

# **Semi-Supervised Kernel Two-Sample Test**

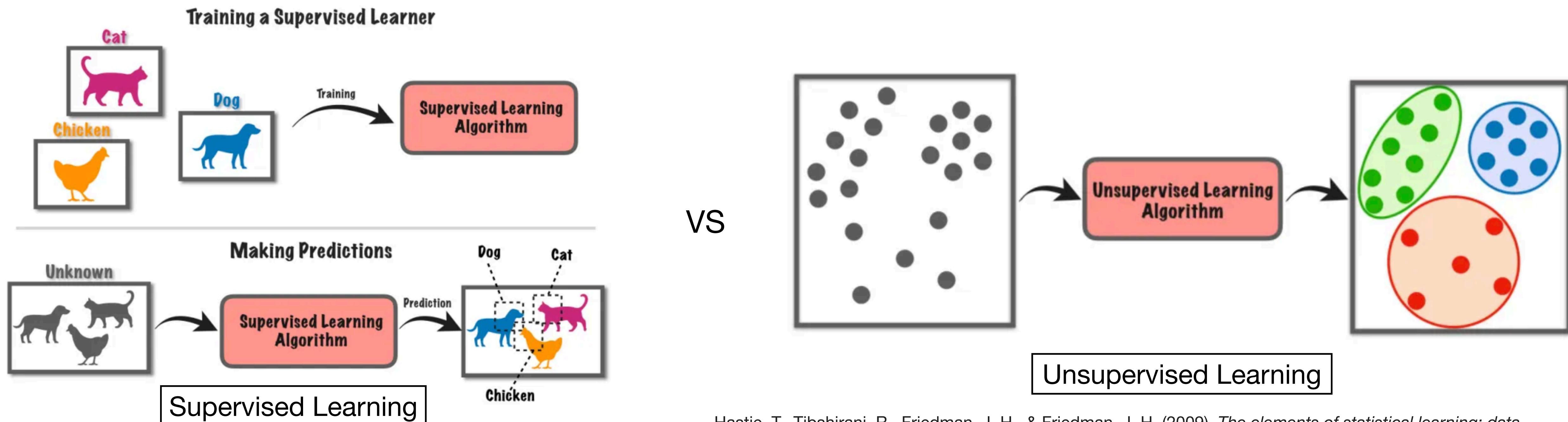
Statistics and Data Analysis  
M.S.  
Gyumin Lee

# Contents

1. Semi-Supervised Inference
2. Kernel Two-Sample Test
3. Semi-Supervised Kernel Two-Sample Test
4. Conclusion

# Semi-Supervised Inference

- Classic machine learning : supervised vs unsupervised
  - We know the outcome variables to guide the learning process?
  - Or, we observe only the features and have no measurements of the outcome?

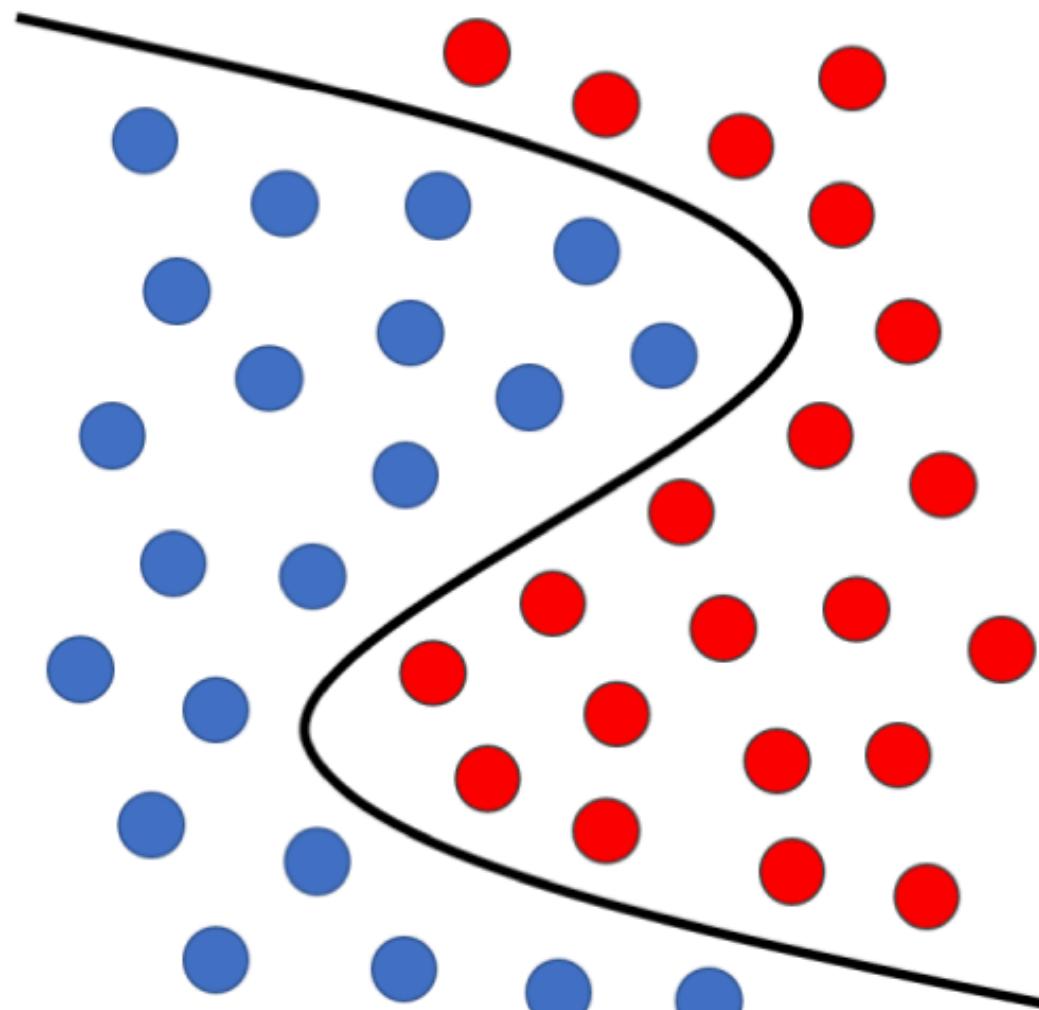


- Hastie, T., Tibshirani, R., Friedman, J. H., & Friedman, J. H. (2009). *The elements of statistical learning: data mining, inference, and prediction* (Vol. 2, pp. 1-758). New York: Springer.
- <https://towardsdatascience.com/supervised-vs-unsupervised-learning-in-2-minutes-72dad148f242>

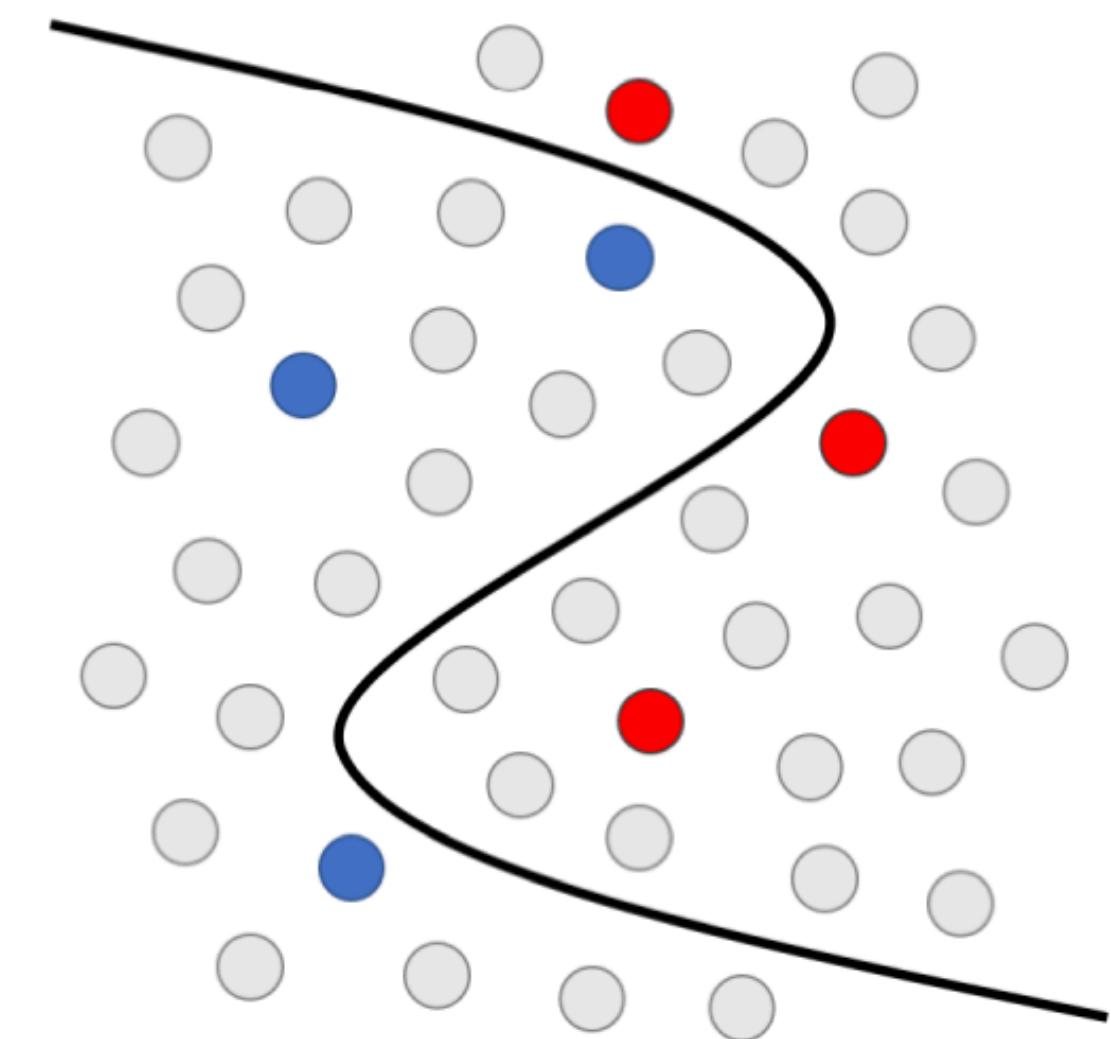
# Semi-Supervised Inference

Then, what is ‘Semi-supervised learning/setting’?

: problem of learning/setting based on a small labeled dataset together with a large unlabeled dataset.



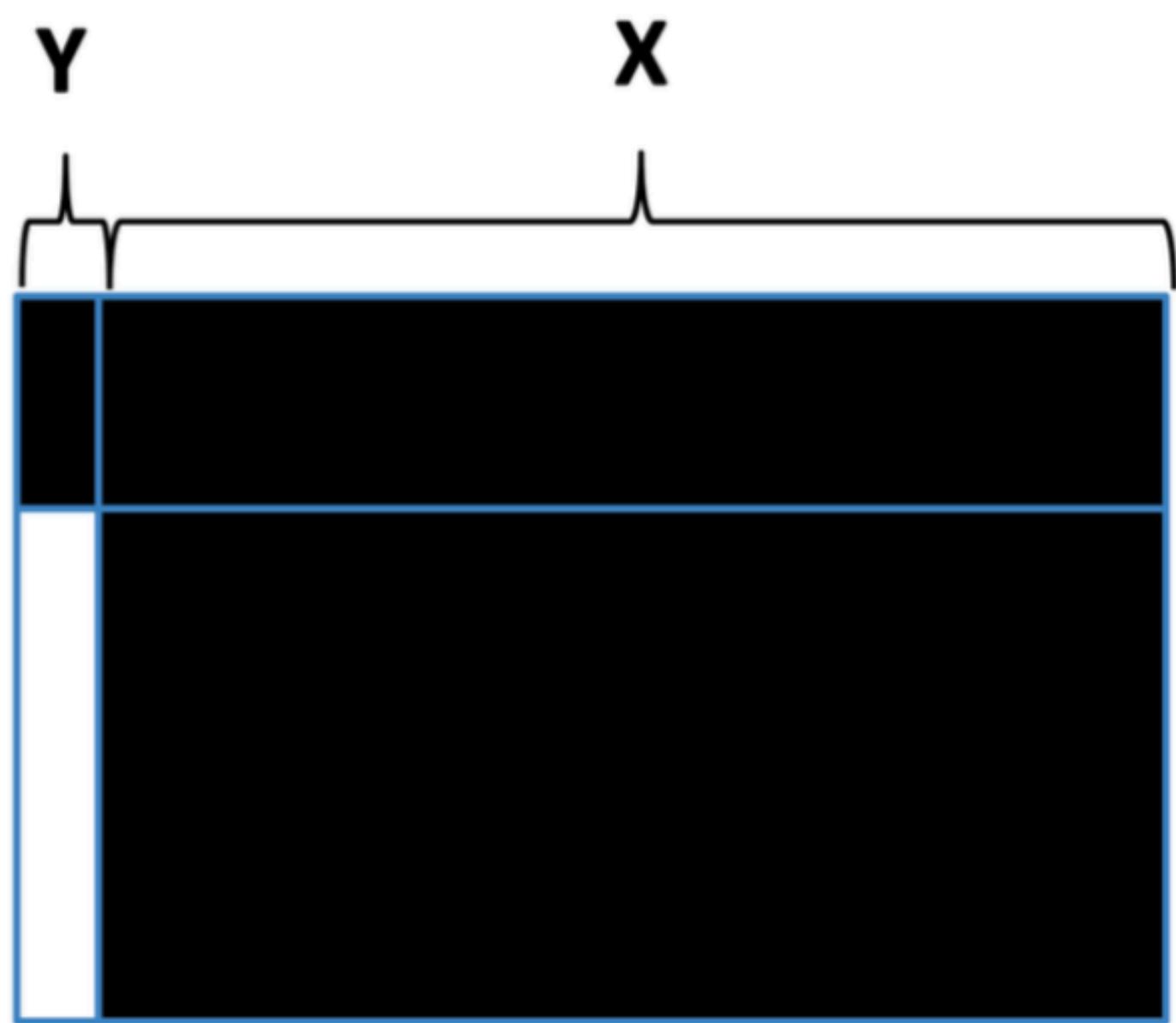
(a) Supervised Learning



(b) Semi-Supervised Learning

**Labeled Samples**

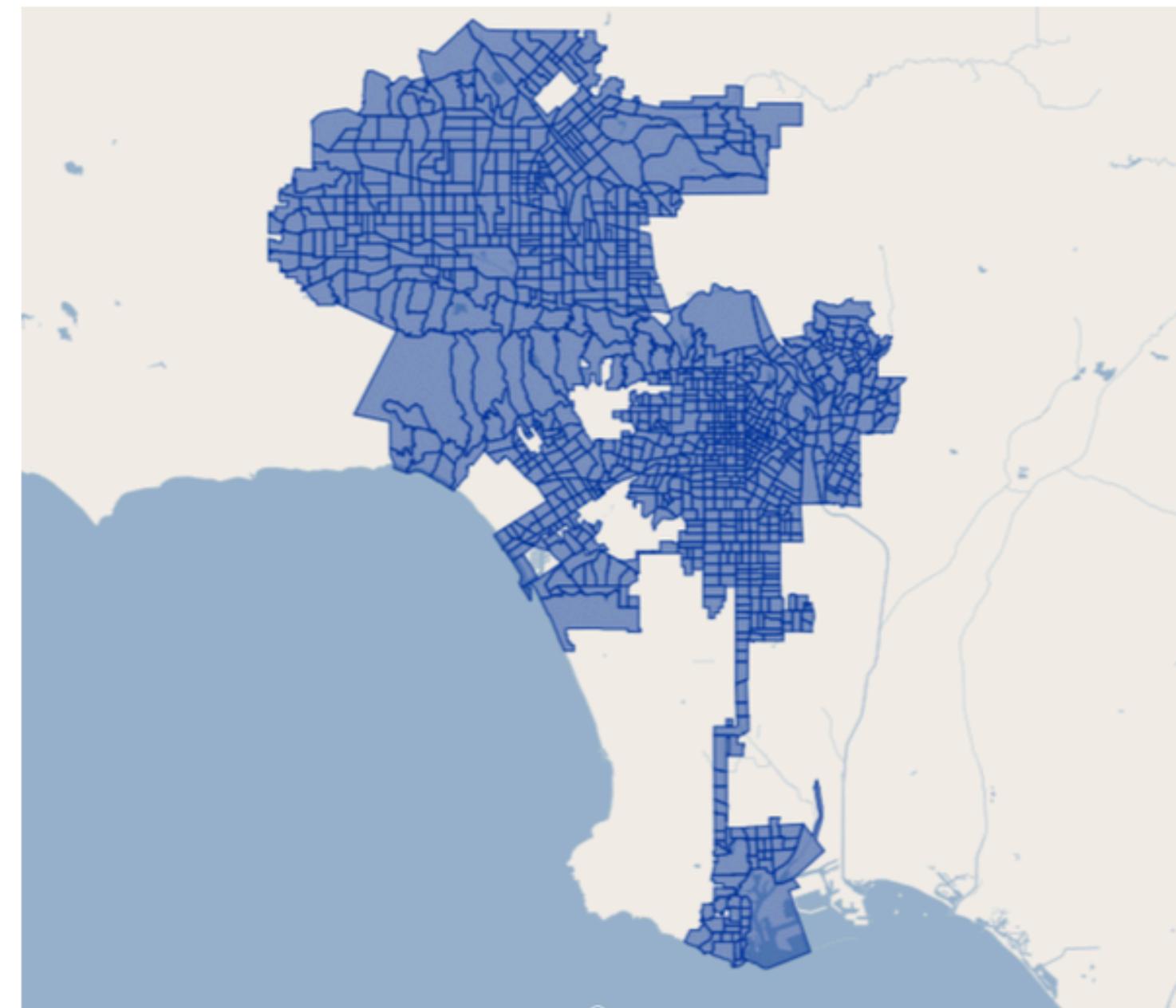
**Unlabeled Samples**



- Jiao, R., Zhang, Y., Ding, L., Xue, B., Zhang, J., Cai, R., & Jin, C. (2023). Learning with limited annotations: a survey on deep semi-supervised learning for medical image segmentation. *Computers in Biology and Medicine*, 107840.
- Zhang, A., Brown, L. D., & Cai, T. T. (2019). Semi-supervised inference: General theory and estimation of means.

