

Exponential Families

Exponential families are perhaps the most important statistical models, and nearly every model we discuss in this course is an exponential family. Today we will examine exponential families carefully, examining carefully their algebraic and combinatorial structure and how we have already seen it appearing

Regular Exponential Families

- X a sample space with σ -algebra \mathcal{A} with (σ -finite) measure ν .
- $T: X \rightarrow \mathbb{R}^k$ a statistic (i.e. a measurable map)
- $h: X \rightarrow \mathbb{R}_{>0}$ a measurable function
- The natural parameter space is

$$N = \left\{ \eta \in \mathbb{R}^k : \int_X h(x) e^{\eta^T T(x)} d\nu(x) < \infty \right\}$$

- For $\eta \in N$ get a probability density function P_η on X :

$$P_\eta(x) := h(x) e^{\eta^T T(x) - \phi(\eta)}, \text{ where}$$

$$\phi(\eta) := \log \int_X h(x) e^{\eta^T T(x)} d\nu(x).$$

- $P_\eta :=$ probability measure on (X, \mathcal{A}) with ν -density P_η .

- set $\nu^T := \nu \circ T^{-1}$, a measure on the Borel σ -algebra \mathcal{B}_k of \mathbb{R}^k :
$$\nu^T(B) = \nu(T^{-1}(B)) \text{ for } B \in \mathcal{B}_k$$

- The support of ν^T is the set

$$\text{supp}(\nu^T) := \bigcap A \quad \left(\begin{array}{l} A \subseteq \mathbb{R}^k \text{ closed:} \\ \nu^T(\mathbb{R}^k \setminus A) = 0 \end{array} \right) \quad \left(\begin{array}{l} \text{Part of } \mathbb{R}^k \text{ carrying} \\ \text{info about } X \dots \end{array} \right)$$

(smallest σ -algebra containing all $T_k(-a_i, b_i)$)

Definition (16) Let k be a positive integer. The probability distributions

form a regular exponential family of order k if N is an open subset in \mathbb{R}^k and the affine dimension of $\text{supp}(\nu^T)$ is equal to k . The statistic T inducing \mathcal{P} is called the canonical sufficient statistic.

- It is clear from the definition that T is a sufficient statistic, but in fact T has a stronger property when \mathcal{P} is regular that corresponds to the condition on the dimension of $\text{supp}(v^T)$:
- The canonical sufficient statistic T of a regular exponential family is minimal: For any sufficient statistic T' and for every $x, y \in \mathcal{X}$

$T(x) = T(y)$ whenever $T'(x) = T'(y)$.
 That is, T is a function of T' : $\exists f$ s.t. $T(x) = f(T'(x))$.

- It turns out that ~~any~~ any exponential family that is not regular can be reduced by reparametrization to a regular exponential family of order $< k$. Here, the statistic for the model is also transformed into a minimal one...
- Today, we want to understand the algebra and geometry of exponential families that are submodels of regular exponential families. We start in the discrete setting.

Discrete Exponential Families

Example 1 (Discrete Random Variables)

- $\mathcal{X} := [r] = \{1, \dots, r\}$.
- $v =$ the counting measure on $[r]$: $v(A) = |A| \quad \forall A \subseteq [r]$.
- $T: [r] \rightarrow \mathbb{R}^{r-1}$ where $T: i \mapsto (\delta_{i=1}, \delta_{i=2}, \dots, \delta_{i=r-1})^t \in \mathbb{R}^{r-1}$,
 where $\delta_{i=j} = \begin{cases} 1 & \text{if } i=j \\ 0 & \text{otherwise} \end{cases}$.

$h(x) \equiv 1$.

Note: $\text{Supp}(v^T) = \{0, e_1, \dots, e_{r-1}\} \subseteq \mathbb{R}^{r-1}$ since
 $v^T(B) = v(T^{-1}(B)) = 0 \iff T^{-1}(B) \cap \{0, e_1, \dots, e_{r-1}\} = \emptyset$.

$\Rightarrow \dim(\text{supp}(v^T)) = r-1$.

Given $\eta \in \mathbb{R}^{r-1}$

$$\phi(\eta) = \log \int_{\mathcal{X}} h(x) e^{\eta^T T(x)} d\nu(x) = \log \left(\sum_{i=1}^r e^{\eta^T T(i)} \right) = \log \left(1 + \sum_{i=1}^{r-1} e^{\eta_i} \right) < \infty.$$

\Rightarrow The natural parameter space is $N = \mathbb{R}^{r-1}$ (and hence N open)

$\Rightarrow \mathcal{P} = \{ p_\eta : p_\eta(i) = e^{\eta^T T(i) - \phi(\eta)} \quad \forall i \in [r] \}$
 is a regular exponential family of order $r-1$.

