Isotonic Distributional Regression (IDR) Leveraging Monotonicity, Uniquely So!

Tilmann Gneiting

Heidelberg Institute for Theoretical Studies (HITS) Karlsruhe Institute of Technology (KIT)

> Alexander Henzi Johanna F. Ziegel Universität Bern

> > MMMS2

June 2020



Heidelberg Institute for Theoretical Studies





▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQで

- 1 What is Regression?
- 2 Mathematical Background
 - 2.1 Calibration and Sharpness
 - 2.2 Proper Scoring Rules
 - 2.3 Partial Orders
- 3 Isotonic Distributional Regression (IDR)
 - 3.1 Definition, Existence, and Universality
 - 3.2 Computing
 - 3.3 Synthetic Example
- 4 Case Study on Precipitation Forecasts

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

- 1 What is Regression?
- 2 Mathematical Background
 - 2.1 Calibration and Sharpness
 - 2.2 Proper Scoring Rules
 - 2.3 Partial Orders
- 3 Isotonic Distributional Regression (IDR)
 - 3.1 Definition, Existence, and Universality
 - 3.2 Computing
 - 3.3 Synthetic Example
- 4 Case Study on Precipitation Forecasts

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Origins of Regression

regression originates from arguably the most notorious priority dispute in the history of mathematics and statistics



between Carl-Friedrich Gauss (1777–1855) and Adrien-Marie Legendre (1752–1833) over the method of least squares

Stigler (1981): "Gauss probably possessed the method well before Legendre, but [...] was unsuccessful in communicating it to his contemporaries"

Current Views: Distributional Regression

Wikipedia notes that

- "commonly, regression analysis estimates the conditional expectation [...] Less commonly, the focus is on a quantile [...] of the conditional distribution [...] In all cases, a function of the independent variables called the regression function is to be estimated"
- "it is also of interest to characterize the variation of the dependent variable around the prediction of the regression function using a probability distribution"
- Hothorn, Kneib and Bühlmann (2014) argue forcefully that the
 - "ultimate goal of regression analysis is to obtain information about the conditional distribution of a response given a set of explanatory variables"
- in a nutshell, distributional regression
 - uses training data

$$\{(x_i, y_i) \in \mathcal{X} \times \mathbb{R} : i = 1, \dots n\}$$

to estimate the conditional distribution of the response variable, $y \in \mathbb{R}$, given the explanatory variables or covariates, $x \in \mathcal{X}$

isotonic distributional regression (IDR) uses monotonicity relations to find nonparametric conditional distributions



bivariate point cloud — regression of Y on X



linear ordinary least squares (OLS; L₂) regression line



linear L_2 regression line with 80% prediction intervals

◆□ > ◆□ > ◆豆 > ◆豆 > ・豆



linear L_1 regression line — median regression



linear quantile regression — levels 0.10, 0.30, 0.50, 0.70, 0.90



linear quantile regression — zoom in



linear quantile regression — beware quantile crossing

▲□ > ▲圖 > ▲目 > ▲目 > ▲目 > ● ④ < ⊙



linear quantile regression



nonparametric isotonic mean (L_2) regression

(日) (四) (日) (日) (日)



nonparametric isotonic median (L_1) regression

(日) (四) (日) (日) (日)



nonparametric isotonic quantile regression



isotonic distributional regression (IDR)

Isotonic Distributional Regression (IDR) ... the Details

isotonic distributional regression (IDR) uses training data of the form

 $\{(x_i, y_i) \in \mathcal{X} \times \mathbb{R} : i = 1, \dots n\}$

to estimate a conditional distribution of the response variable or outcome, $y \in \mathbb{R}$, given the explanatory variables or covariates, $x \in \mathcal{X}$

takes advantage of known or assumed nonparametric monotonicity relations between the covariates, x, and the real-valued outcome, y

has primary uses in prediction and forecasting, where we know the covariates x, but do not know the outcome y

a full understanding relies on a number of (partly, rather recent) mathematical concepts and developments, namely,

