

Геннадий Суворов  
data scientist, финансовый аналитик

# Ловец изменений

## ПРИМЕНЕНИЕ ИСКУССТВЕННОГО ИНТЕЛЛЕКТА В БИРЖЕВОЙ ТОРГОВЛЕ COMMODITY FUTURES И ИХ ПРОИЗВОДНЫМИ

Поставленную проблему решает нейронная сеть, которая с помощью искусственного интеллекта реализуется в виде нелинейной функции. В качестве аргументов нейронная сеть может использовать данные в любом формате (действительные числа, целые числа, строки и так далее) после их обработки data engineer.

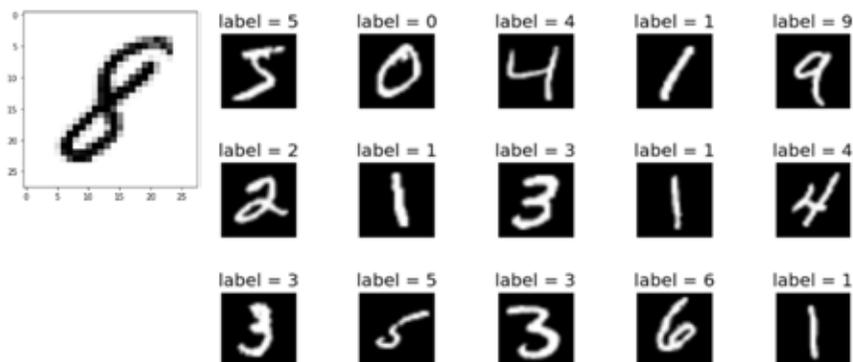
Автор сразу хочет предупредить читателя, что предлагает ему очень техническую и предметную статью об особенностях применения методов искусственного интеллекта (ИИ) к торговле производными финансовыми инструментами на мировых биржах commodities. Статья написана на основании реального проекта применения ИИ в реальной трейдинговой компании, который был реализован более чем два года назад. Но она совершенно не является техническим описанием выполненного проекта, перед вами — концептуальный анализ того, каким образом такой проект может быть реализован сегодня.

Что касается ИИ. Здесь вы не найдете рассуждений о цифровых концлагерях, о том, как

суперкомпьютер SkyNet захватывает мир и как роботы расстреливают ни в чем не повинное гражданское население. Мы будем рассматривать ИИ в узком смысле: в качестве некой компьютерной системы (скрипта, программы), которая специальным образом «обучена» для решения некоторой вполне конкретной практической проблемы.

В дальнейшем изложении придется обильно использовать англоязычные термины. Дело в том, что часть понятий на российском финансовом пространстве просто отсутствует, а другие понятия не имеют полноценного перевода на русский, не говоря уже об их строгой дефиниции.

Рис. 1. КАРТИНКА И НАБОРЫ ДАННЫХ: КАРТИНКА+LABEL



Еще одно замечание относительно понятия «искусственный интеллект». В первой части речь будет идти в основном о нейронных сетях. Поэтому термины «система ИИ» и «нейронная сеть» будут пониматься как синонимы. С помощью технологий ИИ нейронная сеть реализуется в виде нелинейной функции, которая решает поставленную проблему. В качестве аргументов нейронная сеть может использовать действительные числа, целые числа, строки, категории и тому подобное. Практически речь идет о данных в любом формате после их обработки data engineer.

Все нейронные сети работают с некоторой погрешностью. Приемлемы ли параметры точности, аккуратности и так далее этой функции для решаемой проблемы --- в каждом конкретном случае представляет собой отдельный вопрос. Но это не вопрос данной статьи. Поэтому в дальнейшем изложении будем считать, что в нашем случае эти параметры приемлемы.

Далее.

Прежде чем нейронная сеть начнет решать реальные проблемы, ее необходимо обучить. Существует несколько базовых методов обучения: supervised, unsupervised, reinforcement learning.

Supervised learning. При такой форме обучения data scientist готовит большую выборку данных с известными результатами. Data engineer в процессе предварительной обработки данных переводит данные в формат, который может воспринять нейронная сеть.

Пример. В очень известном в мире ИИ-проекте MNIST решается задача распознавания рукописных написаний цифр. Входные данные для MNIST представлены в виде картинок 28x28 пикселей и label (цифра, которую эта картинка представляет). Всего таких наборов в MNIST было 60 тысяч для обучения и 10 тысяч для тестирования нейронной сети. Образцы данных MNIST представлены на рис. 1. От нейронной сети требова-

лось предсказать (распознать), какую именно цифру представляет картинка.

