

today: • exponential family
• sampling

Exponential family

a (flat/canonical) exponential family on X

is a parametric family of distributions defined by two quantities

I) $\underbrace{h(x)}_{\text{reference density}} \underbrace{d\mu(x)}_{\text{base measure}} \rightarrow$ reference measure on X

reference density \swarrow base measure \searrow counting (discrete f.v.)
Lebesgue (cts. f.v.)

II) $T: X \rightarrow \mathbb{R}^P$ called "sufficient statistics" vector
(aka feature vector)

members of family will have dist.:

$$p(x; \eta) d\mu(x) = \exp(\underbrace{\eta^T T(x)}_{\text{"canonical parameter"}} - \underbrace{A(\eta)}_{\text{log normalizer, log partition function or cumulant generating function}}) \underbrace{h(x) d\mu(x)}_{\text{defining pieces } (+ \Omega_X)}$$

If Ω_X is discrete, then $p(x; \eta)$ is a pmf

" " cts. , " " " a pdf

$$\star \text{ want } 1 = \int_X p(x; \eta) d\mu(x) = \int_X \exp(\eta^T T(x)) e^{-A(\eta)} h(x) d\mu(x)$$

$$\Rightarrow A(\eta) \triangleq \log \left(\underbrace{\int_{\mathcal{X}} \exp(\eta^T T(x)) h(x) d\mu(x)}_{Z(\eta)} \right)$$

domain $\Omega \triangleq \{ \eta \in \mathbb{R}^p \mid A(\eta) < \infty \}$
 (set of valid canonical parameters)

* more generally, consider a reparameterization of a subset of the family
 by defining the mapping $\eta : \mathcal{H} \rightarrow \Omega$
 (new set of parameters)

$$\text{consider } p(x; \theta) \triangleq p(x; \eta(\theta)) \text{ for } \theta \in \Theta$$

(get a "curved exponential family" if $\eta(\Theta)$ is a curved manifold in Ω)

* two examples not an exponential family:
 • mixture of Gaussians
 • $\text{Unif}(0, \theta)$

but most standard distributions are (Poisson, Gamma, Gaussian, etc...)

Example: (Multinomial)

$$X \sim \text{Mult}(\pi)$$

$$X = \{0, 1\}^K$$

$$\Omega_X = \Delta_K \cap X$$

parameter form $\pi \in \Delta_K$ suppose $\pi_i > 0 \forall i$

$$p(x; \pi) = \prod_{j=1}^k \pi_j^{x_j} \stackrel{\downarrow}{=} \exp\left(\sum_{j=1}^k x_j \log \pi_j\right)$$

think of this as a "θ"

$$= \exp(\eta(\pi)^T x - \psi(\pi))$$

we have $\eta_j(\pi) = \log \pi_j$

$$T(x) = x$$

$$d\mu(x) = \text{counting measure on } \mathcal{X}$$

$$h(x) = \mathbb{1}_{\{x \in \Omega_x\}} = \mathbb{1}_{\{x \text{ has exactly one entry equal to 1}\}}$$

$$\Theta = \text{int}(\Delta_k)$$

here, $A(\eta(\pi)) = 0 \quad \forall \pi \in \Theta$

but here $\Omega = \mathbb{R}^k$

remark: $\Theta \rightarrow \text{dimension } k-1$

$\eta(\Theta) \rightarrow \text{" } k-1$

$\Omega \rightarrow \text{dimension } k$
 (not be confused with Ω_x)

note:

for any x s.t. $h(x) \neq 0$

$$\sum_{j=1}^k T_j(x) = 1$$

affine linear dep. between components of T

\Rightarrow multiple π 's rep to same distribution "overparametrization"

\hookrightarrow not a minimal exponential family

* for multinomial, minimal exp family

$$T(x) = \begin{pmatrix} x_1 \\ \vdots \\ x_{k-1} \end{pmatrix}$$

$$Z(\eta) = \sum_{x \in \Omega_x} \exp(\eta^T T(x)) = \sum_{j=1}^{k-1} e^{\eta_j} + 1$$

$$p(x; \eta) = \exp\left(\sum_{j=1}^{k-1} \eta_j x_j - \underbrace{\log\left(\sum_{j=1}^{k-1} \exp(\eta_j) + 1\right)}_{A(\eta)}\right)$$

recall: $\nabla_{\eta} A(\eta) = \mathbb{E}_{p(x; \eta)} [T(x)]$ (valid for $\eta \in \text{int}(\Omega)$)

for multinomial, $\frac{\partial A(\eta)}{\partial \eta_j} = \frac{1 \cdot \exp(\eta_j)}{Z(\eta)} = p(x=j | \eta) = \mathbb{E}_{p(x; \eta)} [T_j | x]$ as required

