

# Appendices of “Intentional Control of Type I Error over Unconscious Data Distortion: a Neyman-Pearson Approach to Text Classification”

## A. PROOFS

### A.1 Proof of Theorem 1

*Proof.* Recall that the (classical) oracle classifier regarding the pre-distortion population is  $h^*(x) = \mathbb{I}(\eta(x) > 1/2)$ , where the regression function  $\eta(x) = \mathbb{E}(Y|X = x)$  can be calculated as

$$\eta(x) = \frac{\pi_1 f_1(x)/f_0(x)}{\pi_1 f_1(x)/f_0(x) + \pi_0}.$$

Therefore,  $h^*(x) = \mathbb{I}\left(\frac{f_1(x)}{f_0(x)} > \frac{\pi_0}{\pi_1}\right)$ . When distortion with rates  $\beta_0$  and  $\beta_1$  is applied to class 0 and class 1 respectively, the class proportions become  $\pi_0^{(\beta_0, \beta_1)}$  and  $\pi_1^{(\beta_0, \beta_1)}$  which are defined as

$$\begin{aligned}\pi_0^{(\beta_0, \beta_1)} &= \frac{(1 - \beta_0)\pi_0}{(1 - \beta_0)\pi_0 + (1 - \beta_1)\pi_1}, \\ \pi_1^{(\beta_0, \beta_1)} &= \frac{(1 - \beta_1)\pi_1}{(1 - \beta_0)\pi_0 + (1 - \beta_1)\pi_1},\end{aligned}$$

while class conditional densities remain  $f_0$  and  $f_1$ . Then, the oracle classifier regarding the post-distortion population is to replace  $\pi_0$  and  $\pi_1$  in  $h^*$  by  $\pi_0^{(\beta_0, \beta_1)}$  and  $\pi_1^{(\beta_0, \beta_1)}$  respectively:

$$h_{(\beta_0, \beta_1)}^*(x) = \mathbb{I}\left(\frac{f_1(x)}{f_0(x)} > \frac{\pi_0^{(\beta_0, \beta_1)}}{\pi_1^{(\beta_0, \beta_1)}}\right) = \mathbb{I}\left(\frac{f_1(x)}{f_0(x)} > \frac{1 - \beta_0}{1 - \beta_1} \cdot \frac{\pi_0}{\pi_1}\right).$$

□

## A.2 Proof of Theorem 2

*Proof.* The constrained optimization program (4) in the main text that defines  $\phi_\alpha^*$  does not involve the class priors  $\pi_0 = \mathbb{P}(Y = 0)$  and  $\pi_1 = \mathbb{P}(Y = 1)$ , so  $\phi_\alpha^*$  does not depend on  $\pi_0$  or  $\pi_1$ . Now suppose distortion with rates  $\beta_0$  and  $\beta_1$  is imposed on class 0 and class 1 respectively, then the post-distortion population have class 0 proportion  $[(1 - \beta_0)\pi_0]/[(1 - \beta_0)\pi_0 + (1 - \beta_1)\pi_1]$  and class 1 proportion  $[(1 - \beta_1)\pi_1]/[(1 - \beta_0)\pi_0 + (1 - \beta_1)\pi_1]$ , while keeping the distributions of  $X|Y = 0$  and  $X|Y = 1$  unchanged. Since distortion at rates  $\beta_0$  and  $\beta_1$  only changes class proportion, which NP oracle does not depend upon, the NP oracle is invariant under distortion.  $\square$

## B. COST-SENSITIVE (CS) LEARNING

An insight from studying the classical classification paradigm is that the relative size of classification errors comes largely from the relative weights placed on type I and type II errors in the objective function. So a natural candidate to adjust classification errors is to change the weights. This is the so-called cost-sensitive (CS) learning paradigm, in which users impose costs  $C_0$  and  $C_1$  to type I and type II errors, respectively. On the population level, instead of minimizing the overall classification error  $R(\cdot)$ , one minimizes the CS learning objective:

$$\min_h R^c(h) := C_0\pi_0 R_0(h) + C_1\pi_1 R_1(h), \quad (\text{A.1})$$

or the following variant of (A.1):

$$\min_h R^{\bar{c}}(h) := C_0 R_0(h) + C_1 R_1(h). \quad (\text{A.2})$$

Then, the CS oracle classifier  $h^{c*}$  under the cost-sensitive learning paradigm (A.1) can be calculated by

$$h^{c*}(x) = \mathbb{I} \left( \frac{f_1(x)}{f_0(x)} > \frac{C_0}{C_1} \cdot \frac{\pi_0}{\pi_1} \right),$$

and the CS oracle  $h^{\bar{c}*}$  under (A.2) can be calculated by

$$h^{\bar{c}*}(x) = \mathbb{I} \left( \frac{f_1(x)}{f_0(x)} > \frac{C_0}{C_1} \right).$$

Similar to its counterpart in the classical paradigm, the post-distortion CS oracle classifier is different from the pre-distortion CS oracle, and the pre-distortion CS oracle cannot be recovered in view of an unknown distortion scheme. Lemma 1 follows from arguments similar to the proof of Theorem 1 in the main text.

**Lemma 1.** *Suppose that  $X|(Y = 0)$  and  $X|(Y = 1)$  have probability density functions  $f_0$  and  $f_1$ , and that class priors are  $\pi_0$  and  $\pi_1$  respectively. Let  $\beta_0$  and  $\beta_1$  be the distortion rates of class 0 and class 1 respectively. Then, the oracle classifier under the cost-sensitive learning paradigm (A.1) regarding the post-distortion population is*

$$h_{(\beta_0, \beta_1)}^{c*}(x) = \mathbb{I} \left( \frac{f_1(x)}{f_0(x)} > \frac{1 - \beta_0}{1 - \beta_1} \cdot \frac{C_0}{C_1} \cdot \frac{\pi_0}{\pi_1} \right).$$

Similarly, the oracle classifier under the paradigm (A.2) regarding the post-distortion population is

$$h_{(\beta_0, \beta_1)}^{\bar{c}*}(x) = \mathbb{I} \left( \frac{f_1(x)}{f_0(x)} > \frac{1 - \beta_0}{1 - \beta_1} \cdot \frac{C_0}{C_1} \right).$$

Lemma 1 implies that even if we have the entire post-distortion population, we can only mimic  $h_{(\beta_0, \beta_1)}^{c*}$  or  $h_{(\beta_0, \beta_1)}^{\bar{c}*}$ . However, unless  $\beta_0$  and  $\beta_1$  are known or estimable, there is no hope to mimic  $h^{c*}$  or  $h^{\bar{c}*}$ .

