

LANGUAGE MODELING

Alphabet Σ finite, non-empty set of symbols Vocab words / tokens $\mathcal{E} \notin \Sigma$

String over Σ is a finite sequence of alphabet

$$\text{Kleene Closure } \alpha^* = \bigoplus_{n=0}^{\infty} \alpha^{\otimes n} = \alpha^{\otimes 0} \oplus \alpha^{\otimes 1} \oplus \bigoplus_{n=2}^{\infty} \alpha^{\otimes n} = 1 \oplus \alpha \otimes \bigoplus_{n=0}^{\infty} \alpha^{\otimes n}$$

Language Model Distribution over Σ^* Global Z is infinite sum over $y \in \Sigma^*$

Locally Normalized (Auto-regressive) Collection of conditionals $p(y|y)$

$$\text{String Probability } p(y) = p(y_1 | \text{BOS}) p(y_2 | \text{BOS } y_1) \cdots p(y_N | y_{N-1}) p(\text{EOS} | y)$$

One-hot Encoding Vocab V , $e(w) \in \mathbb{R}^{|V|}$ BagOfWords Sum up one-hot embeddings

n-grams Encode n-grams, $|V| \rightarrow |V|^n$ Pre Processing Tokenization / Lemmatization / Punct Removal

Skip Gram Sentence tokenize corpus, then Build Positive Samples with given

Context window K. Results $N = O(K \cdot C)$, $e_{\text{word}} \in \mathbb{R}^d$ $e_{\text{ctx}}(c) \in \mathbb{R}^d$ log-bilinear model

$$p(c|w) = \gamma_w \exp(e_{\text{word}}(w) \cdot e_{\text{ctx}}(c)) \text{ Training objective: Maximize NLLK}$$

n-gram Assumption $p(y_t | y_{t-1}, \dots, y_{t-n+1})$ Context depends on prev n-1

\hookrightarrow One conditional probability for each possible context $\sum_{j=0}^{n-1} |\Sigma|^j$

\hookrightarrow $|\Sigma|^{n-1}$ free parameters for each conditional probability distribution

Bengio's Model Estimate $p(y|y)$ more efficiently. Share parameters between contexts

$y_{t-1}, \dots, y_{t-n+1}$, Instead of treating words as discrete symbols, treat as embeddings

$$p(y_t | y_{t-1}, \dots, y_{t-n+1}) = \frac{\exp(v(y_t) \cdot h_t)}{\sum_{y \in \Sigma} \exp(v(y) \cdot h_t)} \text{ Word Embeddings Context "Summary"}$$

Advantage: Naive Ngram No long-term dependencies Cons: Slow (softmax $\Sigma \forall$ conditionals)

Solve Curse of Dimensionality 1. Word Embeddings 2. MLP to Combine Embeddings 3. Softmax

$$\text{RNN } h_t = \sigma(Uh_{t-1} + Vv(y_{t-1}) + b_h) \quad p(y_t | y_{t-1}, \dots, y_{t-n+1}) = \exp(v(y_t) \cdot h_t) / \sum_{y \in \Sigma} \exp(v(y) \cdot h_t)$$

CONDITIONAL RANDOM FIELDS

$$p(t|w) = \frac{\exp\{\sum_{n=1}^N \text{score}(t_{n-1}, t_n, w)\}}{\sum_{t' \in \Sigma^N} \exp\{\sum_{n=1}^N \text{score}(t_{n-1}, t_n, w)\}} \quad \text{Naive Z: } O(|T|^N)$$

$$Z = \sum_{t \in \Sigma^N} (\exp\{\sum_{n=1}^N \text{score}(t_{n-1}, t_n, w)\}) \times \prod_{n=1}^N \left(\sum_{t_n \in \Sigma} (\exp\{\sum_{r \in R} \text{score}(t_{n-1}, t_n, r, w)\}) \right) \times \prod_{n=1}^N \exp\{\sum_{r \in R} \text{score}(t_{n-1}, t_n, r, w)\}$$

