# ① Clustering — unsupervised learning
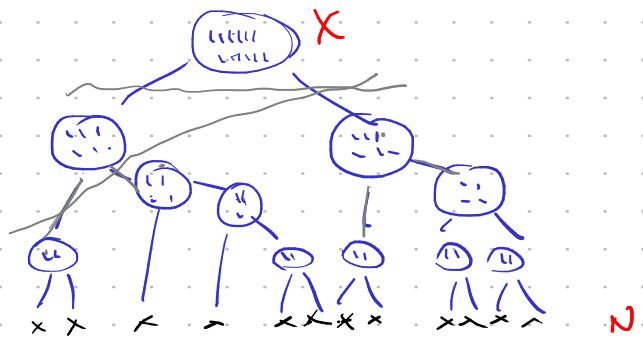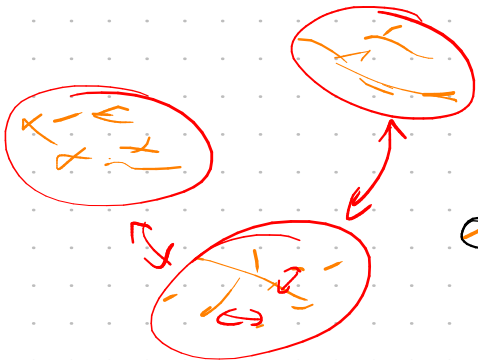
$\bar{x} \in \mathbb{R}^d$

$X = \{\bar{x}_n\}_{n=1}^N$

— Hierarchical clustering
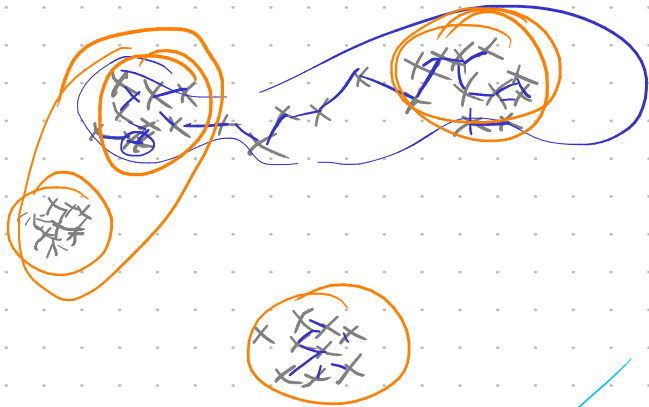
Agglomerative Clustering

X

N

— single-link clustering

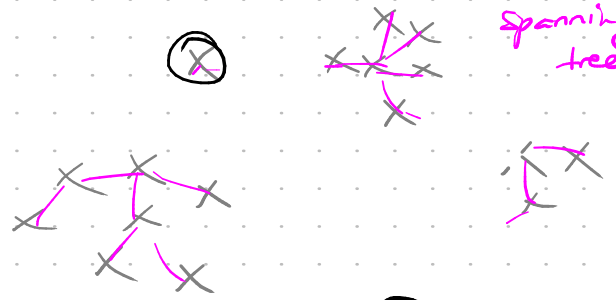$$d(C_1, C_2) = \min_{\substack{x \in C_1 \\ y \in C_2}} d(x, y)$$
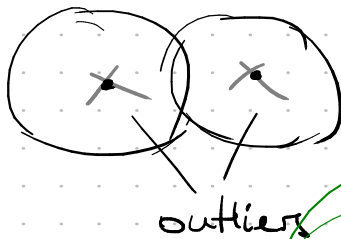
— complete-link clustering

