

# Domain Adaptation



John Blitzer and Hal Daumé III

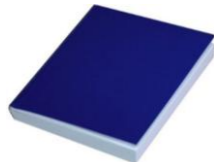


# Classical “Single-domain” Learning



Predict:  $x \rightarrow y$        $(x, y) \sim \text{Pr}[x, y]$

amazon.com



**Running with Scissors**

**Title:** Horrible book, horrible.

This book was horrible. I read half, suffering from a headache the entire time, and eventually i lit it on fire. 1 less copy in the world. Don't waste your money. I wish i had the time spent reading this book back. It wasted my life



So                      the topic of                      ah                      the talk today is online learning



# Domain Adaptation



$$(x, y) \sim \text{Pr}_S[x, y]$$

**Training**

**Source**

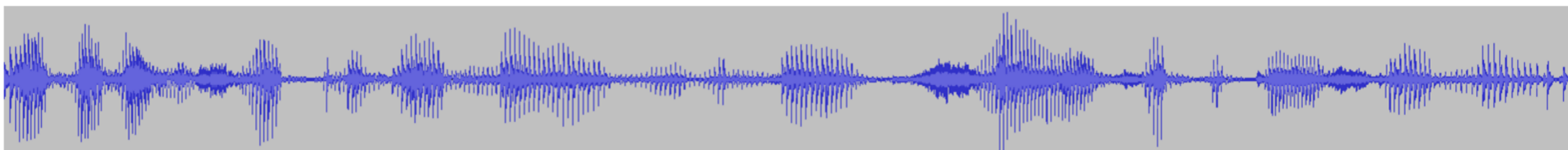


So the topic of ah the talk today is online learning

$$(x, y) \sim \text{Pr}_T[x, y]$$

**Testing**

**Target**



Everything is happening online. Even the slides are produced on-line



# Domain Adaptation



## Natural Language Processing

## Visual Object Recognition



Packed with fascinating info



A breeze to clean up





# Domain Adaptation



## Natural Language Processing



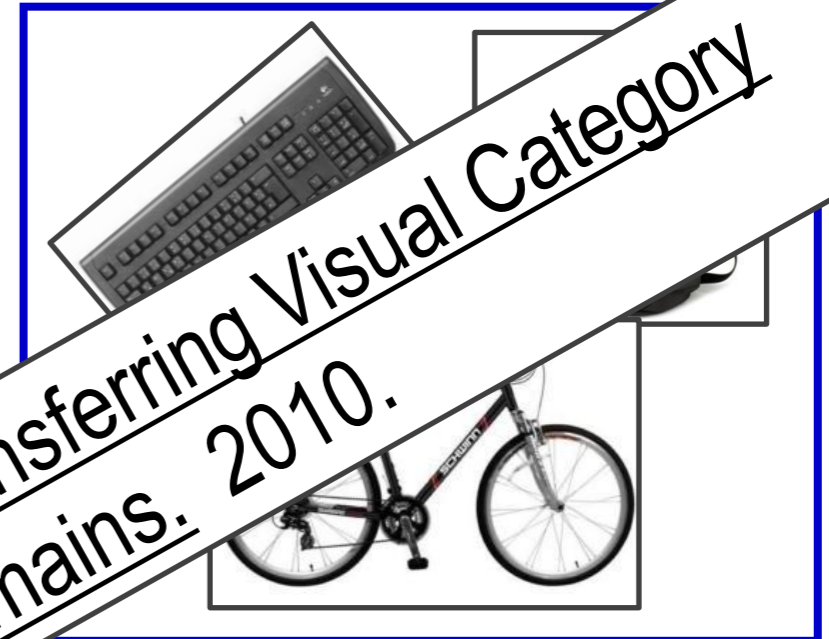
Packed with fascinating info



A breeze to clean up



## Visual Object Recognition



K. Saenko et al. Transferring Visual Category Models to New Domains. 2010.





# Tutorial Outline

---



1. Domain Adaptation: Common Concepts
2. Semi-supervised Adaptation
  - Learning with Shared Support
  - Learning Shared Representations
3. Supervised Adaptation
  - Feature-Based Approaches
  - Parameter-Based Approaches
4. Open Questions and Uncovered Algorithms



# Classical vs Adaptation Error



## Classical Test Error:

$$\epsilon_{\text{test}} \leq \hat{\epsilon}_{\text{train}} + \sqrt{\frac{\text{complexity}}{n}}$$

Measured on the  
same distribution!

## Adaptation Target Error:

$$\epsilon_{\text{test}} \leq ??$$

Measured on a  
new distribution!



# Common Concepts in Adaptation



## Covariate Shift

$$\Pr_S[y|x] = \Pr_T[y|x]$$



understands both



&



## Single Good Hypothesis

$$\exists h^*, \epsilon_S(h^*), \epsilon_T(h^*) \text{ small}$$



understands both

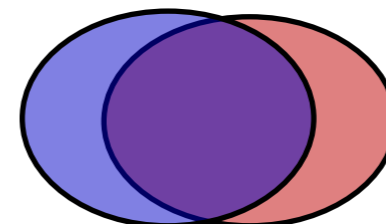


&

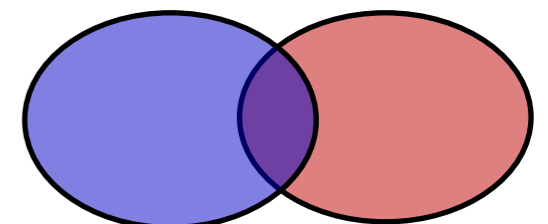


## Domain Discrepancy and Error

Easy



Hard







# Tutorial Outline

---



1. Notation and Common Concepts
2. Semi-supervised Adaptation
  - Covariate shift with Shared Support
  - Learning Shared Representations
3. Supervised Adaptation
  - Feature-Based Approaches
  - Parameter-Based Approaches
4. Open Questions and Uncovered Algorithms



# A bound on the adaptation error

---



Let  $h$  be a binary hypothesis. If  $\Pr_S(y|x) = \Pr_T(y|x)$ , then

$$\epsilon_T(h) \leq \epsilon_S(h) + \int_{\mathcal{X}} |\Pr_T(x) - \Pr_S(x)| dx$$

Minimize the **total variation**



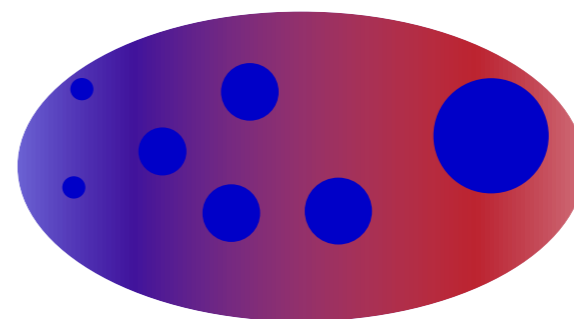
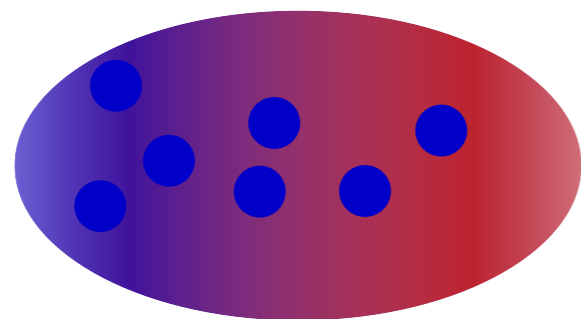
# Covariate Shift with Shared Support



Assumption: Target & Source Share Support

$$\forall x, \Pr_S[x] \neq 0 \text{ iff } \Pr_T[x] \neq 0$$

Reweight source instances to minimize discrepancy





# Source Instance Reweighting



## Defining Error

$$\epsilon_T(h) = \mathbb{E}_{\text{Pr}_T[x]} \mathbb{E}_{\text{Pr}[y|x]} [h(x) \neq y]$$

## Using Definition of Expectation

$$= \sum_x \text{Pr}_T[x] \mathbb{E}_{\text{Pr}[y|x]} [h(x) \neq y]$$

## Multiplying by 1

$$= \sum_x \frac{\text{Pr}_S[x]}{\text{Pr}_S[x]} \text{Pr}_T[x] \mathbb{E}_{\text{Pr}[y|x]} [h(x) \neq y]$$

per-instance weights  $w$

## Rearranging

$$\epsilon_T(h) = \epsilon_S(h, w) = \mathbb{E}_{\text{Pr}_S[x]} \frac{\text{Pr}_t[x]}{\text{Pr}_s[x]} \mathbb{E}_{\text{Pr}[y|x]} [h(x) \neq y]$$



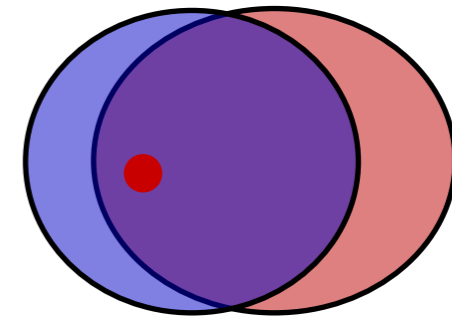
# Sample Selection Bias

---



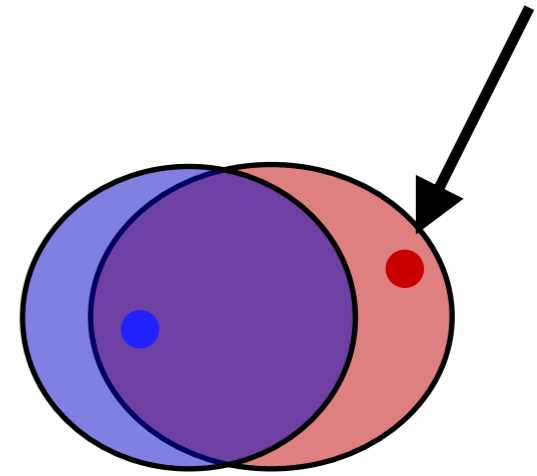
## Another Way to View

- 1) Draw from the target  $\Pr_T[x]$



## Redefine the source distribution

- 1) Draw from the target  $\Pr_T[x]$
- 2) Select into the source with  $\Pr[\sigma = 1|x]$



$$\Pr_S[x] = \frac{\Pr_T[x]\Pr[\sigma = 1|x]}{\Pr[\sigma = 1]}$$



# Rewriting Source Risk



$$\Pr_S[x] = \frac{\Pr_T[x] \Pr[\sigma = 1|x]}{\Pr[\sigma = 1]}$$

## Rearranging

$$\frac{\Pr_T[x]}{\Pr_S[x]} = \frac{\Pr[\sigma = 1]}{\Pr[\sigma = 1|x]}$$

per-instance  
weights  $w$

$\Pr[\sigma = 1]$  not dependent on  $x$

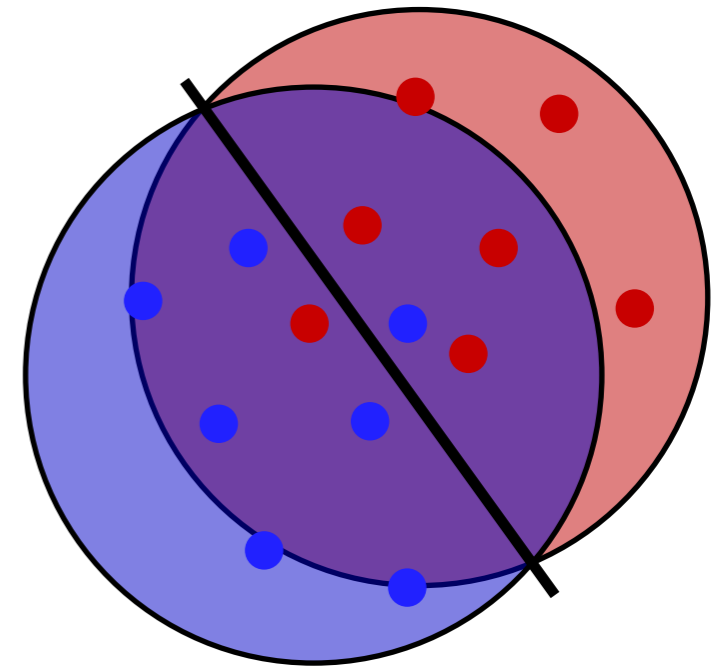
$$\epsilon_S(h, w) \propto \mathbb{E}_{\Pr_S[x]} \left( \frac{1}{\Pr[\sigma = 1|x]} \right) \mathbb{E}_{\Pr[y|x]} [h(x) \neq y]$$



# Logistic Model of Source Selection



$$\Pr[\sigma = 1 | x] = \frac{1}{1 + \exp(\theta^\top x + b)}$$



## Training Data

Source instances,  $\sigma = 1$

Target unlabeled instances,  $\sigma = 0$



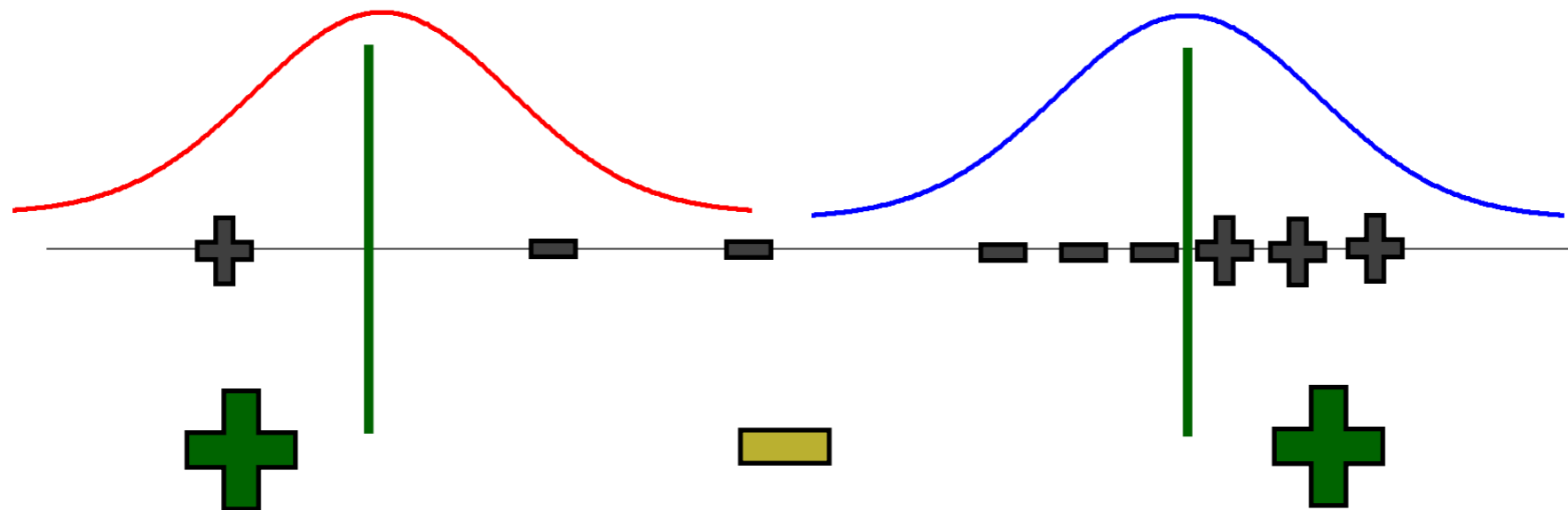


# Selection Bias Correction Algorithm



Input:

Labeled **source** data





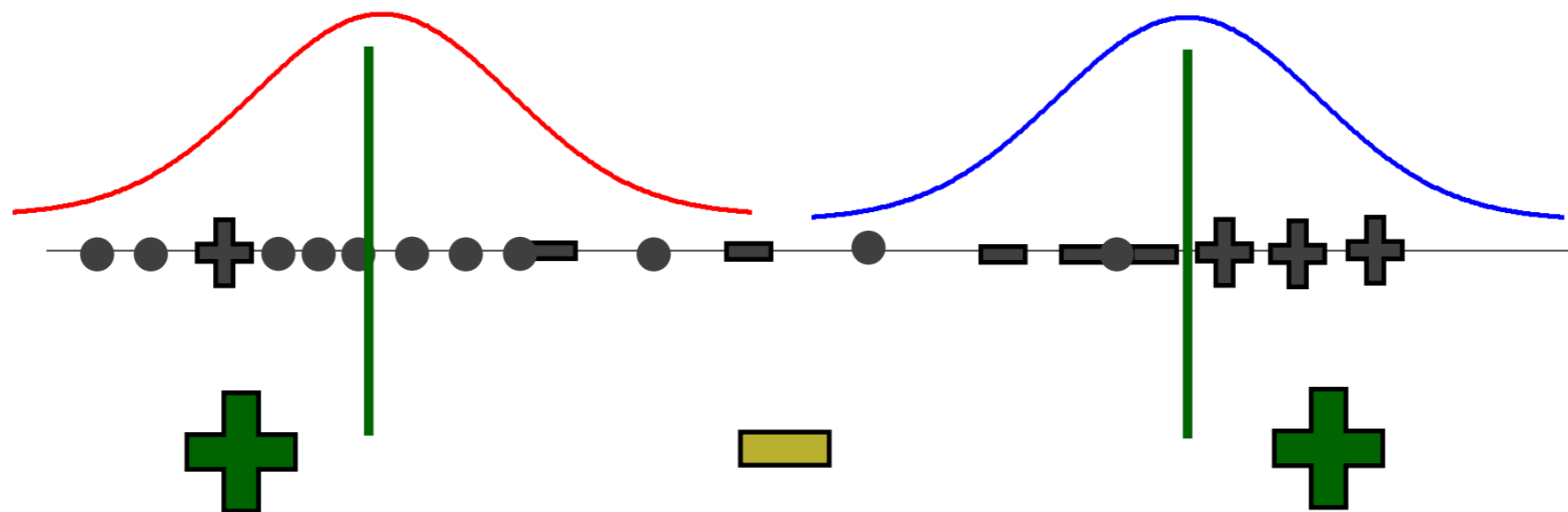
# Selection Bias Correction Algorithm



Input:

Labeled **source** data

Unlabeled **target** data



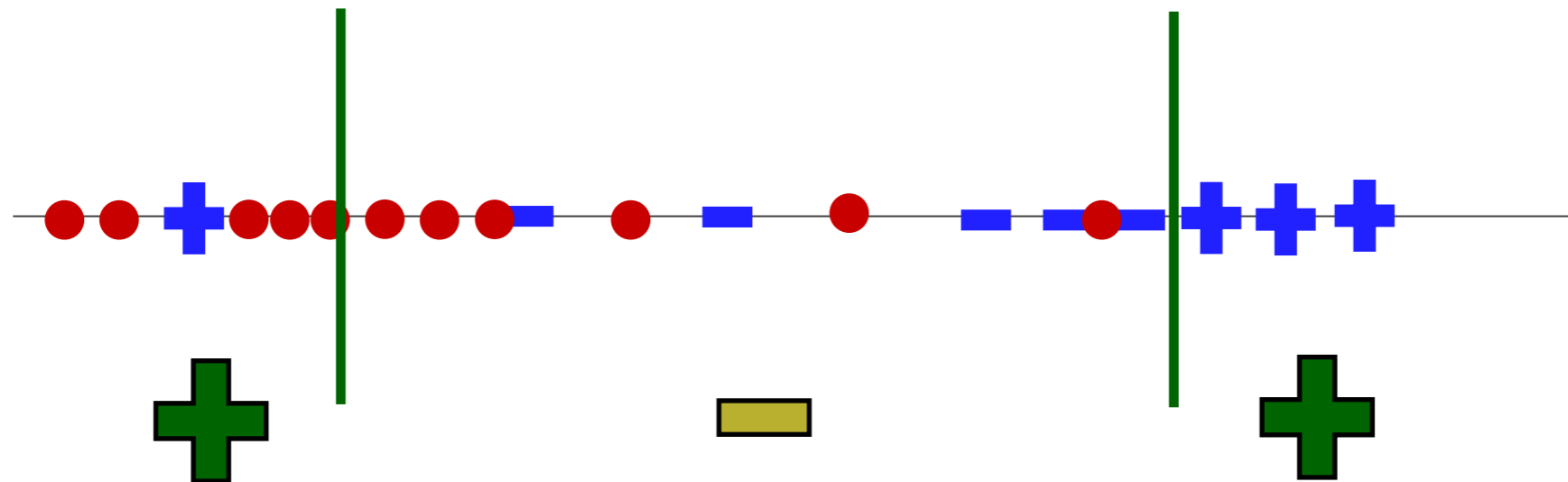


# Selection Bias Correction Algorithm



Input: Labeled **source** and unlabeled **target** data

1) Label source instances as  $\sigma = 1$ , target as  $\sigma = 0$



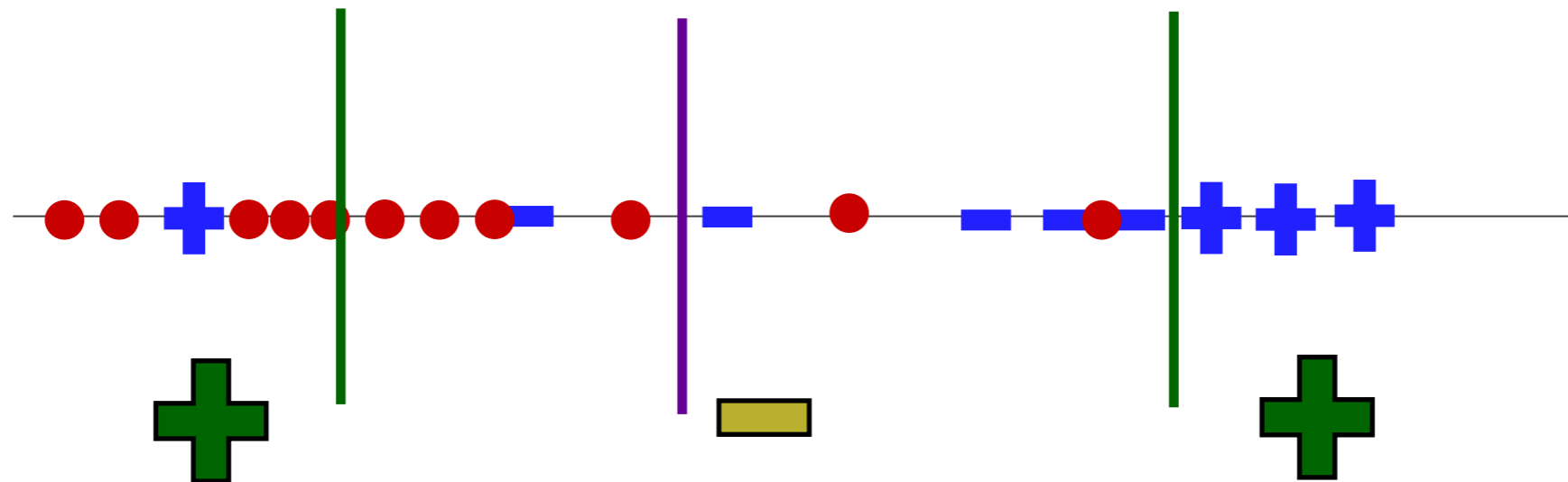


# Selection Bias Correction Algorithm



Input: Labeled **source** and unlabeled **target** data

- 1) Label source instances as  $\sigma = 1$ , target as  $\sigma = 0$
- 2) Train predictor  $\Pr[\sigma = 1|x] = \frac{1}{1+\exp(\theta^\top x+b)}$



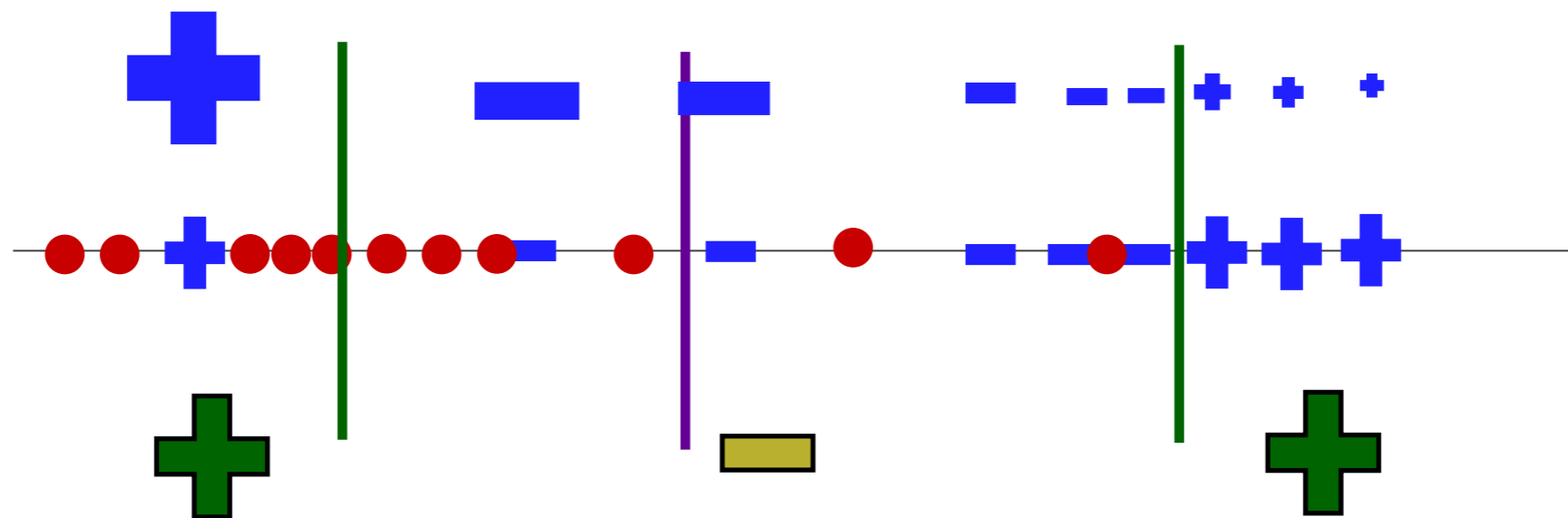


# Selection Bias Correction Algorithm



Input: Labeled **source** and unlabeled **target** data

- 1) Label source instances as  $\sigma = 1$ , target as  $\sigma = 0$
- 2) Train predictor  $\Pr[\sigma = 1|x] = \frac{1}{1 + \exp(\theta^\top x + b)}$
- 3) Reweight source instances



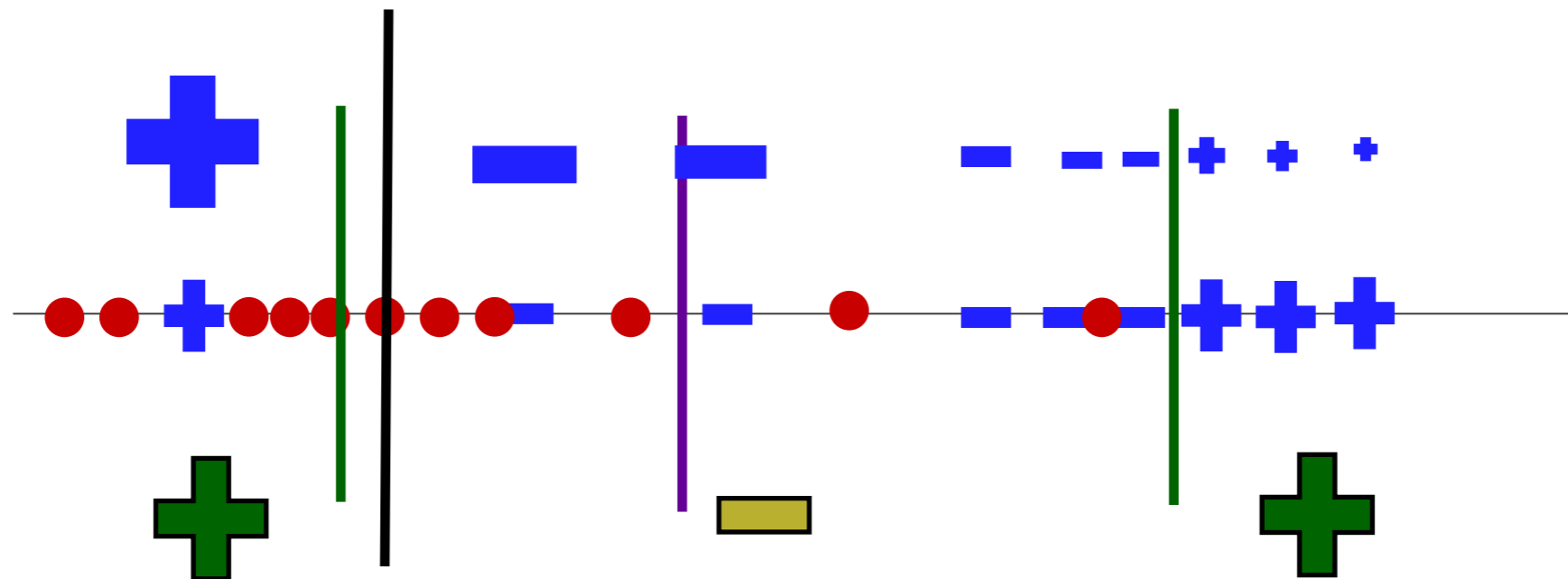


# Selection Bias Correction Algorithm



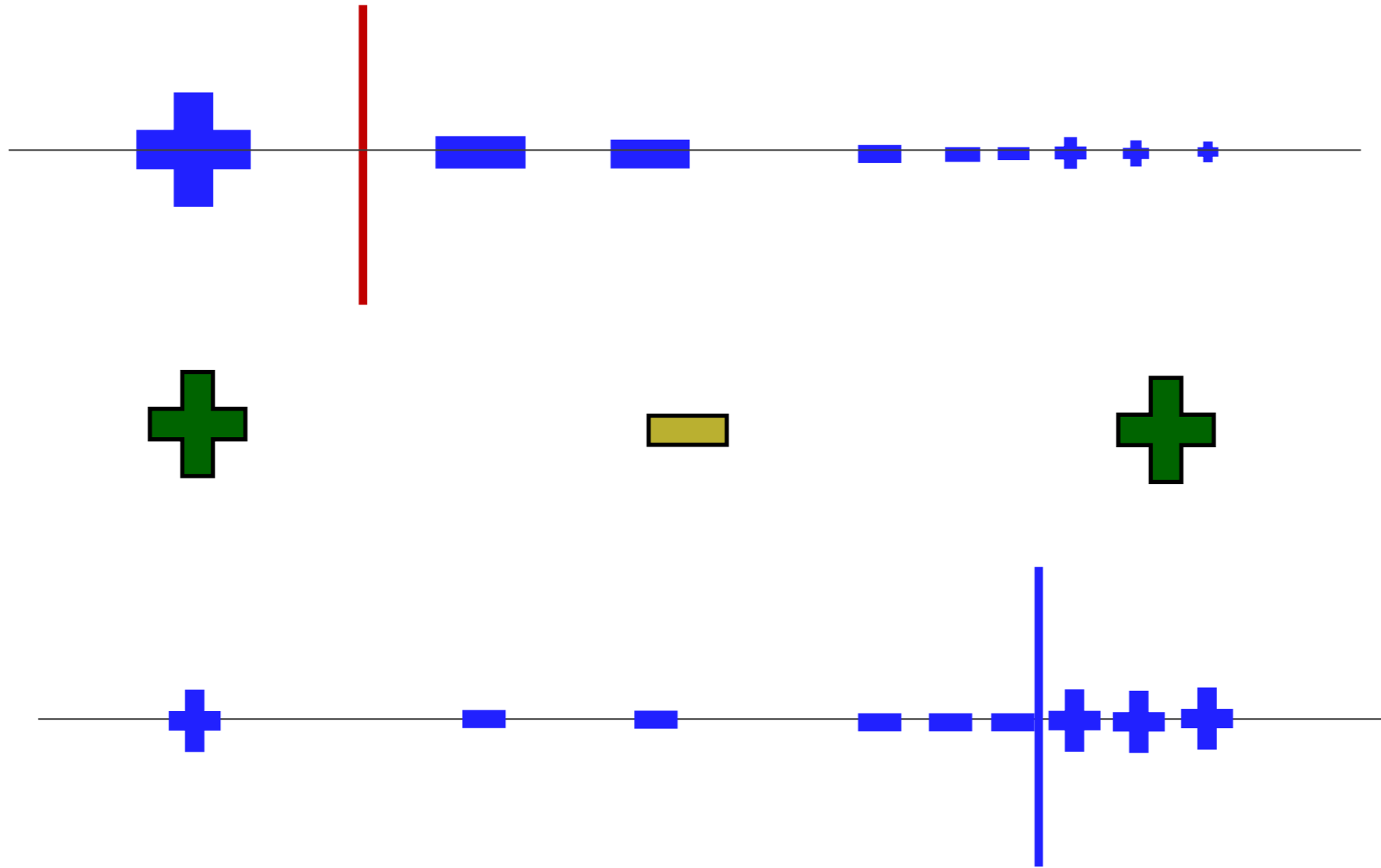
Input: Labeled **source** and unlabeled **target** data

- 1) Label source instances as  $\sigma = 1$ , target as  $\sigma = 0$
- 2) Train predictor  $\Pr[\sigma = 1|x] = \frac{1}{1+\exp(\theta^\top x+b)}$
- 3) Reweight source instances
- 4) Train target predictor





# How Bias gets Corrected





# Rates for Re-weighted Learning



$\hat{\epsilon}_S^n(h, w)$ : weighted source test error on sample of size  $n$

With probability  $1 - \delta$ , for every  $h$

$$|\hat{\epsilon}_S^n(h, w) - \epsilon_T(h)| \leq \sqrt{\frac{O\left(\frac{1}{\delta}\right) + O\left(\max_{x \in \mathcal{X}} w(x)^2\right)}{n}}$$

Adapted from Gretton et al.





