

HMM (continued)

Sriram Sankararaman

The instructor gratefully acknowledges Fei Sha, Ameet Talwalkar, Kai-Wei Chang, Eric Eaton, Andrew Moore and Jessica Wu whose slides are heavily used, and the many others who made their course material freely available online.

Administrivia

- Final exam on Tuesday March 21. Details posted on campuswire.
- Course evaluation open.

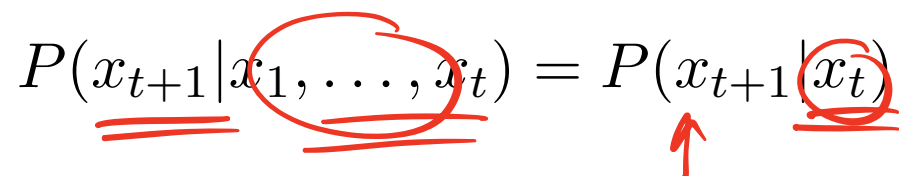
Outline

- 1 Review of last lecture
 - Markov chains
 - Hidden Markov models

Markov Process

Also known as Markov Chain or Markov model

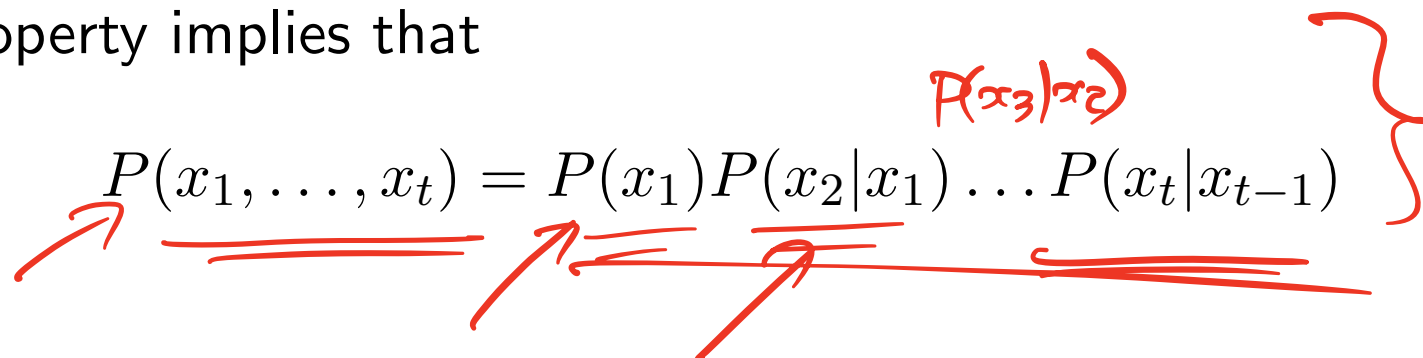
- For a **Markov process**, the next state depends on the current state :

$$\underline{P(x_{t+1} | x_1, \dots, x_t)} = P(x_{t+1} | \underline{x_t})$$


Conditioned on the present, the future is independent of the past

X_{t+1} is independent of $X_i, i < t$ given X_t

- This property implies that

$$\underline{P(x_1, \dots, x_t)} = P(x_1) \overset{P(x_2|x_1)}{P(x_2|x_1)} \dots P(x_t|x_{t-1})$$


Specifying a Markov chain

Initial probability

$$\pi_i = P(\underline{X_1 = i})$$

Transition probability

$$q_{ij} = P(\underline{X_{t+1} = i} | \underline{X_t = j})$$

Computing on Markov chains

$$x_1, \dots, x_T$$

$$P(x_T = i)$$

Compute $P(x_T = i)$

Assume uniform probability of starting in each of the states

$$\Rightarrow Q = \begin{pmatrix} 0.5 & 0.1 & 0.0 \\ 0.3 & 0.0 & 0.4 \\ 0.2 & 0.9 & 0.6 \end{pmatrix}$$

Q_{11}
 Q_{21}

What is the probability of observing a state $X_3 = 3$?

$$\Rightarrow P(\underline{x_1=1}, \underline{x_2=1}, \underline{x_3=3})$$
$$\quad +$$
$$\Rightarrow P(x_1=1, \underline{x_2=2}, \underline{x_3=3})$$
$$\quad +$$
$$\quad \vdots$$

Computing on Markov chains

Compute $P(x_T = i)$

$$x_1 \in \{1 \dots k\}$$

Assume uniform probability of starting in each of the states

$$Q = \begin{pmatrix} 0.5 & 0.1 & 0.0 \\ 0.3 & 0.0 & 0.4 \\ 0.2 & 0.9 & 0.6 \end{pmatrix}$$

What is the probability of observing a state $X_3 = 3$?

