

---

# Model Selection with the Loss Rank Principle

---

**Marcus Hutter**

RSISE @ ANU and SML @ NICTA, Canberra, ACT, 0200, Australia  
marcus@hutter1.net      www.hutter1.net

**Minh-Ngoc Tran**

Department of Statistics and Applied Probability,  
National University of Singapore,    ngoctm@nus.edu.sg

2 March 2010

## Abstract

A key issue in statistics and machine learning is to automatically select the “right” model complexity, e.g., the number of neighbors to be averaged over in  $k$  nearest neighbor (kNN) regression or the polynomial degree in regression with polynomials. We suggest a novel principle - the Loss Rank Principle (LoRP) - for model selection in regression and classification. It is based on the loss rank, which counts how many other (fictitious) data would be fitted better. LoRP selects the model that has minimal loss rank. Unlike most penalized maximum likelihood variants (AIC, BIC, MDL), LoRP depends only on the regression functions and the loss function. It works without a stochastic noise model, and is directly applicable to any non-parametric regressor, like kNN.

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>The Loss Rank Principle</b>	<b>4</b>
<b>3</b>	<b>LoRP for <math>y</math>-Linear Models</b>	<b>8</b>
<b>4</b>	<b>Optimality Properties of LoRP for Variable Selection</b>	<b>11</b>
<b>5</b>	<b>Experiments</b>	<b>13</b>
<b>6</b>	<b>Comparison to Gaussian Bayesian Linear Regression</b>	<b>17</b>
<b>7</b>	<b>Comparison to other Model Selection Schemes</b>	<b>19</b>
<b>8</b>	<b>Loss Functions and their Selection</b>	<b>23</b>
<b>9</b>	<b>Self-Consistent Regression</b>	<b>24</b>
<b>10</b>	<b>Nearest Neighbors Classification</b>	<b>26</b>
<b>11</b>	<b>Conclusion and Outlook</b>	<b>28</b>
	<b>References</b>	<b>30</b>

## Keywords

Model selection, loss rank principle, non-parametric regression, classification, general loss function,  $k$  nearest neighbors.

# 1 Introduction

**Regression.** Consider a regression or classification problem in which we want to determine the functional relationship  $y_i \approx f_{\text{true}}(x_i)$  from data  $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$ , i.e., we seek a function  $r(\cdot|D) \equiv r(D)(\cdot)$  such that  $r(x|D) \equiv r(D)(x)$  is close to the unknown  $f_{\text{true}}(x)$  for all  $x$ . One may define  $r(\cdot|D)$  directly, e.g., “average the  $y$  values of the  $k$  nearest neighbors (kNN) of  $x$  in  $D$ ”, or select  $r(\cdot|D)$  from a class of functions  $\mathcal{F}$  that has smallest (training) error on  $D$ . If the class  $\mathcal{F}$  is not too large, e.g., the polynomials of fixed reasonable degree  $d$ , this often works well.

**Model selection.** What remains is to select the right model complexity  $c$ , like  $k$  or  $d$ . This selection cannot be based on the training error, since the more complex the model (large  $d$ , small  $k$ ) the better the fit on  $D$  (perfect for  $d = n$  and  $k = 1$ ). This problem is called overfitting, for which various remedies have been suggested.

The most popular ones in practice are based on a test set  $D'$  used for selecting the  $c$  for which the function  $r_c(\cdot|D)$  has smallest (test) error on  $D'$ , or improved versions like cross-validation [All74]. Typically  $D'$  is cut from  $D$ , thus reducing the sample size available for regression. Test set methods often work well in practice, but the reduced sample decreases accuracy, which can be a serious problem if  $n$  is small. We will not discuss empirical test set methods any further. See [Mac92] for a comparison of cross-validation with Bayesian model selection.

There are also various model selection methods that allow to use all data  $D$  for regression. The most popular ones can be regarded as penalized versions of Maximum Likelihood (ML). In addition to the function class  $\mathcal{F}_c$  (subscript  $c$  belonging to some set indexing the complexity), one has to specify a sampling model  $P(D|f)$ , e.g., that the  $y_i$  have independent Gaussian distribution with mean  $f(x_i)$ . ML chooses  $r_c(D) = \operatorname{argmax}_{f \in \mathcal{F}_c} P(D|f)$ , Penalized ML (PML) then chooses  $\hat{c} = \operatorname{argmin}_c \{-\log P(D|r_c(D)) + \text{Penalty}(c)\}$ , where the penalty depends on the used approach (MDL [Ris78], BIC [Sch78], AIC [Aka73]). All PML variants rely on a proper sampling model (which may be difficult to establish), ignore (or at least do not tell how to incorporate) a potentially given loss function (see [Yam99, Grü04] for exceptions), are based on distribution-independent penalties (which may result in bad performance for specific distributions), and are typically limited to (semi)parametric models.

**Main idea.** The main goal of the paper is to establish a criterion for selecting the “best” model complexity  $c$  based on regressors  $r_c$  given as a black box without insight into the origin or inner structure of  $r_c$ , that does not depend on things often not given (like a stochastic noise model), and that exploits what is/should be given (like the loss function, note that the criterion can also be used for loss-function selection, see Section 8). The key observation we exploit is that large classes  $\mathcal{F}_c$  or more flexible regressors  $r_c$  can fit more data well than more rigid ones. We define the *loss rank* of  $r_c$  as the number of other (fictitious) data  $D'$  that are fitted better by  $r_c(D')$  than  $D$  is fitted by  $r_c(D)$ , as measured by some loss function. The loss rank is large for regressors fitting  $D$  not well *and* for too flexible regressors (in both

cases the regressor fits many other  $D'$  better). The loss rank has a minimum for not too flexible regressors which fit  $D$  not too bad. We claim that minimizing the loss rank is a suitable model selection criterion, since it trades off the quality of fit with the flexibility of the model. Unlike PML, our Loss Rank Principle (LoRP) works without a noise (stochastic sampling) model, and is directly applicable to any non-parametric regressor, like kNN.

**Related ideas.** There are various other ideas that somehow count fictitious data. In normalized ML [Grü04], the complexity of a stochastic model class is defined as the log sum over all  $D'$  of maximum likelihood probabilities. In the luckiness framework for classification [Her02, Chp.4], the loss rank is related to the level of a hypothesis, if the empirical loss is used as an unluckiness function. The empirical Rademacher complexity [Kol01, BBL02] averages over all possible relabeled instances. Finally, instead of considering all  $D'$  one could consider only the set of all permutations of  $\{y_1, \dots, y_n\}$ , like in permutation tests [ET93]. The test statistic would here be the empirical loss.

**Contents.** In Section 2, after giving a brief introduction to regression, we formally state LoRP for model selection. Explicit expressions for the loss rank for the important class of linear regressors are derived in Section 3; this class includes kNN, polynomial, linear basis function (LBFR), kernel, projective regression, and some others. In Section 4, we establish optimality properties of LoRP for linear regression, namely model consistency and asymptotic mean efficiency. Experiments are presented in Section 5: We compare LoRP to other selection methods and demonstrate the use of LoRP for some specific problems like choosing tuning parameters in kNN and spline regression. In Section 6 we compare linear LoRP to Bayesian model selection for linear regression with Gaussian noise and prior, and in Section 7 to PML, in particular MDL, BIC, and AIC, and then discuss two trace formulas for the effective dimension. Sections 8-10 can be considered as extension sections. In Section 8 we show how to generalize linear LoRP to non-quadratic loss, in particular to other norms. We also discuss how LoRP can be used to select the loss function itself, in case it is not part of the problem specification. In Section 9 we briefly discuss interpolation. LoRP only depends on the regressor on data  $D$  and not on  $x \notin \{x_1, \dots, x_n\}$ . We construct canonical regressors for off-data interpolation from regressors given only on-data, in particular for kNN, Kernel, and LBFR, and show that they are canonical. In Section 10 we derive exact expressions for kNN when  $\{x_1, \dots, x_n\}$  forms a discrete  $d$ -dimensional hypercube, and discuss the limits  $n \rightarrow \infty$ ,  $k \rightarrow \infty$ , and  $d \rightarrow \infty$ . Section 11 contains the conclusions of our work and further considerations that could be elaborated on in the future.

The main idea of LoRP has already been presented at the COLT 2007 conference [Hut07]. In this paper we present LoRP more thoroughly, discover its theoretical properties and evaluate the method through some experiments.

## 2 The Loss Rank Principle

After giving a brief introduction to regression, classification, model selection, over-fitting, and some reoccurring examples, we state our novel Loss Rank Principle for model selection. We first state it for classification (Principle 3 for discrete values), and then generalize it for regression (Principle 5 for continuous values), and exemplify it on two (over-simplistic) artificial Examples 4 and 6. Thereafter we show how to regularize LoRP for realistic regression problems.

**Setup and notation.** We assume data  $D = (\mathbf{x}, \mathbf{y}) := \{(x_1, y_1), \dots, (x_n, y_n)\} \in (\mathcal{X} \times \mathcal{Y})^n =: \mathcal{D}$  has been observed. We think of the  $y$  as having an approximate functional dependence on  $x$ , i.e.,  $y_i \approx f_{\text{true}}(x_i)$ , where  $\approx$  means that the  $y_i$  are distorted by noise from the unknown “true” values  $f_{\text{true}}(x_i)$ . We will write  $(x, y)$  for generic data points, use vector notation  $\mathbf{x} = (x_1, \dots, x_n)^\top$  and  $\mathbf{y} = (y_1, \dots, y_n)^\top$ , and  $D' = (\mathbf{x}', \mathbf{y}')$  for generic (fictitious) data of size  $n$ . A full list of abbreviations and notations used throughout the paper is placed in the appendix.

**Regression and classification.** In regression problems  $\mathcal{Y}$  is typically (a subset of) the real set  $\mathbb{R}$  or some more general measurable space like  $\mathbb{R}^m$ . In classification,  $\mathcal{Y}$  is a finite set or at least discrete. We impose no restrictions on  $\mathcal{X}$ . Indeed,  $\mathbf{x}$  will essentially be fixed and plays only a spectator role, so we will often notationally suppress dependencies on  $\mathbf{x}$ . The goal of regression/classification is to find a function  $f_D \in \mathcal{F} \subset \mathcal{X} \rightarrow \mathcal{Y}$  “close” to  $f_{\text{true}}$  based on the past observations  $D$ . Or phrased in another way: we are interested in a mapping  $r : \mathcal{D} \rightarrow \mathcal{F}$  such that  $\hat{y} := r(x|D) \equiv r(D)(x) \equiv f_D(x) \approx f_{\text{true}}(x)$  for all  $x \in \mathcal{X}$ .

**Example 1 (polynomial regression)** For  $\mathcal{X} = \mathcal{Y} = \mathbb{R}$ , consider the set  $\mathcal{F}_d := \{f_{\mathbf{w}}(x) = w_d x^{d-1} + \dots + w_2 x + w_1 : \mathbf{w} \in \mathbb{R}^d\}$  of polynomials of degree  $d-1$ . Fitting the polynomial to data  $D$ , e.g., by least squares regression, we estimate  $\mathbf{w}$  with  $\hat{\mathbf{w}}_D$ . The regression function  $\hat{y} = r_d(x|D) = f_{\hat{\mathbf{w}}_D}(x)$  can be written down in closed form (see Example 9).  $\diamond$

**Example 2 (k nearest neighbors)** Let  $\mathcal{Y}$  be some vector space like  $\mathbb{R}$  and  $\mathcal{X}$  be a metric space like  $\mathbb{R}^m$  with some (e.g., Euclidean) metric  $d(\cdot, \cdot)$ . kNN estimates  $f_{\text{true}}(x)$  by averaging the  $y$  values of the  $k$  nearest neighbors  $\mathcal{N}_k(x)$  of  $x$  in  $D$ , i.e.,  $r_k(x|D) = \frac{1}{k} \sum_{i \in \mathcal{N}_k(x)} y_i$  with  $|\mathcal{N}_k(x)| = k$  such that  $d(x, x_i) \leq d(x, x_j)$  for all  $i \in \mathcal{N}_k(x)$  and  $j \notin \mathcal{N}_k(x)$ .  $\diamond$

**Parametric versus non-parametric regression.** Polynomial regression is an example of parametric regression in the sense that  $r_d(D)$  is the optimal function from a family of functions  $\mathcal{F}_d$  indexed by  $d < \infty$  real parameters ( $\mathbf{w}$ ). In contrast, the kNN regressor  $r_k$  is directly given and is not based on a finite-dimensional family of functions. In general,  $r$  may be given either directly or be the result of an optimization process.

**Loss function.** The quality of fit to the data is usually measured by a loss function  $\text{Loss}(\mathbf{y}, \hat{\mathbf{y}})$ , where  $\hat{y}_i = \hat{f}_D(x_i)$  is an estimate of  $y_i$ . Often the loss is additive:

$\text{Loss}(\mathbf{y}, \hat{\mathbf{y}}) = \sum_{i=1}^n \text{Loss}(y_i, \hat{y}_i)$ . If the class  $\mathcal{F}$  is not too large, good regression functions  $r(D)$  can be found by minimizing the loss w.r.t. all  $f \in \mathcal{F}$ . For instance,  $r_d(D) = \text{argmin}_{f \in \mathcal{F}_d} \sum_{i=1}^n (y_i - f(x_i))^2$  and  $\hat{y} = r_d(x|D)$  in Example 1.

**Regression class and loss.** In the following we assume a class of regressors  $\mathcal{R}$  (whatever their origin), e.g., the kNN regressors  $\{r_k : k \in \mathbb{N}\}$  or the least squares polynomial regressors  $\{r_d : d \in \mathbb{N}_0 := \mathbb{N} \cup \{0\}\}$ . Each regressor  $r$  can be thought of as a model. Throughout the paper, we use the terms “regressor” and “model” interchangeably. Note that unlike  $f \in \mathcal{F}$ , regressors  $r \in \mathcal{R}$  are not functions of  $x$  alone but depend on all observations  $D$ , in particular on  $\mathbf{y}$ . Like for functions  $f$ , we can compute the empirical loss of each regressor  $r \in \mathcal{R}$ :

$$\text{Loss}_r(D) \equiv \text{Loss}_r(\mathbf{y}|\mathbf{x}) := \text{Loss}(\mathbf{y}, \hat{\mathbf{y}}) = \sum_{i=1}^n \text{Loss}(y_i, r(x_i|\mathbf{x}, \mathbf{y}))$$

where  $\hat{y}_i = r(x_i|D)$  in the third expression, and the last expression holds in case of additive loss.

**Overfitting.** Unfortunately, minimizing  $\text{Loss}_r$  w.r.t.  $r$  will typically *not* select the “best” overall regressor. This is the well-known overfitting problem. In case of polynomials, the classes  $\mathcal{F}_d \subset \mathcal{F}_{d+1}$  are nested, hence  $\text{Loss}_{r_d}$  is monotone decreasing in  $d$  with  $\text{Loss}_{r_n} \equiv 0$  perfectly fitting the data. In case of kNN,  $\text{Loss}_{r_k}$  is more or less an increasing function in  $k$  with perfect regression on  $D$  for  $k=1$ , since no averaging takes place. In general,  $\mathcal{R}$  is often indexed by a “flexibility” or smoothness or complexity parameter, which has to be properly determined. The more flexible  $r$  is, the closer it can fit the data. Hence such  $r$  has smaller empirical loss, but is not necessarily better since it has higher variance. Clearly, too inflexible  $r$  also lead to a bad fit (“high bias”).

