

## ТЕХНОЛОГИЧЕСКАЯ КАРТА ЗАНЯТИЯ

**Тема занятия:** Линейная регрессия.

**Аннотация к занятию:** обучающиеся познакомятся с новой моделью машинного обучения — линейной регрессией. Узнают принцип её работы, геометрический смысл и некоторые свойства. Обсудят, как обучать линейную регрессию и оценивать качество её работы.

**Цель занятия:** сформировать у учеников представление о линейной регрессии, принципах её работы, геометрическом смысле и некоторых свойствах. Научить обучать линейную регрессию и оценивать качество её работы.

**Задачи занятия:**

- познакомить школьников с моделью машинного обучения — линейной регрессией и её геометрической интерпретацией;
- сформировать понятие «линейная регрессия»;
- познакомить с некоторыми свойствами модели: её интерпретируемостью и неспособностью подстраиваться под сложные зависимости в данных;
- обсудить метрику качества для линейной регрессии и способ её обучения.

## Ход занятия

Этап занятия	Время	Деятельность педагога	Комментарии, рекомендации для педагогов
Организационный этап	5 мин.	Добрый день! Приветствую вас на нашем уроке.	Приветствие. Создание в классе атмосферы психологического комфорта.
Постановка цели и задач занятия. Мотивация учебной деятельности обучающихся	7 мин.	<p>Сегодня мы познакомимся с новой моделью машинного обучения — линейной регрессией. Тема урока — «Линейная регрессия».</p> <p>Прочитав тему урока, как вы думаете, что мы будем изучать сегодня?</p> <p><b>Возможные ответы школьников:</b></p> <ul style="list-style-type: none"> <li>узнаем, что такое линейная регрессия, её устройство, принцип работы;</li> <li>поймём, можно ли обучить линейную регрессию.</li> </ul> <p><b>Мы узнаем:</b></p> <ul style="list-style-type: none"> <li>устройство модели линейной регрессии и её геометрический смысл;</li> <li>некоторые особенности линейной регрессии;</li> <li>оценку качества для линейной регрессии и способ её обучения.</li> </ul>	Способствовать обсуждению мотивационных вопросов.

<p><b>Изучение нового материала</b></p>	<p>50 мин.</p>	<p><b>Линейная регрессия</b> — модель машинного обучения, которая используется только для решения задач регрессии. К сожалению, для классификации её не применяют. Скоро мы поймём, почему это так.</p> <p>Идея линейной регрессии — построить линейную поверхность в качестве решающего правила. Давайте на реальном примере поймём, что это значит.</p> <p>Рассмотрим задачу регрессии: предсказание цены квартиры по её характеристикам. На слайде вы видите первые три строки обучающего датасета для этой задачи. В датасете есть пять признаков: площадь квартиры, её удалённость от центра города, год постройки, количество лет после ремонта и тип постройки дома, в котором квартира находится. Признаки пронумерованы от 1 до 5. Целевая переменная — цена квартиры, выраженная в тысячах долларов. Она обозначена как <math>Y</math>.</p> <p>Модель линейной регрессии ставит в соответствие каждому признаку свой коэффициент — некое действительное число. На слайде коэффициенты обозначены как <math>k_i</math>. Первому признаку соответствует коэффициент <math>k_1</math>, второму признаку — коэффициент <math>k_2</math> и так далее. 5 признаков — 5 коэффициентов. Кроме этого, у модели есть свободный член, который не связан ни с одним признаком — <math>k_0</math>. Таким образом, количество коэффициентов линейной регрессии равно количеству признаков датасета плюс один. В нашем случае их 6.</p> <p>Для каждой квартиры наша модель будет предсказывать ответ в виде линейной комбинации признаков, взятых с соответствующими им коэффициентами. Пусть у нас есть элемент <math>x</math> со значениями признаков <math>x_1, x_2, x_3, x_4</math> и <math>x_5</math>. То есть, <math>x_1</math> — значение первого признака этого элемента, <math>x_2</math> — значение второго признака</p>	<p>Для справки: Сайт: <a href="https://habr.com/ru/post/514818/">https://habr.com/ru/post/514818/</a></p> <p>Перед уроком рекомендуется ознакомиться с материалами, представленным и на сайте.</p>
---	----------------	---	--

и так далее. Тогда ответ линейной регрессии на этот элемент будет иметь вид, представленный на слайде.

