

PS 3

②

$$D_{KL}(P||Q) = \sum_{x \in X} P(x) \log \frac{P(x)}{Q(x)}$$

$$= - \sum_x P(x) \log \frac{Q(x)}{P(x)}$$

($\approx E_{\text{exp}}(\log Q/P)$)

$$= - \mathbb{E}_{x \sim P} \left[\log \frac{Q(x)}{P(x)} \right]$$

strict
log is concave
- log is convex

$$-D_{KL} = \mathbb{E}_{x \sim P} [.] \leq \log \mathbb{E}_{x \sim P} \left[\frac{Q(x)}{P(x)} \right] \\ = \log \sum_x P(x) \frac{Q(x)}{P(x)}$$

] ①

$$= \log 1 = 0$$

$$\therefore D_{KL} \geq 0$$

□

$$\text{if } P = Q \quad D_{KL} = \sum_x P(x) \log 1$$

$$D_{KL} = 0$$

$\because \log$ strict concave

$$D_{KL} = \mathbb{E}_{x \sim P} \left[\log \frac{P(x)}{Q(x)} \right] = \log \left(\frac{P(x)}{Q(x)} \right)$$

$\because E(x) = x$
as strict
convex

need strict convex

$$-\log$$

$$-D_{KL} = \sum_x P(x) \log \frac{Q(x)}{P(x)}$$

$$= \mathbb{E}_{x \sim P} \left[\log \frac{Q(x)}{P(x)} \right] = \log \left(\frac{\mathbb{E}_P(Q(x))}{P(x)} \right)$$

① if
antilog $E_P \left[\frac{Q(x)}{P(x)} \right] = \sum P \frac{Q}{P} = 1 \rightarrow \therefore \frac{Q(x)}{P(x)} = 1$ as $E(P) = P$

$$\log E_P [.] = \log 1 = 0$$

$$\therefore P(x) = Q(x)$$

$$\begin{aligned}
 \textcircled{b} \quad D_{KL}(P(x) || Q(x)) &= \sum_x P(x) \log \frac{P(x)}{Q(x)} \\
 + \quad D_{KL}(P(y|x) || Q(y|x)) &= \sum_x P(x) \left(\sum_y P(y|x) \log \frac{P(y|x)}{Q(y|x)} \right) \\
 &= \sum_x P(x) \left(\log \frac{P(x)}{Q(x)} + \sum_y P(y|x) \log \frac{P(y|x)}{Q(y|x)} \right) \\
 &\quad \cancel{\sum_y \sum_x P(y|x) P(x)} \\
 D_{KL}(P(x,y) || Q(x,y)) &= \sum_{x,y} P(x,y) \log \frac{P(x,y)}{Q(x,y)} \\
 &= \sum_{x,y} P(x,y) \log \frac{P(x) P(y|x)}{Q(x) Q(y|x)} \\
 &= \sum_{x,y} P(x,y) \log \frac{P(x)}{Q(x)} + \sum_{x,y} P(x,y) \log \frac{P(y|x)}{Q(y|x)} \\
 &= \sum_x \sum_y P(x|y) P(y) \log \frac{P(x)}{Q(x)} + \sum_{x,y} P(y|x) P(x) \log \frac{P(y|x)}{Q(y|x)} \\
 &\quad \sum_x \log \frac{P(x)}{Q(x)} \quad \sum_y P(x|y) P(y) \\
 &= \sum_x \left(\log \frac{P(x)}{Q(x)} \right) P(x) + \sum_{x,y} P(x) P(y|x) \log \frac{P(y|x)}{Q(y|x)}
 \end{aligned}$$

$$\textcircled{c} \quad \hat{P}(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{x^{(i)} = x\}} \quad \text{uniform dist over training set}$$

$$P_\theta(x) \equiv P(x; \theta) \quad \arg \min_{\theta} D_{KL}(\hat{P} || P_\theta) = \arg \max_{\theta} \sum_{i=1}^n \log P_\theta(x^{(i)})$$