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

- calibration and sharpness,
- proper scoring rules, and
- partial orders

- 1 What is Regression?
- 2 Mathematical Background
 - 2.1 Calibration and Sharpness
 - 2.2 Proper Scoring Rules
 - 2.3 Partial Orders
- 3 Isotonic Distributional Regression (IDR)
 - 3.1 Definition, Existence, and Universality
 - 3.2 Computing
 - 3.3 Synthetic Example
- 4 Case Study on Precipitation Forecasts

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

- 1 What is Regression?
- 2 Mathematical Background
 - 2.1 Calibration and Sharpness
 - 2.2 Proper Scoring Rules
 - 2.3 Partial Orders
- 3 Isotonic Distributional Regression (IDR)
 - 3.1 Definition, Existence, and Universality
 - 3.2 Computing
 - 3.3 Synthetic Example
- 4 Case Study on Precipitation Forecasts

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

What is the Goal in Distributional Regression?

the transition from classical regression to distributional regression poses unprecedented challenges, in that

- the regression functions are conditional predictive distributions in the form of probability measures or, equivalently, cumulative distribution functions (CDFs)
- the outcomes are real numbers
- so, in order to evaluate distributional regression techniques, we need to compare apples and oranges!

guiding principle: the goal is to maximize the sharpness of the conditional predictive distributions subject to calibration

- calibration refers to the statistical compatibility between the conditional predictive CDFs and the outcomes
 - essentially, the outcomes ought to be indistinguishable from random draws from the conditional predictive CDFs
- sharpness refers to the concentration of the conditional predictive distributions

the more concentrated the better, subject to calibration
<□><</p>
<□><</p>
><</p>

Probabilistic Framework

Setting We consider a probability space $(\Omega, \mathcal{A}, \mathbb{Q})$, where the members of the sample space Ω are tuples

$$(X, F_X, Y, V),$$

such that

- the random vector X takes values in the covariate space X (the explanatory variables or covariates),
- F_X is a CDF-valued random quantity that uses information based on X only (the conditional predictive distribution or regression function for Y, given X),
- ▶ Y is a real-valued random variable (the outcome), and
- V is uniformly distributed on the unit interval and independent of X and Y (a randomization device).

Definition The CDF-valued regression function F_X is ideal if $F_X = \mathcal{L}(Y \mid X)$ almost surely.

Notions of Calibration

Definition Let F_X be a CDF-valued regression function with probability integral transform (PIT)

$$Z = F_X(Y-) + V \left[F_X(Y) - F_X(Y-)\right].$$

Then F_X is

(a) probabilistically calibrated if Z is uniformly distributed,

(b) threshold calibrated if

 $\mathbb{Q}(Y \leq y | F_X(y)) = F_X(y)$ almost surely for all $y \in \mathbb{R}$.

Theorem An ideal regression function is both probabilistically calibrated and threshold calibrated.

Remark In practice, calibration is assessed by plotting PIT histograms

- U-shaped PIT histograms indicate underdispersed forecasts with prediction intervals that are too narrow on average
- skewed PIT histograms indicate biased predictive distributions

- 1 What is Regression?
- 2 Mathematical Background
 - 2.1 Calibration and Sharpness
 - 2.2 Proper Scoring Rules
 - 2.3 Partial Orders
- 3 Isotonic Distributional Regression (IDR)
 - 3.1 Definition, Existence, and Universality
 - 3.2 Computing
 - 3.3 Synthetic Example
- 4 Case Study on Precipitation Forecasts

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Scoring Rules

scoring rules seek to quantify predictive performance, assessing calibration and sharpness simultaneously

a scoring rule is a function

that assigns a negatively oriented numerical score to each pair (F, y), where F is a probability distribution, represented by its cumulative distribution function (CDF), and y is the real-valued outcome

a scoring rule S is proper if

 $\mathbb{E}_{Y \sim G}\left[\mathsf{S}(G,Y)\right] \, \leq \, \mathbb{E}_{Y \sim G}\left[\mathsf{S}(F,Y)\right] \quad \text{for all} \quad F,G,$

and strictly proper if, furthermore, equality implies F = G

truth serum: under a proper scoring rule truth telling is an optimal strategy in expectation

characterization results relate closely to convex analysis (Gneiting and Raftery 2007)

