

Las técnicas de N-gramas buscan estimar la probabilidad de secuencias de palabras basadas en datos de entrenamiento.

=> Absolute discounting: Reduce una cantidad fija,  $D$ , de la cuenta de los N-gramas no observados. Esto resuelve el problema de asignar probabilidad a secuencias que no aparecen en el corpus.

Su fórmula (expresión matemática involucrada):

$$P_{AD}(w_i | w_{i-n+1}^{i-1}) = \frac{\max(C(w_{i-n+1}^{i-1}) - D, 0)}{C(w_{i-n+1}^{i-1})} + \lambda P_{\text{backoff}}(w_i | w_{i-n+1}^{i-1})$$

donde  $C$  es la frecuencia observada y  $\lambda$  ajusta la redistribución de la probabilidad.

=> Smoothing Kneser - Ney: Extiende el descuento absoluto incorporando información sobre el número de contextos únicos en los que una palabra ha aparecido, resolviendo mejor la asignación de probabilidad a palabras raras o no vistas.

Su fórmula (expresión matemática involucrada)

$$P_{KN}(w_i | w_{i-n+1}^{i-1}) = \max(C(w_{i-n+1}^{i-1}) - D, 0) / C(w_{i-n+1}^{i-1}) + \lambda P_{KN}(w_i | w_{i-n+1}^{i-1})$$

Ambas técnicas resuelven la problemática de asignar probabilidades a secuencias no vistas en el corpus, evitando probabilidades nulas.

Un ejemplo de Absolute discounting:

Suponiendo que tenemos los siguientes trigramas y sus conteos en un corpus:

- "el gato come": 4
- "gato come pescado": 1
- "come pescado fresco": 0 (no visto)

Con absolute discounting ( $D=0.75$ ), la probabilidad del trigramma "gato come pescado" sería:

$$P_{AD}(\text{pescado} | \text{gato come}) = \max(1 - 0.75, 0) + \lambda P_{\text{backoff}}(\text{pescado} | \text{come})$$

Redistribuimos parte de la probabilidad de los trigramas vistos a los no vistos, como "come pescado fresco".

Y al hacerlo para Kneser, Ney, además del descuento, se consideraría cuántos contextos distintos preceden a "pescado", mejorando la probabilidad para palabras raras en contextos conocidos.

Este enfoque evita frases no vistas como "come pescado fresco" tengan probabilidad cero.