

The next few lectures describe methods for performing inference on MRF's

This lecture introduces free energies and mean field theory approximations to them.

But first, we describe steepest descent. ^{many variables here is the simplob.}

If we have a function $E[x]$ then we can perform steepest descent by

$$\underline{x}^{t+1} = \underline{x}^t - \epsilon \nabla E[\underline{x}]$$

provided $E[\underline{x}]$ is differentiable - i.e. the state variables \underline{x} take continuous values.

In MRF's, the variables are often discrete but can be treated as if they are continuous. For example, in the weak membrane model the ω_i take values $0, 1, \dots, 255$.

But it is convenient to approximate them as continuous values. Alternatively we can think of the ω_i being real values (this can give technical problems for some distributions - but we ignore that.)

$$\text{For the MRF } P(\underline{\omega} | I) = \frac{1}{Z} e^{\sum_i \phi_i(\omega_i | I) + \sum_{i,j} \psi_{ij}(\omega_i, \omega_j)}$$

$$\underline{\hat{\omega}} = \text{ARG MAX}_{\underline{\omega}} P(\underline{\omega} | I)$$

is equivalent to minimizing $E[\underline{\omega}] = \sum_i \phi_i(\omega_i | I) + \sum_{i,j} \psi_{ij}(\omega_i, \omega_j)$

Steepest descent: $\nabla_{\underline{\omega}} E[\underline{\omega}] = \frac{\partial \phi_i}{\partial \omega_i} + \sum_{j \in \text{nbr}(i)} \frac{\partial \psi_{ij}}{\partial \omega_i}$

$$\omega_i^{t+1} = \omega_i^t - \epsilon \nabla_{\omega_i} E[\underline{\omega}]$$

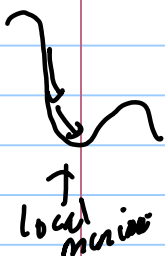
Steepest descent is guaranteed to decrease the energy monotonically - provided ϵ is small enough.

The problems are:

(i) hard to pick the best size of ϵ .

Too small - slow convergence, Too large - unstable algorithm

(ii) algorithm gets the energy surface - it can't make large moves or avoid even small local minima



(2) Variational Free Energies:

$p(\underline{x}|I)$ → suppose we want to estimate the marginals $p(x_i|I) = \sum_{\mathcal{X}_{-i}} p(\underline{x}|I)$.

(later we will extend this to estimate the $\hat{x} = \text{Arg Max } p(\underline{x}|I)$)

Define pseudomarginals $b_i(x_i)$ $b_i(x_i) \geq 0$
these are 'variables' that will $\sum_{x_i} b_i(x_i)$
be used to approximate the $p(x_i|I)$.

Full distribution $B(\underline{x}) = \prod_{i \in \mathcal{D}} b_i(x_i)$

Define a similarity measure between $B(\underline{x})$ and the full distribution $p(\underline{x}|I)$

Kullback-Leibler divergence

$$KL(B, P) = \sum_{\underline{x}} B(\underline{x}) \log \frac{B(\underline{x})}{P(\underline{x}|I)}$$

Proposition $KL(B, P) \geq 0$,

$= 0$, only if $B(\underline{x}) = P(\underline{x}|I)$.

Note: $B^*(\underline{x})$ which minimizes $KL(B, P)$ is the factorizable distribution which best approximates $p(\underline{x}|I)$.

Substituting the form of $B(\underline{x})$ into $KL(B, P)$ yields

$$KL(B, P) = -F_{\text{MFT}}(B) + \log Z.$$

where the 'free energy' $F_{\text{MFT}}(B)$ is:

$$F_{\text{MFT}}(B) = \sum_{i, j \in \mathcal{E}} \sum_{x_i, x_j} b_i(x_i) b_j(x_j) \psi_{ij}(x_i, x_j) \\ + \sum_{i \in \mathcal{V}} \sum_{x_i} b_i(x_i) \phi_i(x_i, I) + \sum_i \sum_{x_i} b_i(x_i) \log b_i(x_i)$$

First term is $\sum_{\underline{x}} B(\underline{x}) E(\underline{x}, I)$ Second is $\sum_{\underline{x}} B(\underline{x}) \log B(\underline{x})$
expected energy negative entropy

Minimizing $F_{\text{MFT}}(B)$ w.r.t. B gives

(i) an upper bound for $\log Z \geq -F_{\text{MFT}}(B^*)$
(relates to model evidence)

(ii) the best factorized approximation to $p(\underline{x}|I)$

(iii) an estimate of the marginals of the distribution
— which relate to the estimate $\hat{x} = \text{Arg Max } p(\underline{x}|I)$

(3) We still need an algorithm to minimize the free energy. Also the free energy function might have local minima so it is possible that the algorithm gets trapped in a local minimum.

