# Introduction to Statistical Machine Translation

David Kauchak

CS457 – Fall 2011

Some slides adapted from

Philipp Koehn          Kevin Knight

CSAIL
Massachusetts Institute of Technology

USC/Information Sciences Institute
USC/Computer Science Department

---

# Admin

- How are the projects going?
- Remaining classes applications
  - 3 MT
    - General overview today
    - Dive into one specific implementation next time
    - MT lab
  - Other applications
    - Information extraction
    - Information retrieval
    - Question answering/summarization

---

# Language translation



Yo quiero Taco Bell

---

# MT Systems

Where have you seen machine translation systems?

AGUA PARA RIEGO DE PLANTAS
AGUA NO POTABLE, PROHIBIDO MOJARSE
NO DRINKING WATER,
IT IS PROHIBITED TO GET WET HERE!

---

## Machine Translation

美国关岛国际机场及其办公室均接获一名自称沙地阿拉伯富商拉登等发出的电子邮件，威胁将会向机场等公众地方发动生化袭击後，关岛经保持高度戒备。

→ The U.S. island of Guam is maintaining a high state of alert after the Guam airport and its offices both received an e-mail from someone calling himself the Saudi Arabian Osama bin Laden and threatening a biological/chemical attack against public places such as the airport .

The classic acid test for natural language processing.

Requires capabilities in both interpretation and generation.

People around the world stubbornly refuse to write everything in English.

---

## Machine Translation

美国关岛国际机场及其办公室均接获一名自称沙地阿拉伯富商拉登等发出的电子邮件，威胁将会向机场等公众地方发动生化袭击後，关岛经保持高度戒备。

Machine translation is becoming very prevalent

Even PowerPoint has translation built into it!

The U.S. island of Guam is maintaining a high state of alert after the Guam airport and its offices both received an e-mail from someone calling himself the Saudi Arabian Osama bin Laden and threatening a biological/chemical attack against public places such as the airport .

United States Guam International Airport and the Office received one claiming to be a wealthy Saudi Arabia an email such as Osama bin Laden, threats to the airport after biochemical attacks in public places such as Guam remain on high alert.

---

## Which is the human?

Beijing Youth Daily said that under the Ministry of Agriculture, the beef will be destroyed after tests.

The Beijing Youth Daily pointed out that the seized beef would be disposed of after being examined according to advice from the Ministry of Agriculture.

**?**

## Which is the human?

Pakistan President Pervez Musharraf Wins Senate Confidence Vote

Pakistani President Musharraf Won the Trust Vote in Senate and Lower House

**?**

## Which is the human?

There was not a single vote against him."

No members vote against him. "

**?**

## Warren Weaver (1947)

**ingcmpnqsnwf cv fpn owoktvcv**

**hu ihgzsnwfv rqcffnw cw owgcnwf**

**kowazoanv ...**

## Warren Weaver (1947)

```
    e      e  e          e
ingcmpnqsnwf cv fpn owoktvcv
         e        e           e
hu ihgzsnwfv rqcffnw cw owgcnwf
      e
kowazoanv ...
```

## Warren Weaver (1947)



```
      e    e   e        the
ingcmpnqsnwf cv fpn owoktvcv
          e       e         e
hu ihgzsnwfv rqcffnw cw owgcnwf
         e
kowazoanv ...
```

## Warren Weaver (1947)



```
      e   he  e        the
ingcmpnqsnwf cv fpn owoktvcv
          e       e       e t
hu ihgzsnwfv rqcffnw cw owgcnwf
         e
kowazoanv ...
```

## Warren Weaver (1947)



```
     e   he  e     of the
ingcmpnqsnwf cv fpn owoktvcv
           e       e        e t
hu ihgzsnwfv rqcffnw cw owgcnwf
          e
kowazoanv ...
```

## Warren Weaver (1947)



```
     e   he  e     of the       fof
ingcmpnqsnwf cv fpn owoktvcv
         e  f   o  e o      oe t
hu ihgzsnwfv rqcffnw cw owgcnwf
          ef
kowazoanv ...
```

## Warren Weaver (1947)

```
 e   he e      the
ingcmpnqsnwf cv fpn owoktvcv
        e     e      e t
hu ihgzsnwfv rqcffnw cw owgcnwf
        e
kowazoanv ...
```

## Warren Weaver (1947)

```
 e   he e    is the     sis
ingcmpnqsnwf cv fpn owoktvcv
        e   s  i e i   ie t
hu ihgzsnwfv rqcffnw cw owgcnwf
        es
kowazoanv ...
```

## Warren Weaver (1947)

```
decipherment is the analysis
ingcmpnqsnwf cv fpn owoktvcv
of documents written in ancient
hu ihgzsnwfv rqcffnw cw owgcnwf
languages ...
kowazoanv ...
```

## Warren Weaver (1947)

Can this be computerized?

The non-Turkish guy next to me is even deciphering Turkish! All he needs is a statistical table of letter-pair frequencies in Turkish …

Collected mechanically from a Turkish body of text, or *corpus*

**Slide 1:**

"When I look at an article in Russian, I say: this is really written in English, but it has been coded in some strange symbols. I will now proceed to decode."
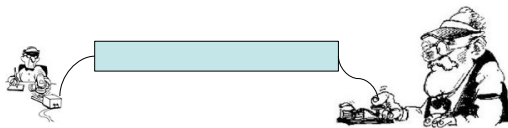- Warren Weaver, March 1947

**Slide 2:**

"When I look at an article in Russian, I say: this is really written in English, but it has been coded in some strange symbols. I will now proceed to decode."
- Warren Weaver, March 1947

"... as to the problem of mechanical translation, I frankly am afraid that the [semantic] boundaries of words in different languages are too vague ... to make any quasi-mechanical translation scheme very hopeful."
- Norbert Wiener, April 1947

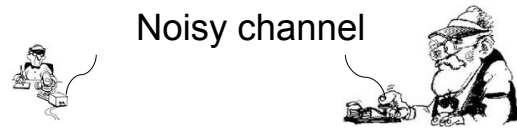**Slide 3:**

## Noisy channel

Some message is sent

along the way things get messed up

What was originally sent?

We have the mutated message, but would like to recover the original

**Slide 4:**

## Noisy channel

known system

If we know something about what goes on inside here, we might be to decode/recover the message.

