Proper scoring rules

Johanna Ziegel

ETH Zurich

CUSO Summer School 7–10 September 2025

Proper scoring rules

Johanna Ziegel

Introduction

Detini

Motivatio

Definition

Divergence and entrop

Examples

Estimation

Classes of scoring

Characterization

ocal scoring rules

ernel scores

Forecast

nformation sets

Decompositions

Summary

Introduction: Probabilistic predictions

- Let $Y \in \mathcal{Y}$ be an unknown future outcome.
 - ▶ Temperature tomorrow at 12:00 in Cambridge. $(Y \in \mathcal{Y} = \mathbb{R})$
 - \blacktriangleright Event of rain tomorrow in London. $(Y \in \mathcal{Y} = \{0,1\})$
 - ▶ Default of credit card client. $(Y \in \mathcal{Y} = \{0, 1\})$
 - \blacktriangleright Amount of precipitation tomorrow in Cambridge and Oxford. ($Y \in \mathcal{Y} = \mathbb{R}^2$)
- ▶ Single valued "best guess" $z \in \mathcal{Y}$ does not quantify uncertainty.

Proper scoring rules

Johanna Ziegel

Introduction

Introduction: Probabilistic predictions

- ▶ Let $Y \in \mathcal{Y}$ be an unknown future outcome.
 - ▶ Temperature tomorrow at 12:00 in Cambridge. $(Y \in \mathcal{Y} = \mathbb{R})$
 - ightharpoonup Event of rain tomorrow in London. $(Y \in \mathcal{Y} = \{0,1\})$
 - ▶ Default of credit card client. $(Y \in \mathcal{Y} = \{0,1\})$
 - lacktriangle Amount of precipitation tomorrow in Cambridge and Oxford. $(Y \in \mathcal{Y} = \mathbb{R}^2)$
- ▶ Single valued "best guess" $z \in \mathcal{Y}$ does not quantify uncertainty.
- ▶ Better: Quantify uncertainty of *Y* by a *probabilistic prediction F*.
 - ightharpoonup F is a distribution on \mathcal{Y} .

Proper scoring rules

Johanna Ziegel

Introduction

Definition

Motivation

Divergence :

Examples

Estimation

Classes of scoring rules

Characterization

Local scoring rules f-scores and gf-scores

rnel scores

Forecast comparison

Information sets

Decompositions

.

Summary

Introduction: Probabilistic predictions

- ▶ Let $Y \in \mathcal{Y}$ be an unknown future outcome.
 - ▶ Temperature tomorrow at 12:00 in Cambridge. $(Y \in \mathcal{Y} = \mathbb{R})$
 - ightharpoonup Event of rain tomorrow in London. $(Y \in \mathcal{Y} = \{0,1\})$
 - ▶ Default of credit card client. $(Y \in \mathcal{Y} = \{0,1\})$
 - lacktriangle Amount of precipitation tomorrow in Cambridge and Oxford. $(Y \in \mathcal{Y} = \mathbb{R}^2)$
- ▶ Single valued "best guess" $z \in \mathcal{Y}$ does not quantify uncertainty.
- ▶ Better: Quantify uncertainty of *Y* by a *probabilistic prediction F*.
 - ightharpoonup F is a distribution on \mathcal{Y} .
- ▶ If X is information available for prediction, F should approximate $\mathcal{L}(Y \mid X)$.

Proper scoring rules

Johanna Ziegel

Introduction

Definition

Motivation

Divergence and en

Examples

Estimation

Classes of scoring rules

haracterization

ocal scoring rules scores and gf-scores

ecast

omparison

formation sets ecompositions

Summary

▶ Temperature tomorrow at 12:00 in Cambridge. ($Y \in \mathcal{Y} = \mathbb{R}$)

ightharpoonup Event of rain tomorrow in London. $(Y \in \mathcal{Y} = \{0,1\})$

▶ Default of credit card client. $(Y \in \mathcal{Y} = \{0,1\})$

lacktriangle Amount of precipitation tomorrow in Cambridge and Oxford. ($Y \in \mathcal{Y} = \mathbb{R}^2$)

▶ Single valued "best guess" $z \in \mathcal{Y}$ does not quantify uncertainty.

Better: Quantify uncertainty of Y by a probabilistic prediction F.

ightharpoonup F is a distribution on \mathcal{Y} .

▶ If X is information available for prediction, F should approximate $\mathcal{L}(Y \mid X)$.

▶ Other possibilities to quantify uncertainty of *Y*: prediction intervals, predictions of some measure of variability, . . .

Proper scoring rules

Johanna Ziegel

Introduction

Definiti

Motivation

Divergence and e

stimation

Classes of scoring rules

Characterization

Local scoring rules

-scores and gf-scores

orecast omparison

Information sets
Decompositions

Lillitation

Summary

▶ Temperature tomorrow at 12:00 in Cambridge. $(Y \in \mathcal{Y} = \mathbb{R})$

lacktriangle Event of rain tomorrow in London. $(Y \in \mathcal{Y} = \{0,1\})$

▶ Default of credit card client. $(Y \in \mathcal{Y} = \{0,1\})$

lacktriangle Amount of precipitation tomorrow in Cambridge and Oxford. $(Y \in \mathcal{Y} = \mathbb{R}^2)$

▶ Single valued "best guess" $z \in \mathcal{Y}$ does not quantify uncertainty.

Better: Quantify uncertainty of Y by a probabilistic prediction F.

ightharpoonup F is a distribution on \mathcal{Y} .

▶ If X is information available for prediction, F should approximate $\mathcal{L}(Y \mid X)$.

▶ Other possibilities to quantify uncertainty of *Y*: prediction intervals, predictions of some measure of variability, . . .

▶ Which loss functions can we use to compare probabilistic predictions?

Proper scoring rules.

Proper scoring rules

Johanna Ziegel

Introduction

Definition

Motivation

Divergence and

Divergence and

Estimation

classes of scoring ules

Characterization
Local scoring rules

cores and gf-scores

orecast omparison

Information sets
Decompositions

Limitatio

Summary

Outline Introduction Definition Motivation Definition Divergence and entropy Examples Estimation Classes of scoring rules Characterization

Local scoring rules *f*-scores and *gf*-scores

Kernel scores
Forecast comparison
Information sets
Decompositions

Limitations Summary

Proper scoring

rules

Johanna Ziegel

Introduction

Proper scoring rules

Johanna Ziegel

Introduction

Definition

Motivation

.

Divergence and entro

xamples

Estimation

Classes of scoring

haracterization

Local scoring rules f-scores and gf-scores

rnel scores

orecast

nformation sets

Decompositions

Summary

References

Definition

Motivation

Definition

Divergence and entropy

Examples

Motivation 1: Forecast comparison

Let \mathcal{P} be a class of distributions on \mathcal{Y} .

Pick a loss function $S: \mathcal{P} \times \mathcal{Y} \to \overline{\mathbb{R}}$, and compare average realized scores: For $(F_1, G_1, Y_1), \ldots, (F_n, G_n, Y_n)$, let

$$\hat{S}_1 = \frac{1}{n} \sum_{i=1}^n S(F_i, Y_i)$$
 and $\hat{S}_2 = \frac{1}{n} \sum_{i=1}^n S(G_i, Y_i)$.

The forecast with the smaller value \hat{S}_i is better.

- ▶ When does this procedure make sense?
- ► *S* needs to be a proper scoring rule.
- ▶ History: In meteorology, Brier (1950) showed that quadratic score for binary outcomes cannot be gamed.

Proper scoring rules

Johanna Ziegel

Introduction

Detinitio

Motivation

Divergence a

Estimation

LStillation

rules

Local scoring rules

f-scores and gf-scores

f-scores and gf-scores Kernel scores

orecast

nformation sets

imitation

ummary

ummary

Motivation 2: McCarthy's forecasting tournament

- Forecasting agent makes probabilistic prediction F for a random variable Y and receives the penalty S(F, Y).
- ▶ Rationally, if the agent believes $Y \sim G$ then it issues

$$\underset{F}{\operatorname{argmin}} \ \mathbb{E}_{G}[S(F,Y)]$$

 $\mathbb{E}_G S(F, Y)$ means: take expectation with $Y \sim G$.

- ▶ Not necessarily equal to *G*.
- ▶ McCarthy's idea. Choose S such that

$$G = \underset{F}{\operatorname{argmin}} \ \mathbb{E}_G[S(F, Y)]$$

- McCarthy characterized such (differentiable) S for $Y \in \{0,1\}$ (McCarthy, 1956).
- lacktriangle Bregman divergences... ~ 10 years before Bregman introduced them in 1967.

Proper scoring rules

Johanna Ziegel

Introduction

Definition

Motivation

Divergence ar

stimation

lasses of scoring

Characterization

ocal scoring rules scores and gf-scores

nel scores

orecast omparison

formation sets

mitations

imitations

Summary

Let \mathcal{P} be a convex class of distributions on \mathcal{Y} .

Definition

A scoring rule is a function $S: \mathcal{P} \times \mathcal{Y} \to \mathbb{R} \cup \{\pm \infty\}$ that is suitably integrable. A scoring rule S is proper if

$$\mathbb{E}_{F}S(F,Y) \leq \mathbb{E}_{F}S(G,Y), \quad F,G \in \mathcal{P}, Y \sim F. \tag{1}$$

S is *strictly* proper if equality implies F = G.

Proper scoring rules

Johanna Ziegel

Introduction

Definition

Motivation

Definition

Divergence and entropy

Estimation

lasses of scoring

Characterization
Local scoring rules
f-scores and gf-scores

orecast omparison

nformation sets Decompositions

imitations

Summary

Let \mathcal{P} be a convex class of distributions on \mathcal{Y} .

Definition

A scoring rule is a function $S: \mathcal{P} \times \mathcal{Y} \to \mathbb{R} \cup \{\pm \infty\}$ that is suitably integrable. A scoring rule S is proper if

$$\mathbb{E}_{F}S(F,Y) \leq \mathbb{E}_{F}S(G,Y), \quad F,G \in \mathcal{P}, Y \sim F. \tag{1}$$

S is strictly proper if equality implies F = G.

Equivalent to (1) is

$$F \in \operatorname*{argmin}_{G} \mathbb{E}_{F} S(G, Y) = \operatorname*{argmin}_{G} S(G, F).$$

Proper scoring rules

Johanna Ziegel

Introduction

Motivation

Definition

Divergence and entropy

Estimation

llasses of scoring ules

Local scoring rules
f-scores and gf-scores

Forecast

nformation sets

imitations

Summary

Let \mathcal{P} be a convex class of distributions on \mathcal{Y} .

