

My first scientific paper

Week 6

Write the theory:
axioms \rightarrow **theorems**,
methods \rightarrow **algorithms**

Vadim Strijov

Moscow Institute of Physics and Technology

2022

Как в моей дипломной работе найти теорему?

Животные делятся на:

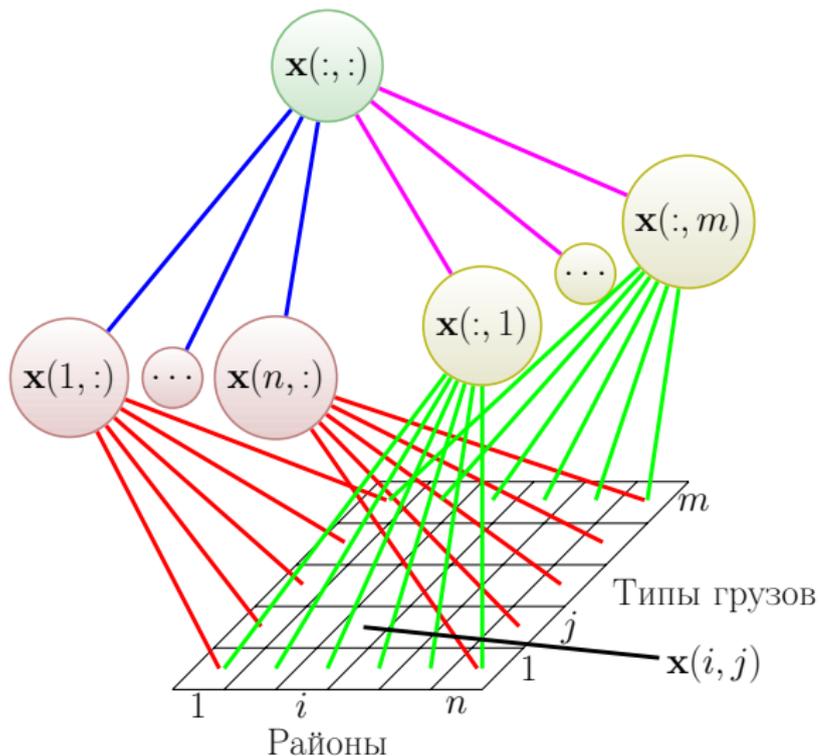
- a) принадлежащих Императору,
- b) набальзамированных,
- c) прирученных,
- d) молочных поросят,
- e) сирен,
- f) сказочных,
- g) бродячих собак,
- h) включённых в эту классификацию,
- i) бегающих как сумасшедшие,
- j) бесчисленных,
- k) нарисованных тончайшей кистью из верблюжьей шерсти,
- l) прочих,
- m) разбивших цветочную вазу,
- n) похожих издали на мух.

¹ Китайская энциклопедия «Божественное хранилище благотворных знаний»

Что доказываем?

- a) существование,
- b) единственность,
- c) эквивалентность множеств или отображений, и даже решений,
- d) сходимость,
- e) корректность по Адамару,
- f) сложность,
- g) включая сложность $\mathcal{O}(n \ln n)$,
- h) устойчивость не только по Ляпунову,
- i) корректность статистического вывода,
- j) и байесовского вывода,
- k) состоятельность, несмещенность, эффективность,
- l) оптимальность (в некотором смысле),
- m) равновесность,
- n) и даже сходимость.

Условие согласованности прогнозов



$$x_t(:, :) = \sum_{i=1}^n x_t(i, :);$$

$$x_t(:, :) = \sum_{j=1}^m x_t(:, j);$$

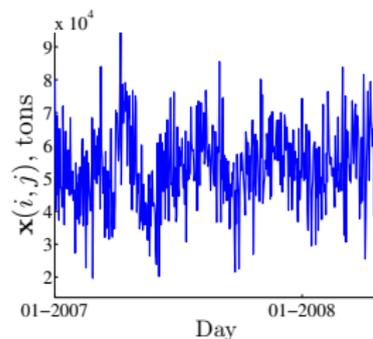
$$x_t(i, :) = \sum_{j=1}^m x_t(i, j),$$

$$i = 1, \dots, n;$$

$$x_t(:, j) = \sum_{i=1}^n x_t(i, j),$$

$$j = 1, \dots, m;$$

$$t = 1, \dots, T.$$



Дано

Матрица связей \mathbf{S} , множества \mathcal{A} , \mathcal{B} и вектор независимых прогнозов $\hat{\chi}$

$$\hat{\chi} \notin \mathcal{A}, \quad \hat{\chi} \in \mathcal{B}.$$

Требуется построить

вектор согласованных прогнозов $\hat{\varphi}$, который удовлетворяет следующим требованиям:

- $\hat{\varphi} \in \mathcal{A}$, $\mathcal{A} = \{\chi \in \mathbb{R}^d \mid \mathbf{S}\chi = \mathbf{0}\}$ — согласованность;
- $\hat{\varphi} \in \mathcal{B}$ — физические ограничения;
- $l_h(\chi_{T+1}, \hat{\varphi}) \leq l_h(\chi_{T+1}, \hat{\chi})$ для любого среза действительных значений $\chi_{T+1} \in \mathcal{A} \cap \mathcal{B}$ — качество.

Согласование как антагонистическая игра

Игрок, выбирающий вектор согласованных прогнозов $\hat{\varphi}$, играет с природой, выбирающей срез иерархии в момент времени $(T + 1)$. Цель игрока — минимизировать свои потери при любом ходе природы.