# Noisy channel

"Hi bob"                    I baab

---

# Noisy channel

"Hi bob"    -'H's often get dropped
            - 'o's go to 'aa' sometimes
            - …                          I baab

---

# Noisy channel

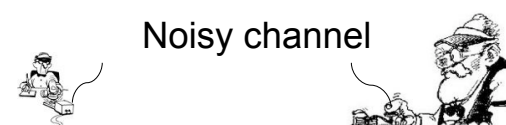"Hi bob"    -'H's often get dropped
            - 'o's go to a sometimes
            - …                          I baab

            Hi bob

---

# Noisy channel

?

Sometime, we don't know what goes on or we
only have a rough guess…

What then?

## Noisy channel

"Hi bob"
"hello"
"banana"
…

?

"I baab"
"ello"
"banana"
…

A data driven approach

Send a bunch of data through where we know what is being sent

## Noisy channel

"Hi bob"
"hello"
"banana"
…

"I baab"
"ello"
"banana"
…

Use this data to train a model

A data driven approach: learn a model of the types of transformations that occur

## Noisy channel

- source (s): what was originally sent
- target (t): what did we get through the noisy channel

We want a probabilistic model of the process:

$$p(s \mid t)$$

What is the probability of a given source sentence, given that we've seen target

## How does a model p(s | t) help us?

"…." (s)

"…." (t)

Decode: $\arg_s \max p(s \mid t)$

## Noisy channel model

$$p(s \mid t) = \frac{p(t \mid s)\,p(s)}{p(t)}$$   Bayes' rule

$p(t)$      how likely is it to receive the target message

$p(s)$      What types of messages are likely to be sent?
What do sent messages look like?

$p(t \mid s)$      What types of transformation happen?
How likely are we to generate t from s?

---

## Noisy channel model

model     $$p(s \mid t) \propto p(t \mid s)\,p(s)$$

**channel model
translation model**
how does the target
get "messed up"
from the source?

**source model
language model**
what types of
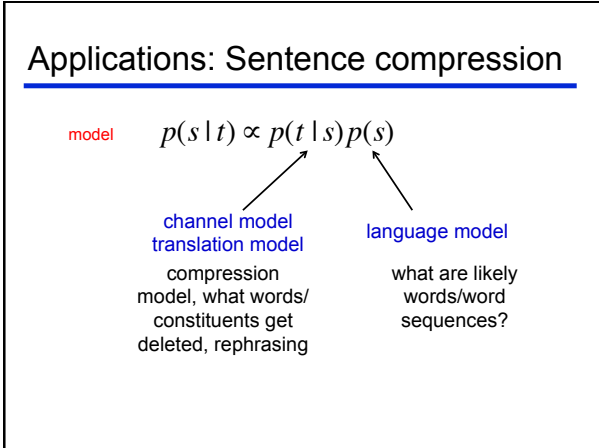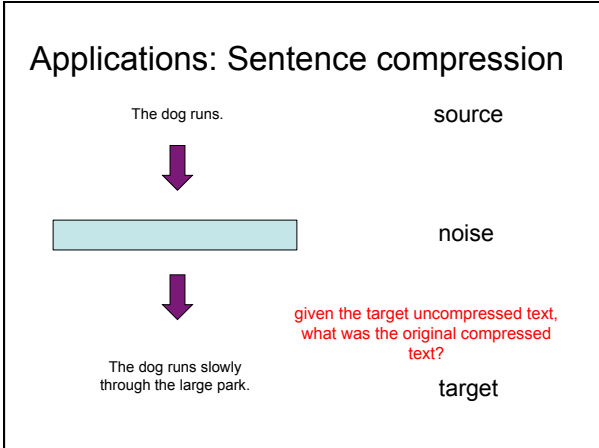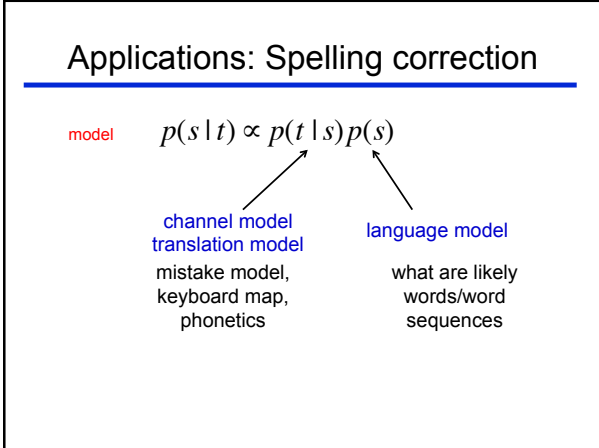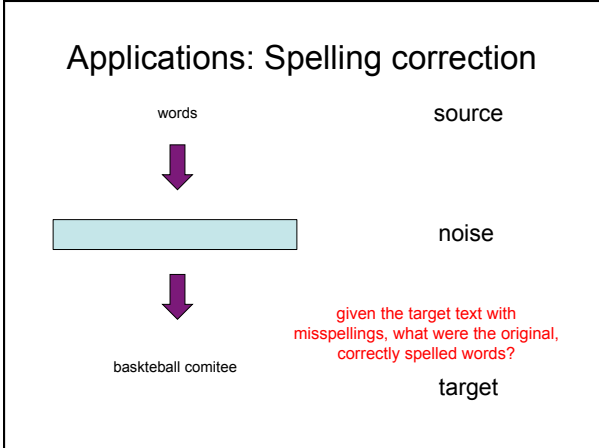things should I
expect?

---

## Applications: Speech recognition

words          source



noise

given the target audio, what
were the original words?

target

---

## Applications: Speech recognition

model     $$p(s \mid t) \propto p(t \mid s)\,p(s)$$

**channel model
translation model**
acoustic model,
how words turn in
to sounds

**language model**
what are likely
words/word
sequences

## Applications: Spelling correction

words                                  source

↓

[ ]                                    noise

↓

given the target text with
misspellings, what were the original,
correctly spelled words?

baskteball comitee

target

## Applications: Spelling correction

model    $p(s \mid t) \propto p(t \mid s) p(s)$

channel model          language model
translation model

mistake model,          what are likely
keyboard map,            words/word
phonetics                sequences

## Applications: Sentence compression

The dog runs.                          source

↓

[ ]                                    noise

↓

given the target uncompressed text,
what was the original compressed
text?

The dog runs slowly
through the large park.

target

## Applications: Sentence compression

model    $p(s \mid t) \propto p(t \mid s) p(s)$

channel model          language model
translation model

compression             what are likely
model, what words/      words/word
constituents get        sequences?
deleted, rephrasing

## Applications: Machine translation

English text      source

noise

given the target Chinese text, what was the original English text?

