

# A Proposal on Evaluation Measures for RTE

Richard Bergmair

University of Cambridge Computer Laboratory  
Natural Language Information Processing

ACL/IJCNLP 2009 Workshop on Applied Textual Inference,  
Aug-6 2009

# Problems with the Current Methodology

- ▶ distribution neither balanced nor representative; so accuracy figures biased.
- ▶ notion of confidence-ranking misleading; accuracy & thresholding contradicts average precision.
- ▶ ENTAILMENT/CONTRADICTION symmetric; average precision doesn't reflect that.

# Problems with the Current Methodology

- ▶ distribution neither balanced nor representative;  
so accuracy figures biased.
- ▶ notion of confidence-ranking misleading;  
accuracy & thresholding contradicts average precision.
- ▶ ENTAILMENT/CONTRADICTION symmetric;  
average precision doesn't reflect that.

# Problems with the Current Methodology

- ▶ distribution neither balanced nor representative; so accuracy figures biased.
- ▶ notion of confidence-ranking misleading; accuracy & thresholding contradicts average precision.
- ▶ ENTAILMENT/CONTRADICTION symmetric; average precision doesn't reflect that.



# Proposal for a New Methodology

- ▶ best: **mutual information**.
- ▶ average precision completely unsuitable!  
confidence-weighted score preferable, but there are still drawbacks;
- ▶ report baselines, be aware of bias with accuracies;  
perhaps use an artificially balanced subset for 3-way task.

# Outline

The Structure of RTE Data

Accuracy

Average Precision

Mutual Information

# Outline

## The Structure of RTE Data

Accuracy

Average Precision

Mutual Information



# RTE Data

N candidate entailments:

$$X = \{x_1, x_2, \dots, x_N\}.$$

gold standard:

$$G : X \mapsto \{\oplus, \diamond, \ominus\}.$$

system output:

$(L, >)$  where

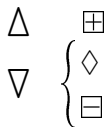
$$L : X \mapsto \{\oplus, \diamond, \ominus\},$$

strict total order  $>$  on  $X$ .

# RTE Data

class labels:

two-way    three way



# Logical Structure of Candidate Entailments

$$(x_{42}) \frac{\text{Socrates is a man and every man is mortal. } (\varphi)}{\therefore \text{Socrates is mortal. } (\psi)}$$

$$\Box(\varphi \rightarrow \psi)$$

$$G(x_{42}) = \Box \quad \neg G(\neg x_{42}) = \Box$$

# Logical Structure of Candidate Entailments

$$(x_{42}) \frac{\text{Socrates is a man and every man is mortal.} \quad (\varphi)}{\therefore \text{Socrates is mortal.} \quad (\psi)}$$

$$\Box(\varphi \rightarrow \psi)$$

$$G(x_{42}) = \Box \quad \neg G(\neg x_{42}) = \Box$$

# Logical Structure of Candidate Entailments

$$(x_{42}) \frac{\text{Socrates is a man and every man is mortal.} \quad (\varphi)}{\therefore \text{Socrates is mortal.} \quad (\psi)}$$

$$\Box(\varphi \rightarrow \psi)$$

$$G(x_{42}) = \Box$$

$$\neg G(\neg x_{42}) = \Box$$

# Logical Structure of Candidate Entailments

$(\neg x_{42}) \frac{\text{Socrates is a man and every man is mortal.} \quad (\varphi)}{\therefore \text{Socrates is **not** mortal.} \quad (\neg\psi)}$

$\square(\varphi \rightarrow \neg\psi)$

$G(x_{42}) = \boxplus$

$\neg G(\neg x_{42}) = \boxminus$

# Logical Structure of Candidate Entailments

(x<sub>43</sub>)  $\frac{\text{Socrates is a man and every man is mortal.} \quad (\varphi)}{\therefore \text{Socrates is (not) wise.} \quad (\psi', \neg\psi')}$

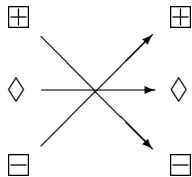
$\diamond(\varphi \rightarrow \psi)$

$\diamond(\varphi \rightarrow \neg\psi)$

$G(x_{43}) = \diamond$

$\neg G(\neg x_{42}) = \diamond$

# Logical Structure of Candidate Entailments





# Outline

The Structure of RTE Data

**Accuracy**

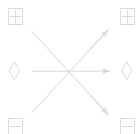
Average Precision

Mutual Information

# Accuracy

$$\mathbb{A}_3(\mathbf{L}; \mathbf{G}) = \frac{1}{N} \sum_{i=1}^N \mathbb{1}([\mathbf{L}(\mathbf{x}_i)]_3 = [\mathbf{G}(\mathbf{x}_i)]_3),$$

$$\mathbb{A}_2(\mathbf{L}; \mathbf{G}) = \frac{1}{N} \sum_{i=1}^N \mathbb{1}([\mathbf{L}(\mathbf{x}_i)]_2 = [\mathbf{G}(\mathbf{x}_i)]_2),$$



