



Intro into Machine Learning

Decision Trees. Bagging. Ensembles. RandomForest. Stacked generalization

Third Machine Learning in High Energy Physics Summer School,
MLHEP 2017, July 17–23

Alexey Artemov^{1,2}

¹ Yandex LLC

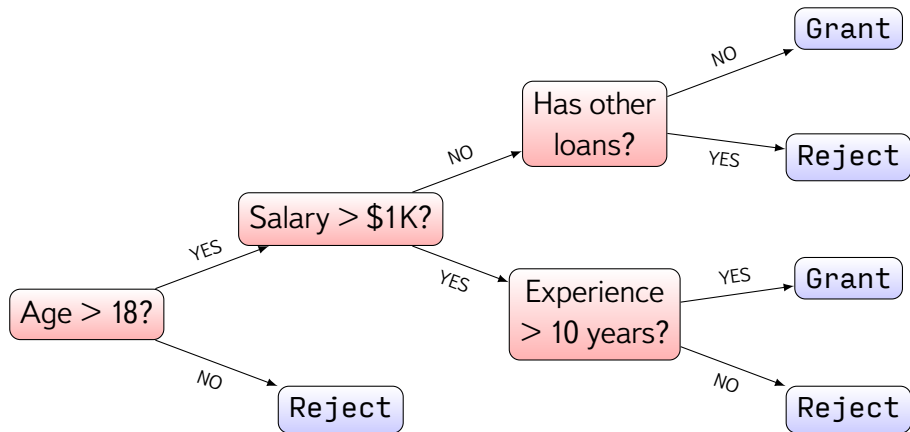
² National Research University Higher School of Economics

Lecture overview

- › Decision Trees
- › Bagging and Random Forests
- › Learning Theory continued
- › Stacked generalization

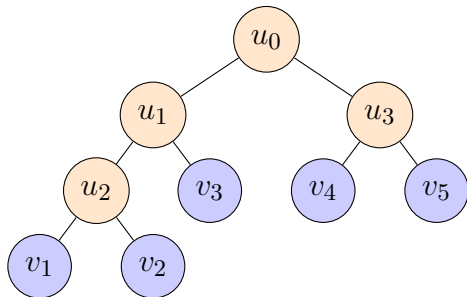
Decision trees

Decision making at a bank



Decision tree formalism

- › Decision tree is a binary tree V
- › Internal nodes $u \in V$: predicates
 $\beta_u : \mathbb{X} \rightarrow \{0, 1\}$
- › Leafs $v \in V$: predictions x
- › Algorithm $h(\mathbf{x})$ starts at $u = u_0$
 - › Compute $b = \beta_u(\mathbf{x})$
 - › If $b = 0$, $u \leftarrow \text{LeftChild}(u)$
 - › If $b = 1$, $u \leftarrow \text{RightChild}(u)$
 - › If u is a leaf, return b
- › In practice: $\beta_u(\mathbf{x}; j, t) = [\mathbf{x}_j < t]$



Greedy tree learning for binary classification

› Input: training set $X^\ell = \{(\mathbf{x}_i, y_i)\}_{i=1}^\ell$

1. Greedily split X^ℓ into R_1 and R_2 :

$$R_1(j, t) = \{\mathbf{x} \in X^\ell | \mathbf{x}_j < t\}, \quad R_2(j, t) = \{\mathbf{x} \in X^\ell | \mathbf{x}_j > t\}$$

optimizing a given loss: $Q(X^\ell, j, t) \rightarrow \min_{(j, t)}$

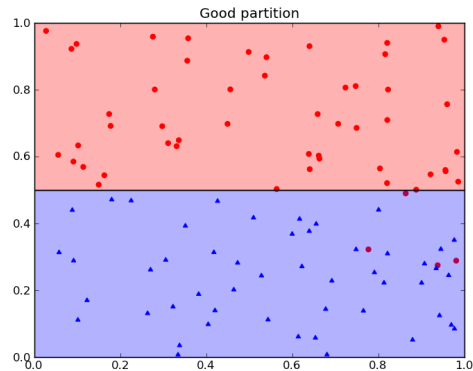
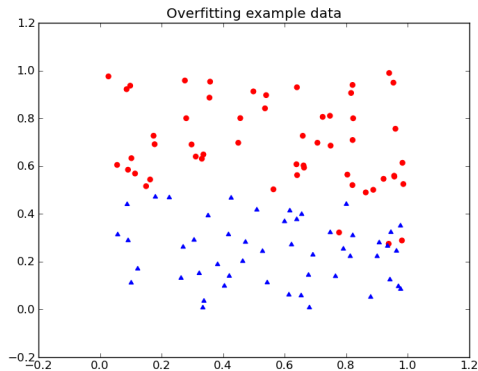
2. Create internal node u corresponding to the predicate $[\mathbf{x}_j < t]$

