

Structured Prediction & Decision Trees

Mengye Ren

(Slides credit to David Rosenberg, He He, et al.)

NYU

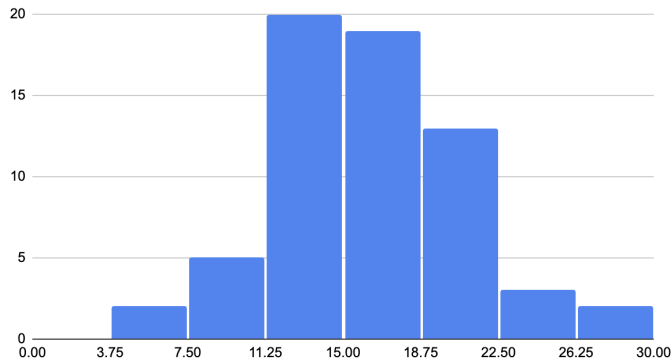
Nov 5, 2024

Final Project

- Project Consultation is mandatory. I am meeting with all 22 teams.
- Every member in the group needs to be present.
- If you have an emergency situation and cannot attend, you need to let me know.
- You may lose participation marks for not showing up in the consultation session without a reason.
- If you plan to show up, be on time. Your team member's time is also valuable.

Midterm

- Out of 30: Average: 16.05, Std: 4.29
- 9 free marks, capped at 30.
- Regrade request — email graders.



Scientific Writing

Introduction

- Introduction is the road map of the full report. Here is a general structure.
- Paragraph 1
 - What are some broad context of the problem? ✓
 - What is the problem that you are studying? ✓
 - Why is the problem interesting? ✓
- Paragraph 2
 - Historically, what have people been doing in the space of similar problems?
 - What is the gap and what extra can this project bring?
- Paragraph 3
 - What technical approach are you taking?
 - What dataset have you experimented with and what are the core results?
 - What are some broader impact of the work?

Introduction

Cybersecurity threats have become a major concern, posing risks to both personal data and organizational assets. In 2023, cyberattacks led to financial damages of over \$10 billion in the United States alone, according to recent reports. Globally, these numbers were even more concerning, with the cost of cybercrime expected to surpass \$8 trillion in 2023, as noted by cybersecurity researchers. In the U.S., approximately 66% of organizations experienced at least one form of cyberattack in 2023. These figures highlight an urgent need for robust strategies to detect, prevent, and mitigate cybersecurity breaches across industries.

What are the issues?

Revised Example

The rapid escalation in cybersecurity threats has led to significant financial and operational impacts, with damages in the United States alone exceeding \$10 billion in 2023, according to the Cybersecurity & Infrastructure Security Agency (CISA) [1]. Despite advancements in cybersecurity solutions, traditional rule-based systems often struggle to detect new or sophisticated attack patterns, especially as cybercriminals evolve their techniques to evade static detection rules. This limitation highlights a critical gap: the need for adaptable, data-driven methods that can dynamically learn and respond to emerging threats.

Continue...

Revised Example

In response to this gap, our project applies machine learning to develop a robust intrusion detection system that leverages a hybrid approach combining supervised and unsupervised learning. Using a labeled dataset of network traffic, our model is first trained to distinguish normal from suspicious activity. To enhance adaptability, we also employ an anomaly detection module using autoencoders, which identifies deviations from typical patterns, even for attack types not present in the training data.

Our results demonstrate the effectiveness of our approach. The model achieved an accuracy of 94% in detecting known attack types and reduced false positives by 30% compared to a traditional rule-based system. Furthermore, the anomaly detection module identified previously unseen attack patterns with an 87% true positive rate. These results suggest that our machine learning-based system can offer a more flexible and accurate defense mechanism against evolving cybersecurity threats.

Abstract

Abstract is the summary of the introduction. Make it into one paragraph.

Cybersecurity threats are increasing in frequency and sophistication, resulting in substantial financial losses and operational disruptions. Traditional rule-based detection systems are limited in their ability to identify novel attack patterns, creating a need for adaptable, data-driven solutions. This project presents a machine learning-based intrusion detection system designed to address this gap by combining supervised learning for known threats with an unsupervised anomaly detection module to identify emerging attack types. Experimental results demonstrate the model's effectiveness, achieving 94% accuracy in detecting known attacks and an 87% true positive rate for novel threats, while reducing false positives by 30% compared to rule-based methods. These findings suggest that our hybrid approach provides a more flexible and accurate defense against evolving cyber threats.

- Survey historical attempts of similar problems (not necessarily the same problem!)
- What categories do historical approaches span across? Which category does your method fall into?
- What is the relation of your work in the context of prior literature?

Example

Intrusion detection has been approached from several angles, with methods broadly categorized into rule-based systems, supervised learning-based intrusion detection, unsupervised anomaly detection, and hybrid models.

Rule-Based Detection: Traditional rule-based systems, such as Snort, use known attack signatures to detect intrusions [1]. While effective for known threats, these systems fail to recognize novel or evolving attacks [2]. Our approach seeks to address this limitation by incorporating machine learning for greater adaptability.

Supervised Machine Learning: Supervised models like support vector machines and neural networks are commonly used due to their high accuracy with labeled data, effectively identifying known attacks [3, 4]. However, their reliance on labeled datasets makes them less adaptable to new threats. Our project builds on supervised methods but adds an anomaly detection layer to handle unknown attacks.

Continue...

