

# Reliable ML

Устойчивость ML-моделей  
после выкатки в прод

**03.06.2023 ReliableML @ DataFest 2023**

**Дмитрий Колодзев**

**ООО Промсофт, Новосибирск**

# Про меня

DE Who is Dmitry Kolodezev in the field of Data Science?

 Finished browsing

I'm sorry, but I'm having trouble finding specific information about Dmitry Kolodezev's contributions to the field of data science. The sources that might have this information are not accessible at this moment.

However, based on the information available from his personal website, Dmitry Kolodezev seems to be an active participant in the field of data science. He has given numerous lectures on topics related to machine learning (ML) and artificial intelligence (AI), such as ML model interpretability, ML in business, and ML system design. He also appears to be involved in teaching and conducting workshops on these subjects. He discusses the impact of language models like GPT on the labor market and the importance of regular expressions in Python [1](#) [2](#) [3](#) [4](#) [5](#) [6](#) [7](#).

If you have a specific aspect of Dmitry Kolodezev's work in data science that you're interested in, such as his research, his teaching, his work in industry, or his involvement in data science communities, I may be able to provide more specific information.



KOLODEZEV.RU

# О чем мы тут

- Reliable ML
- Устойчивость моделей
- Почему модель может работать плохо
- Как узнать, что она уже
- Что делать
- Когда уже пора
- Что почитать и анонсы

# Reliable ML

- Название этой секции на Датафесте
- **Канал в телеграмме**
- Подход к разработке моделей, максимизирующий практическую полезность ML для бизнеса  
→ (мы находимся тут) ←
- **Неплохая книга** (не наша)
- Будущая книга (в основном про ML Design Doc)

# Жизнь после деплоя

- В прошивках оборудования
  - Медицинское оборудование, умные камеры
- В сервисах, построенных вокруг ML
  - Рекомендательные сервисы, поиск
- Как одна из частей бизнес-процесса
  - Скоринг, предсказание продаж, оценка сроков

# Оказывается, его нужно кормить

- Ожидания от ML-разработки
  - Капитальные затраты
  - Повысят эффективность работы
- Реальность
  - Капитальные затраты
  - Операционные затраты
  - Менее гибкие процессы
  - Эффект внедрения нестабилен

# Почему оно сломалось

- Могла сразу работать плохо или не работать вовсе
- Проблемы с данными
  - аномалии, выбросы, редкие значения, пропуски.
  - поменялось распределение данных
- Отказы
  - Программные, аппаратные, организационные
- Проблема с моделью
  - нестабильность, **недоопределенность**, низкое качество
- **Атаки** - см **Атлас**
- Неправильное использование — **model card, datasheets**

# Хрупкость ML-систем

- Модели ищут шаблоны в исторических данных
- Шаблоны фиксированы и скрыты от пользователя
- Меняться будут
  - Бизнес-процессы
  - Классификаторы
  - Схема данных
  - Семантика данных
- Модели не приспособливаются к изменениям
- И пользователь их приспособить тоже не может

