

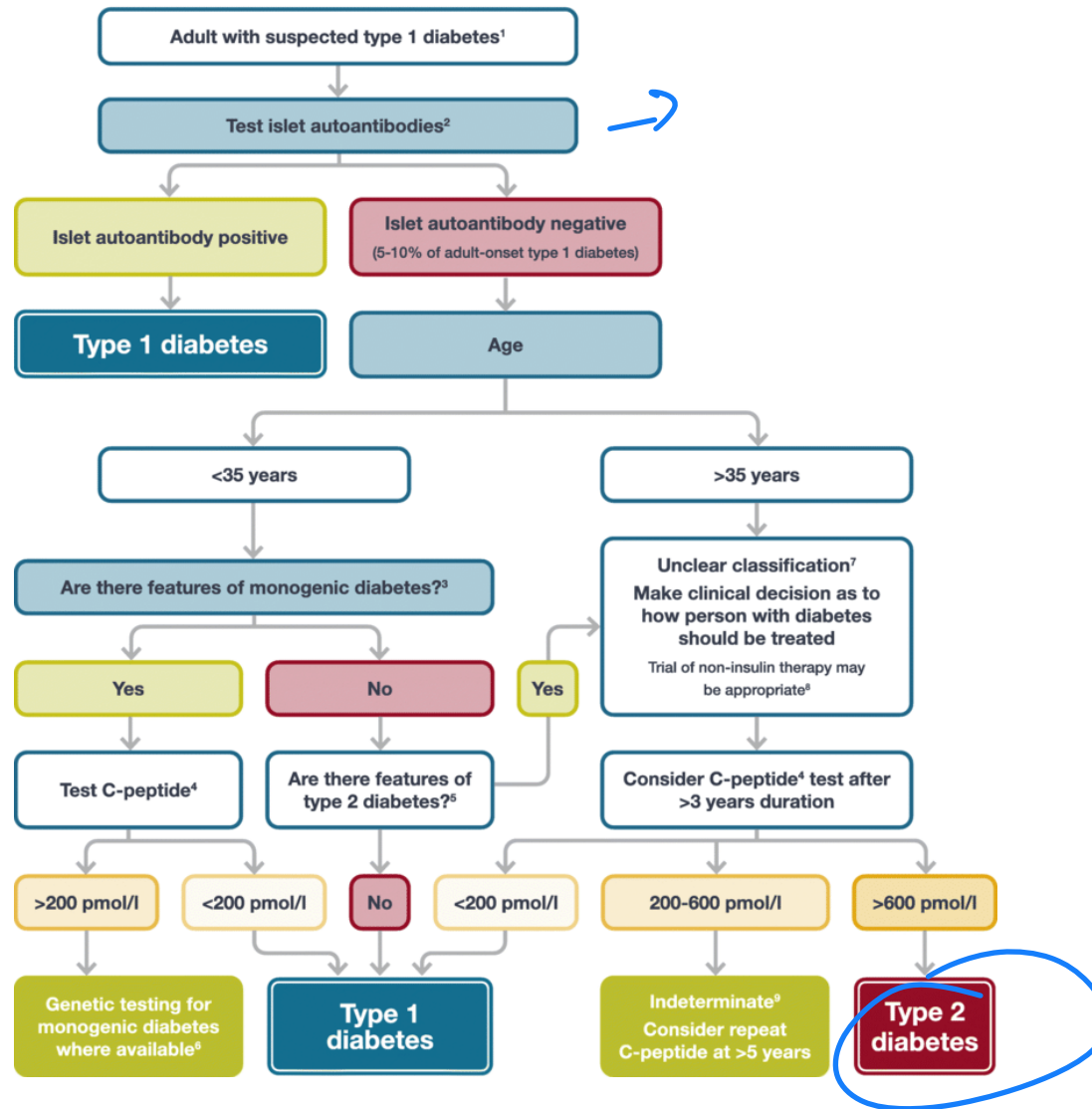
Non-parametric methods

Trees

Natasha Jaques

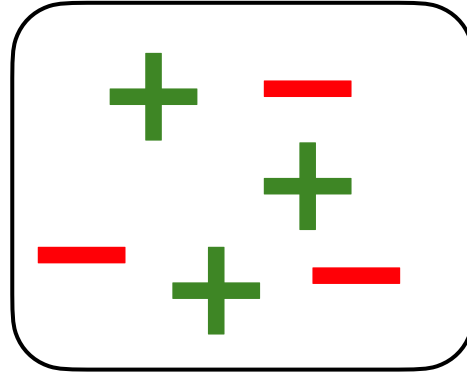


Flow chart for investigation of suspected type 1 diabetes in newly diagnosed adults, based on data from White European populations



Decision trees

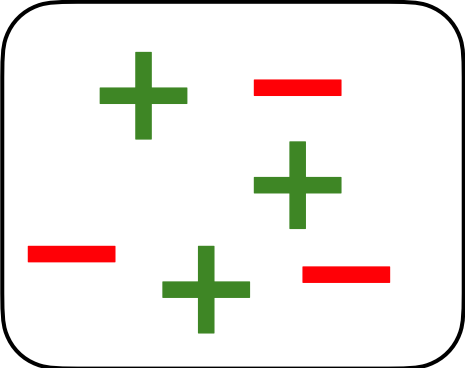
Want to classify spam from not spam using text-based features. What feature gives me the best split?



Full dataset: 50% spam, 50% not
Entropy: 1.0

Decision trees

Want to classify spam from not spam using text-based features. What feature gives me the best split?

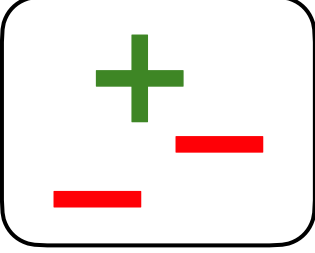
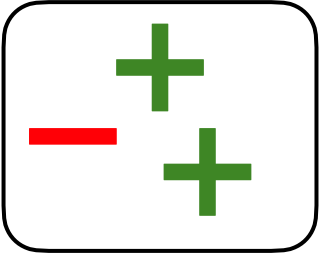


Full dataset: 50% spam, 50% not
Entropy: 1.0

Split on feature: contains “Limited offer”

Yes

No

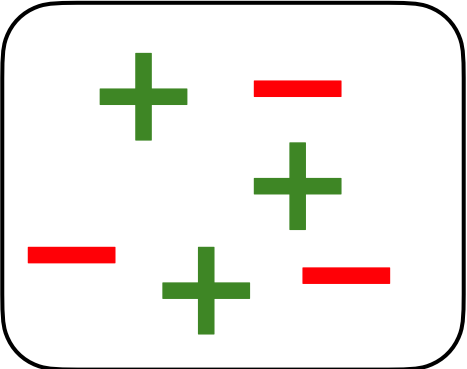


66% spam, 33% not
Entropy: 0.918

33% spam, 66% not
Entropy: 0.918

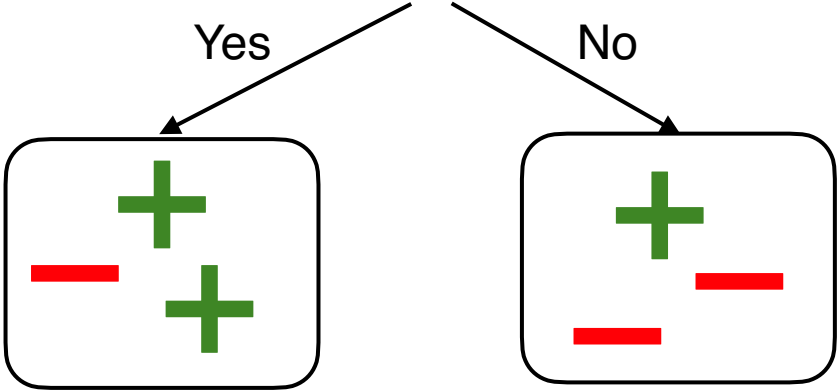
Decision trees

Want to classify spam from not spam using text-based features. What feature gives me the best split?



