



## ТЕХНОЛОГИЧЕСКАЯ КАРТА ЗАНЯТИЯ

**Тема занятия:** Производная, градиент и градиентная оптимизация.

**Аннотация к занятию:** обучающиеся познакомятся с алгоритмом градиентной оптимизации для поиска минимума функции и применением алгоритма для функций многих переменных. Обсудят, как задача минимизации функции связана с машинным обучением.

**Цель занятия:** сформировать у обучающихся представление о задаче поиска минимума функции и познакомить с алгоритмом градиентной оптимизации.

### **Задачи занятия:**

- научить вычислять производную функции многих переменных;
- сформировать представление о частных производных и градиенте;
- обсудить задачу поиска минимума функции;
- познакомить с алгоритмом градиентной оптимизации.

## Ход занятия

Этап занятия	Время	Деятельность педагога	Комментарии, рекомендации для педагогов
<b>Организационный этап</b>	5 мин.	Друзья, здравствуйте! На предыдущем уроке мы познакомились с понятием производной функции, научились её вычислять. На этом уроке мы продолжим работать с производной	Приветствие. Создание в классе атмосферы психологического комфорта
<b>Постановка цели и задач занятия. Мотивация учебной деятельности обучающихся</b>	7 мин.	<p><b>Вопрос для обсуждения</b> Ранее мы рассматривали производные функций одной переменной. Как найти производную функции многих переменных?</p> <p><b>Возможные ответы обучающихся:</b></p> <ul style="list-style-type: none"> <li>• найти по каждой переменной;</li> <li>• её невозможно найти.</li> </ul> <p>Сегодня научимся вычислять производные многих функций. Узнаем правила вычисления производных более сложных, составных функций</p>	Способствовать обсуждению мотивационных вопросов
<b>Изучение нового материала</b>	50 мин.	Напомню, что такое функция многих переменных. На экране вы видите функцию $f$ от двух переменных ( $x_1$ и $x_2$ ) и её график. На графике обозначена точка,	Для справки:

		<p>соответствующая значениям <math>x_1</math> равно минус двадцать и <math>x_2</math> равно минус двадцать. Варьируя значения <math>x_1</math> и <math>x_2</math>, мы будем получать различные значения функции <math>f</math>.</p> <p>Бывают функции от трёх, четырёх и сколько угодно переменных. На самом деле, мы уже сталкивались с функцией многих переменных, когда говорили о линейной регрессии.</p> <p><b>Вопрос для обсуждения</b> Давайте вспомним, что такое линейная регрессия.</p> <p>Это модель для решения задачи регрессии. Она ставит в соответствие каждому признаку датасета свой коэффициент. И затем для входящего элемента вычисляет ответ по формуле, которую вы видите на экране. В этой формуле <math>x_i</math> — <math>i</math>-ый признак элемента.</p> <p>Например, для первого элемента датасета ответ линейной регрессии <math>y_1</math> с шапочкой выглядел бы вот так.</p> <p>Смотрите: формулу ответа модели для первого элемента датасета можно представить в виде функции <math>f</math> от шести переменных: <math>k_0, k_1, k_2, k_3, k_4</math> и <math>k_5</math>. Меняя значения <math>k_i</math>-ых, мы получаем разные модели линейной регрессии, которые выдают разные ответы на первый элемент.</p> <p>Нарисовать графики функций, у которых больше, чем две переменные, мы не можем. Понять, возрастают или убывают они в разных точках и с какой скоростью, можно только с помощью производных.</p>	<p>Сайт: <a href="https://habr.com/ru/post/413853/">https://habr.com/ru/post/413853/</a> Перед уроком рекомендуется ознакомиться с материалами, представленными на сайте.</p>
--	--	---	---

Так давайте научимся вычислять производные функций многих переменных.

Вернёмся к нашей функции  $f$  от  $x_1$  и  $x_2$ . Давайте представим, что  $x_2$  — это не переменная, а константа. Как будто у функции  $f$  только одна переменная —  $x_1$ . У такой функции мы можем взять производную. Она будет равна шесть  $x_1$  плюс два  $x_2$ .

Это называется частной производной функции  $f$  по переменной  $x_1$ . Или проще: производной  $f$  по  $x_1$ . Записывается это как  $f'$  с нижним индексом  $x_1$  или как  $df$  по  $dx_1$ . Обратите внимание, что в записи формулы буква  $d$  должна быть не совсем обычная, а немного закруглённая.