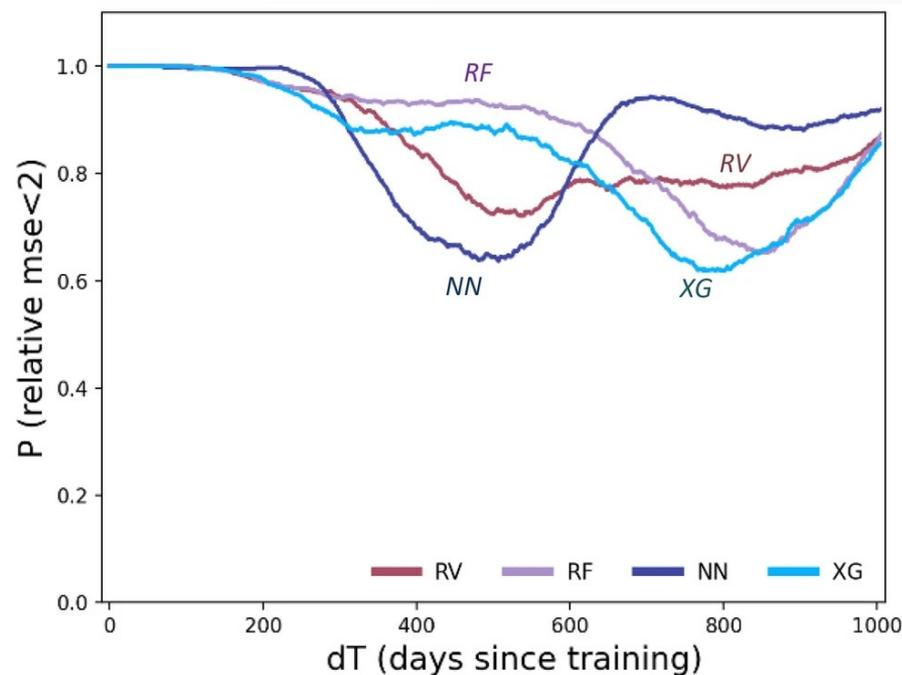
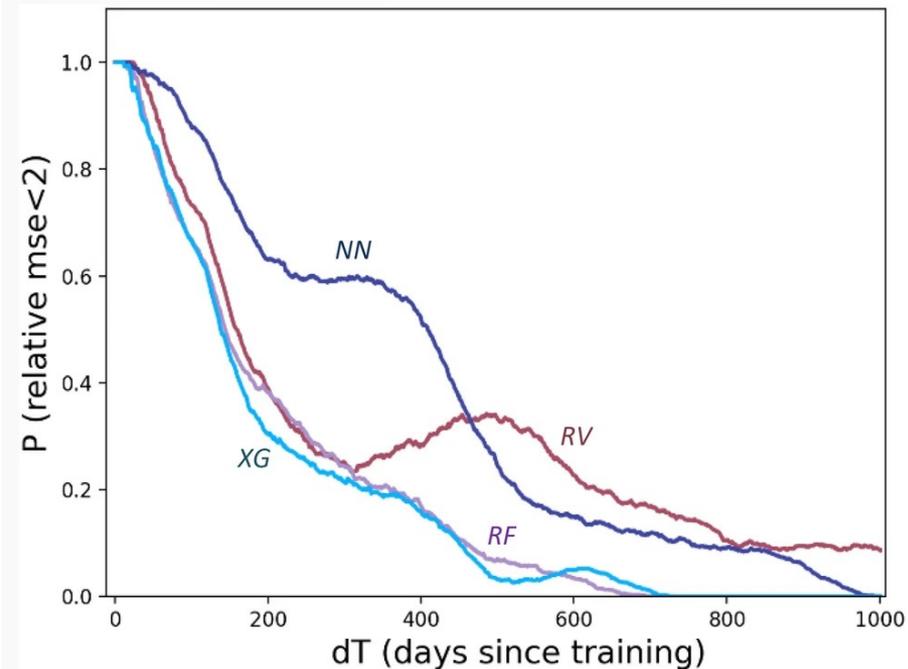
# Гибкость ML-систем

- Настройка порогов и весов
- KNN поверх эмбеддингов
- Регулярные выражения поверх LLM
- Low-Code системы
- Неструктурированные данные
  - Картинки
  - Языковые модели
- Переключение в ручной режим
- Автор в активном поиске антихрупких решений

