$X_1$ $\beta_1$

$X_2$ $\beta_2$

$X_3$ $\beta_3$

label ▨         coeff.

$X_1$

$X_2$

$X_3$

$$|X_1 \cap X_2 \cap X_3| = \frac{1}{64} \cdot n$$

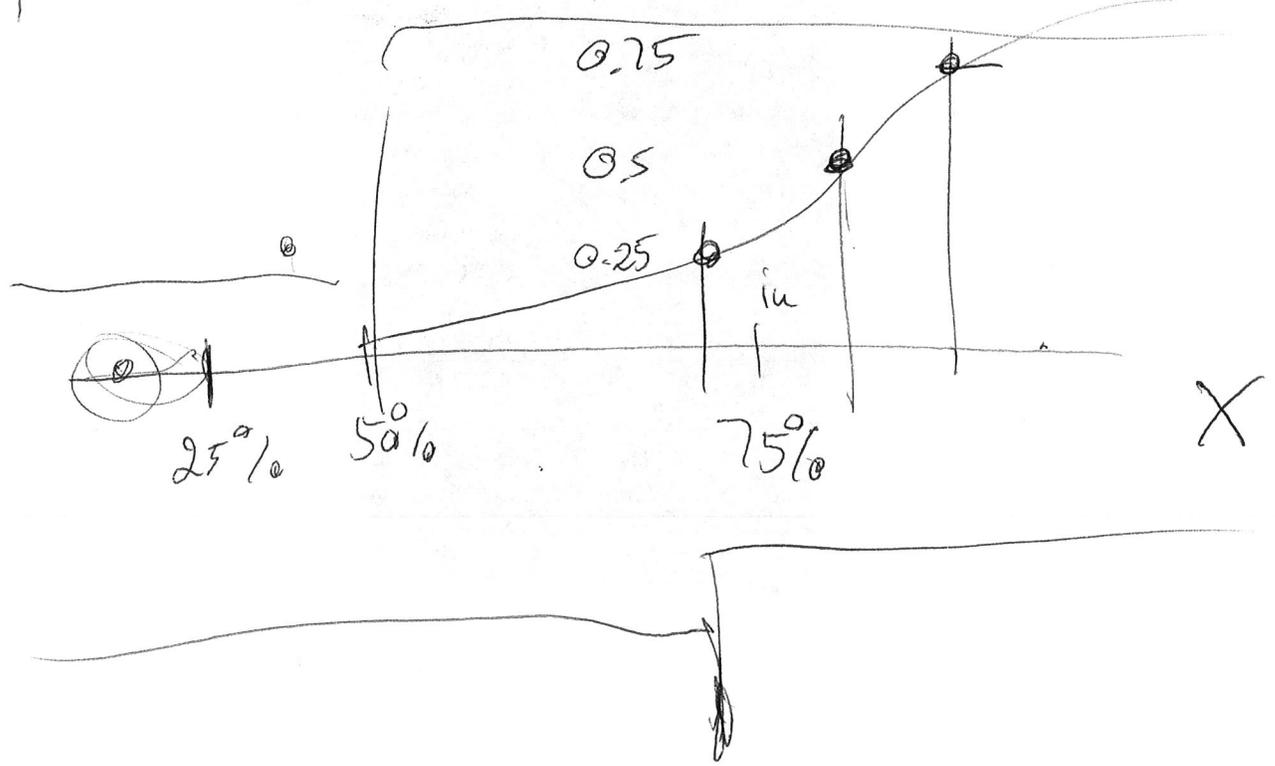| Pred | Obs | | | Obs | rough-cor-ccf md |
|---|---|---|---|---|---|
| 1 | 1 | dop | | 1 | $\begin{pmatrix} 0.3 & 0.7 & 0 \end{pmatrix}$ |
| | | rough. | | | |
| | | ccf | | | |
| | ⋮ | ccf | | | |
| 10 | 1 | rough_cor | | | |
| 1 | 2 | ccf | | | |
| | ⋮ | | | | |
| 10 | 2 | ccf | | | |

