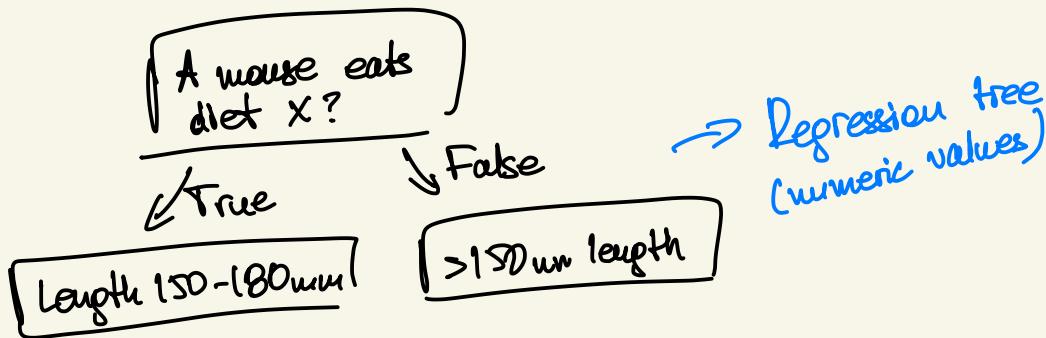
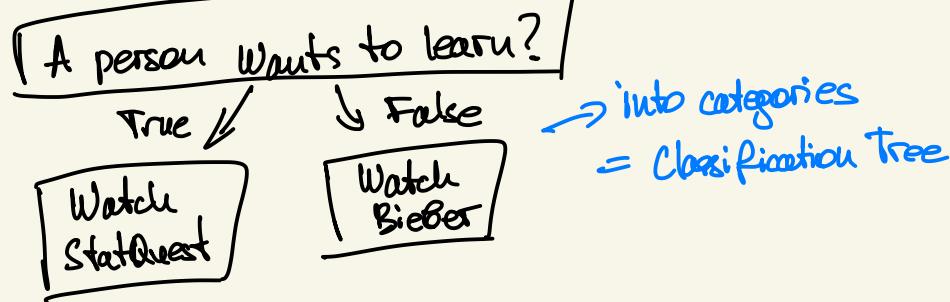


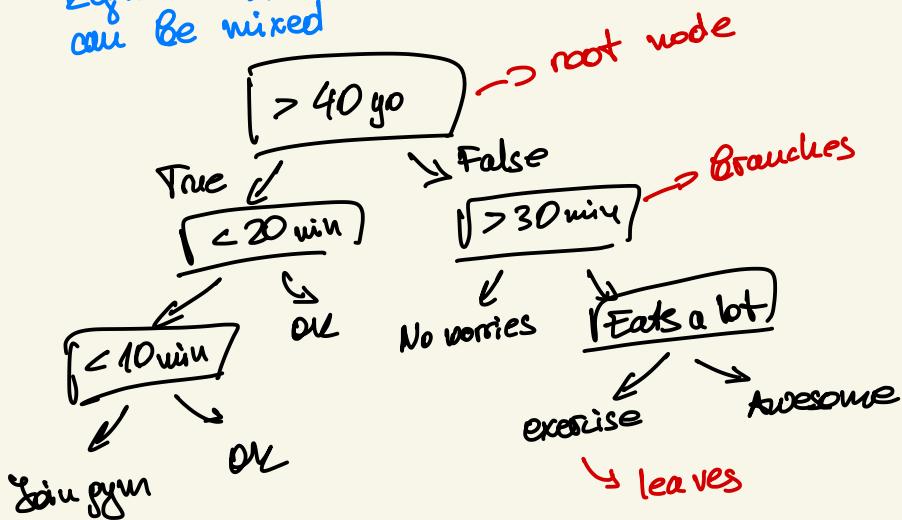
Trees and Forests
(StatQuest)



Decision Trees



Regression & classifications can be mixed



- Build tree from data
How well each feature predict outcome? → Best fit is root
↳ Impure leaves

- Gini Impurity:

$$= 1 - (\text{prob of yes})^2 - (\text{prob No})^2$$

$$= 1 - \left(\frac{1}{4}\right)^2 - \left(\frac{3}{4}\right)^2 = 0,375 \text{ for a leaf}$$

Weighted for both leaves:

$$\frac{4}{(4+3)} \cdot 0,375 + \frac{3}{4+3} \cdot 0,444 = 0,405 \text{ for a feature}$$

soda?

