Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

# Generative models for natural language inference
## DGM4NLP

Miguel Rios
University of Amsterdam

May 12, 2019

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

Applications of Textual Entailment

# Outline

1 **Introduction**
   ● Applications of Textual Entailment

2 Levels of Representation

3 RTE Methods
   ● Evaluation

4 Current Methods

5 Latent Variable Models

6 Uncertainty in Natural Language Inference

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

Applications of Textual Entailment

## Introduction

- Textual entailment is defined as a **directional relation** between pairs of text expressions, the **T** "Text", and the **H** "Hypothesis".

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

Applications of Textual Entailment

## Introduction

- Textual entailment is defined as a **directional relation** between pairs of text expressions, the **T** "Text", and the **H** "Hypothesis".
- We say that T entails H if the meaning of H can be inferred from the meaning of T, as would typically be interpreted by people.
  T → H

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

Applications of Textual Entailment

## Introduction

- Textual entailment is defined as a **directional relation**
  between pairs of text expressions, the **T** "Text", and the **H**
  "Hypothesis".
- We say that T entails H if the meaning of H can be inferred
  from the meaning of T, as would typically be interpreted by
  people.
  T $\rightarrow$ H

> T: The purchase of Houston-based LexCorp by
> BMI for \$2Bn prompted widespread sell-offs by
> traders as they sought to minimize exposure.
> H: BMI acquired an American company.

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

Applications of Textual Entailment

## Recognising Textual Entailment

- Recognition: identification of a thing or person from **previous** encounters or **knowledge**.

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

Applications of Textual Entailment

## Recognising Textual Entailment

- Recognition: identification of a thing or person from **previous** encounters or **knowledge**.

- Physicians are trained in medicine to recognise and treat a disease.

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

Applications of Textual Entailment

## Recognising Textual Entailment

- RTE Challenge (Dagan and Glickman, 2005), provides the first benchmark.

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

Applications of Textual Entailment

## Recognising Textual Entailment

- RTE Challenge (Dagan and Glickman, 2005), provides the first benchmark.
- Participant methods decide for each entailment pair whether T entails H or not.

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

Applications of Textual Entailment

## Recognising Textual Entailment

- RTE Challenge (Dagan and Glickman, 2005), provides the first benchmark.
- Participant methods decide for each entailment pair whether T entails H or not.
- The annotation used for the entailment decision is **TRUE** if T entails H or **FALSE** otherwise.

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

Applications of Textual Entailment

## Recognising Textual Entailment

- RTE Challenge (Dagan and Glickman, 2005), provides the first benchmark.
- Participant methods decide for each entailment pair whether T entails H or not.
- The annotation used for the entailment decision is **TRUE** if T entails H or **FALSE** otherwise.

> RTE can be **framed** as a classification problem, where the entailment relations are the classes, and the RTE benchmark provides the essential evidence to build a **supervised binary classifier** (Dagan et al., 2010)

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

Applications of Textual Entailment

## Applications of Textual Entailment

- RTE has been proposed as a **generic task that captures major semantic inference needs** across natural language processing applications.

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

Applications of Textual Entailment

## Applications of Textual Entailment

- RTE has been proposed as a **generic task that captures major semantic inference needs** across natural language processing applications.

- We can frame natural language processing tasks as recognition.
  Input as T and generated output as H.

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

Applications of Textual Entailment

## Question Answering

- Question Answering system generates as output the best candidate answers. While the top candidate may not be the correct answer, the correct answer is in the set of returned candidates.

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

Applications of Textual Entailment

## Question Answering

- Question Answering system generates as output the best candidate answers. While the top candidate may not be the correct answer, the correct answer is in the set of returned candidates.

  > T/Q: Arabic, for example, is used densely across North Africa and from the Eastern Mediterranean to the Philippines, as the key language of the Arab world.
  > H/A: Arabic is the primary language of the Philippines.

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

Applications of Textual Entailment

# Summarisation

- Identifying if a new sentence contains information already by a summary-in-progress (redundancy detection) can be framed as the current summary as T and the new sentence as H.

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

Applications of Textual Entailment

## Summarisation

- Identifying if a new sentence contains information already by a
  summary-in-progress (redundancy detection) can be framed as
  the current summary as T and the new sentence as H.

  > T/S1: Google and NASA announced a working
  > agreement, Wednesday, that could result in the
  > Internet giant building a complex of up to 1 million
  > square feet on NASA-owned property, adjacent to
  > Moffett Field, near Mountain View.
  > H/S2: Google may build a campus on NASA property.

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

# Outline

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Challenge of RTE

> T: The purchase of Houston-based LexCorp by BMI for $2Bn prompted widespread sell-offs by traders as they sought to minimize exposure.
> H: BMI acquired an American company.

To recognise **TRUE** entailment relation:

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Challenge of RTE

> T: The purchase of Houston-based LexCorp by
> BMI for $2Bn prompted widespread sell-offs by
> traders as they sought to minimize exposure.
> H: BMI acquired an American company.

