

# Spatial Summarization of Image Collections

Diego A. Ballesteros Villamizar

ETH Zürich

January 27th, 2016

## 1 Featurized FLID

## 2 Extending the Model

- 1 In previous presentation, latitude and longitude were normalized to full range, i.e.  $[-90, 90]$  and  $[-180, 180]$  respectively.
- 2 All features normalized to  $[0, 1]$  range **over the data**.
- 3 However, no improvement of score.
- 4 Moreover, the phenomena previously seen on the  $bmW$  weights is still present, i.e. they are the same across dimensions  $d$ .

# Why uniform weights across dimensions?

- ① *Note 1:* Initialization of the  $\mathbf{B}$  weights is obtained from a uniform distribution over  $[0, 0.001]$ .
- ② *Note 2:* The gradient update for  $\mathbf{B}$  is given by equations 1 and 2:

$$\left( \nabla_{\mathbf{B}} \log \frac{1}{\hat{Z}} \tilde{P}(S \mid \mathbf{a}, \mathbf{B}) \right)_{m,l} = x_{i^*,m} - \sum_{i \in S} x_{i,m} \quad (1)$$

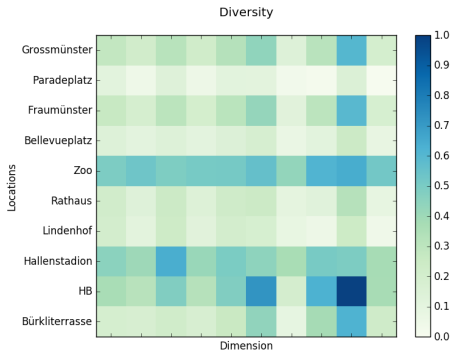
$$i^* = \operatorname{argmax}_{i \in S} \mathbf{x}_i \mathbf{b}_l \quad (2)$$

Which shows that the only difference in the weights across dimensions is given by the value of  $i^*$ .

- ③ Because the initialization is small in comparison with the feature values and the variability is low, the value of  $i^*$  at every timestep is equal in all dimensions  $d$  as it only depends on the feature vector  $\mathbf{x}_i$ .

# Fixing weight $B$ distribution

- 1 To avoid initial updates choosing always the same  $i^*$ , the initialization for  $B$  should be larger. The new range is  $[0, 1]$ .
- 2 The learned weight matrix  $W$  is not uniform with the updated initialization condition.



- 1 However, the scoring remains low. Worse than the score with the identity as feature matrix, i.e. the non-featurized model.

Model	Accuracy	MRR
Modular with features	$17.38 \pm 1.81$	$39.85 \pm 1.52$
<b>FLID (<math>d = 10</math>)</b>	<b><math>29.31 \pm 2.74</math></b>	<b><math>52.19 \pm 1.78</math></b>
FFLID ( $d = 1$ )	$13.60 \pm 1.60$	$37.54 \pm 1.39$
FFLID ( $d = 5$ )	$13.48 \pm 1.61$	$37.51 \pm 1.45$
FFLID ( $d = 10$ )	$12.34 \pm 1.51$	$36.94 \pm 1.47$

1 Featurized FLID

2 Extending the Model

$$P(S \mid \mathbf{a}, \mathbf{B}, \mathbf{C}) = \frac{1}{Z} \exp \left( \sum_{i \in S} u_i + \text{Div}(S, \mathbf{B}) + \text{Coh}(S, \mathbf{C}) \right) \quad (3)$$

$$\text{Div}(S, \mathbf{B}) = \sum_{l=1}^L \left( \max_{i \in S} \mathbf{x}_i \mathbf{b}_l - \sum_{i \in S} \mathbf{x}_i \mathbf{b}_l \right) \quad (4)$$

$$\text{Coh}(S, \mathbf{C}) = \sum_{k=1}^K \left( \sum_{i \in S} \mathbf{x}_i \mathbf{c}_k - \max_{i \in S} \mathbf{x}_i \mathbf{c}_k \right) \quad (5)$$

$$\mathbf{u} = \mathbf{X} \mathbf{a} \quad \mathbf{X} \in \mathbb{R}^{|V| \times M} \quad \mathbf{u} \in \mathbb{R}^{|V|} \quad \mathbf{a} \in \mathbb{R}^M \quad (6)$$

$$\mathbf{W}_B = \mathbf{X} \mathbf{B} \quad \mathbf{B} \in \mathbb{R}^{|M| \times L} \quad (7)$$

$$\mathbf{W}_C = \mathbf{X} \mathbf{C} \quad \mathbf{C} \in \mathbb{R}^{|M| \times K} \quad (8)$$

$$(9)$$



- 1 Supermodular term encourages adding similar items, i.e. values with high  $w_{c_{i,k}}$  close to the  $\max_{i \in S}$ .
- 2 This is a more natural notion in the setting of touristic places, e.g. someone who visits Fraumünster will likely visit Grossmünster as well.
- 3 The learning is analog to the diversity-only case, however more parameters are present so more noise samples should be used. The gradient update is the negative of the equation for the diversity case.