### C. ORACLE CLASSIFIERS WHEN WE RELAX THE FIXED CLASS CONDITIONAL DENSITIES ASSUMPTION

**Proposition 1.** *Suppose that pre-distortion,  $X|(Y = 0)$  and  $X|(Y = 1)$  have probability density functions  $f_0$  and  $f_1$ , and that class priors are  $\pi_0 = \mathbb{P}(Y = 0)$  and  $\pi_1 = \mathbb{P}(Y = 1)$ . Let  $\beta_0$  and  $\beta_1$  be the distortion rates of class 0 and class 1 respectively. Further suppose that the post-distortion class conditional densities of features are  $f'_0$  and  $f'_1$ . Then, the classical oracle classifier regarding the pre-distortion population is*

$$h^*(x) = \mathbb{I} \left( \frac{f_1(x)}{f_0(x)} > \frac{\pi_0}{\pi_1} \right),$$

and that regarding the post-distortion population is

$$h_{(\beta_0, \beta_1)}^{*'}(x) = \mathbb{I} \left( \frac{f'_1(x)}{f'_0(x)} > \frac{1 - \beta_0}{1 - \beta_1} \cdot \frac{\pi_0}{\pi_1} \right).$$

The proof is omitted due to its similarity to that for Theorem 1 in the main text. Note that when  $f'_1/f'_0 = f_1/f_0$ , that is when the ratio of class conditional densities of features is preserved under data distortion, the post-distortion classical oracle classifier  $h_{(\beta_0, \beta_1)}^{*'}(x)$  reduces to  $h_{(\beta_0, \beta_1)}^*(x)$  in Theorem 1, even if the class conditional densities themselves are changed. On the other hand, without assuming any relations between pre and post distortion feature distributions,  $f_1/f_0$  cannot be recovered.

The invariance property (Theorem 2 in the main text) of Neyman-Pearson (NP) oracle classifiers no longer holds in general when the class conditional densities of features are different pre and post distortion. The next proposition illustrates sufficient and necessary conditions under which this invariance property does hold for a fixed  $\alpha$ .

**Proposition 2.** *Denote pre-distortion distributions of  $X|(Y = 0)$  and  $X|(Y = 1)$  by  $f_0$  and  $f_1$  and those post-distortion by  $f'_0$  and  $f'_1$ . When  $f'_1/f'_0 = a \cdot (f_1/f_0)$  and*

$$a \cdot \min\{C \in \mathbb{R} : \mathbb{P}_{f_0}(f_1(X)/f_0(X) > C) \leq \alpha\} = \min\{C \in \mathbb{R} : \mathbb{P}_{f'_0}(f'_1(X)/f'_0(X) > C) \leq \alpha\},$$

*for some  $a > 0$ , the NP oracle classifier  $\phi_\alpha^*$  defined in (4) in the main text is invariant under distortion at various rates  $\beta_0$  (on class 0) and  $\beta_1$  (on class 1), regardless of whether pre-distortion classes are balanced. Moreover, these conditions are also necessary for the invariance property.*

*Proof.* From the NP Lemma, it is easy to see that the two conditions are sufficient for the invariance property of the NP oracles. For the necessary part, again by the NP lemma, the NP oracles pre and post distortion can be written respectively as

$$\phi_\alpha^*(x) = \mathbb{I}(f_1(x)/f_0(x) > C_\alpha), \text{ and } \phi_\alpha^{*'}(x) = \mathbb{I}(f'_1(x)/f'_0(x) > C'_\alpha),$$

for some constants  $C_\alpha$  and  $C'_\alpha$  as determined in the NP Lemma. In other words,

$$C_\alpha = \min\{C \in \mathbb{R} : \mathbb{P}_{f_0}(f_1(X)/f_0(X) > C) \leq \alpha\},$$

$$C'_\alpha = \min\{C \in \mathbb{R} : \mathbb{P}_{f'_0}(f'_1(X)/f'_0(X) > C) \leq \alpha\}.$$

Since  $C_\alpha$  and  $C'_\alpha$  are constants, to have  $\phi_\alpha^*(x) = \phi_\alpha^{*'}(x)$ , it is necessary to have  $f'_1/f'_0 = a \cdot (f_1/f_0)$  for some positive constants  $a$ , and this further demands  $C'_\alpha = a \cdot C_\alpha$ .  $\square$

Note that in general, the constant  $a$  in Proposition 2 depends on  $\alpha$ . In the following, we demonstrate that within certain distribution classes, the more general condition in Proposition 2 falls back to the special case of unchanged class conditional feature distributions, while in others, there are  $a \neq 1$  cases where class conditional feature distributions are different pre and post distortion.

