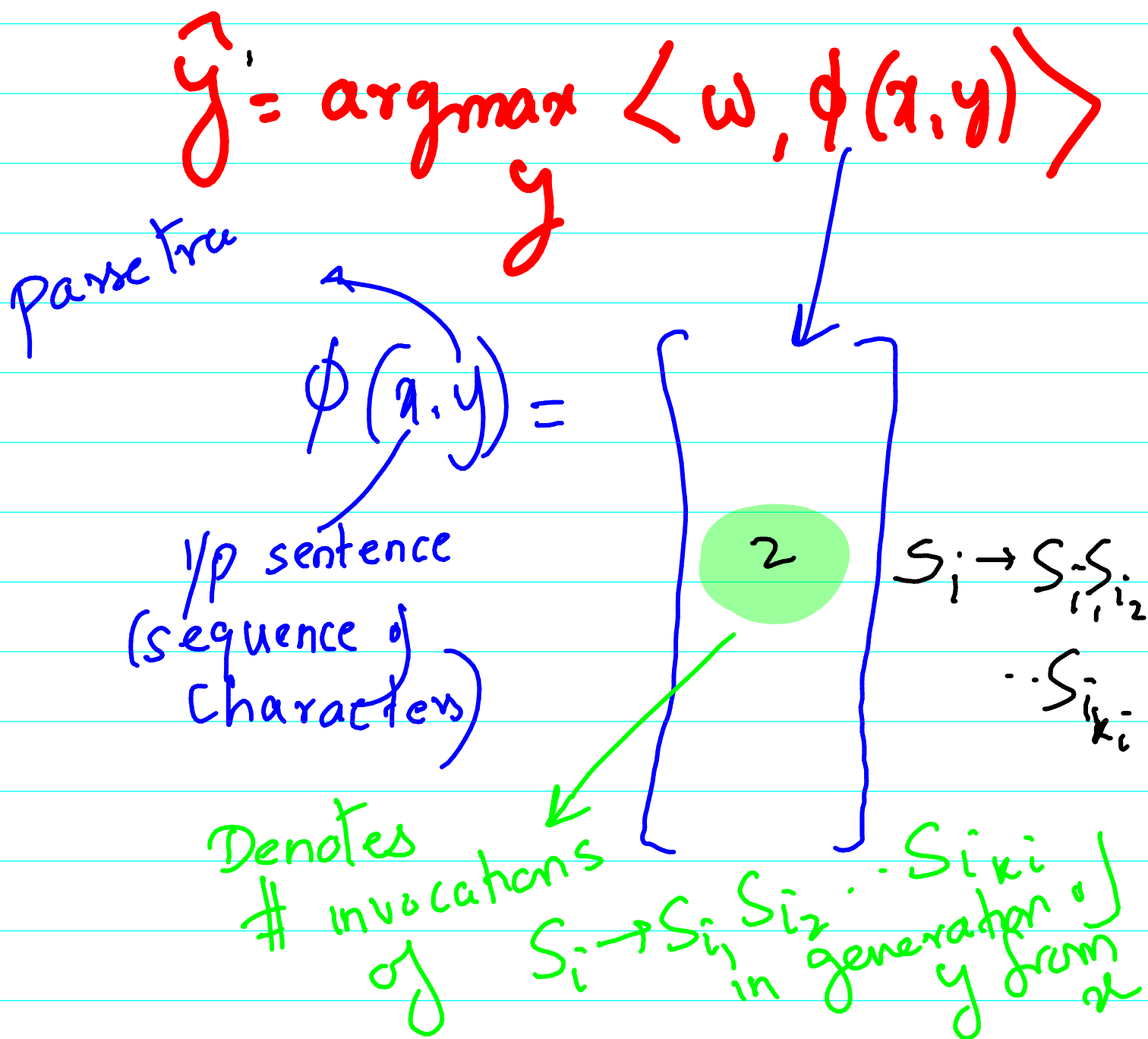


Q: How to compute (efficiently)



Assuming  $w$  is learnt (fixed)  
how to compute  $\hat{y}$  efficiently?  
When  $w = [1 \dots 1]$ , CYK algo.