Definition

A scoring rule is a function $S: \mathcal{P} \times \mathcal{Y} \to \mathbb{R} \cup \{\pm \infty\}$ that is suitably integrable. A scoring rule S is proper if

$$\mathbb{E}_{F}S(F,Y) \leq \mathbb{E}_{F}S(G,Y), \quad F,G \in \mathcal{P}, Y \sim F. \tag{1}$$

S is strictly proper if equality implies F = G.

Equivalent to (1) is

$$F \in \underset{G}{\operatorname{argmin}} \mathbb{E}_F S(G, Y) = \underset{G}{\operatorname{argmin}} S(G, F).$$

- Scoring rules are interpreted as penalties.
- ► Forecasts should be compared with proper scoring rules (Gneiting and Raftery, 2007).
- ▶ Proper scoring rules are also increasingly important in estimation (Dawid et al., 2016).
- ► New review article (Waghmare and Ziegel, 2025).

Proper scoring rules

Johanna Ziegel

ntroduction

definition Motivation

Definition
Divergence and entro

Estimation

llasses of scoring ules

Characterization
Local scoring rules
f-scores and gf-scores

omparison
Information sets

Decompositions

iiiitations

Summary

Examples for
$$\mathcal{Y} = \{0, 1\}$$

Proper scoring rules Johanna Ziegel

Definition

Distributions F on \mathcal{Y} can be identified with a parameter $p \in [0, 1]$.

Brier score

$$S(p,y)=(y-p)^2, \quad p\in[0,1], \ y\in\{0,1\},$$

Logarithmic score

$$S(p,y) = -y \log(p) - (1-y) \log(1-p), \quad p \in [0,1], \ y \in \{0,1\}.$$

Divergence and entropy

For a scoring rule $S: \mathcal{P} \times \mathcal{Y} \to \overline{\mathbb{R}}$, we associate an *entropy*

$$H: \mathcal{P} \to \overline{R}, \quad H(F) = \int S(F,y) \, \mathrm{d}F(y) = \mathbb{E}_F S(F,Y) = S(F,F),$$

and a divergence

$$d: \mathcal{P} \times \mathcal{P} \to \overline{\mathbb{R}}, \quad d(F,G) = S(F,G) - H(G).$$

Proper scoring rules

Johanna Ziegel

ntroduction

efinition

Motivation

Divergence and entropy

Estimation

Classes of scoring rules

Characterization

Local scoring rules
f-scores and gf-scores
Kernel scores

Forecast comparison

nformation sets

imitations

Summary

Logarithmic Score (LogS)

f density of F

$$\mathsf{LogS}(F,y) = -\log f(y)$$

Proper scoring rules

Johanna Ziegel

Introduction

Definition

Motivation

Divergence and entr

Examples

Estimation

Classes of scoring

Characterization Local scoring rules

f-scores and gf-scor

orecast comparison

nformation sets

Limitations

Summary

Logarithmic Score (LogS)

f density of F

$$\mathsf{LogS}(F,y) = -\log f(y)$$

Entropy is Shannon entropy:

$$H(F) = -\int f(x) \log f(x) \,\mathrm{d}\mu(x) = -\mathbb{E}_F \log f(X)$$

Proper scoring rules

Johanna Ziegel

Introduction

Jennitio

Motivation

Divergence ar

Examples

Estimation

Classes of scoring

Characterization

Local scoring rules f-scores and gf-scores

orecast omparison

formation sets

Limitations

Julilliary

Logarithmic Score (LogS)

f density of F

$$\mathsf{LogS}(F,y) = -\log f(y)$$

Entropy is Shannon entropy:

$$H(F) = -\int f(x) \log f(x) d\mu(x) = -\mathbb{E}_F \log f(X)$$

Divergence is Kullback-Leibler divergence: g density of G

$$d(F,G) = \int g(y) \log \left(\frac{g(y)}{f(y)} \right) \, \mathrm{d}\mu(y) = D_{\mathrm{KL}}(G||F)$$

Proper scoring rules

Johanna Ziegel

ntroduction

Petinitio

Motivation

Divergence

Examples

Estimation

Classes of scoring

haracterization

Local scoring rules f-scores and gf-scores

orecast

mparison

ormation sets

imitations

Summary

Logarithmic Score (LogS)

f density of F

$$\mathsf{LogS}(F,y) = -\log f(y)$$

Entropy is Shannon entropy:

$$H(F) = -\int f(x) \log f(x) d\mu(x) = -\mathbb{E}_F \log f(X)$$

Divergence is Kullback-Leibler divergence: g density of G

$$d(F,G) = \int g(y) \log \left(\frac{g(y)}{f(y)} \right) \, \mathrm{d}\mu(y) = D_{\mathrm{KL}}(G||F)$$

Empirical risk minimization with respect to logarithmic score: Maximum likelihood estimation Proper scoring rules

Johanna Ziegel

Introduction

efinition

Motivation

Divergence

Examples

stimation

Classes of scoring ules

haracterization ocal scoring rules

f-scores and gf-scores

orecast omparison offormation sets

mitations

Summary

Continuous Ranked Probability Score (CRPS)

F CDF, finite mean

$$\mathsf{CRPS}(F,y) = \int_{\mathbb{R}} (F(z) - \mathbb{1}\{y \le z\})^2 \, \mathsf{d}z$$

Entropy

$$H(F) = \int F(x)(1 - F(x)) dx$$

Divergence G CDF, finite mean

$$d(F,G) = \int (F(y) - G(y))^2 dy$$

Central role in forecast evaluation in meteorology.

Proper scoring rules

Johanna Ziegel

Introduction

efinition

Motivation

Divergence a

Examples

Estimation

lasses of scoring iles

naracterization ocal scoring rules

cores and gf-sco

recast mparison

ormation sets

imitations

Summary

Proper scoring rules

Johanna Ziegel

Introduction

Definition

D-G-M-

Divergence and entrop

Estimation

Classes of scoring rules

haracterization
ocal scoring rules
-scores and of-score

-scores and gt-sco (ernel scores

orecast comparison

nformation sets

Limitations

Summary

References

Estimation

Motivation 3

Proper scoring rules Johanna Ziegel

Estimation

Let \mathcal{Y} be a Polish space and \mathcal{X} be some measurable covariate space.

Theorem

A scoring rule $S: \mathcal{P} \times \mathcal{Y} \to \overline{\mathbb{R}}$ is strictly proper if and only if for every pair of random variables $(X, Y) \in \mathcal{X} \times \mathcal{Y}$ such that

 \blacktriangleright the conditional distributions $\mathbb{P}_{Y|X=x} \in \mathcal{P}$

we have

$$\{\mathbb{P}_{Y|X=x}\}_{x} = \operatorname*{argmin}_{\{\mathbb{P}^{x}\}_{x}} \mathbb{E}\left[S(\mathbb{P}^{X}, Y)\right].$$

where $\mathbb{P}^{\chi}:\mathcal{X}\to\mathcal{P}$

nformation sets

ecompositions

mitations

ummary

References

Remark

$$\mathbb{P}_{Y|X=}$$
 . = $\operatorname*{argmin}_{\mathbb{P}^{(\cdot)}}\mathbb{E}\left[S(\mathbb{P}^X,Y)\right]$.

is analogous to

$$\mathbb{E}[Y|X = \cdot] = \operatorname*{argmin}_{f:\mathcal{X} \to \mathcal{Y}} \mathbb{E}[(Y - f(X))^2]$$

Proper scoring rules are to conditional distributions what the squared error loss (and Bregman divergences) are to conditional expectations.

Bias-Variance decomposition

Theorem

- ▶ Let $\{\mathbb{P}_{\theta,x}\}_{\theta\in\Theta}$ be a parametric family.
- ▶ Consider an estimator $\mathbb{P}_{\hat{\theta}_{\cdot}}$ of $\mathbb{P}_{Y|X=}$ where $\hat{\theta} = \hat{\theta}(\{(X_j, Y_j)\}_{j=1}^n)$.
- Then,

$$\mathbb{E}[d(\mathbb{P}_{\hat{\theta},X},\mathbb{P}_{Y|X})] = \underbrace{\mathbb{E}[d(\mathbb{P}_{\hat{\theta},X},\bar{\mathbb{P}}_{X})]}_{\text{variance}} + \underbrace{\mathbb{E}[d(\bar{\mathbb{P}}_{X},\mathbb{P}_{Y|X})]}_{\text{bias}}$$

where

$$ar{\mathbb{P}}_{\mathsf{x}} = \operatornamewithlimits{\mathsf{argmin}}_{\mathbb{P}} \mathbb{E}[d(\mathbb{P}_{\hat{ heta},\mathsf{x}},\mathbb{P})]$$

and

$$\mathbb{P}_{Y|X=x} = \operatorname*{argmin}_{\mathbb{P}} \mathbb{E}[S(\mathbb{P},Y) \,|\, X=x]$$

is the best predictor, assuming the two exist.

Proper scoring rules

Johanna Ziegel

ntroduction

efinition

Motivation

Divergence a

Estimation

Classes of scoring ules

Characterization

scores and gf-scores

orecast omparison

formation sets

mitations

Summary

References

(Pfau, 2013)

$$\begin{split} \mathbb{E}[d(\mathbb{P}_{\hat{\theta},X},\mathbb{P}_{Y|X})] &= \underbrace{\mathbb{E}[d(\mathbb{P}_{\hat{\theta},X},\bar{\mathbb{P}}_{X})]}_{\text{variance}} + \underbrace{\mathbb{E}[d(\bar{\mathbb{P}}_{X},\mathbb{P}_{Y|X})]}_{\text{bias}} \\ \bar{\mathbb{P}}_{x} &= \underset{\mathbb{P}}{\text{argmin}} \, \mathbb{E}[d(\mathbb{P}_{\hat{\theta},x},\mathbb{P})] \qquad \mathbb{P}_{Y|X=x} = \underset{\mathbb{P}}{\text{argmin}} \, \mathbb{E}[S(\mathbb{P},Y) \,|\, X=x] \end{split}$$

Proper scoring rules

Johanna Ziegel

Introduction

Definit

Motivation

Definition

Divergence and entrop

Estimation

Classes of scoring

Characterization

Local scoring rules f-scores and gf-score

Kernel scores

orecast comparison

Information sets

Decompositions

Liiiiidaaa

Summary

$$\begin{split} \mathbb{E}[d(\mathbb{P}_{\hat{\theta},X},\mathbb{P}_{Y|X})] &= \underbrace{\mathbb{E}[d(\mathbb{P}_{\hat{\theta},X},\bar{\mathbb{P}}_{X})]}_{\text{variance}} + \underbrace{\mathbb{E}[d(\bar{\mathbb{P}}_{X},\mathbb{P}_{Y|X})]}_{\text{bias}} \\ \bar{\mathbb{P}}_{x} &= \arg\!\min_{\mathbb{E}}\mathbb{E}[d(\mathbb{P}_{\hat{\theta},x},\mathbb{P})] \qquad \mathbb{P}_{Y|X=x} = \arg\!\min_{\mathbb{E}}\mathbb{E}[S(\mathbb{P},Y) \,|\, X=x] \end{split}$$

