$\rightarrow$ Supervised Learning for WSD.

$\quad\hookrightarrow$ Parametric Model ( Naive Bayes )

$$\hat{s} = \underset{s \in S}{\arg\max} \; P(s \mid f)$$

$$= \underset{s \in S}{\arg\max} \; \frac{P(f \mid s) \cdot P(s)}{P(f)}$$

$$= \underset{s \in S}{\arg\max} \; P(s) \cdot P(f \mid s)$$

$$= \underset{s \in S}{\arg\max} \; P(s) \prod_{i=1}^{n} P(f_i \mid s)$$

$$\left( \; [f_1 \cdots f_n] \right.$$
$$\hookrightarrow context$$

$\rightarrow$ Collocation Vector ( Set of words around it)

$\rightarrow$ Setting parameters of Naive Bayes using MLE from training data:

$$P(s_i) = \frac{\text{count}(s_i, v_j)}{\text{count}(v_j)}$$

$$P(f_j | s_i) = \frac{\text{count}(f_j, s_i)}{\text{count}(s_i)}$$

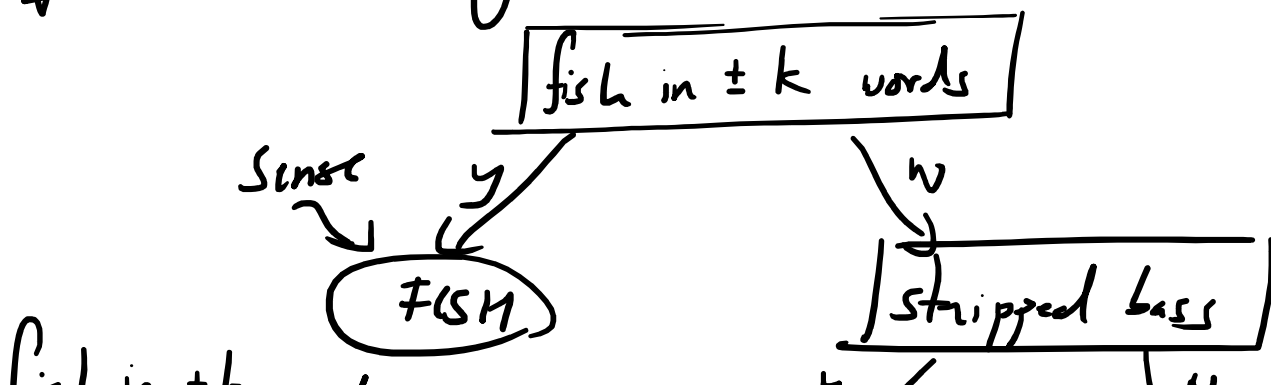→ Non-Parametric Method (Decision List)

   ↳ One Sense per Collocation.

   ↳ Log-Likelihood Ratio:

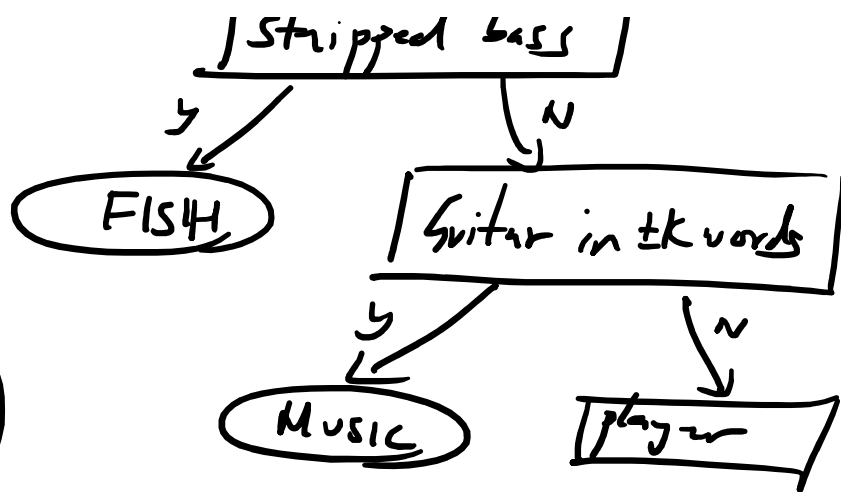$$\log \left\{ \frac{P(\text{Sense-A} | \text{Collocation}_i)}{P(\text{Sense-B} | \text{Collocation}_i)} \right\}$$

