

# Probability and features

# Bayes' Rule

$$p(\textit{English}|\textit{Chinese}) \sim$$

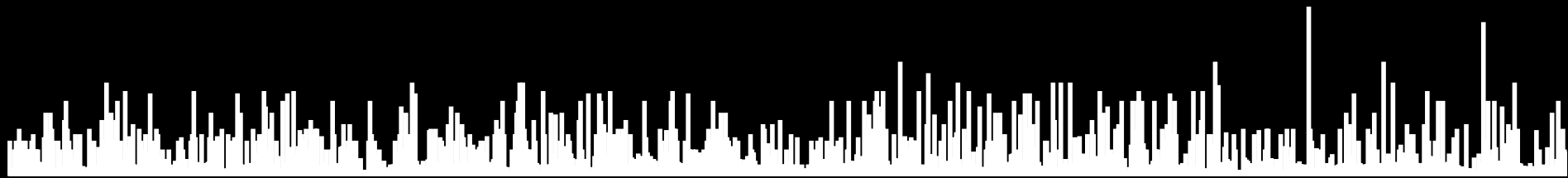
$$p(\textit{English}) \times p(\textit{Chinese}|\textit{English})$$

language model



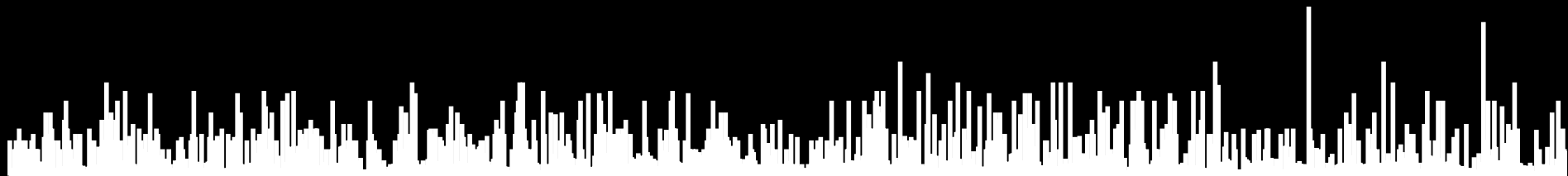
translation model

$p(\textit{Chinese}|\textit{English})$

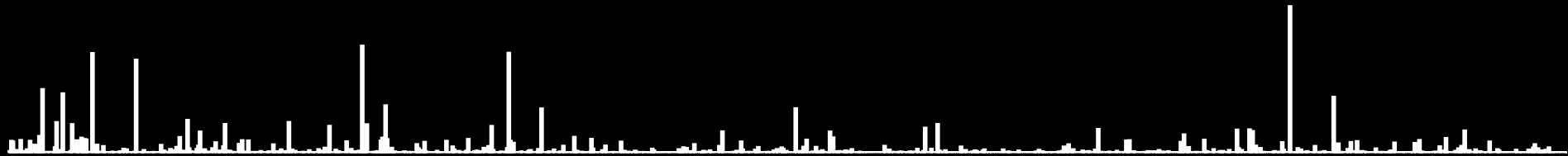


*English*

$p(\textit{Chinese}|\textit{English})$

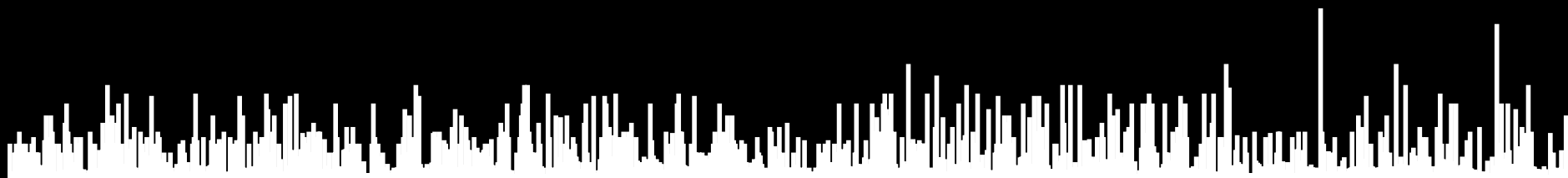


