

(\times) $[\log] \rightarrow @ \text{Exp} \rightarrow b$) Reverse Mode Automatic Differentiation recursively applies Chain-Rule

$Dy_i / Dx_j = \sum_{\theta \in \text{Paths}(i, j)} \prod_{k \in \text{Path}(i, j)} \frac{\partial y_k}{\partial x_j}$ or exponential paths $O(2^n)$ or fully connected $O(\text{Neurons} \cdot \text{hidden layer})$

Backprop (constructive theorem): same asymptotic complexity as original functions, many ineff. algos poss.

$f'(x) = \lim_{h \rightarrow 0} f(x+h) - f(x) / h$ $O(h^2) \frac{\partial}{\partial x} f(g(x)) = f'(g(x)) \cdot g'(x)$ $| \theta^m = \theta^m - \gamma \cdot \theta^m \cdot \nabla L(f(x; \theta^m), y)$

$C = 0 \mid x = 1 \mid ax = a \mid x^2 = 2x \mid \sqrt{x} = (\sqrt{2}) \cdot x^{1/2} \mid e^x = e^x \mid a^x = \ln(a) \cdot a^x \mid \ln(x) = \ln(x) \mid \sin = \cos \mid \cos = -\sin$

Log-linear model very general family of distributions, exponential of lin. function: $p(y|x, \theta) = \exp(\theta \cdot f(x))$

How good is x/y together ≥ 0 | Partition function normalizes $|\log(p(y|x, \theta))| = \theta \cdot f(x, y) + \text{const}$ $E[y \exp(\theta \cdot f(x, y))]$

$\theta^* = \arg\max L(\theta) = \arg\max \prod_{n=1}^N p(y_n|x_n, \theta) = \arg\min -\sum_{i=1}^N \log p(y_n|x_n, \theta) \rightarrow \nabla^2 L(\theta) = -\sum_{i=1}^N \nabla^2_{\theta \theta} p(y_i|x_i, \theta)$ (covariance matrix of random vector)

$\frac{\partial L(\theta)}{\partial \theta} = -\sum_{i=1}^N \frac{\partial \log p(y_i|x_i, \theta)}{\partial \theta} = -\sum_{i=1}^N \frac{\partial}{\partial \theta} (\theta \cdot f(x_i) - \log E[y_i \exp(\theta \cdot f(x_i, y_i))]) = -\sum_{i=1}^N (f_k(x_i, y_i) \cdot \frac{\partial \log E[y_i \exp(\theta \cdot f(x_i, y_i))]}{\partial \theta} - f_k(x_i, y_i) \cdot E[y_i \exp(\theta \cdot f(x_i, y_i))]) = -\sum_{i=1}^N (f_k(x_i, y_i) \cdot E[y_i \exp(\theta \cdot f(x_i, y_i))] - f_k(x_i, y_i))$

$= -\sum_{i=1}^N f_k(x_i, y_i) + \sum_{i=1}^N E[y_i \exp(\theta \cdot f(x_i, y_i))] \cdot \sum_{i=1}^N E[y_i \exp(\theta \cdot f(x_i, y_i))] = -\sum_{i=1}^N f_k(x_i, y_i) + \sum_{i=1}^N E[y_i \exp(\theta \cdot f(x_i, y_i))] \quad (\text{Observed feature count} \neq \text{expected feature count when } \nabla L(\theta) = 0)$

Softmax $h_i, y, T = \exp(h_i/T) / \sum_{j \neq i} \exp(h_j/T)$ | Log-linear is dot product followed by softmax $h_i = \theta \cdot f(x_i)$

$\log(\text{softmax}(h_i, y)) = h_i - \log \sum_{j \neq i} \exp(h_j) \quad \frac{\partial \log(\text{softmax})}{\partial h_i} = \delta_{ij} - \left(\frac{1}{\sum_{j \neq i} \exp(h_j)} \right) \cdot \exp(h_i) = \delta_{ij} - \text{softmax}(h_i, i) \quad (\delta_{ik}=0 \Rightarrow i \neq k \quad \delta_{ik}=1 \Rightarrow i = k)$

