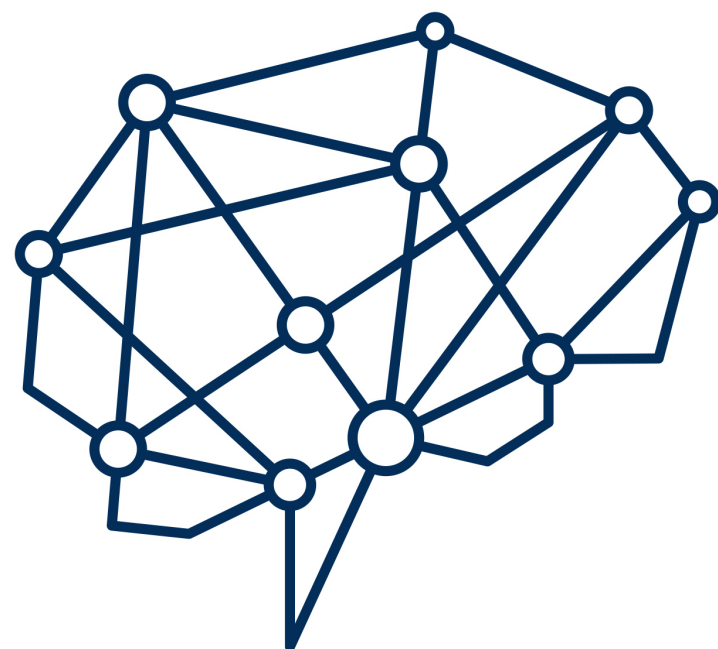


Optimal Kernel Hypothesis Testing

Antonin Schrab
UCL PhD 2020-2024
CDT in Foundational AI
Centre for Artificial Intelligence
Gatsby Computational Neuroscience Unit



Research Papers

MMD Aggregated Two-Sample Test

Schrab, Kim, Albert, Laurent, Guedj & Gretton
JMLR 2022

KSD Aggregated Goodness-of-fit Test

Schrab, Guedj & Gretton
NeurIPS 2022

Efficient Aggregated Kernel Tests using Incomplete U-statistics

Schrab, Kim, Guedj & Gretton
NeurIPS 2022

MMD-FUSE: Learning and Combining Kernels for Two-Sample Testing Without Data Splitting

Biggs*, Schrab* & Gretton
NeurIPS 2023 (Spotlight)

Differentially Private Permutation Tests: Applications to Kernel Methods

Kim* & Schrab*
Under journal review

Robust Kernel Hypothesis Testing under Data Corruption

Schrab* & Kim*
Under conference review

Outline

Part I: Kernel Discrepancies & Estimators

- Maximum Mean Discrepancy (MMD)
- Hilbert-Schmidt Independence Criterion (HSIC)
- Kernel Stein Discrepancy (KSD)
- Efficient estimators as incomplete U-statistics
- Adaptive estimators via kernel pooling

Part II: Optimal Kernel Hypothesis Testing

- Hypothesis testing
- Adaptive testing via kernel aggregation
- Testing constraints: efficiency, privacy & robustness
- Two-sample MMD testing
- Independence HSIC testing
- Goodness-of-fit KSD testing
- Open problems

Part I

Kernel Discrepancies & Estimators

Maximum Mean Discrepancy

Kernel Methods

Kernel: function $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ s.t. $k(x, y)$ measures the similarity between x and y

Inner product of features: feature map ϕ

$$k(x, y) = \langle \phi(x), \phi(y) \rangle$$

Positive definite: for all $x_1, \dots, x_n \in \mathcal{X}$ and all $c_1, \dots, c_n \in \mathbb{R}$

$$\sum_{i=1}^n \sum_{j=1}^n c_i c_j k(x_i, x_j) \geq 0$$

Reproducing Kernel Hilbert Space: \mathcal{H}_k space of real-valued functions on \mathcal{X}

$$k(x, \cdot) \in \mathcal{H}_k \qquad \langle f, k(x, \cdot) \rangle_{\mathcal{H}_k} = f(x)$$

MMD: Maximum Mean Discrepancy

Integral probability metric:

$$\text{MMD}_k(P, Q) = \sup_{f \in \mathcal{H}_k : \|f\|_{\mathcal{H}_k} \leq 1} \mathbb{E}_{X \sim P}[f(X)] - \mathbb{E}_{Y \sim Q}[f(Y)]$$

Difference in mean embeddings: $\langle f, \mu_P \rangle_{\mathcal{H}_k} = \mathbb{E}_P[f(X)]$

$$\text{MMD}_k = \sup_{f \in \mathcal{H}_k : \|f\|_{\mathcal{H}_k} \leq 1} \langle f, \mu_P - \mu_Q \rangle_{\mathcal{H}_k} = \|\mu_P - \mu_Q\|_{\mathcal{H}_k}$$

Expected kernel expression:

$$\text{MMD}_k^2 = \|\mu_P - \mu_Q\|_{\mathcal{H}_k}^2 = \mathbb{E}_{P,P}[k(X, X')] - 2\mathbb{E}_{P,Q}[k(X, Y)] + \mathbb{E}_{Q,Q}[k(Y, Y')]$$

MMD Estimators

Biased V-statistic estimator:

$$V_{\text{MMD}_k^2} = \frac{1}{m^2} \sum_{1 \leq i, i' \leq m} k(X_i, X_{i'}) - \frac{2}{mn} \sum_{i=1}^m \sum_{j=1}^n k(X_i, Y_j) + \frac{1}{n^2} \sum_{1 \leq j, j' \leq n} k(Y_j, Y_{j'})$$

Unbiased U-statistic estimator:

$$U_{\text{MMD}_k^2} = \frac{1}{m(m-1)} \sum_{1 \leq i \neq i' \leq m} k(X_i, X_{i'}) - \frac{2}{mn} \sum_{i=1}^m \sum_{j=1}^n k(X_i, Y_j) + \frac{1}{n(n-1)} \sum_{1 \leq j \neq j' \leq n} k(Y_j, Y_{j'})$$

Core function: $h(x, x'; y, y') = k(x, x') - k(x', y) - k(x, y') + k(y, y')$

$$V_{\text{MMD}_k^2} = \frac{1}{m^2 n^2} \sum_{1 \leq i, i' \leq m} \sum_{1 \leq j, j' \leq n} h(X_i, X_{i'}; Y_j, Y_{j'}) \quad U_{\text{MMD}_k^2} = \frac{1}{m(m-1)n(n-1)} \sum_{1 \leq i \neq i' \leq m} \sum_{1 \leq j \neq j' \leq n} h(X_i, X_{i'}; Y_j, Y_{j'})$$

$$V_{\text{MMD}_k^2} = \frac{1}{n^2} \sum_{1 \leq i, i' \leq n} h(X_i, X_{i'}; Y_i, Y_{i'}) \quad \tilde{U}_{\text{MMD}_k^2} = \frac{1}{n(n-1)} \sum_{1 \leq i \neq i' \leq n} h(X_i, X_{i'}; Y_i, Y_{i'})$$

MMD Kernel Impact

Radial kernel with bandwidth λ :

$$k_{\lambda}(x, y) = f\left(\frac{\|x - y\|_r}{\lambda}\right)$$

MMD behaviour:

$$\text{MMD}_{\lambda}^2 \rightarrow 0 \text{ when } \lambda \rightarrow 0 \text{ or } \lambda \rightarrow \infty$$

$$V_{\text{MMD}_{\lambda}^2} \rightarrow \frac{1}{m} + \frac{1}{n} \text{ when } \lambda \rightarrow 0$$

$$U_{\text{MMD}_{\lambda}^2} \rightarrow 0 \text{ and } \tilde{U}_{\text{MMD}_{\lambda}^2} \rightarrow 0 \text{ when } \lambda \rightarrow 0 \text{ or } \lambda \rightarrow \infty$$

\Rightarrow Kernel bandwidth choice is crucial for having a meaningful metric

Hilbert-Schmidt Independence Criterion

HSIC: Hilbert-Schmidt Independence Criterion

Cross-covariance operator $\mathbf{C}_{P_{XY}}$:

$$\langle f, \mathbf{C}_{P_{XY}} g \rangle_{\mathcal{H}_k} = \mathbb{E}_{P_{XY}} \left[\left(f(X) - \mathbb{E}_{P_X}[f(X')] \right) \left(g(Y) - \mathbb{E}_{P_Y}[g(Y')] \right) \right]$$

$$\left\langle f \otimes g, \mathbf{C}_{P_{XY}} \right\rangle_{\text{HS}} = \mathbb{E}_{P_{XY}} \left[\left\langle f \otimes g, \left(k(X, \cdot) - \mu_{P_X} \right) \otimes \left(\ell(Y, \cdot) - \mu_{P_Y} \right) \right\rangle_{\text{HS}} \right]$$

$$\mathbf{C}_{P_{XY}} = \mathbb{E}_{P_{XY}} \left[\left(k(X, \cdot) - \mu_{P_X} \right) \otimes \left(\ell(Y, \cdot) - \mu_{P_Y} \right) \right]$$

HSIC:

$$\text{HSIC}_{k,\ell}^2(P_{XY}) = \|\mathbf{C}_{P_{XY}}\|_{\text{HS}}^2$$

$$= \mathbb{E}_{P_{XY}, P_{XY}} \left[\left\langle \left(k(X, \cdot) - \mu_{P_X} \right) \otimes \left(\ell(Y, \cdot) - \mu_{P_Y} \right), \left(k(X', \cdot) - \mu_{P_X} \right) \otimes \left(\ell(Y', \cdot) - \mu_{P_Y} \right) \right\rangle_{\text{HS}} \right]$$

$$= \mathbb{E}_{P_{XY}, P_{XY}} \left[k(X, X') \ell(Y, Y') \right] - 2 \mathbb{E}_{P_{XY}} \left[\mathbb{E}_{P_X}[k(X, X')] \mathbb{E}_{P_Y}[\ell(Y, Y')] \right] + \mathbb{E}_{P_X, P_X} \left[k(X, X') \right] \mathbb{E}_{P_Y, P_Y} \left[\ell(Y, Y') \right]$$

HSIC Estimators

Notation: $Z_i = X_i$, $Z_{n+j} = Y_j$, $K_{ij} = k(X_i, X_j)$, $L_{ij} = \ell(Y_i, Y_j)$

Biased V-statistic estimator:

$$V_{\text{HSIC}_{k,\ell}^2} = \frac{1}{N^2} \sum_{1 \leq i, j \leq N} K_{ij} L_{ij} - \frac{2}{N} \sum_{i=1}^N \left(\frac{1}{N} \sum_{j=1}^N K_{ij} \right) \left(\frac{1}{N} \sum_{r=1}^N L_{ir} \right) + \left(\frac{1}{N^2} \sum_{1 \leq i, j \leq N} K_{ij} \right) \left(\frac{1}{N^2} \sum_{1 \leq r, s \leq N} L_{rs} \right)$$

Unbiased U-statistic estimator:

$$U_{\text{HSIC}_{k,\ell}^2} = \frac{1}{|\mathbf{i}_2^N|} \sum_{(i,j) \in \mathbf{i}_2^N} K_{ij} L_{ij} - \frac{2}{|\mathbf{i}_3^N|} \sum_{(i,j,r) \in \mathbf{i}_3^N} K_{ij} L_{ir} + \frac{1}{|\mathbf{i}_4^N|} \sum_{(i,j,r,s) \in \mathbf{i}_4^N} K_{ij} L_{rs}$$

Core function: $h(Z_i, Z_j, Z_r, Z_s) = K_{ij}(L_{ij} - L_{is} - L_{rj} + L_{rs})$

$$V_{\text{HSIC}_{k,\ell}^2} = \frac{1}{N^4} \sum_{1 \leq i, j, r, s \leq N} h(Z_i, Z_j, Z_r, Z_s)$$

$$\tilde{V}_{\text{HSIC}_{k,\ell}^2} = \frac{1}{N^2} \sum_{1 \leq i, j \leq N} h(Z_i, Z_j, Z_{i+N/2}, Z_{j+N/2})$$

$$U_{\text{HSIC}_{k,\ell}^2} = \frac{1}{|\mathbf{i}_4^N|} \sum_{(i,j,r,s) \in \mathbf{i}_4^N} h(Z_i, Z_j, Z_r, Z_s)$$

$$\tilde{U}_{\text{HSIC}_{k,\ell}^2} = \frac{1}{N(N-1)} \sum_{1 \leq i \neq j \leq N} h(Z_i, Z_j, Z_{i+N/2}, Z_{j+N/2})$$

HSIC Kernel Impact

Radial kernels with bandwidths λ, μ :

$$k_{\lambda}(x, y) = f\left(\frac{\|x - y\|_r}{\lambda}\right) \qquad \ell_{\mu}(x, y) = g\left(\frac{\|x - y\|_s}{\mu}\right)$$

HSIC behaviour:

$$\begin{aligned} \text{HSIC}_{\lambda, \mu}^2 &\rightarrow 0 \quad \text{when } \lambda \rightarrow 0 \text{ or } \mu \rightarrow 0 \text{ or } \lambda \rightarrow \infty \text{ or } \mu \rightarrow \infty \\ V_{\text{HSIC}_{\lambda, \mu}^2} &\rightarrow \frac{1}{N} - \frac{1}{N^2} \quad \text{when } \lambda \rightarrow 0 \text{ and } \mu \rightarrow 0 \\ V_{\text{HSIC}_{\lambda, \mu}^2} &\rightarrow 0 \quad \text{when } \lambda \rightarrow \infty \text{ or } \mu \rightarrow \infty \\ U_{\text{HSIC}_{\lambda, \mu}^2} &\rightarrow 0 \quad \text{when } \lambda \rightarrow 0 \text{ or } \mu \rightarrow 0 \text{ or } \lambda \rightarrow \infty \text{ or } \mu \rightarrow \infty \end{aligned}$$

\Rightarrow Kernel bandwidth choice is crucial for having a meaningful metric

HSIC as an MMD

Product kernel:

$$(k \times \ell)((X, Y), (X', Y')) = k(X, X')\ell(Y, Y')$$

HSIC as an MMD:

$$\begin{aligned} & \text{MMD}_{k \times \ell}^2(P_{XY}, P_X \otimes P_Y) \\ &= \mathbb{E}_{P_{XY}, P_{XY}}[k(X, X')\ell(Y, Y')] - 2\mathbb{E}_{P_{XY}, P_X P_Y}[k(X, X')\ell(Y, Y')] + \mathbb{E}_{P_X P_Y, P_X P_Y}[k(X, X')\ell(Y, Y')] \\ &= \mathbb{E}_{P_{XY}, P_{XY}}[k(X, X')\ell(Y, Y')] - 2\mathbb{E}_{P_{XY}}[\mathbb{E}_{P_X}[k(X, X')]\mathbb{E}_{P_Y}[\ell(Y, Y')]] + \mathbb{E}_{P_X, P_X}[k(X, X')]\mathbb{E}_{P_Y, P_Y}[\ell(Y, Y')] \\ &= \text{HSIC}_{k, \ell}^2(P_{XY}) \end{aligned}$$

Why use HSIC? MMD estimator requires data splitting (independent samples)
HSIC estimator uses all paired samples to estimate all three terms (more accurate)

MMD as an HSIC

Setting: Given distributions P and Q , construct joint P_{XY} such that:

$$P_X = w_P P + w_Q Q \qquad Y = 1 \text{ if } X \sim P \text{ and } Y = -1 \text{ if } X \sim Q$$

HSIC as an MMD:

$$\text{HSIC}_{k,1}^2(P_{XY}) = 2w_P^2 w_Q^2 \text{MMD}_k^2(P, Q)$$

$$\text{HSIC}_{k,\ell_\mu}^2(P_{XY}) \rightarrow 2w_P^2 w_Q^2 \text{MMD}_k^2(P, Q) \text{ when } \mu \rightarrow 0$$