• We claim $\mathcal{P} \cong \Delta_{r-1}^0$: p_η is a vector in \mathbb{R}^r . For $1 \leq i \leq r-1$

$$p_\eta(i) = e^{\eta^T T(i) - \phi(\eta)} = \frac{e^{\eta_i}}{1 + \sum_{i=1}^{r-1} e^{\eta_i}}, \text{ and}$$

$$p_\eta(r) = 1 - \sum_{i=1}^{r-1} p_\eta(i) = \frac{1}{1 + \sum_{i=1}^{r-1} e^{\eta_i}}.$$

• Given $(p_1, \dots, p_r) \in \Delta_{r-1}^0$, taking $\eta_i = \log\left(\frac{p_i}{p_r}\right)$ establishes $\mathbb{R}^{r-1} = N \cong \Delta_{r-1}^0$.

• Hence all discrete exponential families are submodels of a regular exponential family.

• To get such a submodel, we consider the following:

• $T: [r] \rightarrow \mathbb{R}^k$ a measurable map. \iff an assignment of vectors $T(i) \in \mathbb{R}^k \quad \forall i \in [r], \dots$

• $h: [r] \rightarrow \mathbb{R}_{>0}$ a vector:
 $h = (h_1, \dots, h_r)^t \in \mathbb{R}_{>0}^r \quad (h_i := h(i))$

• Since measurability is trivial in the discrete case, any collection of vectors works for T and h is any positive vector.

• Each choice $\{T(i) \in \mathbb{R}^k : i \in [r]\}$ and $h \in \mathbb{R}_{>0}^r$ defines an exponential subfamily of Δ_{r-1}^0 .

• Let $Z(\eta) := \sum_{i \in [r]} h_i e^{\eta^T T(i)}$

• Write $T(i) := (a_{i1}, \dots, a_{ik})^t \in \mathbb{R}^k$ and $\eta = (\eta_1, \dots, \eta_k)^t \in \mathbb{R}^k$.

Then $\forall i \in [r]$,

$$P_{\eta}(i) = h(i) e^{\eta^T T(i) - \phi(\eta)} = h_i e^{\eta^T T(i) - \log(Z(\eta))}$$

$$= \frac{h_i e^{\eta^T T(i)}}{Z(\eta)} = \frac{1}{Z(\eta)} h_i \prod_{j \in [k]} e^{\eta_j a_{ji}}$$

• Setting $\Theta_j := \exp(\eta_j) \quad \forall j \in [k]$, we get

$$P_{\Theta}(i) = \frac{1}{Z(\Theta)} h_i \prod_{j \in [k]} \Theta_j^{a_{ji}}$$

• Hence, given any T and h the resulting exponential family is

$$\mathcal{P}_{T,h} := \left\{ P_{\eta} : \eta \in \mathbb{R}^{r-1} \right\} \cong \left\{ P_{\Theta} \in \Delta_{r-1}^{\circ} : P_{\Theta}(i) = \frac{1}{Z(\Theta)} h_i \prod_{j \in [k]} \Theta_j^{a_{ji}} \right\}$$

• Look familiar?

Example 2 Γ a simplicial complex with set of facets $\text{facets}(\Gamma)$.

• $h = \mathbf{1} = (1, \dots, 1)^t \in \mathbb{R}^{|R|}$, where R is state space for (X_1, \dots, X_m) discrete and $E_m = G(\Gamma)$.

• Define the statistic $T: \mathcal{R} \rightarrow \mathbb{R}^{\times |R_F|}$ where $T: \underline{i} \mapsto \sum_{\underline{j} \in R_F, F \in \text{facets}(\Gamma)} \delta_{\underline{i} \neq \underline{j}}$

$$\Rightarrow a_{\underline{j}\underline{i}} = \begin{cases} 1 & \text{if } \underline{i} \text{ restricted to } F \text{ is } \underline{j} \\ 0 & \text{otherwise} \end{cases}$$

• Given $\eta \in \mathbb{R}^{\times |R_F|}$, set $\Theta_{\underline{j}}^{(F)} := \exp(\eta_{\underline{j}, F})$.

