

Royal Institute of Technology

MACHINE LEARNING 2-CONT HMM. EM

Lecture 8

LAST LECTURE

- DGM semantics
- ★ UGM
- ★ De-noising
- ★ HMMs
 - Applications (interesting probabilities)
 - DP for generation probability etc.

INFERENCE IN THE OCCASIONALLY DISHONEST CASINO

Grey regions are states corresponding to biased die







• Filtering: $p(z_t|x_{1:t})$,

online

- Smoothing, MAP state: $p(z_t|x_{1:T})$ offline
- Viterbi, MAP path
 argmax p(z_{1:T}|x_{1:T})

THIS LECTURE

- ★ Backward
- ★ Smoothing
- ★ Sampling
- ★ Viterbi
- ★ K-means (inspiration)
- ★ GMM (towards EM)

MARKOV CHAINS (DISCRETE)

Probabilities on outgoing edges sum to one



THE OCCASIONALLY DISHONEST CASINO



Rolls: 6641532161621152346532143566342616552 Die: LLLLLLLLLLFFFFFFLLLLLLLLLLLFFF



WHAT AN HMM DOES

- \star Starts in the state z₁
- \star When in state z_t
 - outputs $p(\mathbf{x}_t|\mathbf{z}_t)$ B_{x_t,t_t}
 - moves to $p(z_{t+1}|z_t)$ A_{z_{t+1},z_t}
- Stops after a fixed number of steps or when reaching a stop step

The parameters

AN HMM CAN BE SEEN AS A DGM



SPECIAL CASE: HIDDEN MARKOV MODEL (HMM)

Combinations of the transition distributions



Combinations of emission the emission distribution



TRANSITION PROBABILITIES FOR 4 STATES HMM

EMISSION PROBABILITIES - HMM WITH 4 STATES & 3 SYMBOLS

 $Z_1 \to Z_2 \to Z_3 \to \cdots \to Z_T \to Z_{T+1}$ $\begin{bmatrix} x_2 \\ x_2 \end{bmatrix} \begin{bmatrix} x_3 \end{bmatrix}$ $\dot{x_1}$ x_T

All the same

State	1	2	3
1	B ₁₁	B ₂₁	B ₃₁
2	B ₁₂	B ₂₂	B ₃₂
3	B ₁₃	B ₂₃	B ₃₃
4	B ₁₄	B ₂₄	B ₃₄

Sum rule gives Zt+1

Defined by

$$b_t(k) := p(x_{t+1:T} | Z_t = k) = \sum_{l \in [K]} p(x_{t+1:T}, Z_{t+1} = l | Z_t = k)$$

Each term in the sum is a probability of an event

"which is an AND of"

$$\begin{aligned} X_{t} = k \to Z_{t+1} = l & Z_{t+1} = l & Z_{t+1} = l \to ? \to \cdots ? \\ \downarrow & \downarrow & \downarrow \\ X_{t+1} & & x_{t+2} \to \cdots ? \\ \end{bmatrix} \end{aligned}$$

Backward recursion $b_t(k) := p(\boldsymbol{x}_{t+1:T} | \boldsymbol{Z}_t = k)$

• DP also for the backward variable bt

$$b_t(k) = \sum_{l} \underbrace{p(\mathbf{Z}_{t+1} = l | \mathbf{Z}_t = k)}_{\text{transition}} \underbrace{b_{t+1}(l)}_{\text{"smaller"}} \underbrace{p(\mathbf{x}_{t+1} | \mathbf{Z}_{t+1} = l)}_{\text{emission}}$$

 Implementation analogous, complexity same

OFF-LINE SMOOTHING

$p(\mathbf{Z}_t = k | \mathbf{x}_{1:T}) \propto f_t(k) p(\mathbf{x}_t | \mathbf{Z}_t = k) b_t(k)$

emission

OFF-LINE SMOOTHING

$p(\mathbf{Z}_t = k | \mathbf{x}_{1:T}) \propto f_t(k) \underbrace{p(\mathbf{x}_t | \mathbf{Z}_t = k)}_{k} b_t(k)$

emission

Up to a constant

TWO SLICED SMOOTH MARGINALS - MARGINAL OVER PAIRS OF STATES

 $p(\mathbf{Z}_t = k, \mathbf{Z}_{t+1} = l | \mathbf{x}_{1:T})$

Can be computed from forward and backward similarly

SAMPLING FROM POSTERIOR



How much did each previous state contribute to the probability mass of the present state?

