

MARCOV MODELS

$$P(x_{1:T}) = P(x_1) \prod_{t=1}^T P(x_t | x_{t-1}) \quad \text{STATIONARY CHAIN} \rightarrow \text{PARAMETER Tying. TRANSITION MATRIX} = A^N$$

VERILY USED IN LANGUAGE MODELS: $P(x_t = u)$ UNIGRAM STATISTICS, $P(x_t = u | x_{t-1} = v)$ BIGRAM MODELS, ... N-GRAM MODELS ARE NORMALIZED COUNTS. TO HANDLE MISSING N-GRAMS IN TRAINING EITHER MEAN DATA OR ADD-ONE SMOOTHING

DELETED INTERPOLATION: $A_{ju} = (1-\lambda)P_{ju} + \lambda P_{ju}$, CAN DO WITH ADAPTIVE WEIGHTS TOP LANG MODEL: INTERPOLATED KNEISER-NEY

OUT OF VOCABULARY WORDS: USE SYMBOL 'UNK' AND ASSIGN SOME PROBABILITY MASS TO IT; USE DIRICHLET PROCESS (INFINITE STATE SPACE)

STATIONARY DISTRIBUTION: INITIAL ROW STATE π_0 , $\pi_1 = \pi_0 A$, $\pi = \pi A$, BALANCE EQUATIONS. SOLVE $A^T v = v$, $\pi = v^T$, EIGENVECTOR WITH $\lambda = 1$

- IRREDUCIBLE: SIMPLY CONNECTED
- APERIODIC: IF RETURN TO ANY STATE DOESN'T OCCUR AT MULTIPLES OF TIME STEPS LONGER THAN 1
- RECURRENT: YOU WILL RETURN TO EVERY STATE WITH PROBABILITY 1
- NON-NULL: EXPECTED TIME TO RETURN TO STATE IS FINITE
- ERGODIC: APERIODIC, RECURRENT, NON-NULL

• EVERY IRREDUCIBLE ERGODIC MC HAS LIMITING DISTRIBUTION π , STATIONARY AND UNIQUE.

• DETAILED BALANCE EQUATIONS: IF π SATISFIES DBE, π IS STATIONARY. $\sum_i \pi_i A_{ij} = \sum_j \pi_j A_{ji} = \pi_j \sum_i A_{ji} = \pi_j$

HIDDEN MARKOV MODELS

DISCRETE TIME, DISCRETE STATE MC WITH HIDDEN STATES + OBSERVATION MODEL.

$$P(z_{1:T}, x_{1:T}) = P(z_{1:T}) P(x_{1:T} | z_{1:T}) = P(z_1) \prod_{t=1}^T P(z_t | z_{t-1}) \cdot \prod_{t=1}^T P(x_t | z_t)$$

• DISCRETE OBS \rightarrow OBS MATRIX

• CONTINUOUS OBS \rightarrow ~~ISOMORPHIC~~ CONDITIONAL GAUSSIAN
 $P(x_t | z_t = u, \theta) = N(x_t | \mu_u, \Sigma_u)$

• BLACK-BOX DENSITY MODELS ON SEQUENCES. LONG-RANGE DEPENDENCIES. NO MARKOV PROPERTY FOR OBSERVATIONS.

• TIME-SERIES PREDICTIONS

• GENERATING CLASSIFIERS

• APPLICATIONS: AUTO SPEECH RECOGNITION. x_t IS SPEECH RECOGNITION. z_t IS WORDSPOKEN • ACTIVITY RECOGNITION • PART OF SPEECH TAGGING • GENE FINDING • PROTEIN SEQUENCE ALIGNMENT

HMM INFERENCE

INFERR HIDDEN STATE SEQUENCE ASSUMING PARAMS ARE KNOWN. HIDDEN STATE \rightarrow BELIEF STATE

TYPES OF INFERENCE:

- FILTERING: $P(z_t | x_{1:t})$ ONLINE, AS DATA STREAMS IN. REDUCES NOISE BETTER THAN JUST $P(z_t | x_t)$. BAYES RULE IN SEQUENCE.

- SMOOTHING: $P(z_t | x_{1:T})$ OFFLINE, GIVEN ALL EVIDENCE. 'PAST GIVEN FUTURE'

- FIXED-LAG SMOOTHING: $P(z_{t-l} | x_{1:t})$, $l > 0$ LAG. BETTER THAN FILTERING BUT DELAY.

