

Конспект лекции 1. Основные понятия и определения.

Сергей Николенко*

22 мая 2008 г.

Содержание

1 Введение	1
1.1 О чём этот курс	1
1.2 История становления дисциплины	2
1.3 Реальные применения дизайна механизмов	2
2 Теория игр: примеры	3
2.1 Дилемма заключённого	3
2.2 Трагедия общин	4
2.3 Парадокс аукциона за доллар	4
2.4 Winner's curse	5
2.5 Парадокс Браесса	5
3 Основные концепции теории игр	7
3.1 Агенты, стратегии, функция полезности	7
3.2 Равновесие Нэша	8
3.3 Доминантные стратегии и аукцион Викри	9
3.4 Равновесие по Байесу-Нэшу	10
4 Основные понятия дизайна механизмов	11
4.1 Функция социального выбора и механизмы	11
4.2 Оптимальность по Парето	12
4.3 Предположения об агентах	13

1 Введение

1.1 О чём этот курс

Дизайн механизмов (mechanism design) — это раздел *теории игр*, которая, в свою очередь, изучает взаимодействие между агентами, при котором каждый агент пытается выбрать стратегию, максимизирующую его собственную прибыль.

*Законоспектировали Андрей Якушев и Михаил Чураков.

Дизайн механизмов — это конструктивный подход, позволяющий создать такой механизм взаимодействия, при котором эгоистические действия каждого из агентов в сумме приведут к решению, оптимальному с точки зрения общей целевой функции.

Главный пример дизайна механизмов — аукционы. Цели в обычном аукционе:

- либо организатор пытается максимизировать общую прибыль (social welfare);
- либо продавец пытается сделать такой аукцион, чтобы продать дороже.

Кроме того, хочется достичь ситуации, при которой выявляются истинные предпочтения участников (truthfulness), и, конечно, решение должно быть в каком-либо смысле оптимальным и/или устойчивым, иначе оно не сможет реализоваться.

1.2 История становления дисциплины

Слово «mechanism» в этом контексте ввёл Леонид Гурвич (Leonid Hurwicz). Родился в Москве в 1917 году, жил в Польше, в 1940 эмигрировал в США. В 1960 он сформулировал основные положения теории экономических механизмов, в 1972 сформулировал свойство правдивости, а затем и принцип выявления, с которого по сути и началось исследование децентрализованных систем применительно в экономике.

Дальше Эрик Маскин (Eric Maskin) начал implementation theory — то есть, собственно, mechanism design: как сделать такой протокол, чтобы он обладал нужными свойствами.

А потом Роджер Майерсон (Roger Myerson) применил это всё к аукционам и окончательно оформил поле деятельности.

За это им всем троим и дали Нобелевскую премию 2007 года по экономике.

Но ещё раньше - в 1994 - премию дали Нэшу за разработку теории игр, которая, конечно, будет ключевой и ляжет в основу для всей этой науки.

1.3 Реальные применения дизайна механизмов

- Как известно, интернет-компании (Google, Yahoo) зарабатывают большей частью на рекламе, которая продаётся через систему аукционов, использующую последние достижения дизайна механизмов.
- Ebay — крупнейшая система интернет-аукционов.
- Общественно полезные работы — нужно максимизировать social welfare, но участники всё равно эгоистичные.
- Налогообложение — какую систему налогообложения ввести, чтобы максимизировать доход государства и social welfare?
- Аукционы на радиочастоты (3G auctions).

Есть и менее прямые и очевидные примеры применений, например, компьютерные распределённые системы:

- real-time scheduling — к распределённой системе приходят всё новые и новые задачи (заранее неизвестные), нужно как можно больше задач решить в срок;
- *Nobel powered* BitTorrent client — как сделать так, чтобы участникам p2p-сети было выгодно делиться файлами, максимизируя при этом суммарную доступность файлов сети?