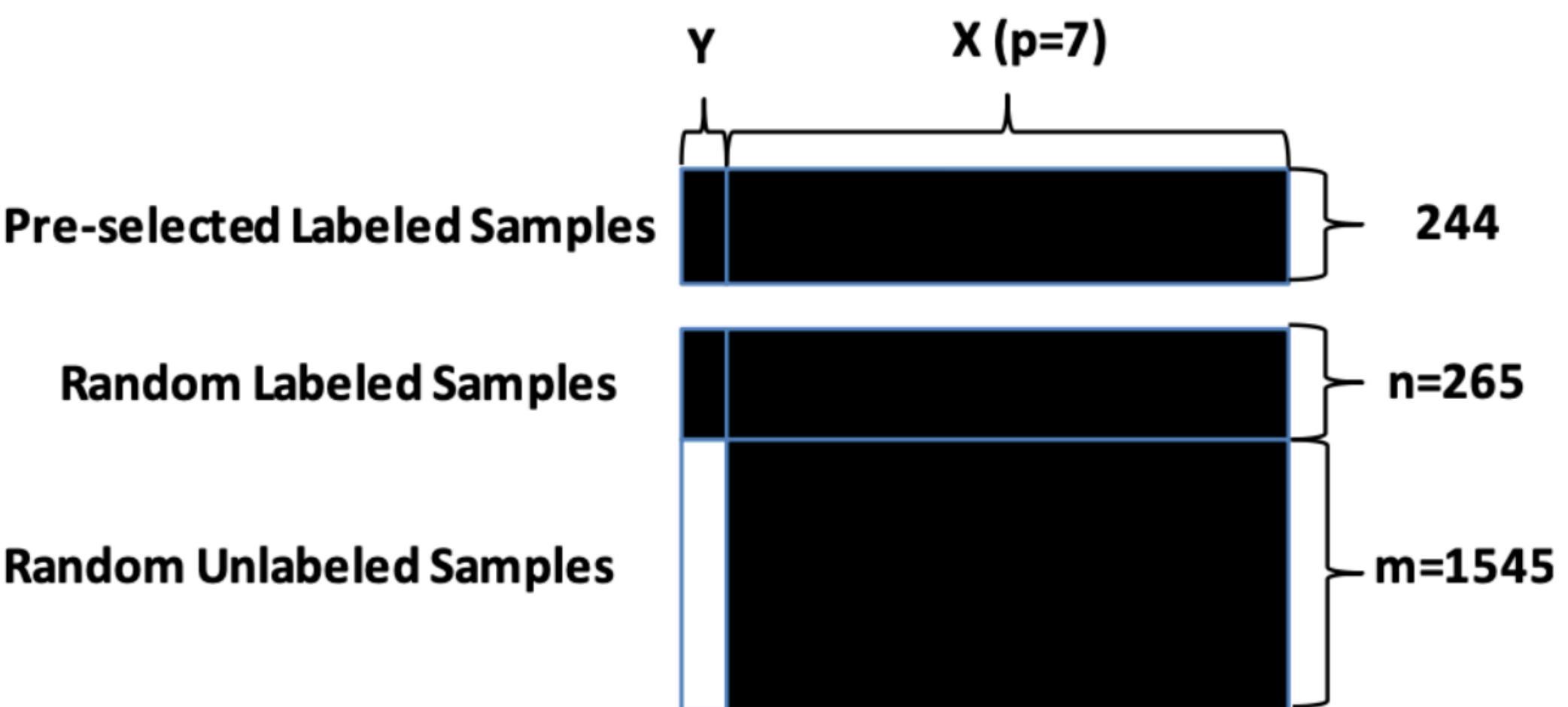
# Semi-Supervised Inference

- Semi-supervised settings often appear across various fields handling regression / classification tasks.
- It plays an important role in many problems in statistics and machine learning, including model fine-tuning, model distillation, self-training, transfer learning and continual learning



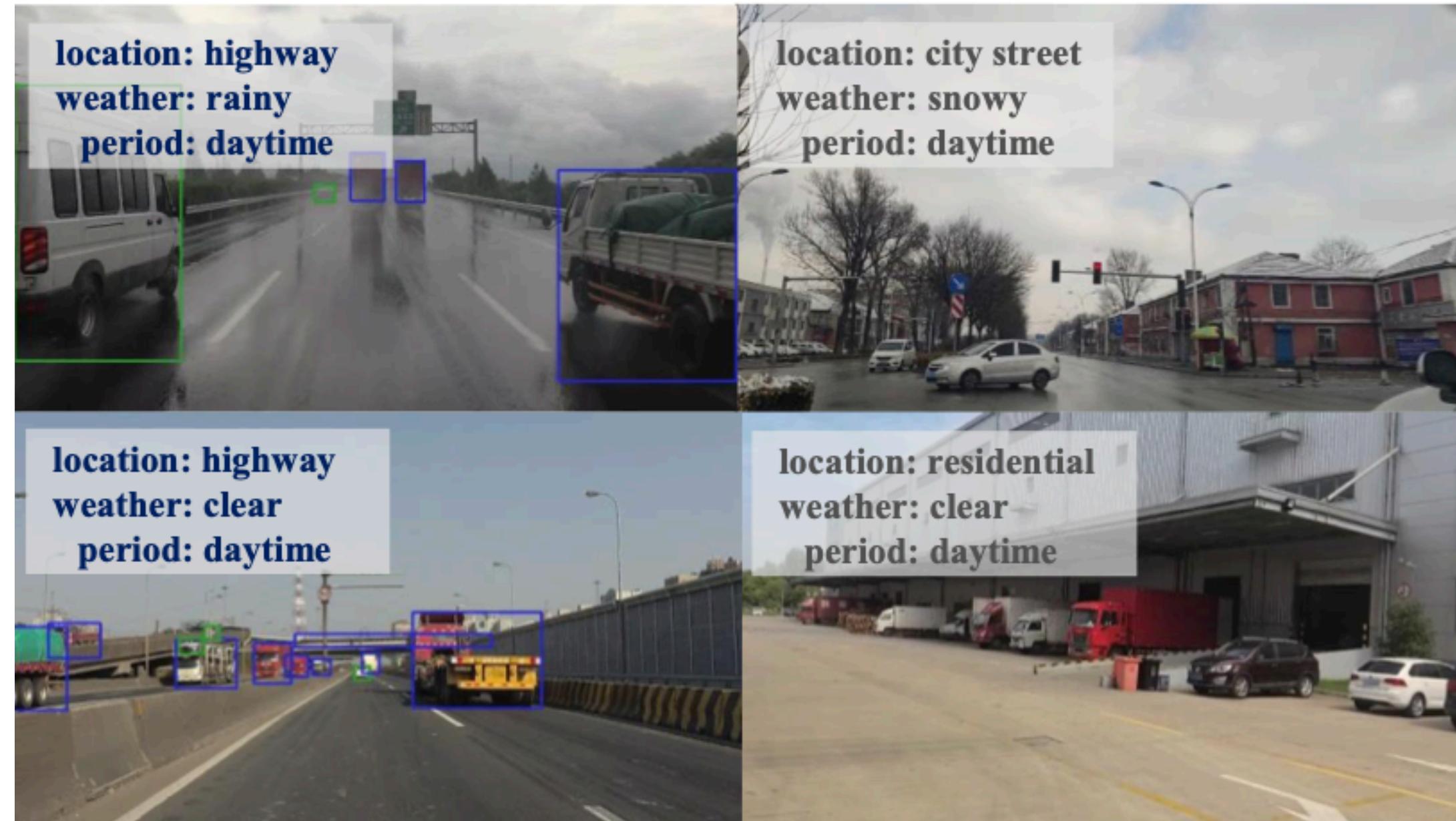
## Ex) Estimating Homeless in Los Angeles County

- inputs: 7 demographic covariates
- labels: the number of homeless people
- 265 labeled data set / 1545 unlabeled data set

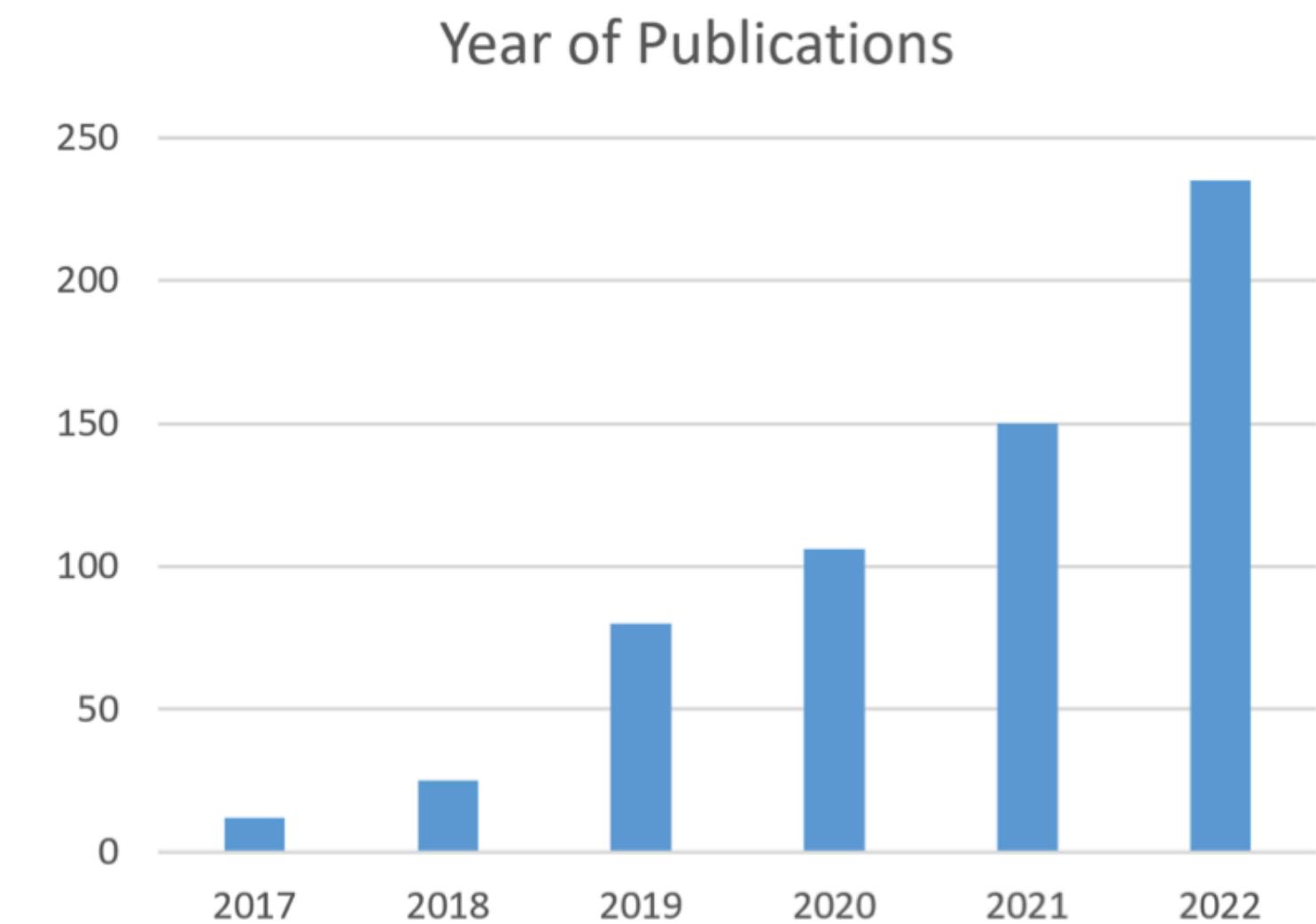
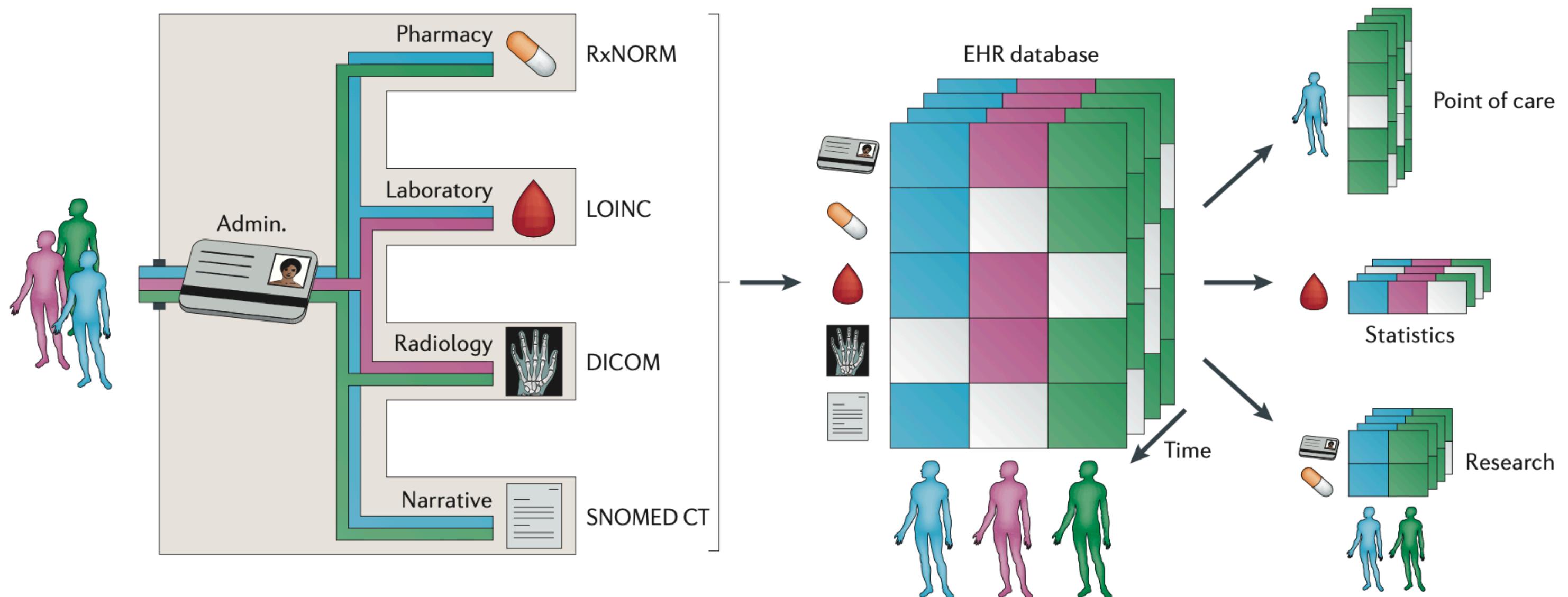


# Semi-Supervised Inference

Ex) Object detection for autonomous driving  
- 20K labeled images, 10M unlabeled images



Ex) Medical image segmentation  
Ex) Electronic health records



- Han, J., Liang, X., Xu, H., Chen, K., Hong, L., Mao, J., ... & Xu, C. (2021). SODA10M: A large-scale 2D self/semi-supervised object detection dataset for autonomous driving. *arXiv preprint arXiv:2106.11118*.
- Jensen, P. B., Jensen, L. J., & Brunak, S. (2012). Mining electronic health records: towards better research applications and clinical care. *Nature Reviews Genetics*, 13(6), 395-405.
- Jiao, R., Zhang, Y., Ding, L., Xue, B., Zhang, J., Cai, R., & Jin, C. (2023). Learning with limited annotations: a survey on deep semi-supervised learning for medical image segmentation. *Computers in Biology and Medicine*, 107840.

# Semi-Supervised Inference

- Basic Setting

$(Y, X_1, X_2, \dots, X_p) : (p + 1)$  dimensional random vector following an unknown joint distribution  $P = P(dy, dx_1, \dots, dx_p)$

$P_X$  : marginal distribution of  $X = (X_1, X_2, \dots, X_p)$

Often, we suppose that there are  $n$  labeled samples from  $P$ ,

$$[\mathbf{Y}, \mathbf{X}] = \{Y_k, X_{k1}, X_{k2}, \dots, X_{kp}\}_{k=1}^n$$

and  $m$  unlabeled samples from  $P_X$ ,

$$\mathbf{X}_{\text{add}} = \{X_{k1}, X_{k2}, \dots, X_{kp}\}_{k=n+1}^{n+m}$$

# Semi-Supervised Inference

- Often, our goal is to estimate some estimand

$$\theta^* = \arg \min_{\theta \in \mathbb{R}^p} \mathbb{E} [\ell_\theta(X_1, Y_1)]$$

which could be applied for many different inference problems

- mean estimation

$$\theta^* = \arg \min_{\theta \in \mathbb{R}} \mathbb{E}[\ell_\theta(Y_1)] = \arg \min_{\theta \in \mathbb{R}} \mathbb{E} \left[ \frac{1}{2}(Y_1 - \theta)^2 \right]$$

- quantile estimation

$$\theta^* = \min \{ \theta : P(Y_1 \leq \theta) \geq q \}$$

- inference for linear regression

$$\theta^* = \arg \min_{\theta \in \mathbb{R}^d} \mathbb{E}[\ell_\theta(X_1, Y_1)] = \arg \min_{\theta \in \mathbb{R}^d} \mathbb{E}[(Y_1 - X_1^\top \theta)^2]$$

and so on...

# Semi-Supervised Inference

- Simple example from ‘Mathematical Statistics 1’

Let our goal of estimation be  $\theta^* = E[Y]$

We have several candidates of estimators (  $f$  is a prediction function )

- sample mean  $\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$  vs semi-supervised mean  $\hat{\theta} = \frac{1}{n} \sum_{i=1}^n (Y_i - f(X_i)) + \frac{1}{n+m} \sum_{i=1}^{n+m} f(X_i)$
- Both are unbiased, so we compare variance...

$$\text{Var}(\bar{Y}) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n Y_i\right) = \frac{1}{n} \text{Var}(Y)$$

$$\begin{aligned}\text{Var}(\hat{\theta}) &= \text{Var}\left(\frac{1}{n} \sum_{i=1}^n (Y_i - f(X_i)) + \frac{1}{n+m} \sum_{i=1}^{n+m} f(X_i)\right) \\ &= \frac{1}{n}(\text{Var}(Y) + \frac{m-n}{m}(\text{Var}(f(X)) - 2\text{Cov}(Y, f(X))))\end{aligned}$$

If  $\text{Var}(f(X)) < 2\text{Cov}(Y, f(X))$ , then semi-supervised mean improves the original one!

