

# Fast multiplication by the Hessian

Consider  $\vec{v}^T H = v_j \nabla_j (\nabla_i E)$

Lecture 13

Goal: calculate  $\vec{v}^T H$  faster than  $\mathcal{O}(w^2)$  required to compute  $H$ .

Define  $R\{\cdot\} = \vec{v}^T \nabla$ , note that

$$R\{\omega_k\} = v_j \nabla_j \omega_k = v_k, \text{ or}$$

$$R\{\omega\} = \underline{\underline{\vec{v}}}$$

Two-layer network:

$$\begin{cases} a_j = \sum_i \omega_{ji} x_i, \\ z_j = h(a_j), \\ y_k = \sum_j \omega_{kj} z_j. \end{cases}$$

$$\text{Now, } \begin{cases} R\{a_j\} = \sum_i R\{\omega_{ji}\} x_i = \sum_i v_{ji} x_i, \\ R\{z_j\} = h'(a_j) R\{a_j\}, \\ R\{y_k\} = \sum_j v_{kj} z_j + \sum_j \omega_{kj} \underbrace{R\{z_j\}}_{h'(a_j) R\{a_j\}} \end{cases}$$

$$\text{Further, } \begin{cases} \delta_k = y_k - t_k, \\ \delta_j = h'(a_j) \sum_k \omega_{kj} \delta_k \end{cases} \text{ as before}$$

↑  
backpropagation

output

Then 
$$\begin{cases} R\{\delta_k\} = R\{y_k\}, \\ R\{\delta_j\} = h''(a_j) R\{a_j\} \sum_k W_{kj} \delta_k + \\ \uparrow \\ \text{hidden} \\ + h'(a_j) \sum_k v_{kj} \delta_k + h'(a_j) \sum_k W_{kj} R\{\delta_k\}. \end{cases}$$

Finally, 
$$\begin{cases} \frac{\partial E}{\partial W_{kj}} = \delta_k z_j, \\ \frac{\partial E}{\partial W_{ji}} = \delta_j x_i. \end{cases}$$

$$\begin{cases} R\left\{\frac{\partial E}{\partial W_{kj}}\right\} = R\{\delta_k\} z_j + \delta_k R\{z_j\}, \\ R\left\{\frac{\partial E}{\partial W_{ji}}\right\} = x_i R\{\delta_j\}. \end{cases} \quad (*)$$

- Algorithm:
1. Do a forward pass, compute  $R\{a_j\}$ ,  $R\{z_j\}$ ,  $R\{y_k\}$ .
  2. Compute  $R\{\delta_k\}$  &  $R\{\delta_j\}$ .
  3. Compute  $W$  elements of  $\vec{v}^T H$  using (\*)

This requires  $O(W)$  operations.

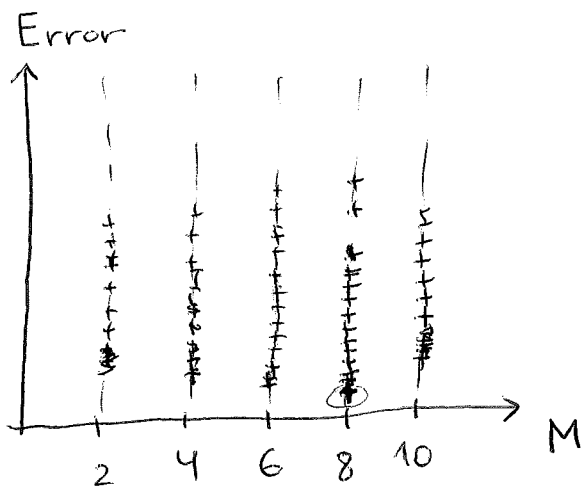
Note that this technique can be used to evaluate  $H$  by choosing  $W$  unit vectors as  $\vec{v}$ :  $\vec{v}^T = \underbrace{00\dots 1\dots 0}_{\substack{\uparrow \\ \text{jth position}}} \Rightarrow O(W^2)$  operations, equivalent to the direct calculation described above

# Regularization in NN

$M$  is a parameter to choose.