Take
$$S(F, y) = (m(F) - y)^2$$
. Then
$$d(F, G) = \mathbb{E}_{G}(m(F) - y)^2$$

$$d(F,G) = \mathbb{E}_G(m(F) - Y)^2 - \text{var}_G(Y) = (m(F) - m(G))^2.$$

$$\mathbb{E}[d(\mathbb{P}_{\hat{\theta},X},\mathbb{P}_{Y\mid X})] = \mathbb{E}[\mathbb{E}[d(\mathbb{P}_{\hat{\theta},X},\mathbb{P}_{Y\mid X})\mid X]] = \mathbb{E}(m_{\hat{\theta}}(X) - Y)^2 - \mathbb{E}\mathrm{var}(Y\mid X)$$

Proper scoring rules

Johanna Ziegel

Introduction

Definit

Motivation

Divergence and e

Examples

Estimation

lasses of scoring

haracterization

Local scoring rules

f-scores and gf-score

rnel scores

orecast comparison

nformation sets

Decompositions

bummary

$$\mathbb{E}[d(\mathbb{P}_{\hat{\theta},X},\mathbb{P}_{Y|X})] = \underbrace{\mathbb{E}[d(\mathbb{P}_{\hat{\theta},X},\bar{\mathbb{P}}_{X})]}_{\text{variance}} + \underbrace{\mathbb{E}[d(\bar{\mathbb{P}}_{X},\mathbb{P}_{Y|X})]}_{\text{bias}}$$

$$\bar{\mathbb{P}}_{X} = \underset{\mathbb{P}}{\operatorname{argmin}} \, \mathbb{E}[d(\mathbb{P}_{\hat{\theta},X},\mathbb{P})] \qquad \mathbb{P}_{Y|X=X} = \underset{\mathbb{P}}{\operatorname{argmin}} \, \mathbb{E}[S(\mathbb{P},Y) \,|\, X=X]$$

Take
$$S(F, y) = (m(F) - y)^2$$
. Then
$$d(F, G) = \mathbb{E}_{G}(m(F) - y)^2$$

$$d(F,G) = \mathbb{E}_G(m(F) - Y)^2 - \operatorname{var}_G(Y) = (m(F) - m(G))^2.$$

$$\mathbb{E}[d(\mathbb{P}_{\hat{\theta}|X}, \mathbb{P}_{Y|X})] = \mathbb{E}[\mathbb{E}[d(\mathbb{P}_{\hat{\theta}|X}, \mathbb{P}_{Y|X}) \mid X]] = \mathbb{E}(m_{\hat{\theta}}(X) - Y)^2 - \mathbb{E}\operatorname{var}(Y \mid X)$$

$$\operatorname*{argmin}_{\mathbb{P}}\mathbb{E}[d(\mathbb{P}_{\hat{\theta},x},\mathbb{P})] \equiv \operatorname*{argmin}_{z}\mathbb{E}(m_{\hat{\theta}}(x)-z)^2 = \mathbb{E}_{\hat{\theta}}m_{\hat{\theta}}(x).$$

Proper scoring rules

Johanna Ziegel

Introduction

Definit

Motivation

Divergence and

Estimation

Sumation

llasses of scoring ules

haracterization

cores and gf-score

recast

comparison

nformation sets

Limitation

Summary

$$\begin{split} \mathbb{E}[d(\mathbb{P}_{\hat{\theta},X},\mathbb{P}_{Y|X})] &= \underbrace{\mathbb{E}[d(\mathbb{P}_{\hat{\theta},X},\bar{\mathbb{P}}_{X})]}_{\text{variance}} + \underbrace{\mathbb{E}[d(\bar{\mathbb{P}}_{X},\mathbb{P}_{Y|X})]}_{\text{bias}} \\ \bar{\mathbb{P}}_{x} &= \operatorname{argmin} \mathbb{E}[d(\mathbb{P}_{\hat{\theta},x},\mathbb{P})] \qquad \mathbb{P}_{Y|X=x} = \operatorname{argmin} \mathbb{E}[S(\mathbb{P},Y) \,|\, X=x] \end{split}$$

Take
$$S(F, y) = (m(F) - y)^2$$
. Then

$$d(F,G) = \mathbb{E}_G(m(F) - Y)^2 - \text{var}_G(Y) = (m(F) - m(G))^2.$$

$$\mathbb{E}[d(\mathbb{P}_{\hat{\theta},X},\mathbb{P}_{Y\mid X})] = \mathbb{E}[\mathbb{E}[d(\mathbb{P}_{\hat{\theta},X},\mathbb{P}_{Y\mid X})\mid X]] = \mathbb{E}(m_{\hat{\theta}}(X) - Y)^2 - \mathbb{E}\mathrm{var}(Y\mid X)$$

$$\operatorname*{argmin}_{\mathbb{P}}\mathbb{E}[d(\mathbb{P}_{\hat{\theta},x},\mathbb{P})] \equiv \operatorname*{argmin}_{z}\mathbb{E}(m_{\hat{\theta}}(x)-z)^2 = \mathbb{E}_{\hat{\theta}}m_{\hat{\theta}}(x).$$

$$\mathbb{E}[d(\mathbb{P}_{\hat{\theta},X},\bar{\mathbb{P}}_X)] = \mathbb{E}[\mathbb{E}[d(\mathbb{P}_{\hat{\theta},X},\bar{\mathbb{P}}_X) \mid X]] = \mathbb{E}(m_{\hat{\theta}}(X) - \mathbb{E}_{\hat{\theta}}m_{\hat{\theta}}(X))^2 \quad \text{variance}$$

Proper scoring rules

Johanna Ziegel

Introduction

Definit

Motivation

Divergence an

Estimation

asses of scoring les

Characterization

scores and gf-score

orecast omparison

nformation sets

Limitatio

Summary

$$\mathbb{E}[d(\mathbb{P}_{\hat{\theta},X},\mathbb{P}_{Y|X})] = \underbrace{\mathbb{E}[d(\mathbb{P}_{\hat{\theta},X},\bar{\mathbb{P}}_{X})]}_{\text{variance}} + \underbrace{\mathbb{E}[d(\bar{\mathbb{P}}_{X},\mathbb{P}_{Y|X})]}_{\text{bias}}$$

$$\bar{\mathbb{P}}_{x} = \operatorname*{argmin}_{\mathbb{P}} \mathbb{E}[d(\mathbb{P}_{\hat{\theta},x},\mathbb{P})] \qquad \mathbb{P}_{Y|X=x} = \operatorname*{argmin}_{\mathbb{P}} \mathbb{E}[S(\mathbb{P},Y) \,|\, X=x]$$

Take
$$S(F, y) = (m(F) - y)^2$$
. Then

$$d(F,G) = \mathbb{E}_G(m(F) - Y)^2 - \text{var}_G(Y) = (m(F) - m(G))^2.$$

$$\mathbb{E}[d(\mathbb{P}_{\hat{\theta},X},\mathbb{P}_{Y\mid X})] = \mathbb{E}[\mathbb{E}[d(\mathbb{P}_{\hat{\theta},X},\mathbb{P}_{Y\mid X})\mid X]] = \mathbb{E}(m_{\hat{\theta}}(X) - Y)^2 - \mathbb{E}\mathrm{var}(Y\mid X)$$

$$\mathop{\rm argmin}_{\mathbb{P}} \mathbb{E}[d(\mathbb{P}_{\hat{\theta},x},\mathbb{P})] \equiv \mathop{\rm argmin}_{z} \mathbb{E}(m_{\hat{\theta}}(x)-z)^2 = \mathbb{E}_{\hat{\theta}} m_{\hat{\theta}}(x).$$

$$\mathbb{E}[d(\mathbb{P}_{\hat{\theta},X},\bar{\mathbb{P}}_X)] = \mathbb{E}[\mathbb{E}[d(\mathbb{P}_{\hat{\theta},X},\bar{\mathbb{P}}_X) \mid X]] = \mathbb{E}(m_{\hat{\theta}}(X) - \mathbb{E}_{\hat{\theta}}m_{\hat{\theta}}(X))^2 \quad \text{variance}$$

$$\mathbb{E}[d(\bar{\mathbb{P}}_X,\mathbb{P}_{Y|X})] = \mathbb{E}[\mathbb{E}[d(\bar{\mathbb{P}}_X,\mathbb{P}_{Y|X}) \mid X]] = \mathbb{E}(\mathbb{E}_{\hat{\theta}} m_{\hat{\theta}}(X) - \mathbb{E}(Y \mid X))^2 \quad \mathsf{bias}$$

Proper scoring rules

Johanna Ziegel

troduction

Definition

Motivation Definition

Divergence an Examples

Estimation

asses of scoring

haracterization

cal scoring rules

orecast omparison

ormation sets compositions

IIIItations

Proper scoring rules

Johanna Ziegel

Introduction

Definition

Motivation

Divergence and entre

xamples

Estimation

Classes of scoring rules

haracterization

Local scoring rules f-scores and gf-scores

rnel scores

orecast omparison

formation sets

Decompositions

Summary

References

Classes of scoring rules

Characterization Local scoring rules f-scores and gf-scores Kernel scores

Characterization of proper scoring rules

Let \mathcal{P} be a convex class of distributions on \mathcal{Y} .

Preliminaries

▶ Concavity: $H: \mathcal{P} \to \overline{\mathbb{R}}$ such that for $F, G \in \mathcal{P}$,

$$H(\alpha F + (1 - \alpha)G) \ge \alpha H(F) + (1 - \alpha)H(G)$$

Entropies H are concave!

▶ Supergradient: $h_F : \mathcal{Y} \to \overline{\mathbb{R}}$ that is *G*-integrable and

$$H(F) + \int h_F(y) d(G - F)(y) \geq H(G)$$

for all $G \in \mathcal{P}$.

▶ Regularity: A scoring rule $S: \mathcal{P} \times \mathcal{Y} \to \overline{\mathbb{R}}$ is called *regular* if H(F) = S(F, F) is finite and $S(F, G) > -\infty$ for every $F, G \in \mathcal{P}$ with $F \neq G$.

