

Bootstrap, Bagging, and Random Forests

Nishant Mehta

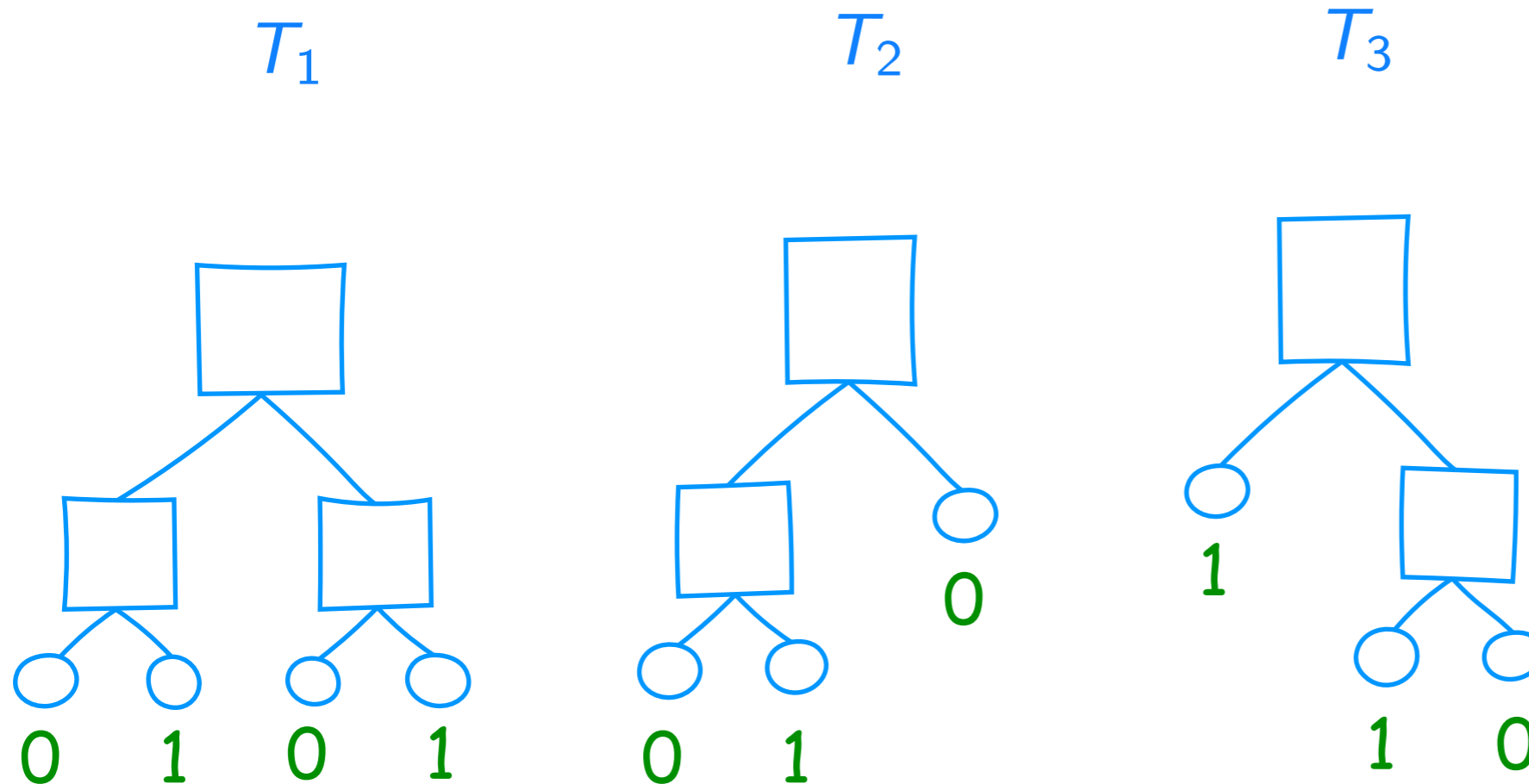
Lecture 3 - Part I

Ensemble methods

An *ensemble method*: a method that predicts by aggregating the predictions of many hypotheses

Bagged Trees and *Random Forests* are both ensemble methods

- Common idea: Take the majority vote (for classification) among a large, diverse set of classifiers, each of which is a decision tree (trained on a random *bootstrap* sample of the training set)

