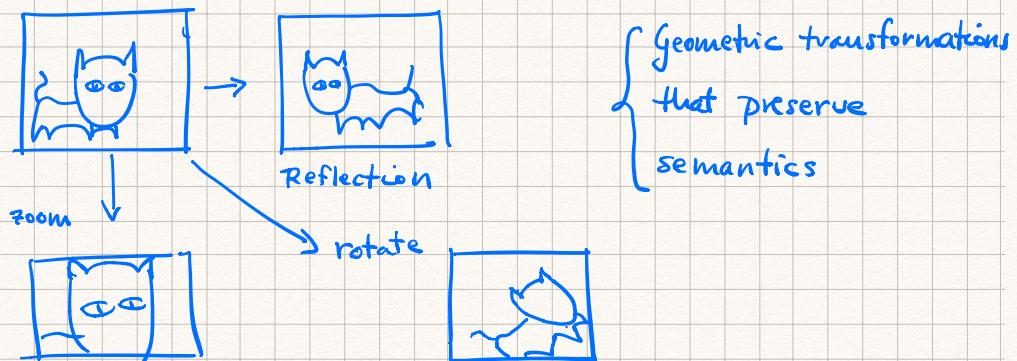


Training neural networks

The dataset

Data augmentation

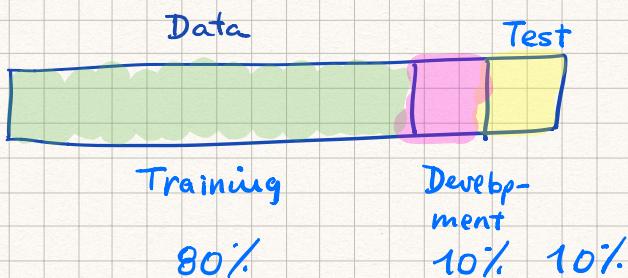
Given a piece of data → replicate



Finding hyperparameters

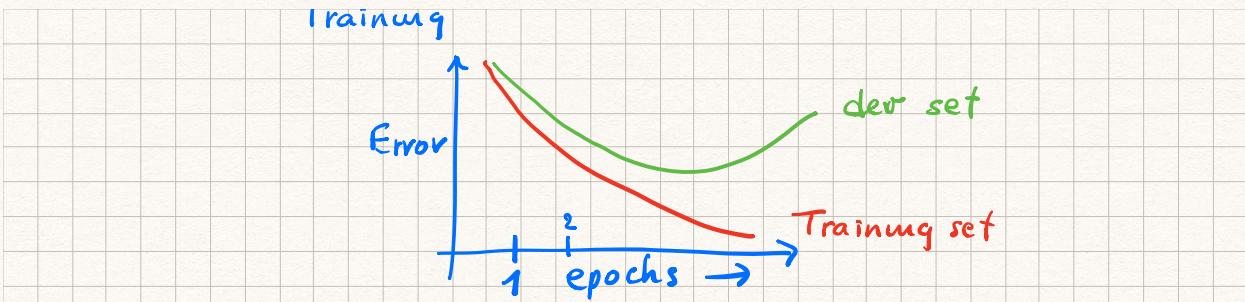
The learning rate δ is a hyperparameter

There are others

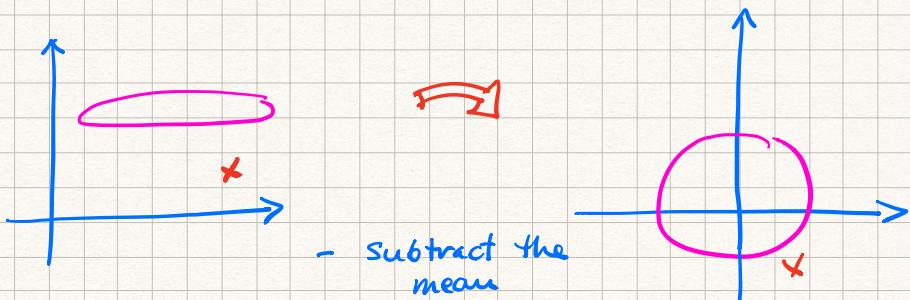


But if we have a lot of data, it can be
98%, 1%, 1%

— . . .