3. If a stopping criterion is satisfied for u ,
declare it a leaf, setting some $c_u \in \mathbb{Y}$ as leaf prediction

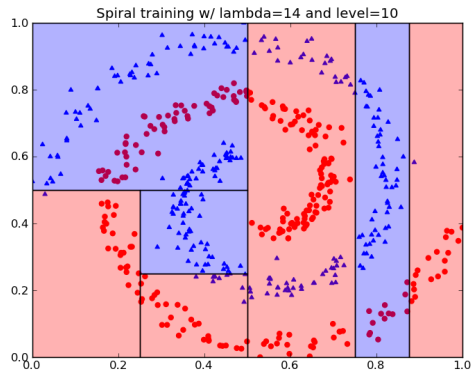
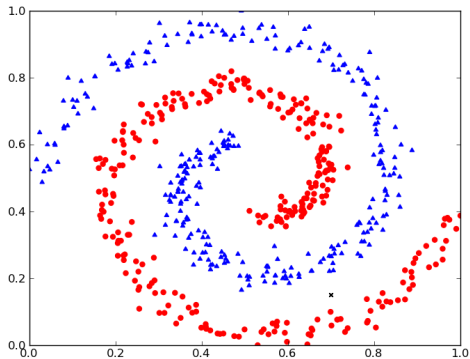
4. If not, repeat 1–2 for $R_1(j, t)$ and $R_2(j, t)$

› Output: a decision tree V

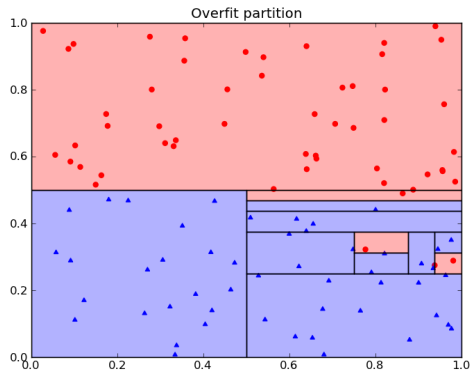
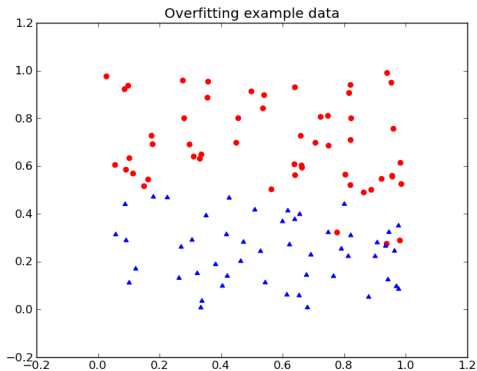
Greedy tree learning for binary classification



Greedy tree learning for binary classification



With decision trees, overfitting is extra-easy!



Design choices for learning a decision tree classifier

- › Type of predicate in internal nodes
 - › The loss function $Q(X^\ell, j, t)$
 - › The stopping criterion
 - › Hacks: missing values, pruning, etc.
-
- › CART, C4.5, ID3

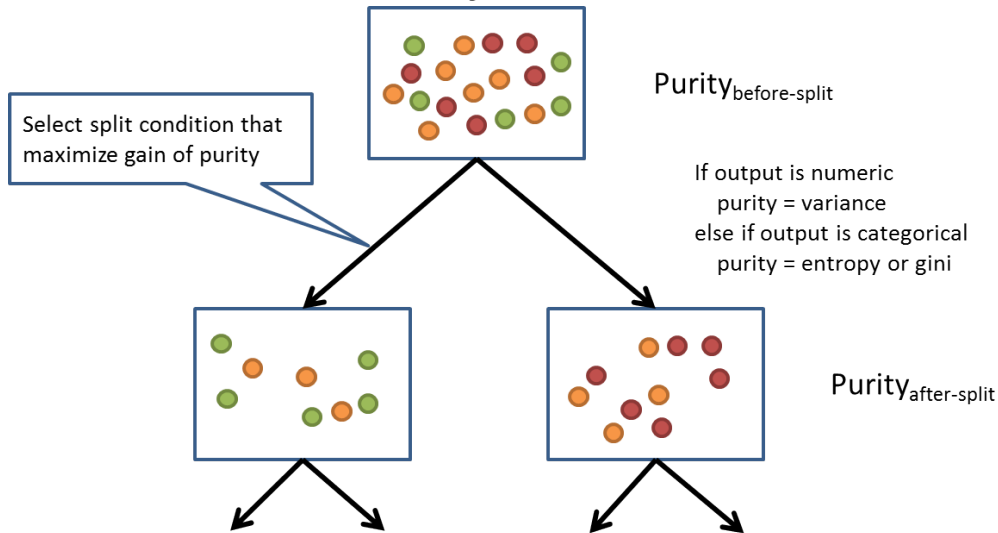
The loss function $Q(X^\ell, j, t)$

- › R_m : the subset of X^ℓ at step m
- › With the current split, let $R_l \subseteq R_m$ go left and $R_r \subseteq R_m$ go right
- › Choose predicate to optimize

$$Q(R_m, j, t) = H(R_m) - \frac{|R_l|}{|R_m|}H(R_l) - \frac{|R_r|}{|R_m|}H(R_r) \rightarrow \max$$

