## Efficient Minimization of Decomposable Submodular Functions

Peter Stobbe and Andreas Krause

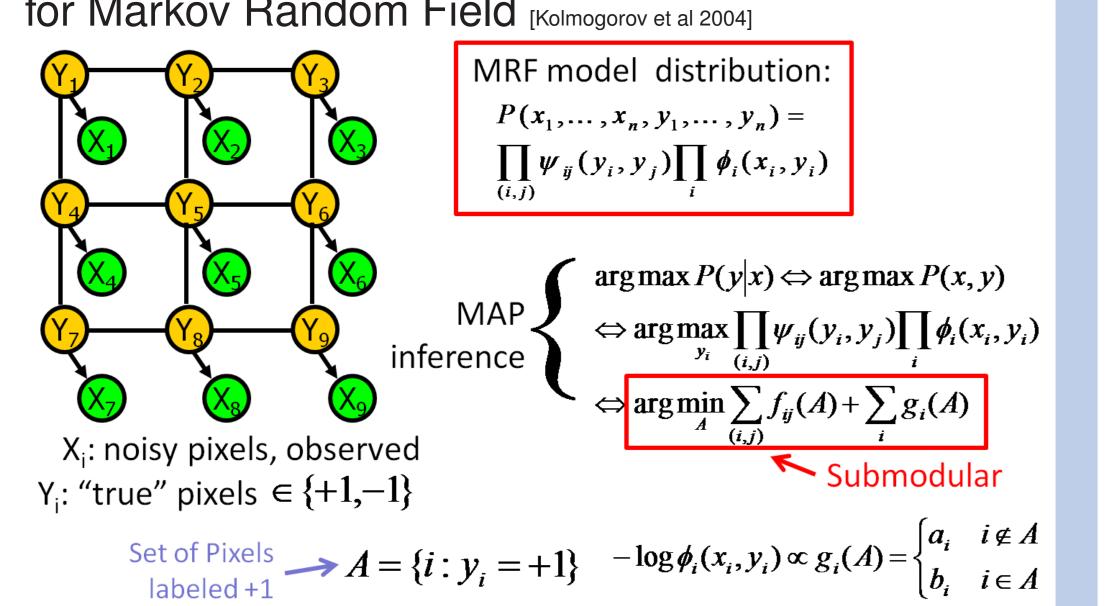
#### Structures in Optimization

- Convexity useful for continuous functions
- $f(\mathbf{x} + \theta \mathbf{h}) f(\mathbf{x}) \leq \theta \left( f(\mathbf{x} + \mathbf{h}) f(\mathbf{x}) \right)$
- Minimization tractable if convex.
- ► Similar Submodular discrete functions:
- Domain of f: subsets of finite set E
- $f(A \cup B \cup C) f(A \cup B) \leq f(A \cup C) f(A)$
- Minimization tractable if submodular.

#### Submodular Minimization Examples

Many important Machine Learning problems require  $A^* \in \text{arg min}_{A \subset E} f(A)$ 

Maximum A Posteri Inference of Hidden Variables for Markov Random Field [Kolmogorov et al 2004]



Factorizing random

variables.



 $c_{ij} \{i\} \cap A \neq \{j\} \cap A$ Mutual Information is submodular:  $f(A) = I(X_A; X_{E \setminus A})$ 

 $\begin{cases}
0 & \{i\} \cap A = \{j\} \cap A
\end{cases}$ 

#### Algorithms

- ► General case:  $O^*(n^5)$  function evaluations. [Iwata Orlin 2009]
- ► Min-norm algorithm. Often practical, unknown

Only pairwise  $-\log \psi_{ij}(y_i, y_j) \propto f_{ij}(A) = \int_{A}^{A} dx$ 

- complexity. [Fujishige et al]
- More efficient special cases:
- Pairwise potentials
- ▶ eg. MAP for Ising model.
- ► Fast mincut algorithms O\*(n²)
- Queyranne's algorithm. ▶ Only symmetric functions  $f(A) = f(E \setminus A)$ .
- ► Running time O\*(n³)
- ► Sum of Submodular Functions [Kolmogorov 2010]
- ► Each term in sum must be relatively low-order (function of few
- Our work: Decomposable functions!