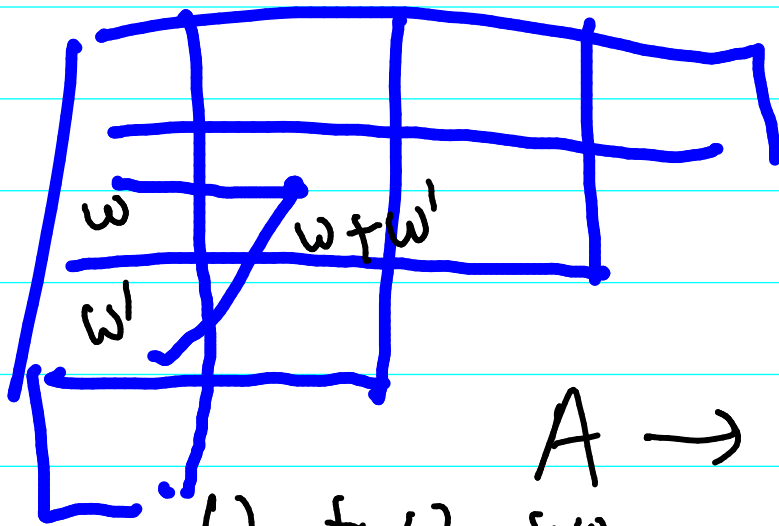
comes to rescue.

(Refer to slides 16 to 27  
of CFL and Parsing.ppt)

How would you modify CYK  
to compute

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

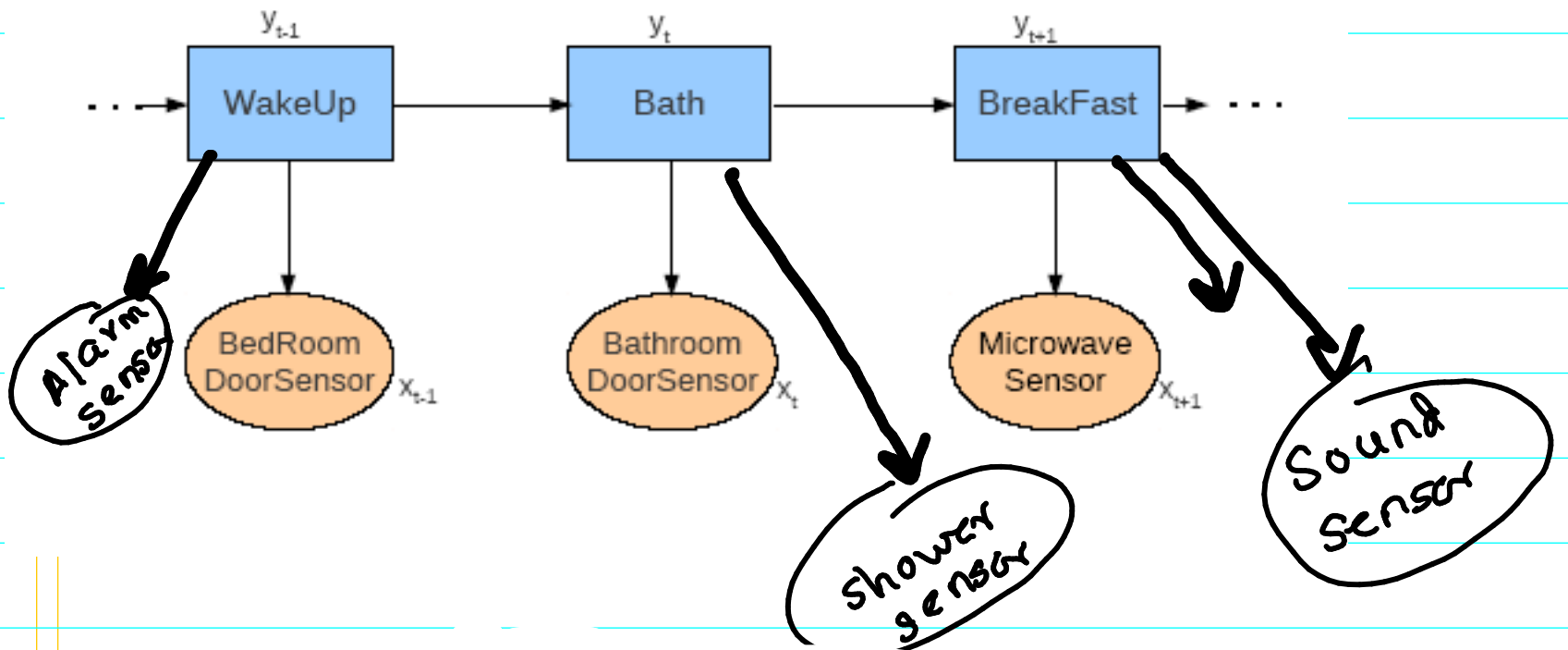
<http://www.cse.iitb.ac.in/~cs717/notes/classNotes/StructuredOutput/CFLandParsing.ppt>