Yes	No	Yes	No
$\frac{1}{3}$	$\frac{1}{3}$	$\frac{0}{3}$	$\frac{3}{3}$

- Gini impurity for numeric value

sort:

$\{ \begin{matrix} 9,5 \\ 15 \\ 26,5 \end{matrix} \} > \{ \begin{matrix} 12 \\ 13 \\ 35 \end{matrix} \}$

for all
→ lowest impurity
of all features
→ root

$\boxed{\text{Age} < 9,5}$

Yes $\frac{0}{6}$ No $\frac{1}{1}$ Yes $\frac{4}{3}$ No $\frac{2}{3}$

$$\text{Imp}_L = 1 - \left(\frac{0}{0+1} \right)^2 - \left(\frac{1}{0+1} \right)^2 = 0$$

$$\text{Imp}_R = 1 - \left(\frac{3}{6} \right)^2 - \left(\frac{3}{6} \right)^2 = 0,5$$

$$\text{Imp}_{\text{TOTAL}} = \frac{1}{1+6} \cdot 0 + \frac{6}{1+6} \cdot 0,5 = 0,429$$

- Nodes after root

reduce with node impure
other features

$\boxed{\text{Soda?}}$

Yes $\frac{3}{3}$ No $\frac{1}{1}$

Yes $\frac{1}{0}$

No $\frac{3}{3}$

$\boxed{\text{Popcorn?}}$

Yes $\frac{1}{1}$ No $\frac{1}{1}$ Yes $\frac{2}{2}$ No $\frac{0}{0}$

$\text{Imp}_{\text{TOTAL}} = 0,25$

$\boxed{\text{Age} < 12,5}$

Yes $\frac{0}{0}$ No $\frac{1}{1}$ Yes $\frac{2}{3}$ No $\frac{0}{0}$

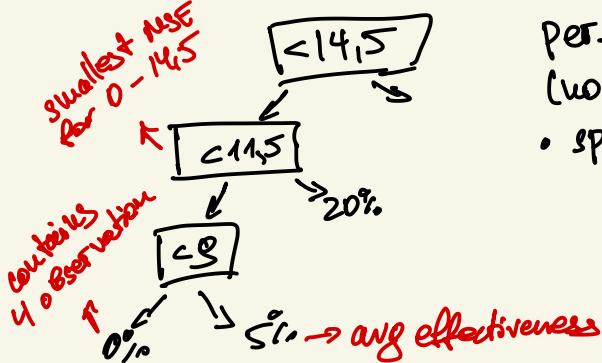
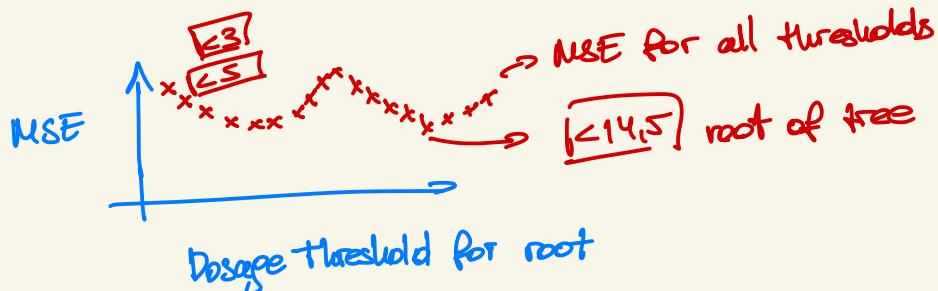
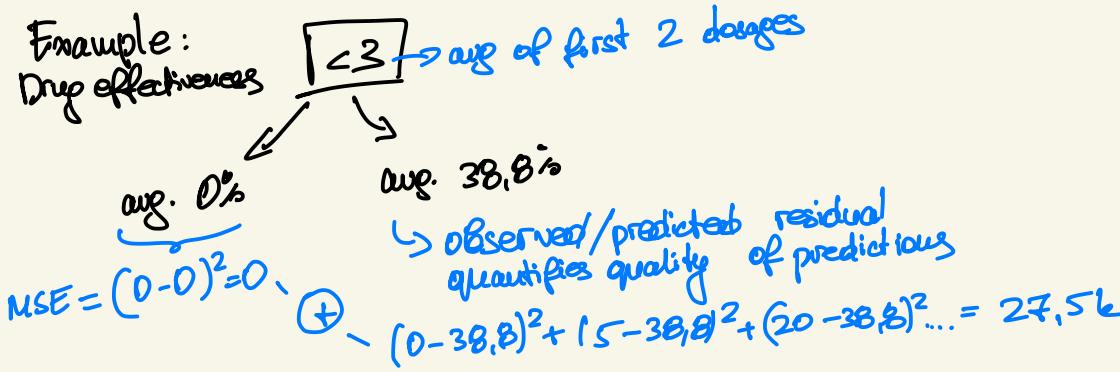
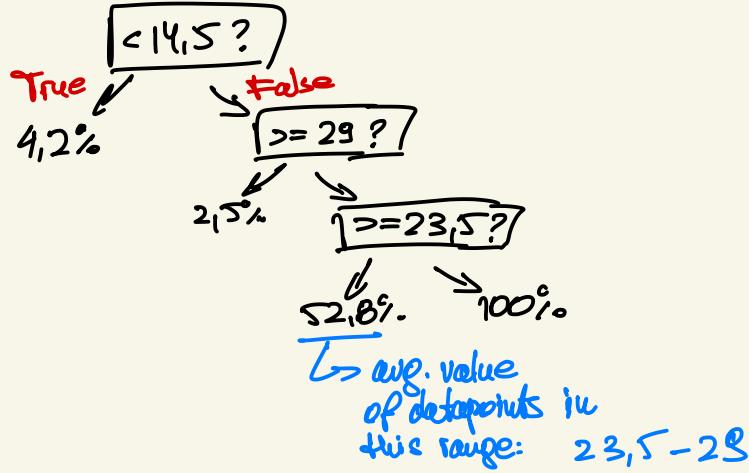
$\text{Imp}_{\text{TOTAL}} = 0 \rightarrow \text{use to split node}$