В проекте MNIST нейронная сеть дает аккуратность предсказания 99%.

После определения наборов входных data scientist определяет тип и конфигурацию нейронной сети и параметры ее обучения. Data engineer преобразует данные для ввода в нейронную сеть. В упрощенном виде тренировка нейронной сети в технологии supervised learning заключается в следующем: на вход нейронной сети последовательно подаются входные данные из наборов данных. Сеть вычисляет, что это за label (с точки зрения сети), сравнивает его с правильным результатом и вычисляет значение ошибки — loss function. Далее сеть подстраивает свои внутренние параметры для того, чтобы минимизировать значение loss function. Процедура подстройки нейронной сети повторяется на всем обучающем наборе данных несколько тысяч раз. После этого на тестовом наборе данных, который не участвовал в «обучении» сети, проводится анализ результатов на предмет точности и аккуратности предсказаний.

Unsupervised learning (самостоятельное обучение). При такой форме обучения data scientist готовит набор данных для обучения, но без классификации и без labels. Система ИИ должна сама определить подобие объектов, классифицировать их и построить плотность распределения вероятностей в классификации. Этот тип обучения в проекте не использовался и поэтому останавливаться на нем мы не будем.

Reinforcement learning (обучение с подкреплением) будет детально рассмотрено во второй части статьи, где обучение этого типа будет предложено для обучения интеллектуального торгового агента.

## Постановка цели

Целью большинства проектов на финансовом рынке является предсказание изменения цены моделируемого

финансового инструмента. Наша цель для первой фазы проекта состоит в том, чтобы предсказать будущее изменение цены инструмента в day-to-day трейдинге.

Опишем постановку задачи проекта.

*Дано:*

Сектор — Natural Gas;

Используемые активы — календарные спреды между месячными фьючерсами;

Горизонт предсказания — до конца дня (точнее, до Daily Settlement Time фьючерса);

Данные для сети — исторические графики изменения цен всех релевантных фьючерсов и спредов за два последних года.

*Требуется:*

Предсказать изменения цены моделируемого актива (спреда) в некотором будущем периоде. Период ограничен горизонтом предсказания и может иметь любую временную протяженность — час, два и так далее.

## Финансовые параметры

Рассмотрим некоторые финансовые параметры активов в нашем проекте.

Ликвидность актива. Спреды в секторе Natural Gas являются активами с очень высокой ликвидностью. Спреды торгуются непрерывно и с большой частотой. Это не значит, что держатель календарного спреда может его мгновенно продать или купить. Обычно в стакане на каждую ask-price и bid-price стоит очередь заявок.

Типы активов. Мы не будем моделировать naked фьючерсы. Цена фьючерса для нашей технологии слишком волатильна. Календарные спреды между фьючерсами намного более устойчивы. Они настолько популярны, что биржа предоставляет трейдеру возможность торговать ими как отдельными инструментами. Покупая/продавая спред между двумя календарными месяцами, трейдер получает в одной транзакции сразу две позиции во фьючерсах: long

и short. В таблице предоставлен пример матрицы календарных спредов.

Например, на пересечении строки JUN14 и столбца SEP14 показаны котировки спреда между фьючерсами июня и сентября, на который трейдер может выставить заявку на покупку или продажу.

Сектор. Мы будем рассматривать торговлю производными инструментами, зависящими от цен фьючерсов в секторе Natural Gas. Это ликвидный, динамичный и относительно волатильный сектор. В секторах Natural Gas существует ярко выраженная сезонность.

Горизонт предсказания. Горизонт предсказания — day trading. Как правило, к концу дня day-trader не имеет открытых позиций или они захеджированы. Это не значит, что к торговле фьючерсами не применим фундаментальный анализ или технический анализ, просто к моменту открытия торгов фьючерсами все новости уже «сидят» в цене фьючерса.

Данные для обучения сети и входные данные. Для обучения сети будут использованы исторические данные OHLC всех релевантных активов. Как было сказано выше, у нас есть доступ к историческим наборам OHLC таких активов за два последних года.

Возникает вопрос — что будет предсказывать нейронная сеть?