美国关岛国际机场及其办公室均接获一
名自称沙地阿拉伯富商拉登等发出的电
子邮件，威胁将会向机场等公众地方发
动生化袭击後，关岛经保持高度戒备。
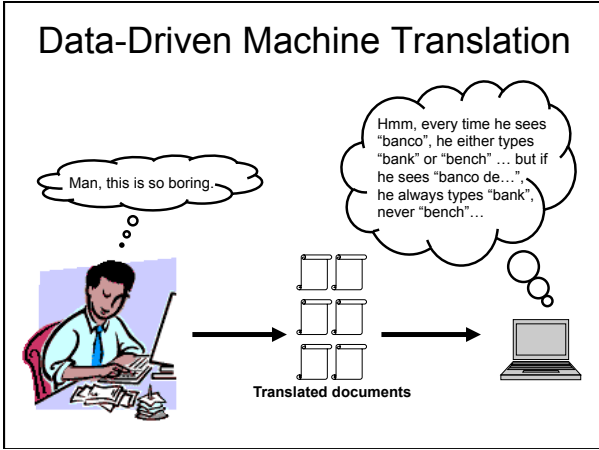
target

---

## Applications: Machine translation

model     $p(s \mid t) \propto p(t \mid s) p(s)$

channel model
translation model
how do English
words/phrases
translate to
Chinese?

language model
what are likely
English words/
word sequences?

---

## Data-Driven Machine Translation

Man, this is so boring.

Hmm, every time he sees "banco", he either types "bank" or "bench" … but if he sees "banco de…", he always types "bank", never "bench"…

**Translated documents**

---

## Welcome to the Chinese Room

Chinese texts with English translations

New Chinese Document

English Translation

You can teach yourself to translate Chinese using *only* bilingual data (without grammar books, dictionaries, any people to answer your questions…)

## Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan: farok crrrok hihok yorok clok kantok ok-yurp

| | |
|---|---|
| 1a. ok-voon ororok sprok . | 7a. lalok farok ororok lalok sprok izok enemok . |
| 1b. at-voon bichat dat . | 7b. wat jjat bichat wat dat vat eneat . |
| 2a. ok-drubel ok-voon anok plok sprok . | 8a. lalok brok anok plok nok . |
| 2b. at-drubel at-voon pippat rrat dat . | 8b. iat lat pippat rrat nnat . |
| 3a. erok sprok izok hihok ghirok . | 9a. wiwok nok izok kantok ok-yurp . |
| 3b. totat dat arrat vat hilat . | 9b. totat nnat quat oloat at-yurp . |
| 4a. ok-voon anok drok brok jok . | 10a. lalok mok nok yorok ghirok clok . |
| 4b. at-voon krat pippat sat lat . | 10b. wat nnat gat mat bat hilat . |
| 5a. wiwok farok izok stok . | 11a. lalok nok crrrok hihok yorok zanzanok . |
| 5b. totat jjat quat cat . | 11b. wat nnat arrat mat zanzanat . |
| 6a. lalok sprok izok jok stok . | 12a. lalok rarok nok izok hihok mok . |
| 6b. wat dat krat quat cat . | 12b. wat nnat forat arrat vat gat . |

## Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan: farok crrrok hihok yorok clok kantok ok-yurp

| | |
|---|---|
| 1a. ok-voon ororok sprok . | 7a. lalok farok ororok lalok sprok izok enemok . |
| 1b. at-voon bichat dat . | 7b. wat jjat bichat wat dat vat eneat . |
| 2a. ok-drubel ok-voon anok plok sprok . | 8a. lalok brok anok plok nok . |
| 2b. at-drubel at-voon pippat rrat dat . | 8b. iat lat pippat rrat nnat . |
| 3a. erok sprok izok hihok ghirok . | 9a. wiwok nok izok kantok ok-yurp . |
| 3b. totat dat arrat vat hilat . | 9b. totat nnat quat oloat at-yurp . |
| 4a. ok-voon anok drok brok jok . | 10a. lalok mok nok yorok ghirok clok . |
| 4b. at-voon krat pippat sat lat . | 10b. wat nnat gat mat bat hilat . |
| 5a. wiwok farok izok stok . | 11a. lalok nok crrrok hihok yorok zanzanok . |
| 5b. totat jjat quat cat . | 11b. wat nnat arrat mat zanzanat . |
| 6a. lalok sprok izok jok stok . | 12a. lalok rarok nok izok hihok mok . |
| 6b. wat dat krat quat cat . | 12b. wat nnat forat arrat vat gat . |

## Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan: farok crrrok hihok yorok clok kantok ok-yurp

| | |
|---|---|
| 1a. ok-voon ororok sprok . | 7a. lalok farok ororok lalok sprok izok enemok . |
| 1b. at-voon bichat dat . | 7b. wat jjat bichat wat dat vat eneat . |
| 2a. ok-drubel ok-voon anok plok sprok . | 8a. lalok brok anok plok nok . |
| 2b. at-drubel at-voon pippat rrat dat . | 8b. iat lat pippat rrat nnat . |
| 3a. erok sprok izok hihok ghirok . | 9a. wiwok nok izok kantok ok-yurp . |
| 3b. totat dat arrat vat hilat . | 9b. totat nnat quat oloat at-yurp . |
| 4a. ok-voon anok drok brok jok . | 10a. lalok mok nok yorok ghirok clok . |
| 4b. at-voon krat pippat sat lat . | 10b. wat nnat gat mat bat hilat . |
| 5a. wiwok farok izok stok . | 11a. lalok nok crrrok hihok yorok zanzanok . |
| 5b. totat jjat quat cat . | 11b. wat nnat arrat mat zanzanat . |
| 6a. lalok sprok izok jok stok . | 12a. lalok rarok nok izok hihok mok . |
| 6b. wat dat krat quat cat . | 12b. wat nnat forat arrat vat gat . |

## Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan: farok crrrok hihok yorok clok kantok ok-yurp

| | |
|---|---|
| 1a. ok-voon ororok sprok . | 7a. lalok farok ororok lalok sprok izok enemok . |
| 1b. at-voon bichat dat . | 7b. wat jjat bichat wat dat vat eneat . |
| 2a. ok-drubel ok-voon anok plok sprok . | 8a. lalok brok anok plok nok . |
| 2b. at-drubel at-voon pippat rrat dat . | 8b. iat lat pippat rrat nnat . |
| 3a. erok sprok izok hihok ghirok . | 9a. wiwok nok izok kantok ok-yurp . |
| 3b. totat dat arrat vat hilat . | 9b. totat nnat quat oloat at-yurp . |
| 4a. ok-voon anok drok brok jok . | 10a. lalok mok nok yorok ghirok clok . |
| 4b. at-voon krat pippat sat lat . | 10b. wat nnat gat mat bat hilat . |
| 5a. wiwok farok izok stok . | 11a. lalok nok crrrok hihok yorok zanzanok . ??? |
| 5b. totat jjat quat cat . | 11b. wat nnat arrat mat zanzanat . |
| 6a. lalok sprok izok jok stok . | 12a. lalok rarok nok izok hihok mok . |
| 6b. wat dat krat quat cat . | 12b. wat nnat forat arrat vat gat . |

# Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan:   farok crrrok hihok yorok clok kantok ok-yurp

| | |
|---|---|
| 1a. ok-voon ororok sprok . | 7a. lalok farok ororok lalok sprok izok enemok . |
| 1b. at-voon bichat dat . | 7b. wat jjat bichat wat dat vat eneat . |
| 2a. ok-drubel ok-voon anok plok sprok . | 8a. lalok brok anok plok nok . |
| 2b. at-drubel at-voon pippat rrat dat . | 8b. iat lat pippat rrat nnat . |
| 3a. erok sprok izok hihok ghirok . | 9a. wiwok nok izok kantok ok-yurp . |
| 3b. totat dat arrat vat hilat . | 9b. totat nnat quat oloat at-yurp . |
| 4a. ok-voon anok drok brok jok . | 10a. lalok mok nok yorok ghirok clok . |
| 4b. at-voon krat pippat sat lat . | 10b. wat nnat gat mat bat hilat . |
| 5a. wiwok farok izok stok . | 11a. lalok nok crrrok hihok yorok zanzanok . |
| 5b. totat jjat quat cat . | 11b. wat nnat arrat mat zanzanat . |
| 6a. lalok sprok izok jok stok . | 12a. lalok rarok nok izok hihok mok . |
| 6b. wat dat krat quat cat . | 12b. wat nnat forat arrat vat gat . |

---

# Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan:   farok crrrok hihok yorok clok kantok ok-yurp

| | |
|---|---|
| 1a. ok-voon ororok sprok . | 7a. lalok farok ororok lalok sprok izok enemok . |
| 1b. at-voon bichat dat . | 7b. wat jjat bichat wat dat vat eneat . |
| 2a. ok-drubel ok-voon anok plok sprok . | 8a. lalok brok anok plok nok . |
| 2b. at-drubel at-voon pippat rrat dat . | 8b. iat lat pippat rrat nnat . |
| 3a. erok sprok izok hihok ghirok . | 9a. wiwok nok izok kantok ok-yurp . |
| 3b. totat dat arrat vat hilat . | 9b. totat nnat quat oloat at-yurp . |
| 4a. ok-voon anok drok brok jok . | 10a. lalok mok nok yorok ghirok clok . |
| 4b. at-voon krat pippat sat lat . | 10b. wat nnat gat mat bat hilat . |
| 5a. wiwok farok izok stok . | 11a. lalok nok crrrok hihok yorok zanzanok . |
| 5b. totat jjat quat cat . | 11b. wat nnat arrat mat zanzanat . |
| 6a. lalok sprok izok jok stok . | 12a. lalok rarok nok izok hihok mok . |
| 6b. wat dat krat quat cat . | 12b. wat nnat forat arrat vat gat . |

---

# Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan:   farok crrrok hihok yorok clok kantok ok-yurp

| | |
|---|---|
| 1a. ok-voon ororok sprok . | 7a. lalok farok ororok lalok sprok izok enemok . |
| 1b. at-voon bichat dat . | 7b. wat jjat bichat wat dat vat eneat . |
| 2a. ok-drubel ok-voon anok plok sprok . | 8a. lalok brok anok plok nok . |
| 2b. at-drubel at-voon pippat rrat dat . | 8b. iat lat pippat rrat nnat . |
| 3a. erok sprok izok hihok ghirok . | 9a. wiwok nok izok kantok ok-yurp . |
| 3b. totat dat arrat vat hilat . | 9b. totat nnat quat oloat at-yurp . |
| 4a. ok-voon anok drok brok jok . | 10a. lalok mok nok yorok ghirok clok . |
| 4b. at-voon krat pippat sat lat . | 10b. wat nnat gat mat bat hilat . |
| 5a. wiwok farok izok stok . | 11a. lalok nok crrrok hihok yorok zanzanok . |
| 5b. totat jjat quat cat . | 11b. wat nnat arrat mat zanzanat . |
| 6a. lalok sprok izok jok stok . | 12a. lalok rarok nok izok hihok mok . |
| 6b. wat dat krat quat cat . | 12b. wat nnat forat arrat vat gat . |

---

# Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan:   farok crrrok hihok yorok clok kantok ok-yurp

| | |
|---|---|
| 1a. ok-voon ororok sprok . | 7a. lalok farok ororok lalok sprok izok enemok . |
| 1b. at-voon bichat dat . | 7b. wat jjat bichat wat dat vat eneat . |
| 2a. ok-drubel ok-voon anok plok sprok . | 8a. lalok brok anok plok nok . |
| 2b. at-drubel at-voon pippat rrat dat . | 8b. iat lat pippat rrat nnat . |
| 3a. erok sprok izok hihok ghirok . | 9a. wiwok nok izok kantok ok-yurp . |
| 3b. totat dat arrat vat hilat . | 9b. totat nnat quat oloat at-yurp . |
| 4a. ok-voon anok drok brok jok . | 10a. lalok mok nok yorok ghirok clok . ??? |
| 4b. at-voon krat pippat sat lat . | 10b. wat nnat gat mat bat hilat . |
| 5a. wiwok farok izok stok . | 11a. lalok nok crrrok hihok yorok zanzanok . |
| 5b. totat jjat quat cat . | 11b. wat nnat arrat mat zanzanat . |
| 6a. lalok sprok izok jok stok . | 12a. lalok rarok nok izok hihok mok . |
| 6b. wat dat krat quat cat . | 12b. wat nnat forat arrat vat gat . |

