
Large-Margin Determinantal Point Processes

Wei-Lun Chao*
 U. of Southern California
 Los Angeles, CA 90089
 weilunc@usc.edu

Boqing Gong*
 U. of Southern California
 Los Angeles, CA 90089
 boqinggo@usc.edu

Kristen Grauman
 U. of Texas at Austin
 Austin, TX 78701
 grauman@cs.utexas.edu

Fei Sha
 U. of Southern California
 Los Angeles, CA 90089
 feisha@usc.edu

In this Supplementary Material, we provide extra details on the following:

- Sec. A: deriving the softmax, eq. (13) in the main text. We also show how to efficiently compute the marginal probability $P_{\{i\}}$ in eq. (13).
- Sec. B: subgradients of the objective function of our large-margin DPP (cf. eq. (14) in the main text).
- Sec. C: extra details on generating oracle video summaries and evaluating summarization results against user summaries.

A CALCULATING THE SOFTMAX

In the main text, we use softmax to deal with the exponential number of large margin constraints and arrive at eq. (13) in the main text. Here we show how to calculate the right-hand side of eq. (13).

First, we compute $\sum_{\mathbf{y} \subseteq \mathcal{Y}_n} \ell_\omega(\mathbf{y}_n, \mathbf{y}) P(\mathbf{y}; \mathbf{L}_n)$ as follows

$$\begin{aligned} & \sum_{\mathbf{y} \subseteq \mathcal{Y}_n} \ell_\omega(\mathbf{y}_n, \mathbf{y}) P(\mathbf{y}; \mathbf{L}_n) \\ &= \sum_{\mathbf{y} \subseteq \mathcal{Y}_n} \left[\sum_{i: i \notin \mathbf{y}} \mathbb{I}(i \notin \mathbf{y}_n) + \omega \sum_{i: i \in \mathbf{y}} \mathbb{I}(i \in \mathbf{y}_n) \right] P(\mathbf{y}; \mathbf{L}_n) \end{aligned} \quad (17)$$

$$\begin{aligned} &= \sum_{i=1}^M \left[\sum_{\mathbf{y}: i \in \mathbf{y}} \mathbb{I}(i \notin \mathbf{y}_n) P(\mathbf{y}; \mathbf{L}_n) \right. \\ & \quad \left. + \omega \sum_{\mathbf{y}: i \notin \mathbf{y}} \mathbb{I}(i \in \mathbf{y}_n) P(\mathbf{y}; \mathbf{L}_n) \right] \end{aligned} \quad (18)$$

$$= \sum_{i=1}^M [\mathbb{I}(i \notin \mathbf{y}_n) P_{n_{\{i\}}} + \omega \mathbb{I}(i \in \mathbf{y}_n) (1 - P_{n_{\{i\}}})] \quad (19)$$

$$= \sum_{i: i \notin \mathbf{y}_n} P_{n_{\{i\}}} + \omega \sum_{i: i \in \mathbf{y}_n} (1 - P_{n_{\{i\}}}) \quad (20)$$

*Equal contribution

$$= \sum_{i: i \notin \mathbf{y}_n} K_{n_{ii}} + \omega \sum_{i: i \in \mathbf{y}_n} (1 - K_{n_{ii}}), \quad (21)$$

where $P_{n_{\{i\}}} = K_{n_{ii}}$ is the marginal probability of selecting item i . Now we are ready to see

$$\begin{aligned} & \text{softmax}_{\mathbf{y} \subseteq \mathcal{Y}_n} \log \ell_\omega(\mathbf{y}_n, \mathbf{y}) + \log P(\mathbf{y}; \mathbf{L}_n) \\ &= \log \sum_{\mathbf{y} \subseteq \mathcal{Y}_n} \ell_\omega(\mathbf{y}_n, \mathbf{y}) P(\mathbf{y}; \mathbf{L}_n) \end{aligned} \quad (22)$$

$$= \log \left(\sum_{i: i \notin \mathbf{y}_n} K_{n_{ii}} + \omega \sum_{i: i \in \mathbf{y}_n} (1 - K_{n_{ii}}) \right). \quad (23)$$