**Case I: Exponential Distribution** Assume that  $f_0(x) = \lambda_0 e^{-\lambda_0 x}$ ,  $f_1(x) = \lambda_1 e^{-\lambda_1 x}$ ;  $f'_0(x) = \lambda'_0 e^{-\lambda'_0 x}$ ,  $f'_1(x) = \lambda'_1 e^{-\lambda'_1 x}$ , where  $x > 0$ . For identifiability concern, let us assume  $\lambda_0 < \lambda_1$ ,  $\lambda'_0 < \lambda'_1$ . Then,

$$\frac{f_1(x)}{f_0(x)} = \frac{\lambda_1}{\lambda_0} e^{-(\lambda_1 - \lambda_0)x},$$

and

$$\frac{f'_1(x)}{f'_0(x)} = \frac{\lambda'_1}{\lambda'_0} e^{-(\lambda'_1 - \lambda'_0)x}.$$

When we demand

$$\frac{f'_1(x)}{f'_0(x)} = a \cdot \frac{f_1(x)}{f_0(x)} \quad \forall x,$$

it follows that

$$\lambda_1 - \lambda_0 = \lambda'_1 - \lambda'_0, \tag{A.3}$$

and

$$\frac{\lambda'_1}{\lambda'_0} = a \cdot \frac{\lambda_1}{\lambda_0}. \tag{A.4}$$

Note that

$$\begin{aligned}
P_{f_0} \left( \frac{f_1(X)}{f_0(X)} > C \right) &= P_{f_0} \left( \frac{\lambda_1}{\lambda_0} e^{-(\lambda_1 - \lambda_0)X} > C \right) \\
&= P_{f_0} \left( e^{-(\lambda_1 - \lambda_0)X} > \frac{\lambda_0}{\lambda_1} C \right) \\
&= P_{f_0} \left( X < -\frac{1}{\lambda_1 - \lambda_0} \ln \left( \frac{\lambda_0}{\lambda_1} C \right) \right) \\
&= 1 - \exp \left\{ -\lambda_0 \cdot \left[ -\frac{1}{\lambda_1 - \lambda_0} \ln \left( \frac{\lambda_0}{\lambda_1} C \right) \right] \right\} \\
&= 1 - \left( \frac{\lambda_0}{\lambda_1} C \right)^{\frac{\lambda_0}{\lambda_1 - \lambda_0}}.
\end{aligned}$$

To choose the minimum  $C$  such that  $P_{f_0} \left( \frac{f_1(X)}{f_0(X)} > C \right) \leq \alpha$ , we get

$$C_\alpha = \frac{\lambda_1}{\lambda_0} (1 - \alpha)^{\frac{\lambda_1 - \lambda_0}{\lambda_0}}.$$

Similarly,

$$C'_\alpha = \frac{\lambda'_1}{\lambda'_0} (1 - \alpha)^{\frac{\lambda'_1 - \lambda'_0}{\lambda'_0}}.$$

Then the condition  $a \cdot C_\alpha = C'_\alpha$  implies that

$$a \cdot \frac{\lambda_1}{\lambda_0} (1 - \alpha)^{\frac{\lambda_1 - \lambda_0}{\lambda_0}} = \frac{\lambda'_1}{\lambda'_0} (1 - \alpha)^{\frac{\lambda'_1 - \lambda'_0}{\lambda'_0}}. \quad (\text{A.5})$$

For any given  $0 < \alpha < 1$ , combining three equations (A.3), (A.4) and (A.5) implies that

$$(1 - \alpha)^{\frac{1}{\lambda_0}} = (1 - \alpha)^{\frac{1}{\lambda'_0}},$$

which implies that  $\lambda_0 = \lambda'_0$ . And then,  $\lambda_1 = \lambda'_1$  and  $a = 1$ . Therefore, we have shown that when the class conditional feature distributions are restricted to the exponential distributions, the invariant property only occurs when  $f_0 = f'_0$  and  $f_1 = f'_1$ .

**Case II: Gaussian Distribution** Assume that  $f_0 : N(\mu_0, \sigma^2)$ ,  $f_1 : N(\mu_1, \sigma^2)$ ,  $f'_0 : N(\mu'_0, \sigma'^2)$ ,

and  $f'_1 : N(\mu'_1, \sigma'^2)$ , where  $\mu_0 < \mu_1$ ,  $\mu'_0 < \mu'_1$ , and  $\sigma \neq \sigma'$ . Then,

$$\frac{f_1(x)}{f_0(x)} = \exp \left\{ \frac{2(\mu_1 - \mu_0)x + \mu_0^2 - \mu_1^2}{2\sigma^2} \right\}$$

and

$$\frac{f'_1(x)}{f'_0(x)} = \exp \left\{ \frac{2(\mu'_1 - \mu'_0)x + \mu'^2_0 - \mu'^2_1}{2\sigma'^2} \right\}.$$

To obtain

$$\frac{f'_1(x)}{f'_0(x)} = a \cdot \frac{f_1(x)}{f_0(x)},$$

the parameters  $\mu_0, \mu_1, \sigma, \mu'_0, \mu'_1, \sigma', a$  must satisfy

$$\frac{2(\mu_1 - \mu_0)}{2\sigma^2} = \frac{2(\mu'_1 - \mu'_0)}{2\sigma'^2}, \quad (\text{A.6})$$

and

$$a = \exp \left\{ \frac{\mu'^2_0 - \mu'^2_1}{2\sigma'^2} - \frac{\mu^2_0 - \mu^2_1}{2\sigma^2} \right\}.$$

Furthermore, denote by  $\Phi(\cdot)$  the cumulative distribution function of standard normal distribution,

$$C_\alpha = \min_C \left\{ C \in R : P_{f_0} \left( \frac{f_1(X)}{f_0(X)} > C \right) \leq \alpha \right\} \text{ and } C'_\alpha = \min_C \left\{ C \in R : P_{f'_0} \left( \frac{f'_1(X)}{f'_0(X)} > C \right) \leq \alpha \right\}.$$

$$\begin{aligned} P_{f_0} \left( \frac{f_1(X)}{f_0(X)} > C \right) &= P_{f_0} \left( \exp \left\{ \frac{2(\mu_1 - \mu_0)X + \mu_0^2 - \mu_1^2}{2\sigma^2} \right\} > C \right) \\ &= P_{f_0} (2(\mu_1 - \mu_0)X + \mu_0^2 - \mu_1^2 > 2\sigma^2 \ln C) \\ &= P_{f_0} \left( X > \frac{2\sigma^2 \ln C + \mu_1^2 - \mu_0^2}{2(\mu_1 - \mu_0)} \right) \\ &= P_{f_0} \left( \frac{X - \mu_0}{\sigma} > \frac{2\sigma^2 \ln C + (\mu_1 - \mu_0)^2}{2(\mu_1 - \mu_0)\sigma} \right). \end{aligned}$$

