

# NEURAL NETWORKS: SHALLOW LEARNING

Evgeny Burnaev

Skoltech, Moscow, Russia

# OUTLINE

- 1 APPROXIMATION PROBLEM AND BASIS EXPANSIONS
- 2 ADDITIVE MODELS AND NEURAL NETWORKS
- 3 SPECIFIC FEATURES OF THE ERM PROBLEM
- 4 RIDGE REGRESSION
- 5 HESSIAN APPROXIMATION

- 1 APPROXIMATION PROBLEM AND BASIS EXPANSIONS
- 2 ADDITIVE MODELS AND NEURAL NETWORKS
- 3 SPECIFIC FEATURES OF THE ERM PROBLEM
- 4 RIDGE REGRESSION
- 5 HESSIAN APPROXIMATION

## PROBLEM STATEMENT

- Let  $y = f(\mathbf{x})$  be some function, which is continuous and defined on a compact  $X \subset R^N$ ,  $N \sim 5 - 50$ .
- The problem is to construct an approximation  $\hat{f}(\mathbf{x})$  using the train sample of size  $m$

$$S_m = \{(\mathbf{x}_i, y_i), \mathbf{x}_i \in X, f(\mathbf{x}_i) = y_i, i = 1, \dots, m\}$$

- Approximation should be accurate in some sense

$$f(\mathbf{x}) \approx \hat{f}(\mathbf{x}), \mathbf{x} \in X \quad (1)$$

Note that (1) should hold for all  $\mathbf{x} \in X$ , not only for  $\mathbf{x} \in S_m$

## LINEAR EXPANSION IN A FUNCTIONAL DICTIONARY

- A model  $\hat{f}(\mathbf{x})$  is composed of functions from some parametric dictionary

$$\hat{f}(\mathbf{x}) = \sum_{j=1}^p \alpha_j \psi_j(\boldsymbol{\theta}_j, \mathbf{x}) + \alpha_0$$

Or, in vector notations

$$\hat{f}(\mathbf{x}) = \boldsymbol{\psi}(\Theta, \mathbf{x}) \boldsymbol{\alpha}, \quad \boldsymbol{\alpha} = \{\alpha_j\}_{j=0}^p, \quad \Theta = \{\boldsymbol{\theta}_j\}_{j=1}^p,$$

where  $\boldsymbol{\psi}(\Theta, \mathbf{x}) = (\psi_1(\boldsymbol{\theta}_1, \mathbf{x}), \dots, \psi_p(\boldsymbol{\theta}_p, \mathbf{x}))$

- Thus  $\hat{f}(\mathbf{x})$  is determined by
  - matrix  $\Theta$  of dictionary functions parameters
  - vector of coefficients  $\boldsymbol{\alpha}$  in the linear expansion

$\psi_j(\mathbf{x})$  can be considered as the  $j$ -th transformation of  $\mathbf{x}$

Once the dictionary  $\{\psi_j(\mathbf{x})\}_{j=1}^p$  is determined, the model  $\hat{f}(\mathbf{x})$  is linear in these new variables, and the fitting is similar to linear regression

# TYPES OF DICTIONARY FUNCTIONS

- Polynomials:  $\psi_j(\boldsymbol{\theta}_j, \mathbf{x}) = \prod_{k=1}^N x^{\theta_{j,k}}, \theta_j \in \{0, 1\}^N$
- Indicators for some regions:

$$\psi_j(\boldsymbol{\theta}_j, \mathbf{x}) = \prod_{k=1}^N 1\{\theta_{k,1} \leq x_k \leq \theta_{k,2}\}$$

- Sigmoid function:

$$\psi_j(\boldsymbol{\theta}_j, \mathbf{x}) = \sigma(\mathbf{x}^T \boldsymbol{\theta}_j^1 + \theta_j^0),$$

$$\text{where } \sigma(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} \text{ or } \sigma(z) = \frac{1}{1 + e^{-z}}$$

- Gaussian function:

$$\psi_j(\boldsymbol{\theta}_j, \mathbf{x}) = \exp\left(-\frac{\|\mathbf{x} - \boldsymbol{\theta}_j^1\|^2}{(\theta_j^0)^2}\right)$$

A dictionary can include other types of functions like trigonometric functions, etc.

