# Deep Learning (BEV033DLE)
# Lecture 14 Recurrent Neural Networks

Czech Technical University in Prague

◆ Recurrent models

◆ Error back propagation through time

◆ Gated recurrent units, GRU and LSTM networks

◆ Recurrent back propagation

**Recurrent models in a nutshell**

◆ input sequence $x = (x_1, \ldots, x_t, \ldots, x_T)$, $x_t \in \mathbb{R}^n$. Similarly: output sequence $y$ with elements $y_t$ and sequence $h$ of (hidden) states with elements $h_t \in \mathbb{R}^d$. Often all three sequences have the same length.

◆ recurrent (dynamic) system with outputs

$$h_t = f(x_t, h_{t-1}, w)$$

$$y_t = g(h_t, v)$$

where $w$ and $v$ are parameters. The model defines sequence mappings $h = F_w(x)$ and $y = G_v(h)$.

◆ loss function $\ell(y, y')$; often locally additive $\sum_t \ell(y_t, y'_t)$

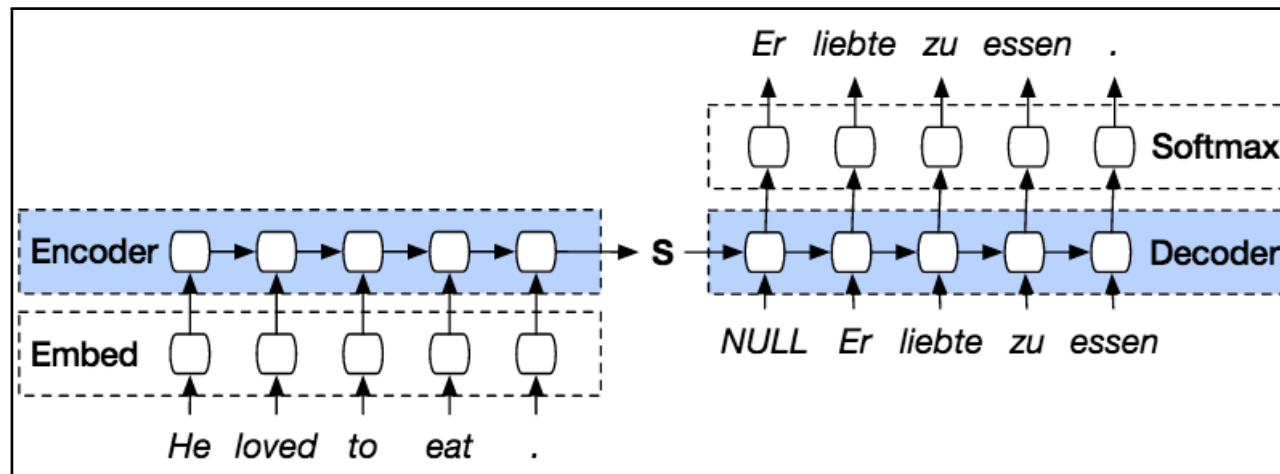**Training goal:** given training data $\mathcal{T}^m = \{(x^j, y^j) \mid j = 1, \ldots, m\}$, learn the model parameters $w$, $v$ by solving

$$\frac{1}{m} \sum_{j=1}^{m} \ell\big(y^j, (G_v \circ F_w)(x)\big) \to \min_{w,v}$$

Incarnations of recurrent models and related tasks

◆ Deep neural network for classification with additional feedback connections. $x$ - input, constant not depending on time. $y$ - output of the network, network head, e.g. $\log \mathrm{softmax}$, $h$ -states of all hidden layers. The loss function depends only on the last output $y_T$.

◆ "infinite state automata": the output space is sufficient for keeping the history, thus $h$ and $y$ can be identified, i.e. $y_t = f(x_t, y_{t-1}, w)$.
Example: landcover type monitoring for a geo-location: $x$ - sequence of spectral satellite measurements, $y$ - sequence of states (e.g. coniferous forest, broadleaf forest, clearcut, bark beetle degradation etc.)

