

Probabilistic Graphical Models

Lecture 8: Information Theory. First Variational Approximation

Volkan Cevher, Matthias Seeger
Ecole Polytechnique Fédérale de Lausanne

24/10/2011



- 1 Why Approximate Inference?
- 2 Some Information Theory
- 3 Variational Mean Field Approximations

Why Approximate Inference?

- Why inference (computing marginal posterior distributions)?
Essential backbone for (almost) anything todo with probabilistic model
 - Answering queries (honest answer: with uncertainties)
 - Learning model parameters
 - Making good decisions
 - Direct further data acquisition
 - Planning strategies (beyond single decisions)

Why Approximate Inference?

- Why inference (computing marginal posterior distributions)?
Essential backbone for (almost) anything todo with probabilistic model
 - Answering queries (honest answer: with uncertainties)
 - Learning model parameters
 - Making good decisions
 - Direct further data acquisition
 - Planning strategies (beyond single decisions)
- Why approximate inference?
Exact inference intractable for almost all real-world models
 - Loops in graphical model: Blow-up of intermediate representations, with no efficient (dynamic programming) way around
 - Potentials not closed under conditioning / marginalization: Blow-up of messages even for tree graphical models
- Bottomline: Bayesian inference powerful, consistent idea.
Without approximate inference: Entirely academic exercise

Conjugate Priors

Transition Part I → Part II

I'll mention some things quickly, without us looking at them in more detail

Sometimes, inference is simple

\mathbf{y} Observation

θ Latent parameters (query)

$P(\mathbf{y}|\theta)$ Likelihood potential (positive function of θ)

Family of distributions $\mathcal{F} = \{P(\theta|\alpha)\}$, α fixed size:

F2

- For every \mathbf{y} : $P(\theta) \in \mathcal{F} \Rightarrow P(\theta|\mathbf{y}) \in \mathcal{F}$
- If $P(\theta) = P(\theta|\alpha_0)$, $P(\theta|\mathbf{y}) = P(\theta|\alpha_1)$: For every (α_0, \mathbf{y}) :
 α_1 easy to find

\Rightarrow Inference a piece of cake! \mathcal{F} **conjugate** to $P(\mathbf{y}|\theta)$ (or to $\{P(\mathbf{y}|\theta)\}$)

Markov Chain Monte Carlo

- General, maybe most flexible framework for approximate inference. Ideas from physics (thermodynamics, statistical mechanics)
- Not covered here (would need own course). I'll just give you cocktail party summary

Markov Chain Monte Carlo

- 1 Inference needs integrals $\int f(\mathbf{x})P(\mathbf{x}) d\mathbf{x}$, \mathbf{x} high-dimensional, $P(\mathbf{x})$ coupled, complicated (posterior)
- 2 Law of large numbers: $\mathbf{x}_1, \dots, \mathbf{x}_N \sim P(\mathbf{x})$ independent:
 $N^{-1} \sum_i f(\mathbf{x}_i) \rightarrow E_P[f(\mathbf{x})]$ almost surely.
Central limit theorem: P, f nice \Rightarrow Convergence as $1/\sqrt{N}$
independent of \mathbf{x} dimensionality.
Catch: Sampling from $P(\mathbf{x})$ hard as well

Markov Chain Monte Carlo

- 1 Inference needs integrals $\int f(\mathbf{x})P(\mathbf{x}) d\mathbf{x}$, \mathbf{x} high-dimensional, $P(\mathbf{x})$ coupled, complicated (posterior)
- 2 Law of large numbers: $\mathbf{x}_1, \dots, \mathbf{x}_N \sim P(\mathbf{x})$ independent:
 $N^{-1} \sum_i f(\mathbf{x}_i) \rightarrow E_P[f(\mathbf{x})]$ almost surely.
 Central limit theorem: P, f nice \Rightarrow Convergence as $1/\sqrt{N}$
independent of \mathbf{x} dimensionality.
 Catch: Sampling from $P(\mathbf{x})$ hard as well
- 3 Let's just do **something**: Start with some \mathbf{x} , draw $\mathbf{x}' \sim K(\mathbf{x}'|\mathbf{x})$, keep doing that. At the very least:

$$P(\mathbf{x}') = \int K(\mathbf{x}'|\mathbf{x})P(\mathbf{x}) d\mathbf{x}$$

Such kernels K exist, need evaluation of $\propto P(\mathbf{x})$ only

Markov Chain Monte Carlo

$$P(\mathbf{x}') = \int K(\mathbf{x}'|\mathbf{x})P(\mathbf{x}) d\mathbf{x}$$

- ④ MCMC magic: Under mild assumptions, that's **all** we need:

