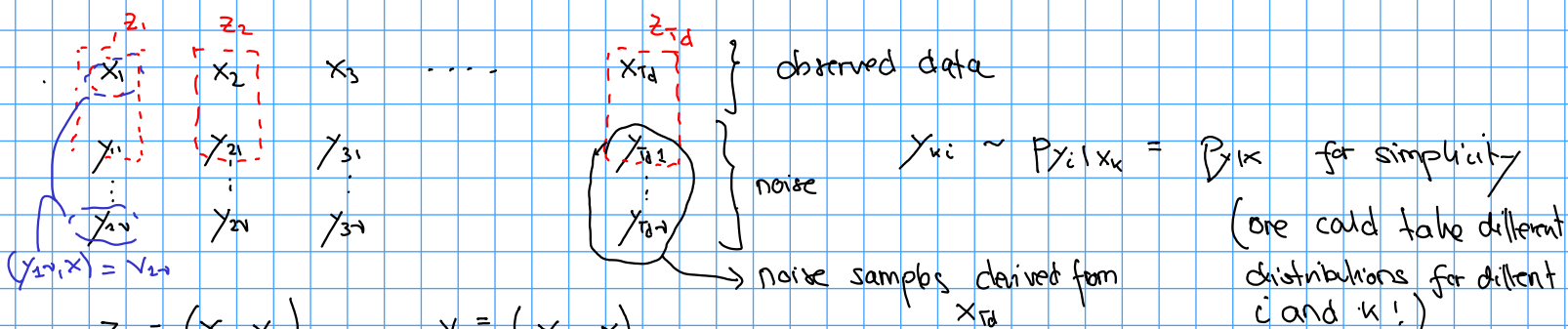


Better formulation of NCE with data-dependent noise distribution

Michael Gutmann June 28 2012

Given sample $(x_1 \dots x_n)$, generate \sim noise samples for each data point:



$$z = (x, y) \quad v = (y, x)$$

$$p_z(u) = p_{y|x}(u_2 | u_1) p_x(u_1) \quad ; \quad p_v(u) = p_{y|x}(u_1 | u_2) p_x(u_2)$$

$$u = (u_1, u_2)$$

Class one: observations $z_1 \dots z_{Td} = (x_1, y_1) \dots (x_{Td}, y_{Td})$; $z_k \sim p_z$

Class zero: observations $v_{1c} \dots v_{Td;c} = (y_{1c}, x_1) \dots (y_{Td;c}, x_{Td})$ $c = 2 \dots \sim$; $v_{kc} \sim p_v$

Class one has T_d members. Class zero has $(\sim - 1) \cdot T_d$ members.

$$P(C=1) = \frac{T_d}{(\sim - 1)T_d + T_d} = \frac{T_d}{\sim T_d} = \frac{1}{\sim}$$

$$P(C=0) = \frac{(\sim - 1)T_d}{\sim T_d} = 1 - \frac{1}{\sim} = \frac{\sim - 1}{\sim} \quad , \quad P(C=0) / P(C=1) = (\sim - 1)$$

$$p(u | C=1; \theta) = p_m(u_1; \theta) p_{y|x}(u_2 | u_1) \quad p(u | C=0; \theta) = p_m(u_2; \theta) p_{y|x}(u_1 | u_2)$$

$$\text{posteriors: } P(C=1 | u) = \frac{p(u | C=1) P(C=1)}{p(u | C=1) P(C=1) + p(u | C=0) P(C=0)}$$

$$= \frac{1}{1 + \frac{p_m(u_2; \theta) p_{y|x}(u_1 | u_2)}{p_m(u_1; \theta) p_{y|x}(u_2 | u_1)} \cdot \left(\frac{P(C=0)}{P(C=1)} \right)^{\sim - 1}}$$

$$P(C=0 | u) = 1 - P(C=1 | u)$$

$$\left(1 - \frac{1}{1+x} = \frac{1+x}{1+x} - \frac{1}{1+x} \right) = \frac{x}{1+x} = \frac{1}{1+1/x}$$

$$= \frac{1}{1 + \frac{p_m(u_1; \theta) p_{y|x}(u_2 | u_1)}{p_m(u_2; \theta) p_{y|x}(u_1 | u_2)} \cdot \frac{1}{\sim - 1}}$$

Logistic regression to distinguish between class 1 and class zero.

$$J_T(\theta) = \frac{1}{T_d} \sum_{t=1}^{T_d} \left[\ln P(c=1 | z_t) + \sum_{i=2}^n \ln P(c=0 | v_{it}) \right]$$

$$= \frac{1}{T_d} \sum_{t=1}^{T_d} \ln P(c=1 | z_t) + \frac{1}{T_d} \sum_{t=1}^{T_d} \sum_{i=2}^n \ln P(c=0 | v_{it})$$

$$\hat{\theta}_T = \operatorname{argmax} J_T(\theta)$$

Relation to NCE:

"Data": $(x_1, y_{11}) \dots (x_{T_d}, y_{1T_d})$

"Noise": $(y_{1i}, x_1) \dots (y_{T_d i}, x_{T_d}) \quad i=2, \dots, n$

} "Data" and "noise" are not independent. This means that the formula for the asymptotic error cannot be read out from the NCE theorems.

