

ДИСЦИПЛИНА

Интеллектуальные системы и технологии

(полное наименование дисциплины без сокращений)

ИНСТИТУТ

информационных технологий

КАФЕДРА

корпоративных информационных систем

полное наименование кафедры

ВИД УЧЕБНОГО
МАТЕРИАЛА

Лекция

(в соответствии с пп. 1-11)

ПРЕПОДАВАТЕЛЬ

Демидова Лилия Анатольевна

(фамилия, имя, отчество)

СЕМЕСТР

1 семестр (осенний), 2024 – 2025 учебный год

(семестр обучения, учебный год)

ЛЕКЦИЯ 1

ПОИСК АССОЦИАТИВНЫХ ПРАВИЛ

1. Основные понятия

Одним из распространенных методов Data Mining является аффинитивный анализ (affinity analysis), целью которого является исследование взаимной связи между событиями, которые происходят совместно.

Анализ рыночной корзины (market basket analysis) является разновидностью аффинитивного анализа.

Целью анализа рыночной корзины является обнаружение ассоциаций между различными событиями посредством поиска ассоциативных правил, позволяющих количественно описать взаимную связь между двумя и более событиями [1, 2].

В контексте задачи поиска ассоциативных правил под транзакцией понимают множество событий, которые произошли одновременно. Таким событием, например, может быть покупка того или иного товара (предмета). Тогда рыночная корзина – это набор товаров (предметов), приобретенных в рамках одной отдельно взятой транзакции.

Транзакционная база данных (база данных транзакций, БД транзакций, transaction database) – двумерная таблица, содержащая номер транзакции и перечень покупок, приобретенных во время этой транзакции.

Пусть число транзакций в базе – N , и каждая транзакция D_i представляет собой некоторый набор предметов ($i = \overline{1, N}$).

В таблице 1 приведен пример набора транзакций (базы данных транзакций) из 10 предметов (товаров). Такой набор может быть сформирован, например, на основе чеков о покупке товаров (предметов) в магазине «Канцтовары». Следует отметить, что этот набор содержит очень мало транзакций (всего 20), что не позволяет считать выводы, сделанные на его основе объективными. Однако этот набор позволяет в полной мере продемонстрировать работу по поиску ассоциативных правил.

Любое ассоциативное правило состоит из двух наборов предметов, называемых условием (antecedent) A и следствием (consequent) B [1, 2].

Ассоциативное правило записывается в виде $A \rightarrow B$, что читается как «из A следует B » или «Если условие A , то следствие B ».

Таблица 1. Пример набора транзакций на основе чеков магазина «Канцтовары»

| Номер чека | Список предметов в чеке | | | | | | | |
|------------|-------------------------|----------|----------|----------|----------|---------|-------|--------|
| 1 | Карандаш | Ластик | Альбом | Линейка | Тетрадь | | | |
| 2 | Ручка | Тетрадь | Карандаш | Ластик | Пенал | Линейка | | |
| 3 | Карандаш | Ластик | Ручка | | | | | |
| 4 | Краски | Альбом | Палитра | Карандаш | Кисть | Ластик | | |
| 5 | Карандаш | Ластик | Тетрадь | Ручка | | | | |
| 6 | Кисть | Краски | Палитра | Альбом | | | | |
| 7 | Краски | Альбом | Карандаш | Ластик | Тетрадь | | | |
| 8 | Карандаш | Ластик | Линейка | | | | | |
| 9 | Тетрадь | Ручка | Пенал | Ластик | Карандаш | | | |
| 10 | Ластик | Карандаш | Тетрадь | Пенал | Кисть | Альбом | Ручка | Краски |
| 11 | Линейка | Карандаш | Ластик | Тетрадь | | | | |
| 12 | Пенал | Ластик | Тетрадь | Карандаш | | | | |
| 13 | Ручка | Пенал | | | | | | |
| 14 | Ручка | Альбом | Линейка | | | | | |
| 15 | Альбом | Тетрадь | Карандаш | Ластик | | | | |
| 16 | Карандаш | Ластик | Ручка | Тетрадь | Линейка | | | |
| 17 | Краски | Палитра | Альбом | Ластик | Карандаш | Тетрадь | | |
| 18 | Карандаш | Краски | Альбом | Ластик | | | | |
| 19 | Пенал | Линейка | | | | | | |
| 20 | Карандаш | Ручка | Тетрадь | Ластик | | | | |

Ассоциативные правила описывают связь между наборами предметов, соответствующими условию A и следствию B .

Эта связь характеризуется такими показателями, как поддержка (support) S и достоверность (confidence) C [1, 2].

Поддержка S – доля транзакций, содержащих одновременно условие A и следствие B :

$$S(A \rightarrow B) = P(A \cap B) = \frac{\text{число транзакций, содержащих } A \text{ и } B}{\text{общее число транзакций}}.$$

Достоверность C – доля транзакций, содержащих и условие A , и следствие B , к числу транзакций, содержащих только условие A .

$$C(A \rightarrow B) = P(B | A) = P(A \cap B) / P(A) = \\ = \frac{\text{число транзакций, содержащих } A \text{ и } B}{\text{число транзакций, содержащих только } A}.$$

Рассмотрим ассоциацию *Карандаш* \rightarrow *Ластик*.

Очевидно, что эти продукты (товары) хорошо совместимы.

Так как число транзакций, содержащих как *Карандаш*, так и *Ластик*, равно 16, а общее число транзакций – 20 (таблица 2), то поддержка этой ассоциации будет:

$$S(\text{Карандаш} \rightarrow \text{Ластик}) = 16/20 = 4/5 = 0,8.$$

Так как число транзакций, содержащих только *Карандаш* (условие), равно 16 (таблица 2), то достоверность этой ассоциации будет:

$$C(\text{Карандаш} \rightarrow \text{Ластик}) = 16/16 = 1.$$

Все наблюдения, содержащие *Карандаш*, содержат и *Ластик*, поэтому такая ассоциация может рассматриваться как **правило**.

Теперь рассмотрим ассоциацию *Кисть* \rightarrow *Ластик*.

Очевидно, что у этих продуктов (товаров) совместимость меньше, чем у продуктов (товаров) *Карандаш* и *Ластик*.

Так как число транзакций, содержащих как *Кисть*, так и *Ластик*, равно 2, а общее число транзакций – 20 (таблица 3), то поддержка этой ассоциации будет:

$$S(\text{Кисть} \rightarrow \text{Ластик}) = 2/20 = 1/10 = 0,1.$$

Так как число транзакций, содержащих только *Кисть* (условие), равно 3 (таблица 3), то достоверность этой ассоциации будет:

$$C(\text{Кисть} \rightarrow \text{Ластик}) = 2/3 = 0,67.$$

**Таблица 2. Пример набора транзакций на основе чеков магазина «Канцтовары»
с выделенными предметами Карандаш и Ластик**

| Номер чека | Список предметов в чеке | | | | | | | |
|------------|-------------------------|----------|----------|----------|----------|---------|-------|--------|
| 1 | Карандаш | Ластик | Альбом | Линейка | Тетрадь | | | |
| 2 | Ручка | Тетрадь | Карандаш | Ластик | Пенал | Линейка | | |
| 3 | Карандаш | Ластик | Ручка | | | | | |
| 4 | Краски | Альбом | Палитра | Карандаш | Кисть | Ластик | | |
| 5 | Карандаш | Ластик | Тетрадь | Ручка | | | | |
| 6 | Кисть | Краски | Палитра | Альбом | | | | |
| 7 | Краски | Альбом | Карандаш | Ластик | Тетрадь | | | |
| 8 | Карандаш | Ластик | Линейка | | | | | |
| 9 | Тетрадь | Ручка | Пенал | Ластик | Карандаш | | | |
| 10 | Ластик | Карандаш | Тетрадь | Пенал | Кисть | Альбом | Ручка | Краски |
| 11 | Линейка | Карандаш | Ластик | Тетрадь | | | | |
| 12 | Пенал | Ластик | Тетрадь | Карандаш | | | | |
| 13 | Ручка | Пенал | | | | | | |
| 14 | Ручка | Альбом | Линейка | | | | | |
| 15 | Альбом | Тетрадь | Карандаш | Ластик | | | | |
| 16 | Карандаш | Ластик | Ручка | Тетрадь | Линейка | | | |
| 17 | Краски | Палитра | Альбом | Ластик | Карандаш | Тетрадь | | |
| 18 | Карандаш | Краски | Альбом | Ластик | | | | |
| 19 | Пенал | Линейка | | | | | | |
| 20 | Карандаш | Ручка | Тетрадь | Ластик | | | | |

**Таблица 3. Пример набора транзакций на основе чеков магазина «Канцтовары»
с выделенными предметами Кисть и Ластик**

| Номер чека | Список предметов в чеке | | | | | | | |
|------------|-------------------------|----------|----------|----------|----------|---------|-------|--------|
| 1 | Карандаш | Ластик | Альбом | Линейка | Тетрадь | | | |
| 2 | Ручка | Тетрадь | Карандаш | Ластик | Пенал | Линейка | | |
| 3 | Карандаш | Ластик | Ручка | | | | | |
| 4 | Краски | Альбом | Палитра | Карандаш | Кисть | Ластик | | |
| 5 | Карандаш | Ластик | Тетрадь | Ручка | | | | |
| 6 | Кисть | Краски | Палитра | Альбом | | | | |
| 7 | Краски | Альбом | Карандаш | Ластик | Тетрадь | | | |
| 8 | Карандаш | Ластик | Линейка | | | | | |
| 9 | Тетрадь | Ручка | Пенал | Ластик | Карандаш | | | |
| 10 | Ластик | Карандаш | Тетрадь | Пенал | Кисть | Альбом | Ручка | Краски |
| 11 | Линейка | Карандаш | Ластик | Тетрадь | | | | |
| 12 | Пенал | Ластик | Тетрадь | Карандаш | | | | |
| 13 | Ручка | Пенал | | | | | | |
| 14 | Ручка | Альбом | Линейка | | | | | |
| 15 | Альбом | Тетрадь | Карандаш | Ластик | | | | |
| 16 | Карандаш | Ластик | Ручка | Тетрадь | Линейка | | | |
| 17 | Краски | Палитра | Альбом | Ластик | Карандаш | Тетрадь | | |
| 18 | Карандаш | Краски | Альбом | Ластик | | | | |
| 19 | Пенал | Линейка | | | | | | |
| 20 | Карандаш | Ручка | Тетрадь | Ластик | | | | |

Невысокая достоверность этой ассоциации дает повод усомниться в том, что она является правилом. При этом поддержка этой ассоциации является совсем маленькой.

При анализе ассоциаций наиболее часто предпочтение отдается тем, которые имеют высокие значения поддержки и достоверности. Хотя в ряде случаев могут выбираться и ассоциации, имеющие только высокую поддержку или только высокую достоверность.

Правила, значения поддержки (достоверности) которых превышают пороговое значение, заданное пользователем, называются *сильными* правилами (strong rules).

Например, при анализе ассоциаций на основе чеков о покупках в магазине с целью формирования рекомендаций по формированию успешного плана его деятельности могут быть интересны ассоциации, для которых значения поддержки и достоверности не ниже пороговых значений в 25% и 75% соответственно.

При анализе ассоциаций с целью выявления мошеннических действий могут быть интересны ассоциации, значение поддержки которых невелико (например, не больше 1%), так как с такими действиями связано малое число транзакций.

Значимость ассоциации может быть вычислена как разность между поддержкой правила в целом и произведением поддержки *только* условия и поддержки *только* следствия.

Если условие и следствие независимы, то поддержка ассоциации примерно соответствует произведению поддержек условия A и следствия B :

$$S_{AB} \approx S_A \cdot S_B.$$

В таком случае можно будет сделать следующий вывод: хотя условие A и следствие B часто встречаются вместе, не менее часто они встречаются и по отдельности.

Например, если товар A встречался в 75 транзакциях из 100, а товар B – в 85, и в 64 транзакциях из 100 они встречаются вместе, то, несмотря на высокое значение поддержки ($S_{AB} = 0,64$), ассоциация $A \rightarrow B$ не обязательно является правилом: товары A и B покупаются независимо друг от друга, но по причине их популярности часто встречаются в одной транзакции.

Поскольку произведение значений поддержек условия A и следствия B вычисляется как $S_A \cdot S_B = 0,75 \cdot 0,85 = 0,6375$, то есть отличается от $S_{AB} = 0,64$ всего на 0,0025, предположение о независимости товаров A и B можно считать достаточно обоснованным.

Если условие A и следствие B независимы, то ассоциация скорее всего не может считаться правилом, даже если она имеет высокие значения поддержки и достоверности.

При оценке значимости ассоциаций применяют объективные и субъективные меры значимости [2]

Объективные меры – поддержка и достоверность. Такие меры значимости могут применяться в любой предметной области.

Субъективные меры – лифт (lift) и левевердж (leverage, плечо, рычаг). Такие меры значимости напрямую связаны с информацией, определяемой контекстом решаемой задачи анализа данных.

Лифт – отношение достоверности ассоциации $A \rightarrow B$ к поддержке следствия B [2]:

$$L(A \rightarrow B) = C(A \rightarrow B) / S(B).$$

Значения лифта, большие чем 1, показывают, что условие чаще появляется в транзакциях, содержащих следствие B , чем в остальных.

Лифт является обобщенной мерой связи двух предметных наборов, фигурирующих в числителе и знаменателе формулы: при значениях лифта, превосходящих 1, связь положительная, при значении лифта, равном 1, связь отсутствует, а при значениях лифта, меньших 1, связь отрицательная.

Рассмотрим ассоциацию *Ластик* \rightarrow *Карандаш* и вычислим для нее лифт.
 $S(\text{Карандаш})=16/20=0,8$; $C(\text{Ластик} \rightarrow \text{Карандаш})=16/16=1$.

Следовательно, $L(\text{Ластик} \rightarrow \text{Карандаш})=1/0,8=1,25$.

Рассмотрим ассоциацию *Ластик* \rightarrow *Кисть* и вычислим для нее лифт.
 $S(\text{Кисть})=3/20=0,15$; $C(\text{Ластик} \rightarrow \text{Кисть})=2/16=0,125$.

Следовательно, $L(\text{Ластик} \rightarrow \text{Кисть})=0,125/0,15=0,833$.

Для первой ассоциации лифт больше 1, для второй – меньше 1, значит, ластик больше влияет на покупку карандаша, чем на покупку кисти.

Следует отметить, что лифт не всегда является удачной мерой значимости ассоциации.

Так, например, ассоциация с меньшим значением поддержкой и большим значением лифта может быть менее значимой, чем ассоциация с большим значением поддержки и меньшим значением лифта, так как что вторая ассоциация используется для большего числа транзакций. Следовательно, увеличение числа транзакций приводит к увеличению связи между условием и следствием ассоциации.

Левередж – разность между наблюдаемой частотой, с которой условие и следствие появляются совместно, и произведением частот появления условия и следствия по отдельности, то есть разность между поддержкой ассоциации и произведением поддержек условия и следствия по отдельности [2]:

$$Lev(A \rightarrow B) = S(A \rightarrow B) - S(A) \cdot S(B).$$

Рассмотрим ассоциации *Карандаш* \rightarrow *Ластик* и *Тетрадь* \rightarrow *Ластик*.

Эти ассоциации имеют одинаковую достоверность $C(A \rightarrow B) = 1$, т.к. карандаш (всего 16 раз) (таблица 2) и тетрадь (всего 12 раз) (таблица 4) всегда продаются вместе с ластиком (16 и 12 соответственно).

Эти ассоциации имеют одинаковые лифты, так как в обеих ассоциациях поддержка следствия $S(\text{Ластик}) = 16/20 = 0,8$ (ластик встретился в транзакциях 16 раз), и, следовательно,

$$L(\text{Карандаш} \rightarrow \text{Ластик}) = L(\text{Тетрадь} \rightarrow \text{Ластик}) = 1/0,8 = 1,25.$$

Вычислим левереджи для этих ассоциаций.

$$S(\text{Карандаш} \rightarrow \text{Ластик}) = 16/20 = 0,8.$$

$$S(\text{Карандаш}) = 16/20 = 0,8. S(\text{Ластик}) = 16/20 = 0,8.$$

$$\text{Следовательно, } Lev(\text{Карандаш} \rightarrow \text{Ластик}) = 0,8 - 0,8 \cdot 0,8 = 0,16.$$

$$S(\text{Тетрадь} \rightarrow \text{Ластик}) = 12/20 = 0,6.$$

$$S(\text{Тетрадь}) = 12/20 = 0,6. S(\text{Ластик}) = 16/20 = 0,8.$$

$$\text{Следовательно, } Lev(\text{Тетрадь} \rightarrow \text{Ластик}) = 0,6 - 0,6 \cdot 0,8 = 0,12.$$

Ассоциация *Карандаш* \rightarrow *Ластик* представляет больший интерес, так как встречается чаще (то есть применяется в большем числе транзакций, чем ассоциация *Тетрадь* \rightarrow *Ластик*).

Таким образом, значимость ассоциации *Карандаш* \rightarrow *Ластик* более высокая.

Иногда вместо лифта используют еще одну субъективную меру – улучшение (improvement), которую вычисляют подобно левереджу, но берут не разность, а отношение между наблюдаемой частотой, с которой условие и следствие появляются совместно, и произведением частот появления условия и следствия по отдельности [2]:

$$I(A \rightarrow B) = \frac{S(A \rightarrow B)}{S(A) \cdot S(B)}.$$

Таблица 4. Пример набора транзакций на основе чеков магазина «Канцтовары»

| Номер чека | Список предметов в чеке | | | | | | | |
|------------|-------------------------|----------|----------|----------|----------|---------|-------|--------|
| 1 | Карандаш | Ластик | Альбом | Линейка | Тетрадь | | | |
| 2 | Ручка | Тетрадь | Карандаш | Ластик | Пенал | Линейка | | |
| 3 | Карандаш | Ластик | Ручка | | | | | |
| 4 | Краски | Альбом | Палитра | Карандаш | Кисть | Ластик | | |
| 5 | Карандаш | Ластик | Тетрадь | Ручка | | | | |
| 6 | Кисть | Краски | Палитра | Альбом | | | | |
| 7 | Краски | Альбом | Карандаш | Ластик | Тетрадь | | | |
| 8 | Карандаш | Ластик | Линейка | | | | | |
| 9 | Тетрадь | Ручка | Пенал | Ластик | Карандаш | | | |
| 10 | Ластик | Карандаш | Тетрадь | Пенал | Кисть | Альбом | Ручка | Краски |
| 11 | Линейка | Карандаш | Ластик | Тетрадь | | | | |
| 12 | Пенал | Ластик | Тетрадь | Карандаш | | | | |
| 13 | Ручка | Пенал | | | | | | |
| 14 | Ручка | Альбом | Линейка | | | | | |
| 15 | Альбом | Тетрадь | Карандаш | Ластик | | | | |
| 16 | Карандаш | Ластик | Ручка | Тетрадь | Линейка | | | |
| 17 | Краски | Палитра | Альбом | Ластик | Карандаш | Тетрадь | | |
| 18 | Карандаш | Краски | Альбом | Ластик | | | | |
| 19 | Пенал | Линейка | | | | | | |
| 20 | Карандаш | Ручка | Тетрадь | Ластик | | | | |

Улучшение $I(A \rightarrow B)$ показывает, полезнее ли ассоциация случайного угадывания. Если улучшение $I(A \rightarrow B)$ больше 1, то это означает, что вероятнее предсказать наличие набора B в следствии с помощью ассоциации, чем угадать случайно.