**Важно:** несмотря на то, что мы вычисляли эту производную, представляя, что  $x_2$  — это не переменная, производная  $f'$  по  $x_1$  — это всё равно функция от двух переменных,  $x_1$  и  $x_2$ . То есть сначала мы берём функцию  $f$ , представляем, что  $x_2$  — это не переменная, а константа, вычисляем так производную  $f$  по  $x_1$ , и дальше снова начинаем считать, что  $x_2$  — это переменная.

Обобщим понятие частной производной функции многих переменных. Пусть у нас есть функция  $f$  от  $n$  переменных. Чтобы вычислить частную производную функции  $f$  по переменной  $x_1$ , нужно:

1. Представить, что все остальные переменные  $x_2$ ,  $x_3$  и так далее — константы.
2. Вычислить производную  $f$  по  $x_1$ .

Для справки:  
[https://ru.wikipedia.org/  
wiki/Частная\\_производная](https://ru.wikipedia.org/wiki/Частная_производная)

		<p>3. Снова начать считать <math>x_2</math>, <math>x_3</math> и так далее переменными.</p> <p>Всё. Так мы получаем функцию <math>f'</math> с индексом <math>x_1</math> от <math>n</math> переменных. Это и будет частная производная <math>f'</math> по <math>x_1</math>.</p> <p>Аналогично вычисляются остальные частные производные <math>f'</math> по <math>x_2</math>, <math>f'</math> по <math>x_3</math> и так далее. Всего у функции с <math>n</math> переменными будет <math>n</math> частных производных.</p> <p>Вернёмся к нашей функции <math>f</math> от двух переменных. Мы уже вычислили для нее <math>f'</math> по <math>x_1</math>.</p> <p>Так как частные производные — это тоже функции, мы можем вычислять значения производных в различных точках.</p> <p>Посмотрим на знакомую нам точку <math>x_1</math> равно <math>x_2</math> равно 20. Обозначим её буквой <math>A</math>. На слайде вы видите её положение на графике функции. Значение частной производной <math>f'</math> по <math>x_1</math> в этой точке равно минус сто шестьдесят. Давайте обсудим, что это значит.</p> <p>Смысл частных производных на самом деле такой же, как и у обычной производной функции одной переменной — о ней мы говорили на прошлых занятиях.</p> <p>Значение частной производной <math>f'</math> по <math>x_1</math> отражает то, что будет происходить со значением функции <math>f</math>, если мы из точки <math>A</math> сдвинемся по <math>x_1</math> на бесконечно малое значение <math>\delta x</math> вправо. То есть перейдём в точку <math>(-20 + \delta x, -20)</math>.</p>	
--	--	--	--

		<p>Ещё проще: смотрите, наше значение производной <math>f'</math> по <math>x_1</math> в точке <math>A</math> равно минус сто шестьдесят. Это значение меньше нуля. Это значит, что функция <math>f</math> в точке <math>A</math> убывает по <math>x_1</math>. Это видно по графику: если мы из точки <math>A</math> при неизменном значении <math>x_2</math> будем двигаться вправо по оси <math>x_1</math>, то значение нашей функции будет уменьшаться.</p> <p>Посмотрим на производную <math>f'</math> по <math>x_2</math> в точке <math>A</math>. <math>f'</math> по <math>x_2</math> равно два <math>x_1</math>. Подставив значение <math>x_1</math> равно минус двадцать, получим, что <math>f'</math> по <math>x_2</math> в точке <math>A</math> равно минус 40. Тоже отрицательное значение. Это значит, что в точке <math>A</math> функция <math>f</math> по переменной <math>x_2</math> убывает. Это заметно и на графике: двигаясь параллельно оси <math>x_2</math> от точки <math>A</math> в сторону увеличения значения <math>x_2</math>, значение функции уменьшается.</p> <p>Замечу, что <math>f'</math> по <math>x_2</math> в точке <math>A</math> меньше по модулю, чем <math>f'</math> по <math>x_1</math> в точке <math>A</math>. Значит, функция <math>f</math> в точке <math>A</math> по <math>x_1</math> убывает быстрее, чем по <math>x_2</math>. Это тоже видно на графике: склон функции при движении вдоль оси <math>x_1</math> гораздо круче, чем при движении вдоль оси <math>x_2</math>.</p> <p>Таким образом, мы увидели, что смысл частных производных такой же, как и смысл производной функции одной переменной. Частные производные отражают характер поведения функции в точке по отдельным переменным.</p> <p>Иначе говоря, частная производная <math>f'</math> по <math>x</math> в точке показывает нам, в какую сторону нам нужно сдвинуться из точки по переменной <math>x</math>, чтобы значение функции уменьшилось.</p>	<p>Для справки:</p>
--	--	---	---------------------

