

# Байесовская оптимизация

Алексей Зайцев  
Руководитель  
лаборатории  
Сколтех

**Skoltech**

A black and white photograph of a modern, multi-story building with a complex, angular roofline and large windows, identified as the Skoltech building.

# Про что лекция

1. Напоминание про регрессию на основе гауссовских процессов
2. Суррогатное моделирование и Байесовская оптимизация
3. Лучшие практики в Байесовской оптимизации

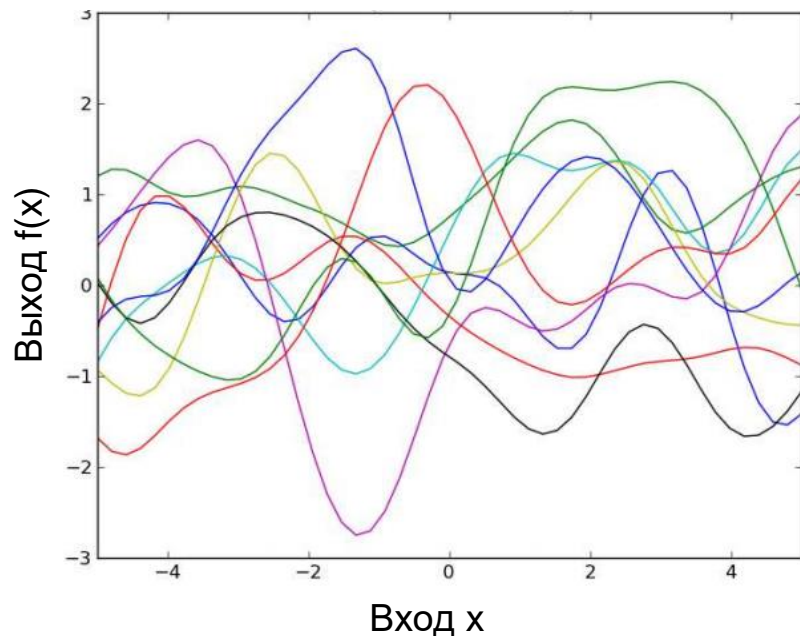
# **Идея регрессии на основе гауссовских процессов**

---

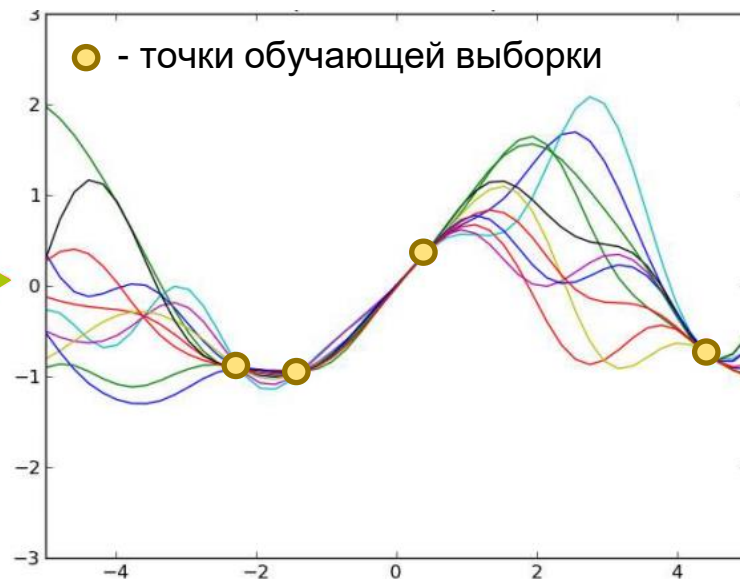


# Идея регрессии на основе гауссовских процессов

Априорное распределение на функциях



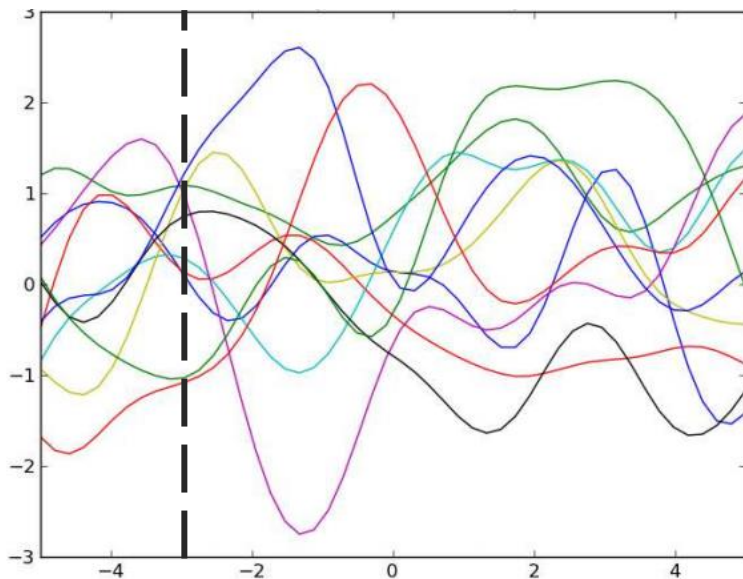
Апостериорное распределение на функциях



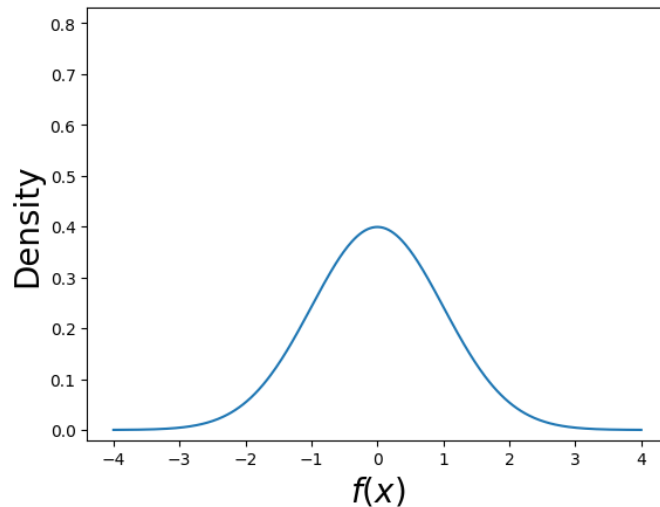
Источник рисунка: <https://stats.stackexchange.com/questions/232959/simulating-the-posterior-of-a-gaussian-process>

