



Machine Learning II: (Graphical Models, EM, Variational Inference)

Prepared by:

Prof. Dr. Visvanathan Ramesh

References and Sources: Nils Bertschinger (ML II lecture slides)
Simon Prince (Learning and Vision), HMM Topology Selection
(Stenger et al), Variational Inference (D. Blei)

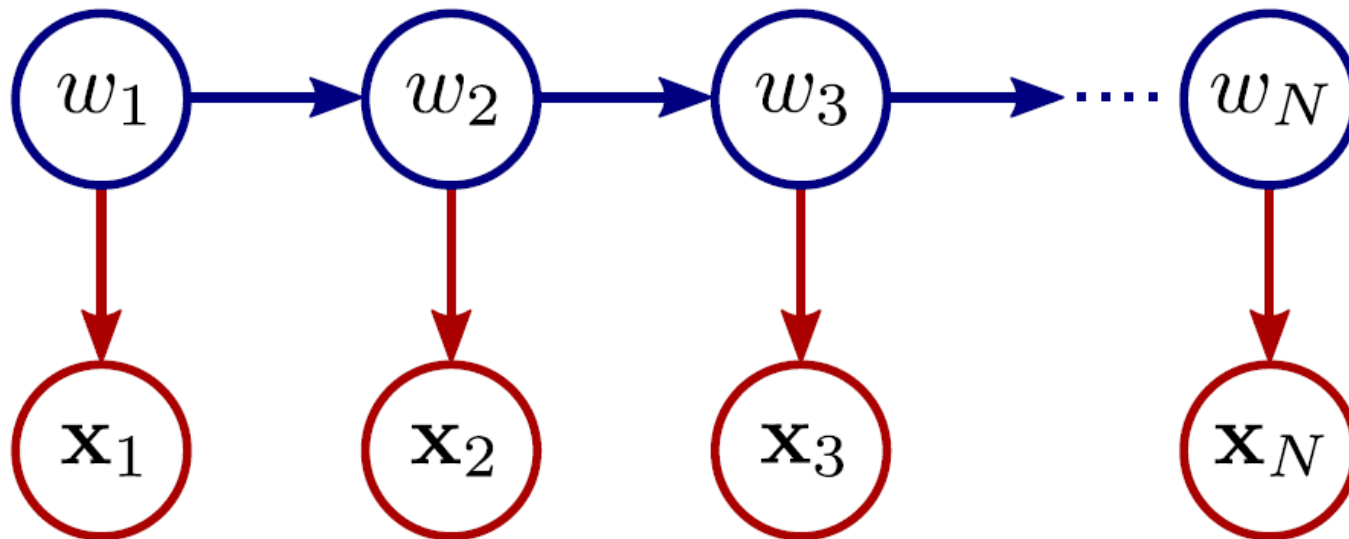


- **Recap – Past Lectures (Based on Simon Prince, Chapters 6, 7, 9, 10, 11)**
 - Regression, Classification, Application examples in Vision
 - Graphical Models and Inference
 - Graphical Models (directed, undirected)
 - Models for Chains and Trees
- Today's Lecture
 - Expectation Maximization Algorithm
 - HMM Model selection with application in Vision
 - Variational Bayes

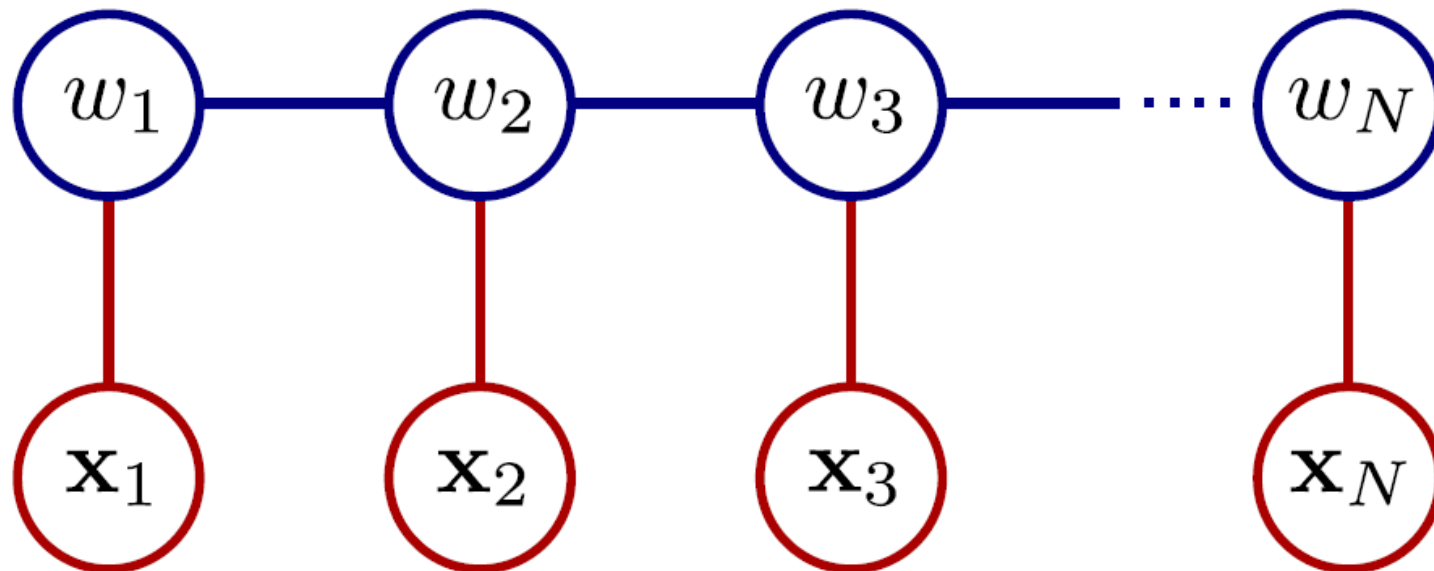


Brief Recap

Directed model for chains (Hidden Markov model)



Undirected model for chains





Supervised learning (where we know world states w_n) - relatively easy.

Unsupervised learning (where we do not know world states w_n) is more challenging. Use the EM algorithm:

- E-step – compute posterior marginals over states
- M-step – update model parameters

For the chain model (hidden Markov model) this is known as the Baum-Welch algorithm.



EM Algorithm



7

- Expectation Maximization Algorithm (Background)
- Convexity
- Jensen's Inequality
- EM Algorithm Formulation
- Short outline of Proofs
- Summary

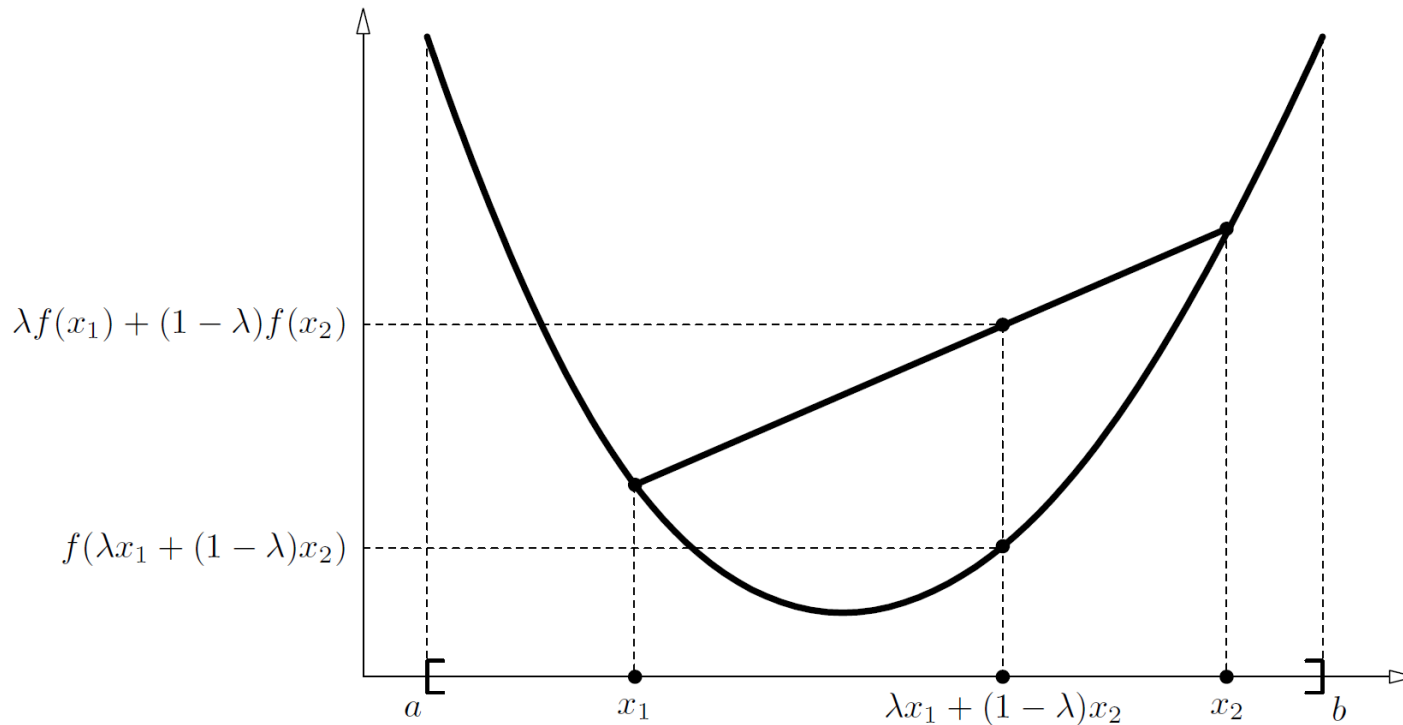


Figure 1: f is *convex* on $[a, b]$ if $f(\lambda x_1 + (1 - \lambda)x_2) \leq \lambda f(x_1) + (1 - \lambda)f(x_2)$
 $\forall x_1, x_2 \in [a, b], \lambda \in [0, 1]$.

Definitions



Definition 1 Let f be a real valued function defined on an interval $I = [a, b]$. f is said to be convex on I if $\forall x_1, x_2 \in I, \lambda \in [0, 1]$,

$$f(\lambda x_1 + (1 - \lambda)x_2) \leq \lambda f(x_1) + (1 - \lambda)f(x_2).$$

f is said to be strictly convex if the inequality is strict. Intuitively, this definition states that the function falls below (strictly convex) or is never above (convex) the straight line (the secant) from points $(x_1, f(x_1))$ to $(x_2, f(x_2))$. See Figure (1).

