

# Лекция 6.

ББББ

(Бутстрепинг, бутстреп,  
бэггинг и бустинг).

Москва  
14.10.2016

Павел Владимирович  
Слипенчук  
[PavelMSTU@stego.su](mailto:PavelMSTU@stego.su)  
ИУ-8

# Бутстрепинг (статистика)

- Генерация выборки размера  $N$  из подвыборки размера  $n \ll N$  путем выбора с повторением.
- Сколько вариантов выбора?

# Бутстрепинг (статистика)

- Генерация выборки размера **N** из подвыборки размера **n** << **N** путем выбора с повторением.
- Вариантов выбора:  $\overline{C}_N^n = C_{N+n-1}^n$
- В бутстреппинге выборка выбирается большое количество раз (например 1000, 10000)
- На каждой выборке высчитывается мат.ожидание, дисперсия.
- Для **N** мат.ожидание и дисперсия принимаются как среднее арифметическое этих величин на совокупности выборок.

# Бутстрепинг (статистика)

- В каких случаях бутстрепинг применим, а когда нет?

# Бутстрепинг (статистика)

- **Центральная Предельная Теорема.**

$X_1, \dots, X_n$  –  $n$  случайных велчин, **независимых** и **одинакого** **Распределенных**, которые **имеют** конечное мат.ожидание и дисперсию.

$\mu, \sigma$  – математическое ожидание и дисперсия

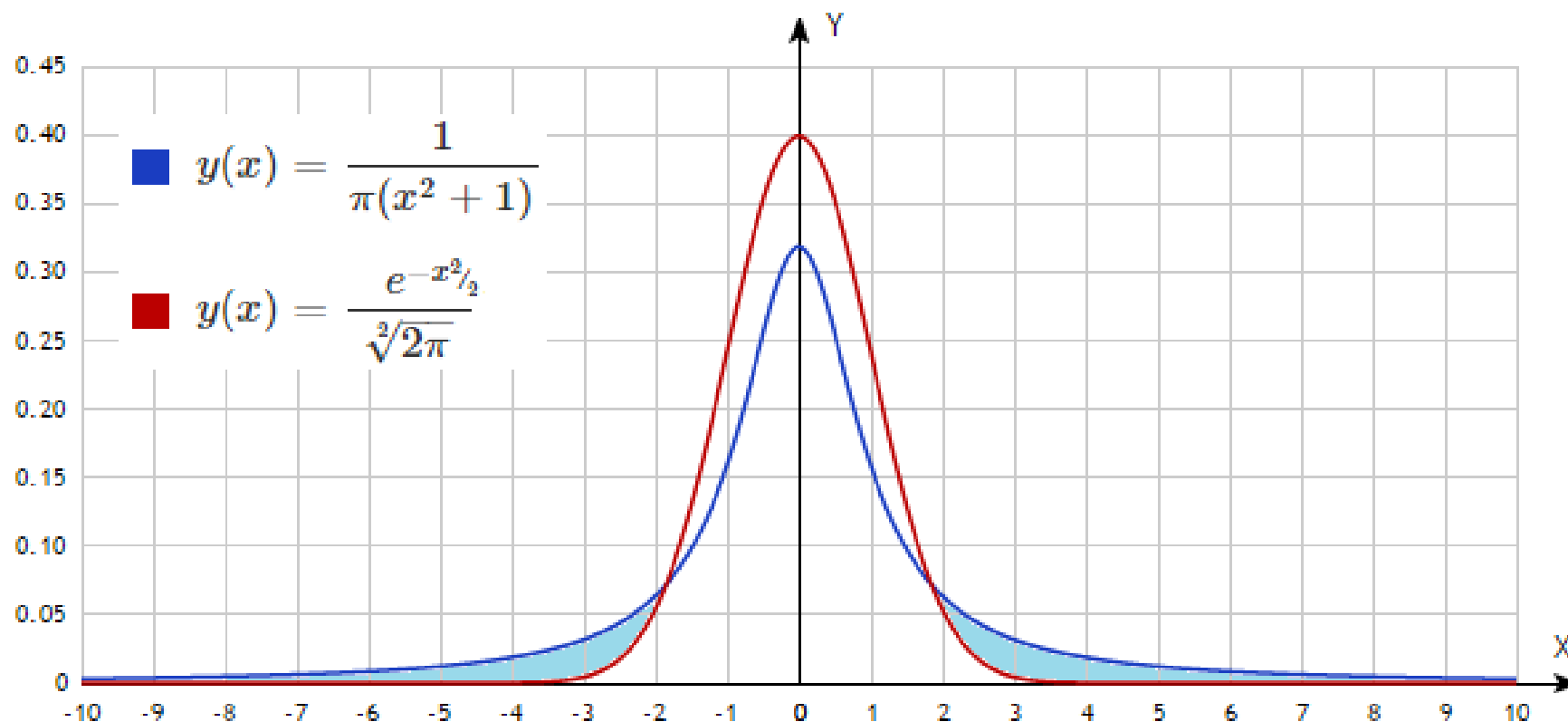
$X_n \stackrel{\text{def}}{=} \sum_{i=1}^n X_i$  – сумма  $n$  величин

