# Learning and Testing Submodular Functions

## Grigory Yaroslavtsev

**With Sofya Raskhodnikova (SODA'13)**

PENNSTATE
1855

## Weizmann Institute

December 30, 2012

# Submodularity

- Discrete analog of convexity/concavity, "law of diminishing returns"
- Applications: combintorial optimization, AGT, etc.

Let $f: 2^X \to [0, R]$:

- **Discrete derivative:**

$$\partial_x f(S) = f(S \cup \{x\}) - f(S), \qquad for \ S \subseteq X, x \notin S$$
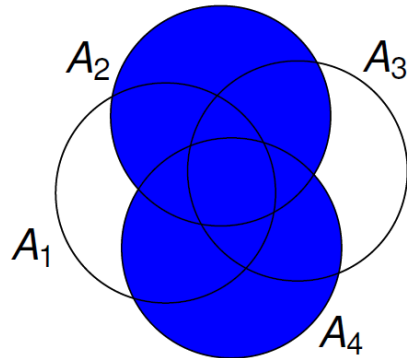
- **Submodular function:**

$$\partial_x f(S) \geq \partial_x f(T), \quad \forall \ S \subseteq T \subseteq X, x \notin T$$
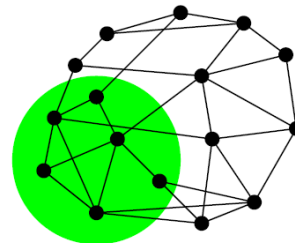
**Coverage function:**
Given $A_1, \ldots, A_n \subset U$,
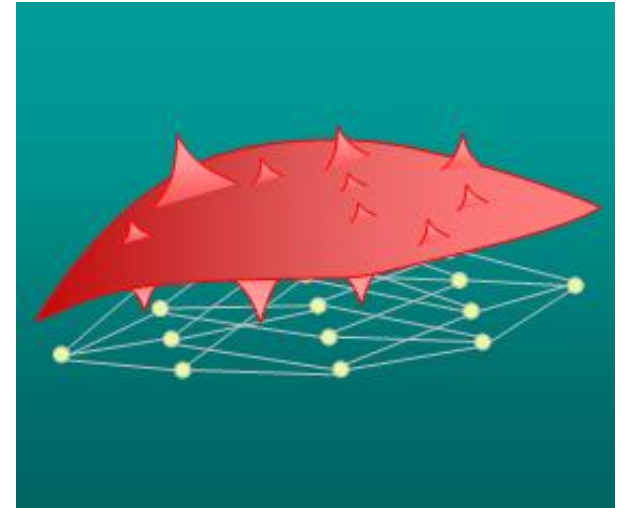
$$f(S) = \left| \bigcup_{j \in S} A_j \right|.$$

$A_2$ $A_3$

$A_1$

$A_4$

**Cut function:**

$$\delta(T) = |e(T, \overline{T})|$$

# Approximating everywhere

- **Q1:** Approximate a submodular $f: 2^X \to [0, R]$ **for all arguments** with only poly($|X|$) queries?

- **A1:** Only $\widetilde{\Theta}\left(\sqrt{|X|}\right)$-approximation (multiplicative) possible [Goemans, Harvey, Iwata, Mirrokni, SODA'09]



- **Q2:** Only for $(1 - \epsilon)$-fraction of arguments (PAC-style learning **with membership queries** under uniform distribution)?

$$\Pr_{randomness\ of\ \boldsymbol{A}}\left[\Pr_{\boldsymbol{S} \sim U(2^X)}[\boldsymbol{A}(\boldsymbol{S}) = \boldsymbol{f}(\boldsymbol{S})] \geq 1 - \epsilon\right] \geq \frac{1}{2}$$

- **A2:** Almost as hard [Balcan, Harvey, STOC'11].

# Approximate learning

- **P**<span style="color:red">**M**</span>**AC**-learning (<span style="color:red">**M**</span>ultiplicative), with poly(|X|) queries :

$$\Pr_{rand.\,of\,A}\left[\Pr_{S\sim U(2^X)}[\,f(S)\leq A(S)\leq \alpha f(S)]\geq 1-\epsilon\right]\geq \frac{1}{2}$$

$$\Omega\left(|X|^{\frac{1}{3}}\right)\leq \alpha \leq O\left(\sqrt{|X|}\right)\ \text{[Balcan, Harvey '11]}$$

- **P**<span style="color:red">**A**</span>**AC**-learning (<span style="color:red">**A**</span>dditive)

$$\Pr_{rand.\,of\,A}\left[\Pr_{S\sim U(2^X)}[\,|f(S)-A(S)|\leq \beta]\geq 1-\epsilon\right]\geq \frac{1}{2}$$

  - Running time: $|X|^{O\left(\frac{R}{\beta}\right)^2\log(\frac{1}{\epsilon})}$ [Gupta, Hardt, Roth, Ullman, STOC'11]

  - Running time: $\text{poly}\left(|X|^{\left(\frac{R}{\beta}\right)^2},\ \log\frac{1}{\epsilon}\right)$ [Cheraghchi, Klivans, Kothari, Lee, SODA'12]

