

# Painless Unsupervised Learning with Features

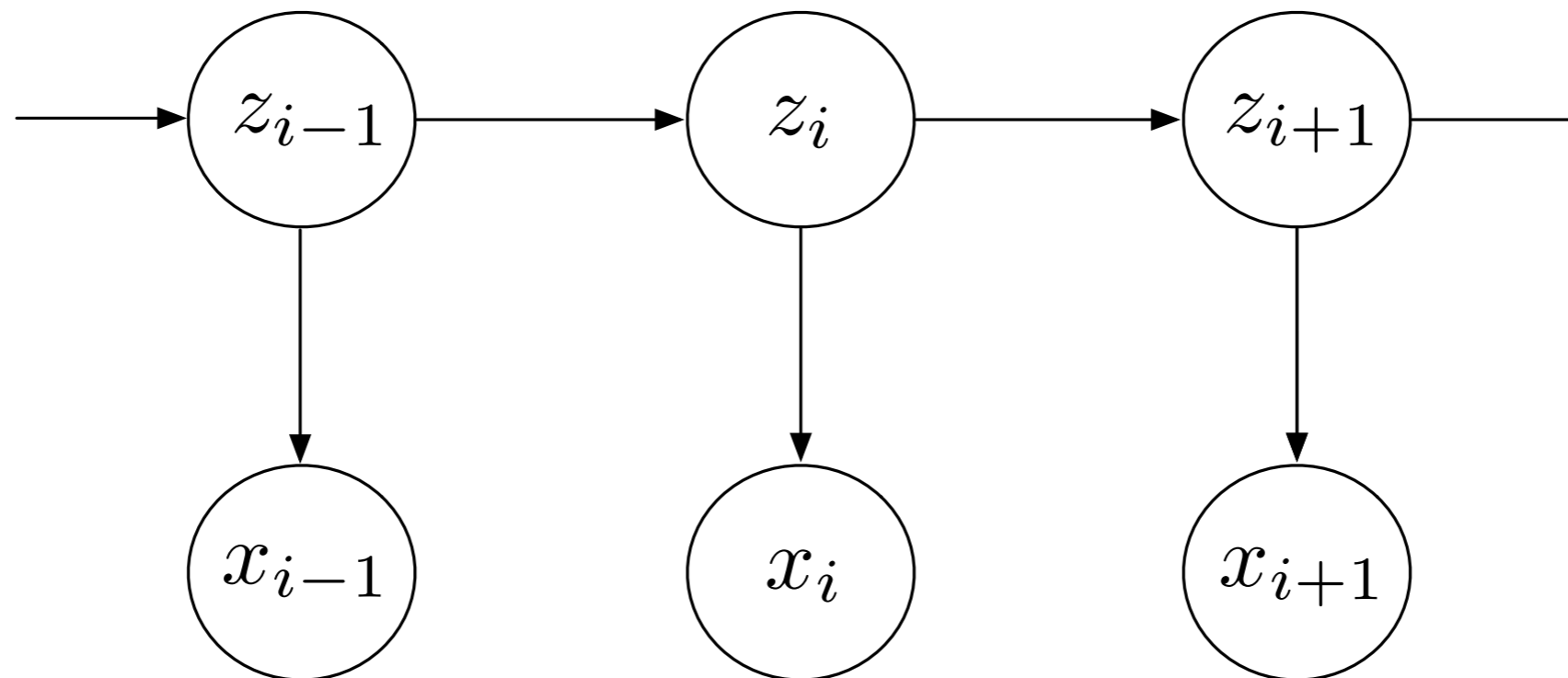


Taylor Berg-Kirkpatrick    Alexandre Bouchard-Côté  
John DeNero    Dan Klein



# Basic HMM for POS Induction

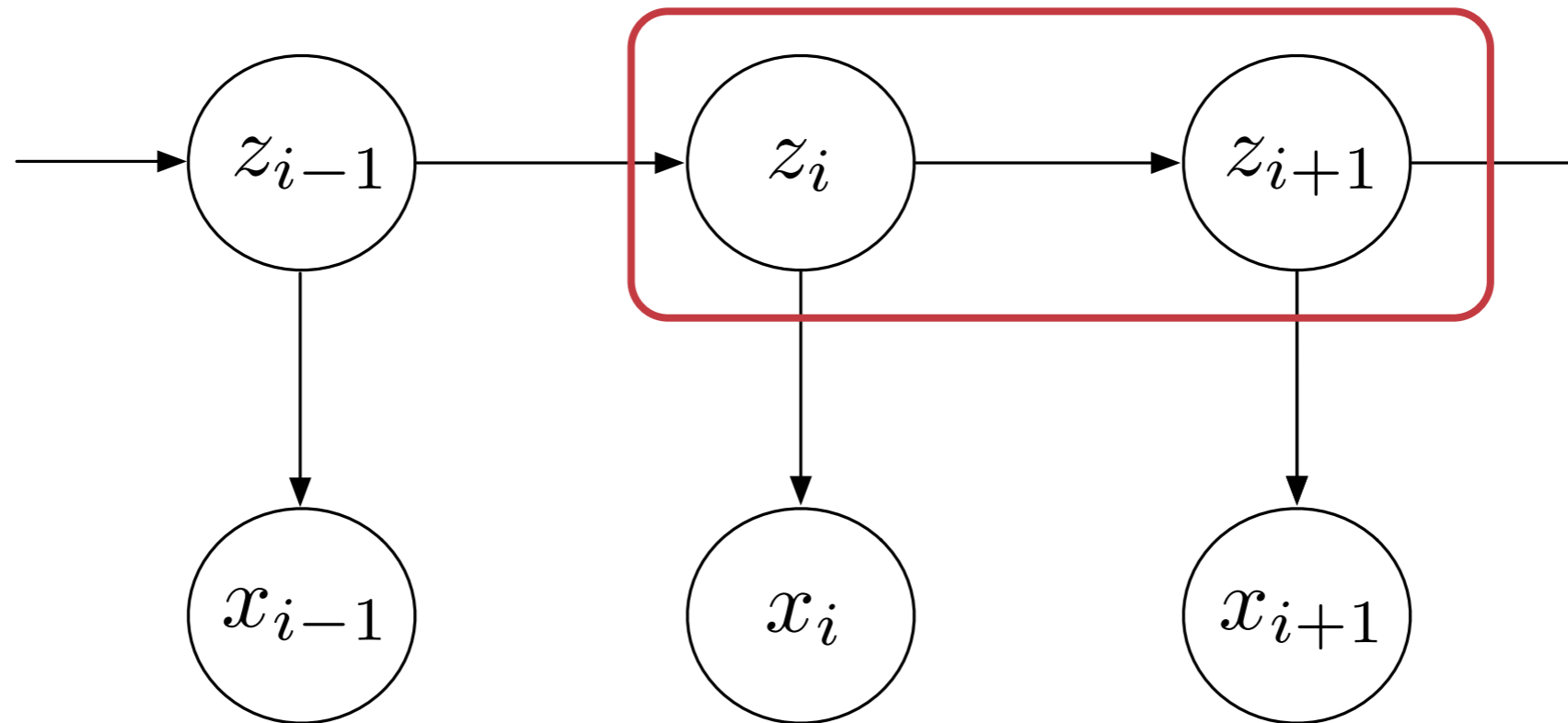
---



# Basic HMM for POS Induction

Transition distribution:

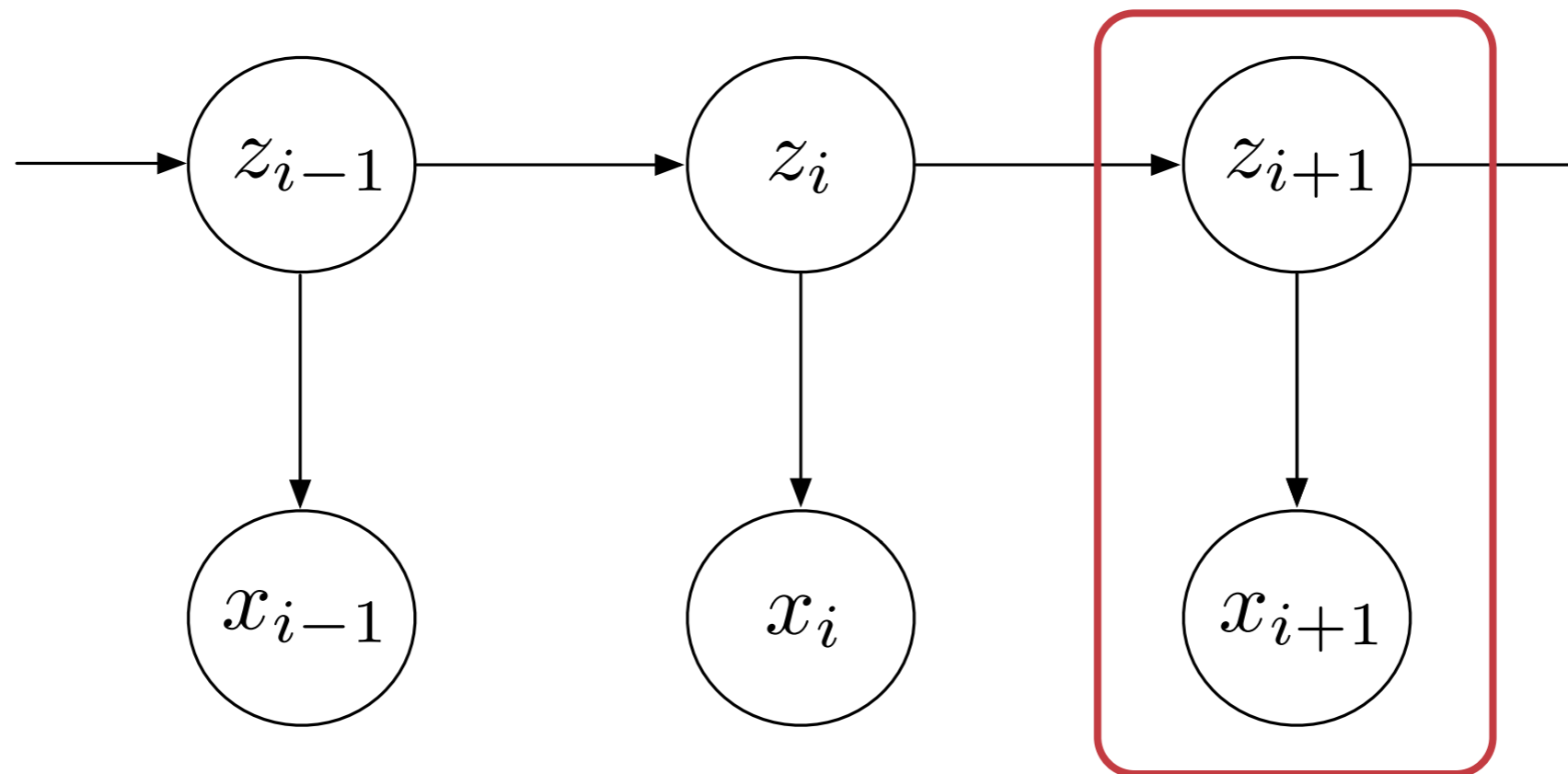
$$P(z' | z)$$



# Basic HMM for POS Induction

Emission distribution:

$$P(x|z)$$





# Parameterization

---

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

$x$

---

John

Mary

running

jumping



# Parameterization

---

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

$\theta_{x \text{NNP}}$	$x$
0.1	John
0.0	Mary
0.2	running
0.0	jumping



# Parameterization

---

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

$\theta_{x \text{NNP}}$	$x$	$\mathbf{f}$
0.1	John	+Cap
0.0	Mary	+Cap
0.2	running	+ing
0.0	jumping	+ing



# Parameterization

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

**W:**

+Cap	+1.2
+ing	-0.3

$\theta_{x \text{NNP}}$	$x$	$\mathbf{f}$	$e^{\mathbf{w}^T \mathbf{f}}$
0.1	John	+Cap	0.3
0.0	Mary	+Cap	0.3
0.2	running	+ing	0.1
0.0	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**

# Unsupervised Learning with Features

---

Main idea: local multinomials become maxents

**EM** + **Maxent M-Step** =  
Unsupervised learning w/ features

Berkeley



# POS Induction Accuracy

---



# POS Induction Accuracy

---

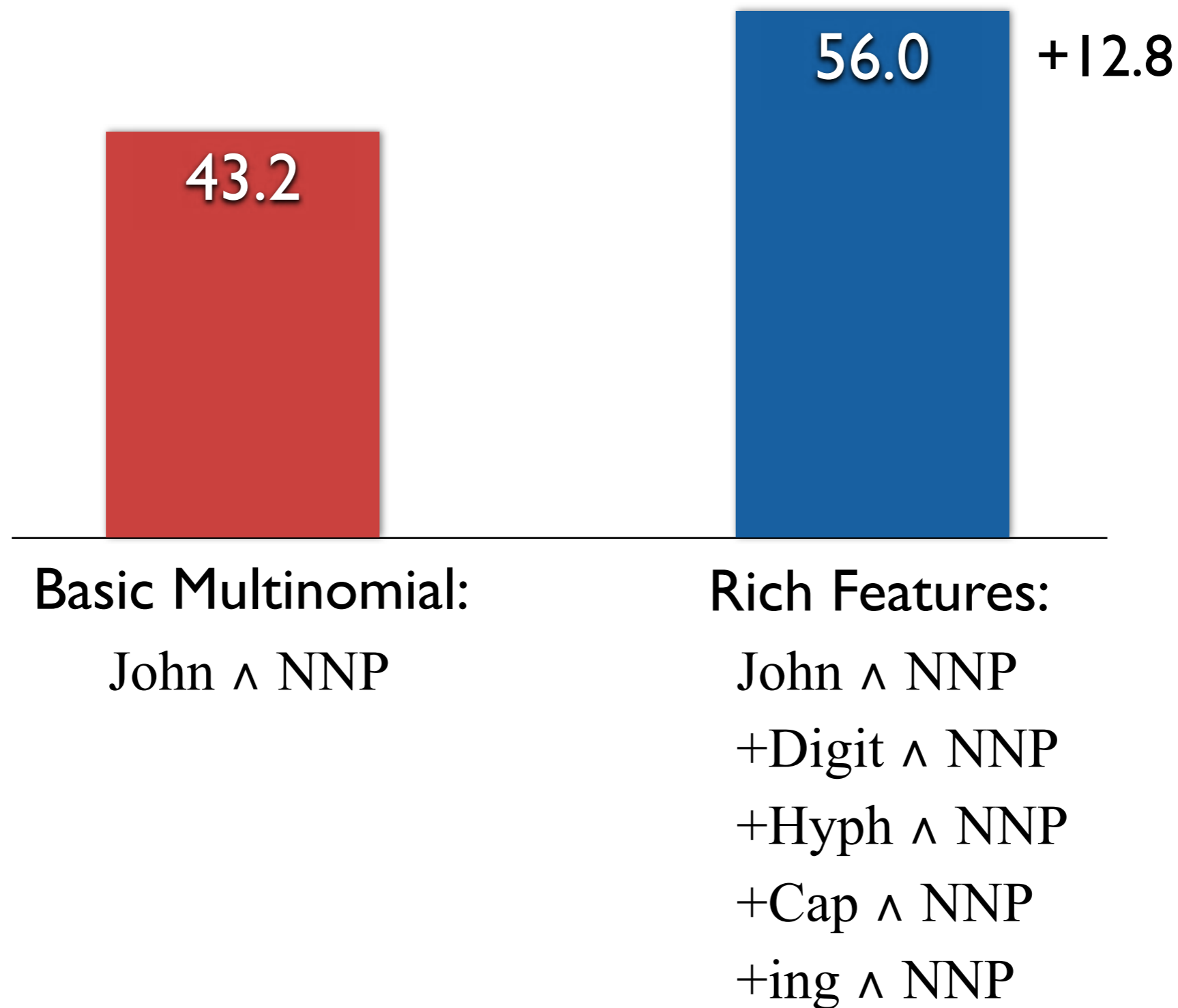
A bar chart with a single red bar representing the accuracy of a Basic Multinomial model. The value '43.2' is printed in white text inside the bar.

Model	Accuracy (%)
Basic Multinomial	43.2

Basic Multinomial:

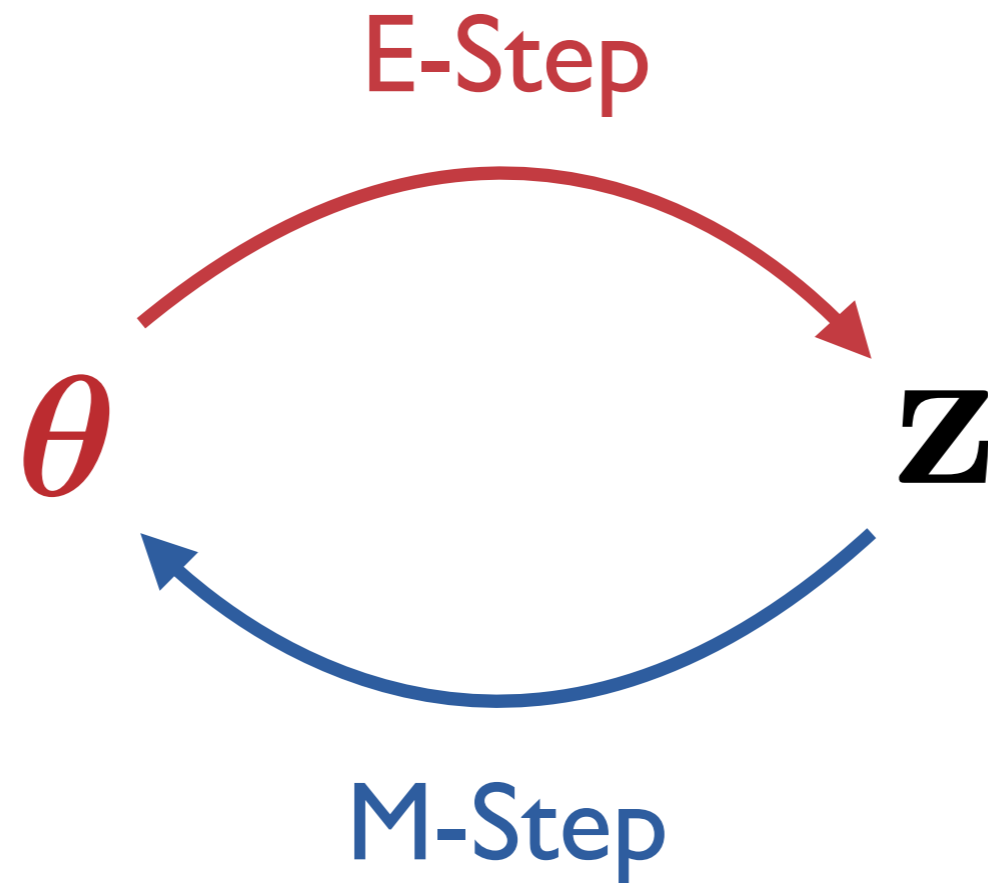
John  $\wedge$  NNP

# POS Induction Accuracy



# Hard EM without Features

---



# Hard EM without Features

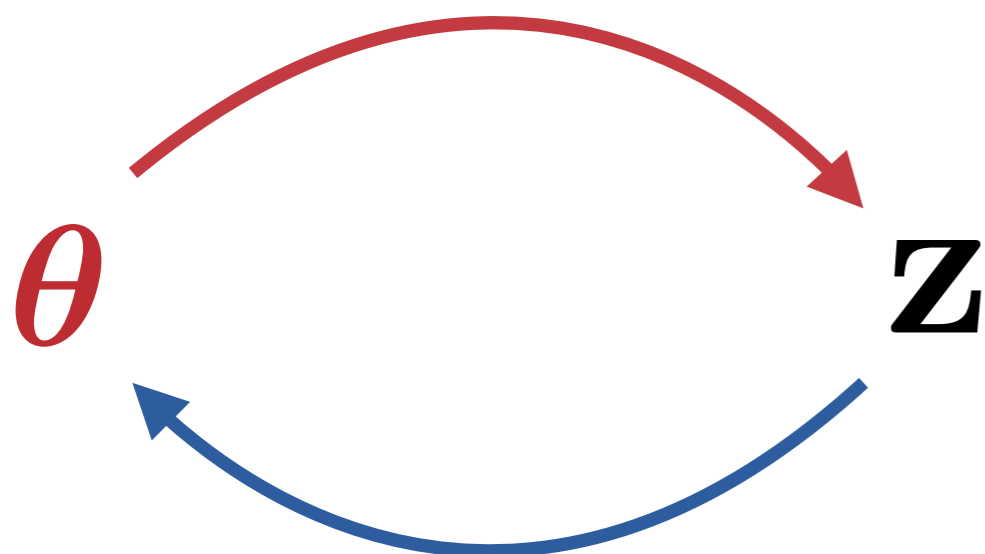
## E-Step: Dynamic Program

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

## M-Step: Divide Counts

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

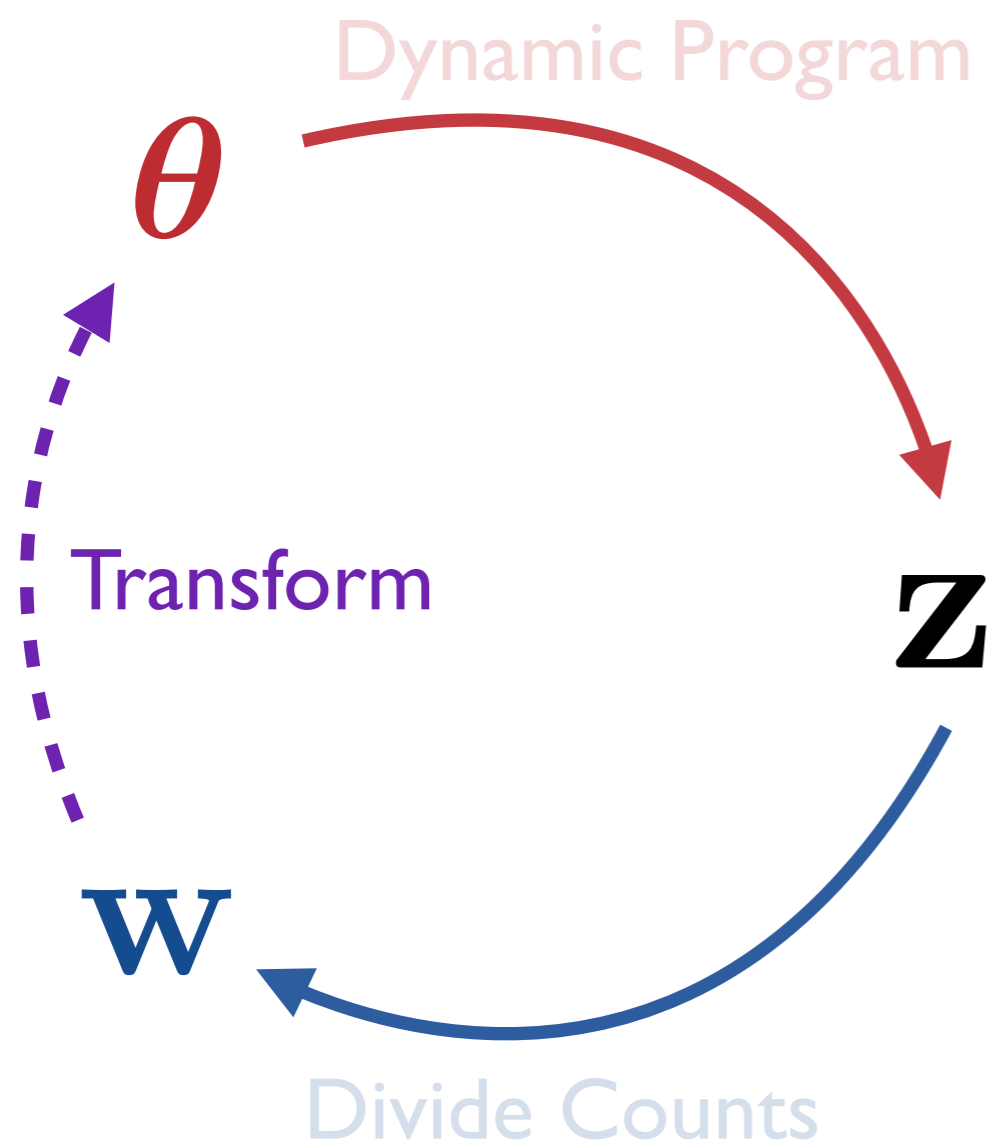
Dynamic Program



Divide Counts



# Hard EM with Features



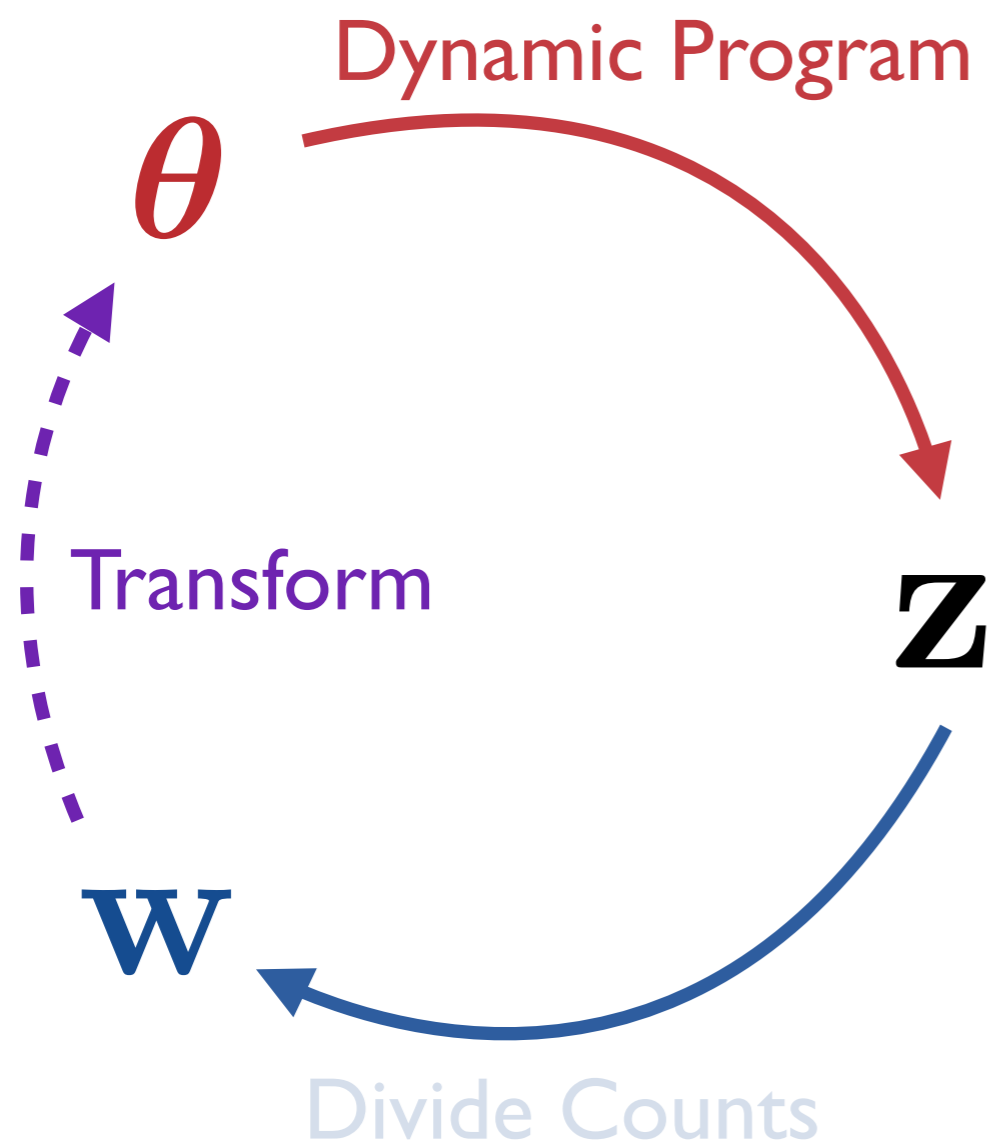
## E-Step: Dynamic Program

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

## M-Step: Divide Counts

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

# Hard EM with Features



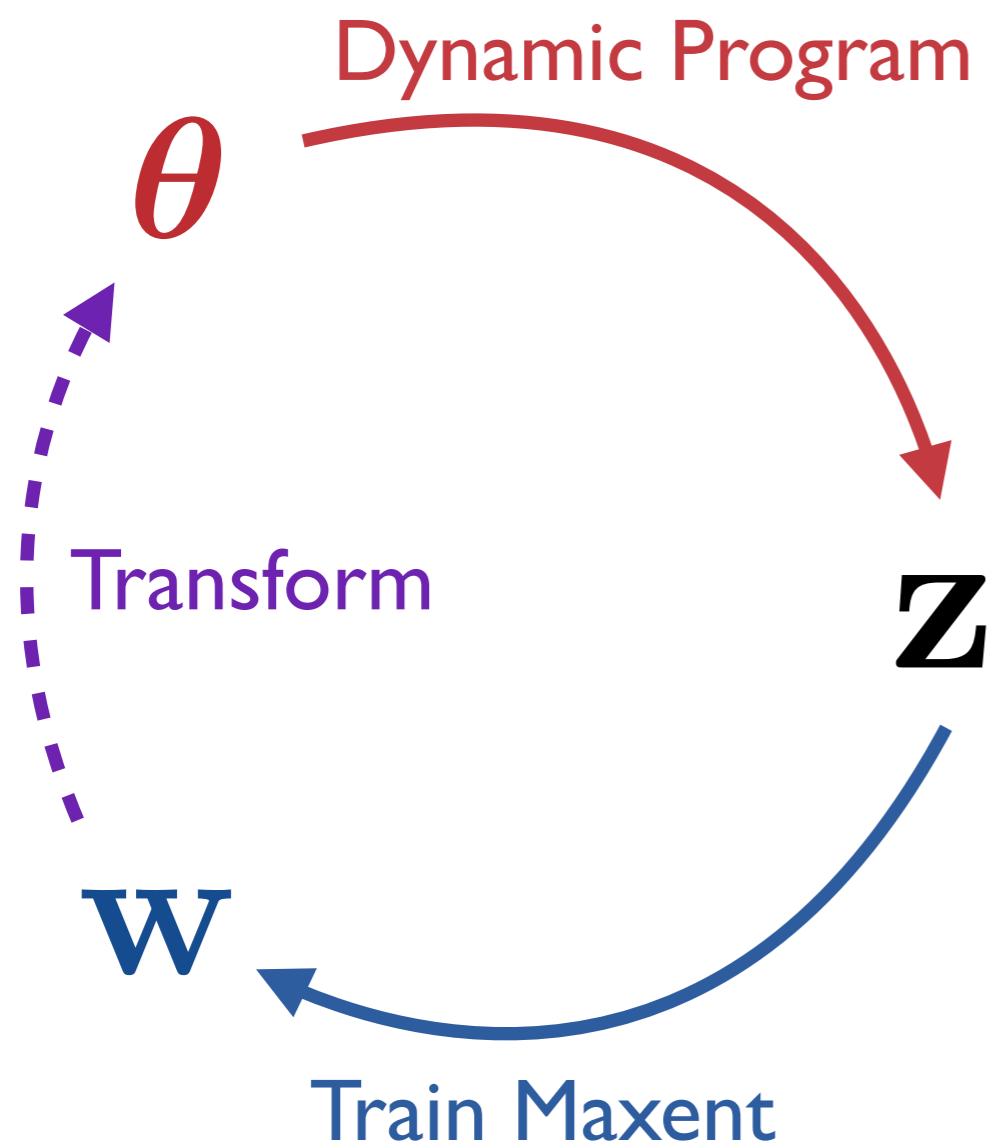
## E-Step: Dynamic Program

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

## M-Step: 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

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

## M-Step: Train Maxent

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



# Hard EM with Features

---

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

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



# 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



# Hard EM with Features

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

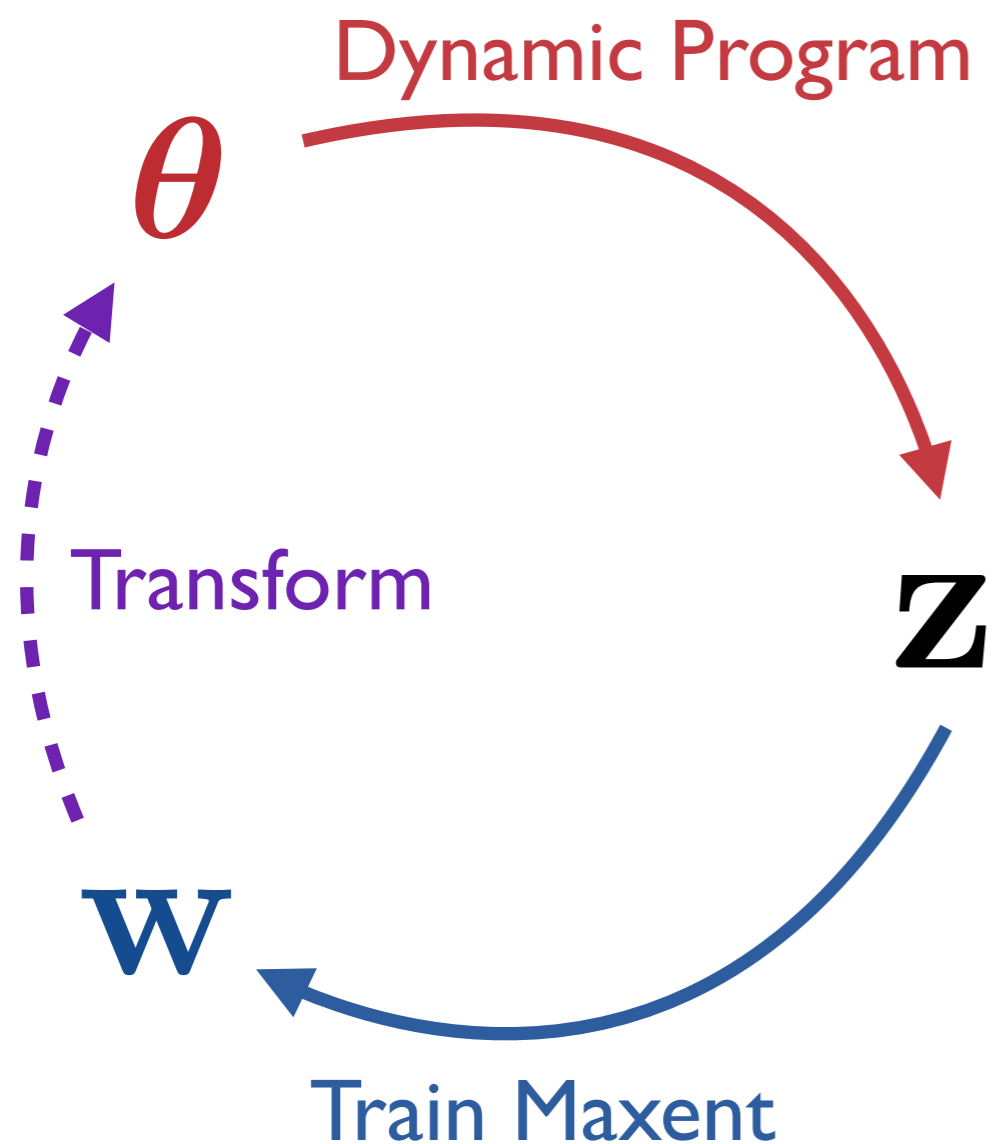
$$= \sum_i \log P(x_i | z_i; \mathbf{w}) + \dots$$

