

# Transformed data

## Lecture 20

Consider a transform with a single parameter:  $\vec{s}(\vec{x}, \xi)$  [ $\vec{s}(\vec{x}, 0) = \vec{x}$ ]

In the infinite data limit, the sum-of-squares error function is given by:

$$E = \frac{1}{2} \iint d\vec{x} dt (y(\vec{x}) - t)^2 p(t|\vec{x}) p(\vec{x})$$

1D output (single output node) for simplicity

Imagine that each  $\vec{x}$  is perturbed many times:  $\vec{x} \rightarrow \vec{s}(\vec{x}, \xi)$ , where  $\xi$  is drawn from  $p(\xi)$ .

Then

$$\tilde{E} = \frac{1}{2} \iiint d\vec{x} dt d\xi (y(\vec{s}(\vec{x}, \xi)) - t)^2 p(t|\vec{x}) p(\vec{x}) p(\xi)$$

↑  
over expanded data set

Now, assume that  $\int d\xi \xi p(\xi) = E(\xi) = 0$ ,

$$E(\xi^2) = \int d\xi \xi^2 p(\xi) = \lambda \ll \text{small variance}$$

s.t. we only consider "small" transformations of  $\vec{x}$ .

Then

$$\begin{aligned} \bar{S}(\vec{x}, \xi) &= \bar{S}(\vec{x}, 0) + \xi \underbrace{\frac{\partial}{\partial \xi} \bar{S}(\vec{x}, \xi)}_{\bar{\tau}} \Big|_{\xi=0} + \\ &+ \frac{\xi^2}{2} \underbrace{\frac{\partial^2}{\partial \xi^2} \bar{S}(\vec{x}, \xi)}_{\bar{\tau}'} \Big|_{\xi=0} + \mathcal{O}(\xi^3) \end{aligned}$$

Next,

$$\begin{aligned} \bar{y}(\bar{S}(\vec{x}, \xi)) &= \bar{y}(\vec{x} + \xi \bar{\tau} + \frac{\xi^2}{2} \bar{\tau}') = \\ &\approx \bar{y}(\vec{x}) + \xi \bar{\tau}_i \frac{\partial \bar{y}(\vec{x})}{\partial x_i} + \frac{\xi^2}{2} \bar{\tau}'_i \frac{\partial \bar{y}(\vec{x})}{\partial x_i} + \\ &+ \frac{\xi^2}{2} \bar{\tau}_i \bar{\tau}_j \frac{\partial^2 \bar{y}(\vec{x})}{\partial x_i \partial x_j} \end{aligned}$$

sums over  $i, j$  implied

Then

$$\begin{aligned} \bar{E} &= \frac{1}{2} \iiint d\vec{x} dt d\xi p(t|\vec{x}) p(\vec{x}) p(\xi) \times \\ &\times \left[ \bar{y}(\vec{x}) + \xi \bar{\tau}_i \frac{\partial \bar{y}}{\partial x_i} + \frac{\xi^2}{2} \bar{\tau}'_i \frac{\partial \bar{y}}{\partial x_i} + \frac{\xi^2}{2} \bar{\tau}_i \bar{\tau}_j \frac{\partial^2 \bar{y}}{\partial x_i \partial x_j} - \right. \\ &\quad \left. - t \right]^2 = \\ &\approx \frac{1}{2} \iiint d\vec{x} dt \underbrace{[ \bar{y}(\vec{x}) - t ]^2}_{E} p(t|\vec{x}) p(\vec{x}) \quad \text{⊕} \end{aligned}$$

↙  $\int d\xi p(\xi) = 1$

$$\oplus \underbrace{E(\xi)}_{\text{"0"}} \frac{1}{2} \int \int d\vec{x} dt \dots + \underbrace{E(\xi^2)}_{\lambda} \frac{1}{2} \int d\vec{x} dt p(t|\vec{x}) p(\vec{x}) \times$$

$$\times \left[ (y(\vec{x}) - t) \left[ \tau_i \frac{\partial y}{\partial x_i} + \tau_j \tau_j \frac{\partial^2 y}{\partial x_i \partial x_j} \right] + \tau_i \frac{\partial y}{\partial x_i} \tau_j \frac{\partial y}{\partial x_j} \right]$$

