

Ridge regression

$$f(x_i) = \beta_0 + \beta_1 x_{i1} + \dots + \beta_m x_{im}, \quad i=1, \dots, n$$

minimise the loss function

$$L_{\beta} = \sum_{i=1}^n \left(y_i - (\beta_0 + \beta_1 x_{i1} + \dots + \beta_m x_{im}) \right)^2$$

$$+ \lambda (\beta_1^2 + \dots + \beta_m^2)$$

a penalty for large coefficients

$$\begin{bmatrix} x_{11} & \dots & x_{1m} \\ \vdots & & \vdots \\ x_{n1} & \dots & x_{nm} \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \rightsquigarrow f(z) = \beta_0 + \beta_1 z_1 + \dots + \beta_m z_m$$

$$z \in \mathbb{R}^m$$

$\lambda \geq 0$ - hyperparameter

Lasso:

$$L_{\beta} = \sum_{i=1}^n \left(-||- \right)^2 + \lambda (|\beta_1| + \dots + |\beta_m|)$$

Lasso has an advantage that it often gives $\beta_j = 0$ for some j 's
interpretability of the model is better

For the usual regression the relative size of the features is irrelevant
(e.g., $x_{.1}$ could be $\in [1, 2)$, $x_{.2} \in [100, 200]$, ...) while for Ridge and Lasso
it is important and usually it is desirable that the features are of the same order of magnitude

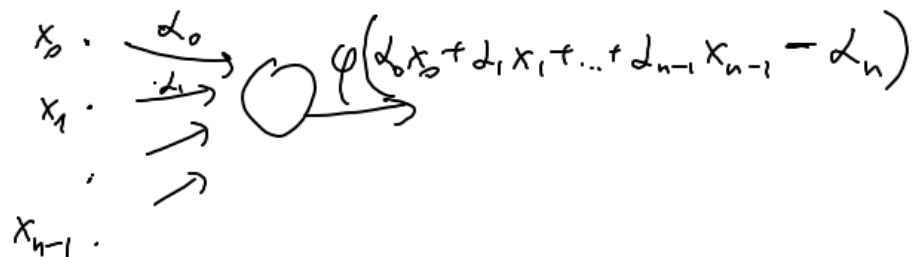
Fix: preprocess the data, for example so that ~~the~~

e.g. each column of X has mean zero and variance 1

or each column of X attains values between 0 and 1

Neural nets

neuron: a function of the form



$$\mathbb{R}^n \ni (x_0, x_1, \dots, x_{n-1}) \xrightarrow{\text{neuron}} \varphi(d_0 x_0 + d_1 x_1 + \dots + d_{n-1} x_{n-1} - d_n) \in \mathbb{R}$$

d_0, \dots, d_n - weights $\in \mathbb{R}$,

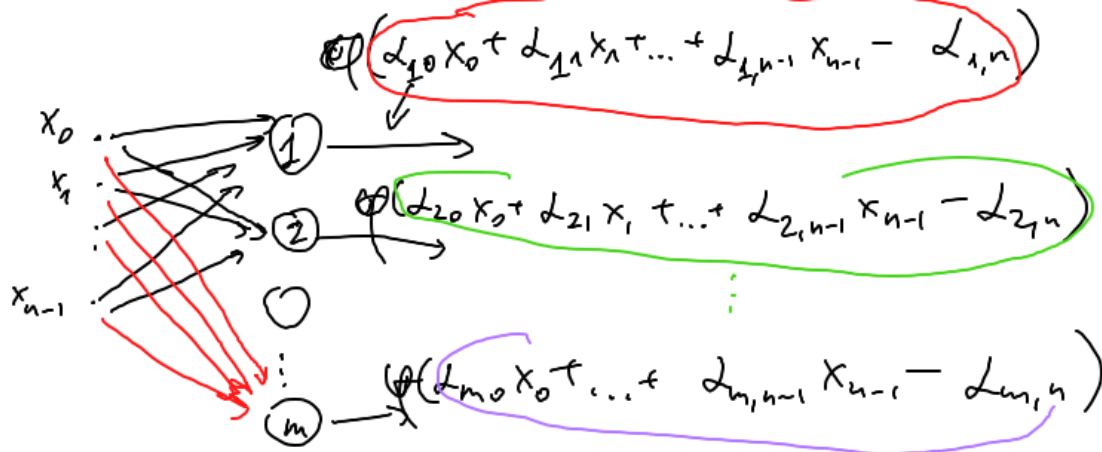
φ - activation function, $\varphi: \mathbb{R} \rightarrow \mathbb{R}$ increasing

