

CS60021: Scalable Data Mining

Subset Selection

Sourangshu Bhattacharya

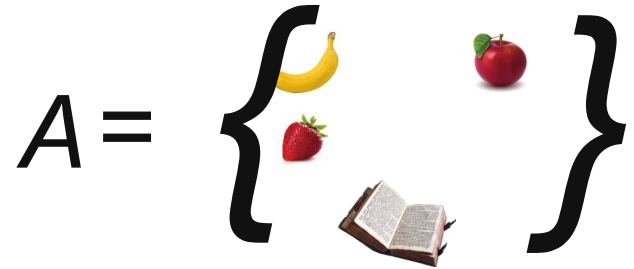
Submodular Subset Selection

Slides taken from IJCAI 2020 tutorial by
Rishabh Iyer and Ganesh Ramakrishnan

Combinatorial Subset Selection Problems



$$f : 2^V \rightarrow \mathbb{R}$$



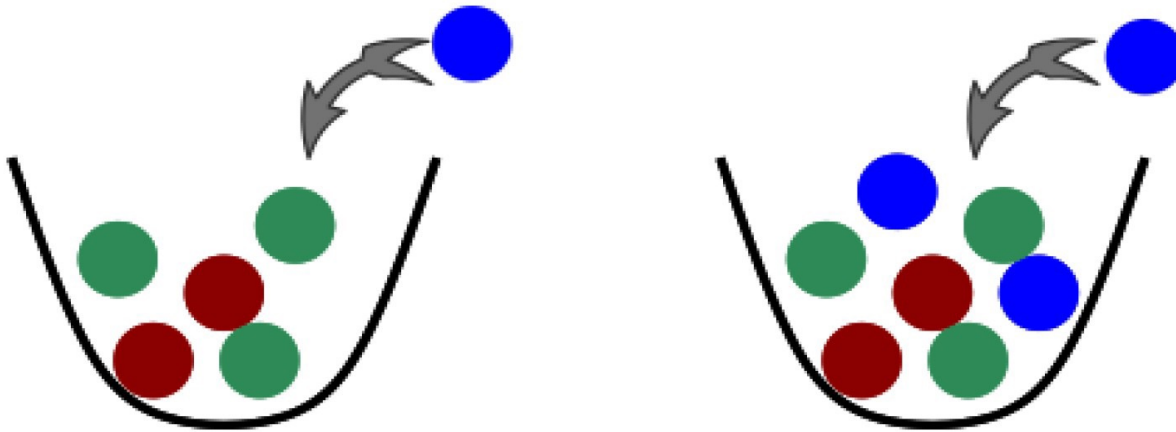
Choose Subset $A \subseteq V$
 $f(A)$ is maximum

General Set function Optimization: very hard!

What if there is some special structure?

Submodular Functions

$$f(A \cup v) - f(A) \geq f(B \cup v) - f(B), \text{ if } A \subseteq B$$



$f = \#$ of distinct colors of balls in the urn.

Negative of a
Submodular
Function is a
Supermodular
Function!

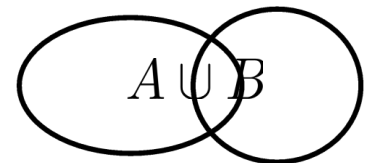
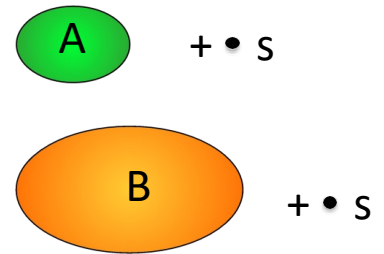
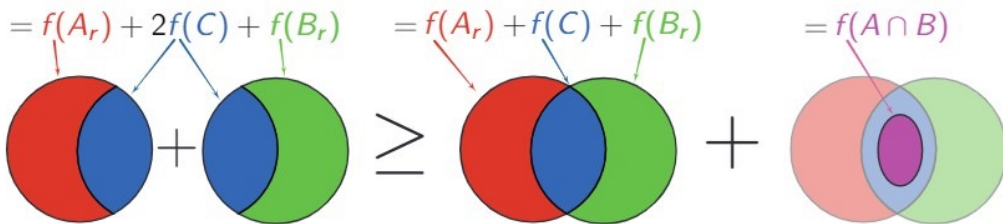
Equivalent Definitions of Submodularity

- **Diminishing gains:** for all $A, B \subseteq V$

$$f(A \cup v) - f(A) \geq f(B \cup v) - f(B), \text{ if } A \subseteq B$$

- **Union-Intersection:** for all $A, B \subseteq V$

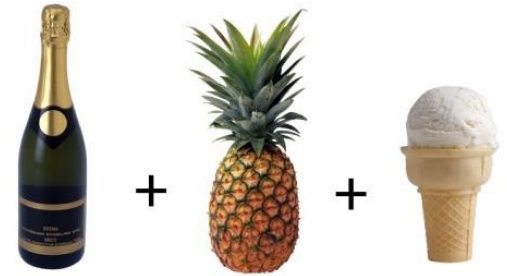
$$f(A) + f(B) \geq f(A \cup B) + f(A \cap B)$$



Modular Functions

- each element e has a weight $w(e)$

$$F(S) = \sum_{e \in S} w(e)$$



$$A \subset B$$

$$F(A \cup e) - F(A) = w(e) = F(B \cup e) - F(B) = w(e)$$

Modular Functions are both submodular and supermodular!

