

Q7 review - Parameter est. uGMs

- (a) check calculations / derivations
- (*) use Jordan 9.3-9.5 + 20 to assist \rightarrow SICM
- (x) skim Cafferty or tutorial

sufficient statisticsfrom Jordan (2003)
ch 9.3

(A1)

- total count and clique counts
- \underline{x}_v - random vector associated with graph
- $\underline{x}_c \subseteq \underline{x}_v$ - random vectors — " subsets of nodes
- IID replicates of obs.
- $\underline{x}_{v,n}$ - n^{th} replicate of subset c
- $D = \{\underline{x}_{v,1}, \underline{x}_{v,2}, \dots, \underline{x}_{v,N}\}$ - completely observed
- Parameterise uGM via clique potentials $\psi_c(\underline{x}_c)$ for $c \in C$
- C - set of cliques
- via Hammersley-Clifford:

Joint prob: $p(\underline{x}_v | \Theta) = \frac{1}{Z} \prod_c \psi_c(\underline{x}_c)$ $\Theta = \{\psi_c(\underline{x}_c), c \in C\}$ (param.)

Z = normalisation

$$Z = \sum_{\underline{x}_v} \prod_c \psi_c(\underline{x}_c)$$

(A) Z - obtained via summing (or integrating) over all configurations
 \underline{x}_v

The total counts

~~def~~
 - no. of times that configuration \underline{x}_v is observed
 in a dataset D

Marginal counts ^(clique counts)
 for clique C

$$m(\underline{x}_v) := \sum_{n=1}^N \delta(\underline{x}_v, \underline{x}_{v,n})$$

$$m(\underline{x}_c) := \sum_{\underline{x}_v \in C} m(\underline{x}_v)$$

$$\cdot N = \sum_{\exists v} m(\exists v) \quad - \text{total no. of observations}$$

(*) log-likelihood for NLMs (tabular clique pot.) (A3)

(**) express log-likelihood in terms of counts (sufficient statistics for discrete models)

- introduce dummy var $\exists v$ (?)

$$p(\exists v, n | \theta) = \prod_{\exists v} p(\exists v | \theta)^{\delta(\exists v, \exists v, n)}$$

- OK- indicator
tickle that switches
on/off
depending on
distri.

$$\textcircled{R} \quad \delta(\exists v, \exists v, n) = \prod \{ \exists v = \exists v, n \}$$

\textcircled{R} dummy variable $\exists v$ ranges across configurations of nodes rather than across data points.

- standard for multinomials (remember Bishop?)

$$\text{Probability of observed data: } - p(D | \theta) = \prod_{n=1}^N p(\exists v, n | \theta)$$

$$= \prod_{n=1}^N \prod_{\exists v} p(\exists v | \theta)^{\delta(\exists v, \exists v, n)}$$

log-likelihood in terms of marginal counts

$$l(\theta; D) = \log p(D | \theta)$$

$$= \sum_{n=1}^N \sum_{\exists v} \delta(\exists v, \exists v, n) \log p(\exists v | \theta)$$

$$= \sum_{\exists v} \sum_n \delta(\exists v, \exists v, n) \log p(\exists v | \theta)$$

factor out terms without sum. avg.

$$= \sum_{\exists v} m(\exists v) \log p(\exists v | \theta)$$

subs. $p(\exists v | \theta)$

$$= \sum_{\exists v} m(\exists v) \log \left(\frac{1}{Z} \prod_C \psi_C(\exists v) \right)$$

$$= \sum_{\exists v} m(\exists v) \sum_C \log \psi_C(\exists v) - \sum_{\exists v} m(\exists v) \log Z$$

$$\text{?} \quad l = \sum_c \sum_{x_c} m(x_c) \log \psi_c(x_c) - N \log Z \quad (9.43)$$

- name
use
@ 0/s 1
- see earlier

- $m(x_c)$ - neg. counts \rightarrow suff. statistics
- $N \log Z$ - appears (not in DGM)

9.2
- Assumptions
and investigate

MLE of l :

(*) $N \log Z$ - coupled, nonlinear set of eq. with implicit poem app.

