

9. Intro

Hypergraph View: Computation graph = labeled acyclic hypergraph. Edges can have multiple sources/targets. **Complexity:** same time as f ; space higher (store intermediates) vec-vec: $O(d)$; mat-vec: $O(nm)$; mat-mat: $O(nml)$

$$\text{NLL } \nabla = \mathbf{0}: \sum_{i=1}^n f(x_i, y_i) = \sum_{i=1}^n \mathbb{E}_{y|x_i} [f(x_i, y)] \\ \text{Observed features} = \text{Expected features} \quad \text{Hessian: } H = \sum_i \text{Cov}_{y|x_i} [f(x_i, y)] \text{ (PSDI)}$$

DAG Properties: Topological order 唯一确定; DP 子问题独立拆分可行; Gradient 反向传播良定义(no cycles) **Hypergraph:** 函数式计算自然表示, multi-inputs → one output

1. Backpropagation

Chain: $\frac{d}{dx}[f(g(x))] = f'(g(x))g'(x)$ **Jacobian:** $f: \mathbb{R}^n \rightarrow \mathbb{R}^m, \frac{dy}{dx} = \begin{bmatrix} \frac{dy_1}{dx_1}, \dots, \frac{dy_m}{dx_n} \end{bmatrix} \in \mathbb{R}^{m \times n}$ **Multivar:**

$$\frac{dy_i}{dx_j} = \sum_{k=1}^m \frac{\frac{dy_i}{dz_k} dz_k}{\frac{dx_j}{dz_k} dz_k}$$

Bauer Path: $\frac{dy_i}{dx_j} = \sum_{p \in \mathcal{P}(j,i)} \prod_{(k,l) \in p} \frac{dz_l}{dz_k} \mathcal{P}(j,i) = \text{all paths } j \rightarrow i; \text{ worst } O(m^n), m \text{ 平均出度, } n \text{ 路径长度}$

Forward vs Reverse: **Forward:** expand $\frac{d}{dx}$ recursively, same flow as fwd **Reverse:** 2 passes—fwd compute vals, bwd compute grads **Complexity:** same time as f ; higher space (store intermediates)

Primitives: **Sum:** $\frac{d(a+b)}{da} = 1$; **Prod:** $\frac{d(ab)}{da} = b$

2. Log-Linear Models

Prob Basics: Bayes: $p(y|x) = \frac{p(x|y)p(y)}{\int p(x|y)p(y)dy}$ Posterior \propto Prior \times Likelihood **Marginal:** $p(x) = \sum_y p(x,y)$

Expectation: $\mathbb{E}[f(x)] = \sum_x f(x)p(x)$

Log-Linear Model: $p(y|x, \theta) = \frac{\exp(\theta \cdot f(x, y))}{Z(\theta)}$ $Z(\theta) = \sum_{y' \in Y} \exp(\theta \cdot f(x, y')) \log p(y|x, \theta) = \theta \cdot f(x, y) - \log Z$ (linear in log space!) **Discrete MLE:** $p(y|x) = \frac{\text{count}(x, y)}{\text{count}(x)}$ (sparse 问题)

MLE $\nabla: \theta_{\text{MLE}} = \arg \min_{\theta} -\sum_n \log p(y_n|x_n, \theta)$

观测 特征 count = 期望 特征 count →

Expectation Matching $\frac{d\mathcal{L}}{d\theta_k} = -\sum_n f_k(x_n, y_n) + \sum_n \sum_{y'} p(y'|x_n; \theta) f_k(x_n, y')$

MAP & Ridge: $\hat{\theta}_{\text{MAP}} = \arg \min_{\theta} [-\log p(\theta) - \log p(D|\theta)]$ Gaussian prior $\mathcal{N}(0, \sigma_p^2 I) \rightarrow \text{L2: } \frac{\lambda}{2} \|\theta\|^2$ Laplace prior → L1 regularization