Monotone Submodular Functions

- A set function is called **monotonic** if

$$A \subseteq B \subseteq V \Rightarrow F(A) \leq F(B)$$

- Examples:

- **Influence** in social networks [Kempe et al KDD '03]
- For discrete RVs, **entropy** $F(A) = H(X_A)$ is monotonic:
Suppose $B = A \cup C$. Then
$$F(B) = H(X_A, X_C) = H(X_A) + H(X_C | X_A) \geq H(X_A) = F(A)$$
- **Information gain**: $F(A) = H(Y) - H(Y | X_A)$

Instantiations of Submodular Functions

Representation Functions

- Facility Location Function (k-medoids clustering)
- Graph Cut Family, Saturated Coverage

Diversity Functions

- Dispersion Functions (Min, Sum, Min-Sum)
- Determinantal Point Processes

Coverage Functions

- Set Cover Function
- Probabilistic Set Cover Function
- Feature Based Functions

Importance Functions

- Modular Functions

Information Functions

- Mutual Information
- Entropy

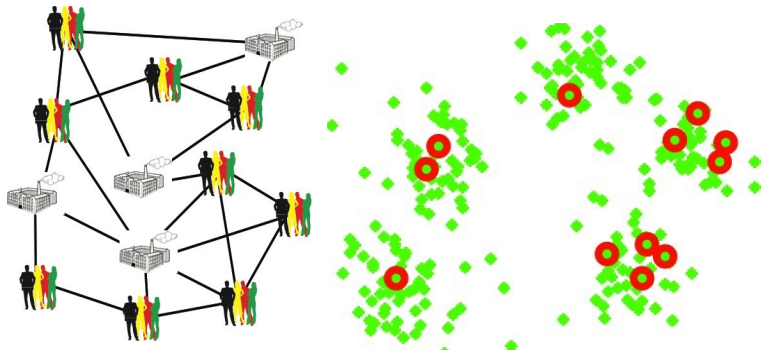
Discounted Cost Functions

- Clustered Concave over Modular Functions
- Cooperative Costs and Saturations

Complexity Functions

- Bipartite Neighborhood Functions

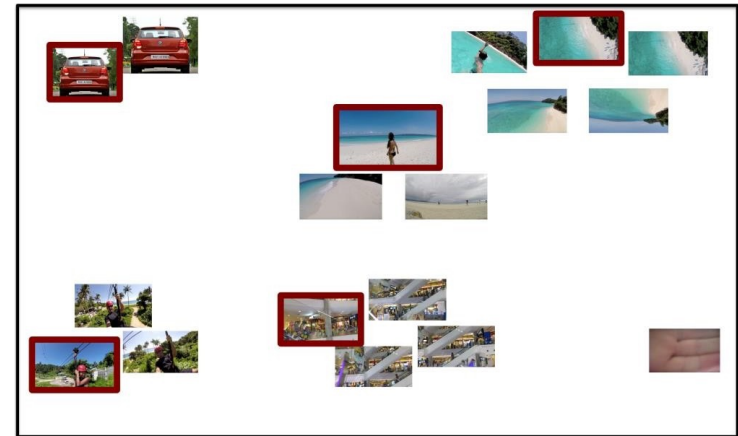
Representation Functions



| | |
|--------------------|---|
| Facility Location | $\sum_{i \in V} \max_{k \in X} s_{ik}$ |
| Saturated Coverage | $\sum_{i \in V} \min\{\sum_{j \in X} s_{ij}, \alpha_i\}$ |
| Graph Cut | $\lambda \sum_{i \in V} \sum_{j \in X} s_{ij} - \sum_{i, j \in X} s_{ij}$ |



Similarity Kernel

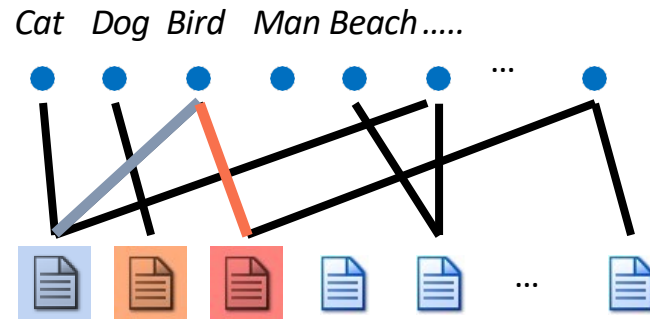


Representation Functions

Picks Centroids

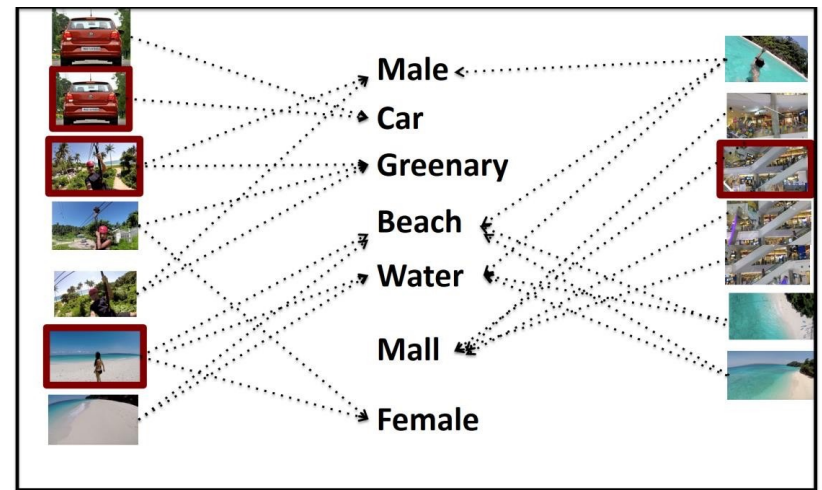
Iyer 2015, Kaushal et al 2019, Tschatchek et al 2014, ...

