



A Kernel-based Approach to Sequence Labeling Problems

(Joint work with *Yuta Tsuboi*)

Hisashi Kashima

IBM Research

Tokyo Research Laboratory

hkashima@jp.ibm.com

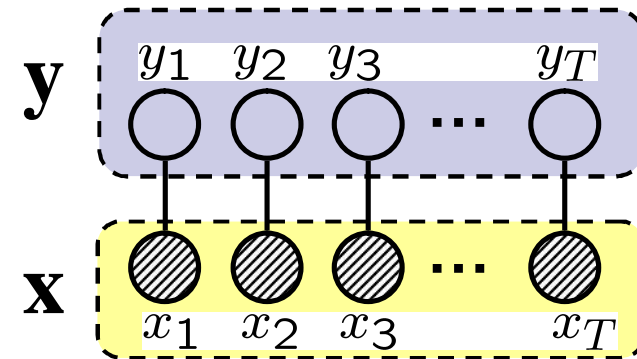


Outline

- Sequence labeling problems
- Hidden Markov model (HMM)
- Hidden Markov (HM) perceptron
- Marginalized labeling perceptron
- Experiments on natural language processing tasks
- Application to bioinformatics

Sequence Labeling Problems

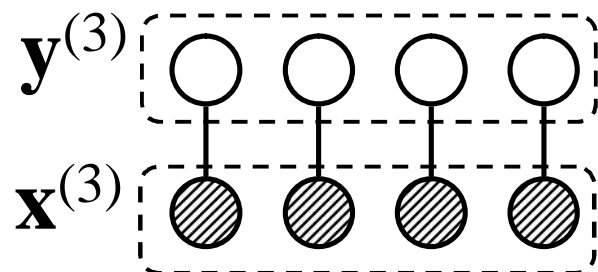
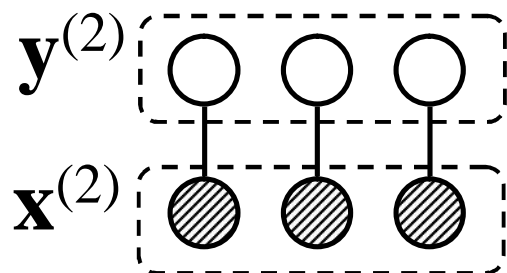
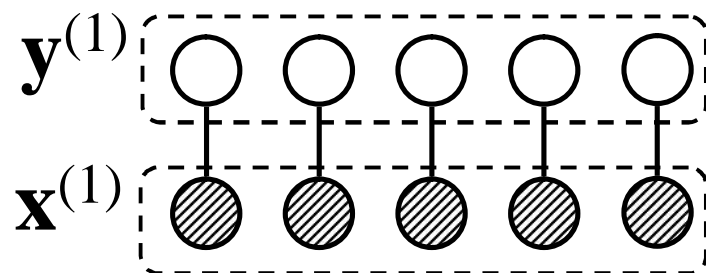
- Mapping $M : \Sigma_x^{|\mathbf{x}|} \rightarrow \Sigma_y^{|\mathbf{x}|}$ ($|\mathbf{x}|=|\mathbf{y}|=T$)
 - from a sequence of observed variables $\mathbf{x} = (x_1, x_2, x_3, \dots, x_T)$
 - to a sequence of hidden variables $\mathbf{y} = (y_1, y_2, y_3, \dots, y_T)$



- Training data
 - Some correctly labeled sequences (pairs of (\mathbf{x}, \mathbf{y}))
 - $(\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), (\mathbf{x}^{(2)}, \mathbf{y}^{(2)}), (\mathbf{x}^{(3)}, \mathbf{y}^{(3)}), \dots$
where $|\mathbf{x}^{(i)}| = |\mathbf{y}^{(i)}| = T^{(i)}$

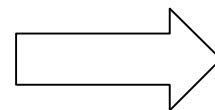
Sequence Labeling Problems

Training examples



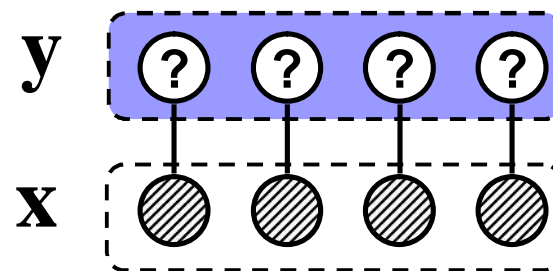
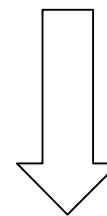
⋮