2 Теория игр: примеры

2.1 Дилемма заключённого

Дилемма заключённого (prisoner's dilemma) — классический пример из теории игр. Двоим заключённым предлагают признаться в преступлении и заложить своего сообщника. Реальных доказательств у обвинения нет, поэтому:

- если оба промолчат, то оба отсидят по полгода за другие грешки;
- если оба признаются, то обоим за примерное поведение дадут по два года;
- если один признается, а другой нет, то признавшегося за сотрудничество отпустят, а упорствующему впаяют по полной, лет десять.

Держать связь заключённые не могут. Как же поступить каждому из них?

Вот какая получается матрица возможных стратегий:

	Промолчать	Сознаться
Промолчать	(0.5, 0.5)	(10, 0)
Сознаться	(0, 10)	(2, 2)

Вне зависимости от выбора первого заключённого, второму в любом случае выгоднее признаться! Получается, что для каждого из них «Сознаться» — доминантная стратегия, и в результате они будут сидеть по 2 года, а не по 0.5.

Реальный пример, в котором возникает именно дилемма заключённого: две фирмы производят один и тот же продукт (других фирм на рынке этого продукта нет). Если рекламы не будет вообще, у них будет одно распределение доходов. Если они обе будут активно рекламироваться, то реклама «взаимно сократится», и относительное потребление их продуктов не изменится, а деньги на рекламу будут потрачены. Но если одна фирма не будет рекламироваться, а вторая будет, то та, что будет, получит большую прибыль от резко увеличившейся доли рынка.

2.2 Трагедия общин

Пример, известный ещё из Фукидида и Аристотеля, возникает, когда у нескольких игроков на рынке есть некий общий ресурс. Выгоды от его использования индивидуальны, а затраты на использование общие, поэтому все пытаются максимизировать своё собственное использование ресурса, и он истощается.

Классическая постановка: на пастбище пасут овец несколько местных овцеводов. Пастбище общее и бесплатное, а каждая дополнительная овца приносит овцеводу прибыль. Поэтому все начинают разводить всё больше и больше овец, и пастбище окончательно вытаптывается. Однако при этом каждый овцевод полностью рационален, потому что лично для него дополнительная овца значит гораздо больше, чем дополнительный ущерб пастбищу от одной овцы.

Такие примеры возникают всё время, где есть общие ресурсы, которые трудно разделить: в загрязнении окружающей среды, использовании воды и воздуха, вырубке лесов, охоте, рыболовстве и т.п. Решение может заключаться только в том, чтобы построить некий общественный механизм (при помощи государства), например, механизм налогообложения или квотирования, при котором общий ресурс не истощится. Вопрос, как сделать это максимально эффективно, — предмет теории механизмов.

2.3 Парадокс аукциона за доллар

Это пример того, к чему может привести дизайн хитрых механизмов. Рассмотрим такой аукцион: лот — один доллар, участники могут перебивать цены друг друга, давший максимальную цену платит её и получает доллар. Но при этом максимальные объявленные цены должны будут уплатить *все* участники аукциона, а не только победитель.

Рассмотрим случай, когда участники действуют рационально.

Первый участник, желая заработать 99 центов, объявляет цену в один цент. Второй перебивает её двумя центами, третий — тремя... Тут первый решает, что заработать 96 центов куда лучше, чем потерять один, и объявляет цену в 4 цента. И так далее. Но рано или поздно цена достигнет 98 центов (пусть такую цену дал первый участник). Второй участник, желая заработать цент, даёт цену в 99 центов. Но для первого даже остаться в нуле гораздо лучше, чем потерять те 98, которые он уже объявлял! И он ставит 100 центов за доллар. А второй... ставит 101!

Адекватного решения у этого парадокса нет. Собственно, и «парадокса» нет — у игры нет равновесия, и игроки могут в конце концов отдать хитрому аукционеру все свои деньги. С другой стороны, конечно, «рациональность» игроков в этом аукционе тоже под вопросом: когда игрок решает, что выгоднее — потерять 98 центов или получить доллар за 100 центов, вторая альтернатива не равна нулю, а должна принимать во внимание вероятность того, что его оппонент не остановится и сделает новую ставку. Ожидание выигрыша составляет бесконечный расходящийся ряд потерь.

