We did a var ~~dist~~ ~~is~~ approx with a
factored dist — BUT we could do others:

eg. Q: a tree structured dist?

$$E_Q = -E_Q \log P - H_Q$$

i) We can compute $H_Q$

recall tree:

$$q(x_1 \cdots x_N) = \left[ \prod_{i \in V} q_i \right] \left[ \prod_{ij \in \varepsilon} \frac{q_{ij}}{q_i q_j} \right]$$

$$= \frac{\prod_{ij \in \varepsilon} q_{ij}}{\prod_{i \in V} q_i^{(d_i - 1)}}$$

degree of edge

So $$H_Q = -\sum_{\substack{values \\ of pairs}} \left[ \sum_{ij \in \varepsilon} q_{ij} \log q_{ij} \right] + \sum_{values} \left[ (d_i - 1) \sum_{i \in V} q_i \log q_i \right]$$

tractable:

We can also compute

$$-E_Q \log P$$

recall:

$$P(H|x) = \frac{1}{Z} \exp\left[ -\sum_{ij} \Theta_{ij}(H_i H_j) - \sum_i \Theta_i(H_i) \right]$$

$$-\log P = \log Z + \sum_{ij} \Theta_{ij}(H_i, H_j) + \sum_i \Theta_i(H_i)$$

Constant – not a prob, cause we've minimizing!

$x_i$ in other hand ω notes ; sorry! vars;

So

$$-E_Q \log P = \log Z + \sum_{\substack{\text{values} \\ \text{of pairs}}} \left[ \sum_{ij \in \mathcal{E}} q_{ij} \, t_{ij} \right]$$

$$+ \sum_{\text{values}} \left[ (d_i - 1) \sum_{i \in V} q_i \, t_i \right]$$

Here I used the change of var

$$t_{ij} = \Theta_{ij} + \Theta_i + \Theta_j$$

$$t_i = \Theta_i$$

to get an expression that looks like entropy.

Now: we want to min

$$E_Q \log Q - E_Q \log P$$

for $Q$ some fixed tree. (which we chose)

Recall the $q_{ij}$, $q_i$, $q_j$ are __marginals__

So we must

$$\min \sum_{values} \left[ \sum_{ij \in \varepsilon} q_{ij} \left[ \log q_{ij} + t_{ij} \right] \right]$$

$$+ \sum_{values} \left[ \sum_i (d_i - 1) \cdot q_i \left[ \log q_i + t_i \right] \right]$$

st. $\sum_i q_{ij} = q_j$ ; $\sum_j q_{ij} = q_i$ ; $\sum q_i = 1$

( cause these are marginals )

write

$\lambda_{\varepsilon_i}$     for LM's assoc. with     $\sum_j q_{ij} = q_i$

$\lambda_{\varepsilon_j}$                 $\sum_i q_{ij} = q_j$

$\lambda_{v_i}$          LM           $\sum_i q_i = 1$

( Notice   $\lambda_{\varepsilon_i}$   is   a   <u>vector</u>

                         <u>scalar</u>

        $\lambda_{v_i}$

write lagrangian $L$.

at stationary point
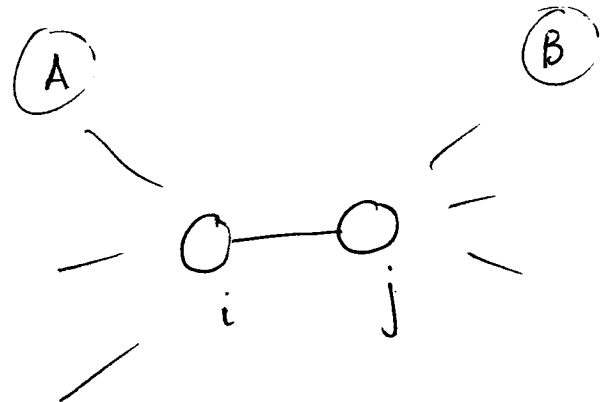
$$\left[\frac{\partial L}{\partial q_{ij}}\right]_{uv} = 0 = \left[\log q_{ij}\right]_{uv} + 1 + \left[\tau_{ij}\right]_{uv} + \sum\left[\lambda_{\varepsilon_i}\right]_u + \sum\left[\lambda_{\varepsilon_j}\right]_v$$

