

INTRO to DATA SCIENCE

LECTURE 12: ENSEMBLE TECHNIQUES

INTRO TO DATA SCIENCE, REGRESSION & REGULARIZATION

DATA SCIENCE IN THE NEWS

DATA SCIENCE IN THE NEWS

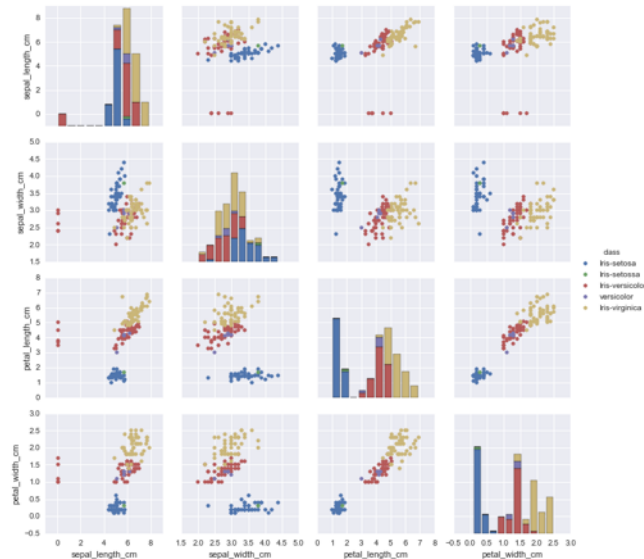
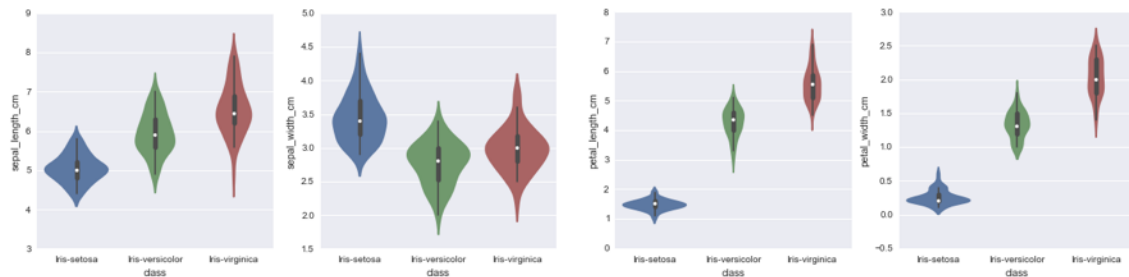
An example machine learning notebook

Notebook by [Randal S. Olson](#)

Supported by [Jason H. Moore](#)

[University of Pennsylvania Institute for Bioinformatics](#)

It is recommended to [view this notebook in nbviewer](#) for the best viewing experience.



Generating Poetry with PoetRNN

Aug 11, 2015

*sunset home
the bird of windows
white the hand*

*rain pond
the light of shadow
of the song*

LAST TIME:

I. DECISION TREES

II. LAB ON DECISION TREES

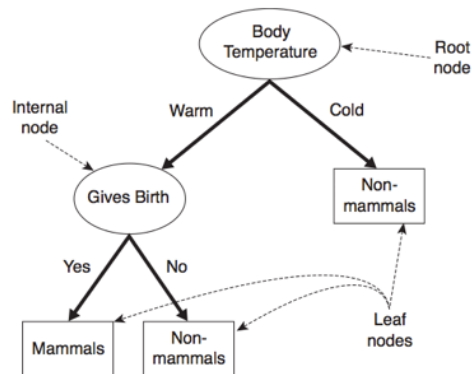


Figure 4.4. A decision tree for the mammal classification problem.

INTRO TO DATA SCIENCE

QUESTIONS?

WHAT WAS THE MOST INTERESTING THING YOU LEARNT?

WHAT WAS THE HARDEST TO GRASP?