$T = \text{Temp. Parameter} \mid T \rightarrow \infty \rightarrow \text{approach uniform categorical dist. (max entropy)} \quad \text{Interpolate softmax} > \text{argmax}$

$T \rightarrow 0 \rightarrow \text{annulling, all mass placed on maximum (min entropy) leading argmax as dist.}$

Exp. Family-specific dist. varies with parameter $|\ p(x|Q) = 1/Z(\theta) h(x) \cdot \exp(\theta \cdot \phi(x))$

Sigmoid: $\sigma(x) = \exp(x) / (1 + \exp(x)) = 1 / (1 + e^{-x}) = 1 / (1 + e^{-\theta}) \mid \text{Sigmoid is binary softmax}$

We can back non-linear feature function into log-linear, but we need to know function \rightarrow Neural Net.

Learn parameters of softmax($\theta \cdot f(x, y)$): f and θ : $h^2 = \sigma^2(w^2 \cdot \sigma^{-1}(w^1 \cdot x + b^1) + d^2) = \text{softmax}(h^1, y) \in (-1, 1)$

Function Engineering: $f(x, y) = [x_1, x_1^2, x_2, x_2^2, -1]^T$ Non linear activation function: $\tanh(x) = 2 \sigma(2x-1) / \text{relu}$

RNN: has time dependant hidden state, language model with infinite context | Cell: $h_t = f(h_{t-1}, e(y_{t-1}))$

Vanilla RNN: Elman Networks: incorporate arbitrary distant contextual information & treat word prediction

as a discriminative learning task: $p(y_t | y_{ct}) = p(y_t | y_1, \dots, y_{t-1}) = \exp(w_{yt} \cdot h_t) / \sum_{y' \in Y} \exp(w_{y't} \cdot h_t)$

Vanilla RNN: $h_t = \sigma(w_1 h_{t-1} + w_2 \cdot e(y_{t-1}))$ (= combination of hidden + current state)

\uparrow Non-linearity e.g. \tanh | Instead of fixed context considered at every step \Rightarrow context

Linguistic Phenomena: long distance dependencies: verb-subject agreement in German subordinate clauses

Language Modeling is structured prediction where we have strings, trees, graphs. It is the intersection of algos. and high-dim. statistics & is weighting of prefix-tree $| y' \in Y$ is extremely large 2^n

Voc V finite, Kleene Closure V infinite (infinite number of paths with finite lengths)

Local Normalization: choose weights of syntax-tree such that $\sum_{y' \in V^*} \prod_{t=1}^n \delta_{y'_t} = 1 \Rightarrow z = 1$

You multiply weights along a path and then sum all paths

N-gram assumption: condition only $n-1$ words: $p(w) = p(w_1) \cdot p(w_2|w_1) \dots p(w_n|w_1, \dots, w_{n-1}) \Leftrightarrow p(w_i | w_1, \dots, w_{i-1}) = p(w_i | U_{i-n+1}, \dots, w_{i-1})$

Bengio neural model: problem: discrete random variables: 10 consecutive words, $|V|=100,000 \Rightarrow 100,000^{10-1}$ params

\hookrightarrow Pass word embedding to MLP and reduce parameters to generalize better — feature function

1. Word embedding like skip gram $\rightarrow e(\text{history})$ | 2. MLP combines embeddings $h_t = f(e(\text{history}))$ | 3. Softmax for prob.

