

# Posterior agreement and stochastic optimization

Statistical learning theory 2020

Maximum entropy  
and  
posterior agreement

## Motivation

$$X = (x_1, \dots, x_n)$$

$$x_i \sim \mathcal{N}(i, \sigma^2)$$

↓ draw

$$X' = (x'_1, \dots, x'_n) \in \mathbb{R}^n$$

instances

$$X'' = (x''_1, \dots, x''_n) \in \mathbb{R}^n$$

Goal:  $\min_i \mathbb{E}[x_i]$

## Formalization

$$\min_c \mathbb{E}_X [R(c, X)]$$

$$X = (X_1, \dots, X_n)$$

$$G = \{1, \dots, n\}$$

$$R(c, X) = X_c$$

# Problem :

$$\min_{c \in C} \mathbb{E}_X [R(c, X)]$$

$C$  : solution space

$X$  : random instance

$R(c, X)$  : cost of solution  $c$  on instance  $X$ .

# Max entropy + PA

Problem: find  $\min_c \mathbb{E}_X [R(c, X)]$ , given  $X', X''$ .

Approach:

1. Compute  $p_{T^*}(\cdot | X') \propto \exp\left(-\frac{1}{T^*} R(\cdot, X')\right)$

$$p_{T^*}(\cdot | X'') \propto \exp\left(-\frac{1}{T^*} R(\cdot, X'')\right)$$

2.  $\hat{c} \leftarrow \frac{1}{Z} p_{T^*}(\cdot | X') \odot p_{T^*}(\cdot | X'')$

$$T^* = \arg \max_T K_T(X', X'')$$

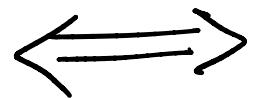
- \* Why ME ?
- \* What's the right value of T?

## Why ME?

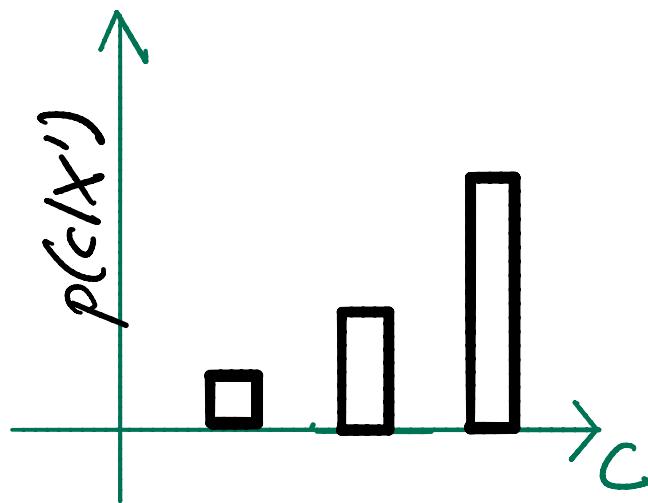
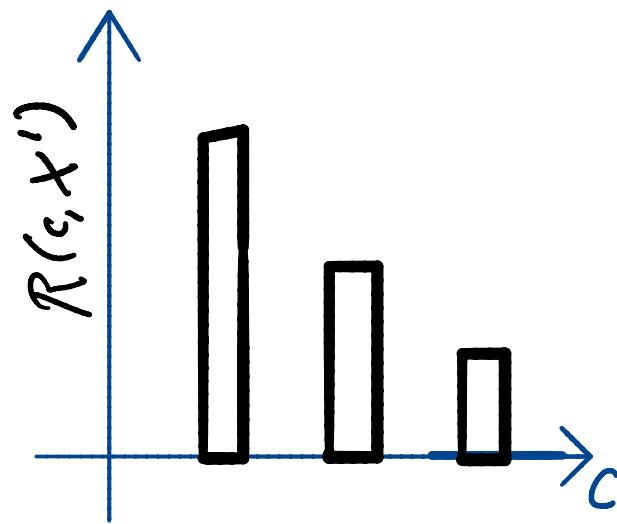
- Algos should output distributions
  - $p(\cdot | X)$  should be the "most general" distribution such that
- \*  $R(c_1, X) \leq R(c_2, X) \Leftrightarrow p(c_1 | X) \geq p(c_2 | X)$
- low cost  $\iff$  high prob.
- \*  $\mathbb{E}_C [R(C, X)] = \text{const} \in [\min R(\cdot, X), \text{avg } R(\cdot, X)]$

$$p(c|x') \propto \exp\left(-\frac{1}{T} R(c, x')\right)$$

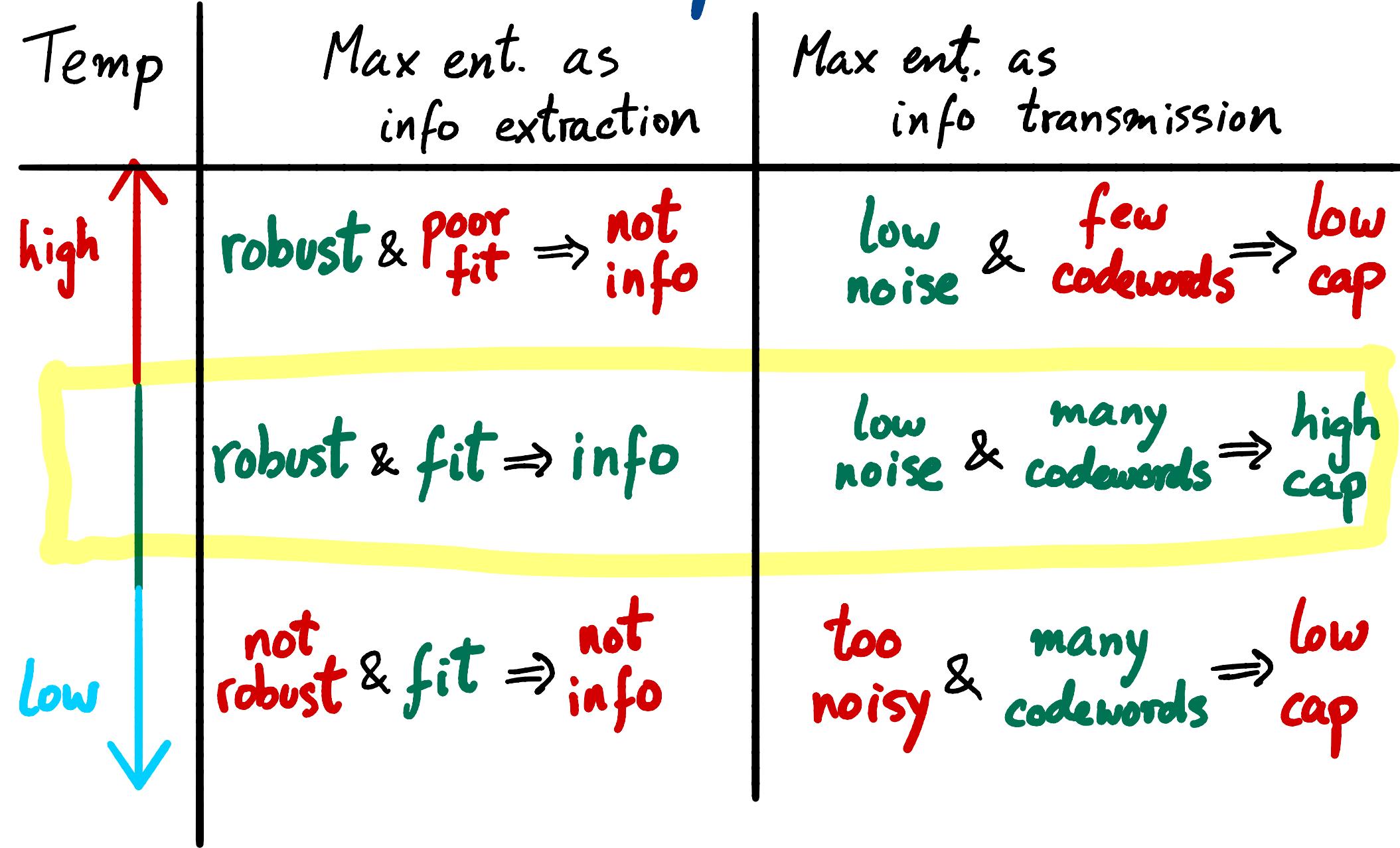
low cost



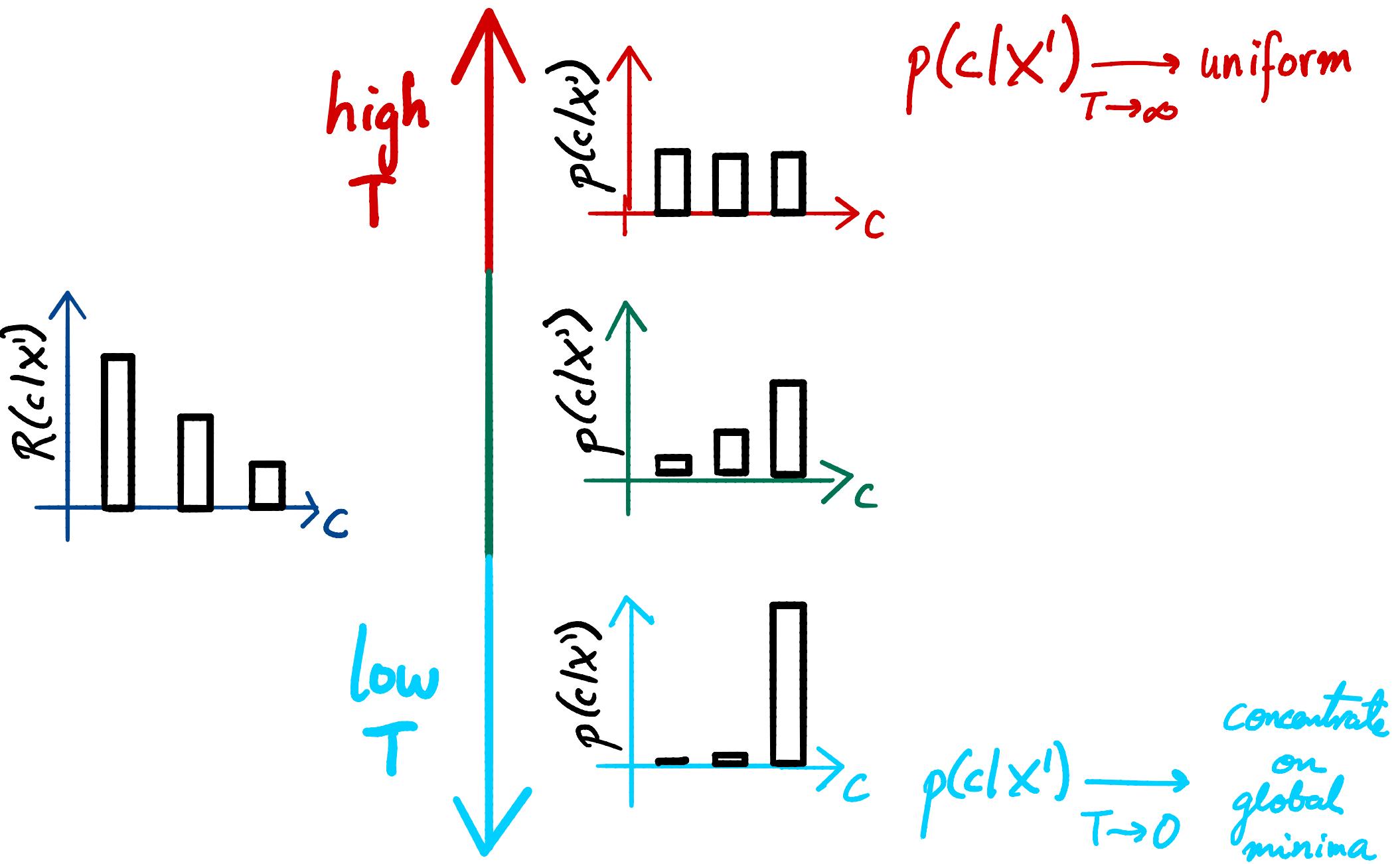
high prob.



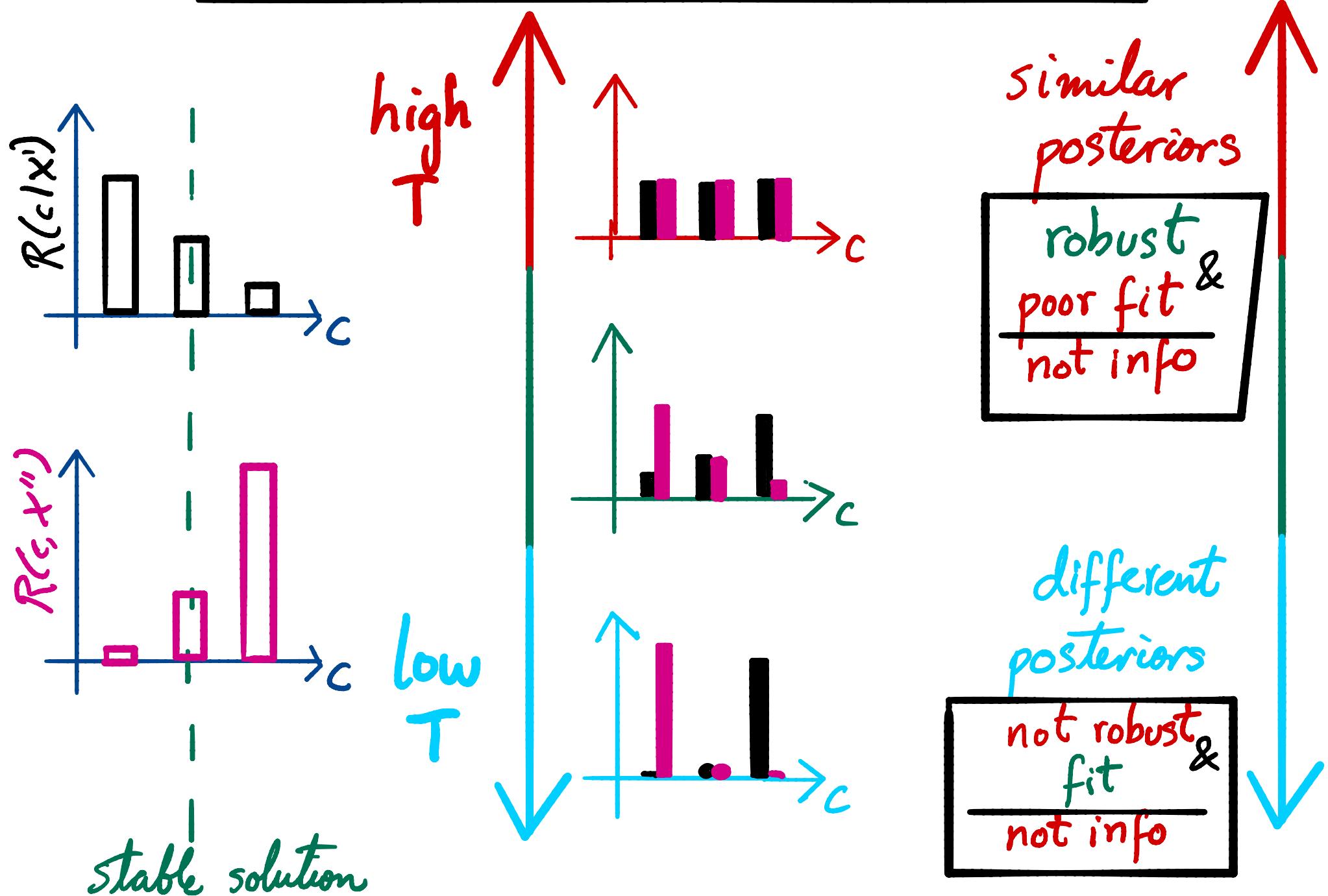
# What is the right value for $T$ ?