# Априорные знания про значения функции

Априорное распределение на функциях

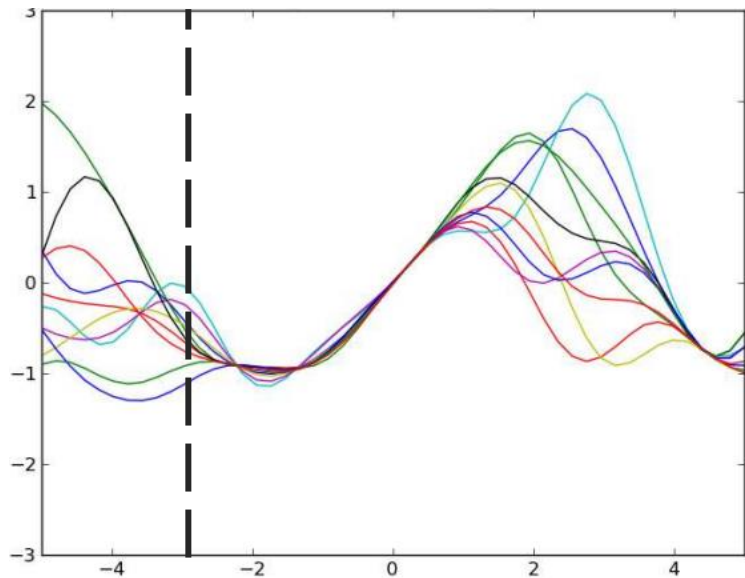


Распределение в точке  $x = -3$

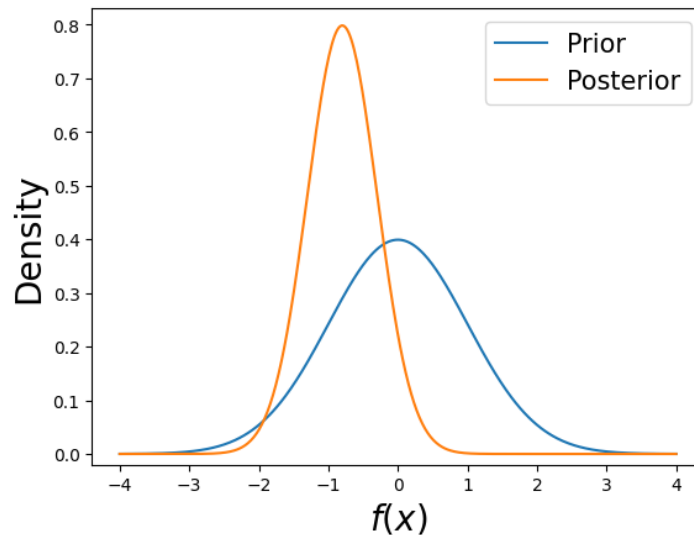


# Апостериорные знания про значения функции

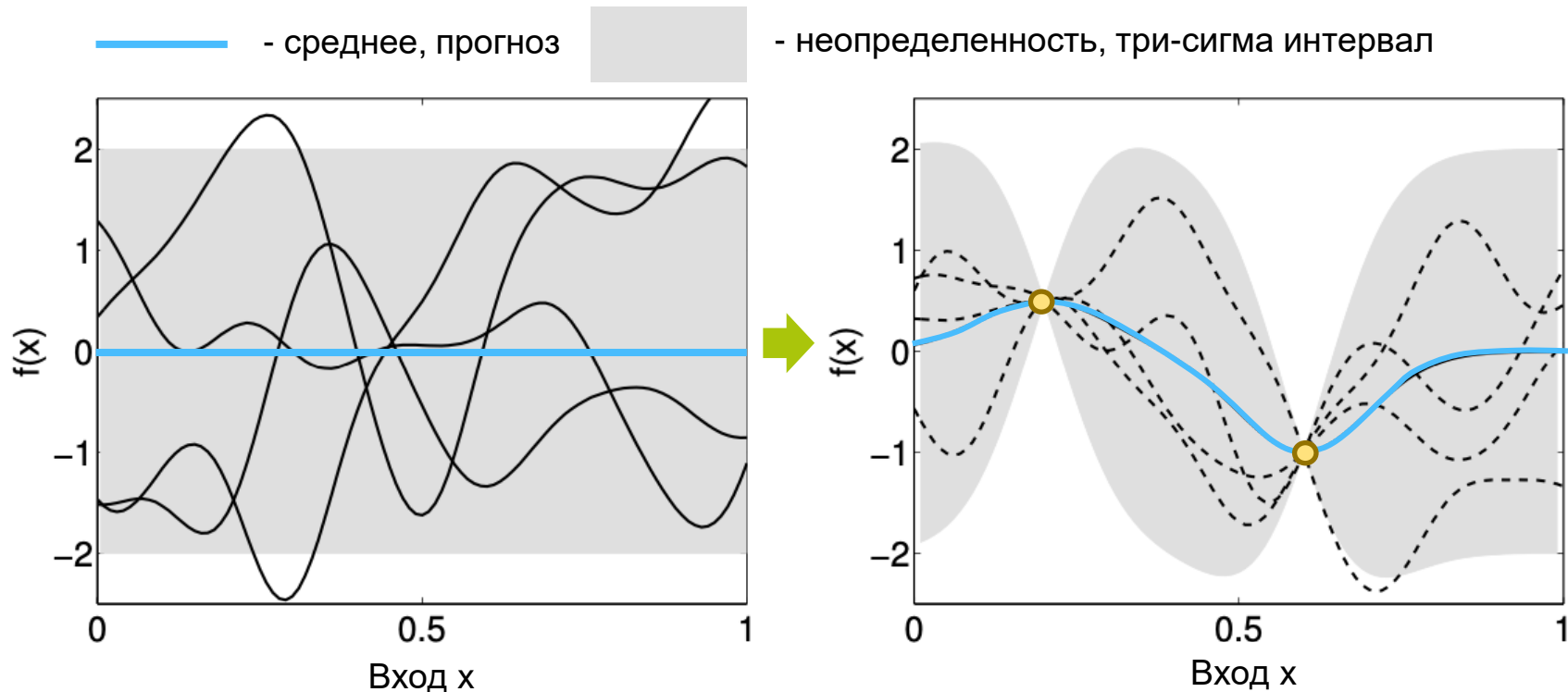
Апостериорное распределение на функциях



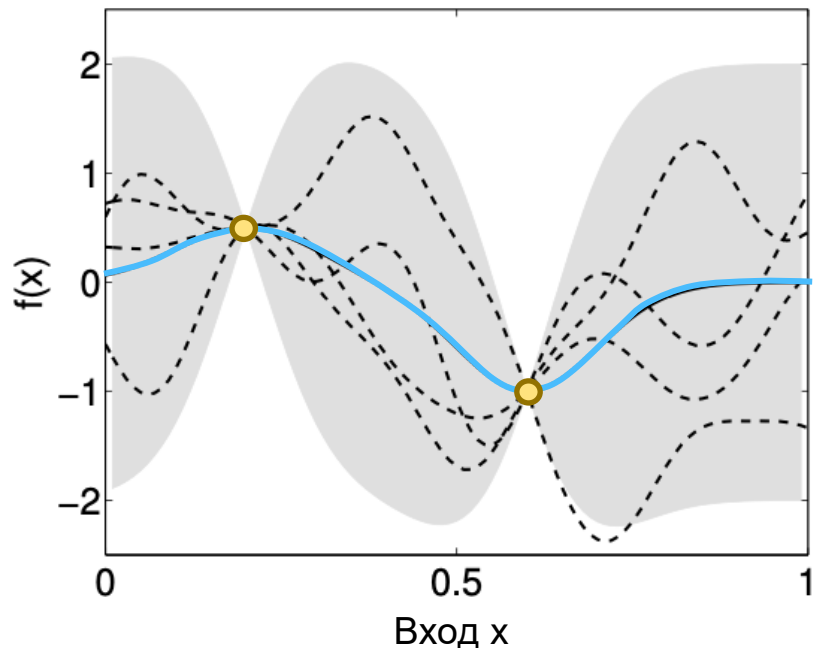
Распределение в точке  $x=-3$



# Усредним по всем функциям



# Свойства регрессии на основе гауссовских процессов



- Нелинейные прогноз
- Интерполяция
- Гладкость (зависит от ковариационной функции)
- Оценка неопределенности
- Выбор параметров с помощью метода максимума правдоподобия
- Явные и быстрые формулы для прогноза



