

I

Canonical form :

$$p(x; \theta) = \frac{1}{Z(\theta)} \exp(\theta^T T(x)) h(x) = \exp(\theta^T T(x) - \underbrace{A(\theta)}_{\log Z(\theta)}) h(x)$$

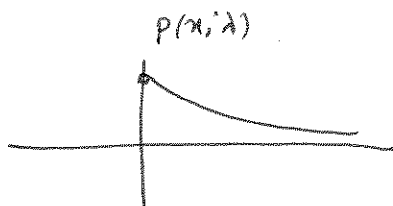
\downarrow Partition function \downarrow Sufficient Stat \swarrow base measure

if θ is function of something else, say η , then $p(x; \eta)$ is called "curved exponential family" ($\dim(\eta) < \dim(\theta(\eta))$)

low dim rep of θ

EX 1: Exponential Distribution

$$p(x; \lambda) = \begin{cases} \lambda e^{-\lambda x} & x > 0 \\ 0 & x \leq 0 \end{cases}$$



$$T(x) = x$$

$$\theta = -\lambda$$

$$p(x; \theta) = \exp\left(\frac{-\lambda x}{\theta} + \log \lambda\right) h(x)$$

$$h(x) = \mathbb{1}_{(0, \infty)}(x)$$

$$A(\theta) = \log\left(+\frac{1}{\lambda}\right)$$

EX 2:

$$\text{Cat}(x; \mu) = \prod_{k=1}^K \mu_k^{x_k} = \exp\left(\sum_{k=1}^K (\log \mu_k) x_k\right)$$

$$\mu \in \Delta^K$$

$$\mu = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_K \end{bmatrix}$$

$$= \exp\left(\sum_{k=1}^{K-1} (\log \mu_k) x_k + \left(1 - \sum_{k=1}^{K-1} x_k\right) \log \left(1 - \sum_{k=1}^{K-1} \mu_k\right)\right)$$

$$= \exp\left(\sum_{k=1}^{K-1} x_k \left[\log \mu_k - \log \left(1 - \sum_{j=1}^{K-1} \mu_j\right)\right] + \log \left(1 - \sum_{k=1}^{K-1} \mu_k\right)\right)$$

one-hot representation $x = \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$ it means $x=k$

II

$$= \exp \left[\sum_{k=1}^{K-1} x_k \log \frac{M_k}{M_K} + \log M_K \right]$$

$$\theta = \left[\log \frac{M_1}{M_K}, \dots, \log \frac{M_{K-1}}{M_K} \right]$$

$$T(x) = \begin{bmatrix} \mathbb{1}(x=1) \\ \vdots \\ \mathbb{1}(x=K-1) \end{bmatrix}$$

$$A(\theta) = -\log M_K$$

$$M_k = \frac{e^{\theta_k}}{1 + \sum_{j=1}^{K-1} e^{\theta_j}}$$

$$A(\theta) = \log \left(1 + \sum_{k=1}^{K-1} e^{\theta_k} \right)$$

softmax function

$$X_i \stackrel{i.i.d}{\sim} P(x; \theta)$$

joint
dist

$$P(x_1, \dots, x_n) = \left(\prod_{i=1}^n h(x_i) \right) \exp \left(\theta^T \left(\sum_{j=1}^n T(x_j) \right) - nA(\theta) \right)$$

joint distr is
also exp. family

III Dirichlet Distribution:

→ Multivariate Gaussian:

$$P(x; \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{k}{2}} \det(\Sigma)^{\frac{1}{2}}} \exp\left\{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)\right\}$$

$$= \frac{1}{(2\pi)^{\frac{k}{2}}} \exp\left[-\frac{1}{2} \overbrace{x^T \Sigma^{-1} x}^{\text{already inner product}} + \overbrace{x^T \Sigma^{-1} \mu}^{\text{already inner product}} - \frac{1}{2} \overbrace{\mu^T \Sigma^{-1} \mu}^{\text{already inner product}} - \log \det(\Sigma)\right]$$

also represented $\langle \Sigma^{-1}, xx^T \rangle$
(inner product)

$$T(x) = \begin{bmatrix} x \\ \text{vec}(xx^T) \end{bmatrix}$$

$$A(\theta) = \log \det \Sigma + \frac{1}{2} \mu^T \Sigma^{-1} \mu$$

$$\theta = \begin{bmatrix} \Sigma^{-1} \mu \\ \frac{1}{2} \text{vec}(\Sigma^{-1}) \end{bmatrix} \begin{matrix} \rightarrow \text{let's call } \theta_1 = \Sigma^{-1} \mu \\ \rightarrow \dots \theta_2 = -\frac{1}{2} \Sigma^{-1} \end{matrix}$$

$$h(x) = (2\pi)^{\frac{k}{2}}$$

then $A(\theta_1, \theta_2) = -\frac{1}{2} \text{tr}(\theta_2 \theta_1 \theta_1^T) - \frac{1}{2} \log(-2\theta_2)$

IV

Sufficient statistics:

Def: $t = T(X)$ is sufficient for underlying parameter θ precisely if conditional probability of data X does not depend on the parameter θ

$$P(x|t, \theta) = P(x|t)$$

→ You already have all info

Alternative way:

$$\cancel{P(\theta, x|t)} \quad P(\theta, x|t) = \underbrace{P(x|t, \theta)}_{\substack{\text{chain} \\ \text{rule}}} P(\theta|t) = P(x|t) P(\theta|t)$$

$\underbrace{P(x|t)}_{\substack{\text{Sufficient} \\ \text{stat}}}$

→ you can easily see

$$P(\theta|t, x) = P(\theta|t)$$

(V)