## Slide 1 (top-left)

# Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan:  farok crrrok hihok yorok clok kantok ok-yurp

| | |
|---|---|
| 1a. ok-voon ororok sprok . | 7a. lalok farok ororok lalok sprok izok enemok . |
| 1b. at-voon bichat dat . | 7b. wat jjat bichat wat dat vat eneat . |
| 2a. ok-drubel ok-voon anok plok sprok . | 8a. lalok brok anok plok nok . |
| 2b. at-drubel at-voon pippat rrat dat . | 8b. iat lat pippat rrat nnat . |
| 3a. erok sprok izok hihok ghirok . | 9a. wiwok nok izok kantok ok-yurp . |
| 3b. totat dat arrat vat hilat . | 9b. totat nnat quat oloat at-yurp . |
| 4a. ok-voon anok drok brok jok . | 10a. lalok mok nok yorok ghirok clok . |
| 4b. at-voon krat pippat sat lat . | 10b. wat nnat gat mat bat hilat . |
| 5a. wiwok farok izok stok . | 11a. lalok nok crrrok hihok yorok zanzanok . |
| 5b. totat jjat quat cat . | 11b. wat nnat arrat mat zanzanat . |
| 6a. lalok sprok izok jok stok . | 12a. lalok rarok nok izok hihok mok . |
| 6b. wat dat krat quat cat . | 12b. wat nnat forat arrat vat gat . |

## Slide 2 (top-right)

# Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan:  farok crrrok hihok yorok clok kantok ok-yurp

| | |
|---|---|
| 1a. ok-voon ororok sprok . | 7a. lalok farok ororok lalok sprok izok enemok . |
| 1b. at-voon bichat dat . | 7b. wat jjat bichat wat dat vat eneat . |
| 2a. ok-drubel ok-voon anok plok sprok . | 8a. lalok brok anok plok nok . |
| 2b. at-drubel at-voon pippat rrat dat . | 8b. iat lat pippat rrat nnat . |
| 3a. erok sprok izok hihok ghirok . | 9a. wiwok nok izok kantok ok-yurp . |
| 3b. totat dat arrat vat hilat . | 9b. totat nnat quat oloat at-yurp . |
| 4a. ok-voon anok drok brok jok . | 10a. lalok mok nok yorok ghirok clok . process of elimination |
| 4b. at-voon krat pippat sat lat . | 10b. wat nnat gat mat bat hilat . |
| 5a. wiwok farok izok stok . | 11a. lalok nok crrrok hihok yorok zanzanok . |
| 5b. totat jjat quat cat . | 11b. wat nnat arrat mat zanzanat . |
| 6a. lalok sprok izok jok stok . | 12a. lalok rarok nok izok hihok mok . |
| 6b. wat dat krat quat cat . | 12b. wat nnat forat arrat vat gat . |

## Slide 3 (bottom-left)

# Centauri/Arcturan [Knight, 1997]

Your assignment, translate this to Arcturan:  farok crrrok hihok yorok clok kantok ok-yurp

| | |
|---|---|
| 1a. ok-voon ororok sprok . | 7a. lalok farok ororok lalok sprok izok enemok . |
| 1b. at-voon bichat dat . | 7b. wat jjat bichat wat dat vat eneat . |
| 2a. ok-drubel ok-voon anok plok sprok . | 8a. lalok brok anok plok nok . |
| 2b. at-drubel at-voon pippat rrat dat . | 8b. iat lat pippat rrat nnat . |
| 3a. erok sprok izok hihok ghirok . | 9a. wiwok nok izok kantok ok-yurp . |
| 3b. totat dat arrat vat hilat . | 9b. totat nnat quat oloat at-yurp . |
| 4a. ok-voon anok drok brok jok . | 10a. lalok mok nok yorok ghirok clok . |
| 4b. at-voon krat pippat sat lat . | 10b. wat nnat gat mat bat hilat . |
| 5a. wiwok farok izok stok . | 11a. lalok nok crrrok hihok yorok zanzanok . cognate? |
| 5b. totat jjat quat cat . | 11b. wat nnat arrat mat zanzanat . |
| 6a. lalok sprok izok jok stok . | 12a. lalok rarok nok izok hihok mok . |
| 6b. wat dat krat quat cat . | 12b. wat nnat forat arrat vat gat . |

## Slide 4 (bottom-right)

# Centauri/Arcturan [Knight, 1997]

Your assignment, put these words in order:  { jjat, arrat, mat, bat, oloat, at-yurp }

| | |
|---|---|
| 1a. ok-voon ororok sprok . | 7a. lalok farok ororok lalok sprok izok enemok . |
| 1b. at-voon bichat dat . | 7b. wat jjat bichat wat dat vat eneat . |
| 2a. ok-drubel ok-voon anok plok sprok . | 8a. lalok brok anok plok nok . |
| 2b. at-drubel at-voon pippat rrat dat . | 8b. iat lat pippat rrat nnat . |
| 3a. erok sprok izok hihok ghirok . | 9a. wiwok nok izok kantok ok-yurp . |
| 3b. totat dat arrat vat hilat . | 9b. totat nnat quat oloat at-yurp . |
| 4a. ok-voon anok drok brok jok . | 10a. lalok mok nok yorok ghirok clok . |
| 4b. at-voon krat pippat sat lat . | 10b. wat nnat gat mat bat hilat . |
| 5a. wiwok farok izok stok . | 11a. lalok nok crrrok hihok yorok zanzanok . zero fertility |
| 5b. totat jjat quat cat . | 11b. wat nnat arrat mat zanzanat . |
| 6a. lalok sprok izok jok stok . | 12a. lalok rarok nok izok hihok mok . |
| 6b. wat dat krat quat cat . | 12b. wat nnat forat arrat vat gat . |

## It's Really Spanish/English

**Clients do not sell pharmaceuticals in Europe => Clientes no venden medicinas en Europa**
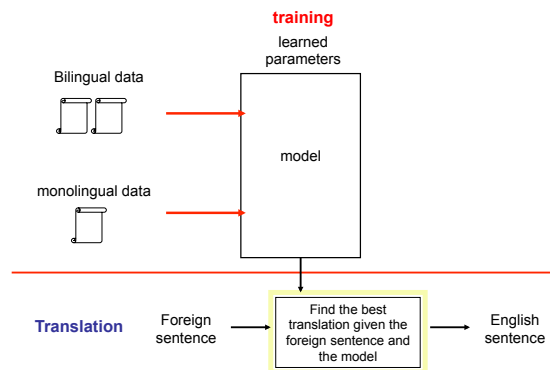
| | |
|---|---|
| 1a. Garcia and associates .<br>1b. Garcia y asociados . | 7a. the clients and the associates are enemies .<br>7b. los clients y los asociados son enemigos . |
| 2a. Carlos Garcia has three associates .<br>2b. Carlos Garcia tiene tres asociados . | 8a. the company has three groups .<br>8b. la empresa tiene tres grupos . |
| 3a. his associates are not strong .<br>3b. sus asociados no son fuertes . | 9a. its groups are in Europe .<br>9b. sus grupos estan en Europa . |
| 4a. Garcia has a company also .<br>4b. Garcia tambien tiene una empresa . | 10a. the modern groups sell strong pharmaceuticals .<br>10b. los grupos modernos venden medicinas fuertes . |
| 5a. its clients are angry .<br>5b. sus clientes estan enfadados . | 11a. the groups do not sell zenzanine .<br>11b. los grupos no venden zanzanina . |
| 6a. the associates are also angry .<br>6b. los asociados tambien estan enfadados . | 12a. the small groups are not modern .<br>12b. los grupos pequenos no son modernos . |



