

Submodular Functions, Optimization, and Applications to Machine Learning

— Spring Quarter, Lecture 1 —

http://www.ee.washington.edu/people/faculty/bilmes/classes/ee563_spring_2018/

Prof. Jeff Bilmes

University of Washington, Seattle
Department of Electrical Engineering
<http://melodi.ee.washington.edu/~bilmes>

Mar 26th, 2018



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

$$-f(A) + 2f(C) + f(B), \quad -f(A) + f(C) + f(B), \quad -f(A \cap B)$$





Announcements

- Welcome to: Submodular Functions, Optimization, and Applications to Machine Learning, EE563.
- Class: An introduction to submodular functions including methods for their optimization, and how they have been (and can be) applied in many application domains.
- Weekly Virtual Office Hours: Mondays, 10:00-11:00pm, via zoom (link posted on canvas).
- EEB 042, class web page is at our web page (http://www.ee.washington.edu/people/faculty/bilmes/classes/ee563_spring_2018/).
- Use our discussion board (https://canvas.uw.edu/courses/1039754/discussion_topics) for all questions, comments, so that all will benefit from them being answered.

Rough Class Outline

- Introduction to submodular functions: definitions, real-world and contrived examples, properties, operations that preserve submodularity, inequalities, variants and special submodular functions, and computational properties. Gain intuition, when is submodularity and supermodularity useful?
- Submodularity is an ideal model for **cooperation, complexity, and attractiveness** as well as for **diversity, coverage, & information**
- Applications in **data science**, **computer vision**, **tractable substructures in constraint satisfaction/SAT and graphical models**, **game theory**, **social networks**, **economics**, **information theory**, **structured convex norms**, **natural language processing**, **genomics/proteomics**, **sensor networks**, **probabilistic inference**, and other areas of **machine learning**.



Rough Class Outline (cont. II)

- theory of matroids and lattices.
- Polyhedral properties of submodular functions, polymatroids generalize matroids.
- The Lovász extension of submodular functions, the Choquet integral, and convex and concave extensions.
- Submodular maximization algorithms under constraints, submodular cover problems, greedy algorithms, approximation guarantees.
- Submodular minimization algorithms, a history of submodular minimization, including both numerical and combinatorial algorithms, computational properties, and descriptions of both known results and currently open problems in this area.
- Submodular flow problems, the principle partition of a submodular function and its variants.



Rough Class Outline (cont. III)

- Constrained optimization problems with submodular functions, including maximization and minimization problems with various constraints. An overview of recent problems addressed in the community.



Classic References

- Jack Edmonds's paper "Submodular Functions, Matroids, and Certain Polyhedra" from 1970.
- Nemhauser, Wolsey, Fisher, "A Analysis of Approximations for Maximizing Submodular Set Functions-I", 1978
- Lovász's paper, "Submodular functions and convexity", from 1983.

Useful Books

- Fujishige, "Submodular Functions and Optimization", 2005
- Narayanan, "Submodular Functions and Electrical Networks", 1997
- Welsh, "Matroid Theory", 1975.
- Oxley, "Matroid Theory", 1992 (and 2011).
- Lawler, "Combinatorial Optimization: Networks and Matroids", 1976.
- Schrijver, "Combinatorial Optimization", 2003
- Gruenbaum, "Convex Polytopes, 2nd Ed", 2003.
- Additional readings that will be announced here.

Recent online material (some with an ML slant)

- Previous video version of this class http://j.ee.washington.edu/~bilmes/classes/ee596a_fall_2014/.
- Stefanie Jegelka & Andreas Krause's 2013 ICML tutorial <http://techtalks.tv/talks/submodularity-in-machine-learning-new-directions-part-i/58125/>
- NIPS, 2013 tutorial on submodularity <http://melodi.ee.washington.edu/~bilmes/pgs/b2hd-bilmes2013-nips-tutorial.html> and <http://youtu.be/c4rBof38nKQ>
- Andreas Krause's web page <http://submodularity.org>.
- Francis Bach's updated 2013 text. http://hal.archives-ouvertes.fr/docs/00/87/06/09/PDF/submodular_fot_revised_hal.pdf
- Tom McCormick's overview paper on submodular minimization <http://people.commerce.ubc.ca/faculty/mccormick/sfmchap8a.pdf>
- Georgia Tech's 2012 workshop on submodularity: <http://www.arc.gatech.edu/events/arc-submodularity-workshop>



Facts about the class

- Prerequisites: ideally knowledge in probability, statistics, convex optimization, and combinatorial optimization these will be reviewed as necessary. The course is open to students in all UW departments. Any questions, please contact me.
- Text: We will be drawing from the book by Satoru Fujishige entitled "Submodular Functions and Optimization" 2nd Edition, 2005, but we will also be reading handouts and research papers that will be posted here on this web page, especially for some of the application areas.
- Grades and Assignments: Grades will be based on a combination of a final project (45%), homeworks (55%). There will be between 3-6 homeworks during the quarter.
- Final project: The final project will consist of a 4-page paper (conference style) and a final project presentation. The project must involve using/dealing mathematically with submodularity in some way or another, and might involve a contest!

Facts about the class

- Homework must be submitted electronically using our assignment dropbox (<https://canvas.uw.edu/courses/1216339/assignments>). PDF submissions only please. Photos of neatly hand written solutions, combined into one PDF, are fine
- Lecture slides - are being updated and improved this quarter. They will likely appear on the web page the night before, and the final version will appear just before class.
- Slides from previous version of this class are at http://www.ee.washington.edu/people/faculty/bilmes/classes/ee596b_spring_2016/.



Cumulative Outstanding Reading

- Read chapter 1 from Fujishige's book.

Class Road Map - EE595 Spring 2016

- L1(3/28): Motivation, Applications, & Basic Definitions
- L2(3/30): Machine Learning Apps (diversity, complexity, parameter, learning target, surrogate).
- L3(4/4): Info theory exs, more apps, definitions, graph/combinatorial examples, matrix rank example, visualization
- L4(4/6): Graph and Combinatorial Examples, matrix rank, Venn diagrams, examples of proofs of submodularity, some useful properties
- L5(4/11): Examples & Properties, Other Defs., Independence
- L6(4/13): Independence, Matroids, Matroid Examples, matroid rank is submodular
- L7(4/18): Matroid Rank, More on Partition Matroid, System of Distinct Reps, Transversals, Transversal Matroid,
- L8(4/20): Transversals, Matroid and representation, Dual Matroids,
- L9(4/25): Dual Matroids, Properties, Combinatorial Geometries, Matroid and Greedy
- L10(4/27): Matroid and Greedy, Polyhedra, Matroid Polytopes,
- L11(5/2): From Matroids to Polymatroids, Polymatroids
- L12(5/4): Polymatroids, Polymatroids and Greedy
- L13(5/9): Polymatroids and Greedy; Possible Polytopes; Extreme Points; Polymatroids, Greedy, and Cardinality Constrained Maximization
- L14(5/11): Cardinality Constrained Maximization; Curvature; Submodular Max w. Other Constraints
- L15(5/16): Submodular Max w. Other Constraints, Most Violated \leq , Matroids cont., Closure/Sat,
- L16(5/18): Closure/Sat, Fund. Circuit/Dep,
- L17(5/23): Min-Norm Point and SFM, Min-Norm Point Algorithm,
- L18(5/25): Proof that min-norm gives optimal, Lovász extension.
- L19(6/1):
- L20(6/6): Final Presentations maximization.

Finals Week: June 6th-10th, 2016.

Class Road Map - EE563

- L1(3/26): Motivation, Applications, & Basic Definitions, Apps (diversity, complexity, parameter, learning target, surrogate).
- L2(3/28):
- L3(4/2):
- L4(4/4):
- L5(4/9):
- L6(4/11):
- L7(4/16):
- L8(4/18):
- L9(4/23):
- L10(4/25):
- L11(4/30):
- L12(5/2):
- L13(5/7):
- L14(5/9):
- L15(5/14):
- L16(5/16):
- L17(5/21):
- L18(5/23):
- L-(5/28): Memorial Day (holiday)
- L19(5/30):
- L21(6/4): Final Presentations maximization.

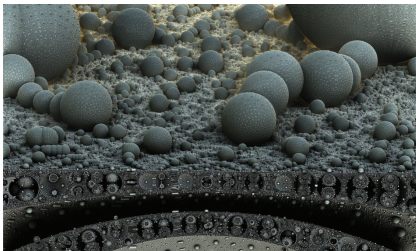
Last day of instruction, June 1st. Finals Week: June 2-8, 2018.

The Ideal Machine Learning Methods

- Simple to define



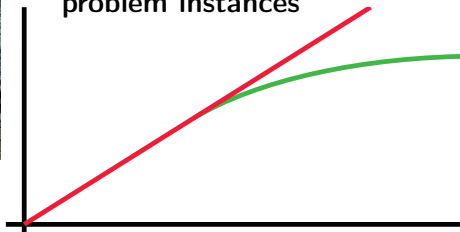
- Mathematically rich



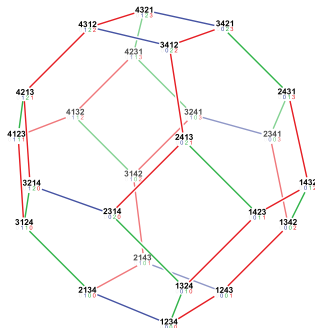
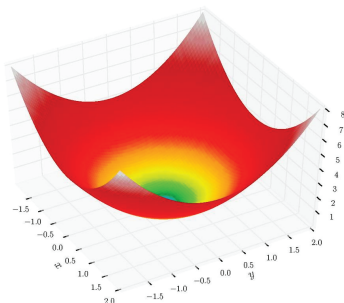
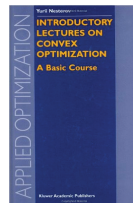
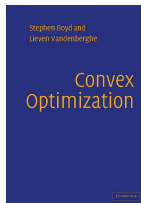
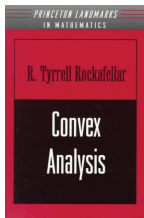
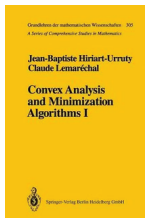
- Naturally suited to many real-world applications



- Efficient & scalable to large problem instances



Convex Analysis in Machine Learning



Successful Convexity in Machine Learning

- Linear and logistic regression, surrogate loss functions.
- Convex sparse regularizers (such as the ℓ_p family and nuclear norms).
- PSD matrices (i.e., positive semidefinite cone) and Gaussian densities.
- Optimizing non-linear and even non-convex classification/regression methods such as support-vector (SVMs) and kernel machines via convex optimization.
- Maximum entropy estimation
- The expectation-maximization (EM) algorithm.
- Ideas/techniques/insight for non-convex methods, convex minimization useful even for non-convex problems, such as Deep Neural Networks (DNNs). Convex analysis for non-convex problems.

A Convexity Limitation: Discrete Problems

Many Machine Learning problems are **inherently discrete**:

- Active learning/label selection.
- MAP & diverse k -best discrete probabilistic inference
- Data Science: data partitioning, clustering, summarization; the science of data management.
- Sparse modeling, compressed sensing, low-rank approximation.
- Probabilistic models: structure learning in graphical models and neural networks. Non-graphical global potentials.
- Variable and feature selection; dictionary selection.
- Natural language processing (NLP): words, phrases, sentences, paragraphs, n -grams, syntax trees, graphs, semantic structures.
- Social choice and voting theory, social networks, viral marketing,
- (Multi-label) image segmentation in computer vision.
- Proteomics: selecting peptides, proteins, drug trial participants
- Genomics: cell-type or assay selection, genomic summarization
- Social networks, influence, viral marketing, information cascades, diffusion networks

Classic Discrete Optimization Problems

- **Operations Research/Industrial Engineering:** facility and factory location, packing and covering.
- **Sensor placement** where to optimally place sensors?
- **Information:** Information theory, sets of random variables.
- **Geometry:** Polytopes and polyhedra
- **Mathematics:** e.g., Monge matrices, efficient dynamic programming, Birkhoff lattice theory
- **Combinatorial Problems:** e.g., sets, graphs, graph cuts, max k coverage, packings, coverings, partitions, paths, flows, matchings, colorings,
- **Algorithms:** Algorithms, and time/space complexity
- **Economics:** markets, economies of scale, mathematics of supply & demand

General Integer Programming (e.g., Integer Linear Programming (ILP), Integer Quadratic Programming (IQP), etc). General case can ignore useful and natural structures common to many problems.

Attractions of Convex Functions

Why do we like Convex Functions? (Quoting Lovász 1983):

- 1 *Convex functions occur in many mathematical models in economy, engineering, and other sciences. Convexity is a very natural property of various functions and domains occurring in such models; quite often the only non-trivial property which can be stated in general.*

Attractions of Convex Functions

Why do we like Convex Functions? (Quoting Lovász 1983):

- 1 *Convex functions occur in many mathematical models in economy, engineering, and other sciences. Convexity is a very natural property of various functions and domains occurring in such models; quite often the only non-trivial property which can be stated in general.*
- 2 *Convexity is preserved under many natural operations and transformations, and thereby the effective range of results can be extended, elegant proof techniques can be developed as well as unforeseen applications of certain results can be given.*

Attractions of Convex Functions

Why do we like Convex Functions? (Quoting Lovász 1983):

- 1 *Convex functions occur in many mathematical models in economy, engineering, and other sciences. Convexity is a very natural property of various functions and domains occurring in such models; quite often the only non-trivial property which can be stated in general.*
- 2 *Convexity is preserved under many natural operations and transformations, and thereby the effective range of results can be extended, elegant proof techniques can be developed as well as unforeseen applications of certain results can be given.*
- 3 *Convex functions and domains exhibit sufficient structure so that a mathematically beautiful and practically useful theory can be developed.*

Attractions of Convex Functions

Why do we like Convex Functions? (Quoting Lovász 1983):

- 1 *Convex functions occur in many mathematical models in economy, engineering, and other sciences. Convexity is a very natural property of various functions and domains occurring in such models; quite often the only non-trivial property which can be stated in general.*
- 2 *Convexity is preserved under many natural operations and transformations, and thereby the effective range of results can be extended, elegant proof techniques can be developed as well as unforeseen applications of certain results can be given.*
- 3 *Convex functions and domains exhibit sufficient structure so that a mathematically beautiful and practically useful theory can be developed.*
- 4 *There are theoretically and practically (reasonably) efficient methods to find the minimum of a convex function.*

Attractions of Submodular Functions

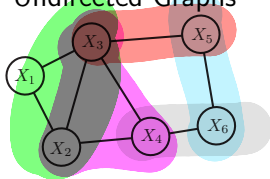
- In this course, we wish to demonstrate that submodular and supermodular functions also possess attractions of these four sorts as well.

Graphical Models and Decomposition

- Let \mathcal{B} be the set of cliques of a graph G . A graphical model prescribes how to write functions $f : \{0, 1\}^n \rightarrow \mathbb{R}$. Let $x \in \{0, 1\}^n$

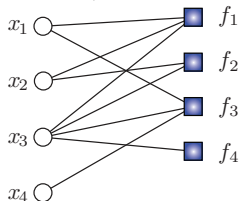
$$f(x) = \sum_{B \in \mathcal{B}} f_B(x_B) \quad (1.1)$$

Example: Undirected Graphs



$$\begin{aligned} f(x_{1:6}) &= f(x_1, x_2, x_3) + f(x_2, x_3, x_4) \\ &\quad + f(x_3, x_5) + f(x_5, x_6) + f(x_4, x_6) \\ f(x_{1:6}) &= f(x_1, x_2) + f(x_2, x_3) + f(x_3, x_1) \\ &\quad + f(x_2, x_3) + f(x_3, x_4) + f(x_4, x_2) \\ &\quad + f(x_3, x_5) + f(x_5, x_6) + f(x_4, x_6) \end{aligned}$$

Example: Factor/Hyper Graphs



$$\begin{aligned} f(x_{1:4}) &= f_1(x_1, x_2, x_3) + f_2(x_2, x_3) \\ &= f_3(x_1, x_3, x_4) + f_4(x_3) \end{aligned}$$

Graphical Models/Decomposition: Real-Object Example

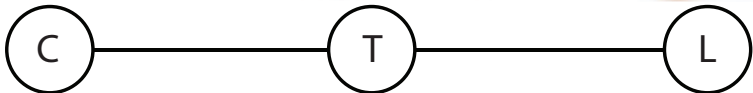
- How to value a set of items?

Graphical Models/Decomposition: Real-Object Example

- How to value a set of items?
- Let C , T , and L be binary variables indicating the presence or absence of items, and we wish to compute $\text{value}(C, T, L)$.

Graphical Models/Decomposition: Real-Object Example

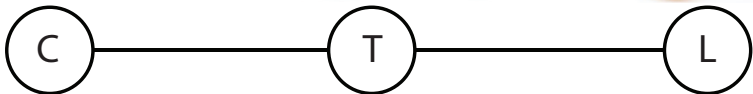
- How to value a set of items?
- Let C , T , and L be binary variables indicating the presence or absence of items, and we wish to compute $\text{value}(C, T, L)$.
- Example: Value of Coffee (C), Tea (T), and Lemon (L).



$$\text{value}(C, T, L) = \text{value}(C, T) + \text{value}(T, L) \quad (1.2)$$

Graphical Decomposition Limitation: Manner of Interaction

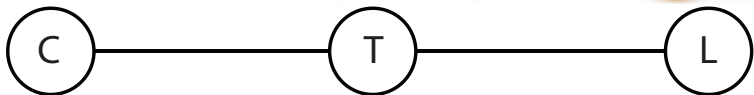
- Value of Coffee (C), Tea (T), and Lemon (L).



$$\text{value}(C, T, L) = \text{value}(C, T) + \text{value}(T, L) \quad (1.3)$$

Graphical Decomposition Limitation: Manner of Interaction

- Value of Coffee (C), Tea (T), and Lemon (L).



