

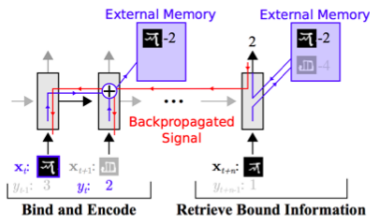
x_1 x_1' x_1''
 train

Model-based Meta-learning (cont) 2

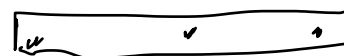
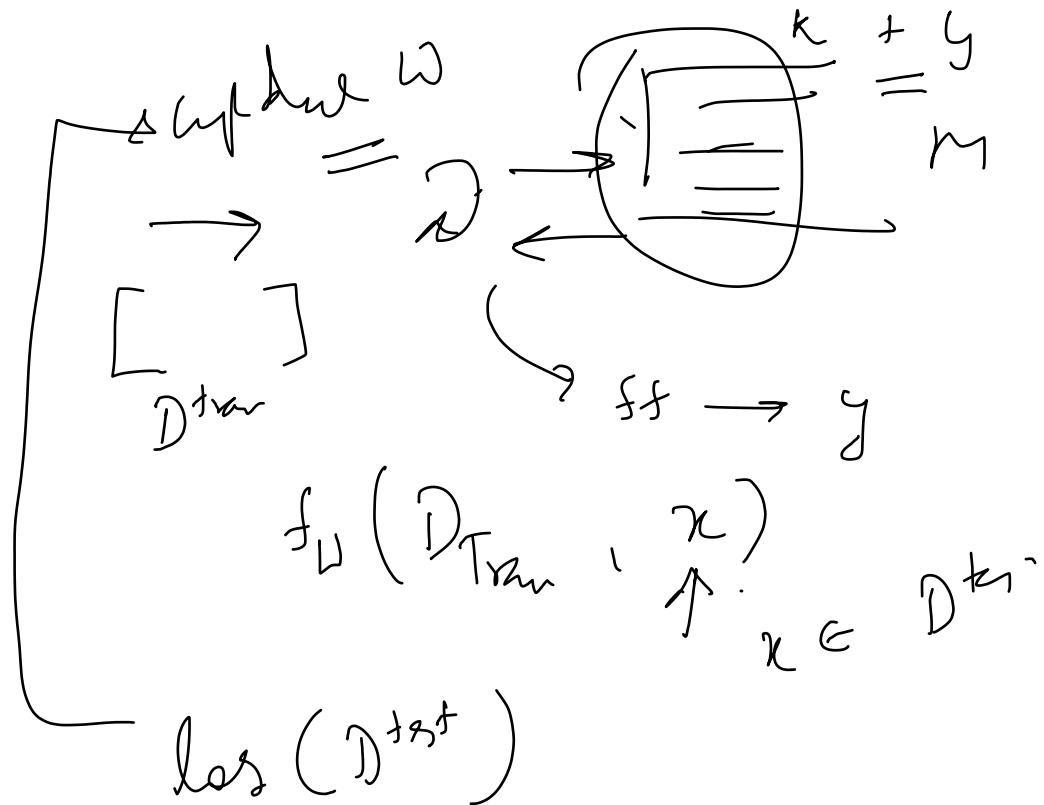
Reference - *Meta-Learning with Memory-Augmented Neural Networks*

Datasets are presented as sequences, $\{(x_t, y_{t-1})\}$ as in the RNN-based approach. A 'memory-augmented LSTM/FF network' is the meta-learner.

Network updates rows of a memory matrix M with keys $k_t = \phi(x_t)^T$:
 $M_t = M_{t-1} + w^w k_t$; w^w = write weights
 Retrieved memory $r_t = w^r M_t$ used to predict \hat{y}_t using a feedforward layer.
 (read weights $w^r = \text{Softmax}(M_t \cdot k_t)$)
 usage weights: $w_t^u = \gamma u_{t-1}^u + w_t^r + w_t^w$
 used to compute w^w as follows:



'least used' $w^u = 1$ for t smallest elements of w^u , 0 otherwise
 $w_t^w = \delta w_{t-1}^w + (1 - \delta) w_t^u$; prior to writing the least used row of M is zeroed.
 Read/write to 'memory' - we can view this as a 'neural Turing machine'.



Model-based Meta-learning (cont) 3

Reference - [A Simple Neural Attentive Meta-learner](#)

Similar to Hochreiter's work above, i.e., dataset is passed as sequence but uses 1D convolutions and attention layers instead of an LSTM. Note: final prediction y_{test} alone is used with target y to backpropagate errors, so train dataset is passed as $z^0 = \{(x_1, y_1) \dots (x_{t-a}, y_{t-1}), (x_t, -)\} \in \mathbb{R}^{d \times t}$ for any query x_t . Both train and test data for each task are used as queries for training.

