

# Multi-Observation Regression

Rafael Frongillo  
CU Boulder

Nishant Mehta  
University of Victoria

Tom Morgan  
Harvard

Bo Waggoner  
CU Boulder

## Multi-Observation Elicitation

**Standard Property Elicitation** (single-observation losses)

$$\arg \min_{r \in \mathbb{R}^d} E_{Y \sim P}[\ell(r, Y)]$$

Squared loss elicits the mean of  $P$

Variance is not elicitable by a single observation loss!

**Multi-Observation Elicitation** (multi-observation losses)

$$\arg \min_{r \in \mathbb{R}^d} E_{(Y_1, \dots, Y_m) \sim P^m}[\ell(r, Y_1, \dots, Y_m)]$$

Variance is now elicitable, using  $m = 2$

$$\ell(r, Y_1, Y_2) = \left( r - \frac{1}{2}(Y_1 + Y_2) \right)^2$$

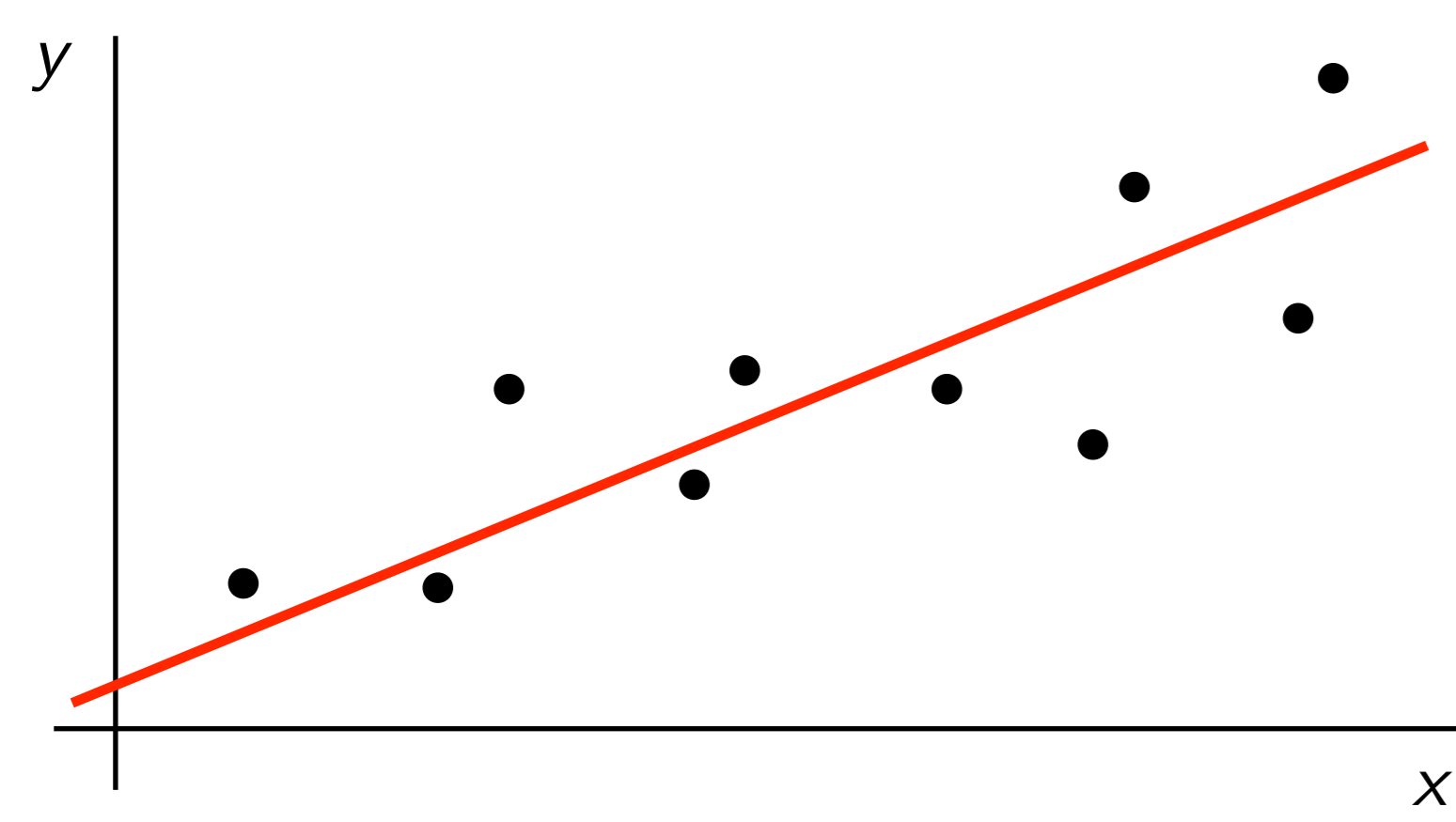
The 2-norm is also elicitable. Let  $\mathcal{Y} = \{1, 2, \dots, K\}$