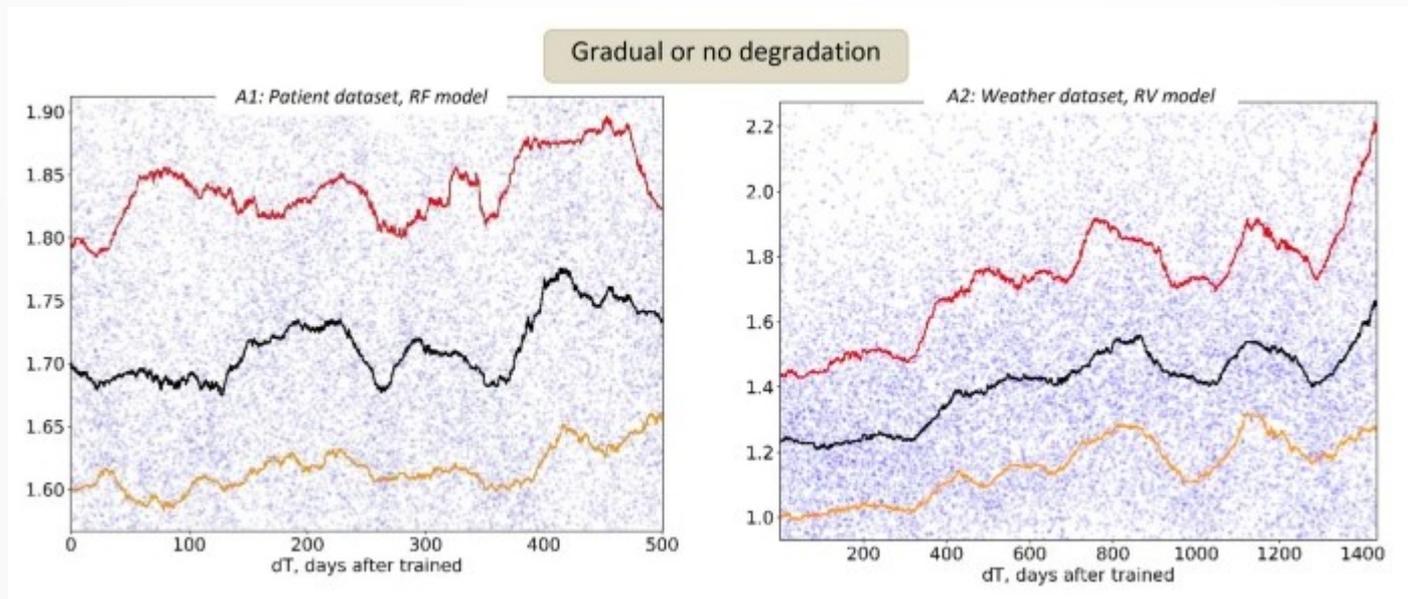
# Как узнаем, что уже?

- Селфчеки т.е. health tests
- Прогноз скорости деградации на исторических данных
- Скорость вызревания разметки
  - Когда мы узнаем, что предсказание было правильным
  - Быстрая: можем позволить деградировать
  - Медленная: не можем позволить деградировать
- Если метки вызревают быстро
  - Отслеживаем качество: SLO, SLI, SLA