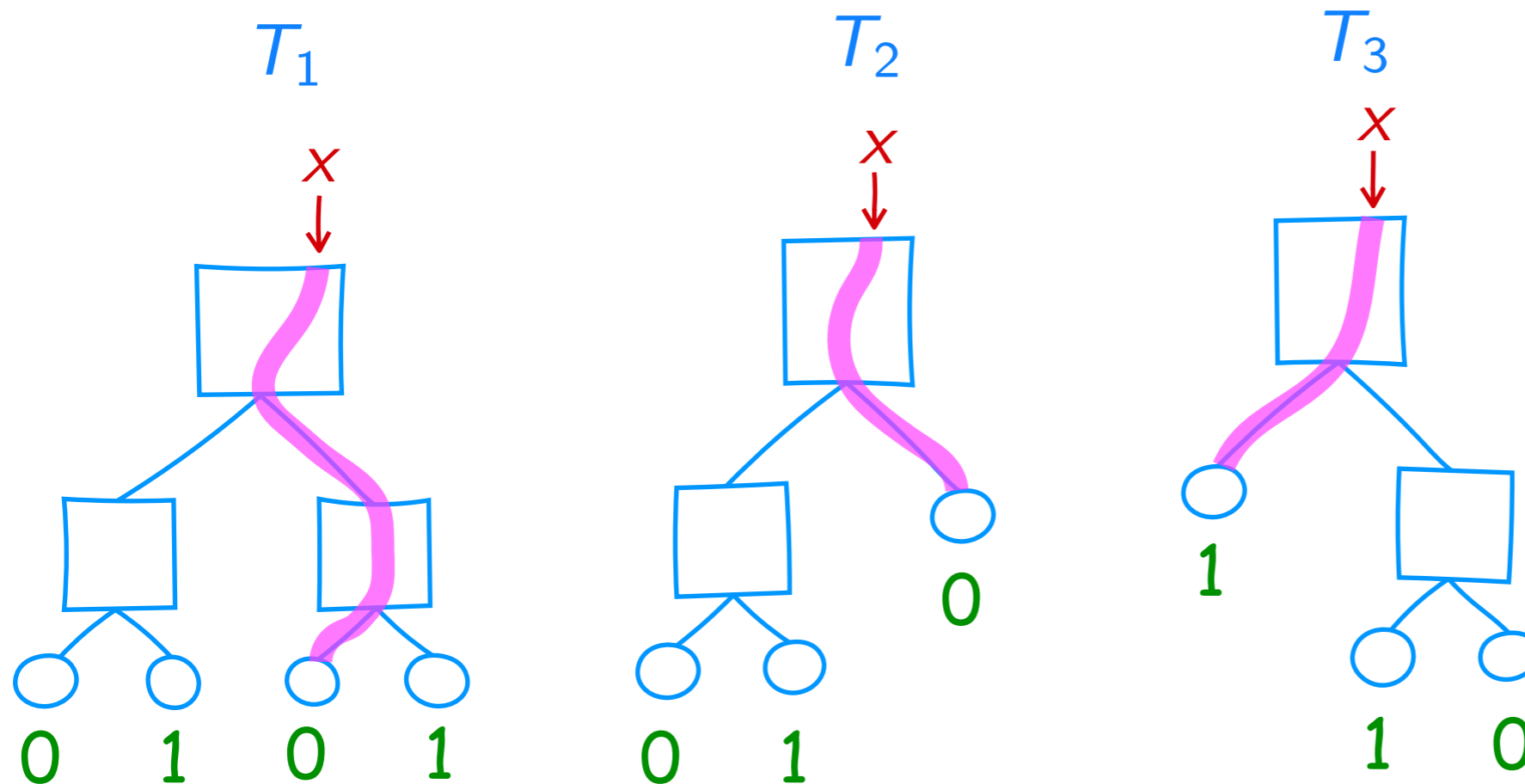


Ensemble methods

An *ensemble method*: a method that predicts by aggregating the predictions of many hypotheses

Bagged Trees and *Random Forests* are both ensemble methods

- Common idea: Take the majority vote (for classification) among a large, diverse set of classifiers, each of which is a decision tree (trained on a random *bootstrap* sample of the training set)



Bootstrap Aggregation (Bagging)

Given a training sample $D = ((X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n))$ of size n , a *bootstrap sample* is obtained by drawing n examples *with replacement* from D .

For a given learning algorithm, *bootstrap aggregation*, or *bagging*, trains the learning algorithm M times as follows:

For $j = 1, 2, \dots, M$

- Draw a new bootstrap sample \tilde{D}_j from D
- Train the algorithm on \tilde{D}_j , giving hypothesis \hat{h}_j

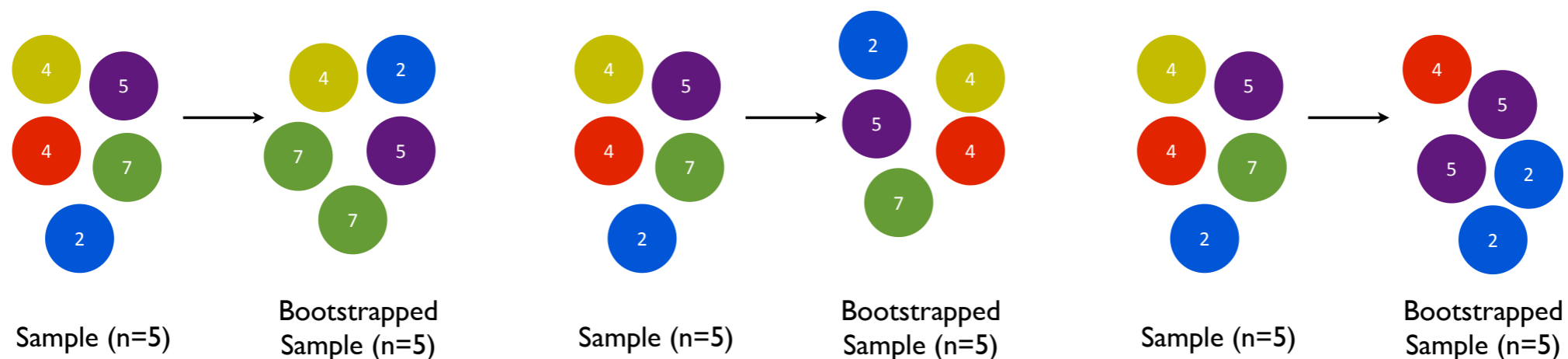
Bootstrap Aggregation (Bagging)

Given a training sample $D = ((X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n))$ of size n , a **bootstrap sample** is obtained by drawing n examples *with replacement* from D .