	Стратегия	Потери
Игрок	$\hat{\varphi} \in \mathcal{A} \cap \mathcal{B}$	$L(\hat{\varphi}, \chi_{T+1}) = I_h(\chi_{T+1}, \hat{\varphi}) - I_h(\chi_{T+1}, \hat{\chi})$
Природа	$\chi_{T+1} \in \mathcal{A} \cap \mathcal{B}$	$-L(\hat{\varphi}, \chi_{T+1})$

Равновесие Нэша в антагонистической игре — это

пара стратегий $(\hat{\varphi}, \chi_{T+1})$, таких что для любых стратегий $\hat{\varphi}'$, χ'_{T+1} выполнено неравенство

$$L(\hat{\varphi}, \chi'_{T+1}) \leq L(\hat{\varphi}, \chi_{T+1}) \leq L(\hat{\varphi}', \chi_{T+1}).$$

Цена игры (Дж. Нэш)

$$V = \min_{\hat{\varphi}} \max_{\chi_{T+1}} L(\hat{\varphi}, \chi_{T+1}) = \max_{\chi_{T+1}} \min_{\hat{\varphi}} L(\hat{\varphi}, \chi_{T+1})$$

определена тогда и только тогда, когда в игре существует равновесие Нэша.

Существование равновесия Нэша и выбор согласованных прогнозов

Теорема 1 (Стенина, 2014)

Пусть $\mathcal{A} \cap \mathcal{B} \neq \emptyset$ и для функции потерь l_h выполнено

- 1 $l_h(\chi_{T+1}, \hat{\chi}) \geq 0$ для произвольных векторов χ_{T+1} , $\hat{\chi}$,
причем $l_h(\chi_{T+1}, \hat{\chi}) = 0 \Leftrightarrow \chi_{T+1} = \hat{\chi}$;
- 2 существует проекция $\chi_{proj} = \arg \min_{\chi \in \mathcal{A} \cap \mathcal{B}} l_h(\chi, \hat{\chi})$;
- 3 для всех $\chi \in \mathcal{B}$ и для всех $\psi \in \mathcal{A} \cap \mathcal{B}$ выполняется
неравенство $l_h(\psi, \chi) \geq l_h(\psi, \chi_{proj}) + l_h(\chi_{proj}, \chi)$.

Тогда

- пара стратегий $(\chi_{proj}, \chi_{proj})$ является равновесием Нэша в антагонистической игре, описывающей задачу согласования прогнозов;
- пара $(\chi_{proj}, \chi_{proj})$ является седловой точкой функции $L(\hat{\varphi}, \chi_{T+1}) = l_h(\chi_{T+1}, \hat{\varphi}) - l_h(\chi_{T+1}, \hat{\chi})$.

Теорема 2: цена игры (Стенина, 2014)

При выполнении требований теоремы 1 цена игры определена и равна

$$V = \min_{\hat{\varphi}} \max_{\chi_{T+1}} L(\hat{\varphi}, \chi_{T+1}) = \max_{\chi_{T+1}} \min_{\hat{\varphi}} L(\hat{\varphi}, \chi_{T+1}) = -l_h(\chi_{proj}, \hat{\chi}) \leq 0.$$

Теорема 3: согласованные прогнозы (Стенина, 2014)

При выполнении требований теоремы 1 использование в качестве вектора согласованных прогнозов $\hat{\varphi}$ вектора

$$\hat{\varphi} = \chi_{proj} = \arg \min_{\chi \in A \cap B} l_h(\chi, \hat{\chi})$$

гарантирует, что вектор согласованных прогнозов будет удовлетворять требованиям согласованности и качества и физическим ограничениям.

Задача согласования прогнозов сводится к решению оптимизационной задачи.

Удовлетворяют условию теоремы 1

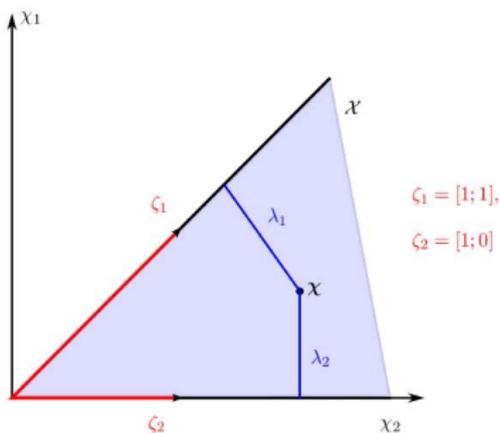
- $l_h(\mathbf{u}, \mathbf{v}) = \|\mathbf{u} - \mathbf{v}\|^2$ — квадрат евклидова расстояния,
- $l_h(\mathbf{u}, \mathbf{v}) = \frac{1}{2}(\mathbf{u} - \mathbf{v})^T Q(\mathbf{u} - \mathbf{v})$ — квадрат расстояния Махаланобиуса,
- $l_h(\mathbf{u}, \mathbf{v}) = \sum_{i=1}^d u_i \log \frac{u_i}{v_i} - \sum_{i=1}^d u_i + \sum_{i=1}^d v_i$ — обобщенная дивергенция Кульбака-Лейблера,
- $l_h(\mathbf{u}, \mathbf{v}) = \sum_{i=1}^d \left(\frac{u_i}{v_i} - \log \frac{u_i}{v_i} - 1 \right)$ — расстояние Itakura-Saito,
- $l_h(\mathbf{u}, \mathbf{v}) = \sum_{i=1}^d |u_i|^a - \sum_{i=1}^d |v_i|^a - a \sum_{i=1}^d \text{sign}(v_i) |v_i|^{a-1} (u_i - v_i)$,
 $a > 1$,
- $l_h(\mathbf{u}, \mathbf{v}) = \frac{2}{a^2} \sum_{i=1}^d (e^{au_i} - e^{av_i}) - \frac{2}{a} \sum_{i=1}^d e^{av_i} (u_i - v_i)$, $a \neq 0$.