$\mathbf{x}^{(j+1)} \sim K(\cdot|\mathbf{x}^{(j)}) \Rightarrow$ Marginal $\mathbf{x}^{(j)} \xrightarrow{D} P(\mathbf{x})$ as $j \rightarrow \infty$

Rough idea why:

- $K(\mathbf{x}'|\mathbf{x})$ contraction of probability mass. Information propagation with K brings marginal distributions closer together
- There is only one fixed point (here: mild assumptions)

Markov Chain Monte Carlo

- MCMC used for many things besides approximate inference
 - Theoretical CS: Counting of combinatorial sets. Volume estimation
 - Statistical physics: Evaluation of thermodynamical numbers (entropy, volume of macrostates). Studying phase transitions of coupled spin systems (magnets, spin glasses)
- Rich theory in the discrete case
- Related to, but different from stochastic optimization

Beware

BEWARE! MCMC sampling can be dangerous!

[OpenBUGS User Manual, page 1]

- MCMC: Simple to code. **Hard** to use properly
- You never exactly know when you're done
 - No definite convergence test in general
 - Hard to spot failures. Very hard to debug
 - Slow convergence can happen even with unimodal distributions, Gaussian tails

Beware

BEWARE! MCMC sampling can be dangerous!

[OpenBUGS User Manual, page 1]

- MCMC: Simple to code. **Hard** to use properly
- You never exactly know when you're done
 - No definite convergence test in general
 - Hard to spot failures. Very hard to debug
 - Slow convergence can happen even with unimodal distributions, Gaussian tails
- MCMC: Black box (in most cases), for good and for bad
 - Easy to code. For some problems, nothing else works. Safe if answers can be checked (search, exploration)
 - Can be very slow, or fail without you noticing. Always compare against something else if you can

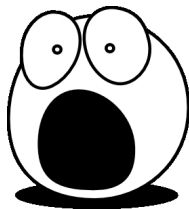
Beware

BEWARE! MCMC sampling can be dangerous!

[OpenBUGS User Manual, page 1]

Dare to find out for yourself?

- Neal: Probabilistic Inference using Markov Chain Monte Carlo Methods (1993)
[<http://www.cs.toronto.edu/~radford/papers-online.html>]
- Gilks *et.al.*: Markov Chain Monte Carlo in Practice (1996)



Elements of Information Theory

Wake Up!

Transition time is over

Elements of Information Theory

Information Theory (Shannon, 1948)

- Narrow sense:
 - Limits of data compression (and how to achieve them)
 - Limits of error-free(!) communication over noisy channel
- Wide sense:
 - Basis of communication (language)
 - What is information? How to best encode it
 - Basis of anything adaptive, of learning
 - Source of great simplifications in number of mathematical domains
 - Information theory \leftrightarrow applied probability / decision theory:
Essentially equivalent in basic concepts, problems, methods



Good luck for students: Amazing textbook available:

- Cover, Thomas: Elements of Information Theory (1991)

Entropy of Distribution

$$H[P(\mathbf{x})] = E_P[-\log P(\mathbf{x})] = - \sum_{\mathbf{x}} P(\mathbf{x}) \log P(\mathbf{x})$$

- Game of questions: I draw $\mathbf{x} \sim P(\mathbf{x})$, give you P but not \mathbf{x} . How many questions [$\mathbf{x} \in \mathcal{E}$] do you need to pin down \mathbf{x} ?