- Output of a leaf = category with most values

- Overfit if only 1 person gets to leaf
require > 3 people per leaf to prevent overfit
- Pruning for overfit
- Cross-validation

Regression Trees

each leaf represents numeric value



perfect fit may overfit
(no bias, but large variance)

- split only when > 20 observations

- Multiple Features

- sort numeric values of feature and calculate MSE for root
- choose lowest MSE
- if categorical, set $\text{Sex} = \text{Female?}$ and calculate MSE
 $\hookrightarrow \text{no } > \text{ or } <$
- stop when observations < 20

Prune Regression Trees

- many versions, here \rightarrow cost complexity pruning or weakest link pruning
- large residuals, best better when testing (overfitting)
- remove leaves, decrease depth

- Calculate SSR for each tree (sum of squared residuals)

Full tree: for each leaf: $320 + 75 + 148,8 + 0 = 543,8$

- 1 leaf: $SSR = 5,5k$

- 2 leaves: $SSR = 18,2k$

- 3 leaves: $SSR = 28,8k$

$$\text{Tree Score} = SSR + \alpha \cdot T$$

α = tuning parameter from cross validation

T = complexity penalty f/number of leaves)

$$\text{Tree}_1 = 543,8 + 10k \cdot 4 = 40,5k$$

$$\text{Tree}_2 = 5,5k + 10k \cdot 3 = 35,5k \rightarrow \text{pick lowest score}$$

$$\text{Tree}_3 = 18,2k + 10k \cdot 2 = 38,2k$$

$$\text{Tree}_4 = 28,8k + 10k = 38,8k$$

- start with $\alpha = 0$ and increase α until removing leaves gives tree lower score

$\alpha = 0$ = full size \rightarrow calculate SSR for testing
for each new tree

$\alpha = 10k = -1$ leaf

$\alpha = 15k = -2$ leaves

$\alpha = 22k = -3$ leaves

with new train/test data, build new trees, calculate α ,
calculate SSR 10 times (10 times cross validation)
lowest SSR α is the final α value

\hookrightarrow use α in the full data

Cross Validation

- divide 25% - 75% test-train (4 Fold Cross Validation)
- use each 25% block as testing block
- keep track of how well method did with test data
 - ↳ Log Regression / SVM / k nearest
- Leave One Out Cross validation (tests each sample individually)
- In practice: 10 Fold
- tuning parameter (guessed) cross valid. can determine it

