

# Information theoretic analysis and coding scheme for learning : regression and more

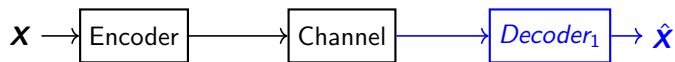
**Jiahui WEI**, Philippe MARY and Elsa DUPRAZ

INSA Rennes, IMT Atlantique

June 7, 2024

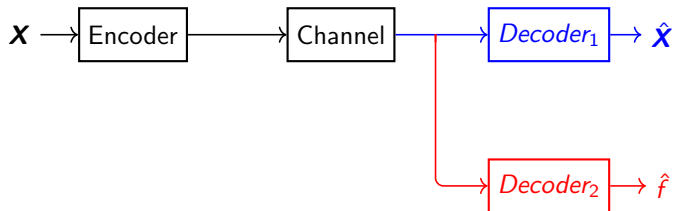


# Context



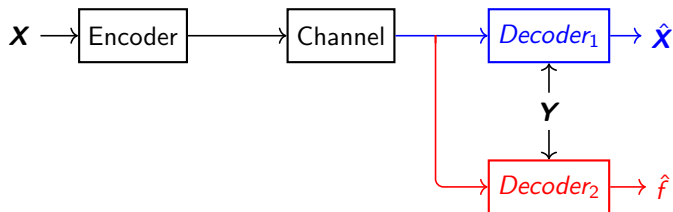
- **Conventional communication** : reconstruct the source even with distortion => Rate-Distortion theory

# Context



- **Conventional communication** : reconstruct the source even with distortion => Rate-Distortion theory
- **Goal-oriented communication** : construct the coding scheme to address some machine learning problems

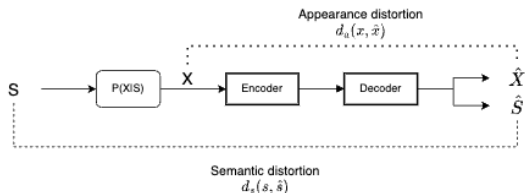
# Context



- **Conventional communication** : reconstruct the source even with distortion => Rate-Distortion theory
- **Goal-oriented communication** : construct the coding scheme to address some machine learning problems
  
- **Question** : Do we need the same method for coding?

# State of the art

- Existing works : **Rate-distortion framework with semantic and appearance distortion**<sup>1</sup>; Rate for parameter estimation<sup>2</sup>; Hypothesis testing<sup>3</sup>; Rate-distortion-perception trade-off<sup>4</sup>.



<sup>1</sup>J. Liu, W. Zhang, and H. V. Poor, "A rate-distortion framework for characterizing semantic information," in 2021 IEEE International Symposium on Information Theory (ISIT). IEEE, 2021, pp. 2894–2899.

<sup>2</sup>M. El Gamal and L. Lai, "Are slepian-wolf rates necessary for distributed parameter estimation?" in 2015 53rd Annual Allerton Conference on Communication, Control, and Computing (Allerton). IEEE, 2015

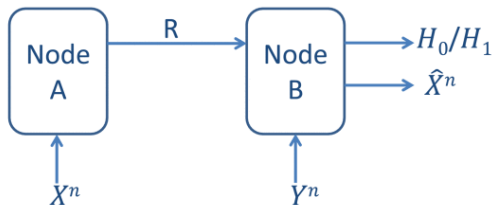
<sup>3</sup>G. Katz, P. Piantanida and M. Debbah, "Distributed Binary Detection With Lossy Data Compression," in IEEE Transactions on Information Theory. Aug. 2017

<sup>4</sup>Y. Blau and T. Michaeli, "Rethinking lossy compression: The rate-distortion-perception tradeoff," in International Conference on Machine Learning. PMLR, 2019

<sup>5</sup>M. Raginsky, "Learning from compressed observations," in IEEE Information Theory Workshop, 2007

# State of the art

- **Existing works** : Rate-distortion framework with semantic and appearance distortion<sup>1</sup>; Rate for parameter estimation<sup>2</sup>; **Hypothesis testing**<sup>3</sup>; Rate-distortion-perception trade-off<sup>4</sup>.



<sup>1</sup>J. Liu, W. Zhang, and H. V. Poor, "A rate-distortion framework for characterizing semantic information," in 2021 IEEE International Symposium on Information Theory (ISIT). IEEE, 2021, pp. 2894–2899.

<sup>2</sup>M. El Gamal and L. Lai, "Are slepian-wolf rates necessary for distributed parameter estimation?" in 2015 53rd Annual Allerton Conference on Communication, Control, and Computing (Allerton). IEEE, 2015

<sup>3</sup>G. Katz, P. Piantanida and M. Debbah, "Distributed Binary Detection With Lossy Data Compression," in IEEE Transactions on Information Theory. Aug. 2017

<sup>4</sup>Y. Blau and T. Michaeli, "Rethinking lossy compression: The rate-distortion-perception tradeoff," in International Conference on Machine Learning. PMLR, 2019

<sup>5</sup>M. Raginsky, "Learning from compressed observations," in IEEE Information Theory Workshop, 2007

# State of the art

- **Existing works** : Rate-distortion framework with semantic and appearance distortion<sup>1</sup>; Rate for parameter estimation<sup>2</sup>; Hypothesis testing<sup>3</sup>; Rate-distortion-perception trade-off<sup>4</sup>.
  - Rate-distortion + **Goal**;
  - **Trade-off** between the task and data reconstruction

<sup>1</sup>J. Liu, W. Zhang, and H. V. Poor, "A rate-distortion framework for characterizing semantic information," in 2021 IEEE International Symposium on Information Theory (ISIT). IEEE, 2021, pp. 2894–2899.