I. ENSEMBLE TECHNIQUES

II. PROBLEMS IN CLASSIFICATION

III. BAGGING

IV. BOOSTING

V. RANDOM FORESTS

EXERCISE:

VI. LAB

KEY OBJECTIVES

- **UNDERSTAND THE POWER OF USING ENSEMBLE CLASSIFIERS**
- **KNOW THE DIFFERENCE BETWEEN A BASE CLASSIFIER AND AN ENSEMBLE CLASSIFIER**
- **UNDERSTAND WHAT BAGGING IS AND BE ABLE TO DESCRIBE IT**
- **UNDERSTAND WHAT BOOSTING IS AND BE ABLE TO DESCRIBE IT**
- **UNDERSTAND WHAT A RANDOM FOREST IS AND BE ABLE TO USE IT**

I. ENSEMBLE TECHNIQUES

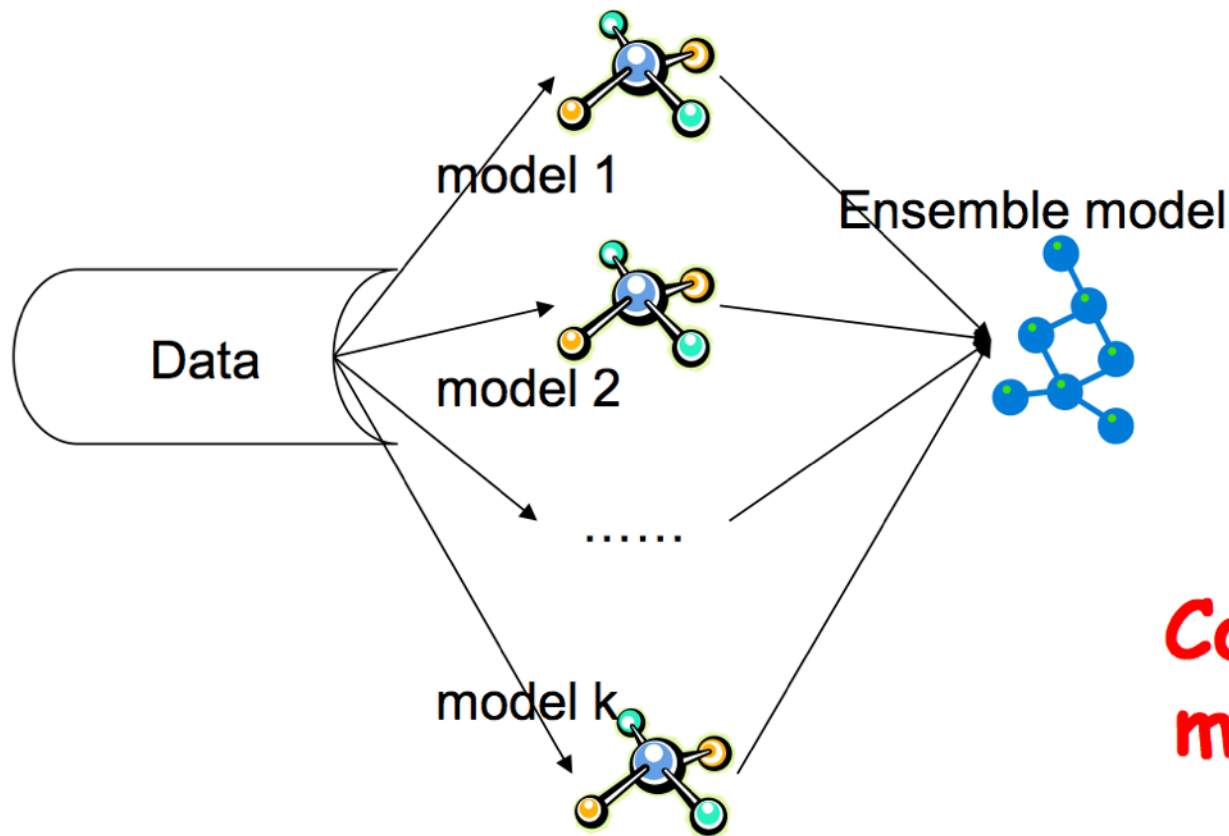
So far, we have only discussed
individual classifiers

What are these?

What is your favorite?

How can we come up with
a better model?

Can we combine multiple classifiers
to produce a better classifier?



**Combine multiple
models into one!**

NETFLIX Watch Instantly | Your Account & Help


Movies, TV shows, actors, directors, genres

[Watch Instantly](#) [Browse DVDs](#) [Your Queue](#) [Movies You'll ♥](#)

Congratulations! Movies we think **You** will ♥

Add movies to your Queue, or **Rate** ones you've seen for even better suggestions.