2.4 Winner's curse

Рассмотрим такую простую ситуацию: есть аукцион, на торги выставлен товар, у каждого участника своё мнение о ценности товара. Участники делают ставки, исходя из своих понятий о ценности. Выигрывает тот, кто сделал самую большую ставку.

Предположим, что мнения участников распределены более-менее нормально вокруг истинной стоимости (т.е. точно её участники не знают, есть отклонения и в большую, и в меньшую сторону). Это нормальная ситуация, например, для аукционов на участки, с которых можно потом качать нефть: информация о количестве нефти общедоступна, но неточна. Тогда, понятное дело, отклонения от настоящей цены будут и в большую, и в меньшую сторону.

Но победит участник с максимальным отклонением в плюс! Иначе говоря, если ты победил на этом аукционе, сам факт твоей победы означает, что ты переплатил. В этом и заключается Winner's curse «парадокс».

2.5 Парадокс Браесса

Рассмотрим две точки, Start и Finish, между которыми есть два пути, проходящие через точки A и B. Если машина едет по незаполненной трассе, она едет со скоростью 100 км/ч. Если трасса заполнилась, то скорость передвижения падает до $\frac{\text{пропускная способность}}{\text{кол-во автомобилей}}$. Водители всё знают и выбирают оптимальный для себя маршрут. Понятно, что в этой симметричной ситуации водители будут выбирать менее загруженную трассу (когда они заполнятся). Пусть проехать должны 2500 машин. Тогда 1250 из них поедут по одной дороге, 1250 — по другой. Путь каждого водителя занимает 75 минут.

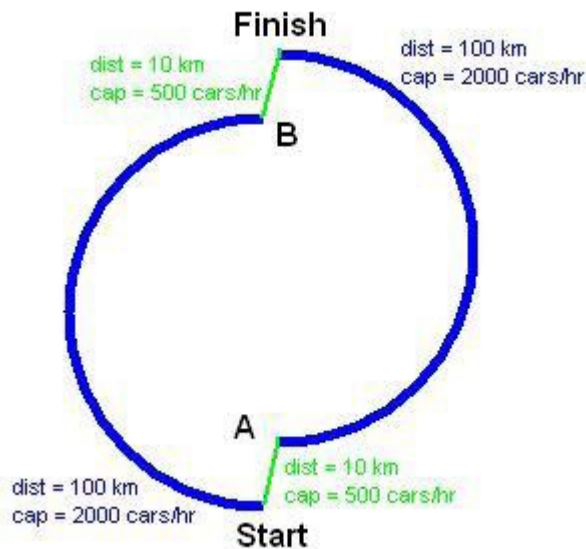


Рис. 1. Парадокс Браесса: до введения новой короткой дороги.

Но вдруг государство решило, что надо бы людям помочь, и построило новую короткую дорогу между *A* и *B*. Эта дорога длиной 60 км супротив 100 км. *Старые дороги никто не закрывает, у водителей просто появляется новый выбор.* Если рассмотреть старое равновесие (1250 на 1250), то при появлении новой дороги по ней ехать будет выгоднее. Новое равновесие (когда все пути одинаковы) достигается, когда из 2500 машин 1500 едут по новой дороге, а по старым — по 500. При этом время в пути окажется равным 84 минутам! Оказывается, что, просто расширив спектр возможностей водителей, мы перевели систему из более эффективного равновесия в менее эффективное. При этом каждый водитель по отдельности действовал рационально: выбирал, где быстрее.



Рис. 2. Парадокс Браесса: после введения новой короткой дороги.

Новая дорога могла бы быть и на пользу, но только если бы в пунктах Start и A сидели регулировщики и распределяли потоки как надо. Это называется *price of anarchy*: иногда регулируемый рынок действительно функционирует эффективнее, чем управляемый лишь невидимой рукой.

Упражнение: Каково оптимальное время проезда в этой системе с 2500 машинами, если регулировщики работают оптимальным образом?