Mix of 1D convolution $C()$ and attention $A()$ layers:

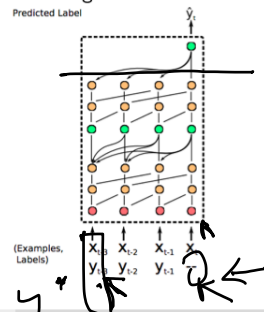
$u = \text{1Dconv}(z, 2^i)$; $a = \tanh(u) * \sigma(u)$; $C(z) = (z, a)$

$k = W_k z + b_k$, $q = W_q z + b_q$, $v = W_v z + b_v$

$p = \text{Softmax}(\frac{qk^T}{\sqrt{d}})$; $A(z) = (z, vp^T)$ $z_{in} \in \mathbb{R}^{d \times t}$

Note: convolutions and softmax are *causal*.

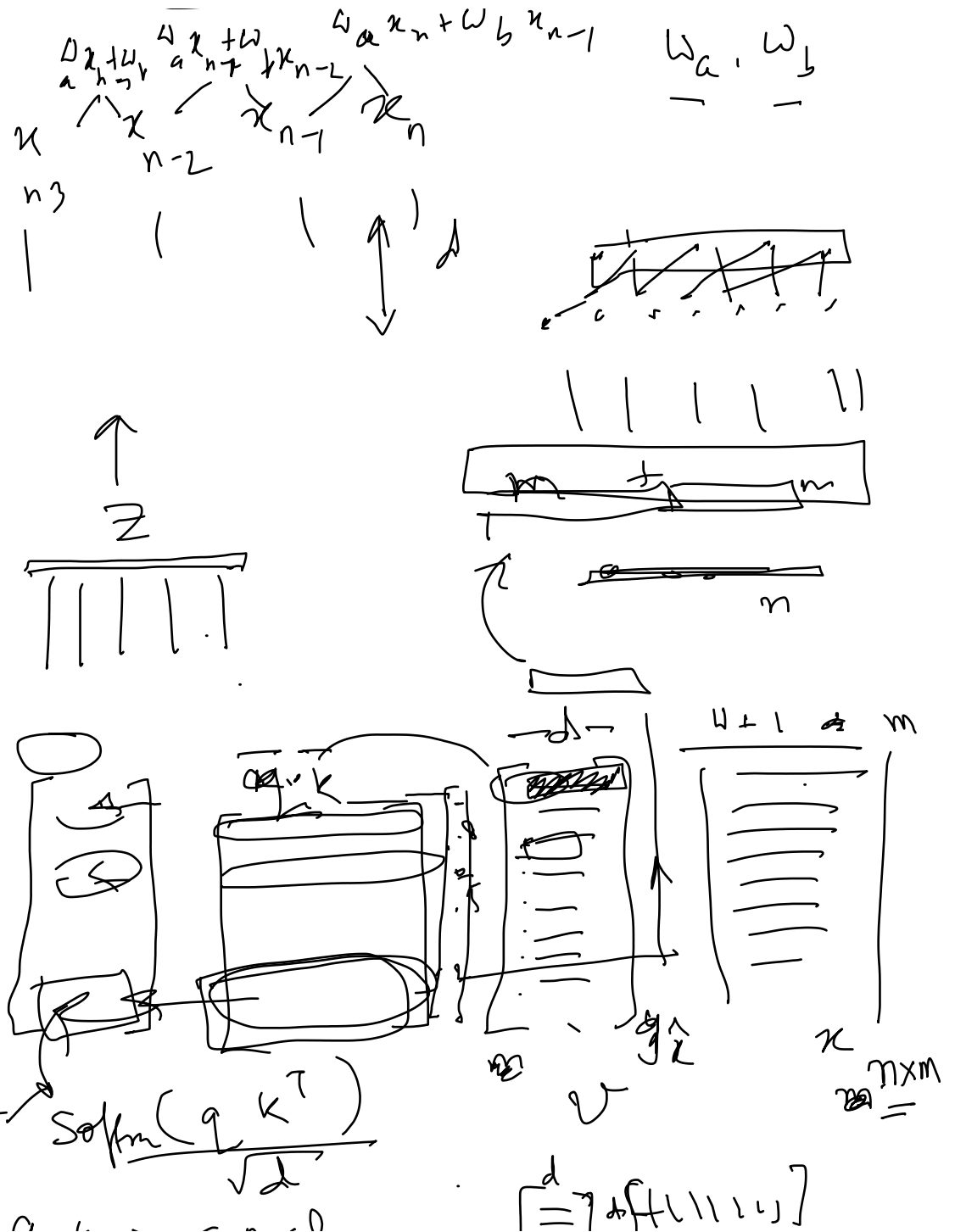
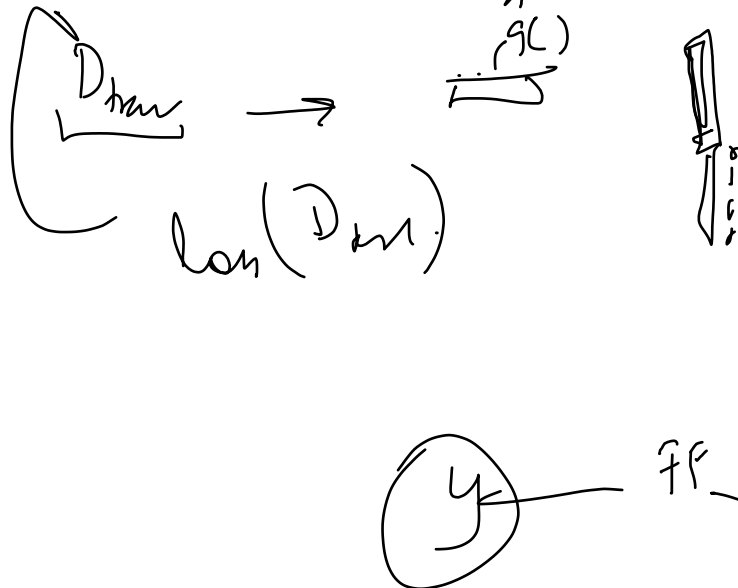
Using convolutions instead of RNNs is faster and easier to train. Attention is able to 'retrieve' inputs in position-independent manner.