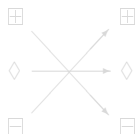
$$\mathbb{A}_3(\mathbf{L}; \mathbf{G}) = \mathbb{A}_3(\neg\mathbf{L}; \neg\mathbf{G})$$

$$\mathbb{A}_2(\mathbf{L}; \mathbf{G}) = \mathbb{A}_2(\neg\mathbf{L}; \neg\mathbf{G})$$

# Accuracy

$$\mathbb{A}_3(\mathbf{L}; \mathbf{G}) = \frac{1}{N} \sum_{i=1}^N \mathbb{1} \left( [\mathbf{L}(\mathbf{x}_i)]_3 = [\mathbf{G}(\mathbf{x}_i)]_3 \right),$$

$$\mathbb{A}_2(\mathbf{L}; \mathbf{G}) = \frac{1}{N} \sum_{i=1}^N \mathbb{1} \left( [\mathbf{L}(\mathbf{x}_i)]_2 = [\mathbf{G}(\mathbf{x}_i)]_2 \right),$$



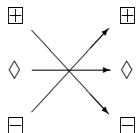
$$\mathbb{A}_3(\mathbf{L}; \mathbf{G}) = \mathbb{A}_3(\neg\mathbf{L}; \neg\mathbf{G})$$

$$\mathbb{A}_2(\mathbf{L}; \mathbf{G}) = \mathbb{A}_2(\neg\mathbf{L}; \neg\mathbf{G})$$

# Accuracy & Logical Symmetry

$$\mathbb{A}_3(\mathbf{L}; \mathbf{G}) = \frac{1}{N} \sum_{i=1}^N \mathbb{1} \left( [\mathbf{L}(\mathbf{x}_i)]_3 = [\mathbf{G}(\mathbf{x}_i)]_3 \right),$$

$$\mathbb{A}_2(\mathbf{L}; \mathbf{G}) = \frac{1}{N} \sum_{i=1}^N \mathbb{1} \left( [\mathbf{L}(\mathbf{x}_i)]_2 = [\mathbf{G}(\mathbf{x}_i)]_2 \right),$$



$$\mathbb{A}_3(\mathbf{L}; \mathbf{G}) = \mathbb{A}_3(\neg \mathbf{L}; \neg \mathbf{G})$$

$$\mathbb{A}_2(\mathbf{L}; \mathbf{G}) = \mathbb{A}_2(\neg \mathbf{L}; \neg \mathbf{G})$$

# Conditioned Accuracy

$$\mathbb{A}'_3(\mathbf{L}; \mathbf{G}, g) = \frac{\sum_{i=1}^N \mathbb{1}([\mathbf{L}(\mathbf{x}_i)]_3 = [\mathbf{G}(\mathbf{x}_i)]_3 = g)}{\sum_{i=1}^N \mathbb{1}([\mathbf{G}(\mathbf{x}_i)]_3 = g)},$$

$$\mathbb{A}'_2(\mathbf{L}; \mathbf{G}, g) = \frac{\sum_{i=1}^N \mathbb{1}([\mathbf{L}(\mathbf{x}_i)]_2 = [\mathbf{G}(\mathbf{x}_i)]_2 = g)}{\sum_{i=1}^N \mathbb{1}([\mathbf{G}(\mathbf{x}_i)]_2 = g)}.$$

$\mathbb{A}'_2(\mathbf{L}; \mathbf{G}, \Delta) \rightsquigarrow$  recall,

$\mathbb{A}'_2(\mathbf{G}; \mathbf{L}, \Delta) \rightsquigarrow$  precision.

# Bias

labels:	RTE-4	RTE-3 PILOT
ENTAILMENT ( $\boxplus$ )	50%	51%
UNKNOWN ( $\diamond$ )	35%	40%
CONTRADICTION ( $\boxminus$ )	15%	9%

$\mathbb{A}_3(L^{\boxplus}; G) = .500$  outperforms **1/3** of all RTE4 participants and **2/3** of all RTE3 PILOT participants!

$\mathbb{A}_3(L^*; G) = .394$ ,  $\mathbb{A}_3(L^{\diamond}; G) = .350$ ,  $\mathbb{A}_3(L^{\boxminus}; G) = .150$

# Outline

The Structure of RTE Data

Accuracy

**Average Precision**

Mutual Information

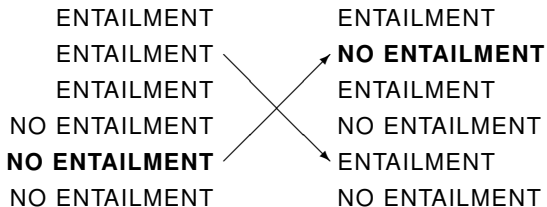
# Average Precision

ENTAILMENT  
ENTAILMENT  
ENTAILMENT  
NO ENTAILMENT  
NO ENTAILMENT  
NO ENTAILMENT

ENTAILMENT  
NO ENTAILMENT  
ENTAILMENT  
NO ENTAILMENT  
ENTAILMENT  
NO ENTAILMENT



# Average Precision



# Confidence-Weighted Score

ENTAILMENT  
ENTAILMENT  
ENTAILMENT  
NO ENTAILMENT  
**NO ENTAILMENT**  
NO ENTAILMENT

ENTAILMENT  
**NO ENTAILMENT**  
ENTAILMENT  
NO ENTAILMENT  
ENTAILMENT  
NO ENTAILMENT



# Average Precision vs. Confidence

- ▶ 2/3 of all RTE-4 participants who submitted confidence-ranked three-way labellings submitted confidence rankings, instead of *AP*-style rankings.
- ▶ Stanford1: 44% → 62% after ranking down negative instances.
- ▶ confusing terminology!
- ▶ more generally: accuracy & average precision have contradictory preferences for rankings.