- › $H(R)$: impurity criterion
- › Generally $H(R) = \min_{c \in \mathbb{Y}} \frac{1}{|R|} \sum_{(\mathbf{x}_i, y_i) \in R} L(y_i, c)$

The idea: maximize purity



Examples of information criteria

› Regression:

- › $H(R) = \min_{c \in \mathbb{Y}} |R|^{-1} \sum_{(\mathbf{x}_i, y_i) \in R} (y_i - c)^2$
- › Sum of squared residuals minimized by $c = |R|^{-1} \sum_{(\mathbf{x}_j, y_j) \in R} y_j$
- › Impurity \equiv variance of the target

› Classification:

- › Let $p_k = |R|^{-1} \sum_{(\mathbf{x}_i, y_i) \in R} [y_i = k]$ (share of y_i 's equal to k)
- › Miss rate: $H(R) = \min_{c \in \mathbb{Y}} |R|^{-1} \sum_{(\mathbf{x}_i, y_i) \in R} [y_i \neq c]$

Minimizing miss rate $1 - p_{k_*}$,

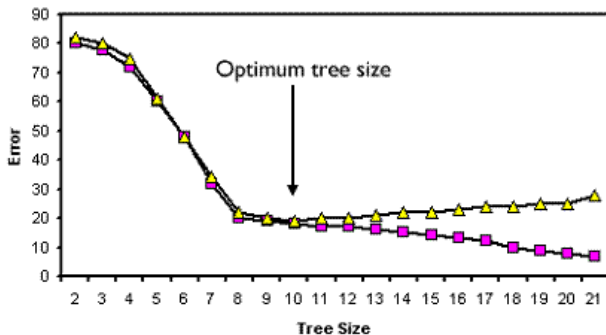
Gini index $\sum_{k=1}^K p_k (1 - p_k)$,

Cross-entropy $-\sum_{k=1}^K p_k \log p_k$

Stopping rules for decision tree learning

- › Significantly impacts learning performance
- › Multiple choices available:
 - › Maximum tree depth
 - › Minimum number of objects in leaf
 - › Maximum number of leafs in tree
 - › Stop if all objects fall into same leaf
 - › Constrain quality improvement
(stop when improvement gains drop below $s\%$)
- › Typically selected via exhaustive search and cross-validation

Decision tree pruning

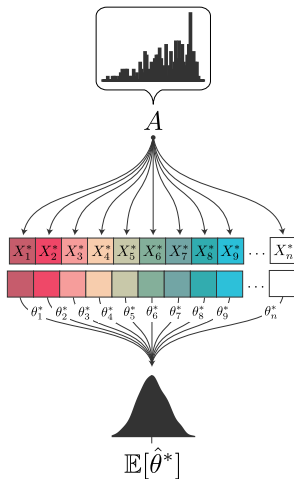


- › Learn a large tree (effectively overfit the training set)
- › Detect overfitting via K -fold cross-validation
- › Optimize structure by removing least important nodes

Bagging and Random Forests

The bootstrapping procedure

- › Input: a sample $X^\ell = \{(x_i, y_i)\}_{i=1}^\ell$
- › **Bootstrapping**: generate new samples X_1^m of (x_i, y_i) drawn from X^ℓ uniformly at random with replacement (replicated (x_i, y_i) possible!)
- › **Ensemble learning idea**:
 1. Generate N bootstrapped samples X_1^m, \dots, X_N^m
 2. Learn N hypotheses h_1, \dots, h_N
 3. Average predictions to obtain
$$h(x) = \frac{1}{N} \sum_{i=1}^N h_i(x)$$
 4. Profit!



Picture credit: <http://www.drbenzen.org/bootstrap-in-picture>

Bagging: bootstrap aggregation (+ demo)

- › Input: a sample

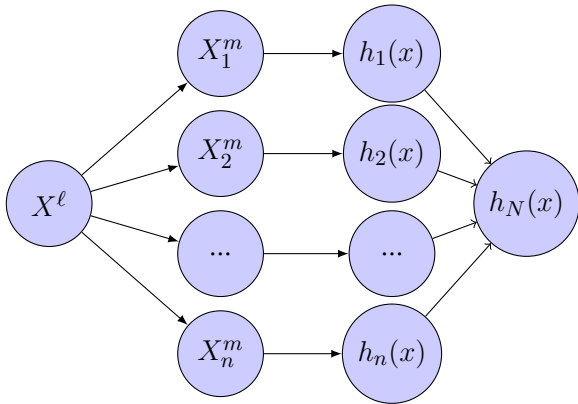
$$X^\ell = \{(x_i, y_i)\}_{i=1}^\ell$$

- › Weak learners via bootstrapping

$$\tilde{\mu}(X^\ell) = \mu(\tilde{X}^\ell)$$

- › Ensemble average

$$\begin{aligned} h_N(x) &= \frac{1}{N} \sum_{i=1}^N h_i(x) = \\ &= \frac{1}{N} \sum_{i=1}^N \tilde{\mu}(X^\ell)(x) \end{aligned}$$



