

# Concavity and Initialization for Unsupervised Dependency Parsing

Kevin Gimpel

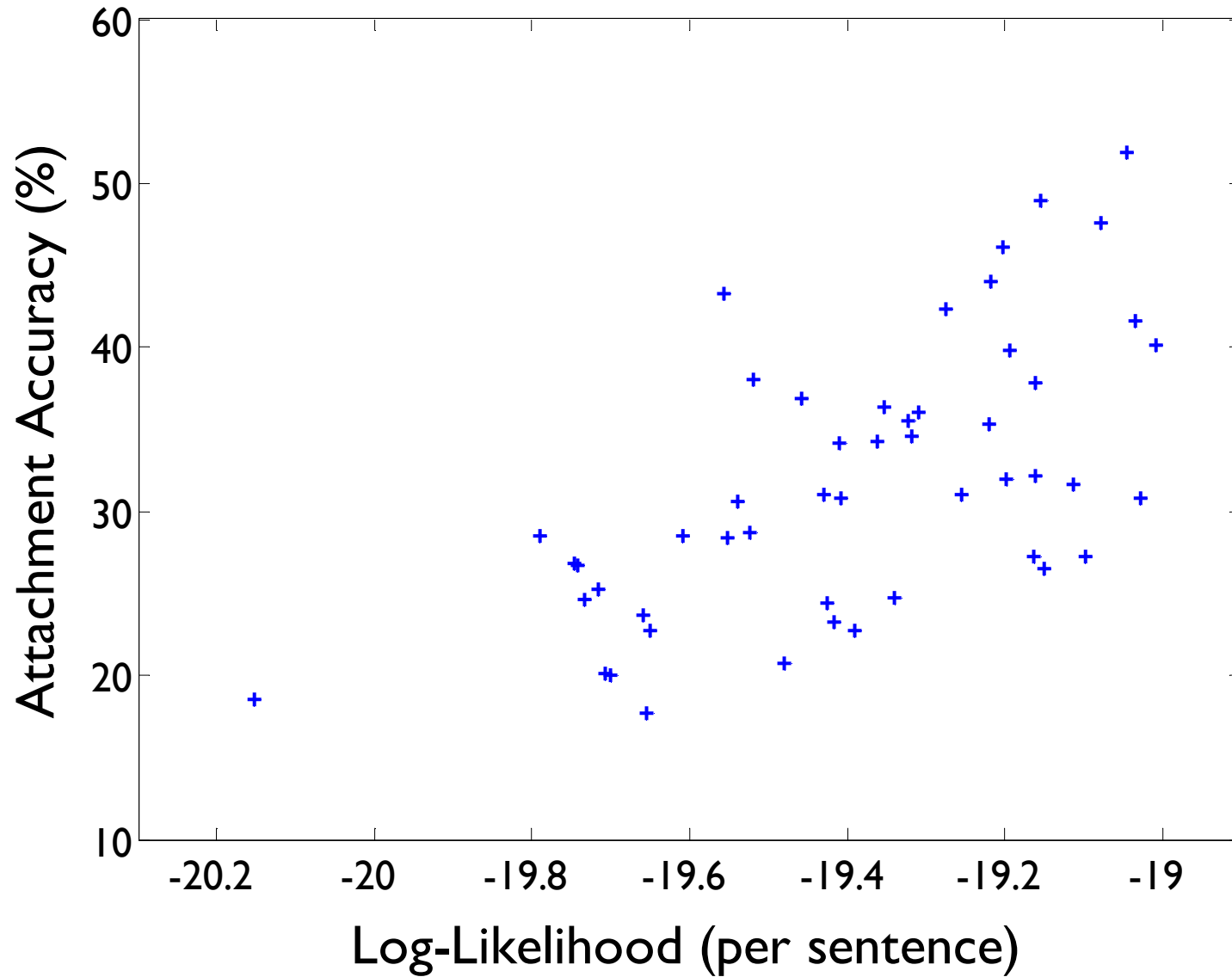
Noah A. Smith



**Carnegie Mellon**

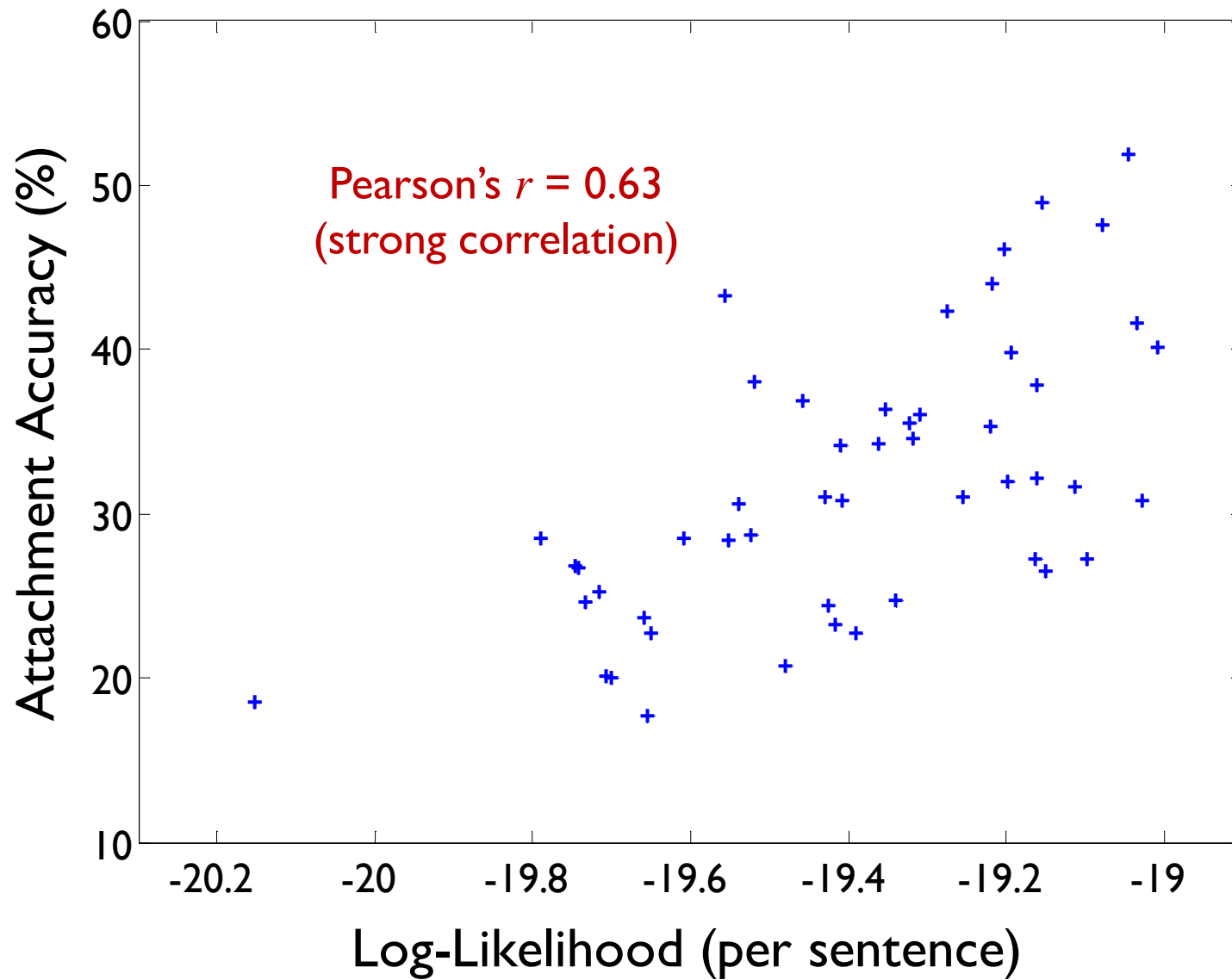
Unsupervised learning in NLP <sup>(typically)</sup> → non-convex optimization

