

SR $(\mathcal{A}, \otimes, \bar{0}, \bar{1})$ | • **Commutative**: $a \otimes b = b \otimes a$ • **Associative**: $(a \otimes b) \otimes c = a \otimes (b \otimes c)$ • **Monoid**: Associative, $a \otimes \bar{1} = a$ where $\bar{1}$ is identity operator | • $(\mathcal{A}, \otimes, \bar{0})$ is a commutative monoid $(\mathcal{A}, \otimes, \bar{1})$ is a monoid • \otimes distributes over \oplus : $(a \oplus b) \otimes c = (a \otimes c) \oplus (b \otimes c)$ • $a \otimes \bar{0} = \bar{0}$ • **Commutative SR**: $a \otimes b = b \otimes a$ | • **Idempotent SR**: $a \otimes a = a$ | • $\bigoplus_{i=1}^K \mathbf{M}^k = (\mathbf{I} + \mathbf{M})^K$ | • **Closed SR**: **Kleene star**: $x^* = \bigoplus_{n=0}^{\infty} x^{\otimes n} = x \otimes x \otimes \dots = \bar{1} \oplus x \oplus x^{\otimes 2} \oplus x^{\otimes 3} \oplus \dots$ which must fulfill $\bar{1} \oplus x \otimes x^* = \bar{1} \oplus x^* \otimes x$, e.g. • **Log-sum-exp SR**: $x^+ = \log(\sum_{i=1}^n e^{x_i}) = \log(\frac{1}{1-e^{-x}})$ for $x < 0$ • **Inside SR**: $x^+ = \sum_{n=0}^{\infty} x^n = \frac{1}{1-x}$ for $x \in (0,1)$ • **Boolean SR**: $x^+ = \bar{1} \oplus x$ • **Boolean bounded SR**: $\bar{0} \oplus x = \bar{1}$ • $x^+ = \bigoplus_{n=0}^{N-1} x^{\otimes n}$ • **Idempotent**: e.g. tropical, arctic SR | **Common SRs**: • **Boolean**: $(\{0,1\}, \vee, \wedge, 0, 1)$ • **Inside**: $(\mathbb{R} \cup \{\infty\}, +, \times, 0, 1)$ • **Log-sum-exp**: $(\mathbb{R} \cup \{-\infty\}, \otimes_{\log}, +, -\infty, 0)$ where $a \otimes_{\log} b = \log(e^a + e^b)$ • **Viterbi**: $([0,1], \max, \times, 0, 1)$ • **Arctic**: $(\mathbb{R} \cup \{-\infty\}, \max, +, -\infty, 0)$ • **Distributions** **Gaussian** — $X \sim N(\mu, \sigma^2)$ | $\frac{1}{\sigma\sqrt{2\pi}} \exp(-\frac{(x-\mu)^2}{2\sigma^2}) = \frac{1}{\sigma\sqrt{2\pi}} \exp(-\frac{x^2}{2\sigma^2} + 2x\frac{\mu}{2\sigma^2} - \frac{\mu^2}{2\sigma^2})$ | $\frac{1}{2\pi n^{1/2}} \frac{1}{|\Sigma|^{1/2}} \exp(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu))$ • **Exponential Families** — $p(x|\theta) = \frac{1}{Z(\theta)} h(x) \exp(\theta \cdot \Phi(x))$ | **Gaussian**: $Z(\sigma, \mu) = \frac{\exp(-\frac{\mu^2}{2\sigma^2})}{\sqrt{2\pi\sigma^2}}$ • $h(x) = 1$ • $\theta = \begin{bmatrix} \mu/\sigma^2 \\ -1/2\sigma^2 \end{bmatrix}$ • $\Phi(x) = \begin{bmatrix} x \\ x^2 \end{bmatrix}$ | **Beta**: • $Z(\theta) = B(\theta_1 + 1, \theta_2 + 1)$ • $h(x) = 1$ • $\theta = \begin{bmatrix} (\alpha-1) \\ (\beta-1) \end{bmatrix}$ • $\Phi(x) = \begin{bmatrix} \ln(x) \\ \ln(1-x) \end{bmatrix}$ | **Binomial**: • $Z(\theta) = e^{n \ln(1+e^\theta)}$ • $h(x) = \binom{n}{x}$ • $\theta = \ln(\frac{1-p}{p})$ • $\Phi(x) = [x]$ • **Hypothesis Testing** • **Errors**: • **TP**: Chose $H_0|H_0$ • **FP**, type I: Chose $H_A|H_0$ • **TN**: Chose $H_A|H_A$ • **FN**, type II: Chose $H_0|H_A$ • **Significance level**: $\alpha = p(\text{FP}) = p(\bar{x} \geq c | H_0) = p(z_n \geq z_\alpha | H_0)$, inflated for K tests $1 - (1 - \alpha)^K$ • **Beta**: $\beta = p(\text{FN})$ • **Critical value**: Threshold c resp. critical value z_α associated with α • **P-value**: $p = P(|z| \geq z_n | H_0)$, smallest significance level, at which we can reject H_0 • **Confidence level**: $1 - \alpha$ • **Power**: $1 - \beta = p(\bar{x} \geq c | H_1) = p(z_n \geq z_\alpha | H_1)$ • **Statistic**: $z = \frac{\bar{x} - \mu_0}{\sigma/\sqrt{n}}$, $z|H_0 \sim N(0,1)$, $z|H_1 \sim N(\frac{\mu_1 - \mu_0}{\sigma/\sqrt{n}}, 1)$ • **Rejection**: $\bar{x} > \mu_0 + \sigma z_\alpha / \sqrt{n}$ • **Other** $S = \sum_{i=0}^{\infty} a_i r^i = \frac{a_1}{1-r}$ for $r < 1$ | $e^{m+n} = e^m e^n$, $\log(mn) = \log m + \log n$ | $\frac{\partial \det x}{\partial x} = a$ | $\det \begin{bmatrix} a & b \\ c & d \end{bmatrix} = ad - bc$