The Random Forest algorithm

- › Bagging over (weak) decision trees
- › Reduce error via **averaging over instances and features**
- › Input: a sample $X^\ell = \{(\mathbf{x}_i, y_i)\}_{i=1}^\ell$, where $\mathbf{x}_i \in \mathbb{R}^d, y_i \in \mathbb{Y}$
- › The algorithm iterates for $i = 1, \dots, N$:
 1. Pick p random features out of d
 2. Bootstrap a sample $X_i^m = \{(\mathbf{x}_i, y_i)\}_{i=1}^\ell$ where $\mathbf{x}_i \in \mathbb{R}^p, y_i \in \mathbb{Y}$
 3. Learn a decision tree $h_i(\mathbf{x})$ using bootstrapped X_i^m
 4. Stop when leafs in h_i contain less than n_{\min} instances

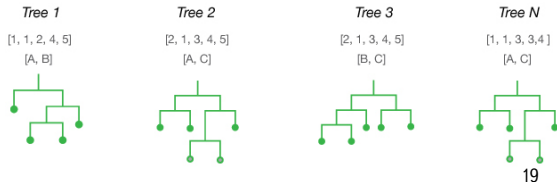
$$\mathbf{x}_i \in \{A, B, C\}$$

$$X^\ell = \{(\mathbf{x}_i, y_i)\}_{i=1}^5$$

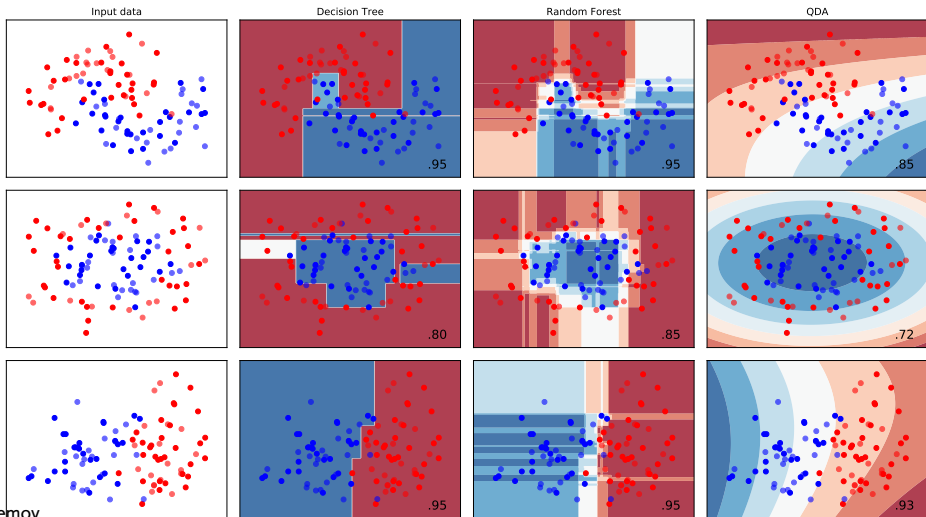
Bootstrap $X_i^m, i \in \{1, 2, 3, 4\}$

Learn $Tree_i(\mathbf{x})$ using X_i^m

Alexey Artemov



Random Forest: synthetic examples



Learning Theory (continued)

Expected risk formalism

- › Input: the training set $X^\ell = \{(x_i, y_i)\}_{i=1}^\ell$
- › Suppose $(x_i, y_i) \in \mathbb{X} \times \mathbb{Y}$ are generated from a distribution $p(x, y)$
- › Consider the MSE loss $Q(y, h) = (y - h(x))^2$
- › Expected risk: average (over $p(x, y)$) squared loss when using h

$$R(h) = \mathbb{E}_{x,y} \left[(y - h(x))^2 \right] = \int_{\mathbb{X}} \int_{\mathbb{Y}} p(x, y) (y - a(x))^2 dx dy.$$

- › **A statement:** the expected risk is minimized by

$$h_*(x) = \mathbb{E}[y | x] = \int_{\mathbb{Y}} yp(y | x)dy = \arg \min_h R(h)$$

