

# CPSC 540: Machine Learning

## Variational Inference, Non-Parametric Bayes

Mark Schmidt

University of British Columbia

Winter 2018

## Previously: Approximate Inference

- We've discussed **approximate inference** in two settings:

- ① Inference in **graphical models** (sum over  $x$  values).

$$E[f(x | w)] = \sum_x f(x)p(x | w)dx.$$

- ② Inference in **Bayesian models** (integrate over posterior values).

$$E[f(\theta)] = \int_{\theta} f(\theta)p(\theta | x)d\theta.$$

- Our previous approach was **Monte Carlo** methods like **MCMC**:
  - Gibbs sampling, Metropolis-Hastings, and so on...
- Alternative class of approximate inference methods is **variational methods**.

## Monte Carlo vs. Variational Inference

Two main strategies for **approximate inference**:

### ① Monte Carlo methods:

- Approximate  $p$  with empirical distribution over samples,

$$p(x) \approx \frac{1}{n} \sum_{i=1}^n \mathcal{I}[x^i = x].$$

- Turns **inference into sampling**.

### ② Variational methods:

- Approximate  $p$  with “closest” **distribution  $q$**  from a tractable family,

$$p(x) \approx q(x).$$

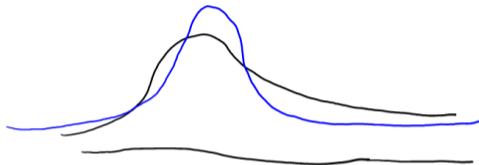
- E.g., Gaussian, independent Bernoulli, or tree UGM.

(or mixtures of these simple distributions)

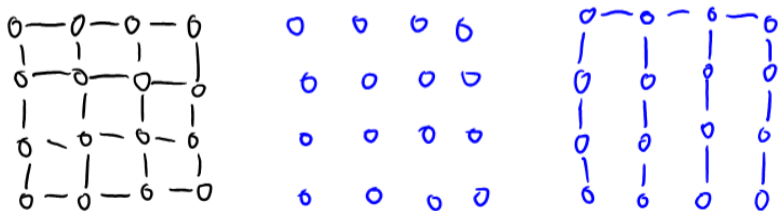
- Turns **inference into optimization**.

## Variational Inference Illustration

- Approximate non-Gaussian  $p$  by a Gaussian  $q$ :



- Approximate loopy UGM by independent distribution or tree-structured UGM:



- Variational methods try to find simple distribution  $q$  that is closest to target  $p$ .
  - This **isn't consistent** like MCMC, but can be **very fast**.

## Laplace Approximation

- A classic variational method is the **Laplace approximation**.

- 1 Find an  $x$  that maximizes  $p(x)$ ,

$$x^* \in \underset{x}{\operatorname{argmin}} \{-\log p(x)\}.$$

- 2 Computer **second-order Taylor expansion** of  $-\log p(x)$  at  $x^*$ .

$$-\log p(x) \approx f(x^*) + \underbrace{\nabla f(x^*)^T}_0 (x - x^*) + \frac{1}{2}(x - x^*)^T \nabla^2 f(x^*) (x - x^*).$$

- 3 Find **Gaussian distribution**  $q$  where  $-\log q(x)$  has **same Taylor expansion**.

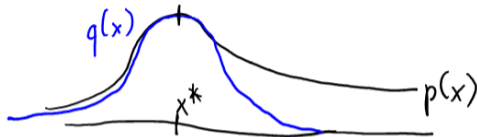
$$-\log q(x) = f(x^*) + \frac{1}{2}(x - x^*)^T \nabla^2 f(x^*) (x - x^*),$$

so  $q$  follows a  $\mathcal{N}(x^*, \nabla^2 f(x^*)^{-1})$  distribution.

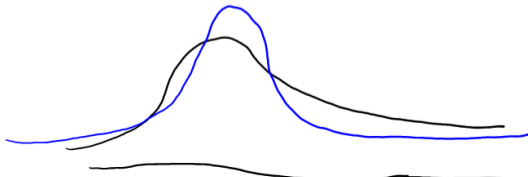
- This is the same approximation used by **Newton's method** in optimization.