<sup>2</sup>M. El Gamal and L. Lai, "Are slepian-wolf rates necessary for distributed parameter estimation?" in 2015 53rd Annual Allerton Conference on Communication, Control, and Computing (Allerton). IEEE, 2015

<sup>3</sup>G. Katz, P. Piantanida and M. Debbah, "Distributed Binary Detection With Lossy Data Compression," in IEEE Transactions on Information Theory. Aug. 2017

<sup>4</sup>Y. Blau and T. Michaeli, "Rethinking lossy compression: The rate-distortion-perception tradeoff," in International Conference on Machine Learning. PMLR, 2019

<sup>5</sup>M. Raginsky, "Learning from compressed observations," in IEEE Information Theory Workshop, 2007

# State of the art

- **Existing works** : Rate-distortion framework with semantic and appearance distortion<sup>1</sup>; Rate for parameter estimation<sup>2</sup>; Hypothesis testing<sup>3</sup>; Rate-distortion-perception trade-off<sup>4</sup>.
  - Rate-distortion + **Goal**;
  - **Trade-off** between the task and data reconstruction
  
- **Regression**
  - A fundamental statistical method;
  - A rate-loss bound for general regression with side information is provided by Raginsky<sup>5</sup>;
  - This bound is **loose** and the **trade-off** is not investigated.

<sup>1</sup>J. Liu, W. Zhang, and H. V. Poor, "A rate-distortion framework for characterizing semantic information," in 2021 IEEE International Symposium on Information Theory (ISIT). IEEE, 2021, pp. 2894–2899.

<sup>2</sup>M. El Gamal and L. Lai, "Are slepian-wolf rates necessary for distributed parameter estimation?" in 2015 53rd Annual Allerton Conference on Communication, Control, and Computing (Allerton). IEEE, 2015

<sup>3</sup>G. Katz, P. Piantanida and M. Debbah, "Distributed Binary Detection With Lossy Data Compression," in IEEE Transactions on Information Theory. Aug. 2017

<sup>4</sup>Y. Blau and T. Michaeli, "Rethinking lossy compression: The rate-distortion-perception tradeoff," in International Conference on Machine Learning. PMLR, 2019

<sup>5</sup>M. Raginsky, "Learning from compressed observations," in IEEE Information Theory Workshop, 2007



# Outline

- 1 Problem statement
- 2 Asymptotic rate-generalization error regions
- 3 Non-asymptotic rate-distortion-generalization error regions
- 4 Practical coding scheme
- 5 Conclusion and perspectives

# Outline

- 1 Problem statement
- 2 Asymptotic rate-generalization error regions
- 3 Non-asymptotic rate-distortion-generalization error regions
- 4 Practical coding scheme
- 5 Conclusion and perspectives

# Problem statement

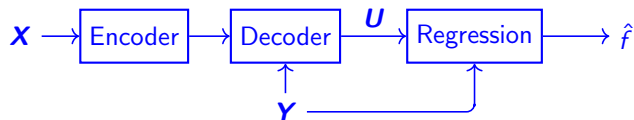


Figure 1: Coding scheme for regression

**Regression** model for  $X$  and  $Y$ :

$$X = f(Y) + N, \quad (1)$$

where  $N \sim \mathcal{N}(0, \sigma^2)$  independent from  $X$  and  $Y$ .

- 1 Training phase with sequence  $\mathbf{Z} = (\mathbf{U}, \mathbf{Y}) \Rightarrow \hat{f}$ ;

# Problem statement

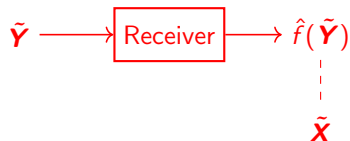


Figure 1: Inference phase

**Regression** model for  $X$  and  $Y$ :

$$X = f(Y) + N, \quad (1)$$

where  $N \sim \mathcal{N}(0, \sigma^2)$  independent from  $X$  and  $Y$ .

- 1 Training phase with sequence  $\mathbf{Z} = (\mathbf{U}, \mathbf{Y}) \Rightarrow \hat{f}$ ;
- 2 Inference phase  $\Rightarrow$  Generalization error.

## Parametric regression : OLS

$$\mathbf{X} = \boldsymbol{\beta}^T \mathbf{Y}^* + N = \sum_{i=0}^{k-1} \beta_i h_i(Y) + N, \quad (2)$$

where  $\mathbf{Y}_j^* = [h_0(Y_j), \dots, h_{k-1}(Y_j)]^T$  and  $\boldsymbol{\beta} = [\beta_0, \dots, \beta_{k-1}]^T$  is unknown. Further define  $\underline{\mathbf{Y}}^* = [\mathbf{Y}_1^*, \dots, \mathbf{Y}_n^*] \in \mathbb{R}^{k \times n}$ , we have

- OLS estimation between  $\mathbf{X}$  and  $\mathbf{Y}$  is given by<sup>6</sup>

$$\hat{\boldsymbol{\beta}} = \left( \underline{\mathbf{Y}}^* \underline{\mathbf{Y}}^{*T} \right)^{-1} \underline{\mathbf{Y}}^* \mathbf{X}. \quad (3)$$

- Properties :

$$\mathbb{E} \left[ \hat{\boldsymbol{\beta}} \right] = \boldsymbol{\beta} \quad \text{and} \quad \mathbb{C} \left[ \hat{\boldsymbol{\beta}} | \mathbf{Y} \right] = \sigma_{X|Y}^2 \left( \underline{\mathbf{Y}}^* \underline{\mathbf{Y}}^{*T} \right)^{-1}, \quad (4)$$

where  $\mathbb{C} \left[ \hat{\boldsymbol{\beta}} | \mathbf{Y} \right]$  is the covariance matrix of  $\hat{\boldsymbol{\beta}}$  given  $\mathbf{Y}$  and  $\sigma_{X|Y}^2$  is the conditional variance of  $X$  given  $Y$ .

