

# Basic HMM for POS Induction

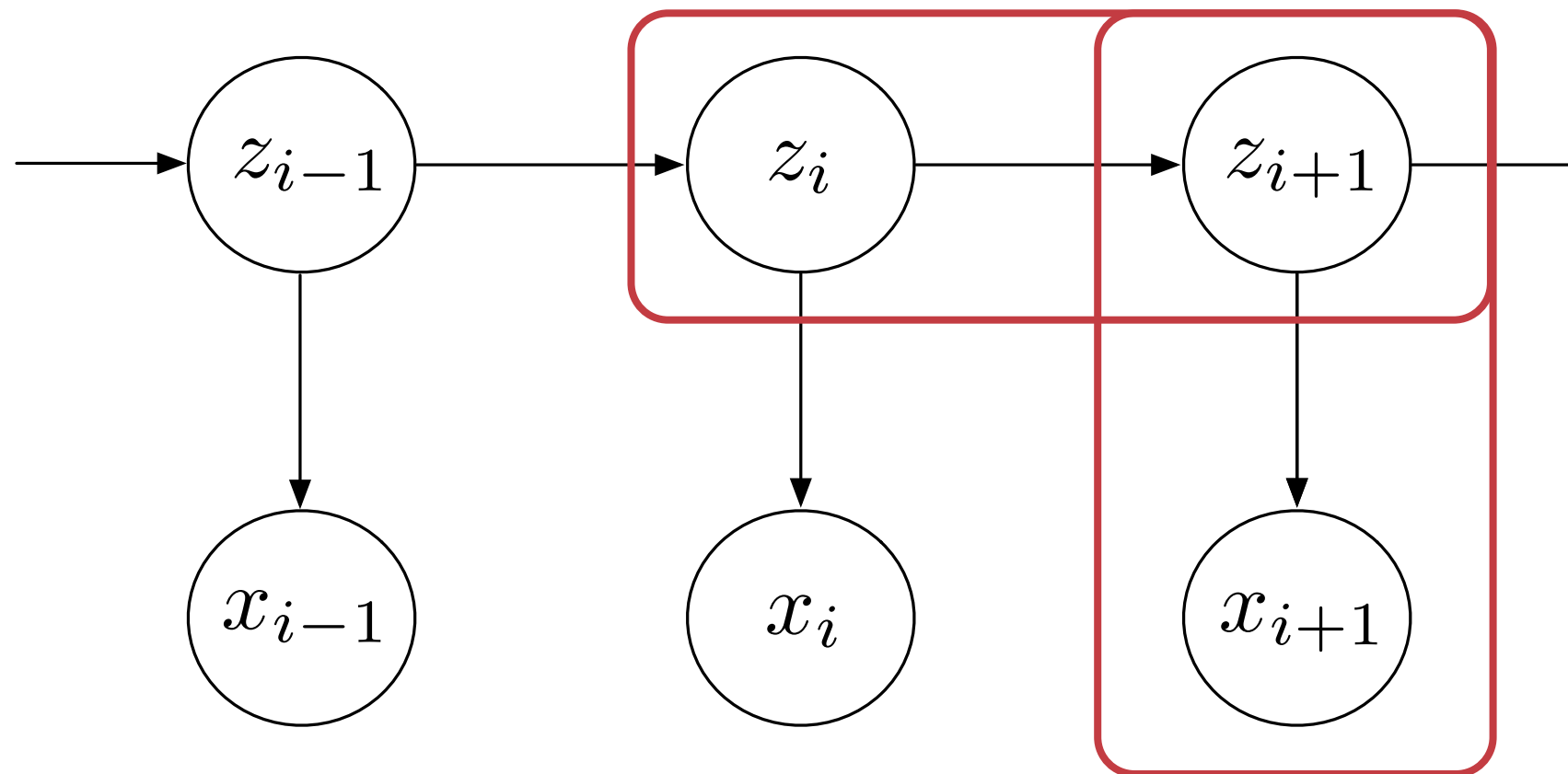
---

Transition distribution:

$$P(z' | z)$$

Emission distribution:

$$P(x | z)$$



# Parameterization

Key distribution:  $P(x|\text{NNP})$

**W:**

+Cap	+2
+ing	-1

$\theta_{x \text{NNP}}$	$x$	$\mathbf{f}$	$e^{\mathbf{w}^T \mathbf{f}}$
0.1	John	+Cap	0.3
0	Mary	+Cap	0.3
0.2	running	+ing	0.1
0	jumping	+ing	0.1

# Parameterization

---



$$\theta_{x|z} = \frac{\exp(\mathbf{w}^\top \mathbf{f}(x, z))}{\sum_{x'} \exp(\mathbf{w}^\top \mathbf{f}(x', z))}$$

# Unsupervised Learning with Features

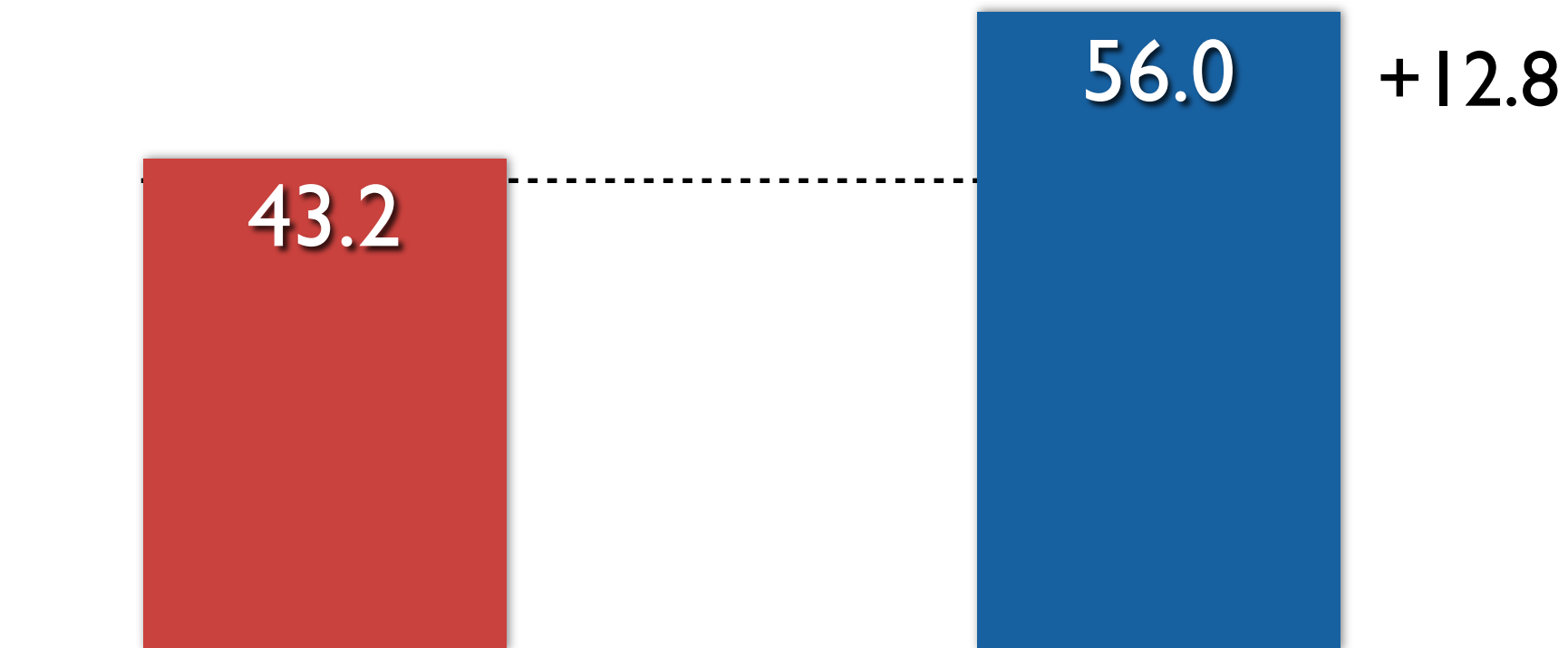
---

- Main idea:  
Local multinomials become maxents
- **EM** + **Maxent M-Step** =  
Unsupervised Learning w/ features

# POS Induction Accuracy

---

I-to-I Accuracy



**Basic Multinomial:**

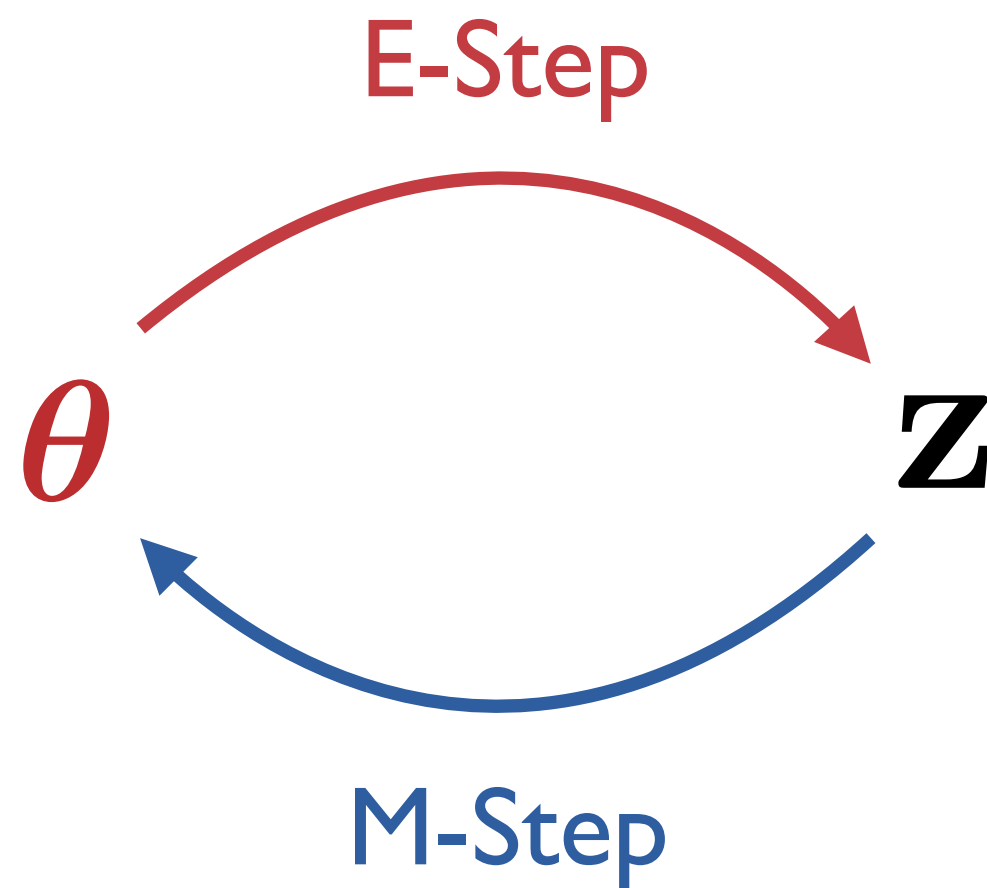
John  $\wedge$  NNP

**Rich Features:**

John  $\wedge$  NNP  
+Cap  $\wedge$  NNP  
+Digit  $\wedge$  NNP  
+Hyphen  $\wedge$  NNP  
+ing  $\wedge$  NNP

# Basic Hard EM

---



# Basic Hard EM

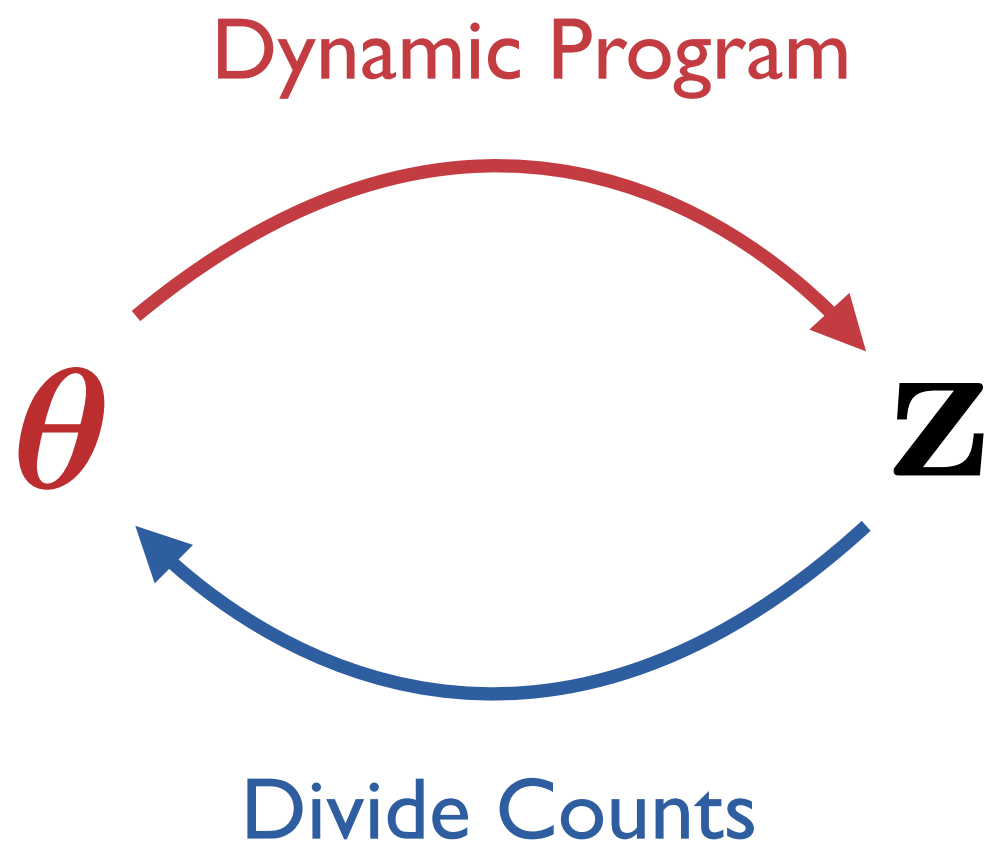
---

## E-Step: Dynamic Program

$$\mathbf{z} \leftarrow \operatorname{argmax}_{\mathbf{z}} P(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta})$$

