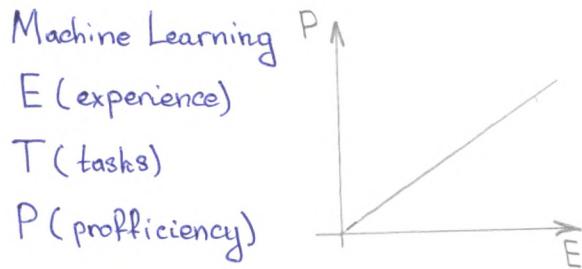


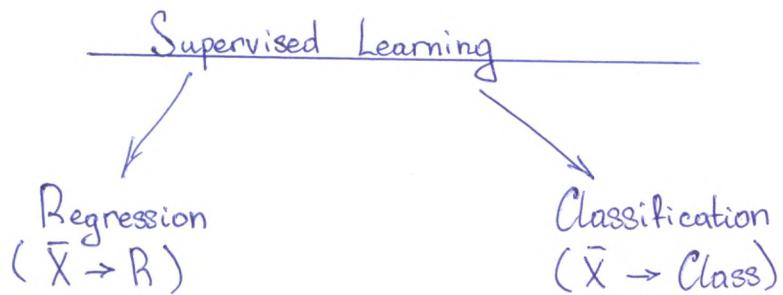
Machine Learning (ML)

- Dynamische Reaktionen (Cognitive Science, Computational Neurobiology)
- Dynamische Psychoanalyse (nur die Psychoanalyse)
- Reichenbachsche Reaktionen (Turing Test "Imitation game" (1950)) 
- Reichenbachsche Psychoanalyse - omnibildende



Statistical approaches : estimation - learning

- Supervised Learning : Regression ($y(\bar{x}) \rightarrow R^c$), Classification ($y(\bar{x}) \rightarrow \vec{p}_{\text{class}}$)
- Unsupervised Learning : Clusterization
- Reinforcement Learning : Agent \rightleftarrows Environment



Data

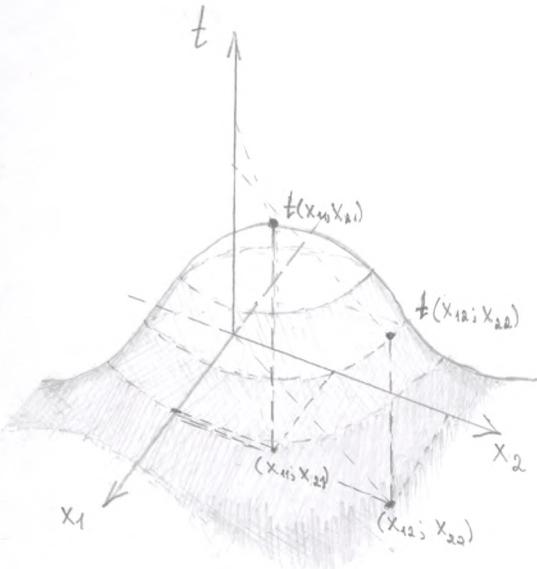
Object features - $\bar{x} \in \mathbb{R}^D$

Training sample

$$\begin{matrix} \bar{x}_1 \\ \bar{x}_2 \\ \bar{x}_3 \\ \vdots \\ \bar{x}_n \end{matrix}$$

Target values

$$\begin{matrix} \bar{t}_1 \\ \bar{t}_2 \\ \bar{t}_3 \\ \vdots \\ \bar{t}_N \end{matrix}$$



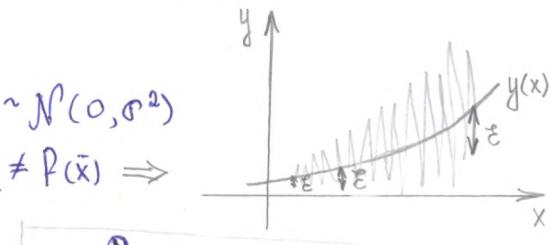
Notation
y - target
\hat{y} - predicted

Linear Regression

$$t(\bar{x}) = y(\bar{x}, \vec{w}) + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2)$$

$\sigma \neq P(\bar{x}) \Rightarrow$
const.

$$y(\bar{x}, \vec{w}) = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_D x_D = w_0 + \sum_{j=1}^D w_j x_j = \vec{w}^T \bar{x}, \quad x_0 = 1$$



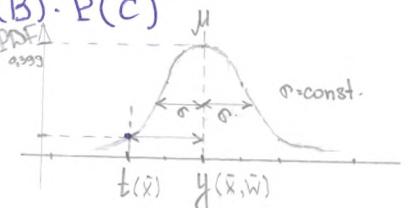
Likelihood

Общая условная плотность вероятности P при $\bar{X}, \vec{w}, \sigma^2$

Events are Independent : $P(A \cap B \cap C) = P(A) \cdot P(B) \cdot P(C)$

$$P(\vec{t} | \bar{X}, \vec{w}, \sigma^2) = \prod_{n=1}^N N(t_n | y(\bar{x}_n, \vec{w}), \sigma^2) \xrightarrow{x} \max$$

$$P(\vec{t} | \bar{X}, \vec{w}, \sigma^2) = \prod_{n=1}^N \frac{1}{\sigma \sqrt{2\pi}} \cdot e^{-\frac{1}{2} \left(\frac{t_n - y(\bar{x}_n, \vec{w})}{\sigma} \right)^2} \xrightarrow{x} \max$$

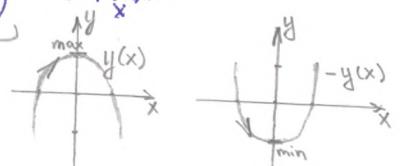


$$\ln P(\vec{t} | \bar{X}, \vec{w}, \sigma^2) = \ln \prod_{n=1}^N \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{t_n - y(\bar{x}_n, \vec{w})}{\sigma} \right)^2} \xrightarrow{x} \max$$

$$= \sum_{n=1}^N \ln \frac{1}{\sigma \sqrt{2\pi}} + \ln \sigma$$

$$\begin{aligned} & P(x) \rightarrow \max \\ & \arg \max_x P(x) = \arg \max_x \log(P(x)) \end{aligned}$$

$$\begin{aligned} \ln P(\vec{t} | \bar{X}, \vec{w}, \sigma^2) &= \sum_{n=1}^N \ln \frac{1}{\sigma \sqrt{2\pi}} + \sum_{n=1}^N \ln e^{-\frac{1}{2} \left(\frac{t_n - y(\bar{x}_n, \vec{w})}{\sigma} \right)^2} = \\ &= \sum_{n=1}^N \ln \frac{1}{\sigma \sqrt{2\pi}} + \sum_{n=1}^N -\frac{1}{2} \left(\frac{t_n - y(\bar{x}_n, \vec{w})}{\sigma} \right)^2 = \\ &= \sum_{n=1}^N \ln \frac{1}{\sigma \sqrt{2\pi}} + \frac{1}{\sigma^2} \underbrace{\left(-\sum_{n=1}^N \frac{1}{2} (t_n - y(\bar{x}_n, \vec{w}))^2 \right)}_{\text{Loss Function}} \xrightarrow{x} \max \end{aligned}$$



Loss Function

$$E(\vec{w}) = \sum_{n=1}^N \frac{1}{2} (t_n - y(\bar{x}_n, \vec{w}))^2 \xrightarrow{x} \min - \text{Mean Squared Error (MSE)}$$

Optimization problem

Adding Non-linearity to Linear Regression

$$t_i(\bar{x}_i) = y(\bar{x}_i, \bar{w}) + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2)$$

$$y(\bar{x}_i, \bar{w}) = \bar{w}^\top \bar{x} = \sum_j x_j w_j = w_0 + x_1 w_1 + x_2 w_2 + \dots + x_D w_D$$

$\bar{x}_i \in \mathbb{R}^D$ - sample

$$\bar{y}(X, \bar{w}) = X^{N \times D} \cdot \bar{w}^{D \times 1} = \begin{bmatrix} \vdots \end{bmatrix}^{N \times 1}$$

$X = \begin{bmatrix} x_{11} & x_{12} & x_{13} & \dots & x_{1D} \\ x_{21} & x_{22} & x_{23} & \dots & x_{2D} \\ \vdots \\ x_{N1} & x_{N2} & x_{N3} & \dots & x_{ND} \end{bmatrix}^{N \times D}$

rows - samples
cols - features

Basis Function

$$\vec{\phi}(\bar{x}_i) = [\phi_0(\bar{x}_i) \ \phi_1(\bar{x}_i) \ \phi_2(\bar{x}_i) \ \dots \ \phi_{M-1}(\bar{x}_i)]^\top$$

$$\vec{\phi} \in \mathbb{R}^M$$

$$y(\bar{x}_i, \bar{w}) = \bar{w}^\top \vec{\phi}(\bar{x}_i) = w_0 + \sum_{j=1}^{M-1} w_j \cdot \phi_j(\bar{x}_i), \quad \phi_0(\bar{x}) = 1$$

Polynomial Basis Functions:

$$\underbrace{[x_1 \ x_2 \ \dots \ x_D]}_{\vec{\phi}_1}, \underbrace{[x_1^2 \ x_2^2 \ \dots \ x_D^2]}_{\vec{\phi}_2}, \underbrace{[x_1^3 \ x_2^3 \ \dots \ x_D^3]}_{\vec{\phi}_3}$$

Design Matrix

$$\Phi(X) = [\phi^\top(\bar{x}_1) \ \phi^\top(\bar{x}_2) \ \dots \ \phi^\top(\bar{x}_N)]^\top$$

$$\Phi(X) = \begin{bmatrix} \phi_0(\bar{x}_1) & \phi_1(\bar{x}_1) & \dots & \phi_{M-1}(\bar{x}_1) \\ \phi_0(\bar{x}_2) & \phi_1(\bar{x}_2) & \dots & \phi_{M-1}(\bar{x}_2) \\ \vdots \\ \phi_0(\bar{x}_N) & \phi_1(\bar{x}_N) & \dots & \phi_{M-1}(\bar{x}_N) \end{bmatrix}^{N \times M}$$

$$\bar{y}(X, \bar{w}) = \Phi(X)^{N \times M} \cdot \bar{w}^{M \times 1} = \begin{bmatrix} \vdots \end{bmatrix}^{N \times 1}$$

\bar{y} нелинейна относительно X .

Linear Regression Optimization.

Normal Equation. Analytical Solution.

$$t(\bar{x}) = y(\bar{x}, \bar{w}) + \varepsilon; \quad \varepsilon \sim N(0, \sigma^2)$$

$\Phi \neq P(\bar{x})$

$$y(\bar{x}, \bar{w}) = \bar{w}^\top \Phi(\bar{x}) = \sum_{i=0}^{M-1} w_i \phi_i(\bar{x})$$

$$P(\vec{t} | X, \bar{w}, \sigma^2) = \prod_{n=1}^N N(t_n | y(\bar{x}_n, \bar{w}), \sigma^2) \rightarrow \max_w$$

$$\bar{x} \in \mathbb{R}^D, \quad \Phi(\bar{x}) \in \mathbb{R}^{M \times D}$$

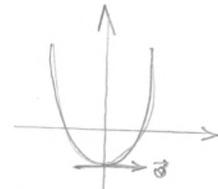
$$\bar{w} \in \mathbb{R}^{N \times D}, \quad \Phi(X) \in \mathbb{R}^{N \times M}$$

Loss function

$$E(\bar{w}) = \sum_{n=1}^N \frac{1}{2} (t_n - y(\bar{x}_n, \bar{w}))^2 \rightarrow \min_w$$

In local minima gradient is equal to zero vector:

$$\nabla E(\bar{w}) = \vec{0}$$



Gradient is vector of partial derivatives: $\nabla y(\bar{x}) = \frac{dy}{dx_1} \vec{i} + \frac{dy}{dx_2} \vec{j} + \frac{dy}{dx_3} \vec{k}$

$$\frac{dE(\bar{w})}{d w_i} = \left(\sum_{n=1}^N \frac{1}{2} (t_n - y(\bar{x}_n, \bar{w}))^2 \right)' = \sum_{n=1}^N 2 \cdot \frac{1}{2} (t_n - y(\bar{x}_n, \bar{w})) \cdot (t_n - y(\bar{x}_n, \bar{w}))'$$

$$\frac{dE(\bar{w})}{d w_i} = -\sum_{n=1}^N (t_n - y(\bar{x}_n, \bar{w})) \cdot \phi_i(\bar{x}_n) \text{ - scalar}$$

$$\nabla E(\bar{w}) = -\sum_{n=1}^N (t_n - y(\bar{x}_n, \bar{w})) \cdot \vec{\Phi}(\bar{x}_n) = \left[\frac{dE(\bar{w})}{d w_0} \quad \frac{dE(\bar{w})}{d w_1} \quad \frac{dE(\bar{w})}{d w_2} \dots \frac{dE(\bar{w})}{d w_{M-1}} \right]^\top$$

vector of partial derivatives

$$\vec{0}^\top = -\sum_{n=1}^N (t_n - y(\bar{x}_n, \bar{w})) \cdot \vec{\Phi}^\top(\bar{x}_n) =$$

$$= -\sum_{n=1}^N t_n \vec{\Phi}^\top(\bar{x}_n) + \sum_{n=1}^N y(\bar{x}_n, \bar{w}) \vec{\Phi}^\top(\bar{x}_n) = -\sum_{n=1}^N t_n \vec{\Phi}^\top(\bar{x}_n) + \sum_{n=1}^N (y) \bar{w}^\top \vec{\Phi}(\bar{x}_n) \vec{\Phi}^\top(\bar{x}_n) =$$

$$= -\sum_{n=1}^N t_n \vec{\Phi}^\top(\bar{x}_n) + (\sum_{n=1}^N \bar{w}^\top \vec{\Phi}(\bar{x}_n)) \vec{\Phi}^\top(\bar{x}_n) =$$

$$= -\vec{t}^\top \vec{\Phi} + \bar{w}^\top \vec{\Phi}^\top \vec{\Phi} = \vec{0}^\top$$

$$\bar{w}^\top \vec{\Phi}^\top \vec{\Phi} = \vec{t}^\top \vec{\Phi}$$

$$\bar{w}^\top = \vec{t}^\top \vec{\Phi} (\vec{\Phi}^\top \vec{\Phi})^{-1}$$

$$\bar{w} = (\vec{t}^\top \vec{\Phi} (\vec{\Phi}^\top \vec{\Phi})^{-1})^\top = ((\vec{\Phi}^\top \vec{\Phi})^{-1})^\top \cdot (\vec{t}^\top \vec{\Phi})^\top = ((\vec{\Phi}^\top \vec{\Phi})^\top)^{-1} \cdot (\vec{\Phi}^\top \vec{t})$$

$$(AB)^\top = B^\top A^\top$$

$$(A^{-1})^\top = (A^\top)^{-1}$$

$$\bar{w} = (\vec{\Phi}^\top \vec{\Phi})^{-1} \vec{\Phi}^\top \vec{t}$$

Moore-Penrose Pseudo inverse

$$\vec{\Phi}^\top \vec{\Phi}^{M \times N} \cdot \vec{\Phi}^{N \times M} = \vec{A}^{M \times M} \text{ - depends on number of features of } \vec{x}$$

$$\vec{A}^{-1} = ?$$

\Rightarrow Gradient Descent

$$\begin{aligned} (f(x) + g(x))' &= f'(x) + g'(x) \\ (f(g(x)))' &= f'(g(x)) \cdot g'(x) \\ (x^n)' &= nx^{n-1} \\ (C)' &= 0 \end{aligned}$$

(Linear Regression Optimization)

$$\vec{t} = \begin{bmatrix} t_1 \\ t_2 \\ \vdots \\ t_n \end{bmatrix}$$

$$\bar{\Phi}(\bar{x}_i) = \begin{bmatrix} \varphi_0(\bar{x}_i) \\ \varphi_1(\bar{x}_i) \\ \varphi_2(\bar{x}_i) \\ \vdots \\ \varphi_{M-1}(\bar{x}_i) \end{bmatrix}$$

$$\Phi(X) = \begin{bmatrix} \varphi_0(\bar{x}_1) & \varphi_1(\bar{x}_1) & \dots & \varphi_{M-1}(\bar{x}_1) \\ \varphi_0(\bar{x}_2) & \varphi_1(\bar{x}_2) & \dots & \varphi_{M-1}(\bar{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \varphi_0(\bar{x}_N) & \varphi_1(\bar{x}_N) & \dots & \varphi_{M-1}(\bar{x}_N) \end{bmatrix}^{N \times M}$$

$$\sum_{n=1}^N t_n \bar{\Phi}^T(\bar{x}_n) = \begin{bmatrix} t_n \varphi_0(x_n) & t_n \varphi_1(x_n) & \dots \\ \vdots & \vdots & \vdots \\ t_n \varphi_0(x_{n+1}) & t_n \varphi_1(x_{n+1}) & \dots \end{bmatrix} = \begin{bmatrix} t_n \varphi_0(x_n) + t_{n+1} \varphi_0(x_{n+1}) & t_n \varphi_1(x_n) + t_{n+1} \varphi_1(x_{n+1}) \\ \vdots & \vdots \\ \vdots & \vdots \end{bmatrix}$$

1) $\vec{t}^T \cdot \Phi^*(X) = [\underline{t_1 \ t_2 \ \dots \ t_n}] \begin{bmatrix} \varphi_0(\bar{x}_1) & \varphi_1(\bar{x}_1) & \dots \\ \varphi_0(\bar{x}_2) & \varphi_1(\bar{x}_2) & \dots \\ \varphi_0(\bar{x}_N) & \varphi_1(\bar{x}_N) & \dots \end{bmatrix}^{N \times M} = \left[\sum_{n=1}^N t_n \varphi_0(\bar{x}_n) \quad \sum_{n=1}^N t_n \varphi_1(\bar{x}_n) \right] \top$

2) $\sum_{n=1}^N \bar{\Phi}(\bar{x}_n) \bar{\Phi}^T(\bar{x}_n) = \sum_{n=1}^N \begin{bmatrix} \varphi_0(\bar{x}_n) \\ \varphi_1(\bar{x}_n) \\ \vdots \\ \varphi_{M-1}(\bar{x}_n) \end{bmatrix} [\varphi_0(\bar{x}_n) \ \varphi_1(\bar{x}_n) \ \dots \ \varphi_{M-1}(\bar{x}_n)] = \sum_{n=1}^N \begin{bmatrix} \varphi_0 \varphi_0 & \varphi_0 \varphi_1 & \dots & \varphi_0 \varphi_{M-1} \\ \varphi_1 \varphi_0 & \varphi_1 \varphi_1 & \dots & \varphi_1 \varphi_{M-1} \\ \vdots & \vdots & \ddots & \vdots \\ \varphi_{M-1} \varphi_0 & \varphi_{M-1} \varphi_1 & \dots & \varphi_{M-1} \varphi_{M-1} \end{bmatrix} =$

$$= \Phi^T \Phi$$

 $M \times N \cdot N \times M = M \times M$

$$\Phi^T \Phi = \begin{bmatrix} \varphi_0(\bar{x}_1) & \varphi_0(\bar{x}_2) & \dots & \varphi_0(\bar{x}_N) \\ \varphi_1(\bar{x}_1) & \varphi_1(\bar{x}_2) & \dots & \varphi_1(\bar{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ \varphi_{M-1}(\bar{x}_1) & \varphi_{M-1}(\bar{x}_2) & \dots & \varphi_{M-1}(\bar{x}_N) \end{bmatrix} \cdot \begin{bmatrix} \varphi_0(\bar{x}_1) & \varphi_1(\bar{x}_1) & \dots & \varphi_{M-1}(\bar{x}_1) \\ \varphi_0(\bar{x}_2) & \varphi_1(\bar{x}_2) & \dots & \varphi_{M-1}(\bar{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \varphi_0(\bar{x}_N) & \varphi_1(\bar{x}_N) & \dots & \varphi_{M-1}(\bar{x}_N) \end{bmatrix}$$

Weight Decay Regularization

$$t(\bar{x}) = y(\bar{x}, \bar{w}) + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2), \quad \sigma^2 \neq f(\bar{x})$$

$$P(\vec{t} | X, \bar{w}, \sigma^2) = \prod_{n=1}^N N(t_n | y(\bar{x}_n, \bar{w}), \sigma^2) \xrightarrow[w]{\text{max}} \text{product - общая условная вероятность}$$

вероятностей \vec{t} при X, \bar{w} и σ^2 .

$$y(\bar{x}_n, \bar{w}) = \bar{w}^\top \bar{\Phi}(\bar{x}_n)$$

$$E(\bar{w}) = \frac{1}{2} \sum_{n=1}^N (t_n - \bar{w}^\top \bar{\Phi}(\bar{x}_n))^2 + \frac{\lambda}{2} \sum_{j=1}^{N-1} |w_j|^q \xrightarrow[w]{\min}$$

- No-bias, zero mean assumption about parameters by regularization.

- standard values for q are 1 and 2.

Регуляризация накладывает давление на минимальное значение параметров модели.

$$E(\bar{w}) = \frac{1}{2} \sum_{n=1}^N (t_n - \bar{w}^\top \bar{\Phi}(\bar{x}))^2 + \frac{\lambda}{2} \bar{w}^\top \bar{w} \rightarrow \min$$

$$\frac{d\left(\frac{\lambda}{2} \sum_{j=1}^{N-1} w_j^2\right)}{d(w_j)} = 2 \cdot \frac{\lambda}{2} \cdot (0 + \dots + w_j) = \lambda w_j \Rightarrow \nabla = \begin{bmatrix} \lambda w_1 \\ \lambda w_2 \\ \vdots \\ \lambda w_{N-1} \end{bmatrix} = \lambda \bar{w}$$

$$\nabla E(\bar{w}) = - \sum_{n=1}^N (t_n - \bar{w}^\top \bar{\Phi}(\bar{x}_n)) \bar{\Phi}(\bar{x}_n) + \lambda \bar{w} \rightarrow \min$$

$$\bar{\Phi}^\top = \nabla E(\bar{w})^\top = -\vec{t}^\top \Phi + \bar{w}^\top \Phi^\top \Phi + \lambda \bar{w}^\top \stackrel{\Phi^\top \Phi \in \mathbb{R}^{M \times M}}{\stackrel{\bar{w}^\top \in \mathbb{R}^M}{\rightarrow}} = -\vec{t}^\top \Phi + \bar{w}^\top (\Phi^\top \Phi + \lambda I)$$

$$\bar{w}^\top = \vec{t}^\top \Phi (\Phi^\top \Phi + \lambda I)^{-1}$$

$$\bar{w} = (\vec{t}^\top \Phi (\Phi^\top \Phi + \lambda I)^{-1})^\top = ((\Phi^\top \Phi + \lambda I)^{-1})^\top \cdot (\vec{t}^\top \Phi)^\top = \underbrace{(\Phi^\top \Phi + \lambda I)^{-1} \cdot (\Phi^\top \vec{t})}_{\text{Normal Equation}}$$

Normal Equation

$$\nabla E(\bar{w}) = - \sum_{n=1}^N (t_n - \bar{w}^\top \bar{\phi}(\bar{x}_n)) \bar{\phi}(\bar{x}_n) + \lambda \bar{w}$$

$$\bar{\Phi} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$n \times 1 \quad n \times 1$

$$\bar{\Phi} = \nabla E(\bar{w}) = - \sum_{n=1}^N (t_n \vec{\phi}(x_n) - (\bar{w}^\top \bar{\phi}(\bar{x}_n)) \bar{\phi}(\bar{x}_n)) + \lambda \bar{w} = -\bar{\phi}^\top \vec{t} +$$

Therefore we should transpose

~~$$\begin{bmatrix} \phi_0(x_1) & \phi_0(x_2) & \dots \\ \phi_1(x_1) & \phi_1(x_2) & \dots \\ \phi_2(x_1) & \phi_2(x_2) & \dots \\ \phi_3(x_1) & \phi_3(x_2) & \dots \end{bmatrix} \quad \begin{bmatrix} \phi_0(x_1) & \phi_1(x_1) & \phi_2(x_1) \\ \phi_0(x_2) & \phi_1(x_2) & \phi_2(x_2) \\ \vdots & \vdots & \vdots \end{bmatrix}$$~~

$$\sum_{n=1}^N \bar{\phi}(\bar{x}_n) \bar{\phi}^\top(\bar{x}_n) = \left(\sum_{n=1}^N \begin{bmatrix} \phi_0(x_n) \\ \phi_1(x_n) \\ \phi_2(x_n) \end{bmatrix} \right)^{\text{3x1}} \cdot \left[\begin{bmatrix} \phi_0(x_1) & \phi_1(x_1) & \phi_2(x_1) \\ \phi_0(x_2) & \phi_1(x_2) & \phi_2(x_2) \\ \vdots & \vdots & \vdots \end{bmatrix} \right]^{\text{1x3}} = \begin{bmatrix} \cdot & \cdot & \cdot \end{bmatrix}^{\text{3x3}} =$$

$$= \sum_{n=1}^N \begin{bmatrix} \phi_0 \phi_0 & \phi_0 \phi_1 & \phi_0 \phi_2 \\ \phi_1 \phi_0 & \phi_1 \phi_1 & \phi_1 \phi_2 \\ \phi_2 \phi_0 & \phi_2 \phi_1 & \phi_2 \phi_2 \end{bmatrix}$$

$$\Phi(\bar{X}) = \begin{bmatrix} \phi_0(\bar{x}_1) & \phi_1(\bar{x}_1) & \phi_2(\bar{x}_1) \\ \phi_0(\bar{x}_2) & \phi_1(\bar{x}_2) & \phi_2(\bar{x}_2) \\ \phi_0(\bar{x}_3) & \phi_1(\bar{x}_3) & \phi_2(\bar{x}_3) \end{bmatrix}^{N \times M}$$

$$\Phi^\top(\bar{X}) = \begin{bmatrix} \phi_0(\bar{x}_1) & \phi_0(\bar{x}_2) & \phi_0(\bar{x}_3) \\ \phi_1(\bar{x}_1) & \phi_1(\bar{x}_2) & \phi_1(\bar{x}_3) \\ \phi_2(\bar{x}_1) & \phi_2(\bar{x}_2) & \phi_2(\bar{x}_3) \\ \vdots & \vdots & \vdots \end{bmatrix}_{M \times N}$$

$$\Phi^\top \cdot \Phi = \begin{bmatrix} \sum_{n=1}^N \bar{\phi}_0(\bar{x}_n) \phi_0(x_n) & \sum_{n=1}^N \bar{\phi}_0(\bar{x}_n) \bar{\phi}_1(\bar{x}_n) & \dots \\ \sum_{n=1}^N \bar{\phi}_1(\bar{x}_n) \phi_0(x_n) & \sum_{n=1}^N \bar{\phi}_1(\bar{x}_n) \bar{\phi}_1(\bar{x}_n) & \dots \\ \sum_{n=1}^N \bar{\phi}_2(\bar{x}_n) \phi_0(x_n) & \sum_{n=1}^N \bar{\phi}_2(\bar{x}_n) \bar{\phi}_1(\bar{x}_n) & \dots \end{bmatrix}_{M \times M}$$

$$w^\top \Phi^\top \Phi w = \begin{bmatrix} w_1 \sum + w_2 \sum + w_3 \sum & -w_1 \sum + w_2 \sum + w_3 \sum & \dots \end{bmatrix}$$

$$w^\top \Phi^\top \Phi w + \lambda w^\top = \begin{bmatrix} w_1 \sum + w_2 \sum + w_3 \sum + w_4 & \dots \end{bmatrix}$$

Методическое решения задачи регрессии

process → Data → Analyze → Estimate (Learning)

estimate unknown params of process generating data.

- Могут не заложить истину.

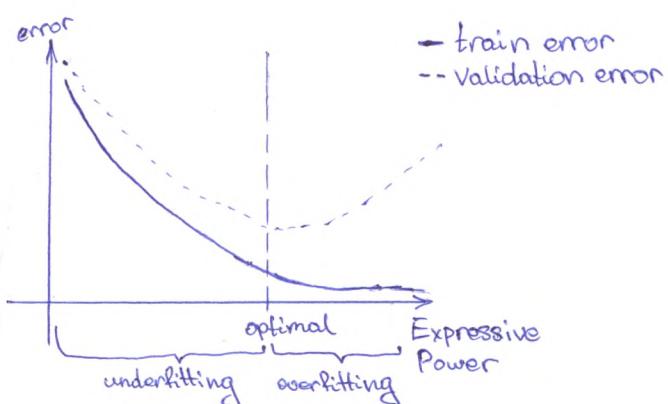
Общее изложение гаусса:

Наше гауссово распределение;
За какой период времени;
Сколько гауссов;
Среднее (μ);
Разброс (σ^2);

Обычные:

Train	Validation	Test
60%	20%	20%

- 1) Try Linear Regression $w^T \bar{x} = \hat{y}$
- 2) Try other basis functions $w^T \phi(\bar{x})$
Test on Validation dataset
- 3) Underfitting vs. Overfitting
- 4) Regularization
- 5) Переупаковка гиперпараметров: φ_i, λ
- 6) Модель с наименьшей ошибкой на вынужденной выборке отбираем.
- 7) Test on test dataset. - реальность.



Пример ненормированных train, validation, test

$10^{1000\dots}$ моделей

модели генерируют 100% accuracy, среди них можно найти модель низкочастотную, и которая подобрана 100% accuracy на validation, но на тестовой выборке имеет ~50% acc.

Classification

$$\text{Class} = \{0, 1, 2, \dots, K\}$$

u - object (cat, patient ...)

$$\bar{x} = \mu(u) - \text{features of object (image; temperature, pain...)} \\ \bar{x} \in \mathbb{R}^D, \quad \bar{x} = [x_1, x_2, \dots, x_D]^T$$

- Binary Classification ($K=2$)

$$y(\bar{x}, \bar{w}) \rightarrow \text{Class or P(Class)} \quad \rightarrow t$$

- Multiclass Classification ($K > 2$)

$$\hat{y}(\bar{x}, \bar{w}) = [p_1, p_2, \dots, p_K]^T \quad \rightarrow \vec{f} - \text{one-hot encoding} \\ \sum_{k=1}^K p_k = 1, \quad p_k \geq 0, \quad k = 1, 2, \dots, K$$

p_k - вероятность принадлежности объекта
к классу k .

Confusion Matrix

$$C = \{c_{ij}\}_{K \times K}$$

$$c_{ij} = |\{u \in C_i \mid y(\bar{x}(u), \bar{w}) = j\}|$$

		y		
		0	1	
t	0	TN	FP	
	1	FN	TP	

В многоклассовом случае можно считать количество
правильных классификаций и ошибок.

Binary Classification Error

$$y(\bar{x}, \bar{w}) = 20, 13$$

- ✓ True Positive = $\{u_i \in C_1 | y(\bar{x}, \bar{w}) = 1\}$
 - ✓ True Negative = $\{u_i \in C_0 | y(\bar{x}, \bar{w}) = 0\}$
 - ✗ False Positive = $\{u_i \in C_0 | y(\bar{x}, \bar{w}) = 1\}$
 - ✗ False Negative = $\{u_i \in C_1 | y(\bar{x}, \bar{w}) = 0\}$

$$\text{Accuracy} = \frac{TP + TN}{N}$$

$$\begin{aligned}TP + TN + FP + FN &= \# \\TP + FN &= \#P \\TN + FP &= \#N\end{aligned}$$

Problem of accuracy is in imbalanced dataset.

Здоровых людей больше, чем больных:
(99%) (1%)

If $y(\bar{x}, \bar{w}) = 0$, then accuracy is 99%, but it is useless.

Они были первыми и вторыми погибшими

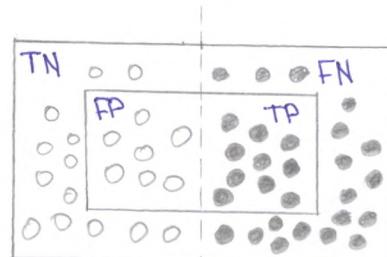
$$d_N = \frac{FP}{FP+TN} \quad (\text{FPR})$$

false positive rate

$$\beta_p = \frac{FN}{FN + TP} \quad (\text{FNR})$$

false negative rate

Recall + $\beta_P = 1$



- Negative
- Positive

Best when:
TN
TP

Precision & Recall

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (\text{TPR})$$

true positive rate

$$F_1\text{-Score} = \frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$$

TN manom stems inconspicuous

$$\text{then } \text{Precision} = \frac{\geq 0}{\geq 0 + \geq 0} \geq 0$$

$$\text{Recall} = \frac{70}{70+0} = 1$$

$$2 \cdot 1 \cdot 20 = \frac{0,00...1}{1,00...1} = 20$$

$$F_1\text{-score} = \frac{2 \cdot 1 \cdot 0,0}{1 + 0,0} = \frac{0,0}{1,00} = 0$$

TN manum Stems nonhegally

F₁ score vs. average

Sensitivity & Specificity

$$\text{Sensitivity (TPR, Recall)} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

$$\text{Specificity} = \frac{\text{TN}}{\text{TN} + \text{FP}} \quad (\text{TNR})$$

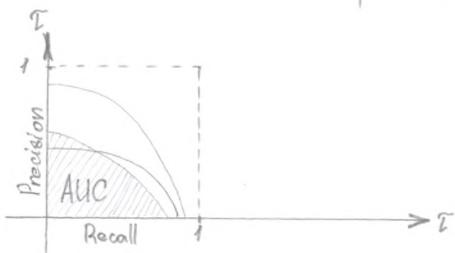
true negative rate

Precision - Recall Curve

Precision & Recall - zuarenuia gaa onpegeññoro noporoboro zuarenuia (2, may).

