

# Приближённый полиномиальный алгоритм для одной задачи разбиения последовательности

Кельманов А. В., Хамидуллин С. А.

*Институт математики им. С. Л. Соболева СО РАН,  
Новосибирск*

10-я Международная конференция  
«Интеллектуализация обработки информации» (ИОИ-2014)

о. Крит, 4–11 октября, 2014 г.

## Предмет исследования —

одна из труднорешаемых задач дискретной оптимизации.

## Цель исследования —

обоснование приближенного полиномиального алгоритма для её решения.

## Мотивация исследования —

отсутствие каких-либо эффективных алгоритмов с гарантированными оценками точности для решения этой задачи.

# Задача 1-MSSC-S-NF разбиения последовательности (Minimum Sum-of-Squares Clustering, for the case of Sequence, NonFixed cardinalities)

## Задача 1-MSSC-S-NF

**Дано:** последовательность  $\mathcal{Y} = (y_1, \dots, y_N)$  векторов из  $\mathbb{R}^q$ ,  
натуральные числа  $T_{\min}$  и  $T_{\max}$ .

**Найти:** набор  $\mathcal{M} = \{n_1, \dots, n_M\} \subseteq \mathcal{N}$  номеров элементов  
последовательности  $\mathcal{Y}$  такой, что

$$F(\mathcal{M}) = \sum_{j \in \mathcal{M}} \|y_j - \bar{y}(\mathcal{M})\|^2 + \sum_{i \in \mathcal{N} \setminus \mathcal{M}} \|y_i\|^2 \rightarrow \min, \quad (1)$$

где  $\bar{y}(\mathcal{M}) = \frac{1}{|\mathcal{M}|} \sum_{i \in \mathcal{M}} y_i$ , при ограничениях

$$1 \leq T_{\min} \leq n_m - n_{m-1} \leq T_{\max} \leq N, \quad m = 2, \dots, M, \quad (2)$$

на элементы искомого набора  $\mathcal{M}$ .

## Имеется

таблица, содержащая упорядоченные по времени результаты многократных измерений набора числовых информационно значимых характеристик некоторого объекта, который может находиться в пассивном и активном состоянии.

## Предполагается, что:

- (1) в пассивном состоянии все числовые характеристики из набора равны нулю, а в любом активном — значение хотя бы одной характеристики не равно нулю;
- (2) в каждом результате измерения, представленном в таблице, имеется ошибка;
- (3) соответствие элементов таблицы какому-либо состоянию объекта неизвестно;
- (4) временной интервал между двумя последовательными активными состояниями объекта ограничен сверху и снизу некоторыми константами.

## Требуется:

- 1 разбить таблицу на подмножества наборов, соответствующих пассивному и активному состояниям объекта, используя критерий минимума суммы квадратов расстояний,
- 2 оценить по результатам измерения наборы характеристик объекта в активном состоянии (учитывая, что данные содержат ошибку измерения).

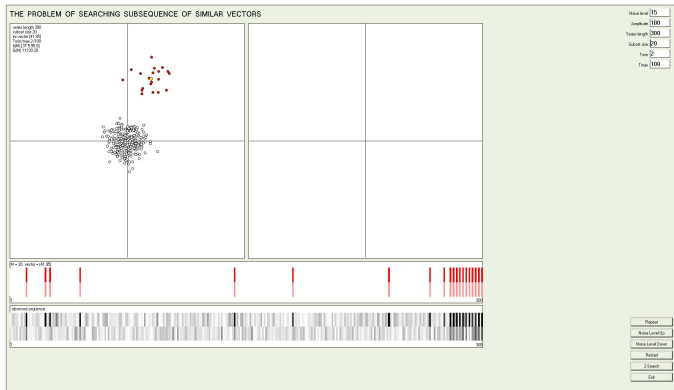
## Замечание

Критерий минимума суммы квадратов расстояний обусловлен оптимизационной моделью помехоустойчивого анализа данных.

## Пример. Задача 1-MSSC-S-NF

300 результатов измерений характеристик объекта, изображенные на плоскости и в виде последовательности.