overbrace: all incoming edges to $i$

$\underset{u,v'th\ entry\ in\ \underline{table}}{\curvearrowleft}$

all inc. to $j$

SO

$$\left[q_{ij}\right]_{uv} \propto \left[e^{-\tau_{ij}}\right]_{uv} \cdot \left[e^{\sum\lambda_{\varepsilon_i}}\right]_u \cdot \left[e^{\sum\lambda_{\varepsilon_j}}\right]_v$$

compare        with        B.P. eqns

Ⓐ                        Ⓑ

$$\left[ q_{ij} \right]_{uv} \propto \left[ \psi_{ij} \, \varphi_i \, \varphi_j \right]_{uv} \cdot \left[ \prod_{\substack{all \; inc \\ to \, i}} M_{ia} \right]_{u} \cdot \left[ \prod_{\substack{all \; inc \\ to \, j}} M_{jb} \right]_{v}$$

## Conclusion

- Messages        =        log $h \cdot M's$

## Two outcomes :

1) We can fit a variational model
   of a single tree (MP as above)

2) loopy BP "like" $\not{z}$ fitting var m
                              of tree
                              without worrying

Now what is happening in terms of M.P. ?

→ fix a tree.

→ Interpret MP $\equiv$ convex hull of all states that can arise in this repn of G.M.

→ ~~we must have that~~
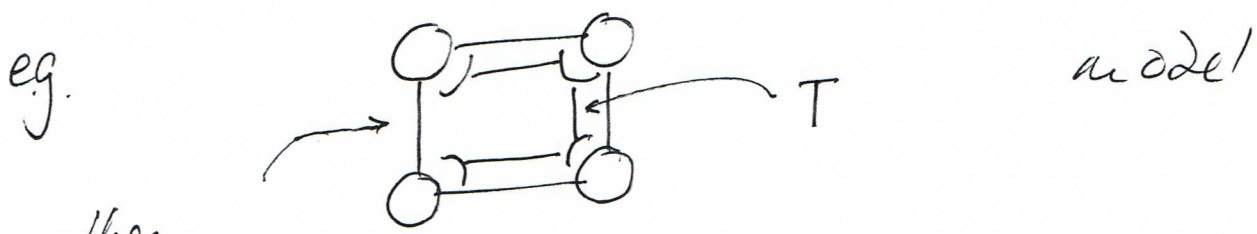
Call the polytope that satisfies

$$\sum_j q_{ij} = q_i \quad ; \quad \sum_{i} q_{ij} = q_j \quad ; \quad \sum q_i = 1$$

the **Local Polytope** $= Lp.$

→ Notice

$$MP \subset Lp.$$

Notice that for choice of tree $T$, a lot of $q_{ij} = q_i \cdot q_j$ (cause the vars are ~~indep~~ cond indep given parents

eg.



$T$     model

then in $T$, these two have $q_{ij} = q_i \cdot q_j$

so we are finding the $q_{ij}, q_i, q_j$

that

- are in $LP$.
- meet indep constraints implied by $T$
- minimize

$$-E_Q \log P \quad + E_Q \log Q$$

(then extract info from $q$ ).

# Extracting info from Q.

- if were lucky, $q_{ij}$ are integer.
  (might be a vert of $M_P$!).
  $\rightarrow$ nothing to do

- else, it's a tree; $\longrightarrow$ <u>max product</u>

# Idea:

rather than

$$\min \quad -E_Q \log P + E_Q \log Q$$

for $Q$ a tree,

do it for $Q \in L_P$.