#### Decomposable submodular functions

Definition:

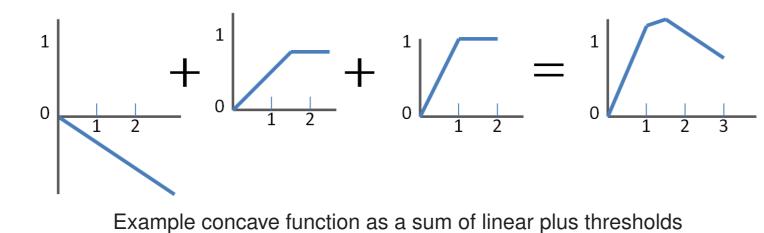
$$f(A) = \sum_{j} \phi_{j}(\mathbf{w}_{j} \cdot \mathbf{e}_{A})$$

 $\phi_j$  concave,  $\mathbf{w}_j \geq \mathbf{0}$ , and  $\mathbf{e}_A[k] = \mathbf{0}$ 

► Key example - Threshold Potentials:

$$f(A) = \min(y, \boldsymbol{w} \cdot \boldsymbol{e}_A)$$

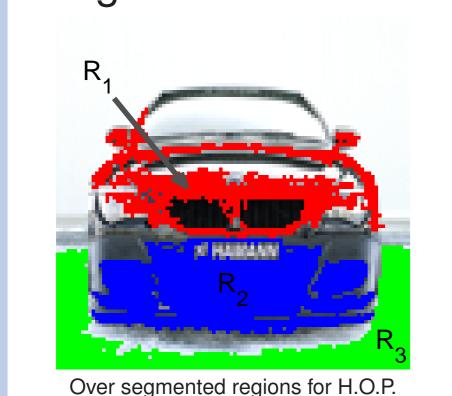
► All decomposable functions can be expressed as a modular part plus a sum/integral of threshold potentials.

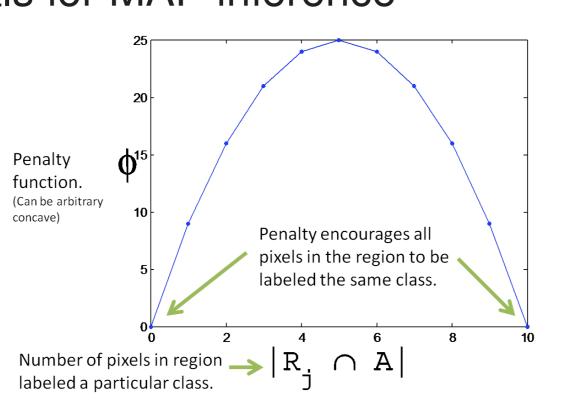


► Concave cardinality functions: A strict subclass of decomposable functions:

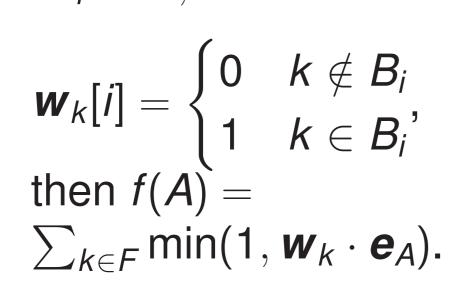
$$f(A) = \sum_{i} \phi_{j}(|R_{j} \cap A|)$$