Такие меры, как лифт $L(A \rightarrow B)$, леввередж $Lev(A \rightarrow B)$ и улучшение $I(A \rightarrow B)$, могут использоваться для последующего сокращения числа рассматриваемых ассоциаций посредством установки порогового значения для используемой меры значимости: ассоциации, у которых значение меры значимости ниже порогового, исключаются из дальнейшего рассмотрения.

2. Алгоритм Apriori

Число возможных ассоциаций растет экспоненциально с увеличением числа предметов, на основе которых могут формироваться ассоциации.

Пусть общее число предметов в базе данных транзакций равно k .

Пусть все ассоциации содержат по одному предмету в условии и в следствии.

В этом случае требуется проанализировать $2 \cdot C_k^2$ ассоциаций.

Пусть, например, $k=100$.

Тогда число возможных ассоциаций равно $2 \cdot C_{100}^2 = 9900$.

Очевидно, что в условии и следствии может содержаться большее число предметов. Следовательно, число возможных ассоциаций будет еще большим.

Исходя из вышесказанного, можно сделать следующий вывод: поиск ассоциативных правил посредством вычисления поддержки и достоверности для всех возможных ассоциаций и сравнения их с заданным пороговым значением будет малоэффективным из-за больших вычислительных затрат.

Обычно при поиске ассоциативных правил используется методика, предполагающая выявление частых предметных наборов. При реализации такой методики анализируются только те ассоциации, которые встречаются часто. При этом задается некоторое пороговое значение, позволяющее отнести предметных набор к частым или к нечастым.

Одним из алгоритмов, основанных на применении методики выявления частых предметных наборов, является алгоритм Apriori [1 – 4].

Этот алгоритм предложили ученые Ракеш Агравал (Rakesh Agrawal) и Рамакришнан Шрикант (Ramakrishnan Srikant) в 1994 году.

Цель алгоритма Apriori – сузить пространство поиска до размеров, обеспечивающих приемлемые временные затраты на поиск ассоциативных правил.

2.1. Выявление частых предметных наборов

В основе алгоритма Apriori лежит понятие частого набора (frequent itemset), который также можно назвать частым предметным набором, для которого определяется частота его появления в базе транзакций.

Под частотой понимается простое число транзакций, в которых содержится рассматриваемый предметный набор.

Частым предметным набором будет тот, который встречается чаще, чем в числе транзакций, описываемом некоторой пороговой частотой Δ .

Кроме того, можно считать, что частый предметный набор – предметный набор с поддержкой, большей заданного порогового значения либо равной этому значению.

Это пороговое значение называется минимальной поддержкой.

Методика поиска ассоциативных правил с использованием частых наборов содержит следующую последовательность шагов.

Шаг 1. Поиск частых предметных наборов.

Шаг 2. Генерация на основе частых предметных наборов ассоциативных правил, удовлетворяющих условиям минимальной поддержки и достоверности.

Алгоритм Apriori использует свойство антимонотонности, чтобы сузить пространство поиска.

Свойство антимонотонности: если предметный набор Z не является частым, то добавление некоторого нового предмета A к набору Z не делает его более частым.

Если набор Z не является частым, то и набор $Z \cup A$ также не будет являться частым.

Использование свойства антимонотонности позволяет значительно уменьшить пространство поиска ассоциативных правил.

Рассмотрим набор транзакций, приведенный в таблице 1, и применим к нему алгоритм Apriori.

Таблица 5. Представление данных из набора транзакций в нормализованном виде

| № чека | Альбом | Карандаш | Кисть | Краски | Ластик | Линейка | Палитра | Пенал | Ручка | Тетрадь |
|------------------|--------|----------|-------|--------|--------|---------|---------|-------|-------|---------|
| 1 | 1 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 1 |
| 2 | 0 | 1 | 0 | 0 | 1 | 1 | 0 | 1 | 1 | 1 |
| 3 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 |
| 4 | 1 | 1 | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 0 |
| 5 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 1 |
| 6 | 1 | 0 | 1 | 1 | 0 | 0 | 1 | 0 | 0 | 0 |
| 7 | 1 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 1 |
| 8 | 0 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 |
| 9 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 1 |
| 10 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 1 | 1 | 1 |
| 11 | 0 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 1 |
| 12 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 1 |
| 13 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 |
| 14 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 |
| 15 | 1 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 |
| 16 | 0 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 1 |
| 17 | 1 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 1 |
| 18 | 1 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 |
| 18 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 |
| 20 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 1 |
| Сумма по столбцу | 9 | 16 | 3 | 6 | 16 | 7 | 3 | 6 | 9 | 12 |

Будем считать частыми такие предметные наборы, которые встречаются в базе транзакций более 5 раз, то есть определим пороговое значение частоты как $\Delta = 5$.

Выполним поиск ассоциативных правил, реализуя следующие шаги.

Шаг 1. Выполним поиск частых однопредметных наборов, представив набор транзакций в нормализованном виде (таблица 5).

На пересечении строки транзакции (строки с номером чека) и столбца с названием предмета поставим число «1», если предмет присутствует в транзакции, и «0» – в противном случае.

Просуммируем значения в каждом столбце для вычисления частоты появления каждого предмета в наборе транзакций.

Запишем частоту появления каждого предмета в наборе транзакций в соответствующую ячейку последней строки таблицы 5.

Как видно из таблицы 5, сумма в каждом столбце, за исключением столбцов с названиями предметов Кисть и Палитра больше порогового значения частоты $\Delta = 5$.

Все предметы, для которых частота появления в наборе транзакций больше или равна $\Delta = 5$, будем считать частыми однопредметными наборами.

Сформируем множество F_1 частых однопредметных наборов:

$F_1 = \{\text{Альбом, Карандаш, Краски, Ластик, Линейка, Пенал, Ручка, Тетрадь}\}.$

Шаг 2. Выполним поиск частых двухпредметных наборов (таблица 6).

В общем случае для формирования множества F_k частых k -предметных наборов реализуют следующий алгоритм:

- создать множество наборов-кандидатов в частные k -предметные наборы посредством связывания множества F_{k-1} с самим собой;
- выполнить редукцию множества F_k , используя свойство антимонотонности.

Оставшиеся предметные наборы формируют искомое множество F_k частых k -предметных наборов.

Так как пороговое значение частоты $\Delta = 5$, то в множество F_2 войдут только те двухпредметные наборы, которые встречаются 5 и более раз (таблица 6).

Информация по найденным частным двухпредметным наборам выделена в таблице 6 жирным шрифтом.

Таблица 6. Двухпредметные наборы-кандидаты в F_2

| Предметный набор | | Частота появления в наборе транзакций |
|------------------|-----------------|--|
| Предмет 1 | Предмет 2 | |
| Альбом | Карандаш | 7 |
| Альбом | Краски | 6 |
| Альбом | Ластик | 7 |
| Альбом | Линейка | 2 |
| Альбом | Пенал | 1 |
| Альбом | Ручка | 2 |
| Альбом | Тетрадь | 5 |
| Карандаш | Краски | 5 |
| Карандаш | Ластик | 16 |
| Карандаш | Линейка | 5 |
| Карандаш | Пенал | 1 |
| Карандаш | Ручка | 7 |
| Карандаш | Тетрадь | 12 |
| Краски | Ластик | 5 |
| Краски | Линейка | 0 |
| Краски | Пенал | 1 |
| Краски | Ручка | 1 |
| Краски | Тетрадь | 3 |

Окончание таблицы 6

| Предметный набор | | Частота появления в наборе транзакций |
|------------------|----------------|--|
| Предмет 1 | Предмет 2 | |
| Ластик | Линейка | 5 |
| Ластик | Пенал | 4 |
| Ластик | Ручка | 7 |
| Ластик | Тетрадь | 12 |
| Линейка | Пенал | 2 |
| Линейка | Ручка | 3 |
| Линейка | Тетрадь | 4 |
| Пенал | Ручка | 4 |
| Пенал | Тетрадь | 4 |
| Ручка | Тетрадь | 6 |

Сформируем множество F_2 частых двухпредметных наборов:

$F_2 \{ \{ \text{Альбом, Карандаш} \}, \{ \text{Альбом, Краски} \}, \{ \text{Альбом, Ластик} \},$
 $\{ \text{Альбом, Тетрадь} \}, \{ \text{Карандаш, Краски} \}, \{ \text{Карандаш, Ластик} \},$
 $\{ \text{Карандаш, Линейка} \}, \{ \text{Карандаш, Ручка} \}, \{ \text{Карандаш, Тетрадь} \},$
 $\{ \text{Краски, Ластик} \}, \{ \text{Ластик, Линейка} \}, \{ \text{Ластик, Ручка} \},$
 $\{ \text{Ластик, Тетрадь} \}, \{ \text{Ручка, Тетрадь} \} \}.$

Шаг 3. Выполним поиск частых трехпредметных наборов (таблица 7).

Для получения множества F_3 свяжем множество F_2 с самим собой.

В общем случае k -предметные наборы являются связываемыми, если у них первые $k - 1$ предметов, следующих в алфавитном порядке, общие.

Так, наборы $\{ \text{Альбом, Карандаш} \}$ и $\{ \text{Альбом, Краски} \}$, для которых $k=2$, чтобы быть связываемыми, должны иметь один ($k-1=1$) общий первый предмет, которым и является Альбом.

В результате связывания получим:

$\{ \text{Альбом, Карандаш} \} + \{ \text{Альбом, Краски} \} = \{ \text{Альбом, Карандаш, Краски} \}.$

Так как пороговое значение частоты $\Delta = 5$, то в множество F_3 войдут только те трехпредметные наборы, которые встречаются 5 и более раз (таблица 7). Информация по найденным частным двухпредметным наборам выделена в таблице 7 жирным шрифтом.

Множество трехпредметных наборов-кандидатов в F_3 сокращается с помощью свойства антимонотонности: поддержка любого набора предметов не может превышать минимальной поддержки любого из его поднаборов.

С увеличением числа предметов в наборе поддержка уменьшается или остается такой же.

В общем случае любой k -предметный набор будет часто встречающимся тогда и только тогда, когда все его $(k - 1)$ -предметные поднаборы будут часто встречающимися.

Рассмотрим набор $\{ \text{Альбом, Карандаш, Краски} \}.$

В нем есть двухпредметные поднаборы: $\{ \text{Альбом, Карандаш} \}$, $\{ \text{Альбом, Краски} \}$, $\{ \text{Карандаш, Краски} \}$, частота появления которых в наборе транзакций соответственно равна 7, 6 и 7. Следовательно, эти поднаборы частые и трехпредметный набор $\{ \text{Альбом, Карандаш, Краски} \}$ не подлежит изъятию из дальнейшего рассмотрения.

Таблица 7. Трехпредметные наборы-кандидаты в F_3

| Предметный набор | | | Частота появления в наборе транзакций |
|------------------|-----------------|----------------|--|
| Предмет 1 | Предмет 2 | Предмет 3 | |
| Альбом | Карандаш | Краски | 5 |
| Альбом | Карандаш | Ластик | 7 |
| Альбом | Карандаш | Тетрадь | 5 |
| Альбом | Краски | Ластик | 5 |
| Альбом | Краски | Тетрадь | 3 |
| Альбом | Ластик | Тетрадь | 5 |
| Карандаш | Краски | Ластик | 5 |
| Карандаш | Краски | Линейка | 0 |
| Карандаш | Краски | Ручка | 1 |
| Карандаш | Краски | Тетрадь | 3 |
| Карандаш | Ластик | Линейка | 5 |
| Карандаш | Ластик | Ручка | 7 |
| Карандаш | Ластик | Тетрадь | 12 |
| Карандаш | Линейка | Ручка | 2 |
| Карандаш | Линейка | Тетрадь | 4 |
| Карандаш | Ручка | Тетрадь | 6 |
| Ластик | Линейка | Ручка | 2 |
| Ластик | Линейка | Тетрадь | 4 |
| Ластик | Ручка | Тетрадь | 6 |

Рассмотрим набор {Альбом, Краски, Тетрадь}.

В нем есть двухпредметные поднаборы: {Альбом, Краски}, {Альбом, Тетрадь}, {Краски, Тетрадь}, частота появления которых в наборе транзакций соответственно равна 6, 5 и 3.

Первые два поднабора – частые, а третий – нечастый. Следовательно, и трехпредметный набор {Альбом, Краски, Тетрадь} не является частым и подлежит изъятию из дальнейшего рассмотрения.

Сформируем множество F_3 частых трехпредметных наборов:

$F_3 \{ \{ \text{Альбом, Карандаш, Краски} \}, \{ \text{Альбом, Карандаш, Ластик} \},$
 $\{ \text{Альбом, Карандаш, Тетрадь} \}, \{ \text{Альбом, Краски, Ластик} \},$
 $\{ \text{Альбом, Ластик, Тетрадь} \}, \{ \text{Карандаш, Краски, Ластик} \},$
 $\{ \text{Карандаш, Ластик, Линейка} \}, \{ \text{Карандаш, Ластик, Ручка} \},$
 $\{ \text{Карандаш, Ластик, Тетрадь} \}, \{ \text{Карандаш, Ручка, Тетрадь} \},$
 $\{ \text{Ластик, Ручка, Тетрадь} \} \}.$

Шаг 4. Выполним поиск частых четырехпредметных наборов (таблица 8).

Для получения множества F_4 свяжем множество F_3 с самим собой и сформируем множество F_4 частых четырехпредметных наборов:

$F_4 \{ \{ \text{Альбом, Карандаш, Краски, Ластик} \},$
 $\{ \text{Альбом, Карандаш, Ластик, Тетрадь} \},$
 $\{ \text{Карандаш, Ластик, Ручка, Тетрадь} \} \}.$

Эти четырехпредметные наборы уже нельзя связать. Поэтому задача поиска частых предметных наборов на исходном множестве транзакций решена.

2.2. Генерация ассоциативных правил

Для генерации ассоциативных правил на основе частых предметных наборов set_j ($j = \overline{1, m}$, m – число выявленных частных предметных наборов) к каждому частому набору set_j необходимо применить следующую процедуру.

1. Сгенерировать все возможные поднаборы $subset_{j,l_j}$ ($j = \overline{1, m}$; $l_j = \overline{1, m_j}$, m_j – число поднаборов на основе набора set_j).
2. Рассмотреть ассоциацию $R: subset_{j,l_j} \rightarrow (set_j - subset_{j,l_j})$, где $subset_{j,l_j}$ – непустой поднабор набора set_j ; $(set_j - subset_{j,l_j})$ – набор set_j без поднабора $subset_{j,l_j}$. Ассоциация R считается ассоциативным правилом, если удовлетворяет условию заданного минимума для поддержки и достоверности.

Таблица 8. Четырехпредметные наборы-кандидаты в F_4

| Предметный набор | | | | Частота появления в наборе транзакций |
|------------------|-----------------|---------------|----------------|--|
| Предмет 1 | Предмет 2 | Предмет 3 | Предмет 4 | |
| Альбом | Карандаш | Краски | Ластик | 5 |
| Альбом | Карандаш | Краски | Тетрадь | 3 |
| Альбом | Карандаш | Ластик | Тетрадь | 5 |
| Карандаш | Ластик | Линейка | Ручка | 2 |
| Карандаш | Ластик | Линейка | Тетрадь | 4 |
| Карандаш | Ластик | Ручка | Тетрадь | 6 |

Рассмотрим предметный набор-кандидат в ассоциативные правила.

$set = \{\text{Альбом, Карандаш, Краски}\}.$

В него входят поднаборы:

$\{\text{Альбом}\}, \{\text{Карандаш}\}, \{\text{Краски}\}, \{\text{Альбом, Карандаш}\}, \{\text{Альбом, Краски}\}, \{\text{Карандаш, Краски}\}.$

Построим ассоциативные правила, в которых условия содержат два предмета или один предмет.

Для ассоциативного правила с двумя предметами в условии предположим, что

$subset = \{\text{Альбом, Карандаш}\}.$

Тогда $(set - subset) = \{\text{Краски}\}.$

Рассмотрим правило:

$R: \{\text{Альбом, Карандаш}\} \rightarrow \{\text{Краски}\}.$

Поддержка, показывающая долю транзакций, которые содержат как условие $\{\text{Альбом, Карандаш}\}$, так и следствие $\{\text{Краски}\}$, в общем наборе из 20 транзакций составляет 25 % (5 из 20 транзакций).

Чтобы найти достоверность, нужно учесть, что набор $\{\text{Альбом, Карандаш}\}$ появляется в 7 из 20 транзакций, 5 из которых также содержат $\{\text{Краски}\}$. Тогда достоверность будет равна $5/7$ (71,4 %).

В таблице 9 приведены ассоциации с двумя предметами в условии для трехпредметного набора $\{\text{Альбом, Карандаш, Краски}\}.$

Если предположить, что минимальная достоверность для ассоциативного правила составляет 60%, то все ассоциации в таблице 9 будут правилами.

Если предположить, что минимальная достоверность для ассоциативного правила составляет 80%, то только вторая и третья ассоциации будут правилами.

