

4.2 Hidden Markov Models

Edwin Simpson

Department of Computer Science,
University of Bristol, UK.

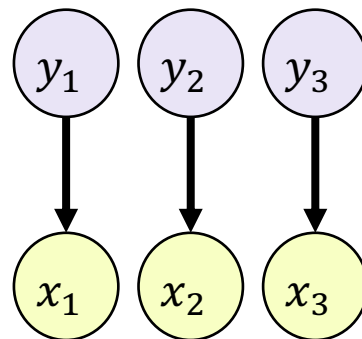
From Naïve Bayes to a Sequential Model: NB for Single Tokens

$$P(y_i|x_i) \propto P(x_i|y_i)P(y_i)$$

Maximum likelihood parameter estimates

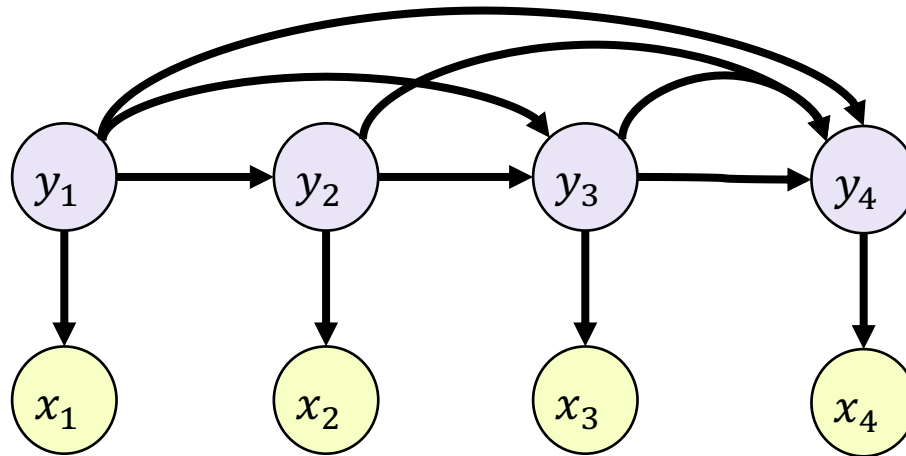
$$P(y_i = c) = \frac{\text{num_tokens_with_tag_c}}{\text{total_num_toks}}$$

$$P(x_i = w|y_i = c) = \frac{\text{count}(w \mid c) + 1}{\sum_{w' \in V} (\text{count}(w' \mid c) + 1)}$$



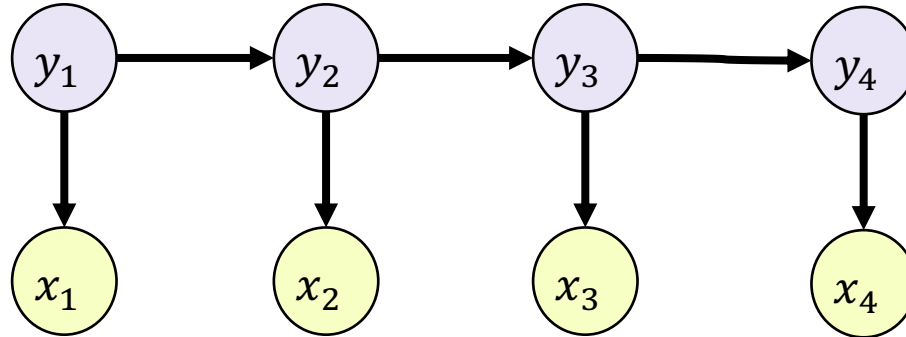
From Naïve Bayes to a Sequential Model:

$$P(y_i|x_{-i}) \propto P(x_i|y_i)P(y_i|y_{-})$$



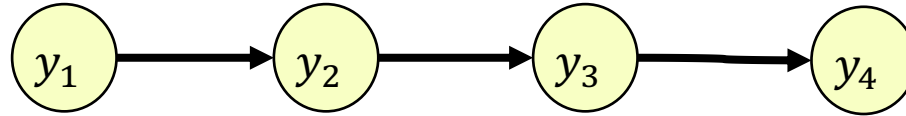
From Naïve Bayes to a Sequential Model:

$$P(y_i|x_{-i}) \propto P(x_i|y_i)P(y_i|\mathbf{y}_{-}) \approx P(x_i|y_i)P(y_i|y_{i-1})$$

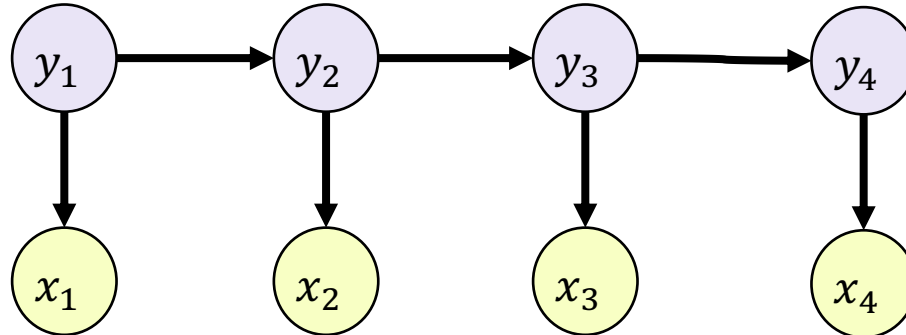


Markov Models and Hidden Markov Models

Markov assumption: $P(y_i | \mathbf{y}_{-}) \approx P(y_i | y_{i-1})$



Hidden Markov model (HMM): the states (y variables) are hidden and we observe x instead.



Transition Matrix

$$P(y_i|x_{-i}) \propto P(x_i|y_i)P(y_i|y_{i-1})$$

Transition matrix, A	Current tag y_i	
Previous tag y_{i-1}	0.5	0.5
	0.9	0.1

Initial probabilities π for $p(y_1)$ instead of the transition matrix	
0.9	0.1

Maximum likelihood parameter estimate

$$P(y_i = c | y_{i-1} = d) = \frac{\text{num_tokens_with_tag_c_preceded_by_d}}{\text{total_num_toks_with_tag_d}}$$

Observation Model

$$P(y_i|x_{-i}) \propto P(x_i|y_i)P(y_i|y_{i-1})$$

- $P(x_i|y_i)$ is defined by the **observation** or **emission model**, **B** ;
- The HMM generalises the observation model to allow any type of word features:
 - Word embeddings – multivariate Gaussian as observation distribution;
 - D binary features, e.g., presence in sentiment lexicons -- $\prod_{d=1}^D P(x_{id}|y_i)$

Learning the Observation Model

- Words as observations (as in naïve Bayes):

Maximum likelihood parameter estimate

$$P(x_i = w | y_i = c) = \frac{\text{count}(w | c) + 1}{\sum_{w' \in V} (\text{count}(w' | c) + 1)}$$

- Word embeddings as Gaussian-distributed observations:

Maximum likelihood parameter estimates

$$P(\vec{x}_i = w | y_i = c) = \mathcal{N}(\vec{x}_i, \mu_c, \Sigma_c)$$

Summary

- The HMM models token likelihoods similarly to naïve Bayes.
- It also models the transition probabilities.
- To make learning and prediction tractable, we use the Markov assumption: the probability of each tag depends only on m predecessors in an order- m Markov model.