Decision Trees and Forests

COMPSCI 371D — Machine Learning

э

< ロ > < 同 > < 回 > < 回 > < 回 > <

Outline

Motivation

- 2 Recursive Splits and Trees
- O Prediction
- 4 Purity
- 6 Splitting
- 6 Forests: Bagging and Randomization
- Forest Training and Inference
- Out-of-Bag Statistical Risk Estimate

・ 同 ト ・ ヨ ト ・ ヨ ト

Linear Predictors \rightarrow Trees \rightarrow Forests

- Linear predictors:
 - + Few parameters \rightarrow Good generalization, efficient training
 - + Convex risk \rightarrow Unique minimum risk, easy optimization
 - + Score-based \rightarrow Measure of confidence
 - Few parameters \rightarrow Limited expressiveness:
 - Regessor is an affine function
 - Classifier is a set of convex regions in X
- Decision trees:
 - Score based (in a sophisticated way)
 - Arbitrarily expressive: Flexible, but generalizes poorly
 - Interpretable: We can audit a decision
- Random decision forests:
 - · Ensembles of trees that vote on an answer
 - Expressive (somewhat less than trees), generalize well

不是我 不是我

Splitting X Recursively



э

< A >

A Decision Tree



э

What's in a Node

- Internal:
 - Split parameters: Dimension $j \in \{1, \ldots, d\}$, threshold $t \in \mathbb{R}$
 - Pointers to children, corresponding to subsets of S:

$$L \stackrel{\text{def}}{=} \{ (\mathbf{x}, y) \in S \mid x_j \le t \}$$
$$R \stackrel{\text{def}}{=} \{ (\mathbf{x}, y) \in S \mid x_j > t \}$$

• Leaf: Distribution of training values *y* in this subset of *X*:

p, discrete for classification, histogram for regression

- At inference time, return a *summary* of **p** as the value for the leaf
 - Mode (majority) for a classifier
 - Mean or median for a regressor (Remember k-NN?)

・ 同 ト ・ ヨ ト ・ ヨ ト

Why Store **p**?

- Can't we just store summary(p) at the leaves?
- With **p**, we can compute a confidence value
- (More important) We need **p** at every node during training to evaluate purity

・ 同 ト ・ ヨ ト ・ ヨ ト

Prediction

```
function y \leftarrow \text{predict}(\mathbf{x}, \tau, \text{summary})
if leaf?(\tau) then
return summary(\tau.p)
else
return predict(\mathbf{x}, \text{split}(\mathbf{x}, \tau), \text{summary})
end if
end function
```

```
function \tau \leftarrow \text{split}(\mathbf{x}, \tau)
if x_{\tau,j} \le \tau.t then
return \tau.L
else
return \tau.R
end if
end function
```

э

< A >

Design Decisions for Training

- How to define (im)purity
- How to find optimal split parameters j and t
- When to stop splitting

< 同 > < 三 > < 三 >

Impurity Measure 1: The Error Rate

- Simplest option: $i(S) = \overline{err}(S) = 1 \max_{y} p(y|S)$
- S: subset of T that reaches the given node
- Interpretation:
 - Put yourself at node τ
 - The distribution of training-set labels that are routed to τ is that of the labels in *S*
 - If the distribution is representative:
 - The best the classifier can do is to pick the label with the highest fraction, maxy p(y|S)
 - err(S) is the probability that the classifier is wrong at τ (empirical risk)

Impurity Measure 2: The Gini Index

- A classifier that always picks the most likely label does best at inference time
- However, it ignores all other labels at training time $\mathbf{p} = [0.5, 0.49, 0.01]$ same error rate as $\mathbf{q} = [0.5, 0.25, 0.25]$
- In **p**, we have almost eliminated the third label
- q closer to uniform, perhaps less desirable
- For evaluating splits (only), consider a *stochastic predictor*. $\hat{y} = h_{\text{Gini}}(\mathbf{x}) = y$ with probability p(y|S)
- The Gini index measures the empirical risk for the stochastic predictor (looks at all of **p**, not just p_{max})
- Says that **p** is a bit better than **q**: **p** is less impure than **q**
- $\mathrm{i}(\mathcal{S}_{p}) pprox 0.51$ and $\mathrm{i}(\mathcal{S}_{q}) pprox 0.62$

