- "Reliability": do strong branching if estimate is "unreliable"
E.g., early in the tree

Default branching rule of SCIP (leading open-source solver):

$$\max\{\tilde{\Delta}_i^+(j), 10^{-6}\} \cdot \max\{\tilde{\Delta}_i^-(j), 10^{-6}\}$$

Estimate of $z(j) - z_i^+(j)$

Estimate of $z(j) - z_i^-(j)$

Learning to rank approaches

- Predict which variable **strong branching** would rank highest
- Using a **linear model** instead of a GNN

- Khalil et al. [AAAI'16]:
 - Use learning-to-rank algorithm **SVM^{rank}** [Joachims, KDD'06]
- Hansknecht et al. [arXiv'18]
 - Use learning-to-rank alg **lambdaMART** [Burges, Learning'10]

Outline (applied techniques)

1. GNNs overview
2. Neural algorithmic alignment
3. Integer programming with GNNs
 - i. Machine learning formulation
 - ii. Baselines
 - iii. Experiments**
 - iv. Additional research

Set covering instances

Always train on “easy” instances

Model	1000 columns, 500 rows				1000 columns, 2000 rows				
	Time	Easy Wins	Nodes		Time	Hard Wins	Nodes		
FSB	17.30 ± 6.1%	0 / 100	17 ± 13.7%		3600.00 ± 0.0%	0 / 0	n/a	± n/a	%
RPB	8.98 ± 4.8%	0 / 100	54 ± 20.8%		1677.02 ± 3.0%	4 / 65	47 299	± 4.9%	
TREES	9.28 ± 4.9%	0 / 100	187 ± 9.4%		2869.21 ± 3.2%	0 / 35	59 013	± 9.3%	
SVMRANK	8.10 ± 3.8%	1 / 100	165 ± 8.2%		2389.92 ± 2.3%	0 / 47	42 120	± 5.4%	
LMART	7.19 ± 4.2%	14 / 100	167 ± 9.0%		2165.96 ± 2.0%	0 / 54	45 319	± 3.4%	
GCNN	6.59 ± 3.1%	85 / 100	134 ± 7.6%		1489.91 ± 3.3%	66 / 70	29 981	± 4.9%	

Set covering instances

Model	Time	Easy Wins	Nodes	Time	Hard Wins	Nodes
FSB	17.30 ± 6.1%	0 / 100	17 ± 13.7%	3600.00 ± 0.0%	0 / 0	n/a ± n/a %
RPB	8.98 ± 4.8%	0 / 100	54 ± 20.8%	1677.02 ± 3.0%	4 / 65	47 299 ± 4.9%
TREES	9.28 ± 4.9%	0 / 100	187 ± 9.4%	2869.21 ± 3.2%	0 / 35	59 013 ± 9.3%
SVMRANK	8.10 ± 3.8%	1 / 100	165 ± 8.2%	2389.92 ± 2.3%	0 / 47	42 120 ± 5.4%
LMART	7.19 ± 4.2%	14 / 100	167 ± 9.0%	2165.96 ± 2.0%	0 / 54	45 319 ± 3.4%
GCNN	6.59 ± 3.1%	85 / 100	134 ± 7.6%	1489.91 ± 3.3%	66 / 70	29 981 ± 4.9%

Set covering instances

- GNN is **faster than SCIP** default VSP (RPB)
- Performance generalizes to **larger instances**
- Similar results for auction design & facility location problems

Model	Time	Easy		Hard					
		Wins	Nodes	Time	Wins	Nodes	Time	Wins	Nodes
FSB	17.30 ± 6.1%	0 / 100	17 ± 13.7%	3600.00 ± 0.0%	0 / 0	n/a ± n/a %			
RPB	8.98 ± 4.8%	0 / 100	54 ± 20.8%	1677.02 ± 3.0%	4 / 65	47 299 ± 4.9%			
TREES	9.28 ± 4.9%	0 / 100	187 ± 9.4%	2869.21 ± 3.2%	0 / 35	59 013 ± 9.3%			
SVMRANK	8.10 ± 3.8%	1 / 100	165 ± 8.2%	2389.92 ± 2.3%	0 / 47	42 120 ± 5.4%			
LMART	7.19 ± 4.2%	14 / 100	167 ± 9.0%	2165.96 ± 2.0%	0 / 54	45 319 ± 3.4%			
GCNN	6.59 ± 3.1%	85 / 100	134 ± 7.6%	1489.91 ± 3.3%	66 / 70	29 981 ± 4.9%			

Max independent set instances

RPB is catching up to GNN on MIS instances

Model	Time	Easy		Hard			
		Wins	Nodes	Time	Wins	Nodes	
FSB	23.58 \pm 29.9%	9 / 100	7 \pm 35.9%	3600.00 \pm 0.0%	0 / 0	n/a \pm n/a %	
RPB	8.77 \pm 11.8%	7 / 100	20 \pm 36.1%	2045.61 \pm 18.3%	22 / 42	2675 \pm 24.0%	
TREES	10.75 \pm 22.1%	1 / 100	76 \pm 44.2%	3565.12 \pm 1.2%	0 / 3	38 296 \pm 4.1%	
SVMRANK	8.83 \pm 14.9%	2 / 100	46 \pm 32.2%	2902.94 \pm 9.6%	1 / 18	6256 \pm 15.1%	
LMART	7.31 \pm 12.7%	30 / 100	52 \pm 38.1%	3044.94 \pm 7.0%	0 / 12	8893 \pm 3.5%	
GCNN	6.43 \pm 11.6%	51 / 100	43 \pm 40.2%	2024.37 \pm 30.6%	25 / 29	2997 \pm 26.3%	

Outline (applied techniques)

1. GNNs overview
2. Neural algorithmic alignment
3. Integer programming with GNNs
 - i. Machine learning formulation
 - ii. Baselines
 - iii. Experiments
 - iv. Additional research**

Additional research

CPU-friendly approaches

Gupta et al., NeurIPS'20

Bipartite representation inspired many follow-ups

Nair et al., '20; Sonnerat et al., '21; Wu et al., NeurIPS'21; Huang et al. ICML'23; ...