Recurrent Neural Language Model: $h_t = f(h_{t-1}, e(y_{t-1})) \mid p(y_t | y_{ct}) = \exp(w_{yt} \cdot h_t) / \sum_{y' \in Y} \exp(w_{y't} \cdot h_t)$

part-of-speech-tagging (=POS): determiner (the), noun, verb, etc. is like path search in graph

Score ($\langle \text{grammar tag sequence}, \text{sentence} \rangle$) | CRF: conditional probabilistic model for sequence labeling

\hookrightarrow higher bdd: can be any function or NN built on logistic regression classifier (log-lin model) of categories

$p(t|w) = \exp(score(t, w)) / \sum_{t' \in T} \exp(score(t', w)) \mid \text{Problem: exp. number of pos assignments } O(|T|^{|V|}) \mid \text{length of sentence}$

Use combinatorial structure of sentence: $score(t, w) = \sum_{n=1}^N score(t_{n-1}, t_n, w) \quad (\text{Every arc gets own score})$

Normalizer: $\sum_{t \in T} \exp(score(t_{n-1}, t_n, w)) = \sum_{t_0 \in T_0} \exp(score(t_0, t_1, w)) \dots \sum_{t_N \in T_N} \exp(score(t_{N-1}, t_N, w))$

Calculate normalizing constant $Z = \beta(w, t_0) \cdot 1 \cdot \beta(w, t_N) \leq 1 \quad (2. \text{ for } n \in N-1, \dots, 0: \quad (= \text{like backpropagation})$

3. $B(w, t_n) = \sum_{t_{n-1} \in T} \exp(score(t_{n-1}, t_n, w)) \cdot B(w, t_{n-1}) \quad (= \text{backward algorithm})$

Compute max path: Utterdi: 1. $\tau(w, t_N) \leq 1 \mid 2. \text{ for } n \in N-1, \dots, 0: \mid 3. \tau(w, t_n) \leq \max_{t_{n+1} \in T} \exp(score(t_n, t_{n+1}))$

Semiring: $R = \langle R^+, +, \times, 0, 1 \rangle: (a+b)+c = a+(b+c) \mid 0+a = a+0 = a \mid a+b = b+a \quad \mathcal{X}(w, t_{n+1})$

$(a \times b) \times c = a \times (b \times c) \mid 1 \times a = a \times 1 = a \parallel a \times (b+c) = (a \times b) + (a \times c) \mid (a+b) \times c = (a \times c) + (b \times c) \mid 0 \times a = a \times 0 = 0$

Utterdi: $[0, 1], \max, \times, 0, 1 \mid \text{Inside: } IR^+, \Sigma, \{+\}, \{0, 1\}, +, \times, 0, 1 \quad (= \text{Semiring } R = \langle IR^+, +, \times, 0, 1 \rangle)$

Normalizer: $Z = \beta(w, t_0): 1. \beta(w, t_n) = 1 \mid 2. \text{ for } n \in N-1, \dots, 0: \mid 3. \beta(w, t_n) = \oplus_{t_{n+1} \in T} \exp(score(t_n, t_{n+1})) \otimes \beta(w, t_{n+1})$

Highest Path: Same as above but with Semiring $R = \langle IR^+, \max, \times, 0, 1 \rangle \quad (= \text{max. this log-likelihood it's softmax})$

Parameters for score with backprop: max. log-likelihood: $\sum_{i=1}^I score(t_i, w_i) - \log \sum_{t' \in T} \exp(score(t', w_i))$

Constituent: word/group that function as single unit: [fruit flies] [like a green banana]

Context Free Grammar: rules can be applied to nonterminal regardless of context | grammars: set of rules

$G = \langle N: \text{non terminal symbols } N_1, \dots, S: \text{start } S, \Sigma: \text{alphabet of terminals } a_1, R: \text{rules } N \rightarrow \alpha \rangle \mid N PVP, \dots$

Probabilistic CFGs: prod. to each rule | multiply all probs. in tree to get prod. of tree

Locally Normalized: prod. in rule sums to 1: $S \rightarrow NP VP \mid N VP \mid N \rightarrow \alpha \Rightarrow \sum_{k=1}^K p(\alpha_k | N) = 1$