# Dependency Model with Valence (Klein & Manning, 2004)

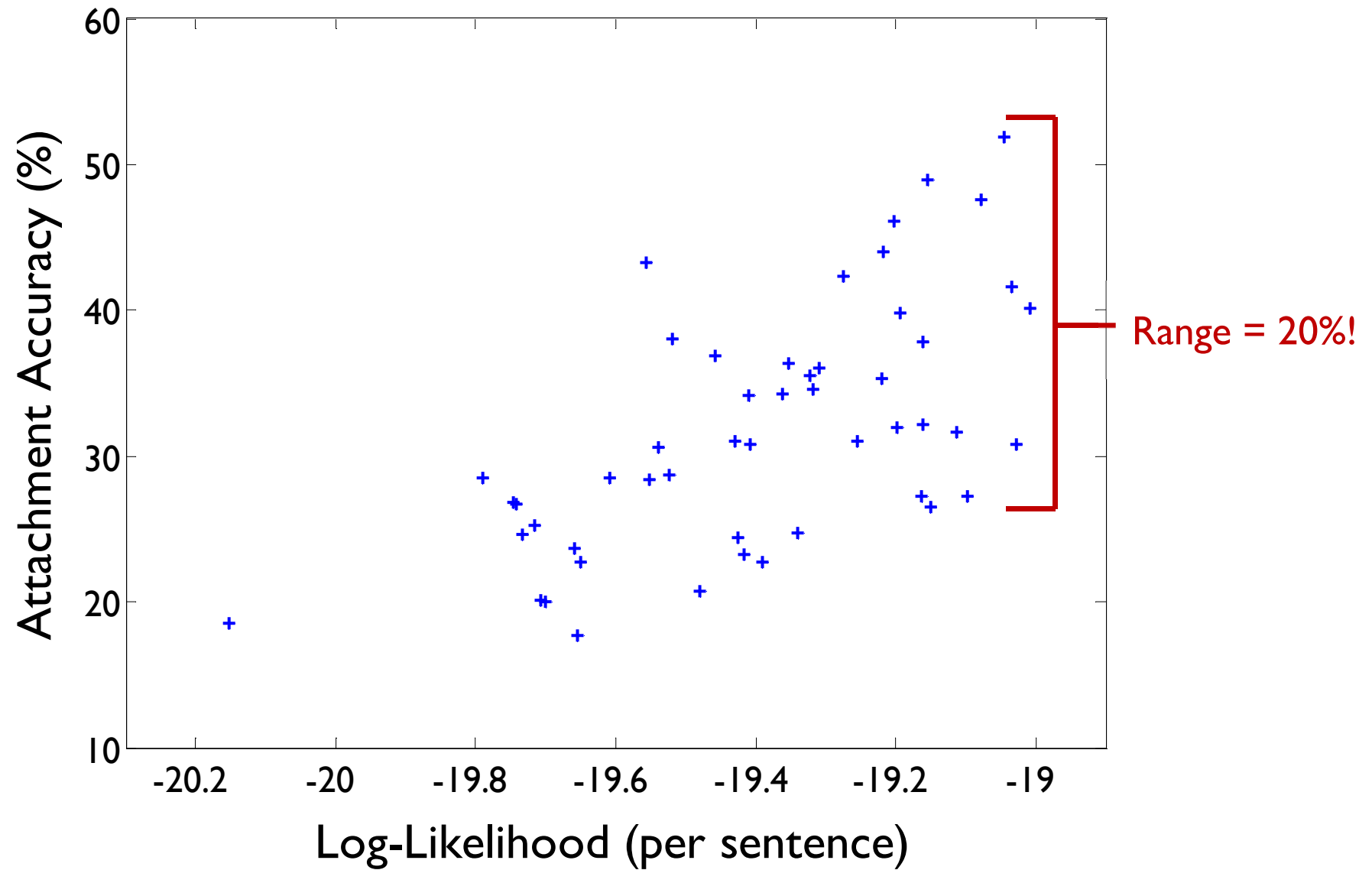


EM with 50  
Random  