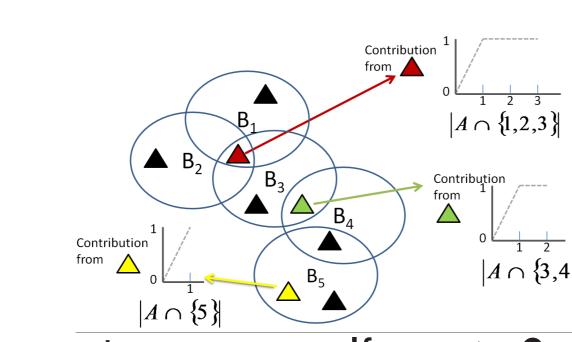
Higher Order Potentials for MAP inference





▶ Set cover functions:  $f(A) = |\cup_{i \in A} B_i|$  where  $B_i \subset F, \forall i \in E$ . Let





▶ Example from queuing systems [Itoko 2007]. If  $u, v \ge 0$ ,  $\phi$ nonincreasing concave

$$f(A) = (oldsymbol{u} \cdot oldsymbol{e}_A) \phi(oldsymbol{v} \cdot oldsymbol{e}_A)$$

#### Overview of method and contributions

- ▶ Reformulate as (nonsmooth) convex minimization
- Use modern technique of smooth minimization of nonsmooth functions
- Novel stopping criterion

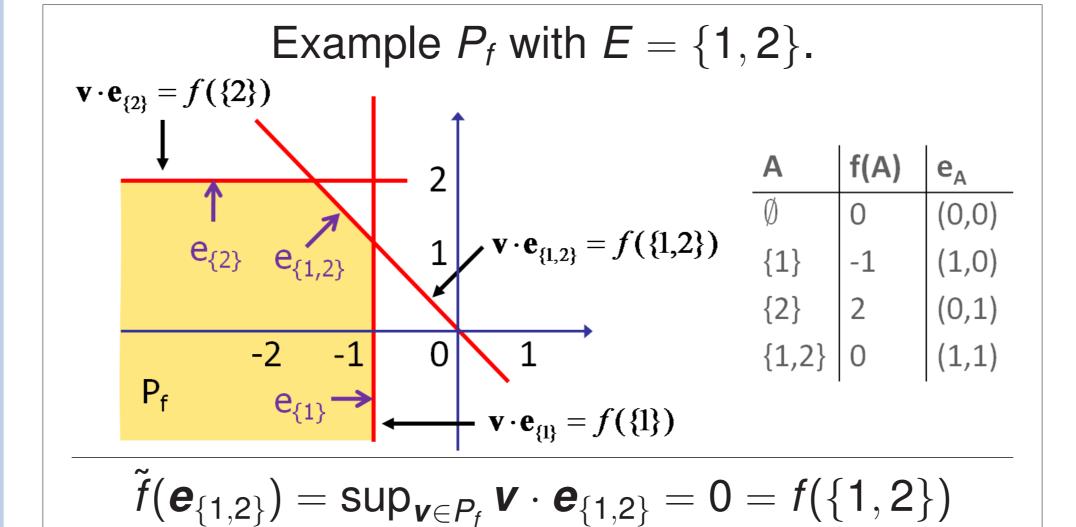
Able to solve problems with 10,000 variables in a minute.

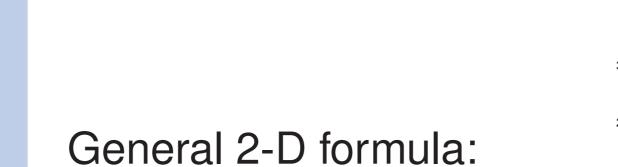
#### **Convex reformulation**

- ▶ Key Properties of *Lovász extension*:  $\tilde{f}:[0,1]^n \to \mathbb{R}$
- Convex
- Agrees with f at corners:  $\tilde{f}(\boldsymbol{e}_A) = f(A)$ ▶ A corner is optimal:  $\{e_A : A \subset E\} \cap \text{arg min}_{x \in [0,1]^n} \tilde{f}(x) \neq \emptyset$ .
- ▶ Definition, assuming  $f(\emptyset) = 0$ .

