# Билет 1

Модель нейрона Маккалока-Питтса.
Перцептрон Розенблатта. Теорема Новикова.
Полносвязные нейронные сети

1) $X = (x_1, ... x_n)$, нейрон $a(x)$ вычисляет n-арную булеву функцию: $a(x) = Heaviside(\sum_{s=1}^n w_s x^s - w_0)$, $w$ - веса ($w_s > 0$ - возб., иначе - тормозящий)

2) Виды элементов: $S$ - сенсорный эл-т, $A$ - ассоциативный эл-т, $R$ - ... $A$ стимулируется, если количество сигналов $S$ на его входе превысило некоторое значение $\Theta$. Затем сигналы ~~идут~~ от возбуждённых $A$ идут на сумматор $R$ с нар. $w_s$ - весом $A$-$R$ связи (веса $S$-$A$ принимают числа -1, 0 или 1) значения порогов $A$ рандомны и ненулевые, общ. вид ф-ции, реализуемой $R$-эл-том, представим как $perceptron(x) = sgn(\sum_{s=1}^n w_s x_s - \Theta)$

3) $\underline{Th}$ Новикова и правило Хэбба: Пусть ли во пределах $Y = \{-1, +1\}$, $x$ - объект обуч. выборки $X^\ell = \{x_j, y_j\}_{j=1}^\ell$, $y_j = y_j^*(x_j) \in Y$ - классов. Алг. класиф. имеет вид $a(x, w) = sgn(\langle x, w \rangle)$, ошибка выделена если знак $\langle , \rangle$ не совпал его данному классу, т.е. $\langle x w \rangle y < 0$. Можно модифицировать веса: $w \to w + y x y$

Собственно теорема: $X = \mathbb{R}^{n+1}$, $Y = \{-1, 1\}$, выборка $X^\ell$ лин. разделима. Т.е. $\exists w^*$ и $\delta > 0$ ч. ч. $\langle x_s, w^* \rangle y_j > \delta$ $\forall_{j=1,...\ell}$. Тогда алгоритм может найти вектор весов, раб. train без ошибок за кон. число итераций $\forall$ нач. прибл. $w_0$ и $\ell > 0$.

\# FCN: input $\to x \to x_w + b \to \sigma \to$ output
dense layer - лин. преобр. входн. данных (обуч. пар-метрица $w$ и вектор $b$): $x \mapsto xW + b$, $W \subset \mathbb{R}^{d \times k}$, $x \in \mathbb{R}^d$, $b \in \mathbb{R}^k$. Слой делает из d-мерных векторов k-мерные. activation - нелин. преобр.

# Билет 2

Аппрокс. теоремы (б.д.):
Лузина, Колмогорова-Арнольда, Цыбенко

1) Лузина: Измеримой $f$ на $[a, b]$ и $\forall \varepsilon > 0$ $\exists f_\varepsilon \in C[a, b]$, такая что $f = f_\varepsilon$ всюду кроме некоторого мн-ва мер $\varepsilon$.

Т.е. $\forall ? x \in [a,b]: f(x) \neq f_\varepsilon(x)\} < \varepsilon$

① Теорема Колмогорова - Арнольда: всякую непр. ф-цию представима в виде композиции непр. ф-ций одной пер.
т.е. $\forall f(x), \; x = (x_1, \ldots x_n) \; \exists$ представл. $\sum\limits_{q=0}^{2n} \Phi_q\left(\sum\limits_{p=1}^{n} \varphi_{q,p}(x_p)\right)$

Усиление Шпрехера: $f = \sum\limits_{q=0}^{2n} \Phi\left(\sum\limits_{p=1}^{n} \lambda_p \varphi(x_p + \eta q) + q\right), \; \eta, \lambda \in \mathbb{R}, \; \Phi \in C,$
$\Phi \colon R \to R, \; \varphi[0,2] \to R.$

③ Цыбенко: $\forall f \in C^0(R), \; f \colon R \to R, \; \forall \varepsilon > 0 \; \exists N, \; w_1, \ldots w_N, \; b_1, \ldots b_N,$
$\lambda_1, \ldots \lambda_N \colon |f - \sum\limits_{i=1}^{N} \lambda_i \sigma(\langle x, w_i \rangle + b_i)| < \varepsilon \; \forall x \in [0,1]^{n} \in \mathbb{R}^m$