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

Continuous Ranked Probability Score (CRPS)

the widely used, proper continuous ranked probability score (CRPS) is defined as

$$CRPS(F, y) = \int_{-\infty}^{\infty} [F(x) - \mathbb{1}(x \ge y)]^2 dx$$
$$= \mathbb{E}_F |X - y| - \frac{1}{2} \mathbb{E}_F |X - X'|,$$

where X and X' are independent with CDF F

for all customary distributions, closed form expressions are available; e.g.,

$$\mathsf{CRPS}(\mathcal{N}(\mu,\sigma^2),y) = \sigma\left(\frac{y-\mu}{\sigma}\left(2\,\Phi\left(\frac{y-\mu}{\sigma}\right) - 1\right) + 2\,\phi\left(\frac{y-\mu}{\sigma}\right) - \frac{1}{\sqrt{\pi}}\right)$$

the CRPS is reported in the same unit as the outcomes, and it generalizes the absolute error, to which it reduces if F is a point measure reduces to the Brier score when the outcome is binary

Mixture (Choquet) Representations of the CRPS

the CRPS can be represented equivalently as

$$CRPS(F, y) = 2 \int_{(0,1)} QS_{\alpha}(F, y) d\lambda(\alpha)$$
$$= 2 \int_{(0,1)} \int_{\mathbb{R}} S^{Q}_{\alpha,\theta}(F, y) d\lambda(\theta, \alpha)$$
$$= \int_{\mathbb{R}} \int_{(0,1)} S^{P}_{z,c}(F, y) d\lambda(c, z)$$

in terms of the asymmetric piecewise linear loss QS_{α} , or the elementary or extremal scoring functions $S_{\alpha,\theta}^{Q}$ for the α -quantile functional, or $S_{z,c}^{P}$ for probability assessments of the binary outcome $\mathbb{1}(y \leq z)$, namely

$$\mathsf{QS}_{\alpha}(F, y) = \begin{cases} (1-\alpha) \left(F^{-1}(\alpha) - y\right), & y \leq F^{-1}(\alpha), \\ \alpha \left(y - F^{-1}(\alpha)\right), & y \geq F^{-1}(\alpha), \end{cases}$$

 $\mathsf{S}^{Q}_{\alpha,\theta}(F,y) = \begin{cases} 1-\alpha, & y \leq \theta < F^{-1}(\alpha), \\ \alpha, & F^{-1}(\alpha) \leq \theta < y, \\ 0, & \text{otherwise}, \end{cases} \qquad \mathsf{S}^{P}_{z,c}(F,y) = \begin{cases} 1-c, & F(z) < c, \; y \leq z, \\ c, & F(z) \geq c, \; y > z, \\ 0, & \text{otherwise}, \end{cases}$

respectively (Ehm et al. 2016)

- 1 What is Regression?
- 2 Mathematical Background
 - 2.1 Calibration and Sharpness
 - 2.2 Proper Scoring Rules
 - 2.3 Partial Orders
- 3 Isotonic Distributional Regression (IDR)
 - 3.1 Definition, Existence, and Universality
 - 3.2 Computing
 - 3.3 Synthetic Example
- 4 Case Study on Precipitation Forecasts

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Partial Orders

a partial order relation \preceq on a general set ${\mathcal X}$

has the same properties as a total order, namely reflexivity, antisymmetry and transitivity

▶ except that the elements need not be comparable, i.e., there might be elements $x \in \mathcal{X}$ and $x' \in \mathcal{X}$ such that neither $x \leq x'$ nor $x' \leq x$

• a key example is the componentwise order on \mathbb{R}^d

of particular importance in our context are partial orders on the set \mathcal{P} of the Borel probability measures on \mathbb{R} , which we identify with their respective CDFs