$$l(\theta, D) = \sum_c \sum_{x_c} m(x_c) \log \psi_c(x_c) - N \log Z \quad (9.43)$$

(*) Partial deriv wrt $\psi_c(x_c)$; with clique c , config x_c held fixed.

$$\begin{aligned}
 \text{(i)} \quad \frac{\partial}{\partial \psi_c(x_c)} m(x_c) \log \psi_c(x_c) &= \frac{m(x_c)}{\psi_c(x_c)} \quad \text{Note that } \theta \text{ and } \hat{x} \text{ are just diff dummy indexing variables} \\
 \text{(ii)} \quad \frac{\partial}{\partial \psi_c(x_c)} \log Z &= \frac{1}{Z} \frac{\partial}{\partial \psi_c(x_c)} \left\{ \sum_{\tilde{x}_D} \prod_D \psi_D(\tilde{x}_D) \right\} \\
 &= \frac{1}{Z} \sum_{\tilde{x}} \delta(\tilde{x}_c, x_c) \frac{\partial}{\partial \psi_c(x_c)} \left(\prod_D \psi_D(\tilde{x}_D) \right) \\
 &= \frac{1}{Z} \sum_{\tilde{x}} \delta(\tilde{x}_c, x_c) \prod_{D \neq c} \psi_D(\tilde{x}_D) \\
 &= \sum_{\tilde{x}} \delta(\tilde{x}_c, x_c) \frac{1}{\psi_c(\tilde{x}_c)} \frac{1}{Z} \prod_D \psi_D(\tilde{x}_D) \\
 &= \frac{1}{\psi_c(x_c)} \sum_{\tilde{x}} \delta(\tilde{x}_c, x_c) p(\tilde{x}) \\
 &= \frac{p_c(x_c)}{\psi_c(x_c)}
 \end{aligned}$$

• note all $\psi_D(\tilde{x}_D)$ where $D \neq c$ are assumed constant wrt $\psi_c(x_c)$

(ds 2)

- not entirely clear crystal

yielding $\frac{\partial l}{\partial \psi_c(x_c)} = \frac{m(x_c)}{\psi_c(x_c)} - N \frac{p(x_c)}{\psi_c(x_c)}$ $\log \psi_c(x_c) > 0$

- set $\frac{\partial l}{\partial \psi_c(x_c)} = 0 \Rightarrow m(x_c) = N p(x_c)$

⑦ define empirical distn $\hat{p}(x) = \frac{m(x)}{N}$, $\hat{p}(x_c) = \frac{m(x_c)}{N}$ is a marg. inde emp. (?)

$\Rightarrow \frac{m(x_c)}{N} = p(x_c) \Rightarrow \hat{p}_{m_c}(x_c) = p(x_c)$

(*) for each clique $c \in C$, the model marginals $p(x_c)$ must be equal to empirical marginals $\hat{p}_{m_c}(x_c)$

Jordan 2003 9.3.3. \rightarrow MLE by inspection for decomposable graphs

- (1) iterative prop. fitting (IPF)

- IPF is not only a fixed point algo, but coordinate ascent algo

(*) use:- $\frac{\hat{p}(x_c)}{\psi_c(x_c)} = \frac{p(x_c)}{\psi_c(x_c)}$

- Hold $\psi_c(x_c)$ fixed in numerator, denom of RHS.
- why? As $\psi_c(x_c)$ appears implicitly through $p(x_c)$
- introduce like marking; - paramest at iter + $\psi_c^{(t)}(x_c)$
- joint prob based on estimates $p^{(t)}(x)$ at iter t

$$\Rightarrow \frac{\psi_c(x_c)}{\hat{p}(x_c)} = \frac{\psi_c(x_c)}{p(x_c)}$$

$$\Rightarrow \psi_c(x_c) = \psi_c(x_c) \frac{\hat{p}(x_c)}{p(x_c)}$$

$$\Rightarrow \psi_c^{(t+1)}(x_c) = \psi_c^{(t)}(x_c) \frac{\hat{p}(x_c)}{p^{(t)}(x_c)} \quad (9.61) \quad (\text{IPF algorithm})$$

(*) cycle through all cliques $c \in C$, applying (9.61); one cycle - one iter.