Forward recursion

 $f_t(k) = \sum_{l} \underbrace{f_{t-1}(l)}_{\text{smaller}} \underbrace{p(\boldsymbol{x}_{t-1} | \boldsymbol{Z}_{t-1} = l)}_{\text{emission } B_{\boldsymbol{x}_{t-1}, l}} \underbrace{p(\boldsymbol{Z}_t = k | \boldsymbol{Z}_{t-1} = l)}_{\text{transition } A_{lk}}$

 $\begin{array}{c|c} \hline Z_1 \to Z_2 &\to Z_3 \to \cdots \to Z_T \to Z_{T+1} \\ \downarrow & \downarrow & \downarrow & \downarrow \\ x_1 & x_2 & x_3 & & x_T \end{array}$

BACKWARDS SAMPLING OF POSTERIOR

Sample
$$z_{1:T+1}^s \sim p(\mathbf{Z}_{1:T+1} = k | \mathbf{x}_{1:T})$$
 by
(i) $z_{T+1}^s \sim p(\mathbf{Z}_{T+1} = k | \mathbf{x}_{1:T}) \propto p(\mathbf{Z}_{T+1} = k, \mathbf{x}_{1:T})$
(ii) $z_t^s \sim p(\mathbf{Z}_t = k | \mathbf{Z}_{t+1} = l, \mathbf{x}_{1:t})$
 $\propto f_t(k) p(\mathbf{Z}_{t+1} = l | \mathbf{Z}_t = k) p(\mathbf{x}_t | \mathbf{Z}_t = k)$

- Sample from posterior
- Sample in order z_T, \dots, z_1
- Start somewhat differently

We want $\operatorname{argmax}_{\boldsymbol{z}_{1:T}} p(\boldsymbol{z}_{1:T} | \boldsymbol{x}_{1:T})$

Not!

 $(\operatorname{argmax}_{\boldsymbol{z}_1} p(\boldsymbol{z}_1 | \boldsymbol{x}_{t+1:T}), \dots, \operatorname{argmax}_{\boldsymbol{z}_T} p(\boldsymbol{z}_T | \boldsymbol{x}_{t+1:T})))$

Viterbi variable

$$v_t(k) := \max_{z_{1:t-1}} p(z_{1:t-1}, Z_t = k, x_{1:t})$$

It gives what we want

 $\max_k v_T(k)$

VITERB

- MAP path
- Viterbi learning: used, as approximation, to speed up parameter learning
- Again DP now with Viterbi variable
- For the path, use back pointers

Chapter 3

EXPECTATION-MAXIMIZATION THEORY

3.1 Introduction

Learning networks are commonly categorized in terms of supervised and unsupervised networks. In unsupervised learning, the training set consists of input training patterns only. In contrast, in supervised learning networks, the training data consist

$$\mathcal{N}(\boldsymbol{x}|\boldsymbol{\mu},\sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(\boldsymbol{x}-\boldsymbol{\mu})^2\right)$$

$$\mathcal{N}(\boldsymbol{x}|\boldsymbol{\mu},\boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2}|\boldsymbol{\Sigma}|^{1/2}} \exp\left(-\frac{1}{2}(\boldsymbol{x}-\boldsymbol{\mu})^T\boldsymbol{\Sigma}^{-1}(\boldsymbol{x}-\boldsymbol{\mu})\right)$$





0.45

-0.05

GAUSSIAN – MVN



TWO DIMENSIONAL NORMAL

80

K-MEANS

★ Data vectors $D=\{x_1,...,x_N\}$

★ Randomly selected classes $z_1,...,z_N$

 \star Iteratively do

$$\boldsymbol{\mu}_c = \frac{1}{N_c} \sum_{n: z_n = c} \boldsymbol{x}_n,$$

where
$$N_c = |\{n : z_n = c\}|$$

$$z_n = \operatorname{argmin}_c || \boldsymbol{x}_n - \boldsymbol{\mu}_c ||_2$$

 \star One step O(NKD), can be improved

ASSIGN \leq (b) (C) (a) 2 2 2 0 0 × X -2 -2_2 -2 2 -22 -22 0 0 0 (d) (e) (f) 2 2 2 ł 0 0 -2 -2-2 -2 0 2 -2 2 -2 2 0 0 (g) (h) (i) 2 2 2 0 0 -2-2-2-2 2 -2 2 -22 0 0 0

