## distributional hypothesis.

↓

## language modeling.

$$X = (x_1, x_2, \cdots, x_\tau) \qquad x_t \in V$$

$$\underbrace{\phantom{X = (x_1, x_2, \cdots, x_\tau)}}_{\text{text}}$$

$$p(X) = p(x_1, x_2, \cdots, x_\tau)$$

① fully factored : $p(x_1, x_2, \cdots, x_\tau) = \prod_{t=1}^{T} p(x_t)$

② <mark>a chain rule of probability</mark>

$$p(x_1, x_2 \cdots, x_\tau) = p(x_{o(1)}) \, p(x_{o(2)} \mid x_{o(1)}) \, p(x_{o(3)} \mid x_{o(1)} o(2)) \cdots p(x_{o(\tau)} \mid -)$$

$$\underline{o : \{1, \cdots, T\} \rightarrow \{1, \cdots, \tau\}}$$

$$o(i) = i$$

③ latent variable models.

$$p(x_1 \cdots, x_\tau) = \int p(z) \prod_{t=1}^{T} p(x_t \mid z) \, dz$$

---

## the goal of LM.

. $p(X; \theta) > p(x' \mid \theta)$   if  <mark>$X \sim D^*$</mark>  &  $x' \not\sim D^*$

. unsupervised lng.

---

## learning.   maximum log-likelihood lng

$$D: \text{data} / \text{corpus}$$

$$L(\theta; D) = -\frac{1}{|D|} \sum_{X \in D} \log p(X; \theta) \qquad \textcolor{red}{-\frac{1}{|D|} \sum_{x,y \in D} \log p(y \mid x; \theta)}$$

## Autoregressive models

$$L(\theta; D) = -\frac{1}{|D|} \sum_{X \in D} \log p(X; \theta)$$

$$= -\frac{1}{|D|} \sum_{n=1}^{|D|} \sum_{t=1}^{T^n} \textcolor{purple}{\log p(X_{o(t)} \mid X_{o(1):o(t-1)}; \theta)},$$

$$\underbrace{\phantom{\log p(X_{o(t)} \mid X_{o(1):o(t-1)}; \theta)}}_{\textcolor{purple}{\text{per-step}}}$$

$$\underbrace{\phantom{XXXXXXXXXXXXXXXX}}_{\textcolor{red}{\text{per-example}}}$$

$$\underbrace{\phantom{XXXXXXXXXXXXXXXX}}_{\textcolor{green}{\text{total}}}$$

$$\log p(X_t \mid X_1 \cdots X_{t-1}; \theta)$$

$$p(X_t \mid X_{<t}) = \frac{p(X_{<t}, X_t)}{p(X_{<t})} = \frac{p(X_{<t}, X_t)}{\sum_{x' \in V} p(X_{<t}, X_t = x')}$$

$$= \frac{\tilde{p}(X_{<t}, X_t) / \cancel{Z}}{\sum_{x' \in V} \tilde{p}(X_{<t}, X_t = x') / \cancel{Z}} = \textcolor{purple}{\frac{c(X_{<t}, X_t)}{\sum_{x' \in V} c(X_{<t}, X_{t'})}}$$

## Count-based estimation

$$N \text{ tosses: } (M \text{ heads} \qquad (N-M) \text{ tails}$$

$$p(H) = \frac{M}{N} \qquad p(T) = \frac{N-M}{N}.$$

Maximum likelihood estimation
(count-based)

$$p(E_i) = \frac{c(E_i)}{\sum_j c(\bar{E}_j)}$$

---

Sparsity.

$$|V|^t \gg |D|$$

0-count

---

n-th order Markov assumption

$$p(x_t | x_1 \cdots, x_{t-1}) = p(x_t | x_{t-n} \cdots x_{t-1})$$

$$\frac{|D|}{O(|V|^n)} \gg \frac{|D|}{O(|V|^T)}$$

$$p(x_t | x_{<t}) = p(x_t | x_{t-n:t-1}) \approx \frac{c(x_{t-n}, \cdots, x_t)}{\sum_{x' \in V} c(x_{t-n}, \cdots, x')}$$

---

Smoothing : MAP version

$$p(x_t | x_{<t}) \approx p(x_t | x_{t-n:t-1}) = \frac{c(x_{t-n}, \cdots, x_t) + \varepsilon}{\sum_{x' \in V} c(x_{t-n}, \cdots, x') + |V|\varepsilon}$$

---

$$d(w, v) = \begin{cases} 0, & \text{if } w = v \\ 1, & \text{otherwise} \end{cases}$$

---

Feedforward language model

$$F(x_{t-n}, x_{t-n+1}, \cdots, x_{t-1}, x) = u_x^T \phi(x_{t-n}, \cdots, x_{t-1}) + b_x$$