



$$p(\bar{x} | \bar{\theta}) = \sum_k \pi_k p_k(\bar{x} | \bar{\theta}_k)$$

$$p(D | \bar{\theta}) = \prod_n \left( \sum_k \pi_k p_k(\bar{x}_n) \right)$$

$\bar{\theta} \rightarrow \max$

$z: z_{nk} = 1 \Leftrightarrow x_n \in C_k$

$$p(D, z | \bar{\theta}) = \prod_n \prod_k \left( \pi_k p_k(\bar{x}_n) \right)^{z_{nk}} \rightarrow \max_{\bar{\theta}}$$

$$E_{\theta^{(m)}} [z_{nk}] \rightsquigarrow \theta^{(m+1)} = \arg \max_{\theta} p(D, E[z_{nk}] | \theta)$$

EM-algorithm

latent vars  $z$

$$l(\theta) = \log p(x | \theta) \rightarrow \max_{\theta}$$

$$p(x, z | \theta)$$

$\log p(\theta | x) \rightarrow \max$

$$\theta^{(0)} \rightarrow \theta^{(1)} \rightarrow \theta^{(2)} \rightarrow \dots \rightarrow \theta^{(m)} \rightarrow \theta^{(m+1)} \rightarrow \dots$$

$$l(\theta^{(m)}) \leq l(\theta^{(m+1)})$$

$$l(\theta) - l(\theta^{(m)}) = \log p(x | \theta) - \log p(x | \theta^{(m)}) =$$

$$= \log \int p(x, z | \theta) dz - \log p(x | \theta^{(m)}) =$$

$$= \log \int p(x, z | \theta) \cdot \frac{p(z | x, \theta^{(m)})}{p(z | x, \theta^{(m)})} dz - \log p(x | \theta^{(m)})$$

$$= \log E_{p(z | x, \theta^{(m)})} \left[ \frac{p(x, z | \theta)}{p(z | x, \theta^{(m)})} \right] - \log p(x | \theta^{(m)}) \geq$$



$$f(ax + (1-a)y) \geq af(x) + (1-a)f(y)$$

$$f(E_{p(x)}[x]) \geq E_{p(x)}[f(x)]$$

Jensen's inequality



$$Q(\theta, \theta^{(m)}) = \int p(z|x, \theta^{(m)}) \log p(x, z|\theta) dz =$$

$$\underset{\theta}{\max} \mathbb{E}_{z|\theta^{(m)}} [\log p(x, z|\theta)]$$

$$X = (\bar{x}_1, \dots, \bar{x}_n) \quad \theta = (\bar{\pi}, \bar{\theta}_1, \dots, \bar{\theta}_K)$$

$$z = (\bar{z}_1, \dots, \bar{z}_n), \quad z_{nk} = 1 \Leftrightarrow \bar{x}_n \in C_k$$

$$\bar{z}_n = (0 \dots 1 \dots 0)$$

$$p(x|\theta) = \prod_n \left( \sum_k \pi_k p_k(\bar{x}_n | \bar{\theta}_k) \right) \rightarrow \max$$

$$\log p = l(\theta)$$

$$p(x, z|\theta) = \prod_n \prod_k \left( \pi_k p_k(\bar{x}_n | \bar{\theta}_k) \right)^{z_{nk}}$$

$$\log p(x, z|\theta) = \sum_n \sum_k z_{nk} (\log \pi_k + \log p_k(\bar{x}_n | \bar{\theta}_k))$$

$$Q(\theta, \theta^{(m)}) = \mathbb{E}_{z|\theta^{(m)}} [\log p(x, z|\theta)] =$$

$$= \mathbb{E}_{p(z|x, \theta^{(m)})} \left[ \sum_n \sum_k z_{nk} (\log \pi_k + \log p_k(\bar{x}_n | \bar{\theta}_k)) \right] =$$

$$= \sum_n \sum_k \left( \mathbb{E}_{z|x, \theta^{(m)}} [z_{nk}] \cdot (\log \pi_k + \log p_k(\bar{x}_n | \bar{\theta}_k)) \right) \rightarrow \max_{\theta}$$

$$Q(\theta, \theta^{(m)}) = \mathbb{E}_{z|\theta^{(m)}} [\log p(x, z|\theta)] \rightarrow \max_{\theta}$$