More generally, use the previous computation to sum over all paths of length t that end in the state i

$$P(X_t = i) = \sum_{x_1, x_2, \dots, x_{t-1}} P(\underline{x_1}, x_2, \dots, x_{t-1}, \underline{x_t = i})$$

Handwritten annotations: A red arrow points to the summation symbol \sum . A red arrow points to the term x_{t-1} in the sequence of states. A red circle highlights the sequence of states x_1, x_2, \dots, x_{t-1} . Below this sequence, a red circle highlights the expression $k \cdot k \dots k$. A red arrow points to the term $x_t = i$ in the probability expression. A red circle highlights the term $x_t = i$. Below the sequence of states, a red circle highlights the expression k .

Computational cost?

$O(K^{T-1})$: Exponential in T !

Computing on Markov chains

Compute $P(x_T = i)$

Define $p_t(i) = P(X_t = i)$.

$$p_1(i) = \pi_i$$

Assume we already know $p_{t-1}(i)$ for all i .

Use inductive definition to compute $p_t(i)$:

$$p_t(i) = P(X_t = i)$$

$$= \sum_j P(X_{t-1} = j, X_t = i)$$

$$= \sum_j P(X_{t-1} = j) P(X_t = i | X_{t-1} = j)$$

$$= \sum_j p_{t-1}(j) q_{ij}$$

$$O(K^{T-1}) \Rightarrow O(TK^2)$$

$$O(K \times K)$$

$$P_t(i) \quad O(K)$$

$$t = 1, 2, \dots, T$$

$$T-1, T$$

$$t=1$$

$$p_1(i) = P(X_1 = i) = \pi_i$$

$$p_{t-1}(i) = P(X_{t-1} = i)$$

$$\Rightarrow p_t(i) = P(X_t = i) ?$$

$$= \sum_j P(X_t = i, X_{t-1} = j)$$

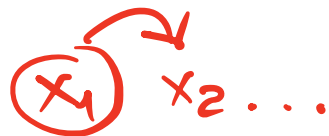
$$= \sum_j \frac{P(X_{t-1} = j)}{P(X_t = i | X_{t-1} = j)}$$

$$= \sum_j p_{t-1}(j) q_{ij}$$

$$O(K)$$

Computing on Markov chains

Compute $P(x_T = i)$



Define $p_t(i) = P(X_t = i)$.

$$p_1(i) = \pi_i$$

$$p_t(i) = \sum_j p_{t-1}(j) q_{ij}$$

Computational cost?

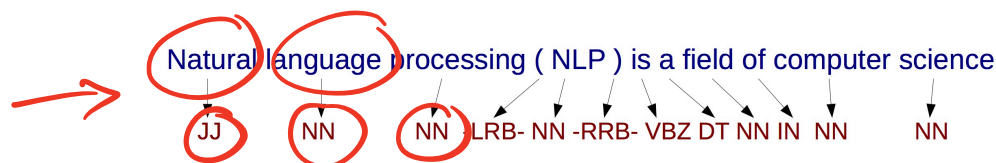
- Computing each $p_t(i)$ for a given t, i takes $O(K)$.
- Computing each $p_t(i)$ for a given t and all i takes $O(K^2)$.
- Computing each $p_t(i)$ for all $t \in \{1, \dots, T\}$ and all i takes $O(K^2(T - 1))$.

Example of Dynamic Programming

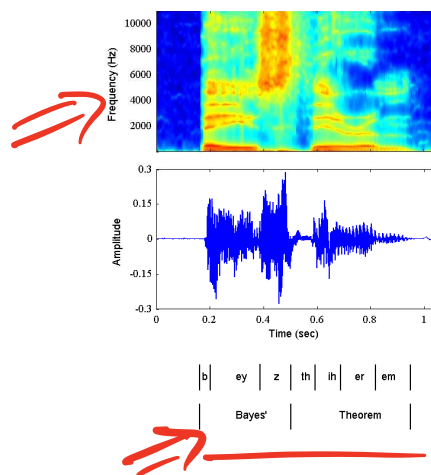
Hidden Markov models

 x_1, \dots, x_t

In many applications, we observe only a noisy or indirect measurement of the state X_t



In POS tagging, the state is the part-of-speech which we do not see. Instead, we observe words.