Normalizing the data set



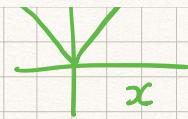
$$\vec{x} \Rightarrow \begin{matrix} 0 \rightarrow \sigma^2 \\ w_0 \quad w_1 \quad w_2 \end{matrix}$$

REGULARIZATION

Like in ridge regression, we add a penalty to the error proportional to the size of the weights

$$m, \times 1 / \|w\|$$

$$E_{\text{new}} = E + \frac{\lambda}{2} \sum_{i=1}^n w_i^2$$



$$\frac{\partial E_{\text{new}}}{\partial w_i} = \frac{\partial E}{\partial w_i} + \lambda w_i$$

Correction

$$\Delta w_i = -\gamma \left(\frac{\partial E}{\partial w_i} + \lambda w_i \right)$$

$\underbrace{\phantom{\Delta w_i = -\gamma (\frac{\partial E}{\partial w_i} + \lambda w_i)}}$ weight decay

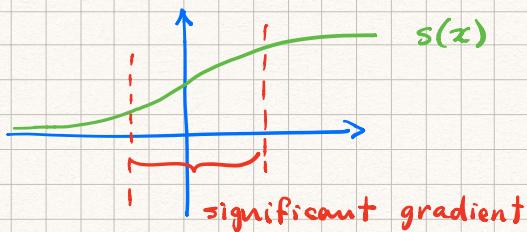
{ the adjustment tries to
 { "pull down" the absolute size of the weight

We can use a regularization parameter λ

$$\Delta w_i = -\gamma \frac{\partial E}{\partial w_i} - \lambda w_i$$

↑ learning rate ↑ weight decay constant

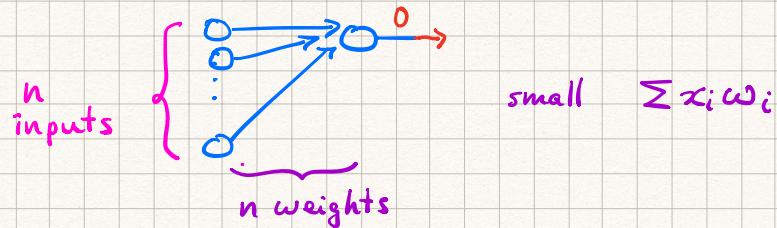
Small weights are better



Initialization

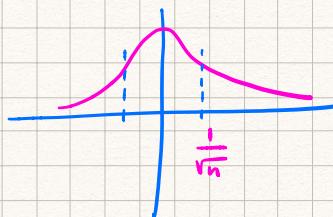
→ Random weights

→ Small weights in terms of the inner product



Xavier initialization

Pick the weights from a Gaussian distribution with zero mean and variance $1/n$



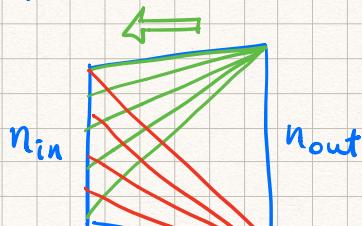
$$\text{Then } \text{Var}(\sum w_i x_i) = n \cdot \text{Var}(w_i) \cdot \text{Var}(x_i)$$