		<p>Теперь мы можем ввести понятие градиента. Градиент функции в точке <math>x</math> — это просто вектор значений частных производных функции в точке <math>x</math>. Сколько у функции переменных, столько и будет элементов в векторе градиента. Например, для нашей функции <math>f</math> её градиент в точке <math>A</math> — это вектор минус сто шестьдесят, минус сорок.</p> <p>Обозначается градиент значком, который называется <math>\nabla f</math>. Он выглядит как перевернутый треугольник.</p> <p>У градиента есть смысл. Он содержит всю информацию о том, возрастает или убывает функция <math>f</math> в точке по всем переменным, и насколько быстро.</p> <p>Получается, вычислив градиент функции в точке, мы можем понять, в какую сторону нам нужно сдвинуться из этой точки по каждой из переменных. Так, чтобы значение функции уменьшилось.</p> <p>Например, для нашей функции <math>f</math> и точки <math>A</math> нам нужно увеличить значения обеих переменных <math>x_1</math> и <math>x_2</math>, чтобы значение функции уменьшилось. На графике движение по обеим переменным <math>x_1</math> и <math>x_2</math> в сторону увеличения показано красной стрелкой. Видно, что, увеличивая <math>x_1</math> и <math>x_2</math>, мы уменьшаем значение <math>f</math>.</p> <p>Градиентная оптимизация — это инструмент, который позволяет решить задачу минимизации функции. Подробнее остановимся на том, что же такое задача минимизации функции.</p>	<p>Сайт: <a href="https://ru.abcdef.wiki/wiki/Gradient">https://ru.abcdef.wiki/wiki/Gradient</a></p>
--	--	--	--

		<p>Пусть у нас есть функция <math>f(x)</math>. Задача минимизации этой функции состоит в нахождении её точки минимума.</p> <p>Мы помним, что минимумы функции бывают локальные и глобальные. Алгоритмы поиска минимума функций чаще всего ищут точки локальных минимумов функций.</p> <p>Для функции многих переменных найти точку минимума значит найти такой набор значений параметров, что значение функции для них будет минимально. На слайде вы видите функцию от двух переменных и её график. Её точка минимума — при значениях <math>x_1</math> и <math>x_2</math>, равных нулю. При любых других значениях <math>x_1</math> и <math>x_2</math> значение функции <math>f</math> будет больше нуля.</p> <p>Задачу минимизации функции называют и по-другому. Например, задачей поиска минимума функции или задачей оптимизации функции.</p> <p>Почему поиск точек минимума функций важен? Дело в том, что обучение многих моделей машинного обучения как раз и состоит в том, чтобы искать точки минимума метрик качества. Так, к примеру, обучаются нейросети. Более подробно об этом вы узнаете в следующих модулях, когда будете изучать логистическую регрессию и нейросети.</p> <p><b>Вопрос для обсуждения</b> Итак, как же решить такую задачу? Как по данной функции <math>f</math> найти её точку минимума?</p> <p><b>Возможный ответ обучающихся</b></p>	
--	--	---	--



		<p>Найти производную.</p> <p>Здесь нам и пригодятся производные. Рассмотрим первый способ нахождения точек минимума — аналитический.</p> <p>Вспомним одно из свойств производной. Производная в точках экстремума функции равна нулю. Поэтому идея поиска точек минимума функции такая:</p> <ul style="list-style-type: none"><li>• берём функцию <math>f</math>, вычисляем её производную <math>f'</math>;</li><li>• затем решаем уравнение <math>f'(x) = 0</math>;</li><li>• находим значения <math>x</math> — корни этого уравнения. Так как при этих значениях <math>x</math> производная <math>f'(x)</math> будет равна нулю, то корни уравнения являются точками экстремумов функции <math>f</math> — то есть, точками минимума и максимума;</li><li>• после этого остаётся последний шаг — понять, какие из найденных точек — точки минимума.</li></ul> <p>Эта идея хороша, но, к сожалению, работает далеко не для всех функций. Убедимся в этом на примере знакомой функции <math>f(x)</math>, график которой вы видите на экране. Её производная <math>f'</math> равна четыре <math>x</math> в третьей плюс пятнадцать <math>x</math> в квадрате минус десять.</p> <p>Получается, чтобы найти точки минимума, нужно решить уравнение <math>f'(x)</math> равно нулю. <math>f'</math> — это многочлен третьей степени, и решить такое уравнение непросто. А что если функция <math>f</math> была бы десятой или двадцатой степени? Не говоря уж о случае, когда <math>f</math> — функция многих переменных.</p>	
--	--	---	--