Здесь  $y$  с шапочкой — это ответ линейной регрессии. Он получается, как  $k_0$  плюс сумма  $k_i \cdot x_i$ , умноженных на  $x_i$ , что просто есть  $k_0$  плюс  $k_1$  умножить на  $x_1$ , плюс  $k_2$  умножить на  $x_2$  и так далее, плюс  $k_5$  умножить на  $x_5$ .

Для наглядности покажем, как выглядит ответ линейной регрессии для первых двух элементов датасета.

Чтобы получить ответ  $y_1$  с шапочкой для первого элемента датасета, нужно подставить вместо  $x_1$ ,  $x_2$  и других  $x$  значения соответствующих признаков элемента. Эти значения — 25 для  $x_1$ , 3 для  $x_2$  и так далее.

Чтобы получить ответ  $y_2$  с шапочкой для второго элемента датасета, нужно подставить вместо  $x$  значения его признаков. Формулы для ответов линейной регрессии на два первых элемента датасета вы видите на экране.

Подчеркну, что линейная регрессия вычисляет ответ на каждый элемент, используя одни и те же значения коэффициентов  $k$ .

В чём же состоит обучение модели линейной регрессии?

**Ответы учеников.**

А вот в чём: на основе входных данных нужно подобрать такие коэффициенты  $k_i$ , чтобы ответы линейной регрессии, полученные с помощью этой формулы, минимально отличались от настоящих. Грубо говоря, чтобы в среднем для всех элементов обучающей выборки линейная регрессия хорошо предсказывала ответы.

Поговорим о геометрическом смысле линейной регрессии. Это поможет нам лучше понять природу модели и процесс подбора её коэффициентов.

Рассмотрим простой случай, когда в датасете всего один признак. Пусть этот признак — площадь квартиры, и пусть в нашем датасете всего 5 квартир. Датасет вы видите на экране. Задача линейной регрессии — по площади квартиры научиться предсказывать её цену.

Изобразим каждую квартиру из датасета в виде точки на декартовой плоскости. По оси  $X$  отложим площадь квартиры, по оси  $Y$  — её цену. Получим рисунок, который вы видите на экране справа.

Так как в датасете всего один признак, то формула предсказания линейной регрессии будет выглядеть так:  
 $Y$  с шапочкой равно  $k_1$  умножить на  $x_1$  плюс  $k_0$ , где  $x_1$  — значение признака входящего элемента.

Допустим, мы подобрали значения коэффициентов  $k_0$  и  $k_1$  для линейной регрессии. Пусть  $k_1$  равен 1,5,  $k_0$  равен 10. Подставим эти коэффициенты в формулу и получим:  $Y$  с шапочкой равно 1,5 умножить на  $x_1$  плюс 10.

Полученная формула — это уравнение прямой. Стандартное уравнение прямой записывается как  $y$  равно  $kx$  плюс  $b$ . У нас вместо  $x$  по оси абсцисс отложен  $x_1$ , а коэффициенты  $k$  и  $b$  равны 1,5 и 10 соответственно.  
Давайте изобразим эту прямую на нашем рисунке.

Как с геометрической точки зрения получить ответ на входящий элемент с помощью линейной регрессии?

#### Ответы учеников

Посмотрим на вторую квартиру в датасете. С точки зрения формулы, ответ модели мы получим, подставив значение площади квартиры — то есть 60 — вместо  $x_1$ . Получим 100. В таблице на экране ответы линейной регрессии на все 5 элементов выборки внесены в правый столбец.

С точки зрения геометрии, чтобы получить ответ на первый элемент датасета, нужно отложить значение его признака — то есть 60 — на оси  $X$ , а затем подняться вверх до пересечения с прямой. Ордината точки пересечения и будет ответом линейной регрессии на первую квартиру.