$\longrightarrow$ <u>How do we get $E_Q \log Q$ ?</u>

Here is one strategy

- Drop the tree
- ~~compu~~ fit $q$ by using the expression for $E_Q \log Q$ that came from tree

$$E_Q \log Q \simeq -\sum_{\substack{edges \\ values \\ of\ pairs}} \sum \left[ q_{ij} \log q_{ij} \right] + \sum_{\substack{\cancel{values} \\ verts}} \sum_{values} \left[ (d_i - 1) \sum_j q_i \log \right.$$

notice I flipped order

$$= H_Q^2$$

This is**n't** the true exp. for $E_Q \log Q$, **but** its easy to eval.

loopy b.p. $= \min E_Q \log P - H_Q^2$
$\text{s.t. } Q \in LP.$

Notice the form of the cost function

$$E_Q \log P \sim H^? Q.$$

↑

linear in Q.

Some property of Q that approx entropy.
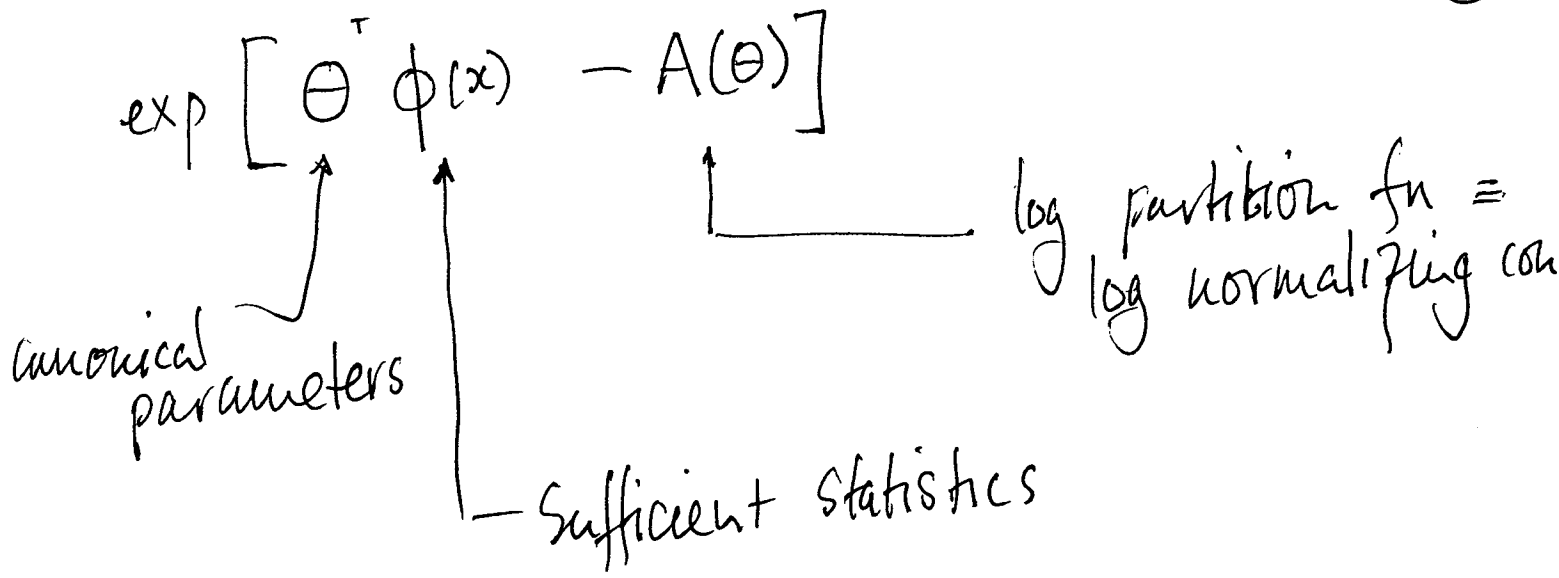
Notice also that we're identifying points in LP (or MP) with probability distributions. It turns out that we can formalize this