→ Higher log-likelihood = More Predictive Evidence.
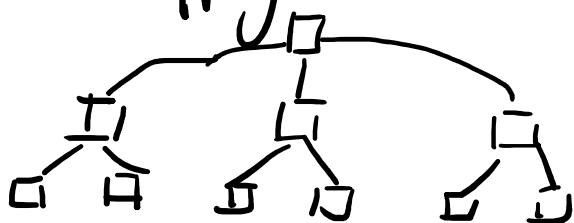
Eg: Discriminating b/w bass (fish / music)



Sense

fish in ± k words

FISH

stripped bass

+(SH)

fish in ±k words
stripped bass
Guitar in ±k words
player

$O(n \log n)$



Stripped bass → [y] FISH
Stripped bass → [N] Guitar in ±k words
Guitar in ±k words → [y] Music
Guitar in ±k words → [N] player

→ Minimally Supervised WSD ~ Janowsky.

   ↳ Bootstrapping or Co-training.

Random Forest

   ↳ Start with small seed, decision list.

   ↳ Use decision list to label corpus.

   ↳ Retain confident labels as annotated
       data to learn new decision list.

⇒ Heuristics (Derived from observations)
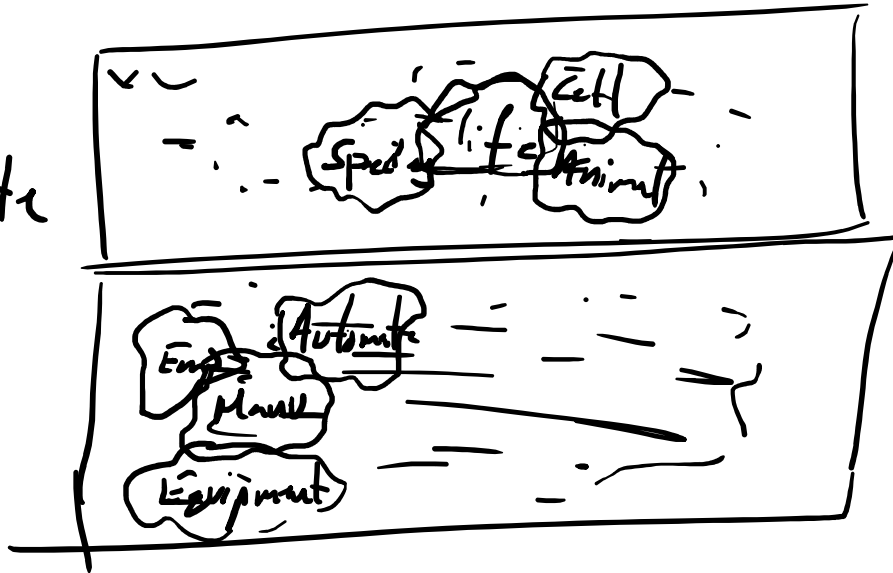
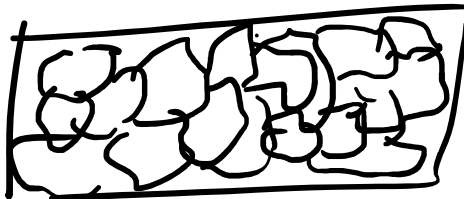   ↳ One sense per discourse

   ↳ One sense per collocation.

Eg:

Eg:

-> Disambiguate plant ( living thing) vs plant
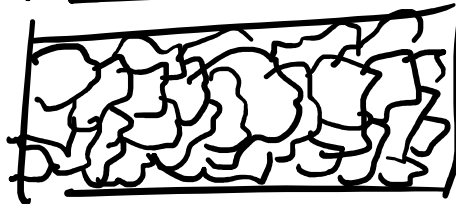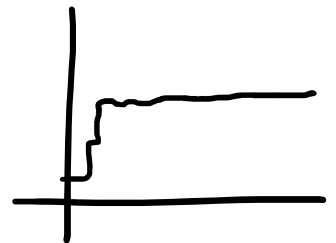
(industrial)

