L5 - Paramete estimation (S.2020)

- EX: Paramete estimation on PGMS (a clean reference point)

---

- Start with named/indexed r.v.s.
- Avol training data (IID); fully observed - every r.v. has instantiation

1.) structure learning
   - In principle, possible to learn structure from g data
   - often experts.

2) Paramete estimation
   UGM - Nos on CPT or potential function values
   DiM

---

EX: learning completely observed UGM - fairly trivial

POGMS (partially obs.)
   - directed - focus on this

estimation principles

MLE → classical setting statistics
EX: ML has 'gone beyond this'; present a universal, standardised view that
   [..]fies these.

eg. reinforcement learning - intrinsic/extrinsic
      adversarial learning - adversarial score
- Traditionally, characterise learning/param est in terms of statistical
   consistency etc.
- in more modern engineering ML → may be compromised

---

- Simplest case
- COGMS where structure is known
- Paramete learning for BN

(*) Analytically write down loss fn → likelihood of data, $y$ as
a function of the parameters
- probability (likelihood)? as product of many local terms

EX: Point you to further reading
- Building blocks of GM:                                    W ⓐ
- single node GM (e.g. root node in tree) → supp slides

  - instances of exponential family distri;

(*) ultimately; parnete (probability) is empirical frequency count

2 node graphical models → ⓌⒶ2
                          ↳ supp slides

EX: Patterns
      modifying parnete learning for other use

exponential family

EX: class interested in building blocks.

  $$p(x|\eta) = h(x) \exp \{\eta^T I(x) - A(\eta)\}$$

  $I(\cdot)$ and $\eta(\cdot)$ most important terms; note dot product

  $A(\cdot)$ - log normaliser
            (canonical)
ex: Estimation of parnetes if only require $I(\cdot)$ (i.e. sufficient statistic)

· examples - MVG                          - canonical, moments related
- $\mu$ and $\Sigma$ - moment parnetes          - ⓌⒶ → review
                                                          exp.
exponential family representation

  vec($\cdot$) - ~~increase~~ / turn the entity into a vector

  $\eta$ and $I(\cdot)$ - of same dimensionality

# Multinomial example

- constrained parameters in Multinomial (degrees of freedom) [sum to unity]
- Hence (K-1) in summation

## Exponential family rep.

Q: Why do we go for exponential family reps?

① - Data and parameters cleanly grouped into 2 terms —— through $\eta$
—— through $T(\cdot)$

② Moment generating function.

$$\frac{dA}{d\eta} = \mathbb{E}[T(x)] \qquad \frac{d^2A}{d\eta^2} = Var[T(x)]$$

- Gives standard operator that yields $n^{th}$ order moments (from derivatives of log-normaliser)

EX: Moments important → characterise

③: Relationship between moment and natural parameters

(*) key-Review moment and canonical param relation.

---

## MLE for exponential family

- only differences between distri in exp. family are form of $\eta$ and $T(\cdot)$   i.e. canonical param and suffic stat.

- IID data → log-likelihood → optimise (set 1st order moment to 0)

- moment matching

(through $T(\cdot)$)

---

EX: Exponential family exposes the relationship between data and parameters through (in a linearly dependent fashion)
(through $\eta$)

- Gives info about transformations of data or forms of data we need to worry about to preserve uniqueness and identity of distri.

- e.g. only store sufficient statistic of data

- 3 ways of conceptualising relation between: (dependencies)

X (data)    T(x) (suff. statistic)    θ (param)

- Bayesian: - draw conclusions on parameters given data
  - dependency of parameters from data
  - use posterior   $p(\theta | T(x), x)$

Frequentist - Data generated from unknown true value of param.
  - Parameters impact data only through suff. stat $p(x | T(x), \theta)$

- ∴ influence flows through $T(\cdot)$ for both Bayesian, frequentist
(due to exponential family)

- Neyman factorisation theorem: (W)(A3) → WIKI - check you understand eq.

- exposes sufficiency of $T(x)$ for parameter $\theta$.