Coverage Functions



$$f(X) = w(\cup_{i \in X} U_i),$$

↑
Concepts Covered by Instance i

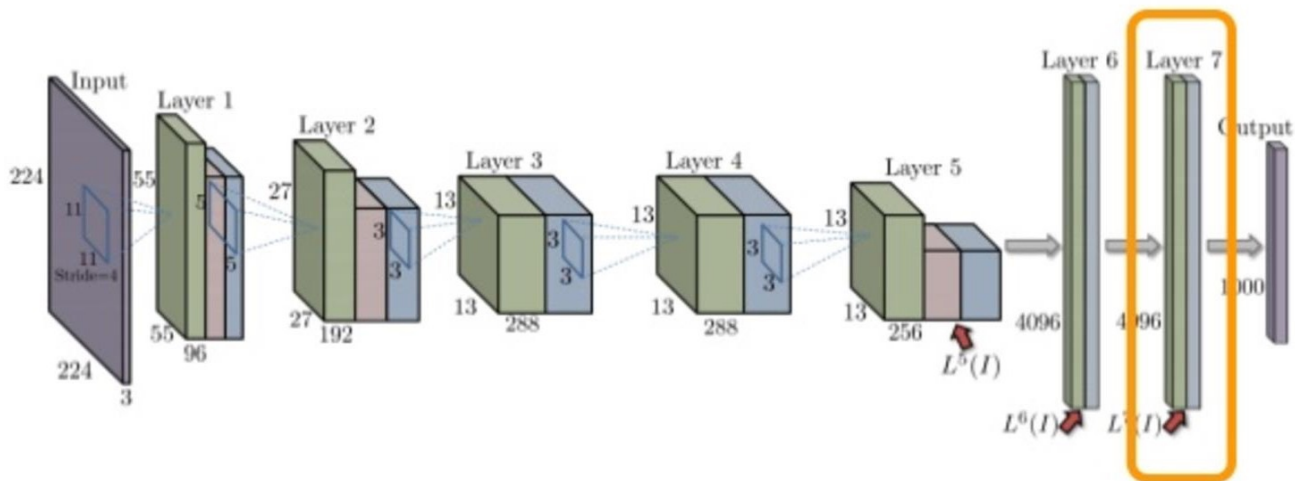


Coverage Functions

Select instances which “cover” all concepts

Wolsey et al 1982, ...

Feature Based Functions



Achieve
Uniformity in
Feature
Coverage

Feature Based Functions

$$f_{\text{fea}}(S) = \sum_{u \in \mathcal{U}} g(m_u(S)).$$

↑

Total Contribution of Feature u in the Set of Images S

Wei-lyer et al 2014...

Information Functions

X_1, \dots, X_n discrete random variables: $X_e \in \{1, \dots, m\}$

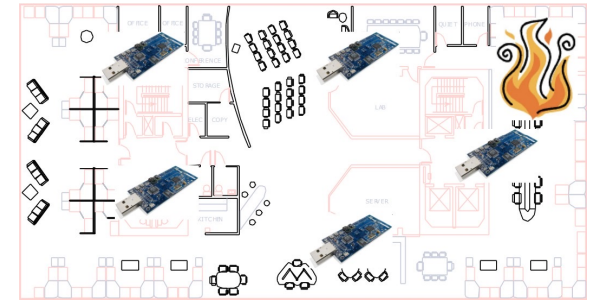
$F(S) = H(X_S) =$ joint entropy of variables indexed by S

$$H(X_e) = \sum_{x \in \{1, \dots, m\}} P(X_e = x) \log P(X_e = x)$$

$$A \subset B, e \notin B \quad F(A \cup e) - F(A) \geq F(B \cup e) - F(B)??$$

$$\begin{aligned} H(X_{A \cup e}) - H(X_A) &= H(X_e | X_A) \\ &\leq H(X_e | X_B) \quad \text{“information never hurts”} \\ &= H(X_{B \cup e}) - H(X_B) \end{aligned}$$

discrete entropy is submodular!



Entropy
Mutual Information
Information Gain

...

Krause et al 2008, ...

Master Optimization Problem

Set Function \rightarrow Selected set

$$\max_{\mathcal{A} \subseteq \mathcal{V}} F(\mathcal{A})$$

Selection cost \rightarrow Budget

$$\text{subject to } C(\mathcal{A}) \leq B$$

F = Monotone Submodular,
Non Monotone Submodular,
Dispersion Functions,
....

F Models:

- Diversity
- Representation
- Coverage
- Information
- Importance
- ...

We shall study this and variants of this Master Optimization Problem!

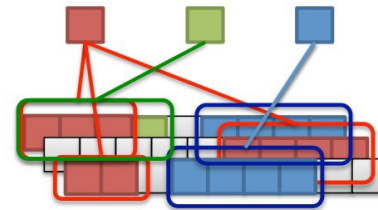
Monotone Submodular Maximization

$$\max_S F(S) \text{ s.t. } |S| \leq k$$

What is the Constraint?
 $C(S) = |S|$

- greedy algorithm:

$$\begin{aligned} S_0 &= \emptyset \\ \text{for } i &= 0, \dots, k-1 \\ e^* &= \arg \max_{e \in \mathcal{V} \setminus S_i} F(S_i \cup \{e\}) \\ S_{i+1} &= S_i \cup \{e^*\} \end{aligned}$$

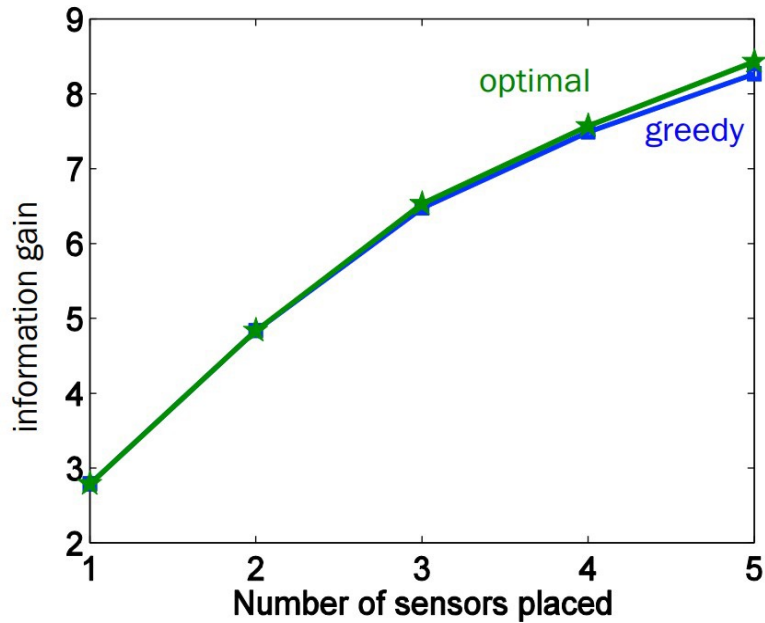


How “good” is S_k ?