Estimators: $w_P = m/(m+n)$, $w_Q = n/(m+n)$

$$V_{\text{HSIC}_{k,1}^2} = \frac{2m^2 n^2}{(m+n)^4} V_{\text{MMD}_k^2} \qquad V_{\text{HSIC}_{k,\ell_\mu}^2} \rightarrow \frac{2m^2 n^2}{(m+n)^4} V_{\text{MMD}_k^2} \text{ when } \mu \rightarrow 0$$

Kernel Stein Discrepancy

Stein Discrepancy

Stein operator: $A_P: \text{Func}(\mathbb{R}^d \rightarrow \mathbb{R}) \rightarrow \text{Func}(\mathbb{R}^d \rightarrow \mathbb{R})$

$$P = Q \iff \mathbb{E}_Q[(A_P \mathbf{f})(X)] = 0 \text{ for all } \mathbf{f} \in \mathcal{F}$$

Stein discrepancy:

$$\text{SD}_{A_P}(P, Q) = \sup_{\mathbf{f} \in \mathcal{F}} \mathbb{E}_Q[(A_P \mathbf{f})(X)] - \mathbb{E}_P[(A_P \mathbf{f})(X)] = \sup_{\mathbf{f} \in \mathcal{F}} \mathbb{E}_Q[(A_P \mathbf{f})(X)]$$

Langevin Stein operator:

$$(A_P \mathbf{f})(x) = \mathbf{f}(x)^\top \nabla \log p(x) + \nabla^\top \mathbf{f}(x) \quad \text{where} \quad \nabla^\top \mathbf{f}(x) = \sum_{i=1}^d \frac{\partial}{\partial x_i} f_i(x)$$

Diffusion Stein operator:

$$(A_P \mathbf{f})(x) = \mathbf{f}(x)^\top \nabla \log p(x) + \nabla^\top \mathbf{f}(x) = \mathbf{f}(x)^\top \left(\frac{\nabla p(x)}{p(x)} \right) + \nabla^\top \mathbf{f}(x) = \frac{1}{p(x)} \left(\mathbf{f}(x)^\top \nabla p(x) + (\nabla^\top \mathbf{f}(x)) p(x) \right) = \frac{1}{p(x)} \left(\nabla^\top (\mathbf{f}(x) p(x)) \right)$$

Stein's identity satisfied:

$$\mathbb{E}_P[(A_P \mathbf{f})(X)] = \int_{\mathbb{R}^d} (A_P \mathbf{f})(x) p(x) dx = \int_{\mathbb{R}^d} \nabla^\top (\mathbf{f}(x) p(x)) dx = \sum_{i=1}^d \int_{\mathbb{R}^d} \frac{\partial}{\partial x_i} (f_i(x) p(x)) dx = 0$$

KSD: Kernel Stein Discrepancy

Kernelise the Stein operator:

$$(\mathbf{A}_P \mathbf{f})(x) = \mathbf{f}(x)^\top \nabla \log p(x) + \nabla^\top \mathbf{f}(x) = \langle \mathbf{f}, \xi_P(x) \rangle_{\mathcal{H}_k^d} \quad \xi_P(x) = \nabla \log p(x) k(x, \cdot) + \nabla k(x, \cdot)$$

Kernel Stein Discrepancy:

$$\text{KSD}_P(Q) = \sup_{\mathbf{f} \in \mathcal{F}_k} \mathbb{E}_Q[(\mathbf{A}_P \mathbf{f})(X)] = \sup_{\mathbf{f} \in \mathcal{F}_k} \langle \mathbf{f}, \mathbb{E}_Q[\xi_P(X)] \rangle_{\mathcal{H}^d} = \left\| \mathbb{E}_Q[\xi_P(X)] \right\|_{\mathcal{H}_k^d}$$

$$\text{KSD}_P^2(Q) = \left\| \mathbb{E}_Q[\xi_P(X)] \right\|_{\mathcal{H}_k^d}^2 = \langle \mathbb{E}_Q[\xi_P(X)], \mathbb{E}_Q[\xi_P(Y)] \rangle_{\mathcal{H}_k^d} = \mathbb{E}_{Q,Q}[\langle \xi_P(X), \xi_P(Y) \rangle_{\mathcal{H}_k^d}] = \mathbb{E}_{Q,Q}[h_P(X, Y)]$$

Stein kernel / core function:

$$\begin{aligned} h_P(x, y) &= \langle \xi_P(x), \xi_P(y) \rangle_{\mathcal{H}_k^d} \\ &= (\nabla \log p(x)^\top \nabla \log p(y)) k(x, y) + \nabla \log p(x)^\top \nabla_y k(x, y) + \nabla \log p(y)^\top \nabla_x k(x, y) + \sum_{i=1}^d \frac{\partial}{\partial x_i y_i} k(x, y) \\ &= \mathbf{A}_{P,2} \mathbf{A}_{P,1} (k(x, y) I_{d \times d}) \end{aligned}$$

Stein's identity:

$$\mathbb{E}_P[h_P(X, \cdot)] = \mathbb{E}_P[\xi_P(X)] = 0$$

KSD Estimators

Kernel Stein Discrepancy:

$$\text{KSD}_P^2(Q) = \mathbb{E}_{Q,Q}[h_P(X, Y)]$$

Stein kernel / core function:

$$h_P(x, y) = \left(\nabla \log p(x)^\top \nabla \log p(y) \right) k(x, y) + \nabla \log p(x)^\top \nabla_y k(x, y) + \nabla \log p(y)^\top \nabla_x k(x, y) + \sum_{i=1}^d \frac{\partial}{\partial x_i y_i} k(x, y)$$

Biased V-statistic estimator:

$$V_{\text{KSD}_k^2} = \frac{1}{n^2} \sum_{1 \leq i, j \leq n} h_P(X_i, X_j)$$

Unbiased U-statistic estimator:

$$U_{\text{KSD}_k^2} = \frac{1}{n(n-1)} \sum_{1 \leq i \neq j \leq n} h_P(X_i, X_j)$$

KSD Kernel Impact

Notation:

$$\mathbf{s}_p(x) = \nabla \log p(x), \quad \delta(x) = \mathbf{s}_p(x) - \mathbf{s}_q(x), \quad (\mathbf{A}_p \mathbf{f})(x) = \mathbf{f}(x)^\top \mathbf{s}_p(x) + \nabla^\top \mathbf{f}(x), \quad k_\lambda(x, y) = f(\|x - y\|_r / \lambda)$$

Useful property:

$$\mathbb{E}_Q[(\mathbf{A}_p \mathbf{f})(X)] = \mathbb{E}_Q[(\mathbf{A}_p \mathbf{f})(X) - (\mathbf{A}_q \mathbf{f})(X)] = \mathbb{E}_Q[\mathbf{f}(X)^\top (\mathbf{s}_p(X) - \mathbf{s}_q(X))] = \mathbb{E}_Q[\mathbf{f}(X)^\top \delta(X)]$$

Alternative KSD expression:

$$\mathbb{E}_{Q,Q} [k(X, Y) \delta(Y)^\top \delta(X)] = \mathbb{E}_{Q,Q} \left[\left(k(X, Y) \mathbf{s}_p(X) + \nabla_X k(X, Y) \right)^\top \delta(Y) \right] = \mathbb{E}_{Q,Q} [h_p(X, Y)] = \text{KSD}_P^2(Q)$$

KSD behaviour:

$$\text{KSD}_\lambda^2 = \mathbb{E}_{Q,Q} [k_\lambda(X, Y) \delta(X)^\top \delta(Y)] \rightarrow \mathbb{E}_Q [\delta(X)^\top \delta(X)] = \text{Fisher}(P, Q) \quad \text{when } \lambda \rightarrow 0$$

$$\text{KSD}_\lambda^2 = \mathbb{E}_{Q,Q} [k_\lambda(X, Y) \delta(X)^\top \delta(Y)] \rightarrow \mathbb{E}_{Q,Q} [\delta(X)^\top \delta(Y)] = \left\| \mathbb{E}_Q [\delta(X)] \right\|_2^2 \quad \text{when } \lambda \rightarrow \infty$$

Link to Fisher divergence:

$$\text{KSD}_P^2(Q) = \mathbb{E}_{Q,Q} [k(X, Y) \delta(X)^\top \delta(Y)] \leq \sqrt{\mathbb{E}_{Q,Q} [k(X, Y)^2] \mathbb{E}_{Q,Q} [(\delta(X)^\top \delta(Y))^2]} \leq \sqrt{\mathbb{E}_{Q,Q} [k(X, Y)^2] \text{Fisher}(P, Q)}$$

KSD as an MMD

Stein's identity:

$$\mathbb{E}_P[h_P(X, \cdot)] = 0$$

KSD as an MMD:

$$\begin{aligned} \text{MMD}_{h_P}^2(P, Q) \\ &= \mathbb{E}_{Q, Q}[h_P(X, X')] - 2\mathbb{E}_{Q, P}[h_P(X, Y)] + \mathbb{E}_{P, P}[h_P(Y, Y')] \\ &= \mathbb{E}_{Q, Q}[h_P(X, X')] \\ &= \text{KSD}_P^2(Q) \end{aligned}$$

Why use KSD?

MMD requires samples from the model P while KSD does not

MMD as a KSD

Simple Stein operator:

$$(\tilde{\mathbf{A}}_P \mathbf{f})(x) = f_1(x) - \mathbb{E}_P[f_1(X)]$$

Stein's identity:

$$\mathbb{E}_P[(\tilde{\mathbf{A}}_P \mathbf{f})(X)] = 0$$

KSD as an MMD:

$$\text{KSD}_{\tilde{\mathbf{A}}_P}(Q) = \sup_{\mathbf{f} \in \mathcal{H}_k^d} \mathbb{E}_Q[(\tilde{\mathbf{A}}_P \mathbf{f})(X)] = \sup_{f_1 \in \mathcal{H}_k} \mathbb{E}_Q[f_1] - \mathbb{E}_P[f_1] = \text{MMD}_k(P, Q)$$

One-sample MMD:

$$\frac{1}{N^2} \sum_{1 \leq i, j \leq N} k(X_i, X_j) - \frac{2}{N} \sum_{i=1}^N \mathbb{E}_P[k(X_i, Y)] + \mathbb{E}_{P, P}[k(Y, Y')]$$

Efficient Kernel Estimators

Efficient Kernel Estimators

Expectation:

$$\mathbb{E}[h(X, X')]$$

Statistic with design \mathcal{D} :

$$\frac{1}{|\mathcal{D}|} \sum_{(i,j) \in \mathcal{D}} h(X_i, X_j)$$

Examples:

$$V_{\text{MMD}_k^2} = \frac{1}{n^2} \sum_{1 \leq i, j \leq n} h(X_i, X_j; Y_i, Y_j),$$

$$h(x, x'; y, y') = k(x, x') - k(x', y) - k(x, y') + k(y, y')$$

$$\tilde{V}_{\text{HSIC}_{k,\ell}^2} = \frac{1}{N^2} \sum_{1 \leq i, j \leq N} h(Z_i, Z_j, Z_{i+N/2}, Z_{j+N/2}),$$

$$h(Z_i, Z_j, Z_r, Z_s) = K_{ij}(L_{ij} - L_{is} - L_{rj} + L_{rs})$$

$$V_{\text{KSD}_k^2} = \frac{1}{n^2} \sum_{1 \leq i, j \leq n} h_p(X_i, X_j), \quad h_p(x, y) = (\mathbf{s}_p(x)^\top \mathbf{s}_p(y)) k(x, y) + \mathbf{s}_p(x)^\top \nabla_y k(x, y) + \mathbf{s}_p(y)^\top \nabla_x k(x, y) + \sum_{i=1}^d \frac{\partial}{\partial x_i y_i} k(x, y)$$

V-statistic (von Mises)

	X_1	X_2	X_3	X_4	X_5	X_6	X_7	X_8	X_9	X_{10}
X_1										
X_2										
X_3										
X_4										
X_5										
X_6										
X_7										
X_8										
X_9										
X_{10}										

