

Об оптимальном пропорционально справедливом распределении ресурсов в сотовых сетях 5G

Д. Г. Колосов¹ Л. С. Городецкий² Д. С. Миненков³

Рассматривается задача оптимального распределения ресурсов в сотовых сетях 5G в предположении полной загрузки (бесконечных потребностей пользователей). В отличие от сетей предыдущего поколения, в 5G пользователи могут делить между собой один и тот же частотный ресурс с небольшой потерей качества, что приводит к более сложной постановке задачи. В работе исследованы свойства пропорционально справедливого алгоритма распределения ресурсов, отвечающего задаче максимизации суммы логарифмов средних скоростей пользователей. На основе изученных свойств предложен вычислительно простой алгоритм, улучшающий рассматриваемый критерий качества по сравнению с другими известными алгоритмами. Сравнение алгоритмов проводилось с использованием реалистичных данных, сгенерированных с помощью библиотеки Sionna.

Ключевые слова: выпуклое математическое программирование, условия Каруша–Куна–Таккера, распределение радиоресурсов, планировщик, сотовые сети 5G, модель полной загрузки, пропорциональная справедливость, Sionna.

1. Введение

Multiple-Input/Multiple-Output (MIMO) — системы беспроводных коммуникаций с несколькими антеннами на передающей станции и несколь-

¹Колосов Дмитрий Григорьевич — студент каф. математической теории интеллектуальных систем мех.-мат. ф-та МГУ; младший инженер, Хуавэй, e-mail: kolosov.dmt@gmail.com.

Kolosov Dmitrii Grigorevich — student, Lomonosov Moscow State University, Faculty of Mechanics and Mathematics, Chair of Mathematical Theory of Intellectual Systems; junior engineer, Huawei Russian Research Institute.

²Городецкий Леонид Сергеевич — студент ф-та математики НИУ ВШЭ; ассистент инженер, Хуавэй, e-mail: ls.gorod.9@gmail.com.

Gorodetskii Leonid Sergeevich — student, Higher School of Economics, Faculty of Mathematics; assistant engineer, Huawei Russian Research Institute.

³Миненков Дмитрий Сергеевич — к.ф.-м.н., с.н.с., Институт проблем механики им. А.Ю. Ишлинского РАН; МГУ им. М.В. Ломоносова, мех.-мат. ф-т; e-mail: minenkov.ds@gmail.com.

Minenkov Dmitrii Sergeevich — PhD, senior researcher, Ishlinsky Institute for Problems in Mechanics RAS, Lomonosov Moscow State University, Faculty of Mechanics and Mathematics.

кими приемными антеннами у пользователей. MIMO — ключевая технология в пятом поколении (5G) мобильной связи. Она позволяет обслуживать несколько пользователей одновременно на одной частоте за счет технологии формирования луча (Beamforming) [1]. У пользователей (UE, User Equipment) может быть по 4 принимающих антенны (с учетом поляризации), а у базовой станции (BS, Base Station) 64 антенны. Сигнал для одного пользователя в этом случае представляется как четыре 64-мерных комплексных вектора. Чтобы избежать интерференции, сигналы, одновременно посылаемые разным пользователям, ортогонализуются по отношению к остальным пользователям, и из-за этого немного уменьшается мощность передаваемого сигнала для каждого конкретного пользователя. При этом некоторые пары пользователей в большей степени мешают друг другу, что сильно снижает качество связи при их одновременном обслуживании. Для решения этой проблемы используется процедура планировщика (Scheduler) — среди всех пользователей нужно выбрать подходящее подмножество для обслуживания. Этот выбор может быть сделан, например, в пользу наибольшей суммарной скорости передачи данных в конкретный момент времени (TTI, Transmission Time Interval) как на одной базовой станции [2, 3, 4, 5, 6], так и на нескольких [7, 8], что позволяет уменьшить межсотовую интерференцию. Большой обзор различных 5G планировщиков представлен в [9].

С другой стороны, важно, чтобы все пользователи были в той или иной степени обслужены за некоторый промежуток времени, что приводит нас к задаче справедливого распределения ресурсов. Один из подходов, связанный с минимизацией суммарного времени ожидания обслуживания, привел к появлению алгоритма SRPTF (Shortest Remaining Processing Time First)[10, 11] для 4G. Другой подход к вопросу справедливого распределения ресурсов — пропорциональная справедливость (PF, Proportional Fairness) [12, 13]. В сетях 4G эту задачу точно решает простой алгоритм градиентного PF-планировщика [14, 15], и при этом пользователи получают данные пропорционально качеству связи. Этот алгоритм был перенесен в сети 5G [5, 16], утратив теоретическое обоснование справедливости алгоритма. В текущей работе получены теоретические свойства PF в 5G MIMO, в том числе условие минимальной справедливости, гарантируемое PF алгоритмом. На основе этих свойств предложен простой алгоритм, решающий задачу пропорционально справедливого распределения ресурсов более успешно, чем хорошо известный алгоритм Baseline PF.

В этой работе мы рассматриваем нисходящую передачу данных от одной базовой станции (BS) с 64 антеннами к N пользовательским приемникам (UE) с 4 антеннами на каждом в модели полной загрузки (Full

Buffer) — когда пользователям всегда есть что получить от базовой станции.

Работа имеет следующую структуру. В пункте 2 даны основные определения, в том числе определение пропорциональной справедливости. В пункте 3 при изменяющихся во времени скоростях пользователей получен градиентный алгоритм, асимптотически решающий задачу пропорционально справедливого распределения ресурсов. В пункте 4 изучаются свойства предельного распределения ресурсов пропорционально справедливого алгоритма в предположении о постоянных во времени скоростях пользователей. В пункте 5 на основании полученных свойств предложен простой эвристический алгоритм и продемонстрирована его эффективность в численных симуляциях. В пункте 6 кратко перечислены основные результаты.

2. Постановка задачи и описание системы

Перед тем, как описывать математическую модель задачи распределения ресурсов в 5G, определим общее свойство пропорциональной справедливости [17].

Определение 1. Пусть $P \subset \mathbb{R}_{>0}^N$ — некоторое множество векторов $X = (X_1, X_2, \dots, X_N) \in P$. Тогда вектор $\hat{X} \in P$ называется пропорционально справедливым (*proportionally fair*) если для любого другого $X \in P$ его суммарный прирост относительно \hat{X} неположителен:

$$\sum_{n=1}^N \frac{X_n - \hat{X}_n}{\hat{X}_n} \leq 0$$

Также известно эквивалентное определение [17]:

Определение 2. Вектор $\hat{X} \in P$ называется пропорционально справедливым тогда и только тогда, когда

$$\hat{X} = \operatorname{argmax}_{X \in P} \sum_{n=1}^N \ln X_n$$

В контексте беспроводных систем нас интересует пропорциональная справедливость вектора X , состоящего из средних скоростей X_n пользователей. Точное определение X_n в этой работе будет даваться в зависимости от сделанных предположений в пунктах 3, 4 и 5. Исходя из определения 2, нам нужно максимизировать функцию

$$F(X) = \sum_{n=1}^N \ln X_n. \quad (2.1)$$

Теперь введем общие для всей работы обозначения. Множество всех пользователей обозначим через $\mathcal{N} = \{1, \dots, N\}$. Множество $\mathcal{A} \subset 2^{\mathcal{N}}$ — набор всех групп пользователей, которых можно одновременно обслуживать на одном частотном ресурсном блоке (RB – Resource Block). Будем считать, что \mathcal{A} задается некоторым критерием совместимости пользователей (например, SUS – Semi-orthogonal User Selection [5]). Через $A = |\mathcal{A}|$ обозначим количество групп в \mathcal{A} .

