

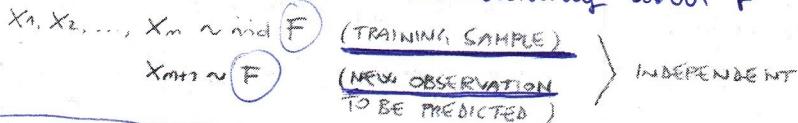
Task: prediction intervals for new observations

thoroughly related to permutation tests

for making probabilistic predictions (not relying on very precise specification of probabilistic models)

CONFORMAL PREDICTION (GAMMERER, VOLK, VAPNIK 1998)

we know abs. nothing about F



VALIDITY
 \downarrow
Coverage

$C(X_1, \dots, X_m)$ IS A VALID PREDICTION SET FOR X_{m+1} IF:

$$P(X_{m+1} \in C(X_1, \dots, X_m)) \geq 1-\alpha$$

We have 2 random objects: the new observation and the prediction set (which we consider as random as well)

COVERAGE / confidence level of the prediction set
 AND NOT THE UNIVARIATE PROBABILITY OF X_{m+1} !

Given a new dataset and a new prediction, the new dataset will generate a random set that will capture the new pred. with probability higher than $1-\alpha$.

Remember: the set is RANDOM

EFFICIENCY
 \downarrow
Pi LENGTH

$C_1(X_1, \dots, X_m)$ IS MORE EFFICIENT THAN $C_2(X_1, \dots, X_m)$ IF:
 or rule for building a pred. set another rule to build a pred. set

$$E(|C_1(X_1, \dots, X_m)|) < E(|C_2(X_1, \dots, X_m)|)$$

on average the length of the prediction intervals provided by the first rule is smaller than the length of the pred. int. provided by the second rule

A PARAMETRIC EXAMPLE: FISHER'S PREDICTION INTERVALS

$$X_1, \dots, X_m \sim \text{ind } N(\mu, \sigma^2) \quad \text{INDEPENDENT, } \mu \text{ AND } \sigma \text{ UNKNOWN}$$

$$\bar{X} = \frac{1}{m} \sum_{i=1}^m X_i \quad \text{POINT-WISE PREDICTION}$$

new observ.

$$\begin{aligned} X_{m+1} - \bar{X} &\sim \text{sample mean of the training samples} \\ S \sqrt{\frac{1}{m}} &\sim t(m-1) \end{aligned}$$

$$P\left(\bar{X} - S \sqrt{\frac{1}{m}} + t_{\alpha/2}(m-1) < X_{m+1} < \bar{X} + S \sqrt{\frac{1}{m}} + t_{\alpha/2}(m-1)\right) = 1-\alpha$$

This is a random set just depending on the training sample s.t. the new observation is included in the random set with prob. $1-\alpha$

formula that allows to build prediction intervals for new observations in a setting in which we know that date are gaussian

This can only be used if the date are gaussian (even if the sample size is large it needs to be gaussian)

For CI we're able to use the asymptotic CI, with PI we don't have asymptotic results

Our goal is to find a rule that provides valid and efficient prediction sets.

How the problem of constructing prediction intervals is tackled in the parametric framework?

in practice they cannot be computed
since we don't know x_{n+1}

Nonparametric counterpart of the Fisher prediction intervals:

CONFORMAL T PREDICTION INTERVALS

$x_1, \dots, x_m \sim \text{ind } F$ \rightarrow INDEPENDENT (CONTINUOUS F)
 $x_{m+1} \sim F$

LET US INDICATE WITH $\bar{x}_{-i} = \frac{1}{m} (x_1 + \dots + x_{i-1} + x_{i+1} + \dots + x_m)$
 $i = 1, \dots, m, m+1$

LET US DEFINE THE $m+1$ NON-CONFORMITY SCORES:

$$|x_i - \bar{x}_{-i}| \quad i = 1, \dots, m+1$$

This quantity measures in some sense the non-conformity of x_i w.r.t. all the other data. It tells how extreme it the i -th date w.r.t. the remaining data.

If this quantity is small then the i -th date is not strange w.r.t. the others.

$$x_1, \dots, x_m, x_{m+1} \sim \text{ind } F \Rightarrow |x_1 - \bar{x}_{-1}|, \dots, |x_{m+1} - \bar{x}_{-(m+1)}|,$$

EXCHANGEABLE

LET US DEFINE THE $m+1$ RANKS OF THE OBSERVATIONS ACCORDING TO THEIR CONFORMITY SCORES.

