

Strategy Free Machine Learning

-

April 9, 2020

1 Short Summary

There are n strategic agents each providing the label of one data point to the principal. The principal is the learner and builds a machine learning model based on the data points provided by the agents. An agent, i , has publicly known feature vector, x_i , and a private discrete label, y_i . The objective of the agent is to maximize the probability that her data point is labeled correctly by the principal's model, and the agent can choose to report y_i^\dagger to achieve the objective, with the possibility of misreporting $y_i^\dagger \neq y_i$. We say a dataset is incentive incompatible with respect to the learner, described by a parametric model, if at least one of the n agents has the incentive to misreport.

The following is the diagram showing a dataset that is incentive incompatible with respect to the multi-class logistic regression model. In the dataset, each of the $n = 18$ agents, i , has a two dimensional feature vector and a private label can take on one of three values: "red", "green", or "blue".

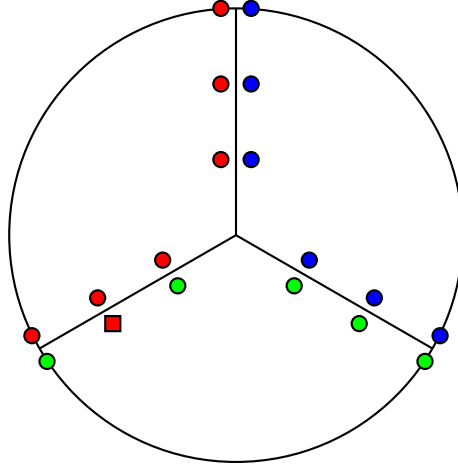


Figure 1: Incentive Incompatible Example

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 form angles of 120 degrees between them. There is one point, labeled by a square in the plot, that is on the "incorrect" side of the boundary. Suppose the point corresponds to the feature vector of an agent i with private label "red", then 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" with probability 0.4966. Therefore, by lying about her label, the agent is able to make the principal learn an incorrect model that classifies her point correctly and with a higher probability.

The same dataset is also incentive incompatible with respect to the one-vs-rest linear support vector machine if the margin is used as the class "probabilities". However, in this case, the agent, with feature vector corresponding to the blue point close to the center and close to the green point, can only improve the margin slightly without making the model switch from classifying her point incorrectly to classifying her point correctly.

However, this dataset is incentive compatible with respect to the Naive Bayes classifiers, and in general, there does not exist any dataset that is incentive incompatible with respect to discrete-valued Naive Bayes classifiers. Misreporting will always lead to lower posterior probability of the agent's true label. In addition, no dataset is incentive incompatible with respect to classifiers that minimizes empirical risk with zero-one loss.

2 Literature Review

Previous work on mechanism design for machine learning with strategic data sources focus 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 machine learner (principal).

- The first group of papers focuses on principal-agent problems 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' 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' feature vectors are also private. Chen, Liu, and Podimata (2019) investigates such problems for linear regressions.
 3. Other models do not distinguish between feature vectors and labels. Each agent has a private valuation. These problems are usually modeled as facility locations problems and the solution involves some variant of the Vickrey-Clarke-Groves or Meyerson auction. These 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).

3 Logistic Regression

3.1 Model and Example

In this section, 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 feature vector, $x_i \in \mathbb{R}^m$, and private discrete label, $y_i \in \{1, 2, \dots, k\}$, has the objective of maximizing the probability that her data point is labeled correctly by the principal's model, 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 $w^\star(y_i^\dagger)$, the agent's objective can be written as,

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

where,

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

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

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

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

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

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

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

Definition 1. A dataset is incentive incompatible with respect to a learner if there exists at least one agent i , and some $y_i^\dagger \neq y_i$ such that,

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

A learner (algorithm) is incentive compatible if there does not exist a dataset that is incentive incompatible.

