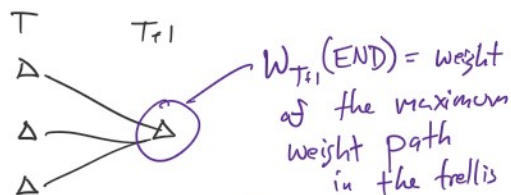
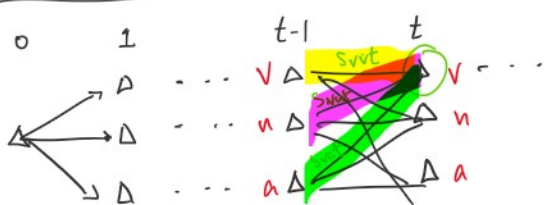


# Section A

Wednesday, March 29, 2023 9:31 AM

## Viterbi Algo (edge weights)



$$W_t(v) = \max (W_{t-1}(v) s_{vvt},$$

$$W_{t-1}(n) s_{vnt},$$

$$W_{t-1}(a) s_{vta})$$

for  $t = 1 \dots T$ :

for  $k = 1 \dots K$ :

$$W_t(k) = \max_{j \in \{1, \dots, K\}} W_{t-1}(j) s_{kjt}$$

$$W_{t-1}(j) s_{kjt}$$

For HMM:  
 $s_{kjt} = p(y_t = k | y_{t-1} = j)$   
 $p(x_t | y_t = k)$

$$b_t(k) = \operatorname{argmax}_{j \in \{1, \dots, K\}} W_{t-1}(j) s_{kjt}$$

$$W_{t-1}(j) s_{kjt}$$

backpointer

Ex: CMU going to the moon?

$N=1 \Rightarrow$  Bloomberg article

$M=1 \Rightarrow$  CMU is going to the moon

$H=1 \Rightarrow$  Elaborate hoax

$A=1 \Rightarrow$  April 1st (or close)

$B=1 \Rightarrow$  CMU researchers are bored

$G=1 \Rightarrow$  building control center in GHC

$C=1 \Rightarrow$  actually just a carnival booth

$S=1 \Rightarrow$  CMU students are amazing

Poll Q1

Q: How to represent  $p(N, M, H, A, B, G, C, S)$ ?

Idea #1: Condition on the observables

$$P(M, H, B, G, C, S \mid N=1, A=1)$$

Idea #2: Chain Rule

$$P(N, M, \dots, S) = P(N | M \dots S) P(M | A \dots S) \dots P(S)$$

pro: looks compact

pro: looks compact  
 con: actually not compact

Idea #3: Full Joint Table

N	M	H	A	B	G	C	S	$P(-)$
0	0	0	0	0	0	0	0	$\theta_1$
0	0	0	0	0	0	0	1	$\theta_2$
0	0	0	0	0	0	1	0	$\theta_3$
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
1	1	1	1	1	1	1	1	$\theta_M$

pro: can encode any distribution  
 con: enormous

Idea #4: Full Independence

$$P(N, M, \dots, S) = P(N) P(M) \dots P(S)$$

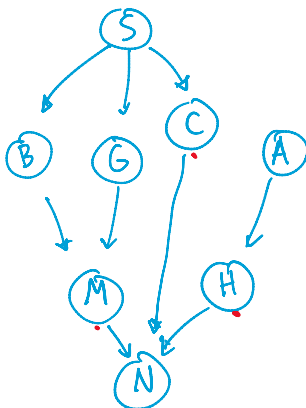
pro: compact  
 con: completely unrealistic

Idea #5: Naive Bayes

$$P(N, M, \dots, S) = P(M) P(N|M) \dots P(S|M)$$

pro: still compact  
 con: missing a lot of dependencies

Idea #6: Bayes net Network



- Write out the chain rule
- Remove same RHS variables

← w/ carefully chosen order

$$\begin{aligned}
 P(N, M, H, C, B, A, G, S) &= P(N | M, H, C, \cancel{B}, \cancel{A}, \cancel{G}, S) \\
 &P(M | \cancel{N}, \cancel{H}, \cancel{C}, \cancel{B}, \cancel{A}, \cancel{G}, S) \\
 &P(H | \cancel{N}, \cancel{M}, \cancel{C}, \cancel{B}, \cancel{A}, \cancel{G}, S) \\
 &P(C | \cancel{N}, \cancel{M}, \cancel{H}, \cancel{B}, \cancel{A}, \cancel{G}, S) \\
 &P(B | \cancel{N}, \cancel{M}, \cancel{H}, \cancel{C}, \cancel{A}, \cancel{G}, S) \\
 &P(A | \cancel{N}, \cancel{M}, \cancel{H}, \cancel{C}, \cancel{B}, \cancel{G}, S) \\
 &P(G | S)
 \end{aligned}$$

**Proof Cond. Indep**

$P(A | X, Z)$   
 $P(G | S)$   
 $P(S)$

Case #1: Cascade



$$\begin{aligned}
 p(x, z | y) &= \frac{p(x, y, z)}{p(y)} \\
 &= \frac{p(x|y) p(y|z) p(z)}{p(y)} \quad \text{BN prob} \\
 &= \frac{p(x|y) p(z|y) p(y)}{p(y)} \\
 &= p(x|y) p(z|y) \\
 &\Rightarrow X \perp\!\!\!\perp Z | Y
 \end{aligned}$$

Case #2:  $\rangle$  left as exercises  
 Case #3:

**OH**

error( $h_{MAP}$ ) > error( $h_{MLE}$ )  $\lambda \sum \frac{\theta_n^2}{n}$   
 $\Rightarrow$  overfit w/ MAP  $\leftarrow L_2: r(\theta) = \lambda (\|\theta\|_2)^2$

