

CADET: Interpretable Parametric Conditional Density Estimation with Decision Trees

Cyrus Cousins

Matteo Riondato

Brown University

Amherst College

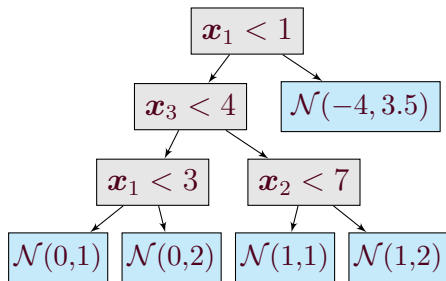
Web: cyrus_cousins@brown.edu Mail: cs.brown.edu/~ccousins/cadet/



ECMLPKDD
Würzburg | 16.–20.09.2019

“Interpretable Parametric Conditional-Density-Estimation”

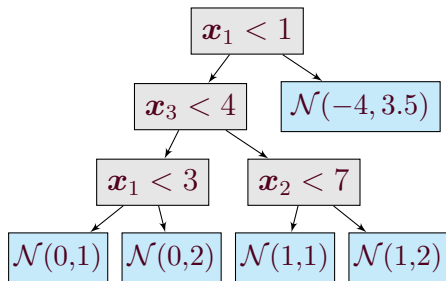
- (1) *Conditional Density Estimation*: predict *distributions* (not point-estimates)
- (2) CADET predicts *parametric densities*, e.g. GAUSSIAN(1, 1) or BETA(3, 2)
- (3) CADET *trees and predictions are interpretable*



“Interpretable Parametric Conditional-Density-Estimation”

- (1) *Conditional Density Estimation*: predict *distributions* (not point-estimates)
- (2) CADET predicts *parametric densities*, e.g. GAUSSIAN(1, 1) or BETA(3, 2)
- (3) CADET *trees and predictions are interpretable*

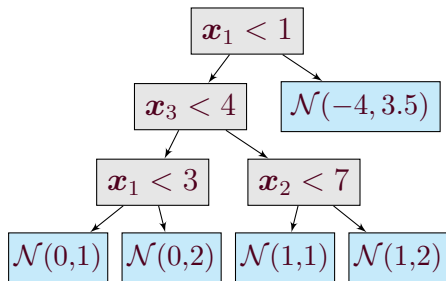
- ▶ Existing CDE tree methods
 - ▶ High training, query, and storage costs
 - ▶ Uninterpretable (non-parametric) estimates
 - ▶ High sample complexity



“Interpretable Parametric Conditional-Density-Estimation”

- (1) *Conditional Density Estimation*: predict *distributions* (not point-estimates)
- (2) CADET predicts *parametric densities*, e.g. GAUSSIAN(1, 1) or BETA(3, 2)
- (3) CADET *trees and predictions are interpretable*

- ▶ Existing CDE tree methods
 - ▶ High training, query, and storage costs
 - ▶ Uninterpretable (non-parametric) estimates
 - ▶ High sample complexity
- ▶ CADET sacrifices *representativeness* for
 - ▶ Efficient training, storage, and querying
 - ▶ Easily understood parametric estimates
 - ▶ Generalizability

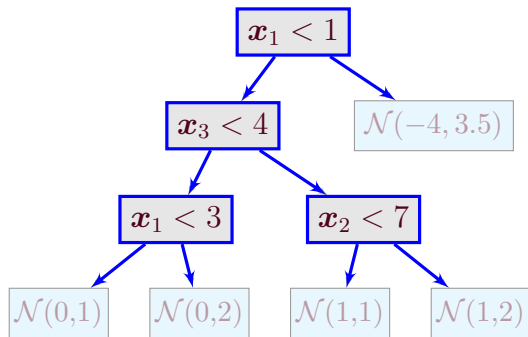


Interpretability in CDE Trees

Interpretability applies to:

(1) *Model*:

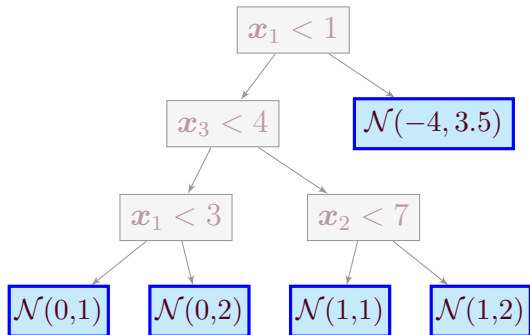
Tree structure easy to *visualize & understand*



Interpretability in CDE Trees

Interpretability applies to:

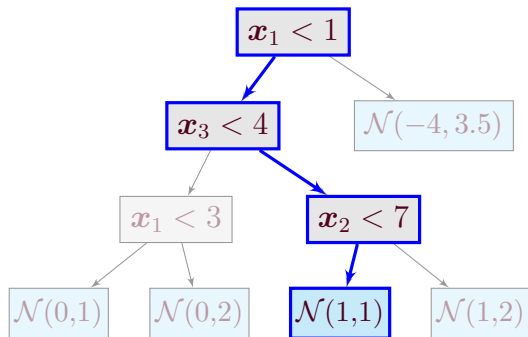
- (1) *Model*:
Tree structure easy to *visualize & understand*
- (2) *Predictions*:
Model *output* must be simple



Interpretability in CDE Trees

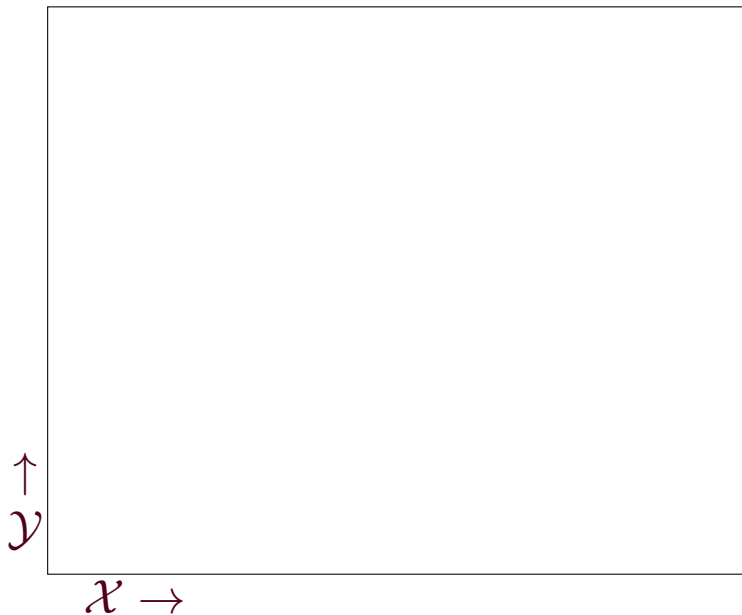
Interpretability applies to:

- (1) *Model*:
Tree structure easy to *visualize & understand*
- (2) *Predictions*:
Model *output* must be simple
- (3) *Decision process*:
Easily audit decision making process



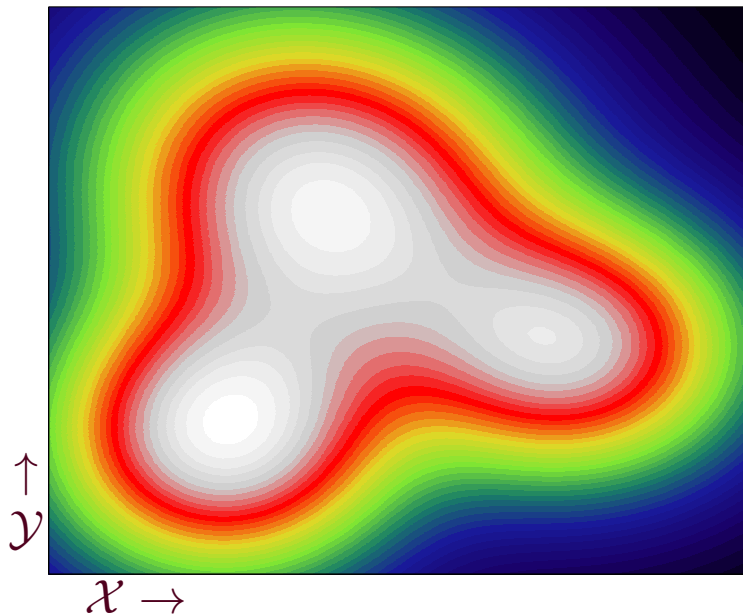
What is Conditional Density Estimation?