# Learning $f: 2^X \to [0, R]$

- For all algorithms $\epsilon = const.$

| | Goemans, Harvey, Iwata, Mirrokni | Balcan, Harvey | Gupta, Hardt, Roth, Ullman | Cheraghchi, Klivans, Kothari, Lee | Raskhodnikova, Y. |
|---|---|---|---|---|---|
| Learning | $\tilde{O}\left(\sqrt{\lvert X \rvert}\right)$- approximation Everywhere | PMAC Multiplicative $\alpha$ $\alpha = O\left(\sqrt{\lvert X \rvert}\right)$ | PAAC Additive $\beta$ | | PAC $f: 2^X \to \{0, \dots, R\}$ (bounded integral range $R \le \lvert X \rvert$) |
| Time | Poly($\lvert X \rvert$) | Poly($\lvert X \rvert$) | $\lvert X \rvert^{O\left(\frac{R}{\beta}\right)^2}$ | | $\lvert X \rvert^3 R^{O(R \cdot \log R)}$ Polylog($\lvert X \rvert$) $R^{O(R \cdot \log R)}$ queries |
| Extra features | | Under arbitrary distribution | Tolerant queries | SQ-queries, Agnostic | |

# Learning: Bigger picture

Subaddittive

$\cup$I

XOS = Fractionally subadditive

$\cup$I

**Submodular**

$\cup$I

Gross substitutes

$\cup$I
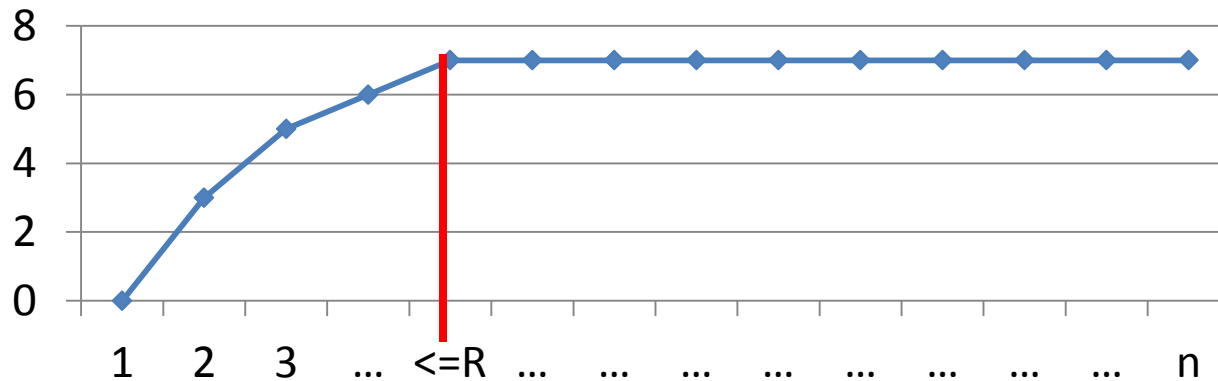
OXS

Additive           Coverage (valuations)
(linear)

} [Badanidiyuru, Dobzinski, Fu, Kleinberg, Nisan, Roughgarden,SODA'12]
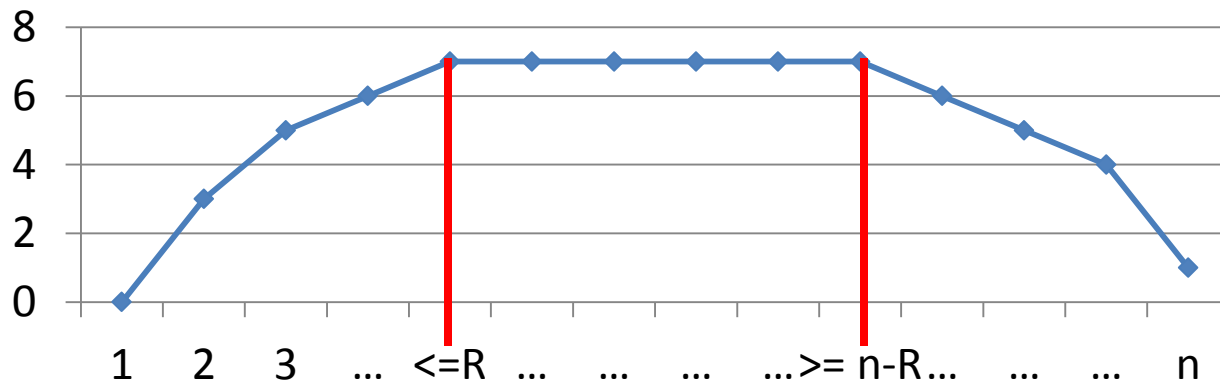
Other positive results:
- Learning valuation functions [Balcan, Constantin, Iwata, Wang, COLT'12]
- $(1 + \epsilon)$ PMAC-learning (sketching) coverage functions [BDFKNR'12]
- $(1 + \epsilon)$ PMAC learning Lipschitz submodular functions [BH'10] (concentration around average via Talagrand)