## M-Step: Divide Counts

$$\begin{aligned} \boldsymbol{\theta} &\leftarrow \operatorname{argmax}_{\boldsymbol{\theta}} P(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) \\ &= \left[ \frac{c(z \rightarrow x)}{c(z \rightarrow \cdot)}, \dots \right] \end{aligned}$$



# Hard EM with Features

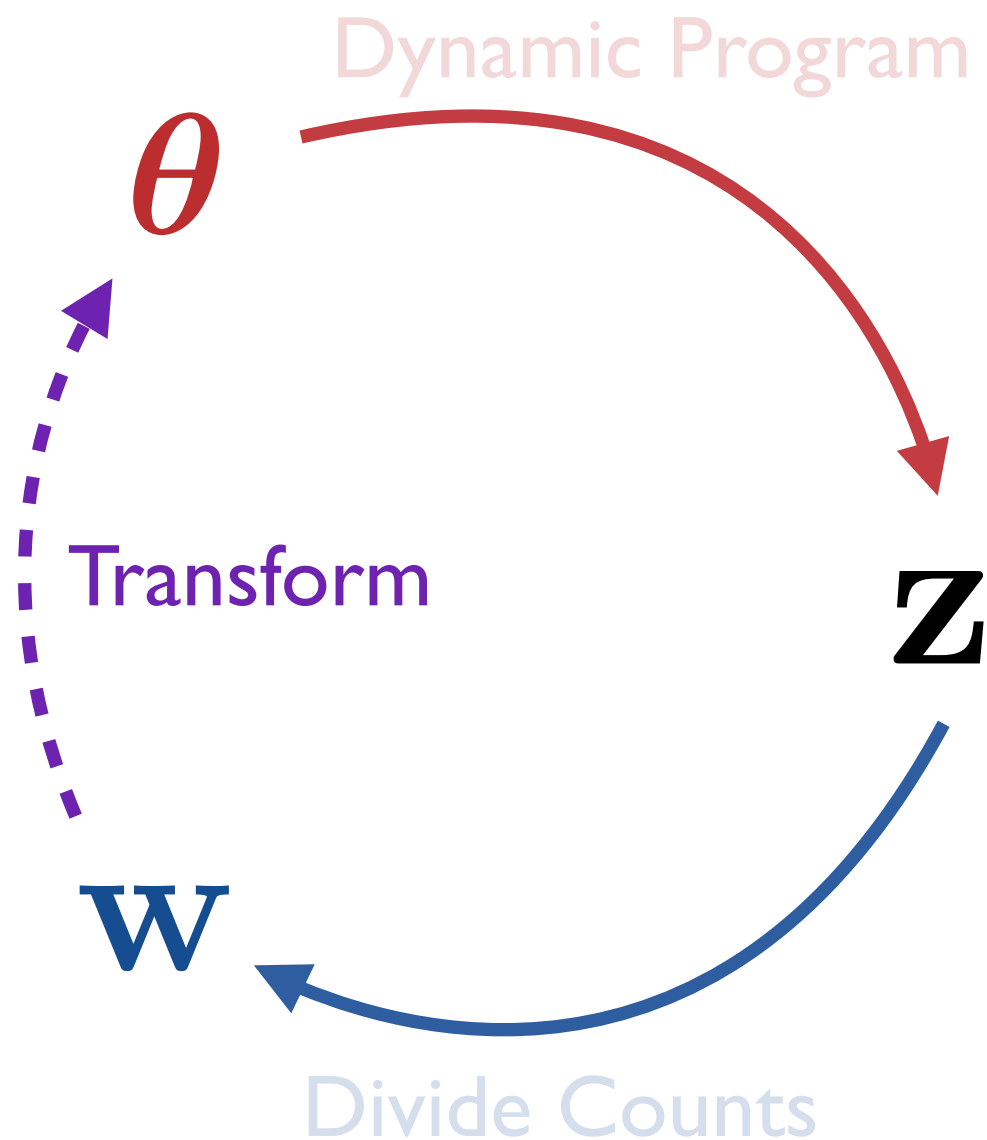
---

## E-Step: Dynamic Program

$$\mathbf{z} \leftarrow \operatorname{argmax}_{\mathbf{z}} P(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta})$$

## M-Step: Divide Counts

$$\begin{aligned} \boldsymbol{\theta} &\leftarrow \operatorname{argmax}_{\boldsymbol{\theta}} P(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) \\ &= \left[ \frac{c(\mathbf{z} \rightarrow \mathbf{x})}{c(\mathbf{z} \rightarrow \cdot)}, \dots \right] \end{aligned}$$





# Hard EM with Features

## E-Step: Dynamic Program

$$\mathbf{z} \leftarrow \operatorname{argmax}_{\mathbf{z}} P(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta})$$

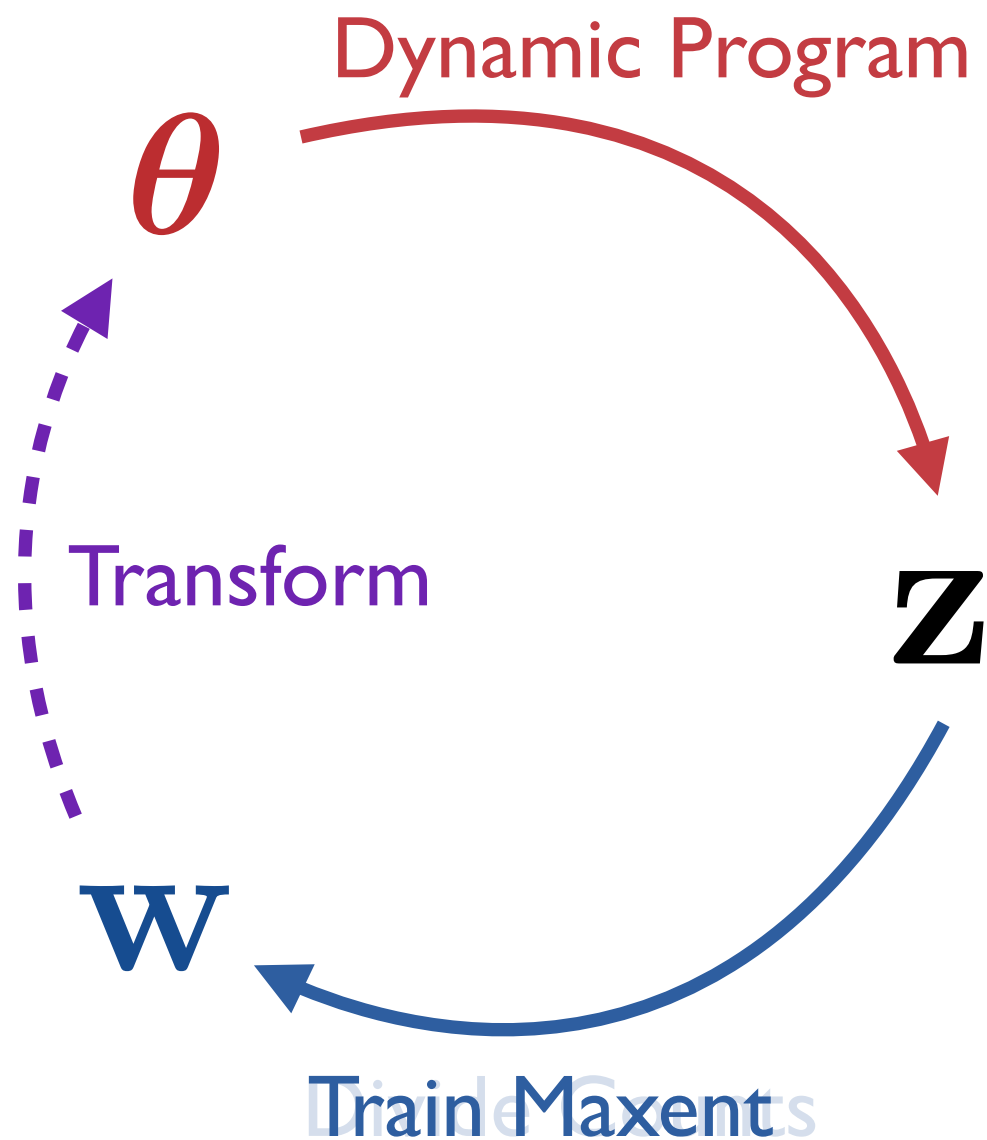
## M-Step: Train Maxent

$$\mathbf{w} \leftarrow \operatorname{argmax}_{\mathbf{w}} \log P(\mathbf{x}, \mathbf{z}; \mathbf{w})$$

## M-Step: Divide Counts

$$\boldsymbol{\theta} \leftarrow \operatorname{argmax}_{\boldsymbol{\theta}} P(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta})$$

$$= \left[ \frac{c(z \rightarrow x)}{c(z \rightarrow \cdot)}, \dots \right]$$



# Hard EM with Features

---

$$\log P(\mathbf{x}, \mathbf{z}; \mathbf{w})$$

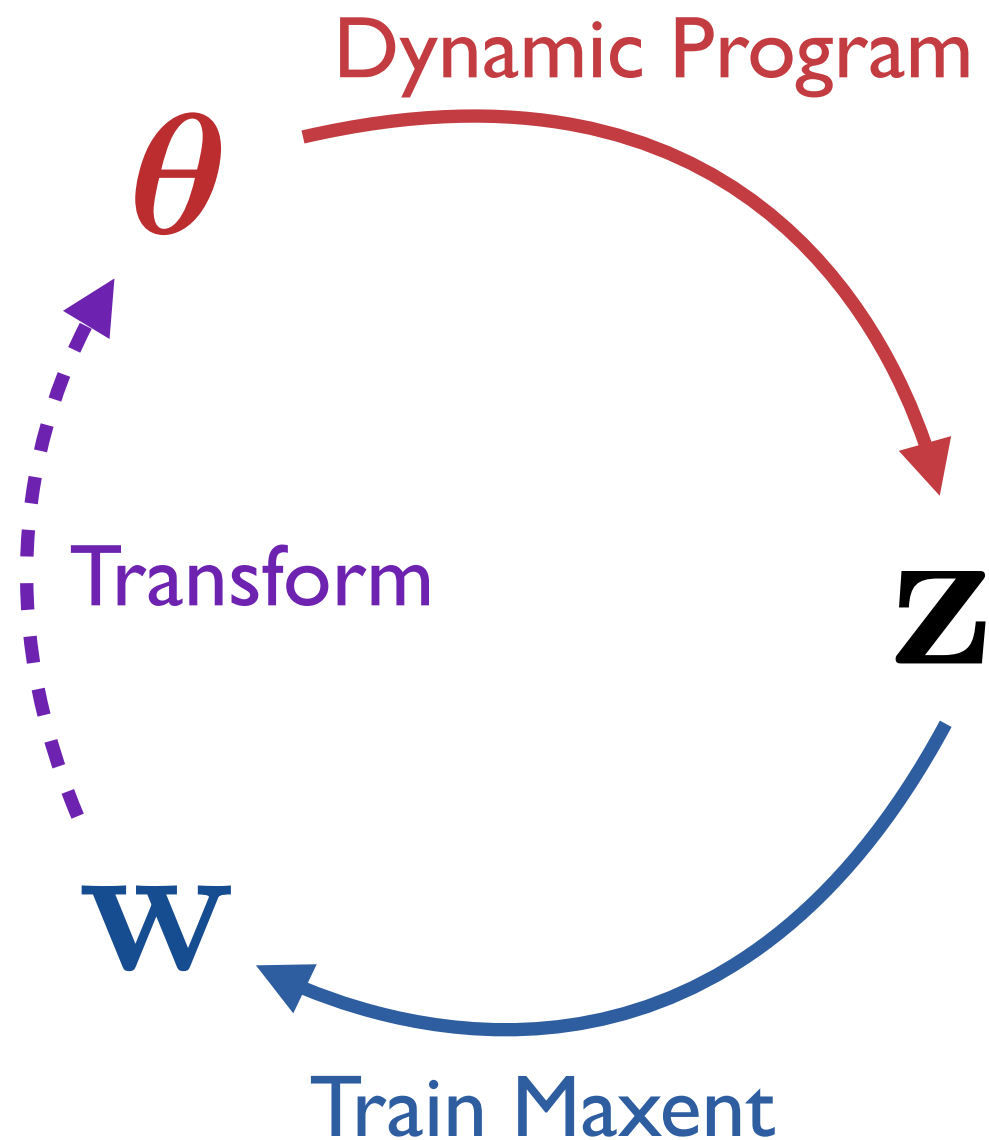
$$= \sum_i \log P(x_i | z_i; \mathbf{w}) + \dots$$

Maxent training example

$$= \sum_{z, x} c(z \rightarrow x) \log P(x | z; \mathbf{w}) + \dots$$

Multiplicity

# Hard EM with Features



## E-Step: Dynamic Program

$$\mathbf{z} \leftarrow \operatorname{argmax}_{\mathbf{z}} P(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta})$$