Прежде чем отвечать на этот вопрос, рассмотрим, какие факторы влияют на цены наших активов. В торговле акциями главными движущими факторами цены являются новости — по акции, по ее сектору в целом, а также квартальные и годовые отчеты. По криптовалюте — новости о рынке криптовалют в целом, о легитимности этого рынка в конкретной стране, степени легальности использования криптовалюты в качестве расчетной валюты.

А что является двигателями наших активов?

Уровень цены торговли фьючерсами определяют новости сектора «Natural

ТАБЛИЦА. ПРИМЕР МАТРИЦЫ КАЛЕНДАРНЫХ СПРЕДО

| Ask    |       | AskQ   |       | Ask    |       | AskQ   |       | Ask    |        | AskQ   |        | Ask    |        | AskQ   |        | Ask   |  | AskQ |  |
|--------|-------|--------|-------|--------|-------|--------|-------|--------|--------|--------|--------|--------|--------|--------|--------|-------|--|------|--|
| 99.735 | 16234 | 99.725 | 233   | 99.705 | 6235  | 99.660 | 12850 | 99.590 | 188    | 99.510 | 5409   | 99.410 | 1997   | 99.290 | 66     |       |  |      |  |
| 99.730 | 9713  | 99.720 | 15231 | 99.700 | 811   | 99.655 | 43    | 99.585 | 4087   | 99.505 | 183    | 99.405 | 613    | 99.285 | 2915   |       |  |      |  |
| MAR14  |       | JUN14  |       | SEP14  |       | DEC14  |       | MAR15  |        | JUN15  |        | SEP15  |        | DEC15  |        |       |  |      |  |
| MAR14  |       | 0.010  | 3060  | 0.035  | 3516  | 0.080  | 783   | 0.145  | 186    | 0.230  | 183    | 0.330  | 387    | 0.450  | 1506   |       |  |      |  |
|        |       | 0.005  | 12844 | 0.030  | 275   | 0.075  | 382   | 0.140  | 677    | 0.220  | 3268   | 0.320  | 1289   | 0.440  | 66     |       |  |      |  |
|        |       |        |       | 0.025  | 44810 | 0.070  | 539   | 0.135  | 93     | 0.220  | 215    | 0.320  | 140    | 0.440  | 140    |       |  |      |  |
|        |       | JUN14  |       | 0.020  | 98014 | 0.065  | 12878 | 0.130  | 4358   | 0.215  | 673    | 0.310  | 1289   | 0.430  | 67     |       |  |      |  |
|        |       |        |       |        |       | 0.050  | 97678 | 0.115  | 4935   | 0.200  | 1184   | 0.300  | 1152   | 0.415  | 1      |       |  |      |  |
|        |       |        |       | SEP14  |       | 0.045  | 60619 | 0.110  | 3851   | 0.195  | 5      | 0.290  | 2582   | 0.410  | 66     |       |  |      |  |
|        |       |        |       |        |       |        |       | 0.070  | 108387 | 0.155  | 17713  | 0.250  | 238    | 0.370  | 824    |       |  |      |  |
|        |       |        |       |        |       | DEC14  |       | 0.065  | 84994  | 0.145  | 19782  | 0.245  | 1288   | 0.365  | 51     |       |  |      |  |
|        |       |        |       |        |       |        |       |        |        | 0.085  | 4970   | 0.185  | 2037   | 0.305  | 5093   |       |  |      |  |
|        |       |        |       |        |       |        |       | MAR15  |        | 0.080  | 119407 | 0.180  | 9606   | 0.295  | 3396   |       |  |      |  |
|        |       |        |       |        |       |        |       |        |        |        |        | 0.100  | 3394   | 0.220  | 9654   |       |  |      |  |
|        |       |        |       |        |       |        |       |        |        | JUN15  |        | 0.095  | 119973 | 0.215  | 254    |       |  |      |  |
|        |       |        |       |        |       |        |       |        |        |        |        |        |        | 0.120  | 100490 |       |  |      |  |
|        |       |        |       |        |       |        |       |        |        |        |        | SEP15  |        | 0.115  | 104904 |       |  |      |  |
|        |       |        |       |        |       |        |       |        |        |        |        |        |        |        |        | DEC15 |  |      |  |

Gas» и сезонность. А вот дневные колебания цены фьючерса и производных от него определяются не так ---они зависят от психологии и стратегий участников рынка.