(1) Properties of IPF

- summarised from Jordan (2003) ch9.3.

- 1) marginal $p^{(t+1)}(x_c)$ is equal to empirical marginal $\hat{p}(x_c)$
- 2) normalisation constant Z - constant across updates

which yields: $p^{(t+1)}(z_v) = p^{(t)}(z_v | z_c) \hat{p}(z_c)$

- each IPF iteration retains previous cond. prob $p(z_v | z_c)$ and replaces prev. mag. prob $p(t)(z_c)$ with new marginal $\hat{p}(z_c)$.

(*) IPF as w-coordinate ascent

- additional derivations in Jordan show this sets gradient of log likelihood to 0; and is coordinate ascent (in log-likelihood)

(*) KL divergence view of IPF (A5)

- A few key equations:

KL-divergence with marginals:

- decompose $p(z_A, z_B) = p(z_B | z_A) p(z_A)$ and $q(z_A, z_B) = q(z_B | z_A) q(z_A)$

KL-div: $D(p(z_A, z_B) || q(z_A, z_B)) = D(p(z_A) || q(z_A)) + \sum_{z_A} p(z_A) D(p(z_B | z_A) || q(z_B | z_A))$

(*) (B): maximizing likelihood equivalent to minimising following KL-div:-

$$D(\hat{p}(z) || p(z | \theta)) = \sum_z \hat{p}(z) \log \frac{\hat{p}(z)}{p(z | \theta)} \quad \hat{p}(z) - \text{empir. dist.}$$

(*) equivalence $\rightarrow -\text{ve log likelihood or KL div diff by } \sum_z \hat{p}(z) \log \hat{p}(z)$
which is independent of θ .

(*) coordinate descent on $D(\hat{p}(z) || p(z | \theta))$ (coordinates - precutes of clique potentials)

- fix a clique C

- adjust clique potential $\psi_C(z_C)$ to min KL-div.

(B)

$$\Rightarrow D(\hat{p}(z) || p(z | \theta)) = D(\hat{p}(z_C) || p(z_C | \theta)) + \sum_{z_C} \hat{p}(z_C) D(\hat{p}(z_{\bar{V}} | z_C) || p(z_{\bar{V}} | z_C, \theta)) \quad (\text{II})$$

(*) A little course lecture also helps with exposition.

- changes in clique pot $\psi_C(z_C)$ does not affect $p(z_{\bar{V}} | z_C, \theta)$

\Rightarrow (I) unaffected by changes to $\psi_C(z_C)$ and minimising KL div

$D(\hat{p}(z) || p(z | \theta))$ amounts to minimising (I)

(*) Minimising (1) achieved by setting $p(x_c|\theta) = \hat{p}(x_c)$ i.e. match to empirical marginal

(*) This what IPF achieves \rightarrow coordinate ascent in log-like coordinate descent in KL-div

(*) Note IPF takes form of scaling algorithm in which potentials are multiplied by a ratio of marginals.

$$\textcircled{A} \max \ell \Leftrightarrow \min \text{KL}(\hat{p}(x) \parallel p(x|\theta)) = \sum_x \hat{p}(x) \log \frac{\hat{p}(x)}{p(x|\theta)} \quad \textcircled{A}$$

(*) Verifying identity (A) :- (KLdiv. decompr.)

$$D(p(x_A, x_B) \parallel q(x_A, x_B)) = \sum_{x_A, x_B} p(x_A) p(x_B|x_A) \log \frac{p(x_A) p(x_B|x_A)}{q(x_A) q(x_B|x_A)}$$

$$= \sum_{x_A, x_B} p(x_A) p(x_B|x_A) \log \frac{p(x_A)}{q(x_A)} + \sum_{x_A, x_B} p(x_A) p(x_B|x_A) \log \frac{p(x_B|x_A)}{q(x_B|x_A)}$$

$$= \sum_{x_A, x_B} p(x_A, x_B) \log \frac{p(x_A)}{q(x_A)} + \sum_{x_A} p(x_A) \sum_{x_B} p(x_B|x_A) \log \frac{p(x_B|x_A)}{q(x_B|x_A)}$$

