

---

# Incentive Incompatibility of Logistic Regression

---

## Abstract

We study the incentive compatibility of multi-class logistic regression. We provide a numerical example in which a strategic data provider has the incentive to misreport her private label to increase the classification probability of her true label. In particular, the model trained given her true label classifies her data point incorrectly, whereas the model trained given her misreported label classifies her data point correctly. We show that this incentive incompatibility disappears for classifiers that satisfy a monotonicity condition and independence of irrelevant alternatives condition. Examples of such classifiers include Bayes classifiers, kernel density estimators, and empirical risk minimization classifiers with zero-one loss.

## 1 Introduction

Consider an insurance company that makes its pricing decisions based on the customers' public observable characteristics, but the decision models are built using the information on private unobservable characteristics the customers' report. If the insurance company is transparent about its models, they might worry that customers have incentives to misreport their private information to get the contract that is the most beneficial to them. Similar examples include other rating systems that depend on the report of private information, such as loan applications, school grades, and employee screening.

In a general mechanism design problem, each of many strategic agents owns one public data point and reports her private label to the principal. The principal is the learner and builds a classifier based on the labels provided by the agents. Each agent chooses a label to report, not necessarily her true label, to maximize the probability that her data point is classified correctly by the principal. We say that a dataset is incentive-incompatible for the classifier if at least one of the agents has the incentive to misreport, and we characterize classifiers that are incentive-compatible with all possible datasets.

We start with an example dataset that is incentive-incompatible for the multi-class logistic regression classifier. In the dataset, each of the 18 agents owns a two-dimensional data point and a private label with one of three values: "red", "green", or "blue".



Figure 1: Incentive-incompatible Example

The 18 points are located inside a unit circle, and each point is 0.004 away from the three line segments through the origin that forms angles of 120 degrees between them. There is one red point, drawn as a square in the diagram, that is on the "incorrect" side of the boundary. For the agent represented by the red square, truthfully reporting her label will lead to a multi-class logistic regression model that classifies her point as "green". The probability that this model classifies her point as "red" is 0.3290. However, if the agent misreports her label as "blue", the resulting model classifies her point as "red", and with a probability of 0.4966. By lying about her label, the agent can make the principal learn an incorrect model that classifies her point correctly and with a higher probability.

The example provides insight into the general incentive incompatibility issue of many classification models used in machine learning. The "red" agent that gets incorrectly classified as "green", does not want to misreport her label as "green" but instead has the incentive to misreport her label as the third alternative "blue" to influence the classifier in a way that changes the decision between "red" and "green". Intuitively, if there are only two classes, then the agent would not be willing to misreport her label if the classifier is monotonic, in the sense that adding a point from one class increases the probability that this point is classified as a member of that class; and if there are more than two classes, then the agent would not be willing to misreport her label if the classifier is independent of irrelevant alternatives, in the sense that adding a point from a third class would not affect the decisions between the two classes.

Previous work on mechanism design for machine learning with strategic data sources focuses on designing robust algorithms to incentivize the data providers to report their private data truthfully. Their models mainly differ in the objective and the possible actions of the data providers (agents) and the learner (principal).

The first group of papers focuses on principal-agent problems similar to our paper in which each agent's private data point is the agent's type that the agent cannot change. The only action the agents can take is whether to report their private information truthfully.

1. Some models assume the agents' data points (or feature vectors) are public, but their labels are private. Perote and Perote-Pena (2004), Chen, Podimata, Procaccia, and Shah (2018), and Gast, Ioannidis, Loiseau, and Roussillon (2013) focus on strategy-proof linear regression algorithms and introduced clockwise repeated median estimators, generalized resistant hyperplane estimators, and modified generalized linear squares estimators. Dekel, Fischer, and Procaccia (2010) investigates the general regression problem with empirical risk minimization and absolute value loss. All the previously mentioned papers assume the labels are continuous variables (regression problems), and Meir, Procaccia, and Rosenschein (2012) assumes the labels are discrete variables (classification problems) and proposes a class of random dictator mechanisms.

2. Some models assume the agents' data points are also private. Chen, Liu, and Podimata (2019) investigates such problems for linear regressions.
3. Other models do not involve labels. Each agent has a private valuation. These problems are usually modeled as facility location problems and the solution involves some variant of the Vickrey-Clarke-Groves or Meyerson auction. They include Dütting, Feng, Narasimhan, Parkes, and Ravindranath (2017), Golowich, Narasimhan, and Parkes (2018), Epasto, Mahdian, Mirrokni, and Zuo (2018), and Procaccia and Tennenholtz (2009).

The second group papers focus on moral-hazard problems in which each agent does not have a type but they can choose an action (with a cost) that affects the probability of obtaining the correct label. Richardson, Rokvic, Filos-Ratsikas, and Faltings (2019) focuses on the linear regression problem in this scenario, and Cai, Daskalakis, and Papadimitriou (2015) and Shah and Zhou (2016) investigates the problem for more general machine learning problems. Mihailescu and Teo (2010) also discusses a similar problem for general machine learning algorithms.

The last group of papers uses machine learning or robust statistics techniques without game-theoretic models. This group of papers include Dekel and Shamir (2009b), Dekel and Shamir (2009a).

## 2 Logistic Regression

### 2.1 Model

In this section, we introduce the model using logistic regression as an example. We assume the principal is training a multi-class logistic (softmax) regression. There are  $n$  strategic agents each providing the label of one data point to the principal. An agent,  $i$ , with public  $x_i \in \mathbb{R}^m$ , and private discrete  $y_i \in \{1, 2, \dots, k\}$ , has the objective of maximizing the probability that her data point is labeled correctly by the principal's classifier, parameterized by the  $m \times (k + 1)$  weights (and bias) matrix  $w$ . The agent can choose to report  $y_i^\dagger$  to achieve the objective, with possibly  $y_i^\dagger \neq y_i$ . Denoting the weights of the model resulting from the false report from agent  $i$  by  $w^\star(y_i^\dagger)$  (and  $w^\dagger$  when the identity and the report of the agent is unambiguous), the agent's objective can be written as,

$$\max_{y^\dagger \in \{1, 2, \dots, k\}} \mathbb{P} \left\{ Y = y_i | x_i; w^\star(y_i^\dagger) \right\},$$

where,

$$\mathbb{P} \{ Y = c | x_i; w \} = \frac{e^{z_{i,c}}}{\sum_{c'=1}^k e^{z_{i,c'}}},$$

$$z_{i,c} = \sum_{j=1}^m w_{j,c} x_{i,j} + b_c, \text{ for } c \in \{1, 2, \dots, k\}.$$