Proper scoring rules

Johanna Ziegel

Introduction

Definitio

Motivation

Divergence and Examples

Estimation

Classes of scoring rules

Characterization

Local scoring rules f-scores and gf-scores

orecast

nformation sets

imitations

Summary

Theorem

A regular scoring rule $S: \mathcal{P} \times \mathcal{Y} \to \overline{\mathbb{R}}$ is (strictly) proper if and only if there is a (strictly) concave function $H: \mathcal{P} \to \mathbb{R}$ such that

$$S(F,y) = H(F) + h_F(y) - \int h_F(x) dF(x)$$

for every $F \in \mathcal{P}$ and $y \in \mathcal{Y}$, where h_F is a supergradient of H at F.

▶ If H is differentiable, $h_F = \nabla_F H$, where $\nabla_F H$ is such that, for all $G \in \mathcal{P}$

$$\lim_{\alpha\downarrow 0}\frac{1}{\alpha}\left[H((1-\alpha)F+\alpha G)-H(F)\right]=\int \nabla_F H(y)\,\mathrm{d}G(y),$$

and

$$d(F,G) = H(F) - H(G) - \int \nabla_F H(y) d(F-G)(y).$$

(Gneiting and Raftery, 2007)

Proper scoring rules

Johanna Ziegel

Introduction

Definition

Motivation

Divergence and e

Estimation

asses of scoring les

Characterization

Local scoring rules
f-scores and gf-scores

orecast

nformation sets

imitation

ummary

Local scoring rules

- Proper scoring rules
 - Johanna Ziegel

- Local scoring rules

- $\triangleright \mathcal{V} \subseteq \mathbb{R}^d$ open
- $\triangleright \mathcal{P}^k$: probability measures on \mathcal{Y} with k-times differentiable densities

Definition

A proper scoring rule $S: \mathcal{P}^k \times \mathcal{V} \to \overline{\mathbb{R}}$ is local of order k if

$$S(F,y) = s(y, f(y), \dots, \nabla_y^k f(y)), \quad F \in \mathcal{P},$$

where f is the density of F.

Example

Local scoring rule of order 0: $S(F, v) = -\log f(v)$.

Uniqueness of logarithmic score

Theorem

Logarithmic score is (essentially) only differentiable local scoring rule of order 0.

Proper scoring rules

Johanna Ziegel

Introduction

Definition

Motivation

Divergence and entropy

Estimation

Classes of scoring ules

naracterization

Local scoring rules

scores and gr-scores ernel scores

orecast omparison

nformation sets

Limitations

Summany

Uniqueness of logarithmic score

Theorem

Logarithmic score is (essentially) only differentiable local scoring rule of order 0.

Proof.

Let S(F,y)=s(y,f(y)) with s(y,z) differentiable. Then $H(F)=\int s(y,f(y))f(y)\,\mathrm{d}y$.

$$\lim_{\alpha \downarrow 0} \frac{1}{\alpha} \left[H((1 - \alpha)F + \alpha G) - H(F) \right] = \int \nabla_F H(y) \, \mathrm{d}G(y)$$

$$= \int \left(s(y, f(y)) + \frac{\partial}{\partial z} s(x, f(y))f(y) \right) g(y) \, \mathrm{d}y - \int s(y, f(y))f(y) \, \mathrm{d}y - \int \frac{\partial}{\partial z} s(y, f(y))(f(y))^2 \, \mathrm{d}y$$

Proper scoring rules

Johanna Ziegel

Introduction

Definition

Motivation Definition

ivergence and entro xamples

Estimation

Classes of scoring rules

racterization

Local scoring rules

cores and gf-scores

orecast omparison

ormation sets

imitations

bummary

Uniqueness of logarithmic score

Theorem

Logarithmic score is (essentially) only differentiable local scoring rule of order 0.

Proof.

Let S(F, y) = s(y, f(y)) with s(y, z) differentiable. Then $H(F) = \int s(y, f(y))f(y) dy$.

$$\lim_{\alpha \downarrow 0} \frac{1}{\alpha} \left[H((1 - \alpha)F + \alpha G) - H(F) \right] = \int \nabla_F H(y) \, \mathrm{d}G(y)$$

$$= \int \left(s(y, f(y)) + \frac{\partial}{\partial z} s(x, f(y))f(y) \right) g(y) \, \mathrm{d}y - \int s(y, f(y))f(y) \, \mathrm{d}y - \int \frac{\partial}{\partial z} s(y, f(y))(f(y))^2 \, \mathrm{d}y$$

Hence, by the characterization theorem

$$S(F,y) = s(y,f(y)) + f(y)\frac{\partial}{\partial y}s(y,f(y)) - \int \underbrace{\frac{\partial}{\partial z}s(y,f(w))(f(w))^{2}}_{=cf(w)} dw$$

$$\Rightarrow \frac{\partial}{\partial z} s(y, f(w)) = \frac{c}{f(w)} \Rightarrow s(y, f(y)) = a \log f(y) + b \text{ for some } a, b \in \mathbb{R}.$$

Proper scoring rules

Johanna Ziegel

ntroduction

efinition

otivation efinition

stimation

lasses of scoring iles

Characterization

Local scoring rules

cores and gf-score

orecast omparison offormation sets

compositions

LIIIILALIOIIS

. .

Score matching

Example (Hyvärinen Score)

Let $\mathcal{Y} = \mathbb{R}^d$ and $\|\nabla_{\mathbf{x}} \log p(\mathbf{x})\| \to 0$ as $\|\mathbf{x}\| \to \infty$ for $p \in \mathcal{P}^2$. Then

$$S(P, \mathbf{y}) = \Delta_{\mathbf{y}} \log p(\mathbf{y}) + \frac{1}{2} \|\nabla_{\mathbf{y}} \log p(\mathbf{y})\|^{2} = \sum_{j} \left\{ \frac{\partial^{2} \log p}{\partial y_{j}^{2}} + \frac{1}{2} \left[\frac{\partial \log p}{\partial y_{j}} \right]^{2} \right\},$$

where $\nabla_{\mathbf{v}} f$ and $\Delta_{\mathbf{v}} f$ are gradient and Laplacian, is a strictly proper scoring rule.

Proper scoring rules

Johanna Ziegel

Local scoring rules

Score matching

Example (Hyvärinen Score)

Let $\mathcal{Y} = \mathbb{R}^d$ and $\|\nabla_{\mathbf{x}} \log p(\mathbf{x})\| \to 0$ as $\|\mathbf{x}\| \to \infty$ for $p \in \mathcal{P}^2$. Then

$$S(P, \mathbf{y}) = \Delta_{\mathbf{y}} \log p(\mathbf{y}) + \frac{1}{2} \|\nabla_{\mathbf{y}} \log p(\mathbf{y})\|^{2} = \sum_{j} \left\{ \frac{\partial^{2} \log p}{\partial y_{j}^{2}} + \frac{1}{2} \left[\frac{\partial \log p}{\partial y_{j}} \right]^{2} \right\},$$

where $\nabla_y f$ and $\Delta_y f$ are gradient and Laplacian, is a strictly proper scoring rule.

Non-normalized Densities

Consider the parametric family $\{P_{\theta}: \theta \in \Theta\}$. Let $dP_{\theta}/d\mu = p_{\theta}(\mathbf{x}) \propto \exp \eta_{\theta}(\mathbf{x})$.

$$S(P_{\theta}, \boldsymbol{x}) = \Delta_{\boldsymbol{x}} \eta_{\theta}(\boldsymbol{x}) + \frac{1}{2} \|\nabla_{\boldsymbol{x}} \eta_{\theta}(\boldsymbol{x})\|^{2}$$

Proper scoring rules

Johanna Ziegel

Introduction

Definition

Definition

Divergence and

Estimation

Classes of scoring rules

haracterization

Local scoring rules

-scores and gf-scores Kernel scores

Forecast comparison

formation sets

Limitation

Summary

Score matching

Example (Hyvärinen Score)

Let $\mathcal{Y} = \mathbb{R}^d$ and $\|\nabla_{\mathbf{x}} \log p(\mathbf{x})\| \to 0$ as $\|\mathbf{x}\| \to \infty$ for $p \in \mathcal{P}^2$. Then

$$S(P, \mathbf{y}) = \Delta_{\mathbf{y}} \log p(\mathbf{y}) + \frac{1}{2} \|\nabla_{\mathbf{y}} \log p(\mathbf{y})\|^{2} = \sum_{j} \left\{ \frac{\partial^{2} \log p}{\partial y_{j}^{2}} + \frac{1}{2} \left[\frac{\partial \log p}{\partial y_{j}} \right]^{2} \right\},$$

where $\nabla_y f$ and $\Delta_y f$ are gradient and Laplacian, is a strictly proper scoring rule.

Non-normalized Densities

Consider the parametric family $\{P_{\theta}: \theta \in \Theta\}$. Let $dP_{\theta}/d\mu = p_{\theta}(\mathbf{x}) \propto \exp \eta_{\theta}(\mathbf{x})$.

$$S(P_{ heta}, \mathbf{x}) = \Delta_{\mathbf{x}} \eta_{ heta}(\mathbf{x}) + \frac{1}{2} \|\nabla_{\mathbf{x}} \eta_{ heta}(\mathbf{x})\|^{2}$$

We have

$$\min_{\theta} \mathbb{E}_{Q} S(P_{\theta}, \mathbf{Y}) = \min_{\theta} \mathbb{E}_{Q} \|\nabla_{\mathbf{y}} \log p_{\theta}(\mathbf{Y}) - \nabla_{\mathbf{y}} \log q(\mathbf{Y})\|^{2}$$

Proper scoring rules

Johanna Ziegel

Introduction

efinition

Definition

Divergence and

stimation

les

racterization

Local scoring rules

cores and gf-scores

mparison
formation sets

ompositions

D.C

 μ measure on \mathcal{Y} , \mathcal{P} all absolutely continuous measures wrt μ $f:[0,\infty)\to\overline{\mathbb{R}}$ concave and differentiable. Define concave entropy

$$H_f(P) = \int f(p(x)) d\mu(x), \quad P \in \mathcal{P};$$

p density of $P \in \mathcal{P}$.

Proper scoring rules

Johanna Ziegel

Introduction

Definiti

Motivation

Divergence and entro

Estimation

Classes of scoring

Characterization

f-scores and gf-scores

ernel scores

orecast omparison

nformation sets

Limitations

Summary

 μ measure on \mathcal{Y} , \mathcal{P} all absolutely continuous measures wrt μ $f:[0,\infty)\to\overline{\mathbb{R}}$ concave and differentiable. Define concave entropy

$$H_f(P) = \int f(p(x)) d\mu(x), \quad P \in \mathcal{P};$$

p density of $P \in \mathcal{P}$.

