

Aggregating Data for Optimal Learning

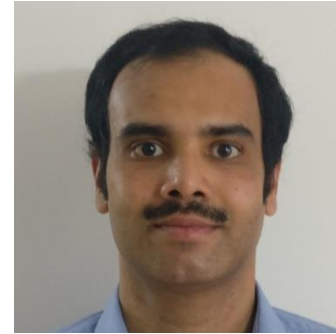


Sushant Agarwal^{*1}



Yukti Makhija²

¹Northeastern University



Rishi Saket²

²Google DeepMind



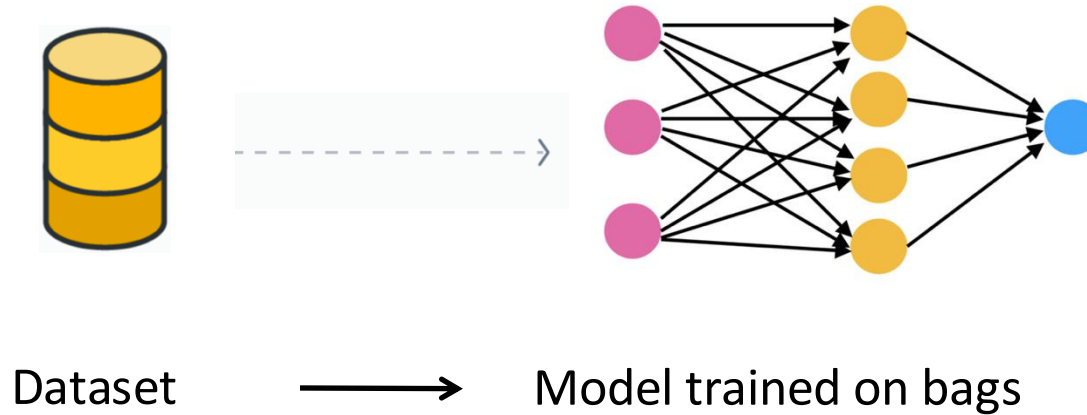
Aravindan Raghuveer²



UAI 2025

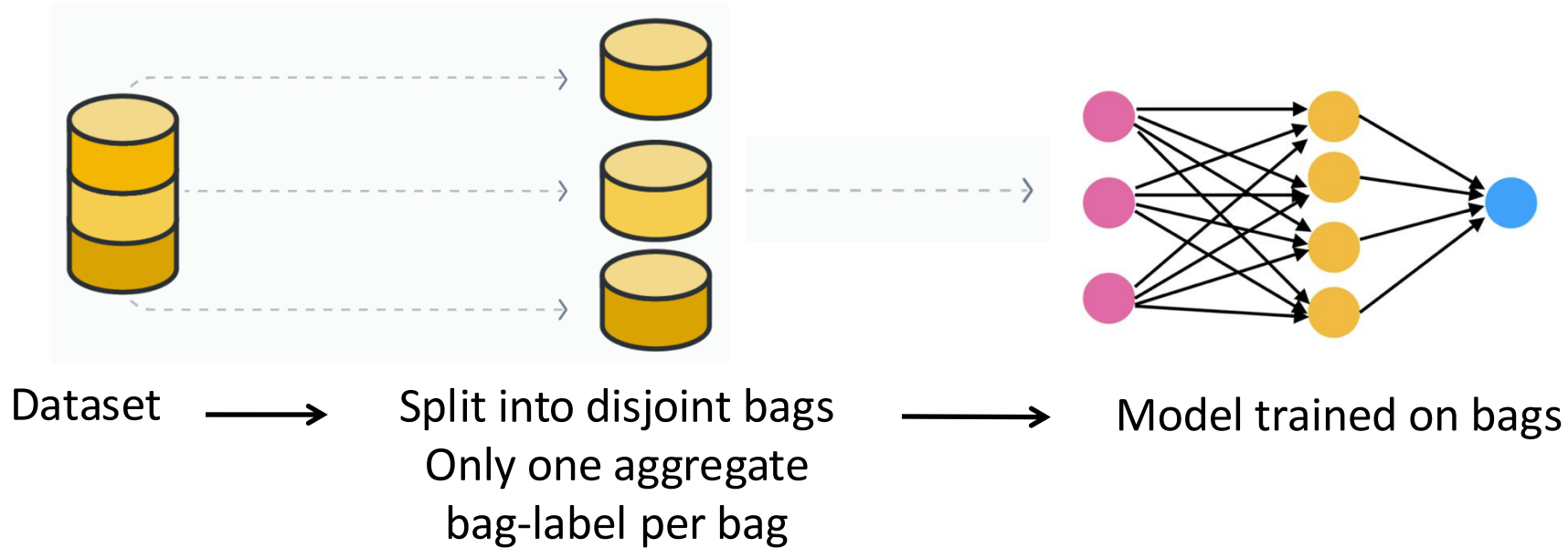
^{*}Work done during an internship at Google DeepMind

Supervised Learning



- Training dataset comprises of n tuples (x_i, y_i) .
 - x_i denotes an instance/feature-vector with label y_i .
 - Denote the sets by X, Y respectively.
- Train a model to predict labels of unseen instances.

