

# Пару слов про сэмплирование

И. Куралёнок

СПб, 2017

# Понятие сэмплирования

*Сэмплирование* — метод исследования множества путём анализа его подмножеств.

Применяется когда:

- 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:

$$P(x), \forall x : P(x) > 0$$

- Невероятностное сэмплирование:

$$P(x), \exists x : P(x) = 0$$

- Без возвратений
- С возвратами

# Типы сэмплирования

- Вероятностное сэмплирование:

$$P(x), \forall x : P(x) > 0$$

*Например: попробуем посчитать соотношение мужчин/женщин*

- Невероятностное сэмплирование:

$$P(x), \exists x : P(x) = 0$$

- Без возвратений
- С возвратениями

# Типы сэмплирования

- Вероятностное сэмплирование:

$$P(x), \forall x : P(x) > 0$$

- Невероятностное сэмплирование:

$$P(x), \exists x : P(x) = 0$$

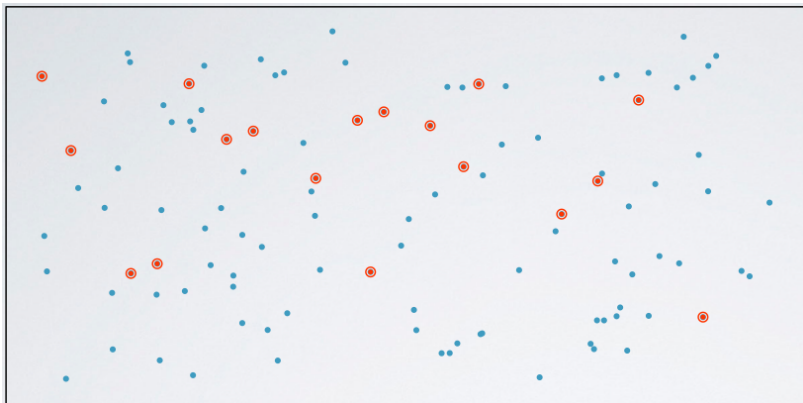
*Например: “по результатам опроса superjob.ru, 100% россиян пользуются интернетом”*

- Без возвращений
- С возвращениями

# Виды сэмплирования

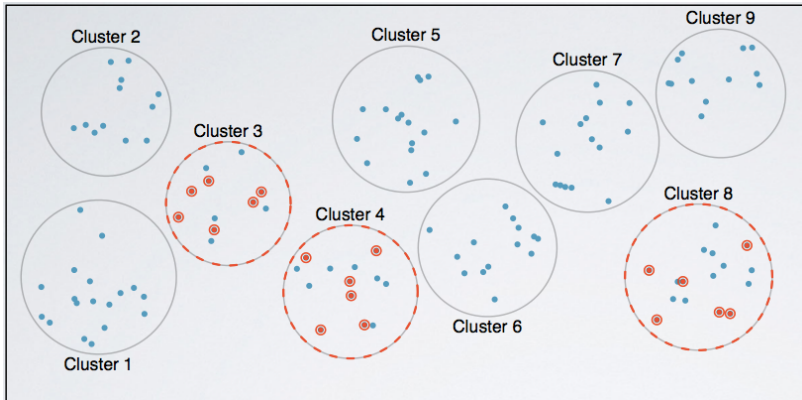
- Вероятностное сэмплирование
  - Простое вероятностное
  - Систематическое
  - Пропорциональное
  - Кластерное
  - Стратифицированное
- Невероятностное сэмплирование
  - Опрос ближайших
  - Панельное сэмплирование

# Простое вероятностное

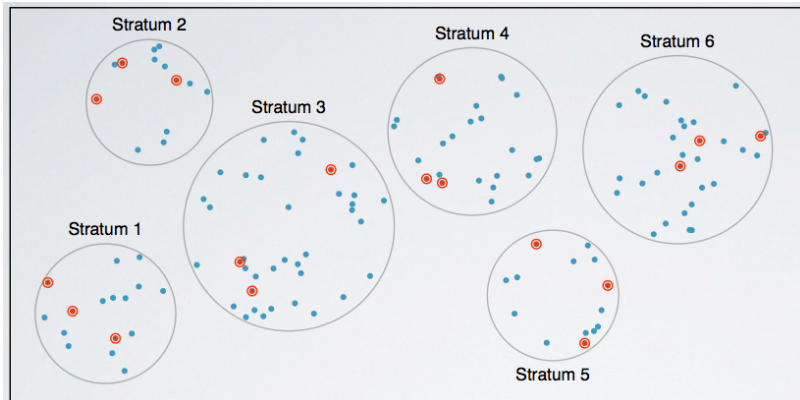




# Кластерное



# Стратифицированное



# Как выбрать нужное?

Надо учитывать:

- природа и размер возможного сэмпла;
- наличие дополнительной информации об элементах;
- необходимая точность измерений;
- точность отдельных измерений в сэмпле;
- стоимость измерений.

