

# CMSC 510 – L14

## Regularization Methods for Machine Learning



---

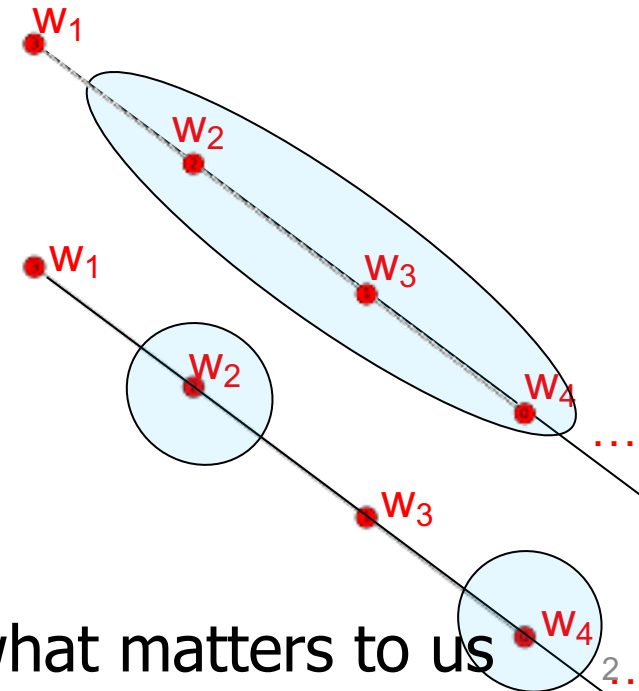
Instructor:  
Dr. Tom Arodz

# Fused Lasso – view involving sets

- If  $f_i$  is important for classification,  $f_{i-1}$  and  $f_{i+1}$  likely to be important, too
- We want a classifier where either both neighboring features are selected (non-zero  $w_i$  and  $w_{i+1}$ ) or both are not selected ( $w_i = w_{i+1} = 0$ ):

$$\Omega(w) = \sum_{f=2}^F |[w_{f-1}] - [w_f]|$$
$$[x] = \text{supp}(x) = |\text{sign}(x)|$$

- $\Omega$  now is essentially a function defined on a set, not on a vector
  - Set of features with non-zero feature weights  $w_f$
  - E.g.  $\Omega(\{f_2, f_3, f_4\}) = 2$   
or  $\Omega(\{f_2, f_4\}) = 4$
  - Detailed values of weights are not what matters to us

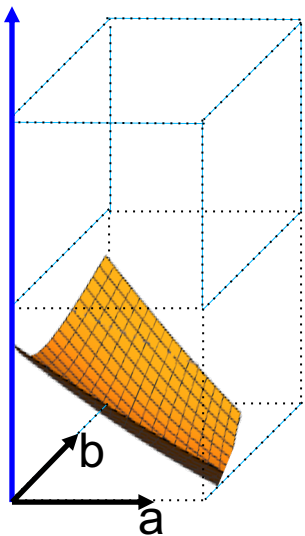


-

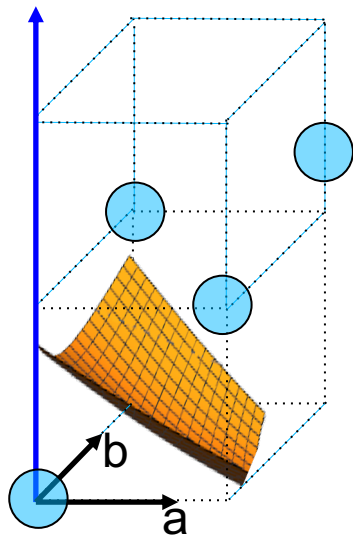
# Big picture

Submodular set functions  
 $\Rightarrow$  convex regularizers!

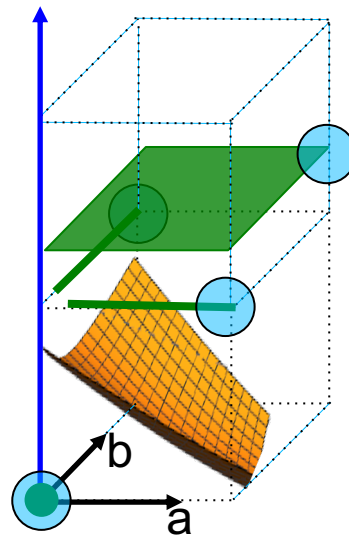
- Example: two features:  $x_a$  and  $x_b$
- Any linear classifier is  $y = \text{sign}(w_a x_a + w_b x_b)$
- What are the weights  $w = (w_a, w_b)$ ?
  - To find out, we optimize  $\text{Risk}(w) + \text{Penalty}(w)$



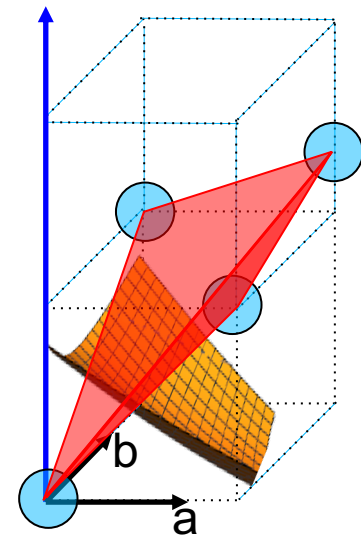
$\text{Risk}(w)$



We can't just add  
penalty on the set  
(of features)  $\Omega(S)$  to  
risk over vectors (of  
feature weights).



What we want:  
 **$\text{Risk}(w) + \text{Penalty}(w)$**   
 $\text{Penalty}(w) = \Omega([w])$   
tough to solve!



What we can:  
 **$\text{Risk}(w) + \text{Penalty}(w)$**   
 $\text{Penalty}(w) = \Omega^-(w)$   
both convex, so  
often easy to solve

# Set functions

- Set function over universe  $V$ , ie.  
 $\Omega: 2^V \rightarrow \mathbb{R}$

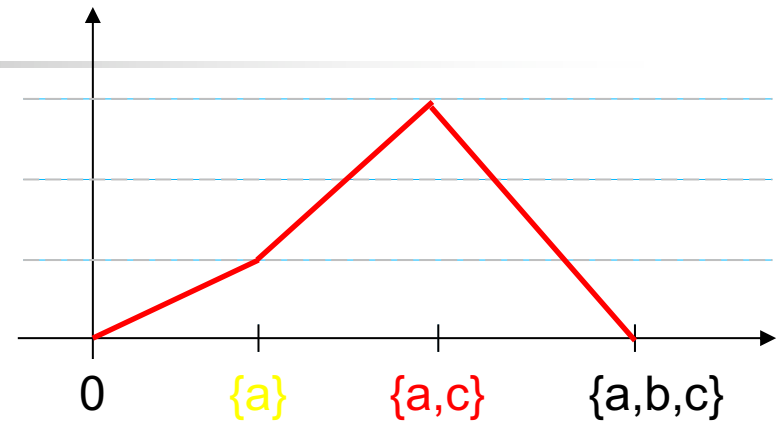
- E.g.  $\Omega_{GC}(S) = \text{graph cut capacity}$   
 $= (S, V-S) = \text{sum of weights of edges between } S \text{ and } V-S$   
 (or just number of edges)

- Assumptions:

- weights are positive
- no edge = zero weight edge

E.g.  $V = \{a, b, c\}$  with this graph:

- $\Omega_{GC}(\{\emptyset\}) = \Omega_{GC}(\{a, b, c\}) = 0$
- $\Omega_{GC}(\{a\}) = \Omega_{GC}(\{b, c\}) = 1$
- $\Omega_{GC}(\{c\}) = \Omega_{GC}(\{a, b\}) = 2$
- $\Omega_{GC}(\{b\}) = \Omega_{GC}(\{a, c\}) = 3$

