

# 基於 Logit 的發音評估多面向統計分析法

超越基本統計量:以高階特徵捕捉細微發音動態

Chieh-Ren Liao, Berlin Chen

# 研究背景:發音評估的機率飽和問題

## 電腦輔助語言學習的核心挑戰

發音好壞度(GOP)評估是CAPT系統的關鍵技術。傳統GOP依賴深度神經網路聲學模型輸出的後驗機率,這些機率透過對原始logits進行softmax歸一化得到。

## 傳統方法的固有缺陷

Softmax函數存在「過度自信」或「機率飽和」問題,將機率分佈推向極端,壓縮了不同音素間的區分度,使細微發音偏差難以偵測。



## 關鍵問題

機率飽和導致細微發音差異被掩蓋



# 研究動機:從機率轉向Logit

## Logit的優勢

直接使用未經處理的Logits計算GOP,避免機率飽和,保留更豐富的鑑別資訊。

## 現有方法侷限

基線指標僅依賴單點統計量或一階動差,忽略Logit序列的複雜動態分佈與時序特性。

## 本研究目標

超越均值與變異數,將Logit序列視為完整統計分佈和時間序列進行建模。

# 基線指標的侷限性分析

基線指標	描述特性	侷限性
GOP_MaxLogit	峰值信心水準	僅依賴單點統計量
GOP_margin	區分程度	僅依賴一階動差(均值)
GOP_variance	信心穩定性	僅依賴一階動差(變異數)

核心問題:這些指標忽略了Logit序列在音素持續時間內的複雜動態分佈與時序特性。

# 五種高階指標體系

01

## 分佈形狀特徵

偏度與峰度:描述Logit分佈的完整形狀

02

## 資訊理論特徵

夏農熵與KL散度:量化整體混淆程度

03

## 分佈擬合特徵

高斯混合模型:捕捉多階段動態

04

## 時序穩定性特徵

自相關分析:評估時間連貫性

05

## 峰值穩健性特徵

Top-k平均:克服單點雜訊敏感性

# 類別1:分佈形狀特徵

## 偏度 (Skewness, G1)

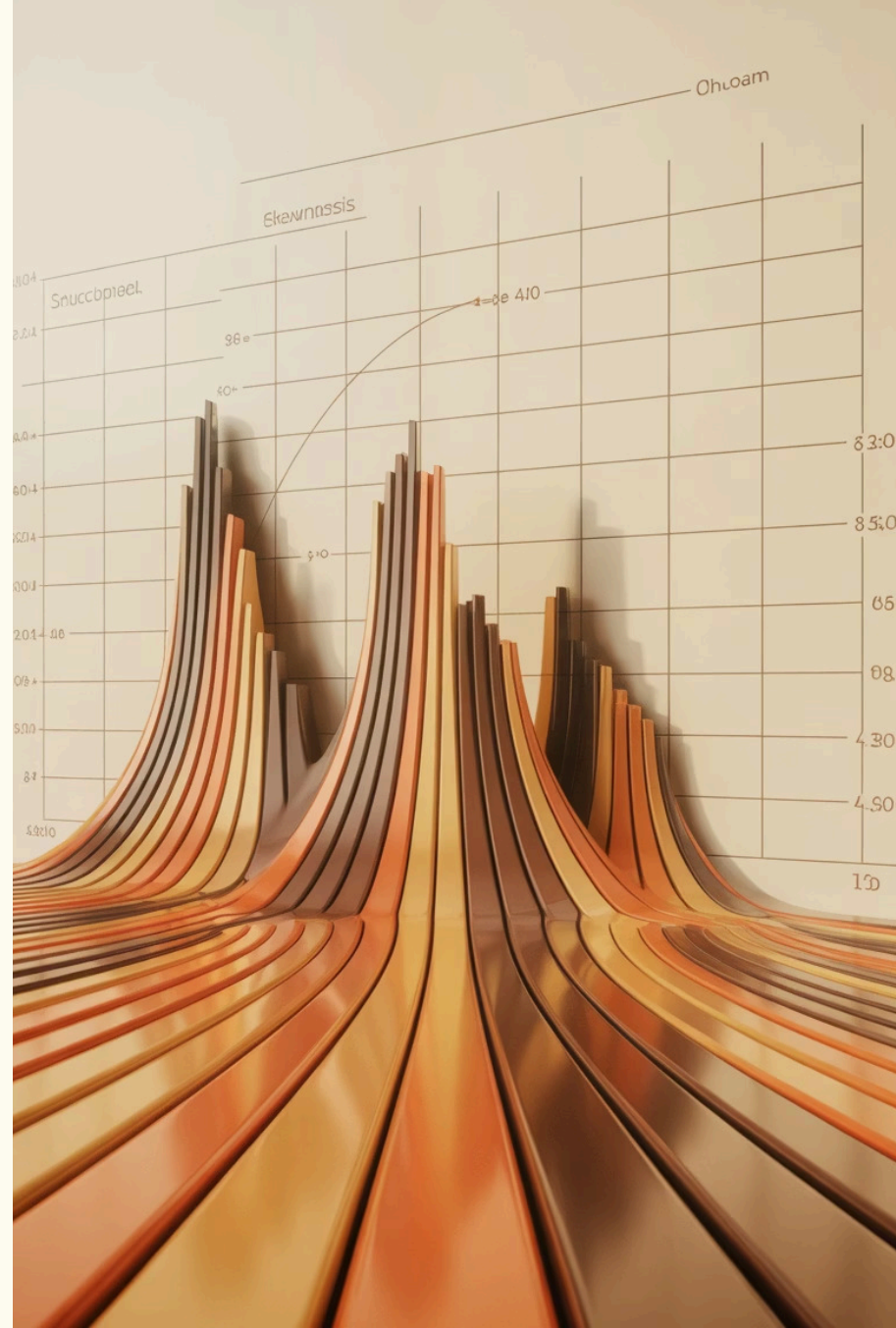
第三階標準化動差,衡量分佈的不對稱性。正偏度可能表示模型信心是逐漸建立後迅速下降的過程。