Spider-Man 3



Add

★★★★☆

☐ Not Interested

300



Add

★★★★☆

☐ Not Interested

The Rundown



Add

★★★★☆

☐ Not Interested

Bad Boys II



Add

★★★★☆

☐ Not Interested

Las Vegas: Season 2
(6-Disc Series)



★★★★☆

☐ Not Interested

The Last Samurai



★★★★☆

☐ Not Interested

Star Wars: Episode III



★★★★☆

☐ Not Interested

Robot Chicken: Season 3
(2-Disc Series)



★★★★☆

☐ Not Interested



award **\$1 million** to anyone
who can improve movie
recommendation by 10%

Supervised learning task

- Training data is a set of users and ratings (1,2,3,4,5 stars) those users have given to movies.

Supervised learning task

- Construct a classifier that given a user and an unrated movie, correctly classifies that movie as either 1, 2, 3, 4, or 5 stars

At first, single-model methods
are developed, and
performances are improved

However, improvements
slowed down

Later, individuals and teams
merged their results,
and significant improvements
are observed

Leaderboard

24

Rank	Team Name	Best Test Score	% Improvement	Best Submit Time
Grand Prize - RMSE = 0.8567 - Winning Team: BellKor's Pragmatic Chaos				
1	BellKor's Pragmatic Chaos	0.8567	10.06	2009-07-26 18:18:28
2	The Ensemble	0.8567	10.06	2009-07-26 18:38:22
3	Grand Prize Team	0.8582	9.90	2009-07-10 21:24:40
4	Opera Solutions and Vandelay United	0.8588	9.84	2009-07-10 01:12:31
5	Vandelay Industries !	0.8591	9.81	2009-07-10 00:32:20
6	PragmaticTheory	0.8594	9.77	2009-06-24 12:06:56
7	BellKor in BigChaos	0.8601	9.70	2009-05-13 08:14:09
8	Dace	0.8612	9.59	2009-07-24 17:18:43
9	Feeds2	0.8622	9.48	2009-07-12 13:11:51
10	BigChaos	0.8623	9.47	2009-04-07 12:33:59

“Our final solution (RMSE=0.8712) consists of blending 107 individual results. “

12	BellKor	0.8624	9.46	2009-07-26 17:19:11
Progress Prize 2008 - RMSE = 0.8627 - Winning Team: BellKor in BigChaos				
13	xiangliang	0.8642	9.27	2009-07-15 14:53:22
14	Gravity	0.8643	9.26	2009-04-22 18:31:32
15	Ces	0.8651	9.18	2009-06-21 19:24:53

“Predictive accuracy is substantially improved when blending multiple predictors. Our experience is that most efforts should be concentrated in deriving substantially different approaches, rather than refining a single technique. “

Progress Prize 2007 - RMSE = 0.8725 - Winning Team: Korben

Cinematix score - RMSE = 0.9525

Q: What are ensemble techniques?

Q: What are ensemble techniques?

A: Methods of improving classification accuracy by aggregating predictions over several base classifiers.

Q: What are ensemble techniques?

A: Methods of improving classification accuracy by aggregating predictions over several base classifiers.

Ensembles are often much more accurate than the base classifiers that compose them.

Q: What are ensemble techniques?

A: Methods of improving classification accuracy by aggregating predictions over several base classifiers.

Ensembles are often much more accurate than the base classifiers that compose them.

NOTE

Base classifiers and ensemble classifiers are sometimes called weak learners and strong learners.

In order for an ensemble classifier to outperform a single base classifier, the following conditions must be met:

In order for an ensemble classifier to outperform a single base classifier, the following conditions must be met:

1) **accuracy**: *base classifiers outperform random guessing*

In order for an ensemble classifier to outperform a single base classifier, the following conditions must be met:

- 1) **accuracy**: *base classifiers outperform random guessing*
- 2) **diversity**: *misclassifications must occur on different training examples*

In order for an ensemble classifier to outperform a single base classifier, the following conditions must be met:

- 1) **accuracy**: *low bias*
- 2) **diversity**: *uncorrelated*

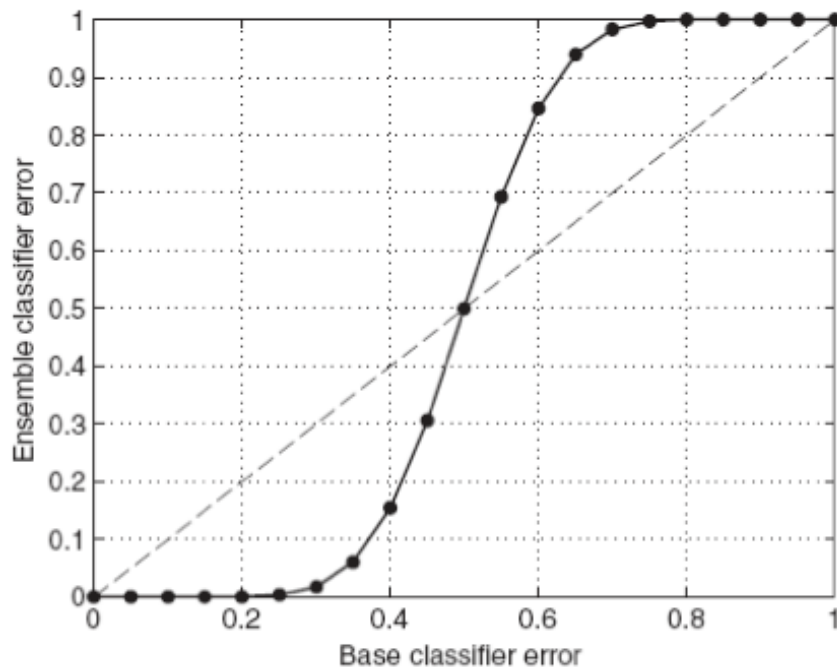
In order for an ensemble classifier to outperform a single base classifier, the following conditions must be met:

- 1) **accuracy**: *low bias*
- 2) **diversity**: *uncorrelated*

NOTE

Ideally, we would also like the base classifiers to be unstable to variations in the training set.

In other words, high variance.

**NOTE**

dashed line = perfectly correlated bc's (no improvement using ensemble)

solid line = perfectly uncorrelated bc's (some improvement for unbiased bc's)

Figure 5.30. Comparison between errors of base classifiers and errors of the ensemble classifier.

Quick check

what base classifiers would you combine to have very different perspectives?

PROBLEMS IN CLASSIFICATION

In any supervised learning task, our goal is to make predictions of the true classification function f by learning the classifier h .

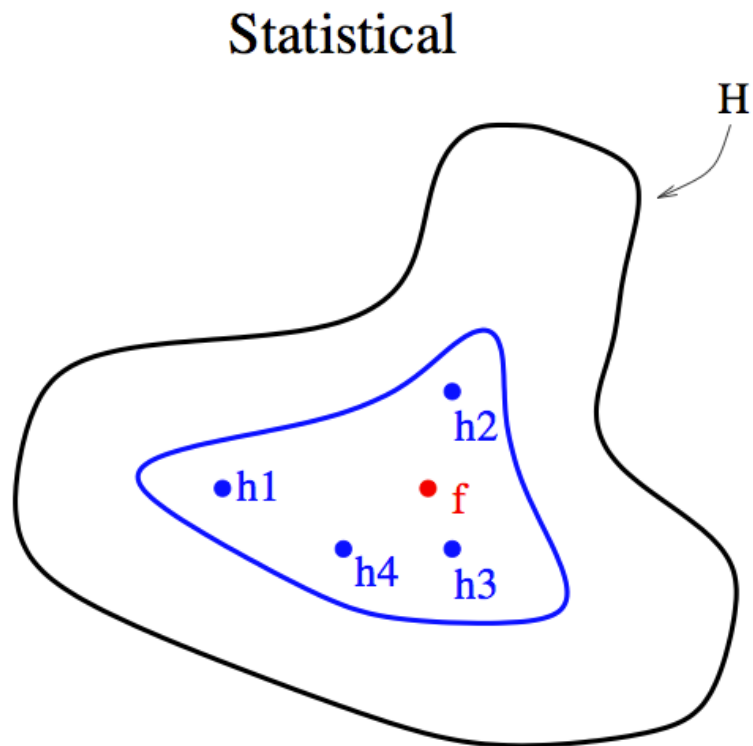
In any supervised learning task, our goal is to make predictions of the true classification function f by learning the classifier h .

There are three main problems that can prevent this:

- statistical problem*
- computational problem*
- representational problem*

If the amount of training data available is small, the base classifier will have difficulty converging to h .

An ensemble classifier can mitigate this problem by “averaging out” base classifier predictions to improve convergence.

**NOTE**

The true function f is best approximated as an average of the base classifiers.

Even with sufficient training data, it may still be computationally difficult to find the best classifier h .

For example, if our base classifier is a decision tree, an exhaustive search of the hypothesis space of all possible classifiers is extremely complex (NP-complete).

Even with sufficient training data, it may still be computationally difficult to find the best classifier h .

For example, if our base classifier is a decision tree, an exhaustive search of the hypothesis space of all possible classifiers is extremely complex (NP-complete).