Definition 2 f is concave (strictly concave) if $-f$ is convex (strictly convex).

Theorem 1 If $f(x)$ is twice differentiable on $[a, b]$ and $f''(x) \geq 0$ on $[a, b]$ then $f(x)$ is convex on $[a, b]$.

Jensen's Inequality



Theorem 2 (Jensen's inequality) *Let f be a convex function defined on an interval I . If $x_1, x_2, \dots, x_n \in I$ and $\lambda_1, \lambda_2, \dots, \lambda_n \geq 0$ with $\sum_{i=1}^n \lambda_i = 1$,*

$$f\left(\sum_{i=1}^n \lambda_i x_i\right) \leq \sum_{i=1}^n \lambda_i f(x_i)$$

- Proof follows by Induction (Trivial for $n = 1$, $n=2 \rightarrow$ follows from convexity, demonstrate for $n+1$ assuming theorem true for n).

Since $\ln(x)$ is concave, we may apply Jensen's inequality to obtain the useful result,

$$\ln \sum_{i=1}^n \lambda_i x_i \geq \sum_{i=1}^n \lambda_i \ln(x_i). \quad (6)$$

This allows us to lower-bound a logarithm of a sum, a result that is used in the derivation of the EM algorithm.

Let \mathbf{X} be random vector which results from a parameterized family. We wish to find θ such that $\mathcal{P}(\mathbf{X}|\theta)$ is a maximum. This is known as the Maximum Likelihood (ML) estimate for θ . In order to estimate θ , it is typical to introduce the *log likelihood function* defined as,

$$L(\theta) = \ln \mathcal{P}(\mathbf{X}|\theta). \quad (7)$$

The likelihood function is considered to be a function of the parameter θ given the data \mathbf{X} . Since $\ln(x)$ is a strictly increasing function, the value of θ which maximizes $\mathcal{P}(\mathbf{X}|\theta)$ also maximizes $L(\theta)$.

The EM algorithm is an iterative procedure for maximizing $L(\theta)$. Assume that after the n^{th} iteration the current estimate for θ is given by θ_n . Since the objective is to maximize $L(\theta)$, we wish to compute an updated estimate θ such that,

$$L(\theta) > L(\theta_n) \quad (8)$$

Equivalently we want to maximize the difference,

$$L(\theta) - L(\theta_n) = \ln \mathcal{P}(\mathbf{X}|\theta) - \ln \mathcal{P}(\mathbf{X}|\theta_n). \quad (9)$$

EM Algorithm (Derivation)



$$L(\theta) - L(\theta_n) = \ln \left(\sum_{\mathbf{z}} \mathcal{P}(\mathbf{X}|\mathbf{z}, \theta) \mathcal{P}(\mathbf{z}|\theta) \right) - \ln \mathcal{P}(\mathbf{X}|\theta_n). \quad (11)$$

$$\begin{aligned} L(\theta) - L(\theta_n) &= \ln \left(\sum_{\mathbf{z}} \mathcal{P}(\mathbf{X}|\mathbf{z}, \theta) \mathcal{P}(\mathbf{z}|\theta) \right) - \ln \mathcal{P}(\mathbf{X}|\theta_n) \\ &= \ln \left(\sum_{\mathbf{z}} \mathcal{P}(\mathbf{X}|\mathbf{z}, \theta) \mathcal{P}(\mathbf{z}|\theta) \cdot \frac{\mathcal{P}(\mathbf{z}|\mathbf{X}, \theta_n)}{\mathcal{P}(\mathbf{z}|\mathbf{X}, \theta_n)} \right) - \ln \mathcal{P}(\mathbf{X}|\theta_n) \\ &= \ln \left(\sum_{\mathbf{z}} \mathcal{P}(\mathbf{z}|\mathbf{X}, \theta_n) \frac{\mathcal{P}(\mathbf{X}|\mathbf{z}, \theta) \mathcal{P}(\mathbf{z}|\theta)}{\mathcal{P}(\mathbf{z}|\mathbf{X}, \theta_n)} \right) - \ln \mathcal{P}(\mathbf{X}|\theta_n) \\ &\geq \sum_{\mathbf{z}} \mathcal{P}(\mathbf{z}|\mathbf{X}, \theta_n) \ln \left(\frac{\mathcal{P}(\mathbf{X}|\mathbf{z}, \theta) \mathcal{P}(\mathbf{z}|\theta)}{\mathcal{P}(\mathbf{z}|\mathbf{X}, \theta_n)} \right) - \ln \mathcal{P}(\mathbf{X}|\theta_n) \quad (12) \end{aligned}$$

$$= \sum_{\mathbf{z}} \mathcal{P}(\mathbf{z}|\mathbf{X}, \theta_n) \ln \left(\frac{\mathcal{P}(\mathbf{X}|\mathbf{z}, \theta) \mathcal{P}(\mathbf{z}|\theta)}{\mathcal{P}(\mathbf{z}|\mathbf{X}, \theta_n) \mathcal{P}(\mathbf{X}|\theta_n)} \right) \quad (13)$$

$$\triangleq \Delta(\theta|\theta_n). \quad (14)$$

$$l(\theta|\theta_n) \triangleq L(\theta_n) + \Delta(\theta|\theta_n)$$

EM Algorithm (Continued)



13

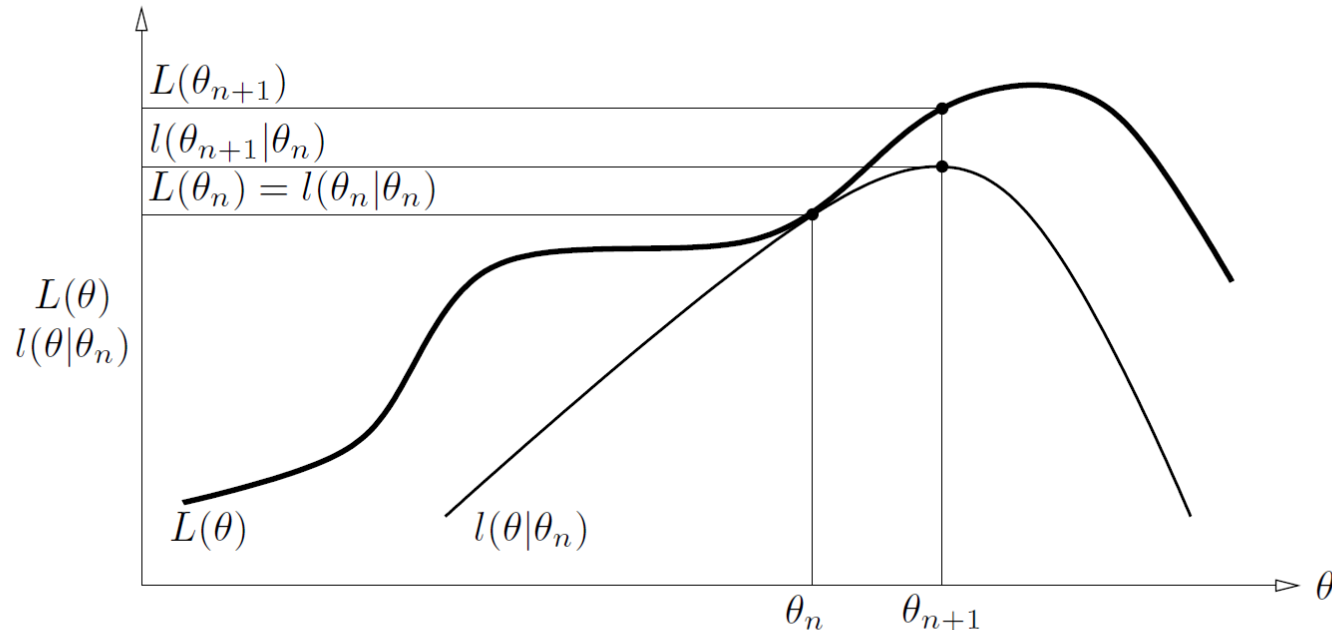


Figure 2: Graphical interpretation of a single iteration of the EM algorithm: The function $l(\theta|\theta_n)$ is bounded above by the likelihood function $L(\theta)$. The functions are equal at $\theta = \theta_n$. The EM algorithm chooses θ_{n+1} as the value of θ for which $l(\theta|\theta_n)$ is a maximum. Since $L(\theta) \geq l(\theta|\theta_n)$ increasing $l(\theta|\theta_n)$ ensures that the value of the likelihood function $L(\theta)$ is increased at each step.

EM Algorithm (Derivation)



14

$$\begin{aligned}\theta_{n+1} &= \arg \max_{\theta} \{l(\theta|\theta_n)\} \\ &= \arg \max_{\theta} \left\{ L(\theta_n) + \sum_{\mathbf{z}} \mathcal{P}(\mathbf{z}|\mathbf{X}, \theta_n) \ln \frac{\mathcal{P}(\mathbf{X}|\mathbf{z}, \theta)\mathcal{P}(\mathbf{z}|\theta)}{\mathcal{P}(\mathbf{X}|\theta_n)\mathcal{P}(\mathbf{z}|\mathbf{X}, \theta_n)} \right\} \\ &\quad \text{Now drop terms which are constant w.r.t. } \theta \\ &= \arg \max_{\theta} \left\{ \sum_{\mathbf{z}} \mathcal{P}(\mathbf{z}|\mathbf{X}, \theta_n) \ln \mathcal{P}(\mathbf{X}|\mathbf{z}, \theta)\mathcal{P}(\mathbf{z}|\theta) \right\} \\ &= \arg \max_{\theta} \left\{ \sum_{\mathbf{z}} \mathcal{P}(\mathbf{z}|\mathbf{X}, \theta_n) \ln \frac{\mathcal{P}(\mathbf{X}, \mathbf{z}, \theta)}{\mathcal{P}(\mathbf{z}, \theta)} \frac{\mathcal{P}(\mathbf{z}, \theta)}{\mathcal{P}(\theta)} \right\} \\ &= \arg \max_{\theta} \left\{ \sum_{\mathbf{z}} \mathcal{P}(\mathbf{z}|\mathbf{X}, \theta_n) \ln \mathcal{P}(\mathbf{X}, \mathbf{z}|\theta) \right\} \\ &= \arg \max_{\theta} \left\{ \mathbb{E}_{\mathbf{Z}|\mathbf{X}, \theta_n} \{ \ln \mathcal{P}(\mathbf{X}, \mathbf{z}|\theta) \} \right\}\end{aligned}\tag{17}$$

EM Algorithm - Summary



15

1. *E-step*: Determine the conditional expectation $E_{\mathbf{Z}|\mathbf{X},\theta_n} \{\ln \mathcal{P}(\mathbf{X}, \mathbf{z}|\theta)\}$
2. *M-step*: Maximize this expression with respect to θ .

