



# LS-Tree: Model Interpretation When the Data Are Linguistic

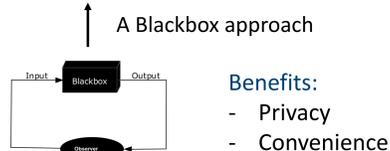
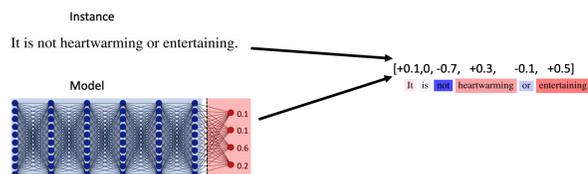
Jianbo Chen, Michael I. Jordan  
UC Berkeley

## ABSTRACT

We study the problem of interpreting trained classification models in the setting of linguistic data sets. Leveraging a parse tree, we propose to assign least-squares-based importance scores to each word of an instance by exploiting syntactic constituency structure. We establish an axiomatic characterization of these importance scores by relating them to the Banzhaf value in coalitional game theory. Based on these importance scores, we develop a principled method for detecting and quantifying interactions between words in a sentence.

## OBJECTIVE

For a given instance, assign a vector of importance scores for each feature.



## MOTIVATION

- Debugging a model



- Increase trust in decision making



## EXISTING METHODS

- LIME (Ribeiro, Singh, and Guestrin 2016)
- Representation Erasure (Li, Monroe, and Jurafsky 2016)
- Quantitative Input Influence (QII) (Datta, Sen, and Zick 2016)
- SHapley Additive exPlanations (SHAP) (Lundberg and Lee 2017)
- L-Shapley and C-Shapley (Chen, et.al. 2018)

### Procedures:

- Step 1: Sample word subsets with a certain scheme
- Step 2: Evaluate target model f on each sampled word subset

A specific example - Shapley value (Shapley 1953):

It is **not** heartwarming or entertaining  
 $f(\text{"not heartwarming"}) - f(\text{"heartwarming"})$   
**It is not** heartwarming or entertaining  
 $f(\text{"It is not"}) - f(\text{"It is"})$   
**It is not** heartwarming or entertaining  
 $f(\text{"It ... not"}) - f(\text{"It"})$

Marginal contribution of i to S:

$$f(S \cup \{i\}) - f(S)$$

where

$$f(S) := f(x_S)$$

- Step 3: Combine model evaluations into attribution scores

A specific example - Shapley value (Shapley 1953):

$$\phi_{f,x}(i) = \frac{1}{d} \sum_{S \subset [d]} \frac{1}{\binom{d-1}{|S|-1}} (f(S \cup \{i\}) - f(S))$$

## LIMITATIONS OF EXISTING METHODS

It is **not** heartwarming or entertaining  $f(\text{"It ... not"}) - f(\text{"It"})$

'It ... not' is not natural language.  
The target model may not respond appropriately.

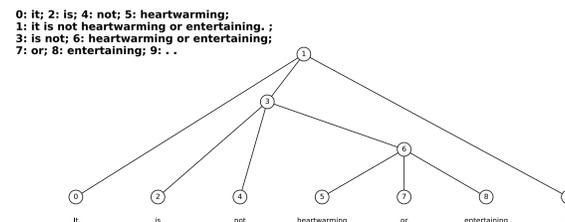
It is **not** heartwarming or entertaining

Is 'not' important as a single word, or because of its interaction with 'heartwarming'

## CONSTITUENCY PARSING FOR LINGUISTIC DATA

What expressions are valid to human?

What interactions are we interested in?



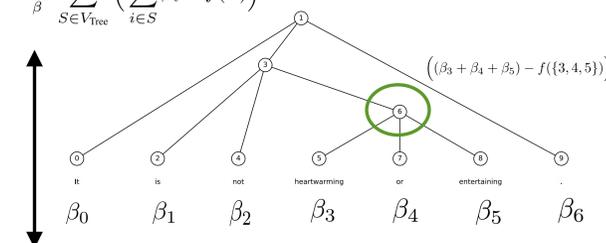
0: It; 2: is; 4: not; 5: heartwarming;  
1: it is not heartwarming or entertaining;  
3: is not; 6: heartwarming or entertaining;  
7: or; 8: entertaining; 9: .

## LS-TREE

Least squares Cook's interaction score

### Step 1: Least squares

$$\min_{\beta} \sum_{S \in V_{Tree}} \left( \sum_{i \in S} \beta_i - f(S) \right)^2$$



An axiomatic framework based on Banzhaf value.

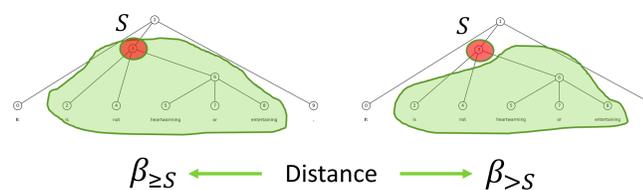
### Cook's distance (Cook 1977)

Capture the influence of instance i:

$$D_i = \text{Const.} \cdot (\hat{\beta}_{(i)} - \hat{\beta})^T X^T X (\hat{\beta}_{(i)} - \hat{\beta})$$

$\hat{\beta}_{(i)}$ : Fit a linear model without data point i.