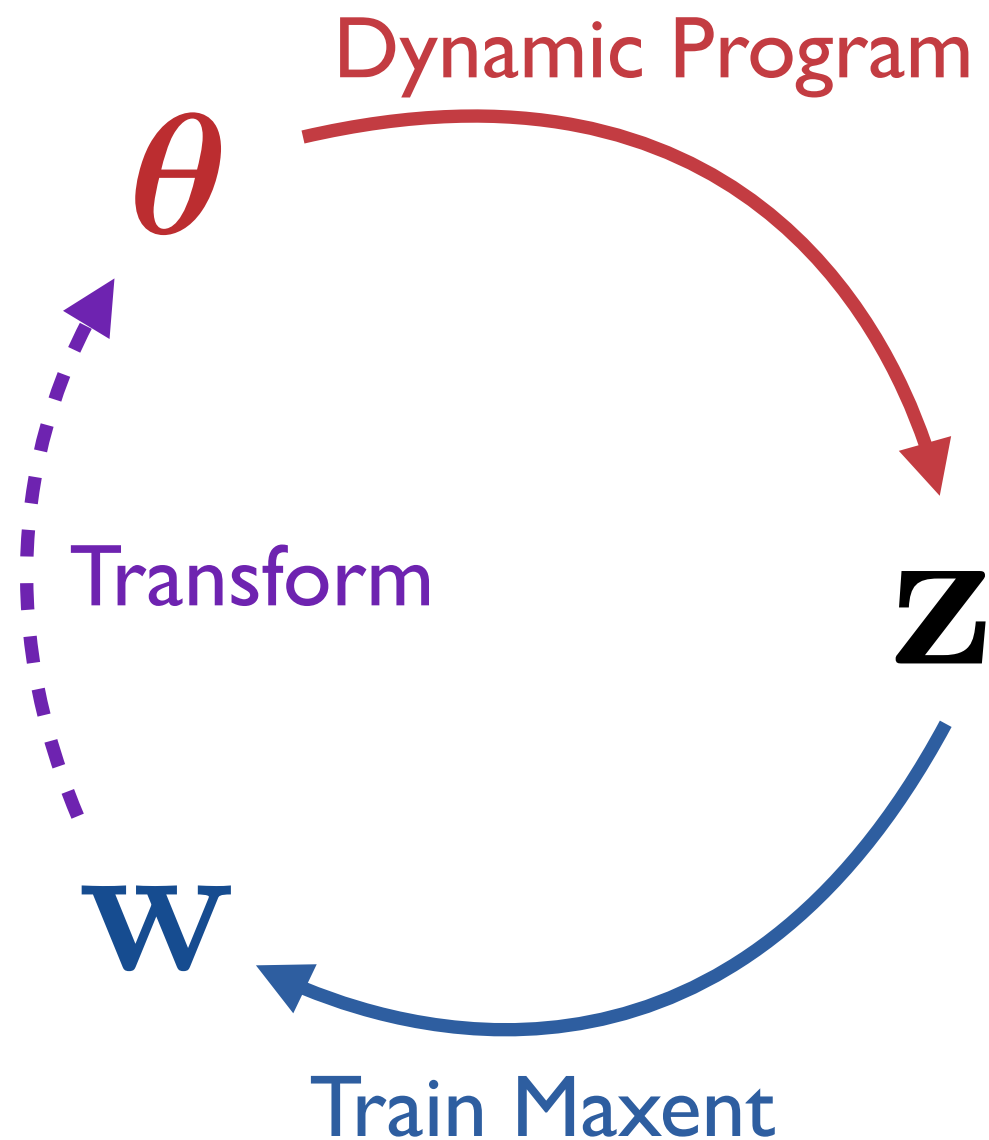
## M-Step: Train Maxent

$$\mathbf{w} \leftarrow \operatorname{argmax}_{\mathbf{w}} \log P(\mathbf{x}, \mathbf{z}; \mathbf{w})$$

## Transform parameters

$$\theta_{x|z} \leftarrow \frac{\exp(\mathbf{w}^T \mathbf{f}(x, z))}{\sum_{x'} \exp(\mathbf{w}^T \mathbf{f}(x', z))}$$

# EM with Features



## E-Step: Dynamic Program

$$e(z \rightarrow x) \leftarrow \mathbb{E}_{\mathbf{z}} [c(z \rightarrow x; \theta)]$$

## M-Step: Train Maxent

$$\mathbf{w} \leftarrow \operatorname{argmax}_{\mathbf{w}} \mathbb{E} [\log P(\mathbf{x}, \mathbf{z}; \mathbf{w})]$$

## Transform parameters

$$\theta_{x|z} \leftarrow \frac{\exp(\mathbf{w}^T \mathbf{f}(x, z))}{\sum_{x'} \exp(\mathbf{w}^T \mathbf{f}(x', z))}$$