Moreover, recall that $\mathbf{K} = \mathbf{L}(\mathbf{L} + \mathbf{I})^{-1}$. Eigendecomposing $\mathbf{L} = \sum_m \lambda_m \mathbf{v}_m \mathbf{v}_m^T$, we have

$$\mathbf{K} = \mathbf{L}(\mathbf{L} + \mathbf{I})^{-1} = \sum_m \frac{\lambda_m}{\lambda_m + 1} \mathbf{v}_m \mathbf{v}_m^T,$$

$$\text{and thus, } K_{ii} = \sum_m \frac{\lambda_m}{\lambda_m + 1} v_{mi}^2. \quad (24)$$

B SUBGRADIENTS OF THE OBJECTIVE FUNCTION

Recall that our objective function in eq. (14) of the main text actually consists of a likelihood term $\mathcal{L}(\cdot)$ and the other term of undesirable subsets. Denote them respectively by

$$\begin{aligned} \mathcal{L}(\boldsymbol{\theta}, \boldsymbol{\alpha}; \mathcal{Y}_n, \mathbf{y}_n) &\triangleq \log P(\mathbf{y}_n; \mathbf{L}_n) \\ &= \log \det(\mathbf{L}_{n_{\mathbf{y}_n}}) - \log \det(\mathbf{L}_n + \mathbf{I}), \end{aligned} \quad (25)$$

$$\mathcal{A}(\boldsymbol{\theta}, \boldsymbol{\alpha}; \mathcal{Y}_n, \mathbf{y}_n) \triangleq \log \left(\sum_{i \notin \mathbf{y}_n} K_{n_{ii}} + \omega \sum_{i \in \mathbf{y}_n} (1 - K_{n_{ii}}) \right). \quad (26)$$

For brevity, we drop the subscript n of \mathbf{L}_n and $K_{n_{ii}}$ and change \mathbf{y}_n to \mathbf{y}^* in what follows.

To compute the overall subgradients, it is sufficient to compute the gradients of the above two terms, \mathcal{L} and \mathcal{A} . Denoting by $\Theta = \{\theta, \alpha, \beta\}$, we have

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial \Theta_k} &= \sum_{i,j} \frac{\partial \mathcal{L}}{\partial L_{ij}} \frac{\partial L_{ij}}{\partial \Theta_k} = \mathbf{1}^T \left(\frac{\partial \mathcal{L}}{\partial \mathbf{L}} \circ \frac{\partial \mathbf{L}}{\partial \Theta_k} \right) \mathbf{1}, \\ \frac{\partial \mathcal{A}}{\partial \Theta_k} &= \mathbf{1}^T \left(\frac{\partial \mathcal{A}}{\partial \mathbf{L}} \circ \frac{\partial \mathbf{L}}{\partial \Theta_k} \right) \mathbf{1},\end{aligned}\quad (27)$$

where \circ stands for the element-wise product between two matrices of the same size. We use the chain rule to decompose $\frac{\partial \mathcal{L}}{\partial \Theta_k}$ from the overall gradients on purpose. Therefore, if we change the way of parameterizing the DPP kernel \mathbf{L} , we only need care about $\frac{\partial \mathcal{L}}{\partial \Theta_k}$ when we compute the gradients for the new parameterization.

B.1 GRADIENTS OF THE QUALITY-DIVERSITY DECOMPOSITION

In terms of the quality-diversity decomposition (c.f. eq. (7) and (8) in the main text), we have

$$\begin{aligned}\frac{\partial \mathbf{L}}{\partial \alpha_k} &= (\mathbf{q}\mathbf{q}^T) \circ S^k, \quad \frac{\partial L_{ij}}{\partial \theta_k} = L_{ij}(x_{ik} + x_{jk}), \\ \text{or } \frac{\partial \mathbf{L}}{\partial \theta_k} &= \mathbf{L} \circ (\mathbf{X}\mathbf{e}_k\mathbf{1}^T + \mathbf{1}\mathbf{e}_k^T\mathbf{X}^T)\end{aligned}\quad (28)$$

where \mathbf{q} is the vector concatenating the quality terms q_i , \mathbf{X} is the design matrix concatenating \mathbf{x}_i^T row by row, and \mathbf{e}_k stands for the standard unit vector with 1 at the k -th entry and 0 elsewhere.