Based on  $P_{f_0} \left( \frac{f_1(X)}{f_0(X)} > C \right) \leq \alpha$ , we get

$$\Phi^{-1}(1 - \alpha) \leq \frac{2\sigma^2 \ln C + (\mu_1 - \mu_0)^2}{2(\mu_1 - \mu_0)\sigma},$$

where  $\Phi^{-1}(\cdot)$  is the inverse function of  $\Phi(\cdot)$ , that is,

$$C \geq \exp \left\{ \frac{2\sigma(\mu_1 - \mu_0)\Phi^{-1}(1 - \alpha) - (\mu_1 - \mu_0)^2}{2\sigma^2} \right\}.$$

Therefore,

$$C_\alpha = \exp \left\{ \frac{2\sigma(\mu_1 - \mu_0)\Phi^{-1}(1 - \alpha) - (\mu_1 - \mu_0)^2}{2\sigma^2} \right\}.$$

Similarly,

$$C'_\alpha = \exp \left\{ \frac{2\sigma'(\mu'_1 - \mu'_0)\Phi^{-1}(1 - \alpha) - (\mu'_1 - \mu'_0)^2}{2\sigma'^2} \right\}.$$

From the relationship  $a \cdot C_\alpha = C'_\alpha$ , we can obtain

$$\begin{aligned} & \frac{\mu_0'^2 - \mu_1'^2}{2\sigma'^2} - \frac{\mu_0^2 - \mu_1^2}{2\sigma^2} + \frac{2\sigma(\mu_1 - \mu_0)\Phi^{-1}(1 - \alpha) - (\mu_1 - \mu_0)^2}{2\sigma^2} \\ &= \frac{2\sigma'(\mu'_1 - \mu'_0)\Phi^{-1}(1 - \alpha) - (\mu'_1 - \mu'_0)^2}{2\sigma'^2}, \end{aligned} \quad (\text{A.7})$$

i.e.,

$$\begin{aligned} & \frac{\mu_0'^2 - \mu_1'^2}{2\sigma'^2} + \frac{(\mu'_1 - \mu'_0)^2}{2\sigma'^2} - \frac{\mu_0^2 - \mu_1^2}{2\sigma^2} - \frac{(\mu_1 - \mu_0)^2}{2\sigma^2} \\ &= \frac{(\mu'_1 - \mu'_0)\Phi^{-1}(1 - \alpha)}{\sigma'} - \frac{(\mu_1 - \mu_0)\Phi^{-1}(1 - \alpha)}{\sigma}, \end{aligned}$$

which is equivalent to,

$$\frac{\mu'_0(\mu'_0 - \mu'_1)}{\sigma'^2} - \frac{\mu_0(\mu_0 - \mu_1)}{\sigma^2} = \left[ \frac{(\mu'_1 - \mu'_0)}{\sigma'} - \frac{(\mu_1 - \mu_0)}{\sigma} \right] \Phi^{-1}(1 - \alpha). \quad (\text{A.8})$$

From equation (A.6)

$$\frac{\mu'_1 - \mu'_0}{\sigma'} = \frac{\sigma'}{\sigma^2}(\mu_1 - \mu_0). \quad (\text{A.9})$$

Putting (A.9) into (A.8),

$$\frac{(\mu_0 - \mu_1)}{\sigma^2}(\mu'_0 - \mu_0) = \left[ \frac{\sigma'}{\sigma^2}(\mu_1 - \mu_0) - \frac{(\mu_1 - \mu_0)}{\sigma} \right] \Phi^{-1}(1 - \alpha),$$



that is,

$$\Phi^{-1}(1 - \alpha) = \frac{\mu_0 - \mu'_0}{\sigma' - \sigma}.$$

Putting the above arguments together, we have shown that under Gaussian distributions, for a given  $\alpha \in (0, 1)$ , the invariance property is satisfied precisely when

$$\frac{(\mu_1 - \mu_0)}{\sigma^2} = \frac{(\mu'_1 - \mu'_0)}{\sigma'^2},$$

$$\Phi^{-1}(1 - \alpha) = \frac{\mu_0 - \mu'_0}{\sigma' - \sigma},$$

and

$$a = \exp \left\{ \frac{\mu_0'^2 - \mu_1'^2}{2\sigma'^2} - \frac{\mu_0^2 - \mu_1^2}{2\sigma^2} \right\}.$$

**Example of Case II:** Let  $f_0 : N(0, 2^2)$ ,  $f_1 : N(1, 2^2)$  and  $f'_0 : N(-4, 4^2)$ , and  $f'_1 : N(0, 4^2)$ . We show that when  $\alpha = 0.023$ , the invariant property holds. First, it is easy to check that the above three equations hold with these density specifications and the choice of  $\alpha$ . In the following, we provide an alternative direct proof.