# Discrete convexity

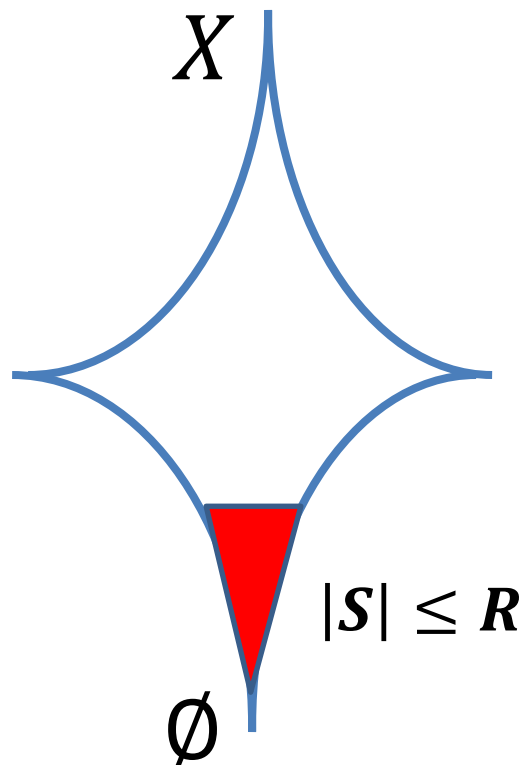- Monotone convex $f : \{1, \dots, n\} \to \{0, \dots, R\}$
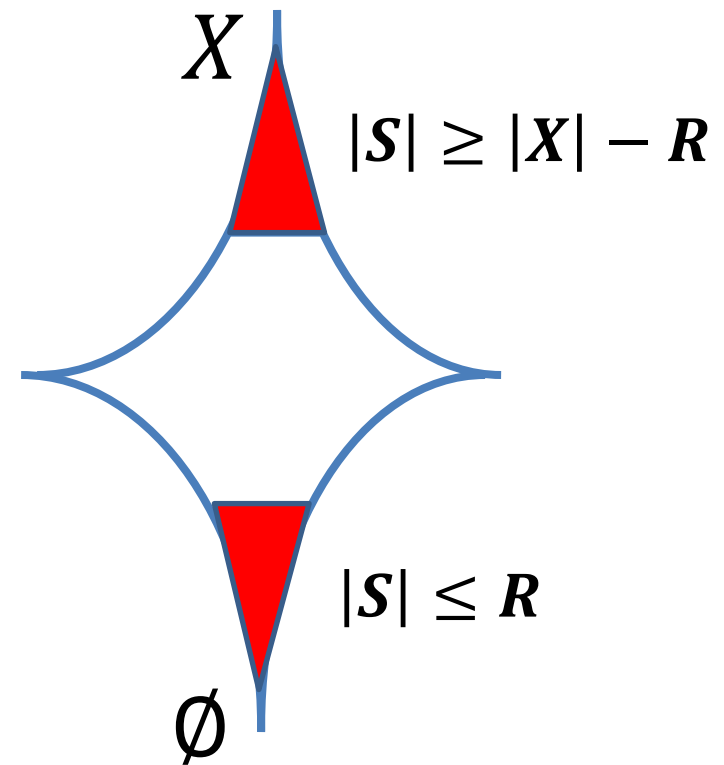