Таблица 9. Ассоциации с двумя предметами в условии для трехпредметного набора $\{\text{Альбом, Карандаш, Краски}\}$

| Если условие, то следствие | Поддержка | Достоверность |
|--|----------------|-----------------|
| Если $\{\text{Альбом и Карандаш}\}$, то $\{\text{Краски}\}$ | $5/20 = 25 \%$ | $5/7 = 71,4 \%$ |
| Если $\{\text{Альбом и Краски}\}$, то $\{\text{Карандаш}\}$ | $5/20 = 25 \%$ | $5/6 = 80 \%$ |
| Если $\{\text{Карандаш и Краски}\}$, то $\{\text{Альбом}\}$ | $5/20 = 25 \%$ | $5/5 = 100 \%$ |

Таблица 10. Ассоциации с одним предметом в условии

| Если условие, то следствие | Поддержка | Достоверность |
|----------------------------|-----------|---------------|
| {Альбом} → {Карандаш} | 0,35 | 0,778 |
| {Карандаш} → {Альбом} | 0,35 | 0,438 |
| {Альбом} → {Краски} | 0,3 | 0,667 |
| {Краски} → {Альбом} | 0,3 | 1 |
| {Альбом} → {Ластик} | 0,35 | 0,778 |
| {Ластик} → {Альбом} | 0,35 | 0,438 |
| {Альбом} → {Тетрадь} | 0,25 | 0,556 |
| {Тетрадь} → {Альбом} | 0,25 | 0,417 |
| {Карандаш} → {Краски} | 0,25 | 0,313 |
| {Краски} → {Карандаш} | 0,25 | 0,833 |
| {Карандаш} → {Ластик} | 0,8 | 1 |
| {Ластик} → {Карандаш} | 0,8 | 1 |
| {Карандаш} → {Линейка} | 0,25 | 0,313 |
| {Линейка} → {Карандаш} | 0,25 | 0,714 |
| {Карандаш} → {Ручка} | 0,35 | 0,438 |
| {Ручка} → {Карандаш} | 0,35 | 0,778 |
| {Карандаш} → {Тетрадь} | 0,6 | 0,75 |
| {Тетрадь} → {Карандаш} | 0,6 | 1 |
| {Краски} → {Ластик} | 0,25 | 0,833 |
| {Ластик} → {Краски} | 0,25 | 0,313 |
| {Ластик} → {Линейка} | 0,25 | 0,313 |
| {Линейка} → {Ластик} | 0,25 | 0,714 |
| {Ластик} → {Ручка} | 0,35 | 0,438 |
| {Ручка} → {Ластик} | 0,35 | 0,778 |
| {Ластик} → {Тетрадь} | 0,6 | 0,75 |
| {Тетрадь} → {Ластик} | 0,6 | 1 |
| {Ручка} → {Тетрадь} | 0,3 | 0,667 |
| {Тетрадь} → {Ручка} | 0,3 | 0,5 |

Таблица 11. Ассоциации с двумя предметами в условии

| Если условие, то следствие | Поддержка | Достоверность |
|--------------------------------|-----------|---------------|
| {Альбом, Карандаш} → {Краски} | 0,25 | 0,714 |
| {Альбом, Краски} → {Карандаш} | 0,25 | 0,833 |
| {Карандаш, Краски} → {Альбом} | 0,25 | 1 |
| {Альбом, Карандаш} → {Ластик} | 0,35 | 1 |
| {Альбом, Ластик} → {Карандаш} | 0,35 | 1 |
| {Карандаш, Ластик} → {Альбом} | 0,35 | 0,438 |
| {Альбом, Карандаш} → {Тетрадь} | 0,25 | 0,714 |
| {Альбом, Тетрадь} → {Карандаш} | 0,25 | 1 |
| {Карандаш, Тетрадь} → {Альбом} | 0,25 | 0,417 |
| {Альбом, Краски} → {Ластик} | 0,25 | 0,833 |
| {Альбом, Ластик} → {Краски} | 0,25 | 0,714 |
| {Краски, Ластик} → {Альбом} | 0,25 | 1 |
| {Альбом, Ластик} → {Тетрадь} | 0,25 | 0,714 |
| {Альбом, Тетрадь} → {Ластик} | 0,25 | 1 |
| {Ластик, Тетрадь} → {Альбом} | 0,25 | 0,417 |
| {Карандаш, Краски} → {Ластик} | 0,25 | 1 |
| {Карандаш, Ластик} → {Краски} | 0,25 | 0,313 |
| {Краски, Ластик} → {Карандаш} | 0,25 | 1 |
| {Карандаш, Ластик} → {Линейка} | 0,25 | 0,313 |
| {Карандаш, Линейка} → {Ластик} | 0,25 | 1 |
| {Ластик, Линейка} → {Карандаш} | 0,25 | 1 |
| {Карандаш, Ластик} → {Ручка} | 0,35 | 0,438 |
| {Карандаш, Ручка} → {Ластик} | 0,35 | 1 |
| {Ластик, Ручка} → {Карандаш} | 0,35 | 1 |
| {Карандаш, Ластик} → {Тетрадь} | 0,6 | 0,75 |
| {Карандаш, Тетрадь} → {Ластик} | 0,6 | 1 |
| {Ластик, Тетрадь} → {Карандаш} | 0,6 | 1 |
| {Карандаш, Ручка} → {Тетрадь} | 0,3 | 0,857 |
| {Карандаш, Тетрадь} → {Ручка} | 0,3 | 0,5 |
| {Ручка, Тетрадь} → {Карандаш} | 0,3 | 1 |
| {Ластик, Ручка} → {Тетрадь} | 0,3 | 0,857 |
| {Ластик, Тетрадь} → {Ручка} | 0,3 | 0,5 |
| {Ручка, Тетрадь} → {Ластик} | 0,3 | 1 |

Таблица 12. Ассоциации с тремя предметами в условии

| Если условие, то следствие | Поддержка | Достоверность |
|--|-----------|---------------|
| {Карандаш, Ластик, Ручка} → {Тетрадь} | 0,3 | 0,857 |
| {Ластик, Ручка, Тетрадь} → {Карандаш} | 0,3 | 1 |
| {Ручка, Тетрадь, Карандаш} → {Ластик} | 0,3 | 1 |
| {Тетрадь, Карандаш, Ластик} → {Ручка} | 0,3 | 0,5 |
| {Альбом, Карандаш, Краски} → {Ластик} | 0,25 | 1 |
| {Карандаш, Краски, Ластик} → {Альбом} | 0,25 | 1 |
| {Краски, Ластик, Альбом} → {Карандаш} | 0,25 | 1 |
| {Ластик, Альбом, Карандаш} → {Краски} | 0,25 | 0,714 |
| {Альбом, Карандаш, Ластик} → {Тетрадь} | 0,25 | 0,714 |
| {Карандаш, Ластик, Тетрадь} → {Альбом} | 0,25 | 0,417 |
| {Ластик, Тетрадь, Альбом} → {Карандаш} | 0,25 | 1 |
| {Тетрадь, Альбом, Карандаш} → {Ластик} | 0,25 | 1 |

Получим все ассоциативные правила на основе выявленных частых предметных наборов *с одним предметом в следствии*.

Ассоциативные правила с двумя и более предметами в следствии могут быть построены аналогичным образом.

В таблице 10 приведены все ассоциации с одним предметом в условии и одним предметом (всего – 28 ассоциаций).

В таблице 11 приведены все ассоциации с двумя предметами в условии (всего – 33 ассоциации).

В таблице 12 приведены все ассоциации с тремя предметами в условии (всего – 12 ассоциаций).

Чтобы оценить значимость сгенерированных ассоциаций необходимо перемножить соответствующие им значения поддержки и достоверности.

Выберем в качестве ассоциативных правил те ассоциации из таблиц 10 – 12, у которых достоверность не ниже 80% (таблица 13).

В результате применения алгоритма Apriori удалось найти 34 ассоциативных правил, которые с достоверностью не ниже 80% показывают, какие продукты из исходного набора предметов чаще всего продаются вместе.

Если понизить пороговое значение достоверности до 60%, то число ассоциативных правил увеличиться на 17 штук и будет равно 51.

Таблица 13. Ассоциативные правила при минимальной достоверности 80%

| Если условие, то следствие | Поддержка S | Достоверность C | $S \cdot C$ |
|--|---------------|-------------------|-------------|
| {Карандаш, Ластик, Ручка} → {Тетрадь} | 0,3 | 0,857 | 0,257 |
| {Ластик, Ручка, Тетрадь} → {Карандаш} | 0,3 | 1 | 0,3 |
| {Ручка, Тетрадь, Карандаш} → {Ластик} | 0,3 | 1 | 0,3 |
| {Альбом, Карандаш, Краски} → {Ластик} | 0,25 | 1 | 0,25 |
| {Карандаш, Краски, Ластик} → {Альбом} | 0,25 | 1 | 0,25 |
| {Краски, Ластик, Альбом} → {Карандаш} | 0,25 | 1 | 0,25 |
| {Ластик, Тетрадь, Альбом} → {Карандаш} | 0,25 | 1 | 0,25 |
| {Тетрадь, Альбом, Карандаш} → {Ластик} | 0,25 | 1 | 0,25 |
| {Альбом, Краски} → {Карандаш} | 0,25 | 0,833 | 0,208 |
| {Карандаш, Краски} → {Альбом} | 0,25 | 1 | 0,25 |
| {Альбом, Карандаш} → {Ластик} | 0,35 | 1 | 0,35 |
| {Альбом, Ластик} → {Карандаш} | 0,35 | 1 | 0,35 |
| {Альбом, Краски} → {Ластик} | 0,25 | 0,833 | 0,208 |
| {Краски, Ластик} → {Альбом} | 0,25 | 1 | 0,25 |
| {Альбом, Тетрадь} → {Ластик} | 0,25 | 1 | 0,25 |
| {Карандаш, Краски} → {Ластик} | 0,25 | 1 | 0,25 |
| {Краски, Ластик} → {Карандаш} | 0,25 | 1 | 0,25 |
| {Карандаш, Линейка} → {Ластик} | 0,25 | 1 | 0,25 |
| {Ластик, Линейка} → {Карандаш} | 0,25 | 1 | 0,25 |
| {Карандаш, Ручка} → {Ластик} | 0,35 | 1 | 0,35 |
| {Ластик, Ручка} → {Карандаш} | 0,35 | 1 | 0,35 |
| {Карандаш, Тетрадь} → {Ластик} | 0,6 | 1 | 0,6 |
| {Ластик, Тетрадь} → {Карандаш} | 0,6 | 1 | 0,6 |
| {Карандаш, Ручка} → {Тетрадь} | 0,3 | 0,857 | 0,257 |
| {Ручка, Тетрадь} → {Карандаш} | 0,3 | 1 | 0,3 |
| {Ластик, Ручка} → {Тетрадь} | 0,3 | 0,857 | 0,257 |
| {Ручка, Тетрадь} → {Ластик} | 0,3 | 1 | 0,3 |
| {Краски} → {Альбом} | 0,3 | 1 | 0,3 |
| {Краски} → {Карандаш} | 0,25 | 0,833 | 0,208 |
| {Карандаш} → {Ластик} | 0,8 | 1 | 0,8 |
| {Ластик} → {Карандаш} | 0,8 | 1 | 0,8 |
| {Тетрадь} → {Карандаш} | 0,6 | 1 | 0,6 |
| {Краски} → {Ластик} | 0,25 | 0,833 | 0,208 |
| {Тетрадь} → {Ластик} | 0,6 | 1 | 0,6 |

Очевидно, что по результатам выявления ассоциативных правил можно разработать адекватную маркетинговую стратегию, оптимизировать закупки и размещение товаров (продуктов) на прилавках и витринах.

3. Алгоритм FPGrowth

Как уже было ранее отмечено, алгоритм Apriori реализует генерацию кандидатов в частые предметные наборы, а это приводит к тому, что вычислительные и временные затраты, которые необходимы для обработки (анализа) кандидатов, могут быть неприемлемыми.

Кроме того, алгоритм Apriori требует многократного сканирования базы данных транзакций: необходимо выполнить столько сканирований, сколько предметов содержит самый длинный предметный набор.

В связи с этим были предложены алгоритмы, позволяющие отказаться от генерации кандидатов и сократить требуемое число сканирований.

Одной из наиболее эффективных альтернатив алгоритму Apriori является алгоритм Frequent Pattern-Growth (FPGrowth), название которого можно интерпретировать как «выращивание часто встречающихся предметных наборов». Алгоритм FPGrowth позволяет избежать затратной процедуры генерации кандидатов, а также – уменьшить необходимое число проходов по базе данных транзакций до двух [4 – 7].

В основе алгоритма FPGrowth лежит предобработка БД транзакций, в процессе которой на основе БД транзакций формируется компактная древовидная структура, а именно – дерево частых предметных наборов (Frequent-Pattern Tree, FP-дерево).

Сжатие БД транзакций в компактную структуру обеспечивает эффективное и полное извлечение частых предметных наборов. При этом при построении FP-дерева применяется парадигма «разделяй и властвуй» (divide and conquer), позволяющий выполнить декомпозицию сложной задачи на некоторое число более простых.

Сравнивая алгоритмы Apriori и FPGrowth, можно отметить следующее.

Алгоритм Apriori использует подход, реализующий генерацию кандидатов в частые предметные наборы и тестирование кандидатов на предмет, являются ли они частыми.

Генерация кандидатов требует больших временных затрат, а также больших затрат памяти для их хранения.

Расчет значений поддержки также требует больших временных затрат, так как сопровождается анализом поднаборов и множественными проходами по БД транзакций.

Алгоритм FPGrowth позволяет выявлять частые предметные наборы без генерации кандидатов в них. При этом алгоритм FPGrowth требует выполнения 2 нижеуказанных этапов [4 – 7].

Этап 1. Построение компактной структуры данных в виде FP-дерева с реализацией 2 проходов по БД транзакций.

Этап 2. Извлечение частых предметных наборов непосредственно из FP-дерева.

При реализации алгоритма FPGrowth предполагается следующее.

1. Узлы FP-дерева соответствуют предметам в БД транзакций и имеют счётчики поддержки.

2. При обращении к БД транзакций считывается одна транзакция, в соответствие которой ставится путь в FP-дереве.

3. Осуществляется упорядочение предметов в каждой транзакции по убыванию частоты появления этих предметов в БД транзакций. Такое упорядочение предметов позволяет сделать так, что пути в дереве перекрываются (накладываются), если имеют общие предметы (т.е. имеют одинаковые префиксы). При наличии наложения путей для транзакций счётчики поддержки соответствующих узлов увеличиваются.

4. Создаются односвязные списки, указатели которых используются перемещения между узлами, соответствующими одному и тому же предмету (на изображении FP-дерева связи между такими узлами помечаются пунктирными линиями со стрелками).

5. Чем больше наложение путей, тем больше сжатие, что позволяет в итоге без больших проблем разместить FP-дерево в памяти.

6. Частые предметные наборы извлекаются из FP-дерева.

Рассмотрим работу алгоритма FPGrowth на примере той же БД транзакций (таблица 1), которая была использована для алгоритма Apriori.

Будем считать (как и в алгоритме Apriori) частыми такие предметные наборы, которые встречаются в базе транзакций более 5 раз, то есть определим пороговое значение частоты как $\Delta = 5$.

В дальнейшем в FPGrowth алгоритме значение частоты будем называть поддержкой.

Таблица 14. Пример набора транзакций на основе чеков магазина «Канцтовары» с упорядочением предметов в транзакциях по убыванию значений их частоты (поддержки)

| Номер чека | Список предметов в чеке | | | | | | | |
|------------|-------------------------|--------|---------|---------|---------|---------|-------|-------|
| 1 | Карандаш | Ластик | Тетрадь | Альбом | Линейка | | | |
| 2 | Карандаш | Ластик | Тетрадь | Ручка | Линейка | Пенал | | |
| 3 | Карандаш | Ластик | Ручка | | | | | |
| 4 | Карандаш | Ластик | Альбом | Краски | Кисть | Палитра | | |
| 5 | Карандаш | Ластик | Тетрадь | Ручка | | | | |
| 6 | Альбом | Краски | Кисть | Палитра | | | | |
| 7 | Карандаш | Ластик | Тетрадь | Альбом | Краски | | | |
| 8 | Карандаш | Ластик | Линейка | | | | | |
| 9 | Карандаш | Ластик | Тетрадь | Ручка | Пенал | | | |
| 10 | Карандаш | Ластик | Тетрадь | Альбом | Ручка | Краски | Пенал | Кисть |
| 11 | Карандаш | Ластик | Тетрадь | Линейка | | | | |
| 12 | Карандаш | Ластик | Тетрадь | Пенал | | | | |
| 13 | Ручка | Пенал | | | | | | |
| 14 | Альбом | Ручка | Линейка | | | | | |
| 15 | Карандаш | Ластик | Тетрадь | Альбом | | | | |
| 16 | Карандаш | Ластик | Тетрадь | Ручка | Линейка | | | |
| 17 | Карандаш | Ластик | Тетрадь | Альбом | Краски | Палитра | | |
| 18 | Карандаш | Ластик | Альбом | Краски | | | | |
| 19 | Линейка | Пенал | | | | | | |
| 20 | Карандаш | Ластик | Тетрадь | Ручка | | | | |

Этап 1. Построение компактной структуры данных в виде FP-дерева.

Число сканирований, т.е. число проходов по БД транзакций, равно 2.

1. Первое сканирование базы данных транзакций.

Произведем сканирование БД транзакций и сформируем множество частых предметов, т.е. предметов, которые встречаются Δ (например, $\Delta = 5$) или более раз. Исключим из дальнейшего рассмотрения нечастые предметы.

Упорядочим выявленные частые предметы в порядке убывания частот их появления в БД транзакций (таблица 5) и получим следующий набор:

(Карандаш, 16), (Ластик, 16), (Тетрадь, 12), (Альбом, 9), (Ручка, 9), (Линейка, 7), (Краски, 6), (Пенал, 6).

Этот порядок будем использовать при построении FP-дерева. В результате общие для транзакций префиксы (пути) в FP-дереве будут накладываться.

Разместим эти наборы в указанном порядке в столбик, отделив название предмета от его частоты (счетчика поддержки) двоеточием (рисунок 1).

В дальнейшем будем размещать узлы дерева, соответствующие одному и тому же предмету, на одном уровне, делая ветви дерева при необходимости более длинными.

Например, все узлы, соответствующие предмету Пенал, должны быть размещены на самом нижнем уровне, т.к. предмет Пенал имеет одну из самых низких частот (наряду с предметом Краски) в списке упорядочения частых предметов.

2. Второе сканирование базы данных транзакций.

Построим FP-дерево, предварительно упорядочив предметы в транзакциях по убыванию значений их частот, то есть по убыванию значений поддержки (таблица 14), удалив из транзакций те предметы, у которых значение поддержки меньше, чем Δ (например, $\Delta = 5$). В таблице 14 нечастые предметы помечены как «вычеркнутые».

Создадим начальный (корневой, root) узел FP-дерева.

Осуществим построение дерева, начиная с транзакции №1 для упорядоченных предметных наборов, т.е. рассмотрим сначала предметный набор (Карандаш, Ластик, Тетрадь, Альбом, Линейка).

При построении дерева будем использовать следующее правило.

Правило. Если для очередного предмета в FP-дереве уже содержится узел, то для предмета не создается новый узел, а счетчик поддержки существующего узла увеличивается на 1. В противном случае для этого предмета создается новый узел и его счётчику поддержки присваивается значение, равное 1.