e.g. $\varphi(x) = \frac{1}{1+e^{-x}} = \sigma(x)$ sigmoid



$\varphi(x) = \max(x, 0)$ ReLU





$$\mathbb{R}^n \rightarrow (x_0, x_1, \dots, x_{n-1})$$

(dense)
The layer of nodes is a function.

$\mathbb{R}^n \rightarrow \mathbb{R}^m$ given by this formula

$$\phi \left(\begin{bmatrix} d_{1,0} & d_{1,1} & \dots & d_{1,n-1} & d_{1,n} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ d_{m,0} & \dots & d_{m,n-1} & d_{m,n} \end{bmatrix} \begin{bmatrix} x_0 \\ \vdots \\ x_{n-1} \\ -1 \end{bmatrix} \right)$$

$$= \phi \left(\begin{bmatrix} \text{red oval} \\ \text{green oval} \\ \vdots \\ \text{purple oval} \end{bmatrix} \right) = \begin{bmatrix} \phi(\text{red oval}) \\ \phi(\text{green oval}) \\ \vdots \\ \phi(\text{purple oval}) \end{bmatrix} \in \mathbb{R}^m$$

Convention

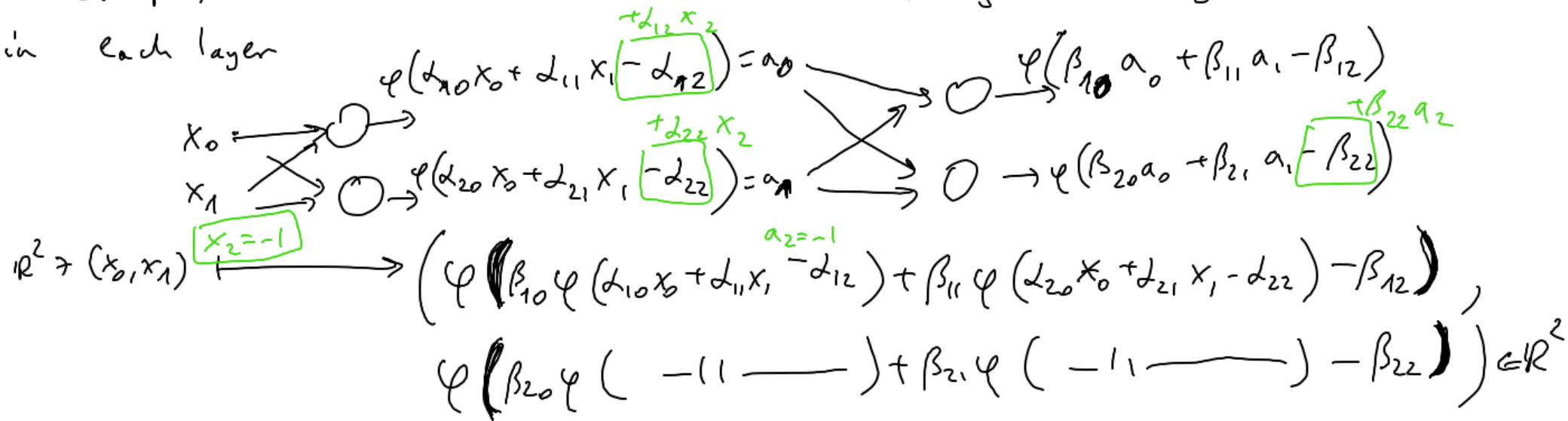
The neural net is a ~~map~~ function which is a composition of layers

