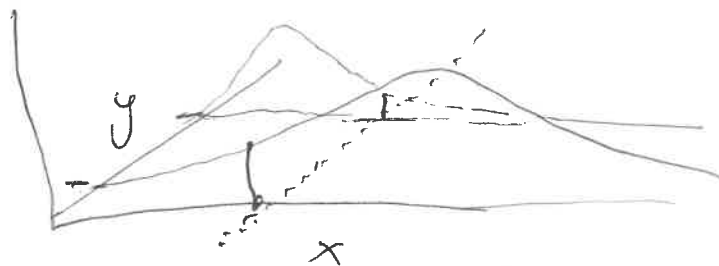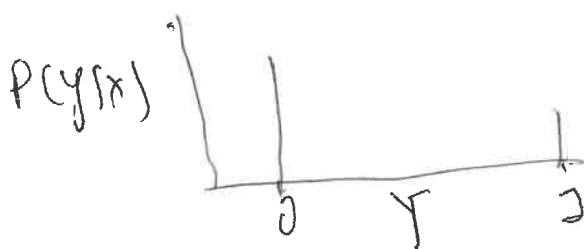# I. Bayes Classification.

$P(X, Y) =$



$P(Y|x) =$



① Pick the most likely one if you dont want to be wrong.

If you made $N$ decisions, avg you'd be wrong on $N(P(\bar{y}|x))$ times.

# II. But now assign costs.

- $C(\text{cold} \to \text{cold}) \to$ take meds.

  $C(\text{cold} \to Nc) \to$ Admitted to ER

  $C(Nc \to \text{cold}) \to$ Poison yourself with meds

  $C(Nc \to Nc) \to$ Nothing.

$E(C(\text{cold})) = \text{Cost}(\text{cold} \to \text{cold}) \cdot P(\text{cold}|x)$
$+ \text{Cost}(Nc \to \text{cold}) \cdot P(Nc|x)$

$E(C(Nc)) = \text{Cost}(\text{cold} \to Nc) \, P(\text{cold}|x) +$
$\text{Cost}(Nc \to Nc) \, P(Nc|x).$

find the min

# Minimum Bayes Risk classification.

$$\hat{Y} = \underset{Y}{\arg\min} \; E(C(Y))$$

$$= \underset{Y}{\arg\min} \; \sum_{y^*} C(y^* \to y) \, P(Y^*|x)$$

$$\underset{Y}{\arg\min} \; \sum_{Y^*} Cost(y^*, y) \, P(Y^*|x).$$

For cost = $1/0$  (0 for $y^* = y$, 1 else)

$$\hat{Y} = \underset{Y}{\arg\min} \; \sum_{Y^* \neq Y} P(Y^*|x) = \underset{Y^* \neq Y}{\arg\min} (1 - P(Y|x))$$

$$= \underset{Y}{\arg\max} \; P(Y|x) \longrightarrow \text{Usual MAP classifier}$$

$$\times \qquad \overset{}{\underset{\longrightarrow}{\quad}} \qquad \overset{}{\underset{\longrightarrow}{\quad}}$$

The above example was simple. Extend to Structured Pred setup.

$$\hat{Y} = \underset{Y}{\arg\min} \; E(Cost(Y))$$

$$= \underset{Y}{\arg\min} \; \sum_{Y^*} Cost(Y^*, Y) \, P(Y^*|x)$$
$$\hookrightarrow \text{Exponential sum.}$$

① Simple Example.    $Y$ = states seq.
   $P(X, Y) =$ HMM.
   Cost = $1/0$
   soln = ?           $\underset{Y}{\arg\max} \; \dfrac{P(Y|x)}{\text{Viterbi}}$

Simple example -

$P(x, y) = $ PCFG,    $Y = $ tree

Cost $= $ 1/0

Soln $= $ ?    $\underline{\underset{Y}{argmax} \; P(Y|x)}$

CYK

More generally, MBR decoding.

HMM variant



MAP decoding is ~~the~~ ~~first~~ MBR with 0/1 cost.

TRIVIAL solution 1
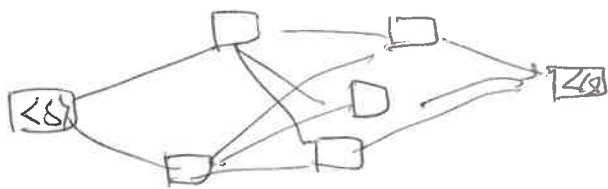
Not defining $C(Y_i^*, Y)$, keeping it generic within our problem setting.

Problem → define a loss between 2 word strings. $L(w^*, w)$.

E.g. Levenshtein distance, or distance which makes some kind of errors more important.

Basic setup ⇒ A word graph.



Nodes & edges have log probs.

Algo 1 : Hack, greedy.