$$\text{Squared 2-norm: } \sum_{j=1}^K P_j^2$$

$$\ell(r, Y_1, Y_2) = (r - 1\{Y_1 = Y_2\})^2$$

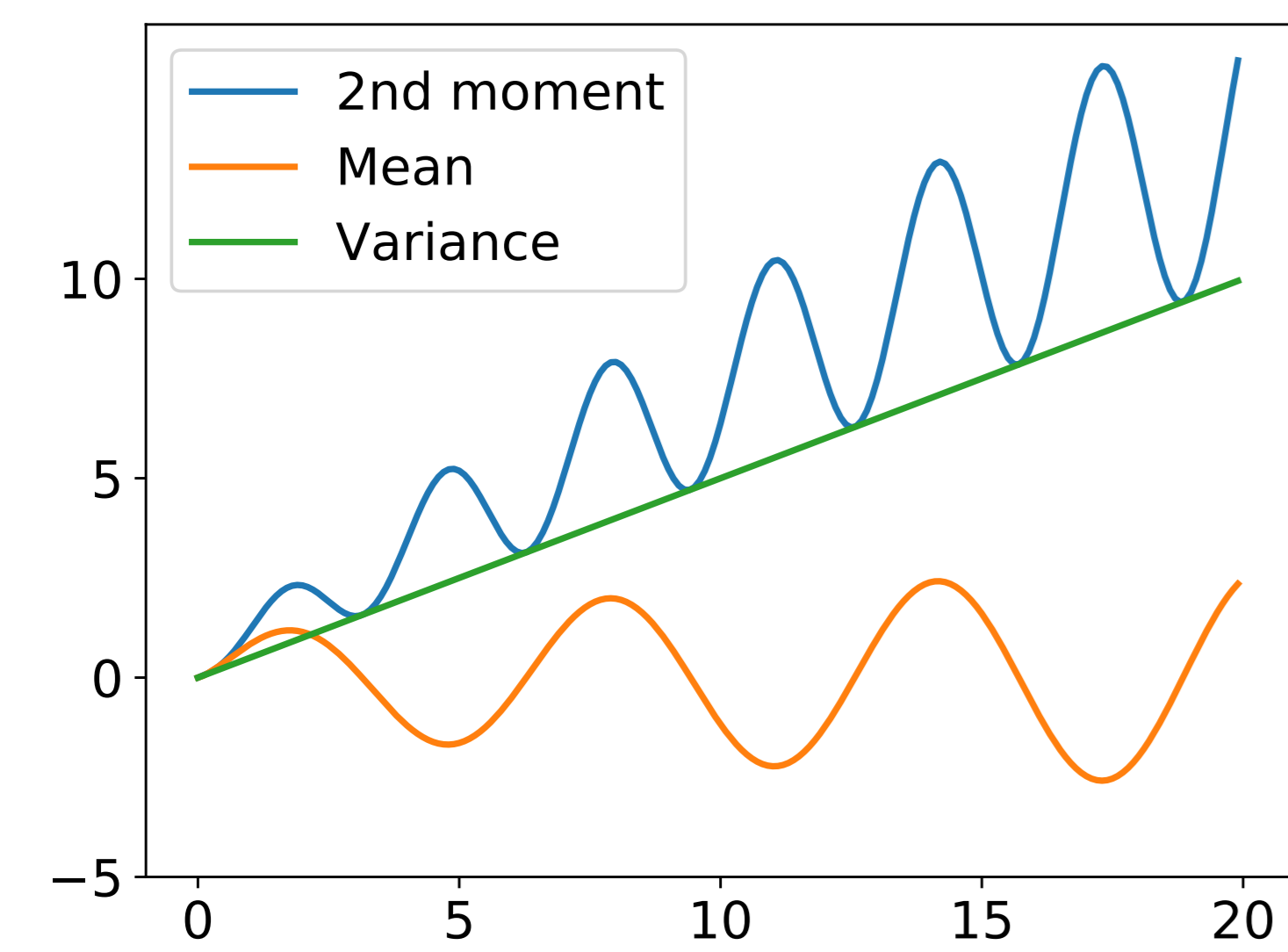
## Multi-Observation Regression

**Classical regression**



Regress on  $y$  using ERM with squared loss

How about e.g. variance or 2-norm?



Traditional approach predicts variance by regressing on  $y$  and  $y^2$  separately.

This might have high sample complexity and is suboptimal when the property varies in a simple way.

