EX: 4th year PHD did not know how to run k-means
- you want to know theorem; but you have to get it working (OT)
- (W): implement these(?)

- Potential functions $\psi_c(x_c)$
- Some issues: - spell-checking example
  - Build affinity models of charact. streams
  - e.g. consec appearances of 3 streams
  - use this criterion to score likelihood
  
  $$\psi_c(x_c) = \psi_c(x_1, x_2, x_3)$$
  
  'zyz' and 'ink'
  - Define potential function over a triplet of characters
  - $26^3$ different features and put a function on it?
  
  $-26^5$    ''

(*) Above shows infeasibility of tabular potentials (akin to enumerat.
                                                                of joint)
Feature based clique pot.              (W): note instantiation
                                            of cliques in practice

(*) Features

(Q): Find a way of appropriately compressing the granularities
- Handcrafted feature design ——> role of human knowledge in A.I
- use feature-engin to save on rep. cost?
- Distinct from current, overcomplete ML models (e.g. gigantic placeholder)
- EX: Against the idea of human knowledge being
              ignored in ML

(*) Micropotentials $f_R$          OR
- Have R features and weights defined over 3 charact. potential

- $\psi_c(c_1, c_2, c_3) = \exp\left\{\sum_{k=1}^{K} \theta_k f_k(c_1, c_2, c_3)\right\}$

(*) overall potential (clique) is exp weighted sum of micropot.

(*) micropotent. distinct from tabular potentials.

(*) $K$ parameters our $K$ features $\rightarrow$ more compact

---

## Combining features

(*) sliding window / overlapping sliding window

(*) note how we can modify standard ctd Gibbs. rep. for exp.

- Allows use of exponential / GLIMS

- ⑦ : Not entirely clear how to apply IPF in this case due to coupling of estimated $\theta_k$ and designed $f_k(c_1, c_2, c_3)$

---

## MLE of feature based UGMs:

- scaled likelihood:- $\hat{\ell}(\theta; D) = \ell(\theta; D)/N = \frac{1}{N}\sum_n \log p(x_n | \theta)$

$$= \sum_x \hat{p}(x) \log p(x|\theta)$$

$$= \underbrace{\sum_x \hat{p}(x) \sum_i \theta_i f_i(x)}_{(i)} - \underbrace{\log \tilde{Z}(\theta)}_{(ii)}$$

(*) calculus+derivatives $\rightarrow$ not fruitful.

Ex: Nonlinearities cause issues (e.g. log/inversion)
- linearise it so argument can be exposed to linear attack

- log has linear upper bound

$\quad$ - $\log \tilde{Z}(\theta) \leq \mu \tilde{Z}(\theta) - \log \mu - 1$

$\quad$ - Bound holds $\forall \mu$ : $\mu = \tilde{Z}^{-1}(\theta^{(t)})$ — fixed point it. strategy
$\quad\quad\quad$ — assume this (this is a previous version of $\tilde{Z}$)

(*) GIS derivation   @ ⑰ : Review

- Define $\Delta\theta_i^{(t)} = \theta_i - \theta_i^{(t)}$ and introduce

- still nasty: → every

(*) Note exp of weighted sum :- $\exp\left\{\sum_i \Delta\theta_i^{(t)} f_i(x)\right\}$

EX: We make distinction between weight $\theta_k$ and $(\Delta\theta_i^{(t)})$ and features $f_i(x)$; ⎫
    algebraically the same.                                                                        ⎪
- Treat $f_i$ as weights; $\Delta\theta_i^{(t)}$ as arguments                                       ⎬ relaxin
- (*) impose prob. constraints (normally applying to weights) to                                   ⎪
    $f_i$ our assumed "weights".                                                                    ⎭

  (*) $\exp(\cdot)$ is convex  ——→ use Jensen's

- Algebraic trick often used in ML

  - getting $\sum_i f_i(x)\exp(\Delta\theta_i^{(t)})$ → only linearly coupled with others.

_____

(*) use low bound of scaled LL :-      GIS

(*) use calculus :-

(*) giving update steps :-

  (*) Note :-   $\dfrac{\sum_x \hat{p}(x) f_i(x)}{\sum_x p^{(t)}(x) f_i(x)}$      " weighted sum of feature:
                                                                                      by emp. prob. $\tilde{p}(x)$ "

                                                                                  " features weighted
                                                                                      by inferred probability $p^t(x)$ "

(*) Iterative re-scaling          ↳ connection with IPF.

_____

(*) Summary        GIS/IPF
- fixed point iterations on LL obj.
- one for tabular, feature based

(*) Where does exponential come from?  (☺: A move from Gauss → Gibbs?)

(☺)(A?): Review exp. family form (LS)

(*) Note of MLE; expectations of sufficient statistics under model match feature average'
                                                    exp.
                                            ^

- Note eq.

(*) Begin with exp. family → get a consequence.

- Reverse rationale          - Impose constraint          arbitrary ones
                             on distr: so you can't give me n - expectat. of
- Maximum entropy                                              feature must
                             - Fixed feat. exp:  $\sum_x p(x) f_i(x) = x_i$   match
                                                                              complexv.
                               (from data)      - encode
                                                 few assumpt. about how a
                                              - entropy as amt             as poss.
                                                 of randomness/ amt of
$$\max_p H(p(x)) = -\sum_x p(x) \log p(x)$$      assumptions made

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

$$\sum_x p(x) = 1 \quad \longrightarrow \quad p(x|\theta) = \frac{1}{Z(\theta)} \exp\left\{ \sum_i \theta_i f_i(x) \right\}$$

(*) variational definition:- define a distrib. as a solution to a constrained optimisation problem.

(☺)(A?): Review Lagrangian sol.

(*) natural consequence gives: an exponential family distri.

(*) Benefit of information theoretic principles on ML  → (☺) explore

- more general max entropy method

- incorporate prior distri on x; replace it. (h(x))
- Estimated distr. has least addit. assumptions from priors
- use KL-divergence rather than entropy

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

$$s.t. \; \sum_x p(x) f_i(x) = \alpha_i$$

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

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

---

(#) constraints from data

ⓠ Where do constraints $\alpha_i$ come from?

- data itself is the constraint

- (h)(g): Automatic consistency?

---

- Geometric step; general process :-

Either :-

1) Assume all exponential family distis as "model" :-

$$\mathcal{E} = \left\{ p(x) : p(x|\theta) = \frac{1}{Z(\theta)} h(x) \exp \left\{ \sum_i \theta_i f_i(x) \right\} \right\}$$

OR

2) Assume all distribs satisfying model constraints

$$M = \left\{ p(x) : \sum_x p(x) f_i(x) = \sum_x \hat{p}(x) f_i(x) \right\}$$      do not acknowledge roots

- Information geometry Pythagorean theorem :-

     → inspires V.I, deep gen. models.

---

(*) Summary

(*) exp family viewed as a sol. to variational exp → maximum entropy

(*)

supplementary → structure learning (see supp.) / 2020 (

case study: CRFs (Lafferty) - at CMU
- insight of experienced modelling

- Lafferty paper → impressive ; clear rationale, motivation

(*) 'local normalisability is a double-edged sword
  ( - makes computing 'simple'
  (?) an - M context of HMM

(*) What you want is global normalisability

- use scores rather than enforcing local normalisibility
  - use potentials

$$\exp\left\{\sum \theta_i f_i\right\}$$

  - features corresp. to nodes
  - use human knowledge for features

- 'Art of modelling' → totally comprehens.