Example

Unsupervised and Anomaly-Based Detection: Unsupervised methods, including autoencoders and clustering, detect unusual patterns in network traffic without labeled data [6, 7]. While flexible, these techniques often produce higher false-positive rates, posing practical challenges [8]. Our model incorporates an autoencoder for anomaly detection, with optimized thresholds to reduce false positives.

Hybrid Approaches: Hybrid models combine supervised and unsupervised techniques to balance accuracy and adaptability. For example, hybrid models like those by Huang et al. (2021) integrate anomaly detection within supervised frameworks, achieving lower false-positive rates [9]. Our work further advances these methods by embedding an autoencoder anomaly module into a supervised model with optimized thresholds, aiming for reliable real-time intrusion detection.

Figures and Tables

- Try to export as PDF instead of JPG/PNG (rasterized).
- Use bigger font size in the figure (same or slightly smaller than the main text).
- When there is a strong trend, use a figure.
- When you need to emphasize a small difference, use a table.
- When the table is too wide, use `resizebox` to fit your table with the text width.
- Keep the same number of significant digits.

Introduction to Structured Prediction

Example: Part-of-speech (POS) Tagging

- Given a sentence, give a part of speech tag for each word:

x	$\underbrace{[\text{START}]}_{x_0}$	$\underbrace{\text{He}}_{x_1}$	$\underbrace{\text{eats}}_{x_2}$	$\underbrace{\text{apples}}_{x_3}$
y	$\underbrace{[\text{START}]}_{y_0}$	$\underbrace{\text{Pronoun}}_{y_1}$	$\underbrace{\text{Verb}}_{y_2}$	$\underbrace{\text{Noun}}_{y_3}$

- $\mathcal{V} = \{\text{all English words}\} \cup \{[\text{START}], ", ."]\}$
- $\mathcal{X} = \mathcal{V}^n, n = 1, 2, 3, \dots$ [Word sequences of any length]
- $\mathcal{P} = \{\text{START, Pronoun, Verb, Noun, Adjective}\}$
- $\mathcal{Y} = \mathcal{P}^n, n = 1, 2, 3, \dots$ [Part of speech sequence of any length]

Example: Action grounding from long-form videos

- Given a long video, segment the video into short windows where each window corresponds to an action from a list of actions.
- E.g. slicing, chopping, frying, washing, etc.
- $\mathcal{V} = \mathbb{R}^D$ image features
- $\mathcal{X} = \mathcal{V}^n$, $n = 1, 2, 3, \dots$ [video frame length]
- $\mathcal{P} = \{\text{Slicing, Chopping, Frying, } \dots\}$
- $\mathcal{Y} = \mathcal{P}^n$, $n = 1, 2, 3, \dots$ [Part of speech sequence of any length]
- Can also be represented with start and end timestamps.

Multiclass Hypothesis Space

- **Discrete** output space: $\mathcal{Y}(x)$
 - Very large but has structure, e.g., linear chain (sequence labeling), tree (parsing)
 - Size depends on input x
- Base Hypothesis Space: $\mathcal{H} = \{h : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}\}$
 - $h(x, y)$ gives **compatibility score** between input x and output y
- Multiclass hypothesis space

$$\mathcal{F} = \left\{ x \mapsto \arg \max_{y \in \mathcal{Y}} h(x, y) \mid h \in \mathcal{H} \right\}$$

- Final prediction function is an $f \in \mathcal{F}$.
- For each $f \in \mathcal{F}$ there is an underlying compatibility score function $h \in \mathcal{H}$.

Structured Prediction

- Part-of-speech tagging

x:	he	eats	apples
y:	pronoun	verb	noun

- Multiclass hypothesis space:

$$h(x, y) = w^T \Psi(x, y) \quad (1)$$

$$\mathcal{F} = \left\{ x \mapsto \arg \max_{y \in \mathcal{Y}} h(x, y) \mid h \in \mathcal{H} \right\} \quad (2)$$

- A special case of multiclass classification
- How to design the feature map Ψ ? What are the considerations?

Unary features

- A **unary feature** only depends on
 - the label at a **single position**, y_i , and x
- Example:

$$\phi_1(x, y_i) = \mathbb{1}[x_i = \text{runs}] \mathbb{1}[y_i = \text{Verb}]$$

$$\phi_2(x, y_i) = \mathbb{1}[x_i = \text{runs}] \mathbb{1}[y_i = \text{Noun}]$$

$$\phi_3(x, y_i) = \mathbb{1}[x_{i-1} = \text{He}] \mathbb{1}[x_i = \text{runs}] \mathbb{1}[y_i = \text{Verb}]$$

Markov features

- A **markov feature** only depends on
 - two **adjacent** labels, y_{i-1} and y_i , and x
- Example:

$$\theta_1(x, y_{i-1}, y_i) = \mathbb{1}[y_{i-1} = \text{Pronoun}] \mathbb{1}[y_i = \text{Verb}]$$

$$\theta_2(x, y_{i-1}, y_i) = \mathbb{1}[y_{i-1} = \text{Pronoun}] \mathbb{1}[y_i = \text{Noun}]$$

- Reminiscent of Markov models in the output space
- Possible to have higher-order features

Local Feature Vector and Compatibility Score

- At each position i in sequence, define the **local feature vector** (unary and markov):

$$\Psi_i(x, y_{i-1}, y_i) = (\phi_1(x, y_i), \phi_2(x, y_i), \dots, \theta_1(x, y_{i-1}, y_i), \theta_2(x, y_{i-1}, y_i), \dots)$$

