





Re-evaluating automatic metrics for image captioning

Mert Kilickaya, Aykut Erdem, Nazli Ikizler-Cinbis, Erkut Erdem









Background



Evaluation

Human



Automatic

- Borrowed from Machine Translation
 O BLEU
 O METEOR
 - ROUGE
- Developed for Image Captioning
 CIDEr
 - SPICE





Candidate



$BLEU \cong \frac{\{cat, on, mat\}}{\{the, cat, sat, on, the, mat\}}$





Candidate



$\begin{array}{ll} ROUGE \cong \frac{\{cat, on, mat\}}{\{an, orange, cat, sitting, on, mat\}} \end{array}$









 $METEOR \cong$



{*an*,*orange*,*cat*,*sitting*,*on*,*mat*}





Candidate



		Word	Frequency
CIDEr		Kitty	0.6
	f(corpus) =	Sit	0.2
		On	0.1
		Mat	0.1









CIDEr

{0.6*kitty,0.2*sit,0.1*on,0.1*mat} {the,cat,sat,on,the,mat} X {0.6*kitty,0.2*sit,0.1*on,0.1*mat} {an,orange,cat,sitting,on,mat}

What is missing?

Metric	Proposed to Evaluate	Idea	Semantic Information
BLEU	Machine Translation	N _{gram} precision	-
ROUGE	Document Summarization	N _{gram} recall	-
METEOR	Machine Translation	N _{gram} with synonyn matching	Synonym Matching
CIDEr	Image Captioning	N _{gram} with corpus re- weighting	-



Visual Complexity

Possible Captions



	Reference	Candidate
parser	{The cat sat on the mat}	{An orange kitty sitting on mat}
pobjects	{ <mark>cat, mat</mark> }	{ <mark>kitty, mat</mark> }
Pattributes	{-}	{ <mark>orange-kitty</mark> }
p _{relations}	{cat-sat, cat-on-mat}	<pre>{kitty-sit, kitty-on-mat}</pre>





Candidate



SPICE \cong fobjects * fattributes * frelations

What if we have a *relevant* caption, with *no* overlapping content?



Reference	A man wearing a lifevest is sitting in a canoe
Candidate 1	A guy with a red jacket is standing on a boat
Candidate 2	A small white ferry rides through water

Word Mover Distance (WMD)





We have BLEU, METEOR, ROUGE, CIDEr, SPICE and WMD

Which one is better?

Correlation Analysis with Williams Significance Test

	Human Judgement	Metric 1	Metric 2
Caption 1		•••	•
Caption 2	•	•	V
Caption 3	•	•••	•
Caption 4	VÝV	•	•••

Corr(Human Judgement, Metric 1)

Corr(Human Judgement, Metric 2)

Spearman Correlation on Datasets

		Flickr-8k	Composite	
	WMD	0.60	0.43	
	SPICE	0.64	0.42	
	CIDEr	0.56	0.42	Is this
	METEOR	0.58	0.44	significant?
	BLEU	0.44	0.38	
	ROUGE	0.44	0.39	
Flickr-8k	SPICE > V	VMD > M	ETEOR > CII	DEr > BLEU > ROUGE
Composite	METEOR	> WMD >	SPICE > CII	DEr > ROUGE > BLEU









Measuring Syntactic Properties

	Description	BLEU	METEOR	ROUGE	CIDEr	SPICE	WMD
Original	A man wearing a red life jacket is sitting in a canoe on a lake	1	1	1	10	1	1

	Description	BLEU	METEOR	ROUGE	CIDEr	SPICE	WMD
Original	A man wearing a red life jacket is sitting in a canoe on a lake	1	1	1	10	1	1
Candidate	A man wearing a life jacket is in a small boat on a lake	0.45	0.28	0.68	2.19	0.40	0.19

	Description	BLEU	METEOR	ROUGE	CIDEr	SPICE	WMD
Original	A man wearing a red life jacket is sitting in a canoe on a lake	1	1	1	10	1	1
Candidate	A man wearing a life jacket is in a small boat on a lake	0.45	0.28	0.68	2.19	0.40	0.19
Synonyms	A guy wearing a life vest is in a small boat on a lake	0.20 (↓)	0.17 (↓)	0.57(↓)	0.65 <mark>(↓)</mark>	0.00 (↓)	0.10(↓)