# Sample Selection Bias Summary



## Two Key Assumptions

- 1) Covariate shift:  $\Pr_S[y|x] = \Pr_T[y|x]$
- 2) Shared support:  $\forall x, \Pr_S[x] \neq 0$  iff  $\Pr_T[x] \neq 0$

Advantage

$$\hat{\epsilon}_S^n(h, w) \xrightarrow{\infty} \epsilon_T(h)$$

Optimal target predictor  
without labeled target data



# Sample Selection Bias Summary



## Two Key Assumptions

- 1) Covariate shift:  $\Pr_S[y|x] = \Pr_T[y|x]$
- 2) Shared support:  $\forall x, \Pr_S[x] \neq 0$  iff  $\Pr_T[x] \neq 0$

Advantage  $\hat{\epsilon}_S^n(h, w) \xrightarrow[n]{\infty} \epsilon_T(h)$

## Disadvantage

Convergence to  $\epsilon_T(h)$  depends on  $\max_x \frac{\Pr_T(x)}{\Pr_S(x)}$



# Sample Selection Bias References

---



<http://adaptationtutorial.blitzer.com/references/>

- [1] J. Heckman. Sample Selection Bias as a Specification Error. 1979.
- [2] A. Gretton et al. Covariate Shift by Kernel Mean Matching. 2008.
- [3] C. Cortes et al. Sample Selection Bias Correction Theory. 2008
- [4] S. Bickel et al. Discriminative Learning Under Covariate Shift. 2009.



# Tutorial Outline

---



1. Notation and Common Concepts
2. Semi-supervised Adaptation
  - Covariate shift
  - Learning Shared Representations
3. Supervised Adaptation
  - Feature-Based Approaches
  - Parameter-Based Approaches
4. Open Questions and Uncovered Algorithms



# Unshared Support in the Real World



## Running with Scissors

**Title:** Horrible book, horrible.

This book was horrible. I read half, suffering from a headache the entire time, and eventually i lit it on fire. 1 less copy in the world. Don't waste your money. I wish i had the time spent reading this book back. It wasted my life

## Avante Deep Fryer; Black

**Title:** lid does not work well...

I love the way the Tefal deep fryer cooks, however, I am returning my second one due to a defective lid closure. The lid may close initially, but after a few uses it no longer stays closed. I won't be buying this one again.





# Unshared Support in the Real World



amazon.com

**Running with Scissors**

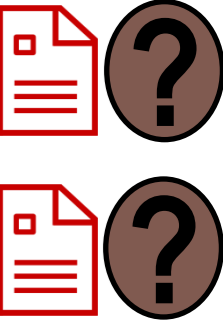
**Title:** Horrible book, horrible.

This book was horrible. I read half

**Avante Deep Fryer; Black**

**Title:** lid **does not work** well...

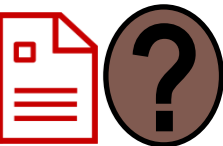
I love the way the Tefal deep fryer



**Error increase: 13% → 26%**

time, and eventually i lit it on fire. 1 less copy in the world. Don't waste your money. I wish i had the time spent reading this book back. It wasted my life

second one due to a **defective** lid closure. The lid may close initially, but after a few uses it no longer stays closed. I won't be buying this one again.





# Linear Regression for Rating Prediction



$$h(x) = \text{sgn}(\theta^\top x) \quad h(x) \in \left\{ \text{👍} \text{👎} \right\}$$



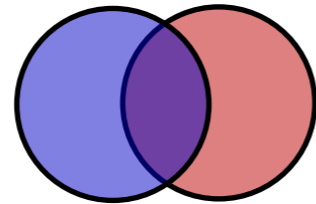
	<i>excellent</i>			<i>great</i>			<i>fascinating</i>		
$x$	3	0	...	0	1	0	...	0	1
$\theta$	0.5	1	...	-1.1	2	-0.3	...	0.1	1.5



# Coupled Subspaces

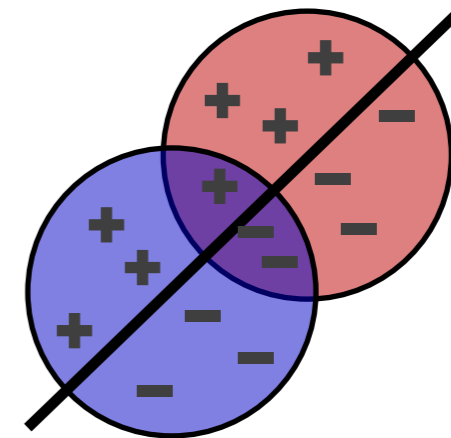


No Shared Support



Single Good Linear Hypothesis

$\exists \theta^*, \epsilon_S(\theta^*) + \epsilon_T(\theta^*)$  small



Stronger than  $\Pr_S[y|x] = \Pr_T[y|x]$

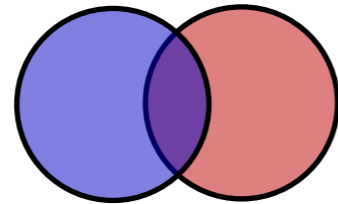




# Coupled Subspaces

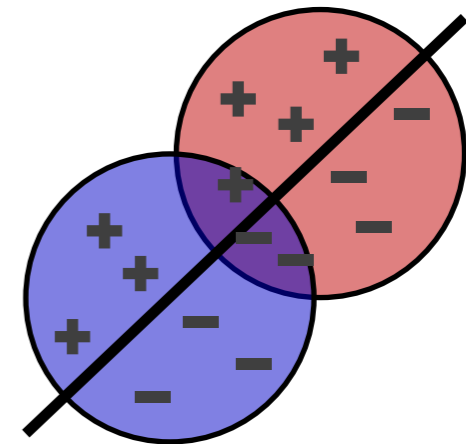


No Shared Support



Single Good Linear Hypothesis

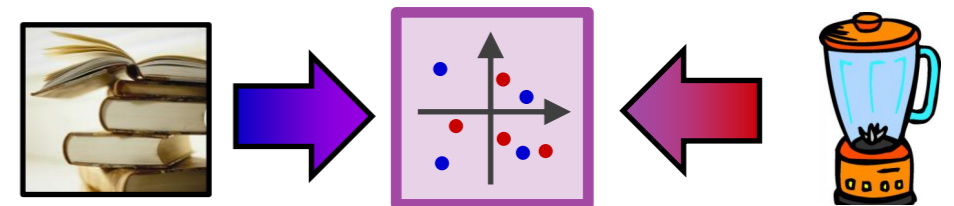
$$\exists \theta^*, \quad \epsilon_S(\theta^*) + \epsilon_T(\theta^*) \quad \text{small}$$



Coupled Representation Learning

$Px$  couples domains

Bound target error  $\epsilon_{P,T}(\theta)$





# Single Good Linear Hypothesis?



$$\exists \theta^*, \quad \epsilon_S(\theta^*) + \epsilon_T(\theta^*) \quad \text{small}$$

## Adaptation Squared Error

<b>Source</b> \ <b>Target</b>	<b>Books</b>	<b>Kitchen</b>
<b>Books</b>	<b>1.35</b>	
<b>Kitchen</b>		<b>1.19</b>
<b>Both</b>		



# Single Good Linear Hypothesis?



$$\exists \theta^*, \quad \epsilon_S(\theta^*) + \epsilon_T(\theta^*) \quad \text{small}$$

## Adaptation Squared Error

<b>Target</b> <b>Source</b>	<b>Books</b>	<b>Kitchen</b>
Books	<b>1.35</b>	
Kitchen		<b>1.19</b>
Both	<b>1.38</b>	<b>1.23</b>



# Single Good Linear Hypothesis?



$$\exists \theta^*, \quad \epsilon_S(\theta^*) + \epsilon_T(\theta^*) \quad \text{small}$$

## Adaptation Squared Error

<b>Target</b> <b>Source</b>	<b>Books</b>	<b>Kitchen</b>
Books	<b>1.35</b>	1.68
Kitchen	1.80	<b>1.19</b>
Both	<b>1.38</b>	<b>1.23</b>



# A bound on the adaptation error



Let  $h$  be a binary hypothesis. If  $\Pr_S(Y|x) = \Pr_T(Y|x)$ , then

$$\epsilon_T(h) \leq \epsilon_S(h) + \int_{\mathcal{X}} |\Pr_T(x) - \Pr_S(x)| dx$$

What if a single good hypothesis exists?

A better discrepancy than total variation?



# A generalized discrepancy distance



Measure how hypotheses make mistakes





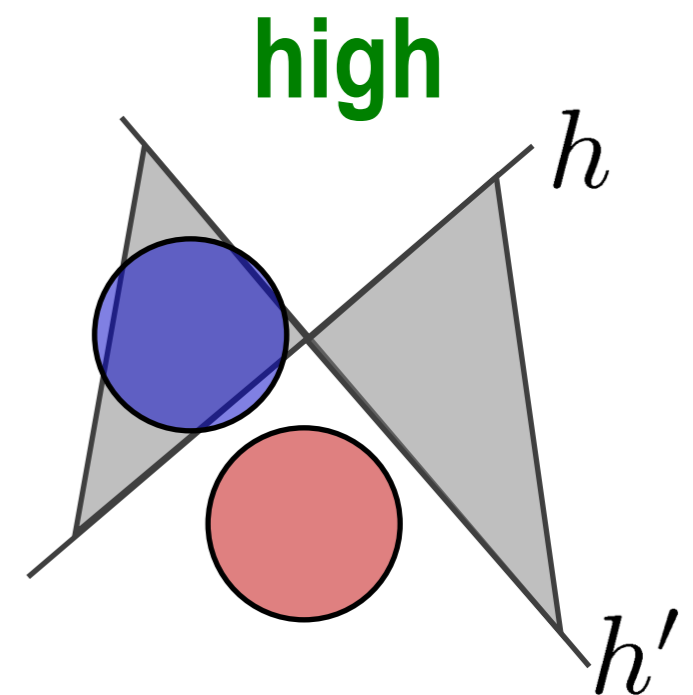
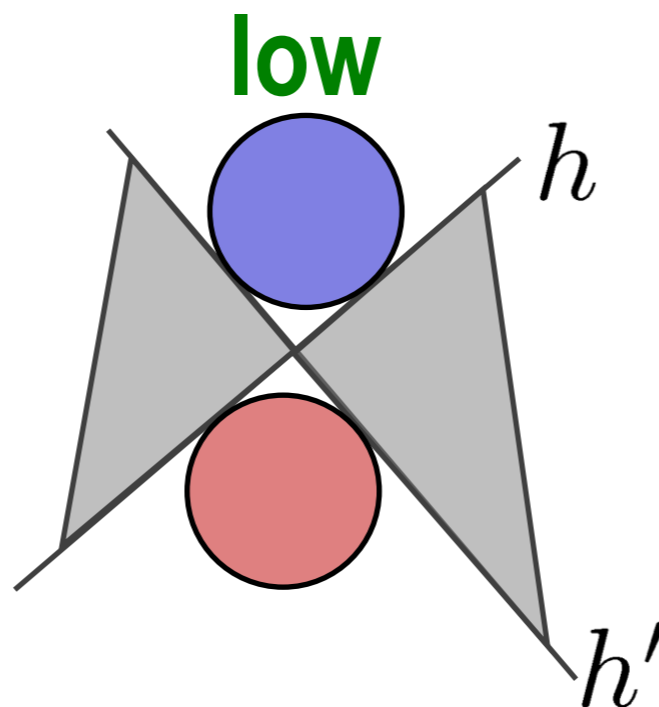
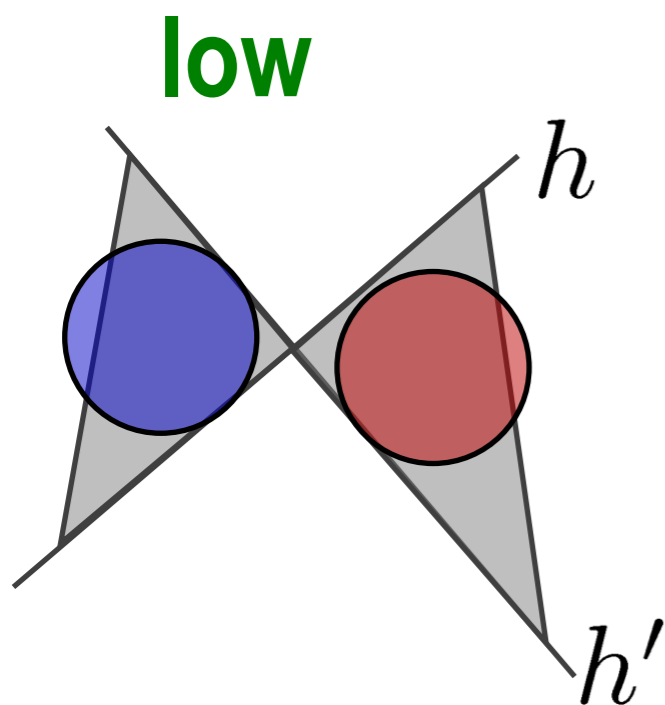
# A generalized discrepancy distance



Measure how hypotheses make mistakes

$$\text{disc}_H(Q, P) =$$

$$\max_{h, h' \in H} |E_Q[h(x) \neq h'(x)] - E_P[h(x) \neq h'(x)]|$$



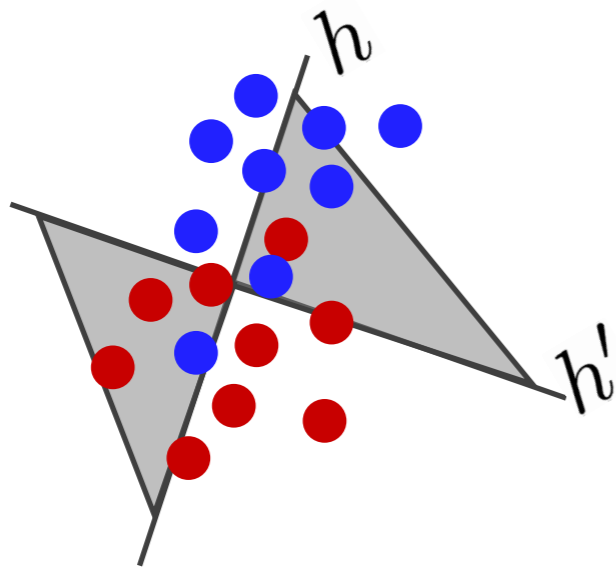


# Discrepancy vs. Total Variation



## Discrepancy

Computable from finite samples.



## Total Variation

Not computable in general



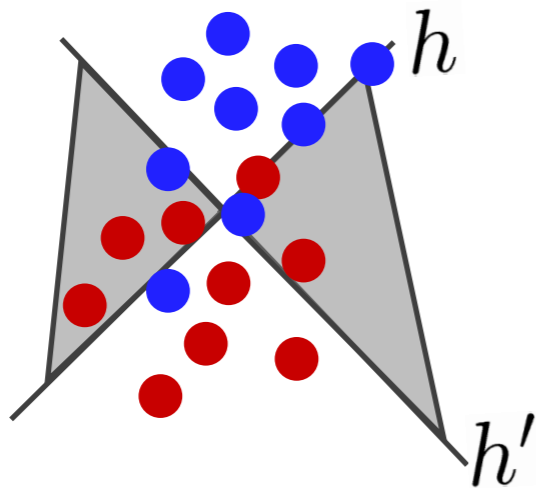


# Discrepancy vs. Total Variation



## Discrepancy

Computable from finite samples.



## Total Variation

Not computable in general

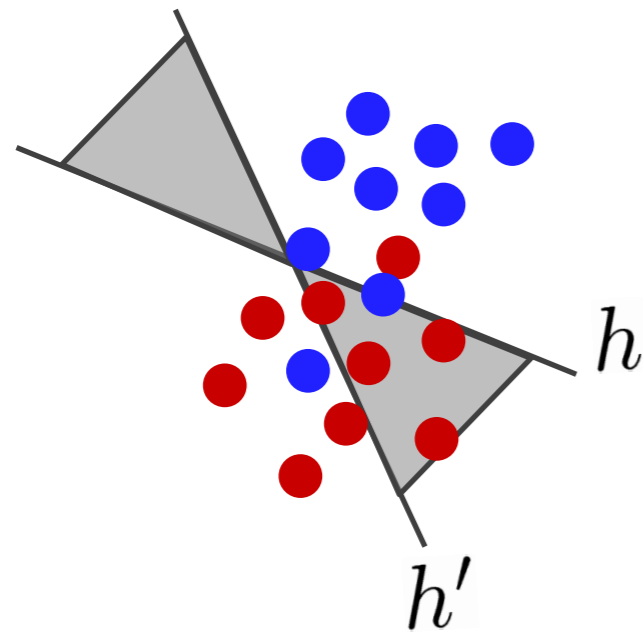


# Discrepancy vs. Total Variation



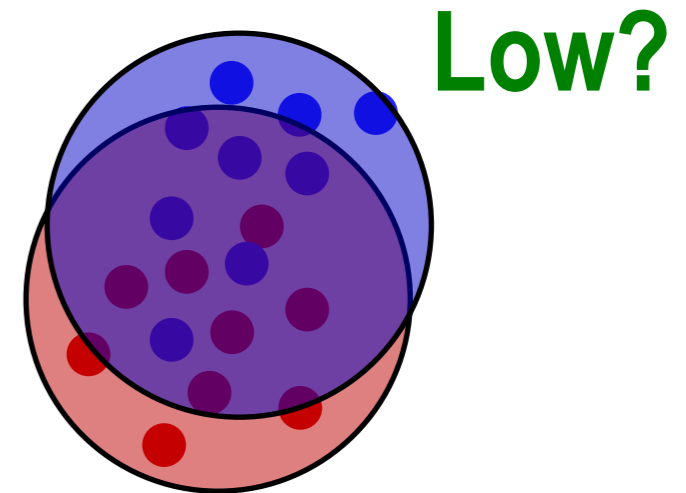
## Discrepancy

Computable from finite samples.



## Total Variation

Not computable in general



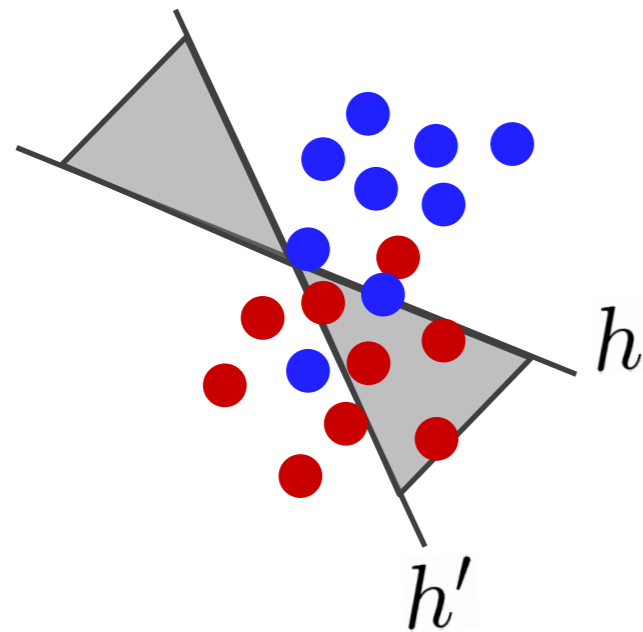


# Discrepancy vs. Total Variation



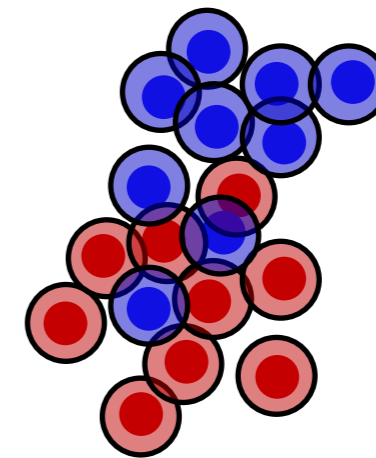
## Discrepancy

Computable from finite samples.



## Total Variation

Not computable in general



High?

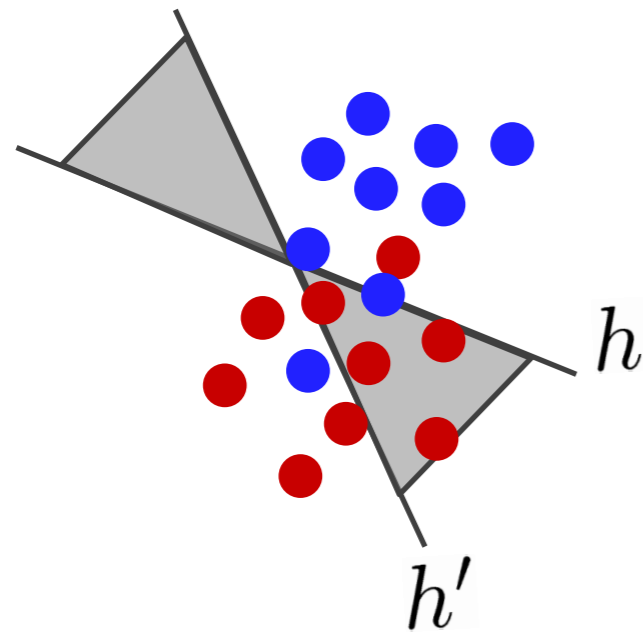


# Discrepancy vs. Total Variation



## Discrepancy

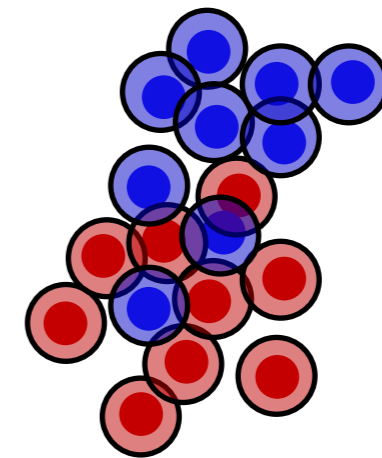
Computable from finite samples.



Related to hypothesis class

## Total Variation

Not computable in general



High?

Unrelated to hypothesis class

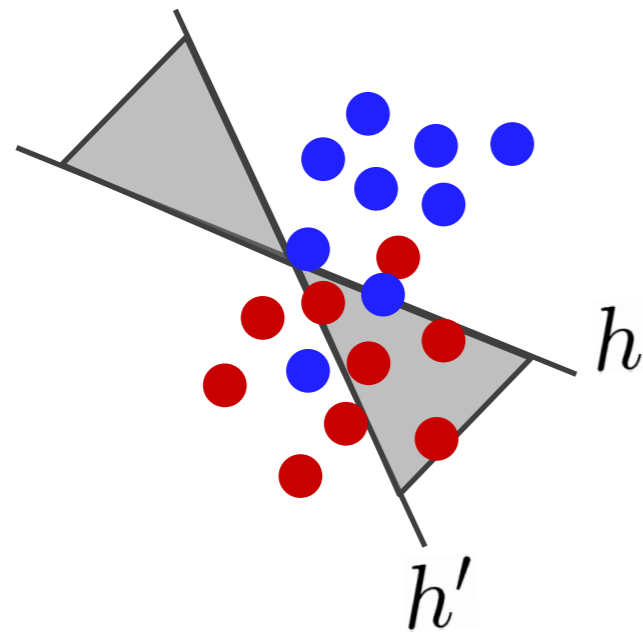


# Discrepancy vs. Total Variation



## Discrepancy

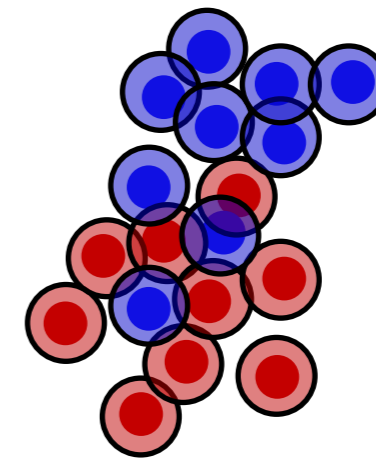
Computable from finite samples.



Related to hypothesis class

## Total Variation

Not computable in general



High?

Unrelated to hypothesis class

Bickel covariate shift algorithm heuristically minimizes both measures



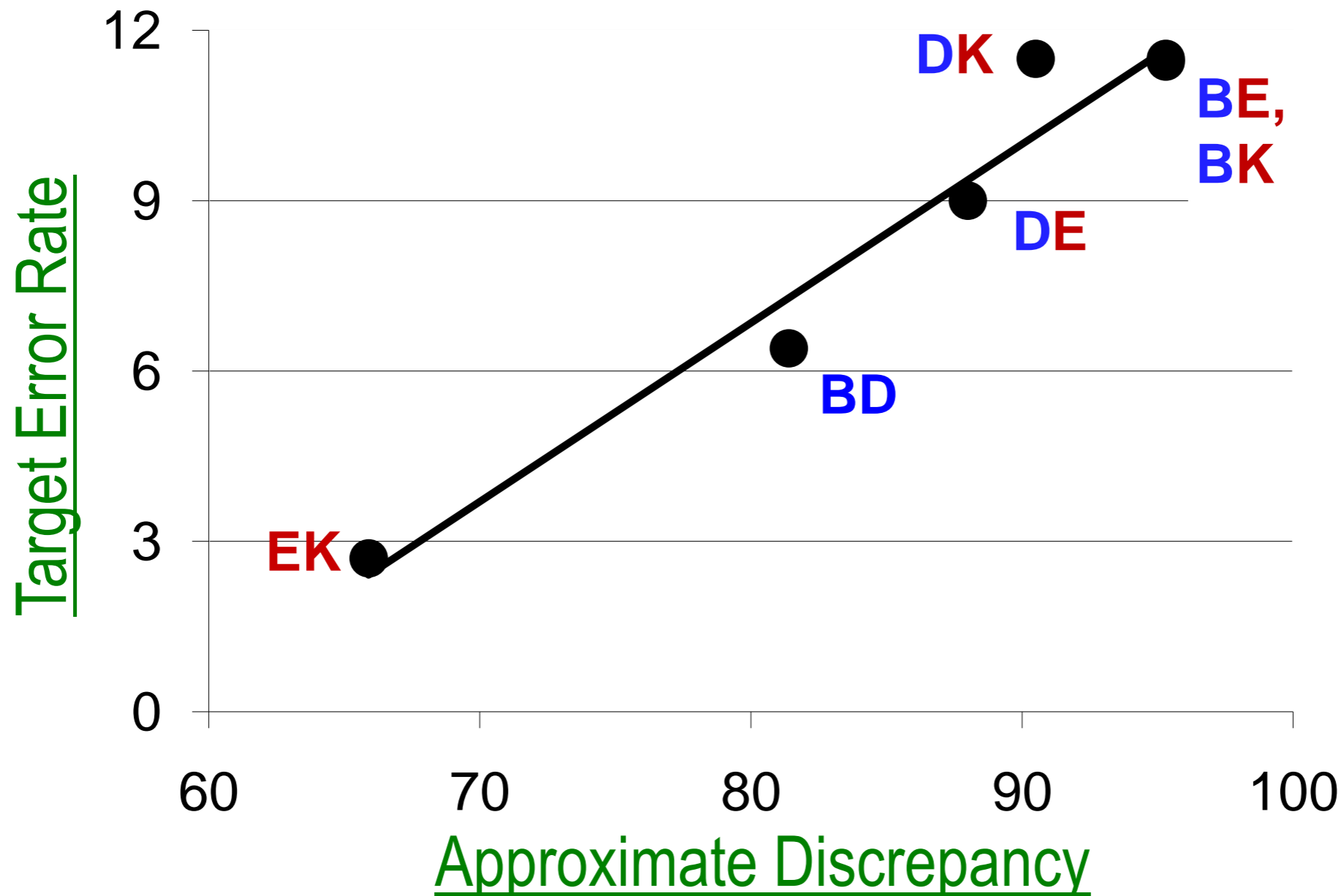
# Is Discrepancy Intuitively Correct?



