

Propagacja i propagacja wsteczna

Weronika Hryniewska

Propagacja

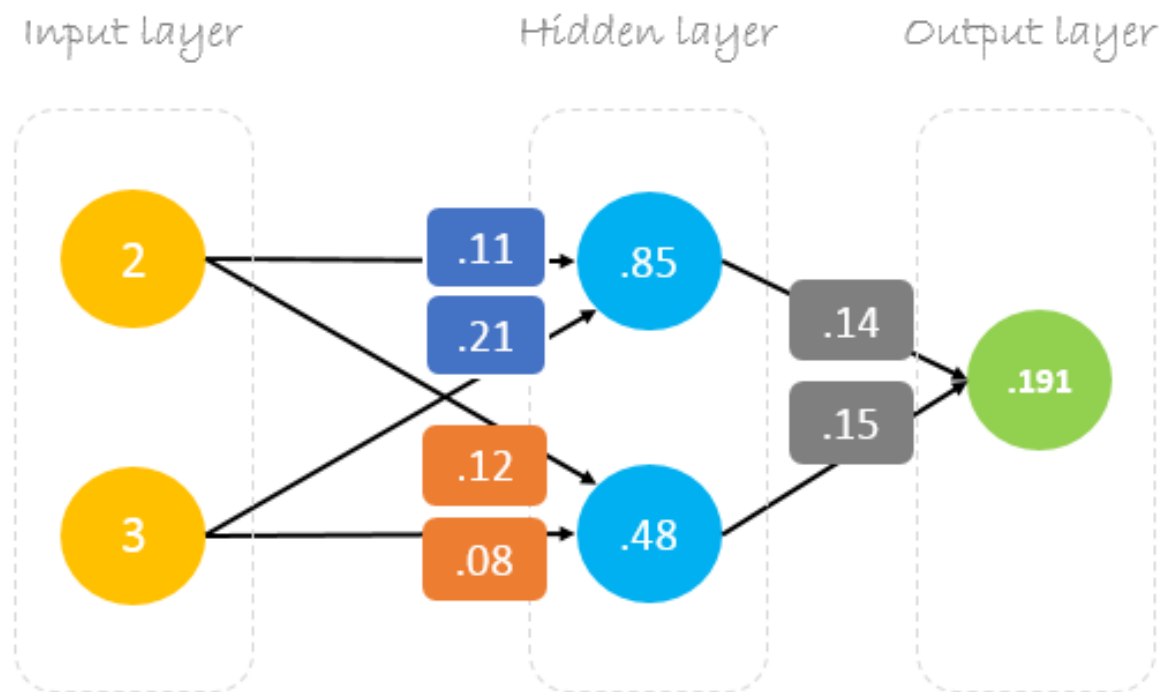
prediction = out

prediction = $(h_1) w_5 + (h_2) w_6$

$$\begin{aligned} h_1 &= i_1 w_1 + i_2 w_2 \\ h_2 &= i_1 w_3 + i_2 w_4 \end{aligned}$$

prediction = $(i_1 w_1 + i_2 w_2) w_5 + (i_1 w_3 + i_2 w_4) w_6$

to change **prediction** value,
we need to change **weights**



Forward Pass

$$\begin{bmatrix} 2 & 3 \end{bmatrix} \cdot \begin{bmatrix} 0.11 & 0.12 \\ 0.21 & 0.08 \end{bmatrix} = \begin{bmatrix} 0.85 & 0.48 \end{bmatrix} \cdot \begin{bmatrix} 0.14 \\ 0.15 \end{bmatrix} = \begin{bmatrix} 0.191 \end{bmatrix}$$