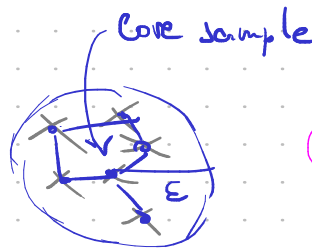$$d(C_1, C_2) = \max_{x \in C_1, y \in C_2} d(x, y)$$
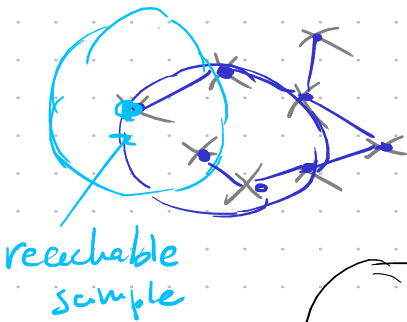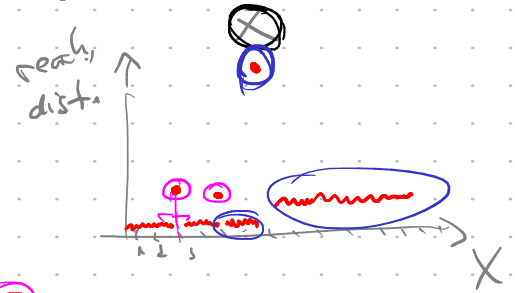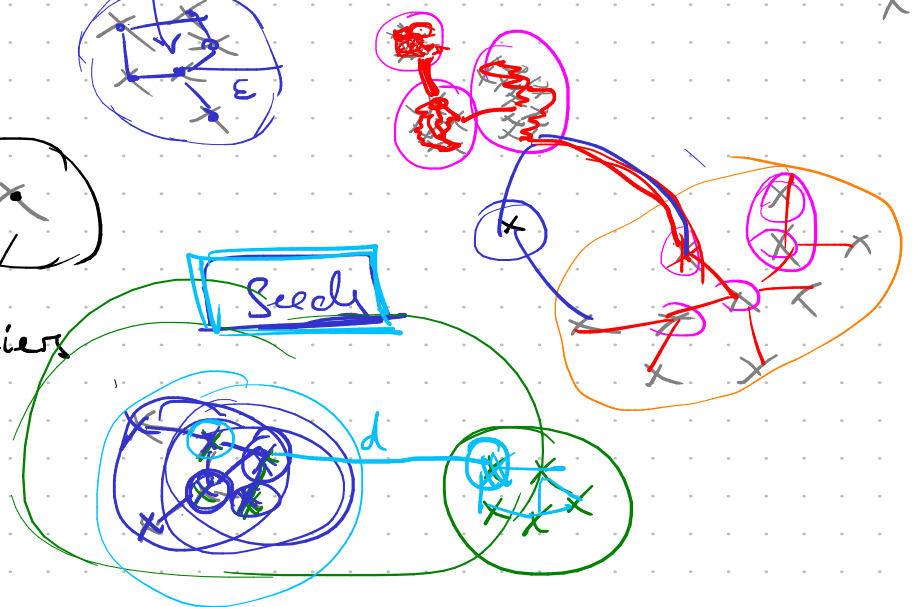
— graph theory

MST
minimal spanning tree

— DBSCAN ($\varepsilon$, min_samples)

Core sample

reach. dist.

reachable sample

outliers

Seeds

$d$

— OPTICS

— BIRCH

# ② Mixture models   GMM



$$\prod_n \left( y_n \cdot p_1(\bar{x}) + (1-y)_n p_2(\bar{x}) \right)$$

$$\prod_n p_1(\bar{x}_n) \, p_2(\bar{x}_n) \quad \boxed{y_n} \quad \boxed{1-y_n}$$

$\bar{x} \in \mathbb{R}^d$

$p(\bar{x})$

$$p(\bar{x}) = \alpha_1 \underbrace{p_1(\bar{x})}_{\bar{\theta}_1} + \alpha_2 \underbrace{p_2(\bar{x})}_{\bar{\theta}_2} + \alpha_3 \underbrace{p_3(\bar{x})}_{\bar{\theta}_3}$$

$$\alpha_1 + \alpha_2 + \alpha_3 = 1$$

$$X = \{\bar{x}_n\}_{n=1}^N$$

$$p(X \mid \bar{\pi}, \bar{\theta}_1 .. \bar{\theta}_k) = \prod_{n=1}^N \left( \pi_1 p_1(\bar{x}_n \mid \bar{\theta}_1) + \pi_2 p_2(\bar{x}_n \mid \bar{\theta}_2) + ... + \pi_k p_k(\bar{x}_n \mid \bar{\theta}_k) \right)$$

$$\xrightarrow{\bar{\pi}, \bar{\theta}_1 .. \bar{\theta}_k} max$$

 $\longrightarrow k - \#кластера \longrightarrow \bar{x} \sim p_k(\bar{x} \mid \bar{\theta}_k)$