<sup>6</sup>Chapter 7, A. C. Rencher and G. B. Schaalje, Linear models in statistics. John Wiley & Sons, 2008

# Non-parametric regression : kernel regression

$$X = f(Y) + N \quad (5)$$

without any prior knowledge of the regression form.

- A one-dim kernel is any smooth and symmetric function  $K : \mathbb{R} \rightarrow \mathbb{R}$  such that  $\forall x \in \mathbb{R}, K(x) \geq 0$ , and the following relations hold

$$\int_{\mathbb{R}} K(x) dx = 1, \quad \int_{\mathbb{R}} xK(x) dx = 0, \quad \text{and} \quad 0 \leq \int_{\mathbb{R}} x^2 K(x) dx < \infty. \quad (6)$$

- The Nadaraya-Watson Kernel regression over  $(\mathbf{X}, \mathbf{Y})$  is defined as <sup>7</sup>:

$$\hat{f}(Y) = \frac{\sum_{j=1}^n K\left(\frac{Y-Y_j}{h}\right) X_j}{\sum_{j=1}^n K\left(\frac{Y-Y_j}{h}\right)}. \quad (7)$$

<sup>7</sup>L. Wasserman, All of Nonparametric Statistics (Springer Texts in Statistics). Berlin, Heidelberg: Springer-Verlag, 2006.

# Non-parametric regression : kernel regression

$$X = f(Y) + N \quad (5)$$

without any prior knowledge of the regression form.

- A one-dim kernel is any smooth and symmetric function  $K : \mathbb{R} \rightarrow \mathbb{R}$  such that  $\forall x \in \mathbb{R}, K(x) \geq 0$ , and the following relations hold

$$\int_{\mathbb{R}} K(x) dx = 1, \quad \int_{\mathbb{R}} xK(x) dx = 0, \quad \text{and} \quad 0 \leq \int_{\mathbb{R}} x^2 K(x) dx < \infty. \quad (6)$$

- The Nadaraya-Watson Kernel regression over  $(\mathbf{X}, \mathbf{Y})$  is defined as <sup>7</sup>:

$$\hat{f}(Y) = \frac{\sum_{j=1}^n K\left(\frac{Y-Y_j}{h}\right) X_j}{\sum_{j=1}^n K\left(\frac{Y-Y_j}{h}\right)}. \quad (7)$$

**Attention :**  $X$  is not available in our setup so the regression needs to be processed with the compressed observation  $U$ .

<sup>7</sup>L. Wasserman, All of Nonparametric Statistics (Springer Texts in Statistics). Berlin, Heidelberg: Springer-Verlag, 2006.

# Definitions

## Definition

A regression scheme at rate  $R$  is defined by a sequence  $\{(e_n, d_n, R, \mathcal{L}_n)\}$  with

an encoder  $e_n : \mathcal{X}^n \rightarrow \{1, 2, \dots, M_n\}$ ,

a decoder  $d_n : \mathcal{Y}^n \times \{1, 2, \dots, M_n\} \rightarrow \mathcal{U}^n$ ,

and the learner  $\mathcal{L}_n : \mathcal{Y}^n \times \mathcal{U}^n \rightarrow \mathcal{F}$ ,

such that

$$\limsup_{n \rightarrow \infty} \frac{\log M_n}{n} \leq R.$$

Loss function  $\ell(x, \hat{x}) = (x - \hat{x})^2$ .

For a fixed function  $f$ , the **expected loss** is defined as

$$L(f) = \mathbb{E} [\ell(X, f(Y))] \tag{8}$$

and the **minimum expected loss** is defined as the

$$L^*(\mathcal{F}) = \inf_{f \in \mathcal{F}} L(f) = \sigma^2 \tag{9}$$



# Generalization error

The **generalization error** is defined as

$$G(\hat{f}^{(n)}, \mathbf{Z}) = \mathbb{E}_{\tilde{X}\tilde{Y}} \left[ \ell \left( \tilde{X}, \hat{f}^{(n)}(\mathbf{Z}, \tilde{Y}) \right) \mid \mathbf{Z} \right]. \quad (10)$$

where  $(\tilde{X}, \tilde{Y}) \sim P_{XY}$  is independent from  $\mathbf{Z}$  (i.e.  $(\mathbf{U}^n, \mathbf{Y}^n)$ ).

## Objective

*Derive the rate-generalization error regions*

2 methods of regression : **Parametric regression using OLS estimator** and **Kernel regression**.

# Generalization error

The **generalization error** is defined as

$$G(\hat{f}^{(n)}, \mathbf{Z}) = \mathbb{E}_{\tilde{X}\tilde{Y}} \left[ \ell \left( \tilde{X}, \hat{f}^{(n)}(\mathbf{Z}, \tilde{Y}) \right) \mid \mathbf{Z} \right]. \quad (10)$$

where  $(\tilde{X}, \tilde{Y}) \sim P_{XY}$  is independent from  $\mathbf{Z}$  (i.e.  $(\mathbf{U}^n, \mathbf{Y}^n)$ ).

## Objective

*Derive the rate-generalization error regions*

2 methods of regression : **Parametric regression using OLS estimator** and **Kernel regression**.

### Our contributions :

- The rate-generalization error regions in both asymptotic and non-asymptotic regime;
- Improvement of the upper bound provided by Raginsky;
- Investigation of reconstruction vs regression trade-off.