This yields the proper scoring rule

$$S_f(P,y)=f'(p(y))+H_f(P)-\int f'(p(x))p(x)\,\mathrm{d}\mu(x).$$

Proper scoring rules

Johanna Ziegel

Introduction

Definition

Motivation

Divergence and entr

Estimation

lasses of scoring

cal scoring rules

f-scores and gf-scores

Kernel scores

orecast

nformation sets

Limitations

C.....

 μ measure on \mathcal{Y} , \mathcal{P} all absolutely continuous measures wrt μ $f:[0,\infty)\to\overline{\mathbb{R}}$ concave and differentiable. Define concave entropy

$$H_f(P) = \int f(p(x)) d\mu(x), \quad P \in \mathcal{P};$$

p density of $P \in \mathcal{P}$.

This yields the proper scoring rule

$$S_f(P,y) = f'(p(y)) + H_f(P) - \int f'(p(x))p(x) d\mu(x).$$

- ▶ Different name: Separable Bregman scores (Grünwald and Dawid, 2004)
- ▶ Only scoring rules of the form $S(P, y) = r(p(y)) + \ell(P)$.

Proper scoring rules

Johanna Ziegel

Introduction

Definitio

Motivation

Divergence and entr Examples

Estimation

lasses of scoring

naracterization ocal scoring rules

f-scores and gf-scores

Kernel scores

orecast comparison

Information sets
Decompositions

Limitation

Summary

$$S_f(P,y) = f'(p(y)) + H_f(P) - \int f'(p(x))p(x) d\mu(x).$$

Examples

Proper scoring rules

Johanna Ziegel

Introductio

Definition

Motivation

Divergence and

Examples

Estimation

Classes of scoring

Characterization

Local scoring rules

f-scores and gf-scores

ernal scores

Kernel scores

Forecast comparison

formation sets

Decompositions

ummary

Consider entropies of the form

$$H(P) = g\left(\int f(p(x)) d\mu(x)\right) = g(H_f(P)),$$

p density of $P \in \mathcal{P}$,

Proper scoring rules

Johanna Ziegel

Introduction

Definiti

Motivation

Divergence and entro

Estimation

Classes of scoring

naracterization ocal scoring rules

f-scores and gf-scores

Kernel scores

orecast comparison

nformation sets

Decompositions

Summary

Consider entropies of the form

$$H(P) = g\left(\int f(p(x)) d\mu(x)\right) = g(H_f(P)),$$

p density of $P \in \mathcal{P}$,

where f, g are such that H is concave, for example

- f concave, g concave and increasing;
- f convex, g concave and decreasing.

If f, g are differentiable, we obtain the proper scoring rule

$$S_{gf}(P,y) = g'(H_f(P))f'(p(y)) + g(H_f(P)) - g'(H_f(P)) \int f'(p(x))p(x) d\mu(x).$$

Proper scoring rules

Johanna Ziegel

Introduction

Definiti

Motivation

Divergence and e

Estimation

Classes of scoring rules

haracterization

f-scores and gf-scores

iner scores

Forecast comparison

formation sets

imitations

Summan.

Consider entropies of the form

$$H(P) = g\left(\int f(p(x)) d\mu(x)\right) = g(H_f(P)),$$

p density of $P \in \mathcal{P}$.

where f, g are such that H is concave, for example

- ► f concave, g concave and increasing;
- f convex, g concave and decreasing.

If f, g are differentiable, we obtain the proper scoring rule

$$S_{gf}(P,y) = g'(H_f(P))f'(p(y)) + g(H_f(P)) - g'(H_f(P)) \int f'(p(x))p(x) d\mu(x).$$

Example (Spherical score)

$$f(u) = u^2$$
, $g(u) = -\sqrt{u}$, $S(P, y) = -p(y)/\sqrt{\int p(x)^2 d\mu(x)}$.

Proper scoring rules

Johanna Ziegel

ntroduction

efinition

otivation

Examples

Estimation

lasses of scoring

Characterization

f-scores and gf-scores

scores and gf-so ernel scores

orecast omparison offormation sets

mitations

ummarv

 \mathcal{P} : class of probability measures on \mathbb{R} with finite mean. Probability measures specified as CDFs F.

Proper scoring rules

Johanna Ziegel

Kernel scores

 \mathcal{P} : class of probability measures on \mathbb{R} with finite mean. Probability measures specified as CDFs F.

Continuous Ranked Probability Score (CRPS)

$$S(F,y) = \int_{\mathbb{R}} (F(x) - \mathbb{1}\{y \le x\})^2 d(x)$$

$$= \int_0^1 (\mathbb{1}\{y \le F^{-1}(\alpha)\} - \alpha) (F^{-1}(\alpha) - y) d\alpha$$

$$= \mathbb{E}_F |X - y| - \frac{1}{2} \mathbb{E}_F |X - X'|$$

- ► Allows to compare discrete, continuous and mixed discrete-continuous distributions.
- Is becoming increasingly popoular also in estimation.

Proper scoring rules

Johanna Ziegel

Introduction

and the table of

Motivation

Divergence a

Estimation

Classes of scoring

Characterization
Local scoring rules

Kernel scores

Forecast comparison

formation sets

Limitations

Summary

Let
$$h(x, y) = |x - y|$$
.

Transformation Models

Let $g_{\theta}(\mathbf{N}) \sim \mathbb{P}_{\theta}$.

$$\hat{S} = \frac{1}{m} \sum_{i=1}^m h(g_{\theta}(\mathbf{N}_i), y) - \frac{1}{2m(m-1)} \sum_{i=1}^m \sum_{i': i' \neq i} h(g_{\theta}(\mathbf{N}_i), g_{\theta}(\mathbf{N}_{i'}))$$

Proper scoring rules

Johanna Ziegel

Introduction

Definition

Motivation

Divergence and entre

Estimation

Classes of scoring rules

Characterization

Local scoring rules

f-scores and gf-scores

Kernel scores

Forecast comparison

Information sets

Limitations

Summary

Let
$$h(x, y) = |x - y|$$
.

Transformation Models

Let $g_{\theta}(\mathbf{N}) \sim \mathbb{P}_{\theta}$.

$$\hat{S} = \frac{1}{m} \sum_{i=1}^{m} h(g_{\theta}(\mathbf{N}_i), y) - \frac{1}{2m(m-1)} \sum_{i=1}^{m} \sum_{i': i' \neq i} h(g_{\theta}(\mathbf{N}_i), g_{\theta}(\mathbf{N}_{i'}))$$

Conditional Transformation Models

Let $g_{ heta}(\mathbf{x},\mathbf{N}) \sim \mathbb{P}_{ heta,\mathbf{x}}$.

$$\hat{S} = \frac{1}{m} \sum_{i=1}^{m} h(g_{\theta}(\mathbf{x}, \mathbf{N}_i), y) - \frac{1}{2m(m-1)} \sum_{i=1}^{m} \sum_{i':i'\neq i} h(g_{\theta}(\mathbf{x}, \mathbf{N}_i), g_{\theta}(\mathbf{x}, \mathbf{N}_{i'}))$$

(Gneiting et al., 2005; Hothorn et al., 2014; Bouchacourt et al., 2016)

Proper scoring rules

Johanna Ziegel

Introduction

efinition

Motivation

vergence and

timation

Classes of scoring ules

ocal scoring rules
scores and gf-score

Kernel scores

orecast omparison

formation sets

mitations

.....

Kernel scores

Instead of h(x, y) = |x - y|...

Proper scoring rules

Johanna Ziegel

Introduction

Definition

Motivation

Divergence and en

Examples

Estimation

Classes of scoring ules

Characterization

Local scoring rules

f-scores and gf-scores

Kernel scores

Forecast comparison

nformation sets

Limitation

Summary

Kernel scores

Instead of h(x, y) = |x - y|...

Take conditionally negative definite kernel h on \mathcal{Y} :

Proper scoring rules

Johanna Ziegel

Introduction

Definition

Deminition

Definition

Divergence and entro

Estimation

Classes of scoring rules

Characterization

Local scoring rules

Kernel scores

orecast comparison

nformation sets

Limitations

Summany

Take conditionally negative definite kernel h on \mathcal{Y} : That is, $h: \mathcal{Y} \times \mathcal{Y} \to [0, \infty)$ is

- ightharpoonup symmetric: h(x,y) = h(y,x);
- $\forall n \geq 1, x_1, \ldots, x_n \in \mathcal{Y}, \alpha_1, \ldots, \alpha_n \in \mathbb{R}$ with $\sum_{j=1}^n \alpha_j = 0$, we have

$$\sum_{i,j=1}^n \alpha_i \alpha_j h(x_i, x_j) \leq 0.$$

Then,

$$H(P) = \frac{1}{2} \int \int h(x, y) dP(x) dP(y).$$

is concave with supergradient

$$\nabla_P H(y) = \int h(x,y) \, \mathrm{d}P(x) - 2H(P).$$

Proper scoring rules

Johanna Ziegel

Introduction

Definitio

Motivation

ivergence and ent

Estimation

Classes of scoring rules

haracterization
ocal scoring rules
-scores and gf-scores

Kernel scores

Forecast comparison

formation sets

iiiitations

ummary

Let
$$\mathcal{P}_h = \{P \mid H(P) < \infty\}.$$

Theorem (Kernel score)

If $h: \mathcal{Y} \times \mathcal{Y} \to [0, \infty)$ is a (strongly) conditionally negative definite kernel, then $S_h: \mathcal{P}_h \times \mathcal{Y} \to \overline{\mathbb{R}}$ given by

$$S_h(P, y) = \int h(x, y) dP(x) - \frac{1}{2} \int \int h(x, y) dP(x) dP(y)$$

is a (strictly) proper scoring rule.

Proper scoring rules

Johanna Ziegel

Introduction

Definitio

Motivation

Divergence and ent

Estimation

Classes of scoring

Characterization

Local scoring rules

F-scores and gf-scores

Kernel scores

Forecast comparison

nformation sets

Limitations

Summary

Let
$$\mathcal{P}_h = \{P \mid H(P) < \infty\}.$$

Theorem (Kernel score)

If $h: \mathcal{Y} \times \mathcal{Y} \to [0, \infty)$ is a (strongly) conditionally negative definite kernel, then $S_h: \mathcal{P}_h \times \mathcal{Y} \to \overline{\mathbb{R}}$ given by

$$S_h(P, y) = \int h(x, y) dP(x) - \frac{1}{2} \int \int h(x, y) dP(x) dP(y)$$

is a (strictly) proper scoring rule.