So,  $\tilde{E} = E + \lambda \Omega$ , where

$$\Omega = \frac{1}{2} \int d\vec{x} p(\vec{x}) \left[ (y(\vec{x}) - \underbrace{E[t|\vec{x}]}_{\int dt p(t|\vec{x}) t}) \left[ \dots \right] + \tau_i \frac{\partial y}{\partial x_i} \tau_j \frac{\partial y}{\partial x_j} \right]$$

↑  $\int dt p(t|\vec{x}) = 1$

Now, note that  $y(\vec{x}) = E[t|\vec{x}]$  minimizes  $E$  and "almost" minimizes

$$\tilde{E}: y(\vec{x}) = E[t|\vec{x}] + \theta(\xi) \quad \text{minimizes } \tilde{E} \text{ in fact}$$

Thus, the 1st term in  $\Omega$  is  $\theta(\xi)$  while the 2nd is  $\theta(\xi^0)$ , so that

$$\Omega \approx \frac{1}{2} \int d\vec{x} p(\vec{x}) \left( \underbrace{\tau_i \frac{\partial y}{\partial x_i}}_{\text{1D Jacobian}} \right)^2 \ll \text{same as before!}$$

Finally, in a special case

$\vec{x} \rightarrow \vec{x} + \vec{\xi}$  we obtain:

↑  
random translations :  $p(\vec{\xi}) = \prod_i p(\xi_i)$

$$\vec{S}(\vec{x}, \vec{\xi}) = \vec{x} + \vec{\xi} \Rightarrow \frac{\partial S_k}{\partial \xi_i} = \delta_{ik} \quad \text{not really needed}$$

Then  $y(\vec{S}) = y(\vec{x}) + \xi_i \nabla_i y(\vec{x}) + \frac{1}{2} \xi_i \xi_j \nabla_i \nabla_j y(\vec{x}) + \dots$

"  $\vec{x} + \vec{\xi}$  "  $\frac{\partial}{\partial x_i}$  "

Then  $E = \frac{1}{2} \int \int \int d\vec{x} dt d\vec{\xi} p(t|\vec{x}) p(\vec{x}) p(\vec{\xi}) \times$

$\times [y(\vec{x}) + \xi_i \nabla_i y(\vec{x}) + \frac{1}{2} \xi_i \xi_j \nabla_i \nabla_j y(\vec{x}) - t]^2 =$

$= E + \underbrace{E(\xi_i)}_{=0} \frac{1}{2} \int \int d\vec{x} dt \dots \quad \oplus$

$\int d\xi_i \xi_i p(\xi_i)$

$\oplus \frac{1}{2} E(\xi_i \xi_j) \int \int d\vec{x} dt p(t|\vec{x}) p(\vec{x}) \times$

$\times [ (y(\vec{x}) - t) \nabla_i \nabla_j y(\vec{x}) + \nabla_i y(\vec{x}) \nabla_j y(\vec{x}) ] =$

$= E + \lambda \Omega$

~~$\int d\xi_i \xi_i p(\xi_i) \int d\xi_j \xi_j p(\xi_j) \dots$~~

$$E(x_i x_j) = \begin{cases} \int \int dx_i dx_j x_i x_j p(x_i) p(x_j) = 0, & i \neq j \\ \int dx_i x_i^2 p(x_i) \equiv \lambda, & i = j \end{cases}$$

$$\text{Here, } \Omega = \frac{1}{2} \int d\vec{x} p(\vec{x}) \left[ \underbrace{(y(\vec{x}) - E[t|\vec{x}])}_{\mathcal{O}(\|\vec{x}\|), \text{ discard}} \underbrace{\nabla^2 y(\vec{x}) + \sum_i \nabla_i^2}_{\text{discard}} y(\vec{x}) + \underbrace{\nabla_i y(\vec{x}) \nabla_i y(\vec{x})}_{\|\nabla y\|^2} \right] =$$

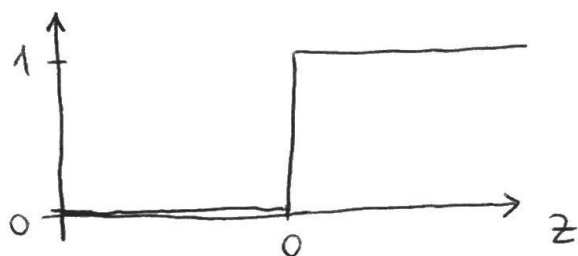
$$= \frac{1}{2} \int d\vec{x} p(\vec{x}) \|\nabla y(\vec{x})\|^2.$$

