# Lecture 16

## Neural networks

<u>Idea</u>: make basis functions weight-dependent and fit those weights.

Previously, we considered

$$y(\vec{x}, \vec{w}) = f\left( \sum_{j=1}^{M} w_j \, \varphi_j(\vec{x}) \right), \text{ where}$$

$f(\cdot)$ is identity for regression and non-linear activation f'n for classification.

Now, consider $(x_1, \ldots, x_D)$ &larr; input vector

$$\underset{\underset{\text{activation}}{\uparrow}}{a_j} = \sum_{i=1}^{D} \underset{\underset{\text{weights}}{\uparrow}}{w_{ji}^{(1)}} x_i + \underset{\underset{\text{bias}}{\uparrow}}{w_{j0}^{(1)}} \qquad j = 1, \ldots, M$$

Then, $z_j = h(\underset{\underset{\text{non-linear activation f'n}}{\uparrow}}{a_j})$

$h(\cdot)$ can be $\sigma(\cdot)$ or $\tanh(\cdot)$

Next, $a_k = \sum_{j=1}^{M} w_{kj}^{(2)} z_j + w_{k0}^{(2)} \qquad k = 1, \ldots, \underset{\uparrow}{K}$

total # outputs

Finally, $\underset{\underset{\text{output vector}}{\nearrow}}{y_k} = \overset{\tilde{h}}{\bullet}(a_k)$, where $\tilde{h}$ may be $\underset{\underset{\downarrow}{\text{or } K>1}}{\overset{K=1}{}}$

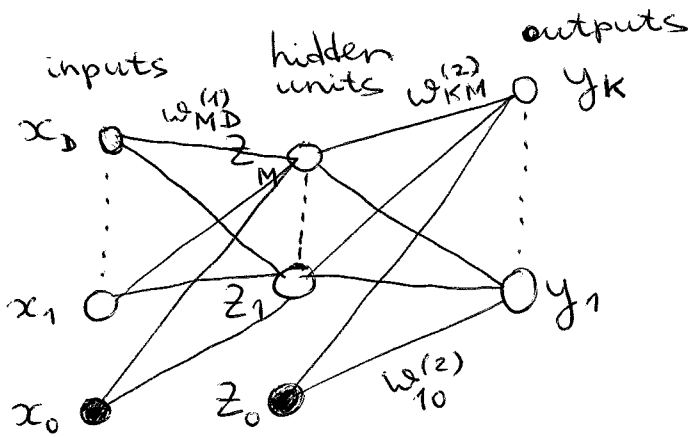identity for regression, $\sigma(\cdot)$ for binary classification ($K=2$), etc.

We have:

$$y_k(\vec{x}, \vec{w}) = \sigma\left( \sum_{j=1}^{M} w_{kj}^{(2)} \overbrace{h\left( \sum_{i=1}^{D} w_{ji}^{(1)} x_i + w_{j0}^{(1)} \right)}^{z_j, \; j=1,\ldots,M} + w_{k0}^{(2)} \right)$$

Note that it does not make sense to have $h(\cdot)$ as identity, since then the argument of the $\sigma$-function is just a linear model with various products of weights.
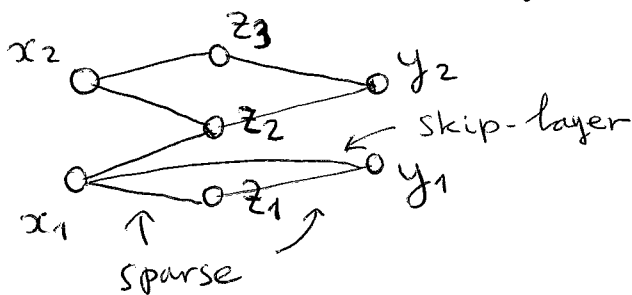


Define $x_0 = 1$ & $z_0 = 1$, then

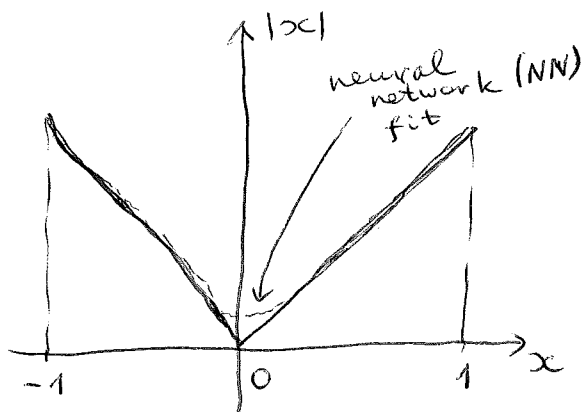$$y_k(\vec{x}, \vec{w}) = \sigma\left( \sum_{j=0}^{M} w_{kj}^{(2)} \underbrace{z_j}_{} \right) = \begin{cases} 1, & j=0 \\ h\left( \sum_{i=0}^{D} w_{ji}^{(1)} x_i \right), & j=1,\ldots,M \end{cases}$$

