



Лаборатория анализа данных физики высоких энергий

Томского государственного университета

Физический анализ данных

Томский Государственный Университет

13.10.2025

Направления работы

Повторение анализа ATLAS Open Data

Воспроизведение измерения сечений tf и Z на данных 2015 года (√s = 13 ТэВ)

Цель: валидация цепочки анализа и инфраструктуры.

✓ Верификация методов анализа

Пример: WVZ в pp-столкновениях

Цель: отработка методологии повышения точности.

✓ Алгоритмы машинного обучения

Пример: tW при 13 ТэВ, ATLAS

Цель: сравнение BDT, NN, Transformers, GNN и др. по метрикам качества классификации.

Изучение рождения и дифференциальных распределений мезонов J/ψ, ω и φ

Протон-протонные столкновения, данные LHC (ATLAS)

Готовится публикация

BBCI метод

BCCI (*Bias-Corrected Bootstrap for Confidence Intervals*) — модифицированный бутстрэп-метод, который улучшает оценку доверительных интервалов параметров за счёт коррекции смещения и нестабильности, возникающих при малых выборках.

Применение к результатам анализа WVZ позволяет проверить, насколько устойчивы оценки силы сигнала и систематических параметров по сравнению с классическим подходом.

Цель:

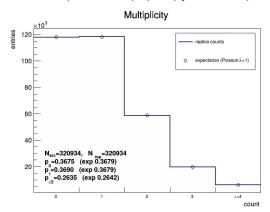
Проверить статистическую устойчивость и надёжность измерения силы сигнала (µ) в условиях ограниченной статистики.

Идея:

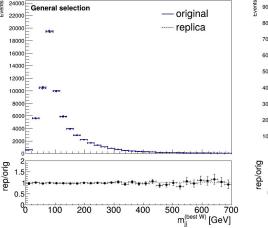
Сгенерировать множество бутстрэп-реплик данных (случайная повторная выборка с заменой) и провести независимую аппроксимацию для каждой из них, чтобы получить распределение оценок µ и его дисперсию.

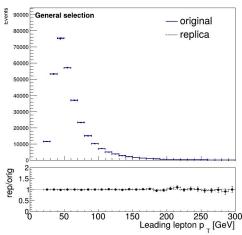
Подготовка реплик

- Сгенерированы 1000 bootstrap-реплик методом случайной повторной выборки с заменой:
 - о для каждой реплики фиксируется seed, что обеспечивает воспроизводимость результатов.
- Каждая реплика сохраняет общее количество событий и все корреляции между переменными:
 - реплики формируются до разбиения на SR/CR/VR.



Распределение числа вхождений событий в репликах согласуется с теоретическим распределением Пуассона (λ = 1). Это подтверждает корректность бутстрэппроцедуры.





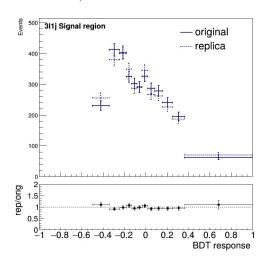
Сравнение форм распределений между оригинальными данными и репликами показывает **согласие** в пределах статистических флуктуаций.

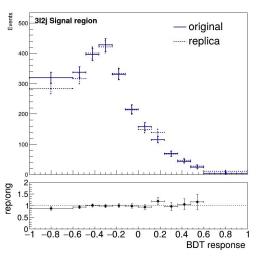
Сигнальный регион

- Рассмотрены события после отбора SR: 322j и 321j.
- Потеря статистики составляет около 1 %, что находится в пределах ожидаемых флуктуаций метода bootstrap.
- Формы распределений **совпадают** с оригиналом в пределах статистических колебаний (*omнowenue replica/original* ≈ 1 ± statistical error).

События:

- **322**j: orig = 2438, repl = 2407
- **3ℓ1j:** orig = 3351, repl = 3311

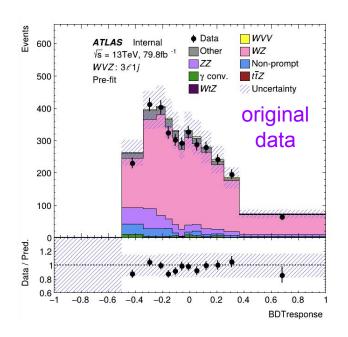


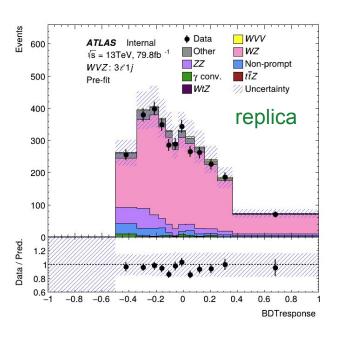


Реплика с seed = 12345 **воспроизводит форму исходных данных**, что подтверждает корректность генерации и сохранение статистических свойств выборки.

Trexfitter

- Обработано 17 реплик по полной цепочке TRExFitter для сигнального региона 321j.
- Показано сравнение оригинальных данных и одной бутстрэп-реплики на уровне pre-fit.
- Наблюдаемые колебания между бинами соответствуют статистическим флуктуациям.
- Реплика воспроизводит структуру сигнала и фоновых компонентов без систематических смещений.

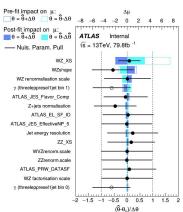




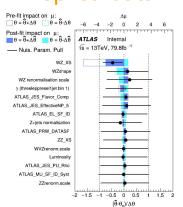
Trexfitter: 311j SR

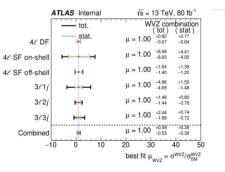
- На данном этапе обработаны реплики для сигнального региона 3€1j.
- Доминирующие неопределённости сохраняются стабильными для всех фитов (оригинал и реплика).
- Следующий шаг: расширение анализа на наиболее чувствительный регион 322j, а затем объединение всех сигнальных областей 321j + 322j + 323j для финальной оценки стабильности метода.