$$= \sum_{x_A} \log \frac{p(x_A)}{q(x_A)} \sum_{x_B} p(x_A, x_B) + \sum_{x_A} p(x_A) \sum_{x_B} p(x_B|x_A) \log \frac{p(x_B|x_A)}{q(x_B|x_A)}$$

$$= \sum_{x_A} p(x_A) \log \frac{p(x_A)}{q(x_A)} + \sum_{x_A} p(x_A) D(p(x_B|x_A) \parallel q(x_B|x_A))$$

$$= D(p(x_A) \parallel q(x_A)) + \sum_{x_A} p(x_A) D(p(x_B|x_A) \parallel q(x_B|x_A))$$

(*) Features and micro-potentials

(*) lecture notes (S2019) are a little clearer on feature design.

(*) Rather than define a table of values to encompass the mapping of $\psi_c(x_c) = \psi_c(x_1, x_2, x_3)$ over e.g. 2^6^3 possibilities

- (6) we instead use domain knowledge (e.g. a vocabulary/dictionary) to reduce representational granularity
- (7) Achieved by assigning a score to common three letter stems; or scores to 'significant' occurrences e.g. fing = 10, fted = 9 etc. \otimes
- (8) remainder → assign an arbitrarily low score to reflect low likelihood of occurrence.
- (9) sense in which 'feature' is a function which is 'vacuous' over all joint settings except a few particular ones.
- mathematical details • define K features $f_k(c_1, c_2, c_3)$ e.g. fingfing
- Assign each of K features a weight θ_k
 - A micropotential is then achieved by exponentiating $e^{\theta_k f_k(c_1, c_2, c_3)}$
- \oplus A clique potential is formed by multiplying together micropotentials

② or use indicators?

- to avoid simult. estim./ setting θ_R AND f_k

• yielding: $\psi_C(x_C) = \psi_C(x_1, x_2, x_3) = e^{\theta_1 f_1} \cdot e^{\theta_2 f_2} \cdots \cdot e^{\theta_K f_K}$

$$= \exp \left\{ \sum_{k=1}^K \theta_k f_k \right\}$$

- Potential is still over 26^3 settings, but only have K parameters if K features are used

(10) we then can estimate the weights θ_k to deal with context. info

(11) θ_k - 'strength' of feature and whether it increases or decreases clique probability

(*) combining features:

- Jordan (2003) → suggests f_k is chosen to be a indicator

- Marginal probability over a clique: $p(c_1, c_2, c_3) \propto \exp \left\{ \sum_{k=1}^K \theta_k f_k \right\}$

- Addit. complexities → overlapping features, microstructure; any fraction of subset of clique variables, overlapping words do not alter anatomy

④ unique potential as weighted sum of expon. features.

$$\psi_c(\underline{x}_c) = \exp \left\{ \sum_{i \in I_c} \theta_i f_k(x_{ci}) \right\}$$

(*) features based model

- usually, $p(\underline{x}_v | \theta) = \frac{1}{z(\theta)} \prod_c \psi_c(\underline{x}_c) = \frac{1}{z(\theta)} \exp \left\{ \sum_c \sum_{i \in I_c} \theta_i f_k(x_{ci}) \right\}$

- But drop explicit assoc.
of features and cigos (?) aduse

$$p(\underline{x} | \theta) = \frac{1}{z(\theta)} \exp \left\{ \sum_i \theta_i f_i(\underline{x}) \right\}$$

where $f_i(\cdot)$ - features
 θ_i - param
 $z(\theta)$ - norm. factor

⑤ exponential family model with features as sufficient statistics

$$z(\theta) = \sum_x \exp \left\{ \sum_i \theta_i f_i(x) \right\}$$

⑥ estimate θ_i from data D.
(param)

MLE of feature based NGMS: (see Jordan 2003) (problem setup)

- ~~as~~ ②
scaled neg.-likelihood: $\tilde{i}(\theta; D) = \frac{l(\theta; D)}{N} = \frac{1}{N} \sum_{n=1}^N \log p(x_n | \theta)$ @?

