# Ensuring More Accurate, Generalisable, and Interpretable Machine Learning Models for Bioinformatics

Wandrille Duchemin, Van Du, Markus Müller

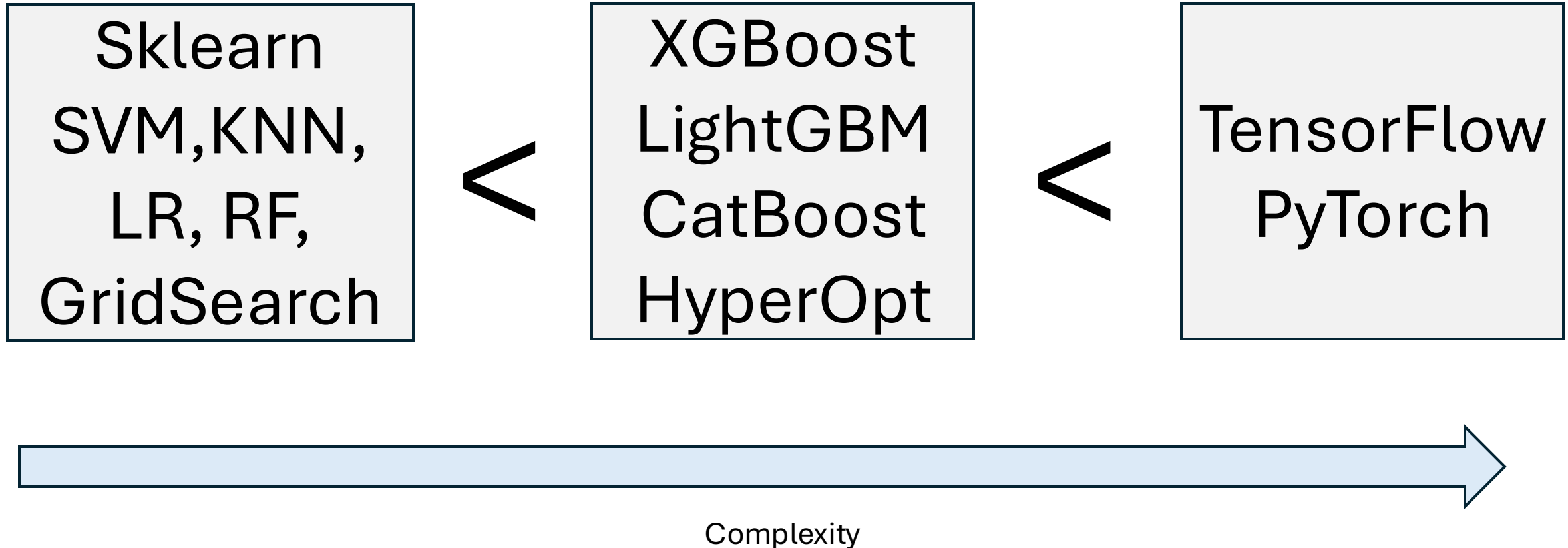Swiss Institute of Bioinformatics

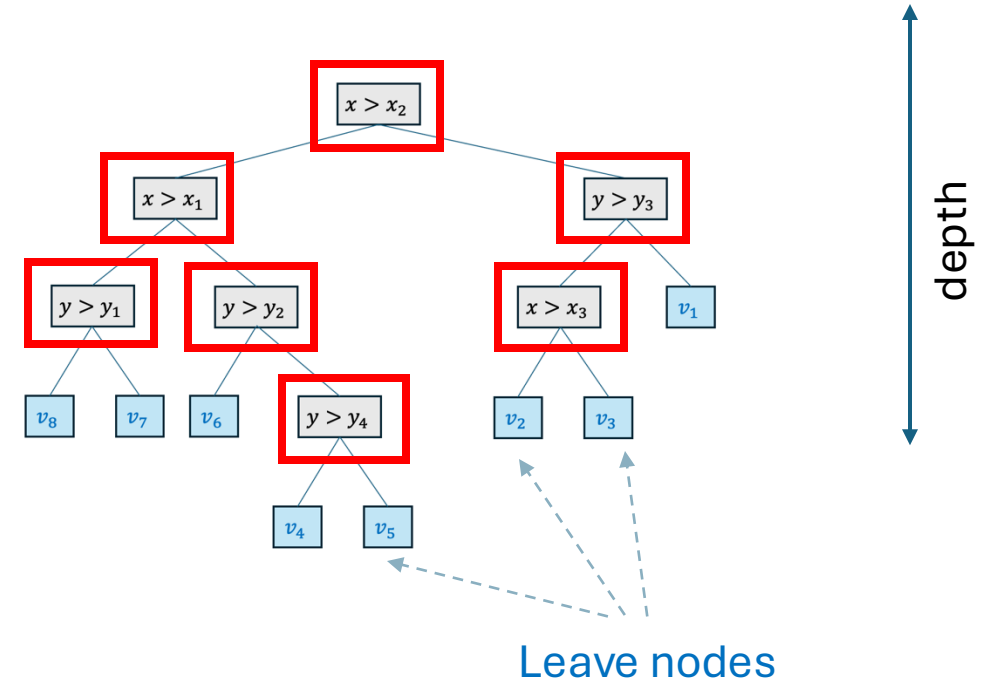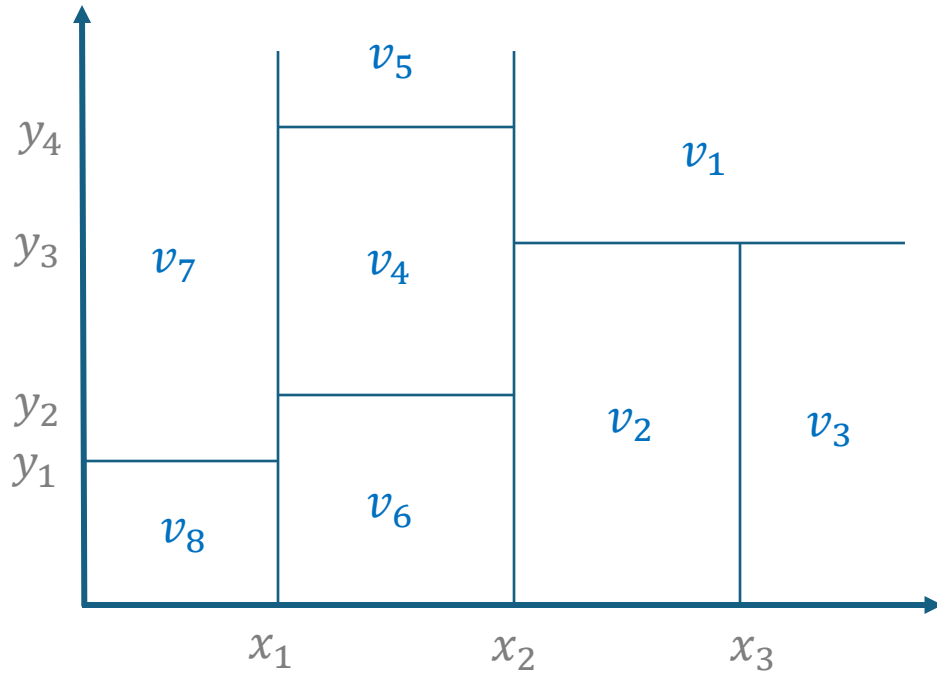Swiss Institute of
Bioinformatics

# Outline

- Intermediate machine learning (ML)