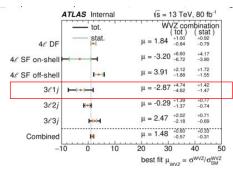
Real data

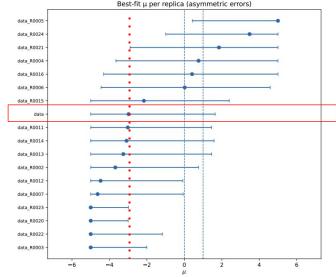


Replica data









17 bootstrap replicas + original data sample.

Сравнение методов ML на НЕР-данных

Цель:

Оценить эффективность современных алгоритмов машинного обучения для задач классификации в физике высоких энергий (HEP). Особое внимание уделено сравнению традиционных методов (BDT, Random Forest) и нейросетевых моделей.

Данные:

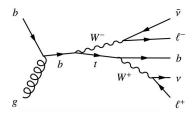
- ATLAS Open Data (2015–2016) открытые данные эксперимента ATLAS.
- HEPData Repository вспомогательные материалы и таблицы с реконструированными событиями.

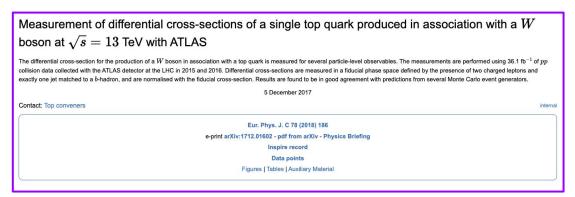
Физический процесс:

Производство топ-кварка в ассоциации с W-бозоном при \sqrt{s} = 13 ТэВ (анализ ATLAS).

Требуется:

- Ровно один b-tag джет,
- И два лептона





Подготовка ntuples

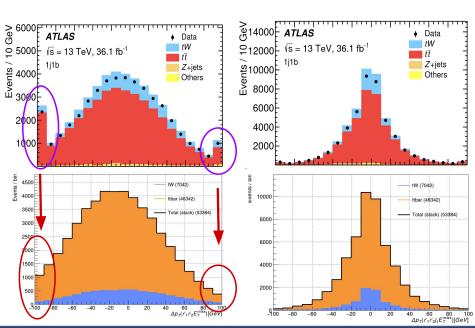
Cutflow:

Все события нормализованы к одинаковой интегральной светимости.

• tW: 7042 событий

tt: 46342 событий

ATLAS Data 5000 vs = 13 TeV, 36.1 fb⁻¹ 1j1b 3000 2000 1000 tW (7042) 1000



Отмечается различие в распределениях переменных, использованных для обучения BDT, в первую очередь в крайних бинах (первом и последнем).

• может означать проблемы с отбором.

События в области сигнала 1j1b

Process	Events 8 300 ± 1 400	
tW		
tī	38400 ± 6600	
Z+jets	620 ± 310	
Diboson	230 ± 58	
Fakes	220 ± 220	
Predicted	47 800 ± 7 300	
Observed	45 273	

Детали тренировки BDT

Цель:

Разделить процессы tW (top+anti-top) и tt в области 1j1b с помощью градиентного BDT.

Используемый инструмент:

Анализ выполнялся в пакете TMVA (ROOT) с использованием метода BDTG — Gradient Boosted Decision Trees.

Основные настройки и особенности:

- Feature scaling: применяется стандартная нормализация z-score, чтобы обучение шло на стандартизированных входных данных.
- Разделение и нормировка: используется случайное разбиение событий TMVA и нормировка по числу событий.
- Оценка модели: для применения обученной модели и построения распределений использовались TMVA Reader + TTreeReader.

Полученное распределение BDT-score для сигнала (красная линия) заметно отличается от представленного в статье:

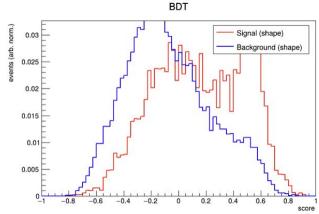
Возможные причины:

- Различие в отборе событий (Selection cuts),
- Использование обновленной реконструкции данных,
- Или различия в весах и нормировках.

Сейчас использовался **85% working point** для b-tagging вместо **77%**, примененного в оригинальной статье:

- 85% WP даёт достаточную статистику и подходит для обучения MLмоделей,
- но для верификации необходимо воспроизвести результаты 77% WP, чтобы убедиться в корректности анализа.

Планируется хранить и анализировать оба набора (77% и 85%) для последующего сравнения эффективности различных ML-алгоритмов.



		score
Variable	Reference Article S (×10 ⁻²)	Produced BDT S (×10 ⁻²)
$pT(\ell\ell E_T^{miss}b)$	4.1	5.38
$\Delta pT(\ell \ell b E_T^{miss})$	2.5	4.19
ΣE_{T}	2.3	0.24
$\eta(\ell\ell E_T^{miss} b)$	1.3	2.85
$\Delta pT(\ell \ell E_T^{miss})$	1.1	0.93
pT(ℓℓb)	1.0	0.62
$C(\ell\ell)$ (centrality)	0.9	0.46
$m(\ell_2,b)$	0.2	0.24
$m(\ell_1,b)$	0.1	0.08

Спасибо за внимание