# Machine Learning Final Project
*Learning Continuous Phrase Representations for Translation Modeling*

Richard Luong (2014403075)
Kerry Zhang (2014403073)

June 25, 2015

Problem introduction
000

Dataset

Method
00000000000000

Results

Summary

# TABLE OF CONTENTS

# PROBLEM INTRODUCTION

### Deep learning

- ▶ Multiple processing layers

- ▶ Neural networks
- ▶ State-of-the-art

- ▶ Statistical machine translation?

PROBLEM INTRODUCTION

**Deep learning**

- ▶ Multiple processing layers
    - ▶ Layers of abstraction

- ▶ Neural networks

- ▶ State-of-the-art

    - ▶ Visual object recognition

    - ▶ Object detection

    - ▶ Speech recognition

    - ▶ Drug discoveries, genomics

- ▶ Statistical machine translation?

## PROBLEM INTRODUCTION

**Deep learning**

- ▶ Multiple processing layers
    - ▶ Layers of abstraction
- ▶ Neural networks
- ▶ State-of-the-art
    - ▶ Visual object recognition
    - ▶ Object detection
    - ▶ Speech recognition
    - ▶ Drug discoveries, genomics
- ▶ Statistical machine translation?

## PROBLEM INTRODUCTION

### Deep learning

- ▶ Multiple processing layers
    - ▶ Layers of abstraction
- ▶ Neural networks
- ▶ State-of-the-art
    - ▶ Visual object recognition
    - ▶ Object detection
    - ▶ Speech recognition
    - ▶ Drug discoveries, genomics
- ▶ Statistical machine translation?

## PROBLEM INTRODUCTION

**Deep learning**

- ▶ Multiple processing layers
    - ▶ Layers of abstraction
- ▶ Neural networks
- ▶ State-of-the-art
    - ▶ Visual object recognition
    - ▶ Object detection
    - ▶ Speech recognition
    - ▶ Drug discoveries, genomics
- ▶ Statistical machine translation?

## Problem introduction

**Deep learning**

- ▶ Multiple processing layers
  - ▶ Layers of abstraction
- ▶ Neural networks
- ▶ State-of-the-art
  - ▶ Visual object recognition
  - ▶ Object detection
  - ▶ Speech recognition
  - ▶ Drug discoveries, genomics
- ▶ Statistical machine translation?

## PROBLEM INTRODUCTION

**Deep learning**

- ▶ Multiple processing layers
    - ▶ Layers of abstraction
- ▶ Neural networks
- ▶ State-of-the-art
    - ▶ Visual object recognition
    - ▶ Object detection
    - ▶ Speech recognition
    - ▶ Drug discoveries, genomics
- ▶ Statistical machine translation?

# PROBLEM INTRODUCTION

**Deep learning**

- ▶ Multiple processing layers
    - ▶ Layers of abstraction
- ▶ Neural networks
- ▶ State-of-the-art
    - ▶ Visual object recognition
    - ▶ Object detection
    - ▶ Speech recognition
    - ▶ Drug discoveries, genomics
- ▶ Statistical machine translation?

## PROBLEM INTRODUCTION

**Deep learning**

- ▸ Multiple processing layers
    - ▸ Layers of abstraction
- ▸ Neural networks
- ▸ State-of-the-art
    - ▸ Visual object recognition
    - ▸ Object detection
    - ▸ Speech recognition
    - ▸ Drug discoveries, genomics
- ▸ Statistical machine translation?