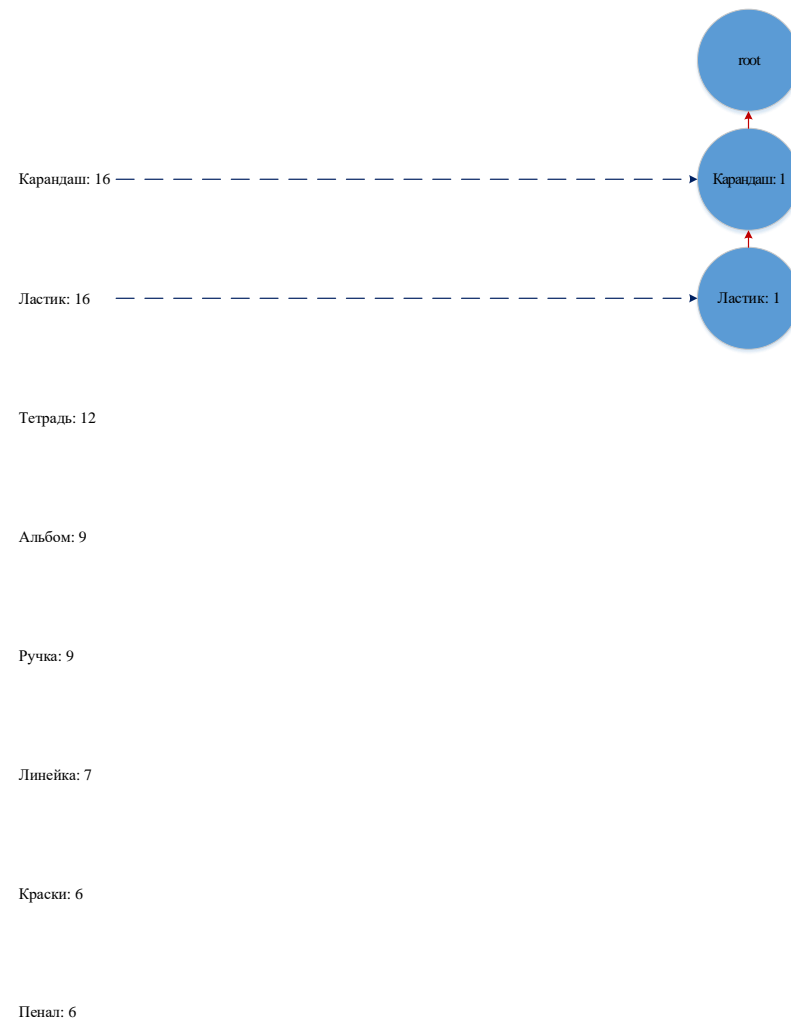
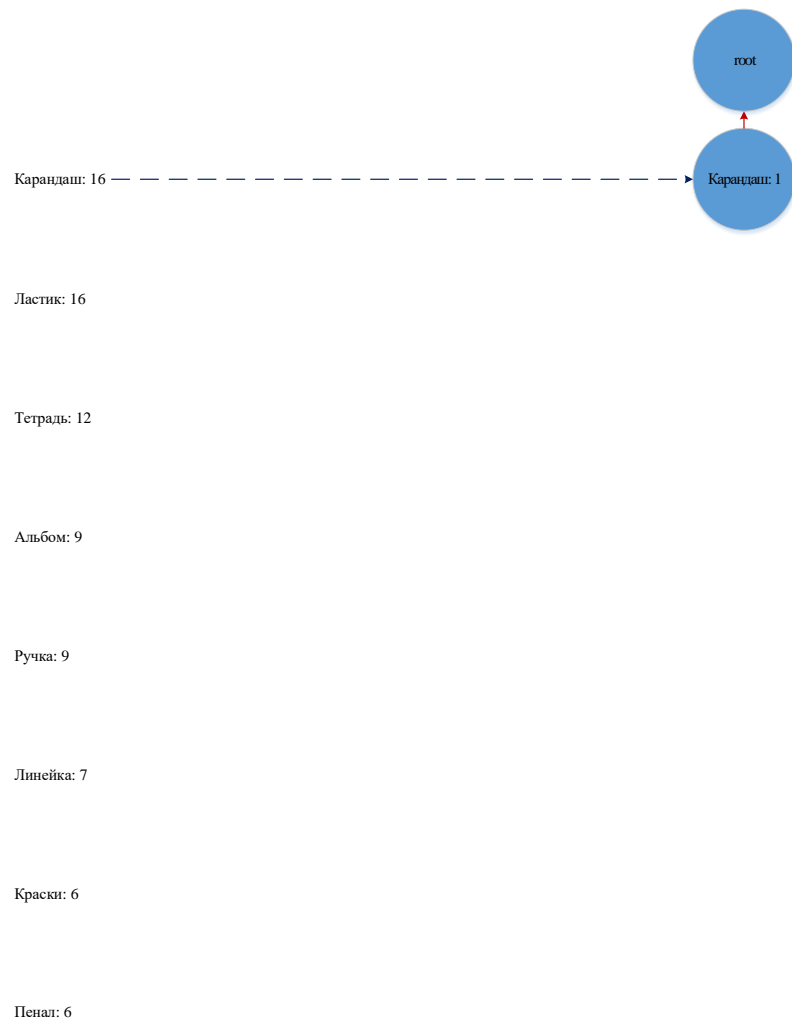
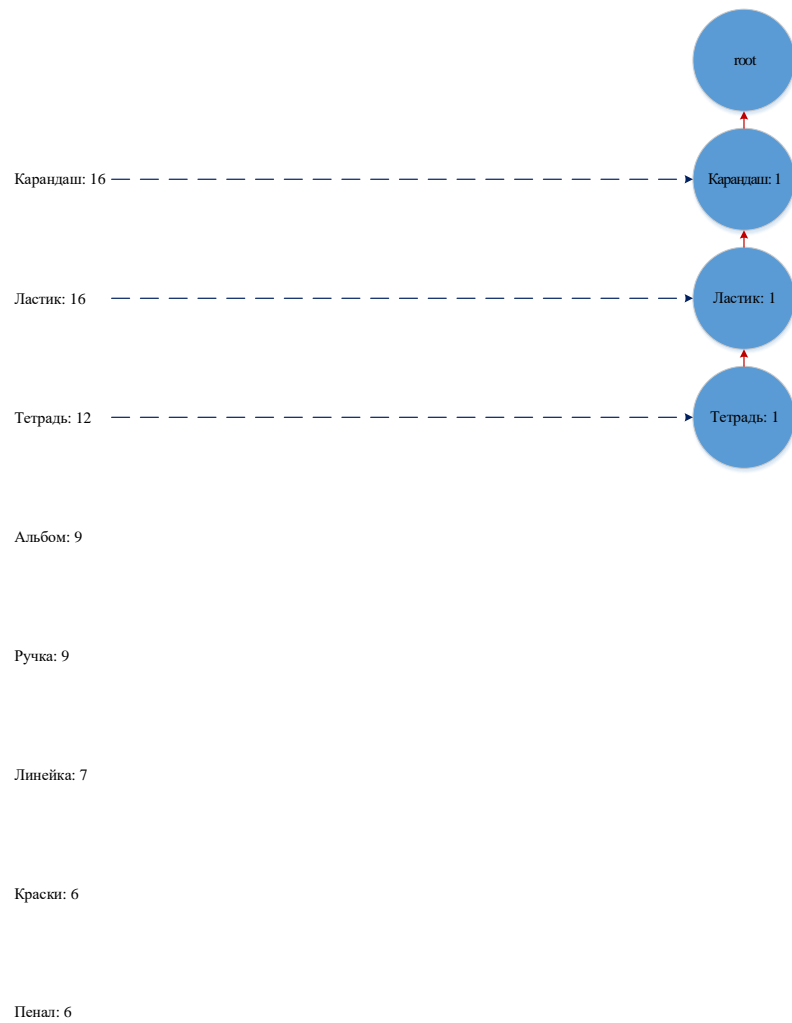
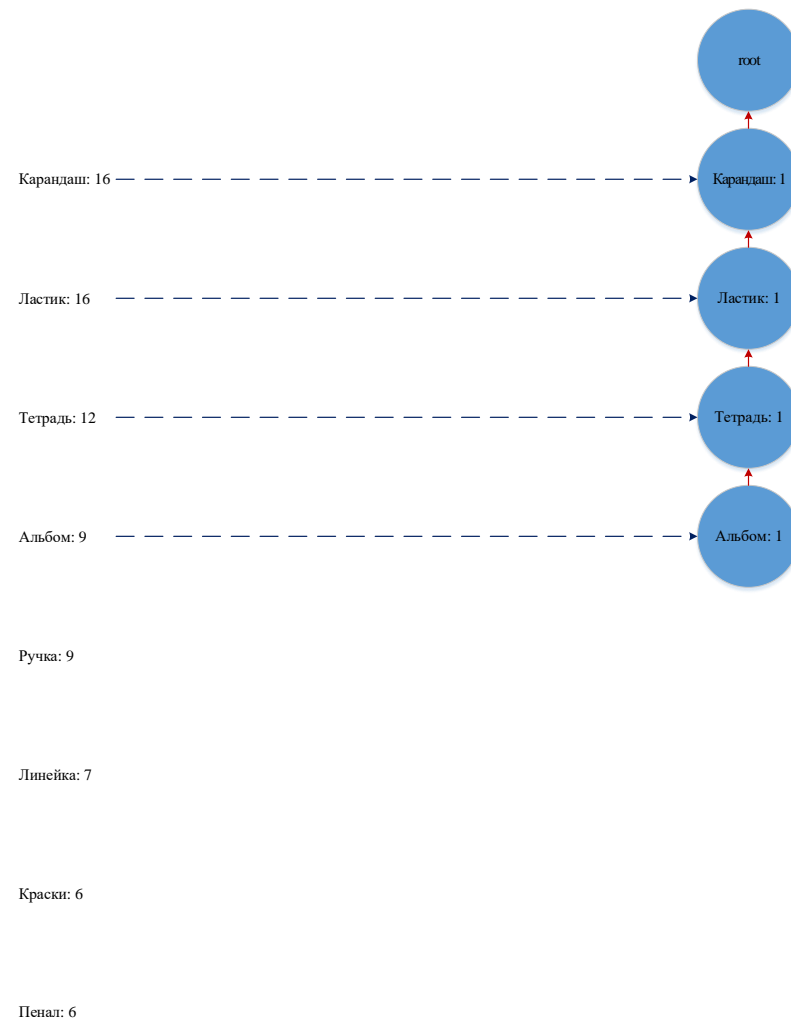


Рисунок 1. Шаги построения FP-дерева на основе транзакции № 1 (начало)



а



б

Рисунок 2. Шаги построения FP-дерева на основе транзакции № 1 (продолжение)

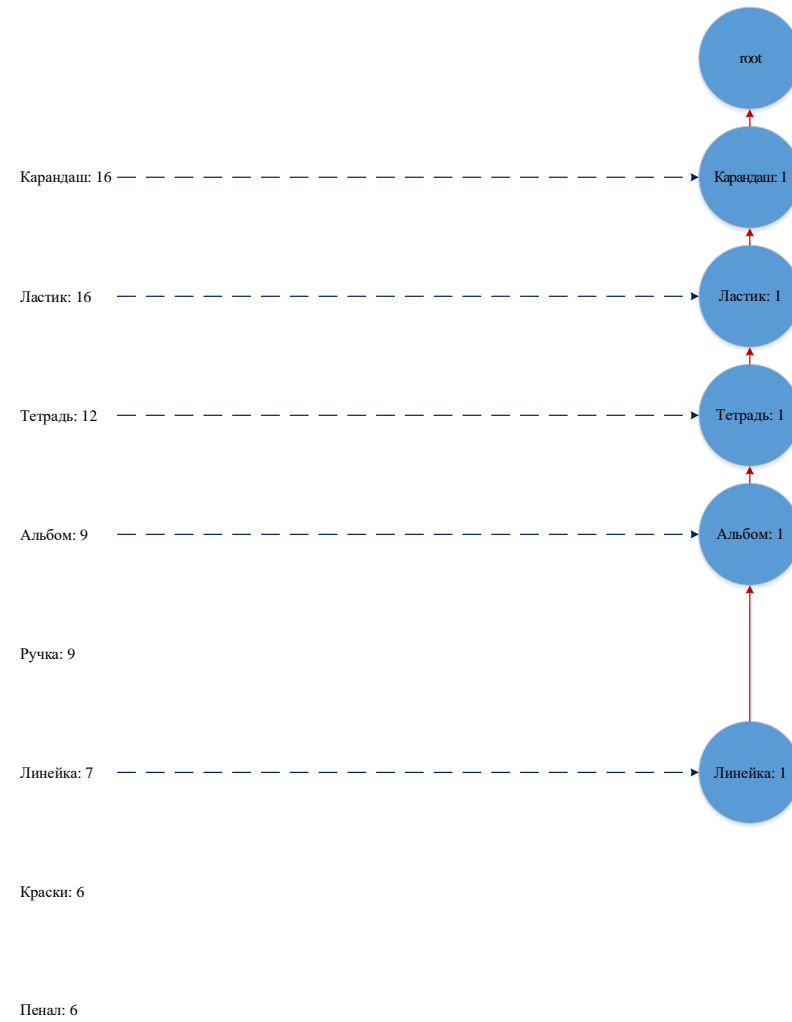


Рисунок 3. Шаги построения FP-дерева на основе транзакции № 1 (окончание)

Сплошными бордовыми линиями со стрелки будем отмечать связи дочерних узлов с родительскими.

Пунктирными темно-синими линиями со стрелками будем отмечать связи между узлами дерева, соответствующими одним и тем же предметам, с целью формирования односвязанных списков.

Сначала возьмем предмет Карандаш из транзакции №1. Поскольку он является первым, то создаем для него узел и соединяем с родительским (корневым) узлом (рисунок 1, а).

Затем возьмем следующий предмет – предмет Ластик. Поскольку FR-дерево пока не содержит узлов с именем Ластик, добавляем в FR-дерево новый узел Ластик, который будет потомком узла Карандаш (рисунок 1, б).

Затем возьмем следующий предмет – предмет Тетрадь.

Рассуждая аналогичным образом, добавляем в FR-дерево новый узел Тетрадь, который будет потомком узла Ластик (рисунок 2, а).

Таким же образом создадим узлы для предметов Альбом (рисунок 2, б) и Линейка (рисунок 3) из транзакции № 1. Для каждого нового узла определим значение его счетчика равным 1.

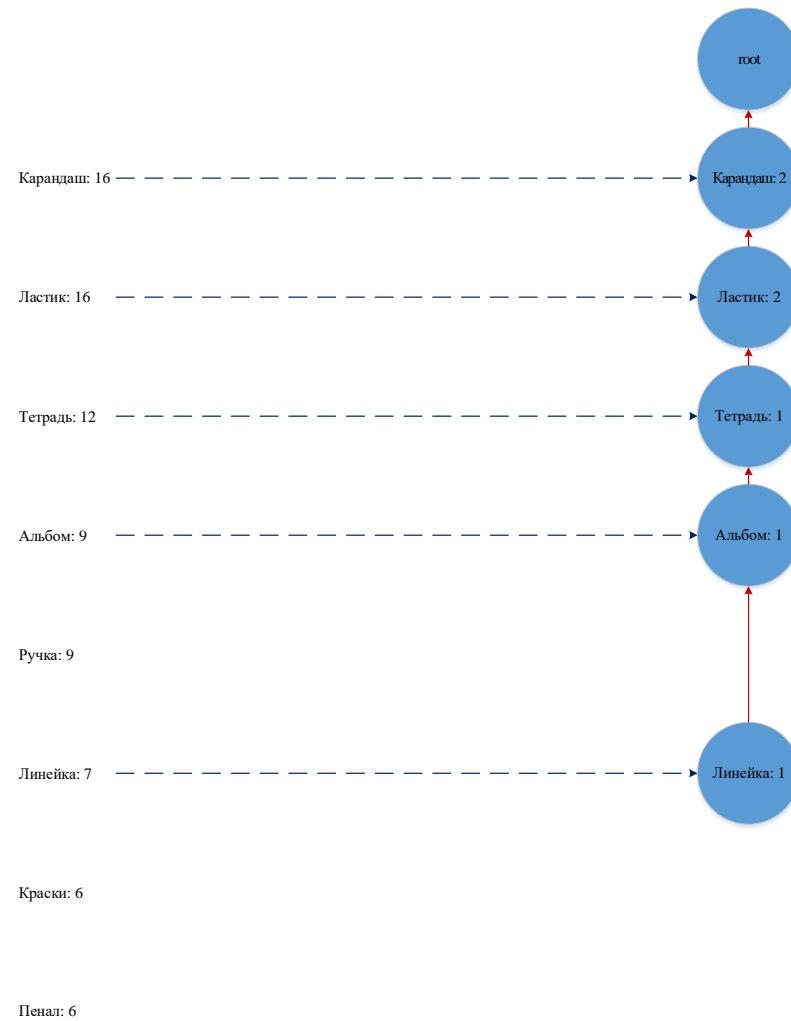
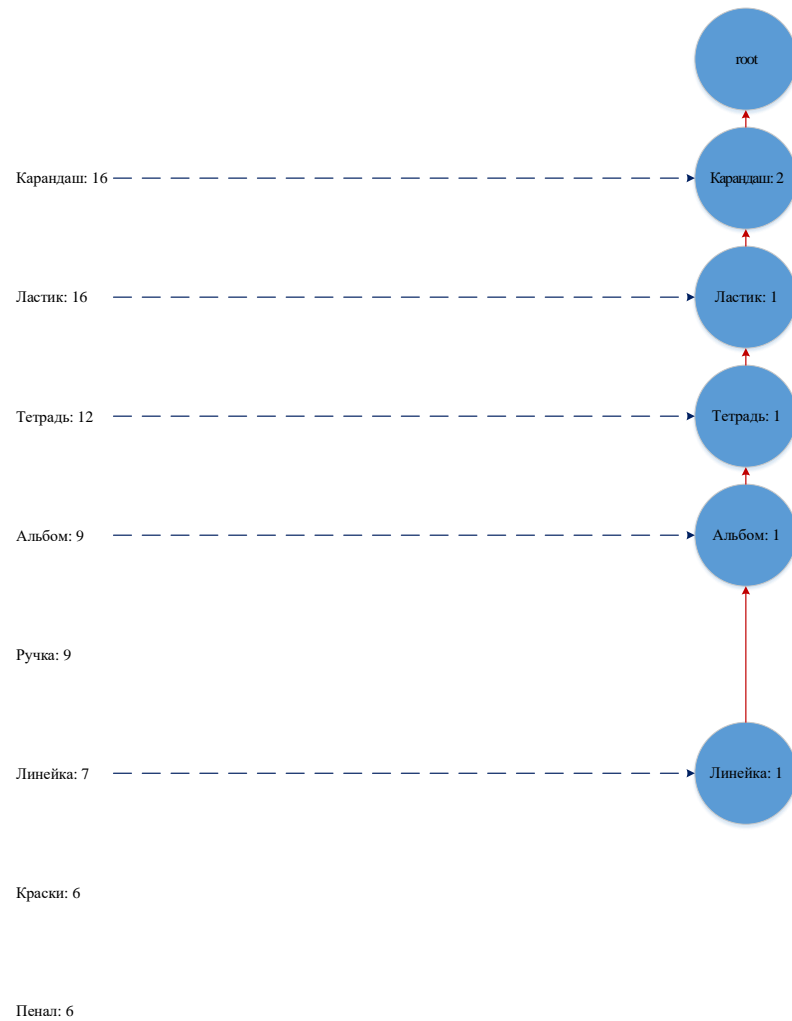
На этом использование транзакции №1 для построения FR-дерева завершено.

Для транзакции №2, содержащей предметы из предметного набора (Карандаш, Ластик, Тетрадь, Ручка, Линейка, Пенал) выбираем первый предмет – предмет Карандаш. Поскольку узел с таким именем уже существует, то в соответствии с правилом построения FR-дерева новый узел не создается, а счетчик поддержки существующего увеличивается на 1 (рисунок 4, а).

При добавлении следующего предмета – предмета Ластик – используем это же правило: поскольку узел Ластик является дочерним по отношению к текущему (т.е. к узлу Карандаш), то новый узел не создается, а счетчик поддержки существующего увеличивается на 1 (рисунок 4, б).

При добавлении следующего предмета – предмета Тетрадь – используем это же правило: поскольку узел Тетрадь является дочерним по отношению к текущему (т.е. к узлу Ластик), то новый узел не создается, а счетчик поддержки существующего увеличивается на 1 (рисунок 5).

Следующий предмет – предмет Ручка – в соответствии с правилом построения FR-дерева породит в FR-дереве новый узел, дочерний к узлу Тетрадь (рисунок 6).



