

1. 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: $O(nml)$

$$\text{NLL } \nabla = \mathbf{0}: \sum_{i=1}^n f(x_i, y_i) = \sum_{i=1}^n \mathbb{E}_{y|x_i} [f(x_i, y)]$$

Observed features = Expected features **Hessian:** $\mathbf{H} = \sum_i \text{Cov}_{y|x_i} [f(x_i, y)]$ (PSD!)

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:**

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

Bauer Path: $\frac{dy_i}{dx_j} = \sum_{p \in \mathcal{P}(j,i)} \prod_{(k,l) \in p} \frac{dz_l}{dx_k}$ $\mathcal{P}(j,i)$ =all paths $j \rightarrow i$; worst $O(m^n)$, m 平均出度, n 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: $p(y|x, \theta) = \frac{\exp(\theta \cdot f(x, y))}{Z(\theta)}$ $Z(\theta) = \sum_{y' \in \mathcal{Y}} \exp(\theta \cdot f(x, y'))$ $\log p(y|x, \theta) = \theta \cdot f(x, y) - \log Z$ (linear in log space!) **Discrete MLE:** $p(y|x) = \frac{\text{count}(x, y)}{\text{count}(x)}$ (sparse 问题)

MLE $\nabla: \theta_{\text{MLE}} = \arg \min_{\theta} -\sum_{n=1}^N \log p(y_n|x_n, \theta)$

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: $\text{sftm}(h, y, T) = \frac{\exp(h_y/T)}{\sum_{y'} \exp(h_{y'}/T)}$ $T \rightarrow 0$:

argmax; $T \rightarrow \infty$: uniform $\log \text{sftm} = h_y - \log \sum_{y'} \exp(h_{y'})$ (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: $\text{sftm}(\theta^\top h_n)$

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) = $\frac{\exp(w_{y_{t-1}})}{\sum_{y' \in V} \exp(w_{y' \cdot h_t})}$, $h_t \in \mathbb{R}^d$ **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 | $\mathbb{R} \cup \{\infty\}$ | 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: $0 \otimes x = x \otimes 0 = 0$

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: $a \otimes b \in \mathbb{K}$
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: $G = \langle \mathcal{N}, \mathcal{S}, \Sigma, \mathcal{R} \rangle$ 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:** $p(t) = \frac{1}{Z} \prod_{r \in t} \exp(\text{score}(r))$ $Z = \sum_{t' \in \mathcal{T}} \prod_{r' \in t'} \exp(\text{score}(r'))$ (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 $i = 1, \dots, N-1$: span: $k \leftarrow i+1$; for $j = i+1, \dots, k-1$; for $X \rightarrow Y Z \in \mathcal{R}$: $C[i, k, X] \oplus \exp(\text{score}) \otimes C[i, j, Y] \otimes C[j, k, Z]$

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

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

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: Arc-factored: 每条边独立打分, 树 score=边 score 之和 优点: global 优化分解为 local 边决策 局限: 无法捕捉 sibling/grandparent effects $\text{score}(t, w) = \sum_{(i \rightarrow j) \in t} \text{score}(i \rightarrow j, w) + \text{score}(r, w)$ $p(t|w) = \frac{1}{Z} \prod_{(i \rightarrow j) \in t} \exp(\text{score}(i, j, w)) \exp(\text{score}(r, w))$

CLE 关键步骤: 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 $O(E + N \log N)$

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

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

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

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], $[\det(\tilde{L})]$, [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:** $(\lambda x.M)N \rightarrow M[x := N]$

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