- Trees

- Bagging and boosting

- Classification and regression with XGBoost

- Hyperparameter optimization with Hyperopt

- Nested cross-validation

- Interpretable machine learning with LIME and SHAP

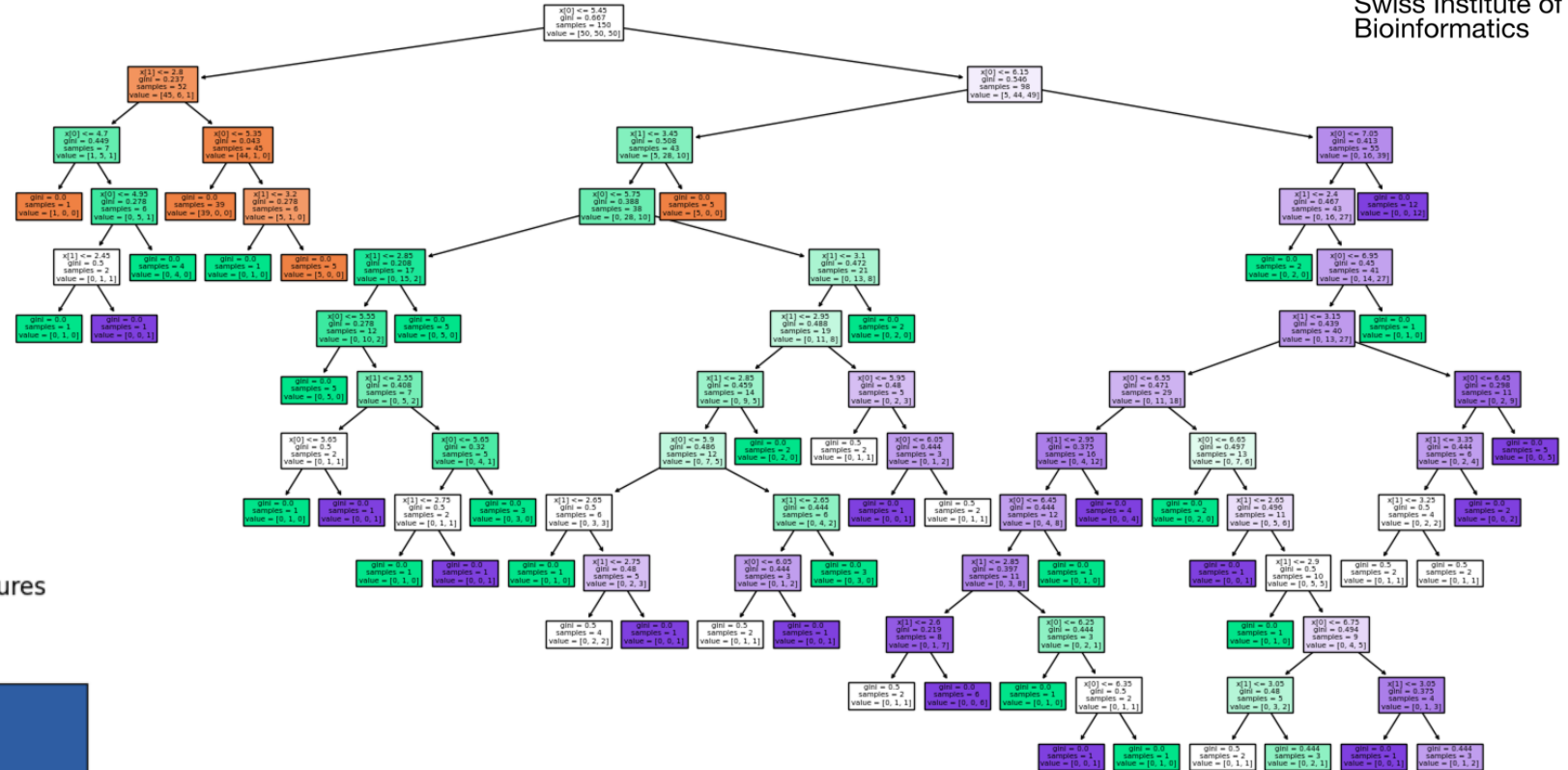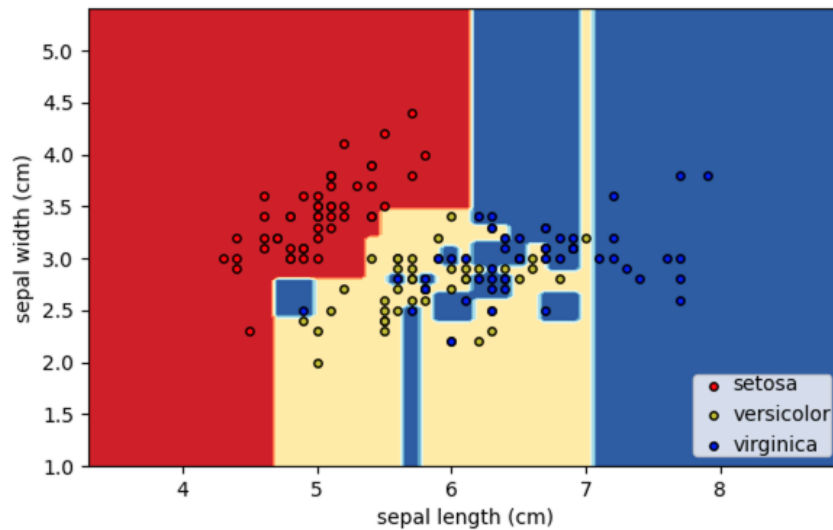# Intermediate ML

# Trees



Leave nodes

# Trees

- Regression and classification trees
- Trees are universal approximators (no bias) if the size is not restricted
- Tree's are naturally able to deal with categorical features, and are not sensitive to the distribution of features values
- Parameters:
  - For each node $k$ we have a feature $F_k$ and threshold $\theta_k$
  - For each leaf node $l$ we have a weight $v_l$ for regression, and a class label $c_l$ for classification
- Tree hyperparameters add constraints to its growth and weights:
  - max_depth: maximal number of internal nodes before leaf
  - max_leaf_nodes: maximal number of leaf nodes in the tree
  - min_samples_leaf: minimal number of data points in the training set within leaf
  - L1 or L2 regularization of weights

# Trees



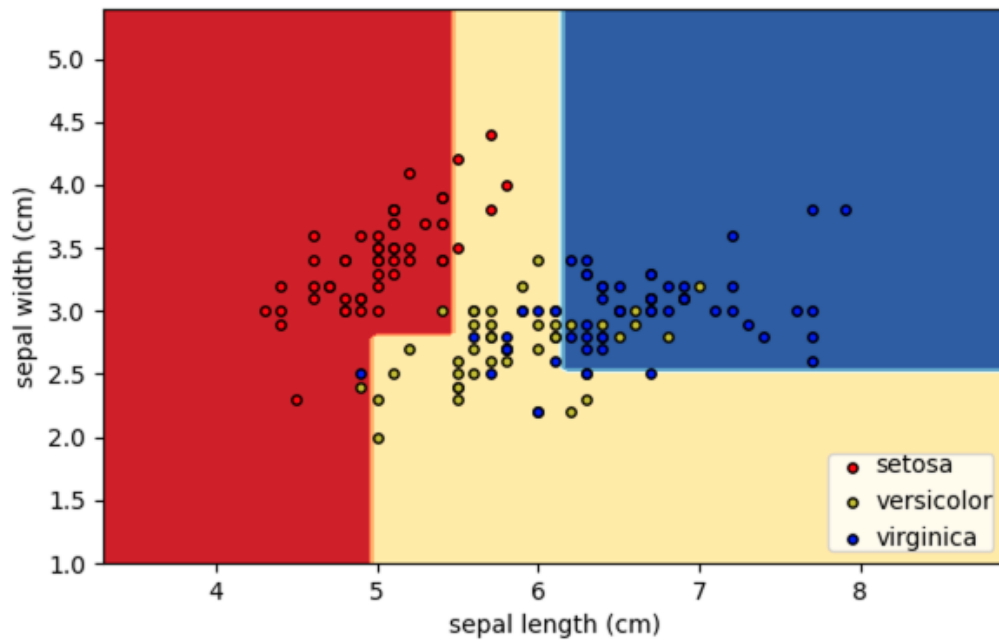Decision tree trained on all the iris sepal length and width