Initializers

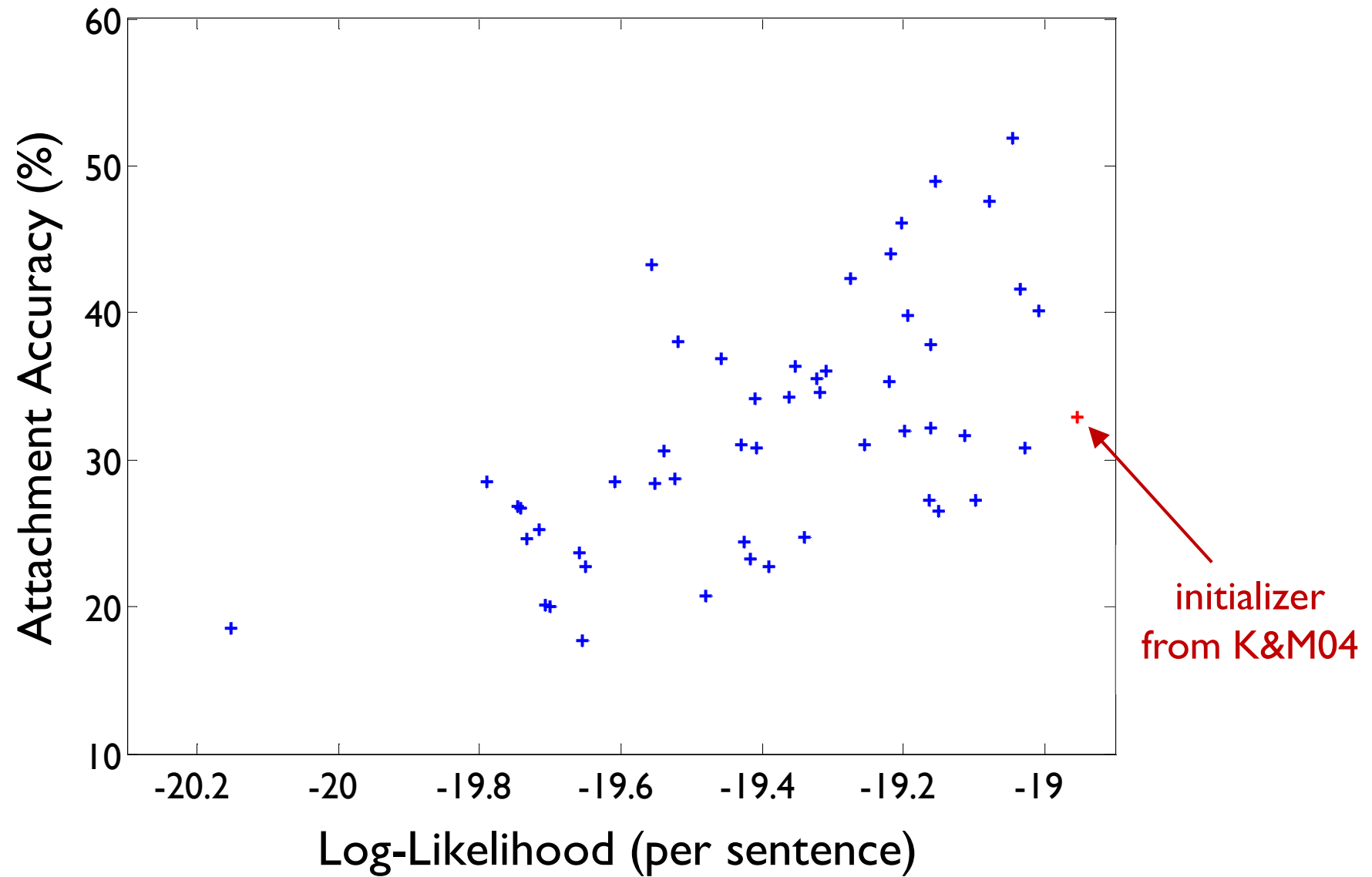
# Dependency Model with Valence (Klein & Manning, 2004)



