

Backpropagation

Linear-time DP for derivatives:

1. Write composite fn as labeled acyclic hypergraph
2. Forward propagation with input
3. Backprop: $\frac{\partial y_i}{\partial x_j} = \sum_{p \in P(j,i)} \prod_{(k \rightarrow \ell) \in p} \frac{\partial z_\ell}{\partial z_k}$
 $\sin'(x) = \cos(x)$, $\cos'(x) = -\sin(x)$, $\log'(x) = \frac{1}{x}$, $\exp'(x) = \exp(x)$

Log-linear Modelling

score(y, x) = $\theta^\top f(x, y)$
NLL gradient = 0: $\sum_{i=1}^n f(x_i, y_i) = \sum_{i=1}^n \mathbb{E}_{y|x_i, \theta}[f(x_i, y)]$

Hessian: $H_\theta(\sum_i -\log p(y_i|x_i)) = \sum_i \text{Cov}_{y|x_i, \theta}[f(x_i, y)]$

Softmax: softmax(\mathbf{h}_y) = $\frac{\exp(h_y/T)}{\sum_{y'} \exp(h_{y'}/T)}$ $T \rightarrow 0$: argmax. $T \rightarrow \infty$: uniform.

Exponential family: $p(x|\theta) = \frac{1}{Z(\theta)} h(x) \exp(\theta^\top \varphi(x))$

Multi-layer Perceptron

Problem: Data must be linearly separable. **Solution:** Learn non-linear feature fn with MLP:

$$h_k = \sigma_k(\mathbf{W}_k^\top \mathbf{h}_{k-1}), \mathbf{h}_1 = \sigma_1(\mathbf{W}_1^\top e(x))$$

Then softmax($\hat{\theta}^\top \mathbf{h}_n$) for prob dist.

Skip-Gram: predict if 2 words in same context. Need good word repr.

Derivative: $\frac{\partial \ell}{\partial \mathbf{W}_k} = \frac{\partial \ell}{\partial y} \frac{\partial y}{\partial \mathbf{h}_n} (\prod_{m=k+1}^n \sigma'_m(\dots) \mathbf{W}_m) \sigma'_k(\dots) \mathbf{h}_{k-1}$

Structured Prediction

$$p(y|x) = \frac{\exp(\text{score}(y, x))}{Z(x)}, Z(x) = \sum_{y' \in \mathcal{Y}} \exp(\text{score}(y', x))$$

Problem: \mathcal{Y} exponentially/ininitely large. **Solution:** Design algorithms using structure of input/output.

Language Modelling

$$p(y) = p(\text{eos}|y) \cdot \prod_{i=1}^N p(y_i | y_{<i}) p(y_i|y_{<i}) = \frac{1}{Z(y_{<i})} \exp(\text{score}(y_{<i}, y_i))$$

Non-tight: Force $p(\text{eos}|y_{<i}) > \xi > 0$

n-gram: $p(y_i|y_{<i}) = p(y_i|y_{i-n+1}, \dots, y_{i-1})$ **Neural n-gram:** Embeddings + MLP

RNN: $h_i = \sigma(\mathbf{W}_h h_{i-1} + \mathbf{W}_x e(y_{i-1}) + b)$

Vanishing gradient: LSTM/GRU

Semirings

Definitions: Monoid $\langle \mathbb{K}, \circ, e \rangle$: assoc, identity Semiring $\langle \mathbb{K}, \oplus, \otimes, \mathbf{0}, \mathbf{1} \rangle$: comm monoid, monoid, distrib, annihilator

Closed: $x^* = \bigoplus_{n=0}^{\infty} x^{\otimes n}$

Boolean, Viterbi $\langle [0, 1], \max, \times, 0, 1 \rangle$, Inside, Real, Tropical, Log, Expectation,

Counting

Part-of-Speech Tagging

Input: $w \in \Sigma^N$. Output: $t \in \mathcal{T}^N$.

CRF: $\text{score}(t, w) = \sum_{n=1}^N \text{score}((t_{n-1}, t_n), w, n) = \text{trans}(t_{n-1}, t_n) +$

emit(w_n, t_n)

Forward: $\alpha_{n, t_n} \leftarrow \bigoplus_{t_{n-1} \in \mathcal{T}} \text{exp}(\text{score}(\dots)) \otimes \alpha_{n-1, t_{n-1}}$ Return $\alpha_{N, \text{eot}}$. Runtime:

$O(N|\mathcal{T}|^2)$

Dijkstra: $O(N|\mathcal{T}|^2 + N|\mathcal{T}| \log(N|\mathcal{T}|))$

Finite-State Automata

WFST: $\Sigma, \Omega, Q, I \subseteq Q, F \subseteq Q, \lambda : I \rightarrow \mathbb{K}, \rho : F \rightarrow \mathbb{K}, \delta$