Decision Trees Feature Selection / Missing Data

- predict missing values with regression
- select features with PCA

Random Forest

- Training

- Original Dataset
 - random samples
 - some samples more than once

Bootstrapped Dataset

- Create Decision Tree, but use random features (columns), not all, for root node
- after root, from the remaining all columns, randomly select subset of features and determine best fit for the node
- Build as usual the tree
- repeat whole process 100+ times

- Testing

- run data through each generated tree
- select class based on most votes

Bootstrapping + aggregate to make decision = Bagging

$\frac{1}{3}$ of original data is not in Bootstrapped dataset

↳ out-of-bag dataset

check which trees label out-of-bag sample falsely \Rightarrow out-of-bag error

↳ falsely classified
correctly classified

- Compare hyperparameters

- change number of vars used per step
- # vars start, then try more or less

Random Forest Missing Data / Sample Clustering

- in original data: make guess based on label, other features
- after running data through a tree, keep similar (ended on same leaf) samples in proximity matrix

Add proximity matrices
after all trees have
been executed

- divide matrix values by number of trees

	1	2	3	4
1				
2				
3				
4				

samples

←

→ ①

→ ①

3/4 ended up on same leaf

Blocked Arteries

	1	2	3	4
1				
2				
3				
4	0,1	0,1	0,8	

Bottle nose 'no' $\Rightarrow 0,8$

$$\frac{0,1}{0,1+0,1+0,8} = 0,1 \rightarrow \text{weight for 'yes'}$$

frequency of 'yes' = $\frac{1}{3}$ of all occurrences

$$\frac{1}{3} \cdot 0,1 = 0,03 \text{ weighted frequency}$$

$$'No': \frac{0,1+0,8}{3} \cdot \frac{2}{3} = 0,6 \text{ weighted frequency}$$

\Rightarrow No is more likely for this sample

- repeat whole process until missing values converge \Rightarrow don't change when we recalculate $\sim 6-7$ times

- Proximity matrices

- largest number in proximity = 100% close
- 1-proximity = distance between samples
- heatmap, mds plot of distances

- Missing data in new samples

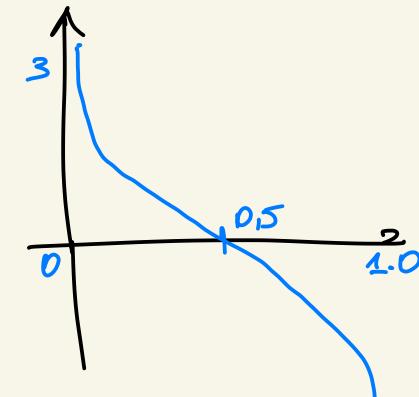
- 2 copies - with/without heart disease (label)
- guess missing features and perform the iterative process until it converges
- run through random forest and see which was labeled correctly most of the times
- the option with most correct wins

Ada Boost

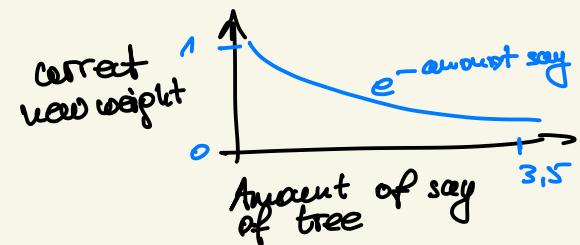
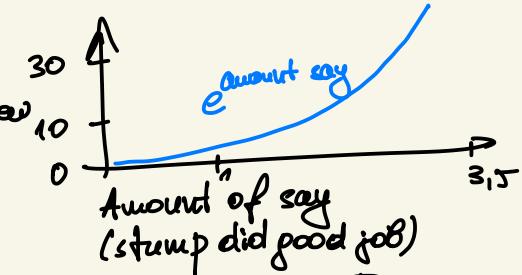
- in forest of AdaBoost trees have 1 node + 2 leaves (stump)
- stumps are weak learners
- some stumps have larger vote
- order is important, errors of prev. stump influence next stump