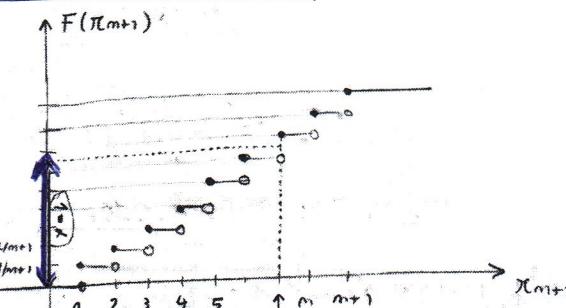
$$R_i = \sum_{j=1}^{m+1} \mathbb{I}(|x_j - \bar{x}_{-j}| \leq |x_i - \bar{x}_{-i}|) \quad i = 1, \dots, m+1$$

$$\text{EXCHANGEABILITY OF SCORES} \Rightarrow R_i \sim \text{UNIF}(\{1, \dots, m+1\}) \quad i = 1, \dots, m+1$$

IN PARTICULAR:

$$R_{m+1} \sim \text{UNIF}(\{1, \dots, m+1\})$$

we have a sequence of non conformity scores. All the sequences obtained by permuting the non conformity scores have exactly the same probability to be observed



$\lceil (1-\alpha)(m+1) \rceil = \text{superior } \alpha \text{ quantile}$

$$P(R_{m+1} \leq \lceil (1-\alpha)(m+1) \rceil) = \frac{\lceil (1-\alpha)(m+1) \rceil}{m+1} \geq 1-\alpha$$

The π_α , π_n
 F

cumulative dist. function of the ranks in variable that is the rank of the new observation

(discrete uniform on $\{1, \dots, n+1\}$)

this event (observed or not) depends only on x_1, \dots, x_{m+1} , not on F

TO BUILD A PREDICTION SET WE JUST NEED TO RESTATE THE EVENT $\{R_{m+1} \leq T(1-\alpha)(m+1)\}$ AS A FUNCTION OF x_1, \dots, x_m AND OF THE POSSIBLE FUTURE VALUES x_{m+1} OF X_{m+1} :

$$C(x_1, \dots, x_m) := \{x_{m+1} \in \mathbb{R} : R_{m+1} \leq T(1-\alpha)(m+1)\}$$

still random because it's depending on the training set

possible future realization of X_{m+1}

MORE EXPLICITLY INTRODUCING THE AUGMENTED SAMPLE $(x_1, \dots, x_m, x_{m+1})$:

$$P(x_{m+1}) := \frac{1}{m+1} \sum_{i=1}^{m+1} (|z_i - \bar{z}_i| \geq |x_{m+1} - \bar{x}_{-(m+1)}|) = \frac{(m+1) - S(m+1)}{m+1}$$

$$\rightarrow P\text{-VALUE OF } x_{m+1} \in \left\{ \frac{1}{m+1}, \frac{2}{m+1}, \dots, \frac{m}{m+1}, \frac{n+2}{n+2} \right\}$$

$$C(x_1, \dots, x_m) = \{x_{m+1} \in \mathbb{R} : p(x_{m+1}) \geq \alpha\}$$

proportion of data in the specific realization having a non-conformative score that is larger or equal than the one associated to the new observation

NB LINK WITH PERMUTATION TESTS FOR TWO INDEPENDENT SAMPLES (x_1, \dots, x_m) AND (x_{m+1})

MORE IN GENERAL NON-CONFORMITY SCORES CAN BE DEFINED THROUGH THE INTRODUCTION OF A NON-CONFORMITY MEASURE $A(\cdot, \cdot)$

$A(\{z_1, \dots, z_m\}, z_{m+1})$ MEASURING THE "NON-CONFORMITY" OF z_{m+1} TO THE SET $\{z_1, \dots, z_m\}$

FORMALLY $A: \mathbb{R}^m \times \mathbb{R} \rightarrow \mathbb{R}$ AND SYMMETRIC WRT THE FIRST m ARGUMENTS (i.e. INVARIANT TO PERMUTATIONS OF THE FIRST m ARGUMENTS)

NON-CONFORMITY SCORES ARE CONSISTENTLY DEFINED AS:

$$A(\{x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_{m+1}\}, x_i) \quad i = 1, \dots, m+1$$

e.g. IN T PREDICTION INTERVALS

$$A(\{z_1, \dots, z_m\}, z_{m+1}) := |z_{m+1} - \bar{z}_{-(m+1)}|$$

NB DIFFERENT NON-CONFORMITY MEASURES CAN BE TO TACKLE THE SAME PREDICTION PROBLEM.

Smart choices of non-conformity measures will allow to obtain smaller prediction sets preferring the same validity.