↑ # hidden units

Train NN on a training set, test on a test set:



← multiple local min's runs from random initializations

↑ choose globally best realization ( $M=8$  here)

More explicitly, one can use

$$\tilde{E}(\vec{w}) = E(\vec{w}) + \frac{\lambda}{2} \vec{w}^T \vec{w} \quad (*)$$

However, (\*) is NOT inv. wrt certain transformation properties obeyed by NNs.

Indeed, consider a two-layer network:

$$\begin{cases} y_k = \sum_j w_{kj} z_j + w_{k0}, & \leftarrow \text{regression} \\ z_j = h\left(\sum_i w_{ji} x_i + w_{j0}\right). \end{cases}$$

Consider  $x_i \rightarrow \tilde{x}_i = ax_i + b$ , then

$$\begin{cases} w_{ji} \rightarrow \tilde{w}_{ji} = \frac{1}{d} w_{ji}, \\ w_{j0} \rightarrow \tilde{w}_{j0} = w_{j0} - \frac{b}{a} \sum_i w_{ji} \quad \text{give} \end{cases}$$

$$\begin{aligned} \tilde{z}_j &= h\left(\sum_i \frac{w_{ji}}{d} (ax_i + b) + w_{j0} - \frac{b}{a} \sum_i w_{ji}\right) = \\ &= h\left(\sum_i w_{ji} x_i + w_{j0}\right) = z_j \Rightarrow \end{aligned}$$

$$\Rightarrow \underline{\underline{\tilde{y}_k = y_k}}$$

Similarly,  $y_k \rightarrow \tilde{y}_k = cy_k + d$  can be "undone" by weight / bias rescaling:

$$\begin{cases} w_{kj} \rightarrow \tilde{w}_{kj} = c w_{kj}, \\ w_{k0} \rightarrow \tilde{w}_{k0} = c w_{k0} + d. \end{cases} \Rightarrow \tilde{x}_i = x_i.$$

Clearly, the  $\lambda$ -term in (\*) breaks this invariance. So, use

$$\frac{\lambda_1}{2} \underbrace{\sum_{w \in W_1} w^2}_{\text{1st layer}} + \frac{\lambda_2}{2} \underbrace{\sum_{w \in W_2} w^2}_{\text{2nd layer}} \quad (\text{biases left unconstrained})$$

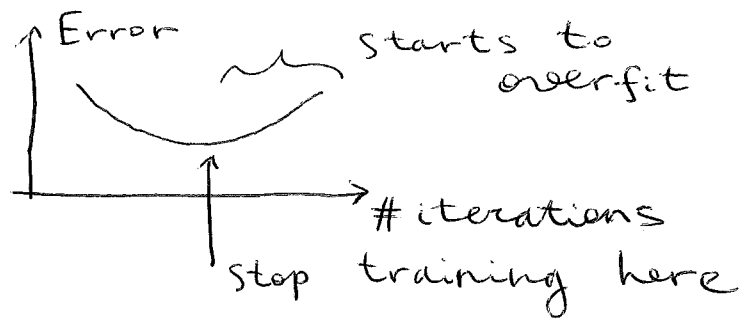
This regularizer will remain inv under  $\tilde{w}_{ji} = d^{-1} w_{ji}$ ,  $\tilde{w}_{kj} = c w_{kj}$  if

$$\lambda_1 \rightarrow d^{1/2} \lambda_1, \quad \lambda_2 \rightarrow c^{-1/2} \lambda_2.$$

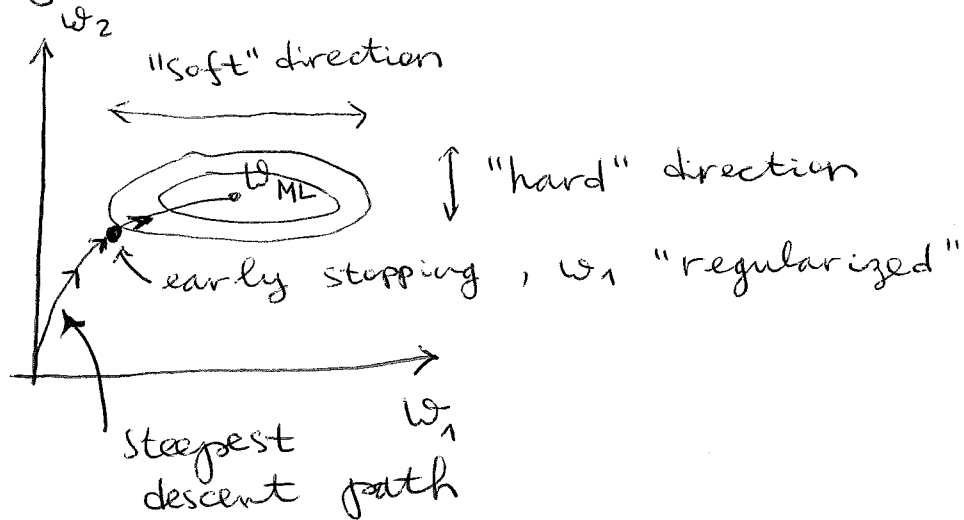
Can divide weights into more groups, not just  $W_1$  &  $W_2$ .