- PREDICTION: 'FUTURE GIVEN PAST' $P(z_{t+h} | x_{1:t})$, h IS PREDICTION HORIZON. EX. $P(z_{t+2} | x_{1:t}) = \sum_{z_{t+1}} \sum_{z_t} P(z_{t+2} | z_{t+1}) P(z_{t+1} | z_t) P(z_t | x_{1:t})$

FOR UP TRANSITION MATRIX AND APPLY TO CUR BELIEF STATE.

\rightarrow PREDICT OBSERVATIONS: $P(x_{t+h} | x_{1:t}) = \sum_{z_{t+h}} P(x_{t+h} | z_{t+h}) P(z_{t+h} | x_{1:t})$ POSTERIOR MARGINAL DENSITY

- MAP $\text{ARGMAX}_{z_{1:T}} P(z_{1:T} | x_{1:T})$ AKA VITERBI DECODING, MOST PROBABLE STATE SEQUENCE

- POSTERIOR SAMPLES: IF ≥ 1 PLAUSIBLE INTERPRETATION OF DATA, SAMPLE $z_{1:T} \sim P(z_{1:T} | x_{1:T})$

- PROBABILITY OF EVIDENCE: $P(x_{1:T}) = \sum_z P(z_{1:T} | x_{1:T})$, CLASSIFY SEQUENCES, MODEL-BASED CLUSTERING.

FORWARDS ALGORITHM

RECURSIVELY COMPUTE FILTERED MARGINALS $P(z_t | x_{1:t})$; ONLINE INFERENCE. EXPLOIT CI OF MARKOV STUFF

• **PREDICTION**: ONE-STEP AHEAD PREDICTION DENSITY. NEW PRIOR FOR t . $P(z_t = j | x_{1:t-1}) = \sum_i P(z_t = j | z_{t-1} = i) P(z_{t-1} = i | x_{1:t-1})$

• **UPDATE**: ABSORB OBSERVATION FOR t WITH BAYES $\alpha_t(j) = \sum_i P(z_t = j | x_{1:t}) = \frac{1}{Z_t} r(x_t | z_t = j) P(z_t = j | x_{1:t-1})$. $Z_t = \sum_j P(z_t = j | x_{1:t-1}) P(x_t | z_t = j)$
 - LOG PROB OF EVIDENCE: $\log r(x_{1:T} | \theta) = \sum_t \log Z_t$ NORM. CONSTANT

- **MATRIX-VECTOR FORM**: $\alpha_t \propto \psi_t \odot (\Psi^T \alpha_{t-1})$. α BELIEF STATE. $\psi_t = r(x_t | z_t = j)$ LOCAL EVIDENCE.

Ψ TRANSITION MATRIX. \odot HADAMARD PRODUCT, ELEMENTWISE MULTIPLICATION

FORWARD/BACKWARDS ALGORITHM

SMOOTHED MARGINALS, OFFLINE INFERENCE

$\alpha_t(j) = P(z_t = j | x_{1:t})$ BELIEF STATE AS BEFORE. $\beta_t(j) = P(x_{t+1:T} | z_t = j)$ CONDITIONAL LIKELIHOOD OF FUTURE EVIDENCE, COMPUTED BACKWARDS.

$\beta_{t-1} = \sum_j \beta_t(j) \psi_t(j) \psi(1,j) \rightarrow \beta_{t-1} = \Psi(\psi_t \odot \beta_t)$

• $\gamma_t(j) = \alpha_t(j) \beta_t(j)$ SMOOTHED POSTERIOR MARGINAL

• IF USING EM TO ESTIMATE TRANSITION MATRIX, NEED TO COMPUTE EXPECTATION OF $1 \rightarrow j$ TRANSITIONS. TWO-SLICE MARGINAL.
 $\xi_{t,t+1} \propto \Psi \odot (\alpha_t \odot \beta_{t+1})^T$

• **COMPLEXITY** BASELINE IS $O(K^2T)$. SHORTCUTS VIA MATRIX SPARSITY $O(TH)$. OR STATE SPACE SPARSITY $O(TH \log K)$. SINCE BOTH DEPEND ON T , TOO.

