



## Paper Reading Seminar

Learning Probabilistic Submodular Diversity Models  
Via Noise Contrastive Estimation

Proceedings of the 19th International Conference on Artificial  
Intelligence and Statistics (AISTATS), 2016

---

HOANG NT - M2

---

2016/10/21

---

# Our Roadmap

3 main topics and 2 off-track topics

**Introduction**  
“Diversity”  
“Summary”

**DPPs**  
Model  
Problems

**FLID**  
Model  
Results Discussion



**Set Theory**  
Quick Review  
Notations



**NCE**  
Z Estimation  
My favorite!

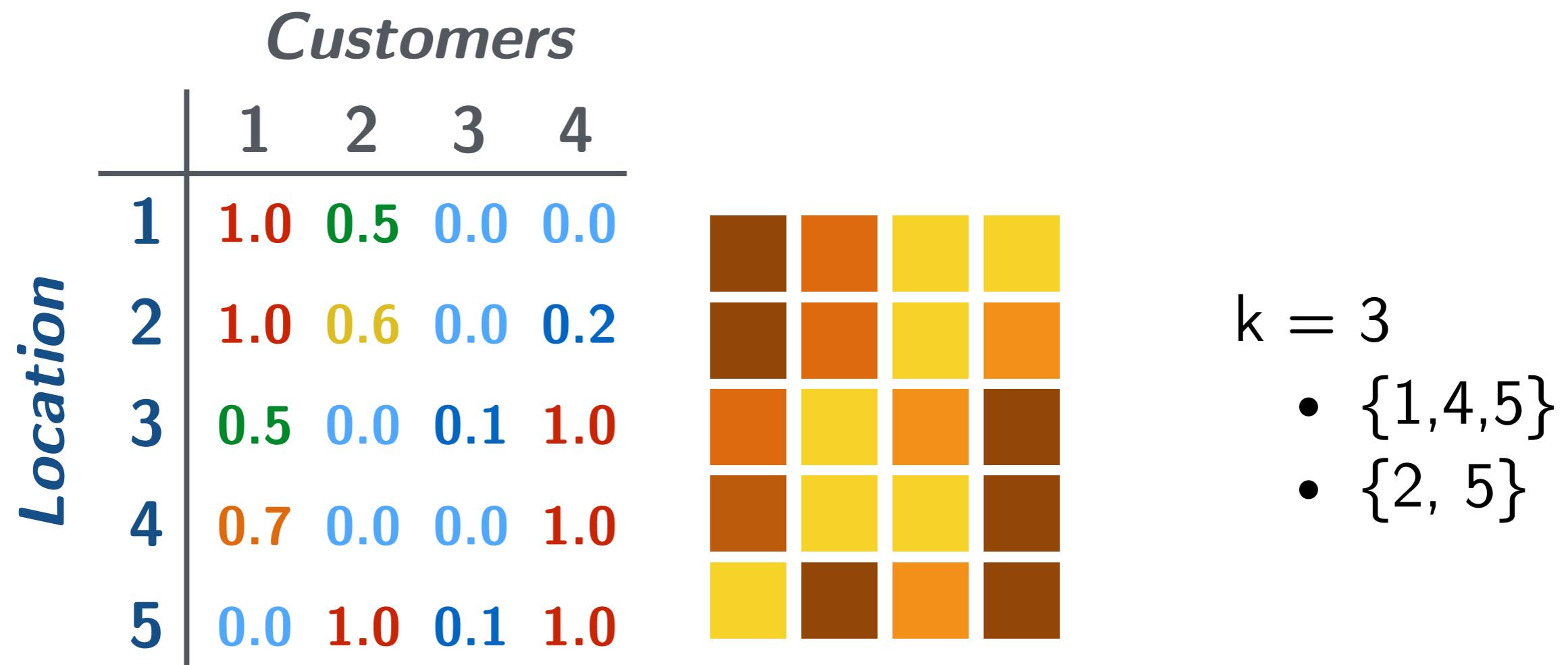
# Diversification and its application

## “Bad” search result example



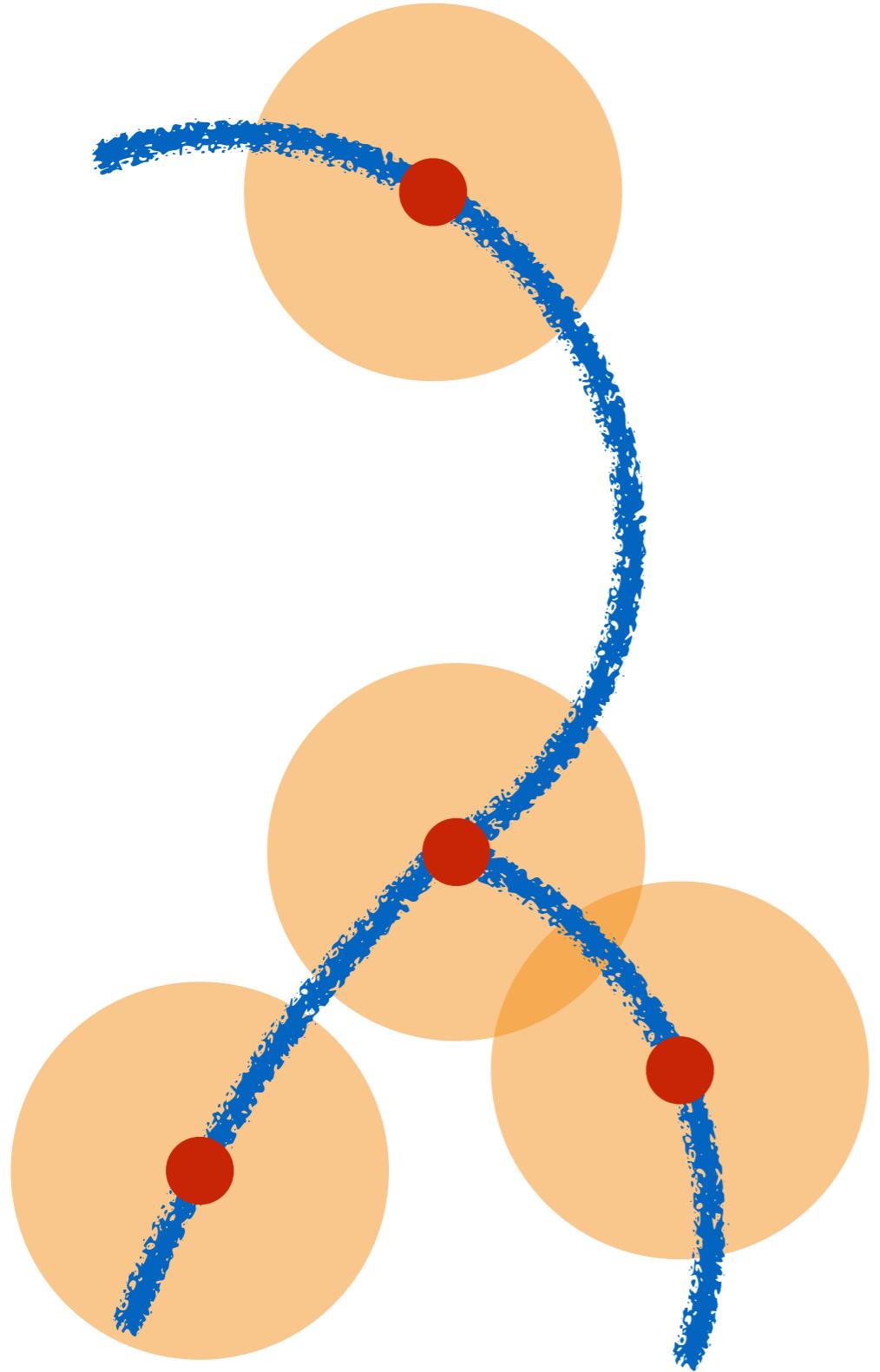
Car





## Diversification and its application

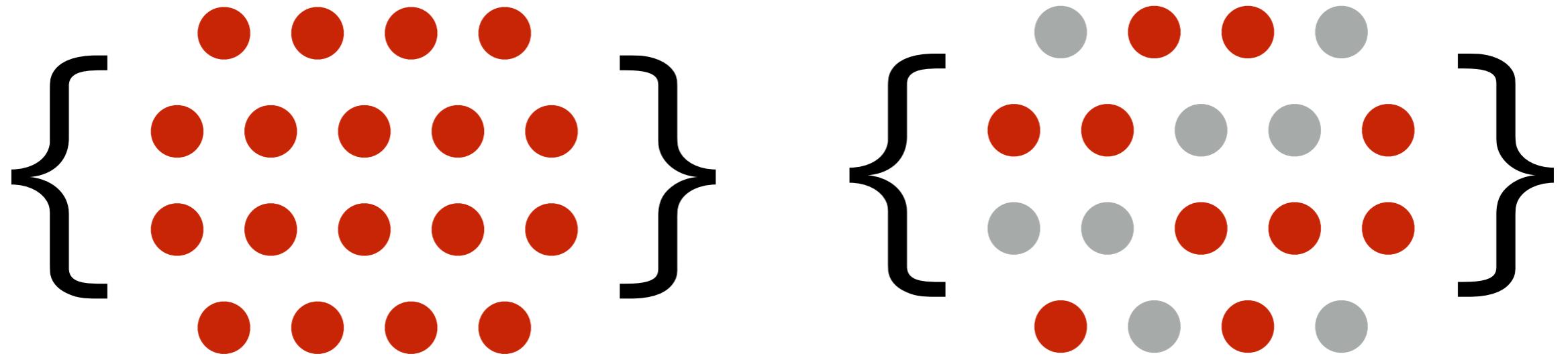
The nature of coverage and submodular functions



$$S = \{ \textcolor{red}{\bullet} \textcolor{red}{\bullet} \textcolor{red}{\bullet} \textcolor{red}{\bullet} \}$$

