

Review + readingsTargeted

(*) coupling in POGMS, FOGMS (*) clarify
 MARK, MCAK etc.

(*) HMM supervised ML estimation ✓

(*) pseudocounts ✓

(*) mixture models

(*) GMMs

(*) EM for GMMs / soft K-means } (~)

(*) EM - compl., incomplete; expected complete; free-energy / entropy / w-o-oid asst. (C)

(*) Jensen, KL div. ✓

(*) Baum-Welch, BWS, CMM (*) (*) (C)

(*) KL div.

- Wikipedia on pseudocounts (thorough exp.)

- pseudocount / additive smoothing

- "Amt added to no. observed cases in order to change the expected probability in a model of those data, when not known to be 0."

- If frequency of item i is x_i out of N samples:

$$p_{i, \text{empirical}} = \frac{x_i}{N}$$

$$p_{i, \alpha\text{-smoothed}} = \frac{x_i + \alpha}{N + \alpha d}$$

Posterior probability after add. smoothing; effectively increases count x_i by α a priori.

- Pseudocount \rightarrow any non-neg. finite value

- rel. value of pseudocounts represent relative prior exp. prob. of their poss. (?)

- motivated in lecture as a way of addressing natural tendency (Bishop) of ML est point estimation to overfit

- used in Naive Bayes (known as Laplace smoothing)

(*) Sum of pseudocounts, may be large, represents estimated weight of prior knowledge compared with all the actual obs. (one for each) when determining expected probability.

(*) HMM supervised estimation;
comfortable with tricky unobserved states setting

Materials (readings)

Jordan (2003) :- Ch 10 mixtures and cond. mixtures

(covers GMM, K-means, mixture; heuristic pres. of EM)
+ CMMs + conj grad, N-R methods

Jordan (2003): Ch 11 EM (*)

- formal pres
+ interpretations

Neal + Hinton (20) : statistical physics interpretation (*)

Borman

Koller (2009) :- Ch 19

Jordan (2003) Ch 11 - EM

- EM central to graphical models
- divide and conquer

⑦: complex dependencies; model "top-down" using latent variables

- unobserved latent variables :

(*) likelihood is marginal probability; obtained by summing/integrating over latent variables

(*) marginalisation couples parameters, obscures underlying structure in likelihood fn.

- E - "Inference" → compute expected sufficient statistics

↳ for multinomial (latent) variables;
reduces to computing probability

↳ OR: calculating probability of latent variables
given observed variables and current parameter values.

- M-step: - update parameters based on inferred latent variables

11.2 - General setting

- \underline{X} - observables \underline{z} - latent
- often $\underline{X}, \underline{z}$ decompose into sets of IID pairs
- $\underline{X} = \begin{pmatrix} x_1 \\ \vdots \\ x_N \end{pmatrix}$ $\underline{z} = \begin{pmatrix} z_1 \\ \vdots \\ z_N \end{pmatrix}$ - Note: the x_i are IID variables
- observations

$$\underline{x} = \begin{pmatrix} x_1 \\ \vdots \\ x_N \end{pmatrix}$$

- \underline{X} - totality of observed variables \underline{x} - entire observed dataset.
- \underline{z} - set of all latent

If \underline{z} observed; MLE estimation:-

complete log likelihood: - $\ell_c(\theta; \underline{x}, \underline{z}) = \log p(\underline{x}, \underline{z} | \theta)$ (11.1)

- MLE maximises (11.1); if $p(\underline{x}, \underline{z} | \theta)$ factors in some way, such that sep. components of θ occur in separate factors ^(*), log has effect of separating likelihood into terms that can be maximised independently

(*) This is what is meant by the est. prob. decouples

- \underline{z} is not observed, hence; (11.2)

incomplete log-likelihood: - $\ell(\theta; \underline{x}) = \log p(\underline{x} | \theta) = \log \sum_{\underline{z}} p(\underline{x}, \underline{z} | \theta)$ (11.2)

- probability of data: \underline{x} - marginal

- No decoupling: log separated from $p(\underline{x}, \underline{z} | \theta)$ by summation.

\underline{z} not observed \Rightarrow complete log likelihood is a random quantity; (note \underline{x} is obs.)

- complete-log like cannot be maximised directly

(*) Average out \underline{z} to remove randomness (using an "averaging distri")

the expected-complete log likelihood (*)

$$\langle \ell(\theta; x, z) \rangle_q = \sum_z q(z|x, \theta) \log p(x, z|\theta) \quad (11.3)$$

(*) This is a deterministic function of parameters θ

⑥: If q is well chosen; perhaps expected complete log like (ECLL) will not be far from the log-likelihood (which one?) (probably ILL); serves as surrogate for ILL*

⑥: Maximising surrogate does not guarantee a value of θ that maximises likelihood.

