

Podział danych na treningowe / testowe (walidacyjne)

MNIST

50k/10k/10k

lrys: 150 : 120/30

• trening sieci na danych treningowych

• sprawdzenie ~~jej~~ dokładności jej przewidzienia na danych testowych (bez użycia)

zwykle dokładność na danych testowych \leq dokładności na danych treningowych

Jestli różnica ^{jest} znacząca, to wskazuje, że dochodzi do przelimitowania sieci (overfitting)

Zwykle prowadzi to o tym, że model jest zbyt skomplikowany i skłonny do imitacji danych ^{uczących}.

Można też użyć regularizacji, tzn. do funkcji kosztu

dołączyć np. d. kwadrat normy L^2 wag (bez użycia w danych)



$$\frac{\partial L}{\partial w} = (\text{jak bez regularizacji}) + 2d \cdot w$$

Wskazywanie wagi: $w^* = (1 - 2d) ; w^- = c \cdot \frac{\partial L}{\partial v}$

'40 pierwsze modele matematyczne neuron (akson?)

'50 pierwsze implementacje + sieci spłotowe

'80 wielka popularność, sieci spłotowe

2000 -

1989

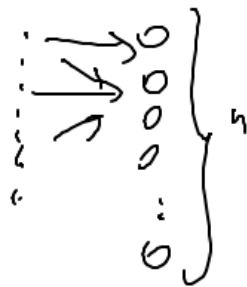
Tw. (Cybenko)

σ - funkcja ciągła, sigmoidalna (tzn. $\sigma(-\infty)=0$, $\sigma(\infty)=1$, $\sigma \uparrow$).

#. Wtedy funkcje postaci

$$G(x) = \sum_{j=1}^n \alpha_j \sigma(y_j^T x + \theta_j) \quad , x \in [0,1]^m$$

je gęste $\hookrightarrow (C([0,1]^m), \|\cdot\|_\infty)$.



Regresja logistyczna

klasyfikacja binarna $X \in \mathbb{R}^n \rightsquigarrow Y \in \{0, 1\}$
klasa

model: $p(x) = P(Y=1|X)$

$\varphi: \mathbb{R} \rightarrow [0, 1]$, $\varphi \nearrow$

$$p(x) = \varphi(\beta_0 + \beta_1 \cdot x) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}} = \frac{e^{\beta_0 + \beta_1 x}}{e^{\beta_0 + \beta_1 x} + 1}, \quad \beta_0 \in \mathbb{R}, \beta_1 \in \mathbb{R}^n$$

$$\varphi(x) = \frac{1}{1 + e^{-x}}$$

obposonywanie współczynników: zwykle maksymalizuje
największą nierozgodność:

$$L(\beta_0, \beta_1) = \prod_{i: y_i=1} p(x_i) \cdot \prod_{i: y_i=0} (1 - p(x_i))$$

$$e^{\beta_0 + \beta_1 x} = p(x) (e^{\beta_0 + \beta_1 x} + 1)$$

$$e^{\beta_0 + \beta_1 x} (1 - p(x)) = p(x)$$

$$e^{\beta_0 + \beta_1 x} = \frac{p(x)}{1 - p(x)}$$

$$\beta_0 + \beta_1 x = \log \frac{p(x)}{1 - p(x)}$$

To odpowiada sieci z jednym neuronem, z funkcją aktywacji $\varphi(x) = \frac{1}{1+e^{-x}}$



z funkcją kosztu $-L(\beta_0, \beta_1)$

Regresja logistyczna dla większej liczby klas:

$X \in \mathbb{R}^n \rightsquigarrow Y \in \{0, 1, \dots, k-1\}$ - klasy

model $\log \frac{P(Y=j|X)}{P(Y=k-1|X)} = \beta_0^{(j)} + \beta_1^{(j)} \cdot X \quad , j=0, \dots, k-2$

$\rightarrow P(Y=j|X) = \frac{e^{\beta_0^{(j)} + \beta_1^{(j)} \cdot X}}{1 + \sum_{i=0}^{k-2} e^{\beta_0^{(i)} + \beta_1^{(i)} \cdot X}} \quad , j=0, \dots, k-2$

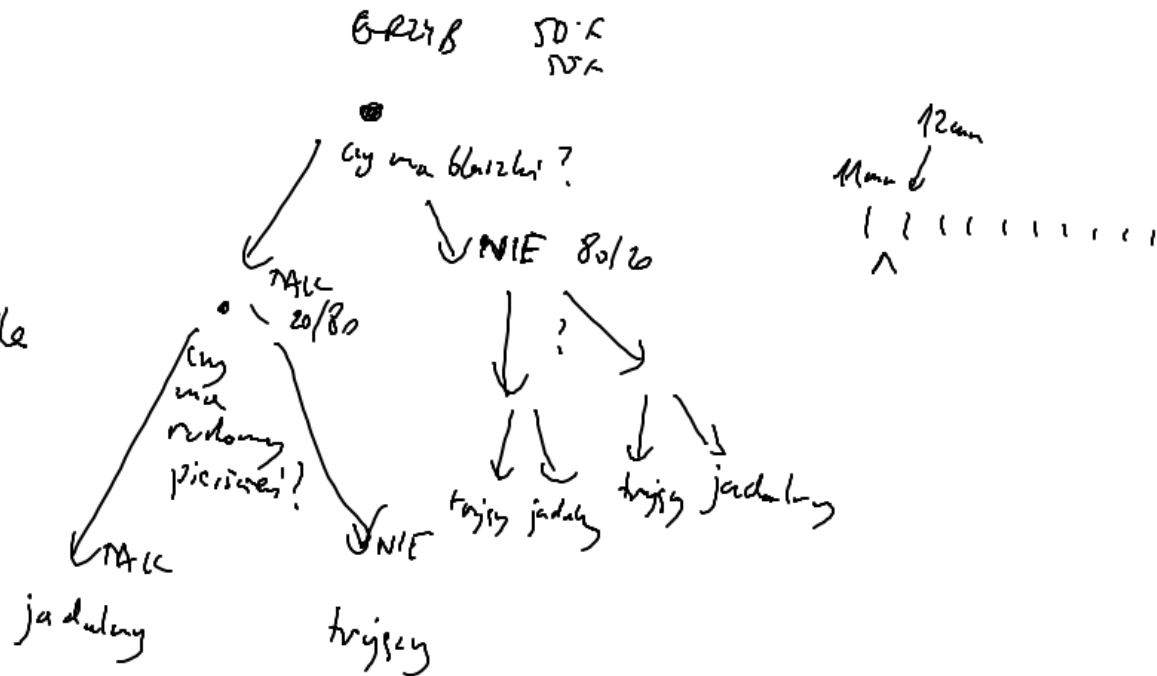
$P(Y=k-1|X) = \frac{1}{1 + \sum_{i=0}^{k-2} e^{\beta_0^{(i)} + \beta_1^{(i)} \cdot X}}$