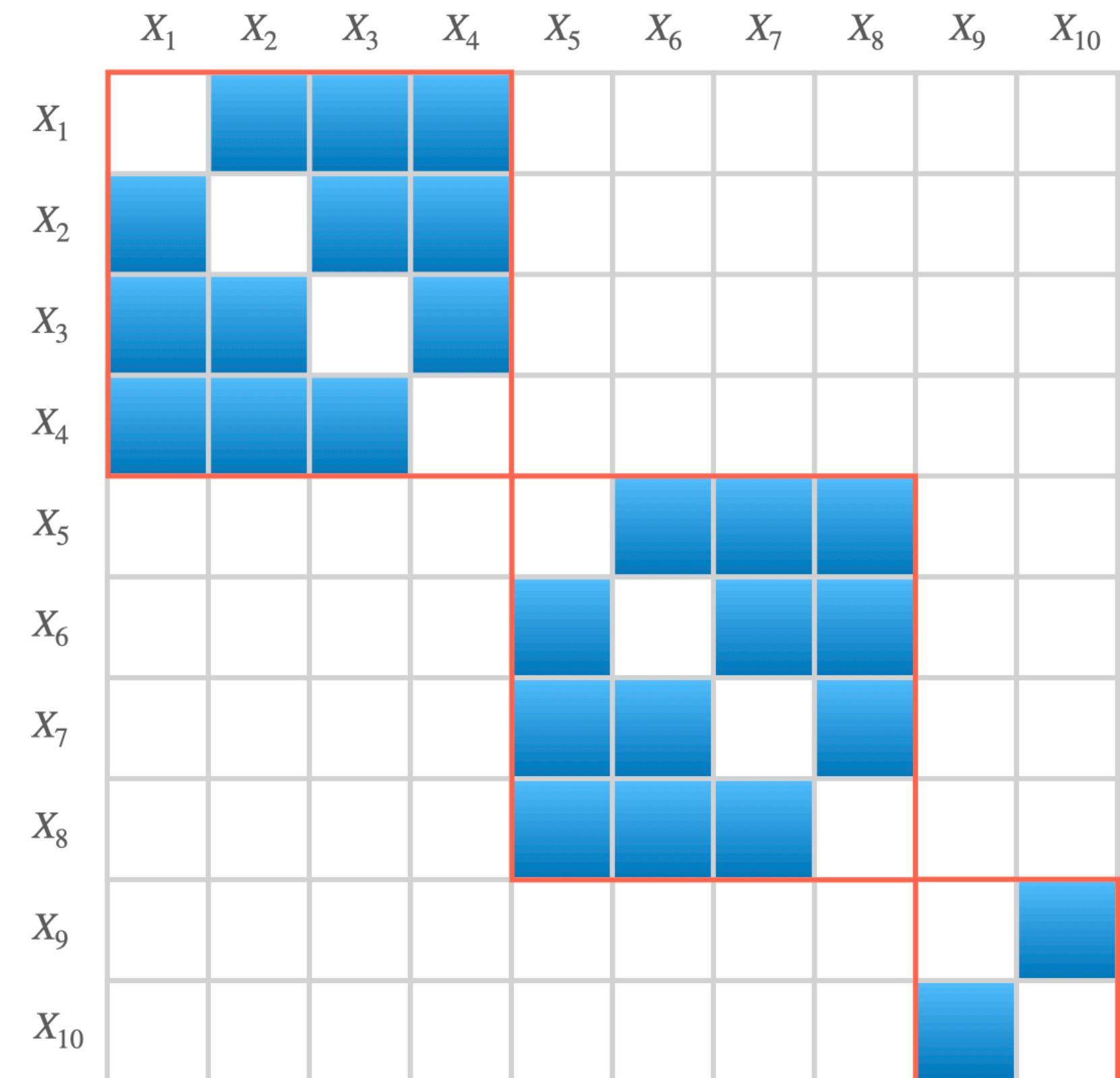
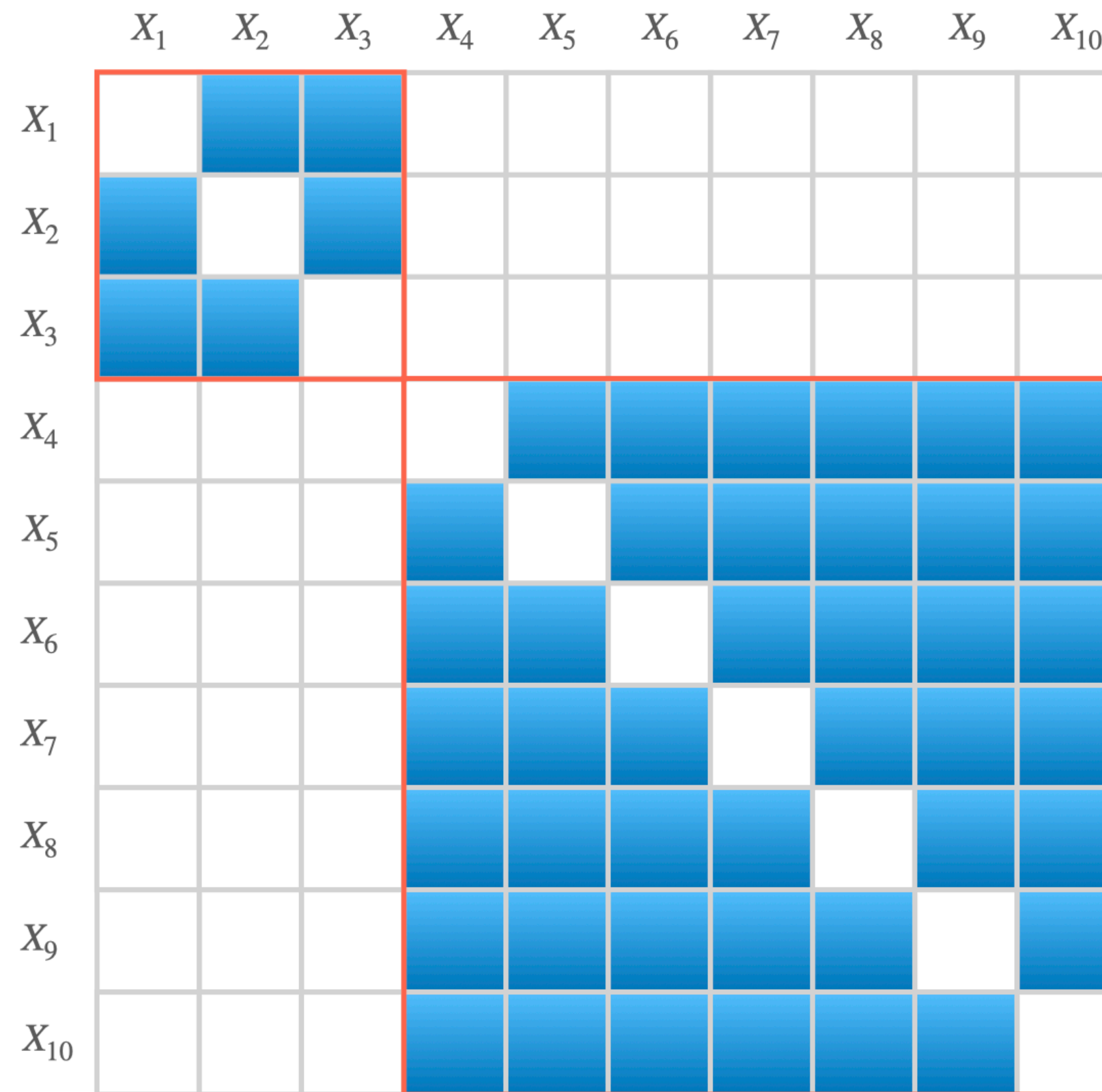
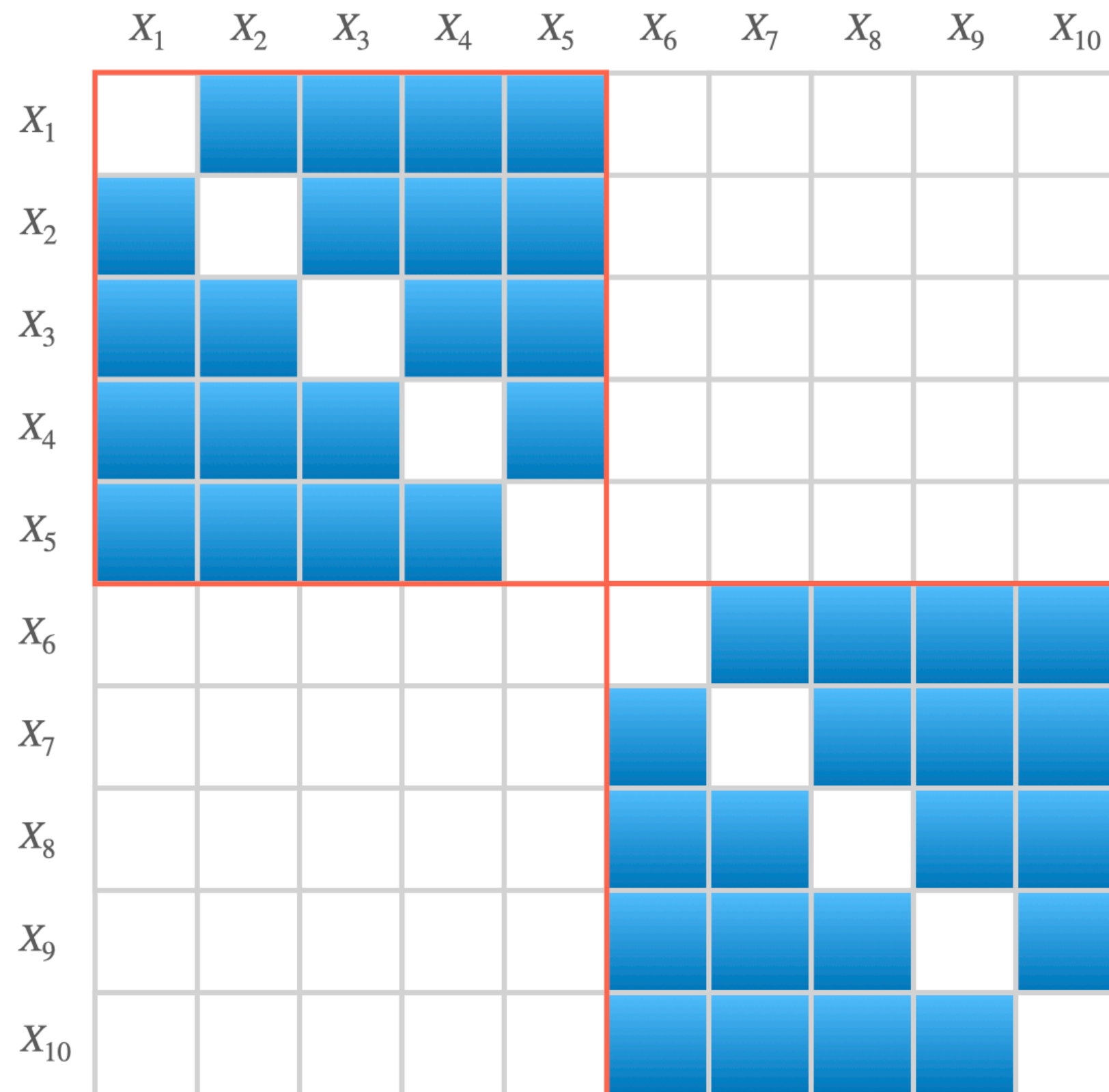
$$V = \frac{1}{n^2} \sum_{1 \leq i, j \leq n} h(X_i, X_j)$$

U-statistic (unbiased)

	X_1	X_2	X_3	X_4	X_5	X_6	X_7	X_8	X_9	X_{10}
X_1										
X_2										
X_3										
X_4										
X_5										
X_6										
X_7										
X_8										
X_9										
X_{10}										

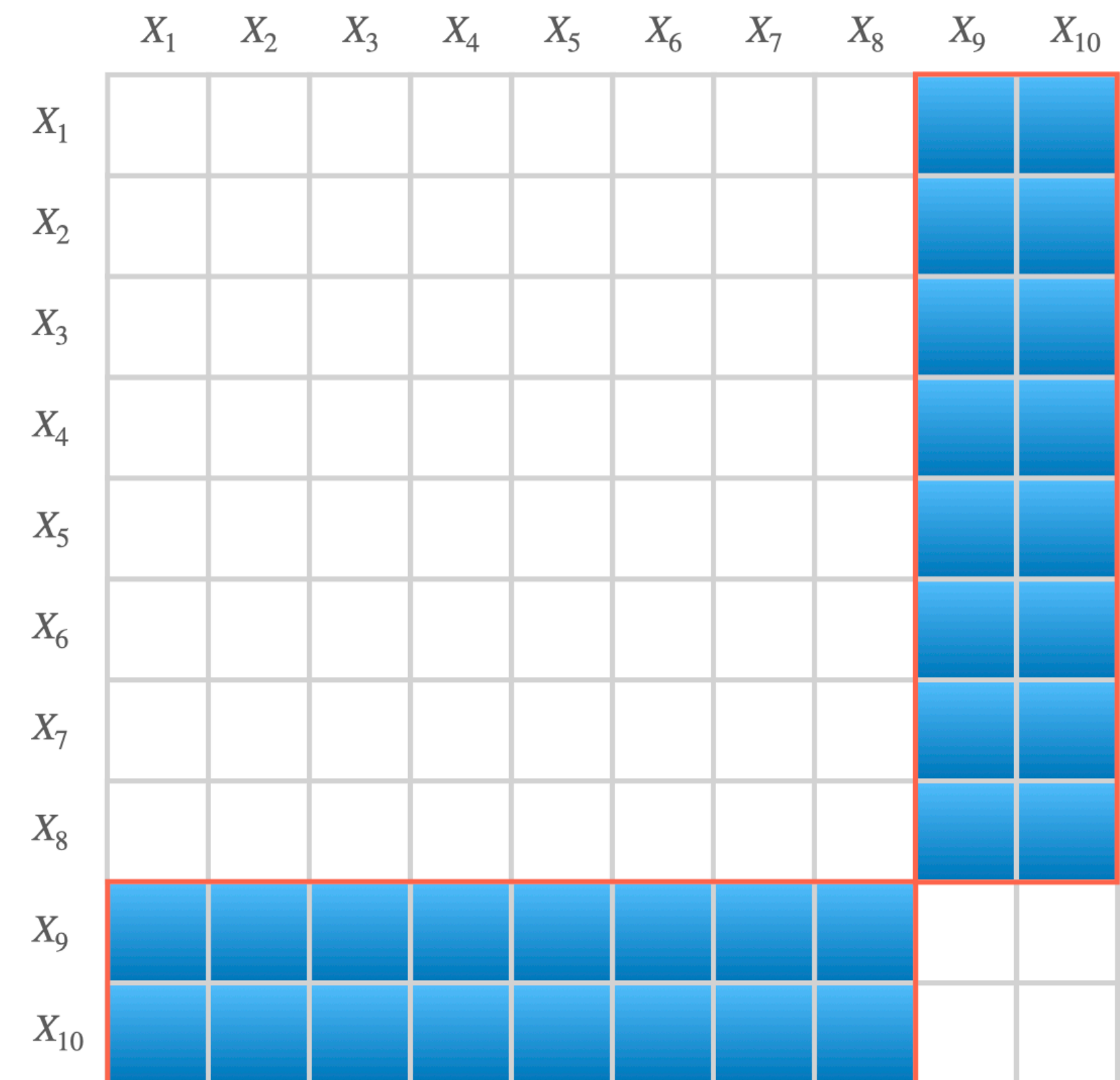
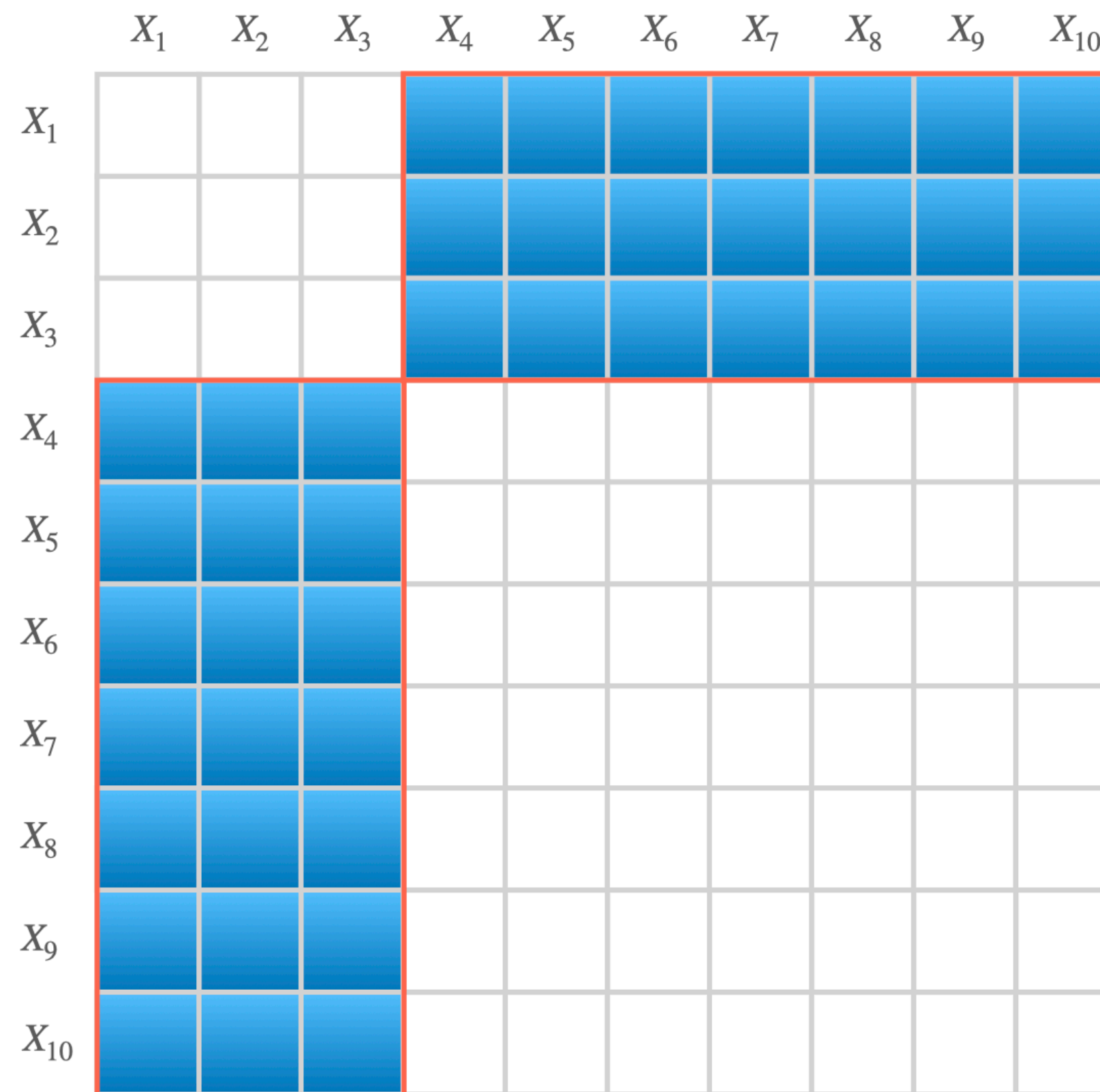
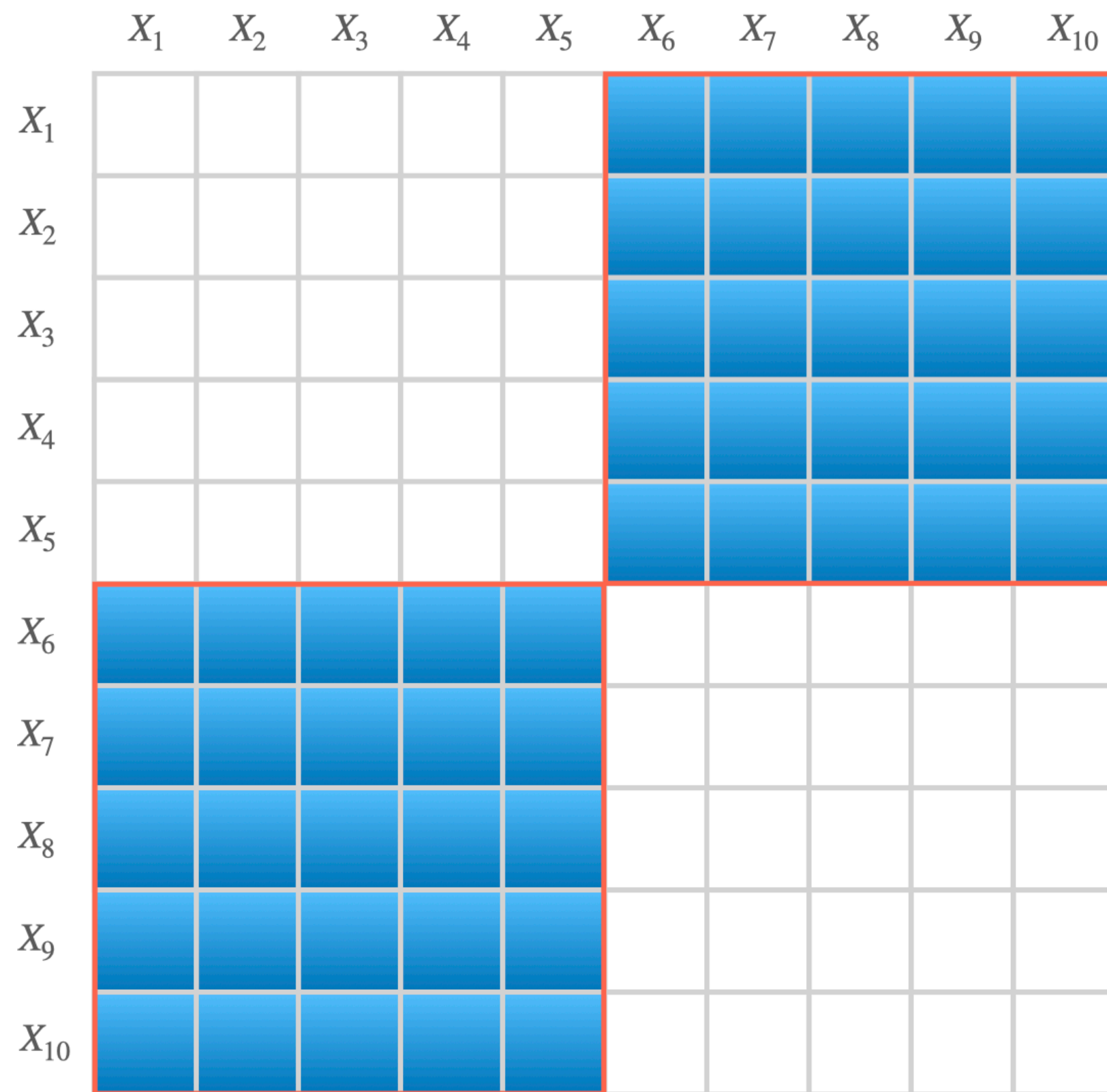
$$U = \frac{1}{n(n-1)} \sum_{1 \leq i \neq j \leq n} h(X_i, X_j)$$