B.2 GRADIENTS WITH RESPECT TO THE DPP KERNEL

In what follows we calculate $\frac{\partial \mathcal{L}}{\partial \mathbf{L}}$ and $\frac{\partial \mathcal{A}}{\partial \mathbf{L}}$ in eq. (27). Noting that eq. (27) sums over all the (i, j) pairs, we therefore do not need bother taking special care of the symmetric structure in \mathbf{L} .

We will need map $\mathbf{L}_{\mathbf{y}^*}$ “back” to a matrix \mathbf{M} which is the same size as the original matrix \mathbf{L} , such that $\mathbf{M}_{\mathbf{y}^*} = \mathbf{L}_{\mathbf{y}^*}$ and all the other entries of \mathbf{M} are zeros. We denote by $\langle \mathbf{L}_{\mathbf{y}^*} \rangle$ such mapping, i.e., $\langle \mathbf{L}_{\mathbf{y}^*} \rangle = \mathbf{M}$. Now we are ready to see,

$$\begin{aligned}\frac{\partial \mathcal{L}}{\partial \mathbf{L}} &= \frac{\partial \log \det(\mathbf{L}_{\mathbf{y}^*})}{\partial \mathbf{L}} - \frac{\partial \log \det(\mathbf{L} + \mathbf{I})}{\partial \mathbf{L}} \\ &= \langle (\mathbf{L}_{\mathbf{y}^*})^{-1} \rangle - (\mathbf{L} + \mathbf{I})^{-1}.\end{aligned}\quad (29)$$

It is a little more involved to compute

$$\begin{aligned}\frac{\partial \mathcal{A}}{\partial \mathbf{L}} &= \frac{1}{\sum_{i \notin \mathbf{y}^*} K_{ii} + \omega \sum_{i \in \mathbf{y}^*} (1 - K_{ii})} \\ &\quad \times \left[\sum_{i \notin \mathbf{y}^*} \frac{\partial K_{ii}}{\partial \mathbf{L}} - \omega \sum_{i \in \mathbf{y}^*} \frac{\partial K_{ii}}{\partial \mathbf{L}} \right],\end{aligned}\quad (30)$$

which involves $\frac{\partial K_{ii}}{\partial \mathbf{L}}$.

In order to calculate $\frac{\partial K_{ii}}{\partial \mathbf{L}}$, we start from the basic identity [Beyer, 1991] of

$$\frac{\partial \mathbf{A}^{-1}}{\partial t} = -\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial t} \mathbf{A}^{-1}, \quad (31)$$

followed by $\frac{\partial \mathbf{A}^{-1}}{\partial A_{mn}} = -\mathbf{A}^{-1} \mathbf{J}^{mn} \mathbf{A}^{-1}$, where \mathbf{J}^{mn} is the same size as \mathbf{A} . The (m, n) -th entry of \mathbf{J}^{mn} is 1 and all else are zeros.

Let $\mathbf{A} = (\mathbf{L} + \mathbf{I})$. Noting that $\mathbf{K} = \mathbf{L}(\mathbf{L} + \mathbf{I})^{-1} = \mathbf{I} - (\mathbf{L} + \mathbf{I})^{-1} = \mathbf{I} - \mathbf{A}^{-1}$ and thus $K_{ii} = 1 - [\mathbf{A}^{-1}]_{ii}$, we have,

$$\begin{aligned}\frac{\partial K_{ii}}{\partial L_{mn}} &= -\frac{\partial [\mathbf{A}^{-1}]_{ii}}{\partial L_{mn}} = -\frac{\partial [\mathbf{A}^{-1}]_{ii}}{\partial A_{mn}} \\ &= [\mathbf{A}^{-1} \mathbf{J}^{mn} \mathbf{A}^{-1}]_{ii} = [\mathbf{A}^{-1}]_{mi} [\mathbf{A}^{-1}]_{ni}.\end{aligned}\quad (32)$$

We can also write eq. (32) in the matrix form,

$$\begin{aligned}\frac{\partial K_{ii}}{\partial \mathbf{L}} &= [\mathbf{A}^{-1}]_{\cdot i} [\mathbf{A}^{-1}]_{\cdot i}^T \\ &= \mathbf{A}^{-1} \mathbf{e}_i \mathbf{e}_i^T \mathbf{A}^{-1} = \mathbf{A}^{-1} \mathbf{J}^{ii} \mathbf{A}^{-1},\end{aligned}\quad (33)$$

where $[\mathbf{A}^{-1}]_{\cdot i}$ is the i -th column of \mathbf{A}^{-1} .

