

Decision Trees + Ensemble Methods

Decision Trees are good to use with ensembles. Tree ensembles

Decision Trees

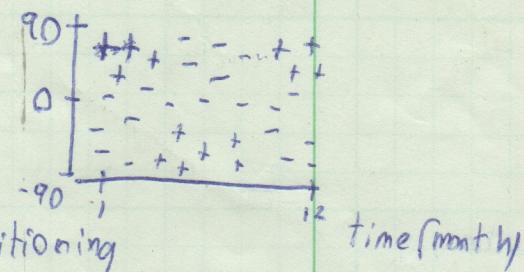
Non-linear (vs. SVMs)

Ski example - linear model bad

④ ⑤ Picking Regions - Split into regions. How?

⑥ Greedy, Top-Down, Recursive ~~partitioning~~

lat($^{\circ}$) win many Kaggle competitions



- 20 questions, draw out tree, branch vs. leaf

⑦ - Formally, given node P (parent node), ~~with~~ covering region R_P , ~~we can define~~ a split s_P on j -th feature with threshold t as

$$s_P(j, t) = \left(\underbrace{\{X | X_j < t, X \in R_P\}}_{R_1}, \underbrace{\{X | X_j \geq t, X \in R_P\}}_{R_2} \right)$$

How to choose splits?

- Intuitively, trying to isolate + and -

- Useful to define $L(R)$ on region

- For now define $L(R)$ as misclassification loss

⑧ Given C total classes, define \hat{p}_c as proportion of examples in R that are of class c

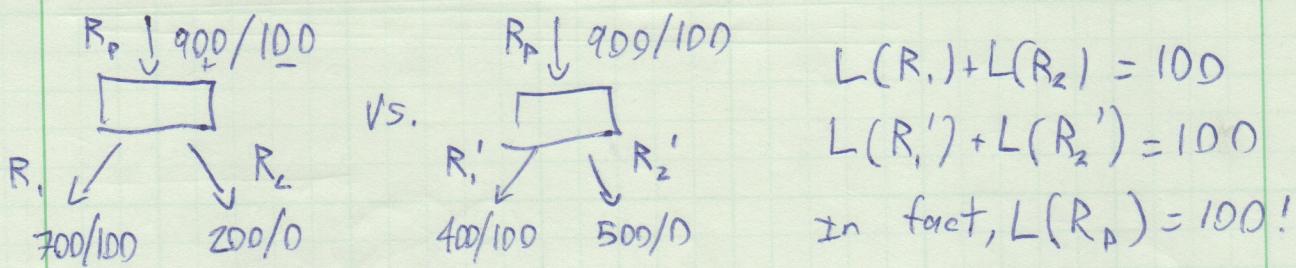
$$L_{\text{misclass}} = 1 - \max_c (\hat{p}_c)$$

⑨ Loss of parent node is : $L(R_p)$

of children : $L(R_1) + L(R_2)$

⑩ Want to maximize reduction in loss ($L(R_p) - [L(R_1) + L(R_2)]$)
 ↳ choose j, t for s_p to do $\Rightarrow 0$

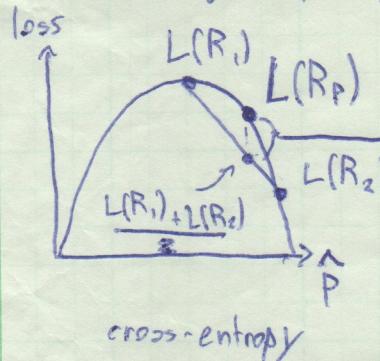
④ Misclassification Loss Has Problems



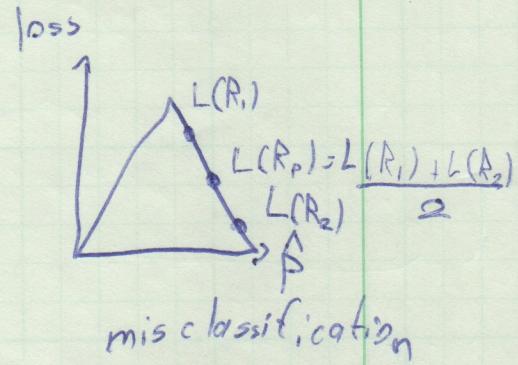
⑤ Instead, define cross-entropy loss

$$L_{\text{cross}} = - \sum_c \hat{p}_c \log_2 \hat{p}_c \quad (\hat{p} \log_2 \hat{p} = 0 \text{ if } \hat{p} = 0)$$

Number of bits needed to specify outcome given distribution is known.



Change in loss
(Information Gain in cross-entropy terms)



Gini loss has similar shape $L_{\text{gini}} = \sum_c \hat{p}_c (1 - \hat{p}_c)$

Now, let's cover some extensions of DTs...