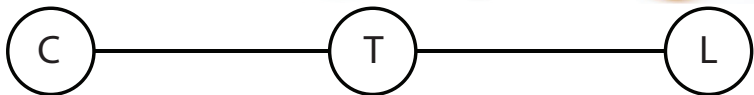
$$\text{value}(C, T, L) = \text{value}(C, T) + \text{value}(T, L) \quad (1.3)$$

- Coffee and Tea are “substitutive”

$$\text{value}(C, T) \leq \text{value}(C) + \text{value}(T) \quad (1.4)$$

Graphical Decomposition Limitation: Manner of Interaction

- Value of Coffee (C), Tea (T), and Lemon (L).



$$\text{value}(C, T, L) = \text{value}(C, T) + \text{value}(T, L) \quad (1.3)$$

- Coffee and Tea are “substitutive”

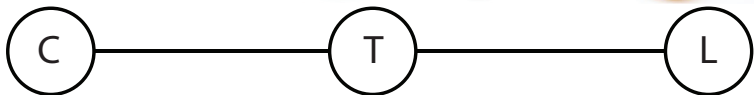
$$\text{value}(C, T) \leq \text{value}(C) + \text{value}(T) \quad (1.4)$$

- Tea and Lemon are “complementary”

$$\text{value}(T, L) \geq \text{value}(T) + \text{value}(L) \quad (1.5)$$

Graphical Decomposition Limitation: Manner of Interaction

- Value of Coffee (C), Tea (T), and Lemon (L).



$$\text{value}(C, T, L) = \text{value}(C, T) + \text{value}(T, L) \quad (1.3)$$

- Coffee and Tea are “substitutive”

$$\text{value}(C, T) \leq \text{value}(C) + \text{value}(T) \quad (1.4)$$

- Tea and Lemon are “complementary”

$$\text{value}(T, L) \geq \text{value}(T) + \text{value}(L) \quad (1.5)$$

- These are distinct non-graphically expressed **manners of interaction!**

Options for Cost Models via Graphical Decomposition

- Three items. Hamburger (H), Fries (F), Soda (S)



Options for Cost Models via Graphical Decomposition

- Three items. Hamburger (H), Fries (F), Soda (S)



- Some graphical model options for $\text{costs}(H, F, S)$:



$$\text{costs}(H, F, S) = \text{cst}_h(H) + \text{cst}_f(F) + \text{cst}_c(S)$$



$$\text{costs}(H, F, S) = \text{cst}_{hf}(H, F) + \text{cst}_{fc}(F, S)$$



$$\text{costs}(H, F, S) = \text{cst}_{hfc}(H, F, S)$$

Decompositions via Manner of Interaction

- costs(H, F, S) of Hamburger (H), Fries (F), Soda (S)



Consider components of cost: consumer-costs (ccs) and health-costs (hcs), each of which is ternary.

$$\text{costs}(H, F, S) = \text{ccs}(H, F, S) + \text{hcs}(H, F, S) \quad (1.6)$$

Decompositions via Manner of Interaction

- costs(H, F, S) of Hamburger (H), Fries (F), Soda (S)



Consider components of cost: consumer-costs (ccs) and health-costs (hcs), each of which is ternary.

$$\text{costs}(H, F, S) = \text{ccs}(H, F, S) + \text{hcs}(H, F, S) \quad (1.6)$$

- Consumer costs

$$\text{CCS} \left(\begin{array}{c} \text{Hamburger} \\ \text{Soda} \end{array} \right) - \text{CCS} \left(\text{Hamburger} \right) \geq \text{CCS} \left(\begin{array}{c} \text{Fries} \\ \text{Hamburger} \\ \text{Soda} \end{array} \right) - \text{CCS} \left(\begin{array}{c} \text{Fries} \\ \text{Hamburger} \end{array} \right)$$

Decompositions via Manner of Interaction

- costs(H, F, S) of Hamburger (H), Fries (F), Soda (S)



Consider components of cost: consumer-costs (ccs) and health-costs (hcs), each of which is ternary.

$$\text{costs}(H, F, S) = \text{ccs}(H, F, S) + \text{hcs}(H, F, S) \quad (1.6)$$

- Consumer costs

$$\text{ccs} \left(\begin{array}{c} \text{Hamburger} \\ \text{Soda} \end{array} \right) - \text{ccs} \left(\text{Hamburger} \right) \geq \text{ccs} \left(\begin{array}{c} \text{Fries} \\ \text{Hamburger} \\ \text{Soda} \end{array} \right) - \text{ccs} \left(\begin{array}{c} \text{Fries} \\ \text{Hamburger} \end{array} \right)$$

- Health costs

$$\text{hcs} \left(\begin{array}{c} \text{Hamburger} \\ \text{Soda} \end{array} \right) - \text{hcs} \left(\text{Hamburger} \right) \leq \text{hcs} \left(\begin{array}{c} \text{Fries} \\ \text{Hamburger} \\ \text{Soda} \end{array} \right) - \text{hcs} \left(\begin{array}{c} \text{Fries} \\ \text{Hamburger} \end{array} \right)$$

Decompositions via Manner of Interaction

- costs(H, F, S) of Hamburger (H), Fries (F), Soda (S)



Consider components of cost: consumer-costs (ccs) and health-costs (hcs), each of which is ternary.

$$\text{costs}(H, F, S) = \text{ccs}(H, F, S) + \text{hcs}(H, F, S) \quad (1.6)$$

- Consumer costs

$$\text{ccs} \left(\begin{array}{c} \text{cup} \\ \text{burger} \end{array} \right) - \text{ccs} \left(\text{burger} \right) \geq \text{ccs} \left(\begin{array}{c} \text{fries} \\ \text{burger} \\ \text{cup} \end{array} \right) - \text{ccs} \left(\begin{array}{c} \text{fries} \\ \text{burger} \end{array} \right)$$

- Health costs

$$\text{hcs} \left(\begin{array}{c} \text{cup} \\ \text{burger} \end{array} \right) - \text{hcs} \left(\text{burger} \right) \leq \text{hcs} \left(\begin{array}{c} \text{fries} \\ \text{burger} \\ \text{cup} \end{array} \right) - \text{hcs} \left(\begin{array}{c} \text{fries} \\ \text{burger} \end{array} \right)$$

- In both cases, graphical-only decompositions fail!

Sets and set functions $f : 2^V \rightarrow \mathbb{R}$

We are given a finite “ground” set V of objects, $2^V \triangleq \{A : A \subseteq V\}$



Also given a set function $f : 2^V \rightarrow \mathbb{R}$ that values subsets $A \subseteq V$.

Ex: $f(V) = 6$

Sets and set functions $f : 2^V \rightarrow \mathbb{R}$

Subset $A \subseteq V$ of objects:



Also given a set function $f : 2^V \rightarrow \mathbb{R}$ that values subsets $A \subseteq V$.

Ex: $f(A) = 1$

Sets and set functions $f : 2^V \rightarrow \mathbb{R}$

Subset $B \subseteq V$ of objects:



Also given a set function $f : 2^V \rightarrow \mathbb{R}$ that values subsets $A \subseteq V$.

Ex: $f(B) = 6$

Set functions are pseudo-Boolean functions

- Any set $A \subseteq V$ can be represented as a binary vector $x \in \{0, 1\}^V$ (a "bit vector" representation of a set).

$$\chi(v) \in \{0, 1\} \quad v \in V$$

Set functions are pseudo-Boolean functions

- Any set $A \subseteq V$ can be represented as a binary vector $x \in \{0, 1\}^V$ (a “bit vector” representation of a set).
- The **characteristic vector** $\mathbf{1}_A \in \{0, 1\}^V$ of a set A is defined one where element $v \in V$ has value:

$$\mathbf{1}_A(v) = \begin{cases} 1 & \text{if } v \in A \\ 0 & \text{else} \end{cases} \quad (1.7)$$

Set functions are pseudo-Boolean functions

- Any set $A \subseteq V$ can be represented as a binary vector $x \in \{0, 1\}^V$ (a “bit vector” representation of a set).
- The **characteristic vector** $\mathbf{1}_A \in \{0, 1\}^V$ of a set A is defined one where element $v \in V$ has value:

$$\mathbf{1}_A(v) = \begin{cases} 1 & \text{if } v \in A \\ 0 & \text{else} \end{cases} \quad (1.7)$$

- Useful to be able to quickly map between $X = X(\mathbf{1}_X)$ and $x(X) \triangleq \mathbf{1}_X$.

Set functions are pseudo-Boolean functions

- Any set $A \subseteq V$ can be represented as a binary vector $x \in \{0, 1\}^V$ (a “bit vector” representation of a set).
- The **characteristic vector** $\mathbf{1}_A \in \{0, 1\}^V$ of a set A is defined one where element $v \in V$ has value:

$$\mathbf{1}_A(v) = \begin{cases} 1 & \text{if } v \in A \\ 0 & \text{else} \end{cases} \quad (1.7)$$

- Useful to be able to quickly map between $X = X(\mathbf{1}_X)$ and $x(X) \triangleq \mathbf{1}_X$.
- $f : \{0, 1\}^V \rightarrow \{0, 1\}$ are known as **Boolean function**.

Set functions are pseudo-Boolean functions

- Any set $A \subseteq V$ can be represented as a binary vector $x \in \{0, 1\}^V$ (a “bit vector” representation of a set).
- The **characteristic vector** $\mathbf{1}_A \in \{0, 1\}^V$ of a set A is defined one where element $v \in V$ has value:

$$\mathbf{1}_A(v) = \begin{cases} 1 & \text{if } v \in A \\ 0 & \text{else} \end{cases} \quad (1.7)$$

- Useful to be able to quickly map between $X = X(\mathbf{1}_X)$ and $x(X) \triangleq \mathbf{1}_X$.
- $f : \{0, 1\}^V \rightarrow \{0, 1\}$ are known as **Boolean function**.
- $f : \{0, 1\}^V \rightarrow \mathbb{R}$ is a **pseudo-Boolean function** (submodular functions are a special case).

Two Equivalent Submodular Definitions

Definition 1.3.1 (submodular concave)

A function $f : 2^V \rightarrow \mathbb{R}$ is **submodular** if for any $A, B \subseteq V$, we have that:

$$f(A) + f(B) \geq f(A \cup B) + f(A \cap B) \quad (1.8)$$

An alternate and (as we will soon see) equivalent definition is:

Definition 1.3.2 (diminishing returns)

A function $f : 2^V \rightarrow \mathbb{R}$ is **submodular** if for any $A \subseteq B \subset V$, and $v \in V \setminus B$, we have that:

$$f(A \cup \{v\}) - f(A) \geq f(B \cup \{v\}) - f(B) \quad (1.9)$$

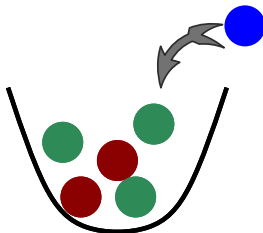
The incremental “value”, “gain”, or “cost” of v decreases (diminishes) as the context in which v is considered grows from A to B .

Example Submodular: Number of Colors of Balls in Urns

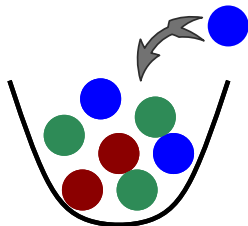
- Consider an urn containing colored balls. Given a set S of balls, $f(S)$ counts the number of distinct colors in S .

Example Submodular: Number of Colors of Balls in Urns

- Consider an urn containing colored balls. Given a set S of balls, $f(S)$ counts the number of distinct colors in S .



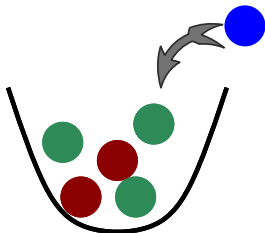
Initial value: 2 (colors in urn).
New value with added blue ball: 3



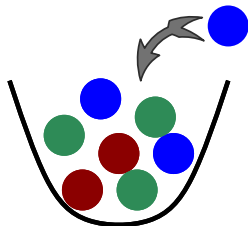
Initial value: 3 (colors in urn).
New value with added blue ball: 3

Example Submodular: Number of Colors of Balls in Urns

- Consider an urn containing colored balls. Given a set S of balls, $f(S)$ counts the number of distinct colors in S .



Initial value: 2 (colors in urn).
New value with added blue ball: 3

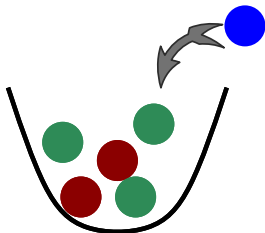


Initial value: 3 (colors in urn).
New value with added blue ball: 3

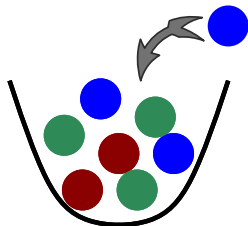
- Submodularity: Incremental Value of Object Diminishes in a Larger Context (diminishing returns).

Example Submodular: Number of Colors of Balls in Urns

- Consider an urn containing colored balls. Given a set S of balls, $f(S)$ counts the number of distinct colors in S .



Initial value: 2 (colors in urn).
 New value with added blue ball: 3



Initial value: 3 (colors in urn).
 New value with added blue ball: 3

- Submodularity: Incremental Value of Object Diminishes in a Larger Context (diminishing returns).
- Thus, f is submodular.

Two Equivalent Supermodular Definitions

Definition 1.3.3 (supermodular)

A function $f : 2^V \rightarrow \mathbb{R}$ is **supermodular** if for any $A, B \subseteq V$, we have that:

$$f(A) + f(B) \leq f(A \cup B) + f(A \cap B) \quad (1.10)$$

Definition 1.3.4 (supermodular (improving returns))

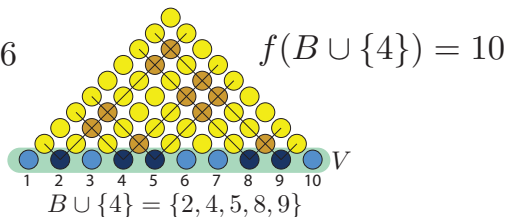
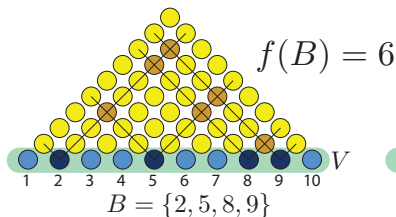
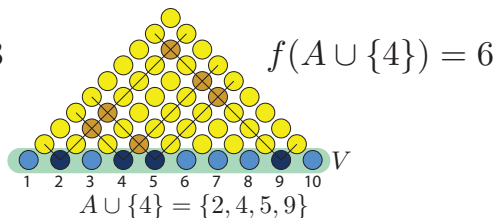
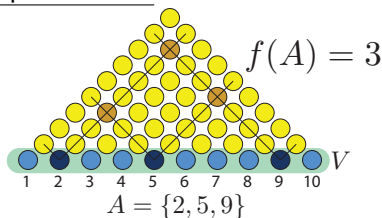
A function $f : 2^V \rightarrow \mathbb{R}$ is **supermodular** if for any $A \subseteq B \subset V$, and $v \in V \setminus B$, we have that:

$$f(A \cup \{v\}) - f(A) \leq f(B \cup \{v\}) - f(B) \quad (1.11)$$

- Incremental “value”, “gain”, or “cost” of v increases (improves) as the context in which v is considered grows from A to B .
- A function f is submodular iff $-f$ is supermodular.
- If f both submodular and supermodular, then f is said to be **modular**, and $f(A) = c + \sum_{a \in A} f(a)$ (often $c = 0$).

Example Supermodular: Number of Balls with Two Lines

Given ball pyramid, bottom row V is size $n = |V|$. For subset $S \subseteq V$ of bottom-row balls, draw 45° and 135° diagonal lines from each $s \in S$. Let $f(S)$ be number of non-bottom-row balls with two lines $\Rightarrow f(S)$ is supermodular.



Scientific Anecdote: Emergent Properties

New York Times column (D. Brooks), March 28th, 2011 on “Tools for Thinking” was about responses to Steven Pinker’s (Harvard) asking a number of scientists “What scientific concept would improve everybody’s cognitive toolkit?”

See <http://edge.org/responses/>

[what-scientific-concept-would-improve-everybodys-cognitive-toolkit](http://edge.org/responses/what-scientific-concept-would-improve-everybodys-cognitive-toolkit)

A common theme was “emergent properties” or “emergent systems”

Emergent systems are ones in which many different elements interact. The pattern of interaction then produces a new element that is greater than the sum of the parts, which then exercises a top-down influence on the constituent elements.

Emergent properties are well modeled by supermodular functions!

Submodular-Supermodular Decomposition

- As an alternative to graphical decomposition, we can decompose a function without resorting sums of local terms.

Submodular-Supermodular Decomposition

- As an alternative to graphical decomposition, we can decompose a function without resorting sums of local terms.

Theorem 1.3.5 (Additive Decomposition (Narasimhan & Bilmes, 2005))

Let $h : 2^V \rightarrow \mathbb{R}$ be **any** set function. Then there exists a submodular function $f : 2^V \rightarrow \mathbb{R}$ and a supermodular function $g : 2^V \rightarrow \mathbb{R}$ such that h may be additively decomposed as follows: For all $A \subseteq V$,

$$h(A) = f(A) + g(A) \tag{1.12}$$

Submodular-Supermodular Decomposition

- As an alternative to graphical decomposition, we can decompose a function without resorting sums of local terms.

Theorem 1.3.5 (Additive Decomposition (Narasimhan & Bilmes, 2005))

Let $h : 2^V \rightarrow \mathbb{R}$ be **any** set function. Then there exists a submodular function $f : 2^V \rightarrow \mathbb{R}$ and a supermodular function $g : 2^V \rightarrow \mathbb{R}$ such that h may be additively decomposed as follows: For all $A \subseteq V$,

$$h(A) = f(A) + g(A) \tag{1.12}$$

- For many applications (as we will see), either the submodular or supermodular component is naturally zero.

Submodular-Supermodular Decomposition

- As an alternative to graphical decomposition, we can decompose a function without resorting sums of local terms.

Theorem 1.3.5 (Additive Decomposition (Narasimhan & Bilmes, 2005))

Let $h : 2^V \rightarrow \mathbb{R}$ be **any** set function. Then there exists a submodular function $f : 2^V \rightarrow \mathbb{R}$ and a supermodular function $g : 2^V \rightarrow \mathbb{R}$ such that h may be additively decomposed as follows: For all $A \subseteq V$,

$$h(A) = f(A) + g(A) \quad (1.12)$$

- For many applications (as we will see), either the submodular or supermodular component is naturally zero.
- Sometimes more natural than a graphical decomposition.

Submodular-Supermodular Decomposition

- As an alternative to graphical decomposition, we can decompose a function without resorting sums of local terms.

Theorem 1.3.5 (Additive Decomposition (Narasimhan & Bilmes, 2005))

Let $h : 2^V \rightarrow \mathbb{R}$ be **any** set function. Then there exists a submodular function $f : 2^V \rightarrow \mathbb{R}$ and a supermodular function $g : 2^V \rightarrow \mathbb{R}$ such that h may be additively decomposed as follows: For all $A \subseteq V$,

$$h(A) = f(A) + g(A) \quad (1.12)$$

- For many applications (as we will see), either the submodular or supermodular component is naturally zero.
- Sometimes more natural than a graphical decomposition.
- Sometimes $h(A)$ has structure in terms of submodular functions but is non additively decomposed (one example is $h(A) = f(A)/g(A)$).

Submodular-Supermodular Decomposition

- As an alternative to graphical decomposition, we can decompose a function without resorting sums of local terms.

Theorem 1.3.5 (Additive Decomposition (Narasimhan & Bilmes, 2005))

Let $h : 2^V \rightarrow \mathbb{R}$ be **any** set function. Then there exists a submodular function $f : 2^V \rightarrow \mathbb{R}$ and a supermodular function $g : 2^V \rightarrow \mathbb{R}$ such that h may be additively decomposed as follows: For all $A \subseteq V$,

$$h(A) = f(A) + g(A) \quad (1.12)$$

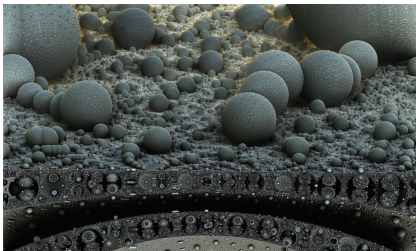
- For many applications (as we will see), either the submodular or supermodular component is naturally zero.
- Sometimes more natural than a graphical decomposition.
- Sometimes $h(A)$ has structure in terms of submodular functions but is non additively decomposed (one example is $h(A) = f(A)/g(A)$).
- Complementary**: simultaneous graphical/submodular-supermodular decomposition (i.e., submodular + supermodular tree).

The Ideal Machine Learning Methods

- Simple to define



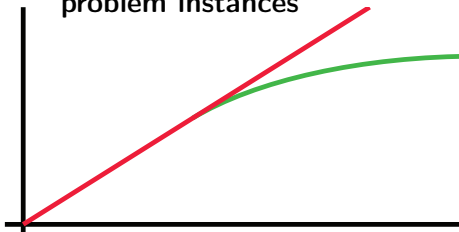
- Mathematically rich



- Naturally suited to many real-world applications



- Efficient & scalable to large problem instances



Discrete Optimization

- Unconstrained minimization and maximization:

$$\min_{X \subseteq V} f(X) \quad (1.13)$$

$$\max_{X \subseteq V} f(X) \quad (1.14)$$

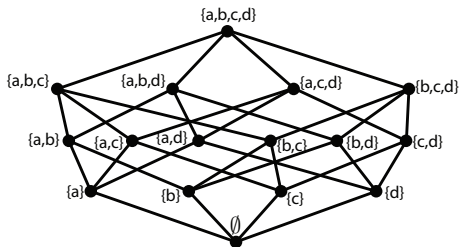
Discrete Optimization

- Unconstrained minimization and maximization:

$$\min_{X \subseteq V} f(X) \quad (1.13)$$

$$\max_{X \subseteq V} f(X) \quad (1.14)$$

- Knowing nothing about f , need 2^n queries for any quality assurance on candidate solution. Otherwise, solution can be unboundedly poor!!