F6

Entropy of Distribution

$$H[P(\mathbf{x})] = E_P[-\log P(\mathbf{x})] = - \sum_{\mathbf{x}} P(\mathbf{x}) \log P(\mathbf{x})$$

- Game of questions: I draw $\mathbf{x} \sim P(\mathbf{x})$, give you P but not \mathbf{x} . How many questions [$\mathbf{x} \in \mathcal{E}$] do you need to pin down \mathbf{x} ?
- Shannon: On average: $\leq H[P(\mathbf{x})] + 1$ questions if you're smart, no less than $H[P(\mathbf{x})]$ even for a genius (log to base 2)
 - \Rightarrow Equivalent: Number bits needed to encode \mathbf{x}
 - \Rightarrow **Amount of uncertainty** in $P(\mathbf{x})$

F6b

Entropy of Distribution

$$H[P(\mathbf{x})] = E_P[-\log P(\mathbf{x})] = - \sum_{\mathbf{x}} P(\mathbf{x}) \log P(\mathbf{x})$$

- Game of questions: I draw $\mathbf{x} \sim P(\mathbf{x})$, give you P but not \mathbf{x} . How many questions [$\mathbf{x} \in \mathcal{E}$] do you need to pin down \mathbf{x} ?
- Shannon: On average: $\leq H[P(\mathbf{x})] + 1$ questions if you're smart, no less than $H[P(\mathbf{x})]$ even for a genius (log to base 2)
 - \Rightarrow Equivalent: Number bits needed to encode \mathbf{x}
 - \Rightarrow **Amount of uncertainty** in $P(\mathbf{x})$

Joint entropy $H[P(\mathbf{y}, \mathbf{x})] = E_P[-\log P(\mathbf{y}, \mathbf{x})]$

Conditional entropy $H[P(\mathbf{y}|\mathbf{x})] = E_P[-\log P(\mathbf{y}|\mathbf{x})]$

F6c

Entropy of Distribution

$$H[P(\mathbf{x})] = E_P[-\log P(\mathbf{x})] = - \sum_{\mathbf{x}} P(\mathbf{x}) \log P(\mathbf{x})$$

- Game of questions: I draw $\mathbf{x} \sim P(\mathbf{x})$, give you P but not \mathbf{x} . How many questions [$\mathbf{x} \in \mathcal{E}$] do you need to pin down \mathbf{x} ?
- Shannon: On average: $\leq H[P(\mathbf{x})] + 1$ questions if you're smart, no less than $H[P(\mathbf{x})]$ even for a genius (log to base 2)
 - \Rightarrow Equivalent: Number bits needed to encode \mathbf{x}
 - \Rightarrow **Amount of uncertainty** in $P(\mathbf{x})$

Joint entropy $H[P(\mathbf{y}, \mathbf{x})] = E_P[-\log P(\mathbf{y}, \mathbf{x})]$

Conditional entropy $H[P(\mathbf{y}|\mathbf{x})] = E_P[-\log P(\mathbf{y}|\mathbf{x})]$

Chain rule of entropy:

$$H[P(\mathbf{x}_1, \dots, \mathbf{x}_n)] = \sum_{i=1}^n H[P(\mathbf{x}_i | \mathbf{x}_{<i})]$$

F6d

Relative Entropy

$$D[P(\mathbf{x}) \parallel Q(\mathbf{x})] = E_P \left[\log \frac{P(\mathbf{x})}{Q(\mathbf{x})} \right]$$

- Game of questions. This time, you get it wrong. You think $\mathbf{x} \sim Q(\mathbf{x})$, but in fact $\mathbf{x} \sim P(\mathbf{x})$. How many questions?

F7

Relative Entropy

$$D[P(\mathbf{x}) \parallel Q(\mathbf{x})] = E_P \left[\log \frac{P(\mathbf{x})}{Q(\mathbf{x})} \right]$$

- Game of questions. This time, you get it wrong. You think $\mathbf{x} \sim Q(\mathbf{x})$, but in fact $\mathbf{x} \sim P(\mathbf{x})$. How many questions?
- On average: $E_P[-\log Q(\mathbf{x})] = H[P(\mathbf{x})] + D[P(\mathbf{x}) \parallel Q(\mathbf{x})]$
 \Rightarrow Number of **additional** bits for using Q instead of true P
 \Rightarrow Natural divergence (distance) measure between distributions F7b

Relative Entropy

$$D[P(\mathbf{x}) \parallel Q(\mathbf{x})] = E_P \left[\log \frac{P(\mathbf{x})}{Q(\mathbf{x})} \right]$$