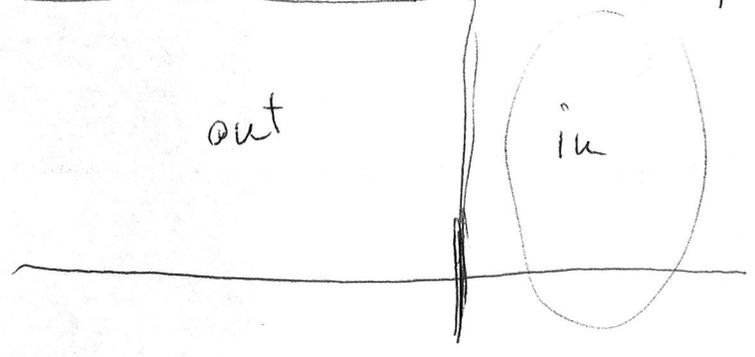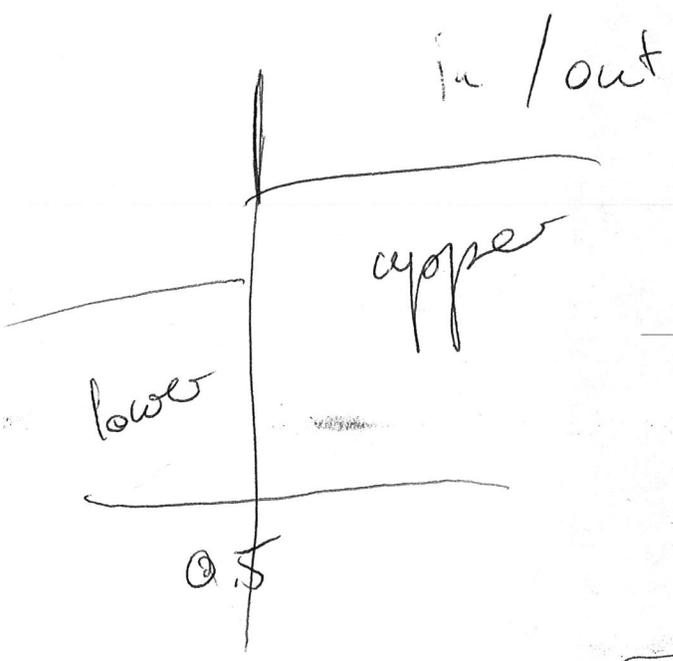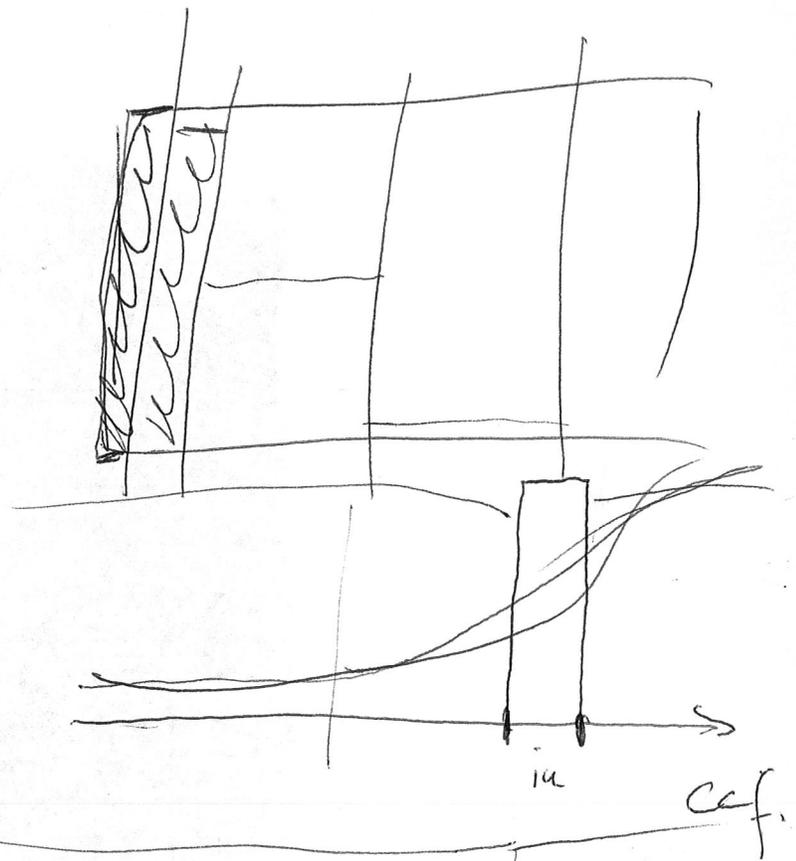$$p_i = (p_{i1} \quad \text{---} \quad p_{iq}) \sim Mult$$