# Average Precision vs. Confidence

- ▶ 2/3 of all RTE-4 participants who submitted confidence-ranked three-way labellings submitted confidence rankings, instead of *AP*-style rankings.
- ▶ Stanford1: 44% → 62% after ranking down negative instances.
- ▶ confusing terminology!
- ▶ more generally: accuracy & average precision have contradictory preferences for rankings.

# Average Precision vs. Confidence

- ▶ 2/3 of all RTE-4 participants who submitted confidence-ranked three-way labellings submitted confidence rankings, instead of *AP*-style rankings.
- ▶ Stanford1: 44% → 62% after ranking down negative instances.
- ▶ confusing terminology!
- ▶ more generally: accuracy & average precision have contradictory preferences for rankings.

# Average Precision vs. Confidence

- ▶ 2/3 of all RTE-4 participants who submitted confidence-ranked three-way labellings submitted confidence rankings, instead of *AP*-style rankings.
- ▶ Stanford1: 44% → 62% after ranking down negative instances.
- ▶ confusing terminology!
- ▶ more generally: accuracy & average precision have contradictory preferences for rankings.

# Average Precision & Logical Symmetry

id		system	gold
223	1	ENTAILMENT	UNKNOWN
	1	ENTAILMENT	UNKNOWN
4	2	ENTAILMENT	ENTAILMENT
	2	ENTAILMENT	CONTRADICTION
313	3	UNKNOWN	UNKNOWN
	3	UNKNOWN	UNKNOWN
534	4	CONTRADICTION	CONTRADICTION
	4	CONTRADICTION	ENTAILMENT
415	5	CONTRADICTION	ENTAILMENT
	5	CONTRADICTION	CONTRADICTION

# Average Precision & Logical Symmetry

id		system	gold
223	1	ENTAILMENT	UNKNOWN
	1	ENTAILMENT	UNKNOWN
4	2	ENTAILMENT	ENTAILMENT
	2	ENTAILMENT	CONTRADICTION
313	3	UNKNOWN	UNKNOWN
	3	UNKNOWN	UNKNOWN
534	4	CONTRADICTION	CONTRADICTION
	4	CONTRADICTION	ENTAILMENT
415	5	CONTRADICTION	ENTAILMENT
	5	CONTRADICTION	CONTRADICTION



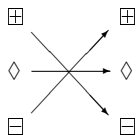
# Average Precision & Logical Symmetry

id		system	gold
223	1	CONTRADICTION	UNKNOWN
	1	ENTAILMENT	UNKNOWN
4	2	CONTRADICTION	ENTAILMENT
	2	ENTAILMENT	CONTRADICTION
313	3	UNKNOWN	UNKNOWN
	3	UNKNOWN	UNKNOWN
534	4	ENTAILMENT	CONTRADICTION
	4	CONTRADICTION	ENTAILMENT
415	5	ENTAILMENT	ENTAILMENT
	5	CONTRADICTION	CONTRADICTION

# Average Precision & Logical Symmetry

id		system	gold
223	5	CONTRADICTION	UNKNOWN
	1	ENTAILMENT	UNKNOWN
4	4	CONTRADICTION	ENTAILMENT
	2	ENTAILMENT	CONTRADICTION
313	3	UNKNOWN	UNKNOWN
	3	UNKNOWN	UNKNOWN
534	2	ENTAILMENT	CONTRADICTION
	4	CONTRADICTION	ENTAILMENT
415	1	ENTAILMENT	ENTAILMENT
	5	CONTRADICTION	CONTRADICTION

# Average Precision & Logical Symmetry

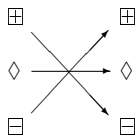


$$A_3(L; G) = A_3(\neg L; \neg G)$$

$$A_2(L; G) = A_2(\neg L; \neg G)$$

$$AP(G; >) \neq AP(\neg G; >')$$

# Average Precision & Logical Symmetry



$$\mathbb{A}_3(L; G) = \mathbb{A}_3(\neg L; \neg G)$$

$$\mathbb{A}_2(L; G) = \mathbb{A}_2(\neg L; \neg G)$$

$$AP(G; >) \neq AP(\neg G; >')$$

# Outline

The Structure of RTE Data

Accuracy

Average Precision

**Mutual Information**

# Mutual Information: Definition

**marginals**

20	25	5	$P(\mathbf{G} = \boxplus)$ = .5
9	18	9	$P(\mathbf{G} = \diamond)$ = .36
1	7	6	$P(\mathbf{G} = \boxminus)$ = .14
$P(\mathbf{L} = \boxplus)$ = .3	$P(\mathbf{L} = \diamond)$ = .5	$P(\mathbf{L} = \boxminus)$ = .2	<b>N</b> <b>= 100</b>
$H(\mathbf{G} \mathbf{L} = \boxplus)$ = 1.0746	$H(\mathbf{G} \mathbf{L} = \diamond)$ = 1.4277	$H(\mathbf{G} \mathbf{L} = \boxminus)$ = 1.5395	

$$\begin{aligned} I(\mathbf{G}; \mathbf{L}) &= H(\mathbf{G}) - H(\mathbf{G}|\mathbf{L}) \\ &= 1.4277 - 1.3441 = 0.0834 \end{aligned}$$

