#### E9: 309 Advanced Deep Learning 12-10-2020

**Instructor:** Sriram Ganapathy sriramg@iisc.ac.in

**Teaching Assistant : Akshara Soman, Prachi Singh, Jaswanth Reddy** aksharas@iisc.ac.in, prachis@iisc.ac.in, jaswanthk@iisc.ac.in

> Schedule - MW - 430-6pm (Microsoft Teams) http://leap.ee.iisc.ac.in/sriram/teaching/ADL2020/









http://phdcomics.com/

E9:309 Advanced Deep Learning

#### www.phdcomics.com



Recap of deep learning







\* Filling the google form in the webpage

- Contents will be made available to the folks in the mailing lists.

with creditors as well as video links.

\* Online registration portal from <u>academics.iisc.ac.in</u>

Your research/faculty advisor may need to approve also before the deadline (Oct. 20th?)

- Announcements regarding evaluations and projects will be shared only
  - **\*** Teams channel interaction and TA session for creditors only.



#### Some notations

 $* \mathbf{x} \in \mathbf{R}^{D}$  - input data.

 $* \mathcal{Y} \in \mathcal{R}^{\mathcal{C}}$ - neural network targets.

 $* \hat{y} \in B^C$  - model outputs.

### $*e, h \in \mathbb{R}^d$ - hidden model representations or embeddings.

\* 🖯 - collection of learnable parameters in the model.

 $* E(y, \hat{y})$  - error function used in the model training.







\*  $\{x_1, ..., x_N, y_1, ..., y_N\}$  - labeled training data

 $*q = \{1...Q\}$  - iteration index.

 $* t = \{1...T\}$  - discrete time index.

 $*^{l} = \{1...L\}$  - layer index

\* - learning rate (hyper-parameter)

\*  $N_b$  - mini-batch size and B is the number of mini-batches.





\* Training data, validation data, test data.



\* Model training data - used for parameter learning.

\* Validation data - used for hyper-parameter tuning (cross validation CV).

E9:309 Advanced Deep Learning

# <section-header><section-header>

#### Unseen Test Data

Evaluation data



#### Last lecture

- \* Feedforward networks
- \* Convolutional networks
- \* Learning in deep networks
  - Parameters, hyper-parameters etc
  - → SGD learning rule
  - Adam optimization, batch normalization.
- \* Reading assignments.





### Module - I Visual and Time Series Modeling





- \* Learn from ordered pairs of  $\mathcal{I}, \mathcal{Y}$ 
  - ✓ All the data samples are treated independently.
    - **\*** Data are shuffled before mini-batch formation
- \* If the input data and output labels are time-series data x(t), y(t)
  - DNNs/CNNs may fail to model the correlation of the data across the time
  - Question how can we build models that capture the time evolution of the data and the labels.





\* An interesting subset of this problem is where the input alone is a time series x(t), y or have different indices x(t), y(u)

- \* Examples
  - ✓ Text sequences Speech and audio Video sequences



#### First order recurrence - hidden layer

\* Making the hidden layer a function of the previous outputs from the hidden layer along with the input

# h(t) = f(h(t-1), x(t)) $\boldsymbol{x}(t)$









#### First order recurrence - output layer

\* Making the hidden layer a function of the previous outputs from the output layer along with the input

 $\boldsymbol{x}(t)$ 

E9:309 Advanced Deep Learning



# $h(t) = f(\hat{y}(t-1), x(t))$







#### First order recurrence - output layer

\* Making the hidden layer a function of the previous outputs from the output layer along with the input

 $\hat{y}(t)$ Training

E9:309 Advanced Deep Learning



# $\boldsymbol{h}(t) = f(\hat{\boldsymbol{y}}(t-1), \boldsymbol{x}(t))$ $\hat{\boldsymbol{x}}(t)$ **Testing**





#### \* Learning in recurrence networks: Back-propagation in time.

#### \* Learning considerations : Issues with forgetting and long-short-term memory networks



#### First order recurrence - hidden layer

\* Making the hidden layer a function of the previous outputs from the hidden layer along with the input.



\* Makes the hidden layer dependent of previous layer outputs in a recurring fashion.





#### First order recurrence - hidden layer

\* Making the hidden layer a function of the previous outputs from the hidden layer along with the input.

# $h(t) = f(h(t-1), \boldsymbol{x}(t))$ $\boldsymbol{x}(t)$ $\hat{\boldsymbol{y}}(t)$

\* Makes the hidden layer dependent of previous layer outputs in a recurring fashion.

E9:309 Advanced Deep Learning



# Model Forward Pass - 1 hidden layer





#### \* Error functions are computed at every time-instant



E9:309 Advanced Deep Learning

 $E(y(T), \hat{y}(T))$  k(t+1) h(T-1) x(T)



## Error back propagation

 Output activations  $\frac{\partial E}{\partial a^2(t)} = \hat{y}(t) - y(t) \quad for \quad t = 1 \dots T$ ✓ Hidden activations at last instant T  $\checkmark$  Hidden activations for previous instances t = T-1,... l  $(W^2)^T \frac{\partial E}{\partial t} + \frac{\partial h^1(t+1)}{\partial t}$ 





 $\checkmark$  Hidden activations for previous instances t = T-1,... 1



E9:309 Advanced Deep Learning

✓ Using the above gradient the rest of the gradients can be computed ...









### Long-term dependency issues









#### Gradients tend to vanish or explode



#### ✓ Intial frames may not have impact in the final predictions.





### Long short term memory (LSTM) idea