# Semi-Supervised Inference

- Semi-supervised mean estimation in regression setting (Zhang, A. et al., 2019)
  - For the ideal semi-supervised inference, where  $m = \infty$  we have ‘least square estimator’  $\hat{\theta}_{LS} = \boldsymbol{\mu}'\hat{\beta} = \hat{\beta}_1 + \boldsymbol{\mu}'\hat{\beta}_{(2)} = \bar{Y} - \hat{\beta}'_{(2)}(\bar{X} - \boldsymbol{\mu})$ .
  - For the ordinary semi-supervised inference, where  $m < \infty$  we have ‘semi-supervised least squared estimator’  $\hat{\theta}_{SSLS} = \hat{\boldsymbol{\mu}}'\hat{\beta} = \bar{Y} - \hat{\beta}'_{(2)}(\bar{X} - \hat{\boldsymbol{\mu}})$ .
  - Both have asymptotic normality

$$\left| \begin{array}{l} \frac{\hat{\theta}_{LS} - \theta}{\tau/\sqrt{n}} \xrightarrow{d} N(0, 1) \\ \frac{\sqrt{n}(\hat{\theta}_{SSLS} - \theta)}{\nu} \xrightarrow{d} N(0, 1) \end{array} \right| \quad \left| \begin{array}{l} MSE/\tau^2 \xrightarrow{d} 1, \quad \text{where } MSE := \frac{\sum_{i=1}^n (\vec{Y}_i - \vec{X}_i^\top \hat{\beta})^2}{n-p-1}. \\ \frac{\hat{\nu}^2}{\nu^2} \xrightarrow{d} 1 \quad \text{where } \hat{\nu}^2 = \frac{m}{m+n} MSE + \frac{n}{m+n} \hat{\sigma}_Y^2 \text{ with } MSE = \frac{1}{n-p-1} \sum_{k=1}^n (\vec{Y}_i - \vec{X}_k^\top \hat{\beta})^2 \text{ and} \\ \hat{\sigma}_Y^2 = \frac{1}{n-1} \sum_{k=1}^n (Y_i - \bar{Y})^2. \end{array} \right.$$

-Zhang, A., Brown, L. D., & Cai, T. T. (2019). Semi-supervised inference: General theory and estimation of means.

# Semi-Supervised Inference

- Semi-supervised mean estimation in regression setting (Zhang, A. et al., 2019)
  - Comparison with respect to  $\ell_2$  Risk

$$n\mathbb{E} \left( \hat{\theta}_{\text{LS}}^1 - \theta \right)^2 = \tau_n^2 + s_n \quad \text{where} \quad s_n = \frac{p^2}{n} A_{n,p} + \frac{p^2}{n^{5/4}} B_{n,p}, \quad \max(|A_{n,p}|, |B_{n,p}|) \leq C$$

$$n\mathbb{E} \left( \hat{\theta}_{\text{SSLS}}^1 - \theta \right)^2 = \tau_n^2 + \frac{n}{n+m} \beta_{(2),n}^\top \Sigma_n \beta_{(2),n} + s_{n,m} \quad \text{where} \quad |s_{n,m}| \leq \frac{Cp^2}{n}.$$

$$\Rightarrow n\mathbb{E} \left( \hat{\theta}_{\text{SSLS}} - \theta \right)^2 = \frac{n}{n+m} n\mathbb{E}(\bar{\mathbf{Y}} - \theta)^2 + \frac{m}{n+m} n\mathbb{E}(\hat{\theta}_{\text{LS}} - \theta)^2 \quad (\text{ } n\mathbb{E}(\bar{\mathbf{Y}} - \theta)^2 = \tau^2 + \beta_{(2)}^\top \Sigma \beta_{(2)})$$

- Other semi-supervised inference works are still ongoing...
  - High-dimensional & mean(Zhang, Y. et al., 2022), high-dimensional & variance (Tony Cai, T. et al., 2020), robustness(Zhu, B. et al., 2024), general framework(Angelopoulos, A. N. et al., 2023), etc

- Zhang, Y., & Bradic, J. (2022). High-dimensional semi-supervised learning: in search of optimal inference of the mean. *Biometrika*, 109(2), 387-403.
- Zhu, B., Ding, M., Jacobson, P., Wu, M., Zhan, W., Jordan, M., & Jiao, J. (2024). Doubly-Robust Self-Training. *Advances in Neural Information Processing Systems*, 36.
- Angelopoulos, A. N., Bates, S., Fannjiang, C., Jordan, M. I., & Zrnic, T. (2023). Prediction-powered inference. *Science*, 382(6671), 669-674.
- Tony Cai, T., & Guo, Z. (2020). Semisupervised inference for explained variance in high dimensional linear regression and its applications. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 82(2), 391-419.

# Kernel Two-Sample Test

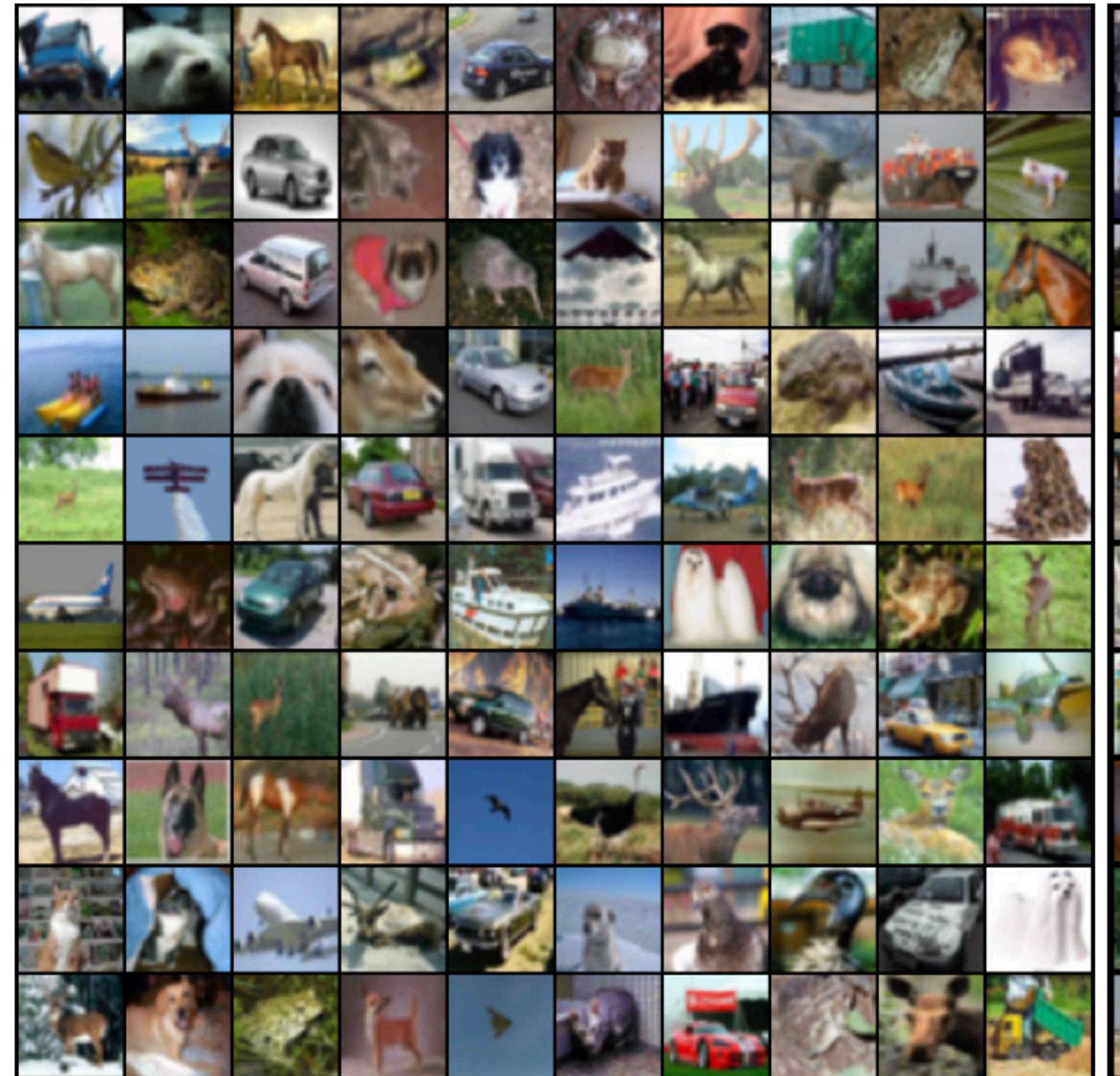
- Two-sample test:  $H_0: \mu_1 - \mu_2 = d_0, H_1: \mu_1 - \mu_2 \neq d_0,$

$\mu_1 - \mu_2 = d_0$	$z = \frac{(\bar{x}_1 - \bar{x}_2) - d_0}{\sqrt{\sigma_1^2/n_1 + \sigma_2^2/n_2}};$ $\sigma_1$ and $\sigma_2$ known	$\mu_1 - \mu_2 < d_0$ $\mu_1 - \mu_2 > d_0$ $\mu_1 - \mu_2 \neq d_0$	$z < -z_\alpha$ $z > z_\alpha$ $z < -z_{\alpha/2}$ or $z > z_{\alpha/2}$
$\mu_1 - \mu_2 = d_0$	$t = \frac{(\bar{x}_1 - \bar{x}_2) - d_0}{s_p \sqrt{1/n_1 + 1/n_2}};$ $v = n_1 + n_2 - 2,$ $\sigma_1 = \sigma_2$ but unknown, $s_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}$	$\mu_1 - \mu_2 < d_0$ $\mu_1 - \mu_2 > d_0$ $\mu_1 - \mu_2 \neq d_0$	$t < -t_\alpha$ $t > t_\alpha$ $t < -t_{\alpha/2}$ or $t > t_{\alpha/2}$
$\mu_1 - \mu_2 = d_0$	$t' = \frac{(\bar{x}_1 - \bar{x}_2) - d_0}{\sqrt{s_1^2/n_1 + s_2^2/n_2}};$ $v = \frac{(s_1^2/n_1 + s_2^2/n_2)^2}{\frac{(s_1^2/n_1)^2}{n_1-1} + \frac{(s_2^2/n_2)^2}{n_2-1}},$ $\sigma_1 \neq \sigma_2$ and unknown	$\mu_1 - \mu_2 < d_0$ $\mu_1 - \mu_2 > d_0$ $\mu_1 - \mu_2 \neq d_0$	$t' < -t_\alpha$ $t' > t_\alpha$ $t' < -t_{\alpha/2}$ or $t' > t_{\alpha/2}$
$\mu_D = d_0$ paired observations	$t = \frac{\bar{d} - d_0}{s_d / \sqrt{n}};$ $v = n - 1$	$\mu_D < d_0$ $\mu_D > d_0$ $\mu_D \neq d_0$	$t < -t_\alpha$ $t > t_\alpha$ $t < -t_{\alpha/2}$ or $t > t_{\alpha/2}$

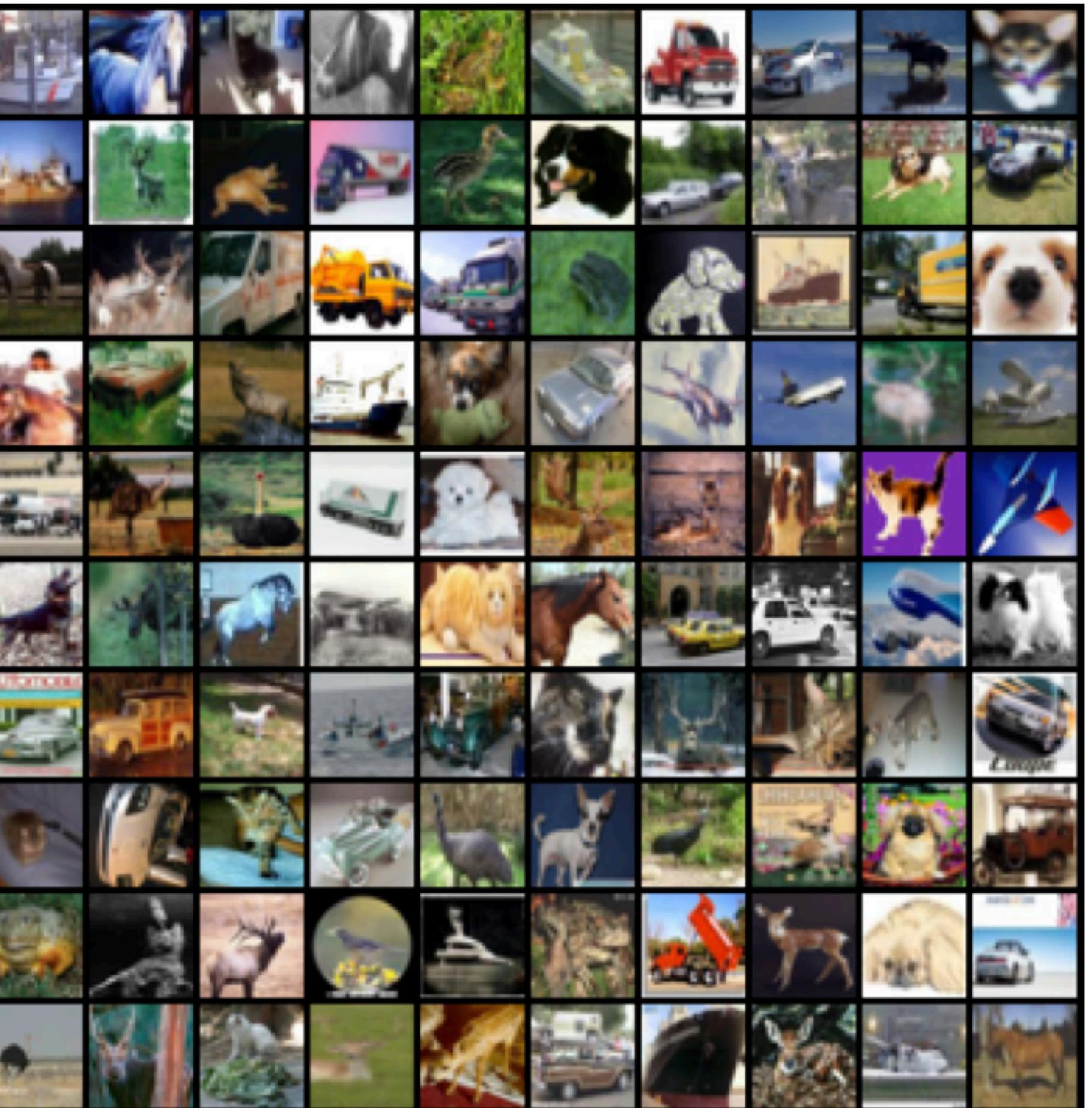
- Walpole, R. E., Myers, R. H., Myers, S. L., & Ye, K. (1993). *Probability and statistics for engineers and scientists* (Vol. 5, pp. 326-332). New York: Macmillan.

# Kernel Two-Sample Test

- Two-sample test:  $H_0 : P_X = P_Y$  vs.  $H_1 : P_X \neq P_Y$



(a) *CIFAR-10* test set



(b) *CIFAR-10.1* test set

$X_1$ : Now disturbing reports out of Newfoundland show that the fragile snow crab industry is in serious decline. First the west coast salmon, the east coast salmon and the cod, and now the snow crabs off Newfoundland.

$X_2$ : To my pleasant surprise he responded that he had personally visited those wharves and that he had already announced money to fix them. What wharves did the minister visit in my riding and how much additional funding is he going to provide for Delaps Cove, Hampton, Port Lorne,

...

?

$$P_X = Q_Y$$

$Y_1$ : Honourable senators, I have a question for the Leader of the Government in the Senate with regard to the support funding to farmers that has been announced. Most farmers have not received any money yet.

$Y_2$ : On the grain transportation system we have had the Estey report and the Kroeger report. We could go on and on. Recently programs have been announced over and over by the government such as money for the disaster in agriculture on the prairies and across Canada.

...

# Kernel Two-Sample Test

- f-divergence

$$D_f(p||q) = \int q(\mathbf{x}) f\left(\frac{p(\mathbf{x})}{q(\mathbf{x})}\right) d\mathbf{x}$$

$$\text{Ex)} D_{\text{KL}}(p \parallel q) = \int p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})} d\mathbf{x}$$

$$D_\alpha^A(p||q) \triangleq \frac{4}{1-\alpha^2} \left( 1 - \int p(\mathbf{x})^{(1+\alpha)/2} q(\mathbf{x})^{(1-\alpha)/2} d\mathbf{x} \right)$$

$$D_H^2(p||q) \triangleq \frac{1}{2} \int \left( p(\mathbf{x})^{\frac{1}{2}} - q(\mathbf{x})^{\frac{1}{2}} \right)^2 d\mathbf{x} = 1 - \int \sqrt{p(\mathbf{x})q(\mathbf{x})} d\mathbf{x}$$

$$\chi^2(p, q) \triangleq \frac{1}{2} \int \frac{(q(\mathbf{x}) - p(\mathbf{x}))^2}{q(\mathbf{x})} d\mathbf{x}$$

Computing divergence in terms of *ratio*

- Integral probability metrics

$$D_{\mathcal{F}}(P, Q) \triangleq \sup_{f \in \mathcal{F}} |\mathbb{E}_{p(\mathbf{x})} [f(\mathbf{x})] - \mathbb{E}_{q(\mathbf{x}')} [f(\mathbf{x}')]|$$

$$\text{Ex)} \mathcal{F} = \{ \|f\|_L \leq 1 \}, \quad \|f\|_L = \sup_{\mathbf{x} \neq \mathbf{x}'} \frac{|f(\mathbf{x}) - f(\mathbf{x}')|}{\|\mathbf{x} - \mathbf{x}'\|}$$

⇒ Wasserstein-1 distance

$$W_1(P, Q) \triangleq \sup_{\|f\|_L \leq 1} |\mathbb{E}_{p(\mathbf{x})} [f(\mathbf{x})] - \mathbb{E}_{q(\mathbf{x}')} [f(\mathbf{x}')]|$$

Computing divergence in terms of *difference*

Maximum Mean Discrepancy (MMD):  $\text{MMD}(P, Q; \mathcal{F}) = \sup_{\substack{f \in \mathcal{F}: \|f\|_L \leq 1 \\ \mathcal{F} \text{ is an RKHS}}} [\mathbb{E}_{p(\mathbf{x})} [f(\mathbf{x})] - \mathbb{E}_{q(\mathbf{x}')} [f(\mathbf{x}')]]$

# Kernel Two-Sample Test

- Reproducing Kernel Hilbert Space (RKHS)

: Hilbert space where the evaluation functional  $L_x(f) = f(x)$  is bounded.