- Creation

- add feature, Sample Weight = (starts at $\frac{1}{\text{total samples}}$) = sample importance
- find which feature best classifies the samples and use it as root (using Gini Index)
- Determine stump Total Error (sum of weights of incorrectly classified samples)
 - ↳ Total Error = 0 - 1
- Amount of stump's say = $\frac{1}{2} \log \left(\frac{1 - \text{Total Error}}{\text{Total Error}} \right)$
- if Total Error $\neq 0/1$, \Rightarrow add small ϵ
- Emphasize wrongly classified samples weights:
 - new weight (incorrect) = sample weight $\cdot e^{\text{amount of say}}$
 - correct sample weight = sample weight $\cdot e^{-\text{amount of say}}$
 - normalize new sample weights to add to 1
- Use new sample weights for new stumps
- Generate random num 0-1 and use weights as distribution to select samples in new collection (more heavy weight samples)
- assign equal weights to all samples again and repeat whole process
- on classify - add amounts-of-say for same classification trees

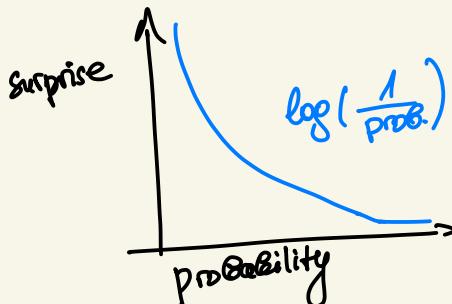


-3 \rightarrow stump gives wrong class constantly



Entropy

- low probability = high surprise
- high prob. = low surprise
- surprise = $-\log \left(\frac{1}{P(X)} \right)$



- for 2 outputs: \log_2

Heads 80% $\Rightarrow \log_2 \left(\frac{1}{0,8} \right) = 0,15$ Surprise

Tails 10% $\Rightarrow \log_2 \left(\frac{1}{0,1} \right) = 3,32$

Total surprise = \sum surprises

$$\underbrace{(0,8 \cdot 100)}_{100 \text{ tosses heads}} \cdot 0,15 + (0,1 \cdot 100) \cdot 3,32 = 4,67 \quad \text{for 100 coin flips}$$
$$\Rightarrow 0,47 \text{ surprise/coin flip} = \text{Entropy} = E(\text{Surprise})$$

$$E(\text{surprise}) = 0,8 \cdot 0,15 + 0,1 \cdot 3,32 = 0,47 = \text{avg. surprise to expect}$$

$$\text{Entropy} = \sum \log \left(\frac{1}{p(x)} \right) p(x) = - \sum p(x) \log(p(x))$$

\uparrow surprise \uparrow probability of surprise

Example:

Area A: 6/7 orange chickens

$$\Rightarrow \text{Entropy} = \frac{6}{7} \cdot \log_2 \left(\frac{1}{\frac{6}{7}} \right) + \frac{1}{7} \cdot \log_2 \left(\frac{1}{\frac{1}{7}} \right) = 0,59$$

Area C: (equal number of chickens)

$$\text{Entropy} = \frac{7}{14} \cdot \log_2 \left(\frac{1}{\frac{7}{14}} \right) + \frac{7}{14} \cdot \log_2 \left(\frac{1}{\frac{7}{14}} \right) = 1$$

Gradient Boost

- continuous values - Gradient Boost for Regression
- makes single leaf \rightarrow first guess is avg. value
- create tree with 8-32 leaves (Fixed size, multiple error-based trees)
- stop when additional trees fail to improve the fit

- Building Model

- start with avg. weight leaf (everyone is 71,2 kg)

- second tree:
 - Error = Observed weight - Predicted weight = Pseudo residual column
 - Use all features to build tree, predicting the residuals (restrict leaves, use avg. of multiple samples)
 - Start with initial 71,2 prediction; run through the tree and sum them
 - Prevent overfit with low bias/high variance: avg. weight + learning rate \times tree predicted = 71,2 + 0,1 \cdot 16,8 = 72,9

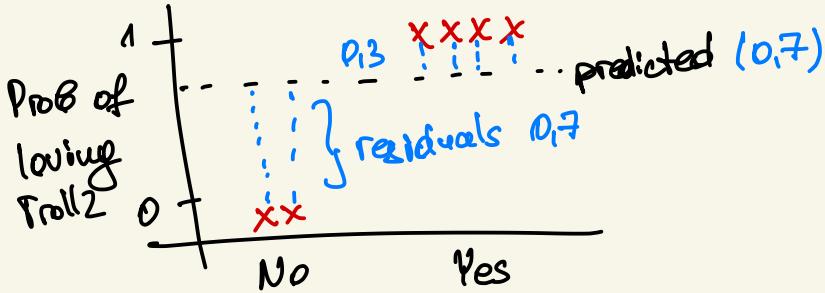
\hookrightarrow scales between 0-1
taking smaller steps