**Main goal.** The main goal of the paper is to establish a selection criterion in order to specify the smallest model to which  $f_{\text{true}}$  belongs or is close to, and simultaneously determine the “best” fitting function  $r(D)$ . The criterion

- is based on  $r$  given as a black box that does not require insight into the origin or inner structure of  $r$ ;
- does not depend on things often not given (like a stochastic noise model); and
- exploits what is or should be given (like the loss function).

**Definition of loss rank.** We first consider discrete  $\mathcal{Y}$  (i.e., classification), fix  $\mathbf{x}$ ,  $\mathbf{y}$  is the observed data and  $\mathbf{y}'$  are fictitious others. The key observation we exploit is that a more flexible  $r$  can fit more data  $D' \in \mathcal{D}$  well than a more rigid one. The more flexible  $r$  is, the smaller the empirical loss  $\text{Loss}_r(\mathbf{y}|\mathbf{x})$  is. Instead of minimizing the unsuitable  $\text{Loss}_r(\mathbf{y}|\mathbf{x})$  w.r.t.  $r$ , we could ask how many  $\mathbf{y}' \in \mathcal{Y}^n$  lead to smaller  $\text{Loss}_r$  than  $\mathbf{y}$ . We define the loss rank of  $r$  (w.r.t.  $\mathbf{y}$ ) as the number of  $\mathbf{y}' \in \mathcal{Y}^n$  with smaller or equal empirical loss than  $\mathbf{y}$ :

$$\text{Rank}_r(\mathbf{y}|\mathbf{x}) \equiv \text{Rank}_r(L) := \#\{\mathbf{y}' \in \mathcal{Y}^n : \text{Loss}_r(\mathbf{y}'|\mathbf{x}) \leq L\} \quad \text{with } L := \text{Loss}_r(\mathbf{y}|\mathbf{x}) \quad (1)$$

We claim that the loss rank of  $r$  is a suitable model selection measure. For (1) to make sense, we have to assume (and will later assure) that  $\text{Rank}_r(L) < \infty$ , i.e., there are only finitely many  $\mathbf{y}' \in \mathcal{Y}^n$  having loss smaller than  $L$ .

Since the logarithm is a strictly monotone increasing function, we can also consider the logarithmic rank  $\text{LR}_r(\mathbf{y}|\mathbf{x}) := \log \text{Rank}_r(\mathbf{y}|\mathbf{x})$ , which will be more convenient.

**Principle 3 (LoRP for classification)** *For discrete  $\mathcal{Y}$ , the best classifier/regressor  $r : \mathcal{D} \times \mathcal{X} \rightarrow \mathcal{Y}$  in some class  $\mathcal{R}$  for data  $D = (\mathbf{x}, \mathbf{y})$  is the one with the smallest loss rank:*

$$r^{best} = \arg \min_{r \in \mathcal{R}} \text{LR}_r(\mathbf{y}|\mathbf{x}) \equiv \arg \min_{r \in \mathcal{R}} \text{Rank}_r(\mathbf{y}|\mathbf{x}) \quad (2)$$

where  $\text{Rank}_r$  is defined in (1).

We give a simple example for which we can compute all ranks by hand to help the reader better grasp how the principle works.

**Example 4 (simple discrete)** Consider  $\mathcal{X} = \{1, 2\}$ ,  $\mathcal{Y} = \{0, 1, 2\}$ , and two points  $D = \{(1, 1), (2, 2)\}$  lying on the diagonal  $x = y$ , with polynomial (zero, constant, linear) least squares regressors  $\mathcal{R} = \{r_0, r_1, r_2\}$  (see Ex.1).  $r_0$  is simply 0,  $r_1$  the  $y$ -average, and  $r_2$  the line through points  $(1, y_1)$  and  $(2, y_2)$ . This, together with the quadratic Loss for generic  $\mathbf{y}'$  and observed  $\mathbf{y} = (1, 2)$  and fixed  $\mathbf{x} = (1, 2)$ , is summarized in the following table

$d$	$r_d(x \mathbf{x}, \mathbf{y}')$	$\text{Loss}_d(\mathbf{y}' \mathbf{x})$	$\text{Loss}_d(D)$
0	0	$y_1'^2 + y_2'^2$	5
1	$\frac{1}{2}(y_1' + y_2')$	$\frac{1}{2}(y_2' - y_1')^2$	$\frac{1}{2}$
2	$(y_2' - y_1')(x - 1) + y_1'$	0	0

From the Loss we can easily compute the Rank for all nine  $\mathbf{y}' \in \{0, 1, 2\}^2$ . Equal rank due to equal loss is indicated by a “=” in the table below. Whole equality groups are actually assigned the rank of their right-most member, e.g., for  $d=1$  the ranks of  $(y_1', y_2') = (0, 1), (1, 0), (2, 1), (1, 2)$  are all 7 (and not 4, 5, 6, 7).

$d$	$\text{Rank}_{r_d}(y_1' y_2'   12)$	$\text{Rank}_{r_d}(D)$
0	$y_1' y_2' = 00 < 01 = 10 < 11 < 02 = 20 < 21 = \mathbf{12} < 22$	8
1	$y_1' y_2' = 00 = 11 = 22 < 01 = 10 = 21 = \mathbf{12} < 02 = 20$	7
2	$y_1' y_2' = 00 = 01 = 02 = 10 = 11 = 20 = 21 = 22 = \mathbf{12}$	9

So LoRP selects  $r_1$  as best regressor, since it has minimal rank on  $D$ .  $r_0$  fits  $D$  too badly and  $r_2$  is too flexible (perfectly fits all  $D'$ ).  $\diamond$

**LoRP for continuous  $\mathcal{Y}$ .** We now consider the case of continuous or measurable spaces  $\mathcal{Y}$ , i.e., normal regression problems. We assume  $\mathcal{Y} = \mathbb{R}$  in the following exposition, but the idea and resulting principle hold for more general measurable

spaces like  $\mathbb{R}^m$ . We simply reduce the model selection problem to the discrete case by considering the discretized space  $\mathcal{Y}_\varepsilon = \varepsilon\mathbb{Z}$  for small  $\varepsilon > 0$  and discretize  $\mathbf{y} \rightsquigarrow \mathbf{y}_\varepsilon \in \varepsilon\mathbb{Z}^n$  (“ $\rightsquigarrow$ ” means “is replaced by”). Then  $\text{Rank}_r^\varepsilon(L) := \#\{\mathbf{y}'_\varepsilon \in \mathcal{Y}_\varepsilon^n : \text{Loss}_r(\mathbf{y}'_\varepsilon | \mathbf{x}) \leq L\}$  with  $L = \text{Loss}_r(\mathbf{y}_\varepsilon | \mathbf{x})$  counting the number of  $\varepsilon$ -grid points in the set

$$V_r(L) := \{\mathbf{y}' \in \mathcal{Y}^n : \text{Loss}_r(\mathbf{y}' | \mathbf{x}) \leq L\} \quad (3)$$

which we assume (and later assure) to be finite, analogous to the discrete case. Hence  $\text{Rank}_r^\varepsilon(L) \cdot \varepsilon^n$  is an approximation of the *loss volume*  $|V_r(L)|$  of set  $V_r(L)$ , and typically  $\text{Rank}_r^\varepsilon(L) \cdot \varepsilon^n = |V_r(L)| \cdot (1 + O(\varepsilon)) \rightarrow |V_r(L)|$  for  $\varepsilon \rightarrow 0$ . Taking the logarithm we get  $\text{LR}_r^\varepsilon(\mathbf{y} | \mathbf{x}) = \log \text{Rank}_r^\varepsilon(L) = \log |V_r(L)| - n \log \varepsilon + O(\varepsilon)$ . Since  $n \log \varepsilon$  is independent of  $r$ , we can drop it in comparisons like (2). So for  $\varepsilon \rightarrow 0$  we can define the log-loss “rank” simply as the log-volume

$$\text{LR}_r(\mathbf{y} | \mathbf{x}) := \log |V_r(L)|, \quad \text{where } L := \text{Loss}_r(\mathbf{y} | \mathbf{x}) \quad (4)$$

**Principle 5 (LoRP for regression)** *For measurable  $\mathcal{Y}$ , the best regressor  $r : \mathcal{D} \times \mathcal{X} \rightarrow \mathcal{Y}$  in some class  $\mathcal{R}$  for data  $D = (\mathbf{x}, \mathbf{y})$  is the one with the smallest loss volume:*

$$r^{\text{best}} = \arg \min_{r \in \mathcal{R}} \text{LR}_r(\mathbf{y} | \mathbf{x}) \equiv \arg \min_{r \in \mathcal{R}} |V_r(L)|$$

where  $\text{LR}$ ,  $V_r$ , and  $L$  are defined in (3) and (4), and  $|V_r(L)|$  is the volume of  $V_r(L) \subseteq \mathcal{Y}^n$ .

For discrete  $\mathcal{Y}$  with counting measure we recover the discrete LoRP (Principle 3).

**Example 6 (simple continuous)** Consider Example 4 but with interval  $\mathcal{Y} = [0, 2]$ . The first table remains unchanged, while the second table becomes

$d$	$V_d(L) = \{\mathbf{y}' \in [0, 2]^2 : \dots\}$	$ V_d(L) $	$\text{Loss}_d(D)$	$ V_d(\text{Loss}_d(D)) $
0	$y_1'^2 + y_2'^2 \leq L$	$\frac{\pi}{4}L$ if $L \leq 4$ ; $4$ if $L \geq 8$ ; $2\sqrt{L-4} + L(\frac{\pi}{4} - \cos^{-1}(\frac{2}{\sqrt{L}}))$ else	5	$\doteq 3.6$
1	$\frac{1}{2}(y_2' - y_1')^2 \leq L$	$4\sqrt{2L-2L}$ if $L \leq 2$ ; $4$ if $L > 2$	$\frac{1}{2}$	3
2	$0 \leq L$	4	0	4

So LoRP again selects  $r_1$  as best regressor, since it has smallest loss volume on  $D$ .  $\diamond$

**Infinite rank or volume.** Often the loss rank/volume will be infinite, e.g., if we had chosen  $\mathcal{Y} = \mathbb{Z}$  in Ex.4 or  $\mathcal{Y} = \mathbb{R}$  in Ex.6. There are various potential remedies. We could modify (a) the regressor  $r$  or (b) the Loss to make  $\text{LR}_r$  finite, (c) the Loss Rank Principle itself, or (d) find problem-specific solutions. Regressors  $r$  with infinite rank might be rejected for philosophical or pragmatic reasons. We will briefly consider (a) for linear regression later, but to fiddle around with  $r$  in a generic (blackbox way) seems difficult. We have no good idea how to tinker with LoRP (c), and also a patched LoRP may be less attractive. For kNN on a grid we

later use remedy (d). While in (decision) theory, the application’s goal determines the loss, in practice the loss is often more determined by convenience or rules of thumb. So the Loss (b) seems the most inviting place to tinker with. A very simple modification is to add a small penalty term to the loss.

$$\text{Loss}_r(\mathbf{y}|\mathbf{x}) \rightsquigarrow \text{Loss}_r^\alpha(\mathbf{y}|\mathbf{x}) := \text{Loss}_r(\mathbf{y}|\mathbf{x}) + \alpha\|\mathbf{y}\|^2, \quad \alpha > 0 \text{ “small”} \quad (5)$$

The Euclidean norm  $\|\mathbf{y}\|^2 := \sum_{i=1}^n y_i^2$  is default, but other (non)norm regularizations are possible. The regularized  $\text{LR}_r^\alpha(\mathbf{y}|\mathbf{x})$  based on  $\text{Loss}_r^\alpha$  is always finite, since  $\{\mathbf{y} : \|\mathbf{y}\|^2 \leq L\}$  has finite volume. An alternative penalty  $\alpha\hat{\mathbf{y}}^\top\hat{\mathbf{y}}$ , quadratic in the regression estimates  $\hat{y}_i = r(x_i|\mathbf{x}, \mathbf{y})$  is possible if  $r$  is unbounded in every  $\mathbf{y} \rightarrow \infty$  direction.

A scheme trying to determine a single (flexibility) parameter (like  $d$  and  $k$  in the above examples) would be of no use if it depended on one (or more) other unknown parameters ( $\alpha$ ), since varying through the unknown parameter leads to any (non)desired result. Since LoRP seeks the  $r$  of smallest rank, it is natural to also determine  $\alpha = \alpha_{\min}$  by minimizing  $\text{LR}_r^\alpha$  w.r.t.  $\alpha$ . The good news is that this leads to meaningful results. Interestingly, as we will see later, a clever choice of  $\alpha$  may also result in alternative optimalities of the selection procedure.

### 3 LoRP for $y$ -Linear Models

In this section we consider the important class of  $y$ -linear regressions with quadratic loss function. By “ $y$ -linear regression”, we mean the linearity is only assumed in  $y$  and the dependence on  $x$  can be arbitrary. This class is richer than it may appear. It includes the normal linear regression model, kNN (Example 7), kernel (Example 8), and many other regression models. For  $y$ -linear regression and  $\mathcal{Y} = \mathbb{R}$ , the loss rank is the volume of an  $n$ -dimensional ellipsoid, which can efficiently be computed in time  $O(n^3)$  (Theorem 10). For the special case of projective regression, e.g., linear basis function regression (Example 9), we can even determine the regularization parameter  $\alpha$  analytically (Theorem 11).

**$y$ -Linear regression.** We assume  $\mathcal{Y} = \mathbb{R}$  in this section; generalization to  $\mathbb{R}^m$  is straightforward. A  $y$ -linear regressor  $r$  can be written in the form

$$\hat{y} = r(x|\mathbf{x}, \mathbf{y}) = \sum_{j=1}^n m_j(x, \mathbf{x})y_j \quad \forall x \in \mathcal{X} \quad \text{and some} \quad m_j : \mathcal{X} \times \mathcal{X}^n \rightarrow \mathbb{R} \quad (6)$$

Particularly interesting is  $r$  for  $x = x_1, \dots, x_n$ .

$$\hat{y}_i = r(x_i|\mathbf{x}, \mathbf{y}) = \sum_j M_{ij}(\mathbf{x})y_j \quad \text{with} \quad M : \mathcal{X}^n \rightarrow \mathbb{R}^{n \times n} \quad (7)$$

where matrix  $M_{ij}(\mathbf{x}) = m_j(x_i, \mathbf{x})$ . Since LoRP needs  $r$  only on the training data  $\mathbf{x}$ , we only need  $M$ .