◆ general sequence segmentation: hidden states $h_t$ are needed for keeping track of longer past and are latent.
Examples: speech recognition, $x$ - audio signal, $y$ -sequence of words. NLP translation:

Learning RNNs is particularly simple in the case that

◆ $h$ and $y$ can be identified, i.e. $y_t = f(x_t, y_{t-1}, w)$ and

◆ the loss is locally additive $\sum_t \ell(y_t, y_t')$

We can split the sequences $(x^j, y^j)$ from training data into triplets $(y_{t-1}^j, x_t^j, y_t^j)$ and train $f$ from

$$\frac{1}{m} \sum_{j=1}^{m} \sum_{t} \ell\big(y_t^j, f_w(x_t^j, y_{t-1}^j)\big) \to \min_w$$

Neither forward nor backward propagation through the sequence are needed.

If the hidden states $h_t$ do not coincide with outputs $y_t$ and are latent, then learning becomes considerably more complicated.

**Assumptions:**

$$h_t = f(x_t, h_{t-1}, w)$$

$$y_t = g(h_t, v)$$

The mappings $f$ and $g$ are implemented by neural networks and are differentiable w.r.t. their inputs and parameters. The loss function $\ell(y, y')$ is differentiable.
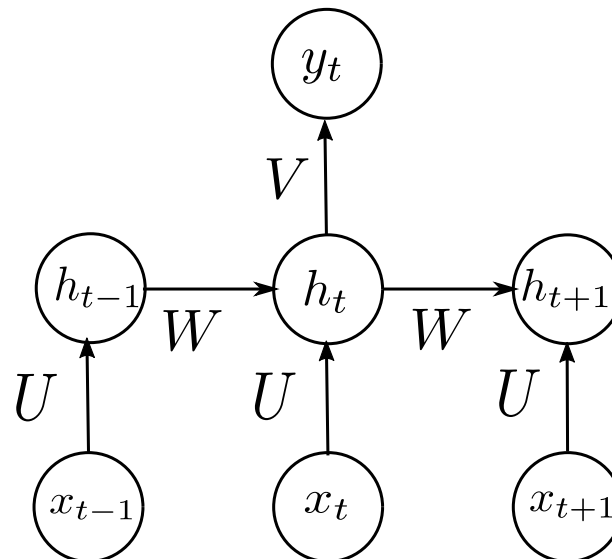
**Example 1.** Both mappings $f$ and $g$ are implemented by one layer networks

$$a_t = W h_t + U x_t + b \qquad\qquad h_t = \tanh(a_t)$$

$$o_t = V h_t + c \qquad\qquad y_t = \mathrm{softmax}(o_t)$$

**Computing the gradients:** Unroll the network in time and apply backpropagation

Let us consider the loss for a single example $(x, y^*)$ from the training data.

Computing the gradient w.r.t. $v$ is easy (see Slide 4.). Let us consider the gradient w.r.t. $w$

$$\partial_w L(y^*, y) = \sum_{t=1}^{T} \partial_w \ell(y_t^*, y_t) = \sum_{t=1}^{T} \partial_{y_t} \ell(y_t^*, y_t) \, \partial_{h_t} g(h_t, v) \, \partial_w h_t$$

The first two terms are simple. For the last one we have the recurrent expression

$$\partial_w h_t = \partial_w f(x_t, h_{t-1}, w) + \partial_{h_{t-1}} f(x_t, h_{t-1}, w) \, \partial_w h_{t-1}$$

This gives

$$\partial_w h_t = \partial_w f(x_t, h_{t-1}, w) + \sum_{i=1}^{t-1} \Big[ \prod_{j=i+1}^{t} \partial_{h_{j-1}} f(x_j, h_{i-1}, w) \Big] \partial_w f(x_i, h_{i-1}, w)$$

**Problems:**

◆ backpropagation through time is computationally expensive

◆ Exploding/vanishing gradients: consider for simplicity the linear recurrence $h_t = W h_{t-1}$. For $\tau$ steps we get $h_\tau = W^\tau h_0$. Suppose that we can write $W = U^{-1} \Lambda U$, where $\Lambda$ is diagonal. We get

$$h_\tau = U^{-1} \Lambda^\tau U h_0.$$

Eigenvalues with magnitude less than one will decay and eigenvalues with magnitude greater than one will explode.