## Laplace Approximation

- So **Laplace approximation** replaces complicated  $p(x)$  with Gaussian  $q(x)$ .
  - Centered at mode and agreeing with 1st/2nd-derivatives of log-likelihood:



- Now you only need to compute Gaussian integrals (linear algebra for many  $f$ ).
  - **Very fast**: just solve an optimization (compared to super-slow MCMC).
  - **Bad approximation** if posterior is heavy-tailed, multi-modal, skewed, etc.
- It might **not even give you the "best" Gaussian** approximation:



## Kullback-Leibler (KL) Divergence

- How do we define “closeness” between a distribution  $p$  and  $q$ ?
- A common measure is **Kullback-Leibler (KL)** divergence between  $p$  and  $q$ :

$$KL(p \parallel q) = \sum_x p(x) \log \frac{p(x)}{q(x)}.$$

- Replace sum with integral for continuous families of  $q$  distributions.
- Also called **information gain**: “information lost when  $p$  is approximated by  $q$ ”.
  - If  $p$  and  $q$  are the same, we have  $KL(p \parallel q) = 0$  (no information lost).
  - Otherwise,  $KL(p \parallel q)$  grows as it becomes hard to predict  $p$  from  $q$ .
- Unfortunately, this **requires summing/integrating over  $p$** .
  - The problem we are trying to solve.

## Minimizing Reverse KL Divergence

- Instead of using KL, most variational methods minimize **reverse KL**,

$$\text{KL}(q \parallel p) = \sum_x q(x) \log \frac{q(x)}{p(x)} = \sum_x q(x) \log \frac{q(x)}{\tilde{p}(x)} Z.$$

which just **swaps all  $p$  and  $q$  values** in the definition (KL is not commutative).

- Not intuitive: “how much information is lost when we approximation  $q$  by  $p$ ”.
- But, **reverse KL only needs unnormalized distribution  $\tilde{p}$** ,

$$\begin{aligned} \text{KL}(q \parallel p) &= \sum_x q(x) \log q(x) - \sum_x q(x) \log \tilde{p}(x) + \sum_x q(x) \log(Z) \\ &= \sum_x q(x) \log \frac{q(x)}{\tilde{p}(x)} + \underbrace{\log(Z)}_{\text{const. in } q}. \end{aligned}$$

- By non-negativity of KL this also gives a **lower bound on  $\log(Z)$** .



## Coordinate Optimization: Mean Field Approximation

- This “**variational lower bound**” still seems difficult to work with.
  - But with appropriate  $q$  we can do **coordinate optimization**.
- Consider minimizing reverse KL with **independent**  $q$ ,

$$q(x) = \prod_{j=1}^d q_j(x_j),$$

where we choose  $q$  to be conjugate (usually discrete or Gaussian).

- If we fix  $q_{-j}$  and optimize the functional  $q_j$  we obtain (see Murphy's book)

$$q_j(x_j) \propto \exp \left( \mathbb{E}_{q_{-j}} [\log \tilde{p}(x)] \right),$$

which we can use to update  $q_j$  for a particular  $j$ .

## Coordinate Optimization: Mean Field Approximation

- Each iteration we choose a  $j$  and set  $q$  based on mean (of neighbours),

$$q_j(x_j) \propto \exp \left( \mathbb{E}_{q_{-j}} [\log \tilde{p}(x)] \right).$$

- This improves the (non-convex) reverse KL on each iteration.
- Applying this update is called:
  - **Mean field** method (graphical models).
  - **Variational Bayes** (Bayesian inference).

## 3 Coordinate-Wise Algorithms

- **ICM** is a coordinate-wise method for approximate decoding:
  - Choose a coordinate  $i$  to update.
  - Maximize  $x_i$  keeping other variables fixed.
- **Gibbs sampling** is a coordinate-wise method for approximate **sampling**:
  - Choose a coordinate  $i$  to update.
  - **Sample**  $x_i$  keeping other variables fixed.
- **Mean field** is a coordinate-wise method for approximate **marginalization**:
  - Choose a coordinate  $i$  to update.
  - **Update**  $q_i(x_i)$  keeping other variables fixed ( $q_i(x_i)$  approximates  $p_i(x_i)$ ).  
for all  $x_i$

