L5 - Parameter estimation for PGMs - review

(*) Review slides → note key eq.
(*) use readings to supplement areas of uncertainty

(*) General parameter est. setup
- Assume structure of graph G is <u>fixed</u> /given in advance
- estimate parameters from IID dataset $D = \{x_1, x_2, \ldots, x_N\}$
  - N training instances
- Each training instance $x_n = \begin{pmatrix} x_{n,1} \\ \vdots \\ x_{n,M} \end{pmatrix}$  $x_n \in \mathbb{R}^M$  - each dim. of $x_n$ corresponds to a realisation of an r.v. i.e. a <u>node</u>

(*) Completely observable
  - $x_{n,i}$ is known $\forall\ n = 1, \ldots, N$ and $i \in 1, \ldots, M$

(*) Partially observable

  $\exists\, i : x_{n,i}$ is not observed.

- log-likelihood (function of param):-   (general BN)          $\pi_i$ - parents of node i

  $$l(\underline{\theta}; D) = \log p(D|\theta) = \log \left( \prod_{n=1}^{N} \left( \prod_{i=1}^{M} p(x_{n,i} | x_n, \pi_i, \theta_i) \right) \right)$$

  $$= \sum_{i=1}^{M} \sum_{i=1}^{M} \log p(x_{n,i} | x_n, \pi_i, \theta_i)$$

(*) Exponential family distr.                    (murphy 2012)
                vector
- for a numeric r.v. X :-  (scalar?)    (*)

  $$p(x|\eta) = h(\underline{x}) \exp \left\{ \eta^T \underline{T}(\underline{x}) - A(\eta) \right\} = \frac{1}{Z(\eta)} h(\underline{x}) \exp \left\{ \eta^T T(\underline{x}) \right\}$$

- is the exponential family distr with eq :-

(*) canonical param $\eta$
(*) sufficient statistic $\underline{T}(\underline{x})$
(*) log normaliser $A(\eta) = \log Z(\eta)$

(*) exp family + GLMS → includes many models instances
                         of this general form

Examples:-

i) MVG:

- $\underline{x} \in \mathbb{R}^k$

$$p(\underline{x} | \underline{\mu}, \underline{\Sigma}) = \frac{1}{(2\pi)^{k/2} |\Sigma|^{1/2}} \exp\left\{ -\frac{1}{2} (\underline{x} - \underline{\mu})^T \Sigma^{-1} (\underline{x} - \underline{\mu}) \right\} \quad \Big\}^{(i)}$$

$$= \frac{1}{(2\pi)^{k/2} |\Sigma|^{1/2}} \exp\left\{ -\frac{1}{2} tr\left( \Sigma^{-1} \underline{x} \underline{x}^T \right) + \underline{\mu}^T \Sigma^{-1} \underline{x} - \frac{1}{2} \underline{\mu}^T \Sigma^{-1} \underline{\mu} - \log |\Sigma| \right\}$$

(i) obtained by expanding; applying the trace trick to quadratic form

$\underline{x}^T \Sigma^{-1} \underline{x}$

- As $\underline{x}^T \Sigma^{-1} \underline{x} = tr(\underline{x}^T \Sigma^{-1} \underline{x}) = tr(\underline{x} \underline{x}^T \Sigma^{-1}) = tr(\Sigma^{-1} \underline{x} \underline{x}^T)$

Exponential family

$$\eta = \begin{bmatrix} \Sigma^{-1} \underline{\mu} \\ -\frac{1}{2} vec(\Sigma^{-1}) \end{bmatrix} = \begin{bmatrix} \eta_1 \\ vec(\eta_2) \end{bmatrix} \quad ; \quad \begin{array}{l} \eta_1 = \Sigma^{-1} \underline{\mu} \\ \eta_2 = -\frac{1}{2} \Sigma^{-1} \end{array} \qquad vec\left( \begin{bmatrix} a & b \\ c & d \end{bmatrix} \right) = \begin{bmatrix} a \\ c \\ b \\ d \end{bmatrix}$$

$$T(\underline{x}) = \begin{bmatrix} \underline{x} \\ vec(\underline{x}\underline{x}^T) \end{bmatrix}$$

$$A(\eta) = \frac{1}{2} \underline{\mu}^T \Sigma^{-1} \underline{\mu} + \log |\Sigma| = -\frac{1}{2} tr(\eta_2 \eta_1 \eta_1^T) - \frac{1}{2} \log(-2\eta_2)$$

$$h(\underline{x}) = (2\pi)^{-k/2}$$

- $k + k^2$ parameters; but less due to symmetric PSD $\Sigma$

ii) Multinomial

- $\underline{x} \in$ Multinomial$(\underline{X} | \underline{\pi})$

$$\underline{X} = \begin{pmatrix} X_1 \\ \vdots \\ X_M \end{pmatrix} \qquad \underline{\pi} = \begin{pmatrix} \pi_1 \\ \vdots \\ \pi_M \end{pmatrix}$$

- Binary(?) (Ex is incorrect) (not sure what lecture notes are referring to)

Jordan (2003): $\underline{X}$ is a collection of integer-valued r.v.s. representing event counts where $X_k$ is no. of times the $k^{th}$ event occurs in $n$ independent trials.

