

# Learning to Rank for Plausible Plausibility

Zhongyang Li<sup>1,2</sup> Tongfei Chen<sup>1</sup> Benjamin Van Durme<sup>1</sup>

<sup>1</sup> Johns Hopkins University; <sup>2</sup> Harbin Institute of Technology

## SUMMARY

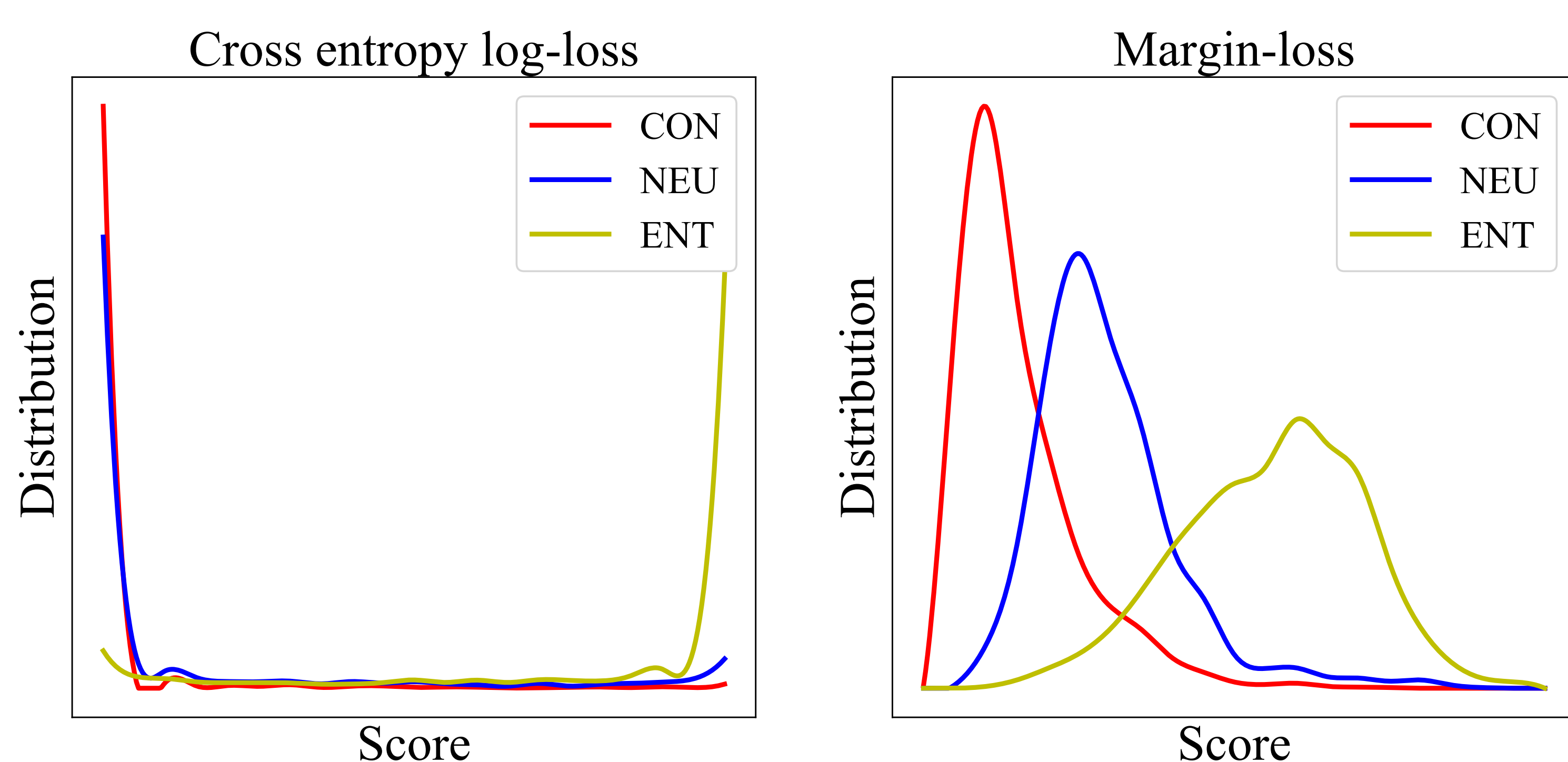
- The categorical nature of RTE / NLI leads to common use of cross-entropy loss for training, even when some data are designed for *plausibility* prediction, not *entailment*
- This loss is intuitively odd: it drives models to assign score 0.0 or 1.0, in contrast to the proposed *margin-based* loss
- Leads to better calibrated predictions to human intuitions
- State-of-the-art result on the Choice of Plausible Alternatives (COPA) task

## CROSS-ENTROPY LOSS

$$P(h_i|p) = \frac{\exp F(p, h_i)}{\sum_i \exp F(p, h_i)}$$

As a *classification* problem – maximizing the posterior probability normalized over all hypothesis alternatives

## PLAUSIBLE PLAUSIBILITY



Margin-based loss gives much more **plausible** scores, leading to a more **plausible** *plausibility* formulation!

## RESULTS ON COPA

Model adapted from the original BERT sentence pair model.

Method	Accuracy
PMI (Jabeen+, 2014)	58.8
PMI_EX (Gordon+, 2011)	65.4
CS (Luo+, 2016)	70.2
CS_MWP (Sasaki+, 2017)	71.2
BERT (cross-entropy)	73.4
<b>BERT (margin)</b>	<b>75.4</b>

## THE PLAUSIBILITY TASK

Given a premise, which hypothesis is *preferred*?

$p$	<i>I just stopped where I was</i>	
$h_E$	<i>I stopped in my tracks</i>	✓
$h_N$	<i>I stopped running right were I was</i>	✗
$h_N$	<i>I stopped running right were I was</i>	✓
$h_C$	<i>I continued on my way</i>	✗

E = Entailed; N = Neutral; C = Contradictory

**In the plausibility task, the correct label depends on the alternatives!**

## MARGIN-BASED LOSS

$$L = \frac{1}{N} \sum_{h>h'} \max\{0, \xi - F(p, h) + F(p, h')\}$$

As a *learning-to-rank* problem – more plausible hypothesis should *rank* higher than other less plausible hypotheses

## EXAMPLE PREDICTIONS

### Recast MultiNLI

$p$		Log	Margin
$p$	<i>I just stopped where I was</i>		
$h_1$	<i>I stopped in my tracks</i>	0.919	0.568
$h_2$	<i>I stopped running right were I was</i>	0.081	0.358
$h_3$	<i>I continued on my way</i>	$1.71 \times 10^{-8}$	0.074

### JOCI

$p$		Log	Margin
$p$	<i>Cheerleaders performs in a lift stunt.</i>		
$h_1$	<i>The stunt is a feat.</i>	0.508	0.304
$h_2$	<i>The stunt is no fluke.</i>	0.486	0.279
$h_3$	<i>The stunt is dangerous.</i>	$2.72 \times 10^{-4}$	0.166
$h_4$	<i>The stunt is remarkable.</i>	$4.13 \times 10^{-3}$	0.153
$h_5$	<i>The stunt backfires.</i>	$2.36 \times 10^{-4}$	0.107

## RELATED DISCUSSIONS

J. Opitz & A. Frank (2018): Addressing the Winograd Schema Challenge as a sequence ranking task. *1st Int'l Workshop on Language Congition and Computational Models*.

## PAPER