Learning from Aggregate Labels



- X is partitioned into disjoint bags $B = \{B_1, B_2, \dots, B_k\}$.
 - Bag B_l has bag-label \bar{y}_l .
- \bar{y}_l is derived from the labels present in B_l via some aggregation function.
- Train a model to predict labels of unseen instances.

LLP and MIR

We focus on two popular paradigms.

- **Learning from Label Proportions (LLP):** \bar{y}_l is the average of the individual instance labels in B_l .
- **Multiple Instance Regression (MIR):** \bar{y}_l is the label of one (undisclosed) instance in B_l , chosen uniformly at random.

LLP and MIR formulations are becoming increasingly prevalent.

- Privacy concerns
 - If the bags are of large size, revealing only the aggregate bag-label to the learner provides privacy protection for individual labels.
- Semi-supervised learning
 - One could partition the data into bags, and query an annotator for the label of one of the instances in each bag.

Problem Statement

- In some cases, bags are fixed.
- In others, there is flexibility in curating the bags.

Main question:

Given a lower bound on bag size, what is the ***optimal bagging strategy***, to maximise utility of models trained on these bags?

The minimum bag size constraint is essential, or else the optimal bagging would be the trivial strategy of putting each point in a separate bag.

Setup

We consider the task of linear regression.

- Assume the existence of an underlying (unknown) θ^* .
 - $y_i = x_i \theta^* + \epsilon_i, \epsilon_i = N(0, \sigma^2)$.
- Given bags and bag-labels, find estimator $\hat{\theta}$ with maximum utility.
 - Utility defined in terms of closeness to θ^* .
- Train a model on bags by minimizing a given loss function.
 - Instance-level loss
 - Bag-level loss
 - Aggregate-level loss

Contributions

- **Optimal bagging:** We provide theoretical utility guarantees, and show that in each case, the optimal bagging strategy reduces to finding the optimal k -means clustering of the feature vectors or the labels.
- **Differential Privacy:** Apart from the inherent privacy that MIR and LLP offer, we can perturb the labels to obtain formal *label differential-privacy* guarantees, incurring an additional utility error.
- **GLMs:** We extend our results for Linear Regression to Generalized Linear Models (GLMs).
- **Experiments:** We experimentally validate our results on both synthetic and real-world data.

Loss Functions

- An estimator $\hat{\theta}$ minimizes **instance-level loss**, if

$$\hat{\theta} := \operatorname{argmin}_{\theta} \frac{1}{n} \sum_{l=1}^m \sum_{i \in B_l} \ell(\bar{y}_l, f_{\theta}(x_i))$$

- An estimator $\hat{\theta}$ minimizes **bag-level loss**, if

$$\hat{\theta} := \operatorname{argmin}_{\theta} \frac{1}{m} \sum_{l=1}^m \ell \left(\bar{y}_l, \frac{\sum_{i \in B_l} f_{\theta}(x_i)}{|B_l|} \right) .$$

- An estimator $\hat{\theta}$ minimizes **aggregate-level loss**, if

$$\hat{\theta} := \operatorname{argmin}_{\theta} \frac{1}{m} \sum_{l=1}^m \ell \left(\bar{y}_l, f_{\theta} \left(\frac{\sum_{i \in B_l} x_i}{|B_l|} \right) \right) .$$

Optimal Bagging

- Intuitively, a bagging provides good utility if the bags are ***homogeneous***, i.e., the feature-vectors and/or labels within a bag are similar.
- By deriving upper bounds on the error, we deduce optimal bagging strategies.

We consider two types of bagging procedures.