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

Proof. The example given previously is a dataset that is incentive incompatible. □

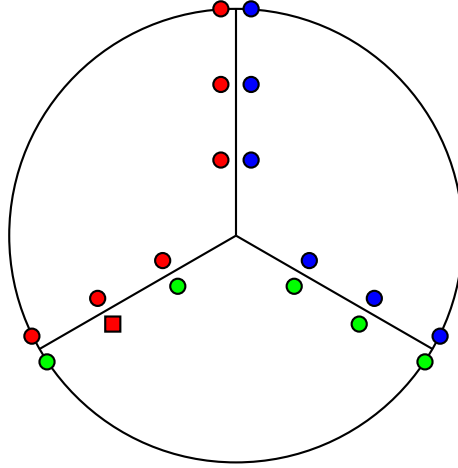


Figure 2: Incentive Incompatible Example

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 correspond to agent 1 with $x_1 = (-1.63, -1.17)$ and $y_1 = 1$.

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

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
2	-0.6053178	104.9925	-181.3391914
3	-0.2852057	209.4190	0.3656777

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

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

3.2 Incentive Incompatibility

To characterize the set of incentive incompatible datasets, we rewrite the principal's choice of optimal weights by,

$$w^* = \arg \max_w \log (\mathbb{P} \{Y = y_i | w, x_i\}) + C_{-i}(w),$$

where the function $C_{-i}(w)$ summarizes the loss due to agents other than i , assuming they are reporting labels truthfully,

$$C_{-i}(w) = \sum_{i'=0, i' \neq i}^n \log (\mathbb{P} \{Y = y_{i'} | w, x_{i'}\}).$$

Since the objective is globally convex and differentiable, as shown in , the problem translates to the first derivative condition,

$$\begin{aligned} \frac{\nabla_w \mathbb{P} \{Y = y_i | w^*, x_i\}}{\mathbb{P} \{Y = y_i | w^*, x_i\}} + \nabla_w (C_{-i}(w^*)) &= 0, \\ \frac{\nabla_w \mathbb{P} \{Y = y_i^\dagger | w^*(y_i^\dagger), x_i\}}{\mathbb{P} \{Y = y_i^\dagger | w^*(y_i^\dagger), x_i\}} + \nabla_w (C_{-i}(w^*(y_i^\dagger))) &= 0. \end{aligned}$$

For logistic regression with weights $w_c, c = 1, 2, \dots, k$, without normalization,

$$\begin{aligned} \mathbb{P} \{Y = c | w, x\} &= \frac{e^{w_c^T x + b_c}}{\sum_{c'} e^{w_{c'}^T x + b_{c'}}}, \\ \nabla_{w_c} \mathbb{P} \{Y = c | w, x\} &= \frac{e^{w_c^T x + b_c} \sum_{c' \neq c} e^{w_{c'}^T x + b_{c'}}}{\left(\sum_{c'} e^{w_{c'}^T x + b_{c'}} \right)^2} x. \\ \nabla_{w_c} \mathbb{P} \{Y = \hat{c}, \hat{c} \neq c | w, x\} &= \frac{e^{w_c^T x + b_c} e^{w_{\hat{c}}^T x + b_{\hat{c}}}}{\left(\sum_{c'} e^{w_{c'}^T x + b_{c'}} \right)^2} x. \end{aligned}$$

The derivative conditions implies,

$$\begin{aligned} (1 - \mathbb{P} \{Y = c | w^*, x_i\}) x_i + \nabla_{w_c} (C_{-i}(w^*)) &= 0, c = y_i, \\ (\mathbb{P} \{Y = c | w^*, x_i\}) x_i + \nabla_{w_c} (C_{-i}(w^*)) &= 0, c \neq y_i, \end{aligned}$$

same for the expression with $w^\star(y_i^\dagger)$.