- ▶ stochastic order (\leq_{st}) $G \leq_{st} H$ if, and only if, $G(y) \geq H(y)$ for $y \in \mathbb{R}$
- increasing convex order (\leq_{icx}) $G \leq_{icx} H$ if, and only if,

$$\mathbb{E}[\phi(X_G)] \leq \mathbb{E}[\phi(X_H)]$$

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

whenever ϕ is increasing and convex and the expectations exist

Partial Orders on \mathbb{R}^d

in our case study $\mathcal{X} = \mathbb{R}^d$, and we consider the

► componentwise order (≤)

$$x \preceq x' \iff x_i \le x_i'$$
 for $i = 1, \dots, d$

- ► empirical stochastic order (≤st) induced by the stochastic order on the associated empirical distributions, and equivalent to the componentwise order on the sorted elements
- ► empirical increasing convex order (≤_{icx}) induced by the increasing convex order on the associated empirical distributions



◆□▶ ◆□▶ ◆ 臣▶ ◆ 臣▶ 三臣 - ∽ � � �

- 1 What is Regression?
- 2 Mathematical Background
 - 2.1 Calibration and Sharpness
 - 2.2 Proper Scoring Rules
 - 2.3 Partial Orders
- 3 Isotonic Distributional Regression (IDR)
 - 3.1 Definition, Existence, and Universality
 - 3.2 Computing
 - 3.3 Synthetic Example
- 4 Case Study on Precipitation Forecasts

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

- 1 What is Regression?
- 2 Mathematical Background
 - 2.1 Calibration and Sharpness
 - 2.2 Proper Scoring Rules
 - 2.3 Partial Orders
- 3 Isotonic Distributional Regression (IDR)
 - 3.1 Definition, Existence, and Universality
 - 3.2 Computing
 - 3.3 Synthetic Example
- 4 Case Study on Precipitation Forecasts

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Isotonic Distributional Regression (IDR): Basic Concepts

basic concepts

we use training data

$$\{(x_i, y_i) \in \mathcal{X} \times \mathbb{R} : i = 1, \dots, n\}$$

to estimate the conditional distribution of the response variable or outcome, $y \in \mathbb{R}$, given the explanatory variables or covariates, $x \in \mathcal{X}$

- ▶ formally, distributional regression generates a mapping from a covariate vector x ∈ X to a probability measure F_x, which serves to model the conditional distribution of the outcome, y, given x
- \blacktriangleright given a partial order \preceq on the covariate space $\mathcal X,$ this mapping is isotonic if

$$x \preceq x' \Rightarrow F_x \leq_{\mathrm{st}} F_{x'},$$

where \leq_{st} denotes the usual stochastic order on the space $\mathcal P$ of the Borel probability measures in $\mathbb R$

IDR: Definition, Existence and Uniqueness

formal setting

- covariate space \mathcal{X} equipped with partial order \preceq
- training data $\{(x_i, y_i) \in \mathcal{X} \times \mathbb{R} : i = 1, \dots n\}$
- ▶ the stochastic order \leq_{st} on the space \mathcal{P} of the Borel probability measures on \mathbb{R}
- proper scoring rule S

Definition (isotonic S-regression) An element $\hat{F} = (\hat{F}_1, \dots, \hat{F}_n) \in \mathcal{P}^n$ is an isotonic S-regression if it is a minimizer of the empirical loss

$$\ell_{\mathsf{S}}(\boldsymbol{F}) = \frac{1}{n} \sum_{i=1}^{n} \mathsf{S}(F_i, y_i)$$

over all $\mathbf{F} = (F_1, \dots, F_n) \in \mathcal{P}^n$, subject to the condition that $F_i \leq_{st} F_j$ if $x_i \leq x_j$, for $i, j = 1, \dots, n$.

Theorem (existence and uniqueness) There exists a unique isotonic CRPS-regression $\hat{F} \in \mathcal{P}^n$.

Terminology We refer to this unique \hat{F} as the isotonic distributional regression (IDR) solution.