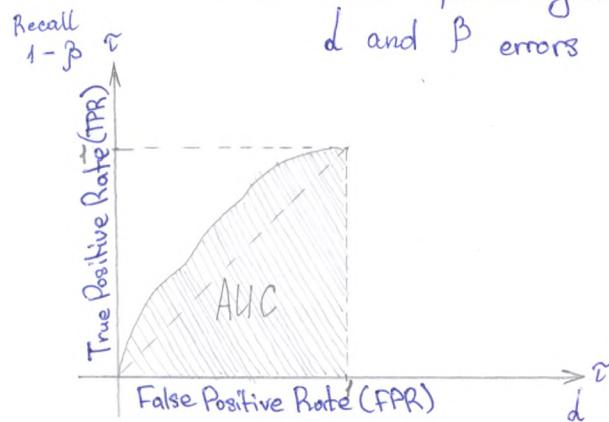
t	P_0	P_1	$\tilde{\tau} = 0,5$	$\tilde{\tau} = 0,6$	$\tilde{\tau} = 0,75$
0	0,3	0,4	1	1	0
0	0,4	0,6	1	1	0
1	0,45	0,55	1	0	0
1	0,2	0,8	1	1	1

$$\uparrow \text{Acc.} = 0,5 \quad \text{Acc.} = 0,25 \quad \text{Acc.} = 0,75$$



AUC - Area Under the Curve

Receiver Operating Characteristic (ROC)



Accuracy:

$$\text{Precision} = \frac{TP}{TP + FP}$$

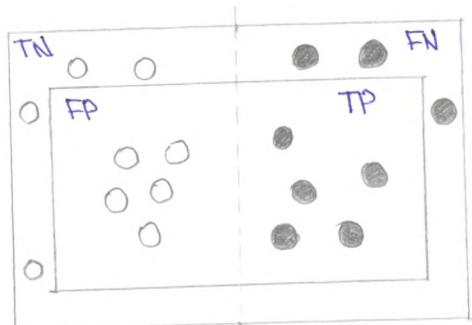
$$\text{Recall} = \frac{TP}{TP + FN} \quad (\text{TPR})$$

Error:

$$\alpha_{\text{FP}} = \frac{FP}{FP + TN} \quad (\text{FPR})$$

$$\beta_{\text{P}} = \frac{FN}{FN + TP} \quad (\text{FNR})$$

● - positive
○ - negative



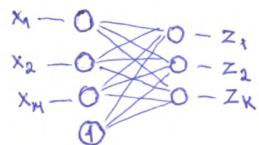
Logistic Regression (Classification)

Data

\vec{x}_i	Class Name	t_i	One-hot encoding vector
\vec{x}_1	Cat	1	[1 0 0]
\vec{x}_2	Dog	2	[0 1 0]
\vec{x}_3	Horse	3	[0 0 1]
\vec{x}_n	"String"	$t_n \in [1; K]$	$\vec{t}_n = [\vec{t}_n^{(1)} \vec{t}_n^{(2)} \dots \vec{t}_n^{(K)}]^T$ $t_n^{(t_n)} = 1, \text{ otherwise } t_n^{(k)} = 0$

Множество наблюдений и модели предсказаний

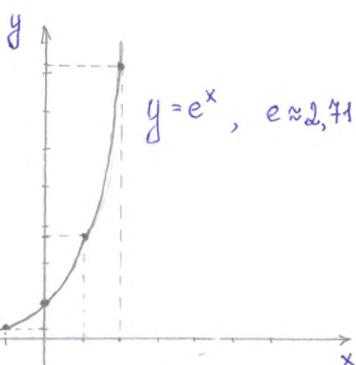
$$\left\{ \begin{array}{l} z_1 = w_{10} + w_{11}x_1 + w_{12}x_2 + \dots + w_{1M}x_M = \vec{w}_1^T \vec{x} + b_1 \\ z_2 = w_{20} + w_{21}x_1 + w_{22}x_2 + \dots + w_{2M}x_M = \vec{w}_2^T \vec{x} + b_2 \\ \dots \\ z_K = w_{K0} + w_{K1}x_1 + w_{K2}x_2 + \dots + w_{KM}x_M = \vec{w}_K^T \vec{x} + b_K \end{array} \right.$$



$$\begin{aligned} \vec{x} &\in \mathbb{R}^M \\ \vec{y} &\in \mathbb{R}^K \end{aligned} \quad \mathbb{W} = \begin{bmatrix} \vec{w}_1^T \\ \vec{w}_2^T \\ \vdots \\ \vec{w}_K^T \end{bmatrix}, \quad \mathbb{W} = (\mathbb{W} \quad \vec{b})$$

$$\mathbb{W}^{K \times (M+1)} \cdot \vec{x}^{(N+1) \times 1} \stackrel{\# \text{ of samples}}{\uparrow} = \vec{z}^{K \times 1}$$

$$\begin{aligned} \vec{y}_i(\vec{x}_i, \mathbb{W}) &= \text{Softmax}(\vec{w}^T \vec{x} + \vec{b}) = \\ &= \text{Softmax}(\mathbb{W}^T \vec{x}) \end{aligned}$$



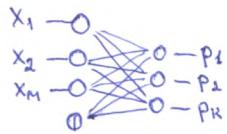
$$\text{Softmax}(\vec{z}) = \left(\frac{e^{z_i}}{\sum_i e^{z_i}} \right)$$

$$1) \sum_i \frac{e^{z_i}}{\sum_i e^{z_i}} = 1$$

2) Softmax is differentiable function

Logistic Regression. Loss Function.

$$\hat{y}_i(\vec{x}_i, \mathbb{W}) = \text{Softmax}(\mathbb{W}\vec{x}_i + \vec{b}), \quad \hat{y} \in \mathbb{R}^{K \times 1}$$



p_i - Softmax probability,
not PDF of Normal Distribution.

Вероятность правильного класса i на 1-ом sample-e и
класса j на 2-ом sample-e.

$$P(y_i, y_j) = y_i(\vec{x}_1, \mathbb{W}) \cdot y_j(\vec{x}_2, \mathbb{W}) - \text{Independent Events:}$$

$$P(A \text{ and } B) = P(A) \cdot P(B)$$

Maximum Likelihood

Одна из основных близких к training sample-ам $X^{M \times N}$,
которые есть в базе данных, называемые правильные
объекты $T^{K \times N}$

$$P(T | X, \mathbb{W}) = \underbrace{y_{t_1}^{(t_1)}(\vec{x}_1, \mathbb{W}) \cdot \dots \cdot y_{t_n}^{(t_n)}(\vec{x}_n, \mathbb{W})}_{y_k \text{ on right class}} \rightarrow \max_{\mathbb{W}}$$

$$P(T | X, \mathbb{W}) = \prod_{n=1}^N \prod_{k=1}^K \underbrace{y_k^{(k)}(\vec{x}_n, \mathbb{W})^{t_n^{(k)}}}_{y_n^{(t_n)}(\vec{x}_n, \mathbb{W}) = \prod_{k=1}^K y_k^{(k)}(\vec{x}_n, \mathbb{W})^{t_n^{(k)}}} \rightarrow \max_{\mathbb{W}}$$

$$\begin{aligned} t_n^{(k)} &\in \{0, 1\} \\ \begin{cases} t_n^{(k)} = 1 & \text{if } k = t_n \\ t_n^{(k)} = 0 & \text{if } k \neq t_n \end{cases} \\ \vec{t}_n &= [0 \ 0 \ 0 \dots 1 \dots 0 \ 0]^\top - \text{one-hot encoding} \\ t_n &- \# \text{of class that's correct.} \end{aligned}$$

$$-\ln(P(T | X, \mathbb{W})) = -\sum_{n=1}^N \sum_{k=1}^K \ln(y_n^{(k)}(\vec{x}_n, \mathbb{W})^{t_n^{(k)}}) \rightarrow \min_{\mathbb{W}}$$

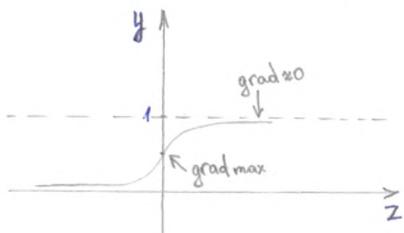
Loss Function (Числова оцінка)

$$E(\mathbb{W}) = -\sum_{n=1}^N \sum_{k=1}^K t_n^{(k)} \ln(y_n^{(k)}(\vec{x}_n, \mathbb{W})) \rightarrow \min_{\mathbb{W}}$$

Weights initialization (Glorot/Xavier)

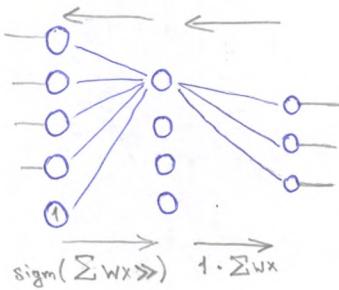
Vanishing / Exploding Gradients Problem ?

$$y(\bar{x}, \bar{w}) = \text{sigm} \left(\sum_i w_i x_i + b \right) , \quad x \sim \mathcal{N}(0, \sigma^2=1)$$



sigm \rightarrow 1) $x \in [0; 1]$
 $w \in [0; \frac{1}{n}]$
 $\sum_i x_i w_i = [0; 1]$

tanh \rightarrow 2) $x \in [-1; 1] \approx \mathcal{N}(0, 1)$
 $w \in [-\frac{1}{n}; \frac{1}{n}]$ or $\mathcal{N}(0, \frac{1}{n})$
 $\sum_i x_i w_i = [-1; 1]$
 we need this



$$w_i \sim \mathcal{N}(\mu=0, \sigma^2 = \left(\frac{2}{n_{in} + n_{out}} \right)^2)$$

↑ ↑
 sigma=1 exploding gradient

Balance between :

- 1) $\sum w_i x_i \gg$
- 2) Exploding gradient

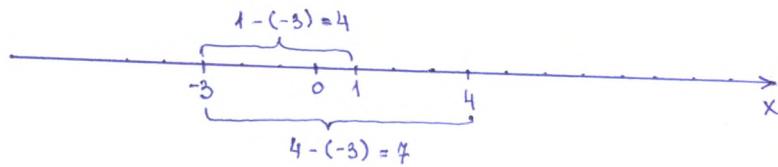
Data Normalization

$$\vec{x} = [x^{(1)} \ x^{(2)} \ \dots \ x^{(M)}]^T$$

$$\tilde{x}^{(j)} = 2 \frac{x^{(j)} - \min_i x_i^{(j)}}{\max_i x_i^{(j)} - \min_i x_i^{(j)}} - 1$$

$$\tilde{x}^{(j)} \in [-1; 1]$$

$x^{(j)}$ нормализуем членко по j гваренум бармопаб \vec{x}_i из $X \in \mathbb{R}^{M \times N}$



$$\tilde{x}^{(j)} = \frac{x^{(j)} - \mu^{(j)}}{\sigma^{(j)}} , \quad \sigma \gg 0$$

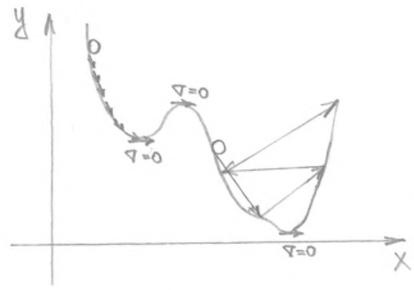
$$\tilde{x}^{(j)} \sim \mathcal{N}(\mu=0, \sigma^2=1)$$

Momentum SGD

1. $\nabla_1 \leftarrow \text{velocity}$

$$2. \boxed{\nabla_2 + \mu \cdot \nabla_1}$$

$$3. \nabla_3 + \mu \nabla_2 + \mu^2 \nabla_1 = \nabla_3 + \mu (\nabla_2 + \mu \nabla_1)$$



$$V_2 = V_1 + \mu V_0 \quad \text{no same on the momentum crazy ostatewumca.}$$

velocity = 0

iter

velocity += momentum

velocity += $\nabla E(w)$

$$w -= lr \cdot \text{velocity}$$

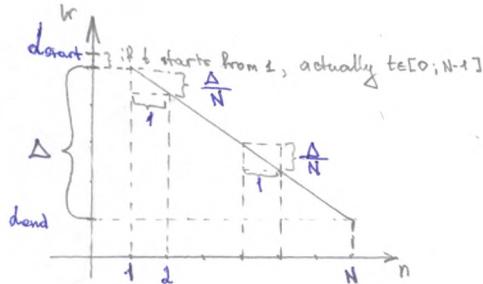
Regularized SGD

$$w -= lr \cdot \nabla E(w) + \lambda w$$

Learning Rate Decay

1. Linear

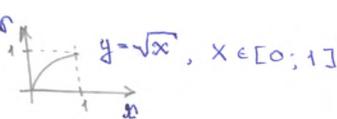
$$d_t = d_{\text{start}} - (d_{\text{start}} - d_{\text{end}}) \cdot \frac{t}{N}$$



2. Exponential

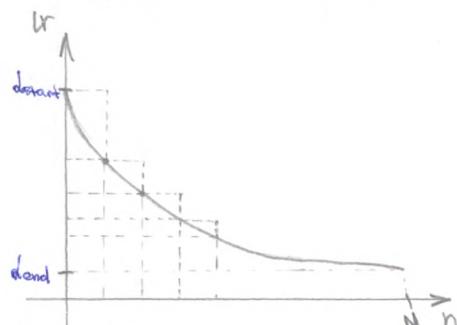
$$d_t = d_{\text{start}} \cdot \sigma^t$$

$$\sigma = \left(\frac{d_{\text{end}}}{d_{\text{start}}} \right)^{\frac{1}{N-1}}$$



$$d_t = d_{\text{start}} \cdot \left(\frac{d_{\text{end}}}{d_{\text{start}}} \right)^{\frac{1}{N-1}}$$

$$\left(\left(d_{\text{start}} \cdot \frac{\sqrt[3]{\sigma}}{\text{const}} \right) \cdot \frac{\sqrt[3]{\sigma}}{\text{const}} \cdot \frac{\sqrt[3]{\sigma}}{\text{const}} \right) = d_{\text{start}} \cdot \frac{d_{\text{end}}}{d_{\text{start}}} = d_{\text{end}}$$



Dropout

- training $0, 1 \rightarrow 0, p$

- inference input $\cdot 0, 1$, we mimic Magnitude $\sim (1-p)$

Logistic Regression. Model parameters estimation.

Gradient Descent

Loss Function

$$E(W) = - \sum_{n=1}^N \sum_{k=1}^K t_n^{(k)} \cdot \ln(y_n^{(k)}(\bar{x}_n, W)) \rightarrow \min_W, \quad W = (W, b)$$

1. Стартовая инициализация весов.

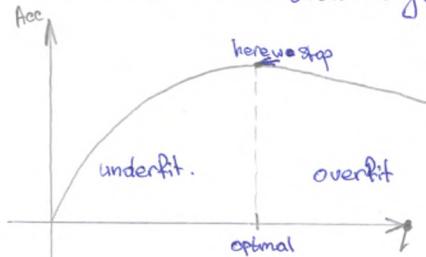
$$w_i \sim \mathcal{N}(\mu=0, \sigma^2 = \frac{1}{M})$$

2. Установка итерации:

$$W_k = W_{k-1} - \gamma \cdot \nabla E(W_{k-1}), \quad \gamma - \text{learning rate, hyperparameter.}$$

3. Метод перебора пробны, когда мотивы на базовом выборке не совпадают с реальными.

Максимум остановки обучения - это же некий hyperparameter.



Logistic Regression. Model parameters estimation.

Loss Function

$$E(\mathbb{W}) = - \sum_{n=1}^N \sum_{k=1}^K t_n^{(k)} \ln y^{(k)}(\bar{x}_n, \mathbb{W}) \rightarrow \min_{\mathbb{W}}, \quad \mathbb{W} = (\mathbb{W}, b)$$

Dua 1-20 saemema fólofren

$$E(\mathbb{W}) = - \sum_{k=1}^K t^{(k)} \ln y^{(k)}(\bar{x}, \mathbb{W}) = - \ln y^{(t)}(\bar{x}, \mathbb{W}) \rightarrow \min_{\mathbb{W}}$$

$$\bar{y}(\bar{x}, \mathbb{W}) = \text{Softmax}(\bar{z}(\bar{x}, \mathbb{W})) = \text{Softmax}(W\bar{x} + b)$$

$$E(\mathbb{W}) = - \ln y^{(t)}(\bar{x}, \mathbb{W}) = - \ln \left(\frac{e^{z_t}}{\sum_{k=1}^K e^{z_k}} \right) = - \ln e^{z_t} + \ln \left(\sum_{k=1}^K e^{z_k} \right) = - z_t + \ln \left(\sum_{k=1}^K e^{z_k} \right) \rightarrow \min_{\mathbb{W}}$$

$$\{ z_i = \sum_j w_j x_j = \bar{w}_i^T \bar{x}_i + b_i \}$$

$$\bar{z} = W\bar{x} + b$$

$$\frac{dE(\mathbb{W})}{d w_{tl}} = - (w_{ti} x_i)^T + \frac{1}{\sum_{k=1}^K e^{z_k}} \cdot e^{z_t} \cdot (w_{ti} x_i)^T = -x_i + y_t^{(t)} \cdot x_i = (y_t^{(t)} - 1) \cdot x_i$$

$$\frac{dE(\mathbb{W})}{d w_{il}} = -0 + y_t^{(i)} \cdot (w_{il} \cdot x_i)^T = y^{(i)} \cdot x_i, \quad i \neq t$$

$$\nabla_{\mathbb{W}} E = \begin{bmatrix} \frac{dE}{dw_{11}} & \frac{dE}{dw_{12}} & \frac{dE}{dw_{13}} & \dots & \frac{dE}{dw_{1M}} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \frac{dE}{dw_{K1}} & \frac{dE}{dw_{K2}} & \frac{dE}{dw_{K3}} & \dots & \frac{dE}{dw_{KM}} \end{bmatrix} = \begin{bmatrix} y_1 x_1 & y_1 x_2 & \dots & y_1 x_N \\ (y_t - 1)x_1 & (y_t - 1)x_2 & \dots & (y_t - 1)x_M \\ \vdots & \vdots & \ddots & \vdots \\ y_K x_1 & y_K x_2 & \dots & y_K x_N \end{bmatrix}$$

$$\nabla_{\mathbb{W}} E = (\bar{y} - t) \cdot \bar{x}^T = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}^{K \times 1} \cdot \begin{bmatrix} \dots & \dots & \dots \end{bmatrix}^{1 \times M} = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}^{K \times M}$$

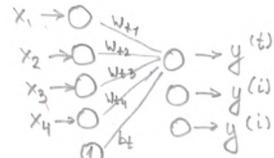
$$\begin{aligned} \frac{dE}{db_t} &= -1 + y_t = y_t - 1 \\ \frac{dE}{db_i} &= -0 + y_i = y_i \end{aligned} \Rightarrow \nabla_b E = (\bar{y} - t)$$

Regularization

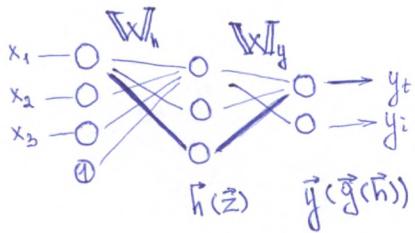
$$\nabla_{\mathbb{W}} E = (\bar{y} - t) \cdot \bar{x}^T + \lambda \mathbb{W}$$

$$\nabla_b E = (\bar{y} - t)$$

- Batch gradient descent - entire dataset
- Mini-batch - sample from dataset
- Stochastic - 1 element from dataset



Multi layer perceptron



$\sum \vec{z} \rightarrow \text{Sigmoid} \rightarrow \vec{h} \rightarrow \text{Sigmoid} \rightarrow \vec{g} \rightarrow \text{Sum} \rightarrow \vec{y}$

$$\vec{y}(g(h(\vec{z}))) = \text{Softmax}(\vec{g}(\vec{h})) = \left(\frac{e^{g_i}}{\sum_{j=1}^K e^{g_j}} \right)$$

$$P(T|X, W_h, W_g) = \prod_{n=1}^N \prod_{k=1}^K y^{(k)} \cdot (\bar{x}_n, W_h, W_g)^{t_n^{(k)}} \rightarrow \max_{W_h, W_g}$$

$$-\ln(P(T|X, W_h, W_g)) = -\sum_{n=1}^N \sum_{k=1}^K \ln y^{(k)} (\bar{x}_n, W_h, W_g)^{t_n^{(k)}} \rightarrow \min_{W_h, W_g}$$

$$-\sum_{k=1}^K \ln y^{(k)} (\bar{x}, W_h, W_g)^{t^{(k)}} = -\ln y^{(t)} (\bar{x}, W_h, W_g)^{t^{(k)}} \rightarrow \min_{W_h, W_g}$$

$$E(W_h) = -\ln y^{(t)} (\bar{x}, W_h, W_g) = -\ln \left(\frac{e^{g_t(\bar{x})}}{\sum_{k=1}^K e^{g_k(\bar{x})}} \right) = -\ln(e^{g_t(\bar{x})}) + \ln \left(\sum_{k=1}^K e^{g_k(\bar{x})} \right) \rightarrow \min_{W_h}$$

$$g_t(\vec{h}(\vec{z})) = \sum_{i=1}^H w_{ti} \cdot h_i(z_i)$$

$$h_i(z_i) = \frac{1}{1+e^{-z_i}}$$

$$z_i(W_h) = \sum_{i=1}^H w_{hi} \cdot x_i$$

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

?

$$X(w) = \sum w_i x_i$$

AutoGrad

$a = \text{NumberWithGrad}(3)$

$$b = a + 4$$

$$c = b + 3$$

$$d = c + (a + 2)$$

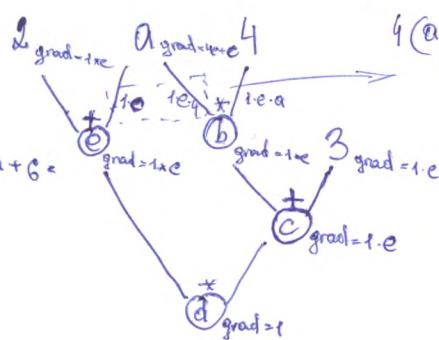
$$df/dx = (b+3)(a+2) = (4a+3)(a+2) = 4a^2 + 8a + 3a + 6 = 8a + 18$$

$$\left[\frac{df}{dg} \frac{dg}{dx} + \frac{df}{dh} \frac{dh}{dx} ? \right]$$

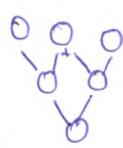
numa $f(g(a), h(c))$?



DL from scratch p.146
code



$$4(a+2) + (b+3) = 4a + 8 + 4a + 3 = 8a + 11$$

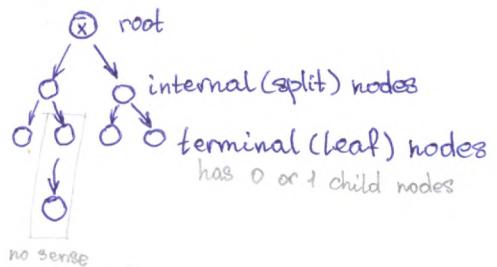


Почему же эти операции не могут производить градиенты для каждого из нейронов?

Decision Tree

Decision Tree ML Algorithm - набор связных классификаторов, организованных в иерархическую структуру - дерево.

Дерево - связной ориентированный граф без циклов, в котором каждая вершина имеет не более 1-го рёбра.

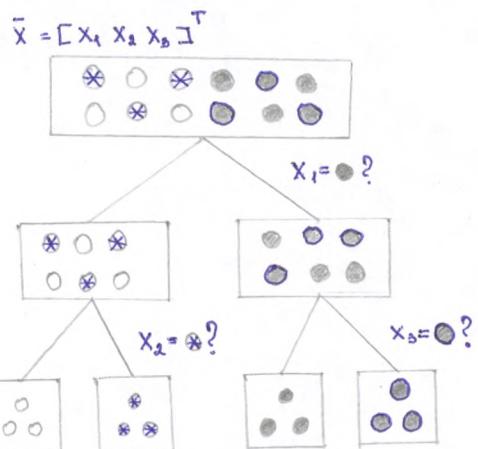


Связный классификатор - это underfit-модель классификаторов, который имеет мало параметров, высокую ошибку, опирается на небольшое кол-во features, быстро вычисляется.

Главное преимущество Decision Tree в том, что для него не нужно весь вектор характеристик (\vec{x}), в отличии от регрессии, то есть когда есть стоимость вычисления каждой характеристики, можно минимизировать издерожки \leftrightarrow максимизировать прибыль.

Например, в случае диагностической регрессии для диагностирования заболевания нужно как будто сдать все возможные лаб. тесты. Но, если все сдали 1 и всё ок, то возможно оставшиеся 500 сдавать не нужно.

Decision Tree легко интерпретируются.



Используют 3 параметра из 3-х.

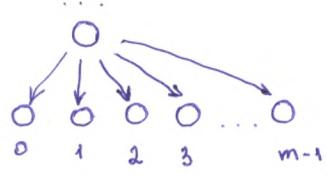
Decision Tree. Internal (split) nodes.

Split function - разделяющая функция.

$$h(\bar{x}, \theta_j) : R^D \times T \rightarrow \{0, m-1\}$$

args Domain Range

$\theta_j \in T$ - параметр разделяющей функции



Split function - функция, определяющая следующий узел из m номеров j-го узла.

If $m=2 \rightarrow$ Binary Decision Tree

$$h(\bar{x}, \theta_j) : R^D \times T \rightarrow \{0, 1\}$$



Feature Selection Function

$$\tilde{\psi}(\bar{x}) : R^D \rightarrow R^{D'}, D' \leq D$$

(ps:)

$$\bar{x} = [x_1 \ x_2 \ x_3 \ x_4 \dots x_D]^T$$

$$\tilde{\psi}(\bar{x}) = [x_2 \ x_3]^T, D' = 2$$

$\tilde{\psi}(\bar{x})$ - гиперпараметр,
своя отдачная для каждого внутреннего узла,
благодаря которому не используется весь feature vector \bar{x} .

$\tilde{\psi}(\bar{x}) = x_i$ - частотный случай.

Количество используемых характеристики - гиперпараметр,
который что имеет лучший результат будем давать выбор всех,
а что будем выбирать все.

Такие unused характеристики - это параметр.

Decision Tree . 3 types of split function.

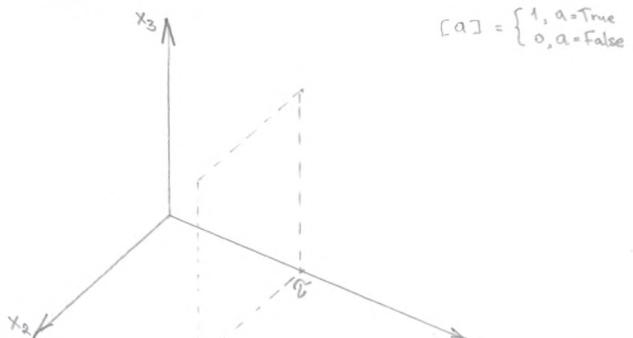
1) Default

$$h(\bar{x}, \theta_j) = [\underbrace{\varphi_1}_{-\infty} > \psi(\bar{x}) > \underbrace{\varphi_2}_{\infty}]$$

$$\theta_j = (\psi, \varphi_1, \varphi_2), \quad \psi(\bar{x}) = x_i$$

Гиперплоскость параллельная всем осям координат, кроме одной.

i.e. $x_i < \varphi$ (brute force)



$$[a] = \begin{cases} 1, & a = \text{True} \\ 0, & a = \text{False} \end{cases}$$

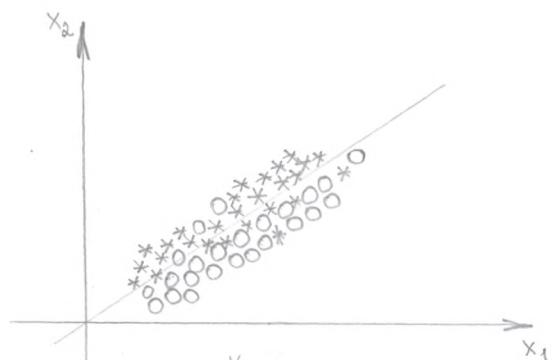
2) Linear

$$h(\bar{x}, \theta_j) = [\varphi_1 > \bar{\psi}(\bar{x})^T \bar{w} > \varphi_2]$$

$$\theta_j = (\bar{\psi}, \bar{w}, \varphi_1, \varphi_2)$$

$$\text{i.e. } x_1 w_1 + x_2 w_2 < \varphi$$

$$x_1 w_1 + \varphi < x_2 \quad (\text{regression})$$



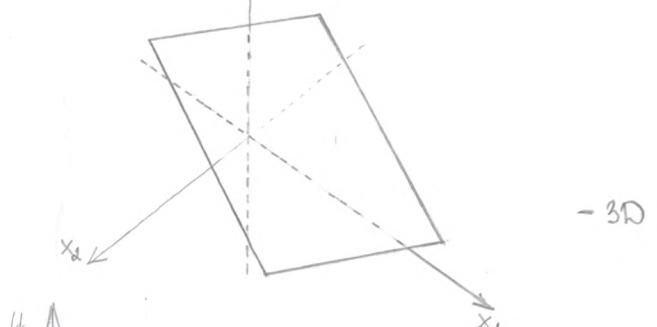
- 2D

3) Non - Linear

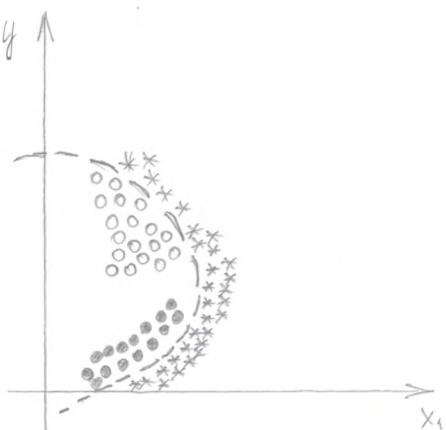
$$h(\bar{x}, \theta_j) = [\varphi_1 > \bar{\psi}(\bar{\psi}(\bar{x}))^T \bar{w} > \varphi_2]$$

$$\theta_j = (\bar{\psi}, \bar{w}, \varphi_1, \varphi_2)$$

ψ - basis functions.



- 3D



Задачи с бинарной меткой могут быть (may).

Decision Tree Training.

Notation

Classification classes $C_1, C_2 \dots C_K$

S_i - множество элементов класса i (часто dataset-a)

$S_i^{(k)}$ - множество элементов класса k , принадлежащих классу i .

$N_i = |S_i|$ - количество элементов множества S_i

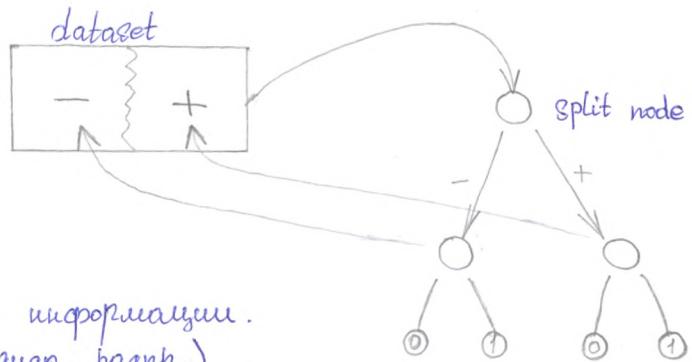
$$N_i^{(k)} = |S_i^{(k)}|$$

S_{ij} - множество элементов подмножества i принадлежащих классу j .

$$S_{ij}^{(k)}$$

$$N_{ij} = |S_{ij}|$$

$$N_{ij}^{(k)} = |S_{ij}^{(k)}|$$

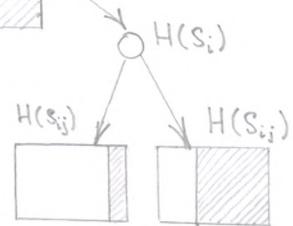
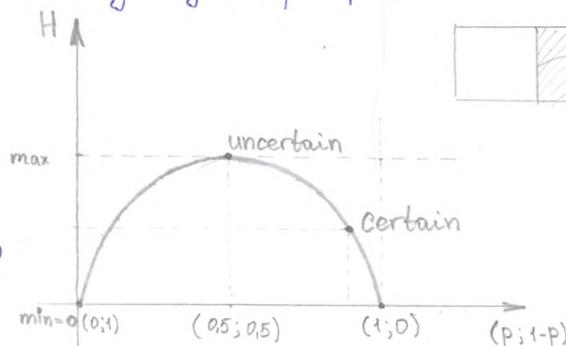


Entropy - оценка количества информации.
(Теория информации Гиббса, Бенара.)

$$H = - \sum_i p_i \log p_i$$

$$\frac{N_i^{(k)}}{N_i} = p^{(k)}$$

$$H(S_i) = - \sum_{k=1}^K \frac{N_i^{(k)}}{N_i} \cdot \log \frac{N_i^{(k)}}{N_i}$$



Information Gain
(Loss Function)

Информация при делении множества на подмножества классов, в худшем случае при uniform распределении.

$$IG = H(S_i) - \sum_j \frac{N_{ij}}{N_i} \cdot H(S_{ij}) \rightarrow \max_{\theta}$$

$$E(\theta) = \sum_j \frac{N_{ij}}{N_i} \cdot H(S_{ij}) \rightarrow \min_{\theta} \in [0; H(S_i)]$$

Функция $\theta = (\psi, \tau)$, которая минимизирует целевую функцию.