$\times p(\textit{English})$

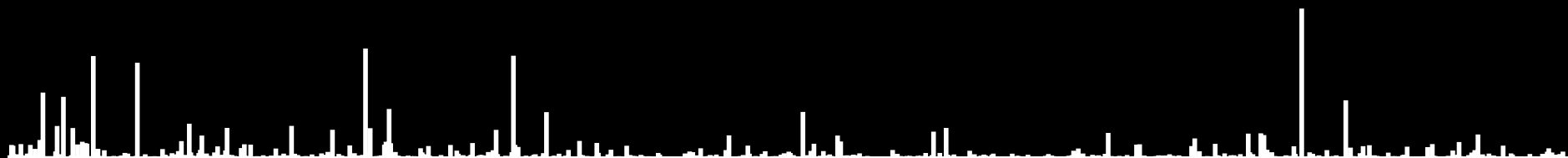


$\textit{English}$

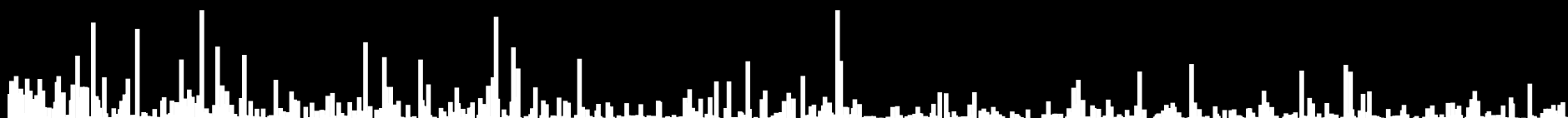
$p(\textit{Chinese}|\textit{English})$



$\times p(\textit{English})$

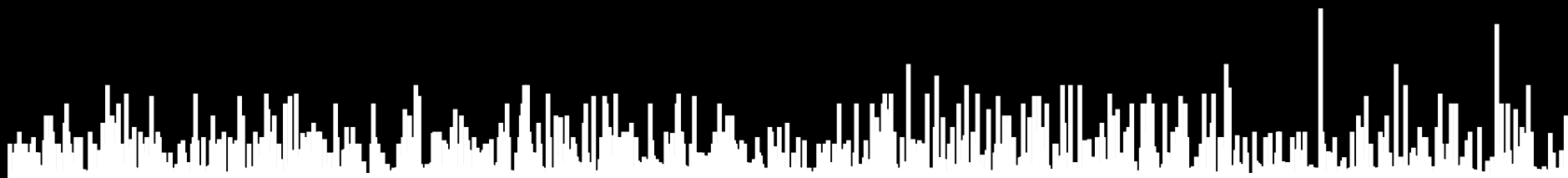


$\sim p(\textit{English}|\textit{Chinese})$

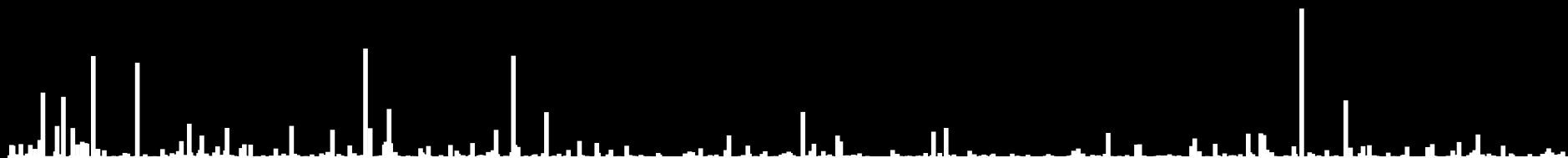


$\textit{English}$

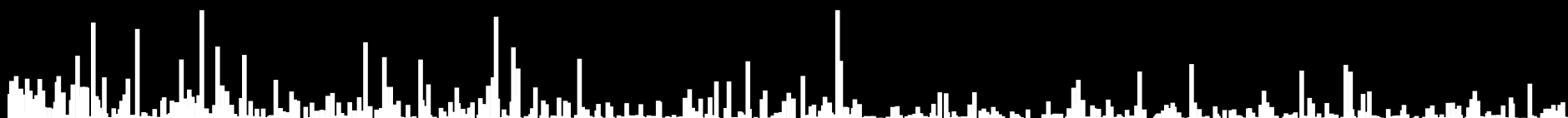
$$p(\textit{Chinese}|\textit{English})^1$$