Gautam Shroff

Meta-learning

Spring 2021 22 / 56

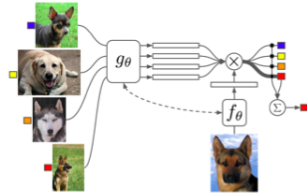


Metric-based Meta-learning

First to introduce the training procedure for few-shot learning, i.e. sampling varieties of few-shot tasks and using test accuracies as an optimisation objective. The network computes a distance kernel ('attention kernel') combining LSTM and attention mechanisms over a given few-shot training set; thus the classifier (based on distance kernel) is 'non-parametric' in that it changes with training set.

$$\hat{y} = \sum_{i=1}^k a(\hat{x}, x_i) y_i \text{ where } a(\hat{x}, x_i) = \frac{e^{c[f_{\theta}(\hat{x}) \cdot g_{\theta}(x_i)]}}{\sum_{i=1}^k e^{c[f_{\theta}(\hat{x}) \cdot g_{\theta}(x_i)]}}$$

Simple case: f, g are MLPs;
full-context case:
 $g(x_i, S)$ takes training examples as input and feeds into $f(\hat{x}, g(S))$.



$$\mathcal{Z}, \mathcal{K}, \mathcal{I} \in \mathbb{R}^{n \times n}$$

$$| \equiv |$$

$$\mathcal{W} \mathcal{X} + \mathcal{b}$$

$$\mathcal{K} \mathcal{W} + \mathcal{b}$$

$$n \times m \uparrow \boxed{m \times d}$$

- ✓ MAML - Opt. based
- ✓ CNP - Model-based
- ✓ MAN - Metric-based

HW 2

compare all three / short ways

speed of training
accuracy

met-learn sample
train

Metric-based Meta-learning (cont)

$D_{Train} = (\{(x_i, y_i)\}, \{\tilde{x}_j\})$; $D_{Test} = \{\tilde{x}_j\}$. D_{Train} has unlabeled examples \tilde{x}_j for semi-supervised learning; y_i s one-hot, and $\tilde{y}(\cdot), \bar{y}(\cdot) = \frac{1}{C}$. D_{Train}, D_{Test} fed to GNN,

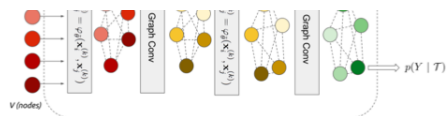
GNN is on a FC graph

$$\mathcal{G}(\{x \in \mathbb{R}^d\} = \{(x, y)\}, \varphi)$$

$$\mathcal{O}^k(x; x) = MI P_{\text{Adv}}(x; -x; \cdot)$$



GC: aggregate from neighbors & concatenate; $W^k, V^k \in \mathbb{R}^{d^{k+1} \times d^k}$:
 $x_i^{k+1} = (\sum_{j \neq i} \varphi_{ij}^k W^k x_j, V^k x_i)$



Finally $p(y|\mathcal{T}) = \text{Softmax}(x^N)$ for all nodes, test and training. Network is trained using available training and test labels across many tasks - each task's data $D_{\text{Train}}, D_{\text{Test}}$ is input as a graph, with possibly different number of examples. This is a mix of metric-based (due to φ s) and model-based approaches.

Score = 1
 mch-hm
 (fix # of training sem)

PAC-Bayes Model: Risk of hypothesis $h(x) = y$, $R(h) = y$ Empirical risk: $R_s(h) = y$, $S = \{x_i, y_i\}$ Gibbs classifier GQ : sample h from distribution Q : Gibbs risks: $R(GQ) = EQR(h)$ and $R_s(GQ) = R_s(h)$ Generalization bounds in PAC-Bayes model: $KL(Q||P) + \log -T \Pr R(GQ) \leq 21-5 \frac{2}{2(m-1)}$ where P is the prior distribution of h "before learning". Note that Q is a function of P and S , i.e., Q is 'learned' starting from P using the samples S ; we bound the true Gibbs risk R in terms of the empirical Gibbs risk R_s , for h sampled from Q post learning. spin. ml