Survey on *Combinatorial Optimization & Reasoning w/ GNNs*:

Cappart, Chételat, Khalil, Lodi, Morris, Veličković, JMLR'23

Conclusions and future directions

Overview

① Theoretical guarantees

- a. Statistical guarantees for algorithm configuration
 - i. Broadly applicable theory for deriving generalization guarantees
 - ii. Proved using connections between primal and dual classes
- b. Online algorithm configuration
 - a. Impossible in the worst cases
 - b. Introduced *dispersion* to provide no-regret guarantees

Overview

① **Theoretical guarantees**

- a. Statistical guarantees for algorithm configuration
- b. Online algorithm configuration

② **Applied techniques:** Graph neural networks

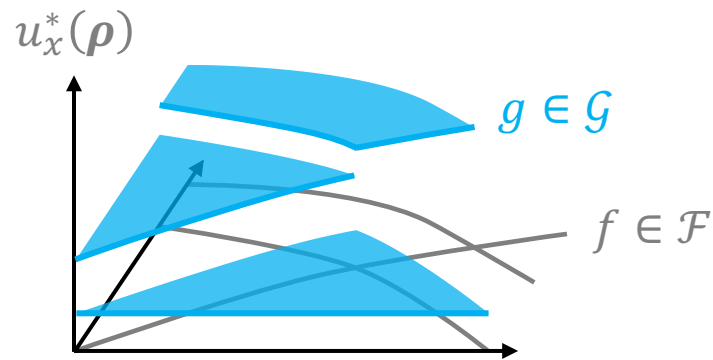
- a. Neural algorithmic alignment
- b. GNNs for variable selection in branch-and-bound

Future work: Tighter statistical bounds

WHP $\forall \rho$, **avg** utility over training set - **exp** utility $\leq \epsilon$

given training set of size $\tilde{O}\left(\frac{1}{\epsilon^2} (\text{Pdim}(\mathcal{G}^*) + \text{VCdim}(\mathcal{F}^*) \log k)\right)$

Number of boundary functions

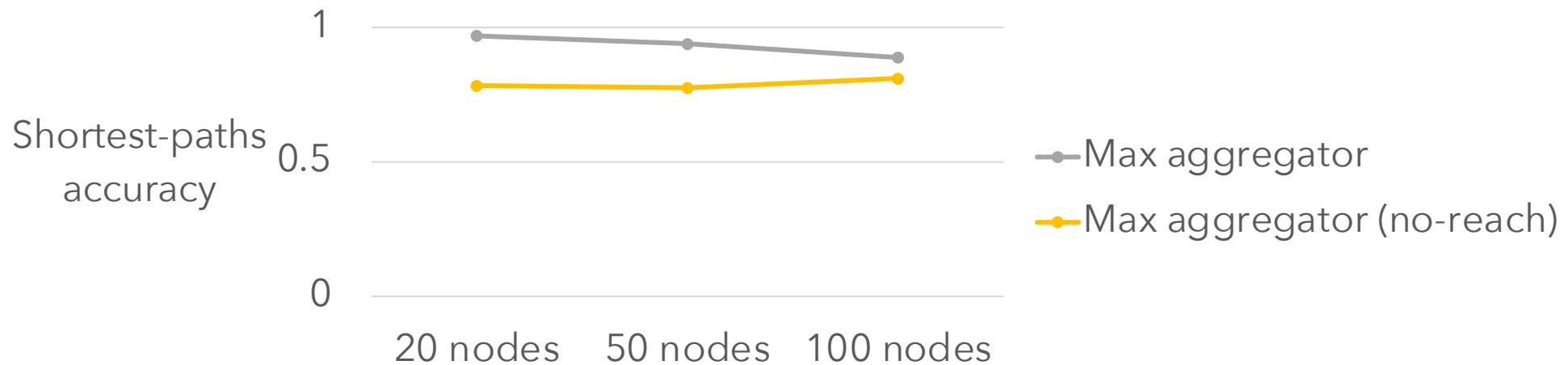


k is often exponential
Can lead to large bounds

I expect this can sometimes be avoided!
Would require more information about duals

Future work: Knowledge transfer

- Training a GNN to solve multiple related problems... can sometimes lead to better **single-task** performance
- E.g., training reachability and shortest-paths (grey line) v.s. just training shortest-paths (**yellow line**)



Future work: Knowledge transfer

- Training a GNN to solve multiple related problems...
can sometimes lead to better **single-task** performance
- Can we understand **theoretically** why this happens?
 - For which sets of algorithms can we expect **knowledge transfer**?

Future work: Size generalization

Machine-learned algorithms can **scale to larger instances**

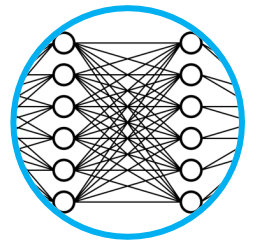
Applied research: Dai et al., NeurIPS'17; Veličković, et al., ICLR'20; ...

Goal: eventually, solve problems **no one's ever been able to solve**

However, size generalization is not immediate! It depends on:

- The **machine-learned algorithm**

Is the algorithm scale sensitive?



