

Network training

Lecture 11

$$n=1, \dots, N : \quad \left\{ \begin{array}{c} \vec{x}_n \\ \uparrow \\ \text{inputs} \end{array} \right\} \quad \left\{ \begin{array}{c} \vec{t}_n \\ \uparrow \\ \text{targets} \end{array} \right\}$$

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

Assume $p(t|\vec{x}, \vec{w}) = \mathcal{N}(t|y(\vec{x}, \vec{w}), \beta^{-1})$
 $\beta = \text{precision} (= \frac{1}{\sigma^2})$

Output unit activation function = identity; single output unit.

as before,

$$\mathcal{J} = \prod_{n=1}^N p(t_n|\vec{x}_n, \vec{w}, \beta), \text{ or}$$

$$-\log \mathcal{J} = \frac{\beta}{2} \sum_{n=1}^N (y(\vec{x}_n, \vec{w}) - t_n)^2 + \frac{N}{2} \log(2\pi) - \frac{N}{2} \log \beta$$

" error f'n

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

then $\frac{\partial E}{\partial \vec{w}} \Big|_{\vec{w}_{ML}} = 0$ gives $\underline{\underline{\vec{w}_{ML}}}$.

ML approach

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

↑ need numerical minimizer

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

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

K=2 classification:

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

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 $\mathcal{L} = \prod_n y_n^{t_n} (1-y_n)^{1-t_n}$, where $y_n = y(\vec{x}_n, \vec{w})$;

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

If we have N separate binary classifications to perform, $t_k = \{0, 1\}$, $k=1, \dots, N$ \leftarrow N output nodes

$$\mathcal{L} = \prod_{n=1}^N \prod_{k=1}^M y_{nk}^{t_{nk}} (1-y_{nk})^{1-t_{nk}}, \text{ where } y_{nk} = y_k(\vec{x}_n, \vec{w}).$$

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

$K > 2$ classification:

1-of- K coding scheme: $\overbrace{00 \dots 1 \dots}^K$
↑ class label

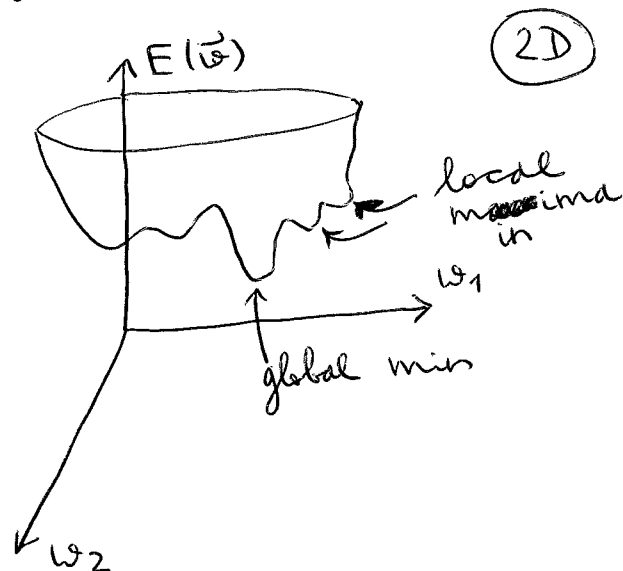
Then $y_k(\vec{x}, \vec{w}) = p(t_k = 1 | \vec{x})$ \Leftarrow K output nodes
interpret

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

$$E(\vec{w}) = - \sum_n \sum_k t_{nk} \log y_{nk}$$

Output unit activation function =
= softmax: $y_k(\vec{x}, \vec{w}) = \frac{e^{d_k(\vec{x}, \vec{w})}}{\sum_{j=1}^K e^{d_j(\vec{x}, \vec{w})}}$

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 $\vec{\nabla} E = 0 \Rightarrow$ steepest descent method.

This will only find a local min in general \Rightarrow 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:

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

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

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

$$\text{Then } \vec{\nabla} E \Big|_{\vec{w}} = \vec{b} + H(\vec{w} - \hat{\vec{w}})$$

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

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

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.

$$\begin{aligned} E(\vec{w}) &\approx E(\vec{w}^*) + \frac{1}{2} \sum_{k,l} (\omega_k - \omega_k^*) H_{kl} (\omega_l - \omega_l^*) = \\ &= E(\vec{w}^*) + \frac{1}{2} \sum_{k,l} \sum_{i,j} \alpha_i \alpha_j u_{i,k} \underbrace{H_{kl} u_{j,l}}_{\lambda_j u_{j,k}} = \\ &= E(\vec{w}^*) + \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j \lambda_j \underbrace{\sum_k u_{i,k} u_{j,k}}_{\delta_{ij}} = \\ &= E(\vec{w}^*) + \frac{1}{2} \sum_i \alpha_i^2 \lambda_i \end{aligned}$$

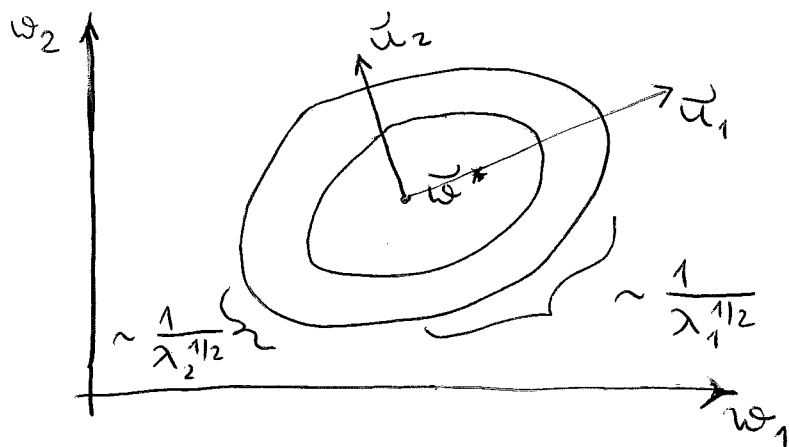
H is pos.-definite iff $\vec{v}^T H \vec{v} > 0$, $\forall \vec{v}$.

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

$$\begin{aligned} \vec{v}^T H \vec{v} &= \sum_{i,j} v_i H_{ij} v_j = \sum_{i,j} \left(\sum_k c_k u_{k,i} \right) \times \\ &\quad \times H_{ij} \left(\sum_l c_l u_{l,j} \right) = \\ &\quad \uparrow \sum_{k,l} c_k c_l \lambda_l \underbrace{\left(\sum_i u_{k,i} u_{l,i} \right)}_{\delta_{kl}} = \sum_l c_l^2 \lambda_l. \\ \sum_j H_{ij} u_{l,j} &= \lambda_l u_{l,i} \end{aligned}$$

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

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

$$\text{If } \vec{w} = \vec{w}^* + \underbrace{d_i}_{\text{along } \vec{u}_i} \vec{u}_i \Rightarrow 2\Delta = \lambda_i d_i^2, \text{ or } d_i \sim \frac{1}{\lambda_i^{1/2}}.$$

Steepest descent:

$$\vec{w}^{(\tau+1)} = \vec{w}^{(\tau)} - \underset{\substack{\uparrow \\ \text{learning rate}}}{\eta} \underset{\substack{\uparrow \\ \text{step \#}}}{\nabla} E(\vec{w}^{(\tau)})$$

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

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

↑ 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.