

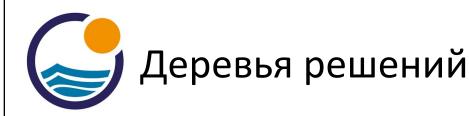
Деревья решений (decision trees)





- Следующие три раздела посвящены методам на основе деревьев.
- Три основных метода:
 - Деревья решений (decision trees)
 - Случайные леса (random forests)
 - Расширяемые деревья (boosted trees)





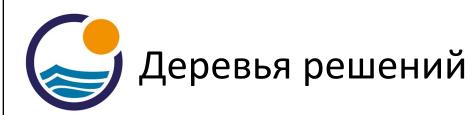
- Все эти три метода основаны на базовом алгоритме деревьев решений.
- Мы пройдём все три метода, каждый в своём разделе курса, и в самом конце уже будут проверочные упражнения.





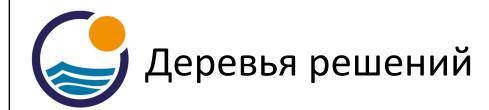
Деревья решений История





- Хотя базовые деревья решений использовались для принятия решений уже очень длительное время, статистические деревья решений были разработаны относительно недавно.
- Обратите внимание, это разные вещи.



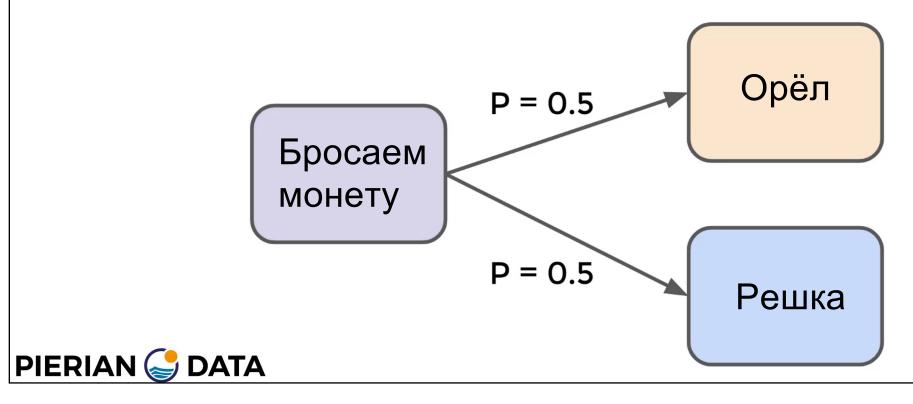


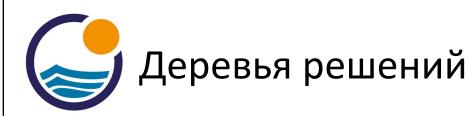
• Общий термин "дерево решений" можно описать как схему, отображающую выбор решения:





 Общий термин "дерево решений" можно описать как схему, отображающую выбор решения:





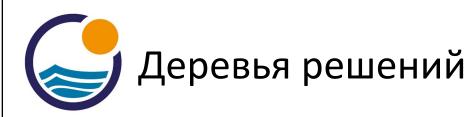
- Обучение деревьев решений это статистическое моделирование с применением деревьев, где решения в каждом узле принимаются на основе некоторого условия (некоторой метрики).
- Давайте рассмотрим шаги, которые привели к возможности создавать предсказания на основе деревьев.



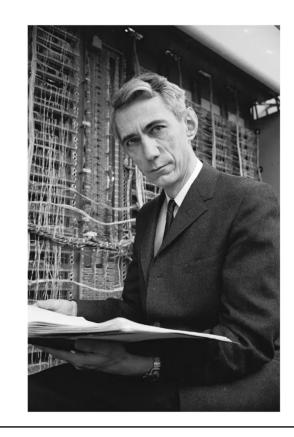


- В двух словах деревья решений и похожие методы разбивают данные на части на основе информации, содержащейся в признаках.
- Это значит, что нам нужно математическое определение термина информация, и способность измерить её.





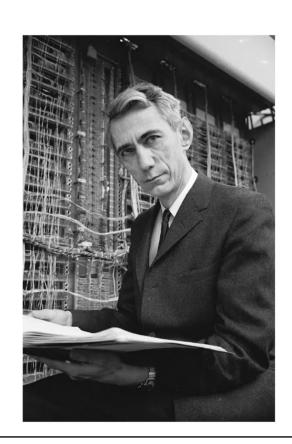
- Клод Шеннон считается отцом теории информации.
- В 1948 году в журнале Bell System Technical Journal опубликовал работу "Математическая теория коммуникаций".



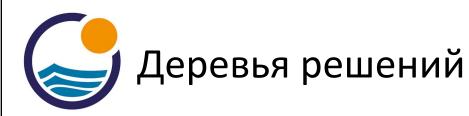




- Работал в разных сферах:
 - Проектирование микросхем
 - Криптография
 - Носимые компьютеры
 - Искусственный интеллект

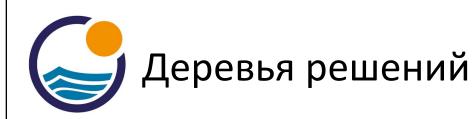






- Способность измерить информацию станет важна, когда мы рассмотрим математические основы построения деревьев решений.
- Мы займёмся этим позже, а пока посмотрим на этапы развития деревьев решений.





 1963: Морган, Санквист - первая публикация регрессионного алгоритма деревьев









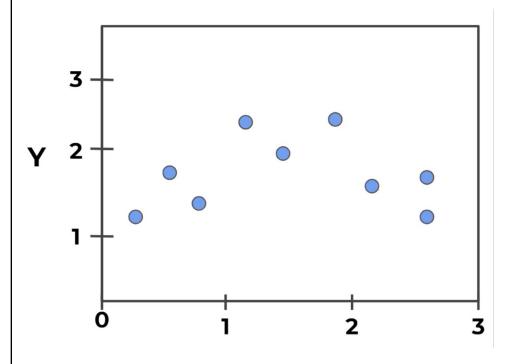
• 1963: piecewise-constant regression tree



X

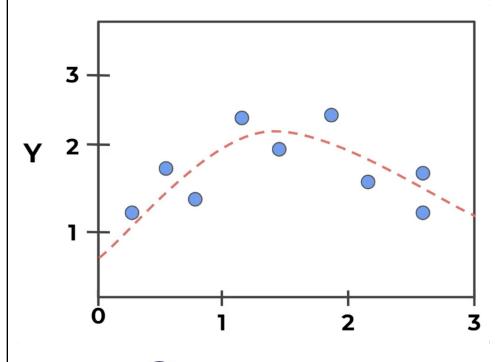






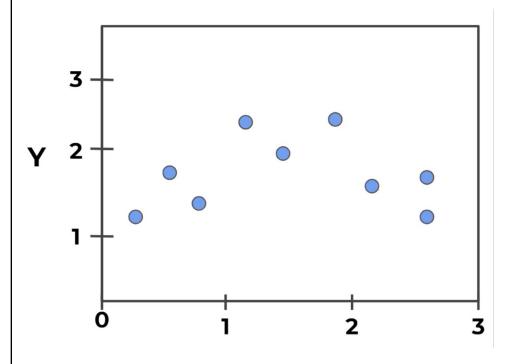








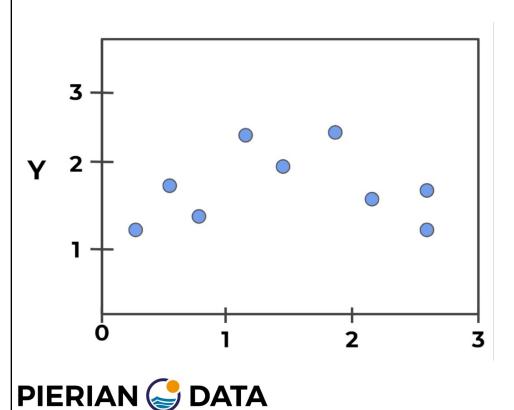






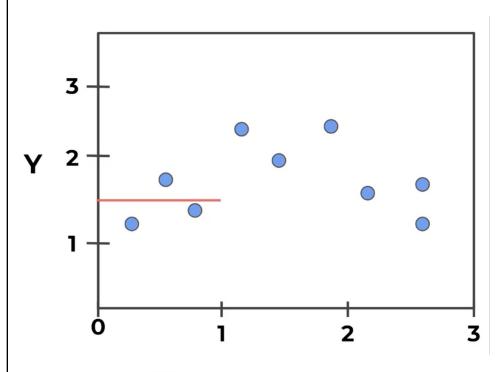


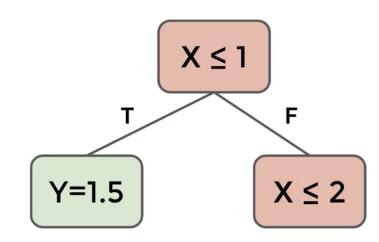
• 1963: piecewise-constant regression tree