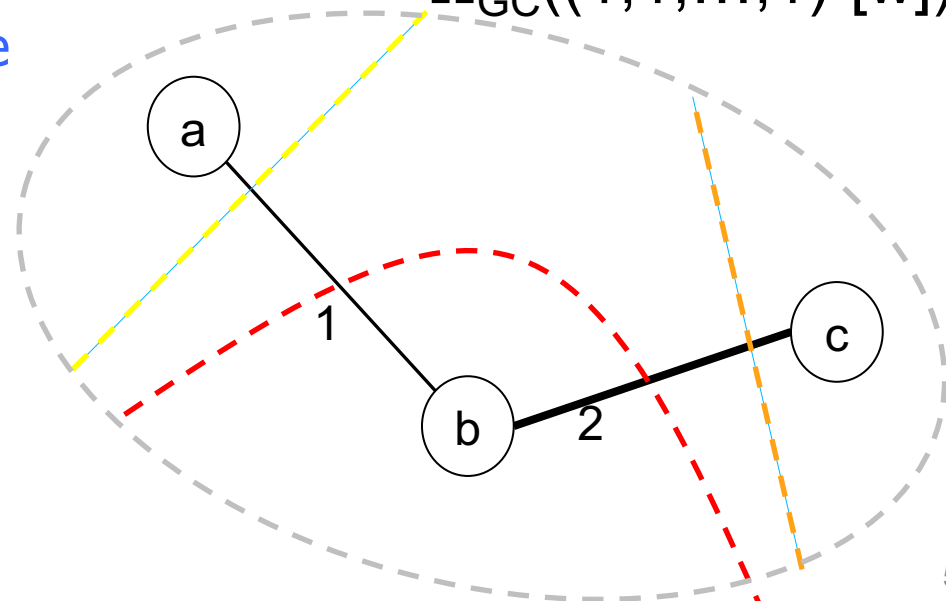


$\Omega_{GC}$  is symmetric:

$$\Omega_{GC}(S) = \Omega_{GC}(V-S)$$

$$\Omega_{GC}([w])$$

$$= \Omega_{GC}((1, 1, \dots, 1) - [w])$$

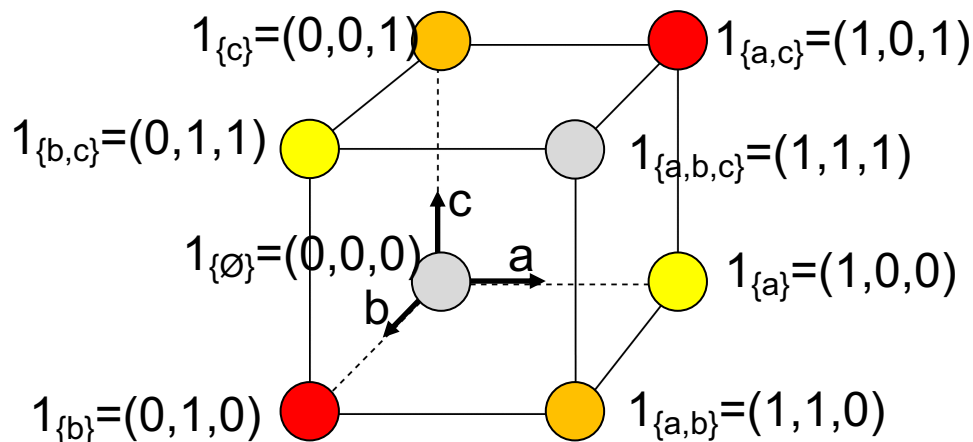


# Set functions

## ■ We can use binary cube to visualize a set function

E.g.  $V=\{a,b,c\}$  with this graph:

- $\Omega_{GC}(\{\emptyset\}) = \Omega_{GC}(\{a,b,c\})=0$
- $\Omega_{GC}(\{a\}) = \Omega_{GC}(\{b,c\}) = 1$
- $\Omega_{GC}(\{c\}) = \Omega_{GC}(\{a,b\}) = 2$
- $\Omega_{GC}(\{b\}) = \Omega_{GC}(\{a,c\}) = 3$



Each dimension corresponds to an element in the universe  $V$

**Notation:**

**vector to set:**

$[w]$ =support of  $w$

vector of 0's and 1's

can be interpreted as a set

**set to vector:**

$1_S$  = vector with 1's for coordinates corresponding to members of  $S$ , and 0's elsewhere

*Often we simplify notation and just use set=vector*

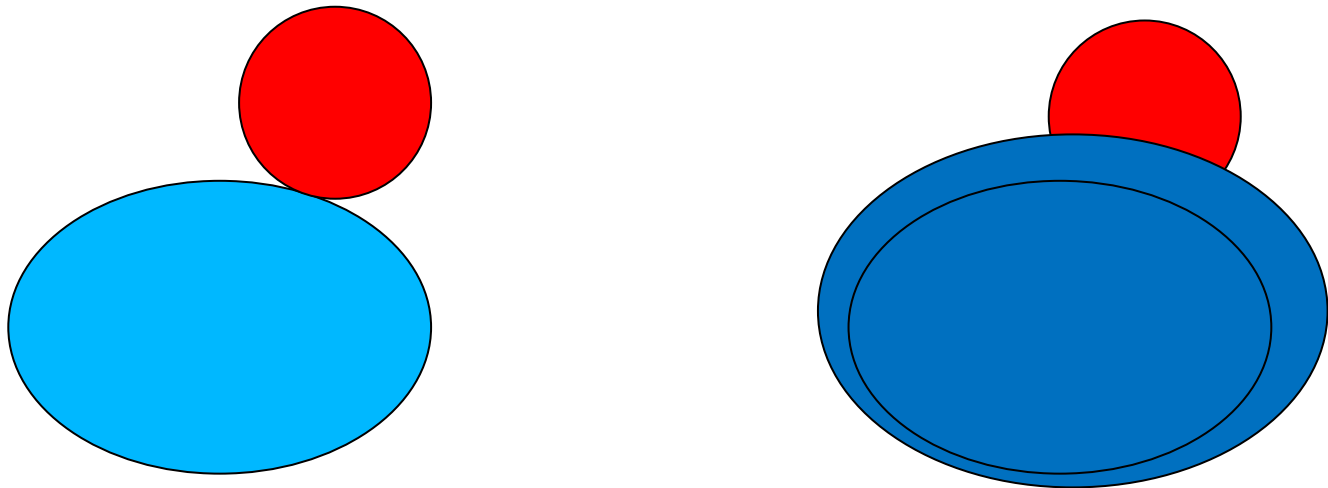
$$\Omega(S) = \Omega(1_S) = \Omega([w])$$

# Submodularity

- Set function over universe  $V$ , ie.  $\Omega: 2^V \rightarrow \mathbb{R}$
- $\Omega_{GC}$  is an example of family of set functions called **submodular set functions**

**Diminishing returns:** 10<sup>th</sup> slice of cake is not as good as the first one!

$$f(\text{red} + \text{blue}) - f(\text{blue}) \geq f(\text{red} + \text{navy}) - f(\text{navy})$$

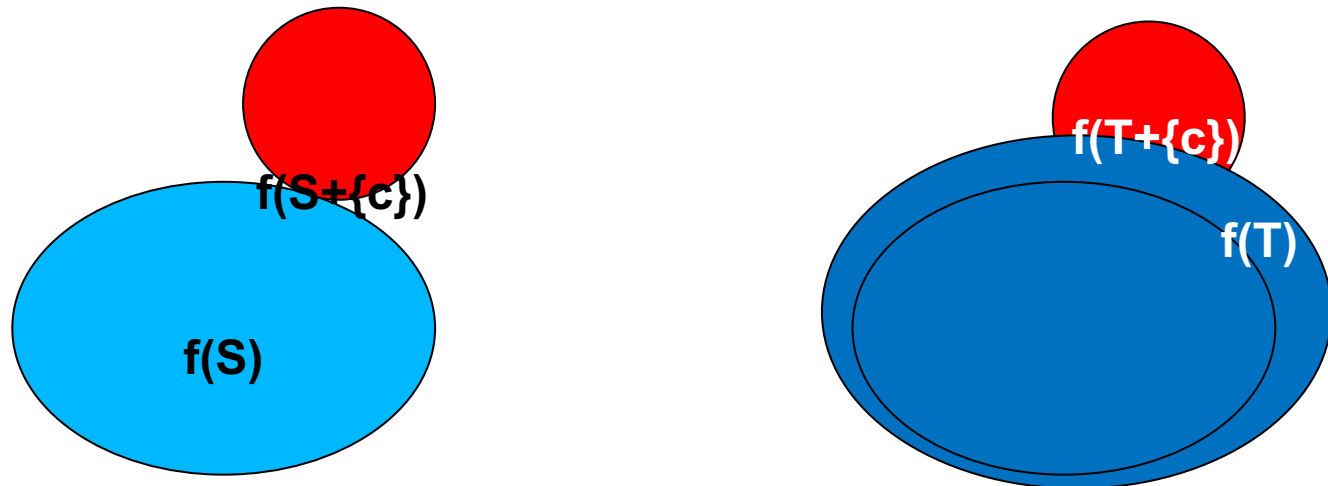


# Submodularity

- Set function over universe  $V$ , ie.  $\Omega: 2^V \rightarrow \mathbb{R}$
- $\Omega_{GC}$  is an example of family of set functions called **submodular set functions**

$$\forall S, T, \{c\} \subseteq V, \quad S \subseteq T, \\ \Omega(S \cup \{c\}) - \Omega(S) \geq \Omega(T \cup \{c\}) - \Omega(T)$$

$$f(\text{red}+\text{blue}) - f(\text{blue}) \geq f(\text{red}+\text{navy}) - f(\text{navy})$$