## Data available

- Many languages
  - Europarl corpus has all European languages
    - http://www.statmt.org/europarl/
    - From a few hundred thousand sentences to a few million
  - French/English from French parliamentary proceedings
  - Lots of Chinese/English and Arabic/English from government projects/interests
    - Chinese-English: 440 million words (15-20 million sentence pairs)
    - Arabic-English: 790 million words (30-40 million sentence pairs)
  - Smaller corpora in many, many other languages
- Lots of monolingual data available in many languages
- Even less data with multiple translations available
- Available in limited domains
  - most data is either news or government proceedings
  - some other domains recently, like blogs

## Statistical MT Overview

## Statistical MT

- We will model the translation process probabilistically

- Given a foreign sentence to translate, for any possible English sentence, we want to know the probability that sentence is a translation of the foreign sentence

- If we can find the most probable English sentence, we're done

  p(english sentence | foreign sentence)

## Noisy channel model

model $p(e \mid f) \propto p(f \mid e) p(e)$

translation model       language model

how do foreign sentences get translated to English sentences?    what do English sentences look like?



## Translation model

- The models define probabilities over inputs
$$p(f \mid e)$$

Morgen fliege ich nach Kanada zur Konferenz

Tomorrow I will fly to the conference in Canada

What is the probability that the English sentence is a translation of the foreign sentence?

## Translation model

- The models define probabilities over inputs

$$p(f \mid e)$$

| Morgen | fliege | ich | nach Kanada | zur Konferenz |

| Tomorrow | I | will fly | to the conference | In Canada |

- What is the probability of a foreign word being translated as a particular English word?
- What is the probability of a foreign foreign phrase being translated as a particular English phrase?
- What is the probability of a word/phrase changing ordering?
- What is the probability of a foreign word/phrase disappearing?
- What is the probability of a English word/phrase appearing?

## Translation model

- The models define probabilities over inputs

$$p(f \mid e)$$

p( Morgen fliege ich nach Kanada zur Konferenz |
Tomorrow I will fly to the conference in Canada )          = 0.1

p( Morgen fliege ich nach Kanada zur Konferenz |
I like peanut butter and jelly )          = 0.0001

## Language model

- The models define probabilities over inputs

$$p(e)$$

Tomorrow I will fly to the conference in Canada

## What is a probability distribution?

- A probability distribution defines the probability over a space of possible inputs
- For the language model, what is the space of possible inputs?
  - A language model describes the probability over ALL possible combinations of English words
- For the translation model, what is the space of possible inputs?
  - ALL possible combinations of foreign words with ALL possible combinations of English words
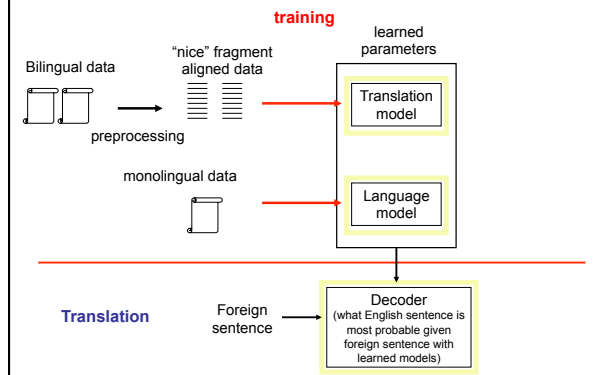
17

## One way to think about it…

**Spanish (foreign)** → [Translation model] → **Broken English** → [language model] → **English**

Que hambre tengo yo → What hunger have I, Hungry I am so, I am so hungry, Have I that hunger … → I am so hungry

---

## Translation

$$p(e \mid f) \propto p(f \mid e)p(e)$$

- Let's assume we have a translation model and a language model
- Given a foreign sentence, what question do we want to ask to translate that sentence into English?

$$\arg_e \max p(e \mid f) \propto p(f \mid e)p(e)$$

---

## Statistical MT Overview

**training**

Bilingual data → preprocessing → "nice" fragment aligned data → learned parameters → Translation model

monolingual data → Language model

**Translation**

Foreign sentence → Decoder (what English sentence is most probable given foreign sentence with learned models)

---

## Problems for Statistical MT

- Preprocessing
  - How do we get aligned bilingual text?
  - Tokenization
  - Segmentation (document, sentence, word)
- Language modeling
  - Given an English string e, assigns P(e) by formula
- Translation modeling
  - Given a pair of strings <f,e>, assigns P(f | e) by formula
- Decoding
  - Given a language model, a translation model, and a new sentence f … find translation e maximizing P(e) * P(f | e)
- Parameter optimization
  - Given a model with multiple feature functions, how are they related? What are the optimal paraeters?
- Evaluation
  - How well is a system doing?  How can we compare two systems?

Caution, butt head against the wall

## Problems for Statistical MT

- Preprocessing
- Language modeling
- Translation modeling
- Decoding
- Parameter optimization
- Evaluation

## Data

We want pairs of aligned sentences/text fragments. How do we get them?

## From No Data to Sentence Pairs

- Easy way 1: Linguistic Data Consortium (LDC)
- Easy way 2: pay $$$
  - Suppose one billion words of parallel data were sufficient
  - At 20 cents/word, that's $200 million
- Hard way: Find it, and then earn it!
  - De-formatting
  - Remove strange characters
  - Character code conversion
  - **Document alignment**
  - **Sentence alignment**
  - **Tokenization (also called Segmentation)**

## If you don't get the characters right…

**ISO-8859-2 (Latin2)**



**ISO-8859-6 (Arabic)**



## Chinese?

- **GB Code**
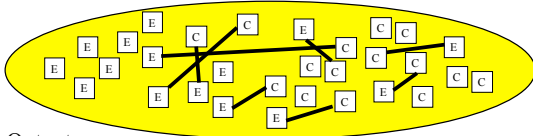- **GBK Code**
- **Big 5 Code**
- **CNS-11643-1992**
- …

## Document Alignment

- Input:
  - Big bag of files obtained from somewhere, believed to contain pairs of files that are translations of each other.