Bias-variance decomposition

- › Input: the training set $X^\ell = \{(x_i, y_i)\}_{i=1}^\ell$
- › Learning method $\mu : (\mathbb{X} \times \mathbb{Y})^\ell \rightarrow \mathbb{H}$
- › Can evaluate quality using average (over possible samples X^ℓ)

$$\begin{aligned} Q(\mu) &= \mathbb{E}_{X^\ell} \left[\mathbb{E}_{x,y} \left[\left(y - \mu(X^\ell)(x) \right)^2 \right] \right] = \\ &= \int_{(\mathbb{X} \times \mathbb{Y})^\ell} \int_{\mathbb{X} \times \mathbb{Y}} \left(y - \mu(X^\ell)(x) \right)^2 p(x, y) \prod_{i=1}^{\ell} p(x_i, y_i) dx dy dx_i dy_i \end{aligned}$$

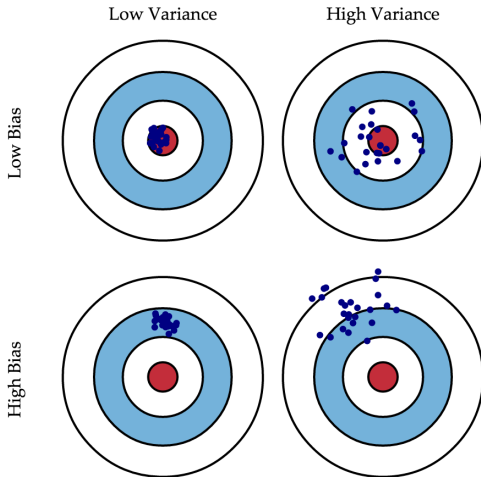
Bias-variance decomposition

- › We arrive at the famous **bias-variance(-noise) decomposition**

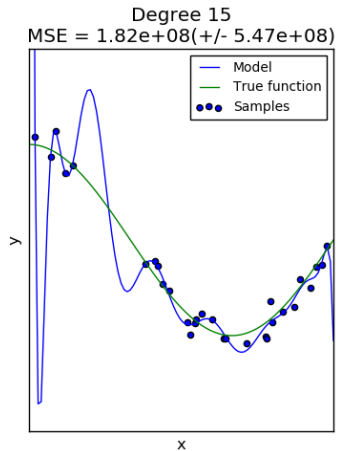
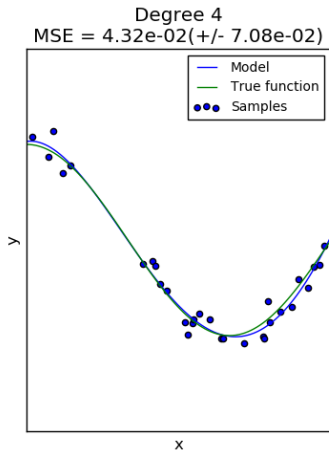
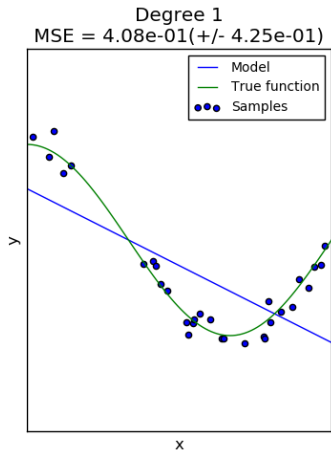
$$\begin{aligned} Q(\mu) = & \underbrace{\mathbb{E}_{x,y} \left[\left(y - \mathbb{E}[y | x] \right)^2 \right]}_{\text{noise}} + \\ & \underbrace{\mathbb{E}_x \left[\left(\mathbb{E}_{X^\ell} [\mu(X^\ell)] - \mathbb{E}[y | x] \right)^2 \right]}_{\text{bias}} + \underbrace{\mathbb{E}_x \left[\mathbb{E}_{X^\ell} \left[\left(\mu(X^\ell) - \mathbb{E}_{X^\ell} [\mu(X^\ell)] \right)^2 \right] \right]}_{\text{variance}} \end{aligned}$$

- › **Noise term**: an error of an ideal learner (nobody can do better!)
- › **Bias term**: learner's approximation of the ideal algorithm
 - › The more complex the learning algorithm, the lower the bias
- › **Variance term**: sensitivity to sample replacement
 - › Simple algorithms have lower variance