Note that

$$\frac{f_1(x)}{f_0(x)} = \exp \left\{ \frac{2x - 1}{8} \right\},$$

and

$$\frac{f'_1(x)}{f'_0(x)} = \exp \left\{ \frac{8x + 16}{32} \right\},$$

Hence,

$$\frac{f'_1(x)}{f'_0(x)} = \exp \left\{ \frac{5}{8} \right\} \cdot \frac{f_1(x)}{f_0(x)}.$$

We can take  $a = \exp \left\{ \frac{5}{8} \right\}$ . Let  $\alpha = 0.023$ . Then  $\Phi^{-1}(1 - \alpha) = 2$ . We solve for  $C_\alpha$  and  $C'_\alpha$  from

$$P_{f_0} \left( \frac{f_1(X)}{f_0(X)} > C_\alpha \right) = \alpha \quad \text{and} \quad P_{f'_0} \left( \frac{f'_1(X)}{f'_0(X)} > C'_\alpha \right) = \alpha.$$

That is,

$$P_{f_0} \left( \exp \left\{ \frac{2X-1}{8} \right\} > C_\alpha \right) = \alpha \quad \text{and} \quad P_{f'_0} \left( \exp \left\{ \frac{8X+16}{32} \right\} > C'_\alpha \right),$$

Or equivalently,

$$P_{f_0} \left( X > \frac{8 \ln C_\alpha + 1}{2} \right) = \alpha \quad \text{and} \quad P_{f'_0} (X > 4 \ln C'_\alpha - 2) = \alpha.$$

That is,

$$\frac{(8 \ln C_\alpha + 1)/2}{2} = 2 \quad \text{and} \quad \frac{4 \ln C'_\alpha - 2 - (-4)}{4} = 2,$$

which implies that

$$C_\alpha = \exp \left\{ \frac{7}{8} \right\} \quad \text{and} \quad C'_\alpha = \exp \left\{ \frac{3}{2} \right\}.$$

Obviously,

$$a \cdot C_\alpha = C'_\alpha,$$

i.e.,

$$a \cdot \min_C \left\{ C \in R : P_{f_0} \left( \frac{f_1}{f_0} > C \right) \leq \alpha \right\} = \min_C \left\{ C \in R : P_{f'_0} \left( \frac{f'_1}{f'_0} > C \right) \leq \alpha \right\}.$$

Therefore, we have constructed a concrete NP oracle invariant example in which  $f_0 \neq f'_0$  and  $f_1 \neq f'_1$ .

#### D. SPARSITY-INDUCING METHODS IN SELECTING MEANINGFUL TOPICS

Among the implemented methods, NP-sLDA performs the best in terms of power and it is a penalized sparsity-inducing method, which means it eliminates certain unimportant features as part of the classifier training process. In this section, we elaborate that such methods are effective in terms of selecting meaningful topics. In particular, we look at results from the first two random repetitions under **Setting 1** in Section 4.5.2 (random seed being set and results are readily available online) with  $K = 10$ . In the first repetition, Table A1 displays the selected ten topics and it's obvious that only topics 4 and 10 are the strike-related topics. Following the common practice of NP umbrella algorithms, we randomly split the training data  $M$  times for training the scoring function and thresholds. Here we use  $M = 7$ , and the final classifier is a majority vote. Figure A1

topic 1	罢工 strike 电话 phone	终于 finally 集体 collective	学校 school 分钟 minute	时间 time 失望 disappoint	一下 a bit 胃 stomach	事件 event 为了 for	彻底 complete 好多 many	哼哼 humph 疑问 question	开 open 多少 how many	对 right 忙 busy
topic 2	罢工 strike 回来 come back	今天 today 太阳 Sun	上班 work 话 words	发生 happen 公交车 bus	年 year 宿舍 dorm	上 go 冷 cold	发现 discover 块 block	问题 problem 过节 festival	种 type 东西 things	衰 decline 思考 think
topic 3	人 people 里 inside	让 let 妈妈 mom	说 speak 老师 teacher	吃 eat 今晚 tonight	时候 time 去 go	事 thing 很多 many	罢课 student strike 找 find	过 pass 出门 go out	哈哈 haha 最近 recent	小 small 班 class
topic 4	年 year 小时 hour	公司 company 后 after	员工 employee 广州 Guangzhou	中 within 抗议 protest	工人 worker 今日 today	工作 work 知道 know	工资 salary 请 please	最后 finally 月日 month-date	月 month 要求 request	鄙视 despise 中国 China
topic 5	抓 clutch 想 think	狂 crazy 潮湿 moist	罢工 strike 下午 afternoon	电脑 computer 结果 result	泪 tear 集 gather	生病 morning 继续 continue	现在 now 部 department	抓狂 go crazy 修 fix	天气 weather 人 people	回家 go home 委屈 be wronged
topic 6	罢工 strike 睡觉 sleep	去 go 鼻屎 mucus	做 do 挖 pick	次 times 第一 first	能 can 今晚 tonight	sick 怒 angry	地 ground 回 back	系 systems 真的 real	偷笑 smirk 叫 shout	没有 without 汗 sweat
topic 7	天 day 出来 out	手机 cellphone 换 exchange	还是 still 已经 already	知道 know 点 bit	能 can 郁闷 depressed	竟然 unexpectedly 鼓掌 applaud	突然 suddenly 听 listen	说 say 一下 a bit	玩 play 真是 really	这个 this 好不容易 hard
topic 8	罢工 strike 星期 week	想 think 然后 therefore	可怜 pity 休息 rest	居然 unexpectedly 家里 home	买 buy 半 half	发 give 悲伤 sad	明天 tomorrow 一直 always	累 tired 本来 originally	点 bit 听说 heard	但是 but 心情 mood
topic 9	罢工 strike 心情 mood	草草 hastily 点 a bit	明天 tomorrow 还有 also	可以 can 刚刚 just	开始 start 这个 this	好好 nicely 之后 after	真的 really 一定 must	新闻 news 为什么 why	爱 love 晚 evening	开 open 上午 morning
topic 10	的士 taxi 营运 operate	汕头 Shantou 三 three	出租车 taxi 原因 reason	司机 driver 政府 government	现在 now 集体 collective	车 car 今日 today	罢工 strike 希望 hope	打 call 四 four	下 get off 市民 citizen	辆 vehicle 月日 month-date

Table A1: top 20 keywords for the ten topics selected from repetition 1.

shows that, over the seven splits, NP-sLDA consistently selects only topics 4 and 10, and all the rest of the topics have corresponding coefficient 0. Similarly, in repetition 2, Table A2 shows that only topics 5 and 6 are the strike-related topics, and Figure A2 shows that NP-sLDA consistently selects topics 5 and 6 over the 7 splits. In summary, these sparsity-inducing methods, such as NP-sLDA, help select meaningful topics.