## Problem introduction

Learning Continuous Phrase Representation for Translation
Modeling - *Gao et. al. 2013*

*This paper tackles the sparsity problem in estimating phrase
translation probabilities by learning continuous phrase
representations, whose distributed nature enables the sharing of
related phrases in their representations.*

# PROBLEM INTRODUCTION

**Data sparsity**

▸ Phrases as translation units

    ▸ Not linguistic notion

    ▸ Sequence of words

▸ Longer phrases

    ▸ Better translations

    ▸ Less likely seen

    ▸ Data sparsity

## PROBLEM INTRODUCTION

### Data sparsity

- ▶ Phrases as translation units
    - ▶ Not linguistic notion
    - ▶ Sequence of words
- ▶ Longer phrases
    - ▶ Better translations
    - ▶ Less likely seen
    - ▶ Data sparsity

# PROBLEM INTRODUCTION

**Data sparsity**

- ▶ Phrases as translation units
    - ▶ Not linguistic notion
    - ▶ Sequence of words
- ▶ Longer phrases
    - ▶ Better translations
    - ▶ Less likely seen
    - ▶ Data sparsity

PROBLEM INTRODUCTION

**Data sparsity**

- ▶ Phrases as translation units
    - ▶ Not linguistic notion
    - ▶ Sequence of words
- ▶ Longer phrases
    - ▶ Better translations
    - ▶ Less likely seen
    - ▶ Data sparsity

# PROBLEM INTRODUCTION

**Data sparsity**

- ▸ Phrases as translation units
  - ▸ Not linguistic notion
  - ▸ Sequence of words
- ▸ Longer phrases
  - ▸ Better translations
  - ▸ Less likely seen
  - ▸ Data sparsity

## PROBLEM INTRODUCTION

**Data sparsity**

- ▶ Phrases as translation units
    - ▶ Not linguistic notion
    - ▶ Sequence of words
- ▶ Longer phrases
    - ▶ Better translations
    - ▶ Less likely seen
    - ▶ Data sparsity

# PROBLEM INTRODUCTION

**Data sparsity**

- ▶ Phrases as translation units
    - ▶ Not linguistic notion
    - ▶ Sequence of words
- ▶ Longer phrases
    - ▶ Better translations
    - ▶ Less likely seen
    - ▶ Data sparsity

DATASET

**TED talk transcripts (English and Chinese)**

## DATASET

**TED talk transcripts (English and Chinese)**

- ▶ Source sentences (Chinese)

# DATASET

**TED talk transcripts (English and Chinese)**

- ▶ Source sentences (Chinese)
- ▶ Target sentences (English)

# Dataset

**TED talk transcripts (English and Chinese)**

- ▶ Source sentences (Chinese)
- ▶ Target sentences (English)
- ▶ N-best-list (N translation hypotheses per sentence)
  - ▶ Generated by Moses: open source statistical machine translation system

# THE LOG-LINEAR MODEL FOR SMT

Standard statistical machine translation (SMT) systems are based on linear combinations of different features:

► Phrase translation table

► Reordering model

► Language model

► Many extensions:

  ► Bidirectional translation probabilities

  ► Word penalty

  ► Phrase penalty

  ► Continuous-space Phrase Translation Model (CPTM)

N-best (translation hypothesis) list generated for each source sentence

# THE LOG-LINEAR MODEL FOR SMT

Standard statistical machine translation (SMT) systems are based on linear combinations of different features:

▶ Phrase translation table

▶ Reordering model

▶ Language model

▶ Many extensions:

    ▶ Bidirectional translation probabilities

    ▶ Word penalty

    ▶ Phrase penalty

    ▶ Continuous-space Phrase Translation Model (CPTM)

N-best (translation hypothesis) list generated for each source sentence

## THE LOG-LINEAR MODEL FOR SMT

Standard statistical machine translation (SMT) systems are
based on linear combinations of different features:

- ▶ Phrase translation table

- ▶ Reordering model

- ▶ Language model

- ▶ Many extensions:

  - ▶ Bidirectional translation probabilities

  - ▶ Word penalty

  - ▶ Phrase penalty

  - ▶ Continuous-space Phrase Translation Model (CPTM)

N-best (translation hypothesis) list generated for each source
sentence

# THE LOG-LINEAR MODEL FOR SMT

Standard statistical machine translation (SMT) systems are based on linear combinations of different features:

- ▶ Phrase translation table
- ▶ Reordering model
- ▶ Language model
- ▶ Many extensions:
    - ▶ Bidirectional translation probabilities
    - ▶ Word penalty
    - ▶ Phrase penalty
    - ▶ Continuous-space Phrase Translation Model (CPTM)

N-best (translation hypothesis) list generated for each source sentence

# THE LOG-LINEAR MODEL FOR SMT

Standard statistical machine translation (SMT) systems are based on linear combinations of different features:

- ▶ Phrase translation table
- ▶ Reordering model
- ▶ Language model
- ▶ Many extensions:
  - ▶ Bidirectional translation probabilities
  - ▶ Word penalty
  - ▶ Phrase penalty
  - ▶ Continuous-space Phrase Translation Model (CPTM)

N-best (translation hypothesis) list generated for each source sentence

# THE LOG-LINEAR MODEL FOR SMT

Standard statistical machine translation (SMT) systems are based on linear combinations of different features:

► Phrase translation table

► Reordering model

► Language model

► Many extensions:
  ► Bidirectional translation probabilities
  ► Word penalty
  ► Phrase penalty
  ► Continuous-space Phrase Translation Model (CPTM)

N-best (translation hypothesis) list generated for each source sentence

# THE LOG-LINEAR MODEL FOR SMT

Standard statistical machine translation (SMT) systems are based on linear combinations of different features:

- ▶ Phrase translation table
- ▶ Reordering model
- ▶ Language model
- ▶ Many extensions:
    - ▶ Bidirectional translation probabilities
    - ▶ Word penalty
    - ▶ Phrase penalty
    - ▶ Continuous-space Phrase Translation Model (CPTM)

N-best (translation hypothesis) list generated for each source sentence

## Continuous-space Phrase Translation Model (CPTM)

- ▶ Addresses the data sparsity problem
- ▶ Score translation pair based on notion of similarity
- ▶ Dimensionality reduction
  - ▶ Raw phrase → Bag-of-words (word vector)
  - ▶ Neural network
    - ▶ Input: raw one-hot word vector of phrase
    - ▶ Share model between source and target languages
    - ▶ Reduce the size to output features
  - ▶ Output: lower dimension feature vector
- ▶ Similarity is measured by dot product of output feature vectors
- ▶

# Continuous-space Phrase Translation Model (CPTM)

- ▶ Addresses the data sparsity problem
- ▶ Score translation pair based on notion of similarity
- ▶ Dimensionality reduction
  - ▶ Raw phrase → Bag-of-words (word vector)
  - ▶ Neural network
    - ▶ Map from input layer to hidden layer
    - ▶ Share similar network across two languages
    - ▶ Learn the non-scaled features
  - ▶ Output: lower dimension feature vector
- ▶ Similarity is measured by dot product of output feature vectors
- ▶

# Continuous-space Phrase Translation Model (CPTM)

- ▶ Addresses the data sparsity problem
- ▶ Score translation pair based on notion of similarity
- ▶ Dimensionality reduction
  - ▶ Raw phrase → Bag-of-words (word vector)
  - ▶ Neural network
    - ▶ Shared vocabulary for both languages
    - ▶ Same neural network for both languages
    - ▶ Helps discovering latent features
  - ▶ Output: lower dimension feature vector
- ▶ Similarity is measured by dot product of output feature vectors
- ▶

## Continuous-space Phrase Translation Model (CPTM)

- ▶ Addresses the data sparsity problem
- ▶ Score translation pair based on notion of similarity
- ▶ Dimensionality reduction
    - ▶ Raw phrase → Bag-of-words (word vector)
    - ▶ Neural network
        - ▶ Shared vocabulary for both languages
        - ▶ Same neural network for both languages
        - ▶ Helps discovering latent features
    - ▶ Output: lower dimension feature vector
- ▶ Similarity is measured by dot product of output feature vectors
- ▶