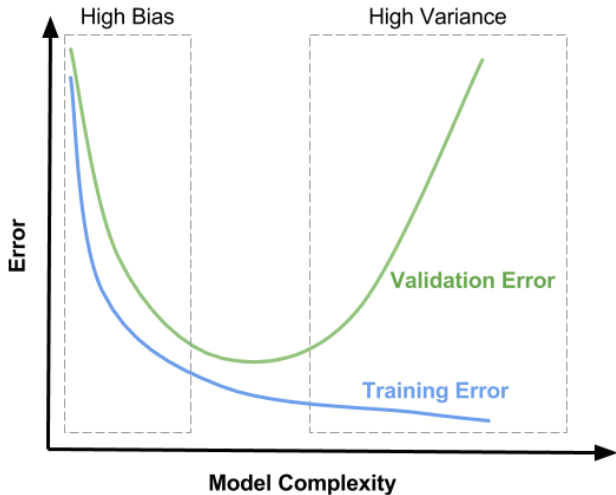
Bias-variance decomposition



Polynomial fits of different degrees



Bias-variance tradeoff



An application: statistical analysis of Bagging

- › **Bias:** not made any worse by bagging multiple hypotheses

$$\begin{aligned}\mathbb{E}_{x,y} \left[\left(\mathbb{E}_{X^\ell} \left[\frac{1}{N} \sum_{n=1}^N \tilde{\mu}(X^\ell)(x) \right] - \mathbb{E}[y | x] \right)^2 \right] &= \\ &= \mathbb{E}_{x,y} \left[\left(\frac{1}{N} \sum_{n=1}^N \mathbb{E}_X^\ell [\tilde{\mu}(X^\ell)(x)] - \mathbb{E}[y | x] \right)^2 \right] = \\ &= \mathbb{E}_{x,y} \left[\left(\mathbb{E}_{X^\ell} [\tilde{\mu}(X^\ell)(x)] - \mathbb{E}[y | x] \right)^2 \right]\end{aligned}$$

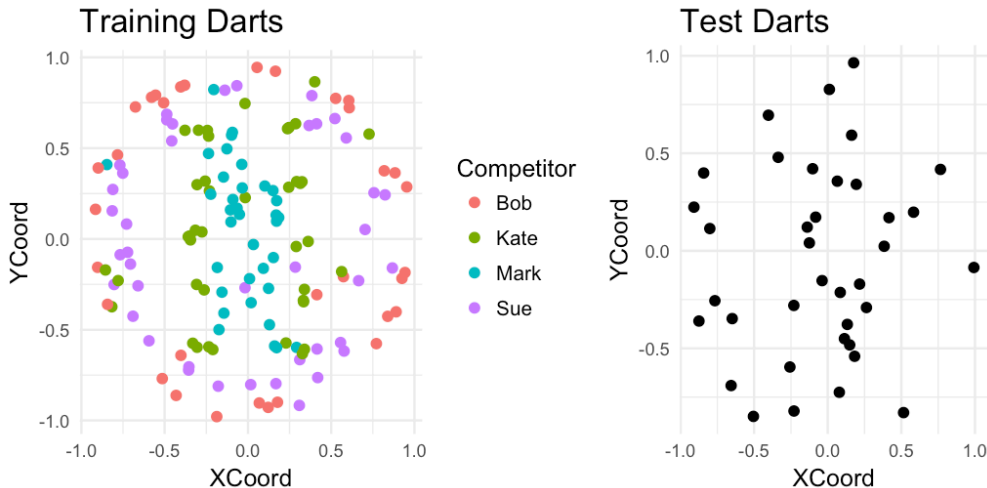
An application: statistical analysis of Bagging

- › **Variance:** N times lower for uncorrelated hypotheses, yet not as much an improvement for highly correlated

$$\begin{aligned} \mathbb{E}_{x,y} \left[\mathbb{E}_{X^\ell} \left[\left(\frac{1}{N} \sum_{n=1}^N \tilde{\mu}(X^\ell)(x) - \mathbb{E}_{X^\ell} \left[\frac{1}{N} \sum_{n=1}^N \tilde{\mu}(X^\ell)(x) \right] \right)^2 \right] \right] &= \\ &= \frac{1}{N} \mathbb{E}_{x,y} \left[\mathbb{E}_{X^\ell} \left[\left(\tilde{\mu}(X^\ell)(x) - \mathbb{E}_{X^\ell} [\tilde{\mu}(X)(x)] \right)^2 \right] \right] + \\ &+ \frac{N(N-1)}{N^2} \mathbb{E}_{x,y} \left[\mathbb{E}_{X^\ell} \left[\left(\tilde{\mu}(X^\ell)(x) - \mathbb{E}_{X^\ell} [\tilde{\mu}(X^\ell)(x)] \right) \times \right. \right. \\ &\quad \left. \left. \times \left(\tilde{\mu}(X^\ell)(x) - \mathbb{E}_{X^\ell} [\tilde{\mu}(X^\ell)(x)] \right) \right] \right] \end{aligned}$$