The principal is not strategic and he maximizes the likelihood of the data,

$$\max_w \sum_{i=1}^n \log \left( \mathbb{P} \{ Y = y_i^\dagger | x_i; w \} \right).$$

We consider the case without a coalition of a group of agents, so only one agent is misreporting at a time, and use the following notations,

$$w^\star = \arg \max_w \sum_{i=1}^n \log (\mathbb{P} \{ Y = y_i | x_i; w \}), \text{ and}$$

$$w^\dagger = w^\star(y_i^\dagger) = \arg \max_w \log \left( \mathbb{P} \{ Y = y_i^\dagger | x_i; w \} \right) + \sum_{i'=0, i' \neq i}^n \log (\mathbb{P} \{ Y = y_{i'} | x_{i'}; w \}).$$

**Definition 1.** A dataset is incentive-incompatible for a classifier if there exists at least one agent  $i$ , and some  $y_i^\dagger \neq y_i$  such that,

$$\mathbb{P}\{Y = y_i | x_i; w^*\} < \mathbb{P}\{Y = y_i | x_i; w^*(y_i^\dagger)\}.$$

A classifier is incentive-compatible if there does not exist a dataset that is incentive-incompatible for the classifier.

**Conjecture 1.** *Multi-class logistic regression is not incentive-compatible.*

The example described in Figure 1 is a dataset that is numerically incentive-incompatible. In this example, agent  $i$  reports  $x_i \in \mathbb{R}^2$  and  $y_i$  is one of 1 (red), 2 (green), or 3 (blue). Suppose the red square point corresponds to agent 1 with  $x_1 = (-1.63, -1.17)$  and  $y_1 = 1$ , then

$$\begin{aligned}\mathbb{P}\{Y = 1 | x_1; w^*\} &= 0.3290, \\ \mathbb{P}\{Y = 1 | x_1; w^*(y_1^\dagger = 3)\} &= 0.4966.\end{aligned}$$

Here, parameter estimation is done using maximum likelihood estimation with BFGS, and  $w^*$  is given by, with class 1 weights normalized to 0,

Class	(Intercept)	x1	x2
1	0	0	0
2	-0.6053178	104.9925	-181.3391914
3	-0.2852057	209.4190	0.3656777

and  $w^*(y_1^\dagger = 3)$  is given by,

Class	(Intercept)	x1	x2
1	0	0	0
2	-0.1915645	3.473426	-5.507418
3	0.8273350	4.309293	-1.200060

Currently, there is no formal proof that the result is not due to numerical instability, therefore, @ (logit) is stated as a conjecture. Numerical experiments indicate that incentive-incompatible datasets are rare. If the data points are two-dimensional and standard normally distributed, and the labels are created using randomly generated weights with a small probability of error, then such a dataset is incentive-incompatible with a probability of 0.005.

## 2.2 Loss Functions

It is, however, possible to change the loss function so that logistic regression is incentive-compatible. Changing the loss function to absolute value  $L^1$  loss is one possibility, due to Dekel, Fischer, and Procaccia (2010). Their result on incentive compatibility of empirical risk minimization in the regression setting is applicable in our model. In addition to absolute value loss, empirical risk minimizers with zero-one loss are also always incentive-compatible.

**Proposition 1.** *Multi-class deterministic empirical risk minimization classifiers with zero-one loss are incentive-compatible.*

*Proof.* For any dataset  $\{(x_i, y_i)\}_{i=1}^n$ , and the hypothesis class  $\mathcal{H}$ , let the optimal classifier be,

$$h^* = \arg \min_{h \in \mathcal{H}} \sum_{i'=1}^n \mathbb{1}_{\{y_{i'} \neq h(x_{i'})\}}.$$

Fix an agent  $i$  with  $x_i$ , and fix the other agents' reports,  $(x_{-i}, y_{-i})$ , define the loss function given the classifier  $h$  and report of agent  $i$ ,  $y_i^\dagger$ , as,

$$\ell(y_i^\dagger; h) = \sum_{i' \neq i} \mathbb{1}_{\{y_{i'} \neq h(x_{i'})\}} + \mathbb{1}_{\{y_i^\dagger \neq h(x_i)\}}.$$

If  $y_i = h^*(x_i)$ , then the classifier is already classifying  $x_i$  correctly, misreporting will not improve the outcome for  $i$ . Assume the prediction is  $h^*(x_i) = y^* \neq y_i$ , and suppose  $h^*$  is making  $q$  mistakes, meaning,

$$q = \min_{h \in \mathcal{H}} \ell(y_i; h^*).$$

Agent  $i$  can misreport in the following two ways: □

1. If agent  $i$  reports  $y_i^\dagger = y^*$ , let the new classifier be  $h^\dagger$ , note that we must have,

$$\ell(y^*; h^\dagger) \leq q - 1,$$

since  $\ell(y^*; h^\dagger) > q - 1 = \ell(y^*; h^*)$  contradicts the optimality of  $h^\dagger$ .

Suppose that agent  $i$  could get her true label with  $h^\dagger$ , meaning  $h^\dagger(x_i) = y_i$ , then,

$$\begin{aligned} \ell(y_i; h^\dagger) &= \ell(y^*; h^\dagger) - 1 \\ &\leq q - 2 \\ &< \ell(y_i; h^*), \end{aligned}$$

which contradicts the optimality of  $h^*$ . Therefore, agent  $i$  cannot improve the outcome by misreporting  $y^*$ .

2. If agent  $i$  reports  $y_i^\dagger = y' \neq y^*$ , let the new classifier be  $h^\dagger$ , note that we must have,

$$\ell(y'; h^\dagger) \leq q,$$

since  $\ell(y'; h^\dagger) > q = \ell(y'; h^*)$  contradicts the optimality of  $h^\dagger$ .

