

# Submodular Functions and Machine Learning

MLSS Kyoto Stefanie Jegelka MIT

### **Set functions**



$$F: 2^{\mathcal{V}} \to \mathbb{R}$$

$$F \left( \begin{array}{c} & & \\ & &$$

We will assume:

• 
$$F(\emptyset) = 0$$

• black box "oracle" to evaluate F

cost of buying items together, or utility, or probability, ... 

### **Discrete Labeling**







Phir

#### **Summarization**



F(S) = relevance + diversity or coverage

Phir

### **Informative Subsets**







- where put sensors?
- which experiments?
- summarization

F(S) = "information"

Mir

## **Sparsity**





Phir





# Formalization

• Formalization: Optimize a set function F(S) (under constraints)



- generally very hard ☺
- submodularity helps: efficient optimization & inference with guarantees!
   ③

# Roadmap

- Submodular set functions
  - definition & basic properties
  - links to convexity
  - special polyhedra
- Minimizing submodular functions coherence, regularization, convexity
- Maximizing submodular functions diversity, repulsion, concavity

## Sensing



V = all possible locations F(S) = information gained from locations in S Mir

# Marginal gain

- Given set function  $F: 2^V \to \mathbb{R}$
- Marginal gain:  $F(s|A) = F(A \cup \{s\}) F(A)$



new sensor s

# **Diminishing marginal gains**



Plii

## **Submodularity**



diminishing marginal costs

#### **Submodular set functions**

- Diminishing gains: for all  $A \subseteq B$ 
  - A +• e  $F(A \cup e) - F(A) \ge F(B \cup e) - F(B)$

• Union-Intersection: for all  $S, T \subseteq \mathcal{V}$ 



# **Supermodular set functions**

• Submodularity: diminishing marginal gains

$$F(A \cup e) - F(A) \ge F(B \cup e) - F(B)$$



Supermodularity: increasing marginal gains

 $F(A \cup e) - F(A) \leq F(B \cup e) - F(B)$ 



 $\emptyset$ 



#### The big picture

Plíř



#### **Examples**

l li î î

• each element e has a weight w(e)

$$F(S) = \sum_{e \in S} w(e)$$

#### $A \subset B$

$$F(A \cup e) - F(A) = w(e) = F(B \cup e) - F(B) = w(e)$$

linear / modular function always submodular!

#### **Examples**



sensing: F(S) = information gained from locations S luir.

#### **Example: cover**

Mir



 $F(A \cup v) - F(A) \ge F(B \cup v) - F(B)$ 

#### More complex model for sensing



Y<sub>s</sub>: temperature at location s

Plii

X<sub>s</sub>: sensor value at location s

 $X_s = Y_s + noise$ 

Joint probability distribution  $P(X_1,...,X_n,Y_1,...,Y_n) = P(Y_1,...,Y_n) P(X_1,...,X_n | Y_1,...,Y_n)$ Prior Likelihood

# **Sensor placement**

Utility of having sensors at subset A of all locations

$$F(A) = H(\mathbf{Y}) - H(\mathbf{Y} \mid \mathbf{X}_A) = I(\mathbf{Y}; \mathbf{X}_A)$$

Uncertainty about temperature Y **before** sensing Uncertainty about temperature Y after sensing



A={1,2,3}: High value F(A)



A={1,4,5}: Low value F(A)

# **Information gain**

 $X_1,\ldots X_n,Y_1,\ldots,Y_m$  discrete random variables

$$F(A) = I(Y; X_A) = H(X_A) - H(X_A|Y)$$
modular

if all  $X_i, X_j$  conditionally independent given Ythen F is submodular! (Exercise: complete the proof)



# Entropy

 $X_1,\ldots,X_n \quad \text{discrete random variables}$   $F(S)=H(X_S)=\text{ joint entropy of variables indexed by }S$   $A\subset B$ 

 $\begin{aligned} H(X_{A\cup e}) - H(X_A) &= H(X_e | X_A) \\ &\leq H(X_e | X_B) \quad \text{``information never hurts''} \\ &= H(X_{B\cup e}) - H(X_B) \end{aligned}$ 

#### discrete entropy is submodular!

# Submodularity and independence

 $X_1, \ldots, X_n$  discrete random variables

 $X_i, i \in S$  statistically independent  $\Leftrightarrow$  H is modular/linear on S  $H(X_S) = \sum_{e \in S} H(X_e)$ 

Similarly: linear independence



## **Maximizing Influence**

F(S) =expected # infected nodes



 $F(S \cup s) - F(S) \ge F(T \cup s) - F(T)$ 

(Kempe, Kleinberg & Tardos 2003)

#### **Graph cuts**



- cut of one edge is submodular!
- large graph: sum of edges

Useful property: sum of submodular functions is submodular

# **Types of submodular functions**

- monotone increasing and integer-valued
   rank functions
- monotone increasing

 $A \subseteq B \implies F(A) \le F(B)$ 

- coverage
- entropy
- spread
- general (non-monotone)
  - graph cuts

#### **Sets and boolean vectors**

any set function with |V| = n.