# Как делать не нужно

Байки про ошибки в создании обучающей выборки

- Соотношение положительных и отрицательных примеров
- Все решения одинаково бесполезны
- Ошибка в данных больше разницы между методами
- Правда в глазах смотрящего (про ошибки в примерах)
- Нужно следить за распределением важных параметров, независимо от способа сбора данных
- Устаревшие данные
- Correlation vs. Causation
- Зависимость результата от одной точки или класса точек
- Закодировать ответ в DS

# Создание DS — тоже оптимизация

Как же понять что построенное множество хорошее? Хотим такого:

$$\begin{aligned}\hat{H} &= \arg \max_H \mu_{\xi \sim D} T(y_\xi, H(x_\xi)) \\ &= \arg \max_H \mu_{\xi \sim \Gamma} T(y_\xi, H(x_\xi))\end{aligned}$$

Но делать будем иначе:

$$D = \arg \max_D \mu_{\eta \sim \Gamma} T(y_\eta, \arg \max_H \mu_{\xi \sim D} T(y_\xi, H(x_\xi)))(x_\eta)$$

Это такая оптимизация

# Заключение о построении DS

Это все область выборочного контроля про которую много написано:

*George E.P. Vox, А. И. Орлов, В.П. Боровиков* Создание хорошего обучающего множества — половина успеха. У меня нету универсальной схемы как такое делать, и есть ощущение, что это на грани искусства.

# Выборка как генеральная совокупность

- LOO методы и jackknife оценки
- Bootstrapping подход (Монте-Карло на выборке)
- Слабая аксиоматика Воронцова

# Leave One Out

Выберем один из примеров и посмотрим насколько поменяется выборка, если его в ней не будет. Усредним по всем.

Эту логику можно применять например так:

$$\text{Var}_{LOO}(\{x_i\}) = \frac{1}{m} \sum_{i=1}^m (x_i - \bar{x}_{-i})^2 = \frac{1}{(m-1)^2} \text{Var}(\{x_i\})$$



# Jackknife оценки

Пусть нас интересует статистика  $\theta$ ,  $\hat{\theta}$ —ее выборочная оценка и  $\bar{\theta} = \frac{1}{m} \sum_i \hat{\theta}_{-i}$ .

$$\begin{aligned}\mathbb{E}(\hat{\theta}) &= \theta + \frac{a}{m} + \frac{b}{m^2} + O(m^{-3}) \\ \Rightarrow \mathbb{E}(\hat{\theta} - \bar{\theta}) &= \frac{a}{m(m-1)} + O(m^{-3}) \\ \Rightarrow b_{jack} &= (m-1)(\bar{\theta} - \hat{\theta})\end{aligned}$$

# Bootstrapping

B. Efron “Bootstrap methods: another look at the jackknife”

Зачем ограничиваться одним примером? Устроим на нашей “генеральной совокупности” Монте-Карло:

- 1 рассматривая выборку как генеральную совокупность, породим повторную выборку  $X'$
- 2 вычислим интересующую нас статистику  $\hat{\theta}$  на  $X'$
- 3 проделаем 1 и 2 достаточное количество раз и построим эмпирическое распределение  $\theta$
- 4 исследуем эмпирическое распределение и сделаем выводы про  $\theta$

# Bootstrapping: простая реализация

- 1 случайно отсортируем исходную выборку
- 2 выбросим  $i$  равномерное от 1 до  $m$
- 3 возьмем  $i$ -й пример в новую выборку
- 4 повторим п. 2 и 3  $m$  раз
- 5 посчитаем на полученной выборке значение  $\theta$
- 6 исследуем поведение  $\theta$  на порожденных выборках

# Bootstrapping: практические реализации

К сожалению в исходной версии bootstrapping часто не эффективен с точки зрения производительности

- Бросить много пуассонов с  $\lambda = 1$
- Байесовский bootstrap (волшебный log)
- Гладкий bootstrap (шумный сигнал)
- etc.