**Solution:** Regress directly on property via multi-observation loss!

So we try to minimize  $R(f) := E_{X \sim \mathcal{D}} [E_{Y \sim \mathcal{D}_X^m} [\ell_f(X, \mathbf{Y})]]$

How? Use ERM:  $\min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n \ell_f(X_i, (Y_{i,1}, \dots, Y_{i,m}))$   
i.i.d.  $\sim \mathcal{D}$  drawn from  $\mathcal{D}_X^m$

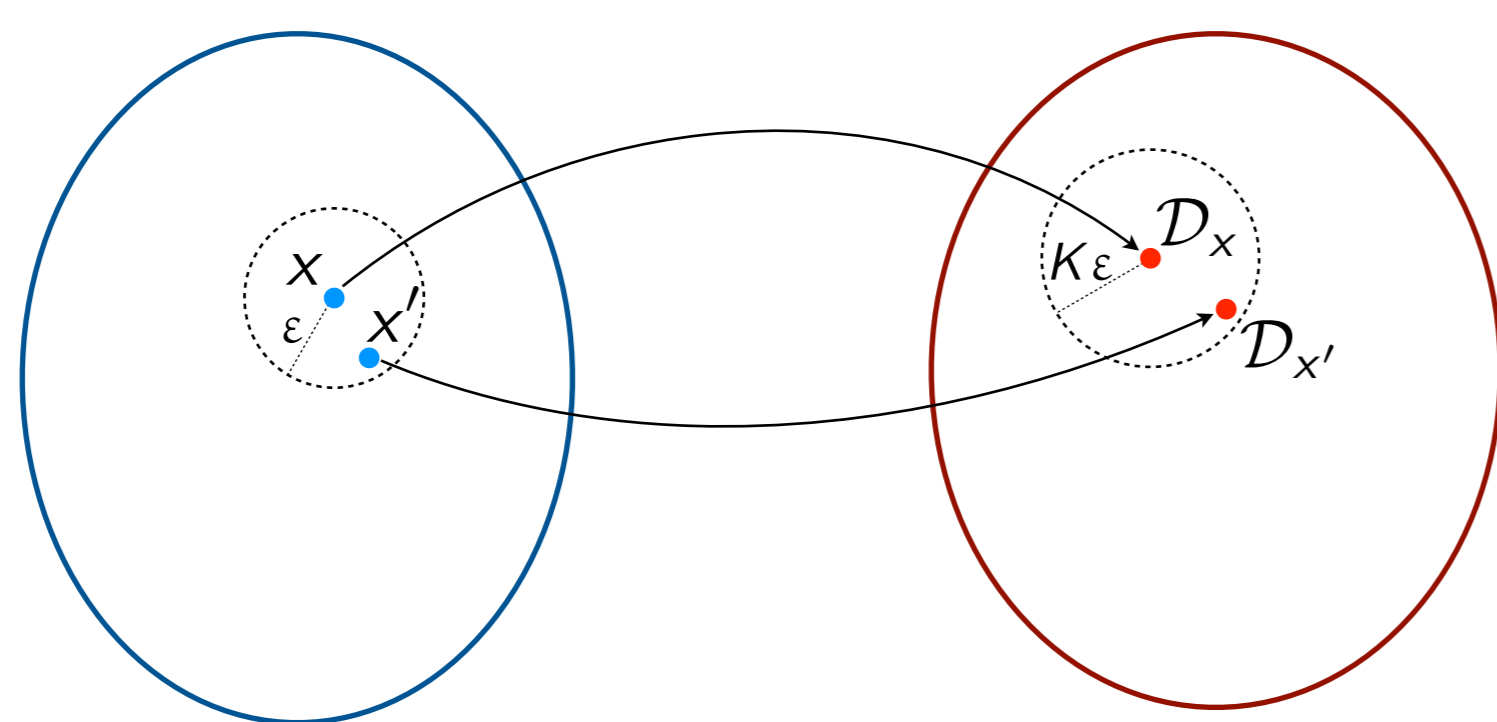
**Problem:** We don't have a multi-observation sampling oracle. So, we don't have meta-samples!