## the exponential family:

any p.d. that is written as

$$P(x) = \exp \left[ \Theta^T \phi(x) - A(\Theta) \right]$$

(for our purposes — other )

$$\exp\left[\theta^{\top}\phi(x) - A(\theta)\right]$$

log partition fn $\equiv$ log normalizing con

canonical parameters

— Sufficient statistics

We will confine attn to case where:

- $\phi(x)$ are linearly independent (no real issue here, just creates a lot of if's, and's, bu

- $\theta$ is such that

$$A(\theta) = \log\int \exp\left[\theta^{\top}\phi(x)\right] < \infty$$

&

# Examples:

### 1D Normal Dist:

$$\exp\left[\cancel{(\alpha ,\ \beta)}\begin{pmatrix}\alpha\\\beta\end{pmatrix}^{T}(x^2, x)\ -A(\Theta)\right]$$

here $\qquad \alpha < 0 \ ; \qquad$ std $= \dfrac{-1}{2\alpha}.$

$$\text{mean} = \left(\dfrac{1}{2\alpha}\right)^{2}\cdot\beta .$$
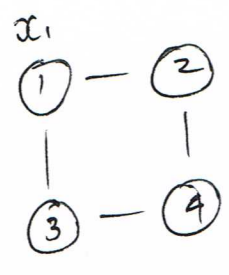
### Multi D ND:

follows easily.

### Poisson dist:
, recall this is dist on ~~the~~ non-neg integers)

.

$$p(K) = \lambda^{K}\ \dfrac{e^{-\lambda}}{K!}$$

rate/intensity.

$$p(K) = \exp\left[\alpha \cdot K - A(\alpha)\right]$$

$$\alpha = \log \lambda \quad , \quad \text{etc.}$$

## Discrete MRF :

$x_1$



etc.

$\mathbb{1}_{\phi}(x_i)$ is 1-hot vector

$\mathbb{1}_{\phi}(x_i, x_j)$ is 1-hot table, straightened into vector

$$p(x) = \exp\left[\Theta^T \begin{bmatrix} \mathbb{1}(x_1) \\ \vdots \\ \vdots \\ \mathbb{1}(x_i, x_j) \\ \vdots \\ \vdots \end{bmatrix} - A(\Theta)\right]$$

$A(\Theta)$ is extremely interesting

$$\nabla_\Theta A = \nabla_\Theta \left[ \log \int e^{\Theta^T \phi} \cdot dx \right]$$

$$= \frac{1}{\int e^{\Theta^T \phi} dx} \cdot \int \phi e^{\Theta^T \phi} dx$$

$$= E_p[\phi]$$

(recall - we've seen something like this before when talking about max-likelihood = max entropy.

Now assume we have some $\phi$
(likely indicator fns in our case).
We can define

$$\Lambda : \Theta \longrightarrow M \longleftarrow \text{marginal polytope}$$

$$\Lambda(\Theta) = E_\Theta[\phi]$$

{ this is in M, cause M is all possible expects of

Thm:

$$\underline{\Delta \qquad \text{is} \qquad 1\text{-}1} \qquad \left(\begin{array}{c}\text{assuming } \phi \text{ are} \\ \text{linearly indep}\end{array}\right).$$

(proof in Wainwright — mildly technical).

Now we want to consider $\underline{\text{dual}}$ of $A(\theta)$.

$$A^*(\mu) \quad = \quad \sup_{\theta \in \Theta} \left[ \langle \mu, \theta \rangle - A(\theta) \right]$$

Note this is a function of $\mu$.

known as a $\underline{\text{conjugate dual}}$

Why is $\wedge$ 1-1 ?