Suppose that agent  $i$  could get her true label with  $h^\dagger$ , then,

$$\begin{aligned} \ell(y_i; h^\dagger) &= \ell(y'; h^\dagger) - 1 \\ &\leq q - 1 \\ &< \ell(y_i; h^*), \end{aligned}$$

which contradicts the optimality of  $h^*$ . Therefore, agent  $i$  cannot improve the outcome by misreporting  $y'$ .

Therefore, no agent can improve the outcome and the dataset is incentive-compatible.

### 2.3 Binary Logistic Regression

Binary logistic regression is always incentive-compatible, and agents with one label do not want to pretend to have the other label. This is not true in general for binary empirical risk minimization classifiers, and additional normalization condition on the loss function is required. The result is discussed in the next section. In the special case of logistic regression, since  $\ell(0; h) + \ell(1; h) = 1$  holds for any  $h \in \mathcal{H}$ , the following Proposition holds.

**Proposition 2.** *Binary logistic classifiers are incentive-compatible.*

*Proof.* For any dataset  $\{(x_i, y_i)\}_{i=1}^n$ , and the logistic hypothesis class  $\mathcal{H}$ , let the optimal classifier in the case every agent reports truthfully be,

$$h^* = \arg \min_{h \in \mathcal{H}} \sum_{i'=1}^n \ell(y_{i'}; h).$$

Fix an agent  $i$  with  $x_i$ , and fix other agents' reports,  $(x_{-i}, y_{-i})$ , define the optimal classifier given the misreport of agent  $i$ ,  $y_i^\dagger = 1 - y_i$  as,

$$h^\dagger = h^* \left( y_i^\dagger \right) = \arg \min_{h \in \mathcal{H}} \sum_{i'=1, i' \neq i}^n \ell(y_{i'}; h) + \ell(y_i^\dagger; h).$$

Now suppose, for a contradiction, that agent  $i$  prefers misreporting,

$$\ell(y_i; h^*) > \ell(y_i; h^\dagger),$$

which implies, since  $y_i^\dagger = 1 - y_i$ , and the special property of the logistic loss function  $\ell(0; h) + \ell(1; h) = 1$ ,

$$\ell(y_i^\dagger; h^*) < \ell(y_i^\dagger; h^\dagger).$$

Note that the above implication only works for binary classification. Due to the optimality of  $h^\dagger$ ,

$$\sum_{i'=1, i' \neq i}^n \ell(y_{i'}; h^\dagger) + \ell(y_i^\dagger; h^\dagger) \leq \sum_{i'=1, i' \neq i}^n \ell(y_{i'}; h^*) + \ell(y_i^\dagger; h^*),$$

using the above inequalities, the comparison can be simplified to,

$$\begin{aligned} \sum_{i'=1, i' \neq i}^n \ell(y_{i'}; h^\dagger) &\leq \sum_{i'=1, i' \neq i}^n \ell(y_{i'}; h^*), \\ \sum_{i'=1}^n \ell(y_{i'}; h^*) &\leq \sum_{i'=1}^n \ell(y_{i'}; h^*), \end{aligned}$$

which is a contradiction to the optimality of  $h^*$ .

□

### 3 Other Classifiers

#### 3.1 Artificial Examples

In this section, we show that classifiers, including Bayes classifiers and kernel density estimators, that satisfy some separability conditions are always incentive-compatible. One intuition behind why some classifiers are incentive-incompatible is that one-vs-one classification decisions are not made independently. Logistic regression has highly interdependent one-vs-one decisions. The following example is one in which 1-vs-2 decisions are completely determined by the 2-vs-3 decisions, and as a result, a class-1 point that is misclassified as class-2 could misreport as class-3 to influence the 2-vs-3 decision boundary and indirectly change the 1-vs-2 decision boundary in its favor.

Consider the three-class threshold classifiers with thresholds symmetric around 0 in the form,

$$\hat{y}(x; t) = \begin{cases} \text{red} & \text{if } x < -t \\ \text{green} & \text{if } -t \leq x \leq t \\ \text{blue} & \text{if } x > t \end{cases}$$

This is effectively a one-vs-one threshold classifier with the red-vs-green threshold of  $-t$ , with the green-vs-blue threshold of  $t$  and the red-vs-blue threshold of any value between  $-t$  and  $t$ . Now define the loss function as the squared error margin,

$$\ell(x_i, y_i; t) = \begin{cases} (\max\{x_i - (-t), 0\})^2 & \text{if } y_i = \text{red} \\ (\min\{\max\{-t - x_i, 0\}, \max\{x_i - t, 0\}\})^2 & \text{if } y_i = \text{green} \\ (\max\{t - x_i, 0\})^2 & \text{if } y_i = \text{blue} \end{cases}$$

Then, the learner's maximization problem is,

$$\min_{t \in \mathbb{R}} \ell(x_i, y_i; t).$$

The agent's problem is to minimize the loss for their point with their true label. The loss can be either the margin loss or zero-one loss, here, for simplicity, assume the agents also want to minimize the margin,

$$\max_{y_i^\dagger \in \{\text{red}, \text{green}, \text{blue}\}} \ell(x_i, y_i; t^*(y_i^\dagger)),$$

where  $t^*(y_i^\dagger)$  is the optimal threshold  $t$  for the learner when agent  $i$  misreports the label as  $y_i^\dagger$ .

For the dataset  $\{(x_i, y_i)\}_{i=1}^5 = \{(-1, \text{red}), (-g, \text{green}), (-r, \text{red}), (g, \text{green}), (1, \text{blue})\}$ , with some  $r, g \in (0, 1)$  and  $r < g < 2r$ , the minimum loss is  $\frac{2}{3}(g-r)^2$  and it occurs when  $t = \frac{2}{3}g + \frac{1}{3}r > r$  when every agent reports truthfully. In this case, agent  $(-r, \text{red})$  is misclassified as green. However, if  $(-r, \text{red})$  misreports as blue, the minimum loss is  $\frac{2}{3}(g+r)^2$  and it occurs when  $t = \frac{2}{3}g - \frac{1}{3}r < r$ . In this case, agent  $(-r, \text{red})$  is classified correctly, and in particular, she improves the loss from  $g-r$  to 0.

