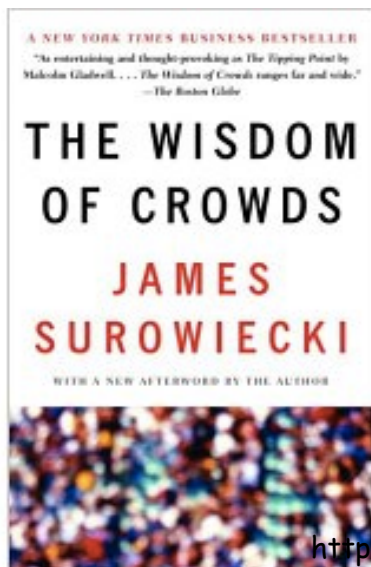
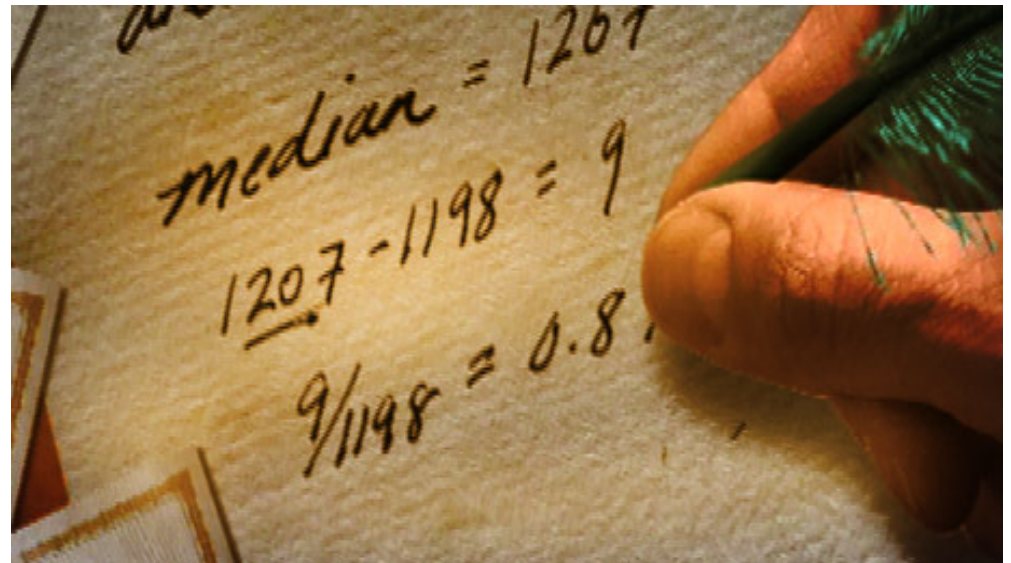

Ensemble learning



<http://unanimous.ai/emergent-intelligence-from-a-jar-of->

The wisdom of the crowds

Sir Francis Galton discovered in the early 1900s that a collection of educated guesses can add up to very accurate predictions!



The paper in which he describes these findings:
Vox Populi. Nature 75, pages 450-451, 1907.
<http://galton.org/essays/1900-1911/galton-1907-vox-populi.pdf>

Image from <http://www.pbs.org/wgbh/nova/physics/wisdom-crowds.html>

Ensemble methods

Intuition: averaging measurements can lead to more reliable estimates

Need: an ensemble of different models from the same training data

How to achieve diversity?

Ensemble methods

Intuition: averaging measurements can lead to more reliable estimates

Need: an ensemble of different models from the same training data

How to achieve diversity?

Training models on random subsets of examples or random subsets of features

Ensemble methods

The general strategy:

- ❖ Construct multiple, diverse predictive models from "tweaked" versions of the training data
- ❖ Combine the predictions

Bootstrap samples

Bootstrap sample: sample with replacement from a dataset

The probability that a given example is not selected for a bootstrap sample of size n :

$$(1 - 1/n)^n$$

This has a limit as n goes to infinity: $1/e = 0.368$

Conclusion: each bootstrap sample is likely to leave out about a third of the examples.

Bagging

Algorithm Bagging(D, T, \mathcal{A})

Input : data set D ; ensemble size T ; learning algorithm \mathcal{A} .