# Mutual Information: Definition

**marginals**

20	25	5	$\mathbb{P}(\mathbf{G} = \boxplus)$ = .5
9	18	9	$\mathbb{P}(\mathbf{G} = \diamond)$ = .36
1	7	6	$\mathbb{P}(\mathbf{G} = \boxminus)$ = .14
$\mathbb{P}(\mathbf{L} = \boxplus)$ = .3	$\mathbb{P}(\mathbf{L} = \diamond)$ = .5	$\mathbb{P}(\mathbf{L} = \boxminus)$ = .2	<b>N</b> = 100
$\mathbb{H}(\mathbf{G} \mathbf{L} = \boxplus)$ = 1.0746	$\mathbb{H}(\mathbf{G} \mathbf{L} = \diamond)$ = 1.4277	$\mathbb{H}(\mathbf{G} \mathbf{L} = \boxminus)$ = 1.5395	

$$\begin{aligned} I(\mathbf{G}; \mathbf{L}) &= \mathbb{H}(\mathbf{G}) - \mathbb{H}(\mathbf{G}|\mathbf{L}) \\ &= 1.4277 - 1.3441 = 0.0834 \end{aligned}$$

# Mutual Information: Definition

**marginals**

20	25	5	$P(\mathbf{G} = \boxplus)$ = .5
9	18	9	$P(\mathbf{G} = \diamond)$ = .36
1	7	6	$P(\mathbf{G} = \boxminus)$ = .14
$P(\mathbf{L} = \boxplus)$ = .3	$P(\mathbf{L} = \diamond)$ = .5	$P(\mathbf{L} = \boxminus)$ = .2	<b>N</b> = 100
$H(\mathbf{G} \mathbf{L} = \boxplus)$ = 1.0746	$H(\mathbf{G} \mathbf{L} = \diamond)$ = 1.4277	$H(\mathbf{G} \mathbf{L} = \boxminus)$ = 1.5395	

$$\begin{aligned} I(\mathbf{G}; \mathbf{L}) &= H(\mathbf{G}) - H(\mathbf{G}|\mathbf{L}) \\ &= 1.4277 - 1.3441 = 0.0834 \end{aligned}$$



# Mutual Information: Definition

**prior entropy**

20	25	5	$\mathbb{P}(\mathbf{G} = \boxplus)$ = .5
9	18	9	$\mathbb{P}(\mathbf{G} = \diamond)$ = .36
1	7	6	$\mathbb{P}(\mathbf{G} = \boxminus)$ = .14
$\mathbb{P}(\mathbf{L} = \boxplus)$ = .3	$\mathbb{P}(\mathbf{L} = \diamond)$ = .5	$\mathbb{P}(\mathbf{L} = \boxminus)$ = .2	$\mathbb{H}(\mathbf{G})$ = 1.4277
$\mathbb{H}(\mathbf{G} \mathbf{L} = \boxplus)$ = 1.0746	$\mathbb{H}(\mathbf{G} \mathbf{L} = \diamond)$ = 1.4277	$\mathbb{H}(\mathbf{G} \mathbf{L} = \boxminus)$ = 1.5395	$\mathbb{H}(\mathbf{G} \mathbf{L})$ = 1.3441

$$\begin{aligned}I(\mathbf{G}; \mathbf{L}) &= \mathbb{H}(\mathbf{G}) - \mathbb{H}(\mathbf{G}|\mathbf{L}) \\ &= 1.4277 - 1.3441 = 0.0834\end{aligned}$$

# Mutual Information: Definition

**prior entropy**

20	25	5	$\mathbb{P}(\mathbf{G} = \boxplus)$ = .5
9	18	9	$\mathbb{P}(\mathbf{G} = \diamond)$ = .36
1	7	6	$\mathbb{P}(\mathbf{G} = \boxminus)$ = .14
$\mathbb{P}(\mathbf{L} = \boxplus)$ = .3	$\mathbb{P}(\mathbf{L} = \diamond)$ = .5	$\mathbb{P}(\mathbf{L} = \boxminus)$ = .2	$\mathbb{H}(\mathbf{G})$ = 1.4277
$\mathbb{H}(\mathbf{G} \mathbf{L} = \boxplus)$ = 1.0746	$\mathbb{H}(\mathbf{G} \mathbf{L} = \diamond)$ = 1.4277	$\mathbb{H}(\mathbf{G} \mathbf{L} = \boxminus)$ = 1.5395	$\mathbb{H}(\mathbf{G} \mathbf{L})$ = 1.3441

$$\begin{aligned}I(\mathbf{G}; \mathbf{L}) &= \mathbb{H}(\mathbf{G}) - \mathbb{H}(\mathbf{G}|\mathbf{L}) \\ &= 1.4277 - 1.3441 = 0.0834\end{aligned}$$