Ответ нашей линейной регрессии далёк от правильного, который равен шестидесяти. Найдём на рисунке точку, соответствующую второй квартире. Она обозначена красным. Её ордината равна 60, т.е. правильному значению цены второй квартиры. Получается, длина линии между красной точкой и точкой пересечения зелёной линии с прямой показывает, насколько линейная регрессия ошиблась в предсказании цены этой квартиры.

Точно так же мы можем изобразить величины ошибок линейной регрессии для других квартир из датасета. Они отмечены зелёным на рисунке. Видно, что предсказание модели для пятой квартиры очень близко к правильному значению, а для второй квартиры ошибка оказалась очень большой. То же самое можно увидеть, сравнив два последних столбца в таблице.

Допустим теперь, что мы подобрали другие коэффициенты для линейной регрессии —  $k_1$  равно 0.5,  $k_0$  равно 10. Прямую, соответствующую этой формуле, вы видите на рисунке. Для такой линейной регрессии ошибки на всех пяти квартирах из датасета довольно велики.

Сравнивая две линейные регрессии, интуитивно мы понимаем, что модель слева на слайде лучше. Величины её ошибок на пяти квартирах из датасета в среднем меньше.

В идеале мы бы хотели, чтобы наша модель не делала ошибок совсем. С геометрической точки зрения это значило бы, что все точки, соответствующие элементам датасета, лежали бы на прямой, соответствующей модели линейной регрессии. В реальности такого, конечно, никогда не будет. В реальности практически никогда не бывает так, что ответ зависит от признаков строго линейно.

Получается, задача обучения линейной регрессии состоит в подборе таких коэффициентов, чтобы точки датасета лежали как можно ближе к прямой. Близость точек к прямой при этом измеряется расстоянием от точки до прямой по оси ординат. О поиске лучших коэффициентов для модели линейной регрессии мы поговорим позже.

Давайте усложним задачу. Пусть в данных теперь два признака: площадь квартиры и удалённость от центра города. В этом случае элементы датасета можно представить в виде точек в трёхмерном пространстве. По оси  $X$  будет отложена площадь квартиры, по оси  $Y$  — удалённость от центра города, по оси  $Z$  — ответ, т.е. стоимость квартиры.

Модель линейной регрессии для решения такой задачи будет иметь 3 коэффициента, и уравнение будет иметь вид, представленный на слайде. Это уравнение будет задавать уже плоскость, а не прямую. И чем ближе точки будут лежать к плоскости по оси Z, тем меньше будет ошибка линейной регрессии на нашем датасете.

По аналогии можно добавить ещё один или несколько признаков в наш датасет. Если признаков  $n$  штук, то у линейной регрессии будет  $n+1$  коэффициент, и она будет строить  $n$ -мерную линейную гиперплоскость в  $n+1$ -мерном пространстве. Если вы не знакомы с понятием гиперплоскость, не пугайтесь: это то же самое, что обычная плоскость, только в пространстве большей размерности. Нарисовать эту гиперплоскость при  $n$  большем или равном 3 мы уже не сможем. Но к линейной регрессии с любым количеством признаков применимы те же рассуждения, что мы проделали выше для  $n$  равного 1 и 2.

На данном этапе должно стать понятно, почему линейную регрессию не получится применить к задаче классификации. Если мы возьмём датасет для задачи бинарной классификации и изобразим элементы из него на плоскости, получим картину, представленную на слайде. Понятно, что мы не можем построить прямую, проходящую близко ко всем точкам. Плюс, ответ линейной регрессии на входящий элемент — это ордината одной из точек прямой, т.е. некое число от минус бесконечности до плюс бесконечности. Непонятно, как такой ответ можно перевести в класс 0 или 1.

Для решения задачи классификации есть другая линейная модель, идейно очень похожая на линейную регрессию. Это логистическая регрессия. С ней вы познакомитесь далее в курсе.



Итак, мы познакомились с устройством линейной регрессии и разобрали её геометрический смысл. Давайте немного поговорим о свойствах линейной регрессии.

Во-первых, линейная регрессия плохо работает на тех датасетах, в которых зависимость между целевой переменной и признаками имеет сильно нелинейный характер. Линейная регрессия предполагает, что между признаками и ответом есть линейная зависимость, и пытается её восстановить, подбирая подходящие коэффициенты. Если же эта зависимость на самом деле не является линейной, как в датасете на слайде, линейная регрессия будет работать плохо.

