

# Conditional Distribution Compression via the Kernel Conditional Mean Embedding

*Dominic Broadbent*<sup>1</sup>

*Nick Whiteley*<sup>1</sup>

*Robert Allison*<sup>1</sup>

*Tom Lovett*<sup>2</sup>



MATHEMATICAL  
INSTITUTE

# Motivation



University of  
BRISTOL



MATHEMATICAL  
INSTITUTE

- Distribution compression seeks to replace large datasets with smaller representative sets that preserve their key statistical properties, reducing the financial, environmental, and time costs of storage and computation.
- Existing methods have been developed for unlabelled data, targeting the distribution  $\mathbb{P}_X$  [1, 2, 3]. However, many real-world datasets are labelled, where preserving relationships between inputs and outputs is essential.
- Depending on the downstream task, one may wish to preserve the joint distribution  $\mathbb{P}_{X,Y}$ , which captures dependencies between features and labels, or the conditional distribution  $\mathbb{P}_{Y|X}$  which governs predictive behaviour.
- Extending distribution compression techniques beyond the marginal case naturally begins with the joint distribution, but directly targeting the conditional distribution is considerably more challenging, requiring new ways to compare and compress conditional structure.

- Distribution compression seeks to replace large datasets with smaller representative sets that preserve their key statistical properties, reducing the financial, environmental, and time costs of storage and computation.
- Existing methods have been developed for unlabelled data, targeting the distribution  $\mathbb{P}_X$  [1, 2, 3]. However, many real-world datasets are labelled, where preserving relationships between inputs and outputs is essential.
- Depending on the downstream task, one may wish to preserve the joint distribution  $\mathbb{P}_{X,Y}$ , which captures dependencies between features and labels, or the conditional distribution  $\mathbb{P}_{Y|X}$  which governs predictive behaviour.
- Extending distribution compression techniques beyond the marginal case naturally begins with the joint distribution, but directly targeting the conditional distribution is considerably more challenging, requiring new ways to compare and compress conditional structure.

- Distribution compression seeks to replace large datasets with smaller representative sets that preserve their key statistical properties, reducing the financial, environmental, and time costs of storage and computation.
- Existing methods have been developed for unlabelled data, targeting the distribution  $\mathbb{P}_X$  [1, 2, 3]. However, many real-world datasets are labelled, where preserving relationships between inputs and outputs is essential.
- Depending on the downstream task, one may wish to preserve the joint distribution  $\mathbb{P}_{X,Y}$ , which captures dependencies between features and labels, or the conditional distribution  $\mathbb{P}_{Y|X}$  which governs predictive behaviour.
- Extending distribution compression techniques beyond the marginal case naturally begins with the joint distribution, but directly targeting the conditional distribution is considerably more challenging, requiring new ways to compare and compress conditional structure.

- Distribution compression seeks to replace large datasets with smaller representative sets that preserve their key statistical properties, reducing the financial, environmental, and time costs of storage and computation.
- Existing methods have been developed for unlabelled data, targeting the distribution  $\mathbb{P}_X$  [1, 2, 3]. However, many real-world datasets are labelled, where preserving relationships between inputs and outputs is essential.
- Depending on the downstream task, one may wish to preserve the joint distribution  $\mathbb{P}_{X,Y}$ , which captures dependencies between features and labels, or the conditional distribution  $\mathbb{P}_{Y|X}$  which governs predictive behaviour.
- Extending distribution compression techniques beyond the marginal case naturally begins with the joint distribution, but directly targeting the conditional distribution is considerably more challenging, requiring new ways to compare and compress conditional structure.

- Extending herding [1] and gradient flow [2] methods, we propose Joint Kernel Herding (JKH) and Joint Kernel Inducing Points (JKIP), which produce compressed sets targeting the joint distribution.
- Extending the distributional metric, Maximum Mean Discrepancy (MMD) [4], we then introduce the Average Maximum Conditional Mean Discrepancy (AMCMD) which we show is a metric on the space of conditional distributions, and derive a closed-form estimate.
- Targeting the AMCMD, we propose Average Conditional Kernel Herding (ACKH) and Average Conditional Kernel Inducing Points (ACKIP), which compress the conditional distribution in linear time.
- Across various downstream tasks, we show that directly compressing the conditional distribution is preferable to joint distribution compression. Moreover, the greedy herding methods (JKH, ACKH) are outperformed by gradient flow methods (JKIP, ACKIP).