Viterbi Algorithm Choose Semiring $R = \langle IR^+, +, \times, 0, 1 \rangle$ and compute

$$\beta(w, t_N) = 1 \quad \text{Works because of the Distributive Property } O(T^2 N)$$

for $n=N-1, \dots, 0$: Correctness: Just Prove Semiring Axioms

$$\beta(w, t_n) = (\bigoplus_{t_{n+1} \in \Sigma} \exp\{\sum_{r \in R} \text{score}(t_n, t_{n+1}, r, w)\}) \otimes \beta(w, t_{n+1}) \quad 0 \oplus a = a$$

Semirings $R = \langle A, \oplus, \otimes, 0, 1 \rangle$ Associativity $(a \oplus b) \oplus c = a \oplus (b \oplus c)$ Identity $a \oplus 0 = a$

Commutativity $a \oplus b = b \oplus a$. \oplus is a commutative Monoid \otimes is a Monoid

Distributivity $(a \oplus b) \otimes c = a \otimes c + b \otimes c$ and $a \otimes (b \oplus c) = a \otimes b \oplus a \otimes c$ Axiom 4 $0 \otimes a = 0$

Dijkstras $O(|W| |T|^2 \log |W| |T|)$ Emission $t \mapsto w$ Transition $t \mapsto t'$

Semirings BOOLEAN $\langle \{0, 1\}, \vee, \wedge, 0, 1 \rangle$ logical deduction, recognition

Viterbi $\langle [0, 1], \max, +, 0, 1 \rangle$ prob. of best derivation

Inside $\langle IR^+ \cup \{+\}, +, \times, 0, 1 \rangle$ prob. of a string Real $\langle IR^+ \cup \{+\}, +, \times, 0, 1 \rangle$ distance

TROPICAL $\langle IR^+ \cup \{+\}, \min, +, \infty, 0, D \rangle$ Shortest Distance with Non-Negative Weights

Counting $\langle \mathbb{N}, +, \times, 0, 1 \rangle$ Number of paths

WEIGHTED FINITE STATE AUTOMATA / TRANSDUCERS (WFSA, WFST)

WFST T over $W = (\mathbb{K}, \oplus, \otimes, 0, 1)$ is an 8-tuple $(\Sigma, \Delta, Q, I, F, \delta, \pi, \rho)$ Where

Σ is a finite input alphabet Δ is a finite output alphabet Q is a finite set of states

$I \subseteq Q$ is a set of Initial states $F \subseteq Q$ is a set of Final states $\delta \subseteq \mathbb{K} \times \Sigma \times \Delta \times Q \times Q$

$\pi: Q \rightarrow \mathbb{K}$ a init weighting function over the set of states $\rho: Q \rightarrow \mathbb{K}$ a final weighting function over set of states Q

$$I = \{q \in Q \mid \pi(q) \neq 0\} \text{ and } F = \{q \in Q \mid \rho(q) \neq 0\}$$

Lehmann's Algorithm Solves all pairs shortest paths given no negative cycles $O(|V|^3)$

$$W \leftarrow |Q| \times |Q| \text{ matrix over a closed semiring } (I + M)^k = \bigoplus_{n=0}^k \binom{k}{n} I^{k-n} \otimes M^n = \bigoplus_{n=0}^k \binom{k}{n} M^n$$

$$R^{(0)} \leftarrow W$$

for $j \leftarrow 1$ up to $|Q|$:

for $i \leftarrow 1$ up to $|Q|$:

$$\|A^* - \sum_{n=0}^k A^M\|_2 \leq \left\| \sum_{n=k+1}^{\infty} A^n \right\|_2$$

$$R_{ik}^{(j)} \leftarrow R_{ik}^{(j-1)} \oplus (R_{ij}^{(j-1)})^* \otimes R_{jk}^{(j-1)}$$

$$\text{return } I \oplus R^{(|Q|)}$$

Composition Let $T_1: U \rightarrow V$ and $T_2: V \rightarrow W$, $T_1 \circ T_2: U \rightarrow W$

$$T(x, y) = \bigoplus_{z \in \Sigma^*} T_1(x, z) \otimes T_2(z, y) \quad \text{Remember to Relabel } \Sigma \text{ transitions}$$