Partition function :

$$\int \exp(\theta^T T(x)) h(x) dx = e^{A(\theta)}$$

$$\nabla_{\theta} \int \exp(\theta^T T(x)) h(x) dx = \nabla_{\theta} A(\theta) e^{A(\theta)}$$

$$\int T(x) \exp(\theta^T T(x)) h(x) dx = \nabla_{\theta} A(\theta) e^{A(\theta)}$$

$$\int \underbrace{T(x) \exp(\theta^T T(x) - A(\theta)) h(x) dx}_{P(x|\theta)} = \nabla_{\theta} A(\theta)$$

$$\mathbb{E}[T(x)] = \nabla_{\theta} A(\theta)$$

(VI)

Posterior estimation

$$p(x|\theta) = h(x) \exp(\theta^T T(x) - A(\theta))$$

$T \in \mathbb{R}^d$



$$p(\theta; \tau, v_0) \propto \exp(v_0 \tau^T T'(\theta) - B(\tau, v_0) h'(\theta))$$

\uparrow
 a vector
 of param

\nwarrow
 scalar
 param

\downarrow
 a different
 sufficient stat
 $T' \in \mathbb{R}^{d'}$

~~There exists a conjugate prior~~

There exists a distribution s.t. $T'(\theta) = \begin{bmatrix} \theta \\ A(\theta) \end{bmatrix}$ ~~$h(\theta) \propto \exp(v_0 A(\theta))$~~
 $h'(\theta) = 1$ ($d' = d+1$)

$$p(\theta | \underbrace{\{x_1, \dots, x_n\}}_D) = p(\theta | D) \propto p(\theta, D)$$

$$= p(D|\theta) p(\theta)$$

$$\propto \left(\prod_{i=1}^n h(x_i) \right) \left(\exp \left(\theta^T \sum_{i=1}^n \overline{T}(x_i) - n A(\theta) \right) \times \exp(v_0 \theta^T \tau_0 - v_0 A(\theta)) \right)$$

$$\propto \exp \left(\theta^T (N \overline{T} + v_0 \tau_0) - A(\theta) (v_0 + n) \right)$$

$$= \exp \left((v_0 + n) \left(\frac{N \overline{T} + v_0 \tau_0}{v_0 + n} \right)^T \theta - A(\theta) (v_0 + n) \right)$$

$$= p \left(\theta; \frac{N \overline{T} + v_0 \tau_0}{v_0 + n}, v_0 + n \right)$$

Let's assume y is ~~the~~ scalar. Remember, the exponential family:

$$P(y; \theta) = \exp[\underbrace{\theta y - A(\theta)}_{T(y)z}] h(y)$$

we are assuming
 $T(y) = y$

I am going to generalize it:

$$P(y; \theta, \sigma^2) = \exp\left[\frac{\theta y - A(\theta)}{\sigma^2}\right] h(y, \sigma^2)$$

Remember:

$$(\text{mean}) \quad \mu \equiv \mathbb{E}[y] = A'(\theta)$$

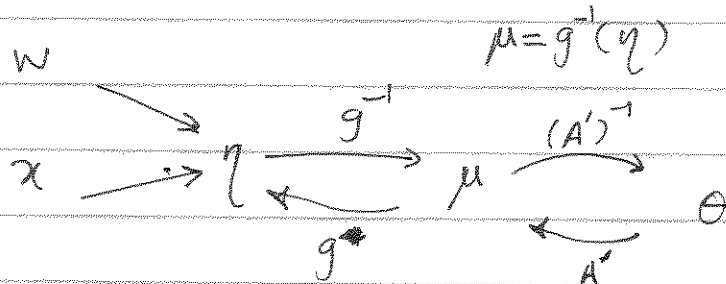
Let's assume σ is known.

We want to predict y using a linear function of features x . In other words $w^T x$. We can also assume any ~~linear~~ invertible ~~function~~ function of mean. Let's call

that function g^{-1} (mean function)

$$\mu = g^{-1}(w^T x)$$

g^{-1} is called mean function and $g(\cdot)$ is called link function. Let's define $\eta = W^T x$



Let's pick g to be $(A')^{-1}$

$$g \equiv (A')^{-1}$$

(Canonical Link function)

then $\mu_i = g^{-1}(\eta) \Rightarrow \theta = \eta$

Now, let's assume we have pairs of (x_i, y_i) .

The ML estimation:

$$\max_w \log p(D; w) = \sum_{i=1}^n \log p(y_i; x_i) = \frac{1}{\sigma^2} \sum_{i=1}^n \theta_i y_i - A(\theta_i)$$

note that

I have indexed

θ_i with i

b/c it's a function
of x_i

let's define

$$l_i \triangleq \theta_i y_i - A(\theta_i)$$

$$\frac{\partial l_i}{\partial w} = \frac{\partial l_i}{\partial \theta_i} \frac{\partial \theta_i}{\partial \mu_i} \frac{\partial \mu_i}{\partial \eta_i} \frac{\partial \eta_i}{\partial w}$$

\downarrow \downarrow
 $y_i - A'(\theta_i)$ x_i
 μ_i

$$\frac{\partial l_i}{\partial w} = (y_i - \mu_i) \frac{\partial \theta_i}{\partial \mu_i} \frac{\partial \mu_i}{\partial \eta_i} x_i$$

For canonical link function $\theta = \eta$ it become

$$\frac{\partial l_i}{\partial w} = (y_i - \mu_i) x_i$$

$$\nabla_w \mathcal{L} = \frac{1}{\sigma^2} \left[\sum_{i=1}^N (y_i - \mu_i) x_i \right]$$