Isotonic Distributional Regression (IDR): Universality

Theorem (universality) The IDR solution \hat{F} is threshold calibrated, and it is an isotonic S-regression under just any scoring rule of the form

$$\mathsf{S}(F,y) = \int_{(0,1)\times\mathbb{R}} \mathsf{S}^{\mathsf{Q}}_{\alpha,\theta}(F,y) \, \mathsf{d}H(\alpha,\theta)$$

or

$$\mathsf{S}(F,y) = \int_{\mathbb{R}\times(0,1)} \mathsf{S}_{z,c}^{P}(F,y) \, \mathrm{d}M(z,c),$$

where $S_{\alpha,\theta}^Q$ and $S_{z,c}^P$ are the elementary quantile and probability scoring functions, and H and M are locally finite Borel measures.

Proof relies on results and techniques in Ehm et al. (2016) and Jordan et al. (2019)

Consequence (theoretical) IDR is optimal under just any proper scoring rule that depends on quantile or binary probability assessments only.

Consequence (practical) IDR subsumes extant approaches to nonparametric isotonic regression as special cases, including but not limited to quantile regression and binary regression.

- 1 What is Regression?
- 2 Mathematical Background
 - 2.1 Calibration and Sharpness
 - 2.2 Proper Scoring Rules
 - 2.3 Partial Orders
- 3 Isotonic Distributional Regression (IDR)
 - 3.1 Definition, Existence, and Universality
 - 3.2 Computing
 - 3.3 Synthetic Example
- 4 Case Study on Precipitation Forecasts

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Estimation

the IDR solution exists and, by definition, is the solution to a constrained optimization problem in \mathcal{P}^n ... but can we actually compute it?

yes — universality and the method of least squares come to the rescue!





by universality (M = δ_z ⊗ λ₁), the IDR solution F̂ satisfies

$$\hat{oldsymbol{F}}(z) = rg \min_{\eta \in [0,1]^n} \sum_{i=1}^n \left(\eta_i - \mathbbm{1}(y_i \leq z)
ight)^2,$$

at every threshold $z \in \mathbb{R}$, subject to the condition that $\eta_i \ge \eta_j$ if $x_i \preceq x_j$, for i, j = 1, ..., n

 at any fixed threshold, the IDR CDFs yield a quadratic programming problem, which we tackle with the OSQP solver (Stellato et al. 2017)

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

- the target function is constant for z inbetween the unique values of y₁,..., y_n, and so it suffices to consider these points only
- the overall computational cost is at least $\mathcal{O}(n^2)$

Prediction

by construction, the IDR solution $\hat{F} = (\hat{F}_1, \dots, \hat{F}_n)$ is defined at the training covariate values $x_1, \dots, x_n \in \mathcal{X}$ only

a key task in practice is to make a prediction at a new covariate value $x \in \mathcal{X}$ where $x \notin \{x_1, \ldots, x_n\}$, for which we proceed as follows

define the sets p(x) and s(x) of the indices of immediate predecessors and successors of x among x₁,..., x_n as

$$p(x) = \{i \in \{1, \dots, n\} : x_i \leq x_j \leq x \implies x_j = x_i, j = 1, \dots, n\}$$

$$s(x) = \{i \in \{1, \dots, n\} : x \leq x_j \leq x_i \implies x_j = x_i, j = 1, \dots, n\},$$

• any predictive CDF F that is consistent with \hat{F} must satisfy

$$\max_{i \in s(x)} \hat{F}_i(z) \le F(z) \le \min_{j \in p(x)} \hat{F}_j(z)$$

at all threshold values $z \in \mathbb{R}$

if both p(x) and s(x) are nonempty, we let F be the pointwise arithmetic average of these bounds, i.e.,

$$F(z) = \frac{1}{2} \left(\max_{i \in s(x)} \hat{F}_i(z) + \min_{j \in p(x)} \hat{F}_j(z) \right)$$

- 1 What is Regression?
- 2 Mathematical Background
 - 2.1 Calibration and Sharpness
 - 2.2 Proper Scoring Rules
 - 2.3 Partial Orders
- 3 Isotonic Distributional Regression (IDR)
 - 3.1 Definition, Existence, and Universality
 - 3.2 Computing
 - 3.3 Synthetic Example
- 4 Case Study on Precipitation Forecasts