$$\times p(\textit{English})^1$$

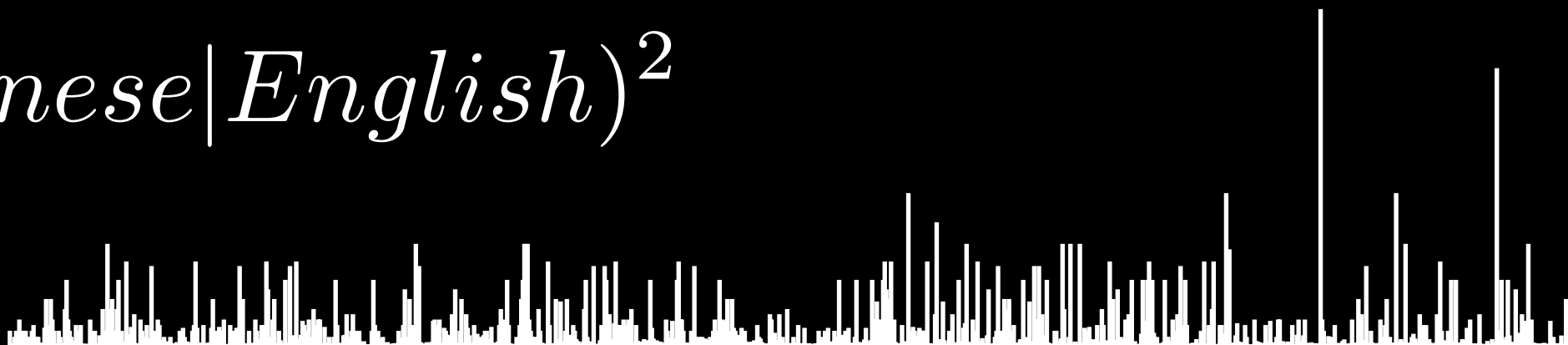


$$\sim p(\textit{English}|\textit{Chinese})$$

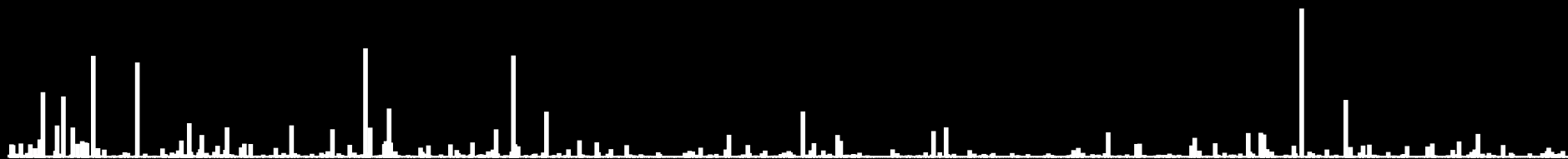


*English*

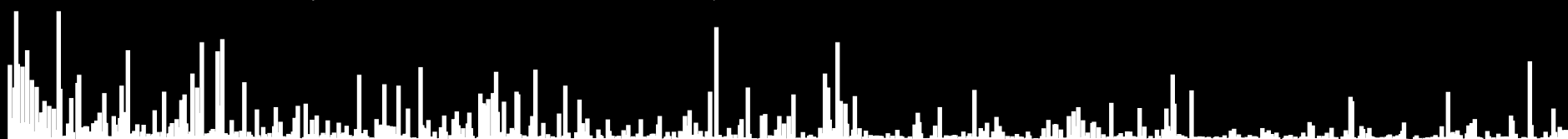
$$p(\textit{Chinese}|\textit{English})^2$$