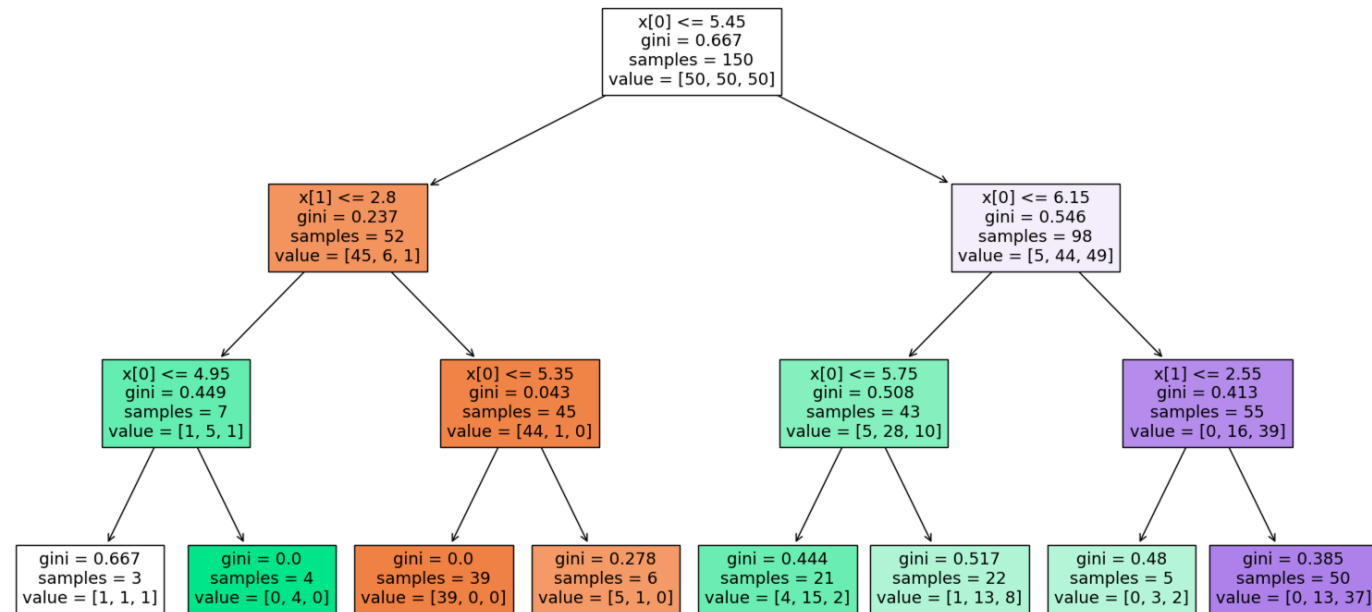Decision surface of decision trees trained on pairs of features

Markus Müller

# Tree hyperparameters

# Bias and Variance

- The *bias* error is an error that occurs if the assumptions in the learning algorithm do not match the data. High bias can cause an algorithm to misclassify a significant number of data items (underfitting). The bias does not disappear if you increase the size of your training data.

- The *variance* is an error from sensitivity of the algorithms to small fluctuations in the training set. High variance may result from an algorithm modeling the noise in the training data (overfitting). The variance becomes smaller if we increase the size of your training data.

# Bagging and Boosting

- Obtain a good classifier by combining many basic, weak classifiers with limited amount of data

- Weak classifiers have either a high bias and/or high variance

- Weak classifiers are for example linear classifiers (high bias for non-linear problems) or tree classifiers (high variance without pruning)

- Weak classifiers have to be better than random labels

# Bagging

- Bootstrap data items, train a base model on bootstrapped data, repeat $N$-times. Average predictions for regression or take a majority vote for classification.

- Reduces variance but not bias

- Bootstrapped models are correlated, which is not optimal. In the Random Forest method, features and data items are sampled in each iteration, which reduces the correlation of base models.

# Boosting

- In each iteration, the boosting method tries to correct the error of the previous iteration by focusing on the wrongly predicted data items.

- The current base model corrects those misclassified data points and is added to the previous models with a weight that is inversely proportional to its error.

- The boosting method can correct the bias of the base classifier.

- AdaBoost implements the boosting algorithm for classification and regression.

- Gradient boosting is a more general framework that includes boosting algorithms.

# XGBoost: a scalable and flexible gradient boosting system

Markus Müller

- Extreme gradient boosting system

- Introduced by Chen & Guestrin in 2016

- Efficient implementation for tree and linear boosting for regression, classification, and ranking.

- Support for parallelization, cluster computing, and GPU.

- Very flexible and configurable API.
  - Callbacks, checkpoints, early stopping, dropout (dart), incremental learning, configurable sampling, regularization, many losses and metrics, choice of algorithm, and pruning, …

- XGBoost provides robust and complementary results when compared to neural network models on tabular data

# XGBoost

- Output after $t$ iterations: $\hat{y}_i^t = \sum_{k=1}^t f_k(\boldsymbol{x}_i)$;
- $f_k(\boldsymbol{x}_i) = $ base classifier at iteration $k$ (here we assume it's a tree)

- Boosting algo: fit $f_t(\boldsymbol{x}_i)$ to error $y_i - \hat{y}_i^{t-1}$ of previous iteration:

$$\mathcal{L}^t\left(\hat{y}_i^{t-1}, f_t\right) = \sum_{i=1}^N l\left(y_i, \hat{y}_i^{t-1} + f_t(\boldsymbol{x}_i)\right) + \Omega(f_t)$$

- loss function $l(\hat{y}_i, y_i)$
- regularization term $\Omega(f)$. For trees $\Omega(f) = \gamma T + \frac{1}{2}\lambda \sum_{i=1}^T v_i^2$
  - $T$: number of leaves; $\boldsymbol{v}$: values at leaves

# XGBoost

- Initialize $\hat{y}_i^0$

- for $t = 1..K$

  - Sample columns of $X$ and/or rows of $X, y$

  - $f_t(\boldsymbol{x}) = \underset{\boldsymbol{v}, Tree}{\operatorname{argmin}} \tilde{\mathcal{L}}^t\left(\hat{y}_i^{t-1}, f_t\right)$

    $\hat{y}_i^t = \hat{y}_i^{t-1} + \eta \, f_t(\boldsymbol{x}_i)$, shrinkage parameter eta: $0 < \eta \leq 1$

- Output $\hat{y}_i^K$

# XGBoost

- Gradient boosting:

  - $\mathcal{L}^t\left(\hat{y}_i^{t-1}, f_t\right) \approx \sum_{i=1}^{N}\left\{l\left(y_i, \hat{y}_i^{t-1}\right) + g_i f_t(\boldsymbol{x}_i) + \frac{1}{2} h_i f_t^2(\boldsymbol{x}_i)\right\} + \Omega(f_t)$

$$= \sum_{i=1}^{N} l\left(y_i, \hat{y}_i^{t-1}\right) + \underbrace{\sum_{i=1}^{N}\left\{g_i f_t(\boldsymbol{x}_i) + \frac{1}{2} h_i f_t^2(\boldsymbol{x}_i)\right\} + \Omega(f_t)}_{\tilde{\mathcal{L}}^t\left(\hat{y}_i^{t-1}, f_t\right)}$$

$$g_i = \frac{\partial}{\partial \hat{y}_i^{t-1}} l\left(y_i, \hat{y}_i^{t-1}\right); \ h_i = \frac{\partial^2}{\partial^2 \hat{y}_i^{t-1}} l\left(y_i, \hat{y}_i^{t-1}\right)$$
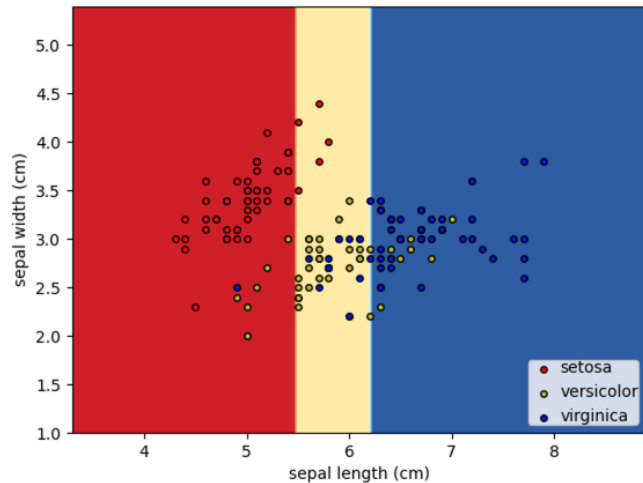
gradients                    hessians