Weighted CFGs: globally normalized: $\exp(score(\text{tree}))$: just structured softmax: $p(t) = \frac{1}{Z} \prod_{r \in \text{tree}} \exp(score(r))$

Parsing: create tree from sentence | Chomsky Normal Form (CNF): $N_1 \rightarrow N_2 N_3, N \rightarrow a \quad Z \in \infty \sum_{t \in T} \exp(score(t))$

CKY-Algorithm: Normalizer of parser \hookrightarrow No longer infinite trees because $N \rightarrow N \rightarrow \dots \rightarrow \emptyset$

Runs in $O(N^3 \cdot \# \text{Rules})$ | Run in (\max, \times) semiring = score of best parse $\{0, 1, \Sigma, \{+, \times, 0, 1\}\}$ check string & grammar

Dependency grammar: link word with its syntactic head | Dependency tree: directed spanning tree: no cycles, one input arc
 Depend. vs. constituency: CFGs not provides info about syntactic head, dep. tree no info about constituency structure
 Projective depend. tree: no crossing arcs, closely related to constituency, DP algos for parsing like CKY
 Non-projective: crossing arcs, closely related to non-projective constituents

Non-projective: crossing arcs, closely related discontinuous constituents
 Prod. dist. over non-projective: $p(t|w) = \frac{1}{Z} \cdot \exp(score(t, w))$ | $Z = \sum_{\text{ET}(w)} \exp(score(t, w))$ ($n-1$) $^{n-2}$
 Edge factored assumption: $p(t|w) = \frac{1}{Z} \cdot \prod_{(i,j) \in t} \exp(score(i, j, w)) \cdot \exp(score(r, w))$ | $A_{ij} = \exp(score(i, j, w))$
 $Z = \sum_{\text{ET}(w)} \prod_{(i,j) \in t} A_{ij} \cdot \exp(r)$ | compute Score of tree score of root | $p_{ij} = \exp(score(i, j, w))$

Kirchhoff's Matrix Tree Theorem: # undirected spanning trees in undirected graph in $O(n^3)$ time. $L_{ij} = p_j$ if $i=1$, $L_{ij} = \sum_{k=1, k \neq i}^n A_{ijk}$ if $j=1$, $L_{ij} = -A_{ij}$ otherwise.

Transliteration: maps string from one character set to another: replace character with approx. phonetic equivalent

Weighted-finite-state-transducers: $T = \langle Q: \text{states}, \Sigma: \text{input vocab}, \Pi: \text{output vocab}, \lambda: \text{state to initial score}, \rho: \text{state to final score}, \delta: \text{transition} \rangle$
ambiguous: more than one accepted path $\prod(x \in \Sigma, y \in \Pi) \mid Wfst \text{ distr.}: p(y|x) = \frac{1}{2} \cdot \exp(score(y, x)) \neq \frac{1}{2} \cdot \exp(\sum_{n=1}^N score(r_n))$
+ structure assumption: $score(\prod: \text{seqn. of state transitions}) = \sum_{n=1}^N score(r_n) \mid \prod(x,y) = \text{all sequences } x \rightarrow y$

Normalizing constant: $Z = \sum_{y \in \text{all words}} \exp(score(y, x)) = \sum_{y' \in \text{all words}} \sum_{\pi \in \text{all paths } x \rightarrow y'} \exp(\sum_{n=1}^{|y'|} score(\pi_n))$

Floyd-Warshall Algo: computes all pairs shortest paths + no negative cycles in $O(n^3)$ (all paths $x \rightarrow y$)
 ↳ same as simple matrix multiplication with denoting $\langle \mathbb{R}^{+} \cup \{\infty\}, \min, +, +\infty, 0 \rangle$ with non-negative weights (tropical ring)
 1 for $n = 1 \rightarrow n/2$ for $i = 1 \rightarrow N/3$. sum $\leq \infty / 4$ for $m = 1 \rightarrow N/5$. sum = min(sum, $1/6 \cdot (sum_1 + sum_2 + sum_3) \cdot 10^7 / 16$). $w^2 / n^2 \cdot 10^7 = sum$

