

Corrected generalized cross-validation for finite ensembles of penalized estimators

Takuya Koriyama

University of Chicago

September 4, 2024

- To appear in Journal of the Royal Statistical Society: Series B (2024)
- Joint work with Pierre C. Bellec (Rutgers), Jin-Hong Du (CMU), Kai Tan (Rutgers), and Pratik Patil (UC Berkeley).

Problem set up

- The response and feature $(y_i, \mathbf{x}_i) \in \mathbb{R} \times \mathbb{R}^p$ ($i = 1, \dots, n$) are i.i.d. distributed.
- Consider the high-dimensional regime

$p/n \rightarrow \text{constant}$ for sample size n and dimension p .

- We are interested in an estimator $\hat{\beta} = \hat{\beta}(\mathbf{y}, \mathbf{X})$ such that the prediction risk

$$\mathbb{E} \left[(y_0 - \mathbf{x}_0^\top \hat{\beta})^2 \mid \mathbf{y}, \mathbf{X} \right] \quad \text{where} \quad (y_0, \mathbf{x}_0) =^d (y_i, \mathbf{x}_i)$$

is small.

- We consider **ensemble estimators** $\tilde{\beta}$ (next slide).

Ensemble estimator $\tilde{\beta}$

We define ensemble estimator $\tilde{\beta}$ as follows:

- 1 Subsampling

$$(I_m)_{m=1}^M \stackrel{iid}{\sim} \text{Uniform}\{I \subset [n] : |I| = k\}$$

for some integers $k \leq n$ and M .

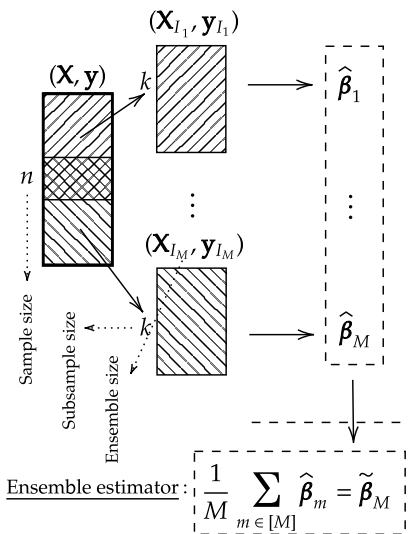
- 2 Fit the penalized least square

$$\hat{\beta}_m \in \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{k} \|\mathbf{y}_{I_m} - \mathbf{X}_{I_m} \beta\|^2 + g(\beta)$$

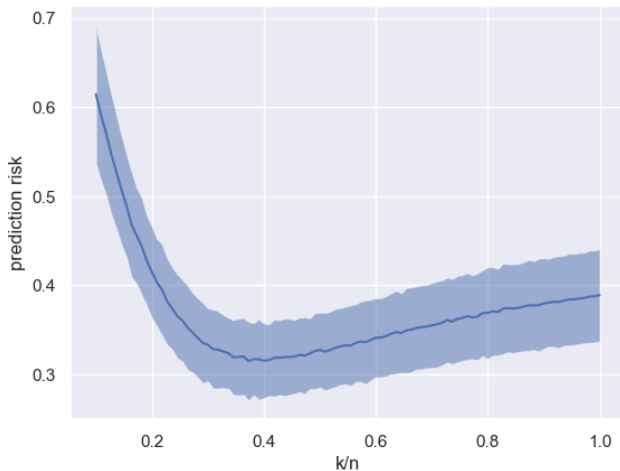
for some convex function $g : \mathbb{R}^p \rightarrow \mathbb{R}$.

- 3 Ensemble $(\hat{\beta}_m)_{m=1}^M$ together

$$\tilde{\beta} = \frac{1}{M} \sum_{m=1}^M \hat{\beta}_m.$$



Prediction risk is U-shape in sub-sample size k



Ensemble of Ridge estimators.

Equivalence between subsampling and regularization

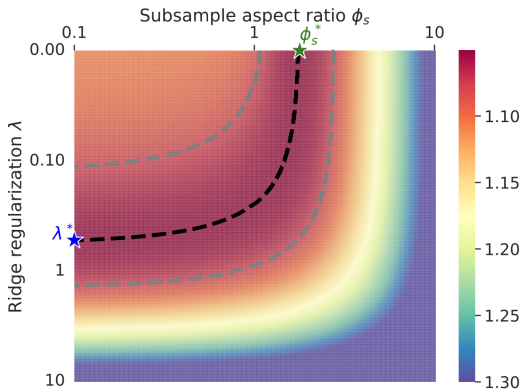


Figure 1 in Du et al. [2023].