Минимизирует uncertainty на child узлах.

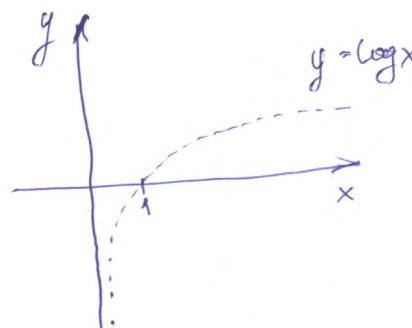
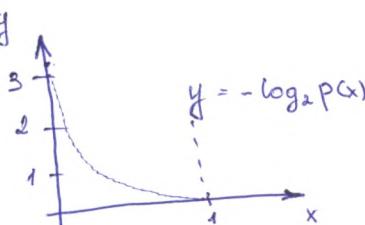


Information theory

Information

$$h(x) = -\log_2(p(x))$$

shannon, e-nats
rare events have more information.
data compression uncertainty.

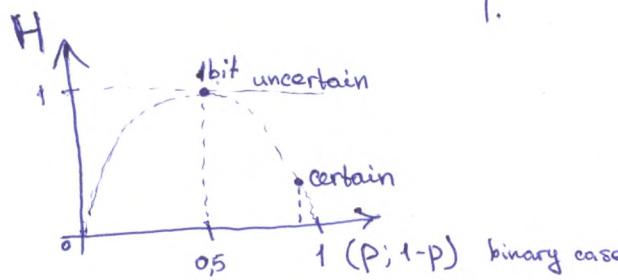


Entropy

$$H(x) = - \sum_i p_i(x) \cdot \log_2(p_i(x))$$

shannon Entropy

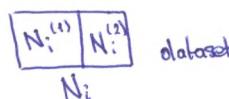
$$\frac{1}{2} \cdot \log_2 \frac{1}{2} = \frac{1}{2} \cdot (\log_2 1 - \log_2 2) = \\ = \frac{1}{2} (0 - 1) = -\frac{1}{2}$$



$$\frac{N_i^{(k)}}{N_i} = P^{(k)}$$

$$H(S_i) = \sum_{k=1}^K \frac{N_i^{(k)}}{N_i} \cdot \log_2 \frac{N_i^{(k)}}{N_i}$$

Деңгээлүү информация,
uncertain - ишүү информация.



$$H_i = \sum_{j=1}^{N_i} \frac{N_{ij}}{N_i} \cdot H(S_{ij})$$

$$E(\theta) = \sum_j \frac{N_{ij}}{N_i} \cdot H(S_{ij}) \rightarrow \min_{\theta}$$

минимизуя информация uncertainty на child уровне, но есть certain (90% ишүү мөнчөк)

$$1) H_i(0,5, 0,5) = 1$$

$$IG = 1 - \sum_i 0,5 \cdot 1 = 0$$



$$2) IG = 1 - (0,1 \cdot 0 + 0,9 \cdot H_{ij}(0,3, 0,7)) = 1 - (0,9 \cdot (0,3 \log_2 0,3 + 0,7 \log_2 0,7)) = 1 - 0,95 \boxed{0,05}$$

~~IG < 0~~ $\Rightarrow IG \in [0; H(S_i)]$

3)

$$\begin{array}{|c|c|} \hline 100 & \\ \hline 0,5 & 0,5 \\ \hline \end{array} - 50/50$$

$$\begin{array}{|c|c|} \hline 20 & 80 \\ \hline 0,7 & 0,3 \\ \hline 14 & 6 \\ \hline \end{array}$$

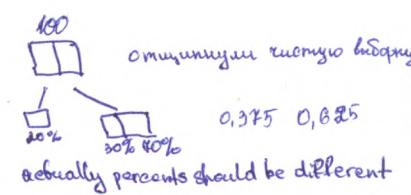
$$\begin{array}{|c|c|} \hline 80 & \\ \hline 0,15 & 0,85 \\ \hline 36 & 44 \\ \hline \end{array}$$

$$H(S_i) = - \sum p(x_i) \log_2 p(x_i) = -\left(\frac{1}{2} \cdot \log_2 \frac{1}{2}\right) = 1$$

$$IG = 1 - (0,2 \cdot H(S_{i1}) + 0,8 \cdot H(S_{i2})) = 1 - 0,9696 \boxed{0,03}$$

$$- 0,2 (0,7 \log_2 0,7 + 0,3 \log_2 0,3) = 0,2 \cdot 0,88$$

$$- 0,8 (0,15 \log_2 0,45 + 0,85 \log_2 0,55) = 0,8 \cdot 0,9928$$



Decision Tree. Terminal Node.

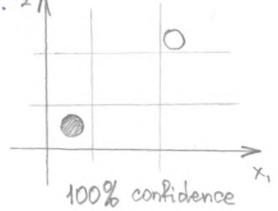
Stopping Criteria (to create terminal node)

- Достигнута заданная глубина d
- $N_i < n$
- $H(S_i) < \epsilon$



If $N_i = N_i^{(k)}$ - Было бы очень глубокое переобученное дерево, с малым кол-вом элементов на терминальных узлах. Мало элементов на внутренних узлах \rightarrow overfitting.

Random Forest: overfit + overfit = ok



$$1) \hat{y} = \underset{k}{\operatorname{argmax}} N_i^{(k)} \rightarrow \{0, 1, \dots, K\}$$

$$2) \vec{\hat{y}} = \left(\frac{N_i^{(1)}}{N_i}, \frac{N_i^{(2)}}{N_i}, \frac{N_i^{(3)}}{N_i}, \dots, \frac{N_i^{(K)}}{N_i} \right)^T$$

Большая уверенность наводит на overfitting.

Random Forest

Random Forest - этоансамбль деревьев решений.

Ensemble classifier - это объединение нескольких классификаторов.

Если мы объединим 6 один алгоритм несколько разных алгоритмов, то это будет называться fusion classifier.

С помощью сильных классификаторов RF может описывать нелинейные данные.

M деревьев: T_1, T_2, \dots, T_M

$$RF(\bar{x}) = RF(T_1(\bar{x}), T_2(\bar{x}), \dots, T_M(\bar{x})) = \frac{1}{M} \sum_{j=1}^M T_j(\bar{x})$$

$$T_i(\bar{x}) = \left(\frac{N_i^{(1)}}{N_i} \quad \frac{N_i^{(2)}}{N_i} \quad \dots \quad \frac{N_i^{(K)}}{N_i} \right), \quad \bar{x} \in S_i \text{ при обучении}$$

Проблема: детерминированный процесс обучения Decision Tree на одинаковой обучающей выборке даёт одинаковые деревья.

$$T_i = T_j, \quad \forall i, j \Rightarrow RF(\bar{x}) = T_1(\bar{x})$$

Bagging (Bootstrap Aggregation)

Метод во взятии случайного подмножества обучающей выборки.

Random Sampling with replacement, то есть можно брать один и тот же элемент несколько раз.

$$\{1, 2, 3, 4, 5, 6, 7\} \rightarrow \{3, 2, 1, 5, 5, 2, 3\}$$

Random Node Optimization

Выбор параметров split function из случайного подмножества параметров разделяния.

$$IG \rightarrow \max_{\Theta_i \in T_i}, \quad T_i \subset T$$

$P = |T_i|$ - "Количество случайности", комбинация параметров Θ_i .

$$\Theta_1 = (x_1, \gamma_1)$$

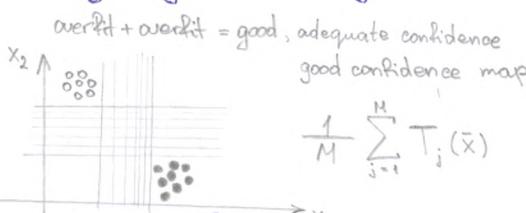
$$\Theta_2 = (x_2, \gamma_2)$$

$$\Theta_3 = (x_1, \gamma_3)$$

$$\Theta_4 = (x_1, \gamma_4)$$

$$\Theta_5 = (x_2, \gamma_5)$$

$\rightarrow \arg \max_{\Theta_i} IG$ - выбираем комбинацию, дающую максимальное значение IG.



$(n+k-1)!$	unique sets
$k! (n-1)!$	
ABC	$\binom{n+k-1}{k}$
ABCA	permutations
4!	$(2+3-1) = \frac{4!}{2! \cdot 2!}$
0000	
00	
000	
0000	

- $P=1$, оптимизация нет, так себе деревья, но разнообразие
- $P=n$, если n большое, деревья будут близже к оптимальному
- $P=M$, если перебирает все значения, то деревья одинаковые.

Чем больше кол-во деревьев, M, тем лучше, но после какого-то момента не будем давать улучшений, а только нагрузжать.

AdaBoost . Adaptive Boosting

Несмотря на Random Forest в том, что он принаследует
к классу деревья. RF не имеет возможности обратить внимание
на то, что есть особо важные примеры, чтобы закрыть пробелы.

Углубленная обучающая выборка

Training samples	Sample weight	Class labels
x_i	w_i	t_i
x_N	w_N	t_N
\bar{x}	\bar{w}	\bar{t}

Числовая функция

Функция:

$$E(\theta) = \frac{1}{2} \sum_{i=1}^N w_i (t_i - y(\bar{x}_i; \theta))^2$$

На всех элементах, на которых значение веса, не является равно нулю.

Классификация:

$$E(\theta) = - \sum_{i=1}^N w_i \sum_{k=1}^K t_i^{(k)} \ln y(\bar{x}_i; \theta)$$

Дерево решений:

$$IG_w = H_w(S_i) - \sum_j \frac{\tilde{N}_{ij}}{N_i} H_w(S_{ij}) \rightarrow_{\theta} \max$$

$$H_w(S_i) = - \sum_k \frac{\tilde{N}_i^{(k)}}{N_i} \log \frac{\tilde{N}_i^{(k)}}{N_i}$$

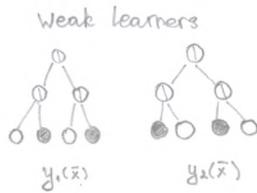
$\tilde{N}_i = \sum_{x_i \in S_i} w_i$ — не количество элементов, а сумма весов элементов.
Например, если вес нулевой, он никак не влияет на числовую функцию.

$$\vec{t}_i = \left(\frac{\tilde{N}_i^{(1)}}{N_i} \quad \frac{\tilde{N}_i^{(2)}}{N_i} \quad \frac{\tilde{N}_i^{(3)}}{N_i} \dots \frac{\tilde{N}_i^{(K)}}{N_i} \right)^T$$

AdaBoost. Рекурсивная загара.

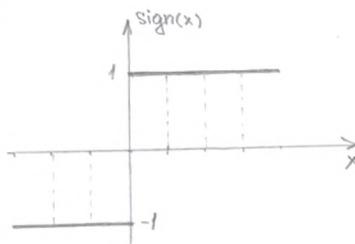
Бинарная загара классификации:

$$C_1, C_2; t_i \in \{-1, 1\}, i=1, 2, \dots, N$$



Следие классификации:

$$y_j(\bar{x}) \Rightarrow \{-1, 1\}, j=1, 2, \dots, M$$



Умоготий классификации:

$$Y(\bar{x}) = \text{sign} \left(\sum_{j=1}^M d_j y_j(\bar{x}) \right)$$

↑
weight of
weak classifier

Если d_j отрицательно, то он переворачивается

Минимизизация бояс обогащений бобошки

На мере:

$$w_i^{(j)} \geq 0, \sum_{i=1}^N w_i^{(j)} = 1$$

$$w_i^{(j)} = \frac{1}{N}, i=1, 2, \dots, N$$

for balanced dataset.

$$\text{or } w_i^{(j)} = \frac{1}{2N_k}, \text{ for imbalanced dataset.}$$

Минимизируемая целевая функция

Основная идея, максимизировать правильные.

$$J_j = \sum_{i=1}^N w_i^{(j)} I(y_j(\bar{x}_i) \neq t_i), \text{ если } y \neq t, \text{ то бояс сущесвует, если нет, то не бояс.}$$

Это не линейное бояс, а бобошка из exponential loss.

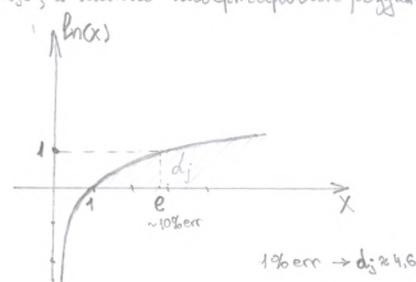
Ошибка каждого классификатора

$$\varepsilon_j = \sum_{i=1}^N w_i^{(j)} \cdot I(y_j(\bar{x}_i) \neq t_i), \text{ если } \varepsilon_j > \frac{1}{2}, \text{ то ошибок много, а много навернется бояс.}$$

best 1) $\ln \frac{1-0}{0} = \ln \infty = \infty$

random 2) $\ln \frac{1-\frac{1}{2}}{\frac{1}{2}} = \ln 1 = 0$

worst 3) $\ln \frac{1-1}{1} = \ln 0 = -\infty$



$$1\% \text{ err} \rightarrow d_j \approx 4.6$$

Обновление бояс обогащений бобошки

$$w_i^{(j)}, i=1, 2, \dots, N$$

$$w_i^{(j+1)} = \frac{w_i^{(j)} \cdot e^{d_j \cdot I(y_j(\bar{x}_i) \neq t_i)}}{Z_j}, i=1, 2, \dots, N$$

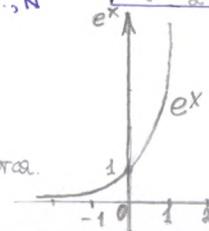
$$Z_j = \sum_{i=1}^N w_i^{(j)} e^{d_j \cdot I(y_j(\bar{x}_i) \neq t_i)}$$

сумма убогащения бояса,

и те изменения на боясе это трабас бояса убогащения.

if $\varepsilon = \frac{1}{2} \rightarrow \text{stop learning, nothing gonna change.}$

$$\boxed{\varepsilon_j \leq \frac{1}{2} \rightarrow d_j \geq 0}$$



$$\ln \frac{1}{\varepsilon} = \ln \frac{1}{\frac{1}{2}} < 0$$

$$\ln \frac{4}{5} = \ln 4 > 0$$

Учебная функция

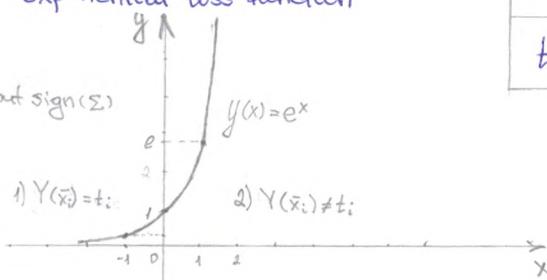
Учебная функция

$$E_0 = \sum_{i=1}^N e^{-\frac{1}{2}t_i Y_N^*(\bar{x}_i)} \quad - \text{exponential loss function}$$

$$Y_N^*(\bar{x}_i) = \sum_{j=1}^M d_j y_j(\bar{x}_i) \quad - \text{without sign} (\Sigma)$$

$$\cdot e^{-\frac{1}{2}t_i Y_N^*(\bar{x}_i)} ; t_i \rightarrow \{-1; 1\}$$

1) $Y(\bar{x}_i) = t_i$
2) $Y(\bar{x}_i) \neq t_i$



1) $\forall y_j(\bar{x}_i) \neq t_i$
 $e^{+\frac{1}{2}|Y_N^*(\bar{x}_i)|}$

y_1, y_2, \dots, y_{m-1} - already given weak learners

$y_m - ?$ - AdaBoost is greedy algorithm, so it cares only about currently last classifier.

$$\begin{aligned} E_0 &= \sum_{i=1}^N e^{-\frac{1}{2}t_i Y_m^*(\bar{x}_i)} = \sum_{i=1}^N e^{-\frac{1}{2}t_i (Y_{m-1}^*(\bar{x}_i) + d_m y_m(\bar{x}_i))} = \sum_{i=1}^N e^{-\frac{1}{2}t_i Y_{m-1}^*(\bar{x}_i) - \frac{1}{2}t_i d_m y_m(\bar{x}_i)} = \\ &= \sum_{i=1}^N \underbrace{e^{-\frac{1}{2}t_i Y_{m-1}^*(\bar{x}_i)}}_{\tilde{w}_i^{(m)}} \cdot e^{-\frac{1}{2}t_i d_m y_m(\bar{x}_i)} = \sum_{i=1}^N \tilde{w}_i^{(m)} \cdot e^{-\frac{1}{2}t_i d_m y_m(\bar{x}_i)} = \underbrace{\sum_{i=1}^N \tilde{w}_i^{(m)} e^{-\frac{1}{2}t_i d_m y_m(\bar{x}_i)}}_{\tilde{Z}_m} \rightarrow \min_{d_m y_m} \quad E \\ \tilde{w}_i^{(m)} &= \frac{\tilde{w}_i^{(m)}}{\tilde{Z}_m} ; \quad \tilde{w}_i^{(m)} > 0 ; \quad \sum_{i=1}^N \tilde{w}_i^{(m)} = 1 \\ \tilde{Z}_m &= \sum_{i=1}^N \tilde{w}_i^{(m)} \Rightarrow \tilde{Z}_m \cdot \tilde{w}_i^{(m)} = \tilde{Z}_m \cdot \frac{\tilde{w}_i^{(m)}}{\tilde{Z}_m} = \tilde{w}_i^{(m)} \end{aligned}$$

$$E = \sum_{i=1}^N w_i^{(m)} \cdot e^{-\frac{1}{2}t_i d_m y_m(\bar{x}_i)} = \sum_{i=1}^N w_i^{(m)} \cdot e^{-\frac{1}{2}t_i d_m y_m(\bar{x}_i)} \underbrace{\mathbb{I}(y_m(\bar{x}_i) \neq t_i)}_{\text{Incorrectly classified}} + \sum_{i=1}^N w_i^{(m)} \cdot e^{-\frac{1}{2}t_i d_m y_m(\bar{x}_i)} \cdot \underbrace{\mathbb{I}(y_m(\bar{x}_i) = t_i)}_{\text{Correctly classified}} =$$

$$= \sum_{i=1}^N w_i^{(m)} \underbrace{e^{+\frac{1}{2}d_m}}_{\text{const}} \cdot \mathbb{I}(y_m(\bar{x}_i) \neq t_i) + \sum_{i=1}^N w_i^{(m)} \underbrace{e^{-\frac{1}{2}d_m}}_{\text{const}} \cdot \mathbb{I}(y_m(\bar{x}_i) = t_i) \pm e^{-\frac{1}{2}d_m} \sum_{i=1}^N w_i^{(m)} \mathbb{I}(y_m(\bar{x}_i) \neq t_i) =$$

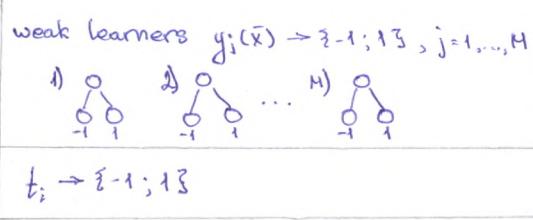
$$= (e^{\frac{1}{2}d_m} - e^{-\frac{1}{2}d_m}) \sum_{i=1}^N w_i^{(m)} \mathbb{I}(y_m(\bar{x}_i) \neq t_i) + e^{-\frac{1}{2}d_m} \sum_{i=1}^N w_i^{(m)} (\mathbb{I}(y_m(\bar{x}_i) = t_i) + \mathbb{I}(y_m(\bar{x}_i) \neq t_i)) =$$

$$= (e^{\frac{d_m}{2}} - e^{-\frac{d_m}{2}}) \sum_{i=1}^N w_i^{(m)} \mathbb{I}(y_m(\bar{x}_i) \neq t_i) + e^{-\frac{d_m}{2}} \rightarrow \min_{d_m y_m}$$

$$\int = \sum_{i=1}^N w_i^{(m)} \cdot \mathbb{I}(y_m(\bar{x}_i) \neq t_i) = \varepsilon_m \Rightarrow E_{dm} = (e^{\frac{d_m}{2}} - e^{-\frac{d_m}{2}}) \varepsilon_m + e^{-\frac{d_m}{2}} \rightarrow \min_{d_m}$$

$$E_{dm} = \left(\frac{1}{2} e^{\frac{d_m}{2}} + \frac{1}{2} e^{-\frac{d_m}{2}} \right) \varepsilon_m + \left(-\frac{1}{2} e^{-\frac{d_m}{2}} \right) = 0 ; \left(e^{\frac{d_m}{2}} + e^{-\frac{d_m}{2}} \right) \varepsilon_m = e^{\frac{d_m}{2}} \cdot e^{\frac{d_m}{2}} ; (e^{\frac{d_m}{2}} + 1) \varepsilon_m = 1 ; e^{\frac{d_m}{2}} + 1 = \frac{1}{\varepsilon_m} ; e^{\frac{d_m}{2}} = \frac{1}{\varepsilon_m} - \frac{\varepsilon_m}{\varepsilon_m}$$

$$d_m = \ln \frac{1 - \varepsilon_m}{\varepsilon_m}$$



Mameuamurecne ocuofor AdaBoost. (2)

$$E_0 = \sum_{i=1}^N e^{-\frac{1}{2}t_i Y_m^*(\bar{x}_i)}$$

$$Y_m^*(\bar{x}_i) = \sum_{j=1}^M d_j y_j(\bar{x}_i)$$

$$y_1, y_2, \dots, y_{m-1}$$

$$y_m - ?$$

$$E_0 = \sum_{i=1}^N e^{-\frac{1}{2}t_i Y_m^*(\bar{x}_i)} = \sum_{i=1}^N e^{-\frac{1}{2}t_i (Y_{m-1}^*(\bar{x}_i) + d_m y_m(\bar{x}_i))} = \sum_{i=1}^N \underbrace{e^{-\frac{1}{2}t_i Y_{m-1}^*(\bar{x}_i)}}_{\tilde{W}_i^{(m)}} \cdot e^{-\frac{1}{2}t_i d_m y_m(\bar{x}_i)} =$$

$$= \sum_{i=1}^N \tilde{W}_i^{(m)} e^{-\frac{1}{2}t_i d_m y_m(\bar{x}_i)}$$

$$\text{when } m=1 \rightarrow \tilde{W}_i^{(m)} = e^0 = 1 \Rightarrow \tilde{W}_i^{(m)} = \frac{1}{\sum_{i=1}^N 1} = \frac{1}{N}$$

$$\left| \begin{array}{l} \tilde{W}_i^{(m)} = \frac{\tilde{W}_i^{(m)}}{\sum_m} \Rightarrow \tilde{W}_i > 0, \sum_{i=1}^N \tilde{W}_i = 1 \\ \sum_m = \sum_{i=1}^N \tilde{W}_i^{(m)} \end{array} \right.$$

$$\bullet \quad \tilde{W}_i^{(m)} = e^{-\frac{1}{2}t_i Y_{m-1}^*(\bar{x}_i)}$$

$$\bullet \quad \tilde{W}_i^{(m+1)} = e^{-\frac{1}{2}t_i Y_m^*(\bar{x}_i)} = e^{-\frac{1}{2}t_i (Y_{m-1}^*(\bar{x}_i) + d_m y_m(\bar{x}_i))} = e^{-\frac{1}{2}t_i Y_{m-1}^*(\bar{x}_i) - \frac{1}{2}t_i d_m y_m(\bar{x}_i)} =$$

$$= \underbrace{e^{-\frac{1}{2}t_i Y_{m-1}^*(\bar{x}_i)}}_{\tilde{W}_i} \cdot e^{-\frac{1}{2}t_i d_m y_m(\bar{x}_i)} = \tilde{W}_i^{(m)} \cdot e^{-\frac{1}{2}t_i d_m y_m(\bar{x}_i)} = \dots$$

$$\left| \begin{array}{l} t_i y_m(\bar{x}_i) = 1 - 2I(y_m(\bar{x}_i) \neq t_i) \\ \begin{cases} t_i \cdot y_m(\bar{x}_i) = 1, \text{ if } t_i = y_m(\bar{x}_i) \\ t_i \cdot y_m(\bar{x}_i) = -1, \text{ if } t_i \neq y_m(\bar{x}_i) \end{cases} \end{array} \right.$$

$$\dots = \tilde{W}_i^{(m+1)} = \tilde{W}_i^{(m)} \cdot e^{-\frac{1}{2}d_m(1 - 2I(y_m(\bar{x}_i) \neq t_i))} = \tilde{W}_i^{(m)} \cdot e^{-\frac{1}{2}d_m + I(y_m(\bar{x}_i) \neq t_i)} = \tilde{W}_i^{(m)} \cdot e^{-\frac{1}{2}d_m} \cdot e^{d_m I(y_m(\bar{x}_i) \neq t_i)}$$

$$\text{Norm: } W_i^{(m+1)} = \frac{\tilde{W}_i^{(m+1)}}{\sum_{l=1}^N \tilde{W}_l^{(m+1)}} = \frac{\tilde{W}_i^{(m)} \cdot (e^{-\frac{1}{2}d_m}) \cdot e^{d_m I(y_m(\bar{x}_i) \neq t_i)}}{\sum_{l=1}^N \tilde{W}_l^{(m)} \cdot (e^{-\frac{1}{2}d_m}) \cdot e^{d_m I(y_m(\bar{x}_i) \neq t_l)}} \cdot \frac{e^{\frac{1}{2}d_m}}{e^{\frac{1}{2}d_m}} = \frac{\tilde{W}_i^{(m)} \cdot e^{d_m I(y_m(\bar{x}_i) \neq t_i)}}{\sum_{l=1}^N \tilde{W}_l^{(m)} \cdot e^{d_m I(y_m(\bar{x}_i) \neq t_l)}} \cdot \frac{1}{\frac{\sum_m}{N}} =$$

$$= \frac{\tilde{W}_i^{(m)}}{\sum_m \frac{\tilde{W}_l^{(m)}}{\sum_m \tilde{W}_l^{(m)}}} \cdot e^{d_m I(y_m(\bar{x}_i) \neq t_i)}$$

$$W_i^{(m+1)} = \frac{W_i^{(m)} \cdot e^{d_m I(y_m(\bar{x}_i) \neq t_i)}}{Z_m}$$

$$1. f = \sum_{i=1}^N w_i^{(m)} \cdot I(y_m(\bar{x}_i) \neq t_i) = \varepsilon_m \rightarrow \min_{y_m}$$

$$2. d_m = \ln \frac{1 - \varepsilon_m}{\varepsilon_m}$$

$$3. W_i^{(m+1)} = \frac{W_i^{(m)} \cdot e^{d_m I(y_m(\bar{x}_i) \neq t_i)}}{Z_m}$$

$$W_i^{(1)} = \frac{1}{N} \text{ or } \frac{1}{2N_k} \text{ for imbalanced dataset.}$$

K-Means Clustering

Задача кластеризации - отнесение элементов выборки к одному из нескольких кластеров так, чтобы элементы из одного кластера были похожи, а из разных непохожи.

Данные: $\bar{x}_1, \bar{x}_2, \dots, \bar{x}_N$; $\bar{x} \in \mathbb{R}^D$

Кластеры: C_1, C_2, \dots, C_k

Центроиды: $\bar{\mu}_i \in \mathbb{R}^D$, $i = 1, 2, \dots, k$, место где каждого кластера свои центроиды.



Целевая функция

$$E = \sum_{k=1}^K \sum_{\bar{x} \in C_k} \|\bar{\mu}_k - \bar{x}\|^2$$

— Euclidean Distance (d) = $\sqrt{(x_1 - x_i)^2 + (y_2 - y_i)^2}$

Инициализация центроидов: Random, min-max range, dataset elements

Назначение элементов кластерам:

$$C_j = \{\bar{x}_i \mid d(\bar{x}_i; \bar{\mu}_j) = \min_k d(\bar{x}_i; \bar{\mu}_k), i = 1, 2, \dots, N\}, j = 1, 2, \dots, N$$

min distance no break.

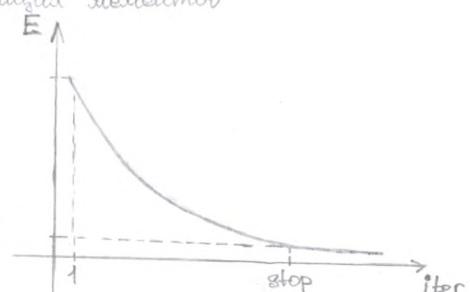
Обновление центроидов

$$\bar{\mu}_j = \frac{1}{|C_j|} \sum_{\bar{x} \in C_j} \bar{x}$$

— среднее значение всех входящих элементов
при-бо

Критерии остановки

- обновление центроидов прекратилось
- падение целевой функции норма прекратилось
- достигнуто заданное кол-во итераций



$$E = \sum_{k=1}^K \sum_{\bar{x} \in C_k} \|\bar{\mu}_k - \bar{x}\|^2 \rightarrow \min_{\mu}$$

1. $\bar{x} \rightarrow \mu$

2. k, C_k решение по центроиду кластера принимается независимо

$$E_k = \sum_{\bar{x} \in C_k} \|\bar{\mu}_k - \bar{x}\|^2 = \sum_{\bar{x} \in C_k} \sum_{i=1}^D (\mu_k^{(i)} - x^{(i)})^2 \rightarrow \min_{\mu}$$

$$\frac{dE_k}{\mu_k^{(i)}} = \sum_{\bar{x} \in C_k} 2(\mu_k^{(i)} - x^{(i)}) = 0$$

$$\sum_{\bar{x} \in C_k} (\mu_k^{(i)} - x^{(i)}) = \sum_{\bar{x} \in C_k} \mu_k^{(i)} - \sum_{\bar{x} \in C_k} x^{(i)} = 0 ; \sum_{\bar{x} \in C_k} \mu_k^{(i)} = |C_k| \mu_k^{(i)} = \sum_{\bar{x} \in C_k} x^{(i)} \Rightarrow$$

$$\Rightarrow \mu_k^{(i)} = \frac{1}{|C_k|} \cdot \sum_{\bar{x} \in C_k} x^{(i)} \Rightarrow \mu_j = \frac{1}{|C_j|} \cdot \sum_{\bar{x} \in C_j} \bar{x}$$

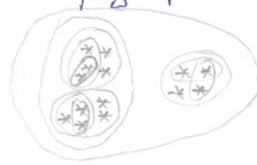
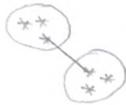
Agglomerative Clustering

$d(\bar{x}_i, \bar{x}_j) = R$ (the larger the norm of the vector difference: the greater the dissimilarity).

Норменное обединение расстояний в расстояние Балмера.