- And **local compatibility score** at position i : $\langle w, \Psi_i(x, y_{i-1}, y_i) \rangle$.
- The compatibility score for (x, y) is the sum of local compatibility scores:

$$\sum_i \langle w, \Psi_i(x, y_{i-1}, y_i) \rangle = \left\langle w, \sum_i \Psi_i(x, y_{i-1}, y_i) \right\rangle = \langle w, \Psi(x, y) \rangle, \quad (3)$$

where we define the **sequence feature vector** by

$$\Psi(x, y) = \sum_i \Psi_i(x, y_{i-1}, y_i). \quad \text{decomposable}$$

Structured perceptron

Given a dataset $\mathcal{D} = \{(x, y)\}$;

Initialize $w \leftarrow 0$;

for $iter = 1, 2, \dots, T$ **do**

for $(x, y) \in \mathcal{D}$ **do**

$\hat{y} = \arg \max_{y' \in \mathcal{Y}(x)} w^T \psi(x, y')$;

if $\hat{y} \neq y$ **then** // We've made a mistake

$w \leftarrow w + \Psi(x, y)$; // Move the scorer towards $\psi(x, y)$

$w \leftarrow w - \Psi(x, \hat{y})$; // Move the scorer away from $\psi(x, \hat{y})$

end

end

end

Identical to the multiclass perceptron algorithm except the $\arg \max$ is now over the structured output space $\mathcal{Y}(x)$.

Structured hinge loss

- Recall the generalized hinge loss

$$\ell_{\text{hinge}}(y, \hat{y}) \stackrel{\text{def}}{=} \max_{y' \in \mathcal{Y}(\mathbf{x})} (\Delta(y, y') + \langle w, (\Psi(\mathbf{x}, y') - \Psi(\mathbf{x}, y)) \rangle) \quad (4)$$

- What is $\Delta(y, y')$ for two sequences?
- Hamming loss** is common:

$$\Delta(y, y') = \frac{1}{L} \sum_{i=1}^L \mathbb{1}[y_i \neq y'_i]$$

where L is the sequence length.

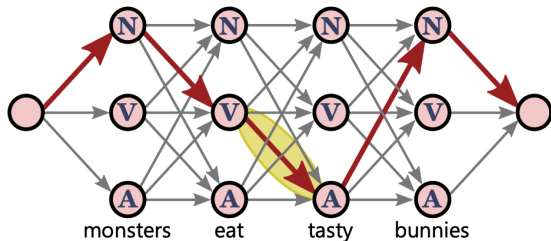
Exercise:

- Write down the objective of structured SVM using the structured hinge loss.
- Stochastic sub-gradient descent for structured SVM (similar to HW3 P3)
- Compare with the structured perceptron algorithm

The argmax problem for sequences

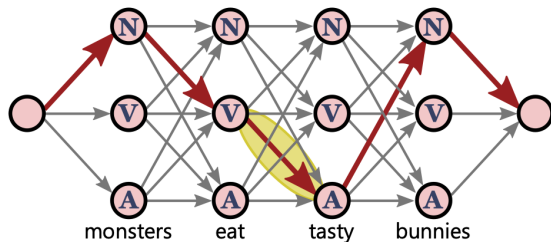
Problem To compute predictions, we need to find $\arg\max_{y \in \mathcal{Y}(x)} \langle w, \Psi(x, y) \rangle$, and $|\mathcal{Y}(x)|$ is exponentially large.

Observation $\Psi(x, y)$ decomposes to $\sum_i \Psi_i(x, y)$.



Solution Dynamic programming
the running time?

Dynamic Programming (MAP inference)



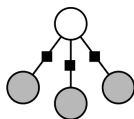
- Initiate $\alpha_j(1) = \exp(w^\top \psi(y_1 = j, x_1))$
- Recursion: $\alpha_j(t) = \max_i \alpha_i(t-1) + w^\top \psi(y_t = j, y_{t-1} = i, x_t)$
- $p_j(t) = \arg \max \dots$ (why?)
- Problem-specific: DP, graph cuts, etc. General: integer linear programming (ILP).

Conditional random field (CRF)

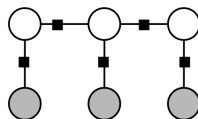
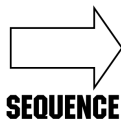
- Recall that we can write logistic regression in a general form:

$$p(y|x) = \frac{1}{Z(x)} \exp(w^\top \psi(x, y)).$$

- Z is normalization constant: $Z(x) = \sum_{y \in Y} \exp(w^\top \psi(x, y))$.
- Example: linear chain $\{y_t\}$
- We can incorporate unary and Markov features: $p(y|x) = \frac{1}{Z(x)} \exp(\sum_t w^\top \psi(x, y_t, y_{t-1}))$



Logistic Regression



Linear-chain CRFs

Conditional random field (CRF)

- Compared to Structured SVM, CRF has a probabilistic interpretation.
- We can draw samples in the output space.
- How do we learn w ? Maximum log likelihood, and regularization term: $\lambda \|w\|^2$
- Loss function:

$$\begin{aligned} l(w) &= -\frac{1}{N} \sum_{i=1}^N \log p(y^{(i)} | x^{(i)}) + \frac{1}{2} \lambda \|w\|^2 \\ &= -\frac{1}{N} \sum_i \sum_t \sum_k w_k \psi_k(y_t^{(i)}, y_{t-1}^{(i)}) + \frac{1}{N} \sum_i \log Z(x^{(i)}) + \frac{1}{2} \sum_k \lambda w_k^2 \end{aligned}$$