$$f(S) = \langle \text{coverage} \rangle$$

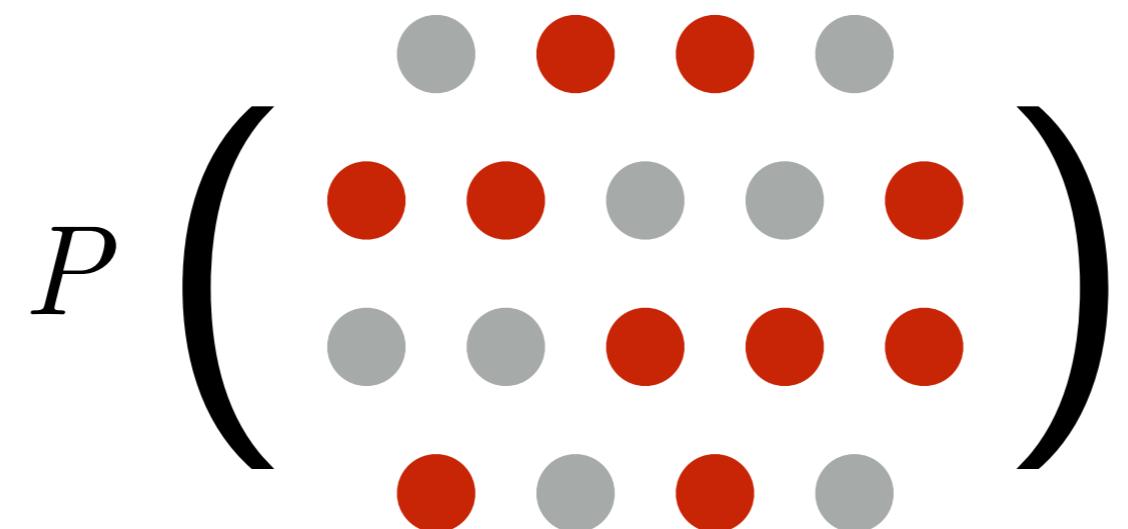
Set S is a collection of “things”.



Ground set (informally) is where we pick elements from.  
(*conf. Matroids and Partially Ordered Sets for more.*)

$2^V$  is the power set, or set of subsets of ground set V.

Given a set  $S$ , a (simple) point process is a random subset  $A$  of  $S$ . (*Terry Tao, Determinantal processes, 2009*)



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

$$F(A \cup \{i\}) - F(A) \geq F(B \cup \{i\}) - F(B),$$

$$A \subseteq B \subseteq V \setminus \{i\}.$$

	1	2	3	4
1	Dark Brown	Orange	Yellow	Yellow
2	Dark Brown	Orange	Yellow	Orange
3	Orange	Yellow	Orange	Dark Brown
4	Dark Brown	Yellow	Yellow	Dark Brown
5	Yellow	Dark Brown	Orange	Dark Brown

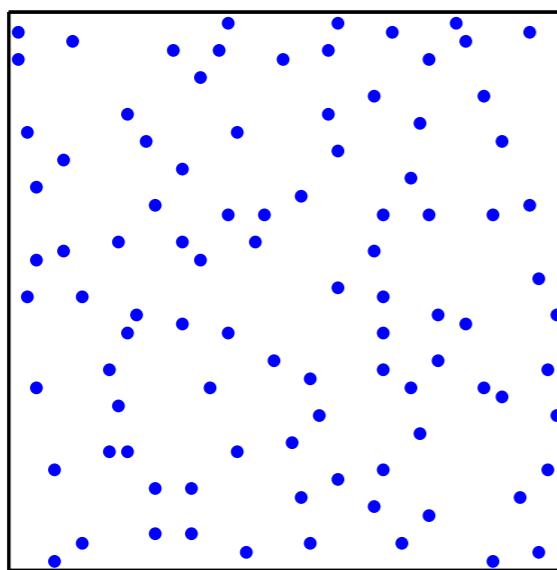
$$F(S) = [ \max(S_{i,d}) \text{ for } S_i \text{ in } S ]$$

$$F(\{2 \mid \{1\}\}) \quad \text{vs.} \quad F(\{2 \mid \{1,3\}\})$$

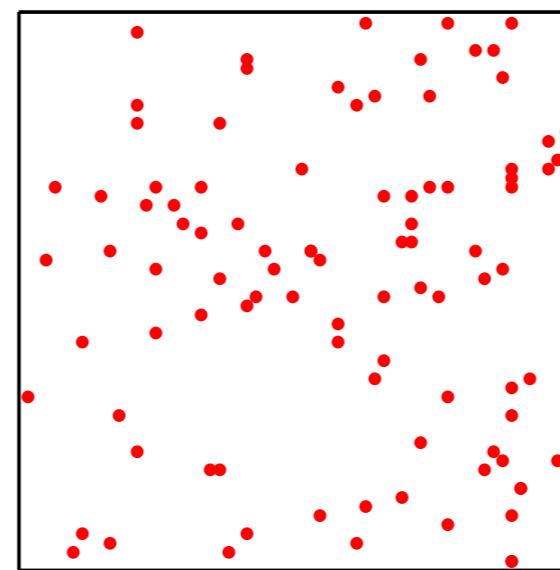
# Determinantal Point Processes

## Solution to diversity sampling problem

$$P \left( \begin{array}{c} \text{red dots} \\ \vdots \\ \text{grey dots} \end{array} \right) = \det(K_A)$$