⑥ Regression Trees

Instead of predicting majority class, predict mean of values in \$R_m\$: $\hat{y}_m = \frac{\sum_{i \in R_m} y_i}{|R_m|}$

⑦ Minimize squared loss

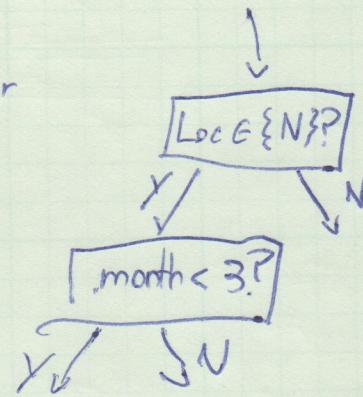
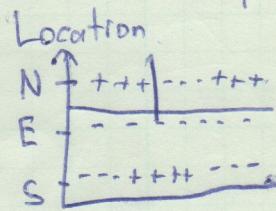
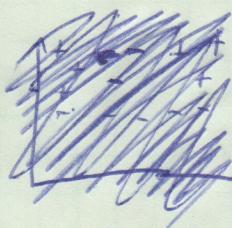
$$L_{\text{squared}} = \frac{\sum_{i \in R_m} (y_i - \hat{y}_m)^2}{|R_m|}$$

ski example snowfall

10.9	0.0	-11.8
0.0	0.0	12.0
0.0	0.0	0.0
0.0	7.6	0.0
0.0	15.0	0.0

④ Categorical Variables

North, South Hemisphere, Equator



- For q categories, 2^q splits possible so does not scale well, except for binary classification case.

⑤ Regularization

- Decision trees if fully grown are high variance, low bias.
Use heuristics for regularization

1) Min leaf size

2) Max depth

3) Max # nodes

4) Min decrease in loss \rightarrow dangerous since might be higher order interactions!

5) Pruning \rightarrow use validation set, measure misclassification or squared error

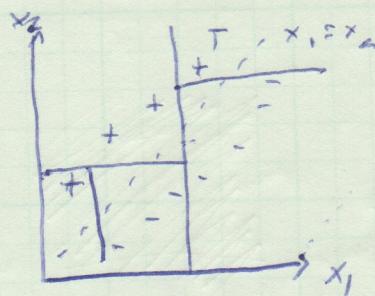
⑥ Runtime - Binary Classification

n examples, f features, d depth

Test time is $O(d)$, if balanced $O(\log n)$

Train time, each point in $O(d)$ nodes, costs $O(f)$ at each node
So $O(nfd)$ n is size of design matrix and d often $\log n$

④⑤ No additive structure



Linear models on other hand are great at additive structure!

⑥⑦ Recap

- | | |
|-------------------|---------------------------|
| ⊕ Easy to explain | ⊖ High Variance |
| ⊕ Interpretable | ⊖ Bad at Additive |
| ⊕ Categorical Var | ⊖ Low Predictive Accuracy |
| ⊕ Fast | |

Will solve these issues via ... ensembling!

⑧⑨⑩ Ensembling

Take independent, identically distributed (iid) random variables (RV)

X_i . Say $\text{Var}(X_i) = \sigma^2$

$$\text{Var}(\bar{X}) = \text{Var}\left(\frac{1}{n} \sum X_i\right) = \frac{\sigma^2}{n}$$

Now drop independence assumption, so only iid.

Say X_i 's are correlated by ρ

$$\text{Then } \text{Var}(\bar{X}) = \rho \sigma^2 + \frac{1-\rho}{n} \sigma^2 \quad \text{as } n \rightarrow \infty \quad \text{Var}(\bar{X}) = \rho \sigma^2$$

⑪ Ways to ensemble

- 1) Use different algorithms
- 2) Use different training sets
- 3) Bagging (RF)
- 4) Boosting (Adaboost, xgboost)

⊗⊗*

Bagging

Stands for bootstrap aggregation, variance reduction method

* Bootstrap - method from statistics to measure uncertainty

* Have true population P , which training set S was drawn from. Write ~~$S \sim P$~~ $S \sim P$.

Ideally, have S_1, S_2, \dots all drawn from P .

Now, assume $S = P$.

* Can now draw new samples Z from S ! $Z \sim S$
(sample w/ replacement, $|Z| = |S|$) Z_1, Z_2, \dots, Z_M

* Train model G_m on Z_m . Can then look at variability in predictions

However, can also define new aggregate predictor

$$\tilde{G}_M(x) = \sum_m \frac{G_m(x)}{M}$$

Bias-Variance Analysis

* Recall $\text{Var}(\bar{X}) = p\sigma^2 + \frac{1-p}{M}\sigma^2$ for M predictors.

* Bagging creates less correlated predictors, $\downarrow p$, thereby reducing Variance.

Note that Bias increases a bit due to reduced sampling), but $\downarrow \text{Var}$ typically outweighs $\uparrow \text{Bias}$

Note too that $\uparrow M$ can't ~~hurt~~, since just decrease Var
make overfit
more

④ Decision Trees + Bagging

- ① Recall DT are high variance, low bias \rightarrow ideal fit for variance reduction of bagging!

Note bagging has another benefit for DT:

If a feature is missing, simply don't use trees in ensemble that contain that feature!

- ② Missing feat \rightarrow Ignore Gm's splitting on it!

