# Masked language modeling

$$J(\theta) = \frac{1}{N} \sum_{n=1}^{N} \mathbb{E}_{M,O \sim \gamma(T_n)} \left[ \sum_{i=1}^{T_m^n} \log p(x_{m_i} | x_{O_1}, \ldots, x_{O_{T_0}}; \theta) \right]$$

① MLM $\Rightarrow$ AR-LM

left to-right ordering

$$p(X) = \prod_{t=1}^{T} p(x_t | x_1, \ldots, x_{t-1}, x_{t+1} = \langle mask \rangle, \ldots, x_T = \langle mask \rangle)$$

② Pseudo-likelihood (Besag, 1971)

$$F(x) = PLL(x) = \sum_{t=1}^{T} \log p(x_t | x_1 \ldots x_{t-1}, x_{t+1} \ldots, x_T)$$

$$p(X) = \frac{\exp\{F(x)\}}{\sum_{x' \in L} \exp\{F(x')\}}, \quad \text{where} \quad L: \text{all possible text snippets.}$$

③ $\ell \sim \text{length}(\ldots)$
$\tilde{X} \sim \text{random}(\ell)$

for $k = 1, \ldots, K$
- uniformly sample an index $i$ at random
- $\tilde{x}_i \sim p(x_i | x_1 \ldots, x_{i-1}, \langle mask \rangle, x_{i+1} \ldots x_T)$.

---

## Classification : constrained optimization

$$\min_{\theta} R(\theta)$$

$\propto x)$.
$\|\theta\|^2$

subject to

$$d(y_n, f(x_n;\theta)) < \varepsilon \quad \text{for all } n=1,\cdots,N$$
per-example loss.

$$\left\{ d'(f'(x'_m;\theta)) < \delta \quad \text{for all } m=1,\cdots,\boxed{M} \right.$$

$$\min_\theta \frac{1}{N} \sum_{n=1}^{N} d(y^a_i, f(x^a;\theta)) + \lambda R(\theta)$$
$$\underset{\geq 0}{\Uparrow}$$

$\underline{SGD}$

$\underline{SGD \text{ cannot move too far.}}$

Obs 1) S.G. : an unbiased estimate of the gradient

$$\left\| \left( \nabla - \mathbb{E}_{B \sim batches} \left[ \sum_{x \in B} \frac{1}{|B|} \nabla_\theta \ell(x;\theta) \right] \right) \right\|^2 = 0$$

$$\nabla \longrightarrow 0$$

Obs 2)

$$Cov_{B \sim batches} [\nabla_B] \gg 0$$

$\underline{\text{Incorporate unlabelled examples for initialization.}}$

stage 1 (pretrain)

$$\theta_0 \overset{(approximate)}{=} \underset{\theta}{arg\,min} \frac{1}{M} \sum_{m=1}^{M} d'(f'(x'_m;\theta))$$
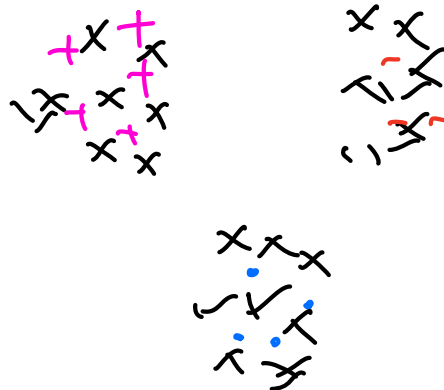$$\underbrace{\qquad}_{\text{density modeling}}$$

stage 2 (finetune)

$$\theta_{t+1} = \theta_t - \eta \frac{1}{|B|} \sum_{(x,y) \in B} \nabla_\theta d(y, f(x;\theta))$$

$$p(y|x) = \frac{p(x|y)\,p(y)}{p(x)} \qquad \overset{= p(x,y)}{\underset{}{\alpha}} \; \underbrace{p(x|y)\,p(y)}_{\substack{\text{Bayes'}\\\text{classifier}}} = \underbrace{p(y)\prod_{i=1}^{d} p(x_i|y)}_{\text{Naive Bayes' classifier}}$$

$$\underbrace{\hspace{3cm}}_{\text{Bayes' rule}}$$

*descriminative*

*Generative*

---



---

Classification **RBM**

$$-\underbrace{\mathbb{E}(x,y)}_{\text{score}} = b_x^T x + b_y^T y + \sum_{k=1}^{K} \log\left(1 + \exp\{w_k^T x + u_k^T y + c_k\}\right).$$
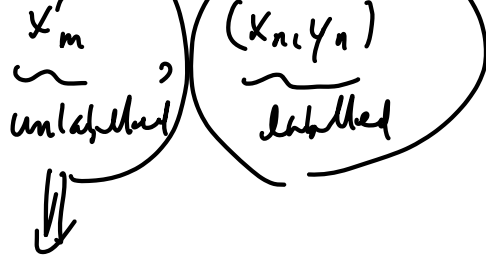
$$\boxed{p(x,y)} = \frac{\exp\{-E(x,y)\}}{\sum_{y\in C}\int_x \exp\{-E(x,y)\}\,dx}$$

losses

$$\boxed{p(y|x)} = \frac{\exp\{-E(x,y)\}}{\sum_{y\in C}\exp\{-E(x,y')\}}$$

$$\boxed{p(x)} = \sum_{y'\in C} p(x,y')$$

*data*

$$\overbrace{x'_m}^{\text{unlabelled}}, \quad \overbrace{(x_n, y_n)}^{\text{labelled}}$$

$$\nabla_\theta \log \sum_y p(x,y) = \mathbb{E}_{y \sim Q(y)} \left[ \nabla_\theta \log p(x,y) \right]$$

$$Q(y) = \frac{\exp\{-E(x,y)\}}{\sum_{y'} \exp\{-E(x,y')\}} = p(y|x)$$

$$\nabla_\theta \log p_\theta(x,y) = -\nabla_\theta E(x,y) + \mathbb{E}_{x',y' \sim p_\theta(x,y)} \left[ \nabla E(x',y') \right]$$

$$-E(x,y,h) = b_x^T x + b_y^T y + \sum_{k=1}^{K} h_k (W_k^T x + U_k^T y + c_k).$$

$\{0,1\}$

$$p(h_k = 1 | x,y) = \overset{\text{sigmoid}}{\sigma} (W_k^T x + U_k^T y + c_k)$$

$$p(y = c | h) \propto \exp\{\bar{u}_c^T h + b_{y,c}\}$$

$$p(x_i = 1 | h) = \sigma(\bar{w}_i^T h + b_x^i).$$

$$\bar{U} = \begin{bmatrix} u_1^T \\ u_2^T \\ \vdots \\ u_K^T \end{bmatrix} \quad \begin{matrix} c^{th} \\ \bar{u}_c \end{matrix}$$

Word embeddings : <u>Continuous bag-of-words (CBoW)</u>

- $p(x_{m_1}, \dots, x_{m_{T_m}} | \text{corrupt}(x))$

$$= \prod_{i=1}^{T_m} p(x_{m_i} | \text{corrupt}(x))$$

$$\approx \prod_{i=1}^{T_m} p(x_{m_i} | x_{m_i - n/2}, \dots, x_{m_i - 1}, <\text{mask}>, x_{m_i + 1}, \dots, x_{m_i + n/2}).$$

softmax

$$F(x_{m_i - n/2} \dots x_{m_i} \dots x_{m_i + n/2}) = e(x_m)^T \left( \sum_{j=m-n/2}^{m-1} e(x_j) + \sum_{j=m+1}^{m+n/2} e(x_j) \right)$$

$$\frac{}{C\,B,W}$$