Maxent training example

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

Multiplicity

# Hard EM with Features



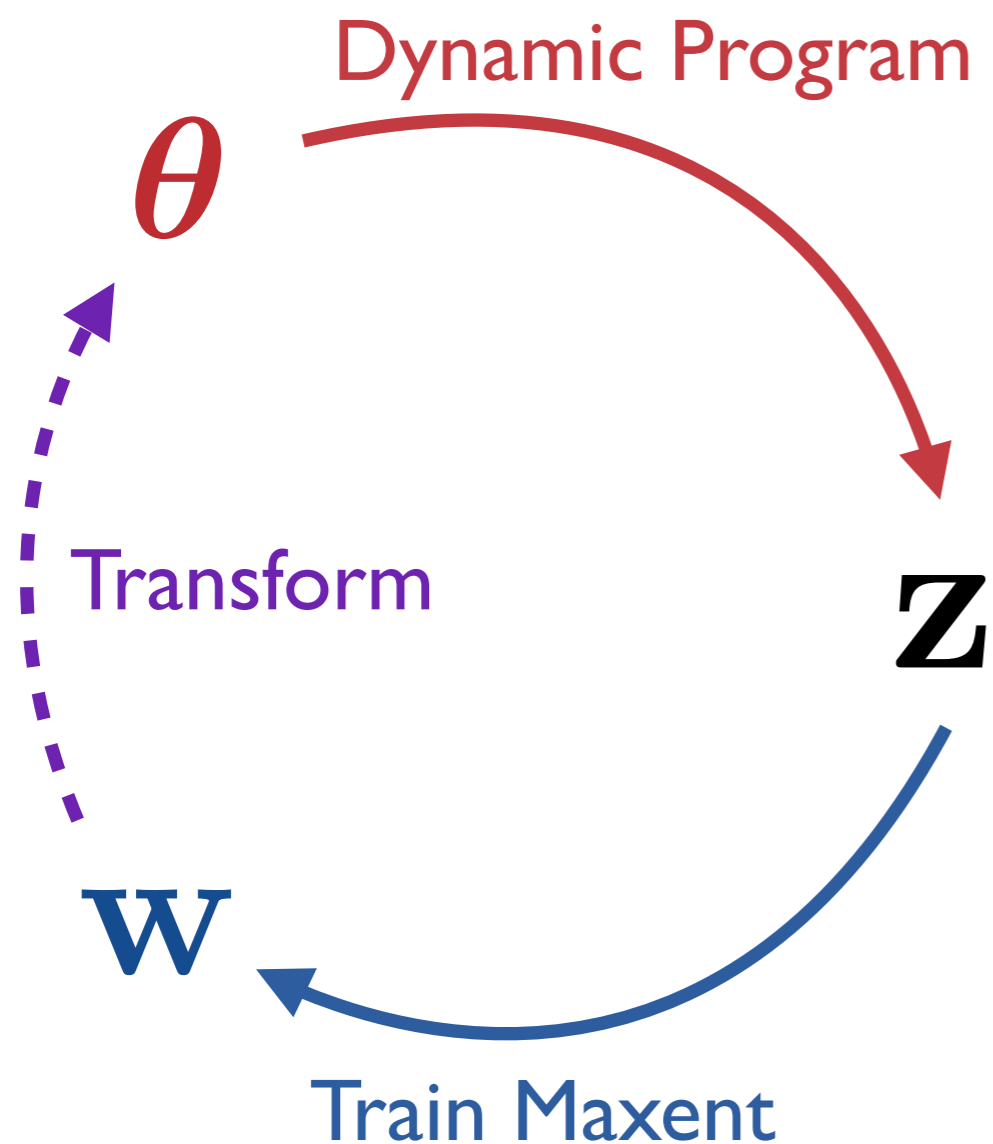
## E-Step: Dynamic Program

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

## M-Step: Train Maxent

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

# Hard EM with Features



## E-Step: Dynamic Program

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

## M-Step: Train Maxent

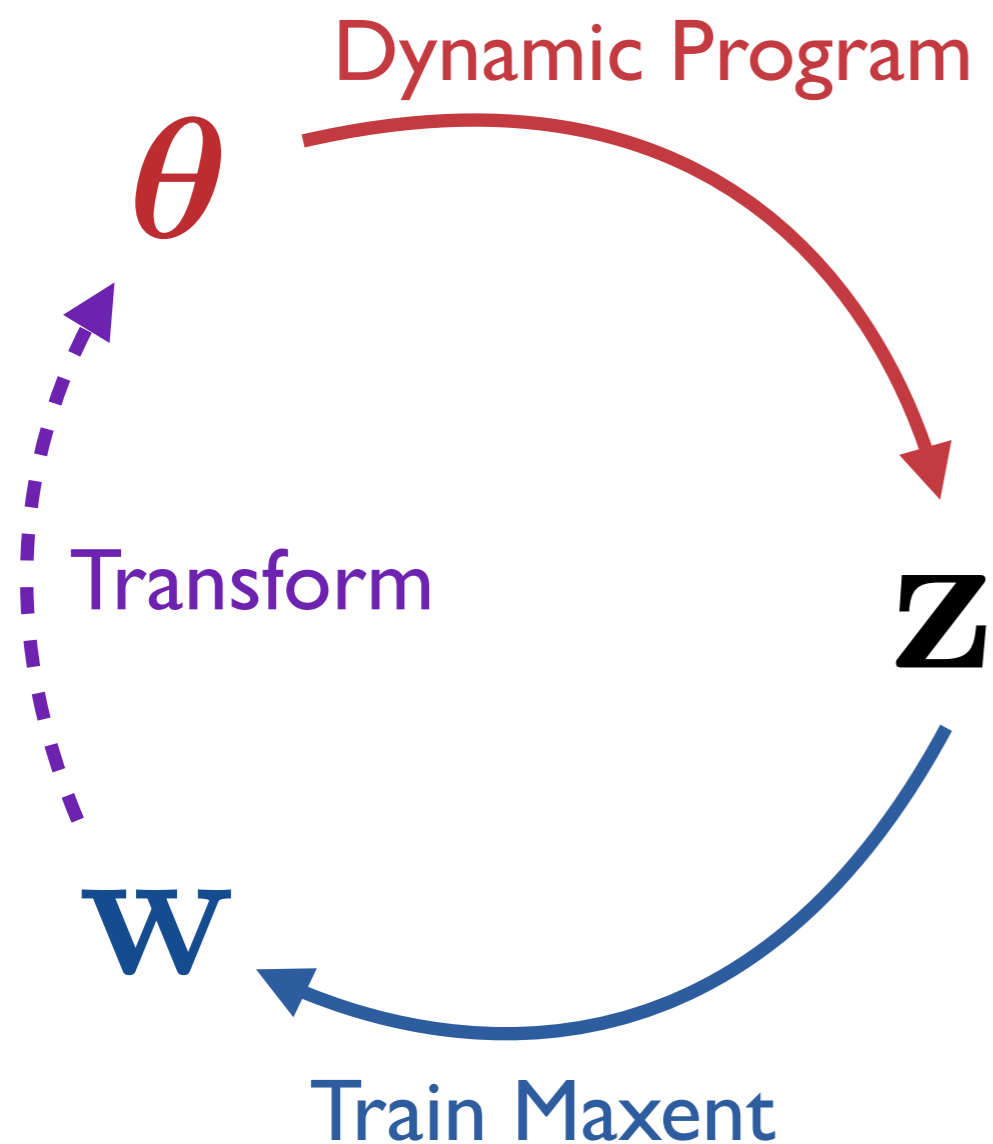
$$\mathbf{w} \leftarrow \underset{\mathbf{w}}{\operatorname{argmax}} \log P(\mathbf{x}, \mathbf{z}; \mathbf{w})$$

## 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

$$z \leftarrow \operatorname{argmax}_z P(z|x; \theta)$$

## M-Step: Train Maxent

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

## 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

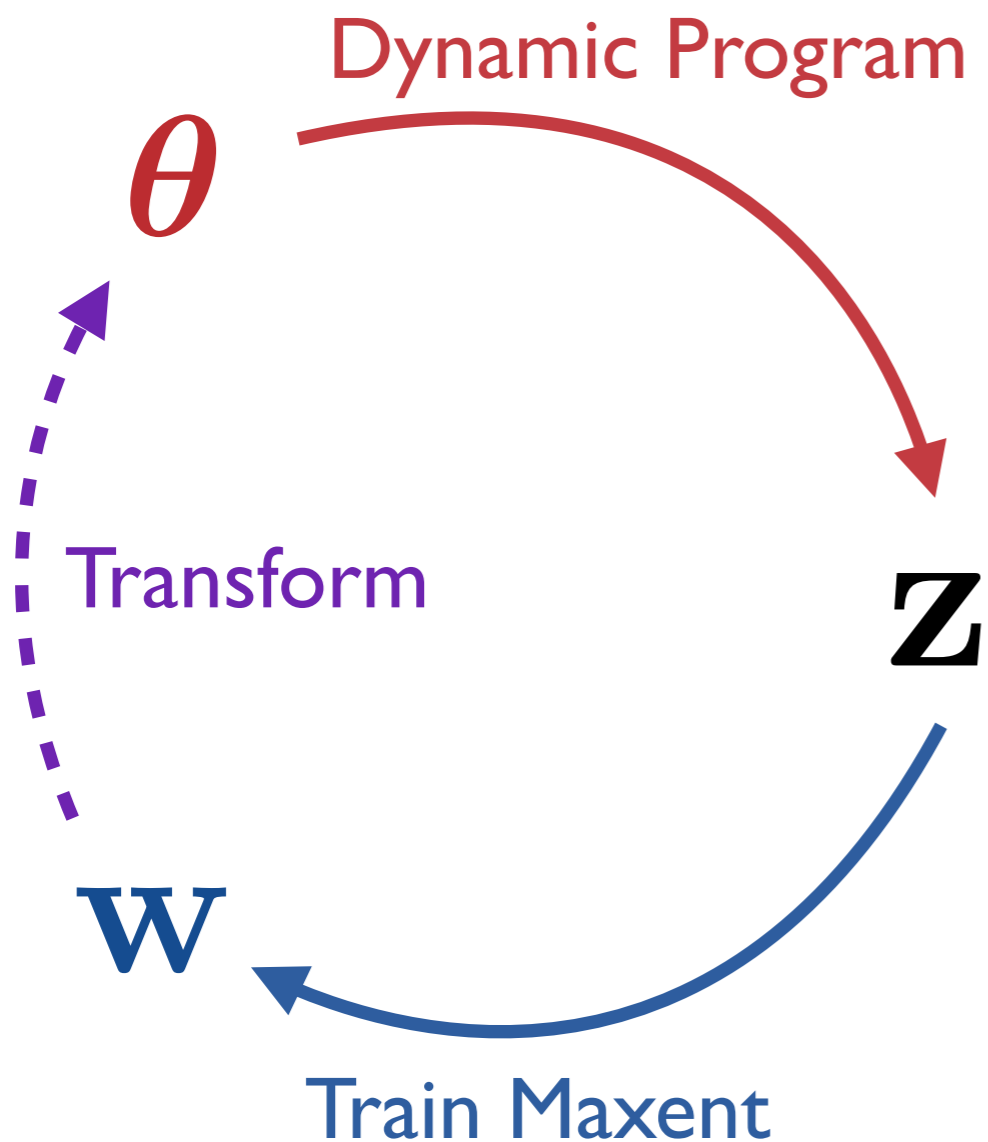
$$e(z \rightarrow x) \leftarrow \mathbb{E}[c(z \rightarrow x)]$$

## M-Step: Train Maxent

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

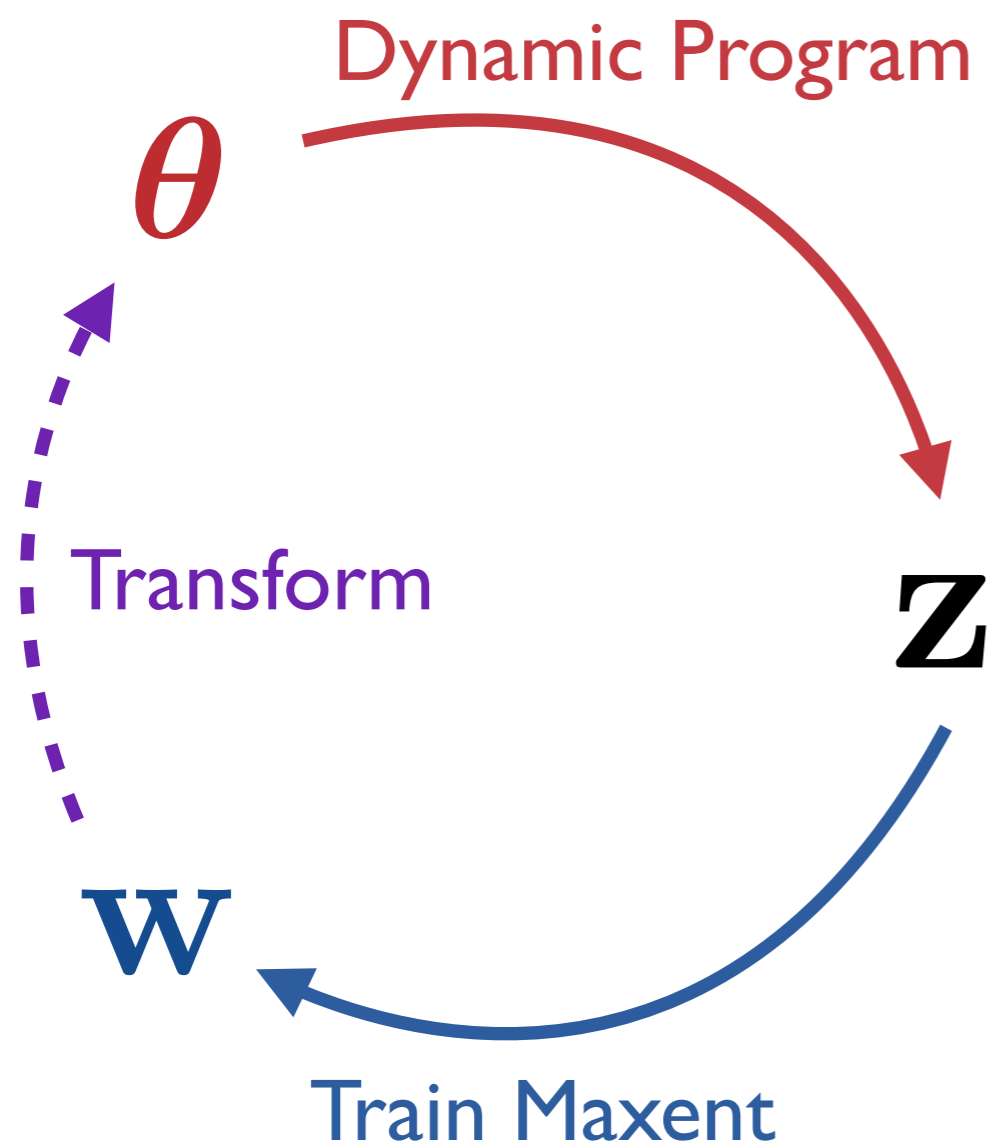
## Transform Parameters

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





# EM with Features



## E-Step: Dynamic Program

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

## M-Step: Train Maxent

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

## Transform Parameters

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



# EM without Features

---

Initialize probabilities  $\theta$

**repeat**

● Compute expected counts  $e$

● Fit parameters  $\theta$

**until** convergence

EM



# EM without Features

---

Initialize probabilities  $\theta$

**repeat**

● Compute expected counts  $e$

● Fit parameters  $\theta$

**until** convergence

A graph showing a smooth, concave-down curve representing the log-likelihood function  $L(\mathbf{w})$ . The curve starts at a low value on the left, rises to a peak, and then falls on the right. The axes are represented by simple black lines.