Решение: Один из оптимальных методов регулирования: разделить машины поровну на путях Start-B-Finish и Start-A-Finish и не пускать машины по пути Start-A-B-Finish. Тогда время движения составит 75 минут.

3 Основные концепции теории игр

3.1 Агенты, стратегии, функция полезности

Постановка задачи:

- В игре участвуют *агенты*.
- У игры есть различные *исходы*.
- У каждого агента есть некий набор *действий*, которые он может предпринимать.

Поставим задачу чуть формальнее.

Во-первых, введём *тип агента* $\theta_i \in \Theta$ для i -го агента (об этом ниже). У игры есть набор исходов \mathcal{O} , и для каждого агента каждый исход означает какую-то прибыль; так появляется *функция полезности* (utility function) $u_i(o, \theta_i)$ для типа θ_i и исхода o . Агент i предпочитает исход o_1 исходу o_2 , если $u_i(o_1, \theta_i) > u_i(o_2, \theta_i)$.

Стратегия агента — это план, который полностью описывает его поведение во всех возможных состояниях окружающего мира. Через Σ_i будем обозначать множество стратегий агента i , через $s_i(\theta_i) \in \Sigma_i$ — его стратегию. Стратегии бывают *чистые* (pure) и *смешанные* (mixed); чистые стратегии жёстко задают поведение в каждом состоянии окружающего мира, смешанные задают распределения вероятностей на множестве возможных действий агента.

В аукционе возрастающей цены состояние мира для агента полностью описывается парой (p, x) , где p — текущая цена, а бит x показывает, является ли агент в текущий момент лидером аукциона. Пусть у агента есть своя (скрытая) оценка лота v , и он готов заплатить любую сумму, которая была бы меньше v . Тогда так называемая best response strategy $s_{BR}(v)$ описывается следующим образом:

$$s_{BR}(p, x, v) = \begin{cases} p, & \text{если } x = 0 \text{ и } p < v, \\ \text{сидеть молча,} & \text{в противном случае.} \end{cases}$$

Здесь b (от слова bid) — это ставка, которую должен сделать агент.

Понятно, что функцию полезности можно с конкретных исходов продолжить на целые стратегии. Если N агентов имеют фиксированные стратегии (s_1, \dots, s_N) , то функция полезности

$$u_i(s_1, \dots, s_N, \theta_i)$$

будет просто равна функции полезности $u_i(o, \theta_i)$ на исходе o , который однозначно задаётся этими стратегиями. Рассмотрим тот же аукцион, в котором участвуют два агента, и оба исповедуют best response strategy. Для агента 2 ценность лота $v_2 = 1$, для агента 1 она равна v_1 . Тогда функция полезности для первого агента будет равна

$$u_1(s_{BR,1}(v_1), s_{BR,2}(1)) = \begin{cases} v_1 - (1 + \epsilon), & \text{если } v_1 > 1, \\ 0, & \text{в противном случае,} \end{cases}$$

где ϵ — минимальное увеличение цены в аукционе.

3.2 Равновесие Нэша

Каждый агент пытается максимизировать свою собственную прибыль. Он решает задачу оптимизации, добиваясь оптимальной стратегии. И в результате система оказывается в каком-нибудь состоянии. Мы будем рассматривать возможные определения *равновесного состояния* системы, к которому она может придти после решения каждым агентом своей локальной задачи.

Обозначим через $\mathbf{s} = (s_1, \dots, s_N)$ профиль всех стратегий участников. Через $\mathbf{s}_{-i} = (s_1, \dots, s_{i-1}, s_{i+1}, \dots, s_N)$ обозначим стратегии всех участников, *кроме* i . Введём также аналогичные обозначения θ и θ_{-i} для типов агентов.

Ключевое понятие всей теории игр — *равновесие Нэша* (Nash equilibrium).