Full dataset: 50% spam, 50% not
Entropy: 1.0

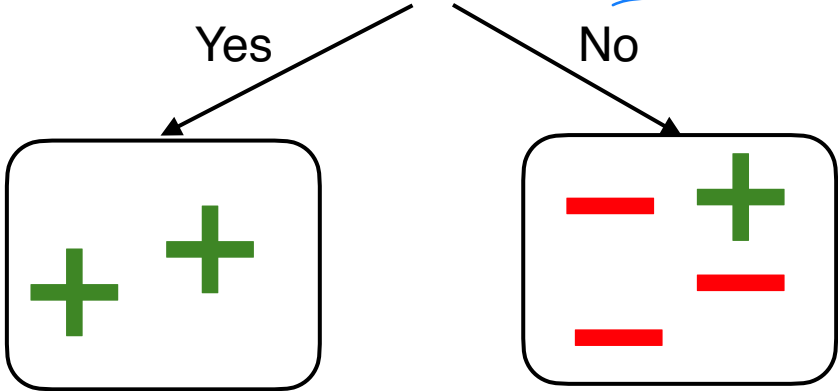
Split on feature: contains "Limited offer"



66% spam, 33% not
Entropy: 0.918

33% spam, 66% not
Entropy: 0.918

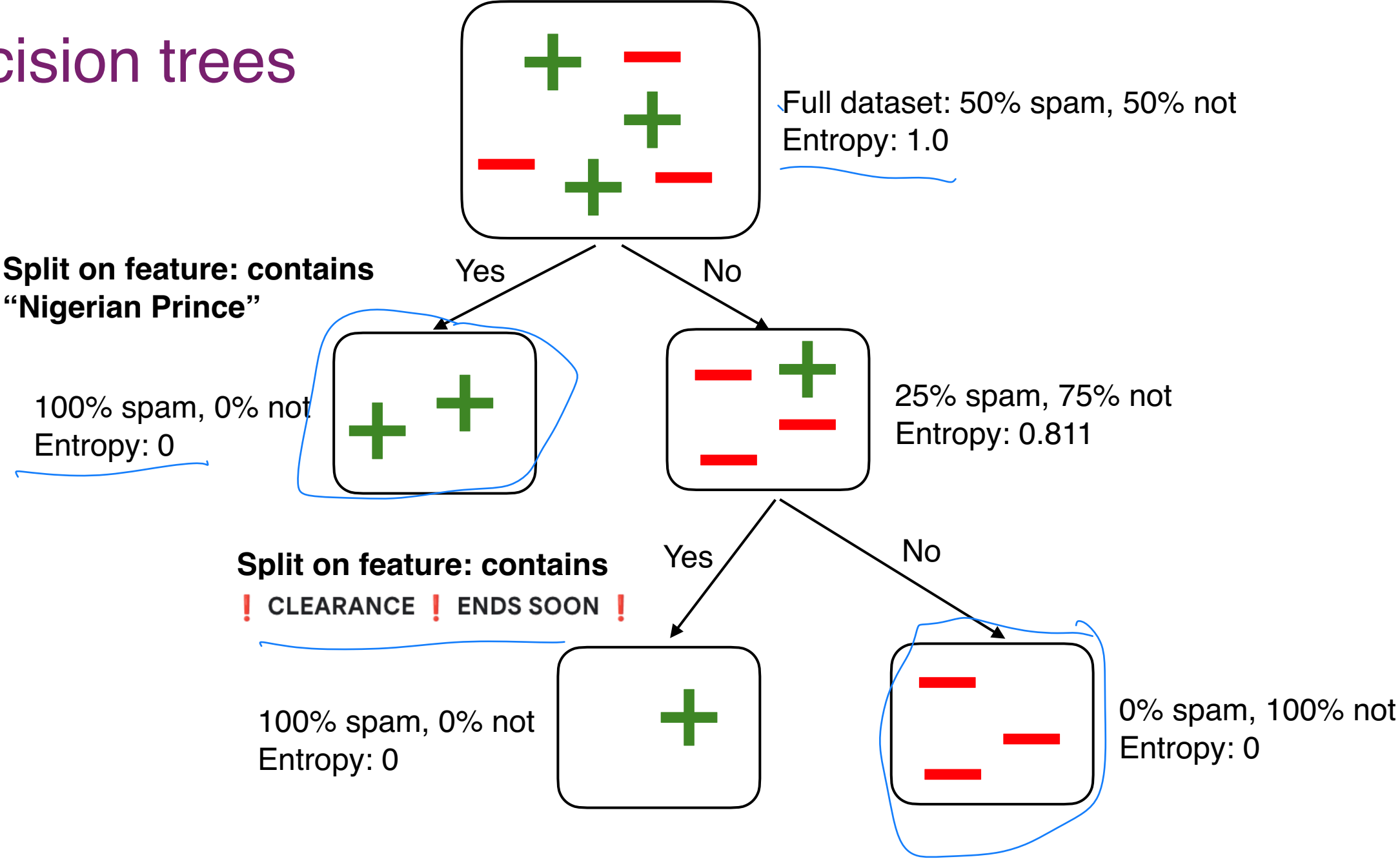
Split on feature: contains "Nigerian Prince"



100% spam, 0% not
Entropy: 0

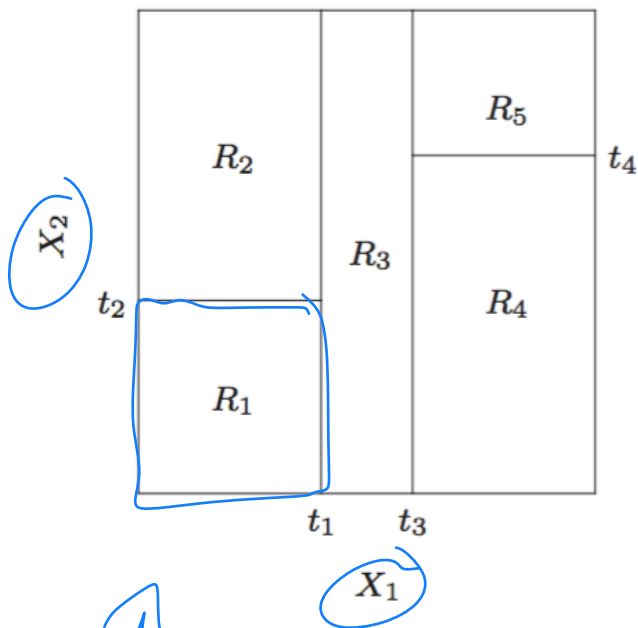
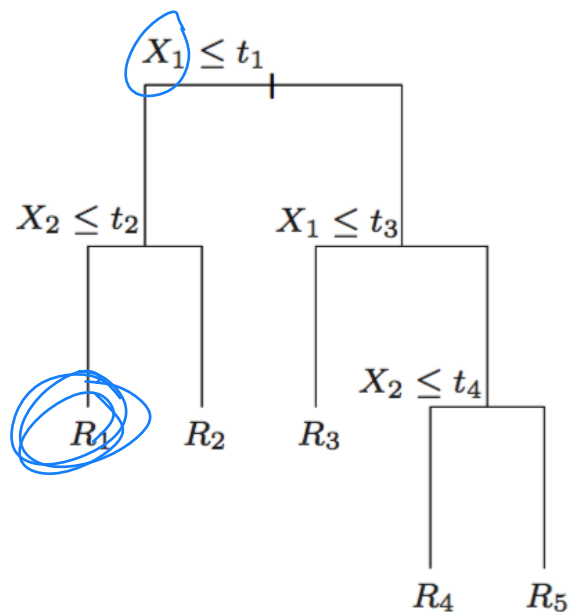
25% spam, 75% not
Entropy: 0.811