Для каждой группы $\alpha \in \mathcal{A}$ и для каждого пользователя $n \in \mathcal{N}$ задана скорость R_n^α пользователя n при обслуживании его в группе α . Если пользователь n не содержится в группе α , то мы полагаем $R_n^\alpha = 0$. Кроме того, для каждого пользователя n есть однопользовательская (SU, Single User) группа $\{n\} \in \mathcal{A}$, скорость пользователя в которой мы называем его SU-скоростью и обозначаем через $R_n = R_n^{\{n\}}$. Мы будем считать, что скорости пользователей не зависят от выбора RB.

3. Дискретно-динамическая модель

В этом пункте рассмотрим следующую 5G систему. В каждый из моментов времени ТТІ $t = 0, \dots, T$ базовая станция должна принять решение об обслуживании некоторого набора пользователей. При этом скорости пользователей зависят от времени: $R_n^\alpha(t)$ — скорость пользователя $n \in \mathcal{N}$ в группе $\alpha \in \mathcal{A}$ в момент времени t .

Так как скорости пользователей одинаковы на всех частотных ресурсах RB, мы будем считать, что на каждом ТТІ все доступные RB выделяются одной группе пользователей. В таком случае без потери общности можно рассматривать модель с одним виртуальным RB у станции, который на каждом ТТІ выделяется одной из групп пользователей.

В качестве управления возьмем индикаторную функцию

$$u_\alpha(t) = \begin{cases} 1, & \text{если группа } \alpha \text{ обслуживается на ТТІ } t, \\ 0, & \text{иначе.} \end{cases} \quad (3.1)$$

Отметим, что $\sum_{\alpha \in \mathcal{A}} u_\alpha(t) = 1$ для любого $t = 0, \dots, T$.

Определим средние скорости X_n пользователей следующим образом. В начальный момент $t = 0$ мы полагаем $X_n(0) = R_n(0)$, после чего X_n изменяются по закону экспоненциального среднего

$$X_n(t+1) = (1 - \beta)X_n(t) + \beta \sum_{\alpha \in \mathcal{A}} u_\alpha(t) R_n^\alpha(t), \quad (3.2)$$

где $0 < \beta \ll 1$ — некоторый коэффициент.

Определение (2) в этой модели перефразируется следующим образом:

Определение 3. *Пропорционально справедливый алгоритм — это алгоритм распределения ресурсов, который в каждый момент времени $t = 0, \dots, T$ максимизирует значение функции*

$$F(X(t+1)) = \sum_{n \in \mathcal{N}} \ln(X_n(t+1)). \quad (3.3)$$

С той же самой целевой функцией $F(X) = \sum_n \ln(X_n)$ можно связать более простой алгоритм, который является приближением пропорционально справедливого алгоритма.

Определение 4. *Градиентный PF — это алгоритм распределения ресурсов, который на каждом ТТИ $t = 0, \dots, T$ обслуживает группу $\alpha(t) \in \mathcal{A}$ с максимальным значением суммы*

$$M_\alpha = \sum_{n \in \alpha} \frac{R_n^\alpha(t)}{X_n(t)}. \quad (3.4)$$

Величина M_α называется PF-метрикой группы α .

Теорема 1. *Зафиксируем некоторый ТТИ t и значения $\{X_n(t)\}$ и $\{R_n^\alpha(t)\}$. Предположим, что PF-метрики M_α для всех групп $\alpha \in \mathcal{A}$ различны (это условие выполнено почти всегда). Тогда для достаточно малого $\beta > 0$ градиентный PF на данном ТТИ будет обслуживать ту же группу, что и пропорционально справедливый алгоритм.*

Доказательство. Возьмем группу $\alpha \in \mathcal{A}$ и найдем прирост целевой функции $F(X)$ в следующий ТТИ при обслуживании этой группы. Так как обслуживается только группа α , то при $\beta \rightarrow 0$ имеем

$$\begin{aligned} \Delta F &= F(X(t+1)) - F(X(t)) = \sum_{n \in \mathcal{N}} \ln X_n(t+1) - \sum_{n \in \mathcal{N}} \ln X_n(t) = \\ &= \sum_{n \in \mathcal{N}} \ln \left(\frac{X_n(t+1)}{X_n(t)} \right) = \sum_{n \in \mathcal{N}} \ln \left(\frac{(1-\beta)X_n(t) + \beta R_n^\alpha(t)}{X_n(t)} \right) = \\ &= \sum_{n \in \mathcal{N}} \ln \left(1 + \left(\frac{R_n^\alpha(t)}{X_n(t)} - 1 \right) \beta \right) = \sum_{n \in \mathcal{N}} \left(\frac{R_n^\alpha(t)}{X_n(t)} - 1 \right) \beta + o(\beta) = \\ &= \beta \sum_{n \in \alpha} \frac{R_n^\alpha(t)}{X_n(t)} - N\beta + o(\beta). \end{aligned} \quad (3.5)$$

Таким образом, при достаточно малом β обслуживание группы $\alpha \in \mathcal{A}$ с наибольшим значением PF-метрики $\sum_{n \in \alpha} \frac{R_n^\alpha(t)}{X_n(t)}$ максимизирует прирост функции $F(X)$. \square

Для каждого пользователя $n \in \mathcal{N}$ определим его PF-метрику

$$M_n = M_{\{n\}} = \frac{R_n(t)}{X_n(t)}.$$

В системе 4G эта метрика дает точное решение задачи пропорционально справедливому распределению ресурсов между пользователями, которое имеет малую вычислительную сложность $O(N)$. В системе 5G нахождение группы с максимальной PF-метрикой может быть задачей экспоненциальной сложности по N . Для того чтобы найти вычислительно простой алгоритм, приближенно решающий эту задачу, мы изучим некоторые свойства пропорционально справедливого алгоритма. Для этого нам потребуется более простая модель.

4. Непрерывно-стационарная модель

С этого момента мы считаем, что скорости пользователей R_n^α постоянны во времени. Такая постановка задачи изучалась в работе [14], и в терминах этой работы градиентный PF является градиентным алгоритмом для утилиты $\sum \ln(X_n)$. Опишем асимптотическое поведение градиентного алгоритма.

Рассмотрим стандартный симплекс

$$\Delta = \left\{ \mu = (\mu_\alpha)_{\alpha \in \mathcal{A}} \mid \mu_\alpha \geq 0 \text{ для всех } \alpha \in \mathcal{A} \text{ и } \sum_{\alpha \in \mathcal{A}} \mu_\alpha \leq 1 \right\}, \quad (4.1)$$

который назовем множеством допустимых распределений ресурсов, и рассмотрим выпуклый многогранник

$$P = \left\{ X = (X_1, \dots, X_N) \mid X_n = \sum_{\alpha \ni n} R_n^\alpha \mu_\alpha; \mu \in \Delta \right\}, \quad (4.2)$$

который назовем множеством векторов долговременных средних скоростей пользователей.

Замечание 1. Поясним, откуда происходят названия этих множеств. Для любого допустимого распределения ресурсов $\mu \in \Delta$ можно рассмотреть соответствующий вероятностный алгоритм распределения ресурсов (Weighted Round Robin): на каждом ТТІ обслуживаемая группа выбирается случайно с распределением μ , то есть вероятность обслуживания группы $\alpha \in \mathcal{A}$ равна μ_α . При таком алгоритме распределения ресурсов при $\beta \rightarrow 0$ и $t\beta \rightarrow \infty$ математическое ожидание вектора $X(t)$ средних скоростей пользователей стремится к вектору X с координатами

$$X_n = \sum_{\alpha \ni n} R_n^\alpha \mu_\alpha. \quad (4.3)$$

Значит многогранник P действительно можно рассматривать как множество возможных долговременных векторов X .