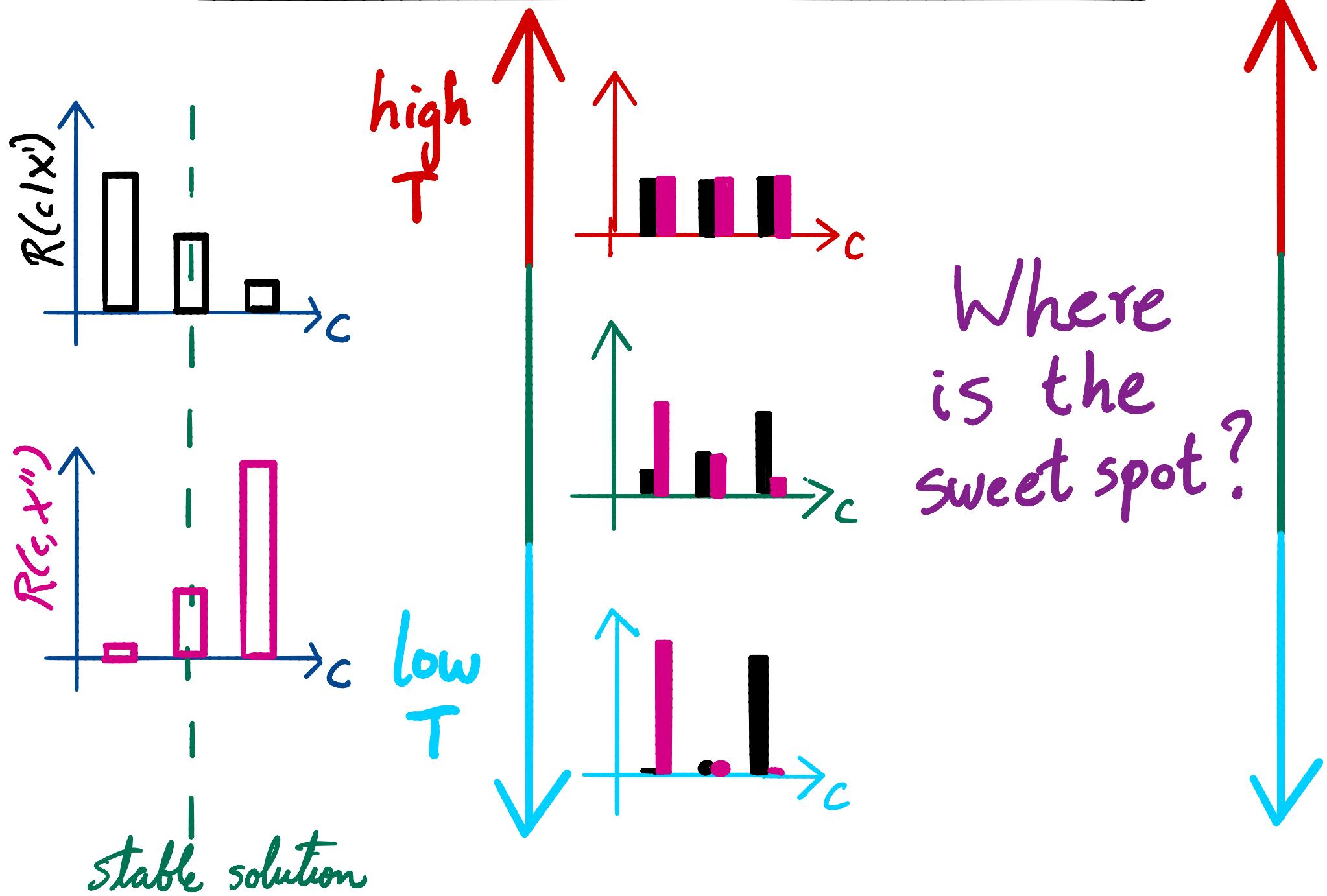
# Max ent as info extraction

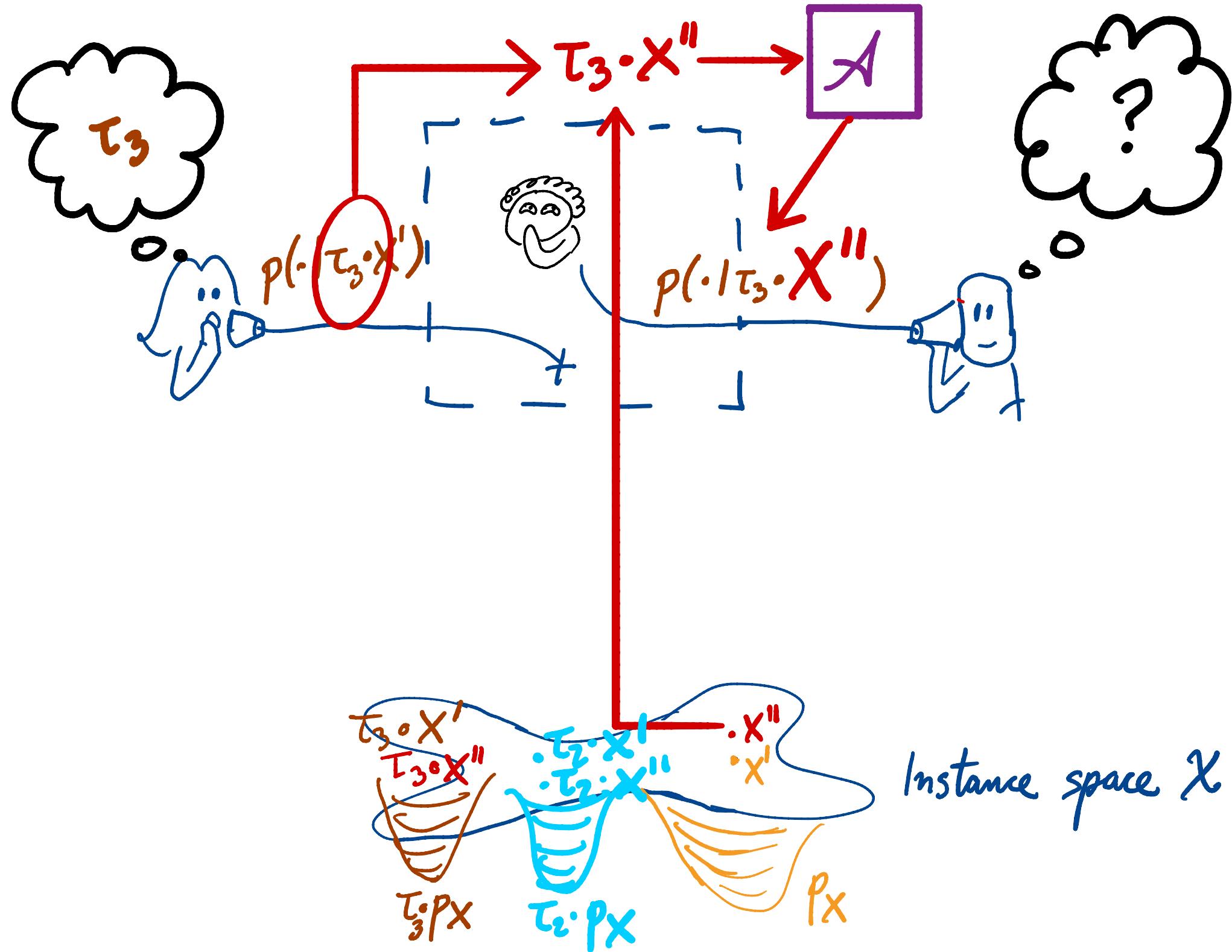


# Max ent as info extraction



# Max ent as info extraction

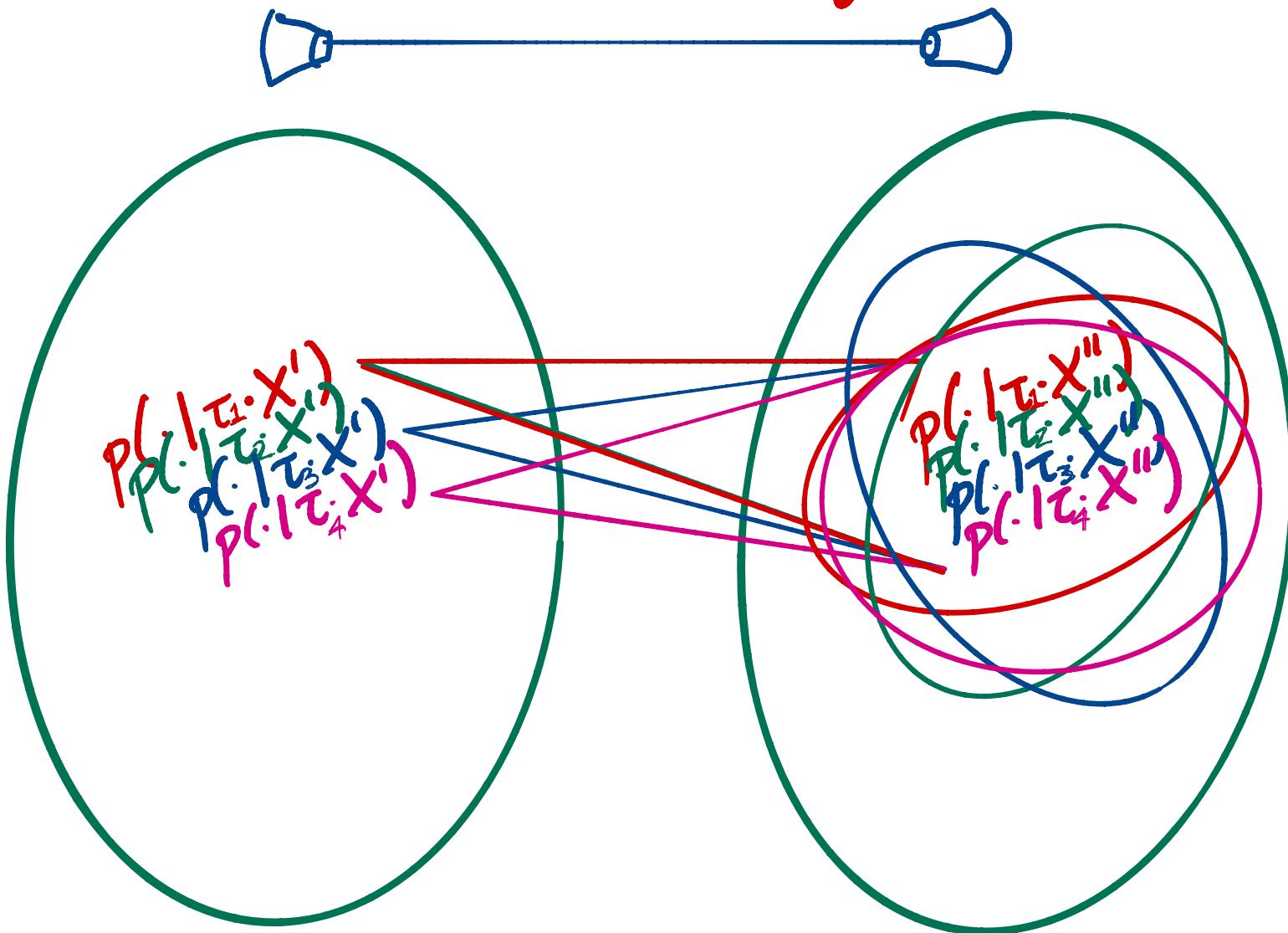




high T: too close codewords and  
narrow typical sets

---

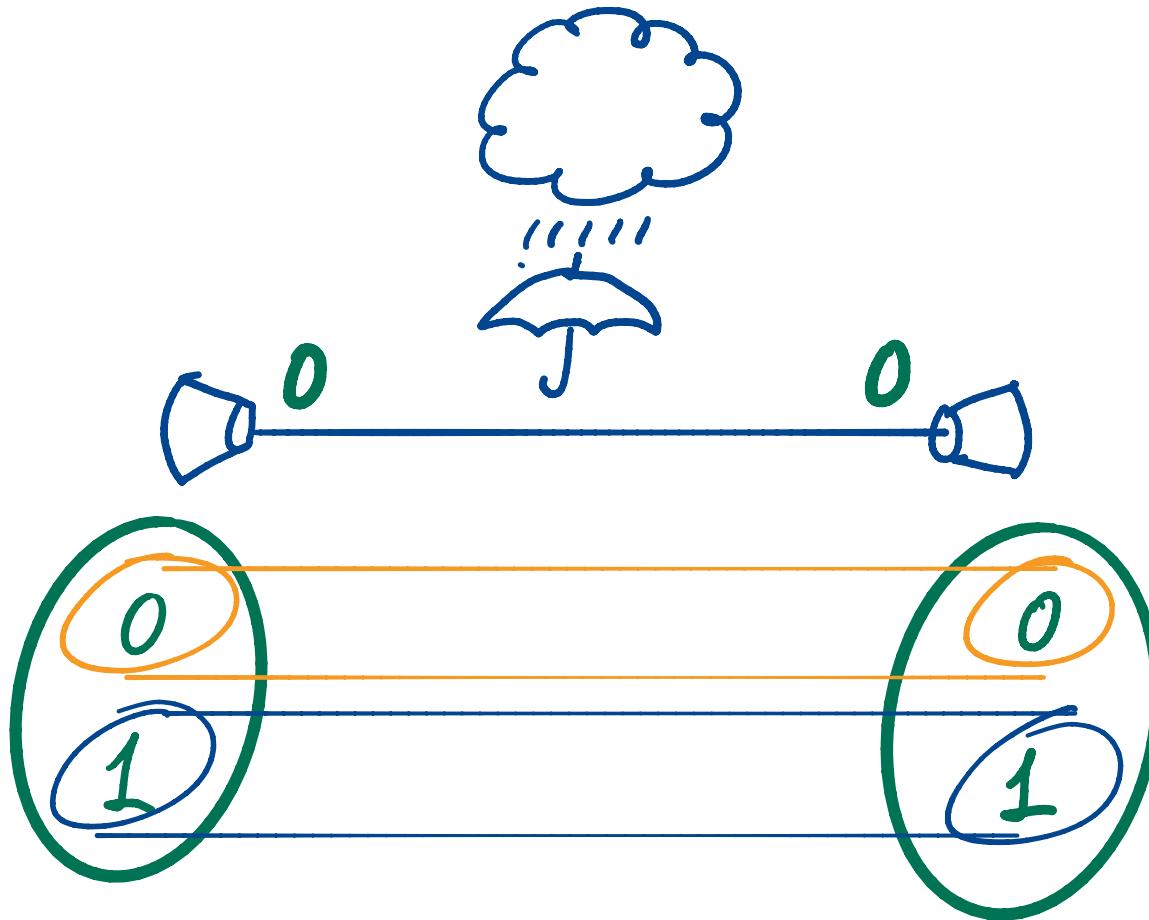
few messages



high T: few codewords and  
narrow typical sets

---

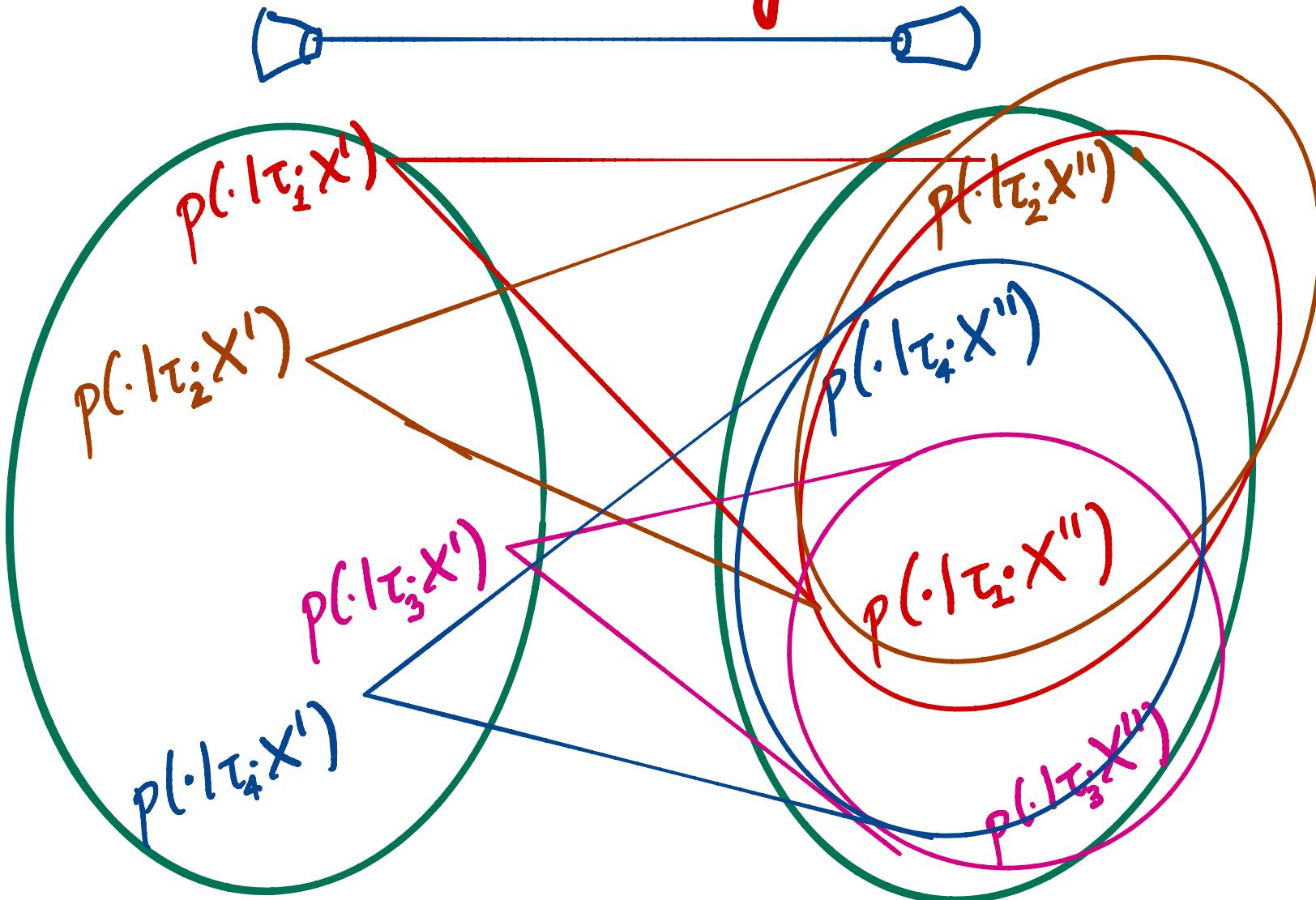
few messages



low T:

distant posteriors and  
too wide typical sets

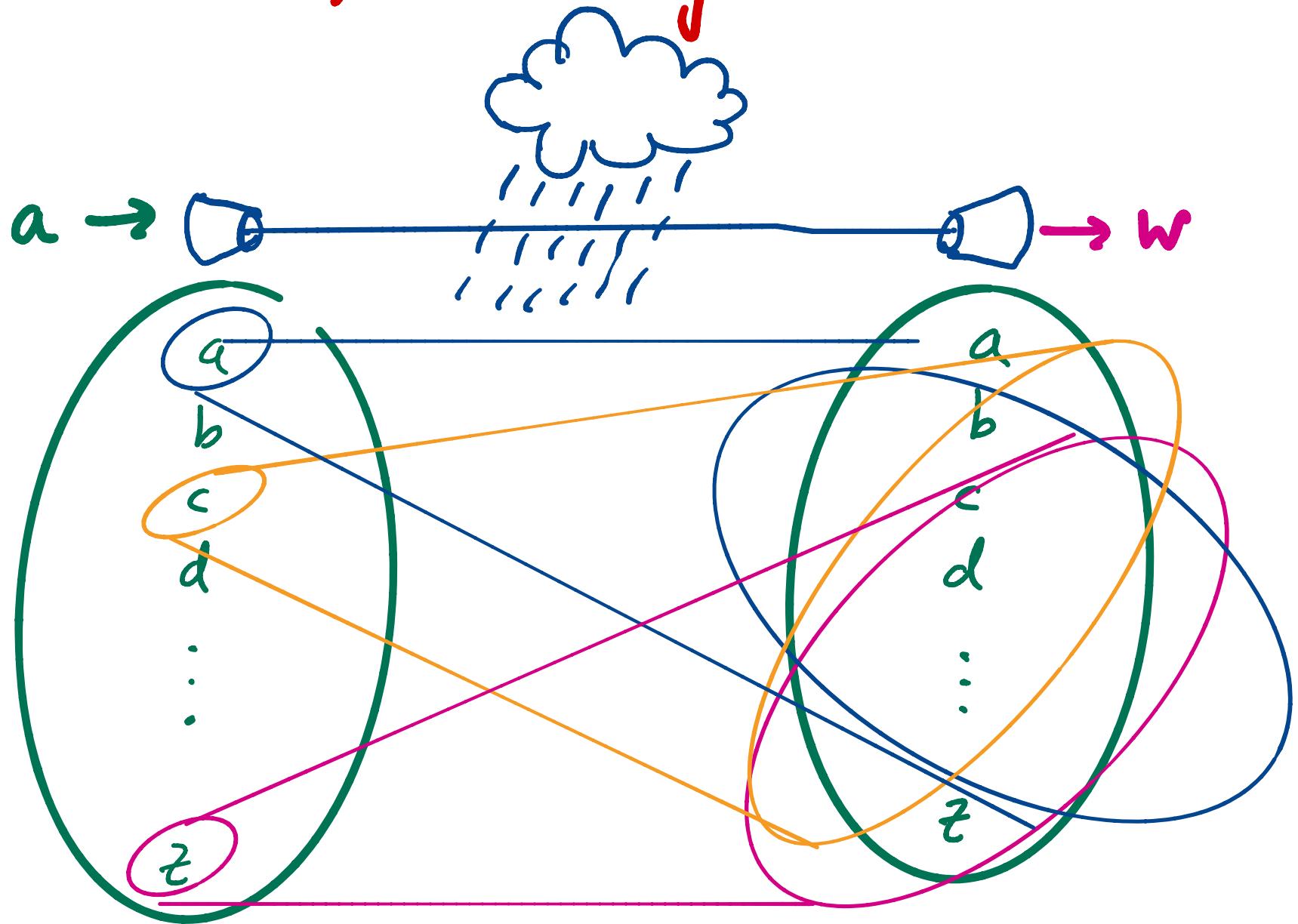
few messages



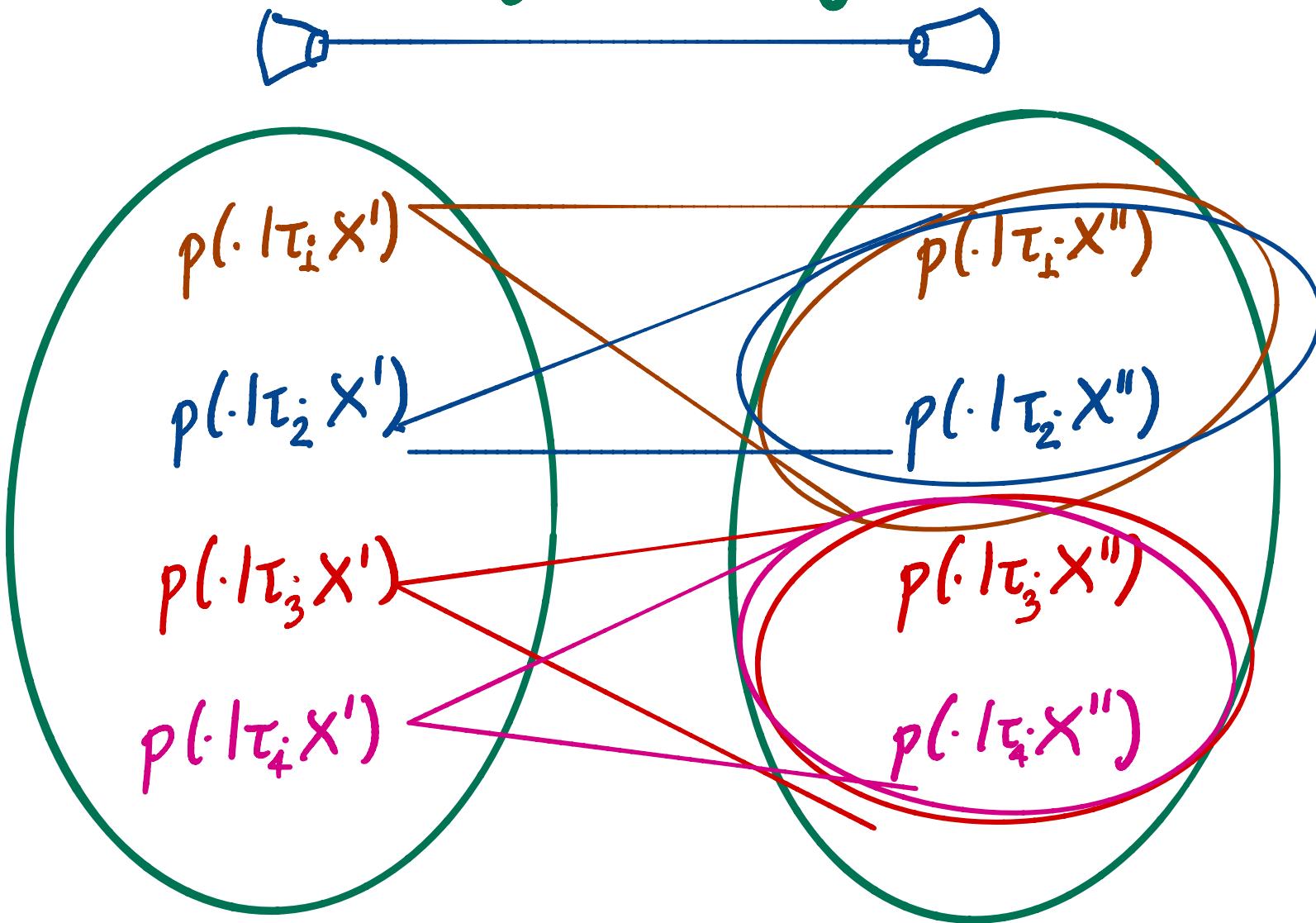
low T:

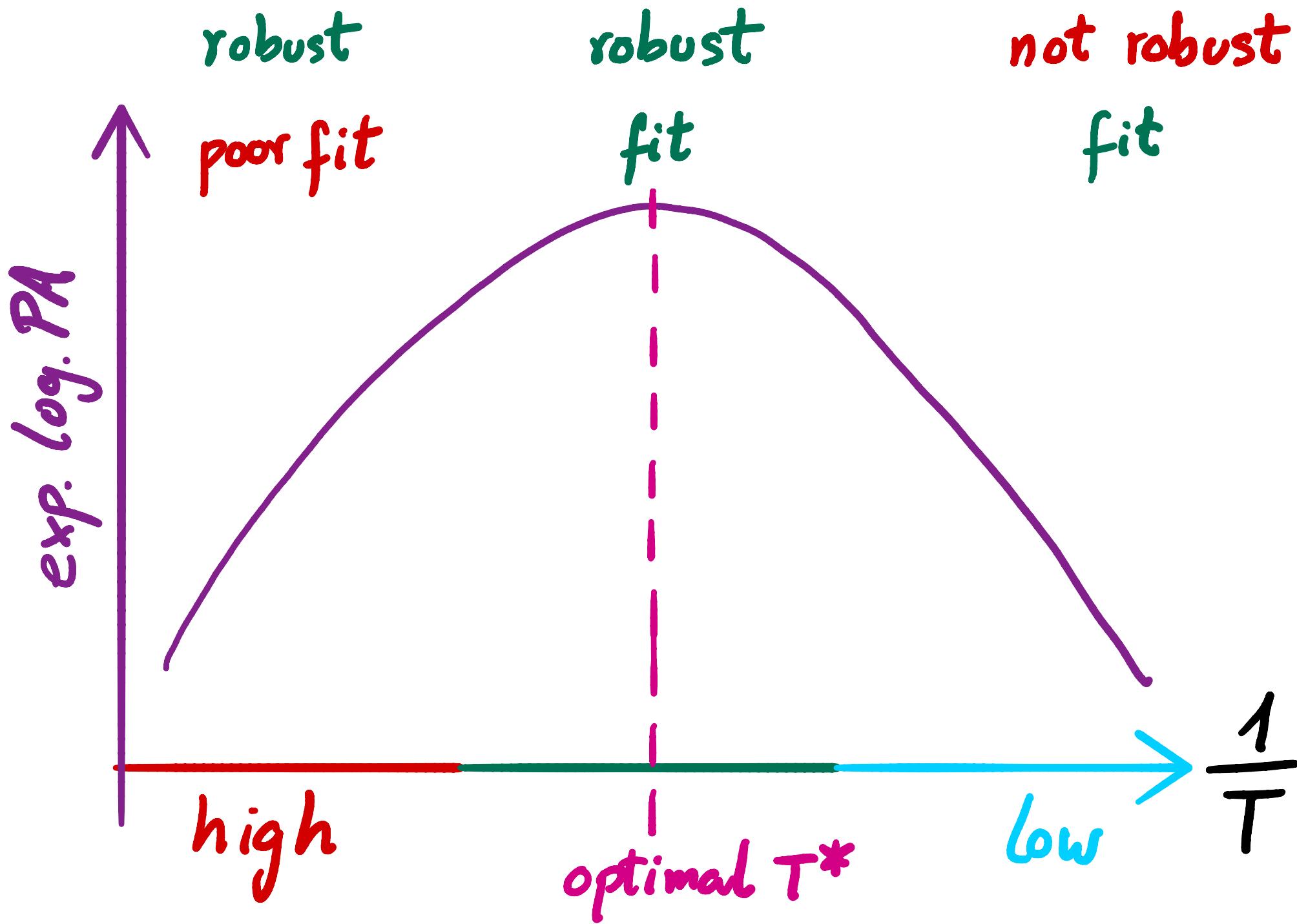
many codewords and  
too wide typical sets

few messages



optimal : distant posteriors and narrow typical sets  
 $\tau$  many messages!





Posterior agreement recommends

$$T^* = \arg \max_T \text{exp. log. PA}$$

$\Sigma E_{x', x''} [\log |C| K(x', x'')] \quad \text{PA}$

$$\approx \arg \max_T \text{emp. log. PA}$$

$$= \arg \max_T \log |C| K(x', x'')$$

$$= \arg \max_T K_T(x', x'')$$

with  $K_T(x', x'') = \langle p_T(\cdot | x'), p_T(\cdot | x'') \rangle$

# Max entropy + PA

Problem: find  $\min_c \mathbb{E}_X [R(c, X)]$ , given  $X', X''$ .

Approach:

1. Compute  $p_{T^*}(\cdot | X') \propto \exp\left(-\frac{1}{T^*} R(\cdot, X')\right)$

$$p_{T^*}(\cdot | X'') \propto \exp\left(-\frac{1}{T^*} R(\cdot, X'')\right)$$

2.  $\hat{c} \leftarrow \frac{1}{Z} p_{T^*}(\cdot | X') \odot p_{T^*}(\cdot | X'')$

$$T^* = \arg \max_T K_T(X', X'')$$

# Experiments comparing ME vs ERM

$$G = \{1, \dots, s, s+1, \dots, s+u\}$$

$$\vec{X} = \left( \underbrace{X_1, \dots, X_s}_{\text{stable}}, \underbrace{\tilde{X}_1, \dots, \tilde{X}_u}_{\text{unstable}} \right)$$