 $F: 2^V \to \mathbb{R}$ 

... is a function on binary vectors!

$$F: \{0,1\}^n \to \mathbb{R}$$



subset selection = binary labeling!

#### **Attractive potentials**





$$\max_{\mathbf{x}\in\{0,1\}^n} \begin{array}{c|c} P(\mathbf{x} \mid \mathbf{z}) \propto \exp(-E(\mathbf{x}; \mathbf{z})) \\ & & & & \\ \text{labels pixel} \\ & & \text{values} \end{array} \Leftrightarrow \min_{\mathbf{x}\in\{0,1\}^n} E(\mathbf{x}; \mathbf{z}) \end{array}$$

Phir

#### **Attractive potentials**



$$E(\mathbf{x};\mathbf{z}) = \sum_{i} E_{i}(x_{i}) + \sum_{ij} E_{ij}(x_{i}, x_{j})$$

spatial coherence:

$$E_{ij}(1,0) + E_{ij}(0,1) \geq E_{ij}(0,0) + E_{ij}(1,1)$$

$$i j \qquad i j \qquad i j$$

$$S = \{i\} \qquad T = \{j\} \qquad S \cap T = \emptyset \qquad S \cup T$$

$$F(S) + F(T) \geq F(S \cup T) + F(S \cap T)$$

29

### **Diversity priors**





luir.

#### $P(S \mid \text{data}) \propto P(S) P(\text{data} \mid S)$

"spread out"

## **Determinantal point processes**



- normalized similarity matrix K
- sample *Y*:

$$P(S \subseteq Y) = \det(K_S)$$

$$P(e_i \in Y) = K_{ii}$$

$$P(e_i, e_j \in Y) = K_{ii}K_{jj} - K_{ij}^2$$

$$= P(e_i \in Y)P(e_j \in Y) - K_{ij}^2$$
repulsion

 $F(S) = \log \det(K_S)$  is submodular

# **Diversity priors**







Phir



(Kulesza & Taskar 10)

### Submodularity and machine learning

distributions over labels, sets log-submodular/ supermodular probability e.g. "attractive" graphical models, determinantal point processes

> submodularity in machine learning!

(convex) regularization submodularity: "discrete convexity" e.g. combinatorial sparse estimation diffusion processes, covering, rank, connectivity, entropy, economies of scale, summarization, ... submodular phenomena

#### **Closedness properties**

 ${\cal F}(S)$  submodular on V. The following are submodular:

• Restriction:  $F'(S) = F(S \cap W)$ 



#### **Closedness properties**

 ${\cal F}(S)$  submodular on V. The following are submodular:

- Restriction:  $F'(S) = F(S \cap W)$
- Conditioning:  $F'(S) = F(S \cup W)$



#### **Closedness properties**

 ${\cal F}(S)$  submodular on V. The following are submodular:

- Restriction:  $F'(S) = F(S \cap W)$
- Conditioning:  $F'(S) = F(S \cup W)$
- Reflection:  $F'(S) = F(V \setminus S)$




## Submodularity ...







... or concavity?

l li î î

#### **Concave aspects**

- submodularity:
- $A \subseteq B, s \notin B$ :  $F(A \cup s) - F(A) \ge F(B \cup s) - F(B)$ +• s В +• s
- concavity:
  - $a \le b, \ s > 0$ :  $f(a+s) - f(a) \ge f(b+s) - f(b)$



## Submodularity and concavity

- suppose  $g: \mathbb{N} \to \mathbb{R}$  and F(A) = g(|A|)
  - F(A) submodular if and only if ... g is concave



## **Minimum of concave functions**



# Minimum of submodular functions

 $F(A) = \min\{ F_1(A), F_2(A) \}$  submodular?

		F <sub>1</sub> (A)	F <sub>2</sub> (A)	
$A \cap B$	{}	0	0	$A \cap B$
A	{a}	1	0	A
B	{b}	0	1	B
$A \cup B$	{a,b}	1	1	$A\cup B$

 $min(F_1,F_2)$  not submodular in general!