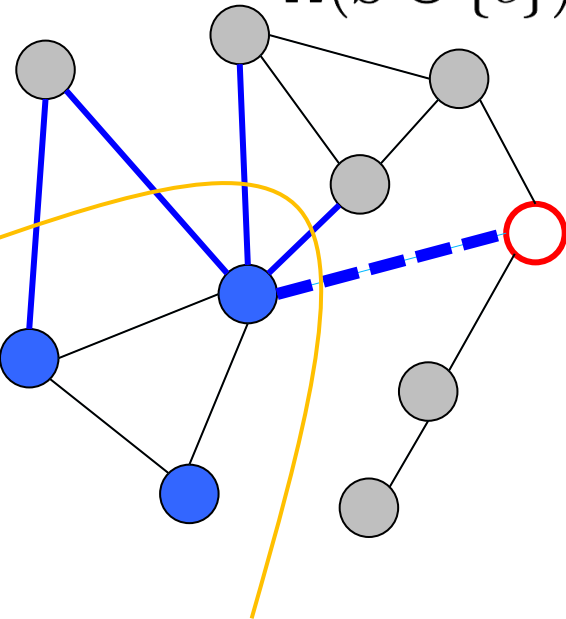
- Adding the same element  $\{c\}$  to a set T increases  $\Omega$  less than adding it to a subset S of T



# Is Graph cut submodular?

c ○ S ● T ●

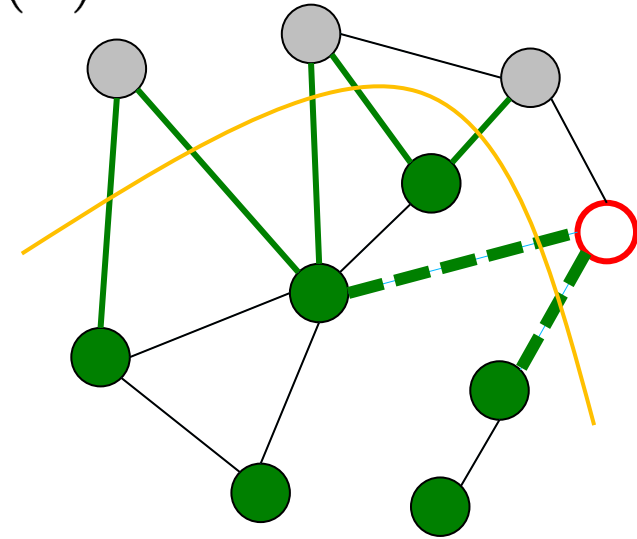
$$\Omega(S \cup \{c\}) - \Omega(S) \geq \Omega(T \cup \{c\}) - \Omega(T)$$



$\Omega(S)=5$   
number of edges  
crossing the cut  
from S to V-S




To check if graph cut is  
submodular, we need:

- Set T
- Set S, a subset of T
- A vertex c  
(not a member of T,  
otherwise we have  
equality)
- Then, we observe what  
happens to  $\Omega(S)$  and  $\Omega(T)$   
when we add  $\{c\}$