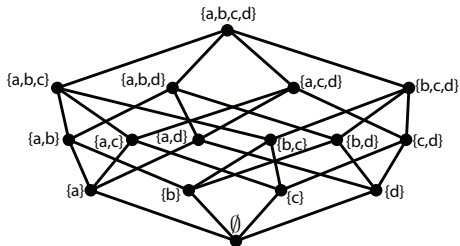
Discrete Optimization

- Unconstrained minimization and maximization:

$$\min_{X \subseteq V} f(X) \quad (1.13)$$

$$\max_{X \subseteq V} f(X) \quad (1.14)$$

- Knowing nothing about f , need 2^n queries for any quality assurance on candidate solution. Otherwise, solution can be unboundedly poor!!



- Alternatively, we may partition V into (necessarily disjoint) blocks $\{V_1, V_2, \dots\}$ that collectively are good in some way.

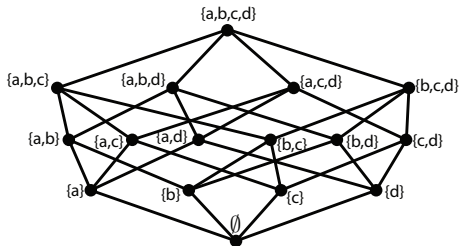
Discrete Optimization

- Unconstrained minimization and maximization:

$$\min_{X \subseteq V} f(X) \quad (1.13)$$

$$\max_{X \subseteq V} f(X) \quad (1.14)$$

- Knowing nothing about f , need 2^n queries for any quality assurance on candidate solution. Otherwise, solution can be unboundedly poor!!



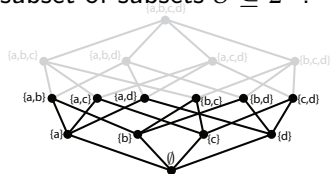
- Alternatively, we may partition V into (necessarily disjoint) blocks $\{V_1, V_2, \dots\}$ that collectively are good in some way.
- When f is submodular, however, Eq. (1.13) is polytime, and Eq. (1.14) is constant-factor approximable. Partitionings are also approximable!

Constrained Discrete Optimization

- Constrained case: interested only in a subset of subsets $\mathcal{S} \subseteq 2^V$.

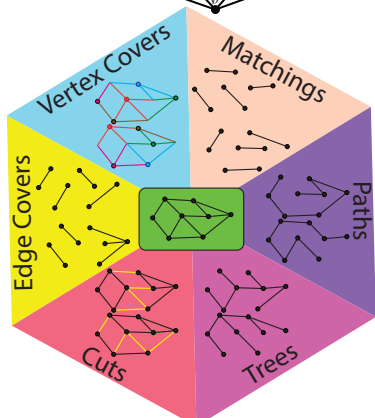
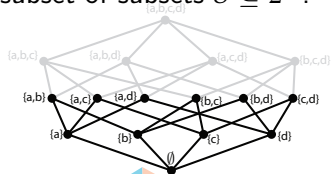
Constrained Discrete Optimization

- Constrained case: interested only in a subset of subsets $\mathcal{S} \subseteq 2^V$.
- Ex: Bounded size $\mathcal{S} = \{S \subseteq V : |S| \leq k\}$, or given cost vector w and budget, bounded cost $\{S \subseteq V : \sum_{s \in S} w(s) \leq b\}$.



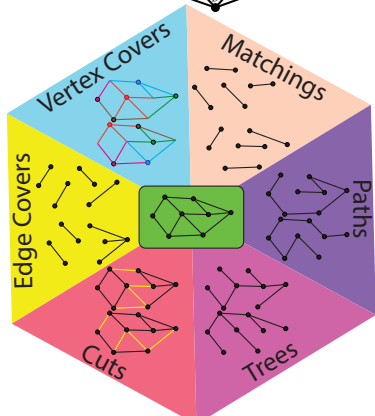
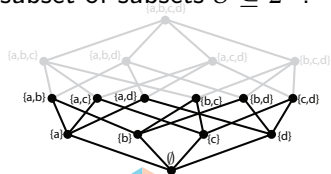
Constrained Discrete Optimization

- Constrained case: interested only in a subset of subsets $\mathcal{S} \subseteq 2^V$.
- Ex: Bounded size $\mathcal{S} = \{S \subseteq V : |S| \leq k\}$, or given cost vector w and budget, bounded cost $\{S \subseteq V : \sum_{s \in S} w(s) \leq b\}$.
- Ex: feasible sets \mathcal{S} as combinatorial objects



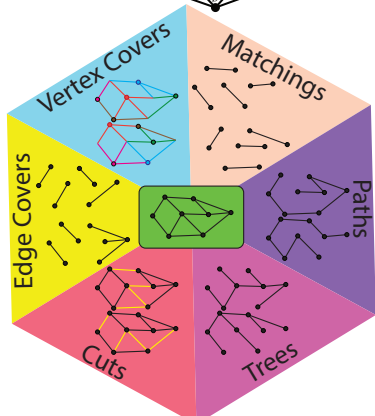
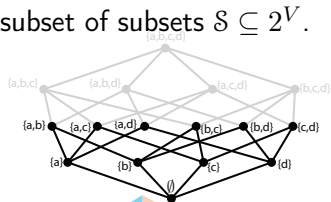
Constrained Discrete Optimization

- Constrained case: interested only in a subset of subsets $\mathcal{S} \subseteq 2^V$.
- Ex: Bounded size $\mathcal{S} = \{S \subseteq V : |S| \leq k\}$, or given cost vector w and budget, bounded cost $\{S \subseteq V : \sum_{s \in S} w(s) \leq b\}$.
- Ex: feasible sets \mathcal{S} as combinatorial objects
- Ex: feasible sets \mathcal{S} as matroids.



Constrained Discrete Optimization

- Constrained case: interested only in a subset of subsets $\mathcal{S} \subseteq 2^V$.
- Ex: Bounded size $\mathcal{S} = \{S \subseteq V : |S| \leq k\}$, or given cost vector w and budget, bounded cost $\{S \subseteq V : \sum_{s \in S} w(s) \leq b\}$.
- Ex: feasible sets \mathcal{S} as combinatorial objects
- Ex: feasible sets \mathcal{S} as matroids.
- Ex: feasible sets \mathcal{S} as sub-level sets of g , $\mathcal{S} = \{S \subseteq V : g(S) \leq \alpha\}$, sup-level sets $\mathcal{S} = \{S \subseteq V : g(S) \geq \alpha\}$



Constrained Discrete Optimization

- Constrained discrete optimization problems:

$$\begin{array}{ll} \text{maximize} & f(S) \\ \text{subject to} & S \in \mathcal{S} \end{array} \quad (1.15)$$

$$\begin{array}{ll} \text{minimize} & f(S) \\ \text{subject to} & S \in \mathcal{S} \end{array} \quad (1.16)$$

where $\mathcal{S} \subseteq 2^V$ is the feasible set of sets.

Constrained Discrete Optimization

- Constrained discrete optimization problems:

$$\begin{array}{ll} \text{maximize} & f(S) \\ \text{subject to} & S \in \mathcal{S} \end{array} \quad (1.15)$$

$$\begin{array}{ll} \text{minimize} & f(S) \\ \text{subject to} & S \in \mathcal{S} \end{array} \quad (1.16)$$

where $\mathcal{S} \subseteq 2^V$ is the feasible set of sets.

- Fortunately, when f (and g) are submodular, these problems can often be solved with guarantees, often very efficiently!

Submodular and Supermodular Applications

- Algorithms: Algorithms can be developed that often are tractable (and as we will see many in this class).
- Applications: There are many seemingly different applications that are strongly related to submodularity.
- Submodularity and supermodularity is as common and natural for discrete problems in machine learning as is convexity/concavity for continuous problems.
- First, let's look at a few more very simple examples of submodular functions.

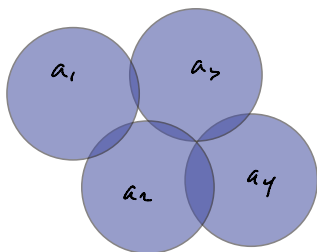
Continuous Set Cover

The area of the union of areas indexed by A

- Let V be a set of indices, and each $v \in V$ indexes a given fixed sub-area of some region in \mathbb{R}^2 .
- Let $\text{area}(v)$ be the area corresponding to item v .
- Let $f(S) = \bigcup_{s \in S} \text{area}(s)$ be the union of the areas indexed by elements in S .
- Then $f(S)$ is submodular, and corresponds to a continuous **set cover function**.

Continuous Set Cover

The area of the union of areas indexed by A — Example

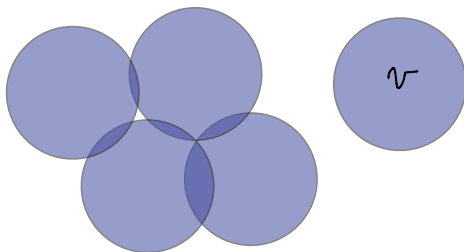


Union of areas of elements of A is given by:

$$f(A) = f(\{a_1, a_2, a_3, a_4\})$$

Continuous Set Cover

The area of the union of areas indexed by A — Example

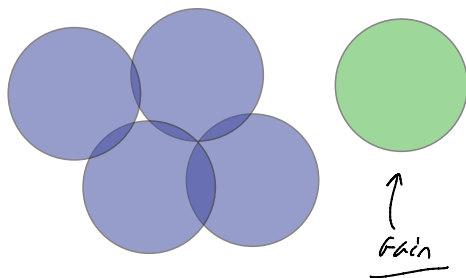


Area of A along with with v :

$$f(A \cup \{v\}) = f(\{a_1, a_2, a_3, a_4\} \cup \{v\})$$

Continuous Set Cover

The area of the union of areas indexed by A — Example



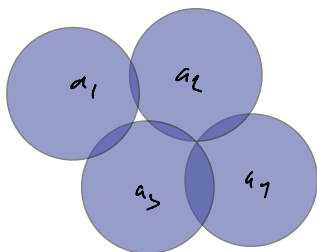
Gain (value) of v in context of A :

$$f(A \cup \{v\}) - f(A) = f(\{v\})$$

We get full value $f(\{v\})$ in this case since the area of v has no overlap with that of A .

Continuous Set Cover

The area of the union of areas indexed by A — Example

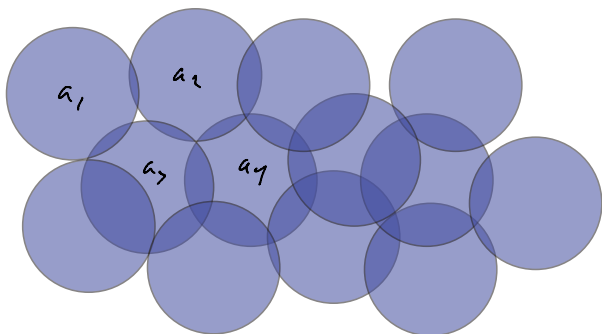


Area of A once again.

$$f(A) = f(\{a_1, a_2, a_3, a_4\})$$

Continuous Set Cover

The area of the union of areas indexed by A — Example

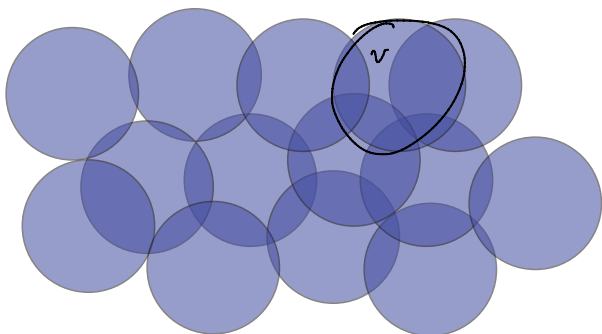


Union of areas of elements of $B \supset A$, where v is not included:

$$f(B) \text{ where } v \notin B \text{ and where } A \subseteq B$$

Continuous Set Cover

The area of the union of areas indexed by A — Example

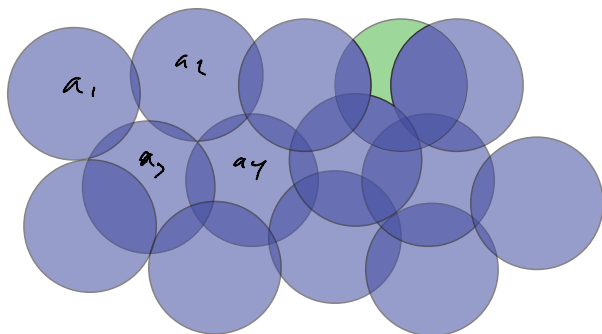


Area of B now also including v :

$$f(B \cup \{v\})$$

Continuous Set Cover

The area of the union of areas indexed by A — Example



Incremental value of v in the context of $B \supset A$.

$$f(B \cup \{v\}) - f(B) < f(\{v\}) = f(A \cup \{v\}) - f(A)$$

So benefit of v in the context of A is greater than the benefit of v in the context of $B \supseteq A$.

Simple Consumer Costs



FUNNYRECIPTS.com

TRADER JOE'S

Store [REDACTED]

OPEN 9:00AM TO 10:00PM DAILY

TJ'S PLAIN SOY MILK	1.69
EGGS BROWN	1.79
VEG TEMPEH ORGANIC 3 GRAIN	1.69
VEG SOY CHORIZO	1.99
PLAIN ORGANIC NONFAT YOGURT 32	2.99
LARGE BABY NON TAXABLE GROCERY	1.99
3 @ 3 FOR 0.49	0.49
SUBTOTAL	\$12.63
TOTAL	\$12.63

- Grocery store: finite set of items V that one can purchase.

Simple Consumer Costs



FUNNYRECIPTS.com

TRADER JOE'S

Store [REDACTED]

OPEN 9:00AM TO 10:00PM DAILY

TJ'S PLAIN SOY MILK	1.69
EGGS BROWN	1.79
VEG TEMPEH ORGANIC 3 GRAIN	1.69
VEG SOY CHORIZO	1.99
PLAIN ORGANIC NONFAT YOGURT 32	2.99
LARGE BABY NON TAXABLE GROCERY	1.99
3 @ 3 FOR 0.49	0.49
SUBTOTAL	\$12.63
TOTAL	\$12.63

- Grocery store: finite set of items V that one can purchase.
- Each item $v \in V$ has a price $m(v)$.

Simple Consumer Costs



FUNNYRECIPTS.com

TRADER JOE'S

Store [REDACTED]

OPEN 9:00AM TO 10:00PM DAILY

TJ'S PLAIN SOY MILK	1.69
EGGS BROWN	1.79
VEG TEMPEH ORGANIC 3 GRAIN	1.69
VEG SOY CHORIZO	1.99
PLAIN ORGANIC NONFAT YOGURT 32	2.99
LARGE BABY NON TAXABLE GROCERY	1.99
3 @ 3 FOR 0.49	0.49
SUBTOTAL	\$12.63
TOTAL	\$12.63

- Grocery store: finite set of items V that one can purchase.
- Each item $v \in V$ has a price $m(v)$.
- Basket of groceries $A \subseteq V$ costs:

$$m(A) = \sum_{a \in A} m(a), \quad (1.17)$$

the sum of individual item costs (no two-for-one discounts).

Simple Consumer Costs



FUNNYRECIPTS.com

TRADER JOE'S

Store [REDACTED]

OPEN 9:00AM TO 10:00PM DAILY

TJ'S PLAIN SOY MILK	1.69
EGGS BROWN	1.79
VEG TEMPEH ORGANIC 3 GRAIN	1.69
VEG SOY CHORIZO	1.99
PLAIN ORGANIC NONFAT YOGURT 32	2.99
LARGE BABY NON TAXABLE GROCERY	1.99
3 @ 3 FOR 0.49	0.49
SUBTOTAL	\$12.63
TOTAL	\$12.63

- Grocery store: finite set of items V that one can purchase.
- Each item $v \in V$ has a price $m(v)$.
- Basket of groceries $A \subseteq V$ costs:

$$m(A) = \sum_{a \in A} m(a), \quad (1.17)$$

the sum of individual item costs (no two-for-one discounts).

- This is known as a modular function.

Discounted Consumer Costs (as we saw earlier)

- Let f be the cost of purchasing a set of items (consumer cost). For example, $V = \{\text{"coke"}, \text{"fries"}, \text{"hamburger"}\}$ and $f(A)$ measures the cost of any subset $A \subseteq V$. We get diminishing returns:

$$f(\text{fries, coke}) - f(\text{fries}) \geq f(\text{fries, coke, hamburger}) - f(\text{fries, hamburger})$$

- Simply rearranging terms, we get the other definition of submodularity:

$$f(\text{fries, coke}) + f(\text{fries, hamburger}) \geq f(\text{fries, coke, hamburger}) + f(\text{fries})$$

- Typical: additional cost of a coke is free only if you add it to a fries and hamburger order.

Shared Fixed Costs (interacting costs)

- Costs often interact in the real world.

Shared Fixed Costs (interacting costs)

- Costs often interact in the real world.
- Ex: Let $V = \{v_1, v_2\}$ be a set of actions with:

$v_1 =$ “buy milk at the store”

$v_2 =$ “buy honey at the store”



Shared Fixed Costs (interacting costs)

- Costs often interact in the real world.
- Ex: Let $V = \{v_1, v_2\}$ be a set of actions with:

$v_1 =$ “buy milk at the store”

$v_2 =$ “buy honey at the store”



- For $A \subseteq V$, let $f(A)$ be the consumer cost of set of items A .

Shared Fixed Costs (interacting costs)

- Costs often interact in the real world.
- Ex: Let $V = \{v_1, v_2\}$ be a set of actions with:

$v_1 =$ “buy milk at the store”

$v_2 =$ “buy honey at the store”



- For $A \subseteq V$, let $f(A)$ be the consumer cost of set of items A .
- $f(\{v_1\}) =$ cost to drive to and from store c_d , and cost to purchase milk c_m , so $f(\{v_1\}) = c_d + c_m$.

Shared Fixed Costs (interacting costs)

- Costs often interact in the real world.
- Ex: Let $V = \{v_1, v_2\}$ be a set of actions with:

v_1 = “buy milk at the store”

v_2 = “buy honey at the store”



- For $A \subseteq V$, let $f(A)$ be the consumer cost of set of items A .
- $f(\{v_1\})$ = cost to drive to and from store c_d , and cost to purchase milk c_m , so $f(\{v_1\}) = c_d + c_m$.
- $f(\{v_2\})$ = cost to drive to and from store c_d , and cost to purchase honey c_h , so $f(\{v_2\}) = c_d + c_h$.

Shared Fixed Costs (interacting costs)

- Costs often interact in the real world.
- Ex: Let $V = \{v_1, v_2\}$ be a set of actions with:

$v_1 =$ “buy milk at the store”

$v_2 =$ “buy honey at the store”



- For $A \subseteq V$, let $f(A)$ be the consumer cost of set of items A .
- $f(\{v_1\}) =$ cost to drive to and from store c_d , and cost to purchase milk c_m , so $f(\{v_1\}) = c_d + c_m$.
- $f(\{v_2\}) =$ cost to drive to and from store c_d , and cost to purchase honey c_h , so $f(\{v_2\}) = c_d + c_h$.
- But $f(\{v_1, v_2\}) = c_d + c_m + c_h < 2c_d + c_m + c_h$ since c_d (driving) is a shared fixed cost.
 $f(v_1) + f(v_2)$

Shared Fixed Costs (interacting costs)

- Costs often interact in the real world.
- Ex: Let $V = \{v_1, v_2\}$ be a set of actions with:

v_1 = “buy milk at the store”

v_2 = “buy honey at the store”



- For $A \subseteq V$, let $f(A)$ be the consumer cost of set of items A .
- $f(\{v_1\})$ = cost to drive to and from store c_d , and cost to purchase milk c_m , so $f(\{v_1\}) = c_d + c_m$.
- $f(\{v_2\})$ = cost to drive to and from store c_d , and cost to purchase honey c_h , so $f(\{v_2\}) = c_d + c_h$.
- But $f(\{v_1, v_2\}) = c_d + c_m + c_h < 2c_d + c_m + c_h$ since c_d (driving) is a shared fixed cost.
- Shared fixed costs are submodular: $f(v_1) + f(v_2) \geq f(v_1, v_2) + f(\emptyset)$

$$c_d + c_m \quad c_d + c_h \quad c_d + c_m + c_h \quad + 0$$

Markets: Supply Side Economies of scale

- **Economies of Scale**: Many goods and services can be produced at a much lower per-unit cost only if they are produced in very large quantities.
- The **profit margin** for producing a unit of goods is improved as more of those goods are created.
- If you already make a good, making a similar good is easier than if you start from scratch (e.g., Apple making both iPod and iPhone).
- An argument in favor of free trade is that it opens up larger markets for firms (especially in otherwise small markets), thereby enabling better economies of scale, and hence greater efficiency (lower costs and resources per unit of good produced).