- recalculate residuals before building next tree
- build tree 3 to predict new residuals (#leaves can be different)
- scale by 0,1 again, residuals, predict residuals on next run again

- make trees until residual improvement is too small
- predict by running through 1st leaf, and all following trees

Gradient Boost for Classification

- start with single leaf: $\log(\text{odds}) \rightarrow \log\left(\frac{4}{2}\right) = 0.7$
- convert to probability: $\frac{e^{\log(\text{odds})}}{1 + e^{\log(\text{odds})}} = 0.7$ 4 people love / 2 people don't
- calculate pseudo residuals ($\text{Residual} = \text{Observed} - \text{Predicted}$)



Loving?

- Build a tree using all other features + limit number of leaves to 8-32
- When multiple values reside in a leaf: $\sum_i \text{residual}_i \rightarrow$ all residuals in the leaf

For 2 residuals in 1 leaf:

$$\frac{0.3 + (-0.7)}{0.7 \cdot (1-0.7) + 0.7 \cdot (1-0.7)} = -1$$

$$\sum_i \text{Previous Prob}_i \times (1 - \text{Prev Prob}_i)$$

- Scale new tree by a learning rate (~ 0.1)

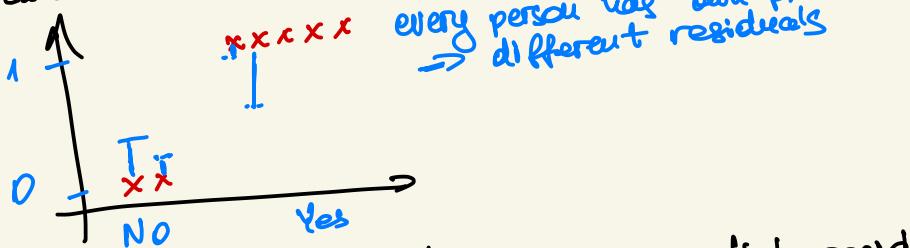
$$0.7 + 0.1 \cdot 1.4 = 1.8 \quad \text{odds + prediction}$$

initial leaf $\rightarrow 1^c$

$$\hookrightarrow \text{prob.} = \frac{e^{1.8}}{1 + e^{1.8}} = 0.9 \rightarrow \text{save in a column}$$

$$\frac{-0.7 \rightarrow \text{current leaf}}{\sum 0.7 \cdot (1-0.7) \rightarrow \text{prob from initial leaf}} = -3.3 \rightarrow \text{output for leaf}$$

- calculate new residuals = Observed - Predicted



- save residuals and build new tree, predict residuals

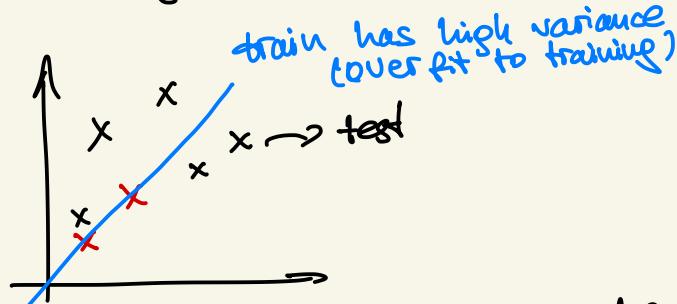
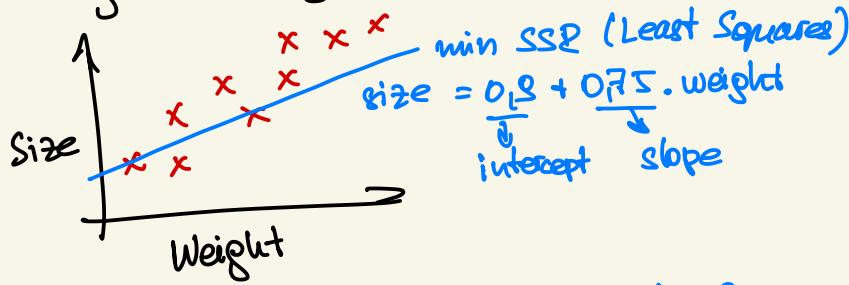
- take last predicted prob. from the chart for each sample