# XGBoost - square loss

- $l(y, \hat{y}) = \frac{1}{2}(y - \hat{y})^2$

- Gradient boosting: $g_i = (\hat{y}_i^{t-1} - y_i); h_i = 1$

  - $\tilde{\mathcal{L}}^t(\hat{y}_i^{t-1}, f_t) = \sum_{i=1}^{N} \left\{ (\hat{y}_i^{t-1} - y_i) f_t(\boldsymbol{x}_i) + \frac{1}{2} f_t^2(\boldsymbol{x}_i) \right\} + \Omega(f_t)$
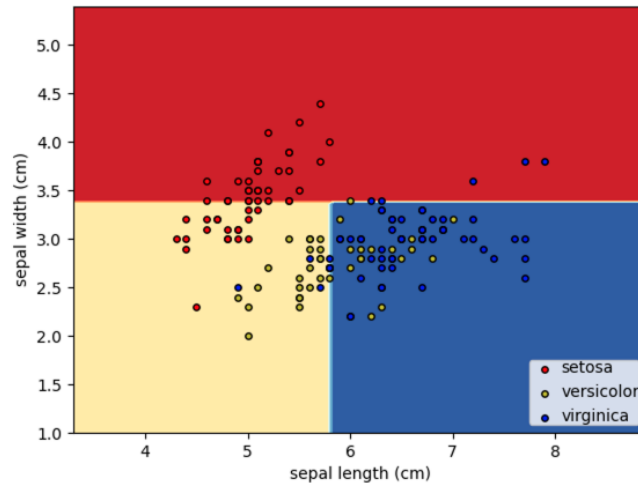
$(\hat{y}_i^{t-1} - y_i) > 0 : f_t(\boldsymbol{x}_i) < 0$

$(\hat{y}_i^{t-1} - y_i) \sim 0 : f_t(\boldsymbol{x}_i) \sim 0 \longleftarrow$

Regularization term

$(\hat{y}_i^{t-1} - y_i) < 0 : f_t(\boldsymbol{x}_i) > 0$

# XGBoost



n_estimators=100, gamma=1, max_depth=2, learning_rate=1, objective='multi:softmax', num_class=3): Iteration 1

n_estimators=100, gamma=1, max_depth=2, learning_rate=1, objective='multi:softmax', num_class=3): Iteration 2

n_estimators=100, gamma=1, max_depth=2, learning_rate=1, objective='multi:softmax', num_class=3): Iteration 3

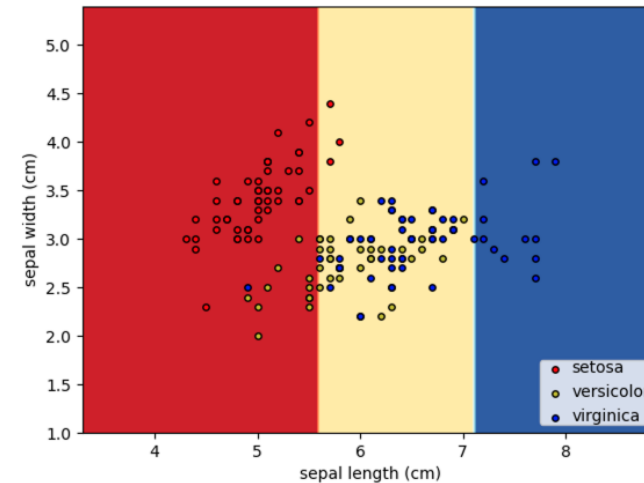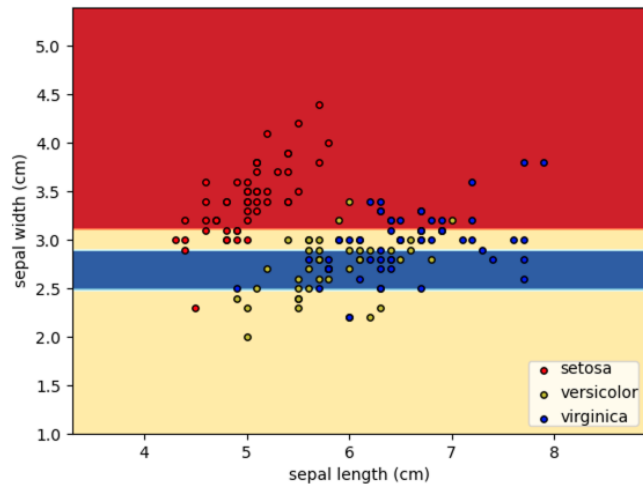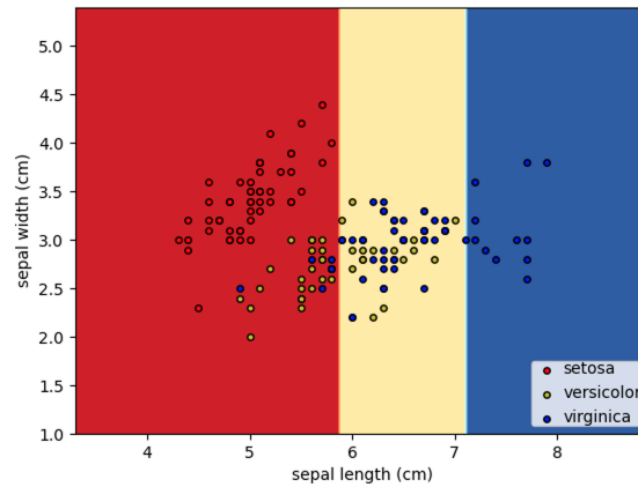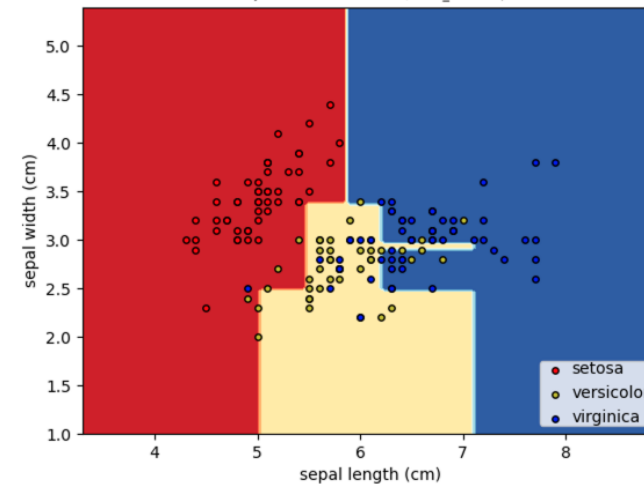n_estimators=100, gamma=1, max_depth=2, learning_rate=1, objective='multi:softmax', num_class=3): Iteration 4

n_estimators=100, gamma=1, max_depth=2, learning_rate=1, objective='multi:softmax', num_class=3): Iteration 5

n_estimators=10, gamma=1, max_depth=2, learning_rate=1, objective='multi:softmax', num_class=3)

# XGBoost

- How does XGBoost find the optimal tree?
- Let's assume the tree $f_t$ is known appart from the leaf values $\boldsymbol{v}$
- Define $I_j = \{i | leaf(\boldsymbol{x}_i) = j\}$, i.e. all points $\boldsymbol{x}_i$ from the training set that belong to leaf $j$
- Then the optimal leaf values are (<u>for square loss</u>):

  - $$v_j = -\frac{\sum_{i \in I_j}(\hat{y}_i^{t-1} - y_i)}{|I_j| + \lambda}$$

  - Then the <u>loss reduction</u> introduced by a split can then be calculated exactly:

$$\mathcal{L}_{split} = \frac{1}{2}\left[\frac{\left(\sum_{i \in I_L}(\hat{y}_i^{t-1} - y_i)\right)^2}{|I_L| + \lambda} + \frac{\left(\sum_{i \in I_R}(\hat{y}_i^{t-1} - y_i)\right)^2}{|I_R| + \lambda} - \frac{\left(\sum_{i \in I}(\hat{y}_i^{t-1} - y_i)\right)^2}{|I| + \lambda}\right] - \gamma$$

- The greedy *gbtree* algorithm finds the feature and the split $I = I_L \cup I_R$ with the largest $\mathcal{L}_{split} > 0$. If always $\mathcal{L}_{split} \leq 0$ then there is no further split.