Conditional random field (CRF)

- Loss function:

$$l(w) = -\frac{1}{N} \sum_i \sum_t \sum_k w_k \psi_k(x^{(i)}, y_t^{(i)}, y_{t-1}^{(i)}) + \frac{1}{N} \sum_i \log Z(x^{(i)}) + \frac{1}{2} \sum_k \lambda w_k^2$$

- Gradient:

$$\frac{\partial l(w)}{\partial w_k} = -\frac{1}{N} \sum_i \sum_t \sum_k \psi_k(x^{(i)}, y_t^{(i)}, y_{t-1}^{(i)}) \quad (5)$$

$$+ \frac{1}{N} \sum_i \frac{\partial}{\partial w_k} \log \sum_{y' \in Y} \exp\left(\sum_t \sum_{k'} w_{k'} \psi_{k'}(x^{(i)}, y'_t, y'_{t-1})\right) + \sum_k \lambda w_k \quad (6)$$

Conditional random field (CRF)

- What is $\frac{1}{N} \sum_i \sum_t \sum_k \psi_k(x^{(i)}, y_t^{(i)}, y_{t-1}^{(i)})$?
- It is the expectation $\psi_k(x^{(i)}, y_t, y_{t-1})$ under the empirical distribution $\tilde{p}(x, y) = \frac{1}{N} \sum_i \mathbb{1}[x = x^{(i)}] \mathbb{1}[y = y^{(i)}]$.

Conditional random field (CRF)

- What is $\frac{1}{N} \sum_i \frac{\partial}{\partial w_k} \log \sum_{y' \in Y} \exp(\sum_t \sum_{k'} w_{k'} \psi_{k'}(x^{(i)}, y'_t, y'_{t-1}))$?

$$\frac{1}{N} \sum_i \frac{\partial}{\partial w_k} \log \sum_{y' \in Y} \exp(\sum_t \sum_{k'} w_{k'} \psi_{k'}(x^{(i)}, y'_t, y'_{t-1})) \quad (7)$$

$$= \frac{1}{N} \sum_i \left[\sum_{y' \in Y} \exp(\sum_t \sum_{k'} w_{k'} \psi_{k'}(x^{(i)}, y'_t, y'_{t-1})) \right]^{-1} \quad (8)$$

$$\left[\sum_{y' \in Y} \exp(\sum_t \sum_{k'} w_{k'} \psi_{k'}(x^{(i)}, y'_t, y'_{t-1})) \sum_t \psi_k(x^{(i)}, y'_t, y'_{t-1}) \right] \quad (9)$$

$$= \frac{1}{N} \sum_i \sum_t \sum_{y' \in Y} p(y'_t, y'_{t-1} | x) \psi_k(x^{(i)}, y'_t, y'_{t-1}) \quad (10)$$

- It is the expectation of $\psi_k(x^{(i)}, y'_t, y'_{t-1})$ under the model distribution $p(y'_t, y'_{t-1} | x)$.

Conditional random field (CRF)

- To compute the gradient, we need to infer expectation under the model distribution $p(y|x)$.
- Compare the learning algorithms: in structured SVM we need to compute the argmax, whereas in CRF we need to compute the model expectation.
- Both problems are NP-hard for general graphs.

CRF Inference

- In the linear chain structure, we can use the forward-backward algorithm for inference, similar to Viterbi.
- Initiate $\alpha_j(1) = \exp(w^\top \psi(y_1 = j, x_1))$
- Recursion: $\alpha_j(t) = \sum_i \alpha_i(t-1) \exp(w^\top \psi(y_t = j, y_{t-1} = i, x_t))$
- Result: $Z(x) = \sum_j \alpha_j(T)$
- Similar for the backward direction.
- Test time, again use Viterbi algorithm to infer argmax.
- The inference algorithm can be generalized to belief propagation (BP) in a tree structure (exact inference).
- In general graphs, we rely on approximate inference (e.g. loopy belief propagation).

- POS tag Relationship between constituents, e.g. NP is likely to be followed by a VP.
- Semantic segmentation
Relationship between pixels, e.g. a grass pixel is likely to be next to another grass pixel, and a sky pixel is likely to be above a grass pixel.
- Multi-label learning
An image may contain multiple class labels, e.g. a bus is likely to co-occur with a car.

Structured Prediction

- Extension to multi-class prediction
- Structured SVM, CRF
- Output space containing structure
- Text and image applications

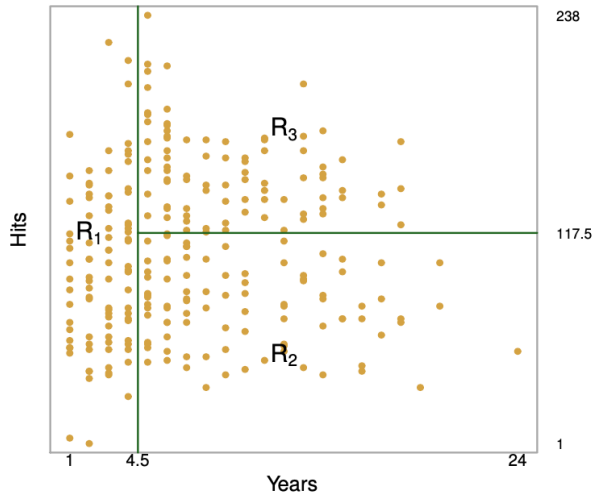
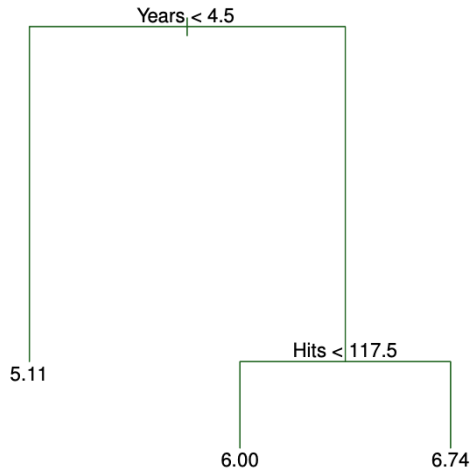
Decision Trees