training



Model
(mapping)
M

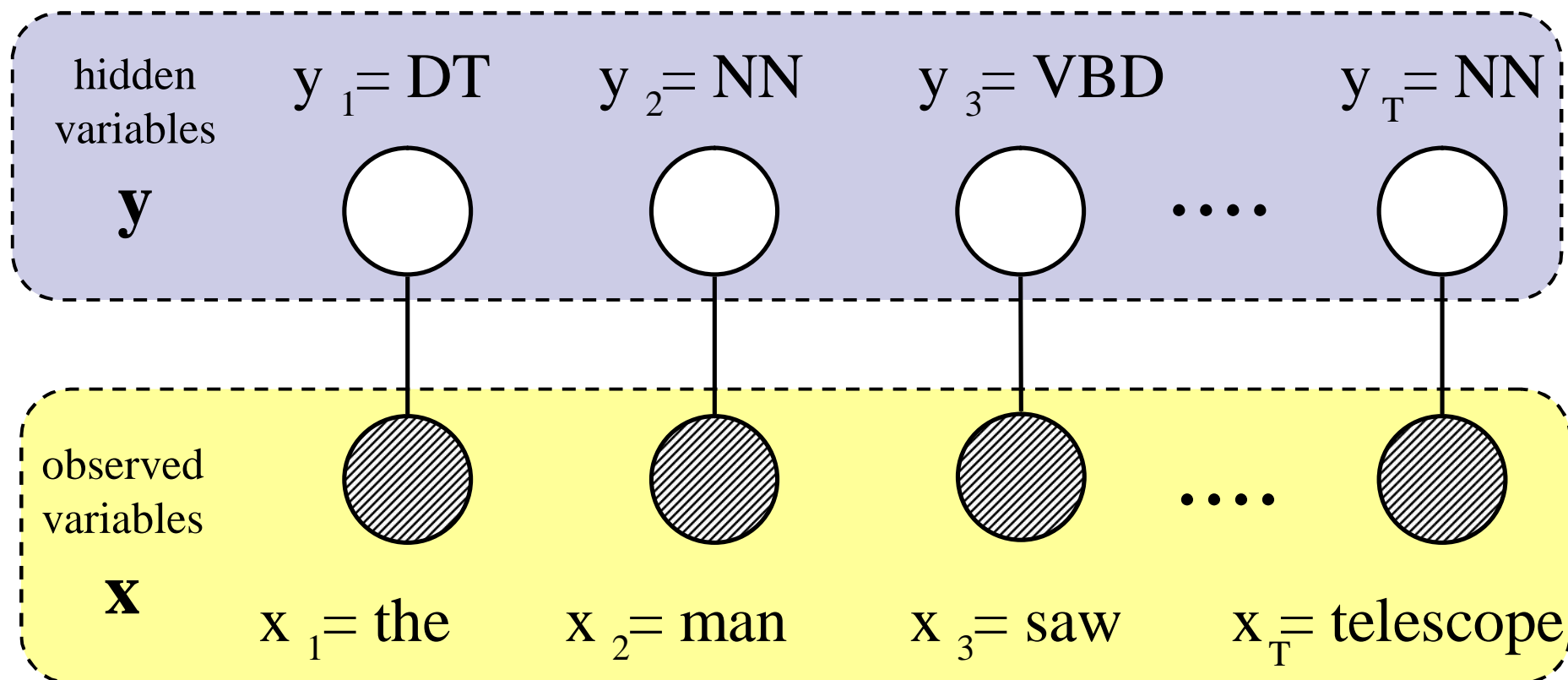
prediction



Unlabeled sequence

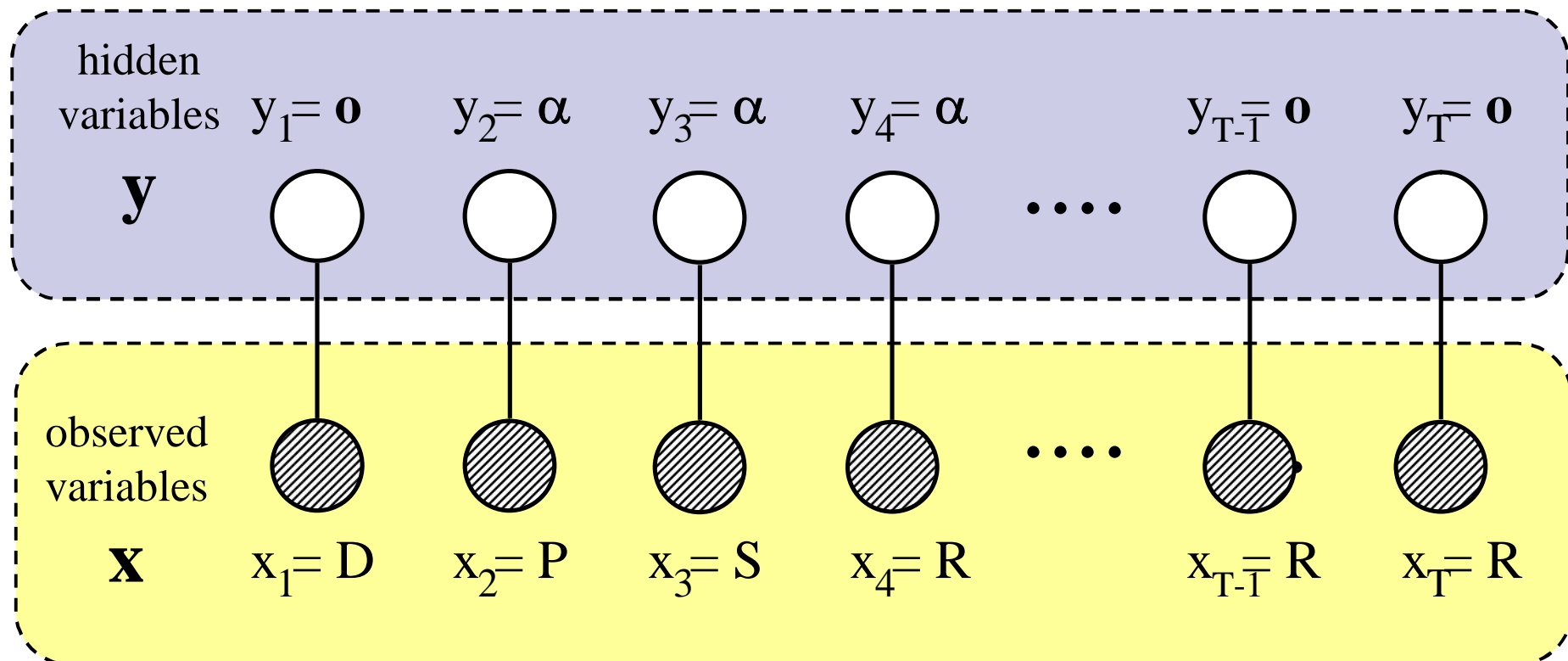
Part-of-Speech Tagging

- Observed variables \mathbf{x} = sequence of words
 - “the man saw the girl with a telescope.”
- Hidden variables \mathbf{y} = part-of-speech tags
 - assign one of part-of-speech tags to each word



Protein Secondary Structure Prediction

- Observed variables \mathbf{x} = protein sequence
 - $\mathbf{x} = (D, P, S, E, \dots, R, R)$
- Hidden variables \mathbf{y} = secondary structure
 - Assign one of “ α ” (α -helix), “ β ” (β -sheet) or “ o ” (other) to each amino acid



Hidden Markov Models for Labeling Sequences

- Joint probability distribution of (\mathbf{x}, \mathbf{y})

$$\Pr(\mathbf{x}, \mathbf{y}) = \Pr \left(\begin{array}{cccc} \text{DT} & \text{NN} & \text{VBD} & \text{NN} \\ \circ & \circ & \circ & \circ \\ | & | & | & | \\ \text{the} & \text{man} & \text{saw} & \text{glasses} \end{array} \right)$$

- Markov assumption (1st-order)

- Two types of probabilities

- Emission probabilities

- Probability of an observed variable emitted from a hidden variable
- e.g. Probability of “the” emitted from “DT”

- Transition probabilities

- Probability of a hidden variable following another hidden variable
- e.g. Probability of “NN” following “DT”

emission $\Pr \left(\begin{array}{c} \text{DT} \\ \circ \\ | \\ \text{the} \end{array} \right)$

transition $\Pr \left(\begin{array}{cc} \text{DT} & \text{NN} \\ \circ & \circ \end{array} \right)$

Markov Assumption

- Likelihood is decomposed into emission probabilities $\Pr(\text{shaded circle} \mid \text{circle})$ and transition probabilities $\Pr(\text{circle}-\text{circle})$

$$\log \Pr \left(\begin{array}{cccc} \text{DT} & \text{NN} & \text{VBD} & \text{NN} \\ \text{circle} & \text{circle} & \text{circle} & \text{circle} \\ | & | & | & | \\ \text{shaded circle} & \text{shaded circle} & \text{shaded circle} & \text{shaded circle} \\ \text{the} & \text{man} & \text{saw} & \text{glasses} \end{array} \right) =$$

$$\log \Pr \left(\begin{array}{c} \text{DT} \\ \text{circle} \\ | \\ \text{shaded circle} \\ \text{the} \end{array} \right) + \log \Pr \left(\begin{array}{c} \text{NN} \\ \text{circle} \\ | \\ \text{shaded circle} \\ \text{man} \end{array} \right) + \log \Pr \left(\begin{array}{c} \text{VBD} \\ \text{circle} \\ | \\ \text{shaded circle} \\ \text{saw} \end{array} \right) + \log \Pr \left(\begin{array}{c} \text{NN} \\ \text{circle} \\ | \\ \text{shaded circle} \\ \text{glasses} \end{array} \right)$$

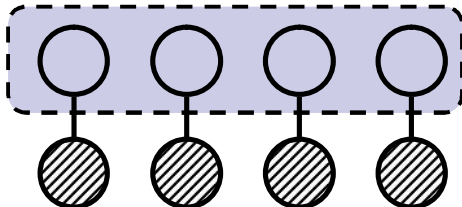
Emission probabilities

$$+ \log \Pr \left(\begin{array}{cc} \text{DT} & \text{NN} \\ \text{circle} & \text{circle} \end{array} \right) + \log \Pr \left(\begin{array}{cc} \text{NN} & \text{VBD} \\ \text{circle} & \text{circle} \end{array} \right) + \log \Pr \left(\begin{array}{cc} \text{VBD} & \text{NN} \\ \text{circle} & \text{circle} \end{array} \right)$$

Transition probabilities

Prediction (sequence labeling) with HMM

- Given \mathbf{x} , output \mathbf{y} that maximizes $\Pr(\mathbf{x}, \mathbf{y})$
 - Among all possible ways of labeling

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y}} \log \Pr \left(\begin{array}{c} \mathbf{y} \\ \text{the} \quad \text{man} \quad \text{saw} \quad \text{glasses} \end{array} \right)$$


Prediction (sequence labeling) with HMM

- Combine emission & transition features (= building blocks) to maximize the total probability
- Viterbi algorithm for finding the best hidden labels
 - Dynamic programming

$$\hat{y} = \operatorname{argmax}_y \log \Pr \left(\begin{array}{c} \text{y} \\ \boxed{y_1 \quad y_2 \quad y_3 \quad y_4} \\ \text{the} \quad \text{man} \quad \text{saw} \quad \text{glasses} \end{array} \right) =$$

$$\log \Pr \left(\boxed{y_1} \right) + \log \Pr \left(\boxed{y_2} \right) + \log \Pr \left(\boxed{y_3} \right) + \log \Pr \left(\boxed{y_4} \right)$$

Emission probabilities

$$+ \log \Pr \left(\boxed{y_1 \quad y_2} \right) + \log \Pr \left(\boxed{y_2 \quad y_3} \right) + \log \Pr \left(\boxed{y_3 \quad y_4} \right)$$

Transition probabilities

Training (Parameter Estimation)

- Maximum likelihood estimation of the joint probability $\Pr(\mathbf{x}, \mathbf{y})$

- Parameters

- Emission probabilities $\Pr(\text{O} \mid \text{shaded circle})$

- Transition probabilities $\Pr(\text{O} \rightarrow \text{O})$

are determined to maximize the joint probability of training data

- Simple

- count the frequencies of features in training data



Training (Parameter Estimation)

■ Hidden Markov models

- ☐ Maximum likelihood estimation of the joint probability of (\mathbf{x}, \mathbf{y})
- ☐ Our purpose: prediction of \mathbf{y} given \mathbf{x}
 - Only needs conditional model $P(\mathbf{y} | \mathbf{x})$
- ☐ MLE of $\Pr(\mathbf{x}, \mathbf{y})$ tackles more difficult problem
 - Tends to need many examples

■ Conditional models

- ☐ More suited to our purpose
- ☐ Estimate conditional model $P(y|x)$ directly
- ☐ Maximum Entropy Markov Model (MEMM) [McCallum et al, 2000]
- ☐ Conditional Random Field (CRF) [Lafferty et al, 2001]
- ☐ HM-SVM [Altun et al, 2003]
- ☐ **HM-Perceptron** [Collins, 2002]
 - Alternative training algorithm for HMM
 - Online algorithm