## Early stopping

Train on a training set, test on a test set:



Early stopping effectively controls model complexity:



## Invariances

Transformations of inputs (e.g. <sup>translating,</sup> rescaling, rotating digits in 2D ~~images~~ images) should not affect predictions.

However, raw data (pixel intensities) change in complex ways.

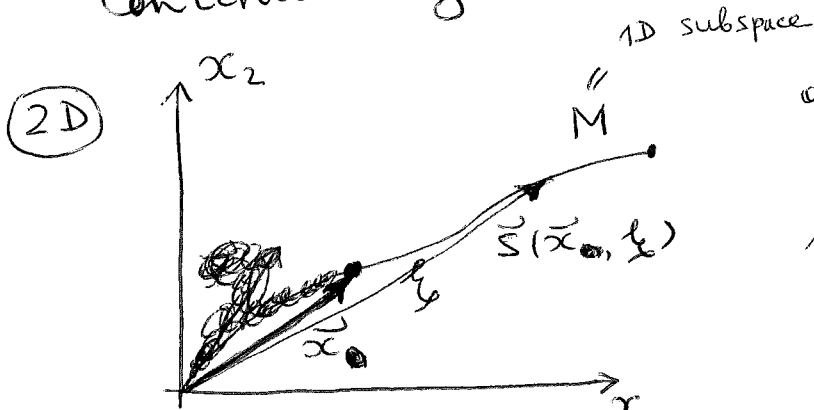
How to deal with this?

1. Add transformed patterns to the dataset, let the NN "sort it out" (i.e. learn the invariances) [data intensive]
2. Add a term to the error function to penalize changes in NN output when input is transformed (tangent propagation)
3. Extract transform-inv features from raw data first [domain knowledge]
4. Build the invariance into structure and weight distributions of the NN

### Tangent propagation

Use regularization to encourage models to be transform-inv.

Continuously transform  $\vec{x}_0$ :



one-parameter continuous transform:  $\xi$

Note that

$$\vec{S}(\vec{x}_0, 0) = \vec{x}_0$$

$$\text{Then } \vec{\tau} = \frac{\partial \vec{S}}{\partial \xi} \Rightarrow \vec{\tau}_0 = \left. \frac{\partial \vec{S}}{\partial \xi} \right|_{\xi=0}$$

tangent vector

Finally,

$$\left. \frac{\partial y_k}{\partial \xi} \right|_{\xi=0} = \sum_{i=1}^{\overset{\# \text{ input nodes}}{\downarrow} D} \frac{\partial y_k}{\partial x_i} \left. \frac{\partial x_i}{\partial \xi} \right|_{\xi=0} = \sum_i \underbrace{J_{ki}}_{\text{Jacobian}} \tau_{0,i}$$

We can introduce  $\tilde{E} = E + \lambda \Omega$ , where

$$\begin{aligned} \Omega &= \frac{1}{2} \sum_{n,k} \left( \left. \frac{\partial y_{nk}}{\partial \xi} \right|_{\xi=0} \right)^2 \\ &= \frac{1}{2} \sum_{n,k} \left( \sum_{i=1}^D J_{n,ki} \tau_{0,i} \right)^2 \end{aligned}$$

This will encourage NN weights to exhibit invariance wrt the transforms. In practice  $\tau$  can be approximated by finite differences for small  $\xi$  [e.g. rotated digits, see Fig. 5.16 in the book]

Also, need to extend the backpropagation technique to  $\tilde{E}$  (that is, backpropagate derivatives of  $\Omega$  wrt weights).