1. Single link

$$d(C_i, C_j) = \min_{x_i \in C_i, x_j \in C_j} d(x_i, x_j)$$



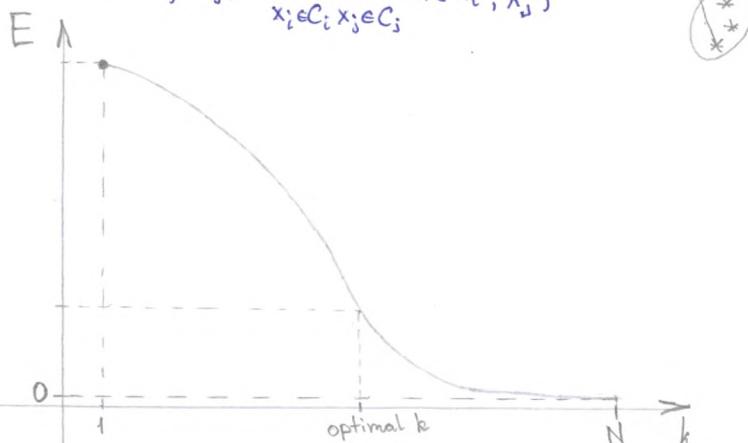
2. Average link

$$d(C_i, C_j) = \frac{1}{|C_i||C_j|} \cdot \sum_{x_i \in C_i} \sum_{x_j \in C_j} d(x_i, x_j)$$



3. Complete link

$$d(C_i, C_j) = \max_{x_i \in C_i, x_j \in C_j} d(x_i, x_j)$$



Basic Linear Algebra

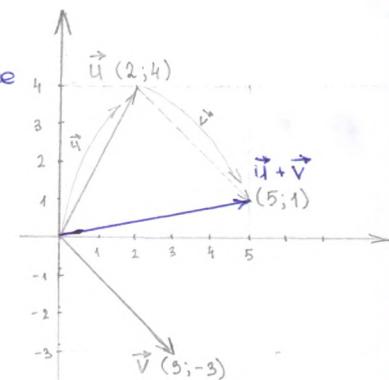
$\vec{u} + \vec{v}$ - elementwise addition of vectors, it's like stepping in space

$c \cdot \vec{v}$ - scaling of vector

The span of \vec{u} and \vec{v} is a set of all linear combinations:

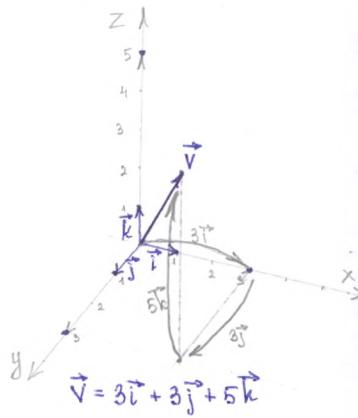
$$c_1 \vec{u} + c_2 \vec{v}$$

Linearly dependent vectors: $\vec{u} = c_1 \vec{v} + c_2 \vec{w}$



Basis vectors of a vector space is a set of linearly independent vectors $\{\vec{i}, \vec{j}, \vec{k}\}$ that span the full space:

$$\vec{v} = c_1 \vec{i} + c_2 \vec{j}$$



Linear Transformations

grid lines remain parallel and evenly spaced

$L(\vec{v}) \rightarrow$ transformed \vec{v}

We can describe any linear transformation for all vectors using the transformed basis vectors:

$$L(\vec{v}) = c_1 L(\vec{i}) + c_2 L(\vec{j}) = \begin{bmatrix} L(\vec{i})_x & L(\vec{j})_x \\ L(\vec{i})_y & L(\vec{j})_y \end{bmatrix} \begin{bmatrix} c_1 \\ c_2 \end{bmatrix}$$

$$L(\vec{v}) = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 1 & -1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$$

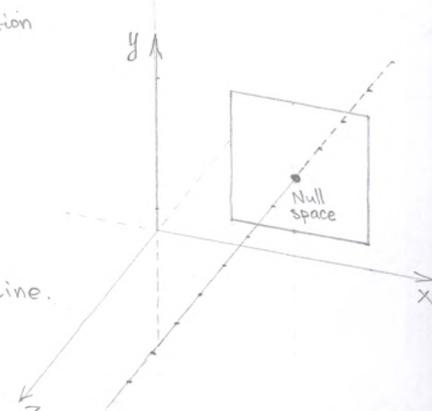
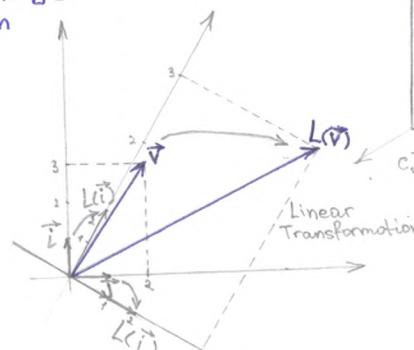
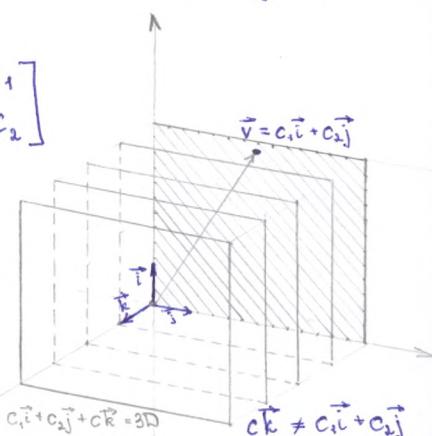
shear rotation composition

$$ABC = A(BC) = (AB)C$$

$$AB \neq BA$$

$$\begin{bmatrix} i_x & j_x \\ i_y & j_y \\ i_z & j_z \end{bmatrix} \quad 2D \rightarrow 3D \quad \text{Full rank 2} \quad \text{Det} = \infty$$

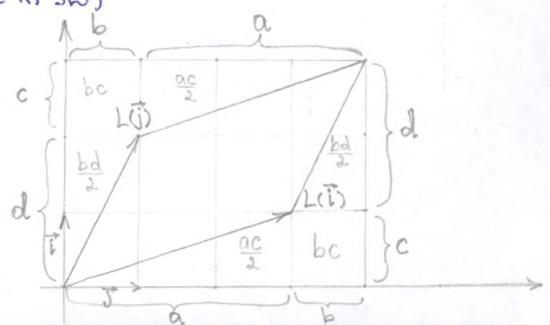
$$\begin{bmatrix} i_x & j_x & k_x \\ i_y & j_y & k_y \end{bmatrix} \quad 3D \rightarrow 2D \quad \text{Rank = 2} \quad \text{Column space = 2} \quad \text{Null space = 1, 1 point can be re-mapped to line. (kernel)}$$



Determinant

determinant - scale of area of linear transformation (or volume in 3D)

$$\det \begin{bmatrix} a & b \\ c & d \end{bmatrix} = ad + bc - ac - bd = (a+b)(c+d) - bc - bc - ac - bd = ad - bc$$



Dot product of vectors

$$\vec{a} \cdot \vec{b} = |\vec{a}| \cdot |\vec{b}| \cdot \cos(\hat{\vec{a}, \vec{b}})$$

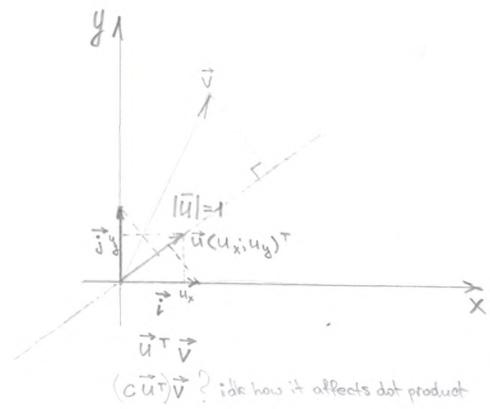
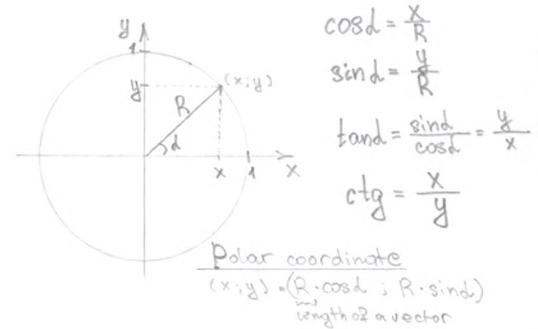
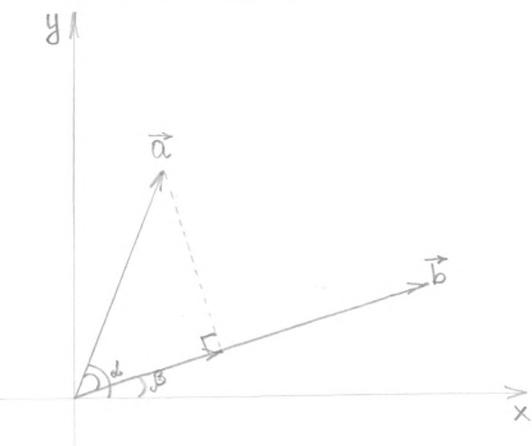
$$\vec{a} \cdot \vec{b} = \vec{a}^T \vec{b} = \vec{b}^T \vec{a} = [i \ j] \begin{bmatrix} x \\ y \end{bmatrix} = xi + yj$$

$$\vec{a} = |\vec{a}| \cdot (\cos d; \sin d)^T$$

$$\vec{b} = |\vec{b}| \cdot (\cos \beta; \sin \beta)^T$$

$$\begin{aligned}\vec{a}^T \vec{b} &= |\vec{a}| |\vec{b}| \cos d \cos \beta + |\vec{a}| |\vec{b}| \sin d \sin \beta = \\ &= |\vec{a}| |\vec{b}| \cos(d - \beta)\end{aligned}$$

$$\cos(d - \beta) = \cos d \cos \beta + \sin d \sin \beta$$



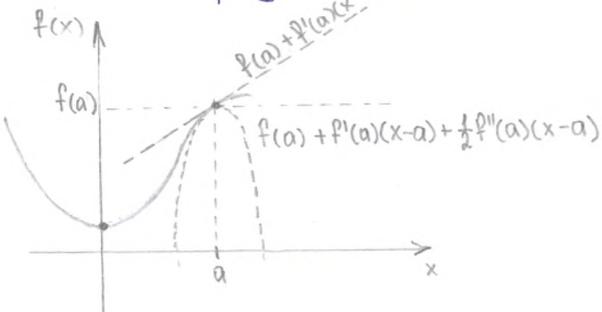
Taylor Series

$$f(x) = P(x) = \sum_{n=0}^{\infty} \frac{f^{(n)}(a)}{n!} (x-a)^n$$

$$P^{(n)}(a) = f^{(n)}(a)$$

Значения функции f и ее производных в точке a и производные высоких степеней обозначают (изменение изменения роста функции)

Approximation of non-polynomial functions (f) with polynomials (P).



$$P(x) = C_0 + C_1(x-a) + C_2(x-a)^2 + C_3(x-a)^3 + C_4(x-a)^4 + \dots$$

$$0) P(a) = C_0 + C_1 \underbrace{(a-a)}_0 = C_0 = f(a)$$

$$1) P^{(1)}(a) = (0 + C_1(1) + 2C_2 \underbrace{(a-a)}_0 + \dots) = C_1 = f'(a)$$

$$2) P^{(2)}(a) = 0 + 0 + 2C_2 + 3 \cdot 2 \cdot C_3 \underbrace{(a-a)}_0 + 4 \cdot 3 \cdot C_4 \underbrace{(a-a)}_0^2 + \dots = 2C_2 = f^{(2)}(a)$$

$$\Rightarrow C_2 = \frac{f^{(2)}(a)}{2}$$

$$3) P^{(3)}(a) = 0 + 0 + 0 + 3 \cdot 2 \cdot 1 C_3 + 4 \cdot 3 \cdot 2 \cdot C_4 \underbrace{(a-a)}_0 + \dots = 3 \cdot 2 \cdot 1 \cdot C_3 = f^{(3)}(a)$$

$$\Rightarrow C_3 = \frac{f^{(3)}(a)}{3!}$$

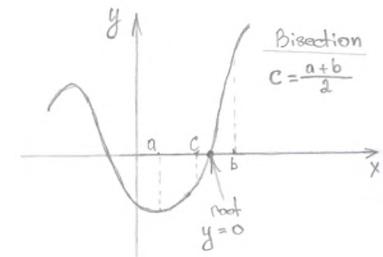
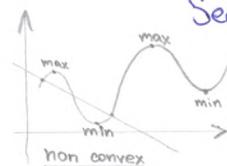
$$P(x) = f(a) + f'(a)(x-a) + \frac{1}{2} f^{(2)}(a)(x-a)^2 + \frac{1}{6} f^{(3)}(a)(x-a)^3 + \frac{1}{24} f^{(4)}(a)(x-a)^4 + \dots$$

Newton's Method

1. Root-finding method. Other: Bisection (Interval halving)

Secant, Brent's methods

2. Optimization method



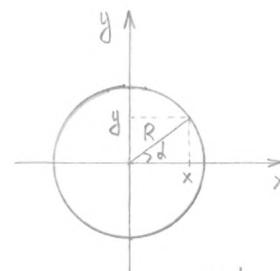
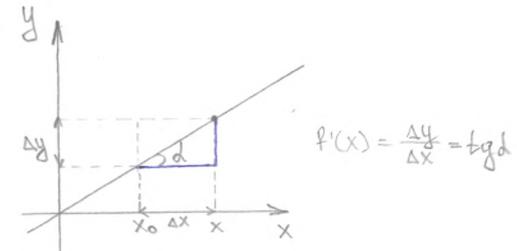
Derivative

$$y = f(x)$$

$$f'(x) = \frac{df}{dx} = \lim_{\Delta x \rightarrow 0} \frac{\Delta y}{\Delta x} = \lim_{\Delta x \rightarrow 0} \frac{f(x+\Delta x) - f(x)}{\Delta x}$$

$$\text{tgd} = \frac{\Delta y}{\Delta x} = \mathcal{V}$$

$$f'(x_0) = \lim_{\Delta x \rightarrow 0} \frac{f(x_0 + \Delta x) - f(x_0)}{\Delta x}$$



$$\begin{aligned}\sin d &= \frac{y}{R} \\ \cos d &= \frac{x}{R} \\ \text{tgd} &= \frac{\sin d}{\cos d} = \frac{y}{x} \\ \text{ctgd} &= \frac{x}{y}\end{aligned}$$

Tangent Line

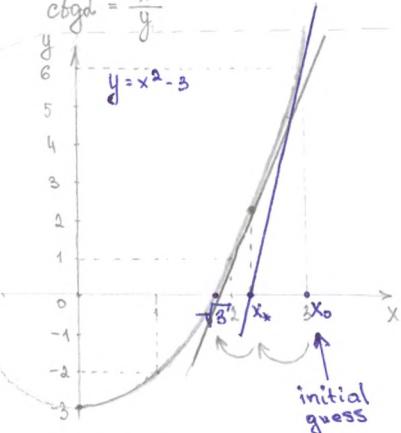
$$f'(x_0) = \lim_{\Delta x \rightarrow 0} \frac{f(x_0 + \Delta x) - f(x_0)}{\Delta x}$$

$$\Delta x = x - x_0$$

$$f'(x_0) = \frac{f(x_0 + (x - x_0)) - f(x_0)}{x - x_0} \Rightarrow f'(x_0)(x - x_0) = f(x) - f(x_0)$$

$$f(x) = f(x_0) + f'(x_0)(x - x_0)$$

$$S_0 + \mathcal{V} \cdot \Delta t$$



Newton's root finding method (iterative)

$$f'(x_0) = \frac{f(x) - f(x_0)}{x - x_0} \Rightarrow x - x_0 = \frac{f(x) - f(x_0)}{f'(x_0)}, \quad f'(x) \neq 0$$

$$f(x) = 0 - \text{root}$$

$$x_* = x_0 - \frac{f(x_0)}{f'(x_0)}, \quad \text{until } |x_{n+1} - x_n| < \epsilon;$$

Newton's optimization method

$y'(x) = 0$ - critical points (minima/maxima),
so we can try to find roots of $y'(x)$, i.e. $y'(x) = 0$.

$$x_* = x_0 - \frac{f'(x_0)}{f''(x_0)}$$

$$\begin{bmatrix} \frac{d^2 f}{d x_1^2} & \frac{d^2 f}{d x_1 d x_2} \dots & \frac{d^2 f}{d x_1 d x_D} \\ \frac{d^2 f}{d x_2 d x_1} & \frac{d^2 f}{d x_2^2} \dots & \frac{d^2 f}{d x_2 d x_D} \\ \dots & \dots & \dots \\ \frac{d^2 f}{d x_D d x_1} & \frac{d^2 f}{d x_D d x_2} \dots & \frac{d^2 f}{d x_D^2} \end{bmatrix}^{-1} \begin{bmatrix} \frac{df}{dx_1} \\ \frac{df}{dx_2} \\ \dots \\ \frac{df}{dx_D} \end{bmatrix} = \begin{bmatrix} \frac{df}{dx_1} \\ \frac{df}{dx_2} \\ \dots \\ \frac{df}{dx_D} \end{bmatrix}$$

$$\begin{bmatrix} \frac{d(f)}{d x_1} \\ \frac{d(f)}{d x_2} \\ \dots \\ \frac{d(f)}{d x_D} \end{bmatrix} = \begin{bmatrix} \frac{d(f)}{d x_1} \\ \frac{d^2 f}{d x_1^2} + \frac{d^2 f}{d x_1 d x_2} + \dots + \frac{d^2 f}{d x_1 d x_D} \\ \frac{d^2 f}{d x_2 d x_1} + \frac{d^2 f}{d x_2^2} + \dots + \frac{d^2 f}{d x_2 d x_D} \\ \dots \\ \frac{d^2 f}{d x_D d x_1} + \frac{d^2 f}{d x_D d x_2} + \dots + \frac{d^2 f}{d x_D^2} \end{bmatrix}$$

Multidimensional case:

$$\vec{x}_1 = \vec{x}_0 - H_{x_0}^{-1} \cdot \nabla f(\vec{x}_0)$$

Hessian Matrix

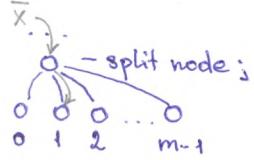
f twice-differentiable.

Decision Tree for Regression Task

- target variables: $t_1, t_2, \dots, t_N \in \mathbb{R}$ $t(\bar{x}) = y(\bar{x}) + \epsilon$

- sample: $\bar{x}_i \in \mathbb{R}^D$

- split function: $h(\bar{x}, \theta_j) : \mathbb{R}^D \times T \rightarrow \{0, 1, \dots, m-1\}$
 - $\theta_j \in T$ - split function parameters
 - usually $m=2$ i.e. Binary Decision Tree



- feature selection function: $\psi(\bar{x}) : \mathbb{R}^D \rightarrow \mathbb{R}^{D'}, D' \leq D$
 - usually $D'=1$ i.e. $\psi(\bar{x}) = x_i$

$$\bar{x} = [x_1, x_2, \dots, x_D]^T$$

$$\bar{\psi}(\bar{x}) = [x_2, x_5]^T, D'=2$$

- Split function types remain the same, for example, default hyperplane:

$$h(\bar{x}, \theta_j) = [\tilde{t}_1 > \psi(\bar{x}) > \tilde{t}_2] - \text{Indicator function}$$

$$\theta_j = (\psi, \tilde{t}_1, \tilde{t}_2)$$

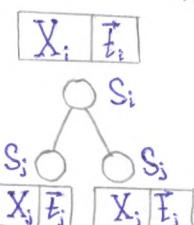
$$\psi(\bar{x}) = x_i$$

$[a] = \begin{cases} 1, & \text{if } a = \text{true} \\ 0, & \text{if } a = \text{false} \end{cases}$

Loss function

$$L = V(S_i) - \sum_j \frac{N_{ij}}{N_i} V(S_j) \rightarrow \min_{\theta} \quad (\text{we use variance instead of entropy})$$

$$V(S_i) = \frac{1}{N_i} \sum_{l=1}^{N_i} (t_i^{(l)} - \underbrace{m_i}_{\substack{\text{prediction} \\ \text{computed only once}}})^2$$



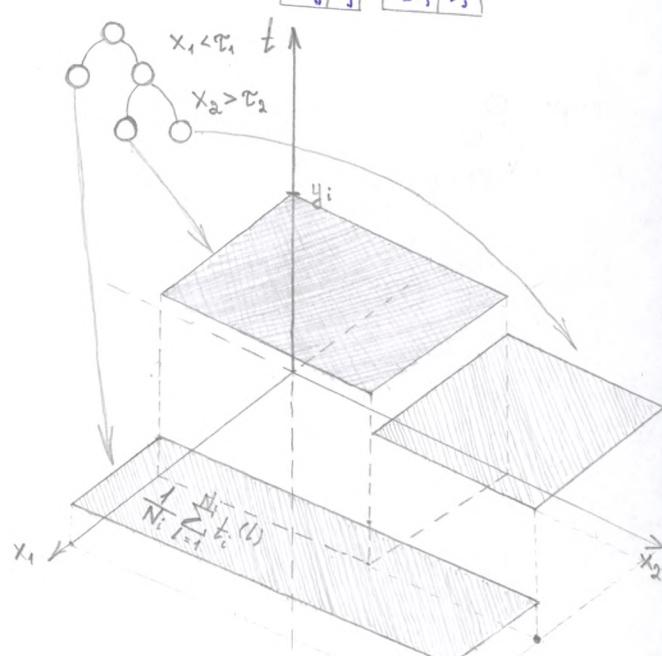
Stopping criteria to create terminal node

- depth d has been reached
- $N_i < n$
- $V(S_i) < \epsilon$

Terminal node (prediction)

$$y_i = \frac{1}{N_i} \sum_{l=1}^{N_i} t_i^{(l)} - \text{average}$$

additionally, we could return variance.



Gradient Boosting: Regression

Объединение нескольких персекторов в единий персектор.

Weak regressor - Decision Tree for Regression Task.
(weak learner in general)

Умственный персектор

$$F(\bar{x}) = \sum_{j=0}^M h_j f_j(\bar{x}) = f_0 + \sum_{j=1}^M h_j f_j(\bar{x}) \quad - \text{сумма умственных персекторов,}$$

$h_0 = 1 \qquad f_0 - \text{const}$

$h_j = h \in [0; 1]$
as learning rate

компьютерное представление остатка (residuals).

Loss Function (any differentiable loss function) $L(F_m(\bar{x}), t) = 0 \Leftrightarrow F_m(\bar{x}) = t, \forall i$
 $L_2 \text{ or } L_1$

$$L_2 = \frac{1}{2} \sum_{i=1}^N (F(\bar{x}_i) - t_i)^2$$

$\exists i : F_m(\bar{x}_i) \neq t_i \Rightarrow L > 0$
but $\bar{x}_i = \tilde{x}_i$ and $t_i \neq \tilde{t}_i \Rightarrow f(\bar{x}_i) = \frac{\tilde{t}_i + t_i}{2}$

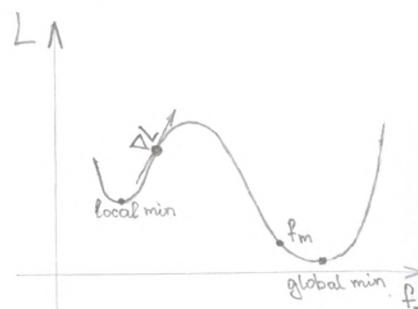
$$L \rightarrow \min_{f_0} ; \left(\frac{1}{2} \sum_{i=1}^N (f_0 - t_i)^2 \right)'_{f_0} = \sum_{i=1}^N (f_0 - t_i) = Nf_0 - \sum_{i=1}^N t_i = 0 \Rightarrow f_0 = \frac{1}{N} \sum_{i=1}^N t_i$$

$$\begin{aligned} L \rightarrow \min_{f_m} ; \quad L &= \frac{1}{2} \sum_{i=1}^N (F_m(\bar{x}_i) - t_i)^2 = \frac{1}{2} \sum_{i=1}^N (F_{m-1}(\bar{x}_i) + f_m(\bar{x}_i) - t_i)^2 = \\ &= \frac{1}{2} \sum_{i=1}^N (f_m(\bar{x}_i) - \underbrace{(t_i - F_{m-1}(\bar{x}_i))}_{\tilde{t}_i^{(m)} - \text{residuals}})^2 = \frac{1}{2} \sum_{i=1}^N (f_m(\bar{x}_i) - \tilde{t}_i^{(m)})^2 \end{aligned}$$

" y насторожа умственного персектора f_m свой набор членов переменных $\tilde{t}_i^{(m)}$ - residuals, совместно насторожа персектор f_m своим разрывом.

Gradient Boosting называется Gradient Descent ($f_m - (-\nabla L)$, $h - lr$), а сума остатков, это f_m называем global minima.

Поменявшись f_m можно overfit-симба, мы же хотим избежать
на ансамбле генерации & Random Forest.



Gradient Boosting: Classification

Умозрение на accuracy

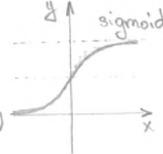
$$y(\bar{x}) = \delta(F(\bar{x})) = \delta\left(\underbrace{f_0}_{\text{const}} + d \sum_{j=1}^M f_j(\bar{x})\right), \quad f_j - \text{regression tree}$$

Sigmoid (δ)

$$\text{Softmax}(\bar{x}) = \left(\frac{e^{x_i}}{\sum_{j=1}^K e^{x_j}} \right), \quad K=2$$

$$\begin{aligned} \text{Softmax}(\bar{x}) &= \left(\frac{e^{x_0}}{e^{x_0} + e^{x_1}}, \frac{e^{x_1}}{e^{x_0} + e^{x_1}} \right) = \left(\frac{1}{1 + e^{x_1 - x_0}}, \frac{1}{e^{x_0 - x_1} + 1} \right) = |x_1 - x_0 = t| = \\ &= \left(\underbrace{\frac{1}{1 + e^t}}_{1 - \delta}, \underbrace{\frac{1}{1 + e^{-t}}}_{\delta} \right) \end{aligned}$$

$$\frac{1}{1 + e^t} = \frac{1 + e^{-t} - e^t}{1 + e^t} = 1 - \frac{e^t}{1 + e^t} = 1 - \frac{1}{e^{-t} + 1} = 1 - \delta(t) = \delta(-t)$$



Loss function

$$P(\vec{t} | X) = \prod_{n=1}^N \prod_{k=1}^K t_n^{(k)} y^{(k)}(\bar{x}_n)^{t_n^{(k)}} \rightarrow \max$$

$\vec{y} = [30\% \ 70\%]$ - probability vector

$$L = - \sum_{n=1}^N \sum_{k=1}^K t_n^{(k)} \ln y^{(k)}(\bar{x}_n) \rightarrow \min, \quad K=2$$

$$L = - \sum_{n=1}^N \sum_{k=1}^2 t_n^{(k)} \ln y^{(k)}(\bar{x}_n) = - \sum_{n=1}^N (t_n^{(0)} \ln y^{(0)}(\bar{x}_n) + t_n^{(1)} \ln y^{(1)}(\bar{x}_n)), \quad \delta(x) \in [0; 1] - \text{scalar}$$

$$L = - \sum_{n=1}^N ((1-t_n) \ln(1-\delta(\bar{x}_n)) + t_n \ln \delta(\bar{x}_n)) \rightarrow \min$$

$\overbrace{t_n = C_0}^1 \quad \overbrace{t_n = C_1}^1$ Одна из них должна быть

$$y(\bar{x}) = \delta(f_0)$$

$$L = - \sum_{n=1}^N ((1-t_n) \ln(1-\delta(f_0)) + t_n \ln \delta(f_0)) \rightarrow \min_{f_0}$$

$$L' = - \sum_{n=1}^N \left((1-t_n) \frac{1}{1-\delta(f_0)} \cdot (-\delta(f_0))' + t_n \frac{1}{\delta(f_0)} \cdot (\delta(f_0))' \right) =$$

$$= - \sum_{n=1}^N \left((1-t_n) \frac{1}{1-\delta(f_0)} \cdot (-\delta(f_0)(1-\delta(f_0))' + t_n \frac{1}{\delta(f_0)} \cdot \delta(f_0)(1-\delta(f_0))' \right) =$$

$$= - \sum_{n=1}^N ((1-t_n) \cdot (-\delta(f_0)) + t_n (1-\delta(f_0))) = \sum_{n=1}^N (-\delta(f_0) + t_n \delta(f_0) + t_n - t_n \delta(f_0)) =$$

$$= - \sum_{n=1}^N (-\delta(f_0) + t_n) = 0 \quad -\text{is it minimum? proof} \rightarrow \text{if } f'' = +C$$

$$0. (a^x)' = a^x \ln a \cdot (x)'$$

$$1. (\log_a f(x))' = \frac{1}{f(x) \ln a} \cdot f'(x)$$

$$\begin{aligned} 2. \delta'(x) &= (1 \cdot (1+e^{-x})^{-1})' = 0 + (-1(1+e^{-x}) \cdot (-e^{-x}))' = \\ &= -\frac{1}{(1+e^{-x})^2} \cdot (e^{-x})(-x)' = \\ &= \frac{e^{-x}}{(1+e^{-x})^2} = \frac{1 \cdot (e^{-x})}{(1+e^{-x})(1+e^{-x})} = \delta(x) \cdot \frac{-e^{-x} + 1 - 1}{1 + e^{-x}} = \\ &= \delta(x) \cdot \frac{1 + e^{-x} - 1}{1 + e^{-x}} = \boxed{\delta(x)(1 - \delta(x))} \end{aligned}$$

$$\uparrow \downarrow \quad \left(\sum_{n=1}^N \delta(f_0) - t_n \right)' = \sum_{n=1}^N \delta(f_0)(1 - \delta(f_0)) = +C$$

$$\sum_{n=1}^N \delta(f_0) = \sum_{n=1}^N t_n = N \delta(f_0)$$

$$\delta(f_0) = \frac{1}{N} \sum_{n=1}^N t_n = \frac{1}{1 + e^{f_0}} = p_1 \Rightarrow$$

bias, dataset balance

$$1 + e^{f_0} = \frac{1}{p_1}$$

$$e^{f_0} = \frac{1}{p_1} - 1 = \frac{1 - p_1}{p_1} = \frac{p_0}{p_1}$$

$$f_0 = \ln \frac{p_0}{p_1}$$

Gradient Boosting: Classification (2)

Умовний наступником

$$y(\bar{x}) = \delta(F(\bar{x})) = \delta(f_0 + d \sum_{j=1}^M f_j(\bar{x})) , \quad f_j - \text{regression tree}$$

$$\delta(f_0) = \frac{1}{N} \sum_{n=1}^N t_n \longleftrightarrow f_0 = \ln \frac{p_0}{p_1}$$

Loss function

$$L = - \sum_{n=1}^N ((1-t_n)(\ln(1-y(\bar{x}_n)) + t_n \ln y(\bar{x}_n))) \rightarrow \min$$

Given: f_0, f_1, \dots, f_{m-1}

f_m - ?

$$\begin{cases} F_m(\bar{x}_i) = F_{m-1}(\bar{x}_i) + d f_m(\bar{x}_i) \\ F_m(\bar{x}_i) = F_{m-1}(\bar{x}_i) - \delta L'(\delta(F_{m-1}(\bar{x}_i)))_{F_{m-1}} \end{cases}$$

$L'(\delta(F_{m-1}(\bar{x}_i)))$ - спосіб пошуку цільової функції L щодо $F_{m-1}(\bar{x}_i)$,
 $\underset{F_{m-1}}{\text{у змозі}} \text{ у змозі} \underset{L}{\text{зменшити}} \underset{F_{m-1}}{\text{показано}} \underset{F_{m-1}(\bar{x}_i)}{\text{справжні}} \underset{\text{напрямлення}}{\text{напрямлення}}.$

$$\begin{aligned} -L'(\delta(F_{m-1}(\bar{x}_i))) &= ((1-t_i) \ln(1-y(\bar{x}_i)) + t_i \ln y(\bar{x}_i))'_{F_{m-1}} = (\delta(F_{m-1}(\bar{x}_i)))'_{F_{m-1}} = 1, \text{ as } x' = 1 \\ &= ((1-t_i) \ln(1-\delta(F_{m-1}(\bar{x}_i))) + t_i \ln \delta(F_{m-1}(\bar{x}_i)))'_{F_{m-1}} = \\ &= (1-t_i) \frac{1}{1-\delta(F_{m-1}(\bar{x}_i))} \cdot (-\delta(F_{m-1}(\bar{x}_i))) (1-\delta(F_{m-1}(\bar{x}_i))) + t_i \frac{1}{\delta(F_{m-1}(\bar{x}_i))} \cdot \delta(F_{m-1}(\bar{x}_i)) (1-\delta(F_{m-1}(\bar{x}_i))) = \\ &= (1-t_i) (-\delta(F_{m-1}(\bar{x}_i))) + t_i (1-\delta(F_{m-1}(\bar{x}_i))) = \\ &= -\delta(F_{m-1}(\bar{x}_i)) + t_i \cdot \cancel{\delta(F_{m-1}(\bar{x}_i))} + t_i - \cancel{t_i \delta(F_{m-1}(\bar{x}_i))} = \\ &= t_i - \delta(F_{m-1}(\bar{x}_i)) = \tilde{t}_i^{(m)} - \text{label for } f_m(\bar{x}_i) \text{ tree} \end{aligned}$$

$$\boxed{\tilde{t}_i^{(m)} = t_i - \delta(F_{m-1}(\bar{x}_i))}$$

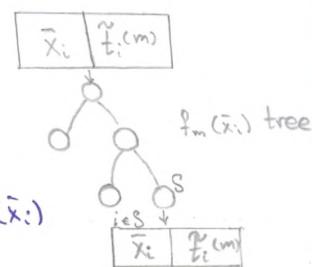
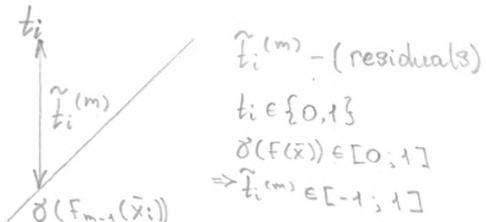
we try to approximate $\tilde{t}_i^{(m)}$ by the underfitting weak regressor tree $f_m(\bar{x}_i)$

but even if it will approximate 100% i.e. $f_m(\bar{x}_i) = \tilde{t}_i^{(m)}, \forall i$

$\boxed{\delta(F_{m-1}(\bar{x}_i) + \tilde{t}_i^{(m)}) \neq t_i}$, it does not guarantee convergence $t_i = \delta(F(\bar{x}_i)), \forall i$

$\tilde{t}_i^{(m)}$ - some направлення. Йона змоза зробити це намагається більше, змоза б зробити
 змоза зробити residuals більше б оголошити terminal value.

$$\boxed{L'(\delta(F_{m-1}(\bar{x}_i) + f_m(\bar{x}_i))) = 0} \quad \text{we can't just return average of } \tilde{t}_i^{(m)}, \text{ too slow convergence}$$



Gradient Boosting: Classification (3)

$$y(\bar{x}) = \delta(F(\bar{x})) = \delta(f_0 + d \sum_{j=1}^M f_j(\bar{x})) , \quad f_j - \text{regression tree}$$

$$\delta(f_0) = \frac{1}{N} \sum_{n=1}^N t_n \longleftrightarrow f_0 = \ln \frac{P_0}{P_1}$$

Given: f_0, f_1, \dots, f_{m-1}

$f_m - ?$

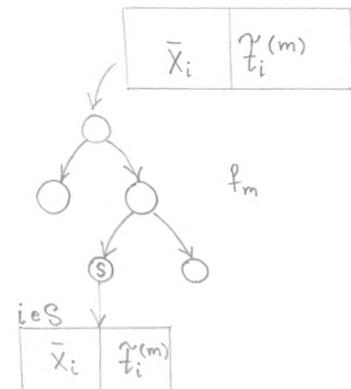
$$\begin{cases} F_m(\bar{x}_i) = F_{m-1}(\bar{x}_i) + d f_m(\bar{x}_i) \\ F_m(\bar{x}_i) = F_{m-1}(\bar{x}_i) - \lambda L'(\delta(F_{m-1}(\bar{x}_i)))_{F_{m-1}} \end{cases}$$

$$f_m(\bar{x}_i) = -L'(\delta(F_{m-1}(\bar{x}_i)))_{F_{m-1}}$$

$$-L'(\delta(F_{m-1}(\bar{x}_i)))_{F_{m-1}} = t_i - \delta(F_{m-1}(\bar{x}_i)) = \tilde{t}_i^{(m)}$$

$\tilde{t}_i^{(m)}$ - direction for minimising L

$\delta(F_{m-1}(\bar{x}_i) + \tilde{t}_i^{(m)}) \neq t_i$, so we can't just return average on terminal node of f_m tree.



residuals, we try to minimize Variance at the nodes of the f_m regression tree at these values

$$L(\delta(F_{m-1}(\bar{x}) + f_m(\bar{x}))) = -\sum_{n=1}^N ((1-t_n) \ln(1-\delta(F_{m-1}(\bar{x}_n) + f_m(\bar{x}_n))) + t_n \ln \delta(F_{m-1}(\bar{x}_n) + f_m(\bar{x}_n))) \rightarrow \min_{f_m}$$

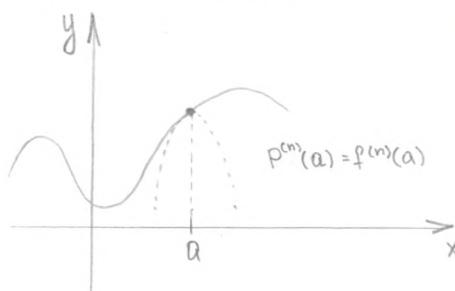
we can't take the derivative of f_m , because we would need f_m in δ but we don't know it.