На рынке торгуют несколько типов трейдеров: индивидуальные трейдеры, роботы (алготрейдинг), крупные proprietary trading houses, глобальные инвестиционные банки (в основном обеспечивают хеджирование). Все они применяют различные тактические приемы, алгоритмы и стратегии. Стратегия обычно состоит из нескольких подготовительных шагов (выставленные ордера), а дальше следует медленный или мгновенный сдвиг текущей ситуации в стакане в направлении, нужном исполнителю стратегии. Трейдеру сложно определить, кто и какую в данный момент применяет стратегию в торговле данным инструментом. А вот нейронная сеть потенциально способна распознать тонкие корреляции между рыночными данными, состоянием стакана и определить, какие стратегии применяются или даже (как мы увидим во второй части) применить свои.

### **Предварительная обработка данных для нейронной сети**

Сформируем набор данных, необходимых для обучения нейронной сети. Предположим, мы пытаемся предсказать движения цены спреда JUN-DEC в секторе Natural Gas commodities.

В этом секторе наблюдается устойчивая сезонность: зимний ноябрь-март (W, сезон высокого потребления) и летний апрель-сентябрь (S, сезон низкого потребления и накопления запасов). Интуитивно ясно и статистически доказано, что «поведение» спреда JUN-DEC (между фьючерсами в летнем и зимнем сезонах) кардинально отличается от поведения спреда JUN-SEP (между двумя фьючерсами в летнем сезоне). Таким образом, мы получаем четыре сезонные группы: WW, WS, SW, SS.

Первая часть данных для рассмотрения — это OHLC, цены релевантных спредов и фьючерсов из группы SW за весь исторический период. Это действительные числа (или Nan если данные отсутствуют), требующие минимальной предварительной обработки (normalization) и удаления Nans (или присваивания им некоторых значений).

Нам также необходимо, чтобы нейронная сеть обучалась на формах графиков изменения цен релевантных спредов и фьючерсов. Нейронная сеть не может принять в качестве входных нейронов графики или последовательности. Остается два варианта: использовать нейронную сеть типа CNN или объединить различные графики в кластеры (по «подобию» формы) и использовать принадлежность кластеру как входной нейрон. Теория гласит, что нужно пробовать оба подхода. Мы же рассмотрим только применение кластеров, так как анализ CNN нейронных сетей не интуитивен по своей природе и уведет нас далеко от финансового характера проблемы.

Итак, мы можем нарезать периоды любой длины из дневного графика изменения цены релевантных инструментов. Назовем их timelet. Встает вопрос, как определить подобия форм timelets в разные периоды торгового дня (как текущего, так и всех исторических). Лобовой подход с использованием Евклидова расстояния (меры) очевидно не срабатывает: мера существенно меняется, когда timelets сдвигаются по времени (рис. 2 с). Существует другой подход к определению меры. Он называется DTW (Dynamic Time Warping) и является одним из алгоритмов для оценки подобия двух временных последовательностей, которые могут быть сдвинуты по фазе и различаться по скорости изменения. Будем использовать DTW расстояния между двумя timelets как меру в смысле подобия по форме.

Следующей задачей, которую нам надо решить, будет объединение близких

по форме timelets в группы (кластеры). В англоязычной литературе это называется cluster analysis (clustering). По сути, это задача группировки множества объектов так, чтобы объекты из одного кластера были больше похожи друг на друга, чем на объекты других кластеров. Эта задача решается применением хорошо известного в ИИ метода K-means clustering, который конструирует кластеры путем объединения объектов вокруг так называемых K centroids (где K — наперед заданное число). Кластеры могут быть построены другими способами (без centroids), но тогда будет необходимо определить centroid для каждого кластера (усреднением или еще как-то) см. рис. 3.

В цели данной статьи не входит детальное рассмотрение метода K-means clustering. Нам достаточно понимать, что существуют методы, которые:

Позволяют распределить наши timelets по кластерам;

В качестве меры близости одного timelet другому могут использовать DTW;