- $T(x)$ d-separates X and θ

(*) Density estimation for single r.v. for many diff distr family
  - use sufficient statistic, moment matching, exp. family

- Move onto 2 nodes

- Generalised instance → GLIM

EX: Builds on knowledge of exponential family

- Discrim - logistic regression ; SVMS

  LDA → No it's generative

- Logistic regression - $p(y=1|x) = \dfrac{1}{1 + e^{-\theta^T x}}$

- These are GLIMS
  - But sigmoid → non linearity (taken care of by blanket function)
  - contains linear rel. ; so we use linear techniques with above.

commonality (!): - $\mathbb{E}_p(Y) = \mu = f(\theta^T x)$

- $p(Y|f(X)) \rightarrow$ conditional density of $X \rightarrow$ use exponential family distri.

- $f(\cdot)$ is a response function $\rightarrow$ the $X$ treated in a linear way (dot product)

- (W)(A4): check you understand formulation of GLIM (*)

(*) Different choices of $p()$ and $f()$ - cover all of 2 node GMs.

Linear Regression

Logistic

MRFs : No y, but exponential family distri.

RBMs

etc. ;

(W)(A5): check that you understand these as GLIMs

GLIM (cont.)

- simple modelling principle: (for many instant.)

(*) Begin from data $x$

Assume set of parameters $\underline{\theta}$ (to be estimated) } interact linearly to form signal $\xi$ $\xi$.

(*) Signal $\xi$ turned into mean parameter of cond. distri of output

via response function $f()$

(*) mean parameters and canonical parameters related by $\psi$ (inverse transform for relating distribution instance params and canonical funct. form)

(*) use exponential family to get y.

- mechanical pipeline

(*) Work in exponential form $\rightarrow$ lots of key results

- EX: clear relationship between $f(\cdot)$ and $\psi(\cdot)$ - allows 'cancellation'

$$f = \psi^{-1}(\cdot)$$

- Analytical simplif (?) w (A6)

MLE for GLIMS → natural response

- $f(\cdot)$ and $\psi(\cdot)$ "cancel" → allows simplistic def. of cond. likelihood
  of output given input.

w(17): check this reasoning:- $f(\cdot)$ and $\psi^{-1}(\cdot)$

Yields online learning for canonical GLIM and stochastic gradient ascent

(Q): How is $\psi(\cdot)$ chosen?          (?)
  - Given a distn instance, use default
  - Can specify any (?)
  - Some don't choose $\psi(\cdot)$

(*) S.G.A can be used for any GLIM model (slow)

- Batch learning for GLIM
- Best is Newton's method          → requires knowledge
- requires computation of Hessian        of $\psi$ function

- we already have a library of canonical response functions

IRLS / Newton-Raphson
- Requires Hessian and gradient of loss.
       (inverse)

- reframe update rule to see how it relates to LS loss.
EX: Spirit is exponential family and GLIMS use universally
    (one node and 2-node building blocks for GM).

w(18): check you understand full GLIM/eq formulation
       of logistic, linear.
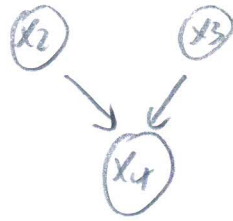
MLE for general BNS:-
(*) Assume global indep. of param; nodes fully observed, decompose BN
    by decomposing log-likelihood into a sum of local terms, one per node.

decomposable likelihood of BN

- graphical illustration

- illustrates analytic decomposition; use of GLIM/exp. methods

EX: How about 2 parents

- Multiplexer function

  $$\prod_{i=3}^{R} X_i^{\partial(X_2, i)}$$

- combine inputs asymmetrically

  $$\log\left(\prod_{i=3}^{R} X_i^{\partial(X_2, i)}\right) \longrightarrow \text{still GLIM}$$

- GLIM/exp?

- Additive $X_1 + X_2$

- Mult.p. $X_1 X_2$

- Multiplexer - (previously pop.)

---

MLE for BNs with tabular CPDS

- Get simple estimators (via same procedure)

---

S. 2020 lecture 5 is a hybrid of S.2019 L5 and L6 material

- so this is technically S.2019 L6 (POGM param est.)

---

- POGMS (partially observed GMs)

EX: too useful in practical e.g. speech HMMs (latent/unobs words)
situations
biolog ev. clustering
(latent clusters)

---

Mixture models

- observe $X$ as 2d co-ord.; no label $(z)$

- estimate $p(z|x)$ generatively via $p(x) \, p(z|x)$

---

- Decided to pause here: concluded run of L6 S.2019