# Continuous-space Phrase Translation Model (CPTM)

- ► Addresses the data sparsity problem
- ► Score translation pair based on notion of similarity
- ► Dimensionality reduction
  - ► Raw phrase → Bag-of-words (word vector)
  - ► Neural network
    - ► Shared vocabulary for both languages
    - ► Same neural network for both languages
    - ► Helps discovering latent features
  - ► Output: lower dimension feature vector
- ► Similarity is measured by dot product of output feature vectors
- ►

# Continuous-space Phrase Translation Model (CPTM)

- ▶ Addresses the data sparsity problem
- ▶ Score translation pair based on notion of similarity
- ▶ Dimensionality reduction
  - ▶ Raw phrase → Bag-of-words (word vector)
  - ▶ Neural network
    - ▶ Shared vocabulary for both languages
    - ▶ Same neural network for both languages
    - ▶ Helps discovering latent features
  - ▶ Output: lower dimension feature vector
- ▶ Similarity is measured by dot product of output feature vectors
- ▶

# Continuous-space Phrase Translation Model (CPTM)

- ► Addresses the data sparsity problem
- ► Score translation pair based on notion of similarity
- ► Dimensionality reduction
  - ► Raw phrase → Bag-of-words (word vector)
  - ► Neural network
    - ► Shared vocabulary for both languages
    - ► Same neural network for both languages
    - ► Helps discovering latent features
  - ► Output: lower dimension feature vector
- ► Similarity is measured by dot product of output feature vectors
- ►

# CONTINUOUS-SPACE PHRASE TRANSLATION MODEL (CPTM)

- ▶ Addresses the data sparsity problem
- ▶ Score translation pair based on notion of similarity
- ▶ Dimensionality reduction
    - ▶ Raw phrase → Bag-of-words (word vector)
    - ▶ Neural network
        - ▶ Shared vocabulary for both languages
        - ▶ Same neural network for both languages
        - ▶ Helps discovering latent features
    - ▶ Output: lower dimension feature vector
- ▶ Similarity is measured by dot product of output feature vectors

- ▶

# CONTINUOUS-SPACE PHRASE TRANSLATION MODEL (CPTM)

- ▶ Addresses the data sparsity problem
- ▶ Score translation pair based on notion of similarity
- ▶ Dimensionality reduction
    - ▶ Raw phrase → Bag-of-words (word vector)
    - ▶ Neural network
        - ▶ Shared vocabulary for both languages
        - ▶ Same neural network for both languages
        - ▶ Helps discovering latent features
    - ▶ Output: lower dimension feature vector
- ▶ Similarity is measured by dot product of output feature vectors
- ▶

# Continuous-space Phrase Translation Model (CPTM)

- ▶ Addresses the data sparsity problem
- ▶ Score translation pair based on notion of similarity
- ▶ Dimensionality reduction
    - ▶ Raw phrase → Bag-of-words (word vector)
    - ▶ Neural network
        - ▶ Shared vocabulary for both languages
        - ▶ Same neural network for both languages
        - ▶ Helps discovering latent features
    - ▶ Output: lower dimension feature vector
- ▶ Similarity is measured by dot product of output feature vectors
- ▶

# Continuous-space Phrase Translation Model (CPTM)

- ▶ Addresses the data sparsity problem
- ▶ Score translation pair based on notion of similarity
- ▶ Dimensionality reduction
    - ▶ Raw phrase → Bag-of-words (word vector)
    - ▶ Neural network
        - ▶ Shared vocabulary for both languages
        - ▶ Same neural network for both languages
        - ▶ Helps discovering latent features
    - ▶ Output: lower dimension feature vector
- ▶ Similarity is measured by dot product of output feature vectors
- ▶

Problem introduction
000

Dataset

Method
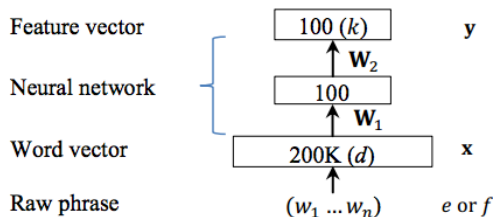00●000000000000

Results

Summary

## OVERVIEW

Learning Continuous Phrase Representation for Translation
Modeling - *Gao et. al. 2013*



Figure 2. A neural network model for phrases
giving rise to their continuous representations.
The model with the same form is used for both
source and target languages.