# Dependency Model with Valence (Klein & Manning, 2004)



# Dependency Model with Valence (Klein & Manning, 2004)



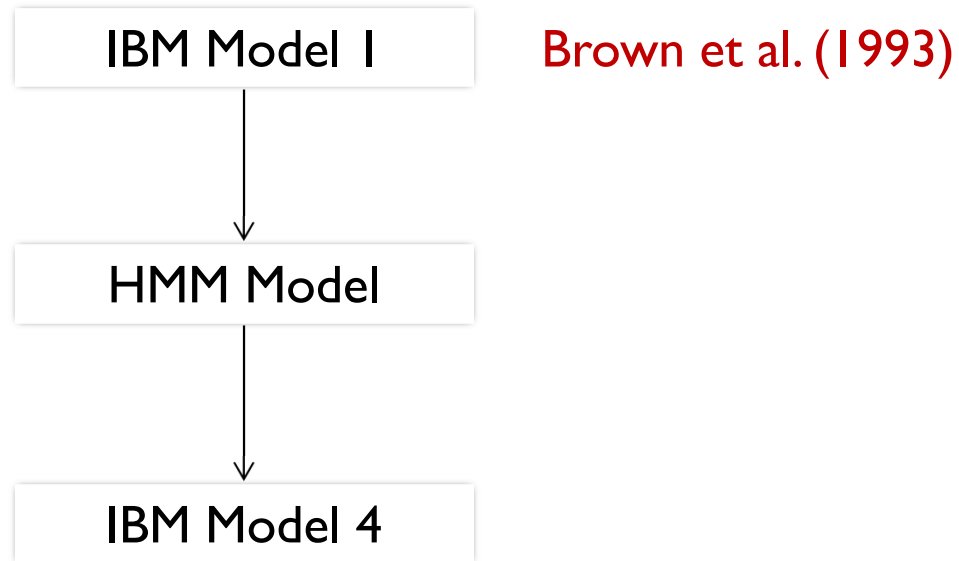
# How has this been addressed?

- Scaffolding / staged training (**Brown et al., 1993**; Elman, 1993; Spitkovsky et al., 2010)
- Curriculum learning (Bengio et al., 2009)
- Deterministic annealing (Smith & Eisner, 2004), Structural annealing (Smith & Eisner, 2006)
- Continuation methods (Allgower & Georg, 1990)



**Carnegie Mellon**

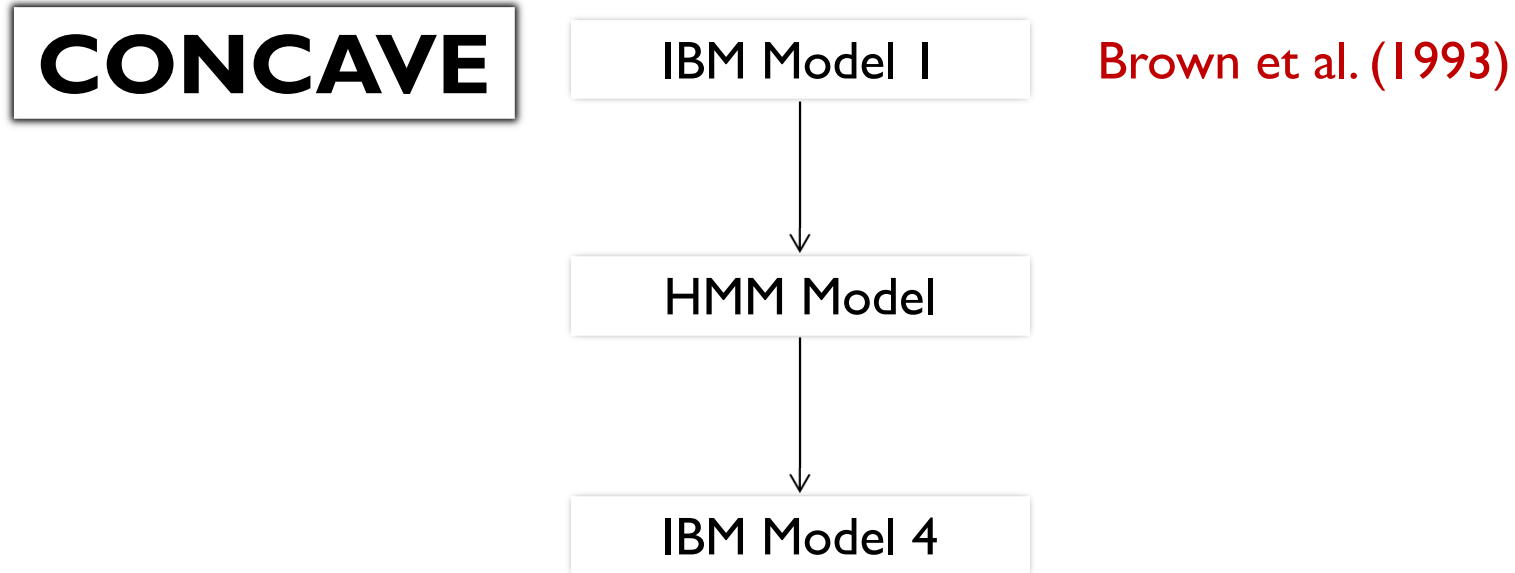
# Example: Word Alignment



Carnegie Mellon



# Example: Word Alignment



Carnegie Mellon

Unsupervised learning in NLP <sup>(typically)</sup> → non-convex optimization

Unsupervised learning in NLP <sup>(typically)</sup> → non-convex optimization

Except IBM Model I for word alignment  
(which has a concave log-likelihood function)

# IBM Model I (Brown et al., 1993)

$$\log p(\mathbf{e} \mid \mathbf{f}) = \log \epsilon + \sum_{j=1}^{|\mathbf{e}|} \log \sum_{i=0}^{|\mathbf{f}|} \frac{1}{|\mathbf{f}| + 1} t(e_j \mid f_i)$$



Carnegie Mellon

# IBM Model I (Brown et al., 1993)

$$\log p(\mathbf{e} \mid \mathbf{f}) = \log \epsilon + \sum_{j=1}^{|\mathbf{e}|} \log \sum_{i=0}^{|\mathbf{f}|} \underbrace{\frac{1}{|\mathbf{f}| + 1}}_{\text{alignment probability}} \underbrace{t(e_j \mid f_i)}_{\text{translation probability}}$$