Softmax: softmax(h, y, T) = $\frac{\exp(h_y/T)}{\sum_{y'} \exp(h_{y'}/T)}$ $T \rightarrow 0$: argmax; $T \rightarrow \infty$: uniform log softmax = $h_y - \log \sum_{y'} \exp(h_{y'})$ (logsumexp)

MLP Architecture: Problem: Log-linear needs linearly separable data **Solution:** Learn non-linear feature fn $h_k = \sigma_k(W_k h_{k-1})$, $h_1 = \sigma_1(W_1^\top e(x))$ Output: softmax($\theta^\top h_n$)

Activations: $\sigma(x) = \frac{1}{1+\exp(-x)}$, $\nabla \sigma = \sigma(1-\sigma)$ **tanh:** $\frac{1-e^{-2x}}{1+e^{-2x}}$, $\nabla = 1 - \tanh^2$ Sigmoid/tanh vanishing gradient → use ReLU **Backprop:** $\frac{\partial \ell}{\partial W_k} = \frac{\partial \ell}{\partial y} \frac{\partial y}{\partial h_n} \left(\prod_{m=k+1}^n \sigma'_m W_m \right) \sigma'_k h_{k-1}$

Exp Family & MaxEnt: $p(x|\theta) = \frac{1}{Z(\theta)} h(x) \exp(\theta \cdot \varphi(x))$ **Max Entropy:** $H(p) = -\sum_x p(x) \log p(x)$ 选

最大熵分布=最少假设=Laplace原则 **优势:** Conjugate priors; Sufficient stats; Convex log-partition → unique MLE

3. Language Models

Structured Prediction: Kleene V^* : infinite set of finite-length strings from V **Language Model:** weighted prefix tree, each sentence=unique path $p(y) = \frac{1}{Z} \prod_{t=1}^{|y|} \text{weight}_{y_t}$

Local Normalization: $Z = 1$ when children edges sum to 1 at each node **Consistency:** $p(\text{EOS}|y_{<t}, V^*) > \varepsilon > 0$ $p(|y| = \infty) \leq \lim_{t \rightarrow \infty} (1-\varepsilon)^t = 0$ (tight)

N-gram Model: $p(y_t|y_{<t}) = p(y_t|y_{t-1}, \dots, y_{t-n+1})$ **Markov:** $P(t_i|t_{1:i-1}) = P(t_i|t_{i-1})$ (1st order) =

$$\frac{\exp(w_{y_t \cdot h_t})}{\sum_{y' \in V} \exp(w_{y' \cdot h_t})}, h_t \in \mathbb{R}^d$$

Bengio: $h_t = f(e(\text{hist}))$, $e(\text{hist}) = [e(y_{t-1}); e(y_{t-2}); \dots]$

RNN: $h_t = f(h_{t-1}, e(y_{t-1}))$ (implicit infinite context)

Vanilla: $h_t = \sigma(W_1 h_{t-1} + W_2 e(y_{t-1}))$ **BPTT:** unroll through time, sum grads over timesteps

4. Word Embeddings

Encoding: One-hot: $v \in O(|V|)$, only word=1 **Bag-of-words:** pooled one-hot (sum/mean/max) **N-grams:** vectors huge—every combo needs slot **Pipeline:** Embedding → Pooling → NN → Softmax

Skip-gram: **Preprocess:** word-context pairs ($k \times C$ many), window k $p(c|w) = \frac{1}{Z(w)} \exp(e_{\text{wrd}(w)} \cdot e_{\text{ctx}(c)})$, $O(2|V|k)$ params **Bilinear:** linear if all-but-one vars held constant **Similarity:** $\cos(u_i, u_j)$

5. CRF & POS Tagging

As Graph: Fully connected graph w/ POS-tag nodes per layer score($(D, N, V, \dots, w), \theta$) = $\theta f(t, w)$ **score** (t, w) = unnormalized log-prob = \sum_n trans+emit Problem: $O(|\mathcal{T}|^N)$ paths in normalizer