Matrix multiplication

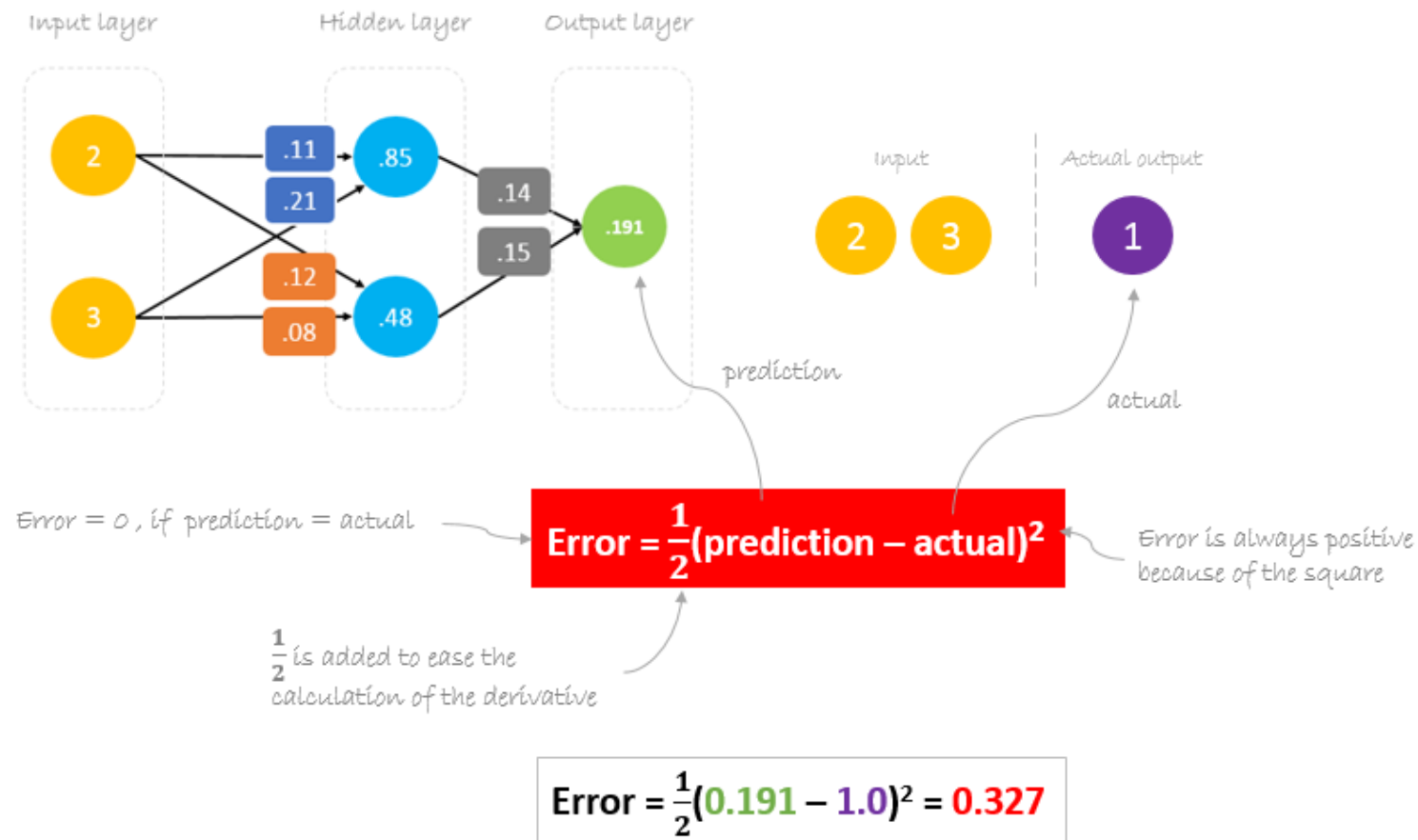
Details

$$2 \times .11 + 3 \times .21 = .85$$

$$.85 \times .14 + .48 \times .15 = .191$$

$$2 \times .12 + 3 \times .08 = .48$$

Obliczanie błędu



Wsteczna propagacja błędu

$$*W_x = W_x - a \left(\frac{\partial \text{Error}}{\partial W_x} \right)$$

Old weight $\rightarrow W_x$
 Derivative of Error with respect to weight $\rightarrow \frac{\partial \text{Error}}{\partial W_x}$
 New weight $\rightarrow *W_x$
 Learning rate $\rightarrow a$

$$\frac{\partial \text{Error}}{\partial W_6} = \frac{\partial \text{Error}}{\partial \text{prediction}} * \frac{\partial \text{prediction}}{\partial W_6} \quad \leftarrow \text{chain rule}$$

$$\text{Error} = \frac{1}{2}(\text{prediction} - \text{actual})^2$$

$$\frac{\partial \text{Error}}{\partial W_6} = \frac{1}{2}(\text{prediction} - \text{actual})^2 * \frac{\partial (i_1 w_1 + i_2 w_2) w_5 + (i_1 w_3 + i_2 w_4) w_6}{\partial W_6}$$

$$\text{prediction} = (i_1 w_1 + i_2 w_2) w_5 + (i_1 w_3 + i_2 w_4) w_6$$

h_2

$$\frac{\partial \text{Error}}{\partial W_6} = 2 * \frac{1}{2}(\text{prediction} - \text{actual}) \frac{\partial (\text{prediction} - \text{actual})}{\partial \text{prediction}} * (i_1 w_3 + i_2 w_4)$$