Another interpretation of HMM

- Log likelihood can be written as the inner product of **weight vector** \mathbf{w} & **feature vector** $\Phi(\mathbf{x}, \mathbf{y})$

$$\log \Pr \left(\begin{array}{c} \text{DT} \quad \text{NN} \quad \text{VBD} \quad \text{NN} \\ \text{the} \quad \text{man} \quad \text{saw} \quad \text{glasses} \end{array} \right) = \langle \mathbf{w} \cdot \Phi(\mathbf{x}, \mathbf{y}) \rangle$$

weight vector
(model parameters)

$$\mathbf{W} = \begin{pmatrix} \log \Pr \left(\begin{array}{c} \text{DT} \\ \text{the} \end{array} \right) \\ \log \Pr \left(\begin{array}{c} \text{NN} \\ \text{man} \end{array} \right) \\ \log \Pr \left(\begin{array}{c} \text{VBD} \quad \text{NN} \end{array} \right) \\ \vdots \end{pmatrix}$$

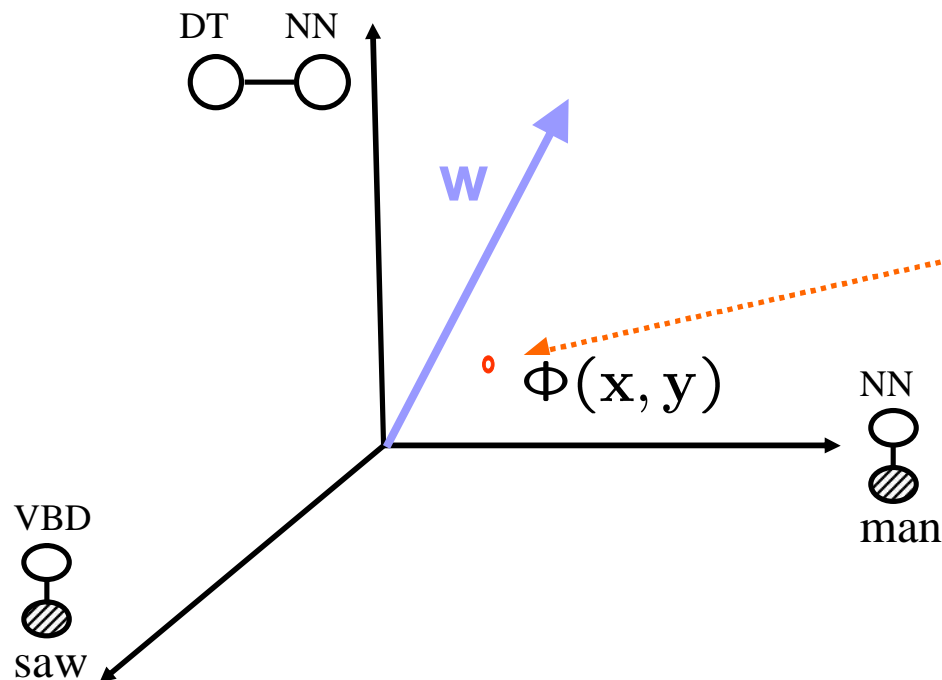
feature vector
(labeled sequence)

$$\Phi(\mathbf{x}, \mathbf{y}) = \begin{pmatrix} \# \left(\begin{array}{c} \text{DT} \\ \text{the} \end{array} \right) \\ \# \left(\begin{array}{c} \text{NN} \\ \text{man} \end{array} \right) \\ \# \left(\begin{array}{c} \text{VBD} \quad \text{NN} \end{array} \right) \\ \vdots \end{pmatrix}$$

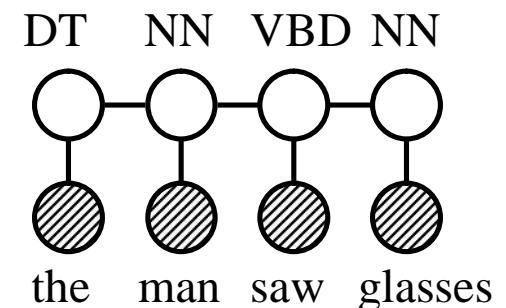
Number of times $\begin{array}{c} \text{VBD} \quad \text{NN} \\ \text{O} \text{---} \text{O} \end{array}$ appearing

Another interpretation of HMM

- A labeled sequence (\mathbf{x}, \mathbf{y}) is mapped into a feature space as a **feature vector** $\Phi(\mathbf{x}, \mathbf{y})$ by using the number of times each feature appears in labeled sequence (\mathbf{x}, \mathbf{y})
- In prediction, feature vectors are projected onto \mathbf{w} , \mathbf{y} with the highest $\langle \mathbf{w} \cdot \Phi(\mathbf{x}, \mathbf{y}) \rangle$ is output
- In training, \mathbf{w} is trained to maximize the joint likelihood of (correctly labeled) training data $\sum_i \langle \mathbf{w} \cdot \Phi(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) \rangle$



labeled sequence
 (\mathbf{x}, \mathbf{y})



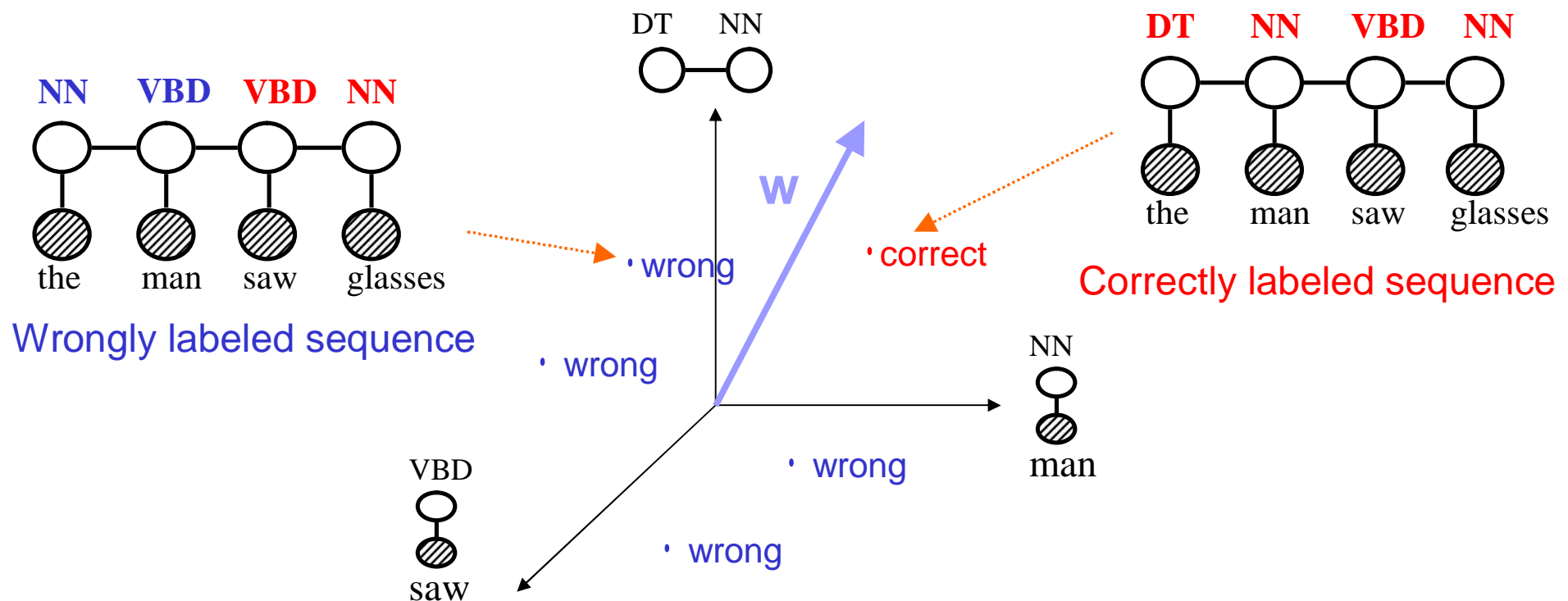
Hidden Markov (HM) Perceptron [Collins, 2002]

- Training algorithm of \mathbf{w} suited for conditional prediction
- \mathbf{w} is trained so that **correctly labeled sequence** has higher score than all **wrongly labeled sequences**

□ For each training example $(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})$,

$$\langle \mathbf{w} \cdot \Phi(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) \rangle > \langle \mathbf{w} \cdot \Phi(\mathbf{x}^{(i)}, \mathbf{y}^{\text{wrong}}) \rangle \text{ for } \forall \mathbf{y}^{\text{wrong}} \neq \mathbf{y}^{(i)}$$

correct labeling wrong labeling



Algorithm of HM-Perceptron

- Online algorithm

- Processes training examples one by one


- Prediction

- Given $\mathbf{x}^{(i)}$, output the prediction for $\mathbf{y}^{(i)}$ with the highest score according to the current \mathbf{w}

- Viterbi algorithm

$$\hat{\mathbf{y}}^{(i)}(\mathbf{x}) = \operatorname{argmax}_{\mathbf{y}} \langle \mathbf{w}, \Phi(\mathbf{x}, \mathbf{y}) \rangle$$

Among
all possible ways
of labeling for \mathbf{y}



- Update \mathbf{w} when the prediction is wrong

- If $\hat{\mathbf{y}}^{(i)} \neq \mathbf{y}^{(i)}$, (prediction for the i -th example is wrong)

$$\mathbf{w}^{new} = \mathbf{w}^{old} + \underbrace{\Phi(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})}_{\text{True labeling}} - \underbrace{\Phi(\mathbf{x}^{(i)}, \hat{\mathbf{y}}^{(i)})}_{\text{Wrong labeling}}$$

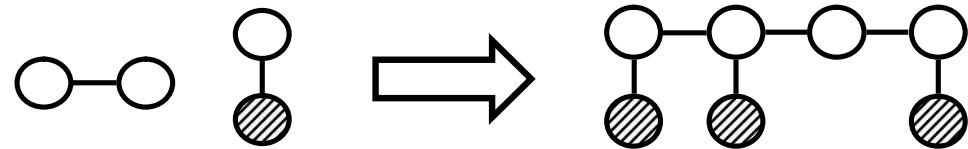
True labeling

Wrong labeling

HM-Perceptron with Large Features

- Arbitrarily **large features** to incorporate wider contexts

- idioms, motifs,...



