

# Байесовское моделирование категориальных признаков

---

Надежда Бугакова

научный руководитель: к. ф.-м. н. И. Е. Кураленок

17 июня 2018 г.

СПб АУ НОЦНТ РАН

# Задача машинного обучения с учителем

Существует:

$$(x, y) \sim P, x \in X, y \in Y$$

Классификация:  $y_i \in \{0, 1\}$

Хотим:

$$\phi : X \rightarrow y; \phi \in \mathbb{F}$$

$$\phi = \arg \min_{\phi \in \mathbb{F}} E_{(x,y) \sim P} L(\phi(x), y)$$

Обучающая выборка:

$\{(x_i, y_i)\}_{i=1}^m$ , реализации  $P$ .

Ищем:

$$\hat{\phi} = \arg \min_{\phi \in \mathbb{F}} E_{\{(x_i, y_i)\}} L(\phi(x), y)$$



Рис. 1: Картинки



Рис. 2: Вещественные признаки



Рис. 3: Категориальные признаки

# Категориальные признаки

- Нет отношения порядка.
- Количество различных значений бывает очень велико.

Методы:

- One-hot-encoding
- Специализированные под предметную область (например, использование вложенности признаков)
- Оценка параметров вероятностной модели (CatBoost)

**Целью** данной работы является изучение эффективности иерархического байесовского моделирования категориальных признаков для задачи классификации на основе обучающего множества.

Для достижения цели необходимо решить следующие **задачи**:

- Исследовать механизм работы с категориальными признаками, используемый в CatBoost.
- Разработать вероятностную модель для категориальных признаков с подбором параметров на основе обучающего множества.
- Реализовать преобразование категориальных признаков в упорядоченные.
- Сравнить на задачах обучения с учителем

# Вероятностная модель для категориальных признаков

**Идея:** моделировать зависимость целевых значений от категориальных признаков.

$$y_c \sim \text{Ber}(\theta_c)$$

**Простейшая модель:**  $\theta_c$  независимы.

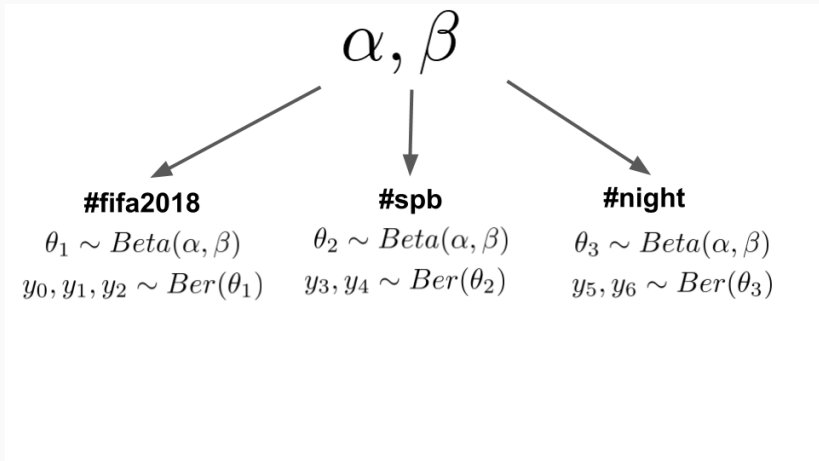
**Проблема:** переобучение.

Для категорий  $y_0 \sim \text{Ber}(0.4)$ ;  $y_1 \sim \text{Ber}(0.5)$

$$\text{Выборка: } \begin{array}{c|c} \text{cat0} & 0 \\ \text{cat1} & 1 \end{array} \Rightarrow \begin{array}{l} \hat{\theta}_0 = \frac{0}{1} = 0 \\ \hat{\theta}_1 = \frac{1}{1} = 1 \end{array}$$