## Билет 3

### Обучение искусств. нейронов:
### stochastic gradient descent, back-propagation.

Минимизируем: $Q(w) = \sum h(w, x_i, y_i) \to \min$

$w = w - \gamma h_i(w), \quad \nabla h_i(w) = \left(\dfrac{\delta h_i(w^t)}{\delta w^k}\right)_{k=1}^{m}, \quad Q = (1-\lambda) Q + \lambda h_i(w)$

Backdrop

Дано $(1 x_i, y_i)_{i=1}^{\ell}$ сети $(He)_{\ell=1}^{h}$, парам. $\gamma, \lambda$. Ищем вектор весов всех слоев $w = (w^1, \ldots w^n)$

1) random $(x_i) \in X^\ell$, $\ell = 1, \ldots h$, $n = 1, \ldots He$

forward: $X_{in}^{\ell} = \sigma_n^{\ell}\left(\sum\limits_{k=0}^{He_{-1}} w_{kn}^{i} x_{in}^{\ell-1}\right)$

$z_{in}^{\ell} = (\sigma_n^{\ell})' \left(\sum\limits_{k=0}^{He_i} w_{kn}^{j} x_{ik}^{\ell-1}\right) \qquad p_{ni}^{h} = \dfrac{\partial h_i(w)}{\partial x_n^{w}}$

Backward: $p_{in}^{\ell_{-1}} = \sum\limits_{h=0}^{He} p_{in}^{\ell} z_{in}^{\ell} w_{kn}^{\ell}$

grad step for all $\ell = 1, \ldots h$, $k = 0, \ldots He_{-1}$, $n = 1, \ldots He$

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

Пока $Q$ и/или весов unstable

Пример:
③ $\dfrac{\partial P}{\partial h_1}, \; \dfrac{\partial P}{\partial h_2}$.

② $\dfrac{\partial P}{\partial z_i} = \dfrac{\partial P}{\partial h_1}\dfrac{\partial h}{\partial z_i} = \dfrac{\partial P}{\partial h_1}\dfrac{\partial h_2}{\partial z_i}; \quad \dfrac{\partial P}{\partial z_2} = \dfrac{\partial P}{\partial h_1}\dfrac{d h_1}{\partial z_2} + \dfrac{\partial P}{\partial h_2}\dfrac{\partial h_2}{\partial z_2}$

① $\dfrac{\partial P}{\partial x_1} = \dfrac{\partial P}{\partial h_1}\dfrac{\partial h_1}{\partial z_1}\dfrac{\partial z_1}{\partial y_1} + \dfrac{\partial P}{\partial h_2}\dfrac{\partial h_2}{\partial z_1}\dfrac{\partial z_1}{\partial y_1} + \dfrac{\partial P}{\partial h_1}\dfrac{\partial h_1}{\partial z_2}\dfrac{\partial z_2}{\partial x_1} + \dfrac{\partial P}{\partial h_2}\dfrac{\partial h_2}{\partial z_1}\dfrac{\partial z_2}{\partial x_1}$
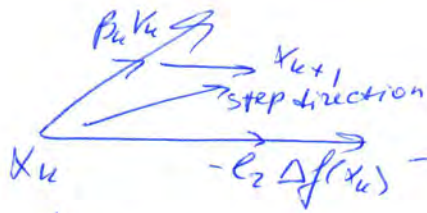
$$> x_1 \longleftarrow z_1 \longleftarrow h_1 \xleftarrow{\frac{\partial P}{\partial h_1}}$$
$$> y_2 \longleftarrow z_2 \longleftarrow h_2 \xrightarrow{\frac{\partial P}{\partial h_2}} P$$

## Билет 4

### Метод обучения Adaptive momentum.
### Метод Prop-ont и его реализация

① $v_0 = 0$
$v_{k+1} = \beta_1 v_k + (1-\beta_1) \overrightarrow{\nabla} f(v_k)$
$G_{k+1} = \beta_2 G_k + (1-\beta_2)(\overrightarrow{\nabla} f(x_k))^2$

$x_{k+1} = x_k - \dfrac{\ell_2}{\sqrt{G_{k+1}} + \varepsilon} v_{k+1}$

$\beta_1 = 0.5$
$\beta_2 = 0.35$
$\varepsilon = 1e-8$
$f(x) = \sum h(x, y_i)$

Получим из momentum и RMS prop.