To recognise **TRUE** entailment relation:

- "company" in the Hypothesis can match "LexCorp",

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Challenge of RTE

> T: The purchase of Houston-based LexCorp by
> BMI for \$2Bn prompted widespread sell-offs by
> traders as they sought to minimize exposure.
> H: BMI acquired an American company.

To recognise **TRUE** entailment relation:

- "company" in the Hypothesis can match "LexCorp",
- "based in Houston" implies "American",

Introduction
**Levels of Representation**
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Challenge of RTE

> T: The purchase of Houston-based LexCorp by
> BMI for \$2Bn prompted widespread sell-offs by
> traders as they sought to minimize exposure.
> H: BMI acquired an American company.

To recognise **TRUE** entailment relation:

- "company" in the Hypothesis can match "LexCorp",
- "based in Houston" implies "American",
- identify the relation "purchase",

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Challenge of RTE

> T: The purchase of Houston-based LexCorp by
> BMI for \$2Bn prompted widespread sell-offs by
> traders as they sought to minimize exposure.
> H: BMI acquired an American company.

To recognise **TRUE** entailment relation:

- "company" in the Hypothesis can match "LexCorp",
- "based in Houston" implies "American",
- identify the relation "purchase",
- determine that "A purchased by B" implies "B acquires A".

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Levels of Representation

- Determining the equivalence or non-equivalence of the meanings of the T-H.

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Levels of Representation

- Determining the equivalence or non-equivalence of the meanings of the T-H.
- The representation (e.g. words, syntax, semantics) of the T-H pair that is used to extract features to train a supervised classifier.

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Lexical level

- Every assertion (word) in the representation of H is contained in the representation T.

Introduction
**Levels of Representation**
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Lexical level

- H and T sentences encode aspects of underlying meaning that cannot be captured by the purely lexical representation.

Introduction
**Levels of Representation**
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Structural level

- Syntactic structure provides cues for the underlying meaning of a sentence.

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Structural level

- If T contains the same structure (i.e, dependency edges), the system will predict TRUE and otherwise FALSE.

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Structural level

- If T contains the same structure (i.e, dependency edges), the system will predict TRUE and otherwise FALSE.
- "John" and "drove," but the two words are **separated** by a sequence of dependency edges.

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Structural level

- If T contains the same structure (i.e, dependency edges), the system will predict TRUE and otherwise FALSE.
- "John" and "drove," but the two words are **separated** by a sequence of dependency edges.
- Given the expressiveness of the dependency representation, many possible sequences of edges that could represent connection, and many other sequences that do not.

Introduction
**Levels of Representation**
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Semantic level

- Semantic role labelling, grouping of words into "arguments" (entity such as a person or place) and "predicates" (a predicate being a verb representing the state of some entity).

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Semantic level

- Semantic role labelling, grouping of words into "arguments" (entity such as a person or place) and "predicates" (a predicate being a verb representing the state of some entity).
- **Immediate connections** between arguments and predicates.

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Semantic level

- Semantic role labelling, grouping of words into "arguments" (entity such as a person or place) and "predicates" (a predicate being a verb representing the state of some entity).
- **Immediate connections** between arguments and predicates.
- "John" is an argument of the predicate "drove"

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Knowledge Acquisition for RTE

- T: The U.S. citizens elected their new president Obama.
  H: Obama was born in the U.S.

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Knowledge Acquisition for RTE

- T: The U.S. citizens elected their new president Obama.
  H: Obama was born in the U.S.

- Assumed **background knowledge**: "U.S. presidents should be naturally born in the U.S."

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Knowledge Acquisition for RTE

- Knowledge is a lexical-semantic relation between two words.

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Knowledge Acquisition for RTE

- Knowledge is a lexical-semantic relation between two words.
- I enlarged my **stock**. and I enlarged my **inventory**.
  **synonym**

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Knowledge Acquisition for RTE

- Knowledge is a lexical-semantic relation between two words.
- I enlarged my **stock**. and I enlarged my **inventory**.
  **synonym**
- I have a **cat**. entails I have a **pet**.
  **hyponymy**

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Knowledge Acquisition for RTE

- Knowledge is a lexical-semantic relation between two words.
- I enlarged my **stock**. and I enlarged my **inventory**.
  **synonym**
- I have a **cat**. entails I have a **pet**.
  **hyponymy**
- But also meaning implication between more complex structures than just lexical terms.
  X *causes* Y $\rightarrow$ Y *is a symptom of* X

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Knowledge Acquisition for RTE

- WordNet specifies **lexical-semantic** relations between lexical items such as hyponymy, synonymy, and derivation.
  chair $\rightarrow$ furniture

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Knowledge Acquisition for RTE

- WordNet specifies **lexical-semantic** relations between lexical items such as hyponymy, synonymy, and derivation.
  chair $\rightarrow$ furniture

- FrameNet is a lexicographic resource for **frames** that are events and includes information on the predicates and argument relevant for that specific event.
  The attack frame, and specifies events: 'assailant', a 'victim', a 'weapon', etc.
  cure X $\rightarrow$ X recover

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Knowledge Acquisition for RTE

- WordNet specifies **lexical-semantic** relations between lexical items such as hyponymy, synonymy, and derivation.
  chair $\rightarrow$ furniture

- FrameNet is a lexicographic resource for **frames** that are events and includes information on the predicates and argument relevant for that specific event.
  The attack frame, and specifies events: 'assailant', a 'victim', a 'weapon', etc.
  cure X $\rightarrow$ X recover

- Wikipedia articles for identifying **is a** relations.
  Jim Carrey $\rightarrow$ actor

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

Knowledge Acquisition for RTE

- **Extended Distributional Hypothesis**: If two paths tend to occur in similar contexts, the meanings of the paths tend to be similar.