		<p>Получается, такой способ нахождения точек минимума функции работает не всегда. Для многих функций <math>f</math> второй пункт этого плана сложно выполним, если вообще выполним.</p> <p>Что же делать? К сожалению, никакого другого способа, который бы точно находил все точки минимума у любой функции, не существует.</p> <p>Мы рассмотрим другой способ, который позволяет находить некоторые точки локального минимума функций. Он основан на производной и называется алгоритмом градиентной оптимизации.</p> <p>Возьмём произвольную точку функции <math>f</math>. Например, точку <math>x</math> равно минус пять. Значение функции в этой точке равно пятидесяти.</p> <p>Вычислим значение производной в этой точке: оно равно минус ста тридцати пяти. Вспомним, о чем говорит нам знак производной: о том, что функция <math>f</math> в точке минус пять убывает.</p> <p>Что это значит: раз функция в точке минус пять убывает, значит, где-то правее этой точки находится локальный минимум. Если мы сдвинемся из точки минус пять вправо, то, скорее всего, подвинемся чуть ближе к точке локального минимума.</p> <p>Давайте сдвинемся вправо на шаг <math>\delta</math> <math>x</math> равно 1 и попадём в точку <math>x</math> равно минус четыре. Значение функции действительно стало меньше — минус двадцать четыре, — а мы стали ближе к локальному минимуму.</p>	
--	--	--	--

		<p>Значение производной в этой точке равно минус двадцати шести: всё ещё меньше нуля. Значит, чтобы найти точку локального минимума, нужно всё ещё двигаться вправо и увеличивать <math>x</math>.</p> <p>Хорошо. Двигаемся дальше: ещё на единицу. Попадаем в точку <math>x</math> равно минус трём. Считаем значение функции и производную: они равны минус двадцати четырём и семнадцати. Смотрите: производная стала больше нуля, а значение функции не изменилось. Это значит, что мы перешагнули точку минимума и попали в место, где функция начала возрастать.</p> <p>Получается, теперь, чтобы попасть в точку минимума, нужно идти левее: уменьшать <math>x</math>. Но если мы уменьшим <math>x</math> на значение <math>\Delta x</math>, равное единице, мы снова перешагнём точку минимума и попадём обратно в точку минус четыре, где уже были. Значит, шаг <math>\Delta x</math> нужно уменьшить.</p> <p>Возьмём <math>\Delta x</math> равно ноль целых пять десятых и сдвинемся влево: получим точку минус три с половиной. Значение функции станет равно около минус двадцати девяти: оно уменьшилось по сравнению с предыдущим шагом, то есть, мы стали ближе к точке минимума. И эта точка, кстати, уже очень близка к реальной точке минимума минус три целых пятьдесят пять сотых.</p> <p>Производная осталась положительной. Значит, мы всё ещё находимся правее точки минимума, и нам нужно снова двигаться влево, чтобы её найти.</p>	
--	--	---	--

		<p>Этот процесс можно продолжать и дальше. Таким образом мы будем все ближе и ближе подходить к точке локального минимума.</p> <p>Обобщим алгоритм, который мы только что придумали. Итак, у нас есть функция <math>f</math>. Чтобы найти её точку минимума, нужно сделать следующее:</p> <ol style="list-style-type: none"><li>1. Выбираем случайную точку функции <math>x</math> и фиксируем значение шага <math>\delta x</math>.</li><li>2. Понимаем, в какую сторону нужно идти из точки <math>x</math>, чтобы приблизиться к точке минимума. Для этого вычисляем производную <math>f'(x)</math> и смотрим на знак. Если производная больше нуля, двигаемся влево, уменьшаем значение <math>x</math>. Если производная меньше нуля, двигаемся вправо, к точке <math>x</math> плюс <math>\delta x</math>.</li><li>3. Снова вычисляем производную функции и смотрим, изменился ли её знак. Если не изменился, всё хорошо: мы идём в верном направлении. Если изменился, значит, мы перешагнули точку минимума: сделали слишком большой шаг <math>\delta x</math>. Тогда уменьшим <math>\delta x</math> и продолжим идти к точке минимума: снова вычислим производную, поймём, в какую сторону двигаться, и так далее.</li></ol> <p>Понятно, что, двигаясь таким образом, мы будем всё ближе и ближе подходить к точке локального минимума. Действительно: как только мы перешагиваем через точку минимума, мы уменьшаем шаг, начинаем двигаться аккуратнее.</p>	
--	--	---	--