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

$$x_{k+1} = x_k - \frac{l_2}{\sqrt{G_{k+1}} + \varepsilon}\,\nabla F(x_k)$$

dropout

При град. шаге $h_i(w) \to \min\limits_{w}$ выключаем $n$-ой нейрон
$l$-го слоя с вер. $p_e$

$$x^{\ell}_{n_i} = \xi^{\ell}_n\, \sigma^{\ell}_n\left(\Sigma\, w^{\ell}_{kn}\, A^{\ell-1}_{k_i}\right) \quad P_2(\xi^{\ell}_n = 0) = p_e$$

выкл. с $(1-p_e)$

на обуч. $x^{\ell}_{n_i} = \frac{1}{1-p_e}\,\xi^{\ell}_n\,\sigma^{\ell}_n\left(\Sigma\, w^{\ell}_{kn}\, x^{\ell-1}_{k_i}\right)$

на валидации $x^{\ell}_{h_i} = \sigma^{\ell}_n\left(\Sigma\, w^{\ell}_{kn}\, x^{\ell-1}_{k_i}\right)$

с учётом $l_2$ gradstep $w = w(1-\gamma\lambda) - l\frac{1}{1-p_e}\,\xi^{\ell}_n\, h'_i(w)$

## Билет 5

Основные слои искусств. нейронов:
Свертка, skip-connection, drop-out, batch norm,
нелинейности.

1) дано 2 тензор. $A\,(n_x \times n_y)$ и $B\,(m_x \times m_y)$ $C = A * B$. C. shape =
$= (n_x - m_x + 1),\ (n_y - m_y + 1)$, т. е. $C_{i,j} \sum\limits_{u=0}^{m_y - 1} \sum\limits_{v \geq 0}^{m_y - 1} A_{i+u,\,j+v}\, B_{u,v}$

2) Skip-connection решает проблему затухания grad,
передача инф. из нижних слоёв в верхние:

$$\to\ layer\ 1 \to activation \to layer\ 2 \to \oplus \to \dots$$

(с пометкой skip-connection над стрелкой)

3) Dropout - случайное выключение некоторого числа нейронов.
Они борятся с переобучением и улучшают обобщающую спо-
собность.

4) Batch norm - считает зависимость от входн. данных
некоторым слоям данные данные с $\mathbb{E} = 0$ $var = 1$
$Batch := \{x_1 \dots x_m\}$ $M_B = \frac{1}{m}\Sigma x_i$ $\sigma^2_B = \frac{1}{m}\Sigma(x_i - M_B)^2$
$\hat{x}_i = \frac{x_i - M_B}{\sqrt{\sigma^2_B + \varepsilon}}$ $y_i = \gamma\hat{x}_i + \beta = Batch\ Norm(x_i)$ - сжатие и сдвиг.

5) ~~нелин~~ нелинейности помогает моделировать сложные
ср-нии.

# Билет 6

## Рекуррентные нейронные сети.
## LSTM.

RNN-сети с учителем, подходящие для обработки последовательностей. Для обучения применяется BPTT-Back error through time. обуч. он так:



Схема слоя RNN и развернутая схема.

Схема с задержкой в скрытом слое:

output
h(t+1)

hidden layer    delay
x(t) ↑

input    h(t)

Виды RNN:
1) 1 вход, много выходов; для итерации аудио
2) Много входов, 1 выход: для оценки исполняемостей
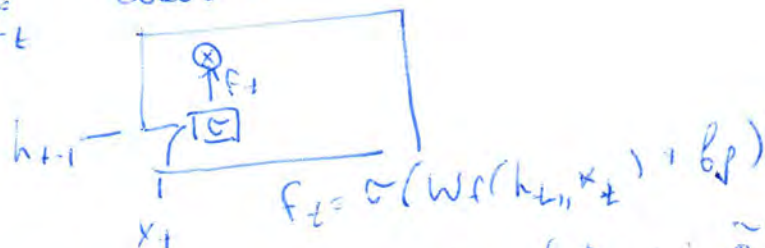3) много входов и выходов: для перевода

## LSTM

помнит данные долгое время



шаг 1

$C_{t-1} \rightarrow \boxed{\times \ +} \rightarrow C_t$

шаг 2

$f_t = \sigma(W_f[h_{t-1}, x_t] + b_f)$

шаг 3

шаг 4

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

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

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

$h_t = \sigma(W_0[h_{t-1}, x_t] + b_0) tgh(C_t)$

# Билет 7

## Метод опорных векторов.
## Kernel trick.

Ставим задачу регрессии - дано train $X = \{(x_i, y_i)\}_{i=1}^{\ell}$, ищем $g(x)$ - аппроксимацию выборки лучшим образом.