### 3 Coordinate-Wise Algorithms

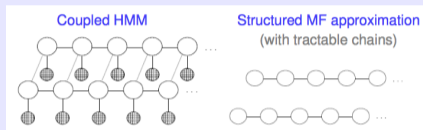
- Consider a **pairwise UGM**:

$$p(x_1, x_2, \dots, x_d) \propto \left( \prod_{i=1}^d \phi_i(x_i) \right) \left( \prod_{(i,j) \in E} \phi_{ij}(x_i, x_j) \right),$$

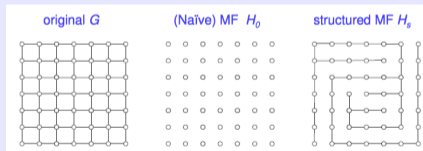
- ICM** for **updating a node  $i$**  with 2 neighbours ( $j$  and  $k$ ).
  - 1 Compute  $M_i(x_i) = \phi_i(x_i)\phi_{ij}(x_i, x_j)\phi_{ik}(x_i, x_k)$  for all  $x_i$ .
  - 2 Set  $x_i$  to the largest value of  $M_i(x_i)$ .
- Gibbs** for **updating a node  $i$**  with 2 neighbours ( $j$  and  $k$ ).
  - 1 Compute  $M_i(x_i) = \phi_i(x_i)\phi_{ij}(x_i, x_j)\phi_{ik}(x_i, x_k)$  for all  $x_i$ .
  - 2 Sample  $x_i$  proportional to  $M_i(x_i)$ .
- Mean field** for **updating a node  $i$**  with 2 neighbours ( $j$  and  $k$ ).
  - 1 Compute  $M_i(x_i) = \exp \left( \sum_{x_j} q_j(x_j) \log \phi_{ij}(x_i, x_j) + \sum_{x_k} q_k(x_k) \log \phi_{ik}(x_i, x_k) \right)$ .
  - 2 Set  $q_i(x_i)$  proportional to  $\phi_i(x_i)M_i(x_i)$ .

# Structure Mean Field

- Common variant is **structured mean field**:  $q$  function includes some of the edges.



<http://courses.cms.caltech.edu/cs155/slides/cs155-14-variational.pdf>



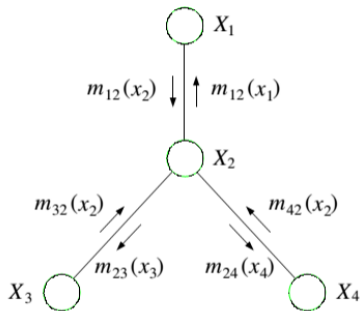
<http://courses.cms.caltech.edu/cs155/slides/cs155-14-variational.pdf>

- Original LDA article proposed a structured mean field approximation.

## Previously: Belief Propagation

- We've discussed **belief propagation** for forest-structured UGMs.

(undirected graphs with no loops, which must be pairwise)



<https://www.quora.com/>

Probabilistic-graphical-models-what-are-the-relationships-between-sum-product-algorithm-belief-propagation-and-junction-tree-

- Defines “messages” that can be sent along each edge.
  - Generalizes forward-backward algorithm.

## Loopy Belief Propagation

- In pairwise UGM, belief propagation “message” from parent  $p$  to child  $c$  is given by

$$M_{pc}(x_c) \propto \sum_{x_p} \phi_i(x_p) \phi_{pc}(x_p, x_c) M_{jp}(x_p) M_{kp}(x_p),$$

assuming that parent  $p$  has parents  $j$  and  $k$ .

- We get marginals by multiplying all incoming messages with local potentials.
- **Loopy belief propagation**: a “hacker” approach to approximate marginals:
  - Choose an edge  $ic$  to update.
  - Update messages  $M_{ic}(x_c)$  keeping all other messages fixed.
  - Repeat until “convergence”.
    - We approximate marginals by multiplying all incoming messages with local potentials.
- Empirically much better than mean field, we’ve spent 20 years figuring out why.