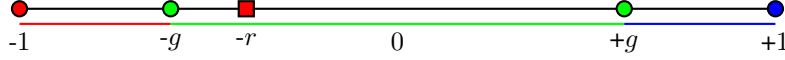


Figure 2: 1D Artificial Incentive-incompatible Example 1 (Truthful)

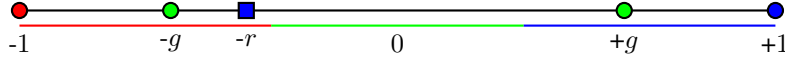


Figure 3: 1D Artificial Incentive-incompatible Example 1 (Misreport)

#### 4 A good example

Consider a 3-way classifier  $h_t : \mathbb{R} \rightarrow \{\text{red}, \text{green}, \text{blue}\}$  parametrized by  $t \geq 0$ :

$$h_t(x) = \begin{cases} \text{red} & \text{if } x < t \\ \text{green} & \text{if } -t \leq x \leq t \\ \text{blue} & \text{if } x > t. \end{cases} \quad (1)$$

Let the hypothesis space be

$$\mathcal{H} = \{h_t : t \geq 0\}. \quad (2)$$

Equivalently, a hypothesis  $h_t$  partitions  $\mathbb{R}$  into three sets:  $X_t^{\text{red}} = (-\infty, -t)$ ,  $X_t^{\text{green}} = [-t, t]$ ,  $X_t^{\text{blue}} = (t, \infty)$ . Given a labeled point  $(x, y)$  with  $y \in \{\text{red}, \text{green}, \text{blue}\}$ , it could be outside the “color region” suggested by  $h_t$ . Accordingly, we define a loss function  $\ell$  based on the distance it takes to move the point to the corresponding color region suggested by  $h_t$ . Concretely,

$$\ell(x, y, h_t) = f(d(x, X_t^y)) \quad (3)$$

where

$$d(x, X_t^y) = \min_{x' \in X_t^y} \|x - x'\| \quad (4)$$

is the shortest distance from the point  $x$  to the set (color region)  $X_t^y$ .  $f \geq 0$  is strictly convex with minimum at 0:  $f(0) = 0$ . For example,  $f$  can be the square function  $f(z) = z^2$ .

Given a training set  $S = \{(x_1, y_1), \dots, (x_n, y_n)\}$ , consider the Empirical Risk Minimizer (ERM)

$$\hat{h} \in \operatorname{argmin}_{h \in \mathcal{H}} \sum_{i=1}^n \ell(x_i, y_i, h). \quad (5)$$

We now exhibit a family of IIC datasets  $S$ , see Figure 2.  $S$  consists of five labeled points:

$$S = \{(x_i, y_i)\}_{i=1}^5 = \{(-1, \text{red}), (-g, \text{green}), (-r, \text{red}), (g, \text{green}), (1, \text{blue})\}, \quad (6)$$

with  $0 < r < g < 2r$ . **Do we need  $g < 1$  or  $2r < 1$ ?**

**Proposition 3.**  $S$  is incentive incompatible with respect to ERM on  $\mathcal{H}$  and  $\ell$ .

*Proof.* The optimal threshold for the original dataset is,

$$t^* = \arg \min_{t \geq 0} \begin{cases} (t-r)^2 & \text{if } t \geq g \\ 2(t-g)^2 + (t-r)^2 & \text{if } r < t < g \\ 2(t-g)^2 & \text{if } t \leq r \end{cases},$$

or,

$$t^* = \arg \min_t \begin{cases} (g-r)^2 & \text{if } t = g \\ \frac{2}{3}(g-r)^2 & \text{if } t = \frac{2}{3}g + \frac{1}{3}r \\ 2(g-r)^2 & \text{if } t = r \end{cases},$$

implying,

$$t^* = \frac{2}{3}g + \frac{1}{3}r \in (r, g),$$

which classifies  $(-r, \text{red})$  as a green point, which is incorrect.

Similarly, the optimal threshold if the agent misreports her label as blue is,

$$t^* = \arg \min_{t \geq 0} \begin{cases} (t+r)^2 & \text{if } t \geq g \\ 2(t-g)^2 + (t+r)^2 & \text{if } r \leq t < g \\ 2(t-g)^2 + (t+r)^2 & \text{if } t < r \end{cases},$$

or,

$$t^* = \arg \min_t \begin{cases} (g+r)^2 & \text{if } t = g \\ > \frac{2}{3}(g+r)^2 & \text{if } t = r \\ \frac{2}{3}(g+r)^2 & \text{if } t = \frac{2}{3}g - \frac{1}{3}r \end{cases},$$

implying,

$$t^* = \frac{2}{3}g - \frac{1}{3}r \in (0, r),$$

which classifies  $(-r, \text{red})$  as a red point, which is correct.

Therefore, the point has the incentive to misreport and the dataset is incentive incompatible. □

**Lemma 1.**  $S_1$  is incentive incompatible for a strictly convex loss margin with thresholds symmetric around 0 if  $r = g - \delta$ , for some small  $\delta > 0$ .

*Proof.* The optimal threshold for the original dataset is,

$$t^* = \arg \min_{t \geq 0} \begin{cases} f(t-r) & \text{if } t \geq g \\ 2f(g-t) + f(t-r) & \text{if } r < t < g \\ 2f(g-t) & \text{if } t \leq r \end{cases},$$

Since  $f(0) = 0$ , the objective is continuous at  $t = g$  and  $t = r$ . Note that  $f(t-r)$  is minimized at  $r < g$  and  $f(t-g)$  is minimized at  $g > r$ , we have,

$$t^* \in (r, g),$$

which classifies  $(-r, \text{red})$  as a green point, which is incorrect.

