



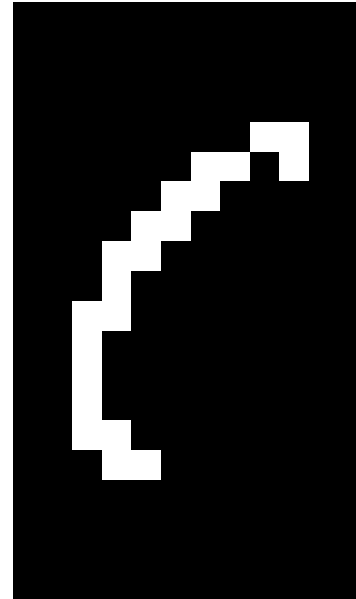
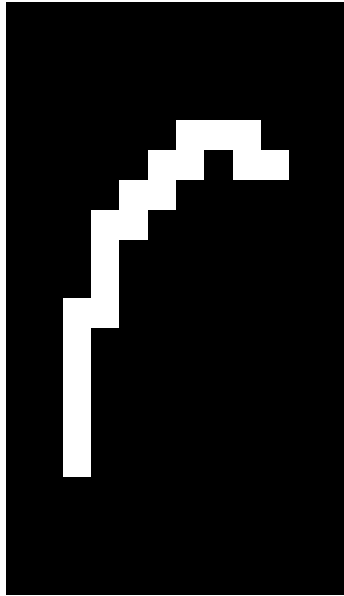
Hidden Markov Models

Machine Learning 10-601
March 17, 2008

Tom M. Mitchell
Machine Learning Department
Carnegie Mellon University

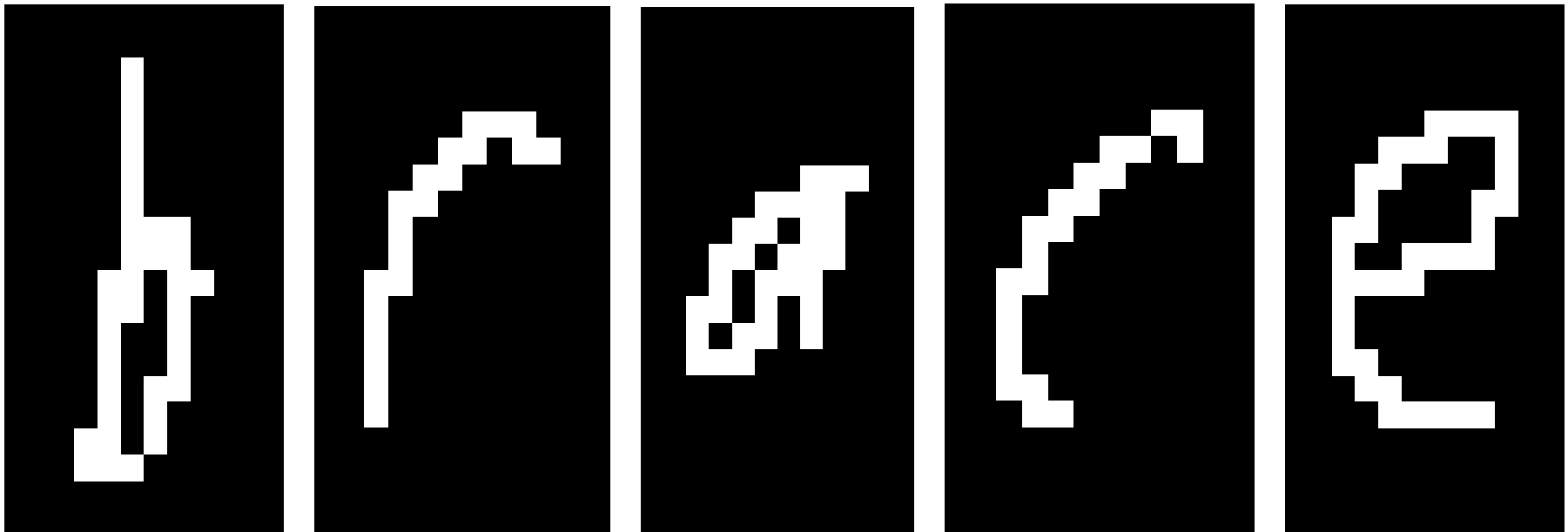
With thanks to Prof. Carlos Guestrin for some of these slides