## Convex functions (Lovász, 1983)

- "occur in many models in economy, engineering and other sciences", "often the only nontrivial property that can be stated in general"
- preserved under many operations and transformations: larger effective range of results
- sufficient structure for a "mathematically beautiful and practically useful theory"
- efficient minimization

"It is less apparent, but we claim and hope to prove to a certain extent, that a similar role is played in discrete optimization by *submodular set-functions*" [...] they share the above four properties.

#### **Convex aspects**

- convex extension
- duality results
- poly-time minimization



## Maximum of submodular functions

•  $F_1(A), F_2(A)$  submodular. What about





 $\max\{F_1, F_2\}$  not submodular in general!

## Roadmap



- Submodular set functions
  - definition & basic properties
  - links to convexity
  - special polyhedra
- Minimizing submodular functions
- Maximizing submodular functions

#### **Submodularity and convexity**

any set function with |V| = n.

 $F: 2^V \to \mathbb{R}$ 

... is a function on binary vectors!

$$F: \{0,1\}^n \to \mathbb{R}$$



subset selection = binary labeling!

#### **Relaxation: idea**







Phir

#### A relaxation (extension)



#### **Examples**

$$f(x) = \sum_{i=1}^{k} \alpha_i F(S_i)$$

• truncation

 $F(S) = \max\{|S|, 1\} \qquad f(x) = 0.5 + 0.5 = \max_{i} x_{i}$ 

• cut function 1 2 
$$f(x) = 0.5 \cdot 0 + (1 - 0.5) \cdot 1$$
  
 $F(S) = \begin{cases} 1 & \text{if } S = \{1\}, \{2\} \\ 0 & \text{if } S = \emptyset, \{1, 2\} \end{cases}$   $= |x_1 - x_2|$   
"total variation"!

#### **Alternative characterization**

$$f(x) = \sum_{i=1}^{k} \alpha_i F(S_i)$$

if *F* is submodular, this is equivalent to:

$$f(x) = \max_{y \in \mathcal{B}_F} y^\top x$$



**Theorem** (Lovász, 1983) Lovász extension is convex  $\Leftrightarrow$  F is submodular. **Theorem** (Lovász, 1983) Lovász extension is convex  $\Leftrightarrow$  F is submodular.

 $f(x) = \sum_{i=1}^{k} \alpha_i F(S_i)$ 

 $\leftarrow$  If F is submodular, then f is the max of linear functions.



convexity and positive homogeneity implies:

$$f(1_A + 1_B) \leq f(1_A) + f(1_B)$$

*Exercise:* this implies that F is submodular

#### **Alternative characterization**

$$f(x) = \sum_{i=1}^{k} \alpha_i F(S_i)$$

if *F* is submodular, this is equivalent to:

$$f(x) = \max_{y \in \mathcal{B}_F} y^\top x$$



**Theorem** (Lovász, 1983) Lovász extension is convex  $\Leftrightarrow$  F is submodular.

## Submodular polyhedra

submodular polyhedron:



## **Base polytopes**

Phir

Base polytope  $\mathcal{B}_F = \{y \in \mathcal{P}_F \mid y(\mathcal{V}) = F(\mathcal{V})\}$ 



## **Base polytope**

$$\mathcal{B}_F = \left\{ y \in \mathbb{R}^n \mid \sum_{a \in S} y_a \le F(S) \text{ for all } S \subseteq \mathcal{V} \\ y(\mathcal{V}) = F(\mathcal{V}) \right\}$$
exponentially many constraints!

$$f(x) = \max_{y \in \mathcal{B}_F} \ y^\top x$$



basis for submodular minimization!



## **Optimization over base polytope**

$$\mathcal{B}_{F} = \left\{ y \in \mathbb{R}^{n} \mid \sum_{a \in S} y_{a} \leq F(S) \\ y(\mathcal{V}) = F(\mathcal{V}) \right\}$$
Edmonds' greedy algorithm:  
1. sort  

