





Pareto Optimization for Subset Selection: Theories and Practical Algorithms

Chao Qian and Yang Yu

LAMDA Group, Nanjing University, China

Email: {qianc, yuy}@lamda.nju.edu.cn



Introduction

Pareto optimization for subset selection

□ Pareto optimization for large-scale subset selection

□ Pareto optimization for noisy subset selection

□ Pareto optimization for dynamic subset selection

Conclusion

Subset selection is to select a subset of size at most *B* from a total set of *n* items for optimizing some objective function



Subset selection has diverse applications, which have different meanings on the item v_i and the objective f

Application - sensor placement

Sensor placement [Krause & Guestrin, IJCAI'09 Tutorial] : select a few places to install sensors such that the information gathered is maximized



Water contamination detection

Fire detection

Item v_i : a place to install a sensor

Objective *f* : entropy

http://www.lamda.nju.edu.cn/qianc/

Application - document summarization

Document summarization [Lin & Bilmes, ACL'11]: select a few sentences to best summarize the documents



http://www.lamda.nju.edu.cn/qianc/

Application - influence maximization

Influence maximization [Kempe et al., KDD'03]: select a subset of users from a social network to maximize its influence spread



Item v_i : a social network user

Objective *f*: influence spread, measured by the expected number of social network users activated by diffusion

Application - sparse regression

Sparse regression [Tropp, TIT'04]: select a few observation variables to best approximate the predictor variable by linear regression



Objective *f*: squared multiple correlation $R_{z,X}^2 =$

http://www.lamda.nju.edu.cn/yuy/

 $\frac{\operatorname{Var}(z) - \operatorname{MSE}_{z,X}}{z}$

Var(z

http://www.lamda.nju.edu.cn/qianc/

Application - maximum coverage

Maximum coverage [Feige, JACM'98]: select at most *B* sets from *n* given sets $V = \{S_1, ..., S_n\}$ to make the size of their union maximal $max_{X \subseteq V}$ $f(X) = ||\bigcup_{S_i \in X} S_i|| s.t. |X| \le B$

Example: $\forall i \leq l, S_i$ contains the same two elements, $\forall i > l, S_i$ contains one unique element; n = 2l, B = 2



Item v_i : a set of elements

Objective *f* : size of the union

Subset selection is to select a subset of size at most *B* from a total set of *n* items for optimizing some objective function

Formally stated: given all items $V = \{v_1, ..., v_n\}$, an objective function $f: 2^V \to \mathbb{R}$ and a budget *B*, to find a subset $X \subseteq V$ such that $max_{X\subseteq V}$ f(X) s.t. $|X| \leq B$ Application v_i sensor placement a place to install a sensor entropy document sum summary quality Many applications, but influence may influence spread **NP-hard in general!** squared multiple sparse regression an observation variable correlation a set of elements size of the union maximum coverage

http://www.lamda.nju.edu.cn/qianc/

Subset selection is to select a subset of size at most *B* from a total set of *n* items for optimizing some objective function

Formally stated: given all items $V = \{v_1, ..., v_n\}$, an objective function $f: 2^V \rightarrow \mathbb{R}$ and a budget B, to find a subset $X \subseteq V$ such that $max_{X \subseteq V} \quad f(X) \quad s.t. \quad |X| \leq B$



http://www.lamda.nju.edu.cn/qianc/

Previous approaches

Greedy algorithm

Process: iteratively select one item maximizing the increment on *f* $v^* = \arg \max_{v \in V \setminus X^{j-1}} f(X^{j-1} \cup \{v\}) - f(X^{j-1})$ v^* $X^1 = \{v^*\}$ Iteration 1: $\{v_1, v_2, \dots, v_n\}$ Run *B* iterations $X^{j} = X^{j-1} \cup \{v^*\}$ v^* Iteration *j*: $V \setminus X^{j-1}$

X^j: the subset obtained after *j* iterations

Approximation guarantees

Subset selection: given all items $V = \{v_1, ..., v_n\}$, an objective function $f: 2^V \rightarrow \mathbb{R}$ and a budget B, to find a subset $X \subseteq V$ such that $max_{X \subseteq V} \quad f(X) \quad s.t. \quad |X| \leq B$

f: monotone and submodular

The approximation guarantee [Nemhauser et al., MP'78]:

 $1 - 1/e \approx 0.632$ by the greedy algorithm

The subset *X* output by the greedy algorithm satisfies

$$f(X) \ge \left(1 - \frac{1}{e}\right) \cdot \text{OPT}_{\text{constrained}}$$
 the optimal function value

A set function $f: 2^V \to \mathbb{R}$ requires a solution to be a subset of V

Monotone: the function value increases as a set extends, i.e.,

$$\forall X \subseteq Y \subseteq V \colon f(X) \le f(Y)$$

Submodular [Nemhauser et al., MP'78]: satisfy the natural diminishing returns property, i.e.,

 $\forall X \subseteq Y \subseteq V, v \notin Y: f(X \cup \{v\}) - f(X) \ge f(Y \cup \{v\}) - f(Y);$ or equivalently,

 $\forall X \subseteq Y \subseteq V: f(Y) - f(X) \le \sum_{v \in Y \setminus X} f(X \cup \{v\}) - f(X);$ or equivalently,

 $\forall X, Y \subseteq V: f(X) + f(Y) \ge f(X \cap Y) + f(X \cup Y)$

Submodular applications

Maximum coverage [Feige, JACM'98] : select at most *B* sets from *n* given sets $V = \{S_1, \dots, S_n\}$ to make the size of their union maximal $max_{X\subseteq V} \quad f(X) = |\bigcup_{S_i \in X} S_i| \quad s.t. \quad |X| \le B$ Monotone: $\forall X \subseteq Y \subseteq V$: $f(X) \leq f(Y)$ f(Y) - f(X) $Y = \{S_1, S_2\}$ S_1 S_2 $X = \{S_1\}$ S_1 Submodular: $\forall X \subseteq Y \subseteq V, v \notin Y$: $f(X \cup \{v\}) - f(X) \ge f(Y \cup \{v\}) - f(Y)$ 53 $X = \{S_1\}$ S_2 S_1 S_1 $Y = \{S_1, S_2\}$ $v = S_3$

http://www.lamda.nju.edu.cn/qianc/

Submodular applications

Maximum coverage [Feige, JACM'98] : select at most *B* sets from *n* given sets $V = \{S_1, ..., S_n\}$ to make the size of their union maximal

$$\max_{X \subseteq V} f(X) = |\bigcup_{S_i \in X} S_i| \quad s.t. \quad |X| \le B$$

More applications:

- Sensor placement
- Document summarization





Their objective functions are all monotone and submodular



http://www.lamda.nju.edu.cn/qianc/

Approximation guarantees

Subset selection: given all items $V = \{v_1, ..., v_n\}$, an objective function $f: 2^V \rightarrow \mathbb{R}$ and a budget B, to find a subset $X \subseteq V$ such that $max_{X \subseteq V} \quad f(X) \quad s.t. \quad |X| \leq B$

f: monotone and submodular

The approximation guarantee [Nemhauser et al., MP'78]: $1 - 1/e \approx 0.632$ by the greedy algorithm

f: monotone

The approximation guarantee [Das & Kempe, ICML'11]:

 $1 - 1/e^{\gamma}$ by the greedy algorithm

Submodular ratio γ : to what extent *f* satisfies the submodular property

Submodular ratio

Submodular [Nemhauser et al., MP'78]:

$$\forall X \subseteq Y \subseteq V, v \notin Y: f(X \cup \{v\}) - f(X) \ge f(Y \cup \{v\}) - f(Y);$$

or $\forall X \subseteq Y \subseteq V$: $f(Y) - f(X) \le \sum_{v \in Y \setminus X} f(X \cup \{v\}) - f(X)$ -----