<u>Generalizations</u>: (1) multiple layers of hidden units

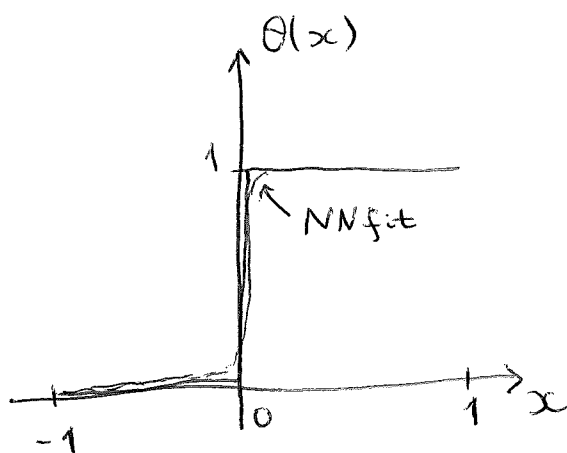(2) sparse network architectures

(3) skip-layer connections:

Note: feed-forward architectures <u>only</u>



skip-layer

sparse

Performance: can fit various functions fairly accurately


|x|, neural network (NN) fit


θ(x), NN fit

Sample $N = 50$ datapoints uniformly in $[-1, 1]$ interval, fit a two-layer network (input layer + hidden layer) discussed above:

3 hidden units,
tanh(.) activation f'n,
linear output units.

## weight-space symmetries:

with tanh(.) activation f's,

$$\tanh(-a) = -\tanh(a) \quad , \text{ and}$$

changing the sign of all weights (and the bias) leading into a unit can be compensated by the change in sign of all weights leading out of that unit.

M hidden units $\rightarrow 2^M$ equivalent weight vectors.

Similarly, can exchange weight values leading in and out of a hidden unit with another hidden unit $\Rightarrow M!$ permutations
So we have $2^M M!$ symmetries for a two-layer network (can be easily generalized to other architectures) and activation functions

# Network training

$$n = 1, \ldots, N : \quad \{\vec{x}_n\} \qquad \{\vec{t}_n\}$$

$\uparrow$ inputs $\qquad \uparrow$ targets

## Regression: (consider scalar targets $\{t_n\}$ for simplicity)

Assume
$$p(t \mid \vec{x}, \vec{w}) = \mathcal{N}(t \mid y(\vec{x}, \vec{w}), \beta^{-1}), \quad \beta$$

$$\beta = \text{precision} \left( = \frac{1}{\sigma^2} \right)$$

Output unit activation function =
= identity ; single output unit.

As before,
$$\mathcal{L} = \prod_{n=1}^{N} p(t_n \mid \vec{x}_n, \vec{w}, \beta) \quad , \quad \text{or}$$

$$-\log \mathcal{L} = \frac{\beta}{2} \sum_{n=1}^{N} (y(\vec{x}_n, \vec{w}) - t_n)^2 + \frac{N}{2} \log(2\pi) -$$
$$\overset{\shortparallel}{\text{error f'n}} \qquad\qquad - \frac{N}{2} \log \beta .$$

Define $\quad E(\vec{w}) = \frac{1}{2} \sum_n (y(\vec{x}_n, \vec{w}) - t_n)^2 ,$

then $\quad \left. \dfrac{\partial E}{\partial \vec{w}} \right|_{\vec{w}_{ML}} = 0 \quad$ gives $\quad \vec{w}_{ML}.$

$\nearrow$ 

ML approach

[or just minimizing $E(\vec{w})$]

$\uparrow$ need numerical minimizer

-4-

given $\vec{w}_{ML}$, $\left.\dfrac{\partial E}{\partial \beta}\right|_{\beta_{ML}} = 0$ yields

$$\beta_{ML}^{-1} = \frac{1}{N} \sum_n \left( y(\vec{x}_n, \vec{w}_{ML}) - t_n \right)^2 .$$

K=2 <u>classification</u>:

$$C_1 : t = 1 \qquad C_2 : t = 0$$

as before, use $\sigma(\cdot)$ on the single output node, and define

$$\begin{cases} p(C_1 | \vec{x}) = \underbrace{y(\vec{x}, \vec{w})}_{[0,1]} \\ p(C_2 | \vec{x}) = 1 - p(C_1 | \vec{x}) \end{cases}$$