But, the free energy formulation replaces a function $E[x; \theta]$ of discrete variables by a function $F[B]$ of continuous variables. This enables us to compute derivatives, do steepest descent and other algorithms (see later).

For the weak membrane model we treat the variables $\{w_{ij}\}$ as if it is continuous, and define pseudomoments over the process variables $\{y_{ij}\}$

$$F_{EM}(B, x) = - \sum_y B(y) \log P(x, y | I) + \sum_y B(y) \log B(y)$$

Also expected energy + negative entropy.

It can be shown (easily) that if we can minimize $F_{EM}(B, x)$ w.r.t. B, x to get $(B^*, x^*) = \text{ARG-MAX}_{(B, x)} F_{EM}(B, x)$

$$\text{then } x^* = \text{ARG-MAX}_x P(x | I) \text{ with } P(x | I) = \sum_y P(x, y | I)$$

This is the theoretical background for the Expectation Maximization (EM) algorithm.

The E-step of EM is equivalent to minimizing $F_{EM}(B, x)$ with respect to B

M-step is equivalent to minimizing $F_{EM}(B, x)$ with respect to x

This free energy formulation shows that the EM algorithm is guaranteed to converge since each step reduces the free energy.

But this formulation is richer than standard EM because it also shows convergence for any algorithm which reduces the free energy at each time step.

Note: this still does not guarantee convergence to the global minimum.

(4) Free energy for the weak membrane.

$$F_r(\beta, \mathbf{x}) = \tau \sum_i (x_i - I_i)^2 + A \sum_{ij} (1 - b_{ij}) (x_i - x_j)^2 + \beta \sum_{ij} b_{ij} + \sum_{ij} \{ b_{ij} \log b_{ij} + (1 - b_{ij}) \log (1 - b_{ij}) \}$$

where $b_{ij} = \alpha_{ij} (y_{ij} = 1)$.

The M-step:

$$\min \text{ wrt. } \langle x_i \rangle \quad \tau \sum_i (x_i - I_i)^2 + A \sum_{ij} (1 - b_{ij}) (x_i - x_j)^2$$

quadratic problem - solution by linear algebra.

note: non-zero b_{ij} weakens the smoothing between x_i & x_j . if $b_{ij} = 1$, then no smoothing occurs

E-step:

min wrt $b_{ij} \rightarrow$ give analytic solution

$$b_{ij} = \frac{1}{1 + e^{-A(x_i - x_j)^2 / \beta}}$$

if $|x_i - x_j|$ is large then $b_{ij} \rightarrow 1$
 $|x_i - x_j|$ is small then $b_{ij} \rightarrow 0$

The iterative algorithm is like smoothing the input $\{I_i\}$, but adapting the weights of the smoothing functions (the b_{ij}) so that they become large (close to 1) if $|x_i - x_j|$ is big.

Note: initial conditions must be well chosen to help converge to the global minimum.

E.g. initialize so that $x_i = I_i$ at time $t=0$, then update the $b_{ij} \dots \rightarrow$ likely to converge (but no guarantees).

Note: it is unusual that the steps of an EM algorithm are so simple.

In particular, the simplicity of the E-step is because the distribution $P(\underline{x}, \underline{y} | I)$ is factorizable in the \underline{y} -variables.

suppose we modify the energy $E(\underline{x}, \underline{y}; I)$ by adding terms like $C \sum_{ij} y_{ij} y_{kl}$, which encourages edges to be continuous spatially.

In this case, we cannot solve the E-step analytically but we can modify the free energy to approximate $P(\underline{y})$ by a factorized distribution.

(5)

Discrete Iterative Algorithms

We give procedures to construct iterative algorithms that can decrease energy functions by making non-local steps and without needing a step size parameter (as for steepest descent).

