② Энтропия и дивергенция (дискр. случай)

Пусть $P$ — распр, т.ч. символы $\{a_1 \dots a_k\}$ имеют вероятности $\{p_1 \dots p_k\}$

$$H(P) = -\sum_{j=1}^{k} p_j \log p_j \qquad - \text{энтропия}$$

**Утв.**  1) $H(P) \geqslant 0$ $\qquad = 0 \quad \Leftrightarrow \quad \exists_j : p_j = 1$ $\qquad (0 \log 0 = 0)$

2) $H(P) \leqslant \log k \qquad = \log k \Leftrightarrow \qquad \forall_j \quad p_j = 1/k$

▲ 1) следует из того, что $-\log p_j \geqslant 0$

2) $H(P) = E \log \dfrac{1}{p(\xi)} \leqslant \log E \dfrac{1}{p(\xi)} = \log \sum_{j=1}^{k} p_j \dfrac{1}{p_j} = \log k$

где $\xi \sim P$ $\qquad$ ↑ по нер-ву Йенсена $\qquad\qquad$ □

$P, Q$ — доминир.

$H(P) = -E \log p(\xi)$ , где $\xi \sim P$     энтропия

$H(P, Q) = -E \log q(\xi)$ , где $\xi \sim P$     кросс-энтропия

$KL(P, Q) = H(P, Q) - H(P) = E \log \dfrac{p(\xi)}{q(\xi)}$ , где $\xi \sim P$     дивергенция

Пример     $P = U(0, 1/2)$     $p(x) = 2 \, I\{x \in (0, 1/2)\}$

$H(P) = -\log 2 < 0$

Свойства KL:

1) $KL(P, Q) \geqslant 0$ $\quad\quad = 0 \iff P = Q$

▲ $-KL(P, Q) = E \log \dfrac{q(\xi)}{p(\xi)} \underset{\text{йенсен}}{\leqslant} \log E \dfrac{q(\xi)}{p(\xi)} = \log \int p(x) \dfrac{q(x)}{p(x)} dx =$

$$= \log 1 = 0$$

т.к. $\log$ — строго выпуклый, то равенство достигается при $P = Q$

2) $KL(P, Q) \neq KL(Q, P)$

3) Пусть $X_1 \dots X_n$ — выборка из $P_\theta \in \{P_\theta \mid \theta \in \Theta\}$ — дискретное

$KL(\hat{P}_n, P_\theta) = E_{\hat{P}_n} \log \dfrac{\hat{p}_n(X)}{P_\theta(X)} = \dfrac{1}{n} \sum\limits_{i=1}^{n} \log \dfrac{1/n}{P_\theta(X_i)} =$

дискр. пл-ти

$$= -\frac{1}{n} \sum_{i=1}^{n} \log p_\theta(x_i) - \log n = -\frac{1}{n} l_x(\theta) - \log n$$

Вывод: $\text{KL}(\hat{P}_n, P_\theta) \longrightarrow \min_{P_\theta} \quad \Longleftrightarrow \quad l_x(\theta) \longrightarrow \max_{\theta}$

**Теорема** Экстремальное свойство правдоподобия [L1 - L3]

$$\forall \theta_0, \theta_1 \in \Theta \quad \text{т.ч.} \quad \theta_0 \neq \theta_1 \qquad P_{\theta_0}\left(L_x(\theta_0) > L_x(\theta_1)\right) \longrightarrow 1 \quad \text{при} \quad n \to \infty$$

① $L_x(\theta_0) > L_x(\theta_1) \quad \Longleftrightarrow \quad \dfrac{L_x(\theta_0)}{L_x(\theta_1)} > 1 \quad \Longleftrightarrow \quad \dfrac{1}{n} \log \dfrac{L_x(\theta_0)}{L_x(\theta_1)} > 0$

$$\frac{1}{n} \log \frac{L_x(\theta_0)}{L_x(\theta_1)} = \frac{1}{n} \log \prod_{i=1}^{n} \frac{p_{\theta_0}(x_i)}{p_{\theta_1}(x_i)} = \frac{1}{n} \sum_{i=1}^{n} \log \frac{p_{\theta_0}(x_i)}{p_{\theta_1}(x_i)} \xrightarrow{\underset{364}{P_{\theta_0}-\text{п.н.}}}$$

$$\longrightarrow E_{\theta_0} \log \frac{p_{\theta_0}(x_i)}{p_{\theta_1}(x_i)} = KL\left(P_{\theta_0}, P_{\theta_1}\right) > 0$$