CRF Model: $p(t|w) = \frac{\exp(\text{score}(t, w))}{\sum_{t' \in \mathcal{T}^N} \exp(\text{score}(t', w))}$ **Decomposition:** $\text{score}(t, w) = \sum_{n=1}^N \text{score}(\langle t_{n-1}, t_n \rangle, w, n)$

$p(t|w) \propto \prod_{n=1}^N \exp\{\text{score}(\langle t_{n-1}, t_n \rangle, w)\}$

DP推导: $O(|\mathcal{T}|^N) \rightarrow O(N|\mathcal{T}|^2)$: Goal: $Z = \sum_{t \in \mathcal{T}^N} \exp \text{score}(t, w)$ **Step1:** 可加 分解 $\text{score} = \sum_n \text{score}_n$ **Step2:** $Z = \sum_t \exp \sum_n \text{score}_n = \sum_t \prod_n \exp \text{score}_n$ (exp) **Step3:** $= \sum_{t_1} \dots \sum_{t_N} \prod_n \exp \text{score}_n$ (展开) **Step4:** $= \sum_{t_1} \exp \text{score}_1 \times \left(\sum_{t_2} \dots \right)$ (distrib 把内层 sum 推进去) **若 3-gram:** 依赖 $t_{n-2}, t_{n-1}, t_n \rightarrow O(N|\mathcal{T}|^3)$

Forward Algorithm: $\alpha[0, t] = \exp(\text{score}(\text{BOS} \rightarrow t))$ (init w/ BOS trans) for $n = 1, \dots, N-1$; for $t_n \in \mathcal{T}$: $\alpha[n, t_n] = \bigoplus_{t_{n-1}} \alpha[n-1, t_{n-1}] \otimes \exp(\text{score})$ return $\bigoplus_t \alpha[N-1, t]$ (sum last column!) **直觉:** prefix之和, 从 seq 开头走到当前状态的所有走法 score 总和

Backward Algorithm: $\forall t_N: \beta[N, t_N] \leftarrow 1$ for $n = N-1, \dots, 0$; for $t_n \in \mathcal{T}$: $\beta[n, t_n] \leftarrow \bigoplus_{t_{n+1}} \exp(\text{score}_{n+1}) \otimes \beta[n+1, t_{n+1}]$ return $\beta[0, \text{BOS}]$ (single value!) **Complexity:** $O(N|\mathcal{T}|^2)$

Fwd vs Bwd Asymmetry: Init: Bwd 直接1; Fwd 需 BOS 转移 Term: Bwd 返回 $\beta[0, \text{BOS}]$ 单值; Fwd 需 \oplus 整列 **原因:** BOS 显式存在, EOS 不显式处理

Viterbi Decoding: $\delta[n, t] = \max_{t_{n-1}} [\delta[n-1, t_{n-1}] + \text{score}(t_{n-1}, t)]$ 每步枚举 t 和 t_{n-1} 的 $|\mathcal{T}|^2$ 种 trans **Backtrack:** 存 argmax 指针 bp, 从 $\arg \max_t \delta[N, t]$ 回溯

Common Semirings:

Name	\mathbb{K}	\oplus	\otimes	0	1	用途
Real	$\mathbb{R}_{\geq 0}$	+	\times	0	1	Z partition
Viterbi	$\mathbb{R} \cup \{-\infty\}$	max	+	$-\infty$	0	最优 path
Log	$\mathbb{R} \cup \{\pm\infty\}$	lse	+	$-\infty$	0	$\log Z$
Boolean	{0,1}	\vee	\wedge	0	1	可达性
Counting	\mathbb{N}	+	\times	0	1	路径数
Tropical	$\mathbb{R} \cup \{\infty\}$	min	+	∞	0	最短路

Semiring Definition: $\langle \mathbb{K}, \oplus, \otimes, 0, 1 \rangle$ where:

1. $(\mathbb{K}, \oplus, 0)$: comm monoid (assoc+comm+identity)
2. $(\mathbb{K}, \otimes, 1)$: monoid (assoc+identity)
3. Distrib: $(x \oplus y) \otimes z = (x \otimes z) \oplus (y \otimes z)$
4. Annihilator: $0 \otimes x = x \otimes 0 = 0$

陷阱: $0 = 1$ 必失败!