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Knowledge Acquisition for RTE

- **Extended Distributional Hypothesis**: If two paths tend to occur in similar contexts, the meanings of the paths tend to be similar.

- X solves Y
  Y is solved by X
  X finds a solution to Y

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

# Knowledge Acquisition for RTE



They had previously bought bighorn sheep from Comstock.

| | | |
|---|---|---|
| (a) | N:subj:V ← buy → V:from:N | ($X$ buys something from $Y$) |
| (b) | N:subj:V ← buy → V:obj:N | ($X$ buys $Y$) |
| (c) | N:subj:V ← buy → V:obj:N → sheep → N:nn:N | ($X$ buys $Y$ sheep) |
| (d) | N:nn:N ← sheep ← N:obj:V ← buy → V:from:N | ($X$ sheep is bought from $Y$) |
| (e) | N:obj:V ← buy → V:from:N | ($X$ is bought from $Y$) |

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

Evaluation

# Outline

1. Introduction
   - Applications of Textual Entailment

2. Levels of Representation

3. RTE Methods
   - Evaluation

4. Current Methods

5. Latent Variable Models

6. Uncertainty in Natural Language Inference

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

Evaluation

## Recognising Textual Entailment Methods

- RTE depend on the representation (e.g. words, syntax, semantics) of the T-H pair that is used to extract features to train a supervised classifier.

**Text**

| John | Smith | drove | to | Seattle | and | bought | a | Honda | Civic |

Drive($E_{T1}$, John Smith, Seattle)
Buy($E_{T2}$, John Smith, a Honda Civic)

**Hypothesis**

| John | Smith | drove | to | Seattle |

Drive($E_{H1}$, John Smith, Seattle)

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

Evaluation

# Recognising Textual Entailment Methods

Introduction
Levels of Representation
**RTE Methods**
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

Evaluation

## Similarity-based approaches

- Pair with a strong similarity score holds a positive entailment relation.

Introduction
Levels of Representation
**RTE Methods**
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

Evaluation

## Similarity-based approaches

- Pair with a strong similarity score holds a positive entailment relation.
- Wordnet similarity.

Introduction
Levels of Representation
**RTE Methods**
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

Evaluation

## Similarity-based approaches

- Pair with a strong similarity score holds a positive entailment relation.
- Wordnet similarity.
- String similarity.

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

Evaluation

## Similarity-based approaches

- Pair with a strong similarity score holds a positive entailment relation.
- Wordnet similarity.
- String similarity.
- Similarity scores computed from different linguistic levels. The goal is to find complementary features.

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

Evaluation

## Alignment-based approaches



- (1,purchase,acquired)
  (3,Hudson-based LexCorp, American company),
  (5,BMI,BMI)

- $\rho_4$ = purchase of $\boxed{X}$ by $\boxed{Y}$ $\rightarrow$ $\boxed{Y}$ acquired $\boxed{X}$

- $\rho_5$ = $\boxed{Z:Noun}$ of $\boxed{X}$ by $\boxed{Y}$ $\rightarrow$ $\boxed{Y}$ $\boxed{Z:Verb}$ $\boxed{X}$

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

Evaluation

# Alignment-based approaches

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

Evaluation

## Edit distance-based approaches

- T entails H if there is a **sequence of transformations** applied to T such that we can obtain H with an overall cost below a certain **threshold**.

Introduction
Levels of Representation
**RTE Methods**
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

Evaluation

## Edit distance-based approaches

- T entails H if there is a **sequence of transformations** applied to T such that we can obtain H with an overall cost below a certain **threshold**.

- Insertion, Substitution, and Deletion.

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

Evaluation

## Edit distance-based approaches

- T entails H if there is a **sequence of transformations** applied to T such that we can obtain H with an overall cost below a certain **threshold**.

- Insertion, Substitution, and Deletion.

- Alternative for expensive theorem provers.

## Evaluation

- Accuracy

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

Evaluation

# Evaluation

- Accuracy
- RTE-3 corpus **1,600** T-H pairs
  information extraction, information retrieval, question
  answering, and summarisation.

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

Evaluation

## Evaluation

- Accuracy

- RTE-3 corpus **1,600** T-H pairs
  information extraction, information retrieval, question
  answering, and summarisation.

- The lexical baseline, between 55% and 58% **accuracy**

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

Evaluation

# Evaluation

- Accuracy

- RTE-3 corpus **1,600** T-H pairs
  information extraction, information retrieval, question
  answering, and summarisation.

