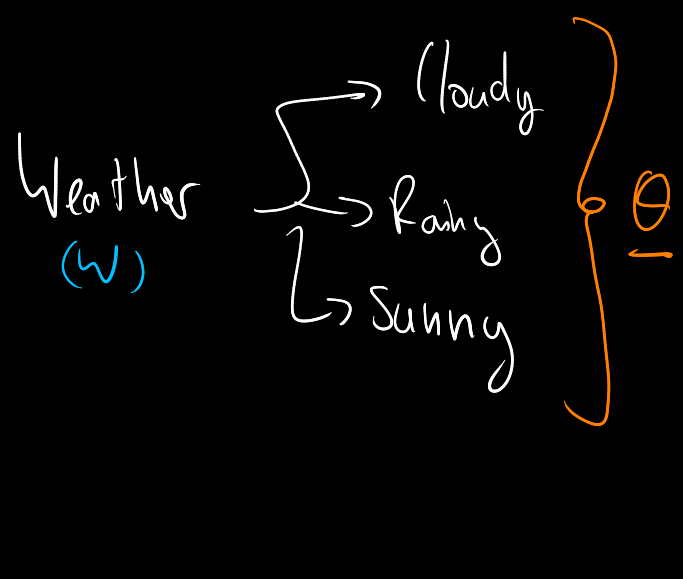


Posterior and MAP Estimate for the Categorical



observe the weather

$$D = \{C, R, S, S, S, R, C, S, \dots\}$$

infer θ from D by MLE

→ is there a robust estimate?

→ can we encode prior knowledge

$$w \sim \text{Cat}(\theta)$$

$$w \in \{C, R, S\}$$

$$\begin{matrix} 0 & 1 & 2 \\ [1, 0, 0]^T & [0, 1, 0]^T & [0, 0, 1]^T \end{matrix}$$

DBM

we have a dataset with N weather obs



$$\theta \sim \text{Dir}(\alpha; \alpha)$$

$$w_i \sim \text{Cat}(w_i | \theta)$$

$$\text{joint } p(\theta, \underline{w}) = p(\theta) \cdot p(\underline{w} | \theta)$$

$$\stackrel{\text{i.i.d.}}{=} p(\theta) \prod_{i=0}^{N-1} p(w_i | \theta)$$

$$p(\theta, \underline{w} = D) = p(\theta) \cdot \prod_{i=0}^{N-1} p(w_i = w^{(i)} | \theta)$$

$p(\theta, D)$

prior

$p(D | \theta)$

↳ likelihood

posterior $p(\theta | D) \stackrel{\text{Bayes' Rule}}{=} \frac{p(D | \theta) p(\theta)}{p(D)}$

$$\sim p(D | \theta) \cdot p(\theta)$$

$$= p(\theta) \cdot \prod_{i=0}^{N-1} p(w_i = w^{(i)} | \theta)$$

all w_i are distributed identically

$$= \text{Dir}(\theta; \alpha) \cdot \prod_{i=0}^{N-1} \text{Cat}(w = w^{(i)} | \theta)$$

$$= \frac{\prod_{d=0}^{D-1} \Gamma(\sum_{i=0}^{N-1} \theta_d^{(i)})}{\prod_{d=0}^{D-1} \Gamma(\alpha_d)} \prod_{d=0}^{D-1} \theta_d^{\alpha_d-1} \cdot \prod_{i=0}^{N-1} \prod_{d=0}^{D-1} \theta_d^{\mathbb{I}(w^{(i)}=d)}$$

(we have D states)

$$\sim \prod_{d=0}^{D-1} \theta_d^{\alpha_d-1} \cdot \prod_{d=0}^{D-1} \prod_{i=0}^{N-1} \theta_d^{\mathbb{I}(w^{(i)}=d)}$$

$$\text{e.g. } D = \{1, 0, 0, 2, 3, 0, \dots\}$$

$$\theta_0^0 \cdot \theta_0^1 \theta_0^1 \theta_0^0 \theta_0^0 \theta_0^1$$

$$= \prod_{d=0}^{D-1} \theta_d^{\alpha_d-1} \cdot \theta_d^{\sum_{i=0}^{N-1} \mathbb{I}(w^{(i)}=d)}$$

$$N_d = \sum_{i=0}^{N-1} \mathbb{I}(w^{(i)}=d)$$

"the number of times we observe the d -th state"

$$= \prod_{d=0}^{D-1} \theta_d^{\alpha_d-1} \cdot \theta_d^{N_d}$$

$$= \prod_{d=0}^{D-1} \theta_d^{\alpha_d + N_d - 1}$$

(can you see why the Dirichlet is the conjugate prior?)