- when predicting: start at initial leaf and follow trees

- stop training when improvement diminish

Regularization

- Ridge (L2) Regression



introduce bias to drop variance (better long-term predictions)

- Ridge regression minimizes $SSR + \lambda \cdot \text{slope}^2 \rightarrow \text{penalty}$

Example

$$\text{size} = 0,4 + 1,3 \cdot \text{weight} = \text{perfect fit}$$

$$\Rightarrow SSR = 0$$

$$\lambda = 1, \text{slope} = 1,3$$

$$\Rightarrow 0 + 1 \cdot 1,3^2 = 1,69 \text{ for } SSR \text{ line}$$

- minimize penalty to create Ridge Regression Line
 - λ determines how small the slope will be
 - 10-fold cross validation to determine λ
 - in log regression: sum of likelihoods + $\lambda \cdot \text{slope}^2$
 - every parameter is scaled by λ
 - more parameters:
- Size = $y\text{-intercept} + \text{slope} \cdot 1 \cdot \text{Weight} + \text{slope} \cdot 2 \cdot \text{Age}$
 ridge regression can solve for all parameters with cross validation

- Lasso (L1) Regression

- $SSR + \lambda \cdot |\text{slope}| \quad \lambda \in [0; +\infty)$ determined by cross-validation

$$\text{size} = y\text{-intercept} + \text{slope} \cdot \text{Weight} + \text{slope} \cdot \text{diet difference} \cdot \text{HighFatDiet}$$

$$\Rightarrow \text{minimize: } SSR + \lambda \times (|\text{slope}| + |\text{diet diff}|)$$

- lasso regression can shrink slope to 0 (ridge only approaches 0)
- eliminates variables without impact

- Elastic Net Regression

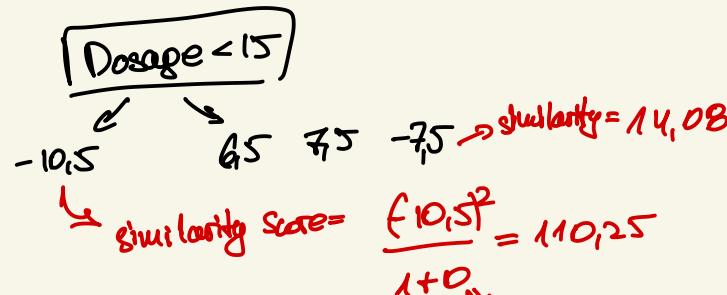
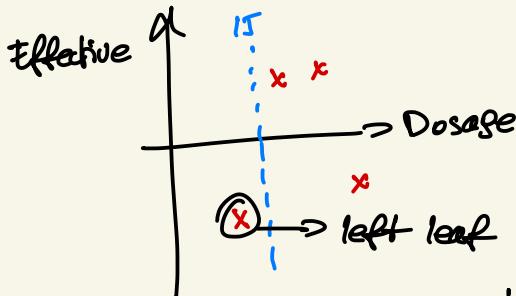
- combines Ridge and Lasso
- $SSR + \lambda_1 (|\text{var1}| + |\text{var2}| \dots) + \lambda_2 (\text{var}_1^2 + \text{var}_2^2 \dots)$

- use cross validation to find λ_1/λ_2
- good if there is correlation between variables

XGBoost Regression

- make initial prediction $\rightarrow 0,5$
- calculate residuals (Observed - Predicted)
- fit XGBoost Tree to the Residuals
- start tree with single leaf (all residuals here)
- calculate Similarity Score = $\frac{(\sum \text{residuals})^2}{\# \text{residuals} + \lambda}$