For a given learning algorithm, **bootstrap aggregation**, or **bagging**, trains the learning algorithm M times as follows:

For $j = 1, 2, \dots, M$

- Draw a new bootstrap sample \tilde{D}_j from D
- Train the algorithm on \tilde{D}_j , giving hypothesis \hat{h}_j



Bootstrap Aggregation (Bagging)

Given a training sample $D = ((X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n))$ of size n , a **bootstrap sample** is obtained by drawing n examples *with replacement* from D .

For a given learning algorithm, **bootstrap aggregation**, or **bagging**, trains the learning algorithm M times as follows:

For $j = 1, 2, \dots, M$

- Draw a new bootstrap sample \tilde{D}_j from D
- Train the algorithm on \tilde{D}_j , giving hypothesis \hat{h}_j

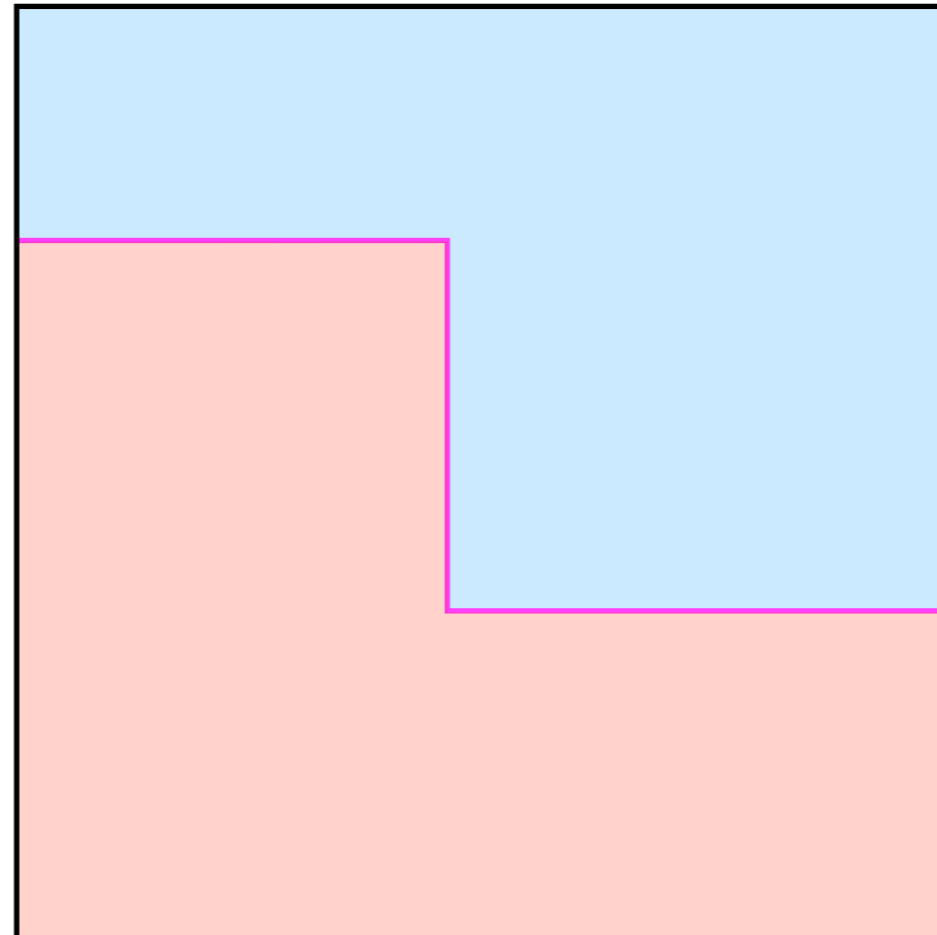
When given a new example, the bagged predictor predicts:

- the **majority vote** among $\hat{h}_1, \hat{h}_2, \dots, \hat{h}_M$ for classification