2 ML Paradigms
MLE • **Log likelihood**: $\hat{\theta} = \arg \max_{\theta} (L) = \prod_{i=1}^n p(y_i | x_i, \theta) = \sum_{i=1}^n \log(p(y_i | x_i, \theta))$ • For binary classification, **log loss** = **binary cross entropy** $= -[y \log(p) + (1-y) \log(1-p)]$ • **Estimating Distributions** • **Likelihood**: $L = (-\frac{1}{\sigma})^n \prod_{i=1}^n \exp(-\frac{1}{2\sigma^2}(x^{(i)} - \mu)^2)$ • **Log likelihood**: $LL = -n \log(\sigma) - \sum_{i=1}^n (-\frac{1}{2\sigma^2}(x^{(i)} - \mu)^2)$ • $\mu_{MLE} = \frac{1}{n} \sum_{i=1}^n x^{(i)}$ • $\sigma_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n (x^{(i)} - \mu)^2$ • **Bayesianism** — • **Prior** $p(\theta)$, likelihood $p(x|\theta)$, posterior $p(\theta|x)$ • Let σ^2 be known, Posterior $p(\mu|x, \mu_0, \sigma_0^2) \propto$ Prior $p(\mu|\mu_0, \sigma_0^2)$ • We get (based on the form of the Gaussian) $\mu_P = \frac{n\bar{x}\sigma_0^2 + \mu_0\sigma^2}{n\sigma_0^2 + \sigma^2}$ and $\sigma_P^2 = \frac{\sigma^2\sigma_0^2}{n\sigma_0^2 + \sigma^2}$ • $p(\mu|x, \mu_0, \sigma_0^2) \sim N(\mu_P, \sigma_P^2)$ • **Conjugate prior** • **Linear Regression** **Formulation** $y^{(i)} = \beta \cdot x^{(i)}$ resp. $y = X\beta$ where $X \in n \times m$ | Transformation: Φ with rows $\phi(x^{(i)})^T$ • **Optimization** **Objective function** — **OLS**: Minimize MSE: $LO = (y - X\beta)^T (y - X\beta)$ | **Optimization** — • $\nabla_{\beta} LO = \frac{1}{2} \nabla_{\beta} (\beta^T X^T X \beta - 2y^T X \beta) = X^T X \beta - X^T y = 0 \Rightarrow \beta = (X^T X)^{-1} X^T y$ • **Ridge L2** Find parameters β subject to $\|\beta\|_2^2 \leq t$ | **Objective function** — Lagrangian: $LO = (y - X\beta)^T (y - X\beta) + \lambda (\|\beta\|_2^2 - t)$ • **Optimization** — $\beta = (X^T X + \lambda I)^{-1} X^T y$ | **Characteristics** — • Strictly convex with global min, unique sol (linearly independent cols), and analytic sol (always invertible) • **Bayesian MAP**, when $\lambda = \frac{\sigma^2}{t}$: Posterior $p(\beta|x, y) \propto N(X\beta, \sigma^2 I_n) \times \text{Prior } p(\beta) \sim N(0, \tau^2 I_m)$ • **Lasso L1** Subject to $|\beta| - t \leq 0$ | **Objective function** — Lagrangian: $LO = \dots + \lambda (|\beta| - t)$ | **Characteristics** — Bayesian MAP, when $\lambda = \frac{\sigma^2}{t}$ • **Log Linear Models**

Formulation $p(y|x, \theta) = \frac{\exp(\theta \cdot f(x, y))}{\sum_{y'} \exp(\theta \cdot f(x, y'))}$ • **Optimization** **Objective function** — • Minimize log loss: $\theta = \arg \min_{\theta} - \sum_{i=1}^n \log(p(y^{(i)} | x^{(i)}, \theta)) = - \sum_{i=1}^n \theta \cdot f(x^{(i)}, y^{(i)}) - \log(\sum_{y'} \exp(\theta \cdot f(x^{(i)}, y')))$ | **Optimization** — First-order derivative: $\frac{\partial}{\partial \theta} \log \text{loss} = - \sum_{i=1}^n (f(x^{(i)}, y^{(i)}) - \frac{\sum_{y'} \exp(\theta \cdot f(x^{(i)}, y'))}{\sum_{y'} \exp(\theta \cdot f(x^{(i)}, y'))}) = - \sum_{i=1}^n (f(x^{(i)}, y^{(i)}) - \sum_{y'} p(y' | x^{(i)}, \theta) \cdot f(x^{(i)}, y')) = - \sum_{i=1}^n (f(x^{(i)}, y^{(i)}) - \sum_{i=1}^n \mathbb{E}(f(x^{(i)}, y')))$ • If we set gradient to 0, we have **expectation matching** **Second-order derivative**: $\frac{\partial^2}{\partial \theta \partial \theta^T} \log \text{loss} = \frac{\partial}{\partial \theta^T} \nabla \log \text{loss} = \sum_{i=1}^n \mathbb{E}[f(x^{(i)}, y') f(x^{(i)}, y')^T] - \sum_{i=1}^n (\mathbb{E}[f(x^{(i)}, y')] \mathbb{E}[f(x^{(i)}, y')^T]) = \sum_{i=1}^n \text{Cov}(f(x^{(i)}, y'))$ • **Logistic** **Formulation** • **Sigmoid**: $\sigma(z) = \frac{1}{1+e^{-z}}$ • $P(y=1|x) = \frac{1}{1+e^{-\beta \cdot x}} = \frac{e^{\beta \cdot x}}{1+e^{\beta \cdot x}}$, $P(y=0|x) = \frac{1}{1+e^{\beta \cdot x}} = \frac{e^{-\beta \cdot x}}{1+e^{-\beta \cdot x}}$ • **Odds**: $\frac{P(y=1|x)}{P(y=0|x)} = e^{\beta \cdot x}$, **Log odds**: $\ln(\frac{P(y=1|x)}{P(y=0|x)}) = \beta \cdot x$ • $z = \beta \cdot x$ is a linear hyperplane: When $z > 0$, then odds > 1 , then predict $y=1$ • **Optimization** **Objective function** — • Likelihood: $L = \prod_{i=1}^n \sigma(z^{(i)})^{y^{(i)}} (1 - \sigma(z^{(i)}))^{1-y^{(i)}}$ • Log likelihood: $LL = \sum_{i=1}^n [y^{(i)} \log \sigma(z^{(i)}) + (1-y^{(i)}) \log(1 - \sigma(z^{(i)}))]$ • $\sum_{i=1}^n [y^{(i)} z^{(i)} - \log(1 + e^{z^{(i)}})]$ | **Optimization** — • $\frac{\partial LL}{\partial \beta} = - \sum_{i=1}^n \frac{\partial}{\partial \beta} [y^{(i)} \log \sigma(z^{(i)}) + (1-y^{(i)}) \log(1 - \sigma(z^{(i)}))]$ • $\sum_{i=1}^n [\sigma(z^{(i)}) - y^{(i)}] x^{(i)}$ • **Multinomial Logistic** **Formulation** **Softmax**: $P(y=k|x) = \frac{e^{f_k(x)/T}}{\sum_{j=1}^K e^{f_j(x)/T}} = \frac{e^{\beta_k \cdot x/T}}{\sum_{j=1}^K e^{\beta_j \cdot x/T}}$ • **Softmax** ($T=1$) vs. **argmax** ($T=0$) • **Optimization** Find parameters β_1, \dots, β_K | **Objective function** — • Likelihood: $L = \prod_{i=1}^n \prod_{k=1}^K (-\frac{e^{\beta_k \cdot x^{(i)}}}{\sum_{j=1}^K e^{\beta_j \cdot x^{(i)}}})^{\delta_{\{y^{(i)}=k\}}}$ • Log likelihood: $LL = \sum_{i=1}^n \sum_{k=1}^K \delta_{\{y^{(i)}=k\}} \log \ell_j[\beta \cdot x^{(i)} - \log(\sum_{j=1}^K e^{\beta_j \cdot x^{(i)}})]$ | **Optimization** — • $\frac{\partial LL}{\partial \beta_k} = - \sum_{i=1}^n \delta_{\{y^{(i)}=k\}} x^{(i)} - P(y=k|x) k | x^{(i)} x^{(i)}$ • For reference:

Softmax derivative $\frac{\partial P(y=\ell|x)}{\partial \beta_k}$: If $\ell=k \Rightarrow P(y=\ell|x) (1 - P(y=\ell|x)) x$, if $\ell \neq k \Rightarrow -P(y=\ell|x) P(y=k|x) x$, can be shown via quotient rule • **Neural Networks** **Formulation** **Architecture** — Inputs: n words $>$ Embedding $e(w_i)$ $>$ Concatenation: $e(x) = \frac{1}{n} \sum w_i e(w_i) >$ Hidden layer: $h^{(k)} = \sigma(w^{(k)} h^{(k-1)}) >$ Activation: **Softmax**: $p(y|x) = \frac{\exp(h_y^{(K)})}{\sum_{y'} \exp(h_{y'}^{(K)})}$ | **Neuron** (j) in layer $[k]$: $h^{(j)}[k] = \varphi(h^{(i)}[k-1] \cdot \beta^{(j)}[k])$ | **Activation functions** — **Sigmoid**: • Not convex • $[0,1]$ • $\varphi(z) = \sigma(z) = \frac{1}{1+e^{-z}} = \frac{e^z}{e^z+1}$ • $\varphi'(z) = \frac{e^{-z}}{(1+e^{-z})^2}$ with maximum at 0.25 | **Hyperbolic tangent**: • Not convex • $[-1,1]$ • $\varphi(z) = \tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} = \frac{1 - e^{-2z}}{1 + e^{-2z}}$ • $\varphi'(z) = 1 - \tanh^2(z) = \frac{4e^z}{(e^z + e^{-z})^2}$ | **ReLU**: • Piecewise convex • $\varphi(z) = \max(0, z)$ • $\varphi'(z) = 1$ if $z > 0$; 0 otherwise • **Optimization** • Forward pass to calculate loss • Backpropagation to calculate gradient: $\frac{\partial L}{\partial \mathbf{B}^{[k]}} = \frac{\partial L}{\partial \mathbf{H}^{[l]}} \frac{\partial \mathbf{H}^{[l]}}{\partial \mathbf{B}^{[k]}} = \mathbf{C}$ • When $l > k+1$, i.e. when going several layers back: $\frac{\partial L}{\partial \mathbf{B}^{[k]}} = \frac{\partial L}{\partial \mathbf{H}^{[l]}} \frac{\partial \mathbf{H}^{[l]}}{\partial \mathbf{S}^{[l-1]}} \frac{\partial \mathbf{S}^{[l-1]}}{\partial \mathbf{H}^{[l-1]}} \frac{\partial \mathbf{H}^{[l-1]}}{\partial \mathbf{B}^{[k]}}$ • When $l=k+1$, i.e. when going one layer back: $\frac{\partial L}{\partial \mathbf{B}^{[k]}} = \frac{\partial L}{\partial \mathbf{H}^{[k+1]}} \frac{\partial \mathbf{H}^{[k+1]}}{\partial \mathbf{S}^{[k]}} \frac{\partial \mathbf{S}^{[k]}}{\partial \mathbf{B}^{[k]}}$ • Gradient descent to find best weights • **Backpropagation** **Formulation** **Computation graph** \mathcal{G} — With n input nodes and $|E| = M$ edges, we have $|V| = M + n$ total nodes | Note: $x^{a+b} = \exp(\log(x^{a+b})) = \exp((a+b) \times \log(x))$ and $\frac{a}{b} = a \times b^{-1}$ • **Bauer's Formula** — $P(j,i)$ is set of all paths from node $j \rightarrow i$ • $\frac{\partial z_i}{\partial z_j} = \sum_{p \in P(j,i)} \prod_{(k,l) \in p} \frac{\partial z_l}{\partial z_k}$ • Challenge: Runtime $O(|P(j,i)|) = O(2^{|E|})$ • **Forward Pass** Randomly initialize values of input nodes and then evaluate function • Runtime $O(|E|)$ • Space $O(|V|)$ • **Backpropagation** 1) Perform forward pass 2) For $i=M, \dots, 1$: $\frac{\partial f}{\partial z_i} = \sum_{j: i \in \text{Pa}(i)} \frac{\partial f}{\partial z_j} \frac{\partial z_j}{\partial z_i}$ 4) Return $[\frac{\partial f}{\partial z_1}, \frac{\partial f}{\partial z_2}, \dots, \frac{\partial f}{\partial z_M}]$ | Runtime: • Partial derivatives on multiple paths are memoized • For first-order derivative runtime and