# Outline

- 1 Problem statement
- 2 Asymptotic rate-generalization error regions**
- 3 Non-asymptotic rate-distortion-generalization error regions
- 4 Practical coding scheme
- 5 Conclusion and perspectives

## Asymptotic Rate-generalization error regions

An  $(\mathbf{n}, \mathbf{M}, \mathbf{G})$  **code** for regression is a code with  $|e_n| = M$  such that  $\mathbb{E}_{\mathbf{Z}} \left[ G(\hat{f}^{(n)}, \mathbf{Z}) \right] \leq G$ .

### Definition

A pair  $(R, \delta)$  is said to be *achievable* if an  $(n, M, G)$ -code exists such as

$$\limsup_{n \rightarrow \infty} \mathbb{E}_{\mathbf{Z}} \left[ G(\hat{f}^{(n)}, \mathbf{Z}) \right] \leq L^*(\mathcal{F}, \mathbf{Z}) + \delta \quad (11)$$

Recall the result of Raginsky<sup>8</sup>

$$L^{\star \frac{1}{2}}(\mathcal{F}, \mathbf{Z}) \leq \limsup_{n \rightarrow \infty} \mathbb{E} \left[ G(\hat{f}^{(n)}, \mathbf{Z})^{\frac{1}{2}} \right] \leq L^{\star \frac{1}{2}}(\mathcal{F}, \mathbf{Z}) + 2\mathbb{D}_{X|Y}(R)^{1/2} \quad (12)$$

<sup>8</sup>M. Raginsky, "Learning from compressed observations," in IEEE Information Theory Workshop, 2007

## Asymptotic Rate-generalization error regions

An  $(\mathbf{n}, \mathbf{M}, \mathbf{G})$  **code** for regression is a code with  $|e_n| = M$  such that  $\mathbb{E}_{\mathbf{Z}} \left[ G(\hat{f}^{(n)}, \mathbf{Z}) \right] \leq G$ .

### Definition

A pair  $(R, \delta)$  is said to be *achievable* if an  $(n, M, G)$ -code exists such as

$$\limsup_{n \rightarrow \infty} \mathbb{E}_{\mathbf{Z}} \left[ G(\hat{f}^{(n)}, \mathbf{Z}) \right] \leq L^*(\mathcal{F}, \mathbf{Z}) + \delta \quad (11)$$

Recall the result of Raginsky<sup>8</sup>

$$L^{\star \frac{1}{2}}(\mathcal{F}, \mathbf{Z}) \leq \limsup_{n \rightarrow \infty} \mathbb{E} \left[ G(\hat{f}^{(n)}, \mathbf{Z})^{\frac{1}{2}} \right] \leq L^{\star \frac{1}{2}}(\mathcal{F}, \mathbf{Z}) + 2\mathbb{D}_{X|Y}(R)^{1/2} \quad (12)$$

### Theorem (Parametric & Kernel regression)

Given any rate  $R > 0$ , the pair  $(\mathbf{R}, \mathbf{0})$  is achievable for the **parametric regression** and **kernel regression** with squared loss, for sources  $(X, Y)$  following the regression model, that is

$$\limsup_{n \rightarrow \infty} \mathbb{E} \left[ G(\hat{f}^{(n)}, \mathbf{Z}) \right] = L^*(\mathcal{F}, \mathbf{Z}) \quad \text{and} \quad \delta = 0. \quad (13)$$

<sup>8</sup>M. Raginsky, "Learning from compressed observations," in IEEE Information Theory Workshop, 2007

## Sketch of proof : Parametric regression

- Idea : **quantization + binning**. Consider a **Gaussian test channel**

$$U = \alpha(X + \Phi), \quad \text{with } \Phi \sim \mathcal{N}(0, \sigma_\Phi^2)$$

**Remark** : The method is based on **prefix transmission of types**<sup>9</sup> of observations + **binning** of conventional WZ coding.

<sup>9</sup>S. C. Draper, "Universal incremental slepian-wolf coding," in Proc. 42nd Allerton Conf. on Communication, Control and Computing. Citeseer, 2004

## Sketch of proof : Parametric regression

- Idea : **quantization + binning**. Consider a **Gaussian test channel**

$$U = \alpha(X + \Phi), \quad \text{with } \Phi \sim \mathcal{N}(0, \sigma_\Phi^2)$$

- For a training sequences  $(\mathbf{y}, \mathbf{u})$ , the OLS estimator  $\hat{\beta}$  becomes

$$\hat{\beta} = \alpha^{-1}(\underline{\mathbf{Y}}\underline{\mathbf{Y}}^T)^{-1}\underline{\mathbf{Y}}\mathbf{u}. \quad (14)$$

where  $\underline{\mathbf{Y}} = \begin{bmatrix} y_1^0 & \dots & y_n^0 \\ \dots & \dots & \dots \\ y_1^{k-1} & \dots & y_n^{k-1} \end{bmatrix}$ .

- The generalization error can be rewritten as

$$G(\hat{f}^{(n)}, \mathbf{Z}) = \sigma^2 + [\beta - \hat{\beta}]^T \mathbb{E}_{\tilde{\mathbf{Y}}} [\tilde{\mathbf{Y}}\tilde{\mathbf{Y}}^T] [\beta - \hat{\beta}] \quad (15)$$

- Let  $\tilde{\Sigma} = \mathbb{E}_{\tilde{\mathbf{Y}}} [\tilde{\mathbf{Y}}\tilde{\mathbf{Y}}^T]$ ,  $\Sigma = \frac{1}{n}\underline{\mathbf{Y}}\underline{\mathbf{Y}}^T$  and  $C = \frac{\lambda_{\max}(\tilde{\Sigma})}{\lambda_{\min}(\tilde{\Sigma})}$ , the **expected generalization error** :

$$\mathbb{E}_{\mathbf{Z}} [G(\hat{f}^{(n)}, \mathbf{Z})] = \sigma^2 + \frac{\sigma^2 + \sigma_\Phi^2}{n} \mathbb{E} [\text{Tr}(\tilde{\Sigma}\Sigma^{-1})] \leq \sigma^2 + \frac{(\sigma^2 + \sigma_\Phi^2)}{n} kC \quad (16)$$

## Sketch of proof : Kernel regression

- The same test channel, we suppose the following conditions for kernel regression :  
 $Y$  **bounded**;  $p_Y$  continuously differentiable and positively bounded;  $\exists f', f''$ ;  
 $n \rightarrow \infty$ ,  $h \rightarrow \mathbf{0}$  and  $nh \rightarrow \infty$ .