$\Omega(T)=(T,V-T)=6$   
number of edges  
crossing the cut  
from T to V-T

# Is Graph cut submodular?

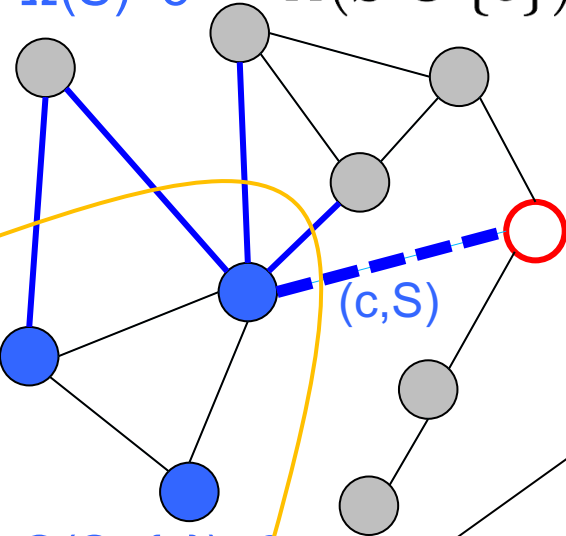
$c$  
 S  
 T

$$\Omega(T) = (T, V-T) = 6$$

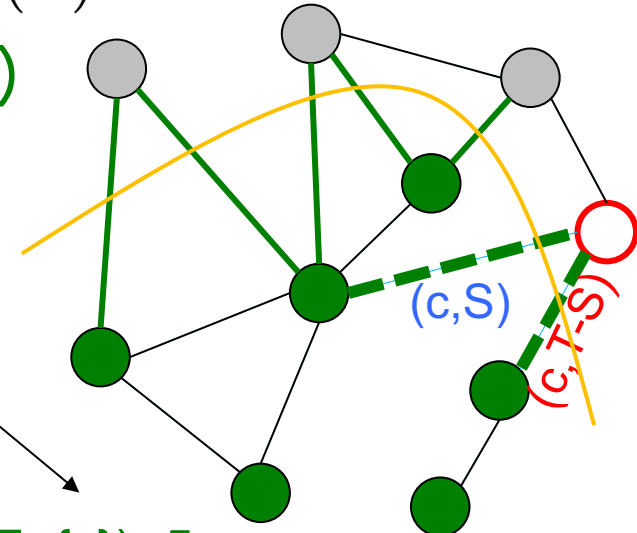
$$\Omega(S) = 5$$

$$\Omega(S \cup \{c\}) - \Omega(S) \geq \Omega(T \cup \{c\}) - \Omega(T)$$

$$\dots + (c, T-S) \geq \dots - (c, T-S)$$



what happens to  $\Omega(S)$  and  $\Omega(T)$  when we add  $\{c\}$

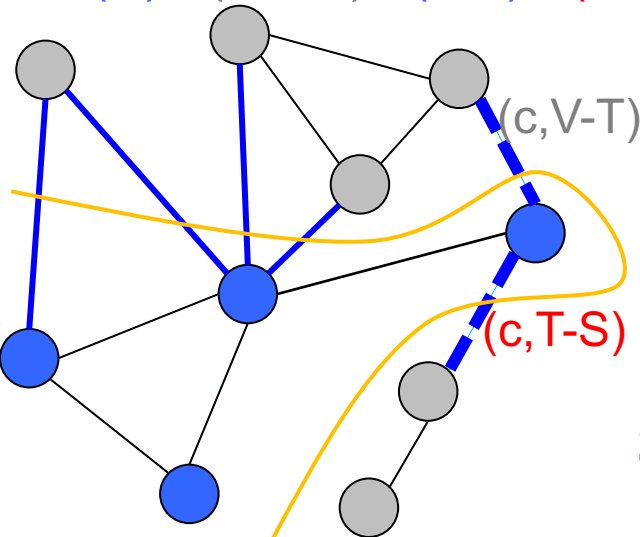


$$\Omega(S + \{c\}) = 6$$

$$= \Omega(S) + (c, V-T) - (c, S) + (c, T-S)$$

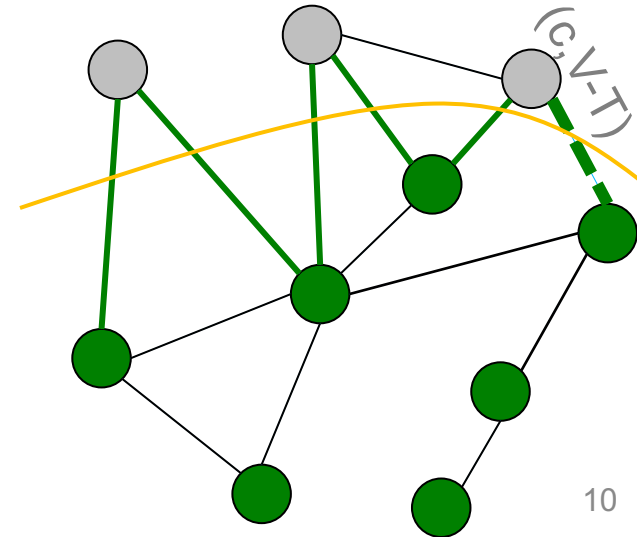
$$\Omega(T + \{c\}) = 5$$

$$= \Omega(T) + (c, V-T) - (c, S) - (c, T-S)$$

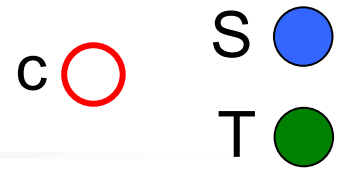


If vertex  $c$  is connected to vertices in  $(T-S)$  with nonzero sum of weights  $(c, T-S)$ , we have " $\geq$ " inequality

So: graph cut capacity is a submodular set function



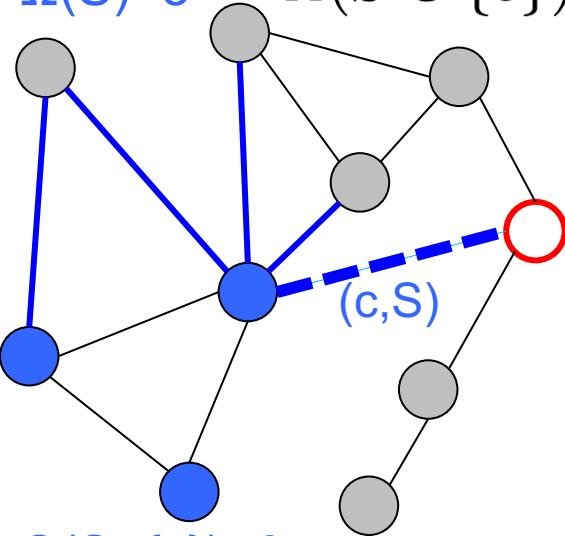
# Is Graph cut submodular?



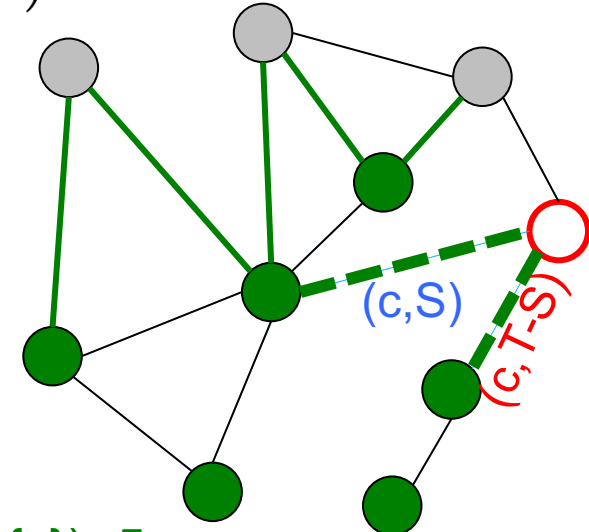
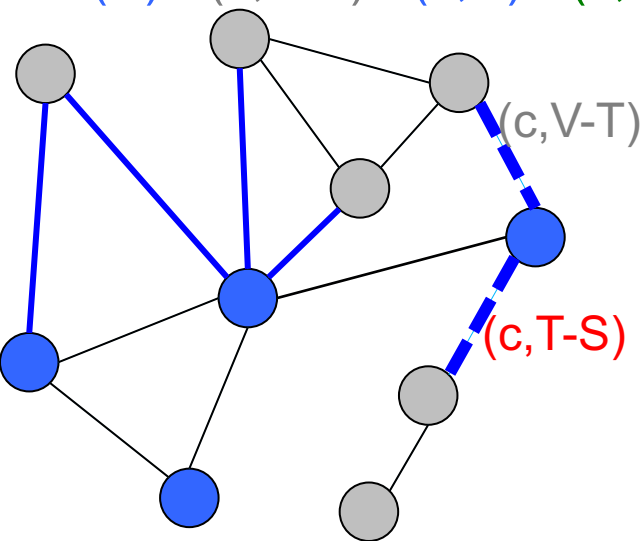
$$\Omega(T) = (T, V-T) = 6$$

$$\Omega(S) = 5$$

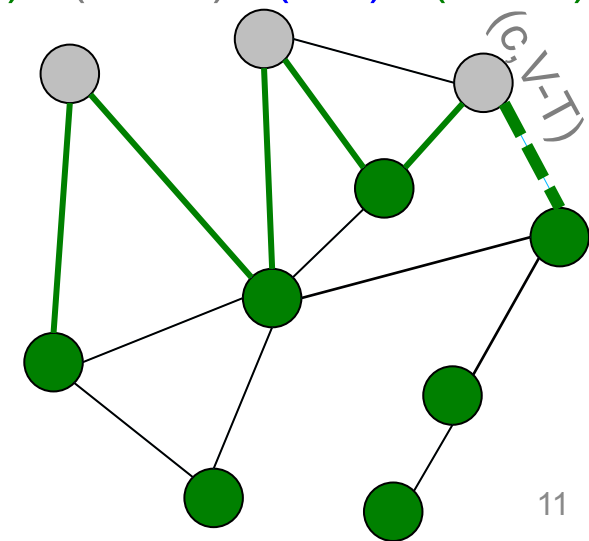
$$\Omega(S \cup \{c\}) - \Omega(S) \geq \Omega(T \cup \{c\}) - \Omega(T)$$



$$\begin{aligned} \Omega(S + \{c\}) &= 6 \\ &= \Omega(S) + (c, V-T) - (c, S) + (c, T-S) \end{aligned}$$

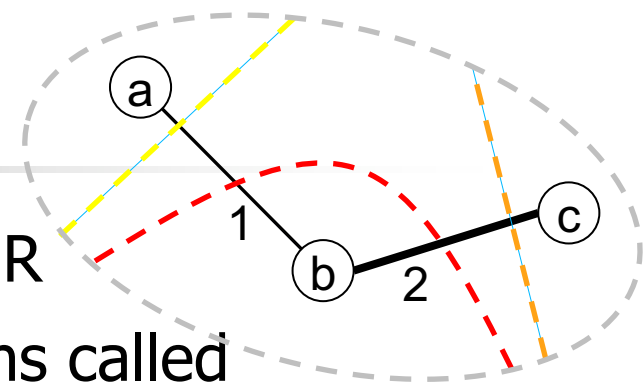


$$\begin{aligned} \Omega(T + \{c\}) &= 5 \\ &= \Omega(T) + (c, V-T) - (c, S) - (c, T-S) \end{aligned}$$



We have equality when  $(c, T-S) = 0$ , i.e., when there are no edges between vertex  $c$  and vertices in  $(T-S)$