(We can define an optimal non-conformity measure, but we would need the distributions)

The nonconformity measure plays the role of the test statistic in permutation tests

A POPULAR WAY TO BUILD NON-COMFORMITY MEASURES IS
RELYING ON THE NOTION OF DISCREPANCY:

$$A((z_1, \dots, z_m), z_{m+1}) := |z_{m+1} - \hat{z}_{m+1}(z_1, \dots, z_m)|$$

in this way we're defining the scores as the PREDICTIVE RESIDUALS

i.e. THE PREDICTIVE RESIDUAL OF A NEW OBSERVATION FROM A PREDICTION OF IT BASED ON THE OTHER DATA

(NB) LINK WITH LEAVE-ONE-OUT CROSSVALIDATION.

(NB) WE HAVE VALID PREDICTION SET FOR ANY CHOICE OF \hat{z}_{m+1}

(NB) EFFICIENCY IS RELATED TO SMART CHOICES FOR \hat{z}_{m+1}

EX1 UNIVARIATE: $z_1, \dots, z_m, z_{m+1} \in \mathbb{R}$

$$A = |z_{m+1} - \bar{z}_{-(m+1)}| \quad (\text{I-INTERVALS})$$

[R-EXAMPLES:
FULL CP.R]

$$A = |z_{m+1} - \text{MED}(z_1, \dots, z_m)|$$

Sample mean (for gaussian data)
this non-conformity measure is the best

$$A = |\log(z_{m+1}) - \frac{1}{m} \log(z_1, \dots, z_m)|$$

useful if we suspect to have outliers
if we suspect that the E is not added but multiplied (!)

EX2 MULTIVARIATE: $z_1, \dots, z_m, z_{m+1} \in \mathbb{R}^p$

$$A = \|z_{m+1} - \bar{z}_{-(m+1)}\|^2$$

$$A = (z_{m+1} - \bar{z}_{-(m+1)}) S_{-(m+1)}^{-1} (z_{m+1} - \bar{z}_{-(m+1)})^\top$$

famous Mahalanobis' distance
of the points to the sample mean

EX3 REGRESSION

$z_1, \dots, z_m, z_{m+1} \in \mathbb{R}^2$ > WE EXPLORE POSSIBLE VALUES y_{m+1} OF y_{m+1}
FOR THE OBSERVED $(x_1, y_1), \dots, (x_m, y_m), x_{m+1}$

$(x_1, y_1) \quad \dots \quad (x_{m+1}, y_{m+1})$ } UNCONDITIONAL VALIDITY

$$A = |y_{m+1} - \hat{y}(x_{m+1})|$$

OLS
RIDGE
LASSO

SPLINES
KERNEL REGRESSION
...

NOT ALL NON-COMFORMITY MEASURES ARE BASED ON DISCREPANCY FROM A POINT-WISE PREDICTION. FOR EXAMPLE:

K-NEAREST NEIGHBOR

$$A(\{z_1, \dots, z_n\}, z_{n+1}) =$$

CONFIDENTIAL PREDICTION SET

$$\sum_{j: z_j \text{ in the first } K \text{ neighbors of } z_{n+1}} d^2(z_{n+1}, z_j)$$

(NB) It can be used in metric spaces
(link with distance based test statistic and to linkages)

ALGORITHMS

IN THE PRACTICE A TOLERANCE ϵ IS FIXED AND x_{m+1} IS EXPLORED ON A Q-SPADED GRID.

COMPUTATIONAL COSTS MAY BE LARGE IF:

- x_{m+1} BELONGS TO HIGH-DIMENSIONAL SPACE
- THE COMPUTATION OF \hat{x}_{m+1} IS COMPUTATIONALLY EXPENSIVE

SOLUTIONS

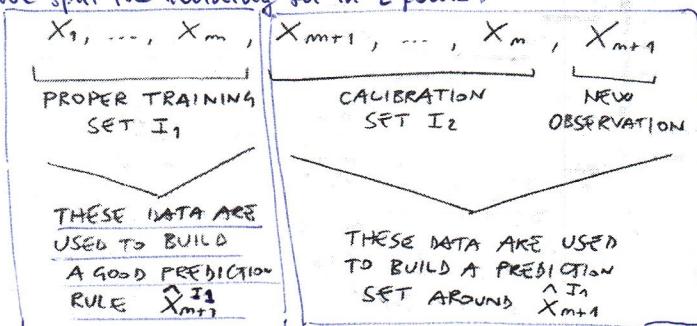
• PARALLELIZATION

• SMART NON-CONFORMITY MEASURES (\downarrow)

FROM FULL CONFORMAL PREDICTION TO SPLIT CONFORMAL PREDICTION

(PAPADOPOULOS, PROEIROU, VOURK, GATHERMAN 2002)

We split the training set in 2 parts:



conformal predictions will be out of this

we'll use a std conformal approach based on discrepancy on the calibration set (we'll

augment the calibration set with the new obs.

and we'll use a non conformal approach

on the calibr. set using a non-conform.

measure which is (kind of) degenerative w.r.t. the calibration set that define the nonconformity score

of one of the observations in the augmented calibration set as the mismatch with the prediction build from the first set.