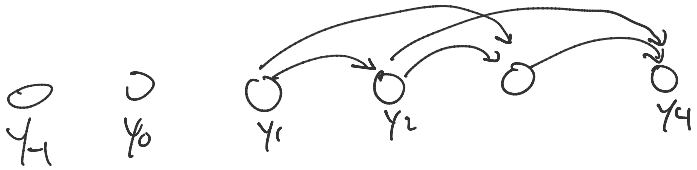
$J'(\theta) = J(\theta) + r(\theta)$   
 $\uparrow$  fit data  $\quad \uparrow$  keep simple

~~$p(\theta) = \text{Gauss}(\mu=0, \sigma^2)$~~   
 $p(\theta) = \prod_{n=1}^M p(\theta_n)$

$J(\theta) = -\log p(D|\theta)$   
 $r(\theta) = -\log p(\theta)$   
 $= -\log \prod_{n=1}^M \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(\theta_n)^2}{2\sigma^2}\right)$   
 $= M \sum_{n=1}^M \left[ \log(\sqrt{2\pi\sigma^2}) + \frac{\theta_n^2}{2\sigma^2} \right]$

$p(\theta_n) \sim \text{Gauss}(\mu=0, \sigma^2)$

$$\begin{aligned} & \text{argmin } J(\theta) \\ &= \text{argmin } J(\theta) + \sum_{n=1}^M \frac{\theta_n^2}{2\sigma^2} \quad \tau = \frac{1}{2\sigma^2} \\ &= \text{argmin } J(\theta) \end{aligned}$$



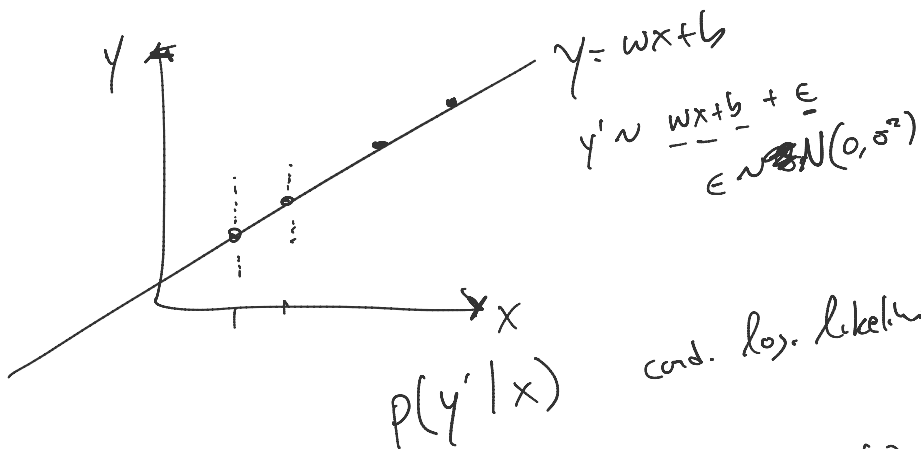
$$p(y_4 | y_3, y_2) p(y_3 | y_2, y_1) p(y_2 | y_1) p(y_1)$$

$x_1$ weather	$x_2$ in. of rain	$y$ Season
rainy	0.6	S
sunny	0.01	S
cloudy	0.05	F
rainy	0.01	F

real-valued  
continuous

$x_1$	$x_1''$	$x_1'''$	$x_2$	$y$
1	0	0	0.6	S
0	1	0	0.01	S
0	0	1	0.05	F
1	0	0	0.01	F



cond. log. likelihood

$$\begin{aligned} x_0 &= \phi_0(x) = x^0 \\ x_1 &= \phi_1(x) = x^1 \\ x_2 &= \phi_2(x) = x^2 \end{aligned}$$

1 - 1





$$x_2 \phi_2(x) = x$$

$$p(y=1|\vec{x}) = \sigma(\underline{w}_1 x_1 + \underline{w}_0 x_0) \quad d=1$$

$$p(y=1|\vec{x}) = \sigma(\underline{w}_1 x_1 + \underline{w}_2 x_2 + \underline{w}_0 x_0) \quad d=2$$

$$\vec{w}' = \underset{\vec{w}}{\text{argmin}} J(\vec{w}) \quad d=1$$

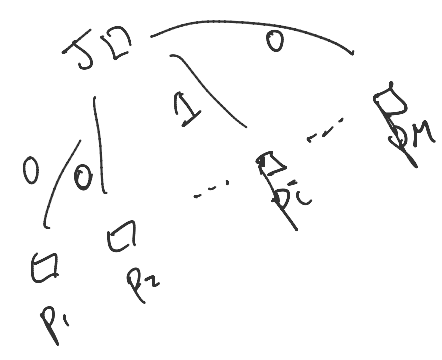
$$\vec{w}'' = \underset{\vec{w}}{\text{argmin}} J_{d=2}(\vec{w})$$

$$\vec{w}'' = \begin{bmatrix} w'_0 & w'_1 & 0 \\ w_0 & w_1 & w_2 \end{bmatrix}$$

$$\vec{p} = f(\vec{h})$$

$$\frac{\partial J}{\partial \vec{p}}$$

$\swarrow$  matrix      $\searrow$  vector  
 $\frac{\partial J}{\partial \vec{h}} = \left( \frac{\partial \vec{p}}{\partial \vec{h}} \right) \left( \frac{\partial J}{\partial \vec{p}} \right)$



$$\frac{\partial J}{\partial p_j} = 0 \quad \forall j \neq i$$

$$\frac{\partial J}{\partial p_j} = 1 \quad \text{for } j=i$$