$$\exp\left(-\sum_{k=1}^{q} p_{ik} \log p_{ik}\right)$$

Shannon entropy

$$\left(1 - \sum_{k=1}^{q} p_{ik}^2\right) \in [0,1]$$

$p(1-p)$

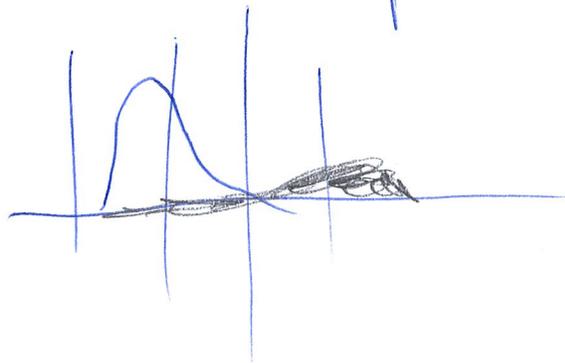$$\sum_{k=1}^{q} p_{ik} = 1$$

$ccf \gtreqless 0.5$

$ccf \geq 0.75$

in / out

upper

lower

0.5

out                    in                    ccf.

0.75

0.5

0.25                    in

25%    50%              75%              X

~~Some~~ Equalbin method

→ skew data will be problematic



Quantile method is ~~better~~

binning + ~~lasso~~ ridge regression
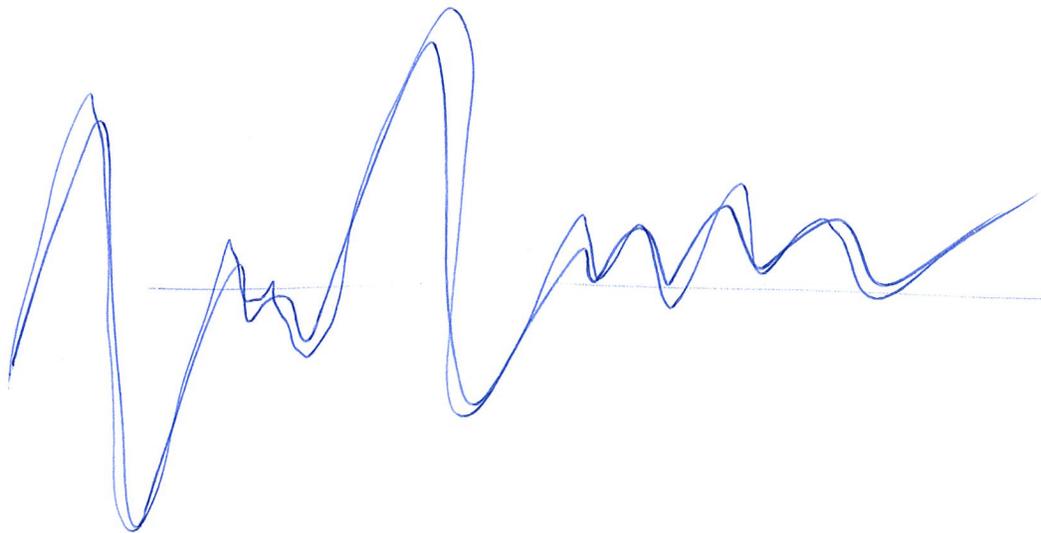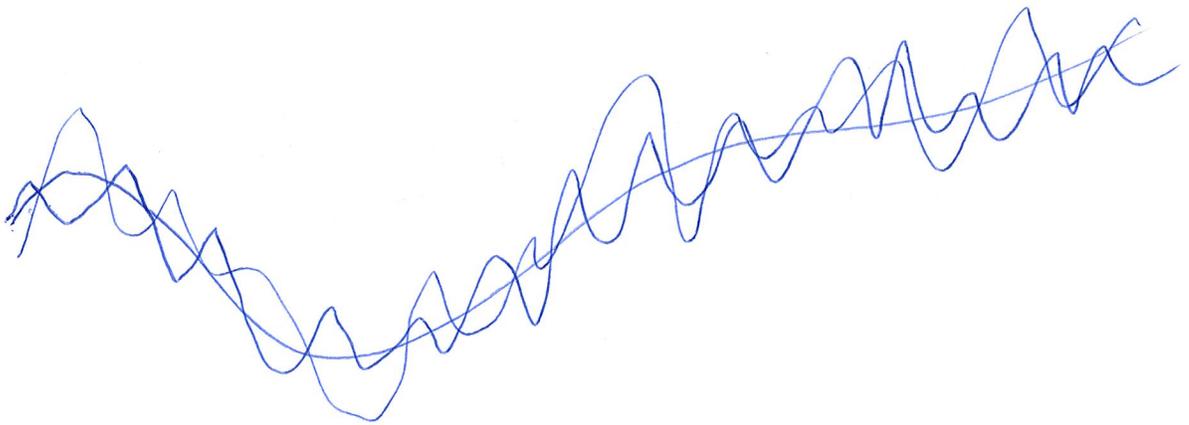
$$y_{ij} = \boxed{\sum \beta_j x_{ij}} + penalty + \varepsilon$$