Submodular ratio [Das & Kempe, ICML'11; Zhang & Vorobeychi, AAAI'16]:

$$\Rightarrow \alpha_f = \min_{\substack{X \subseteq Y, \nu \notin Y}} \frac{f(X \cup \{\nu\}) - f(X)}{f(Y \cup \{\nu\}) - f(Y)}$$

$$\gamma_{U,k}(f) = \min_{\substack{X \subseteq U, Y: |Y| \le k, X \cap Y = \emptyset}} \frac{\sum_{\nu \in Y} f(X \cup \{\nu\}) - f(X)}{f(X \cup Y) - f(X)}$$

Characterize to what extent a set function *f* satisfies the submodular property For example, when *f* is monotone,

- $\forall U, k: \gamma_{U,k}(f) \in [0,1]$, the larger, more close to submodular
- *f* is submodular if and only if $\forall U, k: \gamma_{U,k}(f) = 1$

Submodular ratio [Das & Kempe, ICML'11; Zhang & Vorobeychi, AAAI'16] : characterize to what extent a general set function satisfies the submodular property

$$\alpha_{f} = \min_{X \subseteq Y, v \notin Y} \frac{f(X \cup \{v\}) - f(X)}{f(Y \cup \{v\}) - f(Y)}$$
$$\gamma_{U,k}(f) = \min_{X \subseteq U, \ Y: |Y| \le k, X \cap Y = \emptyset} \frac{\sum_{v \in Y} f(X \cup \{v\}) - f(X)}{f(X \cup Y) - f(X)}$$

Lower bounds on submodular ratio for some non-submodular applications

- Sparse regression: $\gamma_{U,k}(f) \ge \lambda_{min}(C, |U| + k)$ [Das & Kempe, ICML'11]
- Sparse support selection: $\gamma_{U,k}(f) \ge m/M$ [Elenberg et al., Annals of Statistics'18]
- Bayesian experimental design [Bian et al., ICML'17]:

$$\gamma_{U,k}(f) \geq \beta^2 / \left(\|\mathbf{V}\|^2 (\beta^2 + \sigma^{-2} \|\mathbf{V}\|^2) \right)$$

• Determinantal function maximization [Qian et al., IJCAI'18]: $\alpha_f \ge (\lambda_n(A) - 1) / \left((\lambda_1(A) - 1) \prod_{i=1}^{n-1} \lambda_i(A) \right)$

Approximation guarantees



Good optimization performance:

- Good approximation guarantee, i.e., good performance in worst cases
- Practical performance is much better (e.g., close to optima) in most cases
 The greedy nature

Previous approaches (con't)

• Relaxation method

Process: relax the original problem, and then find the optimal solutions to the relaxed problem

Weakness: the optimal solution of the relaxed problem may be distant to the true optimum

$$\max_{X \subseteq V} f(X) \quad s.t. \quad |X| \leq B$$

$$\bigoplus_{max_{w \in R^{n}}} g(w) \quad s.t. \quad |w|_{0} \leq B \quad \text{non-convex}$$

$$\max_{w \in R^{n}} g(w) \quad s.t. \quad |w|_{1} \leq B \quad \text{convex}$$

Variants of subset selection

Subset selection $1 - 1/e^{\gamma}$ $max_{X \subset V}$ f(X) s.t. $|X| \leq B$ [Das & Kempe, ICML'11] General constraints $(\alpha/2)(1-1/e^{\alpha})$ $|X| \leq B \rightarrow c(X) \leq B$ [Zhang & Vorobeychik, AAAI'16] Multiset selection $1 - 1/e^{\beta}$ $(\alpha/2)(1-1/e^{\alpha})$ *X*: a subset \rightarrow a multiset [Alon et al., WWW'12] [Soma et al., ICML'14] k-subsets selection 1/2X: a subset $\rightarrow k$ subsets [Ohsaka & Yoshida, NIPS'15] Sequence selection $1 - \rho^{-1/(2\Delta)}$ *X*: a subset \rightarrow a sequence [Tschiatschek et al, AAAI'17] Ratio optimization $|X^*|$ $(1 + (|X^*| - 1)(1 - \kappa))\gamma$ $min_{X \subseteq V} \quad f(X)/g(X)$ [Bai et al., ICML'16]

http://www.lamda.nju.edu.cn/qianc/

Subset selection: $max_{X\subseteq V}$ f(X) s.t. $|X| \le B$ Two conflicting objectives:

- 1. Optimize the objective $f = \max_{X \subseteq V} f(X)$
- 2. Keep the size small $min_{X\subseteq V} max\{|X| B, 0\}$

Previous theoretical studies have disclosed the advantage of solving single-objective constrained optimization by MOEAs

[Neumann & Wegener, NC'06; Friedrich et al., ECJ'10; Neumann et al., Algorithmica'11; Yu et al., AIJ'12]

Why not optimize the bi-objective formulation?

 $min_{X\subseteq V}$ (-f(X), |X|)



Introduction

Pareto optimization for subset selection

□ Pareto optimization for large-scale subset selection

□ Pareto optimization for noisy subset selection

□ Pareto optimization for dynamic subset selection

Conclusion

A subset $X \subseteq V$ can be naturally represented by a Boolean vector $x \in \{0,1\}^n$

• the *i*-th bit $x_i = 1$ if the item $v_i \in X$; $x_i = 0$ otherwise

•
$$X = \{v_i \mid x_i = 1\}$$

 $V = \{v_1, v_2, v_3, v_4, v_5\}$ a subset $X \subseteq V$ a Boolean vector $x \in \{0,1\}^5$
 \emptyset 00000

 $\{v_1\}$ \longleftrightarrow 10000

 $\{v_2, v_3, v_5\}$ 01101

 $\{v_1, v_2, v_3, v_4, v_5\}$ 11111

Pareto optimization



Pareto optimization



Subset selection with monotone submodular *f*

[Friedrich & Neumann, ECJ'15]

Exclude solutions with size larger than *B*

A simple multi-objective evolutionary algorithm GSEMO [Laumanns et al., TEvC'04]



initialization

updating

Initialization: put a random solution from $\{0,1\}^n$ into the population *P*

Reproduction: pick a solution x randomly from *P*, and flip each bit of $x \in \{0,1\}^n$ with prob. 1/n to generate a new solution

Updating: if the new solution is not dominated by any solution in *P*, put it into *P* and weed out bad solutions

Output: select the best solution with size at most *B*

It can achieve the optimal approximation guarantee of (1 - 1/e)in $O(n^2(B + \log n))$ expected running time

Subset selection with monotone f

The POSS algorithm [Qian, Yu and Zhou, NIPS'15]

$max_{x \in \{0,1\}^n} f(\mathbf{x})$ s.t. $|\mathbf{x}| \le B$ originalTransformation: \mathbf{J} $min_{x \in \{0,1\}^n} (-f(\mathbf{x}), |\mathbf{x}|)$ bi-objective

Algorithm 1 POSS

Input: all variables $V = \{X_1, \dots, X_n\}$, a given objective f and an integer parameter $k \in [1, n]$ **Parameter**: the number of iterations T **Output**: a subset of V with at most k variables Process: 1: Let $s = \{0\}^n$ and $P = \{s\}$. 2: Let t = 0. 3: while t < T do Select *s* from *P* uniformly at random. 4: 5: Generate s' by flipping each bit of s with prob. $\frac{1}{n}$. Evaluate $f_1(s')$ and $f_2(s')$. 6: 7: if $\exists z \in P$ such that $z \prec s'$ then $Q = \{ z \in P \mid s' \preceq z \}.$ 8: $\dot{P} = (P \setminus Q) \cup \{s'\}.$ 9: end if 10:t = t + 1. 11: 12: end while 13: return $\arg\min_{s \in P, |s| \le k} f_1(s)$

Initialization: put the special solution {0}^{*n*} into the population *P*

Exclude solutions with size at least 2B

Reproduction: pick a solution x randomly from P, and flip each bit of x with prob. 1/n to produce a new solution

Updating: if the new solution is not dominated by any solution in *P*, put it into *P* and weed out bad solutions

Output: select the best feasible solution

http://www.lamda.nju.edu.cn/qianc/

Subset selection with monotone f

The POSS algorithm [Qian, Yu and Zhou, NIPS'15]

$$max_{x \in \{0,1\}^n} f(x)$$
 s.t. $|x| \le B$ originalTransformation: $\[min_{x \in \{0,1\}^n} (-f(x), |x|)\]$ bi-objective