4 domains: **B**ooks, **D**VDs, **E**lectronics, **K**itchen    **B&D**, **E&K**    Shared Vocabulary

**B&D**: *fascinating, boring*

**E&K**: *super easy, bad quality*





# A new adaptation bound



$S, T$ : Source and target     $\mathcal{H}$ : Hypothesis class     $n$ : Sample size

$\hat{S}$ : Labeled  $S$  sample     $\hat{T}$ : Unlabeled  $T$  sample

$\mathcal{R}_{\hat{S}}(\mathcal{H}), \mathcal{R}_{\hat{T}}(\mathcal{H})$ : Rademacher complexities

With probability  $1 - \delta$ , for  $h$  the ERM of  $\hat{S}$ :

$$\begin{aligned} \epsilon_T(h) - \epsilon_T(h^*) &\leq \epsilon_{\hat{S}}(h, h^*) + O(\mathcal{R}_{\hat{S}}(\mathcal{H}) + \mathcal{R}_{\hat{T}}(\mathcal{H})) \\ &\quad + O\left(\sqrt{\frac{\log \frac{1}{\delta}}{n}}\right) + \text{disc}_{\mathcal{H}}(\hat{S}, \hat{T}) \end{aligned}$$



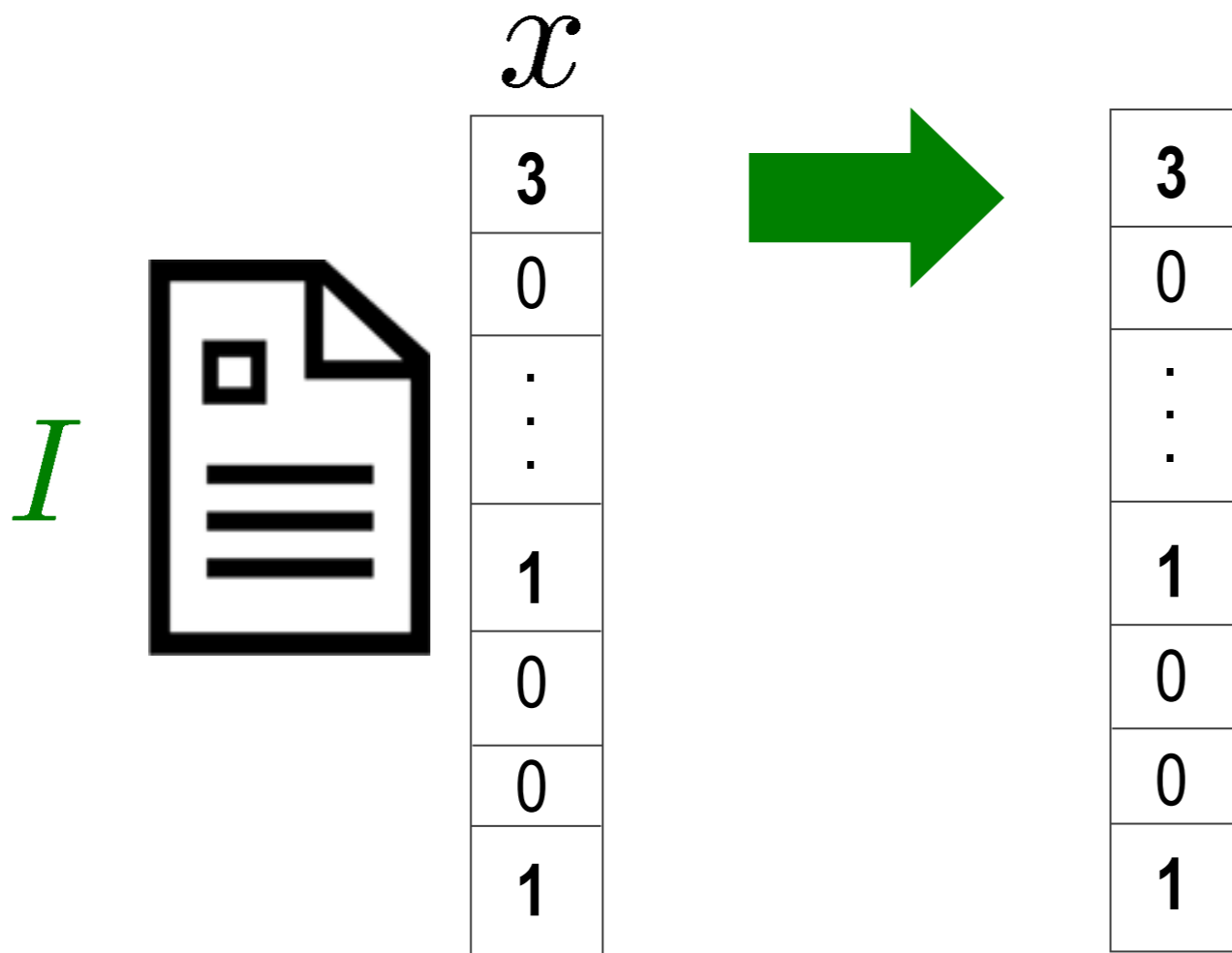
# Representations and the Bound



Linear Hypothesis Class:  $h(x) = \text{sgn}(\theta^\top x)$

$$P = I$$

Hypothesis classes from projections  $P: \theta^\top Px$







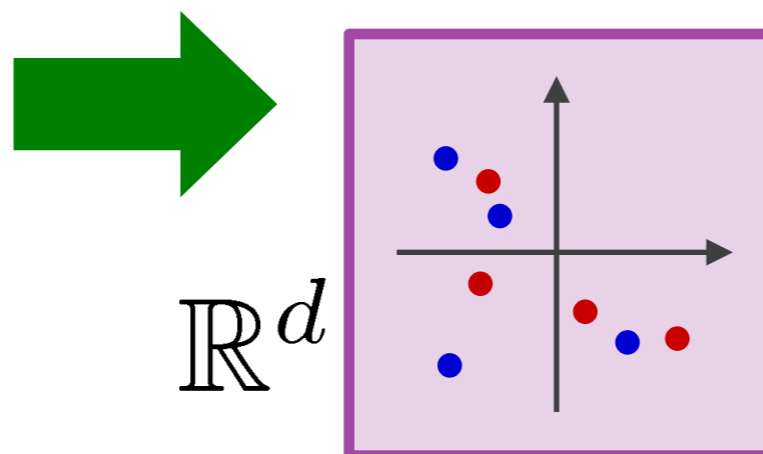
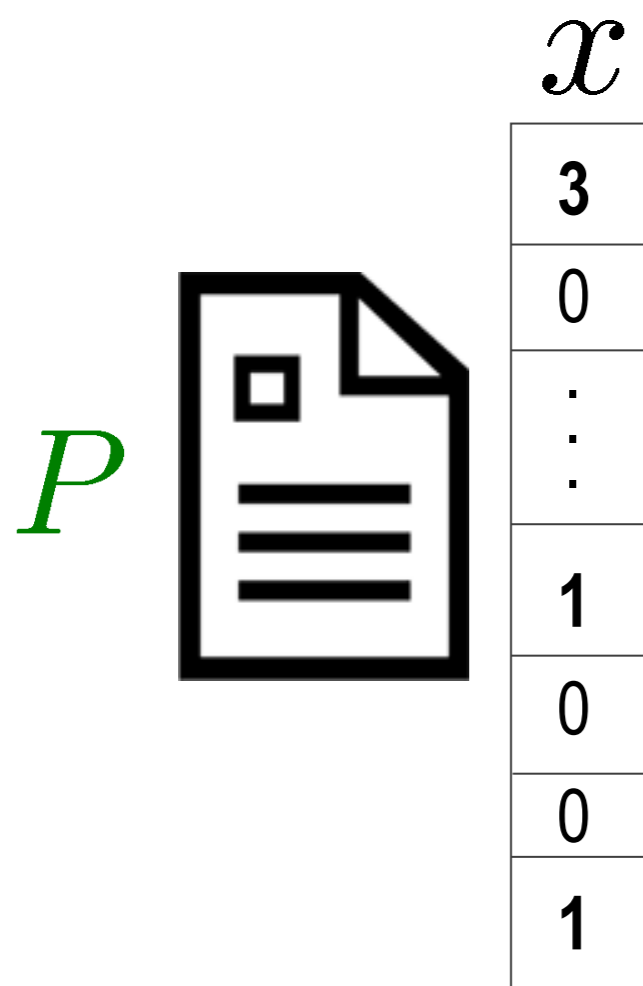
# Representations and the Bound



Linear Hypothesis Class:  $h(x) = \text{sgn}(\theta^\top x)$

$P$

Hypothesis classes from projections  $P$ :  $\theta^\top Px$



Goals for  $P$

- 1) Minimize divergence
- 2)  $\epsilon_{P,T}(\theta^*) - \epsilon_{I,T}(\theta^*)$  small

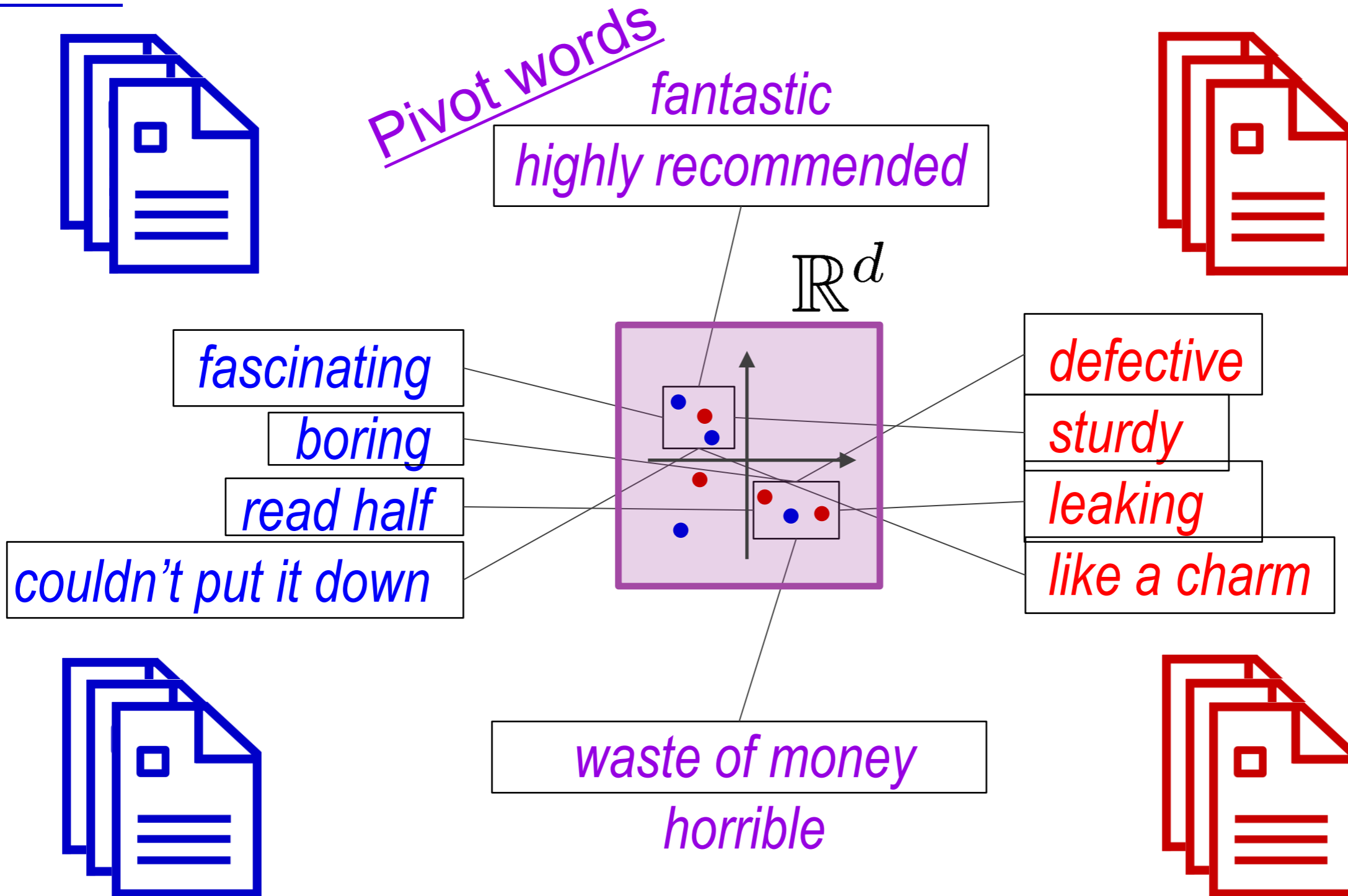


# Learning Representations: Pivots



Source

Target

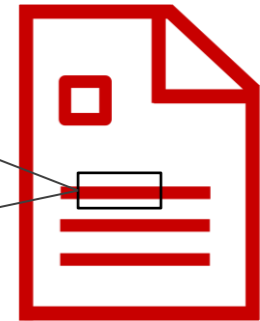




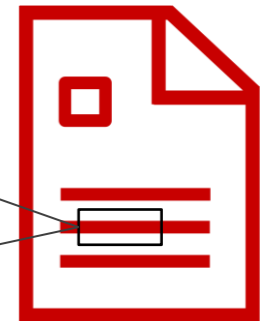
# Predicting pivot word presence



Do **not buy**

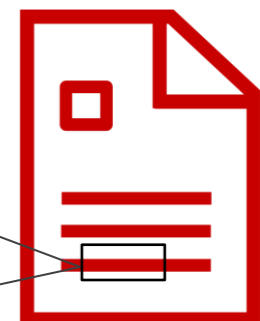


An absolutely **great** purchase



⋮

A **sturdy** deep fryer

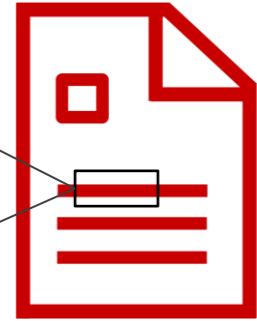




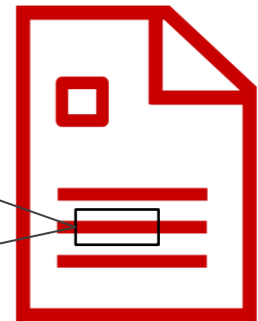
# Predicting pivot word presence



Do **not buy** the Shark portable steamer.  
The trigger mechanism is **defective**.

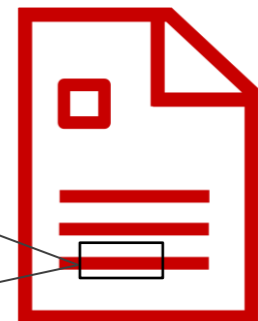


An absolutely **great** purchase



⋮

A **sturdy** deep fryer

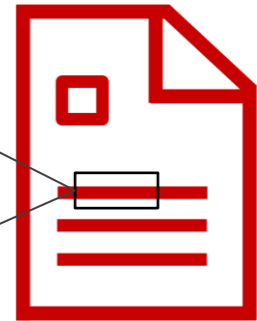




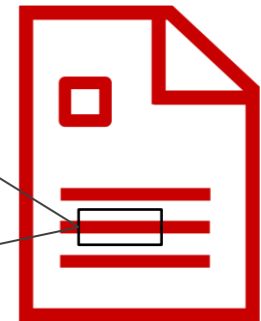
# Predicting pivot word presence



Do **not buy** the Shark portable steamer.  
The trigger mechanism is **defective**.



An absolutely **great** purchase. . . . This  
blender is incredibly **sturdy**.

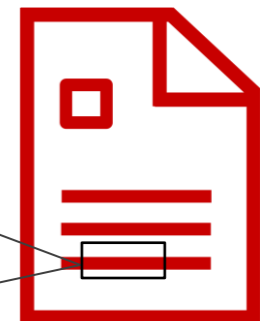


## Predict presence of pivot words

$$P_{w(\textit{great})}(\textit{great} | x) \propto \exp \{ \langle x, w(\textit{great}) \rangle \}$$

⋮

A **sturdy** deep fryer



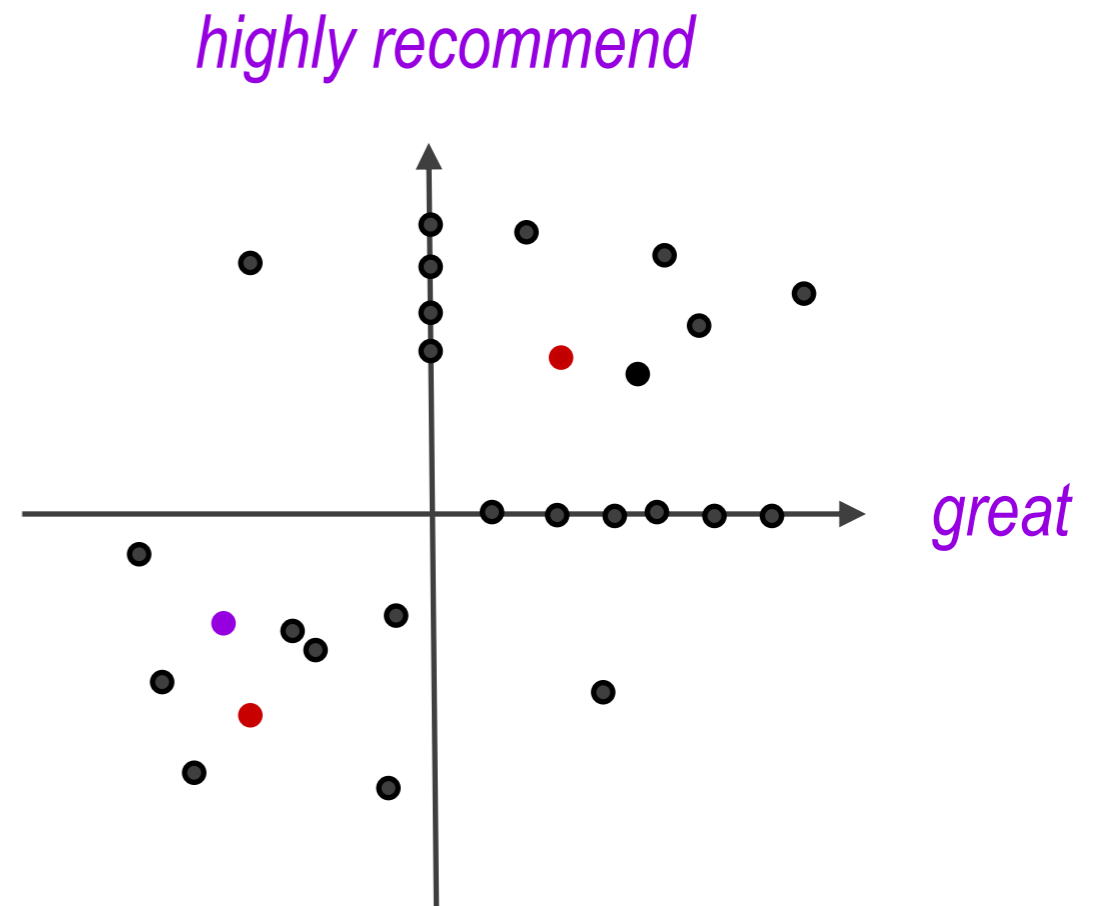


# Finding a shared sentiment subspace



$$W = \begin{bmatrix} | & & | & & | \\ w_1 & \dots & w(\text{highly recommend}) & \dots & w_N \\ | & & | & & | \end{bmatrix}$$

- $p_W(\text{pivots} | x)$  generates N new features
- $p_{w(\text{highly recommend})}(\text{highly recommend} | x)$  : “Did *highly recommend* appear?”
- Sometimes predictors capture non-sentiment information



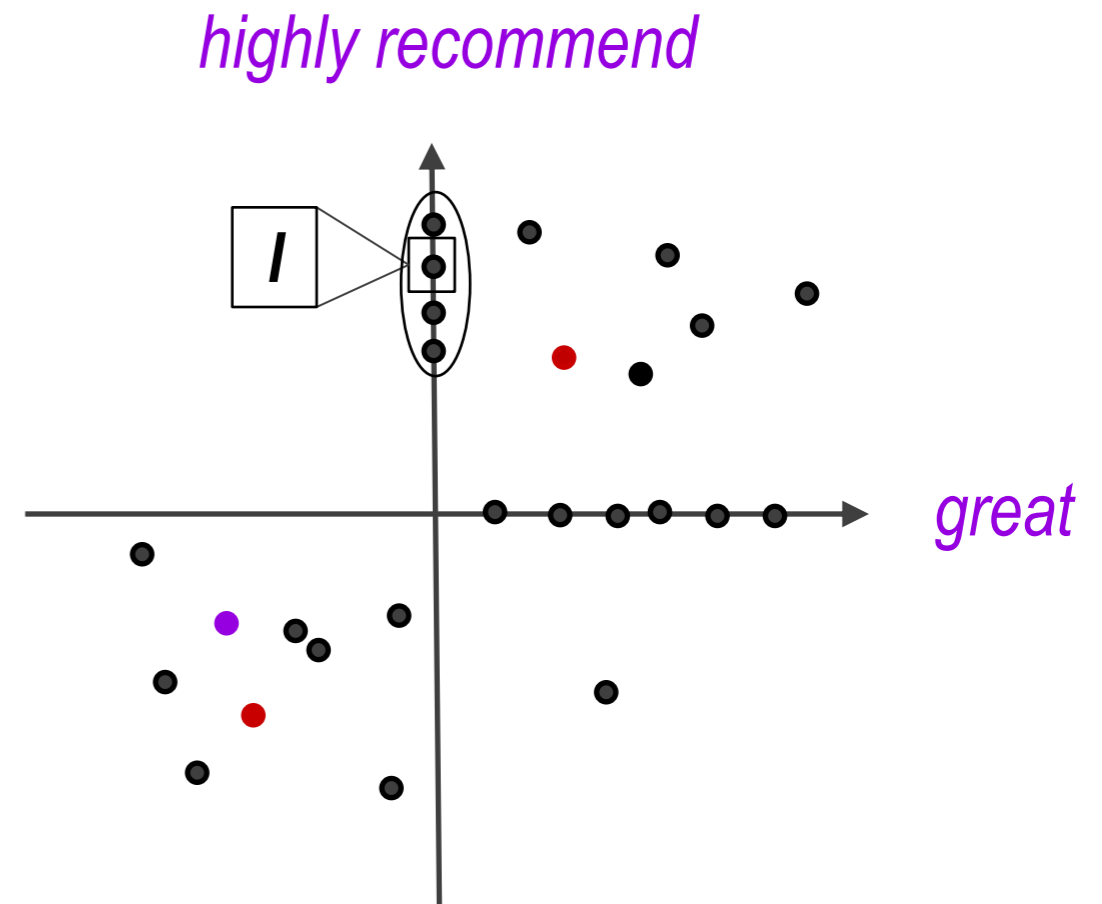


# Finding a shared sentiment subspace



$$W = \begin{bmatrix} | & & | & & | \\ w_1 & \dots & w(\text{highly recommend}) & \dots & w_N \\ | & & | & & | \end{bmatrix}$$

- $p_W(\text{pivots} | x)$  generates N new features
- $p_{w(\text{highly recommend})}(\text{highly recommend} | x)$  : “Did *highly recommend* appear?”
- Sometimes predictors capture non-sentiment information



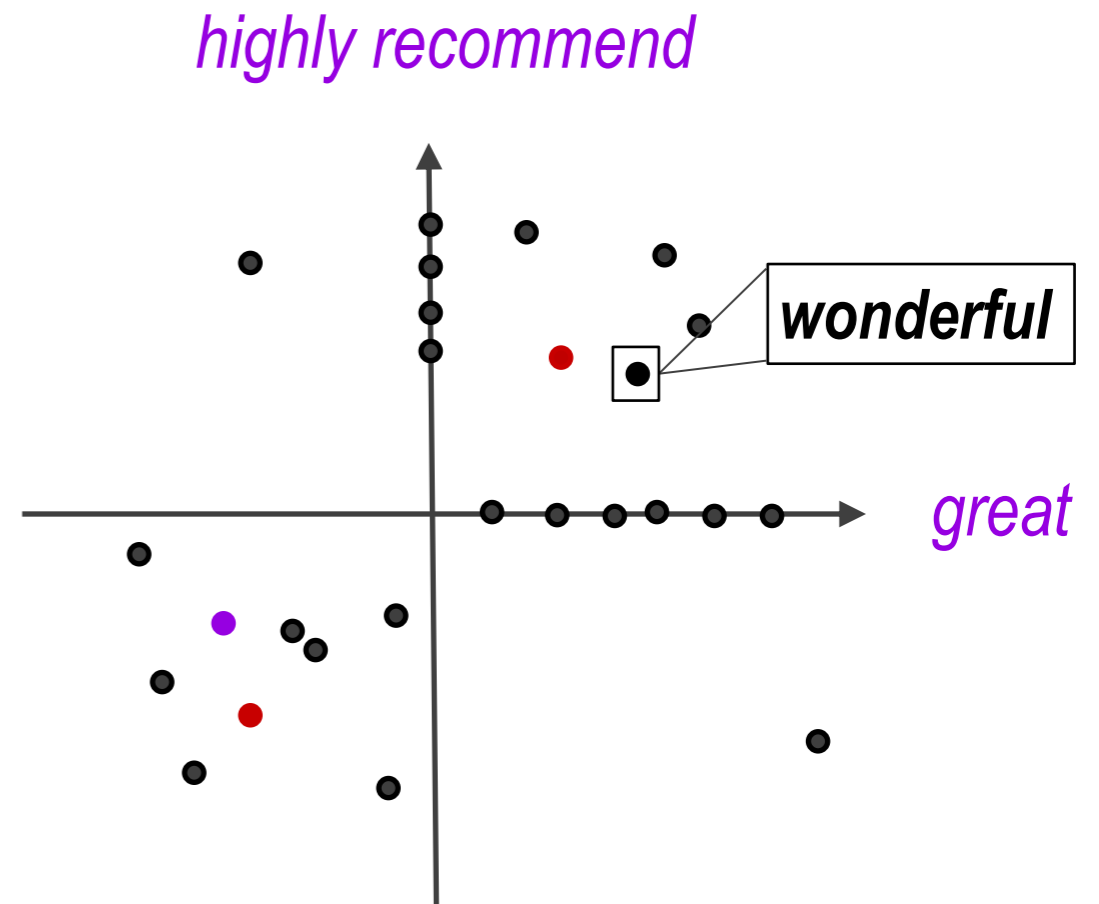


# Finding a shared sentiment subspace



$$W = \begin{bmatrix} | & & | & & | \\ w_1 & \dots & w(\text{highly recommend}) & \dots & w_N \\ | & & | & & | \end{bmatrix}$$

- $p_W(\text{pivots} | x)$  generates N new features
- $p_{w(\text{highly recommend})}(\text{highly recommend} | x)$  : “Did *highly recommend* appear?”
- Sometimes predictors capture non-sentiment information







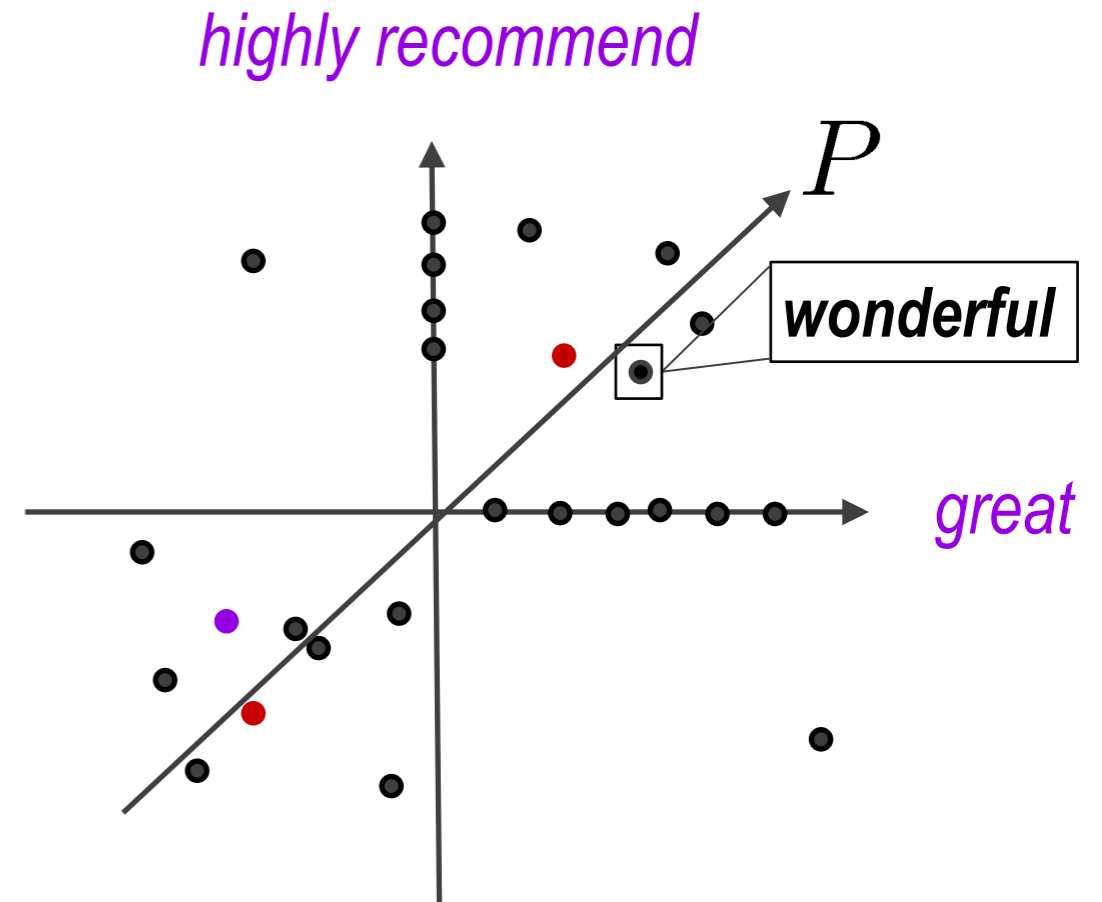
# Finding a shared sentiment subspace



$$W = \begin{bmatrix} | & & | & & | \\ w_1 & \dots & w(\text{highly recommend}) & \dots & w_N \\ | & & | & & | \end{bmatrix}$$

- Let  $P$  be a basis for the subspace of best fit to  $W$

- $p_W(\text{pivots}|x)$  generates  $N$  new features
- $p_{w(\text{highly recommend})}(\text{highly recommend}|x)$  : “Did *highly recommend* appear?”
- Sometimes predictors capture non-sentiment information





# Finding a shared sentiment subspace

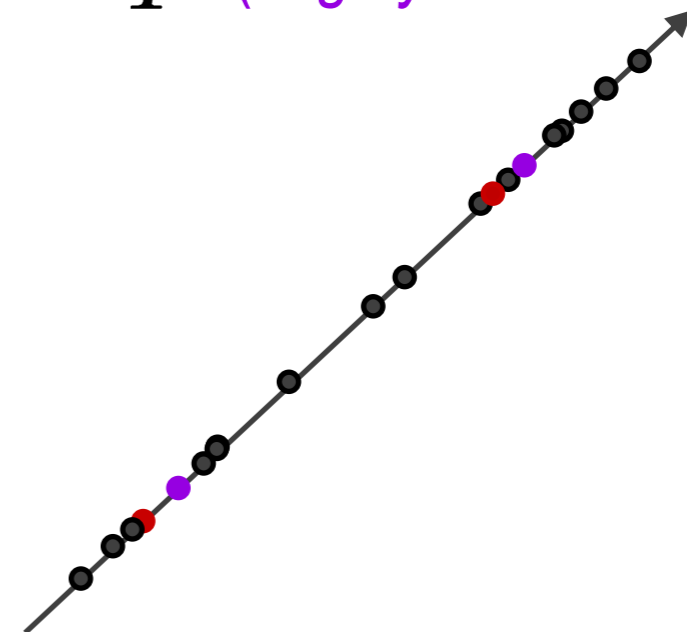


$$W = \begin{bmatrix} | & & | & & | \\ w_1 & \dots & w(\text{highly recommend}) & \dots & w_N \\ | & & | & & | \end{bmatrix}$$

- Let  $P$  be a basis for the subspace of best fit to  $W$
- $P$  captures sentiment variance in  $W$