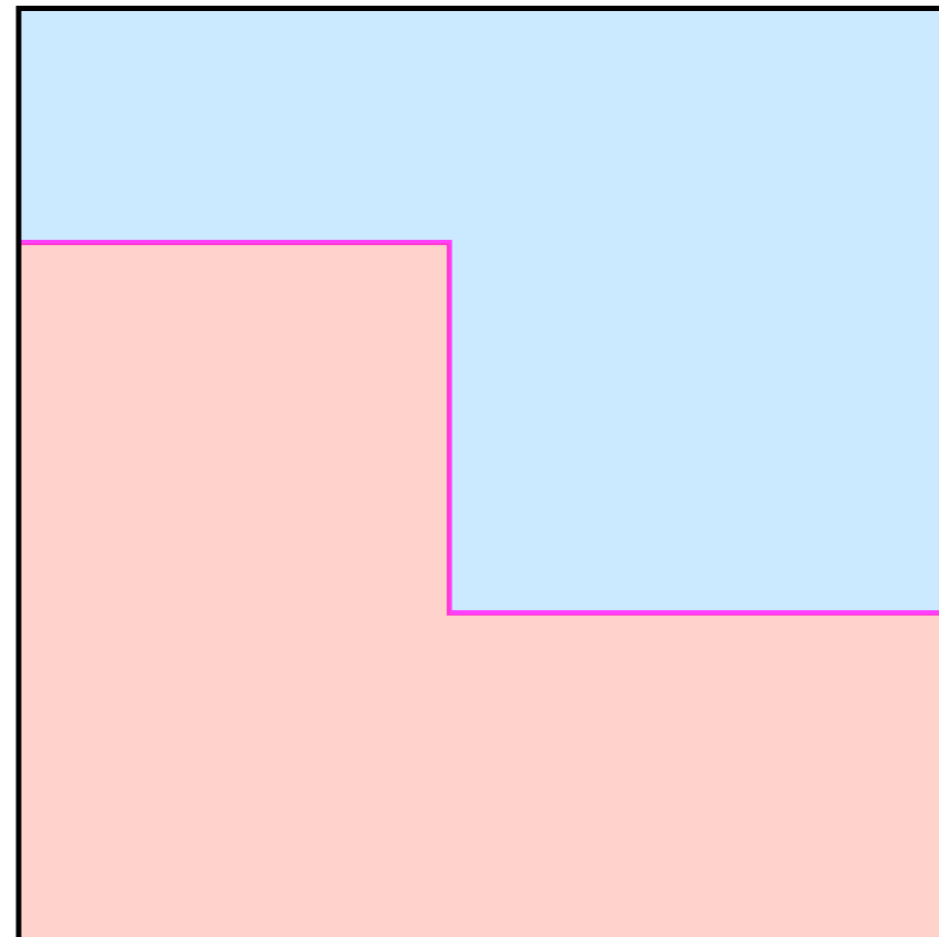
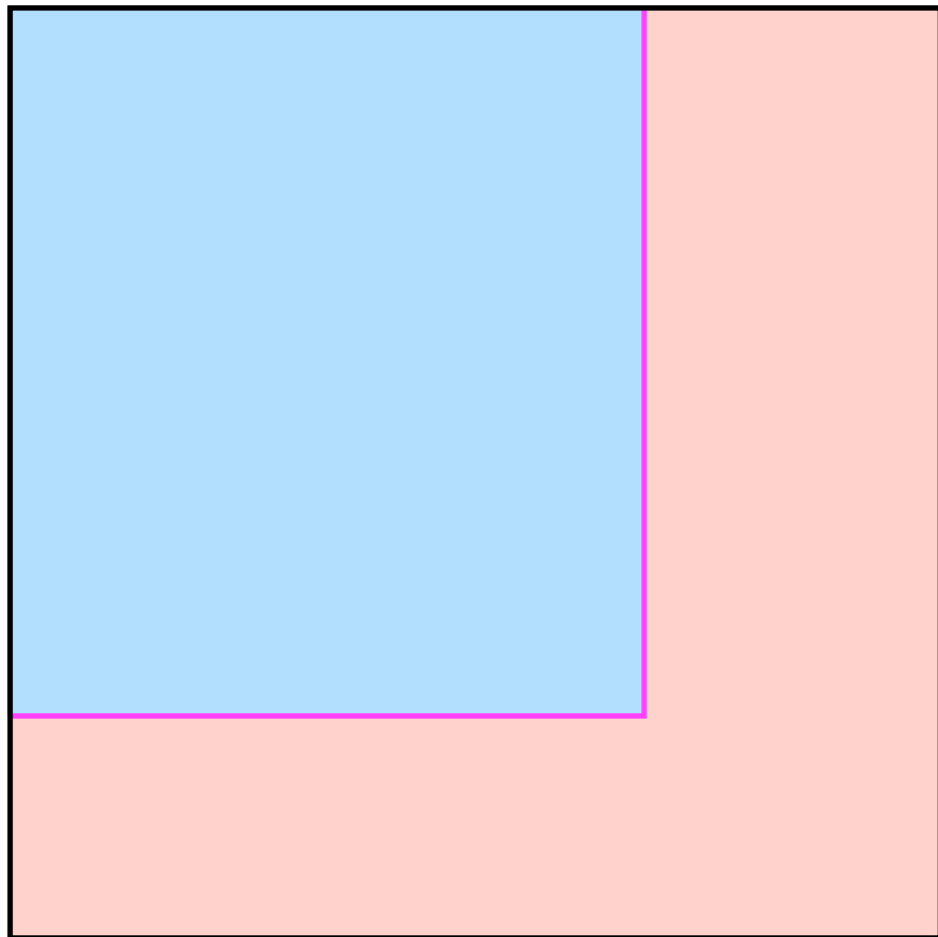
$$\arg \max_{y \in \{1, 2, \dots, k\}} |\{j \in [M] : \hat{h}_j(x) = y\}|$$