Применяя общий результат о сходимости градиентного алгоритма [14], получаем следующее утверждение.

Предложение 1. *При любом начальном значении $X(0)$ вектор $X(t)$, получаемый при обслуживании пользователей при помощи градиентного PF, при $\beta \rightarrow 0$ и $t\beta \rightarrow \infty$ стремится к вектору \hat{X} , который является решением оптимизационной задачи*

$$\begin{aligned} & \underset{X}{\text{maximize}} && F(X) = \sum_{n \in \mathcal{N}} \ln(X_n) \\ & \text{subject to} && X \in P, \end{aligned} \quad (4.4)$$

где P — это выпуклый многогранник из формулы (4.2).

Так как градиентный PF, совпадающий с пропорционально справедливым алгоритмом в пределе при $\beta \rightarrow 0$, асимптотически решает задачу оптимизации (4.4), то для исследования свойств пропорционально справедливого алгоритма распределения ресурсов мы будем изучать свойства решения задачи (4.4). Эту оптимизационную задачу мы будем называть непрерывно-стационарной моделью.

Если в задаче (4.4) ограничение $X \in P$ переписать в виде

$$X_n = \sum_{\alpha \ni n} R_n^\alpha \mu_\alpha, \quad \mu \in \Delta, \quad (4.5)$$

то получится переформулировка задачи в терминах допустимых распределений ресурсов μ :

$$\begin{aligned} & \underset{\mu}{\text{maximize}} && F(X(\mu)) = \sum_{n \in \mathcal{N}} \ln \left(\sum_{\alpha \ni n} \mu_\alpha R_n^\alpha \right) \\ & \text{subject to} && \mu \in \Delta. \end{aligned} \quad (4.6)$$

Пронумеруем каким-либо образом группы $\alpha \in \mathcal{A}$ и будем рассматривать наборы $\mu = (\mu_\alpha)_{\alpha \in \mathcal{A}}$ и $X = (X_n)_{n \in \mathcal{N}}$ как вектор-столбцы. Тогда скорости R_n^α пользователей в группах образуют матрицу $R \in \text{Mat}_{\mathbb{R}}(N, A)$: для группы $\alpha \in \mathcal{A}$ в соответствующем столбце R^α матрицы R стоят скорости пользователей в группе α . Тогда соотношения (4.5) примут вид

$$X = R\mu, \quad \mu \in \Delta. \quad (4.7)$$

Отсюда следует, что P — это выпуклая оболочка столбцов R^α матрицы R .

С использованием введенных обозначений задачу оптимизации (4.6) можно переписать в следующем виде:

$$\begin{aligned} & \underset{\mu}{\text{maximize}} && F(X(\mu)) = \sum_{n \in \mathcal{N}} \ln(X_n) \\ & \text{subject to} && X = R\mu; \mu \in \Delta. \end{aligned} \tag{4.8}$$

4.1. Свойства задачи оптимизации

Отметим, что если $X_n = 0$ для некоторого n (при $\mu \in \Delta$), то функция $F(X(\mu))$ принимает значение $-\infty$, а потому тот факт, что $F(X(\mu))$ определена не на всём Δ , не влияет на свойства максимумов функции.

В силу выпуклости функции суммы логарифмов, выпуклости симплекса Δ и многогранника P , у задачи (4.4) существует глобальный максимум, и любой локальный максимум также является глобальным. При этом у задачи (4.4) максимум единственный, а у задачи (4.8) возможно существование целого множества глобальных максимумов. Кроме того, все максимумы (4.8) соответствуют единственному максимуму (4.4).

Замечание 2. Если предположить, что R — матрица общего вида с неотрицательными компонентами, то для всех $k = 0, \dots, N - 1$ любая k -мерная грань многогранника P является k -мерным симплексом. Отсюда несложно увидеть, что оптимальному вектору $\hat{X} \in P$ соответствует единственное распределение ресурсов $\mu \in \Delta$, такое что $\hat{X} = R\mu$. Следовательно, если матрица R общего вида, то максимум в задаче (4.8) единственный.

Следующая теорема показывает, что несмотря на то, что набор групп \mathcal{A} может состоять из 2^N групп, оптимума можно достичь, используя значительно меньшее число групп.

Теорема 2. *В непрерывно-стационарной модели существует пропорционально справедливое распределение ресурсов μ , такое что не более N переменных μ_α отличны от нуля.*

Доказательство. Пусть $\hat{X} \in P$ — оптимальный (пропорционально справедливый) вектор средних скоростей пользователей. Так как функция $F(X) = \sum_{n \in \mathcal{N}} \ln(X_n)$ монотонна по каждой из переменных, \hat{X} лежит на границе P . Следовательно, \hat{X} лежит на одной из граней P меньшей размерности, а потому можно выбрать гиперплоскость $\Pi \subset \mathbb{R}^N$ так, чтобы она содержала точку \hat{X} , и пересекала P по грани размерности не более $N - 1$. Грань многогранника P , лежащую в гиперплоскости Π и содержащую \hat{X} , обозначим через P' .

Так как размерность гиперплоскости Π равна $N - 1$, то по теореме Каратеодори о выпуклой оболочке, найдется не более N вершин многогранника P' , таких что точка \hat{X} лежит в их выпуклой оболочке. Пусть

эти вершины соответствуют группам $\alpha_1, \dots, \alpha_k$, где $k \leq N$, то есть координаты этих вершин являются столбцами $R^{\alpha_1}, \dots, R^{\alpha_k}$ матрицы R . Так как точка \hat{X} лежит в выпуклой оболочке этих вершин, то

$$\hat{X} = \sum_{i=1}^k \mu_{\alpha_i} R^{\alpha_i}$$

для некоторых коэффициентов $\mu_{\alpha_i} \geq 0$, таких что $\sum_{i=1}^k \mu_{\alpha_i} = 1$. Если положить все остальные μ_{α} равными нулю, то такое распределение ресурсов по группам будет допустимым, и соответствующий вектор средних скоростей пользователей будет в точности равен оптимальному вектору \hat{X} . \square

Отметим, что аналогичным образом из теоремы Каратеодори можно получить более сильное утверждение: существует оптимальное распределение ресурсов μ , такое что столбцы матрицы R , соответствующие ненулевым компонентам μ , линейно независимы.

4.2. Свойства пропорционально справедливого распределения

Так как оптимизационная задача (4.8) выпукла, то её можно решить при помощи условий Каруша–Куна–Таккера [18]. Запишем функцию Лагранжа для задачи с неотрицательными переменными μ :

$$\mathcal{L}(\mu, \lambda_0) = \sum_{n \in \mathcal{N}} \ln \left(\sum_{\alpha \ni n} \mu_{\alpha} R_n^{\alpha} \right) + \lambda_0 \left(1 - \sum_{\alpha \in \mathcal{A}} \mu_{\alpha} \right), \quad (4.9)$$

и с помощью неё из условий Каруша–Куна–Таккера получим

$$\frac{\partial \mathcal{L}}{\partial \mu_{\alpha}} = \sum_{n \in \alpha} R_n^{\alpha} \frac{1}{X_n} - \lambda_0 = 0, \quad \text{если } \mu_{\alpha} > 0; \quad (4.10)$$

$$\frac{\partial \mathcal{L}}{\partial \mu_{\alpha}} = \sum_{n \in \alpha} R_n^{\alpha} \frac{1}{X_n} - \lambda_0 \leq 0, \quad \text{если } \mu_{\alpha} = 0. \quad (4.11)$$

Эти условия являются необходимыми и достаточными условиями максимума для допустимых μ .