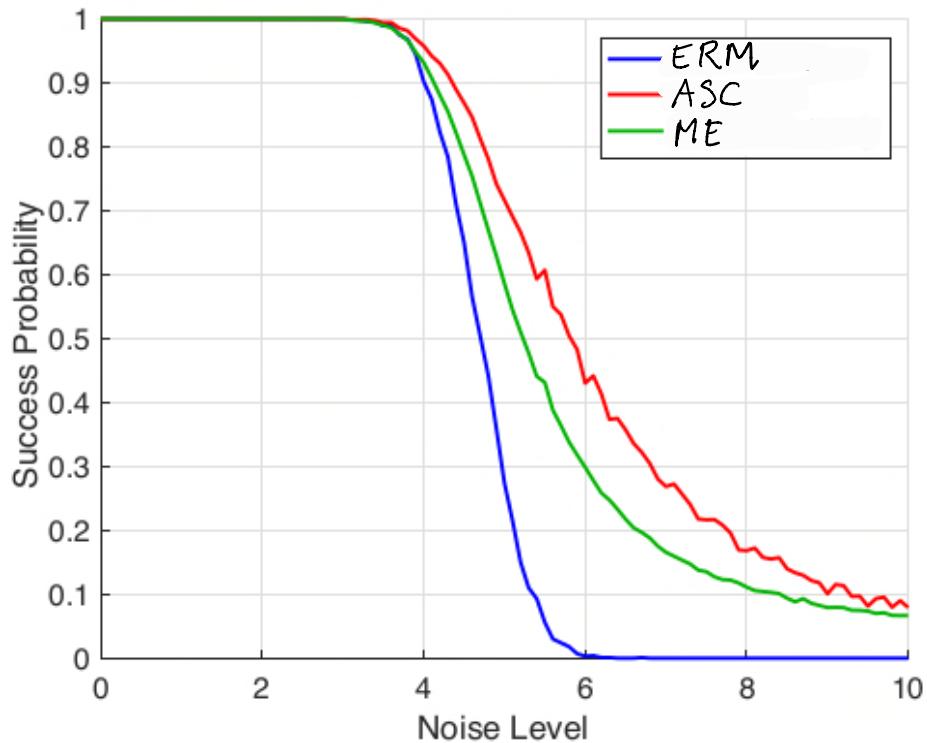
$$X_i \sim \mathcal{N}(1, 1) \quad \tilde{X}_i \sim \mathcal{N}(10, \sigma^2)$$

noise parameter

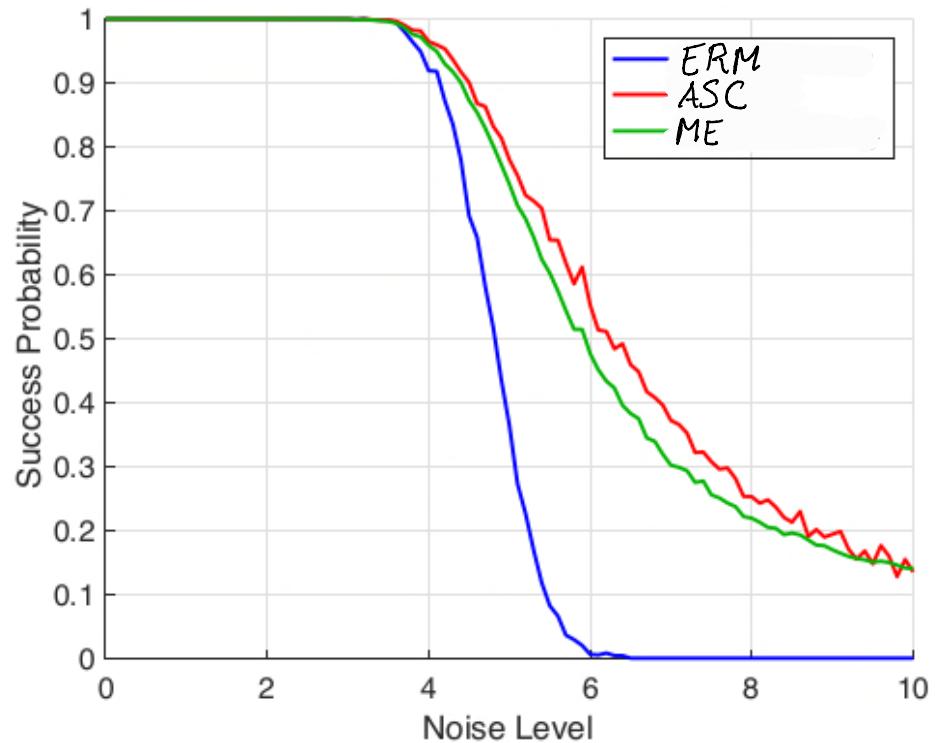
Observe that if  $c \leq s$  then  $\mathbb{E}_{\vec{X}}[R(c, \vec{X})] = \mathbb{E}[X_c] = 1$

if  $s < c \leq s+u$  then  $\mathbb{E}_{\vec{X}}[R(c, \vec{X})] = \mathbb{E}[X_c] = 10$

Alex's experiments on ME vs ERM  
using only two instances  $X'$  and  $X''$ .



(a) 5% of solutions are stable.

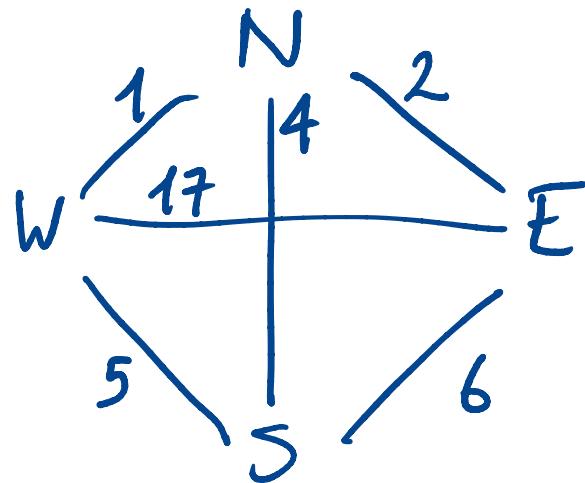


(b) 10% of solutions are stable.

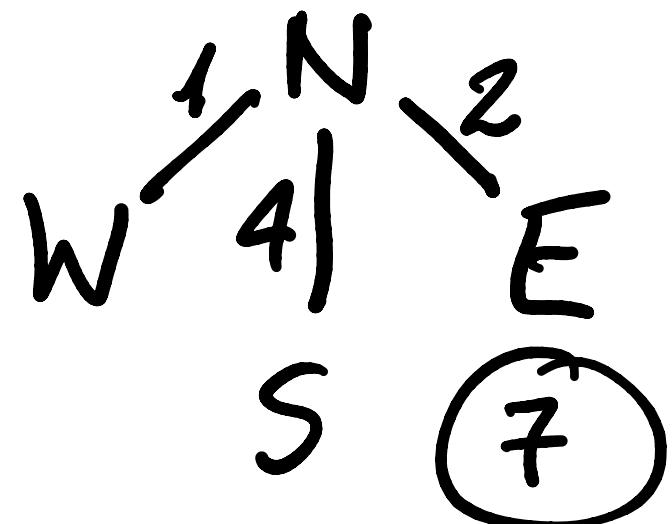
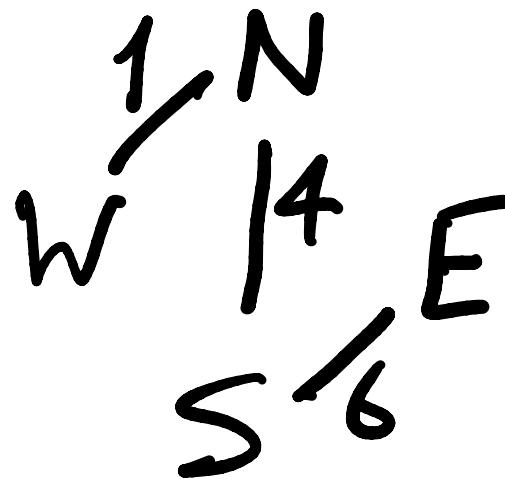
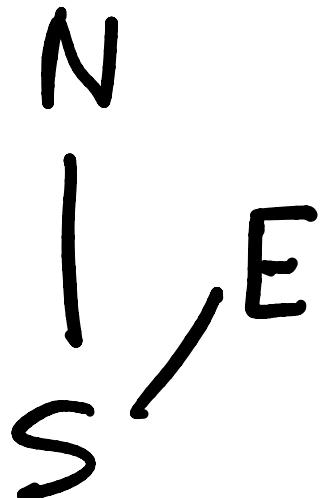
Minimum spanning tree algorithms:  
Regularization by stopping

A. Gronskiy and J. Buhmann

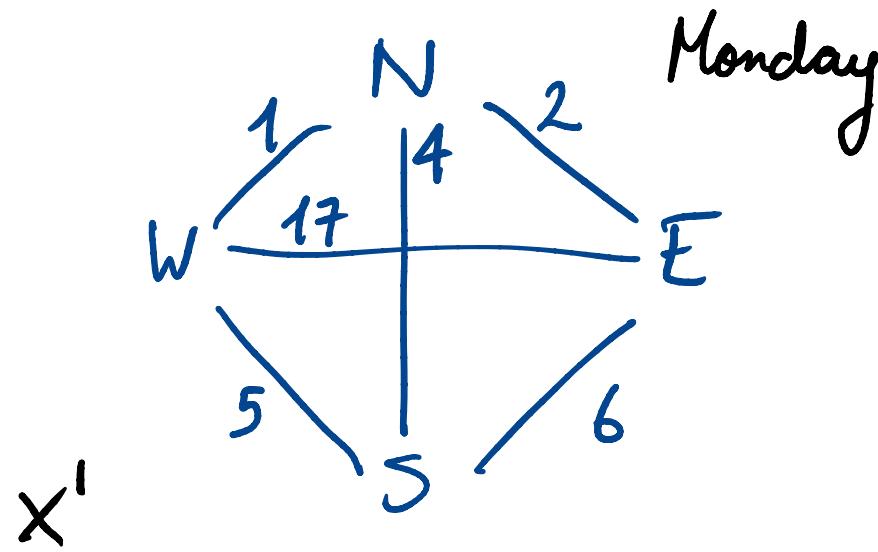
## Motivation



What's the minimum Spanning tree?

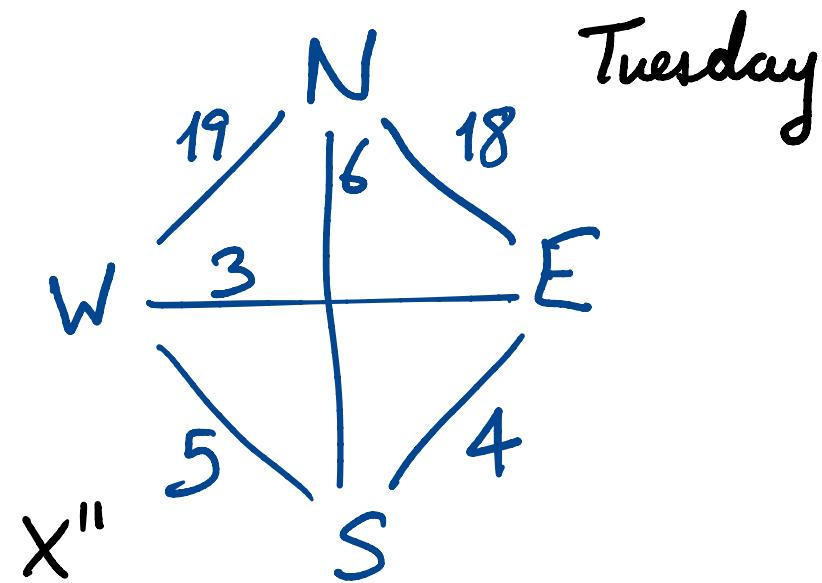
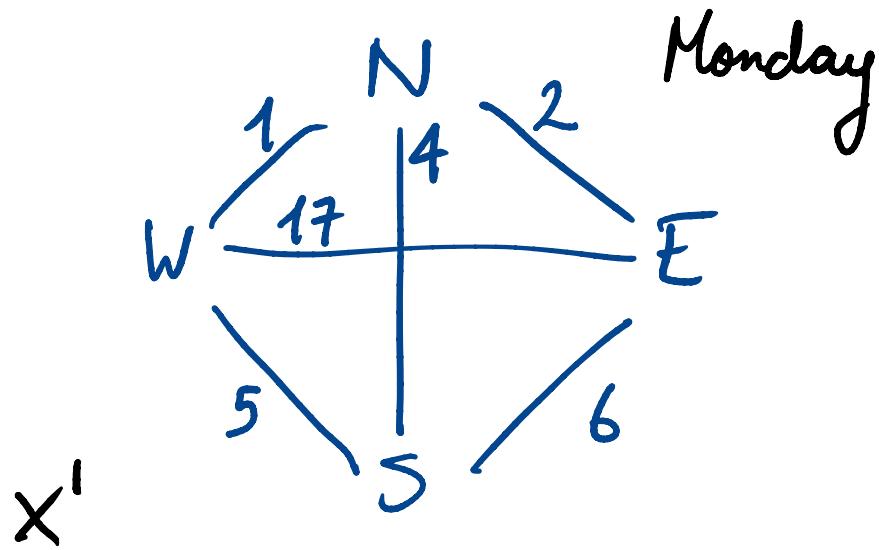


## Motivation

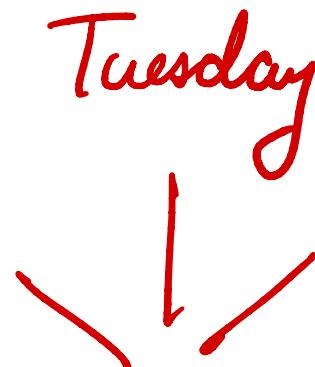
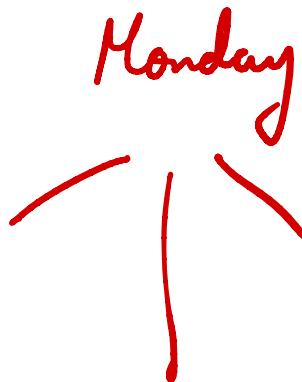


What's the minimum Spanning tree?

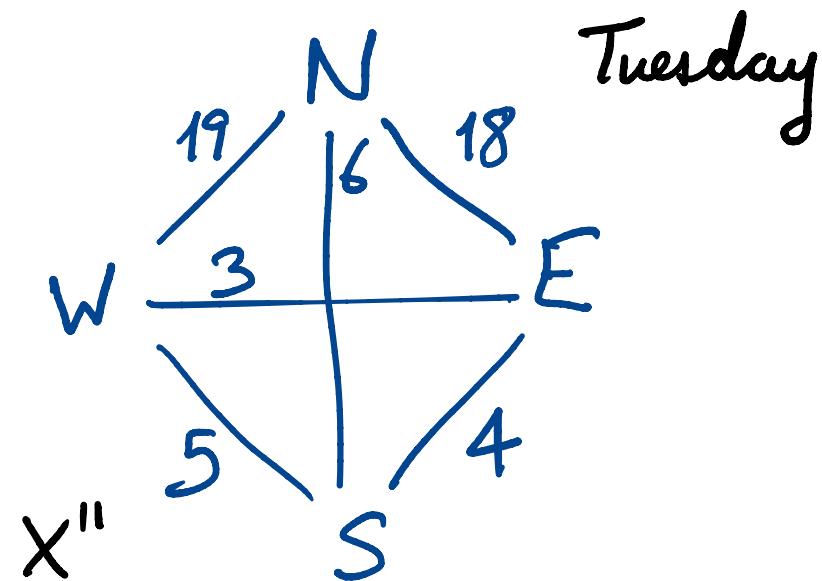
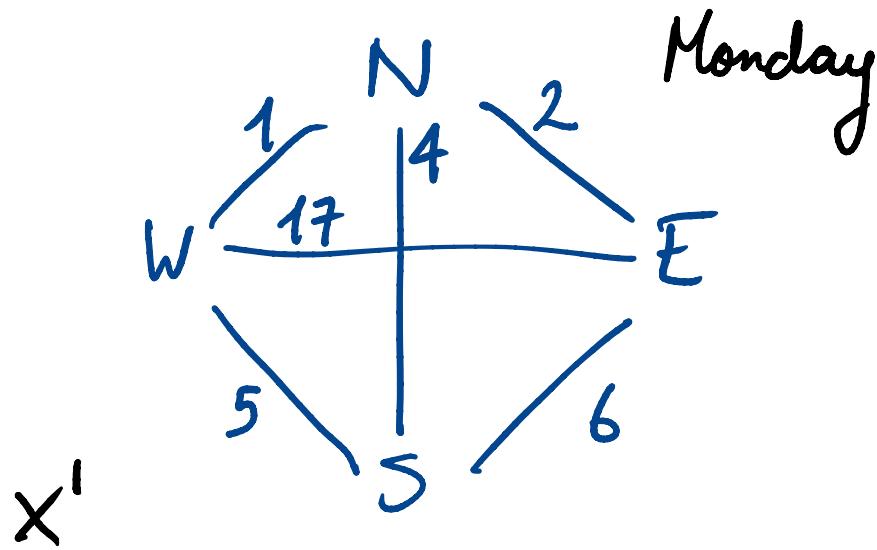
## Motivation



What's the minimum Spanning tree?