Semiring 意义: \oplus : 分治 (split points 合并, OR/MAX/+); \otimes : 连接 (左右子树组合, AND/ $\times/+$); 0: 吸收元, 消除 invalid; 1: 单位元, null 不破坏

Monoid 判定:

1. Closure: $a \oplus b \in \mathbb{K}$
2. Assoc: $(a \otimes b) \otimes c = a \otimes (b \otimes c)$
3. Identity: $\exists e: a \otimes e = e \otimes a = a$

Kleene Star: $a^* = \bigoplus_{n=0}^{\infty} a^{\otimes n} = 1 \oplus a \otimes a^*$ Real 上 $|a| < 1: a^* = \frac{1}{1-a}$ (geometric series) Tropical: $a^* = 0$ if $a \geq 0$ (正环不帮助) 用于 globally normalized LM

6. CFG Parsing

Constituents: Multi-word units as single unit **Tests:** Pronoun substitution, Clefting, Answer ellipsis Ambiguity: PP attachment, modifier scope

CFG Definition: $G = (\mathcal{N}, \mathcal{S}, \Sigma, \mathcal{R})$ Non-terminals, start symbol, terminals, production rules **CNF:** $N_1 \rightarrow N_2 N_3$ or $N \rightarrow a$; $O(4^N)$ trees (Catalan)

Weighted CFG: **Global:** $p(t) = \frac{1}{Z} \prod_{r \in t} \exp(\text{score}(r))$ **Decomposition:** $Z = \sum_{t' \in \mathcal{T}} \prod_{r' \in t'} \exp(\text{score}(r'))$ (可能 ∞ !) **Probabilistic:** local norm $\sum_k p(\alpha_k | N) = 1$

CKY Chart 索引: Position 在 words 之间 : $0|w_1|1|w_2|2| \dots |N$ **Chart[i, k, X]:** span $[i, k]$ 覆盖 w_i, \dots, w_{k-1} **长度:** $k-i$; **对角线:** $k-i=1$ (单词) **Fill order:** 按 span 长度递增 ($\ell = 1, 2, \dots, N$) 同一长度内任意顺序 (topo order 自由度) **Goal:** Chart[0, N, S]

CKY algo: $O(N^3|R|)$, needs CNF **Terminal:** $C[i, i+1, X] = \exp(\text{score}(X \rightarrow w_i))$ for $X \rightarrow w_i \in \mathcal{R}$ **Binary:** for span = 2, ..., N ; for $i = 1, \dots, N$ – span: $k \leftarrow i+span$; for $j = i+1, \dots, k-1$; for $X \rightarrow Y \in \mathcal{R}$: $C[i, k, X] \oplus \exp\{\text{score}\} \otimes C[i, j, Y] \otimes C[j, k, Z]$

CKY Chart 3x3 Example: Sentence: $w_1 w_2 w_3$

	1	2	3
0	$C[0, 1]$	$C[0, 2]$	$C[0, 3] \leftarrow \text{goal}$
1		$C[1, 2]$	$C[1, 3]$
2			$C[2, 3]$

Fill: diag first, then by span length

7. Dependency Parsing

Dependency Tree: Directed spanning tree, root degree 1 **Constraints:** Single head; Connected; Acyclic **Projective:** arcs 不交叉 (嵌套/并列) → CKY 可用 **Non-projective:** arcs 可交叉 → 必须用 CLE/MMT # spanning trees: $O((n-1)^{n-2})$