The Gini Index

• Stochastic predictor:

 $\hat{y} = h_{ ext{Gini}}(\mathbf{x}) = y$ with probability $p(y|\mathcal{S})$ for $y \in Y$

Purity

- What is the empirical risk for *h*_{Gini}?
- Answer y is chosen as \hat{y} with probability p(y|S)
- When the answer is *y*, it is wrong with probability
 ≈ 1 − p(y|S) (fraction of training samples that have true
 answer *y*)
- Therefore, impurity defined as the empirical risk of h_{Gini} is $i(S) = L_S(h_{\text{Gini}}) = \sum_{y \in Y} p(y|S)(1 - p(y|S)) = 1 - \sum_{y \in Y} p^2(y|S)$

3

How to Split

• Split at training time:

If training subset *S* made it to the current node, put all samples in *S* into either *L* or *R* by the split rule

- Split at inference time: Send **x** either to τ .*L* or to τ .*R*
- Either way:
 - Choose (training) or retrieve (inference) a dimension *j* in {1,..., *d*}
 - Choose (training) or retrieve (inference) a threshold t
 - Any data point for which $x_j \leq t$ goes to τL
 - All other points go to τ.R
- How to pick j and t at training time?

-

< ロ > < 同 > < 回 > < 回 > .

How to Pick *j* and *t* at Each Node?

- Try all possibilities and pick the best
- "Best:" Maximizes the decrease in impurity: $\Delta i(S, L, R) = i(S) - \frac{|L|}{|S|}i(L) - \frac{|R|}{|S|}i(R)$
- "All possibilities:" Choices are finite in number
 - Sorted unique values in x_i across T: $x_i^{(0)}, \ldots, x_i^{(u_i)}$
 - Possible thresholds: $t = t_j^{(1)}, \dots, t_j^{(u_j)}$ where $t_j^{(\ell)} = \frac{x_j^{(\ell-1)} + x_j^{(\ell)}}{2}$ for $\ell = 1, \dots, u_j$
- Nested loop: for $j=1,\ldots,d$ for $t=t_j^{(1)},\ldots,t_j^{(u_j)}$
- Efficiency hacks are possible

< ロ > < 同 > < 三 > < 三 > -

Stopping too Soon is Dangerous

Temptation: Stop when impurity does not decrease



< A >

I ≥ ►

When to Stop Splitting

Possible stopping criteria

- Impurity is zero
- Too few samples in either L or R
- Maximum depth reached
- Overgrow the tree, then prune it
- There is no optimal pruning method (Finding the optimal tree is NP-hard) (Reduction from set cover problem, Hyafil and Rivest)
- Better option: Random Decision Forests

< 🗇 > < 🖻 > < 🖻

Summary: Training a Decision Tree

- Use exhaustive search at the root of the tree to find the dimension *j* and threshold *t* that splits *T* with the biggest decrease in impurity
- Store *j* and *t* at the root of the tree
- Make new children with L and R
- · Repeat on the two subtrees until some criterion is met



Summary: Predicting with a Decision Tree

- Use *τ.j* and *τ.t* at the root *τ* to see if **x** belongs in *τ.L* or *τ.R*
- Go to the appropriate child
- Repeat until a leaf is reached
- Return summary(p)
- summary is majority for a classifier, mean or median for a regressor

From Trees to Forests

- Trees are flexible \rightarrow good expressiveness
- Trees are flexible \rightarrow poor generalization
- Pruning is an option, but messy and heuristic
- Random Decision Forests let several trees vote
- Use the bootstrap to give different trees different views of the data
- Randomize split rules to make trees even more independent

Random Forests

- *M* trees instead of one
- Train trees to completion (perfectly pure leaves) or to near completion (few samples per leaf)
- Give tree *m* training bag *B_m*
 - Draw |*T*| training samples independently at random with replacement out of *T*
 - $|B_m| = |T|$
 - About 63% of samples from *T* are in *B_m*
- Make trees more independent by randomizing split dim:

• Original trees: for
$$j = 1, ..., d$$

for $t = t_j^{(1)}, ..., t_j^{(u_j)}$
• Forest trees: $j = random \text{ out of } 1, ..., d$
for $t = t_j^{(1)}, ..., t_j^{(u_j)}$

< 同 > < 回 > < 回 > -

Randomizing Split Dimension

$$j = random \text{ out of } \mathbf{1}, \dots, \mathbf{a}$$

for $t = t_j^{(1)}, ..., t_j^{(u_j)}$

- Still search for the optimal threshold
- Give up optimality for independence
- Dimensions are revisited anyway in a tree
- Tree may get deeper, but still achieves zero training risk
- Independent splits and different data views lead to good generalization when voting
- Bonus: training a single tree is now *d* times faster

・ 同 ト ・ ヨ ト ・ ヨ ト …

Training

 $\begin{array}{ll} \mbox{function } \phi \leftarrow \mbox{trainForest}(T,M) & \triangleright \mbox{ M is the desired number of trees} \\ \phi \leftarrow \ensuremath{\emptyset} & \triangleright \mbox{ The initial forest has no trees} \\ \mbox{for $m = 1, \dots, M$ do} \\ S \leftarrow |T| \mbox{ samples unif. at random out of T with replacement} \\ \phi \leftarrow \phi \cup \{\mbox{trainTree}(S,0)\} & \triangleright \mbox{ Slightly modified trainTree} \\ \mbox{end for} \\ \mbox{end function} \end{array}$

< ロ > < 同 > < 回 > < 回 > .

Inference

function $y \leftarrow$ forestPredict(\mathbf{x}, ϕ , summary) $V = \{\}$ \triangleright A set of values, one per tree, initially empty for $\tau \in \phi$ do $y \leftarrow$ predict(\mathbf{x}, τ , summary) \triangleright The predict function for trees $V \leftarrow V \cup \{y\}$ end for return summary(V) end function

< ロ > < 同 > < 回 > < 回 > .

Out-of-Bag Statistical Risk Estimate

- Random forests have "built-in" training/validation or training/testing splits
- Tree *m*: B_m for training, $V_m = T \setminus B_m$ for testing
- h_{oob} is a predictor that works only for $(\mathbf{x}_n, y_n) \in T$:
 - Let tree *m* vote for *y* only if $\mathbf{x}_n \notin B_m$
 - $h_{\text{oob}}(\mathbf{x}_n)$ is the summary of the votes over participating trees
 - Summary: majority (classification); mean, median (regression)
- Out-of-bag risk estimate:
 - *T*' = {*t* ∈ *T* | ∃ *m* such that *t* ∉ *B_m*} (samples that were left out of *some* bag, so some trees can vote on them)
 - Statistical risk estimate: empirical risk of h_{oob} over T': $L_{T'}(h_{\text{oob}}) = \frac{1}{|T'|} \sum_{(\mathbf{x}, y) \in T'} \ell(y, h_{\text{oob}}(\mathbf{x}))$

-

$T' \approx T$

- *L_{T'}*(*h*_{oob}) can be shown to be an unbiased estimate of the statistical risk
- No separate test set needed if T' is large enough
- How big is *T*'?
- |T'| has a binomial distribution with *N* points, $p = 1 - (1 - 0.37)^M \approx 1$ as soon as M > 20

• Mean
$$\mu = pN$$
, variance $\sigma^2 = p(1 - p)N$

- $\sigma/\mu = \sqrt{\frac{1-p}{pN}} \rightarrow 0$ quite rapidly with growing *M* and *N*
- For reasonably large *N*, the size of *T'* is very predictably close to *N*: All samples in *T* are also in *T'* nearly always

-

Summary of Random Forests

- Random views of the training data by bagging
- Independent decisions by randomizing split dimensions
- Ensemble voting leads to good generalization
- Number M of trees tuned by cross-validation
- OOB estimate can replace final testing
- (In practice, that won't fly for papers)
- More efficient to train than a single tree if *M* < *d*
- Still rather efficient otherwise, and parallelizable
- Conceptually simple, easy to adapt to different problems
- Lots of freedom about split rule
- Example: Hybrid regression/classification problems

э

・ロッ ・ 一 ・ ・ ー ・ ・ ・ ・ ・ ・