## Motivation



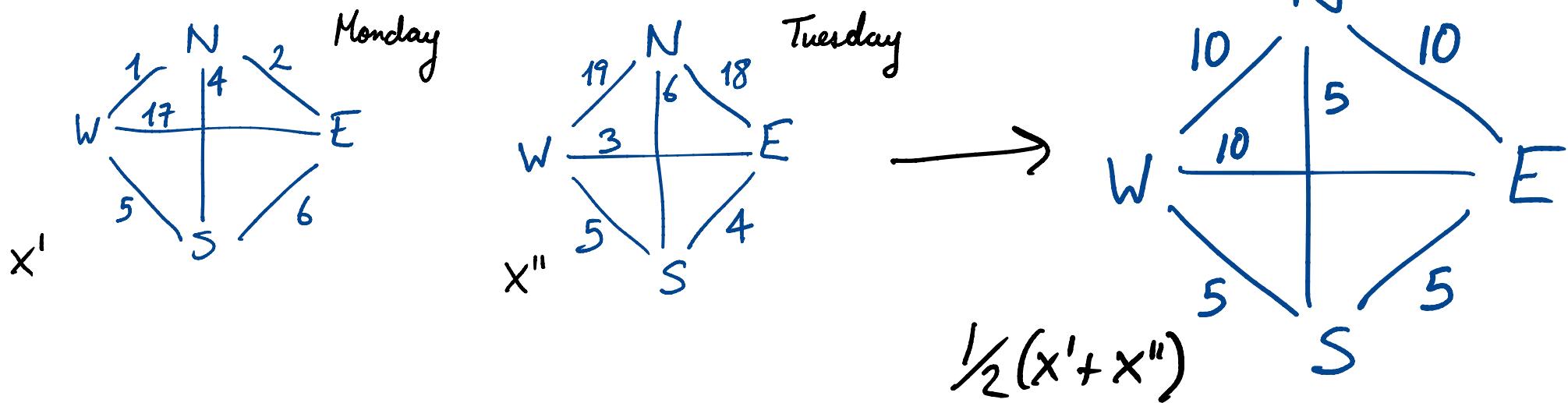
What's the minimum Spanning tree?

$$\min_{c \in C} \mathbb{E}_X [R(c, X)]$$

city graph  
c's weight in X  
all spanning trees

Natural choice (ERM):

$$\mathbb{E}_X [R(c, X)] \approx \frac{1}{2} (R(c, X') + R(c, X'')) \\ = R(c, \frac{1}{2}(X' + X''))$$



Alternative : ME + PA

Let  $p(\cdot | x') \propto \exp\left(-\frac{1}{T^*} R(c, x')\right)$

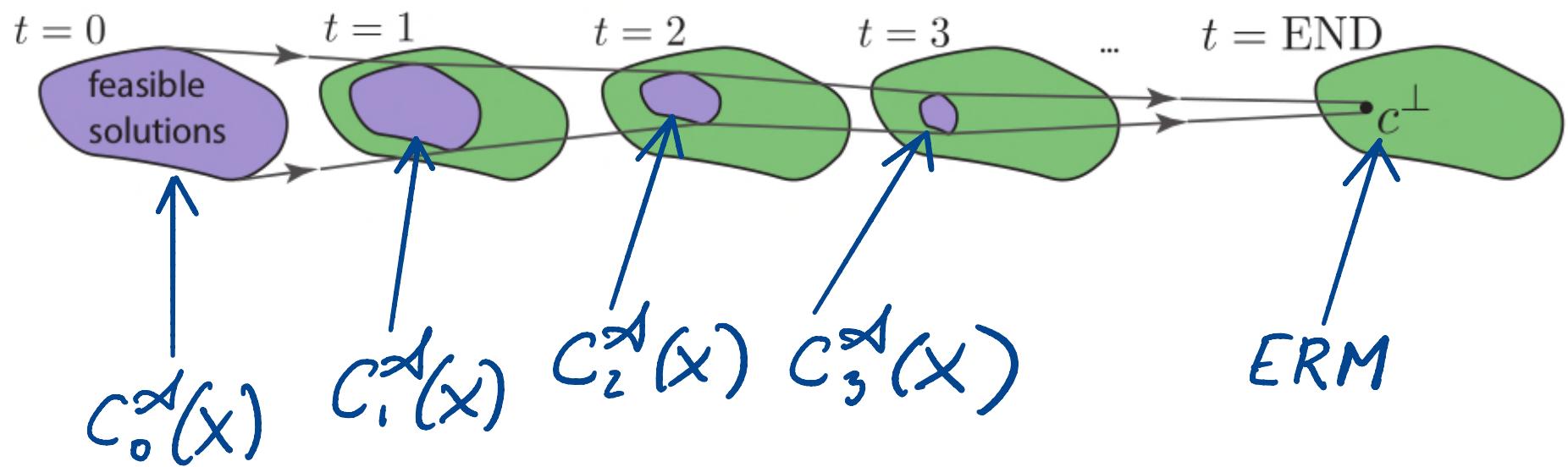
$$p(\cdot | x'') \propto \exp\left(-\frac{1}{T^*} R(c, x'')\right)$$

$$T^* = \arg \max_T K_T(x', x'')$$

Problem: computing  $p(\cdot | x')$  is  
intractable as there are  $n^{n-2}$  trees!

## Robust solving via early stopping + PA

Many MST algs are "contractive" and efficient.



Def: Algo  $\mathcal{A}$  is **contractive** if, given  $X$ ,  
its execution can be expressed as a  
sequence

$$G = C_0^{\mathcal{A}}(X) \supseteq C_1^{\mathcal{A}}(X) \supseteq \dots \supseteq C_T^{\mathcal{A}}(X) = \{C^\perp\}$$

$$G = C_0^{\mathcal{A}}(x) \supseteq C_1^{\mathcal{A}}(x) \supseteq \dots \supseteq C_T^{\mathcal{A}}(x) = \{C^\perp\}$$

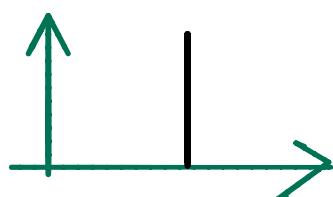
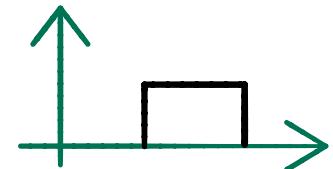
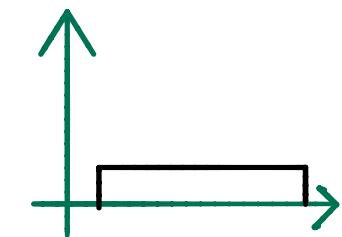
$$P_0^{\mathcal{A}}(\cdot | x), P_1^{\mathcal{A}}(\cdot | x), \dots, P_T^{\mathcal{A}}(\cdot | x)$$

At each step, a uniform distr over  $C_t^{\mathcal{A}}(x)$   
is computed

$$P_0^{\mathcal{A}}(\cdot | x) \leftarrow \text{unif}(G)$$

$$P_t^{\mathcal{A}}(\cdot | x) \leftarrow \text{unif}(G_t^{\mathcal{A}}(x))$$

$$P_T^{\mathcal{A}}(\cdot | x) \leftarrow \text{unif}(\{C^\perp\})$$



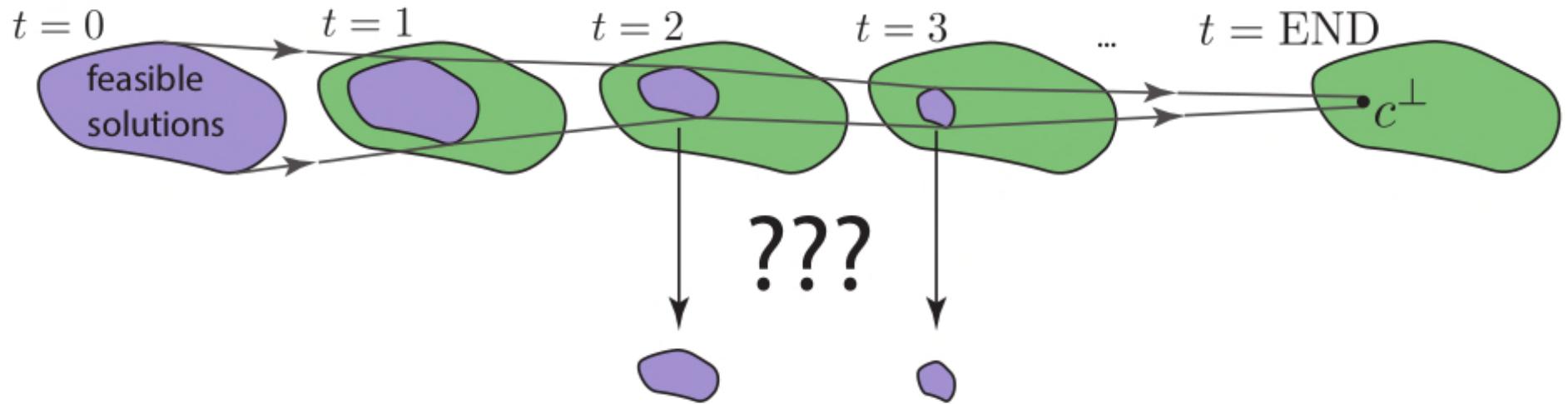
We can apply PA  
to build more robust MST  
algorithms

$$G = C_0^{\mathcal{A}}(x) \supseteq C_1^{\mathcal{A}}(x) \supseteq \dots \supseteq C_T^{\mathcal{A}}(x) = \{c^\perp\}$$

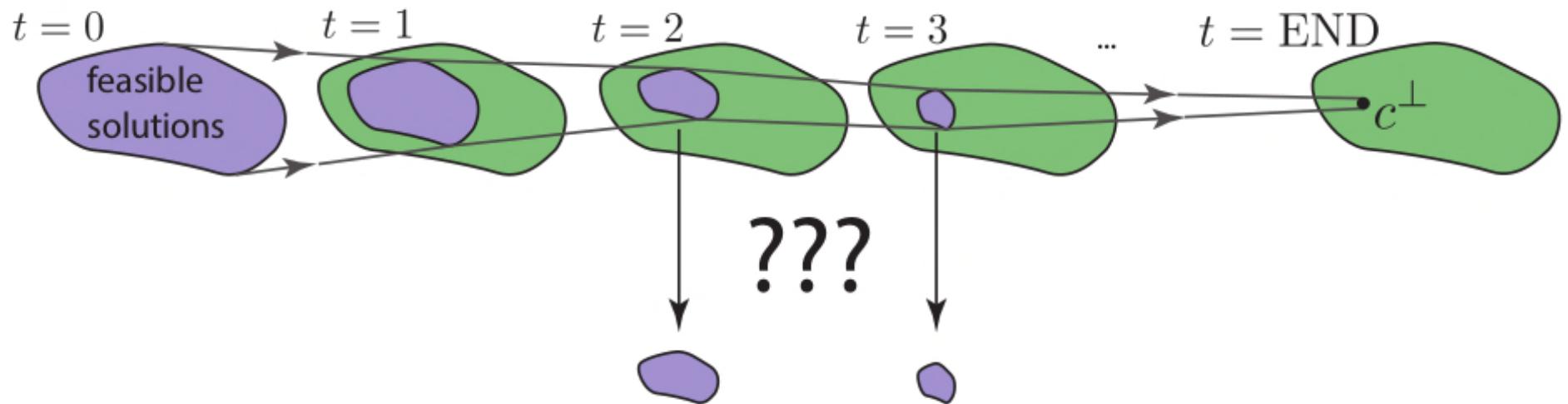
$$P_0^{\mathcal{A}}(\cdot | X), P_1^{\mathcal{A}}(\cdot | X), \dots, P_T^{\mathcal{A}}(\cdot | X)$$

↑  
STOP

Idea: Stop the algorithm earlier at  $t < T$  and output  $P_t^{\mathcal{A}}(\cdot | X)$ .



When do we stop?

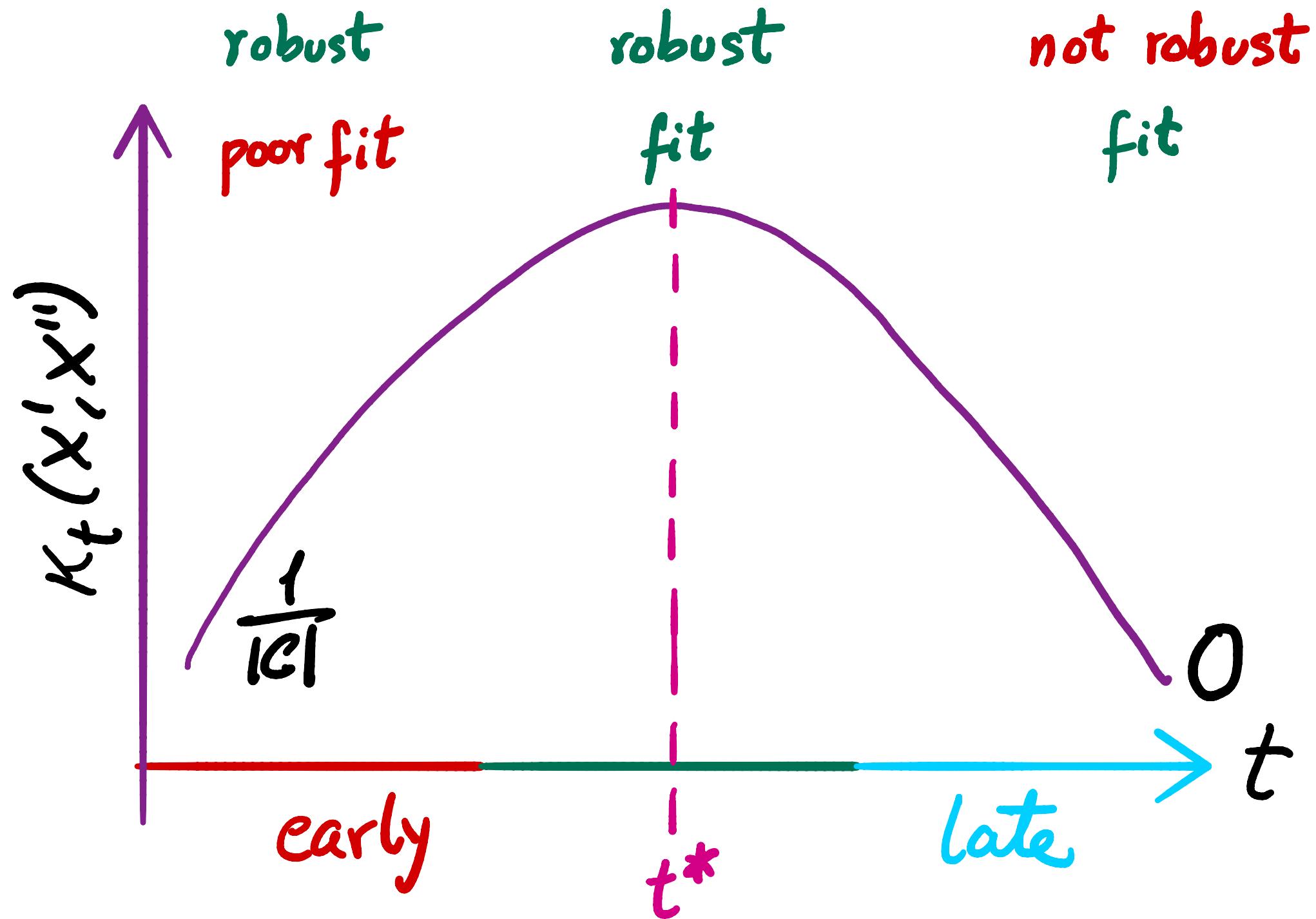


When do we stop?

$$t^* = \arg \max_t K_t^{\mathcal{A}}(x', x'')$$

$$\frac{|C_{t+1}|}{|C_t||C_t|} = \frac{1}{|C_t|}$$

$$K_t^{\mathcal{A}}(x', x'') = \sum_c p_t^{\mathcal{A}}(c|x') p_t^{\mathcal{A}}(c|x'') = \frac{|C_t^{\mathcal{A}}(x') \cap C_t^{\mathcal{A}}(x'')|}{|C_t^{\mathcal{A}}(x')||C_t^{\mathcal{A}}(x'')|}$$



# Regularized MST algo

Input:  $x'$ ,  $x''$ , MST algo  $\mathcal{A}$

1. Run  $\mathcal{A}$  to compute  $C_t^{\mathcal{A}}(x')$  and  $C_t^{\mathcal{A}}(x'')$ , for all  $t$ .

2.  $t^* \leftarrow \arg \max_t \frac{|C_t^{\mathcal{A}}(x') \cap C_t^{\mathcal{A}}(x'')|}{|C_t^{\mathcal{A}}(x')| |C_t^{\mathcal{A}}(x'')|}$

3.  $c^* \leftarrow \text{Unif}(C_{t^*}^{\mathcal{A}}(x') \cap C_{t^*}^{\mathcal{A}}(x''))$

Which MST algo is the best for this?

# Which reg MST algo is the best ?

MST algos : Prim, Kruskal, reverse-delete.

- \* Experimental assessment of PA.
- \* Experimental comparison.  
PA rankings match with performance rankings.
- \* Analysis

Experiments: PA of the reg. MST algos

\* For  $\alpha \in \{0, \dots, 80\}$ , a "random" complete graph with 60 vertices and the exp. log. PA for each algo was calculated.