## Sketch of proof : Kernel regression

- The same test channel, we suppose the following conditions for kernel regression :  
 $Y$  **bounded**;  $p_Y$  continuously differentiable and positively bounded;  $\exists f', f''$ ;  
 $\mathbf{n} \rightarrow \infty$ ,  $\mathbf{h} \rightarrow \mathbf{0}$  and  $\mathbf{nh} \rightarrow \infty$ .
- The kernel regression between the sequence  $(\mathbf{y}, \mathbf{u})$  becomes

$$\hat{f}(y) = \frac{\sum_{i=1}^n K\left(\frac{y-y_i}{h}\right) \frac{u_i}{\alpha}}{\sum_{i=1}^n K\left(\frac{y-y_i}{h}\right)}. \quad (17)$$

- The generalization error can be rewritten as

$$\mathbb{E}_{\mathbf{Z}} \left[ G(\hat{f}^{(n)}, \mathbf{Z}) \right] = \sigma^2 + \int b_n^2(\tilde{y}) p_Y(\tilde{y}) d\tilde{y} + \int V_n(\tilde{y}) p_Y(\tilde{y}) d\tilde{y}. \quad (18)$$

where  $b_n(\tilde{y}) = \mathbb{E} \left[ \hat{f}^{(n)}(\tilde{y}, \mathbf{Z}) - f(\tilde{y}) \right]$  is the bias and  $V_n(\tilde{y}) = \mathbb{V} \left[ \hat{f}^{(n)}(\tilde{y}, \mathbf{Z}) \right]$  is the variance of the estimator  $\hat{f}^{(n)}$  with respect to the training sequence  $\mathbf{Z}$ .

- We proved that

$$b_n(\tilde{y}) = \frac{h^2}{2} \left( 2 \frac{f'(\tilde{y}) p_Y'(\tilde{y})}{p_Y(\tilde{y})} + f''(\tilde{y}) \right) \int_{\mathbb{R}} u^2 K(u) du + o(h^2), \quad (19)$$

$$V_n(\tilde{y}) = \frac{(\sigma^2 + \sigma_{\Phi}^2)}{p_Y(\tilde{y}) nh} \int_{\mathbb{R}} K^2(u) du + o\left(\frac{1}{nh}\right). \quad (20)$$

# Regression-Reconstruction trade-off

## Corollary

There is **no trade-off** in terms of coding rate between distortion and regression generalization error. That is

$$R(G, D) = R(D) \quad (21)$$

where  $R(G, D)$  is the communication rate under reconstruction and regression constraints.

# Regression-Reconstruction trade-off

## Corollary

There is **no trade-off** in terms of coding rate between distortion and regression generalization error. That is

$$R(G, D) = R(D) \quad (21)$$

where  $R(G, D)$  is the communication rate under reconstruction and regression constraints.

- Idea of the proof : show that the scheme provided for regression achieves the optimal rate-distortion region for reconstruction.

# Regression-Reconstruction trade-off

## Corollary

There is **no trade-off** in terms of coding rate between distortion and regression generalization error. That is

$$R(G, D) = R(D) \quad (21)$$

where  $R(G, D)$  is the communication rate under reconstruction and regression constraints.

- Idea of the proof : show that the scheme provided for regression achieves the optimal rate-distortion region for reconstruction.
- For  $X|Y \sim \mathcal{N}$ , by replacing  $\sigma_{\Phi}^2 = \frac{D\sigma^2}{\sigma^2 - D}$ , we obtain

$$R_b(D) = \frac{1}{2} \log \left( \frac{\sigma^2 + \sigma_{\Phi}^2}{\sigma_{\Phi}^2} \right) = \frac{1}{2} \log \left( \frac{\sigma^2}{D} \right) \quad (22)$$

# Outline

- 1 Problem statement
- 2 Asymptotic rate-generalization error regions
- 3 Non-asymptotic rate-distortion-generalization error regions**
- 4 Practical coding scheme
- 5 Conclusion and perspectives

# Non-asymptotic rate-distortion-generalization error regions

In finite block-length  $n$ , **the excess error probability** is  $\mathbb{P} \left[ G(\hat{f}^{(n)}, \beta) \geq G \text{ or } d(X, \hat{X}) \geq D \right]$ .

# Non-asymptotic rate-distortion-generalization error regions

In finite block-length  $n$ , **the excess error probability** is  $\mathbb{P} \left[ G(\hat{f}^{(n)}, \beta) \geq G \text{ or } d(X, \hat{X}) \geq D \right]$ .

