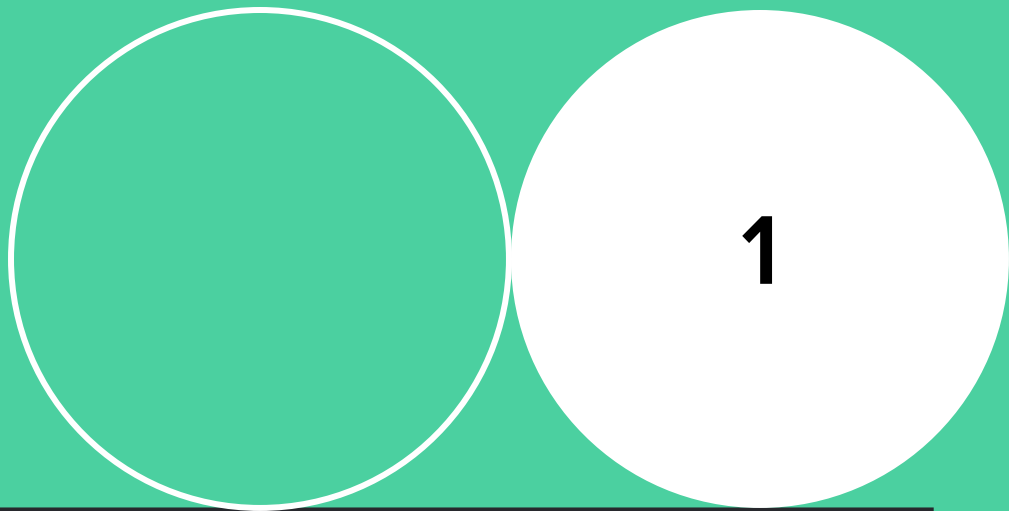


Повышение качества моделей

Повышение качества моделей



Машинное обучение

Улучшение качества модели

1. Работаем с данными.
2. Работа с алгоритмами.
3. Переосмысление проблемы.

Улучшение качества модели

I. Работа с данными.

1. Получить больше данных.
2. Обработка данных.
3. Улучшение качества данных.
4. Придумать больше данных.

II. Работа с алгоритмами.

1. Усложнение при недообучении.
2. Упрощение и/или регуляризация при переобучении.
3. Настройка гиперпараметров
4. Построение ансамбля моделей

III. Переосмыслить проблему

Модель для задачи обучения с учителем

Модель представляем как
функцию с параметрами

где θ - параметры алгоритма
 ε - неустраняемая ошибка

$$y = f(\theta) + \varepsilon$$

$$f(g, w)$$

Параметры модели θ можно разделить на обучаемые w
(просто **параметры**) и необучаемые g (**гиперпараметры**)

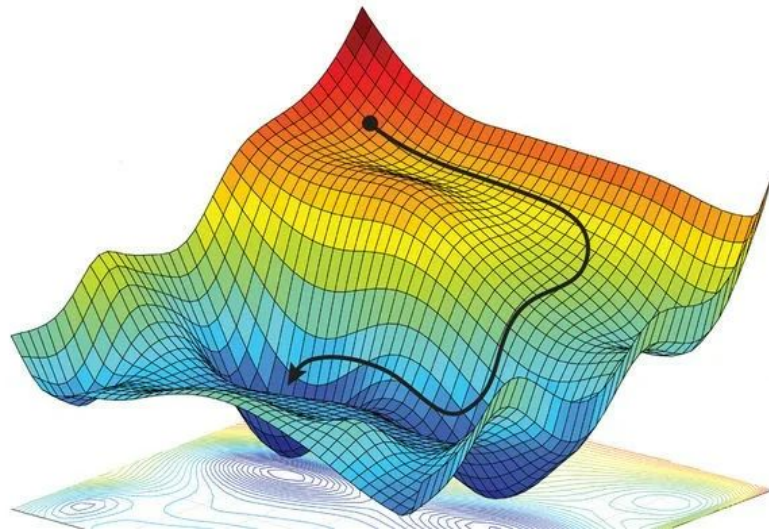
Параметры модели задают **семейство функций**, которые
она может реализовать.

Параметры модели

- Обычные параметры мы получаем в процессе обучения

Подбираем обычные параметры минимизируя функцию потерь L

- Гиперпараметры необходимо подобрать ...