$$h_2 = i_1 w_3 + i_2 w_4$$

$$\frac{\partial \text{Error}}{\partial W_6} = (\text{prediction} - \text{actual}) * (h_2)$$

$$\Delta = \text{prediction} - \text{actual}$$

delta

$$\frac{\partial \text{Error}}{\partial W_6} = \Delta h_2$$

$$\frac{\partial \text{Error}}{\partial W_1} = \frac{\partial \text{Error}}{\partial \text{prediction}} * \frac{\partial \text{prediction}}{\partial h_1} * \frac{\partial h_1}{\partial W_1} \quad \leftarrow \text{chain rule}$$

$$\text{Error} = \frac{1}{2}(\text{prediction} - \text{actual})^2$$

$$\text{prediction} = (h_1) w_5 + (h_2) w_6$$

$$\frac{\partial \text{Error}}{\partial W_1} = \frac{\partial \frac{1}{2}(\text{prediction} - \text{actual})^2}{\partial \text{prediction}} * \frac{\partial (h_1) w_5 + (h_2) w_6}{\partial h_1} * \frac{\partial i_1 w_1 + i_2 w_2}{\partial w_1}$$

$$h_1 = i_1 w_1 + i_2 w_2$$

$$\frac{\partial \text{Error}}{\partial W_1} = 2 * \frac{1}{2}(\text{prediction} - \text{actual}) \frac{\partial (\text{prediction} - \text{actual})}{\partial \text{prediction}} * (w_5) * (i_1)$$

$$\frac{\partial \text{Error}}{\partial W_1} = (\text{prediction} - \text{actual}) * (w_5 i_1)$$

$$\Delta = \text{prediction} - \text{actual}$$

\leftarrow delta

$$\frac{\partial \text{Error}}{\partial W_1} = \Delta w_5 i_1$$

Updated weights

$$\begin{aligned}
 *w_6 &= w_6 - a (h_2 \cdot \Delta) \\
 *w_5 &= w_5 - a (h_1 \cdot \Delta) \\
 *w_4 &= w_4 - a (i_2 \cdot \Delta w_6) \\
 *w_3 &= w_3 - a (i_1 \cdot \Delta w_6) \\
 *w_2 &= w_2 - a (i_2 \cdot \Delta w_5) \\
 *w_1 &= w_1 - a (i_1 \cdot \Delta w_5)
 \end{aligned}$$

$$\begin{bmatrix} w_5 \\ w_6 \end{bmatrix} = \begin{bmatrix} w_5 \\ w_6 \end{bmatrix} - a \Delta \begin{bmatrix} h_1 \\ h_2 \end{bmatrix} = \begin{bmatrix} w_5 \\ w_6 \end{bmatrix} - \begin{bmatrix} a h_1 \Delta \\ a h_2 \Delta \end{bmatrix}$$

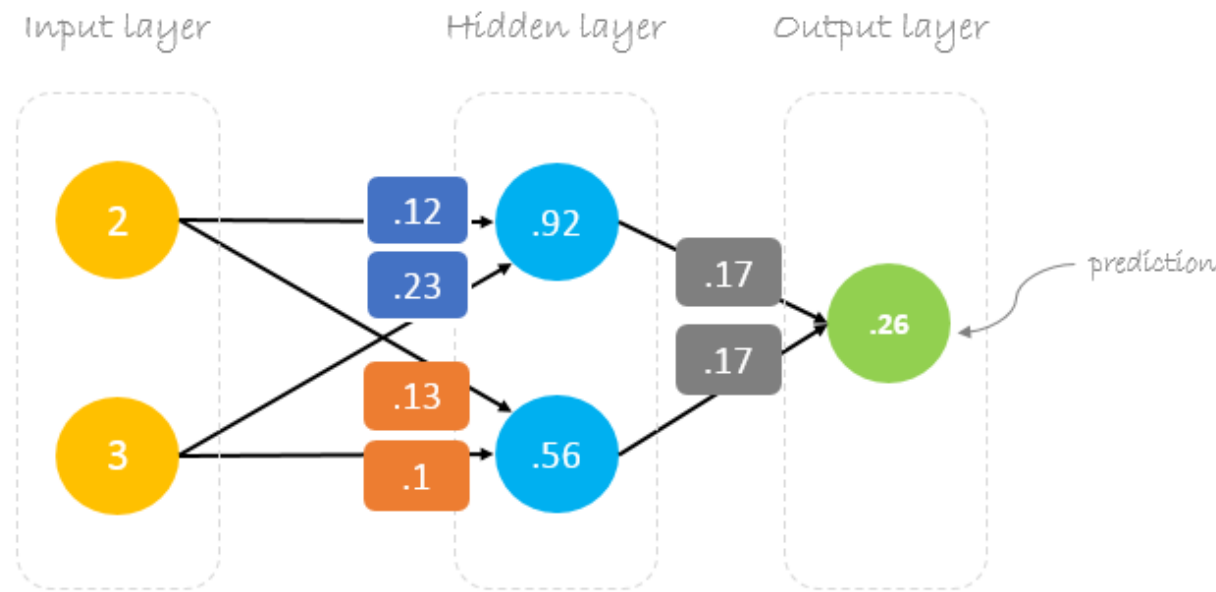
$$\begin{bmatrix} w_1 & w_3 \\ w_2 & w_4 \end{bmatrix} = \begin{bmatrix} w_1 & w_3 \\ w_2 & w_4 \end{bmatrix} - a \Delta \begin{bmatrix} i_1 \\ i_2 \end{bmatrix} \cdot \begin{bmatrix} w_5 & w_6 \end{bmatrix} = \begin{bmatrix} w_1 & w_3 \\ w_2 & w_4 \end{bmatrix} - \begin{bmatrix} a i_1 \Delta w_5 & a i_1 \Delta w_6 \\ a i_2 \Delta w_5 & a i_2 \Delta w_6 \end{bmatrix}$$