Edge-Factored Model: 优点: 全局优化 分解为局部边决策 局限: 无法捕捉 sibling/grandparent effects $\text{score}(t, w) = \sum_{(i,j) \in t} \text{score}(i \rightarrow j, w) + \text{score}(r, w) \quad p(t|w) = \frac{1}{Z} \prod_{(i,j) \in t} \exp(\text{score}(i, j, w)) \exp(\text{score}(r, w))$

Cayley Formula: 无向 K_n : n^{n-2} 棵 spanning trees 有向+固定 root: n^{n-2} 棵 arborescences 有向+任意 root: $n \times n^{n-2} = n^{n-1}$ 棵

Graph Laplacian L : $L_{ij} = \begin{cases} \text{Degree}(i) i=j \text{ (对角线)} \\ -1 i \neq j, i \sim j \text{ (有边)} \\ 0 \text{ otherwise} \end{cases}$

trick: 只看非对角-1判断边存在 MTT: #spanning trees = $\det(\tilde{L})$ (any minor)

Weighted Laplacian (MTT): $A_{ij} = \exp(\text{score}(i \rightarrow j))$, $\rho_j = \exp(\text{score}(j, w)) \quad L_{ij} = \begin{cases} \rho_j & i=1 \text{ (root row)} \\ \sum_{k \neq j} A_{kj} & k \in \text{in-degree} \\ -A_{ij} & \text{else} \\ O(n^3) \end{cases}$

CLE Algorithm: Goal: max spanning arborescence (directed MST)

1. For each node v , pick max incoming edge
2. If no cycle → done (it's a tree)

3. If cycle → **contract** cycle to supernode

4. **Reweighting:** $\omega'(u \rightarrow v) = \omega(u \rightarrow v) - \omega_{\text{in-cycle}(v)}$

5. Recursively solve contracted graph

6. **Expand:** break cycle at min-loss edge

Complexity: $O(N^2)$ or $O(E + N \log N)$

Root Constraint: CLE base 允许多 root outgoing arcs

Naive: 对每条 root arc 分别运行 CLE → $O(N \cdot \text{CLE})$

Clever (Gabow): swap score=next-best - current 删除 swap score 最小的多余 root edge

MTT vs CLE: [维度], [MTT], [CLE], [目标], [Z = $\sum_t \exp(\text{score})$], [$t^* = \arg \max$], [算法], [$\det(L)$], [Greedy+Contract], [复杂度], $[O(N^3)]$, $[O(N^2)]$

8. Semantic Parsing

Syntax vs Semantics: Syntax: structural org (parse tree)

Semantics: underlying meaning **Logical form:** quantifiers, vars, boolean, predicates **Compositionality:** meaning of whole = fn of parts

Lambda Calculus: Terms: 变量 x ; 抽象 $\lambda x.M$; 应用 (MN) **β -reduction:** $(\lambda x.M)N \rightarrow M[x := N]$ α

β -conversion: 重命名 bound 变量 避免 capture β -infinity: $F = \lambda x.((xx)x)$, $FF = \dots$ 不终止

β -reduction 步骤:

1. 找到 $(\lambda x.M)N$ 形式的 redex
2. 在 M 中找所有被该 λx 绑定的 x
3. 将这些 x 替换为 N

注意: 可能需先 α -convert 避免变量捕获!

Free vs Bound Variables: $\text{FV}(x) = \{x\}$

$\text{FV}(\lambda x.M) = \text{FV}(M) - \{x\}$ $\text{FV}(MN) = \text{FV}(M) \cup \text{FV}(N)$

Bound: 在某 λ 的 scope 内 **Free:** 不在任何 abstraction 的 scope 内

Combinatory Logic: $Ix = x$; $Kxy = x$; $Sxyz = xz(yz)$ $Bxyz = x(yz)$ (comp); $Cxyz = xzy$ (flip)

$Txy = yx$ (type-raising) $I = SKK$ (S,K 构成 complete basis)