- Convex $f : \{1, \dots, n\} \to \{0, \dots, R\}$

# Discrete submodularity $f: 2^X \to \{0, \ldots, R\}$

- **Case study**: $R = 1$ (Boolean submodular functions $f: \{0,1\}^n \to \{0,1\}$)
  Monotone submodular = $x_{i_1} \vee x_{i_2} \vee \cdots \vee x_{i_a}$ (monomial)
  Submodular = $(x_{i_1} \vee \cdots \vee x_{i_a}) \wedge (\overline{x_{j_1}} \vee \cdots \vee \overline{x_{j_b}})$ (2-term CNF)
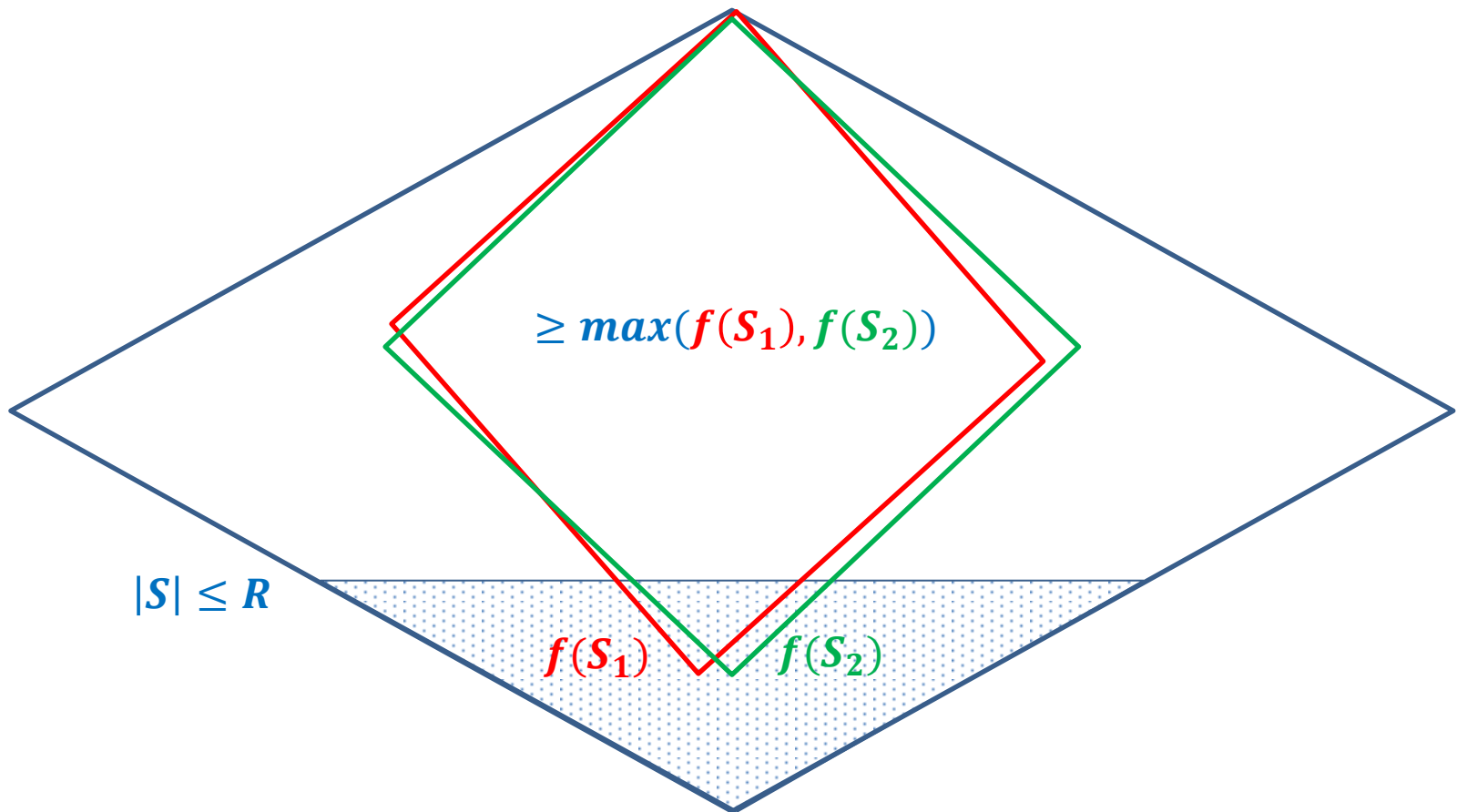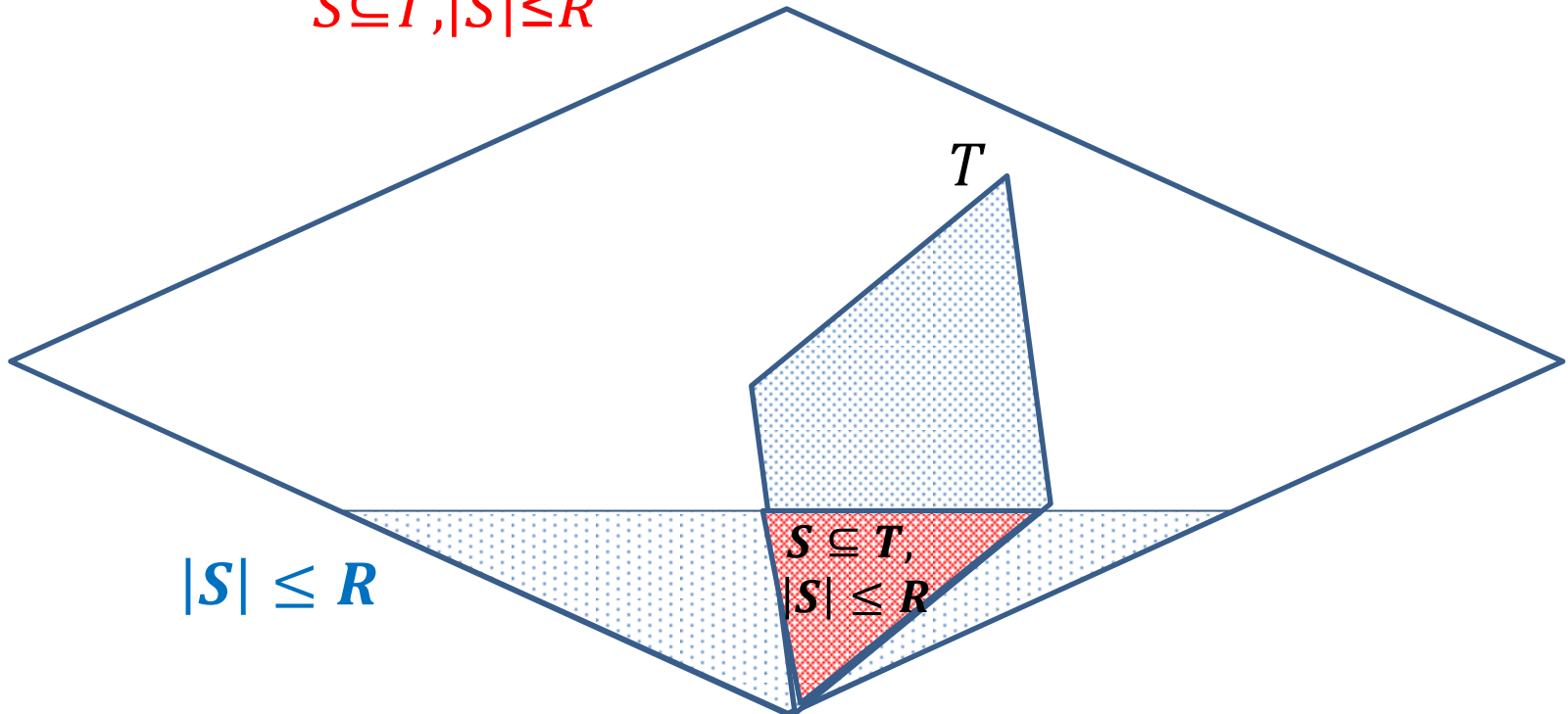
- Monotone submodular