Use $A^*$ to find $N$ most likely decodes.

$$\left.\begin{array}{c} W_1^* \\ \vdots \\ W_N^* \end{array}\right\rangle \text{Candidates.}$$

$$\widehat{w} = \text{argmin}_i \ E\left[\text{Cost}(W_i^*)\right]$$

How do you compute this cost?

Posterior decoding with constraints

$$\hat{S}^{P_k} = \arg\max_{S \in P_A} \prod_{i=1}^{L} P(S_i \mid O) \quad \left\{ \begin{array}{l} \text{I can} \\ \text{do this} \\ \text{buz of} \\ \text{synchronous decoding} \end{array} \right.$$

subject to constraint $\delta(\ell_1, \ell_2)$ (sep constr).

Init:

$$V_{start}(0) = 1 \qquad V_e(0) = 0 \qquad k \neq start$$

Recurse

$$V_k(i) = \max_{S \in Y} \left( V_S(i-1) \, \delta'(S,k) \right) \frac{P(S_{i} = k \mid O)}{r_k(e_i \mid d)}$$

$$P_i(k) = \arg\max_{S \in Y} \left( \qquad \right)$$

Terminati

$$P(\hat{S}^{P_k} \mid O) = \max_S \; V_S[2] \, \delta'(S, END)$$

$$S_L^{P_k} = \arg\max_{S \in S_1}$$

Trace.

$$S_{t-1}^{P_C} = P_t\left(S_t^{P_C}\right).$$

Assignment: $A_i = label\left(S_{t}^{P_C}\right)$

[Draw diagram]

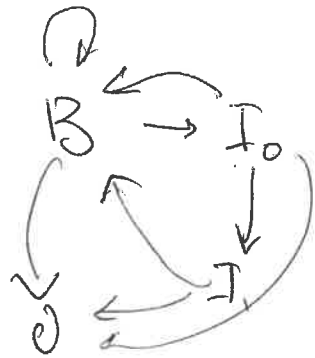If we have multiple states
for the same label, ~~average~~ sum
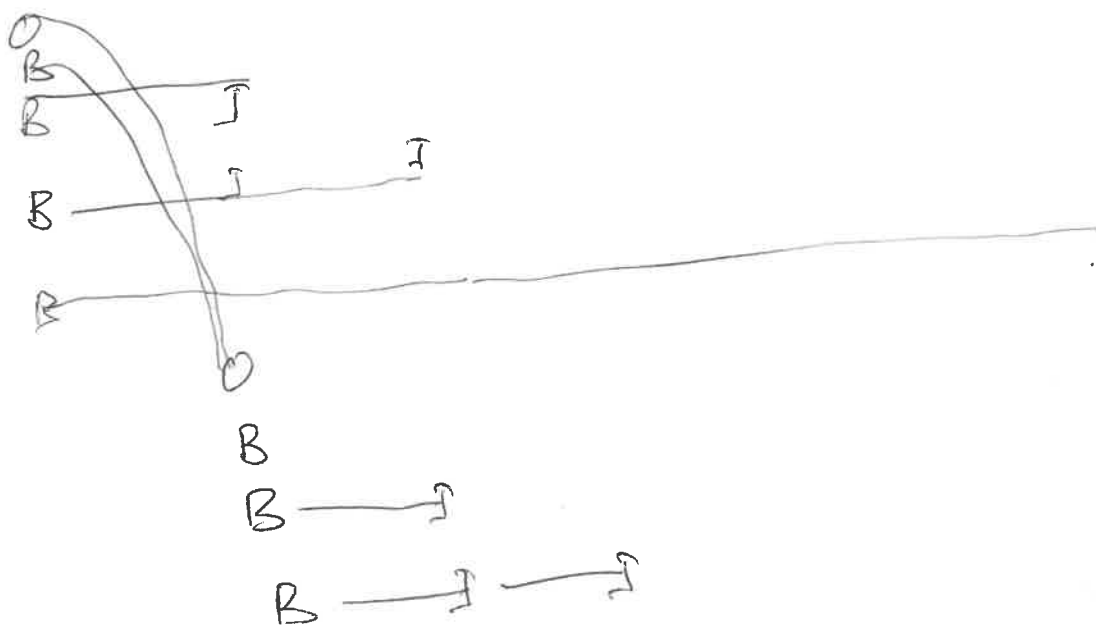first

How about when you have constraints?

E.g BIO tagger



Now you end up with something that looks like Viterbi but has to keep track of sequences two
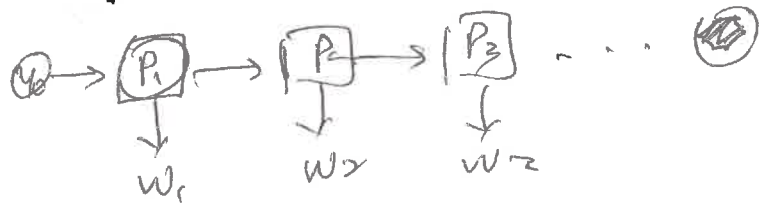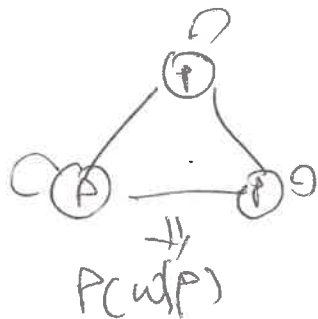
W_1        W_2       W_3       W_4      W_5       W_6



Constrained posterior decoding.