In speech recognition, the state is the word spoken but we only get to see the waveform.

Hidden Markov models

- Previously, in Markov chains, we directly observed the state X_t .
- Now, we observe a new random variable Y_t for each time t that is affected by the state X_t at that time t .
- Can think of Y_t as a noisy version of the true state.

Having observed (Y_1, \dots, Y_T) , we want to ask questions about X_1, \dots, X_T

Hidden Markov models

$$P(x_1, \dots, x_t) = P(x_1) P(x_2 | x_1) P(x_3 | x_2) \dots P(x_t | x_{t-1})$$

- We now define the set of observed states (also called **emission symbols**) $B = \{1, \dots, L\}$: the set of values that Y_t can take.
- Since X_t is not observed, X_t are called **hidden states**.

To link up the hidden and the observed states, we have emission probabilities

$$\Rightarrow P(\underline{Y_t} = \underline{b} | \underline{X_t} = \underline{k}) = \underline{e_k(b)}$$

Constraints: $\sum_b \underline{e_k(b)} = 1.$

Hidden Markov model


Hidden states: $\{1, \dots, K\}$

Observed states: $\{1, \dots, L\}$

Initial probability

$$\pi_i = P(\underline{\underline{X_1}} = i)$$

Transition probability

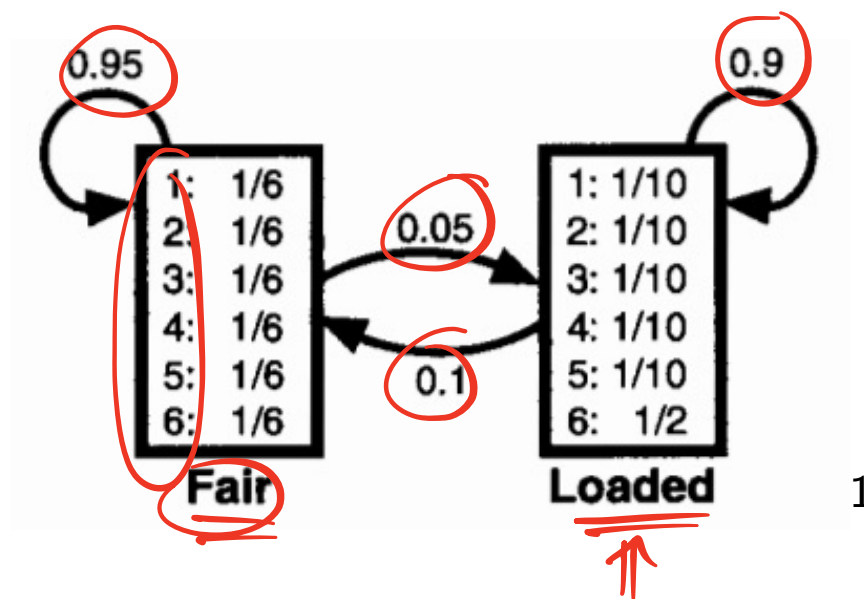
$$q_{ij} = P(X_{t+1} = \underline{i} | X_t = \underline{j})$$


Emission probability

$$e_i(b) = P(\underline{Y_t} = \underline{b} | X_t = i)$$

HMM: Example

The occasionally dishonest casino



- Hidden State: Is the casino using the fair or unfair die?
- Observed State: Roll of die ($\{1, \dots, 6\}$)

Querying an HMM

We observe $(Y_1, \dots, Y_5) = (2, 6, 6, 1, 3)$. Can we infer for which of the throws the casino used the unfair die ?