Decision trees



axis-aligned trees

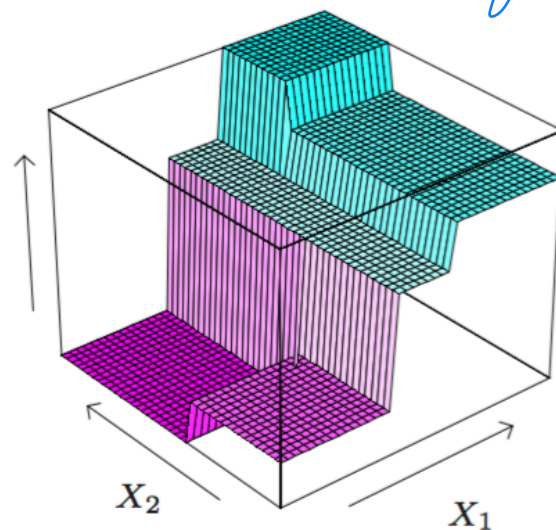
Regression trees



actual output value for regression

$$f(x) = \sum_{m=1}^M c_m I(x \in R_m).$$

$\rightarrow 1$ if $x \in R_m$
0 otherwise



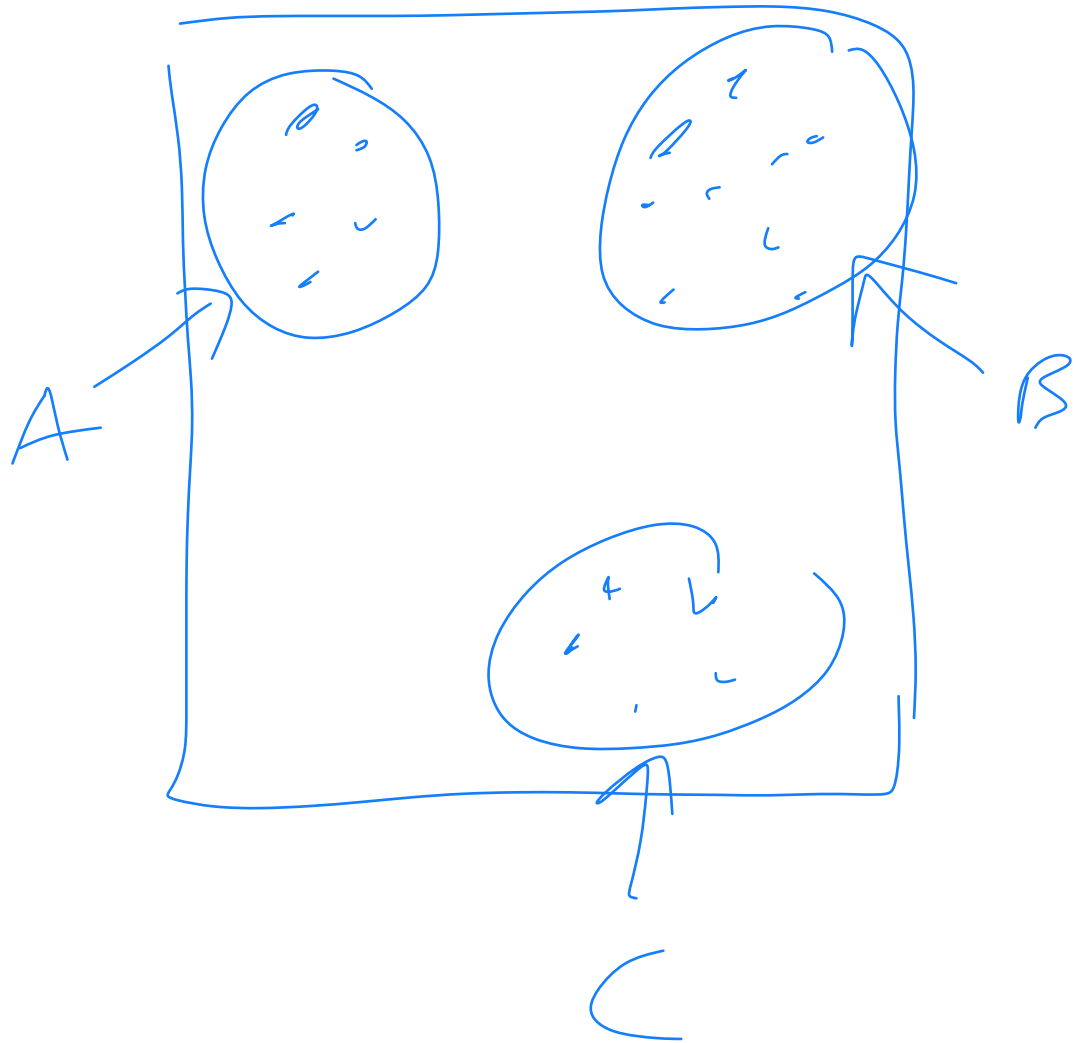
when to use?

(1) interpretable

↑ decision boundary

when not to use?

pixels, super complex high-d data



Generic algorithm for building trees

1. Start from empty decision tree
2. Recursively, for each node:
 - ~~Iterate through~~ all features and compute how good it'd be to split on each feature
 - Split on the "best" feature
3. Prune

Design choices:

- Termination condition
- Tree complexity
- Splitting criterion
- Pruning

max depth
entropy
train & validation error

Splitting regression trees

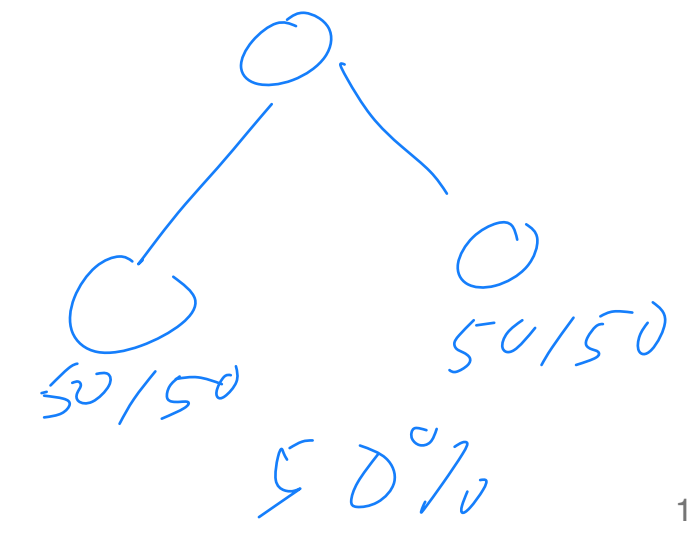
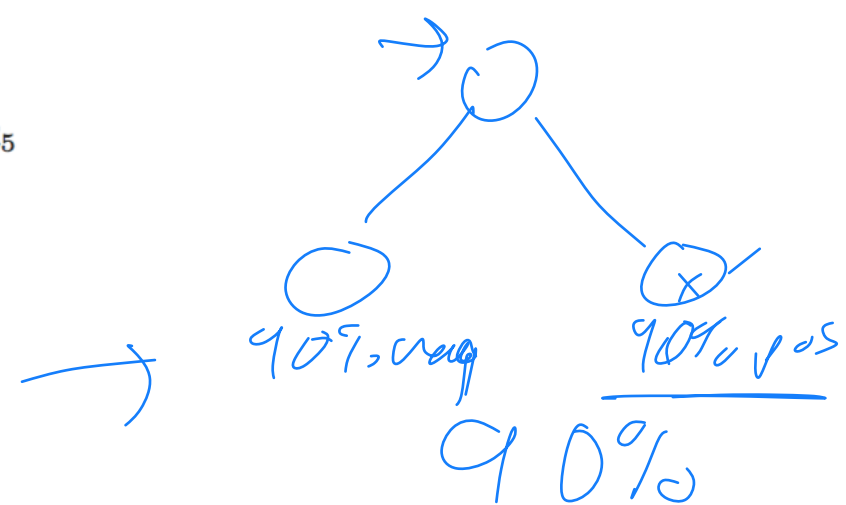
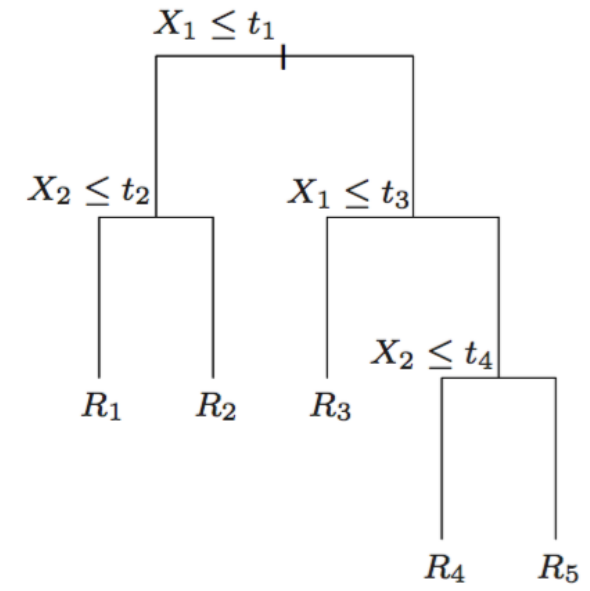
→ 0-1 error / binary