- ▶ Domain \mathcal{X} , codomain \mathcal{Y}



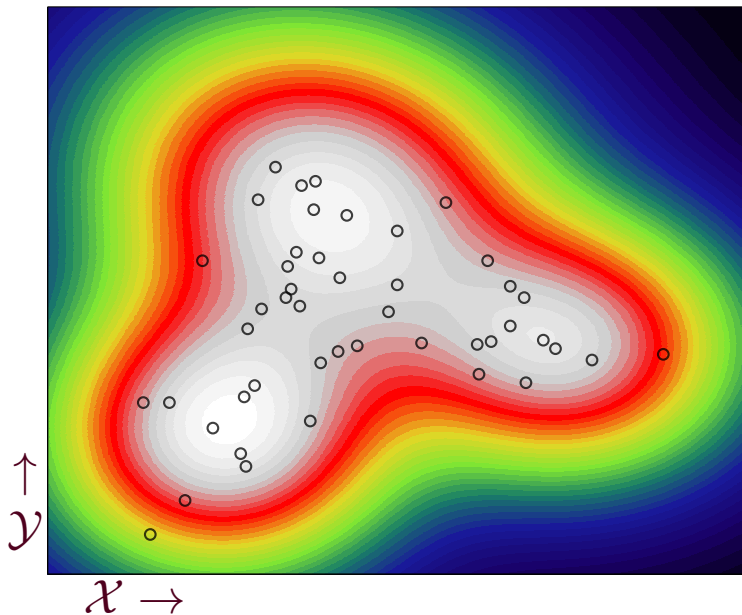
What is Conditional Density Estimation?

- ▶ Domain \mathcal{X} , codomain \mathcal{Y}
- ▶ PDF ρ over $\mathcal{X} \times \mathcal{Y}$



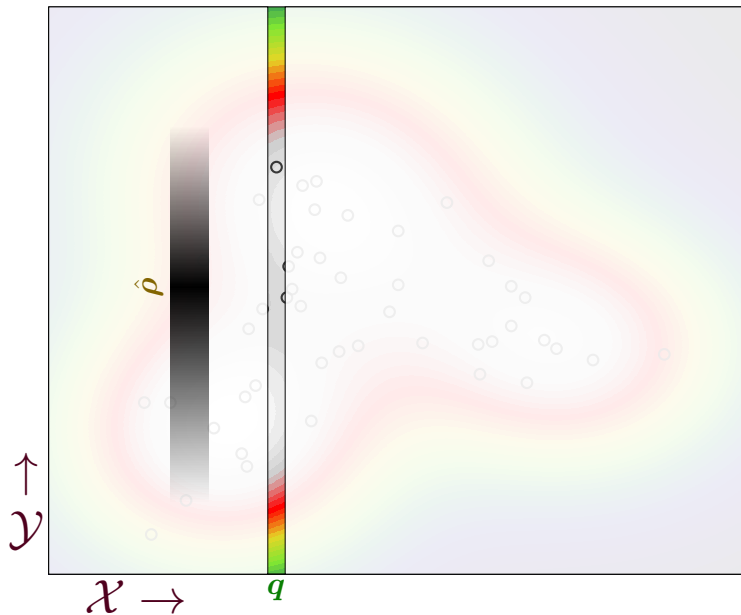
What is Conditional Density Estimation?

- ▶ Domain \mathcal{X} , codomain \mathcal{Y}
- ▶ PDF ρ over $\mathcal{X} \times \mathcal{Y}$
- ▶ Sample (\mathbf{x}, \mathbf{y}) of m points drawn with density ρ



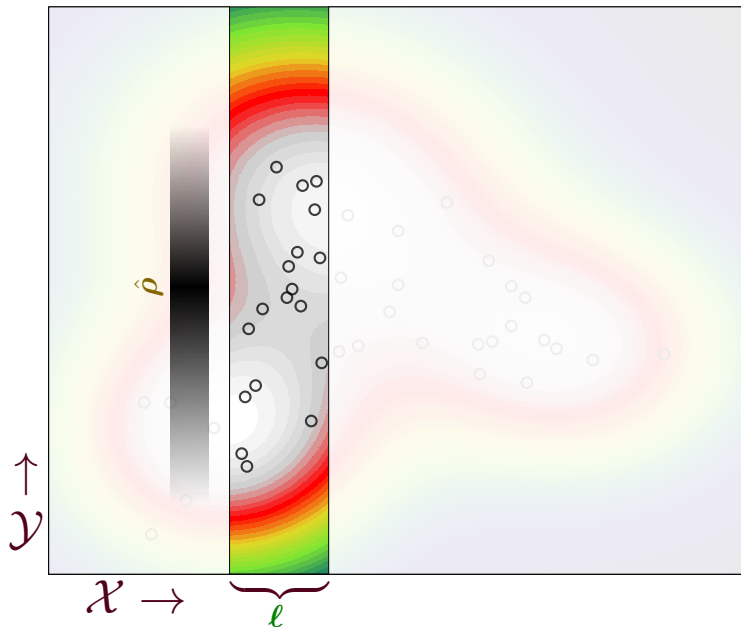
What is Conditional Density Estimation?

- ▶ Domain \mathcal{X} , codomain \mathcal{Y}
- ▶ PDF ρ over $\mathcal{X} \times \mathcal{Y}$
- ▶ Sample (\mathbf{x}, \mathbf{y}) of m points drawn with density ρ
- ▶ Condition on query $\mathbf{q} \in \mathcal{X}$
- ▶ Estimate $\hat{\rho}(\cdot | \mathbf{q}) \approx \rho(\cdot | \mathbf{q})$



What is Conditional Density Estimation?

- ▶ Domain \mathcal{X} , codomain \mathcal{Y}
- ▶ PDF ρ over $\mathcal{X} \times \mathcal{Y}$
- ▶ Sample (\mathbf{x}, \mathbf{y}) of m points drawn with density ρ
- ▶ Condition on query $\mathbf{q} \in \mathcal{X}$
- ▶ Estimate $\hat{\rho}(\cdot | \mathbf{q}) \approx \rho(\cdot | \mathbf{q})$
- ▶ Decision trees:
Fit PDF $\hat{\rho}$ to leaf $\ell \ni \mathbf{q}$



Why Conditional Density Estimation?

Supervised Learning: for any $\mathbf{q} \in \mathcal{X}$, predict statistics of $\rho(\cdot | \mathbf{q})$

Why Conditional Density Estimation?

Supervised Learning: for any $\mathbf{q} \in \mathcal{X}$, predict statistics of $\rho(\cdot | \mathbf{q})$

		Prediction Type	
		Summary	Distribution
\mathcal{Y}	Discrete	Hard Classification $\operatorname{argmax}_y \mathbb{P}(y \mathbf{q})$	Soft Classification $\mathbb{P}(\cdot \mathbf{q})$
	Continuous	Regression $\mathbb{E}_{(x,y) \sim \rho} [y x = \mathbf{q}]$	CDE $\rho(\cdot \mathbf{q})$

Why Conditional Density Estimation?

Supervised Learning: for any $\mathbf{q} \in \mathcal{X}$, predict statistics of $\rho(\cdot | \mathbf{q})$

		Prediction Type	
		Summary	Distribution
\mathcal{Y}	Discrete	Hard Classification $\operatorname{argmax}_y \mathbb{P}(y \mathbf{q})$	Soft Classification $\mathbb{P}(\cdot \mathbf{q})$
	Continuous	Regression $\mathbb{E}_{(x,y) \sim \rho} [y x = \mathbf{q}]$	CDE $\rho(\cdot \mathbf{q})$

- ▶ Regression is a lossy process
- ▶ Only estimate *average outcome*

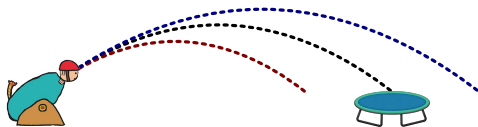


Why Conditional Density Estimation?

Supervised Learning: for any $\mathbf{q} \in \mathcal{X}$, predict statistics of $\rho(\cdot | \mathbf{q})$

		Prediction Type	
		Summary	Distribution
\mathcal{Y}	Discrete	Hard Classification $\operatorname{argmax}_y \mathbb{P}(y \mathbf{q})$	Soft Classification $\mathbb{P}(\cdot \mathbf{q})$
	Continuous	Regression $\mathbb{E} [y x = \mathbf{q}]$ <small>$(x,y) \sim \rho$</small>	CDE $\rho(\cdot \mathbf{q})$

- ▶ Regression is a lossy process
- ▶ Only estimate *average outcome*
- ▶ Want to reason about *many possibilities*

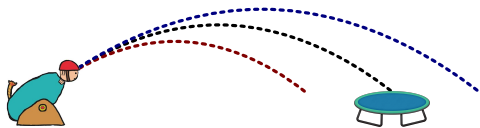


Why Conditional Density Estimation?