Stacked generalization

Stacking motivation: the game of Darts

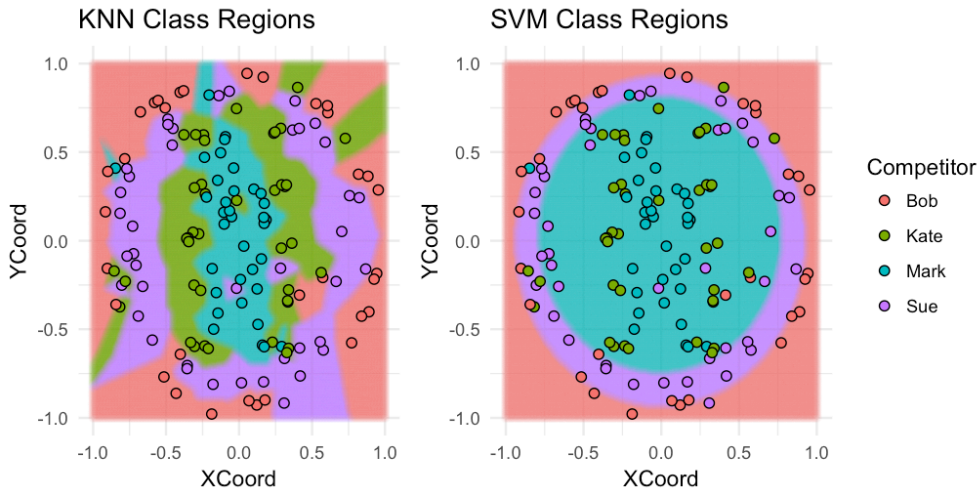


Base model training

- › Select k nearest neighbours as base model 1
- › Fit base model 1 in the most fancy way possible
(grid search for optimal k using K -fold cross-validation, etc.)
- › k -NN accuracy on Test Darts: 70%

- › Select Support Vector Machine as base model 2
- › Fit base model 2 in the most fancy way possible
(different penalizations, grid search for optimal kernel width using K -fold cross-validation, etc.)
- › SVM accuracy on Test Darts: 78%

Results for base models



Stacking base models

1. Partition `train` into 5 folds
2. Create `train_meta/test_meta`: same row/fold Ids as in `train/test`, empty `M1/M2`
3. For each $\text{Fold}_i \in \{\text{Fold}_1, \dots, \text{Fold}_5\}$
 - 3.1 Combine the other 4 folds for training $\rightarrow \text{Fold}_{-i}$
 - 3.2 Fit each base model to `Fold-i`, predict on `Foldi`, save predictions to `M1/M2` in `train_meta`
4. Fit each base model to `train`, predict on `test`, save predictions to `M1/M2` in `test_meta`
5. Fit stacking model `S` to `train_meta`, using `M1/M2` as features
6. Use the stacked model `S` to make final predictions on `test_meta`

Results for base models



Conclusion

- › Decision trees: intuitive and interpretable, yet prone to overfitting
- › Bootstrapping: a general statistical technique for computing sample functionals (and their variance)
- › Bagging: meta-learner over arbitrary weak algorithms via bootstrap aggregation
- › The Random Forest algorithm: Bagging over decision trees
- › Stacked generalization: blend output of weak learners (weak signals) with raw features

Bonus track

Minimum of the expected risk: the proof

- › Transform the loss

$$\begin{aligned}Q(y, h(x)) &= (y - h(x))^2 = (y - \mathbb{E}(y | x) + \mathbb{E}(y | x) - h(x))^2 = \\&= (y - \mathbb{E}(y | x))^2 + 2(y - \mathbb{E}(y | x))(\mathbb{E}(y | x) - h(x)) + \\&+ (\mathbb{E}(y | x) - h(x))^2.\end{aligned}$$

- › Write the expected risk

$$\begin{aligned}R(h) &= \mathbb{E}_{x,y} Q(y, h(x)) = \\&= \mathbb{E}_{x,y} (y - \mathbb{E}(y | x))^2 + \mathbb{E}_{x,y} (\mathbb{E}(y | x) - h(x))^2 + \\&+ 2\mathbb{E}_{x,y} (y - \mathbb{E}(y | x))(\mathbb{E}(y | x) - h(x)).\end{aligned}$$

Minimum of the expected risk: the proof

› Consider $\mathbb{E}_{x,y}(y - \mathbb{E}(y | x))(\mathbb{E}(y | x) - h(x))$ which is essentially some

$$\mathbb{E}_x \mathbb{E}_y [f(x, y) | x] = \int_{\mathbb{X}} \left(\int_{\mathbb{Y}} f(x, y) p(y | x) dy \right) p(x) dx$$

meaning that (as $(\mathbb{E}(y | x) - h(x))$ is independent of y):

$$\begin{aligned} \mathbb{E}_x \mathbb{E}_y \left[(y - \mathbb{E}(y | x)) (\mathbb{E}(y | x) - h(x)) | x \right] &= \\ &= \mathbb{E}_x \left((\mathbb{E}(y | x) - h(x)) \mathbb{E}_y \left[(y - \mathbb{E}(y | x)) | x \right] \right) = \\ &= \mathbb{E}_x \left((\mathbb{E}(y | x) - h(x)) (\mathbb{E}(y | x) - \mathbb{E}(y | x)) \right) = \\ &= 0 \end{aligned}$$