# Mutual Information: Definition

entropy after specific decisions

20	25	5	$\mathbb{P}(\mathbf{G} = \boxplus)$ = .5
9	18	9	$\mathbb{P}(\mathbf{G} = \diamond)$ = .36
1	7	6	$\mathbb{P}(\mathbf{G} = \boxminus)$ = .14
$\mathbb{P}(\mathbf{L} = \boxplus)$ = .3	$\mathbb{P}(\mathbf{L} = \diamond)$ = .5	$\mathbb{P}(\mathbf{L} = \boxminus)$ = .2	$\mathbb{H}(\mathbf{G})$ = 1.4277
$\mathbb{H}(\mathbf{G} \mathbf{L} = \boxplus)$ = 1.0746	$\mathbb{H}(\mathbf{G} \mathbf{L} = \diamond)$ = 1.4277	$\mathbb{H}(\mathbf{G} \mathbf{L} = \boxminus)$ = 1.5395	$\mathbb{H}(\mathbf{G} \mathbf{L})$ = 1.3441

$$\begin{aligned} I(\mathbf{G}; \mathbf{L}) &= \mathbb{H}(\mathbf{G}) - \mathbb{H}(\mathbf{G}|\mathbf{L}) \\ &= 1.4277 - 1.3441 = 0.0834 \end{aligned}$$

# Mutual Information: Definition

entropy after specific decisions

20	25	5	$\mathbb{P}(\mathbf{G} = \boxplus)$ = .5
9	18	9	$\mathbb{P}(\mathbf{G} = \diamond)$ = .36
1	7	6	$\mathbb{P}(\mathbf{G} = \boxminus)$ = .14
$\mathbb{P}(\mathbf{L} = \boxplus)$ = .3	$\mathbb{P}(\mathbf{L} = \diamond)$ = .5	$\mathbb{P}(\mathbf{L} = \boxminus)$ = .2	$\mathbb{H}(\mathbf{G})$ = 1.4277
$\mathbb{H}(\mathbf{G} \mathbf{L} = \boxplus)$ = 1.0746	$\mathbb{H}(\mathbf{G} \mathbf{L} = \diamond)$ = 1.4277	$\mathbb{H}(\mathbf{G} \mathbf{L} = \boxminus)$ = 1.5395	$\mathbb{H}(\mathbf{G} \mathbf{L})$ = 1.3441

$$\begin{aligned} I(\mathbf{G}; \mathbf{L}) &= \mathbb{H}(\mathbf{G}) - \mathbb{H}(\mathbf{G}|\mathbf{L}) \\ &= 1.4277 - 1.3441 = 0.0834 \end{aligned}$$

# Mutual Information: Definition

## relative entropy

20	25	5	$\mathbb{P}(\mathbf{G} = \boxplus)$ = .5
9	18	9	$\mathbb{P}(\mathbf{G} = \diamond)$ = .36
1	7	6	$\mathbb{P}(\mathbf{G} = \boxminus)$ = .14
$\mathbb{P}(\mathbf{L} = \boxplus)$ = .3	$\mathbb{P}(\mathbf{L} = \diamond)$ = .5	$\mathbb{P}(\mathbf{L} = \boxminus)$ = .2	$\mathbb{H}(\mathbf{G})$ = 1.4277
$\mathbb{H}(\mathbf{G} \mathbf{L} = \boxplus)$ = 1.0746	$\mathbb{H}(\mathbf{G} \mathbf{L} = \diamond)$ = 1.4277	$\mathbb{H}(\mathbf{G} \mathbf{L} = \boxminus)$ = 1.5395	$\mathbb{H}(\mathbf{G} \mathbf{L})$ = 1.3441

$$\begin{aligned} I(\mathbf{G}; \mathbf{L}) &= \mathbb{H}(\mathbf{G}) - \mathbb{H}(\mathbf{G}|\mathbf{L}) \\ &= 1.4277 - 1.3441 = 0.0834 \end{aligned}$$

# Mutual Information: Definition

## mutual information

20	25	5	$\mathbb{P}(\mathbf{G} = \boxplus)$ = .5
9	18	9	$\mathbb{P}(\mathbf{G} = \diamond)$ = .36
1	7	6	$\mathbb{P}(\mathbf{G} = \boxminus)$ = .14
$\mathbb{P}(\mathbf{L} = \boxplus)$ = .3	$\mathbb{P}(\mathbf{L} = \diamond)$ = .5	$\mathbb{P}(\mathbf{L} = \boxminus)$ = .2	$\mathbb{H}(\mathbf{G})$ = 1.4277
$\mathbb{H}(\mathbf{G} \mathbf{L} = \boxplus)$ = 1.0746	$\mathbb{H}(\mathbf{G} \mathbf{L} = \diamond)$ = 1.4277	$\mathbb{H}(\mathbf{G} \mathbf{L} = \boxminus)$ = 1.5395	$\mathbb{H}(\mathbf{G} \mathbf{L})$ = 1.3441

$$\begin{aligned}\mathbb{I}(\mathbf{G}; \mathbf{L}) &= \mathbb{H}(\mathbf{G}) - \mathbb{H}(\mathbf{G}|\mathbf{L}) \\ &= 1.4277 - 1.3441 = 0.0834\end{aligned}$$