а

б

Рисунок 4. Шаги построения FP-дерева на основе транзакции № 2 (начало):
добавление предметов Карандаш и Ластик без создания новых узлов

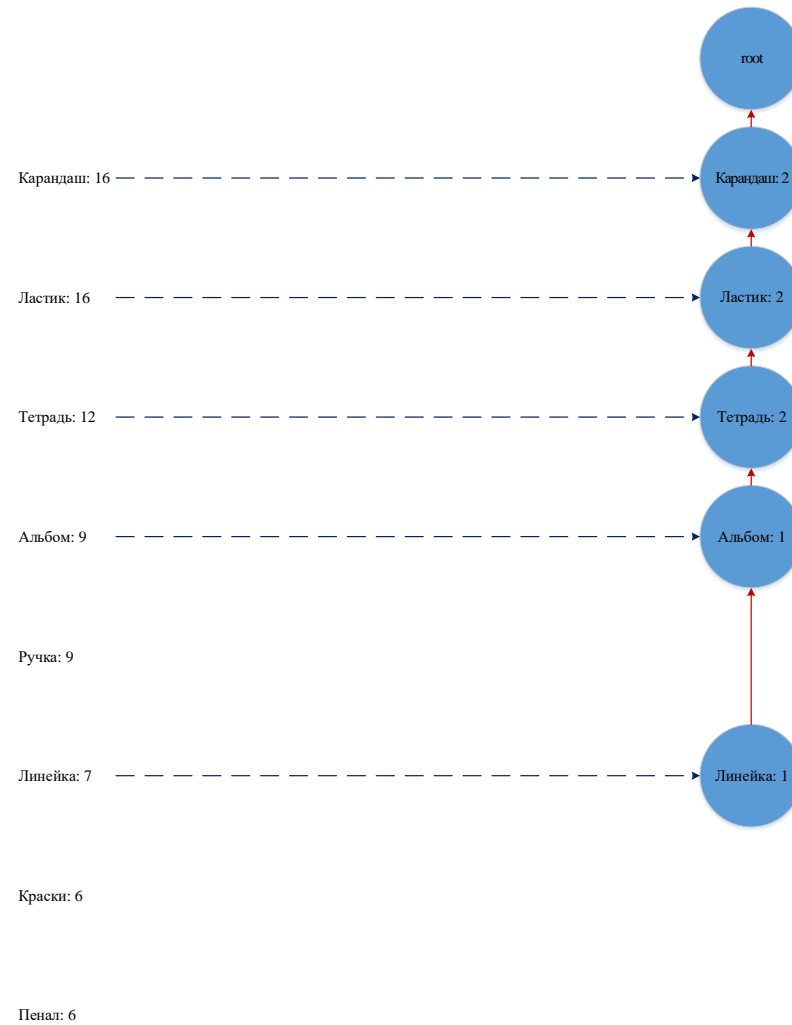


Рисунок 5. Шаги построения FP-дерева на основе транзакции № 2 (продолжение):
добавление предмета Тетрадь без создания нового узла

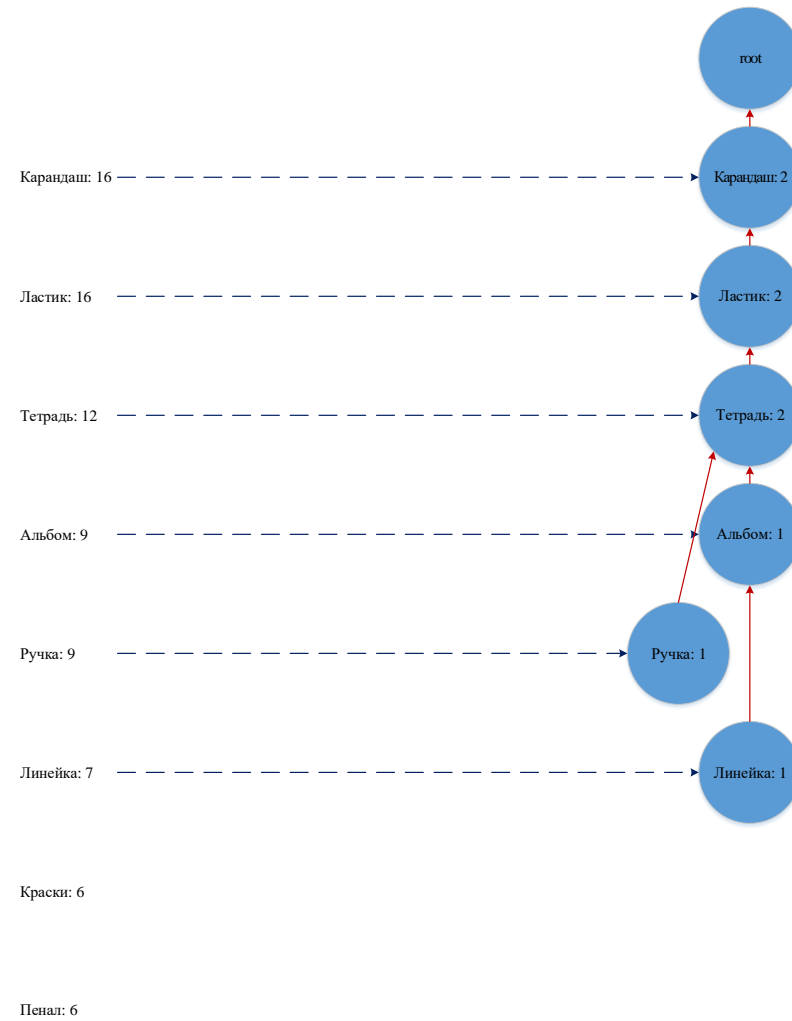


Рисунок 6. Шаги построения FP-дерева на основе транзакции № 2 (продолжение):
добавление предмета Ручка с созданием нового узла

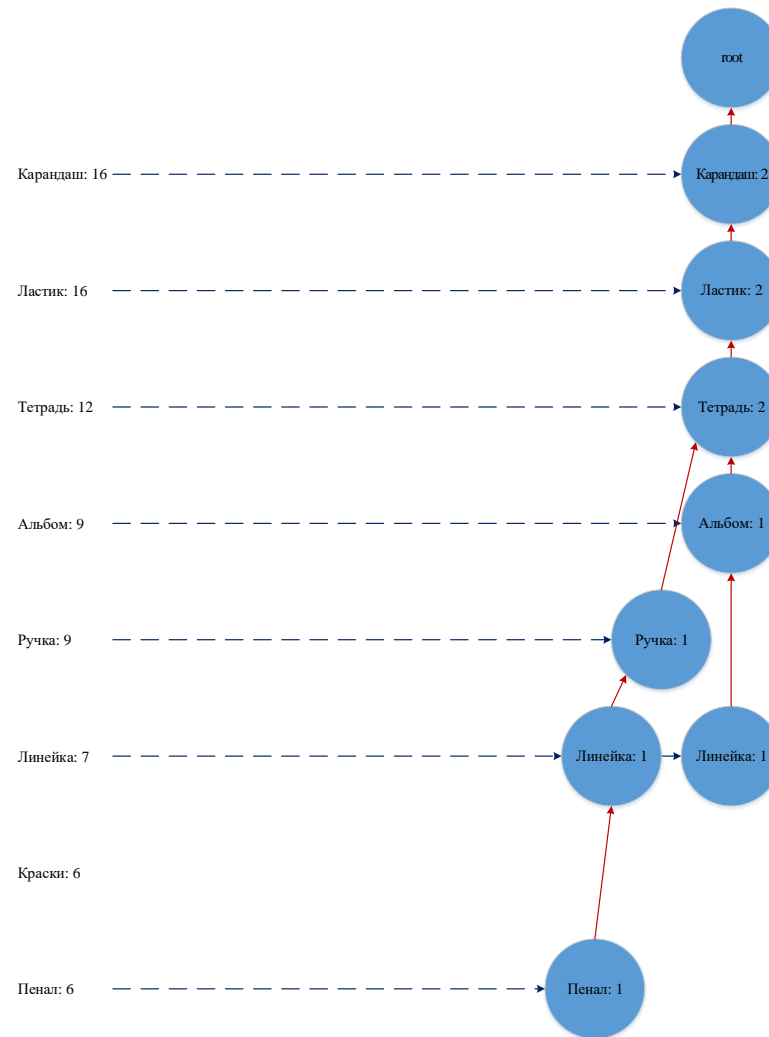


Рисунок 8. Шаги построения FP-дерева на основе транзакции № 2 (окончание):
добавление предмета Пенал с созданием нового узла

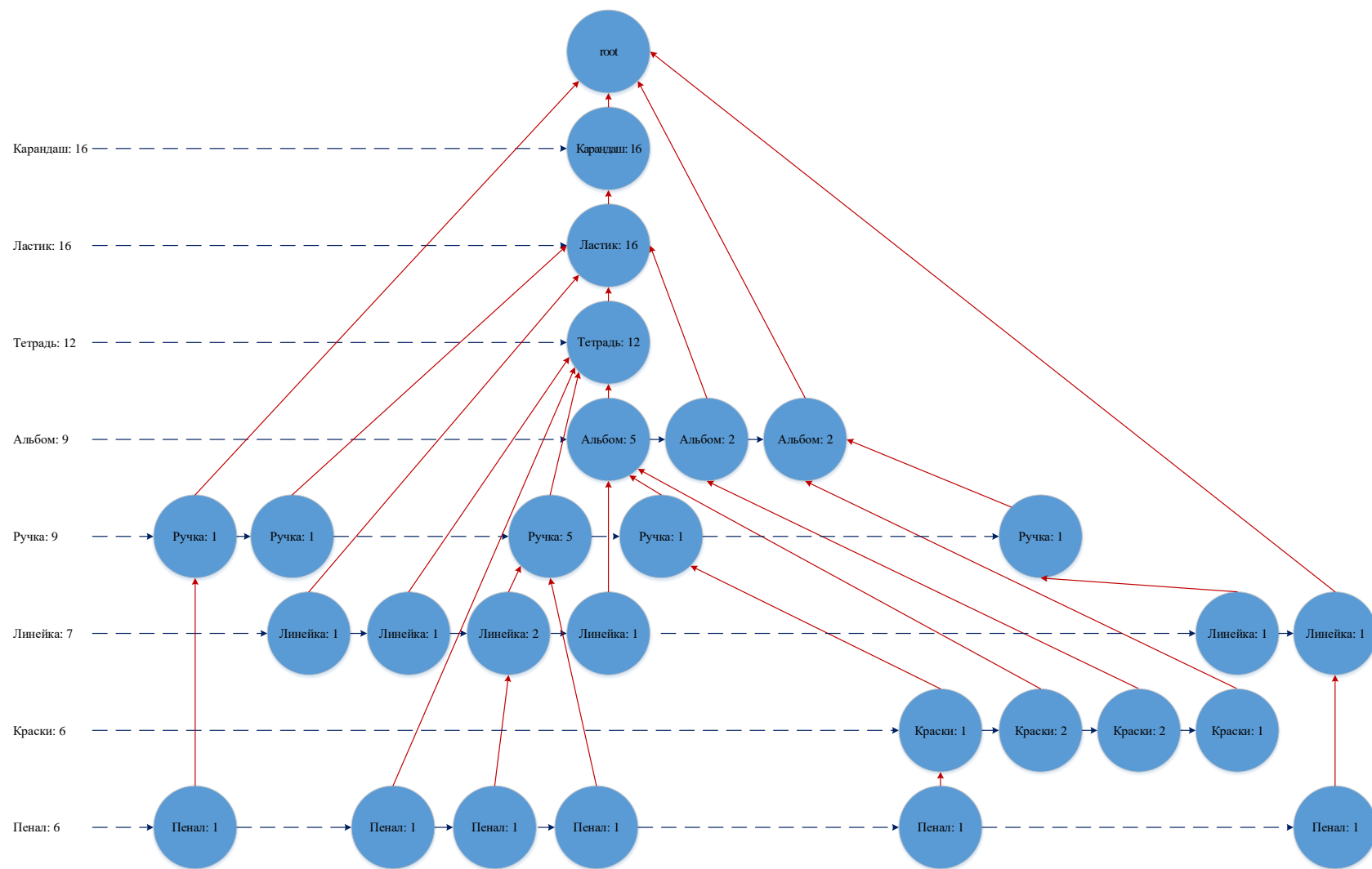


Рисунок 9. Итоговое FP-дерево

В соответствии с правилом построения FР-дерева предмет Линейка породит в FР-дереве новый узел, дочерний к узлу Ручка (рисунок 7), а предмет Пенал породит в FР-дереве новый узел, дочерний к узлу Линейка (рисунок 8).

На этом использование транзакции №2 для построения FР-дерева завершено.

Рассуждая аналогичным образом, можно обработать оставшиеся транзакции под номерами 3 – 20.

На рисунке 9 приведено итоговое FР-дерево, построенное на основе транзакций из таблицы 14.

Можно увидеть, что на любом уровне сумма частот появления одного и того же предмета в разных узлах дерева равна общей частоте этого предмета, указанной через двоеточие от названия предмета (в левой части на рисунке 9).

Таким образом, после двух сканирований БД транзакций и выполнения соответствующей обработки предметных наборов было получено итоговое FР-дерево, которое в компактном виде представляет информацию о частых предметных наборах.

Итоговое FР-дерево обычно имеет меньшие размеры, чем несжатые исходные данные, т.к. многие транзакции имеют общие предметы (и, следовательно, общие префиксы). При этом возможны следующие варианты, характеризующие размеры FР-дерева.

1. Лучший вариант: все транзакции имеют один и тот же набор предметов. Тогда в FР-дереве имеется один единственный путь.

2. Худший вариант: все транзакции содержат уникальные наборы предметов, т.е. у них нет общих предметов. Тогда размеры FР-дерева такие же, как размеры исходных данных, и требования к памяти – велики, т.к. необходимо хранить указатели для узлов и счетчики.

Размеры FР-дерева зависят от того, как упорядочены предметы. Обычно используется упорядочение по убыванию значений поддержки, однако это не всегда приводит к получению FР-дерева наименьших размеров.

Этап 2. Извлечение частых предметных наборов из FР-дерева.

При реализации этого этапа осуществляется движение снизу вверх – от листьев FР-дерева к корню.

Поиск частых предметных наборов выполняется рекурсивно.

При этом используется парадигма «разделяй и властвуй», заключающаяся в рекурсивном разбиении решаемой задачи на две или более подзадачи того же

типа, но меньшего размера, и комбинировании их решений для получения ответа к исходной задаче.

Разбиения выполняются до тех пор, пока все подзадачи не окажутся элементарными.

Для предмета в FR-дереве, представленного узлом, можно указать путь, определяемый последовательностью узлов, через которые необходимо пройти от корневого узла до узла, связанного с этим предметом.

Если предмет представлен в нескольких ветвях FR-дерева, то таких путей будет несколько.

Например, для FR-дерева на рисунке 9 для предмета Линейка можно указать 6 путей:

- {Карандаш, Ластик, Линейка},
- {Карандаш, Ластик, Тетрадь, Линейка},
- {Карандаш, Ластик, Тетрадь, Ручка, Линейка},
- {Карандаш, Ластик, Тетрадь, Альбом, Линейка},
- {Альбом, Ручка, Линейка},
- {Линейка}.

Такой набор путей называется **условным базисом предмета** (*conditional base*).

Каждый путь в базисе состоит из двух частей – префикса и суффикса.

Префикс – последовательность узлов, которые образуют путь.

Префикс – любое начало пути, в том числе и пустое.

Будем говорить, что префикс является собственным, если он не совпадает с пустым путем и полным путем.

Суффикс – любой конец пути, в том числе пустой.

Будем говорить, что суффикс является собственным, если он не совпадает с пустым путем и полным путем.

Пусть суффикс – узел (возможно, лист), к которому «прокладывается» путь.

В условном базисе все пути будут иметь различные префиксы и одинаковый суффикс.

Например, в пути {Карандаш, Ластик, Линейка} в качестве префикса, не совпадающего с полным путем, можно назвать поднабор {Карандаш, Ластик|}, а в качестве суффикса – поднабор из одного предмета, т.е. однопредметный набор {Линейка}.