NOTE

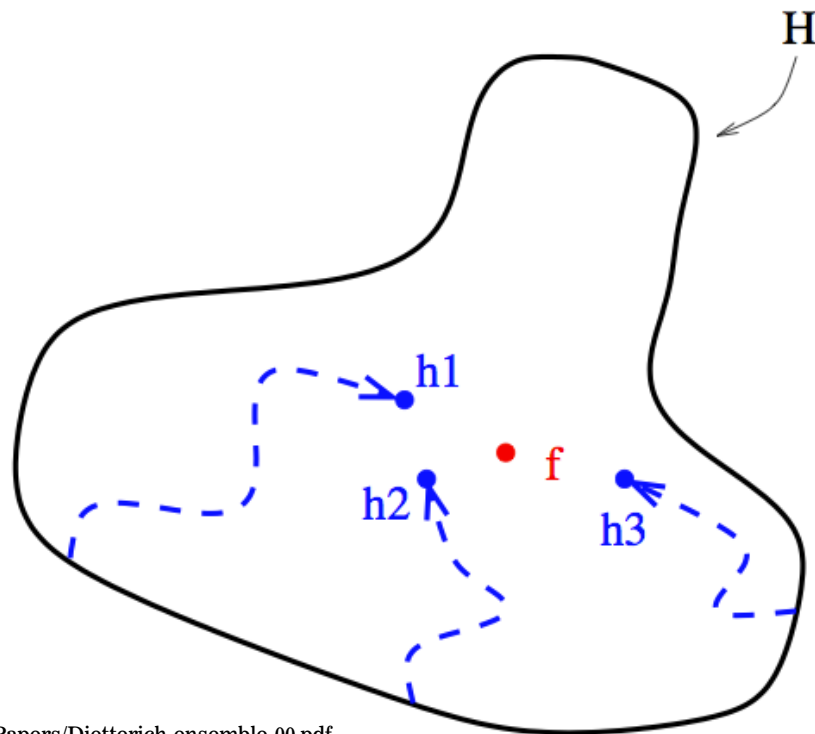
Recall that this is why we used a heuristic algorithm (greedy search).

Even with sufficient training data, it may still be computationally difficult to find the best classifier h .

For example, if our base classifier is a decision tree, an exhaustive search of the hypothesis space of all possible classifiers is extremely complex (NP-complete).

An ensemble composed of several BC's with different starting points can provide a better approximation to f than any individual BC.

Computational

**NOTE**

The true function f is often best approximated by using several starting points to explore the hypothesis space.

Sometimes f cannot be expressed in terms of our hypothesis at all.

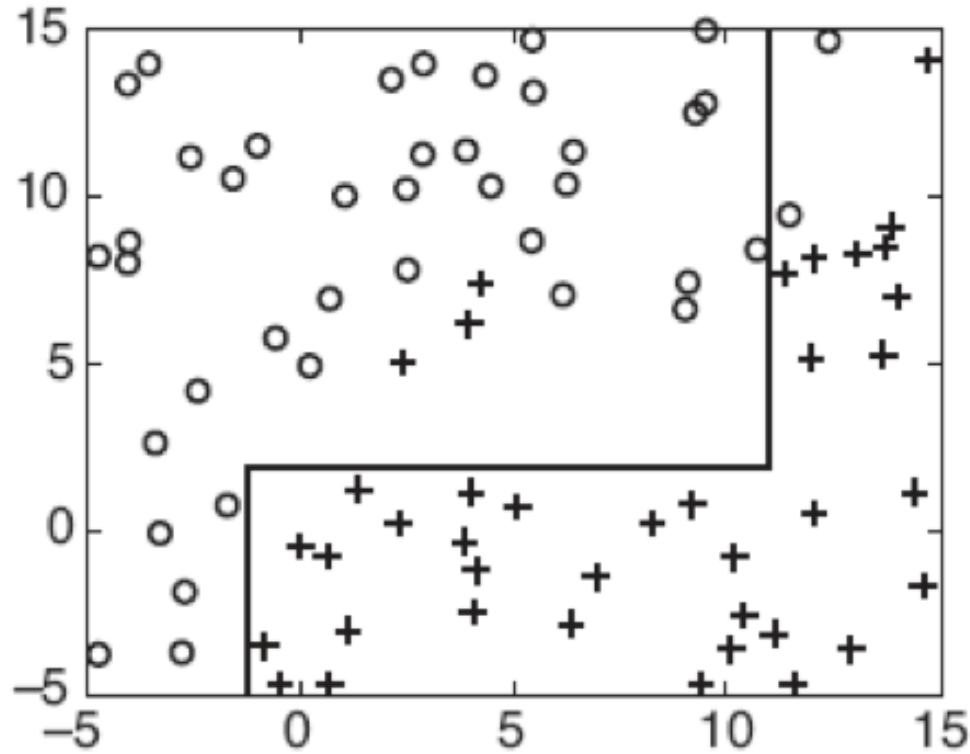
Sometimes f cannot be expressed in terms of our hypothesis at all.