We need to make do with a classical sampling oracle. (classical, single-observation samples)

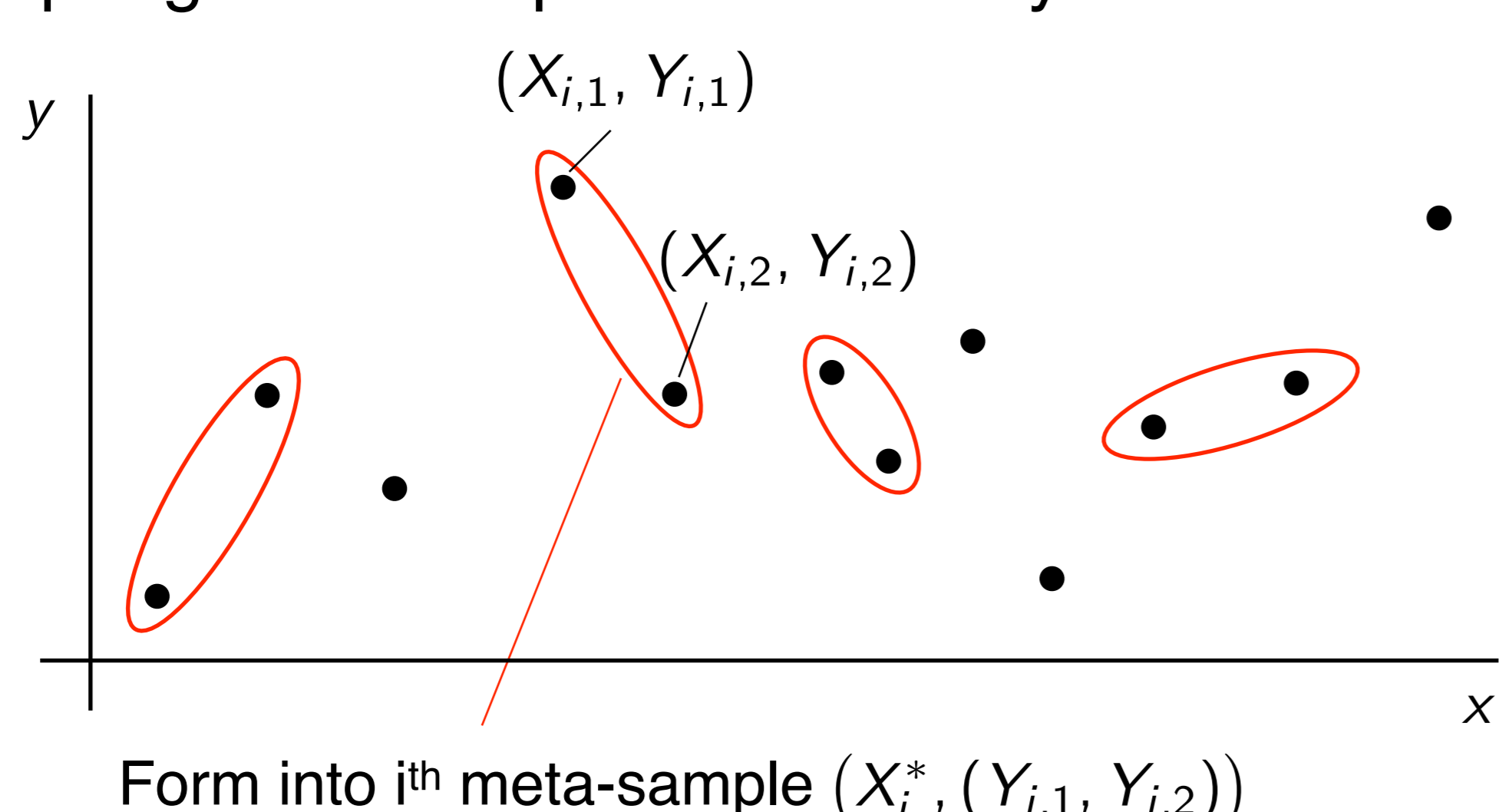
## "Meta" Algorithm for ERM with Meta-Samples