Substitute into the incentive incompatibility condition,

$$\nabla_{w_{y_i,j}}(C_{-i}(w^\star))x_{i,j} \leq \nabla_{w_{y_i,j}}\left(C_{-i}\left(w^\star(y_i^\dagger)\right)\right)x_{i,j}, j = 1, 2, \dots, m.$$

3.3 Continuous Label

The previous formulation does not permit y_i^\dagger to be a continuous variable, but if we rewrite the optimization as the maximization of the cross entropy, then we could treat $y_i^\dagger \in \Delta^{K-1}$ as a continuous multinomial distribution where $y_{i,c}^\dagger \in [0, 1]$ denotes the probability of agent i reporting label $c \in \{1, 2, \dots, K\}$. The principal's problem can be rewritten as,

$$\min_w \sum_{i=1}^n \sum_{c=1}^K -y_{i,c}^\dagger \log(\mathbb{P}\{Y = c|w, x_i\}).$$

Assuming $w^\star(y_i^\dagger)$ is the optimal weights, the objective function becomes,

$$\mathcal{L}(w, y_i^\dagger) = \sum_{c=1}^k -y_{i,c}^\dagger \log(\mathbb{P}\{Y = c|w, x_i\}) - \sum_{i'=1, i' \neq i}^n \log(\mathbb{P}\{Y = y_{i'}|w, x_{i'}\}),$$

and the value function is,

$$\mathcal{L}^\star(y_i^\dagger) = \sum_{c=1}^k -y_{i,c}^\dagger \log(\mathbb{P}\{Y = c|w^\star(y_i^\dagger), x_i\}) - \sum_{i'=1, i' \neq i}^n \log(\mathbb{P}\{Y = y_{i'}|w^\star(y_i^\dagger), x_{i'}\}),$$

and apply the envelope theorem,

$$\begin{aligned} \frac{\partial \mathcal{L}^\star(y_i^\dagger)}{\partial y_i^\dagger} &= -\log(\mathbb{P}\{Y = y_i^\dagger|w^\star(y_i^\dagger), x_i\}) \\ &> 0. \end{aligned}$$

Alternatively, if gradient descent is used in the optimization process, one iteration of the gradient descent with learning rate η is given by,

$$w'_{j,c} = w_{j,c} - \eta x_{i,j} (\mathbb{P}\{Y = c|x_i\} - \mathbb{1}_{y_i^\dagger}).$$

Now fix instance i and define $o_c = \mathbb{P}\{Y = c|x_i\}$, then,

$$\begin{aligned} \frac{\partial o_c}{\partial y_c} &= \frac{\partial o_c}{\partial z_c} \sum_{j=1}^m \frac{\partial z_c}{\partial w_{j,c}} \frac{\partial w_{j,c}}{\partial y_c} \\ &= o_c (1 - o_c) \sum_{j=1}^m x_j^{(i)} x_j^{(i)} \eta \\ &= \eta o_c (1 - o_c) \sum_{j=1}^m \left(x_j^{(i)}\right)^2 \\ &\geq 0. \end{aligned}$$

Similarly,

$$\begin{aligned}
\frac{\partial o_c}{\partial y_{c'}} &= \frac{\partial o_c}{\partial z_{c'}} \sum_{j=1}^m \frac{\partial z_{c'}}{\partial w_{j,c'}} \frac{\partial w_{j,c'}}{\partial y_{c'}} \\
&= -o_c o_{c'} \sum_{j=1}^m x_j^{(i)} x_j^{(i)} \eta \\
&= -\eta o_c o_{c'} \sum_{j=1}^m \left(x_j^{(i)}\right)^2 \\
&\leq 0.
\end{aligned}$$

This implies decreasing $y_{i,c}$ and thus increasing $y_{i,c'}$ for some c' will always increase o_c . Therefore, there should be no incentive to misreport by changing y_i^\dagger slightly from y_i .

3.4 Zero-One Loss Logistic Regression

With zero-one loss, in general, for an agent with label y ,

If $\hat{y} = y$, misreporting could only make $\hat{y} \neq y$.