Определение 1. Профиль стратегий \mathbf{s} находится в равновесии Нэша, если каждый агент при данных стратегиях других агентов выбирает для себя оптимальную стратегию:

$$\forall s'_i \neq s_i \quad u_i(s_i(\theta_i), \mathbf{s}_{-i}(\theta_{-i}), \theta_i) \geq u_i(s'_i(\theta_i), \mathbf{s}_{-i}(\theta_{-i}), \theta_i).$$

В дилемме заключённого только профиль (Сознаться, Сознаться) находится в равновесии Нэша — преступнику всегда выгоднее сознаться, чем промолчать. Бывают игры с несколькими равновесиями Нэша. Бывают игры, где нет равновесий Нэша для чистых стратегий. Но оно всегда есть в смешанных стратегиях.

Доказательство последнего факта следует из теоремы Какутани о неподвижной точке.

Теорема 1. Пусть S — непустое выпуклое компактное подмножество евклидова пространства \mathbb{R}^n , а $\phi : S \rightarrow 2^S$ — многозначная функция на S с замкнутым графиком, такая, что множество $\phi(x)$ непусто и замкнуто для всех $x \in S$. Тогда у ϕ есть неподвижная точка: $\exists x : x \in \phi(x)$.

Упражнение: Доказать, что из теоремы Какутани следует существование равновесия Нэша в играх со смешанными стратегиями.

Решение: Каждая смешанная стратегия есть распределение вероятностей на множестве возможных действий агента, поэтому сумма этих вероятностей равна единице. Но в n -мерном евклидовом пространстве симплекс

$$a = \{(a_1, a_2, \dots, a_n) \mid a_i \geq 0, \sum a_i = 1, i = 1..n\}$$

является выпуклым компактным множеством. Выигрыш игрока в игре со смешанными стратегиями есть математическое ожидание вида:

$$G_i(a_1, a_2, \dots, a_n) = \sum_{i_1=1}^{m_1} \sum_{i_2=1}^{m_2} \dots \sum_{i_n=1}^{m_n} g_i(i_1, i_2, \dots, i_n) a_1^{i_1} a_2^{i_2} \dots a_n^{i_n}$$

Эта функция является линейной и непрерывной по a , при фиксированных остальных аргументах. Следовательно, по теореме Какутани, у этой функции будет неподвижная точка. Это и означает существование равновесия по Нэшу в играх со смешанными стратегиями.

Говорят, в 1949 году Нэш рассказал фон Нейману про равновесие для смешанных стратегий. Фон Нейман в своём стиле ответил: «Это, знаете ли, тривиально. Это же всего лишь теорема о неподвижной точке». Как мы уже знаем, потом Нэшу за это дали Нобелевскую премию.

Равновесие Нэша — фундаментальное понятие, но оно не всегда применимо. Например, оно много чего предполагает о доступной агентам информации. Нужно, чтобы каждый агент знал структуру игры полностью, знал, что другие знают, знал, что все действуют рационально, и, более того, знал, что все выберут одно и то же равновесие Нэша (а их может быть несколько).

3.3 Доминантные стратегии и аукцион Викри

Агент может и не быть уверен, что все остальные всё знают и непременно выберут равновесие Нэша. Но если у него есть *доминантная стратегия*, ему всё равно.

Определение 2. *Стратегия s_i называется доминантной, если она (слабо) максимизирует ожидаемую прибыль агента для всех возможных стратегий других агентов:*

$$\forall s'_i \neq s_i, \mathbf{s}_{-i} \in \Sigma_{-i} \quad u_i(s_i, \mathbf{s}_{-i}, \theta_i) \geq u_i(s'_i, \mathbf{s}_{-i}, \theta_i).$$

Сейчас мы рассмотрим первый пример нетривиального дизайна механизмов — *аукцион Викри* (Vickrey auction). Это аукцион, проводящийся по схеме sealed-bid: участники подают свои заявки в конвертах, потом их вскрывают, и объект продаётся тому, кто предложил самую высокую цену. Например, так обычно проводят тендеры.