# Mutual Information: No Bias!

20 (20)	25 (25)	5 (5)	$\mathbb{P}(\mathbf{G} = \boxplus)$ = .5
9 (9)	18 (18)	9 (9)	$\mathbb{P}(\mathbf{G} = \diamond)$ = .36
1 (1)	7 (7)	6 (6)	$\mathbb{P}(\mathbf{G} = \boxminus)$ = .14
$\mathbb{P}(\mathbf{L} = \boxplus)$ = .3 (.3)	$\mathbb{P}(\mathbf{L} = \diamond)$ = .5 (.5)	$\mathbb{P}(\mathbf{L} = \boxminus)$ = .2 (.2)	$\mathbb{H}(\mathbf{G})$ = 1.4277
$\mathbb{H}(\mathbf{G} \mathbf{L} = \boxplus)$ = 1.0746 (1.0746)	$\mathbb{H}(\mathbf{G} \mathbf{L} = \diamond)$ = 1.4277 (1.4277)	$\mathbb{H}(\mathbf{G} \mathbf{L} = \boxminus)$ = 1.5395 (1.5395)	$\mathbb{H}(\mathbf{G} \mathbf{L})$ = 1.3441 (1.3441)

$$\begin{aligned} \mathbb{I}(\mathbf{G}; \mathbf{L}) &= \mathbb{H}(\mathbf{G}) - \mathbb{H}(\mathbf{G}|\mathbf{L}) \\ &= 1.4277 - 1.3441 = 0.0834 \end{aligned}$$

# Mutual Information: No Bias!

constant choice

50 (20)	0 (25)	0 (5)	$\mathbb{P}(\mathbf{G} = \boxplus)$ = .5
36 (9)	0 (18)	0 (9)	$\mathbb{P}(\mathbf{G} = \diamond)$ = .36
14 (1)	0 (7)	0 (6)	$\mathbb{P}(\mathbf{G} = \boxminus)$ = .14
$\mathbb{P}(\mathbf{L} = \boxplus)$ = 1 (.3)	$\mathbb{P}(\mathbf{L} = \diamond)$ = .0 (.5)	$\mathbb{P}(\mathbf{L} = \boxminus)$ = .0 (.2)	$\mathbb{H}(\mathbf{G})$ = 1.4277
$\mathbb{H}(\mathbf{G} \mathbf{L} = \boxplus)$ = 1.4277 (1.0746)	$\mathbb{H}(\mathbf{G} \mathbf{L} = \diamond)$ ? (1.4277)	$\mathbb{H}(\mathbf{G} \mathbf{L} = \boxminus)$ ? (1.5395)	$\mathbb{H}(\mathbf{G} \mathbf{L})$ = 1.4277 (1.3441)

$$\begin{aligned} \mathbb{I}(\mathbf{G}; \mathbf{L}) &= \mathbb{H}(\mathbf{G}) - \mathbb{H}(\mathbf{G}|\mathbf{L}) \\ &= 1.4277 - 1.4277 = 0.0 \end{aligned}$$



# Mutual Information: No Bias!

constant choice

0 (20)	50 (25)	0 (5)	$\mathbb{P}(\mathbf{G} = \boxplus)$ = .5
0 (9)	36 (18)	0 (9)	$\mathbb{P}(\mathbf{G} = \diamond)$ = .36
0 (1)	14 (7)	0 (6)	$\mathbb{P}(\mathbf{G} = \boxminus)$ = .14
$\mathbb{P}(\mathbf{L} = \boxplus)$ = .0 (.3)	$\mathbb{P}(\mathbf{L} = \diamond)$ = 1 (.5)	$\mathbb{P}(\mathbf{L} = \boxminus)$ = .0 (.2)	$\mathbb{H}(\mathbf{G})$ = 1.4277
$\mathbb{H}(\mathbf{G} \mathbf{L} = \boxplus)$ ? (1.0746)	$\mathbb{H}(\mathbf{G} \mathbf{L} = \diamond)$ = 1.4277 (1.4277)	$\mathbb{H}(\mathbf{G} \mathbf{L} = \boxminus)$ ? (1.5395)	$\mathbb{H}(\mathbf{G} \mathbf{L})$ = 1.4277 (1.3441)

$$\begin{aligned} \mathbb{I}(\mathbf{G}; \mathbf{L}) &= \mathbb{H}(\mathbf{G}) - \mathbb{H}(\mathbf{G}|\mathbf{L}) \\ &= 1.4277 - 1.4277 = 0.0 \end{aligned}$$