Supply Side Economies of Scale

- What is a good model of the **cost** of manufacturing a set of items?

Supply Side Economies of Scale

- What is a good model of the **cost** of manufacturing a set of items?
- Let V be a set of possible items to manufacture, and let $f(S)$ for $S \subseteq V$ be the manufacture costs of items in the subset S .

Supply Side Economies of Scale

- What is a good model of the **cost** of manufacturing a set of items?
- Let V be a set of possible items to manufacture, and let $f(S)$ for $S \subseteq V$ be the manufacture costs of items in the subset S .
- Ex: V might be paint colors to produce: green, red, blue, yellow, white, etc.

Supply Side Economies of Scale

- What is a good model of the **cost** of manufacturing a set of items?
- Let V be a set of possible items to manufacture, and let $f(S)$ for $S \subseteq V$ be the manufacture costs of items in the subset S .
- Ex: V might be paint colors to produce: green, red, blue, yellow, white, etc.
- Producing green when you are already producing yellow and blue is probably cheaper than if you were only producing some other colors.

$$f(\text{green, blue, yellow}) - f(\text{blue, yellow}) \leq f(\text{green, blue}) - f(\text{blue})$$

Supply Side Economies of Scale

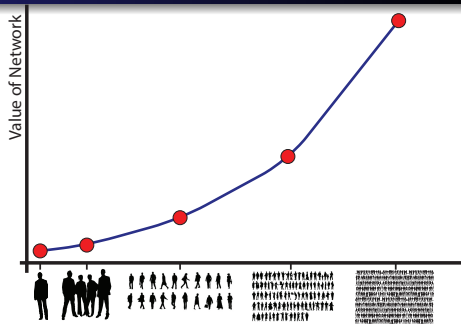
- What is a good model of the **cost** of manufacturing a set of items?
- Let V be a set of possible items to manufacture, and let $f(S)$ for $S \subseteq V$ be the manufacture costs of items in the subset S .
- Ex: V might be paint colors to produce: green, red, blue, yellow, white, etc.
- Producing green when you are already producing yellow and blue is probably cheaper than if you were only producing some other colors.

$$f(\text{green, blue, yellow}) - f(\text{blue, yellow}) \leq f(\text{green, blue}) - f(\text{blue})$$

- So diminishing returns (a submodular function) would be a good model.

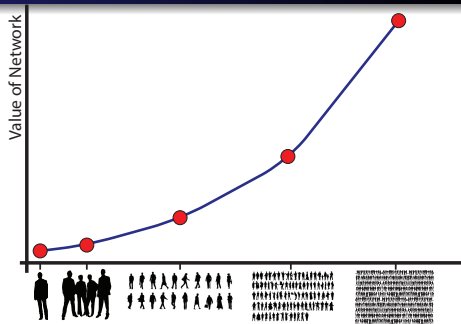
Demand side Economies of Scale: Network Externalities

- Value of a network to a user depends on the number of other users in that network. External use benefits internal use.



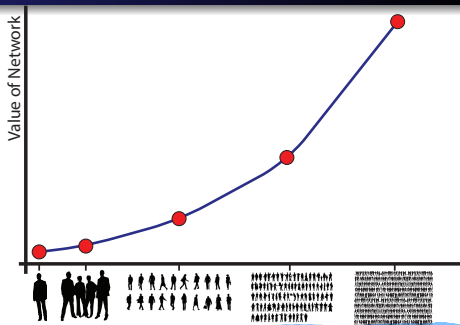
Demand side Economies of Scale: Network Externalities

- Value of a network to a user depends on the number of other users in that network. External use benefits internal use.
- Consumers derive positive incremental value when size of the market for that good increases.



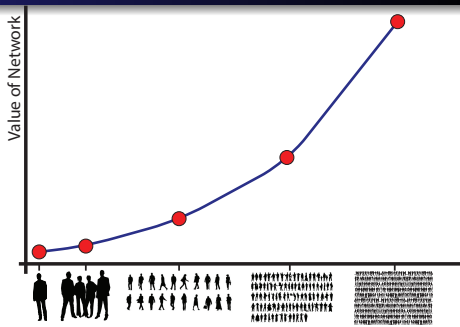
Demand side Economies of Scale: Network Externalities

- Value of a network to a user depends on the number of other users in that network. External use benefits internal use.
- Consumers derive positive incremental value when size of the market for that good increases.
- Called **network externalities** (Katz & Shapiro 1986), or **network effects** and is a form of demand-side economies of scale



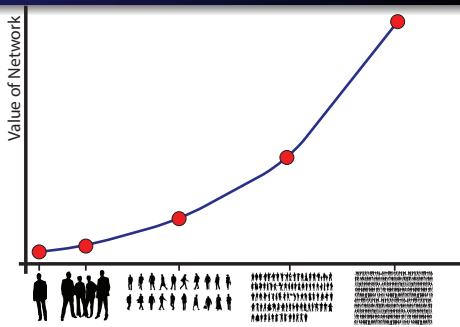
Demand side Economies of Scale: Network Externalities

- Value of a network to a user depends on the number of other users in that network. External use benefits internal use.
- Consumers derive positive incremental value when size of the market for that good increases.
- Called **network externalities** (Katz & Shapiro 1986), or **network effects** and is a form of demand-side economies of scale
- Ex: durable goods (e.g., a car or phone), software (facebook, smartphone apps), and technology-specific human capital investment (e.g., education in a skill), benefit depends on total user base.



Demand side Economies of Scale: Network Externalities

- Value of a network to a user depends on the number of other users in that network. External use benefits internal use.
- Consumers derive positive incremental value when size of the market for that good increases.
- Called **network externalities** (Katz & Shapiro 1986), or **network effects** and is a form of demand-side economies of scale
- Ex: durable goods (e.g., a car or phone), software (facebook, smartphone apps), and technology-specific human capital investment (e.g., education in a skill), benefit depends on total user base.
- Let V be a set of goods, A a subset and $v \notin A$. Incremental gain of good $f(A + v) - f(A)$ gets larger as size of market A grows. This is known as a **supermodular** function.



Examples: Positive Network Effects

- railroad - standard rail format and shared access
- The telephone, who wants to talk by phone only to oneself?
- the internet, more valuable per person the more people use it.
- ebooks (the more people comment, the better it gets)
- social network sites: facebook more valuable with everyone online
- online education, Massive Open Online Courses (MOOCs) such as Coursera, edX, etc. – with many people simultaneously taking a class, all gain due to richer peer discussions due to greater pool of well matched study groups, more simultaneous similar questions/problems that are asked \Rightarrow more efficient learning & training data for ML algorithms to learn how people learn.
- Software (e.g., Microsoft office, smartphone apps, etc.): more people means more bug reporting \Rightarrow better & faster software evolution.
- gmail and web-based email (collaborative spam filtering).
- wikipedia, collaborative documents
- any widely used standard (job training now is useful in the future)
- the “tipping point”, and “winner take all” (one platform prevails)

Examples: Other Network Effects

No Network Externalities

- food/drink - (should be) independent of how many others are eating the type of food.

Examples: Other Network Effects

No Network Externalities

- food/drink - (should be) independent of how many others are eating the type of food.
- Music - your enjoyment should (ideally) be independent of others' enjoyment (but maybe not, see Salganik, Dodds, Watts'06).

Examples: Other Network Effects

No Network Externalities

- food/drink - (should be) independent of how many others are eating the type of food.
- Music - your enjoyment should (ideally) be independent of others' enjoyment (but maybe not, see Salganik, Dodds, Watts'06).

Negative Network Effects

Examples: Other Network Effects

No Network Externalities

- food/drink - (should be) independent of how many others are eating the type of food.
- Music - your enjoyment should (ideally) be independent of others' enjoyment (but maybe not, see Salganik, Dodds, Watts'06).

Negative Network Effects

- clothing

Examples: Other Network Effects

No Network Externalities

- food/drink - (should be) independent of how many others are eating the type of food.
- Music - your enjoyment should (ideally) be independent of others' enjoyment (but maybe not, see Salganik, Dodds, Watts'06).

Negative Network Effects

- clothing
- (Halloween) costumes

Examples: Other Network Effects

No Network Externalities

- food/drink - (should be) independent of how many others are eating the type of food.
- Music - your enjoyment should (ideally) be independent of others' enjoyment (but maybe not, see Salganik, Dodds, Watts'06).

Negative Network Effects

- clothing
- (Halloween) costumes
- perfume?

Review So far

- Machine learning paradigms should be: **easy to define**, **mathematically rich**, **naturally applicable**, and **efficient/scalable**.

Review So far

- Machine learning paradigms should be: **easy to define**, **mathematically rich**, **naturally applicable**, and **efficient/scalable**.
- **Convexity** (continuous structures) and **graphical models** (based on factorization or additive separation) are two such modeling paradigms.

Review So far

- Machine learning paradigms should be: **easy to define**, **mathematically rich**, **naturally applicable**, and **efficient/scalable**.
- **Convexity** (continuous structures) and **graphical models** (based on factorization or additive separation) are two such modeling paradigms.
- **Submodularity/supermodularity** offer a distinct mathematically rich paradigm over discrete space that neither need be continuous nor be additively additively separable,

Review So far

- Machine learning paradigms should be: **easy to define**, **mathematically rich**, **naturally applicable**, and **efficient/scalable**.
- **Convexity** (continuous structures) and **graphical models** (based on factorization or additive separation) are two such modeling paradigms.
- **Submodularity/supermodularity** offer a distinct mathematically rich paradigm over discrete space that neither need be continuous nor be additively additively separable,
- submodularity offers forms of structural decomposition, e.g., $h = f + g$, into potentially global (manner of interaction) terms.

Review So far

- Machine learning paradigms should be: **easy to define**, **mathematically rich**, **naturally applicable**, and **efficient/scalable**.
- **Convexity** (continuous structures) and **graphical models** (based on factorization or additive separation) are two such modeling paradigms.
- **Submodularity/supermodularity** offer a distinct mathematically rich paradigm over discrete space that neither need be continuous nor be additively additively separable,
- submodularity offers forms of structural decomposition, e.g., $h = f + g$, into potentially global (manner of interaction) terms.
- **Set cover, supply and demand side economies of scale,**

Submodularity's utility in ML

- A model of a physical process :

Submodularity's utility in ML

- A model of a physical process :
 - When **maximizing**, submodularity naturally models: diversity, coverage, span, and information.

Submodularity's utility in ML

- A model of a physical process :
 - When **maximizing**, submodularity naturally models: diversity, coverage, span, and information.
 - When **minimizing**, submodularity naturally models: cooperative costs, complexity, roughness, and irregularity.

Submodularity's utility in ML

- A model of a physical process :
 - When **maximizing**, submodularity naturally models: diversity, coverage, span, and information.
 - When **minimizing**, submodularity naturally models: cooperative costs, complexity, roughness, and irregularity.
 - vice-versa for supermodularity.

Submodularity's utility in ML

- A **model of a physical process** :
 - When **maximizing**, submodularity naturally models: diversity, coverage, span, and information.
 - When **minimizing**, submodularity naturally models: cooperative costs, complexity, roughness, and irregularity.
 - vice-versa for supermodularity.
- A submodular function can act as a **parameter** for a machine learning strategy (active/semi-supervised learning, discrete divergence, structured sparse convex norms for use in regularization).

Submodularity's utility in ML

- A model of a physical process :
 - When **maximizing**, submodularity naturally models: diversity, coverage, span, and information.
 - When **minimizing**, submodularity naturally models: cooperative costs, complexity, roughness, and irregularity.
 - vice-versa for supermodularity.
- A submodular function can act as a **parameter** for a machine learning strategy (active/semi-supervised learning, discrete divergence, structured sparse convex norms for use in regularization).
- Itself, as an object or function **to learn**, based on data.

Submodularity's utility in ML

- A **model of a physical process** :
 - When **maximizing**, submodularity naturally models: diversity, coverage, span, and information.
 - When **minimizing**, submodularity naturally models: cooperative costs, complexity, roughness, and irregularity.
 - vice-versa for supermodularity.
- A submodular function can act as a **parameter** for a machine learning strategy (active/semi-supervised learning, discrete divergence, structured sparse convex norms for use in regularization).
- Itself, as an object or function **to learn**, based on data.
- A **surrogate or relaxation strategy** for optimization or analysis

Submodularity's utility in ML

- A **model of a physical process** :
 - When **maximizing**, submodularity naturally models: diversity, coverage, span, and information.
 - When **minimizing**, submodularity naturally models: cooperative costs, complexity, roughness, and irregularity.
 - vice-versa for supermodularity.
- A submodular function can act as a **parameter** for a machine learning strategy (active/semi-supervised learning, discrete divergence, structured sparse convex norms for use in regularization).
- Itself, as an object or function **to learn**, based on data.
- A **surrogate or relaxation strategy** for optimization or analysis
 - An alternate to factorization, decomposition, or sum-product based simplification (as one typically finds in a graphical model). I.e., a means towards tractable surrogates for graphical models.

Submodularity's utility in ML

- A **model of a physical process** :
 - When **maximizing**, submodularity naturally models: diversity, coverage, span, and information.
 - When **minimizing**, submodularity naturally models: cooperative costs, complexity, roughness, and irregularity.
 - vice-versa for supermodularity.
- A submodular function can act as a **parameter** for a machine learning strategy (active/semi-supervised learning, discrete divergence, structured sparse convex norms for use in regularization).
- Itself, as an object or function **to learn**, based on data.
- A **surrogate or relaxation strategy** for optimization or analysis
 - An alternate to factorization, decomposition, or sum-product based simplification (as one typically finds in a graphical model). I.e., a means towards tractable surrogates for graphical models.
 - Also, we can “relax” a problem to a submodular one where it can be efficiently solved and offer a bounded quality solution.

Submodularity's utility in ML

- A **model of a physical process** :
 - When **maximizing**, submodularity naturally models: diversity, coverage, span, and information.
 - When **minimizing**, submodularity naturally models: cooperative costs, complexity, roughness, and irregularity.
 - vice-versa for supermodularity.
- A submodular function can act as a **parameter** for a machine learning strategy (active/semi-supervised learning, discrete divergence, structured sparse convex norms for use in regularization).
- Itself, as an object or function **to learn**, based on data.
- A **surrogate or relaxation strategy** for optimization or analysis
 - An alternate to factorization, decomposition, or sum-product based simplification (as one typically finds in a graphical model). I.e., a means towards tractable surrogates for graphical models.
 - Also, we can “relax” a problem to a submodular one where it can be efficiently solved and offer a bounded quality solution.
 - **Non-submodular problems can be analyzed via submodularity.**

Many different functions are submodular!

- We will see many applications of submodularity in machine learning.
- On next set of slides, we will state (without proof, for now) that many of the functions are submodular (or supermodular).
- In subsequent lectures, we will start showing how to prove submodularity.

Functions to Measure Diversity

Diversity is good, especially when it is high

- Quantitative measurement diversity in data science and ML. **Goal of diversity**: ensure small set properly represents the large.

Functions to Measure Diversity

Diversity is good, especially when it is high

- Quantitative measurement diversity in data science and ML. **Goal of diversity**: ensure small set properly represents the large.
- Web search: given ambiguous search term (e.g., “jaguar”) with no other information, one wants articles more than just about cars.
 - Try google searching for words (e.g., “break”) with many meanings (<http://muse.dillfrog.com/lists/ambiguous>), how well does google’s diversity measure do?

Functions to Measure Diversity

Diversity is good, especially when it is high

- Quantitative measurement diversity in data science and ML. **Goal of diversity**: ensure small set properly represents the large.
- Web search: given ambiguous search term (e.g., “jaguar”) with no other information, one wants articles more than just about cars.
 - Try google searching for words (e.g., “break”) with many meanings (<http://muse.dillfrog.com/lists/ambiguous>), how well does google’s diversity measure do?
 - Overall goal: user quickly finds informative, concise, accurate, relevant, comprehensive information \Rightarrow diversity

Functions to Measure Diversity

Diversity is good, especially when it is high

- Quantitative measurement diversity in data science and ML. **Goal of diversity**: ensure small set properly represents the large.
- Web search: given ambiguous search term (e.g., “jaguar”) with no other information, one wants articles more than just about cars.
 - Try google searching for words (e.g., “break”) with many meanings (<http://muse.dillfrog.com/lists/ambiguous>), how well does google’s diversity measure do?
 - Overall goal: user quickly finds informative, concise, accurate, relevant, comprehensive information \Rightarrow diversity
- Given a set V of items, how do we choose a subset $S \subseteq V$ that is as diverse as possible, with perhaps constraints on S such as its size?
Answer: submodular maximization.

Functions to Measure Diversity

Diversity is good, especially when it is high

- Quantitative measurement diversity in data science and ML. **Goal of diversity**: ensure small set properly represents the large.
- Web search: given ambiguous search term (e.g., “jaguar”) with no other information, one wants articles more than just about cars.
 - Try google searching for words (e.g., “break”) with many meanings (<http://muse.dillfrog.com/lists/ambiguous>), how well does google’s diversity measure do?
 - Overall goal: user quickly finds informative, concise, accurate, relevant, comprehensive information \Rightarrow diversity
- Given a set V of items, how do we choose a subset $S \subseteq V$ that is as diverse as possible, with perhaps constraints on S such as its size?
Answer: submodular maximization.
- How do we choose the smallest set S that maintains a given degree of diversity? Constrained minimization (i.e., $\min |A|$ s.t. $f(A) \geq \alpha$).

Functions to Measure Diversity

Diversity is good, especially when it is high

- Quantitative measurement diversity in data science and ML. **Goal of diversity**: ensure small set properly represents the large.
- Web search: given ambiguous search term (e.g., “jaguar”) with no other information, one wants articles more than just about cars.
 - Try google searching for words (e.g., “break”) with many meanings (<http://muse.dillfrog.com/lists/ambiguous>), how well does google’s diversity measure do?
 - Overall goal: user quickly finds informative, concise, accurate, relevant, comprehensive information \Rightarrow diversity
- Given a set V of items, how do we choose a subset $S \subseteq V$ that is as diverse as possible, with perhaps constraints on S such as its size?
Answer: submodular maximization.
- How do we choose the smallest set S that maintains a given degree of diversity? Constrained minimization (i.e., $\min |A|$ s.t. $f(A) \geq \alpha$).
- Random sample has probability of poorly representing normally underrepresented groups.

Extractive Document Summarization

- The figure below represents the sentences of a document



Extractive Document Summarization

- We extract sentences (green) as a summary of the full document

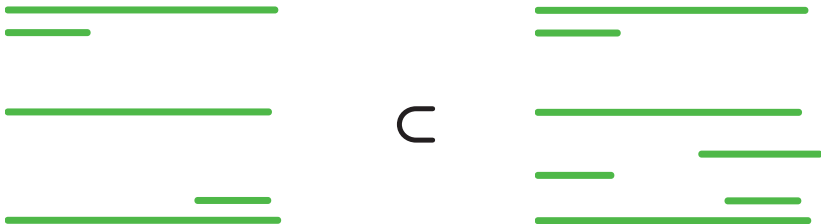


Extractive Document Summarization

- We extract sentences (green) as a summary of the full document

Extractive Document Summarization

- We extract sentences (green) as a summary of the full document



- The summary on the left is a subset of the summary on the right.

Extractive Document Summarization

- We extract sentences (green) as a summary of the full document



- The summary on the left is a subset of the summary on the right.
- Consider adding a new (blue) sentence to each of the two summaries.