Лемма 1. Из условий (4.10), (4.11) следует, что $\lambda_0 = N$.

Доказательство. Если $\hat{\mu}$ — точка глобального максимума, то через \mathcal{A}^+ обозначим множество групп $\alpha \in \mathcal{A}$, для которых $\mu_{\alpha} > 0$. Тогда с исполь-

зованием (4.10) можно записать

$$\begin{aligned}
 \lambda_0 &= \left(\sum_{\alpha \in \mathcal{A}^+} \mu_\alpha \right) \cdot \lambda_0 = \sum_{\alpha \in \mathcal{A}^+} \mu_\alpha \lambda_0 = \sum_{\alpha \in \mathcal{A}^+} \mu_\alpha \sum_{n \in \alpha} \frac{R_n^\alpha}{X_n} = \\
 &= \sum_{\alpha \in \mathcal{A}^+} \sum_{n \in \alpha} \frac{\mu_\alpha R_n^\alpha}{X_n} = \sum_{n \in \mathcal{N}} \sum_{\alpha \ni n} \frac{\mu_\alpha R_n^\alpha}{X_n} = \sum_{n \in \mathcal{N}} \frac{\sum_{\alpha \ni n} \mu_\alpha R_n^\alpha}{X_n} = \\
 &= \sum_{n \in \mathcal{N}} \frac{X_n}{X_n} = \sum_{n \in \mathcal{N}} 1 = N.
 \end{aligned}$$

□

Подставляя значение $\lambda_0 = N$ из леммы в (4.10), (4.11) получаем следующее.

Предложение 2. *Допустимое распределение ресурсов μ является решением оптимизационной задачи (4.8) тогда и только тогда, когда для каждой группы $\alpha \in \mathcal{A}$ выполнено условие*

$$\left\{ \begin{array}{l} \sum_{n \in \alpha} R_n^\alpha \frac{1}{X_n} = N, \text{ если } \mu_\alpha > 0; \\ \sum_{n \in \alpha} R_n^\alpha \frac{1}{X_n} \leq N, \text{ если } \mu_\alpha = 0, \end{array} \right. \quad (4.12)$$

где $X = R\mu$.

По аналогии с пунктом 3 число

$$M_\alpha = \sum_{n \in \alpha} R_n^\alpha \frac{1}{X_n} \quad (4.13)$$

будем называть PF-метрикой группы $\alpha \in \mathcal{A}$. Тогда условие оптимальности (4.12) можно сформулировать так: каждая из обслуживаемых групп имеет PF-метрику равную N , а остальные группы имеют PF-метрику не больше N .

Так как оптимальный вектор \hat{X} единственный, то PF-метрики групп при оптимальном распределении ресурсов определяются однозначно. Будем называть группу $\alpha \in \mathcal{A}$ оптимальной, если её PF-метрика при оптимальном распределении равна N . Из условия (4.12) следует, что для неоптимальной группы $\alpha \in \mathcal{A}$ при любом оптимальном распределении ресурсов μ выполнено $\mu_\alpha = 0$.

Так как неоптимальные группы не получают ресурсов при оптимальном распределении, то временно не будем их рассматривать и предположим, что каждая группа $\alpha \in \mathcal{A}$ имеет PF-метрику равную N . Получаем

систему уравнений на переменные $\mu = (\mu_\alpha)_{\alpha \in \mathcal{A}}$ и $X = (X_n)_{n \in \mathcal{N}}$:

$$\begin{cases} \sum_{\alpha \ni n} R_n^\alpha \mu_\alpha = X_n \text{ для каждого } n \in \mathcal{N} \\ \sum_{n \in \alpha} R_n^\alpha \frac{1}{X_n} = N \text{ для каждой } \alpha \in \mathcal{A}. \end{cases} \quad (4.14)$$

Если через $\frac{1}{X}$ обозначить вектор $(1/X_n)_{n \in \mathcal{N}}$, а через \mathbf{N} — вектор-столбец размера A , все координаты которого равны N , то систему (4.14) можно переписать в матричном виде:

$$\begin{cases} R\mu = X \\ R^T \frac{1}{X} = \mathbf{N}. \end{cases} \quad (4.15)$$

Используя матричную запись, можно найти решение задачи пропорционального распределения ресурсов в частном случае.

Предложение 3. *Если число оптимальных групп A равно числу пользователей N , а матрица R обратима, тогда оптимальное распределение ресурсов μ можно выразить явно, последовательно решая две системы линейных уравнений по следующим формулам:*

$$\begin{aligned} \frac{1}{X} &= (R^T)^{-1} \mathbf{N}, \\ \mu &= R^{-1} X. \end{aligned}$$

Перейдем к изучению других свойств оптимального распределения ресурсов. Определим относительные ресурсы пользователей

$$\nu_n = \frac{X_n}{R_n}, \quad (4.16)$$

где $R_n = R_n^{\{n\}}$ — SU-скорости пользователей. Следующая теорема даёт некоторое условие справедливости, которому удовлетворяет алгоритм PF.

Теорема 3. *При оптимальном распределении ресурсов μ для каждого пользователя $n \in \mathcal{N}$ выполнено условие минимальной справедливости*

$$\nu_n \geq \frac{1}{N}.$$

Доказательство. Для каждого пользователя $n \in \mathcal{N}$ рассмотрим SU-группу $\{n\} \in \mathcal{A}$. Согласно предложению 2, её PF-метрика при оптимальном распределении ресурсов не превосходит N , то есть

$$M_{\{n\}} = \frac{R_n}{X_n} \leq N,$$

откуда сразу следует

$$\nu_n = \frac{X_n}{R_n} \geq \frac{1}{N}.$$

□

4.3. Свойства решения в случае нескольких компонент

Дополнительное свойство пропорциональной справедливости получается, если пользователей можно разбить на несколько компонент $\mathcal{N}_1, \dots, \mathcal{N}_K$, то есть $\mathcal{N} = \bigsqcup_{k=1}^K \mathcal{N}_k$, причем каждая группа $\alpha \in \mathcal{A}$ целиком содержится в одной из компонент \mathcal{N}_k . Множество групп, содержащихся в компоненте \mathcal{N}_k обозначим через \mathcal{A}_k , а общее количество ресурсов k -ой компоненты обозначим через $\eta_k = \sum_{\alpha \in \mathcal{A}_k} \mu_\alpha$. Понятно, что

$$\sum_{k=1}^K \eta_k = \sum_{\alpha \in \mathcal{A}} \mu_\alpha \leq 1. \quad (4.17)$$

Число пользователей в компоненте \mathcal{N}_k обозначим через N_k . В этой ситуации при подходящей нумерации групп матрица R будет иметь блочный вид.

Следующая теорема показывает, как должны распределяться ресурсы между компонентами.

Теорема 4. *Если в непрерывно-стационарной модели пользователи разбиты на компоненты $\mathcal{N}_1, \dots, \mathcal{N}_K$, то при оптимальном (пропорционально справедливом) распределении ресурсов выполняется*

$$\eta_k = \frac{N_k}{N} \quad (4.18)$$

для каждой компоненты $k = 1, \dots, K$. Иными словами, ресурсы распределяются между компонентами пропорционально их размеру.