## OBJECTIVE MEASURE: ERROR FUNCTION

- As mentioned above, the approximation should be close to the original function:

$$f(\mathbf{x}) \approx \hat{f}(\mathbf{x}), \mathbf{x} \in X$$

- Quantitative measure of closeness is the error function

$$\hat{R} = \hat{R}(S_m, \hat{f}) = \frac{1}{2} \sum_{i=1}^m \left( y_i - \hat{f}(\mathbf{x}_i) \right)^2, \quad S_m = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$$

In the sequel we denote by  $\#S$  a cardinality of a set  $S$  (number of points in the sample  $S$ )

## ALGORITHM

- 1 Select dictionary size  $p$
- 2 Initialize dictionary functions parameters  $\Theta$  and linear expansion coefficient  $\alpha$
- 3 Minimize the error function (Empirical Risk Minimization problem)  $\hat{R}(S_m, \hat{f}) = \hat{R}(\Theta, \alpha)$



## ALGORITHM

- 1 Select dictionary size  $p$
- 2 Initialize dictionary functions parameters  $\Theta$  and linear expansion coefficient  $\alpha$
- 3 Minimize the error function (Empirical Risk Minimization problem)  $\hat{R}(S_m, \hat{f}) = \hat{R}(\Theta, \alpha)$

## IN THIS PRESENTATION

- 1 Some connections to well-known models in statistics
- 2 Specific features of the ERM problem
- 3 Incorporate these features into the ERM algorithm to
  - Increase accuracy of approximation
  - Reduce training time

- 1 APPROXIMATION PROBLEM AND BASIS EXPANSIONS
- 2 ADDITIVE MODELS AND NEURAL NETWORKS
- 3 SPECIFIC FEATURES OF THE ERM PROBLEM
- 4 RIDGE REGRESSION
- 5 HESSIAN APPROXIMATION

## ADDITIVE MODELS I

- We try to fit a regression function

$$\hat{f}(\mathbf{x}) = \mathbb{E}(y|\mathbf{x}) = f(x_1, \dots, x_N),$$

in which every level of interaction is potentially present

- It is natural to consider analysis-of-variance (ANOVA) decompositions of the form

$$\hat{f}(x_1, \dots, x_N) = \alpha_0 + \sum_j g_j(x_j) + \sum_{k < r} g_{kr}(x_k, x_r) + \dots$$

- In order to restrict a model class we consider additive models, containing only main effect terms

$$\hat{f}(\mathbf{x}) = \alpha_0 + \sum_{j=1}^N g_j(\mathbf{x}_j)$$

## ADDITIVE MODELS II

Usually backfitting procedure is used:

1. Initialize  $g_1(x_1), \dots, g_N(x_N)$
2. For  $j = 1, \dots, N$ 
  - (A) Calculate residuals  $\epsilon_{j,i} = y_i - \sum_{s \neq j} g_s(x_{s,i})$
  - (B) Fit  $g_j(x_j)$  using the sample  $S_j = \{(x_{j,i}, \epsilon_{j,i})\}_{i=1}^m$
3. If converged STOP. Else, go back to step 2

# PROJECTION PURSUIT REGRESSION I

- Let  $\{\boldsymbol{\theta}_j\}_{j=1}^p$  be unit  $N$ -vectors of unknown parameters
- The projection pursuit regression model has the form

$$\hat{f}(\mathbf{x}) = \sum_{j=1}^p g_j(\boldsymbol{\theta}_j^T \mathbf{x})$$

- This is an additive model, but in the derived features  $z_j = \boldsymbol{\theta}_j^T \mathbf{x}$
- The functions  $g_j$  are unspecified and are estimated. Since

$$g(\boldsymbol{\theta}^T \mathbf{x}) \approx g(\boldsymbol{\theta}_{\text{old}}^T \mathbf{x}) + g'(\boldsymbol{\theta}_{\text{old}}^T \mathbf{x})(\boldsymbol{\theta} - \boldsymbol{\theta}_{\text{old}})^T \mathbf{x},$$

then

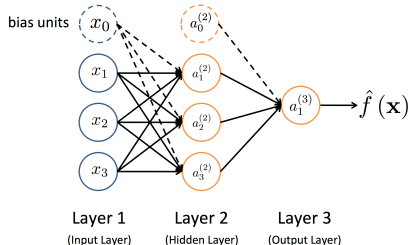
$$\sum_{i=1}^m [y_i - g(\boldsymbol{\theta}^T \mathbf{x}_i)]^2 \approx$$

$$\sum_{i=1}^m g'(\boldsymbol{\theta}_{\text{old}}^T \mathbf{x}_i)^2 \left[ \left( \boldsymbol{\theta}_{\text{old}}^T \mathbf{x}_i + \frac{y_i - g(\boldsymbol{\theta}_{\text{old}}^T \mathbf{x}_i)}{g'(\boldsymbol{\theta}_{\text{old}}^T \mathbf{x}_i)} \right) - \boldsymbol{\theta}^T \mathbf{x}_i \right]^2 \quad (2)$$

