

→ Porter's Algorithm:

1.1:

↳ sses → ss (caresses → caress)
 ↳ ies → i (ponies → poni)
 ↳ ss → ss (caress → caress)
 ↳ s → ∅ (cats → cat)

1.2:

→ (*v*)ing → ∅ (walking → walk, king → king)
 → (*v*)ed → ∅ (played → play)

2:

ational → ate (relational → relate)
 izer → ize (digitizer → digitize)
 ator → ate (operator → operate)

3:

al → ∅ (revival → revive)
 able → ∅ (adjustable → adjust)
 ate → ∅ (activate → activate)

→ Language Models: Probabilistic Model

.

V U
→ N-gram LM: (Content Sensitive Spelling Correction)

→ the office is about fifteen minvets from my house.
minutes = minvets 1700-1800

→ $p(\text{fifteen minutes}) > p(\text{fifteen minvets})$

→ Sentence Plausibility:

→ karma → jeffe

→ $p(\text{high winds}) > p(\text{large winds})$

→ Probabilistic Language Models:

$$P(w) = P(w_1, w_2, w_3, \dots, w_n) \sim \text{Local}$$

$$P(w_i | w_1, w_2, w_3) \sim \text{Related}$$

→ Joint Probability:

$$P(\text{about, fifteen, minutes, from})$$

P (about, fifteen, minutes, from)

→ Basic : Rely on the chain rule of probability.

→ Conditional Prob : $P(B|A) = \frac{P(A, B)}{P(A)}$

Where $P(A, B) = P(A) \cdot P(B|A)$

→ Generalized Rule :

$$P(x_1, x_2, \dots, x_n) = P(x_1) \cdot P(x_2|x_1) \cdot P(x_3|x_1, x_2) \dots \dots P(x_n|x_1, \dots, x_{n-1})$$

$$\therefore P(w_1, w_2, \dots, w_n) = \prod_i P(w_i | w_1, w_2, \dots, w_{i-1})$$

→ $P(\text{about fifteen minutes from}) =$
 $P(\text{about}) \times P(\text{fifteen}|\text{about}) \times \dots$

→ Markov Assumption :

$$P(\text{office} | \text{about fifteen min from}) \approx P(\text{office} | \text{from})$$

→ Markov Assumption : Only k previous words (k^{th} order)

$$P(w_1, w_2, \dots, w_n) \approx \prod_i P(w_i | w_{i-k} \dots w_{i-1})$$

→ Maximum Likelihood Expectation:

$$P(w_i | w_{i-1}) = \frac{\text{count}(w_{i-1}, w_i)}{\text{count}(w_{i-1})}$$

$$\hat{P}(w_i | w_{i-1}) = \frac{c(w_{i-1}, w_i)}{\underline{\underline{c(w_{i-1})}}}$$