$$⑦ \tilde{i}(\theta; D) = \sum_{\underline{x}} \hat{p}(\underline{x}) \log p(\underline{x} | \theta)$$

$$= \sum_{\underline{x}} \hat{p}(\underline{x}) \log \left[\frac{1}{z(\theta)} \left\{ \exp \sum_i \theta_i f_i(\underline{x}) \right\} \right]$$

$$= \sum_{\underline{x}} \hat{p}(\underline{x}) \left\{ \sum_i \theta_i f_i(\underline{x}) \right\} - \log z(\theta)$$

- Jordan (2003) states we use convexity of $\log(\cdot)$ to bound $\log z(\theta)$ ②

- We then get:-

$$\log z(\theta) \leq \mu z(\theta) - \log \mu - 1$$

- 0153: $\log(\cdot)$ is not convex; it's concave
0154: what conseq for presentation?

(*) For now, assume this is okay \rightarrow add to overspill and investigate at the end.

(*) Bound holds for all μ and $\mu = z^{-1}(\theta^{(t)})$

$$\Rightarrow \hat{z}(\theta|D) \geq \sum_x \tilde{p}(x) \sum_i \theta_i f_i(x) - \frac{z(\theta)}{z(\theta^{(t)})} - \log z(\theta^{(t)}) + 1$$

with equality at $z(\theta^{(t)})$ ⑤

(*) Further manipulation of scaled log-like.

- Define $\Delta\theta_i^{(t)} := \theta_i - \theta_i^{(t)}$; then (s.t.s. $z(\theta)$)

$$\hat{z}(\theta|D) \geq \sum_x \tilde{p}(x) \sum_i \theta_i f_i(x) - \frac{1}{z(\theta^{(t)})} \sum_x \exp \left\{ \sum_i \theta_i f_i(x) \right\} - \log z(\theta^{(t)}) + 1$$

$$= \sum_x \tilde{p}(x) \sum_i \theta_i f_i(x) - \frac{1}{z(\theta^{(t)})} \sum_x \exp \left\{ \sum_i (\Delta\theta_i^{(t)} + \theta_i^{(t)}) f_i(x) \right\} - \log z(\theta^{(t)}) + 1$$

$$= \sum_x \tilde{p}(x) \sum_i \theta_i f_i(x) - \frac{1}{z(\theta^{(t)})} \sum_x \exp \left\{ \sum_i \theta_i^{(t)} f_i(x) \right\} \exp \left\{ \sum_i \Delta\theta_i^{(t)} f_i(x) \right\} - \log z(\theta^{(t)}) + 1$$

$$= \sum_x \tilde{p}(x) \sum_i \theta_i f_i(x) - \underbrace{\sum_x \frac{1}{z(\theta^{(t)})} \exp \left\{ \sum_i \theta_i^{(t)} f_i(x) \right\} \exp \left\{ \sum_i \Delta\theta_i^{(t)} f_i(x) \right\}}_{= p(x|\theta^{(t)})} - \log z(\theta^{(t)}) + 1$$

$$= \sum_i \theta_i \sum_x \tilde{p}(x) f_i(x) - \sum_x p(x|\theta^{(t)}) \exp \left\{ \sum_i \Delta\theta_i^{(t)} f_i(x) \right\} - \log z(\theta^{(t)}) + 1$$

copied $f_i(x)$ and $\Delta\theta_i^{(t)}$

assume $f_i(x) \geq 0$ and $\sum_i f_i(x) = 1$

$f_i(\cdot)$ is convex, invoke Jensen's inequal.

$$\exp \left(\sum_i \pi_i x_i \right) \leq \sum_i \pi_i \exp(x_i) \quad \text{for } \sum_i \pi_i = 1$$

0155a

(*) Note; f_i play the role of π_i as they are positive and sum to 1.