## PROJECTION PURSUIT REGRESSION II

- We fit the model by an iterative process:
  - A) Optimize (2) to update  $\theta$  (quadratic optimization!)
  - B) Tune  $g(\cdot)$  by smoothing current residuals
  - C) Repeat steps a)-b) until convergence
- Fit a new term  $g(\theta_{\text{new}}^T \mathbf{x})$  to the residuals, etc.

## NEURAL NETWORK WITH TWO-LAYERS I



- $a_i^{(j)}$  = “activation” of unit  $i$  in layer  $j$
- $\Theta^{(j)}$  = weight matrix controlling function mapping from layer  $j$  to layer  $j + 1$

$$a_1^{(2)}(\mathbf{x}) = \sigma \left( \theta_{10}^{(1)} + \theta_{11}^{(1)} x_1 + \theta_{12}^{(1)} x_2 + \theta_{13}^{(1)} x_3 \right) = \sigma((\boldsymbol{\theta}_1^{(1)})^T \mathbf{x})$$

$$a_2^{(2)}(\mathbf{x}) = \sigma \left( \theta_{20}^{(1)} + \theta_{21}^{(1)} x_1 + \theta_{22}^{(1)} x_2 + \theta_{23}^{(1)} x_3 \right) = \sigma((\boldsymbol{\theta}_2^{(1)})^T \mathbf{x})$$

$$a_3^{(2)}(\mathbf{x}) = \sigma \left( \theta_{30}^{(1)} + \theta_{31}^{(1)} x_1 + \theta_{32}^{(1)} x_2 + \theta_{33}^{(1)} x_3 \right) = \sigma((\boldsymbol{\theta}_3^{(1)})^T \mathbf{x})$$

$$\hat{f}(\mathbf{x}) = a_1^{(3)}(a_1^{(2)}, a_2^{(2)}, a_3^{(2)}) = \theta_{10}^{(2)} + \theta_{11}^{(2)} a_1^{(2)} + \theta_{12}^{(2)} a_2^{(2)} + \theta_{13}^{(2)} a_3^{(2)}$$

Here we assume that bias units  $a_0^{(2)} = 1$  and  $x_0 = 1$ , and that  $\sigma(\cdot)$  is a sigmoid function

## NEURAL NETWORK WITH TWO-LAYERS II

- NN with two layers and  $p$  hidden units coincides with a basis expansion in a dictionary with  $p$  basis functions in case

$$\psi(\boldsymbol{\theta}_j, \mathbf{x}) = \sigma(\boldsymbol{\theta}_j^T \mathbf{x})$$

and  $\alpha_j = \theta_{1,j}^{(2)}$ ,  $j = 0, 1, \dots, p$

- NN with two layers and  $p$  hidden units is similar to the projection pursuit regression when

$$g_j(\boldsymbol{\theta}_j^T \mathbf{x}) = \sigma(\boldsymbol{\theta}_j^T \mathbf{x}),$$

i.e. we consider pre-determined  $g_j(\cdot) = \sigma(\cdot)$



- 1 APPROXIMATION PROBLEM AND BASIS EXPANSIONS
- 2 ADDITIVE MODELS AND NEURAL NETWORKS
- 3 SPECIFIC FEATURES OF THE ERM PROBLEM**
- 4 RIDGE REGRESSION
- 5 HESSIAN APPROXIMATION

## LEARNING OF BASIS EXPANSIONS

Further

- We consider approximation models, defined by basis expansions in nonlinear parametric dictionaries

$$\hat{f}(\mathbf{x}) = \sum_{j=1}^p \alpha_j \psi_j(\boldsymbol{\theta}_j, \mathbf{x}) + \alpha_0$$