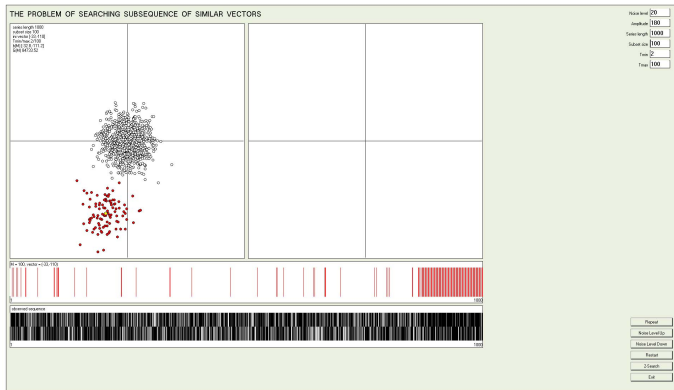
20 раз были измерены характеристики объекта в активном состоянии и 280 — в пассивном.



## Пример. Задача 1-MSSC-S-NF

1000 результатов измерений характеристик объекта, изображенные на плоскости и в виде последовательности.

100 раз были измерены характеристики объекта в активном состоянии и 900 — в пассивном.





## Известные результаты

1. Задача NP-трудна в сильном смысле. Поэтому для этой задачи не существует ни точного полиномиального, ни точного псевдополиномиального алгоритмов, если  $P \neq NP$  (Кельманов, Пяткин, 2009).
2. Параметрический вариант — задача 1-MSSC-S-NF( $T_{\min}, T_{\max}$ ) (Кельманов, Пяткин, 2013)
  - NP-трудна в сильном смысле, когда  $T_{\min} < T_{\max}$ ;
  - разрешима за полиномиальное время при  $T_{\min} = T_{\max}$ .
3. Какие-либо полиномиальные алгоритмы с оценками точности до настоящего времени отсутствовали.

## Новый результат

Обоснован 2-приближенный полиномиальный алгоритм (Кельманов, Хамидуллин, 2014).

## Суть подхода

1. Заменяем решение исходной задачи 1-MSSC-S-NF решением более простой (вспомогательной) задачи.
2. Построим точный полиномиальный алгоритм ее решения.
3. Оценим точность такой замены (приближения).

Подход опирается на записи целевой функции задачи 1-MSSC-S-NF в виде

$$\begin{aligned} F(\mathcal{M}) &= \sum_{j \in \mathcal{M}} \|y_j - \bar{y}(\mathcal{M})\|^2 + \sum_{i \in \mathcal{N} \setminus \mathcal{M}} \|y_i\|^2 \\ &= \sum_{j \in \mathcal{N}} \|y_j\|^2 - \sum_{i \in \mathcal{M}} (2(y_i, \bar{y}(\mathcal{M})) - \|\bar{y}(\mathcal{M})\|^2). \end{aligned}$$

## Задача 1

**Дано:** последовательность  $\mathcal{Y} = (y_1, \dots, y_N)$  векторов из  $\mathbb{R}^q$ , вектор  $b \in \mathbb{R}^q$ , натуральные числа  $T_{\min}$  и  $T_{\max}$ .

**Найти:** набор  $\mathcal{M} = \{n_1, \dots, n_M\} \subseteq \mathcal{N}$  номеров элементов  $\mathcal{Y}$  такой, что

$$G(\mathcal{M}) = \sum_{n \in \mathcal{M}} g(n) \rightarrow \max, \quad (3)$$

где

$$g(n) = 2(y_n, b) - \|b\|^2, \quad n \in \mathcal{N}, \quad (4)$$

при ограничениях (2) на элементы набора  $\mathcal{M}$ .

## Лемма 1

Пусть  $M \leq \lfloor (N - 1) / T_{\min} \rfloor + 1$ . Тогда оптимальное значение  $G_{\max}$  целевой функции задачи 1 находится по формуле

$$G_{\max} = \max_{n \in \mathcal{N}} G_n, \quad (5)$$

а значения функции  $G_n$  вычисляются по следующим рекуррентным формулам

$$G_n = \begin{cases} g(n), & \text{если } n \in \{1, \dots, T_{\min}\}; \\ g(n) + \max\{0, \max_{j \in \gamma(n)} G_j\}, & \text{если } n \in \{T_{\min} + 1, \dots, N\}, \end{cases} \quad (6)$$

где множество  $\gamma$  задается формулой

$$\gamma(n) = \{j \mid \max\{1, n - T_{\max}\} \leq j \leq n - T_{\min}\}.$$

Формулы леммы реализуют схему динамического программирования.