- ▶ Divergence: $d(P,Q) = -\frac{1}{2} \int \int h(x,y) d(P-Q)(x) d(P-Q)(y)$ Symmetric in P and Q!
- ► $S(P, y) = d(P, \delta_y) + \frac{1}{2}h(y, y)$ Divergence is itself a proper scoring rule.
- Divergence is squared maximum mean discrepancy (MMD).

(Gneiting and Raftery, 2007; Dawid, 2007; Steinwart and Ziegel, 2021)

Proper scoring rules

Johanna Ziegel

Introduction

Definition

Motivation

Examples

Estimation

Classes of scoring rules

Characterization

Local scoring rules

f-scores and gf-scores

Kernel scores

Forecast comparison

nformation sets

imitation.

ummarv

Connection to Hilbert space geometry

Theorem

There exists a Hilbert space $\mathcal H$ and a subset $\{\psi_x\}_{x\in\mathcal Y}\subset\mathcal H$ such that the divergence d of S satisfies

$$d(P,Q) = -\frac{1}{2} \iint \|\psi_x - \psi_y\|_{\mathcal{H}}^2 d(P-Q)(x) d(P-Q)(y).$$

Moreover,

$$d(P,Q) = \left\| \int \psi_{\mathsf{x}} \, dP(\mathsf{x}) - \int \psi_{\mathsf{x}} \, dQ(\mathsf{x}) \right\|_{\mathcal{H}}^{2}.$$

Proper scoring rules

Johanna Ziegel

Introduction

Definition

Motivation

Divergence and en Examples

Estimation

Classes of scoring rules

Local scoring rules
f-scores and gf-scores

Kernel scores

Forecast comparison

nformation sets

.iiiiiida tioiis

Summary

Connection to Hilbert space geometry

Theorem

There exists a Hilbert space \mathcal{H} and a subset $\{\psi_x\}_{x\in\mathcal{Y}}\subset\mathcal{H}$ such that the divergence d of S satisfies

$$d(P,Q) = -\frac{1}{2} \iint \|\psi_{x} - \psi_{y}\|_{\mathcal{H}}^{2} d(P-Q)(x) d(P-Q)(y).$$

Moreover,

$$d(P,Q) = \left\| \int \psi_{\mathsf{x}} \, dP(\mathsf{x}) - \int \psi_{\mathsf{x}} \, dQ(\mathsf{x}) \right\|_{\mathcal{H}}^{2}.$$

Kernel score can be constructed on ${\mathcal Y}$ if

- lacktriangle we can construct a conditionally negative definite kernel on \mathcal{Y} ;
- we can construct a positive definite kernel k on \mathcal{Y} (take h = -k);
- ightharpoonup we can embed \mathcal{Y} in a Hilbert space \mathcal{H} .

Proper scoring rules

Johanna Ziegel

Introduction

Definition

Motivation

Divergence and Examples

Estimation

Classes of scoring rules

Local scoring rules
f-scores and gf-scores

Kernel scores

Forecast comparison

Decompositions

Summary

Kernels on \mathbb{R}^d

Bounded continuous kernels

Radial kernels that are strongly positive definite for any *d*:

$$k(x,y) = \varphi(\|x - y\|),$$

where

$$\varphi(t) = \int_0^\infty \exp(-t^2 s) \,\mathrm{d}\nu(s)$$

for a measure μ with supp $\mu \neq \{0\}$.

(Sriperumbudur et al., 2011)

Proper scoring rules

Johanna Ziegel

Introduction

efinition

Motivation

Divergence and e

Estimation

Classes of scoring

Characterization

ocal scoring rules
-scores and gf-scores

Kernel scores

orecast

nformation sets

imitations

Summary

Kernels on \mathbb{R}^d

Bounded continuous kernels

Radial kernels that are strongly positive definite for any d:

$$k(x,y) = \varphi(||x-y||),$$

where

$$\varphi(t) = \int_0^\infty \exp(-t^2 s) \, \mathrm{d}\nu(s)$$

for a measure μ with supp $\mu \neq \{0\}$.

(Sriperumbudur et al., 2011)

Distance kernels

$$h(x,y) = \|x - y\|^{\alpha},$$

for $\alpha \in (0,2)$ are strongly conditionally negative definite for any d.

- ▶ CRPS corresponds to d = 1, $\alpha = 1$.
- Energy scores.

Proper scoring rules

Johanna Ziegel

ntroduction

emnition

Motivation

Divergence a

Estimation

lasses of scoring

haracterization

cores and gf-scores

Kernel scores

recast

rmation sets

ecompositions

nitations

References

Connection to energy distance

Székely and Rizzo (2004), Baringhaus and Franz (2004) introduced the *energy distance* between two distributions P, Q on \mathbb{R}^d with finite first moments

$$\mathbb{E}||Z - W|| - \frac{1}{2}\mathbb{E}||Z - Z'|| - \frac{1}{2}\mathbb{E}||W - W'||,$$

where Z, Z', W, W' are independent with $Z, Z' \sim P$, $W, W' \sim Q$.

Proper scoring rules

Johanna Ziegel

Introduction

Definition

Motivation

Divergence and entrop

Estimation

Classes of scoring ules

Characterization

scores and gf-sco

Kernel scores

Forecast

Information sets

imitations

C.....

Divergence and entre

Estimation

Classes of scoring rules

Characterization

Local scoring rules

scores and gf-scor

Kernel scores

orecast omparison

nformation sets

imitation

ummarv

eferences

Székely and Rizzo (2004), Baringhaus and Franz (2004) introduced the *energy distance* between two distributions P, Q on \mathbb{R}^d with finite first moments

$$\mathbb{E}||Z - W|| - \frac{1}{2}\mathbb{E}||Z - Z'|| - \frac{1}{2}\mathbb{E}||W - W'||,$$

where Z, Z', W, W' are independent with $Z, Z' \sim P$, $W, W' \sim Q$.

- ▶ Energy distance is the divergence of the kernel score with h(x, y) = ||x y|| called the *energy score*.
- ► Energy distance is a squared maximum mean discrepancy between P and Q (Sejdinovic et al., 2013).
- ► Energy score is a popular strictly proper scoring rule for multivariate outcomes.

Distance kernels

When is

$$h(x,y) = ||x - y||$$

strongly conditionally negative definite? Metric of strong negative type

Proper scoring rules

Johanna Ziegel

Introductio

Definition

Motivation

Divergence and

Examples

Estimation

Classes of scoring

haracterization

ocal scoring rules -scores and gf-score

Kernel scores

Forecast

nformation sets

Limitation

Summary

Distance kernels

Proper scoring rules

Johanna Ziegel

Introduction

Definition

Notivation

Divergence and en

Examples

Estimation

Classes of scoring

Characterization

Local scoring rules f-scores and gf-scores

Kernel scores

Forecast

formation sets

Decompositions

Summary

References

When is

$$h(x,y) = ||x - y||$$

strongly conditionally negative definite? Metric of strong negative type

- ► Separable Hilbert spaces
- ▶ Separable L^p -spaces for 1

(Linde, 1986; Lyons, 2013; Sejdinovic et al., 2013)

Characterizations of kernel scores

Proper scoring rules

Johanna Ziegel

Introduction

Definition

Motivation

Divergence a

campies

Estimation

Classes of scoring rules

Characterization

Local scoring rules

f-scores and gf-scores

Kernel scores

orecast omparison

formation sets

imitations

Summary

deferences

Theorem

Let $S: \mathcal{P} \times \mathcal{Y} \to \overline{\mathbb{R}}$ be a proper scoring rule such that $\{\delta_y: y \in \mathcal{Y}\} \subseteq \mathcal{P}$ and the supergradient map $P \mapsto h_P$ is weakly continuous.

Then, the divergence d of S is symmetric if and only if S is a kernel score.

Characterizations of kernel scores

Proper scoring rules

Johanna Ziegel

Kernel scores

Theorem

Let $S: \mathcal{P} \times \mathcal{Y} \to \overline{\mathbb{R}}$ be a proper scoring rule such that $\{\delta_v : v \in \mathcal{Y}\} \subseteq \mathcal{P}$ and the supergradient map $P \mapsto h_P$ is weakly continuous.

Then, the divergence d of S is symmetric if and only if S is a kernel score.

 \blacktriangleright A proper scoring rule with $\{\delta_v : v \in \mathcal{Y}\} \subseteq \mathcal{P}$ corresponds to a squared metric on measures if and only if it is a kernel score!

(Waghmare and Ziegel, 2025)

Then, the divergence d of S is symmetric if and only if S is a kernel score.

- \blacktriangleright A proper scoring rule with $\{\delta_v : v \in \mathcal{Y}\} \subseteq \mathcal{P}$ corresponds to a squared metric on measures if and only if it is a kernel score!
- ▶ This also provides a natural motivation for using kernel-MMD for two sample testing.

(Waghmare and Ziegel, 2025)

Kernel scores

Characterizations of kernel scores

Let S be a scoring rule.

▶ Translation invariance: For $\mathbf{y}, \mathbf{h} \in \mathbb{R}^d$,

$$S(P, \mathbf{y}) = S(P_{\mathbf{h}}, \mathbf{y} + \mathbf{h})$$

where $P_{\mathbf{h}}(A) = P(A + \mathbf{h})$ for Borel sets $A \subset \mathbb{R}^d$.

▶ Homogeneity: For every c > 0, $P \in \mathcal{P}$ and $\mathbf{y} \in \mathbb{R}^d$,

$$S(P_c, c\mathbf{y}) = c^{\alpha}S(P, \mathbf{y})$$

where $P_c(A) = P(c^{-1}A)$ for Borel sets $A \subset \mathbb{R}^d$.

▶ Isotropy: For every rotation matrix $\mathbf{U} \in SO(d)$, $P \in \mathcal{P}$ and $\mathbf{y} \in \mathbb{R}^d$,

$$S(P_{\mathbf{U}}, \mathbf{U}\mathbf{y}) = S(P, \mathbf{y})$$

where $P_{\mathbf{U}}(A) = P(\mathbf{U}^{\top}A)$ for Borel sets $A \subset \mathbb{R}^d$.

Proper scoring rules

Johanna Ziegel

Introduction

finition

lotivation

Divergence a

stimation

lasses of scoring iles

ocal scoring rules

Kernel scores

orecast

formation sets

-14-41---

Lillitations

References

Kererences

Proper scoring rules

Johanna Ziegel

Introduction

.

Motivation

Divergence and ent

Examples

Estimation

Classes of scoring rules

Characterization

ocal scoring rules
-scores and gf-scores

Kernel scores

Forecast

formation sets

imitations

.....

