

# Машинное обучение: Сэмплирование

Кураленок И.Е.

Яндекс

22 марта 2012 г.

# Содержание

- 1 Понятие сэмплирования
  - Процесс сэмплирования
  - Основные методы сэмплирования
- 2 Переборные методы в ML
  - Полный перебор
- 3 Hill climbing
- 4 Сэмплирование марковскими цепями
  - Metropolis-Hastings алгоритм
  - Алгоритм Гиббса
- 5 Построение вероятного пространства для максимизации

# Понятие сэмплирования

*Сэмплирование – метод исследования множества путем анализа его подмножеств.*

Применяется:

- множество слишком велико для перебора;
- каждое дополнительное измерение дорого;
- предварительный анализ;

# Алгоритм сэмплирования

- 1 Понять какое множество мы изучаем
- 2 Осознать что из этого множества мы можем измерить
- 3 Определить количество измерений
- 4 Разработать план сэмплирования
- 5 Провести сэмплирование

# Типы сэмплирования

## Вероятностное сэмплирование

$$p(x), \forall x : p(x) > 0 \quad (1)$$

*Например: попробуем посчитать соотношение мужчин/женщин*

## Невероятностное сэмплирование

$$p(x), \exists x : p(x) = 0 \quad (2)$$

*Например: попробуем посчитать безработных в рабочее время*

Без возвратений

С возвратами

# Виды сэмплирования

- Вероятностное сэмплирование
  - Простое вероятностное
  - Систематическое
  - Стратифицированное (oversampling!)
  - Пропорциональное
  - Кластерное
- Невероятностное сэмплирование
  - Опрос ближайших
  - Панельное сэмплирование

## Как выбрать нужное?

Надо учитывать:

- природа и размер возможного сампла;
- наличие дополнительной информации о элементах;
- необходимая точность измерений;
- точность отдельных измерений в сэмплировании;
- стоимость измерений;

## Возвращаясь к ML

$$F_0 = \underset{F}{\operatorname{argmax}} p(F|X) \quad (3)$$

- + если известны вероятности можно попробовать посэмплировать решения;
- не определено пространство  $F$ ;
- неясно как устроить обход;



# Иногда все просто

$$F_0 = \operatorname{argmax}_{F \in \{f_i\}_{i=1}^n} p(F|X) \quad (4)$$

- 1 введем порядок обхода;
- 2 переберем все возможные решения;
- 3 составим взвешенное решение/выберем лучшее;

## Но чаще все непросто

$$F_0 = \operatorname{argmax}_{F \in \{f_i\}_{i=1}^{\infty}} p(F|X) \quad (5)$$

- 1 введем порядок обхода;
- 2 применим систематическое сэмплирование;
- 3 составим взвешенное решение/выберем лучшее;

## Случайное блуждание

Чтобы построить порядок обхода можно воспользоваться такой схемой:

$$F = F(x, \lambda), \quad \lambda \in \mathbb{R}^n \quad (6)$$

$$F_t = F(x, \lambda_t) \quad (7)$$

$$\lambda_{t+1} = \lambda_t + \xi \quad C(\lambda_{t+1} | \{\lambda_i\}_0^t) \quad (8)$$

- будем блуждать по пространству параметров;
- необходимо определить:
  - 1 способ сделать шаг;
  - 2 условие принятия этого шага;

## Случайное блуждание

На что стоит обратить внимание при построении блуждания:

- размерность  $\lambda$  может быть меньше чем кажется;
- ограничения на  $\lambda$  существенно осложняют процедуру:

## Некоторые виды случайного блуждания

- множество фиксированных шагов  $\xi \sim U(\{\xi_i\}_1^m)$ ;
- гауссовское  $\xi_i \sim N(\mu, \sigma^2)$ ;
- самозависимое;
- etc.

## Simple hill climbing

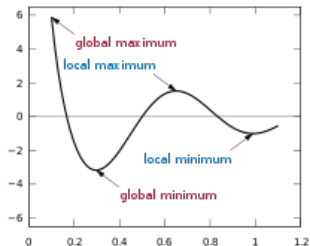
$$\xi \sim U(\{\xi_i\}), \xi_{ii} = \omega, \xi_{ij} = 0, j \neq i, \quad (9)$$

$$C(\lambda_{t+1}|\lambda_t) = \frac{p(F(\lambda_{t+1})|X)}{p(F(\lambda_t)|X)} > 1 \quad (10)$$

Свойства:

- простой;
- быстро сходится;
- зависим от выбора начальной точки.

# Random-restart (shotgun) hill climbing



Проблемы:

- сходится в локальный максимум;
- может долго сходиться, если начало далеко от максимума;
- аллеи.

⇒ Можно рестартить hill climbing из разных начальных точек

# Интуиция

- наверное нельзя всегда ходить “по шерсти”;
- хорошо бы обойти все пространство;
- скорость движения должна меняться.

⇒ Markov Chain Monte-Carlo (MCMC)



# Metropolis-Hastings алгоритм

$$p(\lambda|\lambda_t) \tag{11}$$

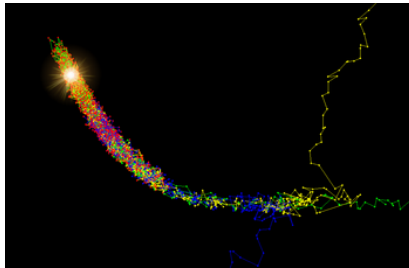
$$\alpha = \frac{p(F(\lambda_t)|X)p(\lambda_t|\lambda)}{p(F(\lambda)|X)p(\lambda|\lambda_t)} \tag{12}$$

$$\psi \sim U(0, 1) \tag{13}$$

$$C(\lambda_{t+1}|\lambda_t) = \begin{cases} 1, & \alpha \geq \psi \\ 0 \end{cases} \tag{14}$$

Например  $p(\lambda|\lambda_t) \sim N(\lambda_t|\sigma^2 E)$

# Свойства



- Доказано:
  - обходит все пространство;
  - это действительно взвешенное сэмплирование;

⇒ точно придем в максимум!

Проблема только с тем, что придем за бесконечное время

# Алгоритм Гиббса

В Метрополисе есть проблема: все зависит от  $p(\lambda|\lambda_t)$ .

$$i \sim U(1, \dots, n) \quad (15)$$

$$p(\xi|\lambda_t) \quad (16)$$

$$\lambda_i = \lambda_{t_i} + \xi, \quad (17)$$

$$\lambda_j = \lambda_{t_j}, j \neq i \quad (18)$$

# Как можно построить $p(F|X)$

Если MSE, то все просто:

$$p(\lambda|X) = \frac{e^{-c\|F(\lambda|X)-Y\|_2}}{Z} \quad (19)$$

$$Z = \int_{\lambda} e^{-c\|F(\lambda|X)-Y\|_2} d\lambda \quad (20)$$

Если максимизируем, то надеемся задрать  $Y$  так, чтобы  $Z$  был определен.