space same as forward pass, for k^{th} order derivative runtime $O(|E|n^{k-1})$ • **RNNs** **Architecture** Input $>$ Hidden layer resp. **memory cell**: **Cell state** h_t , **cell output** $y_t >$ Output | Output of neuron in layer n : $Y_t = \phi(X_t W_x + H_{t-1} W_y + b) V$ • **Optimization** Gradient: • $\nabla_{W_h} L \propto \sum_{k=1}^t (\prod_{i=k+1}^t \frac{\partial h_{i-1}}{\partial h_i}) \frac{\partial h_k}{\partial W_h}$ • $\frac{\partial h_{i+1}}{\partial h_i} = \prod_{j=0}^{k-1} \frac{\partial h_{i+k-j}}{\partial h_{i+k-j-1}}$ • **Attention** **Approach** $Q = E W^q$ resp. $q_i = e_i W^q$ • $E (m \times h)$ • $W_q (h \times d_k) = Q (m \times d_k)$ • $q_i (1 \times d_k)$ • $e_i (1 \times h)$ | $K = E W^k$ • $E (n \times h)$ • $W_k (h \times d_k)$ • $K (n \times d_k)$ • $k_i (1 \times d_k)$ • $e_i (1 \times h)$ | $V = E W^v$ • $E (n \times h)$ • $W_v (h \times d_v)$ • $V (n \times d_v)$ • $v_i (1 \times d_v)$ • $e_i (1 \times h)$ | • $A = \sigma(\frac{QK^T}{\sqrt{d_k}})$ in $(m \times n)$ resp. resp. $z_i = \alpha_i V = \sum_i \alpha_i v_i$ • In CA: Q is decoder input with m , V, K are encoder outputs with n • In SA: Q, V, K are all either encoder or decoder inputs with n or m • In MHA: MHZ = Concat($Z_{\text{head}1}, \dots, Z_{\text{head}_h}$) $W_O + b_O$ where • Concat(...) in $(m \times (n \times n_{\text{heads}}))$ • W_O in $((n_{\text{heads}} \times n) \times d_v)$ • b_O in $1 \times d_v$ | **12 RNNs** **Architecture** • **Encoder**: $h_n^{(d)} = f(W_1^{(e)} h_{n-1}^{(e)} + W_2 w_n)$ • **Decoder**: $h_m^{(d)} = f(W_3^{(d)} h_{m-1}^{(d)} + W_2^{(d)} w'_{m-1} + W_1^{(d)} z_m)$ where $z_m = \sum_{n=1}^N \alpha_{m,n} h_n^{(e)}$ • $h_n^{(e)} = V \circ \alpha_{m,n} = \text{softmax}(h_{m-1}^{(d)} \times [h_1^{(e)}, \dots, h_N^{(e)}])$ with $h_{m-1}^{(d)} = Q$ and $[h_1^{(e)}, \dots, h_N^{(e)}]^T = K$ | Runtime: • Encoder: $O(l_e N d^2)$ from hidden states • Decoder: $O(l_d M d^2 + l_d d N M)$ from hidden states + CA • **Transformers** **Formulation** **Model architecture**: • **Encoder**: • $H_0^{(e)} = X + P \in (N \times d_{\text{model}})$ • Calculate MHA based on $H^{(e)}_{(l-1)}$ • Addition and normalization: $H_l^{(e)} = \text{Layer Norm}(\text{MHA}(Q, K, V) + H^{(e)}_{(l-1)})$ • FFN($H_l^{(e)} = \text{ReLU}(H_l^{(e)} W_1 + b_1) W_2 + b_2$ where $W_1 \in \mathbb{R}^{(d_v \times r)}$ • $b_1 \in \mathbb{R}^{(1 \times r)}$ • $W_2 \in \mathbb{R}^{(r \times d_v)}$ • $b_2 \in \mathbb{R}^{(1 \times d_v)}$ • Add and norm • **Decoder**: • Target sequence inputs are fed into decoder with one time step lag: $H_0^{(d)} = Y + P \in (M \times d_{\text{model}})$ • Calculate masked SA based on $H^{(d)}_{(l-1)}$ • Add and norm on $H^{(d)}_{(l-1)}$

• Encoder outputs are fed into decoder with CA: ■ Calculate CA with $H_l^{(d)}$ for Q and $H^{(e)}_{(N)}$ for K, V • Add and norm • FFN($H_l^{(d)}$) • Add and norm • Linear layer • Softmax layer • Challenge: Output y_t is conditioned on entire history with runtime of $O(|\Sigma|^n)$ • Solution: • **Greedy decoding** • **Beam search**: Calculate local then total probabilities, Keep k -highest-probability tokens at each step, requires normalization • **Language Models** **Formulation** **Prefix tree** | • Challenge: Globally normalized: Since Σ^* is infinite, Z is infinite • Solution: Locally normalized: Assume next word in string is only conditioned on preceding context, normalize over all words in vocabulary • **Tight model**: Locally normalized model, sums to 1, and has finite length paths, enforced by $p(\text{EOS}|\dots) > \delta > 0$ • **N Grams** **Formulation** • String w of length N • Naively: $p(w) = p(w_1) \times p(w_2|w_1) \times p(w_3|w_1, w_2) \times \dots \times p(w_N|w_{<N}) \times p(\text{EOS}|w)$ • **Markov**: Only the last $n-1$ words are considered to predict next word • **N-gram**: • Vocabulary with $|\mathcal{V}|$ words has $|\mathcal{V}|^n$ (distinct) N-grams • Sentence with M words has $M - n + 1$ (potentially overlapping) N-grams • N-gram counts are given by all distinct N-grams and how often they appear in the sentence • Then: $p(w) = (\prod_{i=n}^N p(w_i | w_{1-n+1}, \dots, w_{i-1})) \times p(\text{EOS} | w_{1-n+1}, \dots, w_N)$ • We have: • $|\mathcal{V}|^{n-1} + |\mathcal{V}|^{n-2} + \dots + |\mathcal{V}| + 1 = \sum_{i=0}^{n-1} |\mathcal{V}|^i$ distinct n-grams • $|\mathcal{V}|^n$ distinct conditional probabilities • $|\mathcal{V}| - 1$ parameters for each conditional probability • **Optimization** **Relative frequency counts** — • $p(w_t | \dots) = \frac{\text{count}(w_t - n + 1, \dots, w_t - 1, w_t)}{\text{count}(w_t - n + 1, \dots, w_t - 1)}$ • OOV: • Given a sentence with $M+1$ words and vocabulary \mathcal{V} , probability that a given bigram does not appear in any of the M spots: $(1 - \frac{1}{|\mathcal{V}|^2})^M$ • Solution: **Smoothing**: $p(w_t | w_{1-n+1}, \dots, w_{t-1}) = \frac{\text{count}(w_t - n + 1, \dots, w_{t-1} + 1) + \frac{1}{|\mathcal{V}|}}{\text{count}(w_t - n + 1, \dots, w_{t-1} + 1) + |\mathcal{V}|}$ • **Back-off**: Revert • **Interpolation**: Weight **Neural nets** — $p(w_t | \dots) = \frac{\exp(v(w_t) \cdot h_t)}{\sum_{w' \in \mathcal{V}} \exp(v(w') \cdot h_t)}$ where $v(w)$ are word vectors, h_t are