B-statistic (block)



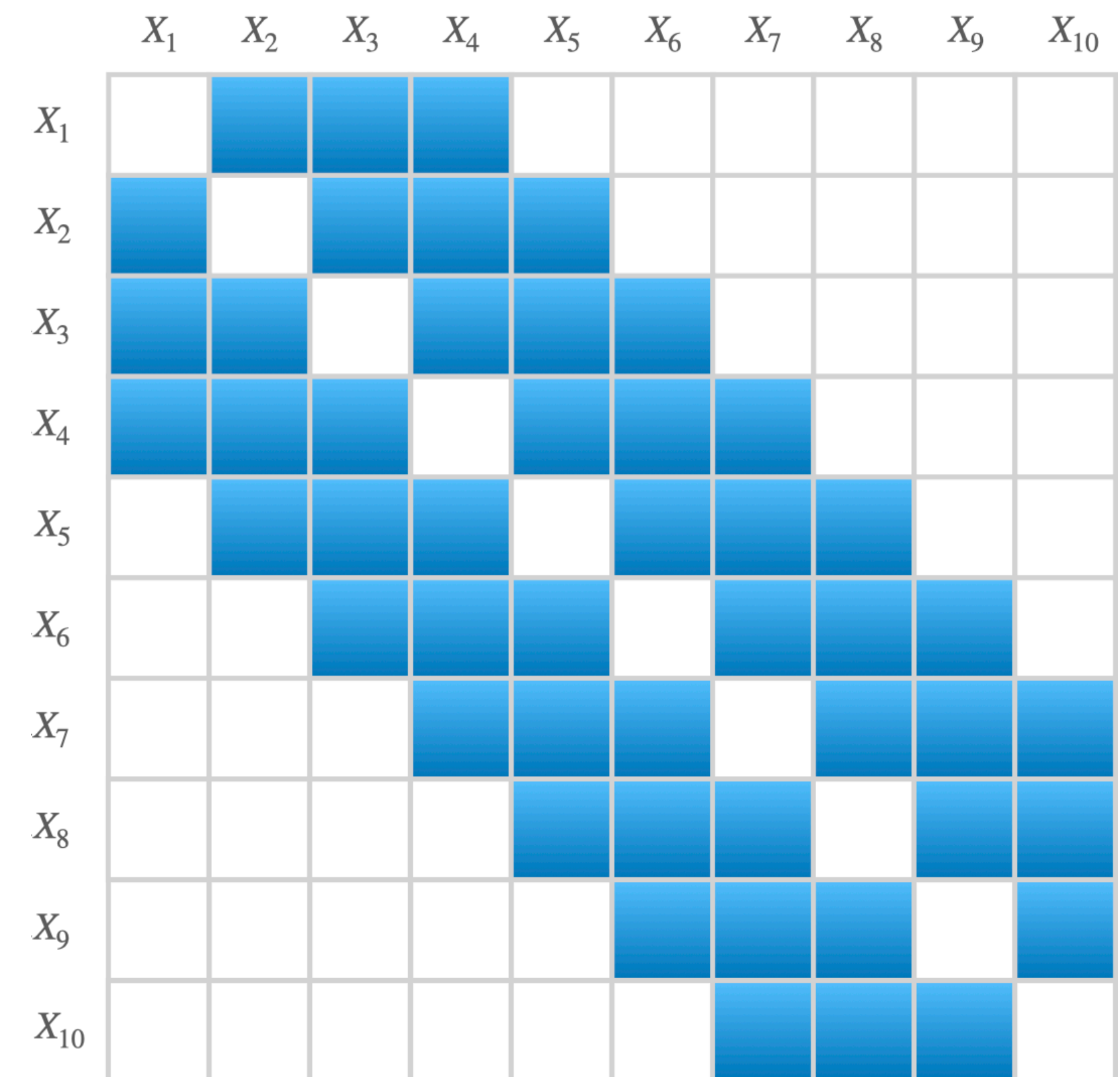
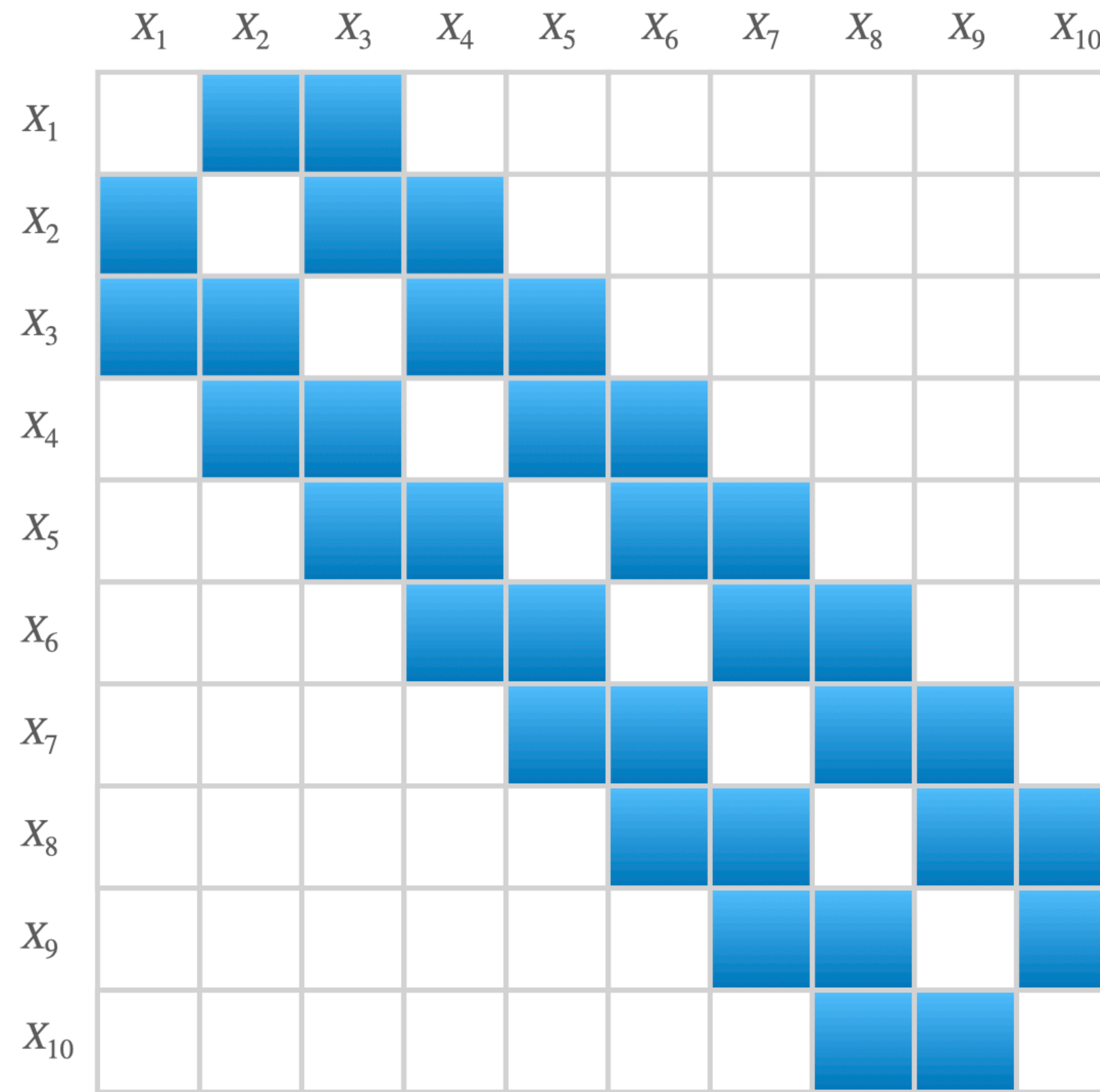
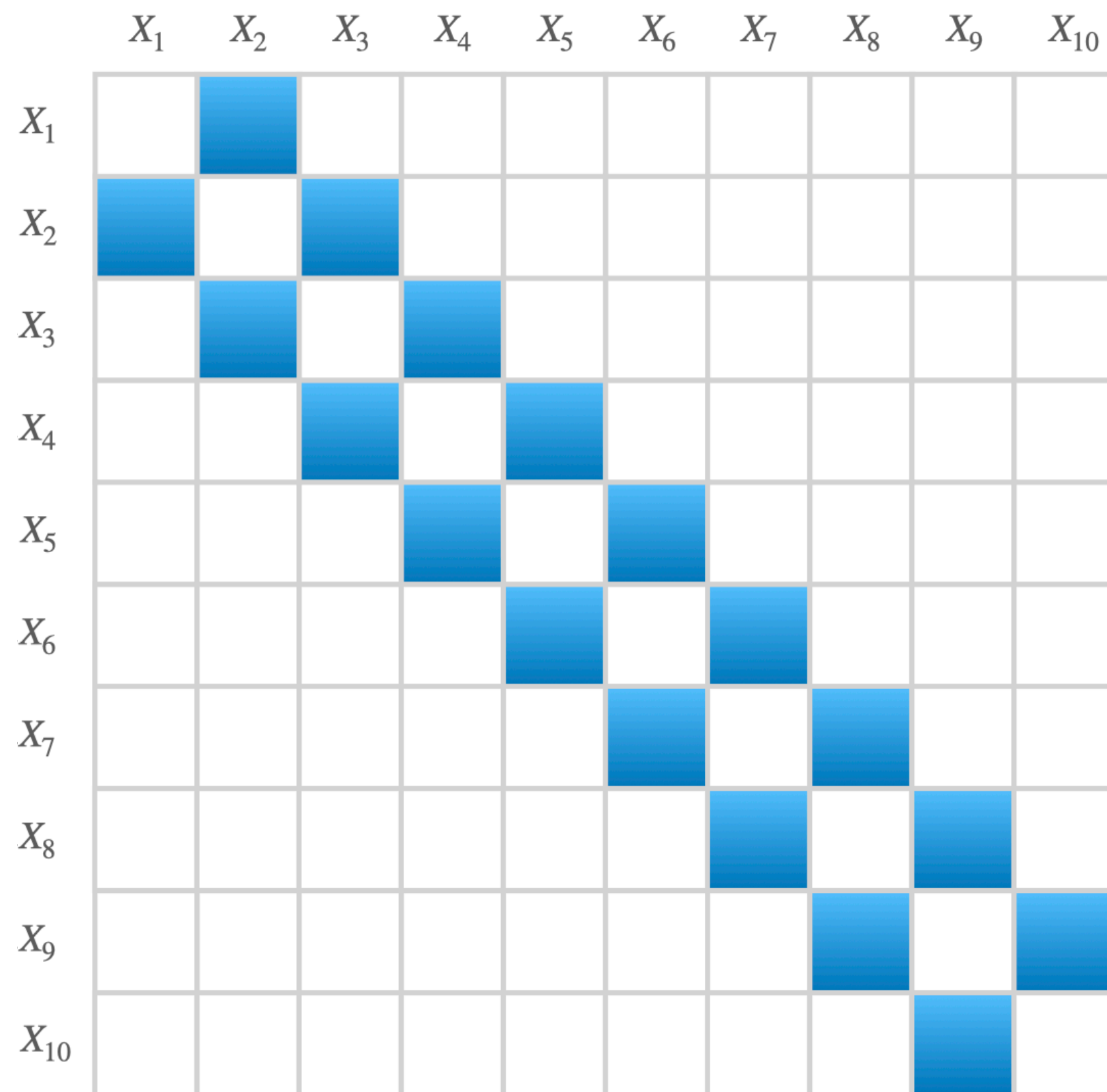
$$B = \sum_{s=1}^b U\left(X_{1+\sum_{t=0}^{s-1} n_t}, \dots, X_{\sum_{t=0}^s n_t}\right)$$

X-statistic (cross)