$$\mathbb{R}^n \ni (x_0, \dots, x_{n-1}) \xrightarrow{\text{layer}} \begin{bmatrix} \varphi(d_{1,0}x_0 + \dots + d_{1,n-1}x_{n-1} - d_{1,n}) \\ \vdots \\ \varphi(d_{m,0}x_0 + \dots + d_{m,n-1}x_{n-1} - d_{m,n}) \end{bmatrix} \xrightarrow{\text{layer}} \begin{matrix} \in \mathbb{R}^p \\ \downarrow \text{layer} \\ \vdots \end{matrix}$$

\uparrow
 \mathbb{R}^m

~~False~~

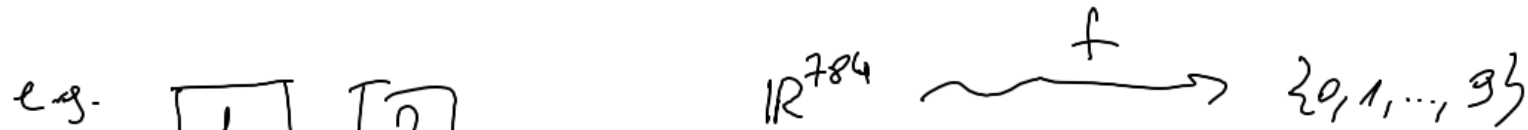
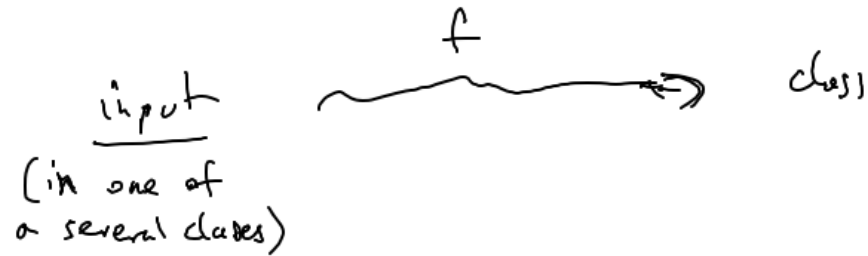
For example, let us consider a neural net consisting of 2 layers with 2 neurons in each layer



If φ is ~~linear~~ affine, then the compositions above are just ~~linear~~ affine functions of (x_0, x_1)

Therefore we do not want to take an affine function as φ .

A common application of neural nets: classification problem



input = an image of a handwritten digit
image \approx vector in $\mathbb{R}^{28 \times 28} = \mathbb{R}^{784}$

one takes a neural net: $\mathbb{R}^{784} \xrightarrow{f} \mathbb{R}^{10}$

the interpretation of the output: $\operatorname{argmax}(f(x)) = \text{class}$

$(0.01, 0.9, 0.1, 0.05, 0.2, 0.1, 0.07, 0.25, 0.1, 0.1) \in \mathbb{R}^{10}$

0 1 2 3 4 5 6 7 8 9

\rightarrow predicted class = 1
 $(0, 1, 0, \dots, 0) \rightarrow$ target

The training procedure:

• take some input $x \in \mathbb{R}^{784}$ and feed it to the network:

we obtain some $y = f(x) \in \mathbb{R}^{10}$

for that input, we know the label, and so we know the target vector:

$$t = (0, 1, 0, \dots, 0) \in \mathbb{R}^{10}$$

↑
this one corresponds to a digit '1'

we introduce some loss function, e.g.