Carnegie Mellon

# IBM Model 1 (Brown et al., 1993)

$$\log p(\mathbf{e} \mid \mathbf{f}) = \log \epsilon + \sum_{j=1}^{|\mathbf{e}|} \log \sum_{i=0}^{|\mathbf{f}|} \underbrace{\frac{1}{|\mathbf{f}| + 1}}_{\text{alignment probability}} \underbrace{t(e_j \mid f_i)}_{\text{translation probability}}$$

# IBM Model 2

$$\log p(\mathbf{e} \mid \mathbf{f}) = \log \epsilon + \sum_{j=1}^{|\mathbf{e}|} \log \sum_{i=0}^{|\mathbf{f}|} a(i \mid j, |\mathbf{f}|, |\mathbf{e}|) t(e_j \mid f_i)$$



Carnegie Mellon

# IBM Model 1 **CONCAVE**

$$\log p(\mathbf{e} | \mathbf{f}) = \log \epsilon + \sum_{j=1}^{|\mathbf{e}|} \log \sum_{i=0}^{|\mathbf{f}|} \underbrace{\frac{1}{|\mathbf{f}| + 1}}_{\text{alignment probability}} \underbrace{t(e_j | f_i)}_{\text{translation probability}}$$

# IBM Model 2

**NOT  
CONCAVE**

$$\log p(\mathbf{e} | \mathbf{f}) = \log \epsilon + \sum_{j=1}^{|\mathbf{e}|} \log \sum_{i=0}^{|\mathbf{f}|} a(i | j, |\mathbf{f}|, |\mathbf{e}|) t(e_j | f_i)$$



Carnegie Mellon

# IBM Model 1 **CONCAVE**

$$\log p(\mathbf{e} | \mathbf{f}) = \log \epsilon + \sum_{j=1}^{|\mathbf{e}|} \log \sum_{i=0}^{|\mathbf{f}|} \underbrace{\frac{1}{|\mathbf{f}| + 1}}_{\text{alignment probability}} \underbrace{t(e_j | f_i)}_{\text{translation probability}}$$

# IBM Model 2

**NOT  
CONCAVE**

$$\log p(\mathbf{e} | \mathbf{f}) = \log \epsilon + \sum_{j=1}^{|\mathbf{e}|} \log \sum_{i=0}^{|\mathbf{f}|} \underbrace{a(i | j, |\mathbf{f}|, |\mathbf{e}|)}_{\text{product of parameters within log-sum}} t(e_j | f_i)$$



Carnegie Mellon



# IBM Model I **CONCAVE**

$$\log p(\mathbf{e} \mid \mathbf{f}) = \log \epsilon + \sum_{j=1}^{|\mathbf{e}|} \log \sum_{i=0}^{|\mathbf{f}|} \underbrace{\frac{1}{|\mathbf{f}| + 1}}_{\text{red}} \underbrace{t(e_j \mid f_i)}_{\text{blue}}$$

## For concavity:

1 parameter is permitted for each atomic piece of latent structure.

No atomic piece of latent structure can affect any other piece.

$$\log p(\mathbf{e} \mid \mathbf{f}) = \log \epsilon + \sum_{j=1}^{|\mathbf{e}|} \log \sum_{i=0}^{|\mathbf{f}|} \underbrace{a(i \mid j, |\mathbf{f}|, |\mathbf{e}|)}_{\text{green}} t(e_j \mid f_i)$$

product of parameters  
within log-sum



Carnegie Mellon

Unsupervised learning in NLP <sup>(typically)</sup> → non-convex optimization

Except IBM Model I for word alignment  
(which has a concave log-likelihood function)

What models can we build without sacrificing concavity?

## For concavity:

- 1 parameter is permitted for each atomic piece of latent structure.
- No atomic piece of latent structure can affect any other piece.



**Carnegie Mellon**

## For concavity:

1 parameter is permitted for each atomic piece of latent structure.

No atomic piece of latent structure can affect any other piece.

single dependency arc



Carnegie Mellon

## For concavity:

1 parameter is permitted for each atomic piece of latent structure.

No atomic piece of latent structure can affect any other piece.

single dependency arc

Every dependency arc must be independent,  
so **we can't use a tree constraint**



Carnegie Mellon

## For concavity:

1 parameter is permitted for each atomic piece of latent structure.

No atomic piece of latent structure can affect any other piece.

single dependency arc

Every dependency arc must be independent,  
so **we can't use a tree constraint**

Only one parameter allowed per dependency arc



Carnegie Mellon

## For concavity:

1 parameter is permitted for each atomic piece of latent structure.