$A \rightarrow BC$   
 $w_A + w_B + w_C$        $w_B$     $w_C$

$A^{(1)} \rightarrow BC$   
 $A^{(2)} \rightarrow EF$  }  $\max \left( \begin{array}{l} w_{A^{(1)}} + w_B + w_C \\ w_A + w_E + w_F \end{array} \right)$

Next we discuss STRUCTSUM  
when  $y$  is a sequence



Can be treated as a special case of parsing to be predicted

①

$$\phi(x, y) =$$

Size = #labels  
+ #labels \* #features  
+ #features

$$\begin{bmatrix} 0 \\ 1 \\ \vdots \\ 2 \\ \vdots \\ 1 \end{bmatrix}$$

- $S \rightarrow S B_f$
- $S \rightarrow S B_a$
- $S \rightarrow S W_a$
- $W_a \rightarrow$  Bedroom sensor
- $W_a \rightarrow$  Bathroom sensor
- $W_a \rightarrow$  Microwave sensor

Will also need

$S \Rightarrow$  every possible feature

In practice, a different representation  $\phi$  is used

To be predicted

2

$\phi(x, y)$   
#labels  
#(labels - 1)  
#labels  
#features

$W_a \rightarrow B_a$   
 $B_a \rightarrow F_a$   
 $F_a \rightarrow W_a$

$W_a \rightarrow$  Bedroom sensor  
 $W_a \rightarrow$  Bathroom sensor  
 $W_a \rightarrow$  Microwave sensor

observations

Since Regular grammar  $\subseteq$  CFG,  
it is possible that with a richer  
( $\&$  therefore larger)  $\phi$ , we could

still compute  $\hat{y} = \arg \max_y \langle \omega, \phi(x, y) \rangle$   
in reasonable time

Use this representation  
for sequence  
labeling

HMMs (Hidden Markov Models)

CRFs (Conditional Random fields)

Struct SVM HMM: Margin

Objective function  
maximized to obtain  $\omega$

Likelihood

Pseudo  
log likelihood  
-0.4

$$\min_{f, \xi} \frac{1}{2} \|f\|^2 + \frac{C}{m} \sum_{i=1}^m \xi_i, \quad \text{s.t. } \forall i: \xi_i \geq 0,$$

$$\forall i, \forall Y \neq Y_i: \langle f, \psi_i^\delta(Y) \rangle \geq 1 - \frac{\xi_i}{\Delta(Y_i, Y)}$$

Alternatives:  $1 - \xi_i + \Delta(Y_i, Y)$   
(discussed in last lecture)

$$\langle f, \psi_i^\delta(Y) \rangle = \langle f, \psi(X_i, Y_i) \rangle - \langle f, \psi(X_i, Y) \rangle$$