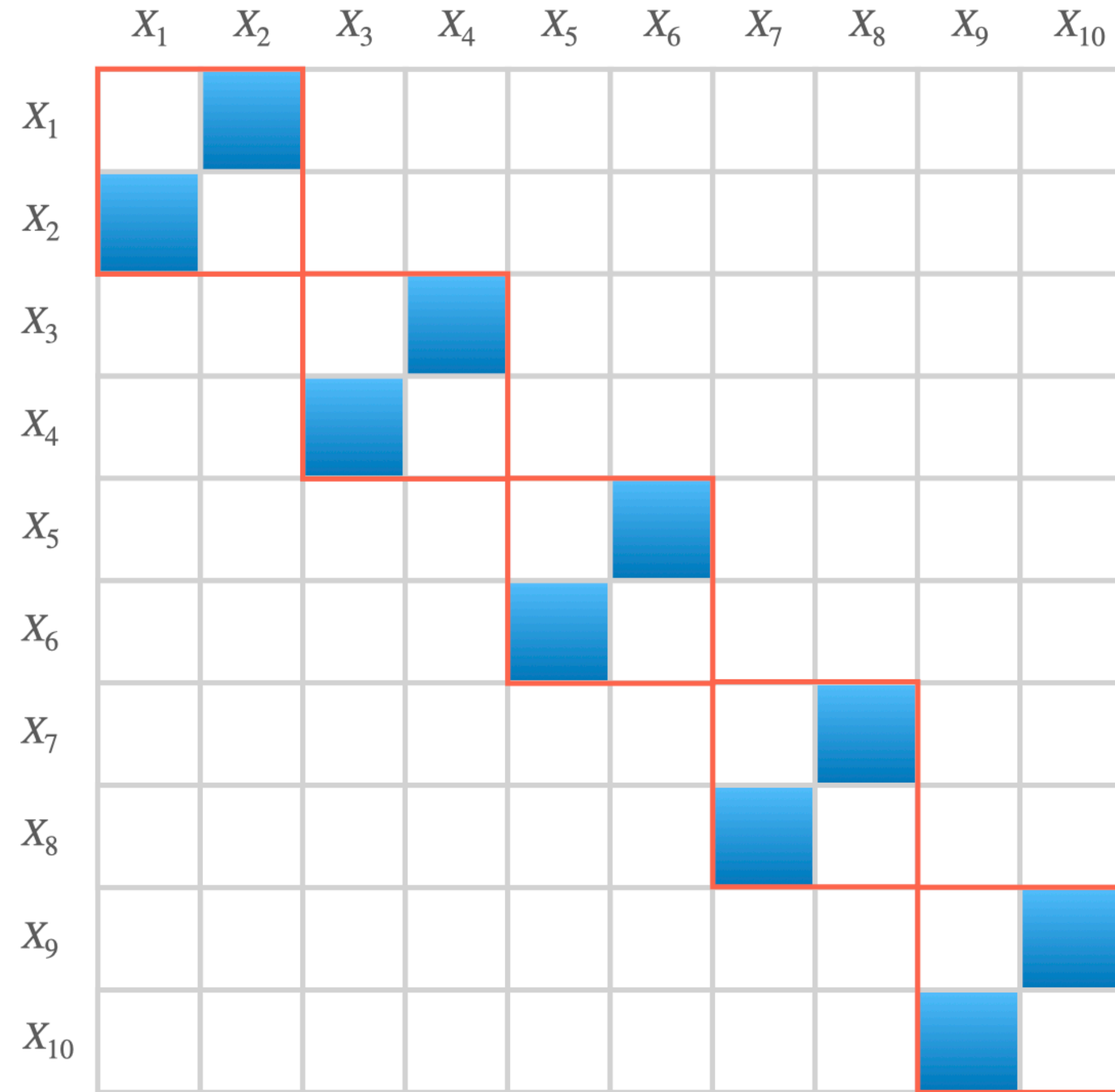
$$X = \frac{1}{n_1(n - n_1)} \sum_{i=1}^{n_1} \sum_{j=n_1+1}^n h(X_i, X_j)$$

D-statistic (diagonals)



$$D = \frac{2}{r(2n - r - 1)} \sum_{j=1}^r \sum_{i=1}^{n-j} h(X_i, X_{i+j})$$

L-statistic (linear)



$$L = \frac{1}{\lfloor n/2 \rfloor} \sum_{1 \leq i \leq \lfloor n/2 \rfloor} h(X_{2i}, X_{2i-1})$$