Output : ensemble of models whose predictions are to be combined by voting or averaging.

```
1 for  $t = 1$  to  $T$  do
2   |   build a bootstrap sample  $D_t$  from  $D$  by sampling  $|D|$  data points with
3   |   replacement;
4   |   run  $\mathcal{A}$  on  $D_t$  to produce a model  $M_t$ ;
5 end
6 return  $\{M_t | 1 \leq t \leq T\}$ 
```

Breiman, Leo (1996). "Bagging predictors". *Machine Learning* 24 (2): 123-140.

Bagging

Comments:

How to combine the classifiers

- ❖ Average the raw classifier scores
- ❖ Convert to probabilities before averaging
- ❖ Voting

When bagging linear classifiers, is the decision boundary of the resulting classifier linear?

Bagging

Comments:

How to combine the classifiers

- ❖ Average the raw classifier scores
- ❖ Convert to probabilities before averaging
- ❖ Voting

The decision boundary of a bagged classifier can be more complex than that of the underlying classifiers.

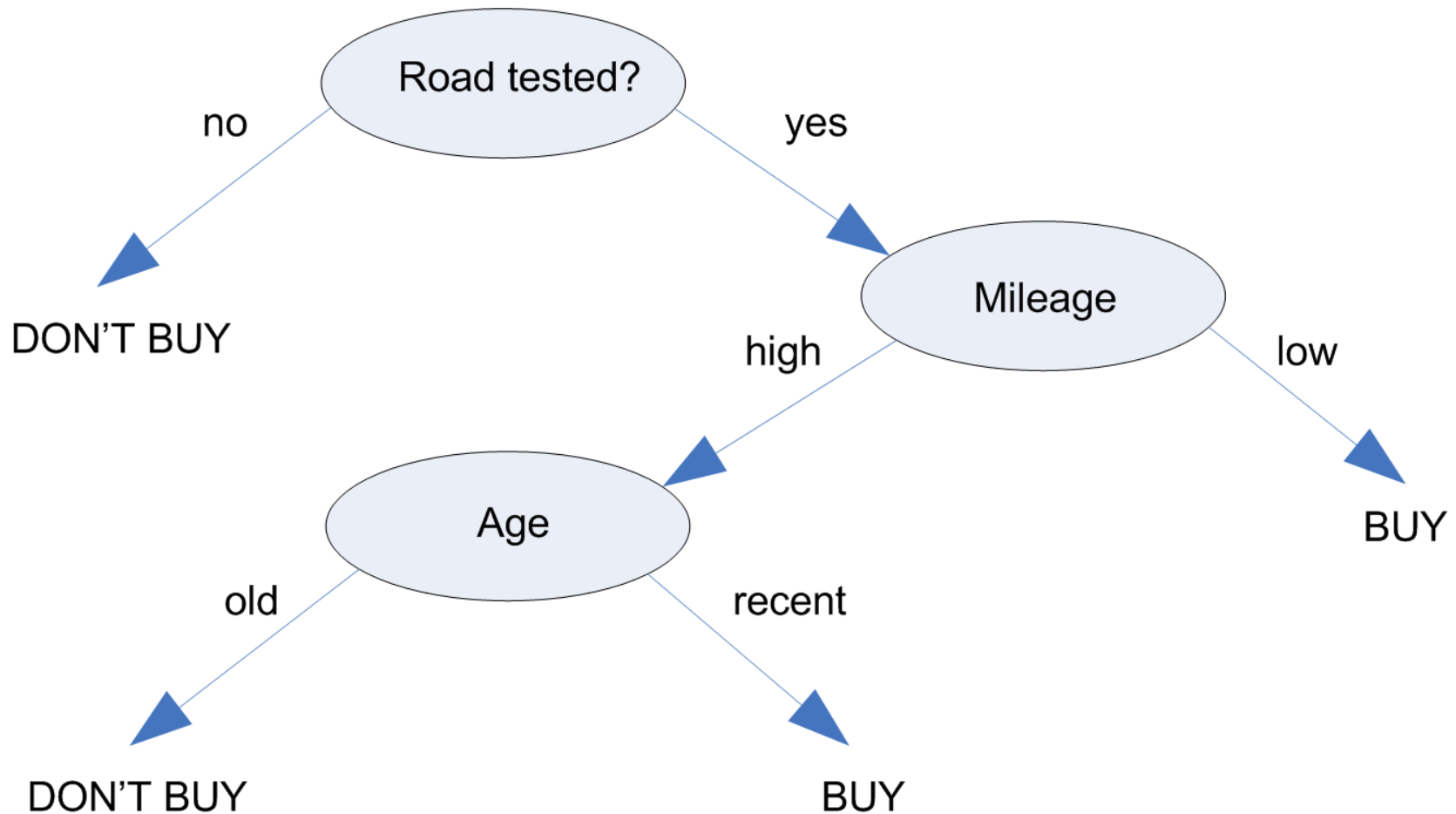
Random forests

Random forests are a special case of bagging with the additional features:

- ❖ Use **decision trees** as the base classifier
- ❖ **Sample the features** for each bootstrapping sample

Random forests. L Breiman. *Machine Learning*, 2001.