# XGBoost gbtree algo



$L_{split} = 9.410$

$x < 8.5$
$I = \{1,2,3,4,5\}$

$\nu = -0.70$
$I = \{1,2,3,4\}$

$\nu = 4.57$
$I = \{5\}$

$\hat{y}^t = \hat{y}^{t-1} + \eta \times f_t$

$\hat{y}^{t-1}$

$\eta = 0.8;\ \lambda = 1; \gamma = 10$

# XGBoost – Approximate algo and missing values

- Instead of testing each value of a numeric feature as a potential threshold, XGBoost has an approximate *histogram*-based algorithm, which only allows splits at certain percentiles.

- In case of missing values, XGBoost calculates a default direction for each node. The default direction gives the largest gain $\mathcal{L}_{split}$ if all missing values are assigned to this direction.

# XGBoost Hyperparameters

- *booster*: base model (gbtree, dart or gblinear)
- *device*: GPU or CPU
- *eta*: step size shrinkage ($\eta$)
- *gamma*: minimum loss reduction required to make a further partition on a leaf node ($\gamma$)
- *max_depth*: maximum depth of a tree
- *min_child_weight*: minimum sum of instance hessians needed in a child
- *subsample*: subsample ratio of the training instances
- *colsample_bytree, colsample_bylevel, colsample_bynode*: parameters for subsampling of columns
- *lambda*: L2 regularization term on weights
- *alpha*: L1 regularization term on weights

- *tree_method*: The tree construction algorithm (exact, approx, and hist)
- *scale_pos_weight*: useful for unbalanced classes
- *process_type*: used to update existing model (transfer learning)
- *max_leaves*: maximum number of leaves
- *monotone_constraints*: constraint of variable monotonicity
- *interaction_constraints*: constraints for interaction representing permitted interactions
- *objective*: loss function
- *eval_metric*: evaluation metrics for validation data
- *max_cat_to_onehot*: how to deal with categorical features

# XGBoost +/-

- When to use XGBoost
  - Classification and regression problems
  - Categorical variables
  - Missing values
  - Tabular data with heterogeneous features
  - Very flexible and comprehensive API
  - Applicable to large data
  - ...

- When not to use XGBoost
  - $n > N$, i.e. more features than observations
  - Regression tasks, where output needs to be continuous
  - Extrapolation tasks
  - For image and language processing NN are preferred

- Other excellent gradient boosting tools are LightGBM, CatBoost, ..

# Hyperparameter (HP) optimization

- ML models have internal parameters, which are learned by the training algorithm, and hyperparameters, which have to be set before training.

- Both internal - and hyperparameters have a strong impact on the performance of the ML model.

- There is a tradeoff between CPU/GPU cycles spent on HP optimization and training. Allocating more cycles for HP optimization can be more rewarding.

- HP optimization is implemented as an outer loop over the CV loop

# Hyperparameter optimization and CV

Hyperparameter space $\Phi$

Dataset $\boldsymbol{X}$, Labels $\boldsymbol{y}$

Number of iterations $M$

Number of CV folds $F$

$S_{max} = -\infty, \boldsymbol{\phi}_{max} = None$

for $i = 1, .., M$:

$\quad \boldsymbol{\phi} = choose(\Phi)$

$\quad$ for $j = 1, .., F$:

$\quad\quad \boldsymbol{X}_{train}, \boldsymbol{X}_{test}, \boldsymbol{y}_{train}, \boldsymbol{y}_{test} = split(\boldsymbol{X}, j)$

$\quad\quad C_{\boldsymbol{\phi},\boldsymbol{\theta}} = \mathcal{C}_{\boldsymbol{\phi}}.\text{train}(\boldsymbol{X}_{train})$

$\quad\quad S_j = S\big(C_{\boldsymbol{\phi},\boldsymbol{\theta}}.\textbf{predict}(\boldsymbol{X}_{test}), \boldsymbol{y}_{test}\big)$

$\quad \bar{S} = \frac{1}{F}\sum_{j=1}^{F} S_j$

$\quad$ if $\bar{S} > S_{max}: S_{max} = \bar{S}, \boldsymbol{\phi}_{max} = \boldsymbol{\phi}$

Hyperparameter search loop

CV loop

Markus Müller

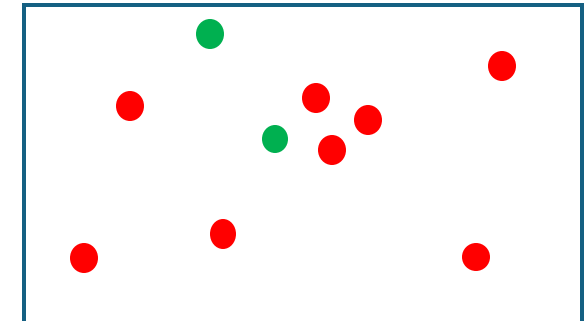SIB
Swiss Institute of
Bioinformatics

# HP searches

Grid search:
Parse all HP grid cells.
Slow if many features
Resolution of grid
limited by CPU time