Handwriting recognition

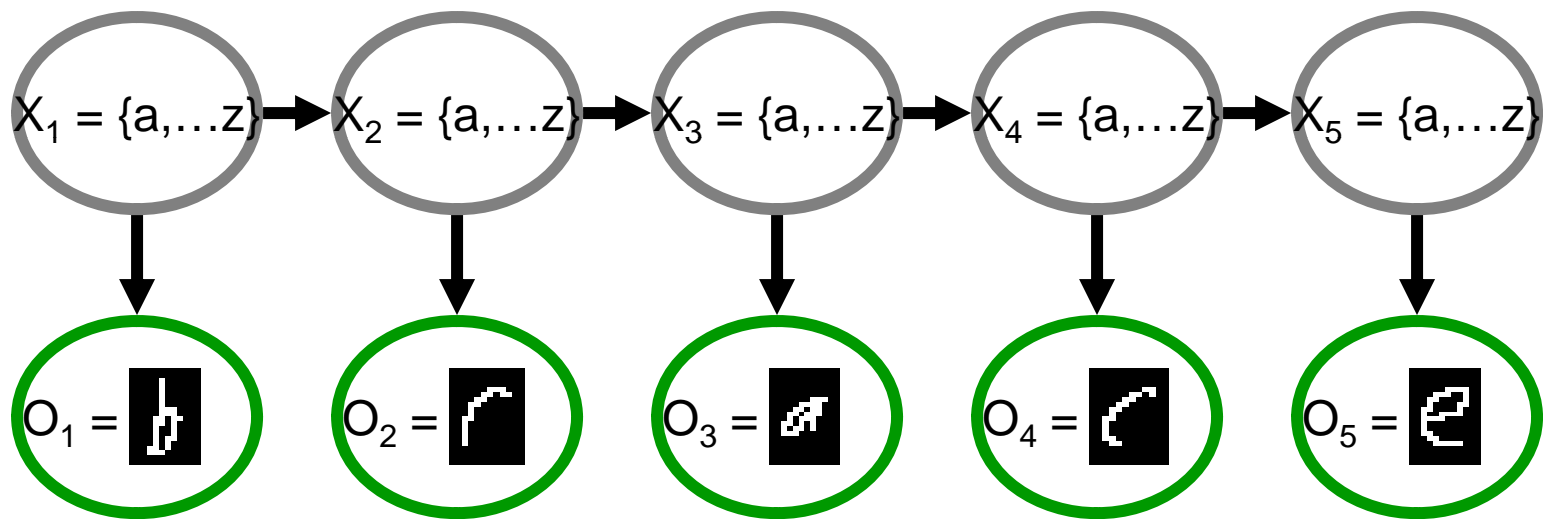


Character recognition, e.g., logistic regression

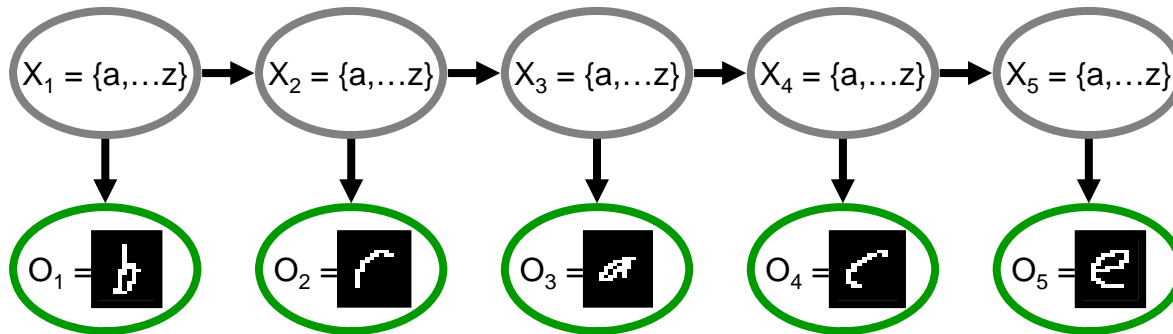
Example of a hidden Markov model (HMM)