Example [Xu et al., ICLR'21]:

- Algorithms represented by GNNs **do generalize**
- Algs represented by MLPs **don't generalize** across size

Future work: Size generalization

Machine-learned algorithms can **scale to larger instances**

Applied research: Dai et al., NeurIPS'17; Veličković, et al., ICLR'20; ...

Goal: eventually, solve problems **no one's ever been able to solve**

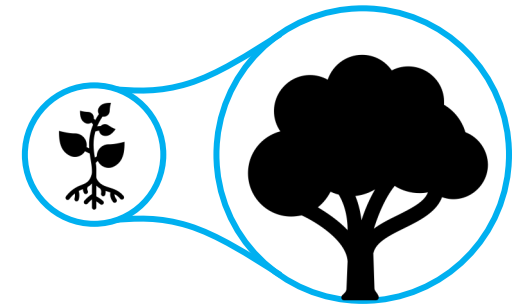
However, size generalization is not immediate! It depends on:

- The **machine-learned algorithm**

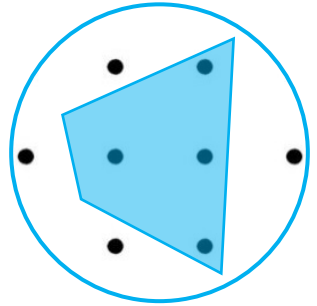
Is the algorithm scale sensitive?

- The **problem instances**

As size scales, what features must be preserved?



Future work: Size generalization



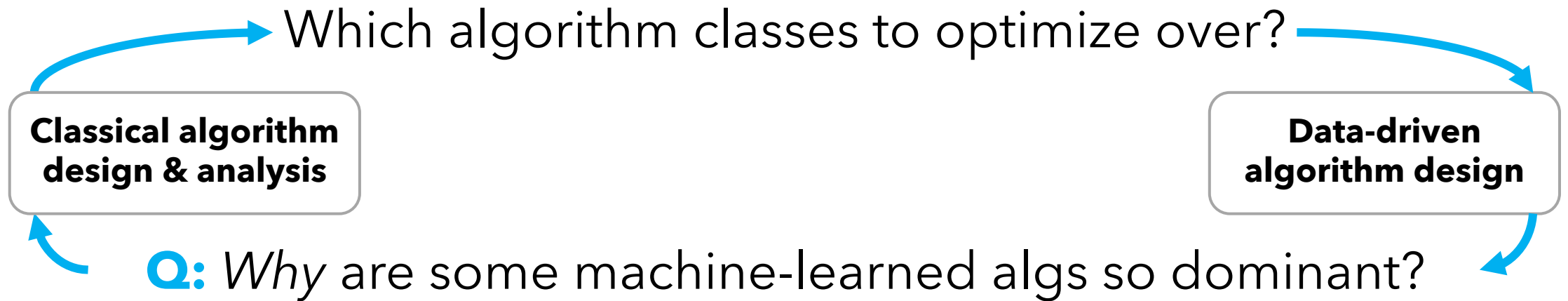
Can you:

1. **Shrink** a set of big integer programs
graphs

...

2. **Learn** a good algorithm on the **small** instances
3. **Apply** what you learned to the **big** instances?

Future work: ML as a toolkit for theory



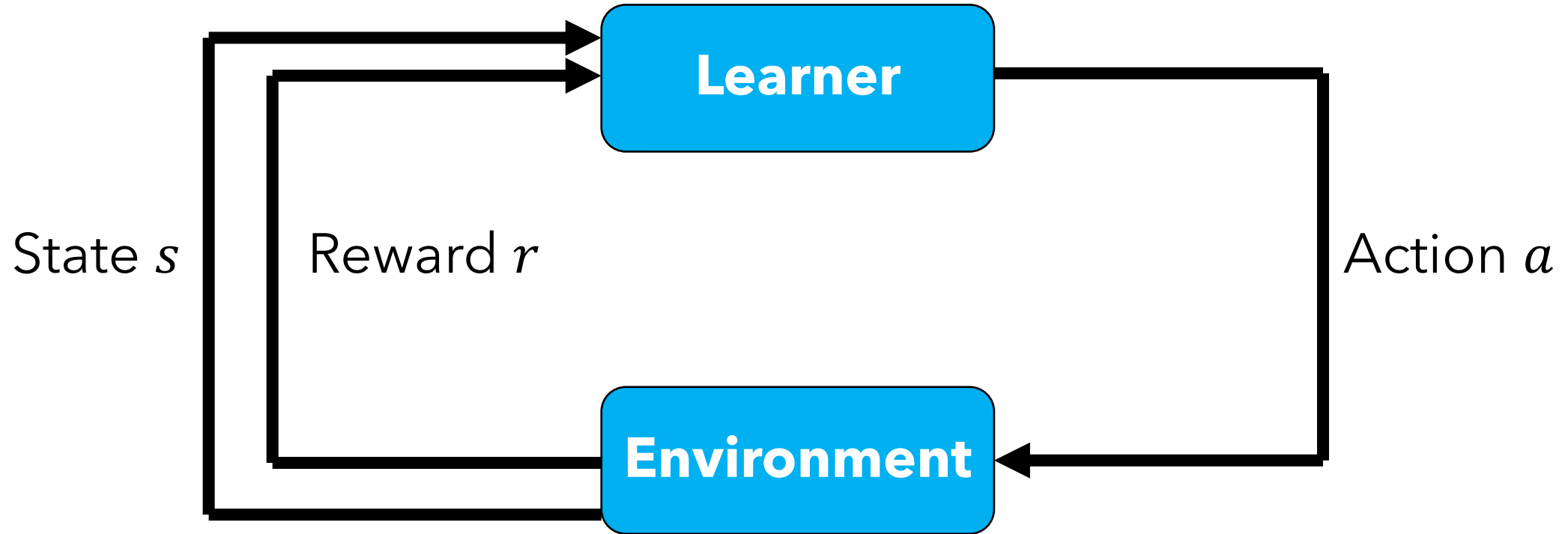
E.g., Dai et al. [NeurIPS'17] write that their RL alg discovered:
"New and interesting" greedy strategies for MAXCUT and MVC
"which **intuitively make sense** but have **not been analyzed** before,"
thus could be a "good **assistive tool** for discovering new algorithms."