Доказательство. Введем относительное количество ресурсов γ_α^k группы $\alpha \in \mathcal{A}$ в компоненте \mathcal{N}_k :

$$\gamma_\alpha^k = \begin{cases} \mu_\alpha / \eta_k, & \text{если } \alpha \in \mathcal{A}_k; \\ 0 & \text{, иначе,} \end{cases}$$

где η_k — ресурсы компоненты \mathcal{N}_k . Тогда

$$\begin{aligned}
F(X) &= \sum_{n \in \mathcal{N}} \ln(X_n) = \sum_{k=1}^K \sum_{n \in \mathcal{N}_k} \ln(X_n) = \sum_{k=1}^K \sum_{n \in \mathcal{N}_k} \ln \left(\sum_{\alpha \ni n} R_n^\alpha \eta_k \gamma_\alpha^k \right) = \\
&= \sum_{k=1}^K \sum_{n \in \mathcal{N}_k} \left(\ln \eta_k + \ln \left(\sum_{\alpha \ni n} R_n^\alpha \gamma_\alpha^k \right) \right) = \\
&= \sum_{k=1}^K N_k \ln \eta_k + \sum_{k=1}^K \sum_{n \in \mathcal{N}_k} \ln \left(\sum_{\alpha \ni n} R_n^\alpha \gamma_\alpha^k \right).
\end{aligned}$$

Так как переменные γ_α^k можно считать независимыми от η_k , то исходная задача оптимизации разбивается на независимые подзадачи — задачу распределения ресурсов по компонентам

$$\begin{aligned}
&\underset{\eta}{\text{maximize}} && F_0(\eta) = \sum_{k=1}^K N_k \ln \eta_k \\
&\text{subject to} && \sum_{k=1}^K \eta_k \leq 1; \eta_k \geq 0 \quad \forall k;
\end{aligned} \tag{4.19}$$

и K независимых задач по распределению ресурсов внутри каждой компоненты:

$$\begin{aligned}
&\underset{\gamma^k}{\text{maximize}} && F_k(\gamma^k) = \sum_{n \in \mathcal{N}_k} \ln \left(\sum_{\alpha \ni n} R_n^\alpha \gamma_\alpha^k \right) \\
&\text{subject to} && \sum_{\alpha \in \mathcal{A}_k} \gamma_\alpha^k = 1; \gamma_\alpha^k \geq 0 \quad \forall \alpha \in \mathcal{A}_k.
\end{aligned} \tag{4.20}$$

Выпишем функцию Лагранжа для задачи (4.19):

$$\mathcal{L}(\eta, \lambda_0) = \sum_{k=1}^K N_k \ln \eta_k + \lambda_0 \left(1 - \sum_{k=1}^K \eta_k \right), \tag{4.21}$$

и из условий Каруша–Куна–Таккера получим

$$\frac{\partial \mathcal{L}}{\partial \eta_k} = \frac{N_k}{\eta_k} - \lambda_0 = 0 \tag{4.22}$$

для всех $k = 1, \dots, K$. Далее

$$\lambda_0 = 1 \cdot \lambda_0 = \left(\sum_{k=1}^K \eta_k \right) \lambda_0 = \sum_{k=1}^K \eta_k \lambda_0 = \sum_{k=1}^K N_k = N,$$

откуда

$$\eta_k = \frac{N_k}{\lambda_0} = \frac{N_k}{N}.$$

□

У доказанной теоремы 4 есть важный частный случай. Предположим, что есть пользователь $n \in \mathcal{N}$, который может обслуживаться только в SU-группе $\{n\} \in \mathcal{A}$ (такой пользователь называется SU-пользователем). Применяя теорему 4 к компоненте, состоящей из одного пользователя n , мы получаем, что при оптимальном распределении ресурсов у этого SU-пользователя

$$\nu_n = \frac{1}{N}. \quad (4.23)$$

5. Численное исследование

В этом пункте будут предложены эвристические алгоритмы, основанные на наших теоретических результатах. Затем мы их проверим в тестовой среде, имитирующей 5G MIMO систему.

5.1. Описание тестовой среды

Тестовый симулятор написан на языке Python для модели Full Buffer с использованием библиотек Numpy — включает в себя различные математические операции и поддержку многомерных массивов, и Sionna (NVIDIA) [19] — продвинутой библиотеки для симуляций телекоммуникационных систем. Этапы работы симулятора следующие:

- 1) Задаются параметры системы: коэффициент усреднения $\beta = 0.01$, количество ТТИ $T = 1500$, количество пользователей $N = 30$, начальное значение вектора средних скоростей пользователей $X(0) = (R_1(0), \dots, R_N(0))$.
- 2) Канал связи между отправляющей антенной и принимающей будем описывать через комплексные числа, показывающие изменение фазы и затухание амплитуды электромагнитной волны при передаче. В 5G MIMO (Multiple-Input/Multiple-Output) системе у базовой станции 64 антенны, а у пользователей по 4 антенны. Следовательно, канал пользователя $n \in \mathcal{N}$ описывается комплексной матрицей $H_n(t)$ размера 4×64 — для каждой принимающей антенны и каждой отправляющей антенны соответствующий элемент матрицы называется канальным числом. Каналы генерируются при помощи библиотеки Sionna в модели городской среды и изменяются со временем согласно случайному гауссовскому блужданию: $H_n(t+1) = H_n(t) + H_n(t) \cdot \mathcal{N}(0, 0.03)$, где $\mathcal{N}(a, \sigma^2)$ — нормальная случайная величина. Значения ОСШ (Отношение Сигнал-Шум) каналов лежат в диапазоне 0dB – 25dB (точное определение ОСШ пользователя будет дано далее). При этом, увеличивая общую пропускную способность системы, будем передавать только по $L = 2$

(эффективных) потоков данных каждому пользователю одновременно, всё еще задействуя все его антенны [20]. Ограничение на количество одновременно обслуживаемых пользователей $|\alpha| \leq 8$.

- 3) На каждом ТТИ $t = 0, \dots, T$ симулятор делает следующее. Используя текущий вектор средних скоростей пользователей $X(t)$ и текущие каналы пользователей, планировщик выбирает группу пользователей $\alpha(t) \subset \mathcal{N}$. Далее для выбранной группы $\alpha(t)$ происходит ММО симуляция передачи данных при помощи RZF (Regularized Zero-Forcing) предварительного кодирования и MMSE-IRC (Minimum Mean Squared Error - Interference Rejection Combiner) распознавания [20], находятся фактические значения $R_n^\alpha(t)$ скоростей пользователей при обслуживании в $\alpha(t)$ и с их помощью обновляется вектор X средних скоростей пользователей:

$$X_n(t+1) = (1 - \beta)X_n(t) + \beta R_n^{\alpha(t)}(t). \quad (5.1)$$

Также обновляется вектор общих объемов переданных данных

$$\text{ТН}_n(t+1) = \text{ТН}_n(t) + R_n^{\alpha(t)}(t), \quad (5.2)$$

который в начальный момент $t = 0$ полагается равным нулю.

- 4) По окончании работы симулятор считает следующие показатели:

- $\text{gMean}(X) = \left(\prod_{n=1}^N X_n \right)^{1/N} = \exp \left(\frac{1}{N} \sum_{n=1}^N \ln(X_n) \right)$,
- $\text{sum}(\text{ТН}) = \sum_{n=1}^N \text{ТН}_n$ — суммарный объём переданных данных.

5.2. Алгоритмы

Все рассматриваемые нами алгоритмы используют следующую конструкцию, упрощающую задачу. Для пользователя n через v_n^1, \dots, v_n^L обозначим первые L сингулярных векторов [21] его матрицы канала H_n . Определим корреляцию пользователей n и k формулой

$$c_{n,k} = |\langle v_n^1, v_k^1 \rangle|. \quad (5.3)$$

Пользователи n и k называются скоррелированными, если их корреляция $c_{n,k}$ больше порогового значения $c_{th} = 0,3$. Для пользователя n определим его отношение сигнал-шум при обслуживании в одиночку (SU) как SINR_n^{su} (Signal to Interference and Noise Ratio):

$$\text{SINR}_n^{su} = \frac{1}{L\sigma^2} \left(\prod_{l=1}^L s_n^l \right)^{2/L}, \quad (5.4)$$

где s_n^1, \dots, s_n^L — первые L сингулярных значений матрицы H_n . Кроме того, если у пользователя n SINR_n^{su} меньше порогового значения $\text{SINR}_{th} = 5\text{dB}$, то он называется SU-пользователем. Остальные пользователи называются MU-пользователями.