- Game of questions. This time, you get it wrong. You think $\mathbf{x} \sim Q(\mathbf{x})$, but in fact $\mathbf{x} \sim P(\mathbf{x})$. How many questions?
- On average: $E_P[-\log Q(\mathbf{x})] = H[P(\mathbf{x})] + D[P(\mathbf{x}) \parallel Q(\mathbf{x})]$
 \Rightarrow Number of **additional** bits for using Q instead of true P
 \Rightarrow Natural divergence (distance) measure between distributions
- Other name: Kullback-Leibler divergence.
 No distance: $D[P \parallel Q] \neq D[Q \parallel P]$

Relative Entropy

$$D[P(\mathbf{x}) \parallel Q(\mathbf{x})] = E_P \left[\log \frac{P(\mathbf{x})}{Q(\mathbf{x})} \right]$$

- Game of questions. This time, you get it wrong. You think $\mathbf{x} \sim Q(\mathbf{x})$, but in fact $\mathbf{x} \sim P(\mathbf{x})$. How many questions?
- On average: $E_P[-\log Q(\mathbf{x})] = H[P(\mathbf{x})] + D[P(\mathbf{x}) \parallel Q(\mathbf{x})]$
 \Rightarrow Number of **additional** bits for using Q instead of true P
 \Rightarrow Natural divergence (distance) measure between distributions
- Other name: Kullback-Leibler divergence.
 No distance: $D[P \parallel Q] \neq D[Q \parallel P]$

Conditional relative entropy:

$$D[P(\mathbf{y}|\mathbf{x}) \parallel Q(\mathbf{y}|\mathbf{x})] = E_P[\log\{P(\mathbf{y}|\mathbf{x})/Q(\mathbf{y}|\mathbf{x})\}]$$

Chain rule of relative entropy:

$$D[P(\mathbf{y}, \mathbf{x}) \parallel Q(\mathbf{y}, \mathbf{x})] = D[P(\mathbf{y}|\mathbf{x}) \parallel Q(\mathbf{y}|\mathbf{x})] + D[P(\mathbf{x}) \parallel Q(\mathbf{x})]$$

F7c

Mutual Information

$$I(\mathbf{x}; \mathbf{y}) = D[P(\mathbf{x}, \mathbf{y}) \parallel P(\mathbf{x})P(\mathbf{y})] = E_P \left[\log \frac{P(\mathbf{x}, \mathbf{y})}{P(\mathbf{x})P(\mathbf{y})} \right]$$

- \mathbf{x} , \mathbf{y} may be dependent. How many additional questions (bits) for ignoring that?

F8

Mutual Information

$$I(\mathbf{x}; \mathbf{y}) = D[P(\mathbf{x}, \mathbf{y}) \parallel P(\mathbf{x})P(\mathbf{y})] = E_P \left[\log \frac{P(\mathbf{x}, \mathbf{y})}{P(\mathbf{x})P(\mathbf{y})} \right]$$

- \mathbf{x} , \mathbf{y} may be dependent. How many additional questions (bits) for ignoring that?
- Mutual information: Reduction in uncertainty of one random variable due to knowledge of other

$$I(\mathbf{x}; \mathbf{y}) = H[P(\mathbf{x})] - H[P(\mathbf{x}|\mathbf{y})] = H[P(\mathbf{y})] - H[P(\mathbf{y}|\mathbf{x})]$$

\Rightarrow Amount of information \mathbf{x} about \mathbf{y} , or \mathbf{y} about \mathbf{x}

Mutual Information

$$I(\mathbf{x}; \mathbf{y}) = D[P(\mathbf{x}, \mathbf{y}) \parallel P(\mathbf{x})P(\mathbf{y})] = E_P \left[\log \frac{P(\mathbf{x}, \mathbf{y})}{P(\mathbf{x})P(\mathbf{y})} \right]$$

