

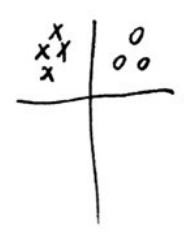
# Natural Language Processing (~1990)

- 1) Spam filtering
- 2) Query Completion
- 3) Document ranking
- 4) Topic extraction

Represent a text

1) Frequency counting

dimension = # words in your text



2) N-grams  
n-gram is just a length n sequence from the text

- predict next word
- decrypting substitution Cyphers
- filling in missing data

Smoothing: "The power of ~~the~~ absolute dis counting"

# Hidden Markov Models

## Markov Chain:

States:  ~~$z_1, \dots, z_n$~~   $s_1, \dots, s_n$

Transition  $P = [P_{ij}]_{i,j=1}^n$

probability  $P_{ij}$  is the probability  $s_i \rightarrow s_j$

$z_t$  state of time  $t$

$$P(z_T = s | z_{T-1}, \dots, z_1) = P(z_T = s | z_{T-1})$$

$$\pi^t = [\frac{1}{3}, \frac{1}{3}, \frac{1}{3}]$$

$\pi \sim$  states, then  $\pi_t = P^t \pi$

$\pi^* = \lim_{t \rightarrow \infty} \pi_t$  exists and is unique independent of  $\pi_0$

$\pi^*$  is limiting / equilibrium / stationary distribution

## Hidden Markov Model:

States:  $s_1, \dots, s_n$

each state has "emission" distribution  $P_i$

Have some variable  $x_1, \dots, x_k$

Goal: Estimate  $P; P_i$

EM algorithm

# Markov Chain Monte Carlo (MCMC)

Idea: Model the distribution we want to estimate on the limiting distribution of MC. simulate the Markov Chain.

---

## Substitution ~~Cyphers~~ Cyphers

$\Sigma$  Alphabet

$$\sigma \in S_{|\Sigma|}$$

Goal: find  $\sigma^{-1}$

1)  $M =$  bigram conditional probability matrix  $= [P_{ij}]$  for English

$$w \in S_{|\Sigma|}$$

$$p(w) = \sum_{x=1}^{|\Sigma|-1} \log \frac{1}{|w(x)|} (P_{w(x)w(x+1)})$$

1) Pick ~~for~~ a random  $w \in S_{|\Sigma|}$

2) loop until for  $N \sim (2000)$

3) Pick  $\tau \in S_{|\Sigma|}$  a transposition

4) if  $p(w.\tau) > p(w)$

$$w = w\sigma\tau$$

else with probability  $\frac{p(w\sigma\tau)}{p(w)}$

$$w = w\sigma\tau$$

# Gibbs Sampling

$D_1, \dots, D_n$  documents

$K$  topics

topic is a set of words and a distribution over the words

$\theta_i$  is a distribution on  $T_1, \dots, T_k$

$(\theta_i)_j$  represents the proportion of  $D_i$  is  $T_j$

1) Assign each word in each document a topic using  $\theta_i$ 's

2) For each document  $D$ , word  $w \in D$

$$P(T|D) = \frac{\# \text{ words in } D \text{ assigned to } T}{\# \text{ words in } D}$$
$$P(w|T) = \frac{\# w \text{ was assigned to } T}{\# \text{ words in } T}$$

$P(T|D) \cdot P(w|T)$