$$\times p(\textit{English})^1$$



$$\sim p(\textit{English}|\textit{Chinese})$$

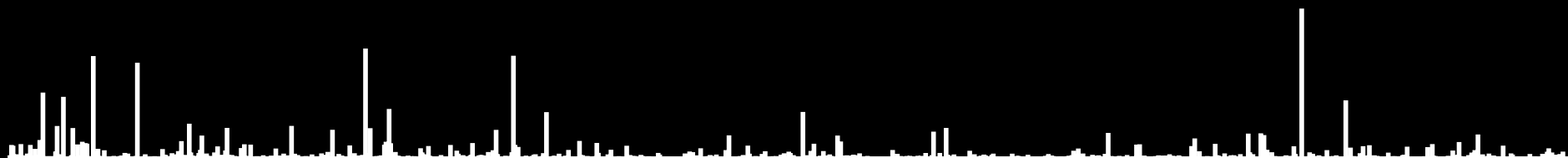


*English*

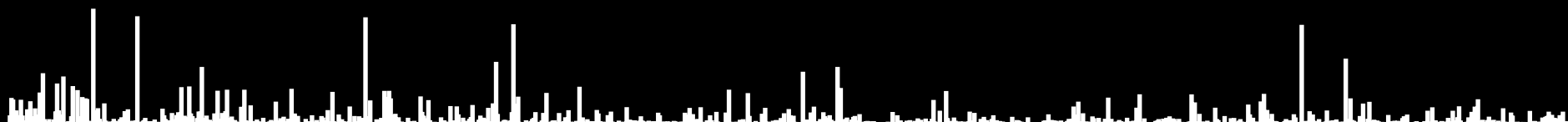
$$p(\textit{Chinese}|\textit{English})^{1/2}$$



$$\times p(\textit{English})^1$$



$$\sim p(\textit{English}|\textit{Chinese})$$



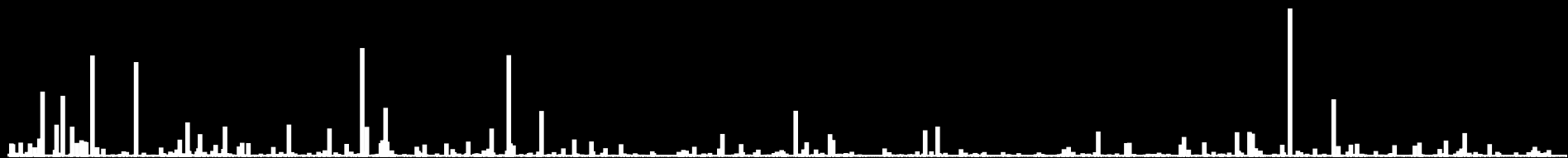
*English*



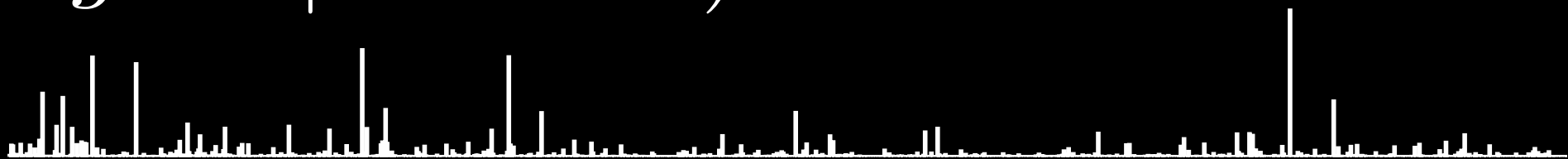
$$p(\textit{Chinese}|\textit{English})^0$$



