

Sparse Conformal Predictor

Mohamed Hebiri
hebiri@math.jussieu.fr

LPMA
Paris 7-Diderot

Rennes - 29 Aout 2008

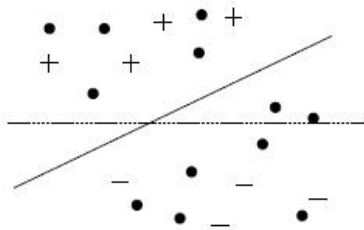
Outline

- 1 Model and Background
 - Transductive setting
 - Background
- 2 Sparse Conformal Predictors
 - Lasso Conformal Predictor
 - Family of conformal predictors
- 3 Simulations
 - Methods and comparison
 - Performances

Transductive setting

References:

- Vapnik '98
- Joachims '99



Linear Regression Model

Observations: $\mathcal{E}_n = \{(x_1, y_1), \dots, (x_n, y_n), x_{new}\}$

$$y_i = x_i \beta^* + \xi_i$$

- Covariates vector: $x_i = (x_{i,1}, \dots, x_{i,p}) \in \mathbb{R}^p$, $i \geq 1$
- New observation: $x_{new} \in \mathbb{R}^p$
- Response: $y_i \in \mathbb{R}$, $i \geq 1$
- Unknown parameter: $\beta^* = (\beta_1^*, \dots, \beta_p^*)' \in \mathbb{R}^p$
- Noise: $\xi_i \sim \mathcal{N}(0, \sigma^2)$

Objectives

Goal I: Given \mathcal{E}_n and $\varepsilon > 0$, provide a *confidence predictor* (interval) Γ^ε in which y_{new} lies with probability $1 - \varepsilon$

Tool: Conformity measure between x_{new} and the previous x_i 's

- distance (geometric, neighborhood,...)
- *distance of similarity*: to be defined later !

Goal II: Exploit the *sparsity* of the model (β^* contains many zero components) if necessary

Tool: Use a variable selection procedure (LASSO,...)

Conformal prediction: Vovk et al. '05

Notation:

- $y \in \mathbb{R}$: potential value of y_{new}
- $|\mathcal{A}|$: cardinal of the set \mathcal{A}

Nonconformity score $\alpha(y) = (\alpha_1(y), \dots, \alpha_n(y), \alpha_{new}(y))'$

- $\alpha_i(y)$: similarity between (x_{new}, y) and (x_i, y_i)
- relative information: *p-value*

$$p(y) = \frac{1}{n+1} |\{i \in \{1, \dots, n, new\} : \alpha_i(y) \leq \alpha_{new}(y)\}|$$

- $p(y) \in [\frac{1}{n+1}; 1]$
- the smaller $p(y)$ is, the less likely the tested pair (x_{new}, y) is (this choice reduces y to be an outlier)

Conformal predictor Γ^ε : labels $y \in \mathbb{R}$ such that $p(y) > \varepsilon$.

LASSO estimator: Tibshirani '96

LASSO estimator

$$\hat{\beta} = \underset{\beta \in \mathbb{R}^p}{\text{Argmin}} \frac{1}{n} \sum_{i=1}^n (y_i - x_i \beta)^2 + \lambda \sum_{j=1}^p |\beta_j|$$

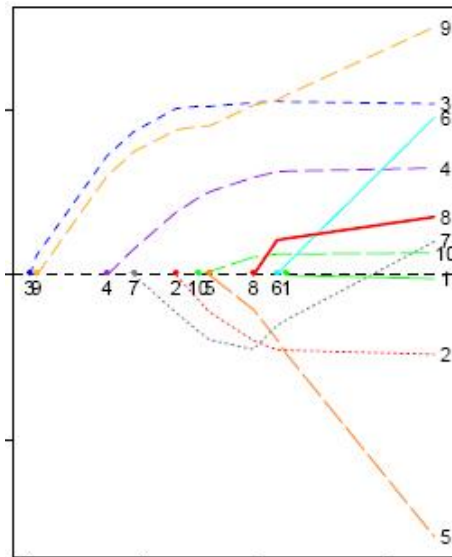
- Tuning parameter: λ

Motivation:

- Sparse solution $\hat{\beta}$ (i.e., many coefficients reduced to 0)
- Interpretable results whenever the model is sparse

Solution approximation: **LARS** algorithm (Efron et al. '04)

- Example: available Diabetes dataset (10 covariates)
- $\hat{\beta}_{\lambda_1}, \dots, \hat{\beta}_{\lambda_K}$: approximations of the LASSO solution at the transition points
 $\lambda = \lambda_1, \dots, \lambda_K$



- Step k : $\hat{\mu}_k = \mathbf{x}_k \hat{\beta}_{\lambda_k} = \mathbf{x}_k (\mathbf{x}'_k \mathbf{x}_k)^{-1} (\mathbf{x}'_k \mathbf{y} - \frac{\lambda_k}{2} \mathbf{s}_k)$
 - response vector: $\mathbf{y} = (y_1, \dots, y_n)'$
 - design matrix: $\mathbf{x} = (x'_1, \dots, x'_n)'$
 - sign vector: \mathbf{s}_k
 - \mathbf{x}_k is the restriction of \mathbf{x} to the columns corresponding to selected covariates

Does not take into account x_{new} while constructing $\hat{\beta}$!