- Two problems

1. Dimension of feature vector $\Phi(\mathbf{x}, \mathbf{y})$ becomes high
2. Computational complexity of label prediction

$$\hat{\mathbf{y}}(\mathbf{x}) = \operatorname{argmax}_{\mathbf{y}} \langle \mathbf{w}, \Phi(\mathbf{x}, \mathbf{y}) \rangle$$

depends **exponentially** on the max number of hidden variables contained in features

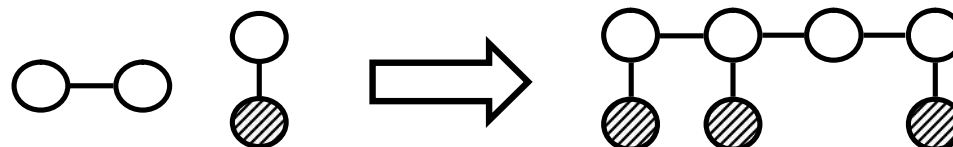
- Viterbi algorithm

⇒ HM-perceptron cannot handle arbitrarily large features

HM-Perceptron with Large Features

- Arbitrarily **large features** to incorporate wider contexts

- idioms, motifs,...



- Two problems

1. Dimension of feature vector $\Phi(\mathbf{x}, \mathbf{y})$ becomes too high

2. Computational complexity of label prediction

$$\hat{\mathbf{y}}(\mathbf{x}) = \operatorname{argmax}_{\mathbf{y}} \langle \mathbf{w}, \Phi(\mathbf{x}, \mathbf{y}) \rangle$$

depends **exponentially** on the max number of hidden variables contained in features

- Viterbi algorithm

⇒ HM-perceptron cannot handle arbitrarily large features

HM-Perceptron in Dual Form

$$\mathbf{w}^{new} = \mathbf{w}^{old} + \Phi(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) - \Phi(\mathbf{x}^{(i)}, \hat{\mathbf{y}}^{(i)})$$

- \mathbf{w} is represented as a linear combination of feature vectors

$$\mathbf{w} = \sum_{j=1}^{|Example|} \sum_{\tilde{\mathbf{y}}} \alpha_j(\tilde{\mathbf{y}}) \Phi(\mathbf{x}^{(j)}, \tilde{\mathbf{y}})$$

New weight

- Prediction

- Given \mathbf{x} , output the prediction of \mathbf{y} with the highest score

$$\hat{\mathbf{y}}(\mathbf{x}) = \operatorname{argmax}_{\mathbf{y}} \langle \mathbf{w}, \Phi(\mathbf{x}, \mathbf{y}) \rangle \text{ (primal form)}$$

$$= \operatorname{argmax}_{\mathbf{y}} \sum_{j=1}^{|Example|} \sum_{\tilde{\mathbf{y}}} \alpha_j(\tilde{\mathbf{y}}) \langle \Phi(\mathbf{x}^{(j)}, \tilde{\mathbf{y}}), \Phi(\mathbf{x}, \mathbf{y}) \rangle \text{ (dual form)}$$

Inner product

- Update

- If $\hat{\mathbf{y}}^{(i)} \neq \mathbf{y}^{(i)}$, (prediction for the i -th example is wrong)

- $\alpha_i^{new}(\mathbf{y}^{(i)}) = \alpha_i^{old}(\mathbf{y}^{(i)}) + 1$

- $\alpha_i^{new}(\hat{\mathbf{y}}^{(i)}) = \alpha_i^{old}(\hat{\mathbf{y}}^{(i)}) - 1,$



Kernel Methods

- Substitute inner products by “kernel function”

$$K((\mathbf{x}, y), (\mathbf{x}', y')) = \langle \Phi(\mathbf{x}, y), \Phi(\mathbf{x}', y') \rangle$$

- Dimension of feature space does not matter
 - As long as we have algorithms to compute inner product (=kernel function) without explicit construction of feature vectors

$$\hat{y}(\mathbf{x}) = \operatorname{argmax}_y \langle \mathbf{w}, \Phi(\mathbf{x}, y) \rangle \text{ (primal form)}$$

$$= \operatorname{argmax}_y \sum_{j=1}^{|Example|} \sum_{\tilde{y}} \alpha_j(\tilde{y}) K(\Phi(\mathbf{x}^{(j)}, \tilde{y}), \Phi(\mathbf{x}, y)) \text{ (dual form)}$$

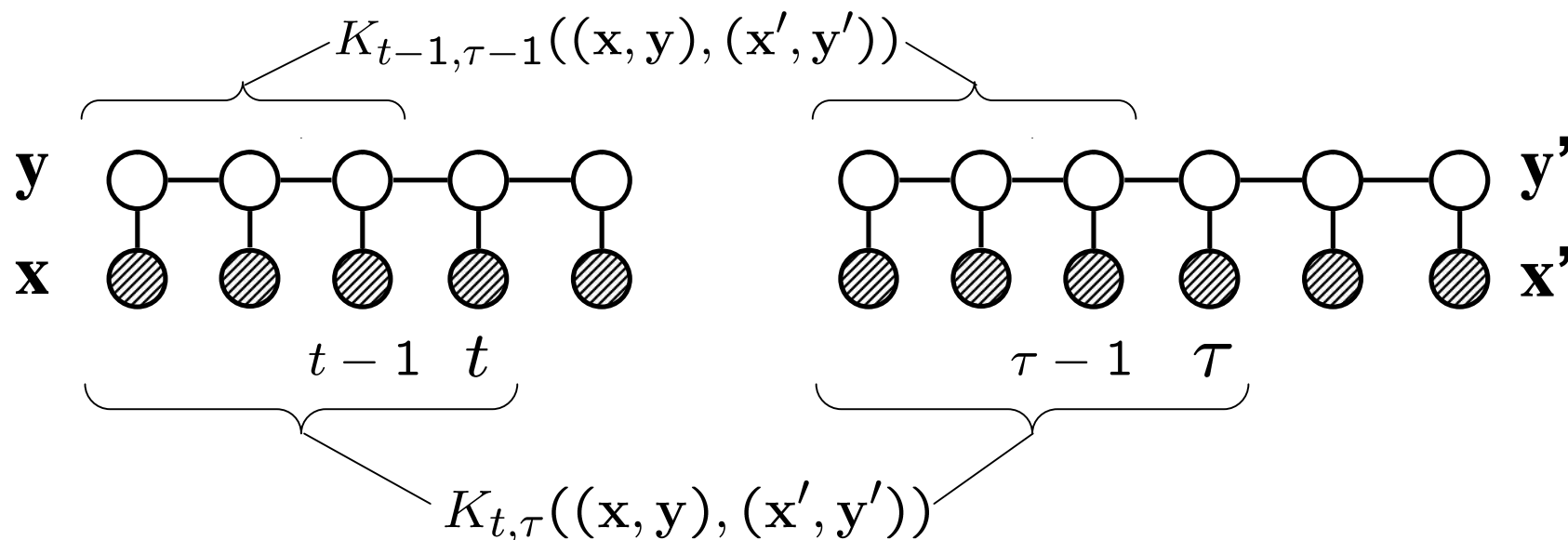
Kernel for Labeled Sequences

- Recursive computation by dynamic programming in $O(|\mathbf{x}||\mathbf{x}'|)$
 - Kernel function is decomposed into kernel functions of pairs of positions

$$K((\mathbf{x}, \mathbf{y}), (\mathbf{x}', \mathbf{y}')) = \sum_{t=1}^{|\mathbf{x}|} \sum_{\tau=1}^{|\mathbf{x}'|} K_{t,\tau}((\mathbf{x}, \mathbf{y}), (\mathbf{x}', \mathbf{y}'))$$

- kernel function of each pair of positions is recursively defined by the kernel function of the previous pair of positions

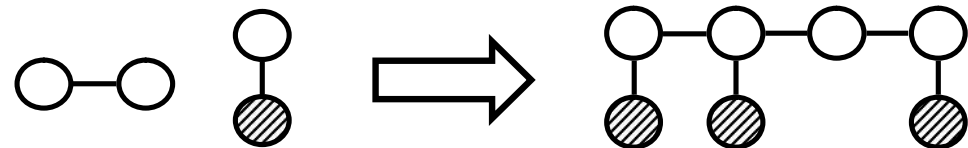
$$K_{t,\tau}((\mathbf{x}, \mathbf{y}), (\mathbf{x}', \mathbf{y}')) = (K_{t-1,\tau-1}((\mathbf{x}, \mathbf{y}), (\mathbf{x}', \mathbf{y}')) + 1) \cdot k(y_t, y'_\tau) \cdot (k(x_t, x'_\tau) + 1)$$



HM-Perceptron with Large Features

- Arbitrarily **large features** to incorporate wider contexts

- idioms, motifs,...