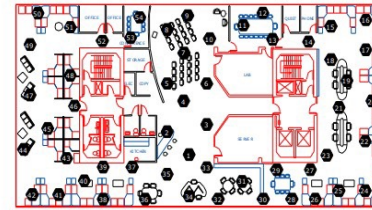
Approximation
Guarantee!

How good is Greedy in Practice?

empirically:



sensor placement



How good is Greedy in Theory?

$$\max_S F(S) \text{ s.t. } |S| \leq k$$

Theorem (Nemhauser, Fisher, Wolsey '78)

F monotone submodular, S_k solution of greedy. Then

$$F(S_k) \geq \left(1 - \frac{1}{e}\right) F(S^*)$$

← optimal solution

No Poly-time algorithm can do better than this in the worst case!

Proof (Nemhauser et al 1978)

Let:

- $A_i = (v_1, v_2, \dots, v_i)$ be the the chain formed by the greedy algorithm, as defined above
- $A^* = (v_1^*, v_2^*, \dots, v_k^*)$ be the optimal solution, in an arbitrary order
- f be a monotone submodular function. Let $f \geq 0$ (*Update on 04/25/2019: I thought this was w.l.o.g., but Andrey Kolobov pointed out that we actually need f to be non negative*)
- $OPT = f(A^*)$, the value of the optimal solution.

We will prove that

$$f(A_k) \geq (1 - 1/e)OPT$$

Source: <https://homes.cs.washington.edu/~marcotcr/blog/greedy-submodular/>

Proof (Nemhauser et al 1978)

For all $i \leq k$, we have:

$$\begin{aligned} f(A^*) &\leq f(A^* \cup A_i) && \text{Monotonicity} \\ &= f(A_i) + \sum_{j=1}^k \Delta(v_j^* | A_i \cup \{v_1^*, v_2^*, \dots, v_{j-1}^*\}) \\ &\leq f(A_i) + \sum_{z \in A^*} \Delta(z | A_i) && \text{Using submodularity} \\ &\leq f(A_i) + \sum_{z \in A^*} \Delta(v_{i+1} | A_i) && v_{i+1} = \operatorname{argmax}_{v \in V \setminus A_i} \Delta(v | A_i) \\ &= f(A_i) + k \Delta(v_{i+1} | A_i) \end{aligned}$$

Rearranging the terms, we have proved that

$$\Delta(v_{i+1} | A_i) \geq \frac{1}{k} (OPT - f(A_i))$$

Source: <https://homes.cs.washington.edu/~marcotcr/blog/greedy-submodular/>

Proof (Nemhauser et al 1978)

Part I

Now we define $\delta_i = OPT - f(A_i)$. This implies
 $\delta_i - \delta_{i+1} = f(A_{i+1}) - f(A_i) = \Delta(v_{i+1}|A_i)$

Plugging this into our previous equation, we have:

$$\Rightarrow \delta_i - \delta_{i+1} \geq \frac{1}{k}(\delta_i)$$

$$\Rightarrow \delta_{i+1} \leq \left(1 - \frac{1}{k}\right)\delta_i$$

Part II

$$\Rightarrow \delta_k \leq \left(1 - \frac{1}{k}\right)^k \delta_0$$

$$\Rightarrow \delta_k \leq \left(1 - \frac{1}{k}\right)^k OPT \leq \frac{1}{e}OPT$$

$$\Rightarrow OPT - f(A_k) \leq \frac{1}{e}OPT$$

$$\Rightarrow f(A_k) \geq \left(1 - \frac{1}{e}\right)OPT$$

Monotone Submodular – Budget Constraints

$$\max F(S) \text{ s.t. } \sum_{e \in S} c(e) \leq B$$

1. run greedy: S_{gr}

2. run a modified greedy: S_{mod}

$$e^* = \arg \max \frac{F(S_i \cup \{e\}) - F(S_i)}{c(e)}$$

3. pick better of S_{gr} , S_{mod}

→ approximation factor:

$$\frac{1}{2} \left(1 - \frac{1}{e} \right)$$

even better but less fast:
partial enumeration
(Sviridenko, 2004) or
filtering (Badanidiyuru &
Vondrák 2014)

Sviridenko 2004:

- Run the cost-sensitive greedy algorithm starting with all possible initial sets $\{l, j, k\}$
- $O(n^3)$ initial complexity
- $1 - 1/e$ approximation!

Sviridenko 2004, Leskovec et al 2007

Summary: Greedy Algorithm Framework

Monotone Submodular Function

$$\max_{S \subseteq V, c(S) \leq \mathcal{B}} f(S)$$

Cost of Summary Subset S (e.g. size)

Problem Formulation

Initialization $S \leftarrow \emptyset$.

repeat

Pick an element $v^* \in \operatorname{argmax}_{v \in V \setminus S} \frac{f(v \cup S) - f(S)}{c(v)}$