- Extending herding [1] and gradient flow [2] methods, we propose Joint Kernel Herding (JKH) and Joint Kernel Inducing Points (JKIP), which produce compressed sets targeting the joint distribution.
- Extending the distributional metric, Maximum Mean Discrepancy (MMD) [4], we then introduce the Average Maximum Conditional Mean Discrepancy (AMCMD) which we show is a metric on the space of conditional distributions, and derive a closed-form estimate.
- Targeting the AMCMD, we propose Average Conditional Kernel Herding (ACKH) and Average Conditional Kernel Inducing Points (ACKIP), which compress the conditional distribution in linear time.
- Across various downstream tasks, we show that directly compressing the conditional distribution is preferable to joint distribution compression. Moreover, the greedy herding methods (JKH, ACKH) are outperformed by gradient flow methods (JKIP, ACKIP).

- Extending herding [1] and gradient flow [2] methods, we propose Joint Kernel Herding (JKH) and Joint Kernel Inducing Points (JKIP), which produce compressed sets targeting the joint distribution.
- Extending the distributional metric, Maximum Mean Discrepancy (MMD) [4], we then introduce the Average Maximum Conditional Mean Discrepancy (AMCMD) which we show is a metric on the space of conditional distributions, and derive a closed-form estimate.
- Targeting the AMCMD, we propose Average Conditional Kernel Herding (ACKH) and Average Conditional Kernel Inducing Points (ACKIP), which compress the conditional distribution in linear time.
- Across various downstream tasks, we show that directly compressing the conditional distribution is preferable to joint distribution compression. Moreover, the greedy herding methods (JKH, ACKH) are outperformed by gradient flow methods (JKIP, ACKIP).



- Extending herding [1] and gradient flow [2] methods, we propose Joint Kernel Herding (JKH) and Joint Kernel Inducing Points (JKIP), which produce compressed sets targeting the joint distribution.
- Extending the distributional metric, Maximum Mean Discrepancy (MMD) [4], we then introduce the Average Maximum Conditional Mean Discrepancy (AMCMD) which we show is a metric on the space of conditional distributions, and derive a closed-form estimate.
- Targeting the AMCMD, we propose Average Conditional Kernel Herding (ACKH) and Average Conditional Kernel Inducing Points (ACKIP), which compress the conditional distribution in linear time.
- Across various downstream tasks, we show that directly compressing the conditional distribution is preferable to joint distribution compression. Moreover, the greedy herding methods (JKH, ACKH) are outperformed by gradient flow methods (JKIP, ACKIP).

- Distribution compression algorithms optimise the compressed set  $\mathcal{C} = \{\mathbf{z}_i\}_{i=1}^m$  to minimise the MMD to the empirical distribution  $\hat{\mathbb{P}}_X$  of the target dataset  $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^n$ :

$$\begin{aligned} \text{MMD}^2(\hat{\mathbb{P}}_X, \hat{\mathbb{P}}_Z) &:= \|\hat{\mu}_X - \hat{\mu}_Z\|_{\mathcal{H}_k}^2 \\ &= \sum_{i,j=1}^n k(\mathbf{x}_i, \mathbf{x}_j) - 2 \sum_{i,j=1}^{n,m} k(\mathbf{x}_i, \mathbf{z}_j) + \sum_{i,j=1}^m k(\mathbf{z}_i, \mathbf{z}_j), \end{aligned}$$

where  $m \ll n$ , and we denote  $\mu_X$  as the *kernel mean embedding* of the distribution  $\mathbb{P}_X$ . The KME  $\mu_X$  lies in the *Reproducing Kernel Hilbert Space* (RKHS)  $\mathcal{H}_k$  induced by the positive definite kernel  $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ , which is defined on the feature space  $\mathcal{X}$ .

