

Билет 1

Смп. 1

Модель нейрона Машинного-Минуса, Перцептрон Розенблатта, Теорема Новикова, Полносвязные нейронные сети

- 1) $X = (x_1, \dots, x_n)$, нейрон $a(x)$ вычисляет n -сумму булеву функцию: $a(x) = \text{Heaviside}(\sum_{i=1}^n w_i x_i - w_0)$, w -веса ($w_i > 0$ -всд, иначе - тормозящий)
- 2) Виды элементов: S -сенсорный эл-т, A -ассоциативный, E -эл, ему совм. несколько S -элементов. A сдвигается, если количество сигналов S на его входе превышает некоторое значение θ . Значим сигналы идут от каждого элемента A идут на сумматор R с весов. w_j - весом A - R связи (веса S - A принимаются ± 1 или 1) значение порога θ случайное и неизменно, общ. вид ф-ции, реализующей R -эл-т, представим как перцептрон $x) = \text{sgn}(\sum_{i=1}^n w_i x_i - \theta)$
- 3) Т. Новикова и правило Хэбба: Пусть мы во пре-дметов $Y = \{-1, +1\}$, x -объект обуч. выборки $X^e = \{x_i, y_i\}_{i=1}^e$, $y_i = y_i^*(x_i) \in Y$ -класс. Алг. максим. имеем вид $a(x, w) = \text{sgn}(\langle x, w \rangle)$, ошибка вычисляется, если $\langle x, w \rangle y_i < 0$. Можно модифицировать веса: $w \rightarrow w + \eta x y$
- Собственно теорема: $X \subset \mathbb{R}^{n+1}$, $Y = \{-1, 1\}$, X^e мин. разделение, т.е. $\exists w^*$ и $\delta > 0$ $\forall x_i, w^* \cdot x_i > \delta \forall i = 1, \dots, e$. Тогда алгоритм может найти вектор весов, разд. train без ошибок за кон. число итераций \forall нач. приближ. w_0 и $\eta > 0$.

* FCN: input $\rightarrow x \rightarrow xw + b \rightarrow \sigma \rightarrow$ output
dense layer - мин. преобр. входн. данных (обуч. пер-мемприца w и вектор b): $x \mapsto xw + b$, $w \in \mathbb{R}^{d \times k}$, $x \in \mathbb{R}^d$, $b \in \mathbb{R}^k$
Слой делает d -мерные векторы k -мерные.
activation - нелин. преобр.

Билет 2

Аппрокс. теорема (Б.Д.). Пушкин, Колмогорова-Арнольда, Унбесенно

- 1) Пушкин: Визмеримой f на $(a, b]$ и $\forall \epsilon > 0 \exists f_\epsilon \in C[a, b]$, такое что $f = f_\epsilon$ везде кроме некоторого мин-ва меры ϵ .

т.е. $\forall x \in [a, b]: f(x) \neq f_2(x) \leq \epsilon$

Смп. 2

Теорема Коши-Вейерштрасса: Всякая непрерывная функция на отрезке $[a, b]$ принимает на нем все значения между $f(a)$ и $f(b)$.

т.е. $\forall f(x), x = (x_1, \dots, x_n) \exists$ предсказ. $\sum_{q=0}^{2^n} \Phi_q(\sum_{p=1}^n \Phi_{q,p}(x_p))$

Универсальная аппроксимация: $f = \sum_{q=0}^{2^n} \Phi(\sum_{p=1}^n \lambda_p \Phi(x_p + \gamma(q) + \theta))$, $\gamma, \lambda \in \mathbb{R}$, $\Phi \in C$, $\Phi: \mathbb{R} \rightarrow \mathbb{R}$, $\Phi[0, 2] \rightarrow \mathbb{R}$.

