

- Maximum Likelihood Estimation (MLE)

⇒ choose model parameter λ to make likelihood
 $P_\lambda(Y|W)$ as large as possible
 ($Y \rightarrow$ acoustic signal, $W \rightarrow$ word sequence)

$$F_{ML}(\lambda) = \sum_u P_\lambda(Y_u|W_u)$$

- Maximum Mutual Information Estimation (MMIE)

⇒ choose model parameter λ to make mutual information
 $I_\lambda(W, Y)$ as large as possible.

$$I_\lambda(W, Y) = \sum_{w, y} P(w, y) \log \left(\frac{P_\lambda(w, y)}{P_\lambda(w) P_\lambda(y)} \right)$$

$$\begin{aligned} \Rightarrow F_{MMIE}(\lambda) &= \sum_u \log \left(\frac{P_\lambda(w, y)}{P_\lambda(w) P_\lambda(y)} \right) \\ &= \sum_u \log \left(\left(\frac{P_\lambda(w, y)}{P_\lambda(w)} \right) \cdot \frac{1}{P_\lambda(y)} \right) = \left[\sum_u \left(\log P_\lambda(y|w) - \log P_\lambda(y) \right) \right] \end{aligned}$$

MLE

⇒ The difference btw MLE, MMIE = $-\log P_\lambda(y)$

(ignore outermost Σ)

$$\begin{aligned} \frac{\partial F_{MMIE}(\lambda)}{\partial \lambda} &\stackrel{\text{(ignore outermost } \Sigma)}{=} \frac{\partial \left\{ \log P_\lambda(y|w) - \log P_\lambda(y) \right\}}{\partial \lambda} = \frac{\partial \log P_\lambda(y|w)}{\partial \lambda} - \frac{\partial \log P_\lambda(y)}{\partial \lambda} \\ &= \frac{\left\{ \frac{\partial P_\lambda(y|w)}{\partial \lambda} \right\}}{P_\lambda(y|w)} - \frac{\left\{ \frac{\partial P_\lambda(y)}{\partial \lambda} \right\}}{P_\lambda(y)} = \frac{\left\{ \frac{\partial P_\lambda(y|w)}{\partial \lambda} \right\}}{P_\lambda(y|w)} - \frac{\left\{ \frac{\partial (\sum_w P_\lambda(y|w) P_\lambda(w))}{\partial \lambda} \right\}}{P_\lambda(y)} \\ &= \frac{\left\{ \frac{\partial P_\lambda(y|w)}{\partial \lambda} \right\}}{P_\lambda(y|w)} - \frac{\left\{ \frac{\partial P_\lambda(y|w)}{\partial \lambda} \right\} P_\lambda(w) + \sum_{\tilde{w} \neq w} \left\{ \frac{\partial P_\lambda(y|\tilde{w})}{\partial \lambda} \right\} P_\lambda(\tilde{w})}{P_\lambda(y)} \\ &= \left\{ \frac{\partial P_\lambda(y|w)}{\partial \lambda} \right\} \left\{ \frac{1}{P_\lambda(y|w)} - \frac{P_\lambda(w)}{P_\lambda(y)} \right\} - \sum_{\tilde{w} \neq w} \left\{ \frac{\partial P_\lambda(y|\tilde{w})}{\partial \lambda} \right\} \frac{P_\lambda(\tilde{w})}{P_\lambda(y)} \end{aligned}$$

derivative of MLE

subtract direction $\frac{\partial P_\lambda(y|\tilde{w})}{\partial \lambda}$
 for incorrect word sequence
 $\tilde{w} \neq w$

$$\rightarrow F_{\text{mmie}}(\lambda) = \sum_u \log \left(\frac{P(w, Y)}{P(w)P(Y)} \right) = \sum_u \left\{ \log \frac{P_\lambda(w, Y)}{P_\lambda(Y)} - \log P_\lambda(w) \right\}$$

$$= \sum_u \left\{ \log \frac{P_\lambda(Y|w)^k P(w)}{\sum_w P_\lambda(Y|w)^k P(w)} - \log P_\lambda(w) \right\}$$

(k = scaling fudge factor)

- Boosted MMI

\Rightarrow boost the likelihood of the sentences have \uparrow errors, generating \uparrow confusable data.

$$\rightarrow F_{\text{mmie}}(\lambda) = \sum_u \left\{ \log \frac{P_\lambda(Y|w)^k P(w)}{\sum_w P_\lambda(Y|w)^k P(w)} e^{-bA(w, w_u)} \right\}$$

(b = boosting factor, $A(w, w_u)$ = accuracy of sentence w given reference w_u)

- Minimum Phone Error (MPE), ^{State-level} minimum Bayes Risk (s-MBR)

\Rightarrow minimize the error corresponding to different granularity of labels (phone / state)

$$\rightarrow F_{\text{mpe}}(\lambda) = \sum_u \left\{ \log \frac{\sum_w P_\lambda(Y|w)^k P(w) A(w, w_u)}{\sum_w P_\lambda(Y|w)^k P(w)} \right\}$$

($A(w, w_u)$ = phone accuracy of w given reference w_u)

$$\rightarrow F_{\text{smbR}}(\lambda) \Rightarrow A(w, w_u) = \text{state-level accuracy}$$