

$$p(x; \theta) = \exp(\underbrace{\theta^T T(x)}_{\text{shift}} - \underbrace{A(\theta)}_{\text{normalizer}}) \underbrace{h(x)}_{\text{base meas.}}$$

Ex: Exponential Dist.
$$p(x; \lambda) = \begin{cases} \lambda e^{-\lambda x} & x > 0 \\ 0 & \text{other} \end{cases}$$

$$p(x; \theta) = \exp\left(\underbrace{-\lambda x}_{\theta} + \underbrace{\log \lambda}_{-A(\theta)}\right) \underbrace{\mathbb{1}(x > 0)}_{h(x)}$$

$$h(x) = \mathbb{1}(x > 0) \quad \theta = -\lambda$$

$$T(x) = x$$

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

Ex: Categorical Dist $x \in \{1, \dots, K\}$ $\mu \in \Delta^K$
$$p(x; \mu) = \prod_{k=1}^K \mu_k \mathbb{1}(x=k) \quad \mu_k = p(x=k)$$

$$= \prod_{k=1}^K \exp(\log \mu_k \mathbb{1}(x=k))$$

$$= \exp\left(\sum_{k=1}^K (\log \mu_k) \mathbb{1}(x=k)\right)$$

$$= \exp\left(\underbrace{(\log \mu)^T \mathbb{1}(x)}_{\sum_{j=1}^K \log \mu_j \mathbb{1}(x=j)} + \underbrace{\left(\sum_{k=1}^K \mathbb{1}(x=k)\right) \left(\frac{1}{\sum_{k=1}^K \log \mu_k}\right)}_{\log(1-\sum \mu_k)}\right)$$

$$= \exp\left(\sum_{j=1}^{K-1} \mathbb{1}(x=j) \log \frac{\mu_j}{\mu_K} + \log \frac{\mu_K}{1-\sum_{j=1}^{K-1} \mu_j}\right)$$

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

$$A(\theta) = \begin{bmatrix} \log \frac{\mu_1}{\mu_K} \\ \vdots \\ \log \frac{\mu_{K-1}}{\mu_K} \end{bmatrix}$$

$$\exp(\underbrace{\theta^T T(x)}_{\text{linear form a non-l. x}} - \underbrace{A(\theta)}_{\text{only } \theta}) \underbrace{h(x)}_{\text{only func of } x}$$

θ :
 $\theta(\eta)$ curved exp

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

$$= \frac{1}{\exp(A(\theta))} \exp(\theta^T T(x)) h(x)$$

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

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

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

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

$$P(\theta, x|t) = P(x|t, \theta) P(\theta|t) \quad \overset{x}{\circ} - \overset{t}{\circ} - \overset{\theta}{\circ}$$

$$= P(x|t) P(\theta|t) \quad \theta \nleftrightarrow x|t$$

$$P(\theta|t, x) = \frac{P(\theta, x|t)}{P(x|t)} = P(\theta|t)$$

$$\frac{T(x)}{\theta} \propto \frac{1}{A(\theta)} \rightarrow \eta$$

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

$$\theta = (\nabla_{\theta} A)^{-1}(\mu)$$

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

$$p(\theta; \tau_0, w) = \exp(\tau_0^T \tilde{T}(\theta) - B(\tau_0, w)) \tilde{h}(\theta)$$

$$T(\theta) = \begin{bmatrix} \theta \\ \tau_0 \cdot A(\theta) \end{bmatrix} \quad B(\cdot) = A(\theta) \quad \exp(\tau_0^T \tilde{T}(\theta) - \tau_0^T A(\theta))$$

$$p(\theta|D) \propto p(D, \theta) = \underbrace{P(D|\theta)}_{\prod_{i=1}^n p(x_i|\theta)} P(\theta) \quad \theta \in \mathbb{R}^d \quad \tilde{T}(\theta) \in \mathbb{R}^{d+1}$$

$$\propto \left(\prod_{i=1}^n h(x_i) \right) \frac{\exp(\theta^T \sum_{i=1}^n T(x_i) - nA(\theta))}{\exp(\tau_0^T \tilde{T}(\theta) - \tau_0^T A(\theta))} \quad \tau_0 = \begin{bmatrix} 0 \\ A(\theta) \end{bmatrix}$$

$$\propto \exp\left(\theta^T \left(\sum_{i=1}^n T(x_i) + \tau_0\right) - A(\theta)(n + \tau_0)\right)$$

$$\propto \exp\left(\theta^T \left(\frac{\sum_{i=1}^n T(x_i) + \tau_0}{n + \tau_0}\right) (n + \tau_0) - A(\theta)(n + \tau_0)\right)$$

$$\tilde{T}_0$$

$$y \quad p(y; \theta) = \exp[\theta y - A(\theta)] h(y)$$

$$\tilde{T}(y) = y$$

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

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

$$(x_i, y_i) \quad x_i \in \mathbb{R}^d \quad y_i \quad \begin{matrix} \circ^x \\ \downarrow \\ \circ^y \end{matrix}$$

$$w^T x_i \xrightarrow{g^{-1}} \mu \quad g(\mu) = w^T x_i$$

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

Param of linear funct \tilde{w}
$$\begin{matrix} \tilde{w} & \xrightarrow{g} & \mu & \xrightarrow{(A')^{-1}} & \theta \\ \uparrow & & \uparrow & & \uparrow \\ x_i & \xrightarrow{g} & \mu & \xrightarrow{A'} & \theta \end{matrix}$$

$$\tilde{w}^T x_i \quad \mu \quad \theta$$

$$A'(\theta) = \mathbb{E}[y]$$

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

$$\max_{\tilde{w}} \log P(D; w) = \sum_{i=1}^n \log p(y_i; w)$$

$$= \frac{1}{\sigma^2} \sum_{i=1}^n \theta_i y_i - A(\theta_i)$$

$$\text{index by } i \text{ how } x_i$$

$$(x_i, y_i) \quad \begin{matrix} \circ^x \\ \downarrow \\ \circ^y \end{matrix} \quad p(y|x)$$

$$\frac{\partial \ell_i}{\partial w} = \frac{\partial \ell_i}{\partial \theta_i} \frac{\partial \theta_i}{\partial w} \quad \frac{\partial \eta}{\partial w}$$

$$(y_i - A'(\theta_i)) \frac{\partial \theta_i}{\partial w} x_i$$

$$\nabla_{\tilde{w}} \mathcal{L} = \frac{1}{\sigma^2} \left[\sum_{i=1}^n (y_i - \mu_i) \left(\frac{\partial \theta_i}{\partial w} x_i \right) \right]$$