$L(\mathbf{w})$



# EM without Features

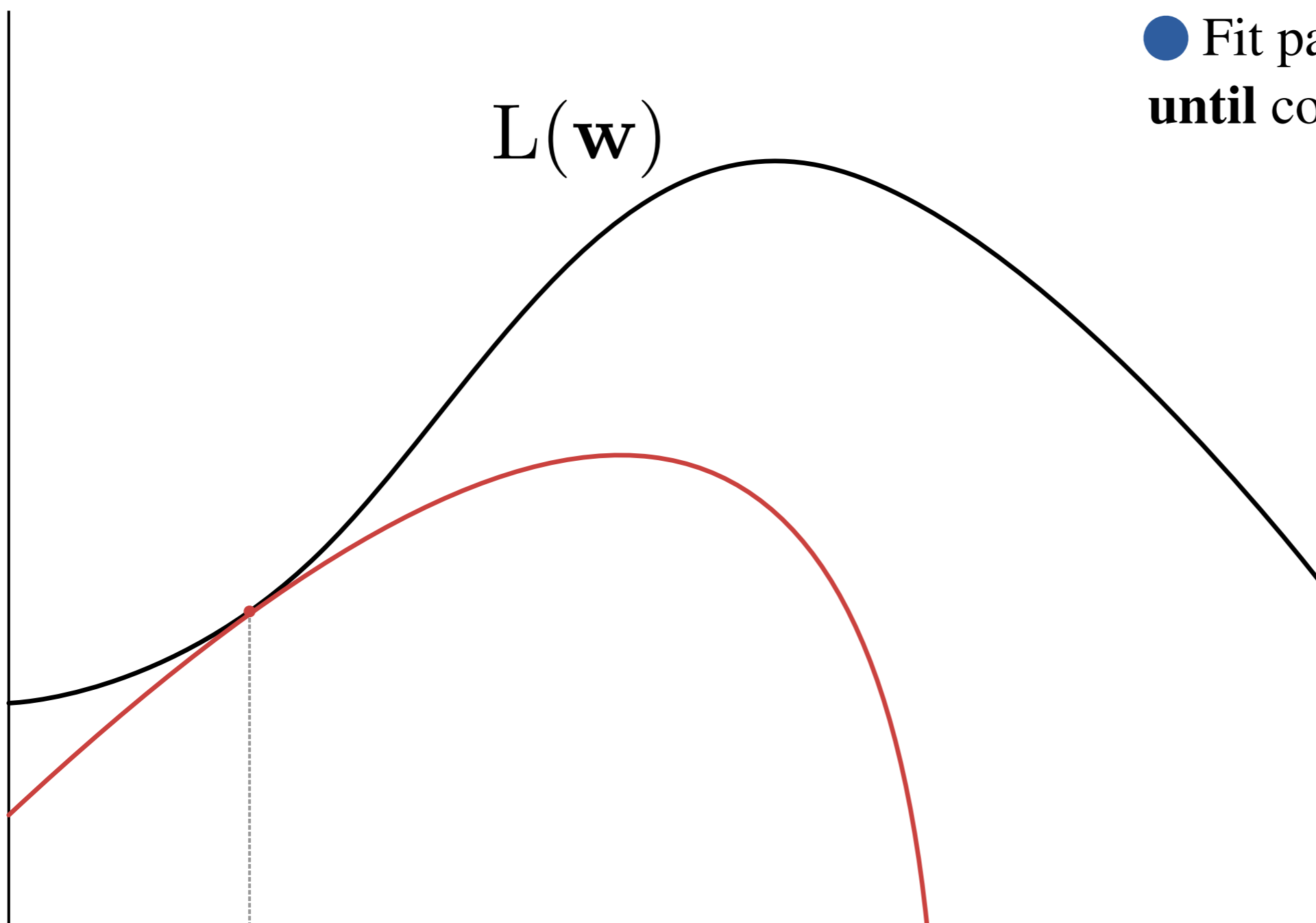
Initialize probabilities  $\theta$

**repeat**

● Compute expected counts  $e$

● Fit parameters  $\theta$

**until** convergence





# EM without Features

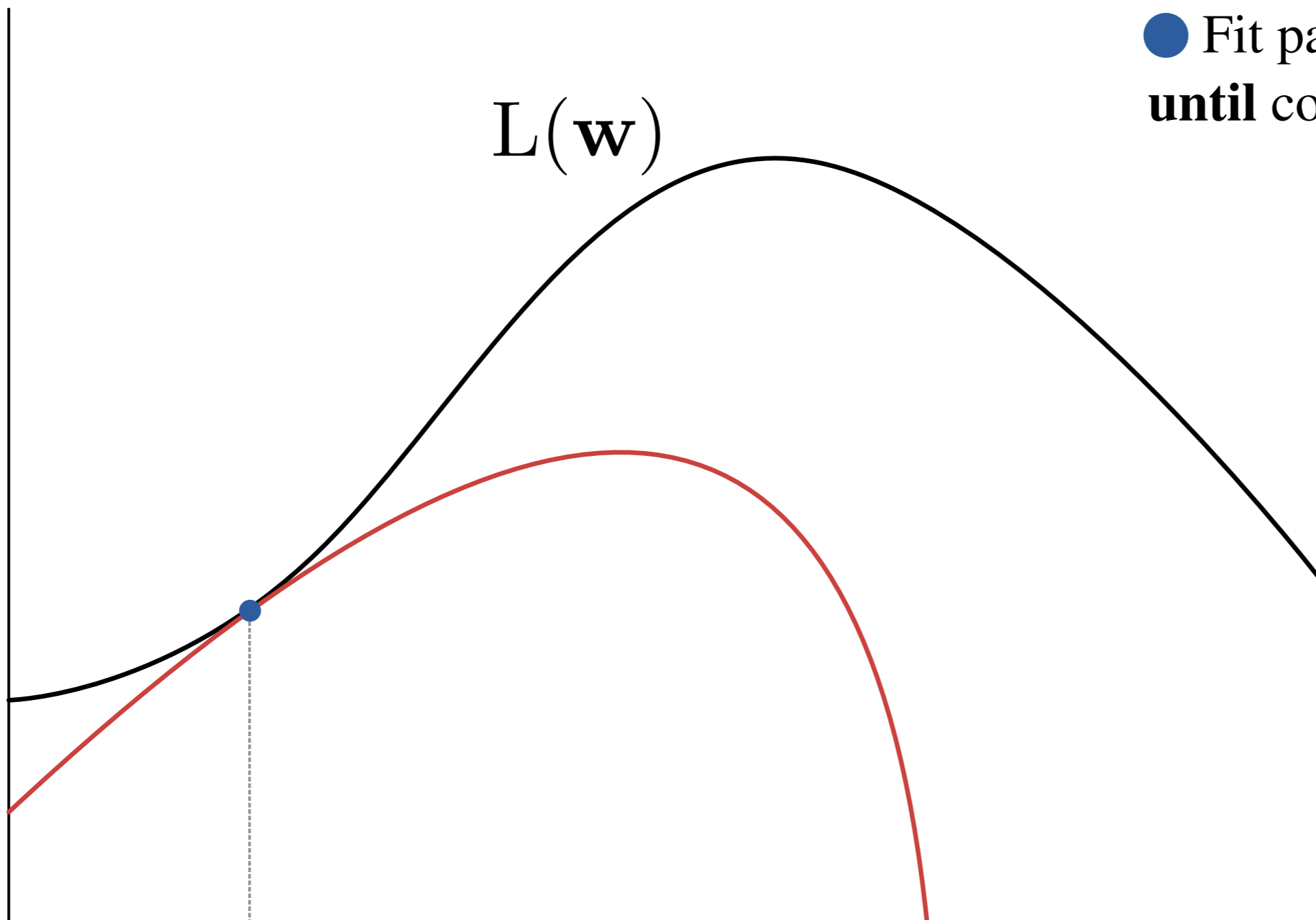
Initialize probabilities  $\theta$

**repeat**

● Compute expected counts  $e$

● Fit parameters  $\theta$

**until** convergence





# EM without Features

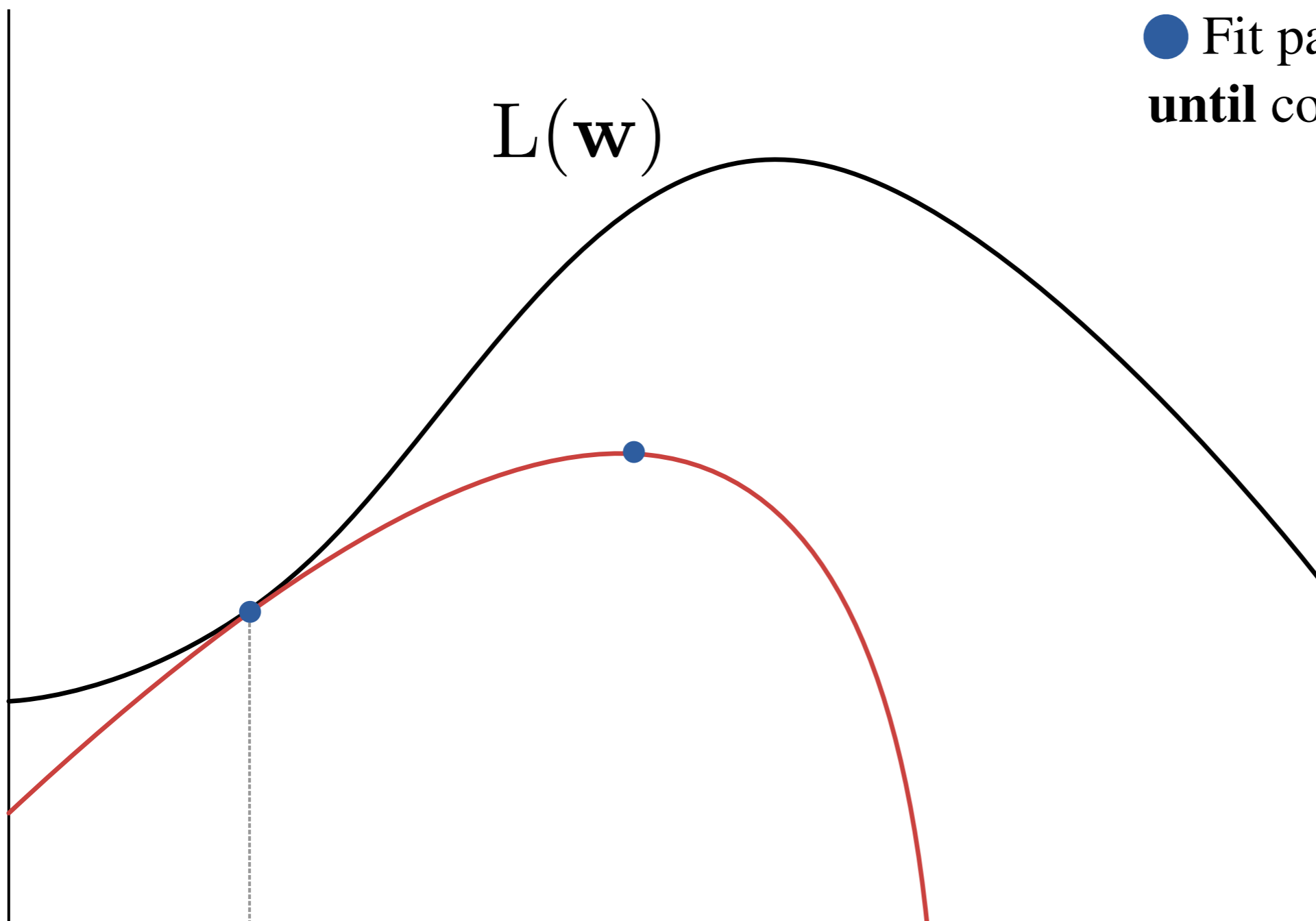
Initialize probabilities  $\theta$

**repeat**

● Compute expected counts  $e$

● Fit parameters  $\theta$

**until** convergence







# EM without Features

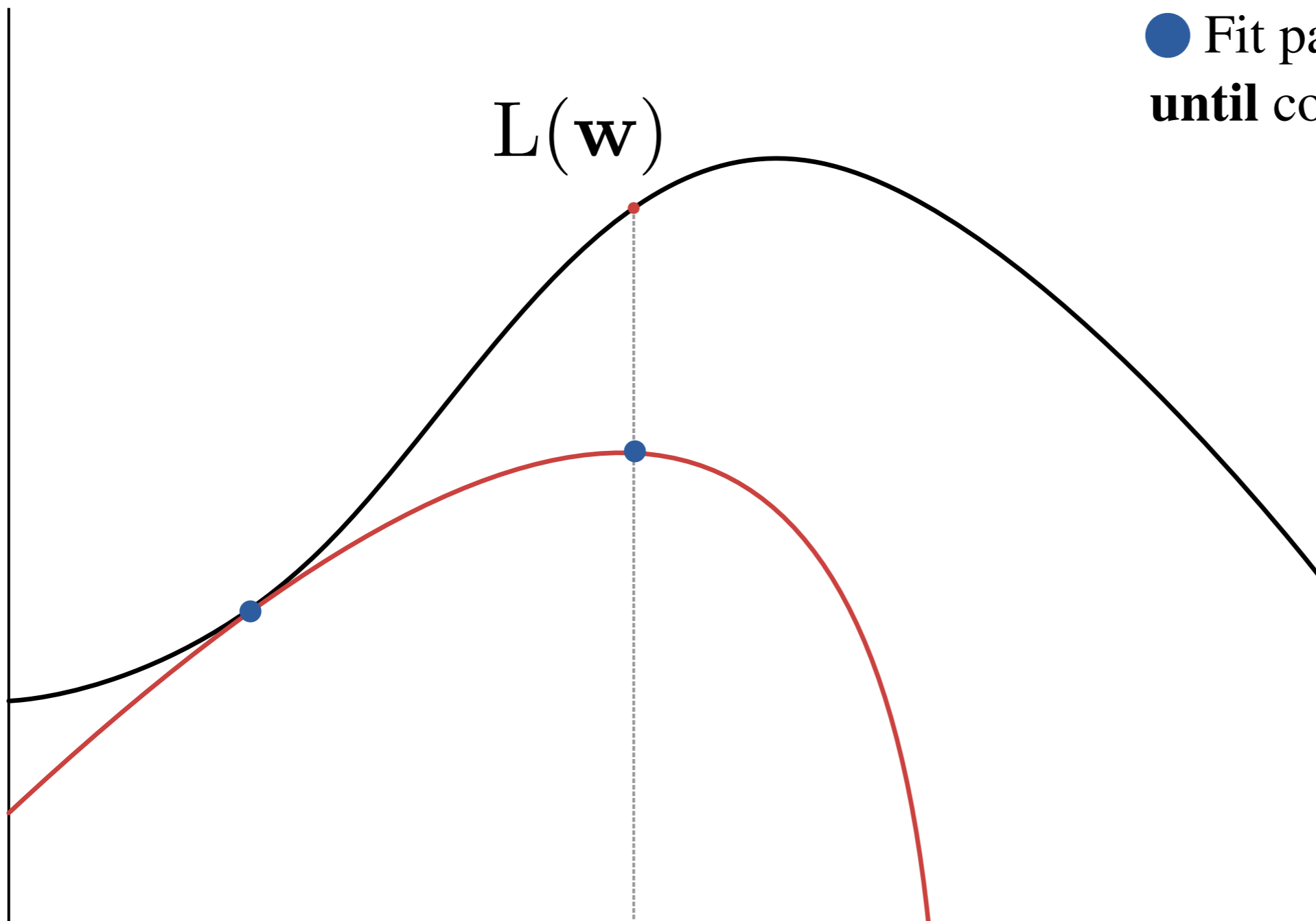
Initialize probabilities  $\theta$

**repeat**

● Compute expected counts  $e$

● Fit parameters  $\theta$

**until** convergence





# EM without Features

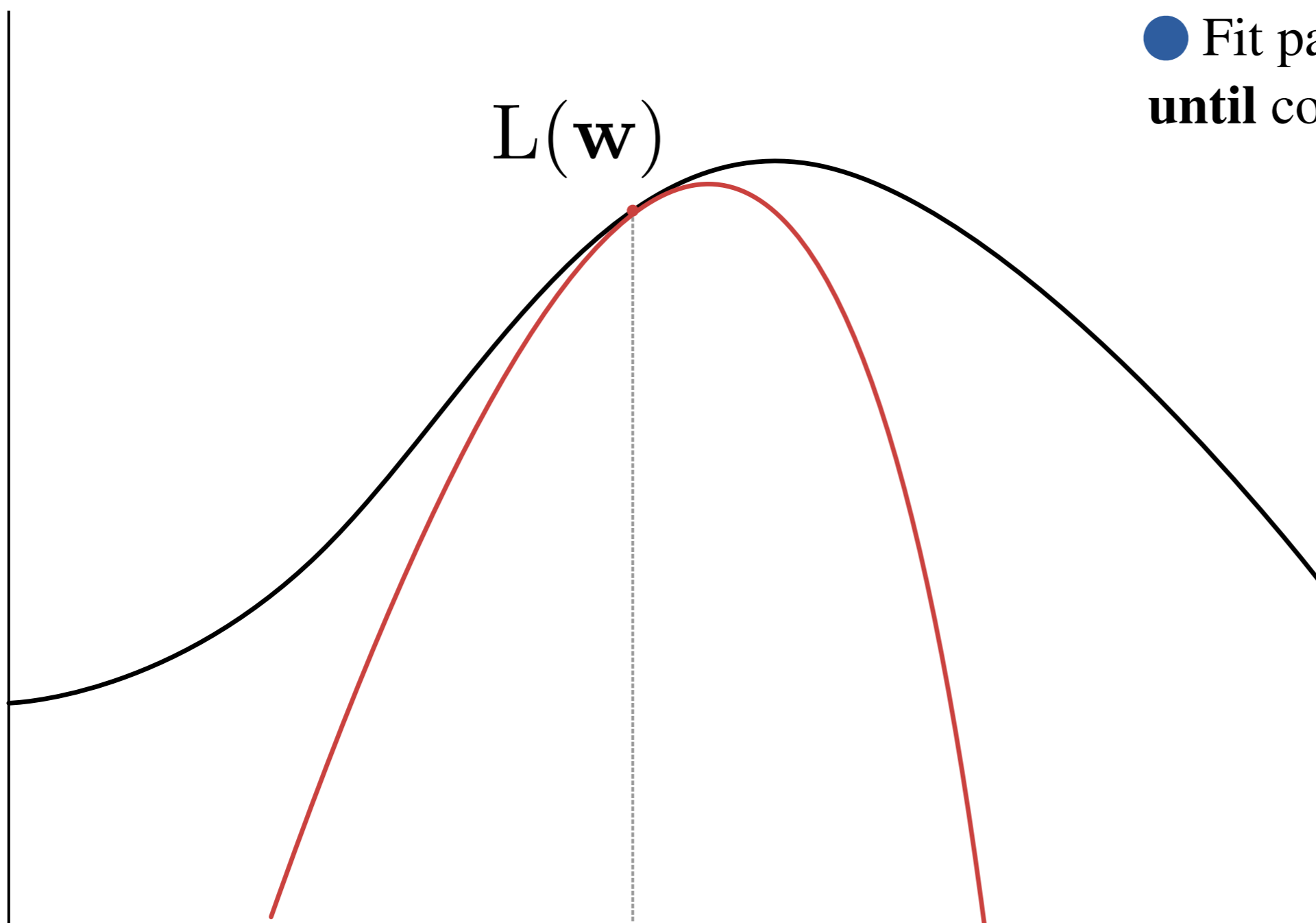
Initialize probabilities  $\theta$

**repeat**

● Compute expected counts  $e$

● Fit parameters  $\theta$

**until** convergence





# EM without Features

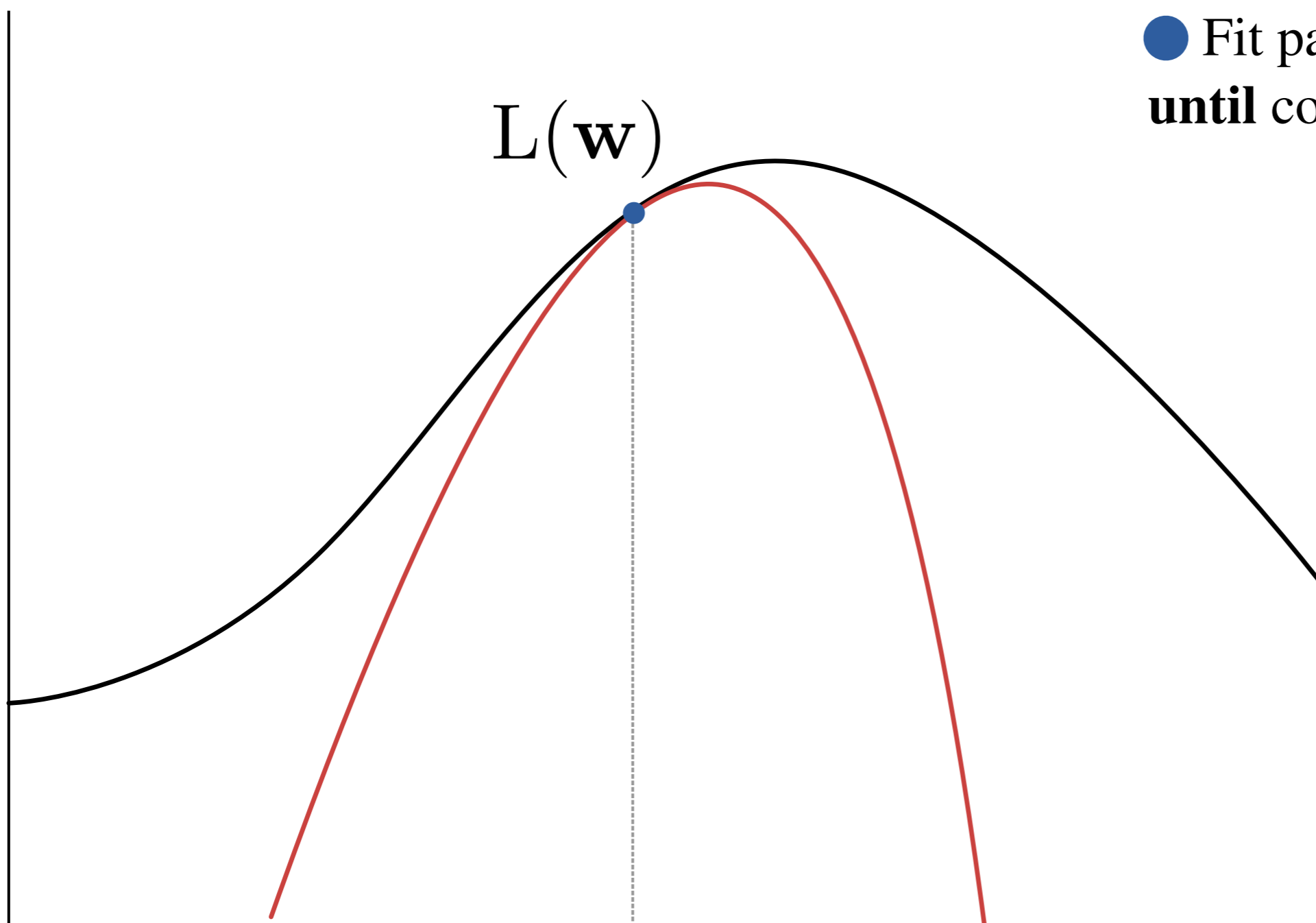
Initialize probabilities  $\theta$

**repeat**

● Compute expected counts  $e$

● Fit parameters  $\theta$

**until** convergence





# EM without Features

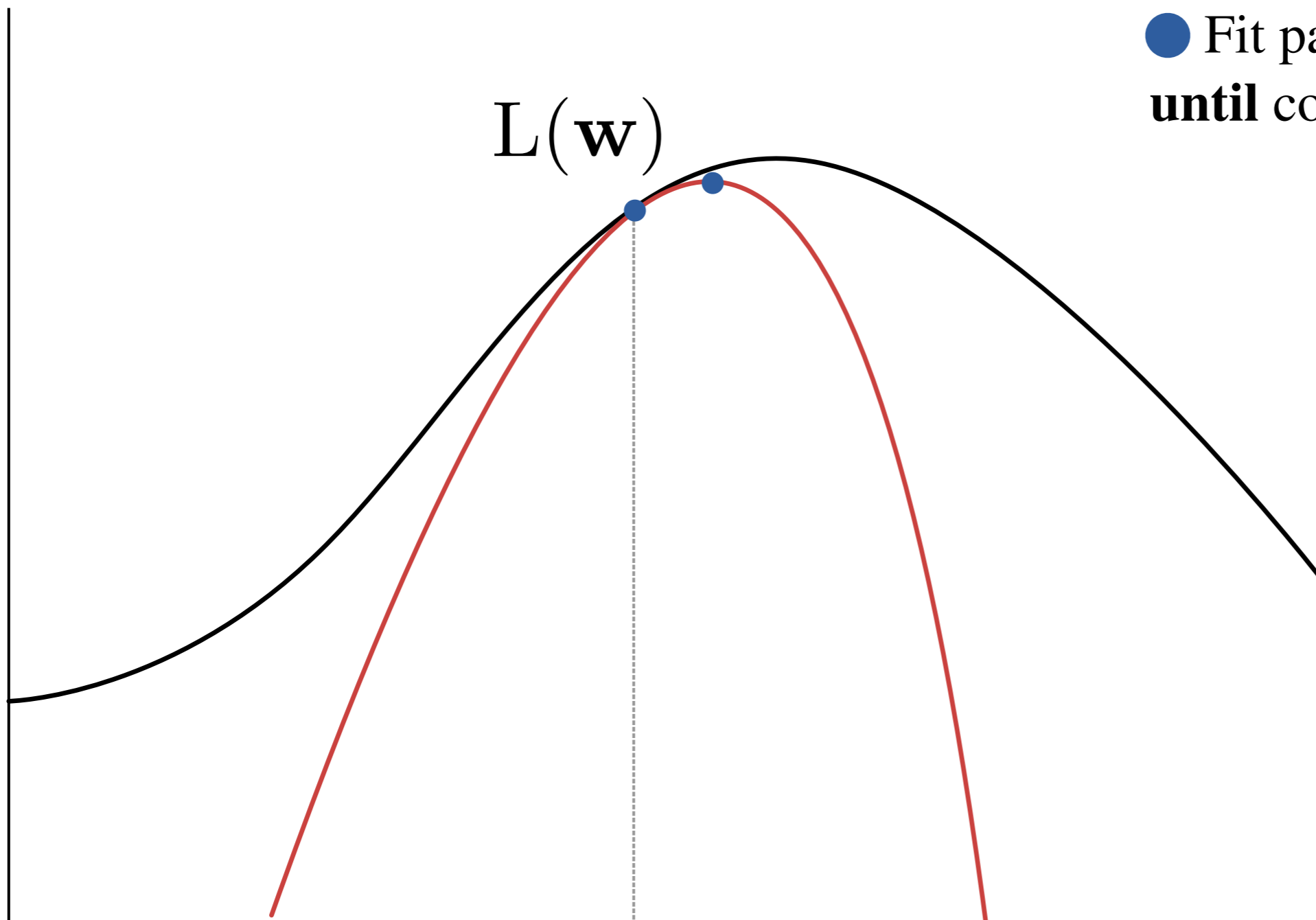
Initialize probabilities  $\theta$

**repeat**

● Compute expected counts  $e$

● Fit parameters  $\theta$

**until** convergence





# EM without Features

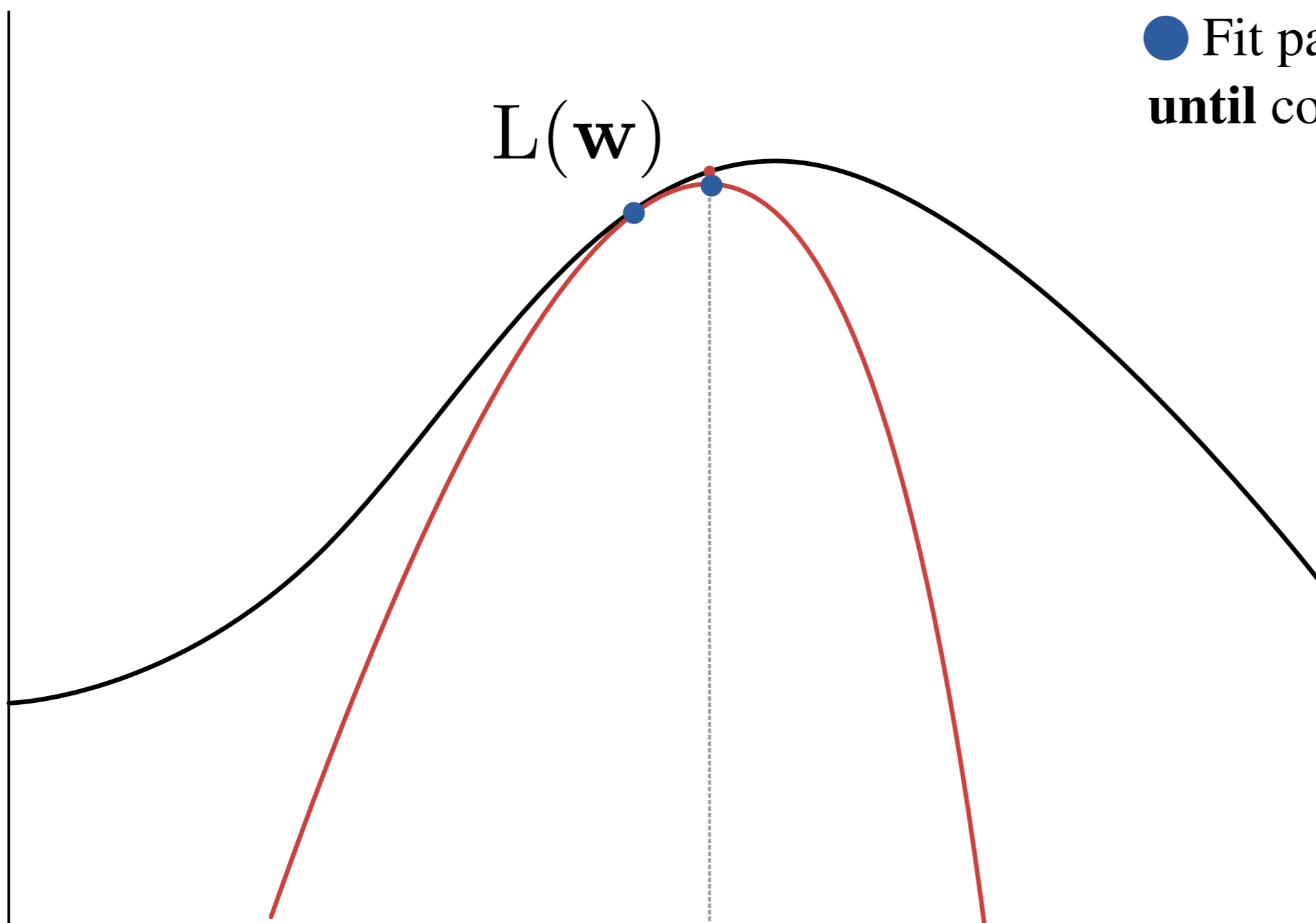
Initialize probabilities  $\theta$

**repeat**

● Compute expected counts  $e$

● Fit parameters  $\theta$

**until** convergence





# EM without Features

---

Initialize probabilities  $\theta$

**repeat**

● Compute expected counts  $e$

● Fit parameters  $\theta$

**until** convergence

EM



# EM with Features

---

Initialize weights  $w$

**repeat**

● Compute expected counts  $e$

● Fit parameters  $w$

● Transform  $w$  to  $\theta$

**until** convergence

EM



# 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  $\ell(\mathbf{w}, \mathbf{e})$

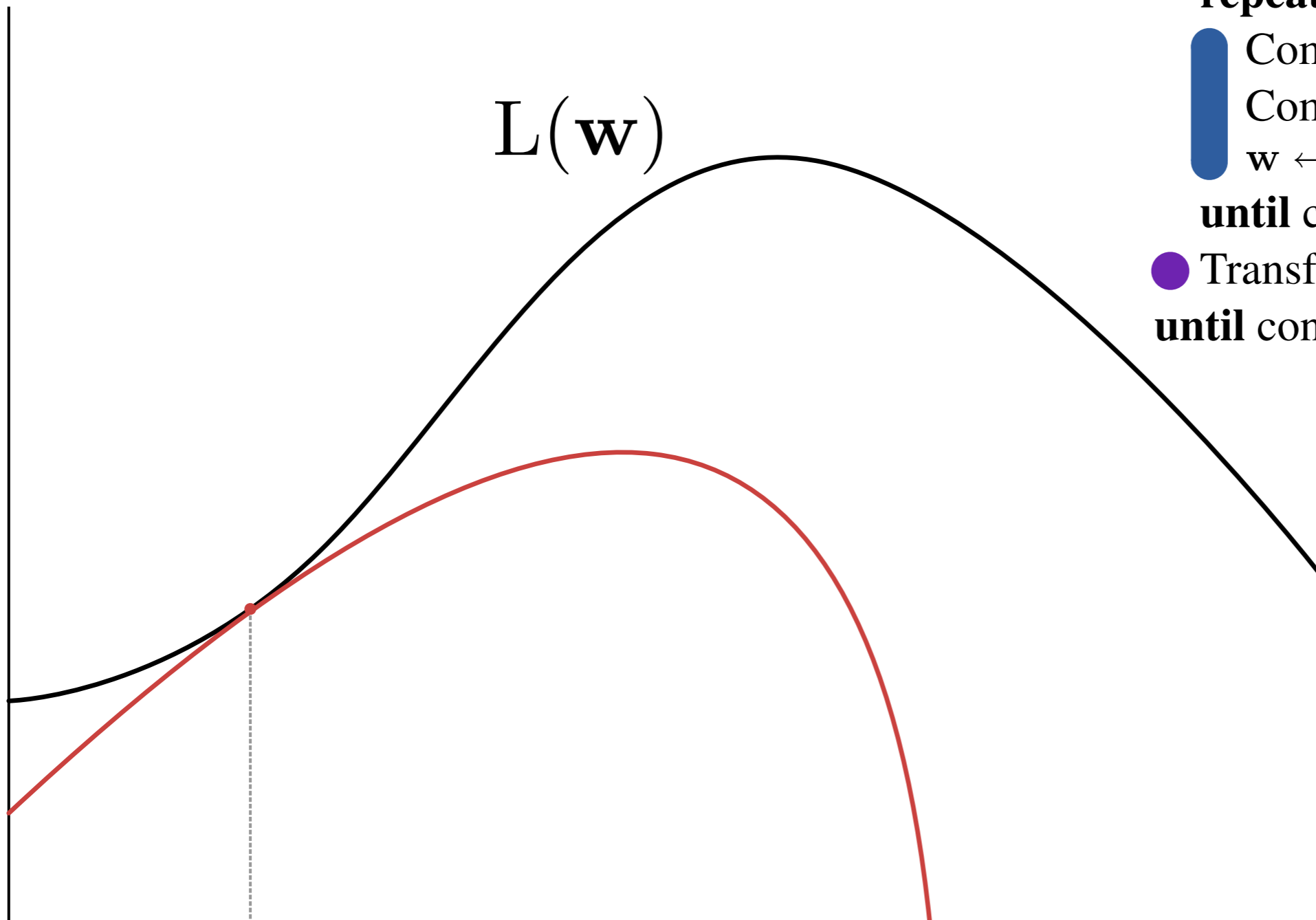
■ Compute  $\nabla \ell(\mathbf{w}, \mathbf{e})$

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

**until** convergence

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

**until** convergence





# EM with Features

Initialize weights  $\mathbf{w}$

**repeat**

● Compute expected counts  $\mathbf{e}$

**repeat**

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

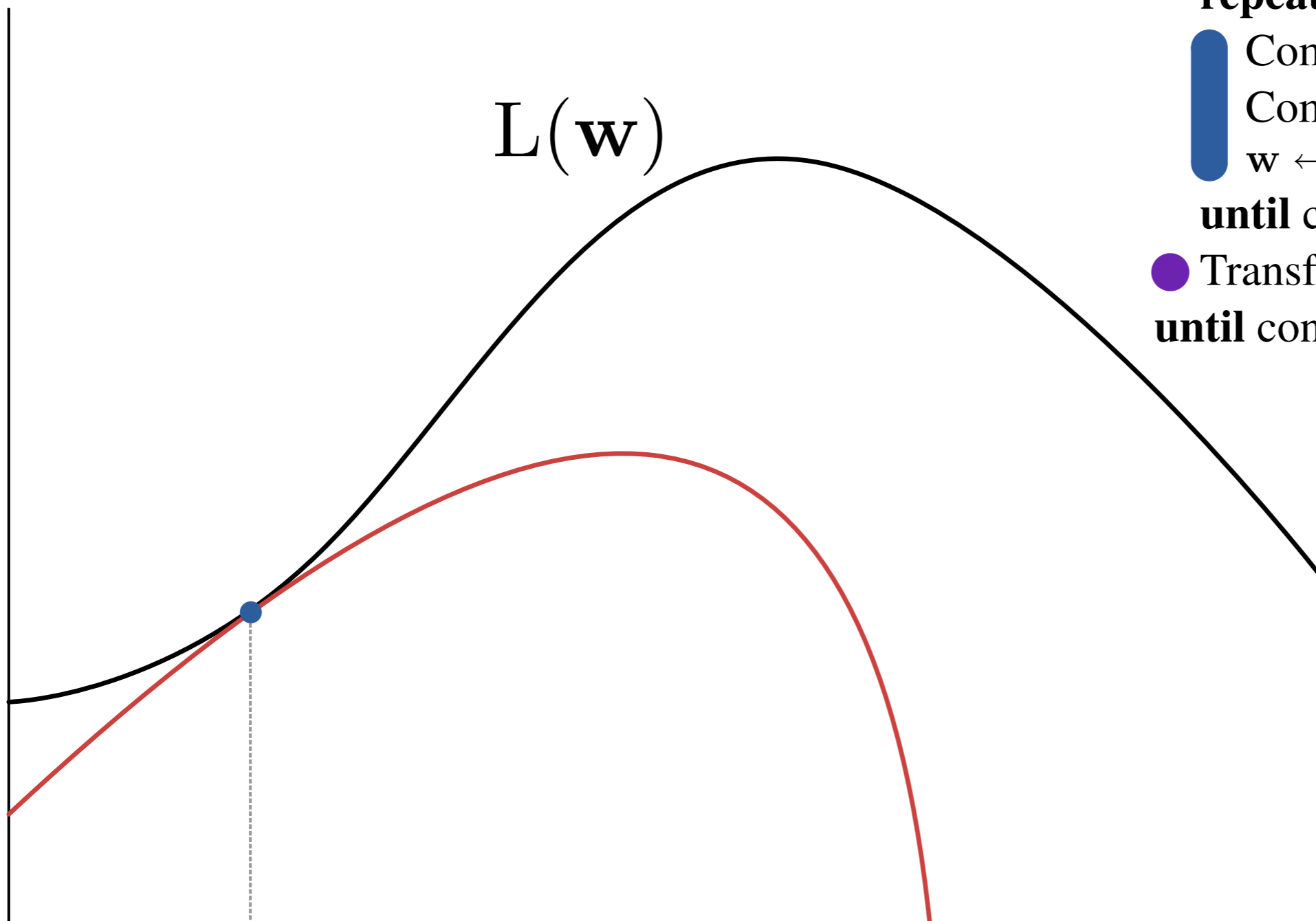
■ Compute  $\nabla \ell(\mathbf{w}, \mathbf{e})$

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

**until** convergence

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

**until** convergence





# EM with Features

Initialize weights  $\mathbf{w}$

**repeat**

● Compute expected counts  $\mathbf{e}$

**repeat**

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

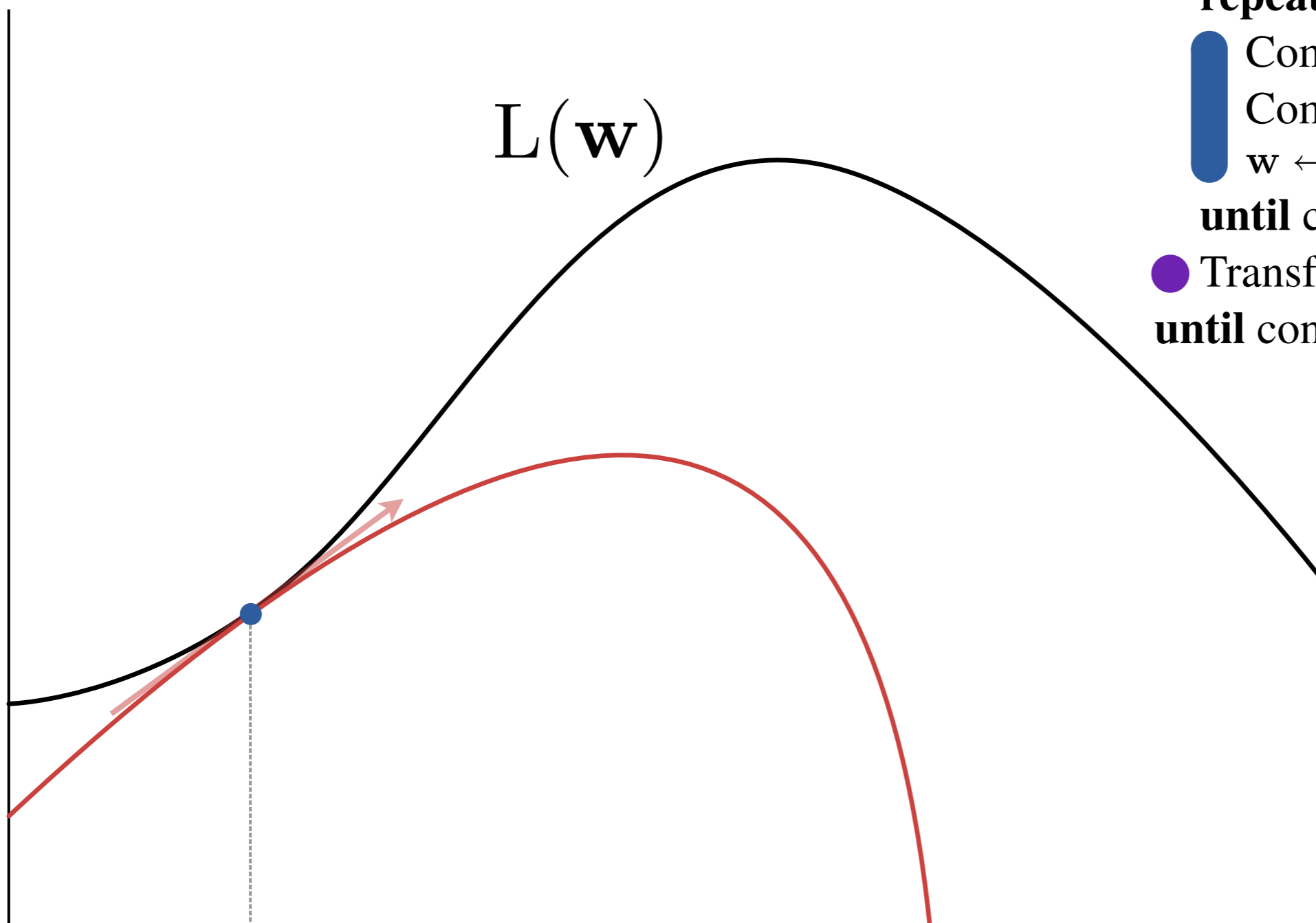
● Compute  $\nabla \ell(\mathbf{w}, \mathbf{e})$

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

**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  $\ell(\mathbf{w}, \mathbf{e})$

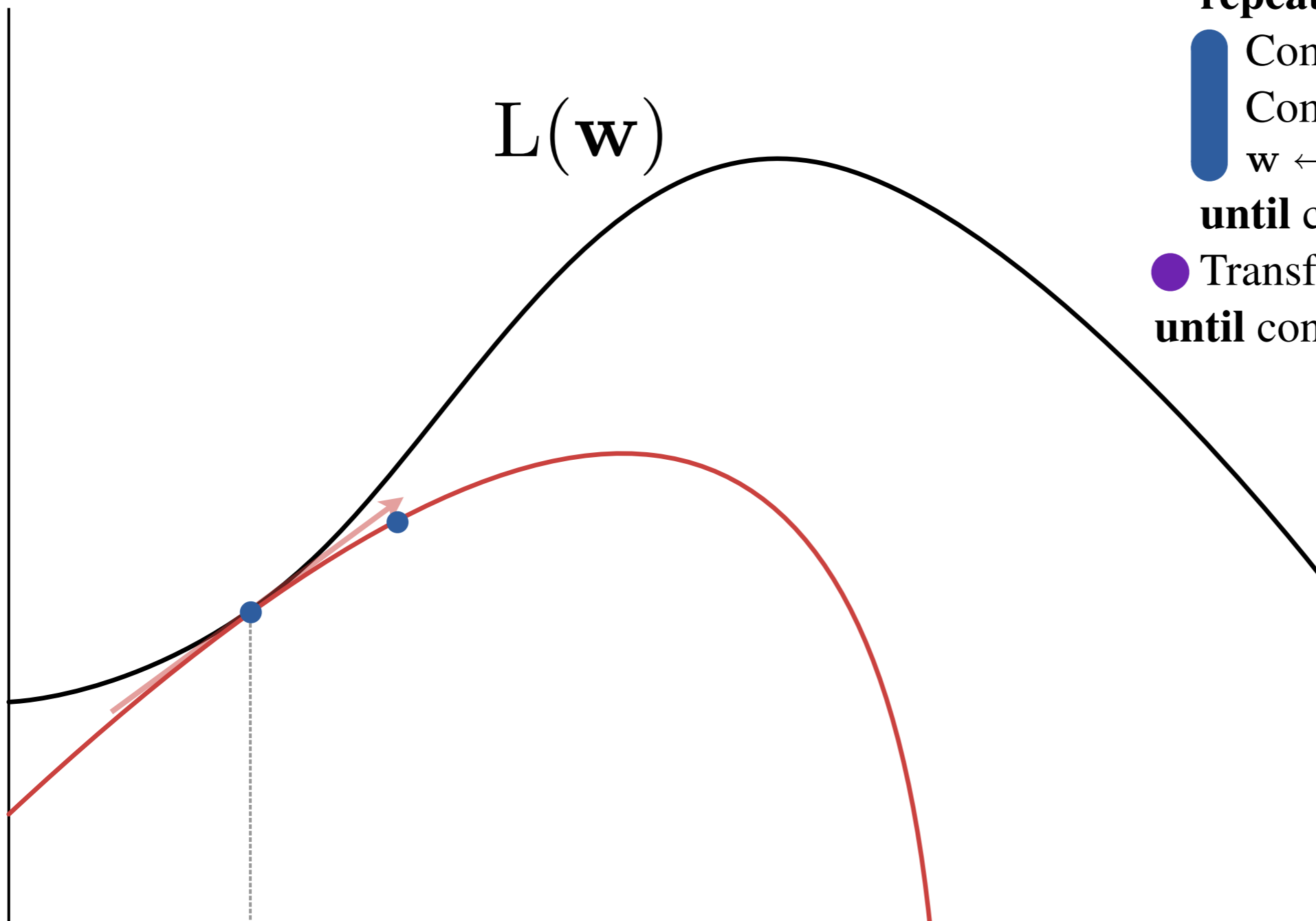
● Compute  $\nabla \ell(\mathbf{w}, \mathbf{e})$