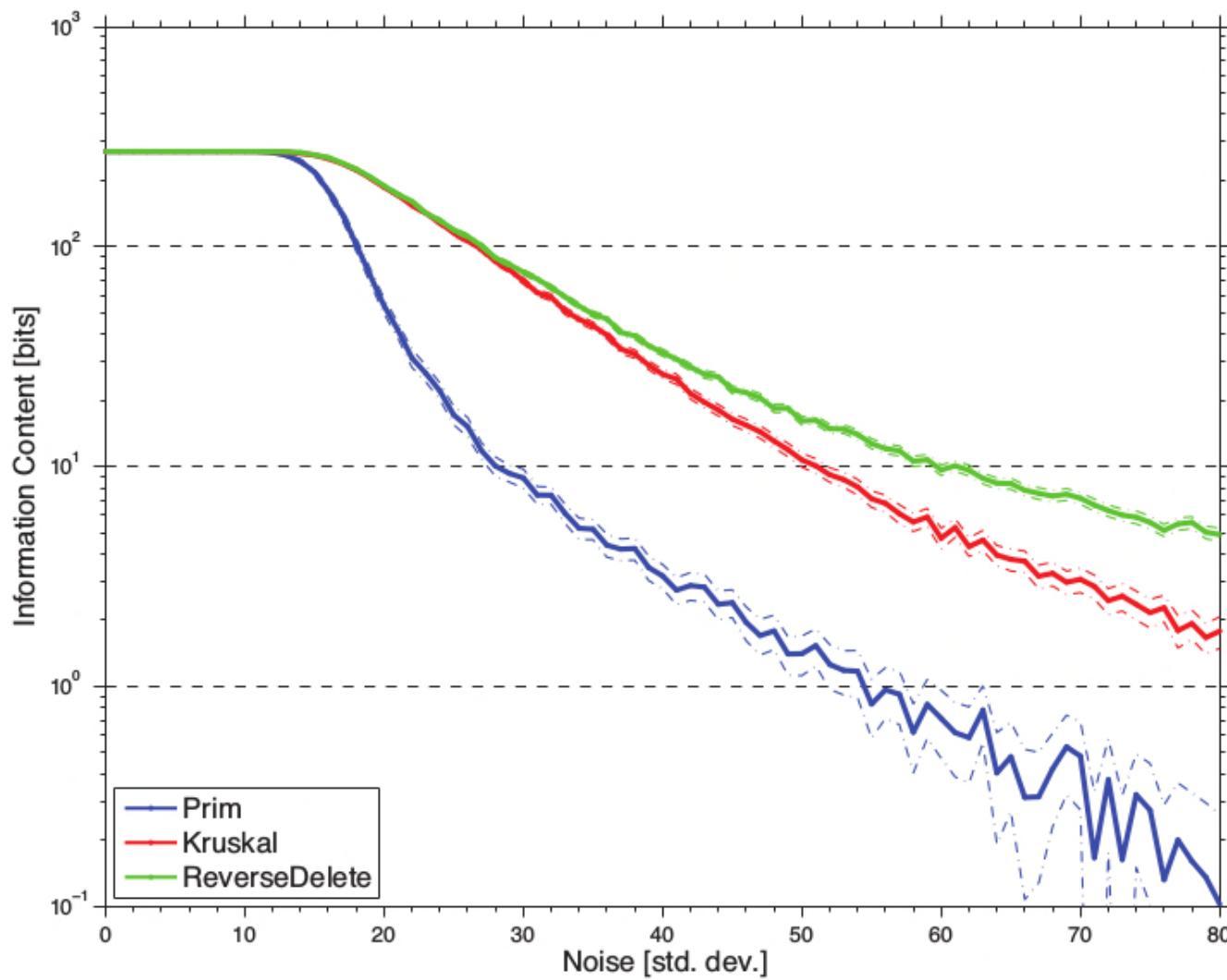
experiment  $\rightarrow G = (V, V^2)$   $V = \{1, \dots, 60\}$

for  $v, w \in V$   $\text{weight}(v, w) \stackrel{\$}{\leftarrow} N(100, 100)$ .

instance  $\rightarrow X = (V, V^2)$

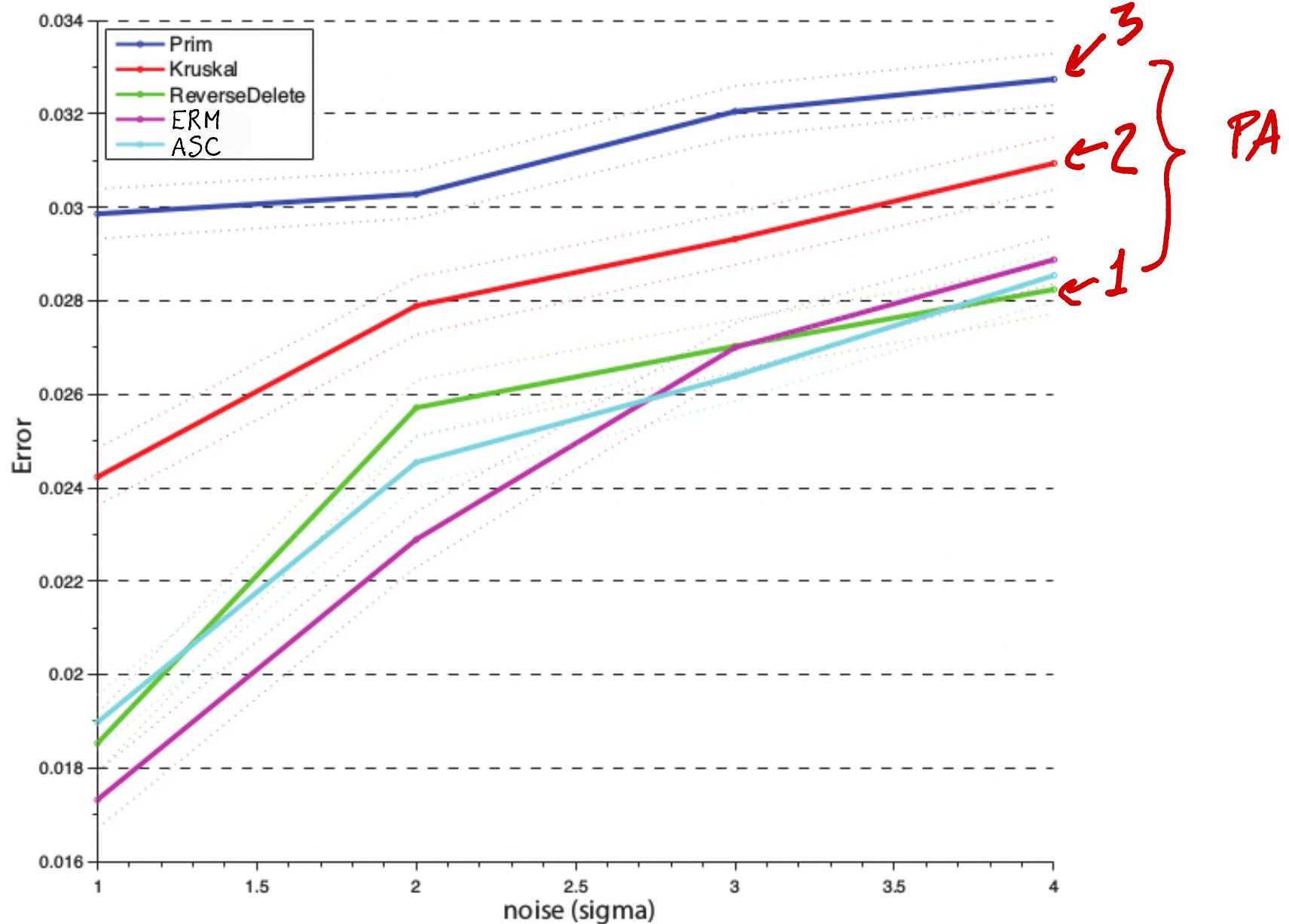
for  $v, w \in V$   $\text{weight}(v, w) \sim N(\text{weight}(u, v), \sigma^2)$

*Exp. log. PA*



Experiments: Performance of the reg. MST algs.

- 400 complete "ground-truth" graphs with 6 vertices
- Each edge weight  $\xleftarrow{\$} \mathcal{N}(\mu, \sigma^2)$
- From each graph, two instances  $X'$  and  $X''$  were generated by adding noise  $\mathcal{N}(0, \sigma_0^2)$  to each edge.



$$\text{Error} = \frac{1}{400} \sum_{i \leq 400} |t_i^* \Delta t_i|$$

↑ Symmetric diff.

$t_i^*$ : MST  
 $t_i$ : ST output by  $\delta$

# MST algorithms and their regularized versions

Prim's algo (growing tree)

- pick an arbitrary vertex  $v_0$ .
- $T \leftarrow \{ (v_0, \arg \min_{v'} \text{weight}(v_0, v')) \} \}$
- while  $T$  not spanning

Expand  $T$  with

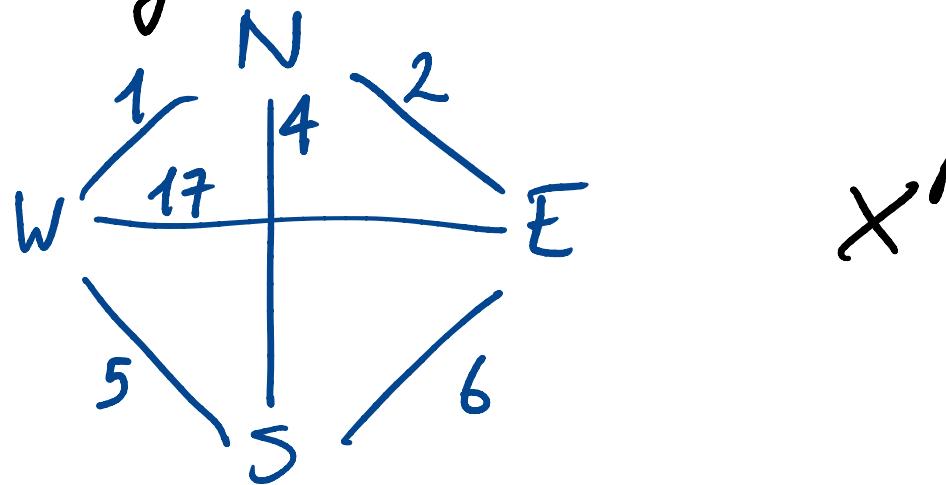
$$\arg \min_{(u,v)} \text{weight}(u, v)$$

s.t.  $u \in T, u \notin T$

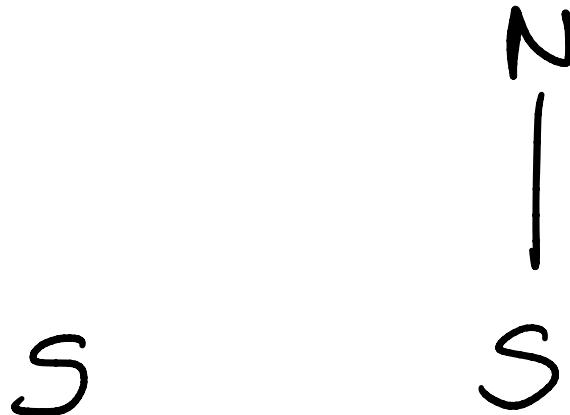
$T \cup \{(u, v)\}$  is a tree.

# MST algorithms and their regularized versions

Prim's algo (growing tree)



$t=0$



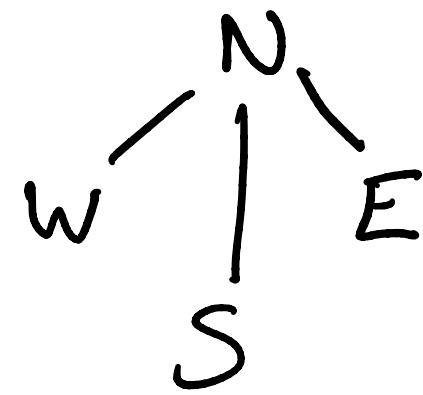
$t=1$



$t=2$



$t=3$

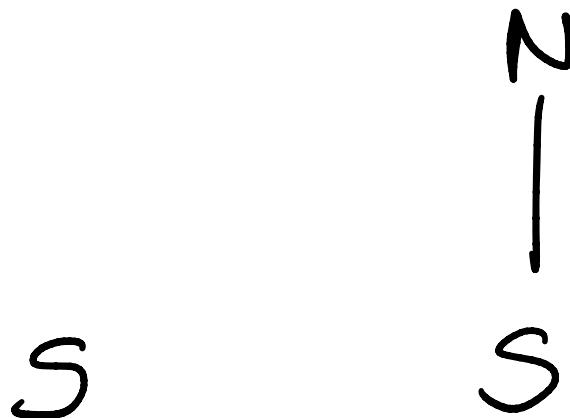


# MST algorithms and their regularized versions

Prim's algo (growing tree)

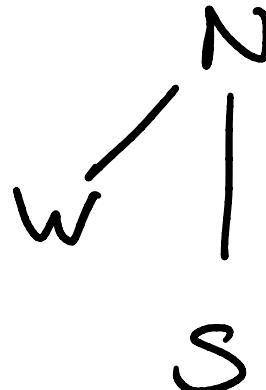
$$\left\{ \text{All trees} \right\} \supseteq \left\{ \begin{matrix} C_0(x') \\ \uparrow \\ 4 \end{matrix} \right\} \supseteq \left\{ \begin{matrix} C_1(x') \\ \uparrow \vee 1 \\ 4 \end{matrix} \right\} \supseteq \left\{ \begin{matrix} C_2(x') \\ \uparrow \\ 4 \end{matrix} \right\} \supseteq \dots \supseteq \left\{ \begin{matrix} C_3(x') \\ \uparrow \\ 1 \end{matrix} \right\}$$

$t=0$

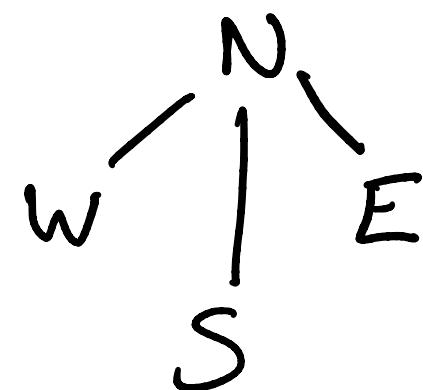


$t=1$

$t=2$

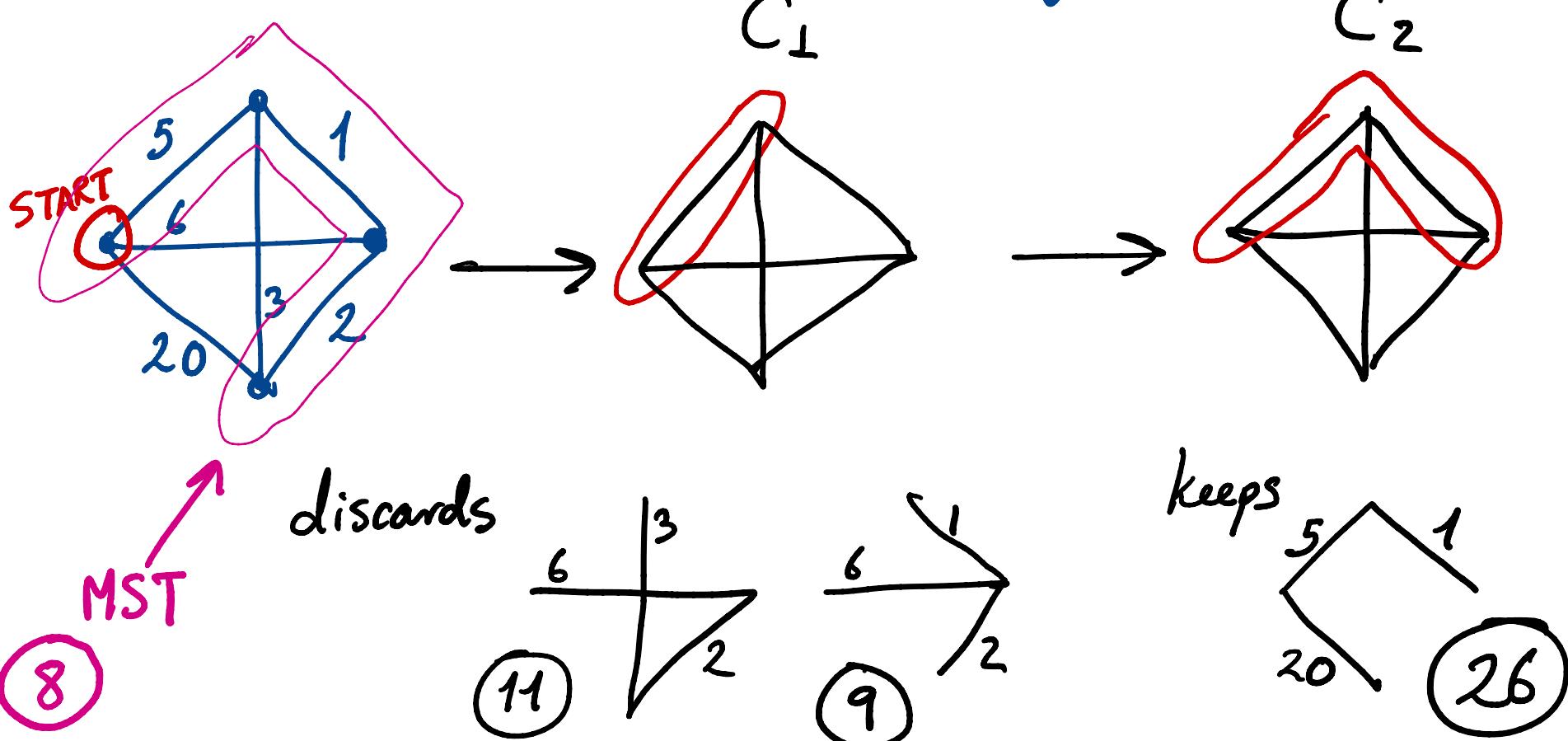


$t=3$



## Prim:

- \* Only adds edges connected to the current tree
- \* If you start wrong, then you discard good candidates too early.