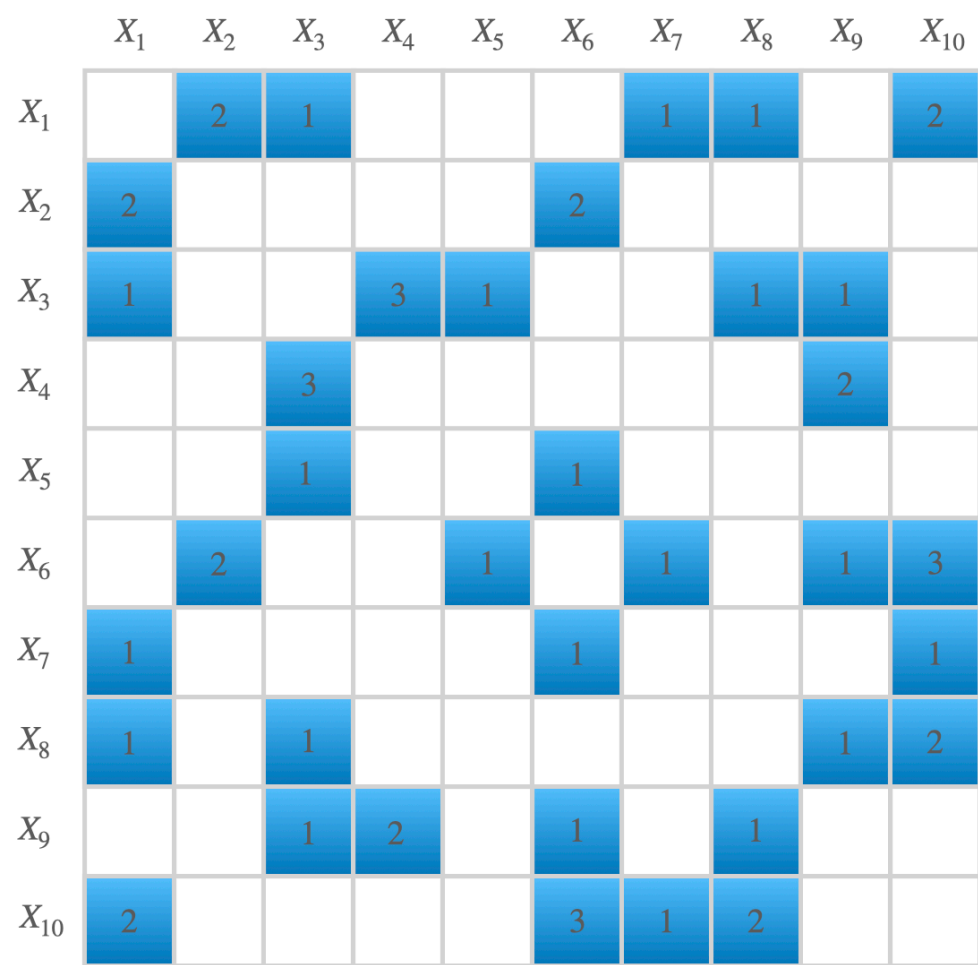
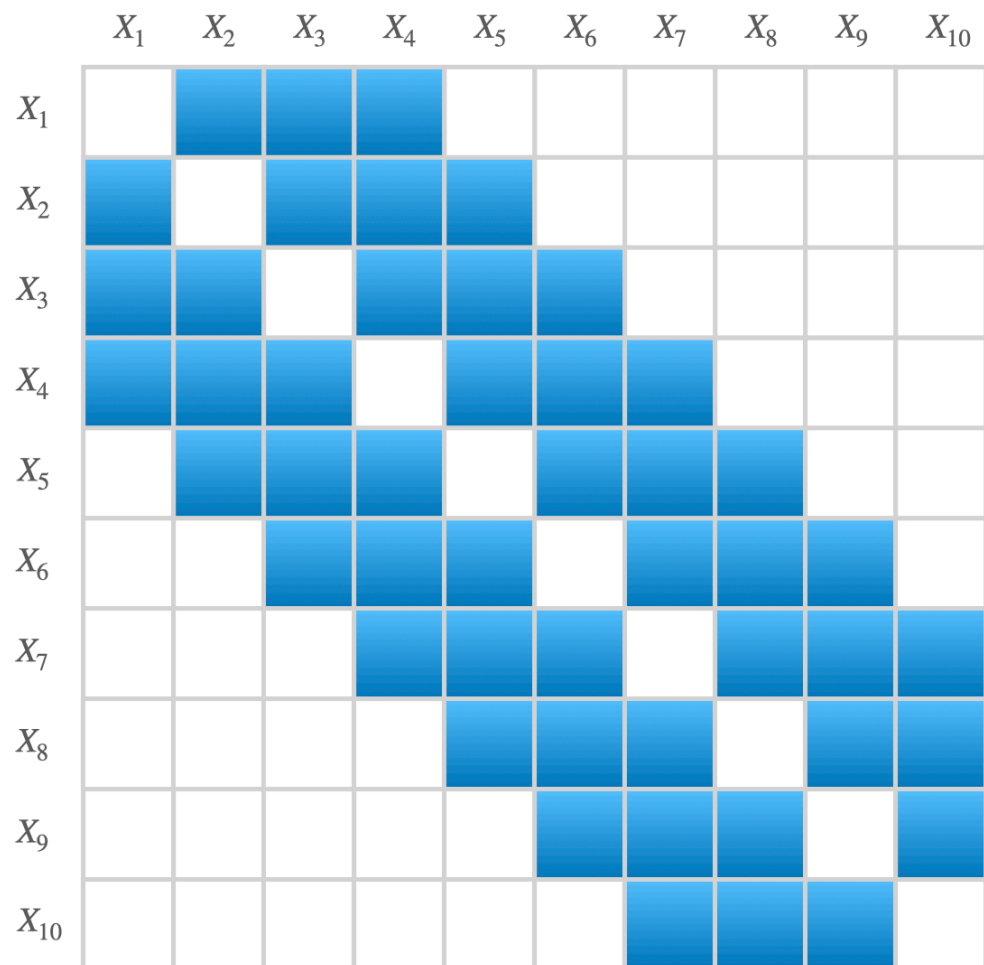
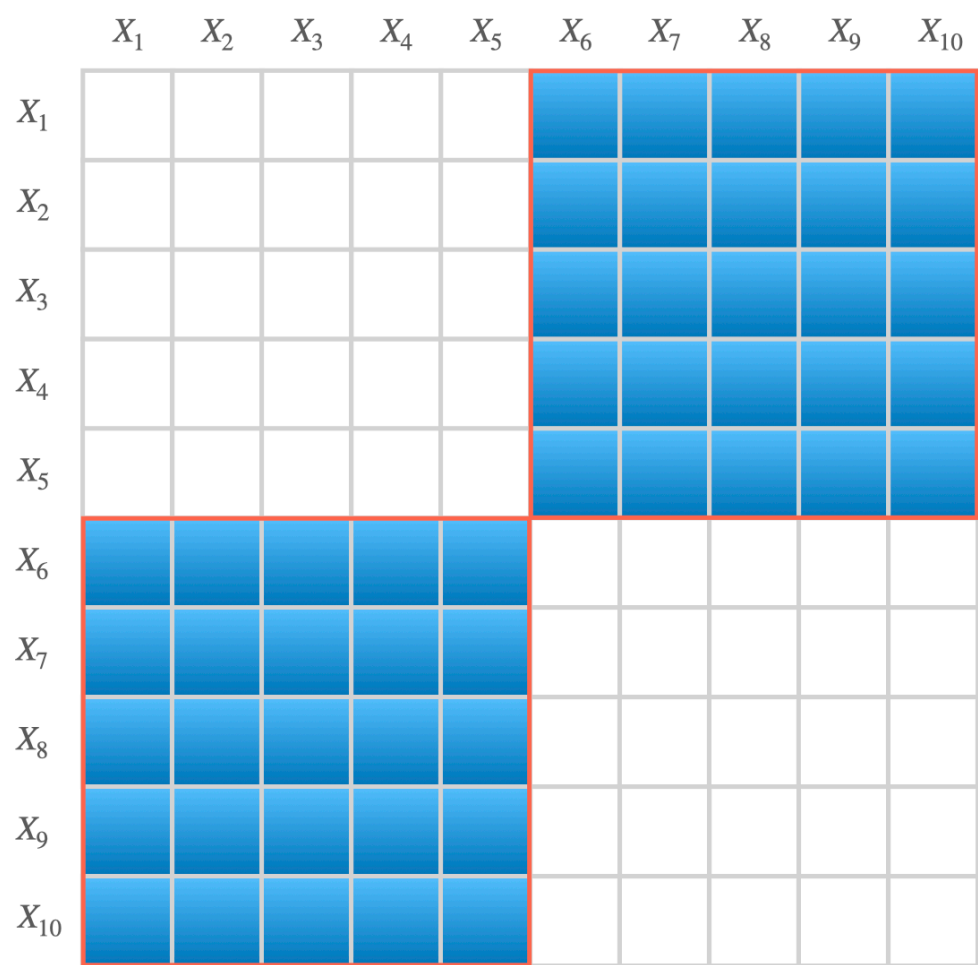
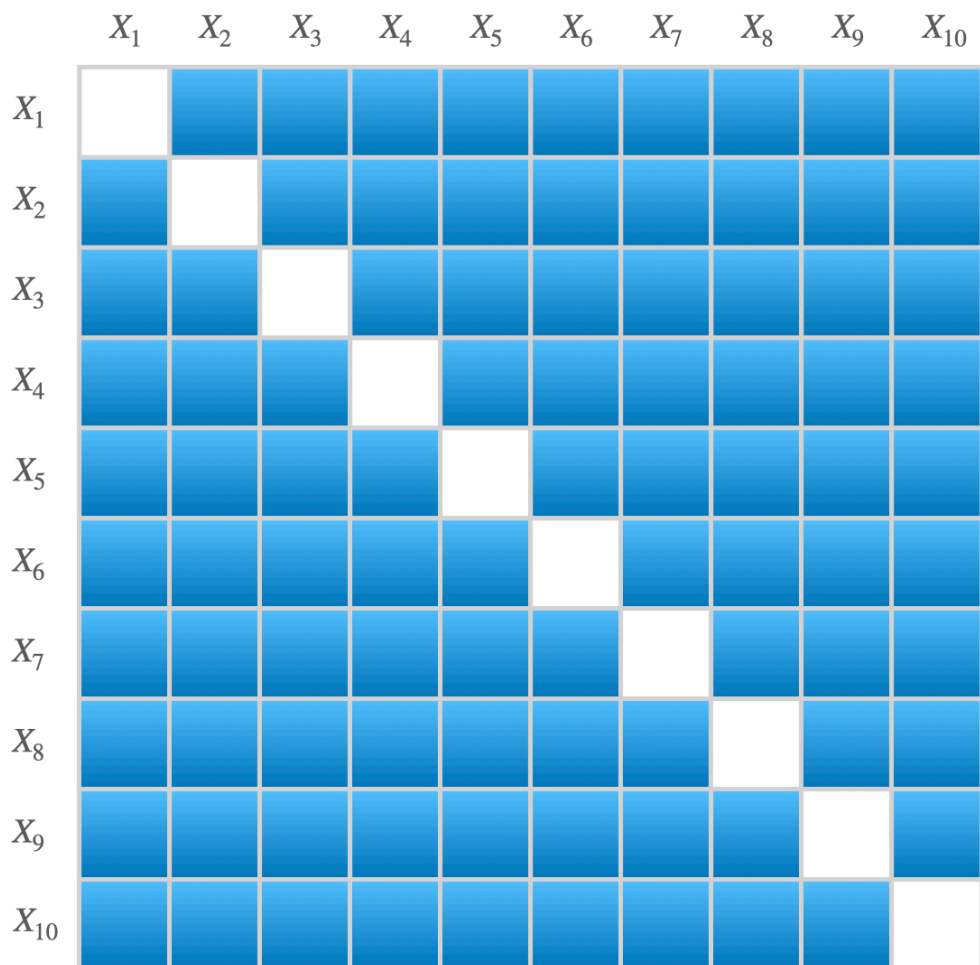
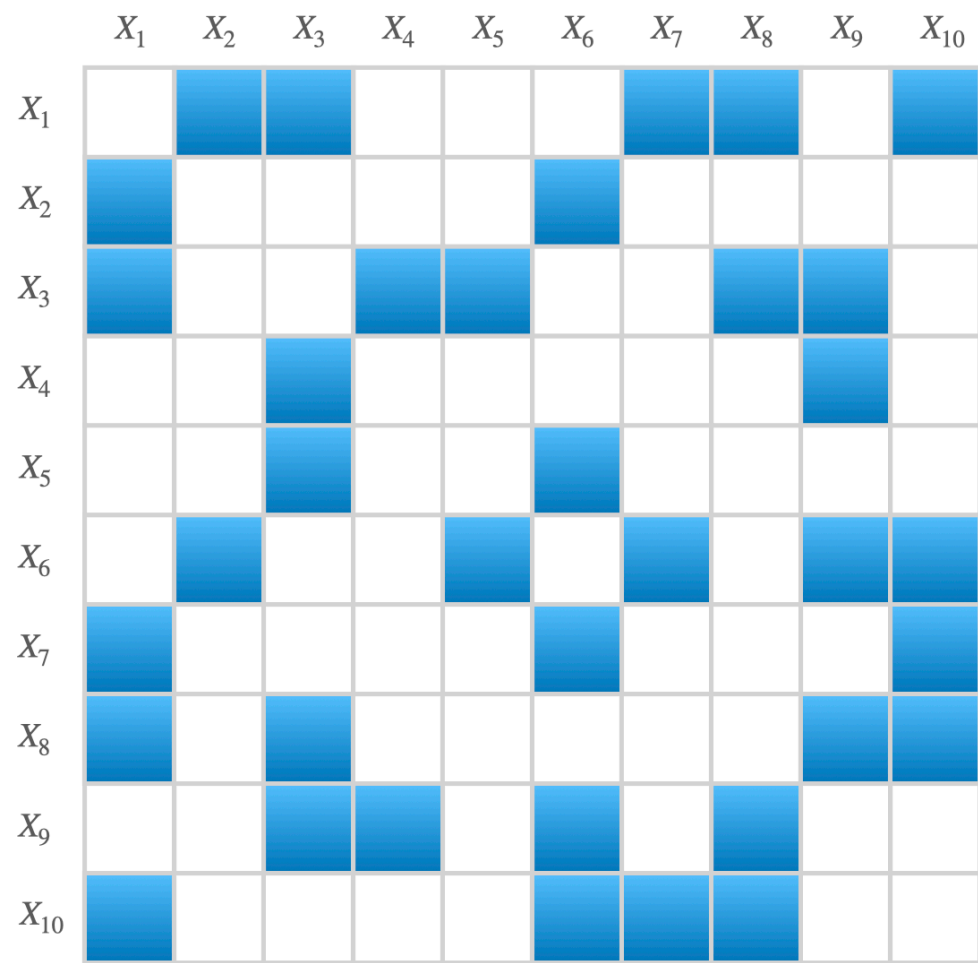
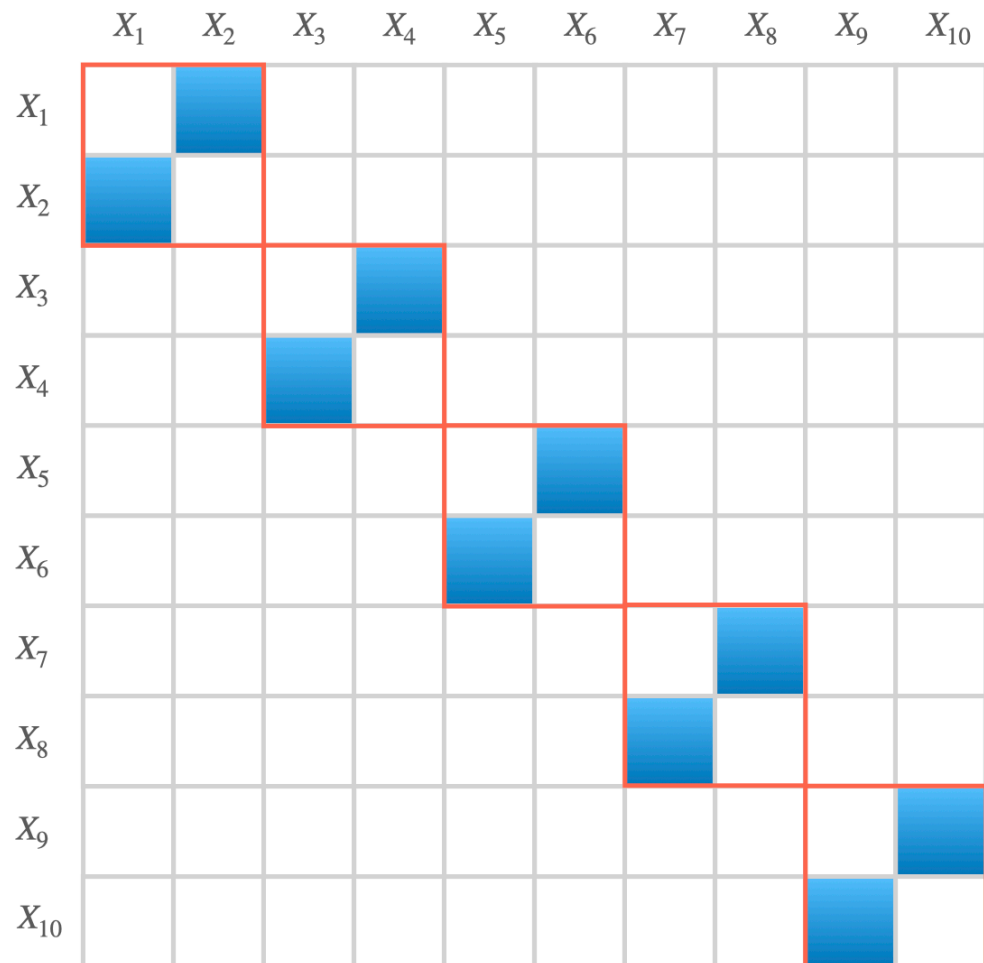
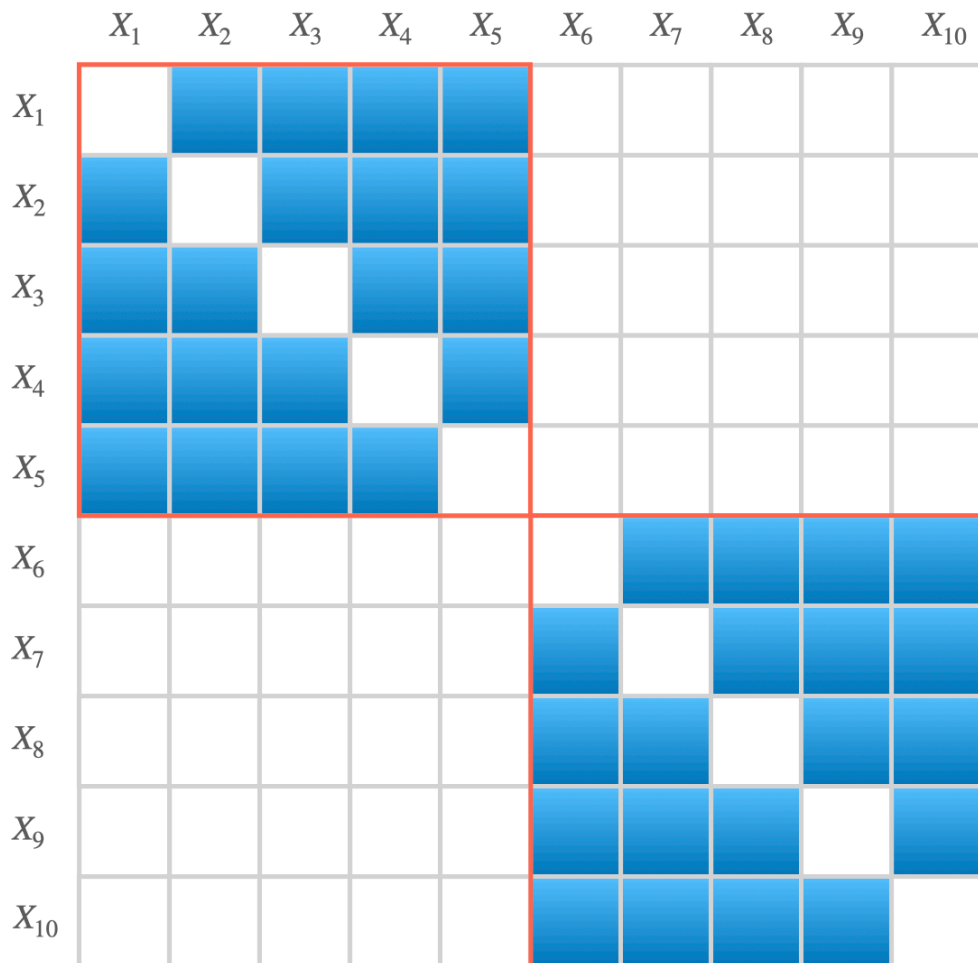
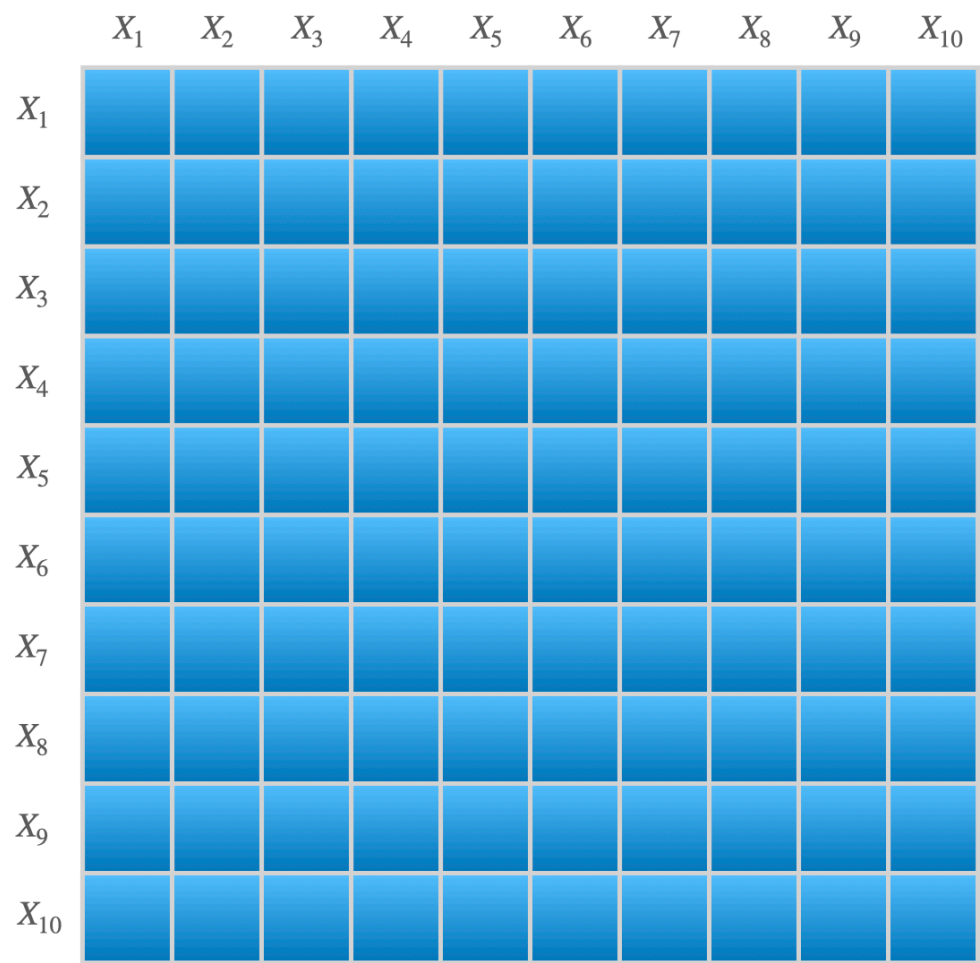
R-statistic (random)

	X_1	X_2	X_3	X_4	X_5	X_6	X_7	X_8	X_9	X_{10}
X_1										
X_2										
X_3										
X_4										
X_5										
X_6										
X_7										
X_8										
X_9										
X_{10}										

	X_1	X_2	X_3	X_4	X_5	X_6	X_7	X_8	X_9	X_{10}
X_1		2	1				1	1		2
X_2	2					2				
X_3	1			3	1			1	1	
X_4			3						2	
X_5			1			1				
X_6		2			1		1		1	3
X_7	1					1				1
X_8	1		1						1	2
X_9			1	2		1		1		
X_{10}	2					3	1	2		

$$R = \frac{1}{|\mathcal{D}_r|} \sum_{(i,j) \in \mathcal{D}_r} h(X_i, X_j)$$

Efficient Kernel Estimators



Adaptive Kernel Estimators

Kernel Pooling for Adaptivity

Collection of statistics / kernels:

$$S_k = \frac{1}{|\mathcal{D}_k|} \sum_{(i,j) \in \mathcal{D}_k} h_k(X_i, X_j), \quad k \in K$$

Normalisation (needed to compare statistics):

$$\frac{S_k}{\sigma_k} \quad \text{where} \quad \sigma_k^2 := \frac{4}{|\mathcal{D}_k^1|} \sum_{i \in \mathcal{D}_k^1} \left(\frac{1}{|\mathcal{D}_k^{2,i}|} \sum_{j \in \mathcal{D}_k^{2,i}} h_k(X_i, X_j) \right)^2 - \left(\frac{2}{|\mathcal{D}_k|} \sum_{(i,j) \in \mathcal{D}_k} h_k(X_i, X_j) \right)^2$$

Kernel pooling:

$$\text{pool}_{k \in K} \frac{S_k}{\sigma_k}$$

Kernel Pooling for Adaptivity

Mean kernel pooling:

$$\text{mean}_{k \in K} \frac{S_k}{\sigma_k} = \frac{1}{|K|} \sum_{k \in K} \frac{S_k}{\sigma_k} \qquad \frac{1}{|K|} \sum_{k \in K} S_k = S_{\frac{1}{|K|} \sum_{k \in K} k}$$

Max kernel pooling:

$$\max_{k \in K} \frac{S_k}{\sigma_k}$$

Fuse kernel pooling:

$$\text{fuse}_{k \in K} \frac{S_k}{\sigma_k} = \frac{1}{\nu} \log \left(\frac{1}{|K|} \sum_{k \in K} \exp \left(\nu \frac{S_k}{\sigma_k} \right) \right)$$
$$\max_{k \in K} \frac{S_k}{\sigma_k} - \frac{\log(|K|)}{\nu} \leq \text{fuse}_{k \in K} \frac{S_k}{\sigma_k} \leq \max_{k \in K} \frac{S_k}{\sigma_k}$$

Kernel Collection

Radial kernel with bandwidth λ :

$$k_\lambda(x, y) = f\left(\frac{\|x - y\|_r}{\lambda}\right)$$

Inter-sample distances:

$$D = \left\{ \|x - x'\|_r : x, x' \in \{X_1, \dots, X_n\} \right\} \setminus \{0\}$$

Bandwidth collection:

$$\Lambda(k, M) = \left\{ q_{5\%}^D + i(q_{95\%}^D - q_{5\%}^D)/M : i = 0, \dots, M \right\}$$

Kernel collection:

$$K = \left\{ k_\lambda : k \in \{\text{Gaussian, Laplace}\}, \lambda \in \Lambda(k, 10) \right\}$$

Part II

Optimal Kernel

Hypothesis Testing

Hypothesis Testing

Hypothesis Testing

Partitioned space of distributions: $\mathcal{P} = \mathcal{P}_0 \cup \mathcal{P}_1$