- The lexical baseline, between 55% and 58% **accuracy**

- RTE-3 higher scores all system entries suggesting an **easier**
  entailment corpus

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

Evaluation

# Evaluation

- Accuracy
- RTE-3 corpus **1,600** T-H pairs
  information extraction, information retrieval, question
  answering, and summarisation.
- The lexical baseline, between 55% and 58% **accuracy**
- RTE-3 higher scores all system entries suggesting an **easier**
  entailment corpus
- RTE-4 and RTE-5 increase the difficulty by adding irrelevant
  signals (additional words, phrases, and sentences).

Introduction
Levels of Representation
RTE Methods
**Current Methods**
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Outline

1 Introduction
   ● Applications of Textual Entailment

2 Levels of Representation

3 RTE Methods
   ● Evaluation

4 Current Methods

5 Latent Variable Models

6 Uncertainty in Natural Language Inference

Introduction
Levels of Representation
RTE Methods
**Current Methods**
Latent Variable Models
Uncertainty in Natural Language Inference
References

## SNLI

- Flickr30k corpus for image captioning domaim.
  Annotated pairs of texts at sentence level

Introduction
Levels of Representation
RTE Methods
**Current Methods**
Latent Variable Models
Uncertainty in Natural Language Inference
References

## SNLI

- Flickr30k corpus for image captioning domaim.
  Annotated pairs of texts at sentence level

- The relations (i.e. 3-way classification labels) are:
  *entailment*, *contradiction*, and *neutral*.

Introduction
Levels of Representation
RTE Methods
**Current Methods**
Latent Variable Models
Uncertainty in Natural Language Inference
References

## SNLI

- Flickr30k corpus for image captioning domaim. Annotated pairs of texts at sentence level

- The relations (i.e. 3-way classification labels) are: *entailment*, *contradiction*, and *neutral*.

- $550, 152$ training, 10K development, and 10k test.

Introduction
Levels of Representation
RTE Methods
**Current Methods**
Latent Variable Models
Uncertainty in Natural Language Inference
References

## SNLI

- Flickr30k corpus for image captioning domaim. Annotated pairs of texts at sentence level

- The relations (i.e. 3-way classification labels) are: *entailment*, *contradiction*, and *neutral*.

- 550, 152 training, 10K development, and 10k test.

- Premise: A soccer game with multiple males playing. Hypothesis: Some men are playing a sport.

Introduction
Levels of Representation
RTE Methods
**Current Methods**
Latent Variable Models
Uncertainty in Natural Language Inference
References

# MNLI

- Multiple genres
  classifiers only learn **regularities** over annotated data, leading
  to **poor generalization** beyond the domain of the training
  data

Introduction
Levels of Representation
RTE Methods
**Current Methods**
Latent Variable Models
Uncertainty in Natural Language Inference
References

# MNLI

- Multiple genres
  classifiers only learn **regularities** over annotated data, leading
  to **poor generalization** beyond the domain of the training
  data

- *matched* (5 in domain genres) $392, 702$ training,
  10k *matched* development,

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

# MNLI

- Multiple genres
  classifiers only learn **regularities** over annotated data, leading
  to **poor generalization** beyond the domain of the training
  data
- *matched* (5 in domain genres) $392, 702$ training,
  10k *matched* development,
- 10k *mismatched* (5 out of domain genres) development.

Introduction
Levels of Representation
RTE Methods
**Current Methods**
Latent Variable Models
Uncertainty in Natural Language Inference
References

# MNLI

- Multiple genres
  classifiers only learn **regularities** over annotated data, leading to **poor generalization** beyond the domain of the training data
- *matched* (5 in domain genres) $392, 702$ training,
  10k *matched* development,
- 10k *mismatched* (5 out of domain genres) development.
- T: 8 million in relief in the form of emergency housing.
  H: The 8 million dollars for emergency housing was still not enough to solve the problem.
  **Government**

Introduction
Levels of Representation
RTE Methods
**Current Methods**
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Drawbacks



| **Premise** | A woman selling bamboo sticks talking to two men on a loading dock. |
|---|---|
| **Entailment** | There are **at least** three **people** on a loading dock. |
| **Neutral** | A woman is selling bamboo sticks **to help provide for her family.** |
| **Contradiction** | A woman is **not** taking money for any of her sticks. |

- Entailment: animal, instrument, and outdoors.

Introduction
Levels of Representation
RTE Methods
**Current Methods**
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Drawbacks



| **Premise** | A woman selling bamboo sticks talking to two men on a loading dock. |
| --- | --- |
| **Entailment** | There are **at least** three **people** on a loading dock. |
| **Neutral** | A woman is selling bamboo sticks **to help provide for her family.** |
| **Contradiction** | A woman is **not** taking money for any of her sticks. |

- Entailment: animal, instrument, and outdoors.
- Neutral: Modifiers (tall, sad, popular) and superlatives (first, favorite, most)

Introduction
Levels of Representation
RTE Methods
**Current Methods**
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Drawbacks



| **Premise** | A woman selling bamboo sticks talking to two men on a loading dock. |
|---|---|
| **Entailment** | There are **at least** three **people** on a loading dock. |
| **Neutral** | A woman is selling bamboo sticks **to help provide for her family.** |
| **Contradiction** | A woman is **not** taking money for any of her sticks. |

- Entailment: animal, instrument, and outdoors.
- Neutral: Modifiers (tall, sad, popular) and superlatives (first, favorite, most)
- Contradiction: Negation words, nobody, no, never and nothing.