$$\begin{aligned}
 D_{KL}(\hat{P} || P_\theta) &= \sum_x \hat{P}(x) \log \frac{\hat{P}(x)}{P(x; \theta)} \\
 &= \sum_x \hat{P}(x) \log \hat{P}(x) - \sum_x \hat{P} \log P(x; \theta) \\
 \arg \min_{\theta} - \sum_x \hat{P} \log P_\theta(x) &= \arg \max_{\theta} \sum \hat{P} \log P_\theta(x)
 \end{aligned}$$

$$\begin{aligned}
 &= \underset{\theta}{\operatorname{argmax}} \sum_{x_i} \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{x_i^{(i)} = x_i\} \log P_\theta(x) \\
 &\quad + \frac{1}{n} \sum_x \sum_{i=1}^n \mathbb{1}\{x_i^{(i)} = x_i\} \log P_\theta(x) \\
 \frac{1}{n} \sum_{i=1}^n \sum_x \mathbb{1}\{x_i^{(i)} = x_i\} P_\theta(x) &= \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^n \log P_\theta(x_i^{(i)}) \\
 \text{MLE} &= \min \text{KL from } \hat{P}
 \end{aligned}$$

Q2 ⑤ 29 bits to 16 colors

$$\begin{aligned}
 (8 \times 3) &\quad \therefore \text{need } \log_2(16) \text{ bits} = 4 \\
 \therefore \frac{29}{4} &= 6 \text{ compression}
 \end{aligned}$$

$$Q3 @ \text{Lsemisup}(\theta^{(t+1)}) = \sum_{i=1}^n \log \sum_z Q_i(z^{(i)}) \frac{P(x^{(i)}, z^{(i)}; \theta)^{(t+1)}}{Q_i(z^{(i)})} + \alpha \text{Lsup}(\theta^{(t+1)})$$

$$\begin{aligned}
 \text{Jensen} &\geq \sum_{i=1}^n \sum_z Q_i(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \theta^{(t+1)})}{Q_i(z^{(i)})} + \alpha \text{Lsup}(\theta^{(t+1)}) \\
 \text{M Step} &\geq \overbrace{\quad \quad \quad}^n \overbrace{\quad \quad \quad}^t; \theta^{(t+1)} \quad \alpha \text{Lsup}(\theta^{(t)})
 \end{aligned}$$

$$\because Q_i(z^{(i)}) = P(z^{(i)} | x^{(i)}; \theta) \quad \& \quad \sum Q_i(z^{(i)}) = 1$$

$$x \wedge z = x | z \cdot z = z | x \cdot x \quad \therefore \frac{P(x \wedge z)}{P(z | x)} = P(x)$$

$$\tilde{x} = \sum_{i=1}^n \log \frac{P(x_i, z^{(i)}; \theta)}{P(z | x_i; \theta)} + \alpha \text{Lsup}$$

$$= \sum_{i=1}^n \log P(x_i^{(i)}; \theta^{(t)}) + \alpha \text{Lsup}(\theta^{(t)})$$

$$= \text{Lsemisup}(\theta^{(t)})$$

⑥

latent is still $z^{(i)}$, i.e. $i \in \{1, \dots, n\}$
 \therefore E step like ~~at~~ unsupervised

$$\therefore w_j^{(i)} = Q_i(z^{(i)}=j) = \frac{p(x^{(i)} | z^{(i)}=j; \cdot) p(z^{(i)}=j)}{\sum_{k=1}^K p(x^{(i)} | z^{(i)}=k; \cdot) p(z^{(i)}=k)}$$

$$g_{jk} = \frac{1}{(2\pi)^{q/2}} |\sum_{j \neq k}|^{1/2}$$

$$\therefore = g_j \frac{\exp \left\{ -\frac{1}{2} (x^{(i)} - u_j)^T \sum_{j=1}^{-1} (x^{(i)} - u_j) \right\}}{\sum_{k=1}^K C_k \exp \left\{ -\frac{1}{2} (x^{(i)} - u_k)^T \sum_{k=1}^{-1} (x^{(i)} - u_k) \right\}}$$

i data pt.