(sketch of proof - details in Wainwright)

$A(\theta)$ is convex

we must show for any $\mu$, there
is some $\theta$ st $E_{p(x;\theta)}[\phi] = \mu$.

BUT $\qquad E_{p(x;\theta)}[\phi] = \nabla_{\theta} A(\theta)$.

under very mild conditions, map

$x \to \frac{df}{dx}$ is 1-1 for $x$ convex

— proof by drawing !

$A(\Theta)$ is a <u>convex</u> function

of $\Theta$.

<u>recall</u> $\quad \dfrac{\partial A}{\partial \Theta_i} = e^{-A} \cdot \int e^{\Theta^T \varphi} \cdot \varphi_i \, dx$

so $\quad \dfrac{\partial^2 A}{\partial \Theta_i \partial \Theta_j} = e^{-A} \cdot \int e^{\Theta^T \varphi} \varphi_i \varphi_j \, dx$

$$- \left[ e^{-A} \cdot e^{-A} \int e^{\Theta^T \varphi} \varphi_j \, dx \right] \left[ \int e^{\Theta^T \varphi} \varphi_j \, dx \right]$$

$$= E_p \left[ \varphi_i \, \varphi_j \right] - E_p \left[ \varphi_i \right] E_p \left[ \varphi_j \right]$$

$$= cov \left( \varphi_i , \varphi_j \right)$$

so $\quad H_A = cov \, mat \left[ \varphi \right]$

this is positive definite under
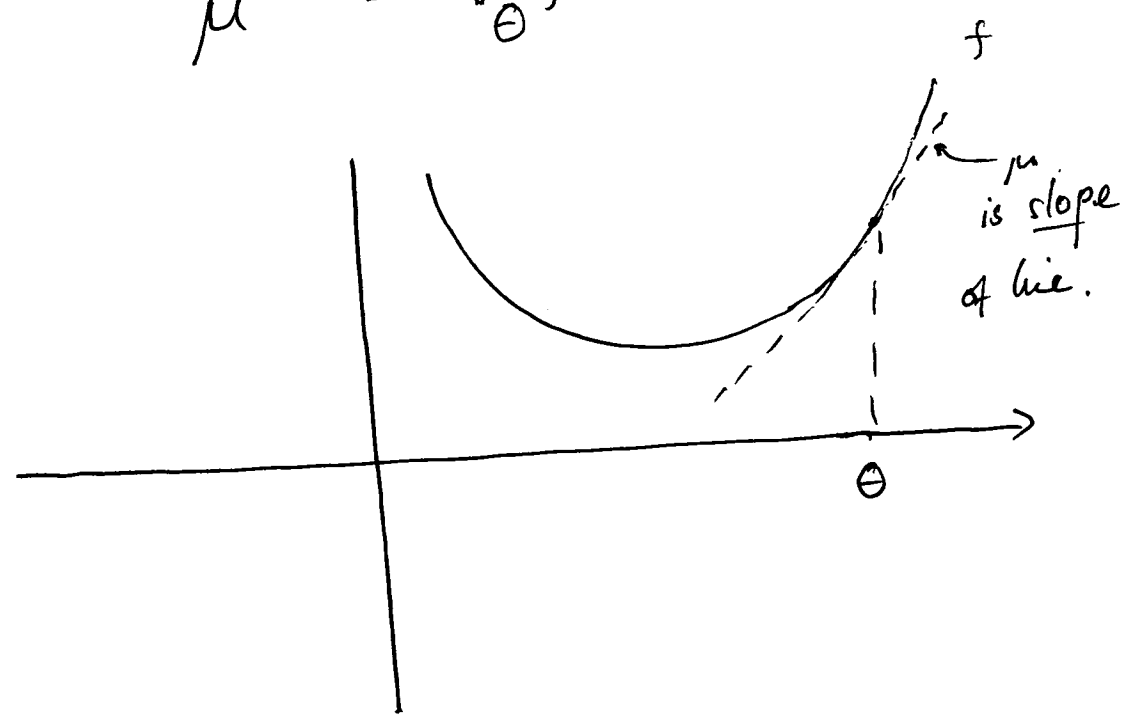our conditions $\left( \varphi \text{ linearly indep.} \right)$

Now consider conjugate dual.

$$f^*(\mu) = \sup_{\theta} \left[ \langle \mu, \theta \rangle - f(\theta) \right]$$

for convex $f$.