Taylor Series

$$f(x) = P(x) = \sum_{n=0}^N \underbrace{\frac{f^{(n)}(a)}{n!}}_{\text{const}} (x-a)^n$$

$$P(x+a) = \sum_{n=0}^N \underbrace{\frac{P^{(n)}(a)}{n!}}_{\text{var}} ((x+a)-a)^n$$

$$f = L, \quad a = F_{m-1}, \quad x = f_m$$



$$L(\delta(F_{m-1}(\bar{x}_i) + f_m(\bar{x}_i))) \approx \sum_{n=0}^2 \frac{L^{(n)}(F_{m-1})}{n!} f_m^n = \underbrace{L(\delta(F_{m-1}(\bar{x}_i)))}_{C_0} + \underbrace{L^{(1)}(\delta(F_{m-1}(\bar{x}_i))) f_m}_{C_1} + \underbrace{\frac{1}{2} L^{(2)}(\delta(F_{m-1}(\bar{x}_i))) f_m^2}_{C_2}$$

$$L'_{f_m}(\delta(F_{m-1}(\bar{x}_i) + f_m(\bar{x}_i))) \approx L^{(1)}(\delta(F_{m-1}(\bar{x}_i))) + L^{(2)}(\delta(F_{m-1}(\bar{x}_i))) \cdot f_m = 0$$

on terminal node

$$f_m = \frac{-L^{(1)}(\delta(F_{m-1}(\bar{x}_i)))_{F_{m-1}}}{L^{(2)}(\delta(F_{m-1}(\bar{x}_i)))_{F_{m-1}}}$$

- Looks like Newton's Optimization Method and it is

$L' = 0$ - minima, we try to find roots of L'

but it's not guaranteed to be global minima, because L is not convex function

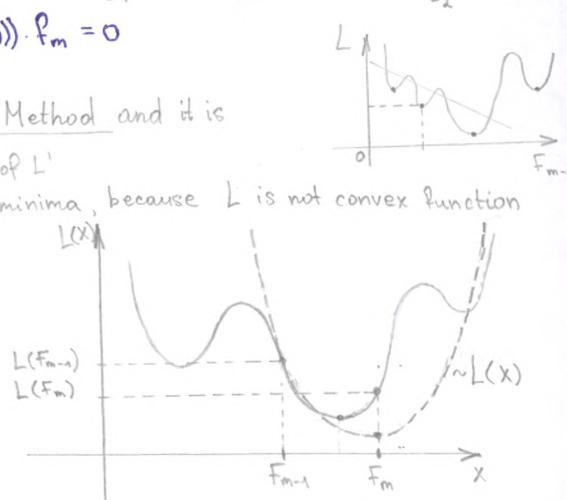
$$F_m(\bar{x}_i) = F_{m-1}(\bar{x}_i) + f_m(\bar{x}_i)$$

$\sim L'(\bar{x}_i) = 0$ local minimum with approximation error

$\sim L(\bar{x}_i) > 0$ but minima

$$L'(\bar{x}_i) = \emptyset$$

$$L(\bar{x}_i) > 0 \quad \text{close to minima}$$



Gradient Boosting: Classification (4)

$$y(\bar{x}) = \delta(F(\bar{x})) = \delta(f_0 + d \sum_{j=1}^M f_j(\bar{x}))$$

$$\delta(f_0) = \frac{1}{N} \sum_{n=1}^N t_n \longleftrightarrow f_0 = \ln \frac{p_0}{p_1}$$

$$\tilde{t}_i^{(m)} = t_i - \delta(F_{m-1}(\bar{x}_i))$$

$$f_m(\bar{x}) = -\frac{L'(\delta(F_{m-1}(\bar{x}_i)))}{L''(\delta(F_{m-1}(\bar{x}_i)))}, \forall i \in S$$

$$L(\delta(F_{m-1}(\bar{x}))) = - \sum_{n=1}^N ((1-t_n) \ln(1-\delta(F_{m-1}(\bar{x}_n))) + t_n \ln \delta(F_{m-1}(\bar{x}_n)))$$

$$\begin{aligned} L'(F_{m-1})(\delta(F_{m-1}(\bar{x}))) &= - \sum_{\substack{n=1 \\ n \in S}}^N ((1-t_n) \frac{1}{1-\delta(F_{m-1}(\bar{x}_n))} \cdot (-\delta(F_{m-1}(\bar{x}_n))(1-\delta(F_{m-1}(\bar{x}_n))) + \\ &\quad + t_n \frac{1}{\delta(F_{m-1}(\bar{x}_n))} \cdot \delta(F_{m-1}(\bar{x}_n))(1-\delta(F_{m-1}(\bar{x}_n)))) = \\ &= - \sum_{n=1}^N ((1-t_n)(-\delta(F_{m-1}(\bar{x}_n))) + t_n(1-\delta(F_{m-1}(\bar{x}_n)))) = \\ &= - \sum_{n=1}^N (-\delta(F_{m-1}(\bar{x}_n)) + t_n \delta(F_{m-1}(\bar{x}_n)) + t_n - t_n \delta(F_{m-1}(\bar{x}_n))) = \\ &= \sum_{n=1}^N (\delta(F_{m-1}(\bar{x}_n)) - t_n) \end{aligned}$$

$$L''(F_{m-1})(\delta(F_{m-1}(\bar{x}))) = \sum_{\substack{n=1 \\ n \in S}}^N \delta(F_{m-1}(\bar{x}_n))(1-\delta(F_{m-1}(\bar{x}_n)))$$

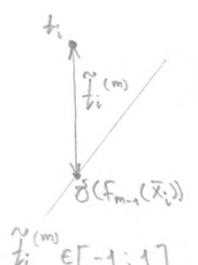
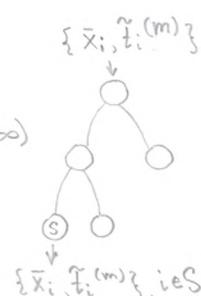
$$f_m(\bar{x}) = \frac{\sum_{n=1}^N (\tilde{t}_n - \delta(F_{m-1}(\bar{x}_n)))}{\sum_{\substack{n=1 \\ n \in S}}^N \delta(F_{m-1}(\bar{x}_n))(1-\delta(F_{m-1}(\bar{x}_n)))} \in [0; \infty) \quad \in (0; 0.25]$$

$$\sim L'(\delta(F_{m-1}(\bar{x}_i)) + f_m(\bar{x}_i)) = 0$$

$\sim L(\delta(F_{m-1}(\bar{x}_i)) + f_m(\bar{x}_i)) > 0$ but minima of $\sim L$

$$L'(\delta(F_{m-1}(\bar{x}_i)) + f_m(\bar{x}_i)) = \emptyset$$

$L(\delta(F_{m-1}(\bar{x}_i)) + f_m(\bar{x}_i)) > 0$ somewhere close to minima according to approximated $\sim L$



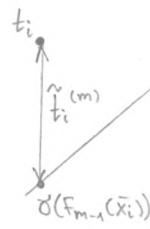
Gradient Boosting: Classification . Summary.

$$y(\bar{x}) = \delta(F(\bar{x})) = \delta(f_0 + d \sum_{j=1}^N f_j(\bar{x}))$$

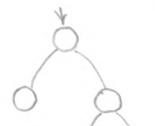
$$f_0 = \frac{p_0}{p_1} \iff \delta(f_0) = \frac{1}{N} \sum_{n=1}^N t_n$$

$$\tilde{t}_i^{(m)} = t_i - \delta(F_{m-1}(\bar{x}_i)) \in [-1, 1]$$

$$f_m(\bar{x}) = \frac{\sum_{i \in S} (t_i - \delta(F_{m-1}(\bar{x}_i)))}{\sum_{i \in S} \delta(F_{m-1}(\bar{x}_i))(1 - \delta(F_{m-1}(\bar{x}_i)))} \in [0; \infty)$$



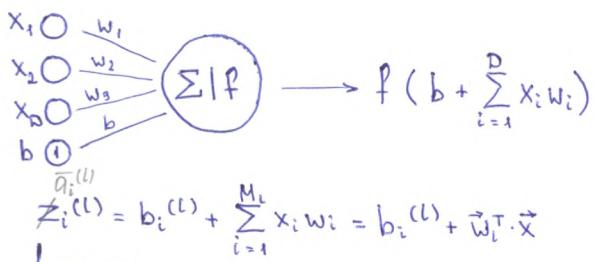
$$\{\bar{x}_i, \tilde{t}_i^{(m)}\}$$



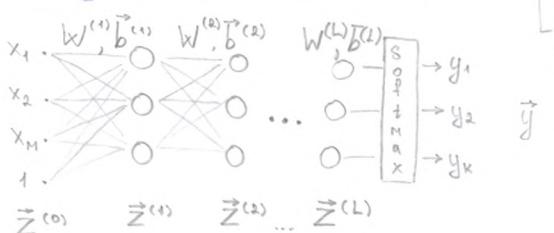
$$\{\bar{x}_i, \tilde{t}_i^{(m)}\}, i \in S$$

Multilayer Perceptron (MLP)

Neuron



Layer



M_L - number of neurons

$$\vec{z}^{(l)} = \vec{f}(\vec{a}^l) = \vec{f}(\vec{a}^{(l)}) = \vec{f}(W^{(l)} \cdot \vec{z}^{(l-1)} + \vec{b}^{(l)}) = \\ = \vec{f}(\vec{b}^{(l)} + W^{(l)} \cdot \vec{f}(\vec{b}^{(l-1)} + W^{(l-1)} \cdot \vec{f}(\dots, \vec{x})))$$

$$Y(\vec{x}, \mathbb{W}) = \text{Softmax}(\vec{z}^{(L)}(\vec{x}, \mathbb{W}))$$

$$z_i > z_j \rightarrow y_i > y_j$$

Loss Function

$$L(W) = - \sum_{i=1}^N \sum_{k=1}^K t_i^{(k)} \ln Y_k(\vec{x}_i, W), \quad t_i^{(k)} = \begin{cases} 1, & \text{if } k=t \\ 0, & \text{otherwise} \end{cases}$$

$$L(\mathbf{W}, \mathbf{X}_{tr}, \mathbf{T}_{tr}) \rightarrow \min_{\mathbf{W}}$$

1. Weights initialization

$$w_i \sim N(\mu=0, \sigma^2 = \frac{2}{n_i})$$

$$2. \quad w_n = w_{n-1} - \gamma \nabla L(w_{n-1})$$

8. Stop criteria

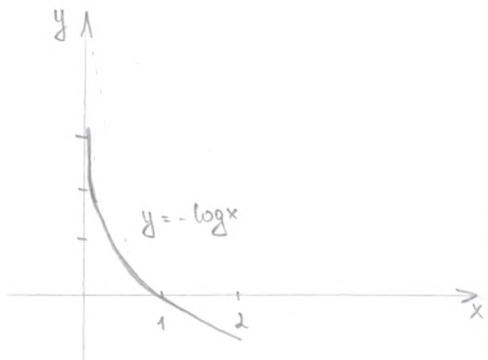
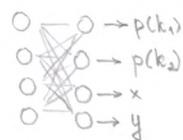
f - activation function :

$$\text{Sigmoid : } f(a) = \frac{1}{1 + e^{-a}}$$

$$\cdot \text{Tanh} : f(a) = \frac{e^{2a} - 1}{e^{2a} + 1}$$

$$\cdot \text{ReLU} : f(a) = \max(0; a)$$

• Leaky ReLU : $f(a) = \max(0.1a; a)$



MLP Backprop (Jacobian)

$$E(\bar{y}(\bar{x})) \quad \nabla_E E = 1$$

$$\bar{y}(\bar{x}), \quad \nabla_y E = \left[\frac{dE}{dy_1}, \frac{dE}{dy_2}, \dots, \frac{dE}{dy_M} \right]^T$$

$$\nabla_x E - ?$$

$$\bar{y}(\bar{x}) = \begin{cases} y_1(\bar{x}) = b_1 + \sum_{j=1}^M w_{j1}x_j \\ y_2(\bar{x}) = b_2 + \sum_{j=1}^M w_{j2}x_j \\ y_3(\bar{x}) = b_3 + \sum_{j=1}^M w_{j3}x_j \end{cases}$$

$\Rightarrow \bar{y} \rightarrow y_1$
 $\Rightarrow \bar{y} \rightarrow y_2$
 $\Rightarrow \bar{y} \rightarrow y_3$

E

$$\nabla_x E = \left[\underbrace{\sum_{i=1}^M \frac{dE}{dy_i} \frac{dy_i}{dx_1}}_{\text{Why sum?}} \quad \sum_{i=1}^M \frac{dE}{dy_i} \frac{dy_i}{dx_2} \dots \right]^T \quad \frac{dE}{dx_1} = \underbrace{\frac{dE}{dy_1} \cdot \frac{dy_1}{dx_1}}_{\frac{dE}{dy_1}} + \underbrace{\frac{dE}{dy_2} \cdot \frac{dy_2}{dx_1}}_{\frac{dE}{dy_2}} + \underbrace{\frac{dE}{dy_3} \cdot \frac{dy_3}{dx_1}}_{\frac{dE}{dy_3}}$$

$$\text{Jacobian} = \frac{d\bar{y}}{dx} = \left[\left[\frac{dy_1}{dx_1}, \frac{dy_1}{dx_2}, \dots \right] \quad \left[\frac{dy_2}{dx_1}, \frac{dy_2}{dx_2}, \dots \right] \quad \left[\frac{dy_3}{dx_1}, \frac{dy_3}{dx_2}, \dots \right] \right]$$

$\frac{df}{dx} = \frac{dy}{dx} \cdot \frac{df}{dy} \rightarrow \frac{df}{dx} = \frac{df}{dy_1} \frac{dy_1}{dx} + \frac{df}{dy_2} \frac{dy_2}{dx}$

$$\nabla_x E = \text{Jacobian}^T \cdot \nabla_y E$$

$x_0 \rightarrow y_1 \rightarrow \text{?}$

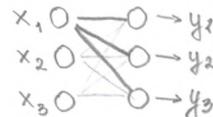
$y_1 \rightarrow y_2 \rightarrow \text{?}$

$y_2 \rightarrow \text{?}$

Матрица
законов
пропорциональности

Why sum?

$$\frac{dE}{dx_1} = \sum_{i=1}^M \underbrace{\frac{dE}{dy_i} \frac{dy_i}{dx_1}}_{\text{Chain rule}}$$



$$E = \sum_{j=1}^M (t^{(j)} - y^{(j)}(x_1, x_2, x_3))^2$$

Сумма виагод $y_i(x_i)$ определяет более изменение.

$$f'(x) = \lim_{\Delta x \rightarrow 0} \frac{\Delta f}{\Delta x} = \lim_{\Delta x \rightarrow 0} \frac{f(x_0 + \Delta x) - f(x_0)}{\Delta x} = \frac{df}{dx}$$

$$f'(x) = \frac{f(x_0 + \Delta x) - f(x_0)}{\Delta x} = \frac{f(x_0 + x - x_0) - f(x_0)}{x - x_0} = \frac{f(x) - f(x_0)}{x - x_0} \Rightarrow f(x) = f(x_0) + f'(x_0)(x - x_0)$$

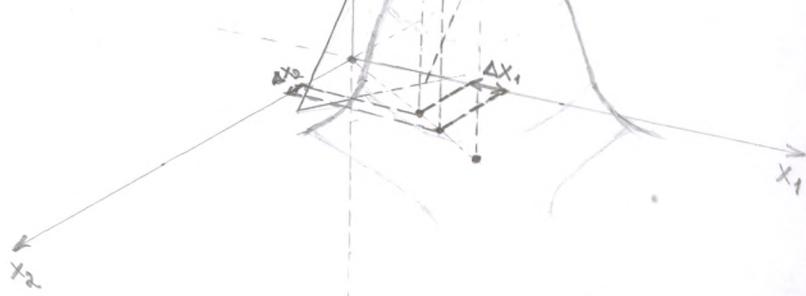
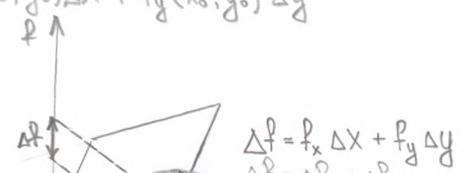
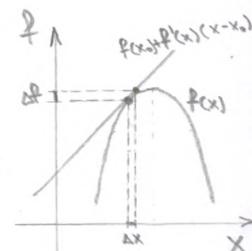
$$f'(x) = \frac{\Delta f}{\Delta x} \Rightarrow \Delta f = f'(x) \Delta x \quad \rightarrow df = f'(x) \cdot dx$$

$$\Delta S = 25 \cdot \Delta t$$

$$\Delta f = [f(x_0, y_0) + f_x(x_0, y_0) \Delta x + f_y(x_0, y_0) \Delta y] - f(x_0, y_0) = f_x(x_0, y_0) \Delta x + f_y(x_0, y_0) \Delta y$$

$$df = \frac{df}{dx} \cdot dx + \frac{df}{dy} \cdot dy$$

$$E(y_1(x_1), y_2(x_1), \dots)$$



Classification MLP Backprop⁽¹⁾ Loss function

$\vec{y} = \text{Softmax}(\vec{z}^{(L)})$, $\vec{t} = [0010\dots0]$ one-hot encoding.

$$P(T|X, W) = \prod_{n=1}^N \prod_{k=1}^K t_n^{(k)} (\vec{x}_n, W)^{z_n^{(k)}} \rightarrow \max$$

$$\ln P(T|X, W) = \sum_{n=1}^N \sum_{k=1}^K t_n^{(k)} \ln y^{(k)}(\vec{x}_n, W) \rightarrow \max$$

$$E = - \sum_{n=1}^N \sum_{k=1}^K t_n^{(k)} \ln y^{(k)}(\vec{x}_n, W) \rightarrow \min$$

Forward pass

$$\vec{y} = \text{Softmax}(\vec{z}^{(L)}) = \text{Softmax}(\vec{f}(\vec{a}^{(L)})) = \text{Softmax}(\vec{f}(\vec{b}^{(L)} + W^{(L)} \vec{z}^{(L-1)}))$$

$$\vec{z}^{(0)} = \vec{x}$$

$$\vec{z}^{(0)} \rightarrow \vec{a}^{(1)} \rightarrow \vec{z}^{(1)} \rightarrow \dots \rightarrow \vec{a}^{(L)} \rightarrow \vec{z}^{(L)} \rightarrow \vec{y} \rightarrow E$$

Backward pass

$$\underbrace{\nabla_E E}_{1} \rightarrow \underbrace{\nabla_y E}_{\ln(y)} \rightarrow \underbrace{\nabla_{z^{(L)}} E}_{\text{Softmax}(z^{(L)})} \rightarrow \underbrace{\nabla_{a^{(L)}} E}_{\nabla_{W^{(L)}} E} \rightarrow \dots \rightarrow \underbrace{\nabla_{a^{(1)}} E}_{\nabla_{W^{(1)}} E} \rightarrow \underbrace{\nabla_{b^{(1)}} E}_{\nabla_{b^{(1)}} E}$$

$$E = - \sum_{k=1}^K t_n^{(k)} \ln y^{(k)}(\vec{x}, W) \rightarrow \min$$

$$\begin{cases} \frac{dE}{dy^{(i)}} = (-\ln y^{(i)}(\vec{x}, W))' y^{(i)} = -\frac{1}{y^{(i)}(\vec{x}, W)}, & i = t_n - \text{ground truth} \\ \frac{dE}{dy^{(i)}} = 0, & i \neq t_n \end{cases}$$

$$\vec{t} = [0 \ 0 \ \dots \underset{t_n}{1} \ \dots \ 0 \ 0 \ 0]^T$$

$$\nabla_y E = [0 \ 0 \ \dots \ -\frac{1}{y^{(t_n)}(\vec{x}, W)} \ \dots \ 0 \ 0]^T$$

$$\vec{y} = \text{Softmax}(\vec{z}^{(L)})$$

$$\frac{dE}{dz_i^{(L)}} = \frac{dE}{dy_i} \cdot \frac{dy_i}{dz_i^{(L)}}$$

$$\nabla_{z^{(L)}} E = [\sum_{i=1}^M \frac{dE}{dy_i} \frac{dy_i}{dz_i^{(L)}}, \sum_{i=1}^M \frac{dE}{dy_i} \frac{dy_i}{dz_2^{(L)}} \dots]^T = [\underbrace{[\frac{dy_i}{dz_1^{(L)}}] \ [\frac{dy_i}{dz_2^{(L)}}]^T}_{\text{Jacobian}} \cdot \nabla_y E]$$

$$\nabla_{z^{(L)}} E = (\frac{dy_i}{dz^{(L)}})^T \cdot \nabla_y E$$

$$\nabla_{z^{(L)}} E = \vec{y} - \vec{t}$$

$$\nabla_z E = \begin{bmatrix} 0,1 \\ 0,5 \\ 0,4 \end{bmatrix} \begin{bmatrix} 0,1 \\ -0,5 \\ 0,4 \end{bmatrix}$$

Несмотря на то что градиент вектора E направлен вправо, градиент вектора $\nabla_z E$ направлен влево.

$$-\nabla_z E = \begin{bmatrix} -0,1 \\ 0,5 \\ -0,4 \end{bmatrix}$$

Уменьшаем коэффициенты ошибки

$$\begin{array}{l} z_1 \rightarrow y_1, \quad t_1 = 0 \\ z_2 \rightarrow y_2, \quad t_2 = 1 \\ z_3 \rightarrow y_3, \quad t_3 = 0 \end{array}$$

Classification MLP Backprop⁽²⁾ (Jacobian)

$$\nabla_y E$$

$$\bar{y} = \text{Softmax}(\bar{z}^{(L)}) = \left(\frac{e^{z_1^{(L)}}}{\sum_{i=1}^M e^{z_i^{(L)}}}, \frac{e^{z_2^{(L)}}}{\sum_{i=1}^M e^{z_i^{(L)}}}, \frac{e^{z_3^{(L)}}}{\sum_{i=1}^M e^{z_i^{(L)}}} \right)$$

For example $z_1^{(L)}$ on $y^{(1)}$

$$i=j \quad \frac{dy^{(j)}}{dz_i^{(L)}} = \left(\frac{e^{z_j^{(L)}}}{\sum_{u=1}^M e^{z_u^{(L)}}} \right)' = \underbrace{(e^{z_j^{(L)}})^1}_{\frac{1}{e^{z_1^{(L)}}}} \underbrace{\left(\frac{1}{\sum_{u=1}^M e^{z_u^{(L)}}} \right)}_{e^{z_1^{(L)}}} + \underbrace{(e^{z_j^{(L)}})}_{\frac{-1}{(\sum_{u=1}^M e^{z_u^{(L)}})^2}} \underbrace{\cdot \partial z_i^{(L)}}_{e^{z_1^{(L)}}} = \\ = y^{(j)} + (-y^{(j)} \cdot y^{(j)}) = \boxed{y^{(j)} - (y^{(j)})^2}$$

$z_1^{(L)}$ on $y^{(2)}$

$$i \neq j \quad \frac{dy^{(j)}}{dz_i^{(L)}} = \left(\frac{e^{z_j^{(L)}}}{\sum_{u=1}^M e^{z_u^{(L)}}} \right)' = \underbrace{(e^{z_j^{(L)}})^1}_{\frac{1}{e^{z_1^{(L)}}}} \underbrace{\left(\frac{1}{\sum_{u=1}^M e^{z_u^{(L)}}} \right)}_{e^{z_1^{(L)}}} + \underbrace{e^{z_j^{(L)}}}_{\frac{-1}{(\sum_{u=1}^M e^{z_u^{(L)}})^2}} \underbrace{\cdot \partial z_i^{(L)}}_{e^{z_1^{(L)}}} = \\ = 0 - y_j^{(j)} \cdot y_i^{(j)} = \boxed{-y_j^{(j)} y_i^{(j)}}$$

$$\left(\frac{u}{v} \right)' = (u \cdot v^{-1})' = u' v^{-1} + u \cdot (-1 v^{-2}) \\ (f(g(x)))' = f'(g) \cdot g'(x)$$

$$\text{Jacobian} = \frac{d\bar{y}}{dz^{(L)}} = \left[\left[\frac{d\bar{y}_1}{dz_1^{(L)}} \right] \left[\frac{d\bar{y}_2}{dz_2^{(L)}} \right] \dots \right] = \left[\begin{array}{c|c|c|c} \frac{d\bar{y}_1}{dz_1^{(L)}} & \frac{d\bar{y}_2}{dz_2^{(L)}} & \dots & \\ \frac{d\bar{y}_1}{dz_2^{(L)}} & \frac{d\bar{y}_2}{dz_3^{(L)}} & & \\ \dots & & & \\ \frac{d\bar{y}_1}{dz_M^{(L)}} & \frac{d\bar{y}_2}{dz_M^{(L)}} & & \end{array} \right]$$

$$\frac{d\bar{y}}{dz^{(L)}} = \left[\begin{array}{cccc} y^{(1)} - (y^{(1)})^2 & -y^{(1)} y^{(2)} & -y^{(1)} y^{(3)} & \dots \\ -y^{(2)} y^{(1)} & y^{(2)} - (y^{(2)})^2 & -y^{(2)} y^{(3)} & \dots \\ -y^{(3)} y^{(1)} & -y^{(3)} y^{(2)} & y^{(3)} - (y^{(3)})^2 & \dots \\ \dots & \dots & \dots & \dots \\ -y^{(M)} y^{(1)} & -y^{(M)} y^{(2)} & -y^{(M)} y^{(3)} & \dots \end{array} \right] \frac{d\bar{y}_1}{dz^{(L)}}$$

$$\nabla_{z^{(L)}} E = \left(\frac{d\bar{y}}{dz^{(L)}} \right)^T \nabla_y E \quad , \quad \nabla_y E = [0 \ 0 \ \dots \ -\frac{1}{y^{(t_k)}} \ \dots \ 0 \ 0]^T$$

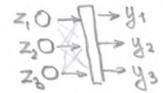
$$\nabla_{z^{(L)}} E = \left[\begin{array}{c} 0 + (-y^{(1)} y^{(t_k)}) \cdot (-\frac{1}{y^{(t_k)}}) = y^{(1)} \\ 0 + (-y^{(2)} y^{(t_k)}) \cdot (-\frac{1}{y^{(t_k)}}) = y^{(2)} \\ \dots \\ 0 + (-y^{(t_k)} y^{(t_k)}) \cdot (-\frac{1}{y^{(t_k)}}) = -1 + y^{(t_k)} \cdot y^{(t_k)-1} \\ y^{(M)} \end{array} \right] = \bar{y} - \bar{t}$$

t_k -th row

$$1) -y^{(1)} y^{(t_k)}, \text{ norga fórumaem } z_i^{(L)}, i \neq t_k$$

$$2) y^{(t_k)} - y^{(t_k)} y^{(t_k)}, \text{ norga fórumaem } z_i^{(L)}, i = t_k$$

Other rows $\times 0$



Regression MLP Backprop

$$t(\bar{x}) = y(\bar{x}, \bar{w}) + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2)$$

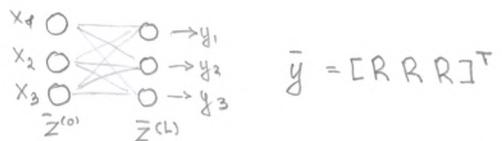
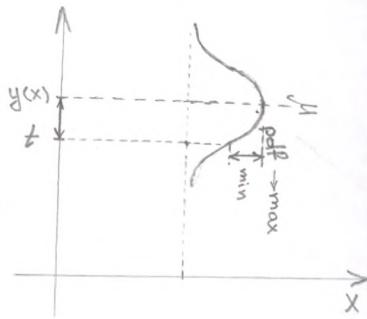
$$P(T|X, W) = \prod_{n=1}^N P(t_n | y(\bar{x}_n, W), \sigma^2) = \prod_{n=1}^N \frac{1}{\sigma \sqrt{2\pi}} \cdot e^{-\frac{1}{2} \left(\frac{t_n - y(\bar{x}_n, W)}{\sigma} \right)^2} \rightarrow \max$$

$$\begin{aligned} \ln P(T|X, W) &= \sum_{n=1}^N \ln \left(\frac{1}{\sigma \sqrt{2\pi}} \right) + \sum_{n=1}^N \left(-\frac{1}{2} \left(\frac{t_n - y(\bar{x}_n, W)}{\sigma} \right)^2 \cdot \ln e \right) = \\ &= \underbrace{N \ln \left(\frac{1}{\sigma \sqrt{2\pi}} \right)}_{\text{const}} - \underbrace{\frac{1}{\sigma^2} \cdot \frac{1}{2} \sum_{n=1}^N (t_n - y(\bar{x}_n, W))^2}_{\min} \rightarrow \max \end{aligned}$$

$$E = \frac{1}{2} \sum_{n=1}^N (t_n - y(\bar{x}_n, W))^2 \rightarrow \min$$

for multiple outputs: $\bar{y} \in \mathbb{R}^M$

$$E = \frac{1}{2} \sum_{n=1}^N \sum_{j=1}^M (t_n^{(j)} - y^{(j)}(\bar{x}_n, W))^2 \rightarrow \min$$