# MST algorithms and their regularized versions

Kruskal's algo (growing forest)

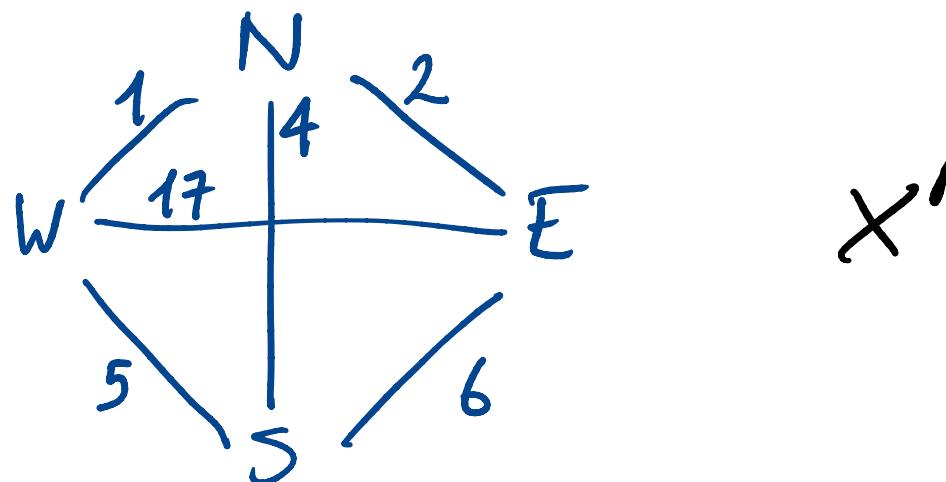
1.  $T \leftarrow \emptyset$

2. While  $T$  not spanning

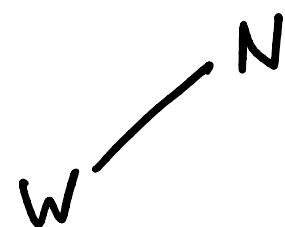
Expand  $T$  with the lightest edge  
that does not induce a cycle

# MST algorithms and their regularized versions

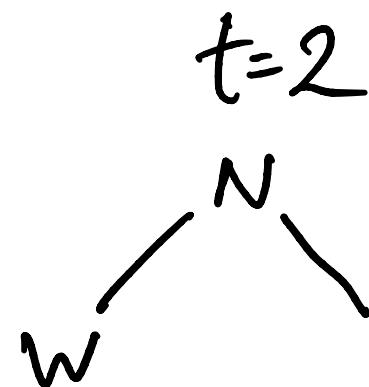
Kruskal's algo (growing forest)



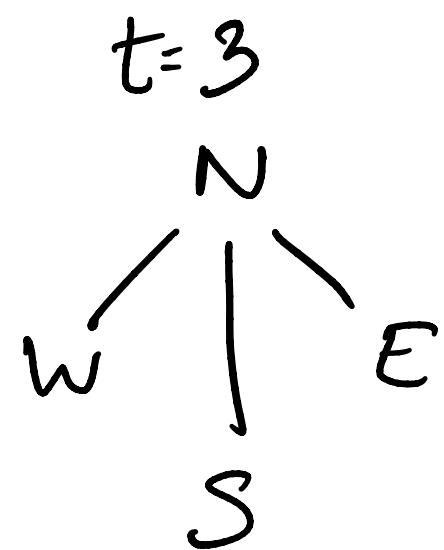
$t=0$



$t=1$



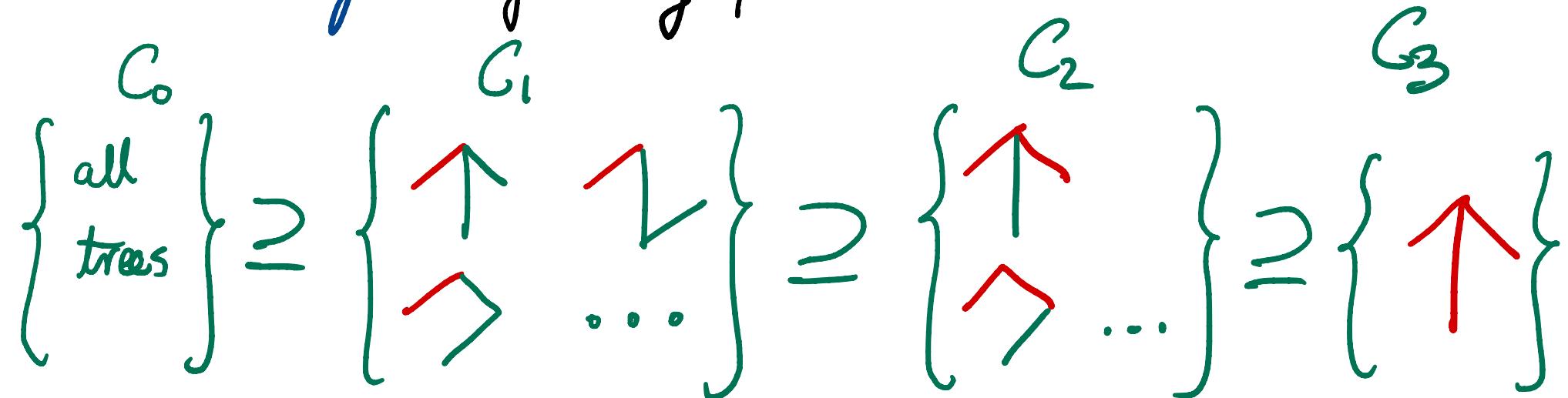
$t=2$



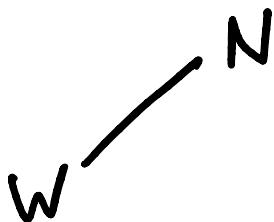
$t=3$

# MST algorithms and their regularized versions

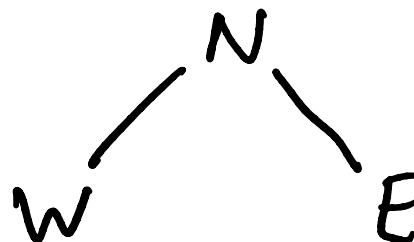
Kruskal's algo (growing forest)



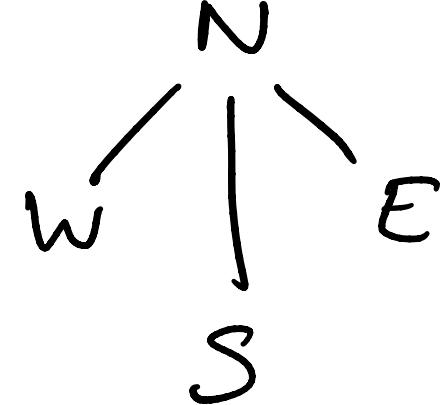
$t=0$



$t=1$



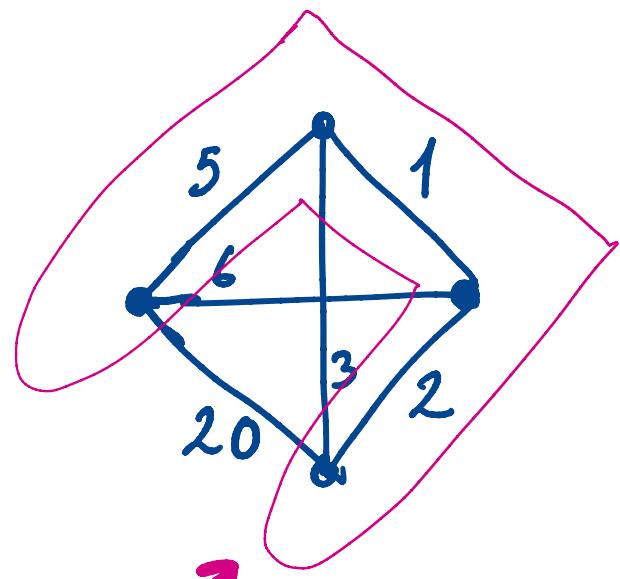
$t=2$



$t=3$

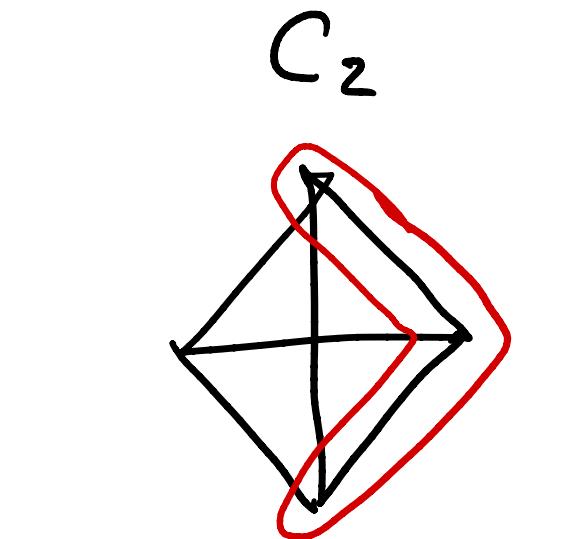
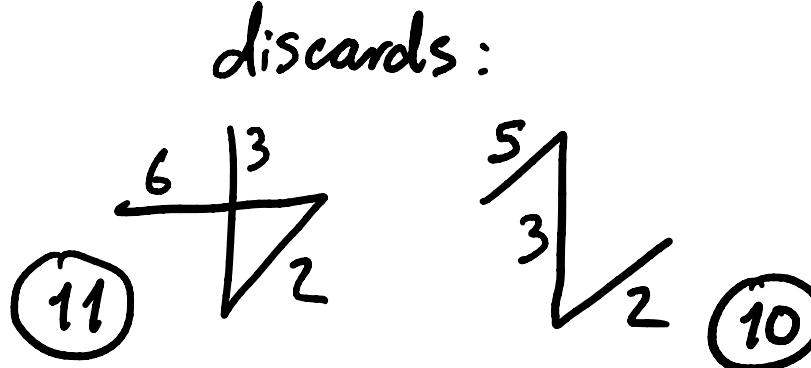
## Kruskal:

- \* Does not require selected edges to be connected.
- \* If you start wrong, you discard good candidates.

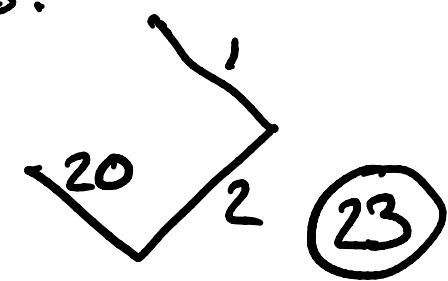


MST

⑧



keeps:



# MST algorithms and their regularized versions

## Reverse-delete (reducing graph)

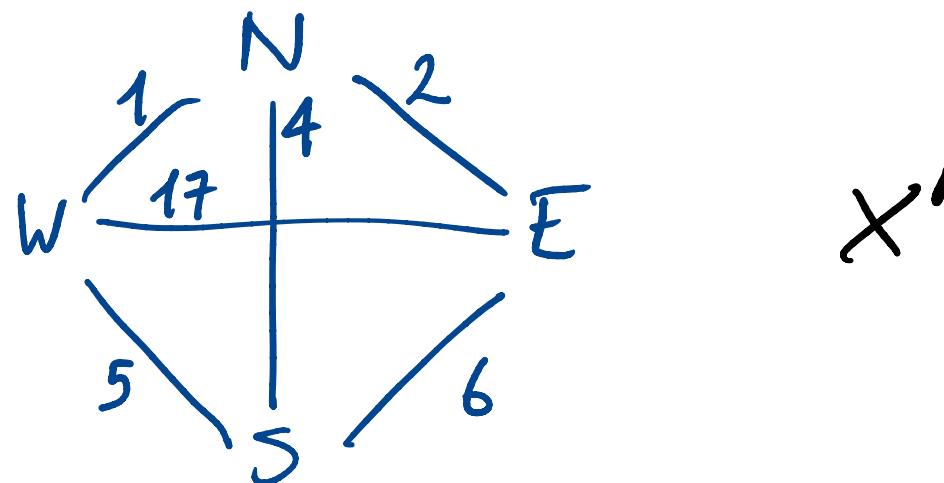
1.  $T \leftarrow$  all edges

2. While  $T$  is not a tree:

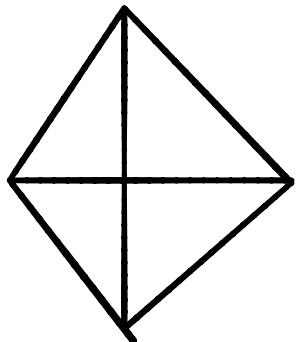
remove the heaviest edge that  
does not disconnect the graph

# MST algorithms and their regularized versions

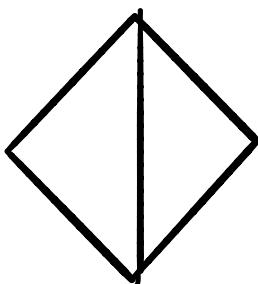
Reverse-delete (reducing graph)



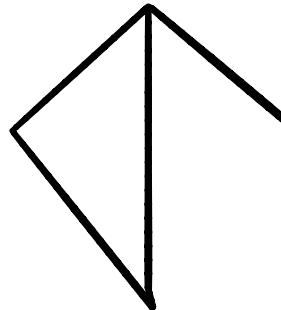
$t=0$



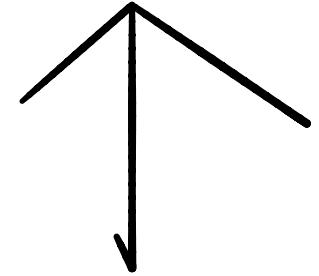
$t=1$



$t=2$

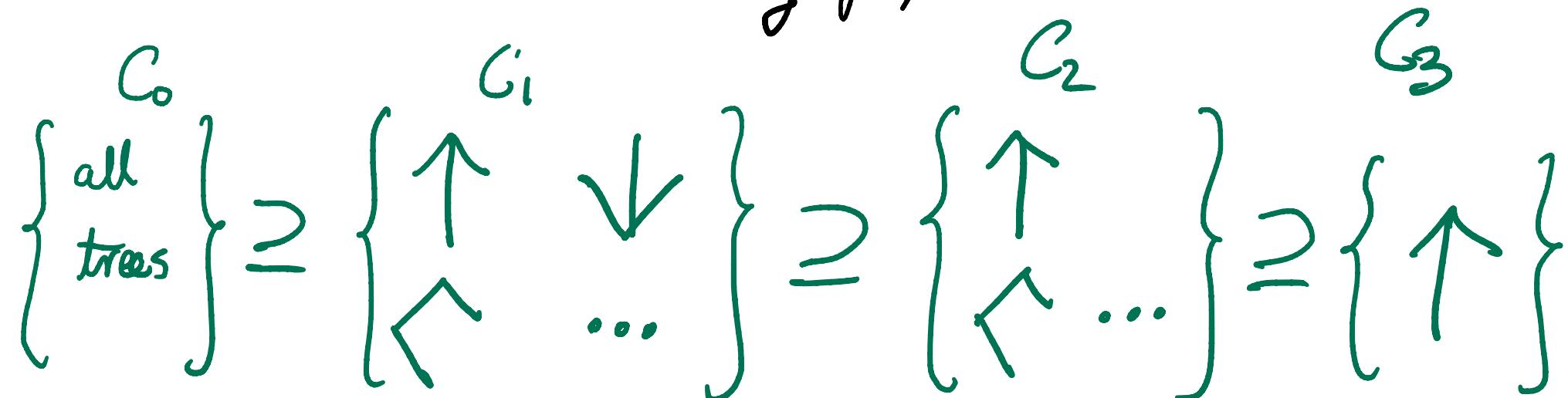


$t=3$

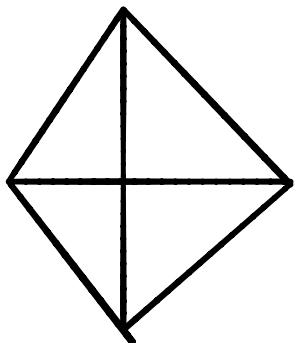


# MST algorithms and their regularized versions

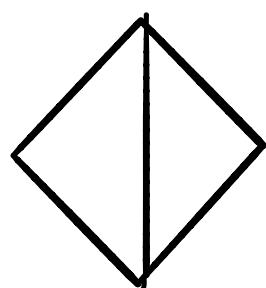
Reverse-delete (reducing graph)