Extractive Document Summarization

- We extract sentences (green) as a summary of the full document



- The summary on the left is a subset of the summary on the right.
- Consider adding a new (blue) sentence to each of the two summaries.
- The marginal (incremental) benefit of adding the new (blue) sentence to the smaller (left) summary is no more than the marginal benefit of adding the new sentence to the larger (right) summary.

Extractive Document Summarization

- We extract sentences (green) as a summary of the full document



- The summary on the left is a subset of the summary on the right.
- Consider adding a new (blue) sentence to each of the two summaries.
- The marginal (incremental) benefit of adding the new (blue) sentence to the smaller (left) summary is no more than the marginal benefit of adding the new sentence to the larger (right) summary.
- **diminishing returns** \leftrightarrow **submodularity**

Large image collections need to be summarized

Many images, also that have a higher level gestalt than just a few, want a summary (subset of images) to represent the diversity in the large image set.

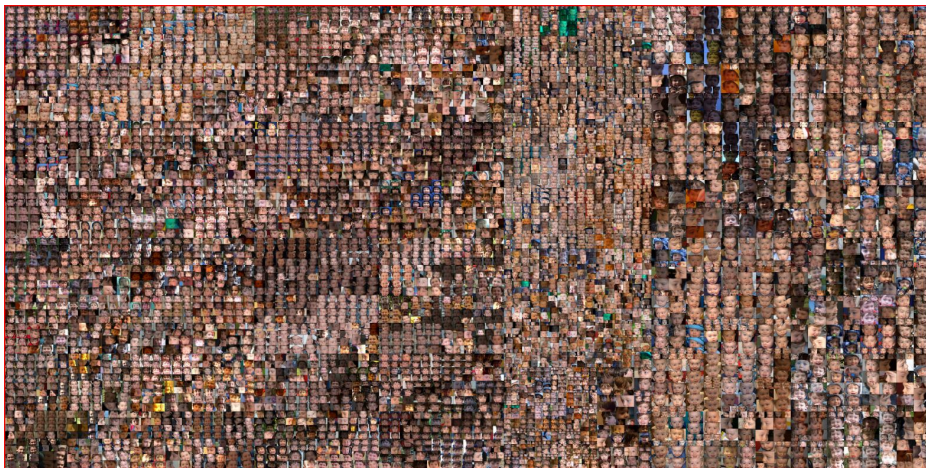


Image Summarization

10×10 image collection:



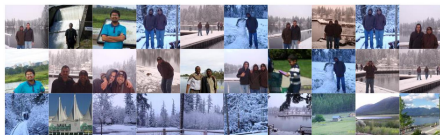
3 good summaries (diverse):



3 ok summaries:



3 poor summaries (redundant):



More Generally: Information and Summarization

- Let V be a set of information containing elements (V might say be any of words, sentences, documents, web pages, or blogs, sensor readings, etc.).
- Each $v \in V$ is one (or a set of) element(s). The total amount of information in V is measure by a function $f(V)$, and any given subset $S \subseteq V$ measures the amount of information in S , given by $f(S)$.
- How informative is any given item v in different sized contexts? Any such real-world information function would exhibit diminishing returns, i.e., the value of v decreases when it is considered in a larger context.
- A submodular function is likely a good model.

Variable Selection in Classification/Regression

- Let Y be a random variable we wish to accurately predict based on at most $n = |V|$ observed measurement variables $(X_1, X_2, \dots, X_n) = X_V$ in a probability model $\Pr(Y, X_1, X_2, \dots, X_n)$.

Variable Selection in Classification/Regression

- Let Y be a random variable we wish to accurately predict based on at most $n = |V|$ observed measurement variables $(X_1, X_2, \dots, X_n) = X_V$ in a probability model $\Pr(Y, X_1, X_2, \dots, X_n)$.
- Too costly to use all V variables. Goal: choose subset $A \subseteq V$ of variables within budget $|A| \leq k$. Predictions based on only $\Pr(y|x_A)$, hence subset A should retain accuracy.

Variable Selection in Classification/Regression

- Let Y be a random variable we wish to accurately predict based on at most $n = |V|$ observed measurement variables $(X_1, X_2, \dots, X_n) = X_V$ in a probability model $\Pr(Y, X_1, X_2, \dots, X_n)$.
- Too costly to use all V variables. Goal: choose subset $A \subseteq V$ of variables within budget $|A| \leq k$. Predictions based on only $\Pr(y|x_A)$, hence subset A should retain accuracy.
- The mutual information function $f(A) = I(Y; X_A)$ ("information gain") measures how well variables A can predicting Y (entropy reduction, reduction of uncertainty of Y).

Variable Selection in Classification/Regression

- Let Y be a random variable we wish to accurately predict based on at most $n = |V|$ observed measurement variables $(X_1, X_2, \dots, X_n) = X_V$ in a probability model $\Pr(Y, X_1, X_2, \dots, X_n)$.
- Too costly to use all V variables. Goal: choose subset $A \subseteq V$ of variables within budget $|A| \leq k$. Predictions based on only $\Pr(y|x_A)$, hence subset A should retain accuracy.
- The mutual information function $f(A) = I(Y; X_A)$ ("information gain") measures how well variables A can predicting Y (entropy reduction, reduction of uncertainty of Y).
- The mutual information function $f(A) = I(Y; X_A)$ is defined as:

$$I(Y; X_A) = \sum_{y, x_A} \Pr(y, x_A) \log \frac{\Pr(y, x_A)}{\Pr(y) \Pr(x_A)} = H(Y) - H(Y|X_A) \quad (1.18)$$

$$= H(X_A) - H(X_A|Y) = H(X_A) + H(Y) - H(X_A, Y) \quad (1.19)$$

Variable Selection in Classification/Regression

- Let Y be a random variable we wish to accurately predict based on at most $n = |V|$ observed measurement variables $(X_1, X_2, \dots, X_n) = X_V$ in a probability model $\Pr(Y, X_1, X_2, \dots, X_n)$.
- Too costly to use all V variables. Goal: choose subset $A \subseteq V$ of variables within budget $|A| \leq k$. Predictions based on only $\Pr(y|x_A)$, hence subset A should retain accuracy.
- The mutual information function $f(A) = I(Y; X_A)$ ("information gain") measures how well variables A can predicting Y (entropy reduction, reduction of uncertainty of Y).
- The mutual information function $f(A) = I(Y; X_A)$ is defined as:

$$I(Y; X_A) = \sum_{y, x_A} \Pr(y, x_A) \log \frac{\Pr(y, x_A)}{\Pr(y) \Pr(x_A)} = H(Y) - H(Y|X_A) \quad (1.18)$$

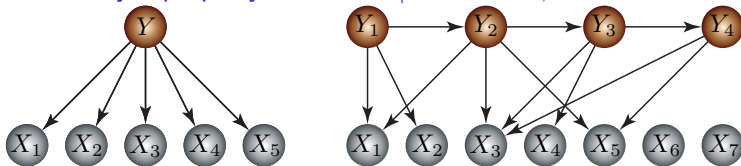
$$= H(X_A) - H(X_A|Y) = H(X_A) + H(Y) - H(X_A, Y) \quad (1.19)$$

- Applicable in pattern recognition, also in sensor coverage problem, where Y is whatever question we wish to ask about environment.

Information Gain and Feature Selection

in Pattern Classification: Naïve Bayes

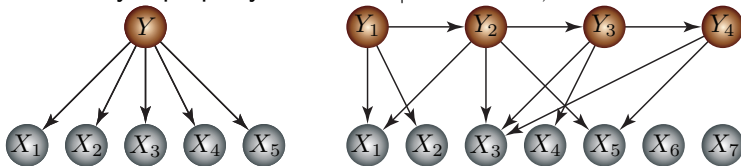
- Naïve Bayes property: $X_A \perp\!\!\!\perp X_B | Y$ for all A, B .



Information Gain and Feature Selection

in Pattern Classification: Naïve Bayes

- Naïve Bayes property: $X_A \perp\!\!\!\perp X_B | Y$ for all A, B .



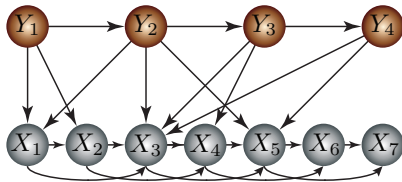
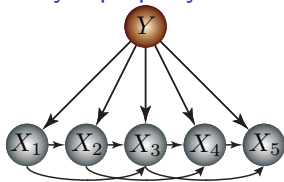
- When $X_A \perp\!\!\!\perp X_B | Y$ for all A, B (the Naïve Bayes assumption holds), then

$$f(A) = I(Y; X_A) = H(X_A) - H(X_A|Y) = H(X_A) - \sum_{a \in A} H(X_a|Y) \quad (1.20)$$

is submodular (submodular minus modular).

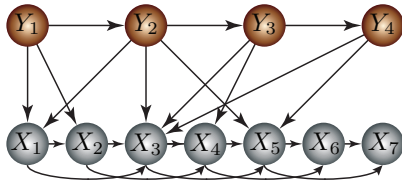
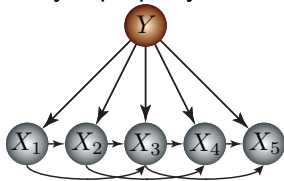
Variable Selection in Pattern Classification

- Naïve Bayes property fails:



Variable Selection in Pattern Classification

- Naïve Bayes property fails:



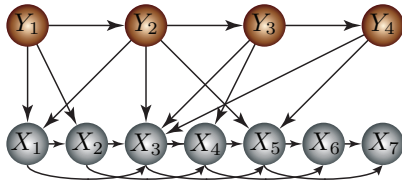
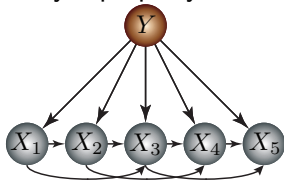
- $f(A)$ naturally expressed as a difference of two submodular functions

$$f(A) = I(Y; X_A) = H(X_A) - H(X_A|Y), \quad (1.21)$$

which is a DS (difference of submodular) function.

Variable Selection in Pattern Classification

- Naïve Bayes property fails:



- $f(A)$ naturally expressed as a difference of two submodular functions

$$f(A) = I(Y; X_A) = H(X_A) - H(X_A|Y), \quad (1.21)$$

which is a DS (difference of submodular) function.

- Alternatively, when Naïve Bayes assumption is false, we can make a submodular approximation (Peng-2005). E.g., functions of the form:

$$f(A) = \sum_{a \in A} I(X_a; Y) - \lambda \sum_{a, a' \in A} I(X_a; X_{a'}|Y) \quad (1.22)$$

where $\lambda \geq 0$ is a tradeoff constant.

Variable Selection: Linear Regression Case

- Next, let Z be continuous. Predictor is linear $\tilde{Z}_A = \sum_{i \in A} \alpha_i X_i$.

Variable Selection: Linear Regression Case

- Next, let Z be continuous. Predictor is linear $\tilde{Z}_A = \sum_{i \in A} \alpha_i X_i$.
- Error measure is the residual variance

$$R_{Z,A}^2 = \frac{\text{Var}(Z) - E[(Z - \tilde{Z}_A)^2]}{\text{Var}(Z)} \quad (1.23)$$

Variable Selection: Linear Regression Case

- Next, let Z be continuous. Predictor is linear $\tilde{Z}_A = \sum_{i \in A} \alpha_i X_i$.
- Error measure is the residual variance

$$R_{Z,A}^2 = \frac{\text{Var}(Z) - E[(Z - \tilde{Z}_A)^2]}{\text{Var}(Z)} \quad (1.23)$$

- $R_{Z,A}^2$'s minimizing parameters, for a given A , can be easily computed ($R_{Z,A}^2 = b_A^\top (C_A^{-1})^\top b_A$ when $\text{Var}Z = 1$, where $b_i = \text{Cov}(Z, X_i)$ and $C = E[(X - E[X])^\top (X - E[X])]$ is the covariance matrix).

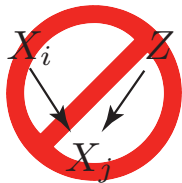
Variable Selection: Linear Regression Case

- Next, let Z be continuous. Predictor is linear $\tilde{Z}_A = \sum_{i \in A} \alpha_i X_i$.
- Error measure is the residual variance

$$R_{Z,A}^2 = \frac{\text{Var}(Z) - E[(Z - \tilde{Z}_A)^2]}{\text{Var}(Z)} \quad (1.23)$$

- $R_{Z,A}^2$'s minimizing parameters, for a given A , can be easily computed ($R_{Z,A}^2 = b_A^\top (C_A^{-1})^\top b_A$ when $\text{Var}Z = 1$, where $b_i = \text{Cov}(Z, X_i)$ and $C = E[(X - E[X])^\top (X - E[X])]$ is the covariance matrix).
- When there are no “suppressor” variables (essentially, no v -structures that converge on X_j with parents X_i and Z), then

$$f(A) = R_{Z,A}^2 = b_A^\top (C_A^{-1})^\top b_A \quad (1.24)$$



is a submodular function (so the greedy algorithm gives the $1 - 1/e$ guarantee). (Das&Kempe).

Data Subset Selection

- Suppose we are given a large data set $\mathcal{D} = \{x_i\}_{i=1}^n$ of n data items $V = \{v_1, v_2, \dots, v_n\}$ and we wish to choose a subset $A \subset V$ of items that is good in some way (e.g., a summary).

Data Subset Selection

- Suppose we are given a large data set $\mathcal{D} = \{x_i\}_{i=1}^n$ of n data items $V = \{v_1, v_2, \dots, v_n\}$ and we wish to choose a subset $A \subset V$ of items that is good in some way (e.g., a summary).
- Suppose moreover each data item $v \in V$ is described by a vector of non-negative scores for a set U of **features** (or “properties”, or “concepts”, etc.) of each data item.

Data Subset Selection

- Suppose we are given a large data set $\mathcal{D} = \{x_i\}_{i=1}^n$ of n data items $V = \{v_1, v_2, \dots, v_n\}$ and we wish to choose a subset $A \subset V$ of items that is good in some way (e.g., a summary).
- Suppose moreover each data item $v \in V$ is described by a vector of non-negative scores for a set U of **features** (or “properties”, or “concepts”, etc.) of each data item.
- That is, for $u \in U$ and $v \in V$, let $m_u(v)$ represent the “degree of u -ness” possessed by data item v . Then $m_u \in \mathbb{R}_+^V$ for all $u \in U$.

Data Subset Selection

- Suppose we are given a large data set $\mathcal{D} = \{x_i\}_{i=1}^n$ of n data items $V = \{v_1, v_2, \dots, v_n\}$ and we wish to choose a subset $A \subset V$ of items that is good in some way (e.g., a summary).
- Suppose moreover each data item $v \in V$ is described by a vector of non-negative scores for a set U of **features** (or “properties”, or “concepts”, etc.) of each data item.
- That is, for $u \in U$ and $v \in V$, let $m_u(v)$ represent the “degree of u -ness” possessed by data item v . Then $m_u \in \mathbb{R}_+^V$ for all $u \in U$.
- Example: U could be a set of colors, and for an image $v \in V$, $m_u(v)$ could represent the number of pixels that are of color u .

Data Subset Selection

- Suppose we are given a large data set $\mathcal{D} = \{x_i\}_{i=1}^n$ of n data items $V = \{v_1, v_2, \dots, v_n\}$ and we wish to choose a subset $A \subset V$ of items that is good in some way (e.g., a summary).
- Suppose moreover each data item $v \in V$ is described by a vector of non-negative scores for a set U of **features** (or “properties”, or “concepts”, etc.) of each data item.
- That is, for $u \in U$ and $v \in V$, let $m_u(v)$ represent the “degree of u -ness” possessed by data item v . Then $m_u \in \mathbb{R}_+^V$ for all $u \in U$.
- Example: U could be a set of colors, and for an image $v \in V$, $m_u(v)$ could represent the number of pixels that are of color u .
- Example: U might be a set of textual features (e.g., ngrams), and $m_u(v)$ is the number of ngrams of type u in sentence v . E.g., if a document consists of the sentence

$v =$ “Whenever I go to New York City, I visit the New York City museum.”

then $m_{\text{the}}(v) = 1$ while $m_{\text{New York City}}(v) = 2$.

Data Subset Selection

- For $X \subseteq V$, define $m_u(X) = \sum_{x \in X} m_u(x)$, so $m_u(X)$ is a modular function representing the “degree of u -ness” in subset X .

Data Subset Selection

- For $X \subseteq V$, define $m_u(X) = \sum_{x \in X} m_u(x)$, so $m_u(X)$ is a modular function representing the “degree of u -ness” in subset X .
- Since $m_u(X)$ is modular, it does not have a diminishing returns property. I.e., as we add to X , the degree of u -ness grows additively.

Data Subset Selection

- For $X \subseteq V$, define $m_u(X) = \sum_{x \in X} m_u(x)$, so $m_u(X)$ is a modular function representing the “degree of u -ness” in subset X .
- Since $m_u(X)$ is modular, it does not have a diminishing returns property. I.e., as we add to X , the degree of u -ness grows additively.
- With g non-decreasing concave, $g(m_u(X))$ grows subadditively (if we add v to a context A with less u -ness, the u -ness benefit is more than if we add v to a context $B \supseteq A$ having more u -ness). That is

$$g(m_u(A + v)) - g(m_u(A)) \geq g(m_u(B + v)) - g(m_u(B)) \quad (1.25)$$

Data Subset Selection

- For $X \subseteq V$, define $m_u(X) = \sum_{x \in X} m_u(x)$, so $m_u(X)$ is a modular function representing the “degree of u -ness” in subset X .
- Since $m_u(X)$ is modular, it does not have a diminishing returns property. I.e., as we add to X , the degree of u -ness grows additively.
- With g non-decreasing concave, $g(m_u(X))$ grows subadditively (if we add v to a context A with less u -ness, the u -ness benefit is more than if we add v to a context $B \supseteq A$ having more u -ness). That is

$$g(m_u(A + v)) - g(m_u(A)) \geq g(m_u(B + v)) - g(m_u(B)) \quad (1.25)$$

- Consider the following class of feature functions $f : 2^V \rightarrow \mathbb{R}_+$

$$f(X) = \sum_{u \in U} \alpha_u g_u(m_u(X)) \quad (1.26)$$

where g_u is a non-decreasing concave, and $\alpha_u \geq 0$ is a feature importance weight. Thus, f is submodular.

Data Subset Selection

- For $X \subseteq V$, define $m_u(X) = \sum_{x \in X} m_u(x)$, so $m_u(X)$ is a modular function representing the “degree of u -ness” in subset X .
- Since $m_u(X)$ is modular, it does not have a diminishing returns property. I.e., as we add to X , the degree of u -ness grows additively.
- With g non-decreasing concave, $g(m_u(X))$ grows subadditively (if we add v to a context A with less u -ness, the u -ness benefit is more than if we add v to a context $B \supseteq A$ having more u -ness). That is

$$g(m_u(A + v)) - g(m_u(A)) \geq g(m_u(B + v)) - g(m_u(B)) \quad (1.25)$$

- Consider the following class of feature functions $f : 2^V \rightarrow \mathbb{R}_+$

$$f(X) = \sum_{u \in U} \alpha_u g_u(m_u(X)) \quad (1.26)$$

where g_u is a non-decreasing concave, and $\alpha_u \geq 0$ is a feature importance weight. Thus, f is submodular.

- $f(X)$ measures X 's ability to represent set of features U as measured by $m_u(X)$, with diminishing returns function g , and importance weights α_u .

Data Subset Selection, KL-divergence

- Let $p = \{p_u\}_{u \in U}$ be a desired probability distribution over features (i.e., $\sum_u p_u = 1$ and $p_u \geq 0$ for all $u \in U$).

Data Subset Selection, KL-divergence

- Let $p = \{p_u\}_{u \in U}$ be a desired probability distribution over features (i.e., $\sum_u p_u = 1$ and $p_u \geq 0$ for all $u \in U$).
- Next, normalize the modular weights for each feature:

$$0 \leq \bar{m}_u(X) \triangleq \frac{m_u(X)}{\sum_{u' \in U} m_{u'}(X)} = \frac{m_u(X)}{m(X)} \leq 1 \quad (1.27)$$

where $m(X) \triangleq \sum_{u' \in U} m_{u'}(X)$.