Порождающее представление конуса

Порождающее представление конуса

Полиэдральный конус \mathcal{X} допускает представление через конечный набор порождающих элементов ζ_1, \dots, ζ_k :

$$\mathcal{X} = \left\{ \sum_{k=1}^r \lambda_k \zeta_k \mid \lambda_k \geq 0 \right\}.$$



Теорема (о порождающем представлении конуса),
[Кузнецов: 2013]

Столбцы матрицы предпочтений $\mathbf{Z}(i, k) = \mathbb{I}[x_i \succcurlyeq x_k]$ являются порождающими элементами конуса предпочтений,

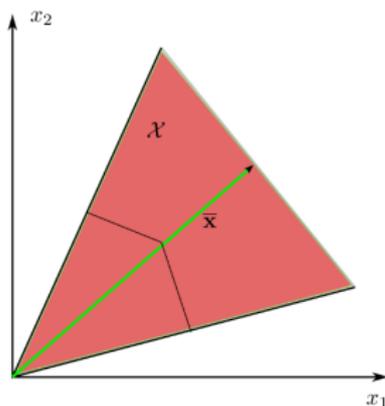
$$\mathcal{X} \supset \{\mathbf{Z}\boldsymbol{\lambda} \mid \boldsymbol{\lambda} \in \mathbb{R}_+^m\}.$$

Регуляризация конусной модели

Линейная конусная модель:

$$\mathbf{f}(X) = \sum_{j=1}^n \mathbf{z}_j \lambda_j, \quad \lambda_j \geq 0.$$

Рассмотрим в конусе \mathcal{X} центральную точку $\bar{\mathbf{x}} = \frac{1}{n} \sum_{j=1}^n \mathbf{z}_j$.



Теорема (о регуляризации конусной модели),
[Кузнецов: 2014]

В случае замены каждого конуса $\mathcal{X}_k = \{\sum \lambda_{jk} \mathbf{z}_{jk} \mid \lambda_k \geq \mathbf{0}\}$ его центральной точкой конусная модель представима в виде

$$\mathbf{f}(\mathbf{x}_1, \dots, \mathbf{x}_m) = \hat{\mathbf{Z}} \boldsymbol{\lambda}, \quad \hat{\mathbf{Z}} = \sum_{j=1}^n w_j \mathbf{z}_j,$$

при ограничениях $w_j \geq 0$, $\sum \lambda_k = 1$, $\boldsymbol{\lambda} \geq \mathbf{0}$.

Приближенная оптимальность релаксированной низкоранговой задачи

Теорема (Мотренко)

Пусть $\underline{\mathbf{A}}^R = \sum_{r=1}^R \mathbf{a}_1^{(r)} \circ \mathbf{a}_2^{(r)} \circ \mathbf{a}_3^{(r)}$ — построенное низкоранговое приближение матрицы $\underline{\mathbf{A}}$. Отклонение функции $F(\underline{\mathbf{A}})$ от $F(\underline{\mathbf{A}}^R)$ линейно зависит от невязки $\|\underline{\mathbf{A}} - \underline{\mathbf{A}}^R\|$:

$$\left| F(\underline{\mathbf{A}}) - F(\underline{\mathbf{A}}^R) \right| = \mathcal{O}(\|\underline{\mathbf{A}} - \underline{\mathbf{A}}^R\|).$$

Теорема (Мотренко)

Пусть $\hat{\underline{\mathbf{A}}} \in [0, 1]^{n_1 \times n_2 \times n_3}$ — некоторая трехиндексная матрица, $\underline{\mathbf{A}}^\varepsilon \in 0, 1^{n_1 \times n_2 \times n_3}$ получена бинаризацией $\hat{\underline{\mathbf{A}}}$:

$$a_{ijk}^\varepsilon = [a_{ijk} > \varepsilon].$$

Отклонение функционала $F(\hat{\underline{\mathbf{A}}})$ от $F(\underline{\mathbf{A}}^\varepsilon)$ линейно зависит от ε :

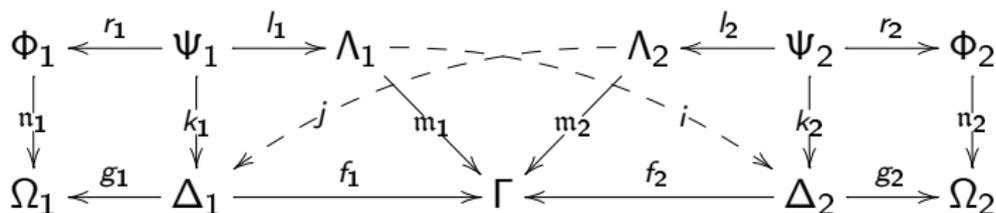
$$\left| F(\hat{\underline{\mathbf{A}}}) - F(\underline{\mathbf{A}}^\varepsilon) \right| = \mathcal{O}(n_1 n_2 n_3 \varepsilon).$$

