

Mathematical Exposition: Conditional Moment Restrictions Library

1. Problem Formulation

1.1 Conditional Moment Restrictions

This library addresses estimation problems defined by **conditional moment restrictions (CMR)** of the form:

$$\mathbb{E}[\psi(X, Y; \theta) | Z] = 0 \quad P_Z\text{-almost surely}$$

where:

- (X, Y, Z) are random variables observed from an i.i.d. sample $\{(x_i, y_i, z_i)\}_{i=1}^n$
- $X \in \mathcal{X} \subseteq \mathbb{R}^{d_x}$ represents treatments or covariates
- $Y \in \mathcal{Y} \subseteq \mathbb{R}^{d_y}$ represents outcomes or responses
- $Z \in \mathcal{Z} \subseteq \mathbb{R}^{d_z}$ represents instrumental variables
- $\theta \in \Theta$ are the parameters of interest
- $\psi : \mathcal{X} \times \mathcal{Y} \times \Theta \rightarrow \mathbb{R}^k$ is a known moment function

1.2 Instrumental Variable Regression

The canonical application is **instrumental variable (IV) regression**, where: - We seek to estimate a structural function $f_\theta : \mathcal{X} \rightarrow \mathcal{Y}$ - The moment function takes the form: $\psi(X, Y; \theta) = f_\theta(X) - Y$ - The CMR becomes: $\mathbb{E}[f_\theta(X) - Y | Z] = 0$

This problem arises when:

- There exists unobserved confounding U between X and Y
- An instrument Z satisfies:
 1. **Relevance:** Z is correlated with X
 2. **Exclusion:** Z affects Y only through X
 3. **Exogeneity:** Z is independent of U

The causal structure is: $Z \rightarrow X \leftarrow U \rightarrow Y \leftarrow X$

1.3 Unconditional Moment Restrictions

When Z is not specified or the conditional restriction reduces to an unconditional one, we have the **moment restriction (MR)** problem:

$$\mathbb{E}[\psi(X, Y; \theta)] = 0$$

This is a special case where Z is trivial or integrated out.

2. Mathematical Framework

2.1 Dual Formulation via f-Divergences

Many estimators in this library are based on **generalized empirical likelihood (GEL)** formulations using f-divergences.

2.1.1 f-Divergences For probability measures P and Q with $P \ll Q$, the f-divergence is:

$$D_f(P\|Q) = \int f\left(\frac{dP}{dQ}\right) dQ$$

where $f : \mathbb{R}_+ \rightarrow \mathbb{R}$ is convex with $f(1) = 0$.

The **Fenchel conjugate** of f is:

$$f^*(t) = \sup_{u \in \mathbb{R}_+} \{tu - f(u)\}$$

Common f-divergences:

- **KL divergence:** $f(u) = u \log u$, so $f^*(t) = e^t$
- **Chi-squared:** $f(u) = (u-1)^2$, so $f^*(t) = \frac{1}{2}(t+1)^2$
- **Log divergence:** $f(u) = -\log u$, so $f^*(t) = -\log(1-t)$ for $t < 1$

2.1.2 GEL Estimator (Unconditional) For unconditional MR, the GEL estimator solves:

$$\min_{\theta \in \Theta} \sup_{\lambda \in \mathbb{R}^k} \left\{ \frac{1}{n} \sum_{i=1}^n f^*(\lambda^\top \psi(x_i, y_i; \theta)) \right\}$$

where λ are dual variables (Lagrange multipliers). This generalizes:

- **GMM:** $f^*(t) = -\frac{1}{2}t^2$ (quadratic)
- **Empirical Likelihood:** $f^*(t) = -\log(1-t)$
- **Exponential Tilting (ET):** $f^*(t) = e^t$
- **Continuous Updating GMM (CUE):** $f^*(t) = \frac{1}{2}(t+1)^2$

2.2 Reproducing Kernel Hilbert Spaces (RKHS)

For conditional moment restrictions, we parameterize the dual functions as elements of an RKHS.

2.2.1 RKHS Basics Let \mathcal{H} be an RKHS of functions $h : \mathcal{Z} \rightarrow \mathbb{R}$ with kernel $k : \mathcal{Z} \times \mathcal{Z} \rightarrow \mathbb{R}$.

Reproducing property: For any $h \in \mathcal{H}$ and $z \in \mathcal{Z}$:

$$h(z) = \langle h, k(z, \cdot) \rangle_{\mathcal{H}}$$

Representer theorem: For finite sample problems, the solution has the form:

$$h^*(z) = \sum_{i=1}^n \alpha_i k(z, z_i)$$

for some coefficients $\alpha_i \in \mathbb{R}$.