Table: Frequency of Item Sets

Set	Frequency
[0, 2]	101
[9, 3]	66
[0, 5]	63
[2, 5]	62
[5, 6]	52
[9, 1]	50
[2, 6]	49
[9, 2]	44
[0, 2, 5]	43
[0, 3]	41

Table: Locations

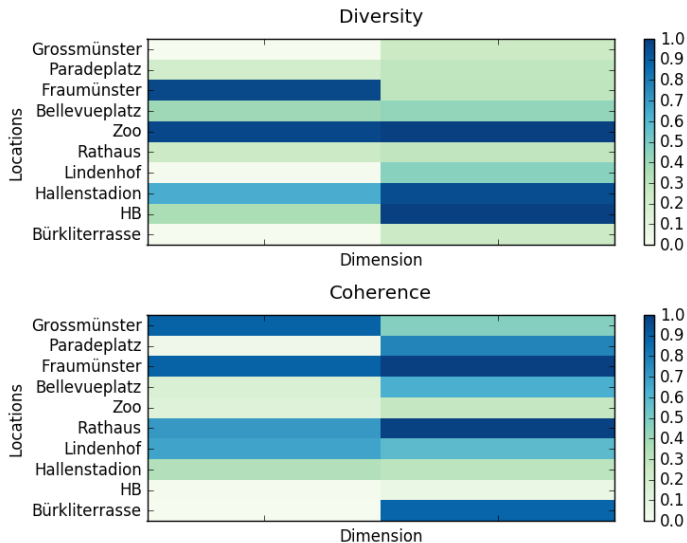
Index	Location
0	Grossmünster
1	Paradeplatz
2	Fraumünster
3	Bellevue
4	Zoo
5	Rathaus
6	Lindenhof
7	Hallenstadion
8	HB
9	Bürkliplatz

# Model without Features

$$\mathbf{X} = \mathbb{I}$$

		K			
		0	2	5	10
L	0	$18.15 \pm 3.08$	$24.85 \pm 7.37$	$31.06 \pm 6.89$	$30.72 \pm 3.22$
	2	$22.39 \pm 2.66$	<b><math>34.07 \pm 2.63</math></b>		
	5	$24.04 \pm 3.22$		$34.35 \pm 2.15$	
	10	$29.31 \pm 2.74$			$35.48 \pm 2.16$

# Model without Features

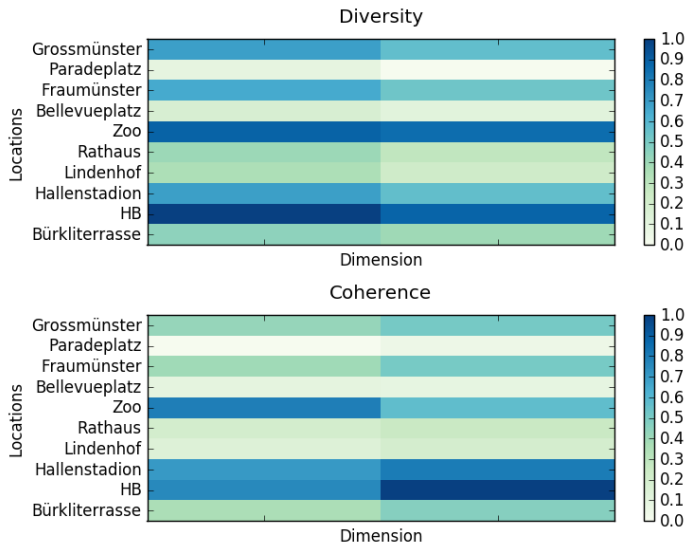


# Model with Features

$$\mathbf{X} \in \mathbb{R}^{10 \times 4}$$

		K			
		0	2	5	10
L	0	$17.38 \pm 1.81$	$16.69 \pm 1.67$	$16.77 \pm 2.13$	$17.90 \pm 1.65$
	2	$13.39 \pm 1.58$	$13.63 \pm 1.64$		
	5	$13.48 \pm 1.61$		$13.21 \pm 1.44$	
	10	$12.34 \pm 1.51$			$12.96 \pm 1.14$

# Model with Features



# Experiment Performance

- ① As the number of parameters increase, what is the cost on running time performance?
- ② As in Tschitschek et al. (2016), let's define  $\kappa = \max_{S \in \mathcal{D} \cup \mathcal{N}} |S|$ .
- ③ The gradient update operations per iteration are:
  - ① Updating the utility vector  $a$ :  $O(\kappa M)$ .
  - ② Updating the weights  $B$ :  $O(\kappa M L)$ .
  - ③ Updating the weights  $C$ :  $O(\kappa M K)$ .
- ④ Then the overall performance is:  $O(|\mathcal{D} \cup \mathcal{N}| \kappa M (L + K))$ .
- ⑤ This is similar to the performance reported on Tschitschek et al. (2016):  $O(|\mathcal{D} \cup \mathcal{N}| \kappa D)$ , as long as  $L + K \approx D$  and  $M \ll \kappa$ .

Tschiatschek, S., Djolonga, J., and Krause, A. (2016). Learning probabilistic submodular diversity models via noise contrastive estimation. In *Proc. International Conference on Artificial Intelligence and Statistics (AISTATS)*.