Sealed-bid highest-price. Что выгодно делать участнику со скрытой ценностью v , если ему продадут вещь по той цене, которую он запросит? Это довольно сложная задача: если его скрытая ценность максимальна из всех участников, ему нужно сделать заявку больше, чем у следующего за ним, но желательно только чуть-чуть больше, чтобы максимизировать свою прибыль. В результате на самом деле никому не лучше — и продавец не максимизирует доход, и social welfare тоже страдает. Мы потом более подробно проанализируем этот случай.

Sealed-bid second-price. В этом типе аукциона (который и называется аукционом Викри) по-прежнему продают тому, кто больше предложил... но продают по цене, которую предложил второй сверху участник! Оказывается, что в нём участникам выгодно говорить правду о своей скрытой ценности! Давайте проверим, что $b_i(v_i) = v_i$ — это действительно доминантная стратегия. Ожидаемая полезность стратегии $b_i(v_i) = v_i$ равна

$$u_i(b_i, b', v_i) = \begin{cases} v_i - b', & \text{если } b_i > b', \\ 0, & \text{в противном случае,} \end{cases}$$

где b' — это наивысшая ставка среди всех остальных агентов.

- Если $b' < v_i$, то оптимальна любая ставка $b_i \geq b'$ (вещь ведь всё равно продадут по цене b').

- Если $b' \geq v_i$, то, опять же, оптимальна любая ставка $b_i \leq v_i$ (всё равно не продадут).
- Ставка $b_i = v_i$ подходит в оба случая и поэтому является доминантной стратегией.

Мы только что на пальцах доказали, что в аукционах Викри каждому участнику выгодно говорить правду. Это очень важное свойство механизмов — *правдивость* (truthfulness). Мы позже увидим, что на самом деле можно ограничиться только правдивыми механизмами.

Оказывается, что доминантные стратегии гораздо удобнее для агентов: им уже не надо ничего предполагать о других агентах, они могут смело пользоваться доминантной стратегией. Поэтому в дизайне механизмов гораздо приятнее получить механизм с доминантными стратегиями у каждого агента, чем просто равновесие Нэша.

3.4 Равновесие по Байесу-Нэшу

Возвращаемся к *типам* агентов; теперь мы предположим, что агент не знает наверняка, каковы типы других агентов, то есть каковы у них функции полезности. Но при этом он знает выплаты для каждого возможного типа, и у него есть некоторое априорное распределение $F(\theta)$ на типах для каждого из других агентов. И, конечно, он пытается максимизировать математическое ожидание своей прибыли в равновесии со такими же оптимизирующими стратегиями других агентов.

Определение 3. *Профиль стратегий \mathbf{s} находится в равновесии по Байесу-Нэшу (Bayesian-Nash equilibrium), если каждый агент при известном ему распределении $F(\theta)$ на типах других агентов выбирает для себя оптимальную стратегию: $\forall s'_i \neq s_i$*

$$\mathbf{E}_{F(\theta)} u_i(s_i(\theta_i), \mathbf{s}_{-i}(\theta_{-i}), \theta_i) \geq \mathbf{E}_{F(\theta)} u_i(s'_i(\theta_i), \mathbf{s}_{-i}(\theta_{-i}), \theta_i).$$

То есть стратегия агента оптимальна по распределению типов других агентов. В одном конкретном эксперименте вполне возможно, что он будет выбирать неоптимальное поведение.

Равновесие по Байесу-Нэшу обобщает обычное — оно делает более естественные предположения о знаниях агентов. Для каждого фиксированного типа $\bar{\theta}_i$ оно тоже должно быть оптимальным: $\forall s'_i \neq s_i$

$$\begin{aligned} \mathbf{E}_{F(\theta)} [u_i(s_i(\bar{\theta}_i), \mathbf{s}_{-i}(\theta_{-i}), \theta_i) \mid \bar{\theta}_i] &\geq \\ &\geq \mathbf{E}_{F(\theta)} [u_i(s'_i(\bar{\theta}_i), \mathbf{s}_{-i}(\theta_{-i}), \theta_i) \mid \bar{\theta}_i]. \end{aligned}$$

Но у него есть другие недостатки равновесия Нэша: например, оно не единственно.

