

Thapar Institute of Engineering & Technology

Computer Science & Engineering Department

EXERCISE (TTS) (2024-25 ODD)

1. If P_θ models a data \mathcal{D} with an intent of generating novel samples, $\mathbf{x} \sim P_\theta$. Is P_θ a posterior distribution? Comment. [2 marks]
2. If $X \equiv \{\mathbf{x}_1, \dots, \mathbf{x}_T\}$ is a temporal sequence and each time step is independently and identically distributed,
 - a) Prove that $\log P(X) = \sum_{i=1}^T \log P(x_i)$;
 - b) If X is a Bernoulli process, i.e. if each x_i is a coin toss, with hit rate p , the log likelihood of k hits is $k \log p + (T - k) \log(1 - p)$.

2 a) $X \equiv \{x_1 \dots x_T\}$

$$\begin{aligned}
 P(X) &= P(x_1 \dots x_T) \\
 &= P(x_T | x_1 \dots x_{T-1}) \cdot P(x_1 \dots x_{T-1}) \\
 &\quad \dots \text{(by Bayes Rule)} \\
 &= \prod_{i=1}^T P(x_i | x_1 \dots x_{i-1}) \quad \text{--- (1)} \\
 &\quad \dots \text{(expanding similarly)}
 \end{aligned}$$

But $x_i \perp x_j \forall i \neq j$ because iid.

hence $P(x_i | x_1 \dots x_{i-1}) = P(x_i)$ --- (2)

from (1) & (2),

$$P(X) = \prod_{i=1}^T P(x_i)$$

$$\log P(X) = \log \prod_{i=1}^T P(x_i) = \sum_{i=1}^T \log P(x_i)$$

2b) For a Bernoulli Process with hit rate ' p ',
The probability of ' k ' hits in ' T ' trials is
given as,

$$P(X) = P(k, T; p) = p^k (1-p)^{T-k}$$

$$\text{or } \log P(X) = k \log p + (T-k) \log (1-p)$$

3. If $X \equiv \{\mathbf{x}_1, \dots, \mathbf{x}_T\}$ is the temporal sequence of a Markov process, prove that $\log P(X) = \sum_{i=1}^T \log P(x_i | x_{i-1})$.

Speech sample X is a temporal sequence of intensities,

$$X \equiv \{x_1 \dots x_T\}$$

Given speech samples \mathcal{D} as evidence, estimate $P_\theta \approx \mathcal{D}$ so that $X \sim P_\theta$ is a valid speech sample.

Let $G_{\theta, \mathcal{N}} : \mathbb{X}^K \rightarrow \mathbb{X}$ represent the model.
where,

- \mathbb{X} is the field of inputs.
 - For a continuous model, $\mathbb{X} \equiv \mathbb{R}$; whereas,
 - For categorical model, $\mathbb{X} \equiv \mathbb{R}_{[0,1]}^{256}$;
 - It may also be a hybrid model,
e.g. $G_{\theta, \mathcal{N}} : \mathbb{R}^K \rightarrow \mathbb{R}_{[0,1]}^{256}$.
- Can you think how?

Let $G_{\theta, \mathcal{N}} : \mathbb{X}^K \rightarrow \mathbb{X}$ represent the model.
where,

- \mathcal{N} is a noise sampler,
 - Generally, implemented as a normal distribution; or
 - Implemented implicitly as dropouts.

Let $G_{\theta, \mathcal{N}} : \mathbb{X}^K \rightarrow \mathbb{X}$ represent the model.
where,

- θ are parameters of the model; and
- $\mathbf{x}_t = G_{\theta, \mathcal{N}}([\mathbf{x}_{t-K} \dots \mathbf{x}_{t-1}])$ models the conditional distribution $P(x_t | x_{t-K} \dots x_{t-1})$

In case of Conditional Generation,

$$\mathbf{x}_t = G_{\theta, \mathcal{N}} \left([\mathbf{x}_{t-K} \dots \mathbf{x}_{t-1}] , \mathbf{h} \right)$$

models the conditional distribution

$$P(x_t \mid x_{t-K} \dots x_{t-1}, \mathbf{h})$$

\mathbf{h} represents the global conditions, *e.g.*

- Text input;
- Speaker;
- Accent;
- and so forth.

Let

$X \equiv [\mathbf{x}_1 \dots \mathbf{x}_T] \sim G_{\theta, \mathcal{N}}$ be the output of the auto-regressive model, so that

$$\mathbf{x}_t = G_{\theta, \mathcal{N}}([\mathbf{x}_{t-K} \dots \mathbf{x}_{t-1}])$$

And,

$Y \equiv [\mathbf{y}_1 \dots \mathbf{y}_T] \sim \mathcal{D}$ be a speech sample from dataset. Recall that for each time step, the sound intensity is pre-processed so that the values are remapped, quantised, and converted to one-hot vectors, so that $\mathbf{y}_t \in \{0, 1\}^{256}; \|\mathbf{y}_t\|_1 = 1$.

The training objective is the cross entropy function, given as,

$$\underset{G}{\text{minimise}} \quad \mathbb{E}_{X \sim G, Y \sim \mathcal{D}} \left[\sum_{i,t} y_{i,t} \log x_{i,t} \right]$$