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Synthetic Example

we compute the IDR solution based on a training sample of size n = 600 from a population where $X \sim \text{Unif}_{(0,10)}$ and

$$Y \mid X \sim \mathsf{Gamma}(\mathsf{shape} = \sqrt{X}, \, \mathsf{scale} = \mathsf{min}\{\mathsf{max}\{X,1\},6\})$$



200

э

Synthetic Example: Subset Aggregation

same setting as before, but now for a training sample of size n = 10000



linear aggregation of IDR estimates on 100 subsamples of size 1 000 each (subagging, panel (b)) is superior to using the full training sample (panel (a)) in terms of both computational costs and estimation accuracy

- 1 What is Regression?
- 2 Mathematical Background
 - 2.1 Calibration and Sharpness
 - 2.2 Proper Scoring Rules
 - 2.3 Partial Orders
- 3 Isotonic Distributional Regression (IDR)
 - 3.1 Definition, Existence, and Universality
 - 3.2 Computing
 - 3.3 Synthetic Example
- 4 Case Study on Precipitation Forecasts

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Numerical Weather Prediction (NWP)

modern weather forecasts rely on numerical weather prediction (NWP) models that represent physical processes in the atmosphere





run operationally on supercomputers, with huge success

nevertheless, major sources of uncertainty remain (initial conditions, representation of sub-grid scale processes, ...)

ensemble simulations seek to quantify uncertainty and provide distributional forecasts

despite continuous improvement, NWP ensemble forecasts remain subject to systematic deficiencies

 $\tt https://celebrating200 years.noaa.gov/breakthroughs/climate_model/AtmosphericModelSchematic.png$

ECMWF Ensemble System

the 52-member ensemble system operated by the European Centre for Medium-Range Weather Forecasting (ECMWF) comprises

- a high-resolution member (x_{hres}) at 9 km horizontal grid spacing
- a control member (x_{ctr}) at 18 km horizontal grid spacing
- ▶ 50 perturbed members (x₁,..., x₅₀) at the same lower resolution but with perturbed initial conditions, to be considered exchangeable



systematic deficiencies call for postprocessing of the raw ensemble output via distributional regression, with covariate vector

$$x = (x_{\text{hres}}, x_{\text{ctr}}, x_1, \dots, x_{50})$$

Case Study: Precipitation Forecasts

our weather data comprise

- 52-member ECMWF ensemble forecasts and associated observations of 24-hour accumulated precipitation
- at prediction horizons of 1 to 5 days ahead
- from 6 January 2007 to 1 January 2017
- > at weather stations on airports in London, Brussels, Zurich and Frankfurt
- ▶ precipitation is a particularly challenging variable, due to its nonnegativity and mixed discrete-continuous character with a point mass at zero and a right skewed component on (0,∞)

we perform an out-of-sample evaluation and comparison of distributional regression forecasts

- years 2015 and 2016 as test period
- prior years serve to provide training data
- generally, IDR uses all available training data, whereas parametric competitors benefit from smaller, rolling training periods

Out-of-sample Comparison of Predictive Performance

systematic deficiencies call for postprocessing of the raw ensemble output via distributional regression, with covariate vector

 $x = (x_{\mathsf{hres}}, x_{\mathsf{ctr}}, x_1, \dots, x_{50})$

we compare IDR to the raw ensemble and state-of-the-art distributional regression techniques developed specifically for the purpose

- ENS ECMWF raw ensemble forecast, i.e., the empirical distribution of the 52 ensemble members
- BMA Bayesian Model Averaging (Sloughter et al. 2007)
 - semi-parametric, based on mixtures of Bernoulli and powertransformed Gamma components
 - plenty of implementation decisions to be made
- EMOS Ensemble Model Output Statistics (Scheuerer 2014)
 - parametric, predictive CDFs from the three-parameter family of left-censored generalized extreme value (GEV) distributions
 - Iocation and scale parameters linked to covariates, numerous implementation decisions to be made

Choice of Partial Order for IDR

IDR applies readily in this setting

without any need for adaptations due to the mixed discrete-continuous character of precipitation, nor requiring data transformations

however, the partial order on the elements $x = (x_{hres}, x_{ctr}, x_1, \dots, x_{50})$ of the covariate space $\mathcal{X} = \mathbb{R}^{52}$ needs to be selected thoughtfully