# Submodularity

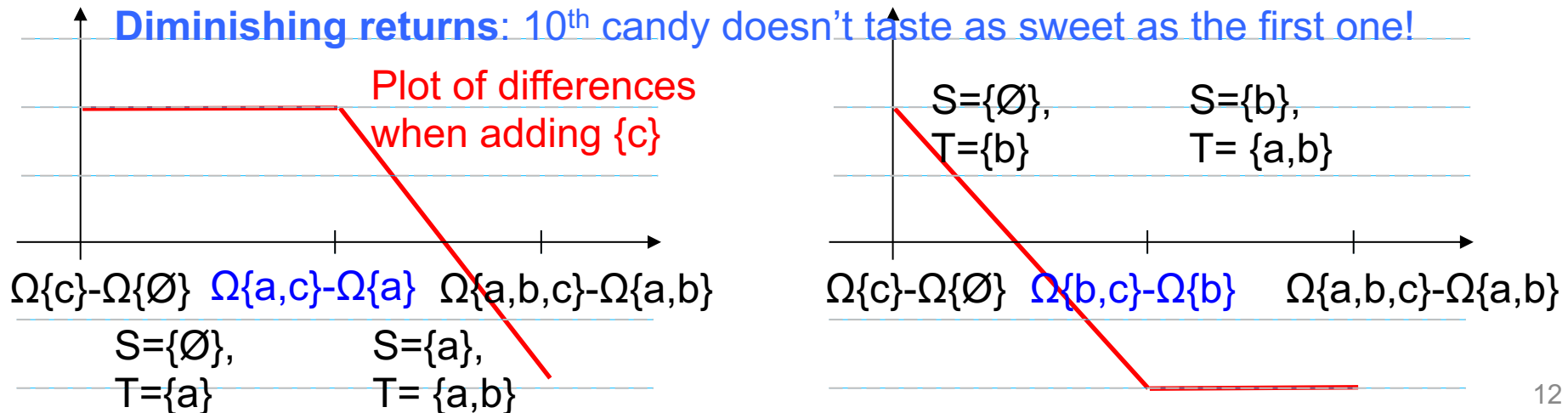


- Set function over universe  $V$ , ie.  $\Omega: 2^V \rightarrow \mathbb{R}$
- $\Omega_{GC}$  is an example of family of set functions called **submodular set functions**

- Set function  $\Omega: 2^V \rightarrow \mathbb{R}$  is **submodular** iff:

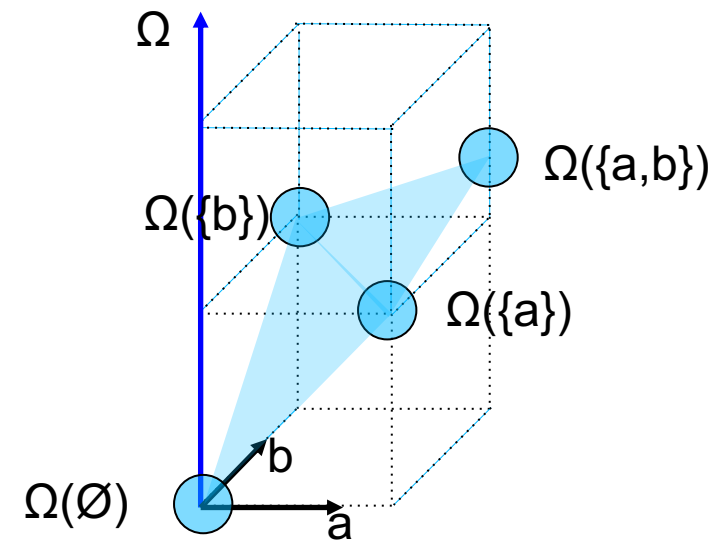
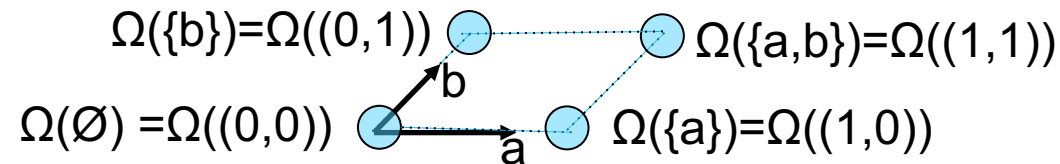
$$\forall S, T, \{c\} \subseteq V, \quad S \subseteq T, \\ \Omega(S \cup \{c\}) - \Omega(S) \geq \Omega(T \cup \{c\}) - \Omega(T)$$

- Adding the same element  $\{c\}$  to a set  $T$  increases  $\Omega$  less than adding it to a subset  $S$  of  $T$



# Submodular set functions

- Submodular set functions are a bit similar to concave functions - they're "bending down"
- we can plot the function on  $\{a,b\}$  on a unit cube
- value  $\Omega(\{a,b\})$  is lower than  $\Omega(\{a\}) + \Omega(\{b\})$ 
  - $S=\emptyset, T=\{a\}, c=b$
  - Assume w.l.o.g.  $\Omega(\emptyset)=0$



submodular  $\Omega$

$$\forall S, T, \{c\} \subseteq V, \quad S \subseteq T,$$

$$\Omega(S \cup \{c\}) - \Omega(S) \geq \Omega(T \cup \{c\}) - \Omega(T)$$

# Alternative definitions

- **Modular** set functions:  $\Omega(\{a,b,\dots,z\}) = \Omega_a + \Omega_b + \dots + \Omega_z$

$$\forall A, B \subseteq V \quad \Omega(A \cup B) + \Omega(A \cap B) = \Omega(A) + \Omega(B)$$

- **Submodular** set functions:  $\Omega(\{a,b,\dots,z\}) \leq \Omega_a + \Omega_b + \dots + \Omega_z$

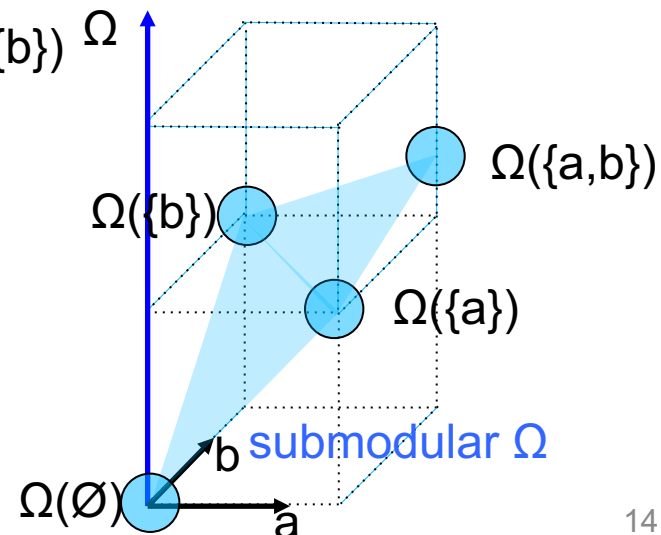
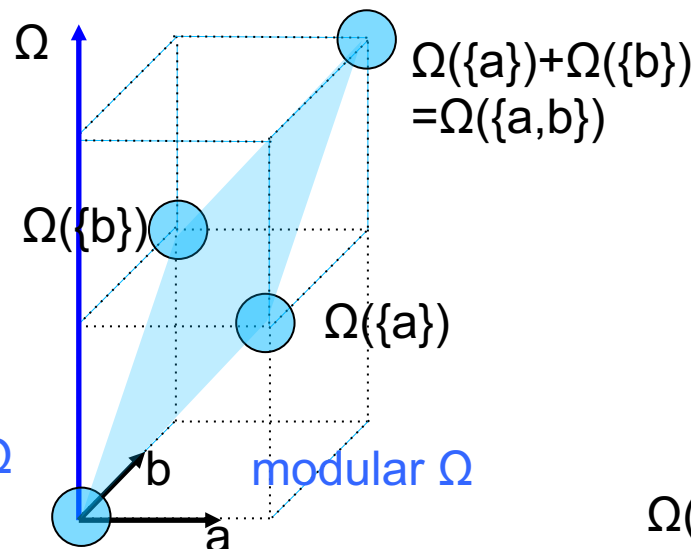
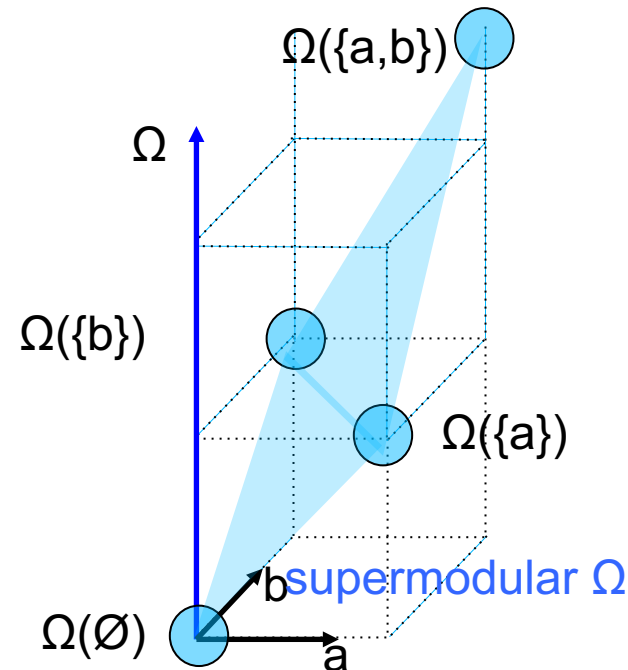
$$\Omega(A \cup B) + \Omega(A \cap B) \leq \Omega(A) + \Omega(B)$$