Theorem

Две трансформации графов $\Gamma \xrightarrow{p_1, m_1} \Omega$ и $\Omega \xrightarrow{p_2, m_2} \Gamma'$ являются последовательно независимыми, если существуют морфизмы $i : \Phi_1 \rightarrow \Delta_2$ и $j : \Lambda_2 \rightarrow \Delta_1$, такие что $f_2 \circ i = n_1$ и $g_1 \circ j = m_2$:

$$\begin{array}{ccccccc}
 \Lambda_1 & \xleftarrow{l_1} & \Psi_1 & \xrightarrow{r_1} & \Phi_1 & \xrightarrow{\quad} & \Lambda_2 & \xleftarrow{l_2} & \Psi_2 & \xrightarrow{r_2} & \Phi_2 & & (3) \\
 \downarrow n_1 & & \downarrow k_1 & & \downarrow n_1 & \text{---} & \downarrow m_2 & & \downarrow k_2 & & \downarrow n_2 & & \\
 \Gamma & \xleftarrow{f_1} & \Delta_1 & \xrightarrow{g_1} & \Omega & \xleftarrow{\quad} & \Delta_2 & \xrightarrow{g_2} & \Gamma' & & & &
 \end{array}$$

The diagram illustrates the commutative relationships between various graph transformations. Solid arrows represent the transformations: $\Lambda_1 \xleftarrow{l_1} \Psi_1 \xrightarrow{r_1} \Phi_1$, $\Lambda_2 \xleftarrow{l_2} \Psi_2 \xrightarrow{r_2} \Phi_2$, $\Gamma \xleftarrow{f_1} \Delta_1 \xrightarrow{g_1} \Omega$, and $\Delta_2 \xrightarrow{g_2} \Gamma'$. Vertical arrows represent morphisms: $\Lambda_1 \downarrow n_1 \Gamma$, $\Psi_1 \downarrow k_1 \Delta_1$, $\Phi_1 \downarrow n_1 \Omega$, $\Lambda_2 \downarrow m_2 \Delta_2$, $\Psi_2 \downarrow k_2 \Delta_2$, and $\Phi_2 \downarrow n_2 \Gamma'$. Dashed arrows represent the morphisms $i : \Phi_1 \rightarrow \Delta_2$ and $j : \Lambda_2 \rightarrow \Delta_1$. The conditions $f_2 \circ i = n_1$ and $g_1 \circ j = m_2$ are indicated by the dashed arrows meeting at the Ω node.



Рассмотрим варианты принадлежности $m_1(v)$.

- 1 Множество $m_1(v) \notin m_2(\Lambda_2)$. Все вершины графа Γ являются образами при применении отображений m_2 или f_2 . Отсюда $m_1(v) \in f_2(\Delta_2)$.
- 2 Множество $m_1(v) \in m_2(\Lambda_2)$. Тогда $m_1(v) \in m_1(\Lambda_1) \cap m_2(\Lambda_2) \subseteq m_1(l_1(\Psi_1)) \cap m_2(l_2(\Psi_2))$. При этом из коммутативной диаграммы следует, что $m_2(l_2(\Psi_2)) = f_2(k_2(\Psi_2))$. Отсюда $m_1(v) \in f_2(\Delta_2)$.

В обоих случаях $m_1(x) \in f_2(\Delta_2)$, из инъективности f_2 : $i(x) = f_2^{-1} \circ m_1(x)$. Аналогично, j из условия $f_1 \circ j = m_2$.

Теорема [Кузьмин, 2016]

Апостериорное распределение параметров q из класса $q(\theta)q(\mathbf{m}, \mathbf{V}, \alpha)$, заданное оптимальными оценками факторов $q^*(\theta)$ и $q^*(\mathbf{m}, \mathbf{V}, \alpha)$, в которых используется верхняя оценка правдоподобия имеет вид

$$q = \mathcal{N}(\alpha_0, a^{-1}\mathbf{I}) \prod_{k=1}^{k_h} \mathcal{N}(\mathbf{m}'_{0k}, (\nu'\mathbf{V}_k)^{-1}) \mathcal{N}(\mathbf{m}_{0k}, (b'\mathbf{V}_k)^{-1}) \mathcal{W}(\mathbf{W}_k, \nu'),$$

$$\mathbf{m}'_{0k} = \mathbf{m}_{0k} + \frac{1}{\nu'} (\mathbf{W}_k^{-1})^T \mathbf{M}_k^T E_{\alpha}[\boldsymbol{\Lambda}] \sum_{n=1}^N \mathbf{x}_n \left(z_{nk} - \frac{\exp(\xi_{nk})}{g(\boldsymbol{\xi}_n)} \right),$$

$$\alpha_0 = \frac{1}{a} \sum_{n=1}^N \sum_{m=1}^{|\mathcal{W}|} x_{nm} \iota_m \sum_{k=1}^{k_h} (\mathbf{M}_k E[\boldsymbol{\theta}_k])_m \left(z_{nk} - \frac{\exp(\xi_{nk})}{g(\boldsymbol{\xi}_n)} \right),$$

$$\mathbf{W}_k^{-1} = (b+1)\mathbf{m}_{0k}\mathbf{m}_{0k}^T + b\mathbf{m}_0\mathbf{m}_0^T + E[\boldsymbol{\theta}_k\boldsymbol{\theta}_k^T] + \mathbf{W}^{-1},$$

$$\mathbf{m}_{0k} = \frac{E[\boldsymbol{\theta}_k] + b\mathbf{m}_0}{1+b}, \quad b' = b+1, \quad \nu' = \nu+1.$$

Операторы релевантности:

R_{MAP} ранжирует кластеры по точечной оценке $p(\tilde{z}_{tk} | \theta_k^{\text{MAP}}, \alpha^{\text{MAP}})$,

$$\theta_k^{\text{MAP}}, \alpha^{\text{MAP}} = \arg \max_{\theta, \alpha} q(\theta, \alpha),$$