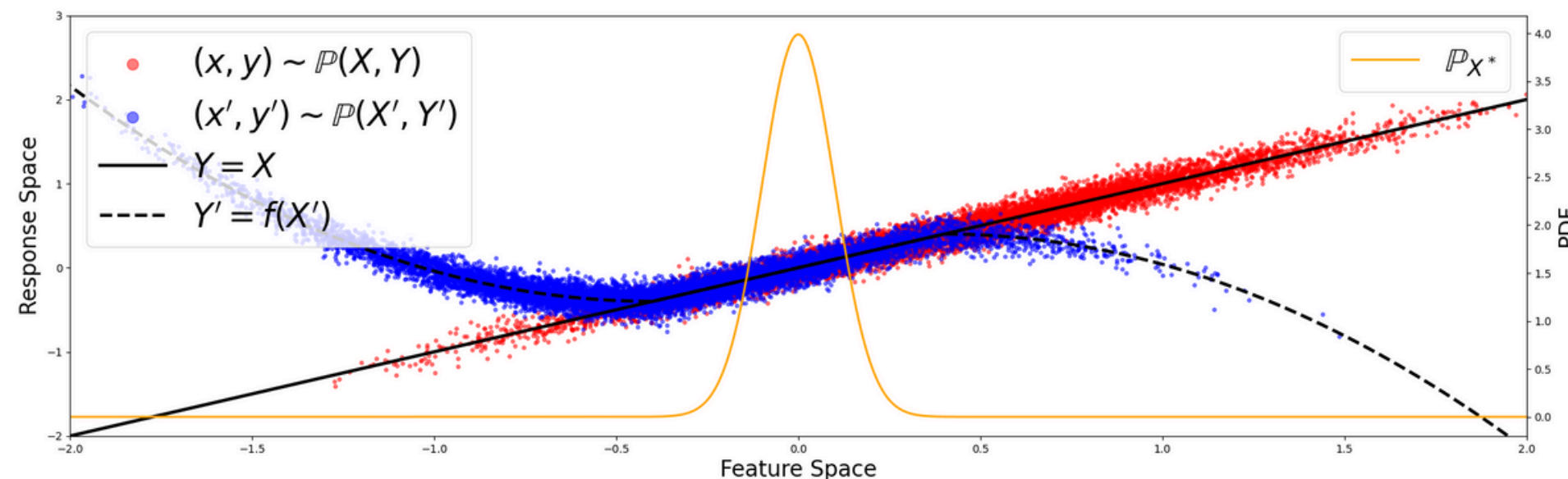
- Given an additional kernel  $l : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$  defined on the response space  $\mathcal{Y}$  we induce the RKHS  $\mathcal{H}_k \otimes \mathcal{H}_l$ . We can then extend existing distribution compression algorithms to optimise a compressed set  $\mathcal{C} = \{(\mathbf{z}_i, \mathbf{w}_i)\}_{i=1}^m$  which minimises the Joint MMD [5] to the empirical distribution of the target dataset  $\mathcal{D} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$ :

$$\begin{aligned} \text{JMMD}^2(\hat{\mathbb{P}}_{X,Y}, \hat{\mathbb{P}}_{Z,W}) &:= \|\hat{\mu}_{X,Y} - \hat{\mu}_{Z,W}\|_{\mathcal{H}_{k \otimes l}}^2 \\ &= \sum_{i,j=1}^n k(\mathbf{x}_i, \mathbf{x}_j) l(\mathbf{y}_i, \mathbf{y}_j) - 2 \sum_{i,j=1}^{n,m} k(\mathbf{x}_i, \mathbf{z}_j) l(\mathbf{y}_i, \mathbf{w}_j) + \sum_{i,j=1}^m k(\mathbf{z}_i, \mathbf{z}_j) l(\mathbf{w}_i, \mathbf{w}_j). \end{aligned}$$

- In order to extend distribution compression to the conditional distribution, we first require a notion of conditional discrepancy, for this we introduce the AMCMD:

$$\text{AMCMD} \left( \mathbb{P}_{X^*}, \mathbb{P}_{Y|X}, \mathbb{P}_{Y'|X'} \right) := \sqrt{\mathbb{E}_{\mathbf{x} \sim \mathbb{P}_{X^*}} \left[ \left\| \mu_{Y|X=\mathbf{x}} - \mu_{Y'|X'=\mathbf{x}} \right\|_{\mathcal{H}_l}^2 \right]}$$

where  $\mathbb{P}_{X^*}$  is a weighting distribution, and  $\mu_{Y|X} : \mathcal{X} \rightarrow \mathcal{H}_l$  is the *kernel conditional mean embedding* (KCME). The KCME is a *vector-valued* function, which takes as inputs conditioning values  $\mathbf{x} \in \mathcal{X}$ , and outputs KMEs  $\mu_{Y|X=\mathbf{x}}$  lying in  $\mathcal{H}_l$ .