- **Supermodular** set functions:  $\Omega(\{a,b,\dots,z\}) \geq \Omega_a + \Omega_b + \dots + \Omega_z$

$$\Omega(A \cup B) + \Omega(A \cap B) \geq \Omega(A) + \Omega(B)$$

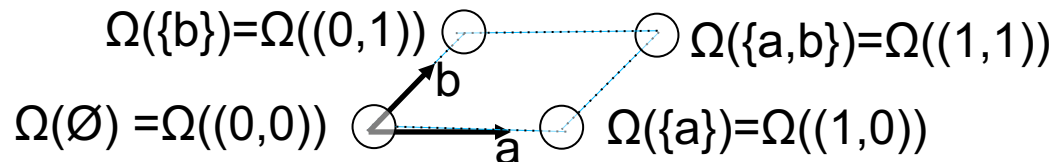
- If  $\Omega$  is modular, it is both submodular and supermodular

- If  $\Omega$  is submodular,  $-\Omega$  is supermodular

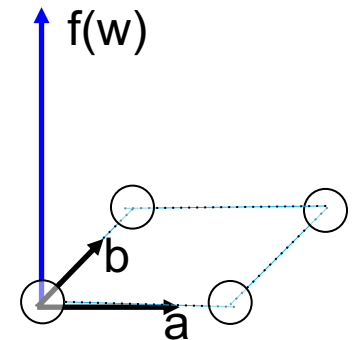


# Submodular set functions

- Set functions: defined vectors with coordinates  $\{0,1\}^F$ , e.g.  $\{a,c\}=(1,0,1)$ 
  - i.e. vectors at the corner of unit cube  $\{0,1\}^F$  where  $F=|V|$  is number of elements in  $V$



- Given a **set** function  $\Omega$  can we define a function  $f(w)$   
 $f: [0,1]^F \rightarrow \mathbb{R}$   
 (or even better,  $\mathbb{R}^F \rightarrow \mathbb{R}$ )?



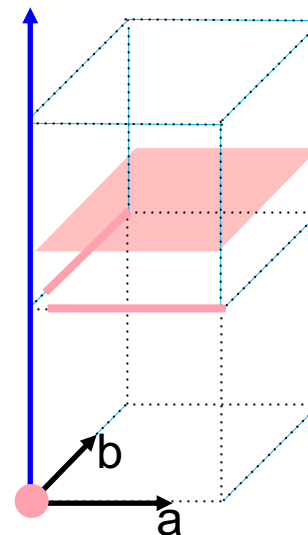
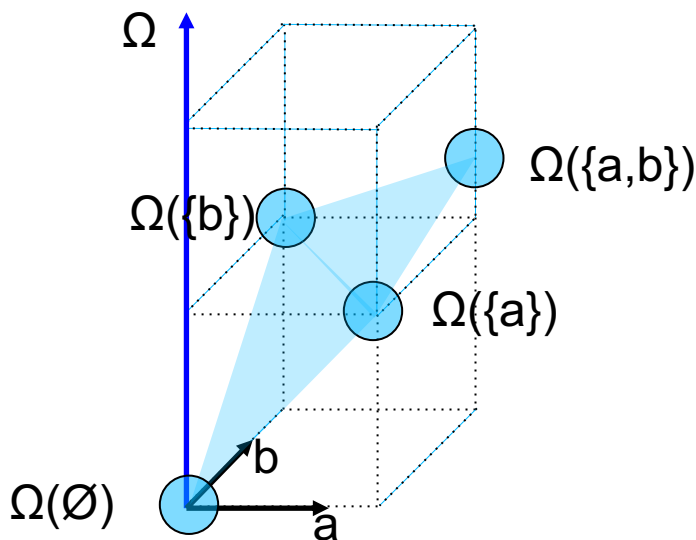
- A function that would be somehow related to  $\Omega$ 
  - We already did!
  - When discussing notation:
  - $f(w) = \Omega([w])$   $f = \# \text{ of non-zero elements in } w$ 
    - For any vector  $\mathbb{R}^F$ , we get vector of 0's/1's, which represents a set, and defines value of  $f(w)$
    - What's wrong with using  $f(w) = \Omega([w])$ ?

# Submodular set functions

- Set functions: defined vectors with coordinates  $\{0,1\}$ , e.g.  $\{a,c\}=(1,0,1)$ 
  - i.e. vectors at the corner of unit cube  $\{0,1\}^F$  where  $F=|V|$  is number of elements in  $V$

$$\begin{array}{ccc} \Omega(\{b\})=\Omega((0,1)) & & \Omega(\{a,b\})=\Omega((1,1)) \\ \Omega(\emptyset)=\Omega((0,0)) & \xrightarrow{a} & \Omega(\{a\})=\Omega((1,0)) \end{array}$$

- Given a set function  $\Omega$  can we define a “nicely behaving” function  $f: [0,1]^F \rightarrow \mathbb{R}$  (or even better,  $\mathbb{R}^F \rightarrow \mathbb{R}$ )?
- A function that would be somehow related to  $\Omega$ 
  - E.g. the light blue piecewise linear function on  $[0,1] \times [0,1]$  we’ve seen before:



$$f(w)=\Omega([w])$$

Not continuous!

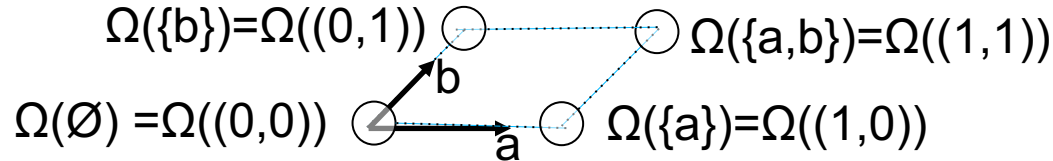
Not convex!

e.g. try  $(1,0)$  and  $(0,1)$  and points in between

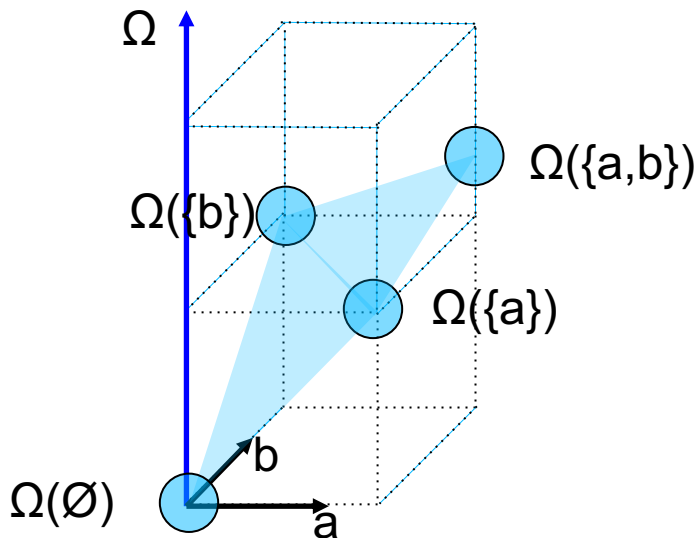


# Submodular set functions

- Set functions: defined vectors with coordinates  $\{0,1\}$ , e.g.  $\{a,c\}=(1,0,1)$ 
  - i.e. vectors at the corner of unit cube  $\{0,1\}^F$  where  $F=|V|$  is number of elements in  $V$