# Basic EM

---

Initialize probabilities  $\theta$

**repeat**

● Compute expected counts  $e$

● Fit parameters  $\theta$

**until** convergence

EM

# Basic EM

---

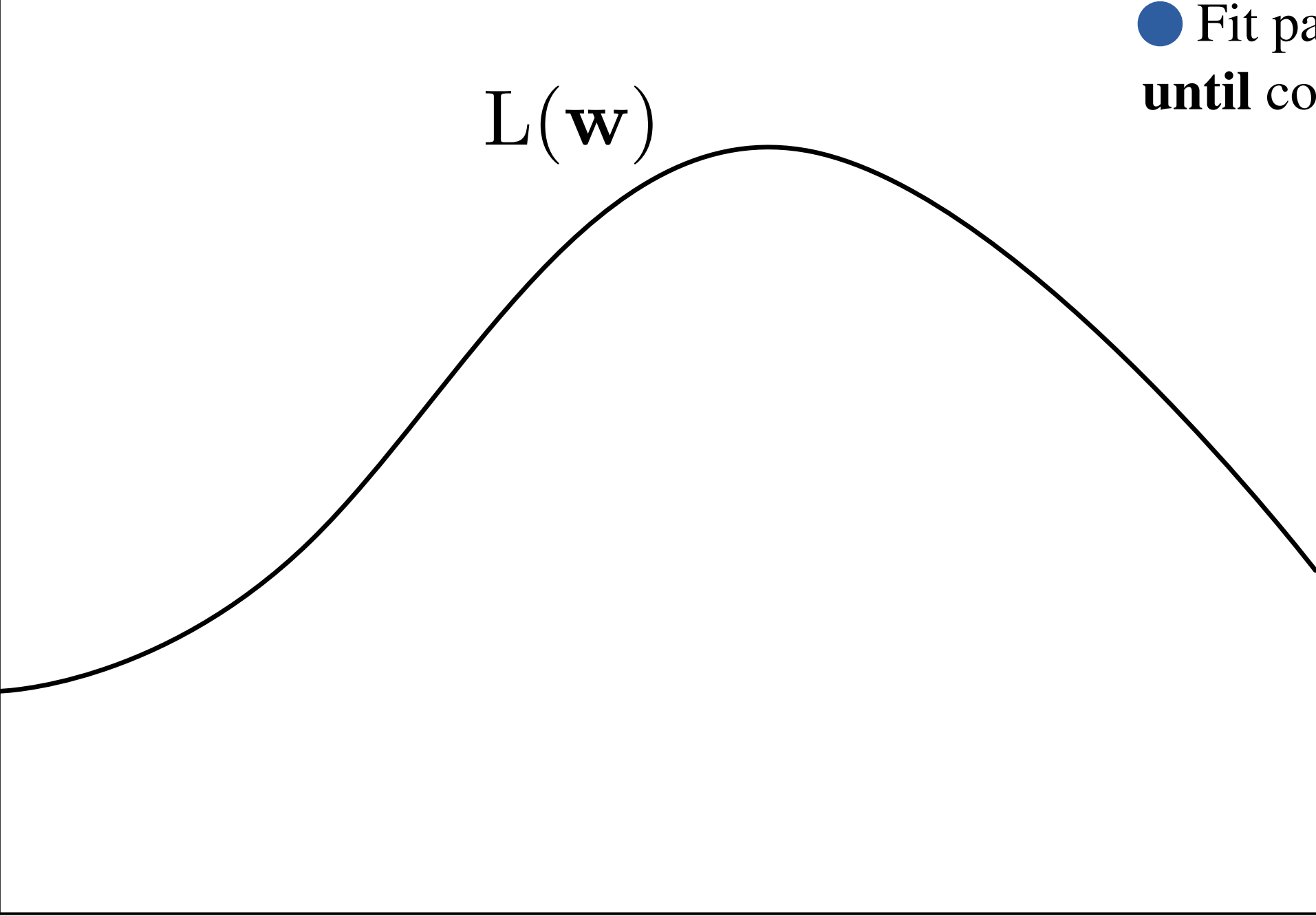
Initialize probabilities  $\theta$

**repeat**

● Compute expected counts  $e$

● Fit parameters  $\theta$

**until** convergence



$L(\mathbf{w})$

# Basic EM

---

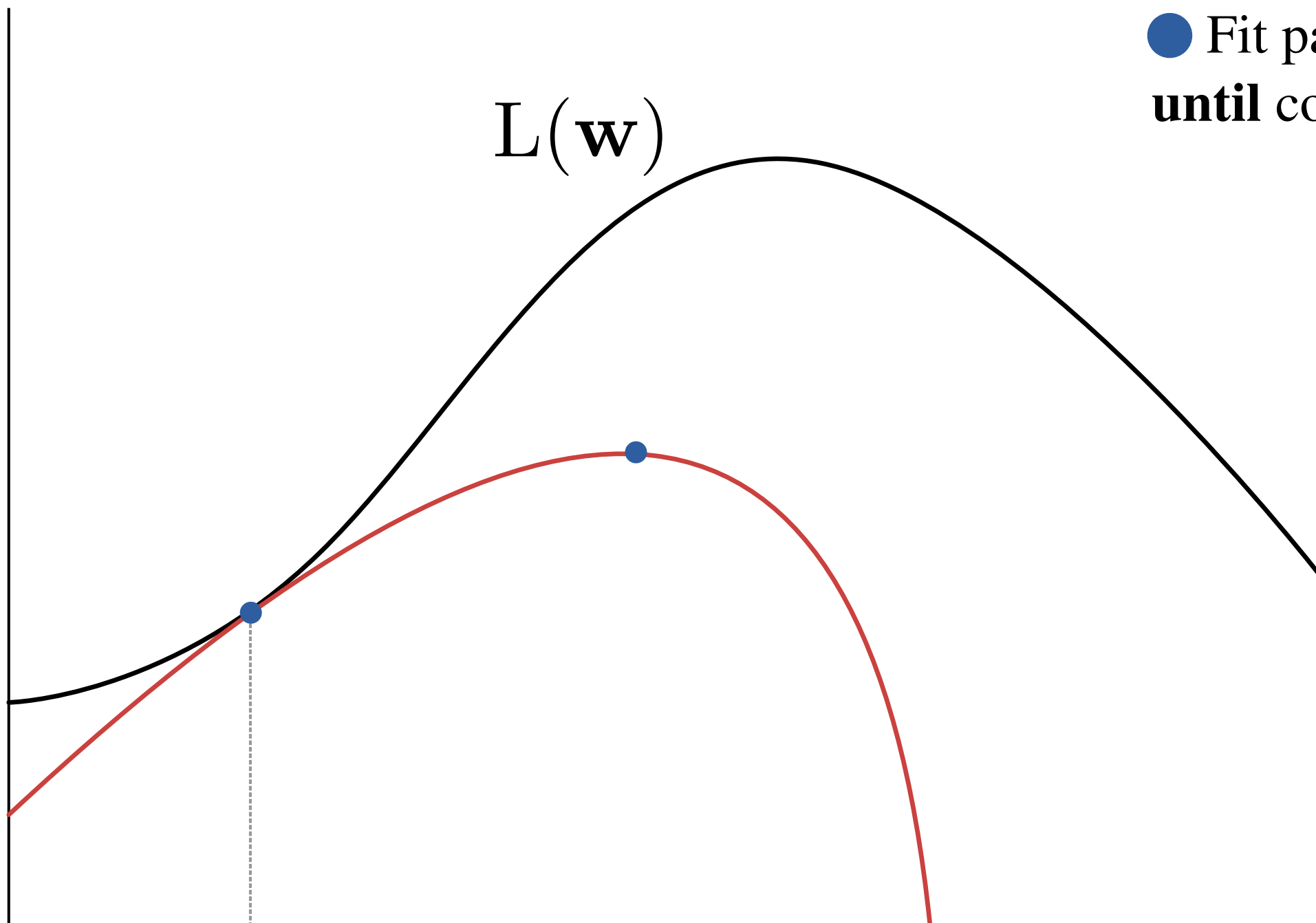
Initialize probabilities  $\theta$

**repeat**

● Compute expected counts  $e$

● Fit parameters  $\theta$

**until** convergence



# Basic EM

---

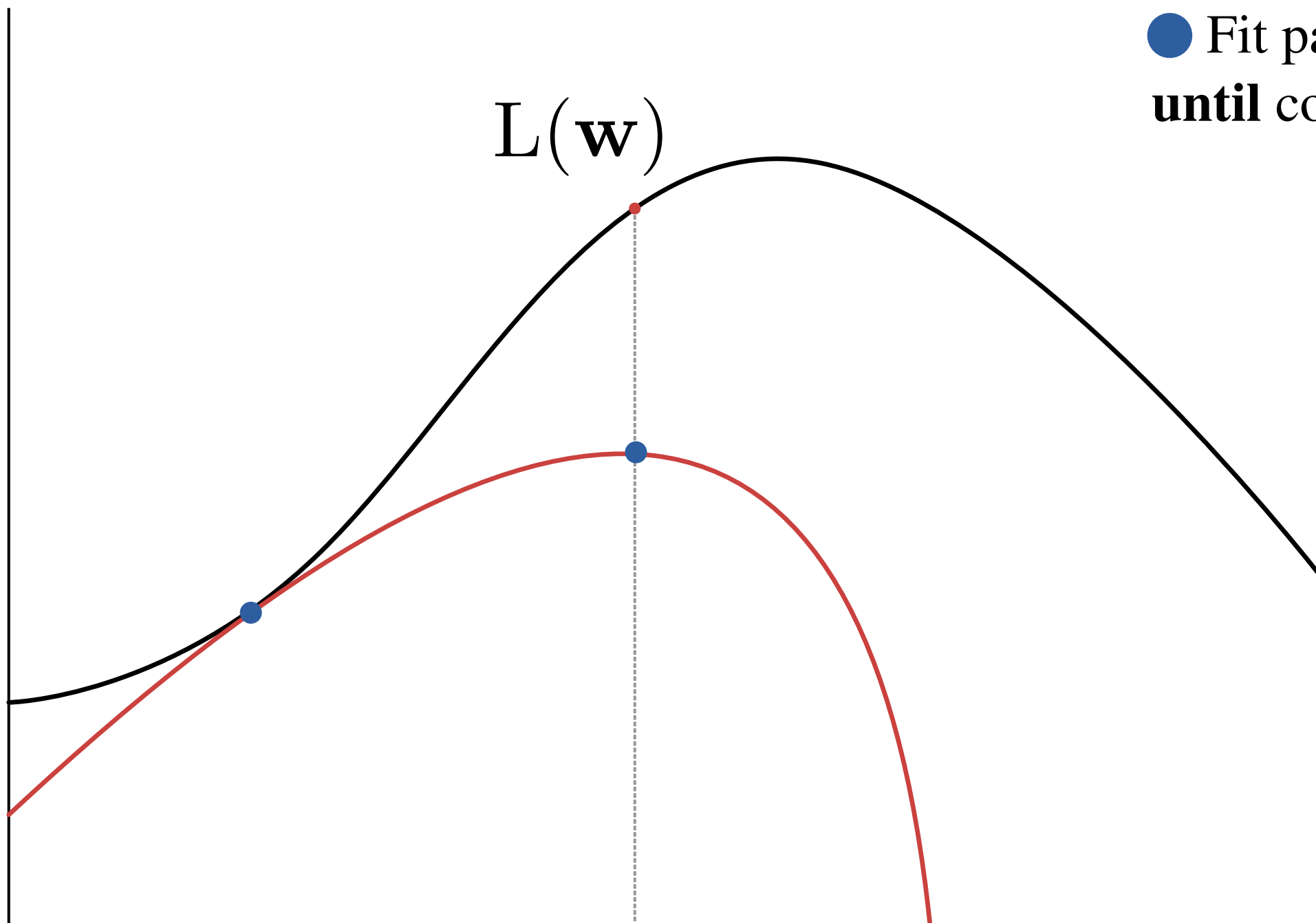
Initialize probabilities  $\theta$

**repeat**

● Compute expected counts  $e$

● Fit parameters  $\theta$

**until** convergence





# Basic EM

---

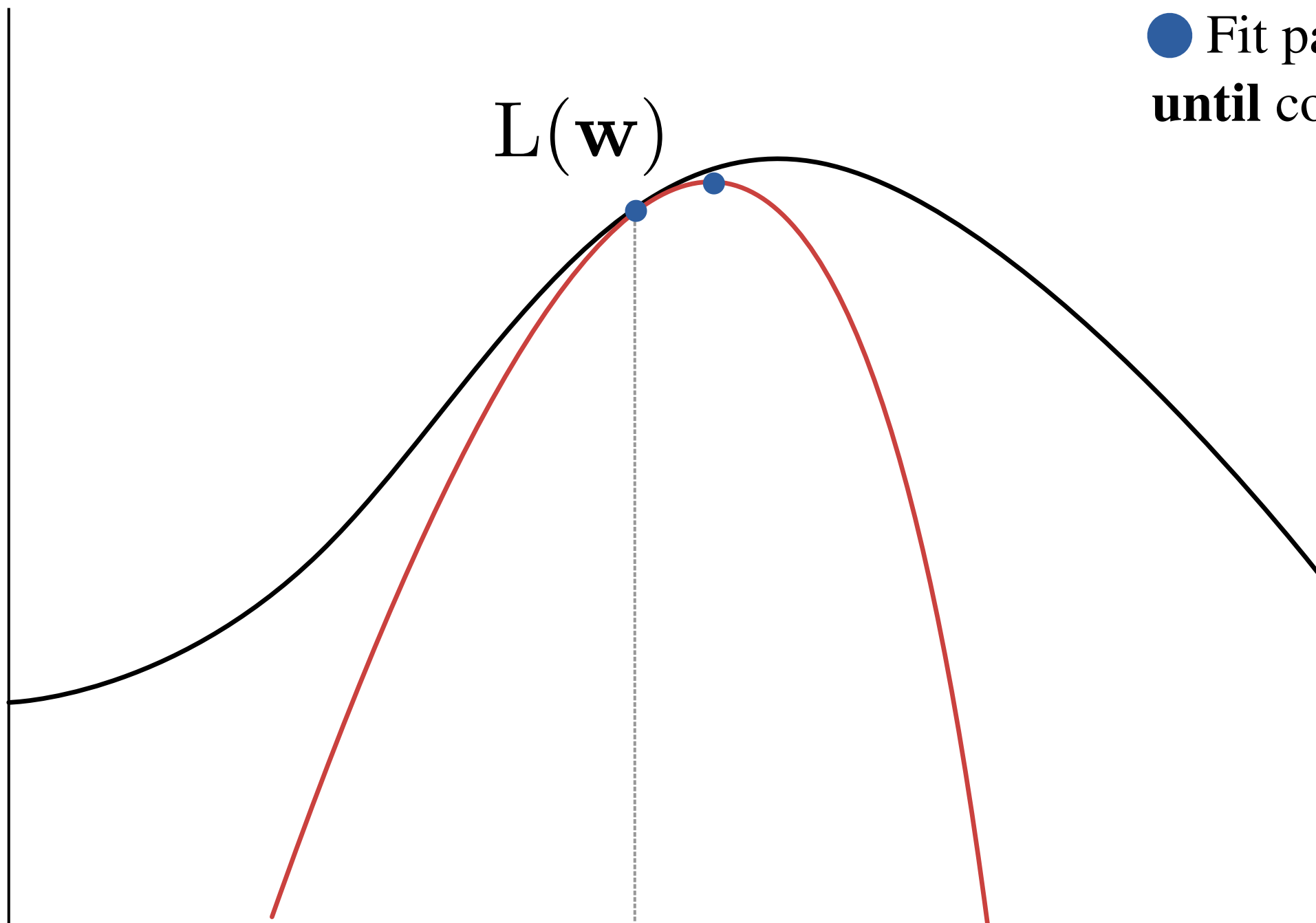
Initialize probabilities  $\theta$

**repeat**

● Compute expected counts  $e$

● Fit parameters  $\theta$

**until** convergence



# Basic EM

---

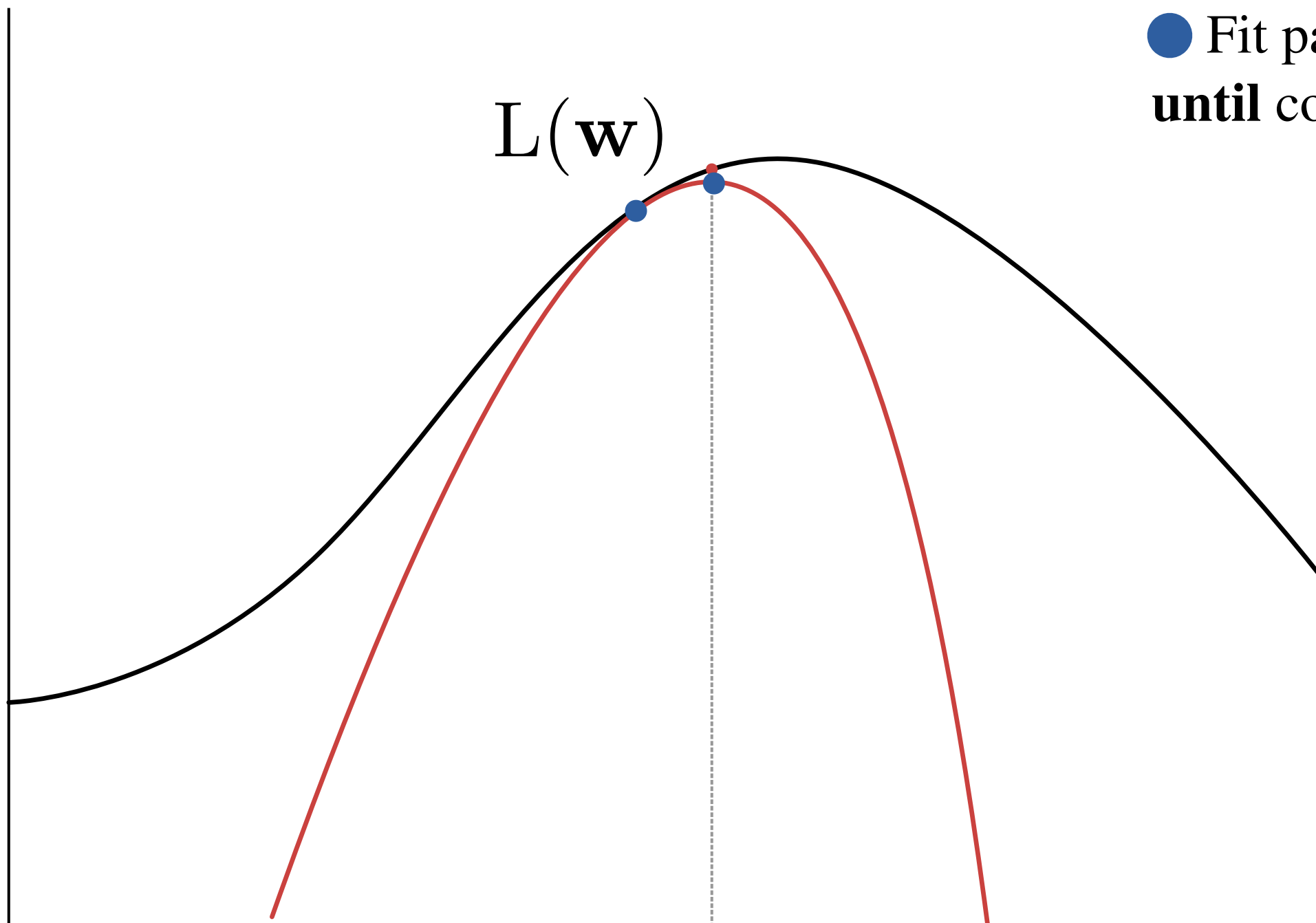
Initialize probabilities  $\theta$

**repeat**

● Compute expected counts  $e$

● Fit parameters  $\theta$

**until** convergence



# Basic EM

---

EM

Initialize probabilities  $\theta$

**repeat**

● Compute expected counts  $e$

● Fit parameters  $\theta$

**until** convergence

# EM with Features

---

Initialize weights  $w$

**repeat**

● Compute expected counts  $e$

● Fit parameters  $w$

● Transform  $w$  to  $\theta$

**until** convergence

$\Sigma$

# EM with Features

---

Initialize weights  $\mathbf{w}$

**repeat**

● Compute expected counts  $\mathbf{e}$

**repeat**

● Compute  $\ell(\mathbf{w}, \mathbf{e})$

● Compute  $\nabla \ell(\mathbf{w}, \mathbf{e})$

●  $\mathbf{w} \leftarrow \text{climb}(\mathbf{w}, \ell(\mathbf{w}, \mathbf{e}), \nabla \ell(\mathbf{w}, \mathbf{e}))$