X < 1

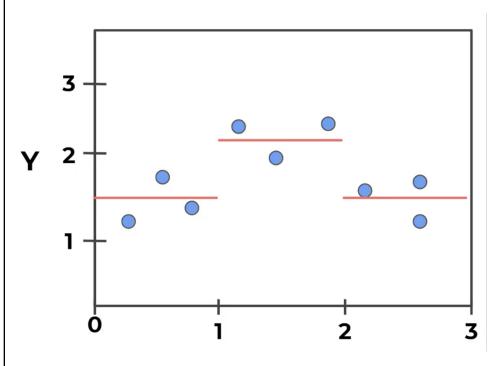


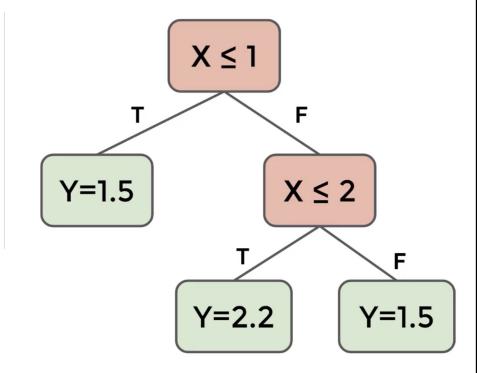






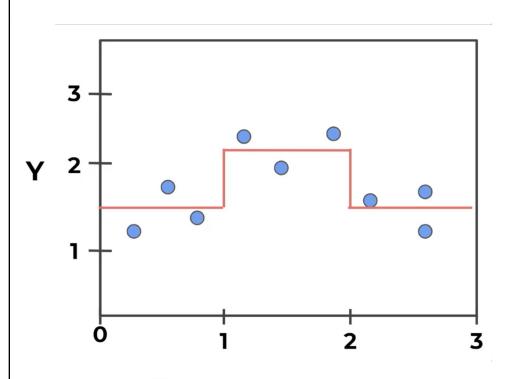


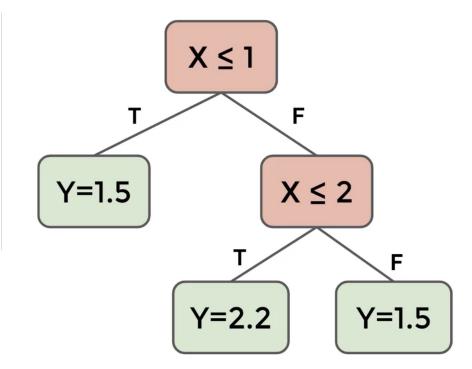












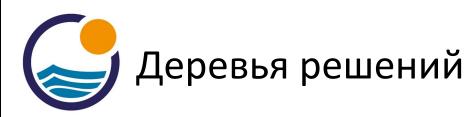




 В статье 1963 года разбиения в каждом узле дерева основывались на "загрязнении узла" (node impurity), которое по сути определялось как метрика ошибки:

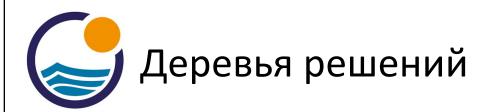
$$\phi(t) = \sum_{i \in t} (y_i - \bar{y})^2$$





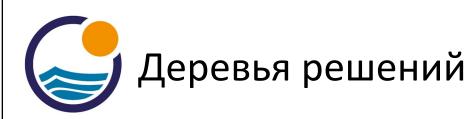
- 1972: Роберт Мессенджер и Льюис Манделл публикуют первый алгоритм классификации с применением деревьев.
- Условие разделения в узлах было названо "Theta Automatic Interaction Detection" (THAID)





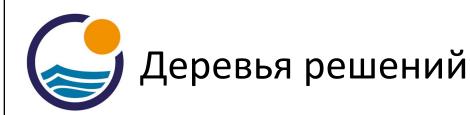
- 1972: Роберт Мессенджер и Льюис Манделл публикуют первый алгоритм классификации с применением деревьев.
- Условие разделения в узлах было названо "Theta Automatic Interaction Detection" (THAID)
- 1980: Гордон Кисс публикует развитие алгоритма CHAID Chi-square automatic interaction detection.





• 1970-е годы: Лео Брейман и Чарльз Стоун из Беркли, Джером Фридман и Ричард Ольшен из Стэнфорда начали разработку алгоритмов CART — Classification and Regression tree.



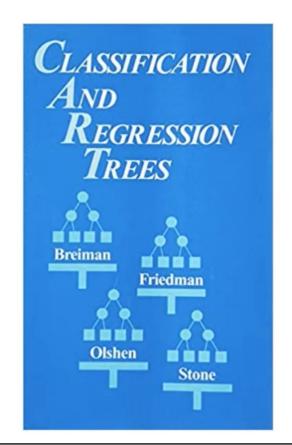


- 1970-е годы: Лео Брейман и Чарльз Стоун из Беркли, Джером Фридман и Ричард Ольшен из Стэнфорда начали разработку алгоритмов CART — Classification and Regression tree.
- 1984: Они опубликовали книгу по алгоритмам CART, включая детали реализации на компьютере.
- Методы на основе CART стали стандартом (и в Scikit-Learn!)

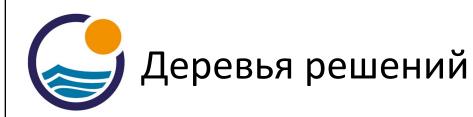




- В книге CART были предложены многие концепции:
 - Кросс-валидация деревьев
 - Усечение деревьев (pruning)
 - Суррогатные разбиения
 - О Оценки важности переменных
 - Поиск линейных разделений

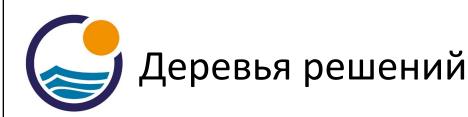






- 1986: Джон Росс Квинлан разработал алгоритм деревьев решений ID3 на основе метрики "gain ratio".
- 1990-е: улучшения ID3 С4.5 (всё ещё популярен)
- 2000-е: выпущена оптимизированная коммерческая версия C5.0 с различными улучшениями.





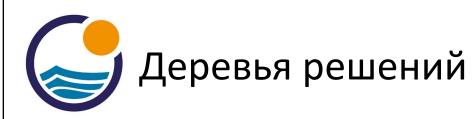
- Многие улучшения основного алгоритма были перенесены на другие методы с применением деревьев – случайные леса и расширяемые деревья.
- Давайте посмотрим на фундаментальную идею деревьев решений!





Деревья решений Терминология



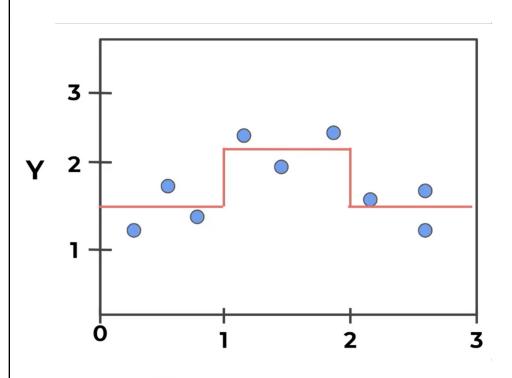


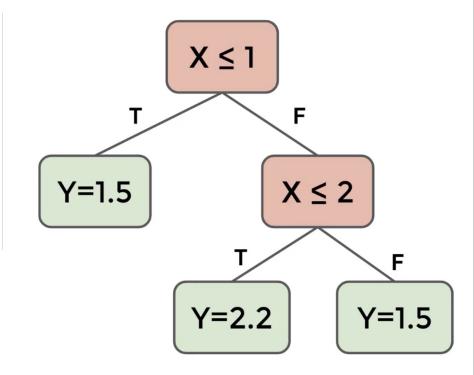
• Для начала обсудим различные термины, которые встречаются при обсуждении деревьев...