Какую бы прямую мы ни построили, ошибка регрессии будет довольно большой.

В реальной жизни практически никогда не бывает так, что зависимость целевой переменной от признаков похожа на линейную. Чаще всего зависимости намного сложнее. Даже цена дома зависит от характеристик дома не совсем линейно. Поэтому модель линейной регрессии чаще всего работает не очень хорошо на реальных датасетах. Таким образом, линейная регрессия — это довольно слабая модель, не способная выражать зависимости сложнее линейных.

Однако даже если зависимость целевой переменной от признаков нелинейна, существует способ преобразовать признаки так, чтобы линейная регрессия работала лучше на этом датасете. Это называется *kernel trick* — по-русски «трюк с ядром».

Следующее и последнее свойство линейной регрессии, о котором мы поговорим, — это интерпретируемость. Линейная регрессия —

хорошо интерпретируемая модель. Формула линейной регрессии явно показывает, как из значений признаков элемента получился ответ. Модули коэффициентов линейной регрессии показывают, насколько большую роль тот или иной признак играет в получении ответа.

Например, пусть линейная регрессия поставила в соответствие четвёртому признаку коэффициент, равный 100, а пятому признаку — коэффициент, равный 1. Получается, что значение четвёртого признака входит в ответ линейной регрессии с коэффициентом 100, а значение пятого признака — с коэффициентом 1. Это говорит о том, что значение четвёртого признака больше влияет на ответ, чем значение пятого признака. Отсюда мы можем сделать вывод, что количество лет, прошедших после ремонта квартиры, влияет на её цену сильнее, чем тип постройки дома, в котором эта квартира находится.

Тут, конечно, нужно принимать во внимание то, насколько велики значения различных признаков. Если третьему признаку модель поставила в соответствие коэффициент 1, а четвёртому признаку — коэффициент 10, то третий признак все равно будет сильнее влиять на ответ, чем четвёртый признак. Потому что третий признак выражается большим по модулю четырёхзначным числом, а четвёртый — маленьким по модулю однозначным. И  $x_3$ , умноженный на  $k_3$ , для любого элемента датасета будет всё же больше, чем  $x_4$ , умноженный на  $k_4$ .

Напомню, что линейная регрессия представляет ответ в качестве линейной комбинации признаков датасета. Если в датасете  $n$  признаков, то у модели будет  $n+1$  коэффициент.

Рассмотрим способы оценки качества линейной регрессии. Для этого мы познакомимся с метриками качества. Это такие функции, которые принимают на вход предсказания линейной регрессии об элементах датасета и правильные ответы, а затем выдают число — некоторую оценку того, насколько предсказания модели близки к правильным ответам.

Для начала разберёмся, зачем вообще нужна метрика качества.

В целом, оценка качества модели позволяет выяснить, насколько хорошо модель справляется с задачей.

Пусть у нас есть две разные модели линейной регрессии с разными наборами коэффициентов для предсказания цены квартиры. Формулы двух моделей вы видите на экране.

Как понять, какая из этих двух моделей работает лучше? Так как в датасете пять признаков, у нас нет возможности визуализировать элементы датасета и модель на графике. Мы не можем увидеть, как хорошо та или иная модель эти точки описывает. И даже если бы мы могли это визуализировать, понять по графику, какая из двух регрессий работает лучше, сложно.

Поэтому для того, чтобы оценить качество предсказаний модели, нужно ввести метрику качества.

Первая метрика качества, которую мы рассмотрим — это Mean Absolute Error, или MAE. Вычисляется она следующим образом.

Пусть у нас есть модель линейной регрессии с некоторыми коэффициентами. Запишем в отдельный столбец предсказания этой модели для всех элементов датасета. На слайде это самый правый столбец.

Для каждого элемента вычислим ошибку, которую модель сделала на этом элементе. Ошибка модели — это модуль разности между правильным значением ответа  $y$  и ответом модели  $y$  с шапочкой. Ответ модели может быть как больше, так и меньше реального значения: именно поэтому разницу нужно брать по модулю.

Запишем ошибку для каждого элемента в отдельный столбец справа.

