

- ◆ 沒有問題，很高興能為您說明這個 $GOP_{\text{classical}}$ 方程式。

這是在GOP (Goodness of Pronunciation) 研究早期（主要基於HMM，隱馬可夫模型）所使用的經典定義。它衡量的是一個「目標音素」與「所有可能的音素」之間的「好壞」對比。

$$GOP_{\text{classical}}(l_i) = \log \left(\frac{p(O_{t_1}^{t_2} | l_i) P(l_i)}{\sum_{q \in Q} p(O_{t_1}^{t_2} | q) P(q)} \right) / (t_2 - t_1)$$

逐項拆解說明

這個公式看起來很複雜，但我們可以把它拆成三個主要部分：

分子 (Numerator)、分母 (Denominator) 和正規化 (Normalization)。

1. 核心概念：一個機率的比例

在 \log 函數內部的那個大分數，本質上是在計算一個後驗機率 (Posterior Probability)，它遵循貝氏定理 (Bayes' Theorem)。

它的白話意思是：

「在聽到了 t_1 到 t_2 這段聲音 (O) 的前提下，說話者
真的在發 l_i (正確音素) 的機率是多少？」

2. 分子 (Numerator): $p(O_{t_1}^{t_2} | l_i)P(l_i)$

這部分代表「這段聲音是 正確音素 l_i 」的機率分數。

- l_i :
 - 代表「標準音素」 (Canonical Phoneme)，也就是這段時間內應該要發的那個正確的音。
- $O_{t_1}^{t_2}$:
 - 代表「聲學觀測值」 (Acoustic Observations)，也就是從時間 t_1 到 t_2 (這段音素的起點和終點) 所擷取到的實際語音訊號特徵 (例如 MFCCs)。
- $p(O_{t_1}^{t_2} | l_i)$:
 - 這是聲學模型 (Acoustic Model) 算出的分數。
 - 它的意思是：「假設說話者打算發 l_i 這個音，那麼他實際發出 $O_{t_1}^{t_2}$ 這段聲音的機率有多高？」
 - 如果發音很標準，這個機率就會很高。如果發音很差，機率就會很低。
- $P(l_i)$:
 - 這是語言模型 (Language Model) 或先驗機率 (Prior)。
 - 它的意思是：「 l_i 這個音素本身在整個語言中出現的機率有多高？」

3. 分母 (Denominator): $\sum_{q \in Q} p(O_{t_1}^{t_2} | q)P(q)$

這部分代表「這段聲音是任何可能音素」的機率分數的總和。

- Q :
 - 代表「所有可能的音素」的集合 (例如 /a/, /b/, /p/, /t/, /k/... 以及 l_i 自己)。
- q :

- 代表 Q 集合中的任何一個音素（可能是正確的 l_i ，也可能是錯誤的其他音）。
- $\sum_{q \in Q} \dots$
 - 這是一個加總 (Summation)。它會把這段聲音 $O_{t_1}^{t_2}$ 分別套用到所有可能的音素 q 上去計算分數（就像分子那樣），然後把這些分數全部加起來。
 - 這代表了這段聲音 $O_{t_1}^{t_2}$ 出現的「總體機率」，而不考慮它到底應該是哪個音。

4. 分數的意義： $\frac{\text{分子}}{\text{分母}}$

這就是貝氏定理的應用：

$$P(l_i | O_{t_1}^{t_2}) = \frac{p(O_{t_1}^{t_2} | l_i) P(l_i)}{\sum_{q \in Q} p(O_{t_1}^{t_2} | q) P(q)}$$

- 如果發音很標準，分子 (l_i 的分數) 會很高，而分母（所有音素的總分）主要由 l_i 貢獻，所以這個比例會接近 1。
- 如果發音很差（例如，把 /l/ 發成了 /n/），那麼 l_i (/l/) 的分數（分子）會很低，而 $q = /n/$ 的分數會很高，導致分母遠大於分子，這個比例就會接近 0。

5. 外圍的計算： $\log(\dots) / (t_2 - t_1)$

這是在做最後的「分數處理」。

- $\log(\dots)$:
 - 取對數 (Logarithm)。這是因為機率值通常都非常小（0 到 1 之間），連乘起來會造成數值下溢 (underflow)。取 log 可以將乘法變加法，將機率轉為「對數機率」 (Log-Likelihood)，在計算上更穩定。
- $\$/ (t_2 - t_1)\$$:
 - 時長正規化 (Duration Normalization)。

- $t_2 - t_1$ 就是這個音素的持續時間（例如 $t_1 = 10$ 帧， $t_2 = 25$ 帧，持續時間就是 15 帧）。
 - 為什麼要除以時長？因為一個發音很長（例如元音 /a:/）的音素，其 log 機率值會比一個發音很短（例如塞音 /t/）的音素低很多（因為是更多 < 0 的值相加）。如果不除以時長，長音素的分數天生就會比較低，這不公平。
 - 除以時長後，GOP 分數就變成了「平均每帧 (per-frame) 的對數機率」，這樣長音素和短音素的分數就可以互相比較了。
-

總結

$GOP_{\text{classical}}$ 的計算步驟：

1. 強制對齊 (Forced Alignment)：首先，系統需要透過語音辨識模型（通常是 HMM）找出 l_i 這個標準音素對應到 t_1 到 t_2 這段實際的聲音。
2. 計算後驗機率：計算「這段聲音是 l_i 的機率」佔「這段聲音是所有可能音素的總機率」的比例。
3. 取 Log 並正規化：將這個比例取 log，然後除以它的持續時間，得到最終的 GOP 分數。

分數意義：一個高的 $GOP_{\text{classical}}$ 分數（通常是接近 0 的負數），意味著與所有其他可能的「競爭音素」相比，這段聲音被模型判斷為「就是那個正確音素 l_i 」的機率非常高。