## Discussion of Loopy Belief Propagation

- Loopy BP decoding is used for “error correction” in WiFi and Skype.
  - Called “turbo codes” in information theory.
- Loopy BP is **not optimizing an objective** function.
  - Convergence of loopy BP is hard to characterize: does not converge in general.
- If it converges loopy BP finds fixed point of “Bethe free energy”:
  - Better approximation than mean field, but not a lower/upper bound.
- Recent works give convex variants that upper bound  $Z$ .
  - **Tree-reweighted belief propagation.**
  - Variations that are guaranteed to converge.
- Messages only have closed-form update for conjugate models.
  - Can approximate non-conjugate models using **expectation propagation.**



## Convex Relaxations

- I've overviewed the “classic” view of variational methods that they minimize KL.
- Modern view: write exact inference as constrained convex optimization (bonus).
  - Different methods correspond to different function/constraints approximations.
  - There are also [convex relaxations](#) that approximate with linear programs.
- For an overview of this and all things variational, see:  
`people.eecs.berkeley.edu/~wainwrig/Papers/WaiJor08_FTML.pdf`

## Variational vs. Monte Carlo

- Monte Carlo vs. variational methods:
  - Variational methods are typically **more complicated**.
  - Variational methods are **not consistent**.
    - $q$  does not converge to  $p$  if we run the algorithm forever.
  - But variational methods typically give **better approximation for the same time**.
    - Although **MCMC is easier to parallelize**.
  - Variational methods typically have similar cost to MAP.
- Combinations of variational inference and stochastic methods:
  - **Stochastic variational inference (SVI)**: use stochastic gradient to speed up variational methods.
  - **Variational MCMC**: use Metropolis-Hastings where variational  $q$  sometimes makes proposals.

# Outline

- 1 Variational Inference
- 2 Non-Parametric Bayes
- 3 GANs and VAEs

## Stochastic Processes and Non-Parametric Bayes

- A **stochastic process** is an infinite collection of random variables  $\{x^i\}$ .
- **Non-parametric Bayesian** methods use priors defined on stochastic processes:
  - Allows extremely-flexible prior, and posterior **complexity grows with data size**.
  - Typically set up so that samples from posterior are finite-sized.
- The two most common priors are **Gaussian processes** and **Dirichlet processes**:
  - Gaussian processes define prior on space of functions (universal approximators).
  - Dirichlet processes define prior on space of probabilities (without fixing dimension).

# Gaussian Processes

- Recall the partitioned form of a multivariate Gaussian

$$\mu = [\mu_x, \mu_y], \quad \Sigma = \begin{bmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{bmatrix},$$

and in this case the marginal  $p(x)$  is a  $\mathcal{N}(\mu_x, \Sigma_{xx})$  Gaussian.

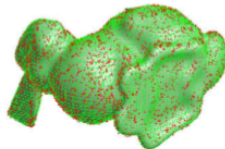
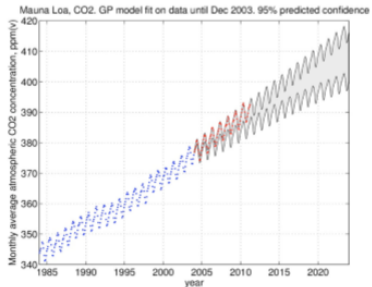
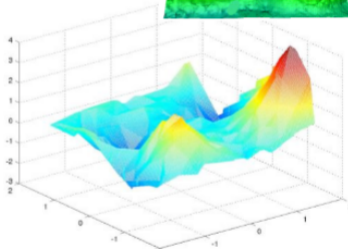
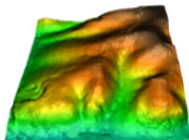
- Generalization of this to **infinite set of variables** is **Gaussian processes** (GPs):
  - Any finite set from collection follows a Gaussian distribution.

# Gaussian Processes

To date kriging has been used in a variety of disciplines, including the following:

- Environmental science<sup>[5]</sup>
- Hydrogeology<sup>[6][7][8]</sup>
- Mining<sup>[9][10]</sup>
- Natural resources<sup>[11][12]</sup>
- Remote sensing<sup>[13]</sup>
- Real estate appraisal<sup>[14][15]</sup>

and many others.