- Output:
  - List of pairs of files that are actually translations.

## Document Alignment

- Input:
  - Big bag of files obtained from somewhere, believed to contain pairs of files that are translations of each other.



- Output:
  - List of pairs of files that are actually translations.

## Sentence Alignment

The old man is happy. He has fished many times. His wife talks to him. The fish are jumping. The sharks await.

El viejo está feliz porque ha pescado muchos veces. Su mujer habla con él. Los tiburones esperan.

## Sentence Alignment

1. The old man is happy.
2. He has fished many times.
3. His wife talks to him.
4. The fish are jumping.
5. The sharks await.

1. El viejo está feliz porque ha pescado muchos veces.
2. Su mujer habla con él.
3. Los tiburones esperan.

## Sentence Alignment

1. The old man is happy.
2. He has fished many times.
3. His wife talks to him.
4. The fish are jumping.
5. The sharks await.

1. El viejo está feliz porque ha pescado muchos veces.
2. Su mujer habla con él.
3. Los tiburones esperan.

## Sentence Alignment

1. The old man is happy. He has fished many times. ——— 1. El viejo está feliz porque ha pescado muchos veces.
2. His wife talks to him. ——— 2. Su mujer habla con él.
3. The sharks await. ——— 3. Los tiburones esperan.

Note that unaligned sentences are thrown out, and sentences are merged in n-to-m alignments (n, m > 0).

## Tokenization (or Segmentation)

- English
  – Input (some byte stream):

            "There," said Bob.
  – Output (7 "tokens" or "words"):

            " There , " said Bob .
- Chinese
  – Input (byte stream):  美国关岛国际机场及其办公室均接获一名自称沙地阿拉伯富商拉登等发出的电子邮件。

  – Output:  美国 关岛国 际机 场 及其 办公 室 均接获 一名 自称 沙地 阿拉 伯富 商拉登 等发 出 的 电子邮件。

## Problems for Statistical MT

- Preprocessing
- Language modeling
- Translation modeling
- Decoding
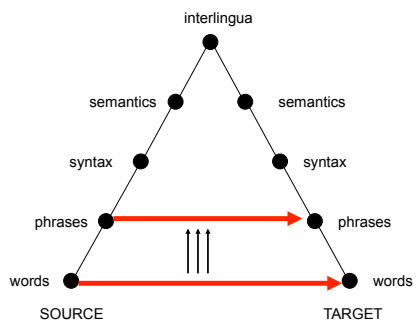- Parameter optimization
- Evaluation

## Language Modeling

- Most common: n-gram language models
- More data the better (Google n-grams)
- Domain is important

---

## Problems for Statistical MT

- Preprocessing
- Language modeling
- Translation modeling
- Decoding
- Parameter optimization
- Evaluation

---

## MT Pyramid

interlingua

semantics — semantics

syntax — syntax

phrases → phrases

words → words

SOURCE — TARGET

---

## Translation Model
### Learn How to Translate from Data

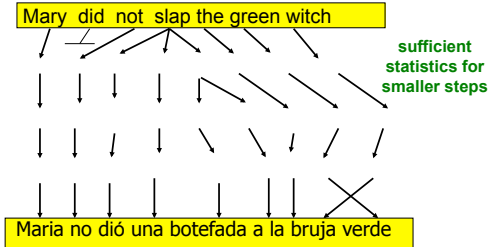**Direct Estimation:**

Mary did not slap the green witch

**not enough data for this
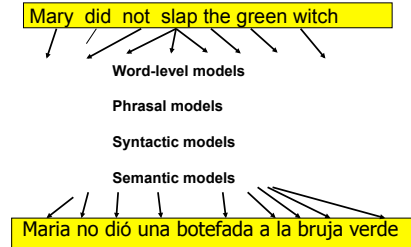(most input sentences unseen)**

Maria no dió una botefada a la bruja verde

---

## Generative Model

**Break up process into smaller steps:**

Mary  did  not  slap the green witch
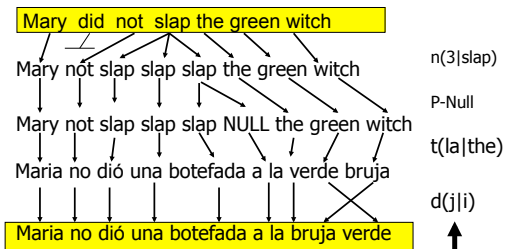
*sufficient statistics for smaller steps*

Maria no dió una botefada a la bruja verde

---

## What kind of Translation Model?

**May use syntactic and semantic representations:**

Mary  did  not  slap the green witch

**Word-level models**

**Phrasal models**

**Syntactic models**

**Semantic models**

Maria no dió una botefada a la bruja verde

---

## The Classic Translation Model

Word Substitution/Permutation [IBM Model 3, Brown et al., 1993]

**Generative story:**

Mary  did  not  slap the green witch

Mary not slap slap slap the green witch $n(3|slap)$

Mary not slap slap slap NULL the green witch $P\text{-Null}$

Maria no dió una botefada a la verde bruja $t(la|the)$

Maria no dió una botefada a la bruja verde $d(j|i)$

**Probabilities can be learned from raw bilingual text.**

---

## Phrase-Based Statistical MT

| Morgen | fliege | ich | nach Kanada | zur Konferenz |

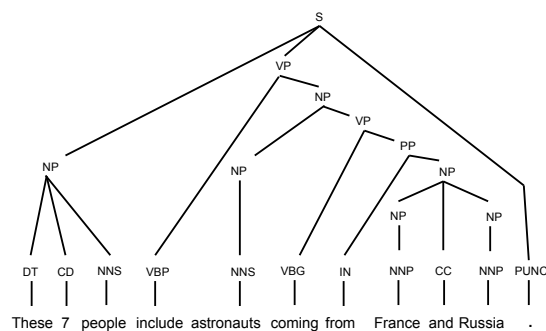| Tomorrow | I | will fly | to the conference | In Canada |

- Foreign input segmented in to phrases
  - "phrase" is any sequence of words
- Each phrase is probabilistically translated into English
  - P(to the conference | zur Konferenz)
  - P(into the meeting | zur Konferenz)
- Phrases are probabilistically re-ordered
- See [Koehn et al, 2003] for an intro.