VITERBI ALGORITHM

MOST PROBABLE SEQUENCE OF STATES IN CHAIN GRAPHICAL MODEL. $Z^* = \text{ARGMAX}_Z P(z_{1:T} | x_{1:T})$ SHORTEST PATH THROUGH TRELLIS DIAGRAM

PATH WEIGHT: $\log \pi_1(z_1) + \log \phi_1(z_1) + \sum_{t=2}^T [\log \psi(z_{t-1}, z_t) + \log \phi_t(z_t)]$

JOINT MOST PROB. SEQUENCE OF STATES

↳ VITERBI
 ↳ GLOBALLY CONSISTENT

↳ NODES ARE STATES AT EACH STEP
 ↳ EDGES ARE LOG PROBABILITIES
 ≠ **MARGINAL** MOST PROB. SEQUENCE OF STATES

• DOES NOT JUST REPLACE SUM WITH MAX IN BACKWARDS PASS

• CHAINS OF MOST PROBABLE / LEAST COST PATH WITH BACKTRACKING.

• LOGS FOR NUMERIC / FINITNESS REASONS. $\log \delta_t(j) = \max_i \log \delta_{t-1}(i) + \log \psi(i,j) + \log \phi_t(j)$

• **COMPLEXITY**: $O(K^2T)$ TIME, $O(KT)$ SPACE

• CAN USE DISCRIMINATIVE METHOD TO RANK PATHS AND RESUME MORE, N-BEST

MPM: $\hat{z} = \text{ARGMAX}_Z P(z_1 | x_{1:T})$
 ↳ MORE ROBUST

FORWARD FILTERING/BACKWARDS SAMPLING

SAMPLE PATHS FROM POSTERIOR $Z_{1:T}^S \sim P(z_{1:T} | x_{1:T})$. FWD - FILTER. COMPUTE 2-SLICE POSTERIOR. COMPUTE CONDITIONALS. SAMPLE.

STRUCTURE: SAMPLE IN BACKPASS.

$P(z_t | z_{t+1}, x_{1:t}) \sim Z_t^S$. $P(z_t = j | z_{t+1} = i, x_{1:t}) = \frac{\phi_{t+1}(i) \psi(i,j) \alpha_t(j)}{\alpha_{t+1}(i)}$ BASIS FOR GIBBS SAMPLING

HMM LEARNING

FOR ESTIMATING $\theta = (\pi, A, B)$, π INITIAL DISTRIBUTION, A TRANSITION MATRIX, B CLASS CONDITIONAL DENSITIES

FULLY OBSERVED DATA

$z_{1:T}$ OBSERVED IN TRAINING. EASY. π, A CLOSED FORM. B LINE FITTING GENERATIVE CLASSIFIER. PSEUDOCOUNTS (MULTIPLY) / GAUSSIAN

NOT FULLY OBSERVED DATA

LINE FITTING MIXTURE MODEL. EM FOR MAP/MLE OR GRADIENT DESCENT. EM IS NOW CALLED BAUM-WELCH

• E STEP SAMPLE DATA LL DEPENDS ON EXPECTED COUNTS, COMPUTED W/ F-D ALGORITHM. $\gamma_{i,t}(j) = P(z_t = j | x_{1:t}, \theta)$, $\xi_{i,t}(j) = P(z_{t-1} = j, z_t = i | x_{1:t}, \theta)$ SMOOTHED MLE, BIC OR MAXIMUM

• M STEP FOR A, π IS JUST NORMALIZING THE EXP. COUNTS

• INITIALIZE W/ FULLY OBSERVED DATA, RANDOM INITIALIZE, **VITERBI TRAINING** \rightarrow APPROX PATH BEFORE W/ MOST PROBABLE PATH

BAYESIAN METHODS VARIATIONAL BAYES EM, MCMC

• DISCRIMINATIVE TRAINING: FOR HMM AS CLASS CONDITIONALS INSIDE GENERATIVE CLASSIFIER. MAXIMIZE CONDITIONAL LIKELIHOOD $\prod_{i=1}^N P(y_i | x_i, \theta)$ WITH GRADIENT METHODS

HMM MODEL SELECTION

• HOW MANY STATES? LINE CHOOSING NO. OF MIXTURES COMPONENTS. GRID SEARCH WITH X VAL LIKELIHOOD, BIC. MCMC. VARIATIONAL BAYES. INFINITE HMM

• STATE TRANSITION DIAGRAM TOPOLOGY? STRUCTURE LEARNING. LEARN SPARSE TRANSITION MATRIX. HEURISTIC METHODS. SPAT MEASURE. INFINITE HMM. MINIMUM ENTROPY PRIOR $P(A_{i,j}) \propto \exp(-H(A_{i,j}))$

HMM VARIANTS

VARIABLE DURATION HMM / SEMI-MARKOV

PROB OF REMAINING IN STATE i FOR D STEPS $P(z_t = i) \propto \exp(D \log A_{ii})$. GEOMETRIC DISTRIBUTION. VARSALISTIC.

