**正向传播**



$$d\,loss = d(sum(y-y\_pred)^2) = sum[2(y\_pred-y)\cdot d(y\_pred)] = 2(y\_pred-y)\times d(y\_pred)^T$$

$$= 2(y\_pred-y)\times d(W_2^T)\times h\_relu^T$$

记 $2(y\_pred-y) = \alpha$，$h\_relu = \beta$，则 $\alpha$，$\beta$ 均为行向量.

$$d\,loss = \alpha\times d(W_2^T)\times\beta^T = sum[(\alpha^T\times\beta)\cdot d(W_2^T)] = sum[(\beta^T\times\alpha)\cdot d(W_2)]$$

故 $W_2$ 的梯度即为 $\beta^T\times\alpha = h\_relu^T\times 2(y\_pred-y)$



同理，有

$$d\,loss = 2(y\_pred-y)\times d(y\_pred^T) = 2(y\_pred-y)\times W_2^T\times d(h\_relu^T)$$

$$= 2(y\_pred-y)\times W_2^T\times diag(h>0)\times d(h^T) = 2(y\_pred-y)\times W_2^T\times diag(h>0)\times d(W_1^T)\times x^T$$

同理可得 $W_1$ 的梯度为 $x^T\times 2(y\_pred-y)\times W_2^T\times diag(h>0) = x^T\times[(2(y\_pred-y)\times W_2^T)\cdot(h>0)]$

故令 $W_1 \mathrel{-}= grad\_W_1 \cdot learning\text{-}rate$　　即完成一次梯度下降过程.

$\quad\quad W_2 \mathrel{-}= grad\_W_2 \cdot learning\_rate$