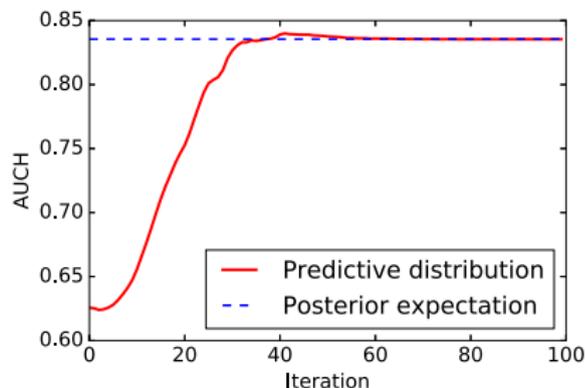
$\hat{R}_{\text{evidence}}$ ранжирует кластеры по оценке обоснованности

$$p(\tilde{z}_{kt} | \tilde{x}_t) = \int p(\tilde{z}_{kt} | \theta, \alpha)_k q(\theta, \alpha) d\theta d\alpha \approx \hat{p}(\tilde{x}_t, \tilde{\xi}_t, \psi_{tk}) \rightarrow \max_{\psi_{tk}} \min_{\tilde{\xi}_t},$$

$\tilde{\xi}_t, \psi_{tk}$ – вариационные параметры оценки.

Теорема [Кузьмин, 2016]

При оптимальных значениях вариационных параметров $\tilde{\xi}_t$ и ψ_t , значение критерия качества $\text{AUCH}(R)$ для операторов R_{MAP} и $\hat{R}_{\text{evidence}}$ совпадают.



Корректная функция сходства s должна быть

- 1 определена в случае несовпадения носителей,
- 2 $s(g_1, g_2) \leq s(g_1, g_1)$,
- 3 $s \in [0, 1]$,
- 4 $s(g_1, g_1) = 1$,
- 5 близка к 1, если $g_2(w)$ — малоинформативное распределение,
- 6 симметрична, $s(g_1, g_2) = s(g_2, g_1)$.

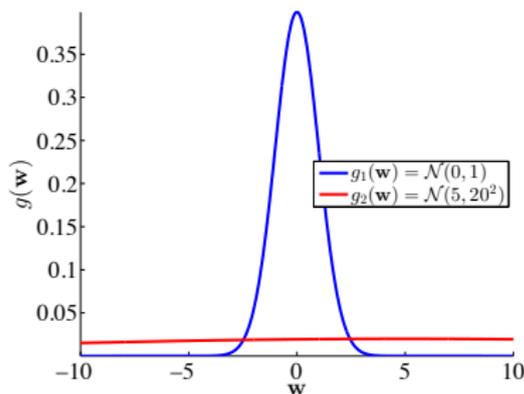
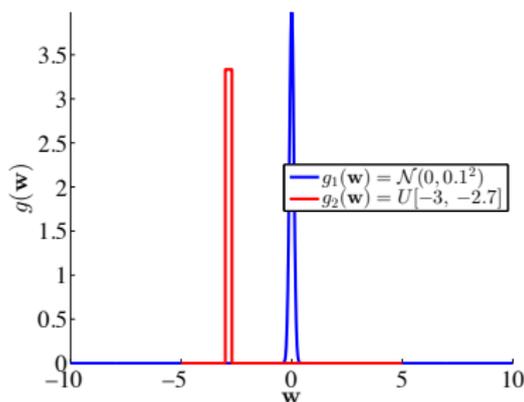
Теорема 3 (Адуенко, 2014)

Функции сходства, порожденные расстояниями Кульбака-Лейблера, Дженсона-Шеннона, Хеллингера, Бхаттачарая, не являются корректными.

Иллюстрация требований к функции сходства

Важно, чтобы значение функции s

было близко к 1, если $g_2(\mathbf{w})$ — малоинформативное распределение.



Теорема 4 (Адуенко, 2014)

Функции сходства, порожденные дивергенциями Брегмана, симметризованными дивергенциями Брегмана и f-дивергенциями, не являются корректными.

В качестве меры сходства распределения предлагается мера сходства s -score:

$$s(g_1, g_2) = \frac{\int_{\mathbf{w}} g_1(\mathbf{w})g_2(\mathbf{w})d\mathbf{w}}{\max_{\mathbf{b}} \int_{\mathbf{w}} g_1(\mathbf{w} - \mathbf{b})g_2(\mathbf{w})d\mathbf{w}}.$$

Теорема 5 (Адуенко, 2014). Предлагаемая функция сходства является корректной.

Примеры:

$g_1(\mathbf{w})$	$g_2(\mathbf{w})$	$s(g_1, g_2)$
$U[0, 1]$	$U[0.5, 1.5]$	0.5
$U[0, 1]$	$U[0, 1]$	1
$\mathcal{N}(0, 1)$	$\mathcal{N}(10, 10^{10})$	1

Выражение для $s(g_1, g_2)$ для пары нормальных распределений

Определение 6. *Обобщенно-линейной моделью с натуральной функцией связи и априорным распределением на вектор параметров $p(\mathbf{w}|\mathbf{A})$ называется вероятностная модель с совместным правдоподобием*

$$p(\mathbf{y}, \mathbf{w}|\mathbf{X}, \mathbf{A}) = p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w}|\mathbf{A}), \text{ где } p(\mathbf{y}|\mathbf{X}, \mathbf{w}) = \prod_{i=1}^m p(y_i|\mathbf{x}_i, \mathbf{w}),$$
$$p(y_i|\mathbf{x}_i, \mathbf{w}) = c(y_i) \exp(\theta_i y_i - b(\theta_i)), \text{ где } \theta_i = \mathbf{w}^\top \mathbf{x}_i.$$

Теорема 6 (Адуенко, 2014).