**Решение:** моделирование со временем: оцениваем  $\theta_c$  только на основе целевых значений ДО текущего момента.

# Оценка параметров вероятностной модели в CatBoost



$\alpha, \beta$

**#fifa2018**

$$\theta_1 \sim \text{Beta}(\alpha, \beta)$$

$$y_0, y_1, y_2 \sim \text{Ber}(\theta_1)$$

0, 1, 1

**#spb**

$$\theta_2 \sim \text{Beta}(\alpha, \beta)$$

$$y_3, y_4 \sim \text{Ber}(\theta_2)$$

**#night**

$$\theta_3 \sim \text{Beta}(\alpha, \beta)$$

$$y_5, y_6 \sim \text{Ber}(\theta_3)$$

$$\hat{\theta}_{1,2} = E_{\theta|\alpha,\beta,y_0,y_1} \theta = \frac{y_0 + y_1 + \alpha}{\alpha + \beta + 2} = \frac{1 + \alpha}{\alpha + \beta + 2}$$

Максимизируем правдоподобие по  $(\alpha, \beta)$ :

$$\prod_c p(D_c | \alpha, \beta) = \prod_c \int_{\theta_c} \prod_{i=1}^{N_c} p(y_{c,i} | \theta_c) dBeta(\alpha, \beta)$$

$$\begin{aligned} \int_{\theta_c} \prod_{i=1}^{N_c} p(y_{c,i} | \theta_c) p(\theta_c | (\alpha, \beta)) d\theta_c &= \int_{\theta_c} \frac{\theta_c^{\sum y_{c,i} + \alpha - 1} (1 - \theta_c)^{\sum (1 - y_{c,i}) + \beta - 1}}{B(\alpha, \beta)} d\theta_c = \\ &= \frac{B(\sum y_{c,i} + \alpha, \sum (1 - y_{c,i}) + \beta)}{B(\alpha, \beta)} \end{aligned}$$



Были реализованы на Python:

- Оптимизация априорных значений с помощью метода Ньютона на основе библиотеки `scipy`.
- Преобразование категориальных признаков в вещественные тремя способами:
  - С автоматическим подбором априорных значений;
  - С фиксированными в CatBoost априорными значениями;
  - С априорными значениями, где  $\alpha$  рассчитывалась как среднее целевого значения по всей выборке и  $\beta = 1 - \alpha$ .
- Воспроизводимые эксперименты.

- Gradient boosting (CatBoost).
- Бутстрепинг тестовой выборки + Wilcoxon signed-rank test.

## Датасеты

Name	Atributes	Categorical	Learn size	Test size
adult	14	8	39074	9768
amazon	9	9	26215	6554
appet	419	38	40000	10000
kick	439	23	58388	14595

Уровень статистической значимости 0.01.

AUC (модель со временем):

Name	CatBoost+time	Auto+time	Simple+time
<b>amazon</b>	<b>0.8564</b>	0.853	0.855
adult	0.9275	<b>0.9281</b>	<b>0.9281</b>
<b>appet</b>	<b>0.8525</b>	0.8475	0.8468
kick	<b>0.7656</b>	<b>0.7657</b>	0.7638

AUC (модель без времени):

Name	CatBoost	Auto	Simple
amazon	0.8053	<b>0.8188</b>	0.81
<b>adult</b>	0.9292	0.9291	0.9292
appet	<b>0.6909</b>	0.6628	0.6615
<b>kick</b>	0.7561	<b>0.7705</b>	0.7623

Время обучения в секундах:

Name	CatBoost	Auto	Simple
amazon	48	38	35
adult	55	48	48
appet	261	226	227
kick	128	105	92

Name	CatBoost+time	Auto+time	Simple+time
amazon	46	38	37
adult	54	47	47
appet	267	235	236
kick	123	90	92

- Исследован механизм преобразования категориальных признаков в CatBoost
- Придуман и изучен автоматический метод подбора априорных значений на основе обучающей выборки.
- Реализованы разные подходы к подбору априорных значений и произведено сравнение на наборе данных

### GitHub (<https://github.com/N-buga/SPBAU-CategBays>)

- CatBoost  
(<https://tech.yandex.com/catboost/>)
- adult  
(<https://archive.ics.uci.edu/ml/datasets/Adult>)
- appet  
(<http://www.kdd.org/kdd-cup/view/kdd-cup-2009/Data>)
- amazon  
(<https://www.kaggle.com/c/amazon-employee-access-challenge>)
- kick  
(<https://www.kaggle.com/c/DontGetKicked>)