Overview: Decision Trees

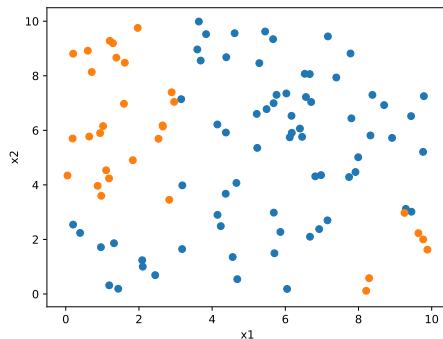
- Our first inherently non-linear classifier: decision trees.
- Ensemble methods: bagging and boosting.

Decision Trees

Regression trees: Predicting basketball players' salaries



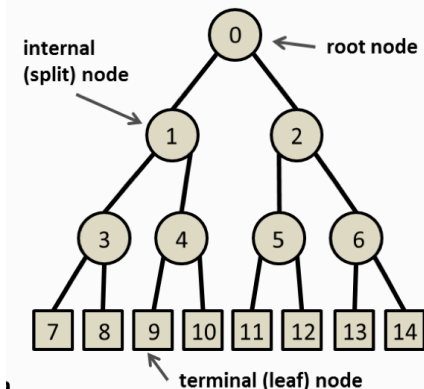
Classification trees



- Can we classify these points using a linear classifier?
- Partition the data into axis-aligned regions **recursively** (on the board)

Decision trees setup

A general tree structure



- We focus on *binary* trees (as opposed to multiway trees where nodes can have more than two children)
- Each node contains a subset of data points
- The data splits created by each node involve only a *single* feature
- For continuous variables, the splits are always of the form $x_i \leq t$
- For discrete variables, we partition values into two sets (not covered today)
- Predictions are made in terminal nodes

Constructing the tree

Goal Find boxes R_1, \dots, R_J that minimize $\sum_{j=1}^J \sum_{i \in R_j} (y_i - \hat{y}_{R_j})^2$, subject to complexity constraints.

Problem Finding the optimal binary tree is computationally intractable.

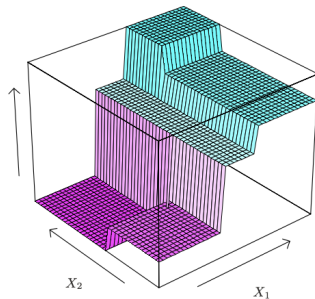
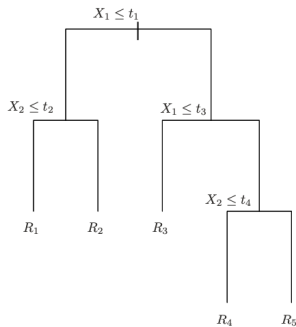
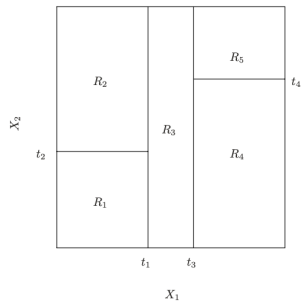
Solution Greedy algorithm: starting from the root, and repeating until a stopping criterion is reached (e.g., max depth), find the non-terminal node that results in the “best” split

- We only split regions defined by previous non-terminal nodes

Prediction Our prediction is the mean value of a terminal node: $\hat{y}_{R_m} = \text{mean}(y_i \mid x_i \in R_m)$

- A greedy algorithm is the one that make the best **local** decisions, without lookahead to evaluate their downstream consequences
- This procedure is not very likely to result in the globally optimal tree

Prediction in a Regression Tree



Finding the Best Split Point

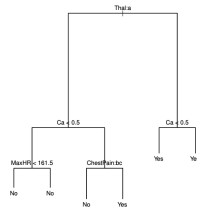
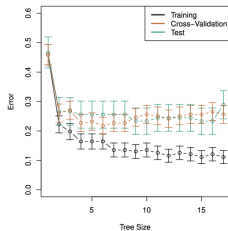
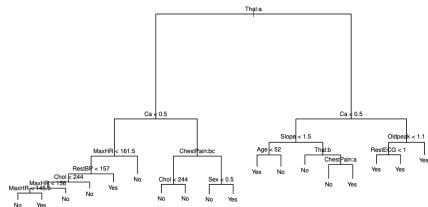
- We enumerate all features and all possible split points for each feature. There are infinitely many split points, but...
- Suppose we are now considering splitting on the j -th feature x_j , and let $x_{j(1)}, \dots, x_{j(n)}$ be the sorted values of the j -th feature.
- We only need to consider split points between two adjacent values, and any split point in the interval $(x_{j(r)}, x_{j(r+1)})$ will result in the same loss
- It is common to split half way between two adjacent values:

$$s_j \in \left\{ \frac{1}{2} (x_{j(r)} + x_{j(r+1)}) \mid r = 1, \dots, n-1 \right\}. \quad n-1 \text{ splits} \quad (11)$$

Decision Trees and Overfitting

- What will happen if we keep splitting the data into more and more regions?
 - Every data point will be in its own region—**overfitting**.
- When should we stop splitting? (Controlling the complexity of the hypothesis space)
 - Limit total number of nodes.
 - Limit number of terminal nodes.
 - Limit tree depth.
 - Require minimum number of data points in a terminal node.
 - **Backward pruning** (the approach used in **CART**; Breiman et al 1984):
 - 1 Build a really big tree (e.g. until all regions have ≤ 5 points).
 - 2 *Prune* the tree back greedily, potentially all the way to the root, until validation performance starts decreasing.