**Solution:** Form approximate meta-samples using extra samples and assume that the conditional distribution is slowly changing:

$$\exists K \text{ such that, for all } x, x' \in \mathcal{X}: \|\mathcal{D}_x - \mathcal{D}_{x'}\|_{TV} \leq K\|x - x'\|_2$$



- Collect large number of classical samples
- Clump together samples with nearby  $x$ -values



- Run ERM on meta-samples  $(X_i^*, (\tilde{Y}_{i,1}, \dots, \tilde{Y}_{i,m}))_{i \in [n]}$

Approximately drawn from  $\mathcal{D}_{X_i^*}^m$

## Algorithm for Constructing Meta-Samples

**ALGORITHM 1**

Given:  $n, m, N, \epsilon$

Sample  $n$  points  $X_1^*, \dots, X_n^*$  i.i.d. from  $\mathcal{D}$

For  $j = 1, \dots, m$

Sample  $k = \frac{N}{m}$  points  $X_1^{(j)}, \dots, X_k^{(j)}$  i.i.d. from  $\mathcal{D}$

Find a maximum matching  $M^{(j)}$  between  $X_1^*, \dots, X_n^*$  and  $X_1^{(j)}, \dots, X_k^{(j)}$ , where  $X_i^*$  and  $X_{i'}^{(j)}$  are adjacent iff  $\|X_i^* - X_{i'}^{(j)}\| \leq \epsilon$

If  $|M^{(j)}| < n$

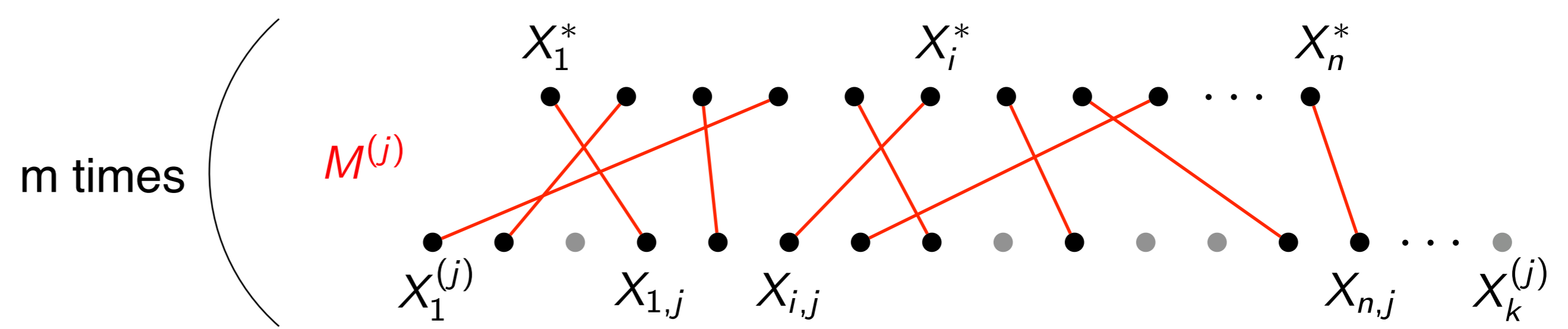
Arbitrarily match remaining  $X_i^*$ 's (ignoring distance constraints)

