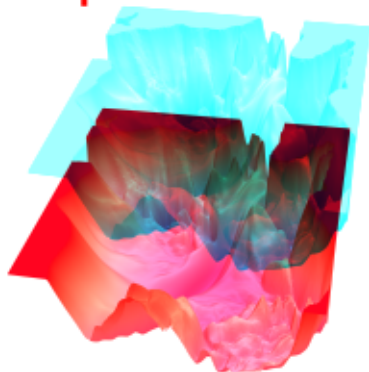


# Introduction to Machine Learning

## Advanced Risk Minimization Properties of Loss Functions



### Learning goals

- Statistical properties

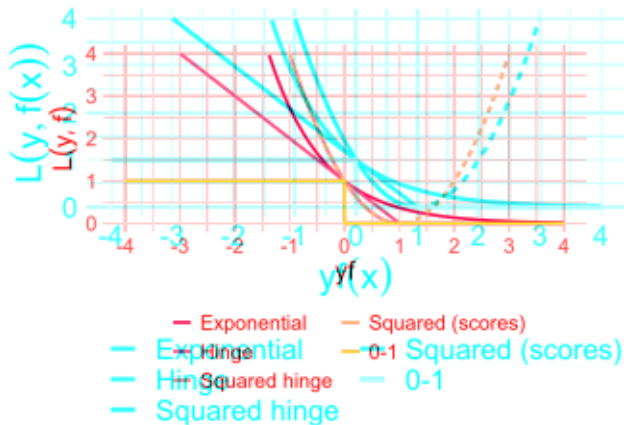
### Learning goals

- Robustness
- Numerical properties
- Some fundamental terminology
- Statistical properties
- Robustness
- Numerical properties
- Some fundamental terminology

## SOME BASIC TERMINOLOGY

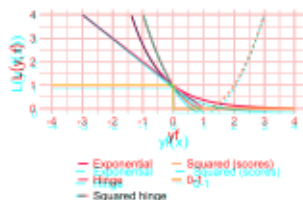
Classification losses are usually expressed in terms of the **margin**:

$$\nu := y \cdot f(\mathbf{x}).$$



# NUMERICAL PROPERTIES: SMOOTHNESS

- **Smoothness** of a function is a property measured by the number of continuous derivatives.
- Derivative-based optimization requires smoothness of the risk  $\mathcal{R}_{\text{emp}}(\theta)$ 
  - If loss is unsmooth, we might have to use derivative-free optimization (or worse, in case of 0-1)
  - Smoothness of  $\mathcal{R}_{\text{emp}}(\theta)$  not only depends on  $L$ , but also requires smoothness of  $f(\mathbf{x})$ !



Squared loss, exponential loss and squared hinge loss are continuously differentiable.  
Hinge loss is continuous but not differentiable.  
0-1 loss is not even continuous.

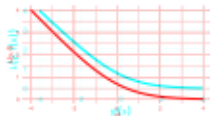


# NUMERICAL PROPERTIES: CONVERGENCE

In case of **complete separation**, optimization might even fail entirely, e.g.:

- Margin-based loss that is strictly monotonically decreasing in  $y \cdot f$ , e.g., **Bernoulli loss**:

$$L(y, f(\mathbf{x})) = \log(1 + \exp(-yf(\mathbf{x})))$$



- $f$  linear in  $\theta$ , e.g., **logistic regression** with  $f(\mathbf{x} | \theta) = \theta^T \mathbf{x}$
- Data perfectly separable by our learner, so we can find  $\theta$ :

$$y^{(i)} f(\mathbf{x}^{(i)} | \theta) = y^{(i)} \theta^T \mathbf{x}^{(i)} > 0 \quad \forall \mathbf{x}^{(i)}$$

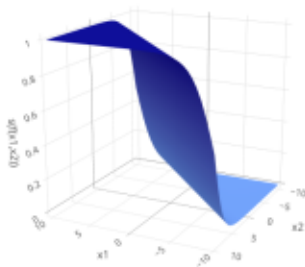
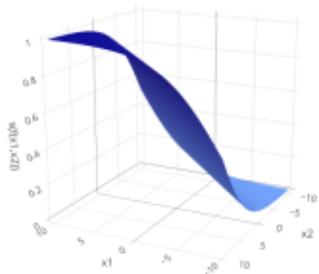
- Can now construct a strictly better  $\theta$

$$\mathcal{R}_{\text{emp}}(2 \cdot \theta) = \sum_{i=1}^n L(2y^{(i)} \theta^T \mathbf{x}^{(i)}) < \mathcal{R}_{\text{emp}}(\theta)$$

- As  $\|\theta\|$  increases, sum strictly decreases, as argument of  $L$  is strictly larger
- We can iterate that, so there is no local (or global) optimum, and no numerical procedure can converge

## NUMERICAL PROPERTIES: CONVERGENCE / 2

- Geometrically, this translates to an ever steeper slope of the logistic/softmax function, i.e., increasingly sharp discrimination:



- In practice, data are seldomly linearly separable and misclassified examples act as counterweights to increasing parameter values.
- Besides, we can use **regularization** to encourage convergence to robust solutions.