# Оценка среднего и дисперсии для регрессии на основе гауссовских процессов

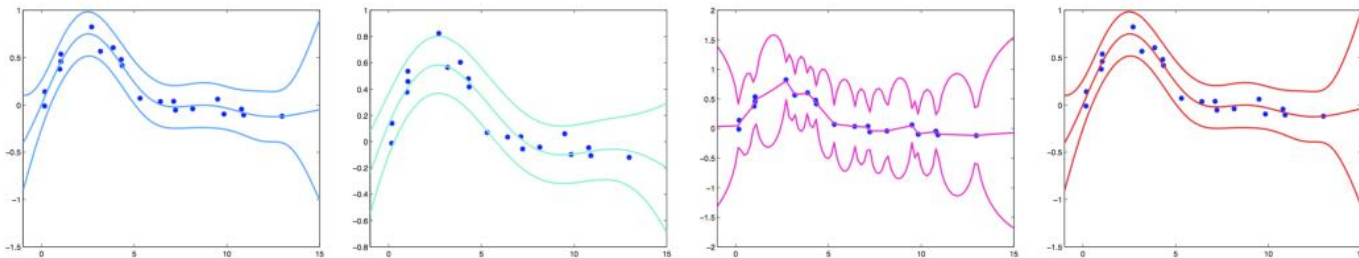
Аналитические формулы для среднего и дисперсии:

$$p(f_*|y) = \mathcal{N}(f_*|\mu_*, \sigma_*^2),$$

$$\mu_* = \mathbf{k}_*^T [\mathbf{K} + \sigma^2 \mathbf{I}_m]^{-1} \mathbf{y}, \quad \sigma_*^2 = K_{**} - \mathbf{k}_*^T [\mathbf{K} + \sigma^2 \mathbf{I}_m]^{-1} \mathbf{k}_*$$

$$\mathbf{K} = \{K(x_i, x_j)\}_{i,j=1}^n \quad \mathbf{k}_* = \{K(\mathbf{x}_*, \mathbf{x}_i)\}_{i=1}^m \quad \text{and} \quad K_{**} = K(\mathbf{x}_*, \mathbf{x}_*)$$

$$\mu_* = \sum_{i=1}^m \alpha_i K(\mathbf{x}_*, \mathbf{x}_i), \quad \alpha = [\mathbf{K} + \sigma^2 \mathbf{I}_m]^{-1} \mathbf{y}$$



# **Pipeline Байесовской оптимизации**

---



# Какую функцию хотим оптимизировать?

Минимизируем тяжелую функцию в многомерном пространстве. Можем сделать ограниченное количество  $B$  запусков.

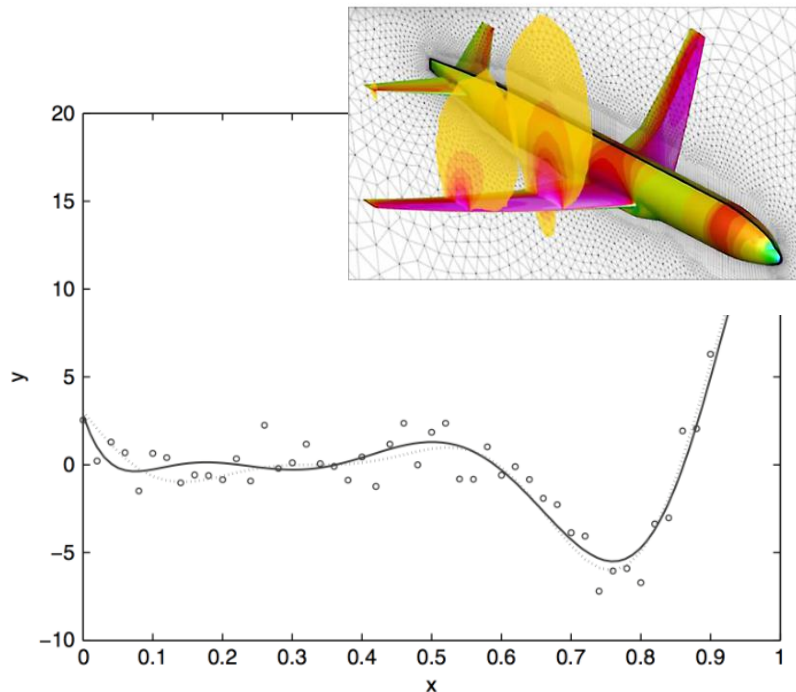
$$f(x) \rightarrow \min_{x \in X}$$

Функция обычно:

- Тяжелая
- Гладкая
- Но наблюдаем мы ее с шумом

Наша задача:

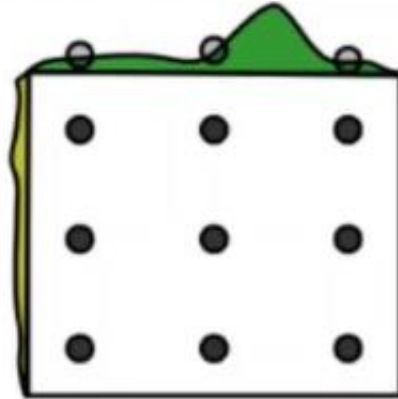
- Оптимизация гиперпараметров модели



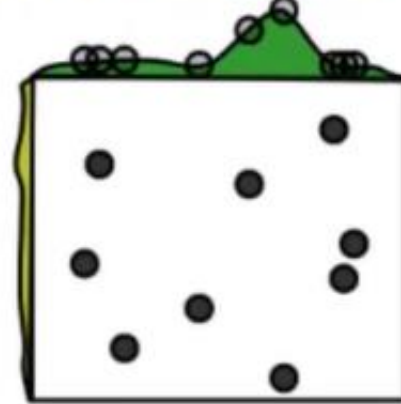
AlphaGo

# Базовые подходы

Grid Search



Random Search



- Не принимают во внимание прошлые вычисления
- В общем случае нужно экспоненциальное по размерности количество вычислений целевой функции

# Алгоритм: как такую оптимизацию устроить

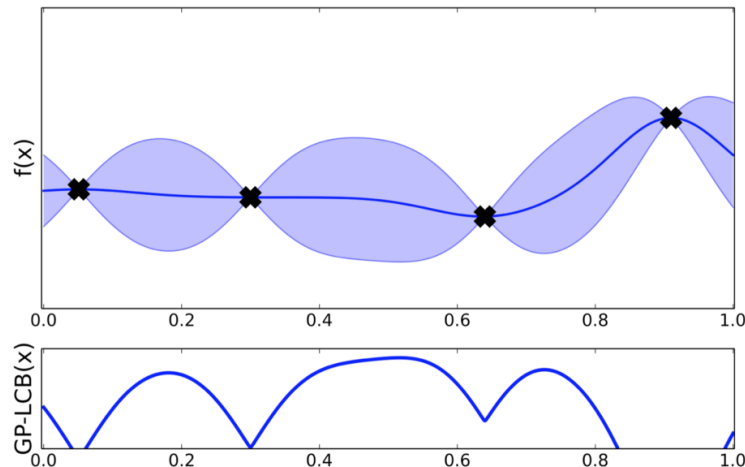
1. Генерируем начальную выборку
2. Учим суррогатную регрессионную модель на этой выборке
3. Оптимизируем суррогатную модель + ее неопределенность вместо исходной тяжелой модели, получили еще одну точку
4. Переучили суррогатную модель
5. Повторяем пока не закончится бюджет

