

1 Объединение и очистка данных

- Интеграция 6 источников → 90,000 строк
- Нормализация форматов, удаление шумовых признаков

2 Feature Engineering 1

- Создание 25 производных признаков
- Out-of-Fold методология с 5-fold кросс-валидацией

3 Feature Engineering 2

- KNN meta-features
- WOE-трансформация для монотонности зависимостей
- Polynomial & Interaction features для нелинейных паттернов

4 Обучение и результаты

- scale_pos_weight решение дисбаланса классов
- Bayesian метод - 30 trials.
- Финальная модель - это ensemble из трёх XGBoost с разными параметрами и оптимизированными весами.
- Добавили Isotonic калибровку для корректных вероятностей

AUC
0.8

Ключевые операции очистки:

Форматирование

- Удаление "\$" и ", " из финансовых полей → конвертация в float
- Нормализация категорий (Full-time, FULL_TIME → Full-Time)
- Преобразование дат в числовые признаки

Обработка пропусков

- JSONL: 89,999 записей → заполнение медианами
- XML: 89,999 записей → интерполяция

Удаление шума

- random_noise_1 – корреляция ~0
- referral_code – 80K+ уникальных значений
- previous_zip_code – низкая предсказательная сила
- Дубликаты: application_id, customer_ref после merge

Feature Engineering 1: Финансовые метрики – 24 признака

1 фин метрики

Колонка	Формула	Описание
calculated_dti	$\text{total_debt_amount} / \text{annual_income}$	Пересчитанный DTI
calculated_pti	$(\text{monthly_payment} + \text{existing_monthly_debt}) / \text{monthly_income}$	Пересчитанный PTI
calculated_utilization	$\text{credit_usage_amount} / \text{total_credit_limit}$	Пересчитанная утилизация
debt_payment_burden	$(\text{monthly_payment} + \text{existing_monthly_debt}) / \text{monthly_income}$	Бремя долга
total_debt_to_income_annual	$\text{total_debt_amount} / \text{annual_income}$	Годовой DTI
total_monthly_debt_payment	$\text{monthly_payment} + \text{existing_monthly_debt}$	Общий месячный платеж
annual_debt_payment	$\text{total_monthly_debt_payment} * 12$	Годовой платеж по долгу
interest_burden	$(\text{interest_rate} / 100) * \text{loan_amount} / \text{monthly_income}$	Бремя процентов

2 ликвидность

Колонка	Формула	Описание
disposable_income	$\text{monthly_income} - \text{total_monthly_debt_payment}$	Располагаемый доход
monthly_free_cash_flow	$\text{monthly_income} - \text{total_monthly_debt_payment} - (\text{loan_amount} / \text{loan_term})$	Свободный денежный поток
free_cash_flow_ratio	$\text{monthly_free_cash_flow} / \text{monthly_income}$	Коэффициент свободного потока
income_to_payment_capacity	$\text{monthly_income} / (\text{monthly_payment} + 0.01)$	Способность платить

4 ДОХОДЫ

Колонка	Формула	Описание
monthly_income_from_annual	$\text{annual_income} / 12$	Месячный доход из годового
income_to_regional_median	$\text{annual_income} / \text{regional_median_income}$	Доход к региональной медиане
income_source_match	$1 \text{ if } \text{abs}(\text{monthly_income} - \text{monthly_income_from_annual}) < 500 \text{ else } 0$	Совпадение источников дохода

3 кредитное поведение

Колонка	Формула/Источник	Описание
account_diversity_index	$\text{num_credit_accounts} / (\text{credit_history_depth} + 1)$	Индекс разнообразия счетов
delinquency_rate	$\text{num_delinquencies_2yrs} / (\text{num_credit_accounts} + 1)$	Уровень просрочек
num_delinquencies_2yrs	Из credit_history	Просрочки за 2 года
num_public_records	Из credit_history	Публичные записи
negative_marks_total	$\text{num_delinquencies_2yrs} + \text{num_public_records}$	Всего негативных отметок
oldest_account_age_months	account_open_years преобразован	Возраст старейшего счета
service_call_intensity	$\text{num_customer_service_calls} / (\text{oldest_account_age_months} + 1)$	Интенсивность звонков