References

Up to positive multiplicative constants:

- 1. CRPS \rightarrow only 1-homogeneous translation invariant kernel score on \mathbb{R} , and
- 2. Energy Scores \rightarrow only homogeneous isotropic translation invariant kernel scores on \mathbb{R}^d

(Waghmare and Ziegel, 2025)

Foundations

Proper scoring rules

Johanna Ziegel

Kernel scores

Where do our loss functions come from?

Logarithmic score. Only proper scoring rule of the form S(P, y) = s(y, p(y)).

- Kernel Maximum Mean Discrepancy. Only proper scoring rules which admit point measures and have symmetric divergences.
- Squared error loss. Only Bregman divergence which is symmetric and isotropic.

Proper scoring rules

Johanna Ziegel

introductio

Definition

Mastration

Definition

Divergence and entrop

Examples

Estimation

Classes of scoring

Characterization

-scores and gf-score

ernel scores

Forecast comparison

Information sets

imitations

Summary

References

Forecast comparison Information sets

Decompositions

References

Pick a strictly proper scoring rule S, compare the average realized score: For $(F_1, G_1, Y_1), \ldots, (F_n, G_n, Y_n)$, let

$$\hat{S}_1 = \frac{1}{n} \sum_{i=1}^n S(F_i, Y_i)$$
 and $\hat{S}_2 = \frac{1}{n} \sum_{i=1}^n S(G_i, Y_i)$.

The forecast with the smaller value \hat{S}_i is better.

- ► Formal tests for differences between expected scores available.

 (Diebold and Mariano, 1995; Giacomini and White, 2006; Lai et al., 2011; Henzi and Ziegel, 2022; Choe and Ramdas, 2024)
- ▶ If (F_i, Y_i) are iid, classical (asymptotic) t-test can be used.

Simulation example

Let $X_1 \sim \mathcal{N}(0,1)$, $X_2 \sim \mathcal{N}(0,2)$ be independent, and

$$\mathbb{P}(Y = 1 \mid X_1, X_2) = \Phi(X_1 + X_2).$$

Predictions:

$$p^{(0)} = 1/2$$
, $p^{(1)} = \Phi\left(\frac{X_1}{\sqrt{3}}\right)$, $p^{(2)} = \Phi\left(\frac{X_2}{\sqrt{2}}\right)$, $p^{(3)} = \Phi(X_1 + X_2)$.

n = 200.

Prediction	Brier Score	Logarithmic score
p_0	0.250	0.693
p_1	0.213	0.613
p_2	0.167	0.499
<i>p</i> ₃	0.116	0.355

Proper scoring rules

Johanna Ziegel

Introduction

Definition

Notivation

Divergence and enti

stimation

Classes of scoring rules

ocal scoring rules
-scores and gf-scores

Forecast comparison

nformation sets

Limitation

Summary

Does forecast comparison depend on the choice of the proper scoring rule *S*? Generally, yes.

- ► Finite samples
- ► Non-nested information sets
- Uncalibrated predicitions/misspecified models

(Patton, 2020; Ziegel et al., 2020)

Proper scoring rules

Johanna Ziegel

Introduction

Definition

.....

Motivation

Divergence an

xamples

Estimation

Classes of scoring

Characterization

Local scoring rules f-scores and gf-scores

nel scores

Forecast comparison

nformation sets
Decompositions

mitations

Summary

Does forecast comparison depend on the choice of the proper scoring rule S? Generally, yes.

- ► Finite samples
- Non-nested information sets
- Uncalibrated predicitions/misspecified models

(Patton, 2020; Ziegel et al., 2020)

Can we avoid the choice of a proper scoring rule *S*?

- For binary outcomes, sometimes yes (Murphy diagrams)
- Otherwise, no.

Proper scoring rules

Johanna Ziegel

Introduction

efinition

NA COLUMN

....

Divergence and

amples

stimation

Classes of scoring

haracterization

Local scoring rules

ores and gf-scores

nel scores

Forecast comparison

Information sets

Decompositions

illitations

Summary

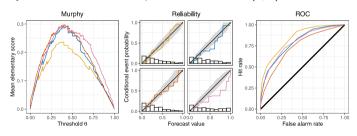
Murphy diagrams

All proper scoring rules for binary outcomes $Y \in \{0,1\}$ can be written as

$$S(p,y) = \int_0^1 S_{\theta}(p,y) \, \mathrm{d}H(\theta)$$

for some measure H on (0,1).

- ▶ One-parameter family $(S_{\theta}(p, y))_{\theta}$ of elementary scores.
- Forecast p is better with respect to all proper scoring rules S if and only if it is better with respect to all S_{θ} , $\theta \in (0,1)$.



Ehm et al. (2016); Krüger and Ziegel (2021); Figure from Dimitriadis et al. (2024)

Proper scoring rules

Johanna Ziegel

Introduction

Definitio

Motivation

Divergence and e

Estimation

Classes of scoring rules

haracterization

Local scoring rules f-scores and gf-scores

Forecast comparison

Information sets

Limitation

Summary

Why should we use proper scoring rules to evaluate predictions?

▶ They incentivize truthful (calibrated) and informative forecasts.

Let S be a (strictly) proper scoring rule.

Theorem

Let
$$F = \mathcal{L}(Y \mid X)$$
, G based on X ($\sigma(X)$ -measurable). Then,

$$\mathbb{E}S(F,Y) \leq \mathbb{E}S(G,Y).$$

Equality implies that F = G almost surely.

Proper scoring rules

Johanna Ziegel

Introduction

Definition

Motivation

Divergence and en

Estimation

Classes of scoring

haracterization

ocal scoring rules scores and gf-scores

recast

Information sets

ecompositions

_imitations

ummary

Why should we use proper scoring rules to evaluate predictions?

▶ They incentivize truthful (calibrated) and informative forecasts.

Let S be a (strictly) proper scoring rule.

Theorem

Let $F = \mathcal{L}(Y \mid X)$, G based on X ($\sigma(X)$ -measurable). Then,

$$\mathbb{E}S(F, Y) \leq \mathbb{E}S(G, Y).$$

Equality implies that F = G almost surely.

Corollary

Let $F = \mathcal{L}(Y \mid X, Z)$, $G = \mathcal{L}(Y \mid X)$. Then,

$$\mathbb{E}S(F,Y) \leq \mathbb{E}S(G,Y).$$

Equality happens if and only if Y and Z are conditionally independent given X.

(Holzmann and Eulert, 2014)

rules

Johanna Ziegel

Johanna Zi

Proper scoring

troduction

efinition

tivation

ivergence and

timation

IIIIation

les

aracterization

cal scoring rule

cores and *gf* rnel scores

parison

Information sets

compositions

imitations

ummary

References

43 / 58

Scoring rules and calibration

Goal

Decompose

$$\hat{S} = \frac{1}{n} \sum_{k=1}^{n} S(F_i, y_i)$$

into interpretable terms quantifying miscalibration (MCB), discrimination ability (DSC), and uncertainty (UNC).

Idea

$$\mathbb{E}S(F,Y) = \underbrace{\mathbb{E}S(F,Y) - \mathbb{E}S(\mathcal{L}(Y\mid F),Y)}_{\text{MCB}} - \left(\underbrace{\mathbb{E}S(\mathcal{L}(Y),Y) - \mathbb{E}S(\mathcal{L}(Y\mid F),Y)}_{\text{DSC}}\right)$$

- \triangleright $\mathcal{L}(Y \mid F)$ is best auto-calibrated forecast given information F.
- \triangleright $\mathcal{L}(Y)$ is uninformative but auto-calibrated.

Proper scoring rules Johanna Ziegel

DSC

Decompositions

Empirical translation of the idea

Data: $(F_1, y_1), \dots, (F_n, y_n)$

Proper scoring rules

Johanna Ziegel

Introduction

Definition

Delilitio

Definition

Divergence and entrop

Estimation

Classes of scoring

Characterization
Local scoring rules
f-scores and gf-scores

orecast omparison

formation sets

Decompositions

imitations

Summary

Empirical translation of the idea

Data: $(F_1, y_1), \dots, (F_n, y_n)$

Uncertainty: UNC = $\mathbb{E}S(\mathcal{L}(Y), Y)$

$$\overline{\mathsf{UNC}} = \frac{1}{n} \sum_{i=1}^{n} S\left(\frac{1}{n} \sum_{k=1}^{n} \delta_{y_k}, y_i\right)$$

Proper scoring rules

Johanna Ziegel

Introduction

Definition

.

Divergence and e

Examples

Estimation

Classes of scoring

Characterization

Local scoring rules

f-scores and gf-score
Kernel scores

orecast omparison

formation sets

Decompositions

bummary

Empirical translation of the idea

Data: $(F_1, y_1), \dots, (F_n, y_n)$

Uncertainty: UNC = $\mathbb{E}S(\mathcal{L}(Y), Y)$

$$\overline{\mathsf{UNC}} = \frac{1}{n} \sum_{i=1}^{n} S\left(\frac{1}{n} \sum_{k=1}^{n} \delta_{y_k}, y_i\right)$$

Miscalibration: $MCB = \mathbb{E}S(F, Y) - \mathbb{E}S(\mathcal{L}(Y \mid F), Y)$

- Need estimate of $\mathcal{L}(Y \mid F)$ that is in-sample auto-calibrated and such that $\overline{\text{MCB}} \geq 0$. Generally not feasible.
- ▶ If S = CRPS, there is possible solution using isotonic distributional regression (Henzi et al., 2021).

Discrimination: DSC = $\mathbb{E}S(\mathcal{L}(Y), Y) - \mathbb{E}S(\mathcal{L}(Y \mid F), Y)$

► Estimate DSC by $\hat{S} - \overline{MCB} - \overline{UNC}$.

Proper scoring rules

Johanna Ziegel

Introduction

Definitio

Motivation

Divergence an

Lampies

Estimation

Classes of scoring ules

haracterization

ocal scoring rules -scores and gf-score

rnel scores

orecast omparison

Information sets

Decompositions

.....

Summary

Data application

Proper scoring rules

Johanna Ziegel

Introduction

Definition

Motivation

Divergence and

Estimation

Classes of scoring rules

Characterization

ocal scoring rules scores and gf-sco

ernel scores

orecast

ormation sets

Decompositions

imitations

ummary

References

Compare probabilistic quantitative precipitation forecasts.