Decision tree: an example



Fast training, but limited accuracy in high dimensions

Random forests

Algorithm RandomForest(D, T, d)

Input : data set D ; ensemble size T ; subspace dimension d .

Output : ensemble of tree models whose predictions are to be combined by voting or averaging.

```
1 for  $t = 1$  to  $T$  do
2   | build a bootstrap sample  $D_t$  from  $D$  by sampling  $|D|$  data points with
   | replacement;
3   | select  $d$  features at random and reduce dimensionality of  $D_t$  accordingly;
4   | train a tree model  $M_t$  on  $D_t$  without pruning;
5 end
6 return  $\{M_t | 1 \leq t \leq T\}$ 
```

Random forests

Random forests are a special case of bagging with the additional features:

- ❖ Use **decision trees** as the base classifier
- ❖ **Sample the features** for each bootstrapping sample

Additional features:

- ❖ **Variable importance:** the values of the i^{th} feature are permuted among the training data and the out-of-bag error is computed on this perturbed data set. The importance score for feature i is computed by averaging the difference in out-of-bag error before and after the permutation over all trees.
- ❖ **Error estimation** during training using out-of-bag data.

Why do committees work?

Let's focus on regression to illustrate this point.

The output produced by the committee:

$$f_{COM}(\mathbf{x}) = \frac{1}{M} \sum_{m=1}^M f_m(\mathbf{x})$$

Let's assume the output of each model can be represented as:

$$f_m(\mathbf{x}) = y(\mathbf{x}) + \epsilon_m(\mathbf{x})$$

The expected error of an individual model:

$$\mathbb{E} \left[(f_m(\mathbf{x}) - y(\mathbf{x}))^2 \right] = \mathbb{E} \left[\epsilon_m(\mathbf{x})^2 \right]$$

Why do committees work?

The average of the errors (corresponds to the situation where each model predicts independently):

$$E_{AV} = \frac{1}{M} \sum_{m=1}^M \mathbb{E} [\epsilon_m(\mathbf{x})^2]$$

Compare that to the expected error from the committee:

$$\begin{aligned} E_{COM} &= \mathbb{E} \left[\left(\frac{1}{M} \sum_{m=1}^M f_m(\mathbf{x}) - y(\mathbf{x}) \right)^2 \right] \\ &= \mathbb{E} \left[\left(\frac{1}{M} \sum_{m=1}^M \epsilon_m(\mathbf{x}) \right)^2 \right] \end{aligned}$$

Why do committees work?

$$E_{AV} = \frac{1}{M} \sum_{m=1}^M \mathbb{E} [\epsilon_m(\mathbf{x})^2]$$
$$E_{COM} = \mathbb{E} \left[\left(\frac{1}{M} \sum_{m=1}^M f_m(\mathbf{x}) - y(\mathbf{x}) \right)^2 \right]$$
$$= \mathbb{E} \left[\left(\frac{1}{M} \sum_{m=1}^M \epsilon_m(\mathbf{x}) \right)^2 \right]$$

Let's assume:

$$\mathbb{E} [\epsilon_m(\mathbf{x})] = 0 \quad \mathbb{E} [\epsilon_m(\mathbf{x})\epsilon_l(\mathbf{x})] = 0$$

Under this assumption: $E_{COM} = \frac{1}{M} E_{AV}$

Mixtures of experts

Bagging produces a result by averaging the predictions of the underlying classifiers.

Alternative: mixture of experts.

have each classifier be responsible for a small part of the input space

Boosting

Similar to bagging, but uses a more sophisticated method for constructing its diverse training sets.

Main ideas:

- ❖ Train the next classifier on examples that previous classifiers made errors on.
- ❖ Assign each classifier a confidence value that depends on its accuracy.

Boosting

Algorithm Boosting(D, T, \mathcal{A})

Input : data set D ; ensemble size T ; learning algorithm \mathcal{A} .