## Riesz Representation Theorem

: For every bounded linear functional  $L$  on a Hilbert space  $\mathcal{H}$ , there exists a unique  $\xi_L \in \mathcal{H}$  such that  $L(f) = (\xi_L, f), \forall f \in \mathcal{H}$ .

Or, we say there exists **kernel function**  $\mathcal{K} : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$  which satisfies reproducing property  $\langle f(\cdot), \mathcal{K}(\cdot, \mathbf{x}') \rangle = f(\mathbf{x}')$  ( $\langle \mathcal{K}(\mathbf{x}, \cdot), \mathcal{K}(\cdot, \mathbf{x}') \rangle = \mathcal{K}(\mathbf{x}, \mathbf{x}')$ )

- Some good properties

: Reproducing kernel is positive definite ( $\sum_{i=1}^n \sum_{j=1}^n a_i a_j K(x_i, x_j) \geq 0$  for every  $n$ , and every  $x_1, \dots, x_n \in \mathcal{X}$ , and every  $a_1, \dots, a_n \in \mathbb{R}$ )

: (Moore-Aronszajn) Suppose  $K$  is a symmetric p.d. kernel. There is a unique RKHS for which  $K$  is a reproducing kernel.

- Yongho, Jeon (2023), Lecture notes for nonparametric function estimation

- Murphy, K. P. (2018). Machine learning: A probabilistic perspective (adaptive computation and machine learning series). The MIT Press: London, UK.

- Gretton, A., Borgwardt, K. M., Rasch, M. J., Schölkopf, B., & Smola, A. (2012). A kernel two-sample test. *The Journal of Machine Learning Research*, 13(1), 723-773.

Feature map  $\varphi(x) \in \mathcal{F}$ ,  $\varphi(x) = [\dots \varphi_i(x) \dots] \in \ell_2$

## Kernel Two-Sample Test

For positive definite  $k$ ,  $k(x, x') = \langle \varphi(x), \varphi(x') \rangle_{\mathcal{F}}$

- Expectations of RKHS functions

Given  $P$  a Borel probability measure on  $\mathcal{X}$ , we define feature map of probability  $P$ , or mean embedding of  $P$  as

$$\mu_P = [\dots \mathbf{E}_P [\varphi_i(X)] \dots] \quad \text{which satisfies } \mathbf{E}_P(f(X)) = \langle f, \mu_P \rangle_{\mathcal{F}}$$

Then, for positive definite  $k(x, x')$ ,

$$\langle \mu_P, \mu_Q \rangle_{\mathcal{F}} = \mathbf{E}_{P,Q} k(\mathbf{x}, \mathbf{y})$$

Does it exist? Yes!

If  $\mathbf{E}_P \sqrt{k(\mathbf{x}, \mathbf{x})} = \mathbf{E}_P \|\varphi(\mathbf{x})\|_{\mathcal{F}} < \infty$  then  $\exists \mu_P \in \mathcal{F}$ .  
 $f(\mathbf{x}) = \langle f, \varphi_{\mathbf{x}} \rangle_{\mathcal{F}}$

Then we obtain its feature map as  $\mu_P(t) = \langle \mu_P, \varphi(t) \rangle_{\mathcal{F}} = \langle \mu_P, k(\cdot, t) \rangle_{\mathcal{F}} = \mathbf{E}_{x \sim P} k(\mathbf{x}, t)$

And MMD is the distance between feature maps

$$\begin{aligned} MMD^2(P, Q) &= \|\mu_P - \mu_Q\|_{\mathcal{F}}^2 \\ &= \underbrace{\mathbf{E}_P k(\mathbf{x}, \mathbf{x}')}_{(a)} + \underbrace{\mathbf{E}_Q k(\mathbf{y}, \mathbf{y}')}_{(a)} - 2 \underbrace{\mathbf{E}_{P,Q} k(\mathbf{x}, \mathbf{y})}_{(b)} \end{aligned}$$

# Kernel Two-Sample Test

- Two-sample test using MMD

We obtain witness functions as

$$f^*(t) \propto \langle \phi(t), \mu_p - \mu_q \rangle_{\mathcal{H}} = \mathbf{E}_x [k(x, t)] - \mathbf{E}_y [k(y, t)],$$

So we take empirical witness function

$$\hat{f}^*(t) \propto \langle \phi(t), \mu_X - \mu_Y \rangle_{\mathcal{H}} = \frac{1}{m} \sum_{i=1}^m k(x_i, t) - \frac{1}{n} \sum_{i=1}^n k(y_i, t).$$

using empirical mean embedding as  $\mu_X := \frac{1}{m} \sum_{i=1}^m \phi(x_i)$  and  $\mu_Y := \frac{1}{n} \sum_{i=1}^n \phi(y_i)$

# Kernel Two-Sample Test

- Two-sample test using MMD

We use unbiased empirical estimator of MMD as

$$\text{MMD}_u^2[\mathcal{F}, X, Y] = \frac{1}{m(m-1)} \sum_{i=1}^m \sum_{j \neq i}^m k(x_i, x_j) + \frac{1}{n(n-1)} \sum_{i=1}^n \sum_{j \neq i}^n k(y_i, y_j) - \frac{2}{mn} \sum_{i=1}^m \sum_{j=1}^n k(x_i, y_j).$$

with asymptotic distribution under the null as

$$t\text{MMD}_u^2[\mathcal{F}, X, Y] \xrightarrow{D} \sum_{l=1}^{\infty} \lambda_l \left[ (\rho_x^{-1/2} a_l - \rho_y^{-1/2} b_l)^2 - (\rho_x \rho_y)^{-1} \right], \quad a_l \sim \mathcal{N}(0, 1) \text{ and } b_l \sim \mathcal{N}(0, 1) \quad \lambda_i \text{ are eigenvalues of}$$

and under the alternative as

$$\int_X \tilde{k}(x, x') \psi_i(x) dp(x) = \lambda_i \psi_i(x').$$

$$m^{\frac{1}{2}} (\text{MMD}_u^2 - \text{MMD}^2[\mathcal{F}, p, q]) \xrightarrow{D} \mathcal{N}(0, \sigma_u^2) \quad \text{where } \sigma_u^2 = 4 \left( \mathbf{E}_z \left[ (\mathbf{E}_{z'} h(z, z'))^2 \right] - [\mathbf{E}_{z, z'} (h(z, z'))]^2 \right)$$

$\Rightarrow$  degenerate U-statistic under the null (intractable limiting distribution!)

So, we use  $(1-\alpha)$  quantile using permutation test.

# Kernel Two-Sample Test

- Permutation-free kernel two-sample test

Instead, we use sample-splitting and studentization to handle this issue.

Original:  $\widehat{\text{MMD}}^2 := \frac{1}{n(n-1)m(m-1)} \sum_{1 \leq i \neq i' \leq n} \sum_{1 \leq j \neq j' \leq m} h(X_i, X_{i'}, Y_j, Y_{j'})$ ,  $h(x, x', y, y') := k(x, x') - k(x, y') - k(y, x') + k(y, y')$

Permutation-free:  $\bar{x}\widehat{\text{MMD}}^2 := x\widehat{\text{MMD}}^2 / \hat{\sigma}$

where  $x\widehat{\text{MMD}}^2 := \frac{1}{n_1 m_1 n_2 m_2} \sum_{i=1}^{n_1} \sum_{i'=1}^{n_2} \sum_{j=1}^{m_1} \sum_{j'=1}^{m_2} h(X_i, X_{i'}, Y_j, Y_{j'}) = \frac{1}{n_1} \sum_{i=1}^{n_1} U_{X,i} - \frac{1}{m_1} \sum_{j=1}^{m_1} U_{Y,j}$ .

$$\hat{\sigma}_X^2 := \frac{1}{n_1} \sum_{i=1}^{n_1} (U_{X,i} - \bar{U}_X)^2, \quad \hat{\sigma}_Y^2 := \frac{1}{m_1} \sum_{j=1}^{m_1} (U_{Y,j} - \bar{U}_Y)^2 \quad \text{and} \quad \hat{\sigma}^2 := \frac{1}{n_1} \hat{\sigma}_X^2 + \frac{1}{m_1} \hat{\sigma}_Y^2.$$

$$U_{X,i} := \langle k(X_i, \cdot), \hat{\mu}_2 - \hat{\nu}_2 \rangle_k \quad U_{Y,j} := \langle k(Y_j, \cdot), \hat{\mu}_2 - \hat{\nu}_2 \rangle_k$$

# Kernel Two-Sample Test

- Permutation-free kernel two-sample test

- Asymptotic normality under the null

: If  $\lim_{n \rightarrow \infty} \frac{\mathbb{E}_P[\bar{k}_n(X_1, X_2)^4]}{\mathbb{E}_P[\bar{k}_n(X_1, X_2)^2]^2} \left( \frac{1}{n} + \frac{1}{m_n} \right) = 0$ , and  $\lim_{n \rightarrow \infty} \frac{\lambda_{1,n}^2}{\sum_{l=1}^{\infty} \lambda_{l,n}^2}$  exists,  
 then  $\widehat{\text{xMMD}}^2 \xrightarrow{d} N(0, 1)$ .

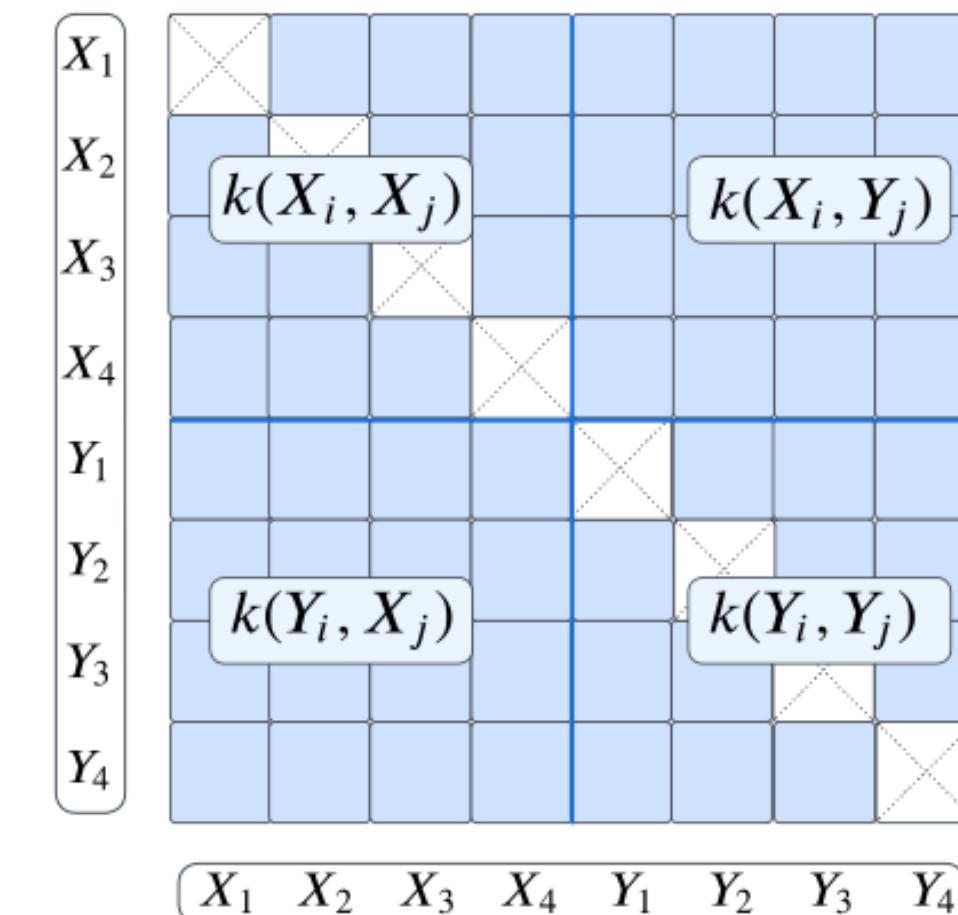
- Power consistency

: If  $0 < \mathbb{E}_P[\bar{k}(X_1, X_2)^4] < \infty$ , and  $0 < \mathbb{E}_Q[\bar{k}(Y_1, Y_2)^4] < \infty$ , then it has asymptotic power of 1.

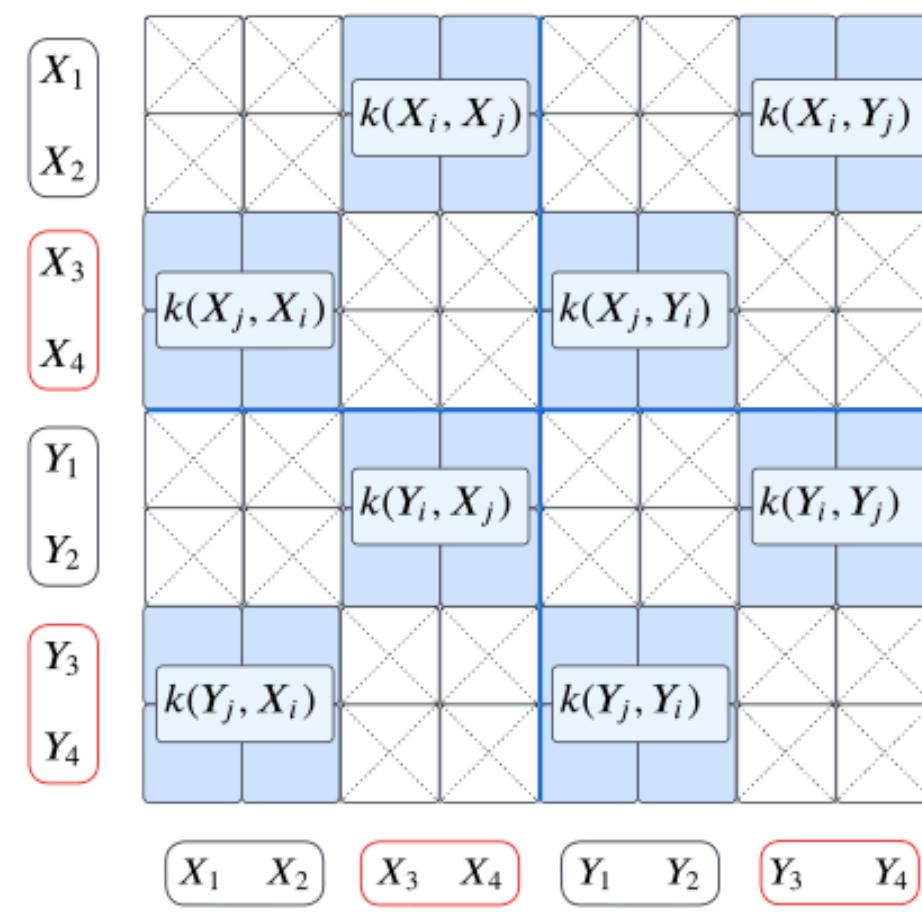
- Minimax rate optimality:  $\lim_{n \rightarrow \infty} \sup_{P_n \in \mathcal{P}_n^{(0)}} \mathbb{E}_{P_n}[\Psi(\mathbb{X}, \mathbb{Y})] \leq \alpha$  and  $\lim_{n \rightarrow \infty} \inf_{(P_n, Q_n) \in \mathcal{P}_n^{(1)}} \mathbb{E}_{P_n, Q_n}[\Psi(\mathbb{X}, \mathbb{Y})] = 1$ .

with  $s_n \asymp n^{4/(d+4\beta)}$ , gaussian kernel and both densities lie in  $\beta$ -Sobolev ball.

Kernel-MMD



Cross-MMD



# Semi-Supervised Kernel Two-Sample Test

- Setting

Now we have samples of both labeled and unlabeled data from two distributions

$$\begin{aligned} \mathcal{D}_{XV,1} &:= \{(X_i, V_i)\}_{i=1}^{n_1} := \{(X_i, V_i)\}_{i \in \mathcal{I}_{XV,1}} \stackrel{\text{iid}}{\sim} P_{XV}, \\ \mathcal{D}_V &:= \{V_i\}_{i=n_1+1}^{n_1+m_1} := \{V_i\}_{i \in \mathcal{I}_V} \stackrel{\text{iid}}{\sim} P_V, \\ \mathcal{D}_{YW,1} &:= \{(Y_i, W_i)\}_{i=1}^{n_2} := \{(Y_i, W_i)\}_{i \in \mathcal{I}_{YW,1}} \stackrel{\text{iid}}{\sim} P_{YW}, \\ \mathcal{D}_W &:= \{W_i\}_{i=n_2+1}^{n_2+m_2} := \{W_i\}_{i \in \mathcal{I}_W} \stackrel{\text{iid}}{\sim} P_W. \end{aligned} \quad \left. \begin{array}{l} \mathcal{D}_{XV,1} \cup \mathcal{D}_V \text{ of size } n_1 + m_1 \\ \mathcal{D}_{YW,1} \cup \mathcal{D}_W \text{ of size } n_2 + m_2 \end{array} \right\}$$

And our goal is to test

$$H_0 : P_X = P_Y \quad \text{vs.} \quad H_1 : P_X \neq P_Y$$

# Semi-Supervised Kernel Two-Sample Test

- Idea

First, estimate the witness function

$$\hat{f}(\cdot) = \frac{1}{n_1} \sum_{i \in \mathcal{I}_{XV,2}} k(X_i, \cdot) - \frac{1}{n_2} \sum_{i \in \mathcal{I}_{YW,2}} k(Y_i, \cdot)$$

From  $\hat{f}(X_1), \dots, \hat{f}(X_{n_1})$  and  $\hat{f}(Y_1), \dots, \hat{f}(Y_{n_2})$ , we compute the oracle estimator

$$\begin{aligned}\hat{\mu}_{X,\hat{f}} &= \frac{1}{n_1} \sum_{i=1}^{n_1} \{\hat{f}(X_i) - \mathbb{E}[\hat{f}(X_i) | V_i, \hat{f}]\} + \frac{1}{n_1 + m_1} \sum_{i=1}^{n_1+m_1} \mathbb{E}[\hat{f}(X_i) | V_i, \hat{f}] \quad \text{and} \\ \hat{\mu}_{Y,\hat{f}} &= \frac{1}{n_2} \sum_{i=1}^{n_2} \{\hat{f}(Y_i) - \mathbb{E}[\hat{f}(Y_i) | W_i, \hat{f}]\} + \frac{1}{n_2 + m_2} \sum_{i=1}^{n_2+m_2} \mathbb{E}[\hat{f}(Y_i) | W_i, \hat{f}].\end{aligned}$$

Then, we may estimate conditional expectations further.