Life

Intermediate
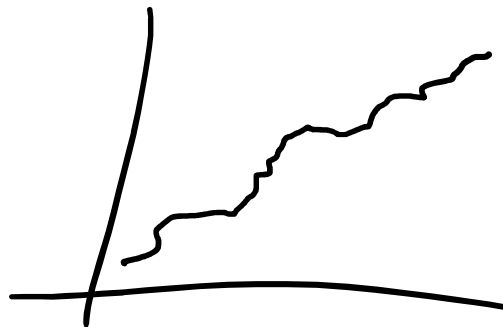State



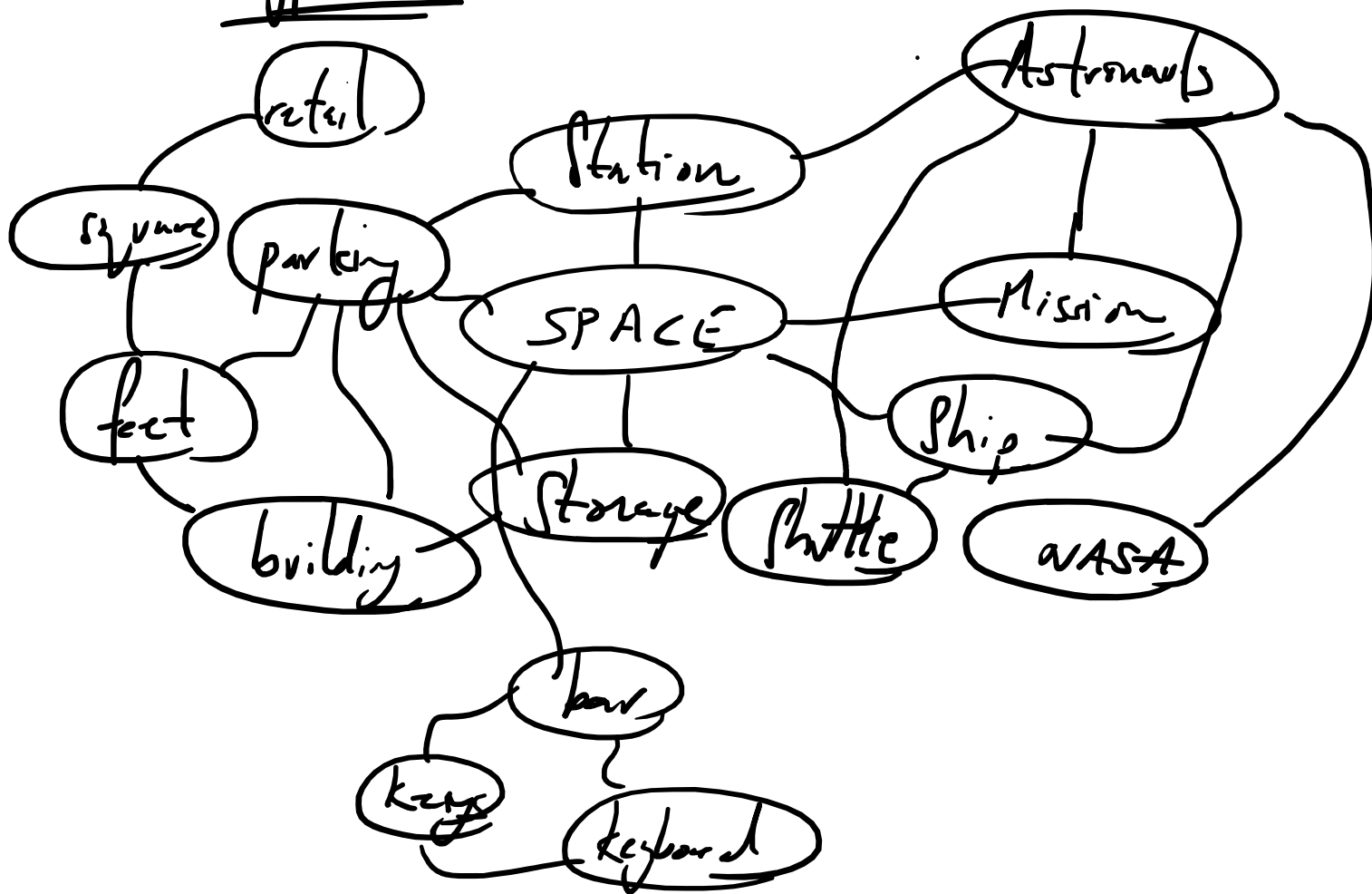Manufacturing

Life



Manufacturing

Termination :

-> Stop when ?