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

**until convergence**

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

**until convergence**





# EM with Features

Initialize weights  $\mathbf{w}$

**repeat**

● Compute expected counts  $\mathbf{e}$

**repeat**

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

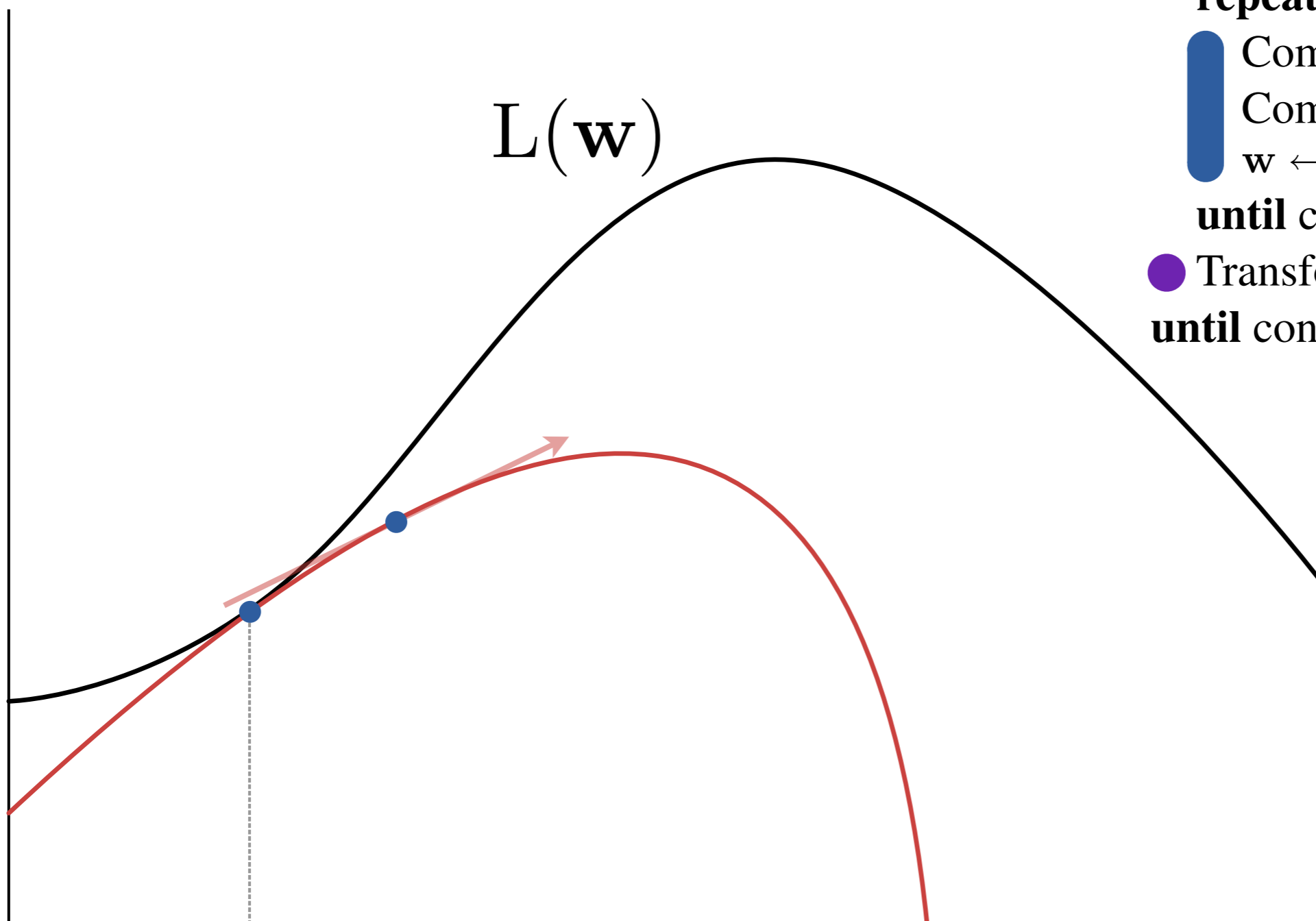
■ Compute  $\nabla \ell(\mathbf{w}, \mathbf{e})$

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

**until** convergence

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

**until** convergence





# EM with Features

Initialize weights  $\mathbf{w}$

**repeat**

● Compute expected counts  $\mathbf{e}$

**repeat**

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

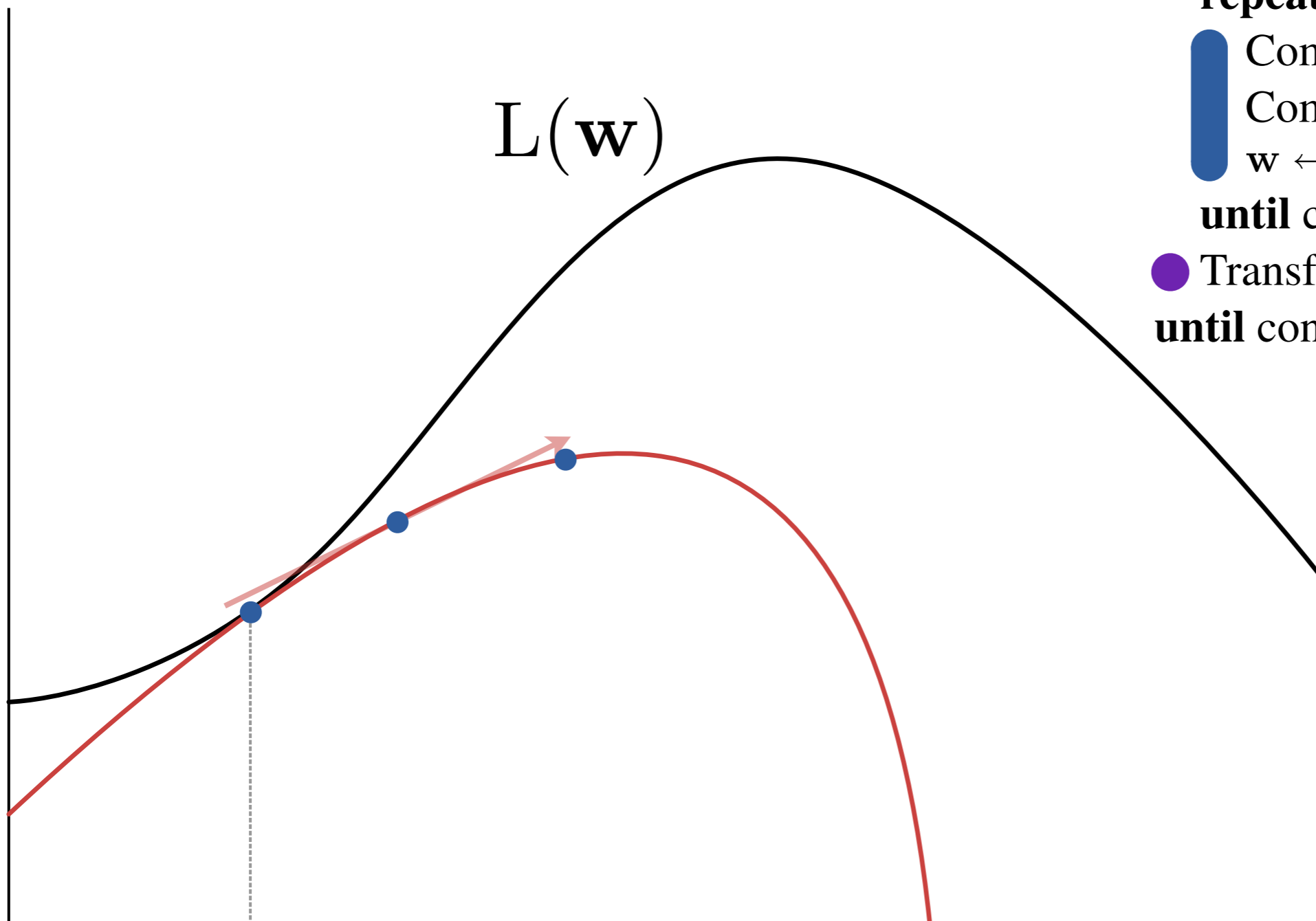
■ Compute  $\nabla \ell(\mathbf{w}, \mathbf{e})$

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

**until** convergence

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

**until** convergence





# EM with Features

Initialize weights  $w$

**repeat**

● Compute expected counts  $e$

**repeat**

■ Compute  $\ell(w, e)$

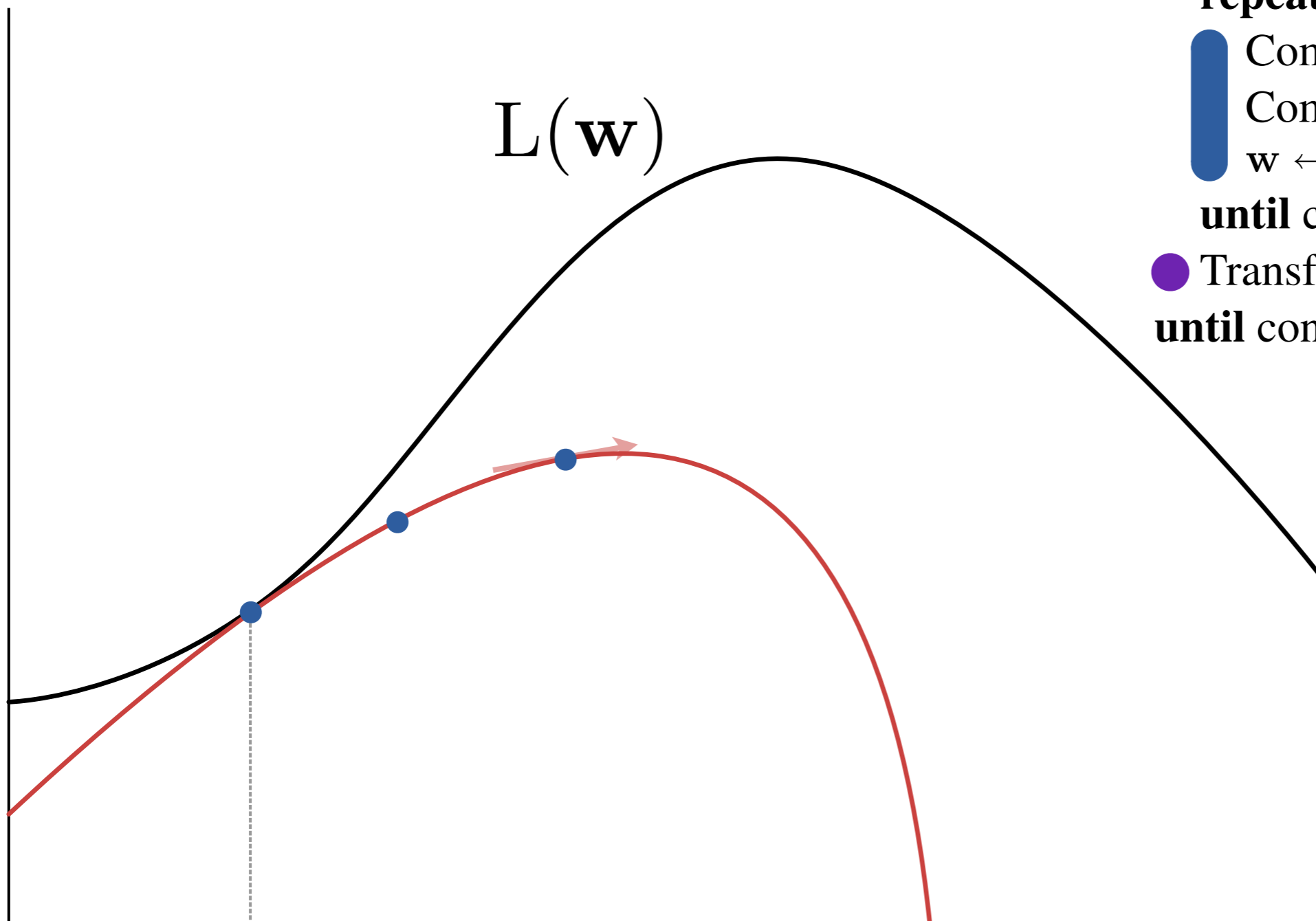
■ Compute  $\nabla \ell(w, e)$

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

**until** convergence

● Transform  $w$  to  $\theta$

**until** convergence





# EM with Features

Initialize weights  $\mathbf{w}$

**repeat**

● Compute expected counts  $\mathbf{e}$

**repeat**

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

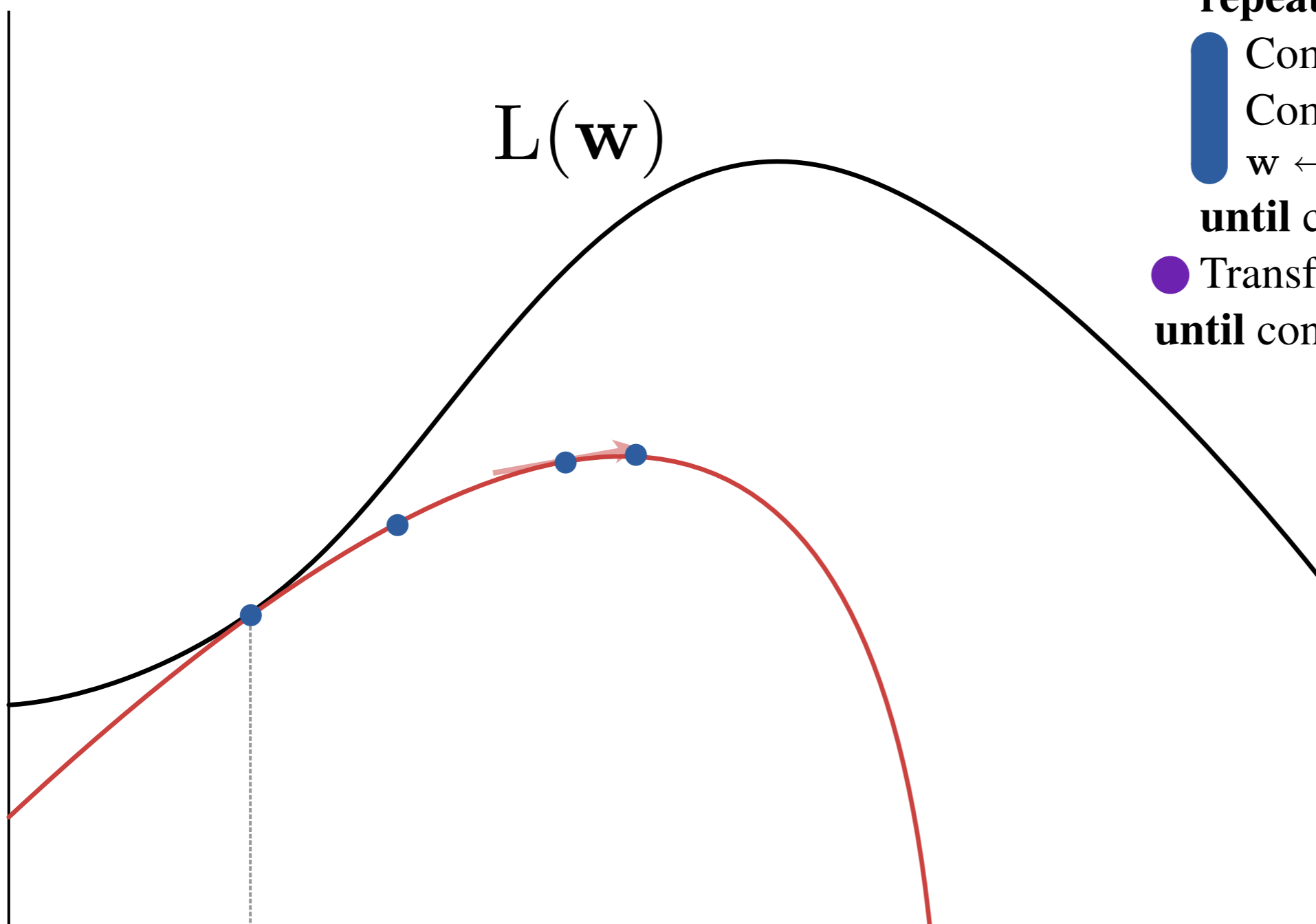
■ Compute  $\nabla \ell(\mathbf{w}, \mathbf{e})$

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

**until** convergence

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

**until** convergence







# EM with Features

Initialize weights  $\mathbf{w}$

**repeat**

● Compute expected counts  $\mathbf{e}$

**repeat**

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

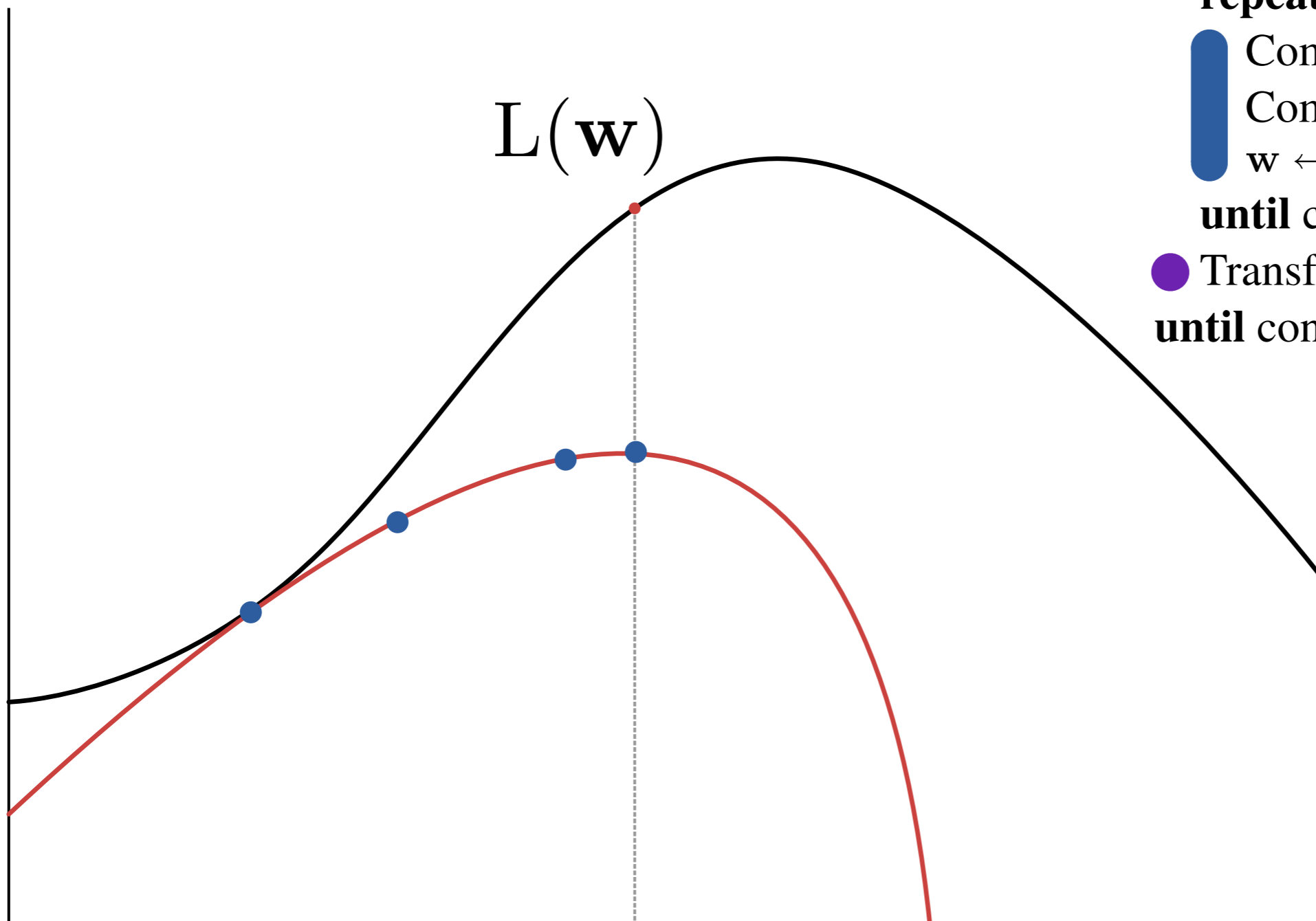
● Compute  $\nabla \ell(\mathbf{w}, \mathbf{e})$

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

**until** convergence

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

**until** convergence





# EM with Features

Initialize weights  $\mathbf{w}$

**repeat**

● Compute expected counts  $\mathbf{e}$

**repeat**

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

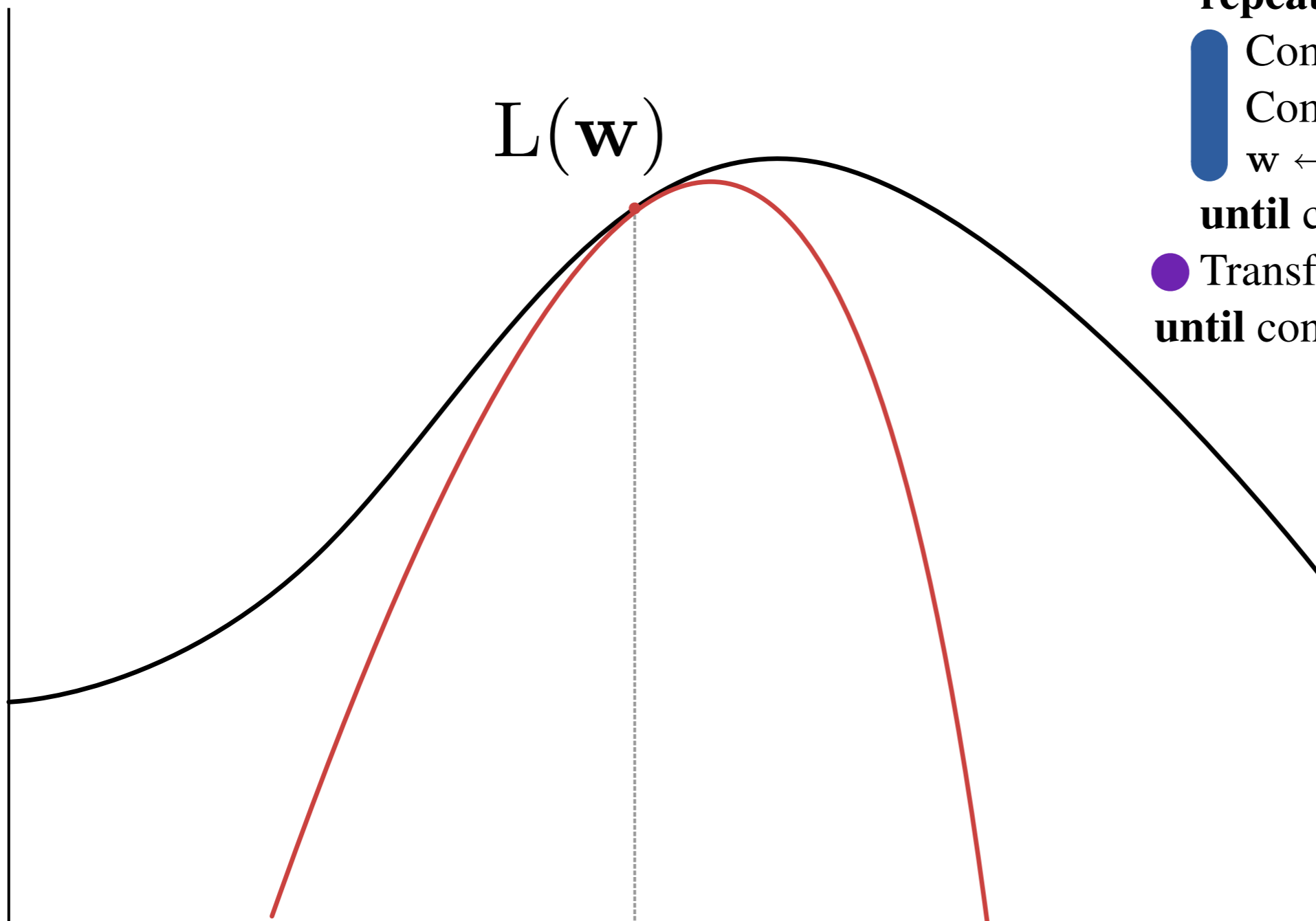
■ Compute  $\nabla \ell(\mathbf{w}, \mathbf{e})$

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

**until** convergence

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

**until** convergence





# EM with Features

Initialize weights  $\mathbf{w}$

**repeat**

● Compute expected counts  $\mathbf{e}$

**repeat**

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

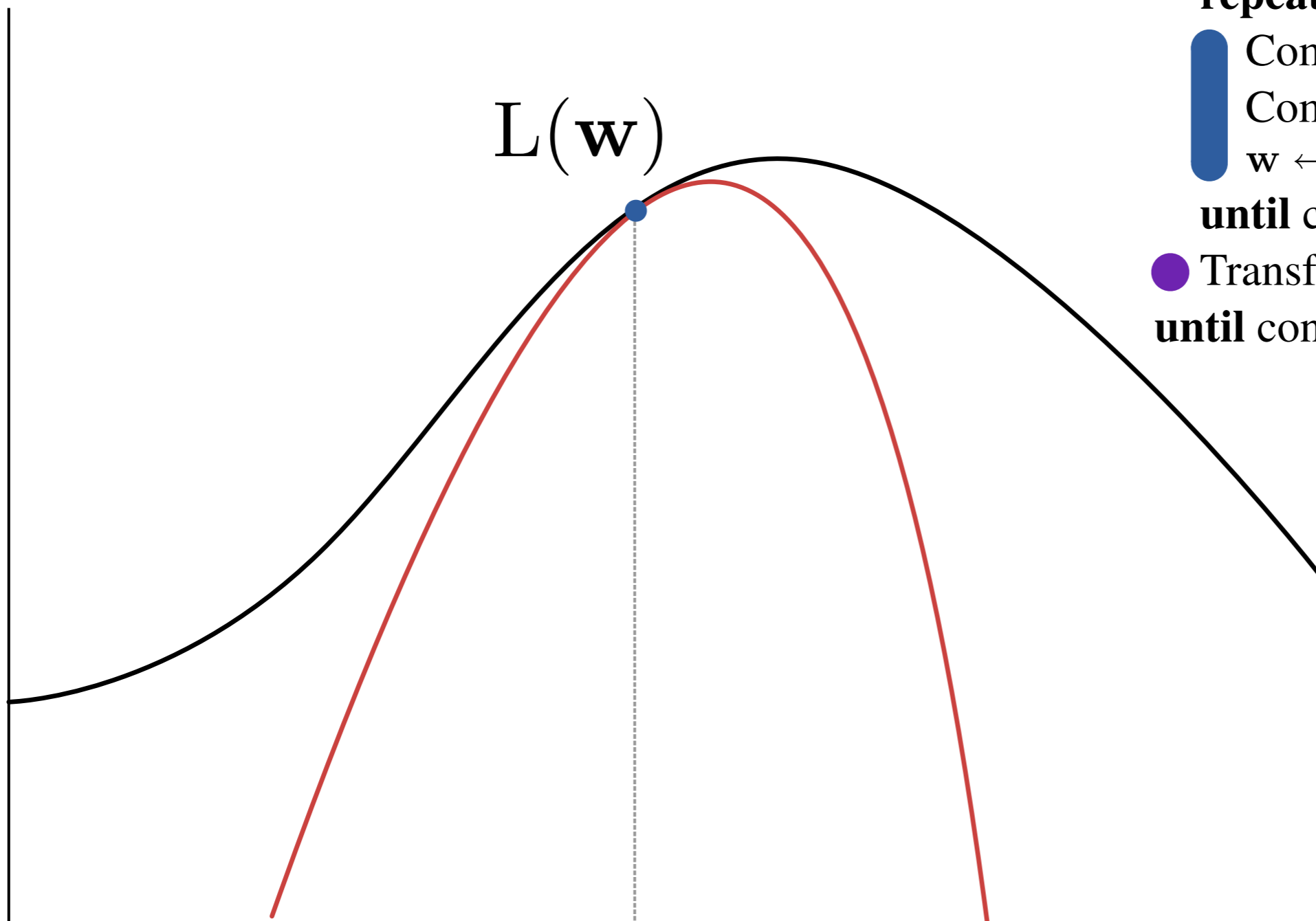
● Compute  $\nabla \ell(\mathbf{w}, \mathbf{e})$

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

**until convergence**

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

**until convergence**





# EM with Features

Initialize weights  $\mathbf{w}$

**repeat**

● Compute expected counts  $\mathbf{e}$

**repeat**

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

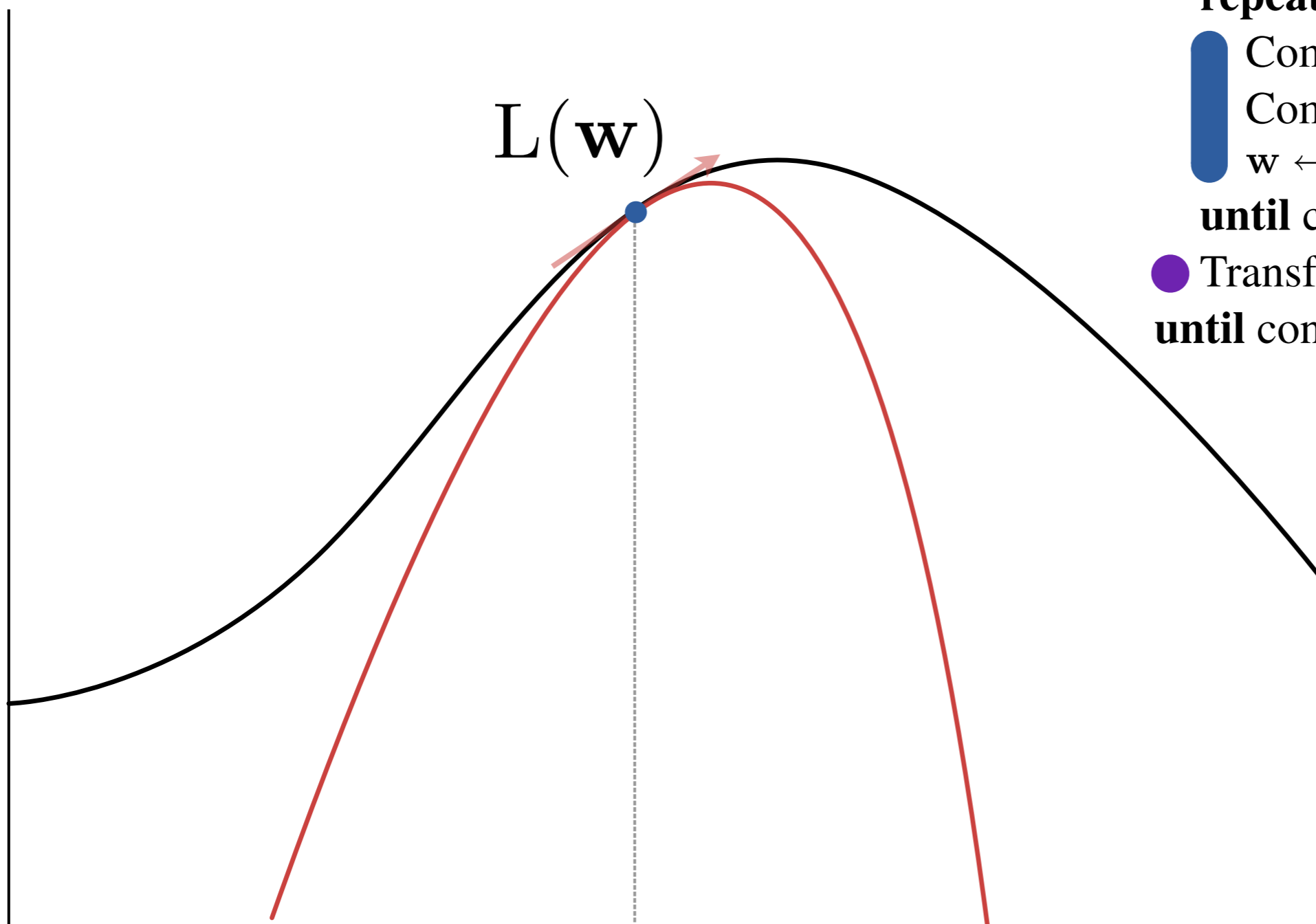
● Compute  $\nabla \ell(\mathbf{w}, \mathbf{e})$