X_1, \dots, X_5

HMM

$$\begin{aligned} x &= x_1, \dots, x_T = x_{1:T} \\ y &= y_1, \dots, y_T = y_{1:T} \end{aligned}$$

The joint probability of the sequence of hidden states and the observed states

$$\begin{aligned} \Rightarrow P(y, x) &= P(y_{1:T}, x_{1:T}) \\ &= \underbrace{P(y_{1:T} | x_{1:T})}_{\text{Markov chain}} \underbrace{P(x_{1:T})}_{\text{Markov chain}} \\ &= \prod_{t=1}^T \underbrace{P(y_t | x_t)}_{\text{Markov chain}} \underbrace{P(x_{1:T})}_{\text{Markov chain}} \\ &= \prod_{t=1}^T \underbrace{e_{x_t}(y_t)}_{\text{Markov chain}} \underbrace{\pi_{x_1}}_{\text{Markov chain}} \underbrace{\prod_{t=1}^{T-1} q_{x_{t+1}x_t}}_{\text{Markov chain}} \end{aligned}$$

Easy to compute but not very useful because we don't know the hidden states.

HMM: The most probable path problem

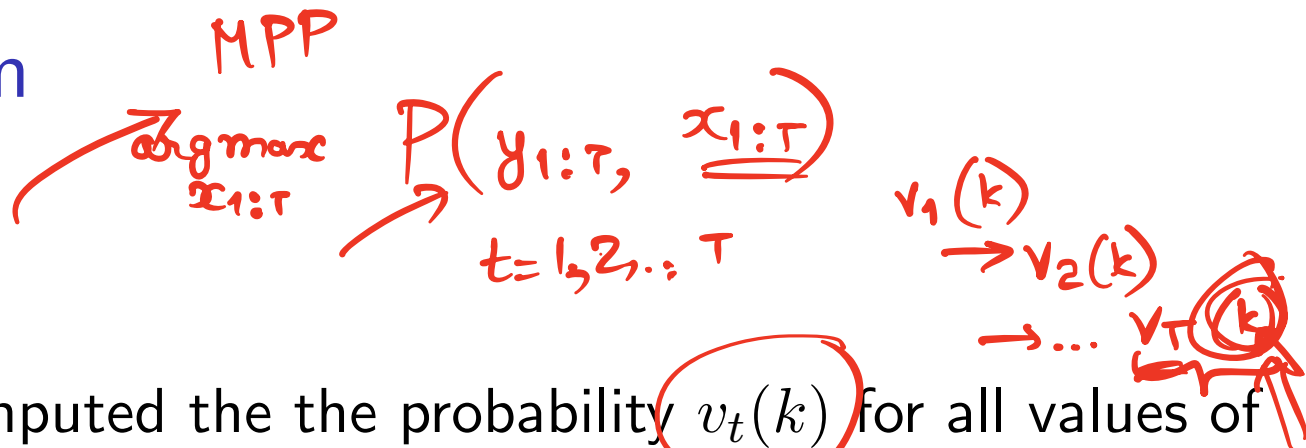
Given a sequence of observations (y_1, \dots, y_T) , what is the most probable sequence of hidden states (x_1, \dots, x_T) ?

$$\arg \max_{x_{1:T}} P(y_{1:T}, x_{1:T})$$

- Often called the **decoding** problem.
- One way to solve this problem is to search over all possible values of $(x_{1:T})$.
- There are exponentially many of them ($O(\underline{K}^T)$)
- Fortunately, it turns out there is an efficient dynamic programming algorithm to solve this problem.

The Viterbi algorithm

A recursive algorithm



- Suppose we have computed the the probability $v_t(k)$ for all values of $t \in \{1, \dots, T\}$ and $k \in \{1, \dots, K\}$:

The probability of the most probable path (MPP) for observations (y_1, \dots, y_t) (so that path has length t) that ends up in state k .

$$\Rightarrow v_t(k) = \max_{x_{1:t-1}} P(y_{1:t}, x_{1:t-1}, x_t = k)$$

$t=1 \rightarrow t=2 \rightarrow \dots \rightarrow t=T$

Why is this a useful quantity?

- If we look at $v_T(k)$, it tells us the probability of the MPP of length T that ends in state k .
- The answer to the MPP problem then is $\max_k v_T(k)$!