DPP



Independent

(Kulesza, UAI'11)

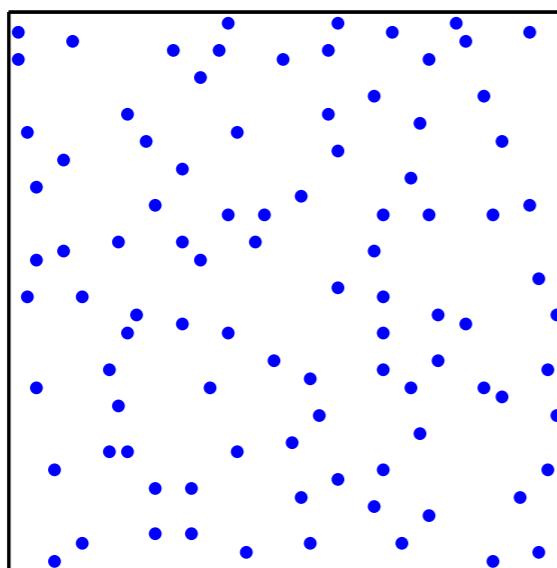
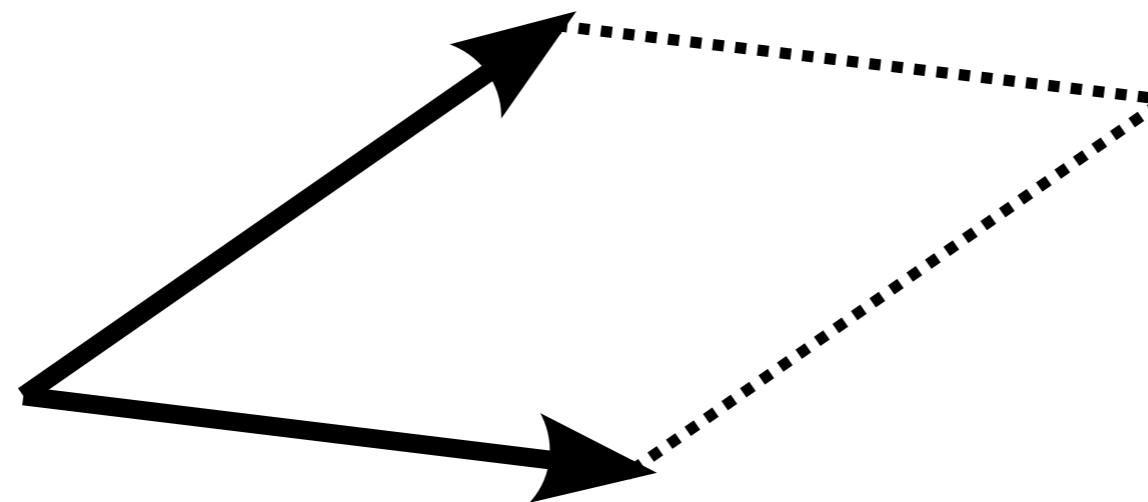
For data modeling, L-ensembles has advantages.

(Borodin, 2009)

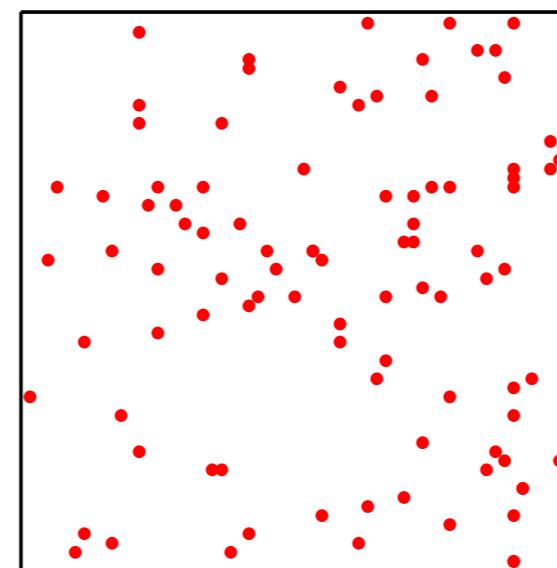
$$\mathcal{P}_L(\mathbf{Y} = Y) = \frac{\det(L_Y)}{\det(L + I)}$$

$$K = (L + I)^{-1} L$$

$$L_{ij} = q_i \phi_i^\top \phi_j q_j$$



DPP



Independent

(Kulesza, UAI'11)

# Noise Contrastive Estimation

## Fast approximation for partition function

(Guntmann, AISTATS'10)

$$\int p_m(\mathbf{u}; \hat{\alpha}) d\mathbf{u} = 1.$$

$$p_m(\cdot; \alpha) = \frac{p_m^0(\cdot; \alpha)}{Z(\alpha)}, \quad Z(\alpha) = \int p_m^0(\mathbf{u}; \alpha) d\mathbf{u},$$