Pruning: Example



What Makes a Good Split for Classification?

Our plan is to predict the **majority label** in each region.

Which of the following splits is better?

Split 1 $R_1 : 8+ / 2-$ $R_2 : 2+ / 8-$

Split 2 $R_1 : 6+ / 4-$ $R_2 : 4+ / 6-$

How about here?

Split 1 $R_1 : 8+ / 2-$ $R_2 : 2+ / 8-$

Split 2 $R_1 : 6+ / 4-$ $R_2 : 0+ / 10-$

Intuition: we want to produce *pure* nodes, i.e. nodes where most instances have the same class.

Misclassification error in a node

- Let's consider the multiclass classification case: $\mathcal{Y} = \{1, 2, \dots, K\}$.
- Let node m represent region R_m , with N_m observations
- We denote the proportion of observations in R_m with class k by

$$\hat{p}_{mk} = \frac{1}{N_m} \sum_{\{i: x_i \in R_m\}} \mathbb{1}[y_i = k].$$

- We predict the majority class in node m :

$$k(m) = \arg \max_k \hat{p}_{mk}$$

Node Impurity Measures

- Three measures of **node impurity** for leaf node m :

- Misclassification error

$$1 - \hat{p}_{mk(m)}.$$

- The Gini index encourages \hat{p}_{mk} to be close to 0 or 1

$$\sum_{k=1}^K \hat{p}_{mk}(1 - \hat{p}_{mk}).$$

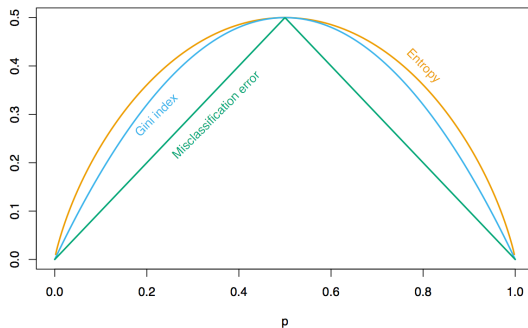
- Entropy / Information gain

$$-\sum_{k=1}^K \hat{p}_{mk} \log \hat{p}_{mk}.$$

- The Gini index and entropy are numerically similar to each other, and both work better in practice than the misclassification error.

Impurity Measures for Binary Classification

(p is the relative frequency of class 1)



Quantifying the Impurity of a Split

Scoring a potential split that produces the nodes R_L and R_R :

- Suppose we have N_L points in R_L and N_R points in R_R .
- Let $Q(R_L)$ and $Q(R_R)$ be the node impurity measures for each node.
- We aim to find a split that minimizes the *weighted average of node impurities*:

$$\frac{N_L Q(R_L) + N_R Q(R_R)}{N_L + N_R}$$

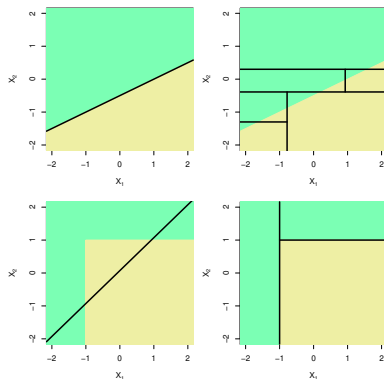
Discussion: Interpretability of Decision Trees



- Trees are easier to visualize and explain than other classifiers (even linear regression)
- Small trees are interpretable – large trees, maybe not so much

Discussion: Trees vs. Linear Models

Trees may have to work hard to capture linear decision boundaries, but can easily capture certain nonlinear ones:



Discussion: Review

Decision trees are:

- Non-linear: the decision boundary that results from splitting may end up being quite complicated
- Non-metric: they do not rely on the geometry of the space (inner products or distances)
- Non-parametric: they make no assumptions about the distribution of the data

Additional pros:

- Interpretable and simple to understand

Cons:

- Struggle to capture linear decision boundaries
- They have high variance and tend to **overfit**: they are sensitive to small changes in the training data (The ensemble techniques we discuss next can mitigate these issues)

Bagging and Random Forests

Recap: Statistics and Point Estimators

- We observe data $\mathcal{D} = (x_1, x_2, \dots, x_n)$ sampled i.i.d. from a parametric distribution $p(\cdot | \theta)$
- A **statistic** $s = s(\mathcal{D})$ is any function of the data:
 - E.g., sample mean, sample variance, histogram, empirical data distribution
- A statistic $\hat{\theta} = \hat{\theta}(\mathcal{D})$ is a **point estimator** of θ if $\hat{\theta} \approx \theta$

Recap: Bias and Variance of an Estimator

- Statistics are random, so they have probability distributions.
- The distribution of a statistic is called a **sampling distribution**.
- The standard deviation of the sampling distribution is called the **standard error**.
- Some parameters of the sampling distribution we might be interested in:

Bias $\text{Bias}(\hat{\theta}) \stackrel{\text{def}}{=} \mathbb{E}[\hat{\theta}] - \theta.$

Variance $\text{Var}(\hat{\theta}) \stackrel{\text{def}}{=} \mathbb{E}[\hat{\theta}^2] - \mathbb{E}^2[\hat{\theta}].$

- Why does variance matter if an estimator is unbiased?
 - $\hat{\theta}(\mathcal{D}) = x_1$ is an unbiased estimator of the mean of a Gaussian, but would be farther away from θ than the sample mean.