CCG Rules: Application: $X/Y Y \Rightarrow X$ (>前向); $Y X \setminus Y \Rightarrow X$ (<后向) Composition: $X/Y Y/Z \Rightarrow X/Z (B_>)$ Type-raising: $X \Rightarrow T/(T \setminus X)$ ($T_>$) rules is universal, language-specific 全在 lexicon

CCG Category 直觉: $S \setminus NP$: 左边要 NP → 产出 S (intransitive) ($S \setminus NP$) / NP: 右边要 NP → $S \setminus NP$ (transitive) **Slash 方向:** / 向右找 arg; \ 向左找 arg

LIG 构造策略: 问题: CFG 无法“计数”($a^n b^n c^n$ 中 n 相等) LIG: 用 stack 记录计数信息. 策略 1: 两端向中间 -先生成首尾, 再生成中间 策略 2: 左向右-前半部分 push, 后半部分 pop Example $a^n b^n c^n d^n$: $S[\sigma] \rightarrow aS[f\sigma]d; S[\sigma] \rightarrow T[\sigma] T[f\sigma] \rightarrow bT[\sigma]c; T[] \rightarrow \varepsilon$

FOL Translation: \forall 配 \Rightarrow : 全称限定条件; \exists 配 \wedge : 存在某具体对象; 否则 \exists 配 \Rightarrow : 往往荒谬; \vee 配 \wedge : 要求满足多个条件.

9. WFST & Lehmann

Transducer Def: $T = \langle Q, \Sigma, \Omega, \lambda, \rho, \delta \rangle$ Q : states; Σ : input; Ω : output $\lambda : Q \rightarrow \mathbb{R}$: initial; $\rho : Q \rightarrow \mathbb{R}$: final $\delta : Q \times (\Sigma \cup \varepsilon) \times (\Omega \cup \varepsilon) \times Q \rightarrow \mathbb{R}$ **ε -transition:** no input/output consumed

FSA vs FST: WFSA (单带): read only, score(π) = $\sum_n \text{score}(\tau_n)$ WFST (双带): read input + write output **Unambiguous:** $|\Pi(x, y)| \leq 1$ **Ambiguous:** $|\Pi(x, y)| > 1 \rightarrow$ need semiring

Path Score: $\text{score}(\pi) = \lambda(q_{\text{start}}) + \sum_{n=1}^{|\pi|} \text{score}(\tau_n) + \rho(q_{\text{end}})$ $p(y|x) = \frac{1}{Z} \sum_{\pi \in \Pi(x, y)} \exp(\text{score}(\pi))$ $Z = \sum_{y' \in \Omega^*} \sum_{\pi'} \exp(\text{score}(\pi'))$ (infinite!)

Matrix Mult View: $C = A \otimes B$: $C_{ij} = \bigoplus_k (A_{ik} \otimes B_{kj})$ **Tropical:** $C_{ij} = \min_k (A_{ik} + B_{kj})$ **Inside:** $C_{ij} = \sum_k (A_{ik} \times B_{kj})$ Naive $W^N: O(N^4) \rightarrow$ Lehmann fixes to $O(N^3)$

Floyd-Warshall: Key: allow 中间 node k incrementally $\text{dist}_k[i][j] = \min(\text{dist}_{k-1}[i][j], \text{dist}_{k-1}[i][k] + \text{dist}_{k-1}[k][j])$ Runtime: $O(N^3)$ FW 是 Lehmann 在 Tropical 的特例 ($a^* = 0$)

Lehmann algo: Generalized FW for any closed semiring:

$$W_{ij}^{(k)} = W_{ij}^{(k-1)} \oplus W_{ik}^{(k-1)} \otimes (W_{kk}^{(k-1)})^* \otimes W_{kj}^{(k-1)}$$

定义: $R_{ik}^{(j)}$ = 从 q_i 到 q_k , 仅经过 $\{q_1, \dots, q_j\}$ 的 paths 的 semiring-sum **直觉:** 经过 $\{1, \dots, j\}$ 的 paths = 不经 j \oplus 经 j 后者分解: $i \rightarrow j$ (不经 j) + j 上 cycles + $j \rightarrow k$ (不经 j) Runtime: $O(|Q|^3)$