Similarly, the optimal threshold if the agent misreports her label as blue is,

$$t^* = \arg \min_{t \geq 0} \begin{cases} f(t+r) & \text{if } t \geq g \\ 2f(g-t) + f(t+r) & \text{if } r \leq t < g \\ 2f(g-t) + f(t+r) & \text{if } t < r \end{cases},$$



Since  $f(t+r)$  is minimized at  $t = -r < g$  and  $f(g-t)$  is minimized at  $g$ ,

$$t^* \in (0, g),$$

and the objective in this range is,

$$2f(g-t) + f(t+r).$$

Now fix small  $\varepsilon > 0$ , there is sufficiently small  $\delta > 0$  such that, due to strict convexity,

$$2 \frac{f(\delta + \varepsilon) - f(\delta)}{\varepsilon} < \frac{f(2r) - f(2r - \varepsilon)}{\varepsilon}.$$

This inequality implies that the objective at  $t = r - \varepsilon$ ,

$$2f(\delta + \varepsilon) + f(2r - \varepsilon) < 2f(\delta) + f(2r),$$

which is the objective at  $t = r$ . Therefore, the objective is decreasing at  $r$ , given that it is the sum of two strictly convex functions thus strictly convex itself, we have,

$$t^* \in (0, r),$$

which classifies  $(-r, \text{red})$  as a red point, which is correct.

Therefore, the point has the incentive to misreport and the dataset is incentive incompatible. □

Intuitively, when the misclassified class-1 point pretends to be a class-3 point with a large loss, the classifier modifies the 2-vs-3 decision boundary to minimize that loss, and the point benefits from the interdependence between the 2-vs-3 decision and the 1-vs-2 decision.

The decision boundaries can be related in more complicated ways. Consider the three-class 1D threshold classifiers in the form,

$$\hat{y}(x; l) = \begin{cases} 1 & \text{if } x < l \\ 2 & \text{if } l \leq x \leq r \\ 3 & \text{if } x > r \end{cases}$$

This is effectively a one-vs-one threshold classifier with the 1-vs-2 threshold of  $l$ , with the 2-vs-3 threshold of  $r$  and the 1-vs-3 threshold of any value between  $l$  and  $r$ . Now define the loss function as the squared distance to the center of the decision region (bounded by the range of the data),

$$\ell(x_i, y_i; t) = \begin{cases} \left(x_i - \frac{l + \min_i x_i}{2}\right)^2 & \text{if } y_i = 1 \\ \left(x_i - \frac{l + r}{2}\right)^2 & \text{if } y_i = 2 \\ \left(x_i - \frac{r + \max_i x_i}{2}\right)^2 & \text{if } y_i = 3 \end{cases}$$

Then, the learner's maximization problem is,

$$\min_{t \in \mathbb{R}} \ell(x_i, y_i; t).$$

The agent's problem is to minimize the loss for their point with their true label. The loss can be either the distance-to-center loss or zero-one loss, here, for simplicity, assume the agents also want to minimize the distance to the center,

$$\max_{y_i^\dagger \in \{1, 2, 3\}} \ell\left(x_i, y_i; l^*(y_i^\dagger), r^*(y_i^\dagger)\right),$$

where  $l^*(y_i^\dagger)$  and  $r^*(y_i^\dagger)$  are the optimal thresholds for the learner when agent  $i$  misreports the label as  $y_i^\dagger$ .

For the same dataset  $\{(x_i, y_i)\}_{i=1}^5 = \{(-5, 1), (-3, 2), (-2, 1), (4, 2), (5, 3)\}$ , the minimum loss

is 29.5 and it occurs when  $l = -2.5$  and  $r = 4$  when every agent reports truthfully. In this case, agent  $(-2, 1)$  is misclassified as  $(-2, 2)$ . However, if  $(-2, 1)$  misreports as  $(-2, 3)$ , the loss from  $l = -2.5$  and  $r = 4$  is above 68 and the minimum loss is 57 and it occurs when  $l = -1$  and  $r = 0$ . In this case, agent  $(-2, 1)$  is classified correctly. In particular, the agent  $(-2, 1)$  improves the loss from 1 to 0.

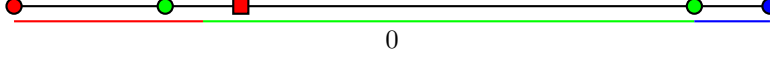


Figure 4: 1D Artificial Incentive-incompatible Example 2 (Truthful)

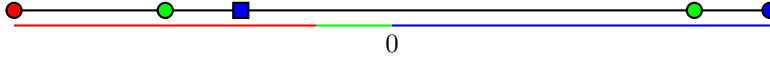


Figure 5: 1D Artificial Incentive-incompatible Example 2 (Misreport)

To generalize the above observations, suppose the learner is a probabilistic classifier with parameters estimated by maximum likelihood, and the agents report their labels to maximize the classification probability of their true labels. Then the following two conditions guarantee that the classification is incentive-compatible.

**Definition 2.** (Monotonic Condition) A multi-class probabilistic classifier is monotonic if, given a training set  $S$ , for any point  $x$  with labels  $a$  and  $b$ ,

$$\frac{\mathbb{P}\{Y = a|x; w^*(S)\}}{\mathbb{P}\{Y = b|x; w^*(S)\}} \geq \frac{\mathbb{P}\{Y = a|x; w^*(S \cup \{(x, y = a)\})\}}{\mathbb{P}\{Y = b|x; w^*(S \cup \{(x, y = a)\})\}}.$$

The assumption says that the probability that  $x$  is classified as  $a$  increases when there is an additional point  $(x, a)$  in the training set.

**Definition 3.** (Independence of Irrelevant Alternatives (IIA) Condition) A multi-class classifier is independent of irrelevant alternatives if, given a training set  $S$ , for any point  $x$  and any pair of labels  $a$  and  $b$ ,

$$\frac{\mathbb{P}\{Y = a|x; w^*(S)\}}{\mathbb{P}\{Y = b|x; w^*(S)\}} = \frac{\mathbb{P}\{Y = a|x; w^*(S \cup \{(x', y' \notin \{a, b\})\})\}}{\mathbb{P}\{Y = b|x; w^*(S \cup \{(x', y' \notin \{a, b\})\})\}}.$$

The assumption says that the ratio between the classification probabilities of  $x$  of any two classes is not changed by adding a point at  $x$  with a third class.

Combining the two assumptions MC and IIA, we have that an agent with label  $a$  cannot change the decision of  $a$  vs  $b$  by misreporting its label as a third class  $c$ . This observation is formalized in the following proposition.

