

2 - BACKPROPAGATION

BACKPROPAGATION IS A LINEAR TIME DYNAMIC PROGRAM FOR COMPUT. ∂_s ; AKA "REVERSE-MODE AUTO-DIFFERENTIATION". GOAL: FIT FUNCTION TO GIVEN DATA BY MINIMIZING $\sum (x_i, y_i) L(f(x_i), y_i)$ (REQUIRES ∇f). BACKPROP CAN COMPUTE A GRADIENT IN THE SAME TIME COMPLEXITY AS COMPUTING f , EXPLOITING THE COMPOSITE NATURE OF COMPLEX FUNCTIONS

$$\frac{\partial y}{\partial x} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}, \quad \frac{\partial y}{\partial x_i} = \left[\frac{\partial y}{\partial x_1}, \dots, \frac{\partial y}{\partial x_m} \right]; \quad \frac{\partial y}{\partial x} = \begin{bmatrix} \frac{\partial y}{\partial x_1} & \dots & \frac{\partial y}{\partial x_m} \\ \vdots & \ddots & \vdots \\ \frac{\partial y}{\partial x_m} & \dots & \frac{\partial y}{\partial x_1} \end{bmatrix} \quad \text{JACOBIAN } (m \times m) \quad \frac{\partial y_i}{\partial x_j} = \sum_{k=1}^m \frac{\partial y_i}{\partial x_k} \cdot \frac{\partial x_k}{\partial x_j}$$

COMPUTATION GRAPHS: LABELED, ACYCLIC HYPERGRAPH WITH VARIABLES AS NODES AND FUNCTION-LABEL EDGES; CAN REPRESENT ANY COMPOSITE FUNCTION; GIVEN A SET OF PRIMITIVES AND THEIR DERIVATIVES

BAUER'S FORMULA: $\frac{\partial y_i}{\partial x_j} = \sum_{k=1}^m \frac{\partial y_i}{\partial z_k} \frac{\partial z_k}{\partial x_j}$ COMPUTING BAUER WITH DP GOES FROM EXPONENTIAL TO LIN.

$\text{PEPL}(i,j) \rightarrow \text{BAUER PATHS}_{\text{PEPL}(i,j)}$ IN THE NUMBER OF NODES OF THE COMPUTATION GRAPH

FORWARD PROP (F, x):

```

 $z_i = x_i$  IF  $i \in M$  ELSE 0
 $z = \text{FORWARD PROP}(F, x)$ 
FOR i = m+1:N
     $z_i = g_i(\langle z_{\leq i-1} \rangle)$ 
     $\frac{\partial F}{\partial z_i} = \sum_{j \in P(i)} \frac{\partial F}{\partial z_j} \frac{\partial z_j}{\partial z_i}$ 
RETURN  $[z_1, \dots, z_n]$ 

```

$y^x = \exp(x \log(y))$; SUM OVER PATHS OF GRAPH!

COMPUTING THE HESSIAN OF A m -VALUED FUNCTION TAKES $O(m \cdot \# \text{NODES GRAP})$

BACUPROPAGATION IS A CONSTRUCTIVE THEOREM SINCE IT GIVES A THEORETICAL GUARANTEE ON THE RUNTIME OF THE BACKWARD PASS AND ALSO AN ALGORITHM TO ACTUALLY CALCULATE IT

C. > IS AN ORDERED SET

3 - LOG-LINEAR MODELS, SOFTMAX

PROBABILITY SPACE P : (SAMPLE SPACE S , EVENT SPACE E , PROBABILITY FUNCTION P) $\xrightarrow{\text{P}} E \in \Omega$

RANDOM VARIABLE: FUNCTION $\Omega \rightarrow T$, TARGET SPACE; $P(x|y) = P(x,y)/P(y)$; $E[\ell(x)] = \sum_x \ell(x) P(x)$

LOG-LINEAR MODEL: $P(y|x, \theta) = \frac{\exp(\theta \cdot F(x, y))}{\sum_{y' \in Y} \exp(\theta \cdot F(x, y'))}$ FEATURE ENG: PREPROCESSING + FEATURE DESIGN

MLE ESTIMATION: $\theta_{MLE} = \arg\max_{\theta} \sum_{(x,y)} \log P(y|x, \theta)$ OBSERVED FEATURE COUNTS: DATASET FEATURES

ARE CONSISTENT: \exists COMPACT SUBSET $\Theta \subset \mathbb{R}^n$ $\sum_{(x,y)} p(y|x, \theta) = \sum_{(x,y)} \sum_{k=1}^m p(y^k|x, \theta) F_k(x, y)$ EXPECTED

$$\Delta \text{NLL}(\theta) = - \sum_{(x,y)} \left(p_k(x, y) - \sum_{y' \in Y} p(y'|x, \theta) \cdot F_k(x, y') \right) = - \sum_{i=1}^m F_k(x_i, y_i) + \sum_{i=1}^m \sum_{y \in Y} p(y^i|x_i, \theta) F_k(x_i, y)$$
 MATCHING

SOFTMAX (\bar{h}, y, T) = $\frac{\exp(\bar{h}_y/T)}{\sum_{y' \in Y} \exp(\bar{h}_{y'}/T)}$ T: TEMPERATURE HYPERPARA. $T \rightarrow \infty \Rightarrow$ UNIFORM CAT. $T \rightarrow 0 \Rightarrow \text{MAX}(\cdot)$ $\bar{h} \mapsto \text{SIMPLEX } \Delta^{n-1}$ LOG-LIN.: DOT PRODUCT FOLLOWED BY SOFTMAX $\bar{h}_y = \theta \cdot \ell(x, y)$