Additional slides

Outline (additional applied techniques)

- 1. Reinforcement learning overview**
2. Learning greedy heuristics with RL

Learner interaction with environment



Markov decision processes

S : set of states (assumed for now to be discrete)

A : set of actions

Transition probability distribution $P(s' | s, a)$

Probability of entering state s' from state s after taking action a

Reward function $R: S \rightarrow \mathbb{R}$

Goal: Policy $\pi: S \rightarrow A$ that maximizes total (discounted) reward

Policies and value functions

Policy is a mapping from states to actions $\pi: S \rightarrow A$

Value function for a policy:

Expected sum of discounted rewards

$$V^\pi(s) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \mid s_0 = s, a_t = \pi(s_t), s_{t+1} \mid s_t, a_t \sim P \right]$$

Discount factor

Optimal policy and value function

Optimal policy π^* achieves the highest value for every state

$$V^{\pi^*}(s) = \max_{\pi} V^{\pi}(s)$$

Value function is written $V^* = V^{\pi^*}$

Several different ways to find π^*

- Value iteration
- Policy iteration

Challenge of RL

MDP (S, A, P, R):

- S : set of states (assumed for now to be discrete)
- A : set of actions
- Transition probability distribution $P(s_{t+1} \mid s_t, a_t)$
- Reward function $R: S \rightarrow \mathbb{R}$

RL twist: We don't know P or R , or too big to enumerate

Q-learning

Q functions:

Like value functions but defined over state-action pairs

$$Q^\pi(s, a) = R(s) + \gamma \sum_{s' \in S} P(s' | s, a) Q^\pi(s', \pi(s'))$$

I.e., Q function is the value of:

1. Starting in state s
2. Taking action a
3. Then acting according to π

Q-learning

$$\begin{aligned} Q^*(s, a) &= R(s) + \gamma \sum_{s' \in \mathcal{S}} P(s' | s, a) \max_{a'} Q^*(s', a') \\ &= R(s) + \gamma \sum_{s' \in \mathcal{S}} P(s' | s, a) V^*(s') \end{aligned}$$

Q^* is the value of:

1. Starting in state s
2. Taking action a
3. Then acting optimally

Q-learning

(High-level) Q-learning algorithm

initialize $\hat{Q}(s, a) \leftarrow 0, \forall s, a$

repeat

Observe current state s and reward r

Take action $a = \operatorname{argmax} \hat{Q}(s, \cdot)$ and observe next state s'

Improve estimate \hat{Q} based on s, r, a, s'

Can use *function approximation* to represent \hat{Q} compactly

$$\hat{Q}(s, a) = f_{\theta}(s, a)$$

Outline (additional applied techniques)

1. Reinforcement learning overview
- 2. Learning greedy heuristics with RL**

RL for combinatorial optimization

Tons of research in this area

Travelling salesman

Bello et al., ICLR'17; Dai et al., NeurIPS'17;
Nazari et al., NeurIPS'18; ...

Bin packing

Hu et al., '17; Laterre et al., '18; Cai et al.,
DRL4KDD'19; Li et al., '20; ...

Maximum cut

Dai et al., NeurIPS'17; Cappart et al.,
AAAI'19; Barrett et al., AAAI'20; ...

Minimum vertex cover

Dai et al., NeurIPS'17; Song et al., UAI'19; ...

This section: Example of a pioneering work in this space

Overview

Goal: use RL to learn new *greedy strategies* for graph problems
Feasible solution constructed by successively adding nodes to solution

Input: Graph $G = (V, E)$, weights $w(u, v)$ for $(u, v) \in E$

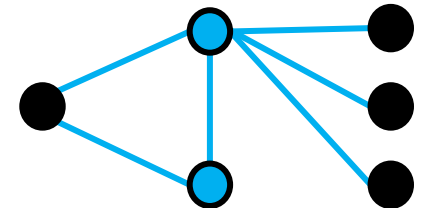
RL state representation: Graph embedding

Outline (additional applied techniques)

1. Reinforcement learning overview
2. Learning greedy heuristics with RL
 - i. Examples: Min vertex cover and max cut**
 - ii. RL formulation
 - iii. Experiments

Minimum vertex cover

Find smallest vertex subset such that each edge is covered

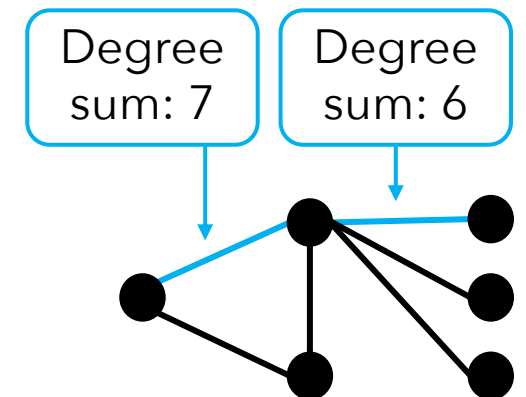


Minimum vertex cover

Find smallest vertex subset such that each edge is covered

2-approximation:

Greedily add vertices of edge with **maximum degree sum**



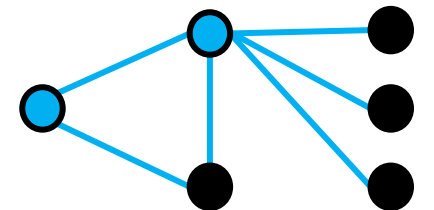
Minimum vertex cover

Find smallest vertex subset such that each edge is covered

2-approximation:

Greedly add vertices of edge with maximum degree sum

Scoring function that guides greedy algorithm



Maximum cut

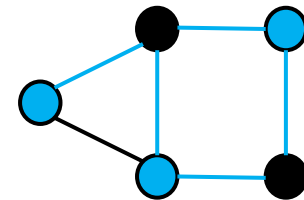
Find partition $(S, V \setminus S)$ of nodes that maximizes

$$\sum_{(u,v) \in C} w(u,v)$$

where $C = \{(u,v) \in E : u \in S, v \notin S\}$

If $w(u,v) = 1$ for all $(u,v) \in E$:

$$\sum_{(u,v) \in C} w(u,v) = 5$$



Maximum cut

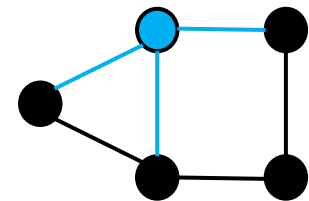
Find partition $(S, V \setminus S)$ of nodes that maximizes

$$\sum_{(u,v) \in C} w(u,v)$$

where $C = \{(u, v) \in E : u \in S, v \notin S\}$

Greedy: move node from one side of cut to the other

Move node that results in the largest improvement in cut weight



Maximum cut

Find partition $(S, V \setminus S)$ of nodes that maximizes

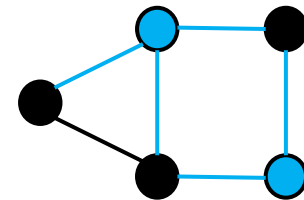
$$\sum_{(u,v) \in C} w(u,v)$$

where $C = \{(u,v) \in E : u \in S, v \notin S\}$

Greedy: move node from one side of cut to the other

Move node that results in the largest improvement in cut weight

Scoring function that guides greedy algorithm



Outline (additional applied techniques)

1. Reinforcement learning overview
2. Learning greedy heuristics with RL
 - i. Example: Min vertex cover and max cut
 - ii. RL formulation**
 - iii. Experiments

Reinforcement learning formulation

State:

- *Goal*: encode partial solution $S = \underline{(v_1, v_2, \dots, v_{|S|})}$, $v_i \in V$

E.g., nodes in independent set, nodes on one side of cut

Reinforcement learning formulation

State:

- *Goal*: encode partial solution $S = (v_1, v_2, \dots, v_{|S|})$, $v_i \in V$
- Use GNN to compute graph embedding μ

$$\text{Initial node features } x_v = \begin{cases} 1 & \text{if } v \in S \\ 0 & \text{else} \end{cases}$$

Action: Choose vertex $v \in V \setminus S$ to add to solution

Transition (deterministic): For chosen $v \in V \setminus S$, set $x_v = 1$

Reinforcement learning formulation

Reward: $r(S, v)$ is change in objective when transition $S \rightarrow (S, v)$

Policy (deterministic): $\pi(v|S) = \begin{cases} 1 & \text{if } v = \operatorname{argmax}_{v' \in S} \hat{Q}(\mu, v') \\ 0 & \text{else} \end{cases}$

Outline (additional applied techniques)

1. Reinforcement learning overview
2. Learning greedy heuristics with RL
 - i. Example: Min vertex cover and max cut
 - ii. RL formulation
 - iii. Experiments**