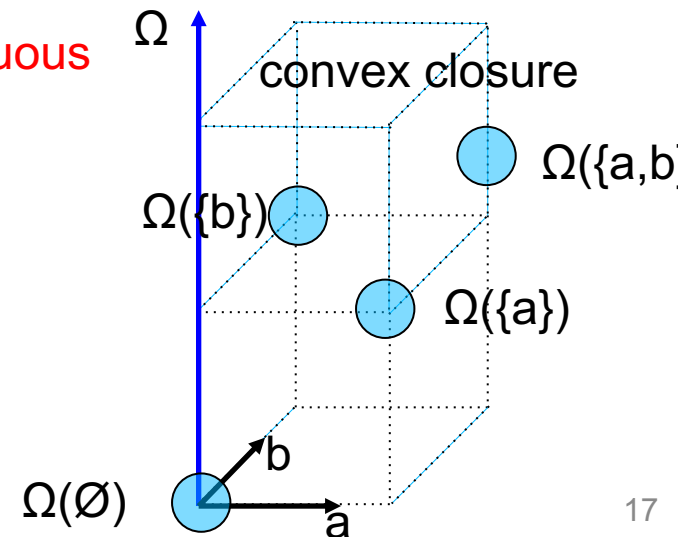


- Given a set function  $\Omega$  can we define a “nicely behaving” function  $f: [0,1]^F \rightarrow \mathbb{R}$  (or even better,  $\mathbb{R}^F \rightarrow \mathbb{R}$ )?
- A function that would be somehow related to  $\Omega$ 
  - E.g. the light blue piecewise linear function on  $[0,1] \times [0,1]$  we’ve seen before:



<- this one is continuous

But is it convex?



# Submodular set functions

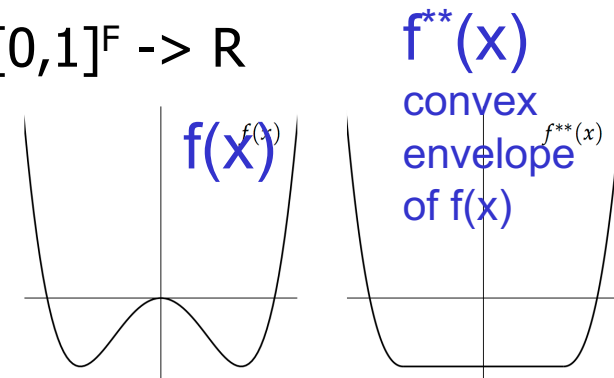
- Given a set function  $\Omega$  we define two functions  $f: [0,1]^F \rightarrow \mathbb{R}$

- Convex closure** of  $\Omega$ :  $\Omega^-$

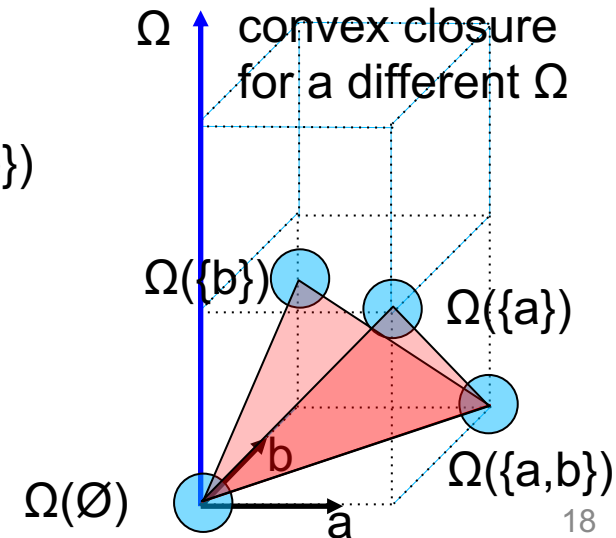
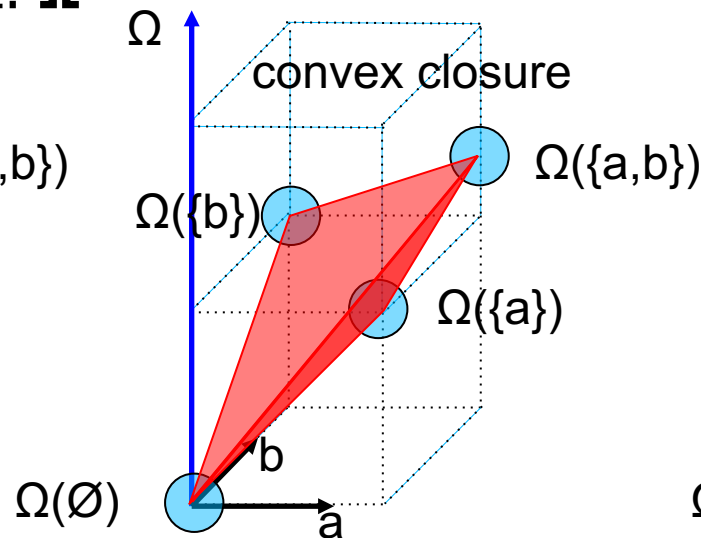
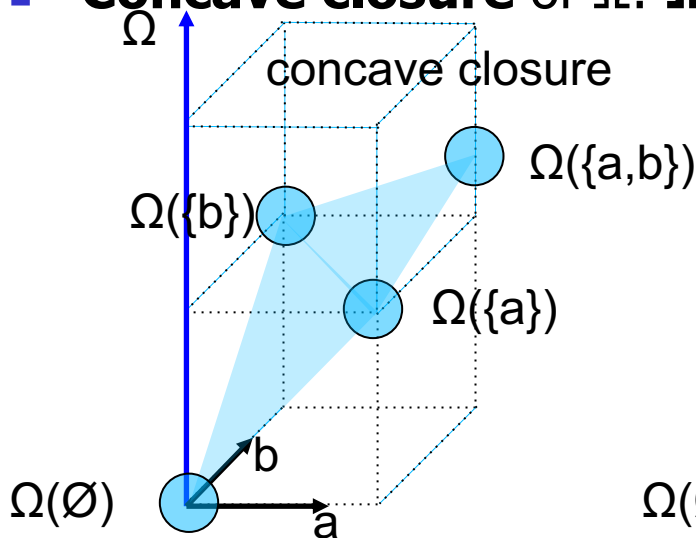
- It is the **convex envelope** ( $\Omega^{**}$ ) of  $\Omega$ : **pointwise highest convex function that bounds  $\Omega$  from below**:  
for any  $w$  in  $\{0,1\}^F$ :  $\Omega^-(w) \leq \Omega(w)$

- Convex closure of set function is **piecewise linear**:

Intuition: Start with any convex function that bounds  $\Omega$  from below, push every point of it up until you can't without losing convexity – you will get maxima of hyperplanes

