Inferring the Why in Images [Pirsiavash et al]

CSC2523 Winter 2015: Paper Presentation Micha Livne

Goals



(b)

Sitting because he wants to watch television Sitting because she intends to see the doctor

Related Work

Most FAVORABLE



Predicted Intents

Energetic Dominant Most COMFORTING Energetic Dominant Most COMPETENT Energetic Dominant Most DOMINANT Energetic Dominant Trustworthy

Visual Persuasion: Inferring Communicative Intents of Images [Joo et al 2014]

(d) Example images and predicted intents

Competent

Comforting

Proposed Solution: Vision Only



[Krizhevsky et al 2012]

Visual classifier



Proposed Solution: Full Solution

Language Potentials

Relationship	Query to Language Model		
action + object + motivation	action the object in order to motivation		
	action the object to motivation		
	action the object because pronoun wants to motivation		
action + object + scene	action the object in a scene		
	in a scene, action the object		
action + scene + motivation	action in a scene in order to motivation		
	action in order to motivation in a scene		
	action because pronoun wants to motivation in a scene		

$$\mathcal{V}_{ij}(y_i, y_j)$$

Dataset



Statistics of Motivations

- Based on PASCAL VOC 2012.
- Only images with a person.
- Annotation of: action, object, scene, and motivation (79).

Proposed Solution: Full Solution

Scoring Function



Proposed Solution: Full Solution

 $\underset{\theta,\xi^n \ge 0}{\operatorname{argmin}} \ \frac{1}{2} ||\theta||^2 + C \sum_n \xi^n$ s.t. $\theta^T \psi(y^n, x^n) - \theta^T \psi(h, x^n) \ge \Delta(y^n, h) - \xi^n \quad \forall_n, \forall_h$

Inference

$$y^* = \underset{y}{\operatorname{argmax}} \ \Omega(y; w, u, x, L)$$

Success



Human Label: sitting on bench in a train station because he is waiting

Top Predictions: 1. sitting on bench in a park because he is waiting

- holding a tv in a park because he wants to take
- 3. holding a seal in a park because he wants to protest
- 4. holding a guitar in a park because he wants to play



Human Label: sitting on chair in a dining room because she wants to eat

Top Predictions: 1. sitting near table in dining room because she wants to eat

- 2. sitting on a sofa in a dining room because she wants to eat
- 3. holding a cup in a dining room because she wants to eat
- 4. sitting on a cup in a dining room because she wants to eat

Failure



Human Label: holding a person in a living room because she wants to show

- **Top Predictions:** 1. sitting on sofa in living room because she wants to pet 2. sitting on sofa in living room because she wants to look 3. sitting on sofa in living room because she wants to read
 - 4. sitting on chair in living room because she wants to pet



Human Label: standing next to table because she wants to prepare

Top Predictions:1. talking to person in dining because she wants to eat
2. standing next to table in dining room because she wants to eat
3. sitting next to table in dining because she wants to eat
4. talking to person in kitchen because she wants to eat

Failure: Vision Only



Human Label: sitting on a bus in a parking lot because he wants to drive

Top Predictions: 1. because he wants to look

- 2. because he wants to ride
- 3. because he wants to drive
- 4. because he wants to eat



Human Label: sitting on chair in living room because she wants to read

Top Predictions: 1. because she wants to eat

- 2. because she wants to look
- 3. because she wants to drink
- 4. because she wants to ride

		Baseline	Our Method
		(Vision Only)	(With Language)
Given Ideal Detectors for:	Action+Object+Scene	13	10
	Action+Object	12	11
	Object+Scene	15	12
	Action+Scene	19	13
	Object	19	13
	Action	18	15
	Scene ¹	37	18
Fully Automatic		23 ²	15

Chance has rank of 39



Point of Strength

- Novel and important problem
- Simple model easy to understand
- Augmenting image with text through data mining was proven to be effective

Point of Weakness

- Results are only ok (qualitatively, failure of visiononly model does not make much more sense)
- Model is linear too simple
- Language queries are simple as well

Contributions

- Introducing the problem of inferring motivation behind people's actions to the computer vision community.
- Propose to use common knowledge mined from web to improve computer vision systems.

Conclusion

- Interesting problem
- The proposed method is more of a baseline
- Future research can extend prediction model, and language model

Thanks!

Questions?