Supervised Learning: for any $\mathbf{q} \in \mathcal{X}$, predict statistics of $\rho(\cdot | \mathbf{q})$

		Prediction Type	
		Summary	Distribution
\mathcal{Y}	Discrete	Hard Classification $\operatorname{argmax}_y \mathbb{P}(y \mathbf{q})$	Soft Classification $\mathbb{P}(\cdot \mathbf{q})$
	Continuous	Regression $\mathbb{E} [y x = \mathbf{q}]$ <small>$(x,y) \sim \rho$</small>	CDE $\rho(\cdot \mathbf{q})$

- ▶ Regression is a lossy process
 - ▶ Only estimate *average outcome*
 - ▶ Want to reason about *many possibilities*
- ▶ CDE quantifies *uncertainty* due to noise or ambiguity
 - ▶ Generalizes *soft classification* to arbitrary \mathcal{Y}
 - ▶ Postprocess to estimate mean, median, ...



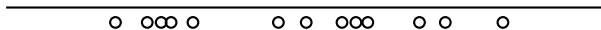
Parametric vs Nonparametric CDE Trees

- ▶ CDE with decision trees:
 - ▶ Tree splits \mathcal{X} into *disjoint cover* (leaves)
 - ▶ Estimate distribution over \mathcal{Y} at each leaf ℓ

Parametric vs Nonparametric CDE Trees

- ▶ CDE with decision trees:
 - ▶ Tree splits \mathcal{X} into *disjoint cover* (leaves)
 - ▶ Estimate distribution over \mathcal{Y} at each leaf ℓ

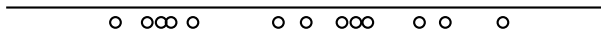
Fitting Labels at ℓ



Parametric vs Nonparametric CDE Trees

- ▶ CDE with decision trees:
 - ▶ Tree splits \mathcal{X} into *disjoint cover* (leaves)
 - ▶ Estimate distribution over \mathcal{Y} at each leaf ℓ
- ▶ **Parametric** CDE trees (CADET)
 - ▶ Quickly converge to parametric approximation
 - ▶ Learn good splits with small samples
 - ▶ Simple predictions, understood at a glance

Fitting Labels at ℓ

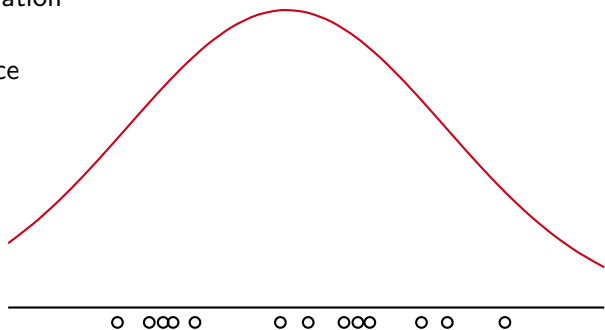


Parametric vs Nonparametric CDE Trees

- ▶ CDE with decision trees:
 - ▶ Tree splits \mathcal{X} into *disjoint cover* (leaves)
 - ▶ Estimate distribution over \mathcal{Y} at each leaf ℓ
- ▶ **Parametric** CDE trees (CADET)
 - ▶ Quickly converge to parametric approximation
 - ▶ Learn good splits with small samples
 - ▶ Simple predictions, understood at a glance

Fitting Labels at ℓ

— Gaussian (Parametric)

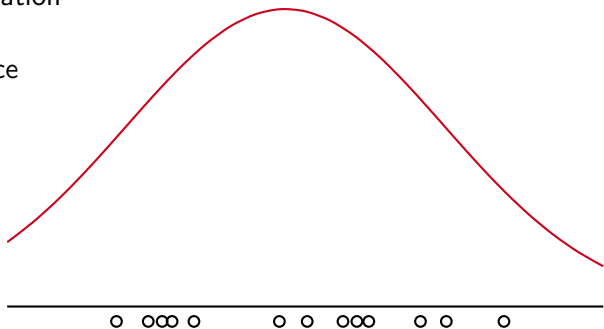


Parametric vs Nonparametric CDE Trees

- ▶ CDE with decision trees:
 - ▶ Tree splits \mathcal{X} into *disjoint cover* (leaves)
 - ▶ Estimate distribution over \mathcal{Y} at each leaf ℓ
- ▶ **Parametric** CDE trees (CADET)
 - ▶ Quickly converge to parametric approximation
 - ▶ Learn good splits with small samples
 - ▶ Simple predictions, understood at a glance
- ▶ **Nonparametric** CDE trees
 - ▶ Asymptotic consistency
 - ▶ *Eventually* get it right
 - ▶ Poor sample complexity
 - ▶ Must fit distribution *at each leaf*
 - ▶ More susceptible to *overfitting*

Fitting Labels at ℓ

— Gaussian (Parametric)

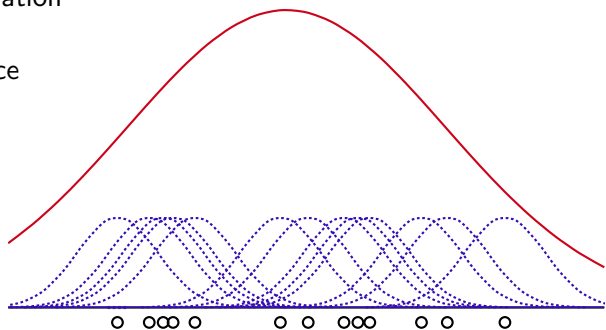


Parametric vs Nonparametric CDE Trees

- ▶ CDE with decision trees:
 - ▶ Tree splits \mathcal{X} into *disjoint cover* (leaves)
 - ▶ Estimate distribution over \mathcal{Y} at each leaf ℓ
- ▶ **Parametric** CDE trees (CADET)
 - ▶ Quickly converge to parametric approximation
 - ▶ Learn good splits with small samples
 - ▶ Simple predictions, understood at a glance
- ▶ **Nonparametric** CDE trees
 - ▶ Asymptotic consistency
 - ▶ *Eventually* get it right
 - ▶ Poor sample complexity
 - ▶ Must fit distribution *at each leaf*
 - ▶ More susceptible to *overfitting*

Fitting Labels at ℓ

- Gaussian (Parametric)
- ⋯ KDE Components

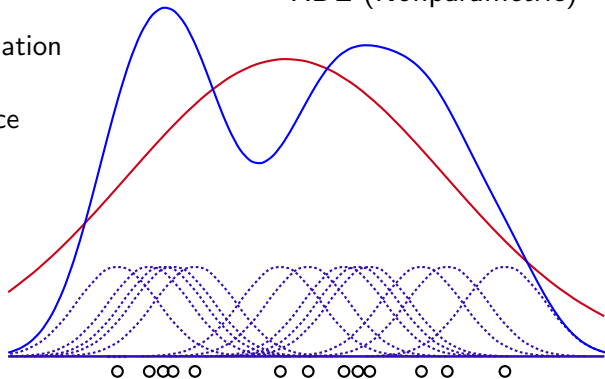


Parametric vs Nonparametric CDE Trees

- ▶ CDE with decision trees:
 - ▶ Tree splits \mathcal{X} into *disjoint cover* (leaves)
 - ▶ Estimate distribution over \mathcal{Y} at each leaf ℓ
- ▶ **Parametric** CDE trees (CADET)
 - ▶ Quickly converge to parametric approximation
 - ▶ Learn good splits with small samples
 - ▶ Simple predictions, understood at a glance
- ▶ **Nonparametric** CDE trees
 - ▶ Asymptotic consistency
 - ▶ *Eventually* get it right
 - ▶ Poor sample complexity
 - ▶ Must fit distribution *at each leaf*
 - ▶ More susceptible to *overfitting*

Fitting Labels at ℓ

- Gaussian (Parametric)
- ⋯ KDE Components
- KDE (Nonparametric)



What are Decision Trees Again?

- ▶ Fitting optimal tree to (\mathbf{x}, \mathbf{y}) is NP-hard
Standard heuristic: impurity reduction
 - (1) Start with a singleton tree, and impurity criterion $I(\cdot)$
 - ▶ $I(\mathbf{y})$ measures *disagreement* among \mathbf{y}
 - ▶ MSE, entropy, GINI, ...

What are Decision Trees Again?

- ▶ Fitting optimal tree to (\mathbf{x}, \mathbf{y}) is NP-hard

Standard heuristic: impurity reduction

- (1) Start with a singleton tree, and impurity criterion $I(\cdot)$

- ▶ $I(\mathbf{y})$ measures *disagreement* among \mathbf{y}
- ▶ MSE, entropy, GINI, ...

- (2) Select split of (\mathbf{x}, \mathbf{y}) into $(\mathbf{x}_L, \mathbf{y}_L)$ and $(\mathbf{x}_R, \mathbf{y}_R)$ to maximize *impurity reduction*:

$$(m_L + m_R) I(\mathbf{y}) - (m_L I(\mathbf{y}_L) + m_R I(\mathbf{y}_R))$$

What are Decision Trees Again?