No atomic piece of latent structure can affect any other piece.

single dependency arc

Our Model:

Like IBM Model I, but we generate the same sentence again,  
aligning words to the original sentence (cf. Brody, 2010)

$$\log p(\mathbf{e} \mid \mathbf{e}') = \log \epsilon + \sum_{j=1}^{|\mathbf{e}|} \log \sum_{i=0, i \neq j}^{|\mathbf{e}|} \frac{1}{|\mathbf{e}| + 1} \text{child}(e_j \mid e'_i)$$



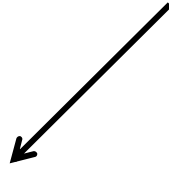
Carnegie Mellon

\$ Vikings came in longboats from Scandinavia in 1000 AD





\$ Vikings came in longboats from Scandinavia in 1000 AD



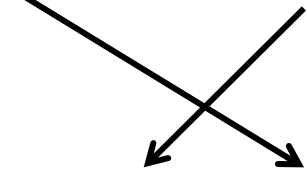
Vikings

\$ Vikings came in longboats from Scandinavia in 1000 AD



**Carnegie Mellon**

\$ Vikings came in longboats from Scandinavia in 1000 AD



Vikings came

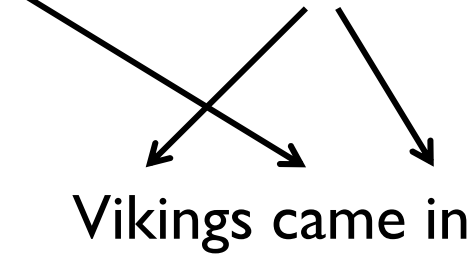


\$ Vikings came in longboats from Scandinavia in 1000 AD



**Carnegie Mellon**

\$ Vikings came in longboats from Scandinavia in 1000 AD

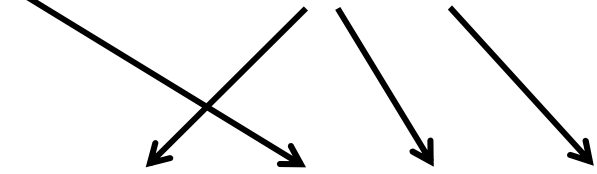


\$ Vikings came in longboats from Scandinavia in 1000 AD



Carnegie Mellon

\$ Vikings came in longboats from Scandinavia in 1000 AD



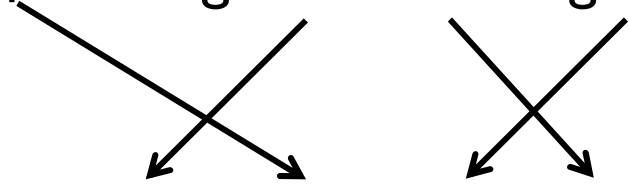
Vikings came in longboats

\$ Vikings came in longboats from Scandinavia in 1000 AD



Carnegie Mellon

\$ Vikings came in longboats from Scandinavia in 1000 AD



Vikings came in longboats



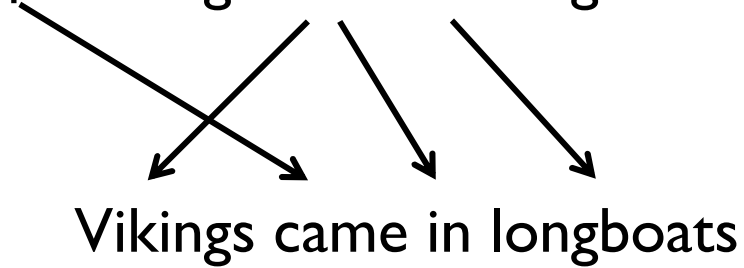
\$ Vikings came in longboats from Scandinavia in 1000 AD

Cycles, multiple roots, and non-projectivity  
are all **permitted** by this model



Carnegie Mellon

\$ Vikings came in longboats from Scandinavia in 1000 AD

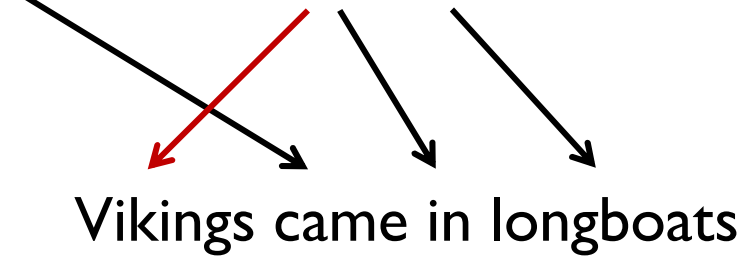


Only one parameter per dependency arc:



Carnegie Mellon

\$ Vikings came in longboats from Scandinavia in 1000 AD

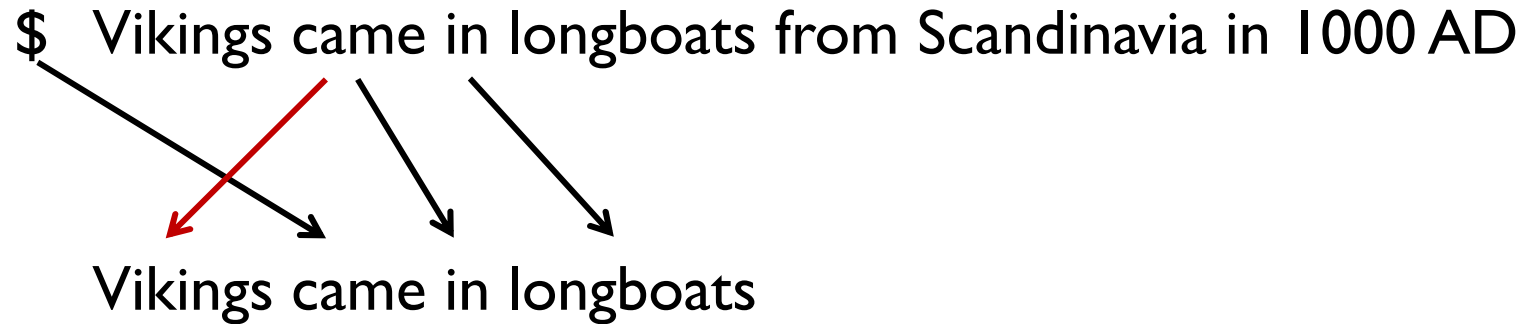


Only one parameter per dependency arc:

$$p(\text{Vikings} \mid \text{came})$$



Carnegie Mellon



Only one parameter per dependency arc:

$$p(\text{Vikings} \mid \text{came})$$

We cannot look at other dependency arcs, but we can condition on (properties of) the sentence:

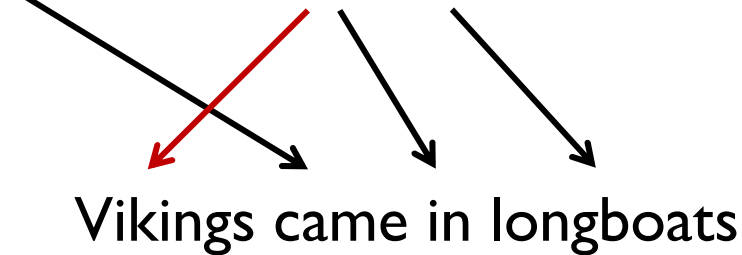
$$p(\text{Vikings} \mid \text{came}, \textit{direction} = \textit{left}, \textit{distance} = 1, \dots)$$



Carnegie Mellon



\$ Vikings came in longboats from Scandinavia in 1000 AD



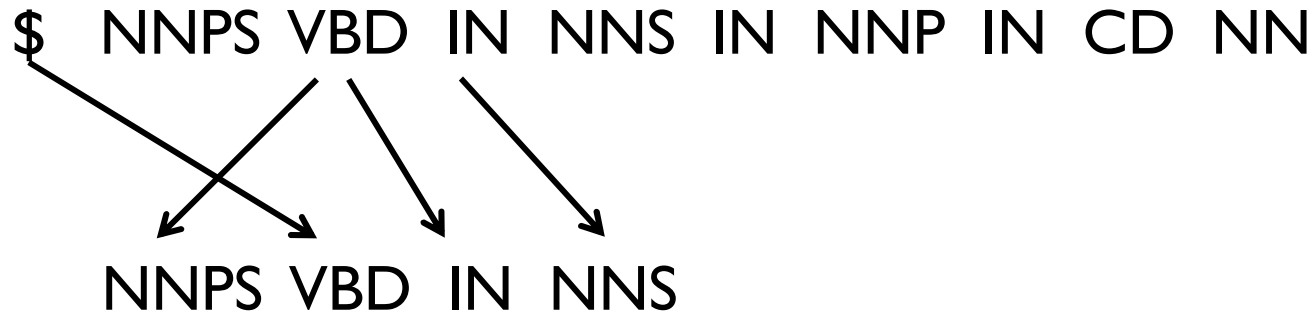
We condition on direction:

$$p(\text{Vikings} \mid \text{came}, \textit{direction} = \textit{left})$$

("Concave Model A")



Carnegie Mellon



Note: we've been using words in our examples, but in our model we follow standard practice and use gold POS tags

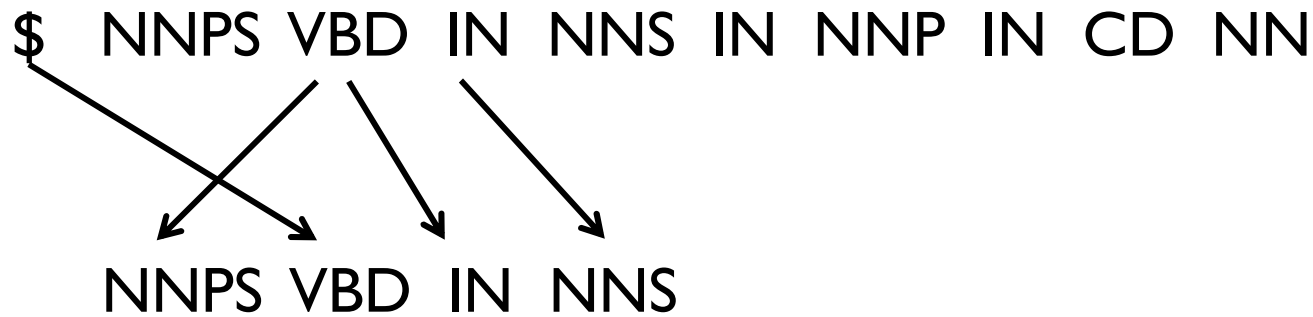
We condition on direction:

$$p(\text{NNPS} \mid \text{VBD}, \text{direction} = \text{left})$$

("Concave Model A")



Carnegie Mellon



Model	Initializer	Accuracy*
Attach Right	N/A	31.7
DMV	Uniform	17.6
DMV	K&M	32.9
Concave Model A	Uniform	25.6

\*Penn Treebank  
test set, sentences  
of all lengths

WSJ10 used for  
training

We condition on direction:

$$p(\text{NNPS} \mid \text{VBD}, \text{direction} = \text{left})$$

(“Concave Model A”)

## Note:

IBM Model 1 is not *strictly* concave  
(Toutanova & Galley, 2011)

Model	Initializer	Accuracy*
Attach Right	N/A	31.7
DMV	Uniform	17.6
DMV	K&M	32.9
Concave Model A	Uniform	25.6