Minimum of the expected risk: the proof

- › We obtain that the expected risk has the form

$$R(h) = \mathbb{E}_{x,y}(y - \mathbb{E}(y | x))^2 + \mathbb{E}_{x,y}(\mathbb{E}(y | x) - h(x))^2.$$

- › Both summands are nonnegative meaning that the sum is minimized when

$$h_*(x) = \mathbb{E}(y | x) = \int_{\mathbb{Y}} yp(y | x)dy.$$

Bias-variance decomposition

- › Input: the training set $X^\ell = \{(x_i, y_i)\}_{i=1}^\ell$
- › Learning method $\mu : (\mathbb{X} \times \mathbb{Y})^\ell \rightarrow \mathbb{H}$
- › Can evaluate quality using average (over possible samples X^ℓ)

$$\begin{aligned} Q(\mu) &= \mathbb{E}_{X^\ell} \left[\mathbb{E}_{x,y} \left[\left(y - \mu(X^\ell)(x) \right)^2 \right] \right] = \\ &= \int_{(\mathbb{X} \times \mathbb{Y})^\ell} \int_{\mathbb{X} \times \mathbb{Y}} \left(y - \mu(X^\ell)(x) \right)^2 p(x, y) \prod_{i=1}^{\ell} p(x_i, y_i) dx dy dx_i dy_i \end{aligned}$$

- › For a fixed sample X^ℓ , we know the expected risk

$$\mathbb{E}_{x,y} \left[\left(y - \mu(X^\ell) \right)^2 \right] = \mathbb{E}_{x,y} \left[\left(y - \mathbb{E}[y | x] \right)^2 \right] + \mathbb{E}_{x,y} \left[\left(\mathbb{E}[y | x] - \mu(X^\ell) \right)^2 \right]$$

Bias-variance decomposition

- › Plugging the expression for fixed X^ℓ into $Q(\mu)$, we obtain

$$Q(\mu) = \mathbb{E}_{X^\ell} \left[\underbrace{\mathbb{E}_{x,y} \left[\left(y - \mathbb{E}[y | x] \right)^2 \right]}_{\text{independent of } X^\ell} + \mathbb{E}_{x,y} \left[\left(\mathbb{E}[y | x] - \mu(X^\ell) \right)^2 \right] \right]$$

- › Transforming the second summand

$$\begin{aligned} \mathbb{E}_{x,y} \left[\mathbb{E}_{X^\ell} \left[\left(\mathbb{E}[y | x] - \mu(X^\ell) \right)^2 \right] \right] &= \\ &= \mathbb{E}_{x,y} \left[\mathbb{E}_{X^\ell} \left[\left(\mathbb{E}[y | x] - \mathbb{E}_{X^\ell} [\mu(X^\ell)] + \mathbb{E}_{X^\ell} [\mu(X^\ell)] - \mu(X^\ell) \right)^2 \right] \right] = \\ &= \mathbb{E}_{x,y} \left[\underbrace{\mathbb{E}_{X^\ell} \left[\left(\mathbb{E}[y | x] - \mathbb{E}_{X^\ell} [\mu(X^\ell)] \right)^2 \right]}_{\text{independent of } X^\ell} + \mathbb{E}_{x,y} \left[\mathbb{E}_{X^\ell} \left[\left(\mathbb{E}_{X^\ell} [\mu(X^\ell)] - \mu(X^\ell) \right)^2 \right] \right] \right. \\ &\quad \left. + 2 \mathbb{E}_{x,y} \left[\mathbb{E}_{X^\ell} \left[\left(\mathbb{E}[y | x] - \mathbb{E}_{X^\ell} [\mu(X^\ell)] \right) \left(\mathbb{E}_{X^\ell} [\mu(X^\ell)] - \mu(X^\ell) \right) \right] \right] \right] \end{aligned}$$

Bias-variance decomposition

› We prove that the last term is zero:

$$\begin{aligned}\mathbb{E}_{X^\ell} \left[\left(\mathbb{E}[y | x] - \mathbb{E}_{X^\ell} [\mu(X^\ell)] \right) \left(\mathbb{E}_{X^\ell} [\mu(X^\ell)] - \mu(X^\ell) \right) \right] &= \\ &= \left(\mathbb{E}[y | x] - \mathbb{E}_{X^\ell} [\mu(X^\ell)] \right) \mathbb{E}_X \left[\mathbb{E}_{X^\ell} [\mu(X^\ell)] - \mu(X^\ell) \right] = \\ &= \left(\mathbb{E}[y | x] - \mathbb{E}_{X^\ell} [\mu(X^\ell)] \right) \left[\mathbb{E}_{X^\ell} [\mu(X^\ell)] - \mathbb{E}_{X^\ell} [\mu(X^\ell)] \right] = \\ &= 0.\end{aligned}$$