N-gram context vectors • **Evaluation** **Perplexity**: $p(w_1, \dots, w_N)^{-\frac{1}{N}}$, if lower, likelihood higher • **POS Tagging** **Formulation** **Point of departure** — • $p(t|w) = \frac{\exp(\text{score}(t, w))}{\sum_{t'} \exp(\text{score}(t', w))}$ • Challenge: Z in $O(|T|^N)$ runtime • Solution: Scoring function that is additively decomposable over tag bigrams: $\text{score}(t, w) = \sum_{n=1}^N \text{score}(\langle t_{n-1}, t_n, w \rangle)$ • Then: $p(t|w) = \frac{\exp(\sum_{n=1}^N \text{score}(\langle t_{n-1}, t_n, w \rangle))}{\sum_{t'} \exp(\sum_{n=1}^N \text{score}(\langle t'_{n-1}, t'_n, w \rangle))}$ • **Backward resp. Viterbi algorithm** — 1) For t_{N-1} (all tags): • $\beta(w, t_{N-1})$ resp. $v(w, t_{N-1}) \leftarrow \exp(\text{score}(\langle t_{N-1}, \text{EOS}, w \rangle))$ • $b(t_{N-1}) \leftarrow \text{EOS}$ since t_N can only be EOS 2) For $n \in N-2, \dots, 1$: For $t_n \in T$ (all tags): • $\beta(w, t_n)$ resp. $v(w, t_n) \leftarrow \bigoplus_{t_{n+1}} [\exp(\text{score}(\langle t_n, t_{n+1}, w \rangle)) \otimes v(w, t_{n+1})]$ • $b(t_n) \leftarrow \text{argmax}_{t_{n+1}} [\exp(\text{score}(\langle t_n, t_{n+1}, w \rangle)) \otimes v(w, t_{n+1})]$ 3) Finally: • $\beta(w, t_0)$ resp. $v(w, t_0) \leftarrow \bigoplus_{t_1} [\exp(\text{score}(\langle t_0, t_1, w \rangle)) \otimes v(w, t_1)]$ 4) For $n=1, \dots, N$: Recover the best sequence $t_n \leftarrow b(t_{n-1})$, starting with $t_1 = b(\text{BOS}), t_2 = b(t_1), \dots$ 5) Return $t_{1:N}$ and $\beta(w, t_0)$ resp. $v(w, t_0)$ | • $\beta(w, t_n)$ are **backward variables** that contain the sum of the scores of all paths starting at EOS and ending at tag t_n • $v(w, t_n)$ are **Viterbi variables** that contain the score of t^* starting at EOS and ending with tag t_n • $b(t_n)$ are **backpointers** • Runtime: • For backpointers: $O(N)$ • For tag N-grams: $O(N|T|^n)$ • With forward and backward, we can compute: $p(t=t_N|w) = \frac{\beta(t_N) \alpha(t_N)}{\sum_{t'} \beta(t') \alpha(t')}$ • Alternatively, **forward algorithm** • Starting from t_1 and going forward towards EOS t_N • Instead of $\langle t_n, t_{n+1} \rangle$ we look at $\langle t_{n-1}, t_n \rangle$ • Variables start at BOS and end with tag t_n • Backpointers: $t_{N-1} \leftarrow b(\text{EOS}), t_{N-2} \leftarrow b(t_{N-1}), \dots$ | Extension of backward algorithm in the expectation SR: • Instead of $\omega = \exp(\text{score}(\langle t_{n-1}, t_n, w \rangle))$, we do: $\omega = \langle \omega, -\omega \log(\omega) \rangle$ • Backward algorithm with these weights in expectation SR computes $\langle Z, H_u \rangle$ where $H_u =$

$\sum_{t \in T} \exp(\text{score}(t, w)) \times \text{score}(t, w)$ is the unnormalized entropy

- Normalized entropy
- $H_n = -\sum_t p(t) \log(p(t))$ as $Z^{-1} H_n + \log(Z)$ • Acts as regularizer in log likelihood: $+ \lambda H_n$

Generalized algorithm —

- Backward algorithm: Computes Z, inside SR
- Viterbi algorithm: Computes the score of i^* and recovers sequence, Viterbi SR
- Challenge: Vanishing probabilities
- Solution: Log-sum-exp SR (backward) resp. arctic SR (Viterbi), returns log

Scoring functions — **Hidden Markov Model**: $\text{sc.}((t_n, t_{n+1}), w) = \text{tr.}(t_n, t_{n+1}) + \text{em.}(t_{n+1}, w_{n+1})$ resp. $\text{sc.}((t_{n-1}, t_n), w) = \text{tr.}(t_{n-1}, t_n) + \text{em.}(t_n, w_n)$ = transition + emission except for $t_N = \text{EOS}$, where emission = 0

17 Syntax — Syntactic Formulation

Probabilistic CFGs (N, S, Σ, R, P) —

- Non-terminals $N = \{N_1, N_2, \dots\}$
- Start non-terminal S
- Terminals $\Sigma = \{a_1, a_2, \dots\}$
- Production rules $r: N \rightarrow \alpha$, where N is non-terminal and $\alpha \in (N \cup \Sigma)^*$
- Probabilities p for each rule, locally normalized over each transition: $\sum_k p(N \rightarrow \alpha_k) = 1$ where $N \rightarrow \alpha_1, \dots, N \rightarrow \alpha_k$ are expansions of node N
- String s of length M
- Probability of a tree: $p(t) = \prod_{r \in T} p(r) = p(S \rightarrow S_1 S_2) M^{-1} \times p(S \rightarrow X) M \times p(X \rightarrow \sigma) M$
- Probability of a string: $p(s) = \sum_{t \in T(s)} p(t)$

Chomsky Normal Form (CNF)

- Production rules: $N_1 \rightarrow N_2 N_3$ are non-terminal productions, $N_3 \rightarrow \varepsilon$ are terminal productions, $S \rightarrow \varepsilon$
- Prohibits cyclic rules
- Transform CFG to CNF:
 - $\alpha \rightarrow \varepsilon$: If CFG contains $A \rightarrow BC, B \rightarrow \varepsilon$, change to $A \rightarrow BC, A \rightarrow C, \alpha \rightarrow B$: ■ If A also appears as first element in other rules: If CFG contains $A \rightarrow BC, B \rightarrow \varepsilon, D \rightarrow \varepsilon$, change to $A \rightarrow DC, D \rightarrow \varepsilon$. For long rules $A \rightarrow B_1 \dots B_k$ with $k > 2$ and mixed rules $A \rightarrow aB$ insert intermediate non-terminals

Weighted CFGs —

- General formulation of PCFGs, globally normalized
- $p(t) = \frac{\prod_{r \in T} \exp(\text{score}(r))}{\sum_{t' \in T} \prod_{r' \in T'} \exp(\text{score}(r'))}$
- Challenge: Z is infinitely large (Σ^*) • Solution:

$p(t|s) = \frac{\prod_{r \in T} \exp(\text{score}(r))}{\sum_{t' \in T(s)} \prod_{r' \in T'} \exp(\text{score}(r'))}$

- Challenge: Z is still potentially infinitely large (cycles)
- Solution: CNF. Then, $|Z|$ is the number of rooted binary trees, i.e. Catalan number C_{M-1}
- Then: $p(t|s) = \frac{\prod_{N \rightarrow N_N} \exp(\text{score}(N_i \rightarrow N_j N_k))}{\sum_{t' \in T(s)} \prod_{N \rightarrow N_N} \exp(\text{score}(N_i \rightarrow N_j N_k)) \times \prod_{N \rightarrow a} \exp(\text{score}(N_i \rightarrow a))}$