Равновесие по Байесу-Нэшу получить лучше, чем обычное, но доминантные стратегии всё равно ещё лучше. В итоге мы ввели и рассмотрели три типа равновесий, которые могут возникнуть в наших механизмах. Перейдём собственно к дизайну.

4 Основные понятия дизайна механизмов

4.1 Функция социального выбора и механизмы

Суть задачи дизайна механизмов заключается в следующем: мы хотим построить механизм, при котором то или иное равновесное состояние системы будет оптимальным относительно той или иной цели. Для этого нужно сначала определить, какая же у нас цель.

Определение 4. Функция социального выбора $f : \Theta_1 \times \dots \times \Theta_N \rightarrow \mathcal{O}$ — это функция, выбирающая тот или иной желаемый результат $f(\theta)$ при данных типах $\theta = (\theta_1, \dots, \theta_N)$.

Функция социального выбора — это то, чего нам бы хотелось получить от механизма, который мы разрабатываем. Но при этом каждый агент будет максимизировать свою собственную прибыль. Надо это примирить.

Теперь наконец можно определить, что же такое *механизм*.

Определение 5. Механизм $\mathcal{M} = (\Sigma_1, \dots, \Sigma_N, g)$ состоит из набора стратегий Σ_i для каждого агента и функция исходов $g : \Sigma_1 \times \dots \times \Sigma_N \rightarrow \mathcal{O}$, которое определяет исход, предусмотренный механизмом для данного профиля стратегий $\mathbf{s} = (s_1, \dots, s_N)$.

Можно проанализировать тот или иной механизм и понять, где у него точки равновесия. При этом может оказаться, что механизм *реализует* ту или иную функцию социального выбора.

Определение 6. Механизм $\mathcal{M} = (\Sigma_1, \dots, \Sigma_N, g)$ реализует функцию социального выбора $f(\theta)$, если для всех $\theta = (\theta_1, \dots, \theta_N) \in \Theta_1 \times \dots \times \Theta_N$

$$g(s_1^*(\theta_1), \dots, s_N^*(\theta_N)) = f(\theta),$$

где профиль стратегий (s_1^*, \dots, s_N^*) находится в равновесии по отношению к игре, индуцированной \mathcal{M} .

Под «равновесием» можно понимать равновесие по Нэшу, по Байесу–Нэшу, по доминантным стратегиям. Обычно нас интересует максимально сильное из возможных равновесий.

Давайте попробуем построить тривиальный механизм, который мог бы реализовывать всевозможные функции социального выбора. Мы просто спросим у каждого агента, какой у него тип (ответы на этот вопрос будут возможными стратегиями агентов), а потом в качестве функции исходов возьмём функцию социального выбора: $g(\theta) = f(\theta)$. Кажется бы, всё работает... но ведь агенты не обязаны говорить нам правду! Агенты будут максимизировать свой доход, сообщая тот тип, который выгоднее, решая (для байесовского равновесия Нэша) задачу оптимизации

$$\max_{\theta' \in \Theta_i} \mathbf{E}_{\theta_{-i}} u_i(\theta', \mathbf{s}_{-i}(\theta_{-i}), \theta_i).$$

Нам нужно построить механизм так, чтобы решение этой задачи для агентов сошлось с желаемым; в частности, в данном случае нам нужно было бы реализовать *правдивый* механизм, при котором агентам было бы выгодно

сообщать свои настоящие типы. Один такой пример мы уже разбирали — это был аукцион Викри.

Есть ряд свойств функций социального выбора, которые могут очень помочь нам при дизайне, а также гарантировать много полезных свойств механизмов, их реализующих. Сейчас мы их рассмотрим и введём (естественные) ограничения на агентов.