**Example 7 (kNN ctd.)** For kNN of Ex.2 we have  $m_j(x, \mathbf{x}) = \frac{1}{k}$  if  $j \in \mathcal{N}_k(x)$  and 0 else, and  $M_{ij}(\mathbf{x}) = \frac{1}{k}$  if  $j \in \mathcal{N}_k(x_i)$  and 0 else.  $\diamond$



**Example 8 (kernel regression)** Kernel regression takes a weighted average over  $\mathbf{y}$ , where the weight of  $y_j$  to  $y$  is proportional to the similarity of  $x_j$  to  $x$ , measured by a kernel  $K(x, x_j)$ , i.e.,  $m_j(x, \mathbf{x}) = K(x, x_j) / \sum_{j=1}^n K(x, x_j)$ . For example the Gaussian kernel for  $\mathcal{X} = \mathbb{R}^m$  is  $K(x, x_j) = e^{-\|x - x_j\|_2^2 / 2\sigma^2}$ . The width  $\sigma$  controls the smoothness of the kernel regressor, and LoRP selects the real-valued “complexity” parameter  $\sigma$ .  $\diamond$

**Example 9 (linear basis function regression, LBFR)** Let  $\phi_1(x), \dots, \phi_d(x)$  be a set or vector of “basis” functions often called “features”. We place no restrictions on  $\mathcal{X}$  or  $\phi: \mathcal{X} \rightarrow \mathbb{R}^d$ . Consider the class of functions linear in  $\phi$ :

$$\mathcal{F}_d = \{f_{\mathbf{w}}(x) = \sum_{a=1}^d w_a \phi_a(x) = \mathbf{w}^\top \phi(x) : \mathbf{w} \in \mathbb{R}^d\}$$

For instance, for  $\mathcal{X} = \mathbb{R}$  and  $\phi_a(x) = x^{a-1}$  we would recover the polynomial regression Example 1. For quadratic loss function  $\text{Loss}(y_i, \hat{y}_i) = (y_i - \hat{y}_i)^2$  we have

$$\text{Loss}_{\mathbf{w}}(\mathbf{y}|\phi) := \sum_{i=1}^n (y_i - f_{\mathbf{w}}(x_i))^2 = \mathbf{y}^\top \mathbf{y} - 2\mathbf{y}^\top \Phi \mathbf{w} + \mathbf{w}^\top B \mathbf{w}$$

where matrix  $\Phi$  is defined by  $\Phi_{ia} = \phi_a(x_i)$  and  $B$  is a symmetric matrix with  $B_{ab} = \sum_{i=1}^n \phi_a(x_i) \phi_b(x_i) = [\Phi^\top \Phi]_{ab}$ . The loss is quadratic in  $\mathbf{w}$  with minimum at  $\mathbf{w} = B^{-1} \Phi^\top \mathbf{y}$ . So the least squares regressor is  $\hat{y} = \mathbf{y}^\top \Phi B^{-1} \phi(x)$ , hence  $m_j(x, \mathbf{x}) = (\Phi B^{-1} \phi(x))_j$  and  $M(\mathbf{x}) = \Phi B^{-1} \Phi^\top$ .  $\diamond$

Consider now a general linear regressor  $M$  with quadratic loss and quadratic penalty

$$\begin{aligned} \text{Loss}_M^\alpha(\mathbf{y}|\mathbf{x}) &= \sum_{i=1}^n \left( y_i - \sum_{j=1}^n M_{ij} y_j \right)^2 + \alpha \|\mathbf{y}\|^2 = \mathbf{y}^\top S_\alpha \mathbf{y}, \\ \text{where}^1 \quad S_\alpha &= (I - M)^\top (I - M) + \alpha I \end{aligned} \quad (8)$$

( $I$  is the identity matrix).  $S_\alpha$  is a symmetric matrix. For  $\alpha > 0$  it is positive definite and for  $\alpha = 0$  positive semidefinite. If  $\lambda_1, \dots, \lambda_n \geq 0$  are the eigenvalues of  $S_0$ , then  $\lambda_i + \alpha$  are the eigenvalues of  $S_\alpha$ .  $V(L) = \{\mathbf{y}' \in \mathbb{R}^n : \mathbf{y}'^\top S_\alpha \mathbf{y}' \leq L\}$  is an ellipsoid with the eigenvectors of  $S_\alpha$  being the main axes and  $\sqrt{L/(\lambda_i + \alpha)}$  being their length. Hence the volume is

$$|V(L)| = v_n \prod_{i=1}^n \sqrt{\frac{L}{\lambda_i + \alpha}} = \frac{v_n L^{n/2}}{\sqrt{\det S_\alpha}}$$

where  $v_n = \pi^{n/2} / \frac{n}{2}!$  is the volume of the  $n$ -dimensional unit sphere,  $z! := \Gamma(z+1)$ , and  $\det$  is the determinant. Taking the logarithm we get

$$\text{LR}_M^\alpha(\mathbf{y}|\mathbf{x}) = \log |V(\text{Loss}_M^\alpha(\mathbf{y}|\mathbf{x}))| = \frac{n}{2} \log(\mathbf{y}^\top S_\alpha \mathbf{y}) - \frac{1}{2} \log \det S_\alpha + \log v_n \quad (9)$$

Since  $v_n$  is independent of  $\alpha$  and  $M$  it is possible to drop  $v_n$ . Consider now a class of linear regressors  $\mathcal{M} = \{M\}$ , e.g., the kNN regressors  $\{M_k : k \in \mathbb{N}\}$  or the  $d$ -dimensional linear basis function regressors  $\{M_d : d \in \mathbb{N}_0\}$ .

<sup>1</sup>The mentioned alternative penalty  $\alpha \|\hat{\mathbf{y}}\|^2$  would lead to  $S_\alpha = (I - M)^\top (I - M) + \alpha M^\top M$ . For LBFR, penalty  $\alpha \|\hat{\mathbf{w}}\|^2$  is popular (ridge regression). Apart from being limited to parametric regression, it has the disadvantage of not being reparametrization invariant. For instance, scaling  $\phi_a(x) \rightsquigarrow \gamma_a \phi_a(x)$  does not change the class  $\mathcal{F}_d$ , but changes the ridge regressor.

**Theorem 10 (LoRP for  $\mathbf{y}$ -linear regression)** For  $\mathcal{Y}=\mathbb{R}$ , the best linear regressor  $M:\mathcal{X}^n \rightarrow \mathbb{R}^{n \times n}$  in some class  $\mathcal{M}$  for data  $D=(\mathbf{x},\mathbf{y})$  is

$$M^{best} = \arg \min_{M \in \mathcal{M}, \alpha \geq 0} \left\{ \frac{n}{2} \log(\mathbf{y}^\top S_\alpha \mathbf{y}) - \frac{1}{2} \log \det S_\alpha \right\} = \arg \min_{M \in \mathcal{M}} \left\{ \frac{\mathbf{y}^\top S_\alpha \mathbf{y}}{(\det S_\alpha)^{1/n}} \right\} \quad (10)$$

where  $S_\alpha = S_\alpha(M)$  is defined in (8).

The last expression shows that linear LoRP minimizes the Loss times the geometric average of the squared axes lengths of ellipsoid  $V(1)$ . Note that  $M^{best}$  depends on  $\mathbf{y}$  unlike the  $M \in \mathcal{M}$ .

**Nullspace of  $S_0$ .** If  $M$  has an eigenvalue 1, then  $S_0 = (I - M)^\top(I - M)$  has a zero eigenvalue and  $\alpha > 0$  is necessary, since  $\det S_0 = 0$ . Actually this is true for most practical  $M$ . Most linear regressors are invariant under a constant shift of  $\mathbf{y}$ , i.e.,  $r(x|\mathbf{x},\mathbf{y}+c) = r(x|\mathbf{x},\mathbf{y}) + c$ , which implies that  $M$  has eigenvector  $(1, \dots, 1)^\top$  with eigenvalue 1. This can easily be checked for kNN (Ex.2), kernel (Ex.8), and LBFR (Ex.9). Such a generic 1-eigenvector effecting all  $M \in \mathcal{M}$  could easily and maybe should be filtered out by considering only the orthogonal space or dropping these  $\lambda_i = 0$  when computing  $\det S_0$ . The 1-eigenvectors that depend on  $M$  are the ones where we really need a regularizer  $\alpha > 0$ . For instance,  $M_d$  in LBFR has  $d$  eigenvalues 1, and  $M_{\text{kNN}}$  has as many eigenvalues 1 as there are disjoint components in the graph determined by the edges  $M_{ij} > 0$ . In general we need to find the optimal  $\alpha$  numerically. If  $M$  is a projection we can find  $\alpha_m$  analytically.

**Numerical approximation of  $(\det S_\alpha)^{1/n}$  and the computational complexity of linear LoRP.** For each  $\alpha$  and candidate model, the determinant of  $S_\alpha$  in the general case can be computed in time  $O(n^3)$ . Often  $M$  is a very sparse matrix (like in kNN) or can be well approximated by a sparse matrix (like for kernel regression), which allows us to approximate  $\det S_\alpha$  sometimes in linear time [Reu02]. To search the optimal  $\alpha$  and  $M$ , the computational cost depends on the range of  $\alpha$  we search and the number of candidate models we have.

**Projective regression.** Consider a projection matrix  $M = P = P^2$  with  $d (= \text{tr} P)$  eigenvalues 1, and  $n - d$  zero eigenvalues. For instance,  $M = \Phi B^{-1} \Phi^\top$  of LBFR Ex.9 is such a matrix. This implies that  $S_\alpha$  has  $d$  eigenvalues  $\alpha$  and  $n - d$  eigenvalues  $1 + \alpha$ , thus  $\det S_\alpha = \alpha^d (1 + \alpha)^{n-d}$ . Let  $\rho = \|\mathbf{y} - \hat{\mathbf{y}}\|^2 / \|\mathbf{y}\|^2$ , then  $\mathbf{y}^\top S_\alpha \mathbf{y} = (\rho + \alpha) \mathbf{y}^\top \mathbf{y}$  and

$$\text{LR}_P^\alpha = \frac{n}{2} \log \mathbf{y}^\top \mathbf{y} + \frac{n}{2} \log(\rho + \alpha) - \frac{d}{2} \log \alpha - \frac{n-d}{2} \log(1 + \alpha). \quad (11)$$

Solving  $\partial \text{LR}_P^\alpha / \partial \alpha = 0$  w.r.t.  $\alpha$  we get a minimum at  $\alpha = \alpha_m := \frac{\rho d}{(1-\rho)n-d}$  provided that  $1 - \rho > d/n$ . After some algebra we get

$$\text{LR}_P^{\alpha_m} = \frac{n}{2} \log \mathbf{y}^\top \mathbf{y} - \frac{n}{2} \text{KL}\left(\frac{d}{n} \parallel 1 - \rho\right), \quad \text{where} \quad \text{KL}(p \parallel q) := p \log \frac{p}{q} + (1-p) \log \frac{1-p}{1-q} \quad (12)$$

is the relative entropy or the Kullback-Leibler divergence. Note that (12) is still valid without the condition  $1 - \rho > d/n$  (the term  $\log((1-\rho)n-d)$  has been canceled

in the derivation). What we need when using (12) is that  $d < n$  and  $\rho < 1$ , which are very reasonable in practice. Interestingly, if we use the penalty  $\alpha \|\hat{\mathbf{y}}\|^2$  instead of  $\alpha \|\mathbf{y}\|^2$ , the loss rank then has the same expression as (12) without any condition<sup>2</sup>.

Minimizing  $\text{LR}_P^{\alpha_m}$  w.r.t.  $P$  is equivalent to maximizing  $\text{KL}(\frac{d}{n} \| 1 - \rho)$ . The term  $\rho$  is a measure of fit. If  $d$  increases, then  $\rho$  decreases and otherwise. We are seeking a tradeoff between the model complexity  $d$  and the measure of fit  $\rho$ , and LoRP suggests the optimal tradeoff by maximizing KL.

**Theorem 11 (LoRP for projective regression)** *The best projective regressor  $P: \mathcal{X}^n \rightarrow \mathbb{R}^{n \times n}$  with  $P = P^2$  in some projective class  $\mathcal{P}$  for data  $D = (\mathbf{x}, \mathbf{y})$  is*

$$P^{best} = \arg \max_{P \in \mathcal{P}} \text{KL}\left(\frac{\text{tr}P(\mathbf{x})}{n} \left\| \frac{\mathbf{y}^\top P(\mathbf{x})\mathbf{y}}{\mathbf{y}^\top \mathbf{y}}\right.\right). \quad (13)$$

## 4 Optimality Properties of LoRP for Variable Selection

In the previous sections, LoRP was stated for general-purpose model selection. By restricting attention to linear regression models, we will point out in this section some theoretical properties of LoRP for variable (also called feature or attribute) selection.

Variable selection is probably the most fundamental and important topic in linear regression analysis. At the initial stage of modeling, a large number of potential covariates are often introduced; one then has to select a smaller subset of the covariates to fit/interpret the data. There are two main goals of variable selection, one is model identification, the other is regression estimation. The former aims at identifying the true subset generating the data, while the latter aims at estimating efficiently the regression function, i.e., selecting a subset that has the minimum mean squared error loss. Note that whether or not there is a selection criterion achieving simultaneously these two goals is still an open question [Yan05, Grü04]. We show that with the optimal parameter  $\alpha$  (defined as  $\alpha_m$  that minimizes the loss rank  $\text{LR}_M^\alpha$  in  $\alpha$ ), LoRP satisfies the first goal, while with a suitable choice of  $\alpha$ , LoRP satisfies the second goal.

Given  $d+1$  potential covariates  $X_0 \equiv 1, X_1, \dots, X_d$  and a response variable  $Y$ , let  $X = \mathbf{x}$  be a non-random design matrix of size  $n \times (d+1)$  and  $\mathbf{y}$  be a response vector respectively (if  $\mathbf{y}$  and  $X$  are centered, then the covariate 1 can be omitted from the models). Denote by  $\mathcal{S} = \{0, j_1, \dots, j_{|\mathcal{S}|-1}\}$  the candidate model that has covariates  $X_0, X_{j_1}, \dots, X_{j_{|\mathcal{S}|-1}}$ . Under a proposed model  $\mathcal{S}$ , we can write

$$\mathbf{y} = X_{\mathcal{S}}\beta_{\mathcal{S}} + \sigma_{\mathcal{S}}\epsilon$$

---

<sup>2</sup>Then  $S_\alpha = (I_n - P)^\top (I_n - P) + \alpha P^\top P = I_n + (\alpha - 1)P$  has  $d$  eigenvalues  $\alpha$  and  $n - d$  eigenvalues 1, thus  $\det(S_\alpha) = \alpha^d$ . The loss rank  $\text{LR}_P^\alpha = \frac{n}{2} \log \mathbf{y}^\top \mathbf{y} + \frac{n}{2} \log(1 + (\alpha - 1)(1 - \rho)) - \frac{d}{2} \log \alpha$  is minimized at  $\alpha_m = \frac{\rho^d}{(1 - \rho)(n - d)}$ . After some algebra we get the same expression of  $\text{LR}_P^{\alpha_m}$  as (12).

where  $\epsilon$  is noise with expectation  $\mathbf{E}[\epsilon] = 0$  and covariance  $\text{Cov}(\epsilon) = I_n$ ,  $\sigma_{\mathcal{S}} > 0$ ,  $\beta_{\mathcal{S}} = (\beta_0, \beta_{j_1}, \dots, \beta_{j_{|\mathcal{S}|-1}})^\top$ , and  $X_{\mathcal{S}}$  is the  $n \times |\mathcal{S}|$  design matrix obtained from  $X$  by removing the  $(j+1)$ st column for all  $j \notin \mathcal{S}$ .