**Theorem 1.** A multi-class probabilistic classifier estimated by maximum likelihood is incentive-compatible if it is monotonic and independent of irrelevant alternatives.

*Proof.* Fix a dataset  $\{(x_i, y_i)\}_{i=1}^n$ , let the maximum likelihood estimates in the case every agent report truthfully be,

$$w^* = \arg \max_w \sum_{i'=1}^n \log(\mathbb{P}\{Y = y_{i'}|x_{i'}; w\}).$$

Fix an agent  $i$ , her feature vector  $x_i$ , and fix other agents' reports,  $(x_{-i}, y_{-i})$ , define the maximum likelihood estimate given the misreport of agent  $i$ ,  $y_i^\dagger$  as,

$$w^\dagger = \arg \max_w \sum_{i'=1, i' \neq i}^n \log (\mathbb{P} \{Y = y_{i'} | x_{i'}, w\}) + \log (\mathbb{P} \{Y = y_i^\dagger | x_i; w\}).$$

Now suppose, for a contradiction, that agent  $i$  prefers misreporting, assume the following incentive inequality,

$$\mathbb{P} \{Y = y_i | x_i; w^\star\} > \mathbb{P} \{Y = y_i | x_i; w^\dagger\}.$$

If there are only two classes, then by symmetry,

$$\mathbb{P} \{Y = y_i^\dagger | x_i; w^\star\} < \mathbb{P} \{Y = y_i^\dagger | x_i; w^\dagger\}.$$

If there are more than two classes, fix a third  $y'_i \notin \{y_i, y_i^\dagger\}$ , and define an intermediate maximum likelihood estimate from removing the point  $(x_i, y_i)$ ,

$$w' = \arg \max_w \sum_{i'=1, i' \neq i}^n \log (\mathbb{P} \{Y = y_{i'} | x_{i'}, w\}),$$

then the Monotonic Condition implies,

$$\frac{\mathbb{P} \{Y = y_i | x_i; w^\star\}}{\mathbb{P} \{Y = y'_i | x_i; w^\star\}} \leq \frac{\mathbb{P} \{Y = y_i | x_i; w'\}}{\mathbb{P} \{Y = y'_i | x_i; w'\}},$$

and the IIA Condition implies,

$$\frac{\mathbb{P} \{Y = y_i | x_i; w'\}}{\mathbb{P} \{Y = y'_i | x_i; w'\}} = \frac{\mathbb{P} \{Y = y_i | x_i; w^\dagger\}}{\mathbb{P} \{Y = y'_i | x_i; w^\dagger\}}.$$

Combining the above two inequalities with the incentive inequality, we have,

$$\mathbb{P} \{Y = y'_i | x_i; w^\star\} > \mathbb{P} \{Y = y'_i | x_i; w^\dagger\}.$$

Note that the above inequality is true for all  $y'_i \notin \{y_i, y_i^\dagger\}$ , summing over all such  $y'_i$  results in,

$$\sum_{y' \notin \{y_i, y_i^\dagger\}} \mathbb{P} \{Y = y'_i | x_i; w^\star\} > \sum_{y' \notin \{y_i, y_i^\dagger\}} \mathbb{P} \{Y = y'_i | x_i; w^\dagger\},$$

given that the class probabilities sum up to 1,

$$1 - \mathbb{P} \{Y = y_i | x_i; w^\star\} - \mathbb{P} \{Y = y_i^\dagger | x_i; w^\star\} > 1 - \mathbb{P} \{Y = y_i | x_i; w^\dagger\} - \mathbb{P} \{Y = y_i^\dagger | x_i; w^\dagger\},$$

and using the incentive inequality again,

$$\mathbb{P} \{Y = y_i^\dagger | x_i; w^\star\} < \mathbb{P} \{Y = y_i^\dagger | x_i; w^\dagger\}.$$

Now, due to the optimality of  $h^\dagger$ ,

$$\begin{aligned} & \sum_{i'=1, i' \neq i}^n \log (\mathbb{P} \{Y = y_{i'} | x_{i'}, w^\dagger\}) + \log (\mathbb{P} \{Y = y_i^\dagger | x_i; w^\dagger\}) \\ & \leq \sum_{i'=1, i' \neq i}^n \log (\mathbb{P} \{Y = y_{i'} | x_{i'}, w^\star\}) + \log (\mathbb{P} \{Y = y_i^\dagger | x_i; w^\star\}), \end{aligned}$$

using the above inequalities, the comparison can be simplified to,

$$\begin{aligned} & \sum_{i'=1, i' \neq i}^n \log (\mathbb{P} \{Y = y_{i'} | x_{i'}, w^\dagger\}) < \sum_{i'=1, i' \neq i}^n \log (\mathbb{P} \{Y = y_{i'} | x_{i'}, w^\star\}), \\ & \sum_{i'=1}^n \log (\mathbb{P} \{Y = y_{i'} | x_{i'}, w^\dagger\}) < \sum_{i'=1}^n \log (\mathbb{P} \{Y = y_{i'} | x_{i'}, w^\star\}), \end{aligned}$$

which is a contradiction to the optimality of  $w^\star$ .

□

**Corollary 1.** *Binary probabilistic classifiers estimated by maximum likelihood are incentive-compatible.*

*Proof.* MC holds due to the optimality conditions and IIA holds since there are only two classes. □

The assumptions Definition 2 (MC) and Definition 3 (IIA) can be significantly simplified for a separable class of the classifiers.

**Definition 4.** A probabilistic classifier is separable if the parameters  $w^*$  can be partitioned into  $k$  classes,  $w_1^*, w_2^*, \dots, w_k^*$ , one set of parameters for each class, such that for given training sets  $S$  and  $S'$  and any label  $a$  if,

$$\{(x_i, y_i) \in S : y_i = a\} = \{(x_i, y_i) \in S' : y_i = a\},$$

then,

$$w_a^*(S) = w_a^*(S'),$$

and if there is a value function  $v_a(x; w_a^*)$  that are independent of  $w_b^*, b \neq a$  such that,

$$\mathbb{P}\{Y = y|x; w^*\} = \frac{v_a(x; w_a^*)}{\sum_{b=1}^K v_b(x; w_b^*)},$$

then the classifier is separable.