Update $S \leftarrow S \cup v^*$

until Reaching the budget, i.e., $c(S) > \mathcal{B}$

Greedy Algorithm

Non-Monotone Submodular Functions

$$\max_S F(S) \text{ s.t. } |S| \leq k$$

Start with $Y_0 = \emptyset$

for $i = 1$ *to* k **do**

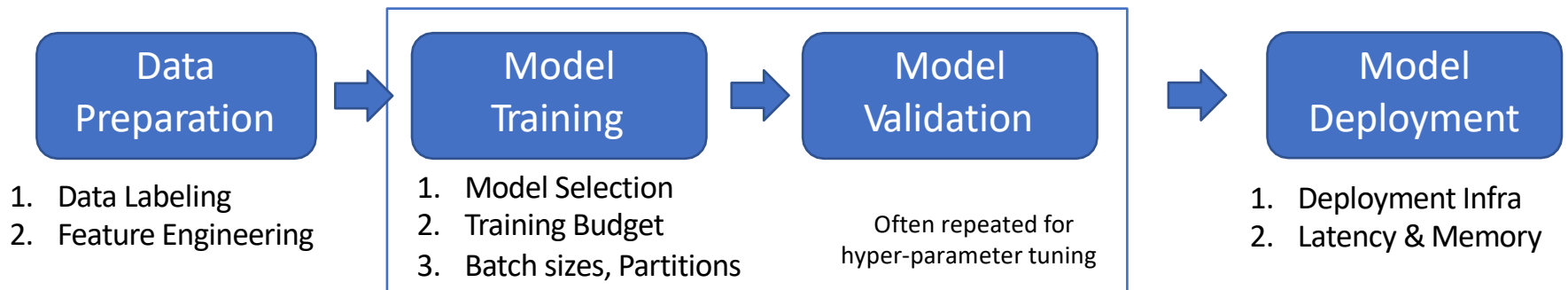
 Let $M_i = \operatorname{argmax}_{X \subseteq V \setminus Y_{i-1}, |X|=k} \sum_{v \in X} f(v|Y_{i-1});$
 Choose y as a uniformly random element in $M_i;$
 $Y_i = Y_{i-1} \cup y;$

return $Y_k.$

Theorem (Buchbinder et al 2014): The Randomized Greedy Algorithm achieves a $1/e$ approximation guarantee for Non-Monotone Submodular Maximization subject to cardinality constraints!

Data subset selection

Make ML Data Efficient and Robust



Production Systems Constraints

1. Data Labeling => Time Consuming, Expensive, Noisy
2. Feature Selection => Latency & Memory
3. Model Training => Compute Intensive and Time Consuming
4. Hyper-Parameter Tuning/NAS => Very Time Consuming
5. Distribution Shift => Deployment vs Training

Can we train Models under these constraints without sacrificing on accuracy?

Data Subset Selection Setup

A Machine Learning model characterized by model parameters θ

Training Data: $\{(x_i, y_i), i \in \mathcal{U}\}$ Training log-likelihood function: $LL_T(\theta, \mathcal{U})$

Training a machine learning model often reduces to finding the parameters that maximizes a log-likelihood function for given training data empirically.

$$\theta^* = \operatorname{argmax}_{\theta} LL_T(\theta, \mathcal{U})$$

Validation Data: $\{(x_i, y_i), i \in \mathcal{V}\}$ Validation log-likelihood function: $LL_V(\theta, \mathcal{V})$

Goal: Select a subset $S \subseteq \mathcal{U}$ such that the resulting model performs the **best!**

Requirements for optimal subset selection

1. The subset selection algorithm needs to be as fast as possible.
 - Subset Selection time \llll Full training time

Example: Subset selection algorithm with negligible time complexity

Training on **10 %** Subset  **10x** Faster training

2. Theoretical guarantees of subset selection algorithm.
 - Can we show theoretical guarantees for subset selection algorithms?

Approaches for Data Subset Selection

- ❑ Several different kinds of approaches studied in literature:
 - ❑ Approach 1: Use Submodular Functions as proxy functions for data subset selection
 - ❑ **Approach 2: Choose data subset which approximates the gradient of the entire dataset**
 - ❑ Approach 3: Choose data subset which approximates the performance on full training dataset (or validation set) as a bi-level optimization!
 - ❑ Approach 4: Choose data subset which minimizes a suitable divergence (e.g. KL divergence) between the distribution induced by the subset and full data!
- ❑ Types of Data Selection
 - ❑ Supervised (Using the labels)
 - ❑ Unsupervised (No access to labels)
 - ❑ Validation based (Access to a validation set for focusing on generalization)

Idea: Gradient Matching/ CoreSets

Can we obtain a weighted gradient of a **subset** of points that approximates the full gradient?

$$\sum_{i \in X_t} w_i^t \nabla_{\theta} L_T^i(\theta) \approx \nabla_{\theta} L(\theta)$$

Gradient Matching/ CoreSets Convergence

Denote L_V as the validation loss, L_T as the full training loss, and L_T^i as the training loss on the i^{th} training example. Furthermore, assume that both losses have gradients bounded by σ_T and σ_V respectively, and that the parameters satisfy $\|\theta^*\|^2 \leq R^2$ (θ^* is the optimal parameter). Then letting L denote either the training or validation loss (with gradient bounded by σ), any data selection algorithm, defined via weights \mathbf{w}^t and subsets X_t for $t = 1, \dots, T$, and run with a learning rate $\alpha = \frac{R}{\sigma_T \sqrt{T}}$ satisfies:

$$\min_{t=1:T} L(\theta_t) - L(\theta^*) \leq \frac{R\sigma_T}{\sqrt{T}} + \frac{R}{T} \sum_{t=1}^{T-1} \text{Err}(\mathbf{w}^t, X_t, L, L_T, \theta_t)$$

where:

$$\text{Err}(\mathbf{w}^t, X_t, L, L_T, \theta_t) = \left\| \sum_{i \in X_t} w_i^t \nabla_{\theta} L_T^i(\theta_t) - \nabla_{\theta} L(\theta_t) \right\|$$

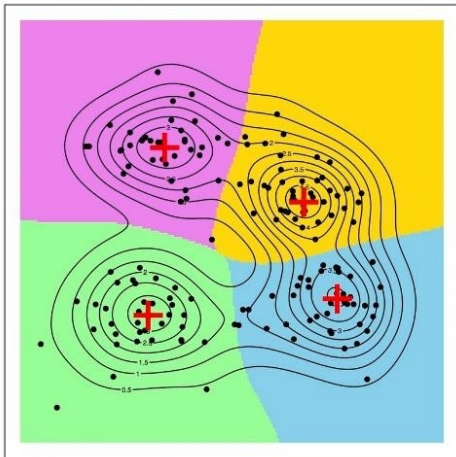
Gradient Matching: Main Idea

The theorem indicates that an effective data selection algorithm should try to have a low error $\text{Err}(\mathbf{w}^t, X_t, L, L_T, \theta_t)$ for $t = 1, \dots, T$. Thus, we can pose the problem as,

$$\begin{aligned}\mathbf{w}^t, X_t &= \min_{\mathbf{w}, X: |X| \leq k} \text{Err}(\mathbf{w}, X, L, L_T, \theta_t) \\ &= \min_{\mathbf{w}, X: |X| \leq k} \left\| \sum_{i \in X_t} w_i^t \nabla_{\theta} L_T^i(\theta_t) - \nabla_{\theta} L(\theta_t) \right\|\end{aligned}$$