$$ELBO_i^{(t)}(\theta) = \mathbb{E}_{\substack{x^{(i)} \\ \sim Q_i(z^{(i)})}} \log \left[\frac{P(x^{(i)} | z^{(i)}; \theta^{(t)}) P(z^{(i)}; \theta^{(t)})}{Q_i(z^{(i)}=j)} \right] \text{unsup}$$

$$+ \alpha \left[\log \frac{P(x^{(i)} | z^{(i)}; \theta^{(t)})}{P(z^{(i)}; \theta^{(t)})} \right] \text{sup}$$

$$\sum_{i=1}^n Q_i(z^{(i)}) \cdot \log [-] + \alpha [-] \underbrace{l_{\text{sup}}(\theta^{(t)})}_{l_{\text{sup}}(\theta^{(t)})}$$

$$\tilde{w}_j^{(i)} = \mathbb{1}\{z^{(i)}=j\} \quad \& \quad w_j^{(i)} = Q_i(z^{(i)}=j)$$

$$GMM: \theta^{(t+1)} = (\mu_j^{(t)}, \sigma_j^{(t)}, \phi_j^{(t)}) \quad Q_i(z^{(i)}) = P(z^{(i)} | x^{(i)}, \theta^{(t)})$$

$$\theta^{(t+1)} = \arg \max_{\theta} l_{\text{demi-sup}}(\theta^{(t)})$$

$$\nabla_{\theta} l_{\text{ss}} = \nabla_{\theta} \sum_{i=1}^n ELBO_i^{(t)}$$

$$= \nabla_{\theta} \left(\sum_{i=1}^n \mathbb{E}_{\substack{x^{(i)} \\ \sim Q_i(z^{(i)})}} \log [-] + \alpha \sum_{i=1}^n l_{\text{sup}}(\theta^{(t)}) \right)$$

$$\begin{cases} \nabla_{\theta} l_{\text{ss}} \\ \nabla_{\theta} M_j^{(t)} \end{cases} \text{ in order}$$

$$\sum_j \nabla_{\theta} M_j^{(t)}$$

$$\arg \min_{\theta} l_{\text{ss}} = \sum Q_i() \log(P(x|z)) + \left(\theta_{\text{true}} \frac{P(z)}{Q(z=j)} \right) + \alpha(x|z) + f(z)$$

param
var all
P=0

$\nabla_{\theta} l_{\text{ss}}$ w/ all param

$\nabla_{\theta} M_j^{(t)}$

$$= \sum_{i=1}^n \mathbb{E}_{\substack{z^{(i)} \sim Q(z^{(i)})}} \left[\nabla_{\mu_j} \log P(x^{(i)} | z^{(i)}; \theta^{(t)}) \right] + \alpha \sum_{i=1}^n \left[\nabla_{\mu_j} \log P(\tilde{x}^{(i)} | \tilde{z}^{(i)}; \theta^{(t)}) \right]$$

$\nabla_{\text{const}} = 0$ of gaussian

$$\therefore \sum_{i=1}^n \mathbb{E}_{\substack{z^{(i)} \sim Q(z^{(i)})}} \left[\nabla_{\mu_j} \left(-\frac{1}{2} (x^{(i)} - \mu_j)^T \Sigma^{-1} (x^{(i)} - \mu_j) \right) \right] + \dots (\tilde{x}^{(i)} - \mu_j)^T \Sigma^{-1} (\tilde{x}^{(i)})$$

$$\partial_i(z^{(i)}) \\ = w_j^{(i)}$$

$$\sum_{i=1}^n w_j^{(i)} \nabla^{-1}(x^{(i)} - \mu_j) + \sum_{i=1}^n \nabla^{-1}(x^{(i)} - \mu_j) = 0$$

$$\therefore \sum_{i=1}^n w_j^{(i)} \nabla^{-1} \mu_j^{(i)} + \alpha \sum_{i=1}^n \nabla^{-1} (\tilde{x}^{(i)})$$

$$\therefore \mu_j = \underbrace{\sum_{i=1}^n w_j^{(i)} x^{(i)}}_{\sum_{i=1}^n w_j^{(i)}} + \alpha \underbrace{\sum_{i=1}^n w_j^{(i)} \tilde{x}^{(i)}}_{\sum_{i=1}^n \tilde{w}_j^{(i)}}$$