Logistic regression satisfies the value function requirement but fails the separability condition since training  $w_a^*$  uses data with labels that are not  $a$ . On the other hand, Bayes-type classifiers are separable. For separable classifiers, Definition 2 (MC) and Definition 3 (IIA) are always satisfied.

**Corollary 2.** *A separable multi-class probabilistic classifier estimated by maximum likelihood is incentive-compatible.*

*Proof.* Due to separability,

$$\begin{aligned} v_a(x; w^*(S)) &= v_a(x; w^*(S \cup \{(x', y' \notin \{a, b\})\})), \text{ and} \\ v_b(x; w^*(S)) &= v_b(x; w^*(S \cup \{(x, y = a)\})) = v_b(x; w^*(S \cup \{(x', y' \notin \{a, b\})\})). \end{aligned}$$

MC follows from the optimality condition of  $w^*(S)$  and IIA follows immediately. □

**Corollary 3.** *Bayes classifiers estimated by maximum likelihood are incentive-compatible.*

*Proof.* Follows from Corollary 2. □

Kernel density estimators are not estimated by maximum likelihood, so the previous results do not hold, although the proof is similar. There are two general approaches to use kernel densities for classification, the first is to use all the points to estimate the density and the second is to estimate the densities for each class separately (see Taylor (1997)). The second approach is similar to a separable classifier. K-Nearest Neighbor is a special case of this with a uniform kernel.

**Corollary 4.** *Kernel density estimators are incentive-compatible.*

*Proof.* The first approach suggests that,

$$\mathbb{P}\{X = x\} = \frac{1}{nh^D} \sum_{i'=1}^n w_{i'}, \text{ where } w_{i'} = K\left(\frac{x - x_{i'}}{h}\right).$$

Define  $w^\star$  as the weight function when all agents report truthfully, and  $w^\dagger = w^\star(y_i^\dagger)$  as the weight if agent  $i$  misreports, then the classification probabilities for agent  $i$  from dividing up the sum based on the class is,

$$\begin{aligned} \mathbb{P}\{Y = y_i | x_i; w^\star\} &= \frac{1}{n} \sum_{i'=1}^n w_{i'}^\star \mathbb{1}_{\hat{y}_{i'}=y_i} \\ &= \frac{1}{n} \sum_{i'=1}^n w_{i'}^\star \mathbb{1}_{\hat{y}_{i'}=y_i, i' \neq i} + \frac{1}{n} w_i^\star \\ &= \mathbb{P}\{Y = y_i | x_i; w^\dagger\} + \frac{1}{n} K(0), \text{ since } y_i^\dagger \neq y_i \\ &\leq \mathbb{P}\{Y = y_i | x = x_i; w^\star\}, \end{aligned}$$

meaning reporting truthfully results in a larger probability compared to reporting  $y_i^\dagger$  instead. Alternatively, the second approach suggests that, if the classification probabilities are computed based on Taylor (1997),

$$\mathbb{P}\{X = x | Y = y\} = \frac{1}{n_y h^D} \sum_{i'=1}^n w_{i'} \mathbb{1}_{\hat{y}_{i'}=y}, n_y = \sum_{i'=1}^n \mathbb{1}_{\hat{y}_{i'}=y}$$

Similar to the above derivation (and also as a special case of a Bayes estimator),

$$\begin{aligned} \mathbb{P}\{Y = y_i | x_i; w^\star\} &= \frac{1}{n_y} \sum_{i'=1}^n w_{i'}^\star \mathbb{1}_{\hat{y}_{i'}=y_i} \\ &= \mathbb{P}\{Y = y_i | x_i; w^\star\} + \frac{1}{n_y} K(0) \\ &\leq \mathbb{P}\{Y = y_i | x_i; w^\star\}. \end{aligned}$$

□

#### 4.1 Empirical Risk Minimization

A similar result can be obtained for empirical risk minimization. We could either add an assumption that the loss function can be normalized so that the sum is constant and it behaves the same way as a probabilistic classifier, or we could use stronger Monotonic and IIA Conditions. Here, we state the additional normalization condition.

**Definition 5.** (Normalized Loss) A loss function  $\ell$  is normalized if given a hypothesis  $h$ , for any point,

$$\sum_y \ell(y; h) = C, \text{ constant.}$$

**Definition 6.** (Monotonic Condition for ERM) Multi-class empirical risk minimization classifiers are monotonic if, given a training set  $S$ , for any point  $x$  with labels  $a$  and  $b$ ,

$$\frac{\ell(y = a; h^\star(S))}{\ell(y = b; h^\star(S))} \geq \frac{\ell(y = a; h^\star(S \cup \{(x, y = a)\}))}{\ell(y = b; h^\star(S \cup \{(x, y = a)\}))}.$$

**Definition 7.** (IIA Condition for ERM) Multi-class empirical risk minimization classifiers are independent of irrelevant alternatives if, given a training set  $S$ , for any point  $x$  and any pair of labels  $a$  and  $b$ ,

$$\frac{\ell(y = a; h^*(S))}{\ell(y = b; h^*(S))} = \frac{\ell(y = a; h^*(S \cup \{(x', y' \notin \{a, b\})\}))}{\ell(y = b; h^*(S \cup \{(x', y' \notin \{a, b\})\}))}.$$

**Corollary 5.** Multi-class empirical risk minimization classifiers with normalized loss functions are incentive-compatible if it is monotonic and independent of irrelevant alternatives.