- ▶ Fitting optimal tree to (\mathbf{x}, \mathbf{y}) is NP-hard

Standard heuristic: impurity reduction

- (1) Start with a singleton tree, and impurity criterion $I(\cdot)$

- ▶ $I(\mathbf{y})$ measures *disagreement* among \mathbf{y}
- ▶ MSE, entropy, GINI, ...

- (2) Select split of (\mathbf{x}, \mathbf{y}) into $(\mathbf{x}_L, \mathbf{y}_L)$ and $(\mathbf{x}_R, \mathbf{y}_R)$ to maximize *impurity reduction*:

$$(m_L + m_R) I(\mathbf{y}) - (m_L I(\mathbf{y}_L) + m_R I(\mathbf{y}_R))$$

- (3) Repeat until termination condition is met

- ▶ Maximum depth, minimum samples per leaf, ...

What are Decision Trees Again?

- ▶ Fitting optimal tree to (\mathbf{x}, \mathbf{y}) is NP-hard

Standard heuristic: impurity reduction

- (1) Start with a singleton tree, and impurity criterion $I(\cdot)$

- ▶ $I(\mathbf{y})$ measures *disagreement* among \mathbf{y}
- ▶ MSE, entropy, GINI, ...

- (2) Select split of (\mathbf{x}, \mathbf{y}) into $(\mathbf{x}_L, \mathbf{y}_L)$ and $(\mathbf{x}_R, \mathbf{y}_R)$ to maximize *impurity reduction*:

$$(m_L + m_R) I(\mathbf{y}) - (m_L I(\mathbf{y}_L) + m_R I(\mathbf{y}_R))$$

- (3) Repeat until termination condition is met

- ▶ Maximum depth, minimum samples per leaf, ...

- ▶ Lower impurity \implies leaf label more accurately describes \mathbf{y}

CADET Trees

- ▶ Hyperparameter: *parametric distribution family* \mathcal{F} over \mathcal{Y}
 - ▶ All CADET predictions (distributions) belong to \mathcal{F}
 - ▶ For this talk, \mathcal{F} in *exponential class*

CADET Trees

- ▶ Hyperparameter: *parametric distribution family* \mathcal{F} over \mathcal{Y}
 - ▶ All CADET predictions (distributions) belong to \mathcal{F}
 - ▶ For this talk, \mathcal{F} in *exponential class*
- ▶ Given query point \mathbf{q} , CADET trees
 - (1) Find the leaf that contains \mathbf{q} , with training labels \mathbf{y}
 - (2) Return $\text{MLE}_{\mathcal{F}}(\mathbf{y})$

CADET Trees

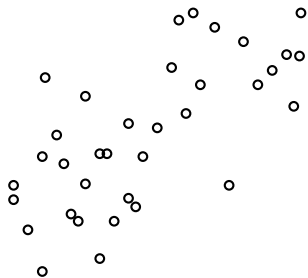
- ▶ Hyperparameter: *parametric distribution family* \mathcal{F} over \mathcal{Y}
 - ▶ All CADET predictions (distributions) belong to \mathcal{F}
 - ▶ For this talk, \mathcal{F} in *exponential class*
- ▶ Given query point \mathbf{q} , CADET trees
 - (1) Find the leaf that contains \mathbf{q} , with training labels \mathbf{y}
 - (2) Return $\text{MLE}_{\mathcal{F}}(\mathbf{y})$
- ▶ CADET trees optimize *cross entropy impurity*
 - ▶ Evaluate CDE with *cross entropy loss* $\ell_{\text{CE}}(\mathbf{y} | \hat{\rho}) \doteq -\ln \hat{\rho}(\mathbf{y})$
 - ▶ *Cross entropy impurity* $I_{\mathcal{F}}(\mathbf{y}) \doteq \frac{1}{m} \sum_{i=1}^m \ell_{\text{CE}}(\mathbf{y}_i | \text{MLE}(\mathbf{y}; \mathcal{F}))$

CADET Trees

- ▶ Hyperparameter: *parametric distribution family* \mathcal{F} over \mathcal{Y}
 - ▶ All CADET predictions (distributions) belong to \mathcal{F}
 - ▶ For this talk, \mathcal{F} in *exponential class*
- ▶ Given query point \mathbf{q} , CADET trees
 - (1) Find the leaf that contains \mathbf{q} , with training labels \mathbf{y}
 - (2) Return $\text{MLE}_{\mathcal{F}}(\mathbf{y})$
- ▶ CADET trees optimize *cross entropy impurity*
 - ▶ Evaluate CDE with *cross entropy loss* $\ell_{\text{CE}}(\mathbf{y} | \hat{\rho}) \doteq -\ln \hat{\rho}(\mathbf{y})$
 - ▶ *Cross entropy impurity* $I_{\mathcal{F}}(\mathbf{y}) \doteq \frac{1}{m} \sum_{i=1}^m \ell_{\text{CE}}(\mathbf{y}_i | \text{MLE}(\mathbf{y}; \mathcal{F}))$
- ▶ CADET organizes computation such that:
 - ▶ Evaluating impurity reduction requires *constant work*
 - ▶ Leaves require *constant storage*

Training CADET Trees

Cross Entropy Impurity:
$$I_{\mathcal{F}}(\mathbf{y}) \doteq \frac{1}{m} \sum_{i=1}^m \ell_{\text{CE}}(\mathbf{y}_i | \text{MLE}(\mathbf{y}; \mathcal{F}))$$



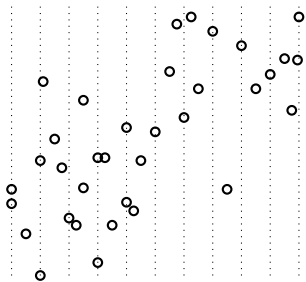
(1) Given training

points $\mathbf{x}_{i=1}^m \in \mathbb{R}^m$ (x axis)

labels $\mathbf{y}_{i=1}^m \in \mathbb{R}^m$ (y axis)

Training CADET Trees

Cross Entropy Impurity:
$$I_{\mathcal{F}}(\mathbf{y}) \doteq \frac{1}{m} \sum_{i=1}^m \ell_{\text{CE}}(\mathbf{y}_i | \text{MLE}(\mathbf{y}; \mathcal{F}))$$



(1) Given training

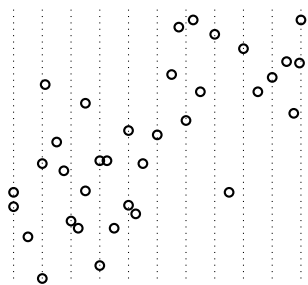
points $\mathbf{x}_{i=1}^m \in \mathbb{R}^m$ (x axis)

labels $\mathbf{y}_{i=1}^m \in \mathbb{R}^m$ (y axis)

(2) Evaluate possible splits

Training CADET Trees

Cross Entropy Impurity:
$$I_{\mathcal{F}}(\mathbf{y}) \doteq \frac{1}{m} \sum_{i=1}^m \ell_{\text{CE}}(\mathbf{y}_i | \text{MLE}(\mathbf{y}; \mathcal{F}))$$

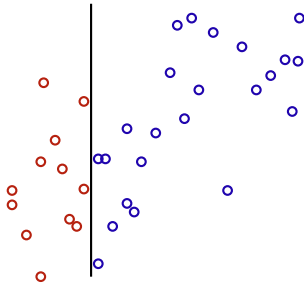


(1) Given training

points $\mathbf{x}_{i=1}^m \in \mathbb{R}^m$ (x axis)

labels $\mathbf{y}_{i=1}^m \in \mathbb{R}^m$ (y axis)

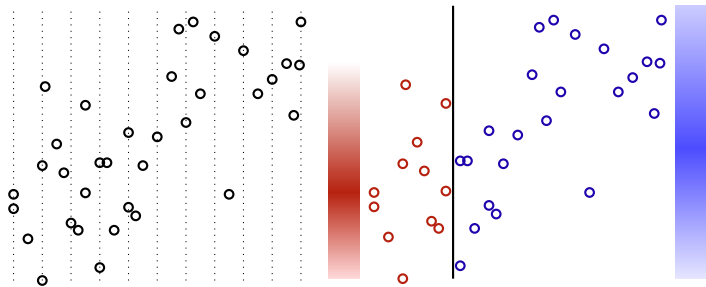
(2) Evaluate possible splits



(1) Consider split $\mathbf{y}_L, \mathbf{y}_R$

Training CADET Trees

Cross Entropy Impurity:
$$I_{\mathcal{F}}(\mathbf{y}) \doteq \frac{1}{m} \sum_{i=1}^m \ell_{\text{CE}}(\mathbf{y}_i | \text{MLE}(\mathbf{y}; \mathcal{F}))$$