• Then $P_{\Theta}(\underline{i}) = \frac{1}{Z(\Theta)} \prod_{F \in \text{facets}(\Gamma)} \Theta_{\underline{i}_F}^{(F)}$

$\Rightarrow \mathcal{P}_{T,h}$ = the log-linear hierarchical model for Γ .

• In fact, all log-linear models are discrete regular exponential families...

Given $A \in [a_{ji}]_{j=1, i=1}^{k,r} \in \mathbb{Z}^{k \times r}$, define the statistic

$$T(i) := (a_{1i}, \dots, a_{ki})^t \text{ to be the columns of } A.$$

Take some $h \in \mathbb{R}_{>0}^r$.

For any $p_\eta \in \mathcal{P}_{T,h}$ we have that

$$p_\eta(i) = h_i e^{\eta^t T(i) - \phi(\eta)}$$

$$\log(p_\eta(i)) = \log(h_i) + \eta^t T(i) - \phi(\eta), \text{ or equivalently,}$$

$$\log(p_\theta(i)) = \log(h_i) + \log(\theta)^t A - \log(\zeta(\theta)).$$

If we assume $\mathbb{1} \in \text{rowspan}(A)$ then $b^t A = \mathbb{1}$ for some b and

$$\mathcal{P}_{T,h} \cong \left\{ p \in \Delta_{r-1}^\circ : \log(p) \in \log(h) + \text{rowspan}(A) \right\}.$$

Definition (17) Let $A \in \mathbb{Z}^{k \times r}$ be an integer matrix such that $\mathbb{1} \in \text{rowspan}(A)$ and let $h \in \mathbb{R}_{>0}^r$. The log-affine model associated to A and h is

$$M_{A,h} := \left\{ p \in \Delta_{r-1}^\circ : \log(p) \in \log(h) + \text{rowspan}(A) \right\}.$$

So we just proved the following:

Proposition (1) Discrete exponential families are the log-affine models.

This gives a geometric interpretation of discrete exponential families as arising from taking linear subspaces of the natural parameter space of the full regular exponential family for discrete random variables:

$$p_\eta : \log(p_\eta) = \log(h) + (\eta + \phi(\eta) \mathbf{b})^t A. \quad (\mathbf{b}^t A = \mathbb{1})$$

Letting $h = \mathbb{1}$, we recover the log-linear models: $M_A := M_{A, \mathbb{1}}$ from chapter 2.

We called log-linear models toric models. Hence, discrete exponential families are the toric models up to rescaling...

- For $A \in \mathbb{Z}^{k \times r}$ and $h \in \mathbb{R}_{>0}^r$ with $\mathbb{1} \in \text{rowspan}(A)$, define the monomial map $\phi^{A,h}: \mathbb{R}^k \rightarrow \mathbb{R}^r$ s.t. $\phi_i^{A,h}(\theta) = h_i \prod_{j \in [k]} \theta_j^{a_{ji}}$.

- The vanishing ideal

$$I_{A,h} := I(\phi^{A,h}(\mathbb{R}^k)) \subseteq \mathbb{R}[\mathbb{R}] := \mathbb{R}[p_1, \dots, p_r]$$

is the toric ideal of $M_{A,h}$.

- By definition, $V(I_{A,h})$ is the Zariski closure of $M_{A,h}$.

- mapping $p_j \mapsto p_j/h_j$ we can recover the generators of $I_{A,h}$ from $I_{A,\mathbb{1}}$. \Rightarrow discrete exponential families are toric models up to rescaling.

- We now prove a claim made in chapter 1 that the toric variety $V(I_L)$ for $L = \ker_{\mathbb{Z}} A$ is the Zariski closure of $M_{A,\mathbb{1}}$:

Proposition (2) | Let $A \in \mathbb{Z}^{k \times r}$. Then the toric ideal $I_{A,\mathbb{1}}$ is a binomial ideal and

$$I_{A,\mathbb{1}} = \langle p^u - p^v : u, v \in \mathbb{N}^r \text{ and } Au = Av \rangle$$

Proof: Any $p^u - p^v$ such that $Au = Av$ is in $I_{A,\mathbb{1}}$ since for any $\theta \in \phi^{A,\mathbb{1}}(\mathbb{R}^k)$ we get $p_i = \prod_{j \in [k]} \theta_j^{a_{ji}}$ and $p^u = \prod_{j \in [k]} \theta_j^{\langle u, (a_{j1}, \dots, a_{jr}) \rangle} = \prod_{j \in [k]} \theta_j^{\langle v, (a_{j1}, \dots, a_{jr}) \rangle} = p^v$. $\leftarrow Au = Av$.