Затем посчитаем среднее значение ошибки на всех элементах датасета. Для этого возьмём сумму всех ошибок и поделим её на число элементов. Для датасета на слайде, в котором всего 3 элемента, средняя ошибка будет составлять 12,66. Это значит, что предсказания нашей модели линейной регрессии отличаются от верных в среднем на 12,66 для каждого элемента набора. Чем меньше величина MAE, тем лучше предсказания модели на этом датасете.

Таким образом можно вычислить значения MAE на одном датасете для двух разных моделей и сравнить. Модель, у которой MAE меньше, имеет меньшую среднюю ошибку на элементах датасета — её можно считать лучшей.

Ранее мы обсуждали геометрический смысл линейной регрессии, где мы рассматривали датасет с одним признаком и изображали элементы датасета в виде точек на плоскости, а линейную регрессию — в виде прямой. Заметим, что MAE — это средняя длина зелёных линий на графике.

На слайде вы видите общую формулу вычисления MAE для датасета с  $n$  элементами.  $Y_i$  — правильный ответ на  $i$ -ый элемент выборки,  $y$  с шапочкой  $i$  — ответ модели на  $i$ -ый элемент выборки.

Для справки:  
Средняя абсолютная масштабированная ошибка (Mean Absolute Error) — показатель ошибки, использующийся для точности модели.

На одной метрике мы не остановимся. Рассмотрим вторую метрику, которая используется для оценки качества работы моделей в задачах регрессии. Это Mean Squared Error, или MSE. По-русски — среднеквадратичная ошибка. Её формулу вы видите на слайде внизу.

Эта метрика отличается от MAE тем, что вычисляет среднее значение не модулей ошибок на элементах, а квадратов ошибок.

Формула вычисления MSE для нашего датасета из трёх элементов показана на экране. Возводим значение каждой ошибки из правого столбца в квадрат и берём среднее значение.

С точки зрения геометрии MSE отражает среднее значение квадратов длин зелёных отрезков.

Заметим, что метрики MAE и MSE можно вычислять не только для ответов линейной регрессии, но и для любой другой модели, например, для KNN. MAE и MSE — это метрики для оценки качества любой модели, решающей задачу регрессии.

Здесь может возникнуть вопрос, зачем нужна метрика MSE. Кажется, что MAE — очень логичная метрика, отражающая природу ошибок линейной регрессии. Действительно, сравнение качества моделей по среднему значению их ошибок на элементах датасета выглядит довольно естественным. Не очень понятно, зачем эти ошибки возводить в квадрат.

На самом деле метрика MSE очень важная и нужная. Зачем она нужна, вы узнаете в одном из следующих видео, когда будете изучать обучение модели линейной регрессии с помощью градиентного спуска.

Для справки:  
MSE  
расшифровывает  
ся как Mean  
Squared Error и  
переводится как  
«средняя  
квадратическая

		<p>Вернёмся к процессу обучения линейной регрессии. Мы ещё не знаем, как линейная регрессия подбирает значения коэффициентов на основе тренировочного датасета.</p> <p>В целом, задача обучения линейной регрессии состоит в том, чтобы найти такие коэффициенты, при которых значение выбранной метрики качества на тренировочном датасете было бы наименьшим. То есть найти такие коэффициенты, при которых ответы линейной регрессии на элементы обучающей выборки были бы лучшими с точки зрения метрики качества.</p> <p>Есть два способа обучить линейную регрессию. Первый — с помощью решения матричного уравнения. Второй — с помощью градиентного спуска. В этом видео мы обсудим первый способ. Второй вы разберёте позже, после того как познакомитесь с понятиями производной и градиентного спуска.</p> <p>Посмотрим на уравнение линейной регрессии. Обозначим через <math>k</math>-вектор столбец всех коэффициентов регрессии, кроме <math>k_0</math>, а через <math>x</math> с верхним индексом 1 — вектор-столбец признаков первого элемента.</p> <p>При таких обозначениях предсказание линейной регрессии для первого элемента можно представить, как умножение транспонированного вектора <math>x_1</math> на вектор <math>k</math>, не забывая при этом прибавить свободный член <math>k_0</math>.</p> <p>Действительно, ведь произведение <math>x_1</math>, транспонированного на <math>k</math>, — это то же, что и <math>k_1</math> умножить на значение первого признака, плюс <math>k_2</math> умножить на значение второго признака и так далее. То есть, в случае первого элемента, — <math>k_1</math> на 25 плюс <math>k_2</math> на 3 и так далее.</p>	<p>ошибка». Суть метода заключается в том, чтобы минимизировать сумму квадратов отклонений фактических значений от расчётных (SSE, Sum of Squared Errors). Если полученную сумму разделить на число наблюдений, то получится та самая MSE.</p>
--	--	--	--