This relation to NCE suggests that $\hat{\theta}_T$ is consistent. On the next page, I consider the nonparametric case, which allows to prove consistency.

Shorthands:

$$f_m(\cdot; \theta) = \ln p_m(\cdot; \theta)$$

$$f_n(u; \theta) = \ln \left[p_{yx}(u_2 | u_1) / p_{yx}(u_2 | u_1) \right]$$

$$h(u; \theta) = 1 + \exp \left(f_m(u_2; \theta) - f_m(u_1; \theta) + f_n(u) \right) \cdot (n-1)$$

$$\check{h}(u; \theta) = 1 + \exp \left(f_m(u_1; \theta) - f_m(u_2; \theta) - f_n(u) \right) \frac{1}{n-1}$$

$$r(u; \theta) = \ln h(u; \theta)$$

$$\check{r}(u; \theta) = \ln \check{h}(u; \theta)$$

Case of large T_0 :

$$\begin{aligned} \bar{J}(\theta) &= \int \ln P(c=1|z) p_z(z) dz + \sum_{c=2}^n \int \ln P(c=0|v_c) p_v(v_c) dv_c \\ &= \int \ln P(c=1|z) p_z(z) dz + (n-1) \int \ln P(c=0|v) p_v(v) dv \end{aligned}$$

Reminder:

$$z = (x, y) \quad v = (y, x)$$

$$p_z(u) = p_{y|x}(u_2|u_1) p_x(u_1) \quad ; \quad p_v(u) = p_{y|x}(u_1|u_2) p_x(u_2)$$

$$\begin{aligned} \bar{J}(\theta) &= - \int \ln \left[1 + \exp \left(f_m(u_2; \theta) - f_m(u_1; \theta) + f_n(u) \right) \cdot (n-1) \right] p_z(u) du \\ &\quad - (n-1) \int \ln \left[1 + \exp \left(f_m(u_1; \theta) - f_m(u_2; \theta) - f_n(u) \right) \frac{1}{n-1} \right] p_v(u) du \end{aligned}$$

Computation of functional derivative

$$\ln(1 + a \exp(v)) = \ln(1 + a \exp(v_0)) + \frac{\frac{1}{1 + a \exp(v_0)}}{1 + a \exp(v_0)} (v - v_0) + \frac{1}{1 + a \exp(v_0)} \frac{1}{1 + a \exp(v_0)} \frac{1}{2} (v - v_0)^2$$

$$\begin{aligned} \frac{d}{dv} \left(\frac{1}{1 + \frac{1}{a} \exp(-v)} \right) &= - \left(\frac{1}{1 + \frac{1}{a} \exp(-v)} \right)^2 (-1) \frac{1}{a} \exp(-v) \\ &= \frac{1}{1 + \frac{1}{a} \exp(-v)} \frac{\frac{1}{a} \exp(-v)}{1 + \frac{1}{a} \exp(-v)} \\ &= \frac{1}{1 + \frac{1}{a} \exp(-v)} \frac{1}{1 + a \exp(v)} \end{aligned}$$

• First integral in \bar{J} : $v_0 = f_m(u_2) - f_m(u_1) + f_n(u)$

$$v - v_0 = \varepsilon \left(\phi(u_2) - \phi(u_1) \right)$$

$$a = n-1$$

$$1 + a \exp(v_0) = h(u)$$

\Rightarrow

$$1 + \frac{1}{a} \exp(-v_0) = \tilde{h}(u)$$

• Second integral in] : $w_0 = f_m(u_1) - f_m(u_2) - f_n(u)$

$$w - w_0 = \varepsilon (\phi(u_1) - \phi(u_2))$$

$$a = 1/(n-1)$$

$$\Rightarrow 1 + a \exp(w_0) = \check{h}(u)$$

$$1 + 1/a \exp(-w_0) = h(u)$$

• $J(f_m + \varepsilon \phi) = J(f_m) + \int \frac{1}{\check{h}(u)} \varepsilon (\phi(u_1) - \phi(u_2)) p_z(u) du$
 $- \int \frac{(n-1)}{h(u)} \varepsilon (\phi(u_1) - \phi(u_2)) p_v(u) du$

this terms are negative
 \Rightarrow negative curvature everywhere
 \Rightarrow critical point is a maximum.

