

1 Basics

1.1 Probability

Bayes rule: $p(y|x) = \frac{p(x|y)p(y)}{p(x)} = \frac{p(x|y)p(y)}{\int p(x|y')p(y')dy'}$
 $p(y|x)$ **posterior**, $p(x|y)$ **likelihood**, $p(y)$ **prior**,
 $p(x)$ or $\int p(x|y')p(y')dy'$ **marginal**

1.2 Activation Functions

ReLU(x): $\max(x, 0)$

ReLU'(x): $\begin{cases} 1 & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}$

Leaky ReLU(x): $\begin{cases} x & \text{if } x \geq 0 \\ 0.01x & \text{if } x < 0 \end{cases}$

Leaky ReLU'(x): $\begin{cases} 1 & \text{if } x \geq 0 \\ 0.01 & \text{if } x < 0 \end{cases}$

Sigmoid: $\sigma(x) = \frac{1}{1+\exp(-x)}$

$\sigma'(x) = \sigma(x)(1 - \sigma(x))$

Hyperbolic tangent: $\tanh(x) = \frac{\exp(x)-\exp(-x)}{\exp(x)+\exp(-x)}$

$\tanh'(x) = 1 - \tanh(x)^2$

2 Semirings

Semiring	Set	\oplus	\otimes	$\bar{0}$	$\bar{1}$	intuition/application
Boolean	$\{0, 1\}$	\vee	\wedge	0	1	logical deduction, recognition
Viterbi	$[0, 1]$	max	\times	0	1	prob. of the best derivation
Inside	$\mathbb{R}^+ \cup \{+\infty\}$	+	\times	0	1	prob. of a string
Real	$\mathbb{R} \cup \{+\infty\}$	min	+	$+\infty$	0	shortest-distance
Tropical	$\mathbb{R}^+ \cup \{+\infty\}$	min	+	$+\infty$	0	with non-negative weights
Counting	\mathbb{N}	+	\times	0	1	number of paths

- $(R, +)$ is a **commutative monoid** with **identity element** 0:
 - $(a + b) + c = a + (b + c)$
 - $0 + a = a = a + 0$
 - $a + b = b + a$
- (R, \cdot) is a **monoid** with identity element 1:
 - $(a \cdot b) \cdot c = a \cdot (b \cdot c)$
 - $1 \cdot a = a = a \cdot 1$
- Multiplication left and right **distributes** over addition:
 - $a \cdot (b + c) = (a \cdot b) + (a \cdot c)$
 - $(a + b) \cdot c = (a \cdot c) + (b \cdot c)$
- Multiplication by 0 **annihilates** R:
 - $0 \cdot a = 0 = a \cdot 0$

3 Log-Linear Models

Log-lin. model: $p(y|x, \theta) = \frac{1}{Z(x, \theta)} \exp(\theta f(x, y))$
 $\log(p(y|x, \theta)) = \theta f(x, y) + \text{const.}$
where $f : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^K$ is the feature function
 $Z(\theta) = \sum_{y' \in \mathcal{Y}} \exp(\theta f(x, y'))$ is the partition function.
MLE: $\mathcal{L}(\theta) = \sum_{n=1}^N \log p(y_n|x_n, \theta)$

$\theta_{\text{MLE}} = \arg \max_{\theta \in \Theta} \mathcal{L}(\theta)$
where Θ is compact subset of \mathbb{R}^K
 $\nabla \mathcal{L}(\theta) = \sum_{i=1}^N f(x_n, y_n) - \sum_{n=1}^N \mathbb{E}_{Y \sim p(\cdot|x_n, \theta)}[f(x_n, Y)]$
 $\sum_{n \leq N} f(x_n, y_n) = \sum_{n=1}^N \mathbb{E}_{Y \sim p(\cdot|x_n, \theta)}[f(x_n, Y)]$