$\delta \log \text{SOFTMAX}(\bar{h}, y) = \delta y_i - \text{SOFTMAX}(\bar{h}, i)$ GUMBLE SOFT MAX: CONTINUOUS RELAXATION OF SOFTMAX FROM δh_i WHICH WE CAN SAMPLE AND BACKPROPAGATE THROUGH

EXPONENTIAL FAMILY: $p(x|\theta) = \frac{1}{Z(\theta)} h(x) \exp(\theta \cdot \phi(x))$ $Z(\theta)$: PARTITION FUNCTION; $h(x)$: SUPPORT; θ : PARAMETERS
SCALED EXPONENTIALIZED DOT PROD.

$$N_{\theta, \sigma} = \frac{1}{\sqrt{2\pi}} \cdot \exp\left(-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right) = \frac{1}{\sqrt{2\pi\sigma^2} \exp(\mu^2/\sigma^2 - 1)} \cdot 1 \cdot \exp\left(-\frac{1}{2\sigma^2} \left(\frac{x-\mu}{\sigma}\right)^2 \cdot [x, x^2]^T\right)$$

CONJUGATE PRIOR: (LIKELIHOOD, PRIOR) PAIR s.t. POSTERIOR = PRIOR; AVOIDS COMPUTING INTRACT. INTEGRAL EXAMPLE: (CATEGORICAL, DIRICHLET); ENTROPY: $H(p(x)) = -\sum x_i p(x_i) \log p(x)$

MAXIMUM ENTROPY PRINCIPLE: THE DISTRIBUTION WHICH BEST REPRESENTS THE CURRENT (PARTIAL) STATE OF KNOWLEDGE IS THE ONE WHICH MAXIMIZES $H(p(x))$; EXPONENTIAL DISTRIBUTIONS MAXIMIZE ENTROPY IN SPECIFIC COND.

LOG-LINEAR \leftrightarrow EXP. FAM: $\log p(y|x) = \theta^T \phi(y) - \log Z(\theta) + \log h(y)$; $\theta(x) = [\theta^T(F(x_1) - F(x, C)), \dots, \theta^T(F(x_n) - F(x, C))]^T$

LOG-LINEAR MODELS DEFINE A FUNCTION FROM VECTORS TO A CATEGORICAL DISTRIBUTION \rightarrow (IN EXP. FAMILY)

CHI-SQUARE AS EXPONENTIAL: $p(x|k) = \frac{1}{2^{k/2} \Gamma(k/2)} \exp(-\frac{k}{2}) \exp((\frac{k}{2}-1) \log(x))$

$\theta(x) = \exp(-x/2)$ $\phi(x) = \log(x)$ $\phi(x) = \log(x)$ $\hookrightarrow x = \exp(\log(x))$ (UNIFORM NOT IN EXP.)

BINARY LOG-LINEAR MODEL IS LOGISTIC REGRESSION WHEN $\theta^0 = [\theta_0, \theta_1, \dots, \theta_d] = \theta_0 - \theta$

START FROM $P(y=1|x, \theta)$ AND GET TO $\sigma(\theta^T x)$ ALGEBRAICALLY | $\theta \in \mathbb{R}^{d+1}$, $F(x, \theta) = [x_1, \dots, x_d, 1, \dots, 0]$

HESSIAN OF LL (LOG-LIN) = $\frac{\partial^2 \text{LL}}{\partial \theta^2} = \sum_{i=1}^n \sum_{y \in Y} p(y|x, \theta) F(x, y)^T - p(y|x, \theta) F(x, y) \sum_{y \in Y} p(y|x, \theta) F(x, y)^T$

= $\sum_{i=1}^n \sum_{y \in Y} p(y|x, \theta) F(x, y) F(x, y)^T - \sum_{i=1}^n (\sum_{y \in Y} p(y|x, \theta) F(x, y) F(x, y)^T) = \sum_{i=1}^n \text{Cov}(F(x, \theta))$

4 - MLPs

HYPERSPACE: MULTIDIMENSIONAL GENERALIZATION OF A LINE $a_1 x_1 + \dots + a_n x_n = b$

LINEAR SEPARABILITY: PROBABILITY OF A SET OF POINTS WHICH CAN BE SEPARATED BY AN HYPERPLANE A LINEAR MODEL CAN EXPRESS A NON-LINEAR DECISION BOUNDARY BY ENCODING A SPECIFIC NON-LINEARITY IN ITS FEATURE FUNCTION; YOU NEED TO KNOW THE NON-LINEARITY A PRIORI

MULTI-LAYER PERCEPTRON: FEED-FORWARD FULLY-CONNECTED NN WHICH ALTERNATES LINEAR PROJECTIONS AND NON-LINEAR ACTIVATION FUNCTIONS; $h^l = \sigma^l(W^l \cdot h^{l-1} + b^l)$ $\sigma(x) \rightarrow$ ENCODING OF x

MLP + SOFTMAX: $p(y|x) = \text{SOFTMAX}(L^m, y) \in \mathbb{R}^{m+1}$, $\sigma'(x) = \sigma(x)(1 - \sigma(x)) \leq 0.25$; $\tanh(x) = 2\sigma(x) - 1$

$\text{tanh}'(x) = 1 - \tanh^2(x)$; $\text{RELU}(x) = \max(0, x)$ $\text{RELU}'(x) = I_{x>0}$; IN AN MLP THE NLL IS NO LONGER CONVEX

MLP WITH AN ARBITRARILY LARGE HIDDEN LAYER AND A NON LINEAR AND UNIVERSAL FUNCTION APPROXIMAT.