$\bar{\pi}$

$$\bar{z} = (0 ... \underset{k}{\textcircled{1}} ... 0) \qquad z_k = [\bar{x} \in C_k]$$

$$p(X, Z \mid \bar{\pi}, \bar{\theta}_1 .. \bar{\theta}_k) = \prod_{n=1}^N p(\bar{x}_n, \bar{z}_n \mid \bar{\pi}, \bar{\theta}_1 ... \bar{\theta}_k) =$$

$$= \prod_{n=1}^N p(\bar{z}_n \mid \bar{\pi}) \, p(\bar{x}_n \mid \bar{z}_n, \bar{\theta}_1 ... \bar{\theta}_k) =$$

$$= \prod_{n=1}^N \prod_{k=1}^k \underbrace{p(z_n = k \mid \bar{\pi})}_{\pi_k}^{z_{nk}} \cdot \prod_{k=1}^k p_k(\bar{x}_n \mid \bar{\theta}_k)^{z_{nk}} =$$

$$= \prod_{n=1}^N \prod_{k=1}^k \left( \pi_k \cdot p_k(\bar{x}_n \mid \bar{\theta}_k) \right)^{z_{nk}}$$

$$\log p(X, Z \mid \bar{\pi}, \bar{\theta}_1 ... \bar{\theta}_k) = \sum_{n=1}^N \sum_{k=1}^k \left[ z_{nk} \log \pi_k + z_{nk} \log p_k(\bar{x}_n \mid \bar{\theta}_k) \right] =$$

$$= \sum_{k=1}^k \underbrace{\left( \underbrace{\sum_{n=1}^N z_{nk}}_{\#[z=k]} \right) \log \pi_k}_{\xrightarrow{\bar{\pi}} max} + \sum_{k=1}^k \underbrace{\left( \sum_{n=1}^N z_{nk} \log p_k(\bar{x}_n \mid \bar{\theta}_k) \right)}_{\xrightarrow{\bar{\theta}_k} max}$$

$$\hat{\pi}_k = \frac{\sum_n z_{nk}}{N}$$

$$\hat{\bar{\theta}}_k = \underset{\bar{\theta}_k}{argmax} \sum_{n: z_{nk} = 1} \log p_k(\bar{x}_n \mid \bar{\theta}_k)$$

(3) Expectation - Maximization Algorithm

$X, \bar{\theta}, Z$ - latent variables:  — $p(X|\bar{\theta}) \xrightarrow{\bar{\theta}} \max$    Трудно

— $p(X,Z|\bar{\theta}) \xrightarrow{\bar{\theta}} \max$    легко

(E-step)  fix $\bar{\theta}$,  find  $\overset{exp}{\boxed{\mathbb{E}[Z]}}$    $\boxed{\text{k-means:} \quad k = \arg\min d(\bar{x}_n, \bar{\mu}_k) \\ z_{nk}=1 \Leftrightarrow k = \arg\max_s p(z_{ns}=1|..)}$

$\mathbb{E}[z_{nk}] = p(z_{nk}=1 | \underbrace{\bar{\pi}, \bar{\theta}_1 - \bar{\theta}_k}, X) = p(z_{nk}=1 | \bar{x}_n, \bar{\pi}, \bar{\theta}_1 - \bar{\theta}_k) =$

$= \dfrac{p(z_{nk}=1, \bar{x}_n | \bar{\pi}, \bar{\theta}_1, \bar{\theta}_k)}{p(\bar{x}_n|\bar{\theta}) = \sum\limits_s p(\bar{x}_n, z_{ns}=1 | \bar{\theta})} = \dfrac{\pi_k \cdot p_k(\bar{x}_n|\bar{\theta}_k)}{\sum\limits_{s=1}^{k} \pi_s p_s(\bar{x}_n | \bar{\theta}_s)}$