case i:

$\boxed{\beta_j x_{ij}}$ for each $j = 1, \rightarrow p$

intercept is 0.8



@ 0.4

$col = rgb(1, 1, 1, alpha = 0.1)$



col      darker-col

# Visual Diagnostics of a model explainer

at the example of LIME

main objective of model explainer:
o understand and explain model performance

LIME does....

Conceptually: models at two levels:
explainer model                          — "simple"
original "black box" model — "complicated"

Usually: model predictions, maybe with
ground truth

     Type I error

     Type II error                          model is wrong

explanation is also a prediction —
    how reliable is that explanation?

Explains model has very low
$R^2$ generally
– probably due to binning

"Local" explanations are not local, but
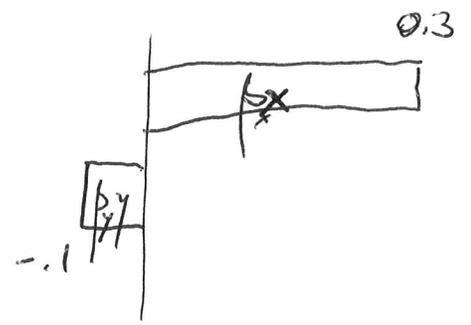are driven by the (global) marginal
distribution of covariates

Describe LIME, including details
on binning and linear regression in
binned features

Motivation

$X = 8$

$Y = 0.5$

0.3

bx

by
H

-.1

$\mu = 0.5$

$X = 8$

$Y = 0.5$

.Sty
Style file

Rmarkdown:

@bibtextag          Author (2010)

[@bibtextag]        (Author 2010 )