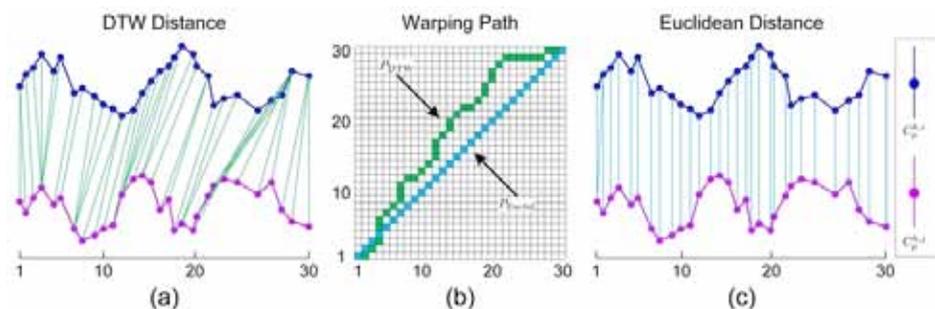
Полученные в процессе формирования кластеров centroids могут служить в качестве усредненных представителей кластеров;

Timelets, которые не участвовали в формировании кластеров, могут сравнивать с centroids кластеров на предмет принадлежности тому или иному кластеру.

Таким образом, мы сформировали начальное множество объектов для обучения нейронной сети: исторические значения OHLC для моделируемого и аналогичных ему активов и набор типовых графиков изменения цен этих активов.

Осталось произвести последнюю операцию над данными — RFE (Recursive Features Elimination). Алгоритм RFE позволяет выбрать наиболее релевантные типы данных и отсеять нерелевантные. Мы можем использовать один из стандартных в

Рис. 2.



ИИ классификаторов (например, SVC) для того, чтобы присвоить веса типам (features) входных данных для нейронной сети. Чем выше вес, тем более релевантным является тип конкретного нейрона. Это делается рекурсивным образом, когда SVC запускается много раз и после каждого запуска из набора данных удаляются наименее релевантные данные. Это интуитивно понятно в нашем случае. Когда мы моделируем спред JUN-DEC, может оказаться, что данные спреда APR-JUN слабо коррелированы с JUN-DEC и классификатор присвоил APR-JUN низкую оценку релевантности. Аналогично спред JUN-NOV может получить высокую оценку релевантности.

Итак, чего мы добиваемся, применяя RFE?

Мы редуцируем количество типов входных данных до оптимального уровня и убираем лишний шум и избыточность из входных данных нейронной сети.

На этом процесс подготовки наборов данных заканчивается. Сформулируем характеристики подготовленных наборов данных для тренировки нейронной сети:

Получен компактный и релевантный набор типов данных (features);

Доступна история «поведения» всех релевантных активов: двухлетние исторические данные OHLC;

Учтена сезонность сектора моделируемого актива. Для каждой сезонной группы SS, SW, WW, WS будет построена отдельная нейронная сеть;

Определены кластеры типовых форм изменения графиков цены активов;

Определен алгоритм отнесения любой части графика изменения цены активов к одному из кластеров.

### Построение и обучение нейронной сети

После того как мы определились с типом нейронной сети и обучающим набором данных, процесс построе-

ния сети становится тривиальным по процедуре и нетривиальным по содержанию.

По процедуре data scientist должен построить внутреннюю структуру нейронной сети. По сути, это практически не формализуемый процесс. Он больше напоминает научное «шаманство» и основан, главным образом, на опыте плюс интуиция. Data scientist пробует большое количество вариантов внутренней структуры нейронной сети, пытаясь добиться максимальных параметров точности и аккуратности. При этом есть опасность «переучить» нейронную сеть (overfitting). Когда она будет показывать блестящие результаты предсказания на обучающем наборе данных, но очень посредственные на новых наборах данных. Но оставим эти проблемы профессиональным data scientists. В результате их работы мы получим нейронную сеть, которая с определенной погрешностью сможет делать некие предсказания о поведении моделируемого актива в обозримом (до конца торгового дня) будущем.

Примеры предсказаний:

в последующие два часа цена спреда JUN-DEC поднимется на 2 тика;

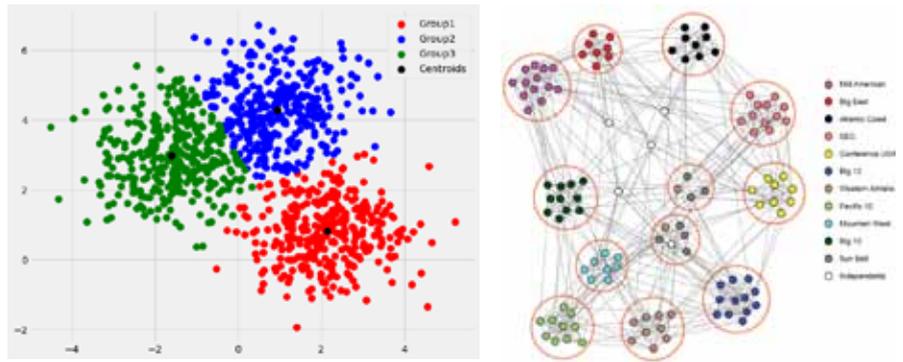
в последующий час график изменения цены спреда JUN-DEC будет меняться по графику centroid кластера 47;

в последующие три часа цена спреда JUN-DEC пойдет вверх.