- Submodular



$|S| \geq |X| - R$

$|S| \leq R$

# Discrete monotone submodularity

- Monotone submodular $f: 2^X \to \{0, \dots, R\}$



$\geq max(f(S_1), f(S_2))$

$|S| \leq R$

$f(S_1)$    $f(S_2)$

# Discrete monotone submodularity

- **Theorem**: for **monotone** submodular $f: 2^X \rightarrow \{0, \dots, R\}$ for all $T$: $f(T) = \max\limits_{S \subseteq T, |S| \leq R} f(S)$

- $f(T) \geq \max\limits_{S \subseteq T, |S| \leq R} f(S)$ (by monotonicity)



$T$

$|S| \leq R$

$S \subseteq T,$
$|S| \leq R$

# Discrete monotone submodularity

- $f(T) \leq \max_{S \subseteq T, |S| \leq R} f(S)$

- S' = **smallest** subset of $T$ such that $f(T) = f(S')$

- $\forall x \in S'$ we have $\partial_x f(S' \setminus \{x\}) > 0$ =>

Restriction of $f$ on $2^{S'}$ is **monotone increasing** =>$|S'| \leq R$
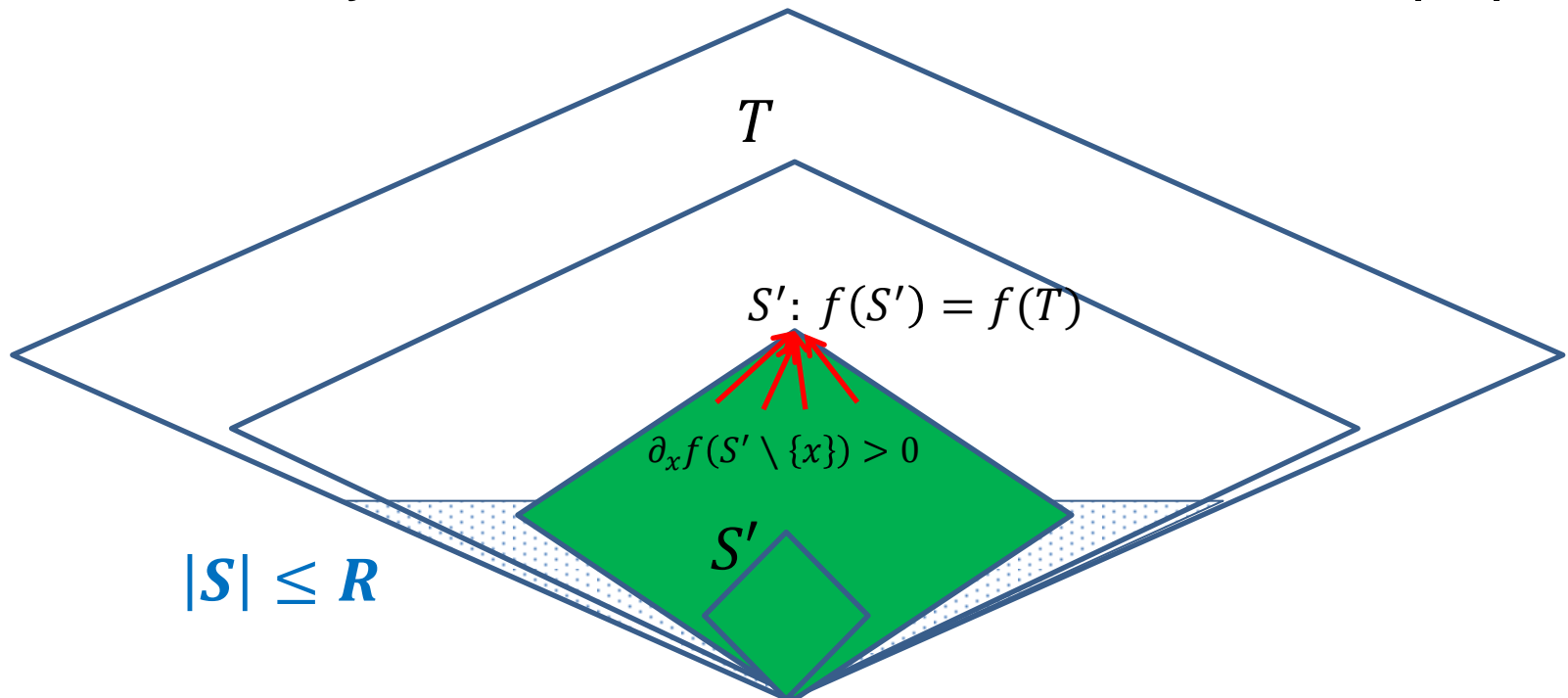
# Representation by a formula

- **Theorem**: for **monotone** submodular $f: 2^X \to \{0, \dots, R\}$ for all $T$:

$$f(T) = \max_{S \subseteq T, |S| \leq R} f(S)$$

- Alternative notation: $|X| \to n$, $2^X \to (x_1, \dots, x_n)$

- **Boolean k–DNF** $= \vee \left( x_{i_1} \wedge \overline{x_{i_2}} \wedge \cdots \wedge x_{i_k} \right)$

- **Pseudo–Boolean k–DNF** ( $\vee \to max$, $A_i = 1 \to A_i \in$ R):

$$max_i \left[ A_i \cdot \left( x_{i_1} \wedge \overline{x_{i_2}} \wedge \cdots \wedge x_{i_k} \right) \right]$$ **(Monotone, if no negations**)

- **Theorem (restated)**:

**Monotone** submodular $f(x_1, \dots, x_n) \to \{0, \dots, R\}$ can be represented as a **monotone** pseudo-Boolean $R$-DNF with constants $A_i \in \{0, \dots, R\}$.
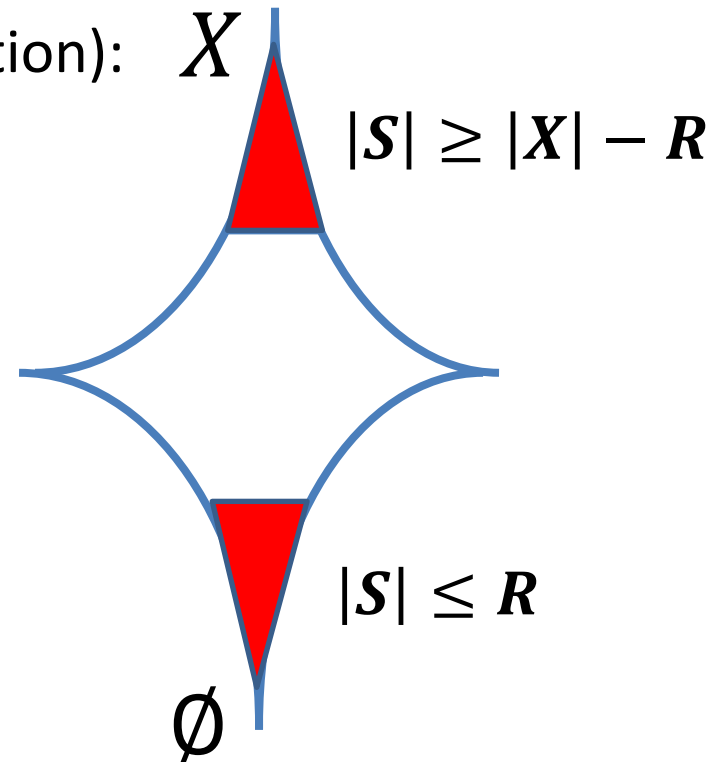
# Discrete submodularity

- Submodular $f(x_1, \ldots, x_n) \to \{0, \ldots, R\}$ can be represented as a pseudo-Boolean **2R**-DNF with constants $A_i \in \{0, \ldots, R\}$.

- Hint [Lovasz] (Submodular monotonization):
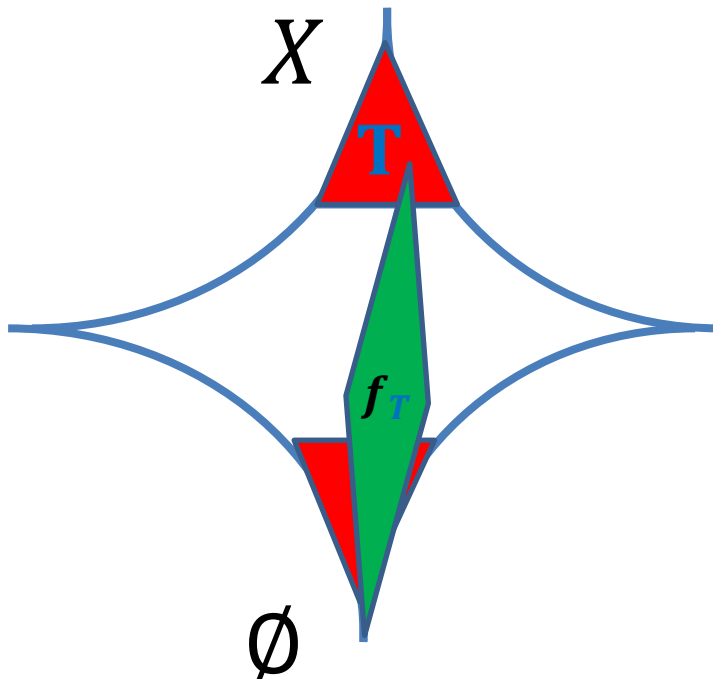
Given submodular $f$, define
$$f^{mon}(S) = min_{S \subseteq T} \, f(T)$$

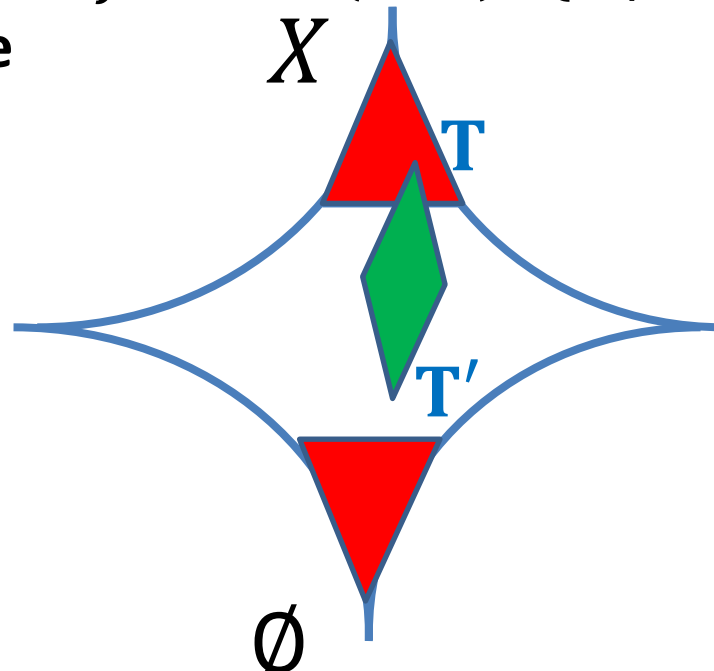Then $f^{mon}$ is monotone and **submodular**.