$$\widetilde{f}(\boldsymbol{x}) = \sup_{\boldsymbol{v} \in P_f} \boldsymbol{v} \cdot \boldsymbol{x}.$$

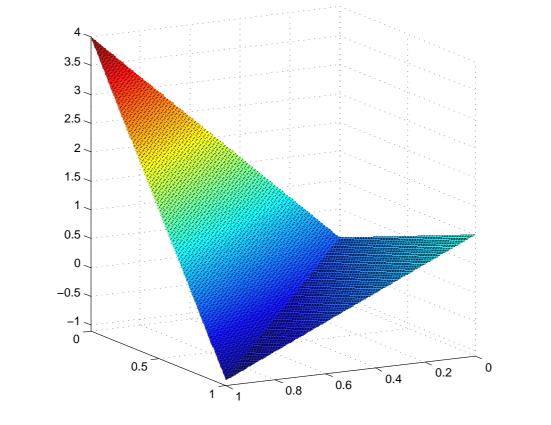
 $P_f = \{ \boldsymbol{v} \in \mathbb{R}^n : \boldsymbol{v} \cdot \boldsymbol{e}_A \leq f(A), \text{ for all } A \in 2^E \}.$  $P_f$  is Submodular polyhedron associated with f.





 $f(x_1,x_2) = a|x_1-x_2| +$ 

 $bx_1 + cx_2$  with  $a \ge 0$ .



#### **Smooth Minimization of Nonsmooth Functions**

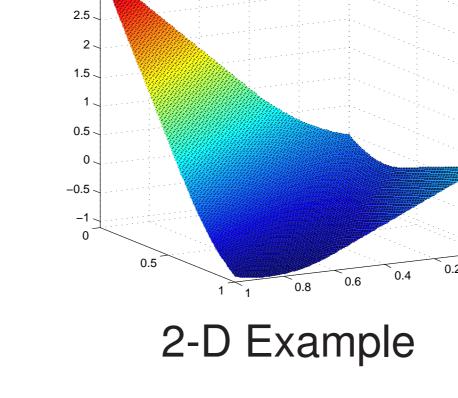
- Groundbreaking work by Nesterov [2004]
- Solves nonsmooth problems in  $\mathcal{O}\left(\frac{1}{\varepsilon}\right)$  iterations
- ► Each iteration neglibly more work than gradient descent
- Not black-box solver; requires problem to have exploitable structure (often true)
- ▶ Example:  $h(x) = \sup_{v \in C} x \cdot y$ . If C easy to project onto, then h can be smoothed and an accelerated gradient descent scheme can be applied.

#### **Smoothed Lovász Extension**



Computing General Smoothed Lovász Gradient ↔ Submodular Minimization Problem. [Bach 2010] (Just as hard as

orginal problem.)



Key Insight: Smoothed Lovász Gradient for a decomposable function is easily computed!

- combinations of threshold potentials, here assume  $f(A) = \min(\mathbf{w} \cdot e_A, y)$
- For **x** in unit cube:

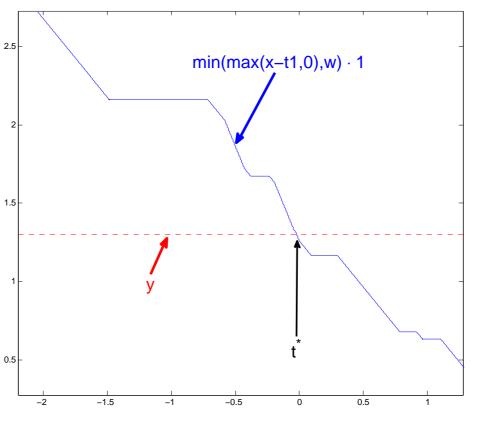
$$\tilde{f}(x) = \max_{\mathbf{0} \leq v \leq w, \ v \cdot \mathbf{1} = y} v \cdot x$$

Smoothed gradient:

$$abla ilde{f}^{\mu}(\mathbf{x}) = \underset{\mathbf{0} \leq \mathbf{v} \leq \mathbf{w}, \ \mathbf{v} \cdot \mathbf{1} = \mathbf{y}}{\arg \min} \|\mathbf{v} - \mathbf{x}/\mu\|$$

$$= \min(\max((\mathbf{x} - t^*\mathbf{1})/\mu, \mathbf{0}), \mathbf{w}).$$

Find root of monotonic linear function. No explicit closed form to compute.