4.2 Оптимальность по Парето

Определение 7. *Функция социального выбора $f(\theta)$ называется оптимальной по Парето, если для всякого набора типов $\theta = (\theta_1, \dots, \theta_i)$ и всякого исхода $\sigma' \neq f(\theta)$*

$$u_i(\sigma', \theta_i) > u_i(f(\theta), \theta_i) \quad \Rightarrow \quad \exists j : u_j(\sigma', \theta_j) < u_j(f(\theta), \theta_j).$$

Оптимальность по Парето значит, что если кому-то стало лучше, чем в предлагаемом функцией f варианте, то кому-то другому обязательно стало хуже. То есть нельзя монотонно улучшить дела сразу всех агентов.

Рассмотрим множество исходов $\mathcal{O} = \{x, y, z\}$ и предположим, что действуют два агента. У первого агента ровно один тип, $\Theta_1 = \{\theta_1\}$, и у этого типа структура предпочтений такова: $x >_1 y >_1 z$. А у второго агента два разных типа $\Theta_2 = \{\theta_2^a, \theta_2^b\}$, и вот их структура предпочтений:

$$z >_2^a y >_2^a x, \quad y >_2^b x >_2^b z.$$

Мы пытаемся реализовать эффективную по Парето (проверьте!) функцию социального выбора:

$$f(\theta_1, \theta_2^a) = y, \quad f(\theta_1, \theta_2^b) = x.$$

Если мы захотим просто спросить у каждого агента его тип, второму будет выгодно соврать: при типе θ_2^b ему будет выгодно сказать, что он θ_2^a и получить в результате исход y , а не x .

Можно теперь ввести вполне естественное определение оптимального механизма.

Определение 8. *Механизм называется оптимальным по Парето, если он реализует оптимальную по Парето функцию социального выбора.*

Это определение на самом деле предполагает, что исход окажется оптимальным по Парето уже для конкретных типов агентов, в итоге, апостериори, *ex post*. Можно рассматривать оптимальность по Парето *ex ante*, когда нет исхода, который бы *в ожидании* строго предпочёл один агент и нестрого — все остальные. Получится более слабое определение.

Чуть отвлечёмся и обобщим предыдущий разговор. Вообще говоря, в литературе о дизайне механизмов есть три разных временных постановки.

- *Ex ante* — до выбрасывания исходов. *Ex ante* агенты знают только распределения (все, включая своё собственное). Информация у всех агентов одинаковая.

- Interim — после выбрасывания исходов, для каждого агента. То есть ситуация рассматривается с точки зрения одного агента, который уже знает свой тип, но не знает типы других агентов (а распределения знает). Информация теперь у агентов разная — каждый знает свой тип.
- Ex post — после того, как типы всех агентов стали известны.

Иначе говоря, о равновесиях или ограничениях можно говорить в трёх случаях.

- Ex ante — в терминах распределений типов агентов.
- Interim — в терминах распределений типов агентов и одного конкретного типа одного агента.
- Ex post — в терминах вектора типов всех агентов.

4.3 Предположения об агентах

Определение 9. Квазилинейная функция полезности агента i с типом θ_i имеет вид

$$u_i(o, \theta_i) = v_i(a, \theta_i) - p_i,$$

где исход o определяет выбор $a \in \mathcal{K}$ из дискретного множества \mathcal{K} и выплату p_i , производимую агентом.

У агента с квазилинейными предпочтениями есть *функция оценки* (valuation function) $v_i(a)$, $a \in \mathcal{K}$. Например, в аукционе, где продаётся одна вещь, $\mathcal{K} = \{0, 1\}$ — агент либо получит эту вещь, либо не получит. А p_i в этом случае — выплата агента продавцу. Это достаточно естественное предположение в случае аукциона.

Есть ещё одно предположение, которое в жизни часто не выполняется. Мы для простоты предполагаем, что агенты нейтральны к риску (risk-neutral agents). То есть если агент может получить возможность с вероятностью p получить вещь ценой в \$100, то он радостно заплатит за это \$100 p . В жизни часто встречаются осторожные агенты (risk-averse agents). Позже мы их рассмотрим и увидим, что меняется в этом случае.