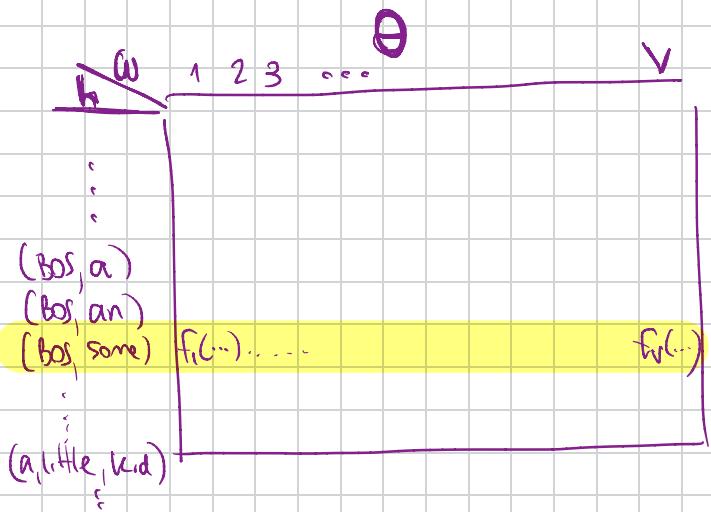


① Categorical Distribution

$$W | H=h \sim \text{Categorical}\left(\underbrace{\mathbf{f}(h; \Theta)}_{\left(f_1(h; \Theta), f_2(h; \Theta), \dots, f_V(h; \Theta)\right)}\right)$$

Tabular Representation



$f \rightarrow$ holds lockup operation

$$f(\underbrace{BOS, \text{some}}_h; \Theta) = \Theta_1^{(h)} / \Theta_2^{(h)} / \Theta_3^{(h)} \dots / \Theta_v^{(h)}$$

② Logistic Representation Categorical CPDs

$$W | H=h \sim \text{Categorical} \left(f(h; \theta) \right)$$

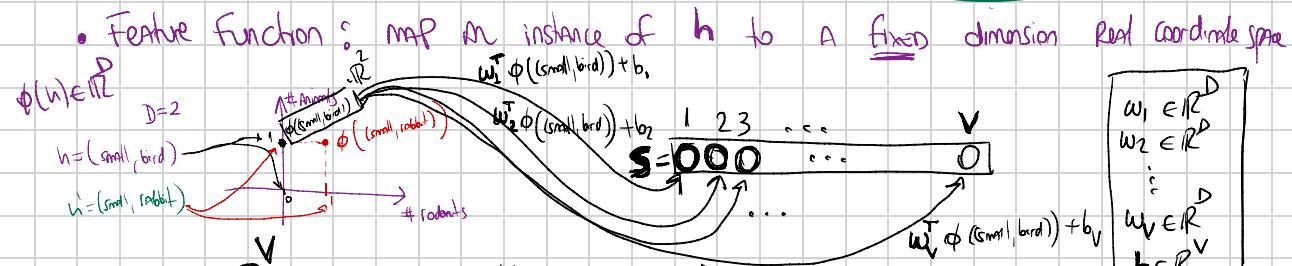
$\underbrace{\quad\quad\quad}_{\left(f_1(h; \theta), f_2(h; \theta), \dots, f_V(h; \theta) \right)}$

\downarrow

$\{1, \dots, V\}$

Instead of storing the cond. prob. masses for every (h, w) condition-outcome pair,

we predict the cond. prob. masses using a log-linear model



$s \in \mathbb{R} \rightarrow \text{probability simplex } \Delta^{V-1}$

$$[\text{softmax}(s)]_i = \frac{\exp(s_i)}{\sum_{j=1}^V \exp(s_j)}$$

1. The model size does not depend on how many different instances of h exist

$$\begin{array}{c} w_1 \in \mathbb{R}^D \dots w_v \in \mathbb{R}^D \\ b_1 \in \mathbb{R} \dots b_v \in \mathbb{R} \end{array} \Rightarrow \underbrace{\vee}_{\vee(D+1)}$$

2. Histories / Conditioning context are no longer treated as unrelated to one another

	cool	dog	cat	dry	positive	animal	positive animal
cool dog	1	1	0	0	1	1	1
cool cat	1	0	1	0	1	1	1
cool dry	1	0	0	1	1	0	0

③ Estimation

$$D = \left\{ \left(h_n, w_n \right)_{n=1}^N \right\}$$

$$W | H=h \sim \text{Categorical} \left(f(h; \theta) \right)$$

$$\phi(h) \in \mathbb{R}_{V \times D}^D$$

$$s = W \phi(h) + b$$

$$W \in \mathbb{R}^V$$

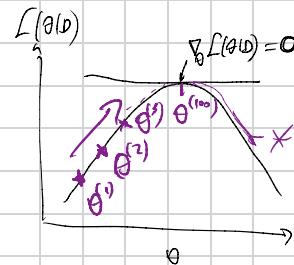
$$f(h; \theta) = \text{softmax}(s)$$

$$b \in \mathbb{R}^V$$

$$L(\theta | D) = \sum_{n=1}^N \log p(w_n | h_n; \theta)$$

$$= \sum_{n=1}^N \log [f(h_n; \theta)]_{w_n}$$

$$= \sum_{n=1}^N \log \left[\text{softmax} (W \phi(h_n) + b) \right]_{w_n}$$



$$\theta^{(t+1)} = \theta^{(t)} + \gamma_{t+1} \nabla_{\theta} L(\theta^{(t)} | D)$$