CRAIG as an upper bound

Facility location can be thought as an upper bound for this. Suppose we define $\pi_t^i \in \arg \min_{j \in X} \|\nabla_{\theta} L^i(\theta) - \nabla_{\theta} L_T^j(\theta)\|$, if $L = L_T$, then $W = U$ otherwise if $L = L_V$ then $W = V$ and $w_j = \sum_{i \in W} [\pi_t^i = j]$ then, for any θ_t we can write



$$\begin{aligned} \text{Err}(\mathbf{w}, X, L, L_T, \theta_t) &= \left\| \nabla L(\theta_t) - \sum_{i \in X} w_i \nabla L_T^i(\theta_t) \right\| \\ &= \left\| \sum_{i \in W} (\nabla_{\theta} L_T^i(\theta_t) - \nabla_{\theta} L^{\pi_t^i}(\theta_t)) \right\| \\ &\leq \sum_{i \in W} \|\nabla_{\theta} L_T^i(\theta_t) - \nabla_{\theta} L^{\pi_t^i}(\theta_t)\|, \end{aligned}$$

Directly Optimizing Gradient Error: GradMatch

Define the regularized version of our objective:

$$E_\lambda(X) = \min_{\mathbf{w}} \underbrace{\left\| \sum_{i \in X_t} w_t^i \nabla_\theta L_T^i(\theta_t) - \nabla_\theta L(\theta_t) \right\|^2}_{E_\lambda(X_t, \mathbf{w}^t)} + \lambda \|\mathbf{w}^t\|^2$$

This problem can be solved efficiently using Orthogonal Matching Pursuit (OMP) described as,

1. Find projection of $r = \nabla_\theta L(\theta_t)$ for each $i \in W$ along $\nabla_\theta L_T^i(\theta_t)$ and chose the i with whom projection is maximum and add it X
2. Solve linear regression problem to find w_t^i for $i \in X$ s.
3. Set $r = \nabla_\theta L(\theta_t) - \sum_{i \in X_t} w_t^i \nabla_\theta L_T^i(\theta_t)$
4. Repeat the steps with new r until the $|r| < \epsilon$ or $|X| < k$ (budget)
5. Return X, w_t

Orthogonal Matching Pursuit

The OMP algorithm

Algorithm 1: OMP(\mathbf{A} , \mathbf{b})

Input: \mathbf{A} , \mathbf{b}

Result: \mathbf{x}_k

- 1 **Initialization** $\mathbf{r}_0 = \mathbf{b}$, $\Lambda_0 = \emptyset$;
 - 2 Normalize all columns of \mathbf{A} to unit L_2 norm;
 - 3 Remove duplicated columns in \mathbf{A} ;
 - 4 **for** $k = 1, 2, \dots$ **do**
 - 5 Step-1. $\lambda_k = \operatorname{argmax}_{j \notin \Lambda_{k-1}} |\langle \mathbf{a}_j, \mathbf{r}_{k-1} \rangle|$;
 - 6 Step-2. $\Lambda_k = \Lambda_{k-1} \cup \{\lambda_k\}$;
 - 7 Step-3. $\mathbf{x}_k(i \in \Lambda_k) = \operatorname{argmin}_{\mathbf{x}} \|\mathbf{A}_{\Lambda_k} \mathbf{x} - \mathbf{b}\|_2$, $\mathbf{x}_k(i \notin \Lambda_k) = 0$;
 - 8 Step-4. $\hat{\mathbf{b}}_k = \mathbf{A} \mathbf{x}_k$;
 - 9 Step-5. $\mathbf{r}_k \leftarrow \mathbf{b} - \hat{\mathbf{b}}_k$;
 - 0 **end**
-

Convex DSS

Aim

- We study the problem of data efficient training of autonomous driving systems.
- Training using many frames on straight road sections may not be necessary. Frames at the turns turn out to be useful.