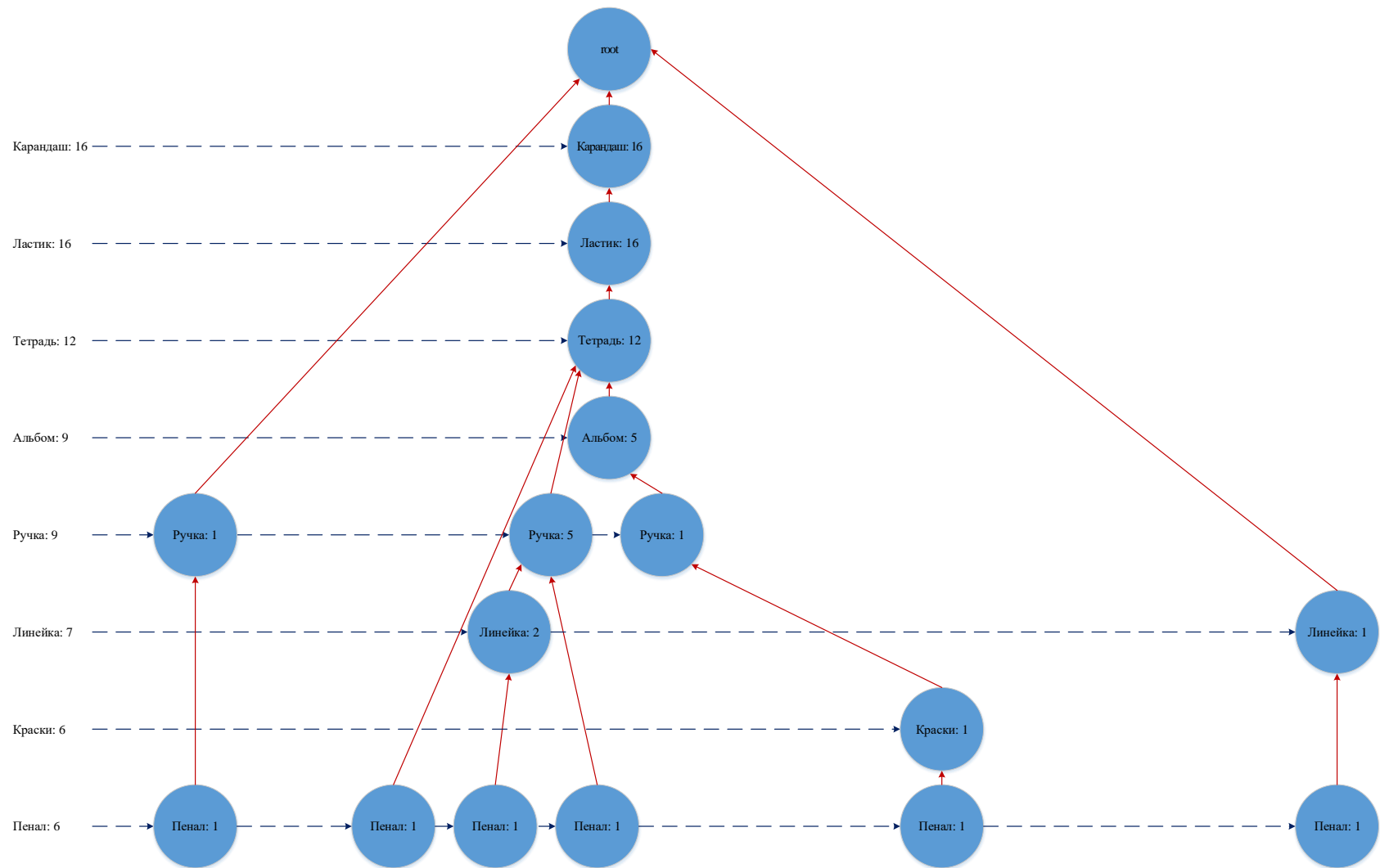


Рисунок 10. Префиксное поддереву путей, заканчивающихся предметом Пенал

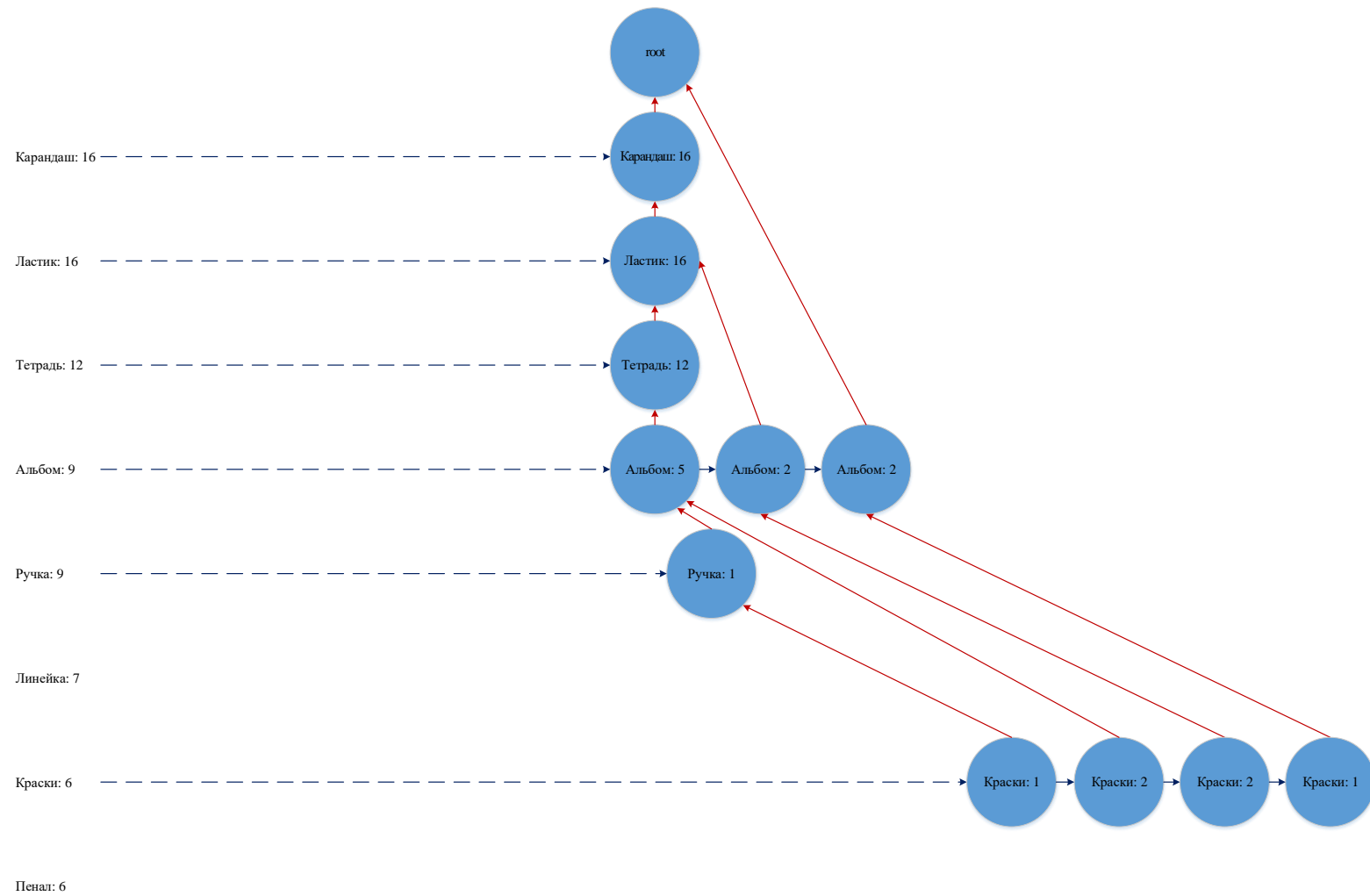


Рисунок 11. Префиксное поддерево путей, заканчивающихся предметом Краски

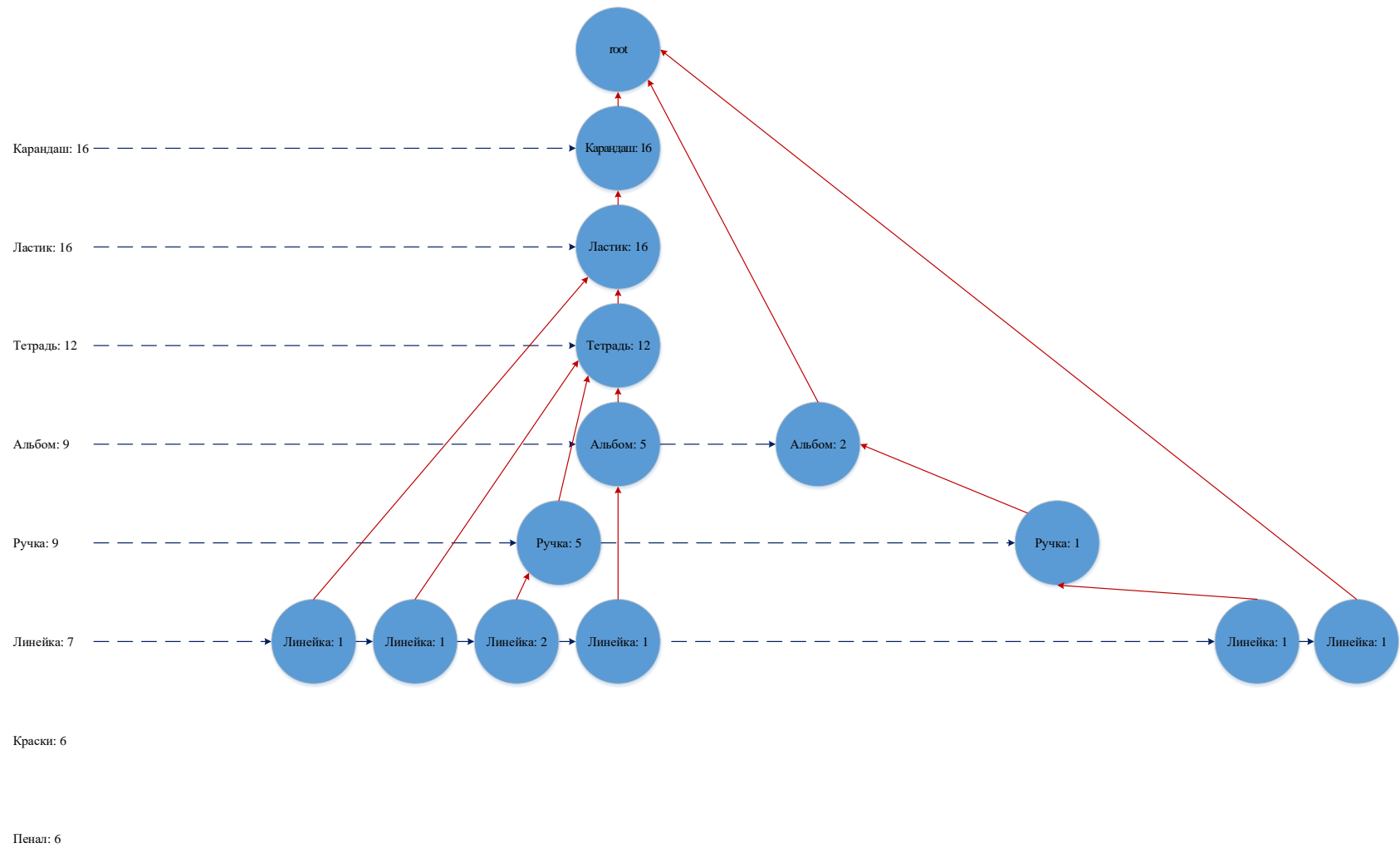


Рисунок 12. Префиксное поддереву путей, заканчивающихся предметом Линейка

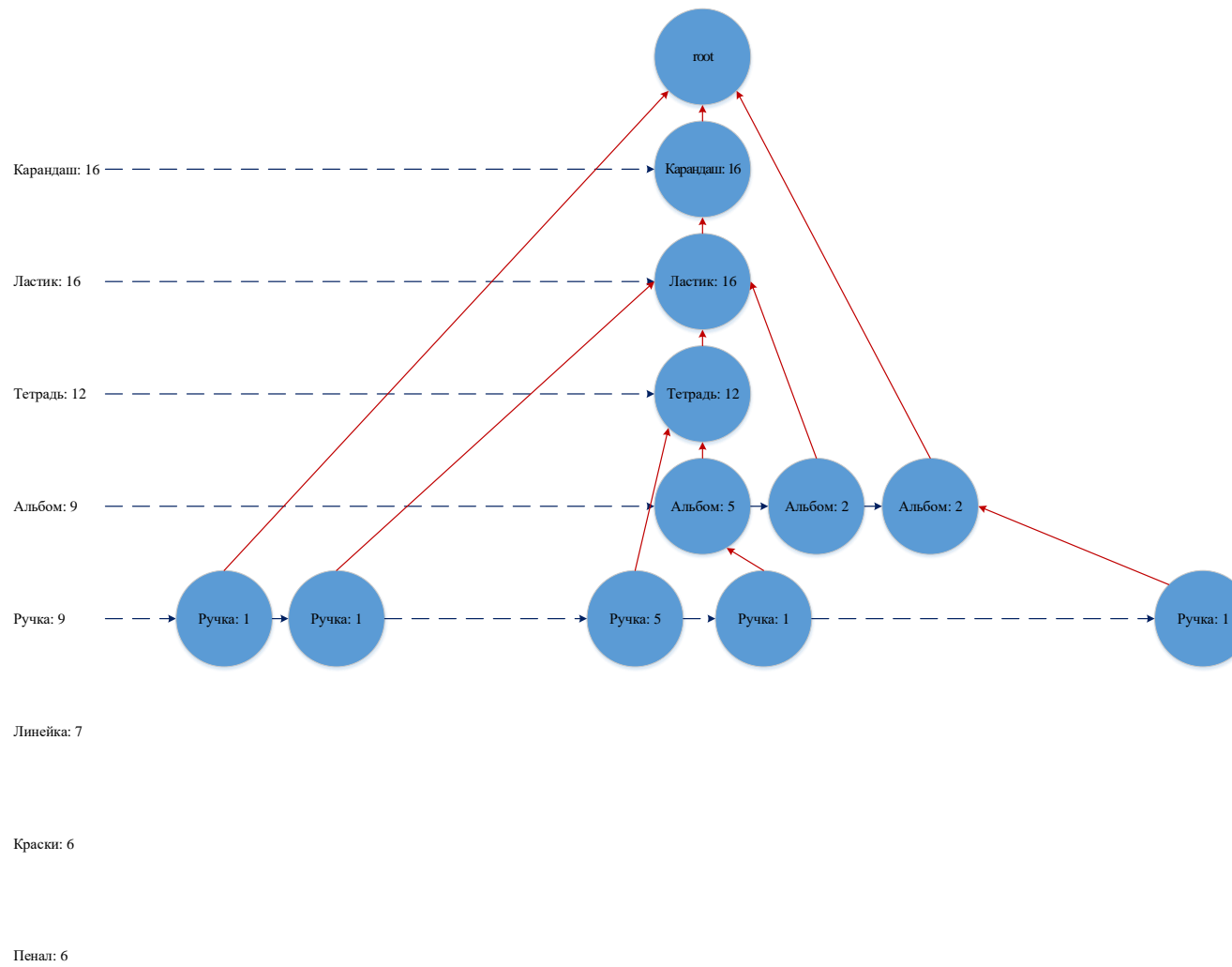


Рисунок 13. Префиксное поддереву путей, заканчивающихся предметом Ручка

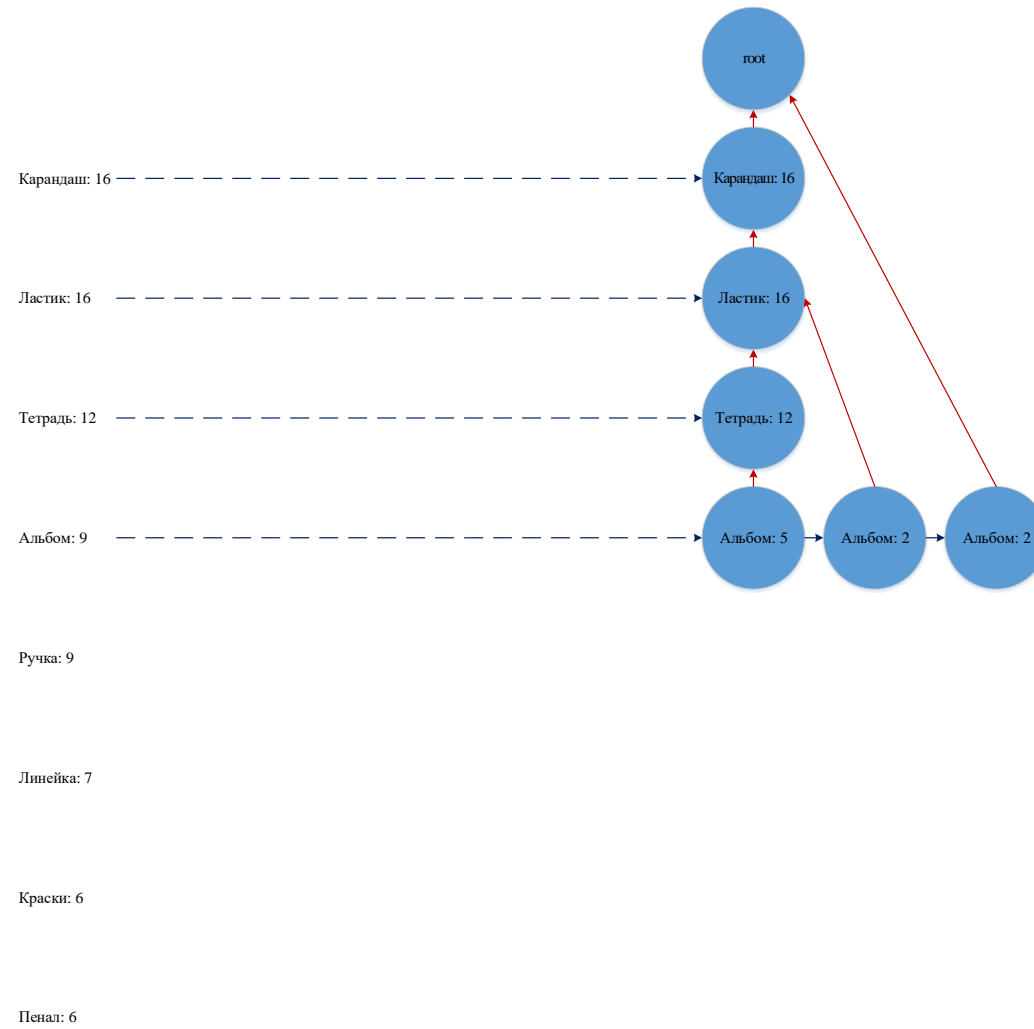


Рисунок 14. Префиксное поддереву путей, заканчивающихся предметом Альбом

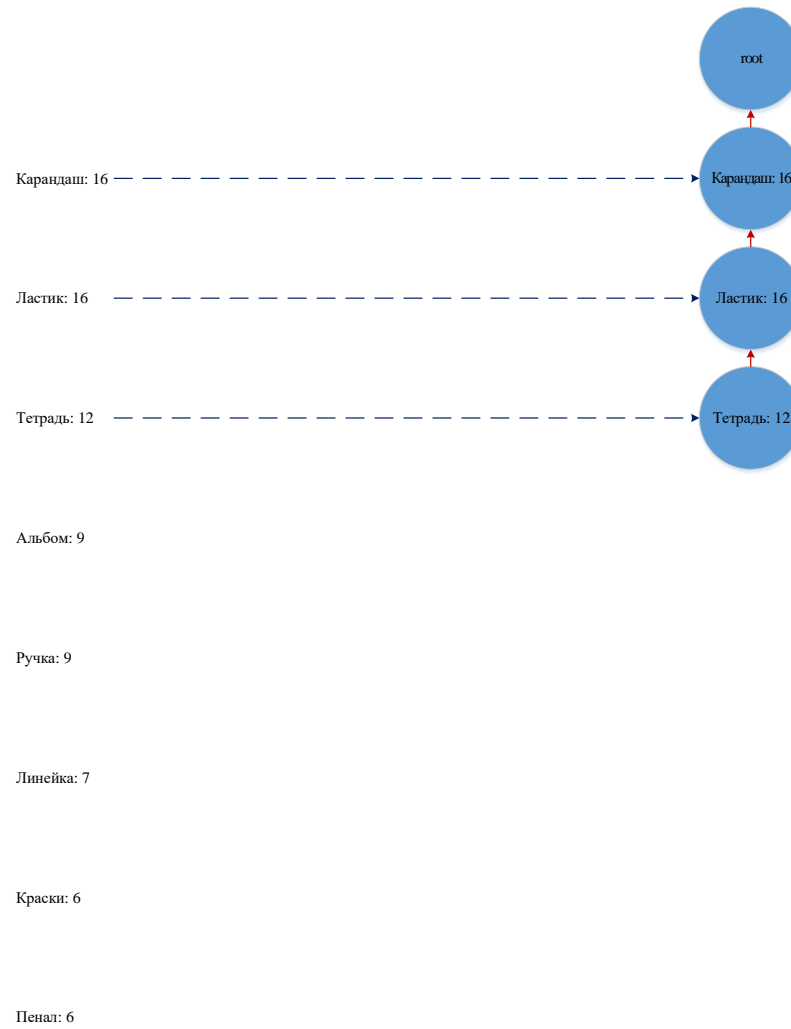


Рисунок 15. Префиксное поддереву путей, заканчивающихся предметом Тетрадь

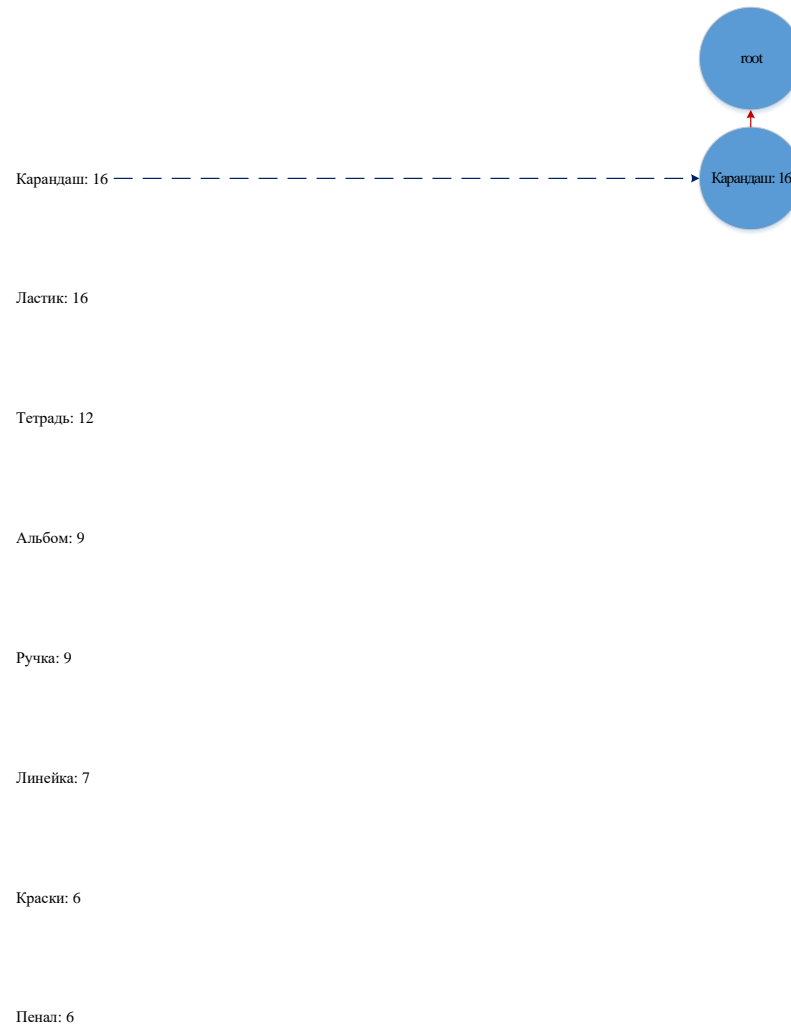


Рисунок 17. Префиксное поддерево путей, заканчивающихся предметом Карандаш

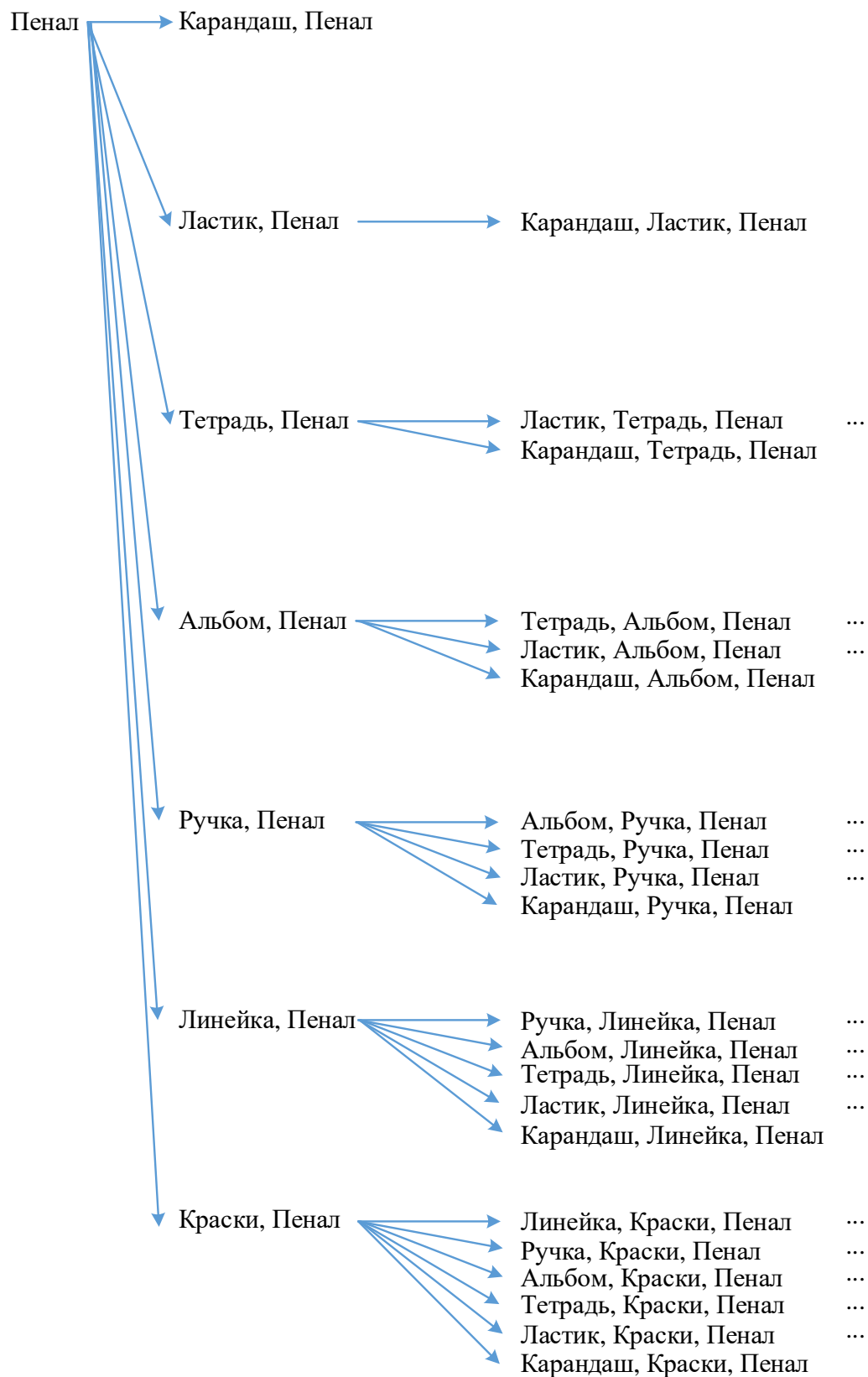


Рисунок 18. Варианты наборов предметов, потенциально определяемых путями, заканчивающимися предметом Пенал и подлежащих проверке на частотность

1. Извлечем из FP-дерева префиксные поддеревья путей, заканчивающиеся предметом или набором предметов.

Сначала осуществим поиск частых предметных наборов среди наборов, определяемых путями, заканчивающимися узлами, расположенными на самом нижнем уровне путями, заканчивающимися связанными узлами, расположенными на двух нижних уровнях и т.п. Аналогичным образом выполним анализ всех путей от листьев к корню.

Например, для FP-дерева, изображенного на рисунке 9, необходимо сначала проанализировать пути, заканчивающиеся предметом Пенал, затем – предметами Краски, Пенал, затем – предметами Линейка, Пенал и т.п.

В результате можно извлечь префиксные поддеревья путей, заканчивающиеся предметом или набором предметов. При этом целесообразно опираться на односвязанные списки, применяемые для связывания узлов дерева, соответствующих одному и тому же предмету.

Так, если в рассматриваемом примере было определено 8 частых предметов, то можно выделить 8 префиксных поддеревьев путей, имеющих в качестве листьев некоторый конкретный предмет (рисунки 10 – 17).

2. Обработаем каждое префиксное поддерево путей рекурсивно для извлечения частых предметных наборов, а затем объединим найденные решения.