◆ We can not apply batch normalisation as simple remedy.

◆ We want the following model ability: events long in the past can trigger changes in conjunction with current measurements.

◆ skip connections?, designate special nodes in $h_t$ for keeping record of events long in the past?

# RNNs with gated recurrent units

LSTM (Hochreiter, Schmidhuber, 1997), GRU (Cho et al., 2014), ...

**Gated recurrent unit (simplified):**

A cell consisting of a recurrent unit $h_t$ and a gate unit $u_t \in [0,1]$

$$h_t = u_{t-1} h_{t-1} + [1 - u_{t-1}] f(x_t, h_{t-1}, w)$$
$$u_t = \mathrm{S}(x_t, h_t, v)$$

The gate unit $u_t$ has sigmoid nonlinearity and "decides" whether to copy $h_t$ from $h_{t-1}$ or to apply the recurrence with $f$.

**Gated recurrent unit (general):**

- $\blacklozenge$   $h$ is a state vector

- $\blacklozenge$   $u$ is a vector of "update" gates

- $\blacklozenge$   $r$ is a vector of "reset" gates

The update equations are

$$h_t = u_{t-1} \odot h_{t-1} + [1 - u_{t-1}] \odot \mathrm{S}\Big(U x_{t-1} + W r_{t-1} \odot h_{t-1}\Big)$$

where $\odot$ denotes the element-wise product of vectors. The gate unit outputs are given by

$$u_t = \mathrm{S}\big(U^u x_t + W^u h_t\big)$$
$$r_t = \mathrm{S}\big(U^r x_t + W^r h_t\big)$$

LSTM cells are somewhat more complicated – they have separate "forget" and "update" gates.

**Recurrent backpropagation:** (Almeida, 1987), (Pineda, 1987) An interesting learning alternative can be applied to deep neural networks for classification with additional feedback connections.

Denote: network input $x$, network output $y_t$ and $h_t$ denoting outputs of all hidden layers.

$$h_t = f(x, h_{t-1}, w) \quad \text{and} \quad y_t = g(h_t, v)$$

**Assumption:** the network configuration $h_t$ converges to a fixpoint $h^*$ if we clamp its input to $x$.

Then we have (implicit function theorem)

$$\frac{\partial h^*}{\partial w} = \left[I - J_F(h^*)\right]^{-1} \frac{\partial F}{\partial w},$$

where $J_F(h^*) = \frac{\partial F(x,w,h^*)}{\partial h}$ is the Jacobian of $F$ w.r.t. $h$.

Now, let us consider the gradient of the loss w.r.t. $w$.

$$\partial_w L = \partial_y L \, \partial_{h^*} y \left[I - J_F(h^*)\right]^{-1} \partial_w f(x, w, h^*)$$

Now, introduce the (column) vector $z$

$$z = \left[I - J_F(h^*)\right]^{-1} \left(\partial_y L\, \partial_{h^*} y\right)^T$$

Multiplying both sides by $\left[I - J_F(h^*)\right]$, we get

$$z = J_F(h^*)^T z + \left(\partial_y L\, \partial_{h^*} y\right)^T.$$

This is a fixpoint equation for $z$ and can be solved by fixpoint iteration. The resulting algorithm for computing the derivative $\frac{\partial L}{\partial w}$ is:

◆ start from $z_0$; iterate

$$z_i = J_F(h^*)^T z_{i-1} + \left(\partial_y L\, \partial_{h^*} y\right)^T$$

until convergence.

◆ Return

$$\frac{\partial L}{\partial w} = z^T \frac{\partial F(x, w, h^*)}{\partial h}$$

This supersedes BPT but requires invertible $\left[I - J_F(h^*)\right]$.