Условие: $Kf \in C^q(\mathbb{R})$, $f: \mathbb{R} \rightarrow \mathbb{R}$, $\forall \epsilon > 0 \exists N, w_1, \dots, w_N, b_1, \dots, b_N$, $\lambda_1, \dots, \lambda_N: |f - \sum_{i=1}^N \lambda_i \sigma(\langle x, w_i \rangle + b_i)| < \epsilon \forall x \in [0, 1]^m \in \mathbb{R}^m$

Бунет 3

Обучение нейронных сетей: stochastic gradient descent, back-propagation.

Минимизация потерь: $Q(w) = \sum h(w, x_i, y_i) \rightarrow \min$

$w = w - \gamma h_i(w)$, $\nabla h_i(w) = \left(\frac{\partial h_i(w)}{\partial w^k} \right)_{k=1}^m$, $Q = (1-\lambda)Q + \lambda h_i(w)$

Back drop

Дано $(x_i, y_i)_{i=1}^n$, сеть $(H_e)_{e=1}^n$, параметры γ, λ . Учим сеть

Берем все веса $w = (w^1, \dots, w^n)$

1) random $(x_i) \in X^e$, $e=1, \dots, h$, $n=1, \dots, H_e$

forward: $x_{in}^e = \sigma_n^e \left(\sum_{k=0}^{H_{e-1}} w_{kn}^e x_{in}^{e-1} \right)$, $p_{ni}^h = \frac{\partial h_i(w)}{\partial x_n^w}$
 $z_{in}^e = \left(\sigma_n^e \right)' \left(\sum_{k=0}^{H_{e-1}} w_{kn}^e x_{in}^{e-1} \right)$

Backward: $p_{in}^{e-1} = \sum_{h=0}^{H_e} p_{in}^e z_{in}^e w_{kn}^e$

grad step for all $e=1, \dots, h$, $k=0, \dots, H_{e-1}$, $n=1, \dots, H_e$

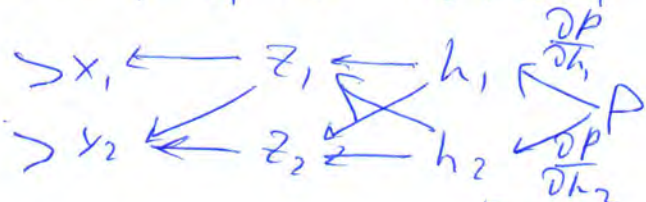
$w_{kn}^e = w_{kn}^e - \gamma p_{in}^e z_{in}^e x_{in}^{e-1}$

Почему Q и/или веса нестабильны

пример:

③ $\frac{\partial P}{\partial h_1}, \frac{\partial P}{\partial h_2}$ ② $\frac{\partial P}{\partial z_i} = \frac{\partial P}{\partial h_1} \frac{\partial h_1}{\partial z_i} = \frac{\partial P}{\partial h_1} \frac{\partial h_2}{\partial z_i}$; $\frac{\partial P}{\partial z_2} = \frac{\partial P}{\partial h_1} \frac{\partial h_1}{\partial z_2} + \frac{\partial P}{\partial h_2} \frac{\partial h_2}{\partial z_2}$

① $\frac{\partial P}{\partial x_1} = \frac{\partial P}{\partial h_1} \frac{\partial h_1}{\partial z_1} \frac{\partial z_1}{\partial x_1} + \frac{\partial P}{\partial h_2} \frac{\partial h_2}{\partial z_1} \frac{\partial z_1}{\partial x_1} + \frac{\partial P}{\partial h_1} \frac{\partial h_1}{\partial z_2} \frac{\partial z_2}{\partial x_1} + \frac{\partial P}{\partial h_2} \frac{\partial h_2}{\partial z_2} \frac{\partial z_2}{\partial x_1}$



Бунет 4

Метод обучения Adaptive momentum.