## Definition

An  $(n, M, G, D, \varepsilon)$  code for the sequence  $\{(e_n, d_n, R, \hat{f}^{(n)})\}$  and  $\varepsilon \in (0, 1)$  is a code with  $|e_n| = M$  such that

$$\mathbb{P} \left[ G(\hat{f}^{(n)}, \beta) \geq G \text{ or } d(X, \hat{X}) \geq D \right] \leq \varepsilon \text{ and } \frac{\log M}{n} \leq R. \quad (23)$$

## Definition

For fixed  $G$  and block-length  $n$ , the finite block-length rate-loss functions with excess loss is defined by:

$$R(n, G, D, \varepsilon) = \inf_R \{ \exists (n, M, G, D, \varepsilon) \text{ code} \} \quad (24)$$

# Non-asymptotic achievable regions

Define the **loss-information density**  $\mathbf{i}$  as  $\mathbf{i}(U, X, Y, \hat{X}) :=$

$$\begin{bmatrix} -\log \frac{P_{U|Y}(U|Y)}{P_U(U)} \\ \log \frac{P_{U|X}(U|X)}{P_U(U)} \\ \mathbb{E}_{\tilde{X}\tilde{Y}} \left[ \ell(\tilde{X}, \hat{f}^{(n)}(\mathbf{Z}, \tilde{Y})) \right] \\ d(X, \hat{X}) \end{bmatrix},$$

Let  $\mathbf{J} = \mathbb{E}_{UXY\hat{X}} [\mathbf{i}]$  and  $\mathbf{V} = \mathbb{C} [\mathbf{i}(U, X, Y, \hat{X})]$ .

The **dispersion region** is defined as  $\mathcal{S}(\mathbf{V}, \varepsilon) := \{\mathbf{b} \in \mathbb{R}^k : \Pr(\mathbf{B} \leq \mathbf{b}) \geq 1 - \varepsilon\}$  with  $\mathbf{B} \sim \mathcal{N}(0, \mathbf{V})$ .



# Non-asymptotic achievable regions

Define the **loss-information density**  $\mathbf{i}$  as  $\mathbf{i}(U, X, Y, \hat{X}) :=$

$$\begin{bmatrix} -\log \frac{P_{U|Y}(U|Y)}{P_U(U)} \\ \log \frac{P_{U|X}(U|X)}{P_U(U)} \\ \mathbb{E}_{\tilde{X}\tilde{Y}} \left[ \ell(\tilde{X}, \hat{f}^{(n)}(\mathbf{Z}, \tilde{Y})) \right] \\ d(X, \hat{X}) \end{bmatrix},$$

Let  $\mathbf{J} = \mathbb{E}_{U, X, Y, \hat{X}} [\mathbf{i}]$  and  $\mathbf{V} = \mathbb{C} [\mathbf{i}(U, X, Y, \hat{X})]$ .

The **dispersion region** is defined as  $\mathcal{S}(\mathbf{V}, \varepsilon) := \{\mathbf{b} \in \mathbb{R}^k : \Pr(\mathbf{B} \leq \mathbf{b}) \geq 1 - \varepsilon\}$  with  $\mathbf{B} \sim \mathcal{N}(0, \mathbf{V})$ .

## Theorem

For every  $0 < \varepsilon < 1$ , and  $n$  sufficiently large, the  $(n, \varepsilon)$ -rate-generalization error function satisfies:

$$R_b(n, \varepsilon, G, D) \leq \inf \left\{ \mathbf{M} \left( \mathbf{J} + \frac{\mathcal{S}(\mathbf{V}, \varepsilon)}{\sqrt{n}} + \frac{2 \log n}{n} \mathbf{1}_4 \right) \right\} \quad (25)$$

with  $\mathbf{M} = [1 \ 1 \ 0 \ 0]$ .

## Sketch of proof

Consider the following sets similar to that of [WKT15]<sup>9</sup>

$$\mathcal{T}_p(\gamma_p) := \left\{ (u, y) : \log \frac{P_{Y|U}(y|u)}{P_Y(y)} \geq \gamma_p \right\}, \quad (26)$$

$$\mathcal{T}_c(\gamma_c) := \left\{ (u, x) : \log \frac{P_{X|U}(x|u)}{P_X(x)} \leq \gamma_c \right\}, \quad (27)$$

$$\mathcal{T}_d(D) := \{(x, \hat{x}) : d(x, \hat{x}) \leq D\}, \quad (28)$$

$$\mathcal{T}_g(G) := \left\{ (\mathbf{u}, \mathbf{y}) : \mathbb{E}_{\tilde{X}\tilde{Y}} \left[ \ell(\tilde{X}, \hat{f}^{(n)}(\mathbf{z}, \tilde{Y})) \right] \leq G \right\}. \quad (29)$$

<sup>9</sup>S. Watanabe, S. Kuzuoka, and V. Y. Tan, "Nonasymptotic and second-order achievability bounds for coding with side-information," IEEE Transactions on Information Theory, vol. 61, no. 4, pp. 1574–1605, 2015.

## Sketch of proof

Consider the following sets similar to that of [WKT15]<sup>9</sup>

$$\mathcal{T}_p(\gamma_p) := \left\{ (u, y) : \log \frac{P_{Y|U}(y|u)}{P_Y(y)} \geq \gamma_p \right\}, \quad (26)$$

$$\mathcal{T}_c(\gamma_c) := \left\{ (u, x) : \log \frac{P_{X|U}(x|u)}{P_X(x)} \leq \gamma_c \right\}, \quad (27)$$

$$\mathcal{T}_d(D) := \{(x, \hat{x}) : d(x, \hat{x}) \leq D\}, \quad (28)$$

$$\mathcal{T}_g(G) := \left\{ (u, y) : \mathbb{E}_{\tilde{X}\tilde{Y}} \left[ \ell(\tilde{X}, \hat{f}^{(n)}(\mathbf{z}, \tilde{Y})) \right] \leq G \right\}. \quad (29)$$