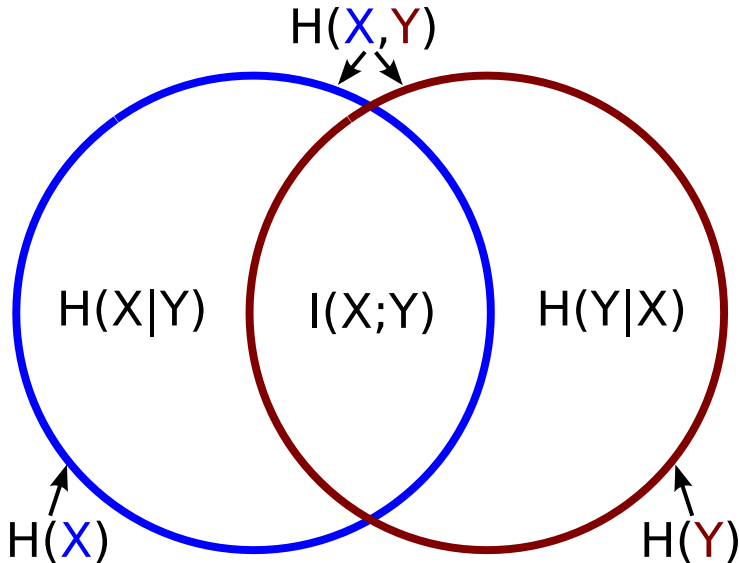
- \mathbf{x} , \mathbf{y} may be dependent. How many additional questions (bits) for ignoring that?
- Mutual information: Reduction in uncertainty of one random variable due to knowledge of other

$$I(\mathbf{x}; \mathbf{y}) = H[P(\mathbf{x})] - H[P(\mathbf{x}|\mathbf{y})] = H[P(\mathbf{y})] - H[P(\mathbf{y}|\mathbf{x})]$$

\Rightarrow Amount of information \mathbf{x} about \mathbf{y} , or \mathbf{y} about \mathbf{x}

- Note: $\mathbf{x} \perp \mathbf{y}$ (independent) $\Rightarrow P(\mathbf{x}, \mathbf{y}) = P(\mathbf{x})P(\mathbf{y}) \Rightarrow I(\mathbf{x}; \mathbf{y}) = 0$. We'll see \Leftarrow . Mutual information: Measure of dependence

Venn Diagram for Information



Information Inequality

- Something missing here

- **More** questions for getting it wrong:

- $H[P(\mathbf{x})] \rightarrow H[P(\mathbf{x})] + D[P(\mathbf{x}) \parallel Q(\mathbf{x})]$

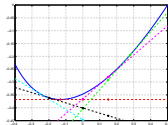
- $I(\mathbf{x}; \mathbf{y})$ measures dependence. $I(\mathbf{x}; \mathbf{y}) = 0$ for $\mathbf{x} \perp \mathbf{y}$

Is $D[P(\mathbf{x}) \parallel Q(\mathbf{x})] \geq 0$? Is $I(\mathbf{x}; \mathbf{y}) \geq 0$?

Information Inequality

- Something missing here
 - **More** questions for getting it wrong:
 - $H[P(\mathbf{x})] \rightarrow H[P(\mathbf{x})] + D[P(\mathbf{x}) \parallel Q(\mathbf{x})]$
 - $I(\mathbf{x}; \mathbf{y})$ measures dependence. $I(\mathbf{x}; \mathbf{y}) = 0$ for $\mathbf{x} \perp \mathbf{y}$
 - Is $D[P(\mathbf{x}) \parallel Q(\mathbf{x})] \geq 0$? Is $I(\mathbf{x}; \mathbf{y}) \geq 0$?
- Convexity comes to the rescue. Jensen's inequality F10

$$E_P[f(\mathbf{x})] \geq f(E_P[\mathbf{x}]), \quad f \text{ convex}$$

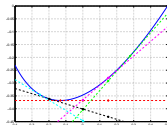


Information Inequality

- Something missing here
 - **More** questions for getting it wrong:
 - $H[P(\mathbf{x})] \rightarrow H[P(\mathbf{x})] + D[P(\mathbf{x}) \parallel Q(\mathbf{x})]$
 - $I(\mathbf{x}; \mathbf{y})$ measures dependence. $I(\mathbf{x}; \mathbf{y}) = 0$ for $\mathbf{x} \perp \mathbf{y}$
- Is $D[P(\mathbf{x}) \parallel Q(\mathbf{x})] \geq 0$? Is $I(\mathbf{x}; \mathbf{y}) \geq 0$?