		<p>Однако ровно в точку минимума мы вряд ли когда-нибудь попадём. Поэтому этот алгоритм ищет точку минимума приближённо: ищет точку, которая будет рядом с точкой минимума. Как, например, на прошлом слайде: мы нашли точку минус три с половиной, которая была довольно близка к минимуму.</p> <p>Понятно, что этот алгоритм находит только одну точку локального минимума: ту, которая ближе всего к начальной точке. Можно попытаться найти несколько точек минимума, проделав этот алгоритм из нескольких разных начальных точек, надеясь, что они приведут к разным точкам минимума. Однако в случае, когда у вас нет доступа к графику функции и вы не знаете, сколько у функции вообще точек минимума, вы не сможете гарантированно найти их все.</p> <p>Несмотря на эти недостатки, алгоритм хорош тем, что может работать с любой функцией, у которой можно взять производную. Всё, что мы делаем, — вычисляем производную функции и смотрим на её знак. Никаких уравнений решать не нужно.</p> <p>Стоит заметить, что с помощью этого алгоритма можно искать и точки максимумов функций. Нужно просто делать шаги на <math>\Delta x</math> в противоположном направлении: не в сторону убывания функции, а в сторону возрастания.</p> <p>Алгоритм такого вида называется градиентной оптимизацией функции. Градиентная она потому, что мы в процессе используем производную — градиент.</p>	
--	--	--	--

### Вопрос для обсуждения

К этому моменту у вас, наверное, возник вопрос: а как понять, когда в этом алгоритме остановиться? Ведь если мы никогда не попадем в точку минимума, мы можем вечно ходить вокруг нее, всё приближаясь к ней, но никогда не достигая.

**Ответ прост:** чаще всего в таких алгоритмах останавливаются тогда, когда значение  $\delta x$  становится маленьким. Например, одна тысячная.

**Смотрите:** если  $\delta x$  маленькое, это значит, что мы очень близко подобралась к точке минимума: настолько, что при чуть большем значении  $\delta x$  мы уже перешагнем через неё. На этом можно остановиться: сказать, что мы нашли точку  $x$ , которая находится очень близко к точке минимума, и нам этого достаточно.

Тогда итоговый алгоритм градиентной оптимизации будет выглядеть так: те же пункты, что на предыдущих слайдах, и к пункту пять добавим оговорку. Мы двигаем точку  $x$  до тех пор, пока  $\delta x$  не станет меньше некоторого числа  $\epsilon$ .

Для начала немного поправим алгоритм градиентной оптимизации. Алгоритм представлен на слайде. Он прекрасно работает и ищет точки рядом с точками минимума, но есть один нюанс.

В этом алгоритме мы до начала движения из точки  $x$  фиксировали начальное значение  $\delta x$ . Например, для этой функции в прошлом мы брали  $\delta x$ , равное

единице. Затем, когда перешагивали через точку минимума, начинали  $\delta x$  уменьшать.

### Вопрос для обсуждения

Какое начальное значение  $\delta x$  стоит брать?

Можно всегда брать  $\delta x$ , равное единице. Но тогда мы с таким  $\delta x$  будем очень-очень долго идти к точке минимума. Например, рассмотрим функцию  $f$  равно  $x$  в четвёртой степени. Её точка минимума — ноль. И пусть мы хотим с помощью нашего алгоритма найти эту точку минимума, двигаясь из точки  $x$  равно двадцать. И пусть  $\delta x$  у нас равно единице. Видно, что прежде чем мы подойдём достаточно близко к точке минимума и уменьшим значение  $\delta x$ , пройдёт как минимум двадцать шагов. Это довольно долго. А что будет, если мы стартуем из точки  $x$  равно двести или  $x$  равно две тысячи?

Кажется,  $\delta x$ , равное единице, было не лучшим выбором для шага в этом случае.