Initialization: put the special solution $\{0\}^n$ into the population P

Reproduction: pick a solution x randomly from P, and flip each bit of x with prob. 1/n to produce a new solution

Updating: if the new solution is not dominated by any solution in *P*, put it into *P* and weed out bad solutions

Output: select the best feasible solution

• Selection: each solution in the population *P* is selected with probability 1/|*P*|

e.g., if *P* contains 10 solutions, each solution is selected with probability 1/10

Exclude solutions with size at least 2B

• Bit-wise mutation: Pr(flip *i* specific bits)= $(1/n)^i(1 - 1/n)^{n-i}$ e.g. the probability of flipping a specific bi

e.g., the probability of flipping a specific bit of a solution is $(1/n)(1 - 1/n)^{n-1}$

Subset selection with monotone f

The POSS algorithm [Qian, Yu and Zhou, NIPS'15]

$$max_{x \in \{0,1\}^n} f(x)$$
 s.t. $|x| \le B$ originalTransformation: $\[min_{x \in \{0,1\}^n} (-f(x), |x|)\]$ bi-objective

Initialization: put the special solution $\{0\}^n$ into the population P

Reproduction: pick a solution x randomly from P, and flip each bit of x with prob. 1/n to produce a new solution

Updating: if the new solution is not dominated by any solution in *P*, put it into *P* and weed out bad solutions

Output: select the best feasible solution

Selection: each solution in the population *P* is selected with probability 1/|*P*|

e.g., if *P* contains 10 solutions, each solution is selected with probability 1/10

Exclude solutions with size at least 2B

• Bit-wise mutation: Pr(flip *i* specific bits)= $(1/n)^i(1 - 1/n)^{n-i}$

e.g., the probability of flipping a specific bit of a solution is $(1/n)(1 - 1/n)^{n-1}$

• The population *P* always contains nondominated solutions generated so-far POSS can achieve the optimal approximation guarantee, previously obtained by the greedy algorithm

Theorem 1. For subset selection with monotone objective function f, POSS using $E[T] \le 2eB^2n$ finds a solution x with $|x| \le B$ and $f(x) \ge (1 - e^{-\gamma}) \cdot OPT$. the expected number of iterations

> the optimal polynomial-time approximation ratio, previously obtained by the greedy algorithm [Das & Kempe, ICML'11]

Lemma 1. For any $X \subseteq V$, there exists one item $\hat{v} \in V \setminus X$ such that

$$f(X \cup \{\hat{v}\}) - f(X) \ge \frac{\gamma}{B} (\text{OPT} - f(X))$$

submodularity ratio [Das & Kempe, ICML'11]

the optimal function value

Roughly speaking, the improvement by adding a specific item is proportional to the current distance to the optimum

Lemma 1. For any
$$X \subseteq V$$
, there exists one item $\hat{v} \in V \setminus X$ such that $f(X \cup {\hat{v}}) - f(X) \ge \frac{\gamma}{B} (\text{OPT} - f(X))$
Main idea:

• consider a solution \mathbf{x} with $|\mathbf{x}| \le i$ and $f(\mathbf{x}) \ge \left(1 - \left(1 - \frac{\gamma}{B}\right)^i\right) \cdot \text{OPT}$

$$i = 0$$

$$i = B$$

$$\downarrow$$
initial solution 00 ... 0
$$1 - \left(1 - \frac{\gamma}{B}\right)^{B} = 1 - \left(1 - \frac{1}{B/\gamma}\right)^{(B/\gamma) \cdot \gamma}$$

$$|00 ... 0| = 0$$

$$let m = \frac{B/\gamma}{A} \ge 1 - e^{-\gamma}$$

$$(1 - 1/m)^{m} \le 1/e$$

http://www.lamda.nju.edu.cn/qianc/



Lemma 1. For any $X \subseteq V$, there exists one item $\hat{v} \in V \setminus X$ such that $f(X \cup {\hat{v}}) - f(X) \ge \frac{\gamma}{B}(\text{OPT} - f(X))$

Main idea:

- consider a solution \mathbf{x} with $|\mathbf{x}| \leq i$ and $f(\mathbf{x}) \geq \left(1 \left(1 \frac{\gamma}{B}\right)^i\right) \cdot \text{OPT}$
- in each iteration of POSS:
 - > select x from the population P

a subset

flip one specific 0-bit of x to 1-bit
 (i.e., add the specific item v in Lemma 1)

$$\Rightarrow |\mathbf{x}'| = |\mathbf{x}| + 1 \leq i + 1 \text{ and } f(\mathbf{x}') \geq \left(1 - \left(1 - \frac{\gamma}{B}\right)^{i+1}\right) \cdot \text{OPT}$$

Lemma 1. For any
$$X \subseteq V$$
, there exists one item $\hat{v} \in V \setminus X$ such that
 $f(X \cup \{\hat{v}\}) - f(X) \ge \frac{\gamma}{B} (\text{OPT} - f(X))$
 $f(x') - f(x) \ge \frac{\gamma}{B} \cdot (\text{OPT} - f(x))$
 $f(x') \ge (1 - \frac{\gamma}{B})f(x) + \frac{\gamma}{B} \cdot \text{OPT}$
 $f(x) \ge (1 - (1 - \frac{\gamma}{B})^i) \cdot \text{OPT}$
 $f(x') \ge (1 - (1 - \frac{\gamma}{B})^i) \cdot \text{OPT} = (1 - (1 - \frac{\gamma}{B})^{i+1}) \cdot \text{OPT}$
Proof

Lemma 1. For any $(X) \subseteq V$, there exists one item $\hat{v} \in V \setminus X$ such that $f(X \cup \{\hat{v}\}) - f(X) \ge \frac{\gamma}{R}(\text{OPT} - f(X))$

Main idea:

• consider a solution **x** with
$$|\mathbf{x}| \le i$$
 and $f(\mathbf{x}) \ge \left(1 - \left(1 - \frac{\gamma}{B}\right)^i\right) \cdot \text{OPT}$

in each iteration of POSS:

a subset

 select *x* from the population *P*, the probability: ¹/_{|P|}
 flip one specific 0-bit of *x* to 1-bit, the probability: ¹/_n (1 - ¹/_n)ⁿ⁻¹ ≥ ¹/_{en} (i.e., add the specific item \hat{v} in Lemma 1)

$$\Rightarrow |\mathbf{x}'| = |\mathbf{x}| + 1 \le i + 1 \text{ and } f(\mathbf{x}') \ge \left(1 - \left(1 - \frac{\gamma}{B}\right)^{i+1}\right) \cdot \text{OPT}$$

$$i \longrightarrow i + 1 \quad \text{the probability: } \frac{1}{|P|} \cdot \frac{1}{en}$$

http://www.lamda.nju.edu.cn/gianc/

http://www.lamda.nju.edu.cn/yuy/

en

Proof

Lemma 1. For any
$$X \subseteq V$$
, there exists one item $\hat{v} \in V \setminus X$ such that
 $f(X \cup \{\hat{v}\}) - f(X) \ge \frac{\gamma}{B} (OPT - f(X))$
Main idea:
a subset
o consider a solution X with $|\mathbf{x}| \le i$ and $f(\mathbf{x}) \ge \left(1 - \left(1 - \frac{\gamma}{B}\right)^i\right) \cdot OPT$
in each iteration of POSS:
 $i \longrightarrow i + 1$ the probability: $\frac{1}{|P|} \cdot \frac{1}{en} \xrightarrow{|P| \le 2B} \frac{1}{2eBn}$
For each size in
 $\{0,1, ..., 2B - 1\},$
there exists at most
one solution in P

http://www.lamda.nju.edu.cn/qianc/

Proof

Lemma 1. For any $X \subseteq V$, there exists one item $\hat{v} \in V \setminus X$ such that $f(X \cup {\hat{v}}) - f(X) \ge \frac{\gamma}{B}(\text{OPT} - f(X))$

Main idea:

• consider a solution **x** with
$$|\mathbf{x}| \le i$$
 and $f(\mathbf{x}) \ge \left(1 - \left(1 - \frac{\gamma}{B}\right)^i\right) \cdot \text{OPT}$

• in each iteration of POSS:

a subset

$$i \longrightarrow i+1$$
 the probability: $\frac{1}{|P|} \cdot \frac{1}{en} \quad |P| \le 2B \quad \frac{1}{2eBn}$

 $i \longrightarrow i + 1$ the expected number of iterations: 2eBn

 $i = 0 \longrightarrow B$ the expected number of iterations: $B \cdot 2eBn$

POSS can achieve the optimal approximation guarantee, previously obtained by the greedy algorithm

Theorem 1. For subset selection with monotone objective function f, POSS using $E[T] \le 2eB^2n$ finds a solution x with $|x| \le B$ and $f(x) \ge (1 - e^{-\gamma}) \cdot OPT$.

the optimal polynomial-time approximation ratio, previously obtained by the greedy algorithm [Das & Kempe, ICML'11]

POSS can do better than the greedy algorithm in cases

Theorem 2. For the Exponential Decay subclass of sparse regression, POSS using $E[T] = O(B^2(n - B)n \log n)$ finds an optimal solution, while the greedy algorithm cannot.

http://www.lamda.nju.edu.cn/qianc/

[Das & Kempe, STOC'08]

Experiments on sparse regression

Sparse regression: given all observation variables $V = \{v_1, ..., v_n\}$, a predictor variable *z* and a budget *B*, to find a subset $X \subseteq V$ such that

$$max_{X\subseteq V} \quad R_{z,X}^2 = \frac{\operatorname{Var}(z) - \operatorname{MSE}_{z,X}}{\operatorname{Var}(z)} \quad s.t. \quad |X| \le B$$

Var(z): variance of z

 $\frac{\text{MSE}_{z,X}}{\text{by using observation variables in } X}$



http://www.lamda.nju.edu.cn/yuy/

http://www.lamda.nju.edu.cn/qianc/

Experimental results - R^2 values

the size constraint: B = 8

the number of iterations of POSS: $2eB^2n$

exhaustive search		greedy algorithms		relaxation methods				
F						*		
Data set	OPT	POSS	FR	FoBa	OMP	RFE	MCP	
housing	.7437±.0297	.7437±.0297	.7429±.0300•	.7423±.0301•	.7415±.0300•	.7388±.0304•	.7354±.0297•	
eunite2001	.8484±.0132	$.8482 \pm .0132$.8348±.0143•	.8442±.0144•	.8349±.0150●	.8424±.0153•	.8320±.0150•	
svmguide3	.2705±.0255	.2701±.0257	.2615±.0260•	.2601±.0279•	.2557±.0270●	.2136±.0325•	.2397±.0237•	
ionosphere	.5995±.0326	.5990±.0329	.5920±.0352•	.5929±.0346•	.5921±.0353•	.5832±.0415•	.5740±.0348•	
sonar	-	$.5365 \pm .0410$.5171±.0440●	.5138±.0432•	.5112±.0425•	.4321±.0636•	.4496±.0482•	
triazines	-	.4301±.0603	.4150±.0592•	.4107±.0600•	.4073±.0591•	.3615±.0712•	.3793±.0584•	
coil2000	-	$.0627 \pm .0076$.0624±.0076•	.0619±.0075•	.0619±.0075•	.0363±.0141•	.0570±.0075•	
mushrooms	_	.9912±.0020	.9909±.0021•	.9909±.0022•	.9909±.0022•	.6813±.1294•	.8652±.0474•	
clean1	-	$.4368 \pm .0300$.4169±.0299•	.4145±.0309•	.4132±.0315•	.1596±.0562•	.3563±.0364•	
w5a	-	.3376±.0267	.3319±.0247•	.3341±.0258•	.3313±.0246•	.3342±.0276•	.2694±.0385•	
gisette	-	$.7265 \pm .0098$.7001±.0116•	.6747±.0145•	.6731±.0134•	.5360±.0318•	.5709±.0123•	
farm-ads	-	$.4217 \pm .0100$.4196±.0101•	.4170±.0113•	.4170±.0113•	-	.3771±.0110•	
POSS: win/tie/loss		_	12/0/0	12/0/0	12/0/0	11/0/0	12/0/0	

• denotes that POSS is significantly better by the *t*-test with confidence level 0.05



POSS is significantly better than all the compared algorithms on all data sets

Experimental results - R^2 values

different size constraints: $B = 3 \rightarrow 8$



POSS tightly follows OPT, and has a clear advantage over the rest algorithms

Experimental results – running time

OPT: n^B/B^B greedy algorithms (FR): Bn POSS: $2eB^2n$



POSS can be much more efficient in practice

Pareto optimization vs. Greedy algorithm

Greedy algorithm:

• Generate a new solution by adding a single item

```
(i.e., single-bit forward search: 0 \rightarrow 1)
```

• Keep only one solution

Pareto optimization:

- Generate a new solution by flipping each bit with prob. 1/n
 - > single-bit forward search : 0 → 1
 - ▶ backward search : $1 \rightarrow 0$
 - ▶ multi-bit search : $00 \rightarrow 11$
- Keep a set of non-dominated solutions due to bi-objective optimization

Pareto optimization may have a better ability of escaping from local optima

Variants of subset selection

- Subset selection
- General constraints
- Multiset selection
- *k*-subsets selection
- Sequence selection

 $max_{x \in \{0,1\}^n} f(x) \quad s.t. \ |x| \le B$ $\max_{\mathbf{x}\in\{0,1\}^n} f(\mathbf{x}) \qquad \text{s.t. } c(\mathbf{x}) \le B$ $max_{\boldsymbol{x}\in\mathbf{Z}_{+}^{n}} f(\boldsymbol{x}) \qquad s.t. \ |\boldsymbol{x}| \leq B$ x_i : the number of times that the item v_i appears $max_{x \in \{0,1,...,k\}^n} f(x) \text{ s.t. } |x| \le B$ x_i : the subset where the item v_i appears s.t. $|x| \leq B$ $max_{x \in \mathcal{S}} f(x)$ *x*: a sequence where the order of items influences *f*

 $\min_{\mathbf{x}\in\{0,1\}^n} f(\mathbf{x})/g(\mathbf{x})$

Ratio optimization

Variants of subset selection

- Subset selection $\max_{x \in \{0,1\}^n} f(x)$ s.t. $|x| \le B$ [Friedrich & Neumann, ECJ'15; Qian et al., NIPS'15]
- General constraints $max_{x \in \{0,1\}^n} f(x)$ s.t. $c(x) \le B$ [Qian et al., IJCAI'17a]
- Multiset selection $\max_{x \in \mathbb{Z}^n_+} f(x)$ $s.t. |x| \le B$ [Qian et al., AAAI'18]
- k-subsets selection $\max_{x \in \{0,1,\dots,k\}^n} f(x)$ s.t. $|x| \le B$ [Qian et al., TEvC'18]
- Sequence selection $\max_{x \in S} f(x)$ s.t. $|x| \le B$ [Qian et al., IJCAI'18]
- Ratio optimization $\min_{x \in \{0,1\}^n} f(x)/g(x)$ [Qian et al., IJCAI'17b]

Pareto optimization can achieve the best-known polynomial-time approximation guarantee, and perform well in practice

Ratio optimization

The PORM algorithm [Qian, Shi, Yu, Tang and Zhou, IJCAI'17]

Transformation:

 $\min_{\mathbf{x}\in\{0,1\}^n} f(\mathbf{x})/g(\mathbf{x})$ original $min_{x \in \{0,1\}^n} (f(x), -g(x))$

bi-objective

Algorithm 2 PORM algorithm

Input: a monotone submodular function $f : \{0, 1\}^n \to \mathbb{R}^+$ and a monotone function $g: \{0,1\}^n \to \mathbb{R}^+$ **Parameter**: the number T of iterations **Output**: a solution $x \in \{0, 1\}^n$ Process: 1: Select x from $\{0,1\}^n$ uniformly at random. 2: Let $P = \{x\}$ and t = 0. 3: while t < T do Select x from P uniformly at random. 4: Generate x' by flipping each bit of x with prob. 1/n. 5: if $\exists z \in P$ such that $z \prec x'$ then 6: $P = (P \setminus \{ z \in P \mid x' \preceq z \}) \cup \{ x' \}.$ 7: $Q = \{ z \in P \mid |z| = |x'| \}.$ 8: $z_1 = \operatorname{arg\,min}_{z \in Q} f(z), \ z_2 = \operatorname{arg\,max}_{z \in Q} g(z),$ 9: $z_3 = \arg\min_{z \in Q} f(z)/g(z).$ $P = (P \setminus Q) \cup \{z_1, z_2, z_3\}.$ 10:end if t = t + 1.12: 13: end while 14: return $\arg\min_{x\in P} f(x)/g(x)$

Initialization: put a random solution from $\{0,1\}^n$ into the population *P*

Reproduction: pick a solution *x* randomly from P, and flip each bit of x with prob. 1/n to produce a new solution

Updating: if the new solution is not dominated by any solution in *P*, put it into *P* and weed out bad solutions; keep at most three solutions for each subset size (smallest f, largest g, smallest ratio f/g)

Output: the solution with the smallest ratio

http://www.lamda.nju.edu.cn/gianc/

Ratio optimization

The PORM algorithm [Qian, Shi, Yu, Tang and Zhou, IJCAI'17] $min_{x \in \{0,1\}^n} f(x)/g(x)$ originalTransformation: \bildlefty $min_{x \in \{0,1\}^n} (f(x), -g(x))$ bi-objective

Theory: PORM achieves the best-known approximation guarantee $\frac{|X^*|}{(1+(|X^*|-1)(1-\kappa))\gamma'}$, previously obtained by GreedRatio [Bai et al., ICML'16]

Application:

F-measure maximization in information retrieval

Always better



Pareto optimization for subset selection

achieve excellent performance on diverse variants of subset selection both theoretically and empirically

The running time (e.g., $2eB^2n$) for achieving a good solution unsatisfactory when the problem size (e.g., *B* and *n*) is large can be reduced to linear time but with performance loss [Crawford & Kuhnle, 2019]

A sequential algorithm that cannot be readily parallelized

How can Pareto optimization be applied to large-scale subset selection problems?



Introduction

Pareto optimization for subset selection

Pareto optimization for large-scale subset selection

Pareto optimization for noisy subset selection

Pareto optimization for dynamic subset selection

Conclusion

Pareto optimization for subset selection



http://www.lamda.nju.edu.cn/qianc/

Parallel Pareto optimization for subset selection



http://www.lamda.nju.edu.cn/qianc/

Parallel Pareto optimization for subset selection



Q: the same solution quality?

Yes!

http://www.lamda.nju.edu.cn/qianc/

Parallel Pareto optimization for subset selection



http://www.lamda.nju.edu.cn/qianc/

Theorem 3. For subset selection with monotone objective function f, the expected number of iterations until PPOSS finds a solution x with $|x| \le B$ and $f(x) \ge (1 - e^{-\gamma}) \cdot \text{OPT}$ is (1) if N = o(n), then $\mathbb{E}[T] \le 2eB^2n/N$: (2) if $N = \Omega(n^i)$ for $1 \le i \le B$, then $\mathbb{E}[T] = O(B^2/i)$; Keep the optimal approximation guarantee (3) if $N = \Omega(n^{\min\{3B-1,n\}})$, then $\mathbb{E}[T] = O(1)$.

• When <u>the number *N* of cores</u> is asymptotically less than <u>the number *n*</u> <u>of items</u>, <u>the expected number E[*T*] of iterations can be reduced linearly w.r.t. the number of cores</u>

Theorem 3. For subset selection with monotone objective function f, the expected number of iterations until PPOSS finds a solution x with $|x| \le B$ and $f(x) \ge (1 - e^{-\gamma}) \cdot \text{OPT}$ is (1) if N = o(n), then $\mathbb{E}[T] \le QeB^2n/N$: (2) if $N = \Omega(n^i)$ for $1 \le i \le B$, then $\mathbb{E}[T] = O(B^2/i)$; Keep the optimal approximation guarantee (3) if $N = \Omega(n^{\min\{3B-1,n\}})$, then $\mathbb{E}[T] = O(1)$.

- When <u>the number *N* of cores</u> is asymptotically less than <u>the number *n*</u> of items, the expected number E[*T*] of iterations can be reduced linearly w.r.t. the number of cores
- With increasing number *N* of cores, the expected number E[*T*] of iterations can be continuously reduced, eventually to a constant

Experiments on sparse regression

Compare the speedup as well as the solution quality measured by R^2 values with different number of cores



http://www.lamda.nju.edu.cn/qianc/

Experiments on sparse regression



(f) *farm-ads* (4143 #inst, 54877 #feat)

PPOSS (blue line): achieve speedup around 8 when the number of cores is 10; the R^2 values are stable, and better than the greedy algorithm **PPOSS-If (red line):** achieve better speedup as expected; the R^2 values are slightly worse

Pareto optimization for subset selection

achieve excellent performance on diverse variants of subset selection both theoretically and empirically

Parallel Pareto optimization for subset selection achieve nearly linear runtime speedup while keeping the solution quality

Require centralized access to the whole data set restrict the application to large-scale real-world problems

Can we make Pareto optimization distributable?

(Parallel) Pareto optimization for subset selection



Require centralized access to the whole data set at each machine

Large-scale data is too large to be stored at one single machine

Distributed Pareto optimization for subset selection



http://www.lamda.nju.edu.cn/qianc/

Experiments on sparse regression

Compare DPOSS with the state-of-the-art distributed greedy algorithm RandGreeDi [Mirzasoleiman et al., JMLR'16] under different number of machines

0.03 DPOSS 0.92 0.86 RANDGREED 0.025 0 0.84 0.88 0.02 0.82 0.86 0.68 0.84 0.015 0.8 DPOSS DPOSS 0.66 0.82 RANDGREED 0.78 0.01 2 10 6 8 10 4 10 2 m m mm MicroMass (n=1, 300)(c) SVHN (n=3,072)(b) colon-cancer (n=2,000)gisette (n = 5, 000)(d) (a) 0.995 0.64 - DPOSS 0.97 -× RANDGREEDI 0.545 0.99 0.62 0.96 0.54 -0.985 0.95 0.535 0.98 0.94 0.53 DPOSS - DPOSS - DPOSS 0.58 RANDGREED × RANDGREED × RANDGREED 0.93 0.975 0.525 10 6 8 10 2 10 6 10 m m m m GHG-Network (n=5, 232)leukemia (n=7, 129)Arcene (n = 10, 000)Dexter (n=20,000)(g) (f)(h) (e)

On regular-scale data sets

DPOSS is always better than RandGreeDi

http://www.lamda.nju.edu.cn/qianc/

Experiments on sparse regression

On regular-scale data sets 0.992 0.994 Approximation ratio 66°0 and 6 0.99 0.992 0.992 0.99 0.986 0.986 the solution quality by DPOSS ratio =the solution quality by POSS **DPOSS** is very close to the centralized POSS 0.95 2 3 4 5 6 7 8 9 10

On large-scale data sets

DPOSS is better than RandGreeDi

Data set	DPOSS	RANDGREEDI
<i>Gas-sensor-flow</i> $(n = 120, 432)$	$.818 {\pm} .005$.710±.017
Twin-gas-sensor $(n = 480, 000)$	$.601 \pm .014$	$.470 \pm .025$
Gas-sensor-sample $(n = 1, 950, 000)$	$.289 {\pm} .029$	$.245 {\pm} .018$

http://www.lamda.nju.edu.cn/qianc/

http://www.lamda.nju.edu.cn/yuy/

т

Experiments on maximum coverage

On regular-scale data sets





On large-scale data sets

Data set	DPOSS	RANDGREEDI
<i>accident</i> $(n = 340, 183)$	175 ± 1	170.6 ± 1.34
kosarak (n = 990, 002)	9263±0	9263±0

DPOSS is very close to the centralized POSS, and is better than RandGreeDi

http://www.lamda.nju.edu.cn/qianc/

Pareto optimization for subset selection

achieve excellent performance on diverse variants of subset selection both theoretically and empirically

Parallel Pareto optimization for subset selection achieve nearly linear runtime speedup while keeping the solution quality

Distributed Pareto optimization for subset selection achieve very close performance to the centralized algorithm

large-scale subset selection

Previous analyses often assume that the exact value of the objective function can be accessed

However, in many applications of subset selection, only a noisy value of the objective function can be obtained





Previous analyses often assume that the exact value of the objective function can be accessed

However, in many applications of subset selection, only a noisy value of the objective function can be obtained



Previous analyses often assume that the exact value of the objective function can be accessed

However, in many applications of subset selection, only a noisy value of the objective function can be obtained



How about the performance for noisy subset selection?