Hypothesis testing: Given i.i.d. samples from $P \in \mathcal{P}$ we test whether

$$\mathcal{H}_0: P \in \mathcal{P}_0 \quad \text{vs} \quad \mathcal{H}_1: P \in \mathcal{P}_1$$

Testing with a discrepancy:

$$\mathcal{H}_0: \text{Disc}(P) = 0 \quad \text{vs} \quad \mathcal{H}_1: \text{Disc}(P) > 0$$

Statistic: estimator using samples from P

$$\widehat{\text{Disc}}$$

Test output: reject \mathcal{H}_0 if $\widehat{\text{Disc}}$ is larger than some threshold

Hypothesis Testing: Level

Test output: reject \mathcal{H}_0 if $\widehat{\text{Disc}}$ is larger than some threshold

Level: type I error

$$\sup_{P_0 \in \mathcal{P}_0} \mathbb{P}_{P_0}(\text{reject } \mathcal{H}_0) \leq \alpha$$

Threshold: $(1 - \alpha)$ -quantile of the **simulated** distribution of $\widehat{\text{Disc}}$ under the null

Permutations: when \mathcal{H}_0 holds if and only if exchangeability holds

$$\widehat{\text{Disc}}(X_{\pi(1)}, \dots, X_{\pi(n)})$$

Exchangeability: $\text{joint}(X_1, \dots, X_n) = \text{joint}(X_{\pi(1)}, \dots, X_{\pi(n)})$ for all permutation π

Wild bootstrap: $\varepsilon_1, \dots, \varepsilon_n$ are i.i.d. Rademacher variables (± 1 with probability 0.5)

$$\frac{1}{|\mathcal{D}|} \sum_{(i,j) \in \mathcal{D}} \varepsilon_i \varepsilon_j h(X_i, X_j)$$

Hypothesis Testing: Power

Type II error:

$$\mathbb{P}_{P_1}(\text{fail to reject } \mathcal{H}_0), P_1 \in \mathcal{P}_1$$

Power:

$$\mathbb{P}_{P_1}(\text{reject } \mathcal{H}_0), P_1 \in \mathcal{P}_1$$

Impossible to control both types of errors uniformly over \mathcal{P}_0 and \mathcal{P}_1 :

$$\sup_{P_0 \in \mathcal{P}_0} \mathbb{P}_{P_0}(\text{reject } \mathcal{H}_0) \leq \alpha \qquad \sup_{P_1 \in \mathcal{P}_1} \mathbb{P}_{P_1}(\text{fail to reject } \mathcal{H}_0) \leq \beta$$

Pointwise power / consistency: for fixed $P_1 \in \mathcal{P}_1$

$$\lim_{N \rightarrow \infty} \mathbb{P}_{P_1}(\text{reject } \mathcal{H}_0) = 1$$

Hypothesis Testing: Power

Uniform power over $\mathcal{S}_1 \subset \mathcal{P}_1$:

$$\sup_{P_1 \in \mathcal{S}_1} \mathbb{P}_{P_1}(\text{reject } \mathcal{H}_0) \geq 1 - \beta \quad (\star)$$

Example of subset \mathcal{S}_1 :

$$\mathcal{S}_1 = \{P_1 \in \mathcal{P}_1 : \text{Disc}(P_1) \geq f(N, \alpha, \beta)\}$$

Upper bound: given a test, find smallest separation rate $f(N, \alpha, \beta)$ s.t. (\star) holds

Lower bound: find largest separation rate $f(N, \alpha, \beta)$ such that no level- α test satisfy (\star)

Minimax optimality: test for which upper bound matches the lower bound

Sobolev regularity: smoothness restriction on P_1

$$\int_{\mathbb{R}^d} \|\xi\|_2^{2s} |\hat{f}(\xi)|^2 d\xi \leq (2\pi)^d$$

Aggregation Kernel Adaptivity

Aggregation (Kernel Adaptivity)

Multiple testing:

reject \mathcal{H}_0 if one of the tests $T_{k_1}, \dots, T_{k_{|K|}}$ rejects

Bonferroni level correction: $\alpha / |K|$

$$\mathbb{P}(\text{reject } \mathcal{H}_0) = \mathbb{P}\left(\bigcup_{i=1}^{|K|} \{T_{k_i} \text{ rejects } \mathcal{H}_0\}\right) \leq \sum_{i=1}^{|K|} \mathbb{P}(T_{k_i} \text{ rejects } \mathcal{H}_0) = \alpha$$

Aggregation: estimate level correction between $\alpha / |K|$ and α such that

$$\mathbb{P}(\text{reject } \mathcal{H}_0) = \alpha$$

with bisection method using null-simulated samples via permutations or wild bootstrap

Kernel Collection

Radial kernel with bandwidth λ :

$$k_\lambda(x, y) = f\left(\frac{\|x - y\|_r}{\lambda}\right)$$

Inter-sample distances:

$$D = \left\{ \|x - x'\|_r : x, x' \in \{X_1, \dots, X_n\} \right\} \setminus \{0\}$$

Bandwidth collection:

$$\Lambda(k, M) = \left\{ q_{5\%}^D + i(q_{95\%}^D - q_{5\%}^D)/M : i = 0, \dots, M \right\}$$

Kernel collection:

$$K = \left\{ k_\lambda : k \in \{\text{Gaussian, Laplace}\}, \lambda \in \Lambda(k, 10) \right\}$$

Testing Constraints:

Efficiency

Privacy

Robustness

Computational Efficiency

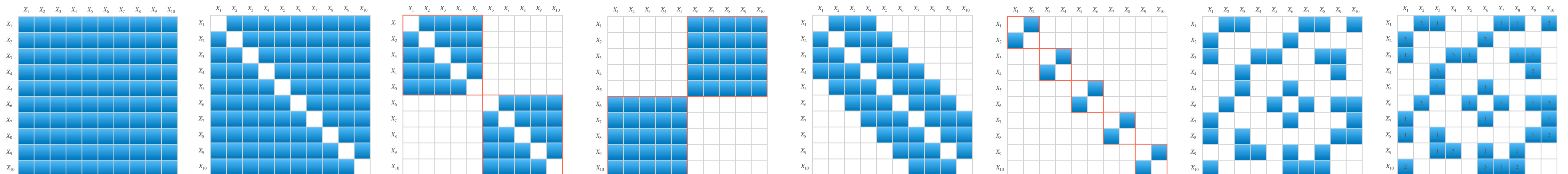
High computational cost: prohibitive for very large datasets

Goal: construct tests with lower computational complexity (e.g. sub-quadratic)

Trade-off: computational efficiency vs test power via uniform separation rates

Our approach: use incomplete statistics as efficient estimators

Others: Nyström approximation, random Fourier features, kernel thinning



Differential Privacy

(ε, δ) -differentially privacy of randomised test:

$$\mathbb{P}(\text{reject } \mathcal{H}_0 \text{ using } \mathbb{X}_N) \leq e^\varepsilon \mathbb{P}(\text{reject } \mathcal{H}_0 \text{ using } \widetilde{\mathbb{X}}_N) + \delta$$

$$\mathbb{P}(\text{fail to reject } \mathcal{H}_0 \text{ using } \mathbb{X}_N) \leq e^\varepsilon \mathbb{P}(\text{fail to reject } \mathcal{H}_0 \text{ using } \widetilde{\mathbb{X}}_N) + \delta$$

for any two datasets \mathbb{X}_N and $\widetilde{\mathbb{X}}_N$ differing only in one entry

Intuition: probability of a test output remains roughly the same when the data of a single user is modified

Privatisation: inject Laplacian noise to original and permuted statistics

$$\widehat{\text{MMD}}_{\pi_i} + \frac{2}{\xi} \frac{\sqrt{2}}{N} \zeta_i \quad \widehat{\text{HSIC}}_{\pi_i} + \frac{2}{\xi} \frac{4(N-1)}{N^2} \zeta_i \quad \xi = \varepsilon + \log \left(\frac{1}{1-\delta} \right)$$

Robustness

Standard testing: Given i.i.d. samples X_1, \dots, X_N from $P \in \mathcal{P}$, determine whether $\mathcal{H}_0: P \in \mathcal{P}_0$ or $\mathcal{H}_1: P \in \mathcal{P}_1$

Robust testing: Given samples X_1, \dots, X_N where

- $N - r$ samples are i.i.d. from $P \in \mathcal{P}$
- r samples have been corrupted (possibly in a non i.i.d. adversarial manner)

determine whether $\mathcal{H}_0: P \in \mathcal{P}_0$ or $\mathcal{H}_1: P \in \mathcal{P}_1$

Intuition: null hypothesis is enlarged with only $N - r$ samples being i.i.d. from P

Robustisation: add shift to permuted quantile (not to original statistic)

$$q_{1-\alpha}^{\text{MMD}} + 2\frac{\sqrt{2}}{N}r$$

$$q_{1-\alpha}^{\text{HSIC}} + 2\frac{4(N-1)}{N^2}r$$

Two-sample Testing

Two-sample testing

Two-sample problem: Given independent

- i.i.d. samples X_1, \dots, X_m from a distribution P ,
- i.i.d. samples Y_1, \dots, Y_n from a distribution Q ,

test whether $\mathcal{H}_0: P = Q$ or $\mathcal{H}_1: P \neq Q$. Let $N = \min(m, n)$.

General notation: $\mathcal{P}_0 = \{(P, Q) : P = Q\}$ and $\mathcal{P}_1 = \{(P, Q) : P \neq Q\}$

Testing with MMD: $\mathcal{H}_0: \text{MMD}(P, Q) = 0$ and $\mathcal{H}_1: \text{MMD}(P, Q) > 0$

Permutations (exchangeability): $Z_i = X_i$, $Z_{m+i} = Y_i$, $(\mathbb{X}_m^\pi, \mathbb{Y}_n^\pi) = (Z_{\pi(1)}, \dots, Z_{\pi(m+n)})$

Wild bootstrap: $\frac{1}{|\mathcal{D}|} \sum_{(i,j) \in \mathcal{D}} \varepsilon_i \varepsilon_j h(X_i, X_j)$, $\varepsilon_1, \dots, \varepsilon_n$ i.i.d. Rademacher variables

Pointwise power / consistency: $\lim_{N \rightarrow \infty} \mathbb{P}_{P,Q}(\text{reject } \mathcal{H}_0) = 1$ for fixed $P \neq Q$

Power in MMD Metric: Fixed Kernel

Standard testing:

$$\text{MMD}_k \gtrsim \sqrt{\frac{\max\{\log(1/\alpha), \log(1/\beta)\}}{N}}$$

Efficient testing:

$$\text{MMD}_k \gtrsim \sqrt{\frac{\max\{\log(1/\alpha), \log(1/\beta)\}}{N/B}}$$

Robust testing:

$$\text{MMD}_k \gtrsim \max\left\{\sqrt{\frac{\max\{\log(1/\alpha), \log(1/\beta)\}}{N}}, \frac{r}{N}\right\}$$

Differentially private testing:

$$\text{MMD}_k \gtrsim \max\left\{\sqrt{\frac{\max\{\log(1/\alpha), \log(1/\beta)\}}{N}}, \frac{\max\{\log(1/\alpha), \log(1/\beta)\}}{N\xi}\right\}$$

Power in MMD Metric: Kernel Pooling

Standard testing:

$$\max_{k \in K} \text{MMD}_k \gtrsim \sqrt{\frac{\max \{ \log(1/\alpha), \log(1/\beta), \log(|K|) \}}{N}}$$

Efficient testing:

$$\max_{k \in K} \text{MMD}_k \gtrsim \sqrt{\frac{\max \{ \log(1/\alpha), \log(1/\beta), \log(|K|) \}}{N/B}}$$

Robust testing:

$$\max_{k \in K} \text{MMD}_k \gtrsim \max \left\{ \sqrt{\frac{\max \{ \log(1/\alpha), \log(1/\beta), \log(|K|) \}}{N}}, \frac{r}{N} \right\}$$

Power in L2 Sobolev Metric: Oracle Kernel

Standard testing:

$$\|p - q\|_{L^2} \gtrsim \left(\frac{\log(1/\alpha) \log(1/\beta)}{N} \right)^{2s/(4s+d)}$$

Efficient testing:

$$\|p - q\|_{L^2} \gtrsim \left(\frac{\log(1/\alpha) \log(1/\beta)}{|\mathcal{D}|/N} \right)^{2s/(4s+d)}$$

Differentially private testing:

$$\|p - q\|_{L^2} \gtrsim N^{-2s/(4s+d)} \quad \text{in low privacy regime } \xi \gtrsim N^{-(2s-d/2)/(4s+d)}$$

$$\|p - q\|_{L^2} \gtrsim \left(N^{3/2} \xi \right)^{-s/(2s+d)} \quad \text{in mid privacy regime } N^{-1/2} \lesssim \xi \lesssim N^{-(2s-d/2)/(4s+d)}$$

$$\|p - q\|_{L^2} \gtrsim \left(N \xi \right)^{-2s/(2s+d)} \quad \text{in high privacy regime } \xi \lesssim N^{-1/2}$$

Power in L2 Sobolev Metric: Kernel Aggregation

Standard testing:

$$\|p - q\|_{L^2} \gtrsim \left(\frac{\log(1/\alpha) \log(1/\beta)}{N / \log(\log(N))} \right)^{2s/(4s+d)}$$

Efficient testing:

$$\|p - q\|_{L^2} \gtrsim \left(\frac{\log(1/\alpha) \log(1/\beta)}{(|\mathcal{D}|/N) / \log(\log(|\mathcal{D}|/N))} \right)^{2s/(4s+d)}$$

Independence Testing

Independence testing

Independence problem: Given

- paired samples $(X_1, Y_1), \dots, (X_N, Y_N)$ drawn i.i.d. from a joint distribution P_{XY} ,

test whether the first and second components of the pairs are independent

$$\mathcal{H}_0: P_{XY} = P_X \otimes P_Y \quad \text{vs} \quad \mathcal{H}_1: P_{XY} \neq P_X \otimes P_Y$$

General notation: $\mathcal{P}_0 = \{P_{XY} : P_{XY} = P_X \otimes P_Y\}$ and $\mathcal{P}_1 = \{P_{XY} : P_{XY} \neq P_X \otimes P_Y\}$

Testing with HSIC: $\mathcal{H}_0: \text{HSIC}(P_{XY}) = 0$ and $\mathcal{H}_1: \text{HSIC}(P_{XY}) > 0$

Permutations (exchangeability): $(X_1, Y_{\pi(1)}), \dots, (X_N, Y_{\pi(N)})$

Wild bootstrap: $\frac{1}{|\mathcal{D}|} \sum_{(i,j) \in \mathcal{D}} \varepsilon_i \varepsilon_j h(X_i, X_j), \quad \varepsilon_1, \dots, \varepsilon_n \text{ i.i.d. Rademacher variables}$

Pointwise power / consistency: $\lim_{N \rightarrow \infty} \mathbb{P}_{P,Q}(\text{reject } \mathcal{H}_0) = 1$ for fixed $P_{XY} \neq P_X \otimes P_Y$

Power in HSIC Metric: Fixed Kernels

Standard testing:

$$\text{HSIC}_{k,\ell} \gtrsim \sqrt{\frac{\max\{\log(1/\alpha), \log(1/\beta)\}}{N}}$$

Efficient testing:

$$\text{HSIC}_{k,\ell} \gtrsim \sqrt{\frac{\max\{\log(1/\alpha), \log(1/\beta)\}}{N/B}}$$

Robust testing:

$$\text{HSIC}_{k,\ell} \gtrsim \max\left\{\sqrt{\frac{\max\{\log(1/\alpha), \log(1/\beta)\}}{N}}, \frac{r}{N}\right\}$$