PERCEPTRON: SINGLE NODE MLP WITH $T \rightarrow 0$, TRAINED WITH SGD (MINIBATCH SIZE OF 1); CAN'T MODEL XOR

WORD PROCESSING: TOKENIZATION, STEMMING, STOP WORD REMOVAL; COMBINE ONE-HOT WORD ENCODING IN BAG

OF WORDS FORMAT, WITH SENTENCE POOLING OR WITH n-GRAM CONCATENATION; EMB. SIZE GOES FROM $|V|$ TO $|V|^n$

SKIP-GRAM: SENTENCIATE AND TOKENIZE CORPUS THEN BUILD POSITIVE SAMPLES WITH GIVEN CONTEXT WINDOW K LOG-BILINEAR MODEL $P(c|w) = Z(w)^{-1} \exp(\text{score}(w) \cdot c)$ USING DIFFERENT CONTEXT AND CENTER EMBDS. TO AVOID COMPUTATION OF $Z(w)$ SAMPLE NEGATIVE SAMPLE PAIRS AND TRAIN FOR A CLASSIFICATION LOGISTIC LOSS OUTPUT: (2) SET OF WORD EMBEDDING; EVALUATION WITH COSINE SIMILARITY OR WORD ANALOGIES (BERT SORTA) MLP FOR SENTIMENT ANALYSIS: EMB + POOLING + MLP + SOFTMAX; WITH MLP FEATURE ENGINEERING IS AVOIDED BUT ARCHITECTURE ENG. CAN ARISE INSTEAD; RELU HAS DEAD NEURONS IF INPUT IS NEGATIVE

5 - LANGUAGE MODELS, RNNs

STRUCTURED PREDICTION: PREDICTING STRUCTURED OBJECTS IN AN EXPONENTIALLY LARGE SPACE KLEENE CLOSURE V^* : SET OF ALL FINITE STRINGS OVER SYMBOLS IN A VOCABULARY V .

LANGUAGE MODEL: NORMALIZED PROB. DISTR. OVER THE KLEENE CLOSURE OF A VOCAB; REPRESENTABLE AS A WEIGHTED PREFIX TREE; $y \in V^* \rightarrow P(y) = Z^{-1} \prod_{t=1}^{|y|} \theta_{yt}$; $Z = \sum_{y \in V^*} \prod_{t=1}^{|y|} \theta_{yt}$; Z CAN BE COMPUTED WITH GLOBAL NORMALIZATION (ARBITRARY WEIGHTS) OR LOCAL NORMALIZATION (CHOOSE 0 S.T. $\sum_{y \in V^*} \theta_{yt} = 1$)

TIGHT LANG MODEL: $P(\text{EOS}|\text{IMENT}) > 0$, ENFORCES FINITE NUMBER OF INFINITELY LONG PATHS; CAN CONDITION ON x M-GRAM ASSUMPTION: $P(y_{t+1}|y_{t-1}, \dots, y_{t-m+1}) = \frac{\exp(w_{y_t \cdot y_{t-1}, \dots, y_{t-m+1}})}{\sum_{y \in V} \exp(w_{y \cdot y_{t-1}, \dots, y_{t-m+1}})}$ | BENGIO ET. AL.: NEURAL LANG.

$w \in V \rightarrow$ WORD VECTOR $h \in \mathbb{R}^d$ M-GRAM CONTEXT VECTOR $\sum_{y \in V} \exp(w_{y \cdot h})$ | MODEL: USES WORD EMBEDDS.

TO SOLVE CURSE OF DIMEN. $(|V|^m \rightarrow d^m)$ | MODEL: GET EMBEDDINGS, COMBINE CONTEXT WITH MLP $\rightarrow h = f(h_{\text{ctx}})$

USE SOFTMAX TO NORMALIZE $P(y_{t-1}, \dots, y_{t-m+1}) \rightarrow$ OVER $|V|$, SO SLOWISH; TRAIN BOTH MODEL PARAM AND EMBDS.

RNN: COMBINE TIME-DEP HIDDEN STATE WITH CURRENT TOKEN; IN THEORY HAS INFINITELY LONG CONTEXT

VANILLA: $h_t = \text{tanh}([W[h_{t-1}; e(y_{t-1})]])$; LSTM CELL: HAS INPUT, OUTPUT AND FORGET GATES + 2 STATES

IS LESS AFFECTED BY VANISHING GRADIENTS DUE TO ADDITIVE UPDATE FUNCTION ON MORE THAN A MATRIX \rightarrow MORE CTX

GRU CELL: ONLY ONE STATE AND TWO GATES (INPUT AND FORGET) \rightarrow LESS PARAMETERS (FOR $t \rightarrow \infty$)

NEURAL RECURRENT MODEL: $h_t = f(h_{t-1}, e(y_{t-1}))$; $P(y_{t+1}|y_{t-1}) = \exp(w_{y_t \cdot h_t}) / (\sum_{y \in V} \exp(w_{y \cdot h_t}))$ | $h_t = f(h_{t-1}, e(y_{t-1}))$

BACKPROP THROUGH TIME: UNROLL RNN AND BACKPROP NORMALLY; VANISHING GRADIENTS ARISE SINCE THE SAME SHARED WEIGHTED MATRIX IS MULTIPLIED FOR EVERY STATE $|V| < 1$ VANISHING $|V| > 1$ EXPLODING GRAD CLIP

RNN ARE PREFERRED OVER MLP FOR NLP SINCE THEY ENCODE A LONGER AND POSITIONAL CONTEXT WHICH IS USEFUL TO ENCODE LONG DISTANCE DEPENDENCIES (E.G. SUBJ. VERB AGREEMENT IN GERMAN SUBORDINATES)

UNIGRAM: $P(\text{HERE YOU ARE EOS}) = P(\text{HERE}) P(\text{YOU}) P(\text{ARE}) P(\text{EOS})$

BIGRAM: $P(\text{BOS HERE YOU EOS}) = P(\text{HERE|BOS}) P(\text{YOU|HERE}) P(\text{EOS|YOU}) = \frac{P(\text{BOS HERE})}{|V|} \frac{P(\text{HERE})}{|V|} \frac{P(\text{YOU})}{|V|}$