1. for $n: 1 \rightarrow N$ | 2. for $p: 1 \rightarrow N$ | 3. sum < too | 4. for $m: 1 \rightarrow N$ | 5. sum = min(sum, $w[n][m] + w[m][p]$) | 6. $w^2[n][p] = \text{sum}$
 w^k encodes the shortest path between two nodes of length k | all-pair shortest paths: $\min(w^0, w^1, \dots, w^{\# \text{nodes}})$

Floyd-Warshall as matrix mult: 1. for each vertex k in N : $W^K[i][j] = 0$ 2. for each vertex k in N : $\text{sum} = \min(\text{sum}, W^{K-1}[n,m] + U^1[m,p])$

$$W^* = \bigoplus_{K=0}^{\infty} W^K = I + W \bigoplus_{K=1}^{\infty} W^K = I + W(W^*) = I \Rightarrow (I - W)W^* = I \Rightarrow (I - W)^{-1} = W^* \quad (\text{Same as LU Decomp.})$$

$$Z = \sum_{y' \in \mathcal{Y}} \exp(score(y', x)) = \sum_{y' \in \mathcal{Y}} \sum_{T_{y'} \in \mathcal{T}(x, y')} \exp(\sum_{n=1}^{|T_{y'}|} score(T_n)) = \alpha^T (\sum_{w \in \mathcal{W}} \nu_w e_w) \cdot \beta$$

We optimize score weights with log-likelihood: $\sum_{i=1}^n \log(p(y_i|x_i)) = \sum_{i=1}^n \text{Score}(y_i|x_i) - \log Z(x_i)$

Sequence-to-Sequence Model: $\tilde{z} = \text{encoder}(x)$ | $y|x = \text{decoder}(\tilde{z})$ | $p(y|x) = \prod_{t=1}^T p(y_t|x, y_1, \dots, y_{t-1})$
 Train model params: $\text{argmax}_{\theta} \sum_{i=1}^N \log p(y_i|x_i; \theta) = \text{argmax}_{\theta} \sum_{i=1}^N \log(p(y_i|x_i, y_{<i}; \theta))$ (= LSTM)

The last hidden state of encoder is first hidden state of decoder $\Rightarrow y_1 = p(\cdot | x) \Rightarrow y_2 = p(\circ, x, y_1) \Rightarrow \dots$

Decoder only receives a single vector from decoder \Rightarrow information bottleneck \Rightarrow Attention Model
Attention: keep hidden vector of decoder at each time step. Weights for each hidden vector

Attention: keep hidden vector of decoder in each time step. weights for each hidden vector in encoder at each time-step (weights computed in query to key-values from training data)

attention weights: $\alpha_{i,j}^{(t)} = \exp(score(q_t, k_i)) / \sum_j \exp(score(q_t, k_j)) \Leftrightarrow \alpha^t = \text{softmax}(score(q_t, K))$
 context vector $c^t = \sum_i \alpha_{i,j}^{(t)} v_i \Leftrightarrow c^t = \alpha^{(t)T} \cdot V | k_i = v_i - h_i^t (\text{hidden state}) | q_t = h_t^d (= \text{decod. at position } t)$