BUT (*) ⑥: If it may yield an improvement from an initial value of θ

(*) If so iterate process and hill-climb

(*) use of averaging distri $q(z|x)$ can provide lower bound on log-likelihood ()

$$\begin{aligned} \ell(\theta; x) &= \log p(x|\theta) \\ &= \log \sum_z p(x, z|\theta) \\ &= \log \sum_z q(z|x) \frac{p(x, z|\theta)}{q(z|x)} \\ &\geq \sum_z q(z|x) \log \frac{p(x, z|\theta)}{q(z|x)} \end{aligned}$$

"ILL (marginal) is a
CLL (joint); marginalised
over latent"
"importance sampling trick"

(*)

$$:= \mathcal{L}(q, \theta) =$$

(*) Jensen: - $\log(\cdot)$ is concave $\Rightarrow \log[\mathbb{E}_q[\cdot]] \geq \mathbb{E}_q[\log[\cdot]]$

• $\mathcal{L}(q, \theta)$ - auxiliary function

• for arbitrary $q(\cdot)$ distri; auxiliary function $\mathcal{L}(q, \theta)$ is a lower bound for the (incomplete) log likelihood.

EM(*) EM as coordinate ascent on auxiliary $\mathcal{L}(q, \theta)$

E-step: $q^{(t+1)} = \arg\max_q \mathcal{L}(q, \theta^{(t)})$

M-step: $\theta^{(t+1)} = \arg\max_{\theta} \mathcal{L}(q^{(t+1)}, \theta)$

(*) Max lower bound on $\mathcal{L}(\theta; x)$ (LL) can maximise $\mathcal{L}(\theta; x)$

(*) M-step as maximisation of ELLL

- recompute auxiliary $\mathcal{L}(q, \theta^{(t+1)})$ / (lower bound)

$$\mathcal{L}(q, \theta) = \sum_z q(z|x) \frac{\log p(x, z | \theta)}{q(z|x)}$$

$$= \sum_z q(z|x) \log p(x, z | \theta) - \sum_z q(z|x) \log q(z|x)$$

$$= \langle \mathcal{L}(\theta; x, z) \rangle_q - H_q$$

(*) $\mathcal{L}(q, \theta)$ is made up of expected-complete log likelihood

(*) And also H_q , which is both indep of param θ

(*) H_q - entropy of $q(\cdot)$; - KL div of distri with itself (?)

(*) E-step: setting $q(\cdot)$ as posterior distri over latent variables, given data and parameters yields maximum: $q^{(t+1)}(z|x) = p(z|x, \theta^{(t)})$

Evaluate $\mathcal{L}(q, \theta^{(t)})$ for $q^{(t+1)} = p(z|x, \theta^{(t+1)})$

$$\mathcal{L}(p(z|x, \theta^{(t)}) | \theta^{(t)}) = \sum_z p(z|x, \theta^{(t)}) \log \frac{p(x, z | \theta)}{p(z|x, \theta^{(t)})}$$

$$= \sum_z p(z|x, \theta^{(t)}) \log p(x | \theta^{(t)})$$

$$= \log p(x | \theta^{(t)})$$

$$= \mathcal{L}(\theta^{(t)}; x)$$

$\theta^{(t)}$?

stackexchange

Bayes rule

$$\mathbb{E}_p(z|x, \theta^{(t)}) [\log p(x | \theta^{(t)})]$$

(*) $\mathcal{L}(\theta; x)$ is an upper bound for $\mathcal{L}(q, \theta^{(t)}) \Rightarrow \mathcal{L}(q, \theta^{(t)})$ maximised by setting $q(\cdot) = p(z|x, \theta^{(t)})$

setting
(*) M-step: $q(\cdot) = p(z|x, \theta^{(t)})$ maximises $l(\theta; x)$ and hence $L(q, \theta^{(t)})$
via KL-divergence

(*) Appendix material on KL-divergence perspective & alternating minimisation

(*) Intuition/Exp: -

(*) $p(z|x, \theta^{(t)})$ as best guess of latents, conditioned on data

(*) Use best guess distn to compute ELL (E-step)

(*) Maximise ELL wrt parameters to yield new $\theta^{(t+1)}$ (M-step)

(*) Given improvement, make better guess $p(z|x, \theta^{(t+1)})$; iterate.

(*) Intuition: -

- effect of EM iteration on log-likelihood $l(\theta; x)$?

- M-step: select θ parameters to increase a lower bound ($L(q^{(t+1)}, \theta)$)

⊗ ⊗ , on likelihood $l(\theta; x)$

- Increasing lower bound on function \nrightarrow increasing function itself

(*) BUT: in E-step, close gap with approp. choice of $q(\cdot)$ distn

$$l(\theta^{(t)}; x) = L(q^{(t+1)}, \theta^{(t)})$$

Ⓜ: some graphics / illustr
would help

(*) EM- will climbing

in log likelihood $l(\theta; x)$

- ~~via~~ Indirect will climbing by co-ordinate ascent in auxiliary $L(q, \theta)$

- rather \rightarrow maximisation of ELL rather than LL