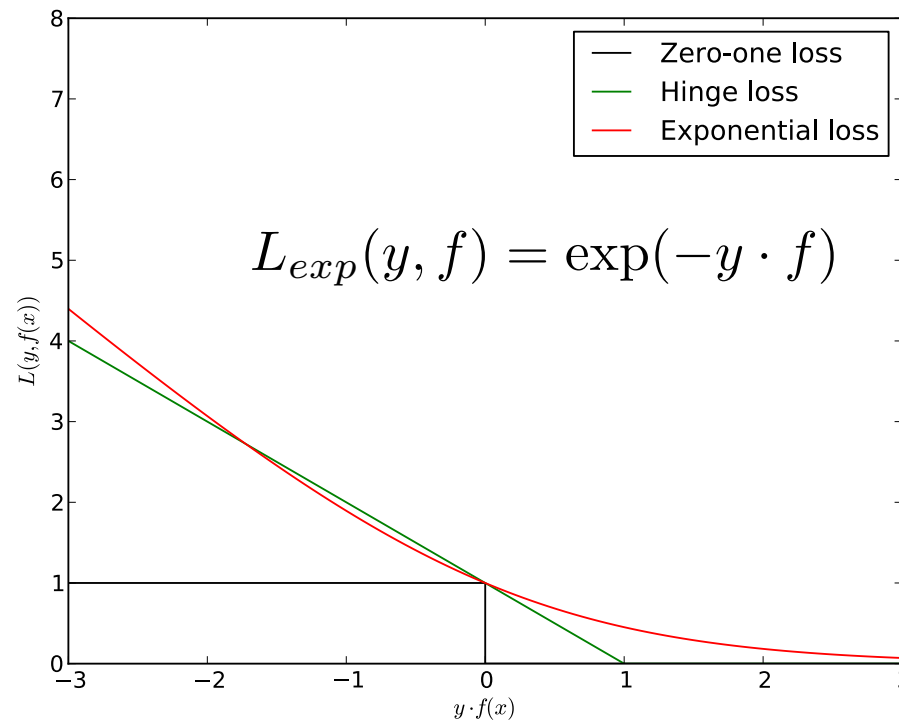
Output : weighted ensemble of models.

```
1  $w_{1i} \leftarrow 1/|D|$  for all  $x_i \in D$  ; // start with uniform weights
2 for  $t = 1$  to  $T$  do
3   run  $\mathcal{A}$  on  $D$  with weights  $w_{ti}$  to produce a model  $M_t$ ;
4   calculate weighted error  $\epsilon_t$ ;
5   if  $\epsilon_t \geq 1/2$  then
6     | set  $T \leftarrow t - 1$  and break
7   end
8    $\alpha_t \leftarrow \frac{1}{2} \ln \frac{1-\epsilon_t}{\epsilon_t}$  ; // confidence for this model
9    $w_{(t+1)i} \leftarrow \frac{w_{ti}}{2\epsilon_t}$  for misclassified instances  $x_i \in D$  ; // increase weight
10   $w_{(t+1)j} \leftarrow \frac{w_{tj}}{2(1-\epsilon_t)}$  for correctly classified instances  $x_j \in D$  ; // decrease
11 end
12 return  $M(x) = \sum_{t=1}^T \alpha_t M_t(x)$ 
```

Boosting and the exponential loss

What does boosting do?

Turns out that boosting minimizes the the exponential loss and leads to large margin solutions.



Robert E. Schapire, Yoav Freund, Peter Bartlett and Wee Sun Lee.

Boosting the margin: A new explanation for the effectiveness of voting methods.

The Annals of Statistics, 26(5):1651-1686, 1998.

Which classifier to boost?

Typically choose a simple classifier such as a “decision stump” or a simple linear classifier.

The bias-variance decomposition

The expected loss can be decomposed as follows:

$$\text{expected loss} = (\text{bias})^2 + \text{variance}$$

bias - the extent to which the average prediction differs from the label

variance - the variability of the classifier when trained on different training sets

Very flexible models have very low bias, but high variance.

Rigid models have high bias and low variance.

The best model will achieve a good balance between the two.

Bagging vs boosting

$$\text{expected loss} = (\text{bias})^2 + \text{variance}$$

bias - the extent to which the average prediction is different from the label

variance - the variability of the classifier when trained on different training sets

Very flexible models have very low bias, but high variance.

Rigid models have high bias and low variance.

Bagging is a variance reduction technique, whereas boosting acts to reduce bias.