For  $i = 1, \dots, n$

Let  $X_{i,j}$  denote the match of  $X_i^*$  in  $M^{(j)}$

Sample a label  $\tilde{Y}_{i,j}$  from  $\mathcal{D}_{X_{i,j}}$

Return  $(X_i^*, (\tilde{Y}_{i,1}, \dots, \tilde{Y}_{i,m}))_{i \in [n]}$



**Guarantee:**

We say  $X_i^*$  is  $\epsilon$ -well-matched by the set of matchings  $M_1, \dots, M_m$  if  $X_i^*$  is matched to an  $\epsilon$ -close point in each matching.

**MATCHING LEMMA**

If  $N = \tilde{\Omega}\left(\frac{m d^{(d+2)/2} n^{(d+1)/2}}{\epsilon}\right)$ , then with probability at least  $1 - \delta$ : all but  $\tilde{O}(\sqrt{n})$  points  $X_i^*$  are  $(1/\sqrt{n})$ -well-matched by  $M_1, \dots, M_m$ .

## General case excess risk bound

**THEOREM**

Assume that the conditional distribution is slowly changing.

Let the loss be  $L$ -Lipschitz, and take  $N = \tilde{\Omega}\left(\frac{m d^{(d+2)/2} n^{(d+1)/2}}{\epsilon}\right)$ .

If Algorithm 1 is run with input  $(n, m, N, 1/\sqrt{n})$ , and if ERM is run on the resulting meta-sample, then with probability at least  $1 - \delta$ ,

$$R(\hat{f}_{\text{ERM}}) - R(f^*) \leq 2LR_n(\mathcal{F}) + 2B \left( 2\sqrt{\log \frac{4}{\delta}} + mK \right) \frac{1}{\sqrt{n}}$$

Rademacher complexity

Upper bound on loss

## Proof Sketch

Idea: Think of ERM as being run on corrupted samples

Let  $\epsilon = 1/\sqrt{n}$

- From Matching Lemma, only  $O(\sqrt{n})$  points  $X_i^*$  fail to be  $\epsilon$ -well-matched. The labels for these points are not even approximately drawn from  $\mathcal{D}_{X_i^*}$ , so we consider the corresponding  $O(\sqrt{n})$  meta-samples as corrupted.

- For each well-matched point  $X_i^*$ :

We have  $\|\mathcal{D}_{X_{i,j}} - \mathcal{D}_{X_i^*}\|_{TV} \leq K\epsilon$  for all  $j \in [m]$

Think of  $\tilde{Y}_{ij}$  as sampled as follows:

First, draw  $Z_{ij}$  from Bernoulli( $K\epsilon$ ).

Then,  $\tilde{Y}_{ij} = \text{sample from } \begin{cases} \mathcal{D}_{X_i^*} & \text{if } Z_{ij} = 0 \\ \mathcal{Q}_{ij} & \text{if } Z_{ij} = 1 \end{cases}$

Arbitrary "bad" mixture component

- With high probability, at most  $O(\sqrt{n})$  of meta-samples from well-matched points have any label coming from a bad mixture component. In all, only  $O(\sqrt{n})$  of meta-samples are corrupted.

## Simulations

**Setup:** Let  $X \sim U([0, 1])$  and  $Y = g(X) + \xi$ , for  $\xi \sim N(0, 1)$

**Goal:** Predict  $\text{Var}[Y | X]$

**Algorithms:**

"2mom linear" - fit linear functions to moments

"2mom quad" - fit quadratic functions to moments

"unbiased" - our algorithm (has theoretical guarantees)

"sliding" / "nearby" - other, non-theoretically rigorous algorithms