To illustrate this, suppose we use a decision tree as our base classifier.

Sometimes f cannot be expressed in terms of our hypothesis at all.

To illustrate this, suppose we use a decision tree as our base classifier.

A decision tree works by forming a rectilinear partition of the feature space.



NOTE

What is a rectilinear decision boundary?

One whose segments are orthogonal to the x & y axes.

But what if f is a diagonal line?

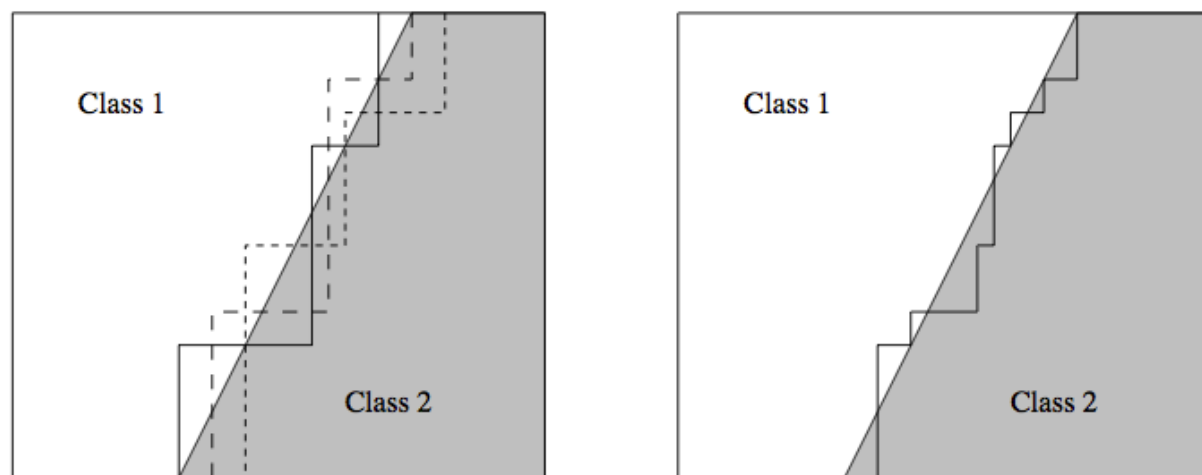
But what if f is a diagonal line?

Then it cannot be represented by finitely many rectilinear segments, and therefore the true decision boundary cannot be obtained by a decision tree classifier.

But what if f is a diagonal line?

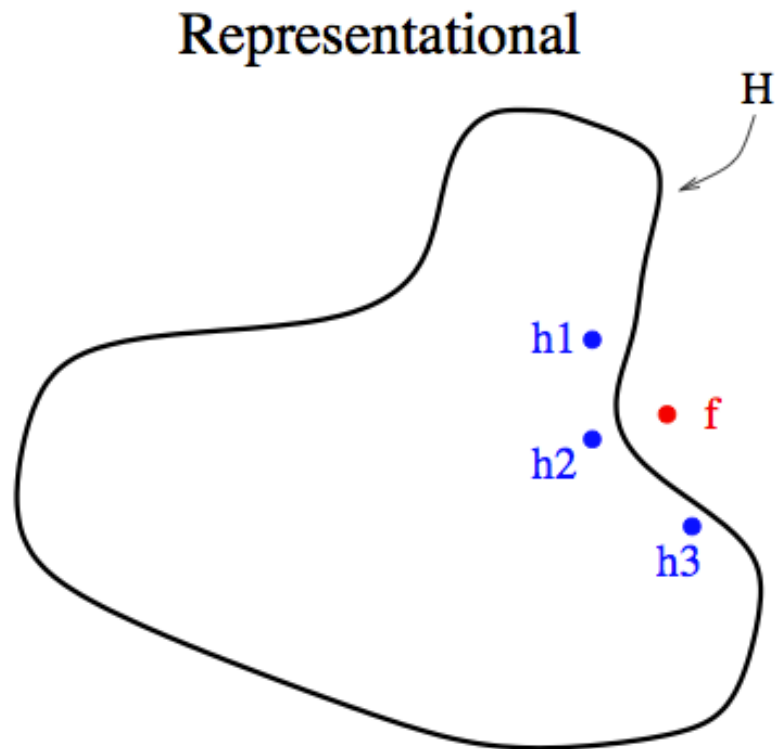
Then it cannot be represented by finitely many rectilinear segments, and therefore the true decision boundary cannot be obtained by a decision tree classifier.