The Viterbi algorithm

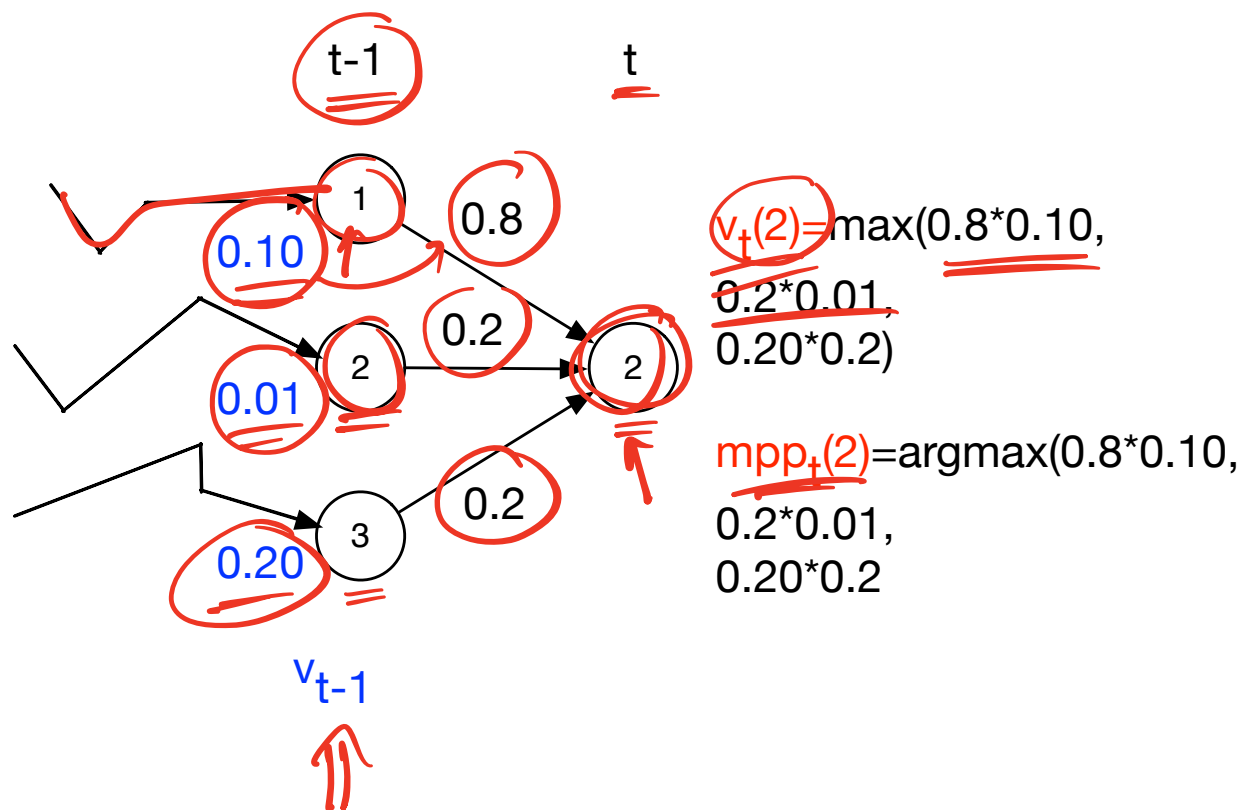
Can we compute $v_t(k)$ efficiently?

Assume we have computed $v_{t-1}(l)$.

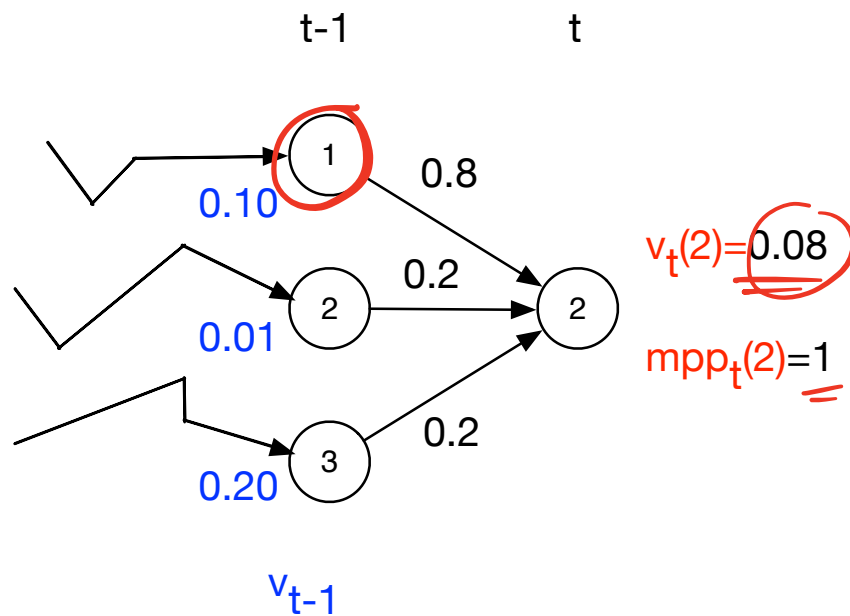
We have computed the probability of the MPP of length $t - 1$ that ends in state l for all values of l