$$\lim_{n \rightarrow \infty} \frac{X_n - \mu \cdot n}{\sigma \cdot \sqrt{n}} = N(0, 1)$$

- Если не конечное мат.ожидание или дисперсия – бутстрепинг не поможет.

# Бутстрепинг (статистика)

- **Пример:** распределение Коши.



- это называют «проблемой тяжелых хвостов»

# Бутсрепинг (Machine Learning)

- Идеи?

# Бутстрепинг (Machine Learning)

- Можно использовать бутстреп для генерации нескольких выборок из одной определенной выборки.
- **Зачем?**



# Бутстрепинг (Machine Learning)

- Можно использовать бутстреп для генерации нескольких выборок из одной определенной выборки.
- **Примеры.**
  - 1) Создание различных выборок для классификаторов ансамбля (например для деревьев случайного леса) – **бутстреп (bootstrap)**.
  - 2) Проверка качества модели (создается несколько тестовых выборок)
  - 3) Требуется  $n \gg N$ . Например для решения проблемы разномощности классов (см. след.слайд)

# Разномощные классы

(~ «несбалансированные классы»)

- Типичная ситуация: количество элементов в одном классе существенно меньше количества элементов в другом классе.
- **Примеры:**
  - 1) фрода существенно меньше, чем легитимных транзакций
  - 2) больных раком меньше, чем здоровых
  - 3) порядочных водителей больше, чем «лихачей»
  - 4) ....
- **Методы решения?**

# Разномощные классы

- **Методы решения**

- 1) Из большего класса взять случайным образом одну подвыборку размерности меньшего класса («классика жанра»)
- 2) Продолжать наблюдения и получать новые события, определять их классы. Затем п.1
- 3) **Бутстрепинг** для меньшего класса ( $n \gg N$ ).
- 4) «экспертная генерация выборки»

# Бутстреп-агрегация или бэггинг (Bootstrap aggregating or bagging)

- По сути это ансамбль над классификатором.
- Бутстреп-агрегация – это обобщение одного и того же классификатора на различных подвыборках или случайных параметрах, которые различны при каждой новой подгонки (fitting).
- Если итоговый классификатор – случайный лес, то беггинг – это усреднение этого леса при  $m > 1$  подгонках.

# Бустинг (boosting)

- **«Общее определение»**

**Бустинг** – это итеративное перевзвешивание наблюдений обучающей выборки. Т.е. это «последовательный» ансамбль, каждый следующий член которого корректирует предыдущий.

- + можно добавлять веса каждому ансамблю в соответствии с ошибками которые он допустил

# Бустинг (boosting)

- **Простой вариант.**

Обучающая выборка подается на вход классификатору №1. Из контрольной выборки выбираются объекты для которых классификатор №1 ошибся. Это – обучающая выборка для классификатора №2. ... ..

# Бустинг (boosting)

## ДЗ. Посмотреть.

- Градиент, градиентный спуск
- AnyBoost
- AdaBoost
- BrownBoost

# Резюме. Who is who?

- Бутстрепинг
- Бутстреп
- Бутстреп агрегация
- Бэггинг
- Бустинг