Adaptive tuning of sub-sample size and penalty

(Recall) Ensemble estimator is $\tilde{\beta} = \frac{1}{M} \sum_{m=1}^M \hat{\beta}_m$ where

$$\hat{\beta}_m \in \arg \min_{\beta \in \mathbb{R}^p} \frac{1}{k} \|\mathbf{y}_{I_m} - \mathbf{X}_{I_m} \beta\|^2 + g(\beta), \quad I_m \sim \text{Uniform}\{I \subset [n] : |I| = k\}$$

for each $m \in [M]$.

- (Goal) Select sub-sample size k and penalty g in a data-driven manner so that the ensemble estimator $\tilde{\beta}$ achieves a small prediction risk

$$\mathbb{E}[(y_0 - \mathbf{x}_0^\top \tilde{\beta})^2 | \mathbf{y}, \mathbf{X}] \quad \text{where} \quad (y_0, \mathbf{x}_0) =^d (y_i, \mathbf{x}_i)$$

- Since the prediction risk is not observable, we need some proxy;
 - ▶ L -fold cross-validation is biased.
 - ▶ Leave one out cross-validation is computationally hard due to high-dimension.
 - ▶ **Generalized cross-validation (GCV).**

Generalized cross validation

For the penalized least square estimator

$$\hat{\beta}(\mathbf{y}, \mathbf{X}) \in \arg \min_{\beta \in \mathbb{R}^p} \left\{ \frac{1}{n} \|\mathbf{y} - \mathbf{X}\beta\|^2 + g(\beta) \right\},$$

Generalized cross-validation (GCV) of $\hat{\beta}$ is defined by

$$(\text{GCV of } \hat{\beta}) := \frac{\|\mathbf{y} - \mathbf{X}\hat{\beta}\|^2}{n(1 - \hat{\text{df}}/n)^2} \quad \text{where} \quad \hat{\text{df}} := \text{tr} \left[\mathbf{X} \frac{\partial \hat{\beta}}{\partial \mathbf{y}} \right].$$

Estimator $\hat{\beta}$	Penalty $g(\beta)$	Degrees of freedom $\hat{\text{df}}$
Lasso	$\lambda \ \beta\ _1$	$ \hat{S} $
Ridge	$\frac{\mu}{2} \ \beta\ _2^2$	$\text{tr} [\mathbf{X} (\mathbf{X}^\top \mathbf{X} + n\mu \mathbf{I}_p)^{-1} \mathbf{X}^\top]$
Elastic net	$\lambda \ \beta\ _1 + \frac{\mu}{2} \ \beta\ _2^2$	$\text{tr} [\mathbf{X}_{\hat{S}} (\mathbf{X}_{\hat{S}}^\top \mathbf{X}_{\hat{S}} + n\mu \mathbf{I}_p)^{-1} \mathbf{X}_{\hat{S}}^\top]$

Example of $\hat{\text{df}}$ for specific penalties. Here, $\hat{S} = \{j \in [p] : e_j^\top \hat{\beta} \neq 0\}$ and $\mathbf{X}_{\hat{S}}$ is the sub-matrix of \mathbf{X} made of columns indexed in \hat{S} .

Consistency of Generalized cross-validation

Theorem (Prediction risk \approx GCV)

$$\mathbb{E}[(y_0 - \mathbf{x}_0^\top \hat{\boldsymbol{\beta}})^2 | \mathbf{y}, \mathbf{X}] \approx \frac{\|\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|^2}{n(1 - \hat{d}f/n)^2}$$

	Penalty	Proof
Patil et al. [2021]	Ridge	Random Matrix Theory
Celentano et al. [2023]	Lasso	Convex Gaussian Min-Max Theorem
Bellec and Shen [2022]	strongly convex	Second order Stein's formula