OBJECTIVE FUNCTION / LOSS FUNCTION

$$L(\theta) = -xBLEU(\theta)$$

*xBLEU* - N-best list expected BLEU score

# BLEU score

**BLEU** - Bilingual Evaluation Understudy

*It is the most widely used automated method of determining the quality of machine translation. The BLEU metric scores a translation on a scale of 0 to 1, but is frequently displayed as a percentage value. The closer to 1, the more the translation correlates to a human translation.*

## Objective function / Loss function

$$L(\theta) = -xBLEU(\theta)$$

*xBLEU* - N-best list expected BLEU score

## Objective function

$$\text{xBleu}(\boldsymbol{\theta}) = \sum_{E \in \text{GEN}(F_i)} P(E|F_i)\text{sBleu}(E_i, E)$$

$sBLEU$ - sentence level BLEU

$P(E|F_i)$ - translation probabilty using $softmax$

$$P(E|F_i) = \frac{\exp(\gamma \boldsymbol{\lambda}^T \mathbf{h}(F_i, E, A))}{\sum_{E' \in \text{GEN}(F_i)} \exp\left(\gamma \boldsymbol{\lambda}^T \mathbf{h}(F_i, E', A)\right)}$$

# SIMILARITY

$$\mathbf{y} \equiv \phi(\mathbf{x}) = \tanh\left(\mathbf{W_2}^\mathrm{T}(\tanh(\mathbf{W_1}^\mathrm{T}\mathbf{x}))\right) \quad (2)$$

$$\mathrm{score}(f, e) \equiv \mathrm{sim}_\boldsymbol{\theta}(\mathbf{x}_f, \mathbf{x}_e) = \mathbf{y}_f^\mathrm{T}\mathbf{y}_e$$

## PARAMETERS

Parameters: $\theta = \{W_1, W_2\}$

# PARAMETERS

Parameters: $\theta = \{W_1, W_2\}$

1. We use a baseline phrase-based SMT system to generate for each source sentence in training data an N-best list of translation hypotheses[4].

2. We set $\lambda$ to that of the baseline system and let $\lambda_{M+1} = 1$, and optimize $\theta$ w.r.t. a loss function on training data[5].

3. We fix $\theta$, and optimize $\lambda$ using MERT (Och 2003) to maximize BLEU on dev data.

# LEARNING $\theta$

▶ We learn $\theta$ by using gradient based numerical optimization
  algorithm (L-BFGS)

▶ We compute the gradient of the loss function

$$\frac{\partial \mathcal{L}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \sum_{(f,e)} \frac{\partial \mathcal{L}(\boldsymbol{\theta})}{\partial \mathrm{sim}_{\boldsymbol{\theta}}(\mathbf{x}_f, \mathbf{x}_e)} \frac{\partial \mathrm{sim}_{\boldsymbol{\theta}}(\mathbf{x}_f, \mathbf{x}_e)}{\partial \boldsymbol{\theta}}$$

$$= \sum_{(f,e)} -\delta_{(f,e)} \frac{\partial \mathrm{sim}_{\boldsymbol{\theta}}(\mathbf{x}_f, \mathbf{x}_e)}{\partial \boldsymbol{\theta}} \quad (7)$$

Computing $\delta sim_\theta(x_f, x_e)/\delta\theta$

$$\frac{\partial \text{sim}_{\boldsymbol{\theta}}(\mathbf{x}_f, \mathbf{x}_e)}{\partial \mathbf{W}_1}$$

$$= \mathbf{x}_f \left( \mathbf{W}_2 \left( \mathbf{y}_e^2 \circ \sigma'(\mathbf{z}_f^2) \right) \circ \sigma'(\mathbf{z}_f^1) \right)^{\mathbf{T}}$$

$$+ \mathbf{x}_e \left( \mathbf{W}_2 \left( \mathbf{y}_f^2 \circ \sigma'(\mathbf{z}_e^2) \right) \circ \sigma'(\mathbf{z}_e^1) \right)^{\mathbf{T}}$$