ASSIGNING POINTS TO MULTIPLE MEANS



K-MEANS AS GMM

 \star Fixed variance, a Gaussian and mean per cluster, i.e., $\theta_c = (\mu_c, \sigma^2)$

 \star Idea: each point can belong to several means (clusters)

 \star Use responsibilities to find means

1

$$r_{nc} = p(z_n = c | \boldsymbol{x}_n, \boldsymbol{\theta}) = \frac{p(z_n = c | \boldsymbol{\theta}) p(\boldsymbol{x}_n | z_n = c, \boldsymbol{\theta})}{\sum_{c=1}^{C} p(z_n = c | \boldsymbol{\theta}) p(\boldsymbol{x}_n | z_n = c, \boldsymbol{\theta})}$$

$$\boldsymbol{\mu}_{c} = \frac{1}{N_{c}} \sum_{n} r_{nc} \boldsymbol{x}_{n}, \qquad \text{where } N_{c} = \sum_{n} r_{nc}$$

IMAGE SEGMENTATION WITH K-MEANS

K=2



K = 3





Original image



GAUSSIAN MIXTURE MODELS (GMM)

$$\mathcal{D} = (oldsymbol{x}_1, \dots, oldsymbol{x}_N)$$

 $\boldsymbol{x}_n = (x_{n1}, \dots, x_{nD})$

 $Z \sim \operatorname{Cat}(\pi)$

 $p(\boldsymbol{X}|Z=c) = p_c(\boldsymbol{X}) = \mathcal{N}(\boldsymbol{X}|\boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)$



1-DIM GAUSSIAN MIXTURE MODELS

 $\mathcal{D} = (oldsymbol{x}_1, \dots, oldsymbol{x}_N)$

 $Z \sim \operatorname{Cat}(\boldsymbol{\pi})$

 $p(\boldsymbol{X}|Z=c) = p_c(\boldsymbol{X}) = \mathcal{N}(\boldsymbol{X}|\boldsymbol{\mu}_c, \sigma_c)$ $\boldsymbol{\theta}_c = (\boldsymbol{\mu}_c, \sigma_c)$

Z hidden $\sim \operatorname{Cat}(\boldsymbol{\pi})$



$$oldsymbol{X} \sim \mathcal{N}(oldsymbol{\mu}_c, \sigma_c)$$

EXAMPLE

 z_n is red with probability 1/2, green with probability 3/10, blue with probability 1/5



So,

 $p(x_n, z_n) = p(z_n)p(x_n|z_n)$

and

$$p(x_n) = \sum_{c=1}^{C} p(z_n = c) p(x_n | z_n = c) = \sum_{c=1}^{C} \pi_c p(x_n | z_n = c)$$

and

$$p(z_n = c | x_n) = \frac{p(z_n = c, x_n)}{p(x_n)} = \frac{\pi_c p(x_n | z_n = c)}{\sum_{c=1}^C \pi_c p(x_n | z_n = c)}$$

MLE -COMPLETE DATA FOR GAUSSIAN $\mathcal{D} = (z_1, x_1), \dots, (z_N, x_N))$ • Maximizing the complete log likelihood

$$l(\theta'; \mathcal{D}) = \sum_{c} N_c \log \pi'_c + \sum_{c} \sum_{n: I(z_n = c)} \log p(\boldsymbol{x}_n | \theta'_c)$$

Boils down to maximizing



MLE FOR GAUSSIAN

...and
$$\sum_{n:I(z_n=c)} \log \frac{1}{\sqrt{2\pi\sigma_c'^2}} \exp \left(-\frac{1}{2\sigma_c'^2}(x_n - \mu_c')^2\right)$$

is maximized by

by
$$\mu'_c = \frac{\sum_{n:I(z_n=c)} x_n}{N_c}$$

where
$$N_c = \sum_n I(z_n = c)$$

EM & EXPECTED LOG LIKELIHOOD (Q-TERM)

- Iteratively maximizing the expected log likelihood in practice always leads to a local maxima
- The expectation is over latent variables given data and current parameters
- We maximize the expression by choosing new parameters.

#