Pathsum & Z: $Z(\mathcal{T}) = \bigoplus_{i,k \in Q} \lambda(q_i) \otimes R_{ik} \otimes \rho(q_k)$ $Z = \alpha^\top \left(\bigoplus_{\omega \in \Sigma^*} W^{(\omega)} \right)^* \beta$ Why Lehmann? Direct sum over infinite paths impossible

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

Transliteration: 3 transducers cascade $\mathcal{T}_x \circ \mathcal{T}_\theta \circ \mathcal{T}_y$

Acyclic WFSA Backward: 前提: DAG 可做 topological sort

1. 按 reverse topo order 遍历 nodes q_M, \dots, q_1
2. $\beta[q_m] \leftarrow \rho(q_m) \oplus \bigoplus_{(q_m, a, w, q') \in \delta} w \otimes \beta[q']$
3. return $\bigoplus_{q \in I} \lambda(q) \otimes \beta[q]$

Complexity: $O(|Q| + |\delta|)$ (linear!)

10. Transformers & MT

Seq2Seq: $z = \text{encoder}(x), y|x \sim \text{decoder}(z)$ $p(y|x) = \prod_{t=1}^T p(y_t|x, y_1, \dots, y_{t-1})$ **Information Bottleneck:** z fixed-length → Attention 解决

Attention: $\alpha^T V = \sum_i \alpha_i v_i^T$ (soft retrieval) $\alpha_i = \text{softmax}(\text{score}(q, k_i))$ $K = V = H^{(e)}$, $q_t = h_t^{(d)}$, $c = \alpha^T V$

Self-Attention: $Q = XW_Q$, $K = XW_K$, $V = XW_V$ SelfAtt = $\text{softmax}\left(\frac{QK^\top}{\sqrt{d_k}}\right)V \sqrt{d_k}$: 防止 dot product 过大 导致 softmax 饱和 **Complexity:** $O(nd^2 + dn^2)$

Positional Encoding: $P_{p,2i} = \sin(p/10000^{2i/d})$ $P_{p,2i+1} = \cos(p/10000^{2i/d})$ motiv: Transformer 无 recurrence, 无法区分位置

Encoder-Decoder 架构: **Encoder:** $+P \rightarrow \text{MHSA} \rightarrow + \rightarrow \text{LN} \rightarrow \text{MLP} \rightarrow + \rightarrow \text{LN}$ **Decoder:** $+ \text{masked self-attn} + \text{cross-attn}$ **Masked:** 只 attend 到左边 positions (causal) **Cross-attn:** Q 来自 decoder, K, V 来自 encoder **Residual:** $x + \text{Layer}(x)$ 缓解 vanishing gradient

Decoding Strategies: $y^* = \arg \max_{y \in \mathcal{Y}} \text{score}(x, y)$ W/o assumptions: $O(|\Sigma|^{n_{\text{max}}})$ paths **Greedy:** 每步 arg max (次优, 快) **Beam:** 保持 k -best candidates **Nucleus/Top-p:** 从累积 prob $\geq p$ 的 tokens 中 sample **Temperature:** $T < 1$ sharper; $T > 1$ uniform **Evaluation:** BLEU (n-gram overlap), METEOR

MT Pipeline:

1. **Tokenize:** subword (BPE/WordPiece)
 2. **Embed:** token → vector + positional
 3. **Encode:** Transformer encoder
 4. **Decode:** autoregressive, $p(y_n|y_{<n}, z)$
 5. **Search:** beam/nucleus sampling
- Train:** MLE, $-\sum \log p(y_n|y_{<n}, x)$

11. Modeling Choices

Prob vs Non-Prob: **Prob:** leverage prob theory, needs assumptions CRF, RNN, N-gram models **Non-Prob:** interpretable, uncertainty unclear Perceptron, SVM, CFG rules