Advantages :

**Printing .**

-> Unsupervised :

    -s Hypertex : Word Sense Induction



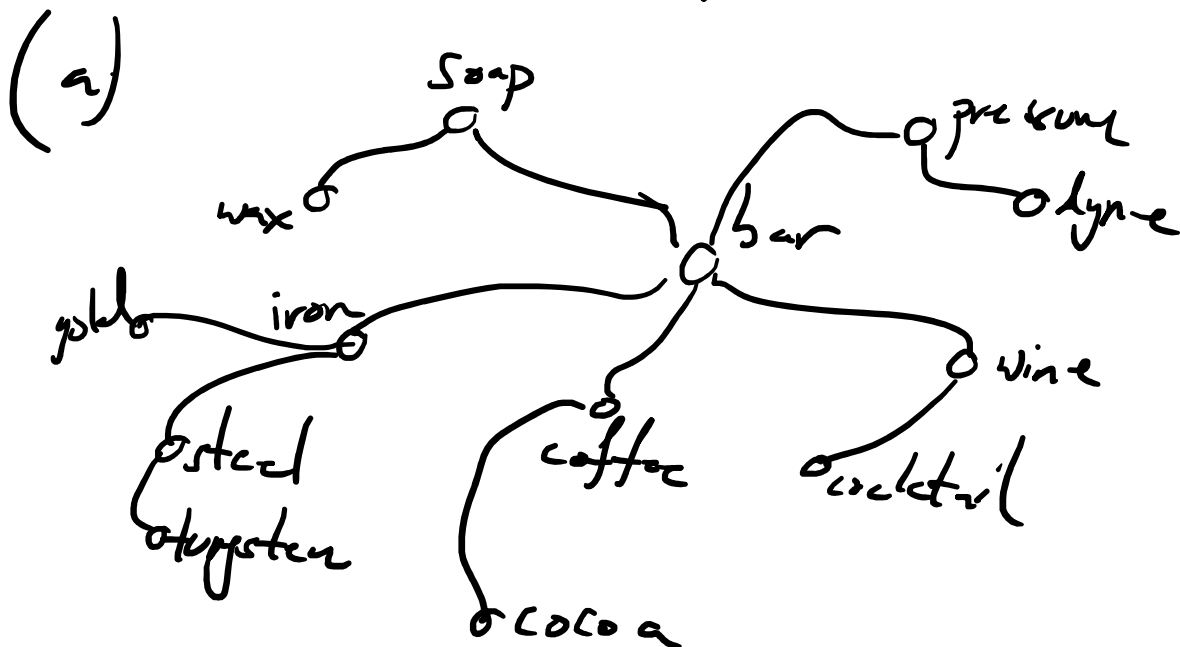-> Detecting Root Hubs :