- Two problems

1. Dimension of feature vector $\Phi(\mathbf{x}, \mathbf{y})$ becomes too high
2. Computational complexity of label prediction

$$\hat{\mathbf{y}}(\mathbf{x}) = \operatorname{argmax}_{\mathbf{y} \in \Sigma_y^{|\mathbf{x}|}} \langle \mathbf{w}, \Phi(\mathbf{x}, \mathbf{y}) \rangle$$

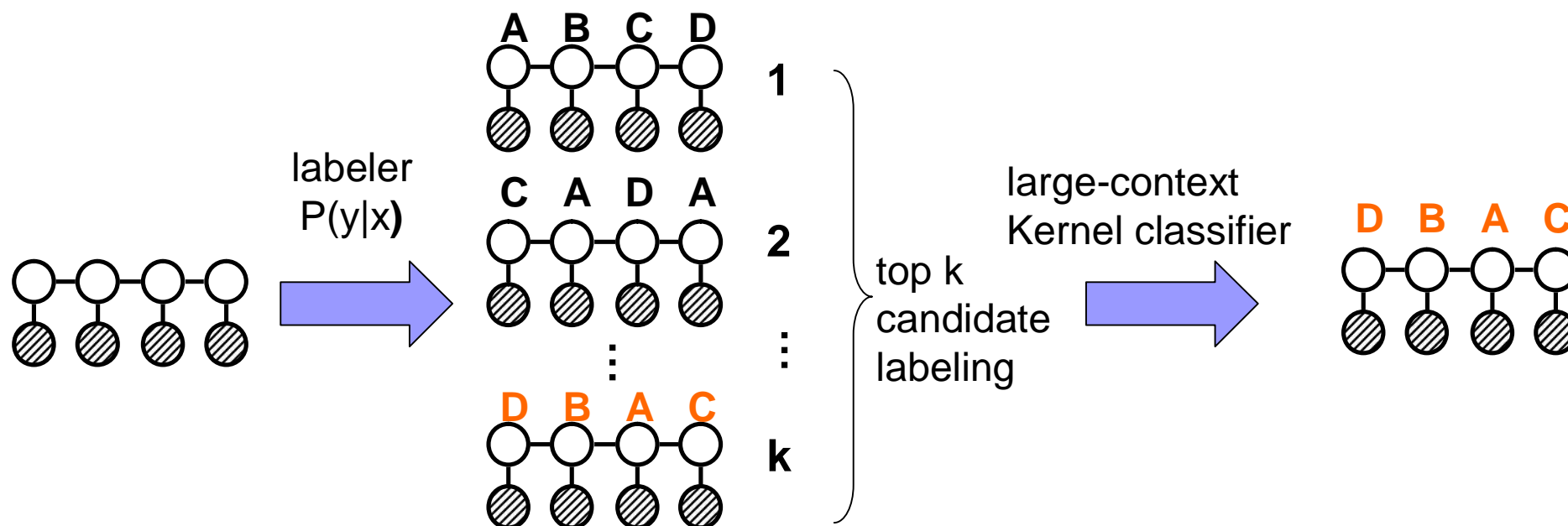
depends **exponentially** on the max number of hidden variables contained in features

- Viterbi algorithm

⇒ HM-perceptron cannot handle arbitrarily large features

Re-ranking [Collins, 2000]

- Avoid to use Viterbi
- Generate top k candidates by standard labeling algorithms (HMM/HM-perceptron)
- Re-rank the candidates by using large features
 - No need to use Viterbi algorithm
- Correct labeling is not guaranteed to be included in the candidates



Point-wise label prediction [Kakade et al. ,2002]

- Avoid to use Viterbi algorithm
- Sequential prediction
 - All labels at once

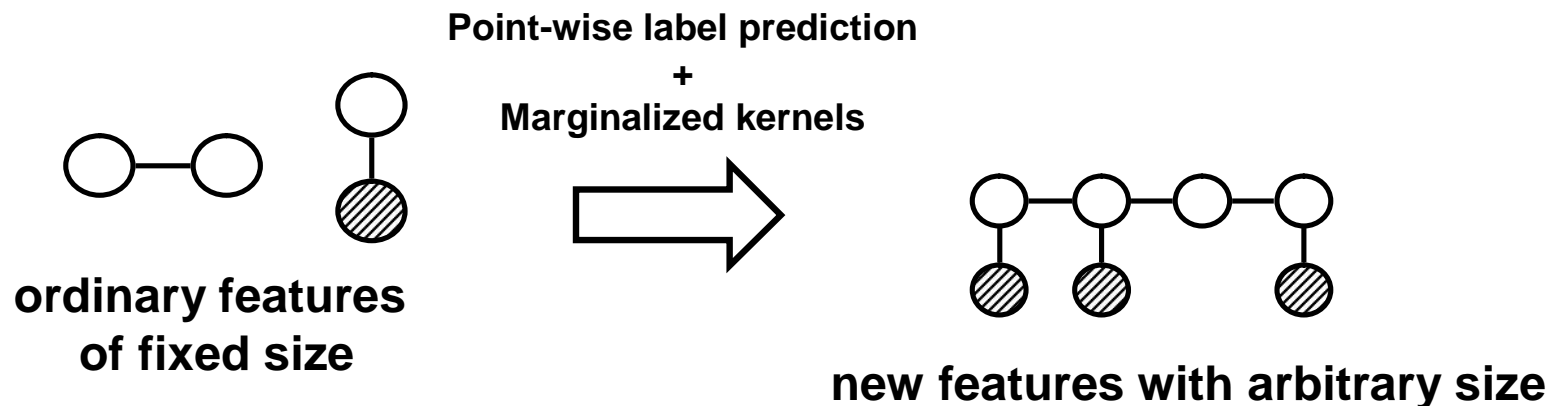
$$\hat{y} = \operatorname{argmax}_y \log \Pr (\begin{array}{c} \text{y} \\ \boxed{y_1 \quad y_2 \quad y_3 \quad y_4} \\ \text{the} \quad \text{man} \quad \text{saw} \quad \text{glasses} \end{array})$$

- Point-wise prediction
 - Each of hidden variables is predicted independently by using marginalized likelihood

$$\hat{y}_1 = \operatorname{argmax}_{y_1} \sum_{y_2, y_3, y_4} \log \Pr (\begin{array}{c} y_1 \quad y_2 \quad y_3 \quad y_4 \\ \boxed{\text{the}} \quad \text{man} \quad \text{saw} \quad \text{glasses} \end{array})$$

Marginalized Labeling Perceptron (Proposed Method)

- Can incorporate arbitrary size features efficiently
 - Dual representation allows to use (marginalized) **kernels**
 - **Point-wise label prediction** [Kakade et al. ,2002]
 - No need to use Viterbi
since each hidden variable is predicted independently



Marginalized Labeling Perceptron (primal)

■ Marginalized Labeling Perceptron

$$\hat{y}_t(\mathbf{x}) = \underset{\tilde{y}_t}{\operatorname{argmax}} \left\langle \mathbf{w}, \sum_{\mathbf{y}: y_t = \tilde{y}_t} P(\mathbf{y}|\mathbf{x}) \Phi(\mathbf{x}, \mathbf{y}) \right\rangle$$

Point-wise prediction
(no Viterbi !)

Marginalized feature vector with fixed label at t
(incorporates all candidates)

Pre-trained labeler with small size features
(e.g. HMM, HM-Perceptron, MEMM, CRF, ...)

Feature vector of
arbitrary size features

■ HM-Perceptron

$$\hat{\mathbf{y}}(\mathbf{x}) = \underset{\mathbf{y}}{\operatorname{argmax}} \langle \mathbf{w}, \Phi(\mathbf{x}, \mathbf{y}) \rangle$$

Find the best hidden sequence
among all possible ways of labeling for \mathbf{y}
(needs Viterbi !)