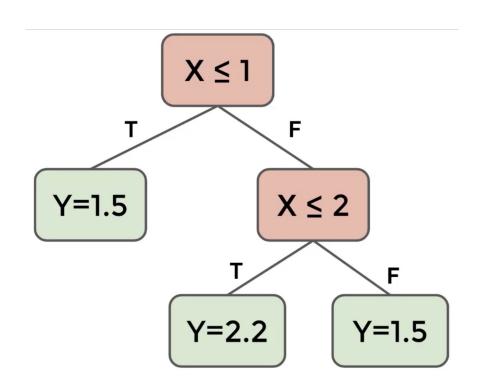
• Вспомним наше простое дерево:





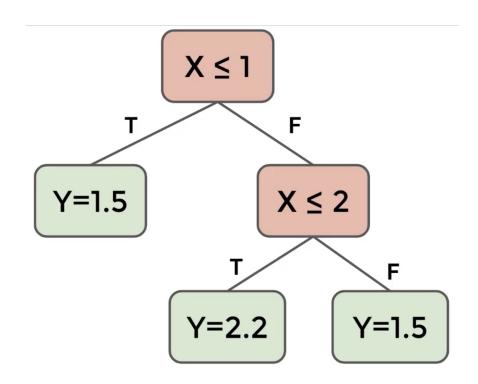


• Вспомним наше простое дерево:

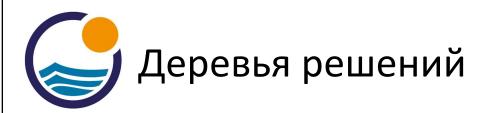




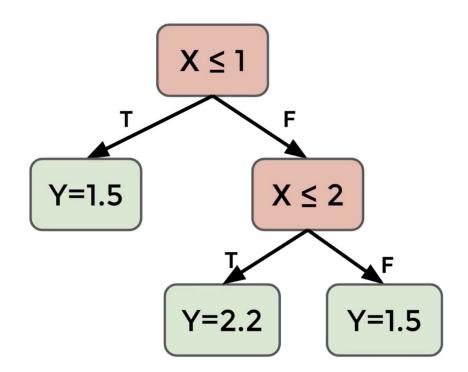
• Разбиение (splitting):







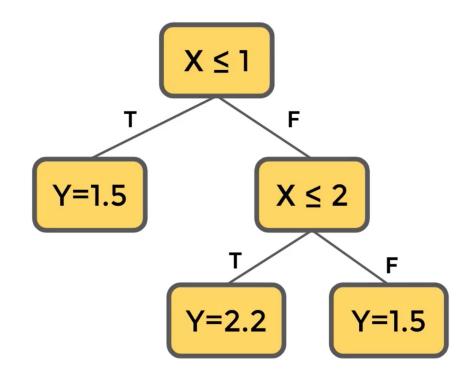
Разбиение (splitting):







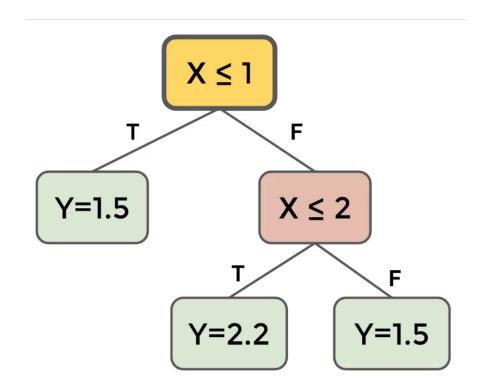
Узлы (nodes):





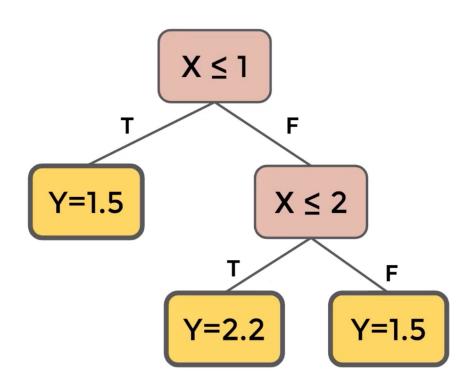


• Корневой узел (root node):





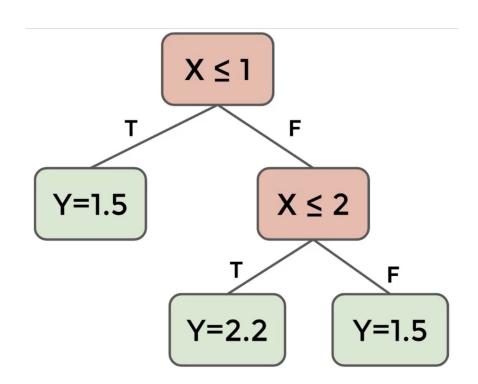
• Листовые / конечные узлы (leaf / terminal nodes):



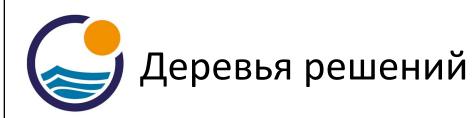




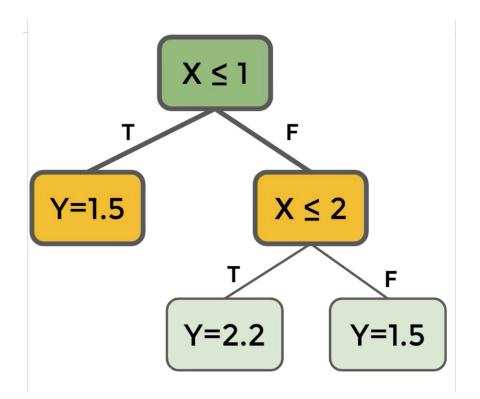
• Родительские/дочерние узлы (parent/children nodes):







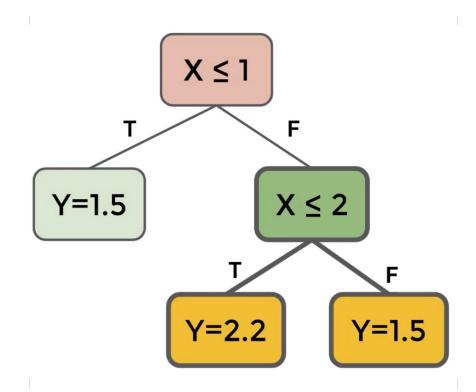
• Родительские/дочерние узлы (parent/children nodes):





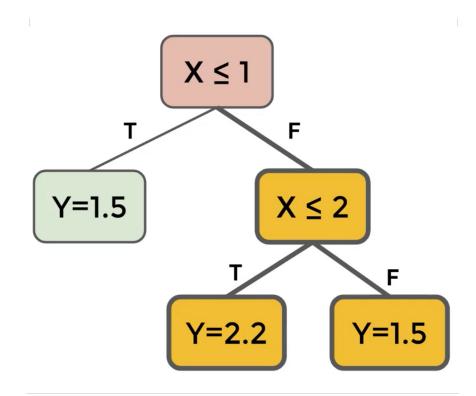


• Родительские/дочерние узлы (parent/children nodes):



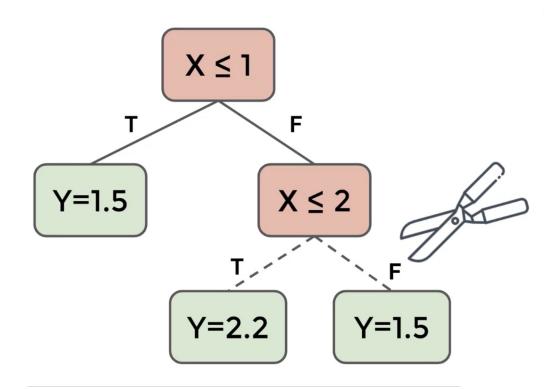


• Ветви / поддеревья (tree branches / sub trees):