Numerical weather prediction models

- ▶ Physical model of the atmosphere is run with current (measured) inital conditions
- ▶ Initial conditions are measured with error: Several model runs with slightly perturbed inital conditions yields *ensemble of forecasts*
- ► Forecast ensembles are interpreted as random draws from the conditional distribution of the outcome
- Ensembles are usually biased and underdispersed: Statistical postprocessing

daily blased and anderdispersed. Statistical postprocessing

Bauer et al. (2015)

Airport station observations at London and Zurich

- ▶ 52 member raw ensemble forecasts from European Centre for Medium-Range Weather Forecasts (ECMWF)
- ► Training data from 2007–2014, evaluation data from 2015-2016
- ▶ Prediction horizons of 1-5 days

Postprocessing methods:

- ► Bayesian model averaging (BMA)
- Ensemble model output statistics (EMOS)
- ► Heteroscedastic censored logistic regression (HCLR)
- ► Isotonic distributional regression (IDR)

Proper scoring rules

Johanna Ziegel

Introduction

Definition

Notivation

Divergence and entro

Estimation

lasses of scoring iles

haracterization

ocal scoring rules -scores and gf-scores

rnel scores

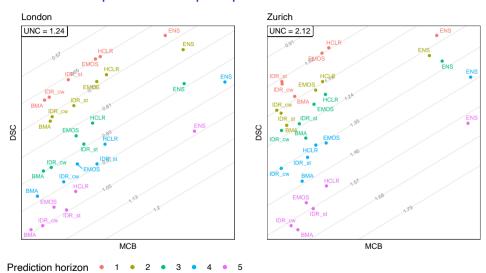
orecast omparison

formation sets

Decompositions

ummary

Probabilistic quantitative precipitation forecasts



Proper scoring rules

Johanna Ziegel

Introduction

efinition

Activation

ivergence and

Estimation

Classes of scoring

Characterization

scores and gf-sco

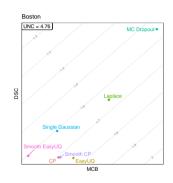
Forecast

formation sets

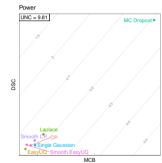
Decompositions

Summary

Uncertainty quantification on benchmark datasets in machine learning







Arnold et al. (2024)

Proper scoring rules

Johanna Ziegel

Introduction

Definition

Jennition

efinition ivergence and

Estimation

Classes of scoring rules

Characterization
Local scoring rules

recast

omparison

Decompositions

Limitations

Summan

Proper scoring rules

Johanna Ziegel

Introduction

Definition

14.00

Dofiniti

Divergence and entrop

Estimation

Classes of scoring ules

haracterization
ocal scoring rules

scores and gf-scor

Forecast comparison

nformation sets
Decompositions

Limitations

Summary

References

50 / 58

Proper scoring rules and extremes

Evaluation of forecast behaviour with respect to extreme events requires care (Lerch et al., 2017): Never condition on the outcomes when evaluating: Forecasters dilemma.

Proper scoring rules

Johanna Ziegel

Proper scoring rules and extremes

- Evaluation of forecast behaviour with respect to extreme events requires care (Lerch et al., 2017): Never condition on the outcomes when evaluating: Forecasters dilemma.
- Proper scoring rules are not necessarily directly useful: Expected scores cannot distinguish different tail behaviour (Brehmer and Strokorb, 2019).

Proper scoring rules

Johanna Ziegel

Proper scoring rules and extremes

- Evaluation of forecast behaviour with respect to extreme events requires care (Lerch et al., 2017): Never condition on the outcomes when evaluating: Forecasters dilemma.
- ▶ Proper scoring rules are not necessarily directly useful: Expected scores cannot distinguish different tail behaviour (Brehmer and Strokorb, 2019).
- ▶ However, evaluating calibration of probabilistic predictions with regards to tails is possible (Allen et al., 2025).

Proper scoring rules

Johanna Ziegel

Proper scoring rules

Johanna Ziegel

Introduction

Definition

D C ...

Divergence and entrop

kamples

Estimation

lasses of scoring

Characterization Local scoring rules

Local scoring rules f-scores and gf-scores Kernel scores

orecast comparison

nformation sets

I incheston

Summary

References

Summary

Summary

- Probabilistic predictions should be calibrated and sharp.
- Ideally, probabilistic predictions should be auto-calibrated.
- Proper scoring rules allow to compare probabilistic predictions simultaneously with respect to calibration and discrimination ability.

Proper scoring rules

Johanna Ziegel

Summary

Divergence and

stimation

Classes of scoring rules

Characterization

scores and gf-scor

orecast

formation sets

Decompositions

Limitation

Summary

References

- S. Allen, J. Koh, J. Segers, and J. Ziegel. Tail calibration of probabilistic forecasts. *Journal of the American Statistical Association*, 2025. To appear.
- S. Arnold, E.-M. Walz, J. Ziegel, and T. Gneiting. Decompositions of the mean continuous ranked probability score. *Electronic Journal of Statistics*, 18:4992–5044, 2024.
- L. Baringhaus and C. Franz. On a new multivariate two-sample test. *Journal of Multivariate Analysis*, 88:190–206, 2004.
- P. Bauer, A. Thorpe, and G. Brunet. The quiet revolution of numerical weather prediction. *Nature*, 525:47–55, 2015.
- D. Bouchacourt, P. K. Mudigonda, and S. Nowozin. DISCO nets: Dissimilarity coefficient networks. In *Advances in Neural Information Processing Systems*, volume 29, pages 352–360, 2016. URL https:

//papers.nips.cc/paper/6143-disco-nets-dissimilarity-coefficients-networks.

- J. Brehmer and K. Strokorb. Why scoring functions cannot assess tail properties. *Electronic Journal of Statistics*, 13:4015–4034, 2019.
- G. W. Brier. Verification of forecasts expressed in terms of probability. *Monthly Weather Review*, 78:1–3, 1950.

References

Y. J. Choe and A. Ramdas. Comparing sequential forecasters. *Operations Research*, 72: 1368-1387, 2024,

A. P. Dawid. The geometry of proper scoring rules. Annals of the Institute of Statistical Mathematics, 59:77–93, 2007.

- A. P. Dawid, M. Musio, and L. Ventura. Minimum scoring rule inference. Scandinavian Journal of Statistics, 43:123-138, 2016.
- F. X. Diebold and R. S. Mariano. Comparing predictive accuracy. Journal of Business & Economic Statistics, 13:253-263, 1995.
- T. Dimitriadis, T. Gneiting, A. I. Jordan, and P. Vogel. Evaluating probabilistic classifiers: The triptych. International Journal of Forecasting, 40:1101-1122, 2024.
- W. Ehm, T. Gneiting, A. Jordan, and F. Krüger. Of quantiles and expectiles: Consistent scoring functions. Choquet representations, and forecast rankings (with discussion). Journal of the Royal Statistical Society: Series B. 78:505–562. 2016.
- R. Giacomini and H. White. Tests of conditional predictive ability. Econometrica, 74: 1545-1578, 2006.
- T. Gneiting and A. E. Raftery. Strictly proper scoring rules, prediction, and estimation. Journal of the American Statistical Association, 102:359-378, 2007.

Divergence and

stimation

Classes of scoring rules

Characterization
Local scoring rules
f-scores and gf-scores

orecast omparison

nformation sets

Limitation

immary

- T. Gneiting, A. E. Raftery, A. H. Westveld, and T. Goldman. Calibrated probabilistic forecasting using ensemble model output statistics and minimum CRPS estimation. *Monthly Weather Review*, 133:1098–1118, 2005.
- A. Gretton, K. M. Borgwardt, M. J. Rasch, B. Schölkopf, and A. Smola. A kernel two-sample test. *Journal of Machine Learning Research*, 13:723–773, 2012.
- P. D. Grünwald and A. P. Dawid. Game theory, maximum entropy, minimum discrepancy and robust Bayesian decision theory. *Annals of Statistics*, 32:1367–1433, 2004.
- A. Henzi and J. F. Ziegel. Valid sequential inference on probability forecast performance. Biometrika, 109:647–663, 2022.
- A. Henzi, J. F. Ziegel, and T. Gneiting. Isotonic distributional regression. *Journal of the Royal Statistical Society: Series B*, 85:963–993, 2021.
- H. Holzmann and M. Eulert. The role of the information set for forecasting with applications to risk management. *Annals of Applied Statistics*, 8:595–621, 2014.
- T. Hothorn, T. Kneib, and P. Bühlmann. Conditional transformation models. *Journal of the Royal Statistical Society: Series B*, 76:3–27, 2014.
- F. Krüger and J. F. Ziegel. Generic conditions for forecast dominance. *Journal of Business & Economic Statistics*, 39:972–983, 2021.

Motivation

Divergence an

stimation

Classes of scoring rules

Characterization

Local scoring rules

f-scores and gf-scores

orecast

nformation sets

Decompositions

mitations

Summary

References

T. L. Lai, S. T. Gross, and D. B. Shen. Evaluating probability forecasts. *Annals of Statistics*, 39:2356–2382, 2011

- S. Lerch, T. L. Thorarinsdottir, F. Ravazzolo, and T. Gneiting. Forecaster's dilemma: Extreme events and forecast evaluation. *Statistical Science*, 32:106–127, 2017.
- W. Linde. On Rudin's equimeasurability theorem for infinite dimensional Hilbert spaces. *Indiana University Mathematics Journal*, 35:235–243, 1986.
- R. Lyons. Distance covariance in metric spaces. Annals of Probability, 41:3284-3305, 2013.
- J. McCarthy. Measures of the value of information. P. Natl. Acad. Sci. USA, 42:654–655, 1956.
- A. J. Patton. Comparing possibly misspecified forecasts. *Journal of Business & Economic Statistics*, 38:796–809, 2020.
- D. Pfau. A generalized bias-variance decomposition for Bregman divergences. Preprint, http://davidpfau.com/assets/generalized_bvd_proof.pdf, 2013. Acessed 2025-02-25.
- D. Sejdinovic, B. Sriperumbudur, A. Gretton, and K. Fukumizu. Equivalence of distance-based and RKHS-based statistics in hypothesis testing. *Annals of Statistics*, 41:2263–2291, 2013.

- B. K. Sriperumbudur, K. Fukumizu, and G. R. G. Lanckriet. Universality, characteristic kernels and RKHS embedding of measures. Journal of Machine Learning Research, 12:2389–2410. 2011.
- I. Steinwart and J. F. Ziegel. Strictly proper kernel scores and characteristic kernels on compact spaces. Applied and Computational Harmonic Analysis, 51:510–542, 2021.
- G. J. Székely and M. Rizzo. Testing for equal distribution in high dimension. *InterStat*, 5, November 2004
- K. Waghmare and J. Ziegel. Proper scoring rules for estimation and forecast evaluation. Annual Review of Statistics and Its Application, 2025. To appear.
- J. F. Ziegel, F. Krüger, A. Jordan, and F. Fasciati, Robust forecast evaluation of expected shortfall. Journal of Financial Econometrics, 18:95–120, 2020.