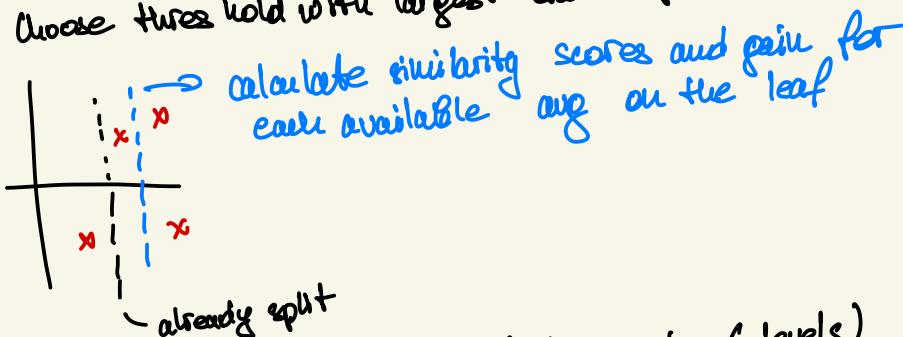
$$\text{regularization param} = \frac{(-10,5 + 6,5 + 7,5 - 7,5)^2}{4 + 0} = 4$$



- split on avg of first 2 dosages
- similar residuals = large score
- How much better leaves cluster similar residuals than root?

$$\text{Grain} = \text{Left similarity} + \text{Right similarity} - \text{Root similarity} = 120,33$$

- shift threshold to between 2/3 calculate similarity / grain (for all remaining datapoints)
- choose threshold with largest Grain for the root



- repeat for other leaves (allow up to 6 levels)

- Pruning XGBoost

- pick number (130) = δ (user defined for pruning)
- (lowest Branch Grain) δ if $< 0 \Rightarrow$ remove Branch

$$140,17 - 130 = 10,17 \Rightarrow \text{Branch remains}$$

- perform for upper Branches as long as last Branch was pruned

- Setting λ
 - set λ to 1, recalculate similarity, similarities decrease
decrease in similarity $\propto \frac{1}{\# \text{residuals}}$
 - larger λ leads to smaller Gain, more pruning
 - if Gain is negative, pruning will be executed
 - Output value for a leaf = sum of residuals
 - When $\lambda > 0$, it reduces influence of $\frac{\# \text{residuals} + \lambda}{\text{observation}}$ on the output
(sensitivity to individual observation)
- Predictions
 - start with initial prediction + learning rate \cdot Output of tree
 - dosage $0,5 + 0,3 \cdot (-10,5) = -2,65$ $\hookrightarrow \epsilon = 0,3$ default
 - takes small steps
 - Build new tree based on the new residuals
 - continue until residuals small, or max # trees is reached
- XGBoost Classification
 - Prob effective vs Dosage
 - make initial prediction (default = 0,5) 50% chance of effective
 - fit tree to the residuals $\frac{(\sum \text{Residual}_i)^2}{\text{similarity score}} = \frac{(\sum \text{Residual}_i)^2}{\sum [\text{PrevProb}_i \cdot (1 - \text{PrevProb}_i)] + \lambda} \rightarrow \text{regularization}$
 - all residuals to the leaf, calculate similarity for the leaf
 - try splitting values into 2 leaves, calculate similarity score for each leaf
 - Calculate Gain = $\text{Left}_{\text{simil}} + \text{Right}_{\text{simil}} - \text{Root}_{\text{simil}}$
 - if no other threshold gives us higher gain \rightarrow choose this
 - Build until reaching max level
 - min number of residuals is determined by Cover
 - for classification: $\text{Cover} = \sum [\text{PrevProb}_i \cdot (1 - \text{PrevProb}_i)] + \lambda$
 - for regression: $\text{Cover} = \text{Number of Residuals}$
 - min value for Cover = 1

• if Cover of the leaf < 1 \Rightarrow remove leaf

- Pruning

- pick δ
- if Gain - $\delta < 0 \Rightarrow$ prune
- $\lambda > 0$ reduce the sensitivity to individual observations

- Convert into probability

$$\frac{P}{1-P} = \text{odds} \quad P=0.5 \Rightarrow \log(\text{odds})=0$$

$$\log\left(\frac{P}{1-P}\right) = \log(\text{odds})$$

• Output of a leaf Output =
$$\frac{\left(\sum \text{Residual}_i\right)}{\sum \left[\text{PrevProb}_i \cdot (1-\text{PrevProb}_i)\right] + \lambda}$$

• Initial prediction: $\log\left(\frac{P}{1-P}\right)$

• Initial prediction + learning rate \cdot Tree output value = $\log(\text{odds})$

$$= \sum \rightarrow \text{Prob} = \frac{e^{\log(\text{odds})}}{1 + e^{\log(\text{odds})}}$$