Then $\quad \mathcal{I} = \prod_n y_n^{t_n} (1 - y_n)^{1 - t_n}$, where

$$y_n = y(\vec{x}_n, \vec{w}) ;$$

$$E(\vec{w}) = -\log \mathcal{I} = - \sum_n \left[ t_n \log y_n + (1 - t_n) \log(1 - y_n) \right].$$

If we have $N$ separate binary classifications to perform, $k = 1, \ldots, N \leftarrow$ $N$ output nodes

$$t_k = \{0, 1\}, \qquad k = 1, \ldots, N$$

$$\mathcal{I} = \prod_{n=1}^{N} \prod_{k=1}^{N} y_{nk}^{t_{nk}} (1 - y_{nk})^{1 - t_{nk}}, \text{ where}$$

$$y_{nk} = y_k(\vec{x}_n, \vec{w}) .$$

$$E(\vec{w}) = -\log \mathcal{I} = - \sum_n \sum_k \left[ t_{nk} \log y_{nk} + (1 - t_{nk}) \log(1 - y_{nk}) \right]$$

## $K > 2$ classification:

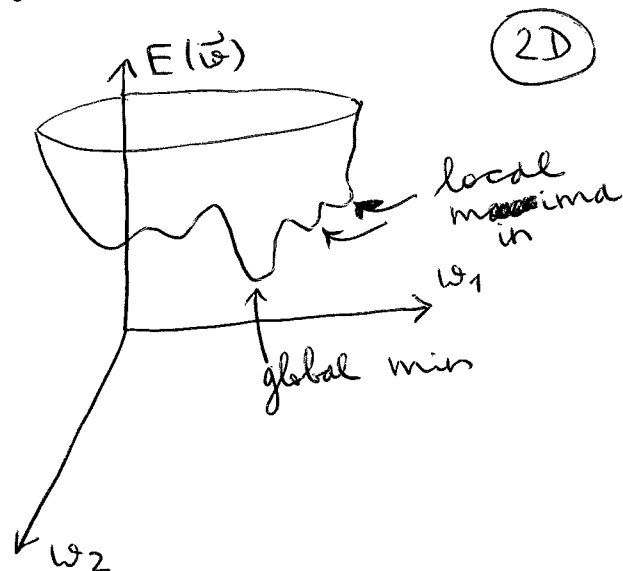1-of-$K$ coding scheme: $\overbrace{\underbrace{00 \cdots 1 \cdots}_{\uparrow \text{class label}}}^{K}$

Then interpret $y_K(\vec{x}, \vec{w}) = p(t_K = 1 | \vec{x}) \iff K$ output nodes

$$\mathcal{L} = \prod_{n=1}^{N} \prod_{k=1}^{K} y_{nk}^{t_{nk}} \quad , \quad \text{or}$$

$$E(\vec{w}) = - \sum_{n} \sum_{k} t_{nk} \log y_{nk}$$

Output unit activation function =
= softmax : $\quad y_K(\vec{x}, \vec{w}) = \dfrac{e^{a_K(\vec{x}, \vec{w})}}{\sum_{j=1}^{K} e^{a_j(\vec{x}, \vec{w})}}$ .

Typically, $E(\vec{w})$ is non-trivial:

# Prm optimization

If $\vec{\nabla} E(\vec{w})$ is available, could make small steps in the $-\vec{\nabla} E$ direction until $\vec{\nabla} E = 0 \implies$ steepest descent method.

This will only find a local min in general $\implies$ could run multiple times & compare the results (i.e., get the best one). Besides, there're symmetries in weight space (e.g. $2^M M!$ equivalent minima in a 2-layer network with M hidden units).

## Local quadratic approximation:

$$E(\vec{w}) \simeq E(\hat{w}) + (\vec{w} - \hat{w})^T \vec{b} + \frac{1}{2}(\vec{w} - \hat{w})^T H (\vec{w} - \hat{w})$$

$\uparrow$ Taylor expansion around $\hat{w}$

$$\vec{b} = \vec{\nabla} E \Big|_{\hat{w}}, \quad H_{ij} = \frac{\partial^2 E}{\partial w_i \, \partial w_j}\Big|_{\hat{w}} \qquad [H^T = H]$$

Then $\quad \vec{\nabla} E_{(\vec{w})} \simeq \vec{b} + H(\vec{w} - \hat{w})$

If $\quad \hat{w} = \vec{w}^*$ is a minimum, s.t. $\vec{b} = 0$, we get:

$$E(\vec{w}) \simeq E(\vec{w}^*) + \frac{1}{2}(\vec{w} - \vec{w}^*)^T H (\vec{w} - \vec{w}^*)$$