Variance of a Mean

- Let $\hat{\theta}(\mathcal{D})$ be an unbiased estimator with variance σ^2 : $\mathbb{E}[\hat{\theta}] = \theta$, $\text{Var}(\hat{\theta}) = \sigma^2$.
- So far we have used a single statistic $\hat{\theta} = \hat{\theta}(\mathcal{D})$ to estimate θ .
- Its standard error is $\sqrt{\text{Var}(\hat{\theta})} = \sigma$
- Consider a new estimator that takes the average of i.i.d. $\hat{\theta}_1, \dots, \hat{\theta}_n$ where $\hat{\theta}_i = \hat{\theta}(\mathcal{D}^i)$.
- The average has the same expected value but smaller standard error (recall that $\text{Var}(cX) = c^2 \text{Var}(X)$, and that the $\hat{\theta}_i$ -s are uncorrelated):

$$\mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n \hat{\theta}_i\right] = \theta \quad \text{Var}\left[\frac{1}{n} \sum_{i=1}^n \hat{\theta}_i\right] = \frac{\sigma^2}{n} \quad (12)$$

Averaging Independent Prediction Functions

- Suppose we have B independent training sets, all drawn from the same distribution ($\mathcal{D} \sim p(\cdot \mid \theta)$).
- Our learning algorithm gives us B prediction functions: $\hat{f}_1(x), \hat{f}_2(x), \dots, \hat{f}_B(x)$
- We will define the average prediction function as:

$$\hat{f}_{\text{avg}} \stackrel{\text{def}}{=} \frac{1}{B} \sum_{b=1}^B \hat{f}_b \quad (13)$$

Averaging Reduces Variance of Predictions

- The average prediction for x_0 is

$$\hat{f}_{\text{avg}}(x_0) = \frac{1}{B} \sum_{b=1}^B \hat{f}_b(x_0).$$

- $\hat{f}_{\text{avg}}(x_0)$ and $\hat{f}_b(x_0)$ have the same expected value, but
- $\hat{f}_{\text{avg}}(x_0)$ has smaller variance:

$$\text{Var}(\hat{f}_{\text{avg}}(x_0)) = \frac{1}{B} \text{Var}(\hat{f}_1(x_0))$$

- **Problem:** in practice we don't have B independent training sets!

The Bootstrap Sample

How do we simulate multiple samples when we only have one?

- A **bootstrap sample** from $\mathcal{D}_n = (x_1, \dots, x_n)$ is a sample of size n drawn *with replacement* from \mathcal{D}_n
- Some elements of \mathcal{D}_n will show up multiple times, and some won't show up at all
- Each x_i has a probability of $(1 - 1/n)^n$ of not being included in a given bootstrap sample
- For large n ,

$$\left(1 - \frac{1}{n}\right)^n \approx \frac{1}{e} \approx .368. \quad (14)$$

- So we expect ~63.2% of elements of \mathcal{D}_n will show up at least once.

The Bootstrap Method

Definition

A **bootstrap method** simulates B independent samples from P by taking B bootstrap samples from the sample \mathcal{D}_n .

- Given original data \mathcal{D}_n , compute B bootstrap samples D_n^1, \dots, D_n^B .
- For each bootstrap sample, compute some function

$$\phi(D_n^1), \dots, \phi(D_n^B)$$

- Use these values as though D_n^1, \dots, D_n^B were i.i.d. samples from P .
- This often ends up being very close to what we'd get with independent samples from P !

Independent Samples vs. Bootstrap Samples

- Point estimator $\hat{\alpha} = \hat{\alpha}(\mathcal{D}_{100})$ for samples of size 100, for a synthetic case where the data generating distribution is known
- Histograms of $\hat{\alpha}$ based on
 - 1000 independent samples of size 100 (left), vs.
 - 1000 bootstrap samples of size 100 (right)

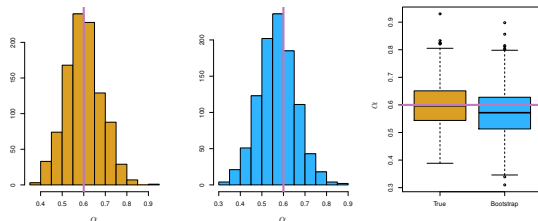


Figure 5.10 from *ISLR* (Springer, 2013) with permission from the authors: G. James, D. Witten, T. Hastie and R. Tibshirani.