*Проблемы нет: теперь работает нормально*

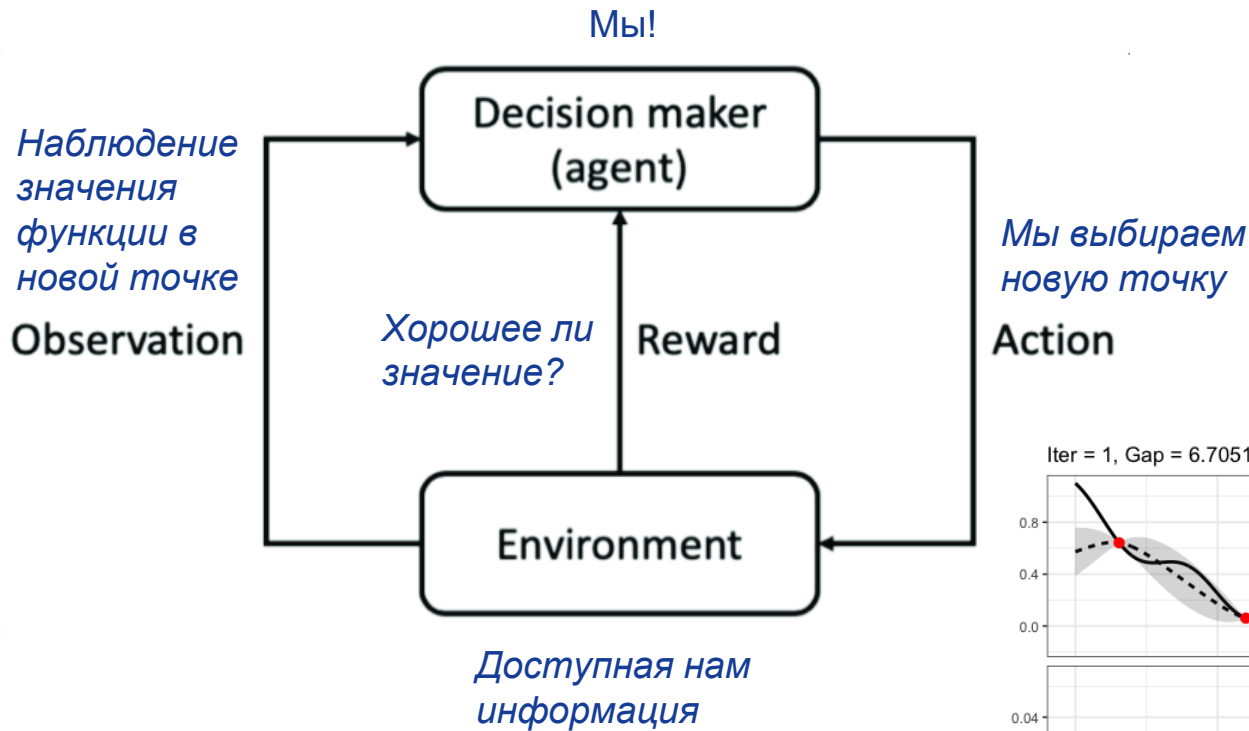
Прогноз  
Оценка неопределенности

$$\alpha_{LCB}(\mathbf{x}) = -\mu_*(\mathbf{x}) + \zeta \cdot \sigma_*(\mathbf{x})$$

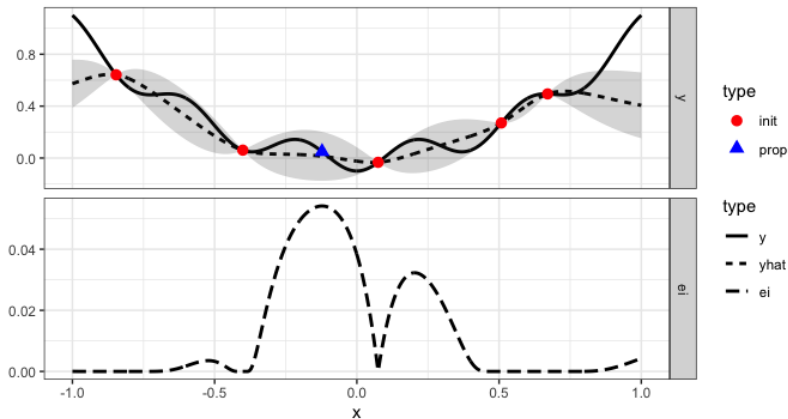
Максимизируем  $\alpha_{LCB}(\mathbf{x})$



# Почти обучение с подкреплением



Iter = 1, Gap = 6.7051e-02



# Acquisition function, функция выгоды LCB

Прогноз                      Оценка  
   неопределенности

↓                                      ↓

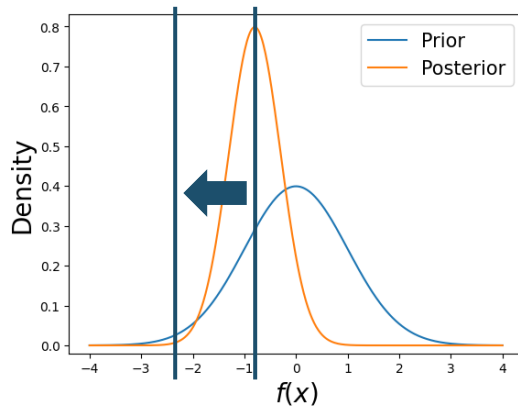
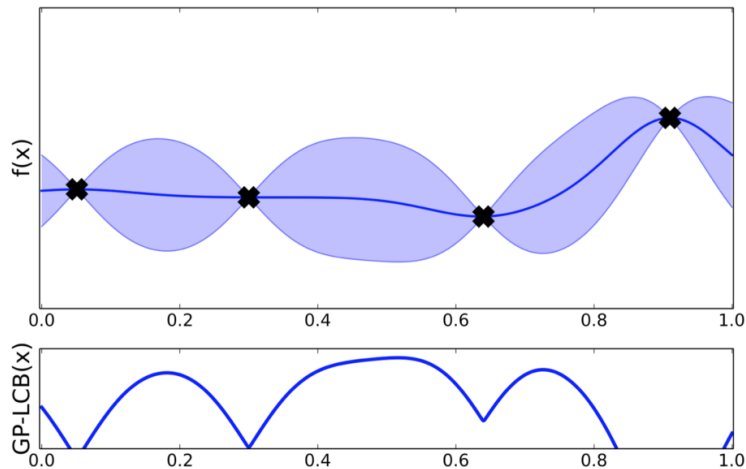
$$\alpha_{LCB}(\mathbf{x}) = -\mu_*(\mathbf{x}) + \zeta \cdot \sigma_*(\mathbf{x})$$

Максимизируем Lower Confidence Bound

$$\alpha_{LCB}(\mathbf{x})$$

Если максимизируем, то смотрим  
на Upper Confidence Bound

**Теорема:** мы достаточно быстро сойдемся к минимуму,  
 $\zeta$  должна расти как корень из количества итераций.

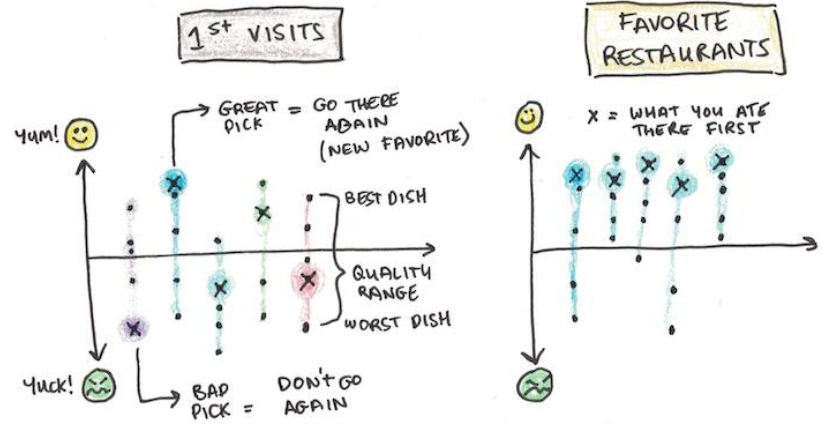


# Функция выгоды балансирует exploration и exploitation

Exploitation      Exploration

↓                      ↓

$$\alpha_{LCB}(x) = -\mu_*(x) + \zeta \cdot \sigma_*(x)$$



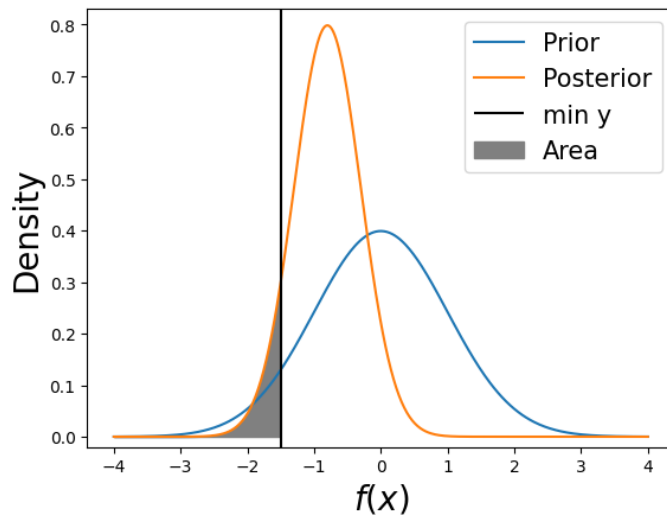


# Функция выгоды Expected Improvement

$$\Delta(\mathbf{x}) = y_{\text{best}} - \mu_*(\mathbf{x}), \quad y_{\text{best}} = \min_{i=1, \dots, m} y_i,$$

$$\begin{aligned} \alpha_{EI}(\mathbf{x}) &= \int \max(0, y_{\text{best}} - y_*) p(y_* | \mathbf{x}) dy_* = & \text{= } \mathbb{E} \max(0, y_{\text{best}} - y_*) \\ &= \Delta(\mathbf{x}) \Phi(-\Delta(\mathbf{x}) / \sigma_*(\mathbf{x})) + \sigma_*(\mathbf{x}) \varphi(\Delta(\mathbf{x}) / \sigma_*(\mathbf{x})) \end{aligned}$$