Key Points:

- Iteratively converges to a local maximum
- Detailed Proof done later demonstrates convergence may not be only to Maxima (e.g. saddle points)
- Method is a unified principle for a number of estimation problems with Hidden variables and/or missing data.
- Several methods followed addressing computational speedups of algorithm



Adaptive Background Modeling Using a Hidden Markov Model

Slide Sources: Thesis by Bjoern Stenger (2000)

SCR Supervisor: Dr. Ramesh Visvanathan

Academic: Prof. J. Buhmann (Uni-Bonn)



Problem of Background adaptation in Video Sequences

Literature Review

An “offline” HMM State-Splitting Algorithm

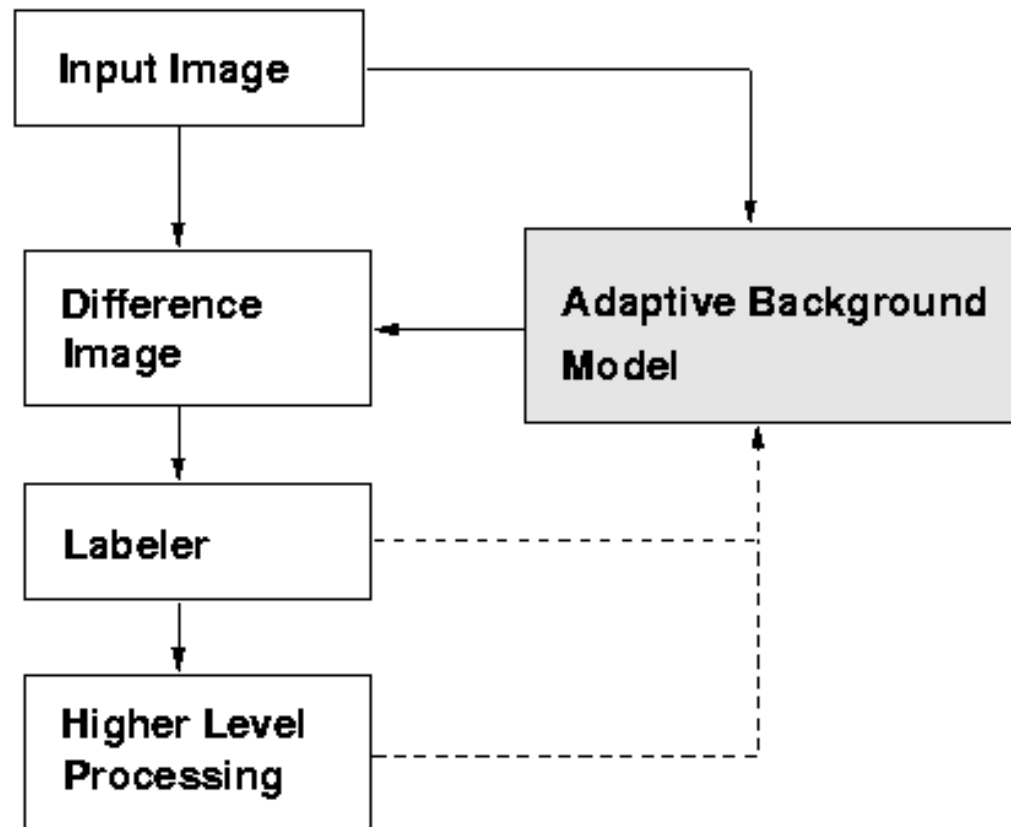
An “online” version for adapting HMM parameters

Example Background adaptation using HMMs

Demonstration

Summary

Module in Video Monitoring – Background Adaption





Statistics Estimation

- Mean and variances per pixel

Dynamic adaptation

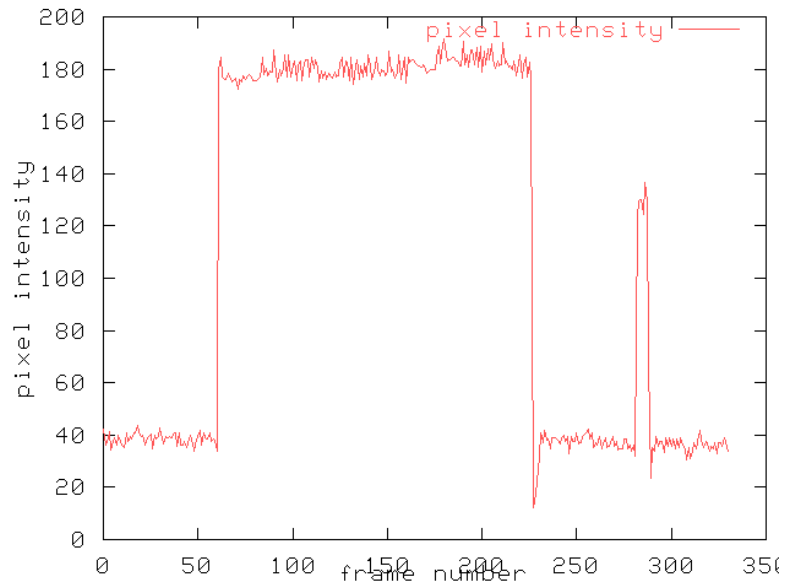
- Linear Prediction
- Kalman-Filter
- “Exponential Forgetting”

Problem: A single global model is not necessarily representative of the statistics. Need a mechanism to handle multiple states and to adapt the model online

Idea: Multi-State model for Global Changes



www.goethe-universitaet.de



State Learning from Video Sequences

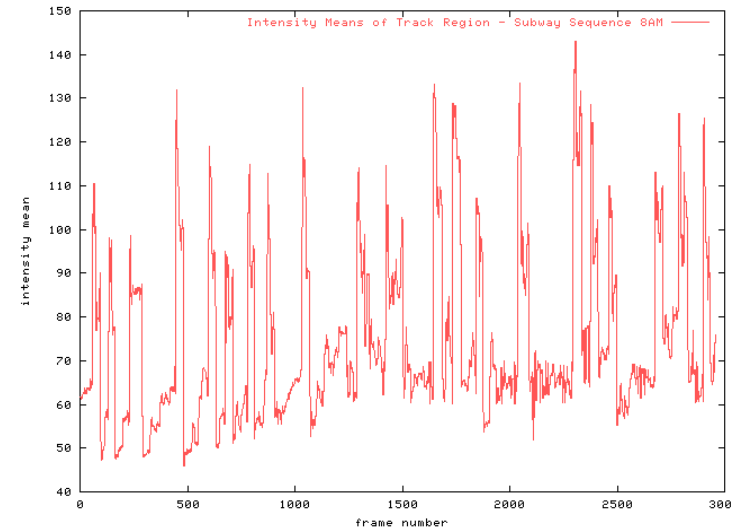


FIAS Frankfurt Institute
for Advanced Studies



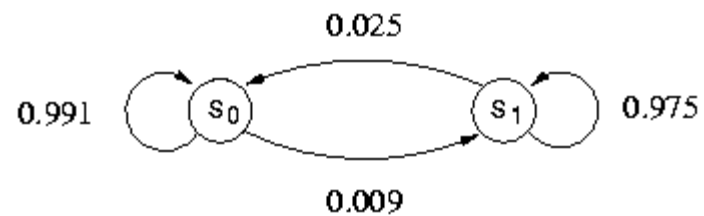
GOETHE
UNIVERSITÄT
FRANKFURT AM MAIN

Example Subway Sequence Observations & HMM Learned



No Train in Station

Train in Station



Relevance to Segmentation



30 Hidden Markov Model: Overview



Set of States S with Cardinality N , State at time index t denoted by s_t

Probability of being in a given state “ i ” at time $t=1$:

$$\pi(i) = P(s_t = i) \quad i \in S.$$

Transition Probability Matrix

$$a_{ij} = P(s_{t+1} = j \mid s_t = i) \quad i, j \in S.$$

Given a hidden state s_t , the observations are assumed to be sampled

from the Alphabet O

$$b_i(k) = P(o_t = k \mid s_t = i) \quad k \in O, i \in S$$

or from a continuous density:

$$b_i(x) = f_{\theta(i)}(x) \quad i \in S.$$



Generation of Sequence

Forward/Backward Algorithms

Useful to evaluate the probability $P(O | \theta)$

Viterbi Algorithm

Used to find the optimal state sequence given the observations and parameters (i.e.) $\operatorname{argmax}_S P(S | O, \theta)$

Baum-Welch Algorithm

To determine the maximum likelihood estimate (i.e) the theta that maximizes: $P(O | \theta)$

State Learning from Video Sequences (offline)



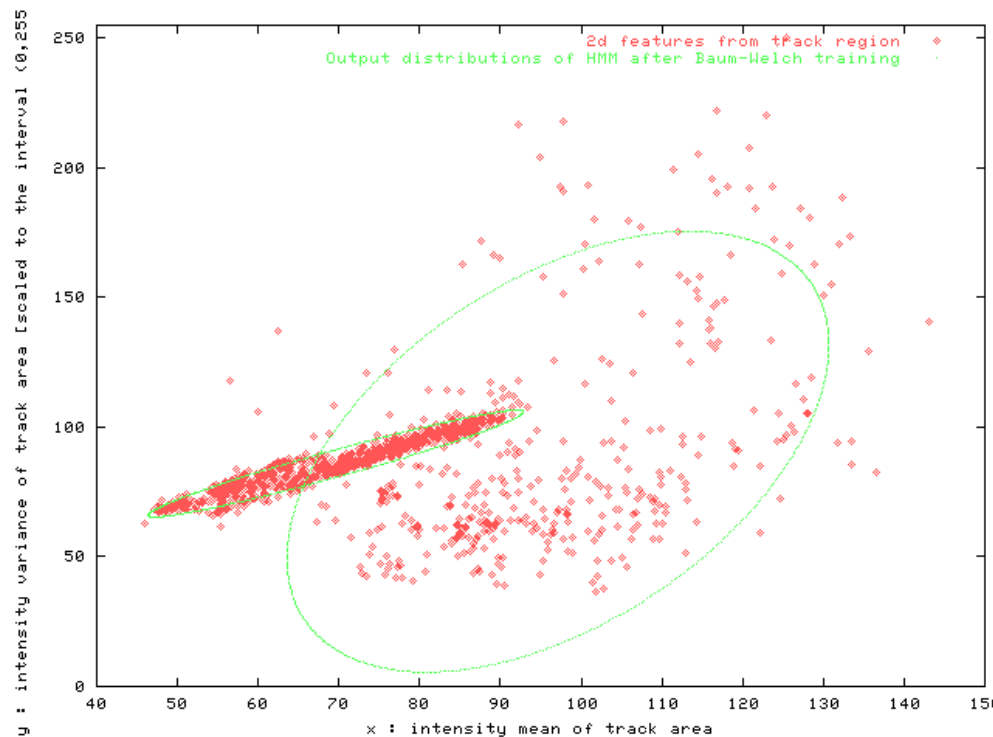
FIAS Frankfurt Institute
for Advanced Studies