Understanding the HMM Semantics



HMMs semantics: Details



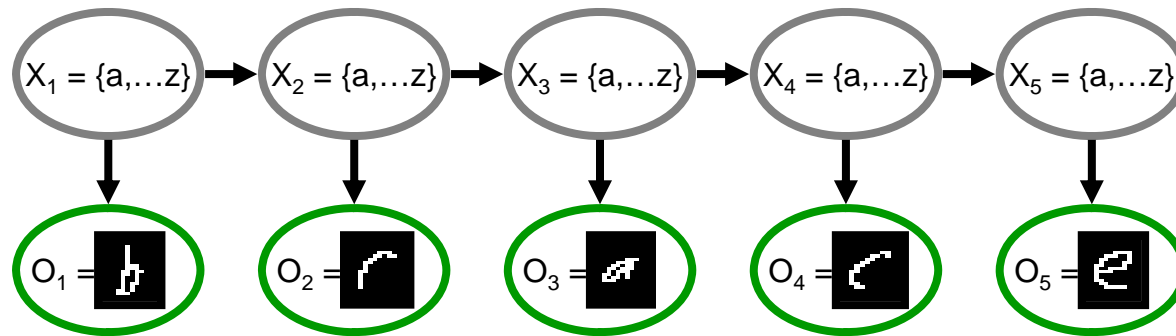
Just 3 distributions:

$$P(X_1)$$

$$P(X_i | X_{i-1})$$

$$P(O_i | X_i)$$

Using and Learning HMM's

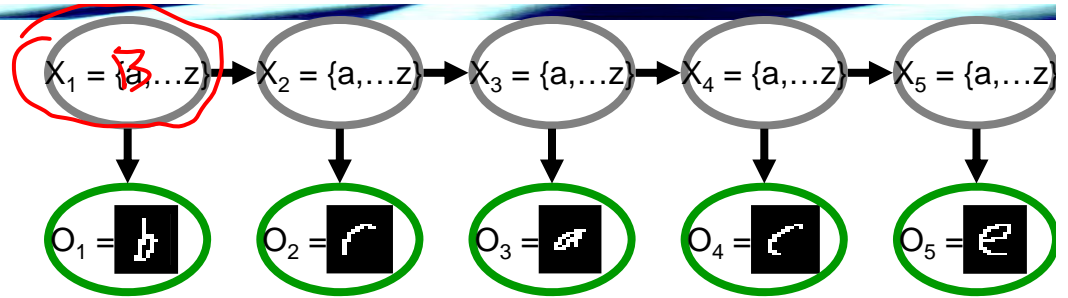


$$P(X_i | o_{1..n})$$

Core HMM questions:

1. How do we calculate $P(o_1, o_2, \dots, o_n)$?
2. How do we calculate argmax over x_1, x_2, \dots, x_n of $P(x_1, x_2, \dots, x_n | o_1, o_2, \dots, o_n)$?
3. How do we train the HMM, given its structure and
 - 3a. Fully observed training examples: $\langle x_1, \dots, x_n, o_1, \dots, o_n \rangle$
 - 3b. Partially observed training examples: $\langle o_1, \dots, o_n \rangle$

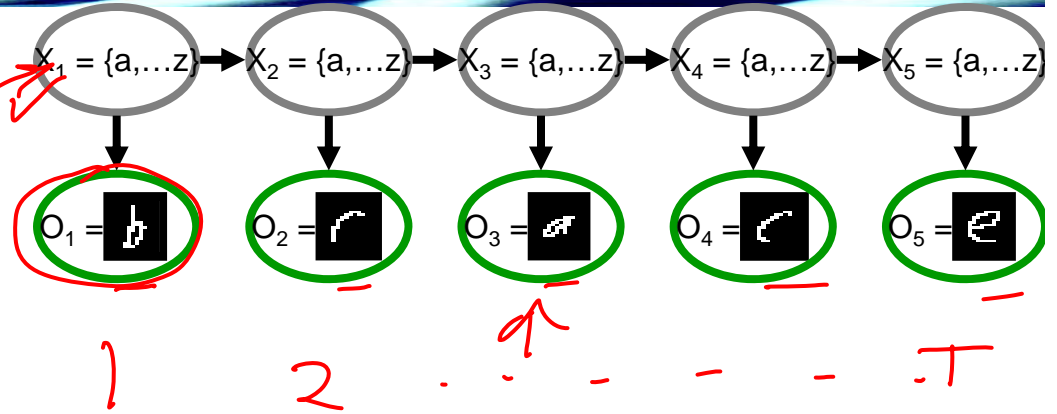
How do we generate a random output sequence following the HMM $P(o_1, o_2, \dots, o_T)$