Метод Drop-out new neurons

① $v_0 = 0$, $v_{k+1} = \beta_1 v_k + (1-\beta_1) \nabla f(x_k)$, $x_{k+1} = x_k - \frac{e_2}{\sqrt{G_{k+1}} + \epsilon} v_{k+1}$, $G_{k+1} = \beta_2 G_k + (1-\beta_2) (\nabla f(x_k))^2$
 $\beta_1 = 0.5$, $\beta_2 = 0.95$, $\epsilon = 1e-8$, $f(x) = \sum h(x, y_i)$

Нормы и momentum - RMS prop.

Задача 3



RMS Prop: $G_{k+1} = \sqrt{G_k + (1-\beta)(\Delta f(x_k)^2)}$

$$x_{k+1} = x_k - \frac{\epsilon_2}{\sqrt{G_{k+1}} + \epsilon} \Delta f(x_k)$$

Dropout

при град. карт $h_i(w) \rightarrow m_i$ единицы n-ой нейрон
l-го слоя с вер. p_c

$$x_{n_i}^l = \xi_n^l \sigma_n^l / (\sum w_{kn}^l x_{k_i}^{l-1}) \quad P_2(\xi_n^l = 0) = p_c$$

Вариант с $(1-p_c)$

$$\text{на одн. } x_{n_i}^l = \frac{1}{1-p_c} \xi_n^l \sigma_n^l (\sum w_{kn}^l x_{k_i}^{l-1})$$

$$\text{на вариациях } x_{n_i}^l = \sigma_n^l (\sum w_{kn}^l x_{k_i}^{l-1})$$

$$\text{в среднем } h_2 \text{ gradstep } w = w(1-\beta) - \beta \frac{1}{1-p_c} \xi_n^l h_i'(w)$$

Задача 5

Основные способы регуляризации:
Л1-норма, skip-connection, drop-out, batch norm,
инварианты.

1) Дано 2 матрицы $A(n_x \times n_y)$ и $B(m_x \times m_y)$, $C = A * B$. C shape:
 $= (n_x - m_x + 1), (n_y - m_y + 1)$, т.е. $C_{i,j} = \sum_{k=0}^{m_x-1} \sum_{l=0}^{m_y-1} A_{i+k,j+l} B_{k,l}$

2) Skip-connection решает проблему затухания град,
передает инфор. из нижних слоев в верхние:

$$\text{Layer 1} \xrightarrow{\text{skip-connection}} \text{Layer 2} \rightarrow \Phi \rightarrow \dots$$

3) Dropout - случайное выключение некоторых из нейронов
для борьбы с переобучением и уменьшения обобщающей способ-
ности

4) Batch norm - стабилизирует дисперсию от входа. Даны
независимые данные x_i с $\mu = 0, \sigma^2 = 1$

$$\text{Batch} := \{x_1, \dots, x_m\} \quad \mu_B = \frac{1}{m} \sum x_i \quad \sigma_B^2 = \frac{1}{m} \sum (x_i - \mu_B)^2$$

$$\hat{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} \quad y_i = f(\hat{x}_i) \quad \beta = \text{Batch Norm}(x_i) - \text{отслеживание и обновление}$$

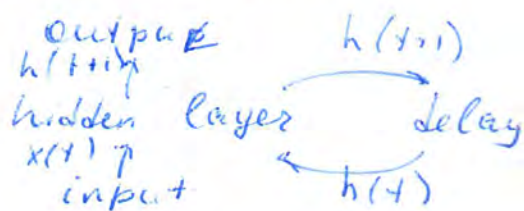
5) ~~онлайн~~ инварианты помогают моделировать сложные
ср-знач.

Рекуррентные нейронные сети LSTM

RNN-сети с циклами, позволяющие им обрабатывать последовательности, где будущее зависит от прошлого. Они решают задачу Backprop Through Time (BPTT).

