





SUBMODULAR HAMMING METRICS JENNIFER GILLENWATER **BETHANY LUSCH** RAHUL KIDAMBI **RISHABH IYER**

{jengi, rkiyer, herwaldt, rkidambi, bilmes}@uw.edu



- First term: Ensures k-th summary is internally diverse.
- Second term: Ensures k-th summary is unlike previous k-1.

JEFF BILMES

SH-MIN AND SH-MAX ALGORITHMS

Note: $f(A \triangle B)$ is submodular in $A \triangle B$, but not necessarily in A. Thus, we cannot trivially use standard submodular optimization algorithms for SH-min and SH-max.

Algorithm 1 UNION-SPLIT						
Define	$(1 (\Lambda))$	C (A)				

Define $f'_i(A) = f_i(A \setminus B_i) + f_i(B_i \setminus A)$ Define $F'(A) = \sum_{i=1}^m f'_i(A)$ **Output:** STANDARD-SUBMOD-OPT (F')

 $A \leftarrow B_1$ for i = 2, ..., m do if $F(B_i) < F(A)$: $A \leftarrow B_i$ **Output**: A

Alg 1—key insight: $A \triangle B = (A \setminus B) \cup (B \setminus A)$ and both $f(A \setminus B)$ and $f(B \setminus A)$ are submodular in A itself.

Alg 2—main idea: If $f_i = f_j \forall i, j$ (call this the "homogeneous" case), then one of the B_i is a reasonable A.

Algorithm 3 MAJOR-MIN

 $A \leftarrow \emptyset$ repeat $c \leftarrow F(A)$ Define $S_j = V \setminus j, T_i = A \triangle B_i$ Set $\boldsymbol{w}_{\hat{F}}(j) = \sum_{i=1}^{m} \begin{cases} f_i(j \mid S_j) \text{ if } j \in A \triangle B_i \\ f_i(j \mid T_i) \text{ otherwise} \end{cases}$ $A \leftarrow \text{MODULAR-MIN}(\boldsymbol{w}_{\hat{F}}, \mathcal{C})$ until F(A) = c**Output**: A

Main idea: Construct \hat{F} , a modular upper bound for F at the current solution A, then (exactly) minimize \hat{F} to get a new A. Iterate until convergence.

THEORETICAL RESULTS

 $n = \text{size of ground set } V, \ m = \# \text{ of } f_i,$ $\kappa_f = \text{curvature (larger values imply } f \text{ is close to modular)}$

Table 1: Hardness for SH-min and SH-max. UC stands for unconstrained, and Card stands for cardinality-constrained. The entry "open" implies that the problem is potentially poly-time solvable.

	SH-min		SH-max	
	homogeneous	heterogeneous	homogeneous	heterogeneous
UC	Open	4/3	3/4	3/4
Card	$ \ \ \Omega\left(\frac{\sqrt{n}}{1+(\sqrt{n}-1)(1-\kappa_f)}\right) $	$\Omega\left(\frac{\sqrt{n}}{1+(\sqrt{n}-1)(1-\kappa_f)}\right)$	1 - 1/e	1 - 1/e

Table 2: Approximation guarantees of algorithms for SH-min and SH-max. '-' implies that no guarantee holds for the corresponding pair. BEST-B only works for the homogeneous case, while all other algorithms work in both cases.

	UNION-SPLIT		BEST-B	MAJOR-MIN	RAND-SET
	UC	Card	UC	Card	UC
SH-min	2	-	2 - 2/m	$\frac{n}{1+(n-1)(1-\kappa_f)}$	_
SH-max	1/4	1/2e	-	-	1/8

SH-MIN EXPERIMENT

Synthetic document clustering: In a setting with 1000 "words" and 100 "documents" assigned to 10 "true" clusters, if each document is a random sampling of 10 words from the 10 word classes W associated with its cluster, then k-means clustering accuracy with kmeans++ initialization is—

Hamming: 57.0% (± 6.8), versus SH: 88.5% (± 8.4)

SH-MAX EXPERIMENT

Diverse k-best image collection summarization: Given 14 image collections, each containing 100 photos, we seek k = 15summaries of size 10 for each. We then use V-ROUGE to measure summary quality.

Table 4: mV-ROUGE (avg over 14 datasets, \pm std dev). Table 5: # of wins (out of 14 datasets).

HM = Hamming optimized via greedy, SP = SH optimized via greedy, TP = SH optimized via UNION-SPLIT

Summaries of the 6th image collection. HM approach (left) and TP approach (right). Images in the green rectangle tend to be more redundant with images from the previous summaries in the HM case than in the TP case; the HM solution contains many images with a "sky" theme, while TP contains more images with other themes. Quality of the individual summaries also tends to become poorer for the later HM sets; considering the images in the red rectangles overlaid on the montage, the HM sets contain many images of tree branches. By contrast, the TP summary quality remains reasonable even for the last few summaries.