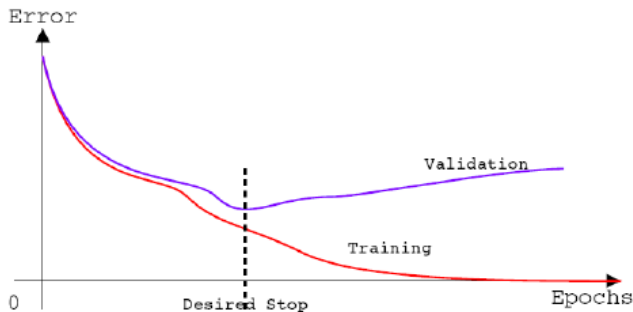
- We provide hints how to take a structure of the ERM problem into account in order to perform minimization efficiently

$$\hat{R}(S_m, \hat{f}) = \hat{R}(\Theta, \boldsymbol{\alpha}) = \frac{1}{2} \sum_{i=1}^m \left( y_i - \hat{f}(\mathbf{x}_i) \right)^2, \quad S_m = \{(\mathbf{x}_i, y_i)\}_{i=1}^m$$

# EARLY STOPPING I

Let  $S_{test}$  be some test sample, generated in  $X$ . Empirical indication of the overfitting is

$$\hat{R}(S_{test}, \hat{f}) \gg \hat{R}(S_m, \hat{f})$$



# EARLY STOPPING II

Let  $S_{test}$  be some test sample, generated in  $X$ . Empirical indication of the overfitting is

$$\hat{R}(S_{test}, \hat{f}) \gg \hat{R}(S_m, \hat{f})$$

Early stopping prevents overfitting:

## ADOPTED ALGORITHM

- ➊ Divide  $S_m$  into  $S_{train}$  and  $S_{val}$
- ➋ Select dictionary size  $p$
- ➌ Initialize dictionary functions parameters  $\Theta$  and linear expansion coefficients  $\alpha$
- ➍ Minimize the error function  $\hat{R}(S_{train}, \hat{f}) = \hat{R}(\Theta, \alpha)$ .
- ➎ Stop the optimization process when  $\hat{R}(S_{val}, \hat{f})$  stops to decrease

## SEPARABILITY OF VARIABLES

Let us consider the error function  $\hat{R}(\Theta, \alpha)$ :

$$\hat{R}(\Theta, \alpha) = (\Psi(\Theta) \alpha - \mathbf{y})^T (\Psi(\Theta) \alpha - \mathbf{y}),$$

where

- $\mathbf{y} = \{y_1, \dots, y_m\}$ ,
- $\Psi(\Theta) = (\psi(\Theta, \mathbf{x}_1), \dots, \psi(\Theta, \mathbf{x}_m))$ . Here and further  $m = \#(S_{train})$  is a size of the training set

We can easily see that

- Dependence of  $\hat{R}$  on  $\Theta$  is nonlinear
- Dependence of  $\hat{R}$  on  $\alpha$  is quadratic

We can find the optimal value of  $\alpha$  using least squares approach:

$$\alpha = \alpha(\Theta)$$

Possible explicit formulas for calculating optimal  $\alpha(\Theta)$

- ① Least squares estimate

$$\alpha(\Theta) = \left( \Psi(\Theta)^T \Psi(\Theta) \right)^{-1} \Psi(\Theta)^T \mathbf{y}$$

- ② Ridge regression

$$\alpha(\Theta) = \left( \Psi(\Theta)^T \Psi(\Theta) + \lambda I_p \right)^{-1} \Psi(\Theta)^T \mathbf{y}$$

- ③ Ridge regression with smoothness penalty:

$$\alpha(\Theta) = \left( \Psi(\Theta)^T \Psi(\Theta) + \lambda \Gamma(\Theta)^T \Gamma(\Theta) \right)^{-1} \Psi(\Theta)^T \mathbf{y},$$

where  $\Gamma(\Theta)$  is a some functional of derivatives of the approximation

## USING OF SEPARABILITY OF VARIABLES

- ❶ Calculate derivatives  $\hat{R}_\Theta$  and  $\hat{R}_{\Theta\Theta}$  of the error function
- ❷ Calculate step size  $\mathbf{d}_k$  by some optimization algorithm using gradient  $\hat{R}_\Theta$  and hessian  $\hat{R}_{\Theta\Theta}$
- ❸ Update the matrix of dictionary functions parameters:  
$$\Theta_{k+1} = \Theta_k + \mathbf{d}_k$$
- ❹ Update linear expansion coefficients  $\alpha(\Theta_{k+1})$ , applying e.g. LSQ formula

We should take into account dependence of  $\alpha(\Theta)$  on  $\Theta$  in step 1!