If $\hat{y} \neq y$,

1. switching to report \hat{y} would improve the zero-one loss, and will not change the prediction \hat{y} . It is impossible to increase the loss by more than 1, because if it were possible, then the original model was not optimal.
2. switching to $y' \neq \hat{y}$ would keep the same zero-one loss, and will either not change the prediction \hat{y} , or will change the prediction to $y' \neq y$, which does not make the agent better off.

4 Naive Bayes Model

No example in which agents have incentive to misreport is found for the Gaussian Naive Bayes estimator.

The prediction is given by:

$$\arg \max_c \mathbb{P}\{Y = c\} \mathbb{P}\{x^{(i)} | Y = c\}$$

Consider agent i with $(x^{(i)}, y^{(i)})$, if she reports truthfully,

$$\begin{aligned}
\mathbb{P}\{Y = y^{(i)}\} &= \left(\frac{n_{y^{(i)}}}{n}\right), \\
\mathbb{P}\{x^{(i)} | Y = y^{(i)}\} &\geq \mathbb{P}\{x^{(i)} | Y = \hat{y}\},
\end{aligned}$$

and if she reports $\hat{y} \neq y^{(i)}$,

$$\mathbb{P}\{Y = y^{(i)}\} = \frac{n_{y^{(i)}} - 1}{n},$$

$$\mathbb{P}\left\{x^{(i)}|Y=y^{(i)}\right\}\leqslant\mathbb{P}\left\{x^{(i)}|Y=\hat{y}\right\},$$

and both terms are smaller.

Given the dataset, the parameters $w = (\mu, \Sigma, \pi)$ (mean, variance, prior) are given by,

$$\begin{aligned}\mu_c &= \frac{\sum_{i'=1}^n x^{(i')} \mathbb{1}_{y^{(i')}=c}}{n_c}, \\ \Sigma_c &= \frac{\sum_{i'=1}^n \left(x^{(i')} - \mu_c\right) \left(x^{(i')} - \mu_c\right)^T \mathbb{1}_{y^{(i')}=c}}{n_c}, \\ \pi_c &= \frac{n_c}{n}.\end{aligned}$$

Then, the classification probability is,

$$\mathbb{P}_c\left\{Y=c|x^{(i)}\right\}\propto\pi_c\frac{1}{|\Sigma_c|}\exp\left(-\frac{1}{2}\left(x^{(i)}-\mu_c\right)^T(\Sigma_c)^{-1}\left(x^{(i)}-\mu_c\right)\right).$$

The condition for incentive compatibility is,

$$\pi'_c\frac{1}{|\Sigma'_c|}\exp\left(-\frac{1}{2}\left(x^{(i)}-\mu'_c\right)^T(\Sigma'_c)^{-1}\left(x^{(i)}-\mu'_c\right)\right)<\pi_c\frac{1}{|\Sigma_c|}\exp\left(-\frac{1}{2}\left(x^{(i)}-\mu_c\right)^T(\Sigma_c)^{-1}\left(x^{(i)}-\mu_c\right)\right),$$

where,

$$\begin{aligned}\mu'_c &= \frac{\sum_{i'\neq i} x^{(i')} \mathbb{1}_{y^{(i')}=c}}{n_c - 1}, \\ \Sigma'_c &= \frac{\sum_{i'\neq i} \left(x^{(i')} - \mu_c\right) \left(x^{(i')} - \mu_c\right)^T \mathbb{1}_{y^{(i')}=c}}{n_c - 1}, \\ \pi'_c &= \frac{n_c - 1}{n}.\end{aligned}$$