### Step 2: Influence of the intersection at node S



All nodes: An  $\Theta(d^3)$  algorithm using the Sherman-Morrison formula.

## ADVERSATIVE RELATIONS

Examples: not, but, yet, though, although, even though, whereas, except, despite, in spite of

Dataset	Model	Avg. Score	not	but	yet	though	although	even though	whereas	except	despite	in spite of
SST	BoW	0.153	0.0006(318)	0.0000(079)	0.0002(105)	0.0000(865)	0.0002(222)	0.0000(000)	-0.1	0.0004(280)	0.0003(519)	0.0000(000)
	CNN	0.634	1.6734(592)	1.6941(444)	0.5680(959)	0.2130(735)	1.6750(462)	0.6260(407)	-0.1	0.9448(175)	1.4824(270)	2.1101(945)
	LSTM	0.79	1.7462(580)	1.5031(531)	1.4602(368)	1.1531(404)	0.3340(162)	1.7460(586)	-0.1	0.8831(325)	1.7621(319)	0.5940(261)
	BERT	0.428	1.7142(383)	2.1481(760)	1.6609(120)	1.8285(268)	1.7441(356)	1.8885(162)	-0.1	1.5563(331)	1.7602(598)	0.8276(252)
IMDB	BoW	0.019	0.0002(283)	0.0000(000)	0.0002(210)	0.0000(000)	0.0000(000)	0.0000(000)	0.0000(000)	0.0001(362)	0.0000(125)	-0.1
	CNN	0.424	1.0500(819)	3.4420(021)	1.4890(295)	0.9230(085)	1.0360(071)	1.1750(467)	0.4091(106)	1.2994(167)	0.6030(434)	-0.1
	LSTM	0.120	0.0003(187)	2.2210(524)	1.5000(248)	0.0110(087)	0.0011(270)	0.9440(483)	1.220(365)	1.2944(400)	0.3600(500)	-0.1
	BERT	1.159	1.4162(207)	3.3901(800)	1.6441(152)	1.2712(061)	1.5812(123)	1.4871(577)	0.2810(149)	1.4212(160)	1.2182(241)	-0.1
Yelp	BoW	0.035	0.0000(488)	0.0001(015)	0.0000(533)	0.0001(664)	0.0001(128)	0.0000(000)	0.0000(530)	0.0000(367)	0.0001(213)	-0.1
	CNN	0.401	2.2870(407)	2.4640(025)	0.5160(043)	0.9800(435)	0.8890(075)	0.7900(011)	0.2600(071)	0.8220(529)	0.4230(889)	-0.1
	LSTM	0.224	2.1735(950)	1.7121(676)	0.9882(165)	0.9844(310)	0.7061(194)	0.3590(483)	1.3981(793)	0.3441(408)	0.3141(135)	-0.1
	BERT	0.746	1.3642(186)	2.4403(658)	0.7810(184)	1.3260(933)	0.5960(615)	1.4190(886)	0.0950(162)	0.3310(074)	1.4410(414)	-0.1

Size of data set: SST (10K) < IMDB (100K) < Yelp (600K)

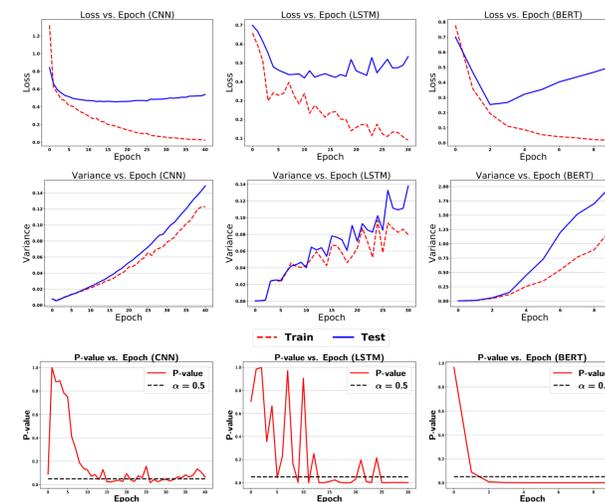
## IS "WHILE" INDICATING A CONTRAST?

Interaction scores of the parent node of "while".

Sentence	Meaning	BoW	CNN	LSTM	BERT
... He said he couldn't help. We had to walk while the snow blue to our faces. When we were almost there, we saw the shuttle pull out with the smoking shuttle driver in it, driving in the opposite direction, away from us. I can not believe how rude they were.	during the time that	0.0000(338)	0.7810(300)	1.761(0.839)	0.062(0.092)
... I ordered a cappuccino. It tasted like milk and no coffee. I was exceptionally disappointed. So while the place has a great reputation, even they can screw it up if they don't pay attention to detail, and at this level they should never screw it up. I had a better cup at Marty's Market for crying out loud!	whereas (indicating a contrast)	0.0000(338)	0.1420(300)	2.155(0.839)	2.167(0.092)
Usually asking the server what is her favorite dish gets you a pretty good recommendation, but in this case, it was crap! The smoked brisket had that discoloration that happens to meat when it's been sitting out for a while. And it wasn't even tender!! Am I asking for too much?	a period of time	0.0000(338)	0.2060(300)	0.4650(839)	0.082(0.092)

## OVERFITTING

Difference between variances of interaction scores between training and test sets as a diagnostic for overfitting.



Permutation test under the null hypothesis of equal average variance