- the mean of $\hat{h}_1, \hat{h}_2, \dots, \hat{h}_M$ for regression

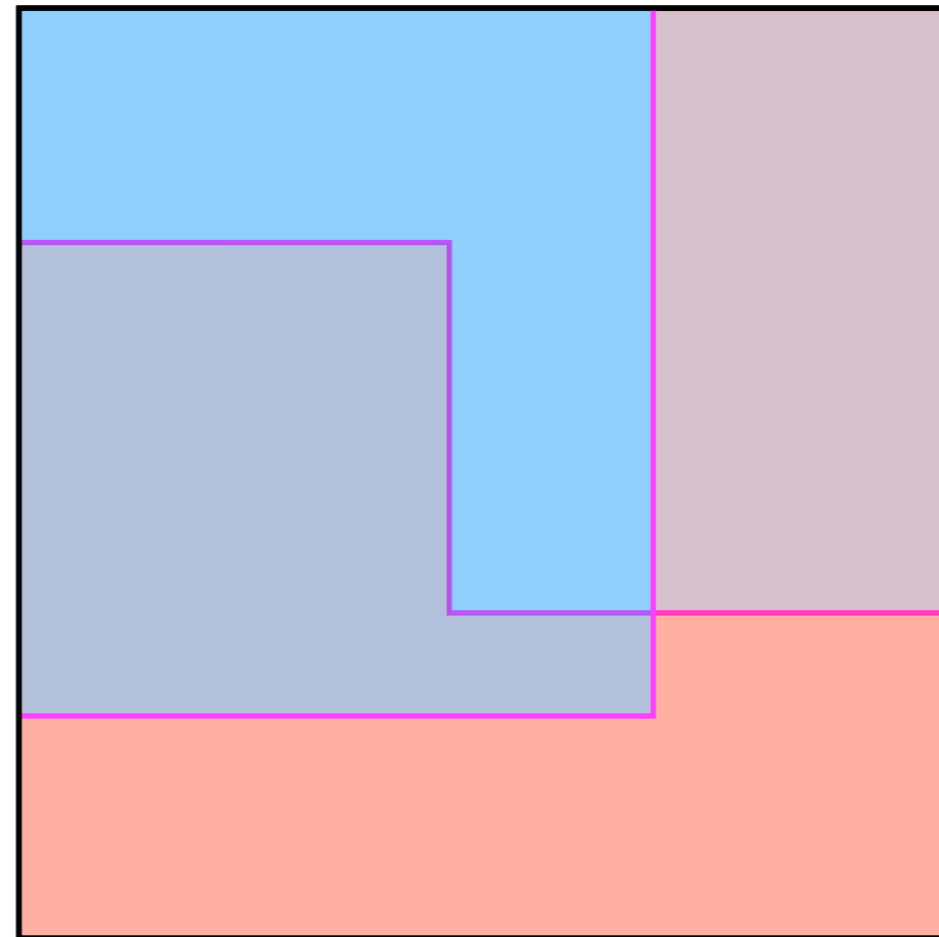
Geometry: One decision tree



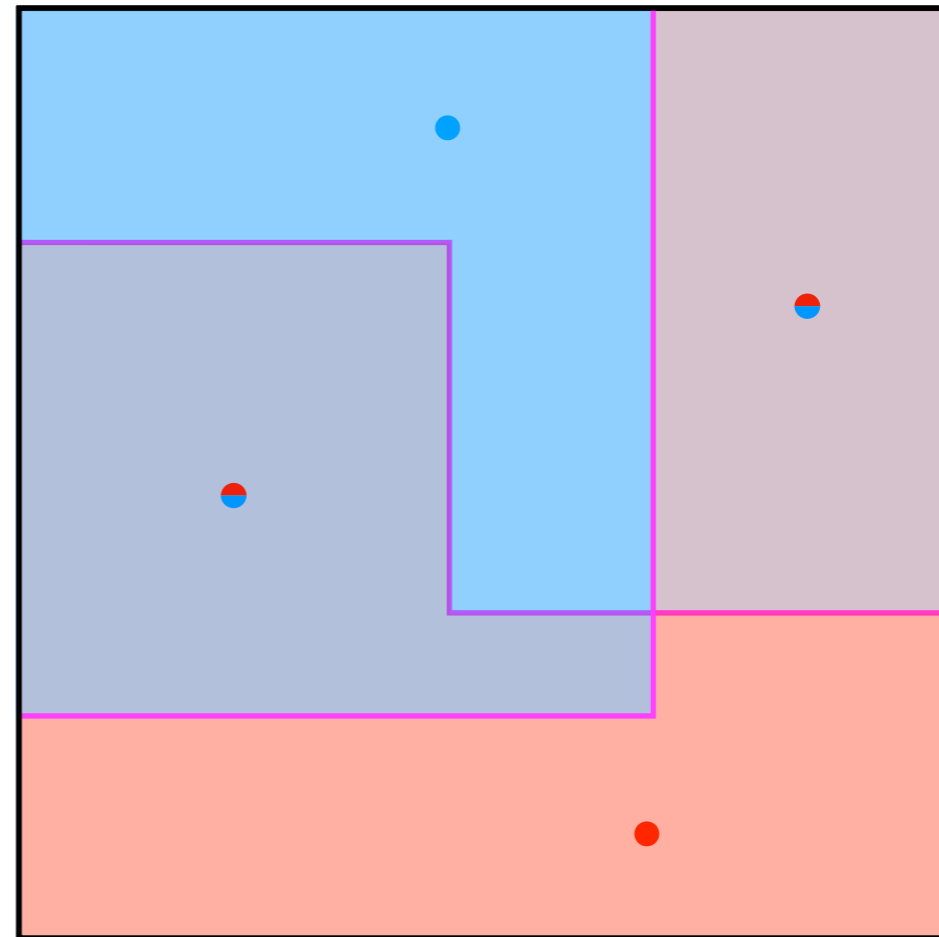
Geometry: Towards majority vote of 2 decision trees