How do we compute

$\rightarrow P(o_1, o_2, \dots, o_T)$

N poss values



1. brute force:

2. Forward algorithm (dynamic progr., variable elimination):

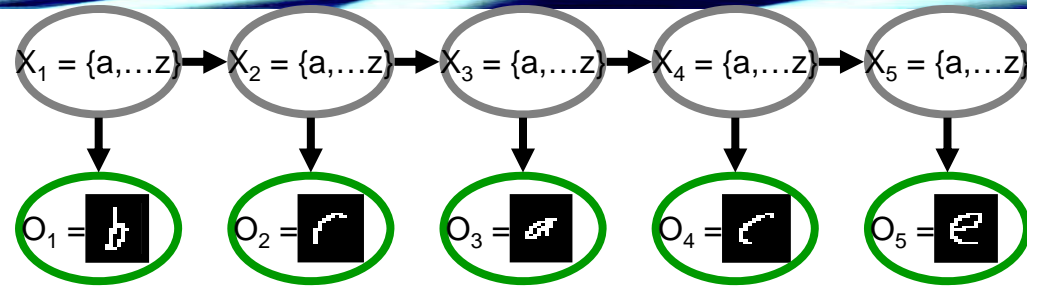
$N^2 T$ define $\alpha_t(k) = P(o_1, o_2, \dots, o_t, X_t = k)$

$\alpha_1(k) = P(X_1 = k) P(o_1 = o_1 | X_1 = k) = P(o_1, X_1 = k)$

$\alpha_{t+1}(k) = \left[\sum_{j=1}^N \alpha_t(j) P(X_{t+1} = k | X_t = j) \right] P(o_{t+1} = o_{t+1} | X_{t+1} = k)$

$P(o_1, o_2, \dots, o_T) = \sum_k \alpha_T(k)$

How do we compute
 $P(X_t = k | o_1, o_2, \dots, o_T)$



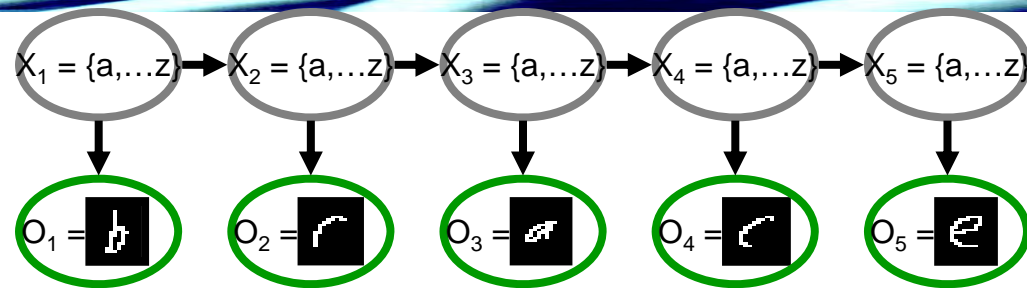
2. Backward algorithm (dynamic progr., variable elimination):

$$\alpha_t(k) = P(o_1, o_2, \dots, o_t, X_t = k)$$

define $\beta_t(k) = P(o_{t+1}, o_{t+2}, \dots, o_T | X_t = k)$

$$P(X_t = k | o_1, o_2, \dots, o_T) = \frac{P(X_t = k, o_1, o_2, \dots, o_T)}{P(o_1, o_2, \dots, o_T)}$$