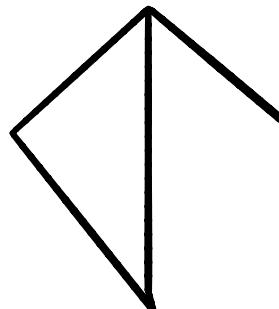
$t=0$



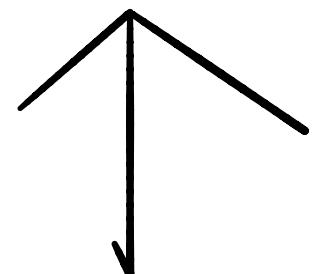
$t=1$



$t=2$

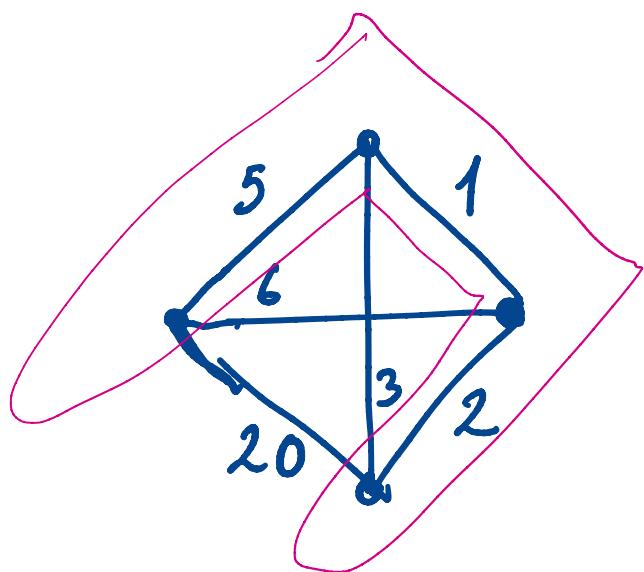


$t=3$

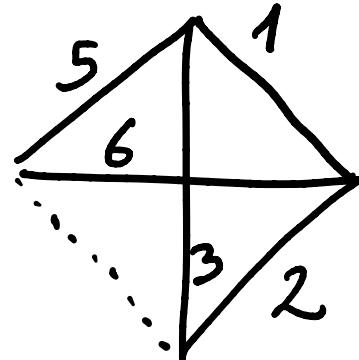


## Reverse-delete

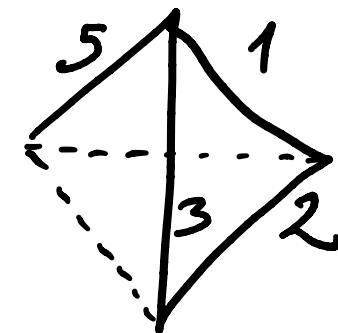
- \* Discards the heaviest edge, while keeping the graph connected
- \* Discards bad candidates early



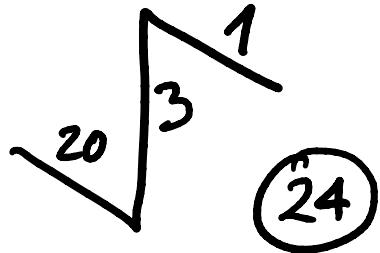
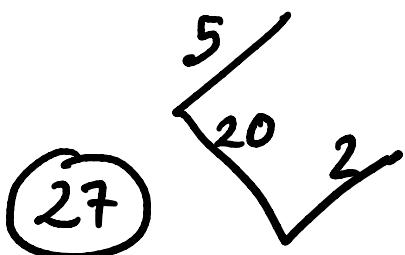
$C_1$



$C_2$



discards



## Conclusion

- \* PA correctly ranked the performance of these algorithms.
- \* PA can guide the choice of hyperparameter values, cost functions, and algorithms.

# Approximate sorting

L. Busse, M. Haghiri Chahregani, J. Buhmann

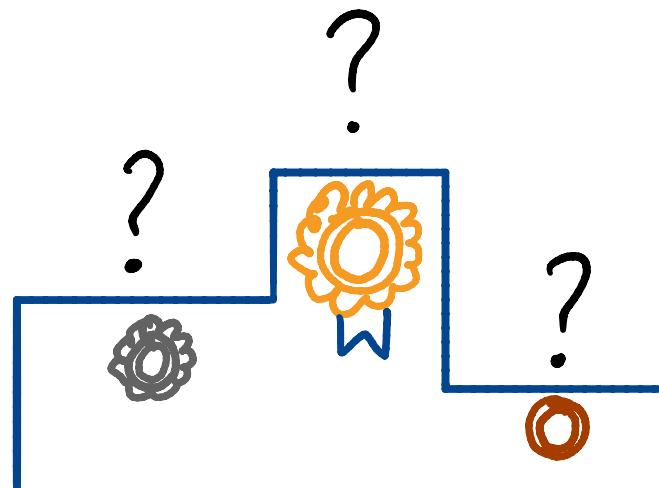
# Problem intro.

	A	B	C
Alice		W	W
Bob	L		W
Charlie	L	L	

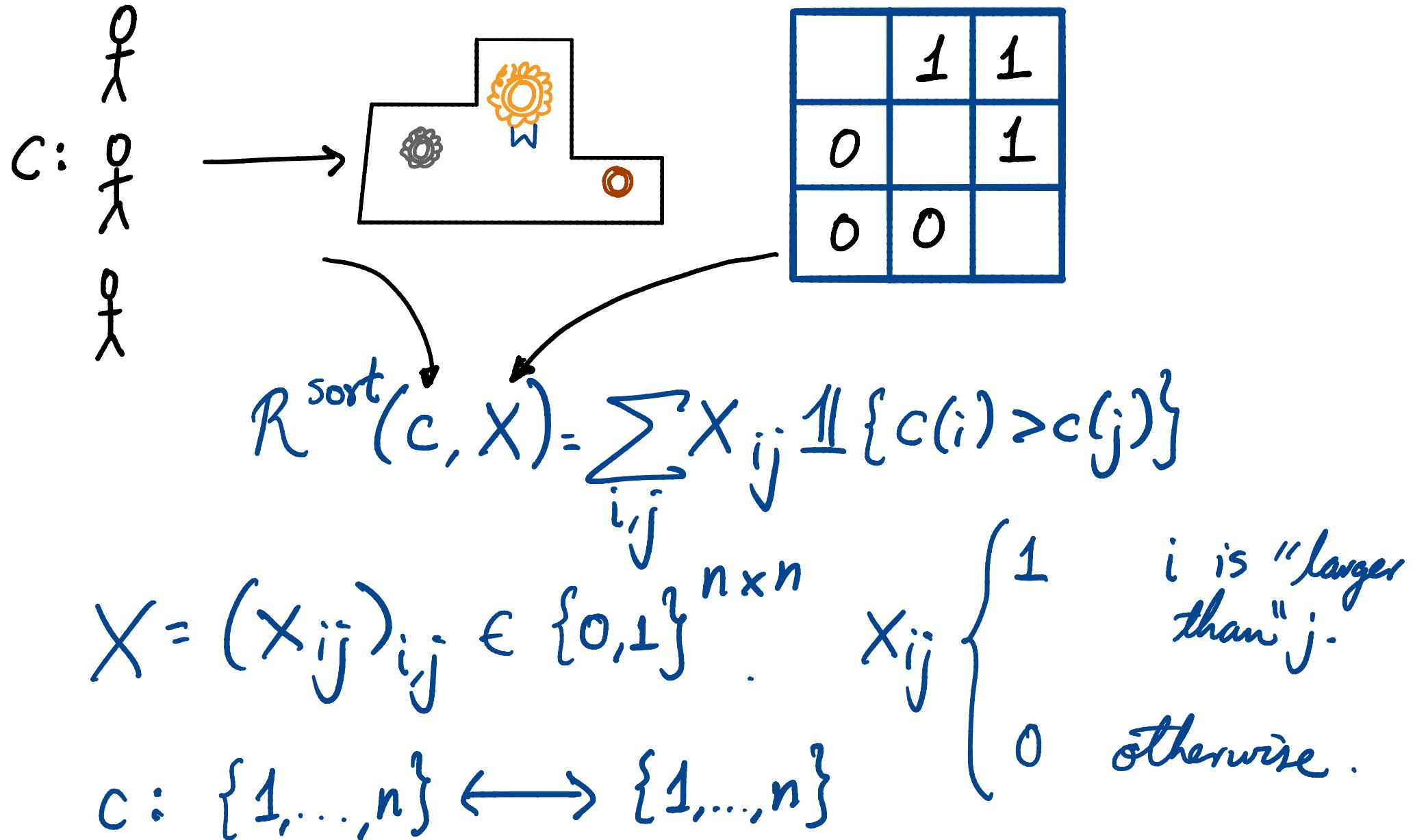
Fall 2019

	A	B	C
Alice		W	W
Bob	L		L
Charlie	L	W	

Spring 2020



# Problem formalization



$$\frac{\text{ME+PA}}{p^{\text{sort}}(c|X) \propto \exp\left(-\frac{1}{T^*} R^{\text{sort}}(c, X)\right)}.$$

$$T^* \leftarrow \arg \max_T K_T(X', X'')$$

Problem:  $p^{\text{sort}}(c|X)$ 's normalization constant  
is intractable.  
( $n!$  different solutions)

Approx. sorting

$$R^{\text{sort}}(c, x) = \sum_{i,j} x_{ij} \mathbb{1}\{c(i) > c(j)\}.$$

$$R^{\text{mf}}(M, \epsilon) = \sum_i \sum_k M_{ik} \epsilon_{ik}.$$

$$R^{\text{mfs}}(M, \epsilon) = \sum_i \sum_k M_{ik} \epsilon_{ik} + \sum_k \mu_k \left( \sum_i M_{ik} - 1 \right)$$

$$= \sum_{i,k} M_{ik} \epsilon_{ik} + \sum_{i,k} M_{ik} \mu_k - \sum_{i,k} \frac{\mu_k}{n}$$

$$= \text{const} + \sum_{i,k} M_{ik} (\epsilon_{ik} + \mu_k)$$

What's the optimal  $\tau$ ?

$$\tau^* = \arg \max_{\tau} K_{\tau}(x', x'')$$

⋮

Tractable!  
 $\frac{1}{C}$

$$= 1 + \frac{1}{\log |G|} \sum_i \log \sum_k \exp\left(-\frac{1}{T}(\varepsilon'_{ik} + \mu'_k + \varepsilon''_{ik} + \mu''_k)\right) \\ - \frac{1}{\log |C|} \sum_i \log \left( \sum_k \exp\left(-\frac{1}{T}(\varepsilon'_{ik} + \mu'_k)\right) \sum_k \exp\left(-\frac{1}{T}(\varepsilon''_{ik} + \mu''_k)\right) \right).$$

## Experiments

Synthetic experiments

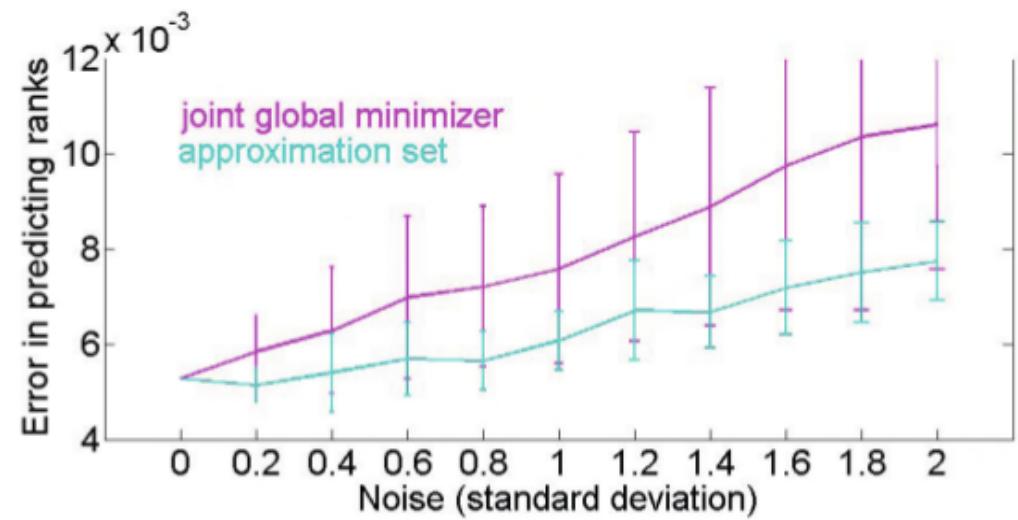
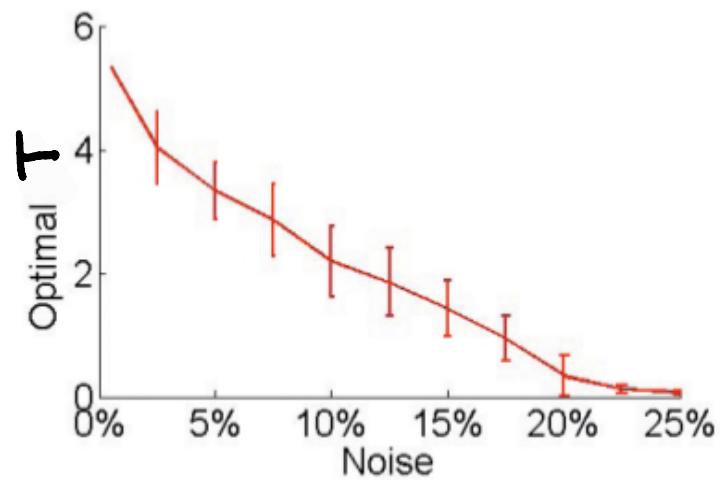
A noisy ranking  $\tilde{\pi}$  was generated

$$\tilde{\pi}_i \sim \mathcal{N}(i, \sigma^2) \quad \text{noise parameter}$$

Then  $X$  was defined by

$$X_{ij} := \mathbb{1}\{\tilde{\pi}_i < \tilde{\pi}_j\}.$$

# Experimental results



## Experiments with real-world data

Tournaments from the German chess league  
competitions ( $n = 16$  teams)

records from wins and losses

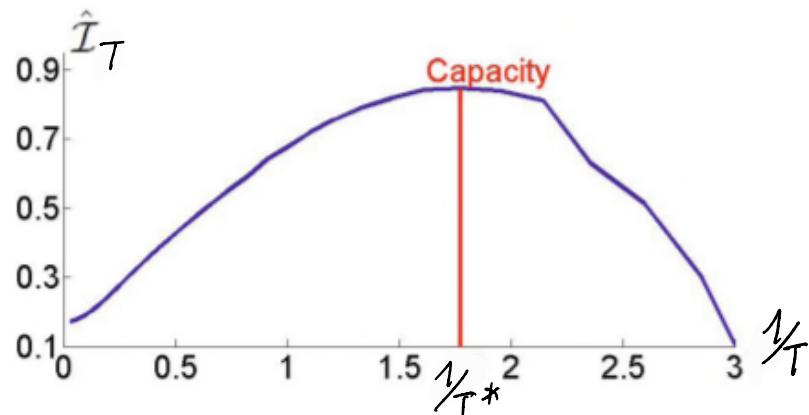
$$X' \leftarrow 2009-2010$$

$$X'' \leftarrow 2010-2011$$

Approach:  $c \leftarrow \underset{\$}{\alpha} \langle p_{T^*}(\cdot | X'), p_{T^*}(\cdot | X'') \rangle$

Baseline:  $\arg \min_c R^{\text{sort}}(c, X') + R^{\text{sort}}(c, X'')$

# Results



Method	Prediction error per rank		
	abs.	rel.	gain
Baseline	0.29	14%	
ME + PA	0.14	6.8%	51%

(b) Prediction of chess league results

# What we learned

- \* PA as a method for validating algorithms.
- \* Algorithms as communication channels.
- \* Asymptotic equipartition property and typicality as tools for simplifying analysis of complex systems.
- \* How to use information theory to develop effective algorithms for robust stochastic optimization.