Linear activation function: $f(a) = a$

Forward pass

$$\bar{x} = \bar{z}^{(0)}$$

$$\begin{aligned} \bar{z}^{(0)} &\rightarrow \bar{a}^{(1)} \rightarrow \bar{z}^{(1)} \rightarrow \bar{a}^{(2)} \rightarrow \bar{z}^{(2)} \rightarrow \dots \rightarrow \bar{z}^{(L)} \rightarrow \bar{y} \rightarrow E \\ \bar{y} &= \bar{z}^{(L)} = f(\bar{a}^{(L)}) = f(\bar{b}^{(L)} + W^{(L)} \cdot \bar{z}^{(L-1)}) \end{aligned}$$

Backward pass

$$\begin{array}{ccccccccc} \nabla_E E & \rightarrow & \nabla_{\bar{y}} E & \rightarrow & \nabla_{\bar{z}^{(L)}} E & \rightarrow & \nabla_{\bar{a}^{(L)}} E & \rightarrow & \nabla_{\bar{z}^{(L-1)}} E \\ \downarrow & & \downarrow & & \downarrow & & \downarrow & & \downarrow \\ \bar{y} = \bar{z}^{(L)} & & & & & & & & \end{array}$$

$$\begin{array}{c} \nabla_{W^{(L)}} E \\ \nabla_{b^{(L)}} E \end{array}$$

$$\begin{array}{c} \nabla_{W^{(L-1)}} E \\ \nabla_{b^{(L-1)}} E \end{array}$$

$$\begin{array}{c} \nabla_{W^{(1)}} E \\ \nabla_{b^{(1)}} E \end{array}$$

$$E = \frac{1}{2} \sum_{n=1}^N \sum_{j=1}^M (t_n^{(j)} - z_j^{(L)})^2 \rightarrow \min$$

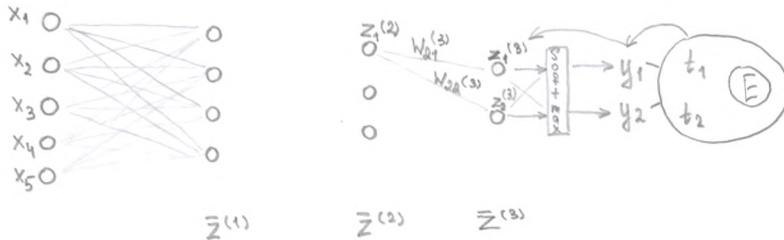
$$\frac{dE}{dz_k^{(L)}} = \frac{1}{2} \sum_{n=1}^N \sum_{j=1}^M 2 \cdot (t_n^{(j)} - z_j^{(L)}) \cdot (t_n^{(j)} - z_j^{(L)})' \Big|_{z_k^{(L)}} = \sum_{n=1}^N \sum_{j=1}^M (z_j^{(L)} - t_n^{(j)}) \mathbf{I}(j=k) = \sum_{n=1}^N (z_k^{(L)} - t_n^{(k)})$$

$\begin{cases} -1, & \text{if } j=k \\ 0, & \text{otherwise} \end{cases}$

$$\nabla_{\bar{z}^{(L)}} E = \bar{z}^{(L)} - \bar{t}$$

$\nearrow \circ \rightarrow z_1^{(L)}$	$- t^{(1)}$
$\nearrow \circ \rightarrow z_2^{(L)}$	$- t^{(2)}$
$\nearrow \circ \rightarrow z_3^{(L)}$	$- t^{(3)}$

Backprop (1)



$$\bar{z}^{(l)} = \bar{f}(\bar{a}^{(l)}) = \bar{f}(b^{(l)} + W^{(l)} \bar{z}^{(l-1)})$$

$$\frac{dE}{da} = \frac{dE}{dy} \frac{dy}{dz} \frac{dz}{df} \frac{df}{da} = \frac{dE}{dz} \frac{df}{da} = \frac{dE}{dz} \frac{df}{da}$$

$$\nabla_{z^{(l)}} E, \quad \bar{z}^{(l)} = \bar{f}(\bar{a}^{(l)})$$

$$z_i^{(l)} = \text{ReLU}(a_i^{(l)})$$

$$\frac{dz}{da} = \frac{dz}{df_1} \frac{df_1}{da} + \frac{dz}{df_2} \frac{df_2}{da} \dots \Rightarrow \nabla_{a^{(l)}} \bar{E} = \left[\sum_{i=1}^M \frac{dz}{df_i} \frac{df_i}{da_1} \quad \sum_{i=1}^M \frac{dz}{df_i} \frac{df_i}{da_2} \dots \right]^T, \text{ but } f_2(a_2) \neq f_2(a_1) \\ \Rightarrow f'_2(a_1) = 0$$

$$\nabla_{a^{(l)}} E = \left(\frac{d\bar{f}}{da} \right)^T \cdot \nabla_{z^{(l)}} E$$

$$\nabla_{a^{(l)}} E = \left[\frac{dz}{da_1} \frac{df_1}{da_1} \quad \frac{dz}{da_2} \frac{df_2}{da_2} \dots \right]^T$$

$$\frac{d\bar{f}}{d\bar{a}} = \begin{bmatrix} \frac{df_1}{da_1} & \frac{df_1}{da_2} & \dots & \frac{df_1}{da_M} \\ \frac{df_2}{da_1} & \frac{df_2}{da_2} & \dots & \frac{df_2}{da_M} \\ \dots & \dots & \dots & \frac{df_M}{da_M} \end{bmatrix} = \begin{bmatrix} \frac{df_1}{da_1} & 0 & \dots & 0 \\ 0 & \frac{df_2}{da_2} & \dots & 0 \\ \dots & 0 & \dots & 0 \\ \dots & \dots & \dots & \frac{df_M}{da_M} \end{bmatrix}$$

$$\left(\frac{d\bar{f}}{d\bar{a}} \right)^T \cdot \nabla_{z^{(l)}} E = \left[\frac{df_1}{da_1^{(l)}} \frac{dE}{dz_1^{(l)}} \quad \frac{df_2}{da_2^{(l)}} \frac{dE}{dz_2^{(l)}} \dots \right]^T$$

$$\nabla_{a^{(l)}} E = \underbrace{\bar{f}'(\bar{a}^{(l)})}_{\text{elementwise}} \circ \nabla_{z^{(l)}} E = \left(\frac{d\bar{f}}{d\bar{a}} \right)^T \nabla_{z^{(l)}} E = \nabla_{a^{(l)}} \bar{E} \quad - \text{activation functions}$$

$$\frac{dE}{db_i^{(l)}} = \frac{dE}{dz_i^{(l)}} \cdot \frac{dz_i^{(l)}}{da_i^{(l)}} \cdot \frac{da_i^{(l)}}{db_i^{(l)}} = \frac{dE}{da_i^{(l)}} \cdot \frac{da_i^{(l)}}{db_i^{(l)}} = \frac{dE}{da_i^{(l)}} \quad a(b) = b + W z^{(l-1)} \\ a'(b) = (b)'_b = 1$$

$$\nabla_{b^{(l)}} E = \left[\frac{dE}{da_1^{(l)}} \quad \frac{dE}{da_2^{(l)}} \quad \frac{dE}{da_3^{(l)}} \dots \right]^T = \nabla_{a^{(l)}} E$$

$$\nabla_{z^{(l)}} E = \nabla_{a^{(l)}} E$$

$\nabla_{z^{(l)}} E$ - loss function
 $\nabla_{a^{(l)}} E$ - activation function

$$\bar{z} = \bar{f}(\bar{a}) = \bar{f}(b + W z^{(l-1)})$$

$$\frac{dE}{da} = \frac{dE}{dz} \cdot \underbrace{\frac{dz}{da}}_{f'(a)} \rightarrow \frac{dE}{db} = \frac{dE}{dz} \frac{dz}{da} \cdot \frac{da}{db} = \boxed{\frac{dE}{dz} \cdot f'(a) \cdot a'(b)}$$

Backprop (2)

$$\bar{z}^{(l)} = \bar{f}(\bar{a}^{(l)}) = \bar{f}(\bar{b}^{(l)} + W^{(l)} \bar{z}^{(l+1)})$$

$$W^{(l)} = \begin{bmatrix} w_{11} & w_{21} & w_{31} & w_{41} \\ w_{12} & w_{22} & w_{32} & w_{42} \\ \dots & \dots & \dots & \dots \\ w_{1M} & w_{2M} & \dots & w_{4M} \end{bmatrix}$$



$$\nabla_{z^{(l)}} E, \nabla_{a^{(l)}} E, \nabla_{b^{(l)}} E$$

$$\nabla_{w^{(l)}} E - ?$$



$$\bar{z}^{(l-1)} = \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_N \end{bmatrix}$$

$$a_i^{(l)} = b_i^{(l)} + \sum_{j=1}^{M-1} w_{ij} \cdot z_j^{(l-1)}$$

$$a_j(w_{ij}) = b_j + \sum_{m=1}^{M-1} w_{mj} \cdot z_m^{(l-1)}, a^t(w_{ij}) = 0 + (0 + (w_{ij} \cdot z_i^{(l-1)})^t) = z_i^{(l-1)}$$

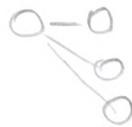
$$\frac{dE}{dw_{ij}} = \frac{dE}{da_j^{(l)}} \cdot \frac{da_j^{(l)}}{d w_{ij}} = \frac{dE}{da_j^{(l)}} \cdot z_i^{(l-1)}$$

$$\nabla_{w^{(l)}} E = \begin{bmatrix} \frac{dE}{dw_{11}} & \frac{dE}{dw_{21}} & \dots & \frac{dE}{dw_{M1}} \\ \frac{dE}{dw_{12}} & \frac{dE}{dw_{22}} & \dots & \frac{dE}{dw_{M2}} \\ \dots & \dots & \dots & \dots \\ \frac{dE}{dw_{1M}} & \frac{dE}{dw_{2M}} & \dots & \frac{dE}{dw_{MM}} \end{bmatrix} = \begin{bmatrix} \frac{dE}{da_1^{(l)}} \cdot \frac{da_1^{(l)}}{d w_{11}} & \frac{dE}{da_1^{(l)}} \cdot \frac{da_1^{(l)}}{d w_{21}} & \dots & \frac{dE}{da_1^{(l)}} \cdot \frac{da_1^{(l)}}{d w_{M1}} \\ \frac{dE}{da_1^{(l)}} \cdot \frac{da_2^{(l)}}{d w_{12}} & \frac{dE}{da_2^{(l)}} \cdot \frac{da_2^{(l)}}{d w_{22}} & \dots & \frac{dE}{da_2^{(l)}} \cdot \frac{da_2^{(l)}}{d w_{M2}} \\ \dots & \dots & \dots & \dots \\ \frac{dE}{da_1^{(l)}} \cdot \frac{da_M^{(l)}}{d w_{1M}} & \frac{dE}{da_M^{(l)}} \cdot \frac{da_M^{(l)}}{d w_{2M}} & \dots & \frac{dE}{da_M^{(l)}} \cdot \frac{da_M^{(l)}}{d w_{MM}} \end{bmatrix} = \begin{bmatrix} \frac{dE}{da_1^{(l)}} \cdot z_1^{(l-1)} & \frac{dE}{da_1^{(l)}} \cdot z_2^{(l-1)} & \dots & \frac{dE}{da_1^{(l)}} \cdot z_N^{(l-1)} \\ \frac{dE}{da_2^{(l)}} \cdot z_1^{(l-1)} & \frac{dE}{da_2^{(l)}} \cdot z_2^{(l-1)} & \dots & \frac{dE}{da_2^{(l)}} \cdot z_N^{(l-1)} \\ \dots & \dots & \dots & \dots \\ \frac{dE}{da_M^{(l)}} \cdot z_1^{(l-1)} & \frac{dE}{da_M^{(l)}} \cdot z_2^{(l-1)} & \dots & \frac{dE}{da_M^{(l)}} \cdot z_N^{(l-1)} \end{bmatrix}$$

$$\nabla_{w^{(l)}} E = \nabla_{a^{(l)}} E \cdot (\bar{z}^{(l)})^T = \begin{bmatrix} \frac{dE}{da_1^{(l)}} \\ \frac{dE}{da_2^{(l)}} \\ \vdots \\ \frac{dE}{da_M^{(l)}} \end{bmatrix} [z_1^{(l)} \ z_2^{(l)} \ \dots \ z_N^{(l)}]$$

$$\nabla_{z^{(l-1)}} E - ?$$

$$\begin{aligned} \frac{dE}{dz_1^{(l-1)}} &= \frac{dE}{dz_1^{(l)}} \frac{da_1^{(l)}}{dz_1^{(l)}} + \frac{dE}{dz_2^{(l)}} \frac{da_2^{(l)}}{dz_2^{(l)}} + \dots \\ &= w_{11} + w_{12} + w_{13} \dots \end{aligned}$$



$$\nabla_{z^{(l-1)}} E = W_l^T \cdot \nabla_{a^{(l)}} E$$

Forward

$$\bar{z}^{(0)} \rightarrow \bar{a}^{(1)} \rightarrow \bar{z}^{(1)} \rightarrow \dots \rightarrow \bar{z}^{(l-1)} \rightarrow \bar{a}^{(l)} \rightarrow \bar{z}^{(l)} \rightarrow y \rightarrow E$$

Backward

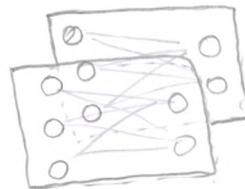
$$\nabla_E E \rightarrow \nabla_y E \rightarrow \nabla_{\bar{z}^{(l)}} E \rightarrow \nabla_{a^{(l)}} E \rightarrow \nabla_{\bar{a}^{(l)}} E \rightarrow \dots \rightarrow \nabla_{\bar{z}^{(1)}} E \rightarrow \nabla_{a^{(1)}} E$$

$$\begin{array}{c} \searrow \\ \nabla_{w^{(l)}} E \\ \searrow \\ \nabla_{b^{(l)}} E \end{array}$$

$$\begin{array}{c} \nearrow \\ \nabla_{w^{(1)}} E \\ \nearrow \\ \nabla_{b^{(1)}} E \end{array}$$

$$W \cdot \vec{v}$$

$$\begin{bmatrix} \textcircled{1} \\ \textcircled{2} \\ \vdots \end{bmatrix} \begin{bmatrix} \textcircled{1} \\ \textcircled{2} \\ \vdots \end{bmatrix} = \begin{bmatrix} \textcircled{1} \\ \textcircled{2} \\ \vdots \end{bmatrix}$$



$$(W \cdot \vec{v})_V^I = W^T \cdot \vec{v}_V \quad \begin{bmatrix} W^T \\ \nabla \end{bmatrix} \begin{bmatrix} \textcircled{1} \\ \textcircled{2} \\ \vdots \end{bmatrix}$$

$$W \cdot V = \underbrace{\begin{bmatrix} \textcircled{1} \\ \textcircled{2} \\ \vdots \end{bmatrix}}_k^{m \times k} \quad m \times k$$

$$\begin{bmatrix} W & 2 \times 3 \\ \nabla & 3 \times 3 \end{bmatrix} \begin{bmatrix} \textcircled{1} \\ \textcircled{2} \\ \vdots \end{bmatrix} = \underbrace{\begin{bmatrix} \textcircled{1} \\ \textcircled{2} \\ \vdots \end{bmatrix}}_{2 \times 3}$$

$$(W \cdot V)_V^I = W^T \nabla, \text{ also } T^{2 \times 3}$$

$$(W \times)_W^I = \nabla \cdot \vec{x}^T$$

$$\begin{bmatrix} \textcircled{1} \\ \textcircled{2} \\ \vdots \end{bmatrix} \begin{bmatrix} \textcircled{1} \\ \textcircled{2} \\ \vdots \end{bmatrix} = \begin{bmatrix} \textcircled{1} \\ \textcircled{2} \\ \vdots \end{bmatrix}$$

$$\begin{bmatrix} \textcircled{1} \\ \textcircled{2} \\ \vdots \end{bmatrix} \begin{bmatrix} \textcircled{1} \\ \textcircled{2} \\ \vdots \end{bmatrix} = \begin{bmatrix} \textcircled{1} \\ \textcircled{2} \\ \vdots \end{bmatrix}$$

$$(W \times)_W^I = \nabla \cdot \vec{x}^T$$

$$\begin{bmatrix} W \\ \nabla \end{bmatrix} \begin{bmatrix} \textcircled{1} \\ \textcircled{2} \\ \vdots \end{bmatrix} \begin{bmatrix} \textcircled{1} \\ \textcircled{2} \\ \vdots \end{bmatrix} = \underbrace{\begin{bmatrix} \textcircled{1} \\ \textcircled{2} \\ \vdots \end{bmatrix}}_{2 \times 2}$$

$$\begin{bmatrix} \textcircled{1} \\ \textcircled{2} \\ \vdots \end{bmatrix} \begin{bmatrix} \textcircled{1} \\ \textcircled{2} \\ \vdots \end{bmatrix} = \underbrace{\begin{bmatrix} \textcircled{1} \\ \textcircled{2} \\ \vdots \end{bmatrix}}_{2 \times 2}$$

$$W \cdot V^I = W^T \cdot \nabla$$

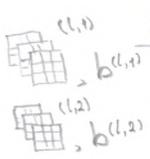
$$W^I \cdot V = \nabla \cdot V^T$$

CNN



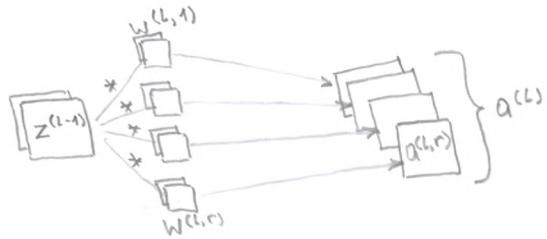
Input Tensor

$$Z^{(L-1)} \in \mathbb{R}^{m_{L-1} \times n_{L-1} \times r_{L-1}}$$



Convolution Filters

$$w^{(l,r)} \in \mathbb{R}^{\tilde{m}_l \times \tilde{n}_l \times r_l}, r=1, 2, \dots, r_l$$



Activation maps

$$a^{(l,r)} \in \mathbb{R}^{m_l \times n_l \times 1}, r=1, 2, \dots, r_l$$

$$a^{(l)} \in \mathbb{R}^{m_l \times n_l \times r_l}$$

$$Z^{(L-1)} * w^{(l)} = a^{(l)}$$

Feature maps ($Z^{(l,r)} = f(a^{(l,r)})$)

$$Z^{(l,r)} \in \mathbb{R}^{m_l \times n_l}$$

$$Z^{(l)} \in \mathbb{R}^{m_l \times n_l \times r_l} \text{ - output tensor}$$

Stride - waż

dilation - skryptowane rozprzestrzenianie



Convolution actually Cross-Correlation (no flip)

$$\text{Image } (I)^{H \times W}$$

Kernel $(K)^{M \times N}$ - convolution filter

$I[i,j]$ - element

$K[m,n]$ - element

$$(I * K)_{[i,j]} = \sum_{m=0}^{M-1} \sum_{n=0}^{N-1} I_{[i+m, j+n]} \cdot K_{[m,n]} + b$$

elementwise

Generally

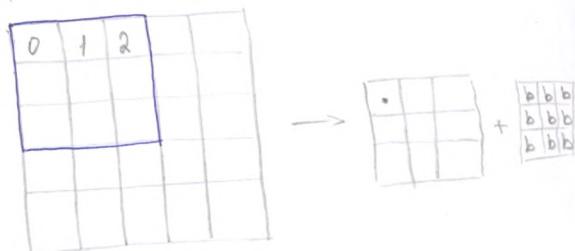
$$\square a^{(l,r)} = Z^{(L-1)} * w^{(l,r)} + b^{(l,r)}, r=1, 2, \dots, r_l$$

$$\square a^{(l)} = Z^{(L-1)} * w^{(l)} + b^{(l)}$$

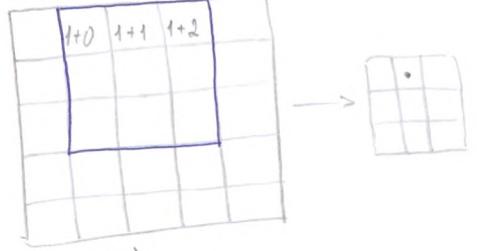
$$Z^{(l)} = f(a^{(l)}), f - \text{ReLU}$$

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \xrightarrow{f} \begin{bmatrix} 0 & 2 \\ 3 & 0 \end{bmatrix}$$

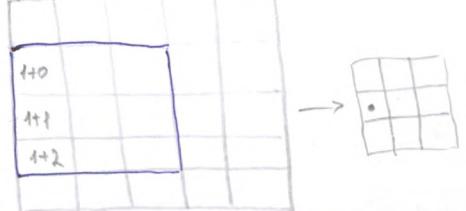
$$(I * K)_{[0,0]}$$



$$(I * K)_{[0,1]}$$

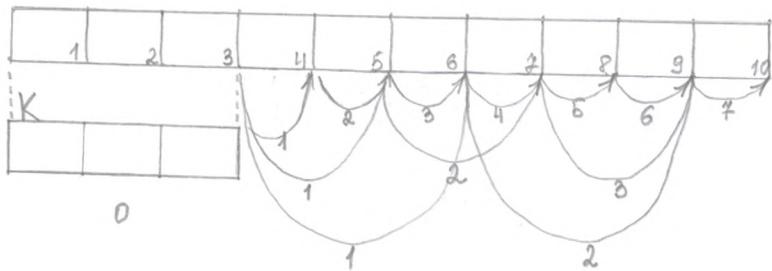


$$(I * K)_{[1,0]}$$



CNN. Output shape

W



$$W = 10, K = 3$$

dilation = 1

- stride = 1 : $W_{out} = 8 = (10-3)+1$
- stride = 2 : $W_{out} = 4 = \lfloor \frac{(10-3)}{2} + 1 \rfloor = \lfloor 3.5 + 1 \rfloor$
- stride = 3 : $W_{out} = 3 = \lfloor \frac{4}{3} + 1 \rfloor = \lfloor 2.3 + 1 \rfloor$

beginning shape as kernel

$$W_{out} = \left\lfloor \frac{(W+2P) - (d(K-1)+1)}{S} + 1 \right\rfloor$$

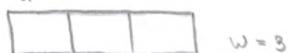
S = stride

2P - padding

d - dilation, default = 1

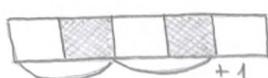
Dilation, как расширение ядра

K=3: d=1



$$W = 3$$

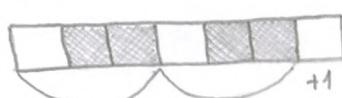
d=2



$$W = 5$$

$$2(3-1)+1 = 4+1$$

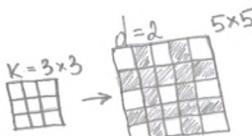
d=3



$$W = 7$$

$$3(3-1)+1 = 6+1$$

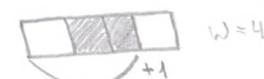
$$d(K-1)+1$$



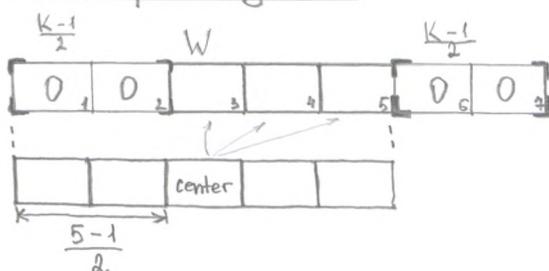
K=2:



$$W = 2$$



Zero padding

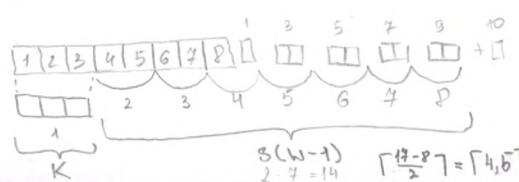


$$K = 2n+1 \in \{3, 5, 7, 9, \dots\}$$

Число нулей $W_{out} = W$ можно заменить $\frac{K-1}{2}$ нулями по краям.

$$P = \frac{K-1}{2}$$

$$P = \left\lceil \frac{S(W-1) - W + K}{2} \right\rceil$$



Без учета гетто $K + S(W-1)$

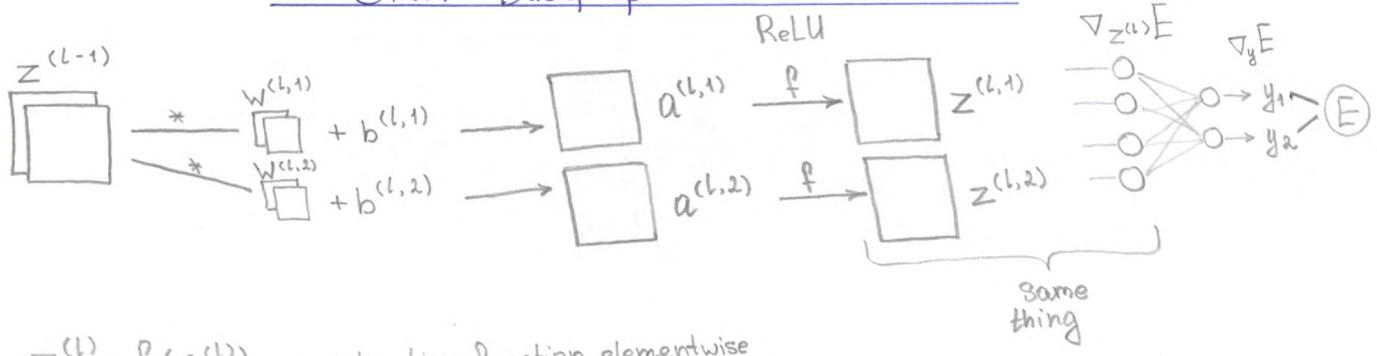
$$\boxed{12345678}$$

$$\boxed{1234}$$

$$4-4$$

$$\left\lceil \frac{4-4}{2} \right\rceil = \lceil 1.5 \rceil = 2$$

CNN Backprop.



$$z^{(l)} = f(a^{(l)}) \quad - \text{activation function elementwise}$$

$$\nabla_{a^{(l)}} E = \left(\frac{d\bar{f}}{da^{(l)}} \right)^T \cdot \nabla_{z^{(l)}} E \quad \frac{dE}{da} = \frac{dE}{dz} \cdot \underbrace{\frac{dz}{df}}_{\bar{f}} \cdot \frac{df}{da}$$

$$\begin{bmatrix} \frac{df_1}{da_1} & \frac{df_2}{da_1} & \dots \\ \frac{df_1}{da_2} & \frac{df_2}{da_2} & \dots \\ \vdots & \vdots & \ddots \end{bmatrix} = \begin{bmatrix} [\bar{f}'(a^{(l,1)})] & 0 & 0 \\ 0 & [\bar{f}'(a^{(l,2)})] & \dots \\ \dots & \dots & \dots \end{bmatrix}$$

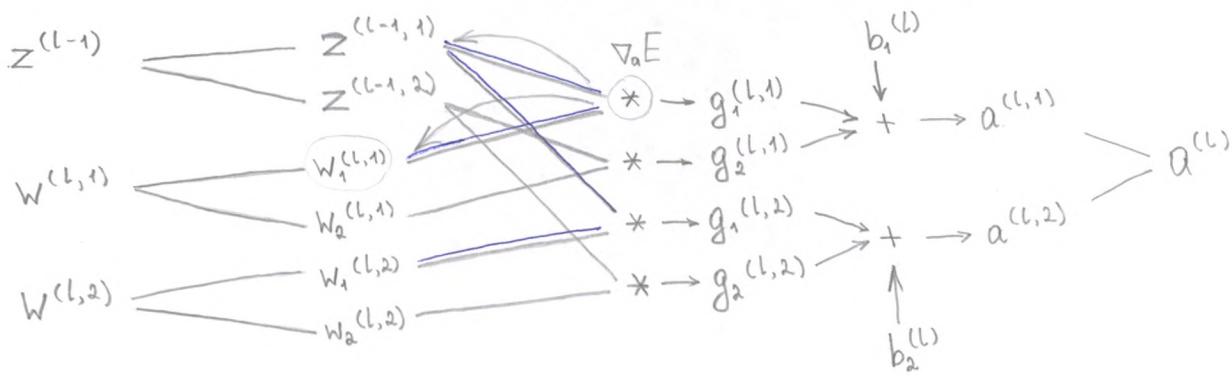
Wegenmehr, wenn wir für zurückfließende $a^{(l)}$ für Fehler rausnehmen

$$\nabla_{a^{(l)}} E = \bar{f}'(a^{(l)}) \circ \nabla_{z^{(l)}} E$$

elementwise

$$\begin{aligned} & \begin{pmatrix} \bar{f}(a_{11}^{(l,1)}) & \bar{f}(a_{12}^{(l,1)}) & \dots \\ \bar{f}(a_{21}^{(l,1)}) & \bar{f}(a_{22}^{(l,1)}) & \dots \\ \vdots & \vdots & \vdots \end{pmatrix} = \\ & = \begin{bmatrix} \frac{dE}{dz^{(l,1)}} \cdot \frac{df}{da_{11}^{(l,1)}} & \dots \\ \frac{dE}{dz^{(l,1)}} \cdot \frac{df}{da_{21}^{(l,1)}} & \dots \end{bmatrix} a \end{aligned}$$

Backprop through convolutions $(Z * W)^T_W$

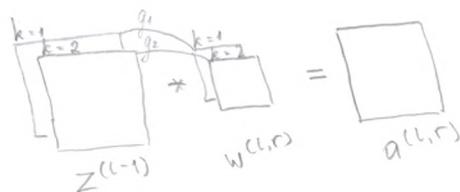


$$g_k^{(l,r)} = Z^{(l-1,k)} * w_k^{(l,r)}$$

$$a^{(l,r)} = \sum_k g_k^{(l,r)} + b_r^{(l)}$$

$$a = \sum_k g_k = \sum (Z * W) \Rightarrow \frac{dE}{dw} = \frac{dE}{da} \cdot \frac{da}{dw}$$

$\nabla_{w_{k(r)}} E - ?$



$$\begin{bmatrix} Z_{11} & Z_{12} & Z_{13} \\ Z_{21} & Z_{22} & Z_{23} \\ Z_{31} & Z_{32} & Z_{33} \end{bmatrix} * \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$$

$$\begin{array}{c} w_{11} \xrightarrow{\cdot Z_{11}} a_{11} \\ \xrightarrow{\cdot Z_{12}} a_{12} \\ \xrightarrow{\cdot Z_{21}} a_{21} \\ \xrightarrow{\cdot Z_{22}} a_{22} \end{array}$$

$$\frac{dE}{dw_{11}} = \frac{dE}{da_{11}} \cdot Z_{11} + \frac{dE}{da_{12}} Z_{12} + \frac{dE}{da_{21}} \cdot Z_{21} + \frac{dE}{da_{22}} \cdot Z_{22}$$

$$\begin{bmatrix} Z_{11} & Z_{12} & Z_{13} \\ Z_{21} & Z_{22} & Z_{23} \\ Z_{31} & Z_{32} & Z_{33} \end{bmatrix} * \begin{bmatrix} \frac{dE}{da_{11}} & \frac{dE}{da_{12}} \\ \frac{dE}{da_{21}} & \frac{dE}{da_{22}} \end{bmatrix}_{3x3}^{2x2} = \begin{bmatrix} \frac{dE}{dw_{11}} & \frac{dE}{dw_{12}} \\ \frac{dE}{dw_{21}} & \frac{dE}{dw_{22}} \end{bmatrix}_{2x2}^{2x2}$$