# COMPUTING $\delta sim_\theta(x_f, x_e)/\delta\theta$

$$\frac{\partial\text{sim}_\theta(\mathbf{x}_f, \mathbf{x}_e)}{\partial\mathbf{W}_1}$$

$$= \mathbf{x}_f \left(\mathbf{W}_2 \left(\mathbf{y}_e^2 \circ \sigma'(\mathbf{z}_f^2)\right) \circ \sigma'(\mathbf{z}_f^1)\right)^\mathrm{T}$$

$$+\mathbf{x}_e \left(\mathbf{W}_2 \left(\mathbf{y}_f^2 \circ \sigma'(\mathbf{z}_e^2)\right) \circ \sigma'(\mathbf{z}_e^1)\right)^\mathrm{T}$$

$$\frac{\partial\text{sim}_\theta(\mathbf{x}_f, \mathbf{x}_e)}{\partial\mathbf{W}_2} = \frac{\partial(\mathbf{y}_f^2)^\mathrm{T}}{\partial\mathbf{W}_2}\mathbf{y}_e^2 + (\mathbf{y}_f^2)^\mathrm{T}\frac{\partial\mathbf{y}_e^2}{\partial\mathbf{W}_2}$$

$$= \mathbf{y}_f^1 \left(\mathbf{y}_e^2 \circ \sigma'(\mathbf{z}_f^2)\right)^\mathrm{T} + \mathbf{y}_e^1 \left(\mathbf{y}_f^2 \circ \sigma'(\mathbf{z}_e^2)\right)^\mathrm{T} \quad (9)$$

## COMPUTING THE ERROR TERM $\delta_{(f,e)}$

$$\delta_{(f,e)}$$
$$= \sum_{(E,A)\in GEN(F_i)} U(\boldsymbol{\theta}, E)P(E|F_i)\lambda_{M+1}N(f, e; A)$$

where $\qquad\qquad\qquad\qquad\qquad$ (14)

$$U(\boldsymbol{\theta}, E) = \text{sBleu}(E_i, E) - \text{xBleu}(\boldsymbol{\theta}).$$

# COMPUTING THE ERROR TERM $\delta_{(f,e)}$

$$\delta_{(f,e)} = \sum_{(E,A) \in GEN(F_i)} U(\boldsymbol{\theta}, E) P(E|F_i) \lambda_{M+1} N(f, e; A)$$

where $\qquad\qquad\qquad\qquad$ (14)

$$U(\boldsymbol{\theta}, E) = \text{sBleu}(E_i, E) - \text{xBleu}(\boldsymbol{\theta}).$$

$N(f, e)$ - number of times $(f, e)$ occurs in $A$

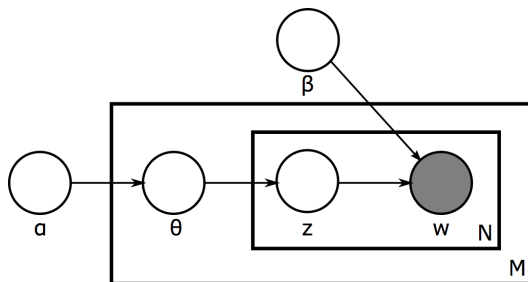$$\text{xBleu}(\boldsymbol{\theta}) = \sum_{E \in \text{GEN}(F_i)} P(E|F_i) \text{sBleu}(E_i, E)$$