Random search:
Random sampling based on prior
knowledge.
Difficult for many features
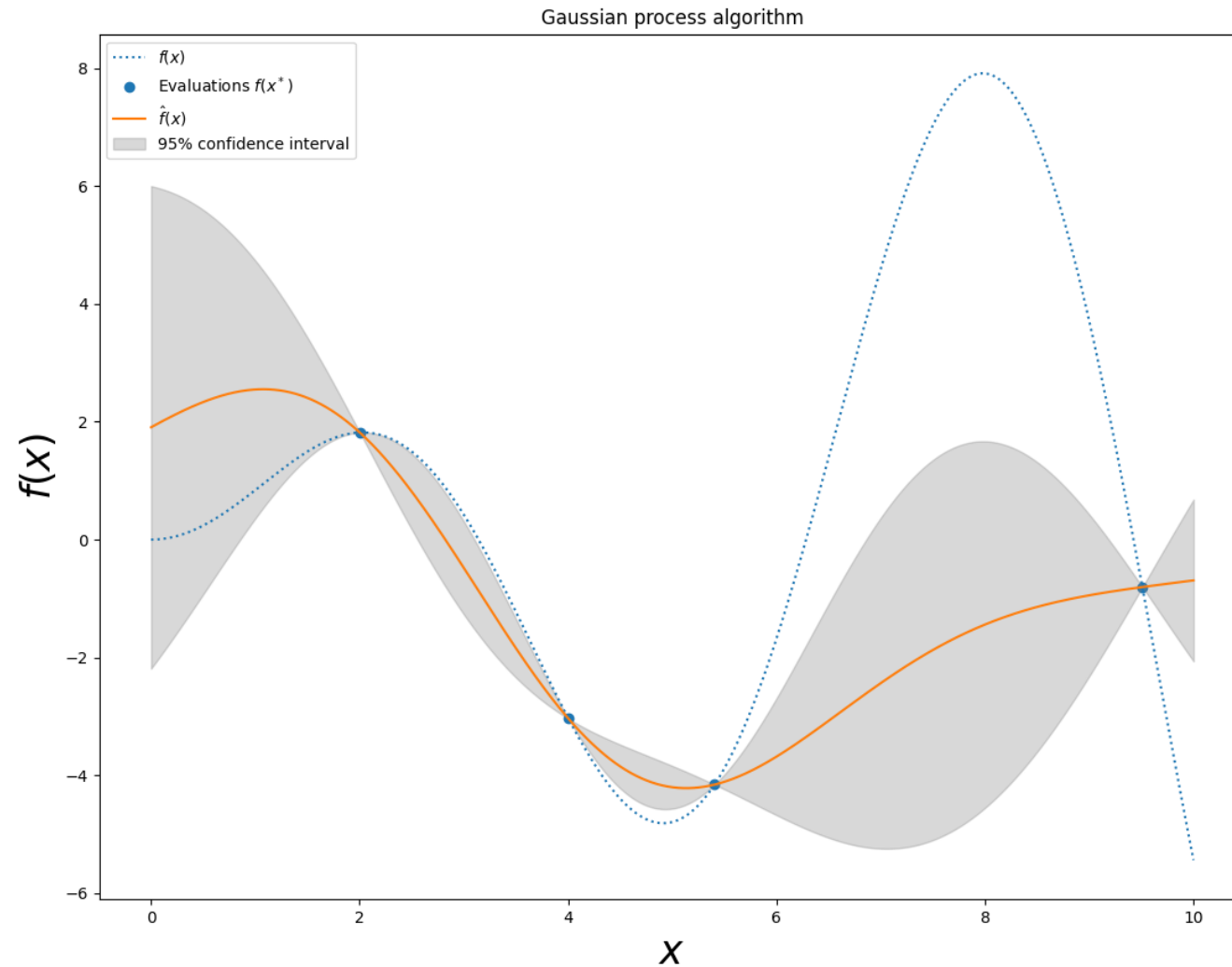Resolution limited by
CPU time

Intelligent searches:
Define new HPs based on
prior and acquired knowledge.
Often more efficient than
grid or random search for
many features

# Hyperopt

- Random algorithms find optimal solutions faster and more accurately than grid searches
- Hyperopt also offers two other model-based search algorithms for python:
  - Random search
  - Tree-structured Parzen window approach (TPE)
- Hyperopt supports nested parameter spaces
  - Hyperopt samples HPs from defined probability distributions or from predefined values
  - The HP space can have a hierarchical structure reflecting dependenies between HPs
- Hyperopt supports warm start using previous results
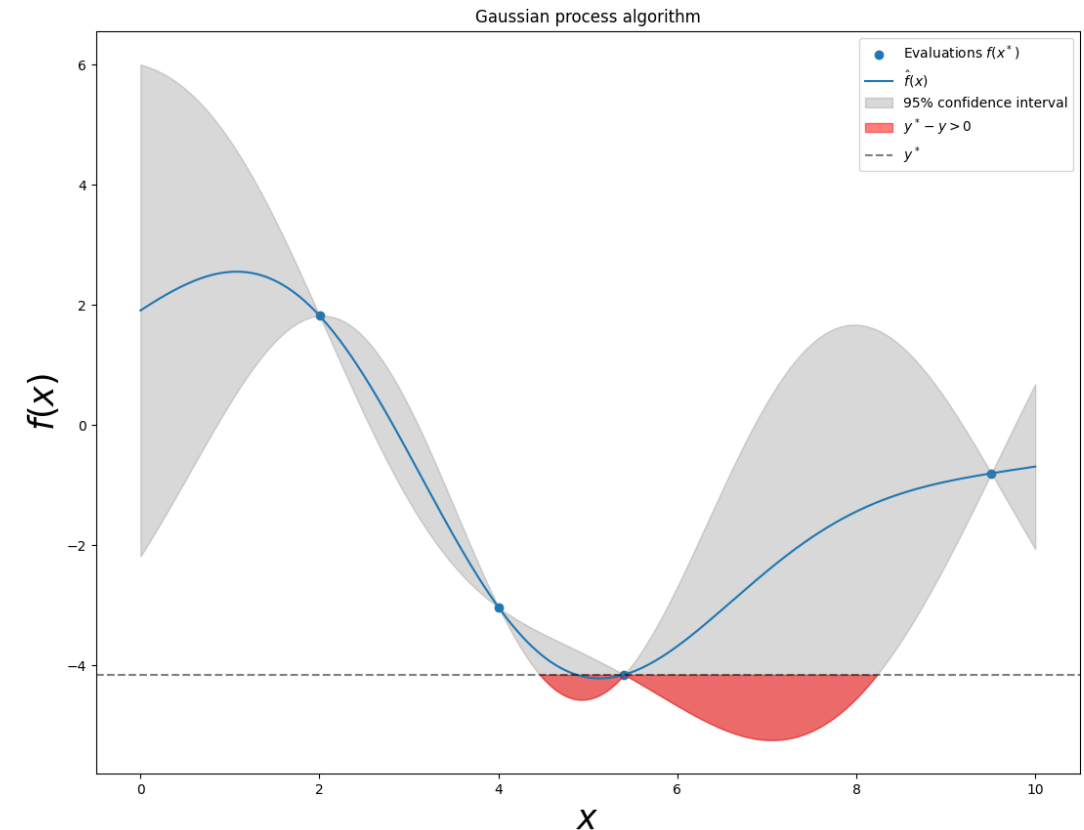
# Hyperopt's bayesian HP exploration

# Hyperopt's optimization score: expected improvement

$EI$: Expected improvement for hyperparameter values $x$

$y = f(\boldsymbol{x})$: objective function to minimize (e.g. loss function $L$)
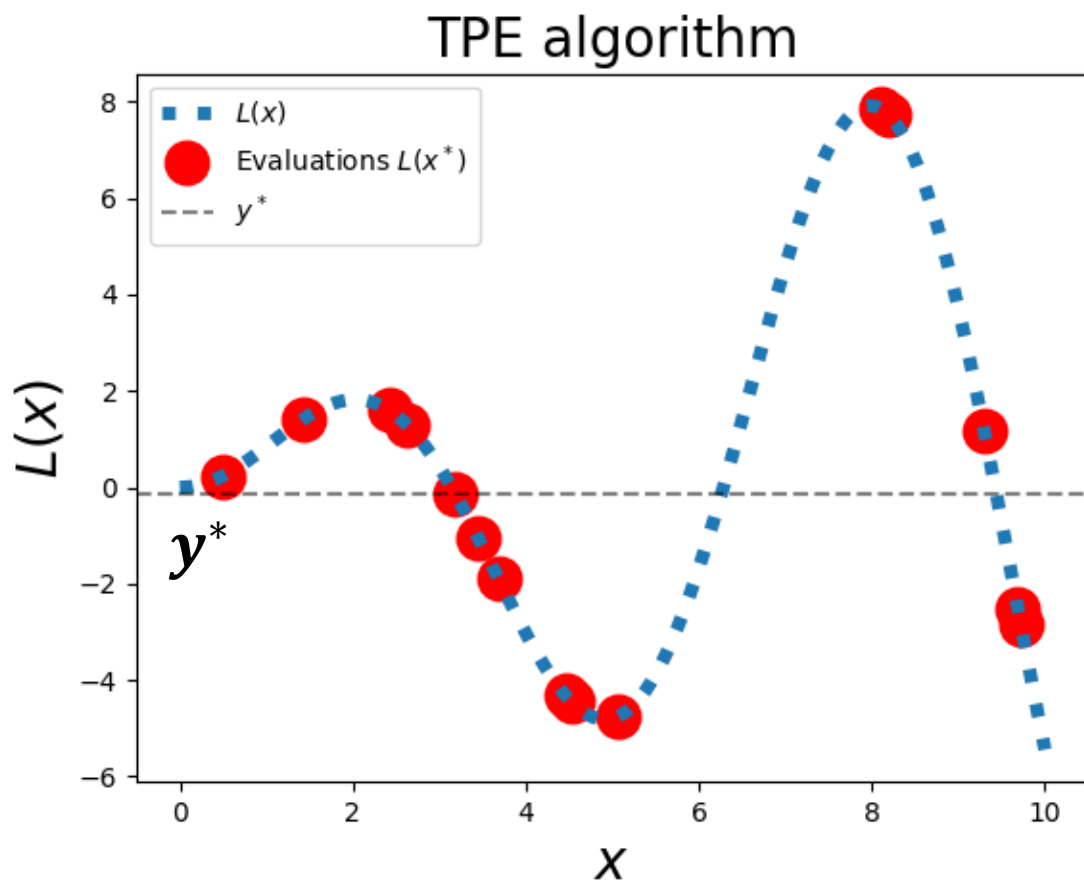
$y^*$: lower percentile of $y$