LIDSTONE SMOOTHING: $P(w_m|w_{m-1}) = \frac{\text{COUNT}(w_{m-1}, w_m) + \lambda}{\sum_{v \in V} \text{COUNT}(w_{m-1}, v) + |V|}$ AVOIDS ASSIGNING PROB. = 0 TO PLATIBLE

LAPLACE SMTH. $\Rightarrow \lambda = 1$

EXPECTED NUMBER OF BIGRAMS WITH ZERO COUNT FROM VOCAB V AND M TOKENS IN CORPUS = $(1 - |V|^{-2})^M$

6 - CONDITIONAL RANDOM FIELDS

PART OF SPEECH: CATEGORY OF WORDS THAT SHOW SIMILAR SYNTACTIC BEHAVIOR $\rightarrow P(t|w) = \frac{\exp(\text{score}(t, w))}{\sum_{t' \in T} \exp(\text{score}(t', w))}$

CONDITIONAL RANDOM FIELD: CONDITIONAL PROBABILISTIC MODEL FOR SEQUENCE LABELING $\sum_{t \in T} \exp(\text{score}(t, w))$

SCORE(t, w) CAN BE A DOT PRODUCT LIKE $\theta^T F(t, w)$ OR NEURALIZED WITH A NEURAL NETWORK $\rightarrow \text{NN}(t, w)$

NAÏVELY COMPUTING Z RUNS IN $O(|T||V|)$ $Z \rightarrow$ SET OF POS TAGS, $|V| = N$, LENGTH OF INPUT SENTENCE

DEFINE ADDITIVELY DECOMPOSABLE SCORE FUNCTION: $\text{SCORE}(t, w) = \sum_{m=1}^M \text{SCORE}([t_{m-1}, t_m], w)$

NORMALIZER: $\oplus_{t \in T} \text{Z}^N \otimes_{m=1}^M \text{EXP}(\text{SCORE}([t_{m-1}, t_m], w)) \oplus_{t \in T} \text{Z}^{N-1} \otimes_{m=1}^M \text{EXP}(\text{SCORE}([t_{m-1}, t_m], w)) = \oplus_{t \in T} \text{Z}^{N-1} \otimes_{m=1}^M \text{EXP}(\text{SCORE}([t_{m-1}, t_m], w)) \oplus_{t \in T} \text{Z}^{N-2} \otimes_{m=1}^M \text{EXP}(\text{SCORE}([t_{m-1}, t_m], w)) = \oplus_{t \in T} \text{Z}^{N-3} \otimes_{m=1}^M \text{EXP}(\text{SCORE}([t_{m-1}, t_m], w)) = \dots = \oplus_{t \in T} \text{Z}^0 \otimes_{m=1}^M \text{EXP}(\text{SCORE}([t_{m-1}, t_m], w)) = \oplus_{t \in T} \text{Z}^0 \otimes_{m=1}^M \text{EXP}(\text{SCORE}([t_{m-1}, t_m], w)) = \text{LINEAR IN } N$

VITERBI ALGORITHM: $O(T^2 N^3)$ THIS ALGORITHM CAN ALSO BE COMPUTED IN A "FORWARD" MANNER, IN THAT CASE

$\gamma(w, t_i) = V_i(t)$ REPRESENTS THE SCORE OF THE BEST TAG SEQUENCE

FOR $m \in N-1, \dots, 0$; ENDING IN TAG t AT POSITION i ; THE ALGO COMPUTES THE BEST SCORE

$\gamma(w, t_m) = \max_{t_{m-1} \in T} \text{Exp}(\text{score}([t_{m-1}, t_m], w)) \times \gamma(w, t_{m-1})$ TAGGING SEQUENCE FOR SENTENCE w

SEMIRING: 5-TUPLE $R = (A, \oplus, \otimes, \bar{0}, \bar{1})$ SO THAT 1) $(A, \oplus, \bar{0})$ IS A COMMUTATIVE MONOID 2) $(A, \otimes, \bar{1})$ IS MONOID

3) \otimes DISTRIBUTES OVER \oplus FROM BOTH LEFT AND RIGHT 4) $\bar{0}$ IS AN ANNihilATOR OF \otimes : $\forall a \in A$, $\bar{0} \otimes a = a \otimes \bar{0} = \bar{0}$

PROOF: 1) CHECK IDENTITY OF $\bar{0}$, ASSOCIATIVITY AND COMMUT. OF \otimes 2) CHECK IDENTITY OF $\bar{1}$ AND ASSOCIATIVITY OF \otimes

3) $(a \otimes b) \otimes c = (a \otimes c) \oplus (b \otimes c)$; $c \otimes (b \otimes a) = (c \otimes a) \oplus (c \otimes b)$ 4) PROPERTY ABOVE

BOOLEAN $\{0, 1\}$ V \wedge 0 1 (THE NORMALIZER CAN BE COMPUTED WITH THE INSIDE SEMIRING

VITERBI [0, 1] MAX X 0 1 (BACKWARD/FORWARD ALGORITHM); THE LOG-Z CAN BE COMPUTED

INSIDE $R^+ U \otimes 0$ + X 0 1 WITH THE LOG SEMIRING (USING LOG ON SCORES AS FIRST STEP)

REAL $R^+ U \otimes 0$ MIN + $\otimes 0$ 0 CRF CAN BE TRAINED BY MAXIMIZING LL WITH BACKPROP ON DATASET

TROPICAL $R^+ U \otimes 0$ MIN + $\otimes 0$ 0 $\sum_{t=1}^T (\text{score}(t, w) - T \log \sum_{t' \in T} \text{Exp}(\text{score}(t', w)/T)) \rightarrow$ SOFTMAX

COUNTING N + X 0 1 ANNEALING A CRF WITH $T \rightarrow 0$ WE OBTAIN A STRUCTURED PERCEP.