Σ_j : using unsup GMM formula for Σ_j
 gives analogous Σ_j for unsup

$$\Sigma_j = \underbrace{\sum_{i=1}^n w_j^{(i)} (x^{(i)} - \mu_j)^T (x^{(i)} - \mu_j)}_{\sum_{i=1}^n w_j^{(i)}} + \alpha \underbrace{\sum_{i=1}^n \tilde{w}_j^{(i)} (\tilde{x}^{(i)} - \mu_j)^T (\tilde{x}^{(i)} - \mu_j)}_{\sum_{i=1}^n \tilde{w}_j^{(i)}}$$

(3) Lagrange $\sum_{j=1}^m \phi_j = 1$

analogy gives $\phi_j = \frac{\sum_{i=1}^n w_j^{(i)}}{n} + \alpha \sum_{i=1}^n \tilde{w}_j^{(i)}$

$$\tilde{w}_j^{(i)} = 1 \quad \{ \tilde{z}^{(i)} = j \}$$

we optimize $\log P(\tilde{z}^{(i)}; \theta^{(t)})$

$$\text{Completeness: } L(\phi, \lambda) = C + \sum_{i=1}^n \sum_{j=1}^k w_j^{(i)} \log \phi_j + \alpha \sum_{i=1}^n \sum_{j=1}^k \tilde{w}_j^{(i)} \times \log \phi_j + \lambda \left(\sum_{j=1}^k \phi_j - 1 \right)$$

(i) SS: fastn (2s iter) compared to US ~ 50

(ii) SS: Stable US: varies with initialization

(iii) US: error - failed single high variance gauss

SS: almost exact as & higher quality underlying

2018

$$\textcircled{1} @ \frac{\partial l}{\partial w_{ij}^{(1)}} = \frac{\partial l}{\partial o} \cdot \frac{\partial o}{\partial h_2} \cdot \frac{\partial h_2}{\partial w_{i,2}^{(1)}}$$

$\underbrace{z_0}_{z_0}$

$$\sigma(z) = \frac{1}{1+e^{-z}} \quad o^{(1)} = \sigma \left(\sum_{j=1}^3 w_j^{(2)} h_j^{(1)} + w_0^{(2)} \right)$$

$$h_j^{(1)} = \sigma \left(\sum_{k=1}^2 w_{kj}^{(1)} x_k^{(1)} + w_{0,j}^{(1)} \right)$$

$$\textcircled{1} \quad \frac{\partial l^{(1)}}{\partial o^{(1)}} = 2(o^{(1)} - y^{(1)})$$

(2) from $\sigma(z)$,

$$\frac{\partial o^{(1)}}{\partial z_0} = o^{(1)}(1-o^{(1)})$$

$$\frac{\partial o^{(1)}}{\partial h_2^{(1)}} = \frac{\partial o^{(1)}}{\partial z_0} \cdot \frac{\partial z_0}{\partial h_2^{(1)}} = o^{(1)}(1-o^{(1)})w_2^{(2)}$$

$$\textcircled{3} \quad \frac{\partial h_2^{(1)}}{\partial z_0} = h_2^{(1)}(1-h_2^{(1)}) \quad \text{sigmoid}$$

$$\therefore \frac{\partial h_2^{(1)}}{\partial w_{i,2}^{(1)}} = \frac{\partial h_2^{(1)}}{\partial z_0} \cdot \frac{\partial z_0}{\partial w_{i,2}^{(1)}} \\ = h_2^{(1)}(1-h_2^{(1)}) \cdot x_i^{(1)}$$

$$\text{Gewichtsupdate: } w_{i,2}^{(1)} := w_{i,2}^{(1)} - \alpha \sum_m \frac{\partial l^{(1)}}{\partial w_{i,2}^{(1)}}$$

$$2(o^{(1)} - y^{(1)}) \cdot o^{(1)}(1-o^{(1)})w_2^{(2)}h_2^{(1)}(1-h_2^{(1)})x_i^{(1)}$$

yes

- ⑥ Hidden layer neurons act as threshold & partitions the input space $y=0.5$ $x=0.5$ $2x+y=1$
- weight affects boundary
- 3 neurons can act as 3 sides of the L to partition
- ⑦ No: if linear, essentially output thresholds on lines combination giving boundary like perceptron. data is not linearly separable