Сделаем небольшой трюк. Добавим к вектору коэффициентов регрессии  $k$  свободный член  $k_0$ . А к вектору признаков элемента в начало добавим значение 1. Смотрите, что получится: при таких векторах  $x_1$  и  $k$  ответ линейной регрессии  $y_1$  с шапочкой на первый элемент представляется в виде произведения транспонированного вектора  $x_1$  на вектор  $k$ . Добавлять к этому произведению дополнительные слагаемые уже не нужно.

Замечу, что произведение транспонированного вектора  $x_1$  на вектор  $k$  — это скалярное произведение векторов  $x_1$  и  $k$ . Таким образом, линейный алгоритм выражает ответ на входящий элемент в виде скалярного произведения вектора коэффициентов на вектор признаков элемента.

Точно так в векторном виде можно представить ответы линейной регрессии и для остальных элементов датасета. Берём вектор признаков элемента и добавляем к нему в начало единицу. Умножаем этот вектор на вектор коэффициентов  $k$  и получаем ответ линейной регрессии на этот элемент.

Добавление единицы к вектору признаков элементов датасета можно представить, как добавление в датасет ещё одного константного нулевого признака. Этот признак для всех элементов будет равен единице. Ему будет соответствовать свободный член  $k_0$ . Тогда вектор  $x_1$  признаков каждого элемента будет состоять из шести значений с единицей в начале. Ответ линейной регрессии на элемент можно представить в виде произведения транспонированного вектора признаков элемента на вектор коэффициентов  $k$ .

Мы научились представлять ответ линейной регрессии на один элемент в виде произведения двух векторов. Пойдём дальше.

Обозначим через  $Y$  большое с шапочкой вектор ответов линейной регрессии на все элементы датасета. Если в датасете 3 элемента,  $y$  вектора тоже будет 3 элемента.

Через  $X$  большое обозначим матрицу признаков всех элементов датасета. То есть, строки матрицы  $X$  — это векторы признаков каждого элемента датасета. Добавим к каждому вектору признаков датасета единицу в начало. Первый столбец матрицы  $X$  будет состоять из единиц.

За  $k$  маленькое обозначим вектор коэффициентов линейной регрессии.

Тогда вектор ответов регрессии  $Y$  с шапочкой можно записать в виде матричного произведения:  $Y$  с шапочкой равно  $k$  умножить на  $X$  транспонированное. Проверьте, что данная формула верна: вычисление ответов в линейной регрессии соответствует матричному произведению, которое вы видите на экране.

Это матричная запись уравнения линейной регрессии. В зависимости от количества признаков и элементов в датасете, в векторе  $k$  и матрице  $X$  будет разное количество элементов, строк и столбцов. Но матричная запись будет иметь тот же самый вид.

Как найти оптимальные значения коэффициентов  $k$ ? Такие, чтобы ошибка линейной регрессии на элементах датасета была минимальна и  $Y$  с шапочкой был как можно ближе к  $Y$ ?