Пусть $g_1 = \mathcal{N}(\mathbf{v}_1, \Sigma_1)$, $g_2 = \mathcal{N}(\mathbf{v}_2, \Sigma_2)$. Тогда выражение для $s(g_1, g_2)$ имеет вид

$$s(g_1, g_2) = \exp\left(-\frac{1}{2}(\mathbf{v}_1 - \mathbf{v}_2)^\top (\Sigma_1 + \Sigma_2)^{-1}(\mathbf{v}_1 - \mathbf{v}_2)\right).$$

Следствие 1. В случае $\Sigma_2 = \mathbf{0}$ выражение для s-score

$$s(g_1, g_2) = \exp\left(-\frac{1}{2}(\mathbf{v}_2 - \mathbf{v}_1)^\top \Sigma_1^{-1}(\mathbf{v}_2 - \mathbf{v}_1)\right).$$

Распределение s-score в условии истинности гипотезы о совпадении моделей

Рассматриваем далее пару обобщенно-линейных моделей с натуральной функцией связи. Обозначим $\mathbf{H}_m(\mathbf{w}) = \mathbf{H}_m$ и введем $O_m^\delta(\mathbf{w}) = \{\mathbf{v} : \|\mathbf{H}_m^{T/2}(\mathbf{v} - \mathbf{w})\| \leq \delta\}$.

Теорема 7 (Адуенко, 2016). Пусть

- Модели f_1 и f_2 совпадают: $\mathbf{w}_1 = \mathbf{w}_2 = \mathbf{w}$, \mathbf{w}_2 известно;
- Априорное распределение: $\mathbf{w} \sim \mathcal{N}(\mathbf{w}|\mathbf{0}, \mathbf{A}_{m^k}^{-1})$, $k = 1, 2$;
- $\exists c_0^k < \infty : \|\mathbf{A}_{m^k}^k\| < c_0^k \forall m^k, k = 1, 2$;
- $\sum_{i=1}^{m^1} \mathbf{x}_i \mathbf{x}_i^\top$ имеет полный ранг для $m^1 \geq m_0$;
- $\lambda_{\min}(\mathbf{H}_{m^1}) \rightarrow \infty$ при $m^1 \rightarrow \infty$;
- $\forall \delta > 0 \max_{\mathbf{v} \in O_{m^1}^\delta(\mathbf{w})} \|\mathbf{H}_{m^1}^{-\frac{1}{2}} \mathbf{H}_{m^1}(\mathbf{v}) \mathbf{H}_{m^1}^{-\frac{1}{2}} - \mathbf{I}\| \rightarrow 0$ при $m^1 \rightarrow \infty$.

Тогда $-2 \log s\text{-score} = (\hat{\mathbf{w}} - \mathbf{w})^\top \tilde{\mathbf{H}}_{m^1}(\hat{\mathbf{w}})(\hat{\mathbf{w}} - \mathbf{w}) \xrightarrow{d} \chi^2(n)$.

Следствие 1. Для случая $n = 2$ s-score имеет асимптотически равномерное распределение на отрезке $[0, 1]$.

Обобщающая задача

Задачу выбора модели \mathbf{h}^*, θ^* назовем обобщающей на множестве

$U_\theta \times U_h \times U_\lambda \subset \mathbb{R}^u \times \mathbb{H} \times \Lambda$, если выполнены условия:

- 1 Для каждого $\mathbf{h} \in U_h$ и каждого $\lambda \in U_\lambda$ решение θ^* определено однозначно.
- 2 **Условие максимизации правдоподобия выборки:** существует $\lambda \in U_\lambda$ и $K_1 \in \mathbb{R}_+$, такие что для любых векторов гиперпараметров $\mathbf{h}_1, \mathbf{h}_2 \in U_h$, $Q(\mathbf{h}_1) - Q(\mathbf{h}_2) > K_1$: матожидания правдоподобия выборок: $E_q \log p(\mathbf{y}|\mathbf{X}, \theta_1, \lambda_{\text{temp}}, \mathbf{f}) > \log E_q p(\mathbf{y}|\mathbf{X}, \theta_2, \lambda_{\text{temp}}, \mathbf{f})$.
- 3 **Условие минимизации сложности модели:** существует $\lambda \in U_\lambda$ и $K_2 \in \mathbb{R}_+$, такие что для любых векторов гиперпараметров $\mathbf{h}_1, \mathbf{h}_2 \in U_h$, $Q(\mathbf{h}_1) - Q(\mathbf{h}_2) > K_2$, $E_q \log p(\mathbf{y}|\theta_1, \lambda_{\text{temp}}, \mathbf{f}) = \log E_q p(\mathbf{y}|\theta_2, \lambda_{\text{temp}}, \mathbf{f})$, количество ненулевых параметров у первой модели меньше, чем у второй.
- 4 **Условие максимизации обоснованности модели:** существует значение гиперпараметров λ , такое что оптимизация задачи эквивалента оптимизации вариационной оценки обоснованности модели:
$$\mathbf{h}^* = \arg \max p(\mathbf{y}|\mathbf{X}, \mathbf{h}', \lambda_{\text{temp}}, \mathbf{f}), \quad \theta^* = \arg \min D_{\text{KL}}(q|p(\mathbf{w}, \Gamma|\mathbf{y}, \mathbf{X}, \lambda_{\text{temp}}, \mathbf{f})).$$
- 5 **Условие перехода между структурами:** Существует константа K_3 , такая что для любых двух векторов $\mathbf{h}_1, \mathbf{h}_2$ и соответствующих векторов $\theta_1^*, \theta_2^* : D_{\text{KL}}(q_{\Gamma_2}, q_{\Gamma_1}) > K_3, D_{\text{KL}}(q_{\Gamma_1}, q_{\Gamma_2}) > K_3$: существуют значения гиперпараметров λ_1, λ_2 , такие что $Q(\mathbf{h}_1, \lambda_1) > Q(\mathbf{h}_2, \lambda_1), Q(\mathbf{h}_1, \lambda_1) < Q(\mathbf{h}_2, \lambda_2)$.
- 6 **Условие непрерывности:** \mathbf{h}^*, θ^* непрерывны по метапараметрам.