GOETHE
UNIVERSITÄT
FRANKFURT AM MAIN

- Data: Subway monitoring sequence

Number of hidden states assumed to be given apriori: Train in station/ not in station

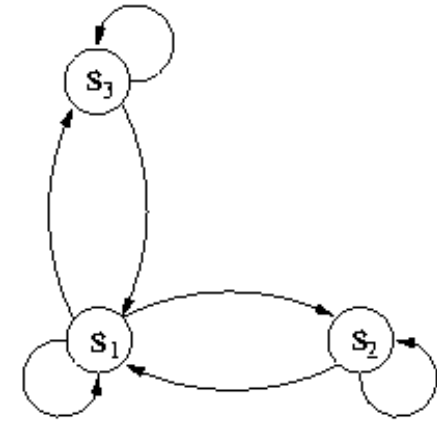
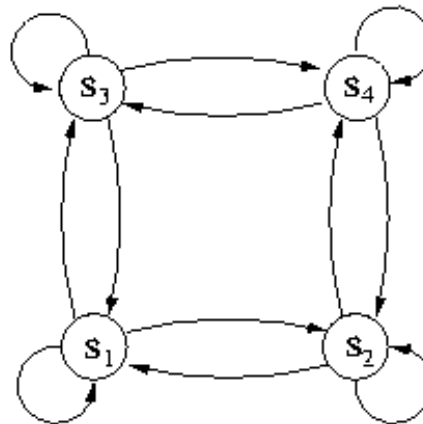
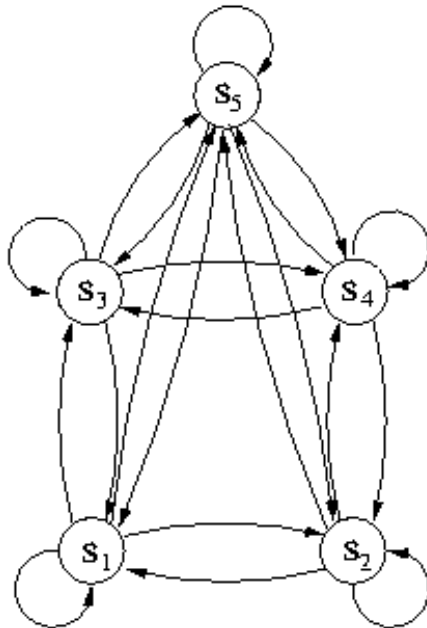


Correct Classification:

EM: 97.2%

HMM: 99.8%

State Splitting and Learning of Topology

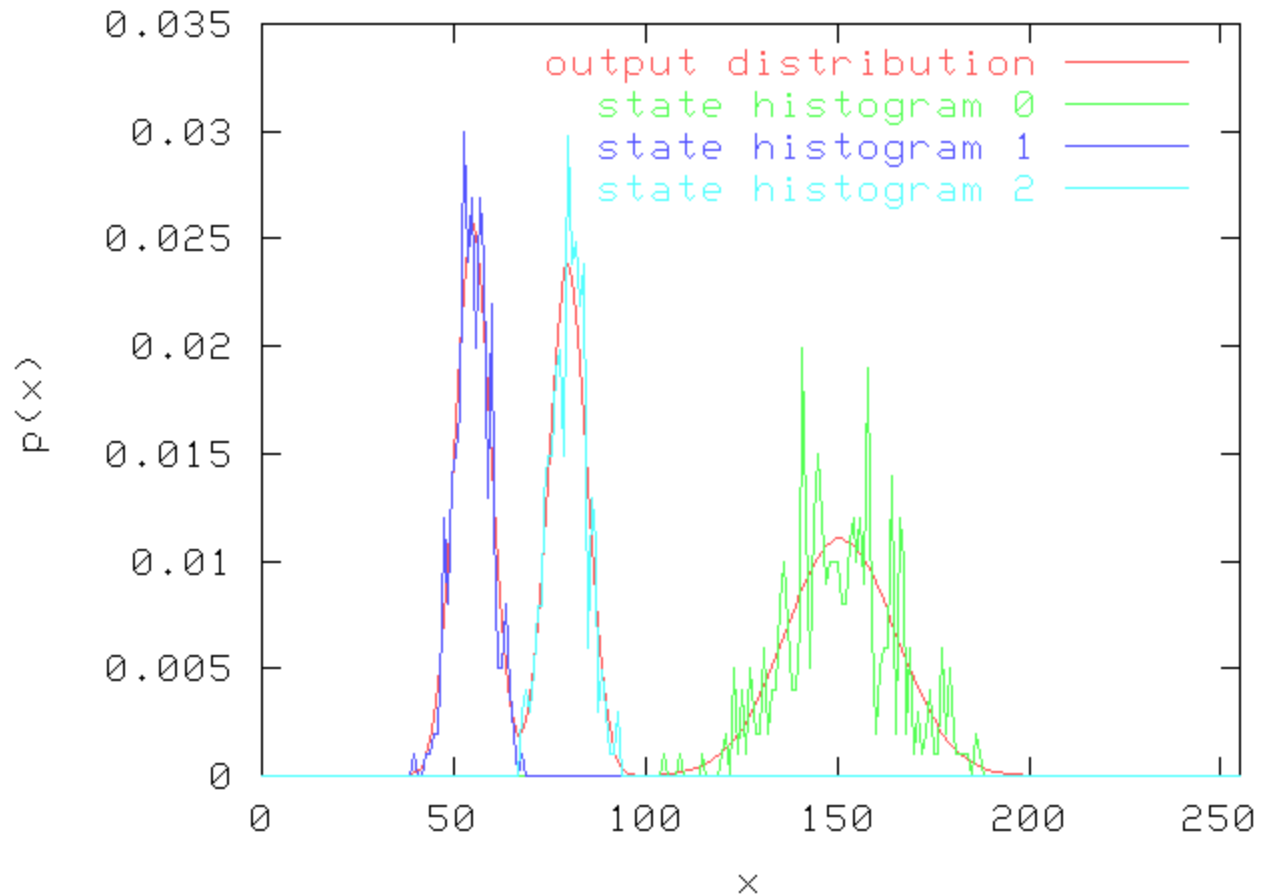


State-Splitting with the χ^2 -goodness of fit test



1. Initialize an HMM with $N=1$ State
2. Train the HMM using the Baum-Welch Algorithm
3. Generate State histograms using the Viterbi Algorithm.
4. Compute the χ^2 -statistic between the state distributions and the state histogram.
5. If at least one difference is statistically significant, split the state with the highest significance.
6. Go to step 2.

Simulation





Limitations of the χ^2 -test state-splitting method

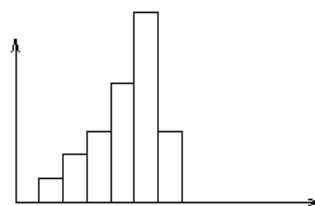
Unless the distributions are indeed Gaussian, the chi-squared test is not robust!

Merging may be needed to deal with false splits

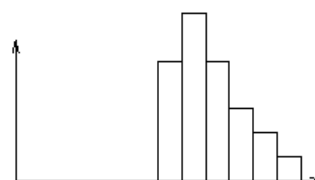
Generalization to higher dimensions is an issue.

Main problem: robustness to false splits.

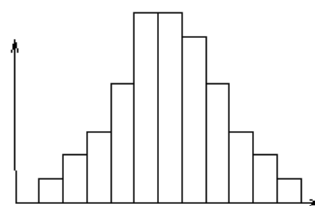
State Splitting (*temporal split*)



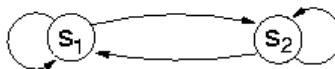
output
in state s_1



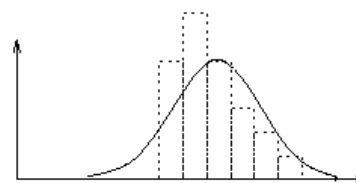
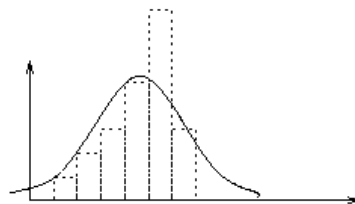
output
in state s_2