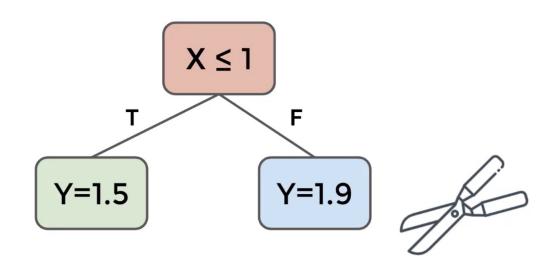


• Усечение (pruning)

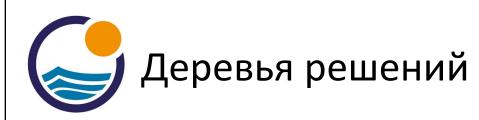




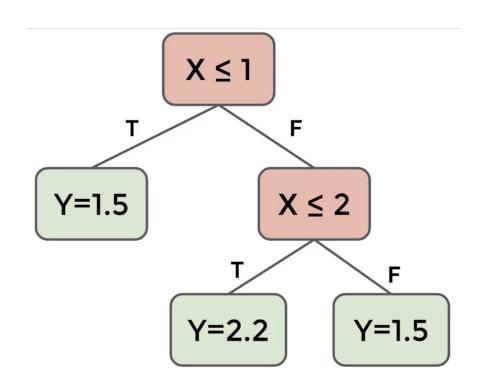
• Усечение (pruning)







• Далее мы перейдём к построению дерева!

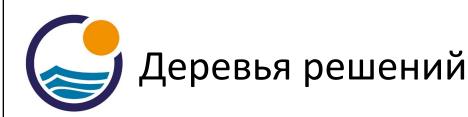






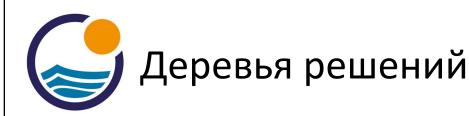
Деревья решений Gini impurity





• Прежде чем мы узнаем, как при построении деревьев выбираются критерии разбиения, давайте взглянем на наиболее часто используемую метрику измерения информации для деревьев решений — gini impurity (загрязнение Джини, индекс Джини, примесь Джини).





- Метрика "gini impurity" измеряет, насколько информация в наборе данных является чистой ("pure").
- С точки зрения задач классификации, мы можем воспринимать это как однородность классов.
- Посмотрим как это выглядит на простейшем примере двух классов...





- Метрика "gini impurity" для классификации
 - Для набора классов С и заданного набора данных Q:

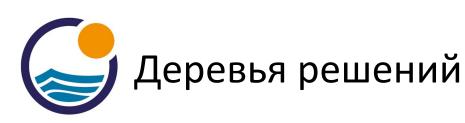
$$G(Q) = \sum_{c \in C} p_c (1-p_c)$$



- Метрика "gini impurity" для классификации
 - Для набора классов С и заданного набора данных Q,
 р это вероятность класса С.

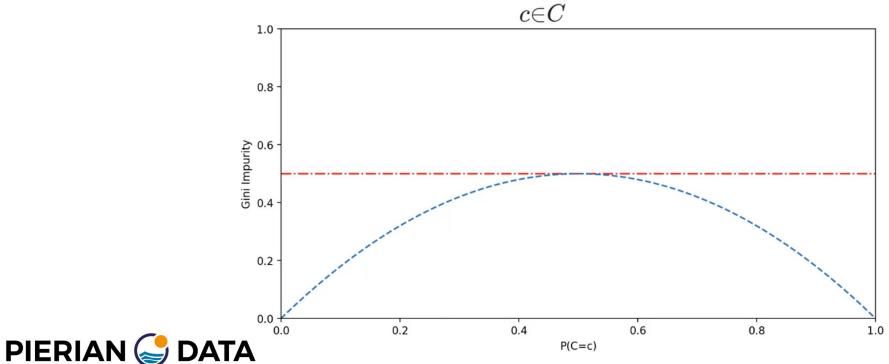
$$p_c = rac{1}{N_Q} \sum_{x \in O} \mathbb{1}(y_{class} = c) \qquad G(Q) = \sum_{c \in C} p_c (1 - p_c)$$





Метрика "gini impurity" для классификации

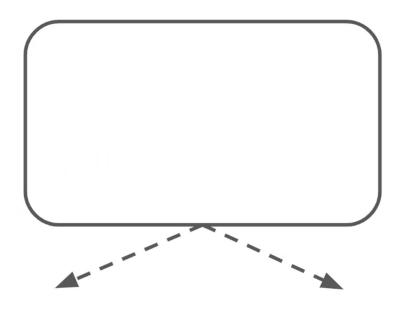
$$G(Q) = \sum_{c \in C} p_c (1-p_c)$$





Метрика "gini impurity" для классификации

$$G(Q) = \sum_{c \in C} p_c (1-p_c)$$

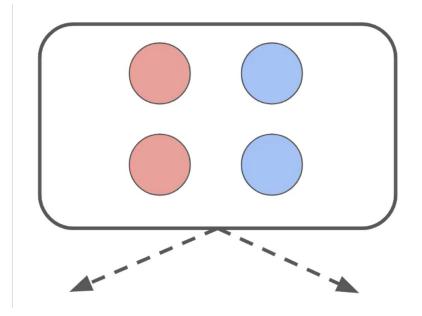






Метрика "gini impurity" для классификации

$$G(Q) = \sum_{c \in C} p_c (1-p_c)$$



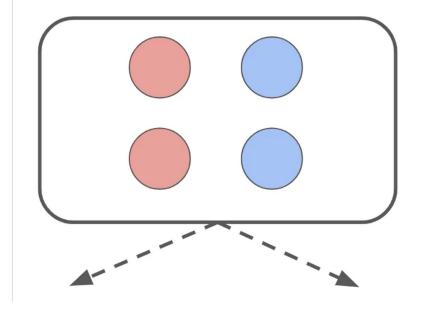




Метрика "gini impurity" для классификации

$$G(Q) = \sum_{c \in C} p_c (1-p_c)$$

Класс "красный" (2/4)(1 - 2/4) = 0.25

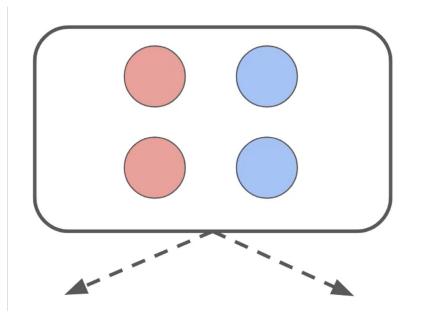






Метрика "gini impurity" для классификации

$$G(Q) = \sum_{c \in C} p_c (1-p_c)$$



Класс "красный" (2/4)(1 - 2/4) = 0.25

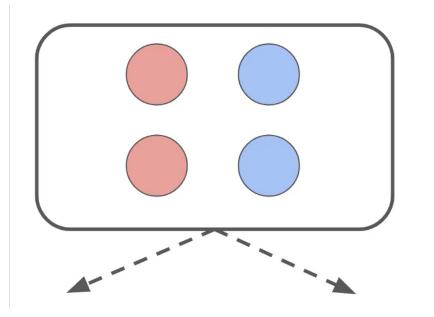
Класс "синий" (2/4)(1-2/4)=0.25





Метрика "gini impurity" для классификации

$$G(Q) = \sum_{c \in C} p_c (1-p_c)$$



Класс "красный" (2/4)(1 - 2/4) = 0.25



Класс "синий" (2/4)(1-2/4)=0.25



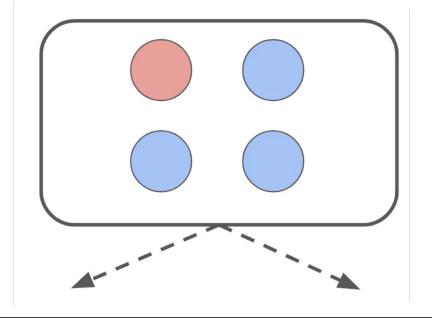
Gini Impurity 0.25 + 0.25 = 0.5





• Данные более чистые (меньше "impurity")

