

11.05.2015

# WYKŁAD

## Algorytm Ada Boost

Yoav Freund, Robert Schapire

Artykuł w 1997, nagroda w 2003

To jest praca z dziedzinie learningu.

Powiedzmy, że chcemy zaprojektować klasyfikator.

Będziemy się zajmować klasyfikatorami binarnymi, czyli takimi, które po prostu odpowiadają TAK lub NIE.

Taki klasyfikator to funkcja  $f: X \rightarrow \{-1, 1\}$ , gdzie  $X$

to jakiś zbiór obiektów. Możemy sobie wyobrazić, że

$X$  to zbiór ~~zawodników~~ <sup>zawodników</sup> i chcemy polecić odpowiedź, czy ~~zawodnik~~ <sup>zawodnik</sup> wygrał następną mecz. Mamy wiele tzw. słabych klasyfikatorów

$f_1, f_2, f_3, \dots = X \rightarrow \{-1, 1\}$  typu: „czy wygrał ostatni mecz”,

„czy ma się dobrego trenera”, „czy z kraju, który zdobył kiedyś mistrzostwo świata” itp. Idea jest taka, by z nich zbudować

dobry klasyfikator, taki, który myli się rzadko.

Można by zorganizować głosowanie między  $f_1, f_2, \dots, f_n$ .

Albo czyli zrobić  $f = f_1 + f_2 + \dots + f_n$  i spojrzeć, czy  $f(x) \geq 0$  dla pewnego  $x \in X$ .

21.05.20.11

Jest jednak lepszy pomysł. Miannoćce można zorganizować głosowanie z wagami.

Ważniejszy  $f = \alpha_1 f_1 + \alpha_2 f_2 + \dots + \alpha_n f_n$

A potem odpowiedź to  $\text{sgn}(f(x))$ .

Pytanie tylko jak dobrać wagi. etc.

Model jest taki: mamy pewien zbiór (stabyd) klasyfikatorsów  $f_1, f_2, \dots, f_n$  oraz zbiór, na którym macemy  $n$  par  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ , gdzie  $x_i \in X, y_i \in \{-1, 1\}$ .

Chcemy znaleźć  $F = \sum_{j=1}^k \alpha_j f_j$   $\alpha_j$  by być w pewnym sensie dobry. Oczekiwane wtedy  $F: X \rightarrow [-1, 1]$

Ada Boost - adaptive boosting. Boosting jest chyba od sumowania, adaptive od tego, że będziemy dobrać po kolei. Okazuje się, że w praktyce to bardzo dobrze działa, choć algorytm jest dość prosty.

Tu tak naprawdę są 2 decyzje do podjęcia:

- które  $f_1, \dots, f_k$  wybrać?
- jakie  $\alpha_1, \dots, \alpha_k$  im dać?

Będziemy za każdym razem dobrać optymalny  $f_j$  i  $\alpha_j$ , w pewnym sensie zachłannie i przydzielad mu wagi.

~~Wzrost~~ Zamrost  $f_{i2}, f_{i1}, \dots, f_{ik}$  bsdziemy  
 oznacz  $f_1, f_2, \dots, f_k$  wybrane klasyfikatory, dla  
 uproszczenia notacji.

$$\text{Niech } F_j(x_i) = \sum_{k=1}^j d_k \cdot f_k(x_i)$$

Powiedzmy, że mamy już  $f_1, \dots, f_{m-1}$  razem z  $d_1, \dots, d_{m-1}$   
 i chcemy wybrać  $f_m$  oraz  $d_m$ .

Bsdziemy chcieć tak wybrać, by bTgd

$$E = \sum_{i=1}^N e^{-y_i \cdot F_m(x_i)}$$

bT jak najmniejszy.