marginal
histogram

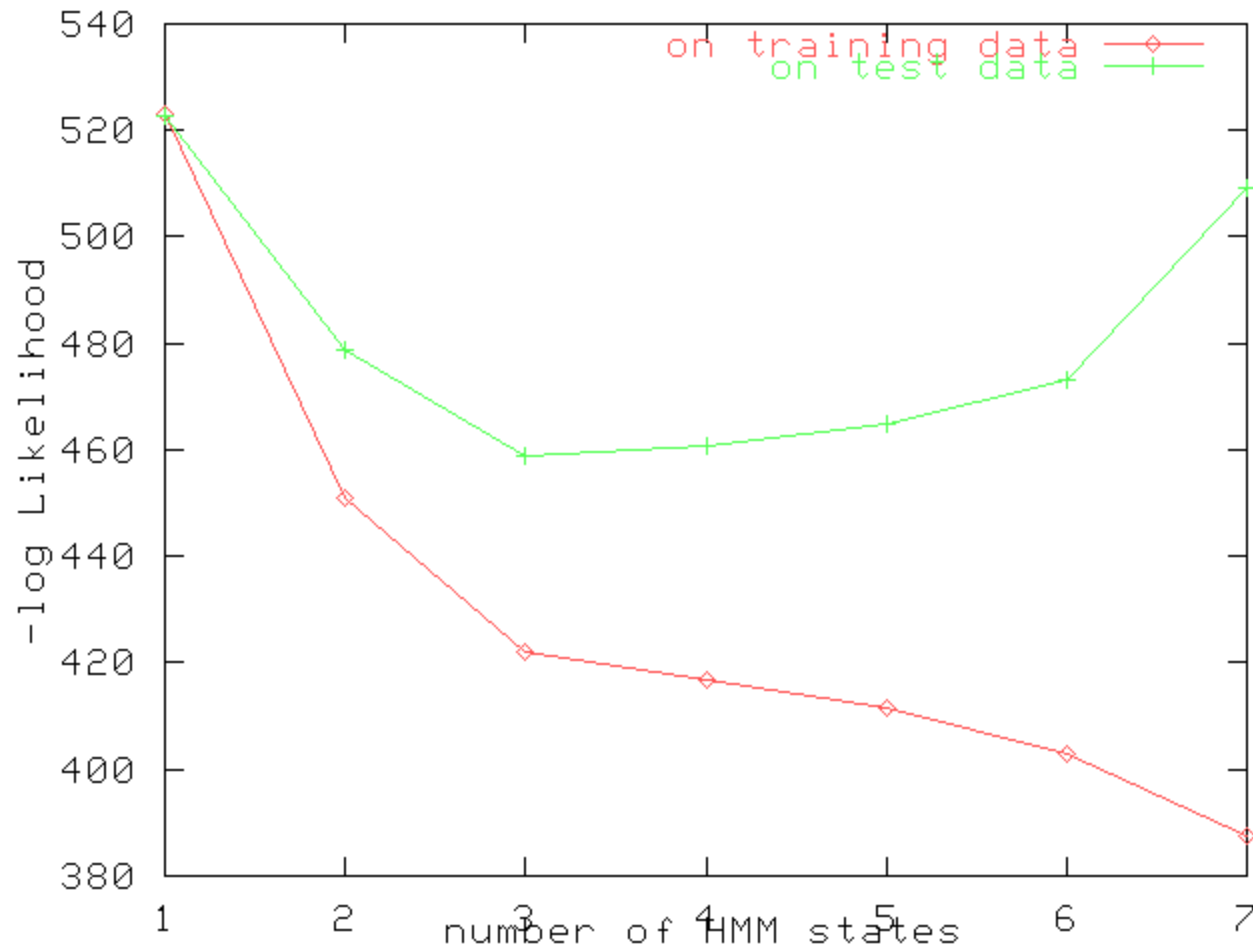


state output
densities after
the split





Overfitting Effect



State-Splitting with Cross-Validation



1. Initialize an HMM with $N=1$ states.
2. Train the model on the training set using the Baum-Welch Algorithm.
3. Compute the data likelihood of the test set, given this model.
4. If the likelihood on the test set decreases with the split, stop.
5. Select a split candidate with a goodness-of-fit test and split this state.
6. Go to step 2.



Model Selection

Given several models M_1, \dots, M_k

And data $x = x_1, \dots, x_n$

Cost function for Model Selection:

$$Q(\theta_k | x) = -\log L(\hat{\theta}_k / x) + C(n, \eta(k)),$$

$\eta(k)$: Number of free model parameters

$\hat{\theta}(k)$: Model parameters (estimated during MLE).

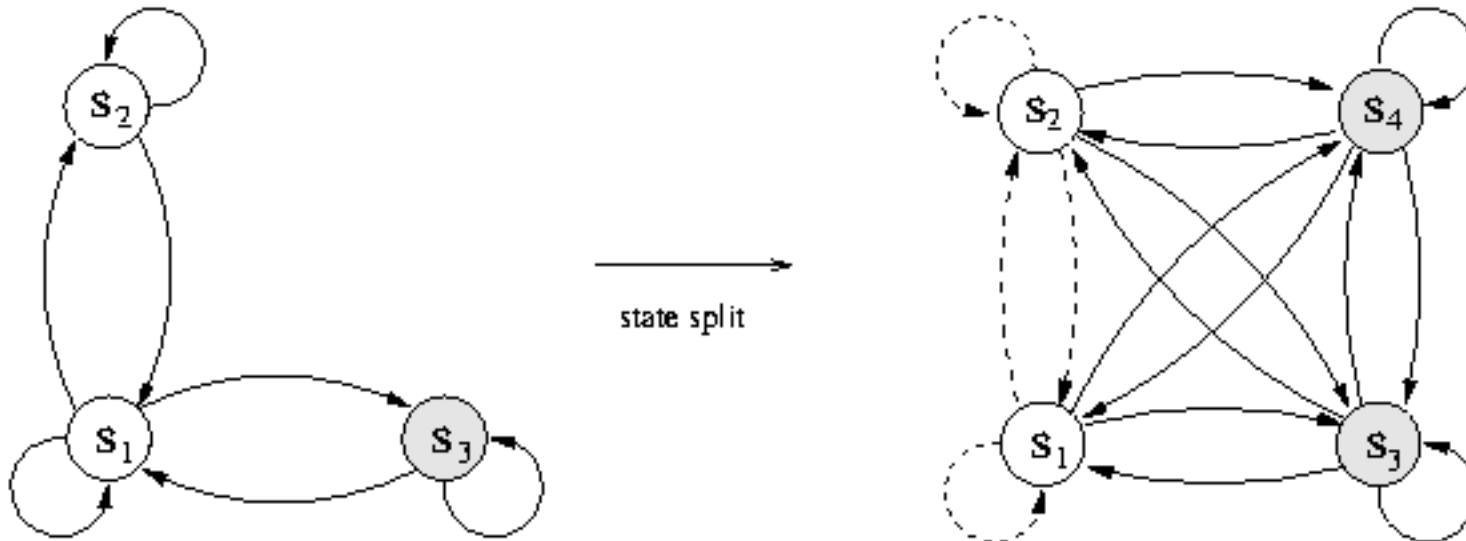
Penalty for model complexity: $C(n, \eta(k))$:

An Information Criterion (AIC): $\eta(k)$

Minimum Description Length (MDL): $\frac{1}{2} \eta(k) \log n$

State Splitting

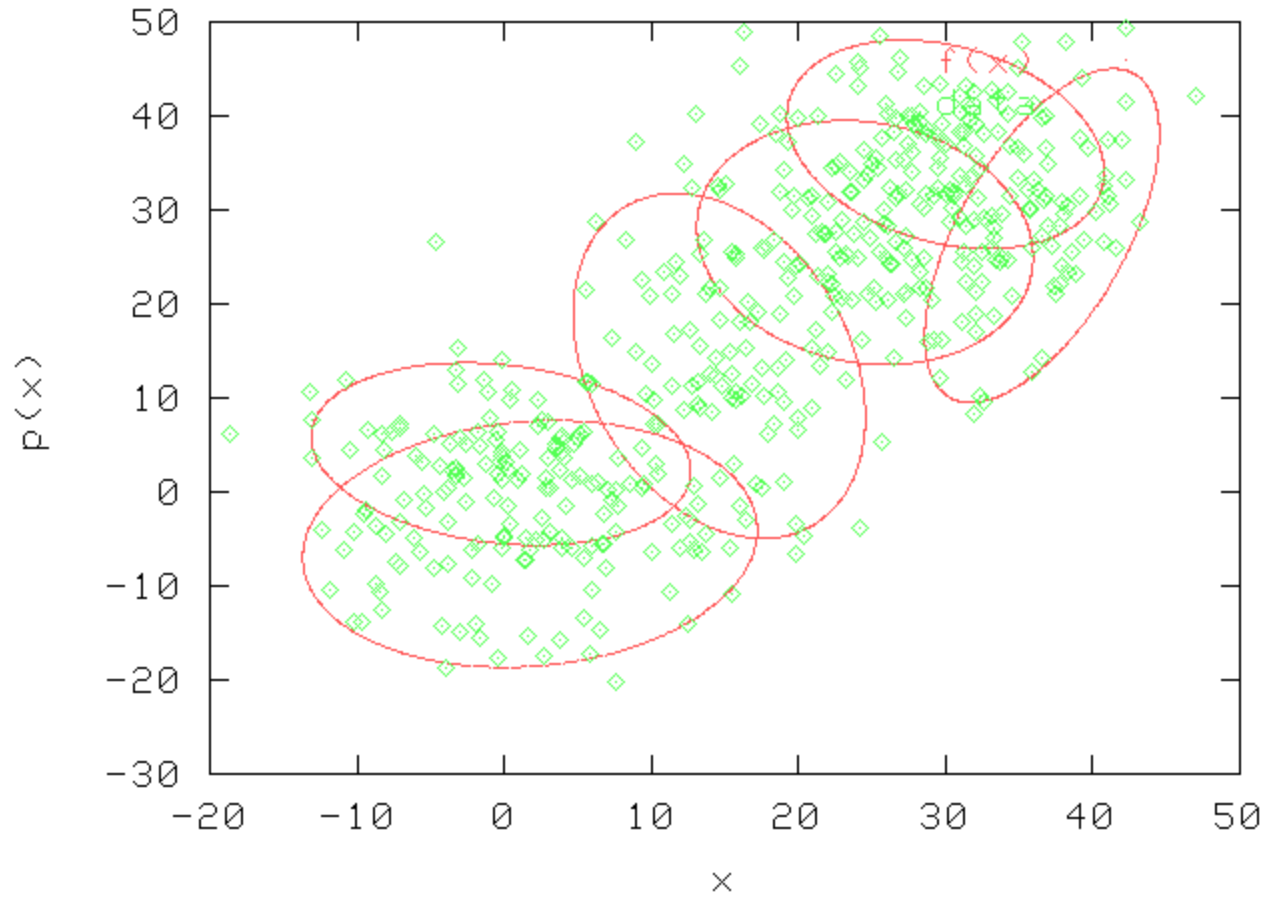
- Split the state which maximally increases the likelihood.
- Compute the likelihood increase only on a constrained subset of the model parameters