## SEPARABILITY OF VARIABLES

## CONSIDER A NEW OBJECTIVE FUNCTION

$$\hat{R}(\Theta) = \hat{R}(\Theta, \alpha(\Theta)), \quad \alpha(\Theta) = \left( \Psi(\Theta)^T \Psi(\Theta) \right)^{-1} \Psi(\Theta)^T \mathbf{y}$$

Let us calculate derivatives of this objective function

- By the definition of the least squares method

$$\hat{R}_\alpha = \mathbf{0}$$

- Gradient of  $\hat{R}(\Theta)$

$$\hat{R}_\Theta = \hat{R}_\Theta + \hat{R}_\alpha \alpha_\Theta = \hat{R}_\Theta + \mathbf{0} \alpha_\Theta = \hat{R}_\Theta$$

- Hessian

$$\hat{R}_{\Theta\Theta} = \hat{R}_{\Theta\Theta} + \hat{R}_{\Theta\alpha} \alpha_\Theta$$



- ① We can obtain  $\hat{R}_{\Theta\alpha}$  and  $\hat{R}_{\Theta\Theta}$  straightforwardly

$$\hat{R}_{\Theta\alpha} = (\mathbf{e}^T \Psi)_{\Theta} = \mathbf{e}^T \odot \Psi_{\Theta} + \mathbf{e}_{\Theta}^T \Psi = \mathbf{e}^T \odot \Psi_{\Theta} + J^T \Psi$$

$$\hat{R}_{\Theta\Theta} = J^T J + \mathbf{e} \odot \hat{f}_{\Theta\Theta}(\mathbf{X})$$

where

- $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$
- $J \stackrel{\text{def}}{=} \mathbf{e}_{\Theta} = \hat{f}_{\Theta}(\mathbf{X})$
- $\odot$  is pseudo-multiplication of 3D matrix:  
 $\mathbf{e}^T \odot \hat{f}_{\theta\theta}(\mathbf{X}) = \sum_{i=1}^m e_i \hat{f}_{\theta\theta}(\mathbf{x}_i)$

- ② Direct computation of  $\alpha_{\Theta}$  is rather lengthy

$$\hat{R}_{\alpha} = 0 \Rightarrow d\hat{R}_{\alpha} = \hat{R}_{\alpha\alpha}d\alpha + \hat{R}_{\alpha\Theta}d\Theta \equiv \mathbf{0},$$

$$\alpha_{\Theta} = \frac{d\alpha}{d\Theta} = -(\hat{R}_{\alpha\alpha})^{-1} \hat{R}_{\alpha\Theta},$$

$$\left( \text{and } \hat{R}_{\alpha\alpha} = \Psi^T \Psi \right)$$

## FINAL FORMULA FOR THE HESSIAN OF $\hat{R}$

- Hessian of  $\hat{R}$

$$\hat{R}_{\Theta\Theta} = \hat{R}_{\Theta\Theta} + \hat{R}_{\Theta\alpha}\alpha_{\Theta}$$

- After substitution:

$$\begin{aligned} \hat{R}_{\Theta\Theta} = J^T J + \mathbf{e} \odot \hat{f}_{\Theta\Theta}(\mathbf{X}) - \\ (\mathbf{e}^T \odot \Psi_{\Theta} + J^T \Psi) (\Psi^T \Psi)^{-1} (\mathbf{e}^T \odot \Psi_{\Theta} + J^T \Psi)^T \end{aligned}$$

- Let us neglect the terms with residuals  $\mathbf{e}$  ( $\mathbf{e} \approx \mathbf{0}$ )

$$\hat{R}_{\Theta\Theta} \approx J^T J - (J^T \Psi) (\Psi^T \Psi)^{-1} (J^T \Psi)^T$$

- 1 APPROXIMATION PROBLEM AND BASIS EXPANSIONS
- 2 ADDITIVE MODELS AND NEURAL NETWORKS
- 3 SPECIFIC FEATURES OF THE ERM PROBLEM
- 4 RIDGE REGRESSION**
- 5 HESSIAN APPROXIMATION