**until** convergence

● Transform  $\mathbf{w}$  to  $\theta$

**until** convergence

EM

Fit Params

# EM with Features

Initialize weights  $\mathbf{w}$

**repeat**

● Compute expected counts  $\mathbf{e}$

**repeat**

■ Compute  $\ell(\mathbf{w}, \mathbf{e})$

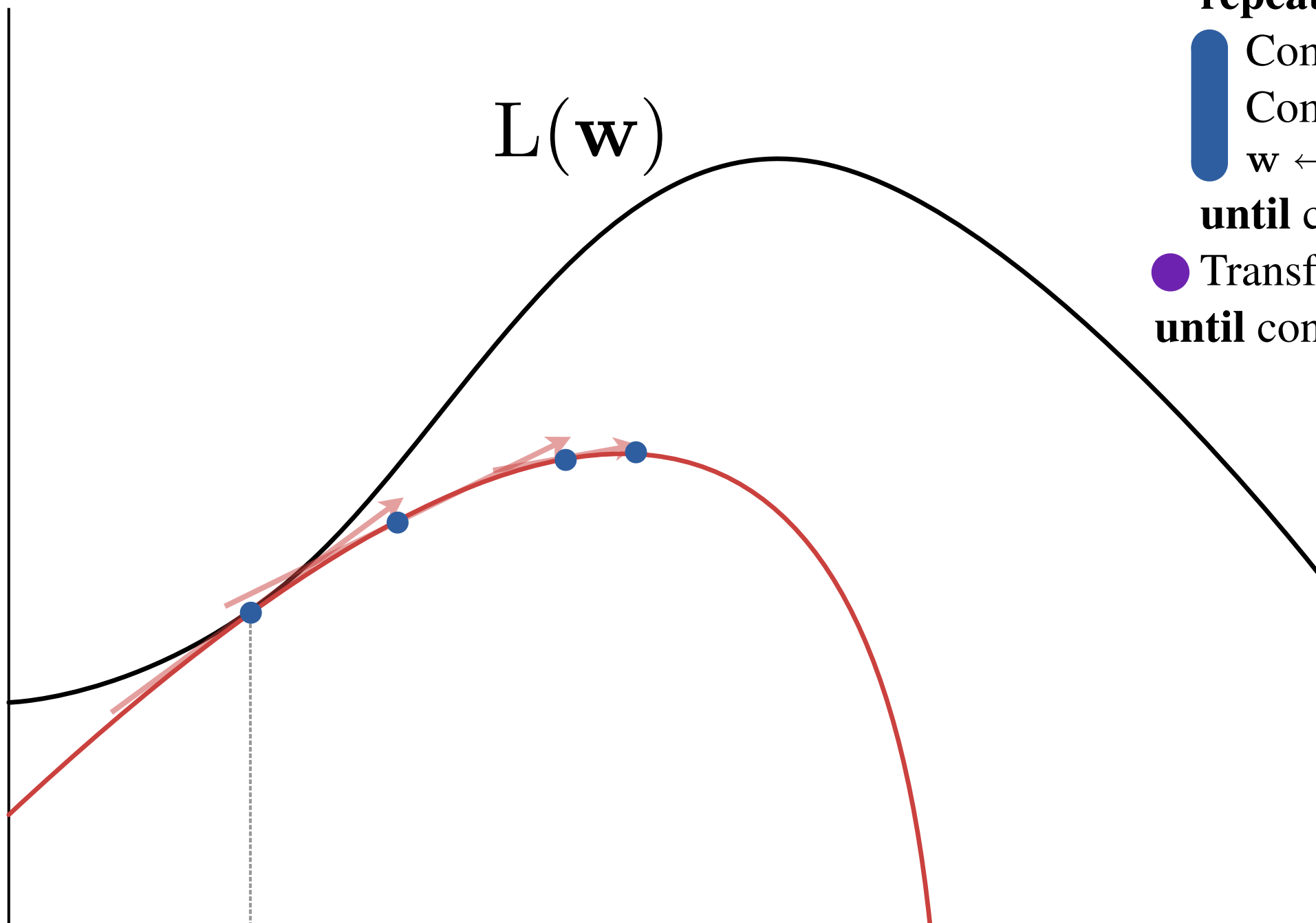
■ Compute  $\nabla \ell(\mathbf{w}, \mathbf{e})$

■  $\mathbf{w} \leftarrow \text{climb}(\mathbf{w}, \ell(\mathbf{w}, \mathbf{e}), \nabla \ell(\mathbf{w}, \mathbf{e}))$

**until** convergence

● Transform  $\mathbf{w}$  to  $\boldsymbol{\theta}$

**until** convergence



# EM with Features

Initialize weights  $\mathbf{w}$

**repeat**

● Compute expected counts  $\mathbf{e}$

**repeat**

■ Compute  $\ell(\mathbf{w}, \mathbf{e})$

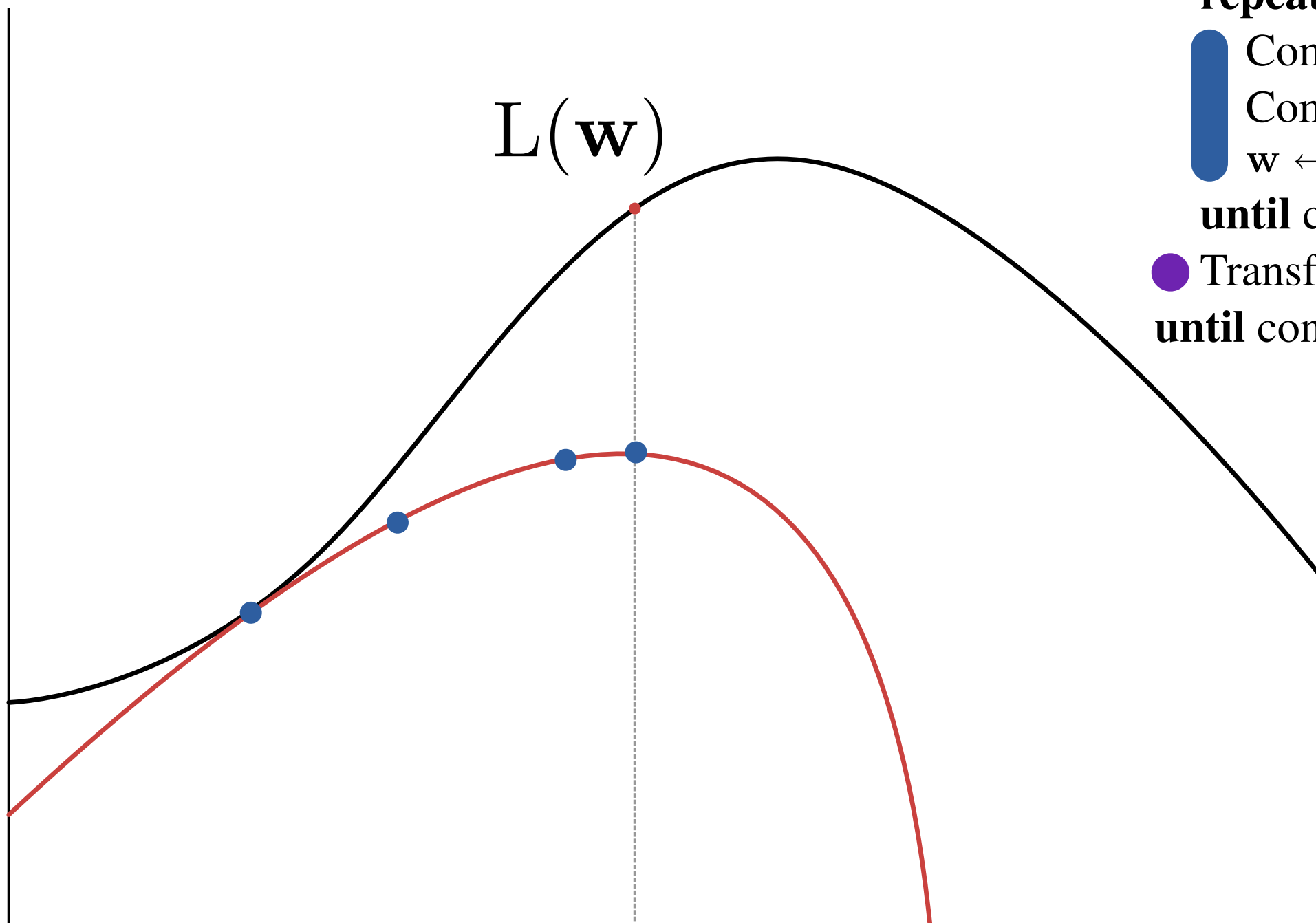
■ Compute  $\nabla \ell(\mathbf{w}, \mathbf{e})$

■  $\mathbf{w} \leftarrow \text{climb}(\mathbf{w}, \ell(\mathbf{w}, \mathbf{e}), \nabla \ell(\mathbf{w}, \mathbf{e}))$

**until** convergence

● Transform  $\mathbf{w}$  to  $\theta$

**until** convergence



# EM with Features

Initialize weights  $\mathbf{w}$

**repeat**

● Compute expected counts  $\mathbf{e}$

**repeat**

■ Compute  $\ell(\mathbf{w}, \mathbf{e})$

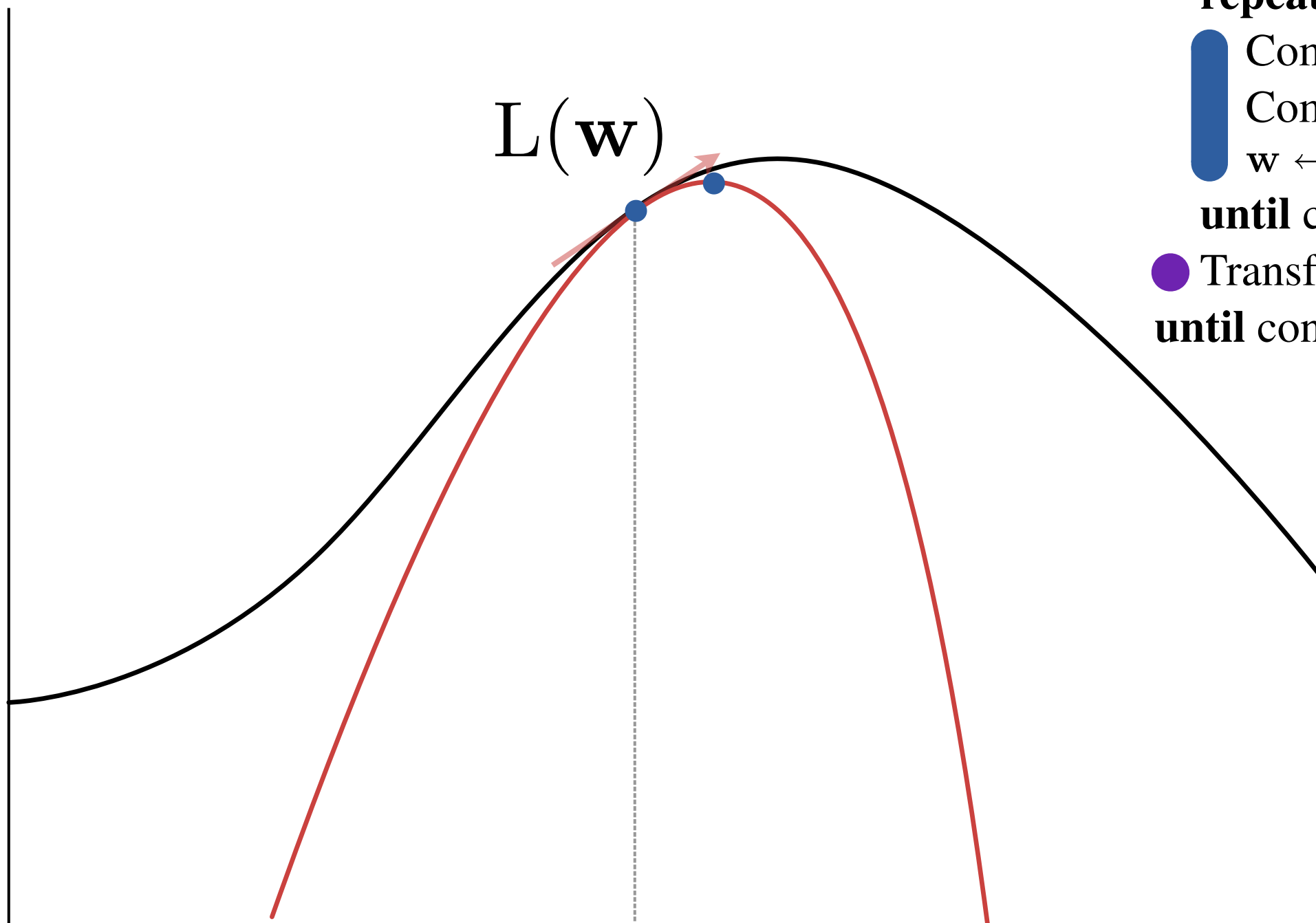
■ Compute  $\nabla \ell(\mathbf{w}, \mathbf{e})$

■  $\mathbf{w} \leftarrow \text{climb}(\mathbf{w}, \ell(\mathbf{w}, \mathbf{e}), \nabla \ell(\mathbf{w}, \mathbf{e}))$

**until convergence**

● Transform  $\mathbf{w}$  to  $\boldsymbol{\theta}$

**until convergence**





# EM with Features

Initialize weights  $\mathbf{w}$

**repeat**

● Compute expected counts  $\mathbf{e}$

**repeat**

■ Compute  $\ell(\mathbf{w}, \mathbf{e})$

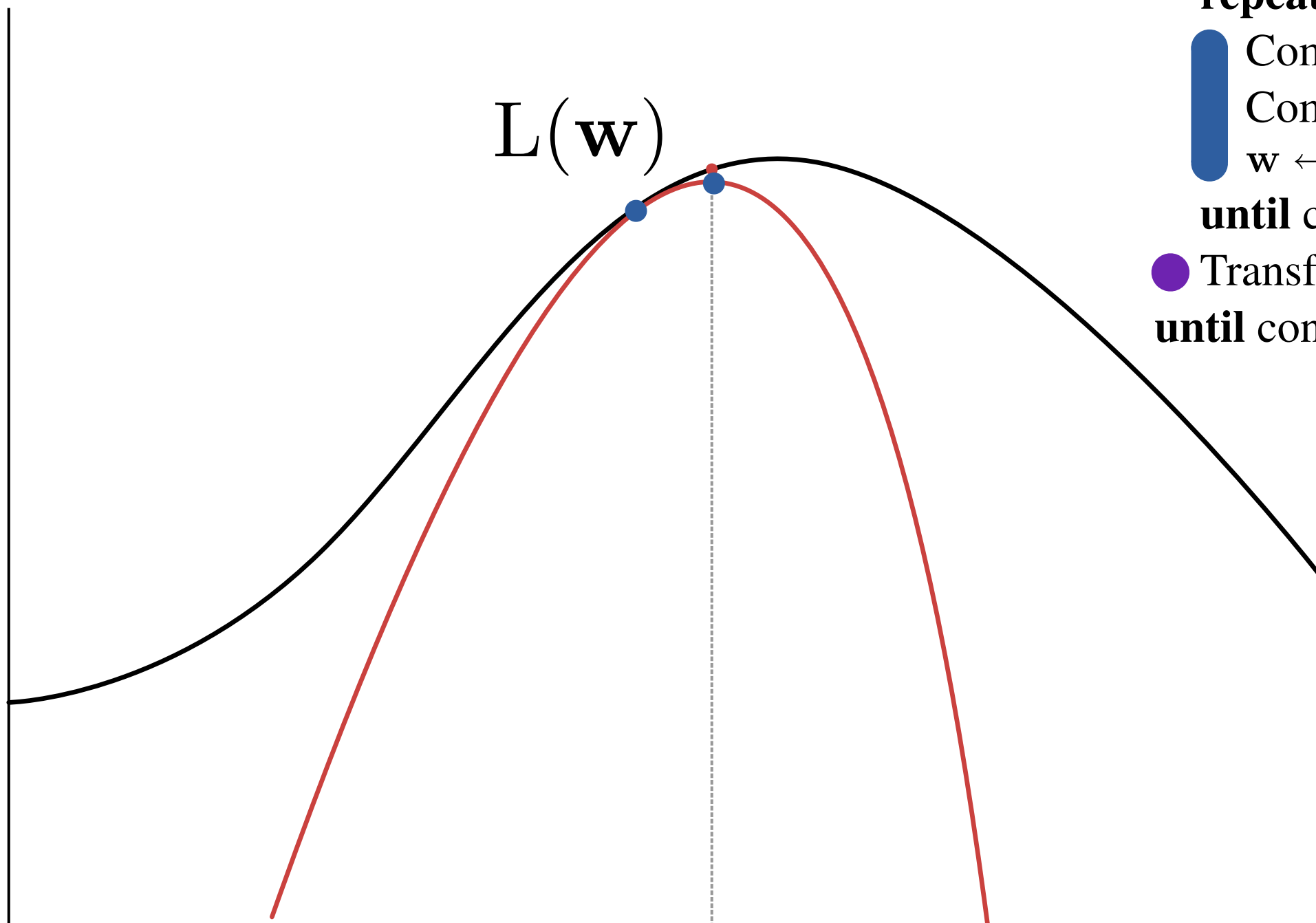
■ Compute  $\nabla \ell(\mathbf{w}, \mathbf{e})$