- assume $f$ <u>differentiable</u>
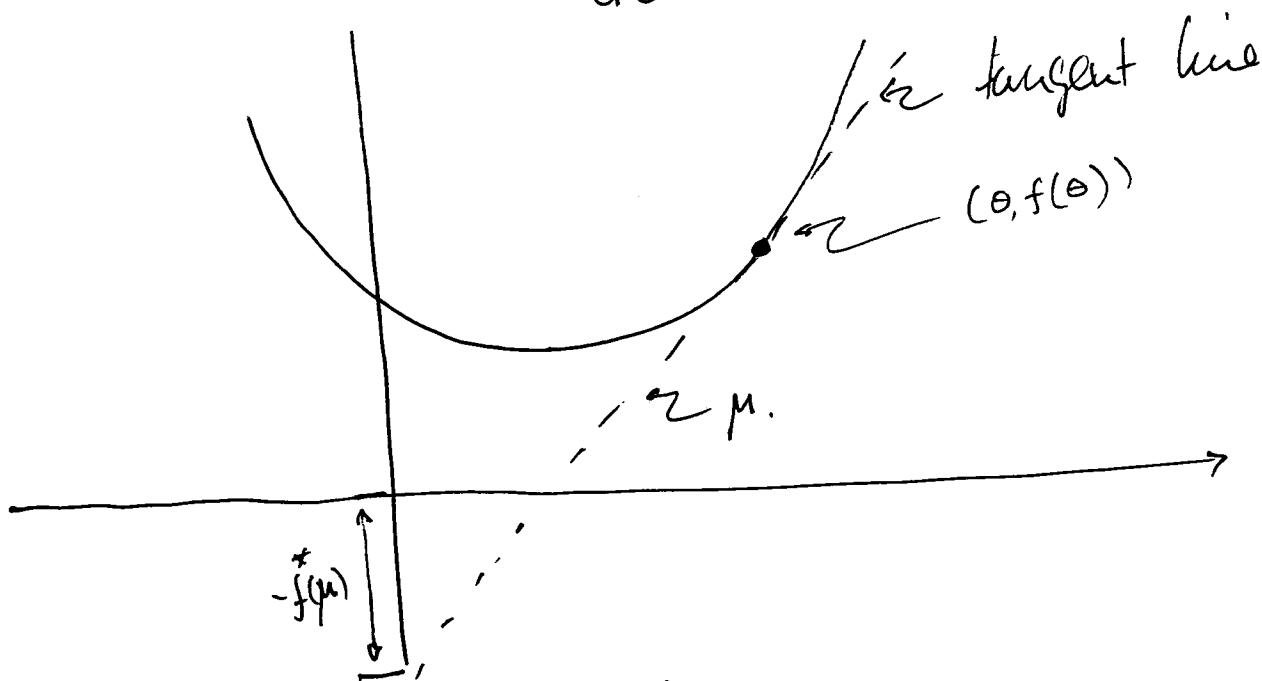
then $\mu = \nabla_\theta f$

$\mu$ is slope of line.

# Common visualization

- choose $\mu \rightarrow$ what is $f^*(\mu)$ ?

- consider $\theta$ such that

$$\mu = \frac{df}{d\theta}$$



tangent line

$(\theta, f(\theta))$

$\mu$

$-f^*(\mu)$

tangent line has slope $\mu$.

passes through $(\theta, f(\theta))$

$$\therefore \quad y = \mu x + (f(\theta) - \mu \theta)$$
$$\therefore \quad \text{at } x = 0, \quad y = -f^*(\mu)$$

for $f$ convex, $f^*$ is convex.