Variational Bounding

Current state \underline{x}^t

define a bounding energy $E_b(\underline{x}; \underline{x}^t)$

such that:

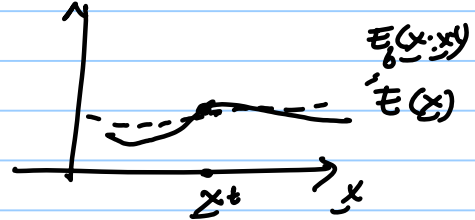
$$E_b(\underline{x}^t; \underline{x}^t) = E(\underline{x}^t)$$

$$\& E(\underline{x}) \leq E_b(\underline{x}; \underline{x}^t) \quad \forall \underline{x}$$

Then choose \underline{x}^{t+1} s.t. $E_b(\underline{x}^{t+1}; \underline{x}^t) \leq E_b(\underline{x}^t; \underline{x}^t)$
which implies $E(\underline{x}^{t+1}) \leq E(\underline{x}^t)$.

Intuition \rightarrow pick the bounding energy $E_b(\underline{x}; \underline{x}^t)$

so that it is easy to minimize



CCCP \rightarrow Concave Convex Procedure.

This is a special case of variational bounding which gives a special set of bounding energies.

Claim: for almost all energies, $E(\underline{x})$, we can write it as $E(\underline{x}) = E_{\text{conv}}(\underline{x}) + E_{\text{conc}}(\underline{x}) \rightarrow$ seems paradoxical

(to prove, consider the Hessian of $E, E_{\text{conv}}, E_{\text{conc}}$)

$$W = U + \Lambda$$

for any W .

$$\text{Then } E_b(\underline{x}; \underline{x}^t) = E_{\text{conv}}(\underline{x}) + E_{\text{conc}}(\underline{x}^t) + (\underline{x} - \underline{x}^t) \cdot \frac{\partial E_{\text{conc}}(\underline{x}^t)}{\partial \underline{x}}$$

this is a convex function - so usually possible to minimize it.

$$\text{corresponds to an update } \frac{\partial E_{\text{conv}}(\underline{x}^{t+1})}{\partial \underline{x}} = -\frac{\partial E_{\text{conc}}(\underline{x}^t)}{\partial \underline{x}}$$

Note: this is a non-local method, so it may avoid local minima that a steepest descent algorithm can get trapped in. But no guarantee that you get to the local minima.

(6) Note: Researchers have invented Discrete Iterative Algorithms for problems on a one by one basis — E.M., Generalized Iterative Scaling, Sinkhorn — but this can be shown to all be equivalent to CCCP.

Example: $E[\omega] = \frac{1}{2} \sum_{ij} T_{ij} \omega_i \omega_j + \sum_i \theta_i \omega_i$ $\omega_i \in [0, 1]$

Mean Field Theory

$$F(b) = \frac{1}{2} \sum_{ij} T_{ij} b_i b_j + \sum_i \theta_i b_i + \sum_i \{ b_i \log b_i + (1-b_i) \log (1-b_i) \}$$

concave
convex

Assume T_{ij} is negative definite matrix

$$F_{\text{max}} = \sum_i \{ b_i \log b_i + (1-b_i) \log (1-b_i) \}$$

$$F_{\text{line}} = \frac{1}{2} \sum_{ij} T_{ij} b_i b_j + \sum_i \theta_i b_i$$

$$\frac{\partial F_{\text{max}}}{\partial b_i} = \log \frac{b_i}{1-b_i}, \quad \frac{\partial F_{\text{line}}}{\partial b_i} = \sum_j T_{ij} + \theta_i$$

yields $b_i^{t+1} = \frac{e^{-\left\{ \sum_j T_{ij} b_j^t + \theta_i \right\}}}{1 + e^{-\left\{ \sum_j T_{ij} b_j^t + \theta_i \right\}}}$

non-local DIA.

Note: steepest descent can be derived as a special case.

$$E(x) = \frac{1}{2} Ax^2 + \underbrace{E(x)} - \underbrace{\frac{1}{2} Ax^2}$$

convex
concave

$$Ax^{t+1} = - \frac{\partial E(x^t)}{\partial x} + Ax^t$$

$$x^{t+1} = x^t - \frac{1}{A} \frac{\partial E(x^t)}{\partial x} \quad \epsilon = \frac{1}{A}$$

Give a size for the ϵ -parameter (to ensure convexity of $E(x) - \frac{1}{2} Ax^2$)
 this can be used to help design stable steepest descent algorithms (e.g. Eyre's method)