LOG $R U - \infty \ln(e^a + e^b)$ + - ∞ 0 TRAINABLE WITH SUBGRADIENT DESCENT AND THE PERCEPTION UPDATE

VITERBI $R U - \infty$ MAX + - ∞ 0 $-T \log \sum_{t' \in T} \text{Exp}(\text{score}(t', w)/T) \geq -\max_{t' \in T} \text{score}(t', w)$

THE FORWARD-BACKWARD ALGORITHM CAN BE SEEN AS FORWARD + BACKPROP WHEN SUBSTITUTING THE BACKWARD

RECURRENCE WITH BACKPROPAGATION. BOTH BACKWARD PASSES USE DYNAMIC PROGRAMMING TO COMPUTE $\Delta \log Z$

VITERBI EXAMPLE: $\max_{t \in T} \text{score}(t, "THEY CAN CAN FISH")$; $v_n(t_n) = \max_{t_n} (\text{score}([t_n, t_n], w) + v_{n-1}(t_{n-1}))$

$v_n(t_n) = \text{score}([t_n, t_n], w)$; $v_n(N) = -2 - 1 = -3$, $v_n(V) = -12$; $v_L(N) = \max(v_n(N) + \text{score}([t_n=N, t_2=N]), v_n(V) + \text{score}([t_n=V, t_2=N])) = \max(-3 - 3, -12 - 3) = -9$; $b_2(N) = N$; $v_L(V) = \max(v_n(N) + \text{score}([t_n=N, t_2=V]), v_n(V) + \text{score}([t_n=V, t_2=V])) = -5$; $b_2(V) = N \rightarrow$ BACKPOINTER, THE TAG THAT MAXIMIZED THE SCORE IS N

EMISSION FEATURES: $t \mapsto w$; TRANSITION FEATURES: $t \mapsto t'$ COMPONENTS OF SCORE FUNCTION

7 - CONTEXT-FREE PARSING

SYNTAX: MATHEMATICAL STUDY OF THE STRUCTURE OF SENTENCES; LANGUAGE IS STRUCTURED HIERARCHICALLY AND CAN BE BROKEN DOWN INTO CONSTITUENTS: MULTI-WORD UNIT THAT FUNCTIONS AS A SINGLE UNIT
KINDS OF AMBIGUITY: ATTACHMENT AMBIGUITY, MODIFIER SCOPE; CONSTITUENCY TESTS: PRONOUN SUBSTITUTION, CLEFTING, ANSWER ELLIPSIS; MANY LINGUISTIC PROPERTIES ARE DEFINED ONLY OVER TREES
GRAMMAR: ORDERED SET OF RULES THAT DEFINES HOW TO FORM STRINGS FROM A LANGUAGE VOCABULARY

AMBIGUOUS GRAMMAR: THERE ARE MORE WAYS TO GENERATE THE SAME STRING;
CFG GRAMMAR: RULES ARE APPLIED REGARDLESS OF THE CONTEXT; $G = \langle N, S, \Sigma, R \rangle$; CHOMSKY NORM.

FORM: ALL PRODUCTION RULES ARE OF THE FORM: $N_i \rightarrow N_j N_k, N_i \rightarrow a$ (IN CFG RULES ARE $N \rightarrow \cdot$)

PARSE-STRUCTURE TREES REPRESENT THE SYNTACTIC STRUCTURE OF THE SENTENCE AND ITS DERIVATION WITH G ; PCFG: ASSIGN A PROB. TO EACH PRODUCTION RULE s.t. $\sum_{k=1}^K p(\cdot|N) = 1$

WCFG: ASSIGN GENERIC NON-NEGATIVE WEIGHTS $\rightarrow p(t) = z^{-1} \prod_{t \in T} \exp(score(t))$ $z = \sum_{t \in T} \exp(score(t))$

THIS z CAN'T BE COMPUTED EASILY DUE TO INFINITE SUM BIGGER THAN Σ^* ; CONDITION ON INPUT SENTENCE s : $p(t|s) = z(s)^{-1} \exp(score(t))$; $z(s) = \sum_{t \in T(s)} \exp(score(t))$; ONLY SUM ON TREES THAT YIELD s ; WE FACTOR THE SCORE OF A TREE ALONG THE PRODUCTIONS; THE NORMALIZER CAN STILL DIVERGE \rightarrow TRANSFORM CFG IN CNF

$z(s) = \sum_{t \in T(s)} \prod_{x \in t} \exp(score(x \rightarrow y)) \cdot \prod_{x \in t} \exp(score(x \rightarrow x))$; IN CNF, THE NUMBER OF ADMISSIBLE TREES IS THE CATALAN NUMBER $C_{M-1} = (M-1+1) \binom{M-1}{M-1}$ FOR A M-WORD LONG SENTENCE $O(4^m)$

CKY ALGORITHM: DP TO COMPUTE z ; TIME COMPLEXITY $O(M^3 / R!)$, SPACE COMPLEXITY $O(M^2 / N!)$

SEMITRING CYK ($S, \langle N, S, \Sigma, R \rangle, score$)
 $N \leftarrow \{s\} \rightarrow (\text{was } M \text{ ABOVE})$

CHART $\leftarrow \emptyset$
FOR $m=1, \dots, N$: \rightarrow FILL IN FIRST ROW OF PYRAMID CHART

FOR $x \rightarrow s_n \in N$:
CHART [$m, m+1, x$] $\oplus = \exp(score(x \rightarrow s_n))$