**Input:** kernels,  $C$ ,  $\epsilon_{margin}$

1.  $S_i \leftarrow \phi \quad \forall i = 1, \dots, m$

2. **repeat**

3.     **for**  $i = 1, \dots, m$  **do**

4.         Define  $H(Y) \equiv [1 - \langle \mathbf{f}, \psi_i^\delta(Y) \rangle] \Delta(Y_i, Y)$

5.         Compute  $\hat{Y} = \arg \max_{Y \in \mathcal{Y}} H(Y)$ .

6.         Compute  $\xi_i = \max\{0, \max_{Y \in S_i} H(Y)\}$ .

7.         **if**  $H(\hat{Y}) > \xi_i + \epsilon_{margin}$ , **then**

8.              $S_i \leftarrow S_i \cup \{\hat{Y}\}$ .

9.              $\alpha \leftarrow$  optimize dual over  $S$ ,  $S = \bigcup_i S_i$ .

10.         **end if**

11.     **end for**

12. **until** no  $S_i$  has changed during the iteration.

//for each example

//Margin Violation

//Max Margin Violation

//Current Max Margin Violation

//adding constraints

// $\mathbf{f}$  can be derived from  $\alpha$

For parse trees  
you use  
CYK  
(modified)

Typically  
for complex models  
inference is also  
a part of learning

For HMM/CRF/StructSVM,  
we discuss Viterbi to  
compute argmax

The same

$\arg \max_y \langle \mathbf{f}, \phi(x, y) \rangle$

will need to be solved for  
inference at runtime as well.

# The Viterbi algo (analog of CYK)

Let us write more expressively

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

$$= \underset{y}{\operatorname{argmax}} \left[ \sum_{i=1}^n w_{y_i x_i} \phi_{y_i x_i}(x, y) \right]$$

node wts  
node features

$$y = [y_1, y_2, \dots, y_n]$$

$$x = [x_1, x_2, \dots, x_n]$$

$$+ \left[ \sum_{i=1}^n w_{y_i y_{i+1}} \phi_{y_i y_{i+1}}(x, y) \right]$$

edge features / wts

Apply this trick:

$$\max_{a, b} g(a) f(a, b) = \max_a g(a) \max_b f(a, b)$$

$\underbrace{\quad}_{[y_1 \dots y_{n-1}]}$       $y_n$

$u(a)$

$$= \max_a g(a) u(a)$$

$$\max_{y_1, \dots, y_n} \left[ \sum_{i=1}^n w_{y_i, x_i} \phi_{y_i, x_i}(x, y) + \sum_{i=1}^n w_{y_{i+1}, y_i} \phi_{y_{i+1}, y_i}(x, y) \right]$$

$$= \max_{y_n} w_{y_n, x_n} \phi_{y_n, x_n}(x, y) + w_{y_n, y_{n-1}} \phi_{y_n, y_{n-1}}(x, y) + \max_{y_1, \dots, y_{n-1}} \left\{ \sum_{i=1}^{n-1} w_{y_i, x_i} \phi_{y_i, x_i}(x, y) + \sum_{i=1}^{n-1} w_{y_{i+1}, y_i} \phi_{y_{i+1}, y_i}(x, y) \right\}$$

$M_{n-1}(y_n)$

$M_n$



$$U_k(y_{k+1}) = \max_{y_k} \left\{ \begin{aligned} & \omega_{x_{k+1}, y_{k+1}} \phi_{x_{k+1}, y_{k+1}}(x, y) \\ & + \omega_{y_k, y_{k+1}} \phi_{y_k, y_{k+1}}(x, y) \\ & + U_{k-1}(y_k) \end{aligned} \right\} (*)$$

Algo: (Viterbi)

$$U_0(y_1) = \omega_{x_1, y_1} \phi_{x_1, y_1}(x, y)$$

for  $k=1 \dots n$  keep applying (\*)

Some additional bookkeeping  
to find  $\text{argmax}$  (H/w)

H/w: Compare run time of

modified CYK & viterbi