Для отыскания оптимального набора определим функцию  $I(n) : \mathcal{N} \rightarrow \{0\} \cup \mathcal{N}$  по следующей формуле

$$I(n) = \begin{cases} 0, & \text{если } n \in \{1, \dots, T_{\min}\}, \\ 0, & \text{если } n \in \{T_{\min} + 1, \dots, N\} \text{ и } \max_{j \in \gamma(n)} G_j < 0, \\ \arg \max_{j \in \gamma(n)} G_j, & \text{если } n \in \{T_{\min} + 1, \dots, N\} \text{ и } \max_{j \in \gamma(n)} G_j \geq 0. \end{cases} \quad (7)$$

Формулы для отыскания оптимального набора

$$\widehat{\mathcal{M}} = (\widehat{n}_1, \dots, \widehat{n}_{\widehat{M}})$$

устанавливает

## Следствие 1

Пусть выполнены условия леммы 1,  $G_n$  и  $I(n)$ ,  $n \in \mathcal{N}$ , — функции, заданные формулами (6) и (7), соответственно. Тогда последняя компонента оптимального набора вычисляется по формуле

$$\hat{n}_{\hat{M}} = \arg \max_{n \in \mathcal{N}} G_n. \quad (8)$$

Если  $I(\hat{n}_{\hat{M}}) > 0$ , то оставшиеся компоненты этого набора находятся по следующим рекуррентным формулам

$$\hat{n}_{m-1} = I(\hat{n}_m), \quad m = \hat{M}, \hat{M} - 1, \dots \quad (9)$$

Если на некотором шаге  $t$  вычислений по формуле (9) имеет место равенство  $I(\hat{n}_m) = 0$ , то  $m = 1$  и найдены все компоненты оптимального набора.

### Алгоритм $\mathcal{A}_1$

Входами алгоритма являются  $\mathcal{Y}$ ,  $b$ ,  $T_{\min}$  и  $T_{\max}$ .

**Шаг 1.** Вычислим значения  $g(n)$ ,  $n \in \mathcal{N}$ , по формуле (4).

**Шаг 2.** Используя формулы (6) и (7), вычислим значения  $G_n$  и  $l(n)$  для каждого  $n \in \mathcal{N}$ . Найдем значение  $G_{\max}$  максимума целевой функции  $G$  по формуле (5).

**Шаг 3.** Вычислим элементы оптимального набора  $\widehat{\mathcal{M}} = \{\widehat{n}_1, \dots, \widehat{n}_{\widehat{M}}\}$  по формулам (8) и (9); выход.

Выходом алгоритма являются значения, зависящие от  $b$ , т.е.  $G_{\max} = G_{\max}(b)$ ,  $\widehat{\mathcal{M}} = \widehat{\mathcal{M}}(b) = \{\widehat{n}_1(b), \dots, \widehat{n}_{\widehat{M}}(b)\}$  и  $\widehat{M} = \widehat{M}(b)$ .



### Теорема 1

Алгоритм  $A_1$  находит оптимальное решение задачи 1 за время  $\mathcal{O}(N(T_{\max} - T_{\min} + q))$ .

**Доказательство.** Оптимальность решения следует из леммы 1.

На **первом шаге** алгоритма требуется  $\mathcal{O}(Nq)$  операций.

Трудоёмкость **второго шага**, выполняемого для каждого  $n = 1, \dots, N$ , определяется мощностью множеств  $\mathcal{N}$  и  $\gamma(n)$ , входящих в определение функции (6). Поэтому трудоёмкость шага 2 есть величина  $\mathcal{O}(N(T_{\max} - T_{\min} + 1))$ .

Из формул (5), (8) и (9) видно, что на **шаге 3** требуется не более  $\mathcal{O}(M_{\max})$  операций, что не превышает  $\mathcal{O}(N)$ , так как  $M_{\max} \leq N$ .

Просуммировав затраты на всех шагах, получим оценку временной сложности, приведенную в формулировке теоремы.

## Замечание 1

*Из ограничений (2) следует двойное неравенство*

$0 \leq T_{\max} - T_{\min} \leq N - 1$ . Поэтому алгоритм  $A_1$  полиномиален по  $N$  и по  $q$ , а его сложность в общем случае можно оценить как  $\mathcal{O}(N(N + q))$  и как  $\mathcal{O}(Nq)$  в частном случае, когда  $T_{\max} = T_{\min}$ .