$$f(x) = \sum_{m=1}^M c_m \mathbb{1}_{\{x \in R_m\}}$$

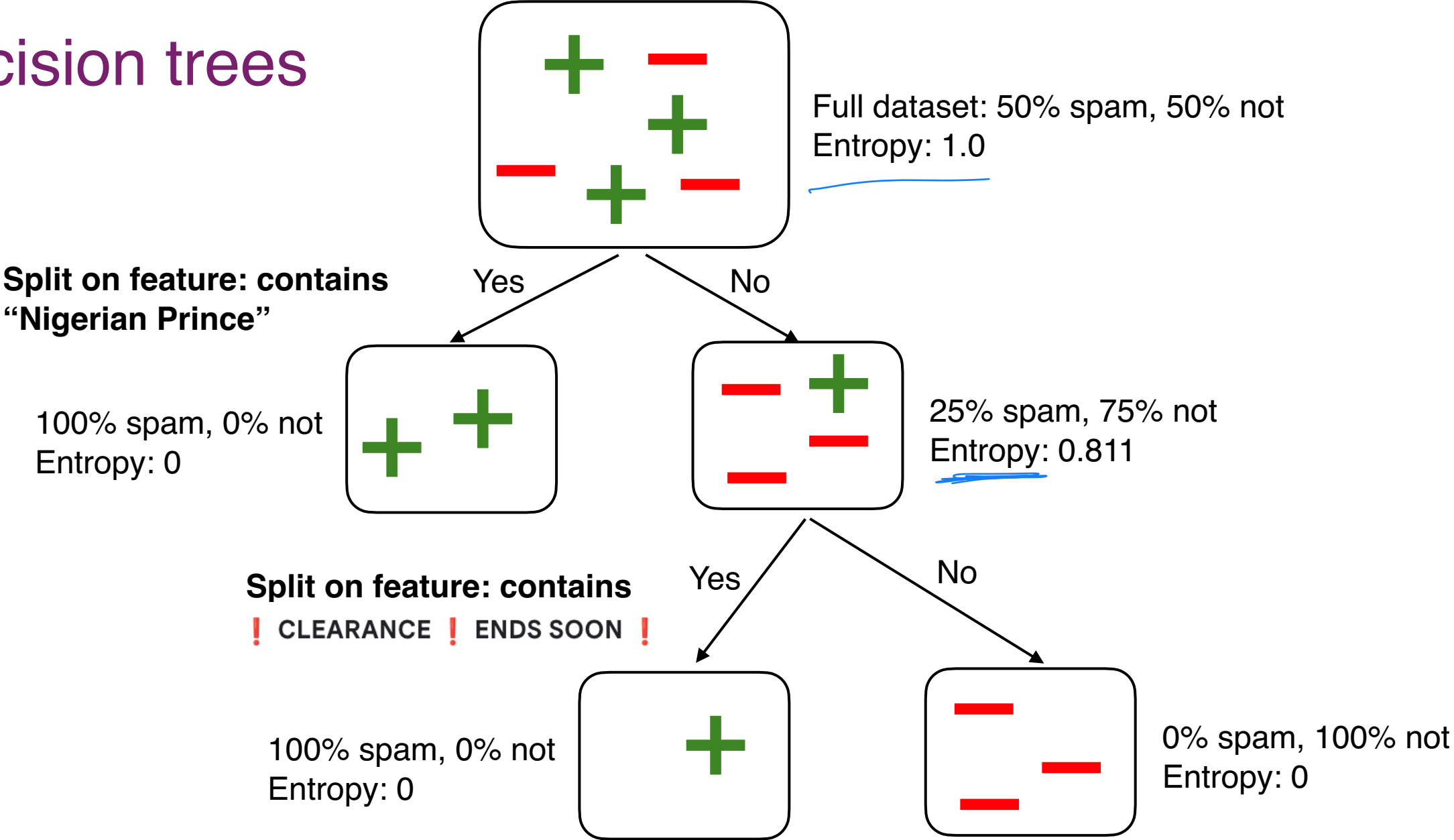
split on $x_i \leq s, x_i > s$

define "impurity" measure I
know well sorted are pos vs neg. exs.

find split that minimizes impurity.



Decision trees

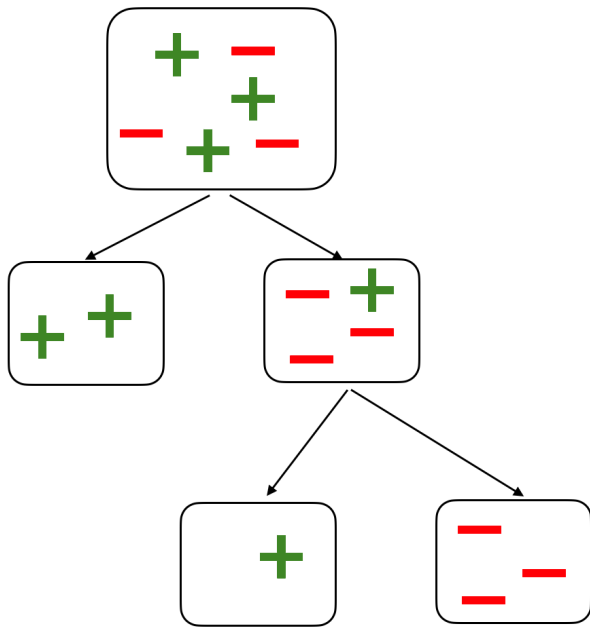


Splitting decision trees

impurity measures

- Gini index
- information gain

* \rightarrow is always $\log_2!$



Entropy $H(X) = - \sum_{x \in X} p(x) \log_2 p(x)$

set $X = \{0, 0, 1\}$ $p_0 = \frac{2}{3}$ $p_1 = \frac{1}{3}$

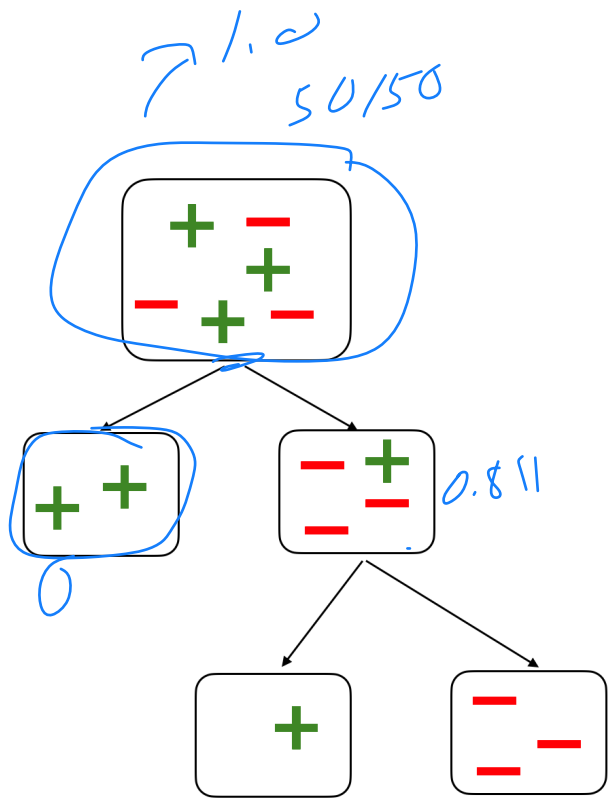
$$H = - \left(\frac{2}{3} \log \frac{2}{3} + \frac{1}{3} \log \frac{1}{3} \right) \approx 0.918$$

from 0 \rightarrow 1
perfectly sorted

\rightarrow uniform distribution

label	A	1.0
	B	0.0
	C	0.0
	D	0.0

Splitting decision trees



Information gain:

$$I(\underbrace{X}_{\text{set}}, \underbrace{a}_{\text{variable}}) = H(X) - H(X|a)$$