$\neq$ т.к. $\theta_0 \neq \theta_1$ и L1

Получаем $\quad P_{\theta_0}\left(L_x(\theta_0) > L_x(\theta_1)\right) \longrightarrow P\left(KL\left(P_{\theta_0}, P_{\theta_1}\right) > 0\right) = 1$

Следствие : Если $|\Theta| < \infty$, то ОМП существует и явл. сост. оц.

$$[L1 - L3]$$

$\triangle$    Пусть $\theta_0$ — истинное знач. парам

Тогда $\quad P_{\theta_0} \left( \forall \theta_1 \neq \theta_0 : L_x(\theta_0) > L_x(\theta_1) \right) =$

$$= P_{\theta_0} \left( \bigcap_{k=1}^{|\Theta| \; \leftarrow \; \text{конечное}} L_x(\theta_0) > L_x(\theta_k) \right) \longrightarrow 1$$

**Теорема**  С вер-тью $\to$ 1  ур-е $\dfrac{\partial L_x(\theta)}{\partial \theta} = 0$  имеет корень, являющийся сост. оц. $\theta$
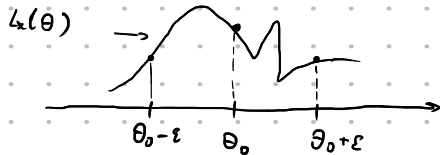
$$[L1 - L5]$$

- ▲  Пусть  $\theta_0$ — истинное  значение  параметра

  $L4 \implies \quad \exists \varepsilon > 0 \quad$ т.т.  $(\theta_0 - \varepsilon, \ \theta_0 + \varepsilon) \subset Ⓗ$

  По т. об экст. св-ве:

  $$P_{\theta_0}\left( L_x(\theta_0) > L_x(\theta_0 - \varepsilon), \ L_x(\theta_0) > L_x(\theta_0 + \varepsilon) \right) \to 1$$



  с вер $\to$ 1

$\Rightarrow$ с вер $\to 1$ на $(\theta_0 - \varepsilon, \theta_0 + \varepsilon)$ имеется корень ур-я правдоподобия.

Пусть $\tilde{\theta}_0$ — ближайший к $\theta_0$ корень ур-я (не факт, что $\tilde{\theta}_0 \in (\theta_0 - \varepsilon; \theta_0 + \varepsilon)$)
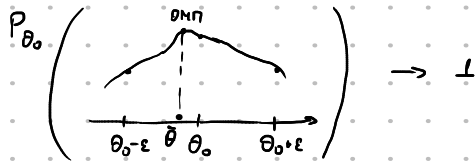
Тогда $P\left( |\tilde{\theta} - \theta_0| > \varepsilon \right) \to 0$

Т.к. это выполнено $\forall \varepsilon > 0$, то получаем, что

$$\tilde{\theta} \xrightarrow{P_{\theta_0}} \theta_0 \quad \Rightarrow \quad \tilde{\theta} - \text{сост. оц. } \theta \qquad \qquad \square$$

Следствие [L1 - L5] Если $\forall n$ $\forall X_1 \ldots X_n$ $\dfrac{\partial l_x(\theta)}{\partial \theta} = 0$ имеет единственный корень,

то $\tilde{\theta}$ — сост. оц. $\theta$

$$P(\tilde{\theta} = \hat{\theta}_{\text{ОМП}}) = 1 \implies \hat{\theta}_{\text{ОМП}} - \text{сост.}$$

▲



$$P_{\theta_0}\left( \right) \to 1$$

## Теорема [LL – L9]

1) $\tilde{\theta}$ – сост. оц. $\theta$, являющаяся решением ур-я правдоподобия

Тогда $\tilde{\theta}$ – а.н.о. $\theta$ с а.д. $i^{-1}(\theta)$

В частности, если $\tilde{\theta}$ – ед. реш ур-я правдоподобия,

Тогда $\tilde{\theta}$ – а.н.о. $\theta$ с а.д. $i^{-1}(\theta)$

2) Пусть $\hat{\theta}$ – а.н.о с а.д. $\sigma^2(\theta)$, т.ч. $\sigma^2(\theta)$ непр. по $\theta$

Тогда $\sigma^2(\theta) \geqslant i^{-1}(\theta)$ <span style="color:red">(д/g)</span>

Следствие [LL – L9] ОМП — сост., а.н.о, асимптотически эффективно
(т.е. имеет наим а.д. среди всех а.н.о. с непр. дисп.)