- Convexity comes to the rescue. Jensen's inequality

$$E_P[f(\mathbf{x})] \geq f(E_P[\mathbf{x}]), \quad f \text{ convex}$$



- **Information inequality:** $D[P(\mathbf{x}) \parallel Q(\mathbf{x})] \geq 0$.
 Since $-\log(\cdot)$ strictly convex (nowhere linear):
 $D[P(\mathbf{x}) \parallel Q(\mathbf{x})] = 0 \Leftrightarrow P(\mathbf{x}) = Q(\mathbf{x})$ P -almost everywhere.

$$I(\mathbf{x}; \mathbf{y}) \geq 0; \quad I(\mathbf{x}; \mathbf{y}) = 0 \Leftrightarrow \mathbf{x} \perp \mathbf{y}$$

Corollaries

Raking in the fruits

- Conditioning reduces entropy (learning always helps)

F11

$$H[P(\mathbf{x}|\mathbf{y})] \leq H[P(\mathbf{x})]$$

Corollaries

Raking in the fruits

- Conditioning reduces entropy (learning always helps)

$$H[P(\mathbf{x}|\mathbf{y})] \leq H[P(\mathbf{x})]$$

- Conditional mutual information:
Measure for **conditional independence**

F11b

$$I(\mathbf{x}; \mathbf{y} | \mathbf{z}) = 0 \quad \Leftrightarrow \quad \mathbf{x} \perp \mathbf{y} | \mathbf{z}$$

Corollaries

Raking in the fruits

- Conditioning reduces entropy (learning always helps)

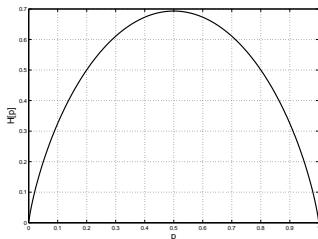
$$H[P(\mathbf{x}|\mathbf{y})] \leq H[P(\mathbf{x})]$$

- Conditional mutual information:
Measure for **conditional independence**

$$I(\mathbf{x}; \mathbf{y} | \mathbf{z}) = 0 \iff \mathbf{x} \perp \mathbf{y} | \mathbf{z}$$

- Entropy: **Concave** function

$$\begin{aligned} & H[\lambda P(\mathbf{x}) + (1 - \lambda)Q(\mathbf{x})] \\ & \geq \lambda H[P(\mathbf{x})] + (1 - \lambda)H[Q(\mathbf{x})] \end{aligned}$$



Remember EM?

One approach to variational approximate inference:

Computations with $P(\mathbf{x}) = Z^{-1} e^{\Psi(\mathbf{x})}$ hard (even $\log Z$)?

\Rightarrow Approximate it by $Q(\mathbf{x})$, for which computations simple

Remember derivation of EM?

F12

Remember EM?

One approach to variational approximate inference:

Computations with $P(\mathbf{x}) = Z^{-1} e^{\Psi(\mathbf{x})}$ hard (even $\log Z$)?

\Rightarrow Approximate it by $Q(\mathbf{x})$, for which computations simple

Remember derivation of EM?

$$\begin{aligned}\log Z &= \log \int e^{\Psi(\mathbf{x})} d\mathbf{x} = \sup_Q E_Q[\log\{e^{\Psi(\mathbf{x})}/Q(\mathbf{x})\}] \\ &= \sup_Q \{E_Q[\Psi(\mathbf{x})] + H[Q(\mathbf{x})]\}\end{aligned}$$

Was called **variational mean field inequality**. Let's see why

Remember EM?

One approach to variational approximate inference:

Computations with $P(\mathbf{x}) = Z^{-1} e^{\Psi(\mathbf{x})}$ hard (even $\log Z$)?