Мы считаем математическое ожидание улучшения



# Функция выгоды Tree Parson Estimator TPE

Пусть у нас есть две генеративных модели

$$P(x \mid \text{bad}) = P(x \mid y > y_{\min})$$

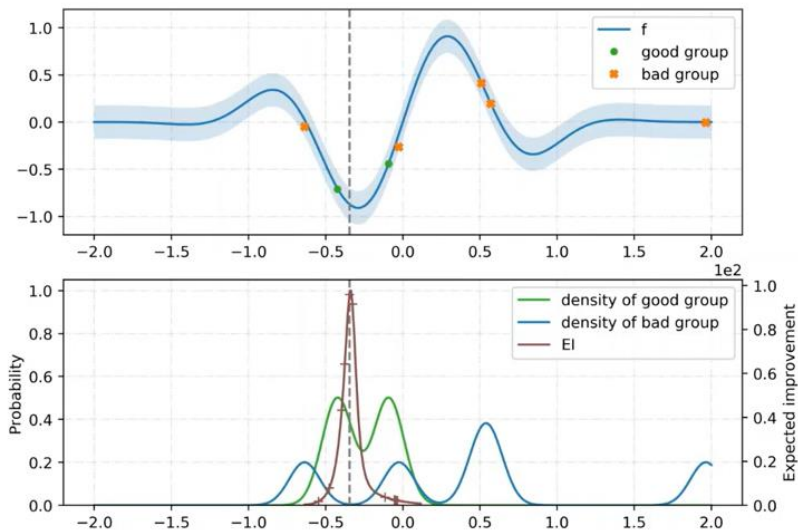
$$P(x \mid \text{good}) = P(x \mid y \leq y_{\min})$$

Перспективность точки:

$$\frac{P(x \mid \text{good})}{P(x \mid \text{bad})}$$

Это пропорционально Expected improvement

Такой метод используется в HyperOpt и Optuna



# Функция выгоды Tree Parson Estimator TPE

Пусть у нас есть две генеративных модели

$$P(x \mid \text{bad}) = P(x \mid y > y_{\min})$$

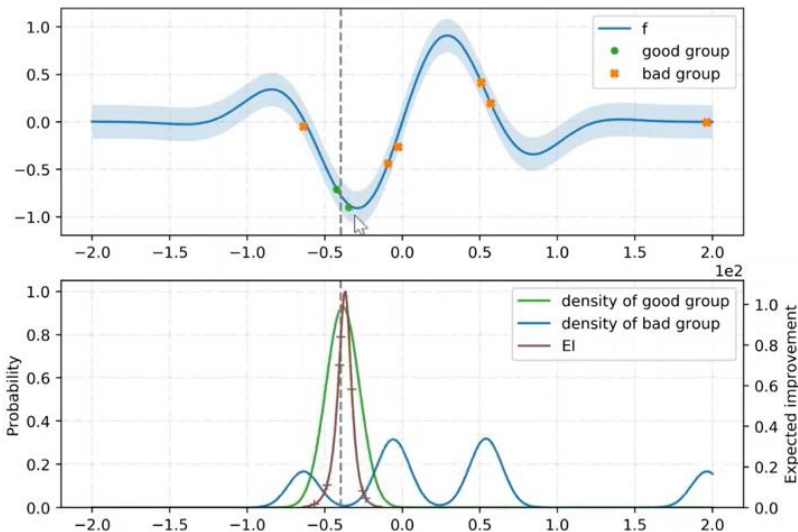
$$P(x \mid \text{good}) = P(x \mid y \leq y_{\min})$$

Перспективность точки:

$$\frac{P(x \mid \text{good})}{P(x \mid \text{bad})}$$

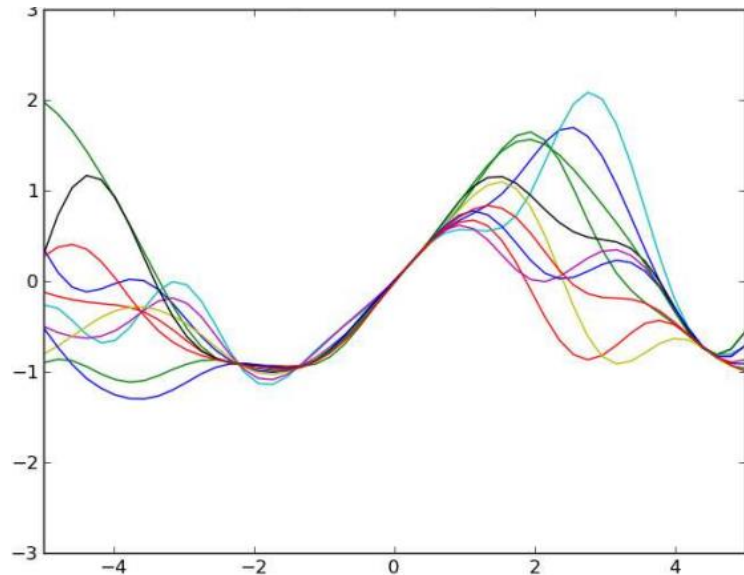
Это пропорционально Expected improvement

Такой метод используется в HyperOpt и Optuna



# Сэмплирование Томпсона

Апостериорное распределение на функциях



На каждой итерации:

- Сэмплируем функцию
- Берем ее минимум

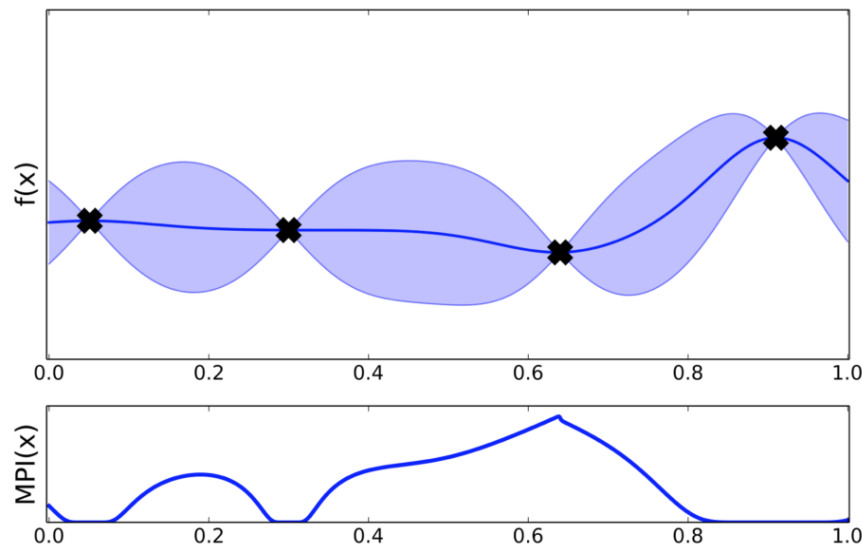
На практике медленней, но не застревает в локальных оптимумах