Introduction

Pareto optimization for subset selection

□ Pareto optimization for large-scale subset selection

Pareto optimization for noisy subset selection

Pareto optimization for dynamic subset selection

Conclusion

Noisy subset selection

Subset selection: given $V = \{v_1, ..., v_n\}$, an objective function $f: 2^V \to \mathbb{R}$ and a budget B, to find a subset $X \subseteq V$ such that $max_{X \subseteq V} \quad f(X) \quad s.t. \quad |X| \leq B$ exact objective value Noise Multiplicative: $(1 - \epsilon) \cdot f(X) \leq F(X) \leq (1 + \epsilon) \cdot f(X)$ Additive: $f(X) - \epsilon \leq F(X) \leq f(X) + \epsilon$

Applications: influence maximization, sparse regression maximizing information gain in graphical models [Chen et al., COLT'15] crowdsourced image collection summarization [Singla et al., AAAI'16]
Greedy algorithm & POSS [Qian et al., NIPS'17]:

Multiplicative noise: ϵ : the noise strength

 $f(X) \ge \frac{1}{1 + \frac{2\epsilon B}{(1 - \epsilon)\gamma}} \left(1 - \left(\frac{1 - \epsilon}{1 + \epsilon}\right)^B \left(1 - \frac{\gamma}{B}\right)^B \right) \cdot \text{OPT}$

Additive noise:

$$f(X) \ge \left(1 - \left(1 - \frac{\gamma}{B}\right)^B\right) \cdot \text{OPT} - \left(\frac{2B}{\gamma} - \frac{2B}{\gamma}e^{-\gamma}\right)\epsilon$$
 constant γ

The noiseless approximation guarantee [Das & Kempe, ICML'11; Qian, Yu and Zhou, NIPS'15]

$$f(X) \ge \left(1 - \left(1 - \frac{\gamma}{B}\right)^B\right) \cdot \text{OPT} \ge (1 - e^{-\gamma}) \cdot \text{OPT} \quad \begin{array}{c} \text{a constant} \\ \text{approximation ratio} \end{array}$$

The performance degrades largely in noisy environments

 $\varepsilon \leq 1/B$ for a constant

PONSS

In our previous work, threshold selection has been theoretically shown to be robust against noise [Qian et al., ECJ'18]

$$\begin{cases} f(X) \ge f(Y) & \longrightarrow & f(X) \ge f(Y) + \theta \\ \text{A solution is better if its objective value is larger} \\ \text{by least a threshold} \\ \end{cases} \begin{array}{l} \text{Exponentially} \\ \text{decrease the} \\ \text{running time} \\ \end{cases} \\ \text{POSS [Qian et al., NIPS'15]} & X \leqslant Y \Leftrightarrow \begin{cases} f(X) \ge f(Y) \\ |X| \le |Y| \\ \end{bmatrix} \\ \\ \text{Reduce the risk} \\ \text{of deleting a} \\ \text{good solution} \\ \end{cases} \\ \\ \text{Multiplicative noise:} & X \leqslant Y \Leftrightarrow \begin{cases} f(X) \ge f(Y) \\ |X| \le |Y| \\ \end{bmatrix} \\ \\ \\ f(X) \ge \frac{1+\theta}{1-\theta}f(Y) \\ |X| \le |Y| \\ \end{cases} \end{cases}$$

PONSS

In our previous work, threshold selection has been theoretically shown to be robust against noise [Qian et al., ECJ'18]

$$f(X) \ge f(Y) \longrightarrow f(X) \ge f(Y) + \theta$$
A solution is better if its objective value is larger
by least a threshold
$$\begin{array}{c} \text{Commate''} \\ \text{Commate''}$$

Multiplicative noise:

$$\begin{array}{ll} \text{PONSS} & f(X) \geq \frac{1-\epsilon}{1+\epsilon} \left(1 - \left(1 - \frac{\gamma}{B}\right)^B\right) \cdot \text{OPT} & \text{better} \\ (\theta \geq \epsilon) \end{array}$$

$$\begin{array}{l} \text{POSS \& Greedy} & f(X) \geq \frac{1}{1 + \frac{2\epsilon B}{(1-\epsilon)\gamma}} \left(1 - \left(\frac{1-\epsilon}{1+\epsilon}\right)^B \left(1 - \frac{\gamma}{B}\right)^B\right) \cdot \text{OPT} \end{array}$$

$$\begin{aligned} \frac{1-\epsilon}{1+\epsilon} \left(1-\left(1-\frac{\gamma}{B}\right)^B\right) &= \frac{1-\epsilon}{1+\epsilon} \cdot \frac{\gamma}{B} \cdot \sum_{i=0}^{B-1} \left(1-\frac{\gamma}{B}\right)^i \geq \frac{1-\epsilon}{1+\epsilon} \cdot \frac{\gamma}{B} \cdot \sum_{i=0}^{B-1} \left(\frac{1-\epsilon}{1+\epsilon} \left(1-\frac{\gamma}{B}\right)\right)^i \\ &= \frac{1-\epsilon}{1+\epsilon} \cdot \frac{\gamma}{B} \cdot \left(1-\left(\frac{1-\epsilon}{1+\epsilon}\right)^B \left(1-\frac{\gamma}{B}\right)^B\right) / \left(1-\frac{1-\epsilon}{1+\epsilon} \left(1-\frac{\gamma}{B}\right)\right) \\ &= \frac{1}{1+\frac{2\epsilon B}{(1-\epsilon)\gamma}} \left(1-\left(\frac{1-\epsilon}{1+\epsilon}\right)^B \left(1-\frac{\gamma}{B}\right)^B\right) \end{aligned}$$

http://www.lamda.nju.edu.cn/qianc/

http://www.lamda.nju.edu.cn/yuy/

Multiplicative noise:

$$\begin{array}{ll} \text{PONSS} & f(X) \geq \frac{1-\epsilon}{1+\epsilon} \left(1 - \left(1 - \frac{\gamma}{B}\right)^B\right) \cdot \text{OPT} & \begin{array}{l} \text{Significantly} \\ \text{better} \end{array}$$

$$\begin{array}{l} \text{POSS \& Greedy} & f(X) \geq \frac{1}{1 + \frac{2\epsilon B}{(1-\epsilon)\gamma}} \left(1 - \left(\frac{1-\epsilon}{1+\epsilon}\right)^B \left(1 - \frac{\gamma}{B}\right)^B\right) \cdot \text{OPT} \end{array}$$

 $\gamma = 1$ (submodular), ϵ is a constant

PONSSa constant approximation ratioPOSS & Greedy $\Theta(1/B)$ approximation ratio

Multiplicative noise:

PONSS
$$f(X) \ge \frac{1-\epsilon}{1+\epsilon} \left(1 - \left(1 - \frac{\gamma}{b}\right)^B \right) \cdot \text{OPT}$$
 better
Greedy $f(X) \ge \frac{1}{1+\frac{2\epsilon B}{(1-\epsilon)\gamma}} \left(1 - \left(\frac{1-\epsilon}{1+\epsilon}\right)^B \left(1 - \frac{\gamma}{B}\right)^B \right) \cdot \text{OPT}$

Additive noise:

POSS &

PONSS
$$f(X) \ge \left(1 - \left(1 - \frac{\gamma}{B}\right)^B\right) \cdot \text{OPT} - 2\epsilon$$
 better
POSS & Greedy $f(X) \ge \left(1 - \left(1 - \frac{\gamma}{B}\right)^B\right) \cdot \text{OPT} - \left(\frac{2B}{\gamma} - \frac{2B}{\gamma}e^{-\gamma}\right)\epsilon$
 $\frac{2B}{\gamma} - \frac{2B}{\gamma}e^{-\gamma} \ge 2$

http://www.lamda.nju.edu.cn/qianc/

http://www.lamda.nju.edu.cn/yuy/

0.

