

$p_1(\bar{x}; \theta_1)$ $p_2(\bar{x}; \theta_2)$ $p(\bar{x}) = \sum_k \pi_k p_k(\bar{x}; \theta_k)$
 $p(D|\pi, \theta) = \prod_n \left(\sum_k \pi_k p_k(\bar{x}_n; \theta_k) \right)$
 $p(D, \underline{z}|\pi, \theta) = \prod_n \prod_k (\pi_k p_k)$
 $\hat{\pi}, \hat{\theta} \rightarrow \max$

$X, \theta \quad p(X|\theta) \rightarrow \max \quad z \quad p(X, z|\theta)$

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

$\log p(\theta|x) = \log p(x|\theta) + \log p(\theta)$
 $\underline{l(\theta)} = \log p(x|\theta) \rightarrow \max$ $\int f(z)q(z)dz$

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

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

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

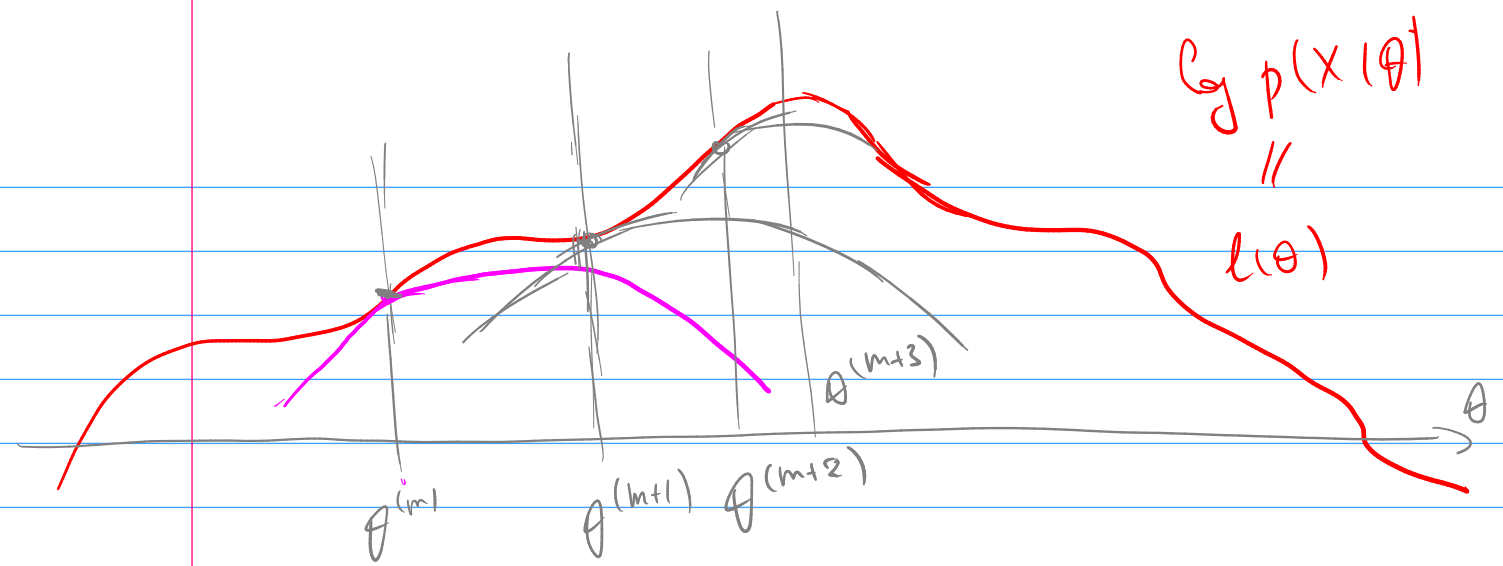
$\geq \int p(z|x, \theta^{(m)}) \log \frac{p(x, z|\theta)}{p(z|x, \theta^{(m)})} dz - \log p(x|\theta^{(m)})$
 $f(\alpha x + (1-\alpha)y)$