$X$

$|S| \geq |X| - R$

$|S| \leq R$

$\emptyset$

# Proof

- We're done if we have a **coverage** $C \subseteq 2^X$ :

  1. All $\mathbf{T} \in C$ have large size: $|\mathbf{T}| \geq |X| - R$

  2. For all $S \in 2^X$ there exists $\mathbf{T} \in C : S \subseteq \mathbf{T}$

  3. For every $\mathbf{T} \in C$ restriction $f_{\mathbf{T}}$ of $f$ on $2^{\mathbf{T}}$ is **monotone**

$X$

$\mathbf{T}$
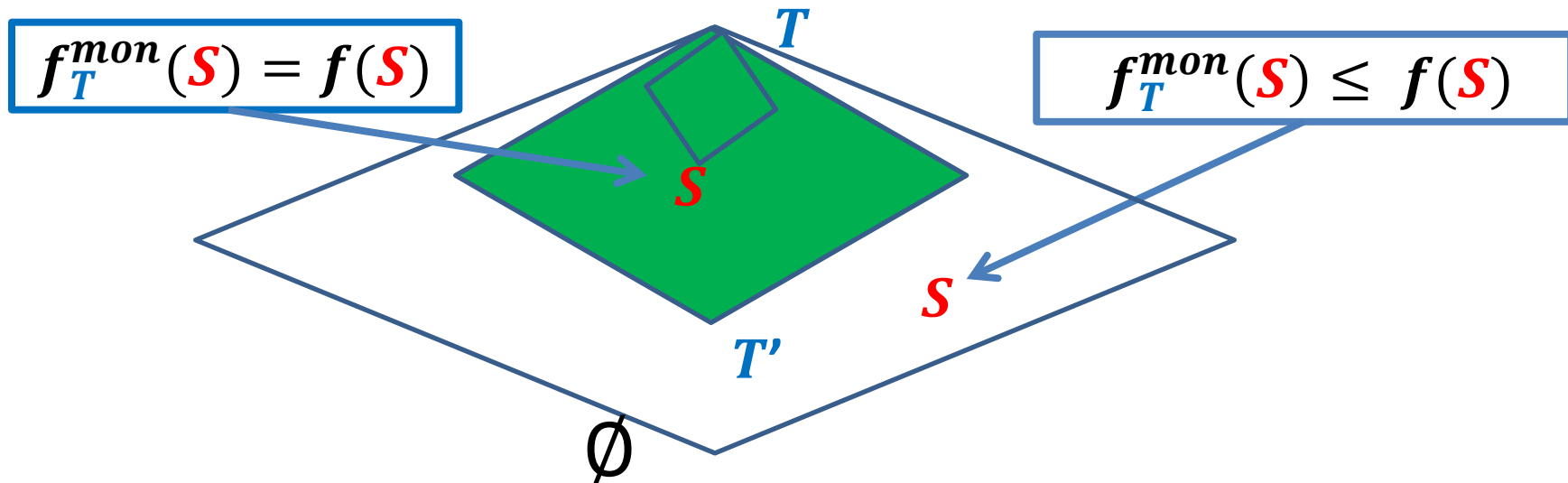
$f_{\mathbf{T}}$

$\emptyset$

- Every $f_{\mathbf{T}}$ is a monotone pB **R**-DNF **(3)**
- Add at most **R** negated variables to every clause to restrict to $2^{\mathbf{T}}$ **(1)**
- $f(S) = \max_{\mathbf{T} \in C} f_{\mathbf{T}}(S)$ **(2)**

# Proof

- There is no such coverage => relaxation [GHRU'11]
  - All $\mathbf{T} \in \boldsymbol{C}$ have large size: $|\mathbf{T}| \geq |X| - \boldsymbol{R}$
  - For all $\boldsymbol{S} \in 2^X$ there exists a pair $\mathbf{T}' \subseteq \boldsymbol{T} \in \boldsymbol{C}$:
    $$\mathbf{T}' \subseteq \boldsymbol{S} \subseteq \boldsymbol{T}$$
  - Restriction of $\boldsymbol{f}$ on all $r(\boldsymbol{T'}, \boldsymbol{T})$: $\{\boldsymbol{S} \,|\, \mathbf{T}' \subseteq \boldsymbol{S} \subseteq \boldsymbol{T}\}$ is **monotone**

$X$

$\mathbf{T}$

$\mathbf{T}'$

$\emptyset$

# Coverage by monotone lower bounds



$$f_T^{mon}(S) = f(S)$$

$$f_T^{mon}(S) \leq f(S)$$

- Let $f_T^{mon}$ be defined as $f_T^{mon}(S) = \min_{S \subseteq S' \subseteq T} f(S')$

  – $f_T^{mon}$ is monotone submodular [Lovasz]

  – For all $S \subseteq T$ we have $f_T^{mon}(S) \leq f(S)$

  – For all $T' \subseteq S \subseteq T$ we have $f_T^{mon}(S) = f(S)$

- $f(S) = \max_{T \in C} f_T^{mon}(S)$ (where $f_T^{mon}$ is a monotone pB R-DNF)