Другой пример: та же функция и стартовая точка  $x$  равно ноль целых две десятых. Казалось бы, эта точка близка к точке минимума: достаточно несколько раз аккуратно сдвинуться на небольшое  $\delta x$  влево, и мы окажемся практически в минимуме. Но если мы будем идти по алгоритму, то сначала сдвинемся в точку минус ноль целых восемь десятых, затем уменьшим  $\delta x$  и попадём в точку минус ноль целых три десятых, затем в точку минус ноль целых две десятых и так далее. Будем блуждать

около точки минимума, постоянно перепрыгивая через неё.

Иногда придуманный нами алгоритм работает не совсем оптимально. Хочется придумать, как выбирать  $\Delta x$  так, чтобы быстрее приходиться в точку минимума: чтобы на каждом шаге сдвигаться как можно ближе к минимуму.

Давайте подумаем, как можно улучшить выбор  $\Delta x$ . Вспомним о ещё одном свойстве производной: модуль производной функции в точке показывает скорость роста функции в этой точке. А теперь посмотрим на график функции  $f$  равно  $x$  в четвёртой. Заметим вот что: чем дальше точка находится от точки минимума, тем быстрее растёт или убывает функция в этой точке.

Получается, для точек, близких к минимуму, модуль производной будет мал. А для точек, далёких от минимума, модуль чаще всего будет большим. Это правило работает не всегда, но для большинства точек это верно. Можно использовать это свойство, чтобы понимать, насколько большой шаг в сторону минимума мы можем совершить, находясь в той или иной точке.

Правило такое: чем больше модуль производной, тем больший шаг по направлению минимума мы можем сделать. Однако модуль производной не говорит нам о том, какой именно по величине шаг нужно выбрать. Шаг, равный модулю производной, не подойдет: он будет слишком большим. Смотрите, в точке  $x$  равно 5 производная функции  $f$  равна пятистам. Если мы



		<p>двинемся на шаг <math>\Delta x</math>, равный пятистам, то улетим далеко по другую сторону от минимума.</p> <p>Решение такое: сдвигаться от точки на величину, равную модулю производной, умноженной на некое маленькое число <math>\alpha</math>. Обычно <math>\alpha</math> берут равным одной тысячной или ещё меньше. Тогда, если мы каждый раз будем сдвигаться из точки <math>x</math> в сторону минимума на шаг <math>\alpha</math> умножить на модуль производной, мы будем двигаться на больший шаг из далёких от минимума точек и на меньший шаг из близких к минимуму точек. По мере приближения к точке минимума модуль производной будет становиться меньше, и шаг <math>\Delta x</math> тоже будет уменьшаться.</p> <p>Например, при <math>\alpha</math> равном одной тысячной, из точки пять мы сдвинемся влево на шаг <math>\Delta x</math>, равный ноль целым пяти десятым, а из точки двадцать — на <math>\Delta x</math>, равный тридцати двум.</p> <p>Конечно, такой выбор <math>\Delta x</math> тоже не идеален. Иногда мы всё ещё будем перепрыгивать точки минимума и тратить большое количество шагов на то, чтобы подойти близко к минимуму. Например, здесь мы из точки двадцать, сдвинувшись на <math>\Delta x</math>, равное тридцати двум, перепрыгнули минимум.</p> <p>На практике получается так, что в среднем для всех функций такой выбор <math>\Delta x</math> ускоряет алгоритм градиентной оптимизации.</p> <p>На этом слайде вы видите процесс работы алгоритма для функции <math>2x</math> умножить на синус <math>x</math> из точки <math>x</math> равной</p>	
--	--	--	--

двум с половиной. Здесь альфа равно пяти тысячным. Альфа достаточно мало, чтобы мы не перескакивали через точку локального минимума, а аккуратно приближались к точке минимума с левой стороны. Чем меньше  $\alpha$ , тем ближе наша функция в итоге сможет подойти к точке минимума. Но при очень малых  $\alpha$  каждый шаг получается маленьким, алгоритм работает долго. Поэтому нужно выбирать  $\alpha$  исходя из того, что важнее: точность или скорость.

Видно, что в местах, где функция убывает быстрее, алгоритм ускоряется: делает большие шаги от точки к точке по направлению минимума. Причина в том, что в этих местах модуль производной функции выше и шаг  $\Delta x$  получается выше.