Context vector $c = \sum_j \alpha_j v_j$ ($\Leftrightarrow c = \alpha^\top V$ | $K_1 \cdot v_i = h_i$ (=hidden state) | $q_t = h_t$ (=decod. at position t)
 $X = V + H^E$ (=Horizontally stacked encoder representation | Best translation $y_t = \operatorname{argmax}_y \delta y$ score(x_t, y))
 $\alpha = \delta c = \delta h_t = \delta(v_1^\top \alpha) = \delta v_1^\top \alpha = \delta v_1^\top (V + H^E) = \delta v_1^\top V + \delta v_1^\top H^E$

score(x, y) = $\log(p(y|x))$ | 1. $E_{x \sim p(x)} \log p(y_t|x_1, y_{t-1})$ (=Non-markovian structure) | $p(y_t|x_1, y_{t-1}, \dots, y_{t-2})$
e.g. language translation, image captioning, — $\hookrightarrow O(|\Sigma|^{n_{\text{max}}})$, summarization $\hookrightarrow O(|\Sigma|^{n_{\text{max}}}) \Rightarrow$ DP

e.g. language translation, image captioning, $\rightarrow O(\mathcal{E})$, summarization $\rightarrow O(\mathcal{E})$, vision $\Rightarrow O(\mathcal{D})$
 Probabilistic Model: make assumption about distr. of data and independence | Non-prob.: more interpretable no uncertainty
 Discrete-valued model: discrete latent variable ($O(\mathcal{D})$) | Continuous latent variable ($O(\mathcal{C})$)

Discriminative? model decision boundary (CRF, RNN) | Generative approach: model dist. (n-gram, Markov Random Field)
 Learned: Perceptron, SVM | Manually Crafted: Context-free grammar | Independence: allow DP_{global}, not needed for NN

Learned: Perception, SVM / Manually crafted: Context-free grammar (Independencies: allow NP_{NPST}, not needed for NN
 Locally normalized: Ngram, RNN: $\frac{1}{2} \cdot \text{score}(x_i, y_i)$ | Globally: CRF, MRF: $\frac{1}{2} \cdot \text{score}(y_i)$ not very realistic for some
 1 - $\frac{1}{2} \cdot \text{score}(x_i, y_i) + \frac{1}{2} \cdot \text{score}(y_i, x_i)$ End-to-end: beam search (faster than backtracking?) End-to-end

$L_1 = l(\theta) + \lambda \cdot \|A\|_1$: coeffs to be 0 | $L_2 = l(\theta) + \lambda \cdot \|B\|_2^2$: small non-zero | Intrinsic: does model improve final task? | problems
 Intrinsic NLP evaluation: based on criteria, more general e.g. log-likelihood | $p(\theta) \cdot p_1(D|\theta)$

$$\text{Cost NN} = \mathcal{D}((CN\lambda)^2) \cdot k_1 + \sum_{l=2}^{L-1} ((k_l+1) \cdot k_{l+1}) + (k_L+1) \cdot C^2 \quad | \quad \text{Prior: } p(\theta) = \frac{p(\theta) \cdot p(D|\theta)}{p(D)} \propto p(\theta) \cdot p(D|\theta)$$

Classification functions: $x_1 - x_2 < 0$, $|x_1| + |x_2| < 2$, $x_1^2 + x_2^2 < 4$ (Sigmoid $\sigma'(x) = \sigma(x) \cdot (1 - \sigma(x))$)
 $\tanh'(x) = 1 - \tanh^2(x)$ / $\text{ReLU}'(x) = 1$ if $x > 0$ else 0 / OR: $y = 1 \cdot \text{ReLU}(1 \cdot x_1 + 1 \cdot x_2) - 2 \cdot \text{ReLU}(x_1 + x_2 - 1)$ / when summing

$\tanh(x) = 1 - \tanh^2(x)$ / $\text{ReLU}'(x) = 1$ if $x > 0$ else 0 / $y = 1 \cdot \text{ReLU}(1 \cdot x_1 + 1 \cdot x_2) - 2 \cdot \text{ReLU}(x_1 + x_2 - 1)$ / "summing p(here or you) = p(you are here) / p(you are) | $z^{-1} = 0,5 | z^{-2} = 0,25 | z^{-3}, z^{-1} = 2^{-9} \Rightarrow \log$, summing: $t = \log(2^9 + 2^6)$