# Bootstrapping: практические реализации

К сожалению в исходной версии bootstrapping часто не эффективен с точки зрения производительности (Почему?)

- Бросить много пуассонов с  $\lambda = 1$
- Байесовский bootstrap (волшебный log)
- Гладкий bootstrap (шумный сигнал)
- etc.

# Несколько рекомендаций

Когда нельзя пользоваться bootstrapping'ом:

- Когда бесконечная дисперсия
- Когда совсем мало примеров и возможны повторы  $X'$

Нужно всегда помнить что:

- Повторные выборки зависимы и степень зависимости варьируется от  $m$
- В повторных выборках значение, например, дисперсии часто ниже чем на генеральной совокупности
- Чтобы понять насколько результат смещен можно применить дополнительные деления выборки

# Домашнее задание

- Сгенерировать последовательность  $\{(x_i, y_i) | x_i, y_i \sim U(0, 1)\}_{i=1}^m$ ,  $m = 10000$
- Выбросить минимально возможное количество примеров так, чтобы  $y$  стало линейно зависимо от  $x$  с уровнем значимости  $\alpha = 0.05$
- Сообщить процент выброшенного

# Переборные методы как частный случай сэмплирования

Сэмплирование можно использовать не только для сбора данных! Представим себе, что процесс оптимизации устроен так:

$$\hat{H} = \arg \max_{\beta \in \mathbb{R}^m} \mathbb{E}_{\xi \sim D}(T(y_\xi, H(x_\xi, \beta)))$$

И нету у нас никаких способов понять свойства зависимости  $T$  от  $\beta$ .



# Формальная постановка

$$\hat{H} = \arg \max_H P(H|X)$$

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

# Иногда все просто

$$\hat{H} = \arg \max_{H \in \{H_i\}_{i=1}^n} P(H|X)$$

- 1 введём порядок обхода;
- 2 переберём все возможные решения;
- 3 составим взвешенное решение/выберем лучшее.

# Но чаще всё непросто

$$\hat{H} = \arg \max_{H \in \{H_i\}_{i=1}^{\infty}} P(H|X)$$

или совсем запущено:

$$\hat{H} = \arg \max_{H(x, \beta), \beta \in \mathbb{R}^k} P(H|X)$$

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

# Случайное блуждание I

$$H = H(x, \beta), \beta \in \mathbb{R}^k$$

Чтобы построить порядок обхода можно воспользоваться такой схемой:

$$H_t = H(x, \beta_t)$$
$$\mathcal{A} = \begin{cases} \beta_{t+1} = \beta_t + \xi \\ C(\beta_{t+1} | \{\beta_i\}_0^t) \end{cases}$$

Для этого необходимо определить:

- 1 начальную точку  $\beta_0$ ;
- 2 способ сделать шаг  $\xi$ ;
- 3 условие принятия этого шага  $C$ .

# Случайное блуждание II

На что стоит обратить внимание при построении блуждания:

- размерность  $\beta$  может быть меньше чем кажется;
- ограничения на  $\beta$  существенно осложняют процедуру.

# Некоторые виды случайного блуждания

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

# Simple hill climbing

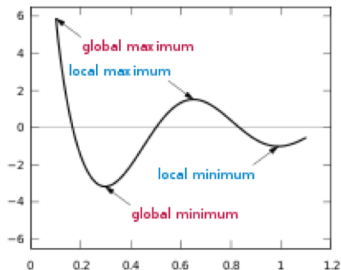
$$\xi \sim U(\{\xi_i\}_1^{2k}), \xi_i = -1^{i \bmod 2} \omega_t,$$

$$C(\beta_{t+1}|\beta_t) = \frac{P(H(\beta_{t+1})|X)}{P(H(\beta_t)|X)} > 1$$

Свойства:

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

# 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

Мы бы хотели получить сэмплирование, а для этого:

- хорошо бы обойти всё пространство;
- нельзя всегда ходить "по шерсти";
- скорость движения должна меняться в зависимости от плотности.