$$G(Q) = \sum_{c \in C} p_c (1-p_c)$$



Класс "красный" (1/4)(1 - 1/4) = 0.1875



Класс "синий" (3/4)(1 - 3/4) = 0.1875



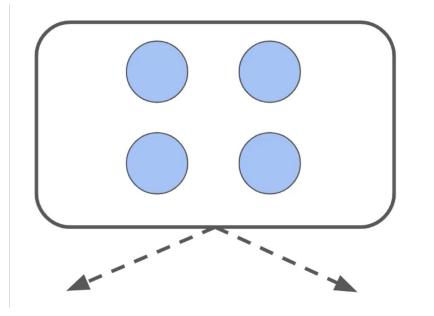
Gini Impurity 0.1875+0.1875 = 0.375





Данные полностью чистые (нулевой "impurity")

$$G(Q) = \sum_{c \in C} p_c (1-p_c)$$



Класс "красный" (0/4)(1 - 0/4) = 0

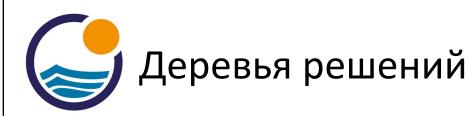
4

Класс "синий" (4/4)(1-4/4)=0



Gini Impurity 0 + 0 = 0





- Если цель дерева решений в том, чтобы отделить классы друг от друга, то мы можем выполнять разбиения данных на основе метрики gini impurity.
- Мы хотим минимизировать gini impurity на листьях.
- Если на листьях gini impurity минимально, то это будет означать, что мы удачно разделили классы.





- В следующей лекции мы построим пример применения gini impurity, вычисляя его на основе набора данных.
- Далее мы рассмотрим различные типы разбиения признаков и определения того, какой признак будет корневым узлом.





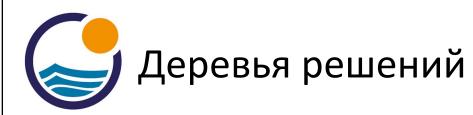
Gini impurity в деревьях – часть 1





- Начнём разбираться в том, как лучше выбрать порядок узлов и как выполнять разбиение внутри дерева.
- Для начала посмотрим, как построить дерево на основе некоторого набора данных, используя метрику gini impurity.





- Начиная строить дерево, нам нужно выбрать признак, который будет использоваться для корневого узла.
- Мы можем применить метрику gini impurity для сравнения информации, содержащейся внутри признаков в наборе данных.
- Рассмотрим этот момент подробнее...





- Метрика "gini impurity" для классификации
 - Для набора классов С и заданного набора данных Q:

$$G(Q) = \sum_{c \in C} p_c (1-p_c)$$



- Метрика "gini impurity" для классификации
 - Для набора классов С и заданного набора данных Q,
 р это вероятность класса С.

$$p_c = rac{1}{N_Q} \sum_{x \in Q} \mathbb{1}(y_{class} = c) \hspace{0.5cm} G(Q) = \sum_{c \in C} p_c (1 - p_c)$$

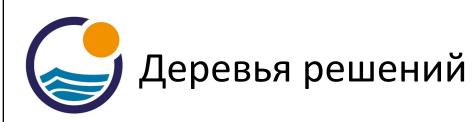




- Метрика "gini impurity" для классификации
 - Для набора классов С и заданного набора данных Q,
 р это вероятность класса С.

$$p_c = egin{pmatrix} rac{1}{N_Q} \sum_{x \in Q} \mathbb{1}(y_{class} = c) & G(Q) = \sum_{c \in C} p_c (1 - p_c) \end{pmatrix}$$





• Рассмотрим следующий набор данных:

X - URL Link	Y-Spam
Yes	Yes
Yes	Yes
No	No
No	No
No	Yes
No	No
Yes	No





• Попытаемся спрогнозировать – спам или нет:

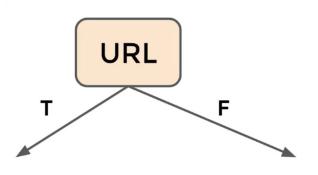
X - URL Link	Y-Spam
Yes	Yes
Yes	Yes
No	No
No	No
No	Yes
No	No
Yes	No







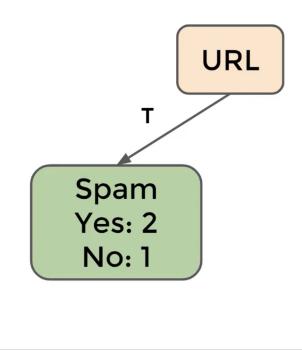
X - URL Link	Y-Spam
Yes	Yes
Yes	Yes
No	No
No	No
No	Yes
No	No
Yes	No







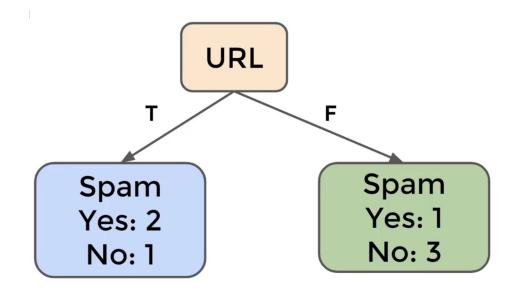
X - URL Link	Y-Spam
Yes	Yes
Yes	Yes
No	No
No	No
No	Yes
No	No
Yes	No







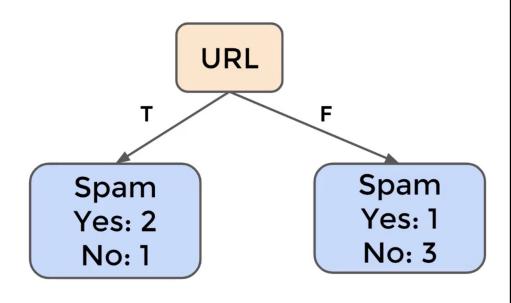
X - URL Link	Y-Spam
Yes	Yes
Yes	Yes
No	No
No	No
No	Yes
No	No
Yes	No







X - URL Link	Y-Spam
Yes	Yes
Yes	Yes
No	No
No	No
No	Yes
No	No
Yes	No

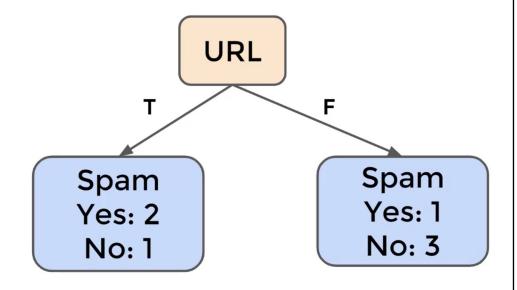






• Вспомним формулу gini impurity:

X - URL Link	Y-Spam	
Yes	Yes	
Yes	Yes	
No	No	
No	No	
No	Yes	
No	No	
Yes	No	



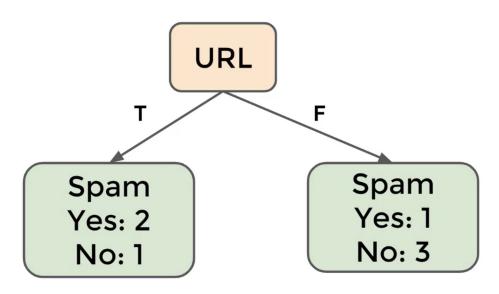
$$G(Q) = \sum_{c \in C} p_c (1-p_c)$$





- Пусть "спам" и "не-спам" будут классами С
- Левый листовой узел:

 $(\frac{2}{3})(1-\frac{2}{3})$



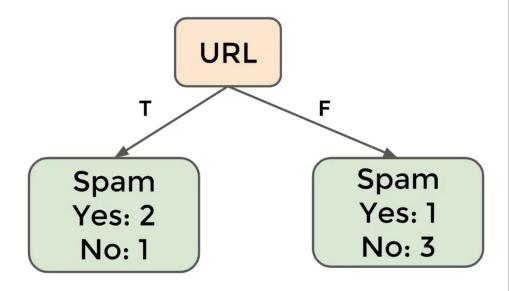
$$G(Q) = \sum_{c \in C} p_c (1-p_c)$$





- Пусть "спам" и "не-спам" будут классами С
- Левый листовой узел:

$$(\frac{2}{3})(1-\frac{2}{3}) + (\frac{1}{3})(1-\frac{1}{3})$$



$$G(Q) = \sum_{c \in C} p_c (1-p_c)$$

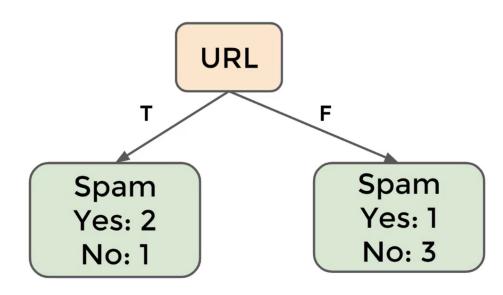




- Пусть "спам" и "не-спам" будут классами С
- Левый листовой узел:

$$(\frac{2}{3})(1-\frac{2}{3}) + (\frac{1}{3})(1-\frac{1}{3})$$

Gini=0.44



$$G(Q) = \sum_{c \in C} p_c (1-p_c)$$



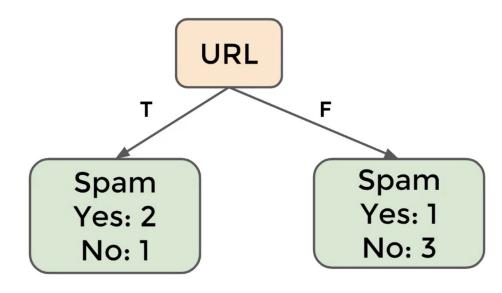


- Пусть "спам" и "не-спам" будут классами С
- Левый листовой узел:

$$(2/3)(1-2/3) + (1/3)(1-1/3)$$

• Правый листовой узел:

$$(1/4)(1-1/4) + (3/4)(1-3/4)$$

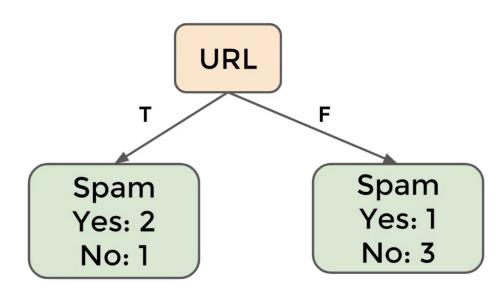


$$G(Q) = \sum_{c \in C} p_c (1-p_c)$$





- Итак, для признака URL метрика gini impurity:
- Левый лист: Gini=0.44
- Правый лист: Gini=0.375

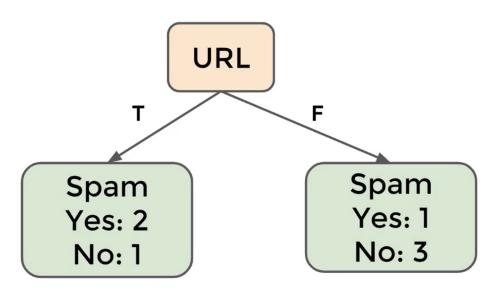


$$G(Q) = \sum_{c \in C} p_c (1-p_c)$$





- Всего писем: (2+1) + (1+3) = 7
- Левый лист: Gini=0.44
- Правый лист: Gini=0.375

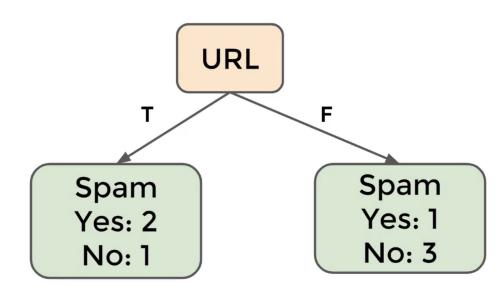


$$G(Q) = \sum_{c \in C} p_c (1-p_c)$$





- Всего писем: (2+1) + (1+3) = 7
- Левый лист: Gini=0.44
- Правый лист: Gini=0.375
- Писем слева: 3
- Писем справа: 4

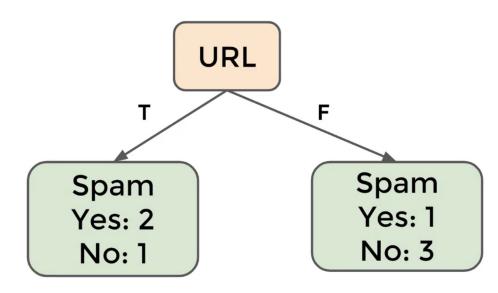


$$G(Q) = \sum_{c \in C} p_c (1-p_c)$$





- Всего писем: (2+1) + (1+3) = 7
- Левый лист: Gini=0.44
- Правый лист: Gini=0.375
- Писем слева: 3
- Писем справа: 4
- \bullet (3/7)*0.44+(4/7)*0.375
- Gini Impurity = 0.403

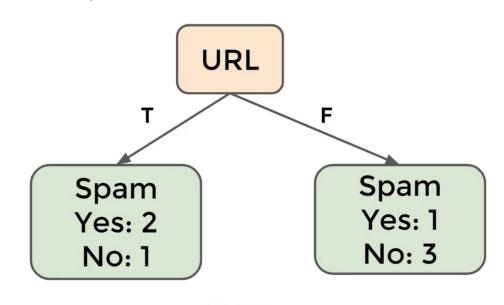


$$G(Q) = \sum_{c \in C} p_c (1-p_c)$$





• Gini Impurity для признака URL равна 0.403



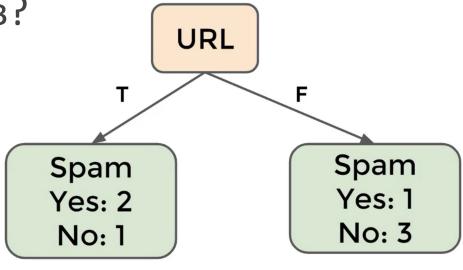
$$G(Q) = \sum_{c \in C} p_c (1-p_c)$$





Gini Impurity для признака URL равна 0.403

• А если несколько признаков?



$$G(Q) = \sum_{c \in C} p_c (1-p_c)$$





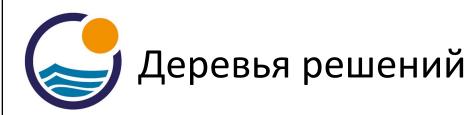
- Нам нужно разобрать ещё несколько вопросов:
 - Несколько признаков
 - Непрерывные признаки
 - Мульти-категориальные признаки
- Мы можем применить gini impurity для каждого из этих случаев, чтобы найти наилучший корневой узел и наилучшие параметры разбиения для листьев.





Gini impurity в деревьях – часть 2





- Ранее мы узнали, как вычислять gini impurity для бинарного категориального признака (состоящего только из двух категорий).
- Нам нужно разобрать ещё несколько вопросов:
 - Непрерывные признаки
 - О Мульти-категориальные признаки (N>2)
 - о Как выбрать корневой узел

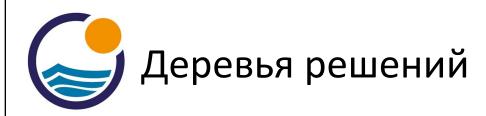




• Представим себе непрерывный признак:

X - Words in Email	Y-Spam	
10	Yes	
40	No	
20	Yes	
50	No	
30	No	

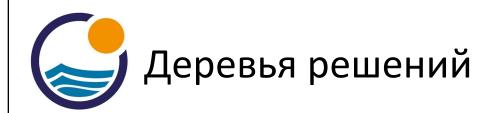




• Вычислим метрику gini impurity

X - Words in Email	Y-Spam	
10	Yes	
40	No	
20	Yes	
50	No	
30	No	





• Сначала отсортируем данные

X - Words in Email	Y-Spam	
10	Yes	
40	No	
20	Yes	
50	No	
30	No	





• Вычислим возможные значения для разделения

X - Words in Email	Y-Spam	
10	Yes	
20	Yes	
30	No	
40	No	
50	No	

Words ≤ N





• Возьмём средние значения "между" строками

X - Words in Email	Y-Spam	
15 10	Yes	
25 20	Yes	
30	No	
35 40	No	
50	No	

Words ≤ N





Вычислим gini impurity для каждого разбиения

X - Words in Email		Y-Spam	
15		Yes	
		Yes	
30		No	
	40	No	
50 No		No	

Words ≤ 15





Вычислим gini impurity для каждого разбиения

X - Words in Email		Y-Spam	
10 20		Yes	
		Yes	
30		No	
	40	No	
	50	No	

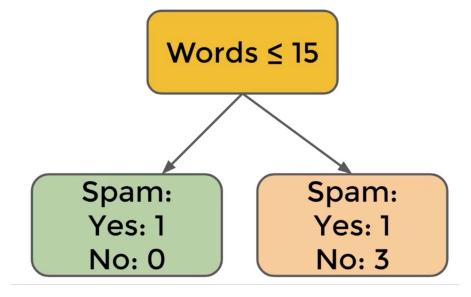
Words ≤ 15





Вычислим gini impurity для каждого разбиения

X - Words in Email		Y-Spam	
10 20		Yes	
		Yes	
30		No	
	40	No	
50		No	



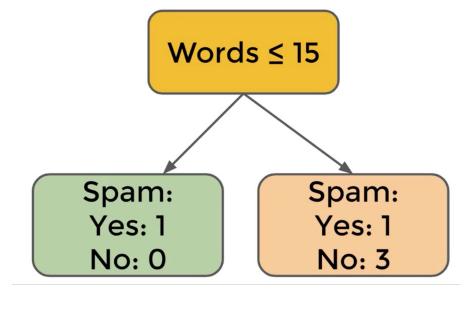
$$G(Q) = \sum_{c \in C} p_c (1-p_c)$$





Вычислим gini impurity для каждого разбиения

X - Words in Email		Y-Spam	
15	10	Yes	
20		Yes	
30		No	
	40	No	
	50	No	



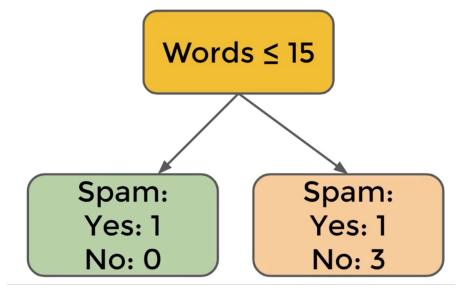
$$G(Q) = {}^{1/5}(0+0) + {}^{4/5}((1/4)(1-1/4)+(3/4)(1-3/4))$$





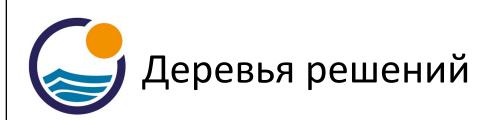
Вычислим gini impurity для каждого разбиения

X - Words in Email		Y-Spam	
10 20		Yes	
		Yes	
30		No	
	40	No	
	50	No	



$$G(Q)= {}^{1/5}(0+0)+{}^{4/5}((1/4)(1-1/4)+(3/4)(1-3/4))$$

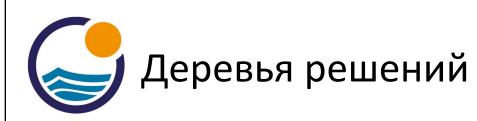




• Вычислим gini impurity для каждого разбиения

X - Wo	rds in Email	Y-Spam	
15	10	Yes	→ Gini=0.3
	20	Yes	U IIII-0.5
	30	No	
	40	No	
	50	No	

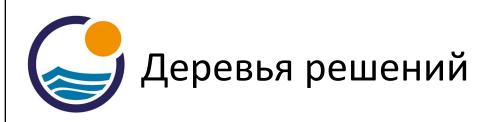




• Вычислим gini impurity для каждого разбиения

X - Words in Email		Y-Spam	
15	10	Yes	→ Gini=0.3
	20	Yes	→ Gini=0.3
25	30	No	
35	40	No	Gini=0.26
45	50	No	Gini=0.4





• Выбираем наименьший gini impurity

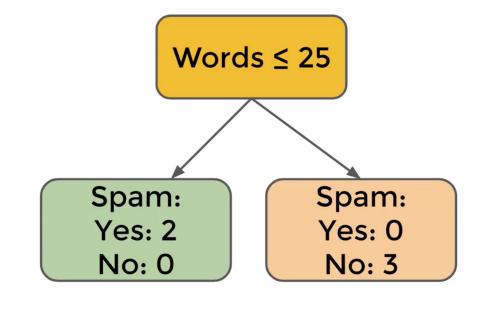
X - Words in Email	Y-Spam	
15 10	Yes	Gini=0.3
20	Yes	Gini=0.3
30	No	
40	No	→ Gini=0.26
50	No	Gini=0.4





• Это значение будет разделителем для узла

X - Words in Email		Y-Spam	
	10	Yes	
25	20	Yes	
	30	No	
	40	No	
	50	No	



$$G(Q) = \mathbf{0}$$





- Мы вычислили gini impurity для следующих случаев:
 - Бинарные признаки
 - Непрерывные признаки
- Далее давайте вычислим эту метрику для мульти-категориальных признаков

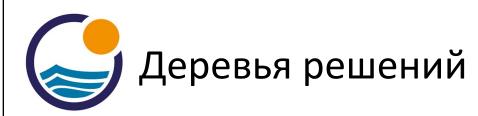




• Мульти-категориальный признак:

X - Sender	Y-Spam	
Abe	Yes	
Bob Yes		
Claire	No	
Abe	No	
Bob	No	





Вычислим gini impurity для всех комбинаций

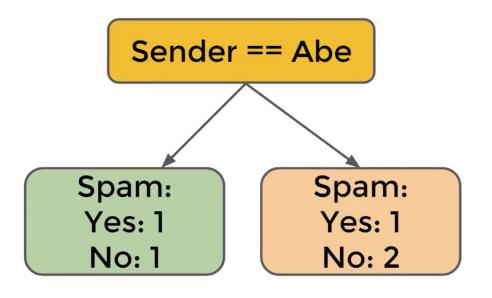
X - Sender	Y-Spam	
Abe	Yes	
Bob	Yes	
Claire	No	
Abe	No	
Bob	No	





• Вычислим gini impurity для всех комбинаций

X - Sender	Y-Spam	
Abe	Yes	
Bob	Yes	
Claire	No	
Abe	No	
Bob	No	

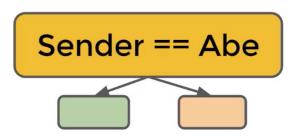






• Вычислим gini impurity для всех комбинаций

X - Sender	Y-Spam	
Abe	Yes	
Bob	Yes	
Claire	No	
Abe	No	
Bob	No	







• Вычислим gini impurity для всех комбинаций

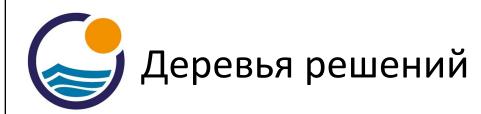
X - Sender	Y-Spam	Sender == Abe	
Abe	Yes		
Bob	Yes	Sender == Bob	
Claire	No	Serider Bob	
Abe	No		
Bob	No	Sender == Claire	
>1.4.1. <i>C</i> = >.4=4			





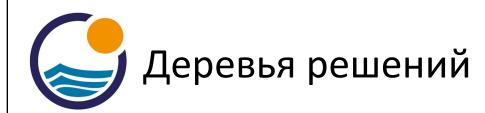
Вычислим gini impurity для всех комбинаций

X - Sender	Y-Spam	Sender == Abe	Sender == Abe или Bob
Abe	Yes		
Bob	Yes	Sender == Bob	Sender ==
Claire	No		Claire или Bob
Abe	No		
Bob	No	Sender == Claire	Sender == Abe или Claire
ΡΙΔΝ 🤼 ΡΔΤΔ			



- Теперь мы умеем вычислить gini impurity для различных типов признаков.
- Но как выбрать признак для корневого узла, когда признаков несколько?
- Нужно вычислить gini impurity для каждого признака, выбрать тот где метрика наименьшая, и взять его первым.