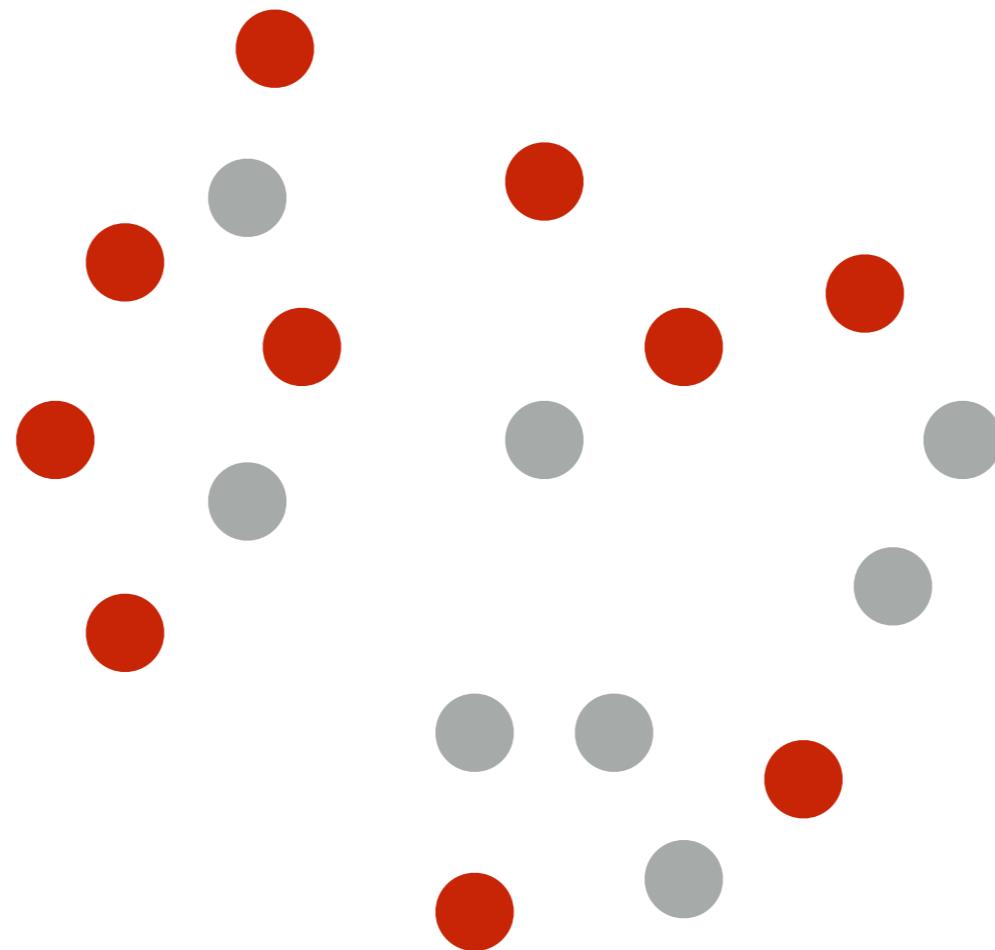
$$\ln p_m(\cdot; \theta) = \ln p_m^0(\cdot; \alpha) + c$$

$$\begin{aligned} P(C = 1 | \mathbf{u}; \theta) &= \frac{p_m(\mathbf{u}; \theta)}{p_m(\mathbf{u}; \theta) + p_n(\mathbf{u})} \\ &= h(\mathbf{u}; \theta) \end{aligned}$$

$$P(C = 0 | \mathbf{u}; \theta) = 1 - h(\mathbf{u}; \theta).$$

# Noise Contrastive Estimation

Classification between “positive” and “negative” samples



(Tschötschek, AISTATS'16)

Log-probabilistic model:

$$P \left( \begin{array}{c} \text{red} \quad \text{red} \quad \text{grey} \quad \text{grey} \quad \text{red} \\ \text{grey} \quad \text{grey} \quad \text{red} \quad \text{red} \quad \text{red} \\ \text{red} \quad \text{grey} \quad \text{red} \end{array} \right) \propto \exp(\text{something})$$

$$P(S | \mathbf{u}, \mathbf{W})$$

$$= \frac{1}{Z} \exp \left( \underbrace{\sum_{i \in S} u_i + \sum_{d=1}^L \left( \max_{i \in S} w_{i,d} - \sum_{i \in S} w_{i,d} \right)}_{\text{Div}(S)} \right),$$
$$\underbrace{\qquad\qquad\qquad}_{\widetilde{P}(S|\mathbf{u}, \mathbf{W})}$$

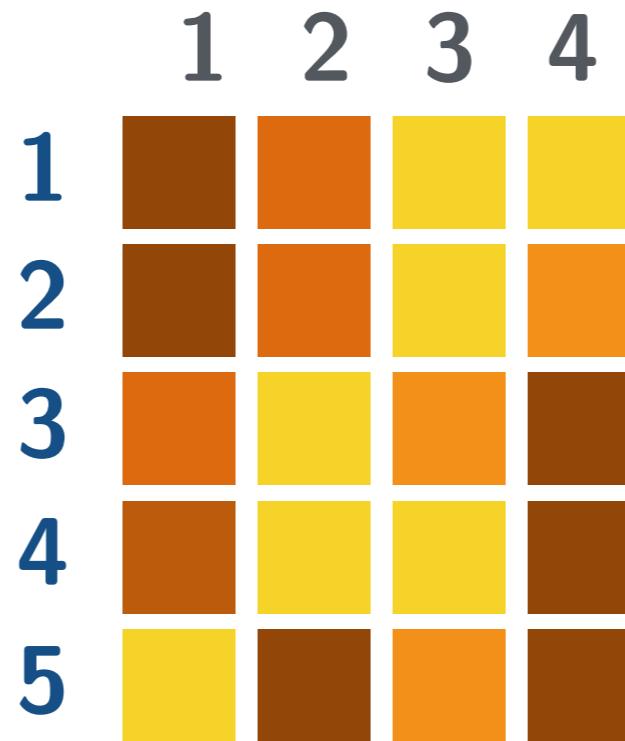
..

