

$$C \cong v_j \rightarrow \alpha \approx 1 \rightarrow \frac{e^{K_{i1}q}}{\sum_{j=1}^n e^{K_{ij}q}} \approx 1 \rightarrow K_j q \gg K_i q \text{ for } j \neq i$$

(1a)

$$C \cong \frac{1}{2} (v_a + v_b) \rightarrow \alpha_a = \alpha_b = \left[\frac{1}{2}, \frac{1}{2}, \dots \right]$$

$\alpha_i \text{ } i \neq b, i \neq a = \{0, 0, \dots\}$

$K_i \Rightarrow$ one hot vectors

$$q = \beta (K_a + K_b), \beta > 0$$

$$\alpha_b = \alpha_a = \frac{e^\beta}{2e^\beta} = \frac{1}{2} \quad \alpha_i \text{ } i \neq b, i \neq a = 0$$

$$C = 0v_1 + 0v_2 + \dots + 0v_i \dots +$$

$$\frac{1}{2}v_a + \frac{1}{2}v_b + \dots + 0v_n$$

$$\rightarrow C = \frac{1}{2} (v_a + v_b)$$

(1b)

$$(1 \quad c \quad l)$$

$$q = \beta (u_a + u_b), \beta \gg 0$$

since $\sum_i = \alpha I$ and α is

close to zero \rightarrow then we can

say K_i is imperfect one hot

because M_i is one hot. So

with large value of β

this problem is approximately
the same as previous one.

$$\{1, 2, \dots, L\}$$

$$K_a \sim \mathcal{N}(\mu_a, \Sigma_a)$$

$$\Sigma_a = \frac{1}{2} (\mu_a \mu_a^T) + \alpha \underbrace{I}_{\approx 0}$$

$$\rightarrow K_a \approx \gamma \mu_a$$

$$\gamma \sim \mathcal{N}(1, \frac{1}{2})$$

ex:

$$\gamma = [0.3, 0.7, 0.5, \dots]$$

$$q = \beta (u_a + u_b)$$

$$K_a \cdot q = \delta_a \beta$$

$$K_b \cdot q = \delta_b \beta$$

$$K_i \cdot q = 0$$

$$i \neq a, i \neq b$$

$$\rightarrow \alpha_a = \frac{e^{\delta_a \beta}}{e^{\delta_a \beta} + e^{\delta_b \beta}}$$

$$\alpha_b = \frac{e^{\delta_b \beta}}{e^{\delta_a \beta} + e^{\delta_b \beta}}$$

$$C \approx \frac{e^{\delta_a \beta}}{e^{\delta_a \beta} + e^{\delta_b \beta}} V_a$$

$$+ \frac{e^{\delta_b \beta}}{e^{\delta_a \beta} + e^{\delta_b \beta}} V_b$$

if $|\delta_a| > |\delta_b|$

↳ c looks like V_a

if $|\delta_b| > |\delta_a|$

↳ c looks like V_b

(1 di)

$$C = \frac{1}{2} (c_a + c_b)$$

$$C \simeq \frac{1}{2} (V_a + V_b)$$

$$\hookrightarrow V_a + V_b \simeq C_a + C_b$$

$$\hookrightarrow q_a = \beta_a n_a \quad q_b = \beta_b n_b$$

$$\beta_a, \beta_b > 0$$

$$\alpha_a = \frac{e^{\beta_a \mu_a}}{0 + 0 + 0 + e^{\beta_a \mu_a}} = 1$$

$$\hookrightarrow C_a = V_a$$

$$\alpha_b = \frac{e^{\beta_b \mu_b}}{0 + 0 + e^{\beta_b \mu_b}} = 1$$

$$\hookrightarrow C_b = V_b$$

$$C = \frac{1}{2} (V_a + V_b)$$

Idil

(like Ici)

$$K_a = \gamma_a \mu_a$$

$$K_b = \gamma_b \mu_b$$

$$q_a = \beta_a \mu_a \quad q_b = \beta_b \mu_b$$

$$\beta_a, \beta_b > 0$$

$$K_a \circ q_a = \gamma_a \beta_a \rightarrow \alpha_a \approx 1$$

$$K_b \circ q_b = \gamma_b \beta_b \rightarrow \alpha_b \approx 1$$

$$C_a \approx V_a \quad C_b \approx V_b$$

$$C \approx \frac{1}{2} (V_a + V_b)$$

lei

$$\alpha_2 = \frac{e^{|v_a|}}{e^0 + e^0 + e^{|v_a|}} \approx 1$$

$$\alpha_i, i \neq 2 = \frac{e^0}{e^0 + e^0 + e^{|v_a|}} \approx 0 \quad \text{since } \beta \text{ is very large}$$

and $|v_a| = |\beta|$

$$C_2 \approx \alpha_2 v_2 \approx v_2 \approx x_2$$

$$\hookrightarrow C_2 \approx v_a = v_a$$

no it cannot approx v_b

unlike u_a and u_c ~~does~~ not connect
any x solely

So for example by
adding m_d , x_i weight
increases but this does so
by increasing v_d and v_b .
the same thing happens for
 m_c .

$$V = (|u_b| - |u_c|) \cdot \frac{1}{\beta^2}$$

$$Q = (u_d^{\cancel{u_a}} + u_c u_d) \cdot \frac{1}{\beta^2}$$

$$K = \overline{I}$$

$$\rightarrow V_1 = u_b \quad \alpha_1 = u_c$$

$$V_2 = 0 \quad \alpha_2 = u_d$$

$$V_3 = u_b - u_c$$

$$\alpha_3 = 0$$

$$\alpha_1 = 0, 0, 1 \quad \alpha_2 = 1, 0, 0$$

$$c_1 \simeq v_3 = u_b - u_c \rightarrow c_2 = v_1 = u_b$$

2d

dev 1.0%

gendon 5%

28

20%

20gi

13.6%

29/12

Syntactic self

attention is not able
to understand the context
of different positions.

CSS in part e
provided more knowledge
the the model. (3a)

① misinformation.

② biased output. (3b)

(It may use some notion
of similarity between name to
purpose the birth place. (3c)
This causes bias (racism, sexism and ...)