① because argmax happens  $\theta$ :  $Q(\theta, \theta^{(m)}) \geq Q(\theta^{(n)}, \theta^{(m)})$   
Generalized EM

②  $Q(\theta, \theta^{(m)}) \approx \frac{1}{R} \sum_{r=1}^R \log p(x, z_r | \theta)$ ,  $z_r \sim p(z|x, \theta^{(m)})$   
 Monte Carlo EM  
 (R=1)  $\rightarrow$  Stochastic EM

Пример Ceppellini et al., 1955

MN - ученик  
 ↑  
 гон.



$$p = \frac{2n_{MM} + n_{MN}}{2(n_{MM} + n_{MN} + n_{NN})}$$

$$q = \frac{2n_{NN} + n_{MN}}{2(n_{MM} + n_{MN} + n_{NN})}$$

Характеристики:

MM, NN - с равной

$$\frac{p^2}{p^2 + 2pq}$$

MN - " - "

$$\frac{2pq}{p^2 + 2pq}$$

$$\theta_{MLE} = \underset{\theta}{\operatorname{argmax}} \log p(\theta|x) = \underset{\theta}{\operatorname{argmax}} (\log p(\theta) + \log p(x|\theta))$$

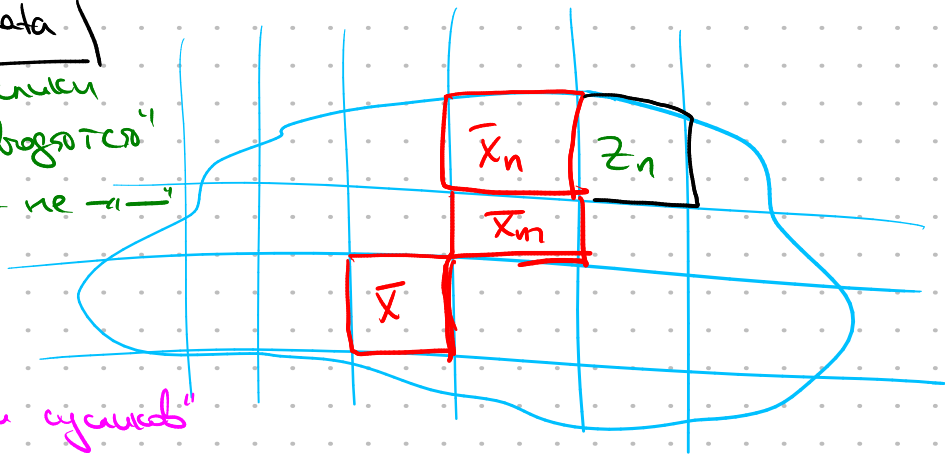
$$\begin{aligned} \log p(\theta|x) - \log p(\theta^{(m)}|x) &= \text{const} + \frac{\log p(x|\theta) + \log p(\theta)}{\log p(x|\theta^{(m)}) - \log p(\theta^{(m)})} \\ &\geq \underbrace{L(\theta, \theta^{(m)}) + \log p(\theta)}_{\text{const}} - \log p(\theta^{(m)}) \text{ const} \end{aligned}$$

$$\theta^{(m+1)} = \underset{\theta}{\operatorname{argmax}} (Q(\theta, \theta^{(m)}) + \log p(\theta))$$

Presence-only data

$z_n = 1 \Leftrightarrow$  "ученик бодрится"

$z_n = 0 \Leftrightarrow$  "не"



$y_n = 1 \Leftrightarrow$  "бегает ученик"