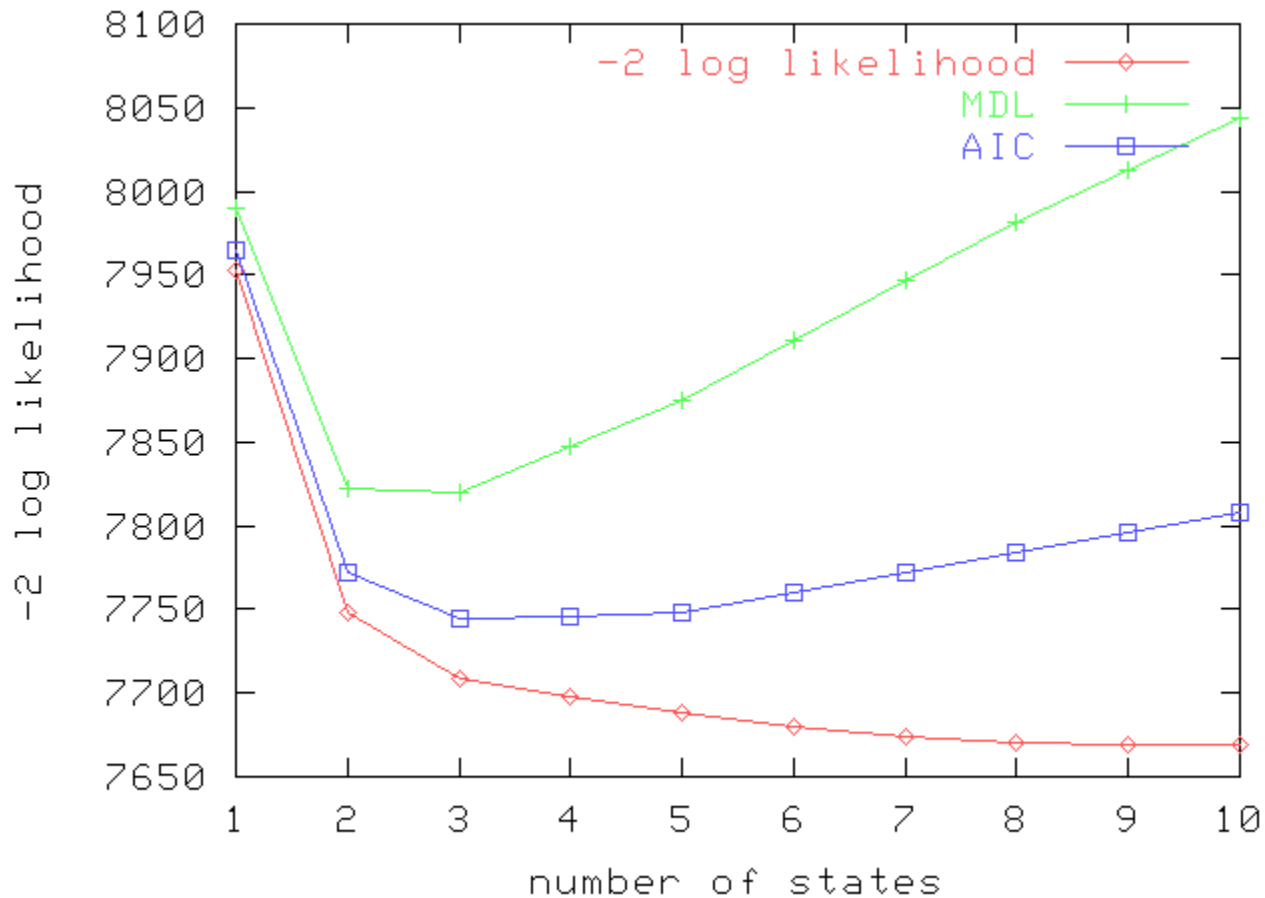


State-Splitting with MDL

1. Initialize a model with $N=1$ states and Train the HMM with the Baum Welch Algorithm.
2. Select the split which maximally increases the likelihood on a constrained subset of parameters.
3. Determine the likelihood increase for the complete model by training a model after state-splitting with Baum-Welch.
4. If the increased likelihood is greater than the MDL penalty difference, stop.
5. Go to step 2.



Cost function for Model Selection





Online Version

Need an online version to handle drifts in data over time. Devised an algorithm based on the online version of the Baum-Welch algorithm.

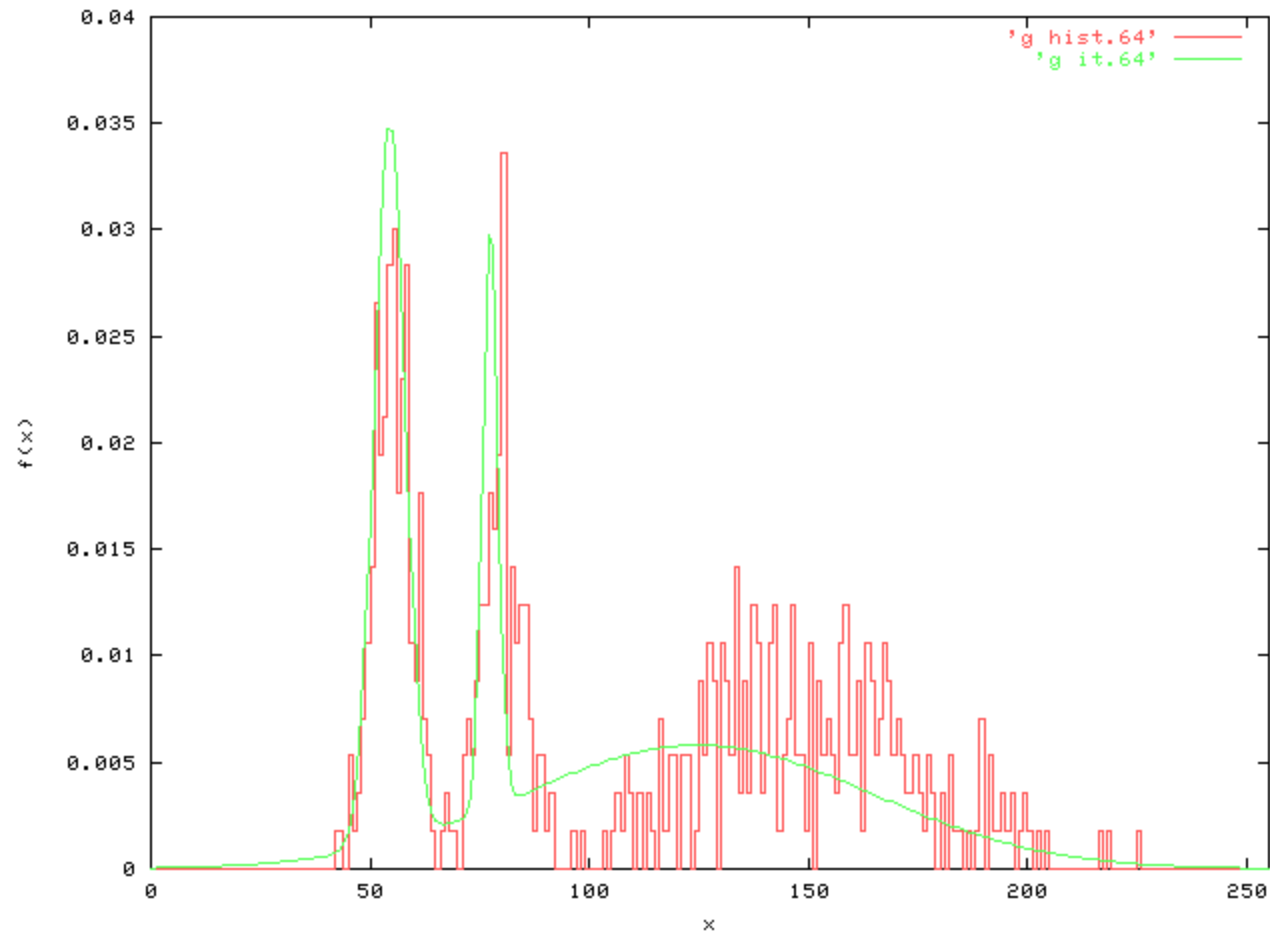
For each data sample x_i the model parameters are updated.

Algorithm Motivated by ‘Incremental EM-Versions’ of Baum-Welch (Nowlan, 1991)(Neal, Hinton).

Converges for stationary processes

Further Advantage: Little memory requirement

Example





Demo



HMM State Model



HMM State-Splitting Algorithms to select HMM topology and parameter estimation.

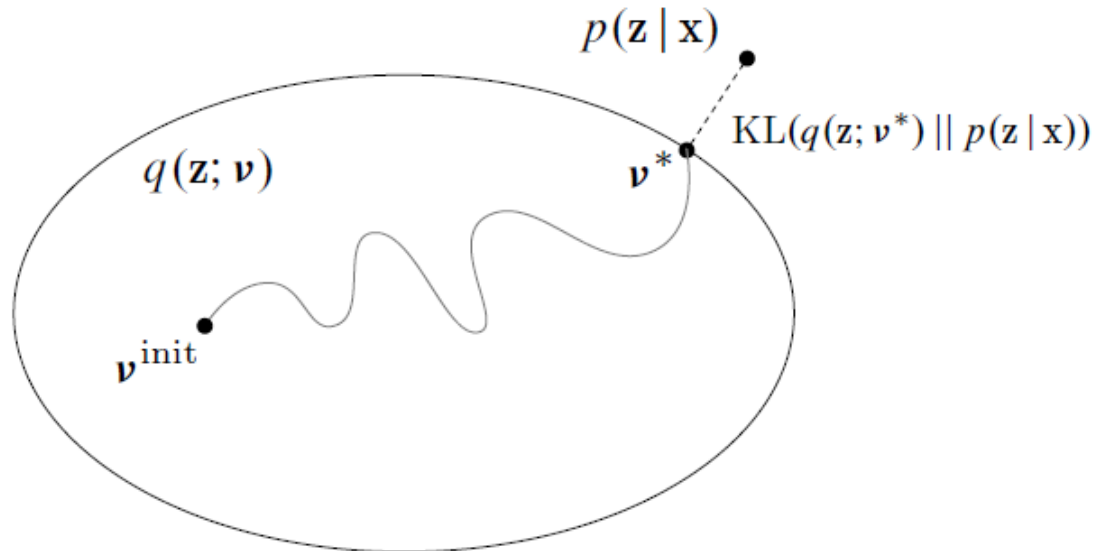
Online-Algorithm for Adapting model parameters.

Application for Global state changes in video analysis.



Variational Inference

Source: David Blei's Tutorial Article



- VI turns **inference** into **optimization**.
- Posit a **variational family** of distributions over the latent variables,

$$q(\mathbf{z}; \nu)$$

- Fit the **variational parameters** ν to be close (in KL) to the exact posterior.
(There are alternative divergences, which connect to algorithms like EP, BP, and others.)

- As usual, we will assume that $x = x_{1:n}$ are observations and $z = z_{1:m}$ are hidden variables. We assume additional parameters α that are fixed.
- Note we are general—the hidden variables might include the “parameters,” e.g., in a traditional inference setting. (In that case, α are the hyperparameters.)
- We are interested in the **posterior distribution**,

$$p(z | x, \alpha) = \frac{p(z, x | \alpha)}{\int_z p(z, x | \alpha)}. \quad (1)$$

- As we saw earlier, the posterior links the data and a model. It is used in all downstream analyses, such as for the predictive distribution.
- (Note: The problem of computing the posterior is an instance of a more general problem that variational inference solves.)