Примечание    если L1-L9 нету, то можно получить более крутые св-ва.

$$X_1 - X_n \sim U[0, \theta]$$

$$n(\theta - X_{(n)}) \xrightarrow{d_\theta} Exp(1) \quad \leftarrow \text{это крутче ас. норм}$$

$$\left(X_{(n)} \pm \frac{c}{n}\right) \qquad\qquad \sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N$$

$$\left(\hat{\theta} \pm \frac{c}{\sqrt{n}}\right)$$

Док-во   п. 1)

Вклад $\quad U_x(\theta) = \dfrac{\partial L_x(\theta)}{\partial \theta}$

разложим по Тейлору $\;U_x(\tilde\theta)$ в точке $\theta$:     (в форме Лагранжа)

$$U_x(\tilde\theta) = U_x(\theta) + U_x'(\theta)(\tilde\theta - \theta) + \frac{1}{2}U_x''(\theta^*)(\tilde\theta - \theta)^2$$

где $\theta^*$ лежит между $\tilde\theta$ и $\theta$

следствие: $\quad \theta^* \xrightarrow{P_\theta} \theta$ $\qquad$ [$\theta^*$ не сост. оц., т.к. не факт, что оц.,
$\qquad\qquad$ т.к. $\tilde\theta \xrightarrow{P_\theta} \theta$ $\qquad\qquad\qquad\qquad$ т.к. может зависеть от $\theta$]

$U_x(\tilde\theta) = 0$, т.к. $\tilde\theta$ — решение ур-я $\;U_x(\theta) = 0$

$\Rightarrow \;\; -U_x(\theta) = (\tilde\theta - \theta)\left(U_x'(\theta) + \frac{1}{2}U_x''(\theta^*)(\tilde\theta - \theta)\right)$

$$\frac{-\sqrt{n} \ \ U_x(\theta)}{U_x'(\theta) + \frac{1}{2} U_x''(\theta^*)(\tilde\theta - \theta)} = \sqrt{n}(\tilde\theta - \theta)$$

$$\frac{\overset{=}{-\sqrt{n} \ \ U_x(\theta) \ \ 1/n} \qquad \to N(0, i(\theta))}{\underset{\to i(\theta)}{\frac{1}{n} U_x'(\theta)} + \underset{\to \infty}{\frac{1}{2n} U_x''(\theta^*)}(\underset{\to 0}{\tilde\theta - \theta})} \qquad \to N\left(0, i^{-1}(\theta)\right)$$

1) $\quad \sqrt{n} \ \frac{1}{n} \ U_x(\theta) = \sqrt{n}\left(\frac{1}{n} \sum_{i=1}^{n} U_{x_i}(\theta) - \overset{0}{\overbrace{E_\theta \ U_{x_i}(\theta)}}\right) \overset{d_\theta}{\underset{ЦПТ}{\longrightarrow}} N\left(0, D_\theta \overset{i(\theta)}{\overbrace{U_{x_i}''(\theta)}}\right)$

2) $\quad \frac{1}{n} U_x'(\theta) = \frac{1}{n} \sum_{i=1}^{n} U_{x_i}'(\theta) \overset{P_\theta}{\underset{ЗБЧ}{\longrightarrow}} E_\theta \ U_{x_i}'(\theta) = -i(\theta)$

3) $\quad \hat\theta - \theta \overset{P_\theta}{\longrightarrow} 0 \quad \Rightarrow \quad хотим \ \left|\frac{1}{n} U_x''(\theta^*)\right| \leq C$

$$\left| \frac{1}{n} U_x^n (\theta^*) \right| \leq \frac{1}{n} \sum_{i=1}^{n} \left| U_{x_i}^n (\theta^*) \right| \leq \frac{1}{n} \sum_{i=1}^{n} H(x_i) \xrightarrow[\text{убзд}]{P_\theta} E_\theta H(x_i) < +\infty$$

$$\text{для} \quad \theta^* \in (\theta - c; \theta + c) \quad \text{по} \quad \text{ЗБЧ}$$

$$\text{т.к.} \quad \theta^* \xrightarrow{P_\theta} \theta, \quad \text{то} \quad P_\theta(\theta^* \in (\theta \pm c)) \longrightarrow 1$$

$$\Rightarrow \quad \frac{1}{n} U_x^n (\theta) \quad \text{огр.} \quad \text{с} \quad \text{вер.} \longrightarrow 1$$

$$\text{т.е.} \quad \exists c(\theta) \quad \forall \theta : \quad P_\theta \left( \left| \frac{1}{n} U_x^n (\theta) \right| \leq c(\theta) \right) \longrightarrow 1$$