Data Subset Selection, KL-divergence

- Let $p = \{p_u\}_{u \in U}$ be a desired probability distribution over features (i.e., $\sum_u p_u = 1$ and $p_u \geq 0$ for all $u \in U$).
- Next, normalize the modular weights for each feature:

$$0 \leq \bar{m}_u(X) \triangleq \frac{m_u(X)}{\sum_{u' \in U} m_{u'}(X)} = \frac{m_u(X)}{m(X)} \leq 1 \quad (1.27)$$

where $m(X) \triangleq \sum_{u' \in U} m_{u'}(X)$.

- Then $\bar{m}_u(X)$ can also be seen as a distribution over features U since $\bar{m}_u(X) \geq 0$ and $\sum_{u \in U} \bar{m}_u(X) = 1$ for any $X \subseteq V$.

Data Subset Selection, KL-divergence

- Let $p = \{p_u\}_{u \in U}$ be a desired probability distribution over features (i.e., $\sum_u p_u = 1$ and $p_u \geq 0$ for all $u \in U$).
- Next, normalize the modular weights for each feature:

$$0 \leq \bar{m}_u(X) \triangleq \frac{m_u(X)}{\sum_{u' \in U} m_{u'}(X)} = \frac{m_u(X)}{m(X)} \leq 1 \quad (1.27)$$

where $m(X) \triangleq \sum_{u' \in U} m_{u'}(X)$.

- Then $\bar{m}_u(X)$ can also be seen as a distribution over features U since $\bar{m}_u(X) \geq 0$ and $\sum_{u \in U} \bar{m}_u(X) = 1$ for any $X \subseteq V$.
- Consider the KL-divergence between these two distributions:

$$D(p \parallel \{\bar{m}_u(X)\}_{u \in U}) = \sum_{u \in U} p_u \log p_u - \sum_{u \in U} p_u \log(\bar{m}_u(X)) \quad (1.28)$$

$$\begin{aligned} &= \sum_{u \in U} p_u \log p_u - \sum_{u \in U} p_u \log(m_u(X)) + \log(m(X)) \\ &= -H(p) + \log m(X) - \sum_{u \in U} p_u \log(m_u(X)) \quad (1.29) \end{aligned}$$

Data Subset Selection, KL-divergence

- The objective once again, treating entropy $H(p)$ as a constant,

$$D(p||\{\tilde{m}_u(X)\}) = \text{const.} + \log m(X) - \sum_{u \in U} p_u \log(m_u(X)) \quad (1.30)$$

Data Subset Selection, KL-divergence

- The objective once again, treating entropy $H(p)$ as a constant,

$$D(p||\{\tilde{m}_u(X)\}) = \text{const.} + \log m(X) - \sum_{u \in U} p_u \log(m_u(X)) \quad (1.30)$$

- But seen as a function of X , both $\log m(X)$ and $\sum_{u \in U} p_u \log m_u(X)$ are submodular functions.

Data Subset Selection, KL-divergence

- The objective once again, treating entropy $H(p)$ as a constant,

$$D(p||\{\tilde{m}_u(X)\}) = \text{const.} + \log m(X) - \sum_{u \in U} p_u \log(m_u(X)) \quad (1.30)$$

- But seen as a function of X , both $\log m(X)$ and $\sum_{u \in U} p_u \log m_u(X)$ are submodular functions.
- Hence the KL-divergence, seen as a function of X , i.e., $f(X) = D(p||\{\tilde{m}_u(X)\})$ is quite naturally represented as a **difference of submodular functions**.

Data Subset Selection, KL-divergence

- The objective once again, treating entropy $H(p)$ as a constant,

$$D(p||\{\bar{m}_u(X)\}) = \text{const.} + \log m(X) - \sum_{u \in U} p_u \log(m_u(X)) \quad (1.30)$$

- But seen as a function of X , both $\log m(X)$ and $\sum_{u \in U} p_u \log m_u(X)$ are submodular functions.
- Hence the KL-divergence, seen as a function of X , i.e., $f(X) = D(p||\{\bar{m}_u(X)\})$ is quite naturally represented as a **difference of submodular functions**.
- Alternatively, if we define (Shinohara, 2014)

$$g(X) \triangleq \log m(X) - D(p||\{\bar{m}_u(X)\}) = \sum_{u \in U} p_u \log(m_u(X)) \quad (1.31)$$

we have a **submodular function** g that represents a combination of its quantity of X via its features (i.e., $\log m(X)$) and its feature distribution closeness to some distribution p (i.e., $D(p||\{\bar{m}_u(X)\})$).

Information Gain for Sensor Placement

- Given an environment, V is set of candidate locations for placement of a sensor (e.g., temperature, gas, audio, video, bacteria or other environmental contaminant, etc.).

Information Gain for Sensor Placement

- Given an environment, V is set of candidate locations for placement of a sensor (e.g., temperature, gas, audio, video, bacteria or other environmental contaminant, etc.).
- We have a function $f(A)$ that measures the “coverage” of any given set A of sensor placement decisions. If a point is covered, we can answer a question about it (i.e., temperature, degree of contaminant).

Information Gain for Sensor Placement

- Given an environment, V is set of candidate locations for placement of a sensor (e.g., temperature, gas, audio, video, bacteria or other environmental contaminant, etc.).
- We have a function $f(A)$ that measures the “coverage” of any given set A of sensor placement decisions. If a point is covered, we can answer a question about it (i.e., temperature, degree of contaminant).
- $f(V)$ is maximum coverage.

Information Gain for Sensor Placement

- Given an environment, V is set of candidate locations for placement of a sensor (e.g., temperature, gas, audio, video, bacteria or other environmental contaminant, etc.).
- We have a function $f(A)$ that measures the “coverage” of any given set A of sensor placement decisions. If a point is covered, we can answer a question about it (i.e., temperature, degree of contaminant).
- $f(V)$ is maximum coverage.
- One possible goal: choose smallest set A such that $f(A) \geq \alpha f(V)$ with $0 < \alpha \leq 1$ (recall the submodular set cover problem)

Information Gain for Sensor Placement

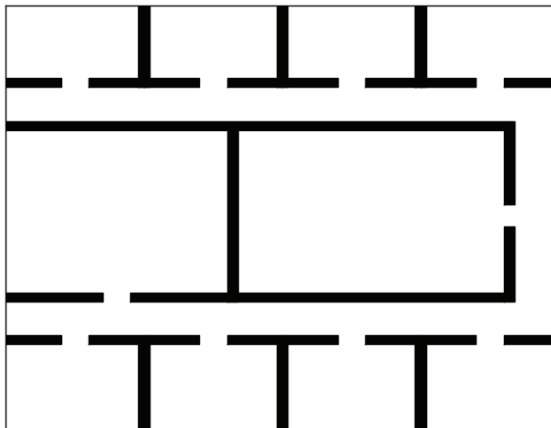
- Given an environment, V is set of candidate locations for placement of a sensor (e.g., temperature, gas, audio, video, bacteria or other environmental contaminant, etc.).
- We have a function $f(A)$ that measures the “coverage” of any given set A of sensor placement decisions. If a point is covered, we can answer a question about it (i.e., temperature, degree of contaminant).
- $f(V)$ is maximum coverage.
- One possible goal: choose smallest set A such that $f(A) \geq \alpha f(V)$ with $0 < \alpha \leq 1$ (recall the submodular set cover problem)
- Another possible goal: choose size at most k set A such that $f(A)$ is maximized.

Information Gain for Sensor Placement

- Given an environment, V is set of candidate locations for placement of a sensor (e.g., temperature, gas, audio, video, bacteria or other environmental contaminant, etc.).
- We have a function $f(A)$ that measures the “coverage” of any given set A of sensor placement decisions. If a point is covered, we can answer a question about it (i.e., temperature, degree of contaminant).
- $f(V)$ is maximum coverage.
- One possible goal: choose smallest set A such that $f(A) \geq \alpha f(V)$ with $0 < \alpha \leq 1$ (recall the submodular set cover problem)
- Another possible goal: choose size at most k set A such that $f(A)$ is maximized.
- Environment could be a floor of a building, water network, monitored ecological preservation.

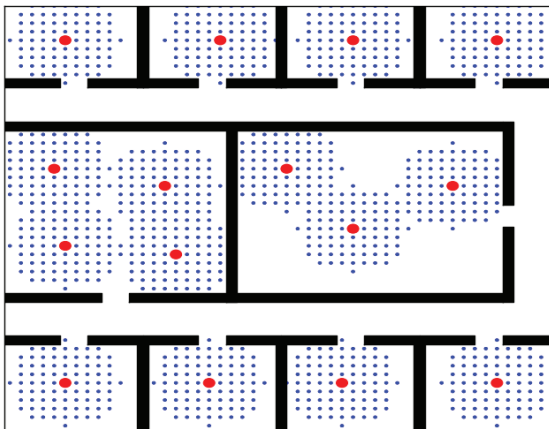
Sensor Placement within Buildings

- An example of a room layout. Should be possible to determine temperature at all points in the room. Sensors cannot sense beyond wall (thick black line) boundaries.



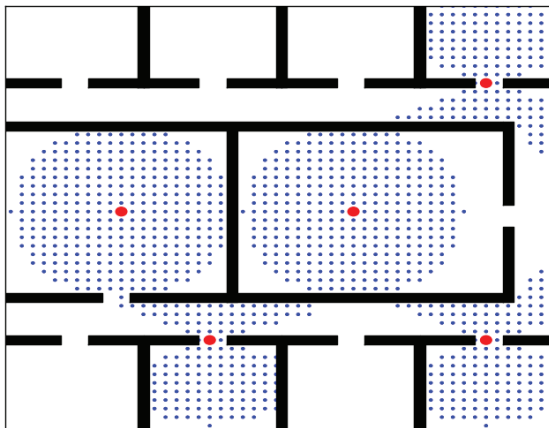
Sensor Placement within Buildings

- Example sensor placement using small range cheap sensors (located at red dots).



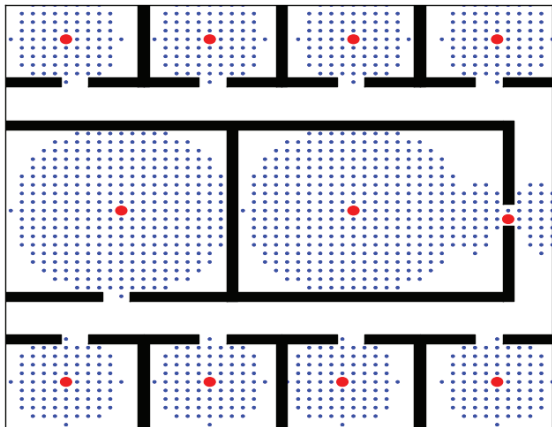
Sensor Placement within Buildings

- Example sensor placement using longer range expensive sensors (located at red dots).



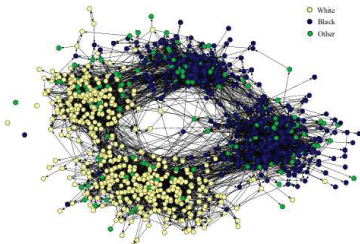
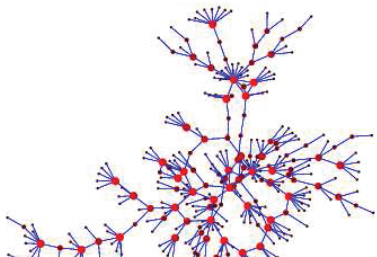
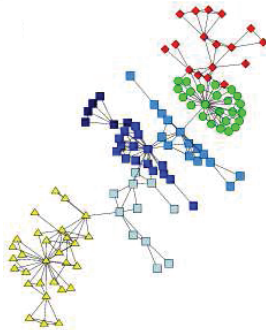
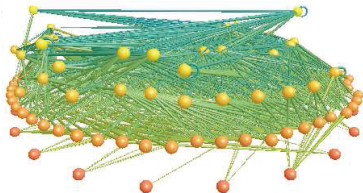
Sensor Placement within Buildings

- Example sensor placement using mixed range sensors (located at red dots).



Social Networks

(from Newman, 2004). Clockwise from top left: 1) predator-prey interactions, 2) scientific collaborations, 3) sexual contact, 4) school friendships.



The value of a friend



- Let V be a set of individuals in a network. How valuable is a given friend $v \in V$?

The value of a friend



- Let V be a set of individuals in a network. How valuable is a given friend $v \in V$? It depends on how many friends you have.

The value of a friend



- Let V be a set of individuals in a network. How valuable is a given friend $v \in V$? It depends on how many friends you have.
- Valuate a group of friends $S \subseteq V$ via set function $f(S)$.

The value of a friend



- Let V be a set of individuals in a network. How valuable is a given friend $v \in V$? It depends on how many friends you have.
- Valuate a group of friends $S \subseteq V$ via set function $f(S)$.
- A submodular model: a friend becomes less marginally valuable as your set of friends grows.

The value of a friend



- Let V be a set of individuals in a network. How valuable is a given friend $v \in V$? It depends on how many friends you have.
- Valueate a group of friends $S \subseteq V$ via set function $f(S)$.
- A submodular model: a friend becomes less marginally valuable as your set of friends grows.
- Supermodular model: a friend becomes **more** valuable the more friends you have.

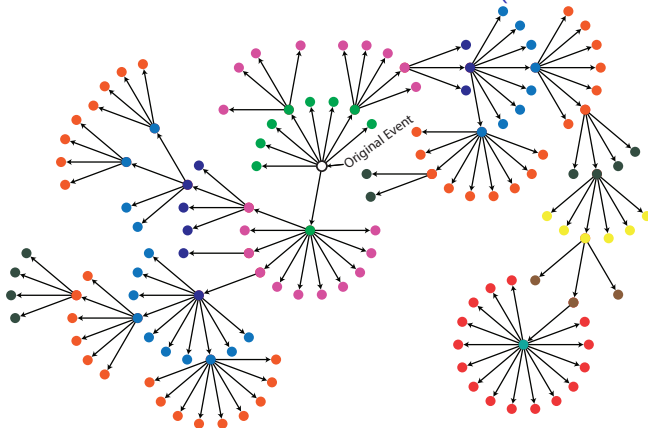
The value of a friend



- Let V be a set of individuals in a network. How valuable is a given friend $v \in V$? It depends on how many friends you have.
- Valuate a group of friends $S \subseteq V$ via set function $f(S)$.
- A submodular model: a friend becomes less marginally valuable as your set of friends grows.
- Supermodular model: a friend becomes **more** valuable the more friends you have.
- Which is a better model?

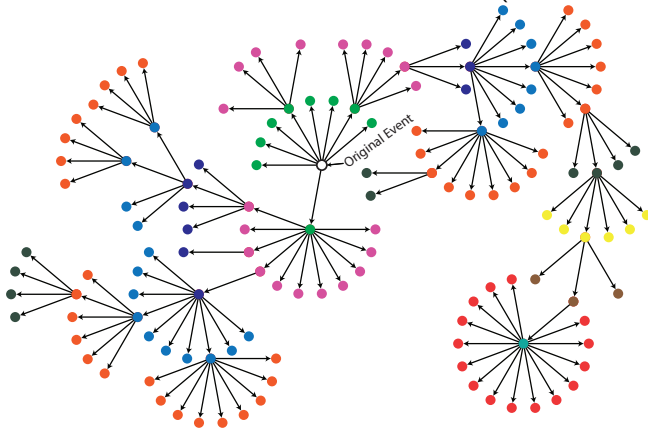
Information Cascades, Diffusion Networks

- How to model flow of information from source to the point it reaches users — information used in its common sense (like news events).



Information Cascades, Diffusion Networks

- How to model flow of information from source to the point it reaches users — information used in its common sense (like news events).



- Goal: How to find the most influential sources, the ones that often set off cascades, which are like large “waves” of information flow?

Diffusion Networks

Where are they useful?

- **Information propagation:** when blogs or news stories break, and creates an information cascade over multiple other blogs/newspapers/magazines.
- **Viral marketing:** What is the pattern of trendsetters that cause an individual to purchase a product?
- **Epidemiology:** who gets sick from whom? What is the infection network of such links? Given finite supply of vaccine, who to inoculate to protect overall population (cut the network)?
 - Infer the connectivity of a network (memes, purchase decisions, viruses, etc.) based only on diffusion traces (the time that each node is “infected”)?
 - How to find the most likely tree or graph?

A model of influence in social networks

- Given a graph $G = (V, E)$, each $v \in V$ corresponds to a person, to each v we have an activation function $f_v : 2^V \rightarrow [0, 1]$ dependent only on its neighbors. I.e., $f_v(A) = f_v(A \cap \Gamma(v))$.
- Goal - Viral Marketing: find a small subset $S \subseteq V$ of individuals to directly influence, and thus indirectly influence the greatest number of possible other individuals (via the social network G).
- Define function $f : 2^V \rightarrow \mathbb{Z}^+$ to model the ultimate influence of an initial infected nodes S . Use following iterative process; at each step:
 - Given previous set of infected nodes S that have not yet had their chance to infect their neighbors,
 - activate new nodes $v \in V \setminus S$ if $f_v(S \cap \Gamma_v) \geq U[0, 1]$, where $U[0, 1]$ is a uniform random number between 0 and 1, and Γ_v are the neighbors of v .
- For many f_v (including simple linear functions, and where f_v is submodular itself), we can show f is submodular (Kempe, Kleinberg, Tardos 1993).

Optimization Problem Involving Network Externalities

- (From Mirrokni, Roch, Sundararajan 2012): Let V be a set of users.

Optimization Problem Involving Network Externalities

- (From Mirrokni, Roch, Sundararajan 2012): Let V be a set of users.
- Let $v_i(S)$ be the value that user i has for a good if $S \subseteq V$ already own the good — e.g. $v_i(S) = \omega_i + f_i(\sum_{j \in S} w_{ij})$ where ω_i is inherent value, and f_i might be a concave function, and w_{ij} is how important $j \in S$ is to i (e.g., a network). Weights might be random.

Optimization Problem Involving Network Externalities

- (From Mirrokni, Roch, Sundararajan 2012): Let V be a set of users.
- Let $v_i(S)$ be the value that user i has for a good if $S \subseteq V$ already own the good — e.g. $v_i(S) = \omega_i + f_i(\sum_{j \in S} w_{ij})$ where ω_i is inherent value, and f_i might be a concave function, and w_{ij} is how important $j \in S$ is to i (e.g., a network). Weights might be random.
- Given price p for good, user i buys good if $v_i(S) \geq p$.

Optimization Problem Involving Network Externalities

- (From Mirrokni, Roch, Sundararajan 2012): Let V be a set of users.
- Let $v_i(S)$ be the value that user i has for a good if $S \subseteq V$ already own the good — e.g. $v_i(S) = \omega_i + f_i(\sum_{j \in S} w_{ij})$ where ω_i is inherent value, and f_i might be a concave function, and w_{ij} is how important $j \in S$ is to i (e.g., a network). Weights might be random.
- Given price p for good, user i buys good if $v_i(S) \geq p$.
- We choose initial price p and initial set of users $A \subseteq V$ who get the good for free.

Optimization Problem Involving Network Externalities

- (From Mirrokni, Roch, Sundararajan 2012): Let V be a set of users.
- Let $v_i(S)$ be the value that user i has for a good if $S \subseteq V$ already own the good — e.g. $v_i(S) = \omega_i + f_i(\sum_{j \in S} w_{ij})$ where ω_i is inherent value, and f_i might be a concave function, and w_{ij} is how important $j \in S$ is to i (e.g., a network). Weights might be random.
- Given price p for good, user i buys good if $v_i(S) \geq p$.
- We choose initial price p and initial set of users $A \subseteq V$ who get the good for free.
- Define $S_1 = \{i \notin A : v_i(A) \geq p\}$ initial set of buyers.