# Semi-Supervised Kernel Two-Sample Test

- Idea
  - $\hat{\mu}_{X,\hat{f}} - \hat{\mu}_{Y,\hat{f}}$  : unbiased estimator of MMD
  - Since  $\text{Var}(\hat{\mu}_{X,\hat{f}} | \hat{f}) := \sigma_{X,\hat{f}}^2 := \frac{1}{n_1} \underbrace{\mathbb{E}[\text{Var}\{\hat{f}(X) | V, \hat{f}\} | \hat{f}]}_{=\sigma_{1,X}^2} + \frac{1}{n_1 + m_1} \underbrace{\text{Var}[\mathbb{E}\{\hat{f}(X) | V, \hat{f}\} | \hat{f}]}_{=\sigma_{2,X}^2}$

we estimate this as

$$\hat{\sigma}_{X,\hat{f}}^2 = \frac{1}{n_1} \hat{\sigma}_{1,X}^2 + \frac{1}{n_1 + m_1} \hat{\sigma}_{2,X}^2 \quad \text{where} \quad \hat{\sigma}_{1,X}^2 = \frac{1}{n_1} \sum_{i=1}^{n_1} \{ \hat{f}(X_i) - \mathbb{E}[\hat{f}(X_i) | V_i, \hat{f}] \}^2 \quad \text{and}$$
$$\hat{\sigma}_{2,X}^2 = \frac{1}{n_1 + m_1} \sum_{i=1}^{n_1+m_1} \left\{ \mathbb{E}[\hat{f}(X_i) | V_i, \hat{f}] - \frac{1}{n_1 + m_1} \sum_{j=1}^{n_1+m_1} \mathbb{E}[\hat{f}(X_j) | V_j, \hat{f}] \right\}^2$$

Then, our oracle test statistic becomes  $T^* = \frac{\hat{\mu}_{X,\hat{f}} - \hat{\mu}_{Y,\hat{f}}}{\sqrt{\hat{\sigma}_{X,\hat{f}}^2 + \hat{\sigma}_{Y,\hat{f}}^2}}$  where we reject the null if  $T^* > z_{1-\alpha} := \Phi^{-1}(1 - \alpha)$

# Semi-Supervised Kernel Two-Sample Test

- Asymptotic normality under the null

1) Test statistic with known variance  $T_\sigma^* = \frac{\hat{\mu}_{X,\hat{f}} - \hat{\mu}_{Y,\hat{f}}}{\sqrt{\sigma_{X,\hat{f}}^2 + \sigma_{Y,\hat{f}}^2}}$

Write numerator as  $\hat{\mu}_{X,\hat{f}} - \hat{\mu}_{Y,\hat{f}} = \sum_{i=1}^{n_1+m_1} G_i - \sum_{i=1}^{n_2+m_2} H_i$  where

$$G_i = \begin{cases} \frac{1}{n_1} \{ \hat{f}(X_i) - \mathbb{E}[\hat{f}(X_i) | \hat{f}] \} - \frac{m_1}{n_1(n_1+m_1)} \{ \mathbb{E}[\hat{f}(X_i) | V_i, \hat{f}] - \mathbb{E}[\hat{f}(X_i) | \hat{f}] \} & \text{if } 1 \leq i \leq n_1, \\ \frac{1}{n_1+m_1} \{ \mathbb{E}[\hat{f}(X_i) | V_i, \hat{f}] - \mathbb{E}[\hat{f}(X_i) | \hat{f}] \} & \text{if } n_1 + 1 \leq i \leq n_1 + m_1, \end{cases}$$

$$H_i = \begin{cases} \frac{1}{n_2} \{ \hat{f}(Y_i) - \mathbb{E}[\hat{f}(Y_i) | \hat{f}] \} - \frac{m_2}{n_2(n_2+m_2)} \{ \mathbb{E}[\hat{f}(Y_i) | W_i, \hat{f}] - \mathbb{E}[\hat{f}(Y_i) | \hat{f}] \} & \text{if } 1 \leq i \leq n_2, \\ \frac{1}{n_2+m_2} \{ \mathbb{E}[\hat{f}(Y_i) | W_i, \hat{f}] - \mathbb{E}[\hat{f}(Y_i) | \hat{f}] \} & \text{if } n_2 + 1 \leq i \leq n_2 + m_2. \end{cases}$$

Define denominator as  $s_{n+m}^2 := \text{Var}(\hat{\mu}_{X,\hat{f}} - \hat{\mu}_{Y,\hat{f}} | \hat{f})$

$$\begin{aligned} &= \sum_{i=1}^{n_1+m_1} \text{Var}(G_i | \hat{f}) + \sum_{i=1}^{n_2+m_2} \text{Var}(H_i | \hat{f}) \\ &= \frac{1}{n_1} \sigma_{1,X}^2 + \frac{1}{n_1 + m_1} \sigma_{2,X}^2 + \frac{1}{n_2} \sigma_{1,Y}^2 + \frac{1}{n_2 + m_2} \sigma_{2,Y}^2 \end{aligned}$$

# Semi-Supervised Kernel Two-Sample Test

- Asymptotic normality

To use Lyapunov's CLT, we prove  $\sum_{i=1}^{n_1+m_1} \mathbb{E}\left[\left|\frac{G_i}{s_{n+m}}\right|^{2+\delta} \middle| \hat{f}\right] + \sum_{i=1}^{n_2+m_2} \mathbb{E}\left[\left|\frac{H_i}{s_{n+m}}\right|^{2+\delta} \middle| \hat{f}\right] \xrightarrow{p} 0$

with additional assumptions like

$$\frac{\mathbb{E}[\bar{k}^4(X_1, X_2)] \frac{n_1+n_2}{n_1 n_2} + \mathbb{E}[\bar{k}^2(X_1, X_3) \bar{k}^2(X_2, X_3)]}{\{\mathbb{E}[\bar{k}^2(X_1, X_2)]\}^2 \left(\frac{n_1 n_2}{n_1+n_2}\right)} \rightarrow 0$$

$$\inf_{\{\alpha_i\}_{i=1}^{\infty}} \left[ (1-c) \sum_{i=1}^{\infty} \lambda_i^2 \alpha_i^2 - \sum_{i,i'=1}^{\infty} \lambda_i \lambda_{i'} \alpha_i \alpha_{i'} \rho_{i,i'} \right] \geq 0$$

for centered kernel  $\bar{k}(x, y) = k(x, y) - \mathbb{E}[k(X, Y) | X = x] - \mathbb{E}[k(X, Y) | Y = y] + \mathbb{E}[k(X, Y)]$

$$= \sum_{i=1}^{\infty} \lambda_i \phi_i(x) \phi_i(y)$$

and  $\rho_{i,i'} = \text{Cov}\{\mathbb{E}[\phi_i(X) | V], \mathbb{E}[\phi_{i'}(X) | V]\}$

# Semi-Supervised Kernel Two-Sample Test

- Asymptotic normality

$$\text{If } \inf_{\{\alpha_i\}_{i=1}^{\infty}} \left[ (1 - c) \sum_{i=1}^{\infty} \lambda_i^2 \alpha_i^2 - \sum_{i,i'=1}^{\infty} \lambda_i \lambda_{i'} \alpha_i \alpha_{i'} \rho_{i,i'} \right] \geq 0$$

$$\begin{aligned} \text{Then, } \sigma_{1,X}^2 &= \mathbb{E} \left[ \left\{ \sum_{i=1}^{\infty} \lambda_i \left( \underbrace{\frac{1}{n_1} \sum_{j=1}^{n_1} \phi_i(X_j)}_{=\bar{\phi}_{i,X}} - \underbrace{\frac{1}{n_2} \sum_{j=1}^{n_2} \phi_i(Y_j)}_{=\bar{\phi}_{i,Y}} \right) (\phi_i(X) - \mathbb{E}[\phi_i(X) | V]) \right\}^2 \mid \mathcal{D}_{XV,2}, \mathcal{D}_{YW,2} \right] \\ &\geq c \sum_{i=1}^{\infty} \lambda_i^2 (\bar{\phi}_{i,X} - \bar{\phi}_{i,Y})^2, \end{aligned}$$

which is used for proof.

# Semi-Supervised Kernel Two-Sample Test

- Asymptotic normality

## 2) Test statistic with unknown variance

Using Slutsky's theorem, it suffices to show that  $\frac{\hat{\sigma}_{X,\hat{f}}^2 - \sigma_{X,\hat{f}}^2}{\sigma_{X,\hat{f}}^2} \xrightarrow{p} 0$  and  $\frac{\hat{\sigma}_{Y,\hat{f}}^2 - \sigma_{Y,\hat{f}}^2}{\sigma_{Y,\hat{f}}^2} \xrightarrow{p} 0$

which implies that  $\frac{\hat{\sigma}_{X,\hat{f}}^2 + \hat{\sigma}_{Y,\hat{f}}^2}{\sigma_{X,\hat{f}}^2 + \sigma_{Y,\hat{f}}^2} \xrightarrow{p} 1 \iff \frac{\hat{\sigma}_{X,\hat{f}}^2 - \sigma_{X,\hat{f}}^2 + \hat{\sigma}_{Y,\hat{f}}^2 - \sigma_{Y,\hat{f}}^2}{\sigma_{X,\hat{f}}^2 + \sigma_{Y,\hat{f}}^2} \xrightarrow{p} 0$

We compute 
$$\frac{|\hat{\sigma}_{X,\hat{f}}^2 - \sigma_{X,\hat{f}}^2|}{\sigma_{X,\hat{f}}^2} \leq \underbrace{\frac{|\hat{\sigma}_{1,X}^2 - \sigma_{1,X}^2|}{\sigma_{1,X}^2}}_{(I)} + \underbrace{\frac{n_1}{n_1 + m_1} \frac{|\hat{\sigma}_{2,X}^2 - \sigma_{2,X}^2|}{\sigma_{1,X}^2}}_{(II)}$$

and both terms are  $o_P(1)$  which completes the proof.

# Semi-Supervised Kernel Two-Sample Test

- Asymptotic normality

## 3) Test statistic using cross-fitting

We estimate  $\mathbb{E}[\hat{f}(X_i) | V_i, \hat{f}]$  as  $\widehat{\mathbb{E}}[\hat{f}(X_i) | V_i, \hat{f}]$  trained based on  $\mathcal{D}_{XV,1,b}$  if  $i \in \{1, 2, \dots, n_1/2\} \cup \{n_1 + 1, \dots, n_1 + m_1/2\}$  and on  $\mathcal{D}_{XV,1,a}$  if  $i \in \{n_1/2 + 1, n_1/2 + 2, \dots, n_1\} \cup \{n_1 + m_1/2 + 1, \dots, n_1 + m_1\}$ .

Then we prove  $T - T^* = o_P(1)$  assuming  $\sup_{i \geq 1} \mathbb{E}[\Delta_i^2] = o(1)$  for  $\Delta_i = \mathbb{E}[\phi_i(X) | V] - \widehat{\mathbb{E}}[\phi_i(X) | V]$

$$\begin{aligned} T - T^* &= \frac{\hat{N}}{\hat{D}} - \frac{\hat{N}}{D} + \frac{\hat{N}}{D} - \frac{N}{D} \\ &= \frac{\hat{N} - N}{D} \left( \frac{D}{\hat{D}} - 1 \right) + \underbrace{\frac{N}{D} \left( \frac{D}{\hat{D}} - 1 \right)}_{=O_P(1)} + \frac{\hat{N} - N}{D}. \end{aligned}$$

# Semi-Supervised Kernel Two-Sample Test

- Power analysis

Define our test function  $\phi^* := \mathbb{1}_{\{T^* > z_{1-\alpha}\}}$

Then we express its power as

$$\begin{aligned} \mathbb{P}(T^* > z_{1-\alpha}) &= \mathbb{P}\left(\frac{\widehat{\mu}_{X,\widehat{f}} - \widehat{\mu}_{Y,\widehat{f}}}{\sqrt{\widehat{\sigma}_{X,\widehat{f}}^2 + \widehat{\sigma}_{Y,\widehat{f}}^2}} > z_{1-\alpha}\right) \\ &= \mathbb{P}\left(\frac{(\widehat{\mu}_{X,\widehat{f}} - \frac{1}{n_1+m_1} \sum_{i=1}^{n_1+m_1} \mathbb{E}[\widehat{f}(X_i)|\widehat{f}]) - (\widehat{\mu}_{Y,\widehat{f}} - \frac{1}{n_2+m_2} \sum_{i=1}^{n_2+m_2} \mathbb{E}[\widehat{f}(Y_i)|\widehat{f}])}{\sqrt{\widehat{\sigma}_{X,\widehat{f}}^2 + \widehat{\sigma}_{Y,\widehat{f}}^2}}\right. \\ &\quad \left. > z_{1-\alpha} - \frac{\frac{1}{n_1+m_1} \sum_{i=1}^{n_1+m_1} \mathbb{E}[\widehat{f}(X_i)|\widehat{f}] - \frac{1}{n_2+m_2} \sum_{i=1}^{n_2+m_2} \mathbb{E}[\widehat{f}(Y_i)|\widehat{f}]}{\sqrt{\widehat{\sigma}_{X,\widehat{f}}^2 + \widehat{\sigma}_{Y,\widehat{f}}^2}}\right) \end{aligned}$$

# Semi-Supervised Kernel Two-Sample Test

- Asymptotic power expression with linear kernel

Assuming normality and denoting  $\Sigma = n_1^{-1}\Sigma_{11} + n_2^{-1}\tilde{\Sigma}_{11}$ ,  $\Lambda = \frac{1-\rho_{XV}^2}{n_1}\Sigma_{22} + \frac{\rho_{XV}^2}{n_1}\Sigma_{11}$ , and

$$\begin{aligned} \begin{pmatrix} X \\ V \end{pmatrix} &\sim N \left( \begin{pmatrix} \mu_X \\ \mu_V \end{pmatrix}, \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} \right) & \tilde{\Lambda} &= \frac{1-\rho_{YW}^2}{n_2}\tilde{\Sigma}_{22} + \frac{\rho_{YW}^2}{n_2}\tilde{\Sigma}_{11} \\ \begin{pmatrix} Y \\ W \end{pmatrix} &\sim N \left( \begin{pmatrix} \mu_Y \\ \mu_W \end{pmatrix}, \begin{pmatrix} \tilde{\Sigma}_{11} & \tilde{\Sigma}_{12} \\ \tilde{\Sigma}_{21} & \tilde{\Sigma}_{22} \end{pmatrix} \right) \end{aligned}$$

we guess the power function of

$$\begin{aligned} \mathbb{P}(T^* > z_{1-\alpha}) &= \Phi \left( z_\alpha + \frac{(\mu_X - \mu_Y)^\top (\mu_X - \mu_Y)}{\sqrt{\text{tr}((\Lambda + \tilde{\Lambda})\Sigma)}} \right) + o_P(1) \\ &= \Phi \left( z_\alpha + \frac{(\mu_X - \mu_Y)^\top (\mu_X - \mu_Y)}{n_1^{-1} \sqrt{(1 - \rho_{XV}^2)\text{tr}(\Sigma_{22}\Sigma_{11}) + \rho_{XV}^2\text{tr}(\Sigma_{11}^2)}} \right) + o_P(1) \end{aligned}$$

# Semi-Supervised Kernel Two-Sample Test

- Power consistency (if we have asymptotic power of 1)

Prove  $\mathbb{E}_{P,Q}[1 - \phi^*] = \mathbb{P}_{P,Q}(T^* \leq z_{1-\alpha}) = \mathbb{P}_{P,Q}(\hat{\mu}_{X,\hat{f}} - \hat{\mu}_{Y,\hat{f}} \leq z_{1-\alpha}\sqrt{\hat{\sigma}_{X,\hat{f}}^2 + \hat{\sigma}_{Y,\hat{f}}^2})$  converges to 0 by showing

$$\frac{\mathbb{E}_{P,Q}[\hat{\sigma}_{X,\hat{f}}^2 + \hat{\sigma}_{Y,\hat{f}}^2]}{\gamma^4 \delta_n} + \frac{\text{Var}_{P,Q}(\hat{\mu}_{X,\hat{f}} - \hat{\mu}_{Y,\hat{f}})}{\gamma^4} \leq \frac{\mathbb{E}_{P,Q}[\hat{\sigma}^2]}{\gamma^4 \delta_n} + \frac{\text{Var}(\widehat{\text{xMMD}})}{\gamma^4} \rightarrow 0$$

Also, using cross-fitted estimator, under the same assumption of  $\sup_{i \geq 1} \mathbb{E}[\Delta_i^2] = o(1)$

for  $\Delta_i = \mathbb{E}[\phi_i(X) | V] - \widehat{\mathbb{E}}[\phi_i(X) | V]$

$$\mathbb{E}_{P,Q}[\widehat{D}^2 - D^2] < \infty \quad \text{and} \quad \text{Var}_{P,Q}(\widehat{N} - N) \rightarrow 0$$

# Semi-Supervised Kernel Two-Sample Test

- Asymptotic normality under the alternative

We need additional assumptions like

$$\max\{\|p/q\|_\infty, \|q/p\|_\infty\} < C$$

$$\frac{\mathbb{E}[\bar{k}_X(X_1, X_2)^4] + n_1 \mathbb{E}[\bar{k}_X(X_1, X_2)^2 \bar{k}_X(X_1, X_3)^2]}{n_1^2 \{\mathbb{E}[\bar{g}_X(X, X)]\}^2} \rightarrow 0$$

$$\frac{\text{MMD}^4 \mathbb{E}[\bar{k}_X(X, X)^2]}{\{n_1 \mathbb{E}[\bar{g}_X(X, X)] + n_1^2 \mathbb{E}[\bar{g}_X(Y_1, Y_2)]\}^2} \rightarrow 0,$$

# Semi-Supervised Kernel Two-Sample Test

- Asymptotic normality under the alternative

Again, we use Lyapunov's CLT showing

$$\frac{1}{\sigma_{n_1, n_2}^4} \left[ \frac{1}{n_1^4} \sum_{i=1}^{n_1} \mathbb{E}[\{\hat{f}(X_i) - \mathbb{E}[\hat{f}(X) | \hat{f}]\}^4 | \hat{f}] + \frac{1}{n_2^4} \sum_{i=1}^{n_2} \mathbb{E}[\{\hat{f}(Y_i) - \mathbb{E}[\hat{f}(Y) | \hat{f}]\}^4 | \hat{f}] \right] \xrightarrow{p} 0$$