THE VITERBI, INSIDE, LOG SEMIRINGS COMPUTE THE BEST PARSE, z AND $\log(z)$ RESPECTIVELY

CFG PARSERS ARE TRAINABLE WITH BACKPROP. BY MAXIMUM LIKELIHOOD ESTIMATION ON A DATASET

FOR SPAN = 2, ..., N: \rightarrow SPAN LENGTH
FOR $i = 1, \dots, N - \text{SPAN}$: \rightarrow START OF SPAN

$k \leftarrow i + \text{SPAN} \rightarrow$ END OF SPAN
FOR $j = i+1, \dots, n-1$: \rightarrow SUBSPAN DIVIDER

FOR $x \rightarrow y \in N$:

CHART [i, k, x] $\oplus = \exp(score(x \rightarrow y))$

\otimes CHART [i, j, y] \otimes CHART [j, k, z]

RETURN CHART [$0, N, S$] \rightarrow CAN ADD BACKPOINTERS

8 - DEPENDENCY PARSING

DEPENDENCY GRAMMAR: LINK EVERY WORD TO ITS SYNTACTIC HEAD; LABELED, DIRECTED HEAD \rightarrow DEPENDENT EDGES FORM A DEPENDENCY TREE \rightarrow DIRECTED SPANNING TREE WITH ROOT CONSTRAINT OF 1; RELATED TO LEXICALIZED CFGs; CONVERT TO AND FROM CONST. TREES WITH GRAMMATICAL HEURISTICS

PROJECTIVE DEP. TREES: HAVE NO CROSSING ARCS AND CAN BE PARSED WITH VERSIONS OF CKY (AND VICE V.)
 $p(t|w) = z^{-1} \exp(score(t, w))$; $z = \sum_{t \in T(w)} \exp(score(t, w)) \rightarrow$ SUM IS $O(m^n)$ NAIVELY ($m = |w|$) \rightarrow ROOT EDGE

ADD EDGE FACTORED ASSUMPTION: $z = \sum_{t \in T(w)} \prod_{i=1}^n \exp(score(i, w)) \exp(score(r, w))$ SCORE

KIRCHHOFF MT: COUNT NUMBER OF UNDIRECTED SPANNING TREES IN $O(m^n)$; GENERALIZED FOR DIRECTED TREES WITH A ROOT CONSTRAINT ST. $z = \delta(t|L)$, $L_i = p_j$ IF $i=1$, $\sum_{j=1, j \neq i}^m A_{ij}$; ELSE $-A_{ij}$

DECODING: FINDING MAXIMUM-WEIGHT TREE WITH ROOT CONSTRAINT, NO CYCLES, ONE INCOMING EDGE PER NODE

$\text{ARGMAX}_{t \in T(w)} \sum_{i \in \text{S}(t)} \text{score}(i, w)$; GREEDY APPROACHES FAIL

CHU-LIU-EDMONDS ALGORITHM: SOLVES THE DECODING PROBLEM IN $O(m^2)$ \rightarrow SPECIAL DATA STRUCTURE MST (G):

\curvearrowright TAKE BEST INCOMING NODE FOR EACH NODE

IF CYCLE IN GREEDY(G): \curvearrowright SELECT MAX \curvearrowright UPDATE ENTER EDGES WEIGHT TO HANDLE CYCLE BREAK

RETURN EXPAND(CONSTRAIN(MST(CONTRACT(G , CYCLE)))

ELSE: RETURN CONSTRAIN(GREEDY(G))

CONSTRAIN(G): \curvearrowright ROOT CONSTRAINT VIOLATION

IF NUM_ROOT_EDGES(GREEDY(G)) > 1: \curvearrowright DELETE NODE WITH LOWER COST

E = ROOT_EDGE_TO_REMOVE(G) COST = EDGE COST - WEIGHT OF NEXT BEST INCOMING EDGE

IF CYCLE IN GREEDY($G-E$): RETURN CONSTRAIN(CONTRACT(G , CYCLE))

ELSE: RETURN CONSTRAIN($G-E$)

ELSE: RETURN GREEDY(G); A NAIVE HANDLING OF THE ROOT CONSTRAINT WOULD ADD $O(m)$ TO THE RUNTIME: FIX EACH EDGE AS THE ONLY ONE EMANATING FROM P AND RE-RUN MST TO SELECT BEST

9 - LAMBDA CALCULUS

ENCODE SEMANTIC REPRESENTATION OF TEXT, REMOVE SEMANTIC AMBIGUITIES

PRINCIPLE OF COMPOSITIONALITY: THE MEANING OF A COMPLEX EXPRESSION IS A FUNCTION OF ITS CONSTITUENT'S LAMBDA OPERATOR: $\lambda x. f(x) \rightarrow$ FUNCTION THAT TAKES x AS INPUT AND PRODUCES $f(x)$ AS OUTPUT \curvearrowright (NOT BOUND)

FREE VARIABLE: VARIABLE THAT DOES NOT OCCUR IN THE SCOPE OF ANY ABSTRACTION WITH ITS NAME

λ -CONVERSION: RENAME A VARIABLE IN AN ABSTRACTION WITH ALL ITS BOUND VARIABLES $\lambda x. (x x) \rightarrow \lambda t. (t t)$

β -REDUCTION: APPLY A LAMBDA TERM TO ANOTHER: $(\lambda x. \lambda y. (x (y x) y)) z \rightarrow \lambda y. ((z (y z) z))$ \rightarrow ONLY FREE CH.