Differentially private testing:

$$\text{HSIC}_{k,\ell} \gtrsim \max\left\{\sqrt{\frac{\max\{\log(1/\alpha), \log(1/\beta)\}}{N}}, \frac{\max\{\log(1/\alpha), \log(1/\beta)\}}{N\xi}\right\}$$

Power in MMD Metric: Kernel Pooling

Standard testing:

$$\max_{k \in K} \max_{\ell \in L} \text{HSIC}_{k,\ell} \gtrsim \sqrt{\frac{\max \{ \log(1/\alpha), \log(1/\beta), \log(|K||L|) \}}{N}}$$

Efficient testing:

$$\max_{k \in K} \max_{\ell \in L} \text{HSIC}_{k,\ell} \gtrsim \sqrt{\frac{\max \{ \log(1/\alpha), \log(1/\beta), \log(|K||L|) \}}{N/B}}$$

Robust testing:

$$\max_{k \in K} \max_{\ell \in L} \text{HSIC}_{k,\ell} \gtrsim \max \left\{ \sqrt{\frac{\max \{ \log(1/\alpha), \log(1/\beta), \log(|K||L|) \}}{N}}, \frac{r}{N} \right\}$$

Power in L2 Sobolev Metric: Oracle Kernels

Standard testing:

$$\|p_{xy} - p_x \otimes p_y\|_{L^2} \gtrsim \left(\frac{\log(1/\alpha) \log(1/\beta)}{N} \right)^{2s/(4s+d)}$$

Efficient testing:

$$\|p_{xy} - p_x \otimes p_y\|_{L^2} \gtrsim \left(\frac{\log(1/\alpha) \log(1/\beta)}{|\mathcal{D}|/N} \right)^{2s/(4s+d)}$$

Differentially private testing:

$$\|p - q\|_{L^2} \gtrsim N^{-2s/(4s+d)} \quad \text{in low privacy regime } \xi \gtrsim N^{-(2s-d/2)/(4s+d)}$$

$$\|p - q\|_{L^2} \gtrsim \left(N^{3/2} \xi \right)^{-s/(2s+d)} \quad \text{in mid privacy regime } N^{-1/2} \lesssim \xi \lesssim N^{-(2s-d/2)/(4s+d)}$$

$$\|p - q\|_{L^2} \gtrsim \left(N \xi \right)^{-2s/(2s+d)} \quad \text{in high privacy regime } \xi \lesssim N^{-1/2}$$

Power in L2 Sobolev Metric: Kernel Aggregation

Standard testing:

$$\|p_{xy} - p_x \otimes p_y\|_{L^2} \gtrsim \left(\frac{\log(1/\alpha) \log(1/\beta)}{N / \log(\log(N))} \right)^{2s/(4s+d)}$$

Efficient testing:

$$\|p_{xy} - p_x \otimes p_y\|_{L^2} \gtrsim \left(\frac{\log(1/\alpha) \log(1/\beta)}{(|\mathcal{D}|/N) / \log(\log(|\mathcal{D}|/N))} \right)^{2s/(4s+d)}$$

Goodness-of-fit Testing

Goodness-of-fit testing

Goodness-of-fit problem: Given

- model distribution P (e.g. via its score function),
- i.i.d. samples X_1, \dots, X_N from a distribution Q ,

test whether $\mathcal{H}_0: P = Q$ or $\mathcal{H}_1: P \neq Q$.

General notation: $\mathcal{P}_0 = \{(P, Q) : P = Q\}$ and $\mathcal{P}_1 = \{(P, Q) : P \neq Q\}$

Testing with KSD: $\mathcal{H}_0: \text{KSD}_P(Q) = 0$ and $\mathcal{H}_1: \text{KSD}_P(Q) > 0$

Permutations do not work: goodness-of-fit does **not** test for exchangeability

Wild bootstrap: $\frac{1}{|\mathcal{D}|} \sum_{(i,j) \in \mathcal{D}} \varepsilon_i \varepsilon_j h(X_i, X_j), \quad \varepsilon_1, \dots, \varepsilon_n \text{ i.i.d. Rademacher variables}$

Pointwise power / consistency: $\lim_{N \rightarrow \infty} \mathbb{P}_Q(\text{reject } \mathcal{H}_0) = 1$ for fixed $P \neq Q$

Power in KSD Metric: Fixed Kernel

Standard testing:

$$\text{KSD}_k \gtrsim \sqrt{\frac{\max\{\log(1/\alpha), \log(1/\beta)\}}{N}}$$

Efficient testing:

$$\text{KSD}_k \gtrsim \sqrt{\frac{\max\{\log(1/\alpha), \log(1/\beta)\}}{N/B}}$$

Power in KSD Metric: Kernel Pooling

Standard testing:

$$\max_{k \in K} \text{KSD}_k \gtrsim \sqrt{\frac{\max\{\log(1/\alpha), \log(1/\beta), \log(|K|)\}}{N}}$$

Efficient testing:

$$\max_{k \in K} \text{KSD}_k \gtrsim \sqrt{\frac{\max\{\log(1/\alpha), \log(1/\beta), \log(|K|)\}}{N/B}}$$

Power in L2 Sobolev Metric: Oracle Kernel

Standard testing:

$$\|(\nabla \log p - \nabla \log q) q\|_{L^2} \gtrsim \left(\frac{\log(1/\alpha) \log(1/\beta)}{N} \right)^{2s/(4s+d)}$$

Efficient testing:

$$\|(\nabla \log p - \nabla \log q) q\|_{L^2} \gtrsim \left(\frac{\log(1/\alpha) \log(1/\beta)}{|\mathcal{D}|/N} \right)^{2s/(4s+d)}$$

Power in L2 Sobolev Metric: Kernel Aggregation

Standard testing:

$$\|(\nabla \log p - \nabla \log q) q\|_{L^2} \gtrsim \left(\frac{\log(1/\alpha) \log(1/\beta)}{N / \log(\log(N))} \right)^{2s/(4s+d)}$$

Efficient testing:

$$\|(\nabla \log p - \nabla \log q) q\|_{L^2} \gtrsim \left(\frac{\log(1/\alpha) \log(1/\beta)}{(|\mathcal{D}|/N) / \log(\log(|\mathcal{D}|/N))} \right)^{2s/(4s+d)}$$

Experiments

Experiments

Synthetic data:

- Perturbed uniform distributions
- Gaussian (mean shift, mixture, etc.)
- Gamma distribution
- Gaussian-Bernoulli Restricted Boltzmann Machine

Real-world data:

- MNIST $d = 28 \times 28 = 784$ (digit distribution, image shift, normalising flow model)
- CIFAR-10 $d = 32 \times 32 = 1024$ (image shift, CIFAR-10 vs CIFAR-10.1)
- Galaxy MNIST $d = 3 \times 64 \times 64 = 12288$ (types of galaxies)
- CelebA $d = 3 \times 178 \times 218 = 116412$ (men/women)
- IMDb movie reviews $d = 3330$ (sentiment analysis)

Level: under the null type I error is at level $\alpha = 0.05$

Power: outperform other state-of-the-art kernel hypothesis tests

Open Problems

Open Problems

- L2 separation upper bound for the HSIIC test with $s \in (0, (d_X + d_Y)/4)$
- L2 separation lower bounds under differential privacy constraint
- L2 and kernel separation lower bounds under efficiency constraint
- L2 separation upper and lower bounds under the robustness constraint
- L2 separation lower bounds for goodness-of-fit testing
- Kernel separation rates for normalised pooling
- Kernel pooling and aggregation procedure for differential private tests
- KSD private and robust test constructions and uniform power separation

Any Questions?

Proof Strategies

Concentration results

Exponential concentration for statistic:

$$\mathbb{P}\left(\left|\widehat{\text{Kdisc}} - \text{Kdisc}\right| > t\right) \leq \exp(-C t^2 N)$$

$$\left|\widehat{\text{Kdisc}} - \text{Kdisc}\right| \lesssim \sqrt{\frac{1}{N} \log\left(\frac{1}{\beta}\right)} \quad \text{with probability } \geq 1 - \beta/2$$

Exponential concentration for bootstrapped statistic (quantile):

$$\mathbb{P}\left(\widehat{\text{Kdisc}}_{\text{boot}} > t\right) \leq \exp(-C t^2 N)$$

$$q_{1-\alpha} \lesssim \sqrt{\frac{1}{N} \log\left(\frac{1}{\alpha}\right)} \quad \text{with probability } \geq 1 - \beta/2$$

Proof Separation in Kernel Metric

Type II error control:

$$\begin{aligned} & \mathbb{P}(\widehat{\text{Kdisc}} \leq q_{1-\alpha}) \\ & \leq \mathbb{P}\left(\text{Kdisc} \leq q_{1-\alpha} + C_1 \sqrt{\frac{1}{N} \log\left(\frac{1}{\beta}\right)}\right) + \frac{\beta}{2} \\ & \leq \mathbb{P}\left(\text{Kdisc} \leq C_2 \sqrt{\frac{1}{N} \log\left(\frac{1}{\alpha}\right)} + C_1 \sqrt{\frac{1}{N} \log\left(\frac{1}{\beta}\right)}\right) + \beta \\ & = \beta \end{aligned}$$

when there is kernel separation:

$$\text{Kdisc} \gtrsim \sqrt{\frac{\max\{\log(1/\alpha), \log(1/\beta)\}}{N}}$$

L2 Testing

L2 testing: $\mathcal{H}_0: \|\psi\|_{L^2} = 0$ vs $\mathcal{H}_1: \|\psi\|_{L^2} > 0$

Two-sample: $\psi = p - q$

Independence: $\psi = p_{xy} - p_x \otimes p_y$

Goodness-of-fit: $\psi = (\nabla \log p - \nabla \log q) q$

Kernel Integral Transform: $(S_\lambda f)(y) = \int_{\mathbb{R}^d} k_\lambda(x, y) f(x) dx$

Linking Kdisc and L2:

$$\text{Kdisc} = \langle \psi, S_\lambda \psi \rangle_{L^2} = \frac{1}{2} (\|\psi\|_{L^2}^2 + \|S_\lambda \psi\|_{L^2}^2 - \|\psi - S_\lambda \psi\|_{L^2}^2)$$

MMD example:

$$\langle \psi, S_\lambda \psi \rangle_{L^2} = \int_{\mathbb{R}^d} \int_{\mathbb{R}^d} k_\lambda(x, y) (p(x) - q(x)) (p(y) - q(y)) dx dy = \text{MMD}_\lambda^2(P, Q)$$

Concentrations with bandwidth dependency

Exponential concentration for statistic:

$$\left| \widehat{\text{Kdisc}} - \text{Kdisc} \right| \lesssim \frac{1}{2} \|S_\lambda \psi\|_{L^2}^2 + \frac{1}{N} \log \left(\frac{1}{\beta} \right)$$

with probability $\geq 1 - \beta/2$

Exponential concentration for bootstrapped statistic (quantile):

$$q_{1-\alpha} \lesssim \frac{1}{N\lambda^{d/2}} \log \left(\frac{1}{\alpha} \right) \log \left(\frac{1}{\beta} \right)$$

with probability $\geq 1 - \beta/2$

Proof Separation in L2 Metric

Type II error control:

$$\begin{aligned} & \mathbb{P}(\widehat{\text{Kdisc}} \leq q_{1-\alpha}) \\ & \leq \mathbb{P}\left(\text{Kdisc} \leq \frac{1}{2}\|S_\lambda\psi\|_{L^2}^2 + C\frac{1}{N\lambda^{d/2}}\log\left(\frac{1}{\alpha}\right)\log\left(\frac{1}{\beta}\right)\right) + \beta \\ & = \beta \end{aligned}$$

when there is kernel separation:

$$\text{Kdisc} \geq \frac{1}{2}\|S_\lambda\psi\|_{L^2}^2 + \frac{C}{N\lambda^{d/2}}\log\left(\frac{1}{\alpha}\right)\log\left(\frac{1}{\beta}\right)$$

equivalently for L2 separation: e.g. use $\text{Kdisc} = \frac{1}{2}(\|\psi\|_{L^2}^2 + \|S_\lambda\psi\|_{L^2}^2 - \|\psi - S_\lambda\psi\|_{L^2}^2)$

$$\|\psi\|_{L^2}^2 \geq \|\psi - S_\lambda\psi\|_{L^2}^2 + \frac{C}{N\lambda^{d/2}}\log\left(\frac{1}{\alpha}\right)\log\left(\frac{1}{\beta}\right)$$

Proof Separation in L2 Metric

Power is guaranteed when:

$$\|\psi\|_{L^2}^2 \geq \|\psi - S_\lambda \psi\|_{L^2}^2 + \frac{C}{N\lambda^{d/2}} \log\left(\frac{1}{\alpha}\right) \log\left(\frac{1}{\beta}\right)$$

Smoothness: Assuming ψ is s -Sobolev smooth: $\int_{\mathbb{R}^d} \|\xi\|_2^{2s} |\widehat{\psi}(\xi)|^2 d\xi \leq (2\pi)^d$

$$\|\psi - S_\lambda \psi\|_{L^2}^2 \lesssim \lambda^{2s}$$

Optimal bandwidth: $\lambda = (\log(1/\alpha)\log(1/\beta)/N)^{2/(4s+d)}$

Power is guaranteed when:

$$\|\psi\|_{L^2} \gtrsim \left(\frac{\log(1/\alpha)\log(1/\beta)}{N} \right)^{2s/(4s+d)}$$

Proofs under testing constraints

Computational efficiency:

Keep track of design size $|\mathcal{D}|$

Differential privacy:

Keep track of privacy parameters ϵ and δ

Robustness:

Keep track of robustness parameter r

Kernel Pooling

Asymptotic Validity

Kernel Pooling Asymptotic Validity (KSD)

Known: $\widehat{\text{Kdisc}}_{\text{wild}}^{(k)} \xrightarrow[n \rightarrow \infty]{D} \widehat{\text{Kdisc}}^{(k)}$ for any fixed kernel k

Cramér–Wold Theorem: a sequence of d -dimensional variables $Z^{(n)} \xrightarrow{D} Z = (Z_1, \dots, Z_d)$ if and only if

$$\sum_{i=1}^d t_i Z_i^{(n)} \xrightarrow[n \rightarrow \infty]{D} \sum_{i=1}^d t_i Z_i \quad \text{for every } (t_1, \dots, t_d) \in \mathbb{R}^d$$

Our setting:

$$\begin{aligned} \sum_{i=1}^K t_i \widehat{\text{Kdisc}}_{\text{wild}}^{(k_i)} &= \sum_{i=1}^K t_i \frac{1}{N(N-1)} \sum_{1 \leq j \neq j' \leq N} \epsilon_j \epsilon_{j'} k_i(X_j, X_{j'}) = \frac{1}{N(N-1)} \sum_{1 \leq j \neq j' \leq N} \epsilon_j \epsilon_{j'} \left(\sum_{i=1}^K t_i k_i(X_j, X_{j'}) \right) \\ &= \widehat{\text{Kdisc}}_{\text{wild}}(\sum_{i=1}^K t_i k_i) \xrightarrow[N \rightarrow \infty]{D} \widehat{\text{Kdisc}}(\sum_{i=1}^K t_i k_i) = \sum_{i=1}^K t_i \widehat{\text{Kdisc}}^{(k_i)} \end{aligned}$$

Cramér–Wold Theorem: $\text{joint}_{k \in K} \left(\widehat{\text{Kdisc}}_{\text{wild}}^{(k)} \right) \xrightarrow[N \rightarrow \infty]{D} \text{joint}_{k \in K} \left(\widehat{\text{Kdisc}}^{(k)} \right)$

Continuous Mapping Theorem: for any continuous pooling function (e.g. mean, max, fuse)

$$\text{pool}_{k \in K} \left(\widehat{\text{Kdisc}}_{\text{wild}}^{(k)} \right) \xrightarrow[N \rightarrow \infty]{D} \text{pool}_{k \in K} \left(\widehat{\text{Kdisc}}^{(k)} \right)$$

Any Questions?