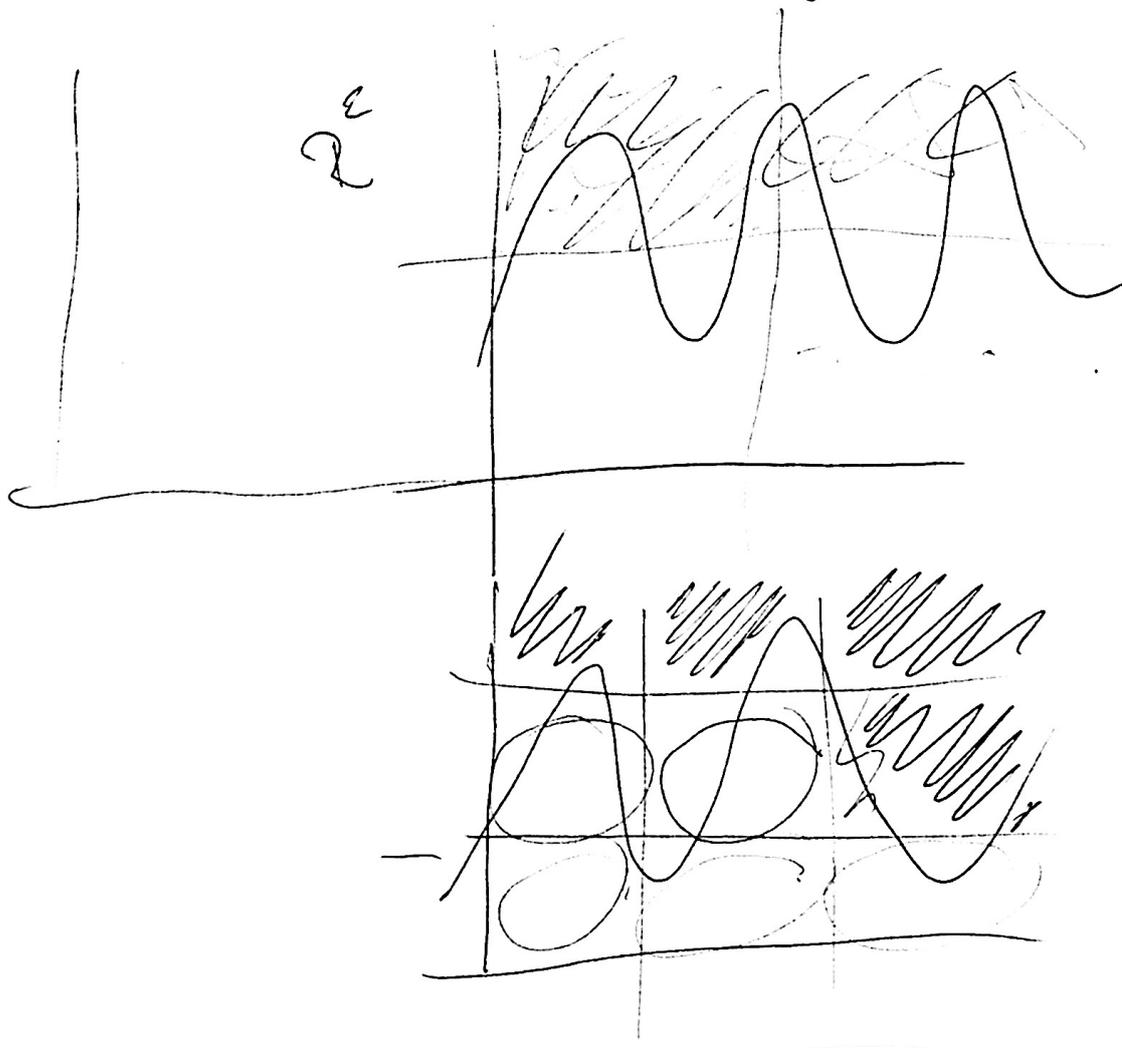[@ref1; @ref2]      (Author1 2010, Author2
                    et al 2013)

yaml:

bibliography: name.bib

$x_2$

$x_1$

low noise          high noise

$\varepsilon^2$

Local
(interpretable)
model agnostic
explanations

$4^2$ is low number
~~simple~~ ✓
$4^P$ might not be very
low
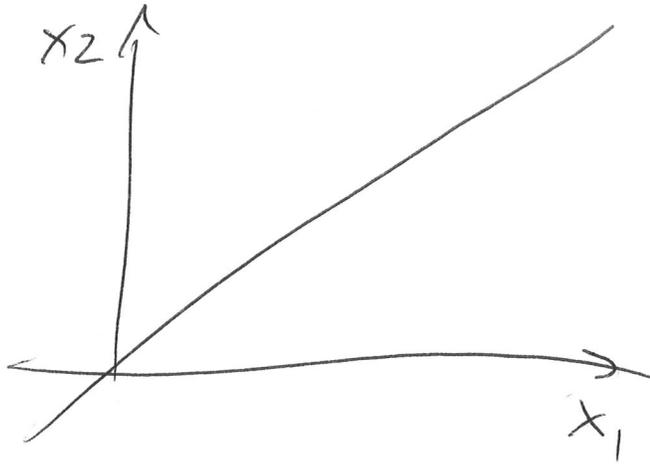
Expectations for explanations
- data driven by relevant features
- "explain": difference in coefficients