However, it may be still be possible to approximate f or even to expand the space of representable functions using ensemble methods.

**NOTE**

An ensemble of decision trees can approximate a diagonal decision boundary.

Fig. 4. The left figure shows the true diagonal decision boundary and three staircase approximations to it (of the kind that are created by decision tree algorithms). The right figure shows the voted decision boundary, which is a much better approximation to the diagonal boundary.



NOTE

Ensemble classifiers can be effective even if the true decision boundary lies outside the hypothesis space.

Q: How do you create an ensemble classifier?

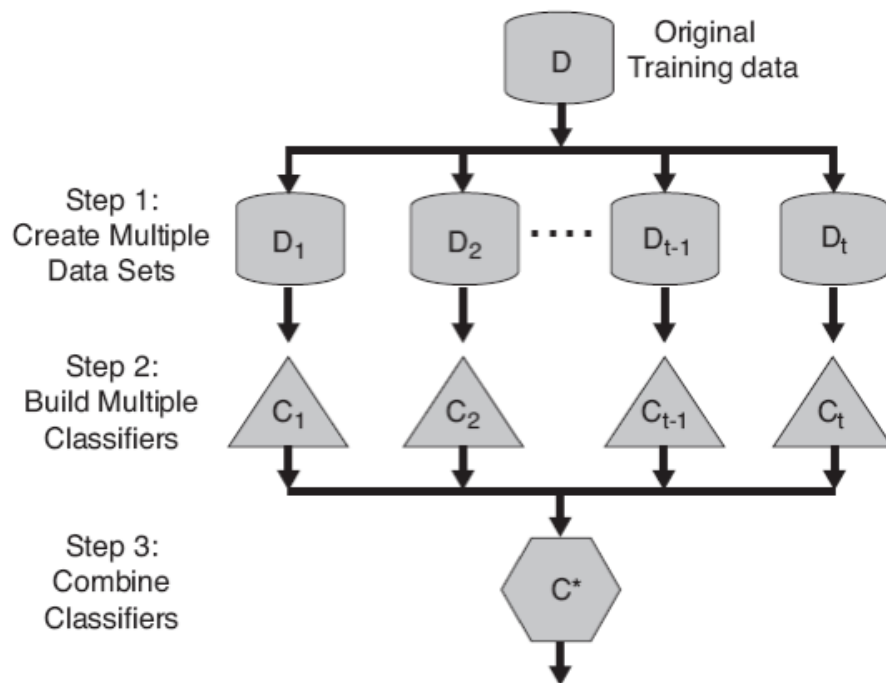


Figure 5.31. A logical view of the ensemble learning method.

Q: How do you generate several base classifiers?

Q: How do you generate several base classifiers?

A: There are several ways to do this:

- manipulating the training set*
- manipulating the output labels*
- manipulating the learning algorithm itself*

We will talk about a few examples of each of these.

Quick check

What could go wrong when you use decision trees for classification?

How could you mitigate the problem?

Discuss with the person next to you

INTRO TO DATA SCIENCE

LAB 1

BAGGING

Bagging (***b**ootstrap **a**ggregating*) is a method that involves manipulating the training set by resampling.

Bagging (*bootstrap aggregating*) is a method that involves manipulating the training  by resampling.

We learn k base classifiers on k different samples of training data.

These samples are independently created by resampling the training data using uniform weights (eg, a uniform sampling distribution).

Bagging (*bootstrap aggregating*) is a method that involves manipulating the training set by resampling.

We learn k base classifiers on k different samples of training data.

These samples are independently created by resampling the training data using uniform weights (eg, a uniform sample distribution).

NOTE

Each training sample is the same size as the original training set.

Bagging (*bootstrap aggregating*) is a method that involves manipulating the training set by resampling.

We learn k base classifiers on k different samples of training data.

These samples are independently created by resampling the training data using uniform weights (eg, a uniform sample distribution).