На рисунке 18 показаны некоторые варианты наборов предметов, потенциально определяемых путями, заканчивающимися предметом Пенал и подлежащих проверке на частоту.

Рассмотрим префиксное поддерево путей для предмета Пенал (рисунок 10) и найдем все частые предметные наборы, содержащие этот предмет.

Предмет Пенал является частым, т.к. встречается в наборе транзакций 6 раз (это число можно также найти, просуммировав частоты, приписанные узлам дерева, соответствующим предмету Пенал).

Т.к. предмет Пенал является частым, проверим на частоту наборы, соответствующие путям, заканчивающимся предметом Пенал. Можно выписать 7 вариантов двухпредметных наборов (рисунок 18).

В таком случае можно говорить о рекурсивной декомпозиции задачи.

Построим *условные FP-деревья* (conditional FP-tree) для частых предметов (или частых наборов предметов).

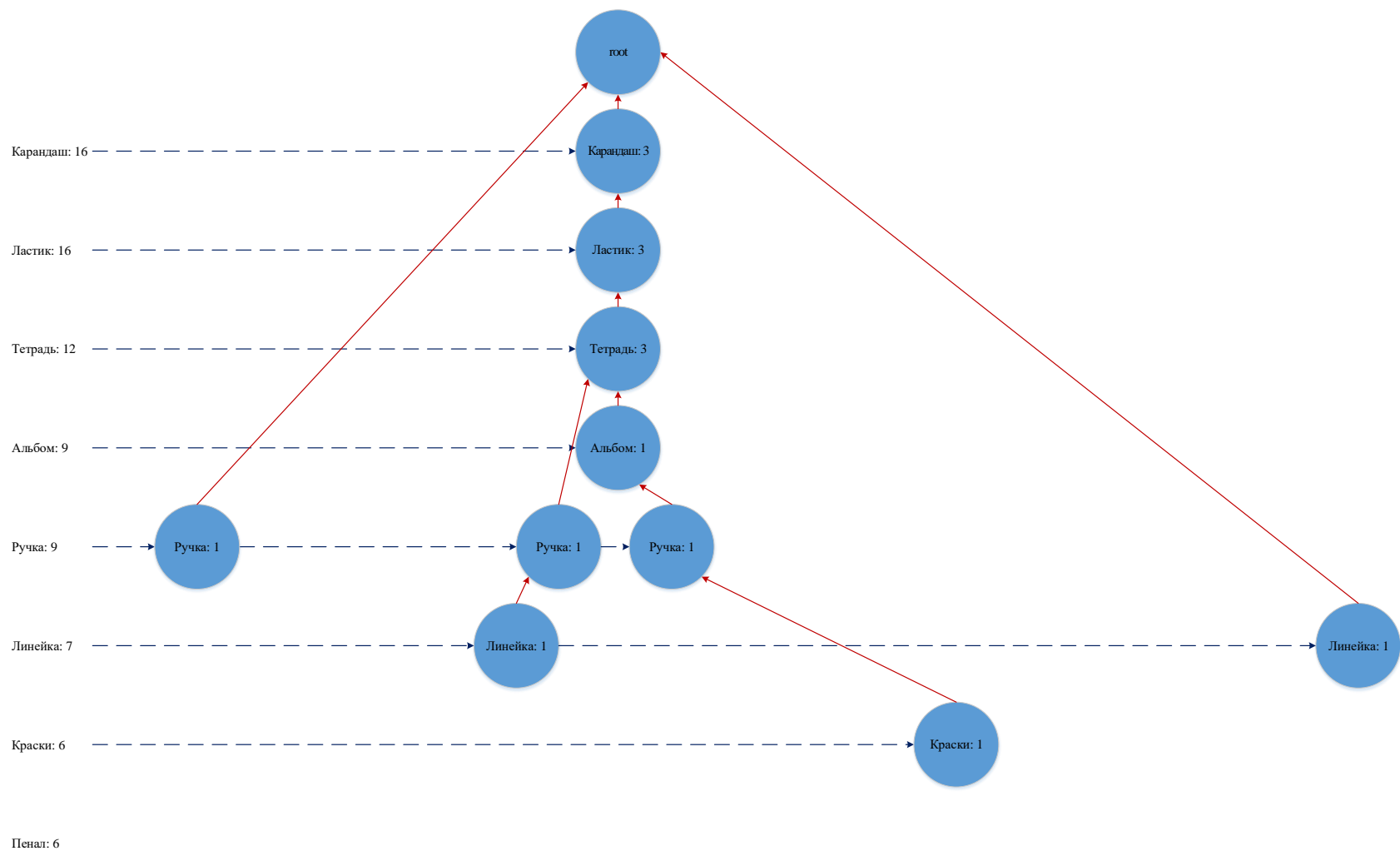


Рисунок 20. Формирование условного FP-дерева для предмета Пенал (продолжение): удаление узлов для предмета Пенал



Рисунок 21. Формирование условного FP-дерева для предмета Пенал (окончание): удаление узлов для нечастых предметов

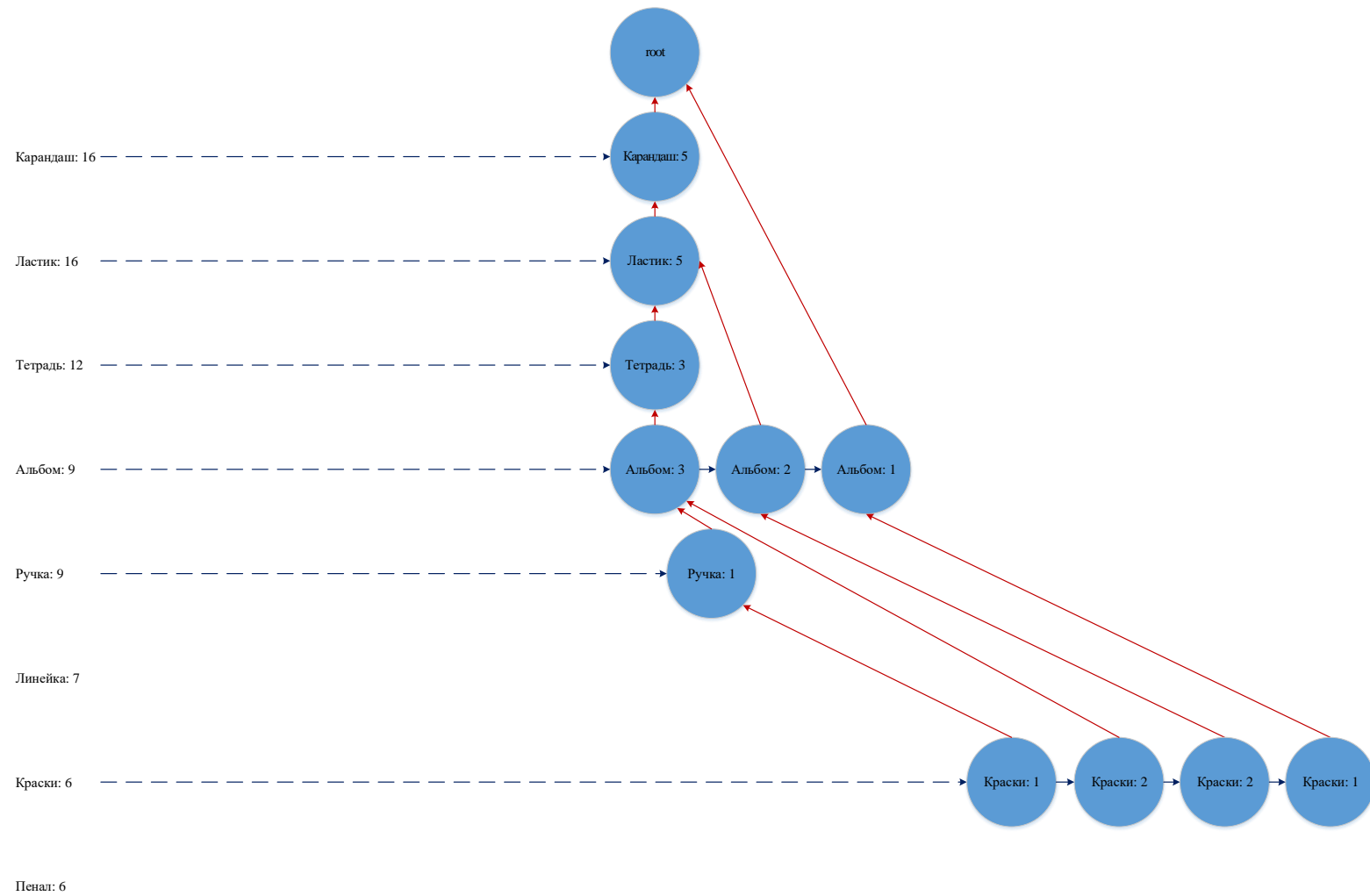


Рисунок 22. Формирование условного FP-дерева для предмета Краски (начало): обновление значений счетчиков поддержки

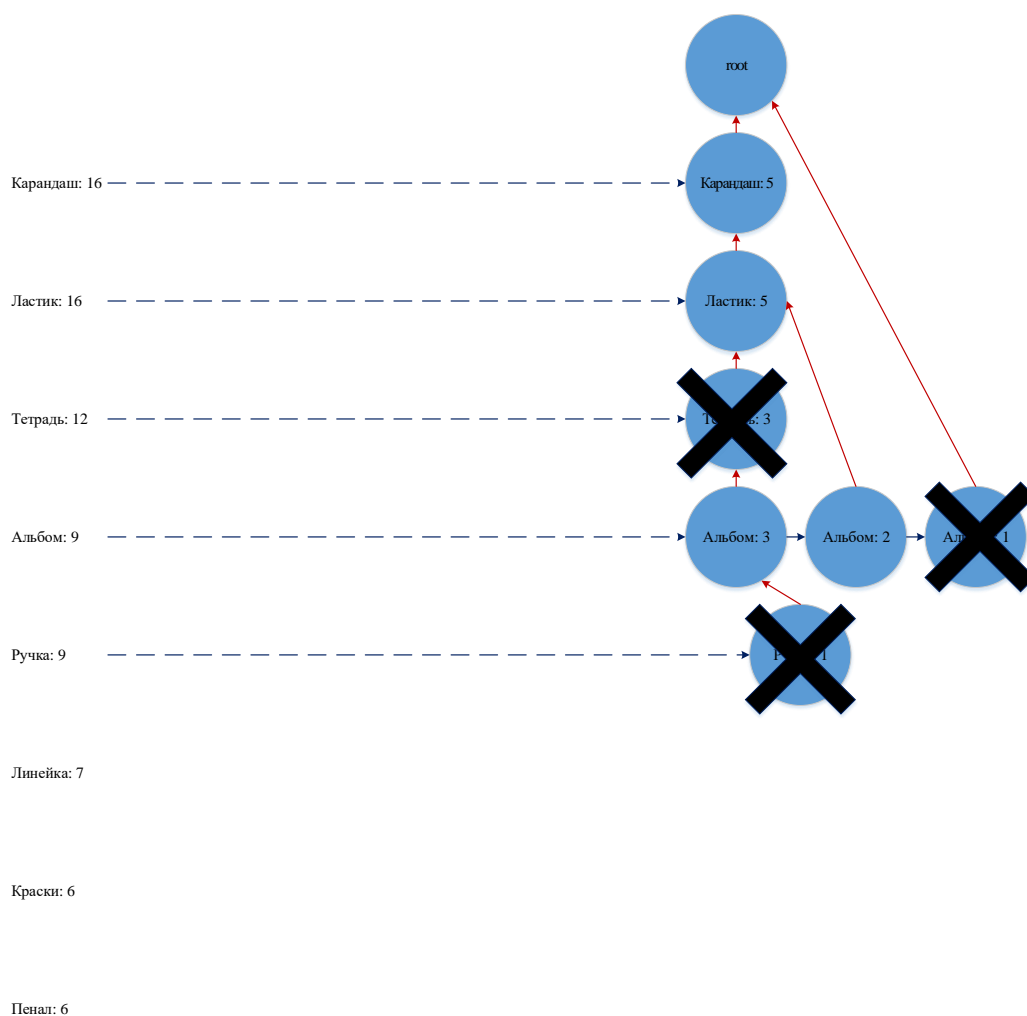


Рисунок 24. Формирование условного FP-дерева для предмета Краски (продолжение): удаление узлов для нечастых предметов

На рисунке 19 приведено префиксное дерева путей предмета Пенал с обновленными значениями счетчиков поддержки (изменены значения счетчиков поддержки у узлов с названиями предметов Карандаш, Ластик, Тетрадь, Альбом, Ручка (второй узел в ряду, считая слева), Линейка (первый узел в ряду, считая слева)).

2.2. Для извлечения условного FP-дерева из префиксного дерева путей необходимо удалить все узлы, соответствующего частому предмету (или частому набору предметов).

На рисунке 20 приведено префиксное дерева путей предмета Пенал с удаленными узлами для предмета Пенал.

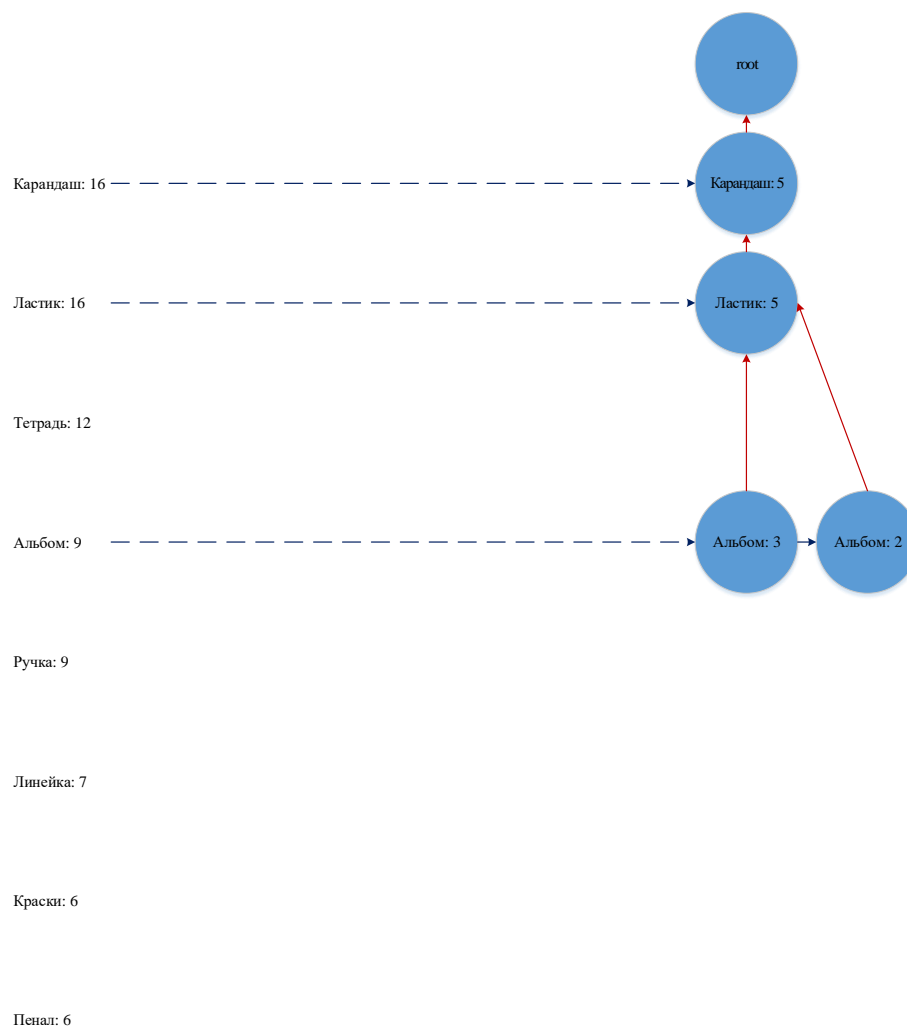


Рисунок 25. Формирование условного FP-дерева для предмета Краски (окончание): итоговое условное FP-дерево

2.3. Для извлечения условного FP-дерева из префиксного дерева путей необходимо удалить все узлы, соответствующего нечастым предметам (или нечастым наборам предметов).