## E. PROOF AND GENERALIZATION OF PROPOSITION 1 IN MAIN TEXT

Proposition 1 in the main text follows as a special case of the next Proposition. Proposition 3 below explores the relationship between type I error  $R_0(\cdot)$ , the distortion rate  $\beta_0$  of class 0 and the class size ratio  $\pi_0/\pi_1$  for the classical post-distortion oracle classifier  $h_{\beta_0, \pi_0}^*$ .

(Intercept)	2.526526	(Intercept)	2.592229	(Intercept)	2.811912	(Intercept)	2.067274
x1	.	x1	.	x1	.	x1	.
x2	.	x2	.	x2	.	x2	.
x3	.	x3	.	x3	.	x3	.
x4	-5.914467	x4	-7.982268	x4	-7.691559	x4	-4.295535
x5	.	x5	.	x5	.	x5	.
x6	.	x6	.	x6	.	x6	.
x7	.	x7	.	x7	.	x7	.
x8	.	x8	.	x8	.	x8	.
x9	.	x9	.	x9	.	x9	.
x10	-19.521468	x10	-17.877815	x10	-20.608779	x10	-16.470045
(Intercept)	2.286663	(Intercept)	3.21628	(Intercept)	2.438294		
x1	.	x1	.	x1	.		
x2	.	x2	.	x2	.		
x3	.	x3	.	x3	.		
x4	-5.279157	x4	-11.65932	x4	-7.726201		
x5	.	x5	.	x5	.		
x6	.	x6	.	x6	.		
x7	.	x7	.	x7	.		
x8	.	x8	.	x8	.		
x9	.	x9	.	x9	.		
x10	-17.663786	x10	-20.64825	x10	-16.653843		

Figure A1: regression coefficients for the 7 splits in NP-sLDA, repetition 1.

topic 1	今天 today 前 forward	天 day 回家 go home	可以 can 还有 also	没有 without 吃饭 eat	开始 start 吃 eat	点 bit 号 day	真的 really 那些 those	日子 day 地铁 subway	明天 tomorrow 哈哈 haha	能 can 玩 play
topic 2	罢工 strike 真是 really	上 go to 三 three	上班 work 为了 for	能 can 生活 life	终于 finally 之后 after	抓狂 go crazy 超级 super	拿 get 只是 just	小时 hour 开心 happy	里 inside 觉得 feel	东西 thing 对 right
topic 3	罢工 strike 搞 do	去 go 过 over	系 be 怒 angry	做 do 公交 public transportation	今日 today 求 beg	地 ground 人 people	睡觉 sleep 甘 willing	后 after 吃 eat	起来 get up 街 street	听 listen 说 speak
topic 4	想 think 次 time	人 people 鄙视 despise	让 let 很多 many	说 speak 新 new	罢课 student strike 但是 but	累 tired 哦 oh	发 happen 感冒 a cold	衰 decline 虽然 although	生病 sick 委屈 be wronged	找 find 竟然 unexpectedly
topic 5	年 year 集体 collective	公司 company 第一 first	月 month 国际 international	员工 employee 次 time	工人 worker 要求 demand	工资 salary 劳动 labor	最后 finally 无法 unable	月日 month-date 机场 airport	工作 work 买 buy	还是 still 法国 France
topic 6	罢工 strike 下 get off	的士 taxi 广州 Guangzhou	汕头 Shantou 出门 go out	现在 now 已经 already	出租车 taxi 政府 government	司机 driver 事件 event	车 car 出 out	打 call 路 street	集体 collective 钱 money	辆 vehicle 问题 problem
topic 7	可怜 pity 生病 sick	小 small 竟然 unexpectedly	结果 result 郁闷 depressed	偷笑 smirk 开 open	发现 find 星期 week	昨天 yesterday 罢工 strike	今晚 tonight 三 three	早上 morning 今天 today	能 can 出来 go out	一直 always 结局 end
topic 8	罢工 strike 今天 today	天 day 时间 time	知道 know 电视 TV	天气 weather 突然 sudden	时候 time 奥特曼 Ultraman	挖 pick 好像 maybe	鼻屎 mucus 应该 should	太阳 Sun 全部 whole	种 type 水 water	周 week 点 bit
topic 9	罢工 strike 事 thing	抓 clutch 搞到 get	狂 crazy 部 department	泪 tear 学生 student	电脑 computer 结 form	居然 unexpectedly 啊啊 ah	今晚 tonight 手机 cellphone	鼓掌 applaud 明天 tomorrow	泪泪 tear 闹 alarm	学校 school 闹钟 alarm clock
topic 10	罢工 strike 女 female	手机 cellphone 怒骂 curse	中 within 分钟 minute	最近 recent 时候 time	哼哼 humph 过 over	一下 a bit 晚 late	哈哈 haha 深圳 Shenzhen	电梯 elevator 第一 first	停播 stop playing 迟到 late	玩 play 下班 off work

Table A2: top 20 keywords for the ten topics selected from repetition 2.

(Intercept)	2.890680	(Intercept)	2.541224	(Intercept)	2.588368	(Intercept)	2.523783
x1	.	x1	.	x1	.	x1	.
x2	.	x2	.	x2	.	x2	.
x3	.	x3	.	x3	.	x3	.
x4	.	x4	.	x4	.	x4	.
x5	-7.087097	x5	-7.949472	x5	-6.654200	x5	-6.538955
x6	-21.213795	x6	-17.069710	x6	-18.605443	x6	-18.279681
x7	.	x7	.	x7	.	x7	.
x8	.	x8	.	x8	.	x8	.
x9	.	x9	.	x9	.	x9	.
x10	.	x10	.	x10	.	x10	.
(Intercept)	2.682391	(Intercept)	2.450597	(Intercept)	2.917165		
x1	.	x1	.	x1	.		
x2	.	x2	.	x2	.		
x3	.	x3	.	x3	.		
x4	.	x4	.	x4	.		
x5	-4.826610	x5	-5.872407	x5	-9.101548		
x6	-21.347954	x6	-18.176870	x6	-19.616616		
x7	.	x7	.	x7	.		
x8	.	x8	.	x8	.		
x9	.	x9	.	x9	.		
x10	.	x10	.	x10	.		