(*) Moment generating properties
   of exp family                                    (scalar form)

· The moments of the relevant distr → obtain via derivatives of
  log normalisation function $A(\eta) = \log(\text{$\int$} He^{\eta\eta}) = \log Z(\eta)$

(*) $\dfrac{dA}{d\eta} = \dfrac{d}{d\eta} \log Z(\eta) = \dfrac{1}{Z(\eta)} \dfrac{d}{d\eta} Z(\eta)$

$\qquad\qquad = \dfrac{1}{Z(\eta)} \dfrac{d}{d\eta} \int h(x) \exp\{\eta T(x)\} \, dx$

(I) $\Big\downarrow$

$\qquad\qquad = \int T(x) \dfrac{h(x) \exp(\eta T(x))}{Z(\eta)} \, dx \qquad\qquad$ (*)

$\qquad\qquad = \mathbb{E}[T(x)] \qquad$ (with respect to?)

(*) $\dfrac{d^2 A}{d\eta^2} = \int T^2(x) \dfrac{h(x) \exp\{\eta^T T(x)\}}{Z(\eta)} \, dx - \int T(x) \dfrac{h(x) \exp\{\eta^T T(x)\}}{Z(\eta)} \, dx \, \dfrac{1}{Z(\eta)} \dfrac{d}{d\eta} Z(\eta)$

$\qquad\qquad = \mathbb{E}[T^2(x)] - \mathbb{E}^2[T(x)]$

$\qquad\qquad = \mathrm{Var}(T(x))$

(0/)

(*) Take derivatives of log-normaliser.
   - $q^{th}$ derivative → $q^{th}$ central moment

$\qquad \dfrac{dA(\eta)}{d\eta}$ - mean $\qquad , \qquad \dfrac{d^2 A(\eta)}{d\eta^2}$ - variance

(*) sufficient statistic - vector ⟹ partial deriv.

(I) clarity on derivation

$\dfrac{dA}{d\eta} = \dfrac{d}{d\eta} \left\{ \log \int \exp\{\eta T(x)\} h(x) \, dx \right\} = \dfrac{\int \dfrac{d}{d\eta} \exp\{\eta T(x)\} h(x) \, dx}{\int \exp\{\eta T(x)\} h(x) \, dx}$

(*) If $p$ If $\pi_k$ the probability of $k^{th}$ event occurring on any given trial:-

$$p(\underline{x}|\underline{\pi}) = \frac{n!}{x_1! x_2! \dots x_m!} \pi_1^{x_1} \pi_2^{x_2} \dots \pi_m^{x_m} \qquad \text{(pmf)}$$

$$\Rightarrow p(\underline{x}|\underline{\pi}) = \exp\left\{\sum_{i=1}^{m} x_i \ln \pi_i\right\} \qquad \text{(trick)}$$

(*) Issue :- factor in probability constraint $\sum_{i=1}^{m} \pi_i = 1$

$\hookrightarrow A(\eta) = 0$?

(*) Parametrise with only $(m-1)$ components of $\underline{\pi}$ :-

$$p(\underline{x}|\underline{\pi}\backslash\pi_m) = \exp\left\{\sum_{i=1}^{m} x_i \ln \pi_i\right\}$$

$$= \exp\left\{\sum_{i=1}^{m-1} x_i \ln \pi_i + \left(1 - \sum_{i=1}^{m-1} x_i\right) \ln\left(1 - \sum_{i=1}^{m-1} \pi_i\right)\right\}$$

$$= \exp\left\{\sum_{i=1}^{m-1} \ln\left(\frac{\pi_i}{1 - \sum_{i=1}^{m-1} \pi_i}\right) x_i + \ln\left(1 - \sum_{i=1}^{m-1} \pi_i\right)\right\}$$

