

TTIC 31230, Fundamentals of Deep Learning

David McAllester, Winter 2019

Connectionist Temporal Classification (CTC)

Connectionist Temporal Classification (CTC)

A Successful Deep Latent Variable Model

A speech signal

$$x = x_1, \dots, x_T$$

is labeled with a phone sequence

$$y = y_1, \dots, y_N$$

with $N \ll T$ and with $y_n \in \mathcal{P}$ for a set of phonemes \mathcal{P} .

The length N of y is not determined by x and the alignment between x and y is not given.

The CTC Model

$$P_{\Phi}(y|x) = \sum_z P_{\Phi}(z|x)P_{\Phi}(y|z).$$

Input Signal: $x = x_1, \dots, x_T$

Latent Label: $z = z_1, \dots, z_T, \quad z_t \in \mathcal{P} \cup \{\perp\}$

Output: $y(z) = y_1, \dots, y_N \quad N \ll T$

$y(z)$ is the result of removing all the occurrences of \perp from z :

$$z \Rightarrow y$$

$$\perp, a_1, \perp, \perp, \perp, a_2, \perp, \perp, a_3, \perp \Rightarrow a_1, a_2, a_3$$

The CTC Model

For $z \in \mathcal{P} \cup \{\perp\}$ we have an embedding $e(z)$. The embedding is a parameter of the model.

$$h_1, \dots, h_T = \text{RNN}_\Phi(x_1, \dots, x_T)$$

$$P_\Phi(z_t | x_1, \dots, x_T) = \underset{z}{\text{softmax}} \ e(z)^\top h_t$$

z_1, \dots, z_T are **all independent** given x .

Dynamic Programming

$$x = x_1, \dots, x_T$$

$$z = z_1, \dots, z_T, \quad z_t \in \mathcal{P} \cup \{\perp\}$$

$$y = y_1, \dots, y_N, \quad y_n \in \mathcal{P}, \quad N \ll T$$

$$y(z) = (z_1, \dots, z_T) - \perp$$

$$\vec{y}_t = (z_1, \dots, z_t) - \perp$$

$$\textcolor{red}{F}[\textcolor{red}{n}, \textcolor{red}{t}] = P(\vec{y}_{\textcolor{red}{t}} = y_1, \dots, y_{\textcolor{red}{n}})$$

$$P(y) = F[N, T]$$

Dynamic Programming

$$\vec{y}_t = (z_1, \dots, z_t) - \perp$$
$$F[n, t] = P(\vec{y}_t = y_1, \dots, y_n)$$

$$F[0, 0] = 1$$

$$\text{For } n = 1, \dots, N \quad F[n, 0] = 0$$

$$\text{For } t = 1, \dots, T$$

$$F[0, t] = P(z_t = \perp)F[0, t - 1]$$

$$\text{for } n = 1, \dots, N$$

$$F[n, t] = P(z_t = \perp)F[n, t - 1] + P(z_t = y_n)F[n - 1, t - 1]$$

Back-Propagation

$$\mathcal{L} = -\ln F[N, T]$$

We can now back-propagate through this computation.

END