ENRICHED LAMBDA CALCULUS: HANDLE SENTENCE SEMANTICS; LOGICAL CONSTANTS, VARIABLES, LITERALS

EXAMPLES: $(\lambda Q. Q(\lambda t. (\forall x. (DOG(x) \Rightarrow LIKES(t, x)))) \lambda P. P(ALEX)) \rightarrow (\lambda P. P(ALEX) \lambda t. (\forall x. (DOG(x) \Rightarrow LIKES(t, x))))$

$\rightarrow (\lambda t. (\forall x. (DOG(x) \Rightarrow LIKES(t, x))) \text{ALEX}) \rightarrow \forall x. (DOG(x) \Rightarrow LIKES(ALEX, x))$; $(\lambda P. \lambda Q. \lambda t. (\forall x. P(\lambda y. LIKES(x, y))) \lambda P. P(ALEX))$

$\rightarrow \lambda Q. Q(\lambda x. (\lambda P. P(ALEX)) \lambda y. LIKES(x, y)) \rightarrow \lambda Q. Q(\lambda x. (\lambda y. LIKES(x, y)) \text{ALEX})$; $(\lambda Q. Q(\lambda x. LIKES(x, ALEX)))$

$\rightarrow \lambda Q. \exists x. (DOG(x) \wedge Q(x)) \rightarrow (\lambda Q. \exists x. (DOG(x) \wedge Q(x))) \lambda x. LIKES(x, ALEX) \rightarrow \exists x. (DOG(x) \wedge (\lambda x. LIKES(x, ALEX))) \lambda x \dots$

10 - WFST

WEIGHTED FINITE STATE TRANSDUCERS: MAP SEQUENCES FROM SOURCE TO TARGET VOCABULARY; HAS FINITE NUMBER OF STATES WEIGHTED WITH PROBS; $T = \langle Q, \Sigma, \Delta, \lambda, p, \delta \rangle$ Q : SET OF STATES; Σ, Δ INPUT AND OUTPUT VOCABULARY
 λ , p MAP STATES TO INITIAL OR FINAL SCORES, δ MAPS TRANSITION ARCS TO SCORES; UNAMBIGUOUS WFST:
EACH PAIR OF STRINGS HAS AT MOST ONE ACCEPTING PATH; SCORE(π) = $\sum_{i=1}^{|T|} \text{score}(\pi_i)$; $p(\pi) = 2^{-|T|} \exp(score(\pi))$
 $Z = \sum_{y \in \Sigma^*} \exp(score(y, x)) = \sum_{y \in \Sigma^*} \sum_{\pi \in T(x,y)} \exp(score(\pi))$

FLOYD-WARSHALL ALGORITHM: SOLVES ALL PAIRS SHORTEST PATHS GIVEN NO NEGATIVE CYCLES $\rightarrow O(|V|^3)$
FOR EACH EDGE (u, v) : DIST MATRIX INIT WITH 0
 $DIST[u][v] = W[u][v] \rightarrow W$

FOR EACH VERTEX v : DIST $[v][v] = 0 \rightarrow W^0$
SEMIRING GENERALIZE FROM MATRIX MULT.

FOR k IN $1, \dots, N$:
FOR i IN $1, \dots, N$: DIST [i][i] = DIST [i][i] \oplus

FOR j IN $1, \dots, N$: DIST [i][j] = DIST [i][k] \otimes DIST [k][j]
IF DIST [i][j] > DIST [i][k] + DIST [k][j] \curvearrowright DIST [i][j] = DIST [i][k] + DIST [k][j]

TRIGRAM: $|V| + |V|$ STATES, $\delta(q_i; w, q_k) = \log p(w = k | w_{m-1} = i, w_{m-2} = j)$ IF $w = k$; m -GRAM: $|V|^{m-2} + |V|^{m-1} + 1$ STATES

11 - SEQUENCE MODELS

TRANSLATION: MAP STRINGS FROM ONE SPACE TO ANOTHER, BUT THERE IS NO CORRECT ANSWER, LOCALITY ASSUMPTIONS ARE NOT REASONABLE; NEURAL MACHINE TRANSLATION MODELS THE MAPPING WITH A COMPLEX FUNCTION;

DIVISIBLE IN MODELING AND DECODING; SEQ-2-SEQ MODEL: $z = \text{ENCODER}(x); y_k = \text{DECODER}(z) \rightarrow$ LOCALLY NORM.

$p(y|x) = \prod_{t=1}^T p(y_t | x_1, \dots, x_{t-1})$; ATTENTION: RATHER THAN ENCODING ALL INPUT INTO A SINGLE CONTEXT VECTOR, PAY VARIABLE ATTENTION TO EACH REPRESENTATION DEPENDING ON THE OUTPUT GENERATION STEP

$d^e = \text{SOFTMAX}(\text{score}(q_t, k))$; $c^t = d^e T V$; $k = v_i = h_i^e$; $q_t = h_t^d$; $K = V = H^p$; SELF-ATTENTION: ATTENTION ON OWN CONTEXT: $k = v = h_i^d$; d^e ; $h_t^d = h_t^d + c^t$; MULTI-HEAD ATTENTION: LEARN MULTIPLE SETS OF ATTENTION WEIGHTS FOR THE SAME INPUT AND CONCATENATE; SCORE CAN BE SCALING DOT PRODUCT $\text{score}(q_t, h_i) = q_t^T h_i \cdot \|h_i\|^2$

TRANSFORMERS: SOLVE SLOW RECURRENCE OF RNNs; USED FOR MOST LANGUAGE GENERATION TASK; POSITIONAL EMBEDDINGS: ENCODE ORDER OF WORDS IN THE SEQUENCE; RESIDUAL CONNECTIONS: MITIGATE VANISHING GRADS.

LAYER NORMALIZATION: HELPS WITH INTERNAL COVARIANCE SHIFT, NORMALIZE INDIVIDUAL LAYER INPUTS

DECODING: $y^* = \arg\max_y y \text{ score}(x, y)$, NO INDEPENDENCE ASSUMPTION, SO HEURISTICS OR SAMPLING METHODS ARE NEEDED; BEAM SEARCH: PRUNED BFS WITH LIMITED BREADTH; NUCLEUS SAMPLING: SAMPLE FROM CONDITIONAL DISTR.