### exp family rep

$$\eta = \begin{bmatrix} \ln\left(\frac{\pi_i}{\pi_m}\right) \\ \\ 0 \end{bmatrix}$$

i.e. $\eta$ has $(m-1)$ components $\eta_i = \ln\left(\frac{\pi_i}{\pi_m}\right)$

and last component $\eta_m = 0$

$$T(\underline{x}) = \underline{x}$$

$$A(\eta) = -\ln\left(1 - \sum_{i=1}^{m-1} \pi_i\right) = \ln\left(\sum_{i=1}^{m} e^{\eta_i}\right)$$

$$h(\underline{x}) = 1$$

(*) Note $\pi_i = \frac{e^{\eta_i}}{\sum_{j=1}^{m} e^{\eta_j}}$  (softmax)  $\pi_i = \text{softmax}_i(\underline{\eta})$

$$= \frac{\int T(x) \exp\{\eta T(x)\} h(x)\, dx}{\int \exp\{\eta T(x)\} h(x)\, dx} = \int T(x) \exp\{\eta^T T(x) - A(\eta)\} h(x)\, dx$$

$$\underbrace{\qquad}_{?}$$

(O/S₁)

$$= \mathbb{E}[T(X)]$$

---

(*) Moment and canonical parameters

(*) Exponential families can have $\begin{cases} \text{canonical parametrisation (via } \eta) \\ \text{moment parametrisation (via } \mu) \end{cases}$

- $\dfrac{dA(\eta)}{d\eta} = \mathbb{E}[T(x)] = \mu$   ;   $\dfrac{d^2 A(\eta)}{d\eta^2} = \text{Var}(T(x)) > 0$   $\left(\begin{array}{l}\text{variance}\\ \text{properties}\\ \rightarrow \text{non-neg}\end{array}\right)$

(*) $A(\eta)$ is convex function

(*) convexity $\Rightarrow$ one-to-one rel between argument $\eta$ and first derivative.

(*) yields an invertible mapping :- Ⓐ

$$\eta = \psi(\mu)$$

---

(*) moment matching, MLE for exponentials

OLS 2-dimensional clarity

- IID data:

- Log likelihood:

$$l(\eta; D) = \log \prod_{n=1}^{N} h(x_n) \exp\{\eta^T T(x_n) - A(\eta)\}$$

$$= \sum_{n=1}^{N} \log h(x_n) + \eta^T \left(\sum_{n=1}^{N} T(x_n)\right) - N A(\eta)$$

$$\nabla_\eta l = \sum_{n=1}^{N} T(x_n) - N \nabla_\eta A(\eta) = 0$$

$$\Rightarrow \nabla_\eta A(\hat\eta) = \frac{1}{N} \sum_{n=1}^{N} T(x_n)$$

- $\mu$ define $\mu = \mathbb{E}[T(x)]$

$$\hat{\mu}_{ML} = \frac{1}{N}\sum_{n=1}^{N} T(x_n)$$

(*) moment matching $\Rightarrow$ can infer canonical param via $\hat{\eta}_{ML} = \psi(\hat{\mu}_{ML})$
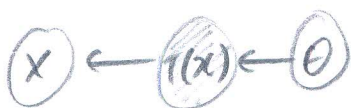not vice versa

---

(*) Sufficiency + Neyman factoriz.

Bayesian:  $\boxed{x} \longrightarrow \circledcirc{T(x)} \longrightarrow \boxed{\theta}$    $p(\theta|T(x), x) = p(\theta|T(x))$  (prob. spec.)

i.e.  $\theta \perp\!\!\!\perp X \mid T(x)$  (C.I rel. spec.)

Freq:  $\boxed{x} \longleftarrow \circledcirc{T(x)} \longleftarrow \boxed{\theta}$    $p(x|T(x), \theta) = p(x|T(x))$

i.e. i.e. $\left( \theta \perp\!\!\!\perp X \mid T(x) \right)$

Neyman fact:

$\boxed{x} \text{---} \circledcirc{T(x)} \text{---} \boxed{\theta}$

⊙ O/S 3

- Reconcile/review
in context of
36-705

- via UGM formalism:- $T(x)$ is <u>sufficient</u> for $\theta$

$$p(x, T(x), \theta) = \psi_1(T(x), \theta)\psi_2(x, T(x))$$

$\Rightarrow p(x|\theta) = g(T(x), \theta)h(x, T(x))$

- proport/norm
constant, absorbed
into one
of $\psi_1/\psi_2$ potentials

- $T(x)$ - determ function
of $x$; drop from LHS; divide
by $p(\theta)$

- for given $g$ and $h$ (functions)

Remarks

(*) can specify either C.I statements or prob. versions
(*) Bayesian intuitions
- $\theta$ is an r.v. $\longrightarrow$ can make C.I. statements involving $\theta$

(*) frequentist intuitions

- Treat $\theta$ as a label rather than r.v.; $T(X)$ is sufficient for $\theta$
if the conditional distr'n $g$ of $X$ given $T(X)$ is <u>not</u> a function of $\theta$.