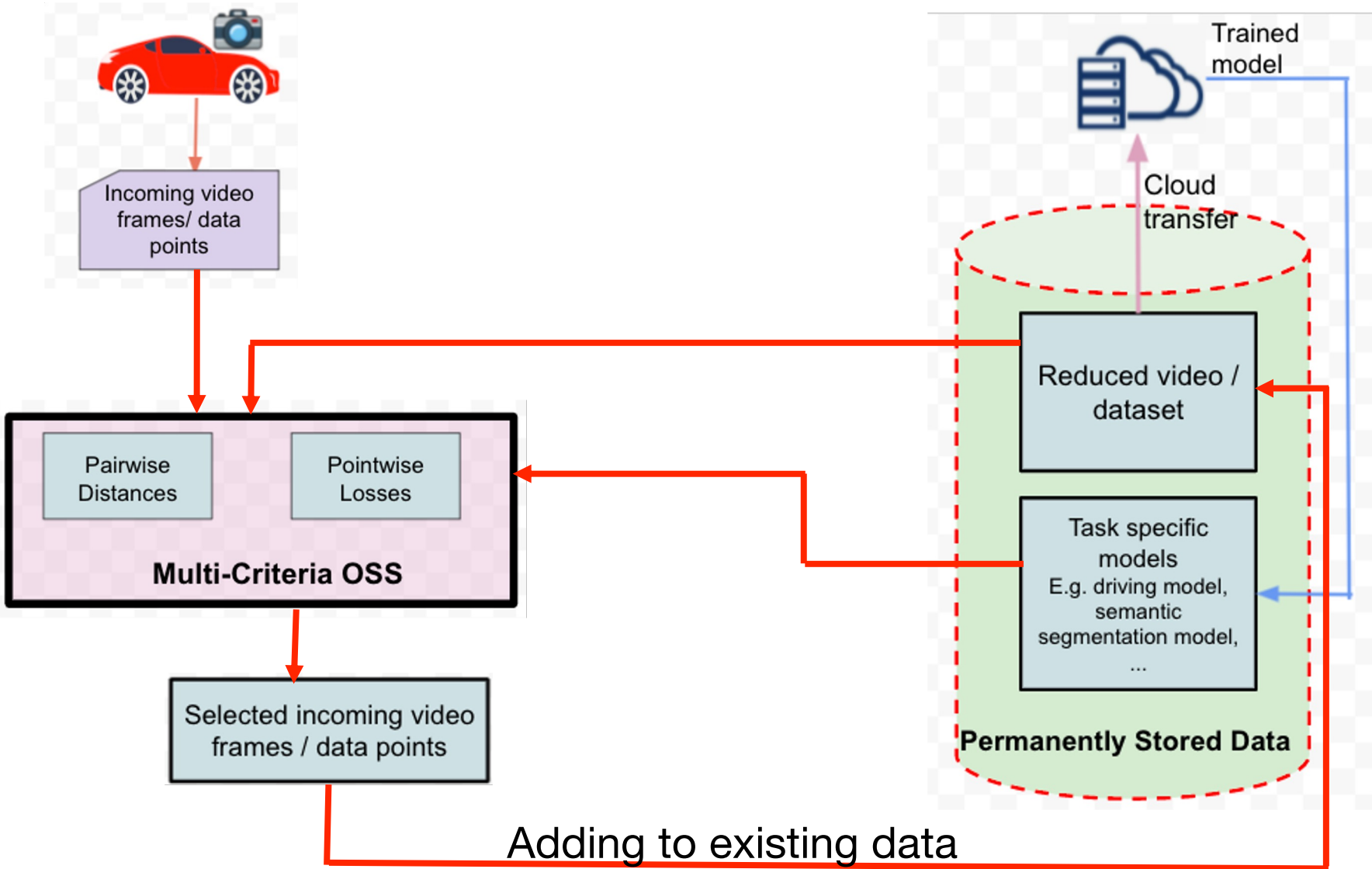
REDUNDANT

INFORMATIVE

| Method | Train One-Turn | Test One-Turn |
|--------------|----------------|---------------|
| Uniform Skip | 3/10 | 5/10 |

In the context of edge device deployment, multi-criteria online subset selection (OSS) framework can be useful in selecting informative frames, essential for an end-task.

Subset selection on Edge devices



High Level Idea

- Given a compression ratio, find out representatives which have the least dissimilarity with the left-out elements besides having the highest task-specific loss.

Problem Setup

- X_t : the set of incoming datapoints at time t (Size m)
- D : set of all data points (Size N)
- R_t : Reduced set of data at time t
- d_{ij} : Distance between data points i and j .
- z_{ij} : Indicator variable indicating that datapoint i is a representative for datapoint j .

Convex Subset Selection

- Original formulation in set notation:

$$\min_{\mathcal{S} \subseteq \mathcal{D}} \lambda |\mathcal{S}| + \sum_{j \in \mathcal{D}} \min_{i \in \mathcal{S}} d_{ij},$$

- Formulation using indicator random variables z_{ij} :

$$\min_{\{z_{ij}\}} \lambda \sum_{i \in \mathcal{D}} \mathbb{I}(\| [z_{i1} \ z_{i2} \ \cdots] \|_p) + \sum_{j \in \mathcal{D}} \sum_{i \in \mathcal{D}} d_{ij} z_{ij}$$

Size regularizer

$$\text{s. t. } z_{ij} \in \{0, 1\}, \quad \sum_{i=1}^N z_{ij} = 1, \quad \forall i, j \in \mathcal{D}.$$

- Convex relaxation:

$$0 \leq z_{ij} \leq 1$$

Online Subset Selection

- At time t :

R_{t-1} : old set (denoted by superscript o)

X_t : in the new set (denoted by superscript n)

R_t : the new reduced set that we are trying to compute using z_{ij}

$$R_t = R_{t-1} \cup \{i \in X_t | Z_{ij} = 1\}$$

- Revised formulation:

$$J'_{\text{enc}} \triangleq \sum_{i \in \mathcal{E}_o} \sum_{j \in \mathcal{D}_n} d_{ij}^{o,n} z_{ij}^{o,n} + \sum_{i \in \mathcal{D}_n} \sum_{j \in \mathcal{D}_n} d_{ij}^{n,n} z_{ij}^{n,n},$$

$$\min_{Z'} J'_{\text{enc}} + \lambda \sum_{i \in \mathcal{D}_n} I(\| [z_{i1}^{n,n} \ z_{i2}^{n,n} \ \dots] \|_p)$$

$$\text{s. t. } z_{ij}^{o,n}, z_{ij}^{n,n} \in \{0, 1\}, \quad \forall i, j,$$

$$\sum_{i \in \mathcal{E}_o} z_{ij}^{o,n} + \sum_{i \in \mathcal{D}_n} z_{ij}^{n,n} = 1, \quad \forall j \in \mathcal{D}_n,$$

$$\begin{aligned} \mathcal{E}_o &= R_{t-1} \\ \mathcal{D}_n &= X_t \end{aligned}$$

High Level Idea

- Given a compression ratio, find out representatives which have the least dissimilarity with the left-out elements besides having the highest task-specific loss.
- Highest task-specific loss ensures having situational tasks needed to be learnt more by the model.