$$\times p(\textit{English})^1$$



$$\sim p(\textit{English}|\textit{Chinese})$$

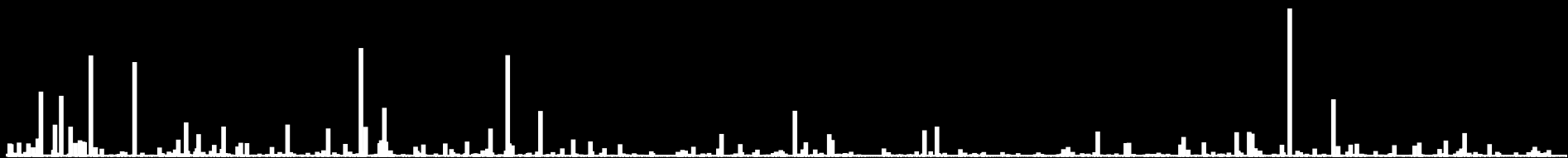


*English*

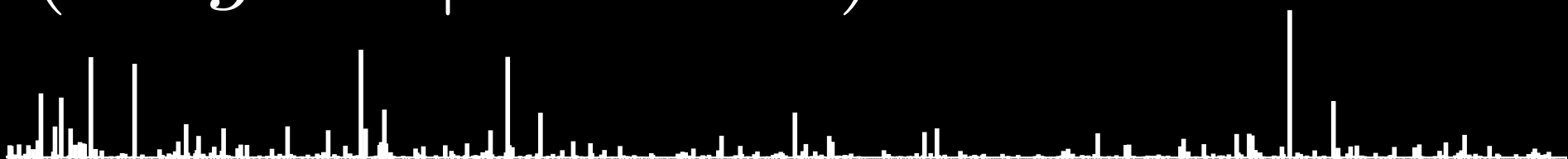
$$0 \cdot \log p(\textit{Chinese}|\textit{English})$$



$$+1 \cdot \log p(\textit{English})$$



$$\sim \log p(\textit{English}|\textit{Chinese})$$



*English*

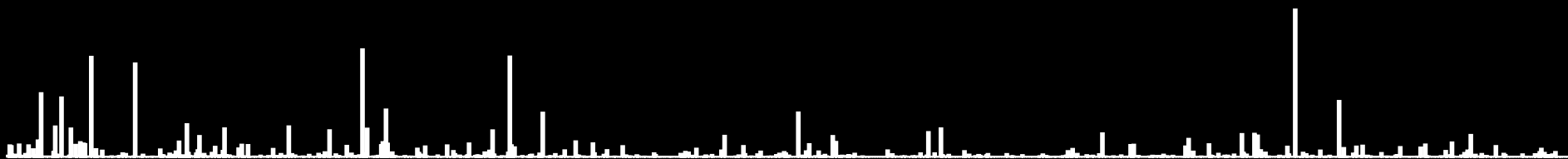
$\log(x)$  is monotonic for positive  $x$ :

$$\log(x) > \log(y) \text{ iff } x > y$$

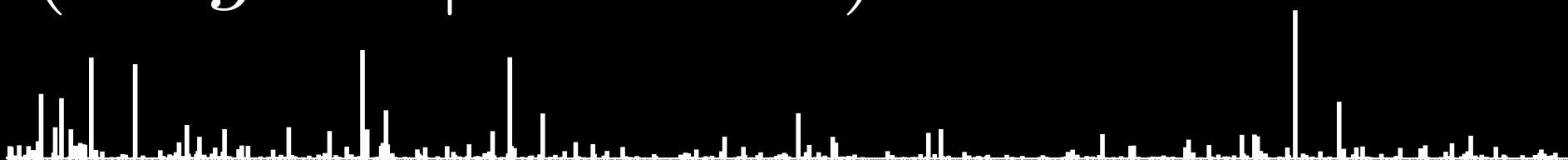
$$0 \cdot \log p(\textit{Chinese}|\textit{English})$$



$$+1 \cdot \log p(\textit{English})$$



$$\sim \log p(\textit{English}|\textit{Chinese})$$



*English*

$\log(x)$  is monotonic for positive  $x$ :

$\log(x) > \log(y)$  iff  $x > y$

$0 \cdot \log p(\text{Chinese} | \text{English})$

Local historic footnote: logarithms were invented in Edinburgh by John Napier, whose ancient family home in Merchiston is now part of Edinburgh Napier University.