5 региональность

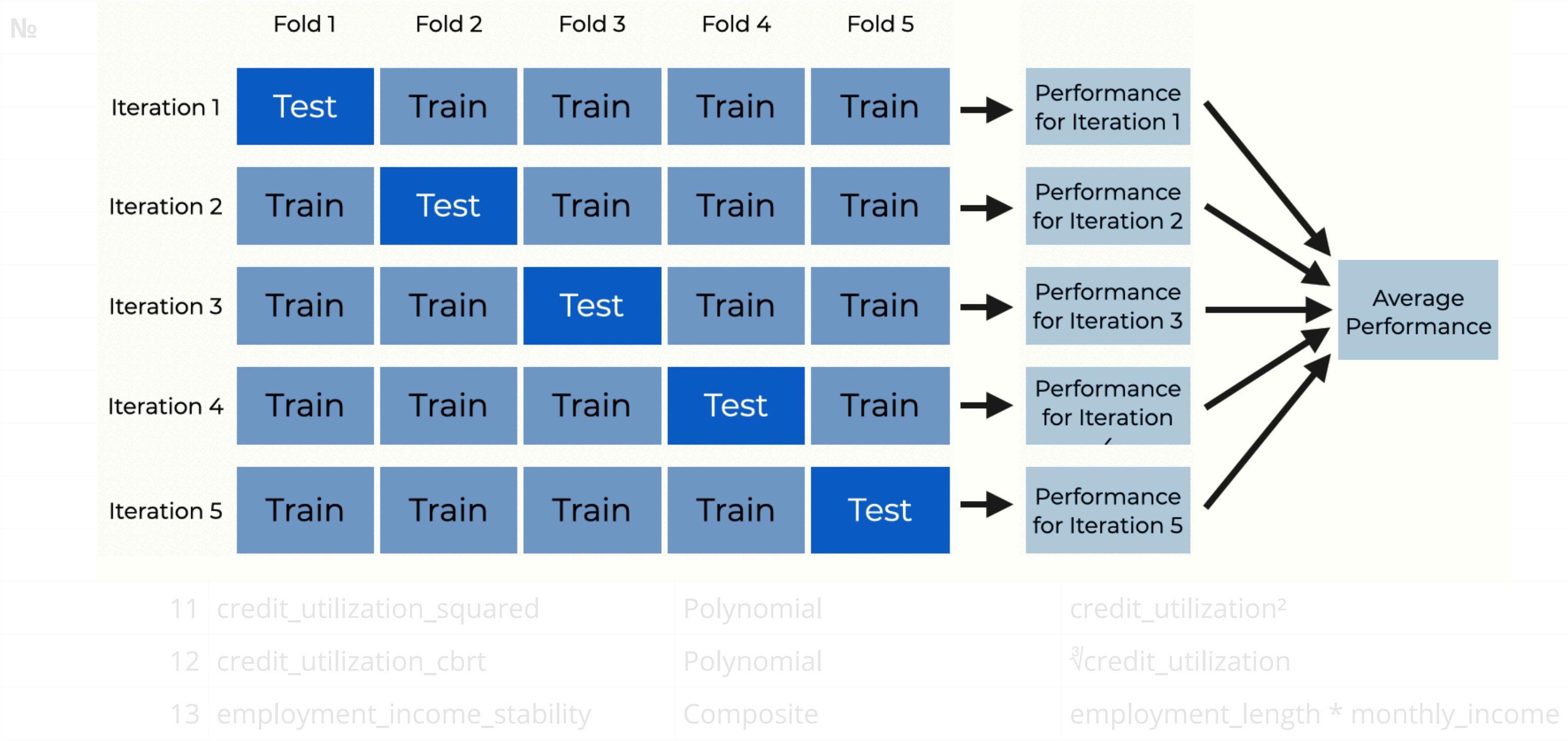
Колонка	Формула	Описание
housing_affordability	$(\text{regional_median_rent} * 12) / \text{annual_income}$	Доступность жилья
regional_stress_index	$\text{regional_unemployment_rate} * \text{cost_of_living_index}$	Региональный стресс