How do we use this to compute the probability of MPP of length t that ends in state k ?

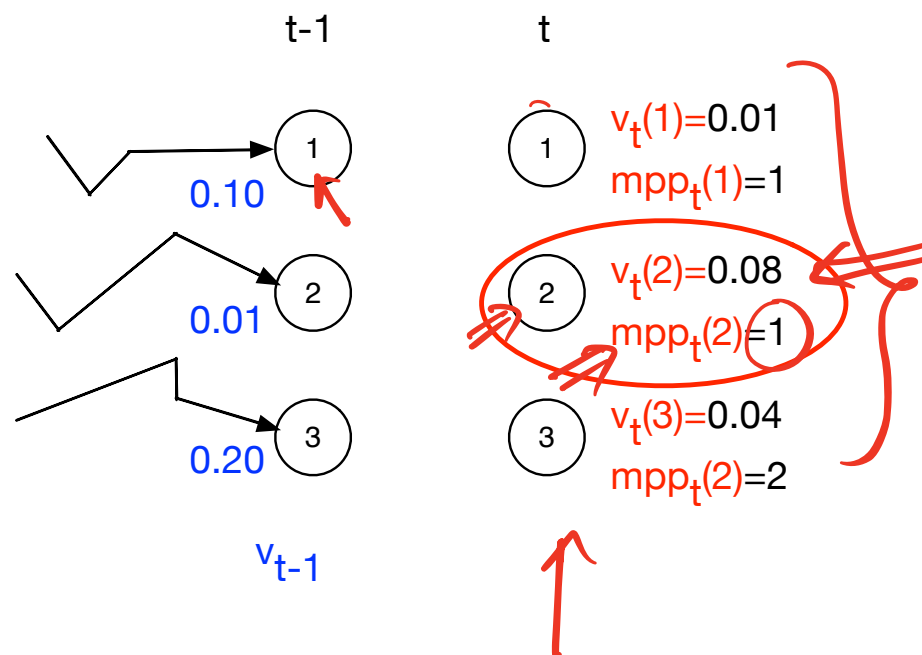
Viterbi example



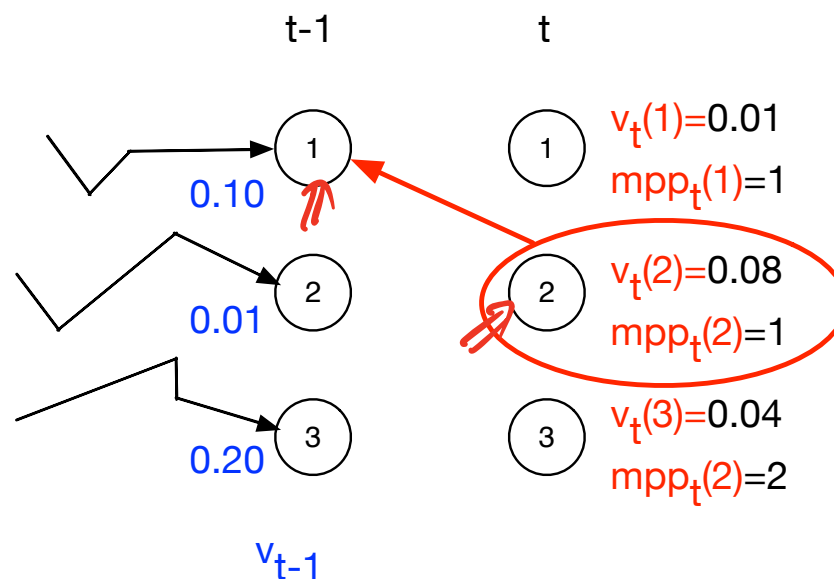
Viterbi example



Viterbi example



Viterbi example



MPP for observations upto time t (backwards):

$(2, 1, mpp_{t-1}(1), mpp_{t-2}(mpp_{t-1}(1)), \dots)$

In other words, state 2 at time t , state 1 at time $t - 1, \dots$

The Viterbi algorithm

Can we compute $v_t(k)$ efficiently?