# Gaussian Processes

- GPs are specified by a **mean function**  $m$  and **covariance function**  $k$ ,

$$m(x) = \mathbb{E}[f(x)], \quad k(x, x') = \mathbb{E}[(f(x) - m(x))(f(x') - m(x'))^T].$$

- We write that

$$f(x) \sim \text{GP}(m(x), k(x, x')),$$

- As an example, we could have a zero-mean and linear covariance GP,

$$m(x) = 0, \quad k(x, x') = x^T x'.$$

## Regression Models as Gaussian Processes

- As an example, **predictions made by linear regression** with Gaussian prior

$$f(x) = \phi(x)^T w, \quad w \sim \mathcal{N}(0, \Sigma),$$

are a Gaussian process with **mean function**

$$\mathbb{E}[f(x)] = \mathbb{E}[\phi(x)^T w] = \phi(x)^T \mathbb{E}[w] = 0.$$

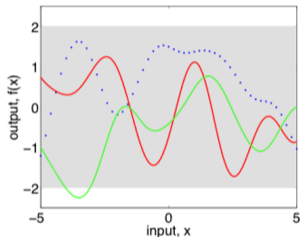
and **covariance function**

$$\mathbb{E}[f(x)f(x')^T] = \phi(x)^T \mathbb{E}[ww^T] \phi(x') = \phi(x)^T \Sigma \phi(x') = k(x, x').$$

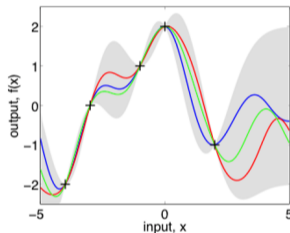


## Gaussian Process Model Selection

- We can view a Gaussian process as a **prior distribution over smooth functions**.



(a), prior



(b), posterior

- Most common choice of covariance is RBF.
- Is this the same as using RBF kernels or the RBFs as the bases?
  - Yes, this is **Bayesian linear regression plus the kernel trick**.

## Gaussian Process Model Selection

- So why do we care?
  - We can get estimate of uncertainty in the prediction.
  - We can use marginal likelihood to learn the kernel/covariance.
- Write kernel in terms of parameters, use empirical Bayes to learn kernel.
- Hierarchical approach: put a hyper-prior of types of kernels.
- Application: Bayesian optimization of non-convex functions:
  - Gradient descent is based on a Gaussian (quadratic) approximation of  $f$ .
  - Bayesian optimization is based on a Gaussian process approximation of  $f$ .
    - Can approximate non-convex functions.

## Dirichlet Process

- Recall the basic mixture model:

$$p(x | \theta) = \sum_{c=1}^k \pi_c p(x | \theta_c).$$

- Non-parametric Bayesian methods allow us to consider **infinite mixture model**,

$$p(x | \theta) = \sum_{c=1}^{\infty} \pi_c p(x | \theta_c).$$

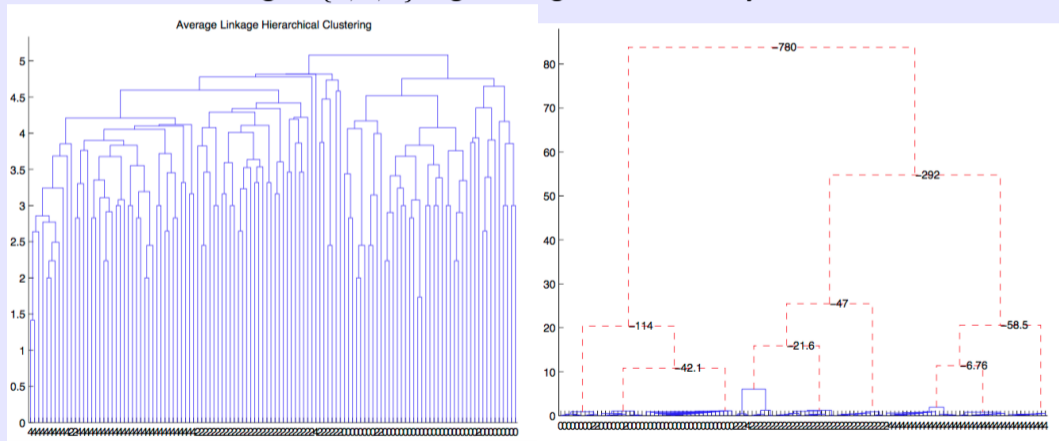
- Common choice for prior on  $\pi$  values is **Dirichlet process**:
  - Also called “Chinese restaurant process” and “stick-breaking process”.
  - For finite datasets, only a fixed number of clusters have  $\pi_c \neq 0$ .
  - But **don't need to pick number of clusters**, grows with data size.