$= \int p(z|x, \theta^{(m)}) \log \frac{p(x, z|\theta)}{p(z|x, \theta^{(m)}) p(x|\theta^{(m)})} dz$
 $\geq \alpha f(x) + (1-\alpha)f(y)$
 $f(E_{p(x)}[x]) \geq E_{p(x)}[f(x)]$

$\underline{l(\theta)} \geq \underline{l(\theta^{(m)})} + \underline{l(\theta, \theta^{(m)})}$

$l(\theta^{(m)}, \theta^{(m)}) = 0$

$\underline{l(\theta, \theta^{(m)})} = E_{z|x} [\log p(x, z|\theta)] - E_z [\log \dots]$



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

$$E_{q(z)}[\dots] - c =$$

$$E_{q(z)}[\dots - c]$$

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

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

$$\theta = (\bar{\theta}_1, \bar{\theta}_2, \dots, \bar{\theta}_K, \bar{\pi})$$

$$z = (\bar{z}_1, \bar{z}_2, \dots, \bar{z}_N)$$

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

$$z_{nk} = 1 \text{ iff } \bar{x}_n \in C_k$$

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

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

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

$$\underbrace{p(z_{nk} | x, \theta^{(m)})}_{E\text{-step}} = p(\bar{x}_n \in C_k | \bar{x}_n, \theta^{(m)}) = \frac{p(C_k) p(\bar{x}_n | C_k)}{\sum_l p(C_l) p(\bar{x}_n | C_l)} = \frac{\pi_k^{(m)} p_k(\bar{x}_n | \theta_k^{(m)})}{\sum_l \pi_l^{(m)} p_l(\bar{x}_n | \theta_l^{(m)})}$$

$$\begin{aligned} Q(\theta, \theta^{(m)}) &= \mathbb{E}_{p(z|x, \theta^{(m)})} [\log p(x, z | \theta)] = \\ &= \sum_n \sum_k \left(\mathbb{E}_{p(z|x, \theta^{(m)})} [z_{nk}] \right) \cdot (\log \pi_k + \log p_k(\bar{x}_n | \theta_k)) \\ &\quad \xrightarrow[\pi_k, \theta]{N\text{-max}} \max \\ &\quad \sum_k \left(\sum_n \mathbb{E} z_{nk} \right) \log \pi_k \xrightarrow{\pi_k} \max \\ &\quad \sum_k \pi_k = 1, \pi_k \geq 0 \end{aligned}$$

$$\pi_k^{(m+1)} = \frac{\sum_n \mathbb{E} z_{nk}}{\sum_l \sum_n \mathbb{E} z_{nl}} = N$$

1955

M N

MM MN NN

$$\begin{aligned} p &= \frac{2n_{MM} + n_{MN}}{2(n_{\cdot\cdot} + n_{\cdot\cdot} + n_{\cdot\cdot})} \\ q &= \frac{2n_{NN} + n_{MN}}{2(\dots)} \end{aligned}$$

Закон Харди-Вайнберга

$$p = p(M), \quad q = 1 - p = p(N)$$

$$\Rightarrow \begin{aligned} p(MM) &= \frac{p^2}{p^2 + 2pq} \\ p(MN) &= \frac{2pq}{p^2 + 2pq} \end{aligned}$$

$$\mathbb{E}[z] = \frac{2p^{(m)}(1-p^{(m)})}{p^{(m)2} + 2p^{(m)}(1-p^{(m)})}$$

$$\begin{aligned} n_{NN}, \mathbb{E} n_{NN} &= (\mathbb{E} z) \cdot (n_{\cdot\cdot} + n_{\cdot\cdot} + n_{\cdot\cdot}), \dots = \mathbb{E} n_{MM} \end{aligned}$$

$$p^{(m+1)} = \frac{2 E \cdot n_{mm} + E n_{mw}}{2 \cdot (n_o + n_e + n_s)}$$

Presence-only data

"бывает наличие"

$y_n = 1 \Rightarrow$ "есть наличие"

$y_n = 0 \Rightarrow ?$

$z_n = 1$ - "есть наличие"

$z_n = 0$ - "нет наличия"

$\pi = p(z_n = 1)$ - prevalence

Если $y_n = 1 \Rightarrow z_n = 1$ с вер-1.