$$- \int \frac{1}{\check{h}(u) h(u)} \frac{1}{2} \varepsilon^2 (\phi(u_1) - \phi(u_2))^2 p_z(u) du$$

$$- \int \frac{(n-1)}{\check{h}(u) h(u)} \frac{1}{2} \varepsilon^2 (\phi(u_1) - \phi(u_2))^2 p_v(u) du + O(\varepsilon^3)$$

First order term is zero if: $u_1 = u_2$ or $u_1 \neq u_2$ and

$$\frac{1}{\check{h}(u)} p_z(u) = \frac{(n-1)}{h(u)} p_v(u)$$

Reminder :

$$h(u; \theta) = 1 + \exp(f_m(u_2; \theta) - f_m(u_1; \theta) + f_n(u)) \cdot (n-1)$$

$$\check{h}(u; \theta) = 1 + \exp(f_m(u_1; \theta) - f_m(u_2; \theta) - f_n(u)) \cdot \frac{1}{n-1}$$

$$p_z(u) = p_{y|x}(u_2|u_1) p_x(u_1) ; \quad p_v(u) = p_{y|x}(u_1|u_2) p_x(u_2)$$

$$\Rightarrow p_z(u) = \frac{(n-1) \check{h}(u)}{h(u)} p_v(u)$$

(1)

$$\frac{p_z(u)}{p_v(u)} = \frac{(n-1) \exp(f_m(u_1) - f_m(u_2) - f_n(u))}{1 + (n-1) \exp(f_m(u_2) - f_m(u_1) + f_n(u))}$$

$$\frac{p_x(u_1) p_{y|x}(u_2|u_1)}{p_x(u_2) p_{y|x}(u_1|u_2)} = \frac{(n-1) \exp(f_m(u_1) - f_m(u_2)) \frac{p_{y|x}(u_2|u_1)}{p_{y|x}(u_1|u_2)}}{1 + (n-1) \exp(f_m(u_2) - f_m(u_1)) \frac{p_{y|x}(u_1|u_2)}{p_{y|x}(u_2|u_1)}}$$

since

$$f_n(u; \theta) = \ln \left[p_{Y|X}(u_1|u_2) / p_{Y|X}(u_2|u_1) \right]$$

Note: For the above equation to make sense, you need $p_{Y|X}(u_2|u_1) \neq 0$ and $p_{Y|X}(u_1|u_2) \neq 0$. Else we find the trivial relation $0=0$ in Eq. (1)

We obtain:

$$\begin{aligned} \frac{p_X(u_1)}{p_X(u_2)} &= \frac{(\gamma-1)p_{Y|X}(u_1|u_2) + \exp\left(\overbrace{f_m(u_1)} - \overbrace{f_m(u_2)}\right) p_{Y|X}(u_2|u_1)}{p_{Y|X}(u_2|u_1) + (\gamma-1)\exp(f_{m_2} - f_{m_1}) p_{Y|X}(u_1|u_2)} \\ &= \exp(f_{m_1} - f_{m_2}) \left[\frac{(\gamma-1)p_{Y|X}(u_1|u_2) \exp(f_{m_2} - f_{m_1}) + p_{Y|X}(u_2|u_1)}{p_{Y|X}(u_2|u_1) + (\gamma-1)\exp(f_{m_2} - f_{m_1}) p_{Y|X}(u_1|u_2)} \right] \\ &= \frac{\exp(f_m(u_1))}{\exp(f_m(u_2))} \cdot 1 \end{aligned}$$

$$\Rightarrow \frac{p_X(u_1)}{\exp(f_m(u_1))} = \frac{p_X(u_2)}{\exp(f_m(u_2))} \quad (2)$$

$$\Rightarrow \frac{p_X(u)}{\exp(f_m(u))} = \text{constant}$$

$$\Rightarrow \ln p_X = f_m + \text{constant}, \text{ or } f_m = \ln p_X + \text{constant}$$

Note Eq (2) means: $\ln p_X(u_1) - \ln p_X(u_2) = f_m(u_1) - f_m(u_2)$ $\forall u_1 \neq u_2$, which can be considered a stronger version of score matching.

This proves that $\tilde{f}(f_m)$ attains its maximum at $\exp(f_m) \propto p_X$.