Рассмотрим неориентированный граф G , вершины которого — пользователи, и две вершины соединены ребром, если они оба являются MU-пользователями и не скоррелированы. Клика (или полный подграф) в графе — это набор вершин, любые две из которых соединены ребром. Каждый из рассматриваемых алгоритмов может обслуживать нескольких пользователей в одной группе, только если они образуют клику в графе G .

Алгоритм перебора клик

Сначала опишем работу алгоритма перебора клик в графе, близкого к оптимальному при малых значениях β . На каждом ТТІ надо

- 1) Для каждой клики C в графе G вычислить её приближенную PF-метрику

$$M_C^{appr} = \sum_{n \in C} \frac{R_n^C}{X_n}, \quad (5.5)$$

где скорость R_n^C пользователя n в клике C оценивается по формуле

$$R_n^C = L \cdot \log_2 \left(1 + \frac{\text{SINR}_n^{su}}{|C|} \right), \quad (5.6)$$

которая выражает максимальную пропускную способность канала [22] при распределении всей мощности базовой станции между $|C|$ пользователями [20]

- 2) На данном ТТІ обслужить клику с максимальной приближенной PF-метрикой.

Формула (5.6) — это приближение скорости пользователя n при обслуживании в клике C , использующее, что пользователи в клике слабо коррелируют друг с другом. Действительно, если пренебречь корреляцией пользователей в клике и считать, что мощность делится между пользователями в группе поровну, то формула (5.6) — это верхняя оценка скорости пользователя в группе, полученная по теореме Шеннона. Таким образом, согласно предложению 1, этот алгоритм будет близок к оптимуму при малых β . Так как количество клик в графе может быть очень большим, данный алгоритм неприменим на практике.

Во всех следующих алгоритмах мы будем использовать приближенные значения SU-скоростей пользователей R_n , которые можно найти по формуле Шеннона [22]

$$R_n = L \cdot \log_2(1 + \text{SINR}_n^{\text{su}}). \quad (5.7)$$

Поиск клики по метрике пользователей

Далее мы будем изучать класс алгоритмов, использующих поиск клики по метрике пользователей. Для удобного описания алгоритмов введем следующее обозначение. Пусть в множестве \mathcal{N} всех пользователей задано некоторое подмножество \mathcal{N}' , а также задана некоторая функция $m: \mathcal{N}' \rightarrow \mathbb{R}$, называемая метрикой пользователей. Определим $C(\mathcal{N}', m)$ как клику в графе, лежащую в \mathcal{N}' , которая набирается по следующему алгоритму:

- 1) Отсортировать пользователей из \mathcal{N}' по убыванию метрик $m(n)$: далее мы считаем, что

$$m(n_1) \geq m(n_2) \geq \dots \geq m(n_{N'}),$$

где N' — число пользователей в \mathcal{N}' , а $n_1, \dots, n_{N'}$ — индексы пользователей из \mathcal{N}' .

- 2) Инициализировать клику C как пустое множество.
- 3) Пройти по набору пользователей $n_1, \dots, n_{N'}$, последовательно добавляя пользователей в клику C по следующему правилу: если пользователь n_i соединен ребрами с каждым из пользователей, которые к данному моменту находятся в клике C , то он добавляется в неё; в противном случае пользователь пропускается. Кроме того, если клика C в какой-то момент стала максимального допустимого размера, то набор пользователей в неё прекращается.

Baseline PF

Референсный алгоритм PF-ZFBF [5], который мы будем называть Baseline PF, устроен следующим образом. На каждом ТТГ он набирает клику $C(\mathcal{N}, \text{PF})$ из множества всех пользователей при помощи PF-метрики

$$M_n = \frac{R_n}{X_n}, \quad (5.8)$$

где R_n — это SU-скорость пользователя n , и обслуживает эту клику. Далее мы увидим, что одна из проблем Baseline PF в том, что он слишком много ресурсов выделяет SU-пользователям.

Компонентный PF

Разделим множество \mathcal{N} всех пользователей на 2 компоненты в смысле теоремы 4: на компоненту \mathcal{N}^{MU} из всех MU-пользователей и на компоненту \mathcal{N}^{SU} из всех SU-пользователей. На каждом ТТІ при помощи PF-метрики (5.8) наберем клики $C(\mathcal{N}^{MU}, PF)$ и $C(\mathcal{N}^{SU}, PF)$ в этих компонентах. В компоненте \mathcal{N}^{SU} клика — это один пользователь с наибольшей PF-метрикой. Далее выберем для обслуживания ту клику из двух найденных, у которой приближенная PF-метрика группы M_C^{appr} наибольшая.

В симуляции мы увидим, что этот алгоритм решает задачу распределения ресурсов между SU-пользователями и MU-пользователями более успешно, чем Baseline PF.

5.3. Результаты симуляции

Результаты симуляции работы трёх описанных алгоритмов на 250 сценариях представлены на рис.1. Для обоих показателей эффективности работы алгоритмов представлен график функции распределения прироста ($gain = (x/x_0 - 1) \cdot 100\%$) значения показателя относительно Baseline PF. По осям абсцисс отображены приросты в процентах показателей относительно Baseline PF, а по осям ординат — проценты от общего числа сценариев.

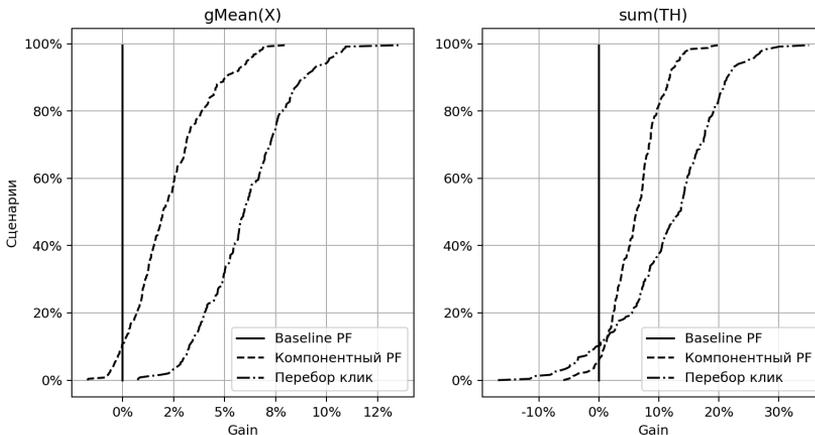


Рис. 1. Графики функций распределения прироста по сценариям относительно Baseline PF значений показателей $gMean(X)$ и $sum(TH)$.