2.2.2 Vector-Valued RKHS For k -dimensional moment functions, we use the vector-valued RKHS:

$$\mathcal{H}^k = \underbrace{\mathcal{H} \times \cdots \times \mathcal{H}}_{k \text{ times}}$$

with norm $\|\mathbf{h}\|_{\mathcal{H}^k}^2 = \sum_{j=1}^k \|h_j\|_{\mathcal{H}}^2$ for $\mathbf{h} = (h_1, \dots, h_k)$.

2.2.3 RBF Kernel The library uses the **radial basis function (RBF) kernel**:

$$k(z, z') = \exp\left(-\frac{\|z - z'\|^2}{2\sigma^2}\right)$$

where $\sigma > 0$ is the bandwidth (automatically selected via median heuristic).

2.3 Maximum Mean Discrepancy (MMD)

The **MMD** between distributions P and Q with respect to kernel k is:

$$\text{MMD}^2(P, Q) = \mathbb{E}_{p, p'}[k(p, p')] - 2\mathbb{E}_{p, q}[k(p, q)] + \mathbb{E}_{q, q'}[k(q, q')]$$

Empirically, with samples $\{p_i\}_{i=1}^n \sim P$ and $\{q_j\}_{j=1}^m \sim Q$:

$$\widehat{\text{MMD}}^2 = \frac{1}{n^2} \sum_{i, i'} K_{ii'} - \frac{2}{nm} \sum_{i, j} K_{ij}^{pq} + \frac{1}{m^2} \sum_{j, j'} K_{jj'}^q$$

2.4 Random Fourier Features (RFF)

For scalability, kernels are approximated using **random Fourier features**.

For the RBF kernel, draw $\omega_1, \dots, \omega_D \sim \mathcal{N}(0, \sigma^{-2}I)$ and define:

$$\phi(z) = \sqrt{\frac{2}{D}} [\cos(\omega_1^\top z), \sin(\omega_1^\top z), \dots, \cos(\omega_D^\top z), \sin(\omega_D^\top z)]^\top$$

Then $k(z, z') \approx \phi(z)^\top \phi(z')$ and the RKHS function becomes:

$$h(z) = \beta^\top \phi(z)$$

for $\beta \in \mathbb{R}^{2D}$, reducing the problem from $O(n^2)$ to $O(nD)$ complexity.

3. Estimation Methods

3.1 Unconditional Moment Restrictions

3.1.1 Ordinary Least Squares (OLS) Objective:

$$\hat{\theta}_{\text{OLS}} = \arg \min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \|\psi(x_i, y_i; \theta)\|^2$$

Use case: Simple baseline, assumes no confounding.

3.1.2 Generalized Method of Moments (GMM) Objective:

$$\hat{\theta}_{\text{GMM}} = \arg \min_{\theta \in \Theta} \bar{\psi}_n(\theta)^\top W_n \bar{\psi}_n(\theta)$$

where:

- $\bar{\psi}_n(\theta) = \frac{1}{n} \sum_{i=1}^n \psi(x_i, y_i; \theta)$
- W_n is a weighting matrix (typically $\hat{\Sigma}^{-1}$ where $\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n \psi_i \psi_i^\top$)

Two-step procedure:

1. Obtain initial estimate $\tilde{\theta}$
2. Compute W_n using $\tilde{\theta}$
3. Re-optimize with the optimal weighting matrix

Parameters: `reg_param` for regularization $W_n = (\hat{\Sigma} + \alpha I)^{-1}$

3.1.3 Generalized Empirical Likelihood (GEL) Primal-dual formulation:

$$\min_{\theta \in \Theta} \max_{\lambda \in \mathbb{R}^k} \left\{ -\frac{1}{n} \sum_{i=1}^n f^*(\lambda^\top \psi_i(\theta)) - \alpha \|\lambda\|^2 \right\}$$

Variants (via `divergence` parameter):

- 'chi2': Chi-squared divergence, $f^*(t) = \frac{1}{2}(t+1)^2$
- 'kl': KL divergence, $f^*(t) = e^t$
- 'log': Log divergence, $f^*(t) = -\log(1-t)$ with constraint $t < 1$

Parameters:

- `divergence`: Choice of f-divergence
- `reg_param`: Regularization parameter $\alpha \geq 0$