Как подбирать гиперпараметры

1. Ручной выбор - на основе своей экспертной оценки вы выбираете гиперпараметры модели
2. Поиск по сетке - перебор всех вариантов
3. Случайный выбор - проверяются случайные семплы из пространства гиперпараметров
4. Байесовская оптимизация - применяется семплирование из пространства гиперпараметров. На основе результатов обучения модели с наборами гиперпараметров строится вероятностная функция отображения из значений гиперпараметра в целевую функцию, которая позволяет более адресно семплировать из пространства гиперпараметров в области максимума.
5. Генетические алгоритмы - применяются эволюционные алгоритмы с мутациями и скрещиванием
6. Градиентные методы и множество других

GridSearch

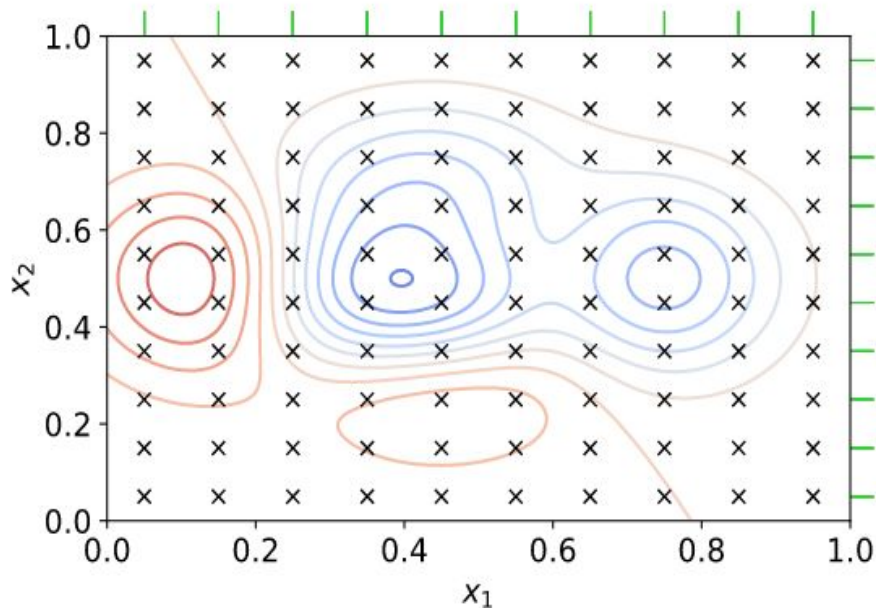
Полный перебор по заданному вручную подмножеству пространства гиперпараметров.

Минусы:

- Долго.
Количество комбинаций $N_1 * N_2 * N_3 * \dots * N_k$,
где N_i - количество возможных значений i -го параметра
- Можно промахнуться мимо минимума

Плюсы:

- Просто и четко
- Можно использовать каскад решеток



RandomizedSearch

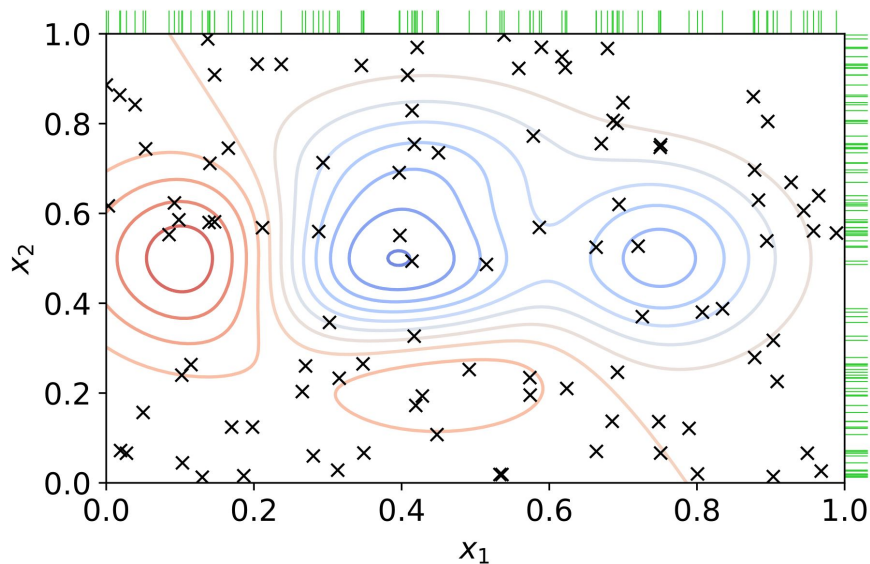
Многократное обучение на случайном сэмплированном наборе гиперпараметров.

Минусы:

- Не гарантируется лучшее решение

Плюсы:

- Достаточно прост
- Обычно быстрее GridSearch находит хорошее решение



Байесовская оптимизация

Основная идея алгоритма – на каждой итерации подбора находится **компромисс** между исследованием **регионов с самыми удачными** из найденных комбинаций гиперпараметров и исследованием **регионов с большой неопределённостью** (где могут находиться ещё более удачные комбинации).