# Функция выгоды Прирост энтропии

$$\alpha_{ES}(\mathbf{x}) = \mathcal{H}[p(\mathbf{x}_{\min}|\mathbf{y})] - \mathbb{E}_{p(\mathbf{y}|\mathbf{y},\mathbf{x})}[\mathcal{H}[p(\mathbf{x}_{\min}|\mathbf{y} \cup \{\mathbf{x}, \mathbf{y}\})]]$$

Мы считаем как улучшается наше знание про локацию минимума

- Явно посчитать уже не получается, нужно использовать методы Монте-Карло для приближенного вычисления интегралов
- Достаточно надежный метод
- Аналогично можно считать прирост знания не про точку минимума, а про значение в этой точке



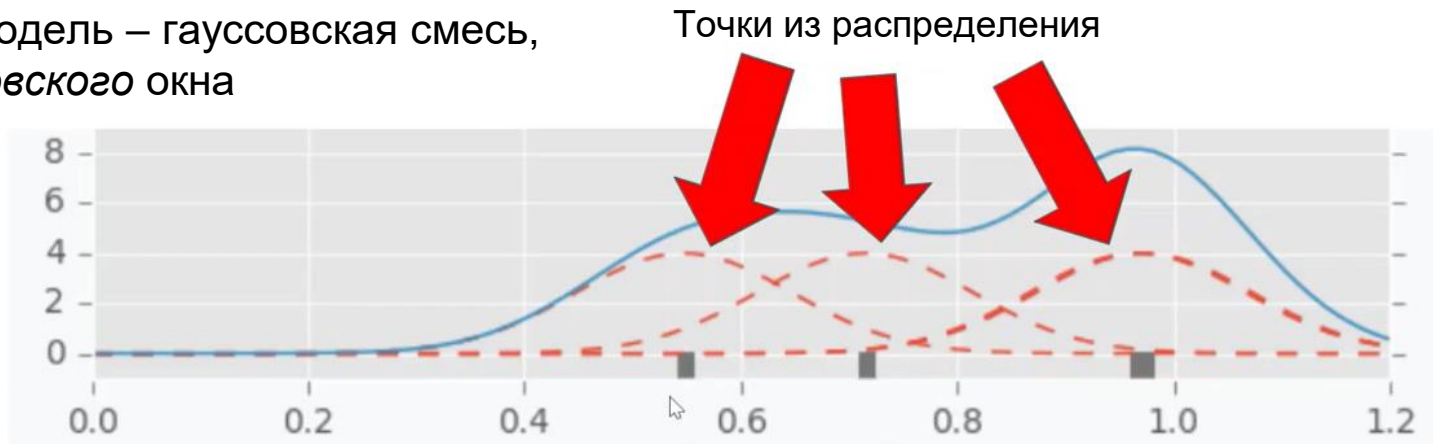
Пример значений энтропийного критерия

# Детали оптимизации

Используем

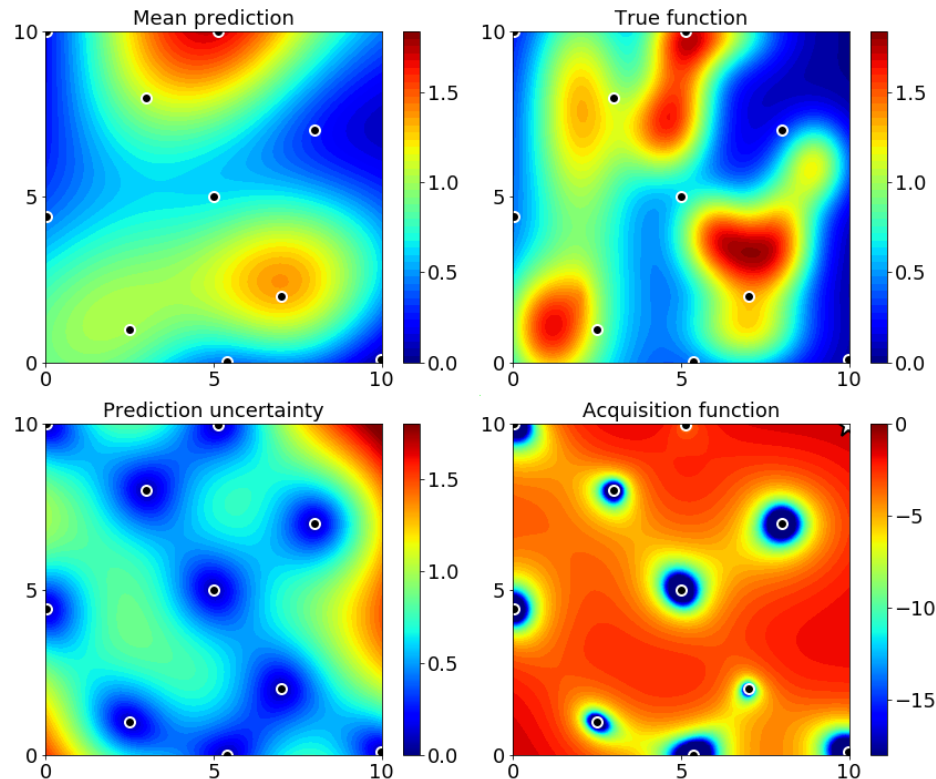
- усеченные (truncated) нормальные распределения для непрерывных гиперпараметров
- категориальные распределения для дискретных

Байесовская модель – гауссовская смесь,  
метод *Парзеновского* окна



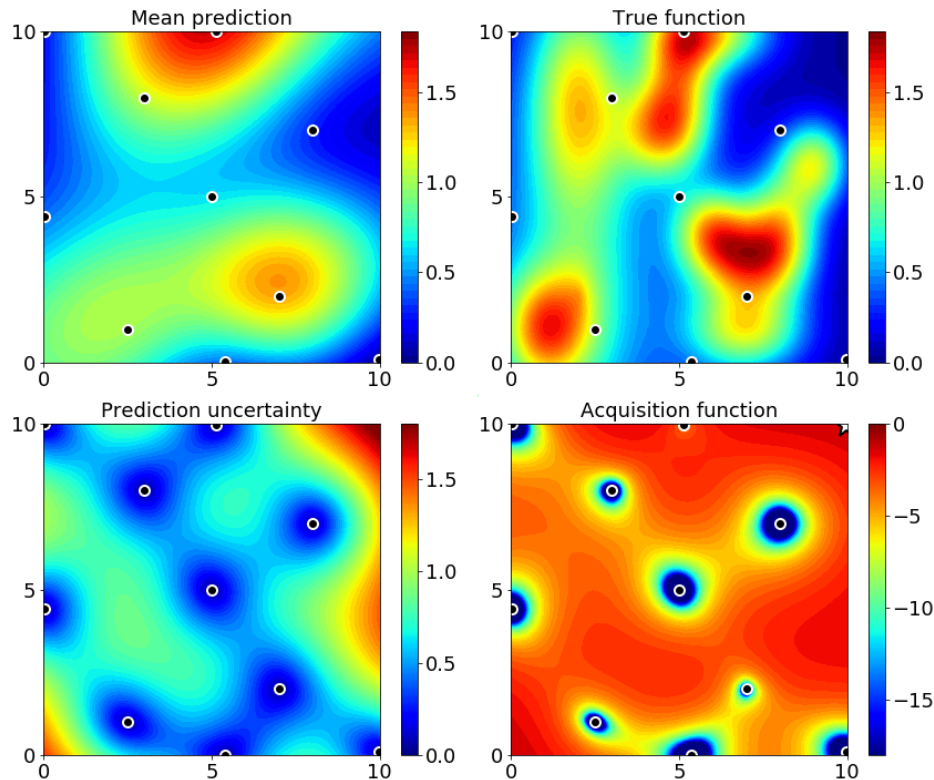
# Как выглядит Expected Improvement

- Очень неприятная функция для оптимизации
- Но ее просто считать и нам не нужен истинный максимум
- Поэтому запускаем мультистарт 5-10 итераций и берем лучший результат