**Theorem** - *The AMCMD is a proper metric*

Suppose the response kernel  $l(\cdot, \cdot)$  is characteristic, that  $\mathbb{P}_X$ ,  $\mathbb{P}_{X'}$ , and  $\mathbb{P}_{X^*}$  are absolutely continuous with respect to each other, and that  $\mathbb{P}(\cdot | X)$  and  $\mathbb{P}(\cdot | X')$  admit regular versions. Then,  $\text{AMCMD}(\mathbb{P}_{X^*}, \mathbb{P}_{Y|X}, \mathbb{P}_{Y'|X'}) = 0$  if and only if, for almost all  $\mathbf{x} \in \mathcal{X}$  wrt  $\mathbb{P}_{X^*}$ ,  $\mathbb{P}_{Y|X=\mathbf{x}}(A) = \mathbb{P}_{Y'|X'}(A)$  for all  $A \in \mathcal{Y}$ .

Moreover, assuming the Radon-Nikodym derivatives  $\frac{d\mathbb{P}_{X^*}}{d\mathbb{P}_X}$ ,  $\frac{d\mathbb{P}_{X'}}{d\mathbb{P}_X}$ , and  $\frac{d\mathbb{P}_{X^*}}{d\mathbb{P}_{X'}}$  are bounded, then the triangle inequality is satisfied, i.e.

$$\text{AMCMD}(\mathbb{P}_{Y|X}, \mathbb{P}_{Y''|X''}) \leq \text{AMCMD}(\mathbb{P}_{Y|X}, \mathbb{P}_{Y'|X'}) + \text{AMCMD}(\mathbb{P}_{Y'|X'}, \mathbb{P}_{Y''|X''}).$$