### Алгоритм $\mathcal{A}$

Входами алгоритма являются  $\mathcal{Y}$ ,  $T_{\min}$ ,  $T_{\max}$  и  $M$ .

**Шаг 1.** Положим  $i = 0$ ,  $M_A = 0$ ,  $\mathcal{M}_A = \emptyset$ ,  $H = -\infty$ .

**Шаг 2.**  $i := i + 1$ ; положим  $b = y_i$ .

**Шаг 3.** Для фиксированного вектора  $b \in \mathcal{Y}$  найдем оптимальное решение  $\widehat{M}(b)$  и значение  $G_{\max}(b)$  целевой функции задачи 1 с помощью алгоритма  $\mathcal{A}_1$ .

**Шаг 4.** Если  $H < G_{\max}(b)$ , то положим  $b_A = b$ ,  $H = G_{\max}(b)$ ,  
 $M_A = \widehat{M}(b)$ ,  $\mathcal{M}_A = \widehat{M}(b)$ .

**Шаг 5.** Если  $i < N$ , то переходим на Шаг 2, иначе — к следующему шагу.

**Шаг 6.** Вычислим вектор  $\bar{y}(\mathcal{M}_A) = (1/M_A) \sum_{n \in \mathcal{M}_A} y_n$  и значение  $F(\mathcal{M}_A)$  целевой функции по формуле (1); положим  $G_{\max}(b_A) = H$ ; выход.

Выходом алгоритма (решением задачи) объявляем набор  $\mathcal{M}_A$ , значение  $F(\mathcal{M}_A)$ , а также векторы  $\bar{y}(\mathcal{M}_A)$  и  $b_A$ . Если максимуму  $G_{\max}(b_A)$  соответствует несколько наборов  $\widehat{M}_A$ , то выбираем любой из них.

## Алгоритм $\mathcal{A}$

Суть алгоритма  $\mathcal{A}$  состоит в решении задачи 1 с помощью алгоритма  $\mathcal{A}_1$  для каждого вектора последовательности  $\mathcal{U}$  и последующего выбора из найденных решений (наборов) наилучшего набора  $\widehat{\mathcal{M}}_{\mathcal{A}}$ , которому соответствует наибольшее значение  $G_{\max}(b_{\mathcal{A}})$  целевой функции задачи 1.

Приведем вспомогательное утверждение.

## Лемма 2

Пусть  $\mathcal{Z}$  — непустое конечное множество векторов из  $\mathbb{R}^q$ , а  $\bar{z}(\mathcal{Z}) = \frac{1}{|\mathcal{Z}|} \sum_{z \in \mathcal{Z}} z$ . Тогда, если вектор  $x \in \mathbb{R}^q$  удовлетворяет условиям

$$\|x - \bar{z}\| \leq \|z - \bar{z}\|, \quad \forall z \in \mathcal{Z},$$

то имеет место неравенство

$$\sum_{z \in \mathcal{Z}} \|z - x\|^2 \leq 2 \sum_{z \in \mathcal{Z}} \|z - \bar{z}\|^2.$$

## Теорема 2

Алгоритм  $\mathcal{A}$  находит 2-приближенное решение задачи 1-MSSC-S-NF за время  $\mathcal{O}(N^2(T_{\max} - T_{\min} + q))$ . Оценка 2 точности алгоритма достижима.

## Доказательство теоремы 2

Пусть  $\mathcal{M}^*$  — оптимальное решение задачи 1-MSSC-S-NF,

$\mathcal{C}^* = \{y_n | n \in \mathcal{M}^*\}$  — подмножество векторов из  $\mathcal{Y}$ , соответствующих оптимальному набору  $\mathcal{M}^*$ ,

$\bar{y}(\mathcal{M}^*) = \frac{1}{|\mathcal{M}^*|} \sum_{n \in \mathcal{M}^*} y_n$  — центр подмножества  $\mathcal{C}^* \subseteq \mathcal{Y}$ ,

$u = \arg \min_{y \in \mathcal{C}^*} \|y - \bar{y}(\mathcal{M}^*)\|$  — вектор, ближайший к центру подмножества  $\mathcal{C}^*$ .