Decision trees

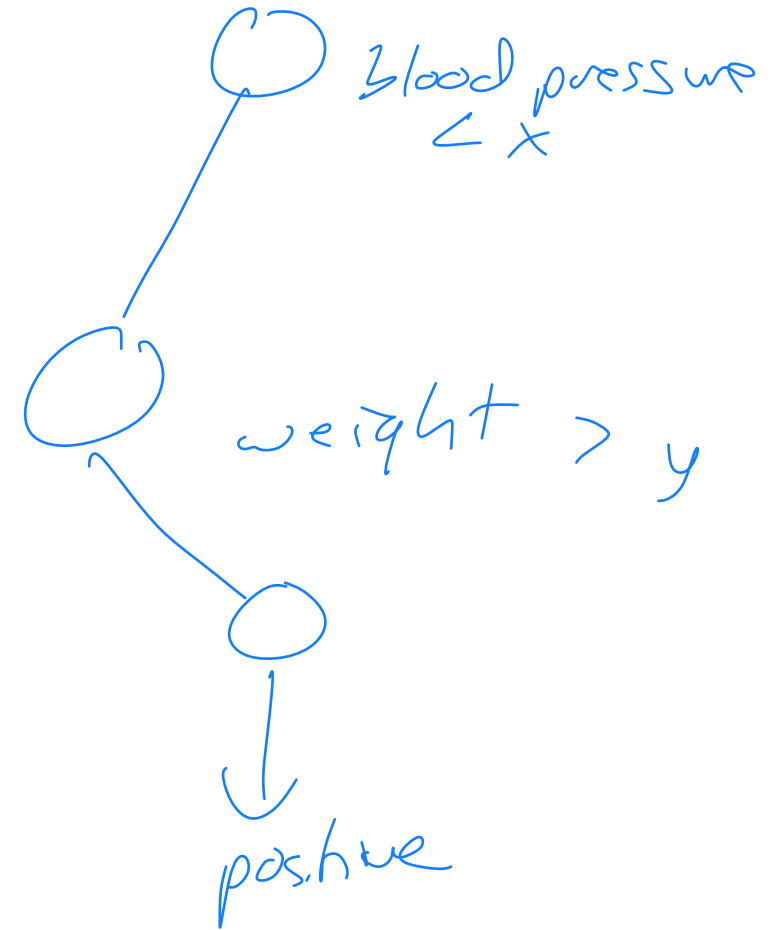
decision criteria for splitting based on $a \rightarrow [X_j = s]$ split n points

$$= H(X) - \sum_{i=1}^k \frac{|R_i(j, s)|}{n} H(R_i(j, s))$$

Interpreting trees

Trees are “easy” to interpret:

- You can explain how the classifier came to the conclusion it did



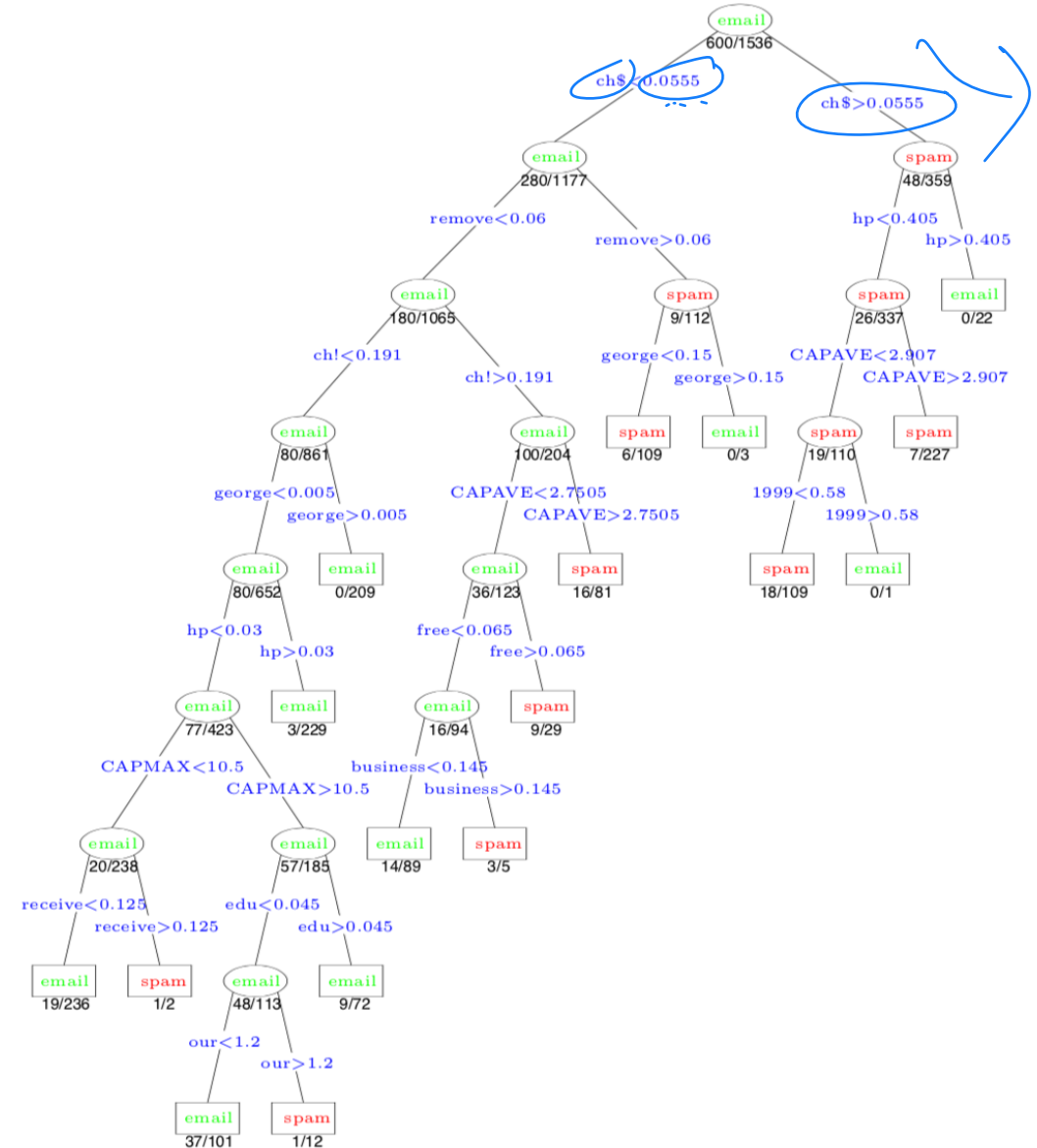
Interpreting trees

Trees are “easy” to interpret:

- You can explain how the classifier came to the conclusion it did

But they can be complex

- Small changes in data can result in large difference in trees



Summary so far

→ if you allow enough depth, fit anything

→ how regularize?
- limit # nodes
- limit depth
- be okay with imperfectly sorted leaves

• Trees have ↓ bias, ↑ variance

→ • Deal with categorical variables well

→ where don't they work?
many related continuous features
eg. pixels

• Intuitive, "interpretable"