Анализ задач выбора моделей

Теорема [Бахтеев, 2019]

Следующие задачи выбора модели не являются обобщающими:

- ① метод максимума правдоподобия: $\max_{\theta} E_q \log p(y|\mathbf{X}, \theta, \lambda_{\text{temp}}, \mathbf{f})$;
- ② метод максимума апостериорной вероятности $\max_{\theta} E_q \log p(y|\mathbf{X}, \theta, \mathbf{f}) p(\theta|\mathbf{h}, \lambda_{\text{temp}}) p(\mathbf{h}|\mathbf{f})$;
- ③ метод максимума вариационной оценки обоснованности модели $\max_{\mathbf{h}} \max_{\theta} E_q \log p(y|\mathbf{X}, \mathbf{w}, \Gamma, \mathbf{f}) - D_{KL}(q(\mathbf{w}, \Gamma) || p(\mathbf{w}, \Gamma, \lambda_{\text{temp}})) p(\mathbf{h}|\mathbf{f})$;
- ④ кросс-валидация $\max_{\mathbf{h}} E_q \log p(y_{\text{valid}}|\mathbf{X}_{\text{valid}}, \theta^*, \lambda_{\text{temp}}, \mathbf{f}) p(\mathbf{h}|\mathbf{f})$,
 $\theta^* = \arg \max_{\theta} E_q \log p(y_{\text{train}}|\mathbf{X}_{\text{train}}, \theta, \lambda_{\text{temp}}, \mathbf{f}) p(\theta|\mathbf{h})$.
- ⑤ AIC: $\max_{\theta} E_q \log p(y|\mathbf{X}, \theta, \lambda_{\text{temp}}, \mathbf{f}) - |\theta_i : \theta_i \neq 0|$;
- ⑥ BIC: $\max_{\theta} E_q \log p(y|\mathbf{X}, \theta, \lambda_{\text{temp}}, \mathbf{f}) - \frac{1}{2} \log(m) |\theta_i : \theta_i \neq 0|$;
- ⑦ перебор структуры модели:
 $\max \Gamma' \max_{\theta} E_q \log p(y|\mathbf{X}, \theta, \lambda_{\text{temp}}, \mathbf{f}) \mathbb{I}(\Gamma = \Gamma')$.

Предлагаемая задача оптимизации

Теорема [Бахтеев, 2018]

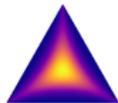
Пусть функции потерь и валидации L, Q являются непрерывно-дифференцируемыми на некоторой области U . Тогда следующая задача является обобщающей на U .

$$\begin{aligned} \mathbf{h}^* &= \arg \max_{\mathbf{h}} Q = & (Q^*) \\ &= \lambda_{\text{likelihood}}^Q E_{q^*} \log p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \Gamma, \mathbf{h}, \lambda_{\text{temp}}, \mathbf{f}) - \\ &\quad - \lambda_{\text{prior}}^Q D_{\text{KL}}(q^*(\mathbf{w}, \Gamma) || p(\mathbf{w}, \Gamma | \mathbf{h}, \lambda_{\text{temp}}, \mathbf{f})) - \\ &\quad - \sum_{p' \in \mathcal{P}, \lambda \in \lambda_{\text{struct}}^Q} \lambda D_{\text{KL}}(\Gamma | p') + \log p(\mathbf{h} | \mathbf{f}), \end{aligned}$$

где

$$\begin{aligned} q^* &= \arg \max_q L = E_q \log p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \Gamma, \mathbf{h}, \lambda_{\text{temp}}, \mathbf{f}) & (L^*) \\ &\quad - \lambda_{\text{prior}}^Q D_{\text{KL}}(q^*(\mathbf{w}, \Gamma) || p(\mathbf{w}, \Gamma | \mathbf{h}, \lambda_{\text{temp}}, \mathbf{f})). \end{aligned}$$

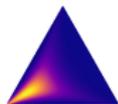
Оптимизационная задача обобщает алгоритмы оптимизации: оптимизация правдоподобия и обоснованности, последовательное увеличение и снижение сложности модели, полный перебор структуры.



$$\lambda_{\text{struct}}^Q = [0; 0; 0].$$



$$\lambda_{\text{struct}}^Q = [1; 0; 0].$$



$$\lambda_{\text{struct}}^Q = [1; 1; 0].$$

Адекватность задачи оптимизации

Теорема, [Бахтеев, 2018]

Пусть задано параметрическое множество вариационных распределений: $q(\theta)$.

Пусть $\lambda_{\text{likelihood}}^L = \lambda_{\text{prior}}^L = \lambda_{\text{prior}}^Q > 0$, $\lambda_{\text{struct}}^Q = 0$. Тогда:

- 1 Задача оптимизации (Q^*) доставляет максимум апостериорной вероятности гиперпараметров с использованием вариационной оценки обоснованности:

$$\log \hat{p}(\mathbf{y}|\mathbf{X}, \mathbf{h}, \lambda_{\text{temp}}, \mathbf{f}) + \log p(\mathbf{h}|\mathbf{f}) \rightarrow \max_{\mathbf{h}}.$$