## Dirichlet Process

- Gibbs sampling in Dirichlet process mixture model in action:  
<https://www.youtube.com/watch?v=0Vh7qZY9sPs>
- We could alternately put a prior on  $k$ :
  - “Reversible-jump” MCMC can be used to sample from models of different sizes.
    - AKA “trans-dimensional” MCMC.
- There a variety of interesting variations on Dirichlet processes
  - Beta process (“Indian buffet process”).
  - Hierarchical Dirichlet process,.
  - Polya trees.
  - Infinite hidden Markov models.

# Bayesian Hierarchical Clustering

- Hierarchical clustering of  $\{0, 2, 4\}$  digits using classic and Bayesian method:

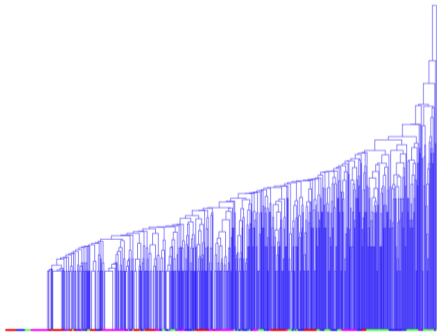


<http://www2.stat.duke.edu/~kheller/bhcnew.pdf> (y-axis represents distance between clusters)

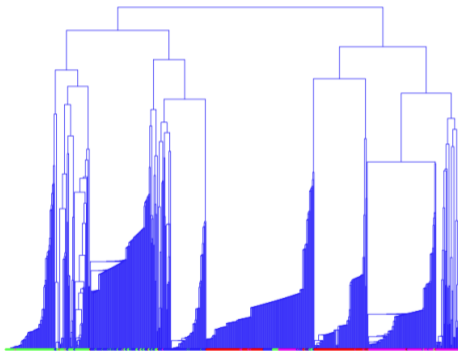
# Bayesian Hierarchical Clustering

- Hierarchical clustering of newgroups using classic and Bayesian method:

4 Newsgroups Average Linkage Clustering



4 Newsgroups Bayesian Hierarchical Clustering



<http://www2.stat.duke.edu/~kheller/bhcnew.pdf> (y-axis represents distance between clusters)

## Summary of Part 1

- **Variational methods** approximate  $p$  with a simpler distribution  $q$ .
  - **Mean field** approximation minimizes KL divergence with independent  $q$ .
  - **Loopy belief propagation** is a heuristic that often works well.
- **Non-Parametric Bayes** puts probabilities over infinite spaces.
  - Gaussian processes are priors over continuous functions.
  - Dirichlet processes are priors over probability mass functions.
- Part 2: new generative deep learning methods.

# Variational Inference: Constrained Optimization View

- Modern view of **variational inference**:
  - Formulate inference problem as constrained optimization.
  - **Approximate the function or constraints** to make it easy.



## Exponential Families and Cumulant Function

- We will again consider log-linear models:

$$P(X) = \frac{\exp(w^T F(X))}{Z(w)},$$

but view them as **exponential family distributions**,

$$P(X) = \exp(w^T F(X) - A(w)),$$

where  $A(w) = \log(Z(w))$ .

- Log-partition  $A(w)$  is called the **cumulant function**,

$$\nabla A(w) = \mathbb{E}[F(X)], \quad \nabla^2 A(w) = \mathbb{V}[F(X)],$$

which implies convexity.

## Convex Conjugate and Entropy

- The **convex conjugate** of a function  $A$  is given by

$$A^*(\mu) = \sup_{w \in \mathcal{W}} \{\mu^T w - A(w)\}.$$

- E.g., in A3 we did this for logistic regression:

$$A(w) = \log(1 + \exp(w)),$$

implies that  $A^*(\mu)$  satisfies  $w = \log(\mu) / \log(1 - \mu)$ .

- When  $0 < \mu < 1$  we have

$$\begin{aligned} A^*(\mu) &= \mu \log(\mu) + (1 - \mu) \log(1 - \mu) \\ &= -H(p_\mu), \end{aligned}$$

**negative entropy of binary distribution with mean  $\mu$ .**

- If  $\mu$  does not satisfy boundary constraint, sup is  $\infty$ .