**Model consistency of LoRP for variable selection.** The ordinary least squares (OLS) fitted vector under model  $\mathcal{S}$  is  $\hat{\mathbf{y}}_{\mathcal{S}} = M_{\mathcal{S}} \mathbf{y}$  with  $M_{\mathcal{S}} = X_{\mathcal{S}}(X_{\mathcal{S}}^\top X_{\mathcal{S}})^{-1} X_{\mathcal{S}}^\top$  being a projection matrix. From Theorem 11 the best subset chosen by LoRP is

$$\hat{\mathcal{S}}_n = \arg \min_{\mathcal{S}} \text{LR}_{\mathcal{S}}^{\alpha_m} = \arg \max_{\mathcal{S}} \{ \text{KL}(\frac{|\mathcal{S}|}{n} \| 1 - \rho_{\mathcal{S}}) \}, \quad \rho_{\mathcal{S}} = \frac{\|\mathbf{y} - \hat{\mathbf{y}}_{\mathcal{S}}\|^2}{\|\mathbf{y}\|^2}.$$

The term  $\rho_{\mathcal{S}}$  is a measure of fit. It will be very close to 0 if model  $\mathcal{S}$  is big, otherwise, it will be close to 1 if  $\mathcal{S}$  is too small. Therefore, it is reasonable to consider only cases in which  $\rho_{\mathcal{S}}$  is bounded away from 0 and 1. In order to prove the theoretical properties of LoRP, we need the following technical assumption.

- (A) For each candidate model  $\mathcal{S}$ ,  $\rho_{\mathcal{S}}$  is bounded away from 0 and 1, i.e., there are constants  $c_1$  and  $c_2$  such that  $0 < c_1 \leq \rho_{\mathcal{S}} \leq c_2 < 1$  with probability 1 (w.p.1).

Let  $\hat{\sigma}_{\mathcal{S}}^2 = \|\mathbf{y} - \hat{\mathbf{y}}_{\mathcal{S}}\|^2/n$  and  $\mathcal{S}_{\text{null}} = \{0\}$ . It is easy to see that for every  $\mathcal{S}$

$$1 - \rho_{\mathcal{S}} = \|\hat{\mathbf{y}}_{\mathcal{S}}\|^2/\|\mathbf{y}\|^2, \quad n\hat{\sigma}_{\mathcal{S}}^2 = \rho_{\mathcal{S}}\|\mathbf{y}\|^2, \quad n\bar{\mathbf{y}}^2 = \|\hat{\mathbf{y}}_{\mathcal{S}_{\text{null}}}\|^2 \leq \|\hat{\mathbf{y}}_{\mathcal{S}}\|^2 \leq \|\mathbf{y}\|^2 \quad (14)$$

where  $\bar{\mathbf{y}}$  denotes the arithmetic mean  $\sum_{i=1}^n y_i/n$ . Assumption (A) follows from

$$(A') \quad 0 < \liminf_{n \rightarrow \infty} (\bar{\mathbf{y}})^2 \leq \limsup_{n \rightarrow \infty} (\frac{1}{n}\|\mathbf{y}\|^2) < \infty \quad \text{and} \quad \forall \mathcal{S}: \hat{\sigma}_{\mathcal{S}}^2 \rightarrow \sigma_{\mathcal{S}}^2 > 0 \text{ w.p.1.}$$

The first condition of (A') is obviously very mild and satisfied in almost all cases in practice. The second one is routinely used to derive asymptotic properties of model selection criteria (e.g., Theorem 2 of [Sha97] and Condition 1 of [WLT07]).

**Lemma 12 (LoRP for variable selection)** *The loss rank of model  $\mathcal{S}$  is*

$$\text{LR}_{\mathcal{S}} \equiv \text{LR}_{\mathcal{S}}^{\alpha_m} = \frac{n}{2} \log(n\hat{\sigma}_{\mathcal{S}}^2) + \frac{n}{2} H(\frac{|\mathcal{S}|}{n}) + \frac{d}{2} \log \frac{1-\rho_{\mathcal{S}}}{\rho_{\mathcal{S}}} \quad (15)$$

where  $\rho_{\mathcal{S}}$  and  $\hat{\sigma}_{\mathcal{S}}^2$  are defined in (14), and  $H(p) := -p \log p - (1-p) \log(1-p)$  is the entropy of  $p$ . Under Assumption (A) or (A'), after neglecting constants independent of  $\mathcal{S}$ , the loss rank of model  $\mathcal{S}$  has the form

$$\text{LR}_{\mathcal{S}} = \frac{n}{2} \log \hat{\sigma}_{\mathcal{S}}^2 + \frac{|\mathcal{S}|}{2} \log n + O_{\mathbb{P}}(1), \quad (16)$$

where  $O_{\mathbb{P}}(1)$  denotes a bounded random variable w.p.1.

**Proof.** Inserting  $\mathbf{y}^\top \mathbf{y} = n\hat{\sigma}_{\mathcal{S}}^2/\rho_{\mathcal{S}}$  into (12) and rearranging terms gives (15). By Assumption (A) the last term in (15) is bounded w.p.1. Taylor expansion  $\log(1-p) = -p + O(p^2)$  implies  $H(p)/p + \log p \rightarrow 1$ , hence  $\frac{n}{2} H(\frac{|\mathcal{S}|}{n}) = \frac{|\mathcal{S}|}{2} \log n + O(1)$ . Finally, dropping the  $\mathcal{S}$ -independent term  $\frac{n}{2} \log n$  from (15) gives (16). ■

This lemma implies that the loss rank  $\text{LR}_{\mathcal{S}}$  here is a BIC-type criterion, thus we immediately can state without proof the following theorem which is the well-known model consistency of BIC-type criteria (interested readers can find the routine proof in, for example, [Cha06]).

**Theorem 13 (Model consistency)** *Under Assumption (A) or (A'), LoRP is model consistent for variable selection in the sense that the probability of selecting the true model goes to 1 for data size  $n \rightarrow \infty$ .*

**The optimal regression estimation of LoRP.** The second goal of model selection is often measured by the (asymptotic) mean efficiency [Shi83] which is briefly defined as follows. Let  $\mathcal{S}_T$  denote the true model (which may contain an infinite number of covariates). For a candidate model  $\mathcal{S}$ , let  $L_n(\mathcal{S}) = \|X_{\mathcal{S}_T}\boldsymbol{\beta}_{\mathcal{S}_T} - X_{\mathcal{S}}\hat{\boldsymbol{\beta}}_{\mathcal{S}}\|^2$  be the squared loss where  $\hat{\boldsymbol{\beta}}_{\mathcal{S}}$  is the OLS estimate, and  $R_n(\mathcal{S}) = \mathbf{E}[L_n(\mathcal{S})]$  be the risk. The mean efficiency of a selection criterion  $\delta$  is defined by the ratio

$$\text{eff}(\delta) = \frac{\inf_{\mathcal{S}} R_n(\mathcal{S})}{\mathbf{E}[L_n(\mathcal{S}_\delta)]} \leq 1$$

where  $\mathcal{S}_\delta$  is the model selected by  $\delta$ .  $\delta$  is said to be asymptotically mean efficient if  $\liminf_{n \rightarrow \infty} \text{eff}(\delta) = 1$ .

By minimizing the loss rank in  $\alpha$  we have shown in the previous paragraph that LoRP satisfies the first goal of model selection. We now show that with a suitable choice of  $\alpha$ , LoRP also satisfies the second goal.

From (11), we have

$$\text{LR}_{\mathcal{S}}^{\alpha}(\mathbf{y}|\mathbf{x}) = \frac{n}{2} \log(\hat{\sigma}_{\mathcal{S}}^2 + \frac{\alpha}{n} \mathbf{y}^{\top} \mathbf{y}) + \frac{n}{2} \log n - \frac{|\mathcal{S}|}{2} \log(\alpha) - \frac{n-|\mathcal{S}|}{2} \log(1 + \alpha).$$

By choosing  $\alpha = \tilde{\alpha} = \exp(-\frac{n(n+|\mathcal{S}|)}{|\mathcal{S}|(n-|\mathcal{S}|-2)})$ , under Assumption (A), the loss rank of model  $\mathcal{S}$  (neglecting the common constant  $\frac{n}{2} \log n$ ) is proportional to

$$\text{LR}_{\mathcal{S}}^{\tilde{\alpha}}(\mathbf{y}|\mathbf{x}) = n \log \hat{\sigma}_{\mathcal{S}}^2 + \frac{n(n+|\mathcal{S}|)}{n-|\mathcal{S}|-2} + o_{\mathbf{P}}(1),$$

which is the corrected AIC of [HT89]. As a result,  $\text{LoRP}(\tilde{\alpha})$  is optimal in terms of regression estimation, i.e., it is asymptotically mean efficient ([Shi83], 1983; [Sha97], 1997).

**Theorem 14 (Asymptotic mean efficiency)** *Under Assumption (A) or (A'), with a suitable choice of  $\alpha$ , the loss rank is proportional to the corrected AIC. As a result, LoRP is asymptotically mean efficient.*

## 5 Experiments

In this section we present a simulation study for LoRP, compare it to other methods and also demonstrate how LoRP can be used for some specific problems like choosing tuning parameters for kNN and spline regression. All experiments are conducted by using MATLAB software and the source code is freely available at <http://www.hutter1.net/ai/lorpcode.zip>.

**Comparison to AIC and BIC for model identification.** Samples are generated from the model

$$\mathbf{y} = \beta_0 + \beta_1 X_1 + \dots + \beta_d X_d + \epsilon, \quad \epsilon \sim N(0, \sigma^2) \quad (17)$$

where  $\boldsymbol{\beta}$  is the vector of coefficients with some zero entries. Without loss of generality, we assume that  $\beta_0 = 0$ , otherwise, we can center the response vector  $\mathbf{y}$  and standardize the design matrix  $X$  to exclude  $\beta_0$  from the model. We shall compare the performance of LoRP to that of BIC and AIC with various factors  $n$ ,  $d$  and signal-to-noise ratio (SNR) which is  $\|\boldsymbol{\beta}\|^2/\sigma^2$  ( $\|\boldsymbol{\beta}\|^2$  is often called the length of the signal).

For a given set of factors ( $n$ ,  $d$ , SNR), the way we simulate a dataset from model (17) is as follows. Entries of  $X$  are sampled from a uniform distribution on  $[-1, 1]$ . To generate  $\boldsymbol{\beta}$ , we first create a vector  $\mathbf{u} = (u_1, \dots, u_d)^\top$  whose entries are sampled from a uniform distribution on  $[-1, 1]$ . The number of true covariates  $d^*$  is randomly selected from  $\{1, 2, \dots, d\}$ , the last  $d - d^*$  entries of  $\mathbf{u}$  are set to zero, then coefficient vector  $\boldsymbol{\beta}$  is computed by  $\beta_i = \{\text{length of signal}\} * u_i / \|\mathbf{u}\|$ . In our simulation, the length of signal was fixed to be 10.  $n$  observation errors  $\epsilon_1, \dots, \epsilon_n$  are sampled from a normal distribution with mean 0 and variance  $\sigma^2 = \|\boldsymbol{\beta}\|^2 / \text{SNR}$ . Finally, the response vector is computed by  $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\epsilon}$ . For each set of factors ( $n$ ,  $d$ , SNR), 1000 datasets are simulated in the same manner to assess the average performance of the methods. For simplicity, a candidate model is specified by its order, i.e., we search the best model among only  $d$  models  $\{1\}, \{1, 2\}, \dots, \{1, 2, \dots, d\}$ . For the general case, an efficient branch-and-bound algorithm [Mil02, Chp.3] can be used to exhaustively search for the best subsets.

Table 1 presents percentages of correctly-fitted models with various factors  $n$ ,  $d$  and SNR. As shown, LoRP outperforms the others. The better performance of LoRP over BIC, which is the most popular criterion for model identification, is very encouraging. This is probably because LoRP is a selection criterion with a data-dependent penalty. This improvement needs a theoretical justification which we intend to do in the future.

Table 1: Percentage of correctly-fitted models over 1000 replications

$n$	$d$	SNR	AIC	BIC	LoRP	$n$	$d$	SNR	AIC	BIC	LoRP
100	5	1	62	62	69	300	5	1	74	82	83
		5	85	85	86			5	78	90	91
		10	80	90	91			10	81	94	94
	10	1	52	42	54		10	1	63	67	71
		5	63	77	77			5	70	85	86
		10	68	84	85			10	74	90	90
	20	1	32	22	36		20	1	54	45	61
		5	55	63	65			5	64	79	80
		10	56	73	74			10	67	85	85

**Comparison to AIC and BIC for regression estimation.** Consider the following model which is from [Shi83]

$$y = y(x) = \log \frac{1}{1-x} + \epsilon, \quad \epsilon \sim N(0, \sigma^2), \quad x \in [0, 1). \quad (18)$$

We approximate the true function by a Fourier series and consider the problem of choosing a good order among models

$$y = \beta_0 + \sum_{l=1}^{k-1} \frac{\cos(\pi l x / \delta)}{l+1} \beta_l + \epsilon, \quad k = 1, \dots, K.$$

In the present context, a model in Section 4 is completely specified by the order  $K$  of the Fourier series. Samples are created from (18) at the points  $x_i = \delta \frac{i}{n+1}$ ,  $i = 1, \dots, n$ . As in [Shi83], we take  $\delta = .99$ , and  $K = 163$  with various  $n$  and  $\sigma$ . The performance is measured by the estimate of mean efficiency over 1000 replications.

Table 2 represents the simulation results. In general, LoRP (with  $\alpha = \tilde{\alpha}$  as in Section 4) outperforms the others, except for cases with unrealistically high noise level. For cases with high noise, mean efficiency of BIC is often larger than that of AIC and LoRP. This was also shown in the simulation study of [Shi83], Table 1. This phenomenon can be explained as follows.

The risk of model  $k$  (the model specified by its order  $k$ ) is  $R_n(k) = \|(I - M_k)\mathbf{y}_{\text{true}}\|^2 + k\sigma^2$  where  $M_k$  is the regression matrix under model  $k$  and  $\mathbf{y}_{\text{true}}$  is the vector of true values  $y(x_i)$ . When  $\sigma \rightarrow \infty$ , the ideal  $k^* = \operatorname{arginf}_k R_n(k) \rightarrow 1$ . Because BIC penalizes the model complexity more strongly than AIC and LoRP do, the order chosen by BIC is closer to  $k^* = 1$  than the ones chosen by AIC and LoRP. As a result, mean efficiency of BIC is larger than that of the others.

Table 2: Estimates of mean efficiency over 1000 replications

$n$	$\sigma$	AIC	BIC	LoRP	$n$	$\sigma$	AIC	BIC	LoRP
400	.001	1.00	.98	.99	600	.001	1.00	.98	1.00
	.01	.93	.68	.90		.01	.99	.67	.92
	.05	.88	.67	.95		.05	.90	.66	.94
	.1	.88	.67	.92		.1	.90	.67	.93
	.5	.81	.66	.85		.5	.82	.66	.83
	1	.79	.63	.82		1	.79	.65	.82
	5	.67	.65	.70		5	.65	.67	.66
	10	.54	.67	.59		10	.54	.59	.54
100	.31	.89	.33	100	.40	.90	.41		

**LoRP for selecting a good number of neighbors in kNN.** Let us now see how LoRP can be applied to select a good parameter  $k$  in kNN regression.