• • •

.1

Experimental results - influence maximization

PONSS (red line) vs POSS (blue line) vs Greedy (black line):

- Noisy evaluation: the average of 10 independent Monte Carlo simulations
- The output solution: the average of 10,000 independent Monte Carlo simulations



Experimental results - sparse regression

PONSS (red line) vs POSS (blue line) vs Greedy (black line):

- Noisy evaluation: a random sample of 1,000 instances
- The output solution: the whole data set



Experimental results – sensitivity to θ

PONSS (red line) vs POSS (blue line) vs Greedy (black line):

"dominate"

$$X \leq Y \Leftrightarrow \begin{cases} f(X) \ge \frac{1+\theta}{1-\theta} f(Y) \\ |X| \le |Y| \end{cases}$$

The performance of PONSS is not very sensitive to θ



http://www.lamda.nju.edu.cn/qianc/

http://www.lamda.nju.edu.cn/yuy/



Introduction

Pareto optimization for subset selection

□ Pareto optimization for large-scale subset selection

□ Pareto optimization for noisy subset selection

Pareto optimization for dynamic subset selection
 Conclusion

Dynamic sensor placement

Sensor placement [Krause & Guestrin, IJCAI'09 Tutorial]: select a few places to install sensors such that the information gathered is maximized



Fire detection

12 sensors (sensor failure)

10 sensors

15 sensors

(more investment)

How about the performance for dynamic subset selection?

Dynamic subset selection

Subset selection with general constraints: given $V = \{v_1, ..., v_n\}$, an objective function $f: 2^V \to \mathbb{R}$, a cost function $c: 2^V \to \mathbb{R}$ and a budget B, to find a subset $X \subseteq V$ such that $max_{X \subseteq V} \quad f(X) \quad s.t. \quad c(X) \leq B$

Dynamic subset selection [Roostapour, Neumann, Neumann and Friedrich, AAAI'19]

The available resources may change over time



The budget *B* may change over time

To examine: Can an algorithm find a good solution quickly for the new problem, when starting from the solutions obtained for the old problem?

Compare Pareto optimization with the greedy algorithm

Both of them achieve the best-known approximation guarantee for the static problem [Zhang & Vorobeychik, AAAI'16; Qian et al., IJCAI'17]

$$max_{X\subseteq V} f(X)$$
 s.t. $c(X) \le B$

Approximation ratio: $(\alpha/2)(1-e^{-\alpha})$ [Zhang & Vorobeychik, AAAI'16]

Process: iteratively select one item making the ratio of the increment on *f* and *c* maximzied

$$v^* = \arg \max_{v \in V \setminus X^{j-1}} \frac{f(X^{j-1} \cup \{v\}) - f(X^{j-1})}{c(X^{j-1} \cup \{v\}) - c(X^{j-1})}$$

 X^{j} : the subset obtained after *j* iterations



:n/qianc/ http://wwv

http://www.lamda.nju.edu.cn/qianc/

http://www.lamda.nju.edu.cn/yuy/

Pareto optimization - POMC

Approximation ratio:
$$(\alpha/2)(1 - e^{-\alpha})$$
 [Qian et al., IJCAI'17] $max_{x \in \{0,1\}^n} f(x)$ s.t. $c(x) \leq B$ originalTransformation: $\begin{subarray}{c} min_{x \in \{0,1\}^n} (-f(x), c(x)) & bi-objective \end{subarray}$

Algorithm 2 POMC Algorithm

Input: a monotone objective function f, a monotone approximate cost function \hat{c} , and a budget B **Parameter**: the number T of iterations **Output**: a solution $x \in \{0, 1\}^n$ with $\hat{c}(x) \leq B$ **Process**: 1: Let $x = \{0\}^n$ and $P = \{x\}$.

1. Let $x = \{0\}$ 2: Let t = 0.

3. while t < T do

4: Select x from P uniformly at random.

5: Generate x' by flipping each bit of x with prob. 1/n.

6: **if**
$$\nexists z \in P$$
 such that $z \succ x'$ **then**

7:
$$P = (P \setminus \{ \boldsymbol{z} \in P \mid \boldsymbol{x'} \succeq \boldsymbol{z} \}) \cup \{ \boldsymbol{x'} \}$$

- 8: **end if**
- 9: t = t + 1.

10: end while

11: return $\arg \max_{\boldsymbol{x} \in P: \hat{c}(\boldsymbol{x}) \leq B} f(\boldsymbol{x})$

Initialization: put the special solution $\{0\}^n$ into the population *P*

Reproduction: pick a solution x randomly from P, and flip each bit of x with prob. 1/n to produce a new solution

Updating: if the new solution is not dominated by any solution in *P*, put it into *P* and weed out bad solutions

Output: select the best feasible solution

http://www.lamda.nju.edu.cn/qianc/

[Roostapour, Neumann, Neumann and Friedrich, AAAI'19]

The greedy algorithm can achieve arbitrarily bad approximation ratios during a sequence of dynamic changes

Theorem 1. For dynamic subset selection, there exist instances of dynamically increasing *B* and decreasing *B* such that the approximation ratios of the greedy algorithm are O(1/n) and $O(1/\sqrt{n})$, respectively.

POMC can maintain good approximation ratios efficiently

Theorem 2. For dynamic subset selection, with a constant probability, POMC achieves an approximation ratio of $(\alpha/2)(1 - e^{-\alpha})$ for any budget $b \in [0, B]$ after $cnP_{max}B/\delta$ iterations. (Already good for decreasing *B*)

Theorem 3. For dynamic subset selection with *B* increasing to B^* , with a constant probability, POMC achieves an approximation ratio of $(\alpha/2)(1 - e^{-\alpha})$ for any budget $b \in [0, B^*]$ after $cnP_{max}(B^* - B)/\delta$ iterations.

Experimental results - influence maximization

GGA: the greedy algorithm starting from scratch for each new budget

AGGA: the greedy algorithm

POMC_{τ}: POMC running τ iterations for each new budget

POMC^{WP}_{τ}: POMC_{τ} with a warm-up phase (running 10,000 iterations for the initial *B*)

Change of the budget *B*:



Changes	GGA		AGGA		$POMC_{1000}$		$POMC_{5000}$		$POMC_{10000}$		$POMC_{1000}^{WP}$		$POMC_{5000}^{WP}$		$POMC_{10000}^{WP}$	
	mean	st	mean	st	mean	st	mean	st	mean	st	mean	st	mean	st	mean	st
1-25	85.0349	12.88	81.5734	14.07	66.3992	17.95	77.8569	18.76	86.1057	17.22	86.3846	10.76	86.9270	12.86	85.8794	14.69
26-50	100.7344	22.16	96.1386	23.99	104.9102	15.50	117.6439	16.71	122.5604	15.54	110.4279	11.08	115.6766	14.21	120.8651	14.97
51-75	118.1568	30.82	110.4893	29.50	141.8249	5.64	155.2126	5.08	158.7228	5.20	140.7838	5.02	149.7658	5.49	157.6169	5.54
76-100	127.3422	31.14	115.2978	27.66	149.0259	3.36	159.9100	3.28	162.7353	3.65	148.3012	3.47	155.1943	4.04	163.1958	3.74
101-125	132.3502	29.62	116.9768	25.45	150.3415	3.17	160.1367	2.81	161.2852	2.68	148.5254	2.67	155.1104	3.05	162.3770	2.81
126-150	134.5256	27.69	118.6962	24.19	147.8998	7.36	154.7319	8.77	154.1470	7.43	143.4908	7.96	150.7567	7.82	156.0363	8.12
151-175	135.7651	25.89	119.4982	22.85	147.2478	4.68	153.1417	5.32	151.2966	3.17	143.2959	4.79	149.5447	4.87	153.2526	3.85
176-200	135.5133	24.41	119.1491	22.04	139.5072	8.08	143.6928	9.16	143.9832	8.67	134.7968	8.72	140.5930	8.61	144.4088	8.08