## Convex Conjugate and Entropy

- More generally, if  $A(w) = \log(Z(w))$  then

$$A^*(\mu) = -H(p_\mu),$$

subject to boundary constraints on  $\mu$  and constraint:

$$\mu = \nabla A(w) = \mathbb{E}[F(X)].$$

- Convex set satisfying these is called **marginal polytope**  $\mathcal{M}$ .
- If  $A$  is convex (and LSC),  $A^{**} = A$ . So we have

$$A(w) = \sup_{\mu \in \mathcal{U}} \{w^T \mu - A^*(\mu)\}.$$

and when  $A(w) = \log(Z(w))$  we have

$$\log(Z(w)) = \sup_{\mu \in \mathcal{M}} \{w^T \mu + H(p_\mu)\}.$$

- We've written **inference as a convex optimization problem**.

## Bonus slide: Maximum Likelihood and Maximum Entropy

- The **maximum likelihood** parameters  $w$  satisfy:

$$\begin{aligned}
 & \min_{w \in \mathbb{R}^d} -w^T F(D) + \log(Z(w)) \\
 &= \min_{w \in \mathbb{R}^d} -w^T F(D) + \sup_{\mu \in \mathcal{M}} \{w^T \mu + H(p_\mu)\} \quad (\text{convex conjugate}) \\
 &= \min_{w \in \mathbb{R}^d} \sup_{\mu \in \mathcal{M}} \{-w^T F(D) + w^T \mu + H(p_\mu)\} \\
 &= \sup_{\mu \in \mathcal{M}} \{ \min_{w \in \mathbb{R}^d} -w^T F(D) + w^T \mu + H(p_\mu) \} \quad (\text{convex/concave})
 \end{aligned}$$

which is  $-\infty$  unless  $F(D) = \mu$  (e.g., maximum likelihood  $w$ ), so we have

$$\begin{aligned}
 & \min_{w \in \mathbb{R}^d} -w^T F(D) + \log(Z(w)) \\
 &= \max_{\mu \in \mathcal{M}} H(p_\mu),
 \end{aligned}$$

subject to  $F(D) = \mu$ .

- Maximum likelihood**  $\Rightarrow$  **maximum entropy + moment constraints.**

## Difficulty of Variational Formulation

- We wrote inference as a convex optimization:

$$\log(Z) = \sup_{\mu \in \mathcal{M}} \{w^T \mu + H(p_\mu)\},$$

- Did this make anything easier?
  - Computing entropy  $H(p_\mu)$  seems as hard as inference.
  - Characterizing marginal polytope  $\mathcal{M}$  becomes hard with loops.
- Practical variational methods:
  - Work with approximation to marginal polytope  $\mathcal{M}$ .
  - Work with approximation/bound on entropy  $A^*$ .
- Notation trick: we put everything “inside”  $w$  to discuss general log-potentials.

## Mean Field Approximation

- Mean field approximation assumes

$$\mu_{ij,st} = \mu_{i,s}\mu_{j,t},$$

for all edges, which means

$$p(x_i = s, x_j = t) = p(x_i = s)p(x_j = t),$$

and that **variables are independent**.

- Entropy is simple under mean field approximation:

$$\sum_X p(X) \log p(X) = \sum_i \sum_{x_i} p(x_i) \log p(x_i).$$

- Marginal polytope is also simple:

$$\mathcal{M}_F = \left\{ \mu \mid \mu_{i,s} \geq 0, \sum_s \mu_{i,s} = 1, \mu_{ij,st} = \mu_{i,s}\mu_{j,t} \right\}.$$

## Entropy of Mean Field Approximation

- Entropy form is from distributive law and probabilities sum to 1:

$$\begin{aligned}\sum_X p(X) \log p(X) &= \sum_X p(X) \log \left( \prod_i p(x_i) \right) \\ &= \sum_X p(X) \sum_i \log(p(x_i)) \\ &= \sum_i \sum_X p(X) \log p(x_i) \\ &= \sum_i \sum_X \prod_j p(x_j) \log p(x_i) \\ &= \sum_i \sum_X p(x_i) \log p(x_i) \prod_{j \neq i} p(x_j) \\ &= \sum_i \sum_{x_i} p(x_i) \log p(x_i) \sum_{x_j \mid j \neq i} \prod_{j \neq i} p(x_j) \\ &= \sum_i \sum_{x_i} p(x_i) \log p(x_i).\end{aligned}$$