# Learning pB-formulas and k-DNF

- $DNF^{k,R}$ = class of pB $k$-DNF with $A_i \in \{0, \ldots, R\}$

- **i-slice** $f_i(x_1, \ldots, x_n) \to \{0,1\}$ defined as

$$f_i(x_1, \ldots, x_n) = 1 \quad \textbf{iff} \quad f(x_1, \ldots, x_n) \geq i$$

- If $f \in DNF^{k,R}$ its **i-slices** $f_i$ are $k$-DNF and:

$$f(x_1, \ldots, x_n) = \max_{1 \leq i \leq R} (i \cdot f_i(x_1, \ldots, x_n))$$

- PAC-learning:

$$\Pr_{rand(A)} \left[ \Pr_{S \sim U(\{0,1\}^n)} [A(S) = f(S)] \geq 1 - \epsilon \right] \geq \frac{1}{2}$$

- Learn every **i-slice** $f_i$ on $(1 - \epsilon / R)$ fraction of arguments => union bound
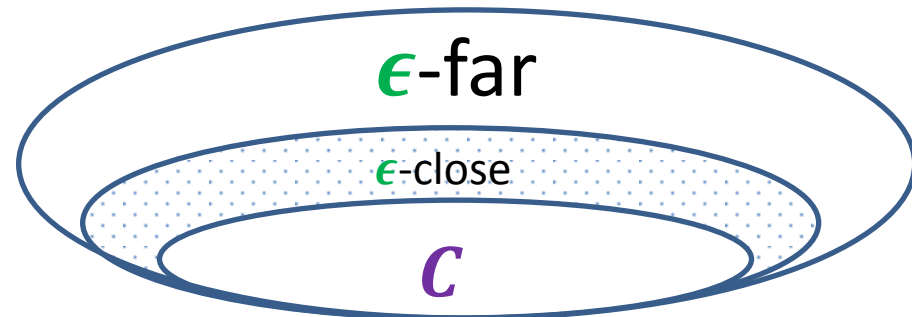
# Learning Fourier coefficients

- Learn $f_i$ ($k$-DNF) on $1 - \epsilon' = (1 - \epsilon / R)$ fraction of arguments

- **Fourier sparsity $S_C(\epsilon)$** = # of largest Fourier coefficients sufficient to PAC-learn every $f \in C$

- $S_{k-\text{DNF}}(\epsilon) = k^{O\left(k \log\left(\frac{1}{\epsilon}\right)\right)}$ [Mansour]: doesn't depend on **n**!
  - Kushilevitz-Mansour (Goldreich-Levin): $poly(n, S_F)$ queries/time.
  - ``Attribute efficient learning'': $\boldsymbol{polylog}(n) \cdot poly(S_F)$ queries
  - Lower bound: $\Omega(2^k)$ queries to learn a random $k$-junta ($\in k$-DNF) up to constant precision.

- $S_{DNF^{k,R}}(\epsilon) = k^{O\left(k \log\left(\frac{R}{\epsilon}\right)\right)}$
  - Optimizations: Do all **R** iterations of KM/GL in parallel by reusing queries

# Property testing

- Let $C$ be the class of submodular $f: \{0,1\}^n \to \{0, \ldots, R\}$
- How to (approximately) test, whether a given $f$ is in $C$?
- Property tester: (randomized) algorithm for distinguishing:

1. $f \in C$

2. ($\epsilon$-far): $\min\limits_{g \in C} \big| f - g \big|_H \geq \epsilon \, 2^n$



- Key idea: $k$-DNFs have small representations:
  - [Gopalan, Meka, Reingold CCC'12]  (using quasi-sunflowers [Rossman'10])

  $\forall \epsilon > 0, \forall k$-DNF formula F there exists:

  $k$-DNF formula F' of size $\leq \left( k \log \frac{1}{\epsilon} \right)^{O(k)}$ such that $\big| F - F' \big|_H \leq \epsilon 2^n$

# Testing by implicit learning

- **Good approximation by juntas => efficient property testing** [Diakonikolas, Lee, Matulef, Onak ,Rubinfeld, Servedio, Wan]

  - $\epsilon$-approximation by $J(\epsilon)$-junta

  - Good dependence on $\epsilon$: $J_{k-\text{DNF}}(\epsilon) = \left(k \log \frac{1}{\epsilon}\right)^{O(k)}$

- For submodular functions $f: \{0,1\}^n \to \{0, \dots, R\}$

  - Query complexity $\left(R \log \frac{R}{\epsilon}\right)^{\tilde{O}(R)}$, independent of **n**!

  - Running time exponential in $J(\epsilon)$

  - $\Omega(k)$ lower bound for testing $k$-DNF (reduction from Gap Set Intersection)

- [Blais, Onak, Servedio, Y.] **exact** characterization of submodular functions

$$J(\epsilon) = \left[O\left(R \log R + \log \frac{1}{\epsilon}\right)\right]^{(R+1)}$$

# Previous work on testing submodularity

$f : \{0,1\}^n \rightarrow [0, R]$ [Parnas, Ron, Rubinfeld '03, Seshadhri, Vondrak, ICS'11]:

- Upper bound $(1/\epsilon)^{O(\sqrt{n})}$.
- Lower bound: $\Omega(n)$

} Gap in query complexity

Special case: coverage functions [Chakrabarty, Huang, ICALP'12].

# Directions

- Close gaps between upper and lower bounds, extend to more general learning/testing settings

- Connections to optimization?

- What if we use $L_1 -$distance between functions instead of Hamming distance in property testing? [Berman, Raskhodnikova, Y.]