# Facility Location Diversity

## Modular and submodular formulation

Modular

$$P(S) \propto \exp \left( \sum_{i \in S} u_i \right)$$



$$\text{Div}(\{1,2\}) = -1.5$$

$$\text{Div}(\{2,4\}) = -0.9$$

Submodular

$$P(S | \mathbf{u}, \mathbf{W})$$

$$= \frac{1}{Z} \exp \left( \underbrace{\sum_{i \in S} u_i + \sum_{d=1}^L \left( \max_{i \in S} w_{i,d} - \sum_{i \in S} w_{i,d} \right)}_{\text{Div}(S)} \right),$$

*index-d of a L-dim vector*

$$\underbrace{\qquad\qquad\qquad}_{\widetilde{P}(S|\mathbf{u},\mathbf{W})}$$

# Facility Location Diversity

Computational complexity for partition function

Straight-forward way to compute  $Z$  is  $\mathcal{O}(|V|^{L+1})$

Therefore, NCE comes to the rescue!

$$\begin{aligned} P(Y_S = 1 \mid S) &= \frac{P(S \mid Y_S = 1)}{P(S \mid Y_S = 1) + \eta P(S \mid Y_S = 0)} \\ &= \frac{P_d(S)}{P_d(S) + \eta P_n(S)}, \end{aligned}$$

$$g(\theta) = \sum_{S \in \mathcal{D}} \log P(Y_S = 1 \mid S, \theta) + \sum_{S \in \mathcal{N}} \log P(Y_S = 0 \mid S, \theta)$$

$$\theta = [\mathbf{u}_{\text{NCE}}, \mathbf{W}_{\text{NCE}}, \hat{Z}]$$

$$\begin{aligned}\nabla \log P(Y_S = Y \mid S) \\ = \left( Y - \frac{1}{1 + \eta \frac{P_n(S)}{\frac{1}{\hat{Z}} \tilde{P}(S \mid \mathbf{u}, \mathbf{W})}} \right) \nabla \log \frac{1}{\hat{Z}} \tilde{P}(S \mid \mathbf{u}, \mathbf{W})\end{aligned}$$

$$\begin{aligned}\nabla \log \left( \frac{Y P_d(S) + (1 - Y) \eta P_n(S)}{P_d(S) + \eta P_n(S)} \right) &= \frac{(2Y - 1) \eta P_n(S)}{(P_d(S) + \eta P_n(S))(Y P_d(S) + (1 - Y) \eta P_n(S))} \nabla P_d(S) \\ &= \frac{(2Y - 1) \eta P_n(S)}{(P_d(S) + \eta P_n(S))(Y + (1 - Y) \eta \frac{P_n(S)}{P_d(S)})} \frac{1}{P_d(S)} \nabla P_d(S)\end{aligned}$$

\*  $Y$  only takes values of 0 or 1



- 29,632 baby registries from [amazon.com](https://www.amazon.com)
- 18 categories
- Filtered out listing with fewer than 5 or more than 100 products.
- Create sub-registries by categories.  
Further filtered out items that did not occur in at least 100 sub-registries.
- Result: 13 categories (avg. 71 products each), 8585 sub-registries. 70% for fitting model. 30% for testing.

Graco Sweet Slumber Sound Machine



Boppy Noggin Nest Head Support



Braun ThermoScan Lens Filters



Aquatopia Bath Thermometer Alarm



Model fitting:

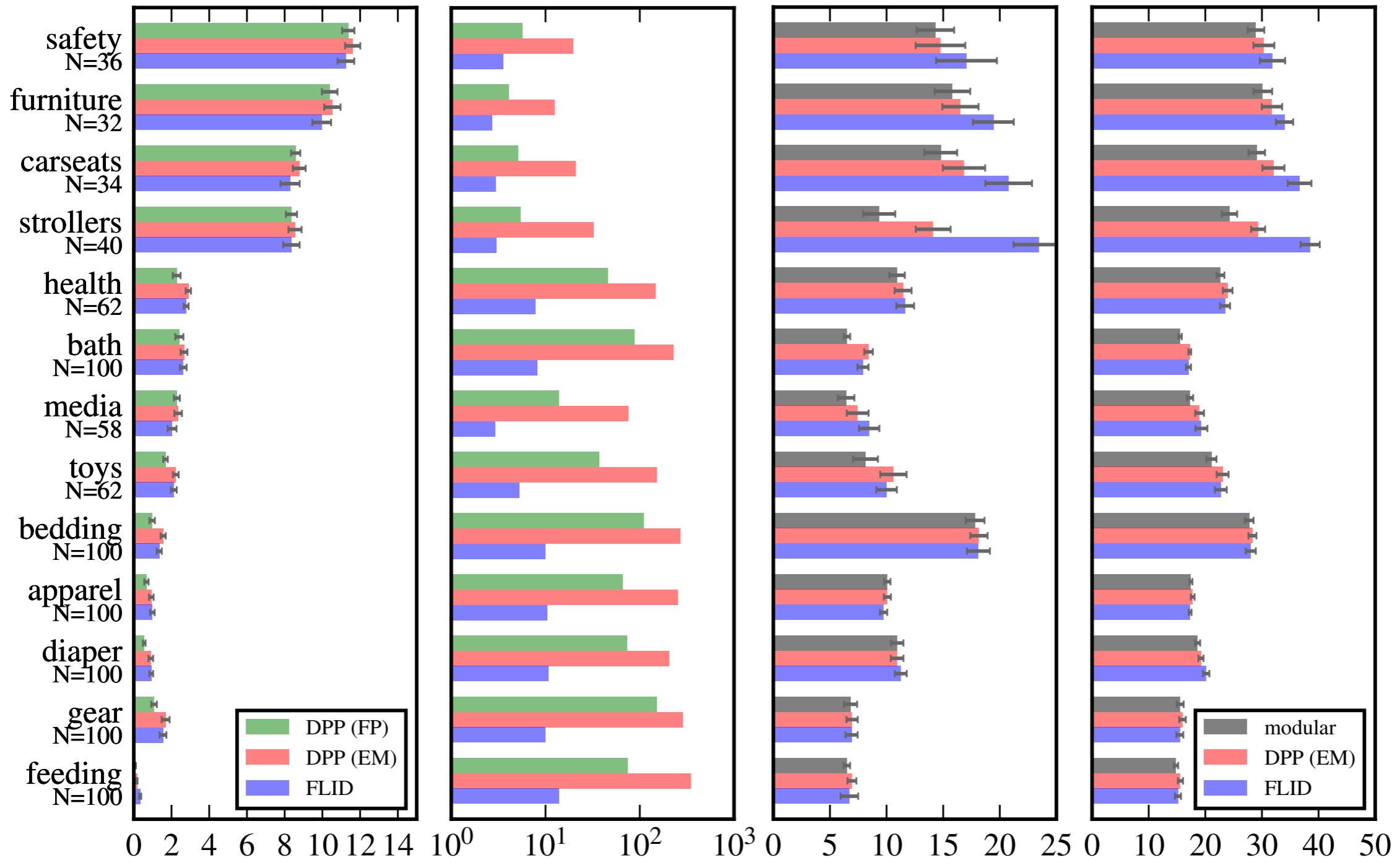
$$\text{LLRI} = 100 \cdot \frac{\mathcal{L}_{\text{method}} - \mathcal{L}_{\text{modular}}}{|\mathcal{L}_{\text{modular}}|},$$

Mean reciprocal rank (MRR)

$$\text{MRR} = \frac{100}{|\mathcal{T}'|} \sum_{\check{S}_i \in \mathcal{T}'} \frac{1}{\text{rank}_i^{\check{S}_i}}.$$

# Facility Location Diversity

## Result for Amazon baby registry



(a) Log-likelihood relative improvement  
(b) Cumulative runtime in seconds  
(c) Accuracy  
(d) MRR