0.5

Ex.

$$\textcircled{b} \quad \begin{aligned} W_{11}x_1 + W_{12}x_2 + b_1 &= z_1 \rightarrow 0 \cdot x_1 + x_2 = 0.5 = z_1 \\ W_{12}x_1 + W_{22}x_2 + b_2 &= z_2 \\ W_{13}x_1 + W_{23}x_2 + b_3 &= z_3 \end{aligned}$$

$$0.5 \quad 0.5 \quad 0$$

$$h_j = f(x_j)$$

$$0 = f(W_{h1}h_1 + W_{h2}h_2 + W_{h3}h_3 + b_0)$$

$$+ \quad 1 \quad 1 \quad -0.5$$

$$1 \quad -1 \quad 1 \quad + \quad 0 \quad 0.5$$

$$\approx 1 - 1 + 1 - 0.5$$

$$h_1 = 1 \quad h_2 = 1 \quad h_3 = -1 \quad b_0 = 0.5$$

Ex pt. 9, 9

$$z_2 = z_1 = 3.5 \quad h_1 = 1 = h_2 \quad z_3 = 9 \quad z \quad h_3 = 1$$

$$0 = f(1 + 1 + (-1)(1) - 0.5)$$

$$= f(0.5) \quad \text{assuming}$$

$$= 1$$

pt 1 1

$$z_1 = z_2 = 0.5 \quad h_1 = h_2 = 1 \quad z_3 = -2 \quad h_3 = 0$$

$$0 = f(1 + 1 + (-1)(0) - 0.5)$$

pt. 0.5, 0.5

$$0, 0, \quad 1, 1 \quad -4 \quad 0$$

$$(1 + 1 + 0 - 1.5) = f(0.5) \therefore 1 \quad \checkmark$$

pt 2, 2

$$z_1 = 1.5 \quad 1 \quad 1 \quad 0 \quad 1$$

$$1 + 1 + (-1)(1) - 0.5 = f(0.5) = 1 \quad x$$

\therefore bias bigger

pt. 0, 0

$$0, 0, 0$$

$$0 + 0 + 0 -$$

$$f(-\infty) = 0 \quad \checkmark$$

$$b_0 = -1.5$$

$$W_{h3} = -1$$

$$W_{h1} = W_{h2} = 1$$

$$\textcircled{1} \quad W_{11} = D \quad b_1 = -0.5$$

$$W_{21} = 1$$

$$\text{similar for } W_{12}, W_{22}$$

$$W_{13} = W_{23} = 1 \quad b_2 = -4$$

pt 1 1

$$z_1 = z_2 = 0.5 \quad h_1 = h_2 = 1 \quad z_3 = -2 \quad h_3 = 0$$

$$0 = f(1 + 1 + (-1)(0) - 0.5)$$

pt. 0.5, 0.5

$$0, 0, \quad 1, 1 \quad -4 \quad 0$$

$$(1 + 1 + 0 - 1.5) = f(0.5) \therefore 1 \quad \checkmark$$

pt 2, 2

$$z_1 = 1.5 \quad 1 \quad 1 \quad 0 \quad 1$$

$$1 + 1 + (-1)(1) - 0.5 = f(0.5) = 1 \quad x$$

\therefore bias bigger

pt. 0, 0

$$0, 0, 0$$

$$0 + 0 + 0 -$$

$$f(-\infty) = 0 \quad \checkmark$$

(2)

$$D_{KL}(p \parallel q) = \int_{-\infty}^{\infty} p(x) \log \frac{p(x)}{q(x)} dx$$