NOTE

Resampling means that some training records may appear in a sample more than once, or even not at all.

Bagging (*bootstrap aggregating*) is a method that involves manipulating the training set by resampling.

We learn k base classifiers on k different samples of training data.

These samples are independently created by resampling the training data using uniform weights (eg, a uniform sampling distribution).

The final prediction is made by taking a majority vote across bc 's.

BAGGING EXAMPLE

Original	1	2	3	4	5	6	7	8
Training set 1	2	7	8	3	7	6	3	1
Training set 2	7	8	5	6	4	2	7	1
Training set 3	3	6	2	7	5	6	2	2
Training set 4	4	5	1	4	6	4	3	8

Bagging reduces the variance in our generalization error by aggregating multiple base classifiers together (provided they satisfy our earlier requirements).

Bagging reduces the variance in our generalization error by aggregating multiple base classifiers together (provided they satisfy our earlier requirements).

If the base classifier is stable, then the ensemble error is primarily due to bc bias, and bagging may not be effective.

Bagging reduces the variance in our generalization error by aggregating multiple base classifiers together (provided they satisfy our earlier requirements).

If the base classifier is stable, then the ensemble error is primarily due to bc bias, and bagging may not be effective.

Since each sample of training data is equally likely, bagging is not very susceptible to overfitting with noisy data.

BOOSTING

Boosting is an iterative procedure that adaptively changes the sampling distribution of training records at each iteration.

Boosting is an iterative procedure that adaptively changes the sampling distribution of training records at each iteration.

The first iteration uses uniform weights (like bagging). In subsequent iterations, the weights are adjusted to emphasize records that were misclassified in previous iterations.

Boosting is an iterative procedure that adaptively changes the sampling distribution of training records at each iteration.

The first iteration uses uniform weights (like bagging). In subsequent iterations, the weights are adjusted to emphasize records that were misclassified in previous iterations.

The final prediction is constructed by a weighted vote (where the weights for a bc depends on its training error).

Boosting is an iterative procedure that adaptively changes the sampling distribution of training records at each iteration.

The first iteration uses uniform weights (like bagging). In subsequent iterations, the weights are adjusted to emphasize records that were misclassified in previous iterations.

NOTE

The bc's focus more and more closely on records that are difficult to classify as the sequence of iterations progresses.

Thus the bc's are faced with progressively more difficult learning problems.

The final prediction is constructed by a weighted vote (where the weights for a bc depends on its training error).

Like in bagging, sampling is done with replacement, and as a result some records may not appear in a given training sample.

Like in bagging, sampling is done with replacement, and as a result some records may not appear in a given training sample.

These omitted records will likely be misclassified, and given greater weight in subsequent iterations once the sampling distribution is updated.

Like in bagging, sampling is done with replacement, and as a result some records may not appear in a given training sample.

These omitted records will likely be misclassified, and given greater weight in subsequent iterations once the sampling distribution is updated.

So even if a record is left out at one stage, it will be emphasized later.

Updating the sampling distribution and forming an ensemble prediction leads to a combination of the base classifiers.

Updating the sampling distribution and forming an ensemble prediction leads to a combination of the base classifiers.

By explicitly trying to optimize the weighted ensemble vote, boosting attacks the representation problem head-on.

INTRO TO DATA SCIENCE

RANDOM FORESTS

RANDOM FORESTS

A random forest is an ensemble of decision trees where each base classifier is grown using a random effect.

A random forest is an ensemble of decision trees where each base classifier is grown using a random effect.

but how?

A random forest is an ensemble of decision trees where each base classifier is grown using a random effect.

One way to do this is to randomly choose one of the top k features to split each node.

A random forest is an ensemble of decision trees where each base classifier is grown using a random effect.

One way to do this is to randomly choose one of the top k features to split each node.

For a small number of features, we can also create linear combinations of features and select splits from the enhanced feature set (Forest-RC).

A random forest is an ensemble of decision trees where each base classifier is grown using a random effect.

One way to do this is to randomly choose one of the top k features to split each node.

For a small number of features, we can also create linear combinations of features and select splits from the enhanced feature set (Forest-RC).

Or, we can select splitting features completely at random (Forest-RI).

Random forests are about as accurate as boosting methods, more robust to noise, and can also have better runtime than other ensemble methods (since the feature space is reduced in some cases).

Quick check

Is there a startup that provides random forests as a service?

If so, what's their competitive advantage?

search the web and find an answer...

INTRO TO DATA SCIENCE

LAB