## LSQ ISN'T A GOOD IDEA!

- Matrix  $\Psi(\Theta)^T \Psi(\Theta)$  can be ill-conditioned
- LSQ estimates are unbiased but can have high variance
- In our case high variance  $\equiv$  unstable optimization process

Let us use a ridge regression to estimate linear expansion coefficients  $\alpha$

$$\alpha(\Theta) = \left( \Psi(\Theta)^T \Psi(\Theta) + \lambda I_p \right)^{-1} \Psi(\Theta)^T \mathbf{y}$$

Estimating  $\alpha$  with this formula is equivalent to minimization of

$$\tilde{R}(\Theta, \alpha) = \frac{1}{2} \left( (\mathbf{y} - \Psi\alpha)(\mathbf{y} - \Psi\alpha) + \lambda \alpha^T \alpha \right)$$

w.r.t.  $\alpha$

Let us apply the same approach to the new objective function

- We consider the modified error function  $\tilde{R}$

$$\tilde{R}(\Theta) = \tilde{R}(\Theta, \alpha(\Theta))$$

- Note that  $\tilde{R}_\alpha = \mathbf{0}$ , therefore we can use a general representation of the hessian matrix

$$\tilde{R}_{\Theta\Theta} = \tilde{R}_{\Theta\Theta} + \tilde{R}_{\Theta\alpha}\alpha_\Theta$$

- Partial derivatives with respect to  $\Theta$  are the same

$$\tilde{R}_{\alpha\Theta} = \hat{R}_{\alpha\Theta}, \quad \tilde{R}_{\Theta\Theta} = \hat{R}_{\Theta\Theta}$$

- A final formula

$$\tilde{R}_{\Theta\Theta} \approx J^T J - (J^T \Psi) (\Psi^T \Psi + \lambda I_p)^{-1} (J^T \Psi)^T$$

Let us compare two objective functions:

- 1 Initial error function:

$$\hat{R}(\Theta, \alpha) = \frac{1}{2} (\mathbf{y} - \Psi \alpha)^T (\mathbf{y} - \Psi \alpha)$$

- 2 Modified error function, based on explicit ridge regression estimate

$$\tilde{R}(\Theta) = \frac{1}{2} (\mathbf{y} - \Psi \alpha(\Theta))^T (\mathbf{y} - \Psi \alpha(\Theta)) + \frac{1}{2} \lambda \alpha(\Theta)^T \alpha(\Theta),$$

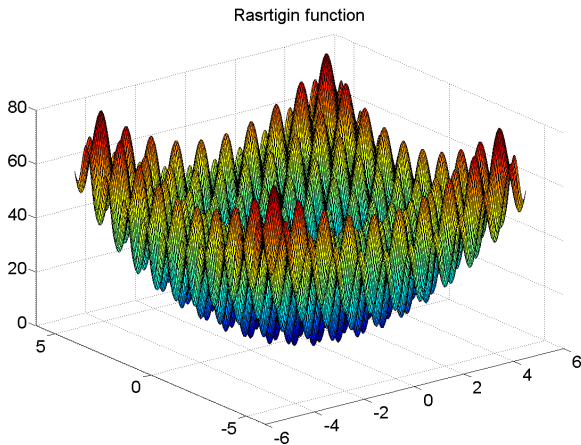
$$\alpha(\Theta) = \left( \Psi(\Theta)^T \Psi(\Theta) + \lambda \mathbf{I}_p \right)^{-1} \Psi(\Theta)^T \mathbf{y}$$

Also we assume that  $\mathbf{e} \approx \mathbf{0}$  in this case

We use the trust region method with a modified More-Sorensen method to solve a restricted quadratic minimization problem

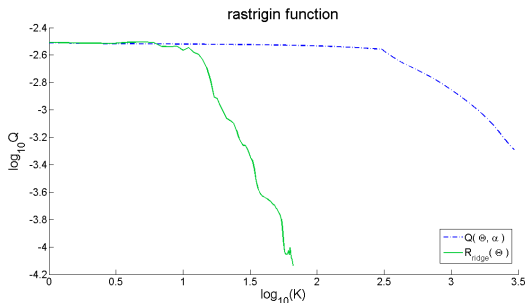
## TOY EXAMPLE: RASTRIGIN FUNCTION I

$$f(\mathbf{x}) = An + \sum_{i=1}^N (x_i^2 - A \cos(2\pi x_i)), \quad A = 10, \quad \mathbf{x} \in [-5.12, 5.12]^N$$