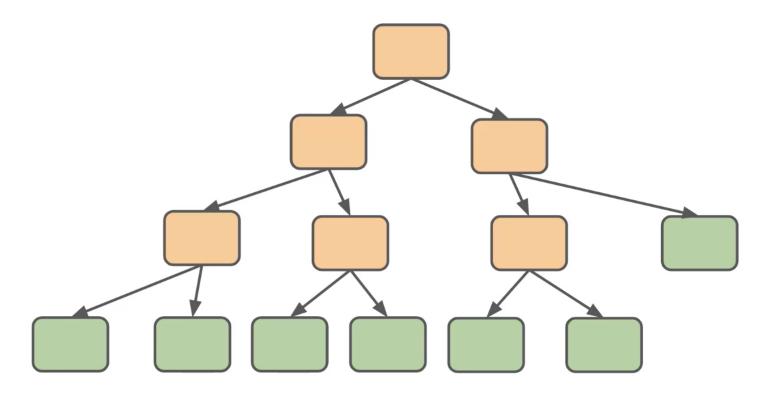


- Мы хотим минимизировать gini impurity на листьях.
- Оставляем листья, когда impurity меньше некоторого порогового значения.
- Выполняем разбиение только в тех случаях, когда impurity больше некоторого порогового значения.

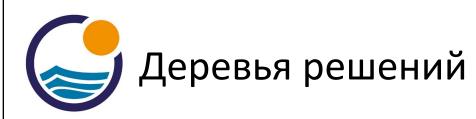




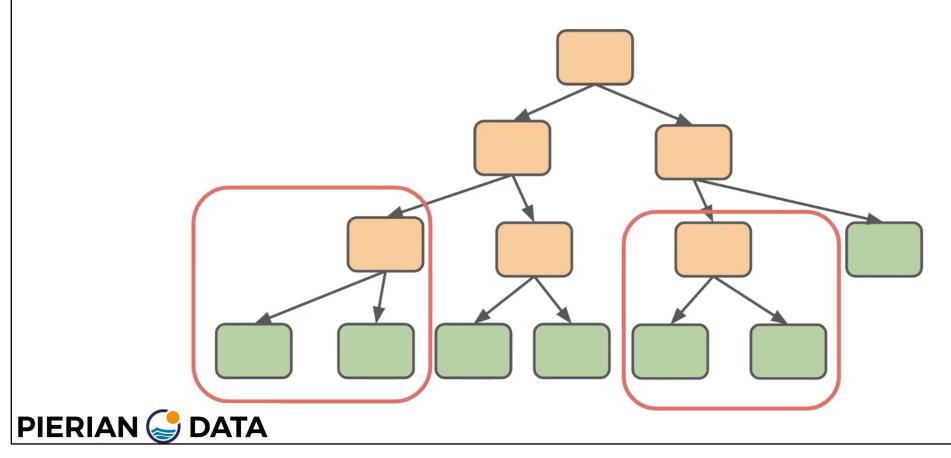
• Слишком большое дерево ("overfitted"):

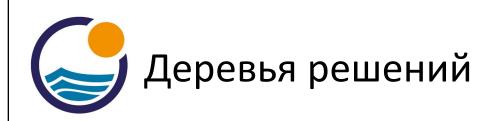




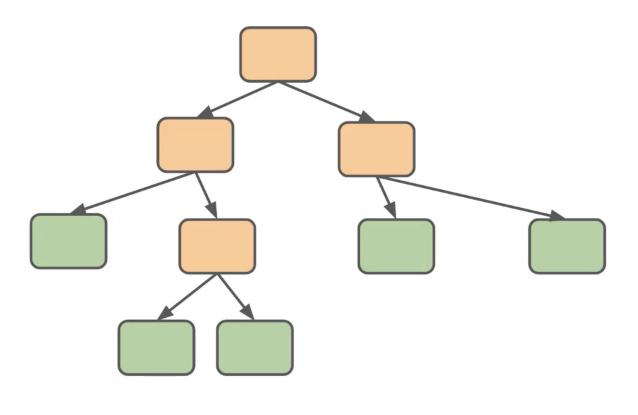


• Добавляем порог уменьшения для gini impurity:





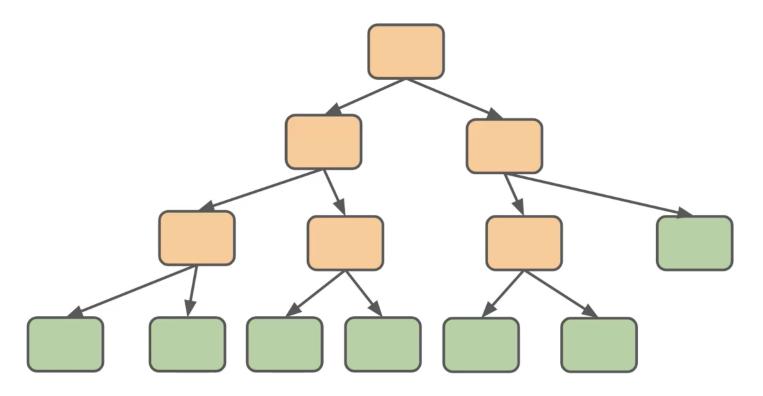
• Добавляем порог уменьшения для gini impurity:







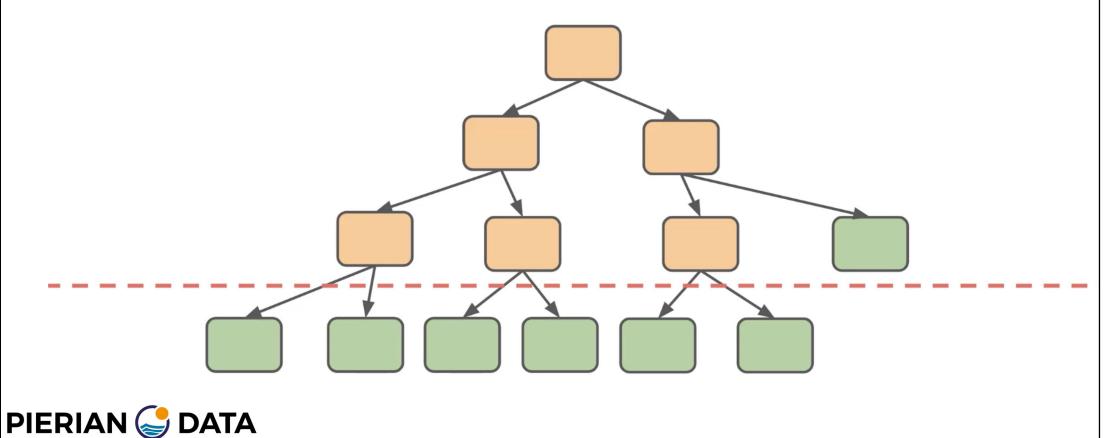
Ещё мы можем установить максимальную глубину

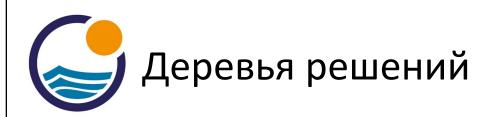




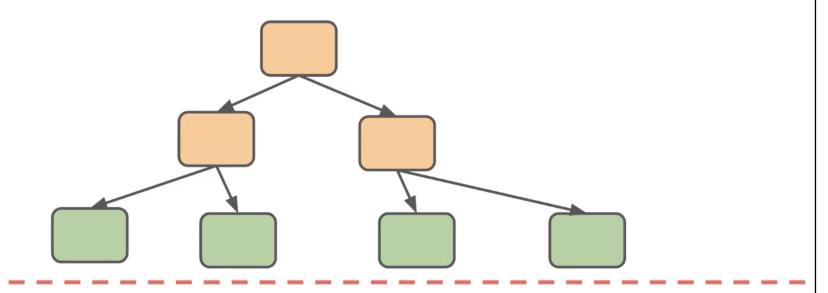


• Ещё мы можем установить максимальную глубину





Ещё мы можем установить максимальную глубину







 Посмотрим различные гиперпараметры во время написания кода!





Пишем код – часть 1 - данные





Пишем код – часть 2 - модель