We created a dataset of  $n = 100$  observations  $(x_i, y_i)$  from the model:

$$y = f(x) + \epsilon, \quad \text{with } f(x) = \frac{\sin(12(x+0.2))}{x+0.2}, \quad x \in [0, 1] \quad (19)$$

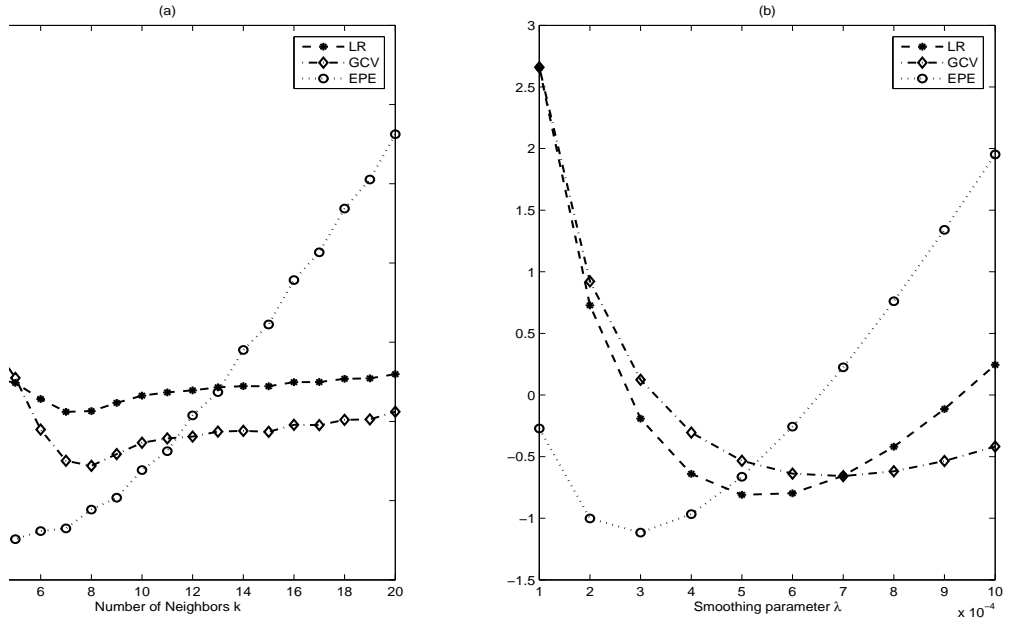


Figure 1: Choosing the tuning parameters in kNN and spline regression. The curves have been scaled by their standard deviations.

where  $\varepsilon \sim N(0, \sigma^2)$  with  $\sigma = 0.5$ . The regression matrix  $M^{(k)}$  for kNN regression is determined by  $M_{ij}^{(k)} = \frac{1}{k}$  if  $j \in \mathcal{N}_k(x_i)$  and 0 else. Then, the loss rank is

$$\text{LR}(k) = \inf_{\alpha \geq 0} \left\{ \frac{n}{2} \log(\mathbf{y}^\top S_\alpha^{(k)} \mathbf{y}) - \frac{1}{2} \log \det S_\alpha^{(k)} \right\},$$

where  $S_\alpha^{(k)} = (I - M^{(k)})^\top (I - M^{(k)}) + \alpha I$ . The most widely-used method to select a good  $k$  is probably Generalized Cross-Validation (GCV) [CW79]:  $\text{GCV}(k) = n \| (I - M^{(k)}) \mathbf{y} \|^2 / [\text{tr}(I - M^{(k)})]^2$ . To judge how well GCV and LoRP work, we compare them to the expected prediction error defined as

$$\text{EPE}(k) = \sum_{i=1}^n \mathbf{E}(y_i - \hat{y}_i)^2 = \sum_{i=1}^n \left[ \sigma^2 + \left( f(x_i) - \frac{1}{k} \sum_{j \in \mathcal{N}_k(x_i)} f(x_j) \right)^2 + \frac{\sigma^2}{k} \right].$$

Figure 1(a) shows the curves  $\text{LR}(k)$ ,  $\text{GCV}(k)$ ,  $\text{EPE}(k)$  for  $k = 2, \dots, 20$  (the trivial case  $k=1$  is omitted), in which  $k=7$ -nearest neighbors is chosen by LoRP and  $k=8$  is chosen by GCV. The “ideal”  $k$  is 5. Both LoRP and GCV do a reasonable job. LoRP works slightly better than GCV.

**LoRP for selecting a good smoothing parameter.** We now further demonstrate the use of LoRP in selecting a good smoothing parameter for spline regression. Consider the following problem: find a function belonging to the class of functions with continuous 2nd derivative that minimizes the following penalized residual sum



of squares:

$$\text{RSS}(f) = \sum_{i=1}^n (y_i - f(x_i))^2 + \lambda \int (f''(t))^2 dt,$$

where  $\lambda$  is called the smoothing parameter. The second term penalizes the curvature of function  $f$  and the smoothing parameter  $\lambda$  controls the amount of penalty. Our goal is to choose a good  $\lambda$ .

It is well known (see, e.g., [HTF01], Section 5.4) that the solution is a natural spline  $f(x) = \sum_{j=1}^n N_j(x)\theta_j$  where  $N_1(x), \dots, N_n(x)$  are the basis functions of the natural cubic spline:

$$N_1(x) = 1, \quad N_2(x) = x, \quad N_{k+2}(x) = d_k(x) - d_{n-1}(x) \quad \text{with} \quad d_k(x) = \frac{(x-x_k)_+^3 - (x-x_n)_+^3}{x_n - x_k}.$$

The problem thus reduces to finding a vector  $\boldsymbol{\theta} \in \mathbb{R}^n$  that minimizes

$$\text{RSS}(\boldsymbol{\theta}) = (\mathbf{y} - N\boldsymbol{\theta})^\top (\mathbf{y} - N\boldsymbol{\theta}) + \lambda \boldsymbol{\theta}^\top \Omega \boldsymbol{\theta}$$

where  $N_{ij} = N_j(x_i)$  and  $\Omega_{ij} = \int N_i''(x)N_j''(x)dx$ . It is easy to see that the solution is  $\hat{\boldsymbol{\theta}}_\lambda = (N^\top N + \lambda\Omega)^{-1}N^\top \mathbf{y}$ , and the fitted vector is  $\hat{\mathbf{y}} = N\hat{\boldsymbol{\theta}}_\lambda = M_\lambda \mathbf{y}$  with  $M_\lambda = N(N^\top N + \lambda\Omega)^{-1}N^\top$ . The fitted vector is linear in  $\mathbf{y}$ , thus the loss rank is

$$\text{LR}(\lambda) = \arg \min_{\alpha \geq 0} \left\{ \frac{n}{2} \log(\mathbf{y}^\top S_\lambda^\alpha \mathbf{y}) - \frac{1}{2} \log \det S_\lambda^\alpha \right\}$$

where  $S_\lambda^\alpha = (I - M_\lambda)^\top (I - M_\lambda) + \alpha I$ .

Let us consider again the dataset generated from model (19). Figure 1(b) shows the curves  $\text{LR}(\lambda)$ ,  $\text{GCV}(\lambda)$  and  $\text{EPE}(\lambda)$ . The derivation of expressions for  $\text{GCV}(\lambda)$  and  $\text{EPE}(\lambda)$  is similar to the previous example.  $\lambda \approx 3 \times 10^{-4}$  is the optimal value selected by the ‘‘ideal’’ criterion EPE.  $\lambda \approx 5 \times 10^{-4}$  and  $\lambda \approx 7 \times 10^{-4}$  are selected by LoRP and GCV, respectively. One again, like the previous example, LoRP selects a better  $\lambda$  than GCV does.

## 6 Comparison to Gaussian Bayesian Linear Regression

We now consider LBFR from a Bayesian perspective with Gaussian noise and prior, and compare it to LoRP. In addition to the noise model as in PML, one also has to specify a prior. Bayesian model selection (BMS) proceeds by selecting the model that has largest evidence. In the special case of LBFR with Gaussian noise and prior and a type II maximum likelihood estimate for the noise variance, the expression for the evidence has a similar structure as the expression of the loss rank.

**Gaussian Bayesian LBFR / MAP.** Recall from Sec.3 Ex.9 that  $\mathcal{F}_d$  is the class of functions  $f_{\mathbf{w}}(x) = \mathbf{w}^\top \boldsymbol{\phi}(x)$  ( $\mathbf{w} \in \mathbb{R}^d$ ) that are linear in feature vector  $\boldsymbol{\phi}$ . Let

$$\text{Gauss}_N(\mathbf{z} | \boldsymbol{\mu}, \sigma) := \frac{\exp(-\frac{1}{2}(\mathbf{z} - \boldsymbol{\mu})^\top \sigma^{-1}(\mathbf{z} - \boldsymbol{\mu}))}{(2\pi)^{N/2} \sqrt{\det \sigma}} \quad (20)$$

denote a general  $N$ -dimensional Gaussian distribution with mean  $\boldsymbol{\mu}$  and covariance matrix  $\sigma$ . We assume that observations  $\mathbf{y}$  are perturbed from  $f_{\mathbf{w}}(x)$  by independent additive Gaussian noise with variance  $\beta^{-1}$  and zero mean, i.e., the likelihood of  $\mathbf{y}$  under model  $\mathbf{w}$  is  $P(\mathbf{y}|\mathbf{w}) = \text{Gauss}_n(\mathbf{y}|\Phi\mathbf{w}, \beta^{-1}I)$ , where  $\Phi_{ia} = \phi_a(x_i)$ . A Bayesian assumes a prior (before seeing  $\mathbf{y}$ ) distribution on  $\mathbf{w}$ . We assume a centered Gaussian with covariance matrix  $(\alpha C)^{-1}$ , i.e.,  $P(\mathbf{w}) = \text{Gauss}_d(\mathbf{w}|\mathbf{0}, \alpha^{-1}C^{-1})$ . From the prior and the likelihood one can compute the evidence and the posterior

$$\text{Evidence:} \quad P(\mathbf{y}) = \int P(\mathbf{y}|\mathbf{w})P(\mathbf{w})d\mathbf{w} = \text{Gauss}_n(\mathbf{y}|\mathbf{0}, \beta^{-1}S^{-1}) \quad (21)$$

$$\text{Posterior:} \quad P(\mathbf{w}|\mathbf{y}) = P(\mathbf{y}|\mathbf{w})P(\mathbf{w})/P(\mathbf{y}) = \text{Gauss}_d(\mathbf{w}|\hat{\mathbf{w}}, A^{-1})$$

$$B := \Phi^T\Phi, \quad A := \alpha C + \beta B, \quad M := \beta\Phi A^{-1}\Phi^T, \quad S := I - M, \quad (22)$$

$$\hat{\mathbf{w}} := \beta A^{-1}\Phi^T\mathbf{y}, \quad \hat{\mathbf{y}} := \Phi\hat{\mathbf{w}} = M\mathbf{y}$$

A standard Bayesian point estimate for  $\mathbf{w}$  for fixed  $d$  is the one that maximizes the posterior (MAP) (which in the Gaussian case coincides with the mean)  $\hat{\mathbf{w}} = \text{argmax}_{\mathbf{w}}P(\mathbf{w}|\mathbf{y}) = \beta A^{-1}\Phi^T\mathbf{y}$ . For  $\alpha \rightarrow 0$ , MAP reduces to Maximum Likelihood (ML), which in the Gaussian case coincides with the least squares regression of Ex.9. For  $\alpha > 0$ , the regression matrix  $M$  is not a projection anymore.

**Bayesian model selection.** Consider now a family of models  $\{\mathcal{F}_1, \mathcal{F}_2, \dots\}$ . Here the  $\mathcal{F}_d$  are the linear regressors with  $d$  basis functions, but in general they could be completely different model classes. All quantities in the previous paragraph implicitly depend on the choice of  $\mathcal{F}$ , which we now explicate with an index. In particular, the evidence for model class  $\mathcal{F}$  is  $P_{\mathcal{F}}(\mathbf{y})$ . BMS chooses the model class (here  $d$ )  $\mathcal{F}$  of highest evidence:

$$\mathcal{F}^{\text{BMS}} = \text{arg max}_{\mathcal{F}} P_{\mathcal{F}}(\mathbf{y})$$

Once the model class  $\mathcal{F}^{\text{BMS}}$  is determined, the MAP (or other) regression function  $f_{\mathbf{w}_{\mathcal{F}^{\text{BMS}}}}$  or  $M_{\mathcal{F}^{\text{BMS}}}$  are chosen. The data variance  $\beta^{-1}$  may be known or estimated from the data,  $C$  is often chosen  $I$ , and  $\alpha$  has to be chosen somehow. Note that while  $\alpha \rightarrow 0$  leads to a reasonable MAP=ML regressor for fixed  $d$ , this limit cannot be used for BMS.

**Comparison to LoRP.** Inserting (20) into (21) and taking the logarithm we see that BMS minimizes

$$-\log P_{\mathcal{F}}(\mathbf{y}) = \frac{\beta}{2}\mathbf{y}^T S \mathbf{y} - \frac{1}{2} \log \det S - \frac{n}{2} \log \frac{\beta}{2\pi} \quad (23)$$

w.r.t.  $\mathcal{F}$ . Let us estimate  $\beta$  by ML: We assume a broad prior  $\alpha \ll \beta$  so that  $\beta \frac{\partial S}{\partial \beta} = O(\frac{\alpha}{\beta})$  can be neglected. Then  $-\frac{\partial \log P_{\mathcal{F}}(\mathbf{y})}{\partial \beta} = \frac{1}{2}\mathbf{y}^T S \mathbf{y} - \frac{n}{2\beta} + O(\frac{\alpha}{\beta}n) = 0 \Leftrightarrow \beta \approx \hat{\beta} := n/(\mathbf{y}^T S \mathbf{y})$ . Inserting  $\hat{\beta}$  into (23) we get

$$-\log P_{\mathcal{F}}(\mathbf{y}) = \frac{n}{2} \log \mathbf{y}^T S \mathbf{y} - \frac{1}{2} \log \det S - \frac{n}{2} \log \frac{n}{2\pi e} \quad (24)$$

Taking an improper prior  $P(\beta) \propto \beta^{-1}$  and integrating out  $\beta$  leads for small  $\alpha$  to a similar result. The last term in (24) is a constant independent of  $\mathcal{F}$  and can be ignored. The first two terms have the same structure as in linear LoRP (10), but the matrix  $S$  is different. In both cases,  $\alpha$  act as regularizers, so we may minimize over  $\alpha$  in BMS like in LoRP. For  $\alpha=0$  (which neither makes sense in BMS nor in LoRP),  $M$  in BMS coincides with  $M$  of Ex.9, but still the  $S_0$  in LoRP is the square of the  $S$  in BMS. For  $\alpha>0$ ,  $M$  of BMS may be regarded as a regularized regressor as suggested in Sec.2 (a), rather than a regularized loss function (b) used in LoRP. Note also that BMS is limited to (semi)parametric regression, i.e., does not cover the non-parametric kNN Ex.2 and kernel Ex.8, unlike LoRP.

Since  $B$  only depends on  $\mathbf{x}$  (and not on  $\mathbf{y}$ ), and all  $P$  are implicitly conditioned on  $\mathbf{x}$ , one could choose  $C = B$ . In this case,  $M = \gamma \Phi B^{-1} \Phi^\top$ , with  $\gamma = \frac{\beta}{\alpha + \beta} < 1$  for  $\alpha > 0$ , is a simple multiplicative regularization of projection  $\Phi B^{-1} \Phi^\top$ , and (24) coincides with (11) for suitable  $\alpha$ , apart from an irrelevant additive constant, hence minimizing (24) over  $\alpha$  also leads to (12).

## 7 Comparison to other Model Selection Schemes

In this section we give a brief introduction to PML for (semi)parametric regression, and its major instantiations, AIC, BIC, and MDL principle, whose penalty terms are all proportional to the number of parameters  $d$ . The *effective number of parameters* is often much smaller than  $d$ , e.g., if there are soft constraints like in ridge regression. We compare MacKay’s trace formula [Mac92] for Gaussian Bayesian LBFR and Hastie’s et al. trace formula [HTF01] for general linear regression with LoRP.