\Rightarrow Approximate it by $Q(\mathbf{x})$, for which computations simple

Remember derivation of EM?

$$\begin{aligned} \log Z &= \log \int e^{\Psi(\mathbf{x})} d\mathbf{x} = \sup_Q E_Q[\log\{e^{\Psi(\mathbf{x})}/Q(\mathbf{x})\}] \\ &= \sup_Q \{E_Q[\Psi(\mathbf{x})] + H[Q(\mathbf{x})]\} \end{aligned}$$

Was called **variational mean field inequality**. Let's see why

- Maximizer: $Q(\mathbf{x}) = P(\mathbf{x})$ itself. Attains $\log Z$
- Any other $Q(\mathbf{x})$: Lower bound. $Q(\mathbf{x})$ closer to $P(\mathbf{x})$?
 \Rightarrow Maximize the lower bound!

F12b

Remember EM?

One approach to variational approximate inference:

Computations with $P(\mathbf{x}) = Z^{-1} e^{\Psi(\mathbf{x})}$ hard (even $\log Z$)?

⇒ Approximate it by $Q(\mathbf{x})$, for which computations simple

Remember derivation of EM?

$$\begin{aligned} \log Z &= \log \int e^{\Psi(\mathbf{x})} d\mathbf{x} = \sup_Q E_Q[\log\{e^{\Psi(\mathbf{x})}/Q(\mathbf{x})\}] \\ &= \sup_Q \{E_Q[\Psi(\mathbf{x})] + H[Q(\mathbf{x})]\} \end{aligned}$$

Was called **variational mean field inequality**. Let's see why

- Maximizer: $Q(\mathbf{x}) = P(\mathbf{x})$ itself. Attains $\log Z$
- Any other $Q(\mathbf{x})$: Lower bound. $Q(\mathbf{x})$ closer to $P(\mathbf{x})$?
⇒ Maximize the lower bound!
- **Relax this problem:** Work with $\mathcal{Q} = \{Q(\mathbf{x})\}$:
 - Lower bound can be evaluated for each $Q \in \mathcal{Q}$
 - Bayesian computations can be done with any $Q \in \mathcal{Q}$ (not with P)

F12c

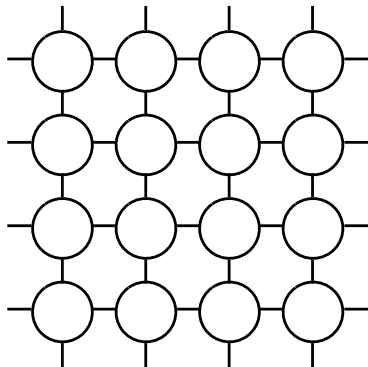
Naive Mean Field for Markov Random Fields

Distributions complicated, because they are **coupled**

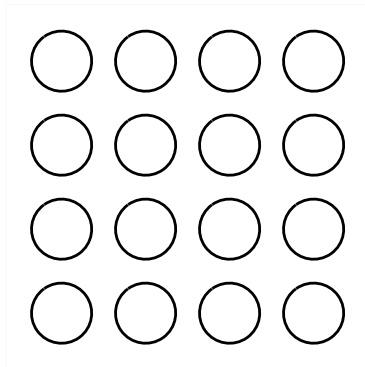
⇒ Mean field: Approximate them by **factorizing** distributions

Naive Mean Field: Drop all edges

True MRF posterior $P(\mathbf{x})$



Approximations $Q(\mathbf{x}) \in \mathcal{Q}$



Naive Mean Field for Markov Random Fields

Variational problem:

$$\operatorname{argmax}_{\{Q(x_k)\}} \left\{ \sum_j E_Q[\Psi_j(\mathbf{x}_{C_j})] + \sum_k H[Q(x_k)] \right\}$$

Our first variational algorithm:

Default-initialize $Q(\mathbf{x}_k)$ (say: uniform)

repeat

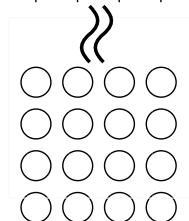
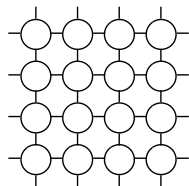
Pick some node k at random

Update $Q(x_k)$, keeping all others fixed

$$Q(x_k) \leftarrow \operatorname{argmax} \left\{ \sum_{j \in \mathcal{N}_k} E_Q[\Psi_j(\mathbf{x}_{C_j})] + H[Q(x_k)] \right\}$$

until Convergence

Prize question: How does that update look like?