(M-step)  fix  $\boxed{\mathbb{E}[Z]}$    $\mathbb{E}\left[ \log p(X,Z|\bar{\theta}) \right] \xrightarrow{\bar{\theta}} \max$

$\mathbb{E}_Z\left[ \log p(X,Z|\bar{\theta}) \right] = \mathbb{E}_Z\left[ \sum\limits_n \sum\limits_k z_{nk} \left( \log \pi_k + \log p_k(\bar{x}_n|\bar{\theta}_k) \right) \right] =$

$= \sum\limits_n \sum\limits_k \mathbb{E}[z_{nk}] \left( \log \pi_k + \log p_k(\bar{x}_n|\bar{\theta}_k) \right)$    $\boxed{\text{k-means:} \\ \bar{\mu}_k := \text{Avg } \bar{x}_n \\ \quad n: z_{nk}=1}$

$\sum\limits_k \left( \sum\limits_n \mathbb{E}[z_{nk}] \right) \log \pi_k \xrightarrow{\bar{\pi}} \max$

$\forall k \quad \sum\limits_n \mathbb{E}[z_{nk}] \log p_k(\bar{x}_n|\bar{\theta}_k) \xrightarrow{\theta_k} \max$

$\sum\limits_n \mathbb{E}[z_{nk}] \overset{\log}{\cdot} N(\bar{x}_n | \bar{\mu}_k, \Sigma_k) \to \max$

(4) EM-алгоритм в общем виде
$X$ — observables  $\bar{\theta}$,  $p(X|\bar{\theta}) \xrightarrow{\bar{\theta}} \max$  — Трудно

$p(X|\bar{\theta}) = \int p(X,Z|\bar{\theta}) dZ \xrightarrow{\bar{\theta}} \max$

EM-algorithm:  — init $\bar{\theta}^{(0)}$
— for $m = 0, 1, \ldots$
$Q(\bar{\theta}, \bar{\theta}^{(m)}) = \mathbb{E}_{p(Z|X,\bar{\theta}^{(m)})}\left[ \log p(X,Z|\bar{\theta}) \right]$;  $\bar{\theta}^{(m+1)} = \arg\max_{\bar{\theta}} Q(\bar{\theta}, \bar{\theta}^{(m)})$

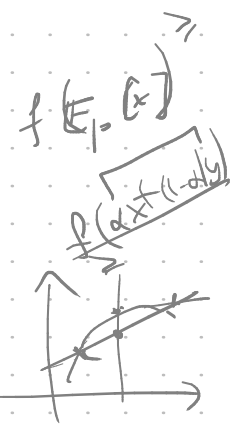<u>Yeas</u>: $\quad p(X|\bar\theta^{(m+1)}) \geq p(X|\bar\theta^{(m)})$



$\approx \sum \quad \sum_{k=1}^{K} \left(\pi_k\right) p_k(\bar x_n|\bar\theta_k)$

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

$$\log p(X|\bar\theta) - \log p(X|\bar\theta^{(m)}) \overset{\log}{=} \int p(x,z|\bar\theta)\,dz \quad \log(x|\bar\theta^{(m)}) = \quad f(E_1, x)$$