# Mutual Information: No Bias!

constant choice

0 (20)	0 (25)	50 (5)	$\mathbb{P}(\mathbf{G} = \boxplus)$ = .5
0 (9)	0 (18)	35 (9)	$\mathbb{P}(\mathbf{G} = \diamond)$ = .36
0 (1)	0 (7)	14 (6)	$\mathbb{P}(\mathbf{G} = \boxminus)$ = .14
$\mathbb{P}(\mathbf{L} = \boxplus)$ = .0 (.3)	$\mathbb{P}(\mathbf{L} = \diamond)$ = .0 (.5)	$\mathbb{P}(\mathbf{L} = \boxminus)$ = 1 (.2)	$\mathbb{H}(\mathbf{G})$ = 1.4277
$\mathbb{H}(\mathbf{G} \mathbf{L} = \boxplus)$ ? (1.0746)	$\mathbb{H}(\mathbf{G} \mathbf{L} = \diamond)$ ? (1.4277)	$\mathbb{H}(\mathbf{G} \mathbf{L} = \boxminus)$ = 1.4277 (1.5395)	$\mathbb{H}(\mathbf{G} \mathbf{L})$ = 1.4277 (1.3441)

$$\begin{aligned} \mathbb{I}(\mathbf{G}; \mathbf{L}) &= \mathbb{H}(\mathbf{G}) - \mathbb{H}(\mathbf{G}|\mathbf{L}) \\ &= 1.4277 - 1.4277 = 0.0 \end{aligned}$$

# Mutual Information: No Bias!

random choice

25 (20)	18 (25)	7 (5)	$\mathbb{P}(\mathbf{G} = \boxplus)$ = .5
18 (9)	13 (18)	5 (9)	$\mathbb{P}(\mathbf{G} = \diamond)$ = .36
7 (1)	5 (7)	2 (6)	$\mathbb{P}(\mathbf{G} = \boxminus)$ = .14
$\mathbb{P}(\mathbf{L} = \boxplus)$ = .5 (.3)	$\mathbb{P}(\mathbf{L} = \diamond)$ = .36 (.5)	$\mathbb{P}(\mathbf{L} = \boxminus)$ = .14 (.2)	$\mathbb{H}(\mathbf{G})$ = 1.4277
$\mathbb{H}(\mathbf{G} \mathbf{L} = \boxplus)$ = 1.4277 (1.0746)	$\mathbb{H}(\mathbf{G} \mathbf{L} = \diamond)$ = 1.4277 (1.4277)	$\mathbb{H}(\mathbf{G} \mathbf{L} = \boxminus)$ = 1.4277 (1.5395)	$\mathbb{H}(\mathbf{G} \mathbf{L})$ = 1.4277 (1.3441)

$$\begin{aligned} \mathbb{I}(\mathbf{G}; \mathbf{L}) &= \mathbb{H}(\mathbf{G}) - \mathbb{H}(\mathbf{G}|\mathbf{L}) \\ &= 1.4277 - 1.4277 = 0.0 \end{aligned}$$

# Mutual Information: No Bias!

(20)	(25)	(5)	$\mathbb{P}(\mathbf{G} = \boxplus)$ = .5
(9)	(18)	(9)	$\mathbb{P}(\mathbf{G} = \diamond)$ = .36
(1)	(7)	(6)	$\mathbb{P}(\mathbf{G} = \boxminus)$ = .14
$\mathbb{P}(\mathbf{L} = \boxplus)$	$\mathbb{P}(\mathbf{L} = \diamond)$	$\mathbb{P}(\mathbf{L} = \boxminus)$	$\mathbb{H}(\mathbf{G})$ = 1.4277
(.3)	(.5)	(.2)	
$\mathbb{H}(\mathbf{G} \mathbf{L} = \boxplus)$	$\mathbb{H}(\mathbf{G} \mathbf{L} = \diamond)$	$\mathbb{H}(\mathbf{G} \mathbf{L} = \boxminus)$	$\mathbb{H}(\mathbf{G} \mathbf{L})$
(1.0746)	(1.4277)	(1.5395)	(1.3441)

$$\begin{aligned} \mathbb{I}(\mathbf{G}; \mathbf{L}) &= \mathbb{H}(\mathbf{G}) - \mathbb{H}(\mathbf{G}|\mathbf{L}) \\ &= 1.4277 - 1.3441 = 0.0834 \end{aligned}$$

# Mutual Information: Degradation Problem

## degradation

20 (20)	25 (25)	5 (5)	$\mathbb{P}(\mathbf{G} = \boxplus)$ = .5
9 (9)	18 (18)	9 (9)	$\mathbb{P}(\mathbf{G} = \diamond)$ = .36
1 (1)	7 (7)	6 (6)	$\mathbb{P}(\mathbf{G} = \boxminus)$ = .14
$\mathbb{P}(\mathbf{L} = \boxplus)$ = .3 (.3)	$\mathbb{P}(\mathbf{L} = \diamond)$ = .5 (.5)	$\mathbb{P}(\mathbf{L} = \boxminus)$ = .2 (.2)	$\mathbb{A}_3(\mathbf{L}; \mathbf{G})$ = .44 (.44)
$\mathbb{H}(\mathbf{G} \mathbf{L} = \boxplus)$ = 1.0746 (1.0746)	$\mathbb{H}(\mathbf{G} \mathbf{L} = \diamond)$ = 1.4277 (1.4277)	$\mathbb{H}(\mathbf{G} \mathbf{L} = \boxminus)$ = 1.5395 (1.5395)	$\mathbb{H}(\mathbf{G} \mathbf{L})$ = 1.3441 (1.3441)

$$\begin{aligned} \mathbb{I}(\mathbf{G}; \mathbf{L}) &= \mathbb{H}(\mathbf{G}) - \mathbb{H}(\mathbf{G}|\mathbf{L}) \\ &= 1.4277 - 1.3441 = 0.0834 \quad (0.0834) \end{aligned}$$