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

	$d=0$	$d=1$
$x=0$	π_{00}	π_{01}
$x=1$	π_{10}	π_{11}
	$p(d=0)$	$p(d=1)$

Retrospective studies

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

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

$$p(d=1 | x=0) = ?$$

Prospective studies

$\pi_{00} \quad \pi_{01}$

$$p(x=0)$$

$$p(x=1)$$

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

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

$$\frac{\sigma(\eta(x))}{1 + e^{-\eta(x)}}$$

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

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

$$p(d=1 | x=0) p(s=1 | d=1, x=0) +$$

$$+ p(d=0 | x=0) p(s=1 | d=0, x=0)$$

$$1 - p(d=1 | x=0)$$

$$\pi_0$$

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

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

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

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

$$\frac{\pi_{01}}{\pi_{00} + \pi_{01}} =$$

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

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

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

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

$$\eta(x) = w_0 + w_1 x$$

$$\eta^*(x) = \eta(x) - \ln \frac{\pi_0}{\pi_1}$$

$$w_0, w_0 - \ln \frac{\pi_0}{\pi_1}$$

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

$S=1$ - болезнь

$\pi = p(z=1)$ - вероятность

n_p - positive
 n_u - unknown

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

логит

логит:

$$z=1: n_p + \pi \cdot n_u$$

$$z=0: (1-\pi)n_u$$

$$p(y=1|\bar{x}, s=1) = \cancel{p(y=1|z=1, s=1, \bar{x})} \cdot p(z=1|s=1, \bar{x}) + \cancel{p(y=1|z=0, s=1, \bar{x})} \cdot p(z=0|s=1, \bar{x})$$

= 0

$$= p(y=1|z=1, s=1) \cdot p(z=1|s=1, \bar{x})$$

$\sigma(\eta_{\text{naive}}(\bar{x}))$

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

$$\eta(\bar{x}) = w_1 x_1 + \dots + w_d x_d$$

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

$$p(y=1|\bar{x}, s=1) = \frac{n_p}{n_p + \pi n_u} \cdot p(z=1|s=1, \bar{x})$$

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

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

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

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

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

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

$$= \prod_i p(y_i=1|s_i=1, \bar{x}_i)^{y_i} (1 - p(y_i=1|s_i=1, \bar{x}_i))^{1-y_i}$$


$$= \prod_i \left(\frac{n_p}{n_p + \pi n_u} \cdot \sigma(\eta(x)) - \log \frac{\pi n_u}{n_p + \pi n_u} \right)^{y_i} \times \left(1 - \frac{n_p}{n_p + \pi n_u} \sigma(\dots) \right)^{1-y_i} \rightarrow \max$$

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

$$Z: \quad \sigma(\underbrace{\eta(\bar{x}_i)} - \log \frac{\pi_0}{\pi_1}) \approx \underline{\underline{z_i}} \quad \eta^{(0)} = \eta_{\text{naive}}$$

M-wei: $\eta^{(m+1)}(\bar{x}) = \eta^*(\bar{x}) + \log \frac{\pi_0}{\pi_1}$, i.e. $\eta^* \circ \text{Syn}$
 Ha $z_i^{(m+1)}$

E-wei: $z_i^{(m+1)} = E[z_i] = \sigma(\eta^{(m)}(\bar{x})) \quad \text{gla } y_i = 0$
 $z_i^{(m+1)} = 1 \quad \text{gla } y_i = 1$

data imputation \equiv  pseudolabels