3.1.4 Kernel Method of Moments (KMM) Formulation:

Minimizes an MMD-based objective with entropy regularization.

$$\min_{\theta} \max_{\nu, \eta, h} \left\{ \mathbb{E}_P[h(X, Y, Z)] + c - \gamma \mathbb{E}_Q[f^*(\frac{1}{\gamma}(h(X, Y, Z) + c - \nu^\top \psi(\theta)))] - \frac{\lambda}{2} \|h\|_{\mathcal{H}}^2 \right\}$$

where:

- P is the empirical distribution
- Q is a reference distribution (estimated via KDE)
- $h \in \mathcal{H}$ is an RKHS function on the joint space
- $\nu \in \mathbb{R}^k$ are dual variables for moments
- $c \in \mathbb{R}$ is a normalization constant
- $\gamma > 0$ is the entropy regularization parameter
- $\lambda > 0$ is the RKHS regularization parameter

Implementation details:

- Uses random Fourier features for h (dimension controlled by `n_random_features`)
- Reference distribution via KDE with bandwidth `kde_bandwidth`
- Samples from reference distribution: `n_reference_samples`
- Product kernel on (X, Y, Z) space

Parameters:

- `entropy_reg_param`: γ (main regularization)
- `rkhs_reg_param`: λ for RKHS norm
- `n_random_features`: Dimension D of RFF approximation
- `n_reference_samples`: Number of samples from reference distribution
- `kde_bandwidth`: Bandwidth for KDE

3.2 Conditional Moment Restrictions

3.2.1 Maximum Moment Restrictions (MMR) **Objective:** Directly minimizes the MMD between $\psi(\theta)$ and 0 conditional on Z .

$$\hat{\theta}_{\text{MMR}} = \arg \min_{\theta \in \Theta} \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n \psi(x_i, y_i; \theta)^\top K_{ij}^Z \psi(x_j, y_j; \theta)$$

where $K_{ij}^Z = k(z_i, z_j)$ is the kernel Gram matrix for instruments.

Advantages:

- Simple, no dual optimization
- Convex in θ for linear models
- Direct interpretation as MMD

Limitations:

- Can be sensitive to kernel choice
- No regularization on moment function

3.2.2 Sieve Minimum Distance (SMD) **Formulation:** Uses sieve basis (polynomials or splines) to approximate the conditional expectation.

Let $\{b_j(z)\}_{j=1}^J$ be a basis for $L^2(P_Z)$. Approximate:

$$\mathbb{E}[\psi(X, Y; \theta) | Z = z] \approx \sum_{j=1}^J \beta_j b_j(z)$$

Objective:

$$\min_{\theta, \beta} \frac{1}{n} \sum_{i=1}^n \left\| \psi(x_i, y_i; \theta) - \sum_{j=1}^J \beta_j b_j(z_i) \right\|^2$$

Properties:

- Classical semiparametric approach
- Theoretical guarantees under regularity conditions
- Basis dimension $J = J_n \rightarrow \infty$ as $n \rightarrow \infty$ (slower than n)

3.2.3 Variational Method of Moments (VMM) Kernel Version (VMM-kernel):

Solves the saddle-point problem:

$$\min_{\theta} \max_{\mathbf{h} \in \mathcal{H}^k} \left\{ \frac{2}{n} \sum_{i=1}^n \psi(x_i, y_i; \theta)^\top \mathbf{h}(z_i) - \mathbf{h}^\top M \mathbf{h} \right\}$$

where M is a block-diagonal matrix with blocks $M_j = K^Z(K^Z K^Z + \alpha I)^{-1} K^Z$ for $j = 1, \dots, k$.

Neural Version (VMM-neural):

Parameterizes $\mathbf{h}(z) = g_\phi(z)$ where g_ϕ is a neural network, leading to:

$$\min_{\theta} \max_{\phi} \left\{ \frac{2}{n} \sum_{i=1}^n \psi(x_i, y_i; \theta)^\top g_\phi(z_i) - \frac{1}{n^2} \sum_{i,j} g_\phi(z_i)^\top K_{ij}^Z g_\phi(z_j) - \alpha \|g_\phi\|^2 \right\}$$

Parameters:

- `reg_param`: Regularization parameter α
- `num_iter`: Number of alternating minimization iterations

Relation to DeepGMM: VMM-neural with specific settings is equivalent to DeepGMM.