Overall, we arrive at a concise form by writing out the right-hand-side of eq. (30) and merging some terms,

$$\begin{aligned}\sum_{i \notin \mathbf{y}^*} \frac{\partial K_{ii}}{\partial \mathbf{L}} - \omega \sum_{i \in \mathbf{y}^*} \frac{\partial K_{ii}}{\partial \mathbf{L}} \\ = \mathbf{A}^{-1} \mathbf{I}_\omega(\overline{\mathbf{y}^*}) \mathbf{A}^{-1} = (\mathbf{L} + \mathbf{I})^{-1} \mathbf{I}_\omega(\overline{\mathbf{y}^*}) (\mathbf{L} + \mathbf{I})^{-1}\end{aligned}\quad (34)$$

where $\mathbf{I}_\omega(\overline{\mathbf{y}^*})$ looks like an identity matrix except that its (i, i) -th entry is $-\omega$ for $i \in \mathbf{y}^*$.

C VIDEO SUMMARIZATION

We provide details on 1) how to generate oracle summaries as the supervised information to learn DPPs and 2) how to evaluate system-generated summaries against user summaries. We also present more results on balancing the precision and recall through our large-margin DPP.

C.1 ORACLE SUMMARY

In the OVP dataset, each video comes along with five user summaries $\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_5$ [de Avila et al., 2011]. Similar to document summarization [Kulesza and Taskar, 2011], we extract an “oracle” summary \mathbf{y}^* from the five user summaries using a greedy algorithm. Initialize $\mathbf{y}^* = \emptyset$. From the frames not in \mathbf{y}^* , we pick out the one i which contributes the most to the marginal gain,

$$\begin{aligned}\text{VSUMM}(\mathbf{y}^* \cup \{i\}, \{\mathbf{y}_1, \dots, \mathbf{y}_5\}) \\ - \text{VSUMM}(\mathbf{y}^*, \{\mathbf{y}_1, \dots, \mathbf{y}_5\}),\end{aligned}\quad (35)$$

where VSUMM is the package developed in [de Avila et al., 2011] to evaluate video summarization results. We postpone to Section C.2 for describing the evaluation scheme of VSUMM. Namely, we select the oracle frames greedily for each video and stop until the marginal gain becomes negative. We evaluate the oracle summaries against users’ and find that they achieve high precision and recalls, 84.1% and 87.7% respectively, validating that the oracle summaries are able to serve as good supervised targets for training DPP models.

The above procedure allows a “user-independent” definition of a good oracle summary for learning. Of course if the application goal were to generate user-specific summaries catering to a particular user’s taste, one would instead simply apply our framework with y^* set to be that particular user’s selection.

C.2 VSUMM: EVALUATING VIDEO SUMMARIZATION RESULTS

We evaluate video summarization results using the VSUMM package [de Avila et al., 2011]. Given two sets of summaries/frames, it searches for the maximum number of matched pairs of frames between them. Two images are viewed as a matched pair if their visual difference is below a certain threshold. VSUMM uses normalized color histograms to compute such difference. Besides, each frame of one set can be matched to at most one frame of the other set, and vice versa. After the matching procedure, one can hence develop different evaluation metrics based on the number of matched pairs. In our experiments, we define F-score, precision, and recall (cf. eq. (15) of the main text).

References

- W. H. Beyer. *CRC standard mathematical tables and formulae*. CRC press, 1991.
- S. E. F. de Avila, A. P. B. Lopes, et al. Vsum: A mechanism designed to produce static video summaries and a novel evaluation method. *Pattern Recognition Letters*, 32(1):56–68, 2011.
- A. Kulesza and B. Taskar. Learning determinantal point processes. In *UAI*, 2011.