Tikhonov regularization

# Non-linear activation functions

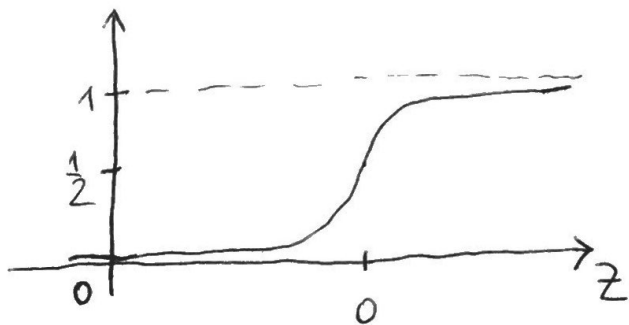
"Classical" activation functions are prone to saturation on the tails, where the gradients become small or vanish completely:

## Perceptron



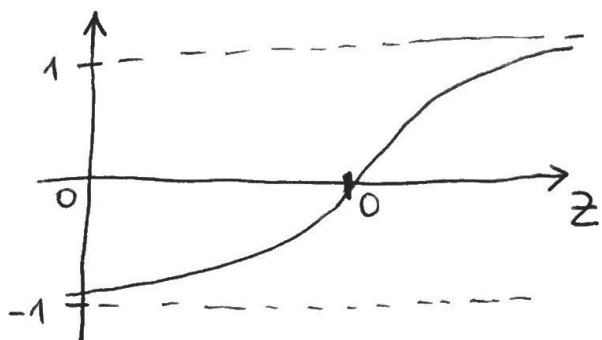
$$\theta(z)$$

## Sigmoid



$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

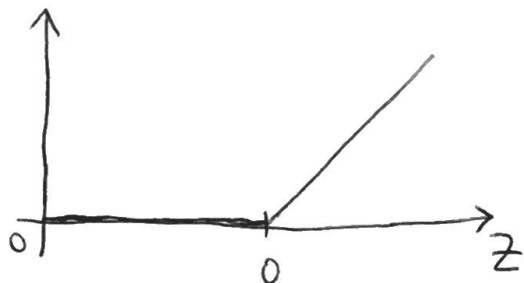
## Tanh



$$\tanh(z)$$

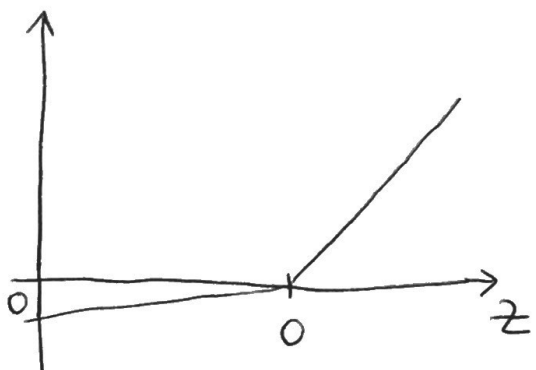
Gradient saturation is widely believed to reduce convergence and/or predictive power. Therefore, other activation functions have been proposed:  
(adopted)

ReLU



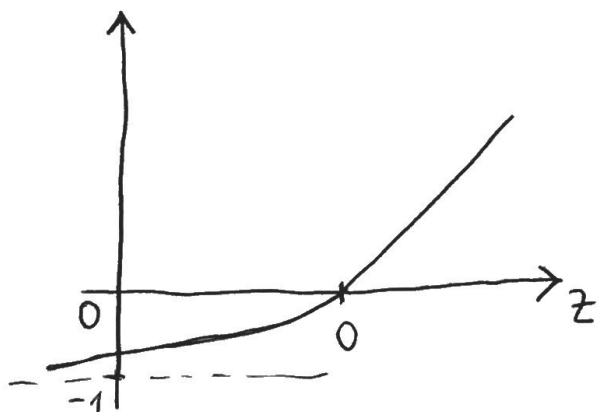
$$\max(0, z)$$

Leaky ReLU



$$\begin{cases} z, z \geq 0 \\ 0.1z, z < 0 \end{cases}$$

ELU



$$\begin{cases} z, z \geq 0 \\ e^z - 1, z < 0 \end{cases}$$