

1) # of clusters

2) Homogeneity may vary across clusters

3) What features to use?

4) What names to assign to clusters?

↳ Dominant?

↳ How about leveraging the wiki topic hierarchy for cluster names?

Suggestions: ① Correlate clusters with category ids?

How abt accounting for category structure in clustering process?

Desirable properties of clusters in DAG

① Independence: Cover disjoint doc set/diverse

Loose Orthogonality

② Exhaustive: High coverage

③ Intra cluster similarity - ④ Depth (Quality Ctrl)

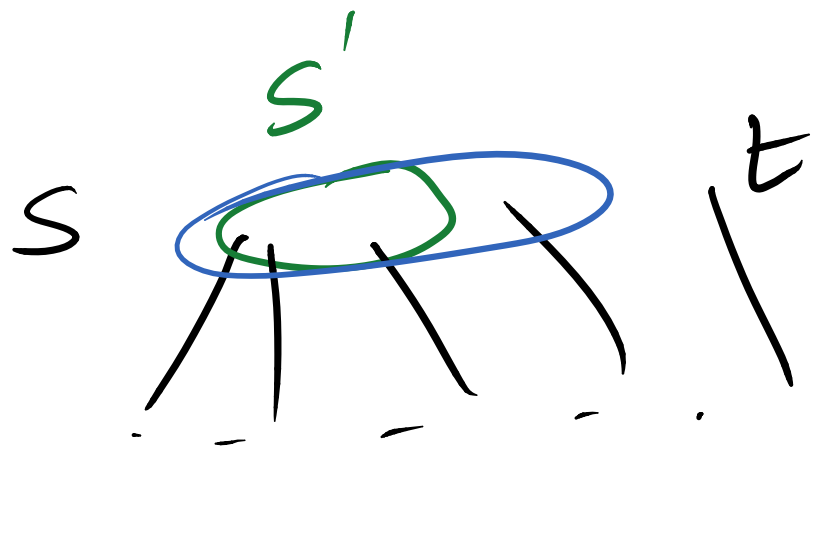
$$\Gamma(S) = \bigcup_{s \in S} \Gamma(s)$$

$$S' \subseteq S$$
$$t \notin S, S'$$

$$|\Gamma(S)| \geq |\Gamma(S')|$$

$$\Gamma(S), \Gamma(t), \Gamma(S')$$

$$\Gamma(S \cup \{t\}) \quad \Gamma(S' \cup \{t\})$$

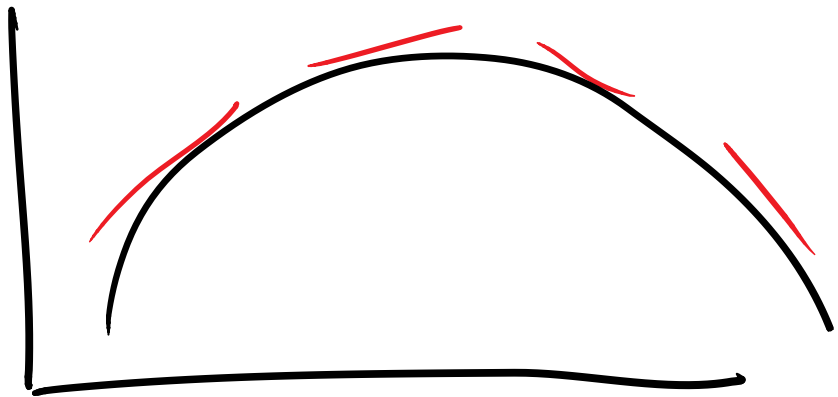


$$\Gamma(S' \cup \{t\}) - \Gamma(S') \geq \Gamma(S \cup \{t\})$$

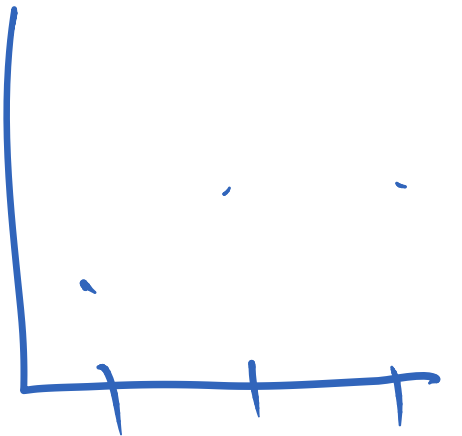
Do submodular fns resemble concave or convex fns in their diminishing returns property?

$$s' \subseteq s \subseteq V \setminus \{v\}$$

$$\underline{f(s' \cup \{v\}) - f(s')} \geq \underline{f(s \cup \{v\}) - f(s)}$$



Slope decreases



$$\Gamma(S) = \bigcup_{s \in S} \Gamma(s)$$

$$\Psi(S) = \sum_{s \in S} |\Gamma(s)|$$



Submodular

Is this submodular?

$$m_u(S) = f\left(\bigcup_{s \in S} s\right)$$

$$\sum_{s \in S} m_u(s)$$

$$\Psi(S) = \sum_{s \in S} \Psi(\{s\})$$

Ans: Modular

$\Psi(m_u(S))$  st  $\Psi$  is concave  $\Rightarrow \Psi(m_u(S)) \geq \sum_{s \in S} \Psi(\{s\})$