# Параллельная оптимизация

- Если мы хотим учить модели параллельно, то мы можем добавить штраф за то, что мы похожи на предыдущую точку
- Все работает неплохо — потому что функция сложная для многомерного входа





# Оптимизация с ограничениями

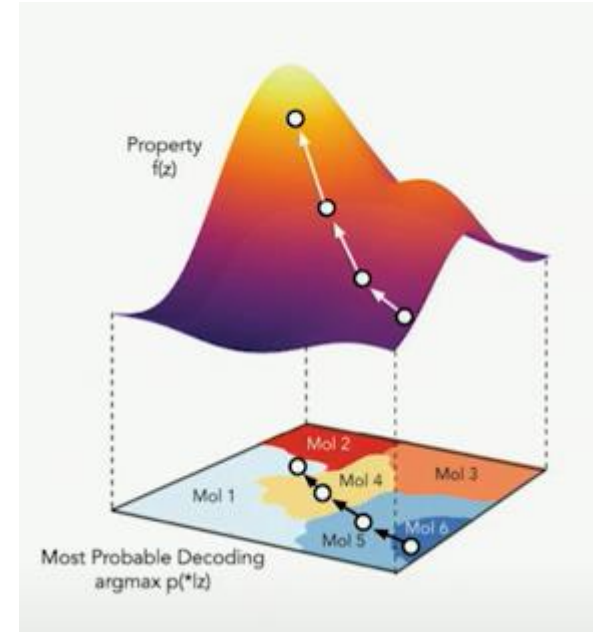
$$f(\mathbf{x}) \rightarrow \min_{\mathbf{x} \in X} \quad c.t. \ c(\mathbf{x}) \geq 0$$

Зависит от ограничений

- В общем – строим и для них модели
- Штрафуем за нарушение ограничений в модели

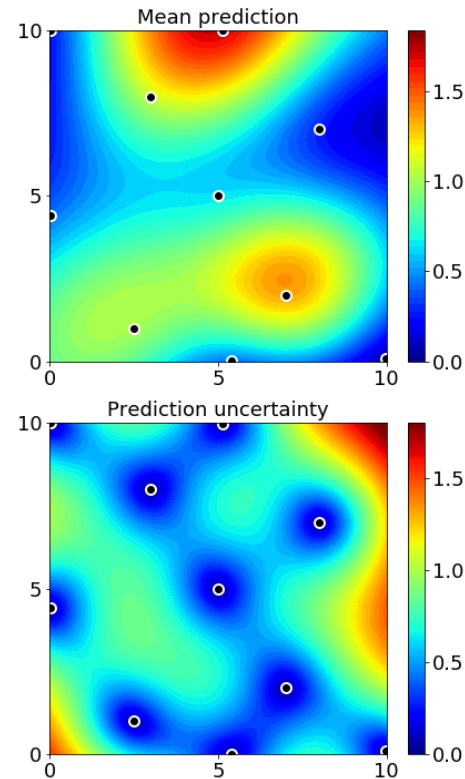
# Приложения Байесовской оптимизации

- Оптимизация гиперпараметров, в том числе AlphaGo
- Суррогатная оптимизация: оптимизация сложных инженерных изделий
- Оптимизация свойств молекул в пространстве представлений
- Оптимизация в роботехнике
- Bayesian Optimization for a Better Desert



# Связанная задача – активное обучение

- Как нам собрать такую выборку, чтобы построенная для нее модель была самой лучшей?
- Будем собирать данные последовательно, на каждом шаге нужно выбрать точку, которую нужно разметить
- У нас тоже возникает функция выгоды – какую точку брать? В этом случае речь только про exploration
- Оптимальный критерий с точки зрения теории – найти точку с наибольшей оценкой неопределенности
- На практике еще хорошо бы «смотреть в будущее»: как нам улучшить интеграл по неопределенности?



$$p(f_*|y) = \mathcal{N}(f_*|\mu_*, \sigma_*^2),$$

$$\mu_* = \mathbf{k}_*^T [\mathbf{K} + \sigma^2 \mathbf{I}_m]^{-1} \mathbf{y}, \quad \sigma_*^2 = K_{**} - \mathbf{k}_*^T [\mathbf{K} + \sigma^2 \mathbf{I}_m]^{-1} \mathbf{k}_*$$

# Выводы



# Выводы

- Байесовская (суррогатная) оптимизация позволяет быстро находить оптимумы тяжелых функций
- Она часто используется для оптимизации гиперпараметров
- Вариант по умолчанию – Ожидаемое улучшение EI или Парзеновская оценка на основе деревьев TPE

# Ссылки

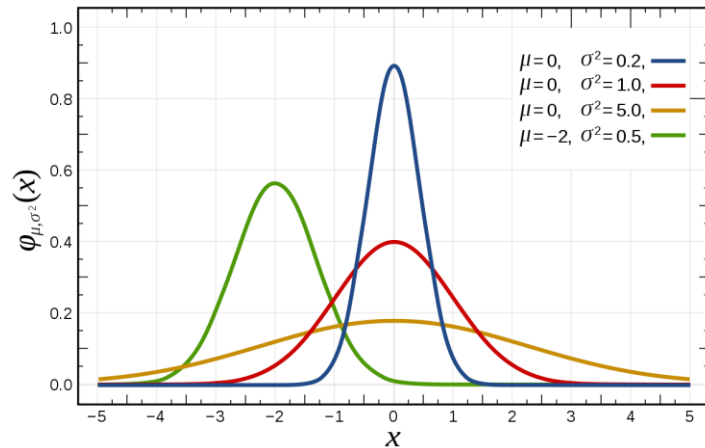
- Вводный материал про регрессию на основе гауссовских процессов  
<https://thegradient.pub/gaussian-process-not-quite-for-dummies/>
- Основная книжка про использование гауссовских процессов в машинном обучении  
<https://gaussianprocess.org/gpml/>
- Книжка про Байесовскую оптимизацию, 2023  
<https://bayesoptbook.com/>
- Лекция про Байесовскую оптимизацию с конференции UAI  
<https://www.youtube.com/watch?v=C5nqEHpdyoE>