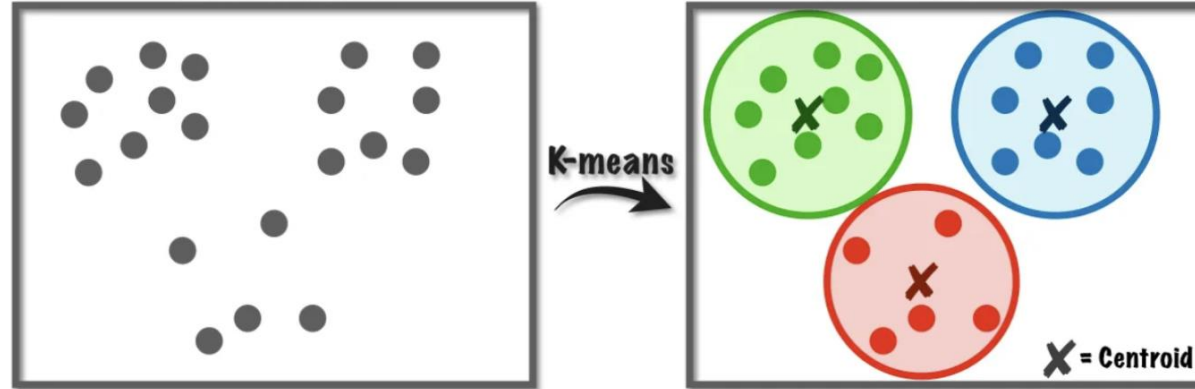
- **Label-dependent bagging:** Individual instance labels are available for the bagging.
- **Label-agnostic bagging:** Individual instance labels are ***not*** available for the bagging.

Optimal Bagging (label-dependent)

	Instance-level loss	Bag/Aggregate-level loss
LLP	1-dimensional k -means clustering of the labels (Javanmard et al. '24)	Minimize the condition number of the covariance matrix of each bag's centroid (Our work)
MIR	1-dimensional k -means clustering of the labels (Our work)	Involves both k -means clustering of labels, and minimizing the condition number (Our work)

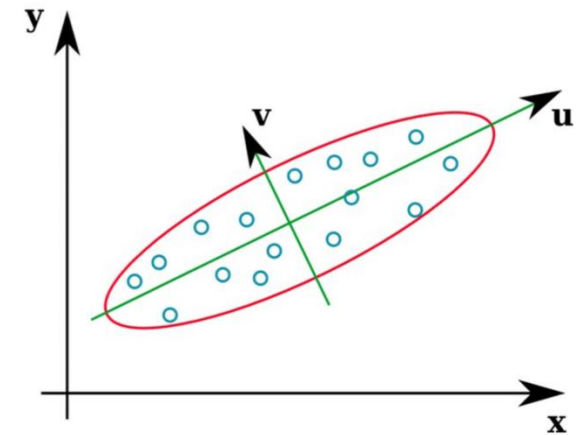
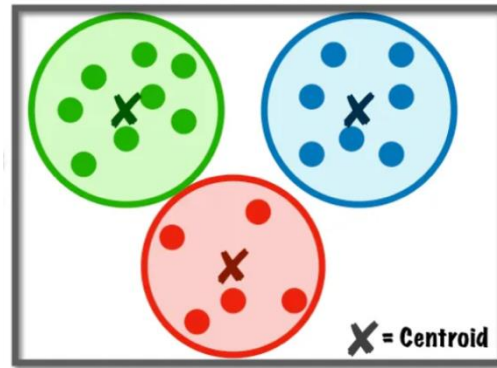
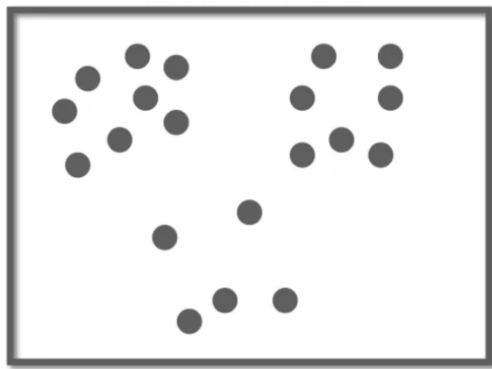
Optimal Bagging (label-dependent)

$$\arg \min_{\mathbf{S}} \sum_{i=1}^k \sum_{\mathbf{x} \in S_i} \|\mathbf{x} - \boldsymbol{\mu}_i\|^2$$



Optimal Bagging (label-dependent)

Minimize the condition number of the covariance matrix of each bag's centroid?



Optimal Bagging (label-agnostic)

- **Instance k -means:** We justify that k -means of the instances X is a good heuristic for each scenario we consider.
 - $y_i \approx x_i \theta^* \Rightarrow k$ -means of instances is a good heuristic for k -means of labels.
 - Maximizing the variance of bag-centroids along a direction \Leftrightarrow finding an optimal k -means clustering of instances projected on that direction.
- **Random bagging:** As a baseline, we also provide a utility analysis of bagging randomly.

Thanks for listening!



For more details,
check out the paper!



Personal Website
agarwal.sus@northeastern.edu