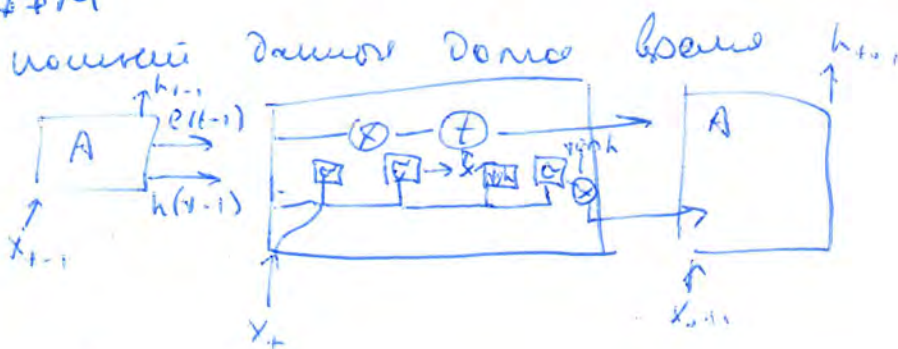


Схема с задержкой в скрытом слое:

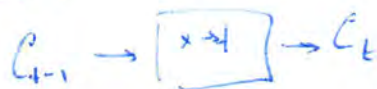


- Виды RNN:
- 1) 1 вход, много выходов: для измерения аудио
 - 2) много входов и выходов: для оценки последовательностей
 - 3) много входов и выходов: для перевода

LSTM



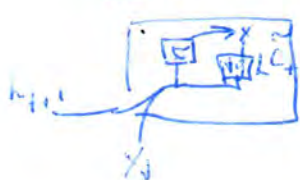
шаг 1



шаг 2



шаг 3



шаг 4



$$C_t = f_t C_{t-1} + i_t \tilde{C}_t$$

$$\tilde{C}_t = \sigma(W_i [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_c [h_{t-1}, x_t] + b_c)$$

$$h_t = \sigma(W_o [h_{t-1}, x_t] + b_o) \tanh(C_t)$$

kernel trick

loss: $h(y, g(x)) = \begin{cases} 0, & |y - g(x)| \leq \epsilon \\ 1, & |y - g(x)| > \epsilon \end{cases}, \epsilon > 0$

$$\text{loss: } q(x_i) = (w, x_i) - w_0 - y_i \quad \forall (x_i, y_i) \text{ where } |z|_2 = \max(0, |z| - \epsilon)$$
$$Q_2(a, X) = \sum_{i=1}^n |(w, x_i) - w_0 - y_i|_2 + \gamma (w, w)^2 \rightarrow \min_{w, w_0}$$

Введем переменные ξ_i^+ и ξ_i^- - дроб. потери при завышении или занижении оценки.

$$\xi_i^+ = (a(x_i) - y_i - \varepsilon)_+, \quad \xi_i^- = (1 - a(x_i) + y_i - \varepsilon)_-, \quad i = 1, \dots, \ell$$

Задана минимизирующая

$$0,5(w, w)^2 + \frac{1}{2\beta} \sum_{i=1}^n (\xi_i' + \xi_i) \rightarrow \min_{w, w_0, \xi_1, \xi_2}$$

$$(w, v_i) \cdot w_i \geq w_i - \epsilon_i$$

$$\xi_i \geq 0, \xi_i^* \geq 0$$

Будем считать двойственную задачу, макс. произв. значением
отрост $k(x_i, x_j)$, d_i^+ , d_i^- - двойств. перем.

$$h(d^+, d^-) = -\varepsilon \sum_{i=1}^p (d_i^+ + d_i^-) + \sum_{i=1}^p (d_i^- - d_i^+) y_i - \frac{1}{\tau} \sum_{i=1}^p (d_i^- - d_i^+) (d_s^- - d_s^+) \cdot k(y_i, x_s) \rightarrow_{d_i^+, d_i^-} \min$$

$$\sum (d_i^+ + d_i^-) \rightarrow 0$$

В рез. все у; доходит на 5 мин

$$1) |a(x_i) - y_i| \leq \epsilon \quad d_i = d_i^+ - \xi_i^+ = \xi_i^- = 0$$

2) $Q(x_i) = y_i + \varepsilon$ $0 < d_i^+ < \varepsilon$, $d_i^- = \xi_i^+ = \xi_i^- = 0$