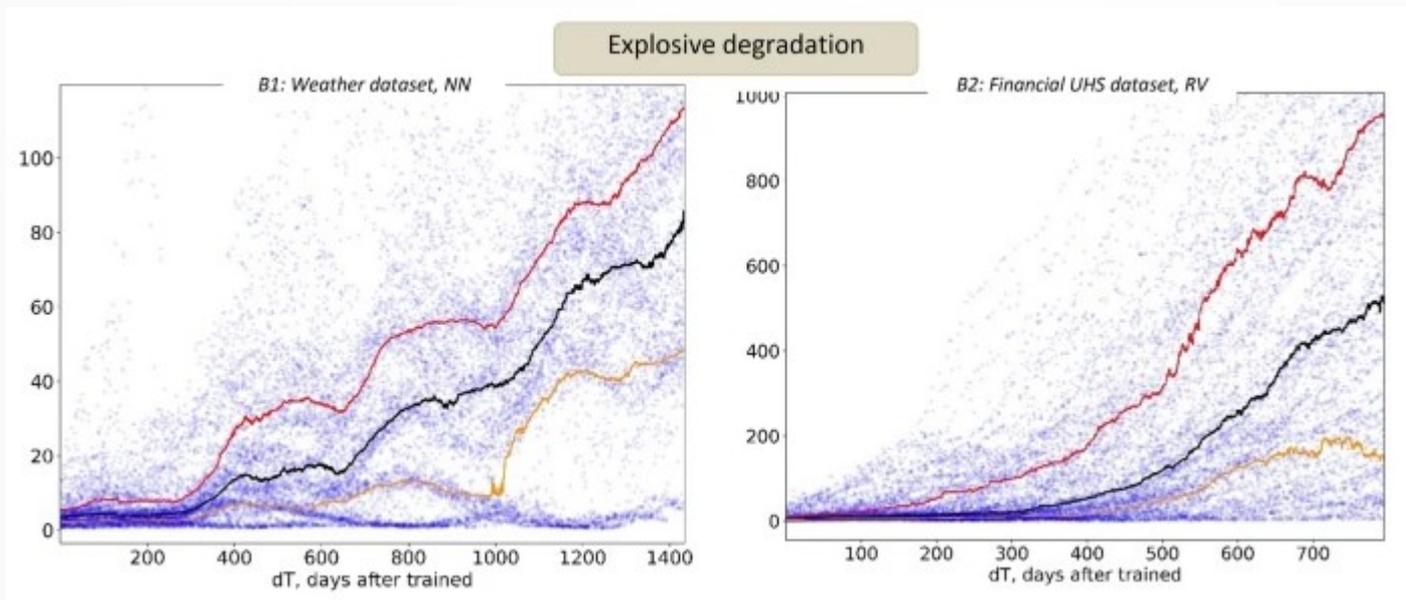
# Деградируем



# Предсказуемая деградация

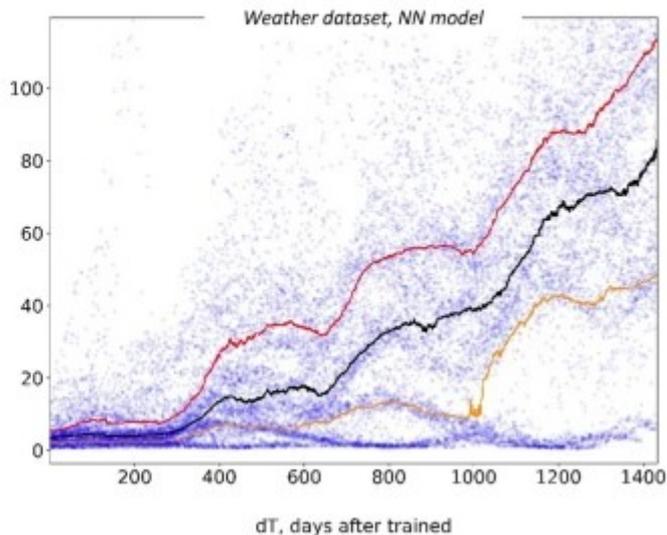
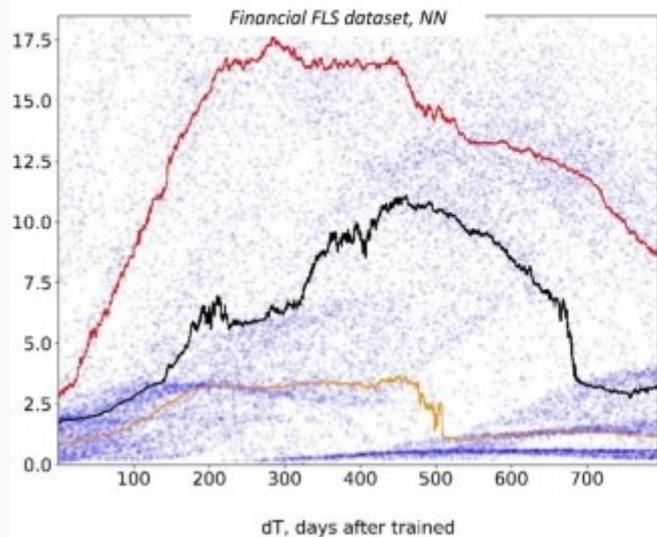


# Взрывная деградация



# Черт-те что

## Strange attractors and chaos



# Медленно вызревающие метки

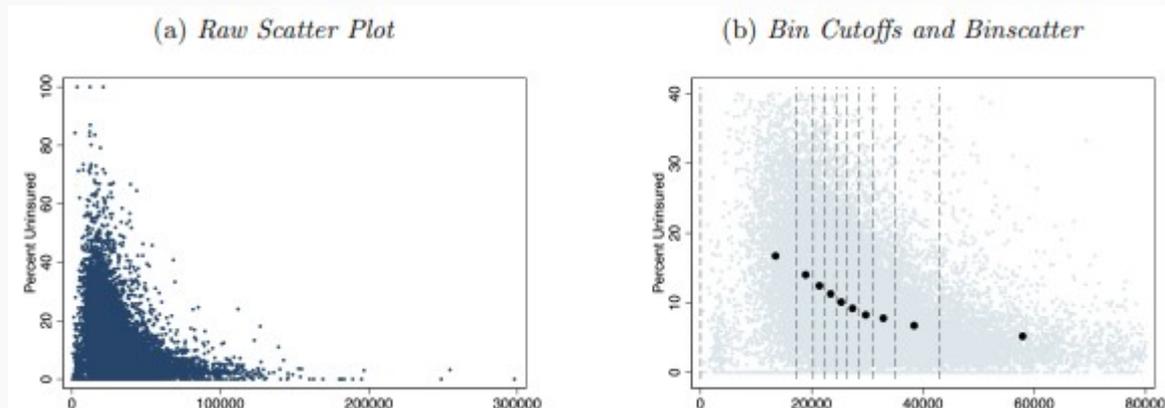
- Прокси-метрики
- Мониторинг распределения предсказаний
- Мониторинг сдвига данных
- Валидация входных данных
- Оценка уверенности модели
- Моделирование невязки
- Слушать пользователей
- Повышать прозрачность решения
- Запасные простые модели