## TOY EXAMPLE: RASTRIGIN FUNCTION II

$m = 1000$  points, a dictionary consists of  $p = 128$  sigmoids

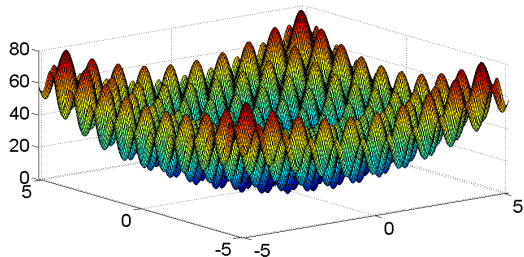


	$\hat{R}(\Theta, \alpha)$	$\tilde{R}(\Theta)$
$\hat{R}(S_{\text{test}}, \hat{f})$	2,99E-02	5,48E-05
<i>Iterations</i>	1532	87

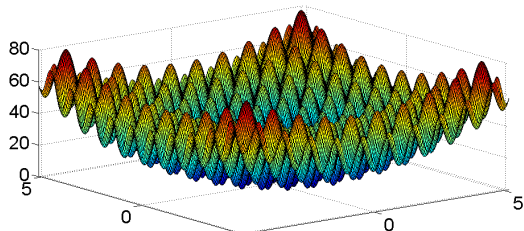


## TOY EXAMPLE: RASTRIGIN FUNCTION III

Approximation by sigmoid decomposition



Original function



## TOY EXAMPLE: RASTRIGIN FUNCTION IV

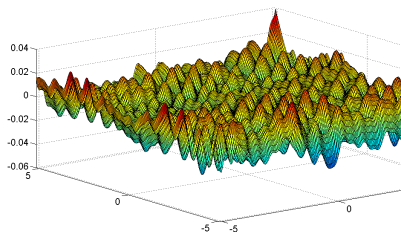


FIGURE : Residuals

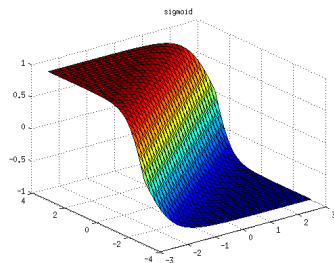


FIGURE : Sigmoid

Range of the original function is  $[0, 80]$ , range of residuals is  $[-0.06, 0.04]$ .

## AVERAGE RESULTS

- “Time” provides values of learning time for approximation construction
- We calculate a square normalized error on some test sample  $S_{test}$
- Ratio (for one task):

$$timeRatio(task) = \frac{\text{basic training algorithm}}{\text{improved training algorithms}}$$

- We consider 26 artificial tasks for each sample size  $m$

$m$	median timeRatio	median errorRatio
160	1.9	1.45
320	2.6	1.39
1000	3.2	1.54

- 1 APPROXIMATION PROBLEM AND BASIS EXPANSIONS
- 2 ADDITIVE MODELS AND NEURAL NETWORKS
- 3 SPECIFIC FEATURES OF THE ERM PROBLEM
- 4 RIDGE REGRESSION
- 5 HESSIAN APPROXIMATION

## HESSIAN ADDITIVITY

Let us consider the error function  $\tilde{R}(\Theta)$  and its hessian

$$\begin{aligned}\tilde{R}(\Theta) = \tilde{R}(\Theta, \alpha(\Theta)) &= \frac{1}{2} \left( (\mathbf{y} - \Psi \alpha(\Theta))^T (\mathbf{y} - \Psi \alpha(\Theta)) \right. \\ &\quad \left. + \lambda \alpha(\Theta)^T \alpha(\Theta) \right)\end{aligned}$$

$$\tilde{R}_{\Theta\Theta}(\Theta) \approx \mathcal{H} \stackrel{\text{def}}{=} J^T J - (J^T \Psi) (\Psi^T \Psi + \lambda I_p)^{-1} (J^T \Psi)^T,$$

$J = J(\Theta)$  are derivatives  $\hat{f}_{\Theta}(\mathbf{x})$  in sample points  $\mathbf{x} \in S_m$