Кроме того, видно, что в начальной точке два с половиной, которая далека от минимума, модуль производной не такой большой, как в более близких точках, например, около  $x$ , равного четырём. Получается, что правило «чем дальше точка от минимума, тем быстрее убывает функция и тем производная больше по модулю» выполняется не всегда. Для многих точек функций это верно, поэтому наш алгоритм градиентной оптимизации в целом работает хорошо и подбирает оптимальные размеры шага  $\Delta x$ .

На этой гифке хорошо видно, что наш алгоритм в процессе работы на каждом шаге вычисляет производную функции и спускается точка за точкой к точке минимума. Поэтому этот алгоритм также называется градиентным спуском.

		<p>Подведём итог и сформулируем обновленный алгоритм градиентного спуска.</p> <p>Суть его осталась той же. Выбираем начальную точку <math>x</math>: мы будем двигаться от неё к точке минимума. Но теперь мы не выбираем значение <math>\Delta x</math>, а фиксируем значение <math>\alpha</math>.</p> <p>Далее на каждом шаге мы будем вычислять производную функции <math>f</math> в текущей точке <math>x</math>, смотреть на знак этой производной и сдвигать точку <math>x</math> в сторону минимума на величину, равную <math>\alpha</math> умножить на модуль производной.</p> <p>Повторяем эти пункты, пока шаг, на который мы сдвигаем нашу точку, не станет очень маленьким. То есть, пока мы не будем близки к минимуму.</p> <p>В алгоритме стало на один пункт меньше, чем было до этого. Нам больше не нужно изменять <math>\Delta x</math> вручную, когда мы перепрыгнули через точку минимума. <math>\Delta x</math> теперь зависит от модуля производной функции в каждой точке.</p> <p>Формулу движения <math>x</math> в пункте 2 можно записать проще, в одну строчку.</p> <p>Запишем это так: новое значение <math>x</math> равно текущее значение <math>x</math> минус <math>\alpha</math> умножить на значение производной в текущей точке. Значение производной без модуля.</p>	
--	--	---	--

Действительно, если  $f'(x)$  больше нуля, то значение  $x$  уменьшится на величину альфа умножить на модуль производной. То есть, мы сдвинем точку влево, где и находится точка минимума, так как знак производной больше нуля. А если  $f'(x)$  меньше нуля, то значение  $x$ , наоборот, увеличится на величину альфа умножить на модуль производной. Такая формула пересчёта координаты точки  $x$  эквивалентна тому, что было написано на предыдущем слайде. Вы можете подумать и убедиться в этом сами.

Мы получили алгоритм градиентной оптимизации, который используется в математике и машинном обучении.

Давайте узнаем, как применять такой алгоритм для функции нескольких переменных. Тут, на самом деле, всё просто. Алгоритм будет точно таким же. Вот как он будет выглядеть для функции двух переменных.

Фиксируем начальную точку  $x$  — координаты  $x_1$  и  $x_2$ . Они могут быть любыми. На каждом шаге вычисляем значения частных производных по  $x_1$  и  $x_2$  в текущей точке. Меняем значения координат  $x_1$  и  $x_2$  как новое значение равно старое минус альфа умножить на значение частной производной в старой точке. Повторяем так до тех пор, пока изменения обеих координат не станут очень малы.

Как мы помним, частная производная функции по переменной  $x_1$  показывает, убывает или возрастает функция по этой переменной. Получается, алгоритм

		<p>градиентной оптимизации для функции двух переменных на каждом шаге понимает для каждой координаты отдельно, в какую сторону и насколько её нужно сдвинуть, чтобы приблизиться к точке минимума функции.</p> <p>Визуализацию градиентного спуска для функции двух переменных вы видите на слайде. Здесь альфа равно пяти сотым. По мере приближения к точке минимума функция начинает убывать медленнее, шаг алгоритма становится всё меньше.</p> <p>Последний шаг: алгоритм градиентного спуска для функции многих переменных можно переписать в терминах градиента. Градиент — это вектор частных производных функции. Поэтому алгоритм можно записать в таком виде, как на слайде.</p> <p>Мы разобрали градиентную оптимизацию или градиентный спуск для функций одной и многих переменных. На практическом занятии мы постараемся применить её на языке Python.</p> <p>В заключение скажу пару слов о том, как градиентный спуск применяется для обучения моделей машинного обучения.</p> <p>Вспомним линейную регрессию и то, как она получает ответ на элемент датасета. Мы уже говорили о том, что функцию получения ответа регрессии можно рассматривать как функцию от <math>n+1</math> переменных, где <math>n</math> — количество признаков в данных. На слайде <math>n</math> равно пяти,</p>	
--	--	---	--