deterministic/
non-deterministic

Random forest
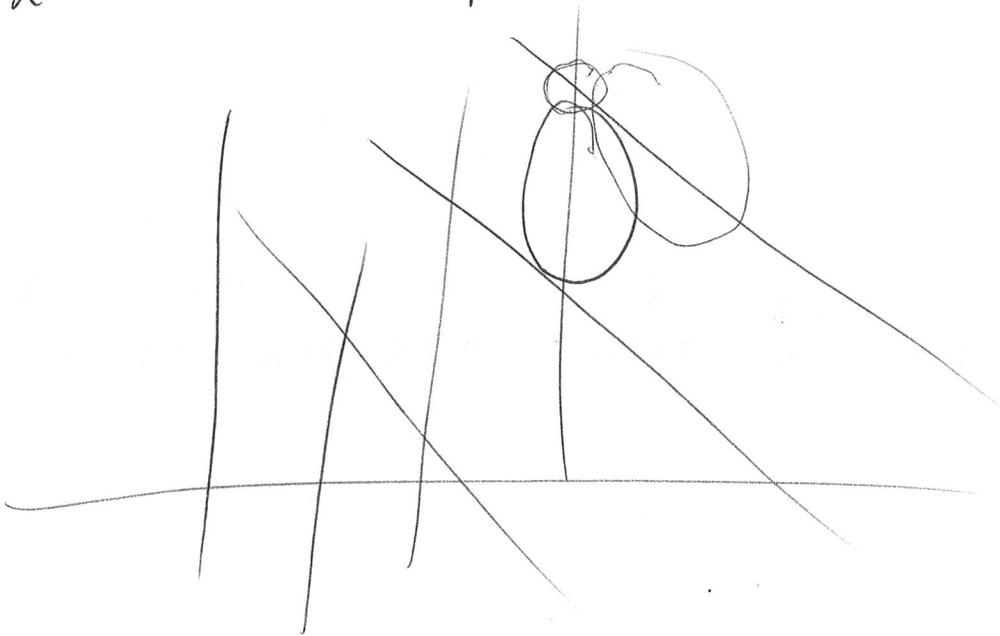variable
importance
for identify
relevant features

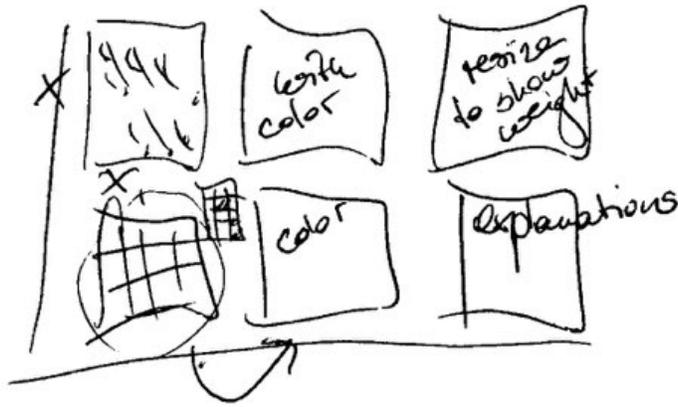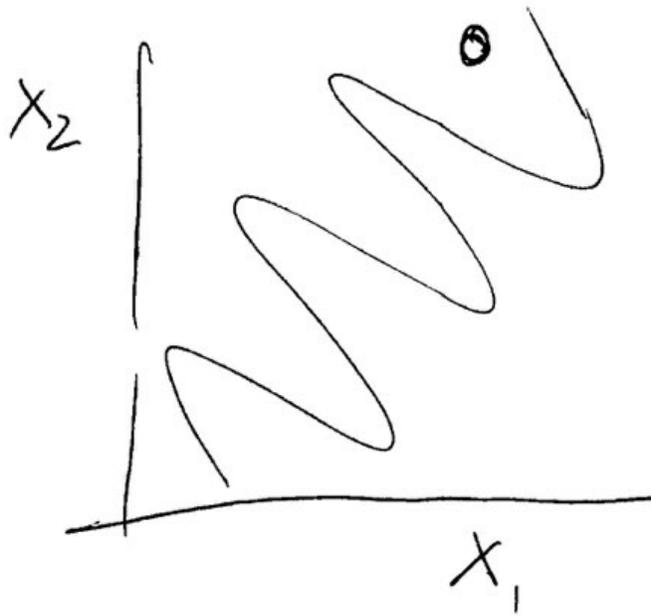Score based Likelihood
Ratio / Bayes factor

global summary



$$z = ax_1 + bx_2$$

local summaries

$X_2$

$X_1$

feature
selection
(not shown)

Create permutations

obtain RF predictions

weight (gower)

bin

weight (not gower)

fit ridge regression

~~get expla~~

Select features

re-fit ridge

get explanations

Suggestion                    Discussion

random Forest                      ✗

   importance features      ✗

$X_1, X_2$    important        ✗

Then a tree on                 E        or
$X_1, X_2$ < on weighted
simulated data of fig 7        binned versions
using rf score as $y$          of $x_1$ and $x_2$

Complexity of tree
directly related to
complexity of explanation
tradeoff between accuracy
and simplicity of explanation
$$\begin{bmatrix} < X_1 < \\ < X_2 < \end{bmatrix}$$

simple versus useful