Optimization Problem Involving Network Externalities

- (From Mirrokni, Roch, Sundararajan 2012): Let V be a set of users.
- Let $v_i(S)$ be the value that user i has for a good if $S \subseteq V$ already own the good — e.g. $v_i(S) = \omega_i + f_i(\sum_{j \in S} w_{ij})$ where ω_i is inherent value, and f_i might be a concave function, and w_{ij} is how important $j \in S$ is to i (e.g., a network). Weights might be random.
- Given price p for good, user i buys good if $v_i(S) \geq p$.
- We choose initial price p and initial set of users $A \subseteq V$ who get the good for free.
- Define $S_1 = \{i \notin A : v_i(A) \geq p\}$ initial set of buyers.
- $S_2 = \{i \notin A \cup S_1 : v_i(A \cup S_1) \geq p\}$.

Optimization Problem Involving Network Externalities

- (From Mirrokni, Roch, Sundararajan 2012): Let V be a set of users.
- Let $v_i(S)$ be the value that user i has for a good if $S \subseteq V$ already own the good — e.g. $v_i(S) = \omega_i + f_i(\sum_{j \in S} w_{ij})$ where ω_i is inherent value, and f_i might be a concave function, and w_{ij} is how important $j \in S$ is to i (e.g., a network). Weights might be random.
- Given price p for good, user i buys good if $v_i(S) \geq p$.
- We choose initial price p and initial set of users $A \subseteq V$ who get the good for free.
- Define $S_1 = \{i \notin A : v_i(A) \geq p\}$ initial set of buyers.
- $S_2 = \{i \notin A \cup S_1 : v_i(A \cup S_1) \geq p\}$.
- This starts a cascade. Let

$$S_k = \{i \notin \cup_{j < k} S_j \cup A : v_i(\cup_{j < k} S_j \cup A) \geq p\},$$

Optimization Problem Involving Network Externalities

- (From Mirrokni, Roch, Sundararajan 2012): Let V be a set of users.
- Let $v_i(S)$ be the value that user i has for a good if $S \subseteq V$ already own the good — e.g. $v_i(S) = \omega_i + f_i(\sum_{j \in S} w_{ij})$ where ω_i is inherent value, and f_i might be a concave function, and w_{ij} is how important $j \in S$ is to i (e.g., a network). Weights might be random.
- Given price p for good, user i buys good if $v_i(S) \geq p$.
- We choose initial price p and initial set of users $A \subseteq V$ who get the good for free.
- Define $S_1 = \{i \notin A : v_i(A) \geq p\}$ initial set of buyers.
- $S_2 = \{i \notin A \cup S_1 : v_i(A \cup S_1) \geq p\}$.
- This starts a cascade. Let

$$S_k = \{i \notin \cup_{j < k} S_j \cup A : v_i(\cup_{j < k} S_j \cup A) \geq p\},$$
- and let S_{k^*} be the saturation point, lowest value of k such that $S_k = S_{k+1}$

Optimization Problem Involving Network Externalities

- (From Mirrokni, Roch, Sundararajan 2012): Let V be a set of users.
- Let $v_i(S)$ be the value that user i has for a good if $S \subseteq V$ already own the good — e.g. $v_i(S) = \omega_i + f_i(\sum_{j \in S} w_{ij})$ where ω_i is inherent value, and f_i might be a concave function, and w_{ij} is how important $j \in S$ is to i (e.g., a network). Weights might be random.
- Given price p for good, user i buys good if $v_i(S) \geq p$.
- We choose initial price p and initial set of users $A \subseteq V$ who get the good for free.
- Define $S_1 = \{i \notin A : v_i(A) \geq p\}$ initial set of buyers.
- $S_2 = \{i \notin A \cup S_1 : v_i(A \cup S_1) \geq p\}$.
- This starts a cascade. Let

$$S_k = \{i \notin \cup_{j < k} S_j \cup A : v_i(\cup_{j < k} S_j \cup A) \geq p\},$$
- and let S_{k^*} be the saturation point, lowest value of k such that $S_k = S_{k+1}$
- Goal: find A and p to maximize $f_p(A) = \mathbb{E}[p \times |S_{k^*}|]$.

Graphical Model Structure Learning

- A probability distribution on binary vectors $p : \{0, 1\}^V \rightarrow [0, 1]$:

$$p(x) = \frac{1}{Z} \exp(-E(x)) \quad (1.32)$$

where $E(x)$ is the energy function.

Graphical Model Structure Learning

- A probability distribution on binary vectors $p : \{0, 1\}^V \rightarrow [0, 1]$:

$$p(x) = \frac{1}{Z} \exp(-E(x)) \quad (1.32)$$

where $E(x)$ is the energy function.

- A graphical model $G = (V, \mathcal{E})$ represents a family of probability distributions $p \in \mathcal{F}(G)$ all of which factor w.r.t. the graph.

Graphical Model Structure Learning

- A probability distribution on binary vectors $p : \{0, 1\}^V \rightarrow [0, 1]$:

$$p(x) = \frac{1}{Z} \exp(-E(x)) \quad (1.32)$$

where $E(x)$ is the energy function.

- A graphical model $G = (V, \mathcal{E})$ represents a family of probability distributions $p \in \mathcal{F}(G)$ all of which factor w.r.t. the graph.
- I.e., if \mathcal{C} are a set of cliques of graph G , then we must have:

$$E(x) = \sum_{c \in \mathcal{C}} E_c(x_c) \quad (1.33)$$

Graphical Model Structure Learning

- A probability distribution on binary vectors $p : \{0, 1\}^V \rightarrow [0, 1]$:

$$p(x) = \frac{1}{Z} \exp(-E(x)) \quad (1.32)$$

where $E(x)$ is the energy function.

- A graphical model $G = (V, \mathcal{E})$ represents a family of probability distributions $p \in \mathcal{F}(G)$ all of which factor w.r.t. the graph.
- I.e., if \mathcal{C} are a set of cliques of graph G , then we must have:

$$E(x) = \sum_{c \in \mathcal{C}} E_c(x_c) \quad (1.33)$$

- The problem of **structure learning in graphical models** is to find the graph G based on data.

Graphical Model Structure Learning

- A probability distribution on binary vectors $p : \{0, 1\}^V \rightarrow [0, 1]$:

$$p(x) = \frac{1}{Z} \exp(-E(x)) \quad (1.32)$$

where $E(x)$ is the energy function.

- A graphical model $G = (V, \mathcal{E})$ represents a family of probability distributions $p \in \mathcal{F}(G)$ all of which factor w.r.t. the graph.
- I.e., if \mathcal{C} are a set of cliques of graph G , then we must have:

$$E(x) = \sum_{c \in \mathcal{C}} E_c(x_c) \quad (1.33)$$

- The problem of **structure learning in graphical models** is to find the graph G based on data.
- This can be viewed as a discrete optimization problem on the potential (undirected) **edges** of the graph $V \times V$.

Graphical Models: Learning Tree Distributions

- Goal: find the closest distribution p_t to p subject to p_t factoring w.r.t. some tree $T = (V, F)$, i.e., $p_t \in \mathcal{F}(T, \mathcal{M})$.

Graphical Models: Learning Tree Distributions

- Goal: find the closest distribution p_t to p subject to p_t factoring w.r.t. some tree $T = (V, F)$, i.e., $p_t \in \mathcal{F}(T, \mathcal{M})$.
- This can be expressed as a discrete optimization problem:

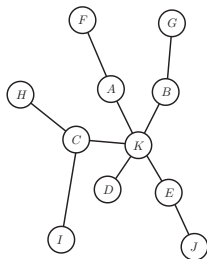
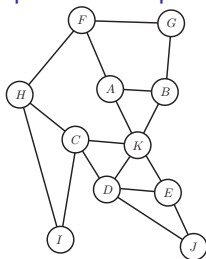
minimize
 $p_t \in \mathcal{F}(G, \mathcal{M})$

$$D(p || p_t)$$

subject to

$$p_t \in \mathcal{F}(T, \mathcal{M}).$$

$T = (V, F)$ is a tree



Graphical Models: Learning Tree Distributions

- Goal: find the closest distribution p_t to p subject to p_t factoring w.r.t. some tree $T = (V, F)$, i.e., $p_t \in \mathcal{F}(T, \mathcal{M})$.
- This can be expressed as a discrete optimization problem:

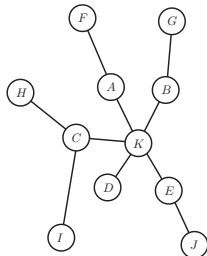
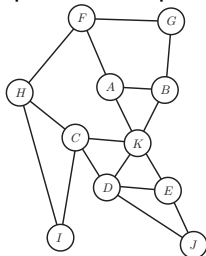
minimize
 $p_t \in \mathcal{F}(G, \mathcal{M})$

$$D(p || p_t)$$

subject to

$$p_t \in \mathcal{F}(T, \mathcal{M}).$$

$T = (V, F)$ is a tree



- Discrete problem: choose the optimal set of edges $A \subseteq E$ that constitute tree (i.e., find a spanning tree of G of best quality).

Graphical Models: Learning Tree Distributions

- Goal: find the closest distribution p_t to p subject to p_t factoring w.r.t. some tree $T = (V, F)$, i.e., $p_t \in \mathcal{F}(T, \mathcal{M})$.
- This can be expressed as a discrete optimization problem:

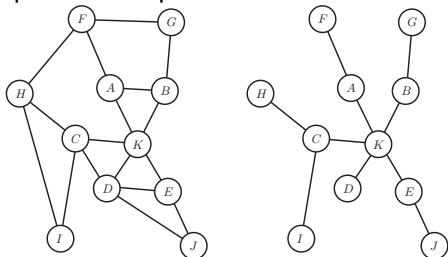
minimize
 $p_t \in \mathcal{F}(G, \mathcal{M})$

$$D(p || p_t)$$

subject to

$$p_t \in \mathcal{F}(T, \mathcal{M}).$$

$T = (V, F)$ is a tree



- Discrete problem: choose the optimal set of edges $A \subseteq E$ that constitute tree (i.e., find a spanning tree of G of best quality).
- Define $f : 2^E \rightarrow \mathbb{R}_+$ where f is a **weighted cycle matroid rank function** (a type of submodular function), with weights $w(e) = w(u, v) = I(X_u; X_v)$ for $e \in E$.

Graphical Models: Learning Tree Distributions

- Goal: find the closest distribution p_t to p subject to p_t factoring w.r.t. some tree $T = (V, F)$, i.e., $p_t \in \mathcal{F}(T, \mathcal{M})$.
- This can be expressed as a discrete optimization problem:

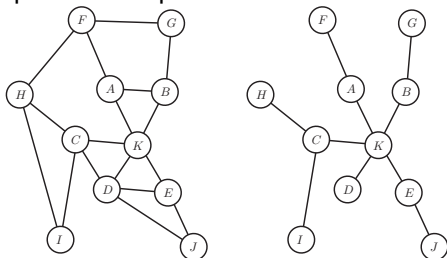
minimize
 $p_t \in \mathcal{F}(G, \mathcal{M})$

$$D(p || p_t)$$

subject to

$$p_t \in \mathcal{F}(T, \mathcal{M}).$$

$T = (V, F)$ is a tree



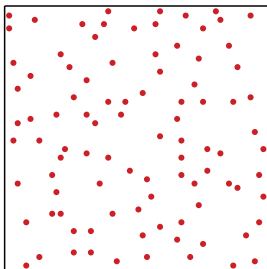
- Discrete problem: choose the optimal set of edges $A \subseteq E$ that constitute tree (i.e., find a spanning tree of G of best quality).
- Define $f : 2^E \rightarrow \mathbb{R}_+$ where f is a **weighted cycle matroid rank** function (a type of submodular function), with weights $w(e) = w(u, v) = I(X_u; X_v)$ for $e \in E$.
- Then finding the maximum weight base of the matroid is solved by the greedy algorithm, and also finds the optimal tree (Chow & Liu, 1968)

Determinantal Point Processes (DPPs)

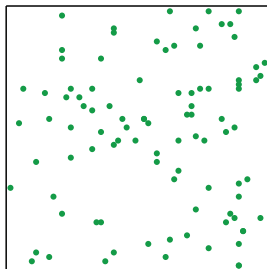
- Sometimes we wish not only to value subsets $A \subseteq V$ but to induce probability distributions over all subsets.

Determinantal Point Processes (DPPs)

- Sometimes we wish not only to valuate subsets $A \subseteq V$ but to induce probability distributions over all subsets.
- We may wish to prefer samples where elements of A are diverse (i.e., given a sample A , for $a, b \in A$, we prefer a and b to be different).



DPP

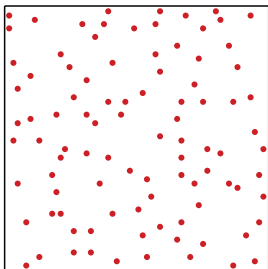


Independent

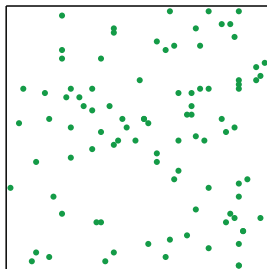
(Kulesza,
Gillen-
water, &
Taskar,
2011)

Determinantal Point Processes (DPPs)

- Sometimes we wish not only to valuate subsets $A \subseteq V$ but to induce probability distributions over all subsets.
- We may wish to prefer samples where elements of A are diverse (i.e., given a sample A , for $a, b \in A$, we prefer a and b to be different).



DPP



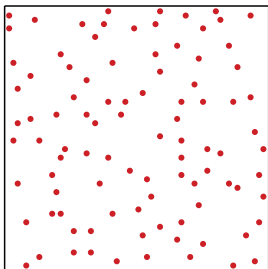
Independent

(Kulesza,
Gillen-
water, &
Taskar,
2011)

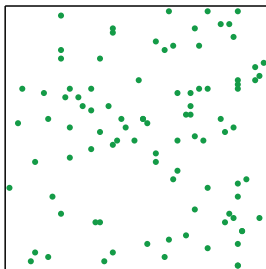
- A Determinantal point processes (DPPs) is a probability distribution over subsets A of V where the “energy” function is submodular.

Determinantal Point Processes (DPPs)

- Sometimes we wish not only to valuate subsets $A \subseteq V$ but to induce probability distributions over all subsets.
- We may wish to prefer samples where elements of A are diverse (i.e., given a sample A , for $a, b \in A$, we prefer a and b to be different).



DPP



Independent

(Kulesza,
Gillen-
water, &
Taskar,
2011)

- A Determinantal point processes (DPPs) is a probability distribution over subsets A of V where the “energy” function is submodular.
- More “diverse” or “complex” samples are given higher probability.

DPPs and log-submodular probability distributions

- Given binary vectors $x, y \in \{0, 1\}^V$, $y \leq x$ if $y(v) \leq x(v), \forall v \in V$.

DPPs and log-submodular probability distributions

- Given binary vectors $x, y \in \{0, 1\}^V$, $y \leq x$ if $y(v) \leq x(v), \forall v \in V$.
- Given a positive-definite $n \times n$ matrix M , a subset $X \subseteq V$, let M_X be $|X| \times |X|$ principle submatrix, rows/columns specified by $X \subseteq V$.

DPPs and log-submodular probability distributions

- Given binary vectors $x, y \in \{0, 1\}^V$, $y \leq x$ if $y(v) \leq x(v), \forall v \in V$.
- Given a positive-definite $n \times n$ matrix M , a subset $X \subseteq V$, let M_X be $|X| \times |X|$ principle submatrix, rows/columns specified by $X \subseteq V$.
- A Determinantal Point Process (DPP) is a distribution of the form:

$$\Pr(\mathbf{X} = x) = \frac{|M_{X(x)}|}{|M + I|} = \exp\left(\log\left(\frac{|M_{X(x)}|}{|M + I|}\right)\right) \propto \det(M_{X(x)}) \quad (1.34)$$

where I is $n \times n$ identity matrix, and $\mathbf{X} \in \{0, 1\}^V$ is a random vector.

DPPs and log-submodular probability distributions

- Given binary vectors $x, y \in \{0, 1\}^V$, $y \leq x$ if $y(v) \leq x(v), \forall v \in V$.
- Given a positive-definite $n \times n$ matrix M , a subset $X \subseteq V$, let M_X be $|X| \times |X|$ principle submatrix, rows/columns specified by $X \subseteq V$.
- A Determinantal Point Process (DPP) is a distribution of the form:

$$\Pr(\mathbf{X} = x) = \frac{|M_{X(x)}|}{|M + I|} = \exp\left(\log\left(\frac{|M_{X(x)}|}{|M + I|}\right)\right) \propto \det(M_{X(x)}) \quad (1.34)$$

where I is $n \times n$ identity matrix, and $\mathbf{X} \in \{0, 1\}^V$ is a random vector.

- Equivalently, defining K as $K = M(M + I)^{-1}$, we have:

$$\sum_{x \in \{0, 1\}^V : x \geq y} \Pr(\mathbf{X} = x) = \Pr(\mathbf{X} \geq y) = \exp\left(\log\left(|K_{Y(y)}|\right)\right) \quad (1.35)$$

DPPs and log-submodular probability distributions

- Given binary vectors $x, y \in \{0, 1\}^V$, $y \leq x$ if $y(v) \leq x(v), \forall v \in V$.
- Given a positive-definite $n \times n$ matrix M , a subset $X \subseteq V$, let M_X be $|X| \times |X|$ principle submatrix, rows/columns specified by $X \subseteq V$.
- A Determinantal Point Process (DPP) is a distribution of the form:

$$\Pr(\mathbf{X} = x) = \frac{|M_{X(x)}|}{|M + I|} = \exp\left(\log\left(\frac{|M_{X(x)}|}{|M + I|}\right)\right) \propto \det(M_{X(x)}) \quad (1.34)$$

where I is $n \times n$ identity matrix, and $\mathbf{X} \in \{0, 1\}^V$ is a random vector.

- Equivalently, defining K as $K = M(M + I)^{-1}$, we have:

$$\sum_{x \in \{0, 1\}^V : x \geq y} \Pr(\mathbf{X} = x) = \Pr(\mathbf{X} \geq y) = \exp\left(\log\left(|K_{Y(y)}|\right)\right) \quad (1.35)$$

- Given positive definite matrix M , function $f : 2^V \rightarrow \mathbb{R}$ with $f(A) = \log |M_A|$ (the logdet function) is submodular.

DPPs and log-submodular probability distributions

- Given binary vectors $x, y \in \{0, 1\}^V$, $y \leq x$ if $y(v) \leq x(v), \forall v \in V$.
- Given a positive-definite $n \times n$ matrix M , a subset $X \subseteq V$, let M_X be $|X| \times |X|$ principle submatrix, rows/columns specified by $X \subseteq V$.
- A Determinantal Point Process (DPP) is a distribution of the form:

$$\Pr(\mathbf{X} = x) = \frac{|M_{X(x)}|}{|M + I|} = \exp\left(\log\left(\frac{|M_{X(x)}|}{|M + I|}\right)\right) \propto \det(M_{X(x)}) \quad (1.34)$$

where I is $n \times n$ identity matrix, and $\mathbf{X} \in \{0, 1\}^V$ is a random vector.

- Equivalently, defining K as $K = M(M + I)^{-1}$, we have:

$$\sum_{x \in \{0, 1\}^V : x \geq y} \Pr(\mathbf{X} = x) = \Pr(\mathbf{X} \geq y) = \exp\left(\log\left(|K_{Y(y)}|\right)\right) \quad (1.35)$$

- Given positive definite matrix M , function $f : 2^V \rightarrow \mathbb{R}$ with $f(A) = \log |M_A|$ (the logdet function) is submodular.
- Therefore, a DPP is a log-submodular probability distribution.

Graphical Models and fast MAP Inference

- Given distribution that factors w.r.t. a graph:

$$p(x) = \frac{1}{Z} \exp(-E(x)) \quad (1.36)$$

where $E(x) = \sum_{c \in \mathcal{C}} E_c(x_c)$ and \mathcal{C} are cliques of graph $G = (V, \mathcal{E})$.