- To show $p^u - p^v$ with $Au = Av$ generate $I_{A,\mathbb{1}}$, it suffices to show that they generate it as a vector space:

- If $f(p) \in I_{A,\mathbb{1}}$ and $c_u p^u$ a monomial in $f(p)$ with $c_u \neq 0$ then for any $\theta \in \mathbb{R}^k$ we have $f(\phi^{A,\mathbb{1}}(\theta)) = 0$.

$$\Rightarrow f(\phi^{A,\mathbb{1}}(\theta)) = 0 \text{ as a polynomial in } \theta.$$

- plugging $\phi^{A,\mathbb{1}}(\theta)$ into $c_u p^u$ we get $c_u \theta^{Au}$ which must cancel.

⇒ There is another monomial $c_\nu p^\nu$ such that plugging in $\phi^{A, \tau}(\theta)$ yields the same monomial $c_\nu \theta^{Au} \Rightarrow Au = Av$.

↪ subtract off $P: f - c_\nu(p^\nu - p^\nu)$, then induct. M

- Aside from being a beautiful correspondence between discrete exponential families and toric varieties, the geometric connection can help us solve important statistical problems...

Theorem 1) The maximum likelihood estimator of a discrete exponential family exists and is unique if and only if the observed value $T(x)$ of the sufficient statistic lies in the interior of the conic hull of $\text{supp}(V^T)$:

• $\text{cone}(V^T) = \left\{ \sum_{V^{(i)} \in \text{supp}(V^T)} \lambda_i V^{(i)} : \lambda_i \geq 0 \right\}$.

- The geometry of a toric variety V_A defined by $A \in \mathbb{Z}^{k \times r}$ is encoded in the polytope $P_A = \text{conv}(A)$.

$\text{cone}(P_A) = \text{cone}(A) = \text{cone}(\text{supp}(V^T))$.

Gaussian Regular Exponential Families

• We can similarly generalize the observation made last week to show that all multivariate Gaussian models are regular exponential families.

• Let $N_m := \{ N_m(\mu, \Sigma) : \mu \in \mathbb{R}^m, \Sigma \in \text{PD}_m \}$.

• $\mathcal{X} = \mathbb{R}^m$ with σ -algebra \mathcal{B}_m and Lebesgue measure ν .

• $T(x) := (x_1, \dots, x_m, -x_1^2/2, \dots, -x_m^2/2, -x_1 x_2, \dots, -x_{m-1} x_m) \in \mathbb{R}^m \times \mathbb{R}^{\binom{m+1}{2}}$.

↳ polynomial components linearly independent $\Rightarrow \dim(\text{supp}(V^T)) = m + \binom{m+1}{2}$.

• $h(x) \equiv 1$
• $\eta \in \mathbb{R}^m \times \mathbb{R}^{\binom{m+1}{2}}$, let η_{cons} be its first m components and

$\eta_{\text{cov}, m \times m} \in \mathcal{S}_{m \times m}$ the symmetric $m \times m$ matrix formed by the last $\binom{m+1}{2}$ components.

[Fact] $x \mapsto e^{\eta^t T(x)}$ integrable $\Leftrightarrow \eta_{[m \times m]}$ is PD.

$\Rightarrow N = \mathbb{R}^m \times PD_m$ is the natural parameter space.

For $\eta \in N$,

$$\phi(\eta) = -\frac{1}{2} (\log \det(\eta_{[m \times m]}) - \eta_{[m]}^t \eta_{[m \times m]} \eta_{[m]} - m \log(2\pi))$$

Hence, the densities P_η are of the form

$$P_\eta(x) = \frac{1}{\sqrt{(2\pi)^m \det(\eta_{[m \times m]}^{-1})}} \exp\left(\eta_{[m]}^t T(x) - \eta_{[m]}^t \eta_{[m \times m]} \eta_{[m]} / 2\right)$$