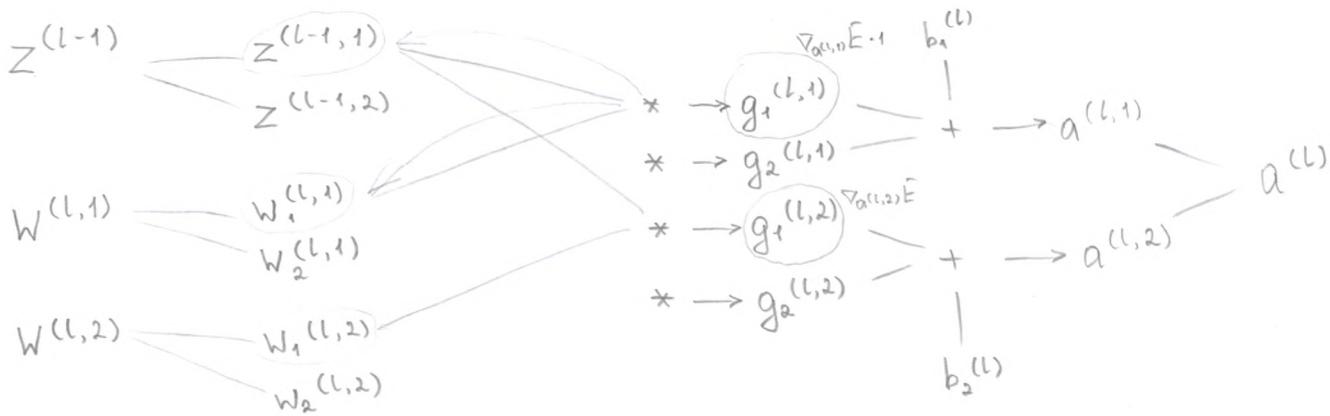
$$Z^{(l-1)} * w^{(l,r)} \quad Z^{(l-1,k)} * \nabla_{a^{(l,r)}} E = \nabla_{w_k^{(l,r)}} E$$

$$\nabla_{w_k^{(l,r)}} E = Z^{(l-1,k)} * \nabla_{a^{(l,r)}} E, \text{ gatunek wykonywania działań podzielony na kanały}$$

$$\nabla_{w^{(l,r)}} E = Z^{(l-1)} * \nabla_{a^{(l,r)}} E$$

□	□	□
□	□	□
□	□	□

Backprop through Convolution $(Z * w)_z'$



$$Z^{(l-1)} * W^{(l)} + b^{(l)}, \quad j = 1, 2, \dots, m_l \times n_l$$

$$\begin{cases} Z^{(l-1,1)} * W_1^{(l,1)} = g_1^{(l,1)} \\ Z^{(l-1,2)} * W_1^{(l,2)} = g_1^{(l,2)} \end{cases} \quad a^l = g^l$$

$$\nabla_{Z^{(l-1,1)}} E - ?$$

$$\tilde{W}_k^{(l,r)} = W_k^{(l,r)} \cdot \text{flip}[E]$$

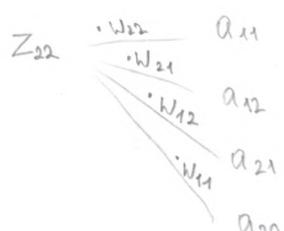
$$\begin{bmatrix} Z_{11} & Z_{12} & Z_{13} \\ Z_{21} & \circled{Z_{22}} & Z_{23} \\ Z_{31} & Z_{32} & Z_{33} \end{bmatrix} * \begin{bmatrix} W_{11} & W_{12} \\ W_{21} & W_{22} \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$$

1 channel

$$\begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & \nabla a_{11} & \nabla a_{12} & 0 \\ 0 & \nabla a_{21} & \nabla a_{22} & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} * \begin{bmatrix} W_{22} & W_{21} \\ W_{12} & W_{11} \end{bmatrix}$$

$\tilde{W}_k^{(l,r)} = \frac{W_k^{(l,r)}}{2}$

$$\sim \nabla_{Z^{(l-1,k)}} E = \text{ZeroPad}(\nabla_{a^{(l,r)}} E) * \tilde{W}_k^{(l,r)}$$



$$\Rightarrow \frac{dE}{dZ_{22}} = \frac{dE}{da_{11}} \cdot W_{22} + \frac{dE}{da_{12}} \cdot W_{21} + \frac{dE}{da_{21}} \cdot W_{12} + \frac{dE}{da_{22}} \cdot W_{11}$$

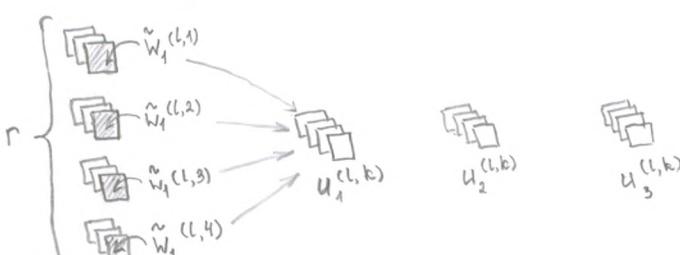
$\nabla_a E \cdot \tilde{W}$ like

Derivative $\square * \square \quad Z^{(l-1,1)} * W_1^{(l,1)}$, no smo ne namne, namny smo on eue zatucum

om smopoi obespmku

$$\boxed{\nabla_{Z^{(l-1,k)}} E = \sum_r \text{ZeroPad}(\nabla_{a^{(l,r)}} E) * \tilde{W}_k^{(l,r)}} = \text{ZeroPad}(\nabla_{a^{(l)}} E) \cdot \underbrace{U^{(l,k)}}_{[W_1^{(l,1)} \ W_1^{(l,2)} \dots \ W_1^{(l,r)}]}$$

$$U_r^{(l,k)} = \tilde{W}_k^{(l,r)}$$

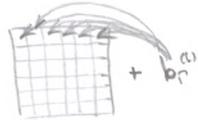


Это интересно б нее нормалю, но на практике все будет наоборот. ZeroPad($\nabla_a E$) * $\tilde{W} = (Z * W)_z'$

Backprop through convolution $(Z^*W + b \cdot 1)_b'$

$$\begin{bmatrix} Z_1^{(l-1)} \\ \vdots \\ Z^{(l-1)} \end{bmatrix} * \begin{bmatrix} W^{(l,r)} \\ b_r^{(l)} \end{bmatrix} = \begin{bmatrix} a_r^{(l)} \end{bmatrix}$$

$$a_r^{(l,r)} = Z^{(l-1)} * W^{(l,r)} + b_r^{(l)}$$

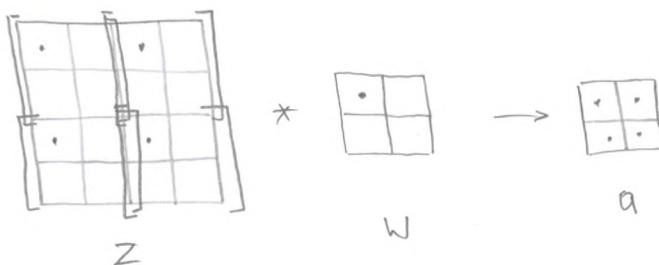


$$\frac{dE}{db_r^{(l,r)}} = \frac{dE}{da_{11}^{(l,r)}} \cdot 1 + \frac{dE}{da_{12}^{(l,r)}} \cdot 1 + \dots$$

$$\Rightarrow \nabla_{b^{(l,r)}} E = \underset{\substack{\uparrow \\ \text{elements of arg.}}}{\text{Sum}} (\nabla_{a^{(l,r)}} E), \quad r = 1, \dots, r^l$$

Backprop through convolution with stride

1)



$$\text{Dilate}(a) = \begin{array}{|c|c|c|} \hline & \cdot & \cdot \\ \hline \cdot & \cdot & \cdot \\ \hline & \cdot & \cdot \\ \hline \end{array}$$

$$\nabla_W E = Z^{(l-1)} * \text{Dilate}(\nabla_a E)$$

2)

$$\begin{array}{ccccc} Z_{11} & Z_{12} & Z_{13} & Z_{14} & Z_{15} \\ Z_{21} & Z_{22} & Z_{23} & Z_{24} & Z_{25} \\ Z_{31} & Z_{32} & Z_{33} & Z_{34} & Z_{35} \\ Z_{41} & Z_{42} & Z_{43} & Z_{44} & Z_{45} \\ Z_{51} & Z_{52} & Z_{53} & Z_{54} & Z_{55} \end{array} * \begin{array}{c} W_{11}^4 \quad W_{12} \quad W_{13}^3 \\ W_{21} \quad W_{22} \quad W_{23} \\ W_{31}^2 \quad W_{32} \quad W_{33}^1 \end{array} = \begin{array}{c} a_{11} \\ a_{12} \\ a_{21} \\ a_{22} \end{array}$$

$$\begin{array}{ll} Z_{33} & \cdot W_{33} \quad a_{11} \\ & \cdot W_{31} \quad a_{12} \\ & \cdot W_{13} \quad a_{21} \\ & \cdot W_{11} \quad a_{22} \end{array}$$

$$\frac{dE}{dZ_{33}} = \underbrace{\frac{dE}{da_{11}} \cdot W_{33}}_{\text{Term 1}} + \underbrace{\frac{dE}{da_{12}} \cdot W_{31}}_{\text{Term 2}} + \underbrace{\frac{dE}{da_{21}} \cdot W_{13}}_{\text{Term 3}} + \underbrace{\frac{dE}{da_{22}} \cdot W_{11}}_{\text{Term 4}}$$

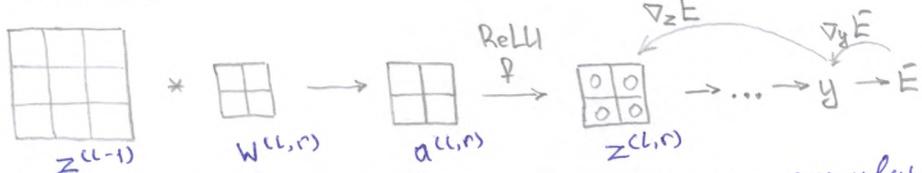
$$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & \frac{dE}{da_{11}} \end{bmatrix} * \begin{array}{c} W_{33} \quad W_{31} \quad W_{13} \\ W_{23} \quad W_{21} \quad W_{11} \\ W_{13} \quad W_{11} \end{array} = \nabla_Z E$$

\tilde{W}

$$\text{ZeroPad}(\text{Dilate}_{S-1}(\nabla_a E))$$

CNN Learnability

1 Channel



Следующие пункты обсуждаются методом активированных feature map.

$\nabla_{z^{(l,r)}} E$ not zero

$$\begin{array}{c} \nabla_{z^{(l,r)}} E \\ \text{GMP} \end{array} \rightarrow \begin{array}{c} E = y + t \\ y = 0, t = 1 \end{array}$$

$$\nabla_{z^{(l,r)}} E = \nabla_y E$$

$\nabla_{a^{(l,r)}} E$

$$\text{ReLU}'(a^{(l,r)}) = \begin{cases} 1, & \text{if } a > 0 \\ 0, & \text{if } a \leq 0 \end{cases}$$

$$Z^{(l-1)} * W^{(l,r)} = a^{(l,r)}$$

$\nabla_{W^{(l,r)}} E$

$$\begin{bmatrix} Z_{11} & Z_{12} & Z_{13} \\ Z_{21} & Z_{22} & Z_{23} \\ Z_{31} & Z_{32} & Z_{33} \end{bmatrix} * \begin{bmatrix} W_{11} & W_{12} \\ W_{21} & W_{22} \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$$

Все бека беком на новый layer

$$W_{11} \begin{bmatrix} Z_{11} & a_{11} \\ Z_{21} & a_{12} \\ Z_{31} & a_{21} \end{bmatrix} \Rightarrow \frac{dE}{dW_{11}} = \frac{dE}{dZ_{11}} \frac{da_{11}}{dZ_{11}} + \frac{dE}{dZ_{21}} \frac{da_{12}}{dZ_{21}} + \dots$$

$$\begin{bmatrix} Z_{11} & Z_{12} & Z_{13} \\ Z_{21} & Z_{22} & Z_{23} \\ Z_{31} & Z_{32} & Z_{33} \end{bmatrix} * \begin{bmatrix} \frac{dE}{da_{11}} & \frac{dE}{da_{12}} \\ \frac{dE}{da_{21}} & \frac{dE}{da_{22}} \end{bmatrix} = \begin{bmatrix} \frac{dE}{dW_{11}} & \frac{dE}{dW_{12}} \\ \frac{dE}{dW_{21}} & \frac{dE}{dW_{22}} \end{bmatrix}$$

ReLU'

Если $\nabla_a E = 0$, то слой не обрабатывается.

Если нормализованное значение активированось $\rightarrow \frac{dE}{da_i} > 0$ то слой обрабатывается.

$\nabla_{z^{(l-1)}} E = \text{ZeroPad}(\nabla_{a^{(l,r)}} E) * W^{(l,r)}$

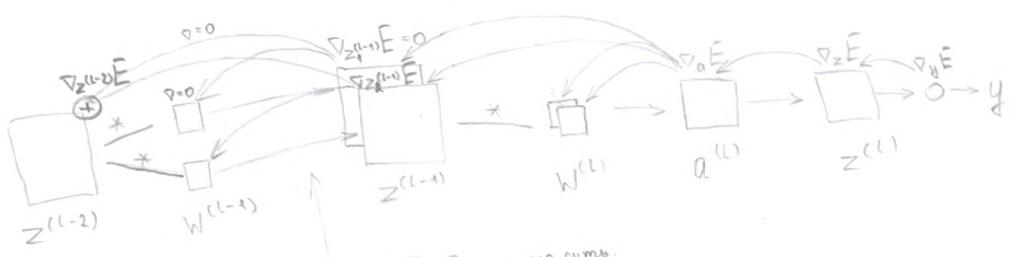
$$\begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & \frac{dE}{da_{11}} & \frac{dE}{da_{12}} & 0 \\ 0 & \frac{dE}{da_{21}} & \frac{dE}{da_{22}} & 0 \\ 0 & 0 & 0 & 0 \end{bmatrix} * \begin{bmatrix} W_{22} & W_{21} \\ W_{12} & W_{11} \end{bmatrix} = \begin{bmatrix} \frac{dE}{dZ_{11}} & \frac{dE}{dZ_{12}} & \frac{dE}{dZ_{13}} \\ \frac{dE}{dZ_{21}} & \frac{dE}{dZ_{22}} & \frac{dE}{dZ_{23}} \\ \frac{dE}{dZ_{31}} & \frac{dE}{dZ_{32}} & \frac{dE}{dZ_{33}} \end{bmatrix}$$

Здесь снова если нормализованное значение активированось, то ненулевое предыдущее слои будут обрабатываться. (?)

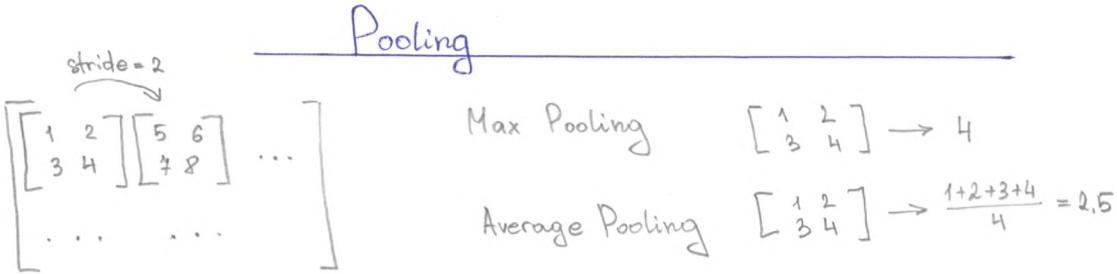
↓ предыдущим можно заменить на каналы, но он тоже будет суммироваться с зрачением групп каналов и обработка пропадет.

$$(a^{(l-1,r)})' = \underbrace{\frac{dE}{dZ^{(l-1)}}}_{\nabla_{Z^{(l-1)}} E} \cdot \text{ReLU}'(a^{(l-1,r)})$$

$\nabla_{Z^{(l-1)}} E$



If my a^{(l-1)} is zero, we're zero.



stride = kernel size

Global Pooling

kernel size = input tensor size

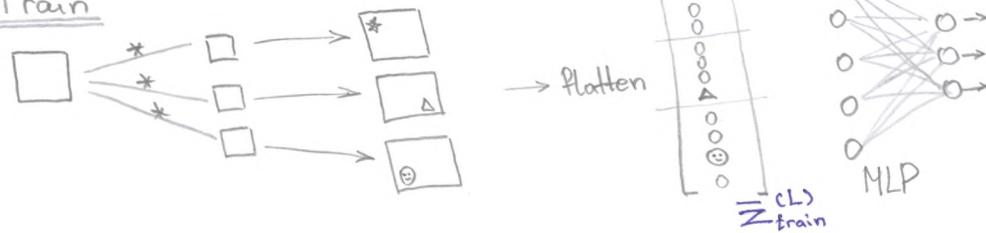
$$\begin{bmatrix} 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \\ 1 & 2 & 3 & 4 \\ 5 & 6 & 7 & 8 \end{bmatrix}$$

Global Max Pooling $\stackrel{(GMP)}{\rightarrow} 8$

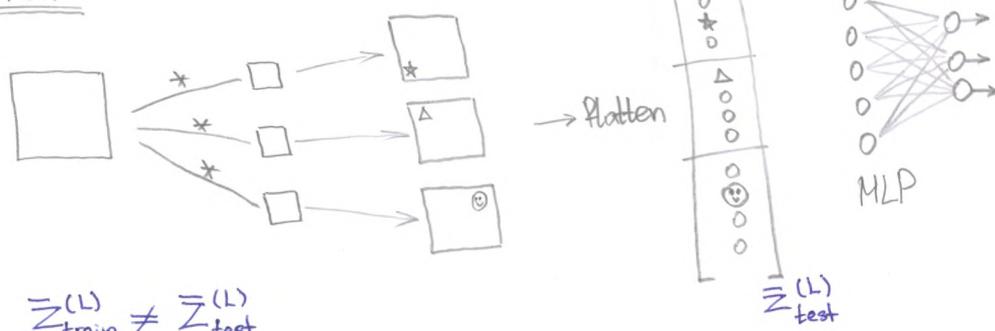
Global Average Pooling $\stackrel{(GAP)}{\rightarrow} \frac{2 \sum_{i=1}^8}{16}$

Why to use GMP

Train



Test

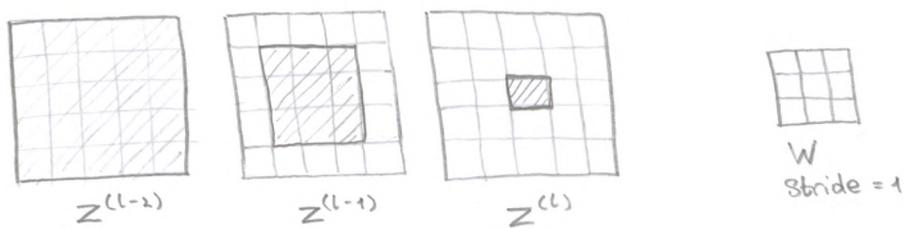


CNN затягивает, искажая информацию в feature maps, но если генерировать flatten, то main taskaa one problema, как Image \rightarrow MLP.

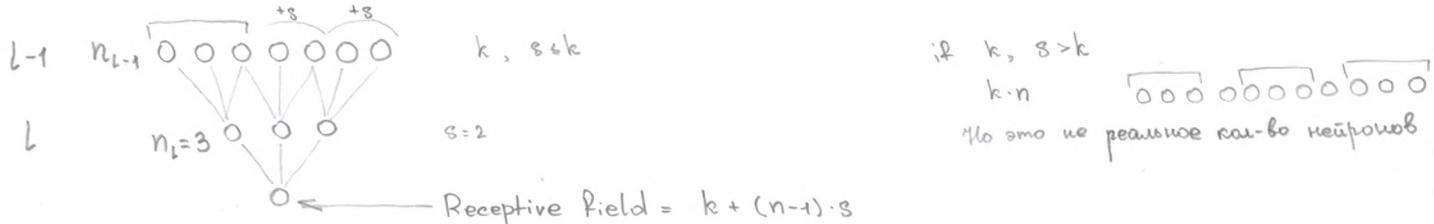
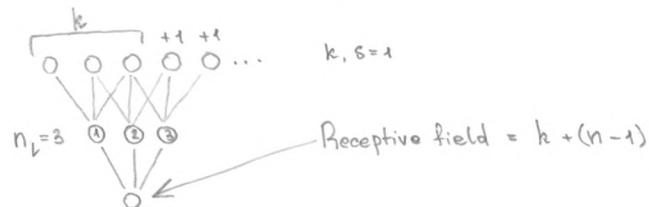
→ Fully Convolutional Neural Network

Receptive Field

$1 \rightarrow 3 \rightarrow 5$

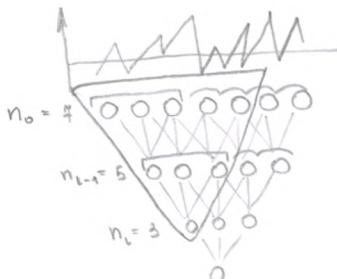


Receptive field - это множество пикселей изображения, от которых зависит активация определенного нейрона.



$$n_{L-1} = k + (n_r - 1) \cdot s$$

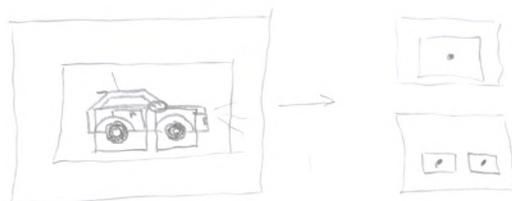
n_r - число нейронов $L-1$ -го слоя, будущий нейрон



знач receptive field нейрона, может ли это определить cf следующего слоя за $O(1)$? Кажется, нет

Dilated convolution
увеличиваём Receptive field
но...

В зависимости от Receptive field нейрон может распознавать предметы определенного размера, или относительное расположение предмета.



CTC for OCR



1. Classify 1M words. Problems: 1) if not in Domain, we can't classify, 2) A lot of images per each word

2. Slice into letters. Problems: 1) Hard task, 2) Annotation for each letter

3. Seq-to-seq. Problem: frame-wise alignment

End-to-end solution - Connectionist Temporal Classification

Уговаривая газа Speech Recognition, Audio → Text
u RNN

Мы используем Fully CNN

- L - конечное множество символов ($ABC, ., :, ?, \dots$ etc.)
- $\mathcal{X} = L^*$ - множество последовательностей символов L
 $\{A\}, \{B\}, \{C\}, \{AB\}, \{BA\}, \{ABC\}, \{ACB\}, \{BAC\}, \{BCA\}, \{CAB\}, \{CBA\}$
Но если можно рассматривать, как слова.
- $\mathcal{X} = (R^m)^*$ - множество последовательностей m -мерных векторов.
 $\{[1], [1]^m, [1]^m, [1]^m, \dots\}$
Но если это изображение $\in R^{m \times n}$, высотой m и произвольной шириной.
- S - множество обучающих примеров из распределения $P_{x, z}$
 $f(x) \rightarrow \mathcal{X}$
 $f((R^m)^*) \rightarrow L^*$
 $f(\text{Image}) \rightarrow L^*$



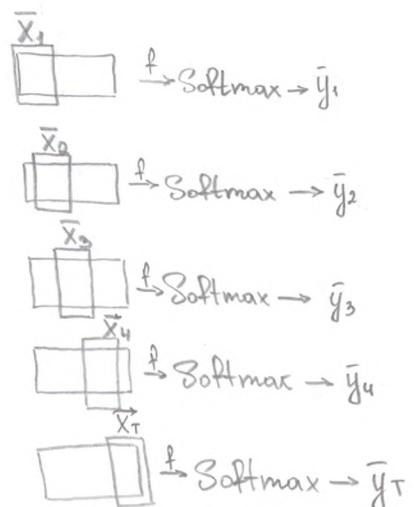
$x = (\bar{x}_1, \bar{x}_2, \bar{x}_3, \dots, \bar{x}_T)$ - Input sequence

$z = (z_1, z_2, z_3, \dots, z_u)$ - Target sequence

$$u \leq T$$

$(x, z) \in S, x \in \mathcal{X}, z \in \mathcal{Z}$

$$f(x) \rightarrow z$$



T - кол-во шагов

- W W - - O O O - R - D D D D -

Z = WORD

Training a Classifier

L - alphabet

$Z = L^*$ - sequence of letters

$X = (R^n)^*$ - sequence of image parts

S - ~~mnogoecmbo~~ $D_{X \times Z}$

$X = (\bar{x}_1, \dots, \bar{x}_T)$ - input sequence (image)

$Z = (z_1, \dots, z_u)$ - target sequence (word)

$u \leq T$

$(X, Z) \in S, X \in \mathcal{X}, Z \in \mathcal{Z}$

Train a classifier on S dataset

$$h : \mathcal{X} \rightarrow \mathcal{Z}$$

images words

Label Error Rate наиболее ошибочную
(LER)

S' - validation dataset

$$\text{LER}(h, S') = \frac{1}{|S'|} \sum_{(x, z) \in S'} \text{ED}(h(x), z) \rightarrow \min$$

$\text{ED}(a, b)$ - минимальное кол-во операций вставки, удаления и замены, необходимых для нарушения последовательности b из a .

$$\text{ED}(a = \text{bank}, b = \text{bank}) = 1$$

ED - Euclidean Distance?

Levenshtein Distance?

CTC

$L' = L \cup \{\text{blank}\}$, blank - нейтральный символ, монодуктивное пространство

L'^* - множество последовательностей.

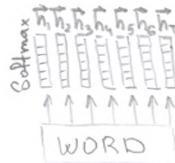
i.g. $I = (l_1, l_2, l_3 \dots, l_u)$ - target последовательность, $I \in L^*$

$I' = (-, l_1, -, l_2, -, l_3, -, \dots, -, l_u, -)$, $I' \in L'^*$

$$|I'| = 2 \cdot |I| + 1$$

CNN $\rightarrow - C - N - N -$

$\bar{x} = (\bar{x}_1, \bar{x}_2, \bar{x}_3, \dots, \bar{x}_T)$ - input sequence



$h(\bar{x}_t) \xrightarrow[\text{Softmax}]{} R^{L'}$ - i.e. предсказание на каждую Symbol апerture + blank

$y_{\pi_t}^t$ π_t - symbol from L'

Для каждого из t Softmax-ов y несёт информацию о labelах, y несёт разные info по Symbol.

Assumption (strong)

$$P(\pi | x) = \prod_{t=1}^T y_{\pi_t}^t, \forall \pi \in L'^T$$

вероятность
символа π_t на шаге t

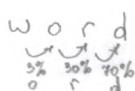
$$\begin{bmatrix} \vec{y}^1 \\ y_1^1 \\ y_2^1 \\ y_3^1 \\ y_4^1 \\ \vdots \end{bmatrix}, \begin{bmatrix} \vec{y}^2 \\ y_1^2 \\ y_2^2 \\ y_3^2 \\ y_4^2 \\ \vdots \end{bmatrix}$$

$$\bar{\pi} = [\pi_1 \ \pi_2 \ \pi_3 \dots]^T$$

Symbol
number

- Так будем y $y_{\pi_t}^t$ с точной концепцией верхний и нижний индекс местами, \vec{y}^t - вектор на шаге t
- $y_{\pi_t}^t$ - конкретное значение вектора

- Проблема: последовательности π - это произведение уверенности модели в символах текущих шагов, которое она предсказывает. Так Independent Events, но это кажется Dependent.



Но, таким образом мы можем легко распознавать слова с ошибками или фрагментами

CTC

$$L' = L \cup \{\text{blank}\} - \text{alphabet} \quad \bar{L} = L'^*$$

L'^* - последовательности символов и blank

$\bar{L} \in L^*$, for example $\bar{L} = \text{hello}_5$

$$\bar{L}' \in L'^* \quad \bar{L}' = -h-e-l-l-o-, |\bar{L}'| = 2|\bar{L}| + 1$$

$$x = (\bar{x}_1, \bar{x}_2, \bar{x}_3, \dots, \bar{x}_T), T - \text{width} \quad x \in \mathcal{X}$$

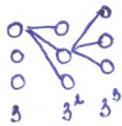
$h(\bar{x}_t) \rightarrow R^{|\bar{L}'|}$ вероятность на каждого символа $\in L'$

$$h(\bar{x}_t) = \bar{y}^t$$

path

$$\bar{\pi} = [\pi_1, \pi_2, \pi_3, \dots, \pi_T]^T, \pi_t - \text{symbol}$$

$$P(\bar{\pi} | x) = \prod_{t=1}^T p_{\pi_t}^t, \forall \pi_t \in L^T$$



$$\bar{L}'^T \in L'^*$$

↑
послед. груп T

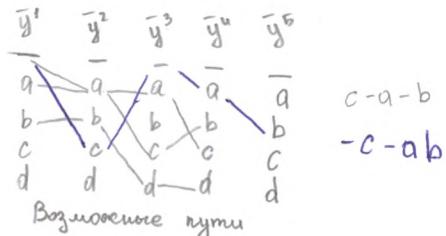
Decoding

$$B: \bar{L} \rightarrow \mathcal{X}$$

1. Удаление нотных: AAA → A

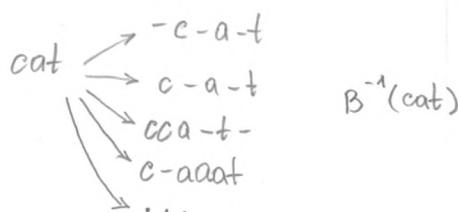
2. Удаление blank: -AA-A → AAA

$$B(-h-e-e-l-l-o) \rightarrow \text{hello}$$



$$P(\bar{L} | x) = \sum_{\bar{\pi} \in B^{-1}(L)} P(\bar{\pi} | x) \quad - \text{Likelihood}$$

Вероятность последовательности \bar{L} при x - это сумма вероятностей всех предсказанных последовательностей, которые заканчиваются в \bar{L} , а не в один конкретный путь.



Возможно количество возможных путей.

$$h(x) = \underset{\bar{L} \in L'^*}{\operatorname{argmax}} P(\bar{L} | x) \quad - \text{наибольшая вероятность последовательности, среди них.}$$

Forward - Backward Algorithm

$$d_t(s) \stackrel{\text{def}}{=} \sum_{\pi \in N^T: \pi_{1:t} = s} \prod_{t'=1}^t y_{\pi_{t'}}^{t'}$$

$B(\pi_{1:t}) = \underbrace{B(L_{1:s}^t)}_{\in L}$

$$N^T = L'^T, \quad L'^T \in L'^*$$

Сума правдоподій всіх підмножин, які складаються з $L_{1:s}^t$

$L = \text{hello}$

$L' = -h-e-l-l-o -$

$L_{1:3} = -h-$

$\pi = h-ee-ll-l-oo -$

$d_t(s)$ - правдоподійність того, що перші t символів S будуть групуватися в перші $\sum_{s=1}^t$ символів L' .

$$B(\bar{y}^1 \bar{y}^2 \bar{y}^3 \dots \bar{y}^T)$$

$$B(-l_1 - l_2 - l_3 - \dots - l_u)$$

\bar{y}^i - softmax i.e. $[P(\text{blank}), P(A), P(B), P(C), \dots, P(Z)]^T$

$$d_1(1) = y_b^1$$

$$d_1(2) = y_{l_1}^1$$

$$\bar{y}^1 \downarrow \\ -l_1 -$$

найлічніша норма всіх відповідей на початку спільної прогноза l_1

$$d_1(s) = 0, \forall s > 2$$

Оскільки прогнозуваннями ми не можем прогнозувати 2 та більше символів

Recursion

Forward - Backward Algorithm

$$\bar{d}_t(s) \stackrel{\text{def}}{=} d_{t-1}(s) + d_{t-1}(s-1)$$

$$d_t(s) = \begin{cases} \bar{d}_t(s) y_{l_s^t}^t, & \text{if } l_s^t = b, \text{ or } l_{s-2}^t = l_s^t \\ (\bar{d}_t(s) + d_{t-1}(s-2)) y_{l_s^t}^t, & \text{otherwise} \end{cases}$$

- ищем запрещающие ноты как blank.
- гомогенность окон.

1) $d_t(s) = \bar{d}_t(s) y_{l_s^t}^t$, if $l_s^t = b$ - H-E-L-L-O-

 $\bar{d}_t(s) = d_{t-1}(s) + d_{t-1}(s-1)$

$d_t(s)$ - вероятность, что ноте t откроется в s и закончится на ноте b .

$\bar{d}_t(s)$ - вероятность того, что $t-1$ откроется в s , и это $s-1$ закончится на b .