Naive GCV for ensemble estimator

For ensemble estimator $\tilde{\beta} = \frac{1}{M} \sum_{m=1}^M \hat{\beta}_m$, we can think of the naive-GCV:

$$\text{naive-GCV} := \frac{\|\mathbf{y} - \mathbf{X}\tilde{\beta}\|^2}{n(1 - \tilde{\text{df}}/n)^2} \quad \text{where} \quad \tilde{\text{df}} = \text{tr}\left[\mathbf{X} \frac{\partial \tilde{\beta}}{\partial \mathbf{y}}\right]$$

Q. Does the naive-GCV consistently estimate the prediction risk?

$$\mathbb{E}[(y_0 - \mathbf{x}_0^\top \tilde{\beta})^2 | \mathbf{y}, \mathbf{X}] \stackrel{?}{\approx} \text{naive-GCV}$$

A. No. The naive-GCV is inconsistent.

Theorem

Under some regularity condition, there exists some positive constant $C \in (0, 1)$ such that

$$\liminf_{n \rightarrow \infty} \mathbb{P}\left(\left|\frac{\mathbb{E}[(y_0 - \mathbf{x}_0^\top \hat{\beta})^2 | \mathbf{y}, \mathbf{X}]}{\text{naive-GCV}} - 1\right| \geq C\right) \geq C.$$

Overview of main result: corrected-GCV (CGCV)

$$\text{CGCV} := \underbrace{\frac{\|\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}}\|^2}{n(1 - \tilde{\text{df}}/n)^2}}_{=\text{naive-GCV}} - \underbrace{\left(\frac{\tilde{\text{df}}}{n - \tilde{\text{df}}}\right)^2 \left(\frac{n}{k} - 1\right) \frac{1}{M^2} \sum_{m=1}^M \frac{\|\mathbf{y}_{I_m} - \mathbf{X}_{I_m}\hat{\boldsymbol{\beta}}_m\|^2}{k(1 - \hat{\text{df}}_m/k)^2}}_{=:\text{correction}}.$$

Theorem (Informal)

Either assumption (a) or (b) below is satisfied.

Assumption	Distribution	Response $y = f(\mathbf{x}, \epsilon)$	Penalty g
(a)	Gaussian	Linear	strongly convex
(b)	Non-Gaussian	Nonlinear	Ridge

Then, we have (Prediction error) \approx CGCV. More precisely,

$$\mathbb{E}[(y_0 - \mathbf{x}_0^\top \tilde{\boldsymbol{\beta}})^2 | \mathbf{y}, \mathbf{X}] = \begin{cases} \text{CGCV} \cdot (1 + O_p(n^{-1/2})) & \text{under (a)} \\ \text{CGCV} + o_p(1) & \text{under (b)} \end{cases}$$

When correction term is small

The theorem implies

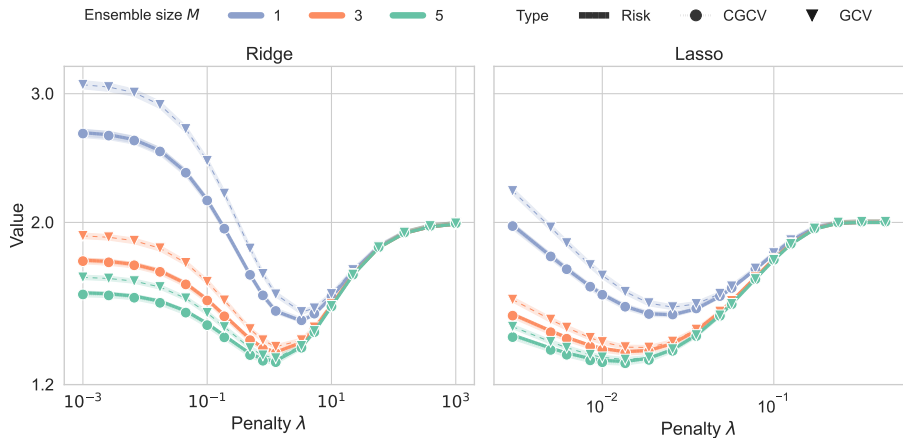
$$(\text{Prediction risk}) \approx \text{CGCV} = \underbrace{\frac{\|\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}}\|^2}{n(1 - \tilde{\text{df}}/n)^2}}_{=\text{naive-GCV}} - \text{correction}$$

where

$$\text{correction} = \left(\frac{\tilde{\text{df}}}{n - \tilde{\text{df}}}\right)^2 \left(\frac{n}{k} - 1\right) \frac{1}{M^2} \sum_{m=1}^M \frac{\|\mathbf{y}_{I_m} - \mathbf{X}_{I_m}\hat{\boldsymbol{\beta}}_m\|^2}{k(1 - \hat{\text{df}}_m/k)^2}.$$

- Naive-GCV overestimates prediction risk.
- Correction term is exactly 0 when sub-sample size k is n .
- Correction term is $O(M^{-1})$.
 \Rightarrow For infinite-ensemble ($M = \infty$), the naive-GCV is consistent.