■  $\mathbf{w} \leftarrow \text{climb}(\mathbf{w}, \ell(\mathbf{w}, \mathbf{e}), \nabla \ell(\mathbf{w}, \mathbf{e}))$

**until convergence**

● Transform  $\mathbf{w}$  to  $\theta$

**until convergence**



# Direct Gradient with Features

---

## EM w/ Features

Initialize weights  $\mathbf{w}$

**repeat**

● Compute expected counts  $\mathbf{e}$

**repeat**

■ Compute  $\ell(\mathbf{w}, \mathbf{e})$

■ Compute  $\nabla \ell(\mathbf{w}, \mathbf{e})$

■  $\mathbf{w} \leftarrow \text{climb}(\mathbf{w}, \ell(\mathbf{w}, \mathbf{e}), \nabla \ell(\mathbf{w}, \mathbf{e}))$

**until** convergence

● Transform  $\mathbf{w}$  to  $\theta$

**until** convergence

## DG w/ Features

Initialize weights  $\mathbf{w}$

**repeat**

● Compute expected counts  $\mathbf{e}$

■ Compute  $L(\mathbf{w})$

■ Compute  $\nabla \ell(\mathbf{w}, \mathbf{e})$

■  $\mathbf{w} \leftarrow \text{climb}(\mathbf{w}, L(\mathbf{w}), \nabla \ell(\mathbf{w}, \mathbf{e}))$

● Transform  $\mathbf{w}$  to  $\theta$

**until** convergence

# Direct Gradient with Features

Initialize weights  $\mathbf{w}$

**repeat**

● Compute expected counts  $\mathbf{e}$

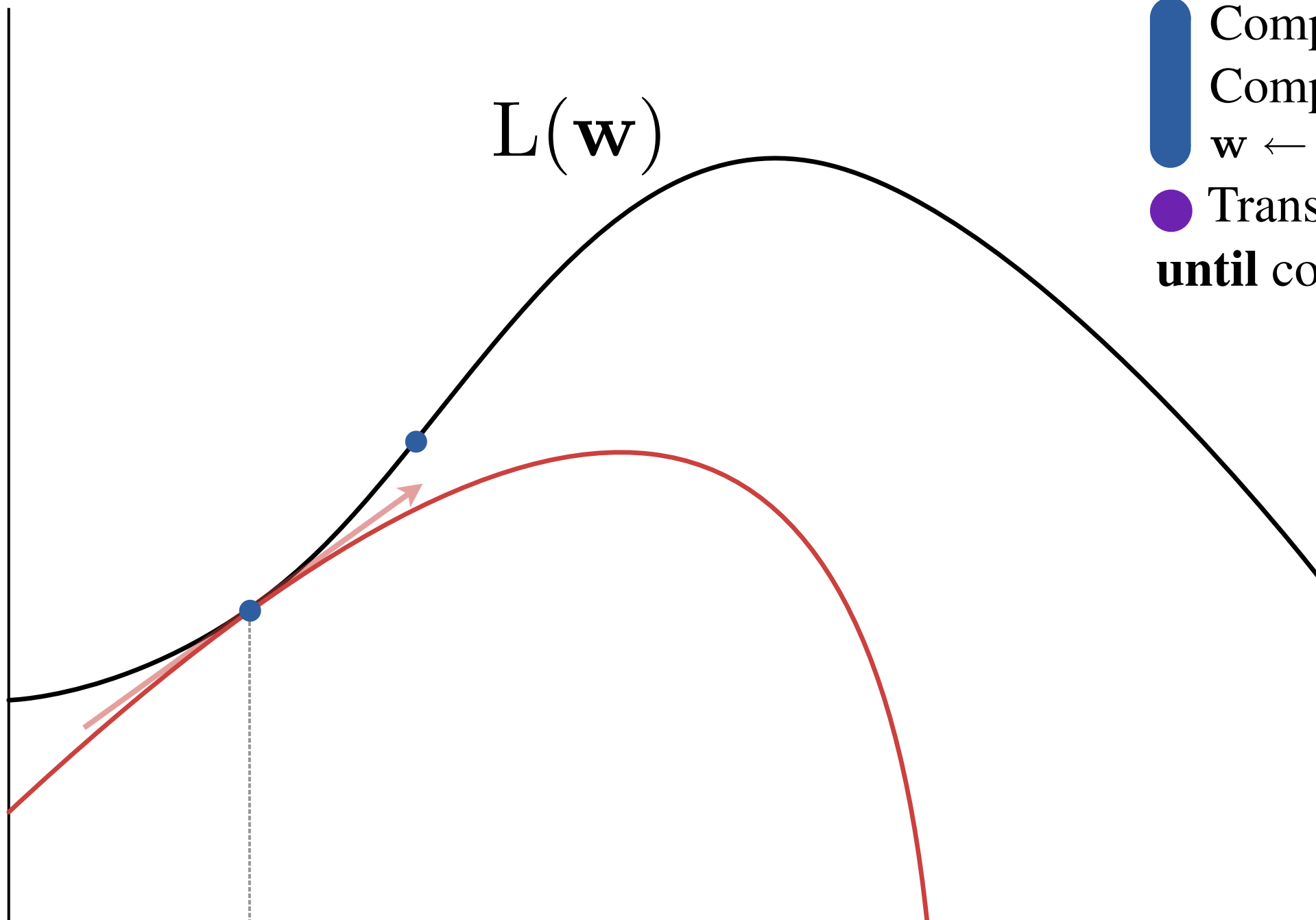
■ Compute  $L(\mathbf{w})$

■ Compute  $\nabla \ell(\mathbf{w}, \mathbf{e})$

■  $\mathbf{w} \leftarrow \text{climb}(\mathbf{w}, L(\mathbf{w}), \nabla \ell(\mathbf{w}, \mathbf{e}))$

● Transform  $\mathbf{w}$  to  $\boldsymbol{\theta}$

**until** convergence



# Direct Gradient with Features

---

Initialize weights  $\mathbf{w}$

**repeat**

● Compute expected counts  $\mathbf{e}$

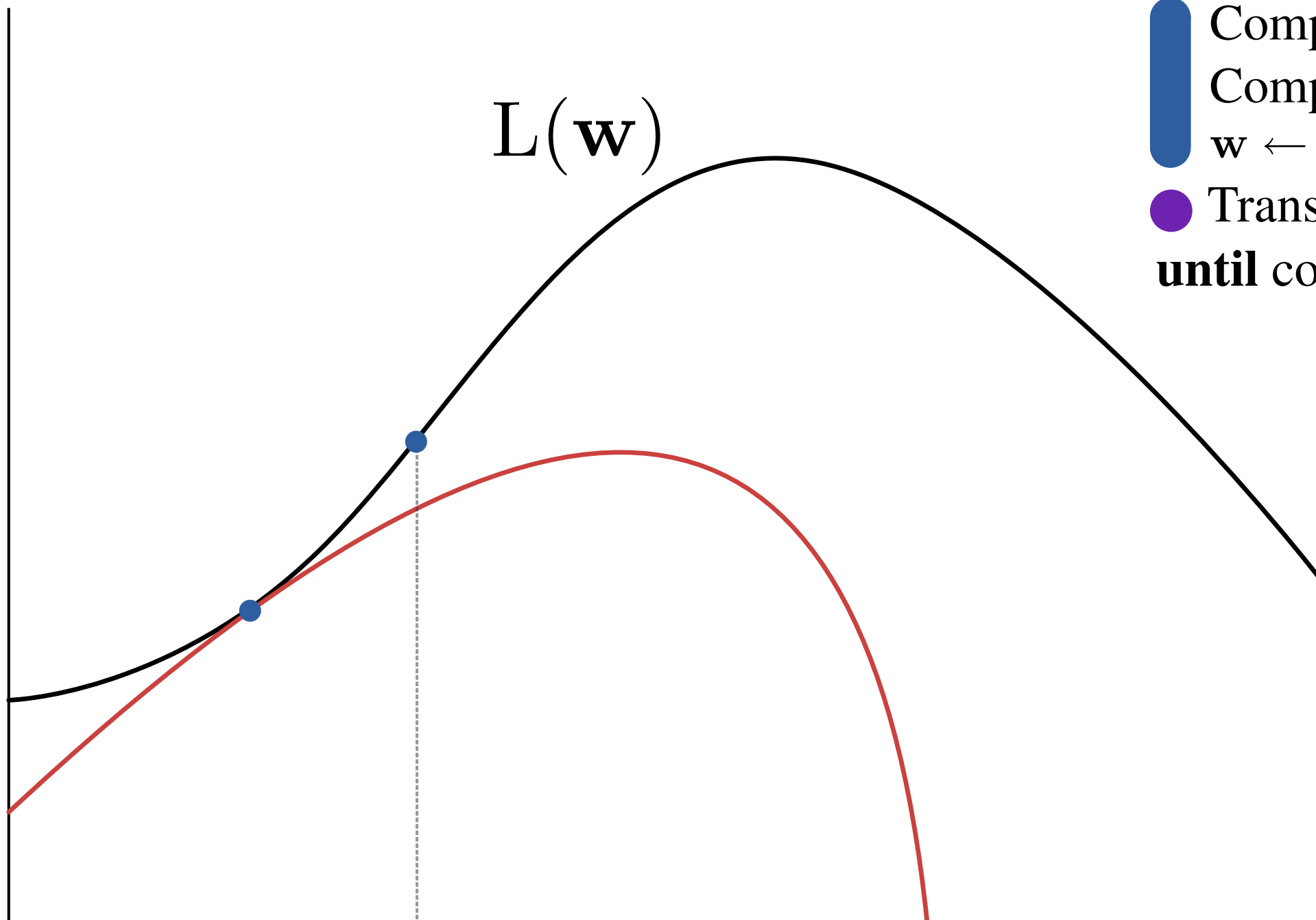
■ Compute  $L(\mathbf{w})$

■ Compute  $\nabla \ell(\mathbf{w}, \mathbf{e})$

■  $\mathbf{w} \leftarrow \text{climb}(\mathbf{w}, L(\mathbf{w}), \nabla \ell(\mathbf{w}, \mathbf{e}))$