2) без док-ва

# Натуральный градиент

**Пример**   $X_1 \ldots X_n \sim N(a, S)$   // $S = \sigma^2$

$$L_x(\theta) = -\frac{n}{2} \log 2\pi - \frac{n}{2} \log S - \frac{1}{2S} \sum_{i=1}^{n} \left( x_i - a \right)^2$$

$$\frac{\partial L_x(\theta)}{\partial a} = -\frac{1}{S} \sum_{i=1}^{n} \left( a - x_i \right) = \frac{n}{S} \left( \overline{x} - a \right)$$

$$\frac{\partial L_x(\theta)}{\partial S} = -\frac{n}{2S} + \frac{1}{2S^2} \sum_{i=1}^{n} \left( x_i - a \right)^2 = \frac{n}{2S^2} \left( \frac{1}{n} \sum_{i=1}^{n} (x_i - a)^2 - S \right)$$

Градиентный спуск:   $a_{t+1} = a_t + \eta_1 \dfrac{\overline{X} - a}{S_t}$

← ответ для $a$

$$s_{t+1} = S_t + \eta_2 \frac{1}{S_t^2} \left( \frac{1}{n} \sum_{i=1}^{n} (x_i - a_t)^2 - S_t \right)$$

← ответ для $S$

что-то ухудшающее сходимость →

Пусть $f: \mathbb{R}^d \to \mathbb{R}$

$$\nabla f(x) = \left(\frac{\partial f(x)}{\partial x_j}\right)_j \qquad \frac{\partial f(x)}{\partial x_j} = \lim_{\Delta x_j \to 0} \frac{f(x + e_j \Delta x_j) - f(x)}{\Delta x_j}$$

$$\nabla f(x) \quad \rightsquigarrow \quad \operatorname*{argmax}_{\substack{\Delta x_j \\ \|\Delta x\| \leq \varepsilon}} \left(f(x + \Delta x) - f(x)\right)$$

<span style="color:red">направление наискорейшего роста</span>

Пусть $P_\theta \in \{P_\theta \mid \theta \in \textcircled{H}\}$

$\qquad f(P_\theta)$ — функционал, напр $L_x(\theta)$
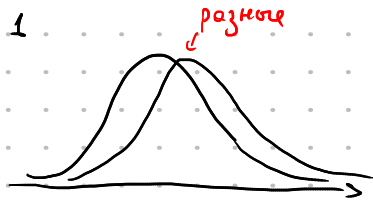
Проблема взятия производной:

$\qquad$ окрестность $\|\Delta\theta\| \leq \varepsilon$ не отражает скместь распределений

Пример <span style="color:red">(на проблему выше)</span>   1)   $N(\theta_1, \theta_2)$

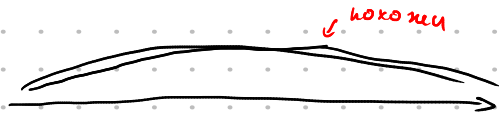$\theta = (0, 1)$   $\Delta\theta = (1, 0)$   $\|\Delta\theta\| = 1$

$P_\theta = N(0, 1)$   $P_{\theta + \Delta\theta} = N(1, 1)$



← разные

2)   $N(\theta_1, \theta_2)$   $\theta = (0, 100)$   $\Delta\theta = (1, 0)$   $\|\Delta\theta\| = 1$

$P_\theta = N(0, 100)$

$P_{\theta + \Delta\theta} = N(1, 100)$



← похожи

Идея: в опр. $grad$ заменим $\|\Delta\theta\| < \varepsilon$ на окр-ть в пр-ве распр

Опр.   Натуральный   градиент

$$\nabla_N f(P_\theta) \sim \underset{\substack{\Delta\theta_j \\ KL(P_\theta, P_{\theta+\Delta\theta}) < \varepsilon}}{argmax} \left[ f(P_{\theta+\Delta\theta}) - f(P_\theta) \right]$$