3) $q(x_i) = y_i$, $0 \leq i \leq n$, $d_1^{(1)} = \{y_i\}_{i=0}^n = 0$
 4) $q(x_i) = y_i$, $0 \leq i \leq n$, $d_1^{(1)} = \{y_i\}_{i=0}^n = 0$

4) $a(x_i) > y_i + \varepsilon$ $d_i^+ = \varepsilon$, $d_i^- = 0$, $\xi_i^+ = a(x_i) - y_i - \varepsilon$, $\xi_i^- = 0$
 5) $a(x_i) < y_i - \varepsilon$ $d_i^+ = 0$, $d_i^- = \varepsilon$, $\xi_i^+ = 0$, $\xi_i^- = y_i - a(x_i) - \varepsilon$

$$3) a(x_i) < y_i, \delta = \varepsilon \quad d_i^+ = 0, \quad d_i^- = \varepsilon, \quad \xi_i^+ = a(x_i) - y_i - \varepsilon, \quad \xi_i^- = 0$$

7.5. оператор U в l_2 задан в базисе $\{e_k\}$ формулой $Ue_k = e_k + \frac{1}{k}e_1$.
 Наблюдать: $(Ux, x) = \|x\|^2 + \frac{1}{2}\|x\|^2 = \frac{3}{2}\|x\|^2$.
 У U нет собственных значений.

$$W_0: (w, x_1) - w_0 = \begin{cases} y + \varepsilon, & \xi_i \in \mathcal{Z} \\ y - \varepsilon, & \xi_i \in \mathcal{S} \end{cases}$$

Маскировка

$\varphi: \mathbb{R}^n \rightarrow H$, H - вейт. пр-во. сразу сведет задачу к мин. раздел. выбора
раздел. φ -числ $f(x) = (w, \varphi(x)) + b$, $w = \sum d_i y_i \varphi(x_i)$, где d_i зависят от y_i и
 $(\varphi(x_i), \varphi(x_j))$, преем. $k(x, y) = (\varphi(x), \varphi(y))$, k неотр. опр. и симм.

Вспомогат. мин. прбл. $f(x)$ на W . $f(x_i) > 0 \forall x_i \in W, f(x_i) < 0$

$$x \in W, y_i = \begin{cases} 1, & x \in W_i \\ -1, & x \in W_j \end{cases}, \quad \text{and } y_i (w_i x_i + b) \geq 0$$

Рассмотрим минимизацию $f(w, x) + \theta = 0$ по T -но θ (вектор-параметр) задана суб. $h(w, \theta, d) = 0$, $s(w, w) = \sum \lambda_i (y_i / (w, x_i) + \theta) - 1 \rightarrow \min_w \max_d$

$\lambda_i \geq 0, d_i (y_i (w, x_i) + b - 1) = 0$, иначе $d_i = 0$ или $y_i ((w, x_i) + b - 1) = 0$, y_i ген. сущ. одной нормы:

$$\begin{cases} \frac{\partial h}{\partial w_s} = w_s - \sum d_i y_i x_{is} = 0 \\ \frac{\partial h}{\partial b} = \sum d_i y_i = 0 \end{cases} \Rightarrow w = \sum d_i y_i x_i \text{ и } \sum d_i y_i = 0$$

$$h(w, b, d) = \sum d_i - 0.5 \sum d_i d_j y_i y_j (x_i, x_j) = \sum d_i - 0.5 \| \sum d_i y_i x_i \|^2 \text{ т.е.}$$