**Penalized ML (AIC, BIC, MDL).** Consider a  $d$ -dimensional stochastic model class like the Gaussian Bayesian linear regression example of Section 6. Let  $P_d(\mathbf{y}|\mathbf{w})$  be the data likelihood under  $d$ -dimensional model  $\mathbf{w} \in \mathbb{R}^d$ . The maximum likelihood (ML) estimator for fixed  $d$  is

$$\hat{\mathbf{w}} = \arg \max_{\mathbf{w}} P_d(\mathbf{y}|\mathbf{w}) = \arg \min_{\mathbf{w}} \{-\log P_d(\mathbf{y}|\mathbf{w})\} \quad (25)$$

Since  $-\log P_d(\mathbf{y}|\mathbf{w})$  decreases with  $d$ , we cannot find the model dimension by simply minimizing over  $d$  (overfitting). Penalized ML adds a complexity term to get reasonable results

$$\hat{d} = \arg \min_d \{-\log P_d(\mathbf{y}|\hat{\mathbf{w}}) + \text{Penalty}(d)\} \quad (26)$$

The penalty introduces a tradeoff between the first and second term with a minimum at  $\hat{d} < \infty$ . Various penalties have been suggested: AIC [Aka73] uses  $d$ , BIC [Sch78] and the (crude) MDL [Ris78, Grü04] use  $\frac{d}{2} \log n$  for  $\text{Penalty}(d)$ . There are at least *three important conceptual differences* to LoRP:

- In order to apply PML one needs to specify not only a class of regression functions, but a full probabilistic model  $P_d(\mathbf{y}|\mathbf{w})$ ,

- PML ignores or at least does not tell how to incorporate a potentially given loss-function,
- PML is mostly limited to selecting between (semi)parametric models.

We discuss two approaches to the last item in the remainder of this section (where AIC, BIC, and MDL are not directly applicable): (a) for non-parametric models like kNN or kernel regression, or (b) if  $d$  does not reflect the “true” complexity of the model. [Mac92] suggests an expression for the effective number of parameters  $d_{eff}$  as a substitute for  $d$  in case of (b), while [HTF01] introduce another expression which is applicable for both (a) and (b).

**The trace penalty for parametric Gaussian LBFR.** We continue with the Gaussian Bayesian linear regression example (see Section 6 for details and notation). Performing the integration in (21), [Mac92, Eq.(21)] derives the following expression for the Bayesian evidence for  $C = I$

$$\begin{aligned} -\log P(\mathbf{y}) &= (\alpha \hat{E}_W + \beta \hat{E}_D) + \left(\frac{1}{2} \log \det A - \frac{d}{2} \log \alpha\right) - \frac{n}{2} \log \frac{\beta}{2\pi} \quad (27) \\ \hat{E}_D &= \frac{1}{2} \|\Phi \hat{\mathbf{w}} - \mathbf{y}\|_2^2, \quad \hat{E}_W = \frac{1}{2} \|\hat{\mathbf{w}}\|_2^2 \end{aligned}$$

(the first bracket in (27) equals  $\frac{\beta}{2} \mathbf{y}^\top S \mathbf{y}$  and the second equals  $-\frac{1}{2} \log \det S$ , cf. (23)). Minimizing (27) w.r.t.  $\alpha$  leads to the following relation:

$$0 = \frac{-\partial \log P(\mathbf{y})}{\partial \alpha} = \hat{E}_W + \frac{1}{2} \text{tr} A^{-1} - \frac{d}{2\alpha} \quad \left(\frac{\partial}{\partial \alpha} \log \det A = \text{tr} A^{-1}\right)$$

He argues that  $\alpha \|\hat{\mathbf{w}}\|_2^2$  corresponds to the effective number of parameters, hence

$$d_{eff}^{\text{McK}} := \alpha \|\hat{\mathbf{w}}\|_2^2 = 2\alpha \hat{E}_W = d - \alpha \text{tr} A^{-1} \quad (28)$$

**The trace penalty for general linear models.** We now return to general linear regression  $\hat{\mathbf{y}} = M(\mathbf{x})\mathbf{y}$  (7). LBFR is a special case of a projection matrix  $M = M^2$  with rank  $d = \text{tr} M$  being the number of basis functions.  $M$  leaves  $d$  directions untouched and projects all other  $n - d$  directions to zero. For general  $M$ , [HTF01, Sec.5.4.1] argue to regard a direction that is only somewhat shrunken, say by a factor of  $0 < \beta < 1$ , as a fractional parameter ( $\beta$  degrees of freedom). If  $\beta_1, \dots, \beta_n$  are the shrinkages = eigenvalues of  $M$ , the effective number of parameters could be defined as [HTF01, Sec.7.6]

$$d_{eff}^{\text{HTF}} := \sum_{i=1}^n \beta_i = \text{tr} M,$$

where HTF stands for Hastie-Tibshirani-Friedman, which generalizes the relation  $d = \text{tr} M$  beyond projections. For MacKay’s  $M$  (22),  $\text{tr} M = d - \alpha \text{tr} A^{-1}$ , i.e.,  $d_{eff}^{\text{HTF}}$  is consistent with and generalizes  $d_{eff}^{\text{McK}}$ .

**Problems.** Though nicely motivated, the trace formula is not without problems. First, since for projections,  $M = M^2$ , one could have argued equally well for  $d_{eff}^{\text{HTF}} =$

$\text{tr}M^2$ . Second, for kNN we have  $\text{tr}M = \frac{n}{k}$  (since  $M$  is  $\frac{1}{k}$  on the diagonal), which does not look unreasonable. Consider now kNN', which is defined as follows: we average over the  $k$  nearest neighbors *excluding* the closest neighbor. For sufficiently smooth functions, kNN' for suitable  $k$  is still a reasonable regressor, but  $\text{tr}M = 0$  (since  $M$  is zero on the diagonal). So  $d_{\text{eff}}^{\text{HTF}} = 0$  for kNN', which makes no sense and would lead one to always select the  $k=1$  model.

**Relation to LoRP.** In the case of kNN',  $\text{tr}M^2$  would be a better estimate for the effective dimension. In linear LoRP,  $-\log\det S_\alpha$  serves as complexity penalty. Ignoring the nullspace of  $S_0 = (I - M)^\top(I - M)$  (8), we can Taylor expand  $-\frac{1}{2}\log\det S_0$  in  $M$

$$-\frac{1}{2}\log\det S_0 = -\text{tr}\log(I - M) = \sum_{s=1}^{\infty} \frac{1}{s}\text{tr}(M^s) = \text{tr}M + \frac{1}{2}\text{tr}M^2 + \dots$$

For BMS (24) with  $S = I - M$  (22) we get half of this value. So the trace penalty may be regarded as a leading order approximation to LoRP. The higher order terms prevent peculiarities like in kNN'.

**Coding/MDL interpretation of LoRP.** The basic idea of MDL is as follows [Grü04]: “The goal of statistical inferences may be cast as trying to find regularity in the data. ‘Regularity’ may be identified with ‘ability to compress’. MDL combines these two insights by *viewing learning as data compression*: it tells us that, for a given set of hypotheses  $\mathcal{H}$  and data set  $D$ , we should try to find the hypothesis or combination of hypotheses in  $\mathcal{H}$  that compress  $D$  most.”

The standard incarnation of (crude) MDL is as follows: If  $H$  is a stochastic model of (discrete) data  $D$ , we can code  $D$  (by Shannon-Fano) in  $[-\log_2\text{P}(D|H)]$  bits. If we have a class of models  $\mathcal{H}$ , we also have to code  $H$  (somehow in, say,  $L(H)$  bits) in order to be able to decode  $D$ . MDL chooses the hypothesis  $H^{\text{MDL}} = \text{argmin}_{H \in \mathcal{H}} \{-\log_2\text{P}(D|H) + L(H)\}$  of minimal two-part code. For instance, if  $\mathcal{H}$  is the class of all polynomials of all degrees with each coefficient coded to  $\frac{1}{2}\log_2 n$  bits (i.e.,  $O(n^{-1/2})$  accuracy) and we condition on  $x$ , i.e.,  $D \rightsquigarrow \mathbf{y}|\mathbf{x}$ , MDL takes the form (25) and (26), i.e.,  $H^{\text{MDL}} = (\hat{\mathbf{w}}, \hat{d})$ .

We now give LoRP (for discrete  $D$ ) a data compression/MDL interpretation. For simplicity, we will first assume that all loss values are different, i.e., if  $\text{Loss}_r(\mathbf{y}'|\mathbf{x}) \neq \text{Loss}_r(\mathbf{y}''|\mathbf{x})$  for  $\mathbf{y}' \neq \mathbf{y}''$  (adding infinitesimal random noise to  $\text{Loss}_r$  easily ensures this). In this case,  $\text{Rank}_r(\cdot|\mathbf{x}): \mathcal{Y}^n \rightarrow \mathcal{N}$  is an order preserving bijection, i.e.,  $\text{Rank}_r(\mathbf{y}'|\mathbf{x}) < \text{Rank}_r(\mathbf{y}''|\mathbf{x})$  iff  $\text{Loss}_r(\mathbf{y}'|\mathbf{x}) < \text{Loss}_r(\mathbf{y}''|\mathbf{x})$  with no gaps in the range of  $\text{Rank}_r(\cdot|\mathbf{x})$ .

Phrased differently,  $\text{Rank}_r(\cdot|\mathbf{x})$  codes each  $\mathbf{y}' \in \mathcal{Y}^n$  as a natural number  $m$  in increasing loss-order. The natural number  $m$  can itself be coded in  $\lceil \log_2 m \rceil$  bits (using plain not prefix coding). Let us call this code of  $\mathbf{y}'$  the *Loss Rank Code* (LRC). LRC has a nice characterization: LRC is the shortest loss-order preserving code. Ignoring the rounding, the *Length* of  $\text{LRC}_r(\mathbf{y}'|\mathbf{x})$  is  $\text{LR}_r(\mathbf{y}'|\mathbf{x})$ :

**Proposition 15 (Minimality property)** *If all loss values are different, i.e., if*

$$\text{Loss}_r(\mathbf{y}'|\mathbf{x}) \neq \text{Loss}_r(\mathbf{y}''|\mathbf{x}) \text{ for all } \mathbf{y}' \neq \mathbf{y}''$$

*then the loss rank (code) of  $\mathbf{y}$  is the smallest/shortest among all loss-order preserving rankings/codes  $C$  in the sense that*

$$\begin{aligned} \text{Rank}(\mathbf{y}) &= \min\{C(\mathbf{y}) : C \in \mathcal{Y}^n \rightarrow \mathbb{N} \wedge (\star)\} \\ \lfloor \text{LR}(\mathbf{y})/\log 2 \rfloor &= \min\{\text{Length}(C(\mathbf{y})) : C \in \mathcal{Y}^n \rightarrow \{0, 1\}^* \wedge (\star)\} \\ (\star) &:= [\text{Loss}(\mathbf{y}') < \text{Loss}(\mathbf{y}'') \Leftrightarrow C(\mathbf{y}') < C(\mathbf{y}''), \forall \mathbf{y}', \mathbf{y}''] \end{aligned}$$

The proof follows from the fact that if a discrete injection (code) is order preserving, there exists a “smallest” one without gaps in the range. So LoRP minimizes the Loss Rank Code, where LRC itself is the shortest among all loss-order preserving codes. From this perspective, LoRP is just a different (non-stochastic, non-parametric, loss-based) incarnation of MDL. The MDL philosophy provides a justification of LoRP (2), its regularization (5), and loss function selection (Section 8). This identification should also allow to apply or adapt the various consistency results of MDL, implying that LoRP is consistent under some mild conditions.

If some losses are equal,  $\text{Rank}_r(\cdot|\mathbf{x}) : \mathcal{Y}^n \rightarrow \mathbb{N}$  still preserves the order  $\leq$ , but the mapping is neither surjective nor injective anymore.

**Large regression classes  $\mathcal{R}$ .** The classes  $\mathcal{R}$  of regressors we considered so far were discrete and “small”, often indexed by an integer complexity index (like  $k$  in kNN or  $d$  in LBFR). But large classes are also thinkable.

As an extreme case, consider the class of *all* regressors. Clearly, there is an  $r=r_D$  which “knows”  $D$  and perfectly fits  $D$  ( $r(x_i|D) = y_i, \forall i$ ), but is the worst possible on all other  $D'$  ( $r(x_i|D') = \infty, \forall i, \forall D' \neq D$ ). This  $r$  has (discrete) Rank 1, so is best according to LoRP. So if  $\mathcal{R}$  is too large, LoRP can overfit too.

Consider a more realistic example by not taking *all* of the first  $d$  basis functions in LBFR, but selecting *some* basis functions  $\phi_{i_1}, \dots, \phi_{i_d}$ , i.e.,  $\mathcal{R}$  is indexed by  $d$  integers, and  $d$  may be variable too.

One solution approach is to group more regressors in  $\mathcal{R}$  into one function class  $\mathcal{F}$ , e.g., the class of functions  $\mathcal{F}_{k,d} = \{w_1\phi_{i_1} + \dots + w_d\phi_{i_d} : \mathbf{w} \in \mathbb{R}^d, 1 \leq i_1 < \dots < i_d \leq k\}$  that are linear in  $d$  of the first  $k$  bases. Now  $\mathcal{R}$  is a small class indexed by  $d$  and  $k$  only.

Looking at the coding interpretation of  $\text{LR}_r$  and the MDL philosophy, suggests to assign a code to  $r \in \mathbb{R}$  in order to get a complete code for  $D$ :

$$r^{best} = \arg \min_r \{\text{LR}_r(\mathbf{y}|\mathbf{x}) + L(r)\}$$

where  $r$  is the length of a code for  $r$  (given  $\mathcal{R}$ ). For  $\mathcal{R} \simeq \mathbb{N}$  a single integer has to be coded, e.g.,  $k$  in  $L(r) = L(k) \approx \log_2 k$  bits, which can usually be safely dropped/ignored. For more complex classes like the (ungrouped) LBFR subset selection above,  $L(r) = L(i_1, \dots, i_d, d) \approx d \log_2 k + \log_2 d$  can become important.