	Description	BLEU	METEOR	ROUGE	CIDEr	SPICE	WMD
Original	A man wearing a red life jacket is sitting in a canoe on a lake	1	1	1	10	1	1
Candidate	A man wearing a life jacket is in a small boat on a lake	0.45	0.28	0.68	2.19	0.40	0.19
Synonyms	A guy wearing a life vest is in a small boat on a lake	0.20 (↓)	0.17 (↓)	0.57(↓)	0.65(↓)	0.00 (↓)	0.10(↓)
Redundancy	A man wearing a life jacket is in a small boat on a lake <mark>at</mark> sunset	0.45	0.28	0.66	2.01	0.36	0.18

	Description	BLEU	METEOR	ROUGE	CIDEr	SPICE	WMD
Original	A man wearing a red life jacket is sitting in a canoe on a lake	1	1	1	10	1	1
Candidate	A man wearing a life jacket is in a small boat on a lake	0.45	0.28	0.68	2.19	0.40	0.19
Synonyms	A guy wearing a life vest is in a small boat on a lake	0.20 (↓)	0.17 (↓)	0.57(↓)	0.65(↓)	0.00 (↓)	0.10(↓)
Redundancy	A man wearing a life jacket is in a small boat on a lake <mark>at</mark> sunset	0.45	0.28	0.66	2.01	0.36	0.18
Word order	In a small boat on a lake a man is wearing a life jacket	0.26 (↓)	0.26 (↓)	0.38(↓)	1.32 (↓)	0.40	0.19

Robustness Analysis

Gold Caption

"Replace Scene"

"Replace Person"

Distractor Type

"Share Person"

"Share Scene"

A man wearing a life jacket is in a small boat on a lake with a ferry in view

A man wearing a life jacket is in a small boat on **takeoff** with a ferry in view

A woman in a blue shirt and headscarf is in a small boat on on a lake with a ferry in view

A man is selecting a chair from a stack under a shady awning

A black and brown dog is playing on the ice at the edge of a lake

Accuracy

Case	#Instances	BLEU	METEOR	ROUGE	CIDEr	SPICE	WMD
Replace Scene	2514	0.62	0.69	0.63	0.83	0.54	0.76
Replace Person	5817	0.73	0.77	0.78	0.78	0.67	0.80
Share Scene	2621	0.79	0.85	0.79	0.81	0.70	0.87
Share Person	4596	0.78	0.85	0.78	0.83	0.67	0.88
Overall	15548	0.73	0.79	0.75	0.81	0.65	0.83

Conclusion

- Word Mover Distance helps as a semantic measure, and opens up opportunity to apply semantic relevancy literature for this task!
- Metrics tend to be effected by small disturbances like scene change, synonym words, or word ordering
- Is current success of image captioning metrics illusion?
- Correlation alone is not enough. It should always coupled with Significance testing over improvements
- There is a big room for improvement for novel metrics, and the best practices proposed in this paper can be useful to show its' contribution

References

- **[BLEU]** Papineni, Kishore, et al. "BLEU: a method for automatic evaluation of machine translation." *Proceedings* of the 40th annual meeting on association for computational linguistics. Association for Computational Linguistics, 2002.
- **[ROUGE]** Lin, Chin-Yew. "Rouge: A package for automatic evaluation of summaries." *Text summarization branches out: Proceedings of the ACL-04 workshop*. Vol. 8. 2004.
- **[METEOR]** Banerjee, Satanjeev, and Alon Lavie. "METEOR: An automatic metric for MT evaluation with improved correlation with human judgments." *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*. Vol. 29. 2005.
- [CIDEr] Vedantam, Ramakrishna, C. Lawrence Zitnick, and Devi Parikh. "Cider: Consensus-based image description evaluation." *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2015.
- **[SPICE]** Anderson, Peter, et al. "Spice: Semantic propositional image caption evaluation." *European Conference on Computer Vision*. Springer International Publishing, 2016.
- [Williams Test] Graham, Yvette. "Improving Evaluation of Machine Translation Quality Estimation." ACL (1). 2015.
- [Robustness Corpus] Hodosh, Micah, and Julia Hockenmaier. "Focused evaluation for image description with binary forced-choice tasks." *Workshop on Vision and Language, Annual Meeting of the Association for Computational Linguistics*. Vol. 3. 2016.
- [Word Mover Distance] Kusner, Matt J., et al. "From Word Embeddings To Document Distances." ICML. Vol. 15. 2015.