# **Дополнительные слайды**

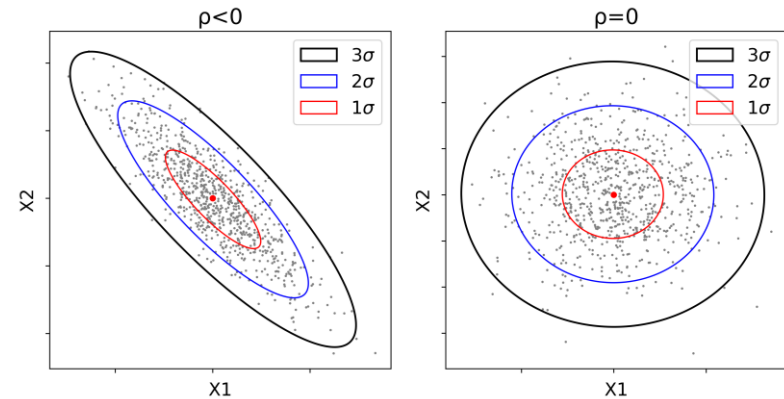
---

# Многомерное гауссовское распределение

Плотность одномерного гауссовского распределения



Плотность многомерного гауссовского распределения



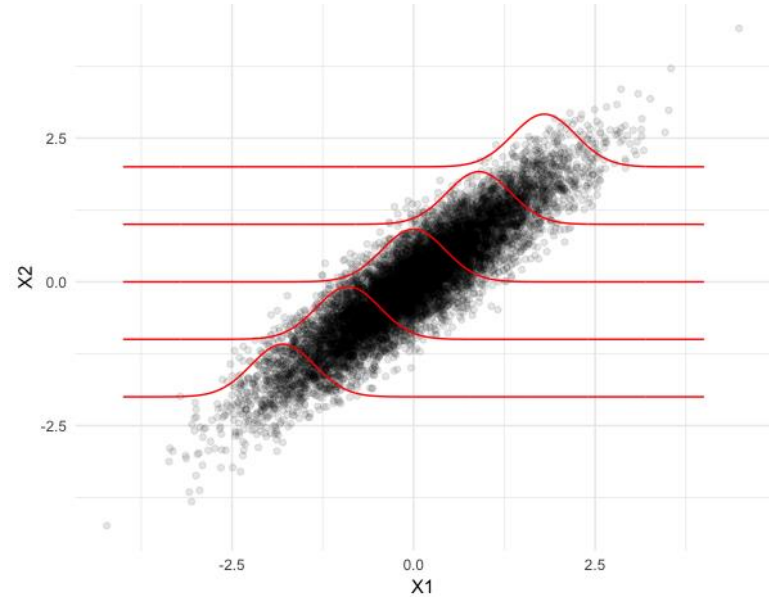
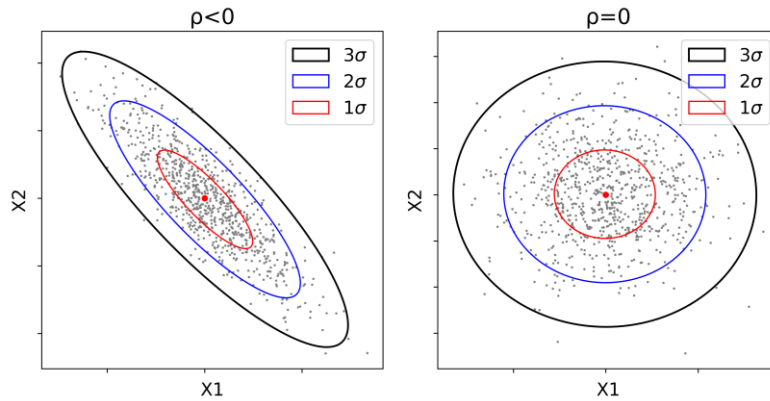
$$p(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\frac{1}{2\sigma^2} (x - \mu)^2\right)$$

$$p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)\right)$$



# Условное гауссовское распределение

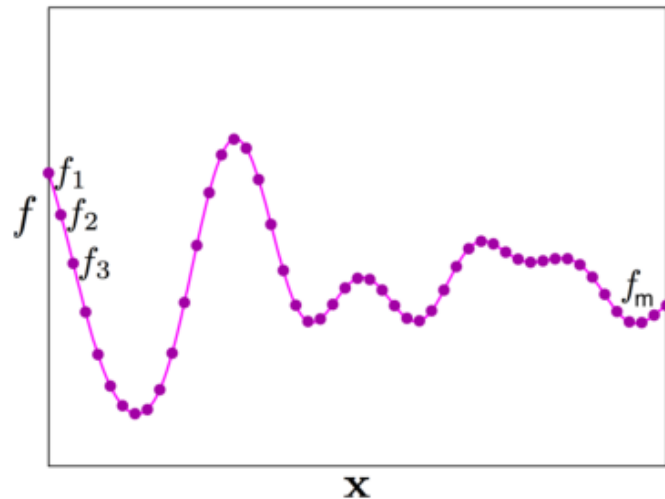
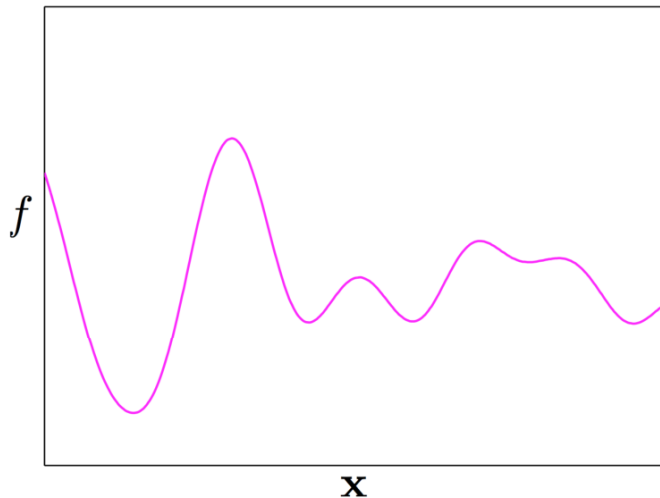
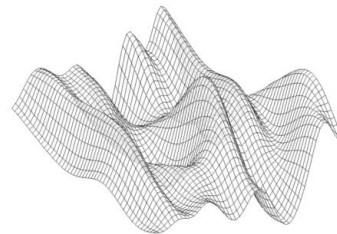
Плотность многомерного  
гауссовского распределения



$$p(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{\frac{d}{2}} |\boldsymbol{\Sigma}|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

# Гауссовский случайный процесс

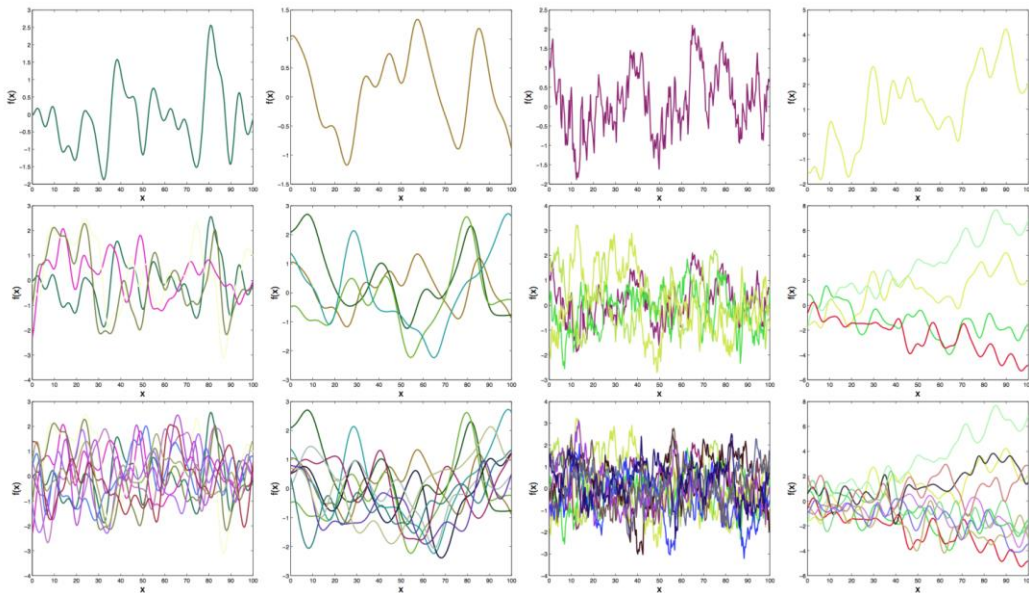
Совместное распределение для любого множества точек - нормальное



$$f(\mathbf{x}) \sim \mathcal{GP}(\cdot | \mu(\mathbf{x}), K(\mathbf{x}, \mathbf{x}'))$$

# Регрессия на основе гауссовских процессов

$$K(\mathbf{x}, \mathbf{x}') = \sigma_f^2 \exp \left\{ - \sum_{i=1}^d \frac{(x_i - x'_i)^2}{2r_i^2} \right\}$$



# Метод максимума правдоподобия для оценки параметров

$$\mathcal{L} = -\log p(\mathbf{y}|\boldsymbol{\theta}) = \underbrace{\frac{1}{2} \log \det \mathbf{C}(\boldsymbol{\theta})}_{\text{regularization}} + \underbrace{\frac{1}{2} \mathbf{y}^T \mathbf{C}^{-1}(\boldsymbol{\theta}) \mathbf{y}}_{\text{data-fit}} + \frac{m}{2} \log(2\pi),$$

- Умеем считать правдоподобие
- Умеем считать производные
- Запускаем градиентный спуск до сходимости, часто глобальный оптимум один