# Интерпретируемость и тд

- Интерпретируем предсказания
- Не должны удивлять экспертов
- Необходимо, но недостаточно — можно накрутить
- Парадокс прозрачности
  - Прозрачные модели надежнее
  - Прозрачные модели проще обманывать
- SHAP и LIME не моделируют причинность
- Посмотрите на контрфактические примеры

# Кстати, EDA

- Идеальная модель выучит шаблоны в данных
- Шаблоны в данных можно посмотреть без модели
- Mean target plot, например
- Мое новое любимое **binsreg** - см **On Binscatter**

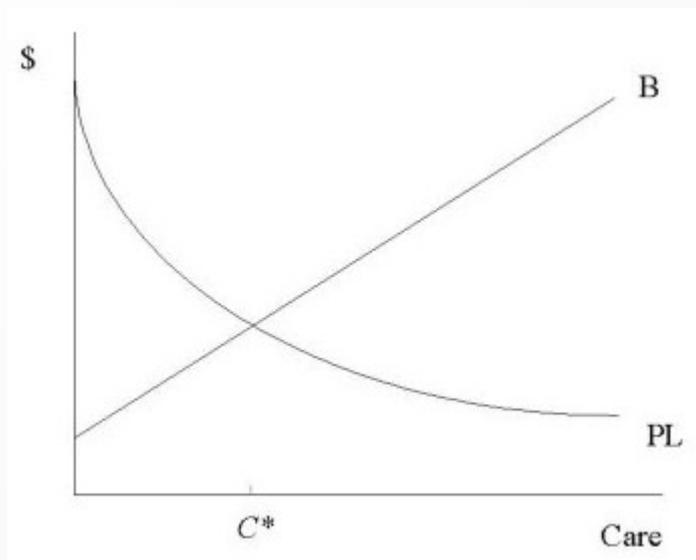


# Предотвратим

- Аудит процесса подготовки данных
  - `df.fillna()` `df.dropna()`
- Детектирование аномалий
  - Клиппирование + флаг
  - Добавлять аномалии специально
- Заполнение пропущенных значений + флаг
  - Выбрасывать признаки специально
- Контролируемая деградация
- Проработанный ML Design Doc — видео, шаблон

# Когда уже пора

- Сколько стоит некачественный предикт?
- Часто не знают
- Сколько стоит мониторить?
- Сколько стоит дообучить?
- Неопределенность
  - Ущерб случаен
  - Успех не гарантирован
  - Например, двойной запас



# Что мешает дообучать модели

- Малый поток данных
  - Синтетика и аугментация
- Медленное вызревание меток
  - Конструируем обратную связь, слабые метки
- Громоздкая подготовка данных
  - Сохранять предикты и признаки
- Трудно сравнивать модели
  - Интерливинг вместо АВ тестов, откат модели

# АНОНС

- Курс по ML System Design, второй запуск осенью
- Переработанный материал
- С ~~блекджеком~~ лабораторными работами
- С промежуточными тестами
- С разделом про устойчивость моделей
- Только онлайн

# Что почитать

- Supervisory Guidance On Model Risk Management
- Machine Learning for High-Risk Applications
- Graceful Degradation and Related Fields
- На DataFest 2023 много докладов про это
- Блог nannyML
- Блог Evidently AI
- Про устойчивость моделей
- Курс ML System Design

# Вопросы

Слайды тут



dkolodezev



promsoft



dmitry\_kolodezev

[https://kolodezev.ru/download/model\\_sustainability\\_2023.pdf](https://kolodezev.ru/download/model_sustainability_2023.pdf)