異常的偏斜可能暗示不自然的發音模式。

## 峰度 (Kurtosis, G2)

第四階標準化動差,衡量分佈的峰銳度與尾部厚度。高峰度表示模型信心高度集中於某個值。

有助於識別發音過程中信心的集中或分散程度。



## 類別2與3:資訊理論與分佈擬合

### 平均夏農熵

計算每一幀後驗機率分佈的夏農熵並取平均。高平均熵意味著模型機率分散在多個候選音素上,是發音含糊或錯誤的強烈信號。

### 平均KL散度

衡量實際後驗機率分佈與理想one-hot向量間的距離。較大的KL散度意味著模型輸出與理想狀態相去甚遠。

### 高斯混合模型

假設Logit序列由音素的多個潛在狀態混合而成。提取各分量的均值、變異數和權重,精細描述發音過程中模型信心的多階段動態。

# 類別4與5:時序穩定性與峰值穩健性

## 自相關分析

Logit序列在延遲為1時的自相關係數,衡量序列隨時間變化的平滑程度與穩定性。

- 高正相關:Logit序列平滑穩定,對應清晰發音
- 接近零或負值:存在劇烈不規則波動,暗示發音不穩定

## Top-k平均值

選取Logit序列中最大的k個值並計算平均。提供更穩健的峰值信心估計,有效平滑單一離群值影響。

克服GOP\_MaxLogit的不穩健性。



# 實驗設計與評估標準



## 資料集

SpeechOcean762 L2英語語音資料庫



## 任務

發音錯誤檢測的分類任務



## 主要指標

馬修斯相關係數(MCC),最適合類別不平衡數據

階段	數據規模	目的
初步實驗	2500筆	進行指標性能的初步探勘
完整實驗	5000筆	驗證指標的泛化能力與穩定性

# 實驗結果:數據規模決定指標選擇

1

## 數據稀疏時(2500筆)

峰度和偏度等描述分佈形狀的高階動差指標表現最佳。形狀是比絕對數值更穩健的錯誤信號。

2

## 數據充足時(5000筆)

發生指標反轉效應。mean\_logit\_margin、evt\_k3和kl\_to\_onehot構成第一梯隊,區分度與峰值強度成為最強分類特徵。

### 📌 核心發現

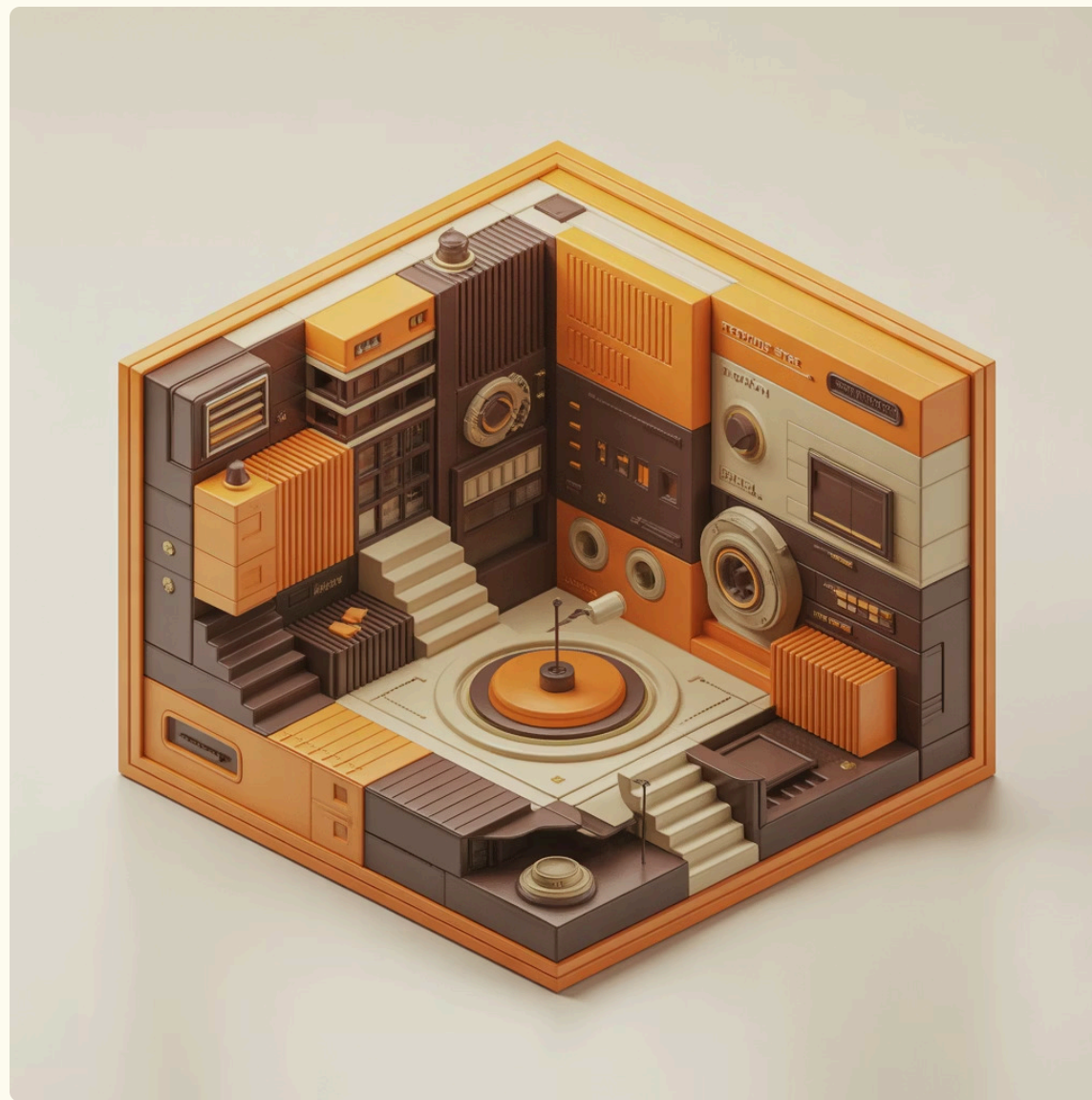