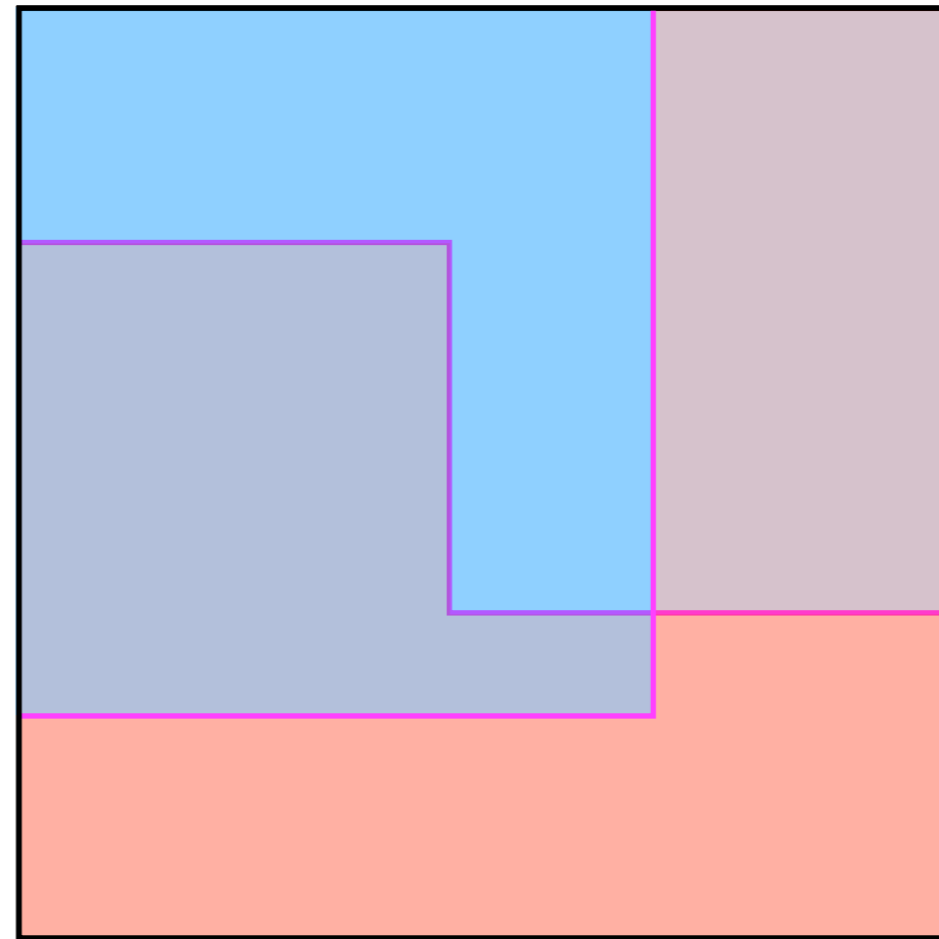
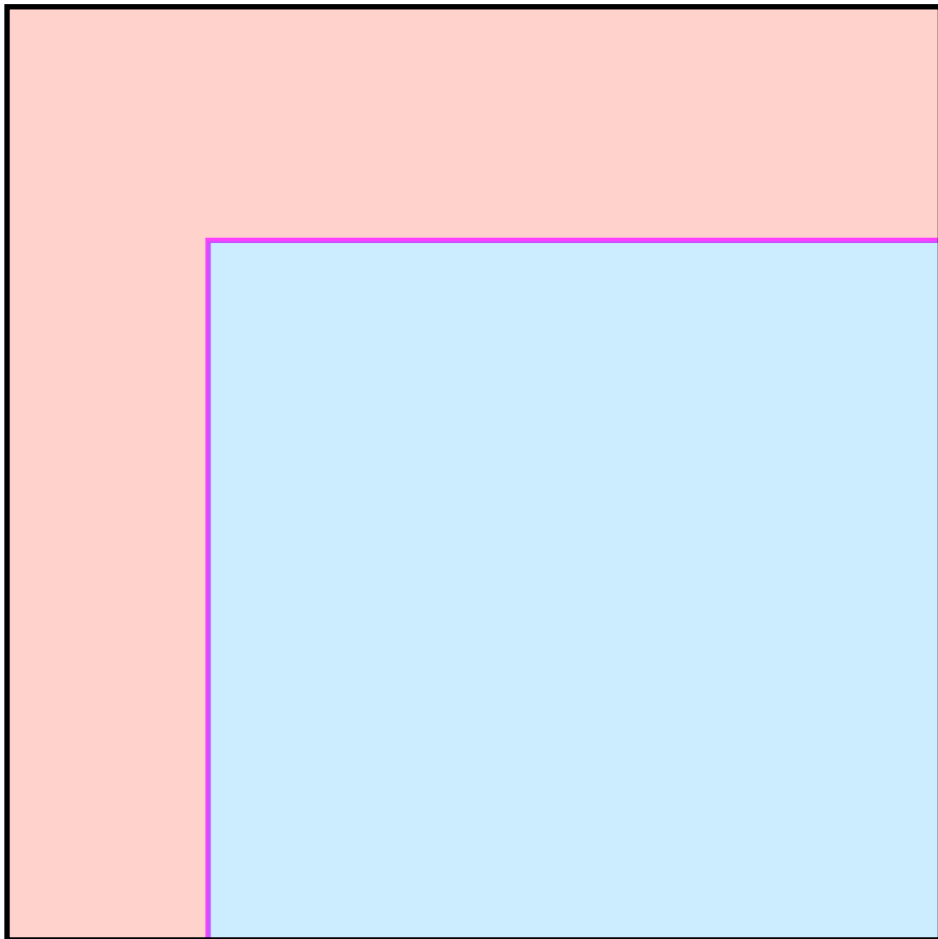
Geometry: Majority vote of 2 decision trees