CONSIDERING ONLY TOP ITEMS COVERING ρ % OF THE PROBABILITY MASS \rightarrow PRECISION OF TRANSLATION EVALUATION METRICS: BLEU: FRACTION OF PREDICTED n -GRAMS THAT APPEAR IN THE GROUND TRUTH TOO METEOR: COMBINATION OF RECALL AND PRECISION OF PREDICTED TRANSLATION W.R.T. THE REFERENCE TRANSFORMERS ARE QUADRATIC IN THE CONTEXT SIZE AND PERFORM BETTER THAN n -GRAM MODELS WITH BIG DATASETS NUMBER OF TERMS DERIVATIVE SEQ-2-SEQ: $O(N-m+m)$, WITH ATTENTION $O((N-m+m) \cdot N)$

$\frac{\partial \ln m!}{\partial \ln x_1} = \frac{\partial \ln m!}{\partial \ln x_2} \dots \frac{\partial \ln x_m}{\partial \ln x_1} \dots \frac{\partial \ln x_m}{\partial \ln x_m}; \frac{\partial \ln m!}{\partial \ln x_1} = \frac{\partial \ln m!}{\partial x_1} + \sum_{i=1}^N \frac{\partial \ln x_i}{\partial \ln x_1} \frac{\partial \ln x_i}{\partial x_1}; h_i^s \rightarrow$ WITH ATTENTION VECTOR

ATTENTION REDUCES THE SHORTEST DECODER \rightarrow ENCODER PATH LENGTH FROM PROPORTIONAL TO THE INPUT (AND OUTPUT) SIZE, TO ONLY CONSTANT

12 - AXES OF MODELING

CLASSIFICATION MODELS CAN BE DIVIDED INTO PROBABILISTIC OR DETERMINISTIC; PROBABILISTIC MODELS CAN BE DISCRIMINATIVE OR GENERATIVE; BOTH CAN BE LOCALLY OR GLOBALLY NORMALIZED; DIFFERENT MODELS HAVE DIFFERENT INDUCTIVE BIASES THAT CAN BE USEFUL TO SOLVE SPECIFIC PROBLEMS; LARGE MODELS OVERFIT DUE TO THE BIAS-VARIANCE TRADEOFF; LOSS FUNCTION: ANY FUNCTION OF MODEL PARAMETERS THAT CLASSIFIES FIT TO DATA

MLE ESTIMATORS ARE CONSISTENT: $\hat{\theta} = \arg\min_{\theta} \mathbb{E}_{x \sim p(x)} \ell(x, \theta)$ AND ASYMPTOTICALLY EFFICIENT, BUT CAN EASILY OVERFIT LOSS SHOULD BE CONVEX, DIFFERENT, LOW COMPUT. COST; HINGE LOSS $\ell(y, \hat{y}) = \max\{0, 1 - y \hat{y}\}$; EXP LOSS: $\ell(y, \hat{y}) = -y \hat{y}$

REGULARIZATION: ADD PRIOR INFORMATION TO PREVENT OVERFITTING TO NOISE \rightarrow NECESSARY FOR GOOD GENERALIZATION MODELS ARE EVALUATED ON EVALUATION METRICS WHICH ARE HIGHLY PROBLEM SPECIFIC AND DIFFERENT FROM LOSS TO WEIGHT ERROR TYPES DIFFERENTLY OR FOR HIGH COMPUTATION COST \rightarrow SOMETIMES THEY COINCIDE (5)

INTRINSIC EVALUATION: EVALUATE BASED ON A CRITERIA; PERPLEXITY (\hat{y}) = $2^{-\frac{\log p(\hat{y})}{|V|}}$ \rightarrow HELD-OUT LOG-LIKELIHOOD EXTRINSIC EVALUATION: EVAL. BASED ON THE PERFORMANCE OF A SUCCESSIVE TASK

MODEL SELECTION: CHOOSING MODEL WITH BEST INFERENCE AND PREDICTION PERFORMANCE; NESTED-CV EVALU.

MODEL STABILITY; P-VALUE OF STATISTICAL TEST $\epsilon: p = 2 \min(P(T \neq H_0), P(T \neq H_1))$, REJECT H_0 IFF $p < \alpha$

STATISTICAL POWER: $P(\text{reject } H_0 | H_1)$, INCREASES WITH INCREASING SAMPLE SIZE; MULTIPLE TESTING: $P(\text{1 false rejection}) > 0 = 1 - (1 - \alpha)^k$; BOFFERONI CORRECTION: $\alpha^* = \alpha / k$; STATISTICAL TESTS CAN BE PARAM. OR NON-PARAMETRIC \rightarrow THE PARAMETER OF THE DISTRIBUTION OF THE TEST IS NOT A FUNCTION OF THE DATA \rightarrow $\chi^2 \sim \chi^2_{k-1}$

MCMILLAN'S TEST: COMPARE DISAGREEMENT PROBABILITY OF TWO CLASSIFIERS THAT SHARE TEST SET $H_0: p_0 = p_1$ \rightarrow $\chi^2 \sim \chi^2_{k-1}$

PERMUTATION TESTS: TEST IF CLASSIFIERS PERFORMS BETTER THAN RANDOM CHANCE ON DATA $P(\text{VAL} = \{1, 0, \dots, 0\}) = 1 - (k+1)^{-k}$

5x2 CV PAIRED T-TEST: TEST TO COMPARE PERFORMANCE OF TWO CLASSIFIERS ON A METRIC P

DIFFERENT TESTS HAVE DIFFERENT STATISTICAL POWERS AND OBSERVED TYPE I ERRORS \rightarrow IMPORT. FOR CHOOSING THE TEST SET CAN BE NOT LARGE ENOUGH TO GIVE ENOUGH POWER TO DETECT DIFFERENCES IN MODEL PERFORMANCE