Introduction
Levels of Representation
RTE Methods
**Current Methods**
Latent Variable Models
Uncertainty in Natural Language Inference
References

# Neural Network Models

- Embeddings like glove or elmo, for fine tuning.



(a) A conventional sequence-based RNN for two sentences.



]

(b) A conventional TreeRNN for two sentences.

Introduction
Levels of Representation
RTE Methods
**Current Methods**
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Neural Network Models

- Embeddings like glove or elmo, for fine tuning.
- Sentence representations.



(a) A conventional sequence-based RNN for two sentences.



(b) A conventional TreeRNN for two sentences.

Introduction
Levels of Representation
RTE Methods
**Current Methods**
Latent Variable Models
Uncertainty in Natural Language Inference
References

# BiLSMT composition

Introduction
Levels of Representation
RTE Methods
**Current Methods**
Latent Variable Models
Uncertainty in Natural Language Inference
References

# ESIM

Introduction
Levels of Representation
RTE Methods
**Current Methods**
Latent Variable Models
Uncertainty in Natural Language Inference
References

## ESIM

$$\mathbf{t}_i = \text{emb}(t_i; \omega_{\text{emb}}) \tag{1a}$$

$$\mathbf{h}_j = \text{emb}(h_j; \omega_{\text{emb}}) \tag{1b}$$

$$\mathbf{s}_1^m = \text{birnn}(\mathbf{t}_1^m; \omega_{\text{enc}}) \tag{1c}$$

$$\mathbf{u}_1^n = \text{birnn}(\mathbf{h}_1^n; \omega_{\text{enc}}) \tag{1d}$$

$$\mathbf{a}_i = \text{attention}(\mathbf{s}_i, \mathbf{u}_1^n) \tag{1e}$$

$$\mathbf{b}_j = \text{attention}(\mathbf{u}_j, \mathbf{s}_1^m) \tag{1f}$$

$$\mathbf{c}_i = [\mathbf{s}_i, \mathbf{a}_i, \mathbf{s}_i - \mathbf{a}_i, \mathbf{s}_i \odot \mathbf{a}_i] \tag{1g}$$

$$\mathbf{d}_j = [\mathbf{u}_j, \mathbf{b}_j, \mathbf{u}_j - \mathbf{b}_j, \mathbf{u}_j \odot \mathbf{b}_j] \tag{1h}$$

$$\mathbf{c}_1^m = \text{birnn}(\mathbf{c}_1^m; \omega_{\text{comp}}) \tag{1i}$$

$$\mathbf{d}_1^n = \text{birnn}(\mathbf{d}_1^n; \omega_{\text{comp}}) \tag{1j}$$

$$\mathbf{q} = [\text{avg}(\mathbf{c}_1^m), \text{maxpool}(\mathbf{c}_1^m), \text{avg}(\mathbf{d}_1^n), \text{maxpool}(\mathbf{d}_1^n)] \tag{1k}$$

$$\mathbf{q} = \tanh(\text{affine}(\mathbf{q}; \omega_{\text{hid}})) \tag{1l}$$

$$f(x) = \text{softmax}(\text{mlp}(\mathbf{q}; \omega_{\text{cls}})) \tag{1m}$$

Introduction
Levels of Representation
RTE Methods
Current Methods
**Latent Variable Models**
Uncertainty in Natural Language Inference
References

## Outline

1. Introduction
   - Applications of Textual Entailment

2. Levels of Representation

3. RTE Methods
   - Evaluation

4. Current Methods

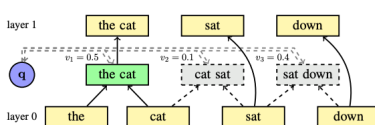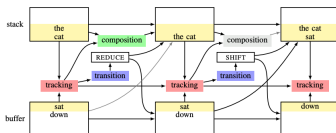5. Latent Variable Models

6. Uncertainty in Natural Language Inference

Introduction
Levels of Representation
RTE Methods
Current Methods
**Latent Variable Models**
Uncertainty in Natural Language Inference
References

# Latent Structure Induction



(a) Two parse trees correspond to two distinct interpretations for the sentence in example (1).

(b) Parses generated by at ST-Gumbel model (left) and the Stanford Parser (right).

Introduction
Levels of Representation
RTE Methods
Current Methods
**Latent Variable Models**
Uncertainty in Natural Language Inference
References

## Deep Generative Models

- Model that generates **hypothesis** and **decision** given a text and a stochastic embedding of the hypothesis-decision pair.

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Deep Generative Models

- Model that generates **hypothesis** and **decision** given a text and a stochastic embedding of the hypothesis-decision pair.
- Models to learn from mixed-domain NLI data
  e.g. by capitalising on lexical domain-dependent patterns.

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Deep Generative Models

- Model that generates **hypothesis** and **decision** given a text and a stochastic embedding of the hypothesis-decision pair.