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

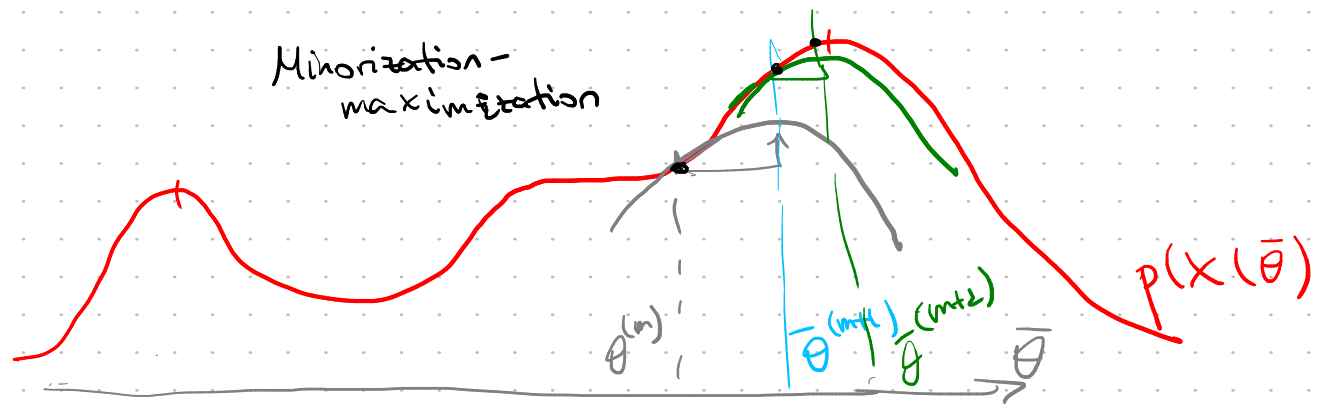$$= \log \mathbb{E}_{p(z|x,\bar\theta^{(m)})}\left[ \frac{p(x,z|\bar\theta)}{p(z|x,\bar\theta^{(m)})} \right] - \log p(x|\bar\theta^{(m)}) \geq \begin{bmatrix} \text{Jensen's} \\ \text{ineq.} \end{bmatrix}$$

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

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

$$\log p(X|\bar\theta) - \log p(X|\bar\theta^{(m)}) \geq \mathbb{E}_{p(z|\theta,\bar\theta^{(m)})}\left[ \log \frac{p(x,z|\bar\theta)}{p(x,z|\bar\theta^{(m)})} \right]$$

$$\log p(X|\bar\theta) \geq L(\bar\theta,\bar\theta^{(m)}) = \log p(x|\bar\theta^{(m)}) + \mathbb{E}[--]$$

Minorization –
maximization



$p(X|\bar\theta)$

$\theta^{(m)}$ $\quad \bar\theta^{(m+1)}$ $\bar\theta^{(m+2)}$ $\quad \bar\theta$

$$h_r(\bar\theta, \bar\theta^{(m)}) = \log p(x \mid \bar\theta^{(m)}) + \mathbb{E}_{p(\bar z \mid x, \bar\theta^m)}\left[\log \frac{p(x, z \mid \bar\theta)}{p(x, z \mid \bar\theta^{(m)})}\right] \to \max_{\bar\theta}$$

$$(=) \quad Q(\bar\theta, \bar\theta^{(m)}) \xrightarrow[\bar\theta]{} \max$$

---

**(5) EM для Mixture models**

$$p(x, z \mid \bar\theta) = \sum \pi_k p_k(\bar x \mid \bar\theta_k)$$

$$p(X, z \mid \bar\theta) = \prod_n p(\bar z_n \mid \bar\theta)\, p(\bar x_n \mid \bar z_n, \bar\theta) = \prod_n \prod_k \left(\pi_k p_k(\bar x \mid \bar\theta_k)\right)^{z_{nk}}$$

$$Q(\bar\theta, \bar\theta^{(m)}) = \mathbb{E}_{p(z \mid x, \bar\theta^{(m)})}\left[\log p(X, z \mid \bar\theta)\right] =$$

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

$$= \sum_n \sum_k \mathbb{E}[z_{nk}] \cdot \left(\log \pi_k + \log p_k(\bar x \mid \bar\theta_k)\right) \xrightarrow[\hat\pi, \hat\theta_1 \ldots \hat\theta_k]{\text{M-step}} \max$$

E-step

---

**(6) Ceppellini et al. 1955**    домин. +
                                    рецесс. —



| a | ++ |
| b | +− / −+ |
| c | −− |

фенотип A
фенотип B

$$q = p(+), \quad 1-q = p(-) \qquad \boxed{q = ?}$$

$$\hat q = \frac{2a + b}{2(a + b + c)} \qquad 1 - \hat q = \frac{b + 2c}{2(a + b + c)}$$