Show for $f \in C^2$, but generally true.

Proof:
$$f^*(p) = \sup_x [xp - f(x)]$$

$f$ diff, convex so

$$p = \frac{df}{dx} \quad \text{at sup.}$$

$\frac{df}{dx}$ is a function, and each is $1-1$

So $g$ st $g \circ \frac{df}{dx} = Id$ exists

$$g(p) = x. \quad \text{at sup}$$

So
$$f^*(p) = p \cdot g(p) - f(g(p))$$

$$\frac{df^*}{dp} = p \frac{dg}{dp} + g - f' \cdot \frac{dg}{dp} = g(p)$$

So
$$\frac{d^2 f^*}{dp^2} = \frac{dg}{dp} = \frac{dx}{dp} \quad ; \quad \text{but} \frac{dp}{dx} = \frac{d^2 f}{dx^2}$$

$$d^2 f^* = \frac{1}{d^2 f/.} > 0$$

all this works in ND as well

(ex: prove it!)

Thm (Fenchel - Moreau).

$$f = (f^*)^*$$

iff

    $f$ is proper, lower semi-continuous.

          and convex

OR

    $f \equiv \infty$

OR

    $f \equiv -\infty$

Now consider $A(\theta) = \log Z(\theta)$.

for an exp. fist.

1) $A^*(\mu) = \sup_{\theta} \left( \langle \theta, \mu \rangle - A(\theta) \right)$

is defined, convex.

2) for $\mu \in \mathcal{M}$

$\quad \underset{\text{marginal polytope.}}{}$

write $\theta(\mu) = \Lambda^{-1}(\mu)$

then

$A^*(\mu) = -H(p(x; \theta(\mu)))$.

Proof of 2 (sketch):

$$\Lambda^{-1}(\mu) = \theta \quad \text{such that}$$

$$E_{p(x;\theta)}\left[\varphi(x)\right] = \mu$$

$$\left(\text{by defn of } \Lambda\right)$$

But if $\mu = E_{p(x;\theta)}[\phi] = \nabla_\theta A(\theta)$

then $\theta$ is sup

so

$$-H(p(x;\theta(\mu))) = E_{p(x;\theta(\mu))}\left[\langle\theta, \phi(x)\rangle - A(\theta)\right]$$

$$= \langle\theta, \mu\rangle - A(\theta)$$

$$= A^*(\mu) \quad \left(\text{cause } \theta \text{ is sup}\right)$$

$$A(\theta) = \sup_{\mu \in M} \left\{ \langle \theta, \mu \rangle - A^*(\mu) \right\}$$

(A is lower semicontinuous — see notes;

then Fenchel-Moreau means

$$\left( A^* \right)^* = A \, .$$

and $\left( A^* \right)^*(\theta) = \sup_{\mu \in M} \left\{ \langle \theta, \mu \rangle - A^*(\mu) \right\}$

Compare:          with

$$E_Q = -E_Q \log p + E_Q \log q$$

Which        we        Minimized      to build var model

marginals of Q

$$E_Q \log p \longrightarrow \langle \theta, \mu \rangle$$

↑

params of

Now we have.

$$\sup_{\mu \in \mathcal{M}} \left\{ \langle \theta, \mu \rangle - A^*(\mu) \right\}. \qquad \longleftarrow \; \text{\textcircled{C}}$$

is attained at

$$\mu = E_{p(x; \theta)}[\phi]$$

So solving \text{\textcircled{C}} gives

- log partition function
- set of mean pars (for our purposes,) arg Max

BUT

$\mathcal{M}$ is hard.

$A^*(\mu)$ is hard.

Now we can unify algs.