Для углубленного анализа алгоритмов каждому пользователю сопоставим показатель

$$\nu_n = \frac{\text{ГН}_n}{\sum_{t=1}^T R_n(t)}, \quad (5.9)$$

отражающий отношение объема переданных пользователю данных к максимально возможному объему данных, который можно было бы передать пользователю за время симуляции. Отметим, что этот показатель согласуется с аналогичным показателем $\nu_n = X_n/R_n$ в непрерывно-стационарной модели. Теорема 3 показывает, что оптимальное распределение ресурсов в пределе должно удовлетворять условию минимальной справедливости

$$\nu_n \geq \frac{1}{N} \approx 0,033. \quad (5.10)$$

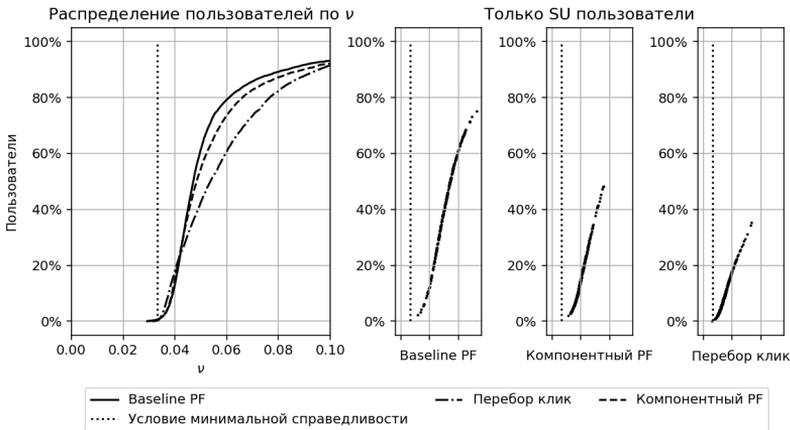


Рис. 2. На левом графике распределение ν_n по всем пользователям во всех сценариях в сравнении с условием минимальной справедливости. На трех правых графике точками выделены SU-пользователи (более 95% времени имели $\text{SINR}_n^{su} < 5\text{dB}$).

Отметим основные особенности этих графиков.

- 1) Алгоритм перебора клик в графе действительно лучше остальных алгоритмов максимизирует полезность $\text{gMean}(X)$.
- 2) Все три описанных алгоритма удовлетворяют минимальному условию справедливости (5.10).
- 3) Исходя из графиков на рис. 2 Baseline PF излишне часто обслуживает SU-пользователей.

- 4) Компонентный PF решает эту проблему, чаще обслуживая МУ-пользователей, за счет чего он имеет более высокие показатели $g\text{Mean}(X)$ и $\text{sum}(TН)$ по сравнению с Baseline PF.

Таким образом, симуляции показывают, что компонентный PF по основным показателям превосходит Baseline PF, при этом обладая практически такой же вычислительной сложностью. Кроме того, продемонстрировано выполнение условия минимальной справедливости.

6. Заключение

В работе рассмотрена проблема распределения радио-ресурсов в модели полной загрузки для одной базовой станции и нескольких приемников. Из задачи максимизации суммы логарифмов средних скоростей пользователей получен градиентный алгоритм обслуживания для одного ресурсного блока, который в каждый момент времени обслуживает группу пользователей с наибольшей групповой PF-метрикой. При предельном переходе по времени и шагу усреднения получена непрерывно-стационарная модель. В этой модели получено несколько свойств оптимального распределения ресурсов, таких как: ограничение количества используемых групп, пропорциональная зависимость количества ресурсов от размера компонент пользователей, оценка снизу необходимого количества выделяемых каждому пользователю ресурсов. На основе полученных свойств предложен алгоритм распределения ресурсов (Компонентный PF), который улучшает основные показатели по сравнению с Baseline PF, не нарушая условие минимальной справедливости. Проведены 5G MIMO симуляции с каналами из библиотеки Sionna, результаты которых показывают эффективность предложенного алгоритма. Этот алгоритм по сложности практически не отличается от референсного, что позволяет легко внедрить предложенный алгоритм в MIMO системы.

Благодарности

Авторы признательны Д.А. Заеву, В.В. Кузнецову, А.М. Миронову, М.Ю. Попеленскому, Д.А. Шмелькину за ценные дискуссии и поддержку. Работа частично выполнена в рамках кооперационного проекта Техкомпании Хуавей и МГУ, частично – в рамках государственного задания ИПМех РАН (№ госрегистрации 123021700044-0).

Список литературы

- [1] Yang, Hong and Marzetta, Thomas L., “Performance of Conjugate and Zero-Forcing Beamforming in Large-Scale Antenna Systems”, *IEEE Journal on Selected Areas in Communications*, **31**:2 (2013), 172–179.
- [2] Chataut, Robin and Akl, Robert., “Channel Gain Based User Scheduling for 5G Massive MIMO Systems”, *2019 IEEE 16th International Conference on Smart Cities: Improving Quality of Life Using ICT & IoT and AI (HONET-ICT)*, 2019, 049–053.
- [3] Naeem, Muddasar and Bashir, Sajid and Ullah, Zaib and Syed, Aqeel A., “A near optimal scheduling algorithm for efficient radio resource management in multi-user MIMO systems”, *Wireless Personal Communications*, **106**:3 (2019), 1411–1427.
- [4] Xia, Xin and Fang, Shu and Wu, Gang and Li, Shaoqian, “Joint User Pairing and Precoding in MU-MIMO Broadcast Channel with Limited Feedback”, *IEEE Communications Letters*, **14**:11 (2010), 1032–1034.
- [5] Taesang Yoo and Goldsmith, A., “On the optimality of multiantenna broadcast scheduling using zero-forcing beamforming”, *IEEE Journal on Selected Areas in Communications*, **24**:3 (2006), 528–541.
- [6] Huang, Xiaoyan and Xue, Guoliang and Yu, Ruozhou and Leng, Supeng, “Joint Scheduling and Beamforming Coordination in Cloud Radio Access Networks With QoS Guarantees”, *IEEE Transactions on Vehicular Technology*, **65**:7 (2016), 5449–5460.
- [7] Tang, Xiaojun and Ramprasad, Sean A. and Papadopoulos, Haralabos, “Multi-Cell User-Scheduling and Random Beamforming Strategies for Downlink Wireless Communications”, *2009 IEEE 70th Vehicular Technology Conference Fall*, 2009, 1–5.
- [8] Li, Min and Collings, Iain B. and Hanly, Stephen V. and Liu, Chunshan and Whiting, Philip, “Multicell Coordinated Scheduling With Multiuser Zero-Forcing Beamforming”, *IEEE Transactions on Wireless Communications*, **15**:2 (2016), 827–842.
- [9] Mamane, Asmae and Fattah, Mohammed and Ghazi, Mohammed El and Bekkali, Moulhime El and Balboul, Younes and Mazer, Said, “Scheduling Algorithms for 5G Networks and Beyond: Classification and Survey”, *IEEE Access*, **10** (2022), 51643–51661.
- [10] Linus Schrage, “Letter to the Editor—A Proof of the Optimality of the Shortest Remaining Processing Time Discipline”, *Operations Research*, **16**:3 (1968), 687–690.