Вид предсказания зависит от результатов взаимодействия data scientist и практикующего трейдера. Не всегда предсказания, даже исполнившиеся, ведут к профиту. Это очень тонкая вещь — построить прибыльную стратегию торгов на базе предсказаний нейронной сети.

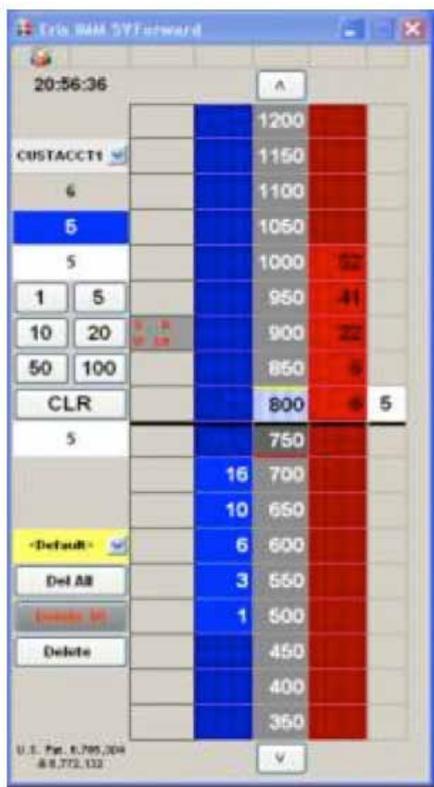
Проиллюстрируем это на примере. Нейронная сеть предсказала, что в ближайшие 2 часа цена спреда XYZ поднимется на 3 тика. Ситуация в стакане показана на рис. 4. Buy-Sell спред (разница между ask и bid) в стакане

Рис. 3.



Слева кластеры с centroids, справа – без centroids.

Рис. 4.



составляет два тика — минимальный ask равен 800, а максимальный bid равен 700. Треjder выставляет ордер на покупку 1 спреда XYZ по 750 и ордер исполняется. Теперь трейдер имеет 1 long позицию в спреде XYZ.

Ожидая, что цена поднимется до 900, трейдер выставляет ордер на продажу 1 спреда XYZ по цене 900. Однако, не все так просто. Правила trading house требуют, чтобы позиция была захеджирована на случай, если рынок пойдет вниз. Хеджирование должно включать в себя два уровня защиты: простой stop sell ордер по цене 650 и hard stop sell ордер по цене 500. Предположим, что в течение первого часа рынок XYZ ушел вниз до 650. Тогда исполнится первый stop sell ордер по цене 650. P&L трейдера будет минус два тика или -100. Если же рынок резко нырнул вниз, то простой stop sell ордер может не сработать и тогда сработает hard stop sell ордер по цене 500. P&L трейдера будет минус пять тиков или -250. Даже если цена спреда, согласно предсказанию нейронной сети, поднимется на втором часу до 900, трейдер УЖЕ потерял деньги.

Рассмотренная выше стратегия не единственная. Можно было попробовать купить call опцион на спред XYZ со strike price ценой 700. При условии, что этот опцион можно купить прямо сейчас. Есть и другие стратегии.

Для эффективного применения первой стратегии нужно будет собирать статистику P&L. Как часто ошибается нейронная сеть, насколько глубоко надо ставить stop sell ордер, какой средневзвешенный P&L при тех или иных параметрах данной стратегии. Это кропотливая совместная домашняя работа трейдера и data scientist.

И последнее. Из психологии известно, что нормальный человек не в состоянии выполнять одновременно более четырех заданий, рабочая память человеческого мозга способна работать

одновременно только в четырех направлениях. Если бы удалось научить ИИ-систему оценивать, мониторить и исполнять трейдинговые стратегии, то это было бы прорывной технологией в цифровизации финансов. Тем более что ИИ система не ограничена четырьмя заданиями. □