Since assuming the alternative,

$$\begin{aligned} \text{Var}[\hat{f}(X) | \hat{f}] &= \mathbb{E}[\langle \bar{\psi}_X - \bar{\psi}_Y, \psi(X) - \mathbb{E}_P[\psi(X)] \rangle^2 | \bar{\psi}_X, \bar{\psi}_Y] \\ &= \mathbb{E}[\langle \bar{\psi}_X - \mathbb{E}[\psi(X)] + \mathbb{E}[\psi(X)] - \bar{\psi}_Y, \psi(X) - \mathbb{E}_P[\psi(X)] \rangle^2 | \bar{\psi}_X, \bar{\psi}_Y] \\ &= \mathbb{E}\left[\left\{ \sum_{i=1}^{\infty} \lambda_i \left( \frac{1}{n_1} \sum_{j=1}^{n_1} \phi_i(X'_j) \right) \phi_i(X) - \sum_{i=1}^{\infty} \lambda_i \left( \frac{1}{n_2} \sum_{j=1}^{n_2} \phi_i(Y'_j) \right) \phi_i(X) \right\}^2 \middle| (X'_j), (Y'_j) \right] \\ &= \mathbb{E}\left[\left\{ \sum_{i=1}^{\infty} \lambda_i (\bar{\phi}_{i,X} - \bar{\phi}_{i,Y}) \phi_i(X) \right\}^2 \middle| (X'_j), (Y'_j) \right] \\ &= \sum_{i=1}^{\infty} \lambda_i^2 (\bar{\phi}_{i,X} - \bar{\phi}_{i,Y})^2. \text{ and, similarly, } \text{Var}[\hat{f}(Y) | \hat{f}] = \sum_{i=1}^{\infty} \mu_i^2 (\bar{\varphi}_{i,X} - \bar{\varphi}_{i,Y})^2. \end{aligned}$$

# Semi-Supervised Kernel Two-Sample Test

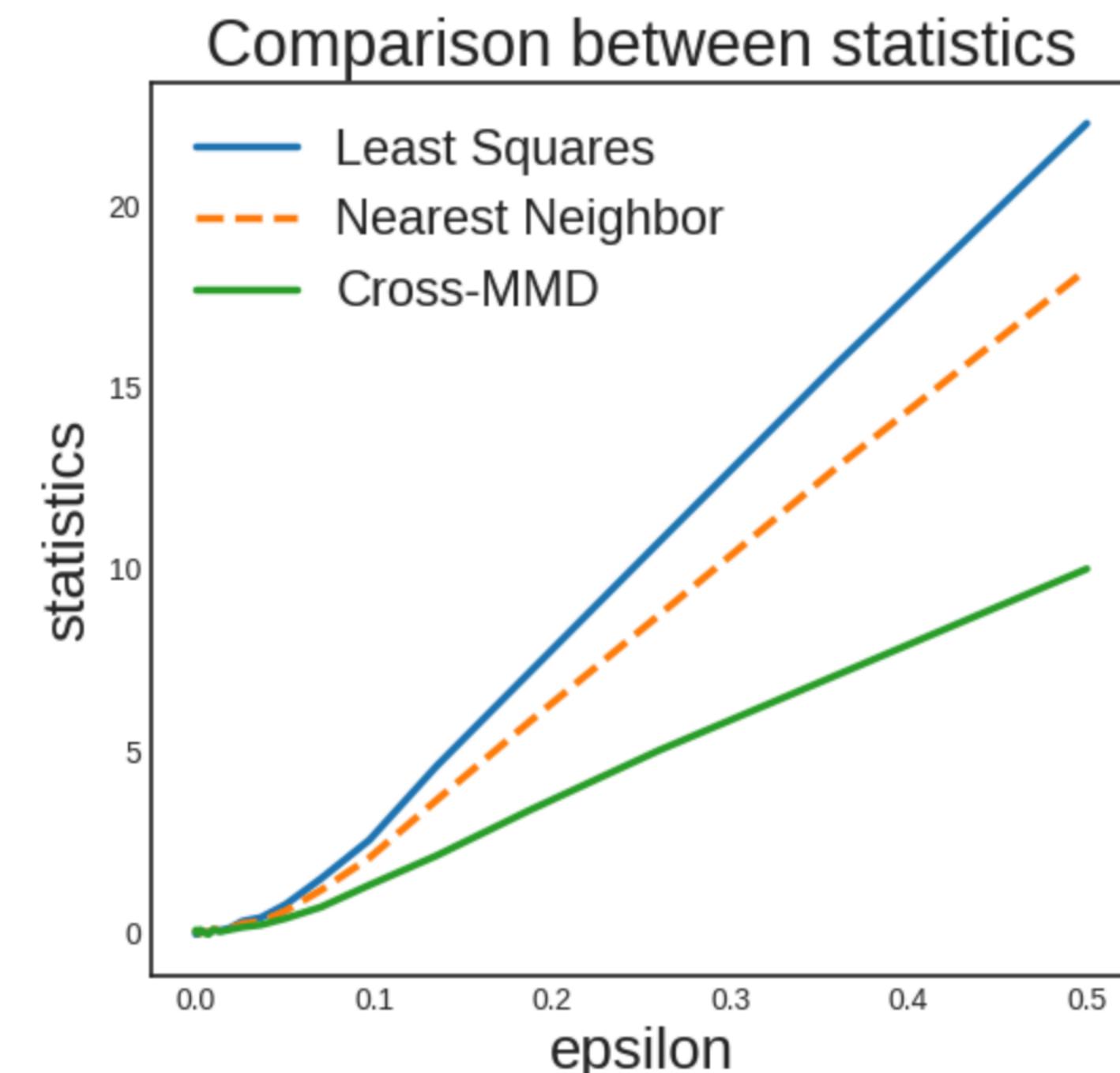
- Numerical Studies ( $d=4$ , RBF kernel)

$$X = V_{(1)} + V_{(2)} + 0.3Z$$

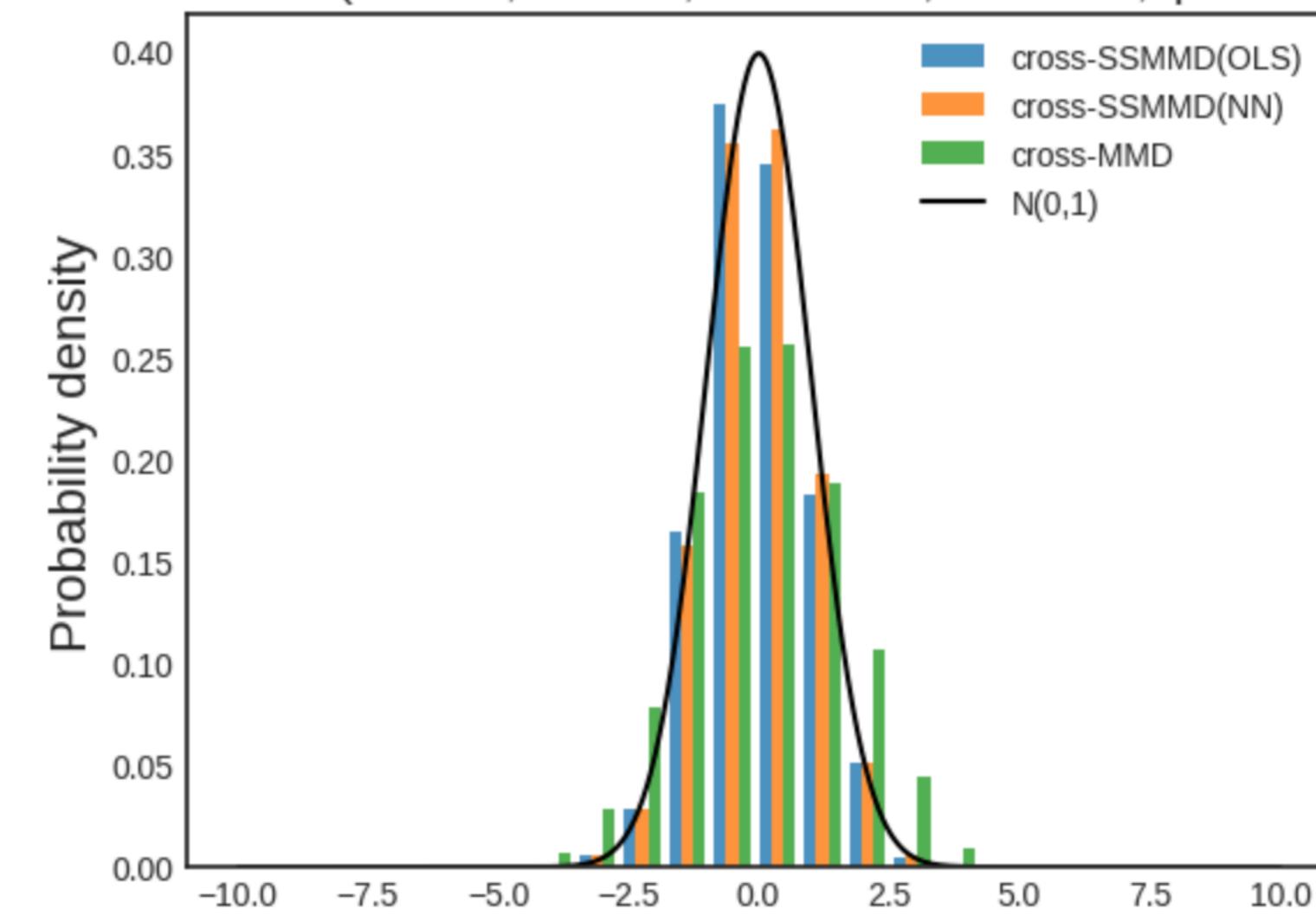
$$Y = W_{(1)} + W_{(2)} + 0.3Z$$

$$V \sim N_d(\mu_1, \Sigma_1), \quad \text{where} \quad \mu_1 = \mathbf{0}_d, \quad \Sigma_1 = 0.3I_d + 0.7\mathbb{1}'_d\mathbb{1}_d$$

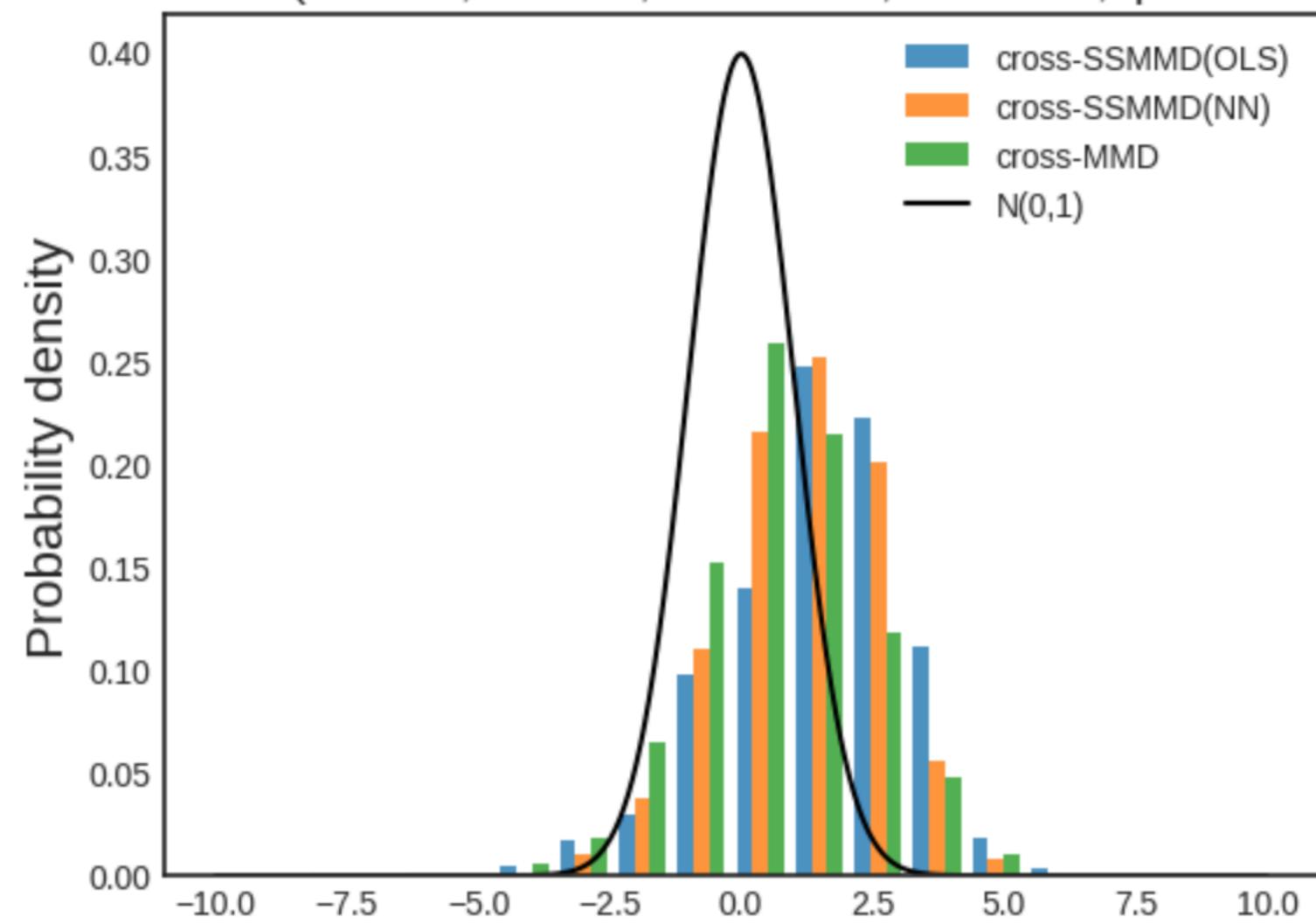
$$W \sim N_d(\mu_2, \Sigma_2), \quad \text{where} \quad \mu_2 = (\epsilon, \dots, \epsilon)', \quad \Sigma_2 = 0.3I_d + 0.7\mathbb{1}'_d\mathbb{1}_d$$



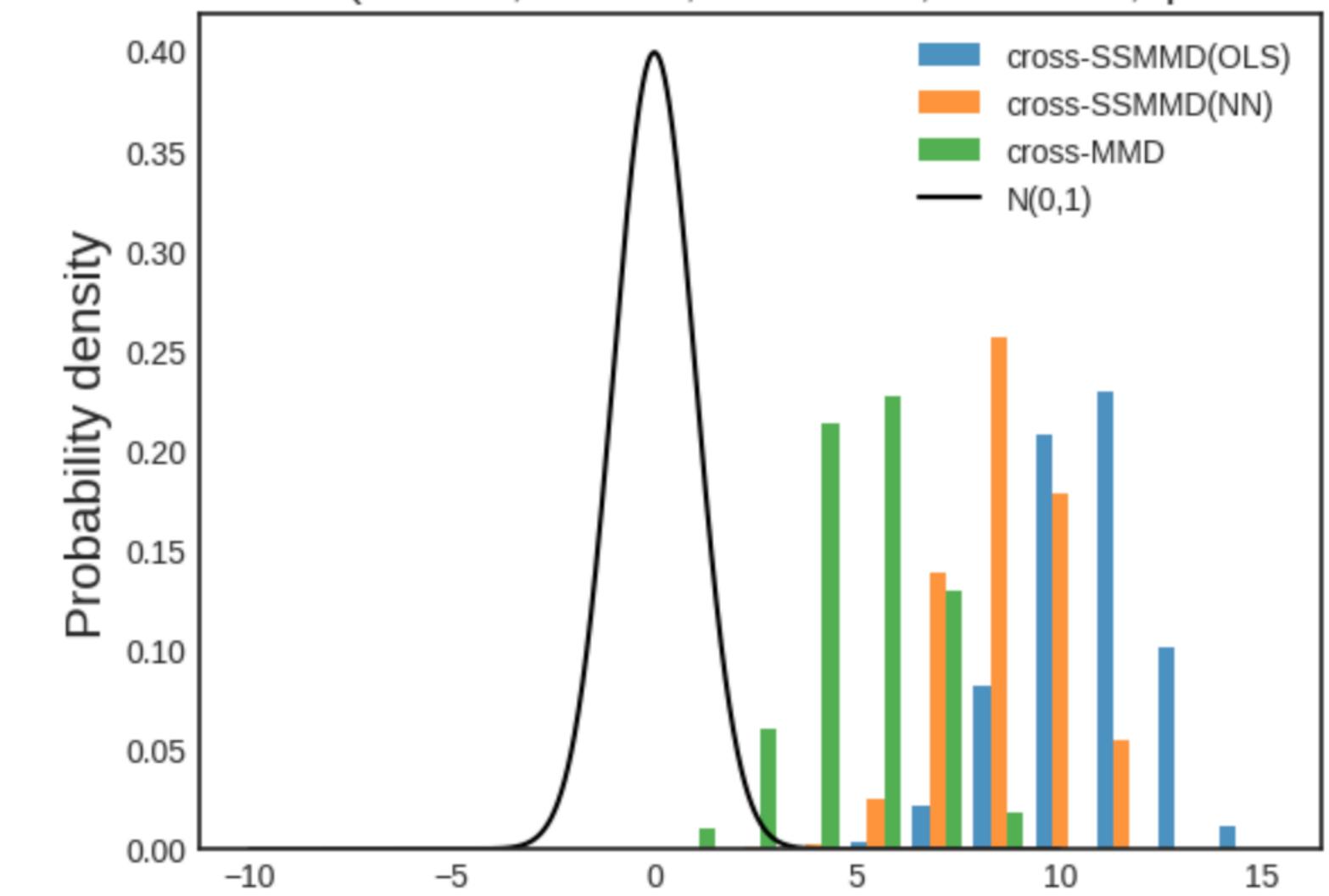
SS-xMMD ( $n_1=500, n_2=400, m_1=10000, m_2=8000, \text{epsilon}=0.0$ )



SS-xMMD ( $n_1=500, n_2=400, m_1=10000, m_2=8000, \text{epsilon}=0.832$ )



SS-xMMD ( $n_1=500, n_2=400, m_1=10000, m_2=8000, \text{epsilon}=5.0$ )