 $POMC_{\tau}$ achieves better performance than GGA and AGGA after 25 changes, and $POMC_{\tau}^{WP}$ can bring improvement in the first 25 changes

http://www.lamda.nju.edu.cn/qianc/

http://www.lamda.nju.edu.cn/yuy/

Conclusion



- Pareto optimization for subset selection
 - Show excellent performance theoretically and empirically
- Pareto optimization for large-scale subset selection
 - Introduce parallel and distributed strategies
- Pareto optimization for noisy subset selection
 - Introduce noise-aware domination relationship
- Pareto optimization for dynamic subset selection
 - Show robustness against dynamic changes

Future work

- Problem issues
 - Non-monotone objective functions [Qian et al., AIJ'19; Do & Neumann, PPSN'20]
 - Continuous submodular objective functions
 - More complex constraints [Neumann & Neumann, PPSN'20; Do & Neumann, PPSN'20]
 - More uncertain environments
 - Other than subset selection [Neumann & Wegener, NC'06; Friedrich et al., ECJ'10;

Neumann et al., Algorithmica'11; Qian et al., IJCAI'15; Crawford, IJCAI'19; Assimi et al., ECAI'20]

- Algorithm issues
 - More complicated MOEAs
- Theory issues
 - Beat the best-known approximation guarantee
- Application issues
 - Attempts on more real-world applications

- F. Neumann and I. Wegener. Minimum spanning trees made easier via multiobjective optimization. *Natural Computing*, 2006, 5(3): 305-319.
- T. Friedrich, J. He, N. Hebbinghaus, F. Neumann and C. Witt. Approximating covering problems by randomized search heuristics using multi-objective models. *Evolutionary Computation*, 2010, 18(4): 617-633.
- F. Neumann, J. Reichel and M. Skutella. Computing minimum cuts by randomized search heuristics. *Algorithmica*, 2011, 59(3): 323-342.
- T. Friedrich and F. Neumann. Maximizing submodular functions under Matroid constraints by evolutionary algorithms. *Evolutionary Computation*, 2015, 23(4): 543-558.
- C. Qian, Y. Yu and Z.-H. Zhou. Pareto ensemble pruning. In: *Proceedings of the 29th AAAI Conference on Artificial Intelligence (AAAI'15)*, Austin, TX, 2015.

- C. Qian, Y. Yu and Z.-H. Zhou. On constrained Boolean Pareto optimization. In: *Proceedings of the 24th International Joint Conference on Artificial Intelligence* (*IJCAI'15*), Buenos Aires, Argentina, 2015.
- C. Qian, Y. Yu and Z.-H. Zhou. Subset selection by Pareto optimization. In: *Advances in Neural Information Processing Systems 28 (NIPS'15)*, Montreal, Canada, 2015.
- C. Qian, J.-C. Shi, Y. Yu, K. Tang and Z.-H. Zhou. Parallel Pareto optimization for subset selection. In: *Proceedings of the 25th International Joint Conference on Artificial Intelligence (IJCAI'16)*, New York, NY, 2016.
- C. Qian, J.-C. Shi, Y. Yu and K. Tang. On subset selection with general cost constraints. In: *Proceedings of the 26th International Joint Conference on Artificial Intelligence (IJCAI'17)*, Melbourne, Australia, 2017.
- C. Qian, J.-C. Shi, Y. Yu, K. Tang and Z.-H. Zhou. Optimizing ratio of monotone set functions. In: *Proceedings of the 26th International Joint Conference on Artificial Intelligence (IJCAI'17)*, Melbourne, Australia, 2017.

- C. Qian, J.-C. Shi, Y. Yu, K. Tang and Z.-H. Zhou. Subset selection under noise. In: *Advances in Neural Information Processing Systems 30 (NIPS'17)*, Long Beach, CA, 2017.
- C. Qian, J.-C. Shi, K. Tang and Z.-H. Zhou. Constrained monotone *k*-submodular function maximization using multi-objective evolutionary algorithms with theoretical guarantee. *IEEE Transactions on Evolutionary Computation*, 2018, 22(4), 595-608.
- C. Qian, Y. Zhang, K. Tang and X. Yao. On multiset selection with size constraints. In: *Proceedings of the 32nd AAAI Conference on Artificial Intelligence (AAAI'18)*, New Orleans, LA, 2018.
- C. Qian, G. Li, C. Feng and K. Tang. Distributed Pareto optimization for subset selection. In: *Proceedings of the 27th International Joint Conference on Artificial Intelligence (IJCAI'18)*, Stockholm, Sweden, 2018.
- C. Qian, C. Feng and K. Tang. Sequence selection by Pareto optimization. In: *Proceedings of the 27th International Joint Conference on Artificial Intelligence* (*IJCAI'18*), Stockholm, Sweden, 2018.

- C. Qian, Y. Yu and Z.-H. Zhou. Analyzing evolutionary optimization in noisy environments. *Evolutionary Computation*, 2018, 26(1): 1-41.
- C. Feng, C. Qian and K. Tang. Unsupervised feature selection by Pareto optimization. In: *Proceedings of the 33rd AAAI Conference on Artificial Intelligence* (AAAI'19), Honolulu, HI, 2019.
- V. Roostapour, A. Neumann, F. Neumann and T. Friedrich. Pareto optimization for subset selection with dynamic cost constraints. In: *Proceedings of the 33rd AAAI Conference on Artificial Intelligence (AAAI'19)*, Honolulu, HI, 2019.
- V. Crawford. An efficient evolutionary algorithm for minimum cost submodular cover. In: *Proceedings of the 28th International Joint Conference on Artificial Intelligence (IJCAI'19)*, Macao, China, 2019.
- V. Crawford and A. Kuhnle. Fast evolutionary algorithms for maximization of cardinality-constrained weakly submodular functions. *arXiv:1908.01230*, 2019.

- C. Qian, Y. Yu, K. Tang, X. Yao and Z.-H. Zhou. Maximizing submodular or monotone approximately submodular functions by multi-objective evolutionary algorithms. *Artificial Intelligence*, 2019, 275: 279-294.
- H. Assimi, O. Harper, Y. Xie, A. Neumann and F. Neumann. Evolutionary biobjective optimization for the dynamic chance-constrained knapsack problem based on tail bound objectives. In: *Proceedings of 24th European Conference on Artificial Intelligence (ECAI'20)*, Santiago de Compostela, Spain, 2020.
- C. Qian, C. Bian and C. Feng. Subset selection by Pareto optimization with recombination. In: *Proceedings of the 34th AAAI Conference on Artificial Intelligence (AAAI'20)*, New York, NY, 2020.
- C. Bian, C. Feng, C. Qian and Y. Yu. An efficient evolutionary algorithm for subset selection with general cost constraints. In: *Proceedings of the 34th AAAI Conference on Artificial Intelligence (AAAI'20)*, New York, NY, 2020.

- A. V. Do and F. Neumann. Maximizing submodular or monotone functions under partition matroid constraints by multi-objective evolutionary algorithms. In: *Proceedings of the 16th International Conference on Parallel Problem Solving from Nature (PPSN'20)*, Leiden, The Netherlands, 2020.
- A. Neumann and F. Neumann. Optimising chance-constrained submodular functions using evolutionary multi-objective algorithms. In: *Proceedings of the 16th International Conference on Parallel Problem Solving from Nature (PPSN'20),* Leiden, The Netherlands, 2020.
- C. Qian. Distributed Pareto optimization for large-scale noisy subset selection. *IEEE Transactions on Evolutionary Computation*, 2020, 24(4): 694-707.

THANK YOU !