CONSTITUENCY PARSING

CFG $G = \langle N, S, \Sigma, R \rangle$ $N \setminus \{S\}$ non-terminal symbols $\{N_1, N_2, N_3, N_4, \dots\}$ S start symbol

Σ Alphabet of terminal Symbols α , R production Rules $N \rightarrow \alpha$ $N \in N \cup \Sigma$

Chomsky Normal Form Grammar in CNF is all productions form of $N_1 \rightarrow N_2 N_3$ $N \rightarrow a$

Locally Normalized Distribution over each production is a valid probability distribution

Weighted CFG Log-linear models over trees in grammar, structured softmax

$$p(t) = \frac{1}{Z} \prod_{r \in t} \exp\{\text{score}(r)\} \quad Z = \sum_{t \in T} \prod_{r \in t} \exp\{\text{score}(r)\} \quad T \text{ is infinite!}$$

Divergence $S \rightarrow S(1.0) \rightarrow \alpha(1.0)$, infinite sum because each summand is 1

CKY (max, \times) Semiring \rightarrow Score of one-best Parse, (Bool) \rightarrow Checks if string in Grammar

Compute $Z(S)$ in $O(N^3(|R|)) \rightarrow \nabla \log Z(S)$ in $O(N^3(|R|))$ Backprop!

DEPENDENCY PARSING

Projective No overlapping Arcs Non-projective Overlapping Arcs

Probability Distr Non-Projective $Z \propto N^N$ Naively $O(N^N)$ Spanning Trees $N^{N-2} + \text{root constraint } (N-1)^{N-2}$

Edge-factored Assumption $p(t|w) = \frac{1}{Z} \prod_{(i,j) \in t} \exp\{\text{score}(i,j; w)\} \exp\{\text{score}(r, w)\}$

$$Z = \sum_{t \in T(w)} \prod_{(i,j) \in t} \exp\{\text{score}(i,j; w)\} \exp\{\text{score}(r, w)\}$$

MTT For an undirected unweighted graph G with N vertices, let L be the graph

Laplacian $L_{ij} = \begin{cases} -A_{ij} & i \neq j \\ p_j + \sum_{k \neq i} A_{ik} & i=j \end{cases}$ Tutte Generalized to Directed Case. Find Z in $O(N^3)$

CLE ρ pseudoroot \rightarrow Incorporate root scores. Need node that can reach all other

nodes in the graph. Greedy Select best incoming edge to each node except ρ , if cycle

Contract and Repeat. Otherwise, Constrain Root to only 1 edge by Comparing cost. $N_{DT}(G) = NN_T(G)$, at root $1 = NN^{-2} (L_i - \lambda I) V^{(P, q)} = 0 \quad \sum_i \delta_{ii} V_i \geq 0$

$x^T Ax \geq 0$, square \Rightarrow PSD Contract: Update each enter edge to cycle by adding

the weights including enter edge but not including edge we break. $O(EN) = O(N^3)$

Tarjan $\rightarrow O(E + N \log N)$

SEQ2SEQ MODELS

y at Inference, generate y_1 according to $p(\cdot|x)$, y_2 according to $p(\cdot|x, y_1)$

Attention Weights $\alpha_i = \text{Softmax}(\text{score}(q, k))_i$ Context $c = \sum_i \alpha_i v_i = \alpha v$

$k_i = v_i = h_i^{(e)}$ Vector representation [hidden state] produced by encoder at

position i $q_t = h_t^{(d)}$ hidden state produced by decoder position t

$K = V = H^{(e)}$ Vertically stacked encoder vector representations

$\alpha = \text{Softmax}(\text{Score}(q, K))$ $MHSA(z) = \sum_{h \in H} \text{Softmax}(A^{(h)} z W_v^{(h)}) w^{(h)}$

Instead of encoding all input into a single context vector, pay variable attention depending on output generation step. MHSA: Learn multiple sets of attention weights for same input then concat.

Scaled: $\text{Score}(q_n, h_t) = \text{Softmax}((q_n^T h_t) / \sqrt{d_k})$ \vee PE Encode order of words

RESIDUAL CONNECTIONS Mitigate Vanishing Gradients.