Learning Marginalized Labeling Perceptron (primal)

■ Marginalized Labeling Perceptron

$$\mathbf{w}^{new} = \mathbf{w}^{old} + \sum_{\mathbf{y}: y_t = y_t^{(i)}} P(\mathbf{y}|\mathbf{x}) \Phi(\mathbf{x}, \mathbf{y}) - \sum_{\mathbf{y}: y_t = \hat{y}_t^{(i)}} P(\mathbf{y}|\mathbf{x}) \Phi(\mathbf{x}, \mathbf{y})$$

$$\Rightarrow \mathbf{w} = \sum_{j=1}^{|Example|} \sum_{\tau=1}^{|\mathbf{x}^{(j)}|} \sum_{\tilde{y}_\tau} \alpha_{j\tau}(\tilde{y}_\tau) \sum_{\mathbf{y}: y_\tau = \tilde{y}_\tau} P(\mathbf{y}|\mathbf{x}) \Phi(\mathbf{x}^{(j)}, \mathbf{y})$$

■ HM-Perceptron

$$\mathbf{w}^{new} = \mathbf{w}^{old} + \Phi(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}) - \Phi(\mathbf{x}^{(i)}, \hat{\mathbf{y}}^{(i)})$$

$$\Rightarrow \mathbf{w} = \sum_{j=1}^{|Example|} \sum_{\tilde{\mathbf{y}}} \alpha_j(\tilde{\mathbf{y}}) \Phi(\mathbf{x}^{(j)}, \tilde{\mathbf{y}})$$

Marginalized Labeling Perceptron (dual)

■ Dual representation

- Kernel methods to handle arbitrary size features efficiently

■ Prediction

$$\hat{y}_t(\mathbf{x}) = \operatorname{argmax}_{\tilde{y}_t} \sum_{j=1}^{|Examples|} \sum_{\tau=1}^{|\mathbf{x}^{(j)}|} \sum_{\tilde{y}_\tau} \alpha_{j\tau}(\tilde{y}_\tau) K(\mathbf{x}^{(j)}, \mathbf{x}, \tau, t, \tilde{y}_\tau, \tilde{y}_t)$$

- Marginalized inner products = Marginalized kernel

$$K(\mathbf{x}, \mathbf{x}', t, \tau, \tilde{y}_t, \tilde{y}'_\tau) = \sum_{\mathbf{y}: y_t = \tilde{y}_t} \sum_{\mathbf{y}': y'_\tau = \tilde{y}'_\tau} P(\mathbf{y}|\mathbf{x}) P(\mathbf{y}'|\mathbf{x}') \langle \Phi(\mathbf{x}, \mathbf{y}), \Phi(\mathbf{x}', \mathbf{y}') \rangle$$

Positions whose labels are fixed

fixed labels at t & τ

Marginalize over all possible ways of labeling of \mathbf{y} & \mathbf{y}' with fixed labels at t & τ

Inner product between two labeled sequences

■ Update

- $\alpha_{it}^{new}(y_t^{(i)}) = \alpha_{it}^{old}(y_t^{(i)}) + 1$
- $\alpha_{it}^{new}(\hat{y}^{(i)}) = \alpha_{it}^{old}(\hat{y}^{(i)}) - 1$

Computation of Marginalized Kernels

- Given two sequences of observed variables \mathbf{x} & \mathbf{x}'
 - Compute marginalized kernels for all pairs of positions t & τ at once
 - Kernel decomposition for efficient computation

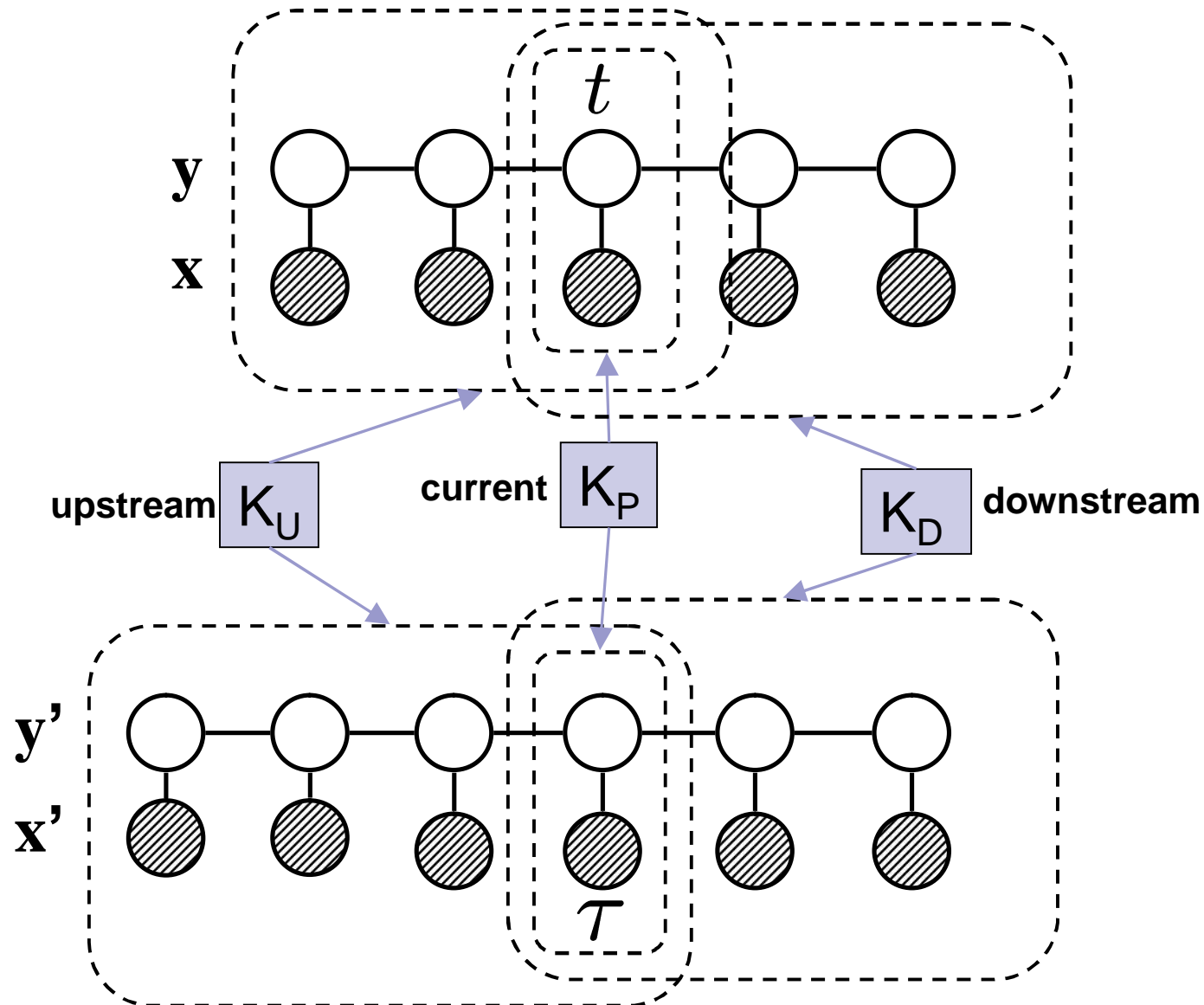
$$\begin{aligned} K(\mathbf{x}, \mathbf{x}', t, \tau, \tilde{y}_t, \tilde{y}'_\tau) &= \sum_{\mathbf{y}: y_t = \tilde{y}_t} \sum_{\mathbf{y}': y'_\tau = \tilde{y}'_\tau} P(\mathbf{y}|\mathbf{x}) P(\mathbf{y}'|\mathbf{x}') \langle \Phi(\mathbf{x}, \mathbf{y}), \Phi(\mathbf{x}', \mathbf{y}') \rangle \\ &= \underbrace{K_U(\mathbf{x}, \mathbf{x}', t, \tau)}_{\text{up-stream kernel}} \cdot \underbrace{K_P(\mathbf{x}, \mathbf{x}', t, \tau, \tilde{y}_t, \tilde{y}'_\tau)}_{\text{current position kernel}} \cdot \underbrace{K_D(\mathbf{x}, \mathbf{x}', t, \tau)}_{\text{down-stream kernel}} \end{aligned}$$

Each kernel is recursively computed by dynamic programming

- Computational Complexity $O(|\mathbf{x}| |\mathbf{x}'|)$

Kernel Decomposition (Sequences)