- $p_W(\text{pivots}|x)$  generates  $N$  new features
- $p_{w(\text{highly recommend})}(\text{highly recommend}|x)$  : “Did *highly recommend* appear?”
- Sometimes predictors capture non-sentiment information

$P$  ( *highly recommend, great* )

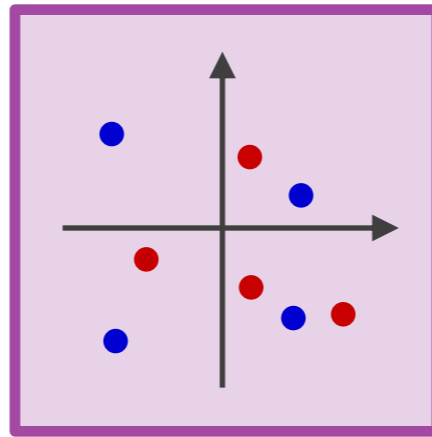




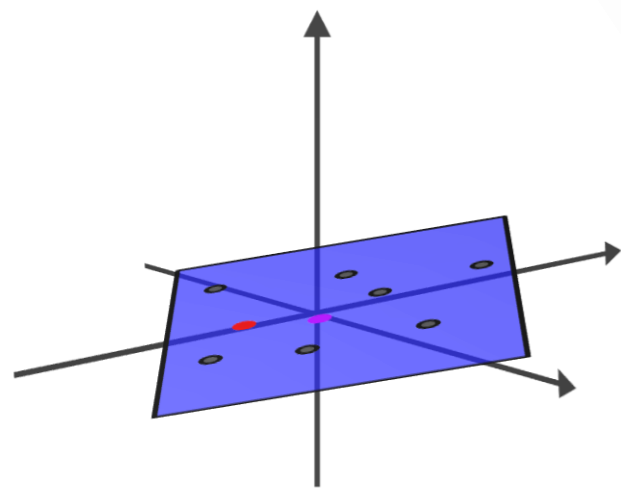
# P projects onto shared subspace



Source

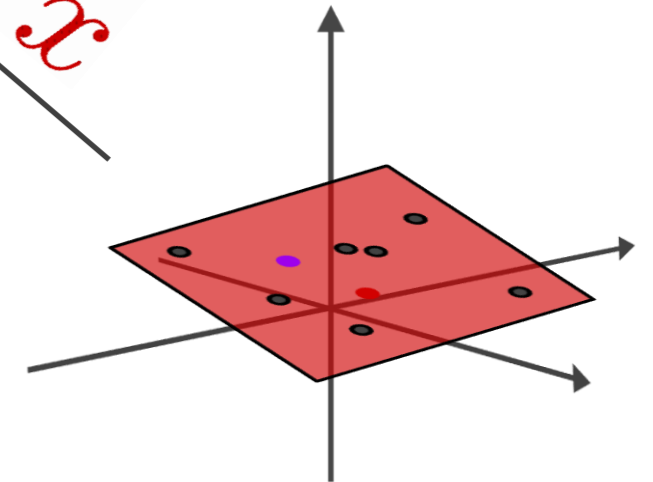


Target



$P x$

$P x$



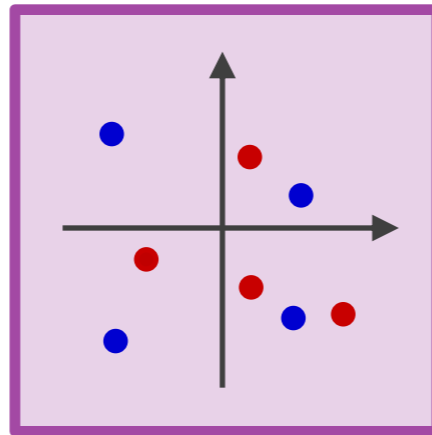
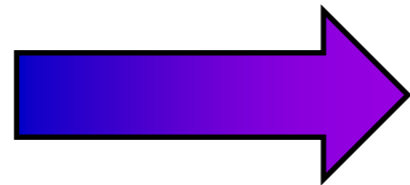
$$p_{\tilde{\theta}}(\text{thumbs up} | x) \propto \exp \left\{ \langle \phi(\text{thumbs up}, P x), \tilde{\theta} \rangle \right\}$$



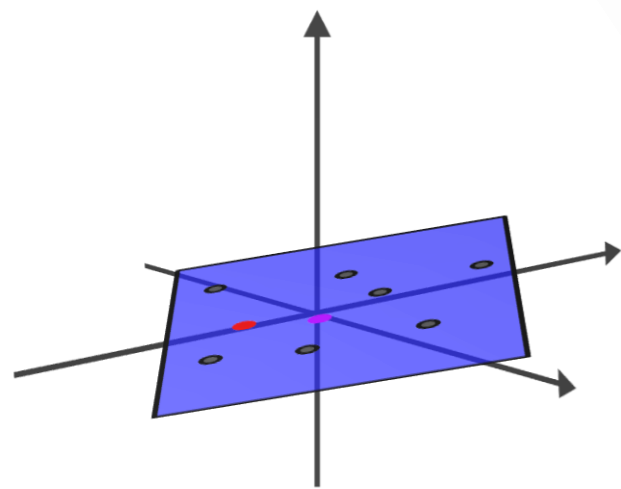
# P projects onto shared subspace



Source

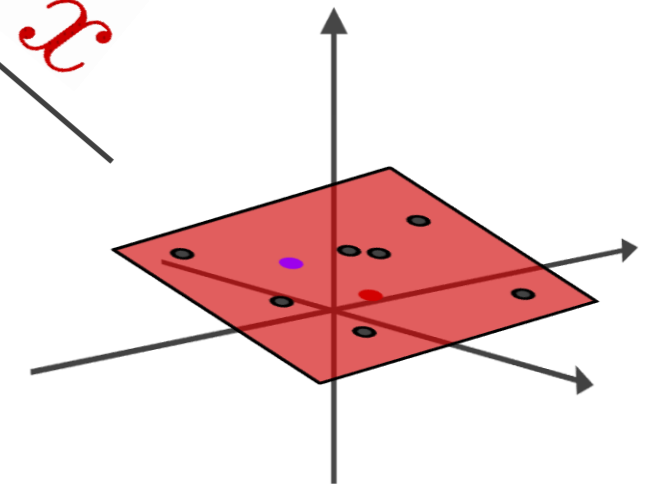


Target



$Px$

$Px$



$$h(x) = \text{sgn}(\theta^T Px)$$



# Correlating Pieces of the Bound



$$\begin{aligned} \epsilon_T(h) - \epsilon_T(h^*) &\leq \epsilon_{\hat{S}}(h, h^*) + O(\mathcal{R}_{\hat{S}}(\mathcal{H}) + \mathcal{R}_{\hat{T}}(\mathcal{H})) \\ &\quad + O\left(\sqrt{\frac{\log \frac{1}{\delta}}{n}}\right) + \text{disc}_{\mathcal{H}}(\hat{S}, \hat{T}) \end{aligned}$$

Component \ Projection	Discrepancy	Source Huber Loss	Target Error
Identity	1.796	0.003	0.253



# Correlating Pieces of the Bound



$$\epsilon_T(h) - \epsilon_T(h^*) \leq \epsilon_{\hat{S}}(h, h^*) + O(\mathcal{R}_{\hat{S}}(\mathcal{H}) + \mathcal{R}_{\hat{T}}(\mathcal{H})) \\ + O\left(\sqrt{\frac{\log \frac{1}{\delta}}{n}}\right) + \text{disc}_{\mathcal{H}}(\hat{S}, \hat{T})$$

<b>Component Projection</b>	<b>Discrepancy</b>	<b>Source Huber Loss</b>	<b>Target Error</b>
Identity	1.796	<b>0.003</b>	0.253
Random	0.223	0.254	0.561



# Correlating Pieces of the Bound

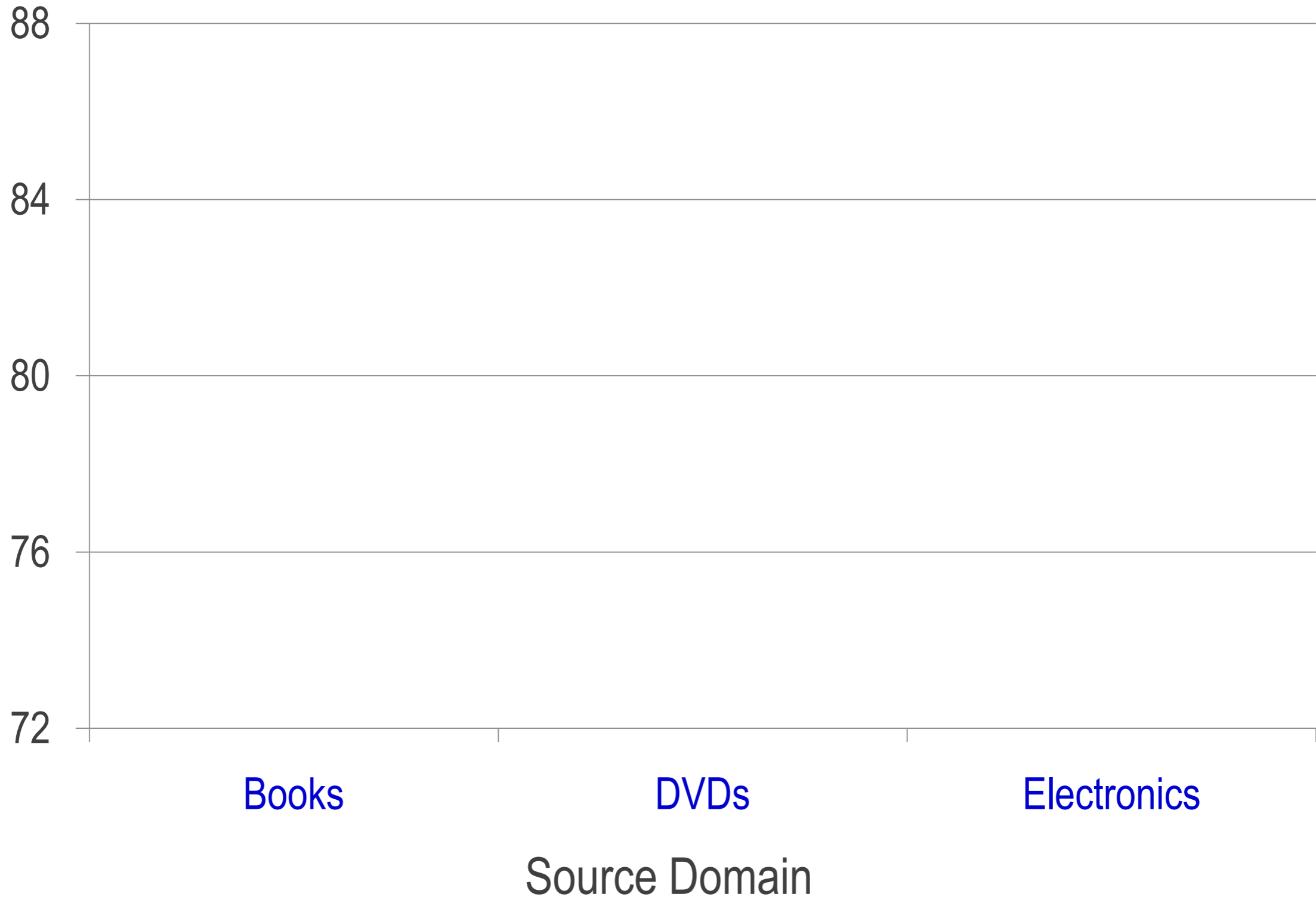


$$\epsilon_T(h) - \epsilon_T(h^*) \leq \epsilon_{\hat{S}}(h, h^*) + O(\mathcal{R}_{\hat{S}}(\mathcal{H}) + \mathcal{R}_{\hat{T}}(\mathcal{H})) \\ + O\left(\sqrt{\frac{\log \frac{1}{\delta}}{n}}\right) + \text{disc}_{\mathcal{H}}(\hat{S}, \hat{T})$$

<b>Component Projection</b>	<b>Discrepancy</b>	<b>Source Huber Loss</b>	<b>Target Error</b>
Identity	1.796	<b>0.003</b>	0.253
Random	0.223	0.254	0.561
Coupled Projection	<b>0.211</b>	0.07	<b>0.216</b>



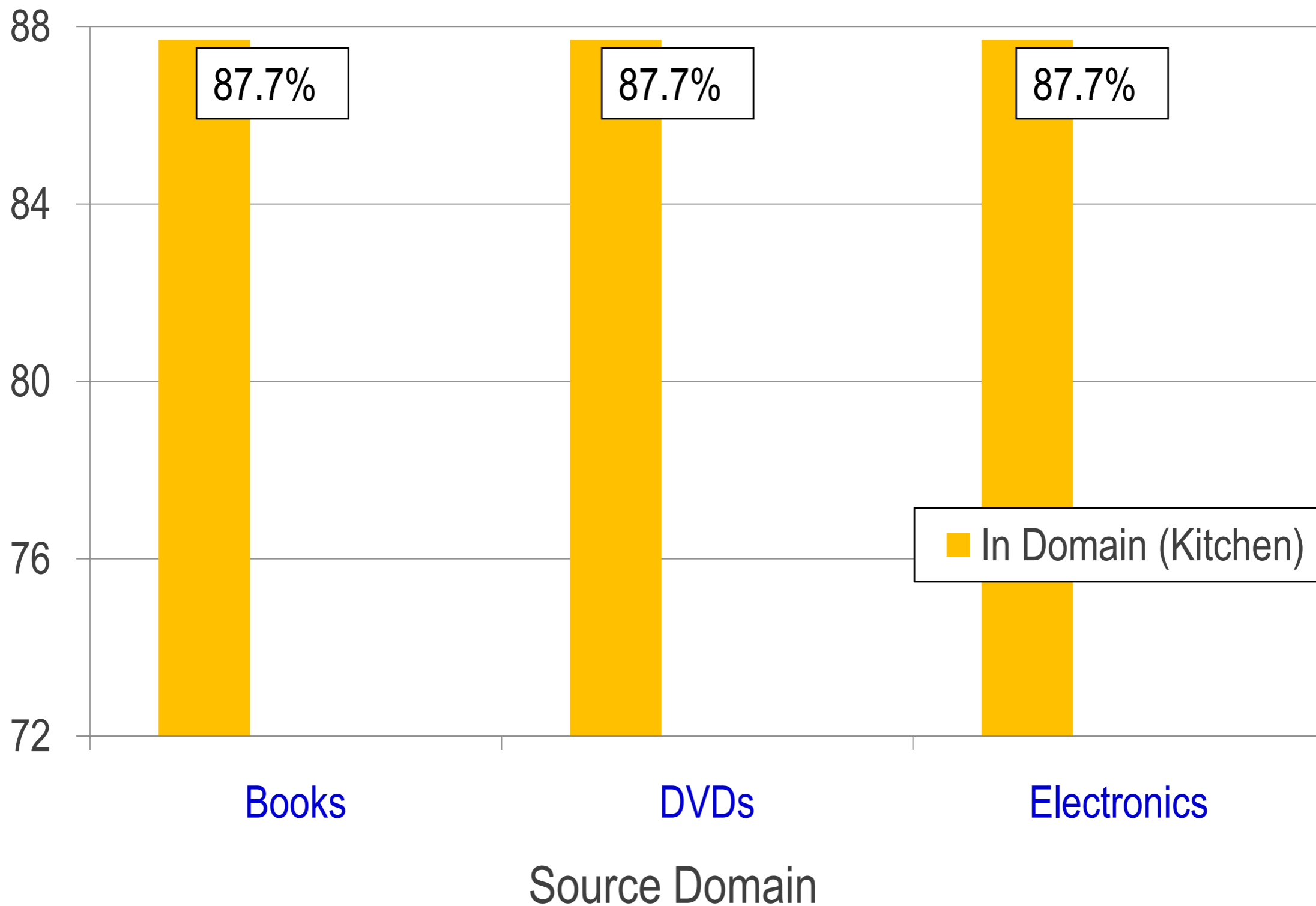
# Target Accuracy: Kitchen Appliances





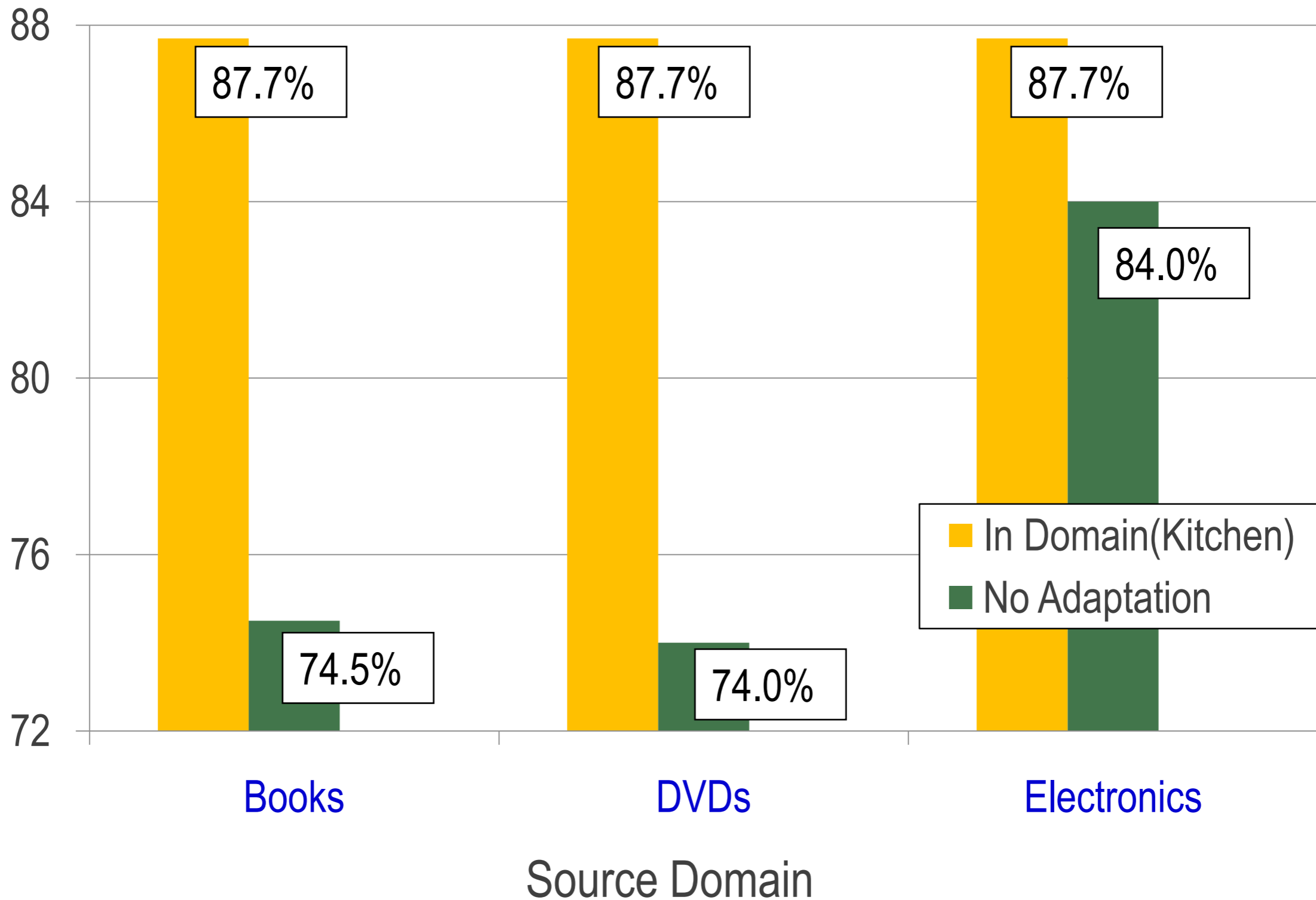


# Target Accuracy: Kitchen Appliances



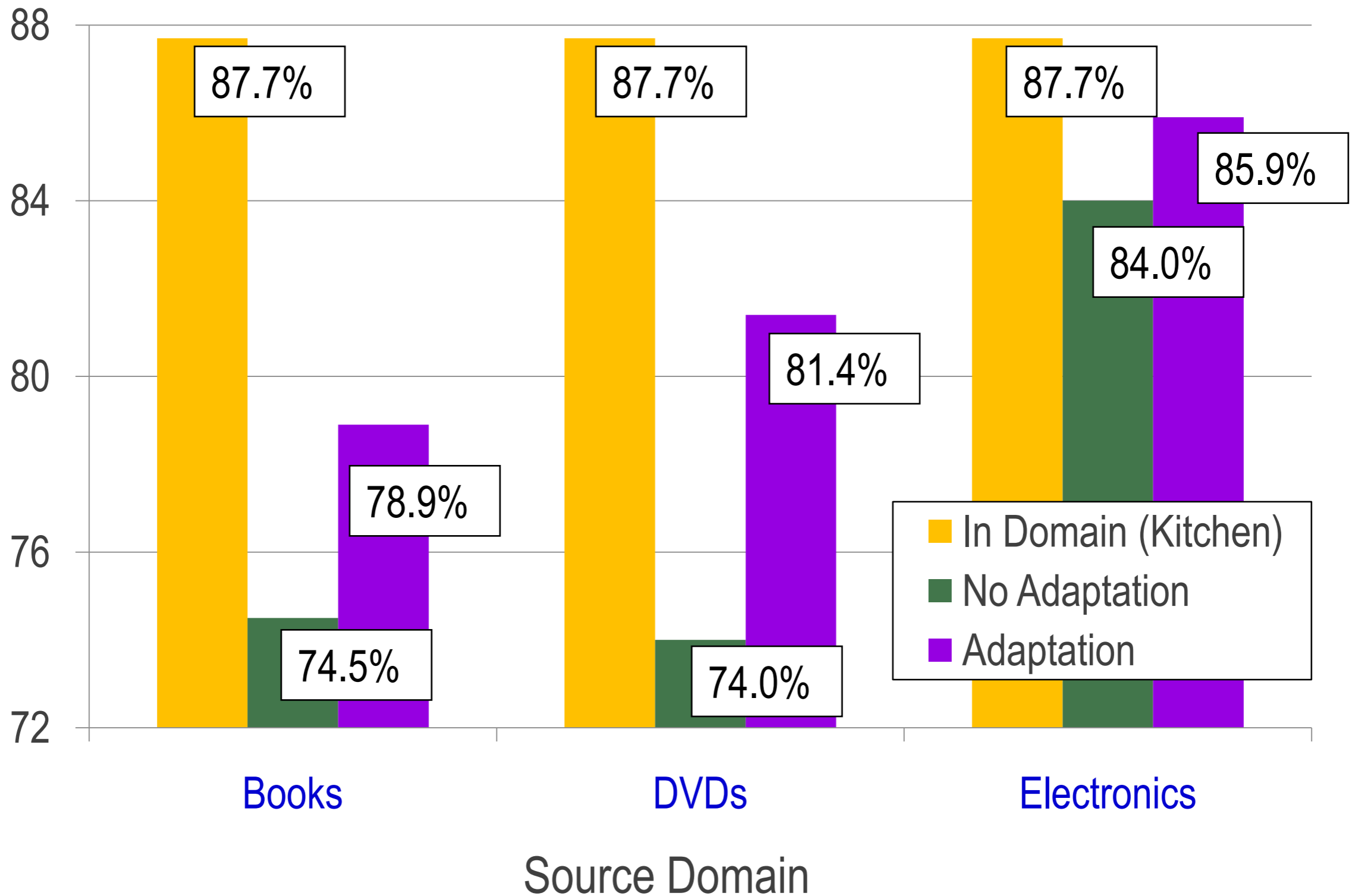


# Target Accuracy: Kitchen Appliances



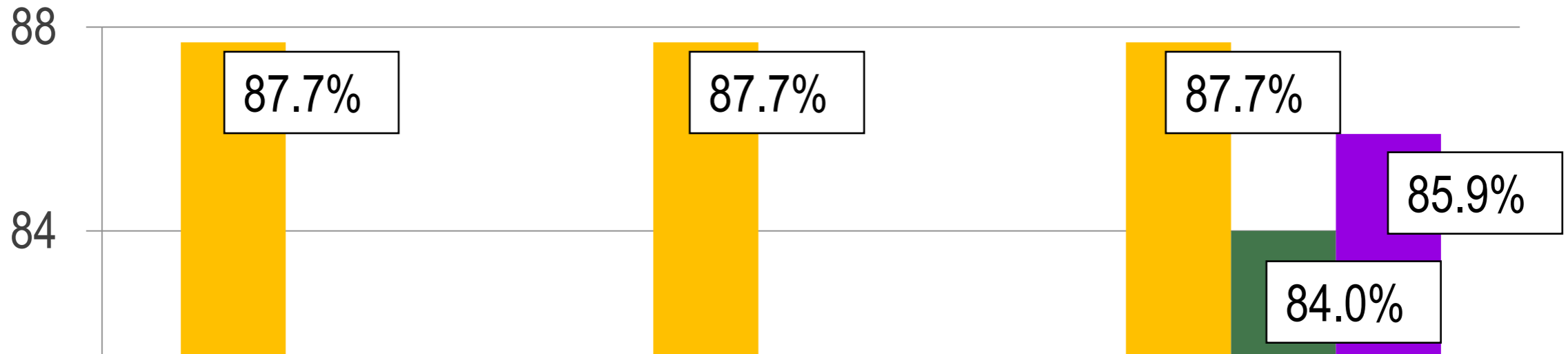


# Target Accuracy: Kitchen Appliances





# Adaptation Error Reduction



**36% reduction in error due to adaptation**



# Visualizing $P$ (books & kitchen)



negative

vs.

positive

*books*

plot

<#>\_pages

predictable

fascinating

engaging

must\_read

grisham

*poorly\_designed*

*awkward\_to*

*espresso*

*years\_now*

*the\_plastic*

*leaking*

*are\_perfect*

*a\_breeze*

*kitchen*



# Representation References

---



<http://adaptationtutorial.blitzer.com/references/>

- [1] Blitzer et al. Domain Adaptation with Structural Correspondence Learning. 2006.
- [2] S. Ben-David et al. Analysis of Representations for Domain Adaptation. 2007.
- [3] J. Blitzer et al. Domain Adaptation for Sentiment Classification. 2008.
- [4] Y. Mansour et al. Domain Adaptation: Learning Bounds and Algorithms. 2009.



# Tutorial Outline

---



1. Notation and Common Concepts
2. Semi-supervised Adaptation
  - Covariate shift
  - Learning Shared Representations
3. Supervised Adaptation
  - Feature-Based Approaches
  - Parameter-Based Approaches
4. Open Questions and Uncovered Algorithms



# Feature-based approaches

---



## Cell-phone domain:

“horrible” is **bad**  
“small” is **good**

## Hotel domain:

“horrible” is **bad**  
“small” is **bad**

### **Key Idea:**

Share some features (“horrible”)  
Don't share others (“small”)

(and let an arbitrary *learning algorithm*  
decide which are which)





F

—

In feature-vector lingo:

$x \rightarrow \langle x, x, 0 \rangle$  (for source domain)

$x \rightarrow \langle x, 0, x \rangle$  (for target domain)

The phone is small

The hotel is small

Original  
Features

W:the  
W:phone  
W:is  
W:small

W:the  
W:hotel  
W:is  
W:small

Augmented  
Features

S:the  
S:phone  
S:is  
S:small

T:the  
T:hotel  
T:is  
T:small



# A Kernel Perspective



In feature-vector lingo:

$x \rightarrow \langle x, x, 0 \rangle$  (for source domain)

$x \rightarrow \langle x, 0, x \rangle$  (for target domain)

$$K^{\text{aug}}(x,z) = \begin{cases} 2K(x,z) & \text{if } x,z \text{ from same domain} \\ K(x,z) & \text{otherwise} \end{cases}$$

**We have *ensured*  
SGH & *destroyed*  
shared support**



# Named Entity Rec.: /bush/



General BC-news Conversations Newswire Weblogs Usenet Telephone

Person

Person	<input type="checkbox"/>	<input type="checkbox"/>				<input type="checkbox"/>	<input checked="" type="checkbox"/>
Geo-political entity	<input type="checkbox"/>	<input checked="" type="checkbox"/>				<input type="checkbox"/>	<input checked="" type="checkbox"/>
Organization	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>				<input type="checkbox"/>	<input checked="" type="checkbox"/>
Location	<input type="checkbox"/>	<input type="checkbox"/>				<input type="checkbox"/>	<input checked="" type="checkbox"/>

Geo-political entity

Organization

Location

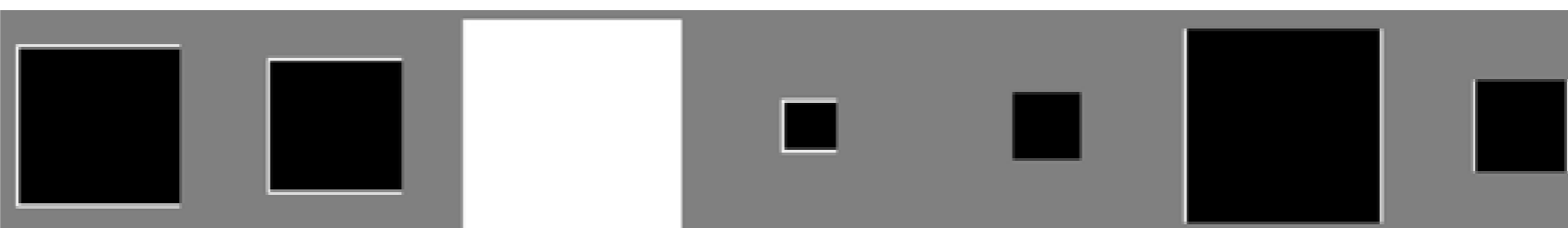


# Named Entity Rec.: p=/the/

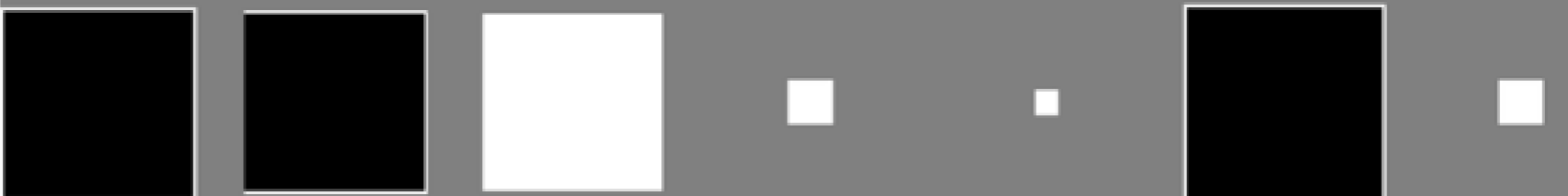


General BC-news Conversations Newswire Weblogs Usenet Telephone

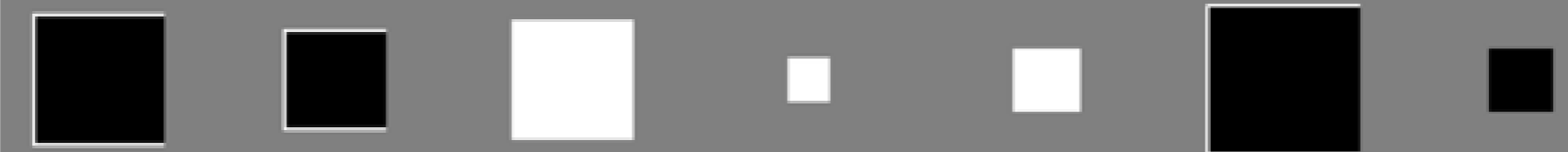
Person



Geo-political entity



Organization



Location





# Some experimental results



Task	Dom	SrcOnly	TgtOnly	Baseline	Prior	Augment
ACE- NER	bn	4.98	2.37	2.11 (pred)	2.06	<b>1.98</b>
	bc	4.54	4.07	3.53 (weight)	<b>3.47</b>	<b>3.47</b>
	nw	4.78	3.71	3.56 (pred)	3.68	<b>3.39</b>
	wl	2.45	2.45	<b>2.12</b> (all)	2.41	<b>2.12</b>
	un	3.67	2.46	2.10 (linint)	2.03	<b>1.91</b>
	cts	2.08	0.46	0.40 (all)	<b>0.34</b>	<b>0.32</b>
CoNLL	tgt	2.49	2.95	<b>1.75</b> (wgt/li)	1.89	<b>1.76</b>
PubMed	tgt	12.02	4.15	3.95 (linint)	3.99	<b>3.61</b>
CNN	tgt	10.29	3.82	3.44 (linint)	<b>3.35</b>	<b>3.37</b>
	wsj	6.63	4.35	4.30 (weight)	4.27	<b>4.11</b>
	swbd3	15.90	4.15	4.09 (linint)	3.60	<b>3.51</b>
Tree bank- Chunk	br-cf	5.16	6.27	<b>4.72</b> (linint)	5.22	5.15
	br-cg	4.32	5.36	<b>4.15</b> (all)	4.25	4.90
	br-ck	5.05	6.32	<b>5.01</b> (prd/li)	5.27	5.41
	br-cl	5.66	6.60	<b>5.39</b> (wgt/prd)	5.99	5.73
	br-cm	3.57	6.59	<b>3.11</b> (all)	4.08	4.89
	br-cn	4.60	5.56	<b>4.19</b> (prd/li)	4.48	4.42
	br-cp	4.82	5.62	<b>4.55</b> (wgt/prd/li)	4.87	4.78
	br-cr	5.78	9.13	<b>5.15</b> (linint)	6.71	6.30
Treebank- brown		6.35	5.75	4.72 (linint)	4.72	<b>4.65</b>