(*) Both imply a factorisation of $p(x|\theta)$

(*) Neyman factorisation $\longrightarrow$ frequentist definition of sufficiency

(*) Sufficiency in
this context means $T(x)$ is sufficient for $\theta$ if :-

$$\theta \perp\!\!\!\perp X \mid T(X)$$

---

Jordan (2003)

- 8.1.8. ML and KL divergence

- A general rel. between ML and KL divergence (not spec. to exp.)
- Necessary for late lec. material L6, L7
- Statistical interp of KL divergence to illus. rel. between KL and exp. family

(*) Empirical distri: $\hat{p}(x)$

- Places a point mass at each data point $x_n$ in $D$ (dataset)          (discrete)

(*) Empirical distri:  $\hat{p}(x) := \frac{1}{N} \sum_{n=1}^{N} \partial(x, x_n) = \frac{1}{N} \sum_{n=1}^{N} \mathbb{I}(x = x_n)$

(*) Sum/Integrate                                                    $\partial$ - kronecker
$\hat{p}(x)$ against a function of $x$; we evaluate                        delta in cont.
$f$ at each point $x_n$                                                      case.
scaled

(*) log likelihood: (also, cross entropy of $\hat{p}(x)$ and $p(x|\theta)$)

$$\sum_{x} \hat{p}(x) \log p(x|\theta) = \sum_{x} \left( \frac{1}{N} \sum_{n=1}^{N} \partial(x, x_n) \log p(x_n|\theta) \right)$$

$$= \frac{1}{N} \sum_{n=1}^{N} \sum_{x} \partial(x, x_n) \log p(x_n|\theta) \qquad \text{antickifying}$$

$$= \frac{1}{N} \sum_{n=1}^{N} \log p(x_n|\theta)$$

$$= \frac{1}{N} \ell(\theta|D)$$

(*) Note, the scaled log likelihood (by factor $\frac{1}{N}$) is equivalent to the ⊗ cross-entropy between the empirical distri and the model

$$(\hat{p}(x)) \qquad \leftrightarrow \qquad p(x|\theta)$$

- same result for continuous

- KL divergence between empirical, model :-

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

$$= \sum_x \hat{p}(x) \log \hat{p}(x) - \sum_x \hat{p}(x) \log p(x|\theta)$$
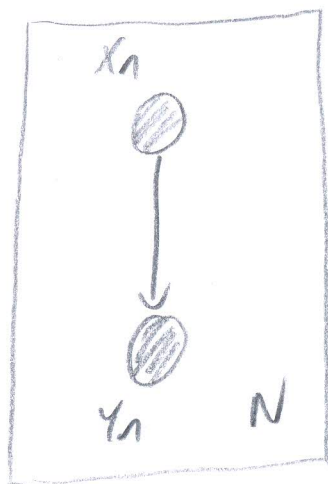
$$= \underbrace{\sum_x \hat{p}(x) \log \hat{p}(x)}_{(I)} - \underbrace{\frac{1}{N} \ell(\theta|D)}_{(II)}$$

(I) - independent of $\theta$

- value of $\theta$ that minimises LHS is the value of $\theta$ that maximises the RHS

⊗ : minimising KL divergence between the empirical distri and model distri is equivalent to maximising the likelihood

_____

(*) Generalised linear Models (GLM) - linear regression/classificatn
　　　　　　　　　　　　§ covers linear reg. / discriminative linear classification.



(*) Both LR/LC → both assume a rep. for conditional expectation of Y.

(*) $\mu = f(\theta^T x) = \mathbb{E}_{y \sim p(Y|f(x))}[Y]$

(*) LR: $f(\cdot)$ - identity
　　 LC: $f(\cdot)$ - sigmoid (logistic)

(*) ALSO: endow Y with a particular cond. prob. distribution, with
$\mu$ as a parameter.

(W): Remember JP $\longrightarrow$ ColumbiaX ML (prob. interp of ML for LR!)

(*) LR - Gaussian V.C. Bernoulli / Multinomial.

_____

(*) Generalised linear model Framework

- 3 assumptions on $p(y|x)$ :-

1. observed input $x$ enters into model via linear comb. $\xi = \theta^T \underline{x}$

2. conditional mean $\mu$ rep as a function $f(\xi)$ of the linear combination
$\xi$ where $f$ is known as the response function

3. observed output $y$ is assumed to be characterised by exp. family
with conditional mean $\mu$.

(a)

$$\theta \quad\quad\quad\quad\quad\quad \overset{f}{\quad} \quad\quad \overset{\psi}{\quad}$$
$$\underline{x} \quad \longrightarrow \xi \longrightarrow \mu \longrightarrow \eta$$