# Semi-Supervised Kernel Two-Sample Test

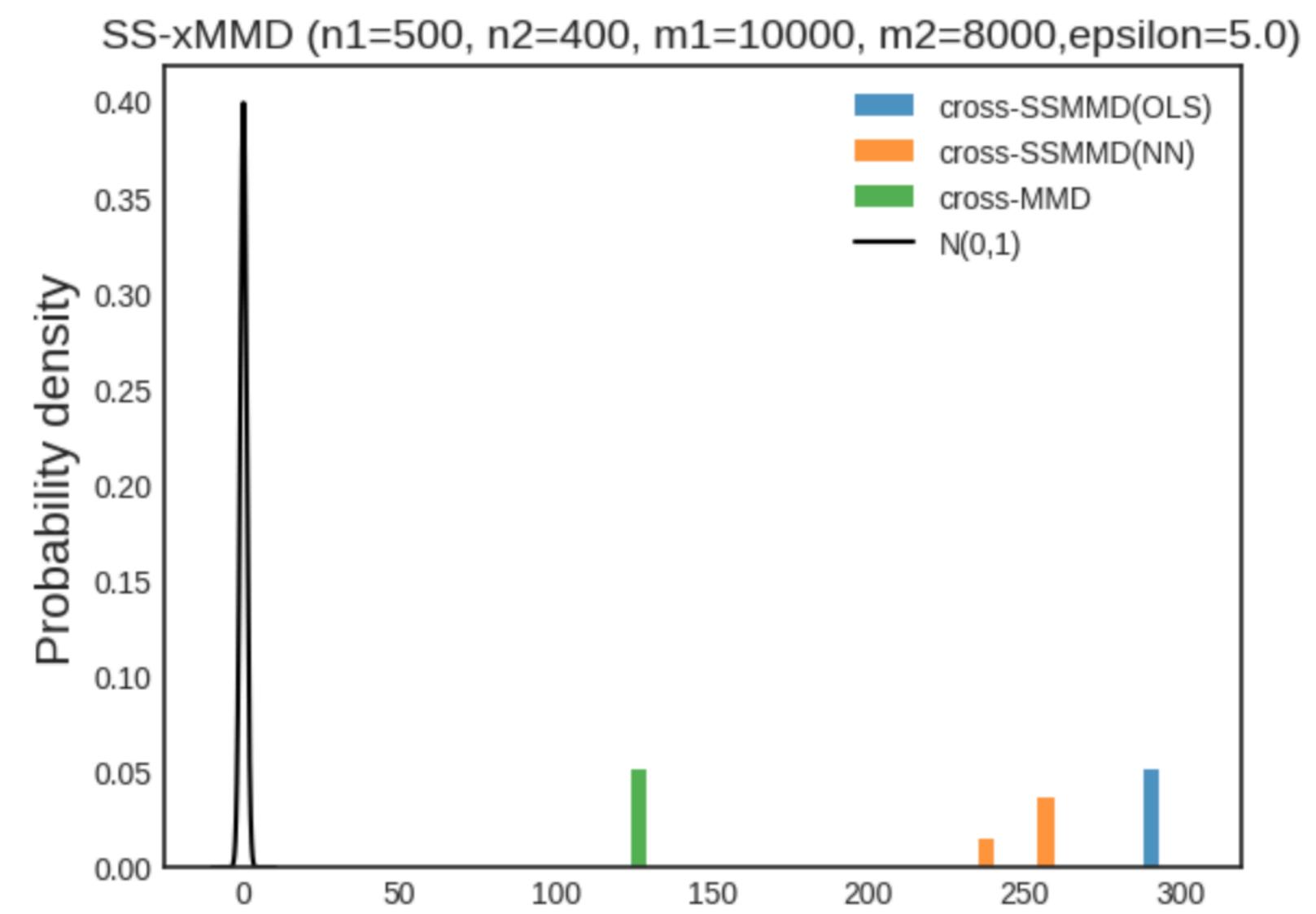
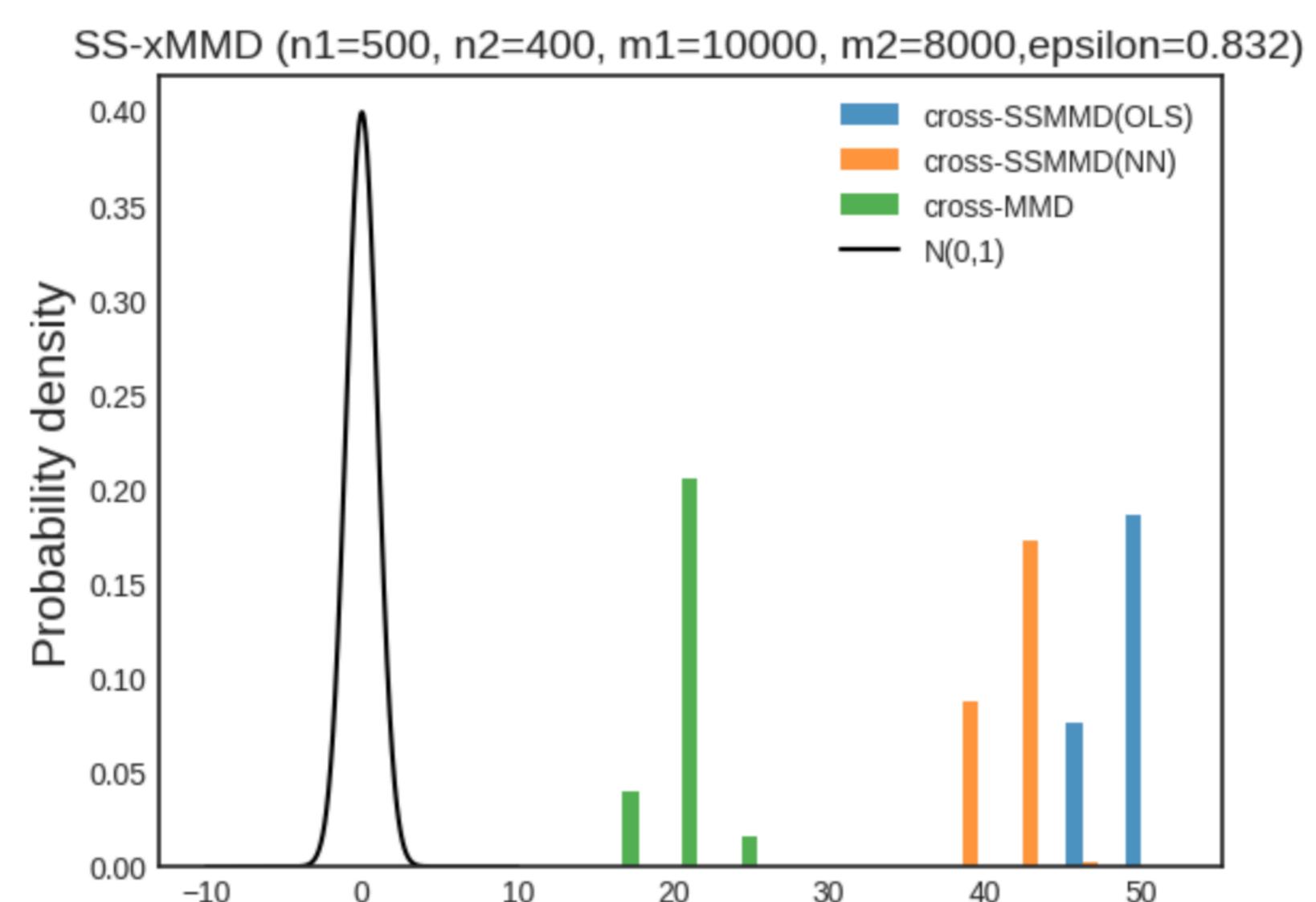
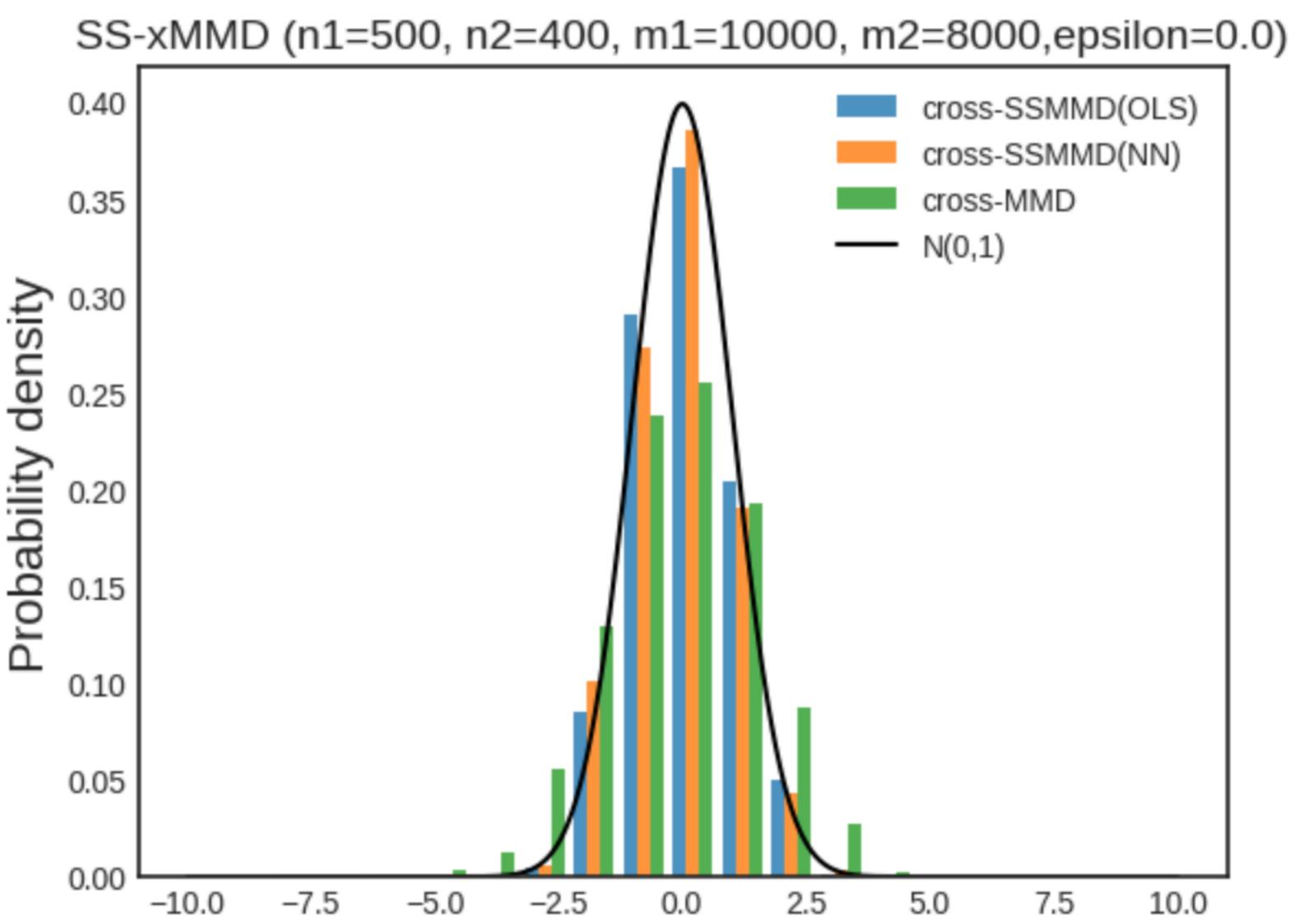
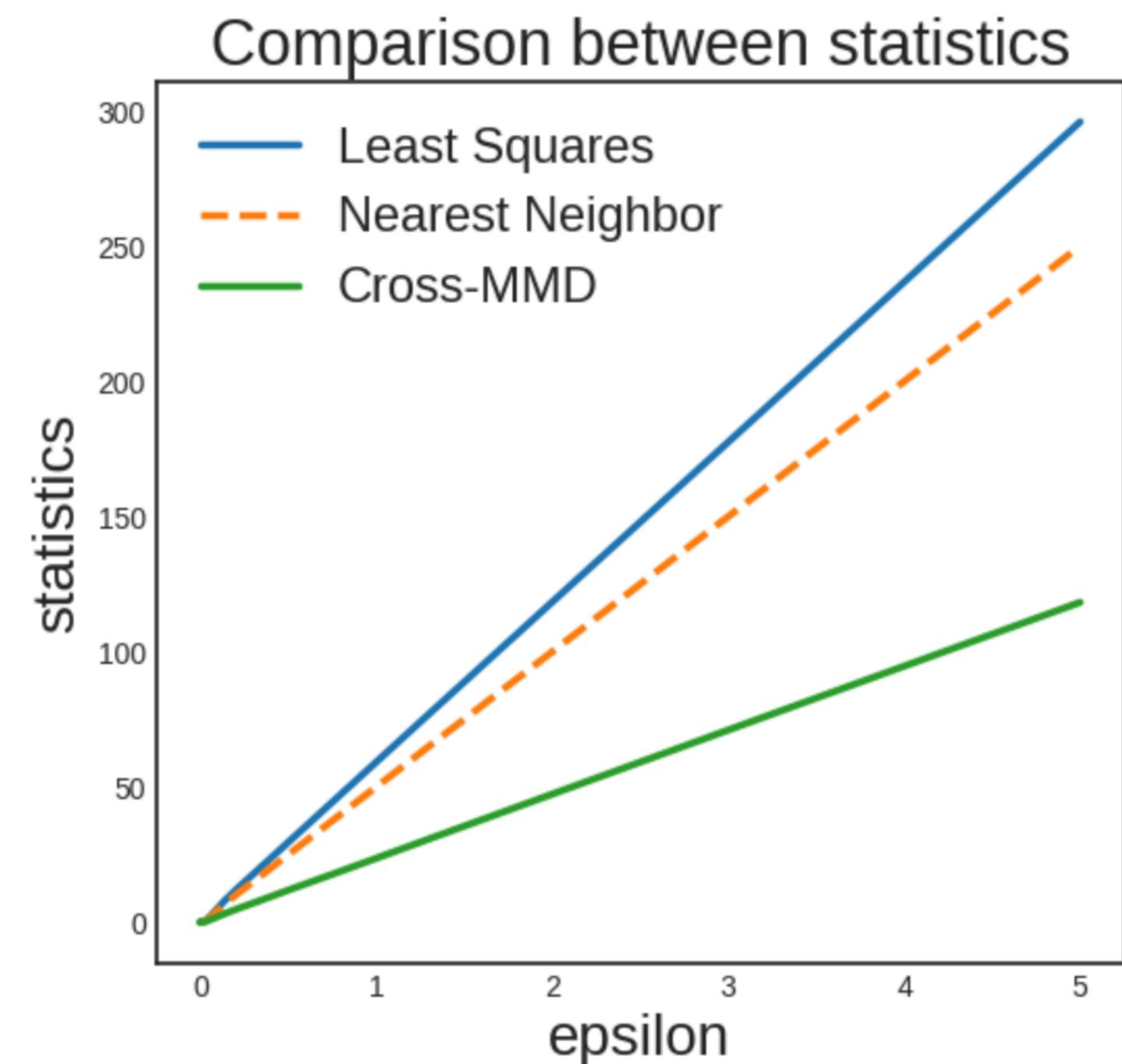
- Numerical Studies ( $d=4$ , linear kernel)

$$X = V_{(1)} + V_{(2)} + 0.3Z$$

$$Y = W_{(1)} + W_{(2)} + 0.3Z$$

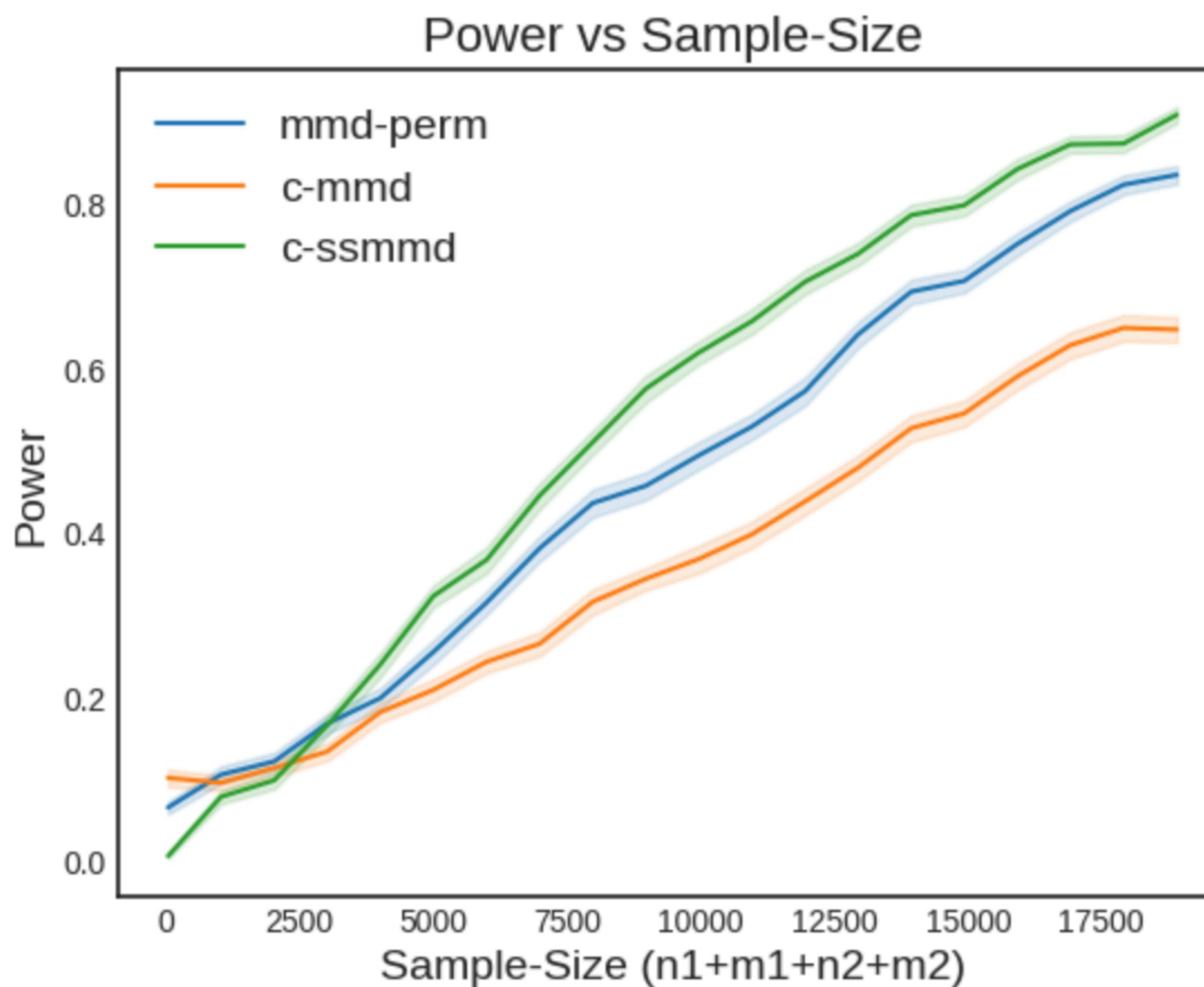
$$V \sim N_d(\mu_1, \Sigma_1), \quad \text{where } \mu_1 = \mathbf{0}_d, \Sigma_1 = 0.3I_d + 0.7\mathbb{1}'_d\mathbb{1}_d$$

$$W \sim N_d(\mu_2, \Sigma_2), \quad \text{where } \mu_2 = (\epsilon, \dots, \epsilon)', \Sigma_2 = 0.3I_d + 0.7\mathbb{1}'_d\mathbb{1}_d$$