Note, that  $\mathcal{H}$  is a sum of “sub-hessians” for training sample points:

$$\mathcal{H} = \sum_{i=1}^m J_i^T J_i - (J_i^T \Psi) (\Psi^T \Psi + \lambda I_p)^{-1} (J_i^T \Psi)^T = \sum_{i=1}^m h(\mathbf{x}_i, \Theta)$$

Computational complexity of

- $\mathcal{H}$  calculation is  $\sim m (pN)^2$ ,
- $\mathcal{H}$  inversion (or Cholesky decomposition) is  $\sim (pN)^3$ .

where

- $m$  is the training sample size,
- $N$  is a dimensionality of  $\mathbf{x}$ ,
- $p$  is a number of functions in the dictionary

If  $m \gg pN$  then computational complexity of the hessian calculation is significantly higher than of the hessian inversion

Let us define subsample  $S_{subtrain} \subset S_{train}$  with size  $\hat{m} \ll m$   
 Now we can calculate approximation of the true hessian  $\mathcal{H}$  using only points from  $S_{subtrain}$

$$\hat{\mathcal{H}} = \sum_{i=1}^{\hat{m}} h(\mathbf{x}_i, \Theta), \mathbf{x}_i \in S_{subtrain},$$

$$\tilde{R}_{\Theta, \Theta} \approx \mathcal{H} \approx \hat{\mathcal{H}}$$

How to select  $S_{subtrain}$  from  $S_{train}$ ?

- ① In order to approximate diagonal of the hessian with a maximal accuracy:

$$\text{diag}(\hat{\mathcal{H}}) \approx \text{diag}(\mathcal{H})$$

- ② Include points with maximal residuals in  $S_{subtrain}$
- ③ Uniformly randomly select  $S_{subtrain}$  (we use a new subsample for each iteration of the learning process)

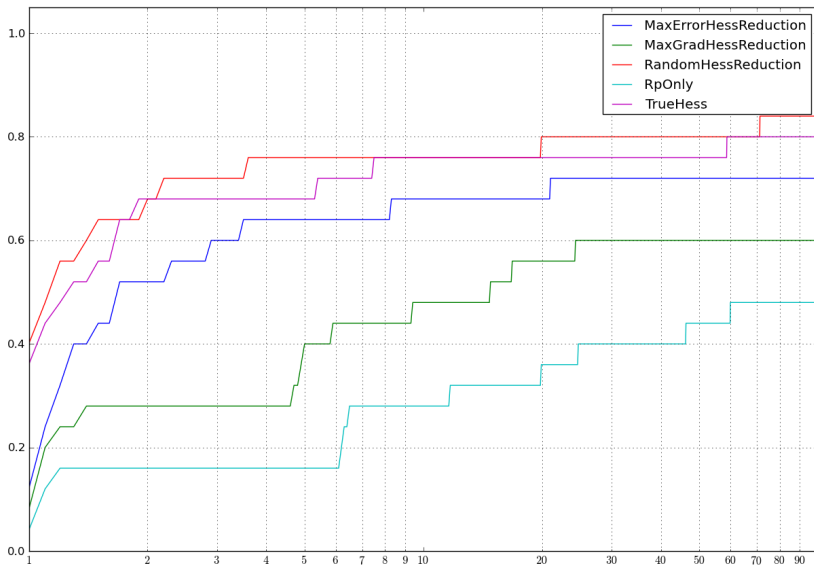
## DOLAN-MORE CURVES

- Let  $\{r_1, \dots, r_n\}$  be the set of compared methods,  $\{S_1, \dots, S_T\}$  be the set of tasks (datasets),  $q_{ti}$  be the quality of the method  $i$  on the dataset  $t$
- For each method  $i$  we introduce  $p_i(\beta)$ , a fraction of datasets, on which the method  $i$  is worse than the best one not more than  $\beta$  times:

$$p_i(\beta) = \frac{1}{T} \left| \left\{ t : q_{ti} \geq \frac{1}{\beta} \max_i q_{ti} \right\} \right|, \quad \beta \geq 1$$

- For example,  $p_i(1)$  is a fraction of datasets where the method  $i$  is the best
- A graph of  $p_i(\beta)$  is called Dolan-More curve for the method  $i$
- This definition implies that the higher the curve, the better the method





- Random approximation works better than others