(1) Given training

points $\mathbf{x}_{i=1}^m \in \mathbb{R}^m$ (x axis)
labels $\mathbf{y}_{i=1}^m \in \mathbb{R}^m$ (y axis)

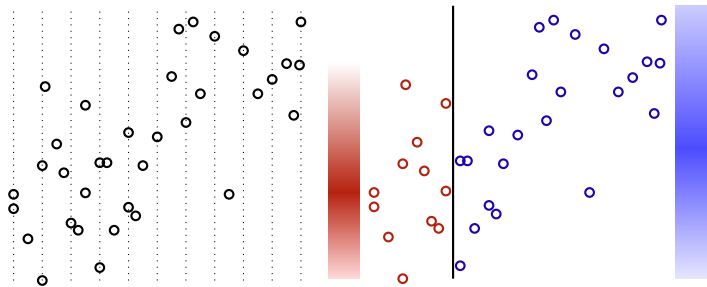
(2) Evaluate possible splits

(1) Consider split $\mathbf{y}_L, \mathbf{y}_R$

(2) Fit $\text{MLE}_{\mathcal{F}}(\mathbf{y}_L), \text{MLE}_{\mathcal{F}}(\mathbf{y}_R)$

Training CADET Trees

Cross Entropy Impurity:
$$I_{\mathcal{F}}(\mathbf{y}) \doteq \frac{1}{m} \sum_{i=1}^m \ell_{\text{CE}}(\mathbf{y}_i | \text{MLE}(\mathbf{y}; \mathcal{F}))$$



(1) Given training

points $\mathbf{x}_{i=1}^m \in \mathbb{R}^m$ (x axis)

labels $\mathbf{y}_{i=1}^m \in \mathbb{R}^m$ (y axis)

(2) Evaluate possible splits

(1) Consider split $\mathbf{y}_L, \mathbf{y}_R$

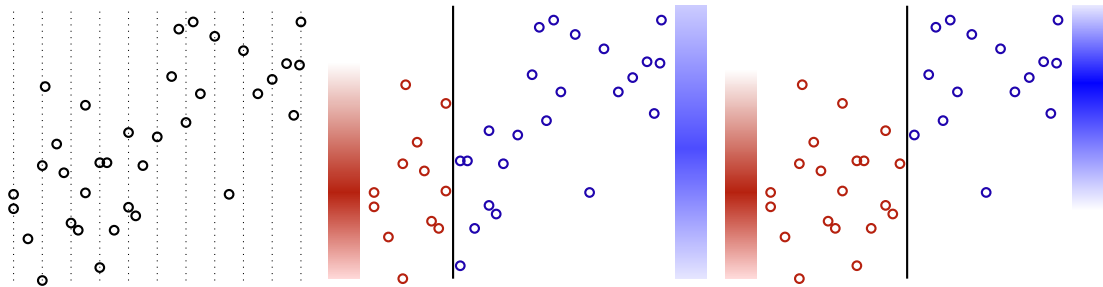
(2) Fit $\text{MLE}_{\mathcal{F}}(\mathbf{y}_L), \text{MLE}_{\mathcal{F}}(\mathbf{y}_R)$

(3) $I_{\mathcal{F}}(\mathbf{y}_L)$: low

$I_{\mathcal{F}}(\mathbf{y}_R)$: high

Training CADET Trees

Cross Entropy Impurity:
$$I_{\mathcal{F}}(\mathbf{y}) \doteq \frac{1}{m} \sum_{i=1}^m \ell_{\text{CE}}(\mathbf{y}_i | \text{MLE}(\mathbf{y}; \mathcal{F}))$$



(1) Given training

points $\mathbf{x}_{i=1}^m \in \mathbb{R}^m$ (x axis)

labels $\mathbf{y}_{i=1}^m \in \mathbb{R}^m$ (y axis)

(2) Evaluate possible splits

(1) Consider split $\mathbf{y}_L, \mathbf{y}_R$

(2) Fit $\text{MLE}_{\mathcal{F}}(\mathbf{y}_L), \text{MLE}_{\mathcal{F}}(\mathbf{y}_R)$

(3) $I_{\mathcal{F}}(\mathbf{y}_L)$: low

$I_{\mathcal{F}}(\mathbf{y}_R)$: high

(1) Now consider $\mathbf{y}'_L, \mathbf{y}'_R$

(2) Fit $\text{MLE}_{\mathcal{F}}(\mathbf{y}'_L), \text{MLE}_{\mathcal{F}}(\mathbf{y}'_R)$

(3) $I_{\mathcal{F}}(\mathbf{y}'_L)$: same

$I_{\mathcal{F}}(\mathbf{y}'_R)$: lower

Why CDE-Specific Impurity Criteria?

CADET sounds complicated, why not just use $I_{\text{MSE}}(\cdot)$?

Why CDE-Specific Impurity Criteria?

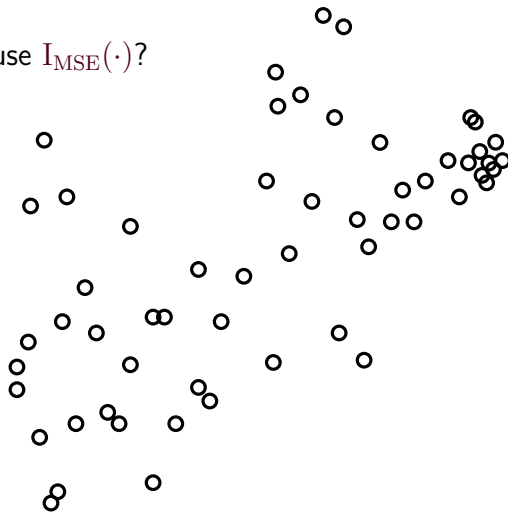
CADET sounds complicated, why not just use $I_{\text{MSE}}(\cdot)$?

- ▶ Consider maximizing $I_{\text{MSE}}(\cdot)$ -reduction
 - ▶ Only sensitive to changes in *expectation*
 - ▶ Insensitive to changes in *variance*
 - ▶ ⚠ Problem even for Gaussian family

Why CDE-Specific Impurity Criteria?

CADET sounds complicated, why not just use $I_{\text{MSE}}(\cdot)$?

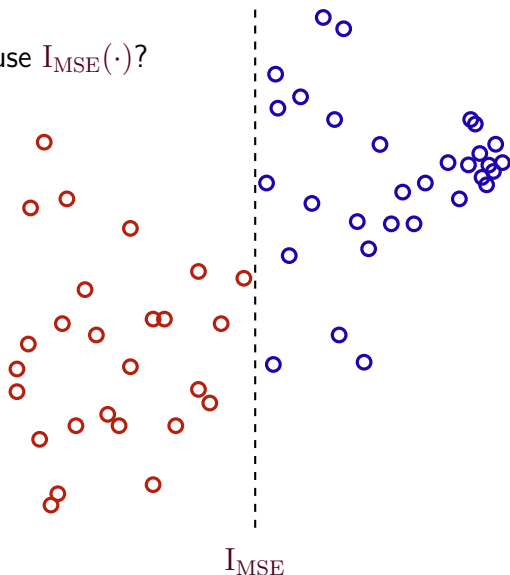
- ▶ Consider maximizing $I_{\text{MSE}}(\cdot)$ -reduction
 - ▶ Only sensitive to changes in *expectation*
 - ▶ Insensitive to changes in *variance*
 - ▶ ⚠ Problem even for Gaussian family



Why CDE-Specific Impurity Criteria?

CADET sounds complicated, why not just use $I_{\text{MSE}}(\cdot)$?

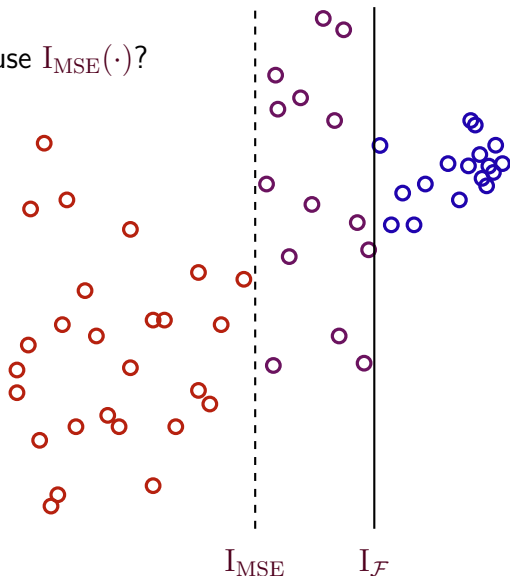
- ▶ Consider maximizing $I_{\text{MSE}}(\cdot)$ -reduction
 - ▶ Only sensitive to changes in *expectation*
 - ▶ Insensitive to changes in *variance*
 - ▶ ⚠ Problem even for Gaussian family



Why CDE-Specific Impurity Criteria?