Признаки feature eng 2

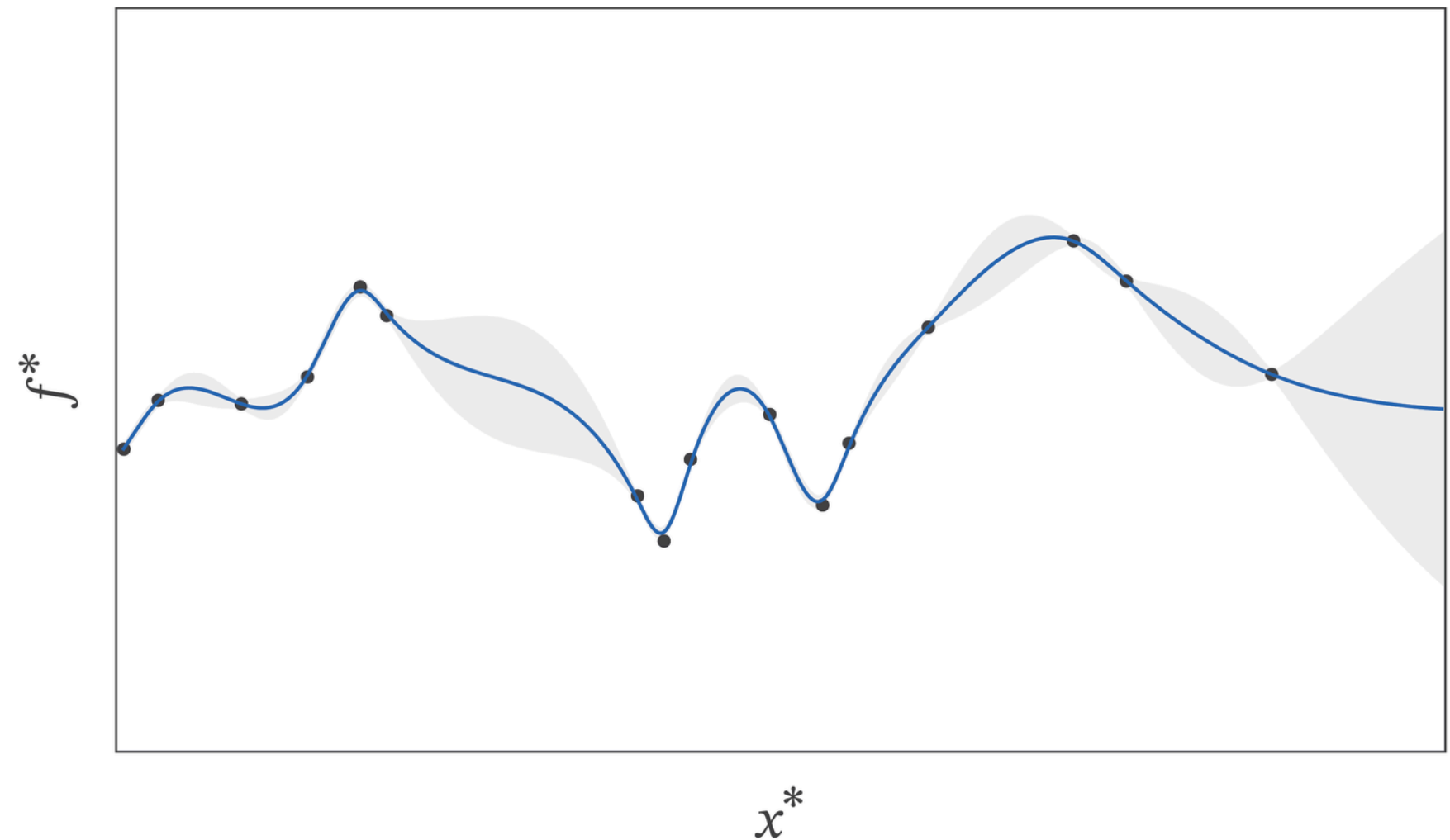
№	Колонка	Категория	Формула/Метод
1	knn_target_prob_50	KNN Meta	KNN(k=50).predict_proba()
2	knn_target_prob_100	KNN Meta	KNN(k=100).predict_proba()
3	knn_target_prob_500	KNN Meta	KNN(k=500).predict_proba()
4	state_target_encoded	Target Encoding	Out-of-fold mean(default) no state
5	debt_credit_interaction	Interaction	debt_to_income_ratio *
6	age_credit_interaction	Interaction	age * credit_score / 100
7	credit_stress_score_squared	Polynomial	credit_stress_score ²
8	credit_stress_score_cbrt	Polynomial	³ √credit_stress_score
9	debt_to_income_ratio_squared	Polynomial	debt_to_income_ratio ²
10	debt_to_income_ratio_cbrt	Polynomial	³ √debt_to_income_ratio
11	credit_utilization_squared	Polynomial	credit_utilization ²
12	credit_utilization_cbrt	Polynomial	³ √credit_utilization
13	employment_income_stability	Composite	employment_length * monthly_income