Распишем   Лагранжиан:

$$L = f(P_{\theta+\Delta\theta}) - f(P_\theta) - \lambda \left( KL(P_\theta, P_{\theta+\Delta\theta}) - \varepsilon \right)$$

$$\nabla_{\Delta\theta} L \Big|_{\Delta\theta=0} = \underbrace{\nabla_{\Delta\theta} f(P_{\theta+\Delta\theta}) \Big|_{\Delta\theta=0}}_{\substack{\nabla_\theta f(P_\theta) \\ \text{т.е. обычный} \\ \text{градиент}}} - \lambda \nabla_{\Delta\theta} KL(P_\theta, P_{\theta+\Delta\theta}) \Big|_{\Delta\theta=0}$$

**Утв.** При $\|\Delta\theta\| \to 0$

$$KL(P_\theta, P_{\theta+\Delta\theta}) = \frac{1}{2}\Delta\theta^T \underbrace{i(\theta)}_{\text{инф. матр. Ришера}} \Delta\theta + o(\|\Delta\theta\|^2)$$

▲ Разложим $KL$ по Тейлору в $\Delta\theta = 0$

$$KL(P_\theta, P_{\theta+\Delta\theta}) = \underbrace{KL(P_\theta, P_\theta)}_{\overset{\shortparallel}{0}} + \underbrace{\nabla_{\Delta\theta}KL(P_\theta, P_{\theta+\Delta\theta})\Big|_{\Delta\theta=0}^T}_{\overset{\shortparallel}{0}, \text{ т.к. это произв. в точке минимума}} \cdot \Delta\theta +$$

$$+ \frac{1}{2}\Delta\theta^T \nabla^2_{\Delta\theta}KL(P_\theta, P_{\theta+\Delta\theta})\Big|_{\Delta\theta=0} \cdot \Delta\theta + o(\|\Delta\theta\|^2)$$

$$\nabla^2_{\Delta\theta}KL(P_\theta, P_{\theta+\Delta\theta})\Big|_{\Delta\theta=0} = \int P_\theta(x)\, \nabla^2_{\Delta\theta}\log\frac{p_\theta(x)}{p_{\theta+\Delta\theta}(x)}\Big|_{\Delta\theta=0} dx =$$

$$= -\int p_\theta(x)\, \underbrace{\nabla^2_{\Delta\theta}\log p_{\theta+\Delta\theta}(x)\Big|_{\Delta\theta=0}}_{\nabla^2_\theta \log p_\theta(x)} dx = -\mathbb{E}_\theta\, \nabla^2_\theta\log p_\theta(x_1) = i(\theta)$$

Вернёмся к Лагранжиану

$$\nabla_{\Delta\theta} L = \nabla_\theta f(P_\theta) - \lambda \, i(\theta) \Delta\theta = 0$$

$$\Rightarrow \quad \Delta\theta = \frac{1}{\lambda} i^{-1}(\theta) \nabla_\theta f(P_\theta)$$

на это можно забить, т.к. натуральный градиент используется в град. спуске $\Rightarrow$ там есть learning-rate

$$\theta_{np} \quad \nabla_N f(\rho_\theta) = i^{-1}(\theta) \nabla_\theta f(P_\theta) \quad - \text{натуральный градиент}$$

Применим к примеру с нормальной выборкой:

Из семинара $i(\theta) = \begin{pmatrix} 1/s & 0 \\ 0 & 1/2s^2 \end{pmatrix} \Rightarrow i^{-1}(\theta) = \begin{pmatrix} s & 0 \\ 0 & 2s^2 \end{pmatrix}$

Расписываем градиентный спуск:

$$a_{t+1} = a_t + \eta_1 (\bar{x} - a_t) = \eta_1 \bar{x} + (1-\eta_1) a_t$$

$$S_{t+1} = S_t + \eta_2 \left( \frac{1}{n} \sum_{i=1}^{n} (x_i - a_t)^2 - S_t \right) = \eta_2 \left( \frac{1}{n} \sum_{i=1}^{n} (x_i - a_t)^2 \right) + (1-\eta_2) S_t$$

# Логистическая регрессия

$$\nabla l_y(\theta) = X^T(y - S(\theta)) \qquad I(\theta) = X^T V(\theta) X$$

$$\theta_{t+1} = \theta_t + \eta \left(X^T V(\theta_t) X\right)^{-1} X^T (y - S(\theta_t))$$

## IRLS

метод 1ого порядка в пр-ве распределений

метод 2ого порядка в пр-ве признаков