- 2 Вариационное распределение q приближает апостериорное распределение $p(\mathbf{w}, \Gamma|\mathbf{y}, \mathbf{X}, \mathbf{h}, \lambda_{\text{temp}}, \mathbf{f})$ наилучшим образом:

$$D_{\text{KL}}(q||p(\mathbf{w}, \Gamma|\mathbf{y}, \mathbf{X}, \mathbf{h}, \lambda_{\text{temp}}, \mathbf{f})) \rightarrow \min_{\theta}.$$

Пусть также распределение q декомпозируется на два независимых распределения для параметров \mathbf{w} и структуры Γ модели \mathbf{f} :

$$q = q_{\mathbf{w}}q_{\Gamma}, q_{\Gamma} \approx p(\Gamma|\mathbf{y}, \mathbf{X}, \mathbf{h}, \mathbf{f}), q_{\mathbf{w}} \approx p(\mathbf{w}|\Gamma, \mathbf{y}, \mathbf{X}, \mathbf{h}, \mathbf{f}).$$

Тогда вариационные распределения $q_{\mathbf{w}}$, q_{Γ} приближают апостериорные распределения $p(\Gamma|\mathbf{y}, \mathbf{X}, \mathbf{h}, \lambda_{\text{temp}}, \mathbf{f})$, $p(\mathbf{w}|\Gamma, \mathbf{y}, \mathbf{X}, \mathbf{h}, \lambda_{\text{temp}}, \mathbf{f})$ наилучшим образом:

$$D_{\text{KL}}(q_{\Gamma}||p(\Gamma|\mathbf{y}, \mathbf{X}, \mathbf{h}, \lambda_{\text{temp}}, \mathbf{f})) \rightarrow \min, \quad D_{\text{KL}}(q_{\mathbf{w}}||p(\mathbf{w}|\mathbf{y}, \mathbf{X}, \mathbf{h}, \mathbf{f})) \rightarrow \min.$$

Анализ обобщающей задачи оптимизации

Теорема, [Бахтеев, 2018]

Пусть $\lambda_{\text{prior}}^L > 0$, $m \gg 0$, $\frac{m}{\lambda_{\text{prior}}^L} \in \mathbb{N}$. Тогда оптимизация функции

$$L = \mathbb{E}_q \log p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \Gamma, \mathbf{h}, \lambda_{\text{temp}}, \mathbf{f}) - \lambda_{\text{prior}}^L D_{\text{KL}}(q || p(\mathbf{w}, \Gamma | \mathbf{h}, \lambda_{\text{temp}}, \mathbf{f})))$$

эквивалентна минимизации $\mathbb{E}_{\hat{\mathbf{X}}, \hat{\mathbf{y}} \sim p(\mathbf{x}, \mathbf{y})} D_{\text{KL}}(q || p(\mathbf{w}, \Gamma | \hat{\mathbf{X}}, \hat{\mathbf{y}}, \mathbf{h}, \lambda_{\text{temp}}, \mathbf{f}))$, где $\hat{\mathbf{X}}, \hat{\mathbf{y}}$ — случайные подвыборки мощностью $\frac{m}{\lambda_{\text{prior}}^L}$ из генеральной совокупности.

Определение

Параметрической сложностью модели назовем минимальную дивергенцию между априорным и вариационным распределением:

$$C_p = \min_{\mathbf{h}} D_{\text{KL}}(q || p(\mathbf{w}, \Gamma | \mathbf{h}, \lambda_{\text{temp}}, \mathbf{f})).$$

Теорема, [Бахтеев, 2018]

При устремлении параметрической сложности модели к нулю относительная плотность параметров модели стремится к единице:

$$C_p \rightarrow 0 \implies \rho(\mathbf{w}) \rightarrow 1, \quad \rho(w) = \frac{q(0)}{q(w)} = \exp\left(-\frac{\mu^2}{\sigma^2}\right).$$

Оптимизация параметрической сложности

Теорема, [Бахтеев, 2018]

Пусть $\lambda_{\text{likelihood}}^Q = \lambda_{\text{prior}}^L > 0$, $\lambda_{\text{struct}}^Q = \mathbf{0}$. Тогда предел оптимизации

$$\lim_{\lambda_{\text{prior}}^Q \rightarrow \infty} \lim_{\eta \rightarrow \infty} T^\eta(Q, \mathbf{h}, T^\eta(L, \theta_0, \mathbf{h}))$$

доставляет минимум параметрической сложности. Существует компактная область U , такая что для любой точки $\theta_0 \in U$ предел данной оптимизации доставляет нулевую параметрическую сложность: $C_p = 0$.

Теорема, [Бахтеев, 2018]

Пусть $\lambda_{\text{likelihood}}^L = 1$, $\lambda_{\text{struct}}^Q = \mathbf{0}$. Пусть $\mathbf{f}_1, \mathbf{f}_2$ — результаты градиентной оптимизации при разных значениях гиперпараметров

$\lambda_{\text{prior}}^{Q,1}, \lambda_{\text{prior}}^{Q,2}, \lambda_{\text{prior}}^{Q,1} < \lambda_{\text{prior}}^{Q,2}$, полученных при начальном значении вариационных параметров θ_0 и гиперпараметров \mathbf{h}_0 . Пусть θ_0, \mathbf{h}_0 принадлежат области U , в которой соответствующие функции L и Q являются локально-выпуклыми. Тогда:

$$C_p(\mathbf{f}_1) - C_p(\mathbf{f}_2) \geq \lambda_{\text{prior}}^L (\lambda_{\text{prior}}^L - \lambda_{\text{prior}}^{Q,1}) \sup_{\theta, \mathbf{h} \in U} |\nabla_{\theta, \mathbf{h}}^2 D_{\text{KL}}(q|p) (\nabla_{\theta}^2 L)^{-1} \nabla_{\theta} D_{\text{KL}}(q|p)|.$$