- **Cannot compute exact posterior distribution for many models**

- Consider the Bayesian mixture of Gaussians,

1. Draw $\mu_k \sim \mathcal{N}(0, \tau^2)$ for $k = 1 \dots K$.

2. For $i = 1 \dots n$:

- (a) Draw $z_i \sim \text{Mult}(\pi)$;

- (b) Draw $x_i \sim \mathcal{N}(\mu_{z_i}, \sigma^2)$.

- Suppressing the fixed parameters, the posterior distribution is

$$p(\mu_{1:K}, z_{1:n} \mid x_{1:n}) = \frac{\prod_{k=1}^K p(\mu_k) \prod_{i=1}^n p(z_i) p(x_i \mid z_i, \mu_{1:K})}{\int_{\mu_{1:K}} \sum_{z_{1:n}} \prod_{k=1}^K p(\mu_k) \prod_{i=1}^n p(z_i) p(x_i \mid z_i, \mu_{1:K})}. \quad (2)$$

- The numerator is easy to compute for any configuration of the hidden variables. The problem is the denominator.



- We return to the general $\{x, z\}$ notation.
- The main idea behind variational methods is to pick a family of distributions over the latent variables with its own **variational parameters**,

$$q(z_{1:m} \mid \nu). \tag{5}$$

- Then, find the setting of the parameters that makes q close to the posterior of interest.
- Use q with the fitted parameters as a proxy for the posterior, e.g., to form predictions about future data or to investigate the posterior distribution of the hidden variables.
- Typically, the true posterior is not in the variational family. (Draw the picture from Wainwright and Jordan, 2008.)



- We measure the closeness of the two distributions with Kullback-Leibler (KL) divergence.
- This comes from **information theory**, a field that has deep links to statistics and machine learning. (See the books “Information Theory and Statistics” by Kullback and “Information Theory, Inference, and Learning Algorithms” by MacKay.)

- The KL divergence for variational inference is

$$\text{KL}(q||p) = E_q \left[\log \frac{q(Z)}{p(Z|x)} \right]. \quad (6)$$

- Intuitively, there are three cases
 - If q is high and p is high then we are happy.
 - If q is high and p is low then we pay a price.
 - If q is low then we don't care (because of the expectation).

Evidence Lower Bound (ELBO)



- We actually can't minimize the KL divergence exactly, but we can minimize a function that is equal to it up to a constant. This is the **evidence lower bound (ELBO)**.
- We use Jensen's inequality on the log probability of the observations,

$$\log p(x) = \log \int_z p(x, z) \tag{8}$$

$$= \log \int_z p(x, z) \frac{q(z)}{q(z)} \tag{9}$$

$$= \log \left(\mathbb{E}_q \left[\frac{p(x, Z)}{q(Z)} \right] \right) \tag{10}$$

$$\geq \mathbb{E}_q[\log p(x, Z)] - \mathbb{E}_q[\log q(Z)]. \tag{11}$$

This is the ELBO. (Note: This is the same bound used in deriving the expectation-maximization algorithm.)

ELBO (continued)



- We choose a family of variational distributions (i.e., a parameterization of a distribution of the latent variables) such that the expectations are computable.
- Then, we maximize the ELBO to find the parameters that gives as tight a bound as possible on the marginal probability of x .
- Note that the second term is the entropy, another quantity from information theory.



- What does this have to do with the KL divergence to the posterior?
 - First, note that

$$p(z | x) = \frac{p(z, x)}{p(x)}. \quad (12)$$

- Now use this in the KL divergence,

$$\text{KL}(q(z) || p(z | x)) = \mathbb{E}_q \left[\log \frac{q(Z)}{p(Z | x)} \right] \quad (13)$$

$$= \mathbb{E}_q[\log q(Z)] - \mathbb{E}_q[\log p(Z | x)] \quad (14)$$

$$= \mathbb{E}_q[\log q(Z)] - \mathbb{E}_q[\log p(Z, x)] + \log p(x) \quad (15)$$

$$= -(\mathbb{E}_q[\log p(Z, x)] - \mathbb{E}_q[\log q(Z)]) + \log p(x) \quad (16)$$

This is the negative ELBO plus the log marginal probability of x .

ELBO (continued)



- Notice that $\log p(x)$ does not depend on q . So, as a function of the variational distribution, minimizing the KL divergence is the same as maximizing the ELBO.
- And, the difference between the ELBO and the KL divergence is the log normalizer—which is what the ELBO bounds.



- In mean field variational inference, we assume that the variational family **factorizes**,

$$q(z_1, \dots, z_m) = \prod_{j=1}^m q(z_j). \quad (17)$$

Each variable is independent. (We are suppressing the parameters ν_j .)

- This is more general than it initially appears—the hidden variables can be grouped and the distribution of each group factorizes.
- Typically, this family does not contain the true posterior because the hidden variables are dependent.
 - E.g., in the Gaussian mixture model all of the cluster assignments z_i are dependent on each other and the cluster locations $\mu_{1:K}$ given the data $x_{1:n}$.
 - These dependencies are often what makes the posterior difficult to work with.

Coordinate Ascent Inference for Optimization of ELBO



- We now turn to optimizing the ELBO for this factorized distribution.
- We will use **coordinate ascent inference**, iteratively optimizing each variational distribution holding the others fixed.



- First, recall the chain rule and use it to decompose the joint,

$$p(z_{1:m}, x_{1:n}) = p(x_{1:n}) \prod_{j=1}^m p(z_j \mid z_{1:(j-1)}, x_{1:n}) \quad (18)$$

Notice that the z variables can occur in any order in this chain. The indexing from 1 to m is arbitrary. (This will be important later.)

- Second, decompose the entropy of the variational distribution,

$$\mathbb{E}[\log q(z_{1:m})] = \sum_{j=1}^m \mathbb{E}_j[\log q(z_j)], \quad (19)$$

where \mathbb{E}_j denotes an expectation with respect to $q(z_j)$.

- Third, with these two facts, decompose the the ELBO,

$$\mathcal{L} = \log p(x_{1:n}) + \sum_{j=1}^m \mathbb{E}[\log p(z_j \mid z_{1:(j-1)}, x_{1:n})] - \mathbb{E}_j[\log q(z_j)]. \quad (20)$$



- Consider the ELBO as a function of $q(z_k)$.
 - Employ the chain rule with the variable z_k as the last variable in the list.
 - This leads to the objective function

$$\mathcal{L} = \mathbb{E}[\log p(z_k | z_{-k}, x)] - \mathbb{E}_j[\log q(z_k)] + \text{const.} \quad (21)$$

- Write this objective as a function of $q(z_k)$,

$$\mathcal{L}_k = \int q(z_k) \mathbb{E}_{-k}[\log p(z_k | z_{-k}, x)] dz_k - \int q(z_k) \log q(z_k) dz_k. \quad (22)$$

Coordinate Ascent Inference – Derivation



- Take the derivative with respect to $q(z_k)$

$$\frac{d\mathcal{L}_j}{dq(z_k)} = \mathbb{E}_{-k}[\log p(z_k | z_{-k}, x)] - \log q(z_k) - 1 = 0 \quad (23)$$

- This (and Lagrange multipliers) leads to the coordinate ascent update for $q(z_k)$

$$q^*(z_k) \propto \exp\{\mathbb{E}_{-k}[\log p(z_k | Z_{-k}, x)]\} \quad (24)$$

- But the denominator of the posterior does not depend on z_j , so

$$q^*(z_k) \propto \exp\{\mathbb{E}_{-k}[\log p(z_k, Z_{-k}, x)]\} \quad (25)$$

- Either of these perspectives might be helpful in deriving variational inference algorithms.



- There is a strong relationship between this algorithm and Gibbs sampling.
 - In Gibbs sampling we sample from the conditional.
 - In coordinate ascent variational inference, we iteratively set each factor to

$$\text{distribution of } z_k \propto \exp\{\mathbb{E}[\log(\text{conditional})]\}. \quad (26)$$

- Easy example: Multinomial conditionals
 - Suppose the conditional is multinomial

$$p(z_j \mid z_{-j}, x_{1:n}) := \pi(z_{-j}, x_{1:n}) \quad (27)$$

- Then the optimal $q(z_j)$ is also a multinomial,

$$q^*(z_j) \propto \exp\{\mathbb{E}[\log \pi(z_{-j}, x)]\} \quad (28)$$

Exponential Family Conditionals



- Suppose each conditional is in the exponential family

$$p(z_j | z_{-j}, x) = h(z_j) \exp\{\eta(z_{-j}, x)^\top t(z_j) - a(\eta(z_{-j}, x))\} \quad (29)$$

- This describes *a lot* of complicated models
 - Bayesian mixtures of exponential families with conjugate priors
 - Switching Kalman filters
 - Hierarchical HMMs
 - Mixed-membership models of exponential families
 - Factorial mixtures/HMMs of exponential families
 - Bayesian linear regression
- Notice that any model containing conjugate pairs and multinomials has this property.