Min vertex cover

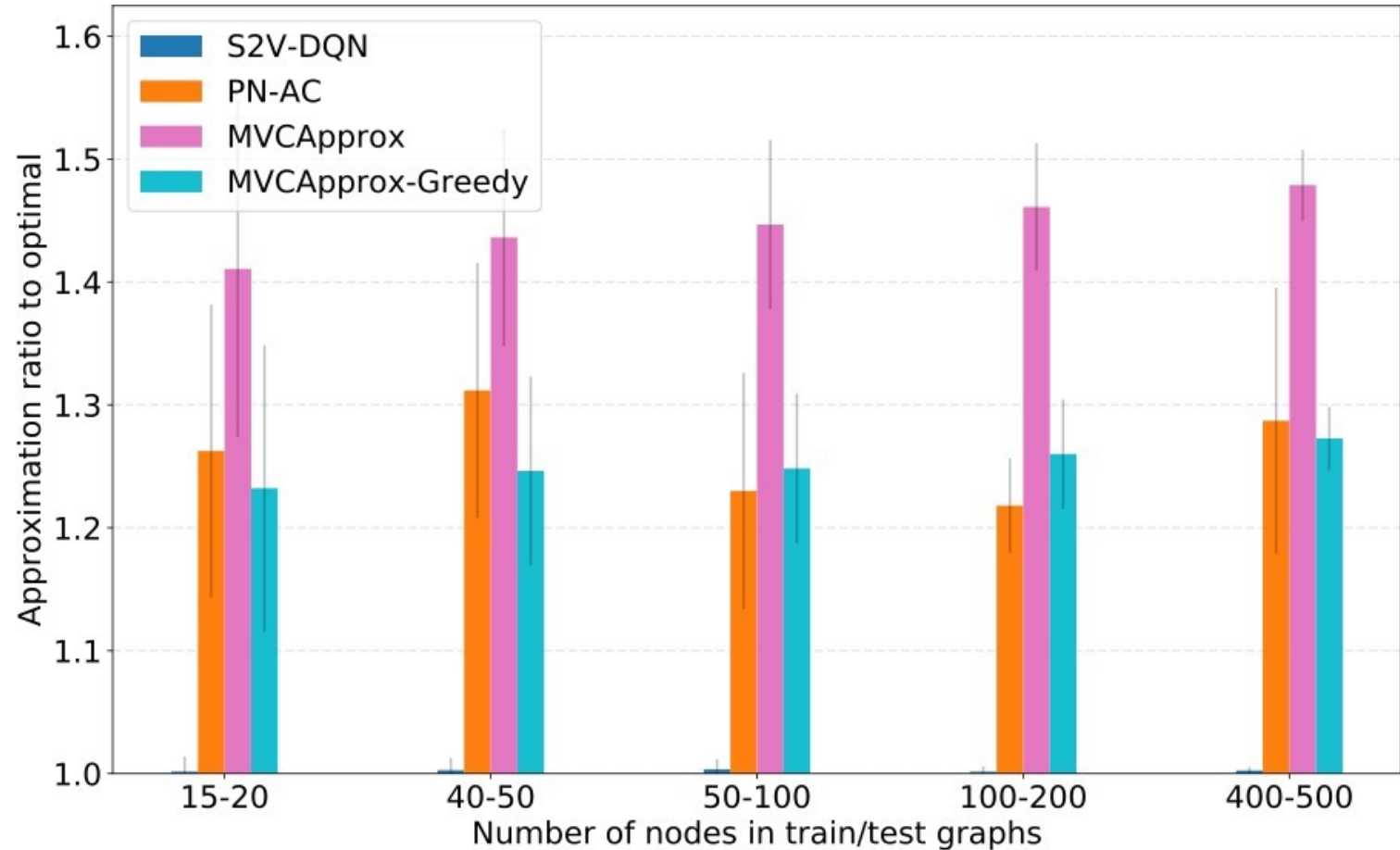
Barabasi-Albert
random graphs

Paper's approach

Another DL approach
[Bello et al., arXiv'16]