⊗ $T(x)$ is called "sufficient" as in statistics:

$T: x \mapsto T(x)$ is "sufficient" for a parametric family \mathcal{P}_{Θ}

iff $\forall \theta \in \Theta \quad p_{\theta}(x) = \underbrace{h(x)}_{\text{fixed for family}} g(T(x); \theta)$

i.e. dependence on θ only happens through $T(x)$

i.i.d. model in exp. family:

$$\log p(x_1, \dots, x_n | \eta) = \sum_{i=1}^n [\log h(x_i) + \eta^T T(x_i) - A(\eta)] \quad \text{Total } h(x_i)$$

\leadsto exp. family on (x_1, \dots, x_n) with reference density $\tilde{h}(x) = \prod h(x_i)$

sufficient statistics $T(x) = \sum_{i=1}^n T(x_i)$
 log-partition $A(\eta) = n A(\eta)$

Example 2: (1D) Gaussian

$$X \sim N(\mu, \sigma^2) \quad X = \mathbb{R}$$

$$p(x; (\mu, \sigma^2)) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

$$= \exp\left(-\frac{x^2}{2} \left[\frac{1}{\sigma^2}\right] + x \left[\frac{\mu}{\sigma^2}\right] - \left[\frac{\mu^2}{2\sigma^2} + \frac{1}{2} \log(2\pi\sigma^2)\right]\right)$$

$$T(x) = \begin{bmatrix} x \\ -\frac{x^2}{2} \end{bmatrix} \quad \eta(\theta) = \begin{bmatrix} \mu/\sigma^2 \\ 1/\sigma^2 \end{bmatrix} = \begin{bmatrix} \eta_1 \\ \eta_2 \end{bmatrix}$$

$$\eta_2 = \frac{1}{\sigma^2} = \text{precision} > 0$$

$$\eta_1 = \eta_2 \cdot \mu \quad \Omega = \{(\eta_1, \eta_2) : \eta_2 > 0\}$$

Example 3: UGM?

let $p \in \mathcal{P}(G)$ G is undirected

with $\psi_c(x_c) > 0$

$$p(x) = \frac{1}{Z} \prod_{c \in \mathcal{C}} \psi_c(x_c) = \exp\left(\sum_c \ln \psi_c(x_c) - \log Z\right)$$

$$= \exp\left(\sum_{c \in \mathcal{C}} \sum_{y_c \in \mathcal{X}_c} \underbrace{\mathbb{1}\{y_c = x_c\}}_{T_c(y_c|x)} \underbrace{\log \psi_c(x_c)}_{\eta_{c,x_c}} - \log Z\right)$$

for given value of c & x_c

for every value of C & y_C

I have $T_{C, y_C}(x) = \mathbb{1}\{y_C = x_C\}$

$$T(x) = \begin{pmatrix} \mathbb{1}\{x_C = y_C\} \\ \vdots \end{pmatrix}_{\substack{y_C \in \mathcal{X}_C \\ C \in \mathcal{C}}} \quad n_{C, y_C} = \log \psi_C(y_C)$$

$$n(C) = \begin{pmatrix} \log \psi_C(y_C) \\ \vdots \end{pmatrix}_{\dots}$$

notes: a) Mult(x) is special case with complete graph (1 big clique)

b) feature perspective: instead of using all indicators $\mathbb{1}\{x_C = y_C\}$
you could choose a subset