## 8 Loss Functions and their Selection

**General additive loss.** Linear LoRP  $\hat{\mathbf{y}} = M(\mathbf{x})\mathbf{y}$  of Section 3 can easily be generalized to non-quadratic loss. Let us consider the  $\rho > 0$  loss

$$\begin{aligned} \text{Loss}_M(\mathbf{y}|\mathbf{x}) &:= (\sum_{i=1}^n (y_i - \hat{y}_i)^\rho)^{1/\rho} = \|\mathbf{y} - \hat{\mathbf{y}}\|_\rho = \|(I-M)\mathbf{y}\|_\rho \\ V(L) &= \{\mathbf{y}' \in \mathbb{R}^n : \|(I-M)\mathbf{y}'\|_\rho \leq L\} = \{(I-M)^{-1}\mathbf{z} \in \mathbb{R}^n : \|\mathbf{z}\|_\rho \leq L\} \\ \text{Let } v_n^\rho &:= |\{\mathbf{z} \in \mathbb{R}^n : \|\mathbf{z}\|_\rho \leq 1\}| = 2^n \prod_{i=1}^{n-1} \frac{i!^{1/\rho}}{\rho^{1/\rho}} / \frac{i+1!}{\rho^{1/\rho}}, \end{aligned}$$

where  $\frac{i!}{\rho} := \Gamma(\frac{i}{\rho} + 1)$ , be the volume of the unit  $d$ -dimensional  $\rho$ -norm ‘‘ball’’. Since  $V(L)$  is a linear transformation of this ball with transformation matrix  $(I-M)^{-1}$  and scaling  $L$ , we have  $|V(L)| = v_n^\rho L^n / \det(I-M)$ , hence

$$\text{LR}_M(\mathbf{y}|\mathbf{x}) = \log |V(\text{Loss}_M(\mathbf{y}|\mathbf{x}))| = n \log \|(I-M)\mathbf{y}\|_\rho - \log \det(I-M) + \log v_n^\rho \quad (29)$$

For the  $\rho=2$  norm, (29) reduces to  $\text{LR}_M^0$  (9). Note that  $\text{Loss}_M := g(\|\mathbf{y} - \hat{\mathbf{y}}\|_\rho)$  leads to the same result (29) for any monotone increasing  $g$ , i.e., only the order of the loss matters, not its absolute value. More generally  $\text{Loss}_M = g(\sum_i h(y_i - \hat{y}_i))$  for any  $h$  implies

$$\begin{aligned} \text{LR}_M(\mathbf{y}|\mathbf{x}) &= n \log v_n^h(\sum_i h(y_i - \hat{y}_i)) - \log \det(I-M), \quad \text{where} \\ v_n^h(l) &:= |\{\mathbf{z} \in \mathbb{R}^n : \sum_i h(z_i) \leq l\}|^{1/n} \end{aligned}$$

is a one-dimensional function of  $l$  (independent  $D$  and  $M$ ), once to be determined (e.g.,  $v_n^h(l) = l \cdot (v_n^\rho)^{1/n} \propto l$  for  $\rho$ -norm loss). Regularization may be performed by  $M \rightsquigarrow \gamma M$  with optimization over  $\gamma < 1$ .

**Loss-function selection.** In principle, the loss function should be part of the problem specification, since it characterizes the ultimate goal. For instance, whether a test should more likely classify a healthy person as sick than a sick person as healthy, depends on the severity of a misclassification (loss) in each direction. In reality, though, having to specify the loss function can be a nuisance. Sure, the loss has to respect some general features, e.g., that it increases with the deviation of  $\hat{y}_i$  from  $y_i$ . Otherwise it is chosen by convenience or rules of thumb, rather than by elicitation of the real goal, for instance preferring the Euclidean norm over  $\rho \neq 2$  norms. If we subscribe to the procedure of *choosing* the loss function, we could ask whether this may be done in a more principled way. Consider a (not too large) class of loss functions  $\text{Loss}^\alpha$ , indexed by some parameter  $\alpha$ . For instance,  $\text{Loss}^\alpha = \|\mathbf{y} - \hat{\mathbf{y}}\|_\alpha$  from the previous paragraph. The regularized loss (5) also constitutes a class of losses. In this case we minimized over the regularization parameter  $\alpha$ . This suggests to choose in general the loss function that has minimal loss rank  $\text{LR}_r^\alpha$ . The justifications are similar to the ones for minimizing  $\text{LR}_r^\alpha$  w.r.t.  $r$ . Note that the term  $\log v_n^\rho$  cannot be dropped anymore, unlike in (10).

## 9 Self-Consistent Regression

So far we have considered only “on-data” regression. LoRP only depends on the regressor  $r$  on data  $D$  and not on  $x \notin \{x_1, \dots, x_n\}$ . We now construct canonical regressors for off-data  $x$  from regressors given only on-data. First, this may ease the specification of the regression functions, second, it is a canonical way for interpolation (LoRP can’t distinguish between  $r$  that are identical on  $D$ ), and third, we show that many standard regressors (kNN, Kernel, LBFR) are self-consistent in the sense that they are canonical. We limit our exposition to linear regression.

**Off-data regression.** A linear regressor is completely determined by the  $n$  functions  $m_j$  (6), but not by the matrix function  $M$  (7). Indeed, two sets  $\{m_j\}$  and  $\{m'_j\}$  that coincide on  $D = (\mathbf{x}, \mathbf{y})$ , i.e.  $m_j(x_i | \mathbf{x}) = m'_j(x_i | \mathbf{x}) \forall i, j$  but possibly differ for  $x \notin \mathbf{x}$ , lead to the same matrix  $M_{ij}(\mathbf{x}) = m_j(x_i | \mathbf{x}) = m'_j(x_i | \mathbf{x})$ . LoRP has the advantage of only depending on  $M$ , but this also means that it cannot distinguish between an  $m_j$  that behaves well on  $x \notin \mathbf{x}$  and one that, e.g., wildly oscillates outside  $\mathbf{x}$ .

Typically, the  $m_j$  are given and, provided the model complexity is chosen appropriately e.g. by LoRP, they properly interpolate  $\mathbf{x}$ . Nevertheless, a canonical extension from  $M$  to  $m_j$  would be nice. In this way LoRP would not be vulnerable to bad  $m_j$ , and we could interpolate  $D$  (predict  $y$  for any  $x \in \mathcal{X}$ ) even without  $m_j$  given a-priori.

We define a self-consistent regression scheme based only on  $M$  (for all  $n$ ). We ask for an estimate  $\hat{y}$  of  $y$  for  $x \notin \mathbf{x}$ . We add a virtual data point  $(x_0, y_0)$  to  $D$ , where  $x_0 = x$ . If we knew  $y_0 = y$  we could estimate  $\hat{y}_0 = r(x_0 | \{(x_0, y_0)\} \cup D)$ , but we don’t know  $y_0$ . But we could require a self-consistency condition, namely that  $\hat{y}_0 = y_0$  for  $x_0 \notin \mathbf{x}$ .

**Definition 16 (canonical and self-consistent regressors)** *Let  $M'_{ij}(\mathbf{x}')_{0 \leq i, j \leq n}$  be the regression matrix for the data set  $D' = \{(x_0, y_0)\} \cup D = ((x_0, \mathbf{x}), (y_0, \mathbf{y})) = (\mathbf{x}', \mathbf{y}')$  of size  $n+1$ .*

- (i) *A linear regressor  $\tilde{y}_0 = \tilde{r}(x_0 | D)$  is called a canonical regressor for  $M'$  if the consistency condition  $\tilde{y}_0 = r(x_0 | D') \equiv \sum_{j=0}^n M'_{0j} y_j$  holds  $\forall x_0, D$ .*
- (ii) *A regressor  $r$  is called self-consistent if  $\tilde{r} = r$ , i.e. if  $r(x_0 | \{(x_0, r(x_0 | D))\} \cup D) = r(x_0 | D) \forall x_0, D$ .*
- (iii) *A class of regressors  $\mathcal{R} = \{r\}$  is called self-consistent if  $\tilde{\mathcal{R}} = \{\tilde{r}\} \subseteq \mathcal{R}$ .*

We denote the solution of the self-consistency condition  $y_0 = \sum_{j=0}^n M'_{0j} y_j$  by  $\tilde{y}_0$ . So we have to solve

$$\tilde{y}_0 = \sum_{j=1}^n M'_{0j} y_j + M'_{00} \tilde{y}_0 \implies \tilde{y}_0 = \frac{\sum_{j=1}^n M'_{0j} y_j}{1 - M'_{00}} = \frac{\sum_{j=1}^n M'_{0j} y_j}{\sum_{j=1}^n M'_{0j}}$$

where the last equality only holds if  $\sum_{j=0}^n M'_{0j} = 1$ , which is often the case, in particular for kNN and Kernel regression, but not necessarily for LBFR.



**Proposition 17 (canonical regressor)** *The linear regressor*

$$y_0 = \tilde{r}(x_0|D) := \sum_{j=1}^n \tilde{m}_j(x_0|\mathbf{x})y_j, \quad \text{where} \quad \tilde{m}_j(x_0|\mathbf{x}) := \frac{M'_{0j}(\mathbf{x}')}{1 - M'_{00}(\mathbf{x}')}$$

is the unique canonical regressor for  $M'$  (if  $M'_{00} < 1$ ).

**Example 18 (self-consistent kNN, ↑Ex.2)**  $M'_{0j}(\mathbf{x}') = \frac{1}{k}$  for  $x_j \in \mathcal{N}'_k(x_0)$  and 0 else. The  $k$  nearest neighbors  $\mathcal{N}'_k(x_0)$  of  $x_0$  among  $\mathbf{x}'$  consist of  $x_0$  and the  $k-1$  nearest neighbors  $\mathcal{N}_{k-1}(x_0) =: J$  of  $x_0$  among  $\mathbf{x}$ , i.e.  $\mathcal{N}'_k(x_0) = \{x_0\} \cup \mathcal{N}_{k-1}(x_0)$ . Hence

$$\tilde{y}_0 = \frac{\sum_{j=1}^n M'_{0j}y_j}{\sum_{j=1}^n M'_{0j}} = \frac{\sum_{j \in J} \frac{1}{k}y_j}{\sum_{j \in J} \frac{1}{k}} = \sum_{j \in J} \frac{1}{k-1}y_j = \sum_{j=1}^n M_{0j}^{(k-1)}y_j = r_{k-1}(x_0|D) = \hat{y}_0$$

Canonical kNN is equivalent to standard  $(k-1)$ NN, so the class of canonical kNN regressors coincides with the standard kNN class.  $\diamond$

**Example 19 (self-consistent kernel)**

$$M'_{0j}(\mathbf{x}') = \frac{K(x_0, x_j)}{\sum_{j=0}^n K(x_0, x_j)} \implies \tilde{y}_0 = \frac{\sum_{j=1}^n K(x_0, x_j)y_j}{\sum_{j=1}^n K(x_0, x_j)} = r(x_0|D) = \hat{y}_0$$

Canonical kernel regression coincides with the standard kernel smoother.  $\diamond$

**Example 20 (self-consistent LBFR)**

$$\begin{aligned} B' &= \sum_{i=0}^n \phi(x_i)\phi(x_i)^\top = B + \phi(x_0)\phi(x_0)^\top \\ \implies M'_{0j} &= \phi(x_0)^\top B'^{-1} \phi(x_j) = \phi(x_0)^\top \left[ B^{-1} - \frac{B^{-1} \phi(x_0)\phi(x_0)^\top B^{-1}}{1 + \phi(x_0)^\top B^{-1} \phi(x_0)} \right] \phi(x_j) \\ &= M_{0j} - \frac{M_{00}M_{0j}}{1 + M_{00}} = \frac{M_{0j}}{1 + M_{00}} \implies 1 - M'_{00} = \frac{1}{1 + M_{00}} \end{aligned}$$

In the first line we used the Sherman-Morrison formula for inverting  $B'$ . In the second line we defined  $M_{0j} = \phi(x_0)^\top B^{-1} \phi(x_j)$ , extending  $M$ .

$$\implies \tilde{y}_0 = \frac{\sum_{j=1}^n M'_{0j}y_j}{1 - M'_{00}} = \sum_{j=1}^n M_{0j}y_j = \sum_{j=1}^n m_j(x_0, \mathbf{x})y_j = \hat{y}_0$$

Canonical LBFR coincides with standard LBFR.  $\diamond$

**Proposition 21 (self-consistent regressors)** *Kernel regression and linear basis function regression are self-consistent. kNN is not self-consistent but the class of kNN regressors  $\mathcal{R} = \{r_{kNN} : k \in \mathbb{N}\}$  is self-consistent.*

To summarize, we expect LoRP to select good regressors with proper interpolation behavior for canonical and self-consistent regressors.

## 10 Nearest Neighbors Classification

We now consider k-nearest neighbors classification in more detail. In order to get more insight into LoRP we seek a case that allows analytic solution. In general, the determinant  $\det S_\alpha$  cannot be computed analytically, but for  $\mathbf{x}$  lying on a hypercube of the regular grid  $\mathcal{X} = \mathbb{Z}^d$  we can. We derive exact expressions, and consider the limits  $n \rightarrow \infty$ ,  $k \rightarrow \infty$ , and  $d \rightarrow \infty$ .

**kNN on one-dimensional grid.** We consider the  $d=1$  dimensional case first. We assume  $\mathbf{x} = (1, 2, 3, \dots, n)$ , a circular metric  $d(x_i, x_j) = d(i, j) = \min\{|i-j|, n-|i-j|\}$ , and odd  $k \leq n$ . The kNN regression matrix

$$M_{ij} = b_{i-j} \quad \text{with} \quad b_{i-j} = \frac{1}{k} \quad \text{if} \quad d(i, j) \leq \frac{k-1}{2} \quad \text{and} \quad 0 \quad \text{otherwise}$$

is a diagonal-constant (Toeplitz) matrix with circularity property  $b_{i-j} = b_{i-j+n}$ . For instance, for  $k=3$  and  $n=5$

$$M = \frac{1}{3} \begin{pmatrix} 1 & 1 & 0 & 0 & 1 \\ 1 & 1 & 1 & 0 & 0 \\ 0 & 1 & 1 & 1 & 0 \\ 0 & 0 & 1 & 1 & 1 \\ 1 & 0 & 0 & 1 & 1 \end{pmatrix}$$

For every circulant matrix, the eigenvectors  $\mathbf{v}^1, \dots, \mathbf{v}^n$  are waves  $v_j^l = \theta^{jl}$  with  $\theta = e^{2\pi\sqrt{-1}/n}$ . The eigenvalues are the fourier transform  $\hat{b}_l = \sum_{j=1}^n b_j \theta^{-jl}$  of  $\mathbf{b}$ , since  $\sum_j M_{ij} v_j^l = \sum_j b_{i-j} \theta^{jl} = \sum_j b_j \theta^{(i-j)l} = v_i^l \sum_j b_j \theta^{-jl} = \hat{b}_l v_i^l$ , where we exploited circularity of  $\mathbf{b}$  and  $\theta^{jl}$ . For  $M_{\text{kNN}}$  in particular we get

$$\hat{b}_l = \frac{1}{k} \sum_{j=-\frac{k-1}{2}}^{\frac{k-1}{2}} \theta^{-jl} = \frac{1}{k} \frac{\theta^{lk/2} - \theta^{-lk/2}}{\theta^{l/2} - \theta^{-l/2}} = \frac{\sin(\pi lk/n)}{k \sin(\pi l/n)} < 1 \quad \text{for} \quad l \neq n$$

$\uparrow$  circularity                       $\uparrow$  geometric sum                       $\uparrow$  insert  $\theta$

and  $\hat{b}_n = 1$ . The only 1-vector  $\mathbf{v}^n = \mathbf{1}$  corresponds to a constant shift  $y_i \rightsquigarrow y_i + c$  under which kNN (like many other regressors) is invariant. Instead of regularizing LoRP with  $\alpha > 0$  we can restrict  $V(L) \subset \mathbb{R}^n$  to the space orthogonal to  $\mathbf{v}^n$ , which means dropping  $\hat{b}_n = 1$  in the determinant. Intuitively, since this invariant direction is the same for all  $k$ , we can drop the same additive infinite constant from LR for every  $k$ , which is irrelevant for comparisons (formally we should compute  $\lim_{\alpha \rightarrow 0} \{\text{LR}_{k_1}^\alpha - \text{LR}_{k_2}^\alpha\}$ ). The exact expression for the restricted log-determinant (denoted by a prime) is

$$-\frac{1}{2} \log \det' S_0 = -\log \det' (\mathbb{1} - M) = -\sum_{l=1}^{n-1} \log(1 - \hat{b}_l) =: \frac{n}{k} c_{nk}^1 = c_{nk}^1 \text{tr} M$$

For large  $n$  (and large  $k$ ) the expression can be simplified. The exact, large  $n$ , and

large  $k \ll n$  expressions are

$$\begin{aligned}
c_{nk}^1 &= -\frac{k}{n} \sum_{l=1}^{n-1} \log \left( 1 - \frac{\sin(\pi lk/n)}{k \sin(\pi l/n)} \right) \\
c_{\infty k}^1 &= -\frac{k}{\pi} \int_{-\pi/2}^{\pi/2} \log \left( 1 - \frac{\sin(kz)}{k \sin(z)} \right) dz && \begin{cases} z = \pi l/n \text{ for } l < \frac{n}{2} \\ z = \pi l/n - \pi \text{ else} \end{cases} \\
c_{\infty \infty}^1 &= -\frac{1}{\pi} \int_{-\infty}^{\infty} \log \left( 1 - \frac{\sin t}{t} \right) dt \doteq 3.202 && (t = kz, \sin(z) \sim z)
\end{aligned}$$

Further,  $c_{\infty 3}^1 = 3 \log 3 \doteq 3.295$ . Since  $c_{\infty k}^1$  is decreasing in  $k$ ,  $c_{\infty k}^1$  equals 3.2 within 3% for all  $k$ .