Для этого алгоритм строит **вероятностную модель функции отображения** значений гиперпараметров на **целевую функцию**.

Т.е. **значения гиперпараметров** в текущей итерации *выбираются* с учётом результатов на **предыдущем шаге**.

Это позволяет во многих случаях найти лучшие значения параметров модели за меньшее количество времени.

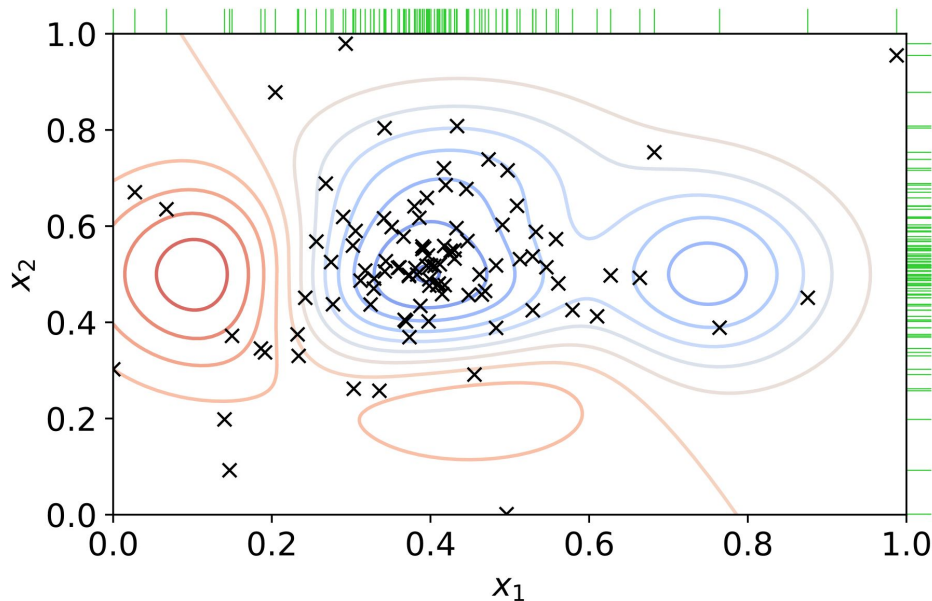
Байесовская оптимизация

Минусы:

- Накладные расходы
- Тяжел для маленьких моделей

Плюсы:

- Хорошо находит оптимум
- Глобальная оптимизация



Генетические алгоритмы оптимизации

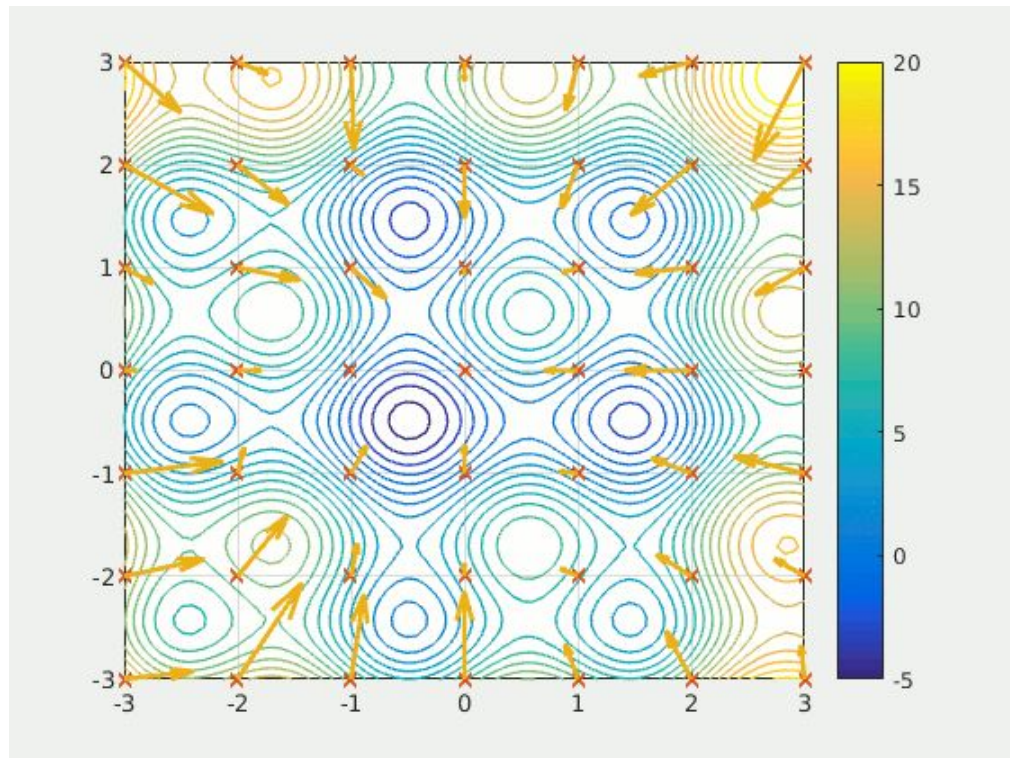
Эволюционная оптимизация гиперпараметров следует процессу, навеянному биологической концепцией эволюции