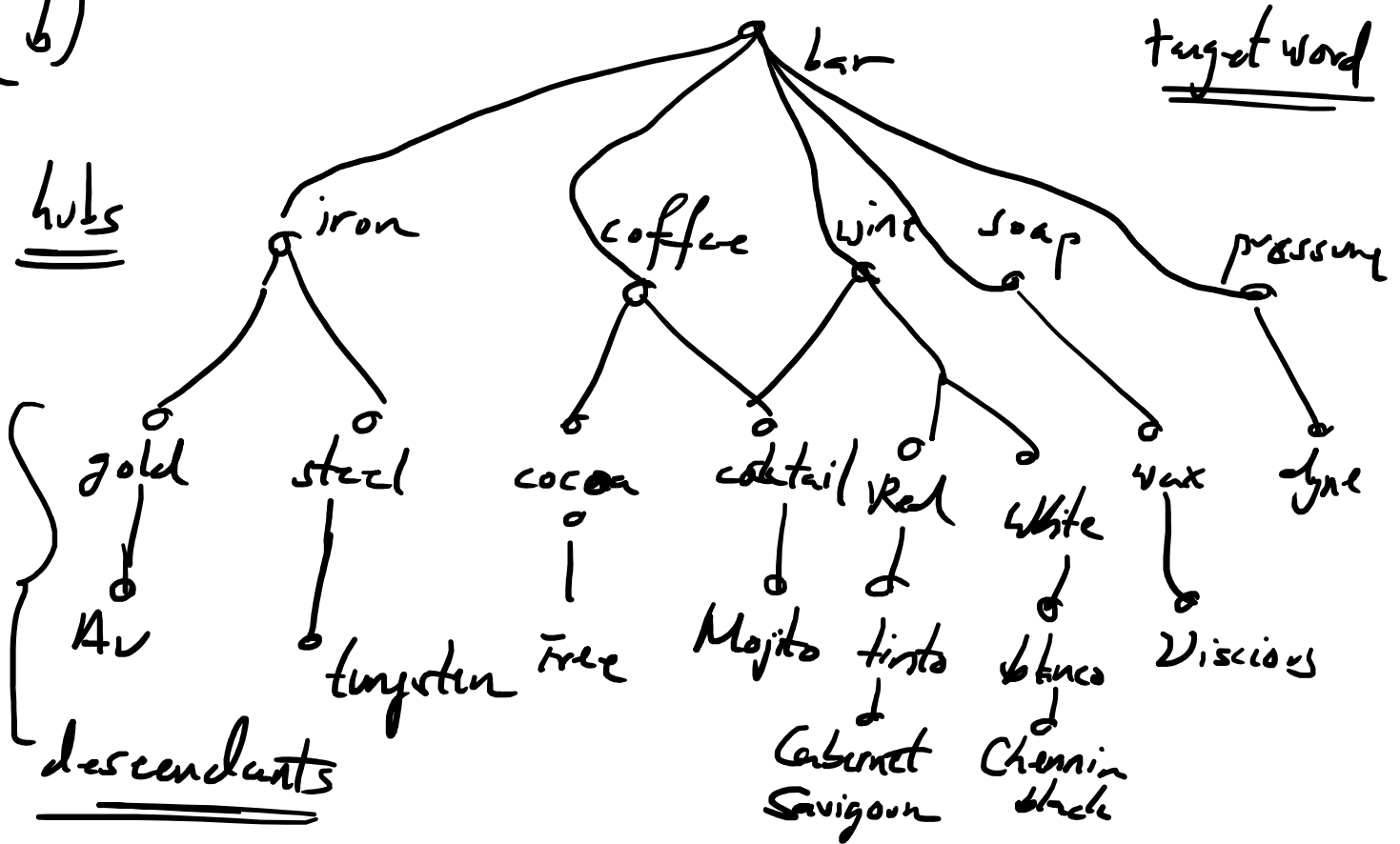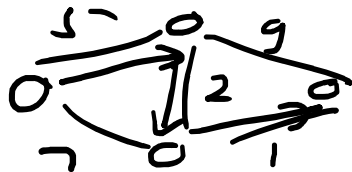    -s Co-occurrence graph :

      Ls Arrange nodes in decreasing order

↳ Arrange nodes in decreasing order of degree.

↳ Select the nodes from graph which has the highest degree. This node will be the hub of the first high density component.

↳ Delete the hub & all its neighbors from Graph.

↳ Repeat steps to detect hubs of other high degree components.

(a)

# (b)

hubs

bar

iron          coffee      wine    soap        pressure

gold    steel    cocoa    cocktail  Red      White      wax      lyne

Au      tungsten  free    Mojito  tinto  blanco   Viscious

descendants

Cabernet      Chennin
Savignon      blanc

→ Delineating Components :



Computing distance b/w 2 nodes
$\omega_i$ & $\omega_j$

$$w_{ij} = 1 - \max \{ P(w_i | w_j), P(w_j | w_i) \}$$

where $P(w_i | w_j) = freq_{ij} / freq_j$

→ <u>Disambiguation</u> : <u>Minimum Spanning Tree</u>

→ Let $w = (w_1, \ldots, w_i, \ldots, w_n)$ be a context
  in which $w_i$ is an instance of over
  target word.

→ Let $w_i$ has $k$ hubs in its MST.

→ A score vector $\underline{s}$ is associated with each
  $w_j \in W (j \neq i)$, such that $S_k$ represents
  the contribution of the $k^{th}$ hub as :

  $S_i = \dfrac{1}{\phantom{(\ldots)}}$

$$S_k = \frac{1}{1 + d(h_k, w_j)}$$

$\quad$ if $h_k$ is an ancestor of $w_j$

$\quad = 0$ otherwise