**kNN on  $d$ -dimensional grid.** We now consider  $\mathbf{x} = \mathcal{X}^d = \{1, \dots, n_1\}^d$  on a  $d$ -dimensional complete hypercube grid with  $n = n_1^d$  points and Manhattan distance  $d(x_i, x_j) = d(\mathbf{i}, \mathbf{j}) = \sum_{a=1}^d d_1(i_a, j_a)$  for all  $x_i = \mathbf{i} \in \mathcal{X}^d$  and  $x_j = \mathbf{j} \in \mathcal{X}^d$ , where  $d_1$  is the one-dimensional circular distance defined above (so actually  $\mathcal{X}^d$  is a discrete torus). For  $k = k_1^d$ , the neighborhood  $\mathcal{N}_k(x)$  of  $x$  is a cube of side-length  $k_1$ . In this case,  $M = M_1 \otimes \dots \otimes M_1$  is a  $d$ -fold tensor product of the 1d  $k_1$ NN matrices  $M_1$  of sample size  $n_1$ . The eigenvectors of  $M$  are  $\mathbf{v}^{l_1} \otimes \dots \otimes \mathbf{v}^{l_d}$  with eigenvalues  $\hat{b}_{l_1} \dots \hat{b}_{l_d}$ . We get

$$\begin{aligned}
-\log \det'(\mathbb{1} - M) &= -\sum_{l_1=1}^{n_1-1} \dots \sum_{l_d=1}^{n_d-1} \log(1 - \hat{b}_{l_1} \dots \hat{b}_{l_d}) && (30) \\
&\xrightarrow{n \gg k \rightarrow \infty} -\frac{1}{\pi^d} \int_{\mathbb{R}^d} \log \left( 1 - \prod_{a=1}^d \frac{\sin t_a}{t_a} \right) d^d \mathbf{t} =: \frac{n}{k} c_{\infty \infty}^d
\end{aligned}$$

For instance, for  $d=2$ , numerical integration gives  $c_{\infty \infty}^2 \doteq 2.2$  compared to 3.2 in one dimension. For higher dimensions, evaluation of the  $d$ -dimensional integral becomes cumbersome, and we resort to a different approximation.

**Taylor series in  $M$ .** We can also (not only for kNN) expand  $\log \det S_0$  in a Taylor series in  $M$ :

$$\begin{aligned}
-\log \det'(\mathbb{1} - M) &= -\text{tr}' \log(\mathbb{1} - M) = \sum_{s=1}^{\infty} \frac{1}{s} \text{tr}'(M^s) \\
&= \sum_{s=1}^{\infty} \frac{1}{s} (\text{tr}' M_1^s)^d = \frac{n}{k} \sum_{s=1}^{\infty} \frac{1}{s} (A_{n_1 k_1 s})^d =: \frac{n}{k} c_{nk}^d
\end{aligned}$$

where we used  $\text{tr}(A \otimes B) = \text{tr}(A) \cdot \text{tr}(B)$  and  $(A \otimes B)^s = A^s \otimes B^s$  and defined

$$A_{n_1 k_1 s} := \frac{k_1}{n_1} \text{tr}'(M_1^s) = \frac{k_1}{n_1} \sum_{l=1}^{n_1-1} (\hat{b}_l)^s \xrightarrow{n \gg k \rightarrow \infty} \frac{1}{\pi} \int_{-\infty}^{\infty} \left( \frac{\sin t}{t} \right)^s dt$$

The one-dimensional integral can be expressed as a finite sum with  $s$  terms or evaluated numerically. For any  $n$  and  $k$  one can show that  $A_{nk1} = A_{nk2} = 1 > A_{nks}$  for

$s > 2$ . So the expansion above is useful for large  $d$ . Note also that  $c_{nk}^d$  is monotone decreasing in  $d$ . For  $d \rightarrow \infty$  we have

$$c_{nk}^\infty = \sum_{s=1}^{\infty} \frac{1}{s} (A_{nks})^\infty = 1 + \frac{1}{2} + 0 + \dots = \frac{3}{2}$$

i.e.  $c_{nk}^d$  decreases monotone in  $d$  from about 3.2 to  $\frac{3}{2}$ .

The practical implication of this observation, though, is limited, since  $k = k_1^d \rightarrow \infty$  is actually not fixed for  $d \rightarrow \infty$ . Indeed, in practical high-dimensional problems,  $k \ll n \ll 3^d$ , but in our grid example  $k = k_1^d \geq 3^d$ . Real data do not form full grids but sparse neighborhoods if  $d$  is large.

## 11 Conclusion and Outlook

We introduced a new method, the Loss Rank Principle, for model selection. The loss rank of a model is defined as the number of other data that fit the model better than the training data. The model chosen by LoRP is the one of smallest loss rank. The loss rank has an explicit expression in case of linear models. Model consistency and asymptotic efficiency of LoRP were considered. The numerical experiments suggest that LoRP works well in practice. A comparison between LoRP and other methods for model selection was also presented.

In this paper, we have only scratched at the surface of LoRP. LoRP seems to be a promising principle with a lot of potential, leading to a rich field. In the following we briefly summarize miscellaneous considerations.

**Comparison to Rademacher complexities.** For a (binary) classification problem, the rank (1) of classifier  $r$  can be re-formulated as the probability that a randomly relabeled sample  $\mathbf{y}'$  behaves better than the actual  $\mathbf{y}$ . The more flexible  $r$  is, the larger its rank is. The Rademacher complexity [Kol01, BBL02] of  $r$  is the expectation of the difference between the misclassifying loss under the actual  $\mathbf{y}$  and the misclassifying loss under a randomly relabeled sample  $\mathbf{y}'$ . The more flexible  $r$  is, the larger its Rademacher complexity is. Therefore, there is a close connection between LoRP and Rademacher complexities. Model selection based on Rademacher complexities has a number of attractive properties and has been attracting many researchers, thus it's worth discovering this connection. Some results have been recently already obtained, however, to keep the present paper not so long, we decide to present the results in another paper.

**Monte Carlo estimates for non-linear LoRP.** For non-linear regression we did not present an efficient algorithm for the loss rank/volume  $\text{LR}_r(\mathbf{y}|\mathbf{x})$ . The high-dimensional volume  $|V_r(L)|$  (3) may be computed by Monte Carlo algorithms. Normally  $V_r(L)$  constitutes a small part of  $\mathcal{Y}^n$ , and uniform sampling over  $\mathcal{Y}^n$  is not feasible. Instead one should consider two competing regressors  $r$  and  $r'$  and compute  $|V \cap V'|/|V|$  and  $|V \cap V'|/|V'|$  by uniformly sampling from  $V$  and  $V'$  respectively e.g., with a Metropolis-type algorithm. Taking the ratio we get  $|V'|/|V|$  and hence the

loss rank difference  $LR_r - LR_{r'}$ , which is sufficient for LoRP. The usual tricks and problems with sampling apply here too.

**LoRP for hybrid model classes.** LoRP is not restricted to model classes indexed by a single integral “complexity” parameter, but may be applied more generally to selecting among some (typically discrete) class of models/regressors. For instance, the class could contain kNN *and* polynomial regressors, and LoRP selects the complexity *and* type of regressor (non-parametric kNN versus parametric polynomials).

**Generative versus discriminative LoRP.** We have concentrated on counting  $y$ 's given fixed  $x$ , which corresponds to discriminative learning. LoRP might equally well be used for counting  $(x, y)$ , as alluded to in the introduction. This would correspond to generative learning. Both regimes are used in practice. See [LJ08] for some recent results on their relative benefit, and further references.

**Acknowledgement.** We would like to thank two anonymous reviewers for their detailed and helpful comments. The second author would like to thank the SML@NICTA for supporting a visit which led to the present paper.

## Appendix: List of Abbreviations and Notations

AIC= Akaike Information Criterion.

BIC= Bayesian Information Criterion.

BMS= Bayesian Model Selection

kNN= k Nearest Neighbors.

LBFR= Linear Basis Function Regression.

LoRP= Loss Rank Principle.

LRC = Loss Rank Code.

MAP= Maximum a Posterior.

MDL= Minimum Description Length.

ML= Maximum Likelihood.

PML= Penalized Maximum Likelihood.

$D = \{(x_1, y_1), \dots, (x_n, y_n)\}$  = observed data.

$\mathcal{D} = \{D\}$  = set of all possible data  $D$ .

$\mathcal{X} \times \mathcal{Y}$  = observation space.

$\mathbf{x} = (x_1, \dots, x_n)$  = vector of  $x$ -observations, similarly  $\mathbf{y}$ .

$f: \mathcal{X} \rightarrow \mathcal{Y}$  = functional dependence between  $x$  and  $y$ .

$\mathcal{F}$  = (“small”) class of functions  $f$ .

$\mathcal{H}$  = class of stochastic hypotheses/models.

$r: \mathcal{D} \rightarrow \mathcal{F}$  = regressor/model.

$\hat{y}_i = r(x_i | D)$  =  $r$ -estimate of  $y_i$ .

$\mathcal{R}$  = (“small”) class of regressors/models.

$\mathbf{w} \in \mathbb{R}^d$  = parametrization of  $\mathcal{F}_d$ .

$\mathcal{N}_k(x)$  = set of indices of the  $k$  nearest neighbors of  $x$  in  $D$ .

$L = \text{Loss}_r(D) = \text{Loss}(\mathbf{y}, \hat{\mathbf{y}})$  = empirical loss of  $r$  on  $D$ .

$\text{Rank}_r(L) = \#\{\mathbf{y}' \in \mathcal{Y}^n : \text{Loss}_r(\mathbf{y}'|\mathbf{x}) \leq L\} =$  loss rank of  $r$ .

$V(L) =$  volume of  $D$  under  $r$ .

$\text{LR}_r(\mathbf{y}|\mathbf{x}) =$  log rank/volume of  $D$ .

$\text{LR}_r^\alpha =$  regularized  $\text{LR}_r$ .

$d_{\text{eff}} =$  effective dimension.

$m_j(x, \mathbf{x}) =$  coefficients of linear regressor.

$M(\mathbf{x}) =$  linear regression matrix or “hat” matrix.

$\log =$  natural logarithm.

$a \rightsquigarrow b:$   $a$  is replaced by  $b$ .

## References

- [Aka73] H. Akaike. Information theory and an extension of the maximum likelihood principle. In *Proc. 2nd International Symposium on Information Theory*, pages 267–281, Budapest, Hungary, 1973. Akademiai Kiadó.
- [All74] D. Allen. The relationship between variable selection and data augmentation and a method for prediction. *Technometrics*, 16:125–127, 1974.
- [BBL02] P. Bartlett, S. Boucheron, and G. Lugosi. Model selection and error estimation. *Machine Learning*, 48:85–113, 2002.
- [Cha06] A. Chambaz. Testing the order of a model. *Ann. Stat.*, 34(3):1166–1203, 2006.
- [CW79] P. Craven and G. Wahba. Smoothing noisy data with spline functions: estimating the correct degree of smoothing by the methods of generalized cross-validation. *Numerische Mathematik*, 31:377–403, 1979.
- [ET93] B. Efron and R. Tibshirani. *An Introduction to the Bootstrap*. Chapman & Hall/CRC, New York, 1993.
- [Grü04] P. D. Grünwald. Tutorial on minimum description length. In *Minimum Description Length: recent advances in theory and practice*, page Chapters 1 and 2. MIT Press, 2004. <http://www.cwi.nl/~pdg/ftp/mdlintro.pdf>.
- [Her02] R. Herbrich. *Learning Kernel Classifiers*. The MIT Press, 2002.
- [HT89] C. M. Hurvich and C. L. Tsai. Regression and time series model selection in small samples. *Biometrika*, 76(2):297–307, 1989.
- [HTF01] T. Hastie, R. Tibshirani, and J. H. Friedman. *The Elements of Statistical Learning*. Springer, 2001.
- [Hut07] M. Hutter. The loss rank principle for model selection. In *Proc. 20th Annual Conf. on Learning Theory (COLT'07)*, volume 4539 of *LNAI*, pages 589–603, San Diego, 2007. Springer, Berlin.
- [Kol01] V. Koltchinskii. Rademacher penalties and structural risk minimization. *IEEE Trans. Inform. Theory*, 47:1902–1914, 2001.

- [LJ08] P. Liang and M. Jordan. An asymptotic analysis of generative, discriminative, and pseudolikelihood estimators. In *Proc. 25th International Conf. on Machine Learning (ICML-2008)*, volume 307, pages 584–591. ACM, 2008.
- [Mac92] D. J. C. MacKay. Bayesian interpolation. *Neural Computation*, 4(3):415–447, 1992.
- [Mil02] A. Miller. *Subset Selection in Regression*. Chapman & Hall/CRC, 2002.
- [Reu02] A. Reusken. Approximation of the determinant of large sparse symmetric positive definite matrices. *SIAM Journal on Matrix Analysis and Applications*, 23(3):799–818, 2002.
- [Ris78] J. J. Rissanen. Modeling by shortest data description. *Automatica*, 14(5):465–471, 1978.
- [Sch78] G. Schwarz. Estimating the dimension of a model. *Annals of Statistics*, 6(2):461–464, 1978.
- [Sha97] J. Shao. An asymptotic theory for linear model selection. *Statistica Sinica*, 7:221–264, 1997.
- [Shi83] R. Shibata. Asymptotic mean efficiency of a selection of regression variables. *Annals of the Institute of Statistical Mathematics*, 35:415–423, 1983.
- [WLT07] H. Wang, R. Li, and C. L. Tsai. Tuning parameter selectors for the smoothly clipped absolute deviation method. *Biometrika*, 3(94):553–568, 2007.
- [Yam99] K. Yamanishi. Extended stochastic complexity and minimax relative loss analysis. In *In Proc. 10th International Conference on Algorithmic Learning Theory - ALT' 99*, pages 26–38. Springer-Verlag, 1999.
- [Yan05] Y. Yang. Can the strengths of aic and bic be shared? a conflict between model identification and regression estimation. *Biometrika*, 92(4):937–950, 2005.