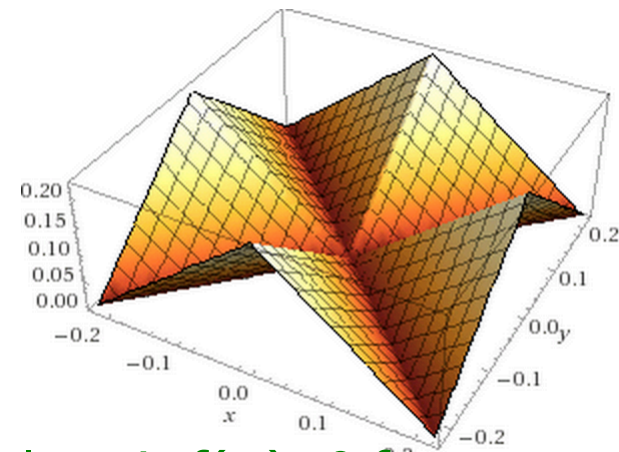


- Concave closure** of  $\Omega$ :  $\Omega^+$



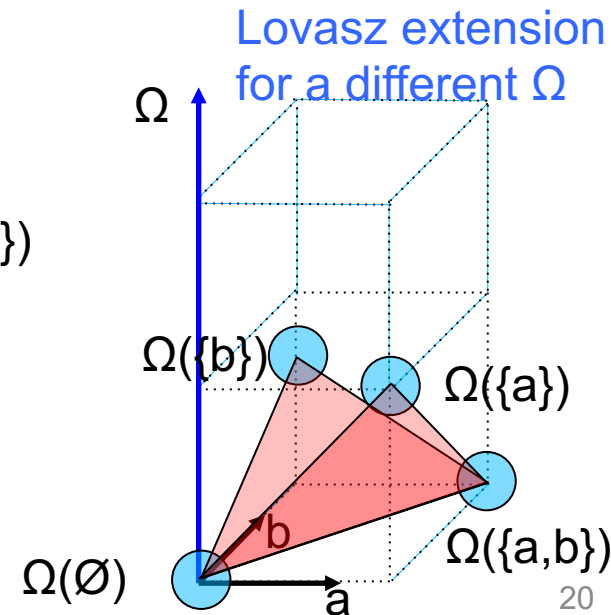
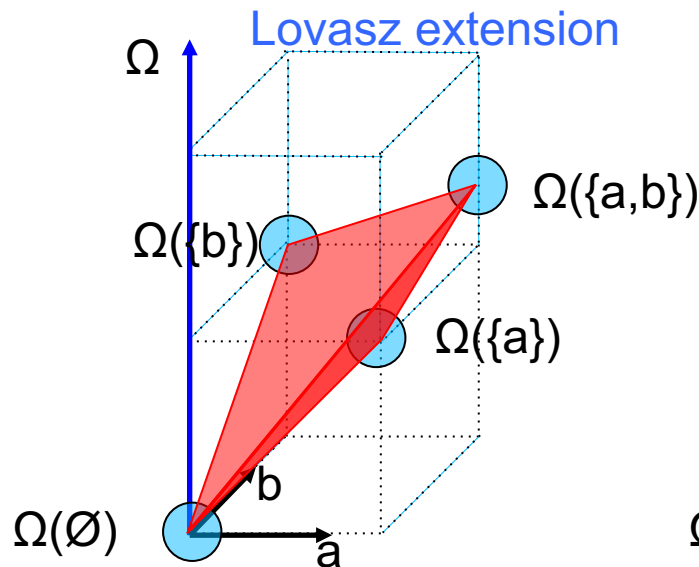
# Extensions vs. envelopes

- **Extension of set function:**
- For any set function  $\Omega$  (submodular or not):  $\Omega^-([w]) = \Omega([w])$
- **Convex closure  $\Omega^-$**  is an **extension** of set function  $\Omega$ , that is, has the same values as  $\Omega$  on the corners of unit cube  $[0,1]^F$
- This is quite remarkable!
  - It's possible because at each dimension, we only have value of  $\Omega$  at two points, 0 and 1, so we can connect them by straight line (which is convex/concave)
  - Imagine we had a function defined over three not two points in every dimension:  $\{-1,0,1\}^F$
  - For most of such functions, convex envelope would not go through values of function at the points  $\{-1,0,1\}^F$
  - Conversely, extension would not be convex
  - E.g. here, extension is not convex, convex envelope is  $f(w)=0$  for any  $w$



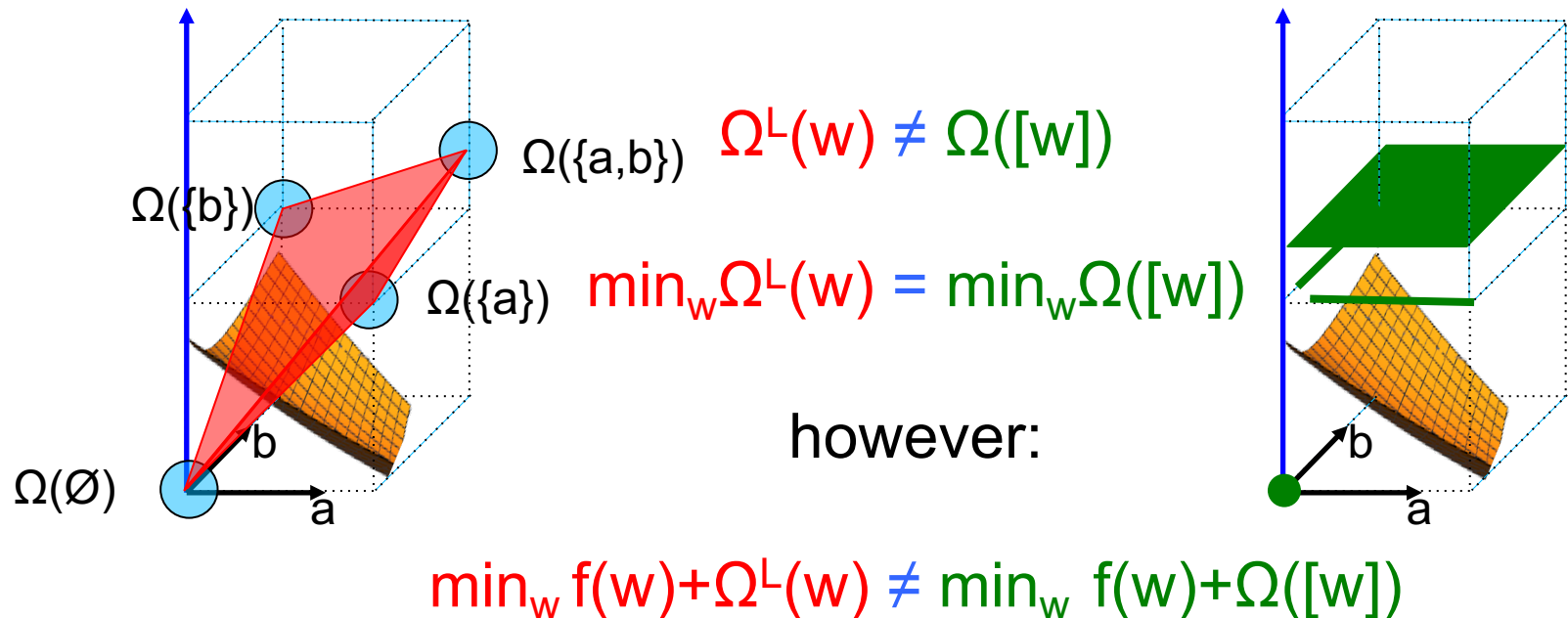
# Submodular minimization

- Given a submodular set function  $\Omega(S)=\Omega(w)$  where  $S=1_w$
- We have a well-defined Lovasz extension  $\Omega^L: [0,1]^n \rightarrow \mathbb{R}$ 
  - For a binary vector  $w$  representing set  $S$ :  $\Omega^L(w)=\Omega(S)$
  - $\Omega^L$  is continuous, convex, piecewise linear
  - $\Omega^L = \Omega^-$  Pointwise highest convex function bounding  $\Omega$  from below
- $\Omega(S)$  and  $\Omega^L(w)$  have the same global minimum
  - And minimum is attained at a corner: at some  $w$  in  $\{0,1\}^F$



# Submodular minimization

- Given a set function  $\Omega(S)=\Omega(w)$  where  $S=1_w$
- We have a well-defined Lovasz extension  $\Omega^L: [0,1]^n \rightarrow \mathbb{R}$



- Consequence:  $\operatorname{argmin}_w R_S(w) + \Omega^L(w) \neq \operatorname{argmin}_w R_S(w) + \Omega([w])$  21

# Overall approach

- For any submodular set function  $\Omega(1_S): \{0,1\}^F \rightarrow \mathbb{R}$   
we can construct its extension  $\Omega^L(w): [0,1]^F \rightarrow \mathbb{R}$ 
  - And then define  $\Omega^L(w): \mathbb{R}^F \rightarrow \mathbb{R}$
- $\operatorname{argmin}_w R_S(w) + \Omega^L(w) \neq \operatorname{argmin}_w R_S(w) + \Omega([w])$ 
  - But still,  $\Omega^L(w)$  captures some aspects of  $\Omega$
  - E.g. **Penalty =  $L_1$  norm** vs. **Penalty = # of features used**
- Unlike  $\Omega([w])$ , the extension  $\Omega^L(w)$  is continuous, convex (though typically non-differentiable)
- We have tool for dealing with non-differentiable, convex terms:
  - Proximal gradient descent, we need:  $\operatorname{argmin}_w \Omega^L(w) + b||w - v||^2$
- So, we can solve problems of the form:
  - **Differentiable risk**  
**+ convex extension of a submodular set function  $\Omega$**



# Big picture

---

- We can solve problems of the form:
  - **Differentiable risk**  
+ **convex extension of a submodular set function  $\Omega$**
- What are some interesting submodular set functions?
  - What are their extensions?
  - Can we solve proximal operator for them?