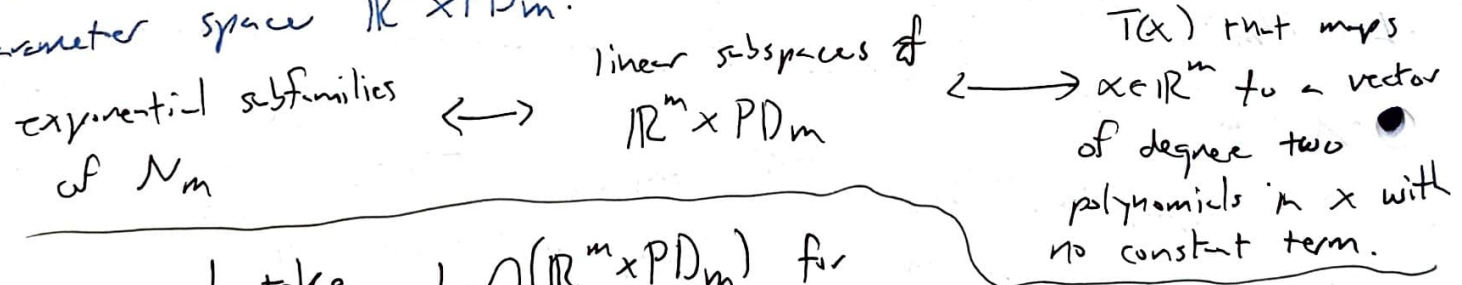
$$\eta^t T(x) = \eta^t (x_1, \dots, x_m, -x_1^2/2, \dots, -x_m^2/2, -x_1 x_2, \dots, -x_{m-1} x_m)^t$$

$$= \eta_{[m]}^t x + \left(-\frac{1}{2}\right) x^t \eta_{[m \times m]} x, \text{ where } x = (x_1, \dots, x_m)^t$$

Setting $\Sigma = \eta_{[m \times m]}^{-1}$ and $\mu = \Sigma^{-1} \eta_{[m]}$ we get

$$P_{\mu, \Sigma} = \frac{1}{\sqrt{(2\pi)^m \det(\Sigma)}} \exp\left(-\frac{1}{2} (x - \mu)^t \Sigma^{-1} (x - \mu)\right)$$

Just as in the discrete setting, the submodels of N_m that are also exponential families arise from taking linear subspaces of the natural parameter space $\mathbb{R}^m \times PD_m$.



In general take $L \cap (\mathbb{R}^m \times PD_m)$ for any linear subspace L of $\mathbb{R}^m \times \mathbb{S}^{m \times m}$.

In practice, $L = L_1 \times L_2$ for linear subspaces $L_1 \subseteq \mathbb{R}^m$ and $L_2 \subseteq \mathbb{S}^{m \times m}$.

Last time: $L_1 = \mathbb{R}^m$ and $L_2 = I_{m \times m}$.

Example (3) Let $L_1 := \mathbb{R}^2$ and L_2 be the linear subspace of

$S_{2 \times 2}$:

$$\left\{ \eta_{11} \begin{pmatrix} 1 & \eta_2 \\ \eta_2 & 0 \end{pmatrix} + \eta_{12} \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} + \eta_{22} \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix} : \eta_{11}, \eta_{22}, \eta_{12} \in \mathbb{R} \right\}.$$

Check that $(L_1 \times L_2) \cap (\mathbb{R}^2 \times PD_2)$ is the collection of parameters for an exponential family given by the statistic

$$T(x) = (x_1, x_2, (-x_1^2/2) + x_1 x_2, -x_2^2/2, -x_1 x_2).$$

• Since the models are linear in the inverse covariance matrix, the vanishing ideals of $N_m(0, \Sigma)$ also have interesting properties...

• If we set $L_1 = \underline{0}$ then the vanishing ideal of the model is

$$I(L_2^{-1}) \subseteq \mathbb{R} [\sigma_{ij} : 1 \leq i, j \leq m] \text{ where}$$

$$L_2^{-1} := \left\{ \Sigma = [\sigma_{ij}] \in PD_m : \Sigma^{-1} \in L_2 \cap PD_m \right\} \text{ is the}$$

inverse linear space of L_2 .

• Since the map $\Sigma \mapsto \Sigma^{-1}$ is rational, $I(L_2^{-1})$ will have nice properties... (examples coming!)

• Today, we didn't say "diffeomorphism." This idea only becomes important when we want to discuss algebraic exponential families: where we

use a nonalgebraic map to transform the natural parameter space and then talk about a semialgebraic subset of the transformed space to define our model:

Definition (18) Let $\{P_\eta : \eta \in N\}$ be a regular exponential family of order k . The subfamily induced by the set $M \subseteq N$ is an

algebraic exponential family if there exists an open subset $\bar{N} \subseteq \mathbb{R}^k$, a diffeomorphism $g: N \rightarrow \bar{N}$, and a semialgebraic set $A \subseteq \mathbb{R}^k$ such that $M = g^{-1}(A \cap \bar{N})$.