		<p>и линейная регрессия выражает функцию от шести переменных: <math>k_1</math>, <math>k_2</math>, <math>k_3</math> и так далее.</p> <p>Дальше. В модуле о линейной регрессии мы говорили о том, что задача модели — подобрать такие коэффициенты <math>k</math>, чтобы значение метрики качества MSE на обучающем датасете было минимально. Давайте запишем формулу метрики MSE для первого элемента датасета на экране.</p> <p>MSE от ответа модели <math>y_1</math> с шапочкой и правильного значения ответа <math>y_1</math> равно <math>y_1</math> минус <math>y_1</math> с шапочкой в квадрате. Подставим вместо <math>y_1</math> с шапочкой формулу для неё. Получим, что MSE — это функция от тех же шести переменных <math>k_1</math>, <math>k_2</math> и так далее. Задача линейной регрессии — минимизировать MSE, то есть найти такие <math>k_1</math>, <math>k_2</math>, <math>k_3</math> и так далее, чтобы значение MSE было минимальным. А это ровно та задача минимизации функции многих переменных, которая решается с помощью алгоритма градиентного спуска.</p> <p>Конечно, здесь есть нюансы. При обучении модели мы хотим, чтобы MSE было минимальным в среднем на всех элементах датасета, а не только на одном, как тут. Как этого добиться, мы узнаем в следующих модулях, когда будем разбирать градиентный спуск для нескольких моделей машинного обучения.</p>	
--	--	---	--

<b>Закрепление изученного материала</b>	15 мин.	<b>Вопросы для обсуждения:</b> <ul style="list-style-type: none"> <li>• Как вычислять производные функций многих переменных?</li> <li>• Что такое градиентная оптимизация?</li> <li>• Как задача минимизации функции связана с машинным обучением?</li> </ul>	Педагог организует беседу по вопросам
<b>Этап подведения итогов занятия (рефлексия)</b>	8 мин.	<b>Вопросы для обсуждения</b> <ul style="list-style-type: none"> <li>• Чему я научился?</li> <li>• С какими трудностями я столкнулся?</li> <li>• Каких знаний мне не хватает для более глубокого понимания изученного материала?</li> <li>• Достиг ли я поставленных целей и задач?</li> </ul>	Педагог способствует размышлению обучающихся над вопросами
<b>Информация о домашнем задании, инструктаж по его применению</b>	5 мин.	<p>В этом домашнем задании вам предстоит поработать с понятием производной и градиента, а также написать градиентный спуск и его вариации.</p> <p><b>Основной пункт задания</b></p> <p>В этом задании Вы должны реализовать функцию <code>grad_descent_v1</code> для нахождения минимума функции с помощью градиентного спуска.</p> <p>Вход функции.</p> <p>Функция <code>func</code>, которую нужно оптимизировать</p> <p>Ее производная <code>deriv</code>. В данном пункте производная вам дана, вычислять ее самостоятельно не нужно. (*)</p> <p>Начальная точка <code>start</code>.</p> <p>Выход функции -- точка локального минимума.</p>	

		<p>Для вашего удобства мы написали функцию для отрисовки траектории градиентного спуска. В реализации градиентного спуска можете предполагать, что на вход подаются функции с единственным, глобальным минимумом. Перед тем, как писать код, ответьте себе на следующие вопросы: Как понять, что пора остановиться? Это может зависеть от градиента или расстояния между двумя соседними шагами алгоритма, так и от числа уже выполненных итераций.</p> <p>Как правильно менять величину шага (learning rate) от итерации к итерации?</p>	
--	--	---	--

### Рекомендуемые ресурсы для дополнительного изучения:

1. Частная производная. [Электронный ресурс] – Режим доступа: [https://ru.wikipedia.org/wiki/Частная\\_производная](https://ru.wikipedia.org/wiki/Частная_производная).
2. Градиент. [Электронный ресурс] – Режим доступа: <https://ru.abcdef.wiki/wiki/Gradient>.
3. Обзор градиентных методов в задачах математической оптимизации. [Электронный ресурс] – Режим доступа: <https://habr.com/ru/post/413853/>.
4. Основы линейной регрессии. [Электронный ресурс] – Режим доступа: <https://habr.com/ru/post/514818/>.