THE $m-m+1$ NON-CONFORMITY SCORES ARE DEFINED AS:

$$A(\{x_{m+1}, \dots, x_m\}, \hat{x}_{m+1}) = |\hat{x}_{m+1} - \hat{x}_{m+1}^{I_1}|$$

AND THE $m-m+1$ RANKS AS:

$$R_i = \sum_{j=m+1}^{m+1} \mathbb{1}(|x_j - \hat{x}_{m+1}^{I_1}| \leq |x_i - \hat{x}_{m+1}^{I_1}|) \quad i = m+1, \dots, m+1$$

AND THE P-VALUE OF x_{m+1} AS:

$$P(x_{m+1}) = \frac{1}{m-m+1} \sum_{j=m+1}^{m+1} \mathbb{1}(|x_j - \hat{x}_{m+1}^{I_1}| \geq |x_{m+1} - \hat{x}_{m+1}^{I_1}|)$$

AND THE PREDICTION SET AS:

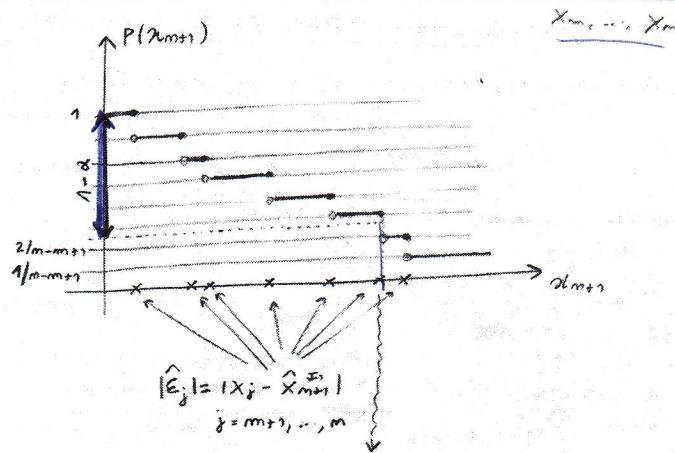
$$C(x_{m+1}, \dots, x_m) = \{x_{m+1} \in \mathbb{R} : P(x_{m+1}) \geq \alpha\}$$

degenerative
because it's not
taking into account
 $\{x_{m+1}, \dots, x_m\}$

ADVANTAGES

- $\hat{X}_{m+1}^{I_1}$ IS COMPUTED ONLY ONCE (COMPUTATIONALLY EXPENSIVE PREDICTION METHODS CAN BE USED)
- NON-CONFORMITY SCORES $|x_j - \hat{X}_{m+1}^{I_1}|$ $j = m+1, \dots, n$ ARE NOT AFFECTED BY THE VALUE OF X_{m+1} (i.e. SIMPLE RESIDUALS)

↳ $p(x_{m+1})$ CAN BE TRIVIALLY COMPUTED BY SORTING THE $m-m$ NON-CONFORMITY SCORES OF THE CALIBRATION SET



THE $\lceil (1-\alpha)(m-m+1) \rceil$ -TH ABSOLUTE RESIDUAL

SO THE PREDICTION SET IS SIMPLY DEFINED AS:

$$(\hat{X}_{m+1}^{I_1} - |\hat{e}_{\lceil (1-\alpha)(m-m+1) \rceil}|; \hat{X}_{m+1}^{I_1} + |\hat{e}_{\lceil (1-\alpha)(m-m+1) \rceil}|)$$

region in which the p-value is $\geq \alpha$

DISADVANTAGES

- POSSIBLE LOSS IN EFFICIENCY WRT FULL CONFORMAL
- EXTRA RANDOMNESS DUE TO THE RANDOM SPLIT

SPLIT PROPORTIONS

LARGER PROPER TRAINING SET → HIGHER QUALITY OF $\hat{X}_{m+1}^{I_1}$
 LARGER CALIBRATION SET → REDUCE THE IMPACT OF CEILING (less discretized p-value function)
 LITERATURE SUGGESTS TO CHOOSE $|I_1|$ AND $|I_2|$ IN BETWEEN

(66%, 33%) AND (75%, 25%) [R-EXAMPLE SPLIT CP.R]

trade-off
(to stabilize m
high or low)