③ ~~Out-of-bag estimation~~ Out-of-bag estimation

Additional benefit of bagging: ~~free validation set.~~

On average, Z will contain $2/3$ of S . Use other $1/3$ to estimate error, call OOB error. In limit as $M \rightarrow \infty$ OOB gives equivalent results to LOOCV.



④ Variable Importance Measure

Do lose some interpretability.

- ⑤ For each feature, find each split that uses it in the ensemble. Measure decrease in loss, average.

Note doesn't measure degradation in perf if didn't have feature, since other features might be able to substitute.

⑥ Random Forest

Intuition: one very strong predictor \rightarrow correlated trees

~~Bagging + Feature selection~~

- ⑦ ~~Bagging~~. Instead, at each split only allow subset of features to be used (ex: $u = \sqrt{f}$). Decreases Var, slight increase in Bias

- ⑧ Again, $\text{Var}(\bar{x}) = p\sigma^2 + \frac{1-p}{M}\sigma^2$ as $\downarrow p \rightarrow \downarrow \text{Var}$
 Just what RF does! Works even better for missing values.. (assuming M large)

④ Recap - Bagging

~~Bagging is variance reducing, random forest is bias reducing~~

- ⊕ ↓ Var (even more so for RF)
- ⊖ ↑ bias (even more so for RF)
- ⊕ Better accuracy
- ⊖ Harder to interpret
- ⊕ Deal with missing values
- ⊖ Still not additive
- ⊖ More expensive

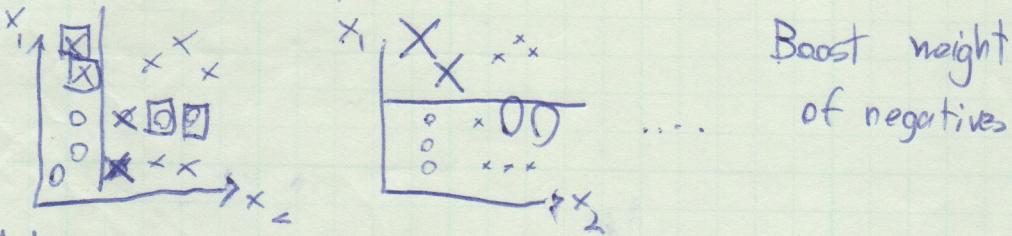
④ Boosting

- ⊗ Bagging was variance reducing, boosting bias reducing
- ~~Boosting is additive~~

Since reducing bias, want high bias, low variance models, weak learners

In terms of DT, just use decision stumps

④ Intuition via Example



④ Adaboost

Training data $(x_1, y_1), \dots, (x_N, y_N)$

1) Set $w_i = 1/N$ for $i = 1 \dots N$

2) for $m = 1 \dots M$:

a) Fit weak classifier G_m to weighted training data

b) Compute weighted err_m = $\frac{\sum w_i I(y_i \neq G_m(x_i))}{\sum w_i}$

c) Compute weight $\alpha_m = \log(\frac{1 - \text{err}_m}{\text{err}_m})$

d) set $w_i = w_i \cdot \exp[\alpha_m I(y_i \neq G_m(x_i))]$

3) Return $f(x) = \text{sign} \left[\sum_m \alpha_m G_m \right]$

④ Adaboost - cont'd

Additive - combine together prediction of many ^{weak} models.
 Oftentimes beats single strong model.
 No longer independent of previous models in sequence
 ↳ can overfit.

⑤ More general framework: Forward Stagewise Additive Modeling (FSAM)

Though derived on its own, Adaboost can be thought of as special case of FSAM.

Still want to create ^{ensemble} classifier f_{es}

1) Initialize $f_0(x) = 0$

2) For $m=1 \dots M$

a) Compute $(\beta_m, \gamma_m) = \underset{\beta, \gamma}{\operatorname{arg\,min}} \sum_{i=1}^N L(y_i, f_{m-1}(x_i) + \beta b(x_i; \gamma))$

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

3) Return $f(x) = f_M(x)$

⑥ Adaboost is FSAM with exponential loss and 2-class classification

$$L(y, f(x)) = \exp(-y f(x)). \text{ Proof in ESL.}$$

⑦ For squared loss for regression, FSAM, is same as fitting

$$L(y, f(x)) = (y - f(x))^2 \quad \text{individual tree to residual } y_i - f_{m-1}(x_i)$$

⑧ Gradient Boosting

For more general losses, can't always write out nice closed-form solution to minimization problem. Turn to numerical optimization
 One way:

Take derivative, do gradient descent.

But ~~can~~ restricted to taking steps that are in model class.

basis of
xgboost

Compute gradient at each training point i ; wrt current predictor f

$$\textcircled{4} \quad g_i := -\frac{\partial L(y_i, f(x_i))}{\partial f(x_i)}$$

New predictor simply tries to train ~~regression~~^{regression} tree to match this gradient. In FSAN:

$$\textcircled{5} \quad \gamma = \operatorname{argmin}_g \sum_{i=1}^N (g_i - b(x_i; \gamma))^2. \text{ Add in, is your gradient step.}$$



Recap - Boosting

⊕ ↓ bias

⊕ Very good accuracy

⊖ ↑ var.

⊖ Prone to overfitting