(*) cross-ref with lecture 7 notes

↳ f_i are being treated as 'weights' and $\Delta\theta_i^{(t)}$ as arguments
yielding the following lower bound on scaled neg-likelihood:- (with poems decoupled)

$$\hat{L}(\theta; D) \geq \sum_i \hat{p}(x) f_i(x) - \sum_x p(x|\theta^{(t)}) \sum_i f_i(x) \exp(\Delta\theta_i^{(t)}) - \log Z(\theta^{(t)}) + 1$$

$$= L(\theta)$$

• we then take derivatives with respect to lower bound $L(\theta)$:-

$$\frac{\partial L}{\partial \theta_i} = \sum_x \hat{p}(x) f_i(x) - \exp(\Delta\theta_i^{(t)}) \sum_x p(x|\theta^{(t)}) f_i(x) = 0$$

0/56

$$\Rightarrow \exp(\Delta\theta_i^{(t)}) = \frac{\sum_x \hat{p}(x) f_i(x)}{\sum_x p(x|\theta^{(t)}) f_i(x)}$$

0/57

• defining (or because?) $p^t(x)$ is an unnormalised version of $p(x|\theta^{(t)})$

$$\Rightarrow \exp(\Delta\theta_i^{(t)}) = \frac{\sum_x \hat{p}(x) f_i(x)}{\sum_x p^t(x) f_i(x)} = z(\theta^{(t)})$$

(0): next bit in Jordan or Xing not clear from notes; attempted to clarify
(fecture notes + slides)

$$(*) \theta_i^{(t+1)} = \theta_i^{(t)} + \Delta\theta_i^{(t)} \Rightarrow p^{(t+1)}(x) = p^{(t)}(x) \left(\prod_i e^{\Delta\theta_i^{(t)} f_i(x)} \right) \quad ? \quad 0/57$$

Jordan (2003):

(*) To upgrade parameters from $\theta^{(t)}$ to $\theta^{(t+1)}$, multiply $p(x|\theta^{(t)})$ by $e^{\Delta\theta_i^{(t)} f_i(x)} \forall i$

$$(*) \text{Note: } \frac{p^{(t)}(x)}{z(\theta^{(t)})} = p(x|\theta^{(t)}) \text{ and } \rightarrow$$

(*) we have:-

$$\begin{aligned}
 p^{(t+1)}(\underline{x}) &= p(\underline{x} | \underline{\theta}^{(t)}) \left(\prod_i e^{\Delta\theta_i^{(t)} f_i(\underline{x})} \right) \\
 &= \frac{p^{(t)}(\underline{x})}{Z(\underline{\theta}^{(t)})} \prod_i \left(\frac{\sum_{\underline{x}} \hat{p}(\underline{x}) f_i(\underline{x})}{\sum_{\underline{x}} p^{(t)}(\underline{x}) f_i(\underline{x})} \right)^{f_i(\underline{x})} \\
 &= \frac{p^{(t)}(\underline{x})}{Z(\underline{\theta}^{(t)})} \prod_i \left(\frac{\sum_{\underline{x}} \hat{p}(\underline{x}) f_i(\underline{x})}{\sum_{\underline{x}} p^{(t)}(\underline{x}) f_i(\underline{x})} \right)^{f_i(\underline{x})} Z(\underline{\theta}^{(t)}) \quad \text{As } \sum_i f_i(\underline{x}) = 1 \\
 &= p^{(t)}(\underline{x}) \prod_i \left(\frac{\sum_{\underline{x}} \hat{p}(\underline{x}) f_i(\underline{x})}{\sum_{\underline{x}} p^{(t)}(\underline{x}) f_i(\underline{x})} \right)^{f_i(\underline{x})} \quad (\text{GIS Algorithm})
 \end{aligned}$$

- Subs for
 $p(\underline{x} | \underline{\theta}^{(t)})$ and $e^{\Delta\theta_i^{(t)}}$

(*) Exponential family, MLE, sufficient statistics

• Import. for understanding
 - drop v subscript method of exponential magnets
 - substitute feat. rep of joint prob.
 (condit. on param)

$$l(\underline{\theta}; D) = \sum_{\underline{x}} m(\underline{x}) \log p(\underline{x} | \underline{\theta})$$

$$= \sum_{\underline{x}} m(\underline{x}) \log \left[\frac{1}{Z(\underline{\theta})} \exp \left\{ \sum_i \theta_i f_i(\underline{x}) \right\} \right]$$

$$= \sum_{\underline{x}} m(\underline{x}) \left\{ \sum_i \theta_i f_i(\underline{x}) - \log Z(\underline{\theta}) \right\}$$