2-approximation
algorithm

Greedy algorithm
from first few slides



Max cut

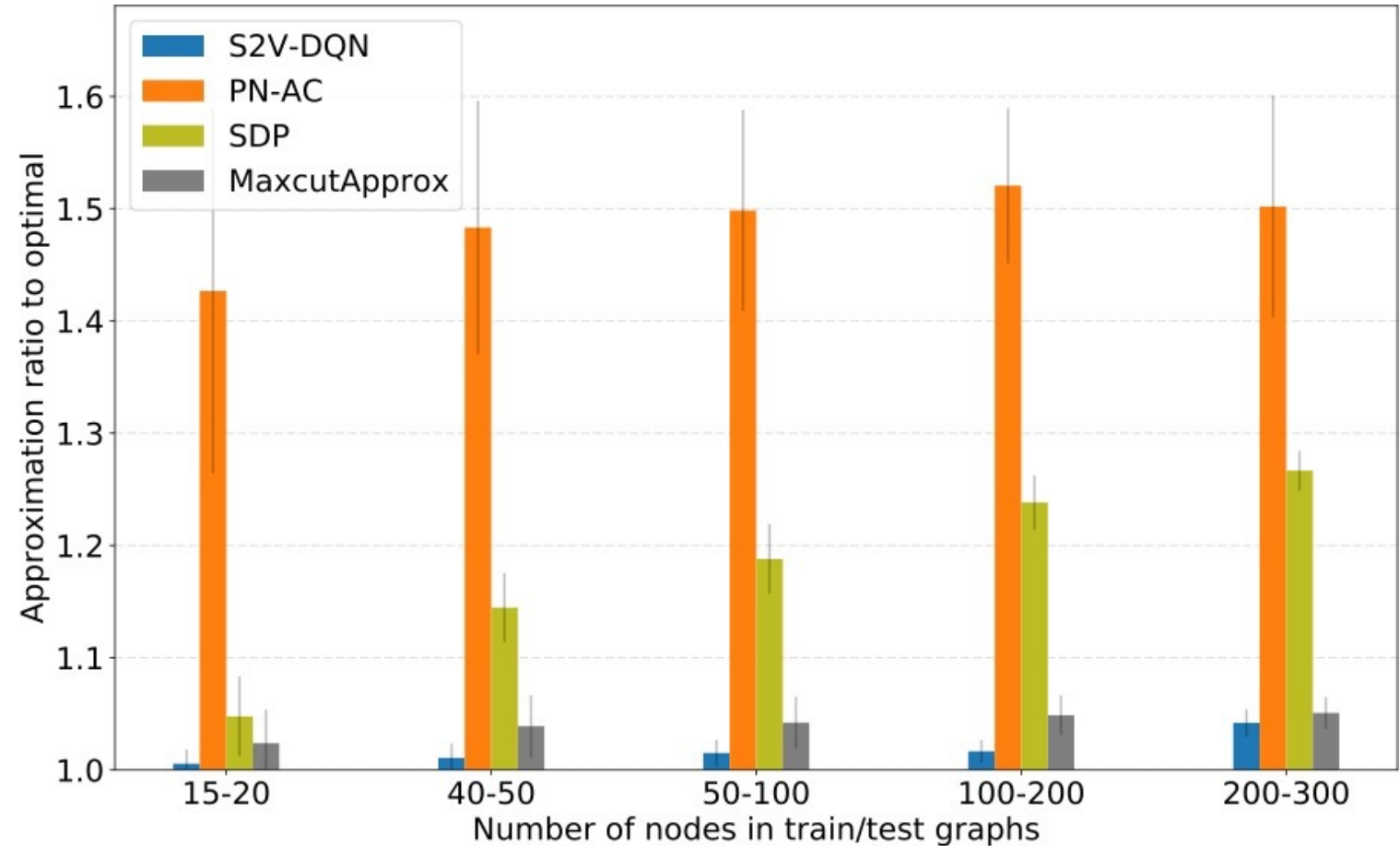
Barabasi-Albert
random graphs

Paper's approach

Another DL approach
[Bello et al., arXiv'16]