HMM for POS tagging



$$P(w|P)$$

Cost = # mislabelled tokens.

$$= \sum_{i=1}^{n} \mathbb{I} (\text{label}(i) \neq \text{label}^*(i))$$

$$= \sum_{i=1}^{n} \mathbb{I}(y_i \neq y_i^*)$$

$$\text{Expected cost} = E \sum_{i=1}^{n} \mathbb{I}(y_i \neq y_i^*)$$

$$\underset{Y}{\text{argmin}} \ \text{Cost} = \underset{y_1 \cdots y_N}{\text{argmin}} \ E \sum_{i=1}^{n} \mathbb{I}()$$

$$= \underset{y_1 \cdots y_N}{\text{argmin}} \sum_{i=1}^{n} E \mathbb{I}(y_i \neq y_i^*)$$

$$\Rightarrow y_i = \underset{y}{\text{argmin}} \ E \mathbb{I}(y \neq y_i^*)$$

$$= \underset{y}{\text{argmin}} \sum_{y^*} P(Y^*|x) \mathbb{I}(y = y_i^*)$$

$$= \underset{y}{\text{argmax}} \ P(y_i = y | x)$$

Its just a series of local decisions

$$T_G = \underset{T}{\text{argmax}} \ E\left(\frac{4}{N_c}\right) = \frac{1}{N_c} \underset{T}{\text{argmax}} \ E\,\mathbb{I}\left((s_i, t_i, R) \in j\right)$$

$$= \underset{T}{\text{argmax}} \ \sum_N P(s_i, t_i, R \mid w_i^n).$$

$$T_G = \underset{T}{\text{argmax}} \ \sum_{T_c} P(T_c \mid w_i^n) \ |T \cap T_c|$$

$$= \underset{T}{\text{argmax}} \ \sum_{T_c} P(T_c \mid w_i^n) \sum_{(s,t,x) \in T} \mathbb{I}\left((s,t,x) \in T_c\right)$$

For a PCFG.

$$P\left(S \rightarrow w_i^{s-1} X w_{t+1}^n \mid w_i^n\right) = \sum_{T_c} P(T_c \mid w_i^n) \,\mathbb{I}\left((s,t,x) \in T_c\right)$$

$$T_G = \underset{T}{\text{argmax}} \ \sum_{(s,t,x) \in T} P\left(S \rightarrow w_i^{s-1} X w_{t+1}^n \mid w_i^n\right)$$

But

$$P\left(S \rightarrow w_i^{s-1} X w_{t+1}^n \mid w_i^n\right)$$

$$= \frac{P(S \rightarrow \cdots , w_i^n)}{P(S \rightarrow w_i^n)} \cdot \frac{P(S \rightarrow \cdots X \cdots) P(X \rightarrow w_s^t)}{P(S \rightarrow w_i^n)}$$

$$= \frac{\beta(S, t) \alpha(S, t)}{P(S \rightarrow w_i^n)}$$

$$\Rightarrow T_G = \underset{T}{\text{argmax}} \ \sum_{(s,t,x) \in T} \gamma(s, t, x)$$

Goodman '96

labeled match $\quad (s,t,e) \rightarrow (s,t,R)$

bracketed match $\quad (s,t)R) \rightarrow (s,r,*)$

CYK etc optimize labeled TREE match.

$$T^* = T \quad , \qquad L/N_c = 1 \implies \begin{array}{c} cost = 0 \\ else = 1 \end{array}$$

$\quad L =$ no of labelled brackets matching.

V. Imp for eg Travel agent

❀ "find me all flights on tuesday}"

If we split it wrong, we'll wait till
tuesday & get a wrong ans.


☒ But in MT "his credentials are nothing
which should be ~~match~~ laughed at"
& MT mis aligns → `His creds are nothing,
which should make ~~you~~ laugh $\neq$ good,
But still helps.
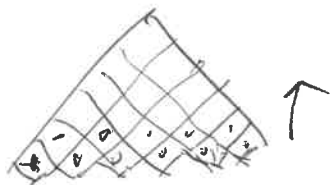Here we want labelled RECALL.

$$T_{G'} = \underset{T'}{argmax} \; L/N_c$$

Algo for max recall rate parse.

$$\text{Max}C\,(S,t) = \max_X \phi(S, t, x)$$

$$+ \max_{r \mid s \le r < t} \Big(\text{Max}C\,(S,r) + \text{Max}C\,(r+1, t)\Big)$$

for $r$ rules
$k$ nts, $O(n^3 + kn^2)$

Dominated by outside prob computation