$$y_{2} \uparrow y_{2} = F(\{e_{1}, e_{2}\}) - F(e_{1})$$

$$x_{\pi(1)} \geq x_{\pi(2)} \geq \dots \geq x_{\pi(n)}$$
2. chain of sets  

$$S_{0} = \emptyset,$$
2. chain of sets  

$$S_{0} = \emptyset,$$

$$S_{1} = \{\pi(1)\} \dots$$

$$S_{i} = \{\pi(1), \dots, \pi(i)\}$$

## **Base polytope**

$$f(x) = \sum_{i=1}^{k} \alpha_i F(S_i)$$

$$f(x) = \max_{y \in \mathcal{B}_F} \ y^\top x$$

#### Remarks:

 chain of sets same as before



Edmonds' greedy algorithm: 1. sort

$$x_{\pi(1)} \ge x_{\pi(2)} \ge \ldots \ge x_{\pi(n)}$$

2. chain of sets  $S_0 = \emptyset, S_i = \{\pi(1), \dots, \pi(i)\}$ 

3. assign values  $y_{\pi(i)} = F(S_i) - F(S_{i-1})$ 

#### **Re-computing our examples**

$$\begin{array}{rcrcr} x \\ \hline \textbf{0.5} \\ 1.0 \end{array} = & 0.5 & \boxed{\begin{array}{c} 1.0 \\ 1.0 \end{array}} + & 0.5 & \boxed{\begin{array}{c} 0 \\ 1.0 \end{array}} \\ \hline 1.0 \end{array}$$

 $F(S) = \max\{|S|, 1\}$ 

sort:  $x_2 \ge x_1 \implies S_1 = \{2\}, S_2 = \{2, 1\}$ 

$$y_2 = F(2) = 1$$
  
 $y_1 = F(2, 1) - F(2) = 1 - 1 = 0$ 

$$f(x) = y^{\top} x = 1 \cdot x_1 + 0 \cdot x_2 = \max_i x_i$$

in general:  $F(S_i) - F(S_{i-1}) > 0$  only for i=1!

$$\Rightarrow f(x) = \max_i x_i$$

### **Re-computing our examples**

$$\begin{array}{cccc} x \\ \hline 0.5 \\ 1.0 \end{array} = & 0.5 \\ \hline 1.0 \\ \hline 1.0 \end{array} + & 0.5 \\ \hline 0 \\ 1.0 \\ \hline 1.0 \\ \hline$$

sort: 
$$x_2 \ge x_1 \implies S_1 = \{2\}, S_2 = \{2, 1\}$$
  
 $y_2 = F(2) = 1$   
 $y_1 = F(2, 1) - F(2) = 0 - 1 = -1$ 

$$f(x) = y^{\top}x = -0.5 + 1 = |x_1 - x_2|$$

### **Convex relaxation**



$$\min_{S \subseteq \mathcal{V}} F(S) = \min_{x \in \{0,1\}^n} F(x) \longrightarrow \min_{x \in [0,1]^n} f(x)$$

1. relaxation: convex optimization (non-smooth)

2. relaxation is exact!

→ submodular minimization in polynomial time! (Grötschel, Lovász, Schrijver 1981)

## Amazing base polytopes

• linear optimization = f(x)Edmonds' greedy algorithm: each vertex determined by a permutation!

$$f(x) = \max_{y \in \mathcal{B}_F} \ y^\top x$$

#### Base polytopes almost everywhere:

- cores of games (Shapley)
- information theory: achievable rates for lossless coding of correlated sources (Slepian-Wolf, Cover, Fujishige)
- matroids



## Base polytopes and spanning trees



ground set: all edges

indicator vectors  $1_T$  of all spanning trees T

 $\mathcal{B} = \text{ convex hull of all tree indicator vectors}$ 

is a base polytope

What is the submodular function? F(S) = size of largest tree within S $= \max\{|T|: T \text{ is a tree and } T \subseteq S\}$ 



## **Greedy algorithm for trees**

 $\max_{y \in \mathcal{B}} y^{\top} x$ 



F(S) = size of largest tree within S

vector x: weight for

each edge

greedy algorithm: y indicator vector of a tree

go through edges in order of their weight

 $e_1 \ e_2 \ e_3 \ e_4 \ e_5 \ e_6$ 

• if edge does not complete a cycle:

 $F(S_i) - F(S_{i-1}) = 1$   $y_i = 1$  "pick i"

if it does:

$$F(S_i) - F(S_{i-1}) = 0$$
  $y_i = 0$  "don't pick i"

finds the maximum weight spanning tree (Kruskal's algorithm)

## **General: Matroids**

Matroid  $\mathcal{M} = (\mathcal{V}, \mathcal{I})$ 

- ground set  $\ \mathcal{V}$
- family of independent sets  $\mathcal{I}$

## Matroids (semi-formally)

S is independent ( = feasible) if ...



• S independent  $\rightarrow$  T  $\subseteq$  S also independent

## Matroids

S is independent (=feasible) if ...