CADET sounds complicated, why not just use $I_{\text{MSE}}(\cdot)$?

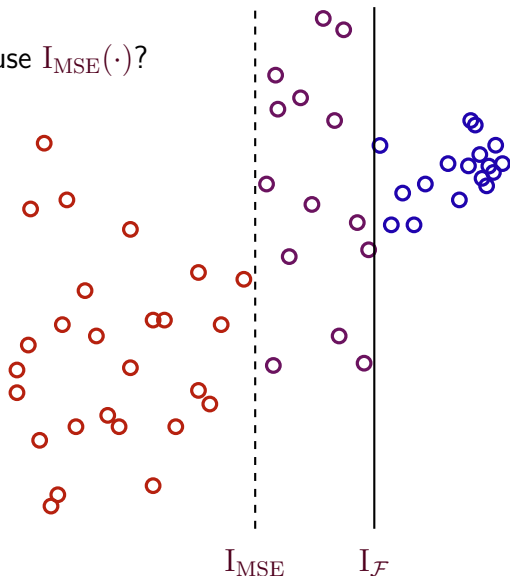
- ▶ Consider maximizing $I_{\text{MSE}}(\cdot)$ -reduction
 - ▶ Only sensitive to changes in *expectation*
 - ▶ Insensitive to changes in *variance*
 - ▶ ⚠ Problem even for Gaussian family



Why CDE-Specific Impurity Criteria?

CADET sounds complicated, why not just use $I_{\text{MSE}}(\cdot)$?

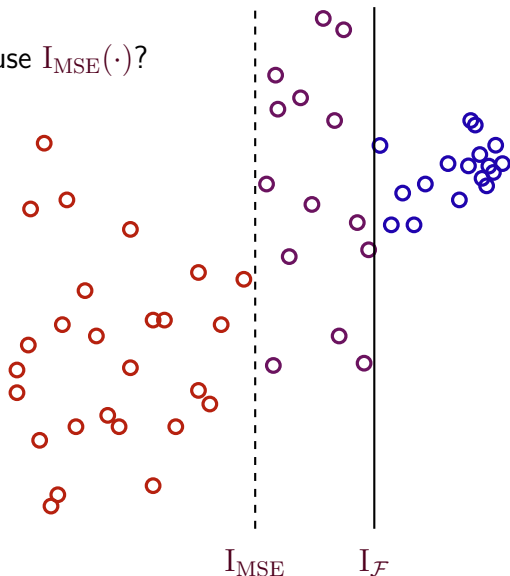
- ▶ Consider maximizing $I_{\text{MSE}}(\cdot)$ -reduction
 - ▶ Only sensitive to changes in *expectation*
 - ▶ Insensitive to changes in *variance*
 - ▶ ⚠ Problem even for Gaussian family
 - ▶ $I_{\text{MSE}}(\cdot)$ undefined for $\mathcal{Y} \neq \mathbb{R}^d$



Why CDE-Specific Impurity Criteria?

CADET sounds complicated, why not just use $I_{\text{MSE}}(\cdot)$?

- ▶ Consider maximizing $I_{\text{MSE}}(\cdot)$ -reduction
 - ▶ Only sensitive to changes in *expectation*
 - ▶ Insensitive to changes in *variance*
 - ▶ ⚠ Problem even for Gaussian family
 - ▶ $I_{\text{MSE}}(\cdot)$ undefined for $\mathcal{Y} \neq \mathbb{R}^d$
- ▶ *Takeaway*: $I_{\mathcal{F}}(\cdot)$ depends on \mathcal{F}
 - ▶ “Pick the splits that improve the fits”



Organizing Computation with Sufficient Statistics

- ▶ A leaf needs training labels \mathbf{y} to select splits in training and answer queries
 - ▶ CART *summarize* \mathbf{y} with *class frequencies* or *means*
 - ▶ “Incremental updates” make split search fast
 - ▶ CADET needs $I_{\mathcal{F}}(\mathbf{y})$ to train and $\text{MLE}_{\mathcal{F}}(\mathbf{y})$ to query

Organizing Computation with Sufficient Statistics

- ▶ A leaf needs training labels \mathbf{y} to select splits in training and answer queries
 - ▶ CART *summarize* \mathbf{y} with *class frequencies* or *means*
 - ▶ “Incremental updates” make split search fast
 - ▶ CADET needs $I_{\mathcal{F}}(\mathbf{y})$ to train and $\text{MLE}_{\mathcal{F}}(\mathbf{y})$ to query
- ▶ *Sufficient statistics* $\mathbf{w}^{(m)}(\mathbf{y})$ w.r.t. \mathcal{F} summarize \mathbf{y} in \mathcal{Y}^m
 - ▶ Compute $\text{MLE}_{\mathcal{F}}(\mathbf{y})$ and minimize $I_{\mathcal{F}}(\mathbf{y})$ from $\mathbf{w}^{(m)}(\mathbf{y})$

Organizing Computation with Sufficient Statistics

- ▶ A leaf needs training labels \mathbf{y} to select splits in training and answer queries
 - ▶ CART *summarize* \mathbf{y} with *class frequencies* or *means*
 - ▶ “Incremental updates” make split search fast
 - ▶ CADET needs $I_{\mathcal{F}}(\mathbf{y})$ to train and $\text{MLE}_{\mathcal{F}}(\mathbf{y})$ to query
- ▶ *Sufficient statistics* $w^{(m)}(\mathbf{y})$ w.r.t. \mathcal{F} summarize \mathbf{y} in \mathcal{Y}^m
 - ▶ Compute $\text{MLE}_{\mathcal{F}}(\mathbf{y})$ and minimize $I_{\mathcal{F}}(\mathbf{y})$ from $w^{(m)}(\mathbf{y})$

Family \mathcal{F}	Suff. Stat. $w^{(m)}(\mathbf{y})$	Log Density $\ln \rho(\mathbf{y})$
GAUSSIAN(μ, σ^2)	$\sum_{i=1}^m \mathbf{y}_i, \sum_{i=1}^m \mathbf{y}_i^2$	$-\frac{1}{2} \ln(2\pi\sigma^2) - \frac{y^2 - y\mu + \mu^2}{2\sigma^2}$

Organizing Computation with Sufficient Statistics

- ▶ A leaf needs training labels \mathbf{y} to select splits in training and answer queries
 - ▶ CART summarize \mathbf{y} with *class frequencies* or *means*
 - ▶ “Incremental updates” make split search fast
 - ▶ CADET needs $I_{\mathcal{F}}(\mathbf{y})$ to train and $\text{MLE}_{\mathcal{F}}(\mathbf{y})$ to query
- ▶ Sufficient statistics $w^{(m)}(\mathbf{y})$ w.r.t. \mathcal{F} summarize \mathbf{y} in \mathcal{Y}^m
 - ▶ Compute $\text{MLE}_{\mathcal{F}}(\mathbf{y})$ and minimize $I_{\mathcal{F}}(\mathbf{y})$ from $w^{(m)}(\mathbf{y})$

Family \mathcal{F}	Suff. Stat. $w^{(m)}(\mathbf{y})$	Log Density $\ln \rho(\mathbf{y})$
GAUSSIAN(μ, σ^2)	$\sum_{i=1}^m \mathbf{y}_i, \sum_{i=1}^m \mathbf{y}_i^2$	$-\frac{1}{2} \ln(2\pi\sigma^2) - \frac{y^2 - y\mu + \mu^2}{2\sigma^2}$
GAMMA(α, β)	$\sum_{i=1}^m \mathbf{y}_i, \sum_{i=1}^m \ln(\mathbf{y}_i)$	$\alpha \ln(\beta) - \ln \Gamma(\alpha) - y + (\alpha - 1) \ln(y)$

Organizing Computation with Sufficient Statistics

- ▶ A leaf needs training labels \mathbf{y} to select splits in training and answer queries
 - ▶ CART summarize \mathbf{y} with *class frequencies* or *means*
 - ▶ “Incremental updates” make split search fast
 - ▶ CADET needs $I_{\mathcal{F}}(\mathbf{y})$ to train and $\text{MLE}_{\mathcal{F}}(\mathbf{y})$ to query
- ▶ Sufficient statistics $w^{(m)}(\mathbf{y})$ w.r.t. \mathcal{F} summarize \mathbf{y} in \mathcal{Y}^m
 - ▶ Compute $\text{MLE}_{\mathcal{F}}(\mathbf{y})$ and minimize $I_{\mathcal{F}}(\mathbf{y})$ from $w^{(m)}(\mathbf{y})$