• considering that the elements of $x_{ptb} = (x_1, \dots, x_{50})$ are exchangeable

we apply IDR in three variants

▶ IDR_{cw} based on x_{hres} , x_{ctr} and $m_{ptb} = \frac{1}{50} \sum_{i=1}^{50} x_i$ and the componentwise order on \mathbb{R}^3 , so that

$$x \preceq x' \iff x_{\mathsf{hres}} \le x'_{\mathsf{hres}}, \, x_{\mathsf{ctr}} \le x'_{\mathsf{ctr}}, \, m_{\mathsf{ptb}} \le m'_{\mathsf{ptb}},$$

- ▶ IDR_{sbg} same as IDR_{cw} but combined with subset aggregation
- ▶ IDR_{icx} invokes the empirical increasing convex order on x_{ptb}, so that

$$x \preceq x' \iff x_{\mathsf{hres}} \le x'_{\mathsf{hres}}, \; x_{\mathsf{ptb}} \preceq_{\mathsf{icx}} x'_{\mathsf{ptb}}$$

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

Example: Predictive CDFs for Brussels, 16 December 2015



prediction horizon: two days

◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ● ○ ○ ○ ○

Calibration Assessed by PIT Histograms



▲□▶ ▲圖▶ ▲園▶ ▲園▶ 三国 - 釣ん(で)

CRPS



◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ▶ ○ □ ○ ○ ○ ○

Brier Score



シック 単 (中本) (中本) (日)

- 1 What is Regression?
- 2 Mathematical Background
 - 2.1 Calibration and Sharpness
 - 2.2 Proper Scoring Rules
 - 2.3 Partial Orders
- 3 Isotonic Distributional Regression (IDR)
 - 3.1 Definition, Existence, and Universality
 - 3.2 Computing
 - 3.3 Synthetic Example
- 4 Case Study on Precipitation Forecasts

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Summary

in regression analysis

- we are witnessing a transition from conditional mean estimation to conditional distribution estimation
- prompted and accompanied by a transition from point forecasts to distributional or probabilistic forecasts (Gneiting and Katzfuss 2014)

isotonic distributional regression (IDR) is a powerful nonparametric technique for estimating conditional distributions under order restrictions

- IDR learns conditional distributions that are calibrated, and simultaneously optimal relative to comprehensive classes of proper scoring rules
- IDR provides a unified treatment of all types of real-valued outcomes
- IDR is entirely generic and fully automated
- code for the implementation of IDR in R is available online, with functions for partial orders, estimation, prediction and evaluation

https://github.com/AlexanderHenzi/isodistrreg

Discussion

IDR might serve as an ideal benchmark technique in distributional regression and probabilistic forecasting problems

- method is entirely generic
- does not require potentially subjective implementation decisions, except for the choice of a partial order
- shows strongly competitive predictive performance in challenging and important applications

deep thinking vs. deep learning?

- IDR requires the a priori selection of a partial order
 - at least for now, this process cannot be automated
 - requires deep thinking about the substantive problem at hand
 - once the partial order has been fixed, IDR is fully automated
- nonparametric distributional regression techniques based on modern neural networks such as CNNs or RNNs (e.g., SQF-RNN, Gasthaus et al. 2019) are attractive alternatives
- ► partly overlapping though largely complementary uses

Selected References

Gneiting, T., Raftery, A. E. (2007), **Strictly proper scoring rules, prediction,** and estimation, *Journal of the American Statistical Association*, 102, 359–378.

Gneiting, T., Katzfuss, M. (2014), **Probabilistic forecasting**, *Annual Review of Statistics and Its Application*, 1, 125–151.

Jordan, A. I., Mühlemann, A., Ziegel, J. F. (2019), **Optimal solutions to the isotonic regression problem**, preprint, https://arxiv.org/abs/1904.04761.

・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・
 ・

Henzi, A., Ziegel, J. F., Gneiting, T. (2019), Isotonic distributional regression, preprint, https://arxiv.org/abs/1909.03725.