loss: $h(y, g(x)) = \begin{cases} 0, & |y - g(x)| \le \varepsilon \\ |y - g(x)| - \varepsilon, & (y - g(x)) \ge \varepsilon \end{cases}$, $\varepsilon > 0$

Ищем решение в линейном виде: $f(x) = (w, x) - w_0$

loss: $a(x_i) = |(w, x_i) - w_0 - y_i|_\varepsilon$ $\forall (x_i, y_i)$ где $|z|_\varepsilon = \max(0, |z| - \varepsilon)$

функционал, потеръ

$Q_\varepsilon(a, X) = \sum_{i=1}^{\ell} |(w, x_i) - w_0 - y_i|_\varepsilon + \gamma (w, w)^2 \to \min_{w, w_0}$

Введём переменные $\xi^+$ и $\xi^-$ - диаг. потороъ = потеръ при завышенном или заниженном ответе.

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

Задача минимизации:

$\begin{cases} 0.5(w, w)^2 + \frac{1}{2\gamma} \sum_{i=1}^{\ell} (\xi_i^+ + \xi_i^-) \to \min_{w, w_0, \xi_i^+, \xi_i^-} \\ (w, x_i) - w_0 \le y_i + \varepsilon + \xi_i^+ \\ (w, x_i) \cdot w_0 \ge y_i - \varepsilon - \xi_i^- \\ \xi_i^- \ge 0, \xi_i^+ \ge 0 \end{cases}$

Будем решать двойственную задачу, скал. произв. заменим скор. $k(x_i, x_j)$, $d_i^+, d_i^-$ - двойств. перем.

$\begin{cases} h(d^+, d^-) = -\varepsilon \sum_{i=1}^{\ell}(d_i^+ + d_i^-) + \sum_{i=1}^{\ell}(d_i^- - d_i^+) y_i - \frac{1}{2} \sum (d_i^- - d_i^+)(d_j^- - d_j^+) \cdot k(x_i, x_j) \to \max_{d_i^+, d_i^-} \\ 0 \le d_i^+ \le c = \frac{1}{2\gamma}, \quad 0 \le d_i^- \le c \\ \sum (d_i^+ + d_i^-) \to 0 \end{cases}$

В рез. все $x_i$ делятся на 5 типов

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

2-5 - опорные, уч. в опр. весов. У-ние регрессии $\sum (d_i^- - d_i^+) k(x_i, x) - w_0$

$w_0$: $(w, x_i) - w_0 = \begin{cases} y + \varepsilon, & \xi_i \in 2 \\ y - \varepsilon, & \xi_i \in 3 \end{cases}$

Классификация

$\varphi : \mathbb{R}^n \to H$. $H$ - вект. пр-во. сразу сведем задачу к лин. раздел. выборке. раздел. ф-ция $f(x) = (w, \varphi(x)) + b$ $w = \sum d_i y_i \varphi(x_i)$, где $d_i$ зависит от $y_i$ и $(\varphi(x_i), \varphi(x_j))$, рассм. $k(x, y) = (\varphi(x), \varphi(y))$, $k$ неотр. опр. и симм.

В случае лин. разд. выборки. ищем $f(x)$ т.ч. $f(x_i) > 0$ $\forall x_i \in \overline{w}$, $f(x_i) < 0$ $\forall x \in \overline{w}$, $y_i = \begin{cases} 1, & x \in \overline{w} \\ -1, & x \in \overline{\overline{w}} \end{cases}$, т.е. $y_i((w, x_i) + b) \ge 0$

Разд. гиперповерхность $(w, x) + b = 0$, по т-ме Куна-Таккера задача эив. $h(w, b, d) = 0.5(w, w) - \sum d_i(y_i((w, x_i) + b) - 1) \to \min_{w, b} \max_{d}$

$\lambda_i \ge 0$, $d_i(y_i(\langle w, x_i\rangle + b) - 1) = 0$, откуда $d_i = 0$ или <inline_ref>Стр. 6</inline_ref>

$y_i(\langle w, x_i\rangle + b) - 1 = 0$, ищу усл. сущ. седловой точки:

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

$L(w, b, d) = \sum d_i - 0.5 \sum d_i d_j y_i y_j \langle x_i, x_j\rangle = \sum d_i - 0.5 \|\sum d_i y_i x_i\|^2$ т.е.