3.2.4 Functional Generalized Empirical Likelihood (FGEL) Kernel Version (FGEL-kernel):

Extends GEL to conditional setting using RKHS dual functions:

$$\min_{\theta} \max_{\mathbf{h} \in \mathcal{H}^k} \left\{ -\frac{1}{n} \sum_{i=1}^n f^*(\mathbf{h}(z_i)^\top \psi(x_i, y_i; \theta)) - \frac{\alpha}{2} \|\mathbf{h}\|_{\mathcal{H}^k}^2 \right\}$$

By representer theorem, $\mathbf{h}(z) = \sum_{i=1}^n A_i k(z, z_i)$ for $A_i \in \mathbb{R}^k$.

Objective in matrix form:

$$\max_{A \in \mathbb{R}^{n \times k}} \left\{ \frac{1}{n} \sum_{i=1}^n f^* \left(\sum_{j=1}^n K_{ij}^Z A_j^\top \psi(x_i, y_i; \theta) \right) - \frac{\alpha}{2} \text{tr}(A^\top K^Z A) \right\}$$

Neural Version (FGEL-neural):

Replaces RKHS function with neural network $g_\phi : \mathcal{Z} \rightarrow \mathbb{R}^k$:

$$\min_{\theta} \max_{\phi} \left\{ -\frac{1}{n} \sum_{i=1}^n f^*(g_\phi(z_i)^\top \psi(x_i, y_i; \theta)) - \alpha \|g_\phi\|^2 \right\}$$

where $\|g_\phi\|^2 = \frac{1}{n} \sum_{i=1}^n \|g_\phi(z_i)\|^2$ is the empirical L2 norm.

Parameters:

- `divergence`: Type of f-divergence ('chi2', 'kl', 'log')
- `reg_param`: RKHS/neural network regularization α
- `dual_func_network_kwargs`: Architecture for neural dual function (neural version)

3.2.5 Kernel Method of Moments - Neural (KMM-neural) Combines MMD-based objective with neural network instrument function:

$$\min_{\theta} \max_{\nu, c, h} \left\{ \frac{1}{n} \sum_{i=1}^n h_\beta(\phi(x_i, y_i, z_i)) + c - \gamma \mathbb{E}_Q \left[f^* \left(\frac{1}{\gamma} (h + c - \nu_\phi(z)^\top \psi(\theta)) \right) \right] - \frac{\lambda}{2} \|\beta\|^2 \right\}$$

where: - $h_\beta(u) = \beta^\top \phi(x, y, z)$ with ϕ being RFF of product kernel - $\nu_\phi(z)$ is a neural network mapping $z \mapsto \mathbb{R}^k$ - Reference samples are drawn from KDE on joint distribution

Key differences from FGEL-neural:

- Uses MMD/entropy regularization instead of pure f-divergence
- Employs reference distribution for better sample efficiency

- Has RKHS function on full (X, Y, Z) space plus neural function on Z

Parameters:

- `entropy_reg_param`: Entropy regularization γ
 - `reg_param`: Neural network regularization
 - `rkhs_reg_param`: RKHS regularization λ
 - `n_random_features`: RFF dimension for h
 - `n_reference_samples`: Samples from reference distribution
 - `dual_func_network_kwargs`: Neural network architecture for ν_ϕ
-

4. API Reference

4.1 High-Level Interface

`estimation()` Main entry point for training any estimator with automatic hyperparameter search.

```
from cmr import estimation

trained_model, stats = estimation(
    model,                      # torch.nn.Module
    train_data,                  # dict: {'t': ndarray, 'y': ndarray, 'z': ndarray}
    moment_function,             # callable: (model_output, y) -> torch.Tensor
    estimation_method,           # str: method name
    estimator_kwargs=None,       # dict: method-specific parameters
    hyperparams=None,            # dict: hyperparameter search space
    sweep_hparams=True,          # bool: whether to search hyperparameters
    validation_data=None,         # dict: validation set
    val_loss_func=None,           # callable or str: validation metric
    normalize_moment_function=True, # bool: normalize moments
    verbose=True                 # bool: print progress
)
```

Arguments:

Parameter	Type	Description
<code>model</code>	<code>torch.nn.Module</code>	Parametric model f_θ containing parameters to estimate
<code>train_data</code>	<code>dict</code>	Training data with keys ' <code>t</code> ' (treatments), ' <code>y</code> ' (outcomes), ' <code>z</code> ' (instruments)
<code>moment_function</code>	<code>callable</code>	Function $(f_\theta(x), y) \mapsto \psi(x, y; \theta) \in \mathbb{R}^k$
<code>estimation_method</code>	<code>str</code>	Method identifier (see table below)
<code>estimator_kwargs</code>	<code>dict</code>	Method-specific parameters (overrides defaults)
<code>hyperparams</code>	<code>dict</code>	Hyperparameter search space as <code>{param: [val1, val2, ...]}</code>
<code>sweep_hparams</code>	<code>bool</code>	If <code>True</code> , performs grid search over <code>hyperparams</code>
<code>validation_data</code>	<code>dict</code>	Validation set (if <code>None</code> , uses <code>train_data</code>)
<code>val_loss_func</code>	<code>callable or str</code>	Validation metric: custom function, ' <code>mmr</code> ', ' <code>hsic</code> ', or ' <code>moment_violation</code> '
<code>normalize_moment_function</code>	<code>bool</code>	If <code>True</code> , pretrains and normalizes moment components to unit variance
<code>verbose</code>	<code>bool or int</code>	Verbosity level: <code>False/0</code> , <code>True/1</code> , or <code>2</code>

Returns:

- `trained_model`: Trained PyTorch model (best from hyperparameter search)
- `stats`: Dictionary containing:
 - `'models'`: List of all trained models
 - `'val_loss'`: Validation losses for each hyperparameter setting
 - `'hyperparam'`: List of hyperparameter configurations
 - `'best_index'`: Index of best model
 - `'train_stats'`: Training statistics for each model

Estimation Methods:

estimation_method	Problem Type	Description
'OLS'	Unconditional	Ordinary least squares
'GMM'	Unconditional	Generalized method of moments
'GEL'	Unconditional	Generalized empirical likelihood
'KMM'	Unconditional	Kernel method of moments
'SMD'	Conditional	Sieve minimum distance
'MMR'	Conditional	Maximum moment restrictions
'VMM-kernel'	Conditional	Variational MM with kernel
'VMM-neural'	Conditional	Variational MM with neural network
'FGEL-kernel'	Conditional	Functional GEL with kernel
'FGEL-neural'	Conditional	Functional GEL with neural network
'KMM-neural'	Conditional	Kernel MM with neural instrument function

4.2 Low-Level Interface

For fine-grained control without hyperparameter search:

```
from cmr.methods.kmm_neural import KMMNeural

estimator = KMMNeural(
    model=model,
    moment_function=moment_function,
    entropy_reg_param=10.0,
    reg_param=1e-2,
    n_random_features=2000,
    verbose=True
)

estimator.train(train_data, validation_data)
trained_model = estimator.model
parameters = estimator.get_trained_parameters()
```

4.3 Data Format

All data must be provided as dictionaries with NumPy arrays:

```
train_data = {
    't': np.array([[...]]),  # shape: (n_samples, d_t)
```

```

'y': np.array([[...]]), # shape: (n_samples, d_y)
'z': np.array([[...]]) # shape: (n_samples, d_z) or None for unconditional
}

```

For unconditional moment restrictions, set 'z': None.

4.4 Moment Function

The moment function must have signature:

```

def moment_function(model_output: torch.Tensor, y: torch.Tensor) -> torch.Tensor:
    """
    Args:
        model_output: Model predictions f_theta(x), shape (n, d_output)
        y: Observed outcomes, shape (n, d_y)

    Returns:
        Moment values psi(x, y; theta), shape (n, k)
    """
    return model_output - y # Example: IV regression

```

Common moment functions:

1. **IV Regression**: lambda pred, y: pred - y
2. **Quantile Regression**: lambda pred, y: (y - pred) * (tau - (y < pred).float())
3. **Vector-valued**: Return tensor with shape[1] = k for k moment conditions

4.5 Validation Metrics

Built-in validation metrics (string identifiers):

- 'moment_violation': $\frac{1}{n} \sum_{i=1}^n \|\psi(x_i, y_i; \theta)\|^2$ (unconditional)
- 'mmr': $\frac{1}{n^2} \sum_{i,j} \psi_i^\top K_{ij}^Z \psi_j$ (conditional, requires Z)
- 'hsic': Hilbert-Schmidt Independence Criterion between $\psi(\theta)$ and Z (conditional)

Custom validation functions:

```

def custom_val_loss(model: torch.nn.Module, val_data: dict) -> float:
    """
    Args:
        model: Trained PyTorch model
        val_data: Validation data dict with keys 't', 'y', 'z'

    Returns:
        Scalar validation loss (lower is better)
    """
    # Your custom validation logic
    return loss_value