Logit-based GOP指標的有效性與最佳選擇,高度依賴於實驗數據的規模。頂尖指標都從不同側面有效捕捉了「模型判斷的確定性」這一核心概念。

# Q1: 為什麼選擇這五類高階指標?

## 設計理念

這五類指標從不同維度全面描述Logit序列的統計特性:

- 形狀特徵捕捉分佈異常
- 資訊理論量化不確定性
- GMM建模多階段動態
- 自相關評估時序連貫性
- 極值理論提供穩健估計



這些指標互補而非重複,共同構建了一個多維度、多層次的發音評估框架。

## Q2: 如何解釋指標反轉效應?

### 數據稀疏階段

信號充滿雜訊時,分佈形狀提供穩健的異常檢測能力,不易受個別數值波動影響。

1

2

3

### 數據充足階段

區分度與峰值強度等直接衡量模型信心的指標成為最有效的分類特徵。

### 轉折點

隨著數據增加,模型判斷趨於穩定,直接的信心度量開始展現優勢。

這種反轉揭示了統計穩健性與判別能力之間的權衡關係。

## Q3: 這些指標的計算複雜度如何?

$O(n)$

基本統計量

偏度、峰度、熵等可在單次遍歷中計算

$O(nk)$

GMM擬合

需要迭代優化,但k值通常很小(2-5)

$O(n)$

自相關

線性時間複雜度,計算高效

所有指標都可在實時系統中高效計算,不會成為CAPT系統的性能瓶頸。相較於深度神經網路的推理時間,這些統計計算的開銷可忽略不計。

# Q4: 如何應用於實際CAPT系統?



## 語音輸入

學習者發音被聲學模型處理,產生逐幀Logit序列

$$\frac{f}{dx}$$

## 特徵提取

根據數據規模選擇適當指標組合進行計算




## 融合判斷

整合多個指標進行綜合評估,提供細緻的發音反饋



## 個性化回饋

根據不同指標的異常模式,提供針對性的改進建議



## Q5: 未來研究方向與展望

### 特徵融合建模

將mean\_logit\_margin、evt\_k3、kl\_to\_onehot等頂尖指標融合,預期多維度綜合判斷將顯著超越單一特徵。

### 數據規模深入探討

探討形狀指標與區分度指標發生性能交叉的數據量級,提供自適應特徵選擇策略。

### 泛化性驗證

應用於不同聲學模型架構(如Whisper)或不同母語背景學習者,驗證結論的普適性。