● Transform  $\mathbf{w}$  to  $\boldsymbol{\theta}$

**until** convergence



# Direct Gradient with Features

Initialize weights  $\mathbf{w}$

**repeat**

● Compute expected counts  $\mathbf{e}$

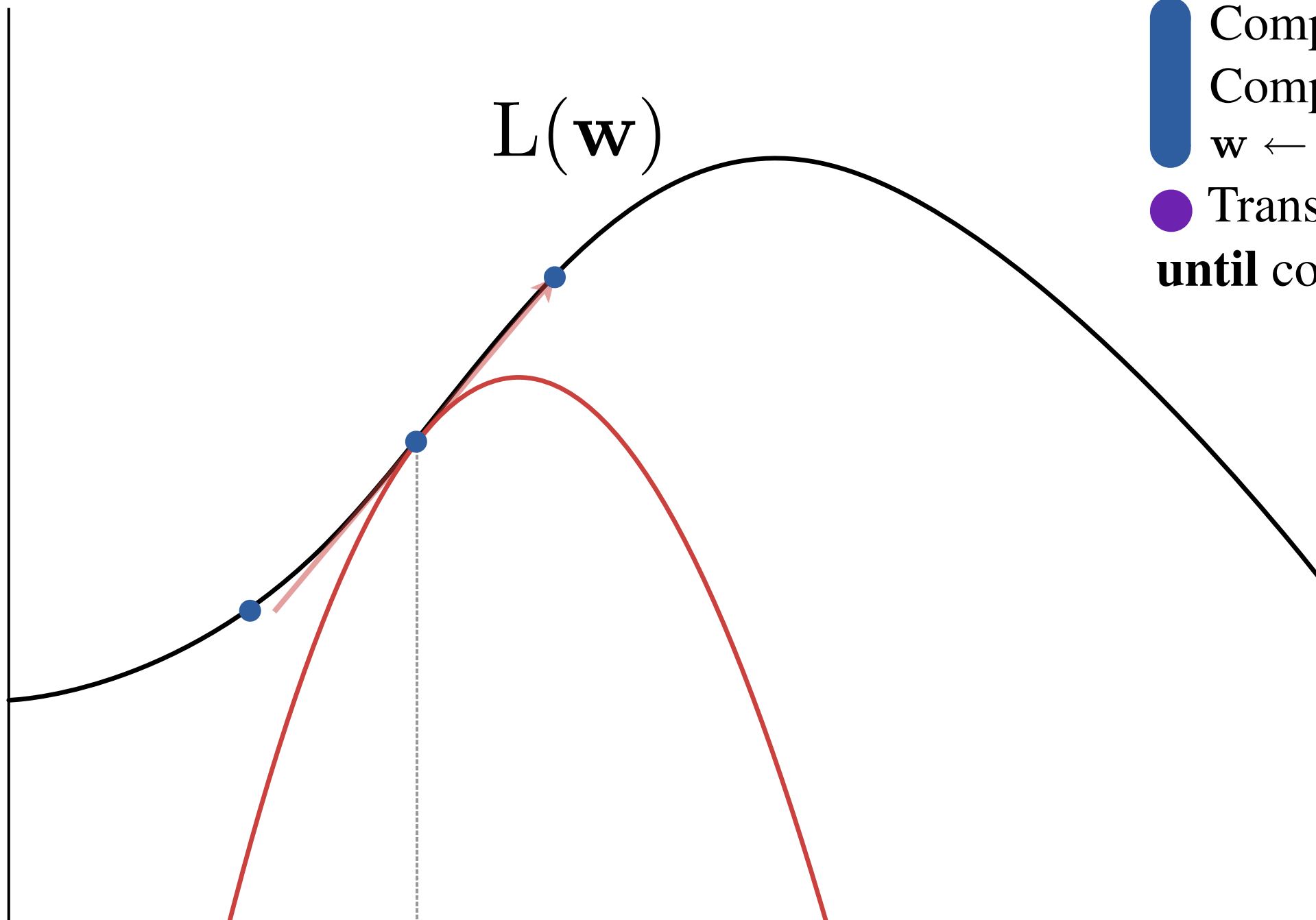
■ Compute  $L(\mathbf{w})$

■ Compute  $\nabla \ell(\mathbf{w}, \mathbf{e})$

■  $\mathbf{w} \leftarrow \text{climb}(\mathbf{w}, L(\mathbf{w}), \nabla \ell(\mathbf{w}, \mathbf{e}))$

● Transform  $\mathbf{w}$  to  $\boldsymbol{\theta}$

**until** convergence



# Direct Gradient with Features

Initialize weights  $\mathbf{w}$

**repeat**

● Compute expected counts  $\mathbf{e}$

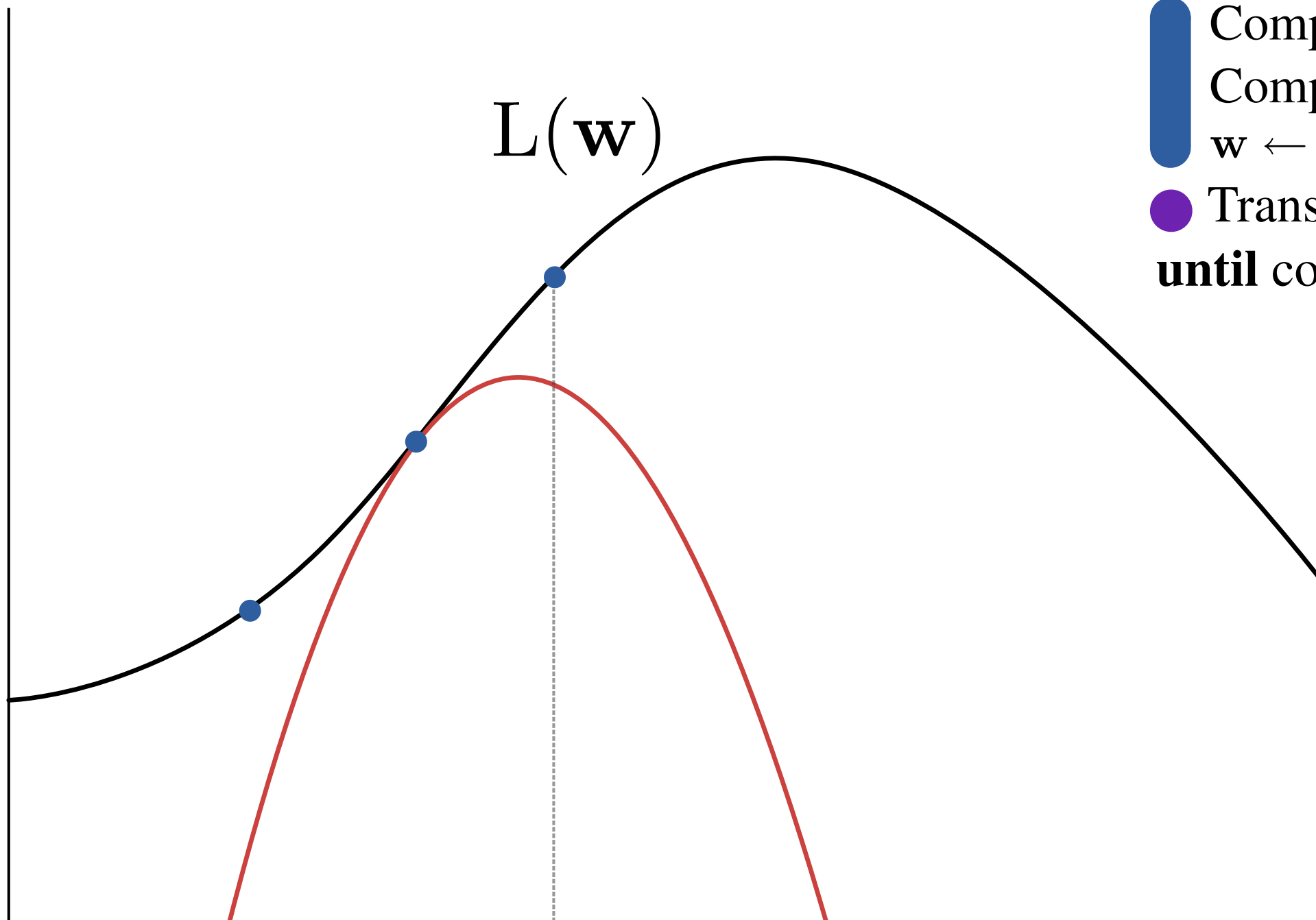
■ Compute  $L(\mathbf{w})$

■ Compute  $\nabla \ell(\mathbf{w}, \mathbf{e})$

■  $\mathbf{w} \leftarrow \text{climb}(\mathbf{w}, L(\mathbf{w}), \nabla \ell(\mathbf{w}, \mathbf{e}))$

● Transform  $\mathbf{w}$  to  $\boldsymbol{\theta}$

**until** convergence



# POS Induction Results

---

DT	JJ	NN	VBZ	IN	NN
The	green	cat	sleeps	at	home.

## Features:

- BASIC:  $\mathbb{1}(y = \cdot, z = \cdot)$
- CONTAINS-DIGIT: Check if  $y$  contains digit and conjoin with  $z$ :  
 $\mathbb{1}(\text{containsDigit}(y) = \cdot, z = \cdot)$
- CONTAINS-HYPHEN:  $\mathbb{1}(\text{containsHyphen}(x) = \cdot, z = \cdot)$
- INITIAL-CAP: Check if the first letter of  $y$  is capitalized:  $\mathbb{1}(\text{isCap}(y) = \cdot, z = \cdot)$
- N-GRAM: Indicator functions for character n-grams of up to length 3 present in  $y$ .

# POS Induction Results

DT	JJ	NN	VBZ	IN	NN
The	green	cat	sleeps	at	home.

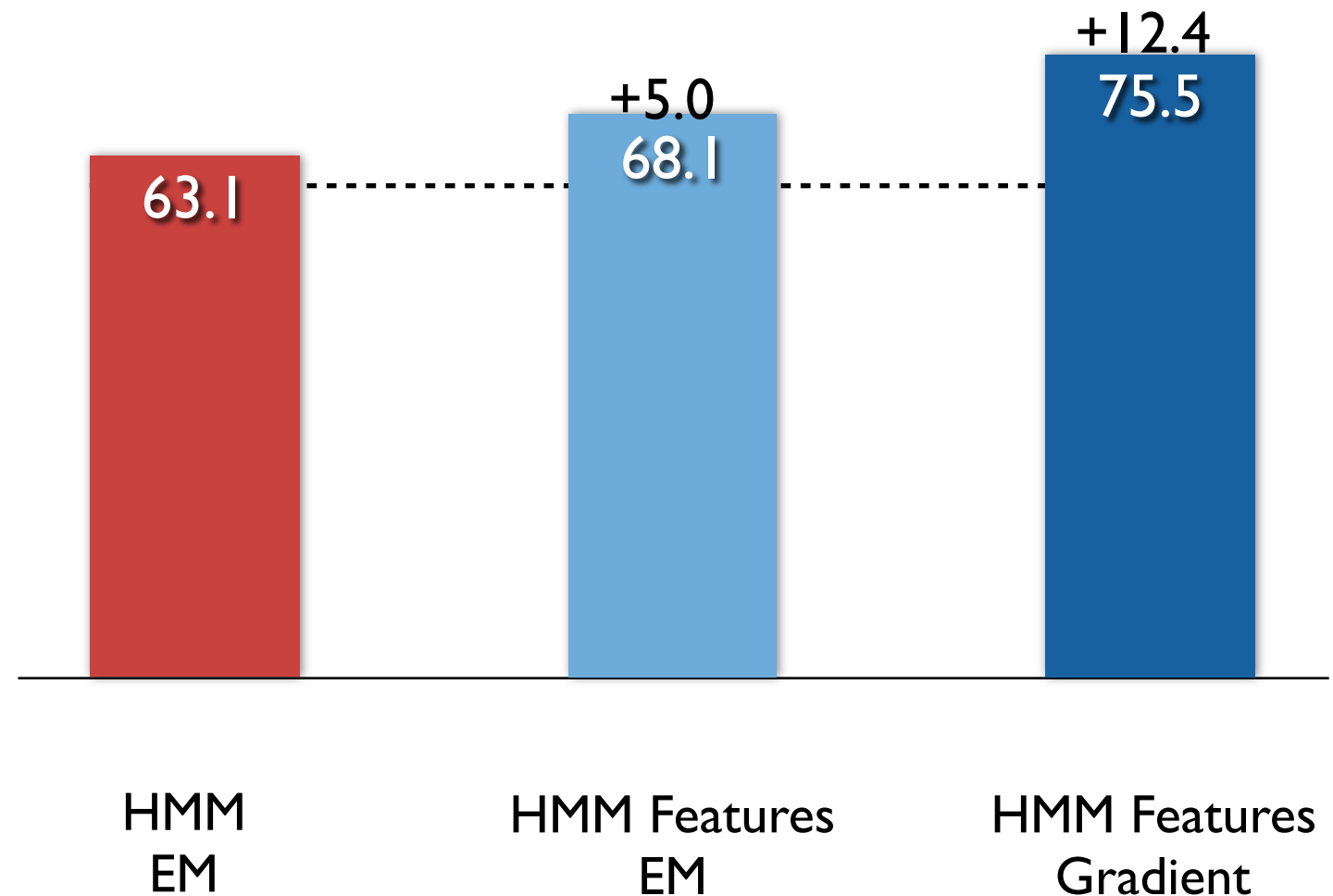
## Features:

BASIC:  $\mathbb{1}(y = \cdot, z = \cdot)$   
CONTAINS-DIGIT: Check if  $y$  contains digit and conjoin with  $z$ :  
 $\mathbb{1}(\text{containsDigit}(y) = \cdot, z = \cdot)$   
CONTAINS-HYPHEN:  $\mathbb{1}(\text{containsHyphen}(x) = \cdot, z = \cdot)$   
INITIAL-CAP: Check if the first letter of  $y$  is capitalized:  $\mathbb{1}(\text{isCap}(y) = \cdot, z = \cdot)$   
N-GRAM: Indicator functions for character n-grams of up to length 3 present in  $y$ .

## Data:

Train and test on entire WSJ  
No tagging dictionary  
45 POS tags

## Many-to-1 Accuracy





# POS Induction Results

DT	JJ	NN	VBZ	IN	NN
The	green	cat	sleeps	at	home.