Comparison of CGCV and naive-GCV



Proof: Second order Stein's fomrula

Theorem (Bellec and Zhang [2021])

For almost surely differentiable function $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ and $\mathbf{z} \sim \mathcal{N}(\mathbf{0}_n, \mathbf{I}_n)$, we have

$$\mathbb{E}\left[\left\{\mathbf{z}^\top \mathbf{f}(\mathbf{z}) - \nabla \cdot \mathbf{f}(\mathbf{z})\right\}^2\right] = \mathbb{E}\left[\|\mathbf{f}(\mathbf{z})\|^2 + \text{tr}\{(\nabla \mathbf{f}(\mathbf{z}))^2\}\right].$$

- Many applications in single index model (Bellec, 2022), multinomial regression (Tan and Bellec, 2023), robust regression (Bellec and Koriyama, 2023).

Summary

- The naive-GCV is inconsistent to the prediction error of ensemble estimators.
- We proposed the corrected GCV and showed its consistency under Gaussian setting and non-Gaussian setting.
- [arXiv:2310.01374](https://arxiv.org/abs/2310.01374)

Reference I

- P. C. Bellec and Y. Shen. Derivatives and residual distribution of regularized M-estimators with application to adaptive tuning. In *Conference on Learning Theory*, 2022.
- P. C. Bellec and C.-H. Zhang. Second-order stein: Sure for sure and other applications in high-dimensional inference. *The Annals of Statistics*, 49(4):1864–1903, 2021.
- M. Celentano, A. Montanari, and Y. Wei. The lasso with general gaussian designs with applications to hypothesis testing. *The Annals of Statistics*, 51(5):2194–2220, 2023.
- J.-H. Du, P. Patil, and A. K. Kuchibhotla. Subsample ridge ensembles: Equivalences and generalized cross-validation. In *International Conference on Machine Learning*, 2023.
- P. Patil, Y. Wei, A. Rinaldo, and R. Tibshirani. Uniform consistency of cross-validation estimators for high-dimensional ridge regression. In *International Conference on Artificial Intelligence and Statistics*, 2021.

Appendix

Consistency of CGCV under assumption (a)

Assumption (a)

- $(y_i, \mathbf{x}_i)_{i=1}^n \in \mathbb{R} \times \mathbb{R}^p$ are iid distributed according to

$$y_i = \mathbf{x}_i^\top \boldsymbol{\beta}^* + \epsilon_i, \quad \mathbf{x}_i \sim \mathcal{N}(\mathbf{0}_p, \boldsymbol{\Sigma}), \quad \epsilon_i \sim \mathcal{N}(0, \sigma^2)$$

for some $\boldsymbol{\beta}^* \in \mathbb{R}^p$, $\boldsymbol{\Sigma} \succ 0$ and $\sigma > 0$.

- g is strongly convex with respect to $\boldsymbol{\Sigma}^a$ (e.g., Ridge, Elastic net).
- $p = O(k)$ for sub-sample size k .

^athe map $\boldsymbol{\beta} \mapsto g(\boldsymbol{\beta}) - \mu \boldsymbol{\beta}^\top \boldsymbol{\Sigma} \boldsymbol{\beta}$ is convex for some $\mu > 0$

Theorem (Prediction risk \approx GCCV)

If the assumption (a) is satisfied, we have

$$\mathbb{E}[(y_0 - \mathbf{x}_0^\top \tilde{\boldsymbol{\beta}})^2 | \mathbf{y}, \mathbf{X}] = [1 + O_P(n^{-1/2})] \cdot \text{CGCV} \quad \text{as } n \rightarrow \infty$$

Consistency of CGCV under Assumption (b)

Assumption (b)