COMPUTING THE ERROR TERM $\delta_{(f,e)}$

$$\delta_{(f,e)} = \sum_{(E,A) \in GEN(F_i)} U(\boldsymbol{\theta}, E) P(E|F_i) \lambda_{M+1} N(f, e; A)$$

where $\qquad\qquad\qquad\qquad$ (14)

$$U(\boldsymbol{\theta}, E) = sBleu(E_i, E) - xBleu(\boldsymbol{\theta}).$$

$N(f, e)$ - number of times $(f, e)$ occurs in $A$

$$xBleu(\boldsymbol{\theta}) = \sum_{E \in GEN(F_i)} P(E|F_i) sBleu(E_i, E)$$

$$P(E|F_i) = \frac{\exp(\gamma \boldsymbol{\lambda}^{\mathbf{T}} \mathbf{h}(F_i, E, A))}{\sum_{E' \in GEN(F_i)} \exp\left(\gamma \boldsymbol{\lambda}^{\mathbf{T}} \mathbf{h}(F_i, E', A)\right)}$$

# LEARNING $\theta$

- The parameters are trained using **Stochastic gradient descent**
- Initialize
  - $W_1$ bilingual topic model trained on parallell data
  - $W_2$ identity matrix

# LATENT DIRICHLET ALLOCATION

LATENT DIRICHLET ALLOCATION

Preprocessing:

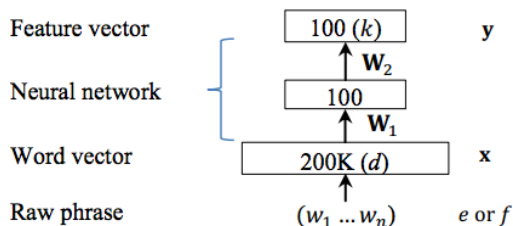- Stop words

## LATENT DIRICHLET ALLOCATION

Preprocessing:

► Stop words

**Output:** $W_1$ - 100 x 50,000

Results

After multiple overnight test runs, our implementation of the CPTM algorithm conducted into satisfiable results.

## SUMMARY



Feature vector — $100\ (k)$ — $\mathbf{y}$

$\mathbf{W}_2$

Neural network — $100$

$\mathbf{W}_1$

Word vector — $200\text{K}\ (d)$ — $\mathbf{x}$

Raw phrase — $(w_1 \dots w_n)$ — $e$ or $f$

$$\theta = \{W_1, W_2\}$$

Gradient descent:

$$\theta_{n+1} = \theta_n - \eta \frac{\partial L(\theta)}{\partial \theta} = \theta_n - \eta \left[ \sum_{(f,e)} -\delta_{(f,e)} \right.$$

# Results

In order to test the theoretical validity of our system, a plot over the loss function was made.
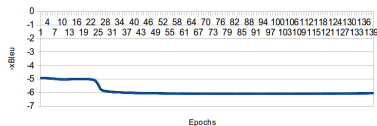


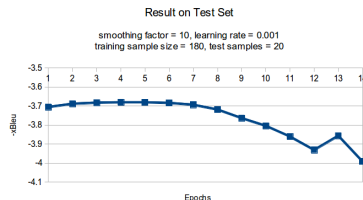Figure: Theoretical validdity of our neural network

# Results



Figure: Result on test set

The results seem to suggest that the system is indeed performing a correct gradient descent.

SUMMARY

► We have been learning continuous phrase representations
   for a SMT system.

# Summary

- ▶ We have been learning continuous phrase representations for a SMT system.
- ▶ Using a neural network, we have projected both source and target sentences into a real-number vector space.

## Summary

- We have been learning continuous phrase representations for a SMT system.
- Using a neural network, we have projected both source and target sentences into a real-number vector space.
- Loss function of the NN has been the negative expected BLEU score for a source sentence and it's possible translation.

## Summary

- We have been learning continuous phrase representations for a SMT system.
- Using a neural network, we have projected both source and target sentences into a real-number vector space.
- Loss function of the NN has been the negative expected BLEU score for a source sentence and it's possible translation.
- We also tested the validity of our system, showing that the loss function decreases.

FINALLY

Thank you for listening! It has been an adventurous journey.