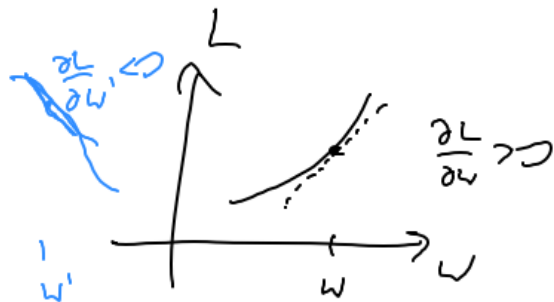
$$L(y, t) = \frac{1}{2} \sum_{i=1}^{10} (y_i - t_i)^2$$

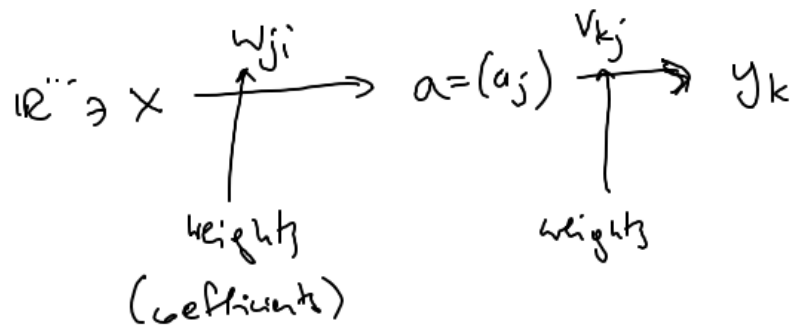
compute $\frac{\partial L}{\partial w}$ for each coefficient w

and adjust the coefficients:

$$w := w - c \frac{\partial L}{\partial w}$$

↑
learning const., e.g. 0.01





2 layers,
 (w_{ji}) - weights in the 1st layer
 (v_{kj}) - " " " " 2nd " " "

$$a_j = \varphi \left(\sum_i w_{ji} x_i \right) \quad \text{(with the convention that } x \text{ has a "1" as the last coordinate)}$$

$$y_k = \varphi \left(\sum_j v_{kj} a_j \right) \quad \text{(with a similar convention for } a \text{)}$$

(e.g. $k=1,2, j=0,1,2, i=0,1,2$ for 2 neurons in each layer)

$x_0 = 0 \rightarrow y_1$
 $x_1 = 0 \rightarrow y_2$
 $x_2 = -1$

$$= \frac{\partial}{\partial v} \left(\sum_j v_{kj} a_j \right)$$

target: (t_k) , $L(t, y) = \frac{1}{2} \sum_k (t_k - y_k)^2 = \frac{1}{2} \sum_k \left(t_k - \varphi \left(\sum_j v_{kj} a_j \right) \right)^2$

$$v = v_{k_0 j_0} : \frac{\partial L}{\partial v} = \sum_k (t_k - \varphi(\sum_j v_{kj} a_j)) \cdot \frac{\partial}{\partial v} \left(t_k - \varphi(\sum_j v_{kj} a_j) \right) = \underbrace{\left(t_{k_0} - \varphi(\sum_j v_{k_0 j} a_j) \right)}_{\frac{\partial L}{\partial y_{k_0}}} \cdot (-1) \cdot \varphi'(\sum_j v_{k_0 j} a_j) \cdot a_{j_0}$$

$$a_j = \varphi\left(\sum_i w_{ji} x_i\right)$$

$$y_k = \varphi\left(\sum_j v_{kj} a_j\right)$$

$$L(t, y) = \frac{1}{2} \sum_k (t_k - y_k)^2$$

$$w = w_{j_0, i_0}$$

$$\frac{\partial L}{\partial w} = \sum_k (t_k - y_k) (-1) \cdot \varphi'\left(\sum_j v_{kj} a_j\right) \cdot \frac{\partial}{\partial w} \left(\sum_j v_{kj} a_j\right) =$$

$$= \sum_k (t_k - y_k) (-1) \varphi'\left(\sum_j v_{kj} a_j\right) \cdot \sum_j v_{kj} \varphi'\left(\sum_i w_{ji} x_i\right) \cdot \frac{\partial}{\partial w} \left(\sum_i w_{ji} x_i\right) =$$

vanishes if $j \neq j_0$

$$= \sum_k \boxed{(t_k - y_k) (-1) \varphi'\left(\sum_j v_{kj} a_j\right)} v_{kj_0} \varphi'\left(\sum_i w_{j_0 i} x_i\right) \cdot x_{i_0}$$