

OpenTag: Open Attribute Value Extraction from Product Profiles

KDD 2018



Get it as soon as Wednesday, Feb. 14 when you choose Two-Day Shipping at checkout Ships from and sold by Cunningham Collective.

by the AAFCO Dog Food Nutrient Profiles for maintenance Product description Variety pack includes: 6 trays of Filet Mignon flavor in meaty juices 6 trays of Porterhouse Steak flavor in meaty juices

the first to review this item

Variety Pack Filet Mignon and

Count) Price: \$92.60 & FREE Shipping

Porterhouse Steak Dog Food (12

6 travs of Filet Mignon flavor in meaty juices

palatability to tempt even the fussiest dog

Cesar pet food has an irresistible taste with exceptiona

Formulated to meet the nutritional levels established

Cesar pet food has an irresistible taste with exceptional palatability to tempt even the fussiest dogs Formulated to meet the nutritional levels established by the AAFCO Dog Food Nutrient Profiles for maintenance Complete & balanced nutrition for small adult dogs Fortified with vitamins and minerals Packaged in convenient feeding trays with no-fuss, peel-away freshness seals Includes 6 Each Chicken & Live

Extracting Structured Values from Unstructured Text

- Given a product title "Variety Pack Fillet Mignon and Porterhouse Steak Dog Food (12 Count)": OpenTag extracts attribute values "size = 12 count" and "flavor = {Porterhouse Streak, Fillet Mignon}"
- It extracts structured multi-word values, and multiple values for an attribute
- Sequence Tagging: It exploits *distributional semantics* to assign tags to tokens to extract attribute values

Sequence	duck	,	fillet	mignon	and	ranch
BIOE	В	0	В	Е	0	В
UBIOE	U	0	В	E	Ο	В
IOB	В	0	В	Ι	Ο	В

OpenTag: Active Learning

- Uses *active learning* to reduce manual annotation effort
- To identify difficult-to-extract instances for the model, ar
- Uses tag flips as a Query Strategy to identify difficult training instances
 - Tag Flip: Change in tag of a token across successive epochs during training
 - Frequent tag flips indicate OpenTag is uncertain about the sample; not stable



Guineng Zheng*, Subhabrata Mukherjee[∆], Xin Luna Dong[∆], FeiFei Li* ^AMazon.com, *University of Utah {subhomj, lunadong}@amazon.com, {guineng, lifeifei}@cs.utah.edu

- Amazon Catalog provides structured information on product attributes. This information is often noisy or missing for a lot of attributes • We develop **OpenTag: a novel deep tagging model** to discover *missing* values of attributes from unstructured product profiles like *title, description, and bullets*
- OpenTag does not rely on dictionaries of values or hand-crafted features
- It discovers *new* attribute values *never* encountered before with an *open world assumption*
- It uses active learning to reduce manual annotation effort by 3.3x

lamb

raised

flavor

()

\frown	_	_	\frown		
()	nen	na	(: om)	nnn	ente
		iug			

- Word Embeddings model distributed attribute-value representations
- Bidirectional LSTM (BiLSTM) captures contextual information from long and short range dependencies via hidden states: $h_t = \sigma([h_t, h_t])$
- BiLSTM captures sequential nature of tokens but not tags
- Conditional Random Field (CRF) considers sequential nature of tags to extract coherent attribute values: $\Pr(y|$

Attention mechanism captures in
and generates interpretable expl

nd	as	k fo	or a	inno	otat	ion

Datasets/Attribute	N
Dog Food: Title	H
Attribute: Flavor	H
	(
Camera: Title	F
Attribute: Brand name	F
	(
Detergent: Title	F
Attribute: Scent	F
	(
Dog Food: Description	H
Attribute: Flavor	H
	(
Dog Food: Bullet	I
Attribute: Flavor	H
	(
Dog Food: Title	F
Multi Attribute:	H
Brand, Flavor, Capacity	(

Extraction via Sequence Tagging

$$x; \Psi) \propto \prod_{t=1}^{T} exp\left(\sum_{k=1}^{K} \psi_k f_k(y_{t-1}, y_t, \langle l_t \rangle)\right)$$

nportance of tags for tokens, important concepts, planation for its verdict: $l_t = \sum \alpha_{t,t'} \cdot h_{t'}$

CRF Layer **Attention Mechanism BiLSTM** Layer Word Embedding





Interpretation via Attention



on of word vectors before attention



(c) Heatmap of attention matrix