# Mutual Information: Degradation Problem

## degradation

45 (20)	0 (25)	5 (5)	$\mathbb{P}(\mathbf{G} = \boxplus)$ = .5
27 (9)	0 (18)	9 (9)	$\mathbb{P}(\mathbf{G} = \diamond)$ = .36
8 (1)	0 (7)	6 (6)	$\mathbb{P}(\mathbf{G} = \boxminus)$ = .14
$\mathbb{P}(\mathbf{L} = \boxplus)$ = .8 (.3)	$\mathbb{P}(\mathbf{L} = \diamond)$ = .0 (.5)	$\mathbb{P}(\mathbf{L} = \boxminus)$ = .2 (.2)	$\mathbb{A}_3(\mathbf{L}; \mathbf{G})$ = .51 (.44)
$\mathbb{H}(\mathbf{G} \mathbf{L} = \boxplus)$ = 1.3280 (1.0746)	$\mathbb{H}(\mathbf{G} \mathbf{L} = \diamond)$ ? (1.4277)	$\mathbb{H}(\mathbf{G} \mathbf{L} = \boxminus)$ = 1.5395 (1.5395)	$\mathbb{H}(\mathbf{G} \mathbf{L})$ = 1.3703 (1.3441)

$$\begin{aligned} \mathbb{I}(\mathbf{G}; \mathbf{L}) &= \mathbb{H}(\mathbf{G}) - \mathbb{H}(\mathbf{G}|\mathbf{L}) \\ &= 1.4277 - 1.3703 = 0.0262 \quad (0.0834) \end{aligned}$$

# Mutual Information: Degradation Problem

## degradation

45 (20)	0 (25)	5 (5)	$\mathbb{P}(\mathbf{G} = \boxplus)$ = .5
27 (9)	0 (18)	9 (9)	$\mathbb{P}(\mathbf{G} = \diamond)$ = .36
8 (1)	0 (7)	6 (6)	$\mathbb{P}(\mathbf{G} = \boxminus)$ = .14
$\mathbb{P}(\mathbf{L} = \boxplus)$ = .8 (.3)	$\mathbb{P}(\mathbf{L} = \diamond)$ = .0 (.5)	$\mathbb{P}(\mathbf{L} = \boxminus)$ = .2 (.2)	$\mathbb{A}_3(\mathbf{L}; \mathbf{G})$ = .51 (.44)
$\mathbb{H}(\mathbf{G} \mathbf{L} = \boxplus)$ = 1.3280 (1.0746)	$\mathbb{H}(\mathbf{G} \mathbf{L} = \diamond)$ ? (1.4277)	$\mathbb{H}(\mathbf{G} \mathbf{L} = \boxminus)$ = 1.5395 (1.5395)	$\mathbb{H}(\mathbf{G} \mathbf{L})$ = 1.3703 (1.3441)

$$\begin{aligned} \mathbb{I}(\mathbf{G}; \mathbf{L}) &= \mathbb{H}(\mathbf{G}) - \mathbb{H}(\mathbf{G}|\mathbf{L}) \\ &= 1.4277 - 1.3703 = 0.0262 \quad (0.0834) \end{aligned}$$

# Mutual Information vs. Accuracy: Degradation

In response to degradation:

- ▶ MI: .0834  $\rightarrow$  .0262
- ▶ Acc: .44  $\rightarrow$  .51



# Outline

The Structure of RTE Data

Accuracy

Average Precision

Mutual Information

# Final Recommendations

- ▶ in addition to accuracy, report MI,
- ▶ use MI for ranking,
- ▶ optimize systems for MI,
- ▶ drop average precision! use confidence-ranked MI and/or bring back CWS
- ▶ be aware of bias and baseline scores when looking at accuracies.

# Final Recommendations

- ▶ in addition to accuracy, report MI,
- ▶ use MI for ranking,
- ▶ optimize systems for MI,
- ▶ drop average precision! use confidence-ranked MI and/or bring back CWS
- ▶ be aware of bias and baseline scores when looking at accuracies.

# Final Recommendations

- ▶ in addition to accuracy, report MI,
- ▶ use MI for ranking,
- ▶ optimize systems for MI,
- ▶ drop average precision! use confidence-ranked MI and/or bring back CWS
- ▶ be aware of bias and baseline scores when looking at accuracies.

# Final Recommendations

- ▶ in addition to accuracy, report MI,
- ▶ use MI for ranking,
- ▶ optimize systems for MI,
- ▶ drop average precision! use confidence-ranked MI and/or bring back CWS
- ▶ be aware of bias and baseline scores when looking at accuracies.

# Final Recommendations

- ▶ in addition to accuracy, report MI,
- ▶ use MI for ranking,
- ▶ optimize systems for MI,
- ▶ drop average precision! use confidence-ranked MI and/or bring back CWS
- ▶ be aware of bias and baseline scores when looking at accuracies.