$y_n = 1 \Rightarrow z_n = 1$

$y_n = 0 \Rightarrow$  ?

$$p(z=1|x) = \sigma(\eta(x)) = \frac{1}{1 + e^{-\eta(x)}}$$

# Prospective vs. retrospective studies

	d=0	d=1	
x=0	$\pi_{00}$	$\pi_{01}$	$\pi_{0x}$
x=1	$\pi_{10}$	$\pi_{11}$	$\pi_{1x}$
$p(d=1)$	$\pi_{x0}$	$\pi_{x1}$	

Prospective

$$p(d=1 | x=0) = \frac{\pi_{01}}{\pi_{00} + \pi_{01}}$$

$$p(d=1 | x=1) = \frac{\pi_{11}}{\pi_{10} + \pi_{11}}$$

$\sigma(\eta(x))$   
 $\approx \sigma(\eta(x))$

## Retrospective

$$\pi_0 = p(s=1 | d=0)$$

$$\pi_1 = p(s=1 | d=1)$$

$$\frac{\pi_{01}}{\pi_{00} + \pi_{01}} = p(d=1 | x=0, s=1) =$$

$$= \frac{p(d=1 | x=0) p(s=1 | d=1, x=0)}{p(d=1 | x=0) p(s=1 | d=1, x=0) + p(d=0 | x=0) p(s=1 | d=0, x=0)}$$

$\pi_1$  (under  $p(s=1 | d=1, x=0)$ )  
 $\pi_1$  (under  $p(d=1 | x=0)$ )  
 $\pi_0$  (under  $p(s=1 | d=0, x=0)$ )

$$\frac{\pi_{01}}{\pi_{00} + \pi_{01}} = \frac{\pi_1 p(d=1 | x=0)}{\pi_1 p(d=1 | x=0) + \pi_0 (1 - p(d=1 | x=0))}$$

$$\pi_{01} (\pi_1 p + \pi_0 - \pi_0 p) = \pi_1 p (\pi_{00} + \pi_{01})$$

$$\pi_0 \pi_{01} - \pi_0 \pi_{01} p = \pi_1 p \pi_{00}$$

$$\sigma(\eta(x)) = p = \frac{\pi_0 \pi_{01}}{\pi_{00} \pi_1 + \pi_{01} \pi_0}$$

$$\frac{\pi_{01}}{\pi_{00} + \pi_{01}} = \frac{\pi_1 \sigma}{\pi_0 + (\pi_1 - \pi_0) \sigma} = \frac{\pi_1 / (1 + e^{-\eta(x)})}{\pi_0 + \frac{\pi_1 - \pi_0}{1 + e^{-\eta(x)}}} =$$

$$= \frac{\pi_1}{\pi_0 (1 + e^{-\eta(x)}) + \pi_1 - \pi_0} = \frac{\pi_1}{\pi_1 + \pi_0 e^{-\eta(x)}}$$

$$\frac{\pi_{01}}{\pi_{00} + \pi_{01}} = \frac{1}{1 + \frac{\tau_0}{\pi_1} e^{-\eta(x)}} = \frac{1}{1 + e^{-(\eta(x) - \ln \frac{\pi_0}{\pi_1})}}$$

$$\frac{\pi_{01}}{\pi_{00} + \pi_{01}} = \sigma \left( \underbrace{\eta(\bar{x})}_{w_0 + w_1 x} + \underbrace{\ln \frac{\pi_0}{\pi_1}} \right)$$

$$\sigma(\eta^*(x)) \approx \frac{\pi_{01}}{\pi_{00} + \pi_{01}} \longrightarrow w_0^*, w_1^*$$

$$\Rightarrow w_1 = w_1^*, w_0 = w_0^* - \ln \frac{\pi_0}{\pi_1}$$

$$\sigma(\eta(\bar{x})) = p(z=1|\bar{x})$$

$s=1$  - выборка  
 $\pi = p(z=1)$  - вероятность

$$\sigma(\eta_{\text{naive}}(\bar{x})) = p(y=1|\bar{x}, s=1)$$

$n_p$  - positive  
 $n_u$  - unknown

выборка:

$$\left[ \begin{array}{l} z=1: \frac{n_p + \pi \cdot n_u}{n_p + \pi \cdot n_u} \\ z=0: \frac{(1-\pi)n_u}{n_p + \pi \cdot n_u} \end{array} \right]$$

$$\begin{aligned} p(y=1|\bar{x}, s=1) &= \\ &= p(y=1|z=1, s=1, \bar{x}) p(z=1|s=1, \bar{x}) \\ &+ p(y=1|z=0, s=1, \bar{x}) p(z=0|s=1, \bar{x}) = \\ &= p(y=1|z=1, s=1, \bar{x}) p(z=1|s=1, \bar{x}) \end{aligned}$$

$$= \sigma(\eta(\bar{x}) - \log \frac{\pi_0}{\pi_1})$$

напрям. выбор.

$$\frac{p(y=1, z=1|s=1)}{p(z=1|s=1)} = \frac{n_p / N}{(n_p + \pi n_u) / N} = \frac{n_p}{n_p + \pi n_u}$$

$$\pi_1 = p(s=1|z=1) = \frac{p(z=1|s=1) p(s=1)}{p(z=1)} = \frac{n_p + \pi n_u}{n_p + n_u} p(s=1)$$

$$\tau_0 = p(s=1|z=0)$$

$$\pi_1 = \frac{n_p + \pi n_u}{\pi(n_p + n_u)} = p(s=1)$$

$$\pi_0 = \frac{(1-\pi)n_u}{(1-\pi)(n_p + n_u)} = p(s=0)$$

$$\frac{\pi_0}{\pi_1} = \frac{\pi n_u}{n_p + \pi n_u}$$

$$\sigma(\eta_{\text{naive}}(\bar{x})) = \frac{n_p}{n_p + \pi n_u} \cdot \sigma(\eta(\bar{x}) - \log \frac{\pi n_u}{n_p + \pi n_u})$$

$$p(y|X, \bar{\eta}) = \prod_{i=1}^{n_p + n_u} p(y_i | s_i=1, \bar{x}_i) =$$

$$= \prod_{i=1}^{n_p + n_u} \left( p(y_i=1 | s_i=1, \bar{x}_i) \right)^{y_i} \left( 1 - p(y_i=1 | s_i=1, \bar{x}_i) \right)^{1-y_i} =$$

$$= \prod_{i=1}^{n_p + n_u} \left( \frac{n_p}{n_p + \pi n_u} \sigma(\eta(\bar{x}) - \log \frac{\pi n_u}{n_p + \pi n_u}) \right)^{y_i} \left( 1 - \frac{n_p}{n_p + \pi n_u} \sigma(\dots) \right)^{1-y_i} \xrightarrow{\max} \bar{\eta}$$

$\stackrel{=}{=} \eta_{\text{naive}}^{(0)}$

EM-algorithm:  $\bar{\eta}^{(0)} \rightarrow \bar{\eta}^{(1)} \rightarrow \bar{\eta}^{(2)} \rightarrow \dots \rightarrow \bar{\eta}^{(m)} \rightarrow \bar{\eta}^{(m+1)} \rightarrow \dots$

M-step:  $\sigma(\eta(\bar{x}_i) - \log \frac{\pi_0}{\pi_1}) \approx z_i^{(m+1)}$ ,  $\eta^{(m+1)} = \eta^* + \log \frac{\pi_0}{\pi_1}$

E-step:  $z_i^{(m+1)} = \mathbb{E}[z_i] = \sigma(\eta^{(m)}(\bar{x}_i))$   $\text{gr} \Rightarrow y_i=0$   
 $= 1, \quad y_i=1$

pseudolabels