- Models to learn from mixed-domain NLI data
  e.g. by capitalising on lexical domain-dependent patterns.

- Performance of standard classifiers tend to vary across domains and especially out of domain.

Introduction
Levels of Representation
RTE Methods
Current Methods
**Latent Variable Models**
Uncertainty in Natural Language Inference
References

## Deep Generative Models



$$Z_i | t_1^m \sim \mathcal{N}(\mu(s_1^m), \sigma^2(s_1^m))$$
$$H_i | z_1^m \sim Cat(f(z_1^m, t_1^m; \theta))$$
$$D_j | z_1^m, h_1^n \sim Cat(g(z_1^m, t_1^m, h_1^n; \theta))$$

Introduction
Levels of Representation
RTE Methods
Current Methods
**Latent Variable Models**
Uncertainty in Natural Language Inference
References

## Deep Generative Models I

- Joint likelihood of y (hypothesis) and d (decision)

$$p(y, d | x, \theta) =$$
$$\int p(z | x, \theta) p(y | x, z, \theta) p(d | x, y, z, \theta) dz. \tag{2}$$

- The *hypothesis generation model*:

$$p(y | x, z, \theta) = \prod_{j=1}^{|y|} p(y_j | x, z, y_{<j}, \theta)$$
$$= \prod_{j=1}^{|y|} \text{Cat}(y_j | f_o(x, z, y_{<j}; \theta)) , \tag{3}$$

Introduction
Levels of Representation
RTE Methods
Current Methods
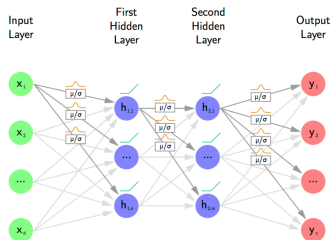**Latent Variable Models**
Uncertainty in Natural Language Inference
References

## Deep Generative Models II

- The *classification model* ESIM:

$$p(d|x, y, z, \theta) = \text{Cat}(d|f_c(x, y, z; \theta)) \qquad (4)$$

- Lowerbound on the log-likelihood function (ELBO)

$$\begin{aligned}
\mathcal{L}(\theta, \phi) = \; &\mathbb{E}_{q(z|x,y,d,\phi)} \left[ \log p(y, d|x, z, \theta) \right] \\
&- \text{KL}(q(z|x, y, d, \phi) || p(z|x, \theta))
\end{aligned} \qquad (5)$$

Introduction
Levels of Representation
RTE Methods
Current Methods
**Latent Variable Models**
Uncertainty in Natural Language Inference
References

## Deep Generative Models

| Model | Dev | |
|---|---|---|
| | matched | mismatched |
| $ESIM_{mnli}$ | $74.39 \pm 0.11$ | $74.05 \pm 0.21$ |
| $+ \mathcal{N}\text{-VAE}_{50z}$ | $74.89 \pm 0.25$ | $74.07 \pm 0.37$ |
| $+ \mathcal{N}\text{-VAE}_{100z}$ | $74.82 \pm 0.28$ | $73.91 \pm 0.59$ |
| $+ \mathcal{N}\text{-VAE}_{256z}$ | $74.87 \pm 0.15$ | $74.08 \pm 0.16$ |

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

# Outline

1. Introduction
   - Applications of Textual Entailment

2. Levels of Representation

3. RTE Methods
   - Evaluation

4. Current Methods

5. Latent Variable Models

6. Uncertainty in Natural Language Inference

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
**Uncertainty in Natural Language Inference**
References

## Bayes by backprop

- NNs perform well with lots of data, however they fail to express uncertainty with little or no data, leading to overconfident decisions.

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Bayes by backprop

- NNs perform well with lots of data, however they fail to express uncertainty with little or no data, leading to overconfident decisions.
- Bayesian neural networks introduce probability distributions over the weights.

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Bayes by backprop

- However, Bayesian inference on the parameters $\omega$ of a neural network is intractable, with data $D$.

$$p(\omega|\mathcal{D}) = \frac{p(\mathcal{D}|\omega)p(\omega)}{p(\mathcal{D})} = \frac{p(\mathcal{D}|\omega)p(\omega)}{\int p(\mathcal{D}|\omega)p(\omega)\mathrm{d}\omega} \qquad (6)$$

(Blundell et al., 2015)

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Bayes by backprop

- However, Bayesian inference on the parameters $\omega$ of a neural network is intractable, with data $D$.

$$p(\omega|\mathcal{D}) = \frac{p(\mathcal{D}|\omega)p(\omega)}{p(\mathcal{D})} = \frac{p(\mathcal{D}|\omega)p(\omega)}{\int p(\mathcal{D}|\omega)p(\omega)\mathrm{d}\omega} \qquad (6)$$

- We need an approximation $q(\omega|\theta)$, over the weights that approximates the true posterior

(Blundell et al., 2015)

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Bayes by backprop

- However, Bayesian inference on the parameters $\omega$ of a neural network is intractable, with data $D$.

$$p(\omega|\mathcal{D}) = \frac{p(\mathcal{D}|\omega)p(\omega)}{p(\mathcal{D})} = \frac{p(\mathcal{D}|\omega)p(\omega)}{\int p(\mathcal{D}|\omega)p(\omega)\mathrm{d}\omega} \tag{6}$$