Remarks

- Does this always converge? Yes. To a unique solution? No
- How to compare different fixed points? Or even different Q ?
You get **lower bound** to $\log Z$

F15

Remarks

- Does this always converge? Yes. To a unique solution? No
- How to compare different fixed points? Or even different Q ?
You get **lower bound** to $\log Z$
- Why “mean field”? $P(\mathbf{x})$: Random field. $Q(\mathbf{x}) \approx P(\mathbf{x})$ has no couplings ($E_Q[x_j x_k] = E_Q[x_j]E_Q[x_k]$).
True means at convergence ($E_Q[x_j] = E_P[x_j]$)? **No**
(remember: $E_P[x_j]$ hard as well!)

Remarks

- Does this always converge? Yes. To a unique solution? No
- How to compare different fixed points? Or even different Q ?
You get **lower bound** to $\log Z$
- Why “mean field”? $P(\mathbf{x})$: Random field. $Q(\mathbf{x}) \approx P(\mathbf{x})$ has no couplings ($E_Q[x_j x_k] = E_Q[x_j]E_Q[x_k]$).
True means at convergence ($E_Q[x_j] = E_P[x_j]$)? **No**
(remember: $E_P[x_j]$ hard as well!)
- General idea here: **Relax** variational problem

$$\sup_Q(\dots) \geq \sup_{Q \in \mathcal{Q}}(\dots)$$
 \mathcal{Q} : Subset of all distributions (factorization constraints). Each $Q(\mathbf{x})$ is distribution.
 \Rightarrow Maximize lower bound over \mathcal{Q}
Note: Might not find maximizer $Q \in \mathcal{Q}$, but local maximum

Variational Mean Field: Minimizing Relative Entropy

$$\log Z = \log \int e^{\Psi(\mathbf{x})} d\mathbf{x} \geq \mathbb{E}_Q[\Psi(\mathbf{x})] + H[Q(\mathbf{x})], \quad P(\mathbf{x}) = Z^{-1} e^{\Psi(\mathbf{x})}$$

- What is the slack in this bound?
Hint: ≥ 0 , and $= 0$ iff $Q(\mathbf{x}) = P(\mathbf{x})$

F16

Variational Mean Field: Minimizing Relative Entropy

$$\log Z = \log \int e^{\Psi(\mathbf{x})} d\mathbf{x} \geq \mathbb{E}_Q[\Psi(\mathbf{x})] + H[Q(\mathbf{x})], \quad P(\mathbf{x}) = Z^{-1} e^{\Psi(\mathbf{x})}$$

- What is the slack in this bound?
Hint: ≥ 0 , and $= 0$ iff $Q(\mathbf{x}) = P(\mathbf{x})$
- Variational mean field: Minimize slack (relative entropy)

$$\min_{Q \in \mathcal{Q}} D[Q(\mathbf{x}) \parallel P(\mathbf{x})]$$

Does that fit relative entropy semantics?

Variational Mean Field: Minimizing Relative Entropy

$$\log Z = \log \int e^{\Psi(\mathbf{x})} d\mathbf{x} \geq \mathbb{E}_Q[\Psi(\mathbf{x})] + H[Q(\mathbf{x})], \quad P(\mathbf{x}) = Z^{-1} e^{\Psi(\mathbf{x})}$$

- What is the slack in this bound?
Hint: ≥ 0 , and $= 0$ iff $Q(\mathbf{x}) = P(\mathbf{x})$
- Variational mean field: Minimize slack (relative entropy)

$$\min_{Q \in \mathcal{Q}} D[Q(\mathbf{x}) \parallel P(\mathbf{x})]$$

Does that fit relative entropy semantics?

- It's the wrong way around! We should minimize $D[P \parallel Q]$.
Alas, even that is hard. For naive mean field, unique solution is

$$Q(x_1, \dots, x_N) = P(x_1) \dots P(x_N)$$

Variational mean field: a **tractable compromise**

Wrap-Up

- Information theory: Fundamental characteristics and limits to compression and faultless information transmission
- Statistical learning, information theory: Different sides of the same coin
- Variational mean field: Tractable approximate inference by factorization assumptions
- Naive mean field: Drop all edges, update node by node