Минусы:

- Сложно и долго

Плюсы:

- Глобальная оптимизация



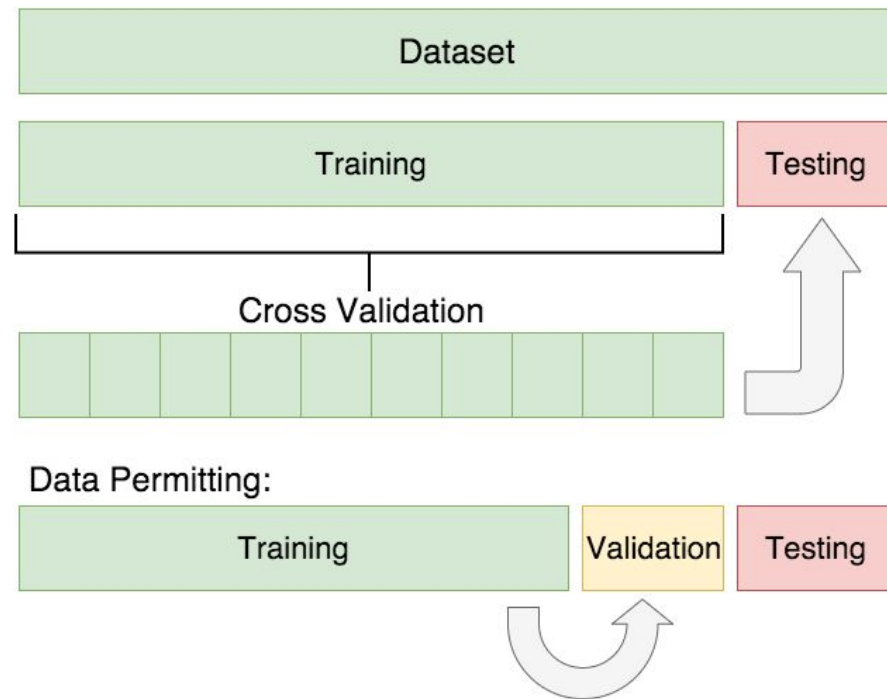
Переобучение

Подбор гиперпараметров для повышения качества можно рассмотреть как метаобучение.

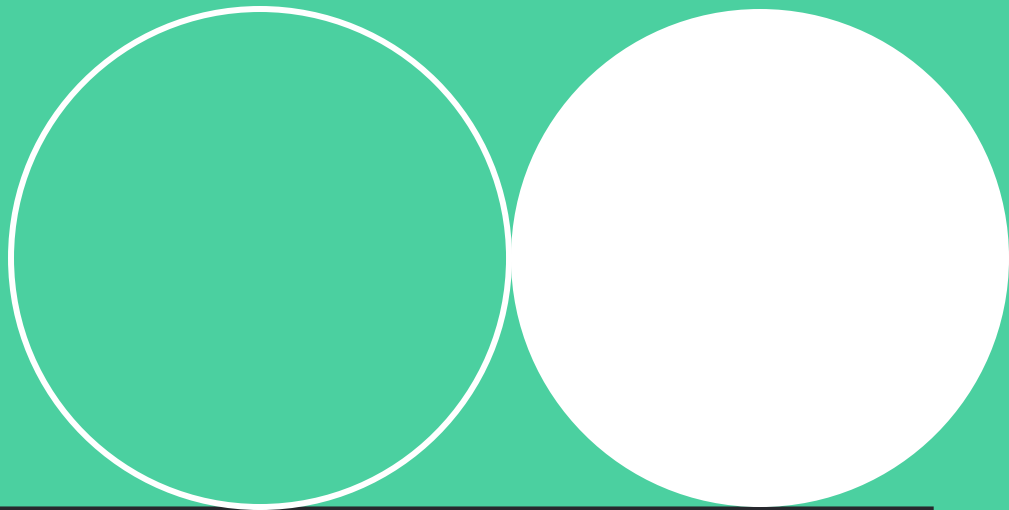
Возможно переобучение на проверочные данные

Решение:

Дополнительный контроль



Практика



Спасибо за внимание!