• $d_{t-1}(s)$

$$\pi_{1:t-1} = \overbrace{\dots \text{HH} \dots \text{EE}}^t -$$

$$l_{1:s}^t = \xleftarrow[s]{\text{H-E}} \quad t \quad \text{nothing changes}$$

• $d_{t-1}(s-1)$

$$\pi_{1:t-1} = \overbrace{\dots \text{HH} \dots \text{EE}}^t -$$

$$l_{1:s-1}^t = \xleftarrow[s-1]{\text{H-E}} \quad s$$

2) $d_t(s) = \bar{d}_t(s) y_{l_s^t}^t$, if $l_s^t = l_{s-2}^t$ - E-M-M-A

$$\bar{d}_t(s) = d_{t-1}(s) + d_{t-1}(s-1)$$

• $d_{t-1}(s)$

$$\pi_{1:t-1} = \overbrace{\dots \text{EE} \dots \text{MM} \dots \text{M}}^t M$$

$$l_{1:s}^t = \xleftarrow[s]{\text{E-M-M}} \quad t \quad \text{nothing changes}$$

• $d_{t-1}(s-1)$

$$\pi_{1:t-1} = \overbrace{\dots \text{EE} \dots \text{MM} \dots}^t M$$

$$l_{1:s-1}^t = \xleftarrow[s-1]{\text{E-M-}} \quad s$$

Forward - Backward Algorithm

$$\bar{d}_t(s) = d_{t-1}(s) + d_{t-1}(s-1)$$

$$d_t(s) = \sum_{\pi \in NT: \pi_{1:t} = l_{1:s}^t} y_{\pi_t}^t$$

$B(\pi_{1:t}) = B(l_{1:s}^t)$

$$d_t(s) = \begin{cases} \bar{d}_t(s) \cdot y_{l_s^t}^t & , \text{ if } l_s^t = b \text{ or } l_{s-2}^t = l_s^t \\ (\bar{d}_t(s) + d_{t-1}(s-2)) \cdot y_{l_s^t}^t & , \text{ otherwise} \end{cases}$$

3) General case

$$d_t(s) = (\bar{d}_t(s) + d_{t-1}(s-2)) \cdot y_{l_s^t}^t = (\underbrace{d_{t-1}(s)}_{\cdot d_{t-1}(s)} + \underbrace{d_{t-1}(s-1)}_{\cdot d_{t-1}(s-1)} + \underbrace{d_{t-1}(s-2)}_{\cdot d_{t-1}(s-2)}) \cdot y_{l_s^t}^t$$

$$\pi_{1:t-1} = \{ \text{HH--EE--L} \}^t$$

$$l_{1:s}^t = \xleftarrow[s=6]{} \text{H-E-L} \xrightarrow[t]{} \text{H-E-L}$$

$$\cdot d_{t-1}(s-1)$$

$$\pi_{1:t-1} = \{ \text{HH--EE--L} \}^t \quad \text{H-E-L}$$

$$l_{1:s-1}^t = \xleftarrow[s=5]{} \text{H-E-} \xrightarrow[t]{} \text{H-E-L}$$

$$\cdot d_{t-1}(s-2)$$

$$\pi_{1:t-1} = \{ \text{HH--EE--L} \}^t$$

$$l_{1:s-2} = \xleftarrow[s=4]{} \text{H-E} \xrightarrow[t]{} \text{H-EL}$$

Дополнительной строкой, когда нет
межем разрешимых промежутков blank.

B некое выражение $B(\pi_{1:t}) = HEL$, то есть наследственность генотипов от $l_{1:s}^t$

$$d_t(s) = \sum_{\pi \in NT: \pi_{1:t} = l_{1:s}^t} y_{\pi_t}^t \cdot \prod_{t'=1}^{t-1} y_{\pi_{t'}}^{t'} = |\pi_t = l_s^t| = y_{l_s^t}^t \cdot (\underbrace{d_{t-1}(s)}_{B(\pi_{1:t}) = B(l_{1:s}^t)} + \underbrace{d_{t-1}(s-1)}_{B(\pi_{1:t}) = B(l_{1:s}^t)} + \underbrace{d_{t-1}(s-2)}_{B(\pi_{1:t}) = B(l_{1:s}^t)})$$

Forward - Backward Algorithm

Likelihood of particular path

$$P(\bar{\pi} | x) = \prod_{t=1}^T y_{\pi_t}^t$$

Likelihood of target sequence \bar{I}

$$P(\bar{I} | x) = \sum_{\pi \in B^{-1}(\bar{I})} P(\bar{\pi} | x) = \sum_{\pi \in B^{-1}(\bar{I})} \prod_{t=1}^T y_{\pi_t}^t = \sum_{\substack{\pi \in N^T : \\ B(\pi_{1:T}) = B(I_1:1^T)}} \prod_{t=1}^T y_{\pi_t}^t$$

сума (беск) нрабгоногобин беск нумеи, генодуруюшиса б \bar{I} .

Likelihood of subsequence $I'_{1:s}$

$$d_t(s) = \sum_{\pi \in N^T : t'=1}^t y_{\pi_{t'}}^t = \sum_{\pi \in N^T : t'=1}^t y_{\pi_t}^t \cdot \prod_{t'=t+1}^{t-1} y_{\pi_{t'}}^t$$

$$B(\pi_{1:t}) = B(I'_{1:s}) \quad B(\pi_{t+1}) = B(I'_{s:s})$$

$$d_1(1) = y_b^1$$

$$d_1(2) = y_{L_1}^1$$

$$d_1(s) = 0, \forall s > 0$$

$$\bar{d}_t(s) = d_{t-1}(s) + d_{t-1}(s-1)$$

$$d_t(s) = \begin{cases} y_{L_s}^t \cdot \bar{d}_t(s), & \text{if } l'_s = b \text{ or } l'_s = l'_{s-1} \\ y_{L_s}^t \cdot (d_t(s) + d_{t-1}(s-1)), & \text{otherwise} \end{cases}$$

-H-	{ -	-L-L	{ L
-H	{ -	-L-	{ L
-H-E-L	{ L		
-H-E-	{ L		
-H-E	{ L		

Likelihood of target sequence \bar{I}'

$$P(\bar{I}' | x) = d_T(|\bar{I}'|) + d_T(|\bar{I}'|-1)$$

\downarrow \downarrow

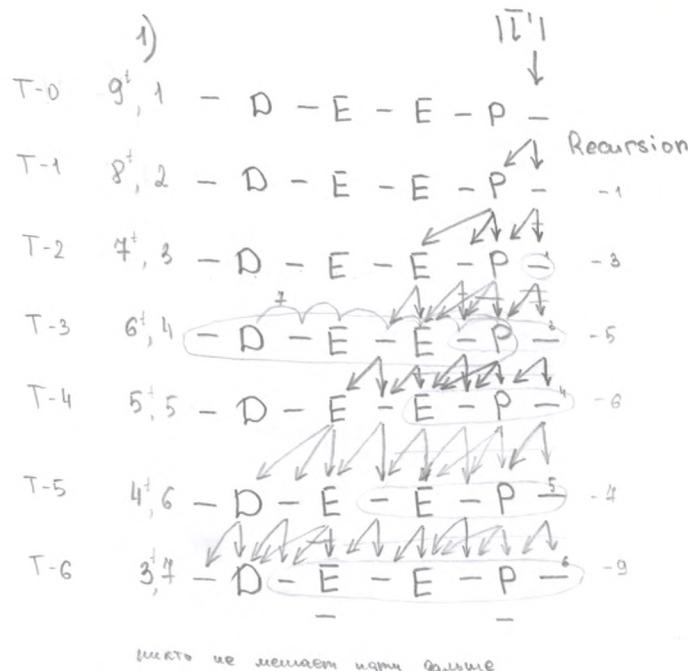
we can skip last blank

Image $\boxed{}$
 $m \times n$
 $n > m$

Forward - Backward Algorithm

$$P(\bar{l}^t | x) = d_t(|\bar{l}'|) + d_t(|\bar{l}'|-1)$$

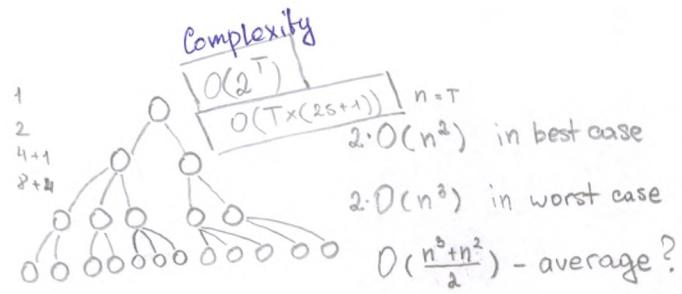
$$d_t(s) = \begin{cases} y_{l'_s}^t \cdot (d_{t-1}(s) + d_{t-1}(s-1)), & \text{if } l'_s = b \text{ or } l'_s = l'_{s-2} \\ y_{l'_s}^t \cdot (d_{t-1}(s) + d_{t-1}(s-1) + d_{t-1}(s-2)), & \text{otherwise} \end{cases}$$



множество не меняет свою генерацию

Если мы будем выбирать blank достаточно, например, то после каждого-то элемента мы не успеем уже свою распознать, поэтому, кажется, можно отбросить некоторые символы.

$$\begin{aligned} d_t(|\bar{l}'|) &= y_b^t \cdot (d_{t-1}(|\bar{l}'|) + d_{t-1}(|\bar{l}'|-1)) \\ 1) d_{t-1}(|\bar{l}'|) &= y_{l'_1}^{t-1} (d_{t-2}(|\bar{l}'|) + d_{t-2}(|\bar{l}'|-1)) \\ 2) d_{t-1}(|\bar{l}'|-1) &= y_{l'_{11}}^{t-1} (d_{t-2}(|\bar{l}'|-1) + d_{t-2}(|\bar{l}'|-2) + d_{t-2}(|\bar{l}'|-3)) \end{aligned}$$



$\frac{l_1 - l_2}{n^2 n^3 n^2 n^2}$
or average just $O(n^3)$
но это не проходит некоторой
 $T \times T - (S-1-\Theta)^2$
С этим приоритетом
беско + раз $d_t(s)$ и все

$$T \leq |\bar{l}'|$$

$$1) \text{ if } d_t(s), t < s-1 \text{ no } 100\% \text{ if } t: 1 2 3 4, s: 1 2 3 4 5, \text{ but } s: 1 2 3 4 5 6$$

2) In example $d_{t-2}(|\bar{l}'|) = d_{t-2}(9)$ не выполняется (это неправильно)

$d_{t-2}(8)$ выполняется, но из-за того что $l'_8 \neq b$ and $l'_8 \neq l'_6$, ищите больше b не упоминаются $d_{t-2}(9)$ не могут выполняться для $s=6$, если $s=6$ не $E-E$

To есть не последовательные символы разбиваются на 2 группы, а не на один

$t + (\text{как-то не надо символов}) < s-1$

$$t < s-1 - (2)$$

$$t < s-3$$

$$7 < 9-3 = 6$$

$$7 < 11-3 = 8$$

Мы берём бинарные, кроме того
когда символ $l_t \neq l_{t-2}$ and $l_t \neq b$

Pruning

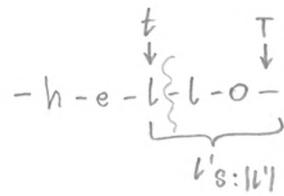
if in $d_t(s)$
 $\Rightarrow t < s-1 - \text{как-то не надо символов, кроме него} \rightarrow$
 завершаем функцию $f(t)$
 не s , а остальные нужны, это может помешать

$$\text{Minimalnoe 1 Symbol} = m \text{ pixels } \boxed{0} \boxed{1} \boxed{0} \boxed{0} \dots \boxed{1} \boxed{0} \boxed{0} \quad \boxed{0} \boxed{1} \boxed{0} \quad \boxed{0} \boxed{1} \boxed{0}$$

Forward-Backward Algorithm

$$\beta_t(s) = \sum_{\pi \in NT} \prod_{t'=t}^T y_{\pi_{t'}}^{t'}$$

$\beta(\pi_{t:T}) = \beta(l'_s:|l'|)$



$$\beta_T(|l'|) = y_b^T$$

- h - e - l - l - o -
↑
 $t = |l'|$

$$\beta_T(|l'|-1) = y_{l_{|l'|}}^T$$

- h - e - l - l - o -
↑
 $t = |l'| - 1$

$$\beta_T(s) = 0, \forall s < |l'| - 1$$

Не имеем озаглав.

$\beta_t(s)$ - бессмысльно озаглавить с кубаров, наименование $t=20$

$$\bar{\beta}_t(s) = \beta_{t+1}(s) + \beta_{t+1}(s+1)$$

$$\beta_t(s) = \begin{cases} \bar{\beta}_t(s) \cdot y_{l'_s}^t, & \text{if } (y_{l'_s}^t) l'_s = b \text{ or } l'_{s+2} = l'_s \\ (\bar{\beta}_t(s) + \beta_{t+1}(s+2)) y_{l'_s}^t, & \text{otherwise} \end{cases}$$

$$\beta_{t+1}(s)$$

- { - l - l - o - }
l'_{s:|l'|}

if $l'_s = b$ or $l'_{s+2} = l'_s$

$$\beta_{t+1}(s+1)$$

- { l - l - o - }
l'_{s+1:|l'|}

не запрещаем пропускать blank

$$\beta_{t+1}(s+2)$$

e { l - l - o - }
l'_{s+2:|l'|}

e { - l - l - o - }
l'_{s+1:|l'|}

e { e - l - l - o - }
l'_s:|l'|

при обмене распределение,

мы разрешаем пропускать blank

Forward - Backward Algorithm

$\prod_x x = 0$, $x < 1$ - underflow

$$C_t = \sum_s d_t(s) \quad \hat{d}_t(s) = \frac{d_t(s)}{C_t}$$

$$D_t = \sum_s \beta_t(s) \quad \hat{\beta}_t(s) = \frac{\beta_t(s)}{D_t} \quad - \text{Normalized}$$

нормированием коэффициентов

$\hat{d}_t(s), \hat{\beta}_t(s)$ - вероятности того, что на t очередь будут символы
(s -й символ)

$\hat{d}_t(s)$ - будущие вероятности s символов l'

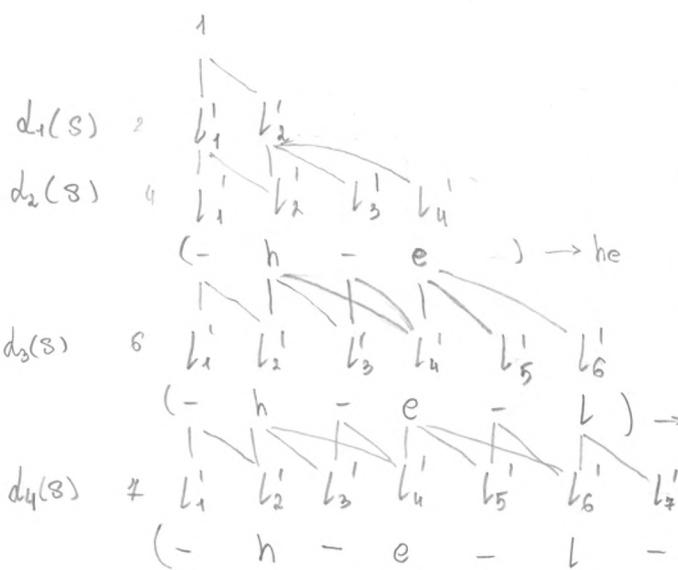
$\hat{\beta}_t(s)$ - будущие ожидания s символа g конца l'

Причина, что в любом случае вероятность $d_t(s)$
и $\beta_t(s)$. Сложность $O(T \cdot S)$.

Но, если $d_1(2), d_1(3) \dots d_1(S) = 0$ и в маине $t=1$.

Последнее $t < s-1-n$ $d_t(s)=0$, то есть фразами, не имеющими собственных

и не b ? Обратно пакетом на 2, но если $l'_s = b$ or $l'_s = l'_{s-2}$,
то пакетом на 1 без



2t - (если $b_t = l_{t+1}$ то b)

Но, если на первом пакете g есть, то на втором есть $d_t(s)$ и $\beta_t(s)$
 $d_t(l'')$

~~β~~

Forward - Backward Algorithm

Loss Function

$$P(S, \mathcal{N}_w) = \prod_{(x,z) \in S} p(z|x) \rightarrow \max$$

$$D^{ML}(S, \mathcal{N}_w) = - \sum_{(x,z) \in S} \ln(p(z|x))$$

$$-\ln(P(S, \mathcal{N}_w)) = - \sum_{(x,z) \in S} \ln(p(z|x)) \rightarrow \min$$

D - objective function

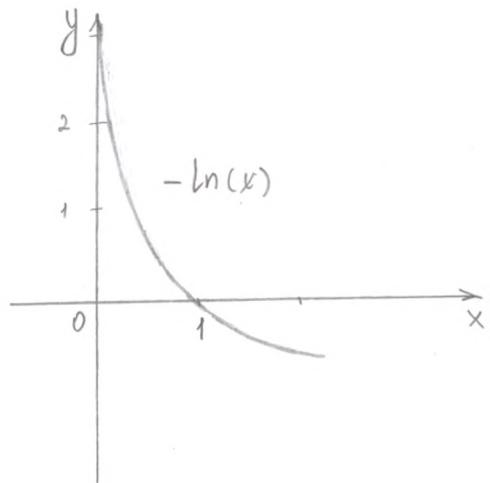
ML - maximum likelihood

S - dataset

\mathcal{N}_w - model

x - input sequence.

z - target sequence.

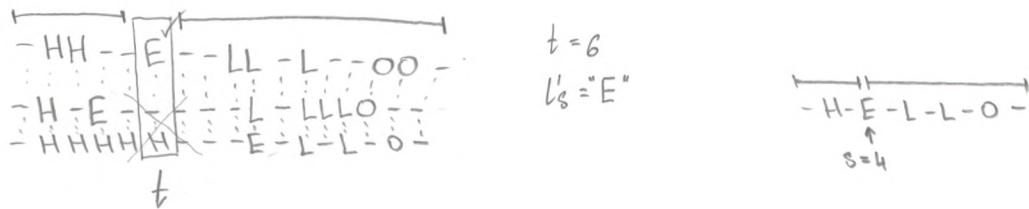


$$\frac{d D^{ML}(z|x, z), \mathcal{N}_w}{d y_k^t} = \frac{-d \ln(p(z|x))}{d y_k^t}$$

partial derivative by some y_k^t

Forward - Backward Algorithm

$$d_t(s) \cdot \beta_t(s) = \sum_{\substack{\pi \in B^{-1}(l) : \\ \pi_t = l'_s}} y_{l'_s}^t \cdot \prod_{t'=1}^T y_{\pi_{t'}}^{t'}$$



$d_t(s)$ - бераамносцы га t үрагамын $g(s)$, мөсөм l'_s :

$\beta_t(s)$ - бераамносцы га t үрагамын $om(s)$ го тооцга, мөсөм $l'_s:l'_1$

То осмыс $\pi_t = l'_s$, маки как бүх $d_t(s)$ и бүх $\beta_t(s)$ нь
нормалуулсан бичсүүлэлт l'_s нь рекурсивно проходжин хо биелийн возможжилсийн
предогоджуулсан нүүрши, иши носигоджуулсан багасгасан $\beta_t(s)$.

$$d_t(s) = \sum_{\pi \in N^T} \prod_{t'=1}^T y_{\pi_{t'}}^{t'}$$

$$B(\pi_{1:t}) = B(l'_{1:s})$$

$$\beta_t(s) = \sum_{\pi \in N^T} \prod_{t'=t}^T y_{\pi_{t'}}^{t'}$$

$$B(\pi_{t:T}) = B(l'_{s:T})$$

$$d_t(s) = \sum_{\pi \in N^T} y_{\pi_t}^t \prod_{t'=1}^{t-1} y_{\pi_{t'}}^{t'}$$

$$B(\pi_{1:t}) = B(l'_{1:s})$$

$$\beta_t(s) = \sum_{\pi \in N^T} y_{\pi_t}^t \prod_{t'=t+1}^T y_{\pi_{t'}}^{t'}$$

$$B(\pi_{t:T}) = B(l'_{s:T})$$

$$\begin{aligned}
 d_t(s) \cdot \beta_t(s) &= (abc + fgh) \cdot (cde + hij) = abc \cdot cde + fgh \cdot hij + abc \cdot hij + fgh \cdot cde = \\
 &= \underbrace{c(\underbrace{abcde}_{\text{full word path}})}_{\text{non-overlapping}} + h(\underbrace{fghij}_{\text{full word path}}) + 0
 \end{aligned}$$

нечастотные
носигоджамжилсийн

$$d_t(s) \cdot \beta_t(s) = \sum_{\pi \in B^{-1}(l)} \underbrace{y_{\pi_t}^t}_{\pi_t = l_s^t} \cdot \underbrace{\prod_{t'=1}^T y_{\pi_{t'}}^t}_{p(\pi | x)} \quad \begin{array}{l} \text{const} \\ \downarrow \\ t = \text{const} \\ s = \text{const} \end{array} \quad \begin{array}{l} d_1(3) = 0 \\ d_2(5) = 0 \\ \dots \end{array}$$

$$p(\pi | x) = \prod_{t=1}^T y_{\pi_t}^t, \forall \pi \in L^{1:T}, L^{1:T} \in L^{1:*$$

$$d_t(s) \cdot \beta_t(s) \cdot \frac{1}{y_{\pi_t}^t} = \sum_{\pi \in B^{-1}(l)} p(\pi | x)$$

$\pi_t = l_s^t$

1) $p(l | x) = \sum_{\pi \in B^{-1}(l)} p(\pi | x) =$, но мы не можем выделить t из π в виде Σ

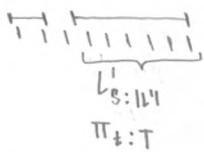
2)
$$p(l | x)_t = \sum_{s=1}^{|L'|} \frac{1}{y_{l_s^t}^t} \cdot d_t(s) \cdot \beta_t(s)$$
 важность каждого $p(l | x)$
но мы не можем $t = \text{const}$, так как $t \neq 1 \neq 2$

$$\sum_{s=1}^{|L'|} \sum_{\pi \in B^{-1}(l)} p(\pi | x), t$$
 $E = O^{ML}(S, \theta_w) = -\sum_{(x, z) \in S} \ln(p(l | x))$

$\pi_t = l_s^t$

Мы отображаем t в π и получаем что не можем выделить t из π ,
 т.е. просто, мы видим $d_t(s) = 0$, видим $\beta_t(s) = 0$,
 соответственно $p(l | x)_t = 0 \dots ?$

Ifem. $d_1(3) = 0 \quad \beta_1(3) > 0$



Для любого t , $p(l | x)_t > 0$, если первый $t = 3-1-n$ где $u \neq g$ для β

π_t

$$\frac{dO^{ML}}{dy_k^t} = \frac{d(-\sum \ln(p(l | x)))}{dy_k^t}$$

$$\nabla E_y^t = \begin{bmatrix} y_1^t \\ y_2^t \\ \vdots \\ y_k^t \\ \vdots \end{bmatrix}$$

for t of 1 word

$$\nabla E_y = \sum_t \nabla E_y^t$$

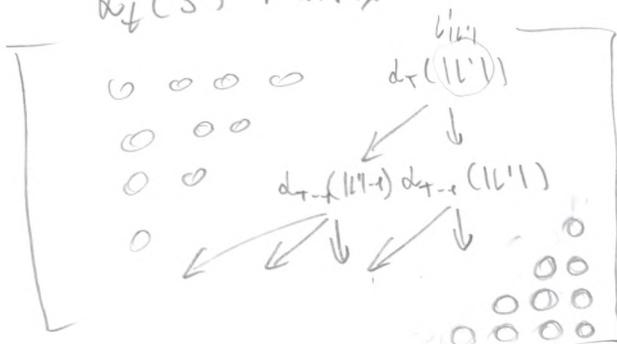
$$T-t \geq (|L'|-s)-1-n?$$

where = $s = \text{символы} - 1$ (- символ?)
 counts noles

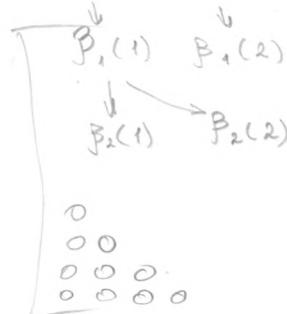


$g-3=8$ - number of words
 $6-1=5$ - blank
 $5-1=4$ - skip

$d_t(s)$ Matrix



$\beta_t(s)$



$$\begin{matrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{matrix}$$

$$\begin{matrix} \cancel{\beta_1(1)} \\ \cancel{\beta_1(2)} \\ \cancel{\beta_2(1)} \\ \cancel{\beta_2(2)} \end{matrix}$$

$$- h - e - l - l - o -$$

\uparrow
 $s=8$

char[]

Рекурсивно берутся векторы $d_t(s)$ и $\beta_t(s)$.

hel - b - o -
h - e l - l o

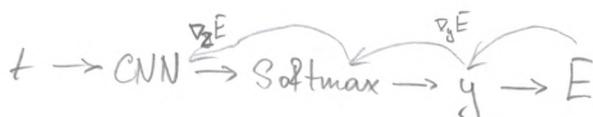
gumbus
 $d_t(4)$
 $d_t(3)$

$$y_k^t = \frac{1}{\sum_{s: l_s^t = k} d_t(s) \beta_t(s)}$$

Как быстро отсеять $s: l_s^t = k$?

Будет некак? можно лучше?

Такая сумма где каждая буква направлена текущей.



- Почему нужны d и β ?
- Вспомогают ли проблемы? 80% да нет, ищите решения RNN.
- В конце нужны levenshtein? или наиболее правдоподобные другие слова типа.
- Потенциальная 'b', 'g' и регуляризация? могут быть распознаваться?

CTC Loss . Forward - Backward algorithm

$$P(\bar{\pi} | x) = \prod_{t=1}^T y_{\pi_t}^t \quad - \text{Likelihood of particular path } \bar{\pi}$$

$\forall \pi \in L^T$

$$h(x = (R^m)^*) = z$$

$$h(\bar{x}_t) = \bar{y}^t, \bar{y}^t \in R^{L'}$$

\bar{y}^t - softmax i.e. $[P(\text{blank}), P(L_1), \dots]^T$

$$P(\bar{l} | x) = \sum_{\bar{\pi} \in B^{-1}(l)} P(\bar{\pi} | x)$$

$$\mathcal{Z} = L^*$$

$$L' = \{\text{blank}\} \cup L, \quad \mathcal{Z} = L'^*$$

$$B: \mathcal{Z} \rightarrow \mathcal{L} \quad - \text{Decoding}$$

$$h(x) = \operatorname{argmax}_{l \in L^T} P(\bar{l} | x) \quad - \text{Classifier (Greedy)}$$

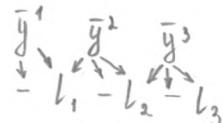
$$d_t(s) = \sum_{\bar{\pi} \in N^T: \bar{\pi}_t = l_s} \prod_{t'=1}^T y_{\pi_{t'}}^{t'} \quad - \text{Сумма независимо от того, сколько символов Sygym генерируются вправо от символа } l_s.$$

$B(\pi_{x,t}) = B(L'_{1:s})$

$$d_1(1) = y_b^1$$

$$d_1(2) = y_{l_1}^1$$

$$d_1(s) = 0, \forall s > 2$$



Однако неправильное не может непрекращаться с символом l' .

$$\bar{d}_t(s) = d_{t-1}(s) + d_{t-1}(s-1) \quad - l_1 - \{ - \quad - l_1 - l_2 \} l_2$$

$$P(\bar{l} | x) = d_T(1|l'|) + d_T(1|l'|-1)$$

$$d_t(s) = \begin{cases} d_{t-1}(s) \cdot y_{l_s}^t & , \text{ if } l_s' = b \text{ or } l_s' = l_{s-1}' \\ (\bar{d}_t(s) + d_{t-1}(s-2)) y_{l_s}^t & , \text{ otherwise} \end{cases} \quad - l_1 - l_2 \} l_2$$

$$\beta_t(s) = \sum_{\bar{\pi} \in N^T: \bar{\pi}_t = l_s} \prod_{t'=1}^T y_{\pi_{t'}}^{t'}, \quad \beta_T(1|l'|) = y_b^T \quad - l_1 \{ l_2 \\ B(\pi_{t:T}) = B(L'_{s+1:T}) \quad \beta_T(1|l'|-1) = y_{l_{s+1}}^T \\ \beta_T(s) = 0, \forall s < |l'| - 1 \quad - l_1 - l_2 \} l_2$$

$$\bar{\beta}_t(s) = \beta_{t+1}(s) + \beta_{t+1}(s+1) \quad - \{ - l - l - 0 - \\ \beta_t(s) = \begin{cases} y_{l_s}^t \bar{\beta}_t(s) & , \text{ if } l_s' = b \text{ or } l_s' = l_{s+2}' \\ y_{l_s}^t (\bar{\beta}_t(s) + \beta_{t+1}(s+2)) & , \text{ otherwise} \end{cases} \quad - \{ l - l - 0 -$$

$$O^{ML}(S, X_w) = - \sum_{(x, z) \in S} \ln(p(z|x))$$

$$d_t(s) \cdot \beta_t(s) = \sum_{\substack{\pi \in B^{-1}(l) : \\ \pi_t = l_s}} y_{l_s}^t \prod_{t'=1}^T y_{\pi_{t'}}^{t'} \quad - H-E-E-L-L-L-O-$$

$$P(l | x) = \sum_{s=1}^{|l'|} \sum_{\substack{\pi \in B^{-1}(l) : \\ \pi_t = l_s}} p(\pi | x) = \sum_{s=1}^{|l'|} \frac{1}{y_{l_s}^T} \cdot d_t(s) \cdot \beta_t(s) \quad P'(l | x) = \sum \frac{1}{y_{l_s}^T} d_t(s) \beta_t(s) \check{y}_{l_s}^t$$

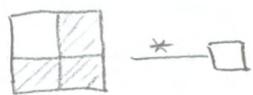
$$\frac{dp(l | x)}{dy_k^t} = \frac{1}{(y_k^t)^2} \cdot \sum_{s \in \text{lab}(l, k)} d_t(s) \beta_t(s)$$

$$\frac{dO^{ML}(S, X_w)}{dy_k^t} = \frac{d(-\ln p(l | x))}{dy_k^t} = -\frac{1}{p(l | x)} \cdot \frac{dp(l | x)}{dy_k^t}$$

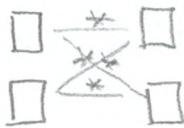
$$\text{lab}(l, k) = \{s : l_s' = k\}$$

most symbol
Set s where

For each word we'll get matrix of derivatives $T \times |l'| \rightarrow [S_1, S_2, \dots]$
 β will be longer than one symbol, one symbol may appear more than once
 $O(T \times S)$



$w/2 \quad h/2$



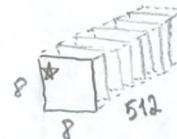
Вся симметрия остается такой же

$$\begin{array}{cccc} 1 & 2 & 3 & 4 \\ 3+2+2+2 & & & \\ 5 & 4 & 9 & \end{array}$$

$$\begin{array}{r} 1 \left(\begin{array}{r} 3 \\ 5 \\ 10 \end{array} \right) \\ 2 \left(\begin{array}{r} 12 \\ 4 \\ 14 \\ 28 \end{array} \right) \\ 3 \left(\begin{array}{r} 30 \\ 8 \\ 32 \\ 64 \end{array} \right) \\ 4 \left(\begin{array}{r} 66 \\ 88 \\ 126 \end{array} \right) \end{array}$$

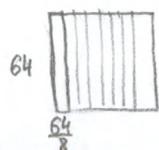
Receptive Field

3	$8 \times 64 \times 64$
5	$64 \times 64 \times 64$
10	
12	$64 \times 32 \times 32$
14	$128 \times 32 \times 32$
28	
30	$128 \times 16 \times 16$
32	$256 \times 16 \times 16$
64	
66	$256 \times 8 \times 8$
68	$512 \times 8 \times 8$

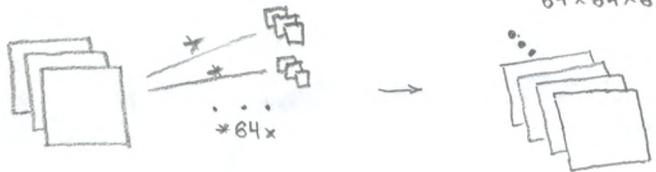


1) GMP $(\frac{w}{8}, \frac{h}{8})$
 $\begin{matrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{matrix} 512$

2) GMP $(1, *)$

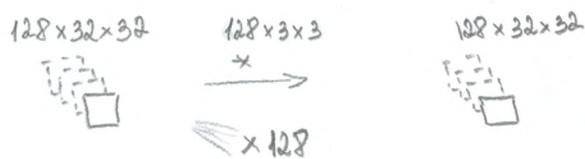
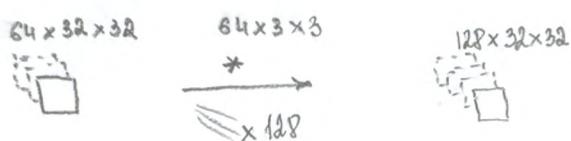
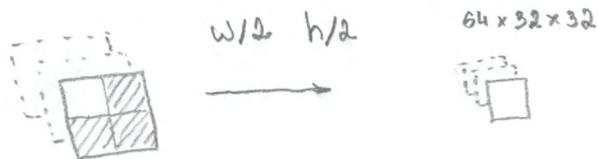


RGB 64×64 + padding

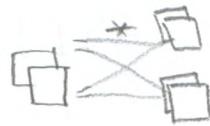


VGG

Max Pool (2×2)



Max Pool



Computational Complexity remains the same

Машындың оғыту 8-жада адалық арамасуның
өзгіліктерін автоматтың түрде талдаудың және
меншік қабылдаудың алгоритмдер мен деяктерді
тәжірибелі жасандық штамшектер көсемшесе.