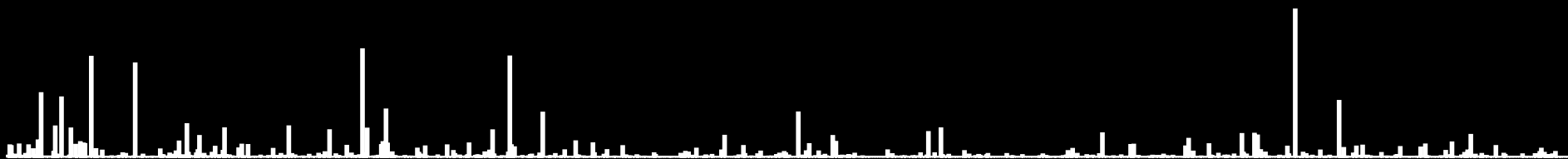
$\log(x)$  is monotonic for positive  $x$ :

$$\log(x) > \log(y) \text{ iff } x > y$$

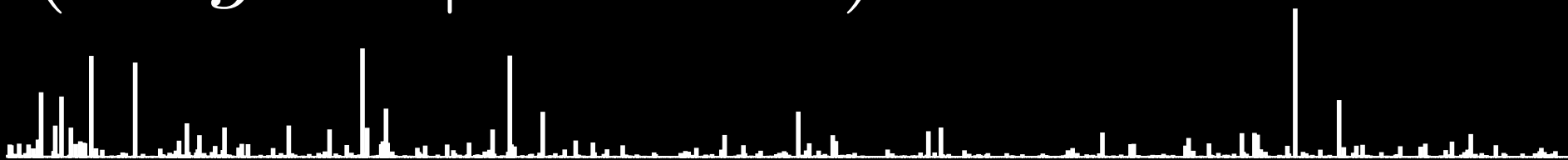
$$0 \cdot \log p(\textit{Chinese}|\textit{English})$$



$$+1 \cdot \log p(\textit{English})$$



$$\sim \log p(\textit{English}|\textit{Chinese})$$

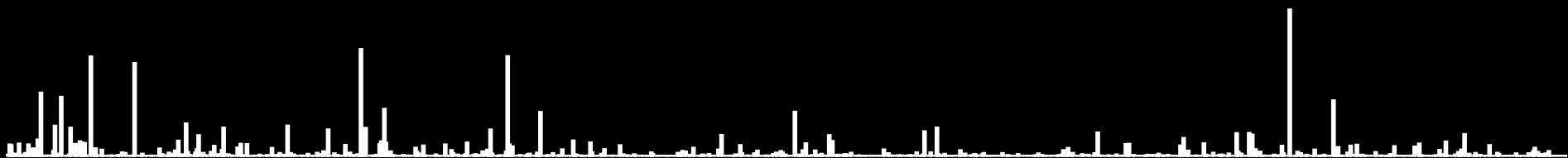


*English*

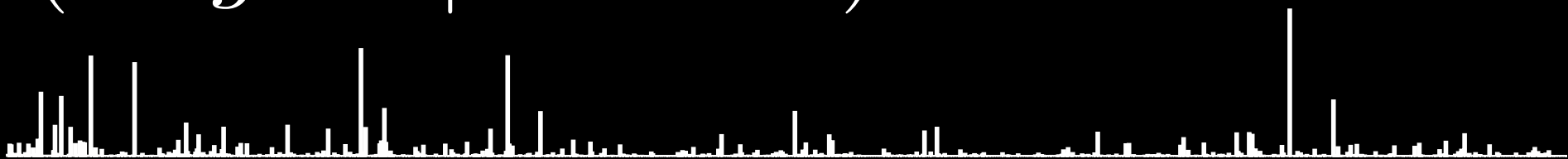
$$0 \cdot \log p(\textit{Chinese}|\textit{English})$$



$$+1 \cdot \log p(\textit{English})$$



$$= \textit{score}(\textit{English}|\textit{Chinese})$$



*English*

$$\begin{aligned} \textit{score}(\textit{English}|\textit{Chinese}) = \\ \lambda_1 \log p(\textit{Chinese}|\textit{English}) + \lambda_2 \log p(\textit{English}) \end{aligned}$$