IN SEMI-MARKOV WE DON'T ONLY CONDITION ON PAST STATE BUT WE ALSO NEED TO KNOW CURRENT PERFORMANCE TIME. THEY MODEL BATCH OF OBSERVATION AT ONE TIME. USEFUL UNBIASED FORM. USED IN GENETICS. STATES HAVE DURATION DISTRIBUTIONS. INFERENCE W/ MLEA VARIABLES OR MARKOVIZATION

• SEMI-MARKOV! EACH STATE REDUCED WITH N STATES. SAME OUTFLOW PROBABILITY. ANY PATH THROUGH SUBSTATES IS $P^{0-N}(1-P)^N$ NEGATIVE BINOMIAL. P IS SELF LOOP PROB. MANY BEHAVIORS BY ADJUSTING $P, 0$ IS FASTER.

HIERARCHICAL HMM

FOR DOMAINS WITH HIERARCHICAL STRUCTURE. GOOD, EFFICIENT INFERENCE. $O(T)$. A TRANSITION AT LEVEL k IS ALLOWED IFF CHAIN AT LEVEL $k-1$ REACHES ITS END STATE

INPUT/OUTPUT HMM

HANDLE INPUTS. CONDITIONAL DENSITY $P(y_{1:T}, z_{1:T} | u_{1:T}, \theta)$ UT INPUT / CTAL. TRANSITION MATRIX IS LOGISTIC REGRESSION MODEL W/ PARAMS DEPEND ON PREVIOUS STATE. HIDDEN VERSION OF MAX. ENTROPY MARKOV MODEL.

AUTOREGRESSIVE HMM

WHERE OBSERVATIONS ARE NOT C.I. GIVEN THE HIDDEN STATE $P(x_t | x_{t-1}, z_{1:t}, \theta) = N(x_t | w_j x_{t-1} + \mu_j, \Sigma_j)$ LINEAR MODEL WHERE PARAMS CHOSEN WRT CURRENT HIDDEN STATE. CAN EXTEND BY CONDITIONING ON LAST L OBSERVATIONS. COMBINE CHAIN ON HIDDEN STATES FOR LONG RANGE DEPENDENCIES AND ONE ON OBSERVATIONS FOR SHORT RANGE.

BURIED HMM

GENERALIZE AR-HMM BY ALLOWING DEPENDENCY STRUCTURE TO CHANGE ON HIDDEN STATE. OF OBSERVABLE NODES BAYESIAN MULTI-NET. MIXTURE OF NETWORKS

FACTORIAL HMM

DISTRIBUTED REPRESENTATION OF HIDDEN STATE VS SINGLE RANDOM VAR. STATES HAVE $z_{1:T} \in \{0, 1\}^T$ IS i TH BIT OF T TH HIDDEN STATE. EXACT ESTIMATION UNTRACTABLE BECAUSE XORING MANY

COUPLED HMM

STATE TRANSITIONS DEPEND ON STATES OF NEIGHBORING CHAINS. $P(z_t | z_{t-1}) = \prod P(z_t | z_{t-1})$ VERY EXPENSIVE $O(CK^4)$ INFLUENCE MODEL NO NEIGHBOUR RESTRICTIONS. CONVEX COMBINATIONS OF PAIRWISE TRANSITION MATRICES. $\alpha(T(u)^2)$

DYNAMIC BAYESIAN NETWORKS

BECAUSE NETWORK REPRESENTATION OF DYNAMICAL SYSTEM. GENERALLY ALL HMM ARE BUILT SPECIFICALLY ARE VERY DOMAIN SPECIFIC IRREGULAR GRAPHS DEFINE FIRST TIME-SLICE. STRUCTURE BETWEEN TIME-SLICES. CPDS. EXACT INFERENCE IS EXPENSIVE.