- $g(\boldsymbol{\beta}) = \lambda \|\boldsymbol{\beta}\|^2$ for some $\lambda > 0$.
- $\mathbb{E}[y_i] = 0$ and $\mathbb{E}[y_i^{4+\delta}] < +\infty$ for some $\delta > 0$.
- $\mathbf{x}_i = \Sigma^{1/2} \mathbf{z}_i$ for some $\Sigma \succ 0$ and $\mathbf{z}_i \in \mathbb{R}^p$ has iid entries such that $\mathbb{E}[z_{ij}] = 0$, $\mathbb{E}[z_{ij}^2] = 1$, and $\mathbb{E}[z_{ij}^{4+\delta}] < +\infty$.
- $p/n \rightarrow \phi \in (0, \infty)$, $p/k \rightarrow \psi \in [\phi, \infty]$.

Theorem

$$\mathbb{E}[(y_0 - \mathbf{x}_0^\top \tilde{\boldsymbol{\beta}})^2 | \mathbf{y}, \mathbf{X}] = \text{CGCV} + o_P(1) \quad \text{as } n \rightarrow +\infty$$

Proof outline

Prediction risk of $\hat{\beta}$, denoted by $R(\hat{\beta})$, can be written as

$$\begin{aligned}R(\hat{\beta}) &= \mathbb{E}[(y_0 - \mathbf{x}_0^\top \hat{\beta})^2 | \mathbf{y}, \mathbf{X}] \\&= \mathbb{E}\left[\{\epsilon_0 - \mathbf{x}_0^\top (\hat{\beta} - \beta^*)\}^2 | \mathbf{y}, \mathbf{X}\right] \quad \text{by } y_0 = \mathbf{x}_0^\top \beta^* + \epsilon_0 \\&= \sigma^2 + (\hat{\beta} - \beta^*)^\top \Sigma (\hat{\beta} - \beta^*) \quad \text{by } \mathbf{x}_0 \sim \mathcal{N}(\mathbf{0}_p, \Sigma), \epsilon_0 \sim \mathcal{N}(0, \sigma^2).\end{aligned}$$

Thus, the prediction risk of the ensemble $\tilde{\beta} = \frac{1}{M} \sum_{m=1}^M \hat{\beta}_m$ is given by

$$\begin{aligned}R(\tilde{\beta}) &= \sigma^2 + \left\{ \left(\frac{1}{M} \sum_{m=1}^M \hat{\beta}_m \right) - \beta^* \right\}^\top \Sigma \left\{ \left(\frac{1}{M} \sum_{m=1}^M \hat{\beta}_m \right) - \beta^* \right\} \\&= \frac{1}{M^2} \sum_{m=1}^M \sum_{\ell=1}^M \left[\sigma^2 + (\hat{\beta}_m - \beta^*)^\top \Sigma (\hat{\beta}_\ell - \beta^*) \right].\end{aligned}$$

Proof outline

The naive-GCV for $\tilde{\beta} = \frac{1}{M} \sum_{m=1}^M \hat{\beta}_m$ is given by

$$\text{naive-GCV} = \frac{\|\mathbf{y} - \mathbf{X}\tilde{\beta}\|^2}{n(1 - \tilde{\text{df}}/n)^2} = \frac{\frac{1}{M^2} \sum_{m=1}^M \sum_{\ell=1}^M (\mathbf{y} - \mathbf{X}\hat{\beta}_m)^\top (\mathbf{y} - \mathbf{X}\hat{\beta}_\ell)}{n(1 - \tilde{\text{df}}/n)^2}$$

Lemma

For all $m, \ell \in [M]$, we have

$$(\mathbf{y} - \mathbf{X}\hat{\beta}_m)^\top (\mathbf{y} - \mathbf{X}\hat{\beta}_\ell) \approx \left[\sigma^2 + (\hat{\beta}_m - \beta^*)^\top \Sigma (\hat{\beta}_\ell - \beta^*) \right] \cdot D_{m\ell},$$

$$\text{where } D_{m\ell} = n - \text{df}_m - \text{df}_\ell + \frac{\hat{\text{df}}_m \hat{\text{df}}_\ell}{|I_m||I_\ell|} |I_m \cap I_\ell|.$$