Zauważmy, że

$$\sum_{i=1}^N e^{-y_i \cdot F_m(x_i)} = \sum_{i=1}^N e^{-y_i \cdot F_{m-1}(x_i)} \cdot e^{-y_i \cdot d_m f_m(x_i)}$$

Możemy więc myśleć, że jest to takie wybranie  $f_m, d_m$ ,  
 by bTgd  $e^{-y_i \cdot d_m f_m(x_i)}$  z wagami  $e^{-y_i \cdot F_{m-1}(x_i)}$

bT jak najmniejszy. Zobaczymy, że  $y_i \cdot F_{m-1}(x_i)$  bliższe 1,

gdz  $F_{m-1}(x_i)$  bliższe  $y_i \in \{-1, 1\}$ . A więc wady waga  $x_i$

to dużo  $e^{-1}$ . Gdy np. natomiast up.  $y_i = 1$ , a  $F_{m-1}(x_i) = -0,8$ ,

to waga to  $e^{0,8}$ , czyli spora.

Oznaczmy wagi  $w_i^{(m)} = e^{-y_i \cdot F_{m-1}(x_i)}$  dla  $m > 1$ ,  $w_i^{(1)} = 1$ .

przypadek skrajny

$$\text{Mamy } E = \sum_{i=1}^N w_i^{(m)} e^{-y_i \cdot d_m f_m(x_i)}$$

Rozpisujemy

$$E = \sum_{i=1}^N w_i^{(m)} e^{-y_i d_m f_m(x_i)} = \sum_{y_i = f_m(x_i)} w_i^{(m)} e^{-d_m} +$$

$$+ \sum_{y_i \neq f_m(x_i)} w_i^{(m)} e^{d_m} = \sum_{i=1}^N w_i^{(m)} e^{-d_m} +$$

$$+ \sum_{y_i \neq f_m(x_i)} w_i^{(m)} (e^{d_m} - e^{-d_m})$$

Zobaczymy, że dla ustalonego  $d_m$  suma zależy tylko od tego jak duża jest  $\sum_{y_i \neq f_m(x_i)} w_i^{(m)}$ .

Cygle po prostu trzeba wybrać  $f_m$ , który ma najmniejszy wagowy błąd względem  $w_i^{(m)}$ .

Jak już mamy  $f_m$  to teraz pytanie brzmi jakie  $d_m$  wybrać. Żeby to zrobić to po prostu zderzamy  $E$  po  $d_m$ .

$$\frac{dE}{dd_m} = - \sum_{y_i = f_m(x_i)} w_i^{(m)} e^{-d_m} + \sum_{y_i \neq f_m(x_i)} w_i^{(m)} e^{d_m}$$

$$\text{Aby } \frac{dE}{dd_m} = 0 \text{ mamy}$$

$$\sum_{y_i = f_m(x_i)} w_i^{(m)} e^{-d_m} = \sum_{y_i \neq f_m(x_i)} w_i^{(m)} e^{d_m}$$

$$e^{2d_m} = \frac{\sum_{y_i = f_m(x_i)} w_i^{(m)}}{\sum_{y_i \neq f_m(x_i)} w_i^{(m)}}$$

Czyli w efekcie dostajemy

$$\alpha_m = \frac{1}{2} \ln \left( \frac{\sum_{g_i = f_m(x_i)} w_i^{(m)}}{\sum_{g_i \neq f_m(x_i)} w_i^{(m)}} \right)$$

Dla ozn.  $\epsilon_m = \sum_{g_i \neq f_m(x_i)} w_i^{(m)} / \sum_{i=1}^N w_i^{(m)}$  mamy

$$\alpha_m = \frac{1}{2} \ln \left( \frac{1 - \epsilon_m}{\epsilon_m} \right)$$

Cyliw po prostu bierzemy szukamy  $f_m$  z min. sumy

$\sum_{g_i \neq f_m(x_i)} w_i^{(m)}$ , a jak to mamy, to bierzemy obliczamy

$\epsilon_m$ , a z tego  $\alpha_m$ .

W praktyce (chyba) możemy wyprodukować wiele losowych (pewnie stabilnych) klasyfikatorów i z nich zrobić dobry używając AdaBoost. Podobno AdaBoost jest dość odporny na tzw. overfitting, czyli zbytwe dopasowanie do danych treningowych.