Examples of Exponential Family



Gaussian	$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\ x-\mu\ ^2/(2\sigma^2)}$	$x \in \mathbb{R}$
Bernoulli	$p(x) = \alpha^x (1 - \alpha)^{1-x}$	$x \in \{0, 1\}$
Binomial	$p(x) = \binom{n}{x} \alpha^x (1 - \alpha)^{n-x}$	$x \in \{0, 1, 2, \dots, n\}$
Multinomial	$p(x) = \frac{n!}{x_1!x_2!\dots x_n!} \prod_{i=1}^n \alpha_i^{x_i}$	$x_i \in \{0, 1, 2, \dots, n\}, \sum_i x_i = n$
Exponential	$p(x) = \lambda e^{-\lambda x}$	$x \in \mathbb{R}^+$
Poisson	$p(x) = \frac{e^{-\lambda}}{x!} \lambda^x$	$x \in \{0, 1, 2, \dots\}$
Dirichlet	$p(x) = \frac{\Gamma(\sum_i \alpha_i)}{\prod_i \Gamma(\alpha_i)} \prod_i x_i^{\alpha_i-1}$	$x_i \in [0, 1], \sum_i x_i = 1$



$$\begin{aligned} p(x) &= \frac{1}{\sqrt{2\pi\sigma^2}} e^{-(x-\mu)^2/(2\sigma^2)} \\ &= \frac{1}{\sqrt{2\pi}} \exp\left(-\log \sigma - \frac{x^2}{2\sigma^2} + \frac{\mu x}{\sigma^2} - \frac{\mu^2}{2\sigma^2}\right) \\ &= \underbrace{\frac{1}{\sqrt{2\pi}}}_{h(x)} \exp\left(\theta^\top T(x) - \underbrace{\log \sigma - \mu^2/(2\sigma^2)}_{A(\theta)}\right) \end{aligned}$$

where

$$T(x) = \begin{pmatrix} x \\ x^2 \end{pmatrix} \quad \theta = \begin{pmatrix} \mu/\sigma^2 \\ -1/(2\sigma^2) \end{pmatrix} \quad \begin{aligned} A(\theta) &= \frac{\mu^2}{2\sigma^2} + \log \sigma \\ &= -\frac{[\theta]_1^2}{4[\theta]_2} - \frac{1}{2} \log(-2[\theta]_2) \end{aligned}$$

Exponential Family Case



- Mean field variational inference is straightforward

- Compute the log of the conditional

$$\log p(z_j | z_{-j}, x) = \log h(z_j) + \eta(z_{-j}, x)^\top t(z_j) - a(\eta(z_{-j}, x)) \quad (30)$$

- Compute the expectation with respect to $q(z_{-j})$

$$\mathbb{E}[\log p(z_j | z_{-j}, x)] = \log h(z_j) + \mathbb{E}[\eta(z_{-j}, x)]^\top t(z_j) - \mathbb{E}[a(\eta(z_{-j}, x))] \quad (31)$$

- Noting that the last term does not depend on q_j , this means that

$$q^*(z_j) \propto h(z_j) \exp\{\mathbb{E}[\eta(z_{-j}, x)]^\top t(z_j)\} \quad (32)$$

and the normalizing constant is $a(\mathbb{E}[\eta(z_{-j}, x)])$.

- So, the optimal $q(z_j)$ is in the same exponential family as the conditional.



- Coordinate ascent algorithm

- Give each hidden variable a variational parameter ν_j , and put each one in the same exponential family as its model conditional,

$$q(z_{1:m} | \nu) = \prod_{j=1}^m q(z_j | \nu_j) \quad (33)$$

- The coordinate ascent algorithm iteratively sets each natural variational parameter ν_j equal to the expectation of the natural conditional parameter for variable z_j given all the other variables and the observations,

$$\nu_j^* = \mathbb{E}[\eta(z_{-j}, x)]. \quad (34)$$



- Let's go back to the Bayesian mixture of Gaussians. For simplicity, assume that the data generating variance is one.
- The latent variables are cluster assignments z_i and cluster means μ_k .
- The mean field family is

$$q(\mu_{1:K}, z_{1:n}) = \prod_{k=1}^K q(\mu_k | \tilde{\mu}_k, \tilde{\sigma}_k^2) \prod_{i=1}^n q(z_i | \phi_i), \quad (35)$$

where $(\tilde{\mu}_k, \tilde{\sigma}_k)$ are Gaussian parameters and ϕ_i are multinomial parameters (i.e., positive K -vectors that sum to one.)

Bayesian Mixtures of Gaussians

- Variational Updates



the coordinate update for $q(z_i)$ is

$$q^*(z_i = k) \propto \exp\{\log \pi_k + x_i \mathbb{E}[\mu_k] - \mathbb{E}[\mu_k^2]/2\}. \quad (40)$$

– For the Gaussian conjugate prior, we map

$$\eta = \langle \mu/\sigma^2, 1/\sigma^2 \rangle. \quad (49)$$

– This gives the variational update in mean parameter form,

$$\mathbb{E}[\mu_k] = \frac{\mu_0/\sigma_0^2 + \sum_{i=1}^n \mathbb{E}[z_i^k] x_i}{1/\sigma_0^2 + \sum_{i=1}^n \mathbb{E}[z_i^k]} \quad (50)$$

$$\text{Var}(\mu_k) = 1/(1/\sigma_0^2 + \sum_{i=1}^n \mathbb{E}[z_i^k]). \quad (51)$$

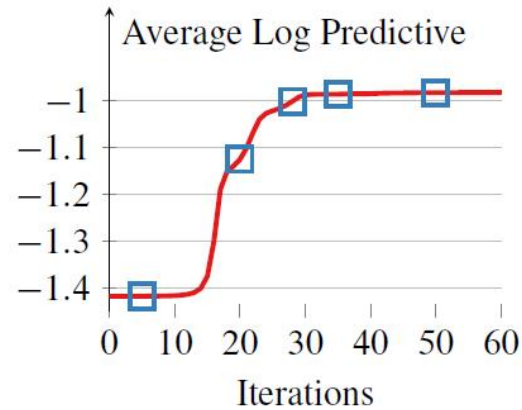
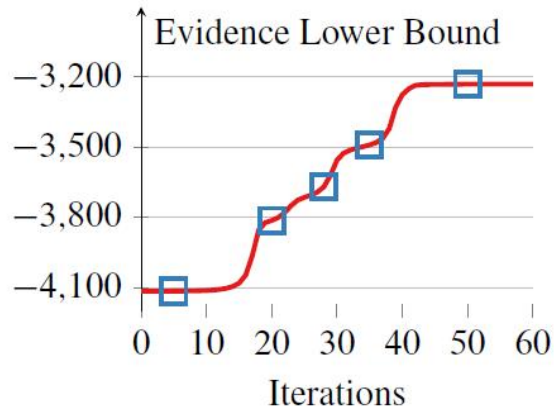
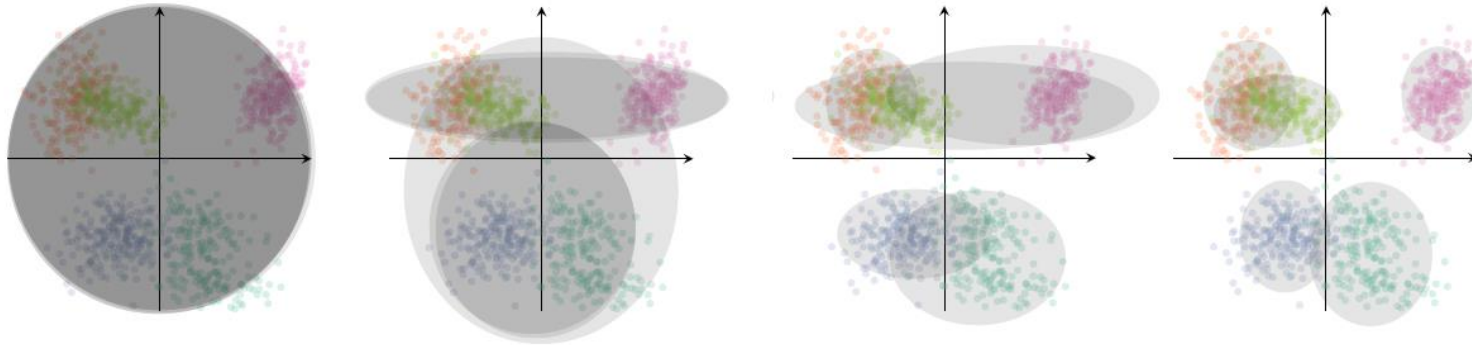
These are the usual Bayesian updates with the data weighted by its variational probability of being assigned to cluster k .



- We are given data $x_{1:n}$, hyperparameters μ_0 and σ_0^2 , and a number of groups K .

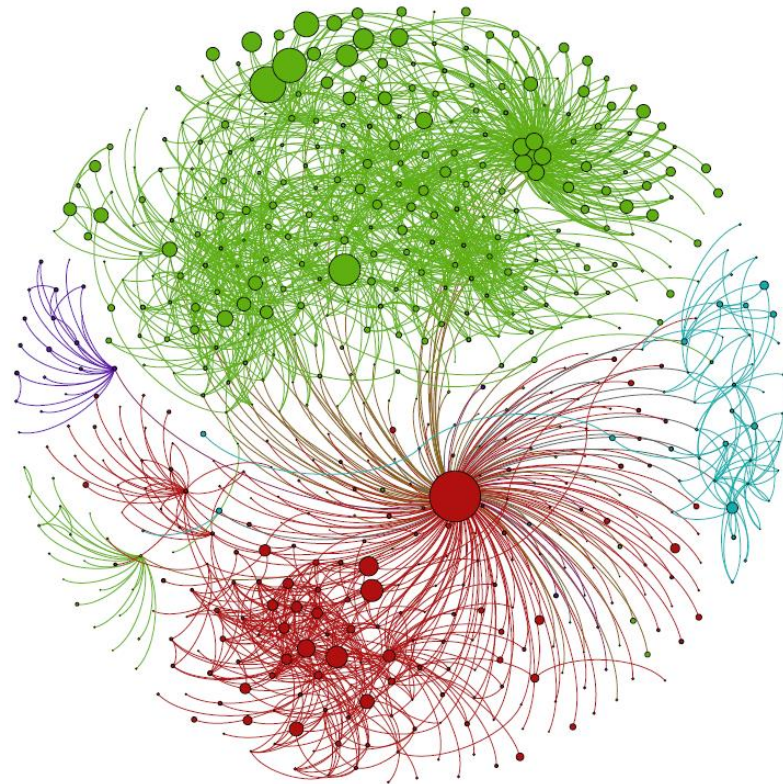
- The variational distributions are
 - * n variational multinomials $q(z_i)$
 - * K variational Gaussians $q(\mu_k | \tilde{\mu}_k, \tilde{\sigma}_k^2)$.
- Repeat until the ELBO converges:
 1. For each data point x_i
 - * Update the variational multinomial $q(z_i)$ from Equation 40.
 2. For each cluster $k = 1 \dots K$
 - * Update the mean and variance from Equation 50 and Equation 51.

Mixtures of Gaussians - Example



[images by Alp Kucukelbir]

Example Applications



Communities discovered in a 3.7M node network of U.S. Patents

[Gopalan and Blei, PNAS 2013]

Applications (Continued)



Topics found in 1.8M articles from the New York Times

[Hoffman, Blei, Wang, Paisley, JMLR 2013]



- alternatives to the optimization of KL as the variational objective function (tighter lower bounds)
- Strong independence assumptions of the mean-field family
- Better understanding of statistical properties of variational inference

- Not Covered in this lecture:
 - Stochastic Gradient Descent combined with Variational Methods
 - Black-box Methods
 - Other methods – e.g. Variational Auto-encoders



Backup