Дано: $(a+b)$, $c$

$$p(++) = q^2, \quad p(--) = (1-q)^2, \quad p(+-) = 2q(1-q)$$

$$q^{(0)}, \quad \mathbb{E}_{q^{(0)}}[a] = (a+b)\cdot p(++\mid \text{фен. } A) = (a+b)\cdot \frac{q^2}{q^2 + 2q(1-q)} = \left.\right]\text{E-step}$$

$$= (a+b)\cdot \frac{q}{2-q}$$

$$\mathbb{E}_{q^{(0)}}[b] = (a+b) - \mathbb{E}[a] = (a+b)\left(1 - \frac{q}{2-q}\right)$$

$$\text{M-step:}\quad q^{(1)} := \mathbb{E}\left[\frac{2a+b}{2n}\right] = \frac{a+b}{2n}\cdot\left(\frac{2q^{(0)}}{2-q^{(0)}} + \left(1 - \frac{q^{(0)}}{2-q^{(0)}}\right)\right) = \frac{a+b}{2n}\cdot\frac{2}{2-q^{(0)}}$$

$$2(a+b) + c$$

$\{A, C, G, T\}$

$L = $ размер alphabet

$\begin{vmatrix} A & C & C & G & T & G & A & A & G & C \\ A & C & G & C & T & G & A & T & G & C \end{vmatrix}$ $M$ — длина строки

$d = \frac{\#[\cdot]}{M}$

$k = 1, \ldots, H, \quad \pi_k = p(C_k)$

$C_k : \quad p(x_m = a_\ell | C_k)$

$= \frac{\theta_{km\ell}}{q}$

$\forall_k \forall_m \sum_\ell \theta_{km\ell} = 1$

$p(\bar{x} | C_k) = \prod_{m=1}^{M} p(x_m | C_k) = \prod_{m=1}^{M} \prod_{\ell=1}^{L} p(x_m = a_\ell | C_k)^{[x_m = a_\ell]}$

$= \prod_{m=1}^{M} \prod_{\ell=1}^{L} \theta_{km\ell}^{[x_m = a_\ell]}$

latent vars: $\quad z_{nk} = [\bar{x}_n \in C_k]$

$p(X, z | \bar{\theta}) = \prod_{n=1}^{N} \prod_{k=1}^{K} \left( \pi_k \, p(\bar{x}_n | C_k) \right)^{z_{nk}} = \prod_n \prod_k \left( \pi_k \prod_m \prod_\ell \theta_{km\ell}^{[x_m = a_\ell]} \right)^{z_{nk}}$

E-war: $\quad \mathbb{E}[z_{nk}] = \dfrac{\pi_k \prod_m \prod_\ell \theta_{km\ell}^{[x_m = a_\ell]}}{\sum_s \pi_s \prod_m \prod_\ell \theta_{sm\ell}^{[x_m = a_\ell]}}$

$\pi, \theta \to$ мат

M-war: $\mathbb{E}_z[\log p(X, z | \bar{\theta})] = \sum_n \sum_k \mathbb{E}[z_{nk}] \cdot \left( \log \pi_k + \sum_m \sum_\ell [x_m = a_\ell] \log \theta_{km\ell} \right)$

$= \underbrace{\sum_k \left( \sum_n \mathbb{E}[z_{nk}] \right) \log \pi_k}_{\bar{\pi} \, \text{max}} + \underbrace{\sum_k \sum_m \sum_\ell \left( \sum_n [x_m = a_\ell] \cdot \mathbb{E}[z_{nk}] \right) \cdot \log \theta_{km\ell}}_{\bar{\theta}_{km*} \, \text{max}}$

$\theta_{km\ell}^{(m+1)} = \dfrac{\sum_n [x_m = a_\ell] \cdot \mathbb{E}[z_{nk}]}{\boxed{\sum_{\ell'} \sum_n [x_m = a_{\ell'}] \, \mathbb{E}[z_{nk}]}} = \sum_n \mathbb{E}[z_{nk}]$

$\log a_1, \log a_2, \ldots, \log a_n$

$\log a_i' = \log \dfrac{e^{\log a_i}}{e^{\log a_1} + \ldots + e^{\log a_n}} = \log a_1 - \log \left( e^{\log a_1} + \ldots + e^{\log a_n} \right)$

logaddexp