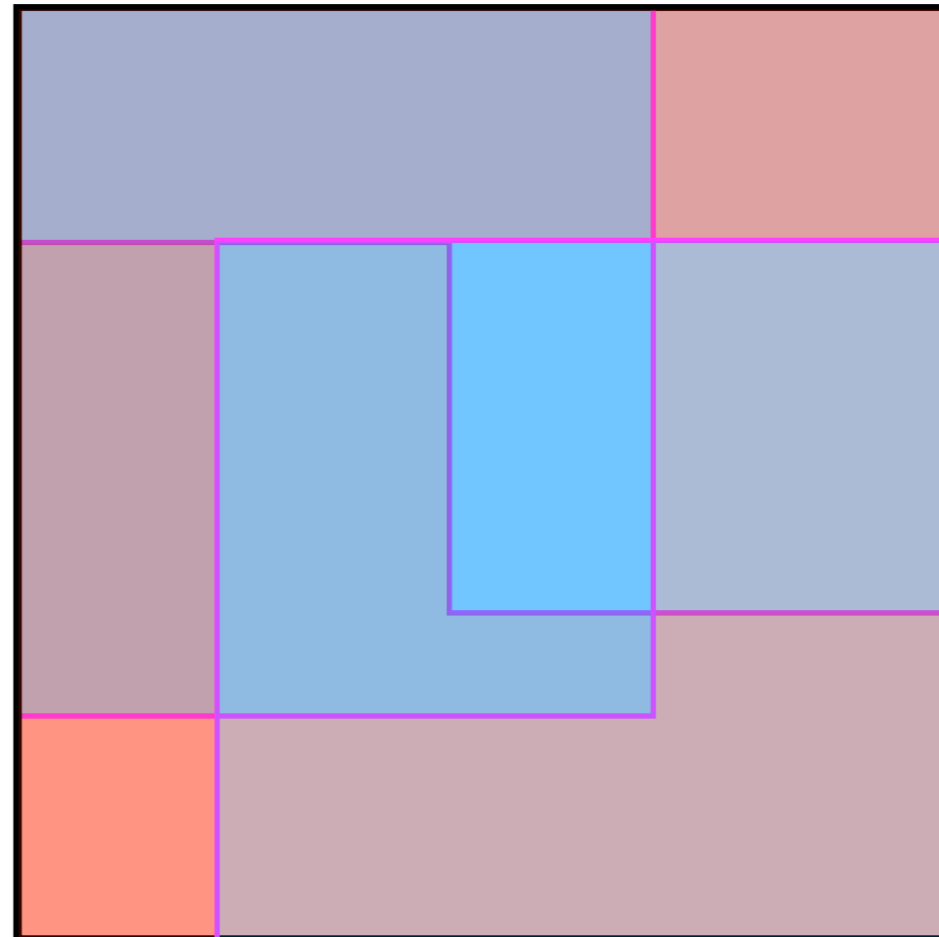
Geometry: Majority vote of 2 decision trees