Вводим и помечу при этом $\Phi(d) = \sum d_i - 0.5 \| \sum d_i y_i x_i \|^2$
 kernel trick
 $w = \sum d_i y_i \phi(x_i)$ $\arg \max f(d_1, \dots, d_n) = \sum d_i - \frac{1}{2} \sum_{i,j} \sum y_i d_i k(x_i, x_j) y_j d_j$
 $\sum d_i y_i = 0, 0 \leq d_i \leq \frac{1}{n \lambda} k_i, \phi(x_i)$ решим на граде в прободр. пр-во
 тогда $b = w^T \phi(x_i) - y_i = \sum y_j d_j k(x_j, x_i) - y_i; z \mapsto \text{sgn}(w^T \phi(z) - b)$

Букет 6

Метод и-близких соседей Оценка
 минимизация ошибки классификации методом kNN
 (метод Ховера-Харфа)

Ран: $\text{train } X^m = \{(x_i, y_i)\}_{i=1}^m$, метрика $\rho(x, x')$ для произв. объектов,
 классифицируем объект x в порядке возр. до $u: \rho(u, x, u) \leq \rho(u, x_m, u)$
 x_i и i -й сосед и тогда алгоритм будет задавать класс
 $a(u) = \arg \max \sum [y_i, u = y] w(i, u)$ где $w(i, u)$ весовый коэффициент
 веса i -го соседа для класса u . Для классификации kNN $w(i, u) =$
 $= [i \leq k]$ при миним. $w(i, u)$ на-во м.б. помин, тогда приме-
 нем одр. пример: $w(i, u) = \frac{1}{k} \left(\frac{\rho(u, x_i, u)}{\rho(u, x_k, u)} \right)$ - пороговое сн
 и-во Ховера-Харфа: верхняя граница ошибки $R^* \leq R_{kNN} \leq R^*$
 $= (1 - \frac{R^*}{M-1})$, R^* - ошибка в R_{kNN} - асимпт. граница, M - число классов

Букет 9

Задача классификации Основание
 метрики (confusion matrix, ROC AUC) и т.д. Оценка
 обобщ. способности модели.

~~Задача~~ Рассм. зад. бинар. классификации $x \rightarrow y, x \in \mathbb{R}^n, y \in \{-1, +1\}$
 $X^m = \{(x_i, y_i)\}_{i=1}^m$, число $a = \text{sgn}(\langle w, x \rangle - w_0)$, ищем параметр w , мин.
 раб. $R(x^m, w, w_0) = \sum [a \neq y_i] = \sum [\langle w, x_i \rangle - w_0 > 0 \wedge y_i < 0]$
 $\arg \min_w = \frac{1}{m} \sum_{i=1}^m [a(x_i) = y_i] = \frac{TP + TN}{FP + FN + TP + TN}$ Матр. ошибок:

При бинарном w_0 число x_i т.т. $a(x_i) = -1$
 добавл. dy-координ, тогда $\text{loss} = d_y [a(x_i) \neq y_i]$
 ROC-AUC по x : $\text{TPR} = \frac{\sum [y_i = -1] \sum [a(x_i) = +1]}{\sum [y_i = -1]}$, по y : $\text{TPR} = \frac{\sum [y_i = +1] \sum [a(x_i) = +1]}{\sum [y_i = +1]}$

Оценки обобщ. способности: $X^e = \{(x_i, y_i)\}_{i=1}^e$
 $y_i = y^*(x_i)$ $A_+ = \{a, x \rightarrow y\}$, $A_- = \{x \rightarrow y\} \xrightarrow{\text{норм.}} A_+$
 $h(a, x) = \text{loss } Q(a, x^e) = \frac{1}{e} \sum \mathcal{L}(a, x_i)$ - ф-ция потерь

Самый common: fold-out $X^e = X_n^e \cup X_{n+1}^e$, $n = 1, \dots, N$

$$\text{cross-val}(u, X^e) = \frac{1}{|N|} \sum_{n=1}^N Q_n(X_n^e, X_{n+1}^e)$$