- Can bound expected target error:

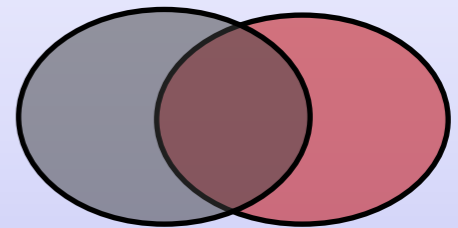
Average training error

$$\epsilon_t \leq \frac{1}{2} (\hat{\epsilon}_s + \hat{\epsilon}_t) + O(\text{complexity})$$

$$+ \left( \frac{1}{\sqrt{N_s}} + \frac{1}{\sqrt{N_t}} \right) O\left(\frac{1}{\delta}\right) + O(\text{disc}_H(S, T))$$

Number of  
source examples

Number of  
target examples

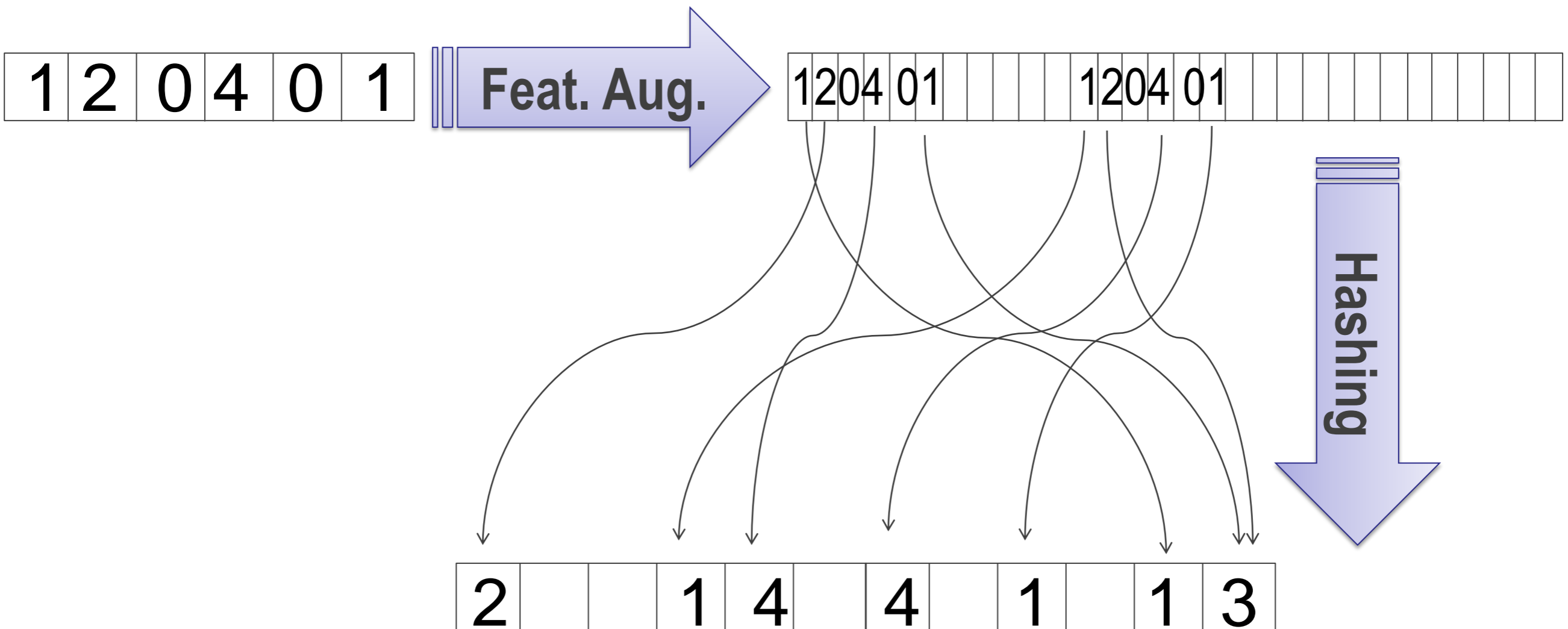




# Feature Hashing

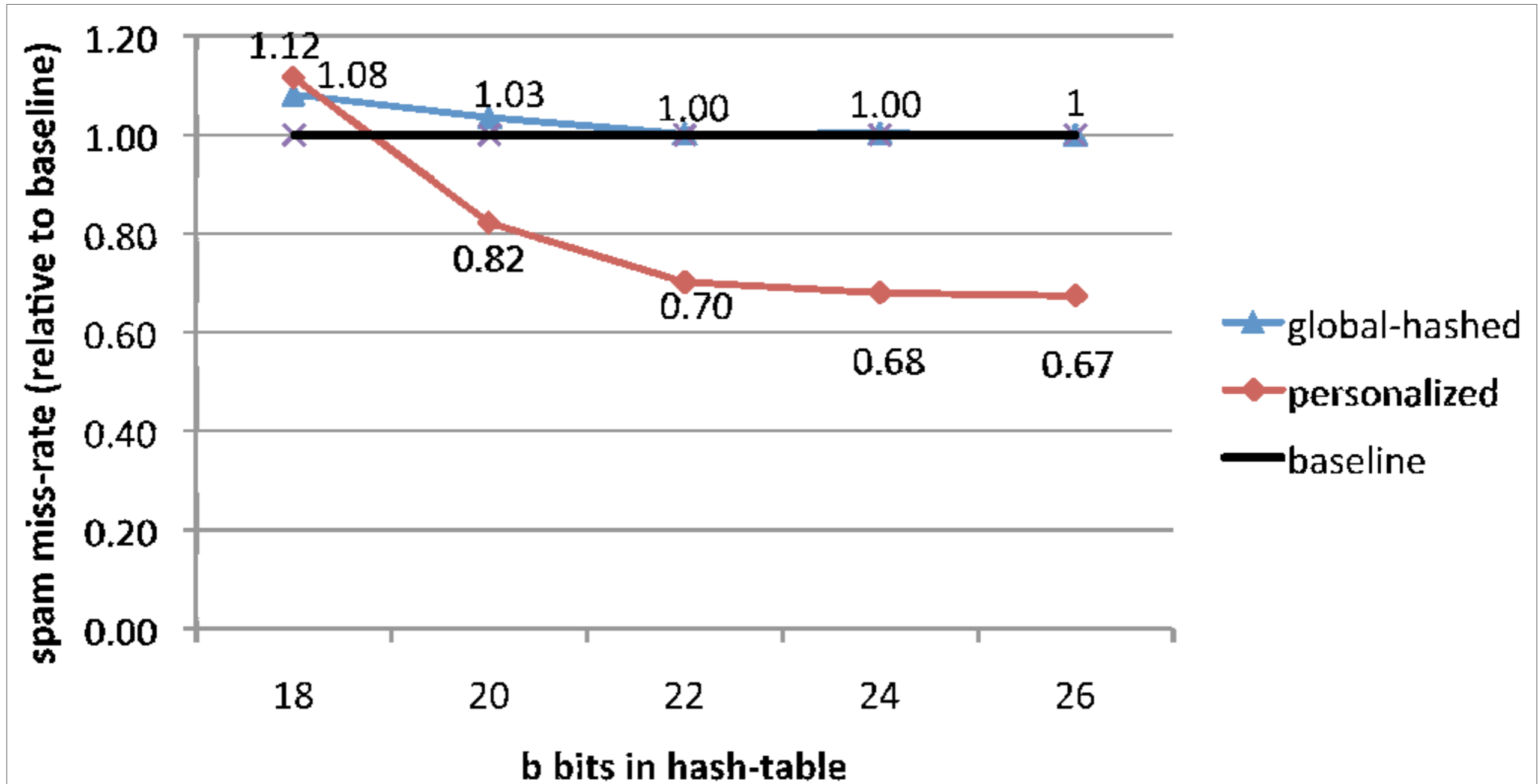


- Feature augmentation creates  $(K+1)D$  parameters
- Too big if  $K \gg 20$ , but *very sparse!*





# Hash Kernels







# Semi-sup Feature Augmentation



- For labeled data:
  - $(y, x) \rightarrow (y, \langle x, x, 0 \rangle)$  (for source domain)
  - $(y, x) \rightarrow (y, \langle x, 0, x \rangle)$  (for target domain)
- What about unlabeled data?
- Encourage agreement:

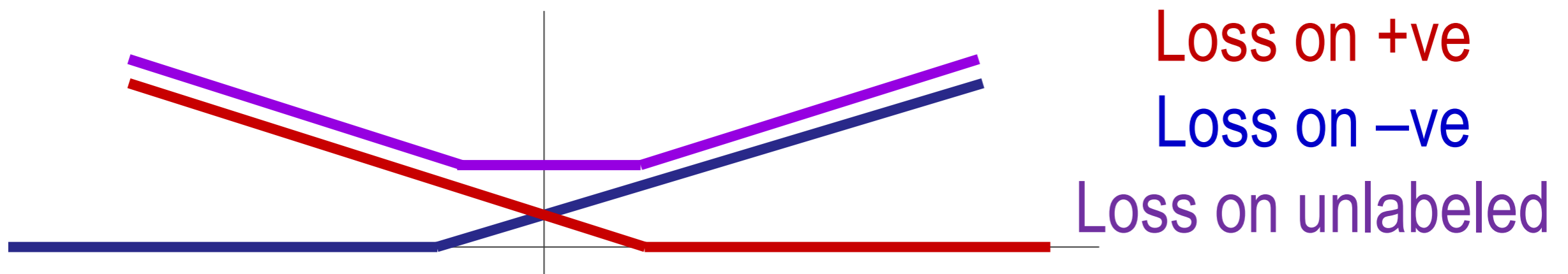
$$[h_t(x) = h_s(x)] \Leftrightarrow [w_t \circ x - w_s \circ x = 0]$$



# Semi-sup Feature Augmentation



- For labeled data:
  - $(y, x) \rightarrow (y, \langle x, x, 0 \rangle)$  (for source domain)
  - $(y, x) \rightarrow (y, \langle x, 0, x \rangle)$  (for target domain)
- What about unlabeled data?
  - $(x) \rightarrow \{ (+1, \langle 0, x, -x \rangle), (-1, \langle 0, x, -x \rangle) \}$



- Encourages *agreement* on unlabeled data
  - Akin to *multiview learning*
  - Reduces generalization bound



# Feature-based References

---



- T. Evgeniou and M. Pontil. Regularized Multi-task Learning (2004).
- H. Daumé III, Frustratingly Easy Domain Adaptation. 2007.
- K. Weinberger, A. Dasgupta, J. Langford, A. Smola, J. Attenberg. Feature Hashing for Large Scale Multitask Learning. 2009.
- A. Kumar, A. Saha and H. Daumé III, Frustratingly Easy Semi-Supervised Domain Adaptation. 2010.



# Tutorial Outline

---



1. Notation and Common Concepts
2. Semi-supervised Adaptation
  - Covariate shift
  - Learning Shared Representations
3. Supervised Adaptation
  - Feature-Based Approaches
  - Parameter-Based Approaches
4. Open Questions and Uncovered Algorithms



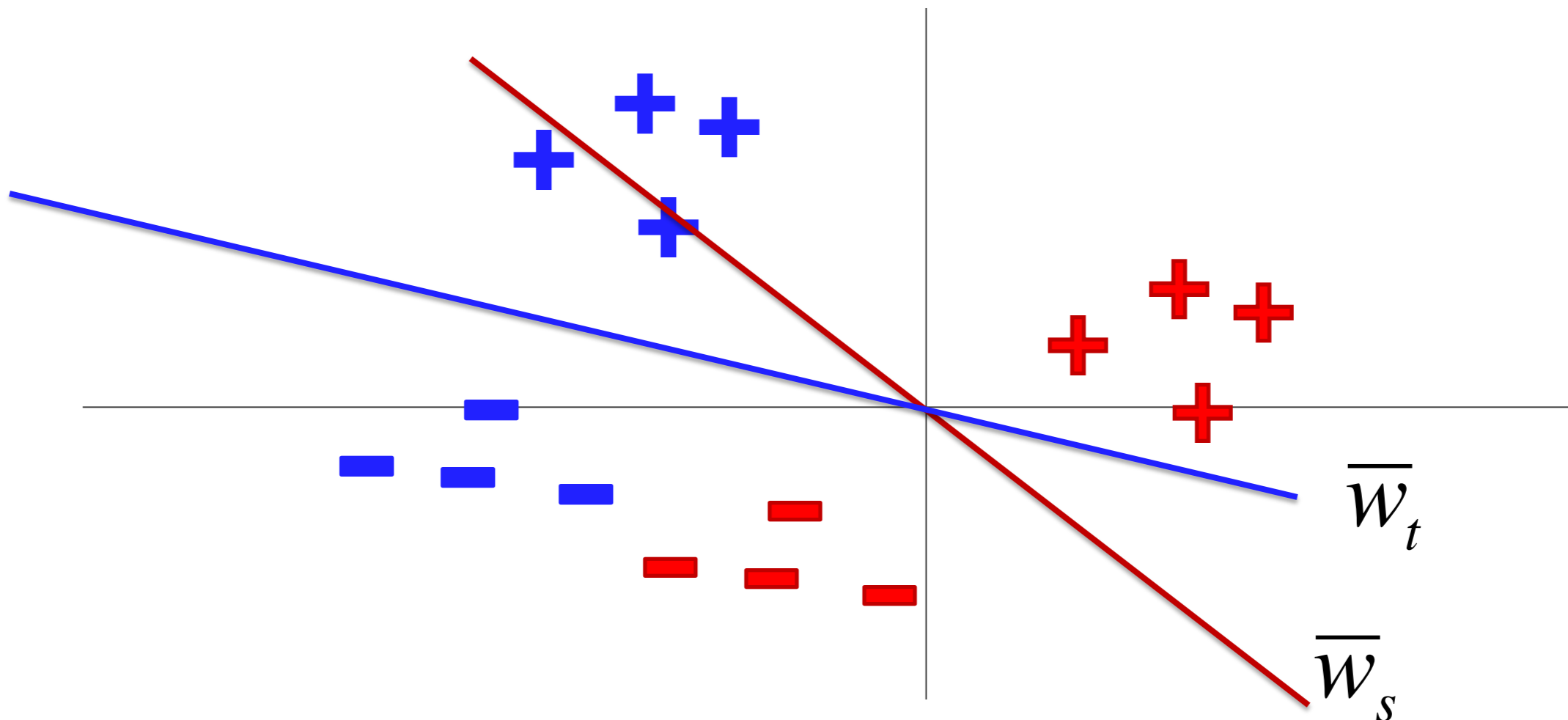
# A Parameter-based Perspective



- Instead of duplicating features, write:

$$\bar{w}_t = \bar{w}_s + \bar{v}$$

- And *regularize*  $\bar{w}_s$  and  $\bar{v}$  toward zero

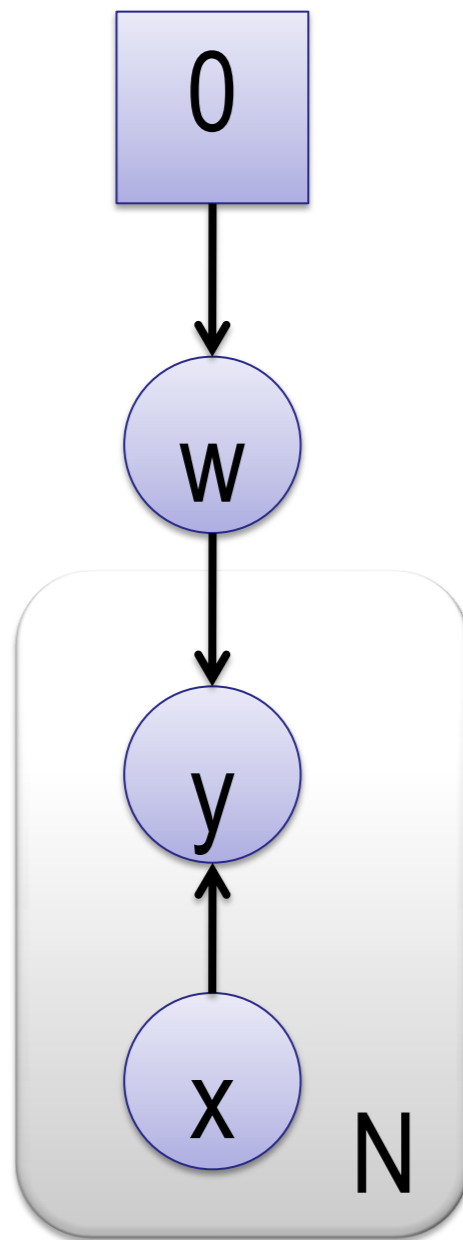


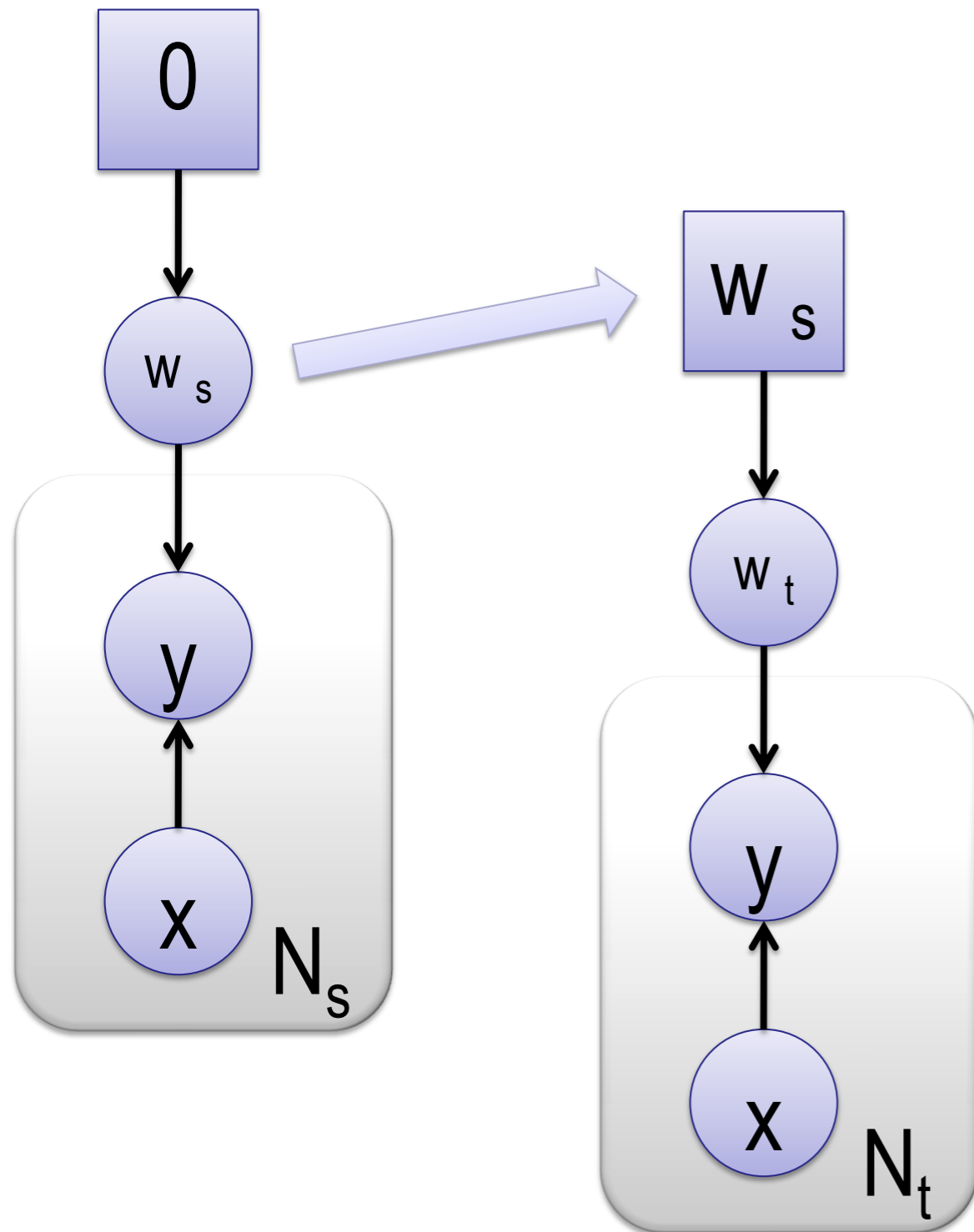


# Sharing Parameters via Bayes

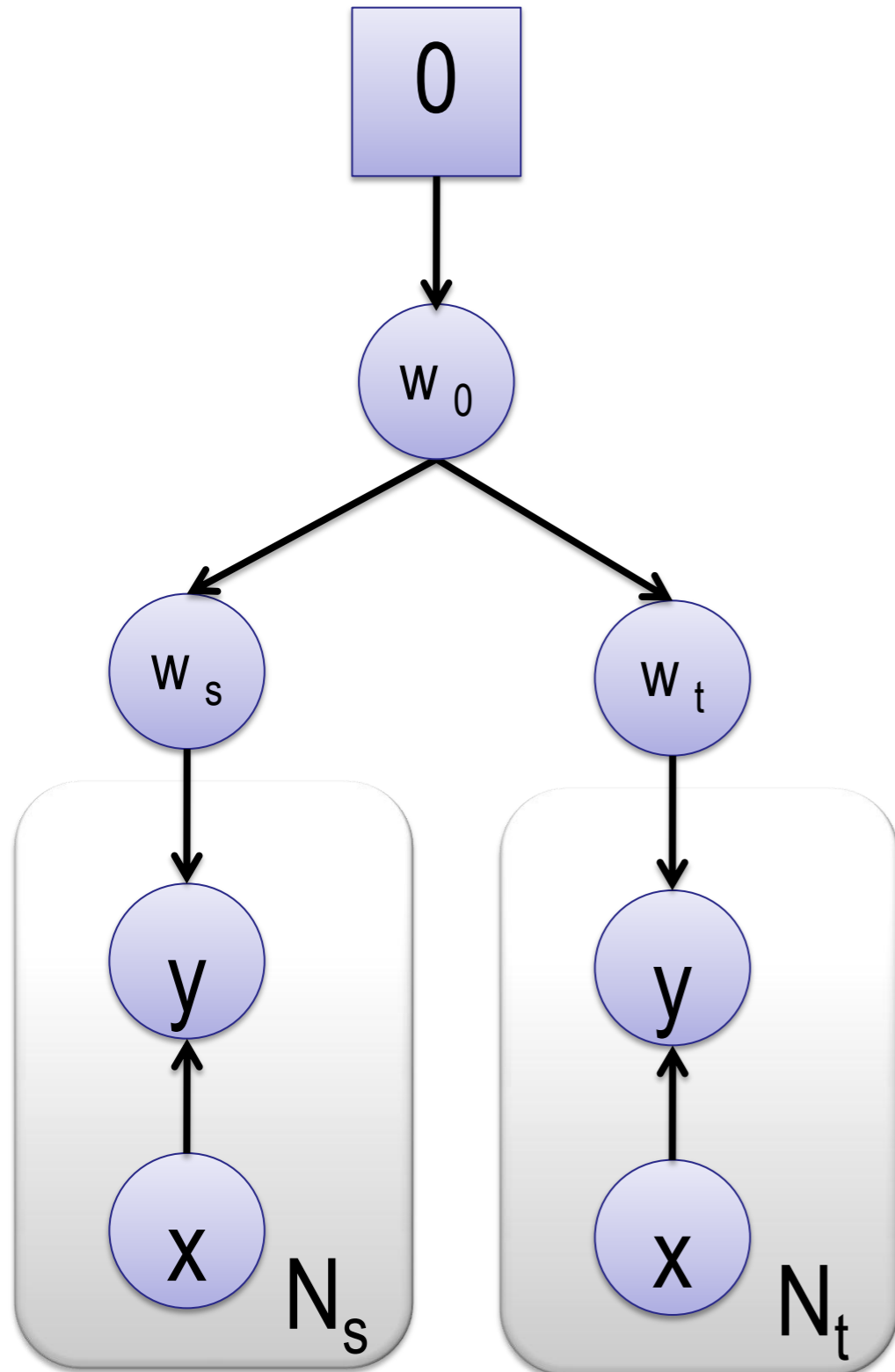


- N data points
- $w$  is regularized to zero
- Given  $x$  and  $w$ , we predict  $y$





- Train model on source domain, regularized toward zero
- Train model on target domain, regularized toward source domain

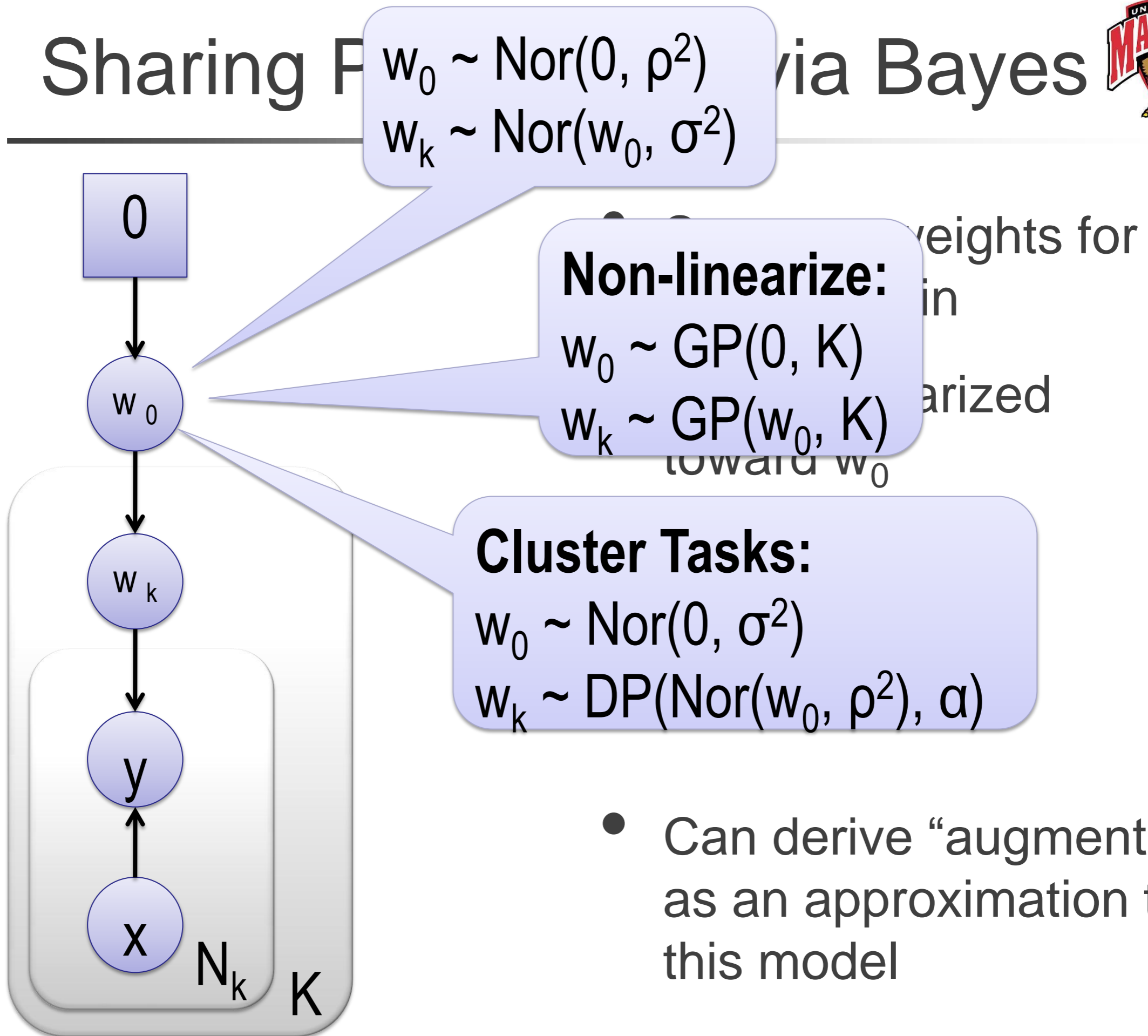


- Separate weights for each domain
- Each regularized toward  $w_0$
- $w_0$  regularized toward zero