Признаки feature eng 2



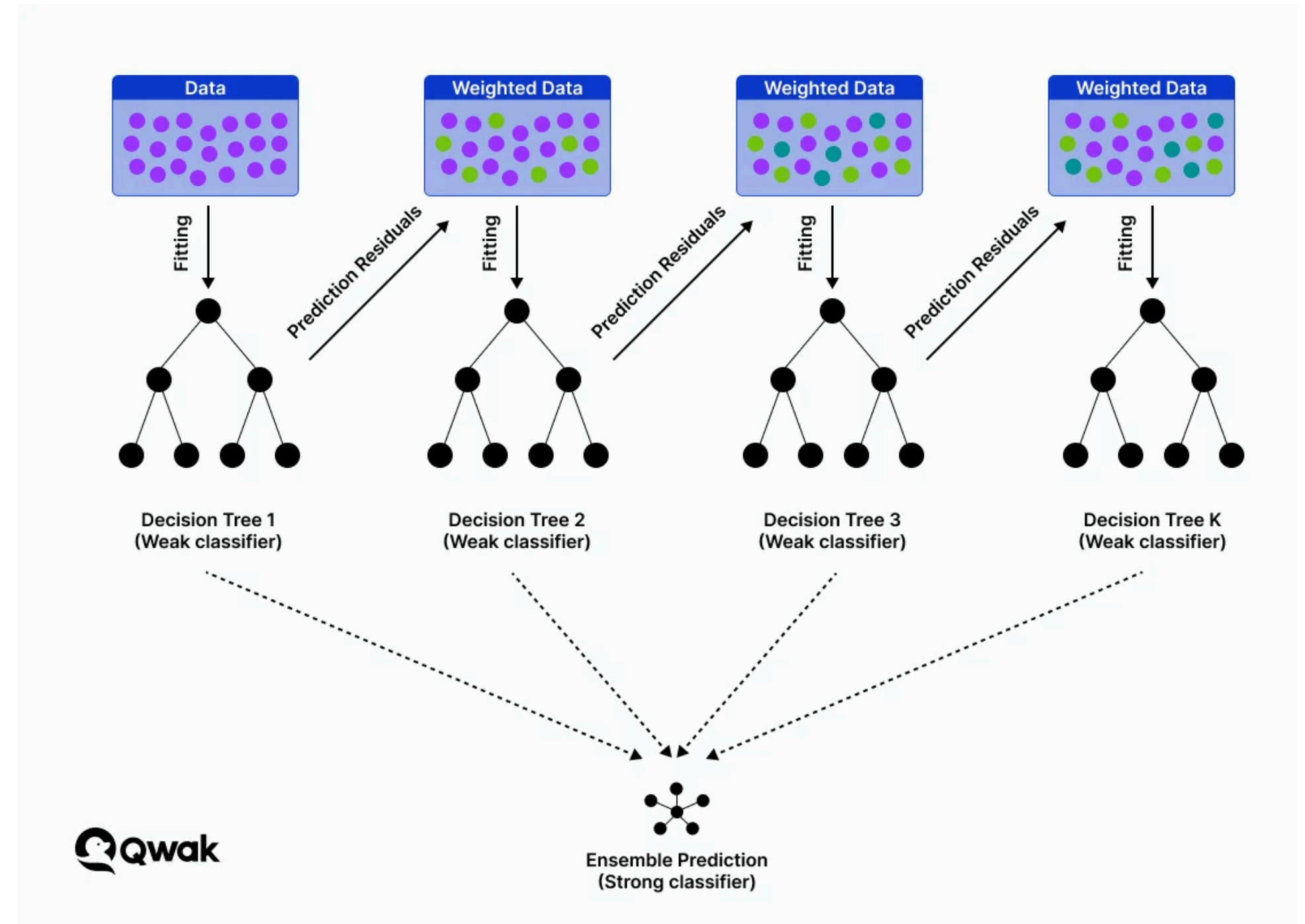
Байесовская оптимизация

это подход к настройке гиперпараметров в моделях машинного обучения, особенно эффективный для ресурсоёмких целевых функций. В отличие от более простых методов, таких как поиск по сетке или случайный поиск, он использует вероятностную модель для управления поиском оптимальных гиперпараметров.



Ensembling method

1. Начальный прогноз:
2. Вычисление остатков (ошибок)
3. Последовательное обучение деревьев: Новое дерево решений обучается на этих остатках (ошибках) с использованием метода оптимизации
4. Обновление прогноза: Прогноз нового дерева добавляется к общему прогнозу ансамбля. Каждое новое дерево вносит небольшой вклад, что контролируется гиперпараметром
5. Повторение: Процесс повторяется, пока не будет достигнут желаемый уровень производительности или максимальное количество деревьев не будет построено.



Isotonic Calibration

isotonic regression or monotonic regression is the technique of fitting a free-form line to a sequence of observations