$$\Delta = 0.191 - 1 = -0.809 \quad \text{Delta} = \text{prediction} - \text{actual}$$

$$a = 0.05 \quad \text{Learning rate, we smartly guess this number}$$

$$\begin{bmatrix} w_5 \\ w_6 \end{bmatrix} = \begin{bmatrix} 0.14 \\ 0.15 \end{bmatrix} - 0.05(-0.809) \begin{bmatrix} 0.85 \\ 0.48 \end{bmatrix} = \begin{bmatrix} 0.14 \\ 0.15 \end{bmatrix} - \begin{bmatrix} -0.034 \\ -0.019 \end{bmatrix} = \begin{bmatrix} 0.17 \\ 0.17 \end{bmatrix}$$

$$\begin{bmatrix} w_1 & w_3 \\ w_2 & w_4 \end{bmatrix} = \begin{bmatrix} .11 & .12 \\ .21 & .08 \end{bmatrix} - 0.05(-0.809) \begin{bmatrix} 2 \\ 3 \end{bmatrix} \cdot \begin{bmatrix} 0.14 & 0.15 \end{bmatrix} = \begin{bmatrix} .11 & .12 \\ .21 & .08 \end{bmatrix} - \begin{bmatrix} -0.011 & -0.012 \\ -0.017 & -0.018 \end{bmatrix} = \begin{bmatrix} .12 & .13 \\ .23 & .10 \end{bmatrix}$$



Forward Pass

$$\begin{bmatrix} 2 & 3 \end{bmatrix} \cdot \begin{bmatrix} 0.12 & 0.13 \\ 0.23 & 0.10 \end{bmatrix} = \begin{bmatrix} 0.92 & 0.56 \end{bmatrix} \cdot \begin{bmatrix} 0.17 \\ 0.17 \end{bmatrix} = \begin{bmatrix} 0.26 \end{bmatrix}$$