## Mean field, single tree, etc

for any $\mu \in M$,

$$A(\theta) \geqslant \langle \mu, \theta \rangle - A^*(\mu).$$

<u>Now</u>  consider  $T \subset M$

$\lefthookdownarrow$ corresponding to
models that are tractable
$\equiv$ can compute $A^*(\mu)$

then  solve

$$\sup_{\mu \in T} \left\{ \langle \theta, \mu \rangle - A^*(\mu) \right\} = A_{MF}(\theta)$$

$\lefthookdownarrow$ same as our exp
but w - sign of max

must have    $A_{MF}(\theta) \leqslant A(\theta)$

# loopy BP

1) recall for a tree structured model

$$H(q) = -\sum_{i \in \text{verts}} \left[ \sum_{x \in \text{values}} q_i(x_i) \log q_i(x_i) \right]$$

$$- \sum_{i,j \in \text{edges}} \left[ \sum_{\substack{x_i, x_j \\ \in \text{values}}} q_{ij}(x_i, x_j) \cdot \log \left[ \frac{q_{ij}(x_i, x_j)}{q_i(x_i) q_j(x_j)} \right] \right]$$

2) approximate

$$A^*(\mu) \approx -H(\mu) \quad \longleftarrow \text{computed using tree expression}$$

$$\underline{\text{Bethe} \quad \text{approx}}$$

3) $$\mathcal{L} \supset M$$

$$\uparrow$$

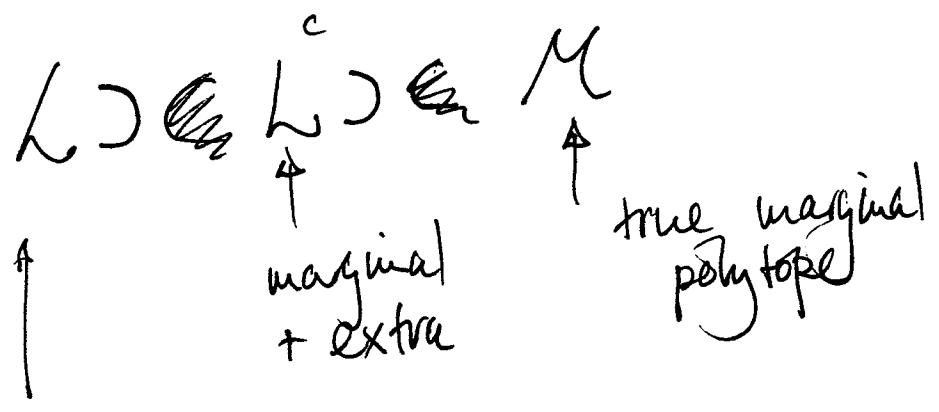local polytope, consistency constraints for pairwise marginals

4) Solve

$$\sup_{\mu \in \mathcal{L}} \left\{ \langle \theta, \mu \rangle + H_B(\mu) \right\} = A_{LBP}(\theta)$$

we must have

$$A_{LBP}(\theta) \geqslant A(\theta)$$

<u>But</u> we can now explore other

approximations :

    - eg. insert constraints so that

$$L \supset \mathcal{L} \supset \mathcal{L}^c \supset \mathcal{L} \supseteq M$$

defined
by marginal
constraints

marginal
+ extra

true marginal
polytope

heres one construction.

Assume some vector $\mu$, which might be in $M$

Construct the matrix

$$M = \begin{array}{|cc|c|cc|}
\hline
1 & \begin{array}{c} \mu_1 2 \\ \text{unary} \\ \text{Marg for first var} \end{array} & & \mu_2 & \cdots \\
\mu_1 & & & & \\
& & & \mu_{12} & \\
\hline
& & & & \\
\hline
\mu_2 & \mu_{21} & & & \\
& & & & \\
\hline
\vdots & & & & \\
\end{array}$$

$M$ is a <u>covariance</u> matrix

so $M \succeq 0$