- [11] Samuli Aalto and Aleksi Penttinen and Pasi Lassila and Prajwal Osti, “Optimal size-based opportunistic scheduler for wireless systems”, *QUEUEING SYSTEMS*, **72**:1–2 (2012), 5–30.
- [12] Kelly, Frank P, “Charging and rate control for elastic traffic”, *European transactions on Telecommunications*, **8**:1 (1997), 33–37.
- [13] Kelly, Frank P and Maulloo, Aman K and Tan, David Kim Hong, “Rate control for communication networks: shadow prices, proportional fairness and stability”, *Journal of the Operational Research society*, **49**:3 (1998), 237–252.
- [14] Alexander L. Stolyar, “On the Asymptotic Optimality of the Gradient Scheduling Algorithm for Multiuser Throughput Allocation”, *Operations Research*, **53**:1 (2005), 12–25.
- [15] Sun, Zhishui and Yin, Changchuan and Yue, Guangxin, “Reduced-Complexity Proportional Fair Scheduling for OFDMA Systems”, *2006 International Conference on Communications, Circuits and Systems*, **2** (2006), 1221–1225.
- [16] Femenias, Guillem and Riera-Palou, Felip and Mestre, Xavier and Olmos, Juan J., “Downlink Scheduling and Resource Allocation for 5G MIMO-Multicarrier: OFDM vs FBMC/OQAM”, *IEEE Access*, **5** (2017), 13770–13786.
- [17] Le Boudec, Jean-Yves, “Rate adaptation, Congestion Control and Fairness: A Tutorial”, 2002.
- [18] Миронов Андрей Михайлович, “Машинное обучение. Часть 1”, <https://is.ifmo.ru/verification/machine-learning-mironov.pdf>.
- [19] Jakob Hoydis and Sebastian Cammerer and Fayçal Ait Aoudia and Avinash Vem and Nikolaus Binder and Guillermo Marcus and Alexander Keller, “Sionna: An Open-Source Library for Next-Generation Physical Layer Research”, 2023, arXiv: <https://arxiv.org/abs/2203.11854>.
- [20] Evgeny Bobrov and Boris Chinyaev and Viktor Kuznetsov and Dmitrii Minenkov and Daniil Yudakov, “Power Allocation Algorithms for Massive MIMO Systems with Multi-Antenna Users”, 2022, arXiv: <https://arxiv.org/abs/2201.08068>.
- [21] Роджер Хорн и Чарльз Джонсон, *Матричный анализ*, **656**, Мир, 1989.
- [22] Shannon, Claude E, “A mathematical theory of communication”, *The Bell system technical journal*, **27**:3 (1948), 379–423.

On the Optimal Proportional Fair Radio Resource Management in 5G Cellular Networks

Kolosov D.G., Gorodetskii L.S., Minenkov D.S.

We consider the problem of the optimal radio resource management in MU-MIMO 5G cellular networks in Full Buffer traffic model. Unlike in the networks of the previous generation, in 5G users can share the same frequency resource with a small loss of quality, which leads to a more complicated problem statement. We study properties of the proportionally fair scheduler that meets the problem of maximizing the sum of logarithms of users' average rates. We propose a computationally simple algorithm based on the studied properties that improves the proportional fairness in comparison with other well-known algorithms. The algorithms were compared using realistic data generated by the Sionna library.

Keywords: convex optimization, Karush–Kuhn–Tucker conditions, radio resource management, scheduler, 5G cellular networks, MU-MIMO, full buffer, proportional fairness, Sionna.

References

- [1] Yang, Hong and Marzetta, Thomas L., “Performance of Conjugate and Zero-Forcing Beamforming in Large-Scale Antenna Systems”, *IEEE Journal on Selected Areas in Communications*, **31**:2 (2013), 172–179.
- [2] Chataut, Robin and Akl, Robert., “Channel Gain Based User Scheduling for 5G Massive MIMO Systems”, *2019 IEEE 16th International Conference on Smart Cities: Improving Quality of Life Using ICT & IoT and AI (HONET-ICT)*, 2019, 049–053.
- [3] Naeem, Muddasar and Bashir, Sajid and Ullah, Zaib and Syed, Aqeel A., “A near optimal scheduling algorithm for efficient radio resource management in multi-user MIMO systems”, *Wireless Personal Communications*, **106**:3 (2019), 1411–1427.
- [4] Xia, Xin and Fang, Shu and Wu, Gang and Li, Shaoqian, “Joint User Pairing and Precoding in MU-MIMO Broadcast Channel with Limited Feedback”, *IEEE Communications Letters*, **14**:11 (2010), 1032–1034.
- [5] Taesang Yoo and Goldsmith, A., “On the optimality of multiantenna broadcast scheduling using zero-forcing beamforming”, *IEEE Journal on Selected Areas in Communications*, **24**:3 (2006), 528–541.
- [6] Huang, Xiaoyan and Xue, Guoliang and Yu, Ruozhou and Leng, Supeng, “Joint Scheduling and Beamforming Coordination in Cloud

- Radio Access Networks With QoS Guarantees”, *IEEE Transactions on Vehicular Technology*, **65**:7 (2016), 5449–5460.
- [7] Tang, Xiaojun and Ramprashad, Sean A. and Papadopoulos, Haralabos, “Multi-Cell User-Scheduling and Random Beamforming Strategies for Downlink Wireless Communications”, *2009 IEEE 70th Vehicular Technology Conference Fall*, 2009, 1–5.
- [8] Li, Min and Collings, Iain B. and Hanly, Stephen V. and Liu, Chunshan and Whiting, Philip, “Multicell Coordinated Scheduling With Multiuser Zero-Forcing Beamforming”, *IEEE Transactions on Wireless Communications*, **15**:2 (2016), 827–842.
- [9] Mamane, Asmae and Fattah, Mohammed and Ghazi, Mohammed El and Bekkali, Moulhime El and Balboul, Younes and Mazer, Said, “Scheduling Algorithms for 5G Networks and Beyond: Classification and Survey”, *IEEE Access*, **10** (2022), 51643–51661.
- [10] Linus Schrage, “Letter to the Editor—A Proof of the Optimality of the Shortest Remaining Processing Time Discipline”, *Operations Research*, **16**:3 (1968), 687–690.
- [11] Samuli Aalto and Aleksi Penttinen and Pasi Lassila and Prajwal Osti, “Optimal size-based opportunistic scheduler for wireless systems”, *QUEUEING SYSTEMS*, **72**:1–2 (2012), 5–30.
- [12] Kelly, Frank P, “Charging and rate control for elastic traffic”, *European transactions on Telecommunications*, **8**:1 (1997), 33–37.
- [13] Kelly, Frank P and Maulloo, Aman K and Tan, David Kim Hong, “Rate control for communication networks: shadow prices, proportional fairness and stability”, *Journal of the Operational Research society*, **49**:3 (1998), 237–252.
- [14] Alexander L. Stolyar, “On the Asymptotic Optimality of the Gradient Scheduling Algorithm for Multiuser Throughput Allocation”, *Operations Research*, **53**:1 (2005), 12–25.
- [15] Sun, Zhishui and Yin, Changchuan and Yue, Guangxin, “Reduced-Complexity Proportional Fair Scheduling for OFDMA Systems”, *2006 International Conference on Communications, Circuits and Systems*, **2** (2006), 1221–1225.
- [16] Femenias, Guillem and Riera-Palou, Felip and Mestre, Xavier and Olmos, Juan J., “Downlink Scheduling and Resource Allocation for 5G MIMO-Multicarrier: OFDM vs FBMC/OQAM”, *IEEE Access*, **5** (2017), 13770–13786.

- [17] Le Boudec, Jean-Yves, “Rate adaptation, Congestion Control and Fairness: A Tutorial”, 2002.
- [18] Mironov A.M., “Machine learning. Part 1” (In Russian), <https://is.ifmo.ru/verification/machine-learning-mironov.pdf>.
- [19] Jakob Hoydis and Sebastian Cammerer and Fayçal Ait Aoudia and Avinash Vem and Nikolaus Binder and Guillermo Marcus and Alexander Keller, “Sionna: An Open-Source Library for Next-Generation Physical Layer Research”, 2023, arXiv: <https://arxiv.org/abs/2203.11854>.
- [20] Evgeny Bobrov and Boris Chinyaev and Viktor Kuznetsov and Dmitrii Minenkov and Daniil Yudakov, “Power Allocation Algorithms for Massive MIMO Systems with Multi-Antenna Users”, 2022, arXiv: <https://arxiv.org/abs/2201.08068>.
- [21] Roger A. Horn and Charles R. Johnson, *Matrix Analysis*, **656**, Mir, 1989 (In Russian).
- [22] Shannon, Claude E, “A mathematical theory of communication”, *The Bell system technical journal*, **27**:3 (1948), 379–423.