Using this lemma,

$$\text{naive-GCV} \approx \frac{1}{M^2} \sum_{m=1}^M \sum_{\ell=1}^M \left[\sigma^2 + (\hat{\beta}_m - \beta^*)^\top \Sigma (\hat{\beta}_\ell - \beta^*) \right] \cdot \frac{D_{m\ell}}{n(1 - \tilde{\text{df}}/n)^2}$$

Proof outline

Lemma (Concentration of $D_{m,\ell}$)

$$\frac{D_{m,\ell}}{n(1 - \tilde{d}f/n)^2} \approx 1 + \mathbf{1}\{m = \ell\} \cdot \left(\frac{n}{k} - 1\right) \frac{(\tilde{d}f/n)^2}{(1 - \tilde{d}f/n)^2}.$$

$$\begin{aligned} \text{naive-GCV} &\approx \frac{1}{M^2} \sum_{m=1}^M \sum_{\ell=1}^M \left[\sigma^2 + (\hat{\beta}_m - \beta^*)^\top \Sigma (\hat{\beta}_\ell - \beta^*) \right] \\ &+ \frac{1}{M^2} \sum_{m=1}^M \left[\sigma^2 + (\hat{\beta}_m - \beta^*)^\top \Sigma (\hat{\beta}_m - \beta^*) \right] \cdot \left(\frac{n}{k} - 1\right) \frac{(\tilde{d}f/n)^2}{(1 - \tilde{d}f/n)^2} \\ &= R(\tilde{\beta}) + \frac{1}{M^2} \sum_{m=1}^M R(\hat{\beta}_m) \cdot \left(\frac{n}{k} - 1\right) \frac{(\tilde{d}f/n)^2}{(1 - \tilde{d}f/n)^2} \end{aligned}$$

Obtain CGCV

We have shown that

$$R(\tilde{\beta}) \approx \text{naive-GCV} - \frac{1}{M^2} \left(\frac{n}{k} - 1 \right) \frac{(\tilde{\text{df}}/n)^2}{(1 - \tilde{\text{df}}/n)^2} \sum_{m=1}^M R(\hat{\beta}_m).$$

Using (prediction risk of $\hat{\beta}_m$) \approx (GCV of $\hat{\beta}_m$ fitted on $(y_i, \mathbf{x}_i)_{i \in I_m}$)

$$R(\hat{\beta}_m) \approx \frac{\|\mathbf{y}_{I_m} - \mathbf{X}_{I_m} \hat{\beta}_m\|^2}{k(1 - \hat{\text{df}}_m/k)},$$

we are left with

$$R(\tilde{\beta}) \approx \underbrace{(\text{naive-GCV}) - \frac{1}{M^2} \left(\frac{n}{k} - 1 \right) \frac{(\tilde{\text{df}}/n)^2}{(1 - \tilde{\text{df}}/n)^2} \sum_{m=1}^M \frac{\|\mathbf{y}_{I_m} - \mathbf{X}_{I_m} \hat{\beta}_m\|^2}{k(1 - \hat{\text{df}}_m/k)^2}}_{=\text{CGCV}}$$

Proof of Lemma 1: Second order Stein's formula

Recall that Lemma 1 claims

$$\left(\sigma^2 + (\hat{\boldsymbol{\beta}}_m - \boldsymbol{\beta}^*)^\top \boldsymbol{\Sigma}(\hat{\boldsymbol{\beta}}_\ell - \boldsymbol{\beta}^*)\right) \cdot D_{m\ell} \approx (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}_m)^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}_\ell),$$

where $D_{m\ell} = n - \text{df}_m - \text{df}_\ell + \frac{\hat{\text{df}}_m \hat{\text{df}}_\ell}{|I_m||I_\ell|} |I_m \cap I_\ell|$.

Theorem (Bellec and Zhang [2021])

For almost surely differentiable function $\mathbf{f} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ and $\mathbf{z} \sim \mathcal{N}(\mathbf{0}_n, \mathbf{I}_n)$, we have

$$\mathbb{E} \left[\left\{ \mathbf{z}^\top \mathbf{f}(\mathbf{z}) - \nabla \cdot \mathbf{f}(\mathbf{z}) \right\}^2 \right] = \mathbb{E} \left[\|\mathbf{f}(\mathbf{z})\|^2 + \text{tr} \{ (\nabla \mathbf{f}(\mathbf{z}))^2 \} \right].$$