## Advantages of Phrase-Based

- Many-to-many mappings can handle non-compositional phrases
- Easy to understand
- Local context is very useful for disambiguating
  - "Interest rate" → …
  - "Interest in" → …
- The more data, the longer the learned phrases
  - Sometimes whole sentences

## Syntax





## Problems for Statistical MT

- Preprocessing
- Language modeling
- Translation modeling
- **Decoding**
- Parameter optimization
- Evaluation

## Decoding

- Of all conceivable English word strings, find the one maximizing P(e) x P(f | e)

- Decoding is an NP-complete problem (for many translation models
  - (Knight, 1999)
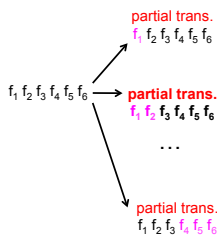
- Several decoding strategies are often available

## Search

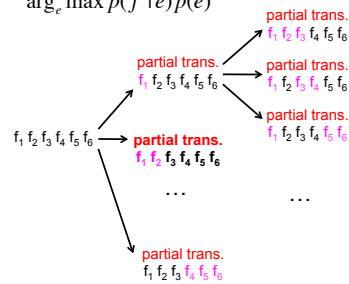$$\arg_e \max p(f \mid e)\, p(e)$$

f$_1$ f$_2$ f$_3$ f$_4$ f$_5$ f$_6$

## Search

$$\arg_e \max p(f \mid e)\, p(e)$$

partial trans.
f$_1$ f$_2$ f$_3$ f$_4$ f$_5$ f$_6$

f$_1$ f$_2$ f$_3$ f$_4$ f$_5$ f$_6$ → **partial trans.**
**f$_1$ f$_2$ f$_3$ f$_4$ f$_5$ f$_6$**

…

partial trans.
f$_1$ f$_2$ f$_3$ f$_4$ f$_5$ f$_6$

## Search

$$\arg_e \max p(f \mid e)\, p(e)$$

partial trans.
f$_1$ f$_2$ f$_3$ f$_4$ f$_5$ f$_6$

partial trans.
f$_1$ f$_2$ f$_3$ f$_4$ f$_5$ f$_6$ → partial trans.
f$_1$ f$_2$ f$_3$ f$_4$ f$_5$ f$_6$

partial trans.
f$_1$ f$_2$ f$_3$ f$_4$ f$_5$ f$_6$

f$_1$ f$_2$ f$_3$ f$_4$ f$_5$ f$_6$ → **partial trans.**
**f$_1$ f$_2$ f$_3$ f$_4$ f$_5$ f$_6$**

…          …

partial trans.
f$_1$ f$_2$ f$_3$ f$_4$ f$_5$ f$_6$

## Search

$$\arg_e \max p(f \mid e)p(e)$$

partial trans.
$f_1\ f_2\ f_3\ f_4\ f_5\ f_6$

partial trans.
$f_1\ f_2\ f_3\ f_4\ f_5\ f_6$

partial trans.
$f_1\ f_2\ f_3\ f_4\ f_5\ f_6$

$f_1\ f_2\ f_3\ f_4\ f_5\ f_6$

**partial trans.**
$f_1\ f_2\ f_3\ f_4\ f_5\ f_6$

partial trans.
$f_1\ f_2\ f_3\ f_4\ f_5\ f_6$

…        …

partial trans.
$f_1\ f_2\ f_3\ f_4\ f_5\ f_6$

- State space search problem
- Huge space
- Evaluate intermediate hypotheses using model and prune partial hypotheses

---

## Problems for Statistical MT

- Preprocessing
- Language modeling
- Translation modeling
- Decoding
- Parameter optimization
- Evaluation

---

## Basic Model, Revisited

argmax  P(e | f)  =
  e

argmax  P(e) x P(f | e) / P(f)  =
  e

argmax  P(e) x P(f | e)
  e

---

## Basic Model, Revisited

argmax  P(e | f)  =
  e

argmax  P(e) x P(f | e) / P(f)  =
  e

argmax  P(e)$^{2.4}$ x P(f | e)      … works better!
  e

## Basic Model, Revisited

argmax  P(e | f)  =
  e

argmax  P(e) x P(f | e) / P(f)
  e

argmax  P(e)$^{2.4}$ x P(f | e) x length(e)$^{1.1}$
  e

Rewards longer hypotheses, since
these are unfairly punished by P(e)

## Basic Model, Revisited

argmax  P(e)$^{2.4}$ x P(f | e) x length(e)$^{1.1}$ x KS $^{3.7}$ …
  e

Lots of knowledge sources vote on any given hypothesis.

"Knowledge source" = "feature function" = "score component".

A feature function simply scores a hypothesis with a real value.

(May be binary, as in "e has a verb").

## The Problem:  Learn Lambdas

$$p(e \mid f) = \frac{p(f \mid e)\, p(e)}{p(f)}$$

$$= \frac{p(f \mid e)^{\lambda_1}\, p(e)^{\lambda_2}}{\sum_{e'} p(f \mid e')^{\lambda_1} \lambda_2\, p(e')^{\lambda_2}}$$

$$= \frac{p(f \mid e)^{\lambda_1}\, p(e)^{\lambda_2}\, p(e \mid f)^{\lambda_3}\, length(e)^{\lambda_4}\ldots}{\sum_{e'} p(f \mid e')^{\lambda_1}\, p(e')^{\lambda_2}\, p(e' \mid f)^{\lambda_3}\, length(e')^{\lambda_4}\ldots}$$

$$= \frac{\exp\left(\lambda_1 \log p(f \mid e) + \lambda_2 \log p(e) + \lambda_3 \log p(e \mid f) + \lambda_4 length(e)\ldots\right)}{\sum_{e'} \exp\left(\lambda_1 \log p(f \mid e') + \lambda_2 \log p(e') + \lambda_3 \log p(e' \mid f) + \lambda_4 length(e')\ldots\right)}$$

$$= \frac{\exp\left(\sum_i \lambda_i h_i(f, e)\right)}{\sum_{e'} \exp\left(\sum_i \lambda_i h_i(f, e')\right)}$$

Given a data set with foreign/English
sentences, find the λ's that:
• maximize the likelihood of the data
• maximize an evaluation criterion