На рисунке 21 приведено условное FP-дерево для предмета Пенал. Оно оказалось пустым, т.к. все предметы, соответствующие вершинам этого условного FP-дерева, не являются частыми.

Построим условные FP-деревья для остальных 7 частых предметов.

Рассмотрим предмет Краски. Построим для него условное FP-дерево.

На рисунках 22 – 25 показаны шаги формирования условного FP-дерева. При этом на рисунке 24 отмечены как вычеркнутые узлы предметов Тетрадь, Ручка, а также – узел предмета Альбом (правая ветвь от корня), т.к. суммарная поддержка которых меньше 5.

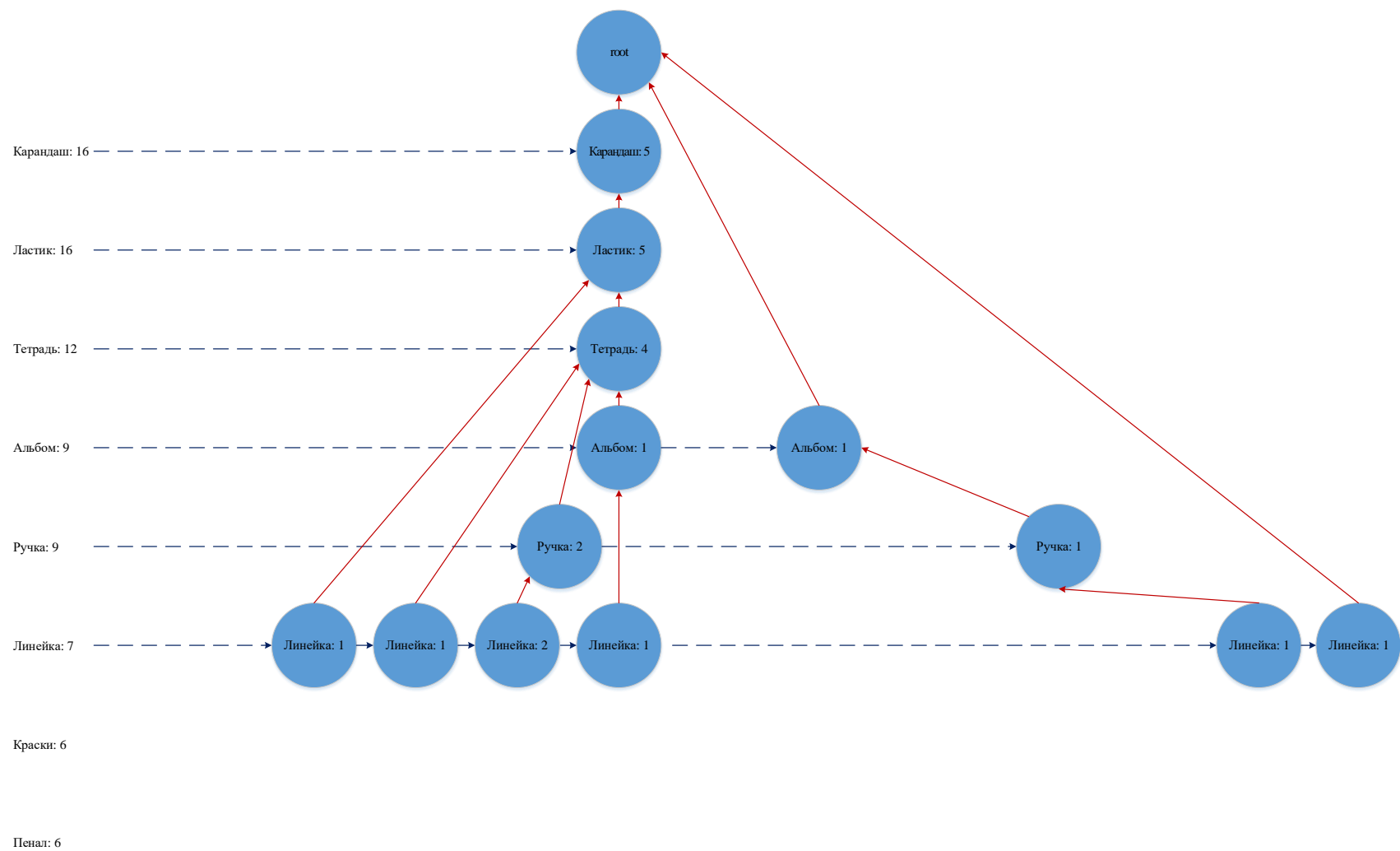


Рисунок 26. Формирование условного FP-дерева для предмета Линейка (начало): обновление значений счетчиков поддержки

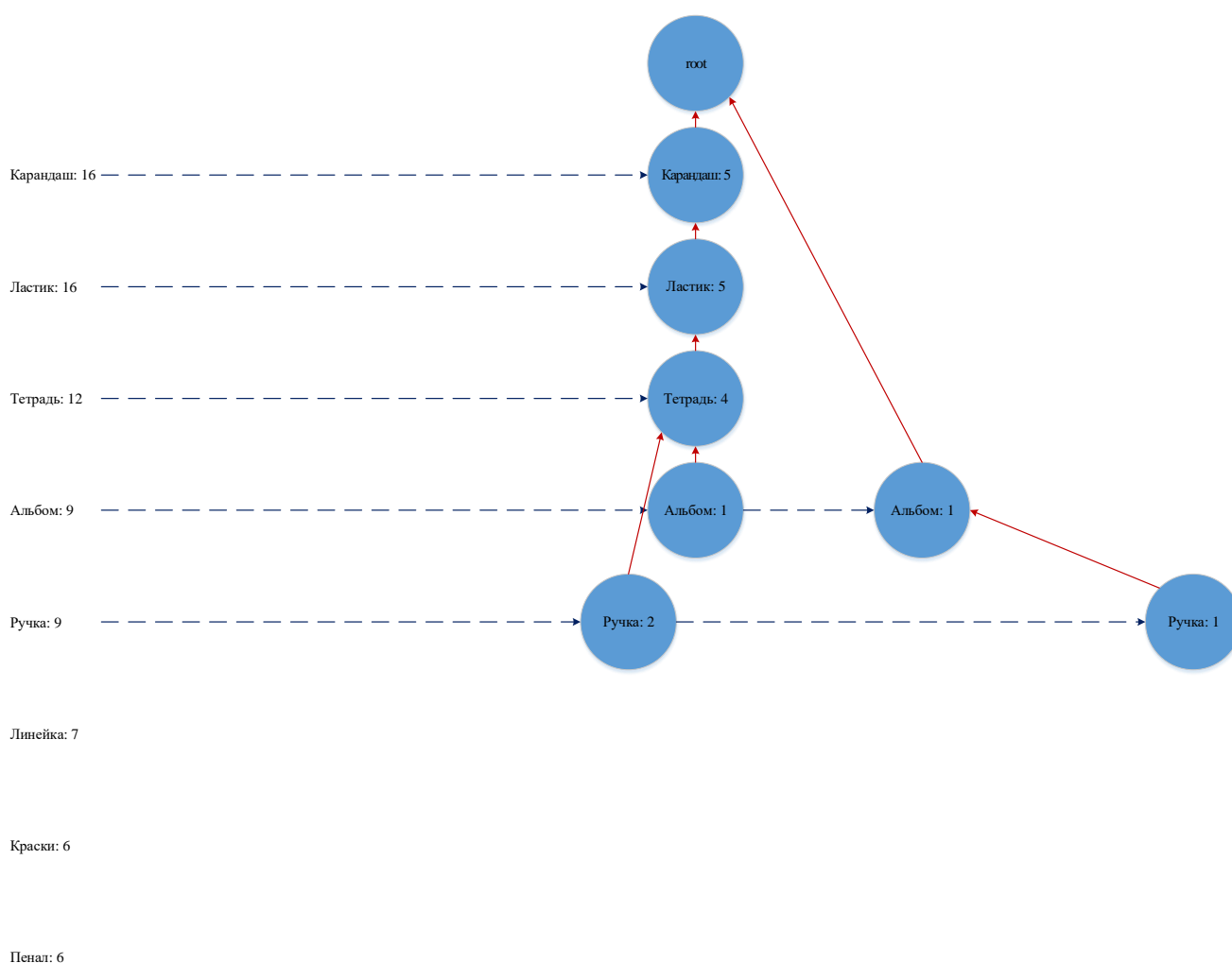


Рисунок 27. Формирование условного FP-дерева для предмета Линейка (продолжение): удаление узлов для предмета Линейка

В итоге частым предметным набором, найденным на основе условного FP-дерева для предмета Краски будет набор {Карандаш, Ластик, Альбом, Краски}.

Самый правый узел предмета Альбом, связанный с корнем дерева, удален потому, что при движении по этому пути (т.е. по пути от корня вправо) предмет Альбом не является частым.

Рассмотрим предмет Линейка. Построим для него условное FP-дерево.

На рисунках 23 – 29 показаны шаги формирования условного FP-дерева.

В итоге частым предметным набором, найденным на основе условного FP-дерева для предмета Линейка будет набор {Карандаш, Ластик, Линейка}.

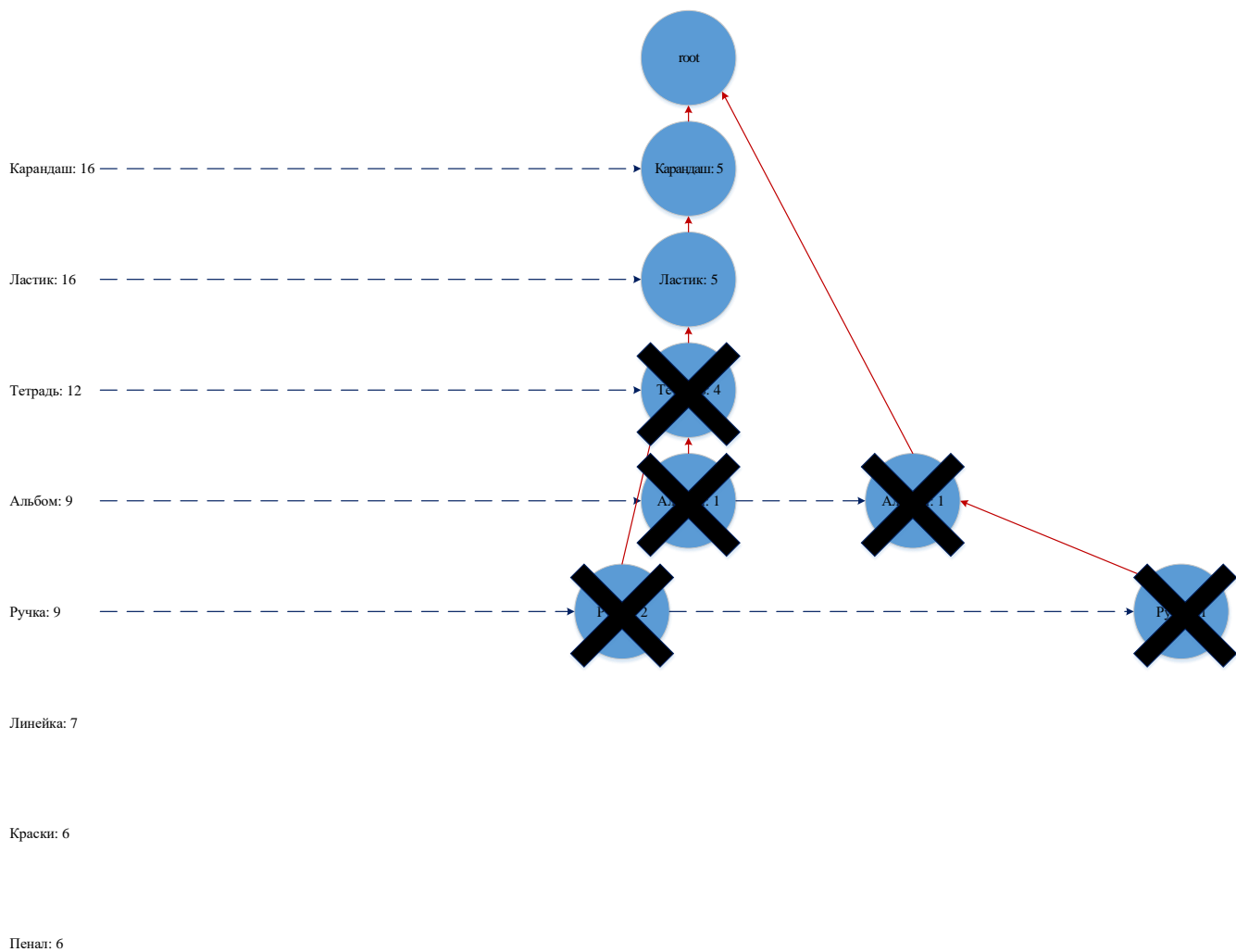


Рисунок 28. Формирование условного FP-дерева для предмета Линейка (продолжение): удаление узлов для нечастых предметов

Аналогичным образом можно выполнить формирование всех условных FP-деревьев сначала для частых однопредметных наборов на основе префиксных деревьев путей (рисунки 10 – 17), затем – для частых двухпредметных наборов и т.д., предполагая, что предметы в этих наборах соответствуют узлам деревьев в конце соответствующих путей.

В результате будут найдены все частые предметные наборы (двухпредметные, трехпредметные и четырехпредметные) в дополнение к уже найденным однопредметным частым наборам. При этом найденные частые наборы совпадут с наборами, найденными с применением алгоритма Apriori.

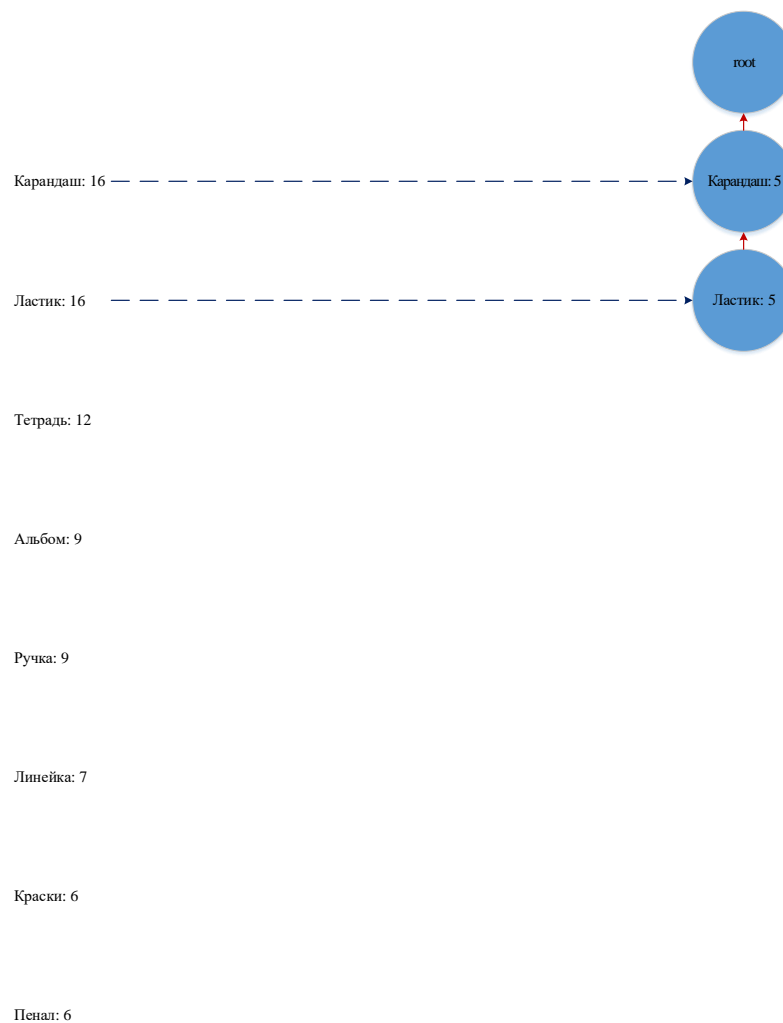


Рисунок 29. Формирование условного FP-дерева для предмета Линейка (окончание): итоговое условное FP-дерево

Множество частых двухпредметных наборов будет иметь вид:

$F_2 \{ \{ \text{Альбом, Карандаш} \}, \{ \text{Альбом, Краски} \}, \{ \text{Альбом, Ластик} \}, \{ \text{Альбом, Тетрадь} \}, \{ \text{Карандаш, Краски} \}, \{ \text{Карандаш, Ластик} \}, \{ \text{Карандаш, Линейка} \}, \{ \text{Карандаш, Ручка} \}, \{ \text{Карандаш, Тетрадь} \}, \{ \text{Краски, Ластик} \}, \{ \text{Ластик, Линейка} \}, \{ \text{Ластик, Ручка} \}, \{ \text{Ластик, Тетрадь} \}, \{ \text{Ручка, Тетрадь} \} \}.$

Множество частых трехпредметных наборов будет иметь вид:

$F_3 \{ \{ \text{Альбом, Карандаш, Краски} \}, \{ \text{Альбом, Карандаш, Ластик} \}, \{ \text{Альбом, Карандаш, Тетрадь} \}, \{ \text{Альбом, Краски, Ластик} \}, \{ \text{Альбом, Ластик, Тетрадь} \}, \{ \text{Карандаш, Краски, Ластик} \}, \{ \text{Карандаш, Ластик, Линейка} \}, \{ \text{Карандаш, Ластик, Ручка} \}, \{ \text{Карандаш, Ластик, Тетрадь} \}, \{ \text{Карандаш, Ручка, Тетрадь} \}, \{ \text{Ластик, Ручка, Тетрадь} \} \}.$

Множество частых четырехпредметных наборов будет иметь вид:

$F_4 \{ \{ \text{Альбом, Карандаш, Краски, Ластик} \},$

$\{ \text{Альбом, Карандаш, Ластик, Тетрадь} \},$

$\{ \text{Карандаш, Ластик, Ручка, Тетрадь} \} \}.$

Сравнительный анализ результатов работы алгоритмов Apriori и FPGrowth показывает, что с увеличением числа транзакций в БД временные затраты на поиск частых предметных наборов растут для алгоритма FPGrowth намного медленнее, чем для алгоритма Apriori.