Consider $H\vec{u}_i = \lambda_i \vec{u}_i$ s.t. $\vec{u}_i^T \vec{u}_j = \delta_{ij}$

Now, expand $\vec{w} - \vec{w}^* = \sum_i \alpha_i \vec{u}_i$ , s.t.

$$E(\vec{w}) \simeq E(\vec{w}^*) + \frac{1}{2} \sum_{k\ell} (w_k - w_k^*) H_{k\ell} (w_\ell - w_\ell^*) =$$

$$= E(\vec{w}^*) + \frac{1}{2} \sum_{k\ell} \sum_{ij} \alpha_i \alpha_j u_{i,k} \underbrace{\sum_\ell H_{k\ell} u_{j,\ell}}_{\lambda_j u_{j,k}} =$$

$$= E(\vec{w}^*) + \frac{1}{2} \sum_{ij} \alpha_i \alpha_j \lambda_j \underbrace{\sum_k u_{i,k} u_{j,k}}_{\delta_{ij}} =$$

$$= E(\vec{w}^*) + \frac{1}{2} \sum_i \alpha_i^2 \lambda_i$$

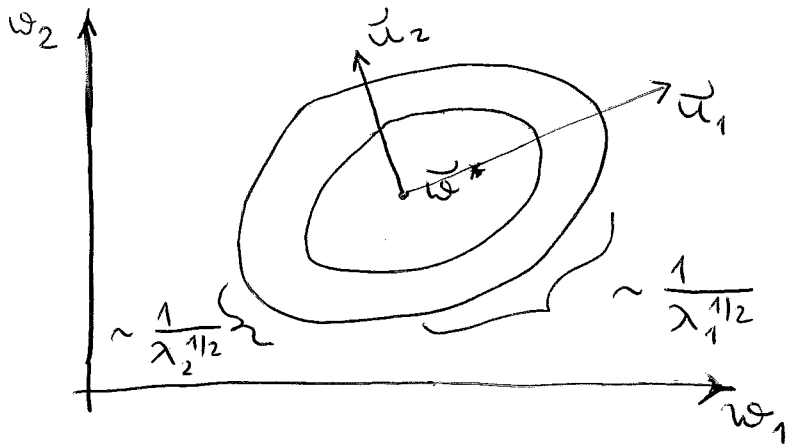$H$ is pos.-definite iff $\vec{v}^T H \vec{v} > 0$ , $\forall \vec{v}$.

Using $\vec{v} = \sum_i c_i \vec{u}_i$ , we obtain:

$$\vec{v}^T H \vec{v} = \sum_{ij} v_i H_{ij} v_j = \sum_{ij} \left( \sum_k c_k u_{k,i} \right) \times$$

$$\times H_{ij} \left( \sum_\ell c_\ell u_{\ell,j} \right) =$$

$$= \sum_{k\ell} c_k c_\ell \lambda_\ell \left( \underbrace{\sum_i u_{k,i} u_{\ell,i}}_{\delta_{k\ell}} \right) = \sum_\ell c_\ell^2 \lambda_\ell$$

$$\overset{\nearrow}{\sum_j H_{ij} u_{\ell,j} = \lambda_\ell u_{\ell,i}}$$

Thus $H$ is pos-def iff $\lambda_i > 0$ , $\forall i$.
This is a requirement for $\vec{w}^*$ to be a minimum, rather than a max or a saddle point.

## Contours of $E(\vec{w})$:



Suppose $E(\vec{w}) - E(\vec{w}^*) = \Delta$, then

$$E(\vec{w}) - E(\vec{w}^*) = \frac{1}{2} \sum_i \lambda_i d_i^2 = \Delta.$$

$$\text{If} \quad \vec{w} = \vec{w}^* + \underbrace{d_i \vec{u}_i}_{\text{along } \vec{u}_i} \quad \Rightarrow 2\Delta = \lambda_i d_i^2, \text{ or}$$

$$d_i \sim \frac{1}{\lambda_i^{1/2}}.$$

---

Steepest descent: $\quad$ learning rate

$$\vec{w}^{(\tau+1)} = \vec{w}^{(\tau)} - \eta \vec{\nabla} E(\vec{w}^{(\tau)})$$

(step #)

Note that $E(\vec{w}) = \sum_{n=1}^{N} E_n(\vec{w})$, so that,

alternatively, we can do

$$\vec{w}^{(\tau+1)} = \vec{w}^{(\tau)} - \eta \vec{\nabla} E_n(\vec{w}^{(\tau)})$$

$\uparrow$ sequential gradient descent
[cycle through datapoints]

This is a diff. algorithm b/c a local min for the whole dataset $\neq$ local min for each individual datapoint.