Semiring check: $\langle A, +, \cdot \rangle$ commutative monoid + $\langle A, \times, 1 \rangle$ monoid + 0 is annihilator for x when $x = + \leftarrow$ when mult. Viterbi: $V_N(t_N) = \max_{t_{N-1}} (\text{score}((t_{N-1}, t_N), w) + V_{N-1}(t_{N-1}))$ $O(T^2 N)$: $V_2(N) = \max(V_1(N) + \text{score}(\dots), V_1(v) + \text{score}(\dots))$

$V_{\text{iterdi}}: V_1(t_n) = \max_{t_{n-1}} (\text{score}(t_{n-1}, t_n) w) + V_{n-1}(t_{n-1})$ $O(T^2N)$: $V_2(N) = \max(V_1(N) + \text{score}(\dots), V_1(v) + \text{score}(\dots))$

α -conversion: rename variables: $\lambda x. (x y) \rightsquigarrow \lambda y. (y y)$ | β -reduction: avoid variables getting bound = outside \rightarrow inside

Termination issue: $f = \lambda x. ((x \times) x) : (F P) = (\lambda x. ((x \times) x)) F \rightarrow \dots$ /n-gram states: $m + |V|^{n-1} + 1$ (m =start state)
 Trigram transition function: $\sigma^-(q_1, \dots, w_i, q_n) = \log p(w_{m+1} = k \mid w_{m-1} = i, w_m = j)$ $m = |V|^{n-2}$ for $n \geq 2$, $0 \leq m \leq |V|^{n-1}$

Trigram transition function: $\pi(q_{ijt}, w_t | q_{ki}) = \log p(w_m=k | w_{m-1}=i, w_{m-2}=j) / m = 1/k$ ($n=2$, $0 \leq n < 1$)
 Initial weight function: $\lambda(q_{ij}) = \log \log(w_m=i | w_{m-1}=j, w_{m-2}=D)$, $w_{m-2} = D$ (softmax-

$$\frac{\exp(f(x_{ij}))}{\exp(f(x_{ij})) + \exp(f(x_{ij})) + \exp(f(x_{ij}))} = \frac{\exp(\theta_1 f(x_{ij}) - f(x_{ij}))}{1 + \exp(\theta_1 f(x_{ij}) - f(x_{ij}))}$$

$$\text{RMN: } h^+ = \text{softmax}(W_{hh} \cdot h^{t-1} + W_{hx} \cdot x^t + b) \quad | \quad \tilde{y}^+ = \text{softmax}(W_{ny} \cdot h^+ + b) \quad | \quad \text{Loss: } L = \sum_{i=1}^T L_i(y^+, y^{i+1})$$

$$\frac{\partial L}{\partial h^+} = \sum_{k=1}^n \frac{\partial L}{\partial h^+} \cdot \frac{\partial h^+}{\partial h_k} + \frac{\partial h^+}{\partial h_k} \cdot \frac{\partial L}{\partial h_k} = \prod_{i=k+1}^T \frac{\partial h^+}{\partial h_i} \cdot \frac{\partial h^+}{\partial h_k} = \prod_{i=k+1}^T \frac{\partial h^+}{\partial h_i} \cdot \frac{\partial h^+}{\partial h^{i-1}} / \text{Attention: } O(dNM)$$

Encoder-Decoder: $h_n^{(s)} = f(W_1^{(s)} \cdot h_{n-1}^{(s)} + W_2^{(s)} \cdot x_n) \Rightarrow h_n^{(s)} = h_{n-1}^{(s)} + f(W_2^{(s)} \cdot x_n)$

$$\text{CRF vs. MEMRA} = \text{locally vs. globally normalized Jacobian: } J = \frac{\partial f}{\partial x_i} = \frac{\partial f}{\partial h_i} \cdot \frac{\partial h_i}{\partial x_i}$$