$\alpha_t(k) \beta_t(k)$
 $\sum_k \alpha_T(k)$



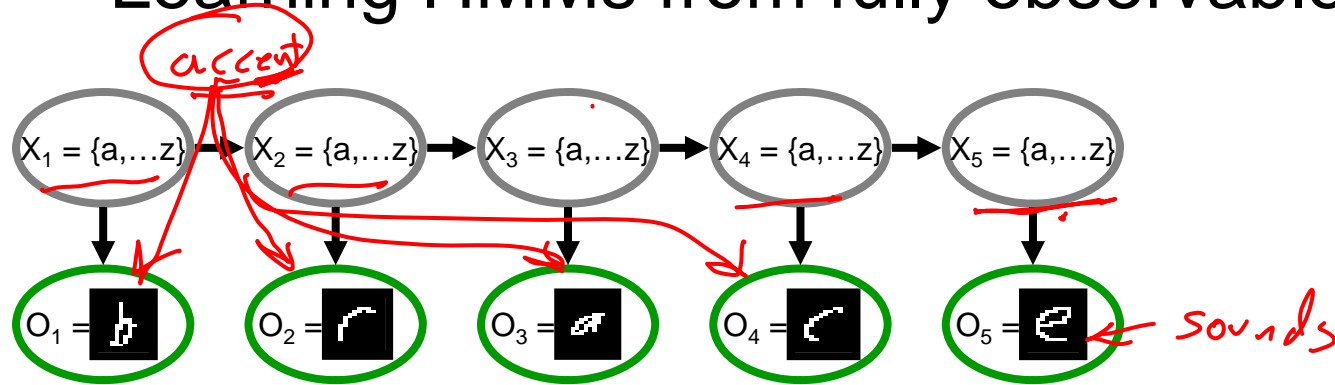
How do we compute

$$\arg \max_{x_1, \dots, x_T} P(x_1, \dots, x_T | o_1, o_2, \dots, o_T)$$

Viterbi algorithm, based on recursive computation of

$$\delta_t(k) = \max_{x_1, \dots, x_{t-1}} P(x_1, \dots, x_{t-1} X_t = k, o_1, o_2, \dots, o_t)$$

Learning HMMs from fully observable data: easy



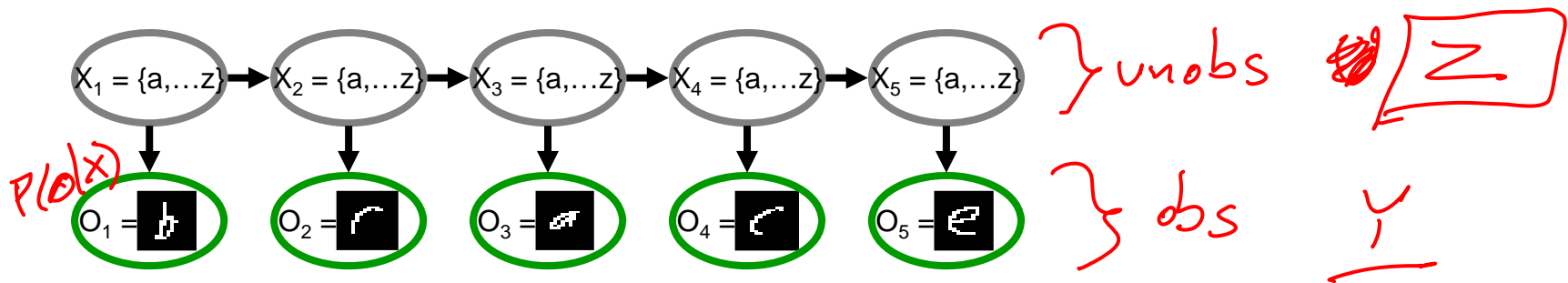
Learn 3 distributions:

$$P(X_1)$$

$$P(O_i | X_i)$$

$$P(X_i | X_{i-1})$$

Learning HMMs when only observe $o_1 \dots o_T$



EM

initial model params θ

Loop to convergence.

E step: use current θ to calc. expected values of Z
 (eg. Viterbi alg)

- $E[\# \text{ times trans from } x_i \text{ to } x_j]$
- $E[\# \text{ times in } x_i]$
- $E[\# \text{ times in } x_i \text{ and obs } o_j]$

M step re-maximize likelihood by choosing θ'

$$\theta' \leftarrow \underset{\theta'}{\operatorname{argmax}} E \left[P(o_1, o_2, \dots, o_T, x_1, \dots, x_T | \theta') \right]$$

\uparrow
 $P(z|Y, \theta)$



What you need to know

- Hidden Markov models (HMMs)
 - Very useful, very powerful!
 - Speech, OCR, time series, ...
 - Parameter sharing, only learn 3 distributions
 - Trick reduces inference from $O(n^2)$ to $O(n)$
 - Special case of Bayes net
 - Dynamic Bayesian Networks

Thanks to Carlos Guestrin for many slides