Diverse Sequential Subset Selection for Supervised Video Summarization

Boging Gong^{+*}, Wei-Lun Chao^{+*}, Kristen Grauman[‡], and Fei Sha[†] +University of Southern California, +University of Texas at Austin

Highlight

- Pose video summarization as a supervised learning problem for subset selection
- Propose sequential determinantal point process (seqDPP) as the underlying probabilistic model
- Evaluate on three video summarization tasks and obtain state-of-the-art performance

Introduction

Video summarization: pressing need - 100 hours of new Youtube video per min 422,000 CCTV cameras in London 24/7 Summaries by three users









Challenges

- Heterogeneous subjects/categories
- Various temporal changing rates
- Subjective, disparate, and noisy labels

Previous work

- Criteria: representativeness vs. diversity
- Largely unsupervised, frame clustering
- Require sophisticated handcrafting

Our main idea

- Supervised learning from human supplied annotations
- Summarization as subset selection
- Modeling temporal cue & diversity

Approach

Sequential DPP (seqDPP)



- $P(Y_t = \boldsymbol{y}_t | Y_{t-t})$
- Ω_t : kernel over

Parameterization of DPP kernel

- Neural networks (NN)

Inference

 $\boldsymbol{y}_1^* = rg \operatorname{ma}$ $oldsymbol{y}_2^* = rg \, \mathrm{ma}$

Learning via

- through gradient descent

In contrast, bag DPPs: Model permutable items (no temporal info) Often use quality-diversity kernel (limited) Inference NP hard

1. Partition video into T disjoint segments 2. Introduce subset selection (of frames) variable Y_t for each segment 3. Condition Y_t on $Y_{t-1} = y_{t-1}$ by DPP

- Linear embedding (L): $f_i^T W^T W f_j$

$$\begin{aligned} & \operatorname{ax}_{\boldsymbol{y}\in\mathcal{Y}_1} P(Y_1 = \boldsymbol{y}) \\ & \operatorname{ax}_{\boldsymbol{y}\in\mathcal{Y}_2} P(Y_2 = \boldsymbol{y}|Y_1 = \boldsymbol{y}_1^*) \\ & \mathsf{MLE} \end{aligned}$$

Generating target summaries

User study on inter-annotator agreement

- Data: 100 videos from Open Video Project and Youtube
- Annotation: 5 user summaries per video
- **Observation:** high inter-annotator agreement

Generate target summaries by greedy search

 $y^* \leftarrow y^* \cup rg \max$

Experiments

Setup

- Data: OVP (50), Youtube (39), Kod
- Feature: Fisher vector, saliency, co
- Evaluation: Precision, Recall, F-sco
- Comparison: bag DPP and previou (unsupervised) DT, STIMO, VSUMM

Results on Youtube and

	VSUMM2			Ours (L)			Ours (NN)		
	F	Р	R	F	Ρ	R	F	Ρ	R
Youtube	55.7	59.7	58.7	57.8	54.2	69.8	60.3	59.4	64.9
Kodak	68.9	75.7	80.6	75.3	77.8	80.4	78.9	81.9	81.1

Target summary

seqDPP (ours) (F=70, P=60, R=88)

VSUMM1 (F=59, P=65, R=55)









a novel evaluation method". Pattern Recognition Letters, 32(1):56–68, 2011. Machine Learning, 5(2-3):123–286, 2012.



$$\mathbf{x}_i \sum_u F_{\mathbf{y}^* \cup i, \mathbf{y}_u}$$



dak (18)
ontext
ore
JS

	F	Р	
DT	57.6	67.7	53.2
STIMO	63.4	60.3	72.2
VSUMM1	70.3	70.6	75.8
	60 0	70 1	60

Results on OVP

VSUMM2	68.2	73.1	69.1
bag DPP	70.8±0.3	71.5±0.4	74.5±0.3
Ours + Q/D	68.5±0.3	66.9 ±0.4	75.8 ±0.5
Ours (L)	75.5±0.4	77.5 ±0.5	78.4±0.5
Ours (NN)	77.7 ±0.4	75.0 ±0.5	87.2 ±0.3



[1] S. Avila, A. Lopes, A. Luz Jr, A. Araujo. "VSUMM: A mechanism designed to produce static video summaries and [2] A. Kulesza and B. Taskar. "Determinantal point processes for machine learning". Foundations and Trends® in