softmax: $\begin{matrix} \rightarrow y_0 \\ \rightarrow \\ \rightarrow y_{k-1} \end{matrix}$ $\frac{e^{y_i}}{\sum_{i=0}^{k-1} e^{y_i}}$ (?)

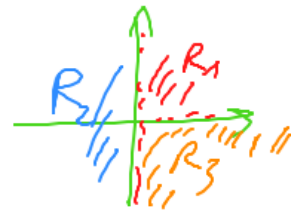
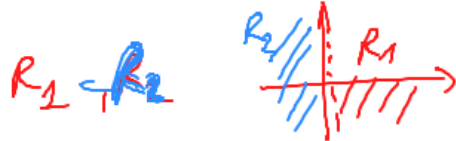
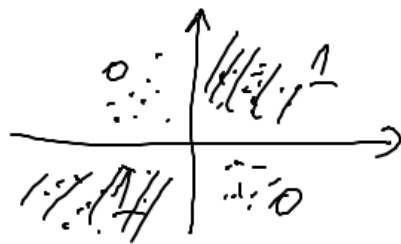
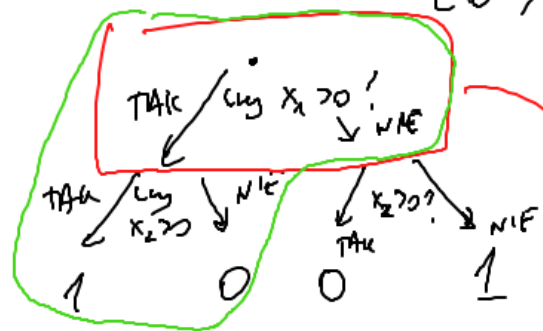
$\begin{matrix} \rightarrow 0 \\ \rightarrow 0 \\ \vdots \\ \rightarrow 0 \\ \rightarrow 0 \end{matrix} \left. \begin{matrix} \\ \\ \\ \\ \end{matrix} \right\} \begin{matrix} k-1 \\ \\ \\ \end{matrix}$ $\varphi(x) = \frac{1}{1+e^{-x}}$
 $\rightarrow 0 \quad k \quad \frac{1}{2}$

Drzewa decyzyjne

Drzewa ^{decyzyjne} zwykle buduje się za pomocą algorytmu rekursywnego:
 - Spisuje możliwych pytań ^{o danych węzła} i decyduje
 wybrai takie, które minimalizują
 pewną funkcję kosztu.



np. $X \in \mathbb{R}^2$, $Y = \begin{cases} 1, & X_1 X_2 > 0 \\ 0, & X_1 X_2 \leq 0 \end{cases}$



$X = (x_1, \dots, x_p) \rightsquigarrow Y$ klasa lub $\rightsquigarrow Y \in \mathbb{R}$

dane: $X^{(1)}, X^{(2)}, \dots, X^{(m)}$
 $y^{(1)}, y^{(2)}, \dots, y^{(m)}$

algorytm:

1) Porównujemy, ile X wstak podrobony na pełną linię obrotów R_1, \dots, R_J ;
 powyzkowo $J=1, R_1=X$.

2) Rozwiązujemy wszystkie $R_j, j=1, \dots, J$. Niech $R=R_J$.

Rozwiązujemy "wzrostkie pytania, które możemy zadać", tzn. dla każdego $k=1, \dots, p$

~~rozwiązujemy~~: a) jeśli X_k ma wartości w \mathbb{R} , to rozwiązujemy wszystkie $s \in \mathbb{R}$ i zbieramy postacie

$$R_1(j, s) = R \cap \{x_j < s\}, \quad R_2(j, s) = R \cap \{x_j \geq s\}$$

b) jeśli X_k ma wartości w pewnym zbiorze skończonym, to rozwiązujemy wszystkie jego podzbiory S i

$$R_1(j, s) = R \cap \{x_j \in S\}, \quad R_2(j, s) = R \cap \{x_j \notin S\}$$

3) wybieramy podzbiór $R : (j, S)$ tak, aby zminimalizować pewną funkcję kosztu.

Np. dla regresji:

$$\sum_{i: x^{(i)} \in R_1(j, S)} (y_i - \hat{y}_{R_1(j, S)})^2 + \sum_{i: x^{(i)} \in R_2(j, S)} (y_i - \hat{y}_{R_2(j, S)})^2 + \sum_{R \in \{R_1, \dots, R_j\} | R} \sum_{i: x^{(i)} \in P} (y_i - \hat{y}_P)^2$$