Pathsum: $Z(\mathcal{T}) = \bigoplus_{i, k \in Q} \lambda(q_i) \otimes R_{ik} \otimes \rho(q_k)$

Lehmann: $R_{ik}^{(j)} \leftarrow R_{ik}^{(j-1)} \oplus R_{ij}^{(j-1)} \otimes (R_{jj}^{(j-1)})^*$ $\otimes R_{jk}^{(j-1)}$ Runtime: $O(|Q|^3)$

Composition: $\mathcal{T}(x, y) = \bigoplus_{z \in \Omega^*} \mathcal{T}_1(x, z) \otimes \mathcal{T}_2(z, y)$

Transliteration

Map $\Sigma^* \rightarrow \Omega^*$. Three transducers:

1. \mathcal{T}_x : maps $x \rightarrow x$

2. \mathcal{T}_θ : maps $\Sigma^* \rightarrow \Omega^*$

3. \mathcal{T}_y : maps $y \rightarrow y$

Compose for $Z(x)$ and score(y, x)

Constituency Parsing

CFG: $N, S, \Sigma, \mathcal{R}$ (rules $N \rightarrow \alpha$) PCFG: locally normalized. WCFG: globally normalized.

CNF: $N_1 \rightarrow N_2 N_3$ or $N \rightarrow a$ (no cycles)

CKY: $C_{i,k,X} \leftarrow \bigoplus_{X \rightarrow YZ} \text{exp}(\text{score}(X \rightarrow YZ)) \otimes C_{i,j,Y} \otimes C_{j,k,Z}$ Return

$C_{1,N+1,S}$. Runtime: $O(N^3 |\mathcal{R}|)$

Dependency Parsing

$(N-1)^{N-2}$ spanning trees with single-root.

score(t, w) = $\rho_r + \sum_{(i \rightarrow j) \in t} A_{ij}$

Koo MTT: Laplacian $L_{ij} = \begin{cases} \rho_j & \text{if } i=1 \\ -A_{ij} & \text{if } i \neq j \\ \sum_{k \neq i} A_{kj} & \text{otherwise} \end{cases}$ $Z(w) = \det(L)$. Runtime: $O(N^3)$

Chu-Liu-Edmonds: Greedy graph \rightarrow contract cycles \rightarrow swap loss \rightarrow expand.

$O(N^2)$

Semantic Parsing

Lambda calculus: $x, y, z; (\lambda x. f(x)); (MN) \beta\text{-reduction}: ((\lambda x. M)N) \rightarrow M[x := N]$

$\alpha\text{-conversion}: \lambda x. M[x] \rightarrow \lambda y. M[y]$

CCG rules: $X/Y Y \Rightarrow X (>), Y X \setminus Y \Rightarrow X (<) X/Y Y/Z \Rightarrow X/Z (B_>) X \Rightarrow T/(T \setminus X) (T_>)$

LIG: CFG with stacks. Push/pop rules.

Transformers

Self-attention: Learn $\mathbf{W}_Q, \mathbf{W}_K, \mathbf{W}_V \in \mathbb{R}^{d \times d}$ SelfAtt(X) =

softmax $\left((\mathbf{W}_Q^\top X)^\top \frac{\mathbf{W}_K^\top X}{\sqrt{d_q}} \right) (\mathbf{W}_V^\top X)^\top$ Runtime: $O(nd^2 + dn^2)$

Positional encoding: $P_{pi} = \sin(p/10000^{i/d})$ or cos

Encoder: $\oplus P \rightarrow \text{MHS} \rightarrow \oplus \rightarrow \text{LN} \rightarrow \text{MLP} \rightarrow \oplus \rightarrow \text{LN}$ **Decoder:** + linear + softmax

Beam search, nucleus sampling

Axes of Modelling

Bias-variance: High bias = underfit. High variance = overfit. **Regularization:**

$\ell(\theta) + \lambda \|\theta\|_2^2$

MLE: $\hat{\theta} = \arg \min_{\theta} -\log \prod_{(x,y) \in \mathcal{D}} p_\theta(y|x)$

Precision = TP/PP, **Recall** = TP/(TP+FN), **F1** = $2 \cdot \frac{P \cdot R}{P + R}$

Locally norm: efficient, label bias **Globally norm:** needs normalizer

Tips

Gradient: Sum over paths, product within paths. **Reuse terms** in backprop for efficiency.

Complexities: vec-vec $O(d)$, mat-vec $O(nm)$, mat-mat $O(nm\ell)$

Activations:

• $\sigma(x) = \frac{1}{1+\exp(-x)}$, $\sigma'(x) = \sigma(x)(1-\sigma(x))$

• ReLU(x) = $\max\{0, x\}$, ReLU'(x) = $1\{x > 0\}$

• tanh(x) = $\frac{\exp(x)-\exp(-x)}{\exp(x)+\exp(-x)}$, tanh'(x) = $1 - \tanh^2(x)$