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

**until convergence**

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

**until convergence**





# EM with Features

Initialize weights  $w$

**repeat**

● Compute expected counts  $e$

**repeat**

● Compute  $\ell(w, e)$

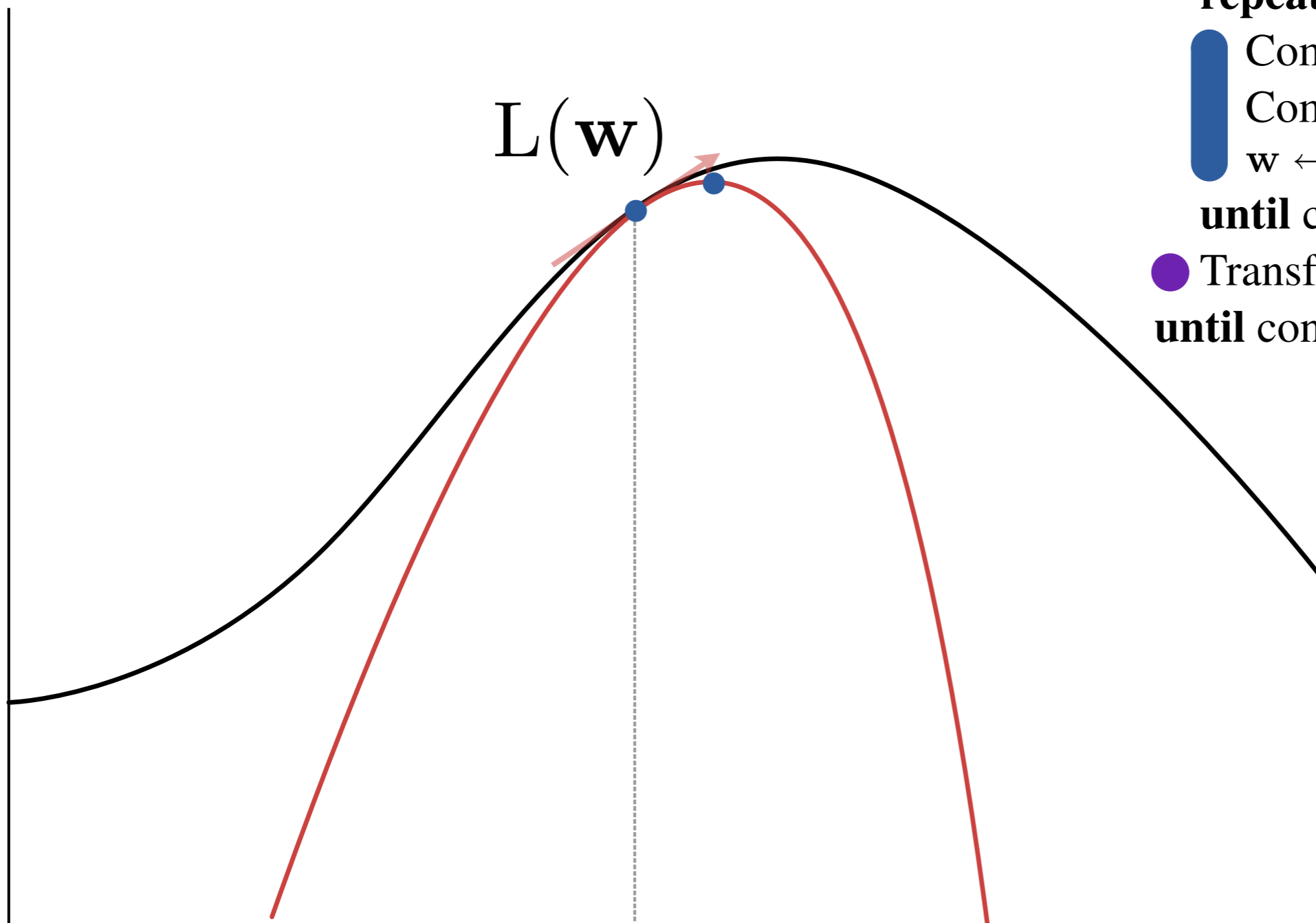
● Compute  $\nabla \ell(w, e)$

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

**until convergence**

● Transform  $w$  to  $\theta$

**until convergence**





# EM with Features

Initialize weights  $\mathbf{w}$

**repeat**

● Compute expected counts  $\mathbf{e}$

**repeat**

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

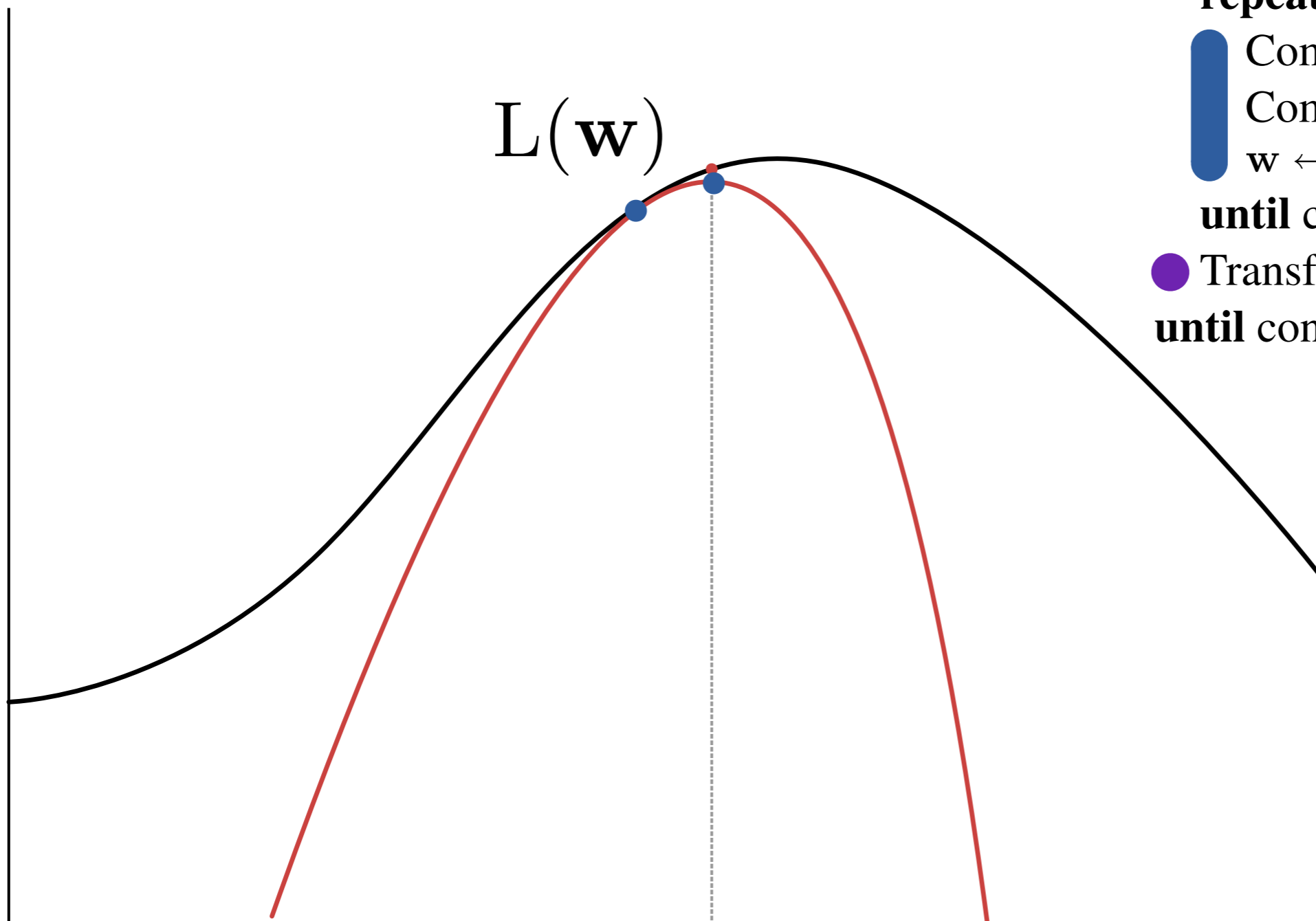
■ Compute  $\nabla \ell(\mathbf{w}, \mathbf{e})$

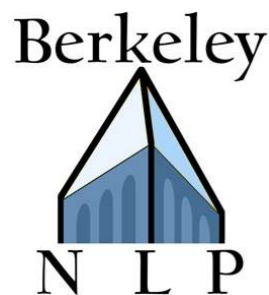
■  $\mathbf{w} \leftarrow \text{climb}(\mathbf{w}, \ell(\mathbf{w}, \mathbf{e}), \nabla \ell(\mathbf{w}, \mathbf{e}))$