4 Viterbi

$t^* = \arg \max_{t \in \mathcal{T}^N} \exp(\text{score}(t, w)) = \arg \max_{t \in \mathcal{T}^N} \prod_{n=1}^N \exp\{\text{score}(\langle t_{n-1}, t_n \rangle, w)\}$

Algorithm 5.2

def VITERBI ALGORITHM(w, T, N):
 for t_{N-1} ∈ T:
 v(w, t_{N-1}, N - 1) ← exp(score((t_{N-1}, eos), w))
 end for
 for n ∈ N - 2, ..., 1:
 for t_n ∈ T:
 v(w, t_n, n) ← max_{t_{n+1} ∈ T} exp(score((t_n, t_{n+1}), w)) × v(w, t_{n+1}, n + 1)
 b(t_n, n) ← argmax_{t_{n+1} ∈ T} exp(score((t_n, t_{n+1}), w)) × v(w, t_{n+1}, n + 1) ▷ Keeps track of the best tags
 end for
 end for
 v(w, eos, 0) ← max_{t₁ ∈ T} (v(w, eos, 0), exp(score((eos, t₁), w)) × v(w, t₁, 1))
 b(eos, 0) ← argmax_{t₁ ∈ T} (v(w, eos, 0), exp(score((eos, t₁), w)) × v(w, t₁, 1))
 for n ∈ 1, ..., N:
 t_n ← b(t_{n-1}, n - 1) ▷ Recovers the best tagging sequence using backpointers
 end for
 return t_{1:N}, v(w, eos, 0)

Runtime: $\mathcal{O}(N|T|^2), \mathcal{O}(N|T|^3)$ when considering triplets.

5 Grammar

Definition 6.2

A context-free grammar G is a quadruple $\langle \mathcal{N}, S, \Sigma, \mathcal{R} \rangle$ consisting of:

- A finite set of non-terminal symbols \mathcal{N} ; written in upper-case letters, e.g. N_1, N_2, N_3
- A distinguished start non-terminal S
- An alphabet of terminal symbols Σ ; written as lower-case letters, e.g. a_1, a_2, a_3
- A set of production rules \mathcal{R} of the form $N \rightarrow \alpha$, where $N \in \mathcal{N}$ and $\alpha \in (\mathcal{N} \cup \Sigma)^*$ (Kleene closure of $\mathcal{N} \cup \Sigma$)

5.1 Chomsky Normal Form

Grammar is in Chomsky Normal Form (CNF) if RHS of every production rule includes either two non-terminals or a single terminal symbol: $N_1 \rightarrow N_2N_3$ or $N \rightarrow a$

5.2 PCFG

Probabilistic Context Free Grammar (PCFG) $\langle \mathcal{N}, S, \Sigma, \mathcal{R}, \mathcal{P} \rangle$, where \mathcal{P} are probabilities assigned to each production rule.

5.3 WCFG

Weighted Context Free Grammar (WCFG) $\langle \mathcal{N}, S, \Sigma, \mathcal{R}, \mathcal{W} \rangle$, where \mathcal{W} are non-negative weights assigned to each production rule. PCFG is special case of WCFG.

6 CKY

Algorithm 6.1

def WEIGHTEDCKY(s, $\langle \mathcal{N}, S, \Sigma, \mathcal{R} \rangle$, score):
 N ← |s|
 chart ← 0
 for n = 1, ..., N:
 for X → s_n ∈ R:
 chart[n, n + 1, X] += exp(score(X → s_n)) ▷ Handles single word tokens
 end for
 end for
 for span = 2, ..., N:
 for i = 1, ..., N - span + 1: ▷ i marks the beginning of the span
 k ← i + span ▷ k marks the end of the span
 for j = i + 1, ..., k - 1: ▷ j marks the breaking point of the span
 for X → Y Z ∈ R:
 chart[i, k, X] += exp(score(X → Y Z)) × chart[i, j, Y] × chart[j, k, Z]
 end for
 end for
 end for
 end for
 return chart[1, N + 1, S]