# Задача 1-MSSC-S-NF. Алгоритм решения

Оценки точности и трудоёмкости

Согласно пошаговой записи алгоритм находит вектор

$$b_A = \arg \max_{b \in \mathcal{Y}} G_{\max}(b) \quad (10)$$

из множества  $\mathcal{Y}$ , набор  $\mathcal{M}_A = \widehat{\mathcal{M}}(b_A) = \{\hat{n}_1(b_A), \dots, \hat{n}_{\widehat{M}}(b_A)\}$ , вектор  $\bar{y}(\mathcal{M}_A)$ , значение  $F(\mathcal{M}_A)$  целевой функции задачи 1-MSSC-S-NF, а также максимум

$$G_{\max}(b_A) = \sum_{n \in \mathcal{M}_A} \{2(y_n, b_A) - \|b_A\|^2\} = \max_{b \in \mathcal{Y}} \max_{\mathcal{M}} \sum_{n \in \mathcal{M}} \{2(y_n, b) - \|b\|^2\} \quad (11)$$

функции  $G_{\max}(b)$ ,  $b \in \mathcal{Y}$ .

# Задача 1-MSSC-S-NF. Алгоритм решения

Оценки точности и трудоёмкости

Справедливость утверждения теоремы вытекает из следующей цепочки равенств и неравенств

$$\begin{aligned} F(\mathcal{M}_A) &= \sum_{n \in \mathcal{M}_A} \|y_n - \bar{y}(\mathcal{M}_A)\|^2 + \sum_{n \in \mathcal{N} \setminus \mathcal{M}_A} \|y_n\|^2 \\ &\stackrel{(1)}{\leq} \sum_{n \in \mathcal{M}_A} \|y_n - b_A\|^2 + \sum_{n \in \mathcal{N} \setminus \mathcal{M}_A} \|y_n\|^2 \\ &= \sum_{n \in \mathcal{N}} \|y_n\|^2 - \sum_{n \in \mathcal{M}_A} \{2(y_n, b_A) - \|b_A\|^2\} = \end{aligned}$$

Неравенство 1 справедливо, так как для любого конечного множества  $\mathcal{Z}$  векторов из  $\mathbb{R}^q$  (в частности, для  $\mathcal{Z} = \{y_n | n \in \mathcal{M}_A\}$ ) минимум суммы квадратов  $\sum_{z \in \mathcal{Z}} \|z - c\|^2$  по  $c$  достигается в точке  $c = \bar{z}(\mathcal{Z}) = \frac{1}{|\mathcal{Z}|} \sum_{z \in \mathcal{Z}} z$  (т.е. в точке  $c = \bar{y}(\mathcal{M}_A)$ ).



# Задача 1-MSSC-S-NF. Алгоритм решения

Оценки точности и трудоёмкости

$$\begin{aligned} &=_{(2)} \sum_{n \in \mathcal{N}} \|y_n\|^2 - \max_{b \in \mathcal{Y}} \max_{\mathcal{M}} \sum_{n \in \mathcal{M}} \{2(y_n, b) - \|b\|^2\} \\ &=_{(3)} \min_{b \in \mathcal{Y}} \min_{\mathcal{M}} \left( \sum_{n \in \mathcal{N}} \|y_n\|^2 - \sum_{n \in \mathcal{M}} \{2(y_n, b) - \|b\|^2\} \right) \\ &= \min_{b \in \mathcal{Y}} \min_{\mathcal{M}} \left( \sum_{n \in \mathcal{M}} \|y_n - b\|^2 + \sum_{n \in \mathcal{N} \setminus \mathcal{M}} \|y_n\|^2 \right) \leq \end{aligned}$$

Равенство 2 следует из теоремы 1 и формул (10), (11).

Равенство 3 справедливо, т.к.  $\sum_{n \in \mathcal{N}} \|y_n\|^2$  — константа.