**until convergence**

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

**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  $\mathbf{e}$

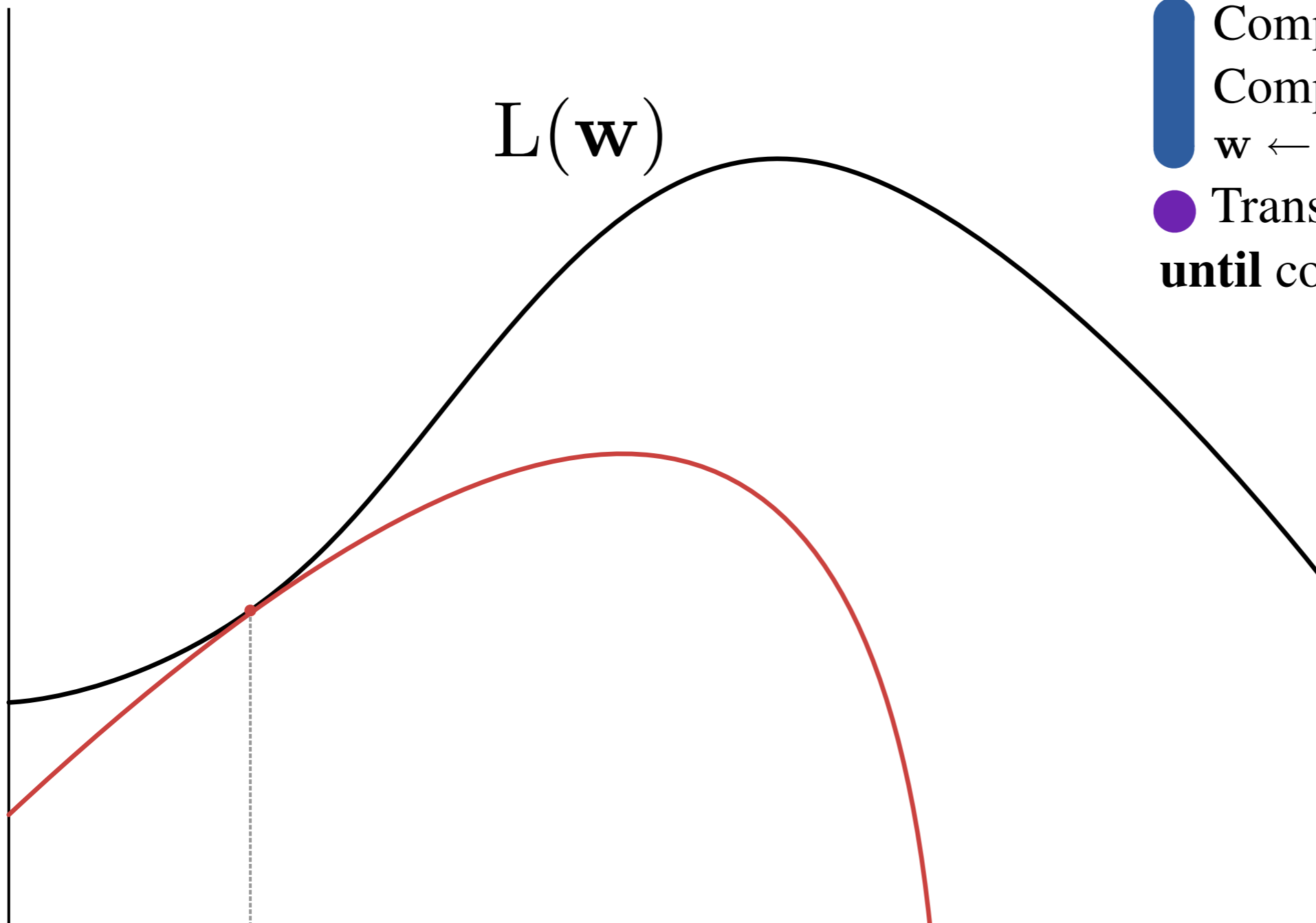
■ Compute  $L(\mathbf{w})$

■ 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  $\theta$

**until** convergence







# Direct Gradient with Features

Initialize weights  $\mathbf{w}$

**repeat**

● Compute expected counts  $\mathbf{e}$

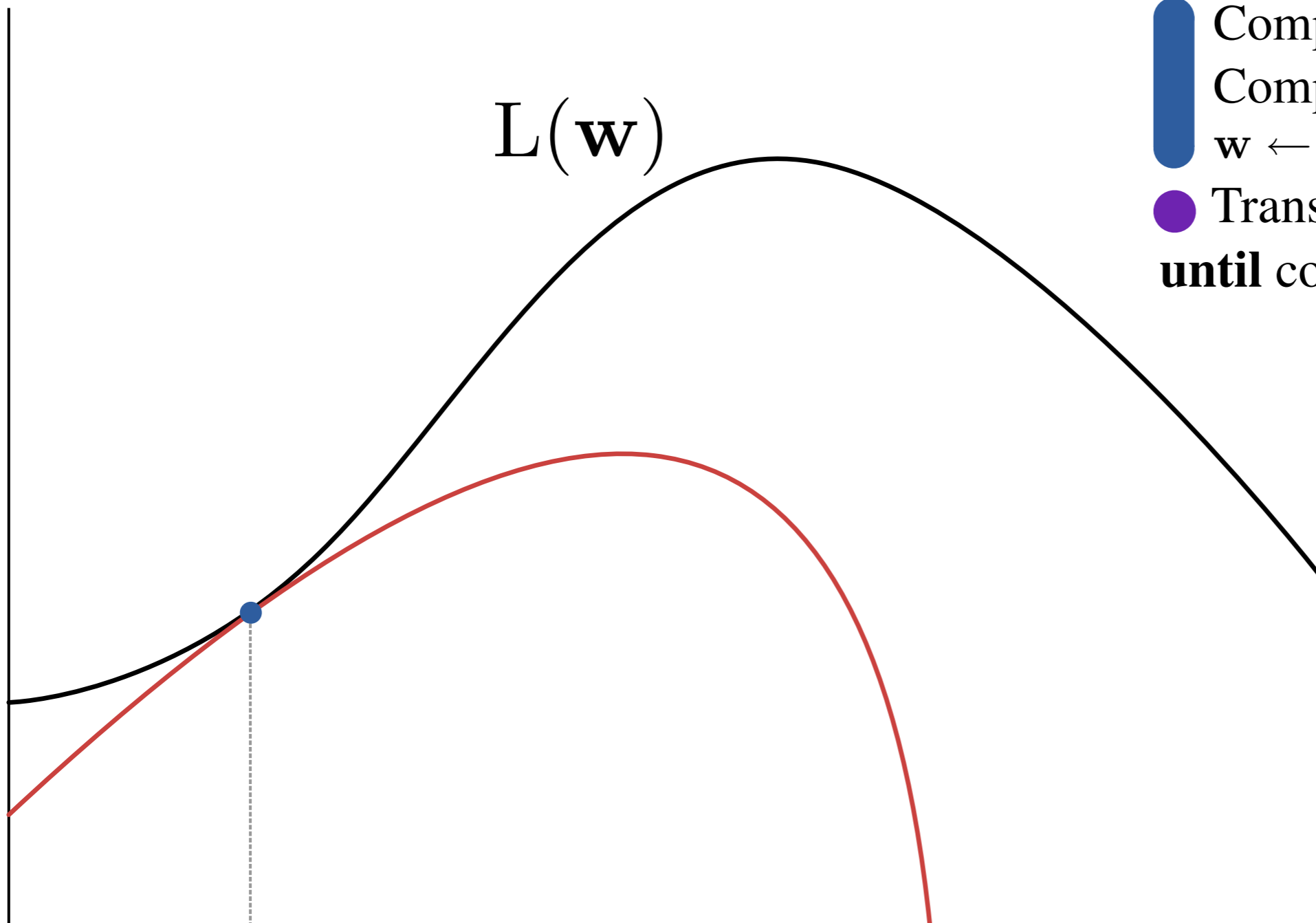
■ Compute  $L(\mathbf{w})$

■ 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  $\theta$

**until** convergence





# Direct Gradient with Features

Initialize weights  $\mathbf{w}$

**repeat**

● Compute expected counts  $\mathbf{e}$

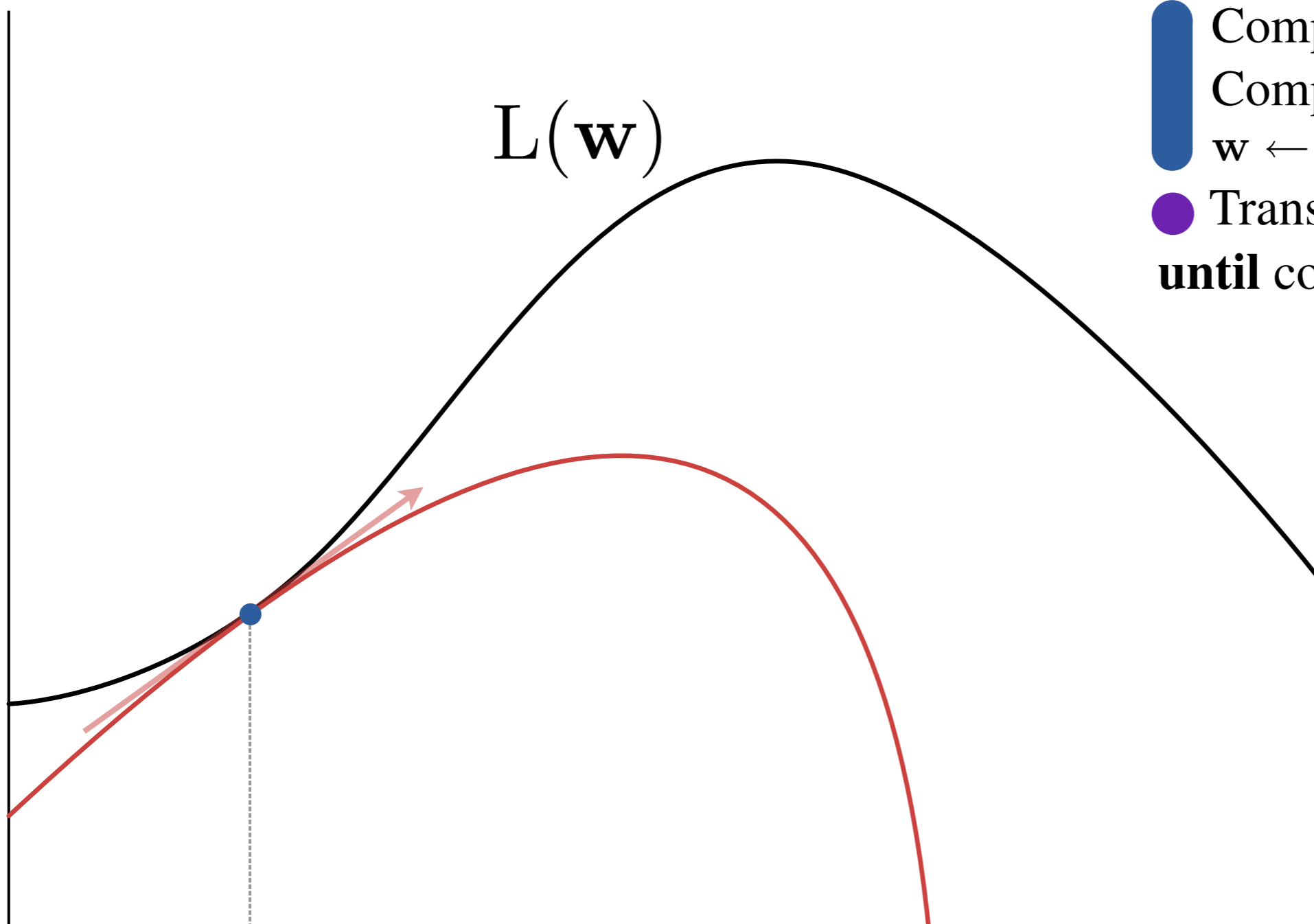
■ Compute  $L(\mathbf{w})$

■ 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  $\mathbf{e}$

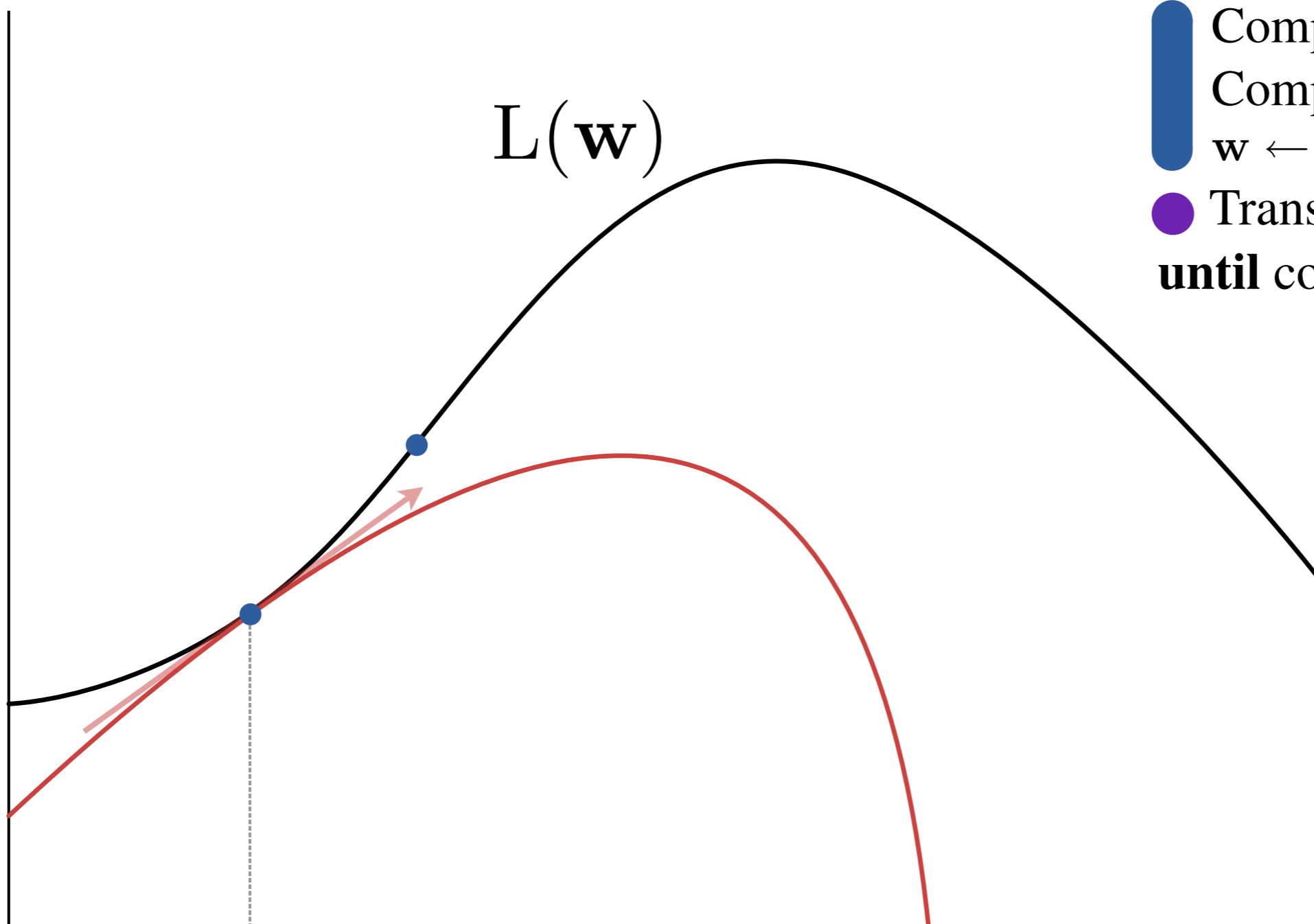
■ Compute  $L(\mathbf{w})$

■ 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  $\mathbf{e}$

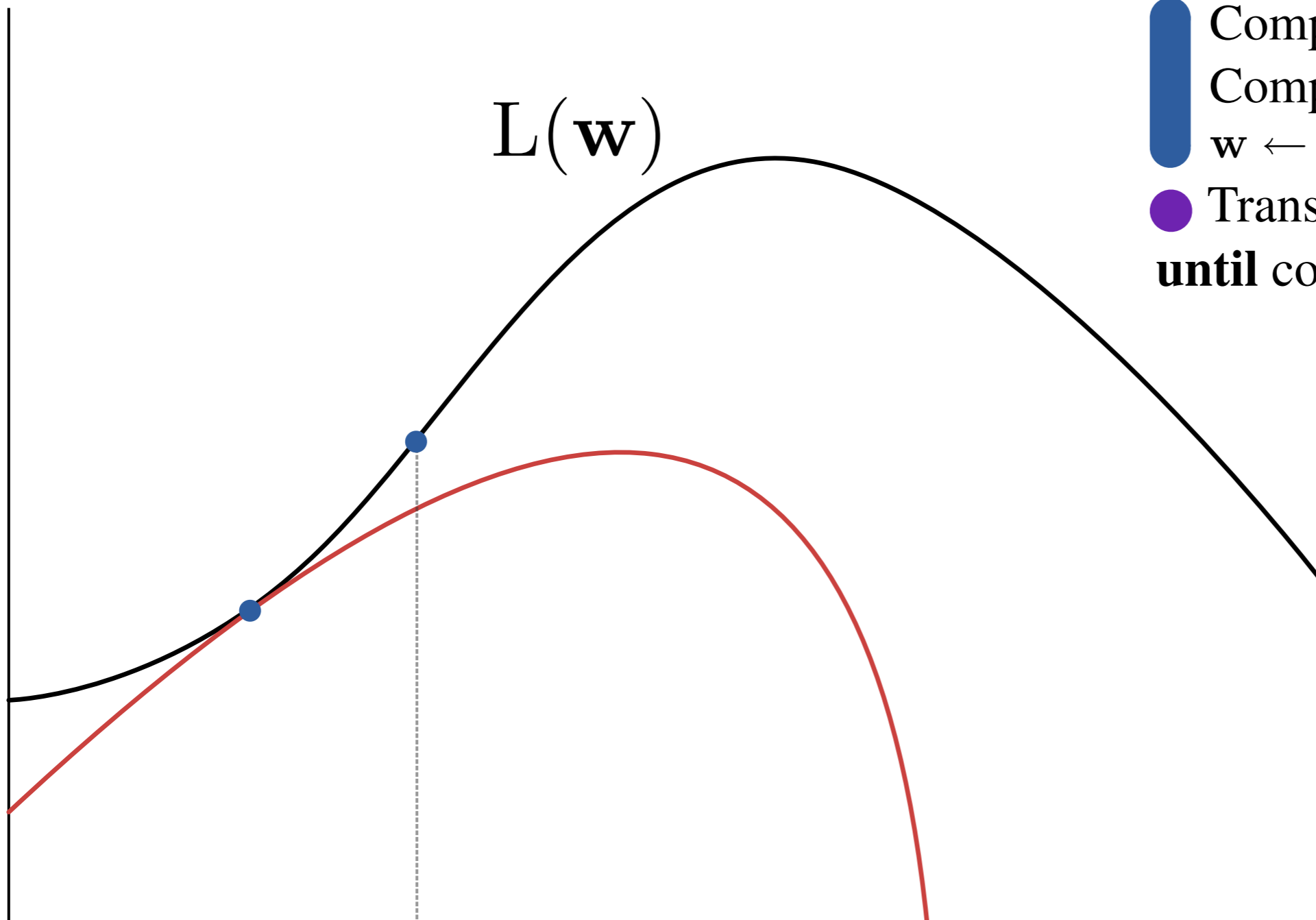
■ Compute  $L(\mathbf{w})$

■ 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  $\mathbf{e}$

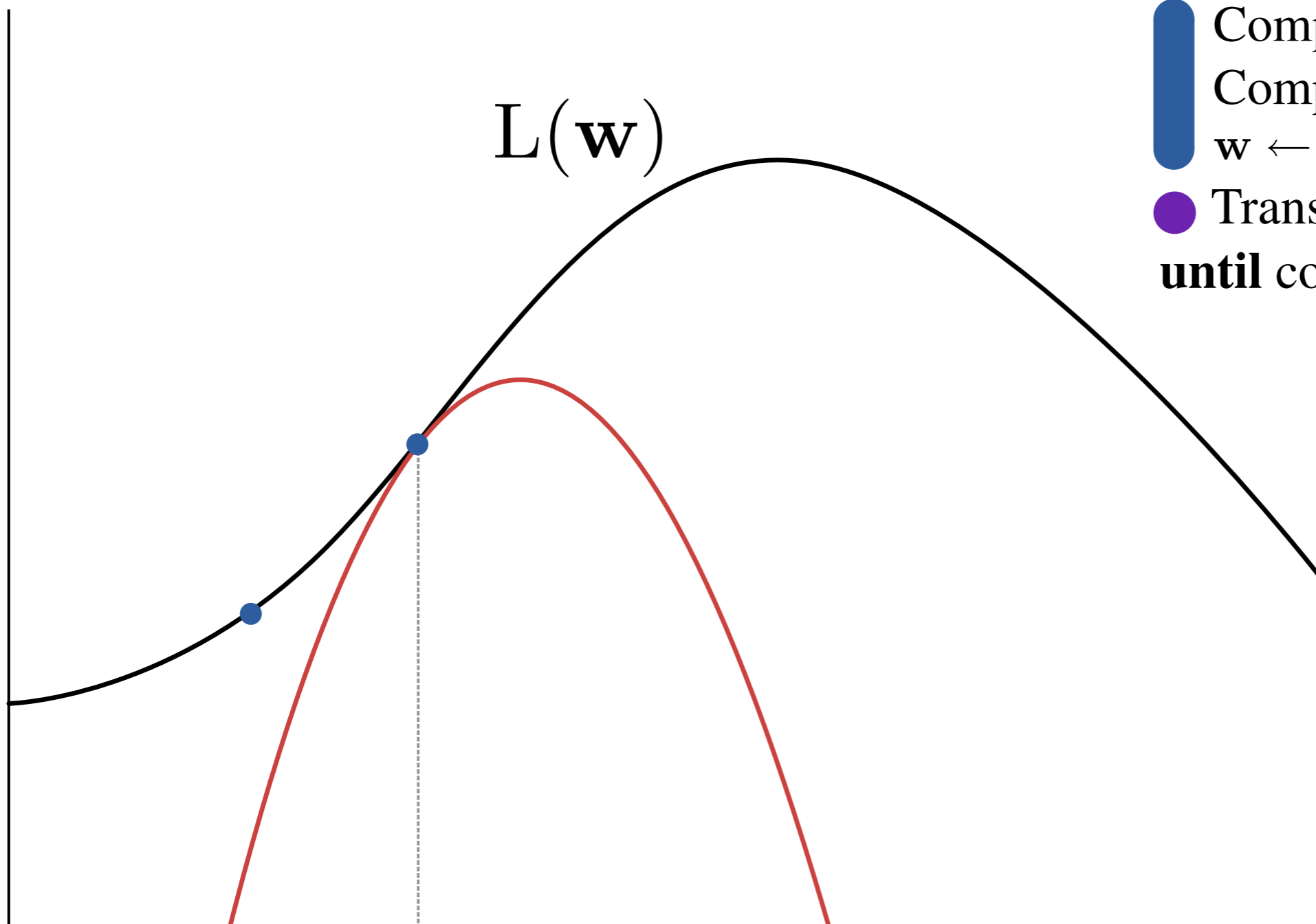
■ Compute  $L(\mathbf{w})$

■ 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





# Direct Gradient with Features

Initialize weights  $\mathbf{w}$

**repeat**

● Compute expected counts  $\mathbf{e}$

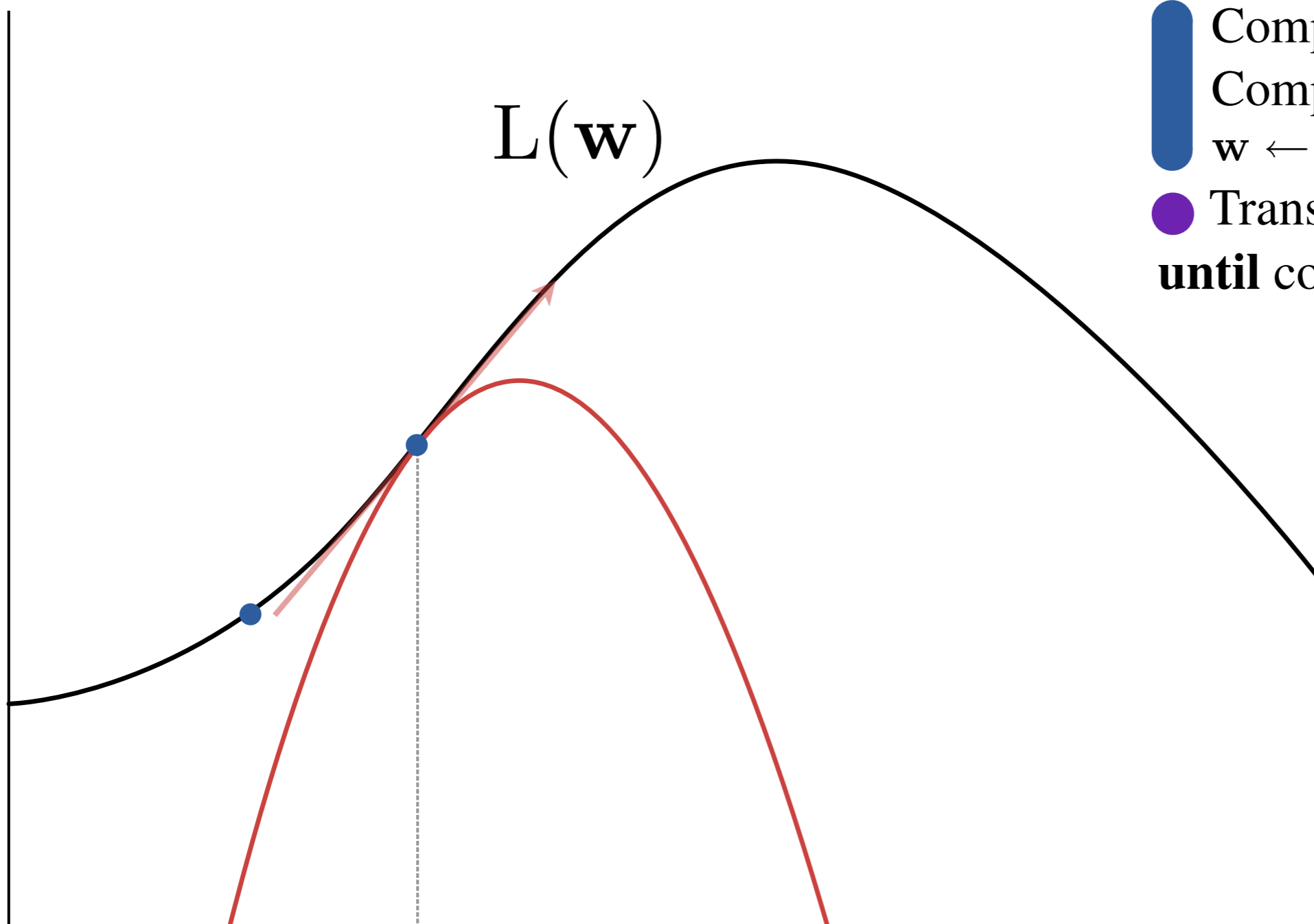
■ Compute  $L(\mathbf{w})$

■ 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  $\mathbf{e}$

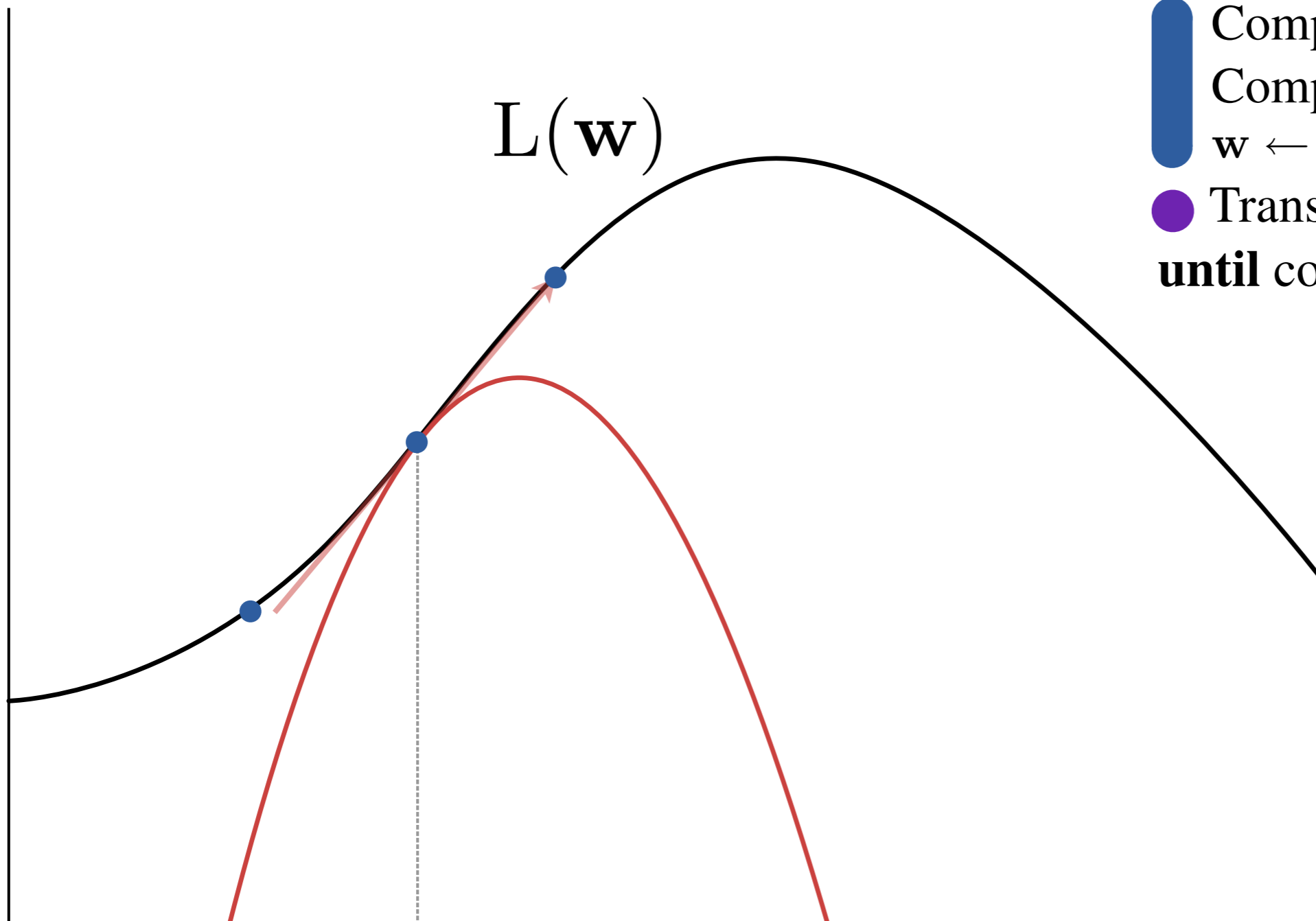
■ Compute  $L(\mathbf{w})$

■ 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  $\theta$

**until** convergence





# Direct Gradient with Features

Initialize weights  $\mathbf{w}$

**repeat**

● Compute expected counts  $\mathbf{e}$

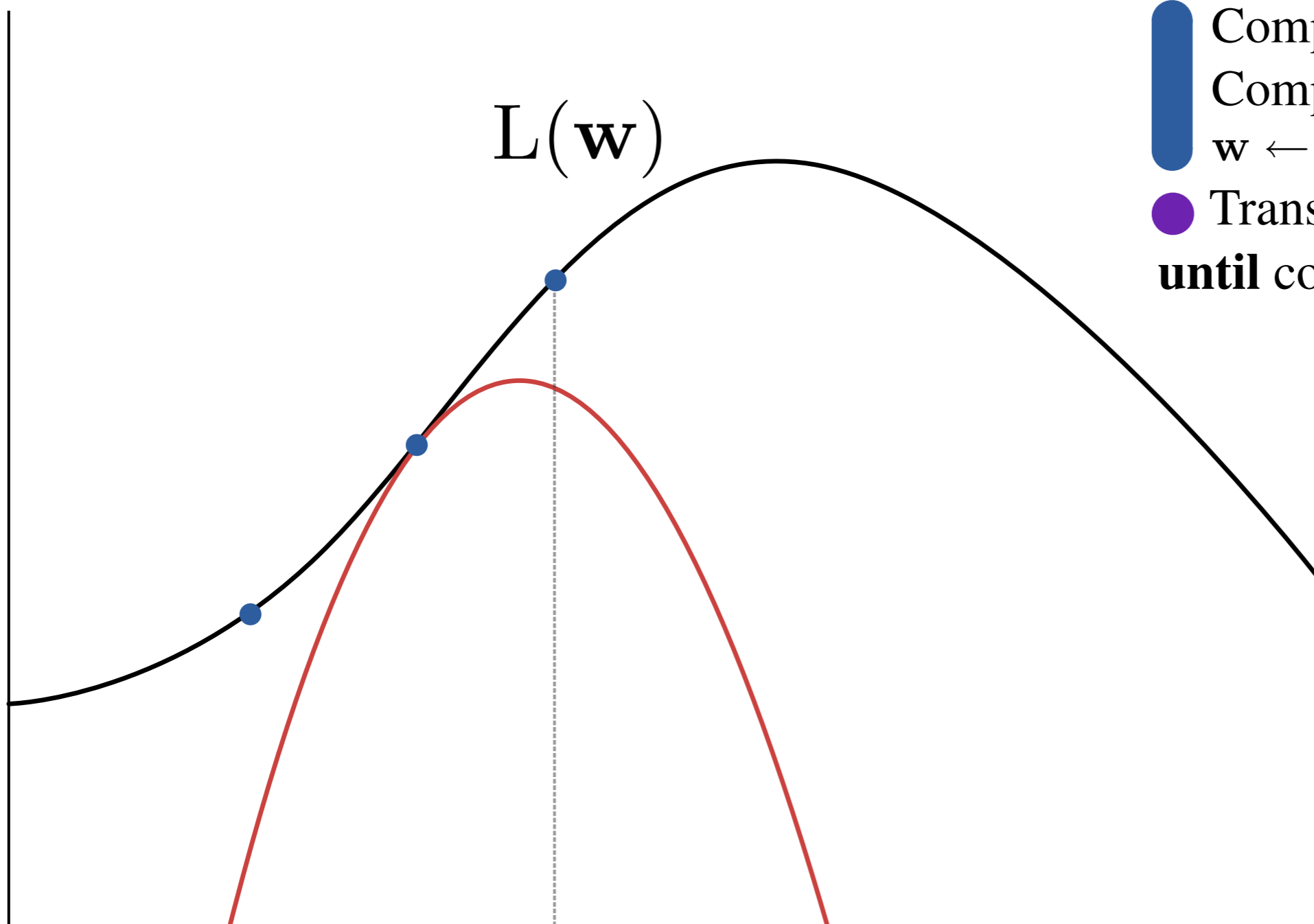
■ Compute  $L(\mathbf{w})$

■ 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  $\theta$

**until** convergence





# Unsupervised Induction Tasks

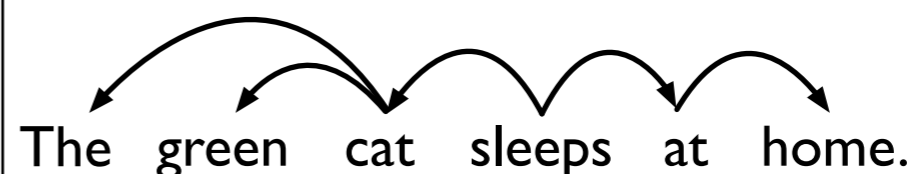
---

POS Induction:

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

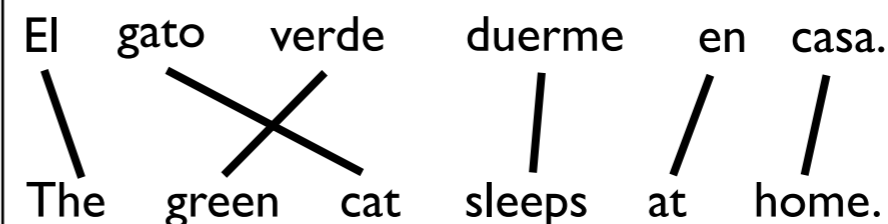
Grammar Induction:

The green cat sleeps at home.



Word Alignment:

El gato verde duerme en casa.  
 The green cat sleeps at home.



Word Segmentation:

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



# POS Induction Results

---

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



# POS Induction Results

---

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

Key distribution:  $P(\text{John}|\text{NN})$



# POS Induction Results

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

Key distribution:  $P(\text{John}|\text{NN})$

Features:

Basic:	$\text{John} \wedge \text{NN}$
Contains-Digit:	$+\text{Digit} \wedge \text{NN}$
Contains-Hyphen:	$+\text{Hyph} \wedge \text{NN}$
Initial-Capital:	$+\text{Cap} \wedge \text{NN}$
Suffix:	$+\text{ing} \wedge \text{NN}$



# POS Induction Results

---

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

Many-to-1 Accuracy

## Features:

Basic:	John $\wedge$ NNP
Contains-Digit:	+Digit $\wedge$ NNP
Contains-Hyphen:	+Hyph $\wedge$ NNP
Initial-Capital:	+Cap $\wedge$ NNP
Suffix:	+ing $\wedge$ NNP

## Data:

Train and test on entire WSJ

No tagging dictionary

45 POS tags

---



# POS Induction Results

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

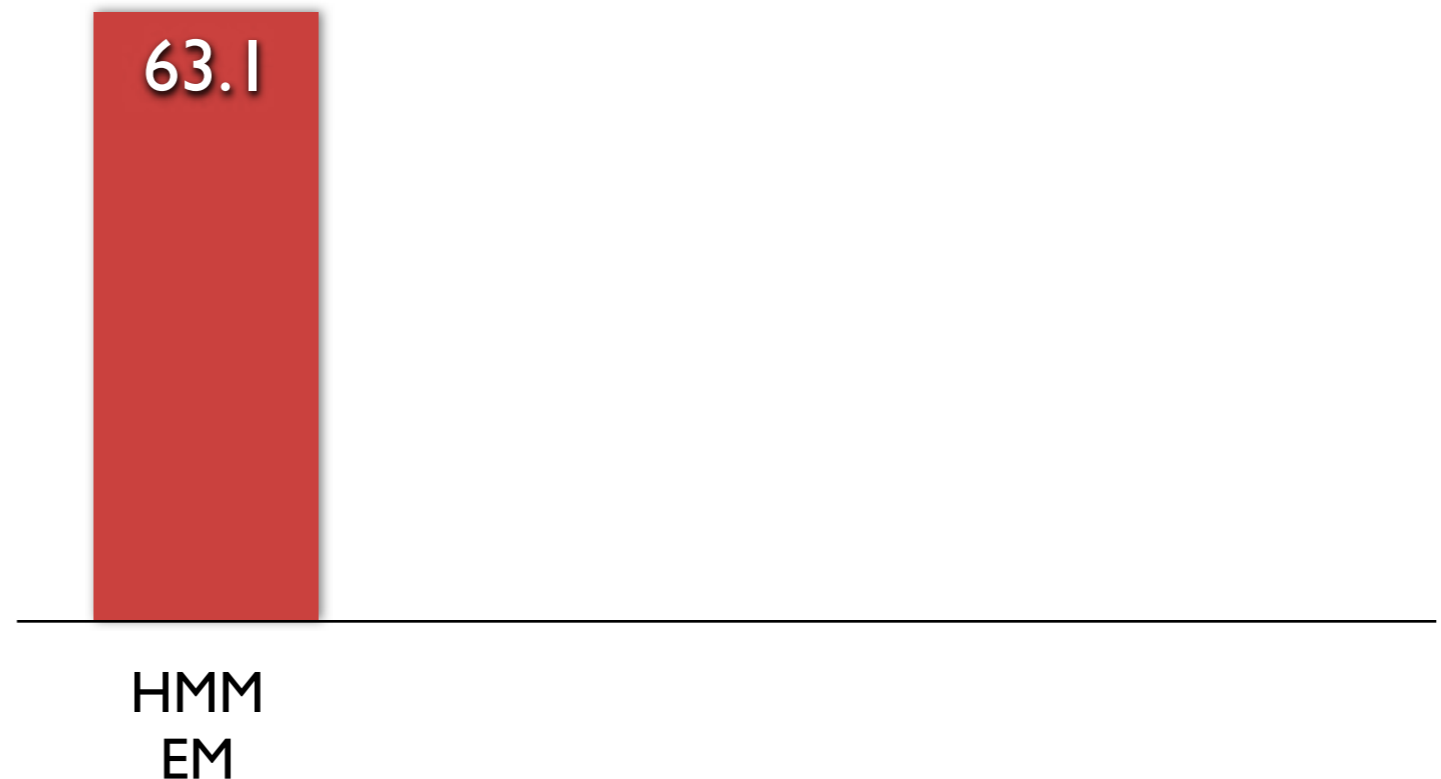
Many-to-1 Accuracy

## Features:

Basic:	John $\wedge$ NNP
Contains-Digit:	+Digit $\wedge$ NNP
Contains-Hyphen:	+Hyph $\wedge$ NNP
Initial-Capital:	+Cap $\wedge$ NNP
Suffix:	+ing $\wedge$ NNP

## Data:

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





# POS Induction Results

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

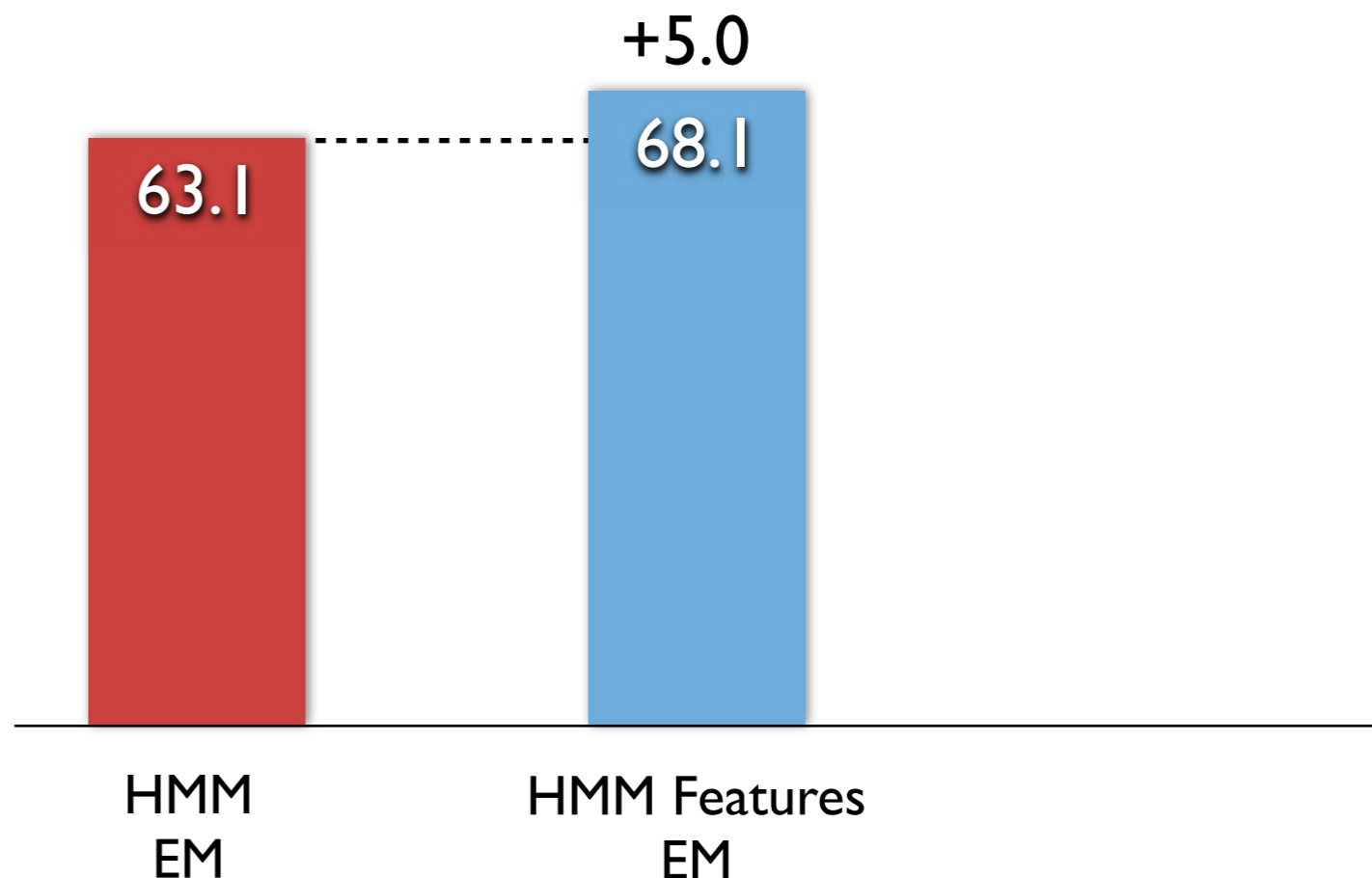
## Features:

Basic:	John $\wedge$ NNP
Contains-Digit:	+Digit $\wedge$ NNP
Contains-Hyphen:	+Hyph $\wedge$ NNP
Initial-Capital:	+Cap $\wedge$ NNP
Suffix:	+ing $\wedge$ NNP

## 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.

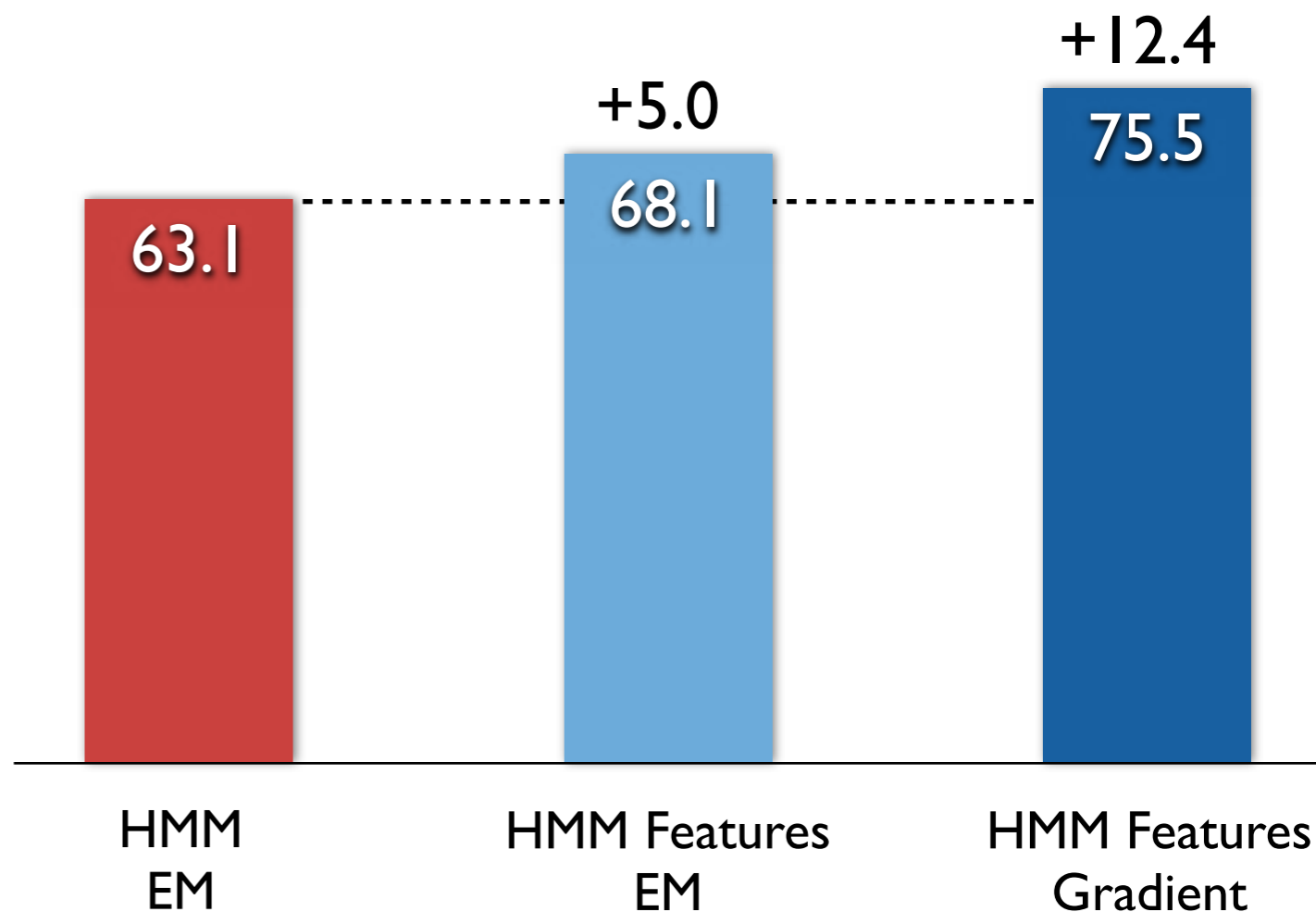
## Features:

Basic:	John $\wedge$ NNP
Contains-Digit:	+Digit $\wedge$ NNP
Contains-Hyphen:	+Hyph $\wedge$ NNP
Initial-Capital:	+Cap $\wedge$ NNP
Suffix:	+ing $\wedge$ NNP

## 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.

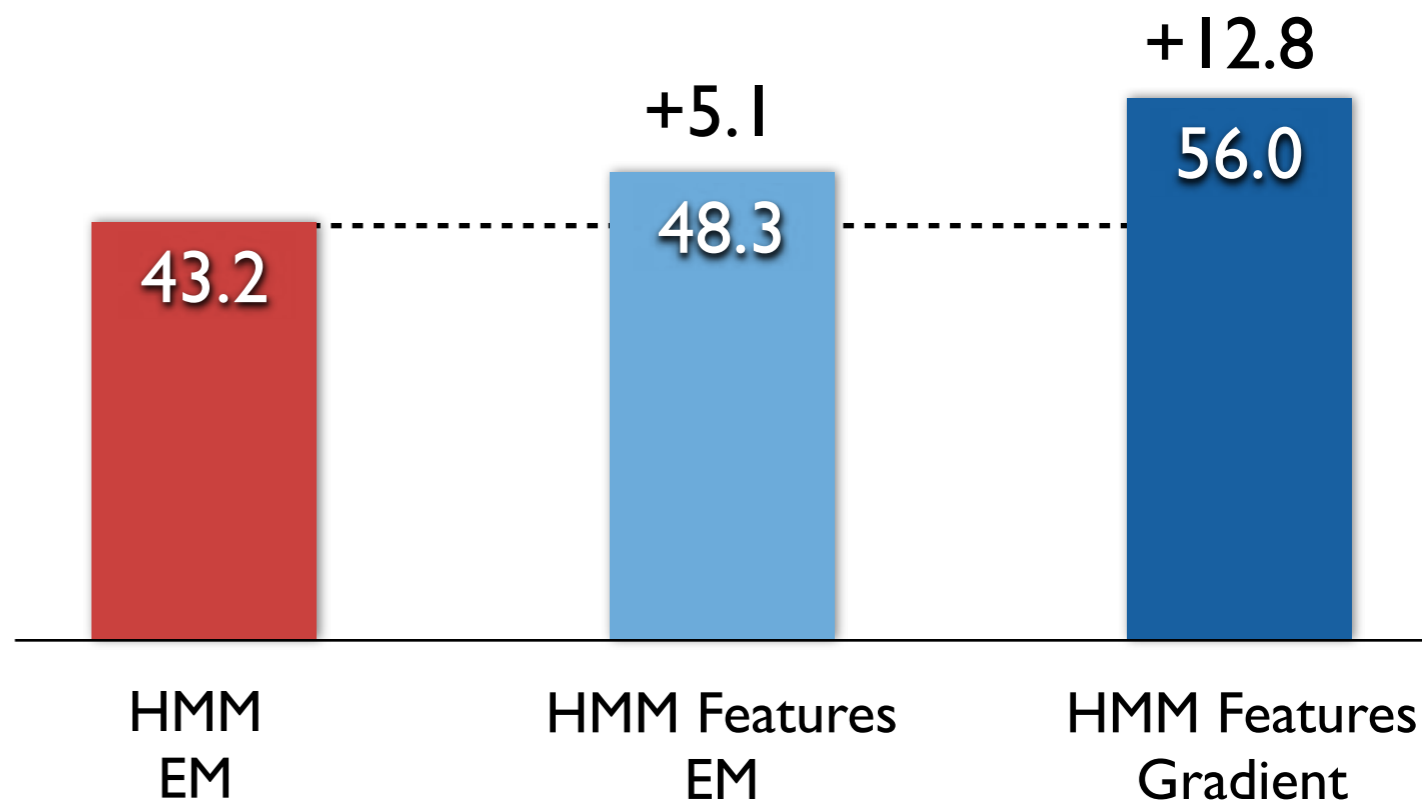
## Features:

Basic:	John $\wedge$ NNP
Contains-Digit:	+Digit $\wedge$ NNP
Contains-Hyphen:	+Hyph $\wedge$ NNP
Initial-Capital:	+Cap $\wedge$ NNP
Suffix:	+ing $\wedge$ NNP

## Data:

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

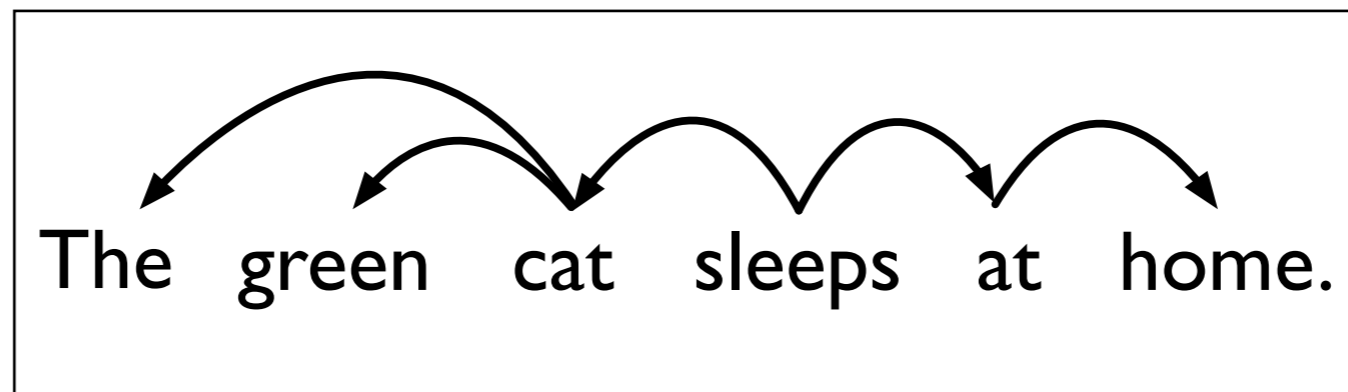
## I-to-I Accuracy





# Grammar Induction Results

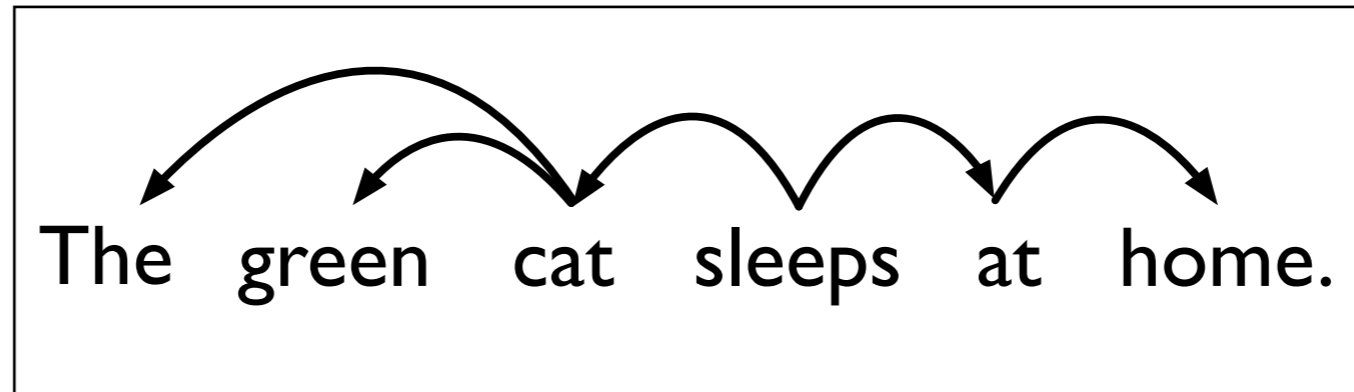
---





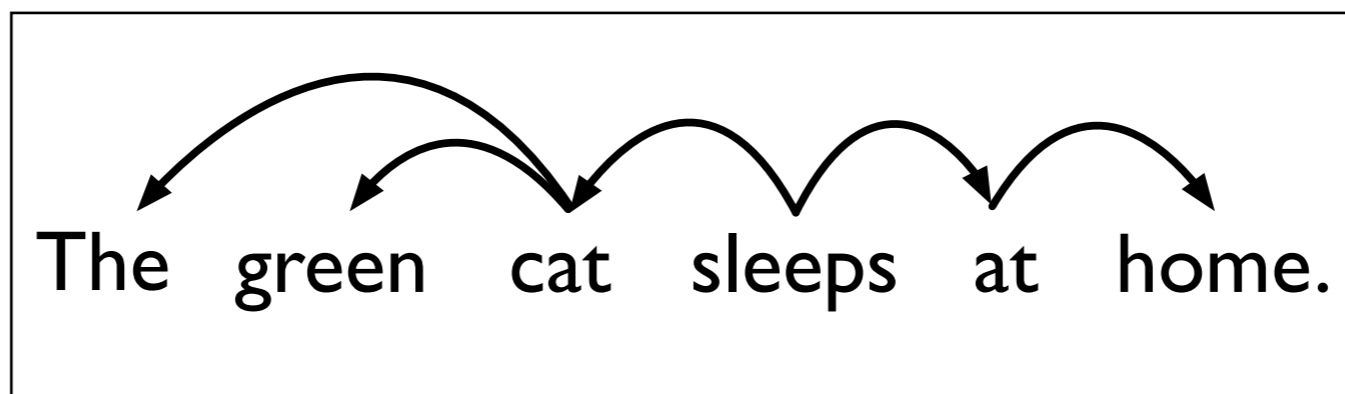
# Grammar Induction Results

---



Key distributions:  $P(\text{JJ}|\text{NN})$   $P(\text{stop}|\text{NN})$

# Grammar Induction Results



Key distributions:  $P(\text{JJ}|\text{NN})$   $P(\text{stop}|\text{NN})$

Features:

Basic:  $\text{JJ} \wedge \text{NN}, \text{JJ} \wedge \text{NNS}$

Noun:  $\text{JJ} \wedge \text{Noun}$

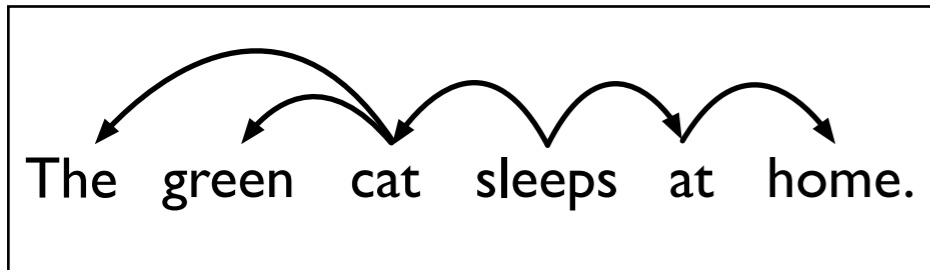
Verb:  $\text{JJ} \wedge \text{Verb}$

Noun-verb:  $\text{JJ} \wedge \text{NounOrVerb}$

# Grammar Induction Results

---

English Directed Accuracy



## Features:

Basic: JJ  $\wedge$  NN, JJ  $\wedge$  NNS

Noun: JJ  $\wedge$  Noun

Verb: JJ  $\wedge$  Verb

Noun-verb: JJ  $\wedge$  NounOrVerb

---

Chinese Directed Accuracy

## Data:

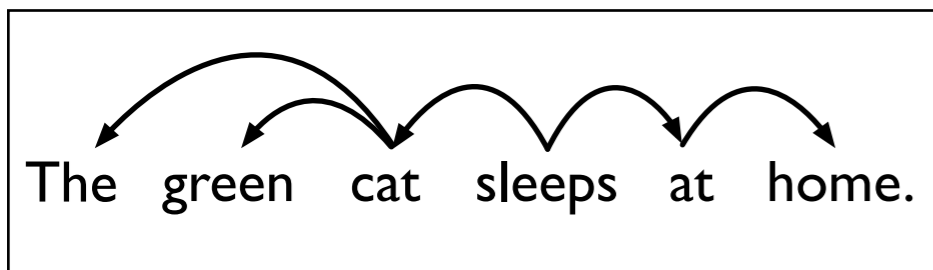
Train WSJ10 Sec. 2-21  
CTB10 Sec. 1-270

Tune WSJ10 Sec. 22  
CTB10 Sec. 400-454

Test WSJ10 Sec. 23  
CTB10 Sec. 271-300

---

# Grammar Induction Results



## Features:

Basic: JJ  $\wedge$  NN, JJ  $\wedge$  NNS

Noun: JJ  $\wedge$  Noun

Verb: JJ  $\wedge$  Verb

Noun-verb: JJ  $\wedge$  NounOrVerb

## 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

English Directed Accuracy

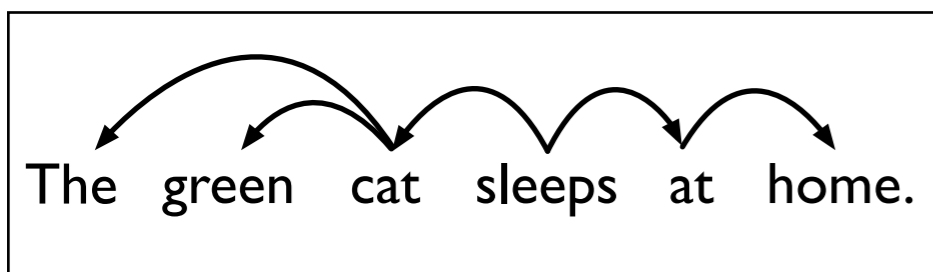
47.8

Chinese Directed Accuracy

42.5

DMV  
EM

# Grammar Induction Results



## Features:

Basic: JJ  $\wedge$  NN, JJ  $\wedge$  NNS

Noun: JJ  $\wedge$  Noun

Verb: JJ  $\wedge$  Verb

Noun-verb: JJ  $\wedge$  NounOrVerb

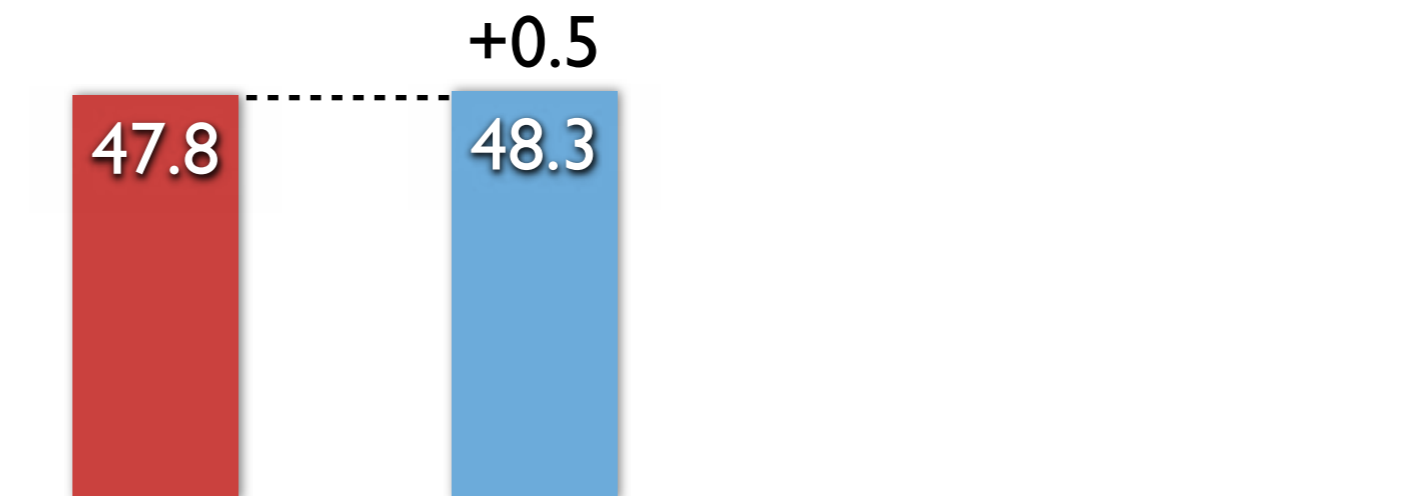
## Data:

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

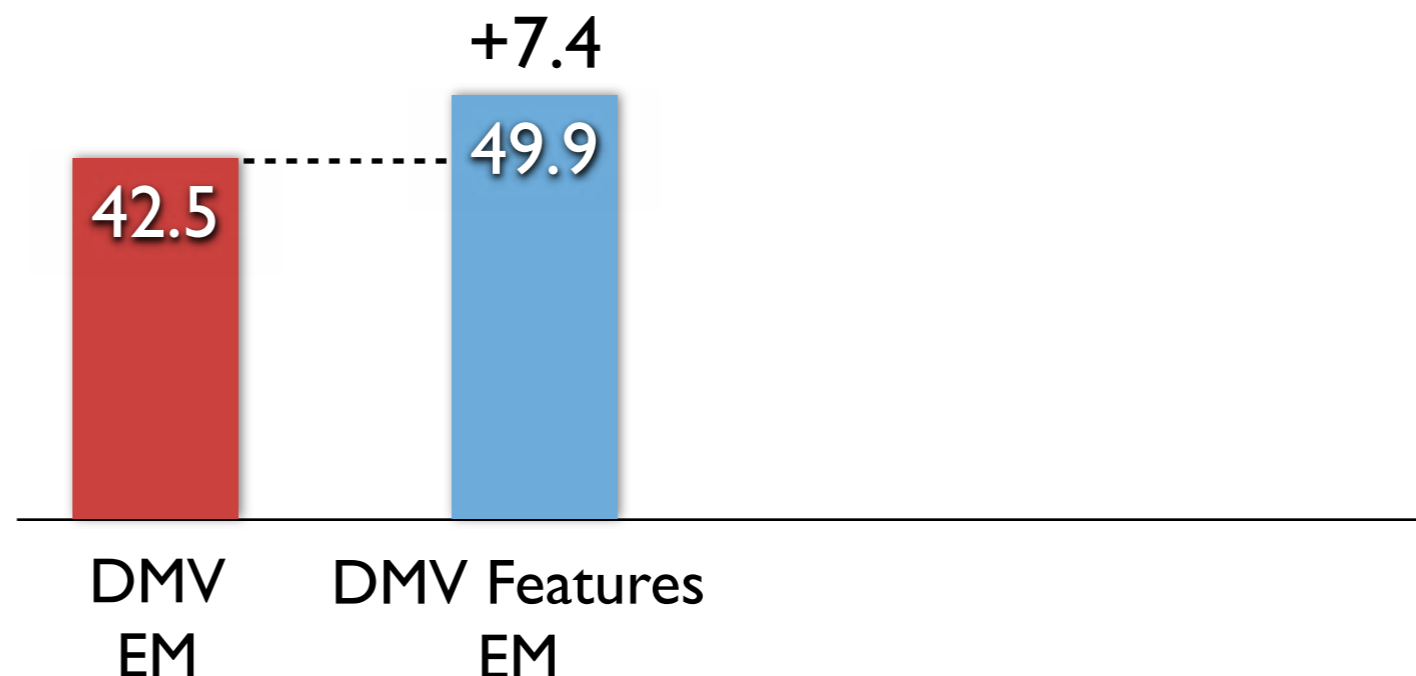
Tune WSJ10 Sec. 22  
 CTB10 Sec. 400-454

Test WSJ10 Sec. 23  
 CTB10 Sec. 271-300

English Directed Accuracy

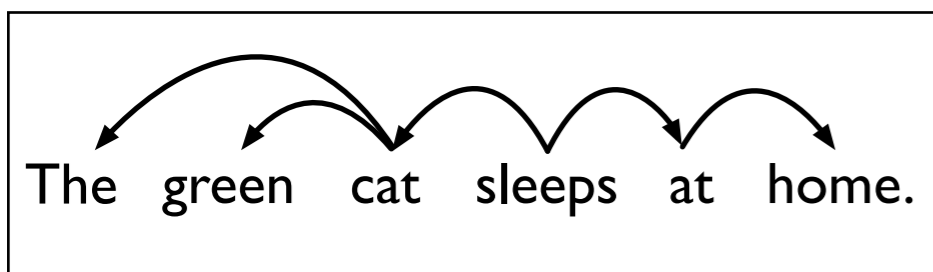


Chinese Directed Accuracy





# Grammar Induction Results



## Features:

Basic: JJ  $\wedge$  NN, JJ  $\wedge$  NNS

Noun: JJ  $\wedge$  Noun

Verb: JJ  $\wedge$  Verb

Noun-verb: JJ  $\wedge$  NounOrVerb

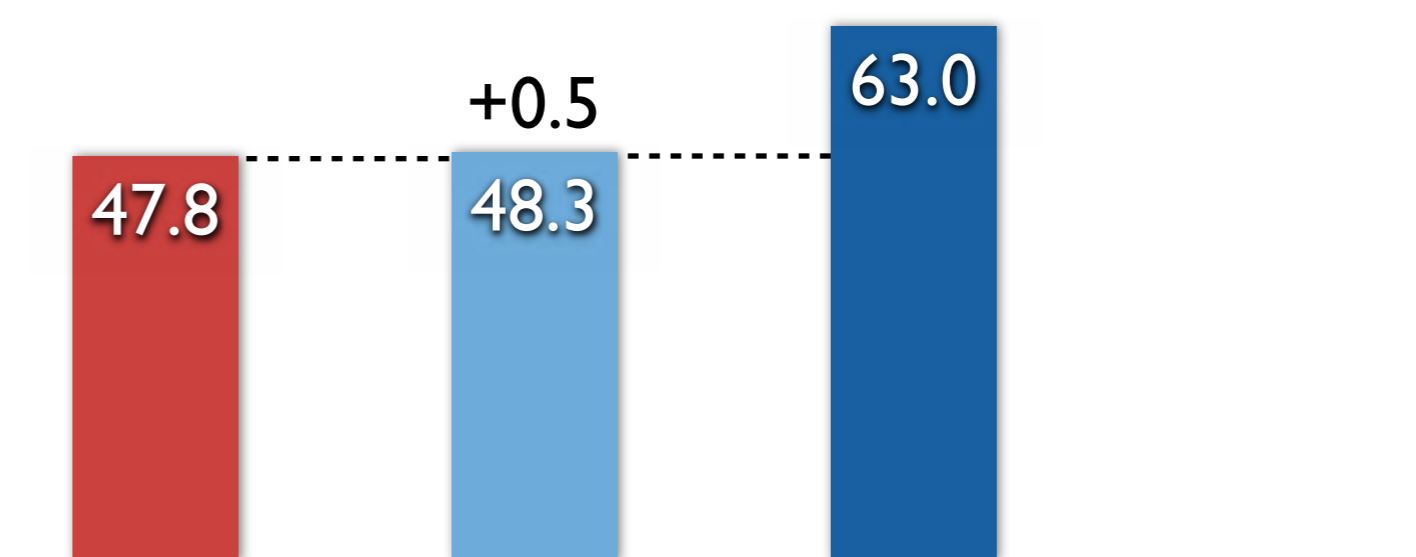
## Data:

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

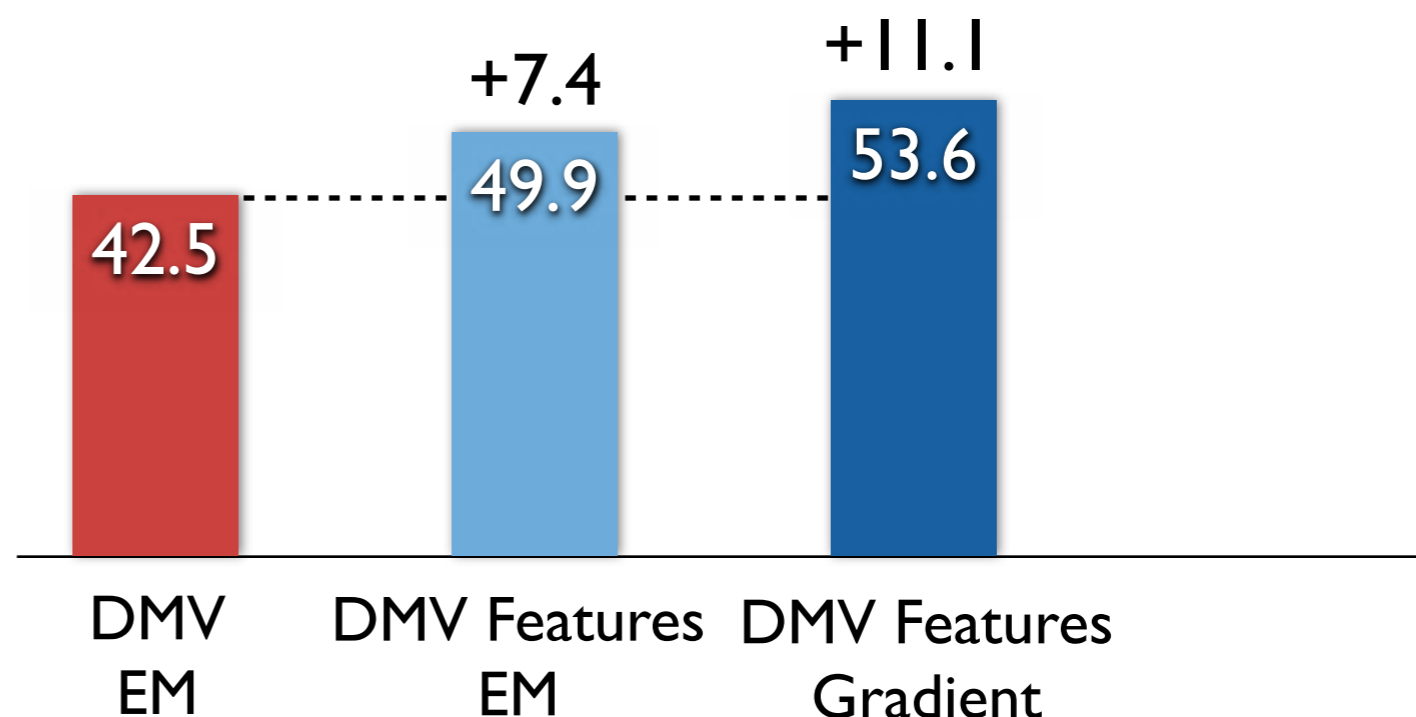
Tune WSJ10 Sec. 22  
CTB10 Sec. 400-454

Test WSJ10 Sec. 23  
CTB10 Sec. 271-300

English Directed Accuracy  
+15.2

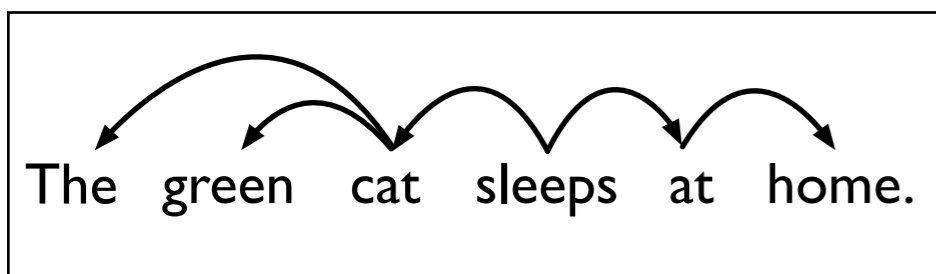


Chinese Directed Accuracy  
+11.1





# Grammar Induction Results



## Features:

Basic: JJ  $\wedge$  NN, JJ  $\wedge$  NNS

Noun: JJ  $\wedge$  Noun

Verb: JJ  $\wedge$  Verb

Noun-verb: JJ  $\wedge$  NounOrVerb

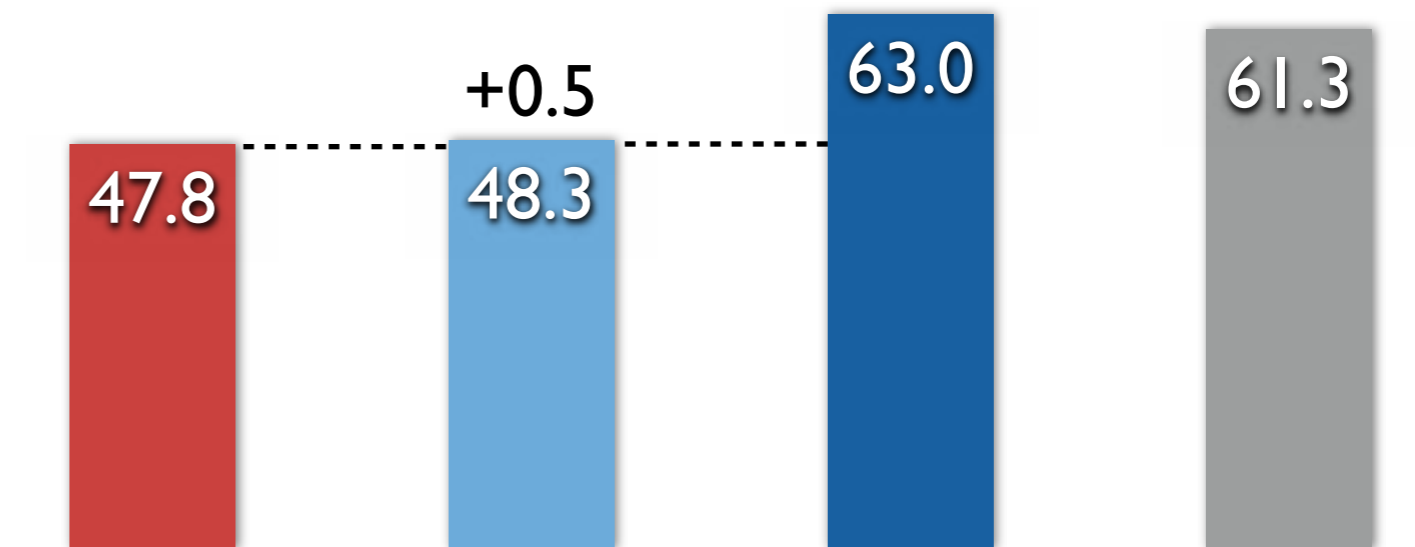
## Data:

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

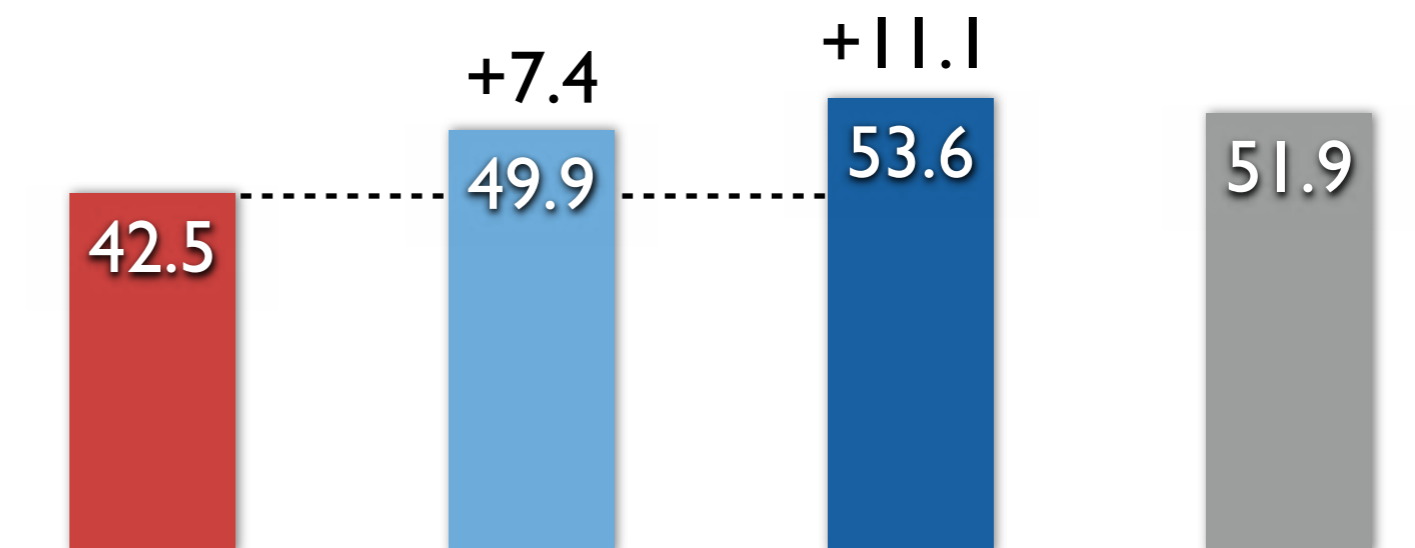
Tune WSJ10 Sec. 22  
CTB10 Sec. 400-454

Test WSJ10 Sec. 23  
CTB10 Sec. 271-300

English Directed Accuracy  
+15.2



Chinese Directed Accuracy  
+11.1



DMV  
EM

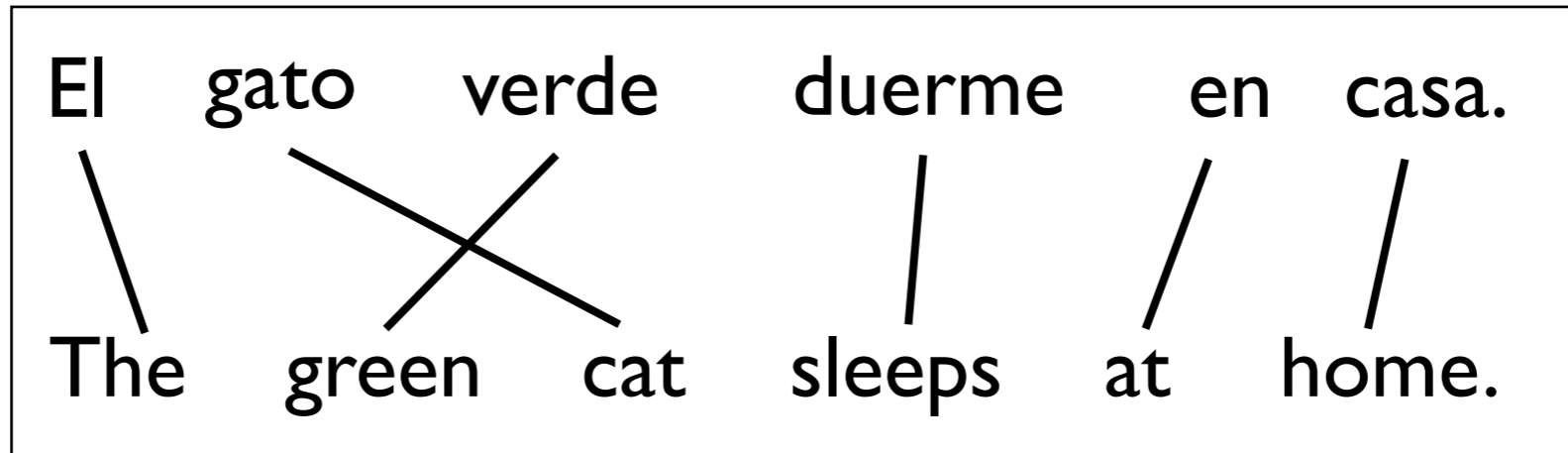
DMV Features  
EM

DMV Features  
Gradient

Cohen and  
Smith '09  
SLN DMV

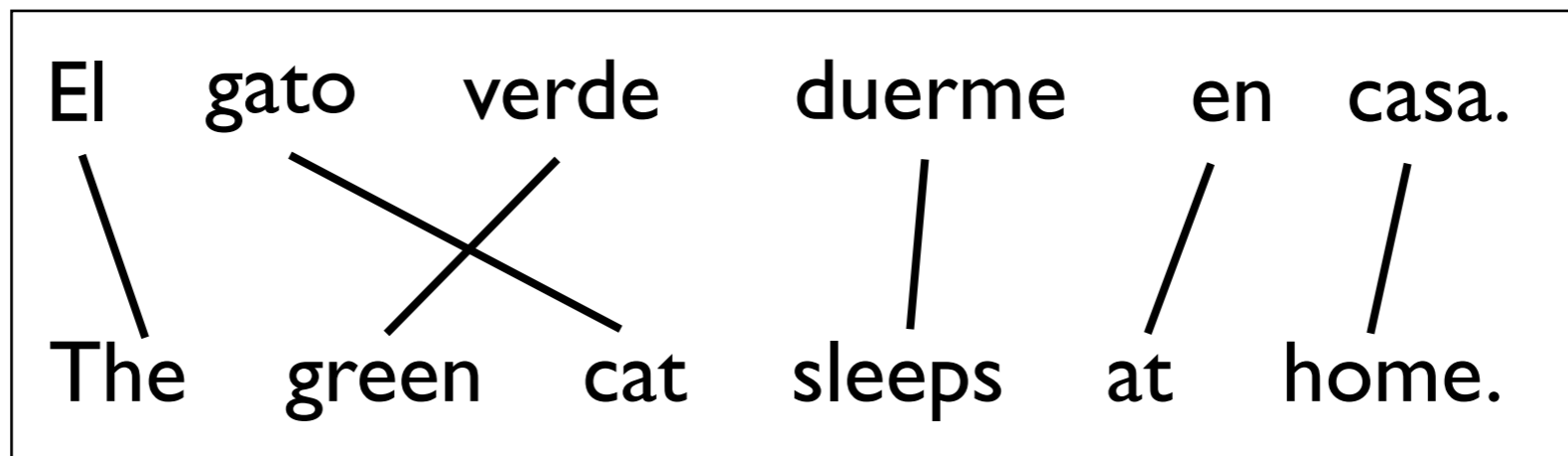


# Word Alignment Results



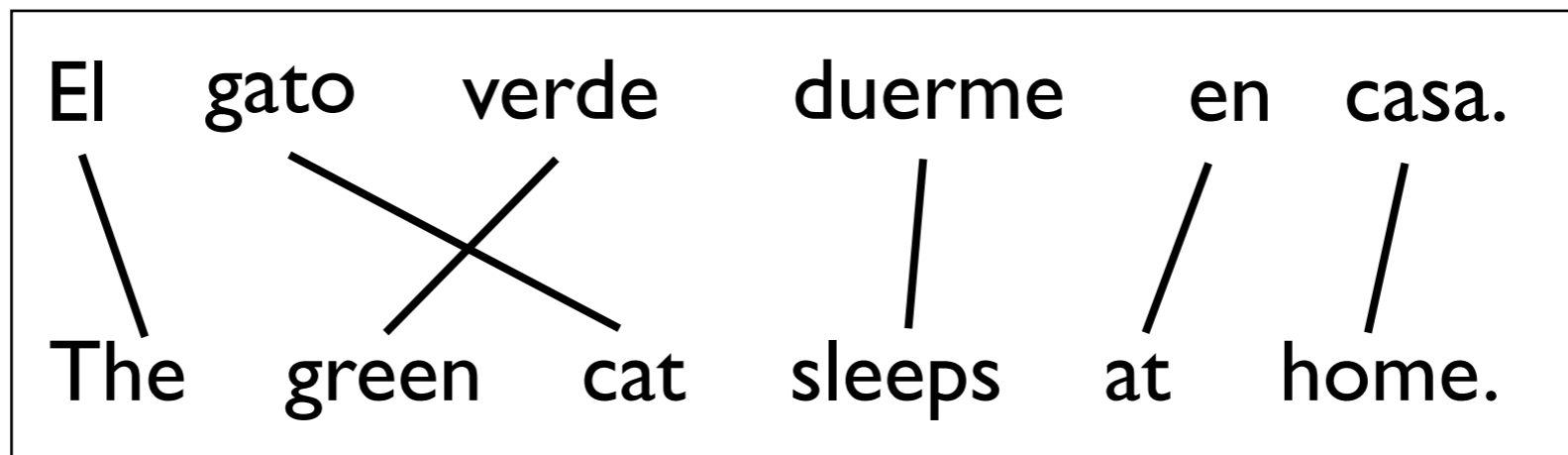


# Word Alignment Results



Key distribution:  $P(\text{gato}|\text{cat})$

# Word Alignment Results



Key distribution:  $P(\text{gato}|\text{cat})$

Features:

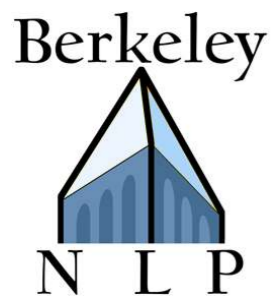
Basic: gato  $\wedge$  cat

Edit-Distance:  $\text{edit}(\text{gato}, \text{cat}) = 2$

Dictionary:  $(\text{gato}, \text{cat}) \in \text{Dict}$

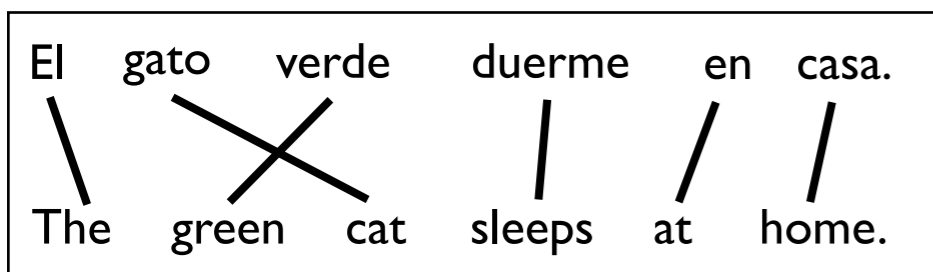
Stem: gato  $\wedge$  +stem(cat)

Prefix: gato  $\wedge$  +ca



# Word Alignment Results

---



Alignment Error Rate

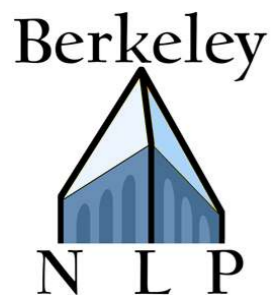
## Features:

- Basic:  $\text{gato} \wedge \text{cat}$
- Edit-Distance:  $\text{edit}(\text{gato}, \text{cat}) = 2$
- Dictionary:  $(\text{gato}, \text{cat}) \in \text{Dict}$
- Stem:  $\text{gato} \wedge +\text{stem}(\text{cat})$
- Prefix:  $\text{gato} \wedge +\text{ca}$

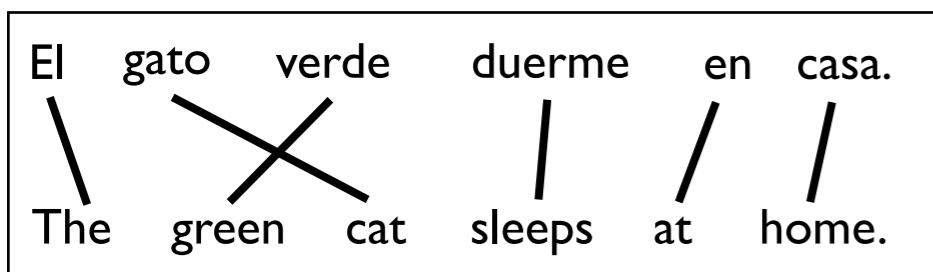
## Data:

Train 10K sentences of FBIS  
Chinese-English newswire

Test NIST 2002 Chinese-English dev set



# Word Alignment Results



Alignment Error Rate

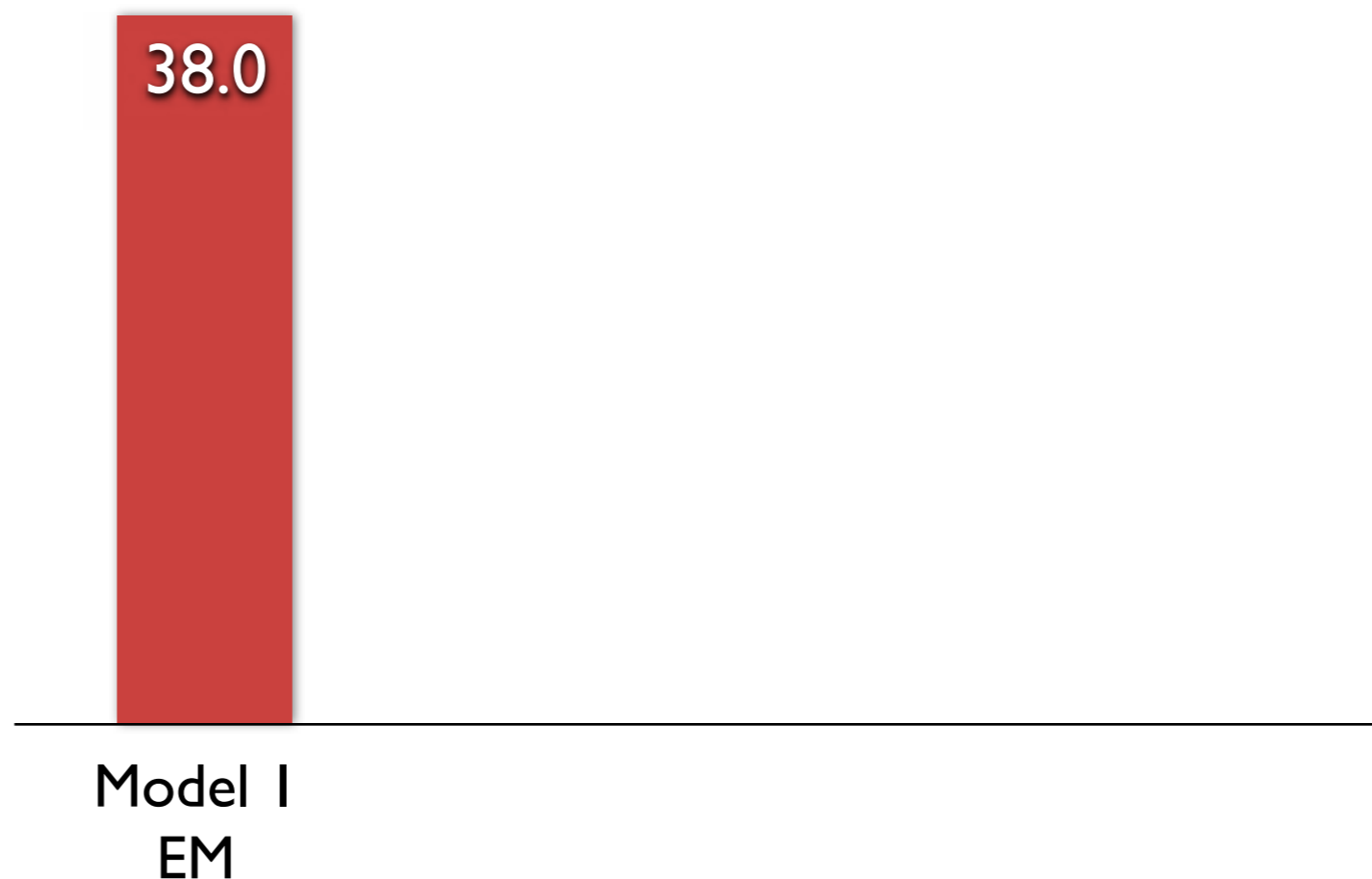
## Features:

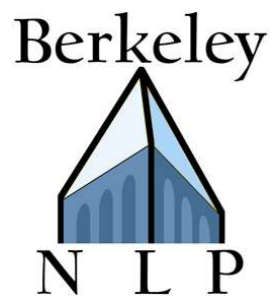
- Basic:  $\text{gato} \wedge \text{cat}$
- Edit-Distance:  $\text{edit}(\text{gato}, \text{cat}) = 2$
- Dictionary:  $(\text{gato}, \text{cat}) \in \text{Dict}$
- Stem:  $\text{gato} \wedge +\text{stem}(\text{cat})$
- Prefix:  $\text{gato} \wedge +\text{ca}$

## Data:

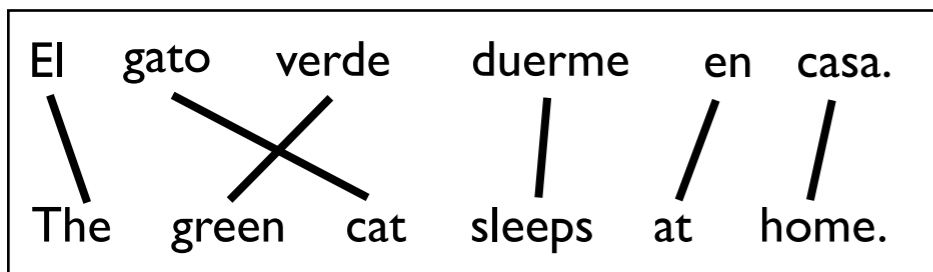
Train 10K sentences of FBIS  
Chinese-English newswire

Test NIST 2002 Chinese-English dev set





# Word Alignment Results



## Alignment Error Rate

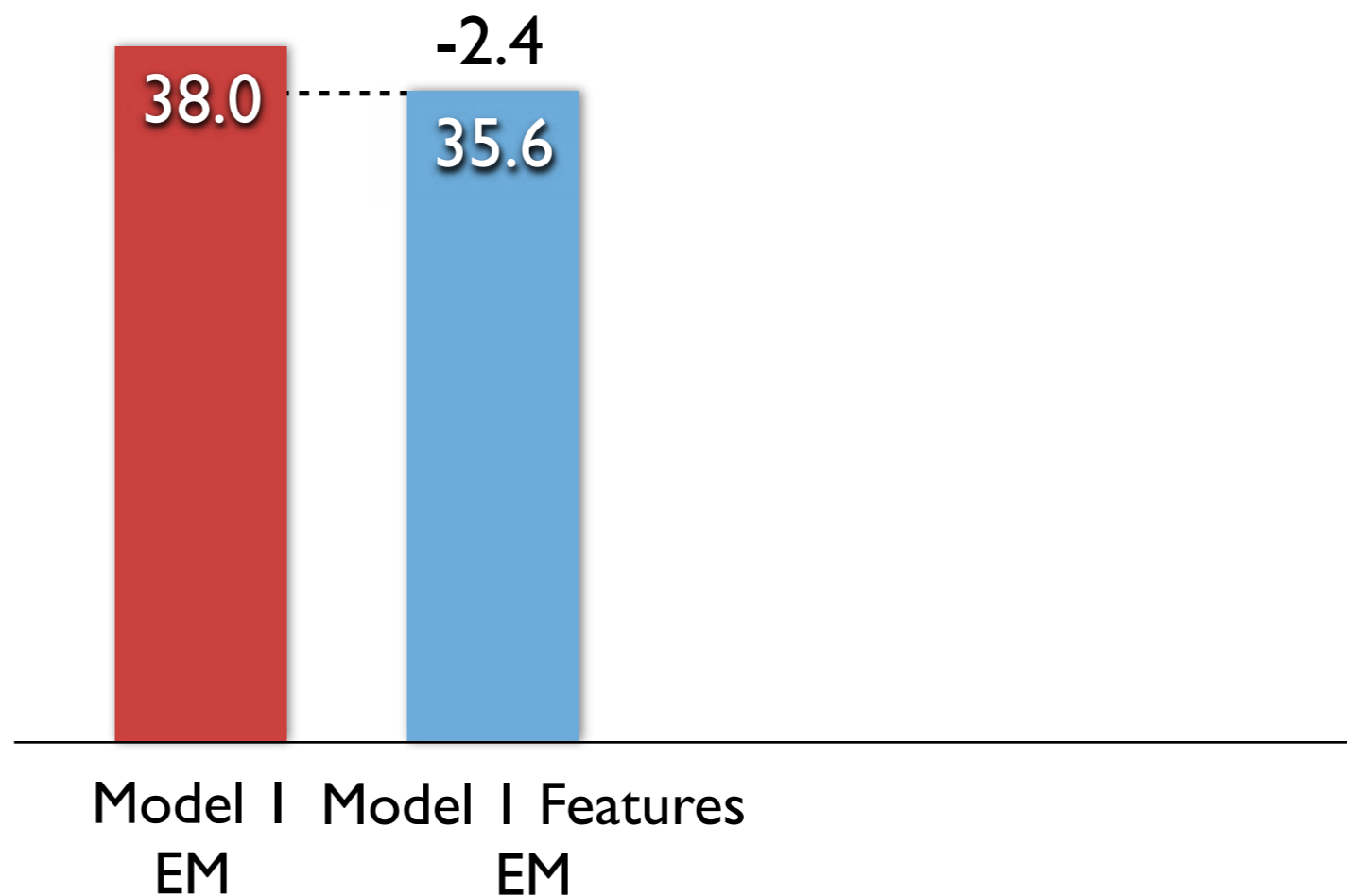
### Features:

- Basic:  $gato \wedge cat$
- Edit-Distance:  $edit(gato, cat) = 2$
- Dictionary:  $(gato, cat) \in Dict$
- Stem:  $gato \wedge +stem(cat)$
- Prefix:  $gato \wedge +ca$

### Data:

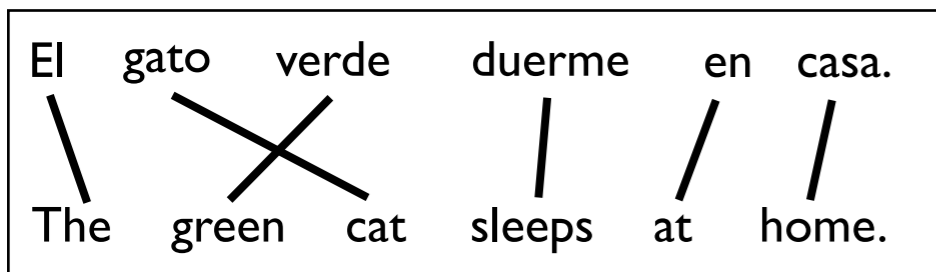
Train 10K sentences of FBIS  
Chinese-English newswire

Test NIST 2002 Chinese-English dev set





# Word Alignment Results



## Alignment Error Rate

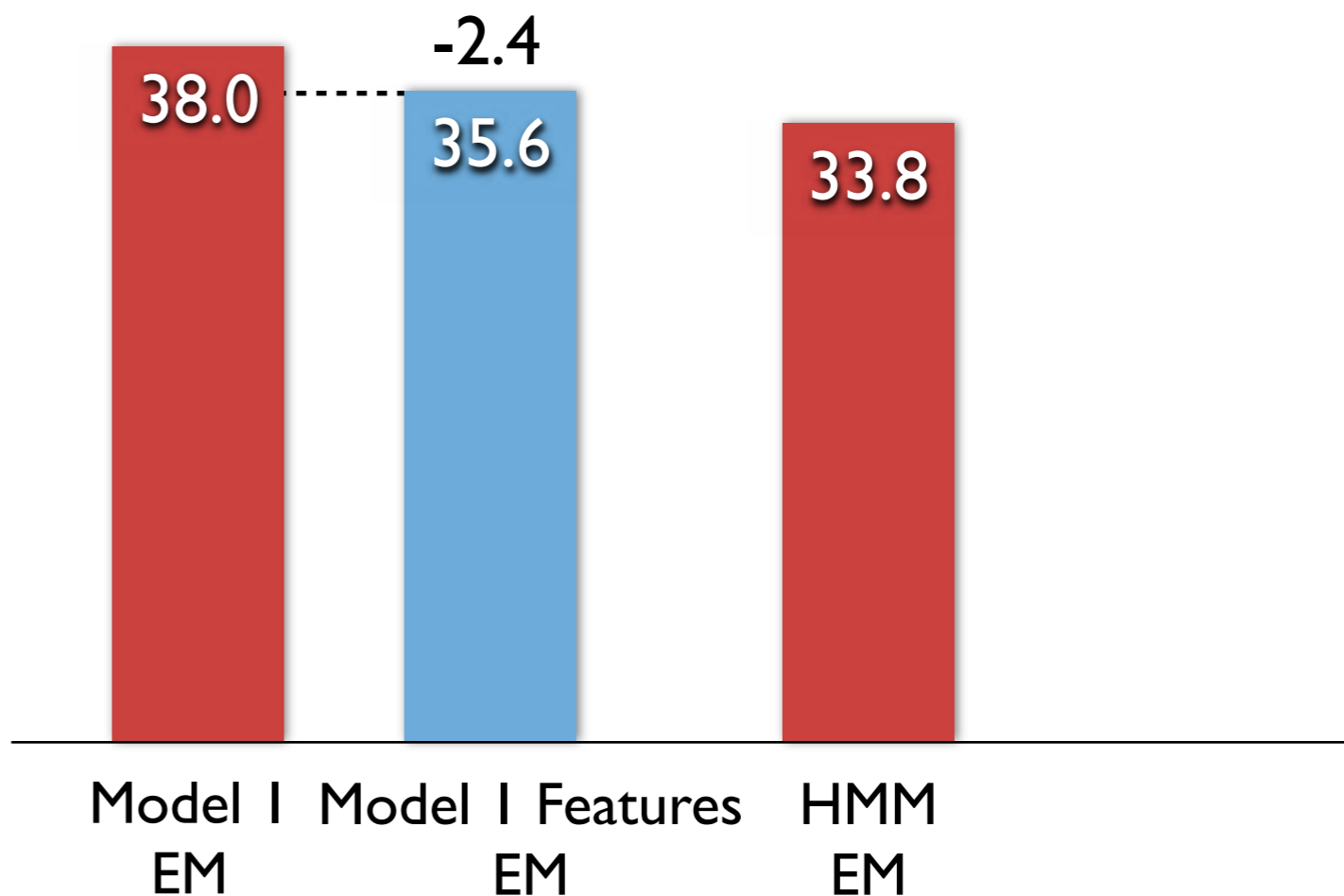
### Features:

- Basic:  $\text{gato} \wedge \text{cat}$
- Edit-Distance:  $\text{edit}(\text{gato}, \text{cat}) = 2$
- Dictionary:  $(\text{gato}, \text{cat}) \in \text{Dict}$
- Stem:  $\text{gato} \wedge +\text{stem}(\text{cat})$
- Prefix:  $\text{gato} \wedge +\text{ca}$

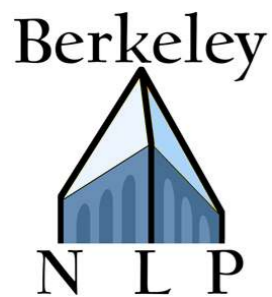
### Data:

Train 10K sentences of FBIS  
Chinese-English newswire

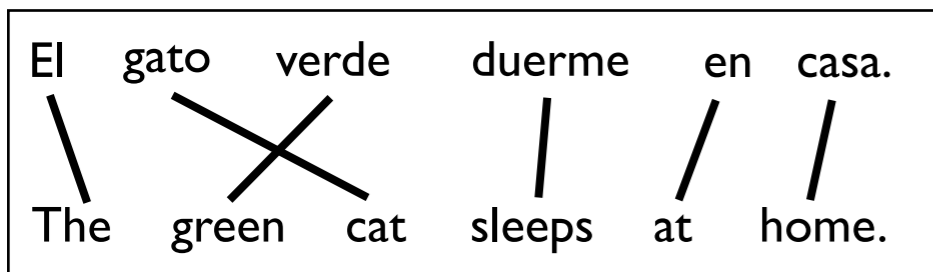
Test NIST 2002 Chinese-English dev set







# Word Alignment Results



## Alignment Error Rate

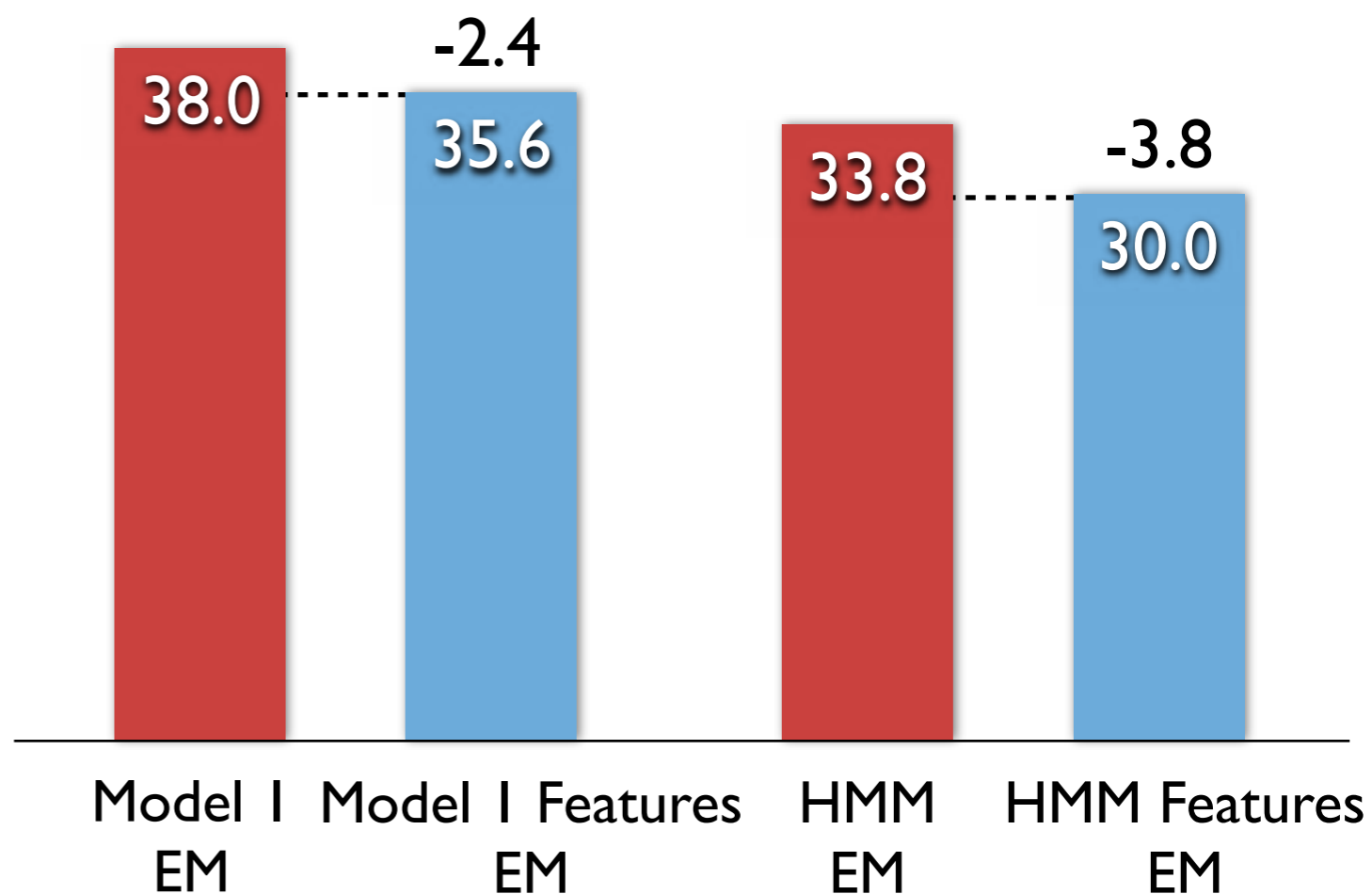
### Features:

- Basic:  $gato \wedge cat$
- Edit-Distance:  $edit(gato, cat) = 2$
- Dictionary:  $(gato, cat) \in Dict$
- Stem:  $gato \wedge +stem(cat)$
- Prefix:  $gato \wedge +ca$

### Data:

Train 10K sentences of FBIS  
Chinese-English newswire

Test NIST 2002 Chinese-English dev set





# Word Segmentation Results

---

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



# Word Segmentation Results

---

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

Key distribution:  $P(\text{running})$

# Word Segmentation Results

---

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

Key distribution:  $P(\text{running})$

Features:

Basic: running

Length:  $\text{length}(\text{running}) = 7$

Num-Vowels:  $\text{numV}(\text{running}) = 2$

Coarse-Phono-Prefix: +rAn

Coarse-Phono-Suffix: +IN



# Word Segmentation Results

---

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

Token F1

## Features:

Basic: running  
Length: length(running) = 7  
Num-Vowels: numV(running) = 2  
Coarse-Phono-Prefix: +rAn  
Coarse-Phono-Suffix: +IN

## Data:

---

Train and test on phonetic version  
of Bernstein-Ratner corpus

# Word Segmentation Results

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

Token F1

## Features:

Basic: running  
Length: length(running) = 7  
Num-Vowels: numV(running) = 2  
Coarse-Phono-Prefix: +rAn  
Coarse-Phono-Suffix: +IN

## Data:

Train and test on phonetic version  
of Bernstein-Ratner corpus

76.9

Unigram  
EM

# Word Segmentation Results

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

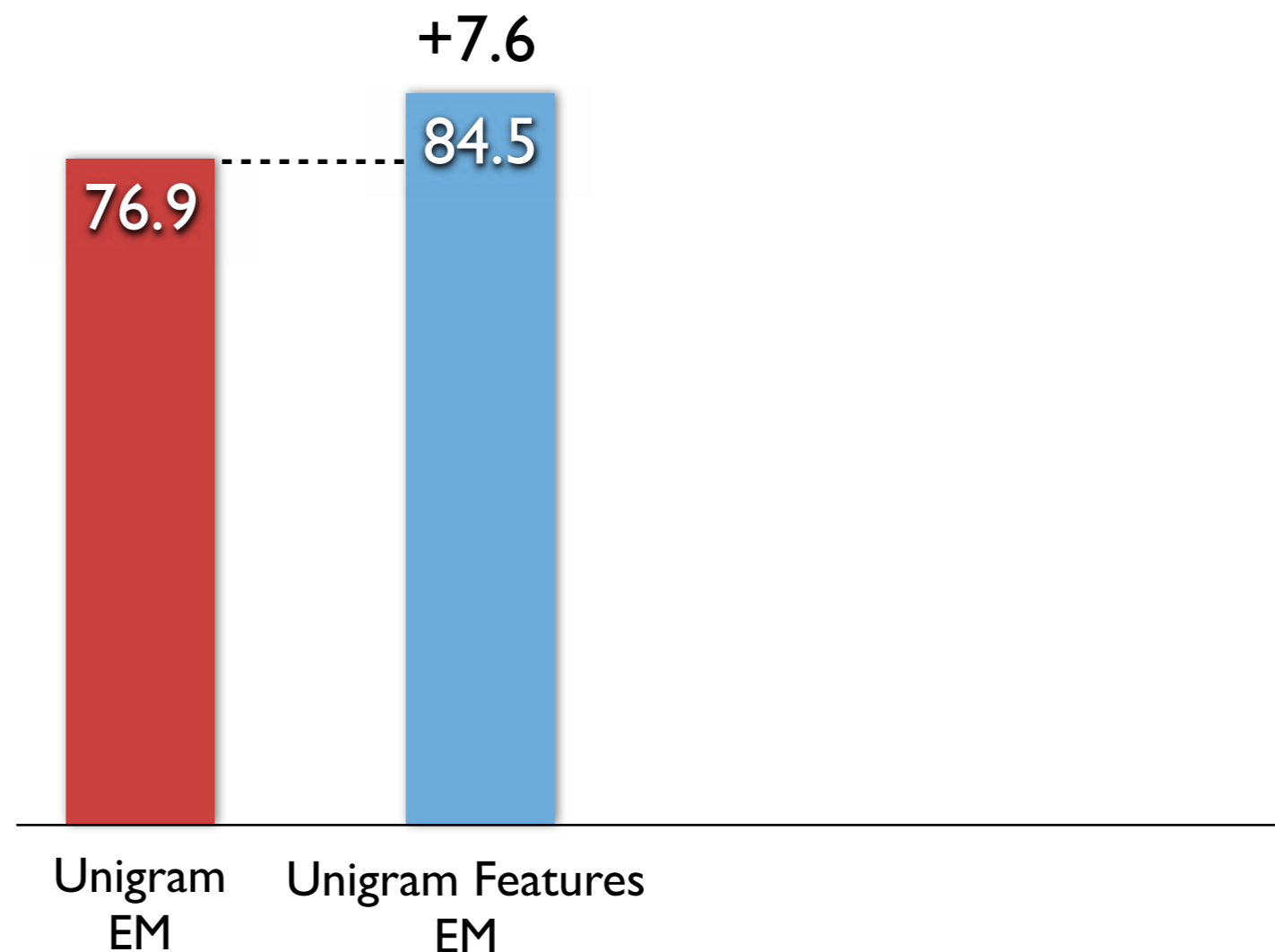
Token F1

## Features:

Basic: running  
Length: length(running) = 7  
Num-Vowels: numV(running) = 2  
Coarse-Phono-Prefix: +rAn  
Coarse-Phono-Suffix: +IN

## Data:

Train and test on phonetic version  
of Bernstein-Ratner corpus



# Word Segmentation Results

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

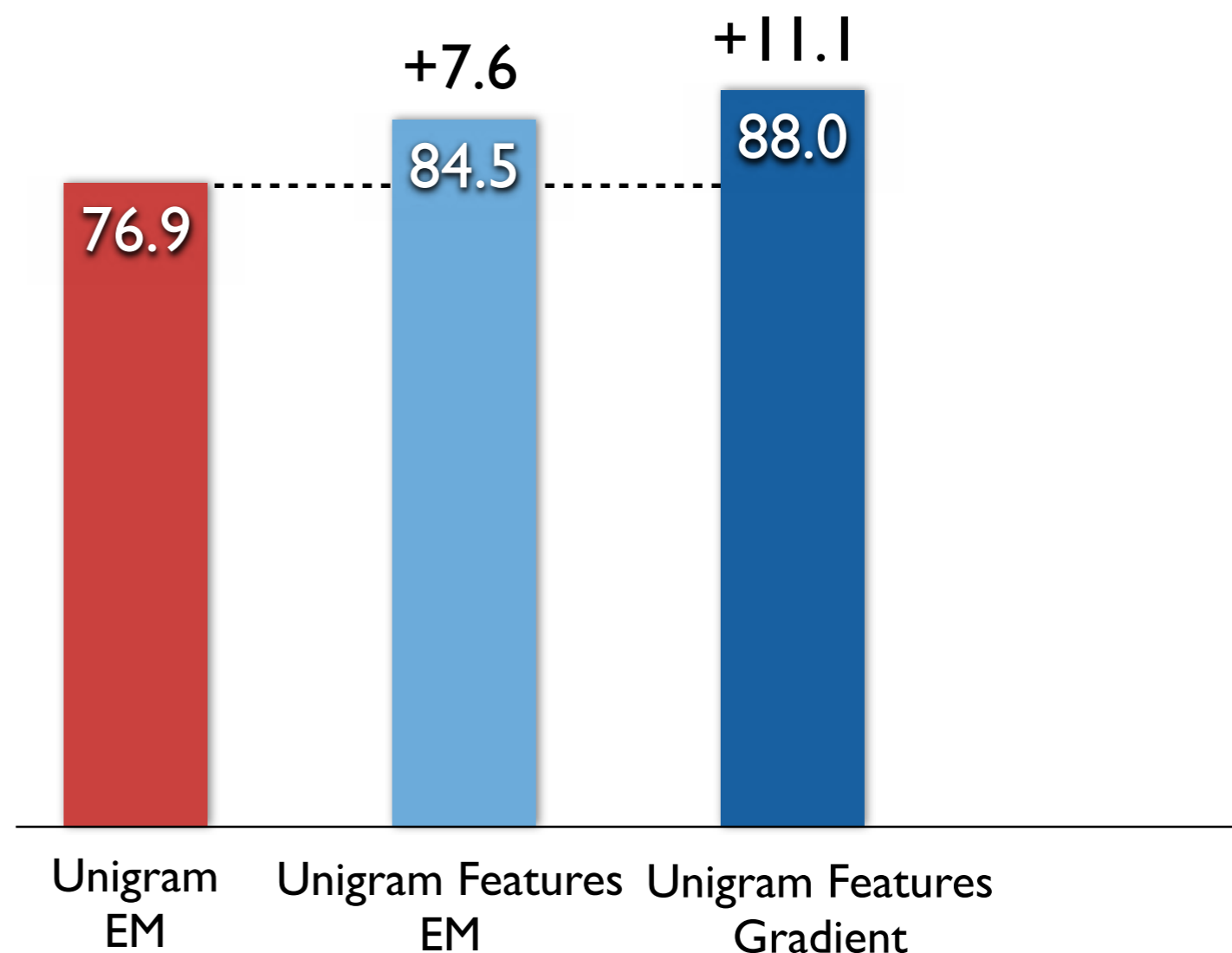
## Features:

Basic: running  
Length:  $\text{length}(\text{running}) = 7$   
Num-Vowels:  $\text{numV}(\text{running}) = 2$   
Coarse-Phono-Prefix: +rAn  
Coarse-Phono-Suffix: +IN

## Data:

Train and test on phonetic version  
of Bernstein-Ratner corpus

## Token F1





# Word Segmentation Results

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

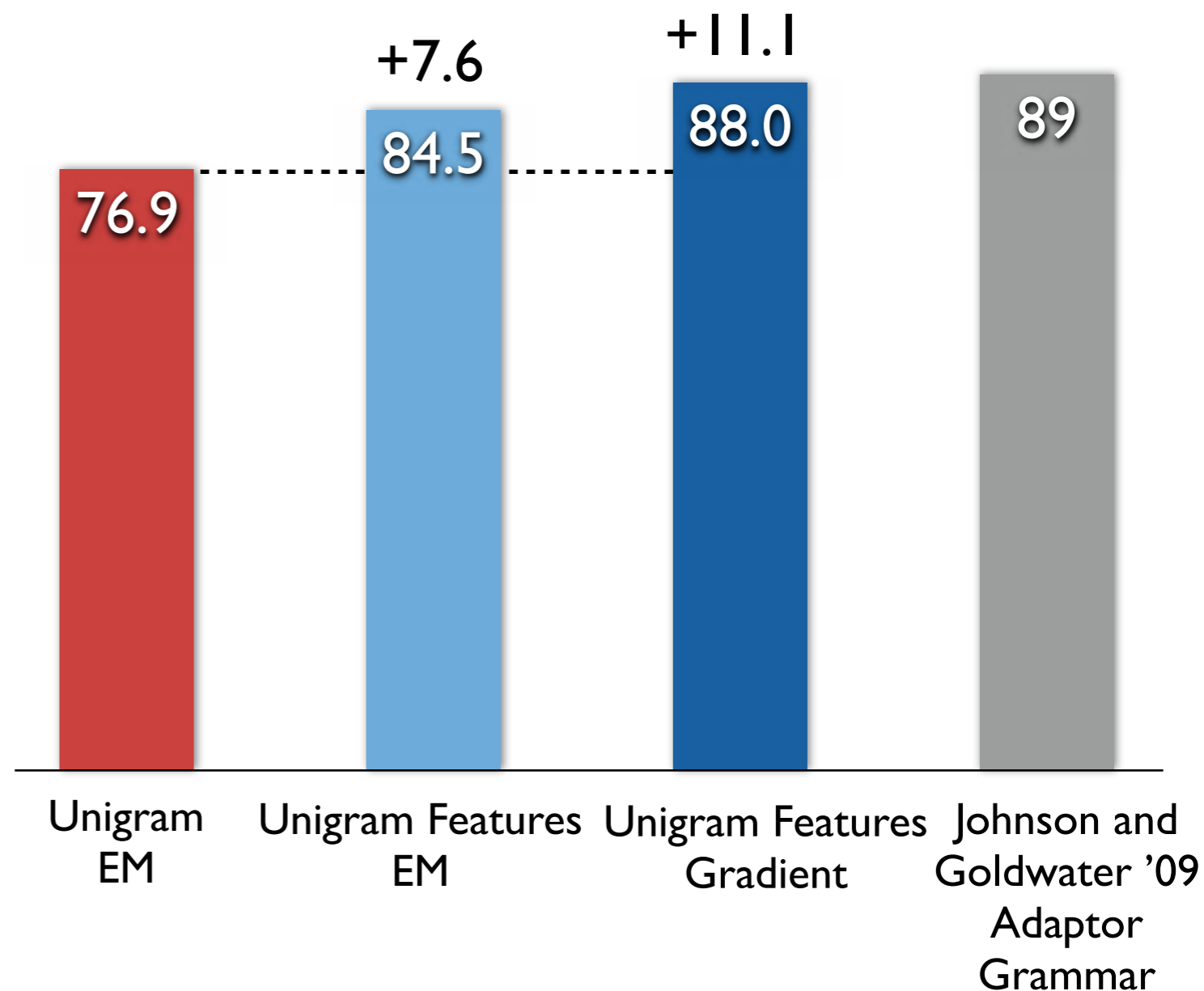
## Features:

Basic: running  
 Length: length(running) = 7  
 Num-Vowels: numV(running) = 2  
 Coarse-Phono-Prefix: +rAn  
 Coarse-Phono-Suffix: +IN

## Data:

Train and test on phonetic version  
 of Bernstein-Ratner corpus

## Token F1





# Apply to New Models

---

- I. Take a generative model



# Apply to New Models

---

1. Take a generative model
2. Brainstorm features local to the component multinomials

# Apply to New Models

---

1. Take a generative model
2. Brainstorm features local to the component multinomials
3. Run this algorithm

# Apply to New Models

---

1. Take a generative model
2. Brainstorm features local to the component multinomials
3. Run this algorithm
4. Crush your baseline



# Conclusion

---

- State-of-the-art results



# Conclusion

---

- State-of-the-art results
- Can implemented using off-the-shelf NLP tools

# Conclusion

---

- State-of-the-art results
- Can implemented using off-the-shelf NLP tools
- Directly optimizing data-likelihood can outperform EM



# Conclusion

---

- State-of-the-art results
- Can implemented using off-the-shelf NLP tools
- Directly optimizing data-likelihood can outperform EM
- Works on a wide range of induction tasks

Berkeley



# Conclusion

---

**Thanks!**