$$2 \times .12 + 3 \times .23 = .85$$

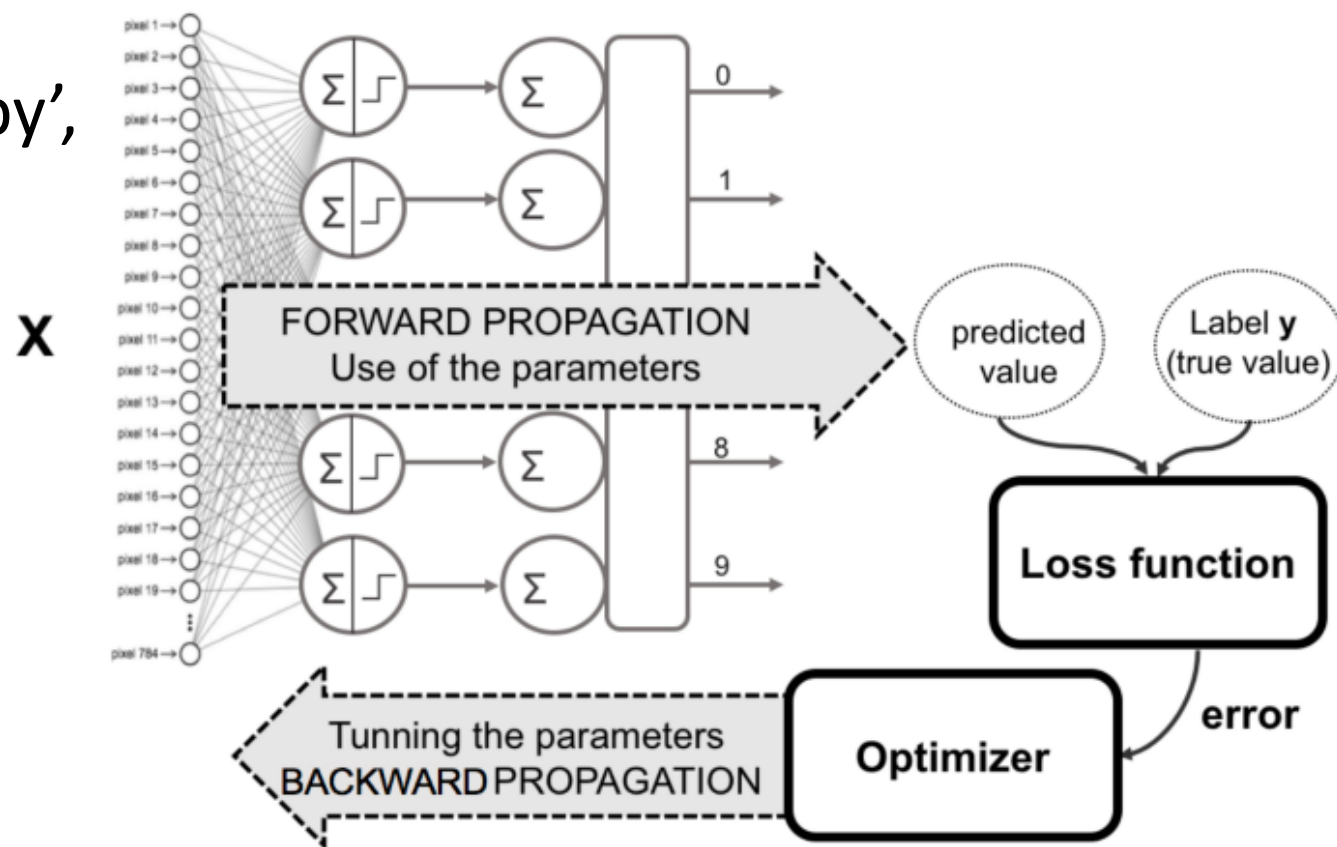
$$.92 \times .17 + .56 \times .17 = .26$$

$$2 \times .13 + 3 \times .10 = .48$$

Komponenty propagacji wstecznej

Weronika Hryniewska


```
model.compile(  
    loss='categorical_crossentropy',  
    optimizer='sgd',  
    metrics=['accuracy']  
)
```



Funkcja straty

- Wyliczana przez ostatnią warstwę sieci
- Stara się minimalizować błąd poprzez zmianę wag neuronów
- Musi być różniczkowalna

Przykłady funkcji straty

- Regression Loss Functions
 - Mean Squared Error Loss
 - Mean Squared Logarithmic Error Loss – wartość docelowa ma szeroki zakres wartości
 - Mean Absolute Error Loss – duże lub małe wartości oddalone od wartości średniej
- Binary Classification Loss Functions
 - Binary Cross-Entropy – preferowana, wartości docelowe w zbiorze $\{0, 1\}$.
 - Hinge Loss – głównie do SVM, wartości docelowe w zbiorze $\{-1, 1\}$.
 - Squared Hinge Loss
- Multi-Class Classification Loss Functions
 - Multi-Class Cross-Entropy Loss – preferowana, wartości docelowe w zbiorze $\{0, 1, \dots, n\}$, gdzie każdej klasie przypisana jest unikalna wartość całkowita
 - Sparse Multiclass Cross-Entropy Loss – klasyfikacja z dużą ilością etykiet
 - Kullback Leibler Divergence Loss - rekonstrukcja oryginalnych danych wejściowych

Optymalizator

- Algorytm aktualizujący wagi
- uwagi mogą być aktualizowane po każdym przykładzie (online) lub po grupie przykładów (offline, batch)

Przykłady optymalizatorów

- SGD
- Adagrad
- RMSprop
- AdaDelta
- Adam (RMSprop + Momentum)
- AdaMax
- Nadam