$$= \sum_{\underline{x}} m(\underline{x}) \sum_i \theta_i f_i(\underline{x}) - \log Z(\underline{\theta}) \sum_{\underline{x}} m(\underline{x})$$

$$= \sum_{\underline{x}} m(\underline{x}) \sum_i \theta_i f_i(\underline{x}) - N \log Z(\underline{\theta})$$

$$\sum_{\underline{x}} m(\underline{x}) = N$$

(*) Taking derivatives wrt θ_i ; setting to 0, i.e. going for an MLE est.:-

$$\frac{\partial}{\partial \theta_i} l(\underline{\theta}; D) = \sum_{\underline{x}} m(\underline{x}) f_i(\underline{x}) - N \frac{\partial}{\partial \theta_i} \log Z(\underline{\theta}) \quad \textcircled{2}$$

$$= \sum_{\underline{x}} m(\underline{x}) f_i(\underline{x}) - N \sum_{\underline{x}} p(\underline{x} | \underline{\theta}) f_i(\underline{x}) = 0$$

$$\Rightarrow \sum_{\underline{x}} p(\underline{x} | \underline{\theta}) f_i(\underline{x}) = \sum_{\underline{x}} \frac{m(\underline{x})}{N} f_i(\underline{x}) = \sum_{\underline{x}} \hat{p}(\underline{x} | \underline{\theta}) f_i(\underline{x})$$

0/59/2

(*) At ML estimate:-

- Expectations of the sufficient statistics under the model must match the empirical feature average
- Presumably former is $\sum_x p(x|\theta) f_i(x)$; and latter is $\sum_x \hat{p}(x|\theta) f_i(x)$.

(*) (2)- Information theoretic and statistical physics in ML

- Have to review (15) on exponentials and GLMS

(2)- MLEs with 36-705:

- Sufficiency, exponentials

Fittman-Koopman-Denjoy theorem

(*) "Among families of probability distributions whose domain does not vary with the parameter being estimated, only in exponential families is there a sufficient statistic whose dimension remains bounded as sample size ↑"

(*) "Sufficiency sharply restricts possible forms of distri"

(i): Leads to some interesting discussions by Peri Diaconis

(*) Info-theory/statistical physics principles in ML

(i) Ex makes point for modelling:- (and 10 why exponential family comes up)

(ii) Review exponential family as a solution to a constrained, variational optims. problem in which objective is entropy or KL divergence; and constraints are that expectations under the distri are matched to expectations under the empirical distri.

i) variational in the sense of min/max or functional by choosing fn (distri)

ii) ~~obj~~

- ex: Begin with fixed feature exp. $\sum_x p(x) f_i(x) = x_i$

(*) Assuming consistent exp; choose a distri:-

$$\max_P H(p(x)) = - \sum_x p(x) \log p(x)$$

$$\text{s.t. } \sum_x p(x) f_i(x) = x_i \Rightarrow p^*(x|\theta) = \frac{1}{Z(\theta)} \exp \left\{ \sum_i \theta_i f_i(x) \right\}$$

$$\sum_x p(x) = 1$$

(*) more generally;

$$\min_p \text{KL}(p(x) || h(x)) = \min_p \sum_x p(x) \log \frac{p(x)}{h(x)} = -H(p) - \sum_x p(x) \log h(x)$$

$$\text{s.t. } \sum_x p(x) f_i(x) = \alpha_i$$

$$\sum_x p(x) = 1$$

$$\Rightarrow p^*(x|\theta) = \frac{1}{Z(\theta)} h(x) \exp\left\{\sum_i \theta_i f_i(x)\right\}$$

(*) interpretation

(*) maximum entropy principle:-

- From amongst distributions consistent with the data, select the distri^{*} whose shannon entropy is maximal

~ amounts to choosing distri with maximal uncertainty, as defined by the entropy functional

(*) for KL divergence

- incorporates prior $h(x)$
- choose distri that contains 'least added ass' above priors.

(*) Additional details on

info-theory - stat mech \rightarrow Jaynes (1967)

~ Jordan (2008) - Foundations & Trends ch 3.