Family \mathcal{F}	Suff. Stat. $w^{(m)}(\mathbf{y})$	Log Density $\ln \rho(\mathbf{y})$
GAUSSIAN(μ, σ^2)	$\sum_{i=1}^m \mathbf{y}_i, \sum_{i=1}^m \mathbf{y}_i^2$	$-\frac{1}{2} \ln(2\pi\sigma^2) - \frac{y^2 - y\mu + \mu^2}{2\sigma^2}$
GAMMA(α, β)	$\sum_{i=1}^m \mathbf{y}_i, \sum_{i=1}^m \ln(\mathbf{y}_i)$	$\alpha \ln(\beta) - \ln \Gamma(\alpha) - y + (\alpha - 1) \ln(y)$

- ▶ Always exist $w(\cdot)$ for \mathcal{F} in the *exponential class* s.t.

$$w^{(m_L+m_R)}(\mathbf{y}_L \circ \mathbf{y}_R) \doteq w^{(m_L)}(\mathbf{y}_L) + w^{(m_R)}(\mathbf{y}_R)$$

Organizing Computation with Sufficient Statistics

- ▶ A leaf needs training labels \mathbf{y} to select splits in training and answer queries
 - ▶ CART summarize \mathbf{y} with *class frequencies* or *means*
 - ▶ “Incremental updates” make split search fast
 - ▶ CADET needs $I_{\mathcal{F}}(\mathbf{y})$ to train and $\text{MLE}_{\mathcal{F}}(\mathbf{y})$ to query
- ▶ Sufficient statistics $w^{(m)}(\mathbf{y})$ w.r.t. \mathcal{F} summarize \mathbf{y} in \mathcal{Y}^m
 - ▶ Compute $\text{MLE}_{\mathcal{F}}(\mathbf{y})$ and minimize $I_{\mathcal{F}}(\mathbf{y})$ from $w^{(m)}(\mathbf{y})$

Family \mathcal{F}	Suff. Stat. $w^{(m)}(\mathbf{y})$	Log Density $\ln \rho(\mathbf{y})$
GAUSSIAN(μ, σ^2)	$\sum_{i=1}^m \mathbf{y}_i, \sum_{i=1}^m \mathbf{y}_i^2$	$-\frac{1}{2} \ln(2\pi\sigma^2) - \frac{y^2 - y\mu + \mu^2}{2\sigma^2}$
GAMMA(α, β)	$\sum_{i=1}^m \mathbf{y}_i, \sum_{i=1}^m \ln(\mathbf{y}_i)$	$\alpha \ln(\beta) - \ln \Gamma(\alpha) - y + (\alpha - 1) \ln(y)$

- ▶ Always exist $w(\cdot)$ for \mathcal{F} in the *exponential class* s.t.

$$w^{(m_L+m_R)}(\mathbf{y}_L \circ \mathbf{y}_R) \doteq w^{(m_L)}(\mathbf{y}_L) + w^{(m_R)}(\mathbf{y}_R)$$

- ▶ Time to compute impurity reduction:

- ▶ With $w(\cdot)$: $\mathcal{O}(1)$ amortized time

- ▶ Without $w(\cdot)$: $\mathcal{O}(m)$ time

CADET as a Unifying Framework

- ▶ Cross-entropy impurity criterion $I_{\mathcal{F}}(\cdot)$ tailored to \mathcal{F}

CADET as a Unifying Framework

- ▶ Cross-entropy impurity criterion $I_{\mathcal{F}}(\cdot)$ tailored to \mathcal{F}

\mathcal{F}		$I_{\mathcal{F}}(\mathbf{y}) \equiv$	Tree Model
GAUSSIAN($\cdot, 1$)		$I_{\text{MSE}}(\mathbf{y})$	Regression Tree

CADET as a Unifying Framework

- ▶ Cross-entropy impurity criterion $I_{\mathcal{F}}(\cdot)$ tailored to \mathcal{F}

\mathcal{F}	$I_{\mathcal{F}}(\mathbf{y}) \equiv$	Tree Model
GAUSSIAN($\cdot, 1$)	$I_{\text{MSE}}(\mathbf{y})$	Regression Tree
CATEGORICAL(\cdot)	$I_{\text{H}}(\mathbf{y})$	Information-Gain Tree

CADET as a Unifying Framework

- ▶ Cross-entropy impurity criterion $I_{\mathcal{F}}(\cdot)$ tailored to \mathcal{F}

\mathcal{F}	$I_{\mathcal{F}}(\mathbf{y}) \equiv$	Tree Model
GAUSSIAN($\cdot, 1$)	$I_{\text{MSE}}(\mathbf{y})$	Regression Tree
CATEGORICAL(\cdot)	$I_{\text{H}}(\mathbf{y})$	Information-Gain Tree

- ▶ We've reconstructed two models from the 80s...

CADET as a Unifying Framework

- ▶ Cross-entropy impurity criterion $I_{\mathcal{F}}(\cdot)$ tailored to \mathcal{F}

\mathcal{F}	$I_{\mathcal{F}}(\mathbf{y}) \equiv$	Tree Model
GAUSSIAN($\cdot, 1$)	$I_{\text{MSE}}(\mathbf{y})$	Regression Tree
CATEGORICAL(\cdot)	$I_{\text{H}}(\mathbf{y})$	Information-Gain Tree

- ▶ We've reconstructed two models from the 80s...

- ▶ Two underlying philosophies for split selection

- (1) *Maximum likelihood*, maximize *sum-log-likelihood* of \mathbf{y}
- (2) *Minimax entropy*, minimize *uncertainty* of predictions

$$I_{\mathcal{F}}(\mathbf{y}) = H(\mathbf{y}, \text{MLE}_{\mathcal{F}}(\mathbf{y}))$$
$$I_{\text{H}, \mathcal{F}}(\mathbf{y}) = H(\mathbf{y})$$

CADET as a Unifying Framework

- ▶ Cross-entropy impurity criterion $I_{\mathcal{F}}(\cdot)$ tailored to \mathcal{F}

\mathcal{F}	$I_{\mathcal{F}}(\mathbf{y}) \equiv$	Tree Model
GAUSSIAN($\cdot, 1$)	$I_{\text{MSE}}(\mathbf{y})$	Regression Tree
CATEGORICAL(\cdot)	$I_{\text{H}}(\mathbf{y})$	Information-Gain Tree

- ▶ We've reconstructed two models from the 80s...

- ▶ Two underlying philosophies for split selection

(1) *Maximum likelihood*, maximize *sum-log-likelihood* of \mathbf{y}

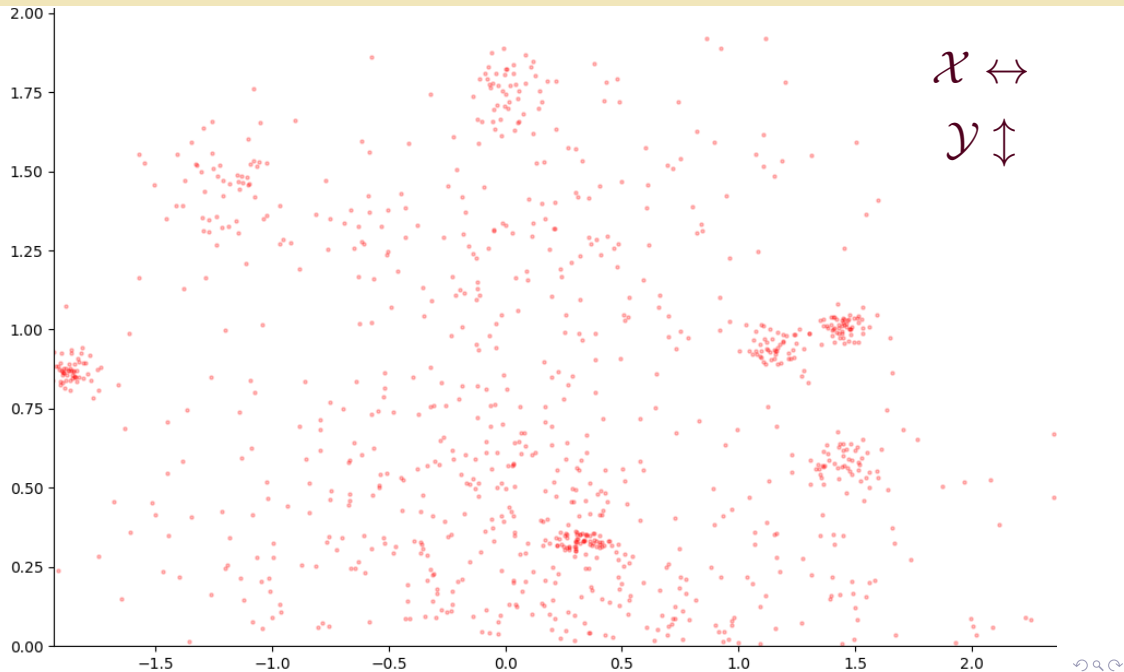
(2) *Minimax entropy*, minimize *uncertainty* of predictions

$$I_{\mathcal{F}}(\mathbf{y}) = H(\mathbf{y}, \text{MLE}_{\mathcal{F}}(\mathbf{y}))$$

