$$a @ b$$

$$\text{np.tensordot}( \ldots ) \iff (a.T @ b).T$$
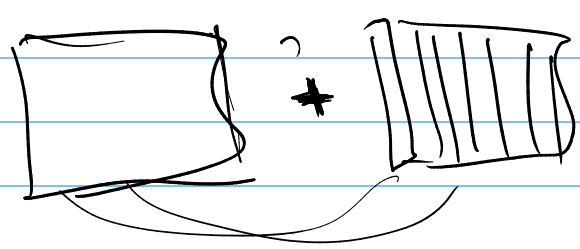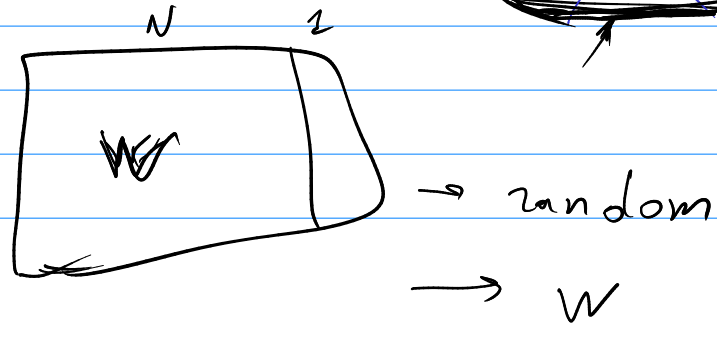
$$(Ax)'_x = A^T = \Sigma$$

np.dot( )

SGD

```
for p in params:
    p.data -= τ p.grad.data
    p.data = p.data - ...
```

Model → [ self.w1, self.bias1, self.w2,...]



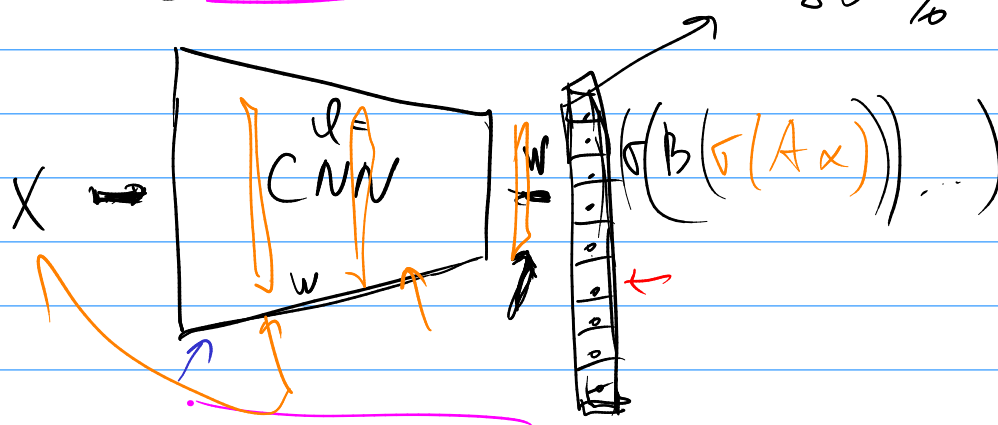$$x_1 W_1 + \ldots + x_N W_N \quad \#1 \cdot W_{N+1} = X \cdot W$$

$$28^2 + 1$$

$$x = (\ldots, 1)$$

N        1

W → random

→ W

Transfer learning

Image Net           $1000 \times 1000 = 1M$
                    80% accuracy

$X \to$ [CNN $W$] $\sigma\left(B\left(\sigma(Ax)\right)\right)...)$

                    2048        $2048 \times 1000$

$logits = \mathcal{U}(X) \cdot W$

5  классов
            $2048 \times 5$

$X \to$ [ $\mathcal{U}(X)$  W ]
                random

[DT] , [RF] , [GB]

# Finetuning

$$\underbrace{\psi(X)}_{} \cdot \underbrace{W}_{}$$

$$l_2 = 10^{-4} \qquad l_2 = 10^{-3}$$

---

## $L2$ -regularization

$$f_w(x) = XW$$

$$LOSS = \mathcal{L}(XW, Y) + \lambda \|W\|_2^2 \longrightarrow \min_w$$

$$\|W\|_2^2 = \sum_{ij} w_{ij}^2 \longrightarrow \min$$

$$\frac{\partial}{\partial x} \begin{vmatrix} 10^6 \cdot X \end{vmatrix}$$

$$f - \text{nnn.}$$

$$L_2 - \text{regularizatin} \iff (early\ stopping)$$

$$(best\text{-}ep - ep) \geq es.$$

$$\frac{\partial \|W\|^2}{\partial W} = \cdots$$

$$_B \boxed{\overset{D}{X}} \ _D \boxed{\overset{E}{W}}$$

$$\frac{\partial x}{\partial x}, \frac{\partial c}{\partial x}$$

$$\frac{\partial}{\partial W}(X W) \qquad \frac{\partial}{\partial W_{\ell m}}\left(\sum_k X_{ik} W_{kj}\right) = \qquad \delta_{ij} = \begin{cases} 1, & i = j \\ 0, & i \neq j \end{cases}$$

$$= \sum_k X_{ik} \frac{\partial W_{kj}}{\partial W_{\ell m}} = \boxed{\sum_k X_{ik} \delta_{k\ell} \delta_{jm}}$$

$$\left(\frac{\partial (XW)_{ij}}{\partial W_{\ell m}}\right).$$