### Theorem

For arbitrary constants  $\gamma_p, \gamma_c, G, D \geq 0$ , and positive integer  $N$ , there exists an  $(n, M, G, D, \varepsilon)$  code satisfying

$$\begin{aligned} \varepsilon \leq & P_{UXY\hat{X}} \left[ (\mathbf{u}, \mathbf{y}) \in \mathcal{T}_p(\gamma_p)^c \cup (\mathbf{u}, \mathbf{x}) \in \mathcal{T}_c(\gamma_c)^c \cup (\mathbf{u}, \mathbf{y}) \in \mathcal{T}_g(G)^c \cup (\mathbf{x}, \hat{\mathbf{x}}) \in \mathcal{T}_d(D)^c \right] \\ & + \frac{N}{2^{\gamma_p |\mathcal{M}|}} + \frac{1}{2} \sqrt{\frac{2\gamma_c}{N}}. \end{aligned} \quad (30)$$

<sup>9</sup>S. Watanabe, S. Kuzuoka, and V. Y. Tan, "Nonasymptotic and second-order achievability bounds for coding with side-information," IEEE Transactions on Information Theory, vol. 61, no. 4, pp. 1574–1605, 2015.

# Numerical results

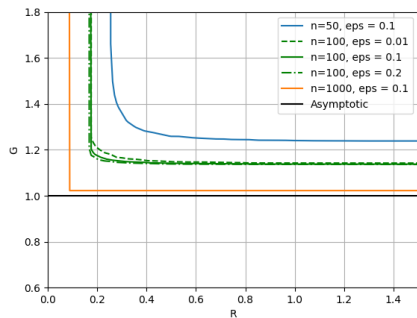


Figure 2: Rate-generalization error region for polynomial regression labeled on the block-length  $n$  and the excess probability  $\varepsilon$ .

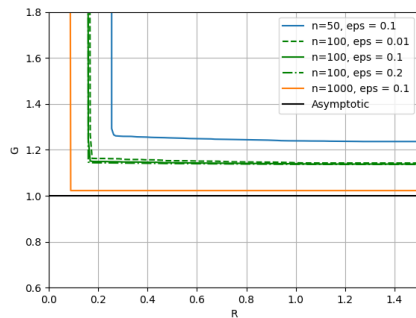


Figure 3: Rate-generalization error region for kernel regression labeled on the block-length  $n$  and the excess loss probability  $\varepsilon$ .

# Non-asymptotic trade-off

## Corollary

For  $0 < \varepsilon < 1$  and  $n$  sufficiently large, there exists an achievable rate-distortion-generalization error region such that

$$R_b(n, G, D, \varepsilon) > \max\{R_b(n, G, \varepsilon), R_b(n, D, \varepsilon)\}. \quad (31)$$

And there is **no trade-off** between generalization error of regression and reconstruction.

By Gaussian approximation, the dispersion matrix  $\mathcal{S}(\underline{\mathbf{V}}, \varepsilon)$  is determined by the correlation matrix  $\underline{\mathbf{V}}$ . We show that

$$\text{Cov} \left( \mathbb{E}_{\tilde{\mathbf{X}}\tilde{\mathbf{Y}}} \left[ \ell(\tilde{\mathbf{X}}, \hat{f}^{(n)}(\mathbf{Z}, \tilde{\mathbf{Y}})) \right], d(\mathbf{X}, \hat{\mathbf{X}} | \mathbf{Z} = \mathbf{z}) \right) = 0 \quad (32)$$

# Distortion-generalization error region

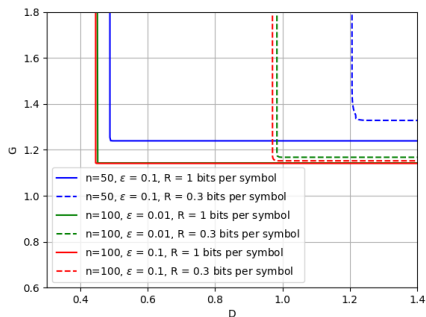


Figure 4: Distortion-generalization error region for polynomial regression on the block-length  $n$ , the excess loss probability  $\varepsilon$  and rate  $R$ .

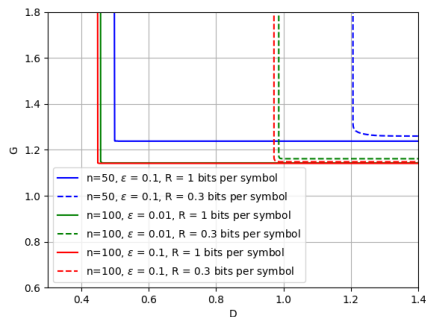


Figure 5: Distortion-generalization error region for kernel regression on the block-length  $n$ , the excess loss probability  $\varepsilon$  and rate  $R$ .

# Outline

- 1 Problem statement
- 2 Asymptotic rate-generalization error regions
- 3 Non-asymptotic rate-distortion-generalization error regions
- 4 Practical coding scheme**
- 5 Conclusion and perspectives

# Practical coding scheme for parametric regression

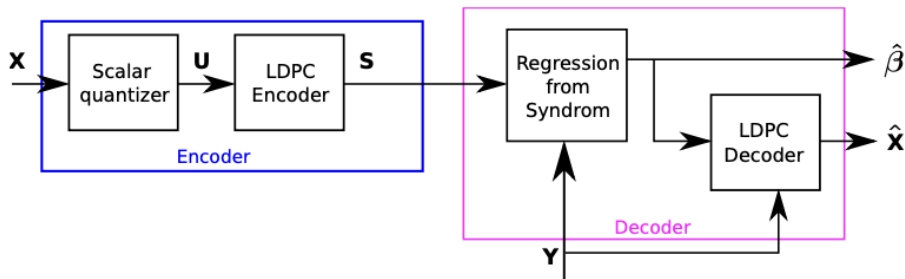


Figure 6: Practical coding scheme using LDPC codes

- **Encoder**

- Scalar quantizer over  $2^q$  levels (Lloyd-max or uniform)
- LDPC encoder :  $\mathbf{s} = \mathbf{H}\mathbf{u}$  with LDPC codes in  $GF(q)$