$$= \text{var}(x_i)$$

$$\text{Var}(w_i) = \frac{1}{n}$$

Glorot & Bengio

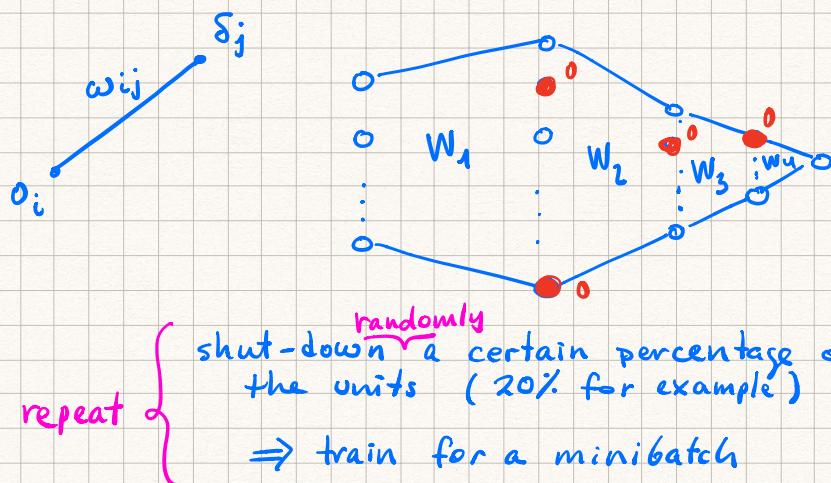
Consider backprop



$$\text{Var}(w_i) = \frac{2}{n_{\text{in}} + n_{\text{out}}}$$

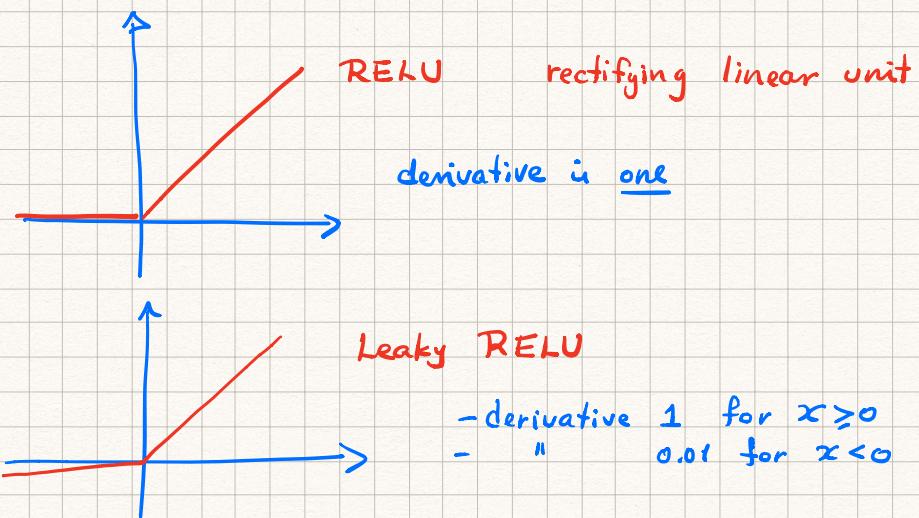
Initialize in
this layer so

Dropout regularization



- The network learns to be robust
- The weights are "spread" through the network

Activation functions



Backprop with momentum



if a ball is rolling
down, it gains
"momentum"

In gradient descent we can "remember" the previous correction in step $i-1$ and then:

$$\Delta \vec{w}^{(i)} = \underbrace{\beta \Delta w^{(i-1)}}_{\text{previous correction}} - (1-\beta) \underbrace{\vec{\nabla} E}_{\text{gradient}}$$

β is a new hyperparameter

for example $\beta = 0.9$

The result is an exponentially weighted average

$$\begin{aligned}\Delta \vec{w}^{(i)} &= \beta \Delta w^{(i-1)} + (1-\beta) \vec{\nabla} E^{(i)} \\ &= \beta^2 \Delta w^{(i-2)} + \beta (1-\beta) \vec{\nabla} E^{(i-1)} + (1-\beta)^2 \vec{\nabla} E^{(i)}\end{aligned}$$