• Good software exists

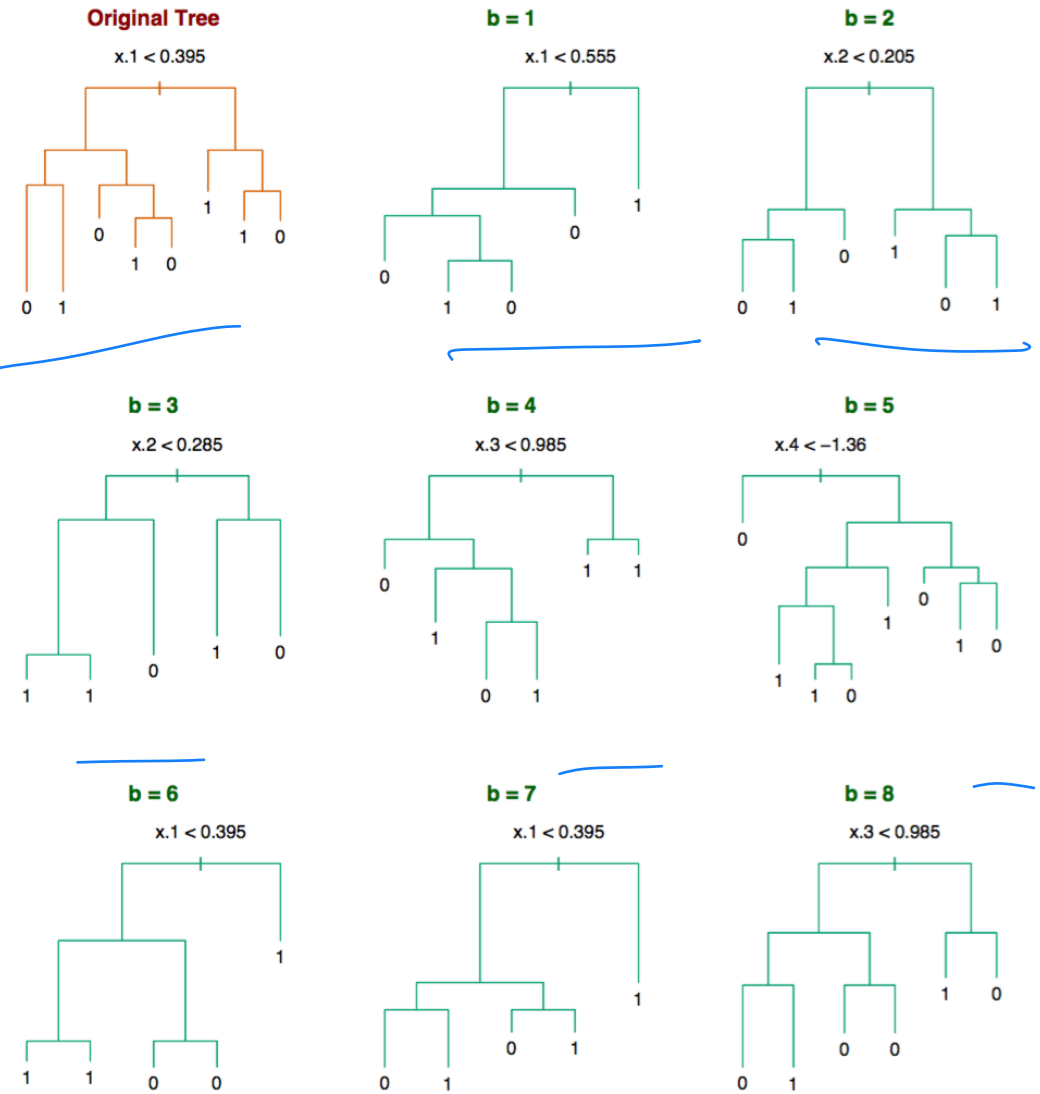
• Some theoretical guarantees

Random forests

- Forest = many trees
- Tree methods have low bias but high variance
- We can reduce variance by constructing many “lightly correlated” trees and averaging them

- Bagging: Bootstrap aggregating

each tree trained on new bootstrap sample of data



Random forests

Algorithm 15.1 *Random Forest for Regression or Classification.*

1. For $b = 1$ to B :

→ with replacement

(a) Draw a bootstrap sample \mathbf{Z}^* of size N from the training data.

(b) Grow a random-forest tree T_b to the bootstrapped data, by recursively repeating the following steps for each terminal node of the tree, until the minimum node size n_{min} is reached.

i. Select m variables at random from the p variables.

ii. Pick the best variable/split-point among the m .

iii. Split the node into two daughter nodes.

frac tree

2. Output the ensemble of trees $\{T_b\}_1^B$.

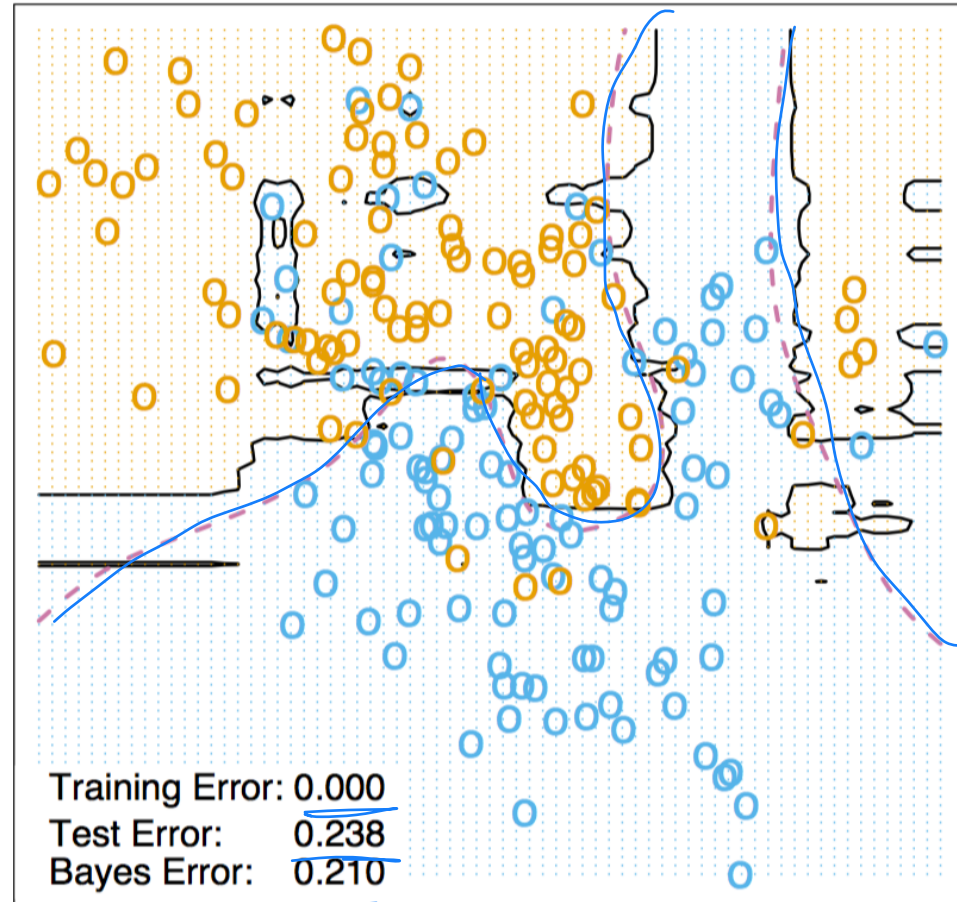
use a random subset of the feats

To make a prediction at a new point x :

Regression: $\hat{f}_{\text{rf}}^B(x) = \frac{1}{B} \sum_{b=1}^B T_b(x)$. *a average*

Classification: Let $\hat{C}_b(x)$ be the class prediction of the b th random-forest tree. Then $\hat{C}_{\text{rf}}^B(x) = \text{majority vote } \{\hat{C}_b(x)\}_1^B$.

Random forests: Decision boundary example



→ Bayes optimal

Random forests

Average of lightly correlated variables has lower variance → trees

Given random variables Y_1, Y_2, \dots, Y_B with
 $\mathbb{E}[Y_i] = y$, $\mathbb{E}[(Y_i - y)^2] = \sigma^2$, $\mathbb{E}[(Y_i - y)(Y_j - y)] = \rho\sigma^2$

Sampled

variance

correlation b/w predictors

correlation coefficient $0 < \rho = 1$
 ↓
 no relationship duplicate

$$\mathbb{E}\left[\left(\frac{1}{B} \sum_{i=1}^B Y_i - y\right)^2\right] = \mathbb{E}\left[\left(\frac{1}{B} \sum_{i=1}^B (Y_i - y)\right)^2\right]$$

→ # of trees

$$= \mathbb{E}\left[\frac{1}{B^2} \sum_{i=1}^B (Y_i - y)^2\right] + \mathbb{E}\left[\frac{1}{B^2} \sum_{i \neq j} (Y_i - y)(Y_j - y)\right]$$

$$= \frac{\sigma^2}{B} + \frac{B-1}{B} \rho \sigma^2$$

if $\rho < 1$:
 can reduce variance by increasing B

if $\rho = 1$
 $\frac{\sigma^2}{B} + \frac{(B-1)\sigma^2}{B}$
 $\frac{B}{B} \sigma^2$