- **Decoder**

- Maximum Likelihood parameter estimation over the syndrom  $\mathbf{s}$  :  $\hat{\beta} = \arg \max_{\beta \in \mathbb{R}^p} \mathbb{P}(\mathbf{s}|\mathbf{y})$
- LDPC decoder from estimated  $\hat{\beta}$



## Numerical results for polynomial regression

- $Y \sim \mathcal{U}[-1, 1]$ ,  $X = \beta_0 + \beta_1 Y + \beta_2 Y^2 + N$
- $n = 100$ , regular (3, 6)-LDPC code in  $\text{GF}(4)$ , with rate  $r = \frac{1}{2}$
- $R = 1$  bit/symbol
- Lower bound :  $\sigma^2$ ; Upper bound :  $\sigma^2 + \frac{kC(\sigma^2 + \sigma_\Phi^2)}{n}$

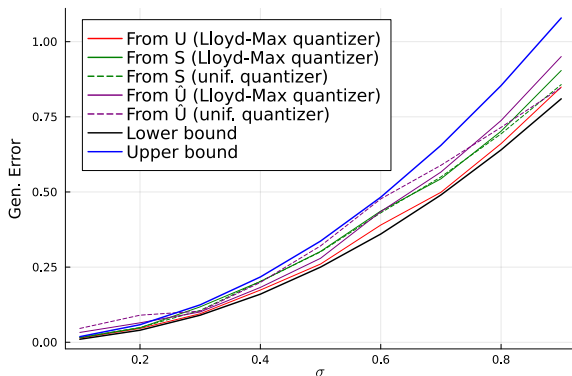


Figure 7: Generalization error with respect to  $\sigma$  for polynomial regression of order 3 with rate  $R = 1$  bit/symbol

## Numerical results for logistic regression

- $Y \sim \mathcal{U}[-1, 1]$ ,  $X = \beta_0 + \frac{\beta_1}{1+e^{-2Y}} + \frac{\beta_2}{1+e^{-4Y}} + N$
- $n = 100$ , regular (3, 6)-LDPC code in  $\text{GF}(4)$ , with rate  $r = \frac{1}{2}$
- $R = 1$  bit/symbol

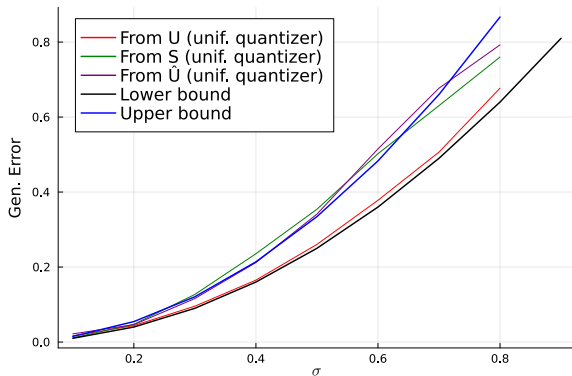


Figure 8: Generalization error with respect to  $\sigma$  for logistic regression of order 3 with rate  $R = 1$  bit/symbol

# Outline

- 1 Problem statement
- 2 Asymptotic rate-generalization error regions
- 3 Non-asymptotic rate-distortion-generalization error regions
- 4 Practical coding scheme
- 5 Conclusion and perspectives

# Conclusion

- Conclusions :
  - In asymptotic regime:  $R(D, G) = R(D)$ ;
  - In non-asymptotic regime: an achievable region with excess probability is provided;
  - **No trade-off** between regression and reconstruction;
  - Learning over compressed data without any prior decompression is possible;
  - **Same coding method can be used for regression.**
  
- Ongoing work :
  - Semantic communication with decoder's side information;
  - Classification as an example.

## Current work : Classification

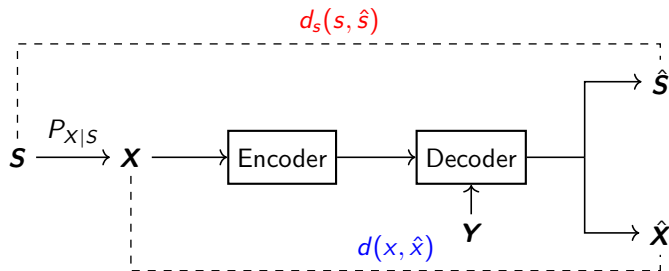


Figure 9: Coding scheme for semantic WZ coding

We consider the following problem

$$R(D, D_s) = \inf_{p(u|x):} I(X; U|Y). \quad (33)$$

$$\mathbb{E} \left[ d(X, \hat{X}) \right] \leq D$$

$$\mathbb{E} \left[ d'(X, \hat{S}) \right] \leq D_s$$

with  $d'(x, \hat{s}) = \frac{1}{p(x)} \sum_{s \in \mathcal{S}} p(x, s) d_s(s, \hat{s})$ ,  $d_s(\mathbf{s}^n, \hat{\mathbf{s}}^n) = \frac{1}{n} \sum_{i=1}^n d_s(s_i, \hat{s}_i)$  and  $d(\mathbf{x}^n, \hat{\mathbf{x}}^n) = \frac{1}{n} \sum_{i=1}^n d(x_i, \hat{x}_i)$ .

# Thank You!

J. Wei, E. Dupraz, and P. Mary, “Distributed source coding for parametric and non-parametric regression,” arXiv preprint arXiv:2404.18688, 2024