# Задача 1-MSSC-S-NF. Алгоритм решения

Оценки точности и трудоёмкости

$$\begin{aligned} &\leq_{(4)} \sum_{n \in \mathcal{M}^*} \|y_n - u\|^2 + \sum_{n \in \mathcal{N} \setminus \mathcal{M}^*} \|y_n\|^2 \\ &\quad \leq_{(5)} 2 \sum_{n \in \mathcal{M}^*} \|y_n - \bar{y}(\mathcal{M}^*)\|^2 + \sum_{n \in \mathcal{N} \setminus \mathcal{M}^*} \|y_n\|^2 \\ &\leq 2 \left( \sum_{n \in \mathcal{M}^*} \|y_n - \bar{y}(\mathcal{M}^*)\|^2 + \sum_{n \in \mathcal{N} \setminus \mathcal{M}^*} \|y_n\|^2 \right) = 2F(\mathcal{M}^*). \quad (12) \end{aligned}$$

Неравенство 4 следует из того, что подмножество  $\mathcal{M}^*$  и вектор  $u \in \mathcal{Y}$  — допустимое решение задачи  $\min_{b \in \mathcal{Y}} \min_{\mathcal{M}}(\cdot)$ .

Справедливость неравенства 5 вытекает из леммы 2.

# Задача 1-MSSC-S-NF. Алгоритм решения

Оценки точности и трудоёмкости

Таким образом, из (12) следует, что  $\frac{F(\mathcal{M}_A)}{F(\mathcal{M}^*)} \leq 2$ , т.е. алгоритм  $\mathcal{A}$  находит 2-приближенное решение.

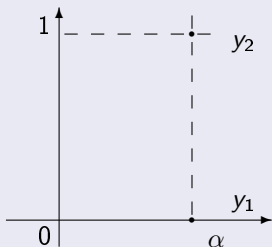
Покажем, что оценка 2 точности алгоритма достижима.

Для этого, приведем пример, показывающий существование таких входных данных задачи, для которых отношение  $F(\mathcal{M}_A)/F(\mathcal{M}^*)$  может быть сколь угодно близко к 2 и равно 2.

# Задача 1-MSSC-S-NF. Алгоритм решения

Оценки точности и трудоёмкости

## Пример достижимости



Пусть  $q = 2$ ,  $N = 2$ ,  $T_{\min} = 1$ ,  $T_{\max} = 2$ ,  $\mathcal{Y} = (y_1, y_2)$ , где  $y_1 = (0, \alpha)$ ,  $y_2 = (1, \alpha)$ .

Тогда если  $0 < \alpha < 1$ , то  $\mathcal{M}_A = \{2\}$ ,  $\mathcal{M}^* = \{1, 2\}$ ,  $F(\mathcal{M}_A) = \alpha^2$ ,  $F(\mathcal{M}^*) = 1/2$ . Таким образом, отношение  $F(\mathcal{M}_A)/F(\mathcal{M}^*) = 2\alpha^2$  может быть сколь угодно близко к 2 при  $\alpha \rightarrow 1$ .

Если же  $\alpha = 1$ , то имеем два алгоритмических решения: либо  $\mathcal{M}_A = \{1, 2\}$ , либо  $\mathcal{M}_A = \{2\}$ . При этом для второго решения  $F(\mathcal{M}_A) = 1$ ;  $\mathcal{M}^* = \{1, 2\}$ ,  $F(\mathcal{M}^*) = 1/2$  и  $F(\mathcal{M}_A)/F(\mathcal{M}^*) = 2$ , т.е. оценка точности алгоритма достижима.

# Задача 1-MSSC-S-NF. Алгоритм решения

Оценки точности и трудоёмкости

Оценим временную сложность алгоритма. Время вычислений определяется трудоёмкостью шага 3. На этом шаге  $N$  раз решается вспомогательная задача 1 с помощью алгоритма  $\mathcal{A}_1$ , трудоёмкость которого оценена в теореме 1. Отсюда следует оценка сложности. Теорема 2 доказана.

## Замечание 2

*В соответствии с замечанием 1 алгоритм  $\mathcal{A}$  полиномиален по  $N$  и по  $q$ , а его сложность можно оценить как  $\mathcal{O}(N^2(N + q))$ .*

Обоснован 2-приближенный полиномиальный алгоритм для решения NP-трудной в сильном смысле задачи разбиения последовательности на два кластера, которая индуцируется, в частности, оптимизационной моделью одной из актуальных проблем анализа данных.

Рассмотренная задача относится к числу практически неизученных в алгоритмическом плане. Поэтому исследование вопросов её аппроксимируемости, а также обоснование алгоритмов другого типа (асимптотически точных, рандомизированных и др.) для её решения представляется делом ближайшей перспективы.

Спасибо за внимание!