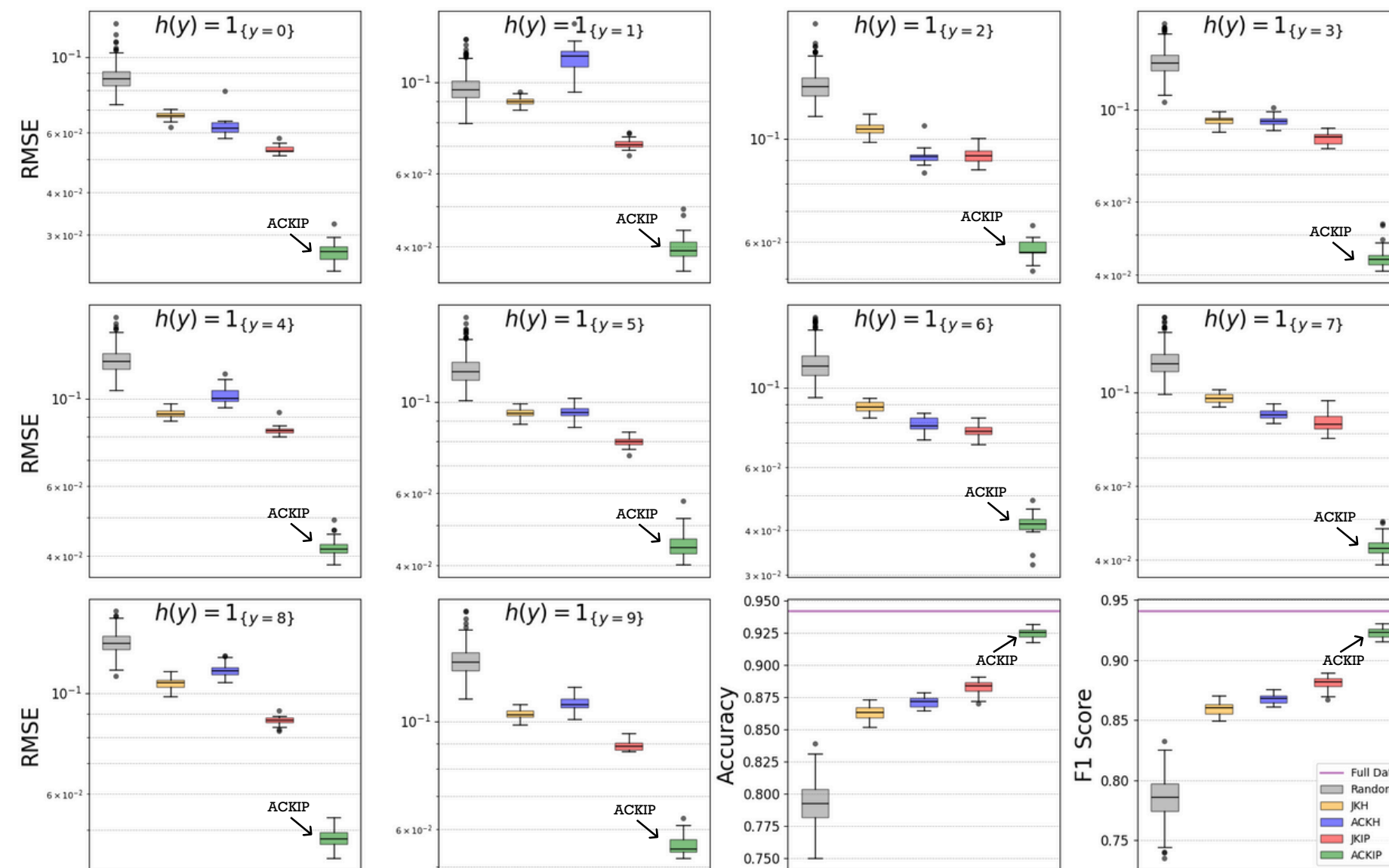
- We can now optimise a compressed set  $\mathcal{C} = \{(z_i, w_i)\}_{i=1}^m$  which minimises the AMCMD to the empirical conditional distribution of the target dataset  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$  :

$$\text{AMCMD}^2 \left( \hat{\mathbb{P}}_{X^*}, \hat{\mathbb{P}}_{Y|X}, \hat{\mathbb{P}}_{Z|W} \right) = \frac{1}{q} \sum_{i=1}^q \left\| \hat{\mu}_{Y|X=x_i^*} - \hat{\mu}_{Z|W=x_i^*} \right\|_{\mathcal{H}_l}^2.$$

- We can obtain a closed-form representation of this, however it has  $\mathcal{O}(n^3)$  cost. For distribution compression, it is natural to choose  $\mathbb{P}_{X^*} = \mathbb{P}_X$ , then by applying the tower property, we can reduce to  $\mathcal{O}(n)$  cost, enabling linear-time conditional distribution compression.



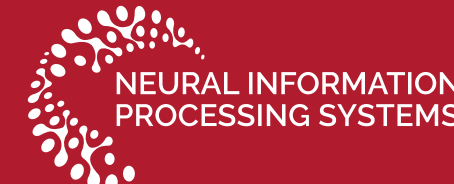
- The KCME has many important applications. In particular it may be used as a regressor and classifier. In our work, we investigate how compression effects these downstream tasks. Below, we show results on MNIST after 98% compression:



# References



University of  
BRISTOL



MATHEMATICAL  
INSTITUTE

- [1] - Yutian Chen, Max Welling, and Alex Smola. Super-samples from kernel herding. UAI 2010
- [2] - Michael Arbel, Anna Korba, Adil Salim, and Arthur Gretton. Maximum mean discrepancy gradient flow, NeurIPS 2019
- [3] - Raaz Dwivedi and Lester Mackey. Kernel thinning. COLT 2021
- [4] - Arthur Gretton, Karsten M. Borgwardt, Malte J. Rasch, Bernhard Schölkopf, and Alexander Smola. A kernel two-sample test, JMLR 2012
- [5] - Mingsheng Long, Han Zhu, Jianmin Wang, and Michael I. Jordan. Deep transfer learning with joint adaptation networks, ICML 2017