Begin with $t = 1$

$$v_t(k) = v_1(\underline{k}) = \underline{P(y_1, x_1 = k)} \\ = \underline{P(y_1 | x_1 = k)} \underline{P(x_1 = k)}$$

$$\begin{aligned} v_1(k) &= P(y_1, x_1 = k) \\ &= P(y_1 | x_1 = k) P(x_1 = k) \\ &= \underline{e_k(y_1) \pi_k} \end{aligned}$$

The Viterbi algorithm

Can we compute $v_t(k)$ efficiently?

Assume we have computed $v_{t-1}(l)$.

We have computed the probability of the MPP of length $t - 1$ that ends in state l for all values of l

How do we use this to compute the probability of MPP of length t that ends in state k ?

The Viterbi algorithm

Can we compute $v_t(\underline{k})$ efficiently?

The most probable path with last two states $\underline{l, k}$ is the most probable path with state l at time $t - 1$ followed by a transition from state l to state k and emitting the observation at time t .

What is the probability of this path?

$$\begin{aligned} & \left[\underline{v_{t-1}(l)} P(\underline{X_t = k} | \underline{X_{t-1} = l}) P(\underline{y_t} | \underline{X_t = k}) \right] \\ & = \underline{v_{t-1}(l)} \underline{q_{lk}} \underline{e_k(y_t)} \end{aligned}$$

So the most probable path that ends in state k at time t is obtained by maximizing over all possible states l in the previous time $t - 1$.

$$\underline{v_t(k)} = \max_l \underline{v_{t-1}(l)} \underline{q_{kl}} \underline{e_t(y_t)}$$

$O(K^2T)$

Also keep a pointer to the state that lead to the current state

$$mpp_t(k) = l^*$$

$$l^* = \arg \max_l v_{t-1}(l) q_{kl} e_t(y_t)$$

The Viterbi algorithm

Can we compute $v_t(k)$ efficiently?

Continue till we compute $v_T(k)$, $k \in \{1, \dots, K\}$. Let:

$$k^* = \arg \max_k v_T(k)$$

To obtain the MPP, follow the pointers defined by $mpp_t(k)$.

The Viterbi algorithm

Can we compute $v_t(k)$ efficiently?

$$P(X_t = k | X_{t-1} = l)$$

$$v_t(k) = \max_l v_{t-1}(l) \underline{q_{kl}} e_t(y_t)$$

- The cost of computing this is $O(K)$ for a given k .
- The total cost of computing this is $O(K^2)$ for all k .
- Total cost of computing $v_T(k)$ is $O(\underline{TK^2})$.

Other HMM computations

$$P(X_2 = k | y_1, \dots, y_5) \quad T=5, y_1, \dots, y_5$$
$$\neq P(X_2 = k | y_1, y_2)$$

Given a sequence of observations (y_1, \dots, y_T) , what is the probability that state at time t is k ?

$$P(X_t = k | y_{1:T}) \neq P(X_t = k | y_{1:t})$$

Given a sequence of observations (y_1, \dots, y_T) , what is the probability of the observations ?

$$\underline{\underline{P(y_{1:T})}}$$

Can also be computed efficiently using dynamic programming.

Learning HMMs

We assume the parameters are known. Can we learn parameters from data?

Parameters of the HMM

$$\theta = (\underline{\pi}, \underline{Q}, \underline{E})$$

Here \mathbf{E} is the matrix of emission probabilities. $E_{kb} = e_k(b)$.

Given training data of observed states $\underline{y_{1:T}}$, find parameters θ that maximizes the log likelihood.

Learning HMMs

We assume the parameters are known. Can we learn parameters from data?

Our model contains observed and unobserved random variables and hence is **incomplete**.