сводим к поиску при помощи проекции $\Phi(d) = \sum d_i - 0.5\|\sum d_i y_i x_i\|^2$

kernel trick

$w = \sum d_i y_i \varphi(x_i)$   $\arg\max f(d_1 \ldots d_n) = \sum d_i - \frac{1}{2}\sum_i \sum_j y_i d_i k(x_i, x_j) y_j d_j$

$\sum d_i y_i = 0$,  $0 \le d_i \le \frac{1}{2n\lambda}$ $\forall i$,  $\varphi(x_i)$ лежит на гран. в преобр. пр-ве

(пусть $b = w^T \varphi(x_i) - y_i = \sum y_j d_j k(x_j, x_i) - y_i$;  $z \mapsto \operatorname{sgn}(w^T \varphi(z) - b)$

## Билет 8

### Метод K-ближайших соседей. Оценки оптимальности классификации методом kNN

#### (нер-во Ковера-Харта)

Для $X^m = \{(x_i, y_i)\}_{i=1}^m$, метрика $\rho(x, x')$ для процев. объекта, располагаем объекты $x_i$ в порядке возр. до и $\rho(u, x_{i;u}) \le \rho(u, x_{m;u})$

$x_{i;u}$ — $i$-й сосед и тогда алгоритм будем задавать поиск

$a(u) = \arg\max\limits_{y \in Y} \sum [y_{i;u} = y] w(i, u)$ где $(i, u)$ оценивает степень важности $i$-го соседа для класс. $u$. Для классического kNN $w(i, u) =$

$= [i \le k]$ при линейном $w(i, u)$ нар-во м. б. многими, тогда прим-

няют ядро. Пример: $w(i, u) = k\left(\dfrac{\rho(u, x_{i;u})}{h}\right)$ — парзеновское окно

H-во Ковера-Харта: верхняя граница частота ошибки $R^* \le R_{kNN} \le R^*$

$= (2 - \dfrac{MR^*}{M-1})$, $R^*$ — частота ошибок БК $R_{kNN}$ — асимпт. частота, $M$ — число классов

## Билет 9

### Задача классификации. Основные

#### метрики (confusion matrix, ROC AUC) и их оценки

#### Общая способности модели.

Расем. зад. бин. классификации $X \to Y$ $X = \mathbb{R}^n$, $Y = \{-1, +1\}$

$X^m = (x_i, y_i)_{i=1}^m$, класс. $a = \operatorname{sgn}(\langle w, x\rangle - w_0)$, ищем параметры $w$, мин-

риск. $R(x^m, w, w_0) = \sum [a \ne y_i] = \sum [(\langle w, x_i\rangle - w_0) y_i < 0]$

$\text{accuracy} = \frac{1}{\ell} \sum_{i=1}^{\ell} [a(x_i) = y_i] = \dfrac{TP + TN}{FP + FN + TP + TN}$

Матр. ошибок:

| | pred + | pred − |
|---|---|---|
| + | TP | FN |
| − | FP | TN |

При большом $w_0$ число $x_i$ т.т. $a(x_i) = -1$

Добавл. $d_y$ -шараги, тогда loss $= d_y [a(x_i) \ne y_i]$

ROC-AUC: по $x$: $TPR = \dfrac{\sum [y_i = -1][a(x_i) = +1]}{\sum [y_i = -1]}$, по $y$: $TPR =$

$= \dfrac{\sum [y_i = +1][a(x_i) = +1]}{\sum [y_i = +1]}$. Оценки общую способности $X^\ell = \{(x, y, \xi)$

$y_i = y^*(x_i)$   $A_\mu = \{a, X \to Y\}$, $\mu_y = (x \times y)^\ell \xrightarrow[\text{обуч}]{\text{метод}} A_\mu$

$h(a, x)$ — Loss $Q(a, x^\ell) = \frac{1}{\ell}\sum h(a, x_i) = \varphi$-ционал качества

Самый common: hold-out $X^\ell = X_n^\ell \sqcup X_n^h$ $n = 1 \ldots N$

$\text{cross-val}(\mu, X^n) = \dfrac{1}{|N|}\sum_{n \in N} Q_\mu(X_n^\ell, X_n^h)$