Disc vs Generative: **Discriminative:** model boundary $p(y|x)$ **Generative:** model own dist $p(x, y)$

Local vs Global Norm: **Local:** efficient train, biased predictions **Global:** needs Z , unbiased Independence assumptions control complexity

Regularization: **LogLoss:** $\ell(y, y') = \log(1 + e^{-y \cdot y'})$

Exp-Loss: $\ell(y, y') = e^{-y \cdot y'}$ **L1/L2:** weight penalties (Laplace/Gaussian prior)

Evaluation Metrics: **Prec:** $P_{\text{true}}/P_{\text{all}}$; **Recall:** $P_{\text{true}}/(P_{\text{true}} + N_{\text{false}})$ **Acc:** $(P_{\text{true}} + N_{\text{true}})/N$ **F-score:** $((1 + \beta^2)(\text{prec} \cdot \text{recall})) / (\beta^2 \text{prec} + \text{recall})$

Statistical Tests: $p = 2 \min(P(T \geq t|H_0), P(T \leq t|H_0))$; Rej if $p < \alpha$ **Power:** $P(\text{reject } H_0|H_1)$ **Multiple tests:** $P(|\text{FalseRej}| > 0) = 1 - (1 - \alpha)^K$ **Bonferroni:** $\alpha^* = \alpha/K$ **McNemar:** $\chi^2 = \frac{(b-c)^2}{b+c} \sim \chi_1^2$

12. Bias & Fairness

Bias Sources: **Labeling:** reproduce annotator bias **Sample selection:** training fits certain profile

Task definition: excludes certain groups **Imbalanced test:** loss ignores minorities

Ethical Frameworks: **Consequentialism:** best consequence **Utilitarianism:** hedonistic/preference/welfare

Deontology: rules must be kept **Social Contract:** natural equality **Anti-subordination:** positive discrimination for equality

Quick Ref: Chain: $\frac{d}{dx}[f(g(x))] = f'(g)g'(x)$; Bauer: sum over all paths Softmax: $\exp(h_y) / \sum \exp(h_{y'})$; $T \rightarrow 0 = \arg\max$ Log-Linear: $p(y|x) = \exp(\theta \cdot f)/Z$; MLE matches expected features DP: distrib 把 $O(|T|^N) \rightarrow O(N|T|^2)$; 3-gram 则 $O(N|T|^3)$ Fwd/Bwd: Fwd init BOS+sum last col; Bwd init 1+single value Viterbi: max instead of sum + backpointer CKY: $O(N^3|R|)$; CNF: diag first, span 递增 MTT: $Z = \det(L)$ in $O(n^3)$; CLE: greedy+contract $O(n^2)$ Lehmann: $R^{(j)} = R^{(j-1)} \oplus R \otimes R^* \otimes R$; $O(|Q|^3)$ Kleene: Inside $\frac{1}{1-a}$; Tropical 0 if $a \geq 0$ CCG: $X/Y Y \Rightarrow X(>)$; $Y X \setminus Y \Rightarrow X(<)$ β -reduce: $(\lambda x.M)N \rightarrow M[x := N]$; 先 α -convert 避免捕获 Self-Attn: softmax $(QK^\top/\sqrt{d})V$; $O(nd^2 + dn^2)$ Cayley: 固定 root n^{n-2} ; 任意 root n^{n-1}

Abbrev: BOS/EOS: Begin/End of Sentence; CCG: Combinatory Categorical Grammar; CFG: Context-Free Grammar; CKY: Cocke-Kasami-YOUNGER; CNF: Chomsky Normal Form; CRF: Conditional Random Field; DP: Dynamic Programming; LIG: Linear Indexed Grammar; MLE: Max Likelihood Est; MST: Min Spanning Tree; MTT: Matrix-Tree Theorem; POS: Part-of-Speech; RNN: Recurrent Neural Network; WFST: Weighted Finite State Transducer;