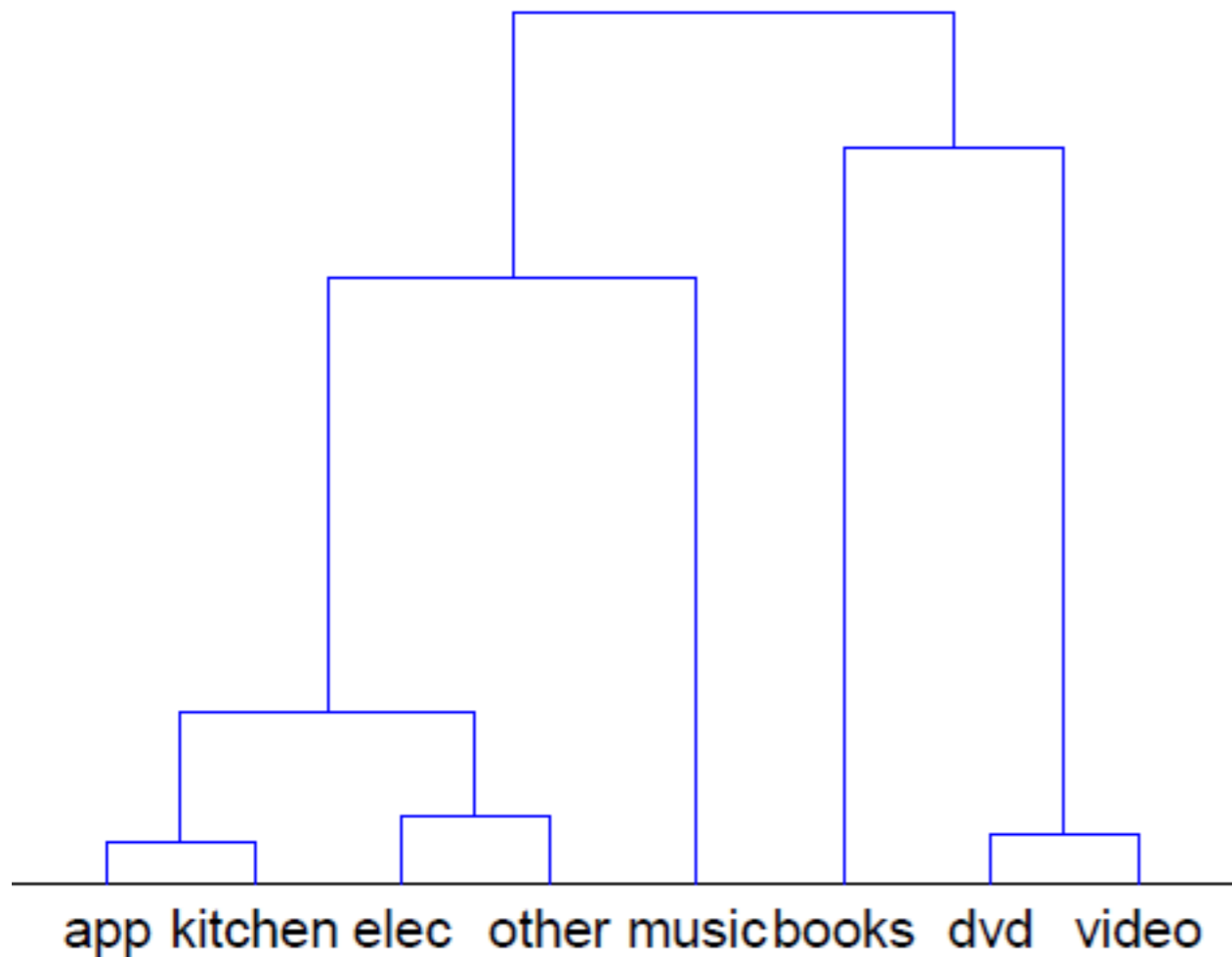
# Sharing Priors via Bayes



- Can derive “augment” as an approximation to this model



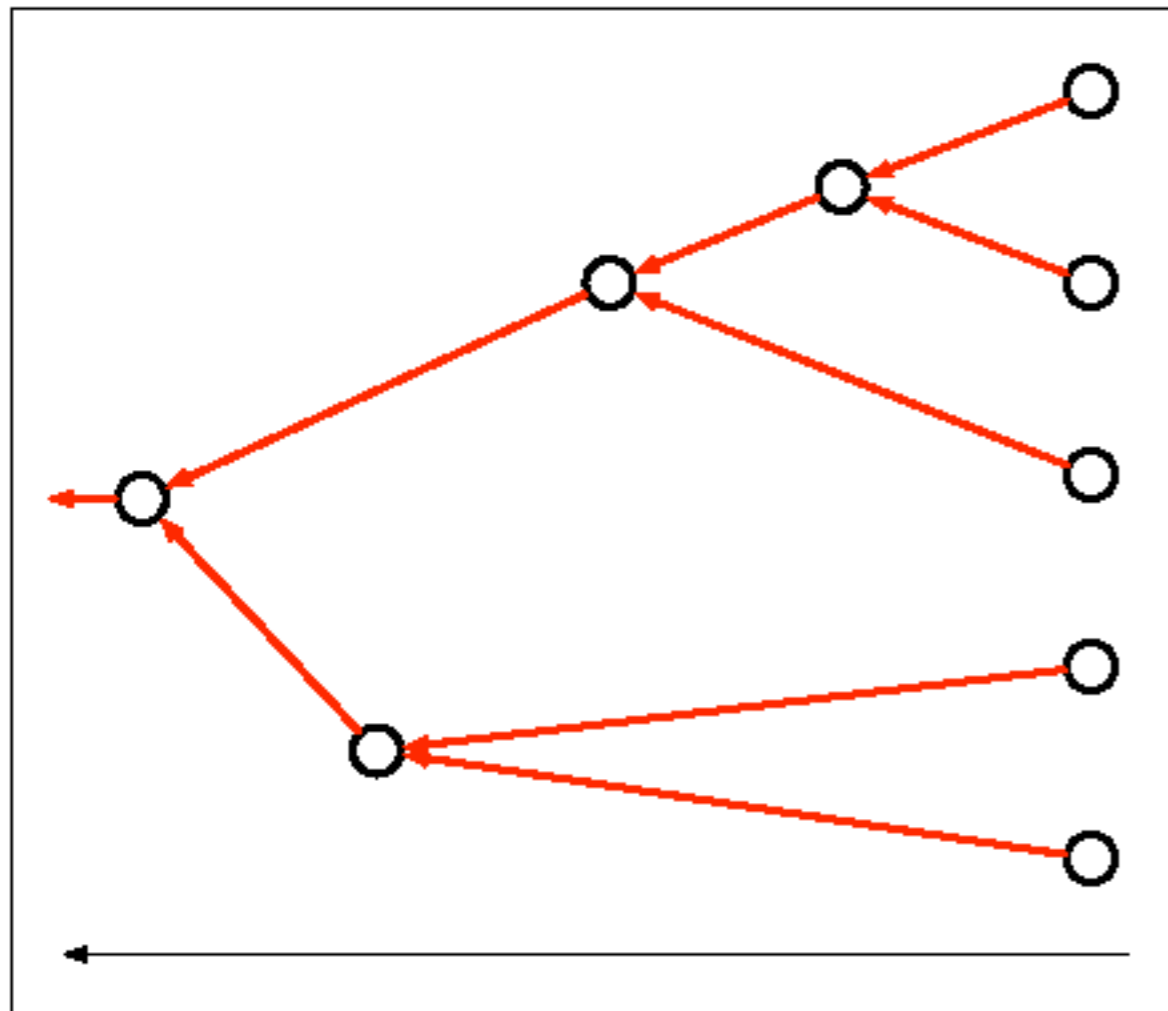
# Not all domains created equal



- Would like to infer tree structure automatically
- Tree structure should be good for the *task*
- Want to simultaneously infer tree structure and parameters



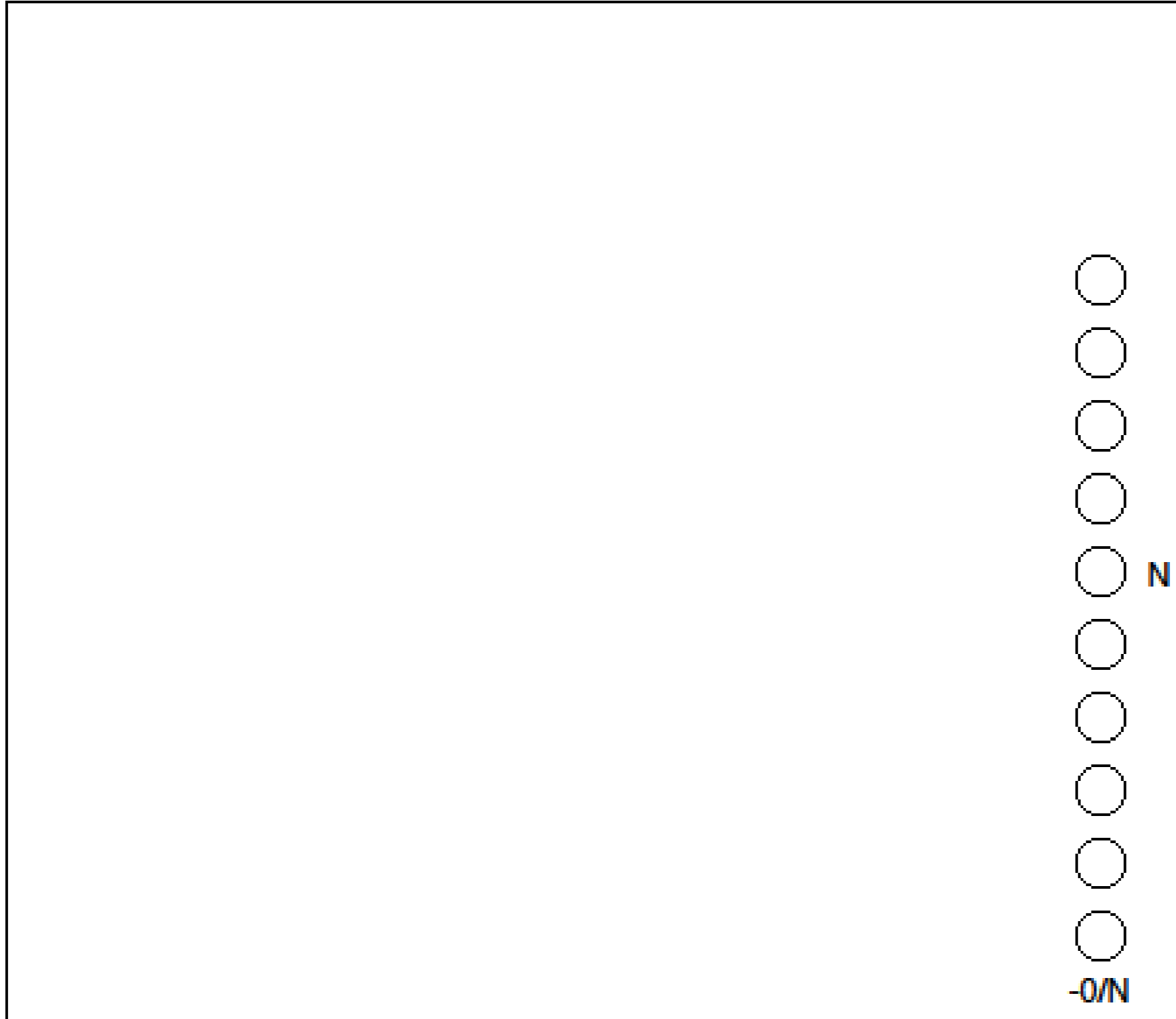
# Kingman's Coalescent



- A standard model for the genealogy of a population
- Each organism has exactly one parent (haploid)
- Thus, the genealogy is a tree