Goemans-Williamson
algorithm

Greedy algorithm
from first few slides



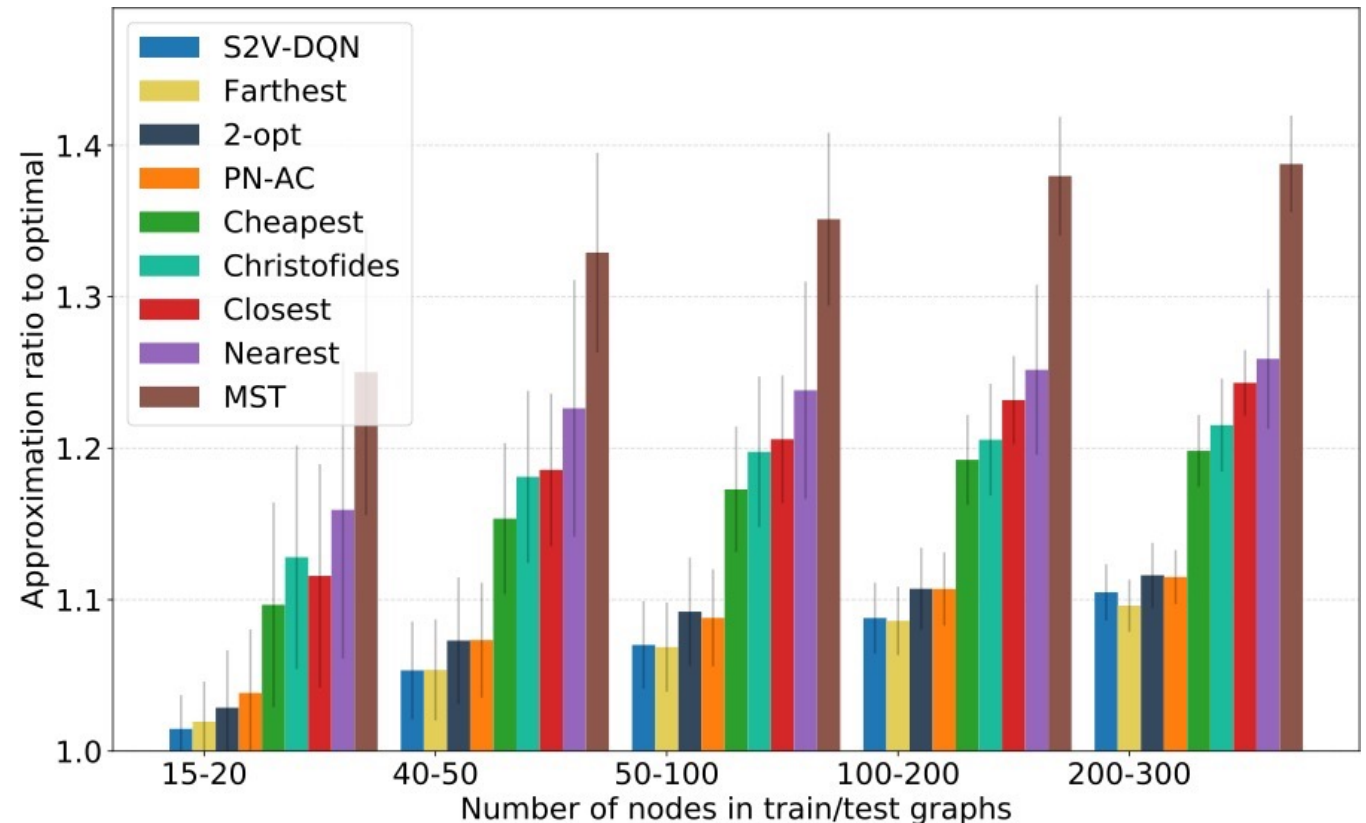
TSP

Uniform random points on 2-D grid

Paper's approach

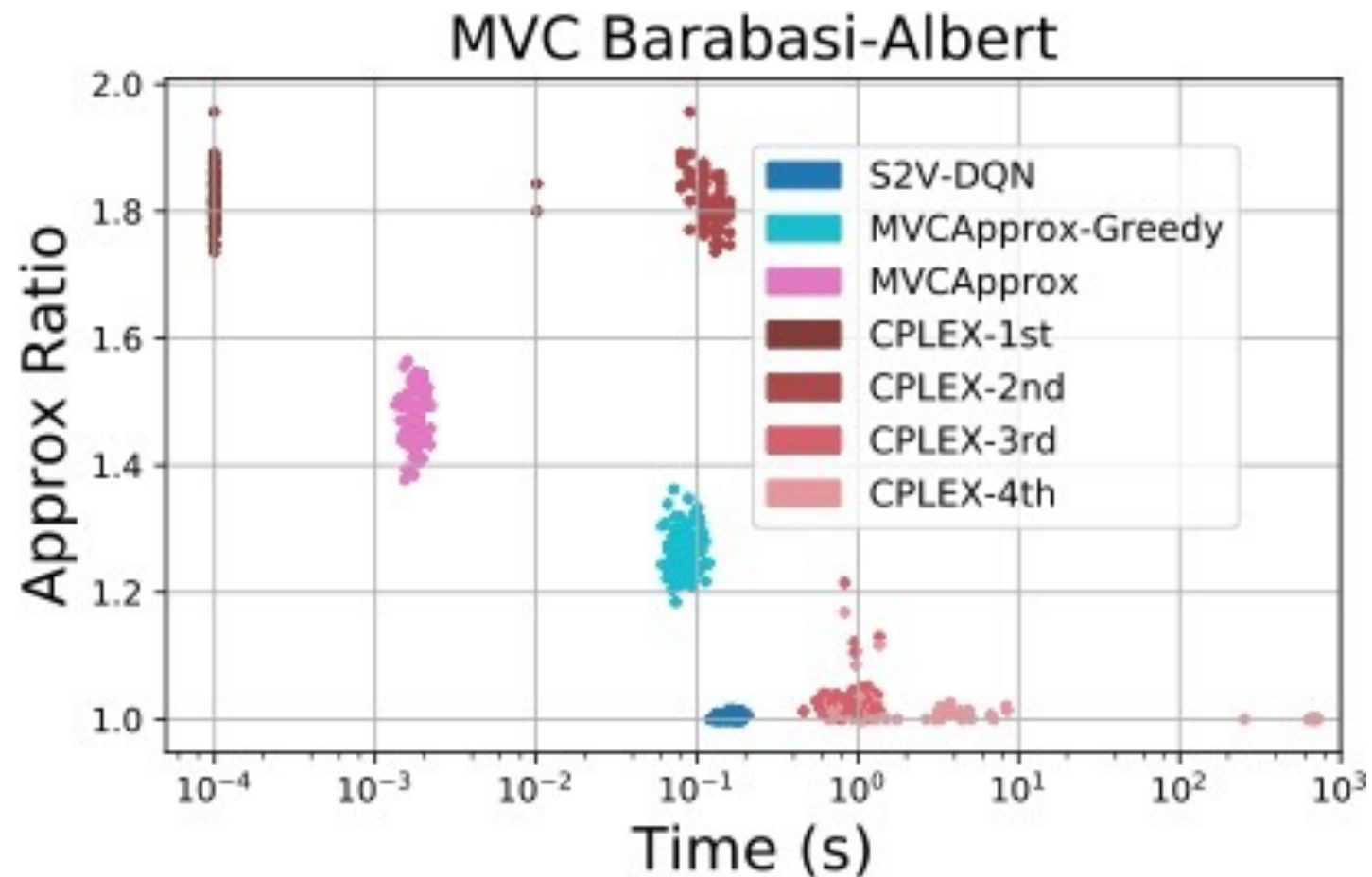
- Initial subtour: 2 cities that are farthest apart
- Repeat the following:
 - Choose city that's *farthest* from any city in the subtour
 - Insert in position where it causes the smallest distance increase

[Rosenkrantz et al., SIAM JoC'77]

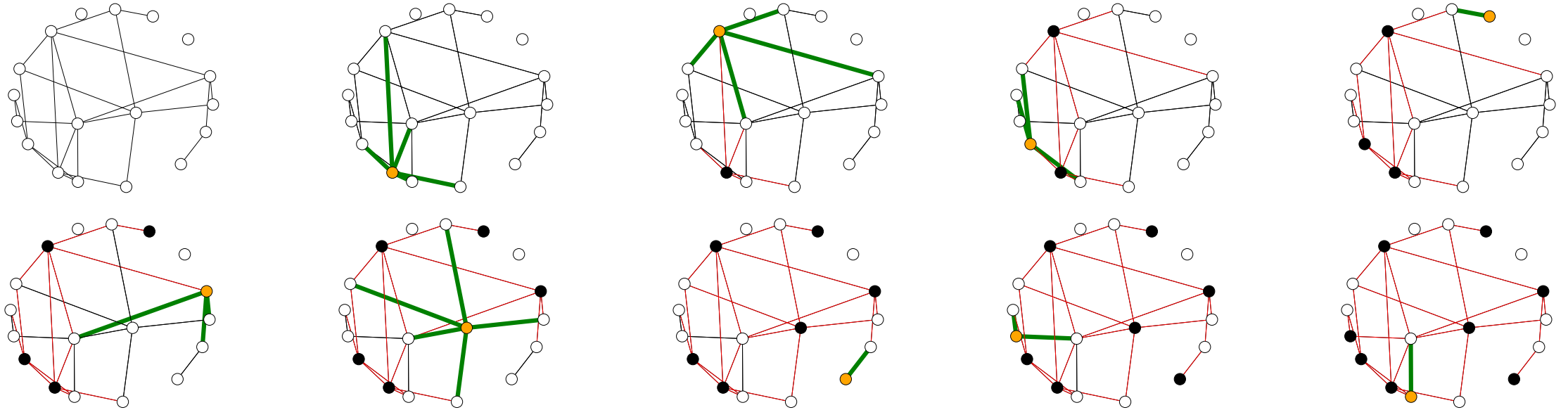


Runtime comparisons

CPLEX-1st: 1st feasible solution found by CPLEX



Min vertex cover visualization



Nodes seem to be selected to balance between:

- Degree
- Connectivity of the remaining graph