- S independent  $\rightarrow$   $T \subseteq$  S also independent
- Exchange property: S, U independent, |S| > |U|
   → some e ∈ S can be added to U: U ∪ e independent

## **General: Matroids**

Matroid  $\mathcal{M} = (\mathcal{V}, \mathcal{I})$ 

- ground set  $\mathcal{V}$
- family of independent sets  $\mathcal{I}$
- rank function:

 $F(S) = \max\{ |T| : T \subseteq S \text{ and } T \in \mathcal{I} \}$ 

always submodular and increasing

• another special case: matrix rank

### **Convex relaxation**



$$\min_{S \subseteq \mathcal{V}} F(S) = \min_{x \in \{0,1\}^n} F(x) = \min_{x \in [0,1]^n} f(x)$$

1. relaxation: convex optimization (non-smooth)

2. relaxation is exact!

→ submodular minimization in polynomial time! (Grötschel, Lovász, Schrijver 1981)

## Minimizing the Lovász extension

$$\min_{x \in [0,1]^n} f(x)$$

- subgradient method
- combinatorial algorithms: dual

#### **Subgradients**



## subgradient: $g_x \in \arg \max_y y^\top x$ s.t. $y \in \mathcal{B}_F$

## **Projected subgradient method**



### Convergence

#### Theorem

Let  $D = \sqrt{n}$  and  $L = \max_{g \in \mathcal{B}_F} ||g|| \le 3 \max_S |F(S)|$ . With step size  $\alpha_t = \frac{D}{L\sqrt{t}}$ , the error decreases as

$$\min_{\tau \le t} f(x^{\tau}) - f(x^*) \le \frac{4DL}{\sqrt{t}}$$

- D: diameter of [0,1]<sup>n</sup>
  - L: Lipschitz constant
- for error  $\leq \epsilon$  need  $O(\frac{1}{\epsilon^2})$  iterations
# **Submodular minimization**

#### convex optimization

- ellipsoid method (Grötschel-Lovasz-Schrijver 81)
- subgradient method
- minimum-norm point / Fujishige-Wolfe algorithm

### combinatorial methods

- first polynomial-time: (Schrijver 00, Iwata-Fleischer-Fujishige 01)
- currently fastest:

 $O(n^4T + n^5\log M)$  (Iwata 03)

 $O(n^6 + n^5 T) \qquad \qquad \text{(Orlin 09)}$ 

ult:  $O(n^2 T \log nM + n^3 \log^c nM)$  $O(n^3 T \log^2 n + n^4 \log^c n)$  (Lee-Sidford-Wong 15)

## **Convex duality**

l li î î

.

.

$$\min_{S \subseteq \mathcal{V}} F(S) = \min_{x \in [0,1]^n} f(x)$$
  
$$= \min_{x \in [0,1]^n} \max_{y \in \mathcal{B}_F} y^\top x$$
  
$$= \max_{y \in \mathcal{B}_F} \min_{x \in [0,1]^n} x^\top y \qquad = \max_{y \in \mathcal{B}_F} \left( \sum_{i=1}^n \min\{y_i, 0\} \right)$$

Optimality conditions:  $(S^*, y^*)$  optimal primal-dual pair if

1. 
$$y^* \in \mathcal{B}_F$$
  
2.  $\{y^* < 0\} \subseteq S^* \subseteq \{y^* \le 0\}$   $\checkmark$   
3.  $y^*(S^*) = F(S^*)$ 

## **Combinatorial algorithms**

solve 
$$\max_{y \in \mathcal{B}_F} \left( \sum_{i=1}^n \min\{y_i, 0\} \right)$$

- remove "negative mass"
- challenges:
  - need to stay in polytope:  $\sum y_i \leq F(S)$
  - cannot test feasibility
  - → network flow algorithms







# **Submodular minimization**

#### convex optimization

- ellipsoid method (Grötschel-Lovasz-Schrijver 81)
- subgradient method
- minimum-norm point / Fujishige-Wolfe algorithm

### combinatorial methods

- first polynomial-time: (Schrijver 00, Iwata-Fleischer-Fujishige 01)
- currently fastest:

 $O(n^4T + n^5\log M)$  (Iwata 03)

 $O(n^6 + n^5 T) \qquad \qquad \text{(Orlin 09)}$ 

ult:  $O(n^2 T \log nM + n^3 \log^c nM)$  $O(n^3 T \log^2 n + n^4 \log^c n)$  (Lee-Sidford-Wong 15)

### Some fun 🙂

- Complete the proof that Lovasz extension convex → set function is submodular (slide 51)
- 2. Submodular oder not? Let F be increasing and submodular, and define  $G(S) = \min\{F(S), c\}$ for a constant c.

slides and pointers to literature: people.csail.mit.edu/stefje/mlss

## **Example: costs**



### **Example: costs**

Phir



## **Shared fixed costs**



marginal cost: #new shops + #new items

decreasing  $\rightarrow$  cost is submodular!

- shops: shared fixed cost
- economies of scale