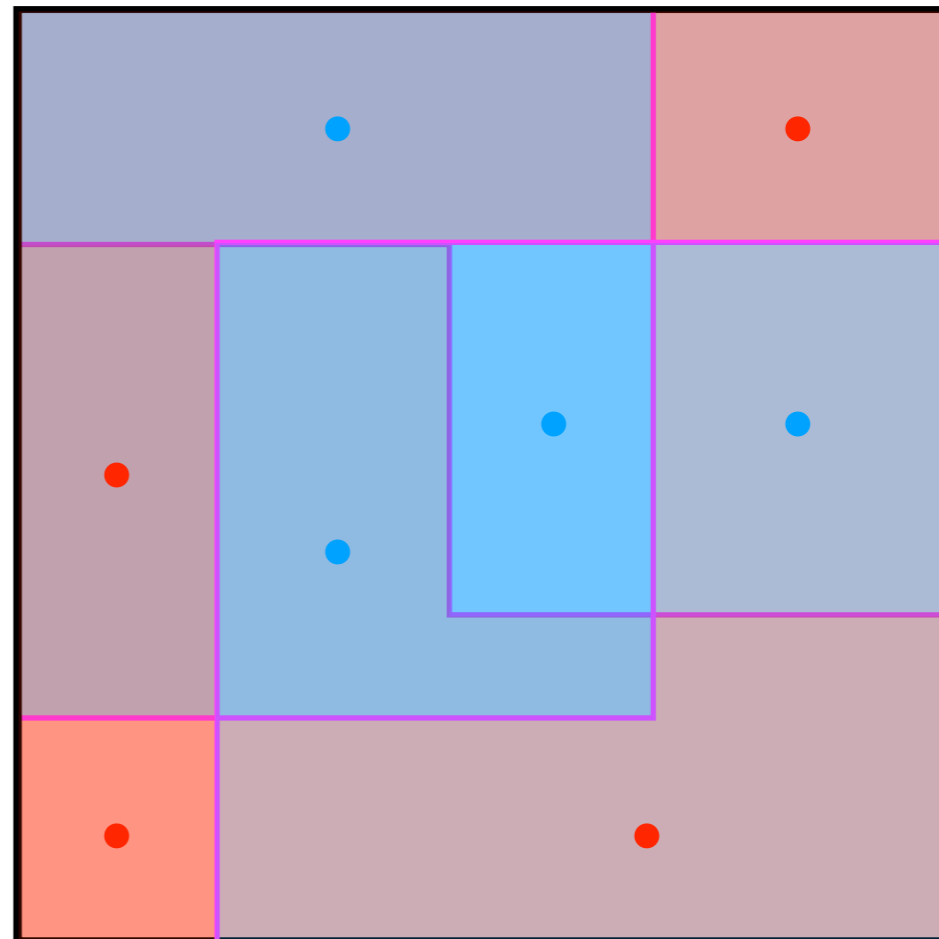
Geometry: Towards majority vote of 3 decision trees



Geometry: Majority vote of 3 decision trees



Geometry: Majority vote of 3 decision trees



Bagged Trees

Each tree is constructed as follows:

Input

- a training set of n examples with d features
- Number of trees, M

Algorithm

For each of the M trees:

- Sample (with replacement) n examples from the training set
- Train the tree using the sample

Does bagging help? Yes for regression (variance reduction).

For classification, could be better or worse than single decision tree.

Random Forest Construction

Each tree is constructed as follows:

Input

- a training set of n examples with d features
- n' ($\leq n$) - number of examples in each decision tree's training set
- d' ($< d$) - number of features used to determine the decision at each node of the tree
- Number of trees, M

Algorithm

For each of the M trees:

- Sample (with replacement) n' examples from the training set
- During tree construction, for each node of the tree, randomly select d' distinct features on which to base the decision at the node. So, the best split will be computed based on these d' features.

Random Forest Construction

Each tree is constructed as follows:

Input

- a training set of n examples with d features
- n' ($\leq n$) - number of examples in each decision tree's training set
- d' ($< d$) - number of features used to determine the decision at each node of the tree
- Number of trees, M

Algorithm

For each of the M trees:

- Sample (with replacement) n' examples from the training set
- During tree construction, for each node of the tree, randomly select d' distinct features on which to base the decision at the node. So, the best split will be computed based on these d' features.

Rules of Thumb

Try $n' = n$ (bootstrap samples)

Classification: try $d' = \sqrt{d}$

Regression: try $d' = d/3$

Prediction using a Random Forest

- (1) Get a prediction for each tree T in the forest \mathcal{T}
- (2) Take the majority vote of the predictions (for classification)

$$\text{Predict } \arg \max_{y \in \{1, 2, \dots, k\}} |\{T \in \mathcal{T} : T(x) = y\}|$$




set of potential labels

Random Forests

Quite old, and yet one of the best methods out there

Intuition for why they work well:

- Trees have low bias (high representational power) but can have high variance (high sensitivity to particular draw of training sample). Achieve some variance reduction by aggregating over many trees trained on bootstrapped data
- If errors of individual models tend not to be correlated (different models make different types of errors), the majority vote can far outperform the best individual model



Random forests have this advantage,
but bagged trees do not. Why?

Random Forests

Good news: Good performance in practice

Bad news: A random forest hypothesis is not easily interpretable (whereas a single decision tree typically is interpretable)