Sparse Conformal Predictor

- Consider the augmented dataset: $\tilde{\mathbf{x}} = (x'_1, \dots, x'_n, x'_{new})'$ and $\tilde{\mathbf{y}} = (y_1, \dots, y_n, y)'$
- For **each** transition point λ_k , define the LASSO estimator $\hat{\mu}_k$ based on $\tilde{\mathbf{x}}_k$ and $\tilde{\mathbf{y}}$

Define the **Nonconformity score**

$$\alpha^k(y) := |\tilde{\mathbf{y}} - \hat{\mu}_k| = |\mathbf{A}_k + \mathbf{C}_k + \mathbf{B}_k y|$$

where $|\cdot|$ is meant here componentwise and

$$\left\{ \begin{array}{l} \mathbf{A}_k = (\mathbf{a}_1^k, \dots, \mathbf{a}_n^k, \mathbf{a}_{new}^k)' := (\mathbf{I} - \mathbf{H}_k) (y_1, \dots, y_n, 0)' \\ \mathbf{B}_k = (\mathbf{b}_1^k, \dots, \mathbf{b}_n^k, \mathbf{b}_{new}^k)' := (\mathbf{I} - \mathbf{H}_k) (0, \dots, 0, 1)' \\ \mathbf{C}_k = (\mathbf{c}_1^k, \dots, \mathbf{c}_n^k, \mathbf{c}_{new}^k)' := \frac{\lambda_k}{2} \tilde{\mathbf{x}}_k (\tilde{\mathbf{x}}_k' \tilde{\mathbf{x}}_k)^{-1} \mathbf{s}_k \end{array} \right.$$

- The $\alpha^k(y)$ are piecewise linear

Sparse Conformal Predictor

- *p-value*: $p^k(y) = \frac{1}{n+1} |\{i : \alpha_i^k(y) \leq \alpha_{new}^k(y)\}|$
- Predictor at step k : $\Gamma_k^\varepsilon = \{y \in \mathbb{R} : p^k(y) > \varepsilon\}$.

Proposition

Points y such that $\alpha_i^k(y) = \alpha_{new}^k(y)$ exist

i) if $b_i^k \neq b_{new}^k$: when y equals

$$-\frac{a_i^k - a_{new}^k + c_i^k - c_{new}^k}{b_i^k - b_{new}^k} \quad \text{and} \quad -\frac{a_i^k + a_{new}^k + c_i^k + c_{new}^k}{b_i^k + b_{new}^k}.$$

ii) if $b_i^k = b_{new}^k \neq 0$: when y equals

$$-\frac{a_i^k + a_{new}^k + c_i^k + c_{new}^k}{2b_i^k}$$

Conformal Lasso Predictor Γ_ε^F : the smallest Γ_ε^k

Extension

Estimator of the form:

$$\hat{\mu} = u(\tilde{\mathbf{x}}, \mathbf{s})\tilde{\mathbf{y}} + v(\tilde{\mathbf{x}}, \mathbf{s})$$

where $u(\cdot)$ and $v(\cdot)$ are piecewise constant functions w.r.t. $\tilde{\mathbf{y}}$

We are interested in

- CoLP: $u(\tilde{\mathbf{x}}, \mathbf{s}) = \tilde{\mathbf{x}}_k(\tilde{\mathbf{x}}_k'\tilde{\mathbf{x}}_k)^{-1}\tilde{\mathbf{x}}_k'$
 $v(\tilde{\mathbf{x}}, \mathbf{s}) = -\lambda_k\tilde{\mathbf{x}}_k(\tilde{\mathbf{x}}_k'\tilde{\mathbf{x}}_k)^{-1}\mathbf{s}_k$
- CoRP: $u(\tilde{\mathbf{x}}, \mathbf{s}) = \tilde{\mathbf{x}}(\tilde{\mathbf{x}}'\tilde{\mathbf{x}} + \mu\mathbf{I}_p)^{-1}\tilde{\mathbf{x}}'$ and $v = 0$
- CENeP: $u(\tilde{\mathbf{x}}, \mathbf{s}) = \tilde{\mathbf{x}}_k(\tilde{\mathbf{x}}_k'\tilde{\mathbf{x}}_k + \mu_k\mathbf{I}_k)^{-1}\tilde{\mathbf{x}}_k'$
 $v(\tilde{\mathbf{x}}, \mathbf{s}) = -\lambda_k\tilde{\mathbf{x}}_k(\tilde{\mathbf{x}}_k'\tilde{\mathbf{x}}_k)^{-1}\mathbf{s}_k$