Key ideas:

- In general, **ensemble methods** combine multiple weak models into a single, more powerful model
- Averaging i.i.d. estimates reduces variance without changing bias
- We can use bootstrap to simulate multiple data samples and average them
- Parallel ensemble (e.g., bagging): models are built independently
- Sequential ensemble (e.g., boosting): models are built sequentially
 - We try to find new learners that do well where previous learners fall short

Bagging: Bootstrap Aggregation

- We draw B bootstrap samples D^1, \dots, D^B from original data \mathcal{D}
- Let $\hat{f}_1, \hat{f}_2, \dots, \hat{f}_B$ be the prediction functions resulting from training on D^1, \dots, D^B , respectively
- The **bagged prediction function** is a *combination* of these:

$$\hat{f}_{\text{avg}}(x) = \text{Combine} \left(\hat{f}_1(x), \hat{f}_2(x), \dots, \hat{f}_B(x) \right)$$

Bagging: Bootstrap Aggregation

- Bagging is a general method for variance reduction, but it is particularly useful for decision trees
- For classification, averaging doesn't make sense; we can take a **majority vote** instead
- Increasing the number of trees we use in bagging does not lead to overfitting
- Is there a downside, compared to having a single decision tree?
- Yes: if we have many trees, the bagged predictor is much less interpretable

Aside: Out-of-Bag Error Estimation

- Recall that each bagged predictor was trained on about 63% of the data.
- The remaining 37% are called **out-of-bag (OOB)** observations.
- For i th training point, let

$$S_i = \{b \mid D^b \text{ does not contain } i\text{th point}\}$$

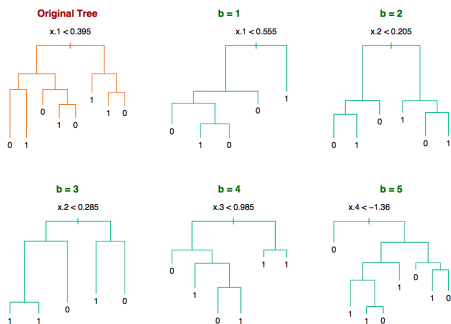
- The **OOB prediction** on x_i is

$$\hat{f}_{\text{OOB}}(x_i) = \frac{1}{|S_i|} \sum_{b \in S_i} \hat{f}_b(x_i)$$

- The OOB error is a good estimate of the test error
- Similar to cross validation error: both are computed on the training set

Applying Bagging to Classification Trees

- Input space $\mathcal{X} = \mathbb{R}^5$ and output space $\mathcal{Y} = \{-1, 1\}$. Sample size $n = 30$.



- Each bootstrap tree is quite different: different splitting variable at the root!
- High variance:** small perturbations of the training data lead to a high degree of model variability
- Bagging helps most when the base learners are relatively unbiased but have high variance (exactly the case for decision trees)

Motivating Random Forests: Correlated Prediction Functions

Recall the motivating principle of bagging:

- For $\hat{\theta}_1, \dots, \hat{\theta}_n$ *i.i.d.* with $\mathbb{E}[\hat{\theta}] = \theta$ and $\text{Var}[\hat{\theta}] = \sigma^2$,

$$\mathbb{E}\left[\frac{1}{n}\sum_{i=1}^n \hat{\theta}_i\right] = \mu \qquad \text{Var}\left[\frac{1}{n}\sum_{i=1}^n \hat{\theta}_i\right] = \frac{\sigma^2}{n}.$$

- What if $\hat{\theta}$'s are correlated?
- For large n , the covariance term dominates, limiting the benefits of averaging
- Bootstrap samples are
 - independent samples from the training set, but
 - **not** independent samples from $P_{\mathcal{X} \times \mathcal{Y}}$
- Can we reduce the dependence between \hat{f}_i 's?

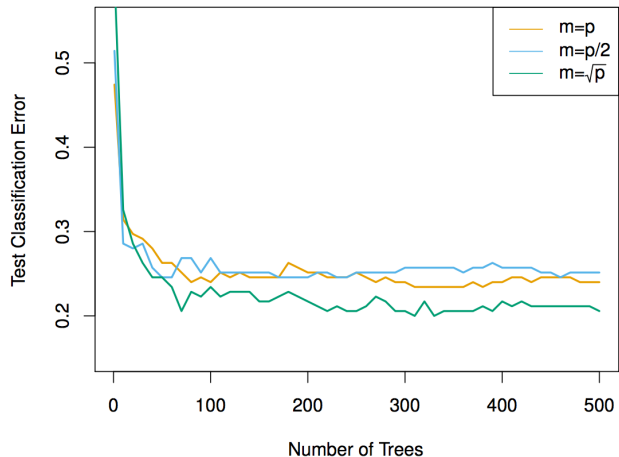
Random Forests

Key idea

Use bagged decision trees, but modify the tree-growing procedure to reduce the dependence between trees.

- Build a collection of trees independently (in parallel), as before
- When constructing each tree node, restrict choice of splitting variable to a randomly chosen subset of features of size m
 - This prevents a situation where all trees are dominated by the same small number of strong features (and are therefore too similar to each other)
- We typically choose $m \approx \sqrt{p}$, where p is the number of features (or we can choose m using cross validation)
- If $m = p$, this is just bagging

Random Forests: Effect of m



From *An Introduction to Statistical Learning, with applications in R* (Springer, 2013) with permission from the authors: G. James, D. Witten, T. Hastie and R. Tibshirani.

- The usual approach is to build very deep trees—low bias but **high variance**
- Ensembling many models reduces variance
 - Motivation: Mean of i.i.d. estimates has smaller variance than single estimate
- Use bootstrap to simulate many data samples from one dataset
 - \implies Bagged decision trees
- But bootstrap samples (and the induced models) are correlated
- Ensembling works better when we combine a diverse set of prediction functions
 - \implies Random forests: select a random subset of features for each decision tree