$$= \underset{y \sim p}{E} [\log p(y)] - \underset{y \sim p}{E} [\log q(y)]$$

(a) Show $\underset{y \sim p(y; \theta)}{E} [\nabla_{\theta} \log p(y; \theta')|_{\theta'=\theta}] = 0$

$$\text{Score of } p(y; \theta) = \nabla_{\theta} \log p(y; \theta)$$

$$= \frac{\nabla_{\theta} p(y; \theta)}{p(y; \theta)}$$

$$\begin{aligned} & \underset{y \sim p(y; \theta)}{E} \left[\frac{\nabla_{\theta} p(y; \theta)}{p(y; \theta)} \right] \\ &= \int_{-\infty}^{\infty} p(y; \theta) \cdot \frac{\nabla_{\theta} p(y; \theta)}{p(y; \theta)} dy \end{aligned}$$

$$= \int_{-\infty}^{\infty} \nabla_{\theta} p(y; \theta) dy$$

Applying Leibniz rule

$$= \nabla_{\theta} \int_{-\infty}^{\infty} p(y; \theta) dy = \nabla_{\theta} 1 = 0$$

(b) Fisher I(\theta) = $\underset{y \sim p(y; \theta)}{E} [\nabla_{\theta} \log p(y; \theta')|_{\theta'=\theta}]$ Cov of score

Show $\underset{y \sim p(y; \theta)}{E} [\nabla_{\theta} \log p(y; \theta) \nabla_{\theta} \log p(y; \theta)^T] = 0$

$$\text{Cov} = \frac{1}{n} \mathbf{x} \mathbf{x}^T \approx \mathbb{E}[\mathbf{x} \mathbf{x}^T] = \mathbb{E}[(\mathbf{x} - \mathbb{E}(\mathbf{x})) (\mathbf{x} - \mathbb{E}(\mathbf{x}))^T]$$

\therefore Using $I = \int_{-\infty}^{\infty} \nabla_{\theta} \log p(y; \theta) \cdot p(y; \theta) dy$

in discrete, $\text{Cov}(\nabla \log p) = \frac{1}{n} \log p \nabla \log p^T$

③ Show $\mathbb{E}_{y \sim p(y; \theta)} [-\nabla_{\theta}^2 \log p(y; \theta)] = I(\theta)$ [E[-Hessian]]

$$\nabla_{\theta} \log p = \frac{\nabla_{\theta} p}{p}$$

$$\nabla_{\theta}(\cdot) = \nabla_{\theta}\left(\frac{P}{P''}\right) = (\nabla_{\theta} p)\left(\frac{-\nabla_{\theta}^2 p}{P^2}\right) + \frac{p}{P^2} \nabla_{\theta}^2 p$$

$$-\nabla_{\theta}(\cdot) = \frac{(\nabla_{\theta} p)^2 - p \nabla_{\theta}^2 p}{P^2} - \cancel{\frac{\nabla_{\theta}^2 p}{P^2}} \cancel{(\nabla_{\theta} p)^2} + \frac{p}{P^2} \cancel{\nabla_{\theta}^2 p}$$

$$J(\theta) = \mathbb{E}\left[\frac{\nabla_{\theta} p}{p} \frac{\nabla_{\theta} p}{p}\right] = \int_{-\infty}^{\infty} \frac{(\nabla_{\theta} p)^2}{p^2} p dy$$

$$\mathbb{E}(-\nabla_{\theta}(\cdot)) = \int_{-\infty}^{\infty} p \cancel{(\nabla_{\theta} p)^2 dy} - \int_{-\infty}^{\infty} \nabla_{\theta}^2 p dy - \cancel{\int_{-\infty}^{\infty} \nabla_{\theta}^2 p dy}$$

$$\therefore \mathbb{E}(-\nabla_{\theta}(\cdot)) = J(\theta)$$

④ $\log p(y; \bar{\theta}) \approx \log p(y; \theta) + (\hat{\theta} - \theta)^T \nabla_{\theta} \log p(y; \theta) \Big|_{\theta=\theta}$
 Taylor around θ