Span —

- Admissible if $X \rightarrow w[i:j]$
- For tree where nodes expand along right diagonal of tree only, only 1 admissible span per span size: $[M - \text{span size} + 1, M + 1]$

Optimization

CKY Algorithm —

- Can: Compute Z (inside SR or log-sum-exp SR), find best parse and its probability (Viterbi or arctic SR), determine if a given string is admissible by the grammar (Boolean SR)
- Requires grammar in CNF
- For Z: Chart $[i, j, X]$ is probability that non-terminal X generates the subtree resp. substring s_i, \dots, s_j
- 1) For $m=1, \dots, M$: Terminal productions: For $X \rightarrow s_m$: Chart $[m, m+1, X] \oplus = \exp(\text{score}(X \rightarrow s_m))$
- 2) For span $= 2, \dots, M$: For $i=1, \dots, M - \text{span} + 1$: $k \leftarrow i + \text{span} - 1$ For $j=i+1, \dots, k-1$: Non-terminal productions: For $X \rightarrow YZ$: Chart $[i, k, X] \oplus = \exp(\text{score}(X \rightarrow YZ)) \otimes \text{Chart}[i, j, Y] \otimes \text{Chart}[j, k, Z]$
- 3) Return Chart $[1, M+1, S]$
- Runtime: $O(M^3 |R|)$ for tree where nodes expand along right diagonal of tree only, $i=M - \text{span size} + 1, j=i+1$, then $O(M |R|)$

18 Syntax — Dependency

Formulation String w | **Spanning tree**: N nodes + 1 root node

- $N-1$ edges, of which 1 edge is fixed (root) • **Cayley's**: Undirected N^{N-2} vs. directed $N^{(N-1)}$ spanning trees • Directed, labeled spanning trees: $N^{(N-1) \times N}$

Probability of spanning tree

- $p(t|w) = \frac{\exp(\text{score}(t, w))}{\sum_{t' \in T} \exp(\text{score}(t', w))}$
- Challenge: Naively: $(N-1)^{N-2}$ spanning trees with one root
- Solution: **Edge factored assumption** $p(t|w) = \prod_{(i \rightarrow j) \in T} \exp(\text{score}(i, j, w)) \exp(\text{score}(r, w))$ where $(i \rightarrow j)$ is an edge, r is the root
- Matrix Tree Theorem** —
- Adjacency matrix**: One entry for each node i to node j , if they were connected via an edge $i \rightarrow j$, otherwise $\bar{0}$ $A_{ij} = \exp(\text{score}(i, j, w))$
- Root vector**: One entry for each node j , if it were the root node,

$\bar{0}_{p=j} = \exp(\text{score}(j, w))$

- Construct **Laplacian matrix**, accounting for constraint that there is only one root node:

$$L_{ij} = \begin{cases} \rho_j & \text{if } i=1 \\ \sum_{t'=1, t' \neq j} A_{t'j} & \text{if } i=j \\ -A_{ij} & \text{otherwise} \end{cases} \quad \text{i.e.}$$

$(-1)^{st}$ row of L contains root scores, diagonal of L (except $L_{1,1}$) contains sum within each column of A (except A $(-1)^{st}$ row) contains elements of A \times with -1

- According to matrix tree theorem: $|L| = \det(L) = Z = \text{number of directed spanning trees}$

- Runtime $O(N^3)$

Optimization

CLE algorithm:

- Greedy algorithm selects best incoming edge for each node
- Contract cycles
- Break cycle: For each enter edge, break cycle by removing edges that are also incoming at the node where the enter edge is incoming
- Re-weight: To enter edge, add weights of remaining edges on cycle
- If there are multiple edges from the root: For each edge, calculate cost of deleting it: Cost = Weight of root edge - weight of next-best incoming edge to target node
- Prelim remove edge with lowest cost, but keep target node intact
- Prelim repeat steps 5,6 as needed
- Re-run greedy algorithm
- If cycle: Undo removal, contract, then re-expand. Otherwise: Re-expand
- Runtime $O(N^2)$

19 Semantic Parsing

Lambda Calculus

Basic components:

- Logical constants
- Variables: Undetermined logical constants, free (does not occur in the scope of any abstraction that holds its name) vs. bound
- Objects x, y, z, \dots in $\lambda x.f(x)$ • Relations P, Q, R, \dots in $\lambda P.P(\dots)$ • Relation has **arity**, which is number of objects it relates, e.g. $P(x, y)$ has arity 2
- Abstraction: Let M, N be terms and x be a variable. λx is an abstraction with argument x
- $\lambda x.M$ is a function with input x , which replaces every free occurrence of x in M with whatever the function is applied to $\lambda x.MN$ is a function applied to N and body resp. **scope** MN
- Output of $\lambda x.MN$: $M[x:=N]$
- Application: $M, N \rightarrow (MN)$
- Parentheses are left-associative, e.g. $(\lambda x.\lambda y.\lambda z.z) = ((\lambda x.\lambda y.\lambda z.z).z)$
- Alpha conversion** — Renaming
- We can rename a variable in an abstraction together with all its

occurrences in the scope, which are bound to the same abstraction, if the renamed variable remains bound to the same abstraction and the remaining variables remain free

- resp. bound
- Note: Only variables are renamed, not functions
- Beta reduction** — Application We can apply one lambda term to another if the free variables in N remain free in $M[x:=N]$

LIG resp. CCGLIG

$\langle N, S, I, \Sigma, R \rangle$ —

- N : Non-terminals
- N_1, N_2, N_3, \dots
- T : Indices f, g, h, \dots
- Σ : Alphabet of terminals
- a_1, a_2, a_3, \dots
- R : Rules: $N[\sigma] \rightarrow aM[\sigma]\beta$, $N[\sigma] \rightarrow aM[f\sigma]\beta$, $N[f\sigma] \rightarrow aM[\sigma]\beta$
- Can generate languages, which cannot be generated by CFG, e.g. $\{a^n b^n \mid c^n d^n : n \geq 0\}$ with non-terminals $N = \{S, T\}$, indices $I = \{f, @, \}$, and rules $R = \{S[\sigma] \rightarrow aS[f\sigma]d, S[\sigma] \rightarrow T[\sigma@], T[\sigma f] \rightarrow bT[\sigma]c, T[\varepsilon] \rightarrow \varepsilon\}$