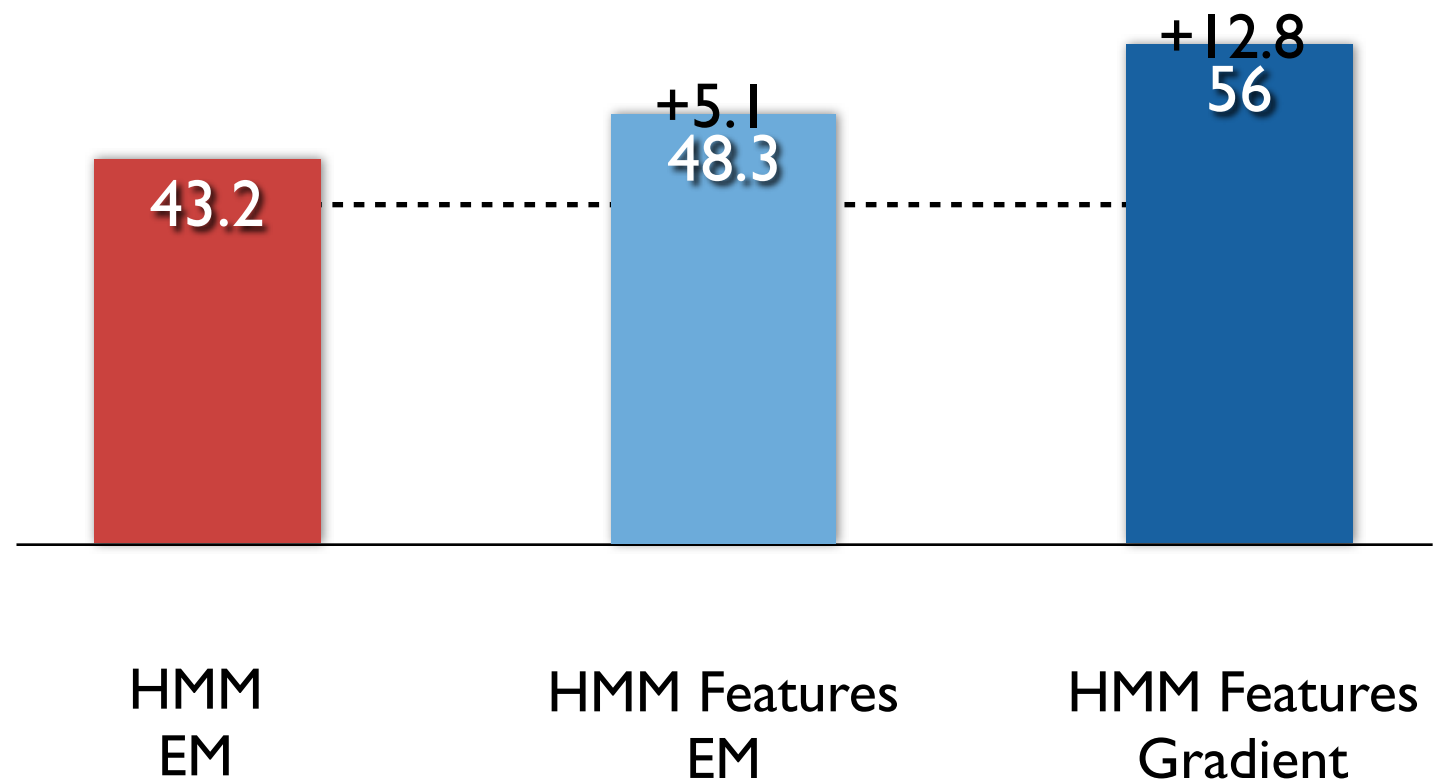
## Features:

BASIC:  $\mathbb{1}(y = \cdot, z = \cdot)$   
CONTAINS-DIGIT: Check if  $y$  contains digit and conjoin with  $z$ :  
 $\mathbb{1}(\text{containsDigit}(y) = \cdot, z = \cdot)$   
CONTAINS-HYPHEN:  $\mathbb{1}(\text{containsHyphen}(x) = \cdot, z = \cdot)$   
INITIAL-CAP: Check if the first letter of  $y$  is capitalized:  $\mathbb{1}(\text{isCap}(y) = \cdot, z = \cdot)$   
N-GRAM: Indicator functions for character n-grams of up to length 3 present in  $y$ .

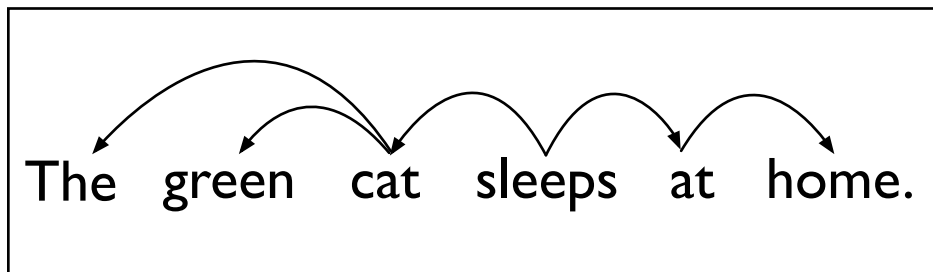
## Data:

Train and test on entire WSJ  
No tagging dictionary  
45 POS tags

## I-to-I Accuracy



# Grammar Induction Results



## Data:

Train WSJ10 Sec. 2-21  
CTB10 Sec. 1-270

Tune WSJ10 Sec. 22  
CTB10 Sec. 400-454

Test WSJ10 Sec. 23  
CTB10 Sec. 271-300

## Features:

BASIC:  $\mathbb{1}(a = \cdot, h = \cdot, \delta = \cdot)$

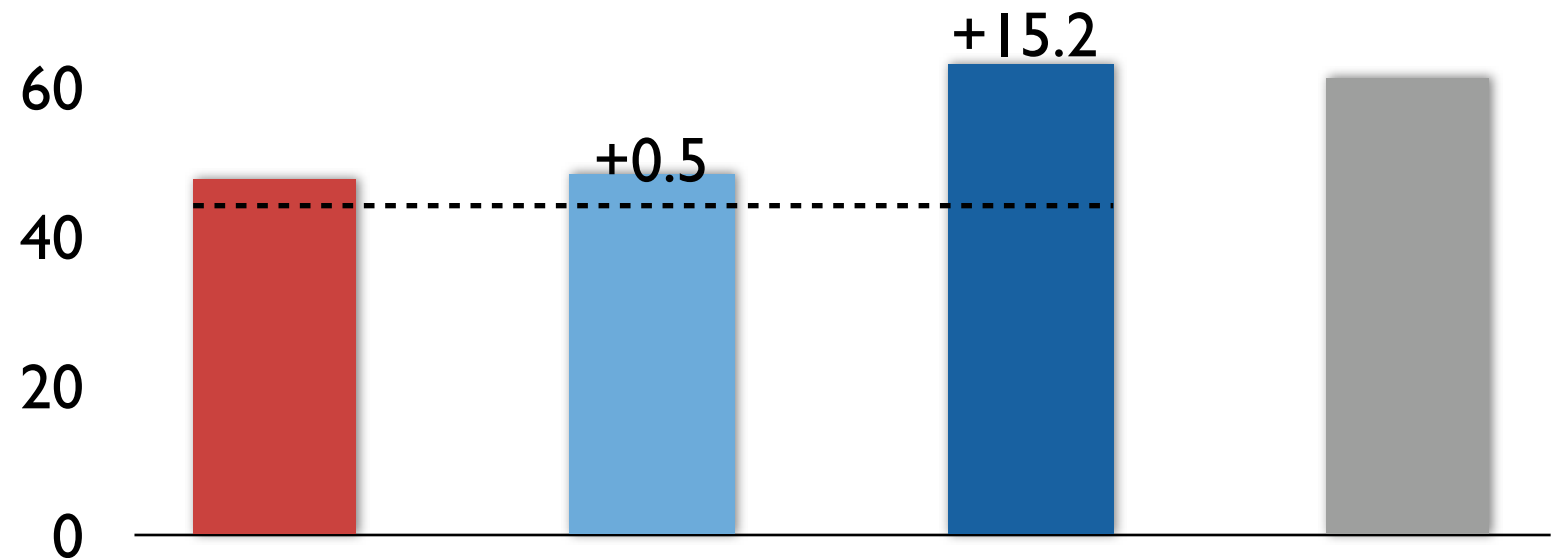
NOUN: Generalize the morphological variants of nouns by using  $\text{isNoun}(\cdot)$ :  
 $\mathbb{1}(a = \cdot, \text{isNoun}(h) = \cdot, \delta = \cdot)$   
 $\mathbb{1}(\text{isNoun}(a) = \cdot, h = \cdot, \delta = \cdot)$   
 $\mathbb{1}(\text{isNoun}(a) = \cdot, \text{isNoun}(h) = \cdot, \delta = \cdot)$

VERB: Same as above, generalizing verbs instead of nouns by using  $\text{isVerb}(\cdot)$

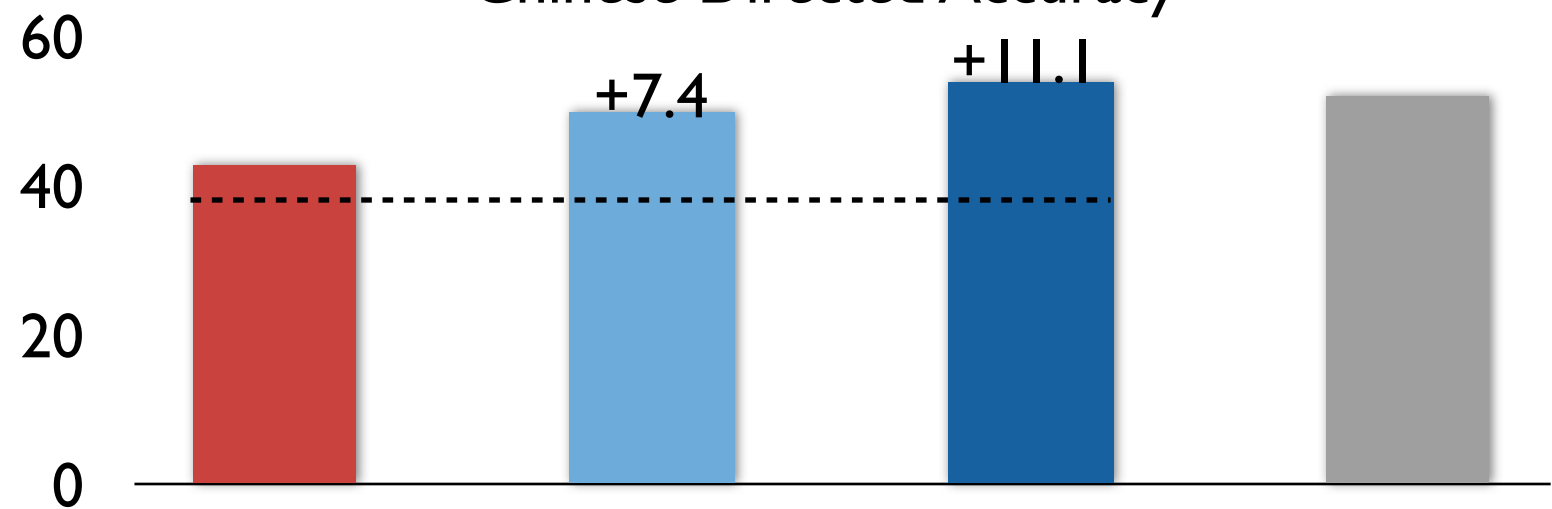
NOUN-VERB: Same as above, generalizing with  $\text{isVerbOrNoun}(\cdot) = \text{isVerb}(\cdot) \vee \text{isNoun}(\cdot)$

BACK-OFF: We add versions of all other features that ignore direction or adjacency.

English Directed Accuracy



Chinese Directed Accuracy



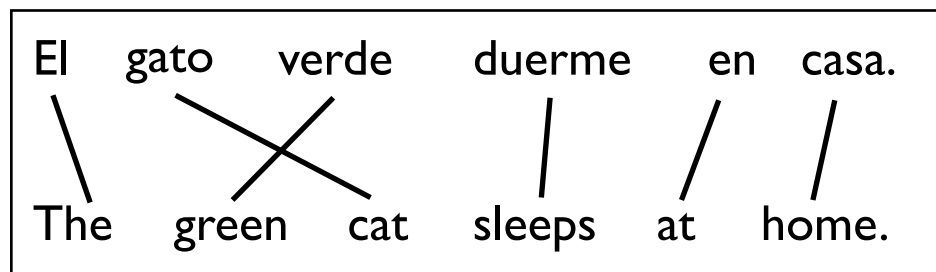
DMV  
EM

DMV Features  
EM

DMV Features  
LBFGS

Cohen and  
Smith '09  
SLN DMV

# Word Alignment Results



## Data:

Train 10K sentences of FIBIS  
Ch-En newswire

Test NIST 2002 Ch-En dev set

## Features:

BASIC:  $\mathbb{1}(e = \cdot, y = \cdot)$   
EDIT-DISTANCE:  $\mathbb{1}(\text{dist}(y, e) = \cdot)$   
DICTIONARY:  $\mathbb{1}((y, e) \in D)$  for dictionary  $D$ .  
STEM:  $\mathbb{1}(\text{stem}(e) = \cdot, y = \cdot)$  for Porter stemmer.  
PREFIX:  $\mathbb{1}(\text{prefix}(e) = \cdot, y = \cdot)$  for prefixes of length 4.  
CHARACTER:  $\mathbb{1}(e = \cdot, \text{charAt}(y, i) = \cdot)$  for index  $i$  in the Chinese word.

50

AER

40

-2.4

30

-3.8

20

10

0

Model I  
EM

Model I Features  
EM

HMM  
EM

HMM Features  
EM

# Word Segmentation Results

[T h e][g r e e n][c a t]

## Data:

Train and test on phonetic version  
of Bernstein-Ratner corpus

## Features:

BASIC:  $\mathbb{1}(z = \cdot)$   
LENGTH:  $\mathbb{1}(\text{length}(z) = \cdot)$   
NUMBER-VOWELS:  $\mathbb{1}(\text{numVowels}(z) = \cdot)$   
PHONO-CLASS-PREF:  $\mathbb{1}(\text{prefix}(\text{coarsePhonemes}(z)) = \cdot)$   
PHONO-CLASS-PREF:  $\mathbb{1}(\text{suffix}(\text{coarsePhonemes}(z)) = \cdot)$