TMCOSS

Adopts a facility location objective involving multiple criteria

$$\min_{z_{ij}^o, z_{ij}^n} \mathcal{G}(z_{ij}^o, z_{ij}^n) \text{ s.t. } \sum_{j=1}^{|R_t|} z_{i,j}^o + \sum_{j=1}^m z_{i,j}^n = 1; z_{i,j}^n, z_{i,j}^o \in [0,1]; \sum_{j=1}^m \|[z_{1,j}^n \dots z_{m,j}^n]\|_p \leq \text{frac} * m$$

Objective function

Constraint 1

Constraint 2

Compression Ratio

$z_{ij}^o = 1$ Denotes j from existing set o is a representative of element i from incoming set n

$z_{ij}^n = 1$ Denotes j from incoming set n is a representative of element i from incoming set n

$$\mathcal{G}(z_{ij}^o, z_{ij}^n) = \rho \left(\sum_{i=1}^m \sum_{j=1}^{|R_t|} z_{ij}^o d_{ij}^o(t) + \sum_{i,j=1}^m z_{ij}^n d_{ij}^n(t) \right) - (1 - \rho) \left(\sum_{j=1}^{|R_t|} S_j^o * L_j^o + \sum_{j=1}^m S_j^n * L_j^n \right) \text{ where, } S_j^o = \frac{1}{\epsilon} \min(\epsilon, \sum_{i=1}^m z_{ij}^o), S_j^n = \frac{1}{\epsilon} \min(\epsilon, \sum_{i=1}^m z_{ij}^n)$$

Dissimilarity

Task specific Loss


Representative power of element j thresholded by ϵ

Justification for thresholding

Theorem 1 Let z_{ij}^o and z_{ij}^n be the optimal solution for formulation 1. A new frame $j \in X_{t+1}$ is selected as a representative frame for at least one incoming frame $i \in X_{t+1}$, i.e. $z_{ij}^n = 1$, only if BOTH these conditions hold:

- For some incoming frame $i \in X_{t+1}$, $Q_{ij}^n < Q_{ij}^o$, for all $j' \in X_{t+1}$ and $j' \neq j$
- For some incoming frame $i \in X_{t+1}$, $Q_{ij}^n < \frac{\sum_{i'=1}^m z_{i',k}^o Q_{i',k}^o + \lambda \| [z_{1,j}^n \dots z_{m,j}^n] \|_p}{\| \mathbf{z}_j^n \|_1}$

where $k = \operatorname{argmin}_j \sum_{i=1}^m z_{i,j}^o Q_{i,j}^o$, and $\| \mathbf{z}_j^n \|_1 = \sum_{i'=1}^m z_{i',j}^n$

$$\rho = 0$$


Corollary 1.1 Let z_{ij}^o and z_{ij}^n be the optimal solution for formulation 1. A new frame $j \in X_{t+1}$ is selected as a representative frame for at least one incoming frame $i \in X_{t+1}$, i.e. $z_{ij}^n = 1$, only if BOTH these conditions hold:

- $L_j^n > L_{j'}^n$ for all $j' \in X_{t+1}$ and $j' \neq j$
- $L_j^n > \frac{\sum_{i=1}^m z_{i,k}^o L_k^o - \lambda \| [z_{1,j}^n \dots z_{m,j}^n] \|_p}{\| \mathbf{z}_j^n \|_1}$

where $k = \operatorname{argmin}_j \sum_{i=1}^m z_{i,j}^o Q_{i,j}^o$, and $\| \mathbf{z}_j^n \|_1 = \sum_{i'=1}^m z_{i',j}^n$

Multi-criteria OSS (MCOSS)¹

$$Q_{ij}^n = \rho d_{ij}^n - (1 - \rho) L_j^n; Q_{ij}^o = \rho d_{ij}^o - (1 - \rho) L_j^o$$

$$\min_{z_{ij}^o, z_{ij}^n} \sum_{i=1}^m \sum_{j=1}^{|R_t|} z_{ij}^o Q_{ij}^o + \sum_{i,j=1}^m z_{ij}^n Q_{ij}^n + \lambda \sum_{j=1}^m \| [z_{1,j}^n \dots z_{m,j}^n] \|_p$$

$$s . t . \sum_{j=1}^{|R_t|} z_{i,j}^o + \sum_{j=1}^m z_{i,j}^n = 1, \forall i \in X_{t+1} z_{i,j}^n, z_{i,j}^o \in [0,1], \forall i, j$$

¹. Soumi Das, Sayan Mondal, Ashwin Bhojar, Madhumita Bharde, Niloy Ganguly, Suparna Bhattacharya, Sourangshu Bhattacharya, "Multi-criteria onlineframe-subset selection for autonomous vehicle videos." *Pattern Recognition Letters* 133 (2020): 349-355.

References

Submodular Optimization

1. Andreas Krause and Daniel Golovin. Submodular Function Maximization. 2012.
2. IJCAI 2021 tutorial by Rishabh Iyer and Ganesh Ramakrishnan
3. Buchbinder, Niv, et al. Submodular maximization with cardinality constraints. SODA 2014.
4. Mirzasoleiman, Baharan, Jeff Bilmes, and Jure Leskovec. Coresets for data-efficient training of machine learning models. ICML 2020.

Orthogonal Matching Pursuit

1. Killamsetty, Krishnateja, S. Durga, Ganesh Ramakrishnan, Abir De, and Rishabh Iyer. Grad-match: Gradient matching based data subset selection for efficient deep model training. ICML 2021.
2. Tropp, Joel A., and Anna C. Gilbert. Signal recovery from random measurements via orthogonal matching pursuit. IEEE Transactions on information theory, (2007).
3. Cai, T. Tony, and Lie Wang. "Orthogonal matching pursuit for sparse signal recovery with noise." IEEE Transactions on Information theory, (2011).

Convex Optimization

1. Elhamifar, Ehsan, and M. Clara De Paolis Kaluza. Online summarization via submodular and convex optimization. CVPR 2017.
2. Das, Soumi, Harikrishna Patibandla, Suparna Bhattacharya, Kshounis Bera, Niloy Ganguly, and Sourangshu Bhattacharya. "TMCOS: Thresholded Multi-Criteria Online Subset Selection for Data-Efficient Autonomous Driving." CVPR 2021.