CCG

$\langle V_T, V_N, S, f, R \rangle$ —

- V_N : Non-terminals
- V_T : Terminals
- R : Rules: $f: V_T \rightarrow C(V_N)$
- Lexicon**: $C(V_N)$: Categories
- Categories:
 - Atomic**: Terminals
 - Complex**:
 - Built from atomic categories via operators
 - Function with pattern: $X|Y \rightarrow X$, where $X|Y$ is function, $|$ is operator, Y is argument, and X is output, and applying function to Y yields X
 - Have **arity** which is number of arguments it relates, e.g. $X|Y|Z$ has arity 2
 - Operators: Backward and forward slash
 - Note: Operators are read from outside in
 - Rules:
 - Function application**:
 - Forward** ($>$): $X: X|Y \rightarrow X$
 - Backward** ($<$): $Y: X|Y \rightarrow X$

CCGs that only have application rules, have power of CFG

- Function composition**: **Forward** (\triangleright): $(X|Y) (Y|Z) \rightarrow (X|Z)$, etc.
- Higher-order rules**: **Forward** ($>^n$): $X|Y|Y_1 \dots Y_n \dots |Y_1 \rightarrow X|_n Y_n \dots |Y_1$, etc.
- Type raising**: **Forward** (\triangleright): $X \rightarrow T/(T|X)$, etc.
- 2^{V_N} forward and 2^{V_N} backward rules, but infinitely many rule instances
- Lexicon: Associates terminals with categories • E.g. Atomic: Harry: $= NP$ • E.g. Complex: walks: $= (S \setminus NP)$

CCG parsing: $\frac{A_1 \dots A_k}{B}$ where B is a consequence of A_1, \dots, A_k

is a lexicon entry • **Inference rules**: $\frac{[X|Y, i, j] [Y\beta, j, k]}{[X\beta, i, k]} X|Y \rightarrow X\beta$ resp. $\frac{[Y\beta, i, j] [X|Y, j, k]}{[X\beta, i, k]} Y\beta \rightarrow X\beta$

- Polynomial time algorithm:
 - Arity bounded by grammar
 - constant C_g : $\text{ar}(X) \leq C_g$
 - Challenge: Categories with arity $> C_g$ can no longer be derived
 - Solution: Decompose longer derivations into smaller pieces: **Derivation contexts** c
 - $C_g \geq \max\{\ell, a+n\}$: At least as large as maximal Y (determined by largest arity a in lexicon) and β (determined by the maximum degree n of composition rules)

CCG can be paired with Lambda calculus

SK-Calculus • Alternative to Lambda calculus • **Variables**: x, y, z, \dots • **Primitive functions** resp. **combinators**:

- S : $Sxyz = (xz)(yz) = ((xz)(yz)) \circ K$
- K : Constant function: $Kxy = ((Kx)y) = x$
- I : Identity function: $Ix = x$
- SKK and I are equivalent: $S(KK)x = SKKx = K(Kx) = x$
- Parentheses are left-associative, e.g. $(Kxyz) = (((Kx)y)z)$

Lambda to SK-Calculus

- $T(x) = x$ for every variable x
- $T([E_1 E_2]) = (T[E_1] T[E_2])$ for all lambda terms E_1, E_2
- $T[\lambda x.E] = (KT[E])$ for every lambda term E where x is either bound or absent within the term
- $T[\lambda x.x] = (SKK) = I$
- $T[\lambda x.\lambda y.E] = T[\lambda x.T[\lambda y.E]]$ for every lambda term E where x is free within the term
- $T[\lambda x.(E_1 E_2)] = (ST[\lambda x.E_1] T[\lambda x.E_2])$ for all lambda terms E_1, E_2 where x is free within at least one of the two terms

20 WFSFA

FSA $\mathcal{A} = (\Sigma, Q, I, F, \delta)$ • Alphabet Σ with a, b, c, \dots • States Q , initial states I , final states F

- Transitions $\delta \subseteq Q \times (\Sigma \cup \{\varepsilon\}) \times Q$ resp. $q \xrightarrow{a \in \Sigma \cup \{\varepsilon\}} q'$
- Sequentially reads individual symbols of an input string s and transitions from state q to state q' iff $(q, a, q') \in \delta$
- If ends up in a state $q_f \in F$, automaton **accepts** the string

WFSFA

$\mathcal{A} = (\Sigma, Q, I, F, \delta, \lambda, \rho)$

- Transitions weighted with SR
- $\delta \subseteq Q \times (\Sigma \cup \{\varepsilon\}) \times \mathbb{K} \times Q$ resp. $q \xrightarrow{a \in \Sigma \cup \{\varepsilon\} / w} q' : \lambda: Q \rightarrow \mathbb{K}$ resp. $p: Q \rightarrow \mathbb{K}$ are initial resp. final weighting functions, $\bar{0}$ for non-initial resp. non-final q

WFSFA Terminology Paths: π

- $p(\pi) = q_1$ is the beginning and $q(\pi) = q_N$ the ending state of the path
- Yield** the concatenation of symbols • **Inner path weight**: $w_I(\pi) = \bigotimes_{n=1}^N w_n$
- Path weight**: $w(\pi) = \lambda(p(\pi)) \otimes w_{w_I}(\pi) \otimes \rho(q(\pi))$
- A path is **accepting** iff $w(\pi) \neq \bar{0}$
- States:
 - Accessible** iff $q \in I$ or path from I to q with $w(\pi) > 0$
 - Co-accessible** iff $q \in F$ or path from q to F with $w(\pi) > 0$

WFST

$\mathcal{T} = (\Sigma, \Omega, Q, I, F, \delta, \lambda, \rho)$ • Σ is the input alphabet, Ω is the output alphabet

- $\delta \subseteq Q \times (\Sigma \cup \{\varepsilon\}) \times (\Omega \cup \{\varepsilon\}) \times \mathbb{K} \times Q$ resp. $q \xrightarrow{a \in \Sigma \cup \{\varepsilon\}; b \in \Omega \cup \{\varepsilon\} / w} q'$
- From $\mathcal{T}_1 = (\Sigma, \Omega, Q_1, I_1, F_1, \delta_1, \lambda_1, \rho_1)$ and $\mathcal{T}_2 = (\Sigma, \Omega, Q_2, I_2, F_2, \delta_2, \lambda_2, \rho_2)$: $\mathcal{T} = (\Sigma, \Gamma, Q, I, F, \delta, \lambda, \rho)$ such that: $\mathcal{T}(x, y) = \bigoplus_{z \in \Sigma^*} \mathcal{T}_1(x, z) \otimes \mathcal{T}_2(z, y)$ where $\mathcal{T}(i, j)$ is weight assigned to mapping from input i to an output j
- Naive algorithm:
 - Initial state $(q_1^f, q_2^f) \in F_1 \times F_2$
 - Final state $(q_1^f, q_2^f) \in F_1 \times F_2$
 - For $q_1, q_2 \in Q_1 \times Q_2$: For $(q_1 \xrightarrow{a/b/w_1} q_1'), (q_2 \xrightarrow{c/d/w_2} q_2') \in E_{\mathcal{T}_1} \times E_{\mathcal{T}_2}$: If $b=c$, (q_1, q_2) and (q_1', q_2') are accessible and co-accessible (reachable from (q_1^f, q_2^f) in \mathcal{T} : Add new (q_1, q_2) , (q_1', q_2') and $(q_1, q_2) \xrightarrow{a/d/w_1 w_2} (q_1', q_2') \rightarrow \mathcal{T}$
- For $(q_1^f, q_2^f): \lambda_{\mathcal{T}} = \lambda_1(q_1^f) \otimes \lambda_2(q_2^f)$
- For $(q_1^f, q_2^f): \rho_{\mathcal{T}} = \rho_1(q_1^f) \otimes \rho_2(q_2^f)$
- Return \mathcal{T} • Challenge: Runtime $O(|Q_1| |Q_2|)$