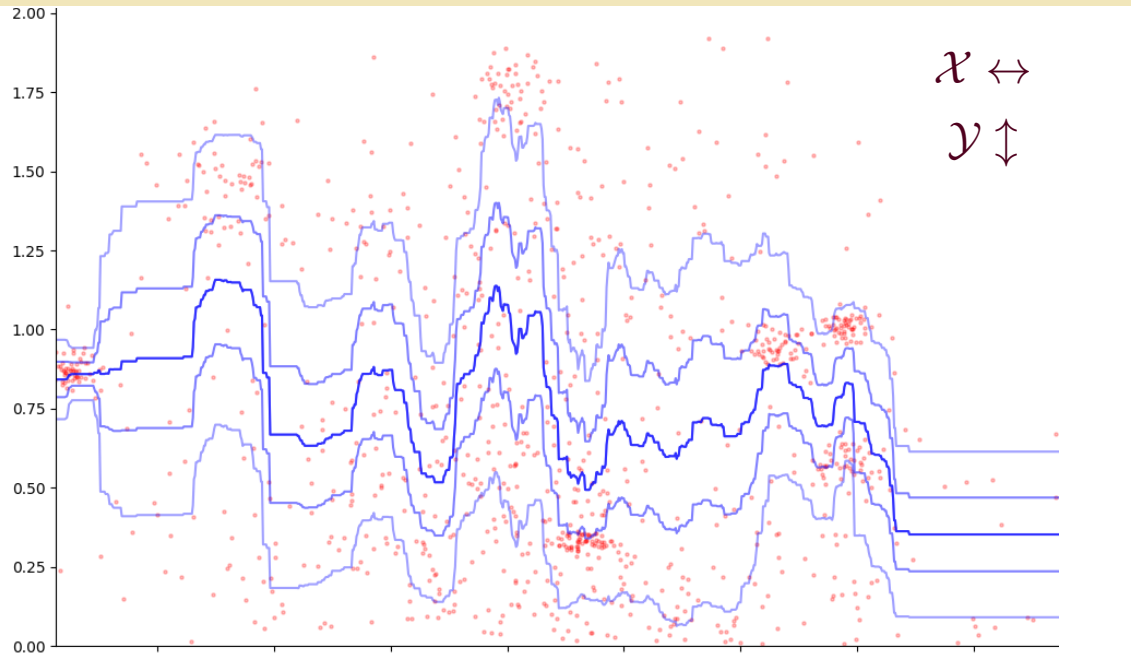
$$I_{\text{H}, \mathcal{F}}(\mathbf{y}) = H(\mathbf{y})$$

- ▶ Lemma 1: conditions on \mathcal{F} under which $I_{\mathcal{F}}(\cdot) = I_{\text{H}, \mathcal{F}}(\cdot)$

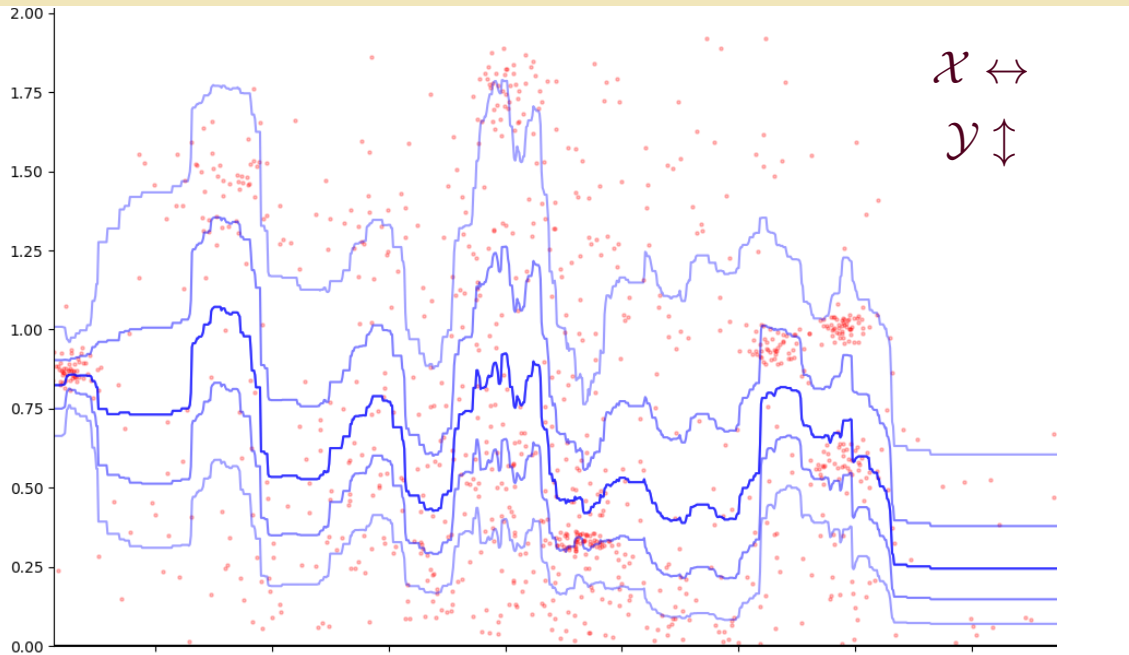
Visualizing Gaussian-CADET Forests



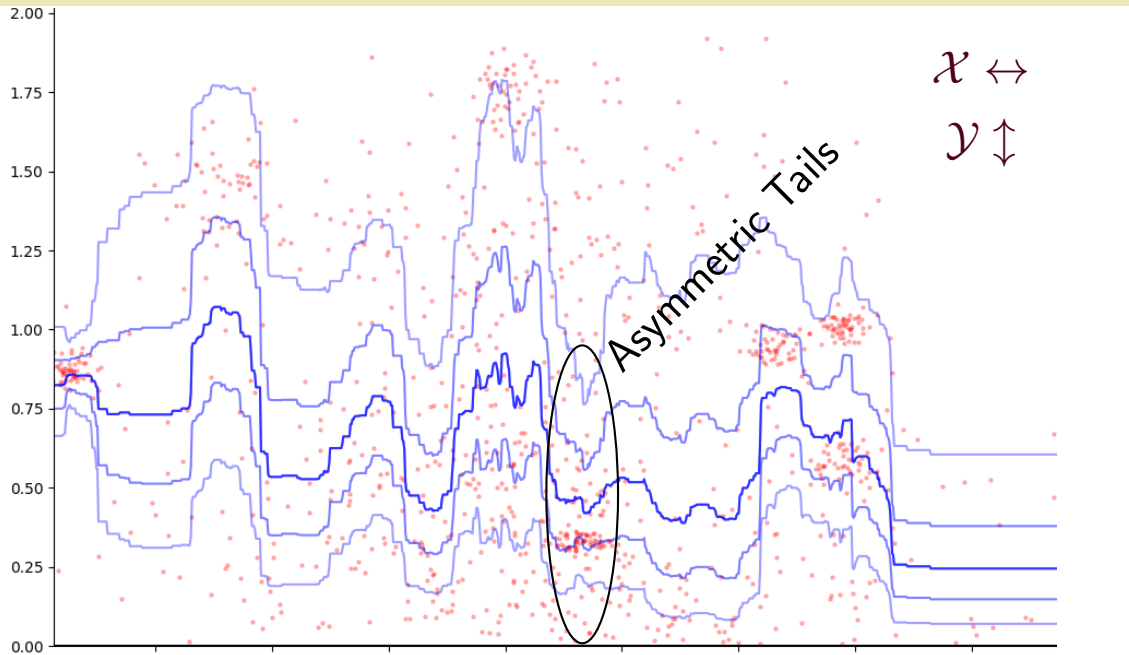
Visualizing Gaussian-CADET Forests



Visualizing Gamma-CADET Forests



Visualizing Gamma-CADET Forests



A Brief Recapitulation

- ▶ CADET: simple, interpretable, parametric CDE trees
 - ▶ Nonparametric methods uninterpretable
- ▶ Efficient training, query, and storage costs with
 - ▶ Additive sufficient statistics
 - ▶ Efficiency matches CART
 - ▶ $\Omega(m)$ speedup over nonparametric CDE trees
- ▶ Generalize existing tree methods
 - ▶ Information-gain classification trees
 - ▶ MSE regression trees