- Observed: $\mathcal{D} = y_{1:T}$
- Unobserved (hidden): $x_{1:T}$

$$\begin{aligned}\hat{\theta} &= \arg \max_{\theta} \ell(\theta) \\ &= \arg \max_{\theta} \log P(y_{1:T} | \theta) \\ &= \arg \max_{\theta} \log \sum_{x_{1:T}} P(y_{1:T}, x_{1:T} | \theta)\end{aligned}$$

The objective function $\ell(\theta)$ is called the incomplete log likelihood.

We can optimize this function using the EM algorithm (we won't get into the details in this course)

Summary

HMM

- Allows us to model dependencies.
- Can perform computations efficiently using dynamic programming.
- Can learn the parameters using the EM algorithm.

Summary of the course

Types of learning problems

- Supervised, unsupervised and reinforcement
- Labeled vs unlabeled data
- Labeled: Supervised learning. Type of label: classification (categorical) vs regression (quantitative)
- Unlabeled: Unsupervised learning.

Summary of the course

Supervised learning

- Model/hypotheses
- Loss function
- Regularizer
- Algorithm to solve optimization problem

Summary of the course

Supervised learning

- Key goal is to pick hypothesis h that minimizes risk for some loss function:

$$\mathcal{R}[h(\mathbf{x})] = \sum_{\mathbf{x}, y} \ell(\underline{h(\mathbf{x})}, \underline{y}) \underline{p(\mathbf{x}, y)}$$

Difficulty: we don't know the data generating distribution $p(\mathbf{x}, y)$.

- Instead pick h that minimizes empirical risk (a.k.a training error)

$$\mathcal{R}^{\text{EMP}}[h(\mathbf{x})] = \frac{1}{N} \sum_n \ell(\underline{h(\mathbf{x}_n)}, \underline{y_n})$$

Summary of the course

Supervised learning

- Setup: given a training dataset $\{\mathbf{x}_n, y_n\}_{n=1}^N$, learn a function $h(\mathbf{x})$ to predict y given \mathbf{x} .
 - ▶ Choose hypothesis space/models.
 - ▶ Define a loss function.
 - ▶ Define a cost function (typically loss function evaluated over the training data + regularizer).
 - ▶ Algorithm to solve optimization problem.

Summary of the course

Supervised learning

- Hypotheses
 - ▶ Decision trees, Nearest neighbors
 - ▶ Linear models: $h(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$
 - ▶ Kernels to extend to non-linear functions.
 - ▶ Neural Networks: jointly learn ϕ and \mathbf{w} .
 - ▶ Ensembles as a way to combine classifiers.
- Loss functions
 - ▶ Squared loss: least squares for regression
 - ▶ 0-1 loss for binary classification and surrogates for 0-1 loss.
 - ▶ Logistic loss, Exponential loss, Hinge loss
- Main principles *logistic regression* *Adaboost* *SVM*
 - ▶ Many of these learning algorithms can be thought of as solving the problem of finding "good" parameters for some probabilistic model.
 - ▶ Principles for finding good parameters: Maximum likelihood, regularize likelihood
 - ▶ Generative vs discriminative models.

Summary of the course

Supervised learning

- Optimization
 - ▶ Convex vs non-convex optimization problems
 - ▶ Methods: gradient descent (batch vs stochastic)
 - ▶ Constrained optimization. Lagrange function. Primal vs dual formulations.
- Concepts
 - ▶ Training error vs generalization error
 - ▶ Overfitting vs underfitting
 - ▶ The role of inductive bias
- Practical issues
 - ▶ How to tune hyperparameters, how to estimate generalization error
 - ▶ Importance of train-validation-test setup and cross-validation

Unsupervised learning

- Finding structure in data.
- Dimensionality reduction, Clustering and mixture models, Modeling dependencies
- Clustering
 - ▶ K-means. Requires solving a non-convex problem
 - ▶ Can be viewed as a probabilistic model with hidden variable (GMM)
 - ▶ EM algorithm: iterative algorithm to estimate MLE
- PCA
 - ▶ Linear Dimensionality reduction
 - ▶ Finds projections that maximize variance, minimize reconstruction error
 - ▶ Obtained by computing the top eigenvectors
- Hidden Markov Models (HMM)
 - ▶ Model dependency among observations
 - ▶ Use dynamic programming to efficiently query the HMM

- Thank you for your participation!
- All the best for the exam!