#### Smoothed Lovász Gradient (SLG) Algorithm

▶ Input: Accuracy  $\varepsilon$ ; decomposable function

- For t = 0, 1, 2, ...
- $ho oldsymbol{g}_t = 
  abla \widetilde{f}^{\mu}(oldsymbol{x}_{t-1})/L$
- $ullet oldsymbol{z}_t = P_{[0,1]^n} \left( oldsymbol{z}_{-1} \sum_{s=0}^t \left( rac{s+1}{2} 
  ight) oldsymbol{g}_s 
  ight)$
- $\boldsymbol{y}_t = P_{[0,1]^n}(\boldsymbol{x}_t \boldsymbol{g}_t)$
- ▶ If gap<sub>t</sub>  $\leq \varepsilon/2$  stop.
- $ightharpoonup m{x}_t = (2m{z}_t + (t+1)m{y}_t)/(t+3)$
- $ullet oldsymbol{x}_arepsilon = oldsymbol{y}_t$
- ▶ Output:  $\varepsilon$ -optimal  $\boldsymbol{x}_{\varepsilon}$  to min $_{\boldsymbol{x}\in[0,1]^n}\tilde{f}(\boldsymbol{x})$

#### Theorem

SLG is guaranteed to provide an  $\varepsilon$ -optimal solution after running for  $\mathcal{O}(\frac{D}{c})$  iterations.

Given  $\varepsilon$ -optimal  $x_{\varepsilon} \in [0, 1]^n$ , can round to find corresponding  $\varepsilon$ -optimal set:

- ▶ Input:  $\mathbf{x} \in [0, 1]^n$ ; submodular function f(A).
- ▶ Choose  $\sigma$ :  $\boldsymbol{x}[\sigma(1)] \geq \ldots \geq \boldsymbol{x}[\sigma(n)]$ .
- $\triangleright S_k = \{\sigma(1), \ldots, \sigma(k)\}.$
- $ightharpoonup k^* = \operatorname{arg\;min}_k f(S_k), A = S_{k^*}$
- ▶ Output: Set *A* satisfying  $f(A) \leq \tilde{f}(x)$

#### **Early Stopping**

- ▶ By rounding, may find optimal set before continous convergence.
- ▶ Use  $g \in \partial \tilde{f}(e_A)$  to bound optimality gap of A:

$$f(A) - f^* \le \sum_{k \in A} \max(0, \boldsymbol{g}[k]) + \sum_{k \in E \setminus A} \max(0, -\boldsymbol{g}[k])$$

$$If \boldsymbol{g}[k] \le 0 \text{ for } k \in A \text{ and } \boldsymbol{g}[k] \ge 0 \text{ for } k \in E \setminus A, \text{ then }$$

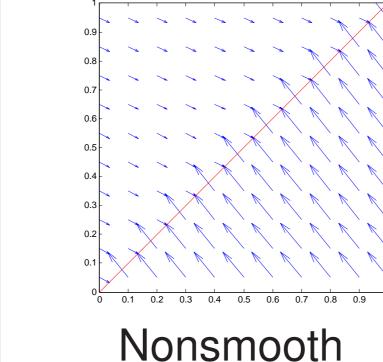
- A is optimal! ▶ Choose  $\mathbf{g} \in \partial \tilde{f}(\mathbf{e}_A)$  from smoothed gradient