## Mean Field as Non-Convex Lower Bound

- Since  $\mathcal{M}_F \subseteq \mathcal{M}$ , yields a **lower bound** on  $\log(Z)$ :

$$\sup_{\mu \in \mathcal{M}_F} \{w^T \mu + H(p_\mu)\} \leq \sup_{\mu \in \mathcal{M}} \{w^T \mu + H(p_\mu)\} = \log(Z).$$

- Since  $\mathcal{M}_F \subseteq \mathcal{M}$ , it is an **inner approximation**:

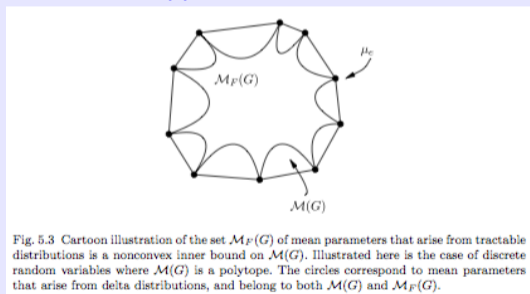


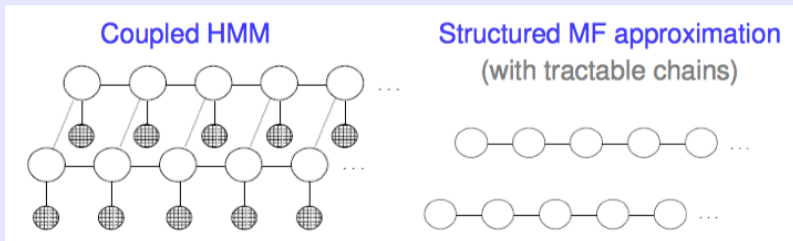
Fig. 5.3 Cartoon illustration of the set  $\mathcal{M}_F(G)$  of mean parameters that arise from tractable distributions is a nonconvex inner bound on  $\mathcal{M}(G)$ . Illustrated here is the case of discrete random variables where  $\mathcal{M}(G)$  is a polytope. The circles correspond to mean parameters that arise from delta distributions, and belong to both  $\mathcal{M}(G)$  and  $\mathcal{M}_F(G)$ .

- Constraints  $\mu_{ij,st} = \mu_{i,s}\mu_{j,t}$  make it **non-convex**.
- Mean field algorithm is **coordinate descent** on  $w^T \mu + H(p_\mu)$  over  $\mathcal{M}_F$ .



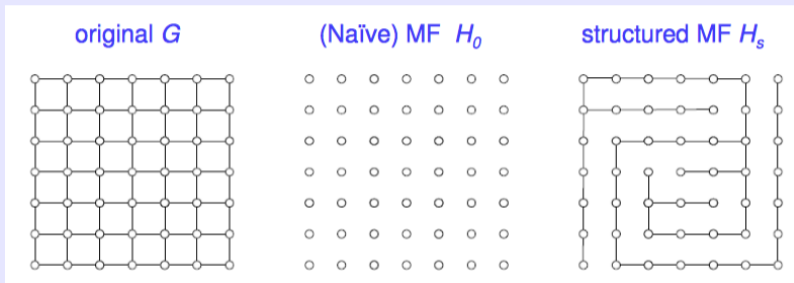
## Discussion of Mean Field and Structured MF

- Mean field is weird:
  - Non-convex approximation to a convex problem.
  - For learning, we want **upper** bounds on  $\log(Z)$ .
- Structured mean field:
  - Cost of computing entropy is similar to cost of inference.
  - Use a subgraph where we can perform exact inference.



## Structured Mean Field with Tree

- More edges means better approximation of  $\mathcal{M}$  and  $H(p_\mu)$ :



<http://courses.cms.caltech.edu/cs155/slides/cs155-14-variational.pdf>

- Fixed points of loopy correspond to using “Bethe” approximation of entropy and “local polytope” approximation of “marginal polytope”.
- You can design better variational methods by constructing better approximations.