\*Penn Treebank  
test set, sentences  
of all lengths

WSJ10 used for  
training

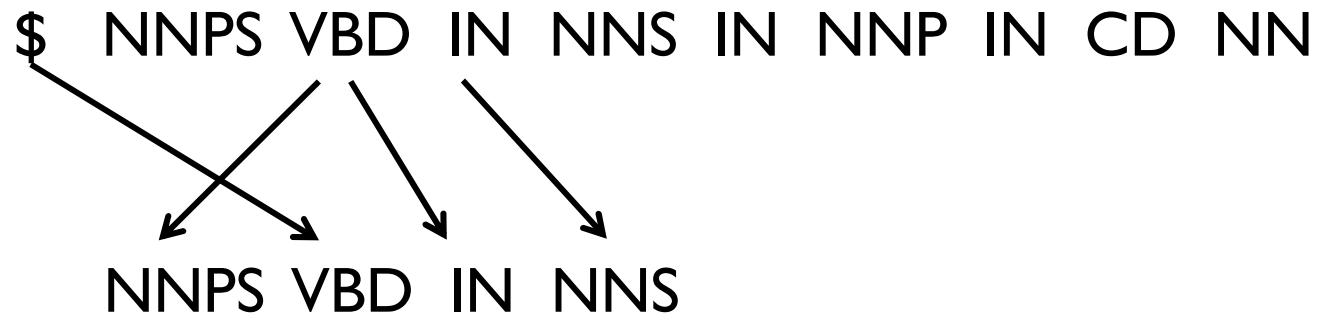
We condition on direction:

$$p(\text{NNPS} \mid \text{VBD}, \text{direction} = \text{left})$$

(“Concave Model A”)



Carnegie Mellon



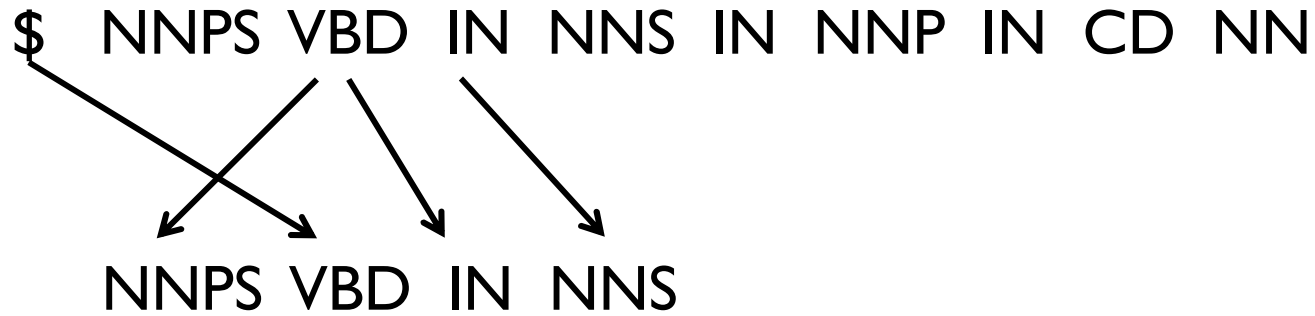
We can also use hard constraints while preserving concavity:

The only tags that can align to \$ are verbs  
(Marecček & Žabokrtský, 2011; Naseem et al., 2010)

(“Concave Model B”)



Carnegie Mellon



Model	Initializer	Accuracy*
Attach Right	N/A	31.7
DMV	Uniform	17.6
DMV	K&M	32.9
Concave Model A	Uniform	25.6
Concave Model B	Uniform	28.6

\*Penn Treebank  
test set, sentences  
of all lengths

WSJ10 used for  
training



Carnegie Mellon

Unsupervised learning in NLP <sup>(typically)</sup> → non-convex optimization

Except IBM Model I for word alignment  
(which has a concave log-likelihood function)

What models can we build without sacrificing concavity?

Can these concave models be useful?



**Carnegie Mellon**

As IBM Model I is used to initialize other word alignment models,  
we can use our concave models to initialize the DMV





As IBM Model I is used to initialize other word alignment models,  
we can use our concave models to initialize the DMV

Model	Initializer	Accuracy
Attach Right	N/A	31.7
DMV	Uniform	17.6
DMV	K&M	32.9
DMV	Concave Model A	34.4
DMV	Concave Model B	<b>43.0</b>

\*Penn Treebank  
test set, sentences  
of all lengths

WSJ10 used for  
training

As IBM Model I is used to initialize other word alignment models,  
we can use our concave models to initialize the DMV

Model	Initializer	Accuracy*
DMV, trained on sentences of length $\leq 20$	Concave Model B	53.1
Shared Logistic Normal (Cohen & Smith, 2009)	K&M	41.4
Posterior Regularization (Gillenwater et al., 2010)	K&M	53.3
LexTSG-DMV (Blunsom & Cohn, 2010)	K&M	55.7
Punctuation/UnsupTags (Spitkovsky et al., 2011), trained on sentences of length $\leq 45$	K&M'	59.1

\*Penn Treebank test set, sentences of all lengths



Carnegie Mellon

# Multilingual Results

(averages across 18 languages)

Model	Initializer	Avg. Accuracy*	Avg. Log-Likelihood †
DMV	Uniform	25.7	-15.05
DMV	K&M	29.4	-14.84
DMV	Concave Model A	30.9	-14.93
DMV	Concave Model B	<b>35.5</b>	<b>-14.45</b>

\* Sentences of all lengths from each test set

† Micro-averaged across sentences in all training sets  
(used sentences  $\leq 10$  words for training)



Carnegie Mellon

Unsupervised learning in NLP <sup>(typically)</sup> → non-convex optimization

Except IBM Model I for word alignment  
(which has a concave log-likelihood function)

What models can we build without sacrificing concavity?

Can these concave models be useful?

Like word alignment, we can use simple, concave models to initialize more complex models for grammar induction



**Carnegie Mellon**

Thanks!



**Carnegie Mellon**