$$EI_{y^*}(x) = \int_{-\infty}^{y^*} (y^* - y)\, p_M(y|x)\, dy$$

$EI$ is high where $\hat{f}(x) = \int_{-\infty}^{\infty} y\, p_M(y|x)\, dy$ is low and/or

where the variance $\sigma^2(x) = \int_{-\infty}^{\infty} \left(y - \hat{f}(x)\right)^2 p_M(y|x)\, dy$

is high.
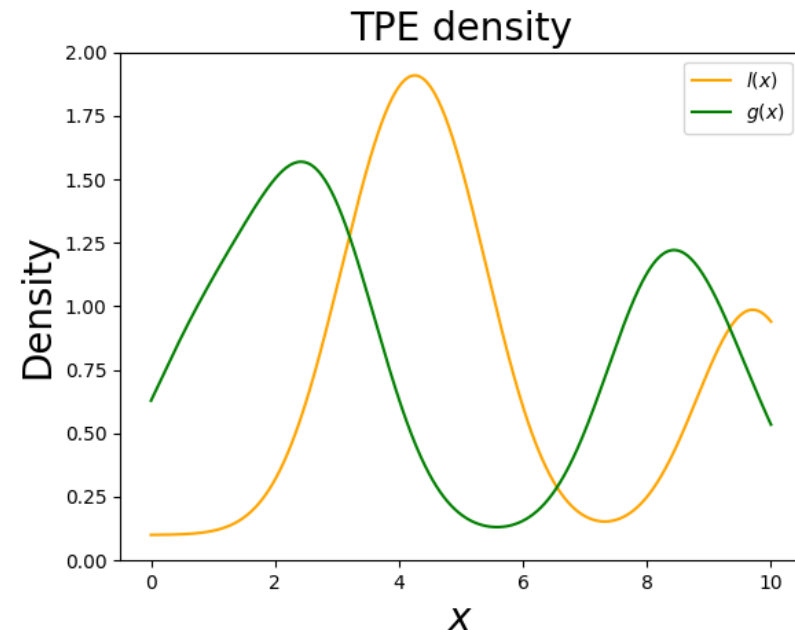
=> We need to estimate $p_M(y|x)$



Gaussian process algorithm

# Hyperopt's TPE algorithm

TPE: Tree structured Parzen estimator

$$y^* = \gamma - \text{quantile}: p(y < y^*) = \gamma$$



$$p_M(x|y) = \begin{cases} l(x) \text{ if } y \leq y^* \\ \text{g}(x) \text{ if } y > y^* \end{cases}$$

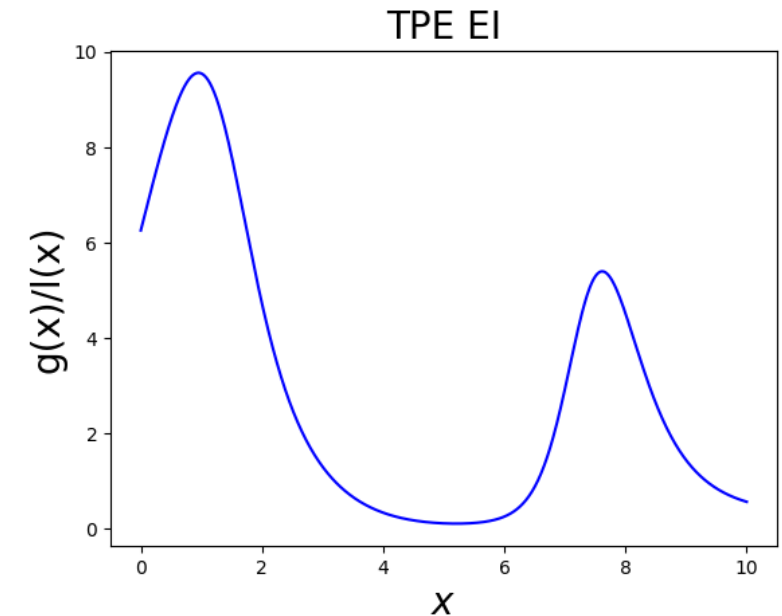$$p_M(x) = \gamma l(x) + (1 - \gamma)\text{g}(x)$$

$$p_M(y|x) = \frac{p_M(x|y)p(y)}{p(x)}$$

# Hyperopt's TPE algorithm

- $EI_{y^*}(x) = \int_{-\infty}^{y^*} (y^* - y) \frac{p_M(x|y)p(y)}{p(x)} dy$

$$\propto 1 / \left( \gamma + \frac{g(x)}{l(x)} (1 - \gamma) \right)$$



TPE EI

- The TPE algorithm will sample hyperscores from $l(x)$ and return the $x^*$ with the smallest $g(x)/l(x)$ value.

- $g(x)$ and $l(x)$ are hierarchical as defined in the HP space variable.

- $g(x)$ and $l(x)$ are initialized by the uniform or gaussian prior distributions specified in HP space variable. After each iteration, these distributions are updated with the newly sampled $x^*$ values.

# Hyperopt's model-based search

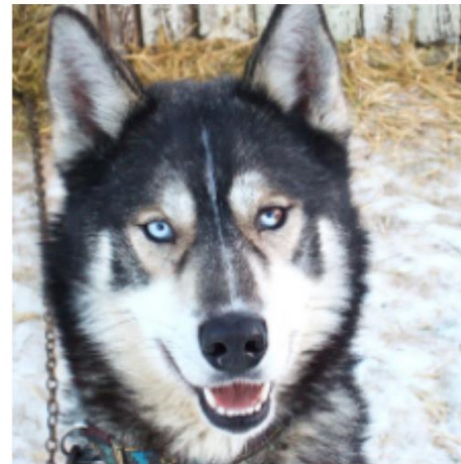$$M_t(x) \Rightarrow p_M(y|x)$$

$$S(x, M) = -EI_{y^*}(x)$$

$\text{SMBO}(f, M_0, T, S)$

| | |
|---|---|
| 1 | $\mathcal{H} \leftarrow \emptyset,$ |
| 2 | For $t \leftarrow 1$ **to** $T,$ |
| 3 | $x^* \leftarrow \text{argmin}_x \ S(x, M_{t-1}),$ |
| 4 | Evaluate $f(x^*),$ $\quad \triangleright$ *Expensive step* |
| 5 | $\mathcal{H} \leftarrow \mathcal{H} \cup (x^*, f(x^*)),$ |
| 6 | Fit a new model $M_t$ to $\mathcal{H}.$ |
| 7 | **return** $\mathcal{H}$ |

Figure 1: The pseudo-code of generic Sequential Model-Based Optimization.

# Interpretable Machine Learning (IML)

- ML tools are optimized for metrics such as accuracy, which can lead to accurate, but wrong models.

- Many models are black boxes, i.e. they do not reveal in a human readable form how they make their decision.