In one-dimensional case, if agent i with $x^{(i)}$ and label $y^{(i)} = c$ reports truthfully,

$$\begin{aligned}\mathbb{P}\left\{Y=c|x^{(i)}\right\} &= \frac{1}{\sigma_c\sqrt{2\pi}}\exp\left(-\frac{\left(x^{(i)}-\mu_c\right)^2}{2\sigma_c^2}\right), \\ \mathbb{P}'\left\{Y=c|x^{(i)}\right\} &= \frac{1}{\sigma'_c\sqrt{2\pi}}\exp\left(-\frac{\left(x^{(i)}-\mu'_c\right)^2}{2\left(\sigma'_c\right)^2}\right),\end{aligned}$$

where,

$$\begin{aligned}\mu'_c &= \frac{n_c\mu}{n_c - 1} - \frac{x^{(i)}}{n_c - 1}, \\ \left(\sigma'_c\right)^2 &= \frac{n_c\sigma_c^2}{n_c - 1} - \frac{\left(x^{(i)} - \mu_c\right)\left(x^{(i)} - \mu'_c\right)}{n_c - 1}\end{aligned}$$

$$= \frac{n_c \sigma_c^2}{n_c - 1} - \frac{n_c (x^{(i)} - \mu_c)^2}{(n_c - 1)^2}.$$

If $x^{(i)} > \mu_c$, then $\mu'_c < \mu_c$, but it is not clear how σ'_c and σ_c compares? Is it possible to have the following scenario?

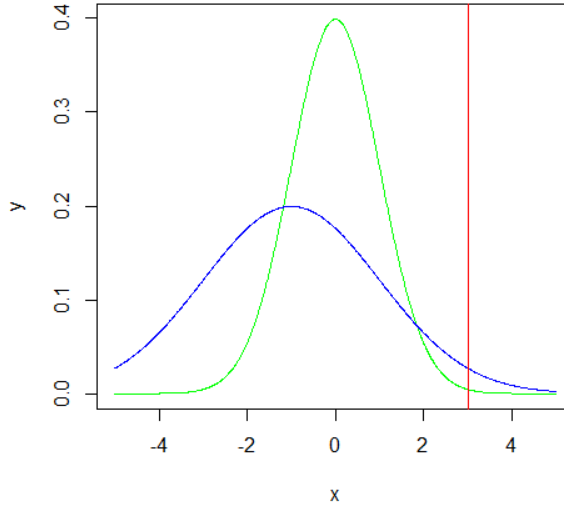


Figure 3: Normal Means

-

5 Support Vector Machines

5.1 One-vs-One

Since binary SVM is incentive compatible, no agent can gain from misreporting in any of the one-vs-one SVMs. Therefore, there will be no incentive to misreport in the multi-class SVM.

5.2 One-vs-Rest

If margin is used as the prediction probabilities, then it is possible to improve the margin by misreporting the third class label, for example on the 18-point data set.

5.3 Tree-Based

Since binary SVM is incentive compatible, no agent can gain from misreporting in any stage. Therefore, there will be no incentive to misreport in the multi-class SVM.

6 Numerical Results

6.1 Structured Examples

A dataset is incentive incompatible with respect to the model parameterized by w if at least one agent has the incentive to report $\hat{y}^{(i)} \neq y^{(i)}$. Formally, let the model estimated with the true label be,

$$w^* = \arg \max_w \sum_{i'=1^n} \log \left(\mathbb{P} \left\{ Y = y^{(i')} | w, x^{(i')} \right\} \right) + \lambda \|w\|,$$

and the model estimated with the misreported label be,

$$\hat{w} = \arg \max_w \left(\sum_{i' \neq i} \log \left(\mathbb{P} \left\{ Y = y^{(i')} | w, x^{(i')} \right\} \right) \right) + \log \left(\mathbb{P} \left\{ Y = \hat{y}^{(i)} | w, x^{(i)} \right\} \right) + \lambda \|w\|.$$

The parameter λ is the regularization parameter.

The dataset is incentive incompatible if there is i such that,

$$\mathbb{P} \left\{ Y = y^{(i)} | \hat{w}, x^{(i)} \right\} > \mathbb{P} \left\{ Y = y^{(i)} | w^*, x^{(i)} \right\}.$$

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.