$$\begin{aligned} \textit{score}(\textit{English}|\textit{Chinese}) = \\ \exp(\lambda_1 \log p(\textit{Chinese}|\textit{English}) + \lambda_2 \log p(\textit{English})) \end{aligned}$$



$$p(\textit{English}|\textit{Chinese}) = \frac{\exp(\lambda_1 \log p(\textit{Chinese}|\textit{English}) + \lambda_2 \log p(\textit{English}))}{\sum_{n_{\textit{English}}} \exp(\lambda_1 \log p(\textit{Chinese}|\textit{English}) + \lambda_2 \log p(\textit{English}))}$$

$$p(\textit{English}|\textit{Chinese}) = \frac{\exp(\lambda_1 \log p(\textit{Chinese}|\textit{English}) + \lambda_2 \log p(\textit{English}))}{\sum_{n_{\textit{English}}} \exp(\lambda_1 \log p(\textit{Chinese}|\textit{English}) + \lambda_2 \log p(\textit{English}))}$$

log-linear model  
maximum entropy model  
conditional random field  
undirected model

$$p(\textit{English}|\textit{Chinese}) =$$

$$p(\textit{English}) \times p(\textit{Chinese}|\textit{English})$$

Note: Original model is a special case of this model!

log-linear model

maximum entropy model

conditional random field

undirected model

$$p(\textit{English}|\textit{Chinese}) = \frac{\exp(\lambda_1 \log p(\textit{Chinese}|\textit{English}) + \lambda_2 \log p(\textit{English}))}{\sum_{n_{\textit{English}}} \exp(\lambda_1 \log p(\textit{Chinese}|\textit{English}) + \lambda_2 \log p(\textit{English}))}$$

log-linear model  
maximum entropy model  
conditional random field  
undirected model

$$p(\textit{English}|\textit{Chinese}) =$$

$$\exp \left\{ \sum_k \lambda_k h_k(\textit{English}, \textit{Chinese}) \right\}$$


---

$$\sum_{\textit{English}'} \exp \left\{ \sum_k \lambda_k h_k(\textit{English}', \textit{Chinese}) \right\}$$

log-linear model

maximum entropy model

conditional random field

undirected model

$$p(\textit{English}|\textit{Chinese}) = \frac{1}{Z} \exp \left\{ \sum_k \lambda_k h_k(\textit{English}, \textit{Chinese}) \right\}$$

log-linear model  
maximum entropy model  
conditional random field  
undirected model

$$p(\textit{English}|\textit{Chinese}) = \frac{1}{Z} \exp \left\{ \sum_k \lambda_k h_k(\textit{English}, \textit{Chinese}) \right\}$$

Z is the normalization term or *partition function*

log-linear model  
maximum entropy model  
conditional random field  
undirected model

$$p(\textit{English}|\textit{Chinese}) = \frac{1}{Z} \exp \left\{ \sum_k \lambda_k h_k(\textit{English}, \textit{Chinese}) \right\}$$

$Z$  is the normalization term or *partition function*

The functions  $h_k$  are *features* or *feature functions*

They are deterministic (fixed) functions of the  
input/output pair.

The parameters of the model are the  $\lambda_k$  terms.



# What's a Feature?

# What's a Feature?

A feature can be *any* function in the form:

$$h_k : \textit{English} \times \textit{Chinese} \rightarrow \mathbb{R}_+$$

# What's a Feature?

A feature can be *any* function in the form:

$$h_k : \textit{English} \times \textit{Chinese} \rightarrow \mathbb{R}_+$$

- Language model:  $p(\textit{English})$

# What's a Feature?

A feature can be *any* function in the form:

$$h_k : \textit{English} \times \textit{Chinese} \rightarrow \mathbb{R}_+$$

- Language model:  $p(\textit{English})$
- Translation model:  $p(\textit{Chinese} \mid \textit{English})$

# What's a Feature?

A feature can be *any* function in the form:

$$h_k : \textit{English} \times \textit{Chinese} \rightarrow \mathbb{R}_+$$

- Language model:  $p(\textit{English})$
- Translation model:  $p(\textit{Chinese} \mid \textit{English})$
- Reverse translation model:  $p(\textit{English} \mid \textit{Chinese})$

# What's a Feature?

A feature can be *any* function in the form:

$$h_k : \textit{English} \times \textit{Chinese} \rightarrow \mathbb{R}_+$$

- Language model:  $p(\textit{English})$
- Translation model:  $p(\textit{Chinese} \mid \textit{English})$
- Reverse translation model:  $p(\textit{English} \mid \textit{Chinese})$
- The number of words in the English sentence.