- We need an approximation $q(\omega|\theta)$, over the weights that approximates the true posterior

- The ELBO is:

$$\begin{aligned} \mathcal{L}(\mathcal{D},\theta) &= \int q(\omega|\theta)\log\frac{q(\omega|\theta)}{p(\omega)} - q(\omega|\theta)\log p(\mathcal{D}|\omega)\mathrm{d}\omega \\ &= \mathrm{KL}[q(\omega|\theta)\|p(\omega)] - \mathbb{E}_{q(\omega|\theta)}[\log p(\mathcal{D}|\omega)] \end{aligned} \tag{7}$$

(Blundell et al., 2015)

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## MC dropout I

- On NLI training inputs $X = \langle (t_1, h_1), \ldots, (t_n, h_n) \rangle$ are premise ($t$) and hypothesis ($h$) pairs, and the corresponding outputs $Y = \langle y_1, \ldots, y_n \rangle$ over $N$ instances.

- The likelihood for classification is defined by:

$$p(y|x, \omega) = \text{Cat}(y|f(x; \omega)), \tag{8}$$

over $y$ entailment relations computed by mapping from the input to the class probabilities with a neural network $f$ parameterised by $\omega$.

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## MC dropout II

- A Bayesian NN (MacKay, 1992) is defined by placing a prior distribution over the model parameters $p(\omega)$, where this prior is often a Gaussian distribution $p(\omega) \sim \mathcal{N}(0, I)$.

- The Bayesian NN formulation leads to a posterior distribution over the parameters given our observed data, instead of a single estimate.

- We are interested on estimating the posterior distribution over the parameters $p(\omega|\mathcal{D})$, given our observed data $X, Y$.

- The goal is to predict a new input instances by marginalising over the parameters:

$$p(y^*|x^*, \mathcal{D}) = \int p(y^*|x^*, \omega)p(\omega|\mathcal{D})d\omega. \qquad (9)$$

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## MC dropout III

- However, the true posterior $p(\omega|\mathcal{D})$ is intractable, and Gal and Ghahramani (2016a) use variational inference to approximate this posterior.

- We define an approximate distribution $q_\theta(\omega)$, to minimise the KL divergence between the approximation and the true posterior.

- The objective for optimisation is a lower-bound on the log-likelihood function (ELBO):

$$\mathcal{L} = \mathbb{E}_{q(\omega)} \left[ \sum_{i=1}^{N} \log p(y_i | f(x_i; \omega)) \right] \tag{10}$$
$$- KL(q_\theta(\omega) || p(\omega)),$$

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## MC dropout IV

where the KL term is approximated with $L_2$ regularisation.

- Gal and Ghahramani (2016a) show that the use of dropout in NNs before each weight layer is an approximation to variational inference in Bayesian NNs.

- By replacing the true posterior $p(\omega|\mathcal{D})$ with the approximate posterior $q_\theta(\omega)$, we obtain a Monte Carlo (MC) estimate for future predictions :

$$
\begin{aligned}
p(y^*|x^*, \mathcal{D}) &\approx \int p(y^*|x^*, \omega) q_\theta(\omega) d\omega \\
&\approx \frac{1}{T} \sum_t^T p(y^*|x^*, \hat{\omega}_t),
\end{aligned}
\tag{11}
$$

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## MC dropout V

where $\hat{\omega}_t \sim q_\theta(\omega)$

- In practice, the approximation to the predictive distribution is based on performing $T$ stochastic forward passes through the network and averaging the results.

- In other words, this is achieved by performing **dropout at test time** (MC dropout).

- Finally, for classification, a way to quantify uncertainty is by computing the entropy of the output probability vector $\mathcal{H}(p) = -\sum_{c=1}^{C} p_c \log p_c$ over $c$ classes.

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
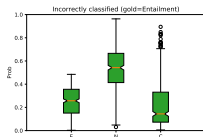Uncertainty in Natural Language Inference
References

## Uncertainty in natural language inference

- ESIM for classification (without syntactic parses)

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
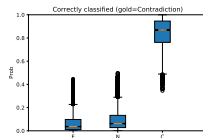Uncertainty in Natural Language Inference
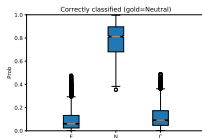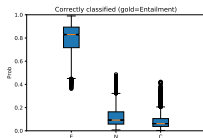References

## Uncertainty in natural language inference

- ESIM for classification (without syntactic parses)
- The word embedding and the bidirectional LSTMs are shared between the pair of texts. A single (tanh) hidden layer MLP with a softmax output predicts the class probabilities.

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Uncertainty in natural language inference

- ESIM for classification (without syntactic parses)
- The word embedding and the bidirectional LSTMs are shared between the pair of texts. A single (tanh) hidden layer MLP with a softmax output predicts the class probabilities.
- We use dropout on both the LSTM (variational RNN), and the word embedding.