⇒ Markov Chain Monte-Carlo (MCMC)

# Metropolis-Hastings алгоритм

Введем  $p(\beta_1|\beta_2)$ , отвечающую за локальность.

$$\alpha = \frac{P(H_{\beta_{t+1}}|X)P(\beta_{t+1}|\beta_t)}{P(H_{\beta_t}|X)P(\beta_t|\beta_{t+1})}$$

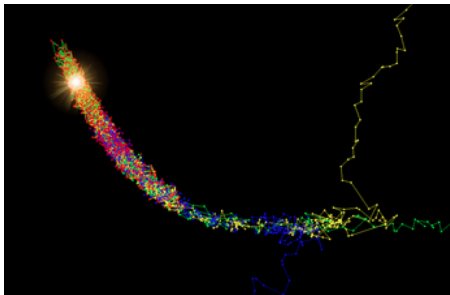
$$\psi \sim U(0, 1)$$

$$C(\beta_{t+1}|\beta_t) = \begin{cases} 1, & \alpha \geq \psi \\ 0 & \end{cases}$$

Например,  $P(\beta_{t+1}|\beta_t) \sim N(\beta_t, \sigma^2 E)$

Если  $P(\beta_1|\beta_2) = P(\beta_2|\beta_1)$  — это Metropolis

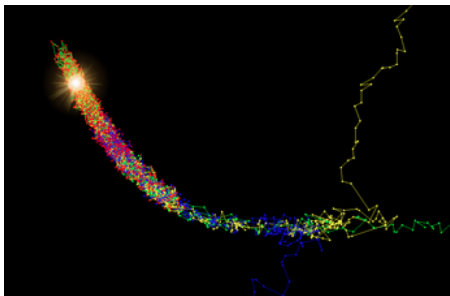
# Свойства



- + Обходит всё пространство
- + Это действительно взвешенное сэмплирование
- Последовательные сэмплы похожи друг на друга
- На этапе разогрева показывает что-то странное

⇒ *Точно придём в максимум!*

# Свойства



- + Обходит всё пространство
- + Это действительно взвешенное сэмплирование
- Последовательные сэмплы похожи друг на друга
- На этапе разогрева показывает что-то странное

⇒ *Точно придём в максимум!*

**Проблема только с тем, что придём за бесконечное время**

# Сложности в использовании

- Сходимость зависит от выбора  $P(\beta_{t+1}|\beta_t)$
- Если хотим использовать разумное распределение, оно многомерное  $\Rightarrow$  его сложно реализовывать

# Пример с дартс I

Вася и Петя повесели на стене мишень и поиграли в дартс. На следующий день пришел их бригадир Юра и задался вопросом кто из подчиненных будет заделывать дырки в стенах.

- Будет честнее, чтобы “косой” платил больше
- Дырки все одинаковые
- Каждый говорит: “да я токо разок кинул!”
- Один признался, что: “ну вот это — моя дырка”
- Играли без употребления, так что можно считать, что кидали  $\sim N(m_i, \Sigma_i)$ ,  $i \in \{\text{Вася, Петя}\}$  и параметры не менялись во времени

**Надо помочь Юре!**

# Пример с дартс II

При фиксированных параметрах вероятность увидеть дырки  $\{x_j\}$ ,  $x_j \in \mathbb{R}^2$

$$LL = \sum_j \log(\pi N(m_A, \Sigma_A)(x_j) + (1 - \pi)N(m_B, \Sigma_B)(x_j))$$

где  $A$  — Вася,  $B$  — Петя,  $\pi$  — доля Васиных бросков,  $m_i$  — сбитость прицела,  $\Sigma_i$  — разброс стрелка.

Максимизировать такое добро трудно из-за суммы под логарифмом. Было бы классно точно знать, что тот или иной бросок точно сделал Вася  $A_t \in \{0, 1\}$ :

$$LL = \sum_t A_t \log(N(m_A, \Sigma_A)(x_j)) + (1 - A_t) \log(N(m_B, \Sigma_B)(x_j))$$

В таком варианте все совсем просто считается (если мы еще не забыли формулу плотности нормального многомерного распределения :))

# Пример с дартс III

Если чуть более формально, то мы ввели скрытую/ненаблюдаемую переменную  $I(A)$  (такое называется data augmentation) и хотим:

- 1 Прикинуть начальные  $m_i, \Sigma_i, \pi$
- 2 Посчитать ожидание скрытого параметра  $A_j = \mu((m_i, \Sigma_i, \pi))$
- 3 В условиях найденных  $A_j$  максимизировать  $LL$
- 4 Перейти к второму шагу до сходимости