Learning rate decay

Start with γ_0

Let γ decay

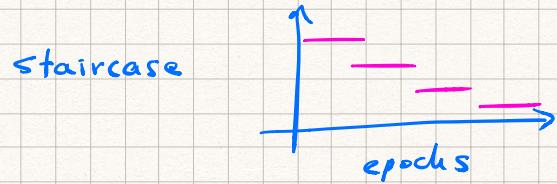
$$\gamma = \frac{\gamma_0}{1 + \underbrace{\text{decay-rate} * \text{epoch}}_{\text{one pass through data}}}$$

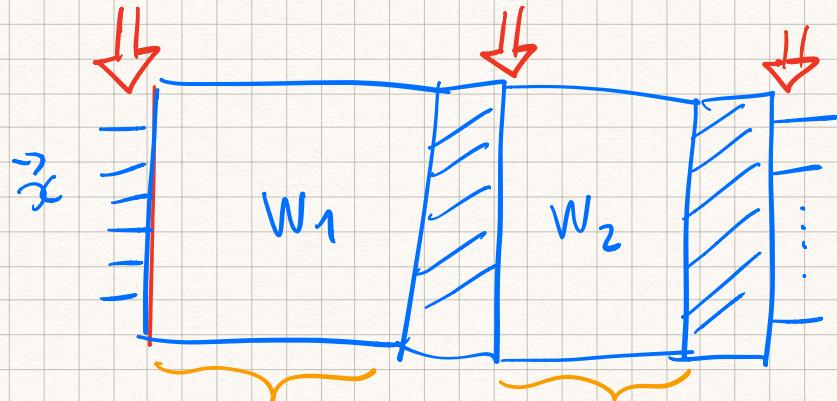
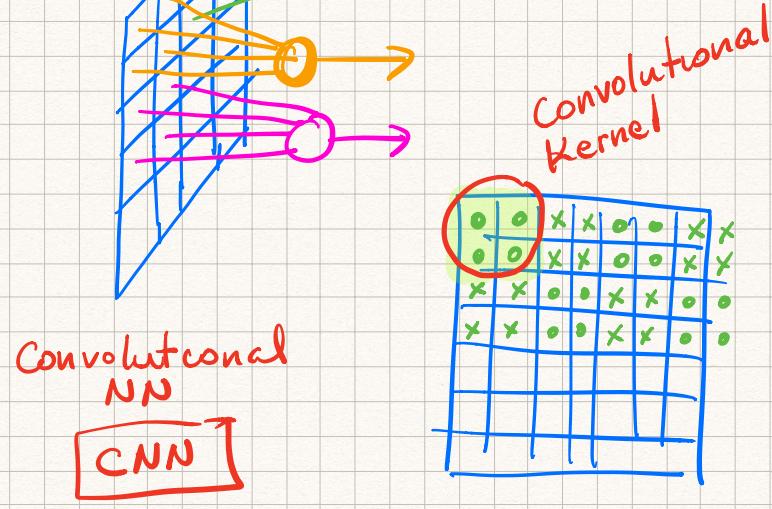
$$\gamma = \gamma_0 \cdot 0.95^{\text{epoch}}$$

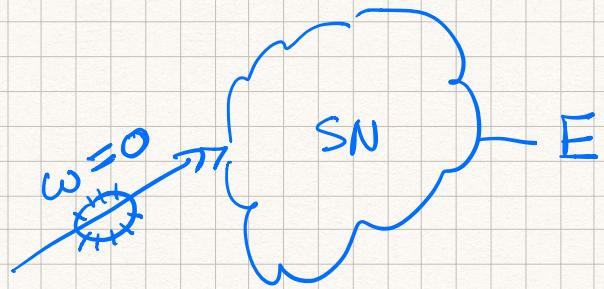
$$\gamma = \frac{k}{\sqrt{\text{epoch}}} \cdot \gamma_0$$

manual decay

staircase







$$\Delta w = -\gamma \frac{\partial E}{\partial w}$$

$$\frac{\partial E}{\partial w}$$