Graphical Models and fast MAP Inference

- Given distribution that factors w.r.t. a graph:

$$p(x) = \frac{1}{Z} \exp(-E(x)) \quad (1.36)$$

where $E(x) = \sum_{c \in \mathcal{C}} E_c(x_c)$ and \mathcal{C} are cliques of graph $G = (V, \mathcal{E})$.

- MAP inference problem is important in ML: compute

$$x^* \in \operatorname{argmax}_{x \in \{0,1\}^V} p(x) \quad (1.37)$$

Graphical Models and fast MAP Inference

- Given distribution that factors w.r.t. a graph:

$$p(x) = \frac{1}{Z} \exp(-E(x)) \quad (1.36)$$

where $E(x) = \sum_{c \in \mathcal{C}} E_c(x_c)$ and \mathcal{C} are cliques of graph $G = (V, \mathcal{E})$.

- MAP inference problem is important in ML: compute

$$x^* \in \operatorname{argmax}_{x \in \{0,1\}^V} p(x) \quad (1.37)$$

- Easy when G a tree, exponential in k (**tree-width** of G) in general.

Graphical Models and fast MAP Inference

- Given distribution that factors w.r.t. a graph:

$$p(x) = \frac{1}{Z} \exp(-E(x)) \quad (1.36)$$

where $E(x) = \sum_{c \in \mathcal{C}} E_c(x_c)$ and \mathcal{C} are cliques of graph $G = (V, \mathcal{E})$.

- MAP inference problem is important in ML: compute

$$x^* \in \operatorname{argmax}_{x \in \{0,1\}^V} p(x) \quad (1.37)$$

- Easy when G a tree, exponential in k (**tree-width** of G) in general.
- Even worse, NP-hard to find the tree-width.

Graphical Models and fast MAP Inference

- Given distribution that factors w.r.t. a graph:

$$p(x) = \frac{1}{Z} \exp(-E(x)) \quad (1.36)$$

where $E(x) = \sum_{c \in \mathcal{C}} E_c(x_c)$ and \mathcal{C} are cliques of graph $G = (V, \mathcal{E})$.

- MAP inference problem is important in ML: compute

$$x^* \in \operatorname{argmax}_{x \in \{0,1\}^V} p(x) \quad (1.37)$$

- Easy when G a tree, exponential in k (**tree-width** of G) in general.
- Even worse, NP-hard to find the tree-width.
- Tree-width can be large even when degree is small (e.g., regular grid graphs have low-degree but $\Omega(\sqrt{n})$ tree-width).

Graphical Models and fast MAP Inference

- Given distribution that factors w.r.t. a graph:

$$p(x) = \frac{1}{Z} \exp(-E(x)) \quad (1.36)$$

where $E(x) = \sum_{c \in \mathcal{C}} E_c(x_c)$ and \mathcal{C} are cliques of graph $G = (V, \mathcal{E})$.

- MAP inference problem is important in ML: compute

$$x^* \in \operatorname{argmax}_{x \in \{0,1\}^V} p(x) \quad (1.37)$$

- Easy when G a tree, exponential in k (**tree-width** of G) in general.
- Even worse, NP-hard to find the tree-width.
- Tree-width can be large even when degree is small (e.g., regular grid graphs have low-degree but $\Omega(\sqrt{n})$ tree-width).
- Many approximate inference strategies utilize additional factorization assumptions (e.g., mean-field, variational inference, expectation propagation, etc).

Graphical Models and fast MAP Inference

- Given distribution that factors w.r.t. a graph:

$$p(x) = \frac{1}{Z} \exp(-E(x)) \quad (1.36)$$

where $E(x) = \sum_{c \in \mathcal{C}} E_c(x_c)$ and \mathcal{C} are cliques of graph $G = (V, \mathcal{E})$.

- MAP inference problem is important in ML: compute

$$x^* \in \operatorname{argmax}_{x \in \{0,1\}^V} p(x) \quad (1.37)$$

- Easy when G a tree, exponential in k (**tree-width** of G) in general.
- Even worse, NP-hard to find the tree-width.
- Tree-width can be large even when degree is small (e.g., regular grid graphs have low-degree but $\Omega(\sqrt{n})$ tree-width).
- Many approximate inference strategies utilize additional factorization assumptions (e.g., mean-field, variational inference, expectation propagation, etc).
- Can we do exact MAP inference in polynomial time regardless of the tree-width, without even knowing the tree-width?

Order-two (edge) graphical models

- Given G let $p \in \mathcal{F}(G, \mathcal{M}^{(f)})$ such that we can write the **global energy** $E(x)$ as a sum of **unary** and **pairwise** potentials:

$$E(x) = \sum_{v \in V(G)} e_v(x_v) + \sum_{(i,j) \in E(G)} e_{ij}(x_i, x_j) \quad (1.38)$$

Order-two (edge) graphical models

- Given G let $p \in \mathcal{F}(G, \mathcal{M}^{(f)})$ such that we can write the **global energy** $E(x)$ as a sum of **unary** and **pairwise** potentials:

$$E(x) = \sum_{v \in V(G)} e_v(x_v) + \sum_{(i,j) \in E(G)} e_{ij}(x_i, x_j) \quad (1.38)$$

- $e_v(x_v)$ and $e_{ij}(x_i, x_j)$ are like local energy potentials.

Order-two (edge) graphical models

- Given G let $p \in \mathcal{F}(G, \mathcal{M}^{(f)})$ such that we can write the **global energy** $E(x)$ as a sum of **unary** and **pairwise** potentials:

$$E(x) = \sum_{v \in V(G)} e_v(x_v) + \sum_{(i,j) \in E(G)} e_{ij}(x_i, x_j) \quad (1.38)$$

- $e_v(x_v)$ and $e_{ij}(x_i, x_j)$ are like local energy potentials.
- Since $\log p(x) = -E(x) + \text{const.}$, the smaller $e_v(x_v)$ or $e_{ij}(x_i, x_j)$ become, the higher the probability becomes.

Order-two (edge) graphical models

- Given G let $p \in \mathcal{F}(G, \mathcal{M}^{(f)})$ such that we can write the **global energy** $E(x)$ as a sum of **unary** and **pairwise** potentials:

$$E(x) = \sum_{v \in V(G)} e_v(x_v) + \sum_{(i,j) \in E(G)} e_{ij}(x_i, x_j) \quad (1.38)$$

- $e_v(x_v)$ and $e_{ij}(x_i, x_j)$ are like local energy potentials.
- Since $\log p(x) = -E(x) + \text{const.}$, the smaller $e_v(x_v)$ or $e_{ij}(x_i, x_j)$ become, the higher the probability becomes.
- Further, say that $D_{X_v} = \{0, 1\}$ (binary), so we have binary random vectors distributed according to $p(x)$.

Order-two (edge) graphical models

- Given G let $p \in \mathcal{F}(G, \mathcal{M}^{(f)})$ such that we can write the **global energy** $E(x)$ as a sum of **unary** and **pairwise** potentials:

$$E(x) = \sum_{v \in V(G)} e_v(x_v) + \sum_{(i,j) \in E(G)} e_{ij}(x_i, x_j) \quad (1.38)$$

- $e_v(x_v)$ and $e_{ij}(x_i, x_j)$ are like local energy potentials.
- Since $\log p(x) = -E(x) + \text{const.}$, the smaller $e_v(x_v)$ or $e_{ij}(x_i, x_j)$ become, the higher the probability becomes.
- Further, say that $D_{X_v} = \{0, 1\}$ (binary), so we have binary random vectors distributed according to $p(x)$.
- Thus, $x \in \{0, 1\}^V$, and finding MPE solution is setting some of the variables to 0 and some to 1, i.e.,

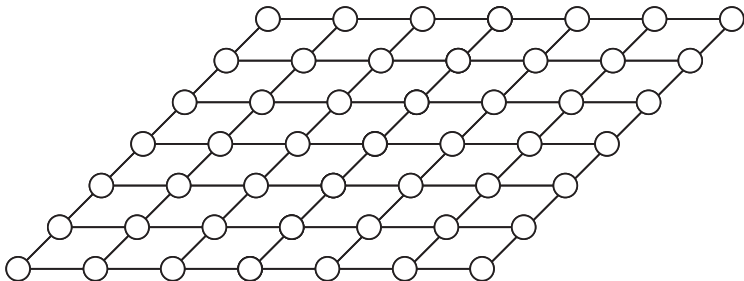
$$\min_{x \in \{0,1\}^V} E(x) \quad (1.39)$$

MRF example

Markov random field

$$\log p(x) \propto \sum_{v \in V(G)} e_v(x_v) + \sum_{(i,j) \in E(G)} e_{ij}(x_i, x_j) \quad (1.40)$$

When G is a 2D grid graph, we have



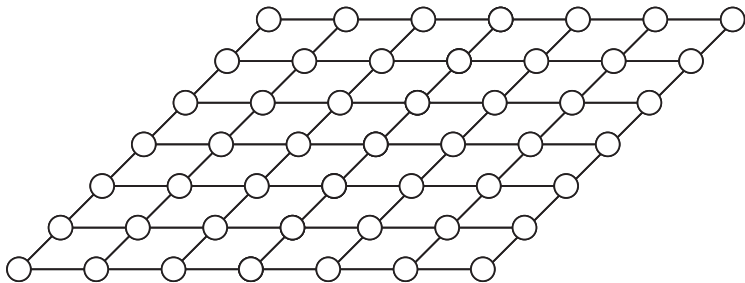
Create an auxiliary graph

- We can create auxiliary graph G_a that involves two new “terminal” nodes s and t and all of the original “non-terminal” nodes $v \in V(G)$.
- The non-terminal nodes represent the original random variables $x_v, v \in V$.
- Starting with the original grid-graph amongst the vertices $v \in V$, we connect each of s and t to all of the original nodes.
- I.e., we form $G_a = (V \cup \{s, t\}, E + \cup_{v \in V} ((s, v) \cup (v, t)))$.

Transformation from graphical model to auxiliary graph

Original 2D-grid graphical model G and energy function

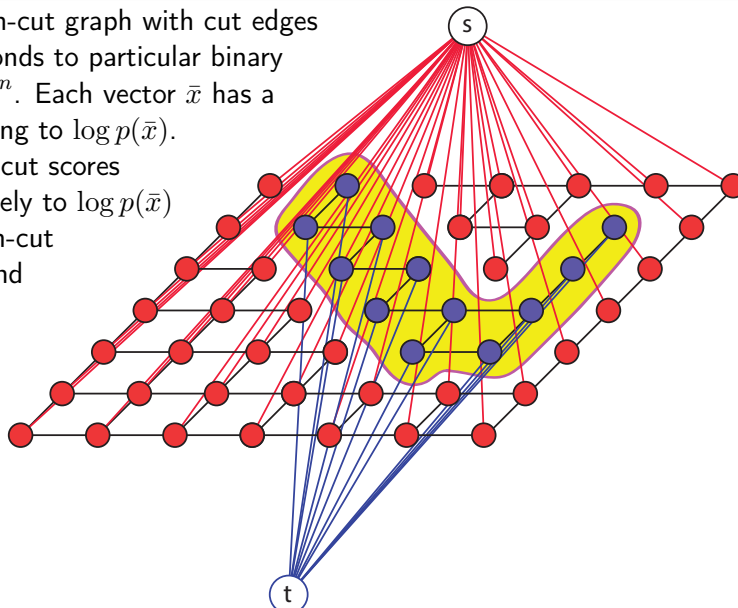
$E(x) = \sum_{v \in V(G)} e_v(x_v) + \sum_{(i,j) \in E(G)} e_{ij}(x_i, x_j)$ needing to be minimized over $x \in \{0, 1\}^V$. Recall, tree-width is $O(\sqrt{|V|})$.



Transformation from graphical model to auxiliary graph

Augmented graph-cut graph with cut edges removed corresponds to particular binary vector $\bar{x} \in \{0, 1\}^n$. Each vector \bar{x} has a score corresponding to $\log p(\bar{x})$.

When can graph cut scores correspond precisely to $\log p(\bar{x})$ in a way that min-cut algorithms can find minimum of energy $E(x)$?



Setting of the weights in the auxiliary cut graph

- Any graph cut corresponds to a vector $\bar{x} \in \{0, 1\}^n$.
- If weights of all edges, except those involving terminals s and t , are non-negative, graph cut computable in polynomial time via max-flow (many algorithms, e.g., Edmonds&Karp $O(nm^2)$ or $O(n^2m \log(nC))$; Goldberg&Tarjan $O(nm \log(n^2/m))$), see Schrijver, page 161).
- If weights are set correctly in the cut graph, and if edge functions e_{ij} satisfy certain properties, then graph-cut score corresponding to \bar{x} can be made equivalent to $E(x) = \log p(\bar{x}) + \text{const.}$.
- Hence, poly time graph cut, can find the optimal MPE assignment, regardless of the graphical model's tree-width!
- In general, finding MPE is an NP-hard optimization problem.

Submodular potentials

submodularity is what allows graph cut to find exact solution

- Edge functions must be **submodular** (in the binary case, equivalent to “associative”, “attractive”, “regular”, “Potts”, or “ferromagnetic”): for all $(i, j) \in E(G)$, must have:

$$e_{ij}(0, 1) + e_{ij}(1, 0) \geq e_{ij}(1, 1) + e_{ij}(0, 0) \quad (1.48)$$

Submodular potentials

submodularity is what allows graph cut to find exact solution

- Edge functions must be **submodular** (in the binary case, equivalent to “associative”, “attractive”, “regular”, “Potts”, or “ferromagnetic”): for all $(i, j) \in E(G)$, must have:

$$e_{ij}(0, 1) + e_{ij}(1, 0) \geq e_{ij}(1, 1) + e_{ij}(0, 0) \quad (1.48)$$

- This means: on average, preservation is preferred over change.

Submodular potentials

submodularity is what allows graph cut to find exact solution

- Edge functions must be **submodular** (in the binary case, equivalent to “associative”, “attractive”, “regular”, “Potts”, or “ferromagnetic”): for all $(i, j) \in E(G)$, must have:

$$e_{ij}(0, 1) + e_{ij}(1, 0) \geq e_{ij}(1, 1) + e_{ij}(0, 0) \quad (1.48)$$

- This means: **on average, preservation is preferred over change.**
- As a set function, this is the same as:

$$f(X) = \sum_{\{i,j\} \in \mathcal{E}(G)} f_{i,j}(X \cap \{i, j\}) \quad (1.49)$$

which is submodular if each of the $f_{i,j}$'s are submodular!

Submodular potentials

submodularity is what allows graph cut to find exact solution

- Edge functions must be **submodular** (in the binary case, equivalent to “associative”, “attractive”, “regular”, “Potts”, or “ferromagnetic”): for all $(i, j) \in E(G)$, must have:

$$e_{ij}(0, 1) + e_{ij}(1, 0) \geq e_{ij}(1, 1) + e_{ij}(0, 0) \quad (1.48)$$

- This means: **on average, preservation is preferred over change.**
- As a set function, this is the same as:

$$f(X) = \sum_{\{i,j\} \in \mathcal{E}(G)} f_{i,j}(X \cap \{i, j\}) \quad (1.49)$$

which is submodular if each of the $f_{i,j}$'s are submodular!

- A special case of more general submodular functions – unconstrained submodular function minimization is solvable in polytime.

On log-supermodular vs. log-submodular distributions

- Log-supermodular distributions.

$$\log \Pr(x) = g(x) + \text{const.} = -E(x) + \text{const.} \quad (1.50)$$

where g is supermodular ($E(x) = -g(x)$ is submodular). MAP (or high-probable) assignments should be “regular”, “homogeneous”, “smooth”, “simple”. E.g., attractive potentials in computer vision, ferromagnetic Potts models statistical physics.

On log-supermodular vs. log-submodular distributions

- Log-supermodular distributions.

$$\log \Pr(x) = g(x) + \text{const.} = -E(x) + \text{const.} \quad (1.50)$$

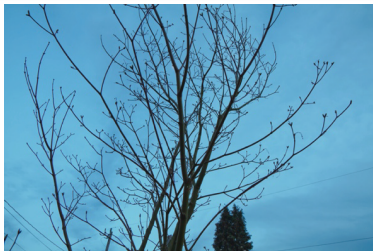
where g is supermodular ($E(x) = -g(x)$ is submodular). MAP (or high-probable) assignments should be “regular”, “homogeneous”, “smooth”, “simple”. E.g., attractive potentials in computer vision, ferromagnetic Potts models statistical physics.

- Log-submodular distributions:

$$\log \Pr(x) = f(x) + \text{const.} \quad (1.51)$$

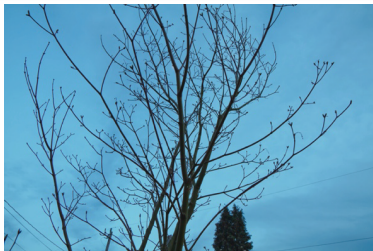
where f is submodular. MAP or high-probable assignments should be “diverse”, or “complex”, or “covering”, like in determinantal point processes.

Shrinking bias in graph cut image segmentation



What does graph-cut based image segmentation do with elongated structures (top) or contrast gradients (bottom)?

Shrinking bias in graph cut image segmentation



Addressing shrinking bias with edge submodularity

- Standard graph cut, uses a **modular** function $w : 2^E \rightarrow \mathbb{R}_+$ defined on the edges to measure cut costs. Graph cut node function is submodular.

$$f_w(X) = w\left(\{(u, v) \in E : u \in X, v \in V \setminus X\}\right) \quad (1.52)$$

Addressing shrinking bias with edge submodularity

- Standard graph cut, uses a **modular** function $w : 2^E \rightarrow \mathbb{R}_+$ defined on the edges to measure cut costs. Graph cut node function is submodular.

$$f_w(X) = w\left(\{(u, v) \in E : u \in X, v \in V \setminus X\}\right) \quad (1.52)$$

- Instead, we can use a submodular function $g : 2^E \rightarrow \mathbb{R}_+$ **defined on the edges** to express cooperative costs.

$$f_g(X) = g\left(\{(u, v) \in E : u \in X, v \in V \setminus X\}\right) \quad (1.53)$$

Addressing shrinking bias with edge submodularity

- Standard graph cut, uses a **modular** function $w : 2^E \rightarrow \mathbb{R}_+$ defined on the edges to measure cut costs. Graph cut node function is submodular.

$$f_w(X) = w\left(\{(u, v) \in E : u \in X, v \in V \setminus X\}\right) \quad (1.52)$$

- Instead, we can use a submodular function $g : 2^E \rightarrow \mathbb{R}_+$ **defined on the edges** to express cooperative costs.

$$f_g(X) = g\left(\{(u, v) \in E : u \in X, v \in V \setminus X\}\right) \quad (1.53)$$

- Seen as a node function, $f_g : 2^V \rightarrow \mathbb{R}_+$ is not submodular, but it uses submodularity internally to solve the shrinking bias problem.

Addressing shrinking bias with edge submodularity

- Standard graph cut, uses a **modular** function $w : 2^E \rightarrow \mathbb{R}_+$ defined on the edges to measure cut costs. Graph cut node function is submodular.

$$f_w(X) = w\left(\{(u, v) \in E : u \in X, v \in V \setminus X\}\right) \quad (1.52)$$

- Instead, we can use a submodular function $g : 2^E \rightarrow \mathbb{R}_+$ **defined on the edges** to express cooperative costs.

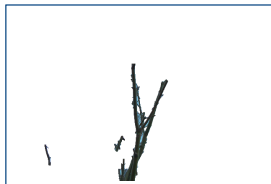
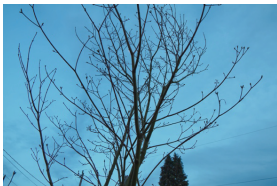
$$f_g(X) = g\left(\{(u, v) \in E : u \in X, v \in V \setminus X\}\right) \quad (1.53)$$

- Seen as a node function, $f_g : 2^V \rightarrow \mathbb{R}_+$ is not submodular, but it uses submodularity internally to solve the shrinking bias problem.
- \Rightarrow cooperative-cut (Jegelka & B., 2011).

Graph-cut vs. cooperative-cut comparisons

Graph Cut

Cooperative Cut



(Jegelka&Bilmes,'11). There are fast algorithms for solving as well.