- Analogous to Forward-Backward algorithm for HMM



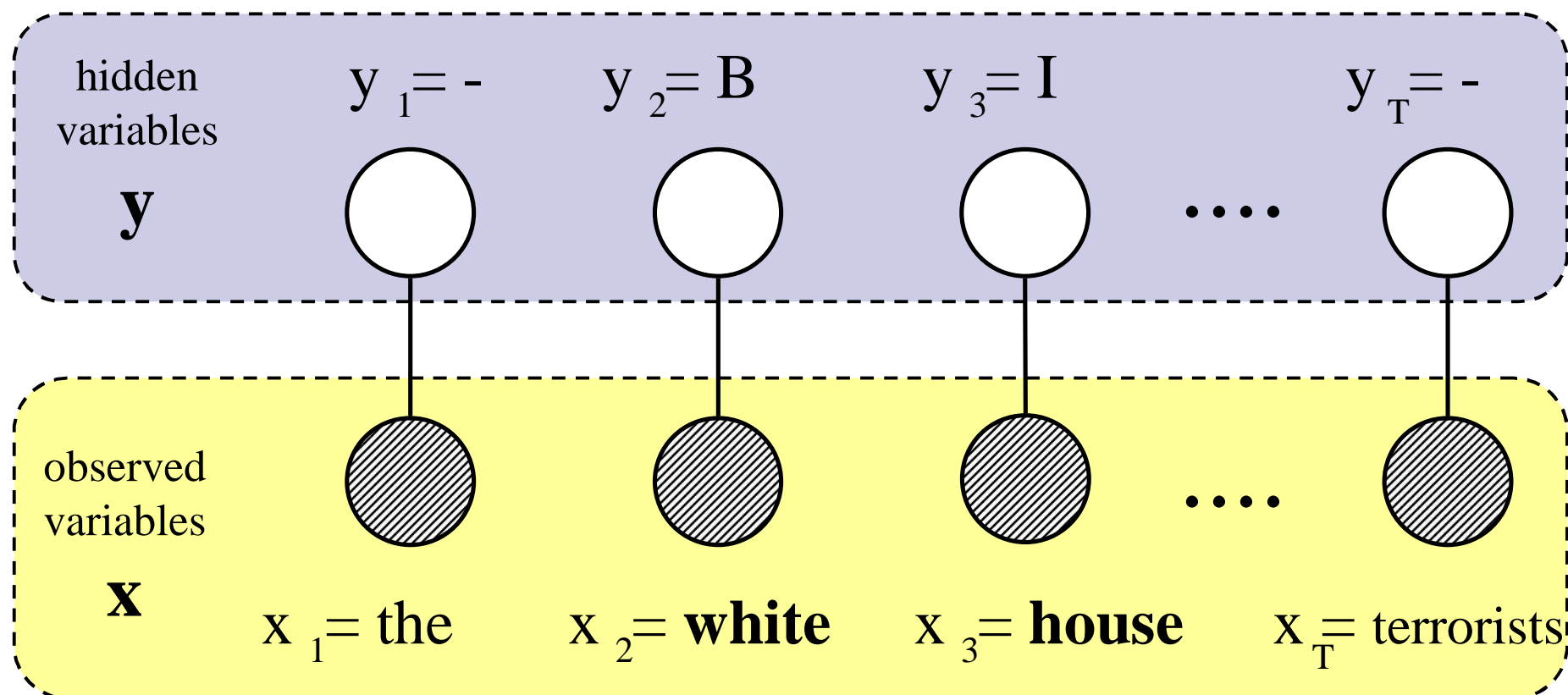


Experiments: Named Entity Recognition

- NLP version of secondary structure prediction
- CoNLL2002 data
 - 300 sentences (8,500 words)
 - 9 labels to indicate *person name, organization name, place, ...*
- Sales log data
 - 184 sentences (3,500 words)
 - a sales log written by sales representatives (written in Japanese)
 - “They ported their server to a **linux cluster** for its **availability**”
 - 12 labels to indicate *product name, company name, reason for buying the product,...*
- Comparison of two methods
 - HM-Perceptron
 - window size = 3
 - Marginalized labeling perceptron with sequence kernel

Named Entity Recognition

- Given \mathbf{x} = sequence of words
 - “The white house was attacked by the terrorists”
- predict \mathbf{y} = named entities
 - beginning (B), inside (I), other (-)



Results (5-fold cross validation)

- Uniform distributions as the priors $P(\mathbf{y}|\mathbf{x})$ for marginalized kernels
- CoNLL2002 data

	PRECISION	RECALL	F1
HM-PERCEPTRON	23.8% (14.6)	17.9% (3.0)	18.6 (5.2)
SEQUENCE KERNEL	49.0% (6.0)	23.1% (8.1)	30.5 (6.7)

(standard deviation)

- Sales log

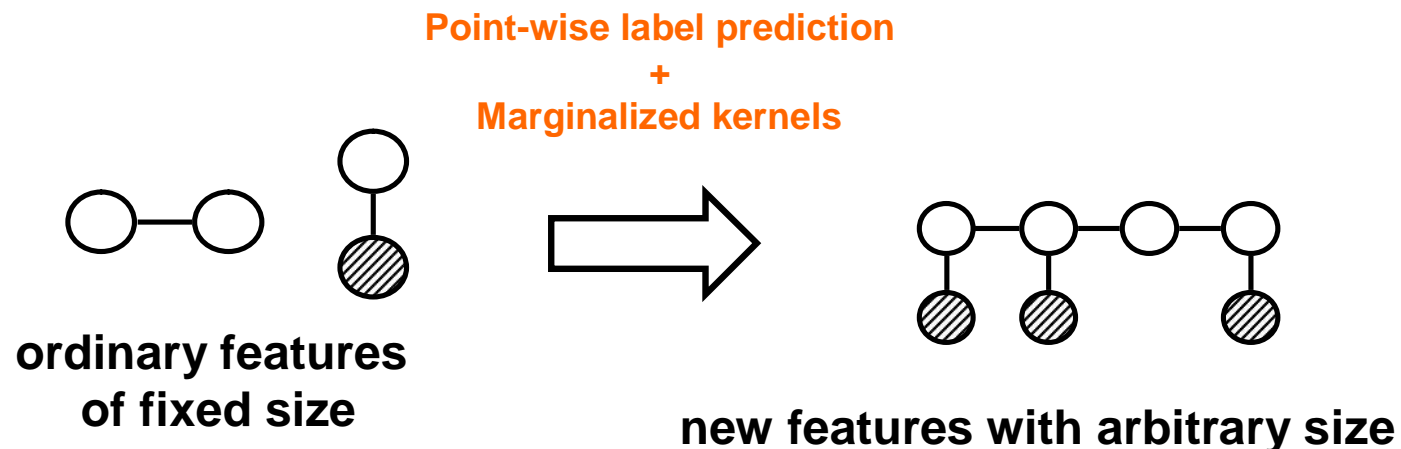
	PRECISION	RECALL	F1
HM-PERCEPTRON	51.5%(8.5)	24.0%(21.4)	28.9(20.3)
SEQUENCE KERNEL	52.2% (9.5)	29.6% (4.0)	37.5 (4.1)

(standard deviation)

- Precision = $\#(\text{correctly predicted NE labels}) / \#(\text{predicted NE labels})$
- Recall = $\#(\text{correctly predicted NE labels}) / \#(\text{true NE labels})$
- F1 = $2 \cdot \text{Specificity} \cdot \text{Sensitivity} / (\text{Specificity} + \text{Sensitivity})$

Summary

- Marginalized labeling perceptron
 - Labeling learning algorithm for sequence labeling
 - Can handle features with arbitrarily many hidden variables by using
 - Point-wise label prediction avoids Viterbi
 - Marginalized kernels
avoids explicit handling of high dimensional feature vectors





Application in Bioinformatics

- Protein secondary structure prediction
 - Predict alpha-helix/beta-strand regions in protein sequences
- Gene finding
 - Predict gene regions in DNA sequences
- Phosphorylation site prediction
 - Predict phosphorylated sites in protein sequences
- Information extraction from biomedical texts
 - Protein name recognition

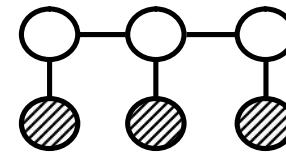
Marginalized Kernels for Labeling Structured Data

$$\blacksquare K(\mathbf{x}, \mathbf{x}', t, \tau, \tilde{y}_t, \tilde{y}'_\tau) = \sum_{\mathbf{y}: y_t = \tilde{y}_t} \sum_{\mathbf{y}': y'_\tau = \tilde{y}'_\tau} P(\mathbf{y}|\mathbf{x}) P(\mathbf{y}'|\mathbf{x}') \langle \Phi(\mathbf{x}, \mathbf{y}), \Phi(\mathbf{x}', \mathbf{y}') \rangle$$

■ 3 kernels

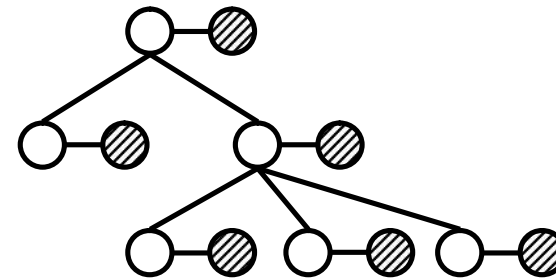
□ Sequence labeling

■ Sequence features



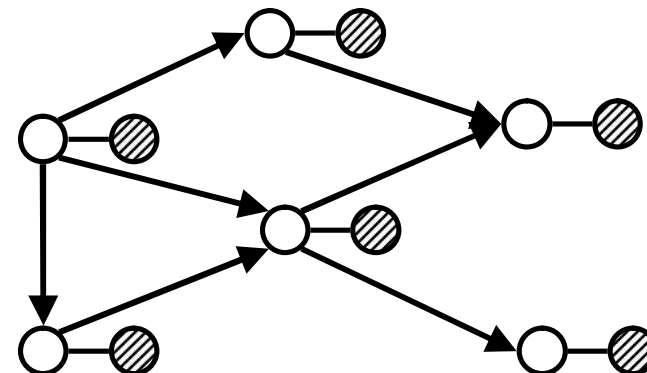
□ Ordered tree labeling

■ Tree features



□ Directed acyclic graph (DAG) labeling

■ Path features





Outline

- Sequence labeling problems
- Hidden Markov model (HMM)
- Hidden Markov (HM) perceptron
- Marginalized labeling perceptron
- Experiments on natural language processing tasks
- Application to bioinformatics