# Expectation maximization алгоритм

Фокус, который мы проделали называется *EM*-алгоритм. Пусть данные у нас состоят из видимой части  $L$  и скрытой  $Z$ ,  $X = (L, Z)$ . Мы хотим найти оптимальные параметры  $\beta$ :

$$\log P(L|\beta) = \log P(X|\beta)$$

- 1 Возьмем какой-то  $\beta_0$
- 2 Посчитаем ожидание (expectation step):

$$\mathbb{E}(\beta|\beta_t) = \mathbb{E}_Z(\log P(X, \beta, Z|\beta_t))$$

- 3 Проведем оптимизацию (maximization step):

$$\beta_{t+1} = \arg \max_{\beta} \mathbb{E}(\beta|\beta_t)$$

- 4 Перейти к второму шагу до сходимости

# Немного про Гиббса/Больцмана

Есть такое распределение:

$$p(x) = \frac{e^{\frac{e(x)}{kT}}}{Z = \int_x e^{\frac{e(x)}{kT}} dx}$$

Про него известно, для фиксированной функции  $e(x) \geq 0$  :  $Z < \infty$  и условия на общую энергию системы:

$$\int_x e(x) dP(x) \leq \text{const}$$

оно доставляет в максимум энтропию. А еще это добро эквивалентно MRF.

# Алгоритм Гиббса

В Метрополисе есть проблема: многомерное распределение  $P(H|X) = P(\beta|X)$ . Можно попробовать рассматривать его по частям.

- 1 Начнем с какого-то  $\beta$
- 2 Сгенерируем  $\beta_{t+1}$  по правилам

$$P(\beta_{t+1,i}|X, \beta_{t+1,1}, \dots, \beta_{t+1,i-1}, \beta_{t,i+1}, \dots, \beta_{t,n})$$

- 3 Будем бегать по  $i = 1 \dots n$ , пока не сойдется
- 4 Полученная  $\beta$  — следующая точка сэмплирования
- 5 Пока хочется идем в п.2

# Пример с дартс IV

На языке товарищей Геман:

- 1 Прикинуть начальные  $m_i, \Sigma_i, \pi$
- 2 Для каждого сэмпла сгенерировать  $A_j$  из текущего  $\pi$
- 3 В условиях найденных  $A_j$  максимизировать  $LL$
- 4 Перейти к второму шагу до сходимости того момента пока распределение параметров не перестанет меняться

# Как можно построить $P(F|X)$ для RMSE

$$P(\beta|X) = \frac{e^{-c\|H(\beta|X)-Y\|_2}}{Z}$$

$$Z = \int_{\beta} e^{-c\|H(\beta|X)-Y\|_2} d\beta$$

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

Есть только одна проблема:

# Как можно построить $P(F|X)$ для RMSE

$$P(\beta|X) = \frac{e^{-c\|H(\beta|X)-Y\|_2}}{Z}$$

$$Z = \int_{\beta} e^{-c\|H(\beta|X)-Y\|_2} d\beta$$

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

Есть только одна проблема: мы хотим максимизировать, а не считать среднее :)

# NFL: формальная формулировка

$$d_m = \{(d_m^x(1), d_m^y(1)), \dots, (d_m^x(m), d_m^y(m))\}$$
$$f : \mathcal{X} \rightarrow \mathcal{Y}$$
$$\mathcal{F} = \mathcal{Y}^{\mathcal{X}}$$
$$p(f) = \frac{1}{\mathcal{F}}$$

Theorem (David Wolpert and William G. Macready (1997))

*Для любых двух способов обхода  $a_1$  и  $a_2$ :*

$$\sum_f P(d_m^y | f, m, a_1) = \sum_f P(d_m^y | f, m, a_2)$$

# NFL: следствия

- Нам всем хватит работы :)
- Полоса белая, полоса черная
- Нужно искать близкие задачи



# Заключение

Сэмплирование очень полезная вещь:

- собирать данные;
- брать интегралы;
- обучаться.

Посмотрите на `bayesian inference` и библиотеку Гиббса для него: `BUGS`.

# Задание на дом

- Датасет, как обычно, в svn
- Суть - учимся предсказывать размер аудитории сервиса
- $N(t) = \left(\frac{1}{1+e^{at+b}}\right)N_0$
- Ищем  $a$ ,  $b$ ,  $N_0$
- Можно прямо hill climbing
- Можно подумать и воспользоваться алгоритмом Гиббса (это будет плюсом)
- Дедлайн 17 октября