Pathsum and Algorithms

Pathsum: $Z(\mathcal{A}) = \bigoplus_{\pi \in \Pi(\mathcal{A})} w(\pi)$

- Backward algorithm: Computes finite $Z(\mathcal{A})$ in acyclic WFSFA in $O(|\delta|)$ time:
 - For $q \in \text{Rev-Top}(\mathcal{A})$ (starts with final state):
 - If $q \in F$: $\beta(q) = \rho(q)$ Else: $\beta(q) = \bigoplus_{a/w} a/w \otimes \beta(q')$ resp. $q \xrightarrow{a/w} q'$
- Return: $\bigoplus_{q \in I} \lambda(q^i) \otimes \beta(q^i)$
- Forward: Exchange I and F and use $q' \xrightarrow{a/w} q$
- Lehmann's**: Computes infinite $Z(\mathcal{A})$ in cyclic WFSFA:
 - First WFSFA as a matrix:
 - Adjacency matrix** for all $\sigma \in \Sigma \cup \{\varepsilon\}$: $\mathbf{W}^{(a)} \in \mathbb{R}^{(|Q| \times |Q|)}$
 - From states q_n in rows, entries as weights w for $q_n \xrightarrow{a/w} q_m$ if it exists, else $\bar{0}$
 - Can be collapsed: $\mathbf{W} = \bigoplus_{a \in \Sigma \cup \{\varepsilon\}} \mathbf{W}^{(a)}$
 - Pathsum for paths of length exactly l : $\mathbf{W}^l = \mathbf{W} \otimes \mathbf{W} \otimes \dots$
 - Compute $Z(\mathcal{A})$

via Kleene closure of matrix R over closed SR: **1** $R \in \mathbb{R}^{(|Q| \times |Q|)}$ with entries as \oplus -sum of the weights for $q_i \rightarrow q_k$: $R_{ik} = \bigoplus_{\pi \in \Pi(q_i, q_k)} w_I(\pi)$ where cyclical terms are denoted with $*$

- $Z(\mathcal{A}) = \bigoplus_{i,k=1}^{|Q|} \lambda(q_i) \otimes R_{ik} \otimes \rho(q_k)$ where SR retains only pairs (i, k) where i is initial and k is final state
- Runtime $O(|Q|^3)$
- Compute R :
 - $R^{(0)} \leftarrow W$
 - For $j \leftarrow 1$ up to $|Q|$: For $i \leftarrow 1$ up to $|Q|$: For $k \leftarrow 1$ up to $|Q|$: $R_{ik}^{j-1} \oplus (R_{ij}^{j-1} \otimes (R_{jj}^{j-1})^* \otimes R_{ik}^{j-1})$ where R_{ik}^{j-1} does not traverse nodes with index $> j$
 - Return: $I \oplus R^{(|Q|)}$

Transliteration

$p(y|x)$ is probability of all paths in $\mathcal{T}_x \circ \mathcal{T}_y$ that align x to y : $= \frac{\text{pathsum of } \mathcal{T}_x \circ \mathcal{T}_y}{\text{pathsum of } \mathcal{T}_x \circ \mathcal{T}}$

- Training: For $p(y|x)$ use Lehmann's with the log-sum-exp SR
- Inference: For highest scoring y of $\mathcal{T}_x \circ \mathcal{T}$ use arctic SR
- Conditioning on x in $\mathcal{T}_x \circ \mathcal{T}$ makes WFST a WFSFA
- Lehmann's applications**
- Floyd-Warshall**:
 - Aim: Find shortest distances • Lehmann over tropical SR • We can drop $(R_{jj}^{j-1})^*$ in Lehmann
- Gauss-Jordan**:
 - Aim: For $M^{D \times D}$, find M^{-1}
 - Run Lehmann on $(\mathbf{I} - M)$
- (W)FSFA applications** **N-gram models** — As WFSFA:
 - $\Sigma = \bigcup \{(\text{BOS}), (\text{EOS})\}$
 - $Q = \bigcup_{n=0}^N \{ \{(\text{BOS})\}^{N-n} \times V^{n-1} \times \{(\text{EOS})\} \}$ represent $N-1$ histories, in total $|V|^{N-1} + 2$ states, incl. start and end
 - $I = \{(\text{BOS}, \dots, \text{BOS})\}$ (N times) is a subset of Q at $n=0$
 - $F = Q$
 - $\delta = \{y_N \dots y_1, y_0, p(y_0) \mid y_{N+1}, \dots, y_{N-1}, y_{N-1} \dots y_0 \text{ for } (y_{N-1} \dots y_1) \in \bigcup_{n=0}^{N-1} \{(\text{BOS})\}^{N-1-n} \times V^n\} \text{ for } N \geq 2, \delta(q_{ij}, y_n, l, q_{ki}) = \begin{cases} \log(p(y_n = k) \mid y_{n-1} = i, y_{n-2} = j, \dots, y_n = N) & \text{if } y_n = k \\ -\infty & \text{else} \end{cases}$
 - $\lambda = (\text{BOS}) \dots (\text{BOS}) \rightarrow \bar{1}$ (N times, full padding)
 - $\lambda(q_{ij}) = \log(p(y_n = i \mid y_{n-1} = \text{BOS}, \dots))$
 - $\rho = y_N \dots y_{N-1} \dots y_1 (\text{EOS}) \rightarrow \bar{1}$
 - $\rho(q_{ij}) = \log(p(y_n = \text{EOS} \mid y_{n-1} = i, \dots))$
 - $p(y) = \lambda(q_1) + \sum_{n=1}^L \delta(q_{y_N - N, \dots, y_{n-1}, y_n, q_{y_{N-1} - N + 1, \dots, y_1, y_0})$
 - $\rho(q_F) = |V|^{N-2} + |V|^{N-1} + 1$ start states + intermediate states + end state
- String edit FSA** — Edit distance $\leq d$ for string of length N has $(d+1)(N+1)$ states