Ensembles are powerful

pc 1

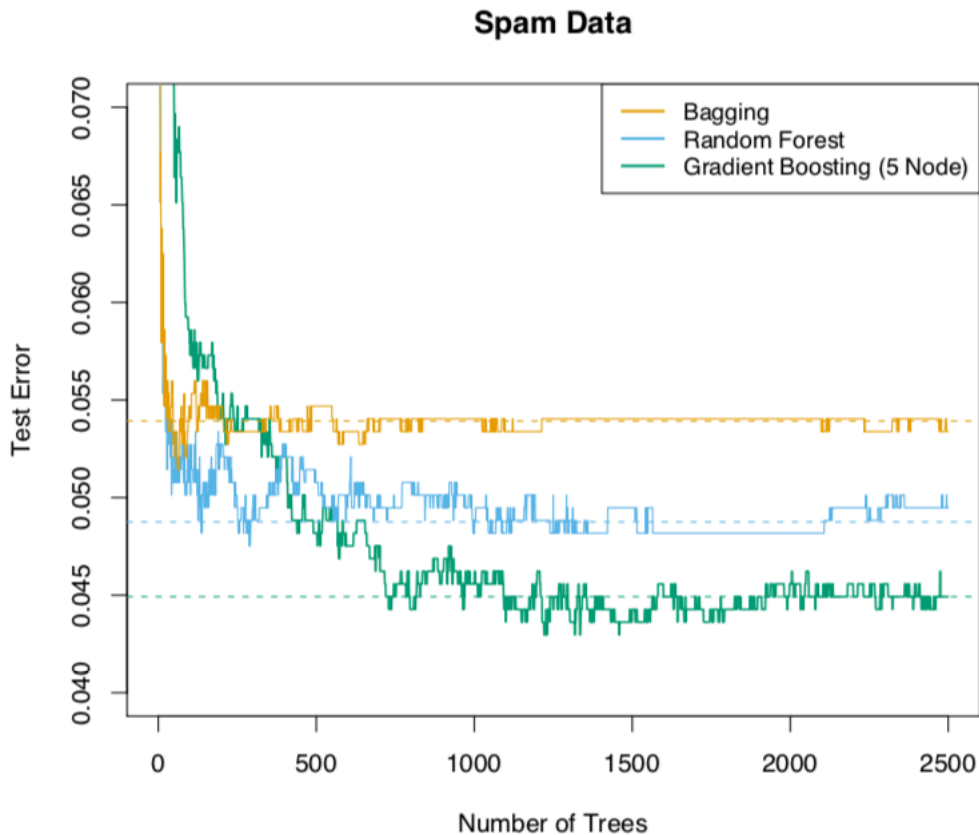
For binary classification, if you have N weak classifiers, but each one is slightly better than random chance (gets right answer with $p > 0.5$), what happens if you take the majority vote?

As $N \rightarrow \infty$, what is the probability that the majority vote gets the right answer?

$\rightarrow 100\%$

- Marquis de Condorcet, "Essay on the Application of Analysis to the Probability of Majority Decisions" (1785). Known as the Condorcet Jury Theorem.
- Schapire, "The Strength of Weak Learnability" (1990)

The power of weakly correlated predictors



→ Bagging: Averaged trees on bootstrapped datasets using all d features

Random forest: Averaged trees on bootstrapped datasets using m randomly selected features

↳ less correlated

Takeaways:

- Reducing correlation improves performance
- Ensembles are powerful

Summary so far

- Random forests have ↓ bias, ↓ variance
- ✓ • Deal with categorical variables well
- ✗ • Not that intuitive nor “interpretable”
- ✓ • Gives some notion of confidence estimates ✓✓✓
- Good software exists
- Some theoretical guarantees

Boosting and additive models

Instead of ensembling bootstrapped models, can we:

- Keep the idea of ensembling / combining simpler models, but
- Not necessarily have the models be identically distributed?

Key idea: Given a current collection of models, add a new model that focuses on what the previous models got wrong

Boosting

Additive models

classification

Given $\{(x_i, y_i)\}_{i=1}^n$ $x_i \in \mathbb{R}^d$, $y_i \in \{+1, -1\}$

$b_m: \mathbb{R}^d \rightarrow \mathbb{R}$ $m = 1 \dots M$ // a model

$\beta_m \in \mathbb{R}$ // weights ~~of a particular model~~
parameters

$$\hat{\beta}_1, \hat{b}_1, \dots, \hat{\beta}_M, \hat{b}_M = \underset{\dots}{\operatorname{arg\,min}} \sum_{i=1}^n d(y_i, \sum_{m=1}^M \beta_m b_m(x_i))$$

// too hard!
Simultaneous
opt. of so
many params

Sequentially learn each new model (greedily)
one at a time.

Forward stagewise additive models

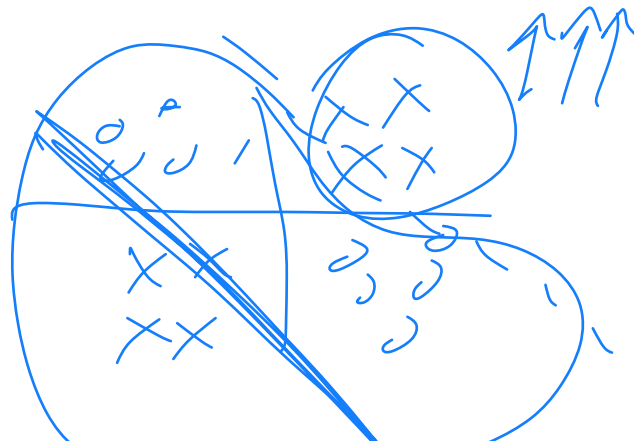
Algorithm 10.2 Forward Stagewise Additive Modeling.

1. Initialize $f_0(x) = 0$. *→ ensemble*
2. For $m = 1$ to M : *→ for all models in ensemble*
 - (a) Compute

$$(\beta_m, \gamma_m) = \arg \min_{\beta, \gamma} \sum_{i=1}^N L(y_i, f_{m-1}(x_i) + \beta b(x_i; \gamma)).$$

← h param
← previous

- (b) Set $f_m(x) = f_{m-1}(x) + \beta_m b(x; \gamma_m)$.
-



Adaboost *→ binary classifier*

$$b(x, \gamma) : \mathbb{R}^d \rightarrow \{-1, 1\}$$

$$L(y, f(x)) = e^{-y f(x)}$$

// upweight points the previous ensemble got wrong

↳ "boosts" the weight on mistakes

Forward stagewise additive models

Algorithm 10.2 *Forward Stagewise Additive Modeling.*

1. Initialize $f_0(x) = 0$.
2. For $m = 1$ to M :
 - (a) Compute

$$(\beta_m, \gamma_m) = \arg \min_{\beta, \gamma} \sum_{i=1}^N L(y_i, f_{m-1}(x_i) + \beta b(x_i; \gamma)).$$

- (b) Set $f_m(x) = f_{m-1}(x) + \beta_m b(x; \gamma_m)$.
-

Boosted regression
trees

$b(x; \gamma) =$ regression tree

$$L(y, f(x)) = (y - f(x))^2$$

$$L(\underbrace{f_{m-1}(x_i)}_{\text{previous ensemble}} + \underbrace{\beta b(x_i, \gamma)}_{\text{new node}}) = \underbrace{(y_i - f_{m-1}(x_i) - \beta b(x_i, \gamma))}_{\text{residual}}^2$$

~~residual~~

$$(r_{im} - \beta b(x_i, \gamma))^2$$

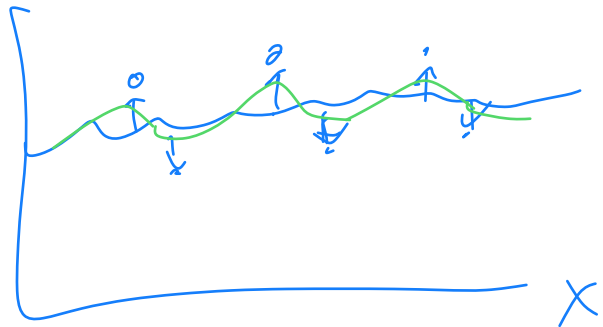
Gradient boosting

take gradient steps to minimize loss in "function space"

$$L(y_i, f(x_i)) = L(y_i, \underline{f_i}) = (y_i - f_i)^2$$
$$\frac{\partial L(y_i, f_i)}{\partial f_i} = -2(y_i - f_i)$$

gradient descent in function space

$f(x)$



$$f_i := \underline{f} + \beta \underline{b} \rightarrow \text{step size} \rightarrow \text{gradient direction}$$