Figure A2: regression coefficients for the 7 splits in NP-sLDA, repetition 2.

**Proposition 3.** Suppose probability densities of class 0 ( $X|Y = 0$ ) and class 1 ( $X|Y = 1$ ) follow distributions  $\mathcal{N}(\mu_0, \Sigma)$  and  $\mathcal{N}(\mu_1, \Sigma)$  respectively; class 0 composes  $\pi_0 \in (0, 1)$  proportion of the population and  $\beta_0 \in (0, 1)$  is the censorship rate of class 0 (i.e., the proportion of class 0 posts that were removed from some government censorship scheme). Suppose class 1 is not distorted (i.e.,  $\beta_1 = 0$ ). Let  $h_{\beta_0, \pi_0}^*$  be the classical oracle classifier in the post-distortion population. Then the type I error of  $h_{\beta_0, \pi_0}^*$  (regarding either the pre-distortion or the post-distortion population) is calculated as :

$$R_0(h_{\beta_0, \pi_0}^*) = \Phi \left( \frac{-\frac{1}{2}C - \log((1 - \beta_0)p)}{\sqrt{C}} \right), \quad (\text{A.10})$$

where  $C = (\mu_0 - \mu_1)^\top \Sigma^{-1} (\mu_0 - \mu_1)$  and  $p = \pi_0 / (1 - \pi_0)$ . Equation (A.10) implies that

1. Keeping  $\pi_0$  fixed (hence  $p$  is fixed),  $R_0(h_{\beta_0, \pi_0}^*)$  is a monotone increasing function of the class 0 censorship rate  $\beta_0 \in (0, 1)$ . Moreover, we have i). if  $pe^{3C/2} \leq 1$ ,  $R_0(h_{\beta_0, \pi_0}^*)$  is a concave function of  $\beta_0 \in (0, 1)$ ; and ii). if  $pe^{3C/2} > 1$ ,  $R_0(h_{\beta_0, \pi_0}^*)$  is a convex function of  $\beta_0$  for  $\beta_0 \in \left(0, 1 - \frac{1}{pe^{3C/2}}\right)$ , and a concave function for  $\beta_0 \in \left(1 - \frac{1}{pe^{3C/2}}, 1\right)$ .
2. Keeping  $\beta_0$  fixed,  $R_0(h_{\beta_0, \pi_0}^*)$  is a monotone decreasing function of the class ratio  $p = \pi_0 / (1 - \pi_0)$ . In other words, the larger the proportion of class 0 in the uncensored population, the smaller the type I error of  $h_{\beta_0, \pi_0}^*$ . Moreover,  $R_0(h_{\beta_0, \pi_0}^*)$  is a convex function of  $p$  for  $p > \frac{1}{(1 - \beta_0)e^{3C/2}}$ , and it is a concave function of  $p$  for  $p \leq \frac{1}{(1 - \beta_0)e^{3C/2}}$ .

*Proof.* Since equation (2) in the main text is the decision boundary of  $h_{\beta_0, \pi_0}^*$ , we have

$$R_0(h_{\beta_0, \pi_0}^*) = P_{X \sim \mathcal{N}(\mu_0, \Sigma)} \left\{ X^\top \Sigma^{-1} (\mu_0 - \mu_1) - \frac{1}{2} (\mu_0 - \mu_1)^\top \Sigma^{-1} (\mu_0 + \mu_1) + \log \left( \frac{(1 - \beta_0)\pi_0}{\pi_1} \right) \leq 0 \right\}.$$

For  $X$  in class 0,  $X^\top \Sigma^{-1} (\mu_0 - \mu_1) =: Z' \sim \mathcal{N}(\mu_0^\top \Sigma^{-1} (\mu_0 - \mu_1), (\mu_0 - \mu_1)^\top \Sigma^{-1} (\mu_0 - \mu_1))$ . Therefore,

$$\begin{aligned} R_0(h_{\beta_0, \pi_0}^*) &= P_{Z' \sim \mathcal{N}(\mu_0^\top \Sigma^{-1} (\mu_0 - \mu_1), (\mu_0 - \mu_1)^\top \Sigma^{-1} (\mu_0 - \mu_1))} \left\{ Z' \leq \frac{1}{2} (\mu_0 - \mu_1)^\top \Sigma^{-1} (\mu_0 + \mu_1) - \log \left( \frac{(1 - \beta_0)\pi_0}{\pi_1} \right) \right\} \\ &= \Phi \left( \frac{-\frac{1}{2} (\mu_0 - \mu_1)^\top \Sigma^{-1} (\mu_0 - \mu_1) - \log \left( \frac{(1 - \beta_0)\pi_0}{\pi_1} \right)}{\sqrt{(\mu_0 - \mu_1)^\top \Sigma^{-1} (\mu_0 - \mu_1)}} \right). \end{aligned}$$

Regarding part 1, for fixed  $\pi_0$ , let  $f(\beta_0) = R_0(h_{\beta_0, \pi_0}^*)$ .

$$f'(\beta_0) = \phi\left(\frac{-\frac{1}{2}C - \log((1 - \beta_0)p)}{\sqrt{C}}\right) \cdot \frac{1}{\sqrt{C}(1 - \beta_0)},$$

where  $\phi(\cdot)$  is the probability density function of the standard normal random variable. This implies that for  $\beta_0 \in (0, 1)$ ,  $f'(\cdot)$  is positive, so  $R_0(h_{\beta_0, \pi_0}^*)$  is a monotone increasing function of  $\beta_0$  for fixed  $\pi_0$ . Taking the second derivative of  $f$ , we have

$$f''(\beta_0) = \phi'\left(\frac{-\frac{1}{2}C - \log((1 - \beta_0)p)}{\sqrt{C}}\right) \cdot \frac{1}{C(1 - \beta_0)^2} + \phi\left(\frac{-\frac{1}{2}C - \log((1 - \beta_0)p)}{\sqrt{C}}\right) \cdot \frac{1}{\sqrt{C}(1 - \beta_0)^2}.$$

Let  $g(w) = \phi'(w) + \sqrt{C}\phi(w)$ . Then

$$g(w) = \frac{1}{\sqrt{2\pi}}e^{-\frac{w^2}{2}} \cdot (-w) + \frac{\sqrt{C}}{\sqrt{2\pi}}e^{-\frac{w^2}{2}}.$$

Note that  $g(w) > 0$  iff  $w < \sqrt{C}$ .