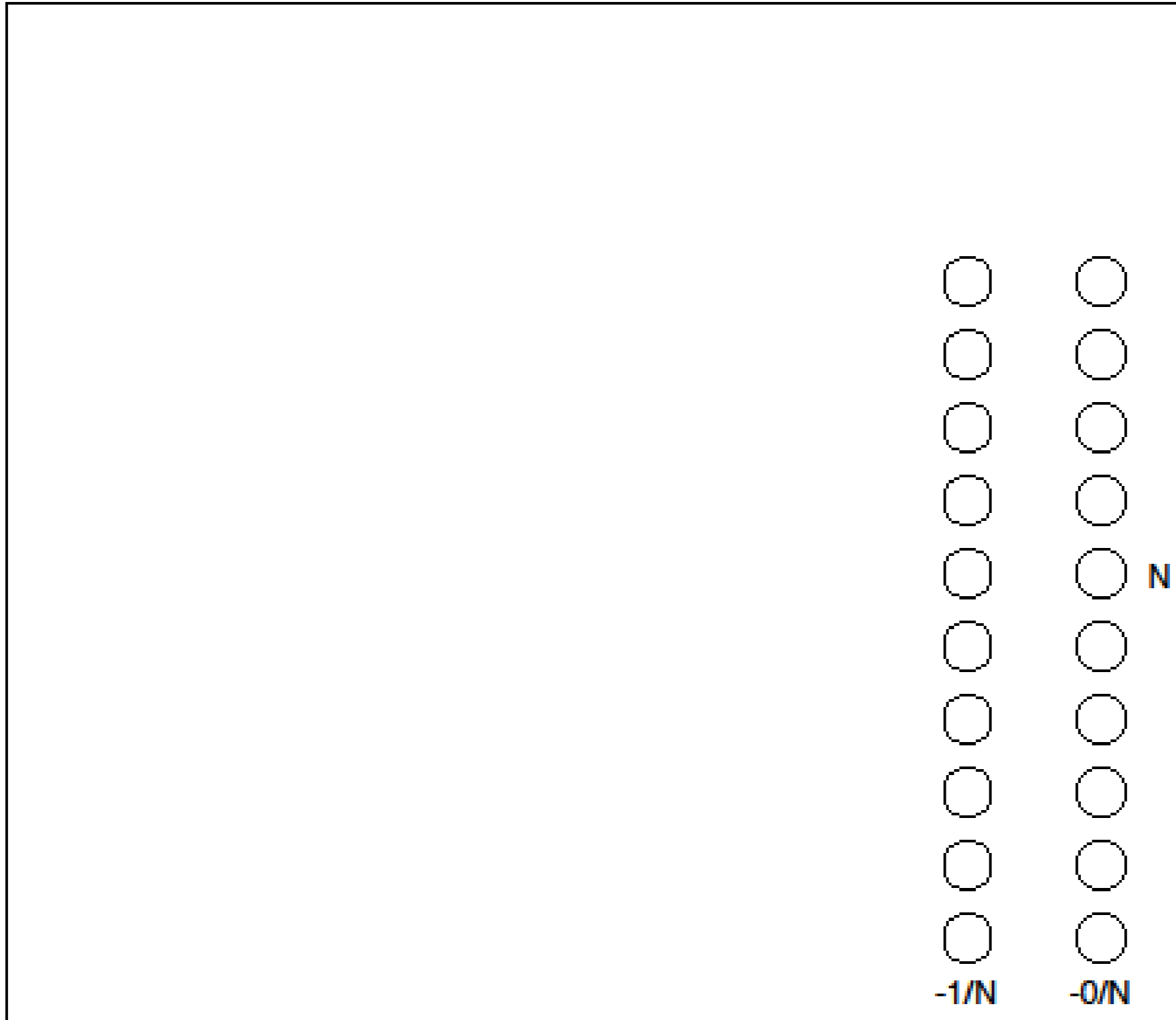


# A distribution over trees



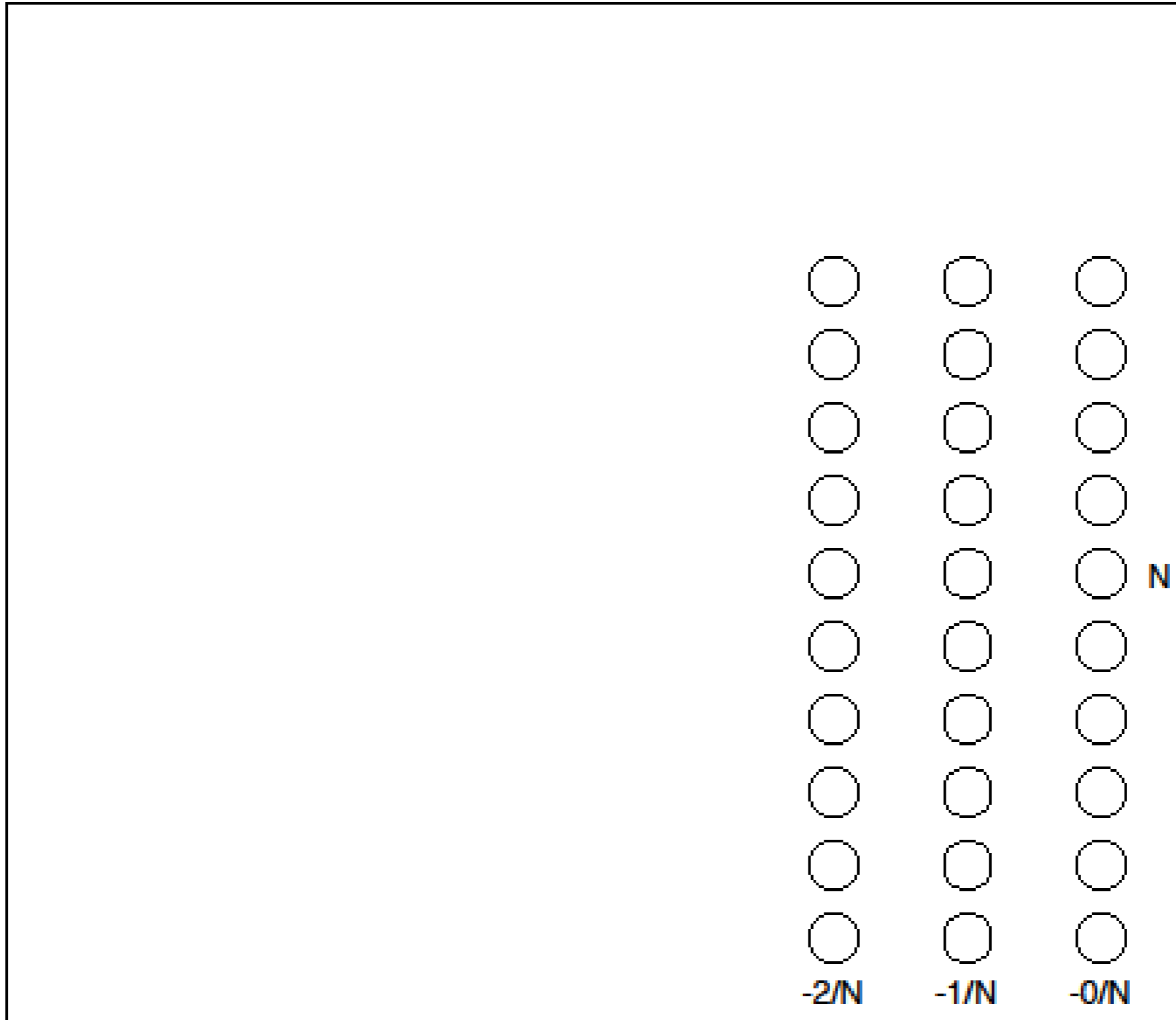


# A distribution over trees



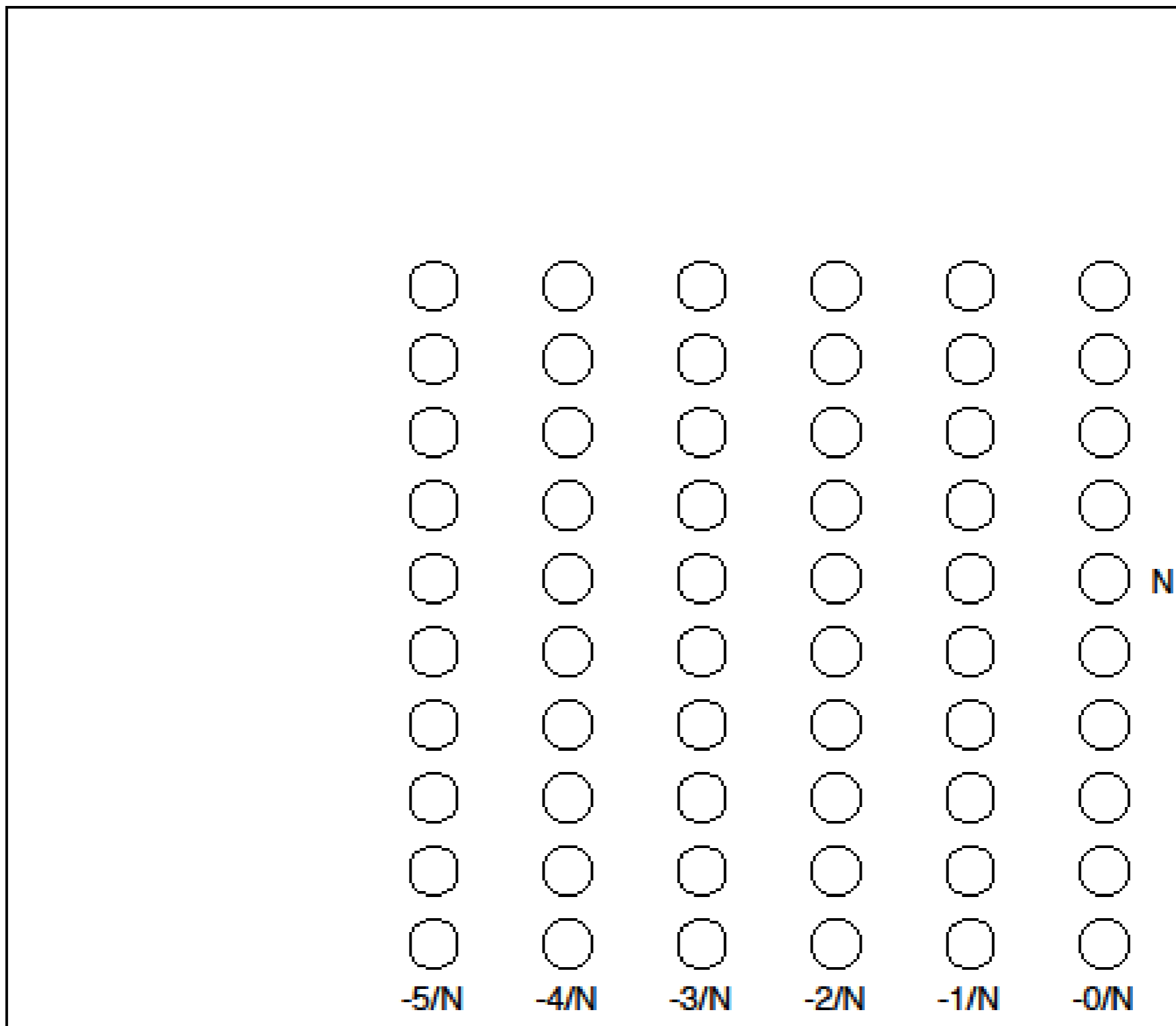


# A distribution over trees



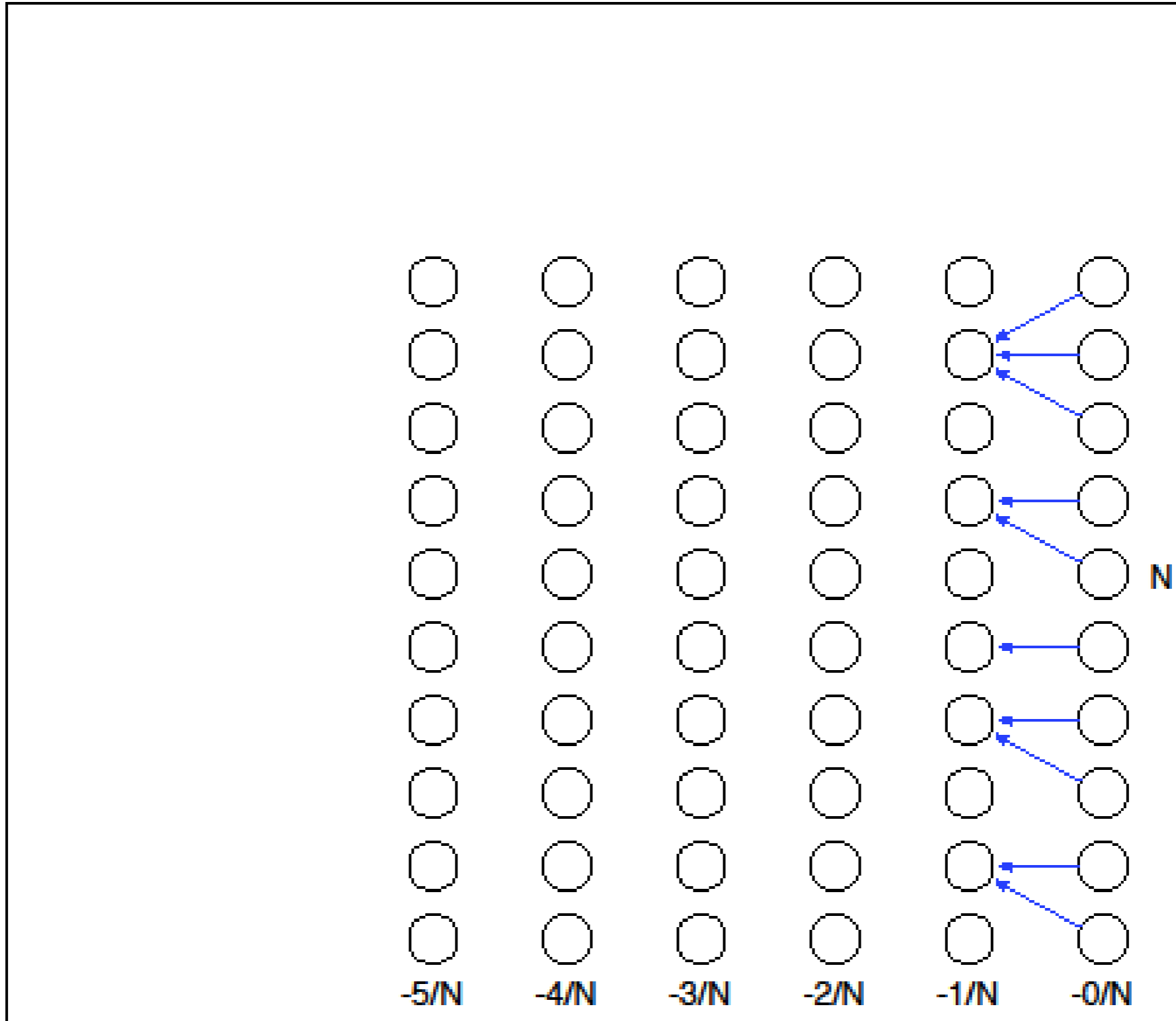


# A distribution over trees





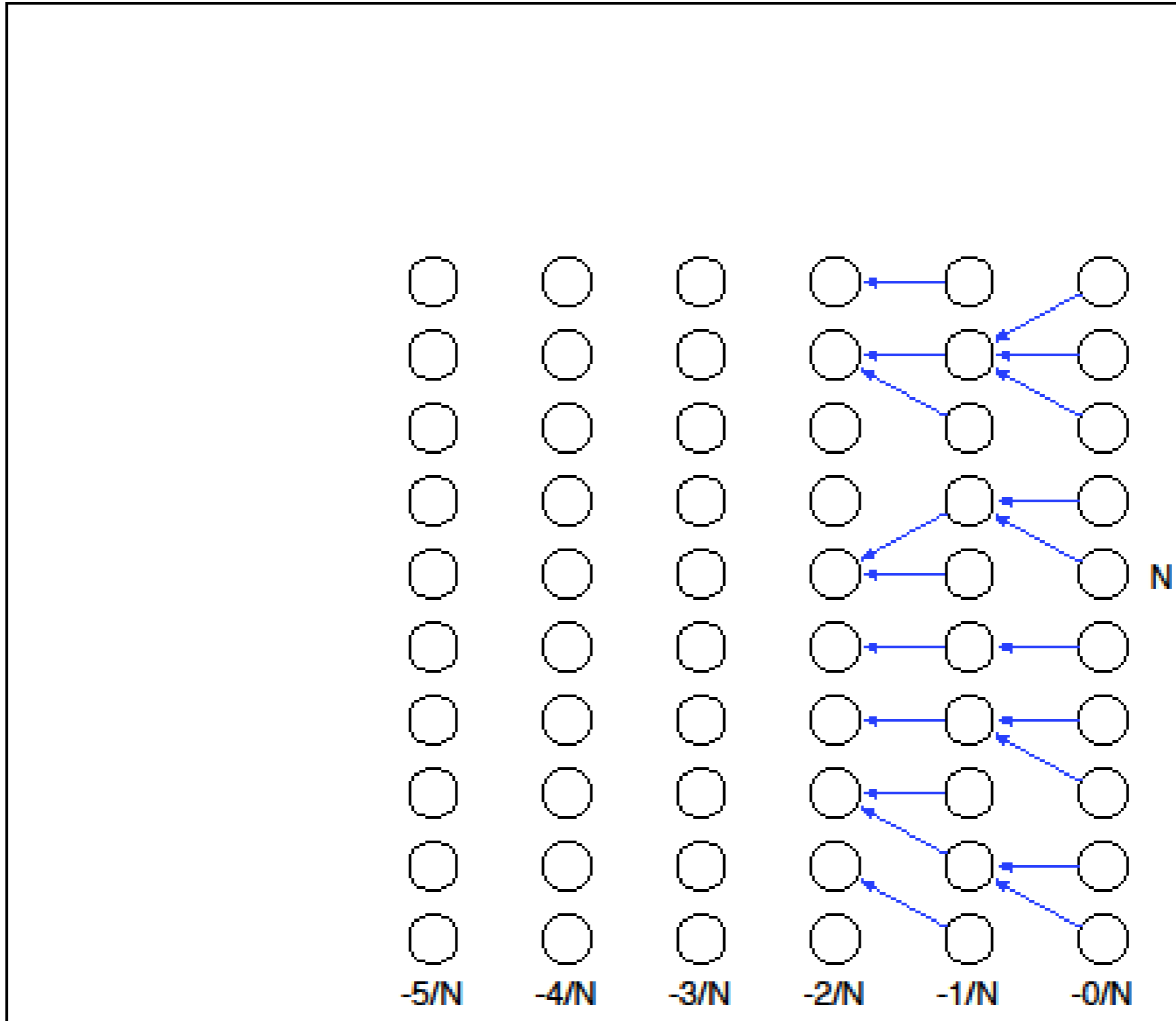
# A distribution over trees





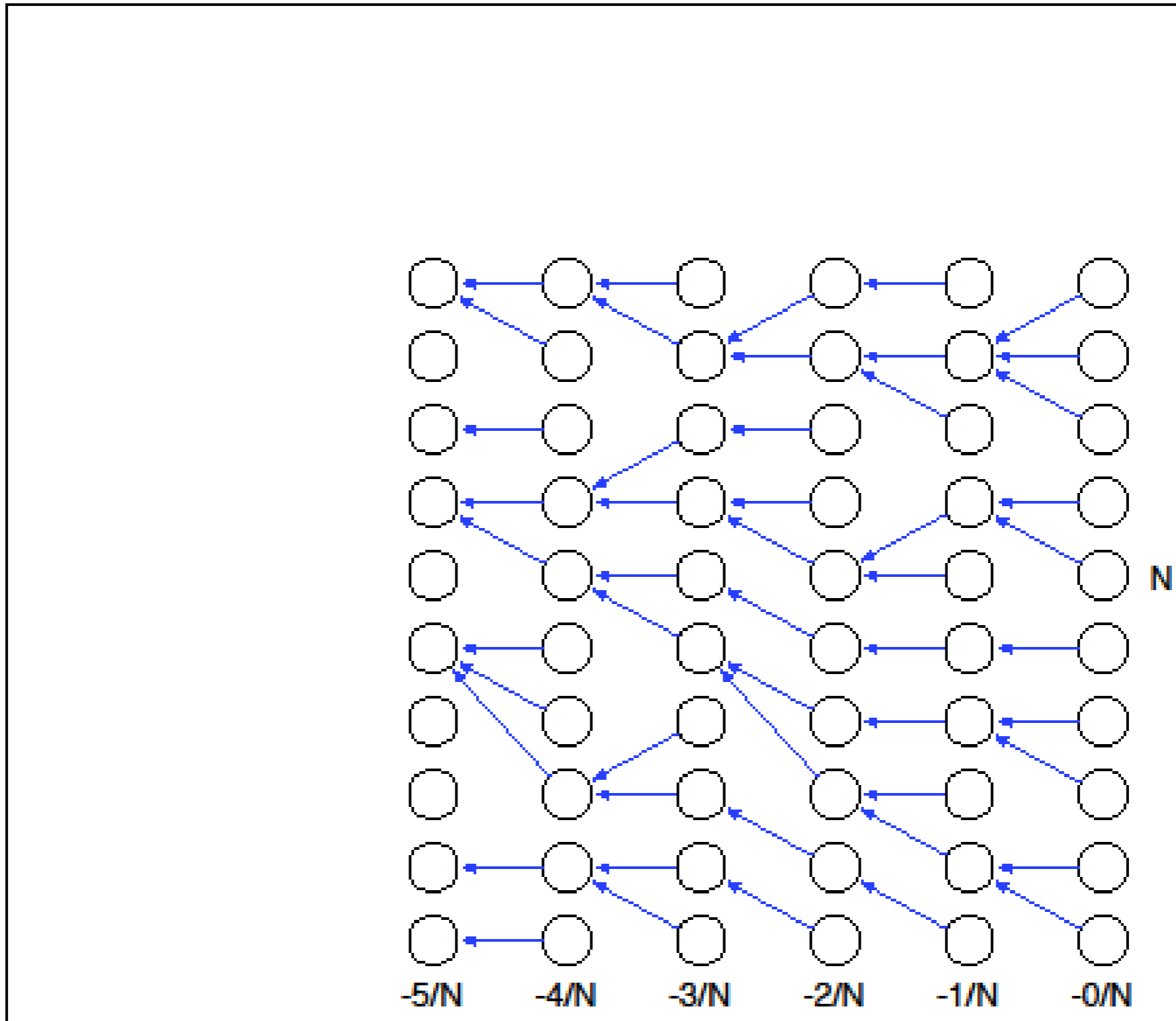


# A distribution over trees



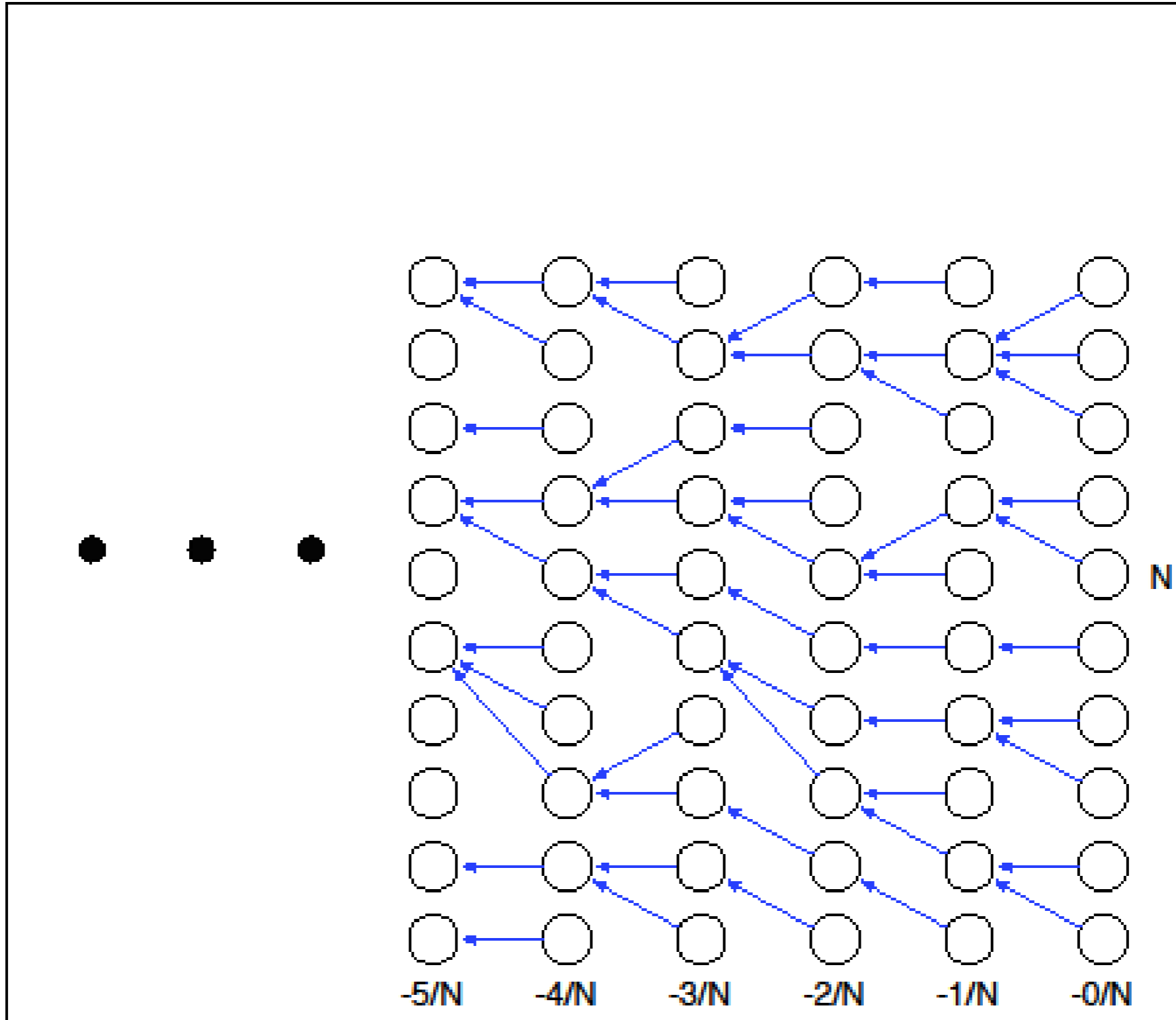


# A distribution over trees



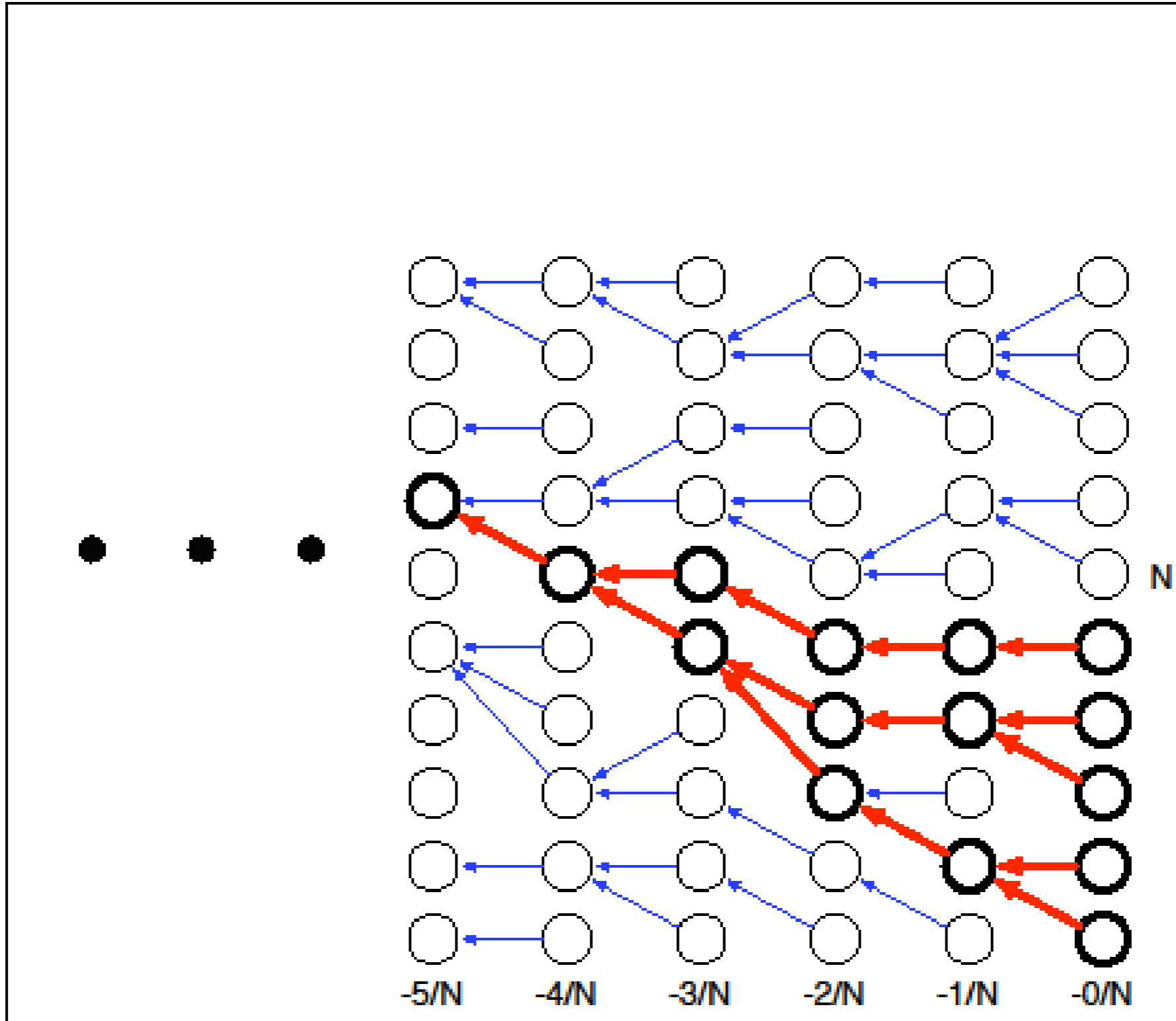


# A distribution over trees



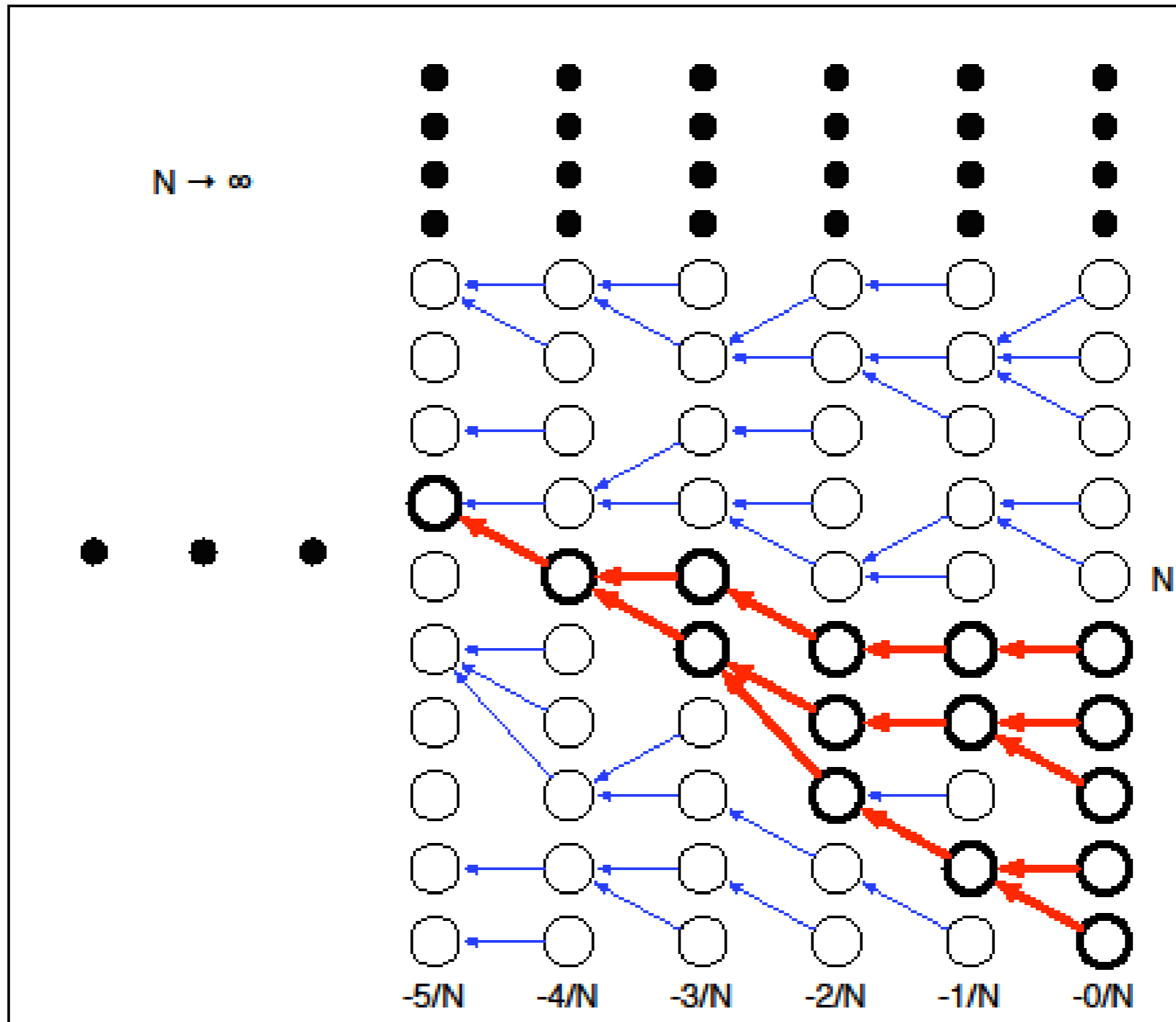


# A distribution over trees



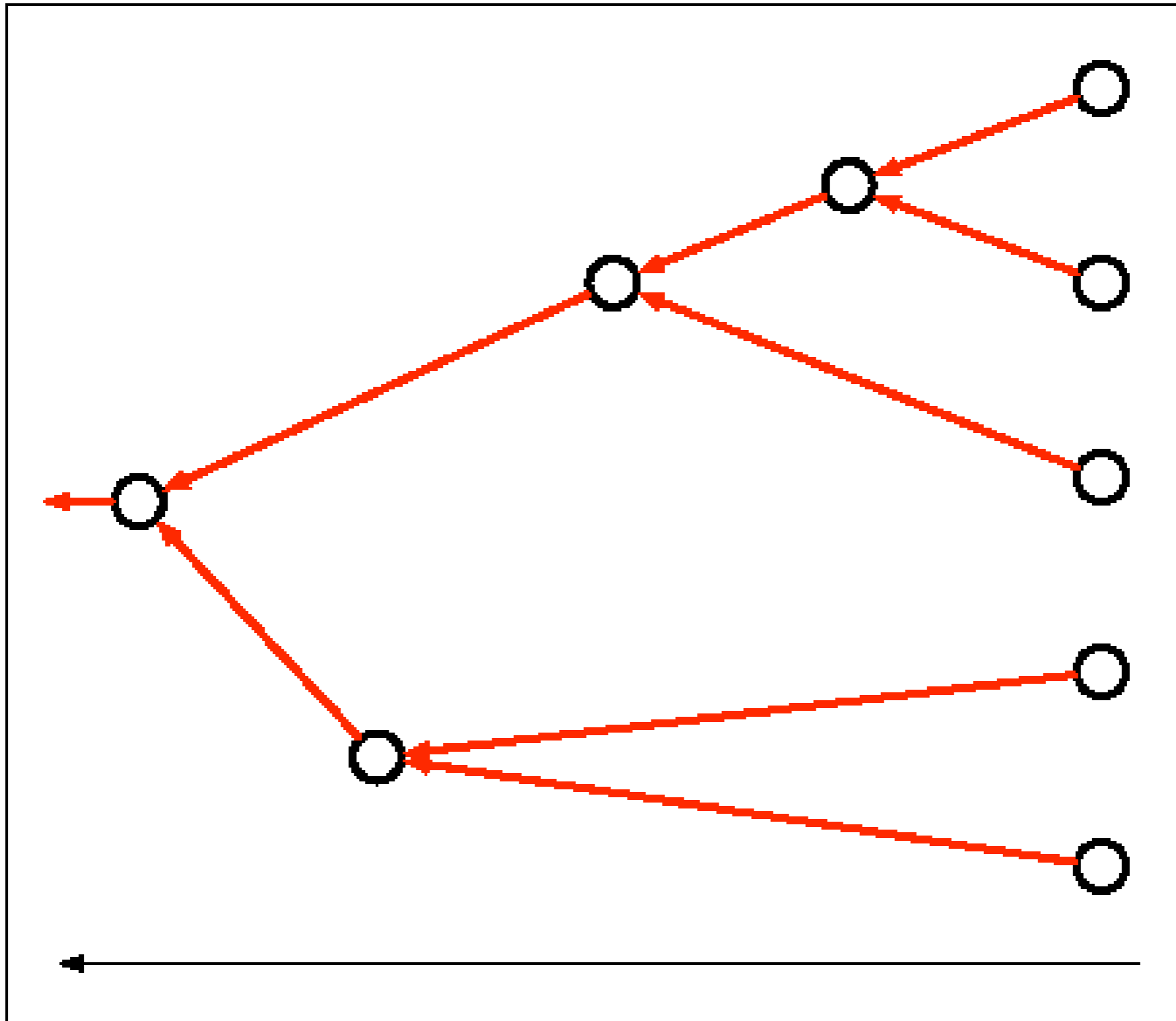


# A distribution over trees



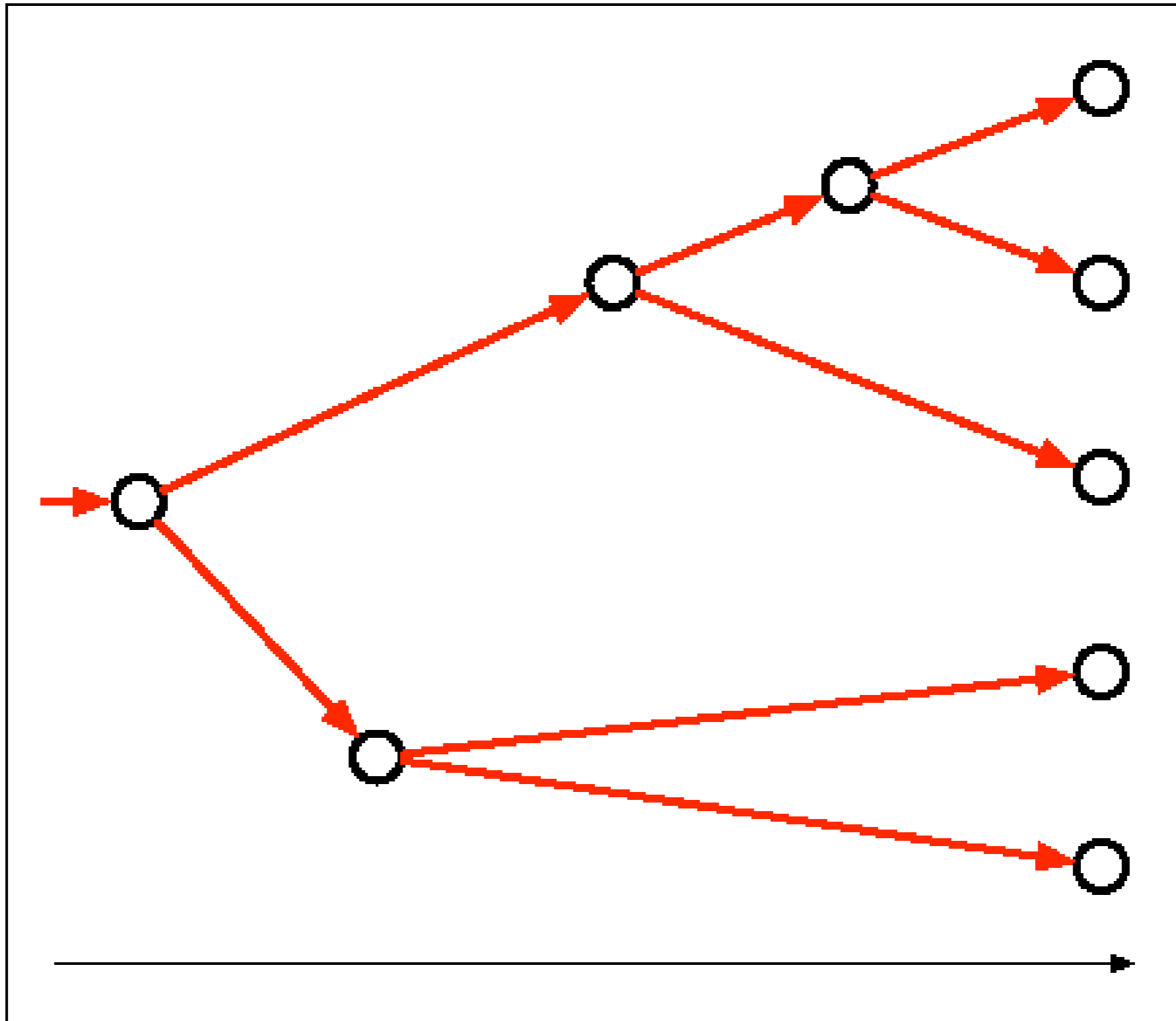


# Coalescent as a graphical model



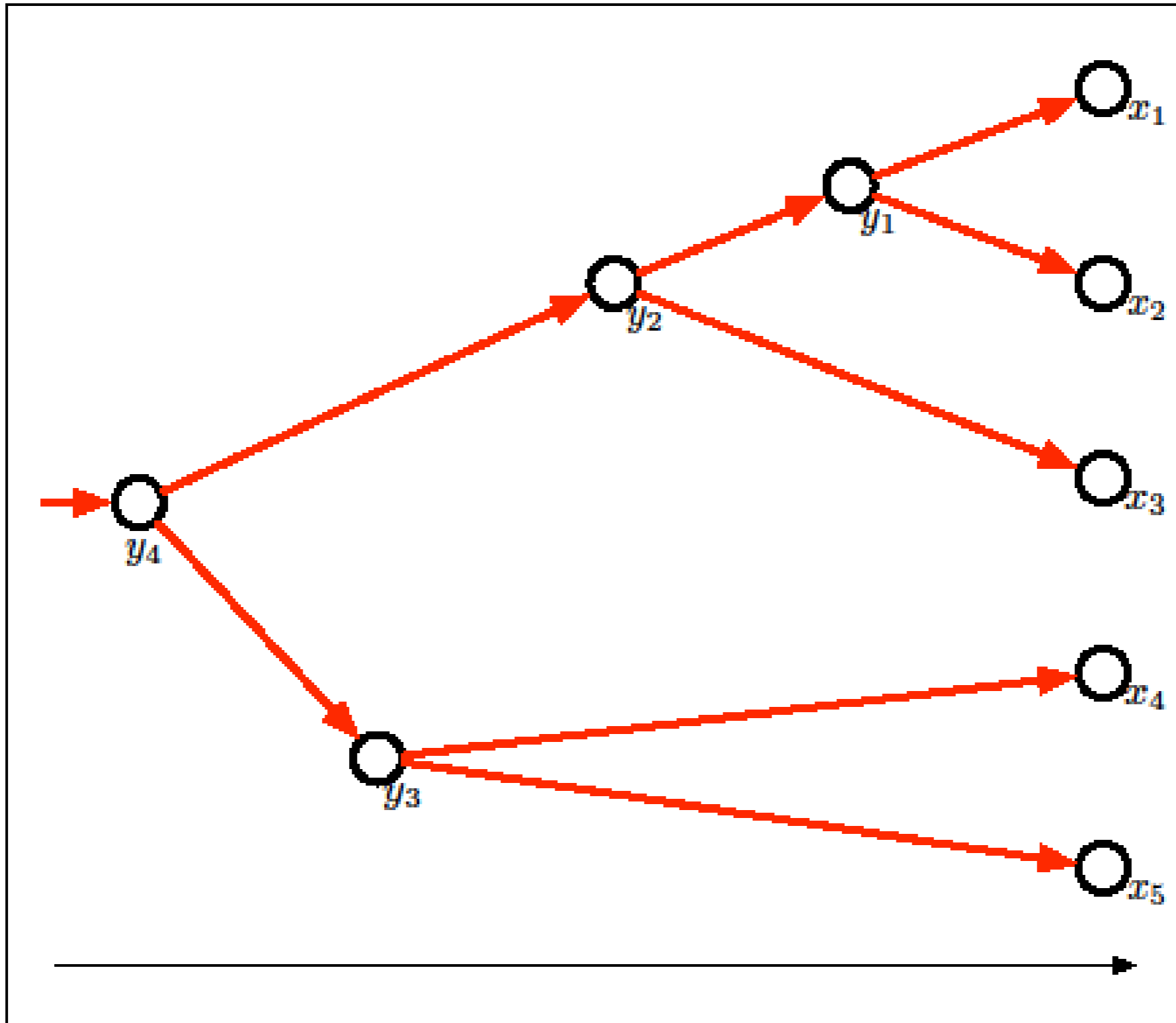


# Coalescent as a graphical model





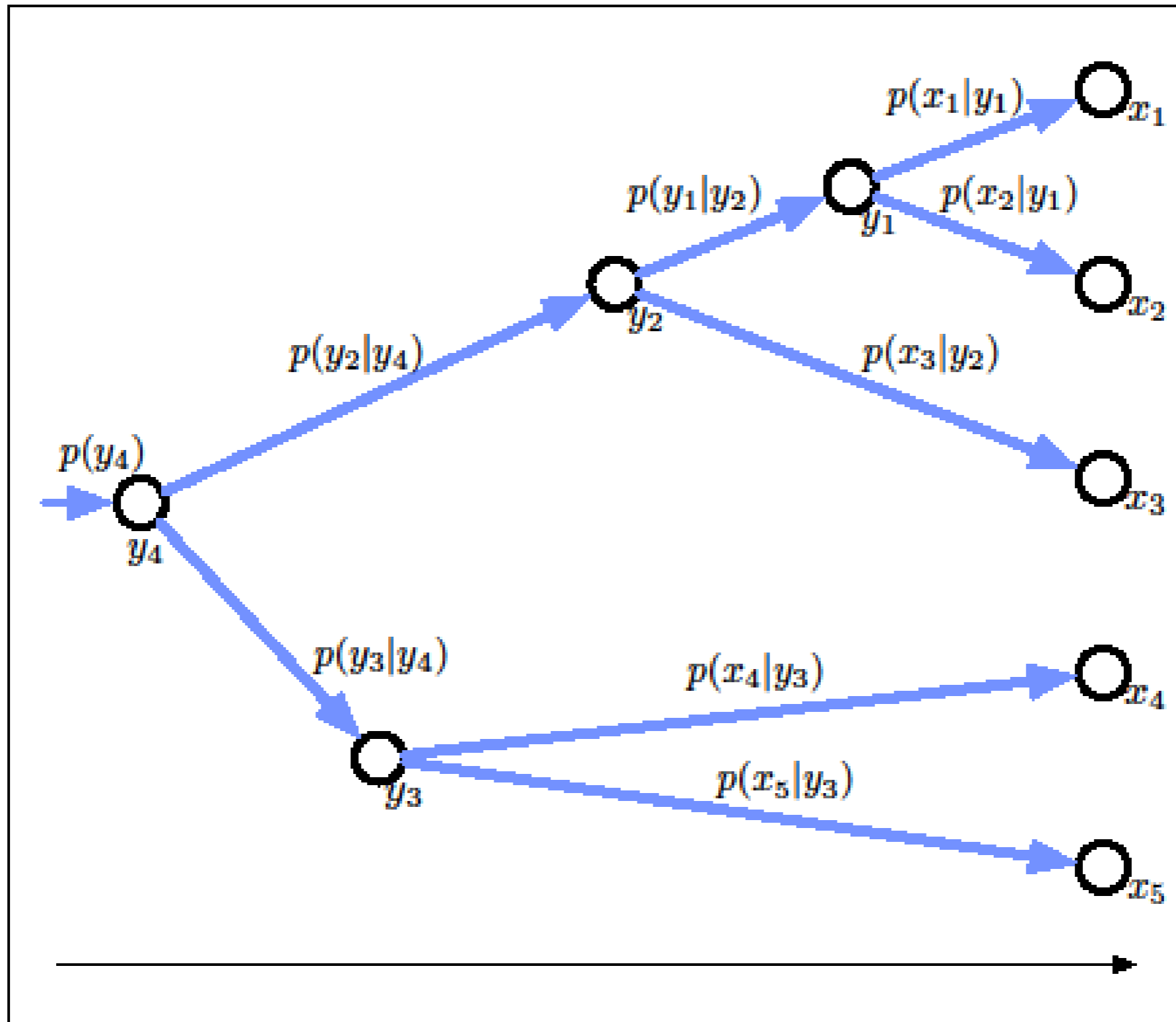
# Coalescent as a graphical model







# Coalescent as a graphical model

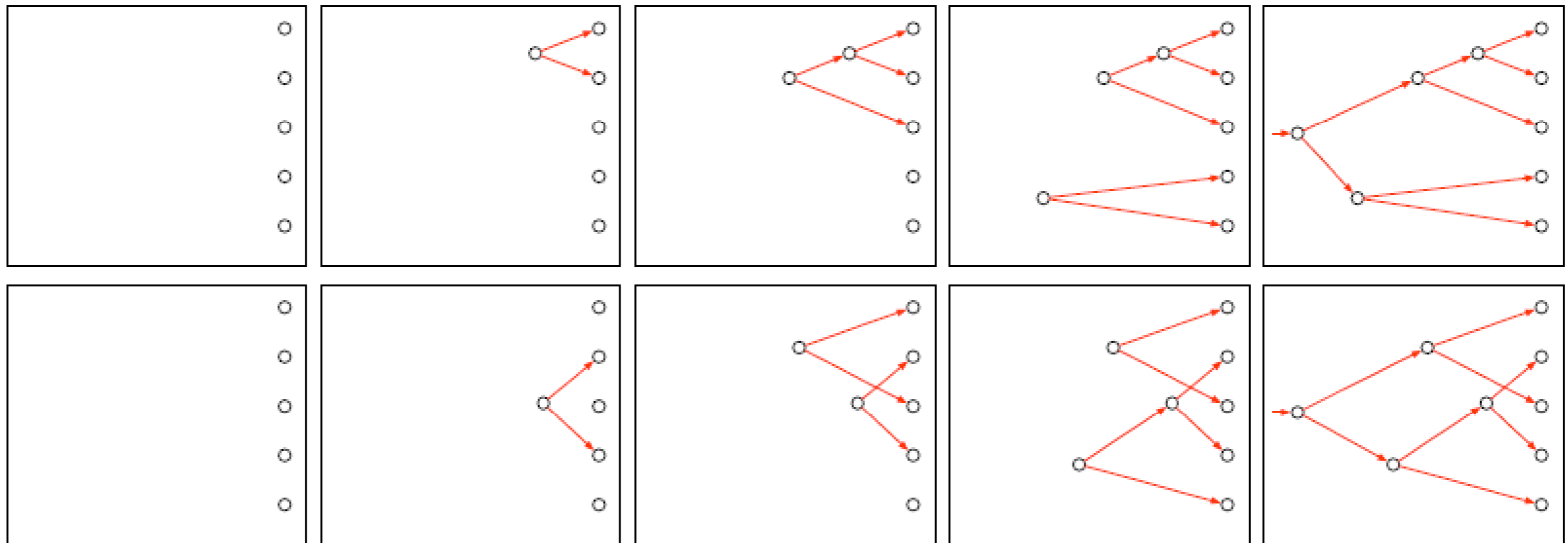




# Efficient Inference



- Construct trees in a bottom-up manner



- Greedy: At each step, pick optimal pair (maximizes joint likelihood) and time to coalesce (branch length)
- Infer values of internal nodes by belief propagation