$$\therefore \mathcal{M}(D; \theta)$$

"posterior likelihood"

Full posterior

→ full distribution

$$p(\theta | D) \sim \mathcal{M}(D; \theta)$$

$$\alpha_d' := \alpha_d + N_d$$

$$p(\theta | D) \sim \prod_{d=0}^{D-1} \theta_d^{\alpha_d' - 1}$$

isn't it just a Dirichlet? *differentiation*

$$p(\theta | D) = \text{Dir}(\theta; \alpha')$$

$$p(\theta | D) = \frac{\prod_{d=0}^{D-1} \Gamma(\sum_{i=0}^{N-1} \alpha_d^{(i)})}{\prod_{d=0}^{D-1} \Gamma(\alpha_d)} \prod_{d=0}^{D-1} \theta_d^{\alpha_d-1}$$

MAP estimate
Maximum A Posteriori

→ point estimate

$$\theta_{\text{MAP}}^* = \underset{\theta}{\text{argmax}} (\mathcal{M}(D; \theta))$$

$$\text{s.t. } \sum_{d=0}^{D-1} \theta_d = 1$$

→ "log posterior likelihood"

$$\ln(D; \theta) = \ln(\mathcal{M}(D; \theta))$$

$$= \sum_{d=0}^{D-1} (\alpha_d + N_d - 1) \cdot \ln(\theta_d)$$

Build Lagrangian to include constraint

$$\tilde{m}(D; \theta, \lambda) = \ln(D; \theta) + \lambda (1 - \sum_{d=0}^{D-1} \theta_d)$$

$$= \sum_{d=0}^{D-1} ((\alpha_d + N_d - 1) \cdot \ln(\theta_d)) + \lambda \cdot (1 - \sum_{d=0}^{D-1} \theta_d)$$

↳ Take derivative and set to 0

$$\frac{\partial \tilde{m}}{\partial \theta_e} = (\alpha_e + N_e - 1) \cdot \frac{1}{\theta_e} - \lambda \stackrel{!}{=} 0$$

e -th component of θ

$$\Leftrightarrow \lambda = \frac{\alpha_e + N_e - 1}{\theta_e} \quad (*)$$

(also holds for $(e-1), (e+1)$)

$$\lambda = \frac{\alpha_d + N_d - 1}{\theta_d}$$

$$\frac{\partial \tilde{m}}{\partial \lambda} = 1 - \sum_{d=0}^{D-1} \theta_d \stackrel{!}{=} 0$$

$$\Leftrightarrow \sum_{d=0}^{D-1} \theta_d = 1$$

How to eliminate λ

$$\lambda = \lambda = \lambda \cdot 1 = \lambda \cdot \sum_{d=0}^{D-1} \theta_d = \sum_{d=0}^{D-1} \lambda \theta_d$$

$$= \sum_{d=0}^{D-1} \frac{\alpha_d + N_d - 1}{\theta_d} \cdot \theta_d$$

$$= \sum_{d=0}^{D-1} (\alpha_d + N_d - 1)$$

$$= \sum_{d=0}^{D-1} (\alpha_d) + \sum_{d=0}^{D-1} (N_d) - \sum_{d=0}^{D-1} (1)$$

$$= N - D + \sum_{d=0}^{D-1} \alpha_d = \lambda$$

plug into (*)

$$N - D + \sum_{d=0}^{D-1} \alpha_d = \frac{\alpha_e + N_e - 1}{\theta_e}$$

MAP

$$\theta_e = \frac{N_e + \alpha_e - 1}{N - D + \sum_{d=0}^{D-1} \alpha_d}$$

"pseudo-counts"

$$\text{with } N_e = \sum_{i=0}^{N-1} \mathbb{I}(w^{(i)} = e)$$

(for One-Hot-Categorical $N_e = \sum_{i=0}^{N-1} w_e^{(i)}$)