- In order to understand a model, we need to approximate the model's decision with a simple human readable model.

- If this simple approximation corresponds to human intuition and knowledge and is faithful (i.e. truly reflects the underlying ML model), we can better trust the ML model.

# Interpretable Machine Learning (IML)

Model → Data and Prediction → Explainer (LIME) → Explanation → Human makes decision

sneeze | Flu
weight
headache
no fatigue
age

sneeze
headache
no fatigue

Ribeiro et al., ACM SIGKDD, 2016

(a) Husky classified as wolf          (b) Explanation

# Interpretable Machine Learning (IML)



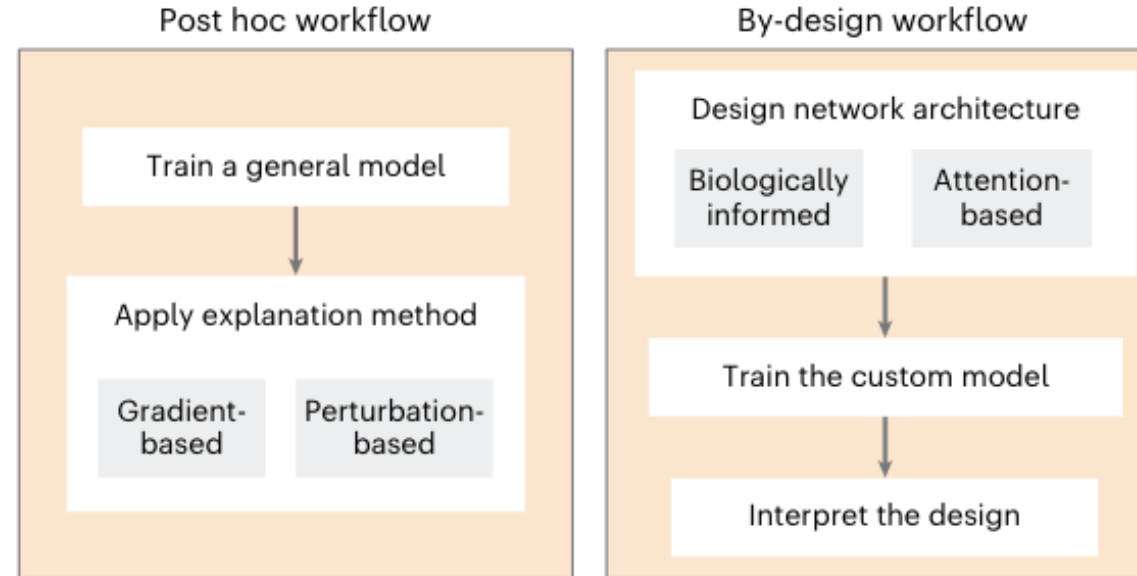**Fig. 1 | The two main IML approaches used to explain prediction models are post hoc explanations and by-design explanations.** Each approach has its canonical workflows and popular types of IML methods: post hoc explanations are model agnostic and are applied after a model is trained, while by-design explanations are typically built into or inherent to the model architecture.

Chen et al, Nature Methods, 2024

# IML Guidelines

- Use several IML methods and compare the results

- Use several datasets if possible, or randomly sample from your data, to show that your results are robust.

- Use different hyperparameter choices for the underlying model. Use more than one model.

- Evaluate the plausibility and faithfulness of the IML interpretation:
  - Do they agree with prior or expert knowledge?
  - Do they represent the model's reasoning process?

- Avoid cherry-picking nice examples, but always report full results.

# IML

- Local interpretation:
  - An interpretation is calculated for a single data vector $x$.
  - A local interpretation tells you, how a specific data vector $x$ was predicted by the ML model. It tells you which features are predominantly involved in the prediction.

- Global interpretation:
  - An interpretation is calculated on a whole data set $\{x_k\}, k = 1, .., N$
  - The global interpretation tells you, which features are predominantly involved in the prediction of a whole dataset.
  - A global interpretation can be obtained by *averaging* over local interpretations.

# LIME

- LIME: Local Interpretable Model-Agnostic Explanations
- Approximates any ML model by a simple local model
- Provides control over the complexity of the local model
- Provides global interpretation of the ML model
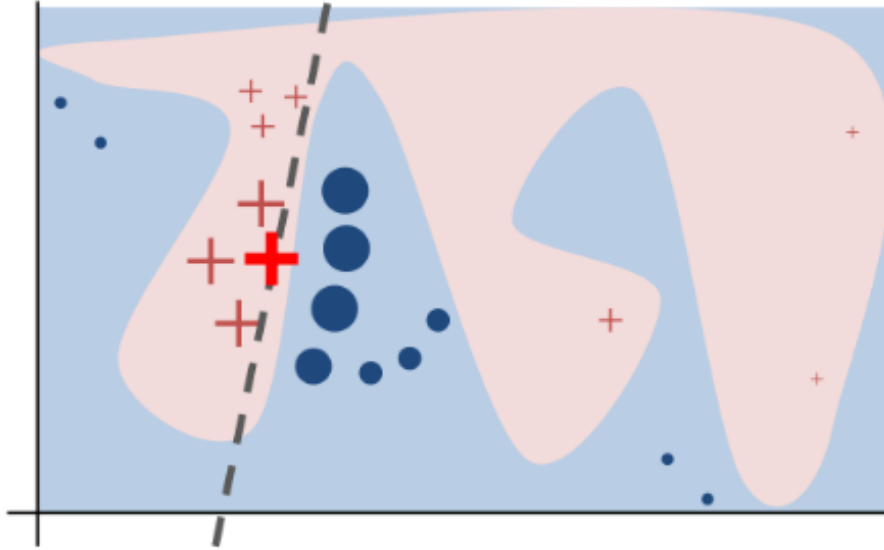
# LIME: local interpretation

Figure 3: Toy example to present intuition for LIME. The black-box model's complex decision function $f$ (unknown to LIME) is represented by the blue/pink background, which cannot be approximated well by a linear model. The bold red cross is the instance being explained. LIME samples instances, gets predictions using $f$, and weighs them by the proximity to the instance being explained (represented here by size). The dashed line is the learned explanation that is locally (but not globally) faithful.

$$\xi(x) = \operatorname*{argmin}_{g \in G} \mathcal{L}(f, g, \pi_x) + \Omega(g)$$

# Shapley Values

- Loyd Shapley (1951): How do you evaluate the importance of a single player in a team of $n$ players $\mathcal{N}$?

- Idea: the importance of a player is evaluated relative to all subsets of players $S$.

# Shapley Values

- The Shapley value $\varphi_i$ of player $i$ is defined by:

$$\varphi_i(f) = \frac{1}{n!} \sum_P \left( f\left(S_i^P \cup \{i\}\right) - f\left(S_i^P\right) \right)$$

Sum over all permutations $P$ of $\{1, \ldots, n\}$

Gain in $f$ by including player $i$ in team
$S_i^P = \{P_j \mid j < P(i)\}$

- $E_{\boldsymbol{x}}\left(f(\boldsymbol{x})\right) + \sum_{i=1}^{n} \varphi_i(\boldsymbol{x}) = f(\boldsymbol{x})$

# Shapley Values for ML

- Let's assume we have a <u>trained</u> ML model, which predicts the output $f(x_j)$ of a vector $x_j = (x_{j1}, \ldots, x_{jn}) \in \mathbb{R}^n$ from a dataset $\mathcal{D} = \{x_j\}, 1 \leq j \leq N$

- $S \subset \mathcal{N} = \{1, \ldots, n\}$. How do we define $f(x, S)$?

  - $f(x, S)$: keeps values $x_i$ for $i \in S$ and replaces values $x_i$ for $i \notin S$ by some backgound values
  - In the Shap package, you can define this background with the *masker* parameter.
  - You can also group the features before masking (e.g. if features are strongly dependent)

# IML with Shapley values

Markus Müller



Good prediction

Bad prediction