$$\tilde{\theta} = \theta + d$$

$$+ \frac{1}{2} (\hat{\theta} - \theta)^T (\nabla_{\theta}^2 \log p(y; \theta') \Big|_{\theta'=\theta}) (\hat{\theta} - \theta)$$

$$\mathbb{E}_{y \sim p(y; \theta)} [\log p] = \mathbb{E}_{\theta} [\log p] + d^T \underbrace{\nabla_{\theta} \log p \Big|_{\theta=\theta}}_0 + \frac{1}{2} d^T \underbrace{\mathbb{E} [\nabla_{\theta}^2 \log p] \Big|_{\theta=\theta}}_0 + \frac{1}{2} d^T \mathbb{E} [\nabla_{\theta}^2 \log p] (d)$$

$$\therefore \mathbb{E} [\log p] = \mathbb{E} [\log p] + \frac{1}{2} d^T I(\theta) d \quad I(\theta)$$

$$D_{KL}(p_{\theta} || p_{\hat{\theta}}) = \mathbb{E}_{x \sim p(y; \theta)} [\log p(y; \theta)^2] - \mathbb{E}_{y \sim p(y; \hat{\theta})} [\log p(y; \hat{\theta})]$$

$$\approx \frac{1}{2} d^T I(\theta) d$$

$$\begin{aligned}
 \textcircled{C} \quad l(\theta + d) &\approx l(\theta) + d^T \nabla_{\theta} l(\theta') \Big|_{\theta'=\theta} \\
 &= \log p(y; \theta) + d^T \frac{\nabla_{\theta'} p(y; \theta') \Big|_{\theta'=\theta}}{p(y; \theta)} \\
 D_{KL}(p_0 \parallel p_{\theta}) &\approx \frac{1}{2} d^T I(\theta) d
 \end{aligned}$$

$$\begin{aligned}
 L(d, \lambda) &= l(\theta + d) - \lambda [D_{KL}(p_0 \parallel p_{\theta}) - c] \\
 &\approx \log p(y; \theta) + d^T \frac{\nabla_{\theta'} p(y; \theta') \Big|_{\theta'=\theta}}{p(y; \theta)} - \lambda \left[\frac{1}{2} d^T I(\theta) d - c \right] \\
 \nabla_d L(d, \lambda) &= 0 + \frac{\nabla_{\theta'} p(y; \theta') \Big|_{\theta'=\theta}}{p(y; \theta)} = \lambda I(\theta) d \stackrel{!}{=} 0
 \end{aligned}$$

$$d = \frac{\nabla_{\theta'} p(y; \theta') \Big|_{\theta'=\theta}}{p(y; \theta)} \cdot \frac{1}{\lambda I(\theta)} = \frac{1}{\lambda} (2)$$

$$\nabla_{\lambda} L(d, \lambda) = 0 = \frac{1}{2} d^T I(\theta) d - c$$

$$\begin{aligned}
 2c &= d^T I(\theta) d \\
 &= \frac{1}{\lambda} z^T I(\theta) z
 \end{aligned}$$

$$\lambda = \pm \sqrt{z^T I(\theta) z / 2c}$$

$$d^* = \frac{1}{\lambda} z = \frac{z}{\sqrt{z^T I(\theta) z}} \sqrt{2c} \quad \text{when } z > \frac{\nabla_{\theta'} p(y; \theta') \Big|_{\theta'=\theta}}{p(y; \theta)} \frac{1}{I(\theta)}$$

$$\begin{aligned}
 \textcircled{F} \quad \text{Newton} : \theta &:= \theta - H^{-1} \nabla_{\theta} l(\theta) \\
 \text{Natural grad} : I(\theta) &= \mathbb{E}_{y \sim p(y; \theta)} \left[-\nabla_{\theta'}^2 \log p(y; \theta') \Big|_{\theta'=\theta} \right]
 \end{aligned}$$

$$\theta := \theta + \tilde{d} = \theta + \frac{1}{\lambda} I(\theta)^{-1} \nabla_{\theta} l(\theta)$$

$$-\nabla_{\theta'}^2 l(\theta) = H$$