Layer Norm Helps with internal covariate shift, normalize individual layer inputs INPUT \rightarrow PE \rightarrow $N \times [\text{MHSA}, \text{LayerNorm}, \text{FF}, \text{LayerNorm}]$

TRANSFORMER OUTPUT \rightarrow PE \rightarrow $N \times [\text{MHSA}, \text{LN}, \text{MHSA}, \text{LN}, \text{FF}, \text{LN}]$ - Linear-Softmax

DECODING $y^* = \text{Argmax} \text{Score}(x, y)$. Beam Search, Sampling where words are sampled iteratively from $p(y_t|x, y_{\leq t})$ Nucleus Sampling, only top words p^*

SEMANTIC PARSING

Compositionality Meaning of Complex expression is a function of the meanings of that expressions Constituent Parts. \rightarrow Build Models over Components

$\lambda x.f(x)$ Function that takes in x as input and $f(x)$ as output.

λ -Conversion $\lambda x.(xx)y \rightarrow \lambda t.(tt)y$ Rename Var in abstraction with all bound vars

β -Reduction Apply lambda term to another $(\lambda x.\lambda y.(x((\lambda x.z)x)y))z \rightarrow \lambda y.z((\lambda x.x.z)y)$

AXES OF MODELING

Implications of Structural Constraints Determine Complexity of the model.

A model with too many constraints (too simple) might underfit the data.

Too few constraints (too complex) \rightarrow Overfit. Domain specific structures, Sequential data.

Cost/Benifits of Independence Assumptions Simplifies computation and model understanding. Benifical with high dimensional data. Loss of Accuracy.

Regularization Increases models generalization performance.

Considerations Data Balance, Robustness to Noise, Interpretability, Transferability

Inflated Confidence Overfitting to training data, Data Leakage, Improper Cross Val, Ignoring Model assumptions, Selective Reporting of Results.

OTHER IMPORTANT STUFF

Log-Linear Models Very general family of probability distributions that can be defined as $p(y|x, \theta) = Y_{z(\theta)} \exp(\theta \cdot f(x, y))$ If we take log \rightarrow Linear function Can be thought as dot-product followed by the Softmax

Softmax T : Non-negative temperature. $T \rightarrow \infty$ Uniform (Maximum Entropy)

$T \rightarrow 0$ (Annealing) all mass placed on maximum (minimum Entropy)

Exponential Family $p(x|\theta) = Y_{z(\theta)} h(x) \exp(\theta, \phi(x))$ generalizes Softmax

FF Neural Network Connection between nodes dont form Cycles.

$\nabla L(\theta) = -\sum_{n=1}^N \log p(y_n|x_n, \theta) = -\sum_{n=1}^N f(x_n, y_n) - \sum_{n=1}^N \sum_{y' \neq y_n} p(y'|x_n, \theta) \cdot f(x_n, y')$

Expectation Matching $\nabla L(\theta) = 0$ Observed features counts = Expected Feature Counts

$p(w:k|x) = \sum_{Y, Z \in N} p(x \rightarrow YZ) \left[\sum_{j=1}^{k-1} P_{\text{inside}}(w_{-j}|Y) P_{\text{pre}}(w_{j+1:k}|Z) + P_{\text{pre}}(w:k|Y) \right]$

$M^* = I + M + M^2 + \dots = \sum_{n=0}^{\infty} M^n = I + M \sum_{n=0}^{\infty} M^n = I + MM^* = (I - M)^{-1}$

Matrix Inversion $O(N^3) \Rightarrow P_{\text{lc}}(YZ|x) \in O(N^4)$ $P_{\text{inside}}(w_{i:j}|Y)$ by CKY

$\in O(N^3 | R_l|^3) = O(N^3 | N^l |^3)$ $P_{\text{pre}}(w_k|x) \in O(N^3 | N^k |)$ $P_{\text{pre}}(w) \in O(N^3 | N^l |^3 + N^l | N^k |)$

All prefix probability $\rightarrow N$ times $\in O(N^4 | N^l |^3 + N^l | N^k |)$