*Proof.* For a fixed dataset  $\{(x_i, y_i)\}_{i=1}^n$ , and the hypothesis class  $\mathcal{H}$ , let the optimal classifier in the case every agent report truthfully be,

$$h^* = \arg \min_{h \in \mathcal{H}} \sum_{i'=1}^n \ell(y_{i'}; h).$$

Fix an agent  $i$ , her feature vector  $x_i$ , and fix other agents' reports,  $(x_{-i}, y_{-i})$ , define the optimal classifier given the classifier  $h$  and the misreport of agent  $i$ ,  $y_i^\dagger$  as,

$$h^\dagger = \arg \min_{h \in \mathcal{H}} \sum_{i'=1, i' \neq i}^n \ell(y_{i'}; h) + \ell(y_i^\dagger; h).$$

Now suppose, for a contradiction, that agent  $i$  prefers misreporting, assume the following incentive inequality,

$$\ell(y_i; h^*) > \ell(y_i; h^\dagger).$$

If there are only two classes, then by symmetry,

$$\ell(y_i^\dagger; h^*) < \ell(y_i^\dagger; h^\dagger).$$

If there are more than two classes, fix a third  $y'_i \notin \{y_i, y_i^\dagger\}$ , and define an intermediate maximum likelihood estimate from removing the point  $(x_i, y_i)$ ,

$$h' = \arg \min_{h \in \mathcal{H}} \sum_{i'=1, i' \neq i}^n \ell(y_{i'}; h),$$

then the Monotonic Condition for ERM implies,

$$\frac{\ell(y_i; h^*)}{\ell(y'_i; h^*)} \leq \frac{\ell(y_i; h')}{\ell(y'_i; h')},$$

and the IIA Condition implies,

$$\frac{\ell(y_i; h')}{\ell(y'_i; h')} = \frac{\ell(y_i; h^\dagger)}{\ell(y'_i; h^\dagger)}.$$

Combining the above two inequalities with the incentive inequality, we have,

$$\ell(y'_i; h^*) > \ell(y'_i; h^\dagger).$$

Note that the above inequality is true for all  $y'_i \notin \{y_i, y_i^\dagger\}$ , summing over all such  $y'_i$  results in,

$$\sum_{y' \notin \{y_i, y_i^\dagger\}} \ell(y'; h^*) > \sum_{y' \notin \{y_i, y_i^\dagger\}} \ell(y'; h^\dagger),$$

and given the losses are normalized,

$$C - \ell(y_i; h^*) - \ell(y_i^\dagger; h^*) > C - \ell(y_i; h^\dagger) - \ell(y_i^\dagger; h^\dagger),$$

and using the incentive inequality again,

$$\ell(y_i^\dagger; h^*) < \ell(y_i^\dagger; h^\dagger).$$

Now, due to the optimality of  $h^\dagger$ ,

$$\sum_{i'=1, i' \neq i}^n \ell(y_{i'}; h^\dagger) + \ell(y_i^\dagger; h^\dagger) \leq \sum_{i'=1, i' \neq i}^n \ell(y_{i'}; h^*) + \ell(y_i^\dagger; h^*),$$

using the above inequalities, the comparison can be simplified to,

$$\begin{aligned} \sum_{i'=1, i' \neq i}^n \ell(y_{i'}; h^\dagger) &\leq \sum_{i'=1, i' \neq i}^n \ell(y_{i'}; h^*), \\ \sum_{i'=1}^n \ell(y_{i'}; h^\dagger) &\leq \sum_{i'=1}^n \ell(y_{i'}; h^*), \end{aligned}$$

which is a contradiction to the optimality of  $h^*$ .

□

## References

- Cai, Yang, Constantinos Daskalakis, and Christos Papadimitriou (2015), “Optimum statistical estimation with strategic data sources.” In *Conference on Learning Theory*, 280–296.
- Chen, Yiling, Yang Liu, and Chara Podimata (2019), “Grinding the space: Learning to classify against strategic agents.” *arXiv preprint arXiv:1911.04004*.
- Chen, Yiling, Chara Podimata, Ariel D Procaccia, and Nisarg Shah (2018), “Strategyproof linear regression in high dimensions.” In *Proceedings of the 2018 ACM Conference on Economics and Computation*, 9–26.
- Dekel, Ofer, Felix Fischer, and Ariel D Procaccia (2010), “Incentive compatible regression learning.” *Journal of Computer and System Sciences*, 76, 759–777.
- Dekel, Ofer and Ohad Shamir (2009a), “Good learners for evil teachers.” In *Proceedings of the 26th annual international conference on machine learning*, 233–240.
- Dekel, Ofer and Ohad Shamir (2009b), “Vox populi: Collecting high-quality labels from a crowd.” In *COLT*.
- Dütting, Paul, Zhe Feng, Harikrishna Narasimhan, David C Parkes, and Sai Srivatsa Ravindranath (2017), “Optimal auctions through deep learning.” *arXiv preprint arXiv:1706.03459*.
- Epasto, Alessandro, Mohammad Mahdian, Vahab Mirrokni, and Song Zuo (2018), “Incentive-aware learning for large markets.” In *Proceedings of the 2018 World Wide Web Conference*, 1369–1378.
- Gast, Nicolas, Stratis Ioannidis, Patrick Loiseau, and Benjamin Roussillon (2013), “Linear regression from strategic data sources.” *arXiv preprint arXiv:1309.7824*.
- Golowich, Noah, Harikrishna Narasimhan, and David C Parkes (2018), “Deep learning for multi-facility location mechanism design.” In *IJCAI*, 261–267.
- Meir, Reshef, Ariel D Procaccia, and Jeffrey S Rosenschein (2012), “Algorithms for strategyproof classification.” *Artificial Intelligence*, 186, 123–156.
- Mihailescu, Marian and Yong Meng Teo (2010), “Strategy-proof dynamic resource pricing of multiple resource types on federated clouds.” In *International Conference on Algorithms and Architectures for Parallel Processing*, 337–350, Springer.
- Perote, Javier and Juan Perote-Pena (2004), “Strategy-proof estimators for simple regression.” *Mathematical Social Sciences*, 47, 153–176.
- Procaccia, Ariel D and Moshe Tennenholtz (2009), “Approximate mechanism design without money.” In *Proceedings of the 10th ACM conference on Electronic commerce*, 177–186.

Richardson, Adam, Ljubomir Rokvic, Aris Filos-Ratsikas, and Boi Faltings (2019), “Privately computing influence in regression models.”

Shah, Nihar B and Dengyong Zhou (2016), “Double or nothing: Multiplicative incentive mechanisms for crowdsourcing.” *The Journal of Machine Learning Research*, 17, 5725–5776.

Taylor, Charles (1997), “Classification and kernel density estimation.” *Vistas in Astronomy*, 41, 411–417.