# What's a Feature?

A feature can be *any* function in the form:

$$h_k : \textit{English} \times \textit{Chinese} \rightarrow \mathbb{R}_+$$

- Language model:  $p(\textit{English})$
- Translation model:  $p(\textit{Chinese} \mid \textit{English})$
- Reverse translation model:  $p(\textit{English} \mid \textit{Chinese})$
- The number of words in the English sentence.
- The number of verbs in the English sentence.

# What's a Feature?

A feature can be *any* function in the form:

$$h_k : \textit{English} \times \textit{Chinese} \rightarrow \mathbb{R}_+$$

- Language model:  $p(\textit{English})$
- Translation model:  $p(\textit{Chinese} \mid \textit{English})$
- Reverse translation model:  $p(\textit{English} \mid \textit{Chinese})$
- The number of words in the English sentence.
- The number of verbs in the English sentence.
- 1 if the English sentence has a verb, 0 otherwise.



# What's a Feature?

A feature can be *any* function in the form:

$$h_k : \textit{English} \times \textit{Chinese} \rightarrow \mathbb{R}_+$$

# What's a Feature?

A feature can be *any* function in the form:

$$h_k : \textit{English} \times \textit{Chinese} \rightarrow \mathbb{R}_+$$

- A word-based translation model:  $p(\textit{Chinese} \mid \textit{English})$

# What's a Feature?

A feature can be *any* function in the form:

$$h_k : \textit{English} \times \textit{Chinese} \rightarrow \mathbb{R}_+$$

- A word-based translation model:  $p(\textit{Chinese} \mid \textit{English})$
- Agreement features in the English sentence.

# What's a Feature?

A feature can be *any* function in the form:

$$h_k : \textit{English} \times \textit{Chinese} \rightarrow \mathbb{R}_+$$

- A word-based translation model:  $p(\textit{Chinese} | \textit{English})$
- Agreement features in the English sentence.
- Features over part-of-speech sequences in the English sentence.

# What's a Feature?

A feature can be *any* function in the form:

$$h_k : \textit{English} \times \textit{Chinese} \rightarrow \mathbb{R}_+$$

- A word-based translation model:  $p(\textit{Chinese} \mid \textit{English})$
- Agreement features in the English sentence.
- Features over part-of-speech sequences in the English sentence.
- How many times the sentence pair includes the English word *north* and Chinese word 北.

# What's a Feature?

A feature can be *any* function in the form:

$$h_k : \textit{English} \times \textit{Chinese} \rightarrow \mathbb{R}_+$$

- A word-based translation model:  $p(\textit{Chinese} \mid \textit{English})$
- Agreement features in the English sentence.
- Features over part-of-speech sequences in the English sentence.
- How many times the sentence pair includes the English word *north* and Chinese word 北.
- Do words *north* and 北 appear in a dictionary?

# Probability again

$$p(\textit{English}|\textit{Chinese}) = \frac{1}{Z} \exp \left\{ \sum_k \lambda_k h_k(\textit{English}, \textit{Chinese}) \right\}$$

# Probability again

$$p(y|x) = \frac{1}{Z} \exp \left\{ \sum_k \lambda_k h_k(x, y) \right\}$$



# Probability again

$x$  and  $y$  can be anything here.

Suppose:  $y = w_i$   
 $x = w_{i-1} \dots w_{i-n+1}$

This is a *maximum  
entropy language  
model*

$$p(y|x) = \frac{1}{Z} \exp(\lambda^\top \cdot h_k(x, y))$$

Now design  $h$  so that it is sensitive to morphology