→ gradient of the func at the parts where you have data; move func towards those points

Gradient boosting generalizes boosting to arbitrary loss funcs

A brief history of boosting

a rare case of
theory to practice

- 1988 Kearns and Valiant: “Can weak learners be combined to create a strong learner?”
- 1990 Schapire: “Yup, in theory” or 1785?
- 1995 Schapire and Freund: “Practical for 0/1 loss” -> AdaBoost
- 2001 Friedman: “Practical for arbitrary losses”
- 2014 Tianqi Chen: “Scale it up!” -> XGBoost

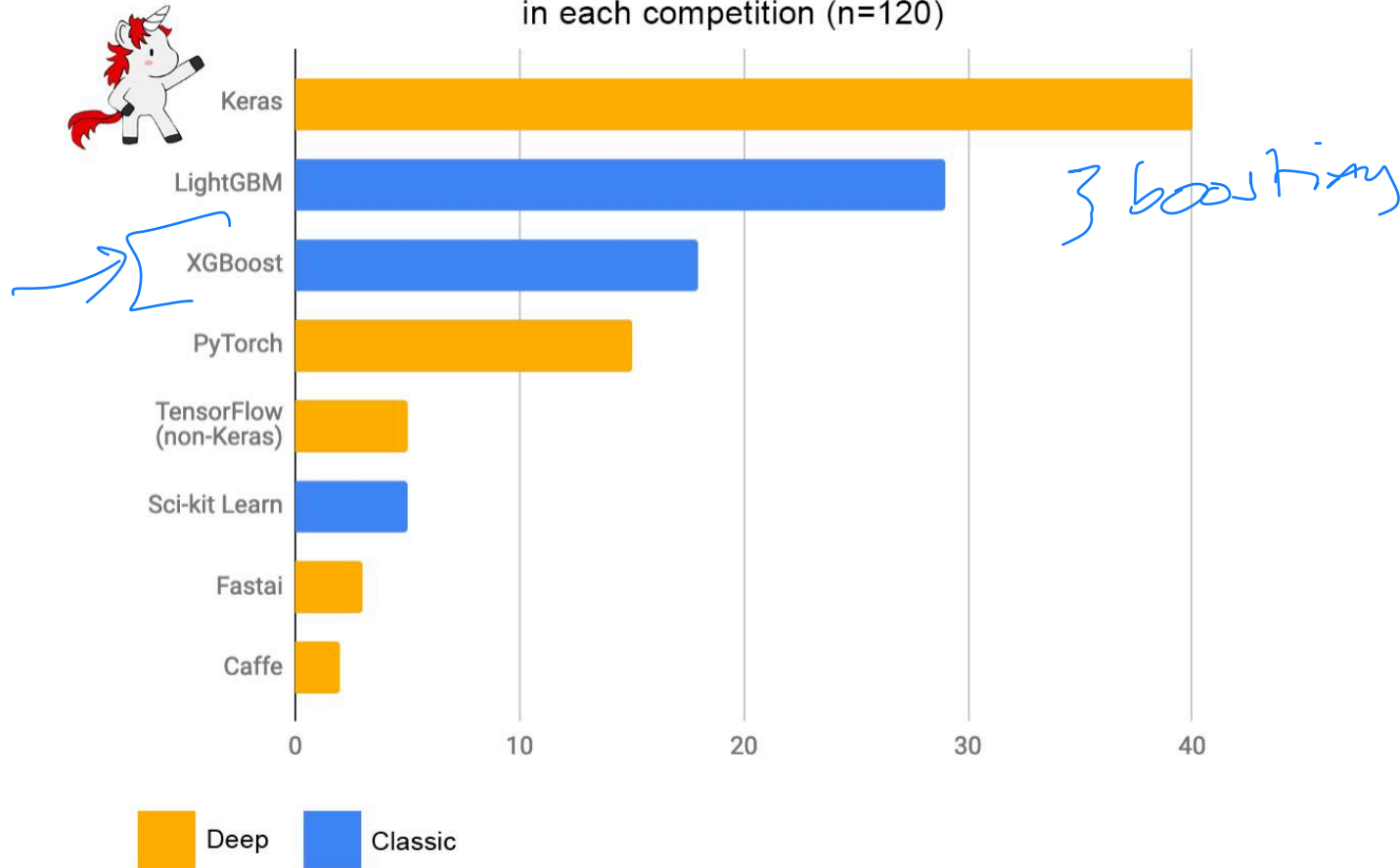
→ aw a lot!

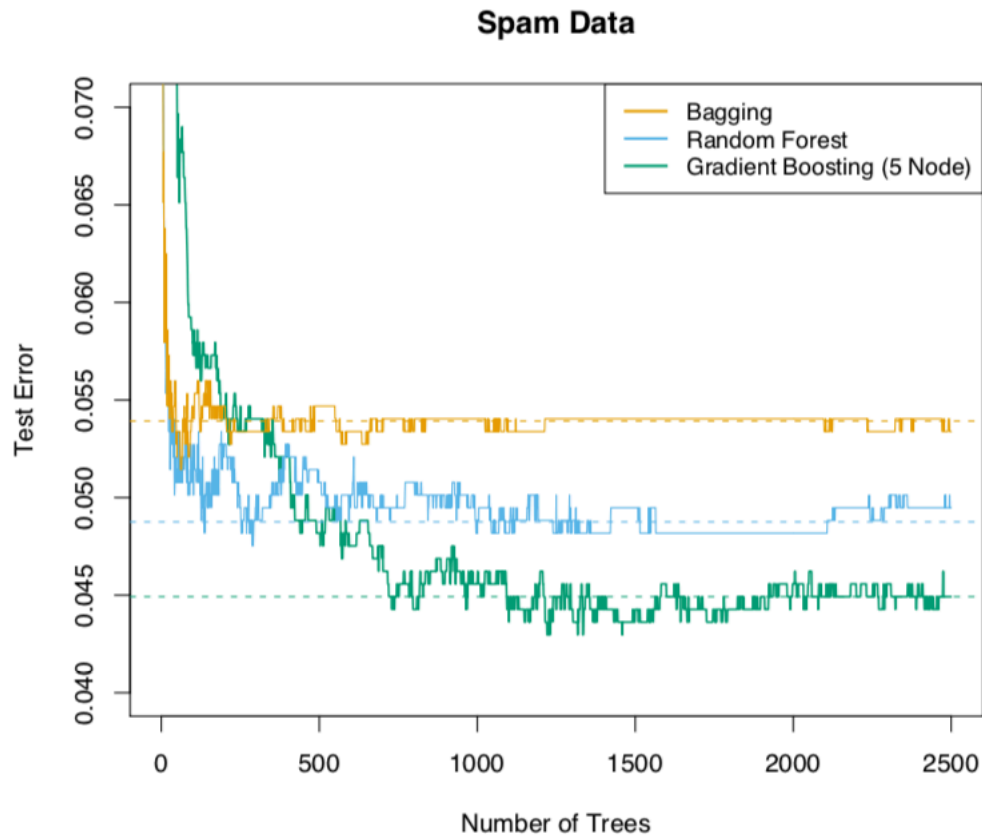


François Chollet ✓ @fchollet · Apr 3, 2019

What machine learning tools do Kaggle champions use? We ran a survey among teams that ranked in the *top 5* of a competition since 2016.

Primary ML software tool used by **top-5 teams** on Kaggle in each competition (n=120)





Bagging: Averaged trees on bootstrapped datasets using all d features

Random forest: Averaged trees on bootstrapped datasets using m randomly selected features

Boosting: Learned combinations of trees

Takeaways

- Single trees: low bias, high variance
- Ensembles: low bias, (relatively) low variance
- Bagging averages many lightly dependent models to reduce variance
 - Random forests: same but with random subset of features
- Boosting learns a linear combination of high bias, highly dependent classifiers to reduce error
- Gradient boosted trees are commonly used for categorical data