Thank you



Conditional Models

■ conditional model

$$\Pr(y|\mathbf{x}) = \frac{\exp \langle \mathbf{w}, \Phi(\mathbf{x}, y) \rangle}{\sum_{\tilde{y}} \exp \langle \mathbf{w}, \Phi(\mathbf{x}, \tilde{y}) \rangle}$$

■ MLE of conditional models

□ (sequential objective function)

$$\hat{\mathbf{w}} = \operatorname{argmax}_{\mathbf{w}} \sum_i \Pr(\mathbf{y}^{(i)} | \mathbf{x}^{(i)})$$

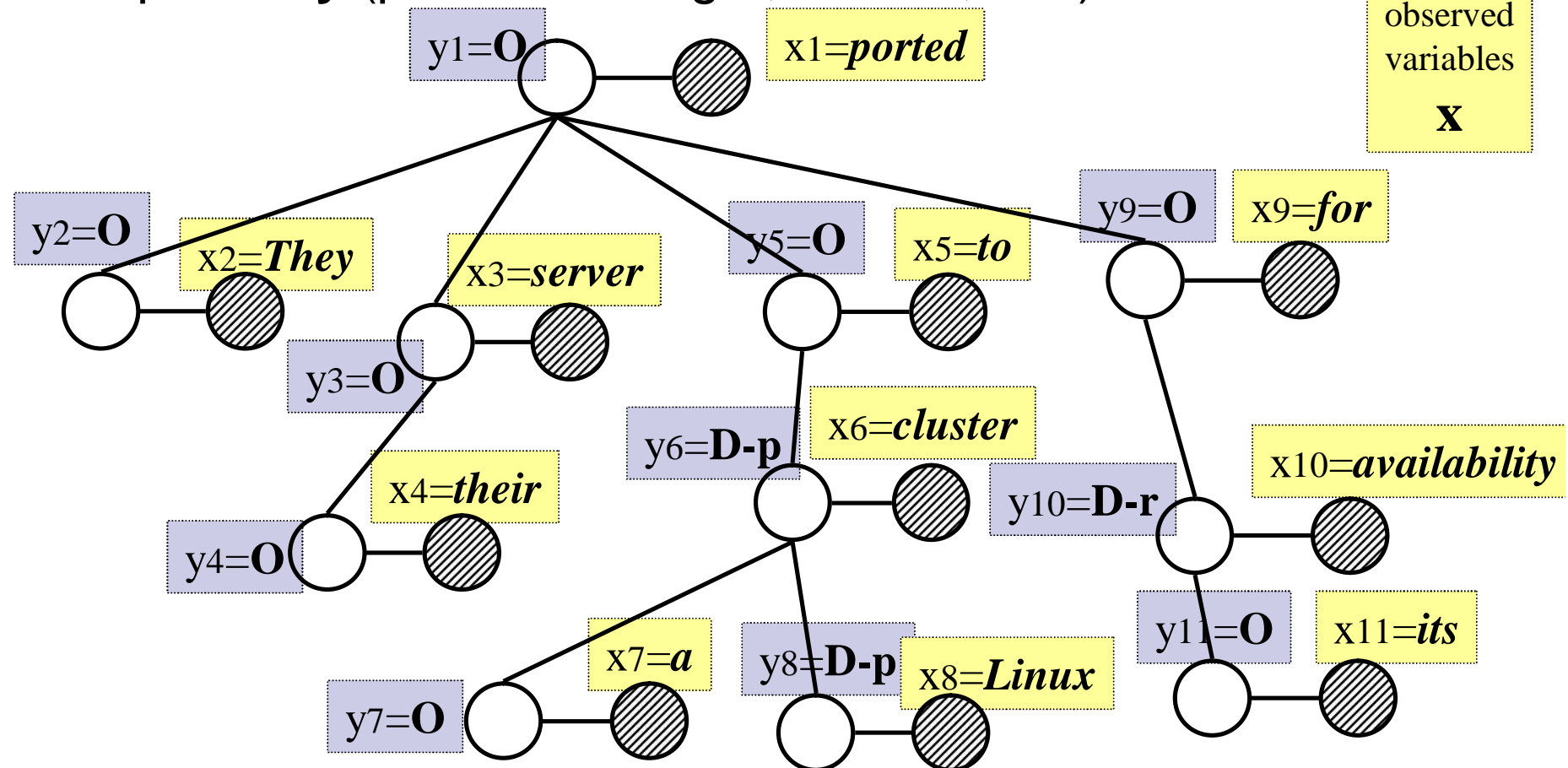
□ (point-wise objective function)

$$\hat{\mathbf{w}} = \operatorname{argmax}_{\mathbf{w}} \sum_i \sum_{t=1}^{T(i)} \sum_{\mathbf{y}: y_t = y_t^{(i)}} \Pr(\mathbf{y} | \mathbf{x}^{(i)})$$

Tree Labeling Problem

■ Information Extraction from parse trees

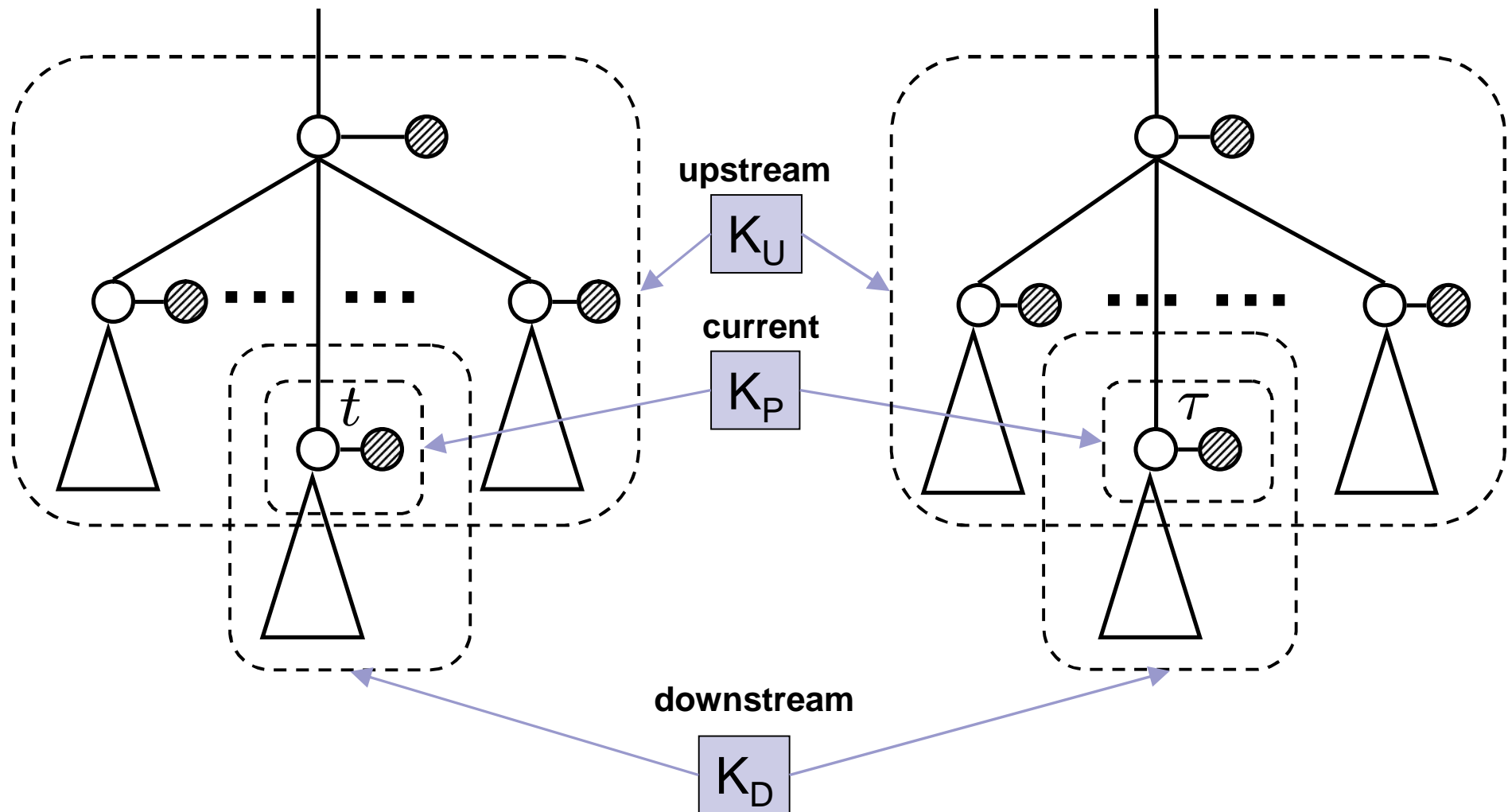
- Given \mathbf{x} (parse tree for text),
predict \mathbf{y} (products bought, reason, etc.)



*“They ported their server to a **linux cluster** for its **availability**”*

Kernel Decomposition (Ordered Trees)

- Analogous to Inside-Outside algorithm for probabilistic CFG



Kernel Decomposition (DAG)