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Uncertainty in natural language inference

- ESIM for classification (without syntactic parses)
- The word embedding and the bidirectional LSTMs are shared between the pair of texts. A single (tanh) hidden layer MLP with a softmax output predicts the class probabilities.
- We use dropout on both the LSTM (variational RNN), and the word embedding.
- In the word embedding $\omega_{\mathrm{emb}} \in R^{V \times D}$, with $V$ vocabulary and $D$ dimensionality, the dropout masks types (rows) instead of words in a sequence .

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Uncertainty in natural language inference

- ESIM for classification (without syntactic parses)
- The word embedding and the bidirectional LSTMs are shared between the pair of texts. A single (tanh) hidden layer MLP with a softmax output predicts the class probabilities.
- We use dropout on both the LSTM (variational RNN), and the word embedding.
- In the word embedding $\omega_{\text{emb}} \in \mathsf{R}^{V \times D}$, with $V$ vocabulary and $D$ dimensionality, the dropout masks types (rows) instead of words in a sequence .
- Finally, for the additional $L_2$ regularisation, we use a separate weight decay: for weights $\lambda_\omega = \frac{1 - p_{\text{drop}}}{N}$ with $p_{\text{drop}}$ dropout, and for biases ($b$): $\lambda_{\text{b}} = \frac{1}{N}$.

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Results

| Training | Model | SNLI | Breaking NLI |
|----------|-------|------|--------------|
| SNLI | $ESIM^{\dagger}$ | 87.9 | 65.6 |
| | $ESIM_{ours}$ | $86.4 \pm 0.09$ | $57.6 \pm 1.9$ |
| | $ESIM_{MC}$ | $86.5 \pm 0.13$ | $68.9 \pm 1.7$ |
| MNLI+SNLI | $ESIM^{\dagger}$ | 86.3 | 74.9 |
| | $ESIM_{ours}$ | $86.8 \pm 0.05$ | $68.8 \pm 3.5$ |
| | $ESIM_{MC}$ | $86.6 \pm 0.16$ | $75.2 \pm 1.3$ |

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

# Results SNLI



(d) Gold label entailment.

(e) Gold label neutral.

(f) Gold label contradiction.

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

# Results SNLI and Breaking



(g) SNLI



(j) Breaking

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

# Results

P: The little girl is riding in the car with her dad.
H: The small girl is riding in the car with her dad.

P: The little girl is riding in the car with her dad.
H: The little girl is riding in the car with her father.

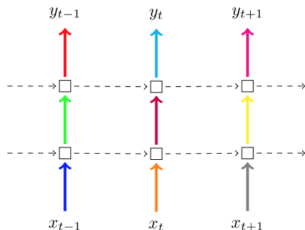P: The little girl is riding in the car with her dad.

H: The tiny girl is riding in the car with her dad.



Breaking (Category=synonyms, Gold=Entailment)

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
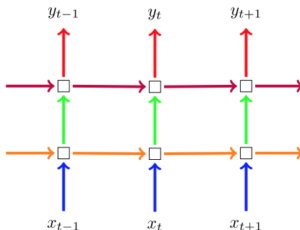Uncertainty in Natural Language Inference
References

## Homework!!

- Dropout in Recurrent Networks (Gal and Ghahramani, 2016b)
- Use the same dropout mask at each time step for both inputs, outputs, and recurrent layers
- The RNN can be framed as a probabilistic model.



(a) Naive dropout RNN          (b) Variational RNN

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Literature I

Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan
    Wierstra. Weight uncertainty in neural networks. In *Proceedings of the
    32Nd International Conference on International Conference on Machine
    Learning - Volume 37*, ICML'15, pages 1613–1622. JMLR.org, 2015.
    URL http://dl.acm.org/citation.cfm?id=3045118.3045290.

Qian Chen, Xiaodan Zhu, Zhen-Hua Ling, Si Wei, Hui Jiang, and Diana
    Inkpen. Enhanced lstm for natural language inference. In *Proceedings
    of the 55th Annual Meeting of the Association for Computational
    Linguistics (Volume 1: Long Papers)*, volume 1, pages 1657–1668,
    2017.

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Literature II

Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In Maria Florina Balcan and Kilian Q. Weinberger, editors, *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pages 1050–1059, New York, New York, USA, 20–22 Jun 2016a. PMLR. URL http://proceedings.mlr.press/v48/gal16.html.

Yarin Gal and Zoubin Ghahramani. A theoretically grounded application of dropout in recurrent neural networks. In *Advances in neural information processing systems*, pages 1019–1027, 2016b.

Introduction
Levels of Representation
RTE Methods
Current Methods
Latent Variable Models
Uncertainty in Natural Language Inference
References

## Literature III

Max Glockner, Vered Shwartz, and Yoav Goldberg. Breaking NLI systems with sentences that require simple lexical inferences. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 650–655, Melbourne, Australia, July 2018. Association for Computational Linguistics. URL https://www.aclweb.org/anthology/P18-2103.

David J. C. MacKay. A practical bayesian framework for backpropagation networks. *Neural Comput.*, 4(3):448–472, May 1992. ISSN 0899-7667. doi: 10.1162/neco.1992.4.3.448. URL http://dx.doi.org/10.1162/neco.1992.4.3.448.