Оказывается, это просто: оптимальные значения  $k$  получаются с помощью формулы справа внизу на слайде. Чтобы получить значение  $k$ , мы должны взять матрицу датасета  $X$ , транспонировать её и умножить справа на  $X$ . Затем от полученной матрицы взять



		<p>обратную и снова умножить её на <math>X</math> транспонированное справа. Полученное произведение справа умножить на вектор коэффициентов <math>k</math>. Вы можете самостоятельно убедиться в том, что все матричные умножения в формуле корректны.</p> <p>Линейная регрессия с коэффициентами, полученными по этой формуле, обладает наименьшим возможным значением метрики MSE на элементах датасета. Какие бы другие коэффициенты для регрессии мы ни выбрали, линейная регрессия с ними будет иметь большую ошибку MSE на элементах нашего датасета.</p> <p>Получается, обучение линейной регрессии происходит так: мы берём тренировочный датасет, получаем из него матрицы <math>X</math> и <math>Y</math> и вычисляем значение вектора коэффициентов <math>k</math> с помощью формулы на слайде. Далее на любой элемент из тренировочного или тестового датасета линейная регрессия выдаёт ответ, умножая вектор <math>k</math> на вектор признаков элемента.</p> <p>Такой матричный способ нахождения коэффициентов <math>k</math> — первый способ обучения линейной регрессии. Со вторым способом — обучением с помощью градиентного спуска — вы познакомитесь вна одном из следующих уроков.</p>	
<p><b>Закрепление изученного материала</b></p>	<p>15 мин.</p>	<p><b>Вопросы для обсуждения</b></p> <ul style="list-style-type: none"> <li>• Что такое линейная регрессия?</li> <li>• В чём смысл геометрической интерпретации линейной регрессии?</li> <li>• Как оценить качество линейной регрессии?</li> <li>• Как обучить линейную регрессию?</li> </ul>	<p>Педагог организует беседу по вопросам.</p>

<p><b>Этап подведения итогов занятия (рефлексия)</b></p>	<p>8 мин.</p>	<p><b>Вопросы для обсуждения</b></p> <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>	<p>Педагог способствует размышлению обучающихся над вопросами.</p>
<p><b>Информация о домашнем задании, инструктаж по его применению</b></p>	<p>5 мин.</p>	<p>В домашнем задании вам предстоит написать свой класс линейной регрессии и протестировать его для решения задачи регрессии. Работать мы будем с тем же датасетом пингвинов, что и на семинаре. Ссылка на скачивание датасета: <a href="#">датасет</a>.</p> <p>Ваша задача — написать код для методов класса MyLinearRegression. Несколько комментариев к заданию:</p> <ul style="list-style-type: none"> <li>• для каждого метода класса (fit, predict) описано, что этот метод принимает на вход и какой функционал реализует. По сути, fit — это аналог метода fit модели линейной регрессии из sklearn, predict — аналог метода predict модели линейной регрессии из sklearn;</li> <li>• в методе fit при получении коэффициентов линейной регрессии вам нужно получить две переменные: self.coef_ и self.intercept_. В self.coef_ должен получиться массив коэффициентов, которые модель поставила в соответствие признакам датасета. В self.intercept_ должно получиться одно число — коэффициент-свободный член, который выучила модель. По сути, self.coef_ и self.intercept_ должны быть аналогами этих же переменных модели Linear Regression из sklearn.</li> </ul>	

		<p><b>Алгоритм реализации метода fit:</b></p> <ol style="list-style-type: none"><li>1. Перевести X и y в numpy array (для удобства). Это уже реализовано;</li><li>2. Добавить к X первый столбец из единиц. Подсказка: чтобы это сделать, удобно использовать np.hstack. Подумайте, как именно;</li><li>3. Получить массив коэффициентов k по формуле;</li><li>4. Разбить полученный массив коэффициентов k на self.coef_ и self.intercept_</li></ol> <p><b>Алгоритм реализации метода predict:</b></p> <ol style="list-style-type: none"><li>1. Получить значения y_pred, используя значения выученных коэффициентов и входящих признаков X по формуле.</li></ol>	
--	--	--	--

### Рекомендуемые ресурсы для дополнительного изучения:

1. Основы линейной регрессии. [Электронный ресурс] – Режим доступа: <https://habr.com/ru/post/514818/>.
2. Оценка результатов линейной регрессии. [Электронный ресурс] – Режим доступа <https://habr.com/ru/post/195146/>.
3. Основы линейной регрессии. [Электронный ресурс] – Режим доступа: <http://statistica.ru/theory/osnovy-lineynoy-regressii/>.
4. Методы оценки качества прогноза. [Электронный ресурс] – Режим доступа: <https://habr.com/ru/post/19657/>.