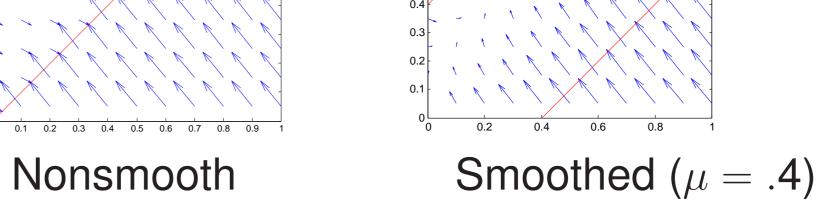
#### Lemma

 $\mathsf{min}_{k\in\mathcal{A},l\in\mathcal{E}\setminus\mathcal{A}}oldsymbol{x}[k]-oldsymbol{x}[l]\geq 2\mu\Rightarrow
abla ilde{f}^{\mu}(oldsymbol{x})\in\partial ilde{f}(oldsymbol{e}_{\mathcal{A}})$ 

#### Example 2-D negative gradients:

of  $A = \{1, 2\}$ 





# SFM3 LEX2

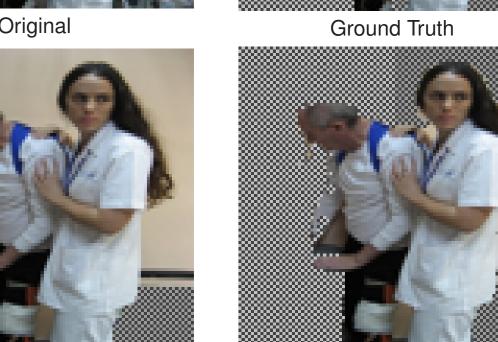
#### **Segmentation Results**

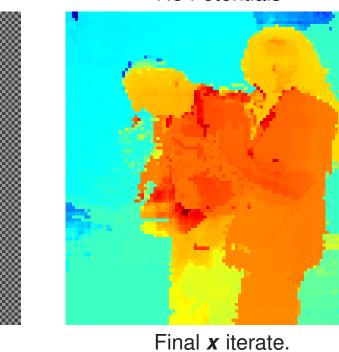
Textonboost classification, regularized with submodular potential functions:

**Synthetic Results Comparision** 



Pairwise Potentials





Minimization of H. O. Potentials takes 70 sec. with SLG algorithm vs. 2 hrs. for min-norm.

#### Conclusions

- A new class of submodular functions that can be efficiently minimized
- Apply Nesterov smoothing technique to Lovász extension
- Novel way to stop early
- Can solve larger-scale problems than previously possible

#### General Lovász Extension Properties

- ▶ Piecewise linear.
- Nonsmooth at points with equal components. (eg.
- Defined by LP with exponentially many constraints.
- ▶ Computable in  $O^*(n)$  time: Sort x components. Choose  $\sigma$ :  $\boldsymbol{x}[\sigma(1)] \geq \ldots \geq \boldsymbol{x}[\sigma(n)]$ . Let  $S_k = \{\sigma(1), \ldots, \sigma(k)\}.$ 
  - $\tilde{f}(\mathbf{x}) = \sum_{k} \mathbf{x}[\sigma(k)](f(S_k) f(S_{k-1})).$
- Can compute extension and also a subgradient in linear time:
- Can use projected subgradient descent, but slow convergence  $\mathcal{O}\left(\frac{1}{\epsilon^2}\right)$  iterations to achieve  $\varepsilon$  accuracy. Impractical.

 $\partial \tilde{f}(\boldsymbol{x}) \ni \sum \boldsymbol{e}_{\sigma(k)}(f(S_k) - f(S_{k-1})).$ 

### Computation of Smoothed Gradient

- ► Since all decomposable functions will be linear

With  $t^*$  chosen so  $\nabla \tilde{f}^{\mu}(\mathbf{x}) \cdot \mathbf{1} = \mathbf{y}$ 

continuous piecewise expression, but simple