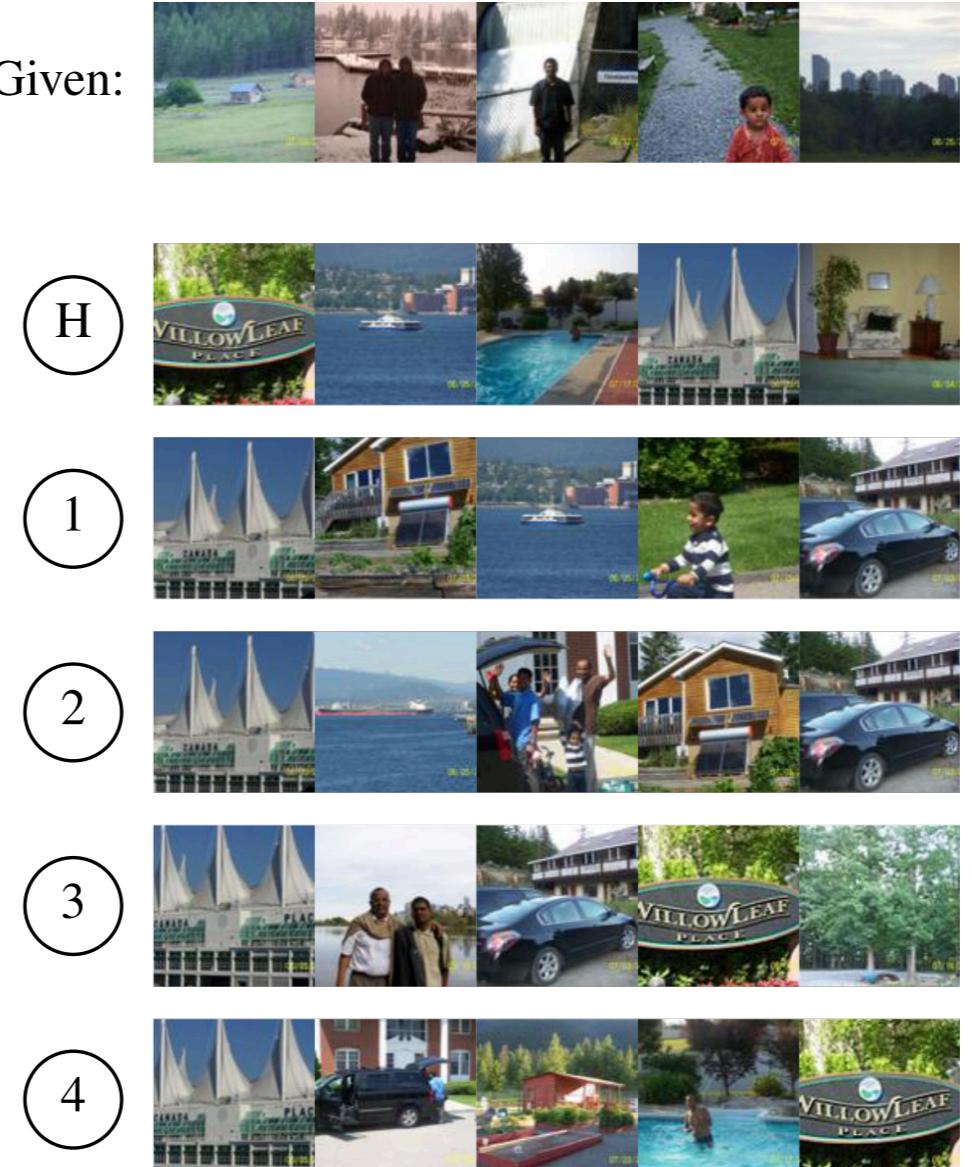
# Facility Location Diversity

## Result for Image summarization



(a) Image collection

Given:

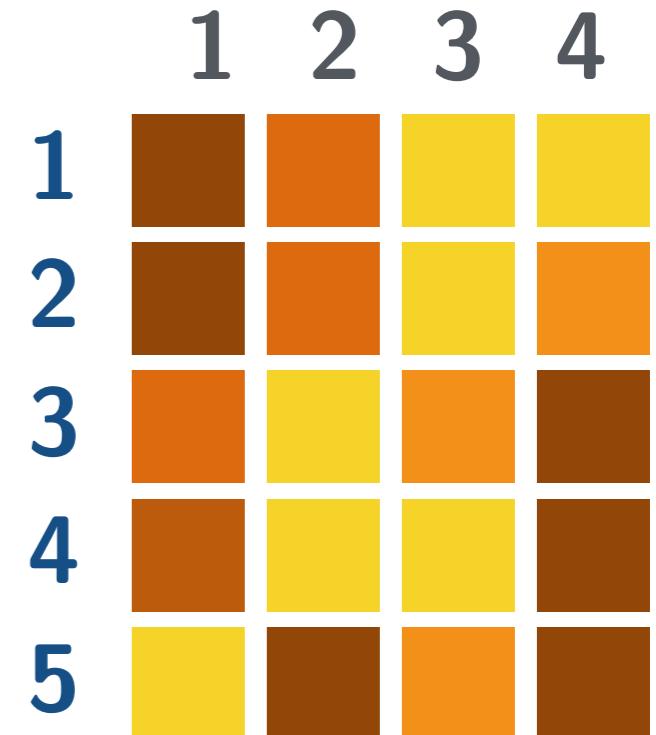


(b) Completed summaries

# Facility Location Diversity

## Discussion

$$\begin{aligned}
 P(S | \mathbf{u}, \mathbf{W}) &= \frac{1}{Z} \exp \left( \sum_{i \in S} u_i + \underbrace{\sum_{d=1}^L \left( \max_{i \in S} w_{i,d} - \sum_{i \in S} w_{i,d} \right)}_{\text{Div}(S)} \right), \\
 &\quad \underbrace{\qquad\qquad\qquad}_{\tilde{P}(S|\mathbf{u},\mathbf{W})}
 \end{aligned}$$



$$g(\theta) = \sum_{S \in \mathcal{D}} \log P(Y_S = 1 | S, \theta) + \sum_{S \in \mathcal{N}} \log P(Y_S = 0 | S, \theta)$$

The importance / sensitivity of the parameter eta.  
 Transitivity of FLID compared to DPPs.  
 Other application / usable dataset for FLID.

(Kulesza, UAI'11) *Learning Determinantal Point Processes.*

(Guntmann, AISTATS'10) *Noise-contrastive estimation: A new estimation principle for unnormalized statistical models.*

(Tschischek, AISTATS'16) *Learning Probabilistic Submodular Diversity Models Via Noise-Contrastive Estimation.*

(Terry Tao, 2009) *Determinantal Point Processes.*