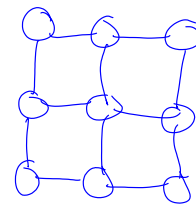
for example; suppose x_i is a word

feature on x_i & x_{i+1} $\mathbb{1}\{x_i \text{ is a verb} \& x_{i+1} \text{ is a noun}\}$

c) binary Ising model

$$x_i \in \{0, 1\} \quad |C| \leq 2$$

suppose use nodes & edges as cliques



\Rightarrow dimension of $T(x)$ is $2|V| + 4|E| \rightarrow$ "overparameterized"

a minimal representation:

$$T(x) = \begin{pmatrix} (x_i)_{i \in V} \\ (x_i x_j)_{\{i,j\} \in E} \end{pmatrix} \rightarrow \dim = |V| + |E|$$

$$\text{ie. } p(x; \eta) = \frac{1}{Z} \exp \left(\sum_{\{i,j\} \in E} \eta_{ij} x_i x_j + \sum_{i \in V} \eta_i x_i \right)$$

def: $p(x; \eta)$ is minimal exp family iff there does not exist $\{a_j\}$ s.t. $\sum_{j=1}^p a_j T_j(x) + a_{p+1} = 0 \quad \forall x \text{ s.t. } h(x) \neq 0$

Properties of A :

- $A(\cdot)$ is C^∞ for $\eta \in \text{int}(\Omega)$
- $A(\cdot)$ is convex (Ω is a convex set)
- $\nabla_\eta A(\eta) = \mathbb{E}_{p(x; \eta)} [T(x)] \triangleq \mu(\eta)$ "moment vector"
for $\eta \in \text{int}(\Omega)$
- $\nabla^2 A(\eta) = \mathbb{E}_{p(x; \eta)} [(T(x) - \mu(\eta))(T(x) - \mu(\eta))^T] = \text{cov}(T(x))$
(Hessian)
(proof as exercise)
- $Z(\eta)$ is the mgf for $p(x; \eta)$ $t \mapsto Z(\eta + t)$ for $t \in \mathbb{R}^p$
 $\hookrightarrow \mathbb{E} \exp(t^T T(x))$

$A(\eta) = \log \text{mgf} \rightarrow$ cumulant generating functions
 \rightarrow why derivatives of A give cumulants

1st derivative give mean
 2nd " " covariance
 3rd " " 3rd-cumulant etc...

Approximate inference \rightarrow sampling

why sampling?

$$X = (x_1, \dots, x_p)$$

a) simulation: $X^{(i)} \sim p$

b) approximate $f(x_j)$

\rightarrow more generally,

consider $f: \mathbb{R}^p \rightarrow \mathbb{R}^d$, approximate $\mu = \mathbb{E}_p[f(X)]$

e.g. if $f(X) = \mathbb{1}\{X_A = x_A\}$ $\mathbb{E}_p[f(X)] = p(X_A = x_A)$

Monte-Carlo integration / estimation \rightarrow appears in physics, applied math, ML, etc...

to approximate $\mu = \mathbb{E}_p[f(X)]$

alg:
 • n samples $X^{(i)} \stackrel{\text{iid}}{\sim} p$
 • estimate: $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n f(X^{(i)})$

this is true even if
 $X^{(i)}$ are dependent

properties: 1) unbiased $\mathbb{E}[\hat{\mu}] = \frac{1}{n} \sum_{i=1}^n \mathbb{E} f(X^{(i)}) = \mu$

$\textcircled{2}$ expected error: $\mathbb{E}[\|\hat{\mu} - \mu\|^2] = \mathbb{E}\left[\frac{1}{n^2} \sum_{i,j} \langle f(x^{(i)}) - \mu, f(x^{(j)}) - \mu \rangle\right]$
 $\text{tr}(\text{cov}(\hat{\mu}, \hat{\mu}))$

$\xrightarrow{\text{by independence}} \text{off-diagonal terms are zero}$
 $= \frac{1}{n^2} \sum_{i=j} \mathbb{E}[\langle f(x^{(i)}) - \mu, f(x^{(i)}) - \mu \rangle]$
 $\mathbb{E}[\|f(x) - \mu\|^2]$
 $\triangleq \sigma^2 = \text{tr}(\text{cov}(f(x), f(x)))$

$$\boxed{\mathbb{E}[\|\hat{\mu} - \mu\|^2] = \frac{\sigma^2}{n}}$$

[aside: asymptotic analysis LLN $\hat{\mu} \xrightarrow{as} \mu$

CLT

$$\sqrt{n}(\hat{\mu} - \mu) \xrightarrow{d} N(0, \sigma^2)$$