```

4.6 Common Workflow

```

import torch
import numpy as np
from cmr import estimation

# 1. Generate/load data
train_data = {'t': X_train, 'y': Y_train, 'z': Z_train}
val_data = {'t': X_val, 'y': Y_val, 'z': Z_val}

```

```

test_data = {'t': X_test, 'y': Y_test, 'z': Z_test}

# 2. Define model architecture
model = torch.nn.Sequential(
    torch.nn.Linear(dim_t, 50),
    torch.nn.LeakyReLU(),
    torch.nn.Linear(50, 20),
    torch.nn.LeakyReLU(),
    torch.nn.Linear(20, dim_y)
)

# 3. Define moment function
def moment_function(pred, y):
    return pred - y

# 4. Train with hyperparameter search
trained_model, stats = estimation(
    model=model,
    train_data=train_data,
    moment_function=moment_function,
    estimation_method='KMM-neural',
    validation_data=val_data,
    verbose=True
)

# 5. Evaluate
predictions = trained_model(torch.Tensor(test_data['t']))
mse = torch.mean((predictions - torch.Tensor(test_data['y'])))**2

```

5. Theoretical Guarantees

5.1 Identification

For the CMR to identify θ , we require:

- **Completeness:** The conditional expectation operator $\mathbb{E}[\cdot | Z]$ is injective
- **Global identification:** θ is the unique solution to $\mathbb{E}[\psi(X, Y; \theta) | Z] = 0$

Under completeness, the CMR is equivalent to:

$$\int \psi(x, y; \theta) p(x, y | z) dx dy = 0 \quad \text{for all } z \in \mathcal{Z}$$

5.2 Consistency and Convergence Rates

Under regularity conditions:

MMR:

- Consistency: $\hat{\theta}_n \xrightarrow{p} \theta_0$
- Rate: $\|\hat{\theta}_n - \theta_0\| = O_p(n^{-r})$ where r depends on smoothness and dimension

FGEL-kernel:

- With RKHS regularization $\alpha_n = O(n^{-\gamma})$ for appropriate γ
- Achieves minimax optimal rates under smoothness assumptions

KMM:

- Combines advantages of MMD testing and moment restriction estimation
- Reference distribution improves finite-sample performance
- Convergence rate depends on RFF approximation quality and reference sample size

5.3 Asymptotic Normality

For semiparametric efficient estimation, under regularity:

$$\sqrt{n}(\hat{\theta}_n - \theta_0) \xrightarrow{d} \mathcal{N}(0, V)$$

where V is the semiparametric efficiency bound. Methods like FGEL and VMM can achieve this bound with appropriate choices of divergence and regularization.

5.4 Computational Complexity

Method	Per-Iteration Complexity	Memory
OLS/GMM	$O(nk^2)$	$O(nk)$
GEL	$O(nk)$	$O(nk)$
MMR	$O(n^2k)$	$O(n^2)$
VMM-kernel	$O(n^3k)$ (matrix inversion)	$O(n^2k)$
FGEL-kernel	$O(n^2k^2)$	$O(nk)$
VMM-neural	$O(nk \cdot \text{NN})$	$O(nk)$
FGEL-neural	$O(nk \cdot \text{NN})$	$O(nk)$
KMM	$O(n^2)$ (full kernel)	$O(n^2)$
KMM-neural	$O(nD + n \cdot \text{NN})$	$O(nD)$

where NN denotes neural network forward/backward pass cost and D is the RFF dimension.

Scalability recommendations:

- $n < 1000$: Any method works well
- $1000 \leq n < 5000$: Use neural versions or KMM with RFF
- $n \geq 5000$: Mandatory to use neural networks or RFF approximations
- Use mini-batch training (`batch_size` parameter) for $n > 10000$

References

1. Kremer, H., Nemmour, Y., Schölkopf, B., & Zhu, J. J. (2023). Estimation Beyond Data Reweighting: Kernel Method of Moments. arXiv:2305.10898
2. Kremer, H., Zhu, J. J., Muandet, K., & Schölkopf, B. (2022). Functional Generalized Empirical Likelihood Estimation for Conditional Moment Restrictions. ICML 2022.
3. Bennett, A., Kallus, N., & Schnabel, T. (2020). Deep Generalized Method of Moments for Instrumental Variable Analysis. NeurIPS 2019.
4. Muandet, K., Fukumizu, K., Sriperumbudur, B., & Schölkopf, B. (2017). Kernel Mean Embedding of Distributions: A Review and Beyond. Foundations and Trends in Machine Learning.