Therefore,  $f''(\beta_0) > 0$  iff  $g\left(\frac{-\frac{1}{2}C - \log((1 - \beta_0)p)}{\sqrt{C}}\right) > 0$  iff  $\frac{-\frac{1}{2}C - \log((1 - \beta_0)p)}{\sqrt{C}} < \sqrt{C}$  iff  $\beta_0 < 1 - \frac{1}{pe^{3C/2}}$ . Similarly  $f''(\beta_0) < 0$  iff  $\beta_0 > 1 - \frac{1}{pe^{3C/2}}$ .

Regarding part 2, for fixed  $\beta_0$ , let  $k(p) = R_0(h_{\beta_0, \pi_0}^*)$ , then

$$k'(p) = \phi\left(\frac{-\frac{1}{2}C - \log((1 - \beta_0)p)}{\sqrt{C}}\right) \cdot \frac{-1}{\sqrt{C}p}.$$

Clearly,  $k'(p) < 0$  for all  $p > 0$ .

$$k''(p) = \phi'\left(\frac{-\frac{1}{2}C - \log((1 - \beta_0)p)}{\sqrt{C}}\right) \cdot \frac{1}{Cp^2} + \phi\left(\frac{-\frac{1}{2}C - \log((1 - \beta_0)p)}{\sqrt{C}}\right) \cdot \frac{1}{\sqrt{C}p^2}.$$

Note that  $k''(p) > 0$  iff  $\frac{-\frac{1}{2}C - \log((1 - \beta_0)p)}{\sqrt{C}} < \sqrt{C}$  iff  $p > \frac{1}{(1 - \beta_0)e^{3C/2}}$ . □

The constant  $C$  can be considered as a measure of separability of the two classes. Note that when  $p = 1$ , that is when  $\pi_0 = 1 - \pi_0 = 1/2$ , if  $C$  is large (i.e., it is easy to separate the two classes),  $1/(pe^{3C/2}) \approx 0$ , then  $R_0(h_{\beta_0, \pi_0}^*)$  is a convex function of  $\beta_0 \in (0, 1)$ . On the other hand, when  $C$  is so small (i.e., two classes are hard to separate) that  $pe^{3C/2} \leq 1$ ,  $R_0(h_{\beta_0, \pi_0}^*)$  is a concave function of  $\beta_0 \in (0, 1)$ .



## F. NEYMAN-PEARSON LEMMA

The oracle classifier under the NP paradigm (NP oracle) arises from its close connection to the Neyman-Pearson Lemma in statistical hypothesis testing. Hypothesis testing bears strong resemblance to binary classification if we assume the following model. Let  $P_1$  and  $P_0$  be two *known* probability distributions on  $\mathcal{X} \subset \mathbb{R}^d$ . Assume that  $Y \sim \text{Bern}(\zeta)$  for some  $\zeta \in (0, 1)$ , and the conditional distribution of  $X$  given  $Y$  is  $P_Y$ . Given such a model, the goal of statistical hypothesis testing is to determine if we should reject the null hypothesis that  $X$  was generated from  $P_0$ . To this end, we construct a randomized test  $\phi : \mathcal{X} \rightarrow [0, 1]$  that rejects the null with probability  $\phi(X)$ . Two types of errors arise: type I error occurs when  $P_0$  is rejected yet  $X \sim P_0$ , and type II error occurs when  $P_0$  is not rejected yet  $X \sim P_1$ . The Neyman-Pearson paradigm in hypothesis testing amounts to choosing  $\phi$  that solves the following constrained optimization problem

$$\text{maximize } \mathbb{E}[\phi(X)|Y = 1], \text{ subject to } \mathbb{E}[\phi(X)|Y = 0] \leq \alpha,$$

where  $\alpha \in (0, 1)$  is the significance level of the test. A solution to this constrained optimization problem is called *a most powerful test* of level  $\alpha$ . The Neyman-Pearson Lemma gives mild sufficient conditions for the existence of such a test.

**Lemma 2** (Neyman-Pearson Lemma). *Let  $P_1$  and  $P_0$  be two probability measures with densities  $f_1$  and  $f_0$  respectively, and denote the density ratio as  $r(x) = f_1(x)/f_0(x)$ . For a given significance level  $\alpha$ , let  $C_\alpha$  be such that  $P_0\{r(X) > C_\alpha\} \leq \alpha$  and  $P_0\{r(X) \geq C_\alpha\} \geq \alpha$ . Then, the most powerful test of level  $\alpha$  is*

$$\phi_\alpha^*(X) = \begin{cases} 1 & \text{if } r(X) > C_\alpha, \\ 0 & \text{if } r(X) < C_\alpha, \\ \frac{\alpha - P_0\{r(X) > C_\alpha\}}{P_0\{r(X) = C_\alpha\}} & \text{if } r(X) = C_\alpha. \end{cases}$$

Under mild continuity assumption, we take the *NP oracle classifier*

$$\phi_\alpha^*(x) = \mathbb{I}\{f_1(x)/f_0(x) > C_\alpha\} = \mathbb{I}\{r(x) > C_\alpha\}, \quad (\text{A.11})$$

as our plug-in target for NP classification.