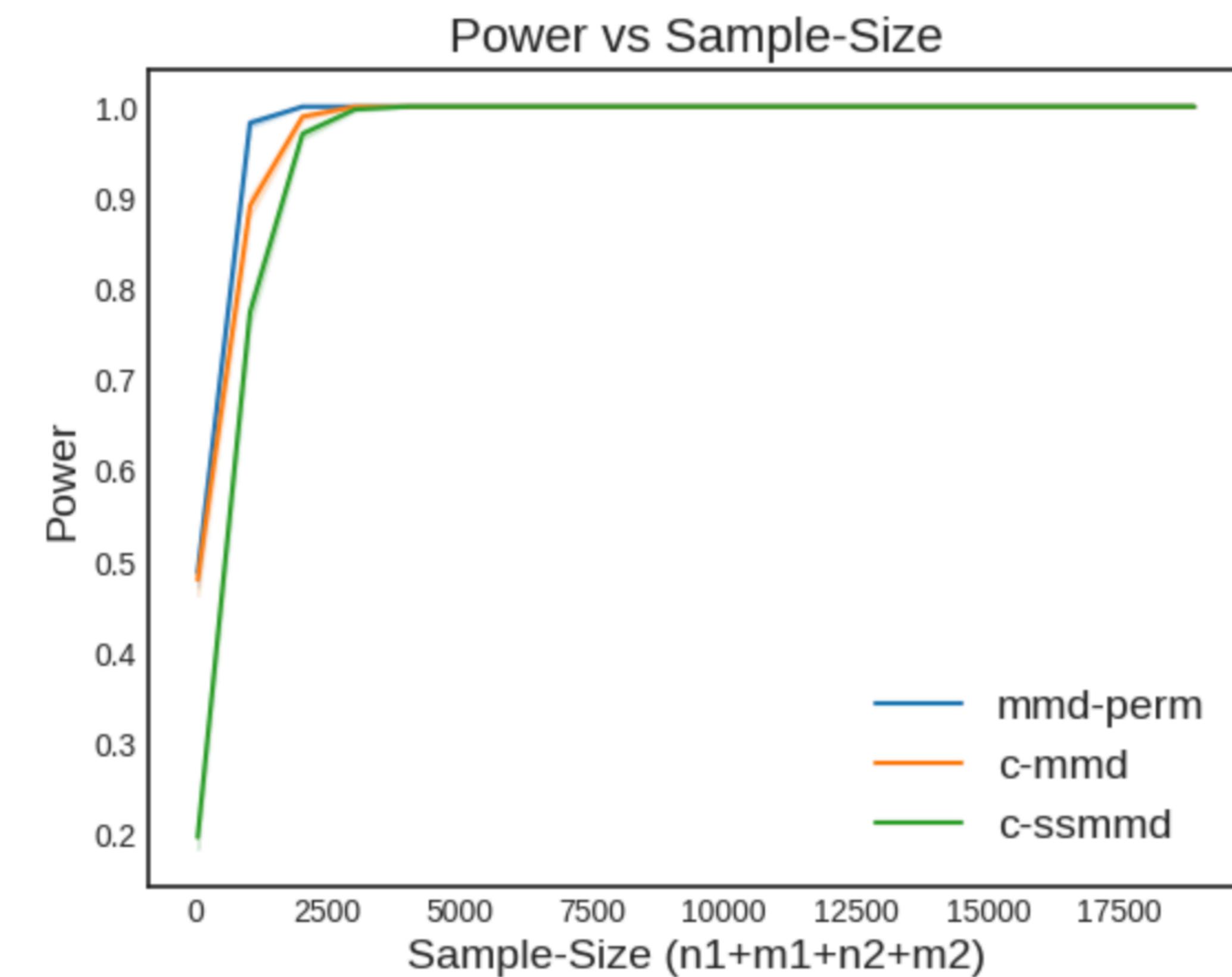


# Semi-Supervised Kernel Two-Sample Test

- Numerical Studies  $X = V_{(1)} + V_{(2)} + V_{(d)} + 0.3Z$   $V \sim N_d(\mu_1, \Sigma_1)$ , where  $\mu_1 = \mathbf{0}_d$ ,  $\Sigma_1 = I_d$   
 $Y = W_{(1)} + W_{(2)} + W_{(d)} + 0.3Z$   $W \sim N_d(\mu_2, \Sigma_2)$ , where  $\mu_2 = (\underbrace{\epsilon, \dots, \epsilon}_{p \text{ elements}}, 0, \dots, 0)'$ ,  $\Sigma_2 = I_d$



$d=10, p=3, \text{eps}=0.3, n_1=400, m_1=8000, n_2=500, m_2=10000$



$d=50, p=50, \text{eps}=1, n_1=400, m_1=8000, n_2=500, m_2=10000$

# Semi-Supervised Kernel Two-Sample Test

- Numerical Studies

$$X = V_{(1)} + V_{(2)} + V_{(d)} + 0.3Z$$

$$Y = W_{(1)} + W_{(2)} + W_{(d)} + 0.3Z$$

$$V \sim N_d(\mu_1, \Sigma_1), \text{ where } \mu_1 = \mathbf{0}_d, \Sigma_1 = I_d$$

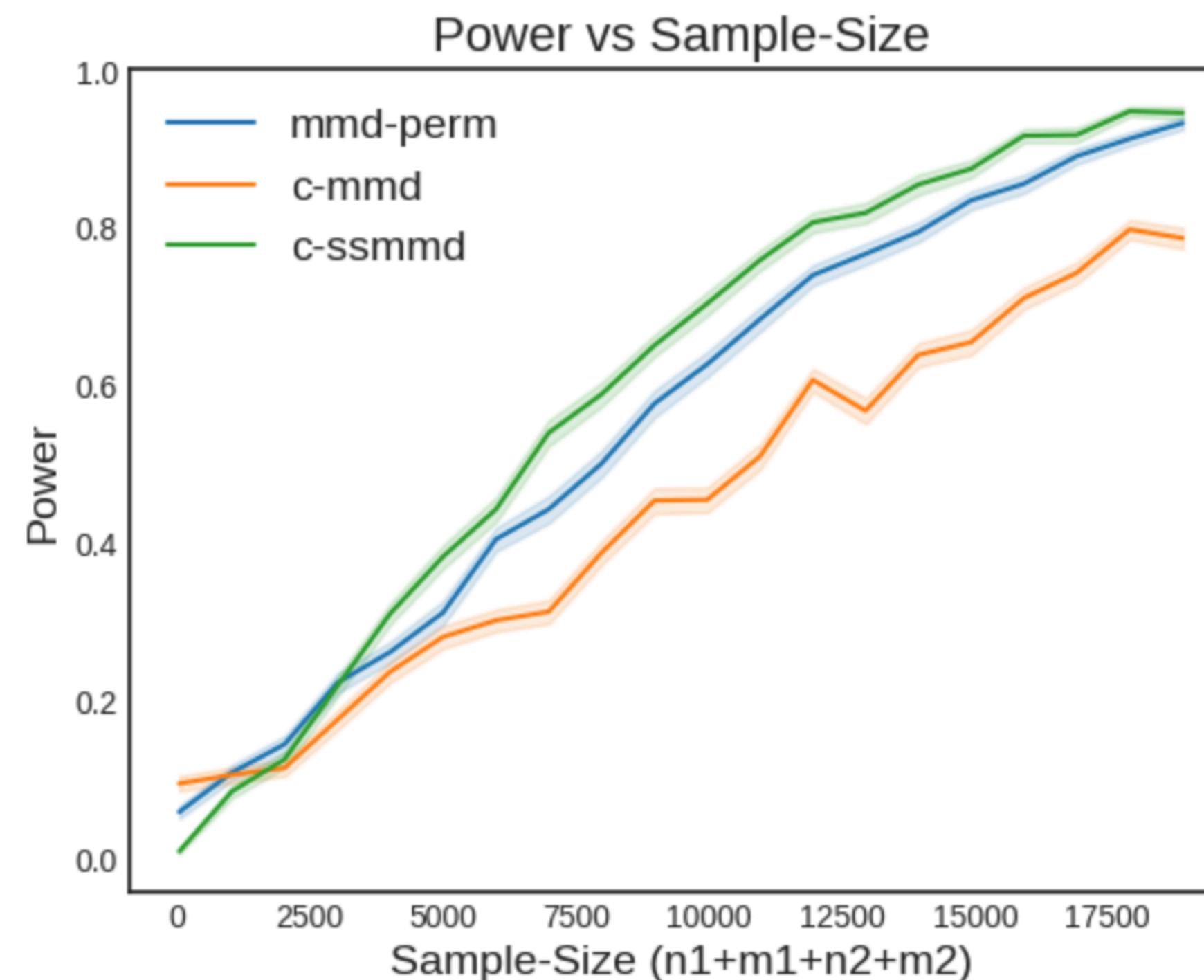
$$W \sim N_d(\mu_2, \Sigma_2), \text{ where } \mu_2 = \mathbf{0}_d, \Sigma_2 = \text{diag}(\underbrace{\epsilon, \dots, \epsilon}_{p \text{ elements}}, 1, \dots, 1)$$

$$X = V_{(1)} + V_{(2)} + V_{(d)} + 0.3Z$$

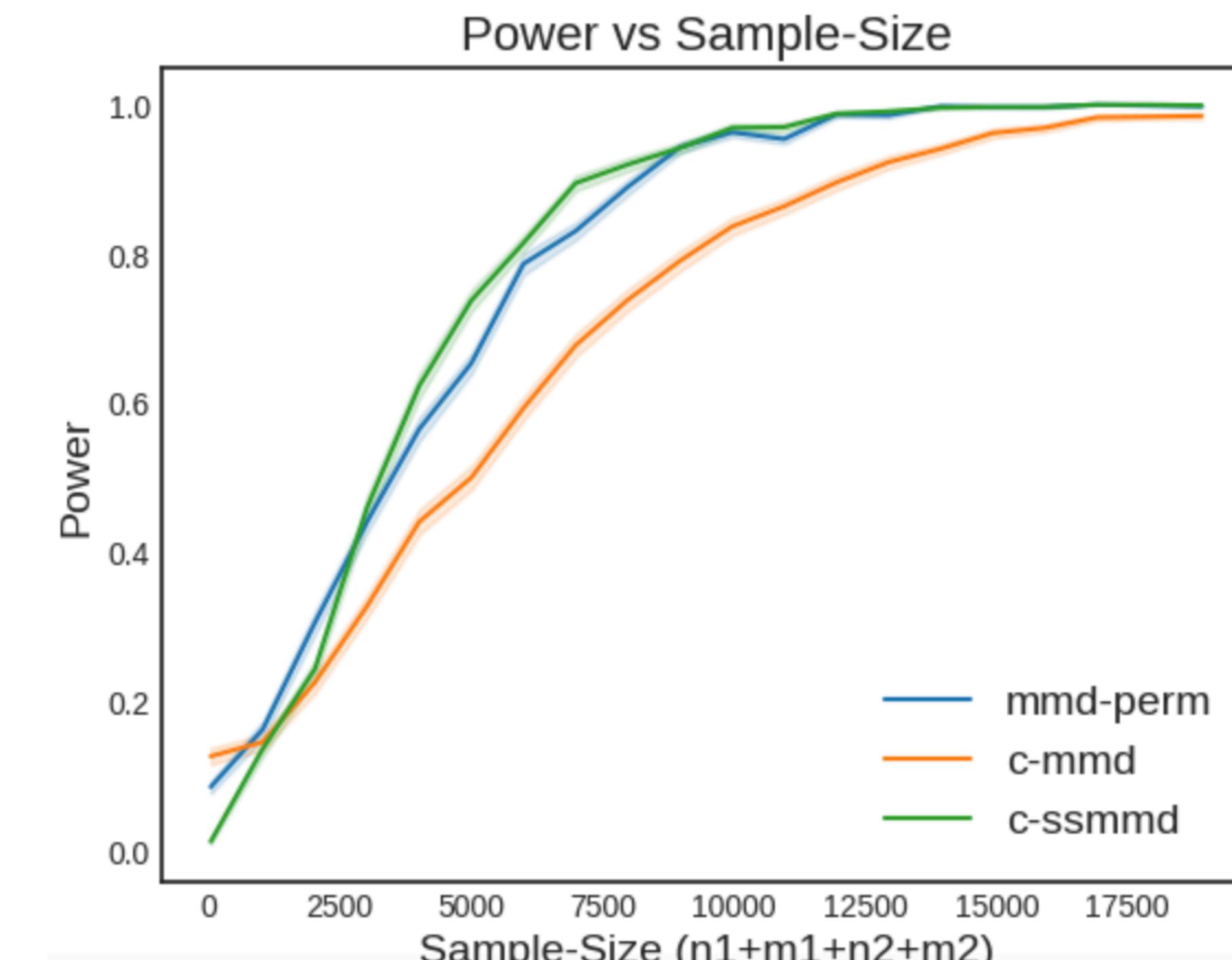
$$Y = W_{(1)} + W_{(2)} + W_{(d)} + 0.3Z$$

$$V \sim N_d(\mu_1, \Sigma_1), \text{ where } \mu_1 = \mathbf{0}_d, \Sigma_1 = I_d$$

$$W \sim N_d(\mu_2, \Sigma_2), \text{ where } \mu_2 = (\underbrace{\epsilon, \dots, \epsilon}_{p \text{ elements}}, 0, \dots, 0)', \Sigma_2 = \text{diag}(\underbrace{\epsilon, \dots, \epsilon}_{p \text{ elements}}, 1, \dots, 1)$$



d=10,p=3,eps=0.3,n1=400,m1=8000,n2=500,m2=10000



d=10,p=3,eps=0.3,n1=400,m1=8000,n2=500,m2=10000

# Conclusion

- Proposed semi-supervised kernel two-sample test statistic using sample splitting, and cross-fitting.
- Proved asymptotic normality of the test statistic under the null and the alternative and figured out under which condition it shows the properties.
- Derived asymptotic power expression using linear kernel and showed power consistency with proper assumptions.

# Future work

- What if not using cross-fitting? What if it is not linear?
- Are all the assumptions necessary?
- What if the distributions, so MMD are not fixed w.r.t. n?

# Reference

- Angelopoulos, A. N., Bates, S., Fannjiang, C., Jordan, M. I., & Zrnic, T. (2023). Prediction-powered inference. *Science*, 382(6671), 669-674.
- Gretton, A., Borgwardt, K. M., Rasch, M. J., Schölkopf, B., & Smola, A. (2012). A kernel two-sample test. *The Journal of Machine Learning Research*, 13(1), 723-773.
- Han, J., Liang, X., Xu, H., Chen, K., Hong, L., Mao, J., ... & Xu, C. (2021). SODA10M: A large-scale 2D self/semi-supervised object detection dataset for autonomous driving. *arXiv preprint arXiv:2106.11118*.
- Hastie, T., Tibshirani, R., Friedman, J. H., & Friedman, J. H. (2009). The elements of statistical learning: data mining, inference, and prediction (Vol. 2, pp. 1-758). New York: springer.
- Jensen, P. B., Jensen, L. J., & Brunak, S. (2012). Mining electronic health records: towards better research applications and clinical care. *Nature Reviews Genetics*, 13(6), 395-405.
- Jiao, R., Zhang, Y., Ding, L., Xue, B., Zhang, J., Cai, R., & Jin, C. (2023). Learning with limited annotations: a survey on deep semi-supervised learning for medical image segmentation. *Computers in Biology and Medicine*, 107840.
- Ilmun Kim, (2023), Lecture notes “Selective Topics in Mathematical Statistics” for graduate mathematical statistics course
- Liu, F., Xu, W., Lu, J., Zhang, G., Gretton, A., & Sutherland, D. J. (2020, November). Learning deep kernels for non-parametric two-sample tests. In *International conference on machine learning* (pp. 6316-6326). PMLR.
- Murphy, K. P. (2018). Machine learning: A probabilistic perspective (adaptive computation and machine learning series). *The MIT Press: London, UK*.
- Shekhar, S., Kim, I., & Ramdas, A. (2022). A permutation-free kernel two-sample test. *Advances in Neural Information Processing Systems*, 35, 18168-18180.
- Tony Cai, T., & Guo, Z. (2020). Semisupervised inference for explained variance in high dimensional linear regression and its applications. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 82(2), 391-419.
- Walpole, R. E., Myers, R. H., Myers, S. L., & Ye, K. (1993). *Probability and statistics for engineers and scientists* (Vol. 5, pp. 326-332). New York: Macmillan.
- Yongho, Jeon (2023), Lecture notes for nonparametric function estimation
- Zhang, A., Brown, L. D., & Cai, T. T. (2019). Semi-supervised inference: General theory and estimation of means.
- Zhang, Y., & Bradic, J. (2022). High-dimensional semi-supervised learning: in search of optimal inference of the mean. *Biometrika*, 109(2), 387-403.
- Zhu, B., Ding, M., Jacobson, P., Wu, M., Zhan, W., Jordan, M., & Jiao, J. (2024). Doubly-Robust Self-Training. *Advances in Neural Information Processing Systems*, 36.
- <https://towardsdatascience.com/supervised-vs-unsupervised-learning-in-2-minutes-72dad148f242>
- [http://www.gatsby.ucl.ac.uk/~gretton/papers/columbia23/columbia23\\_21](http://www.gatsby.ucl.ac.uk/~gretton/papers/columbia23/columbia23_21)