Experimental setting

- All simulations are based on $M = 1000$ replications
- Performance measures:
 - Accuracy: length of intervals
 - Validity criterion: $\text{VAL}^\epsilon = M^{-1} \sum_{m=1}^M \mathbb{I}(y_{new} \in (\Gamma_\epsilon^F)_m)$
 - Variable selection: recovery of the support of β^*

Benchmarks:

- Variable selection: Original LASSO (Tibshirani '96) and the Original Elastic-Net (Zou & Hastie '05) (based on the BIC criterion)
- Accuracy and validity: CoRP (Vovk et al. '05)

Data $p = 50$

$\mathcal{A}^* = \{j : \beta_j^* \neq 0\}$: set of relevant covariates

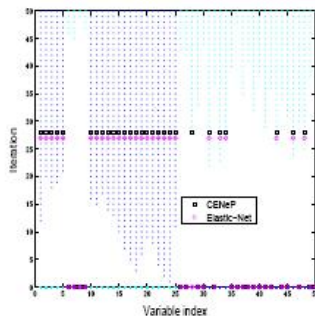
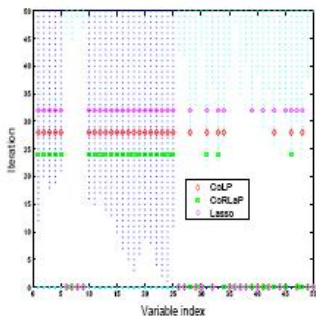
- Example(a): $\mathcal{A}^* = \{1\}$; exponentially decreasing correlations between successive covariates $\{15, \dots, 35\}$
- Example(b): $\mathcal{A}^* = \{1, \dots, 5\} \cup \{10, \dots, 25\}$; correlations as in Example(a)
- Example(c): $\mathcal{A}^* = \{1, \dots, 15\}$; three groups of highly correlated variables: $G_1 = \{1, \dots, 5\}$, $G_2 = \{6, \dots, 10\}$ and $G_3 = \{11, \dots, 15\}$
- Example(d): $\mathcal{A}^* = \{1, \dots, p\}$; exponentially decreasing correlations between successive covariates $\{1, \dots, p\}$

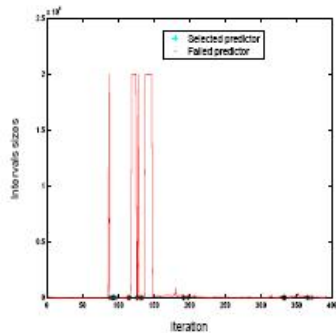
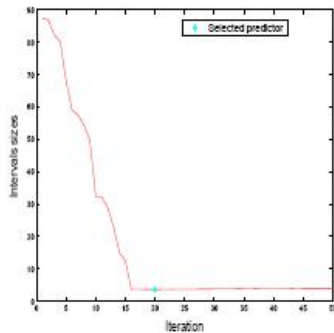
Validity

Table: VAL^ε evaluations

EXAMPLE[n/σ]	CoRP	CoLP	CoLARP	CENeP
EX (A)[300/1]	0.90± 0.02	0.88± 0.02	0.85± 0.02	0.88± 0.02
EX (A)[300/7]	0.89± 0.02	0.91± 0.02	0.89± 0.02	0.90± 0.02
EX (A)[300/15]	0.89± 0.02	0.89 ± 0.02	0.88± 0.02	0.89± 0.02
EX (B)[300/1]	0.90± 0.02	0.88± 0.02	0.87± 0.02	0.87± 0.02
EX (C)[300/1]	0.90± 0.02	0.90± 0.02	0.89± 0.02	0.90± 0.02
EX (D)[300/1]	0.89± 0.02	0.90± 0.02	0.90± 0.02	0.90± 0.02
EX (A)[50/3]	0.89± 0.02	0.67± 0.03	0.41± 0.03	0.79± 0.02
EX (A)[20/3]	0.86± 0.02	0.60± 0.03	0.30± 0.03	0.69± 0.03
EXAMPLE[n/σ]	CoRP	CoLP	STOPPED-CoLP	2-PN-CoLP
EX (A)[50/7]	0.85± 0.02	0.62± 0.03	0.82± 0.02	0.88± 0.02
EX (B)[50/1]	0.88± 0.02	0.56± 0.03	0.82± 0.02	0.91 ± 0.02
EX (C)[20/15]	0.88± 0.02	0.61± 0.03	0.77± 0.03	0.90± 0.02
EX (D)[20/1]	0.90± 0.02	0.60± 0.03	0.79± 0.02	0.89± 0.02

Variable selection: Example(b)[300/5]



Accuracy: Example(b)[$n/5$]

Conclusion

- Theoretical validity (Vovk et al. '05)
- Theoretical variable selection consistency when selection is based on the accuracy criterion.