M

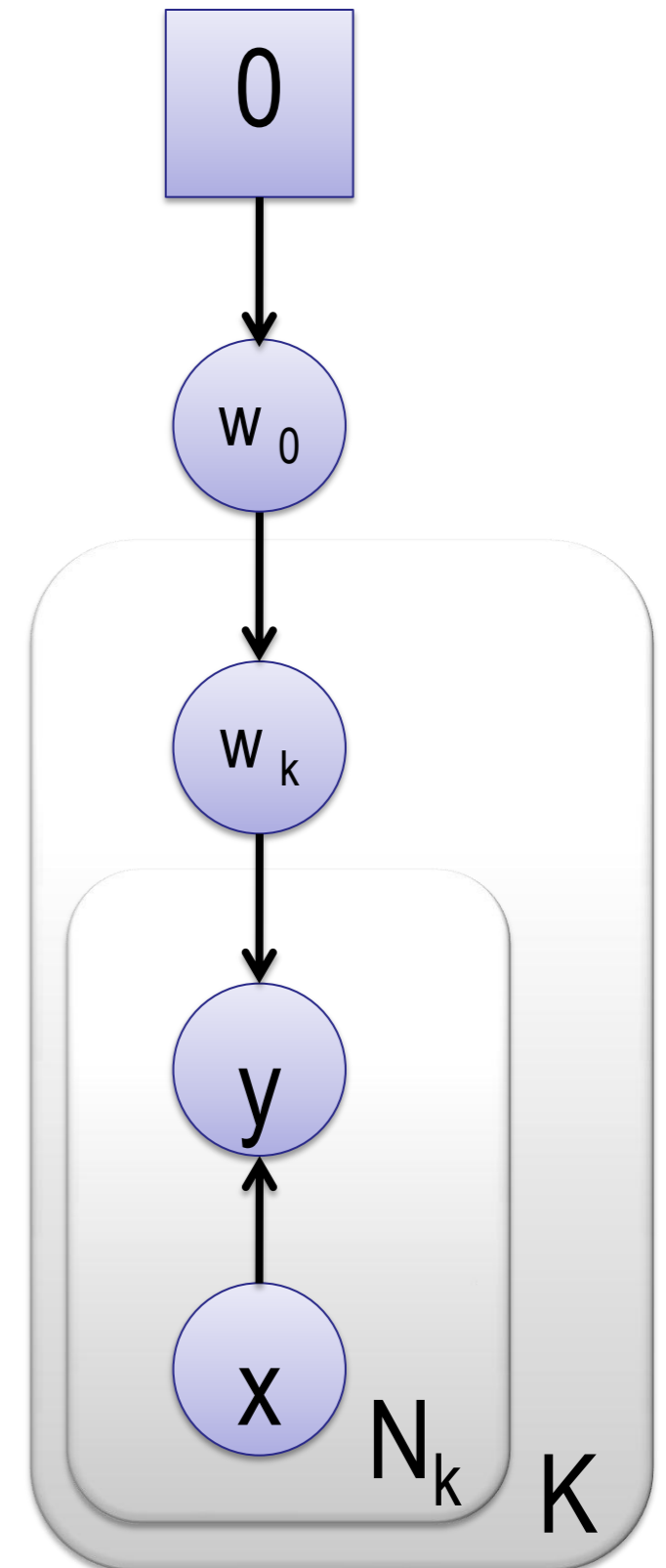
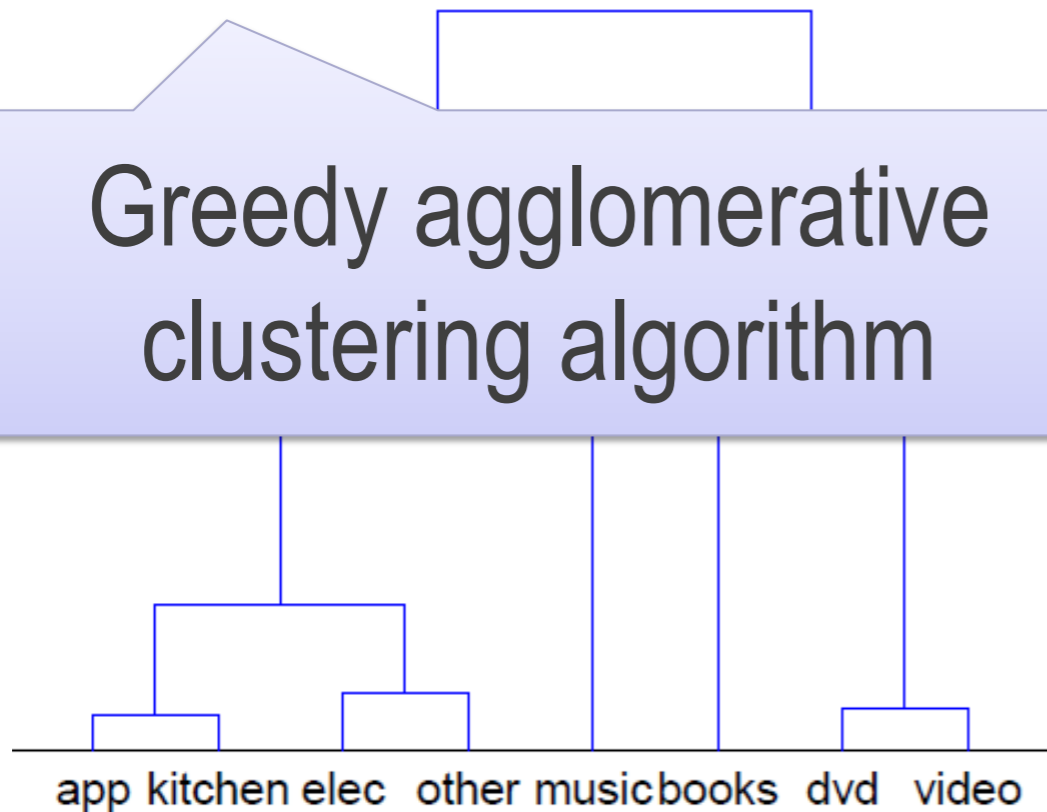
Message passing on coalescent tree; efficiently done by belief propagation

ated equal



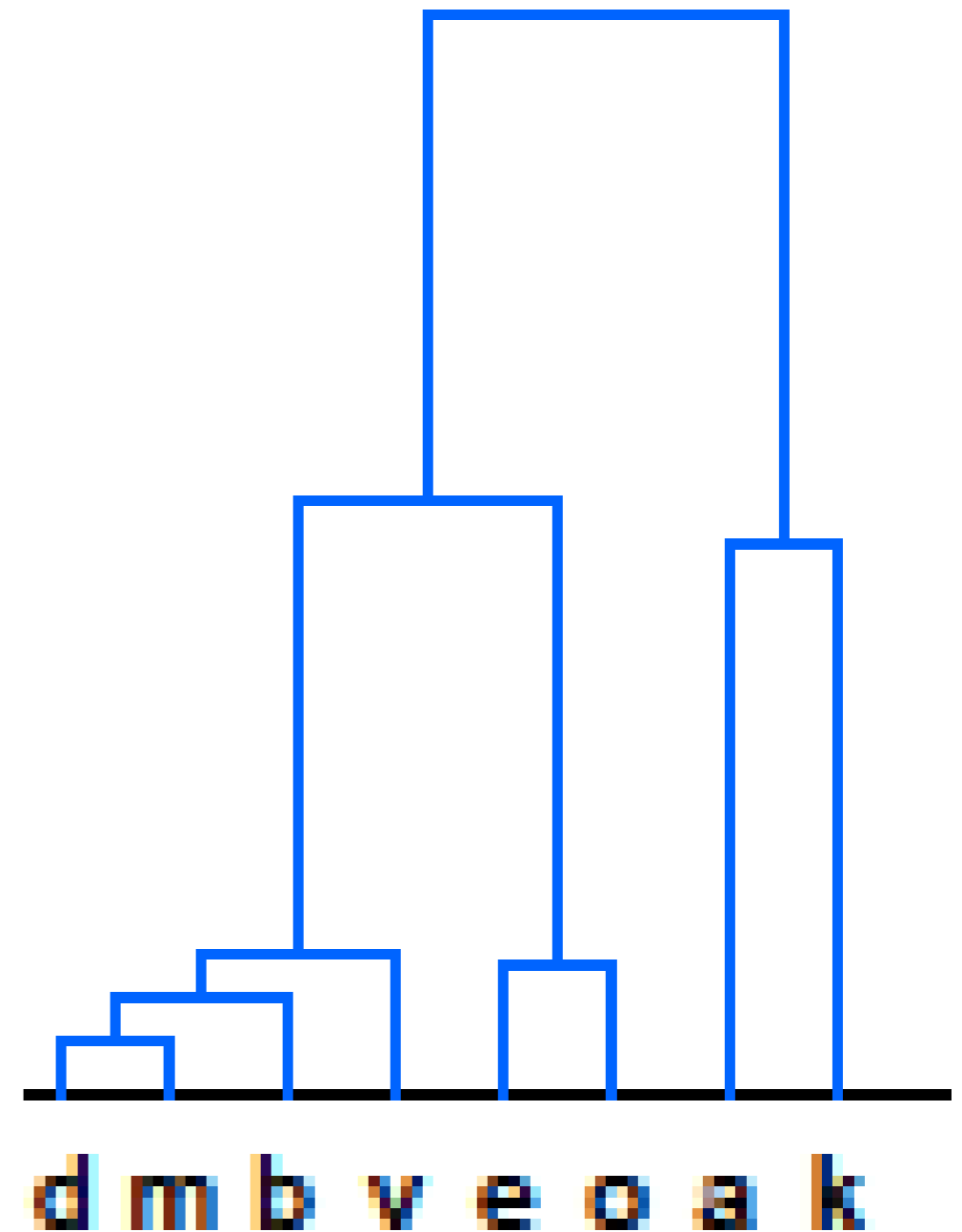
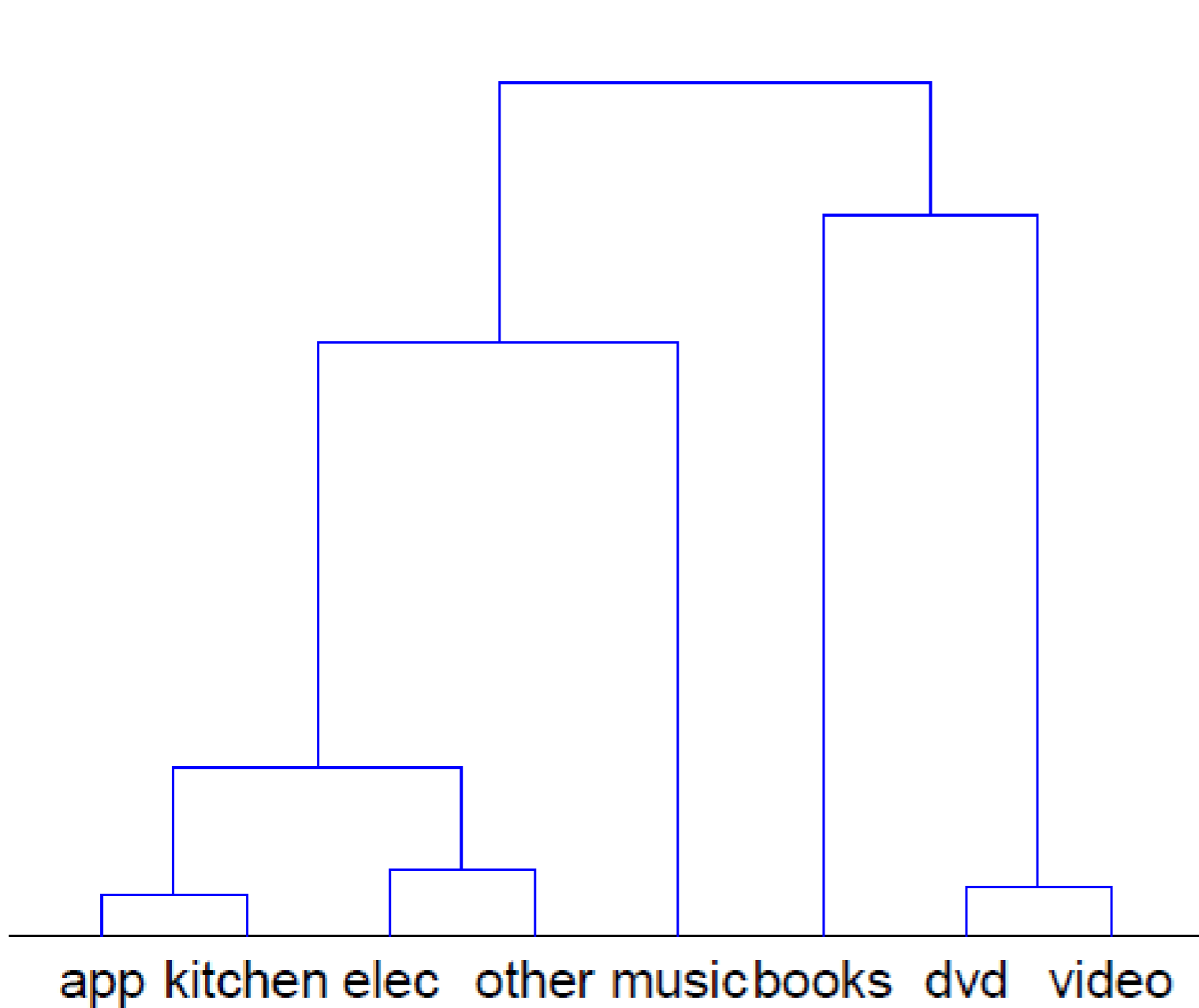
- Inference by EM:
- E: compute expectations over weights
- M: maximize tree structure

Greedy agglomerative clustering algorithm



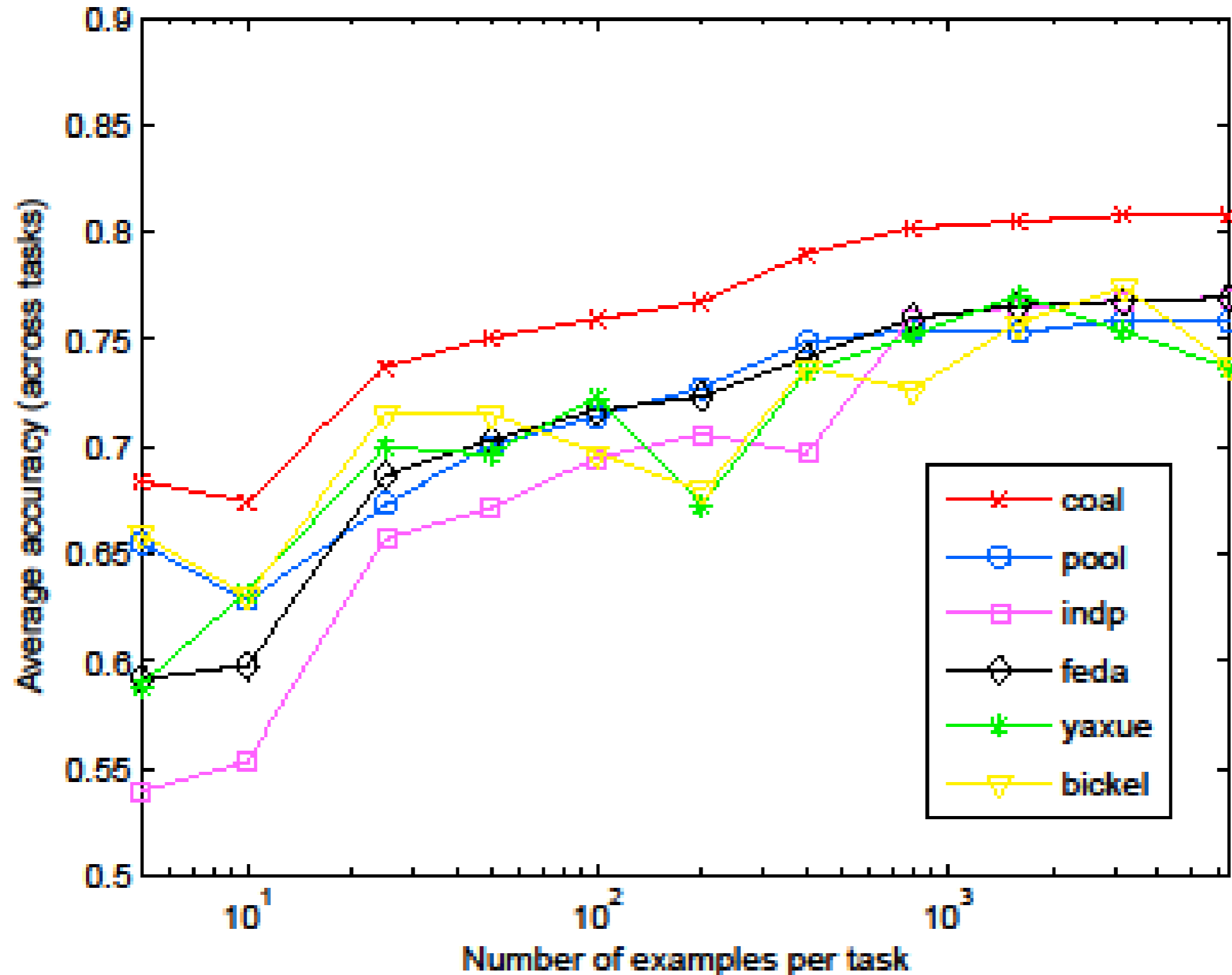


# Data tree *versus* inferred tree





# Some experimental results





# Parameter-based References

---



- O. Chapelle and Z. Harchaoui. A machine learning approach to conjoint analysis. NIPS, 2005.
- K. Yu, V. Tresp, and A. Schwaighofer. Learning Gaussian processes from multiple tasks. ICML, 2005.
- Y. Xue, X. Liao, L. Carin, and B. Krishnapuram. Multi-task learning for classification with Dirichlet process priors. JMLR, 2007.
- H. Daumé III. Bayesian Multitask Learning with Latent Hierarchies. UAI 2009.



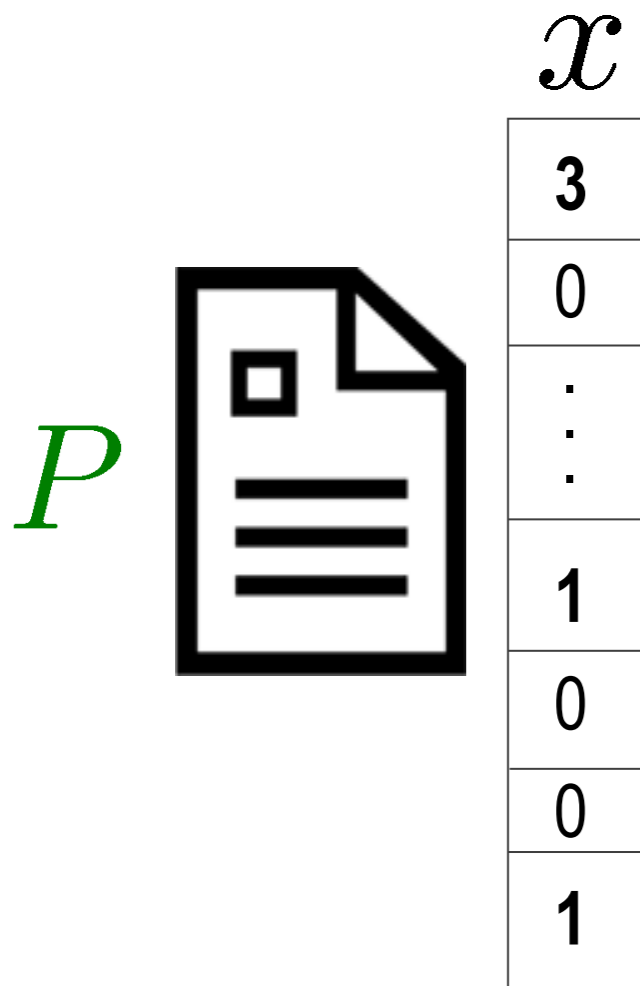
# Tutorial Outline

---

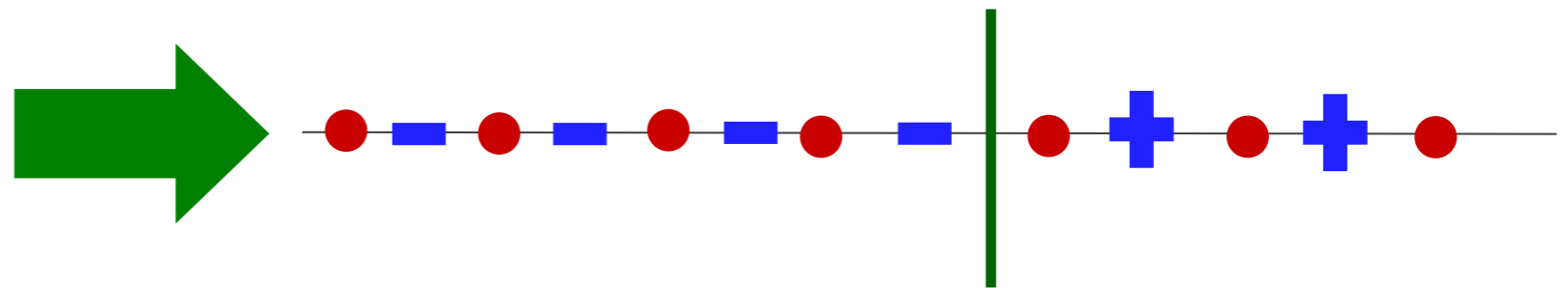


1. Notation and Common Concepts
2. Semi-supervised Adaptation
  - Covariate shift
  - Learning Shared Representations
3. Supervised Adaptation
  - Feature-Based Approaches
  - Parameter-Based Approaches
4. Open Questions and Uncovered Algorithms

Hypothesis classes from projections  $P : \theta^\top P x$



$$P = \hat{\theta} \hat{\theta}^\top \quad \hat{\theta} = \text{source ERM}$$



- 1) Minimize divergence
- 2)  $\hat{\epsilon}_S(\theta)$  small

~~$\epsilon_{P,T}(\theta^*) - \epsilon_{I,T}(\theta^*)$  small~~





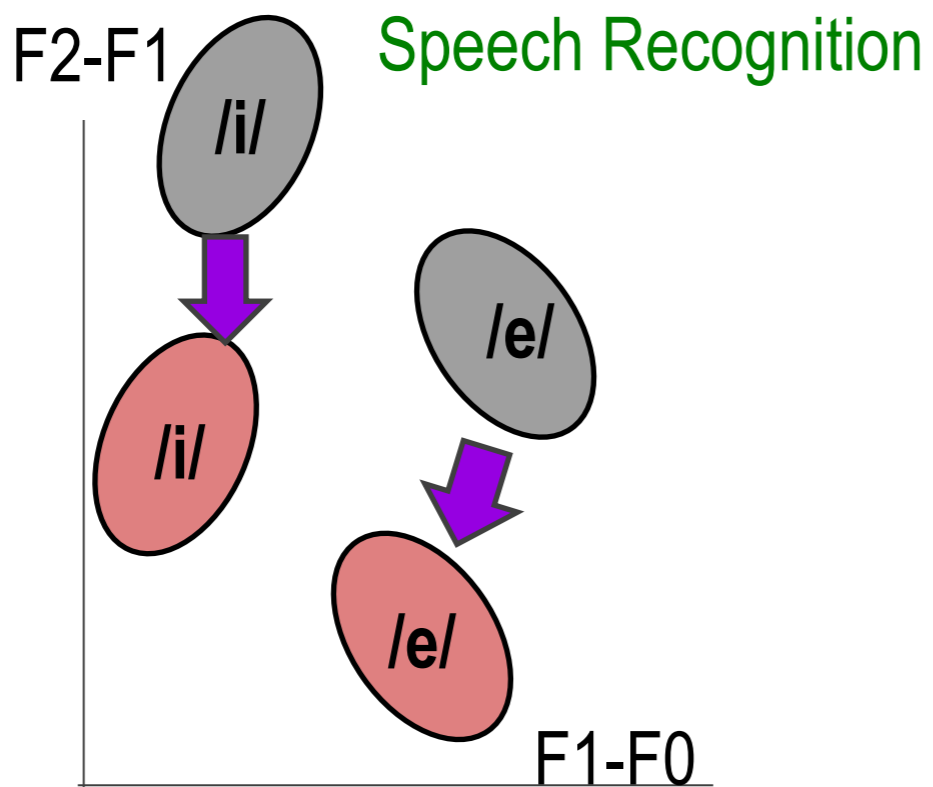
# Open Questions



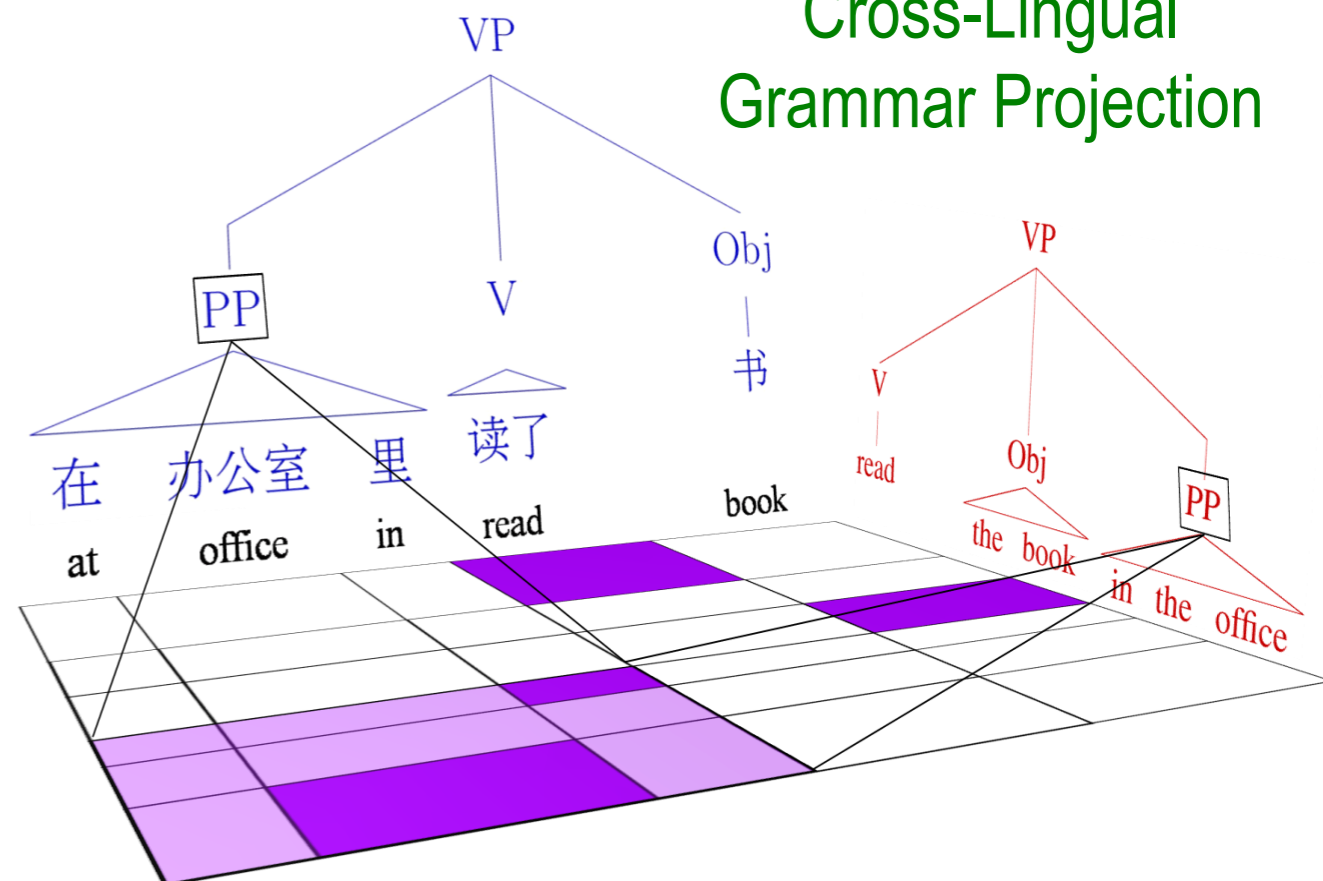
## Matching Theory and Practice

Theory does not exactly suggest what practitioners do

## Prior Knowledge



## Cross-Lingual Grammar Projection





# More Semi-supervised Adaptation

---



<http://adaptationtutorial.blitzer.com/references/>

## Self-training and Co-training

- [1] D. McClosky et al. Reranking and Self-Training for Parser Adaptation. 2006.
- [2] K. Sagae & J. Tsuji. Dependency Parsing and Domain Adaptation with LR Models and Parser Ensembles. 2007.

## Structured Representation Learning

- [3] F. Huang and A. Yates. Distributional Representations for Handling Sparsity in Supervised Sequence Labeling. 2009.



# What is a domain anyway?



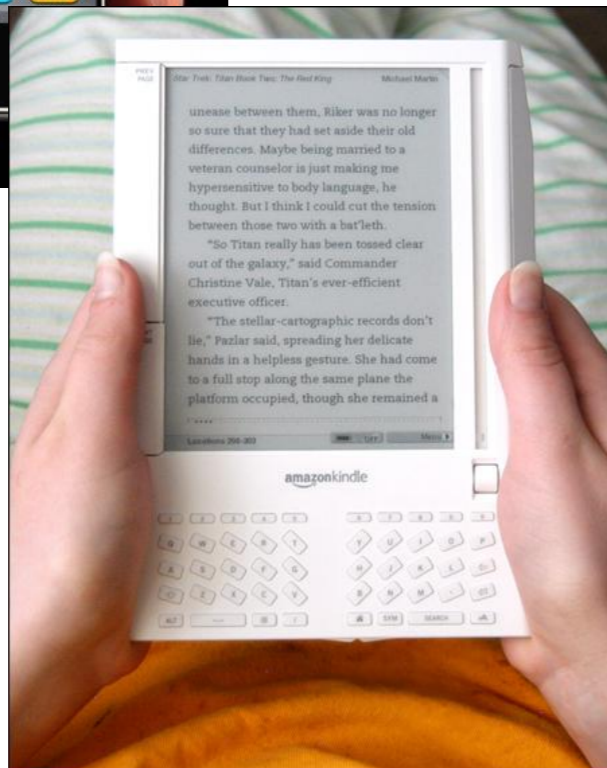
- Time?
  - News the day I was born vs news today?
  - News yesterday vs news today?
- Space?
  - News back home vs news in Haifa?
  - News in Tel Aviv vs news in Haifa?
- Do my data even come with a domain specified?

Suggest a continuous structure

Stream of  $\langle x, y, d \rangle$  data with  $y$  and  $d$  sometimes hidden?



# We're *all* domains: personalization



- adapt learn across millions of “domains”?
- share enough information to be useful?
- share little enough information to be safe?
- avoid negative transfer?
- avoid DAAM (domain adaptation spam)?



Thanks

---



Questions?