

Lecture 17 - Mar 6

Gradient Boosting
XGBoost

References

- [*Elements of Statistical Learning*](#)
 - Ch 10: Boosting and Additive Trees
- T. Chen and C. Guestrin. 2016. XGBoost: A Scalable Tree Boosting System. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD '16)*, 785–794. <https://arxiv.org/abs/1603.02754>
- <https://medium.com/syncedreview/tree-boosting-with-xgboost-why-does-xgboost-win-every-machine-learning-competition-ca8034c0b283>

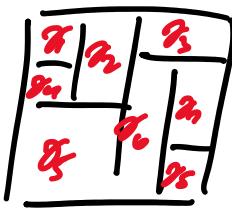
Upcoming Deadlines

Homework 4 (Mar 15)

Boosting Trees

Trees split space into R_1, R_2, \dots, R_J such that prediction is

$$x \in R_j \Rightarrow f(x) = \gamma_i$$



A tree can be expressed

$$T(x; \Theta) = \sum_{j=1}^J \gamma_j \mathbb{1}_{\{x \in R_j\}}$$

with parameters $\Theta = \{R_j, \gamma_j\}_1^J$ found by minimizing the empirical risk

$$\hat{\Theta} = \arg \min_{\Theta} \sum_{j=1}^J \sum_{x \in R_j} L(y_i, \gamma_i)$$

Difficult... so we usually go for suboptimal solutions by

① Find γ_j given R_j : easy, \bar{y}_j or model class

② Find R_j : hard, approximate. Typical strategy is to use a greedy, recursive partitioning. Sometimes, estimate $\hat{\Theta}$ with a more convenient criterion:

$$\hat{\Theta} = \arg \min_{\Theta} \sum_{i=1}^N \tilde{L}(y_i, T(x_i, \Theta))$$

Then estimate γ_j more precisely.

... , ... see in 9.2 for decision tree classifier via Gini

Similar process to 9.2 for decision tree classifier via Gini index instead of loss to identify R_j 's.

Boosted tree model

$$f_m(x) = \sum_{m=1}^M T(x; \Theta_m)$$

Induced by forward-stagewise approach. At each stage, solve

$$\hat{\Theta}_m = \arg \min_{\Theta_m} \sum_{i=1}^N L(y_i, f_{m-1}(x_i) + T(x_i, \Theta_m))$$

Solving is slow,
so we need
approximations
frequently

for the region set and constants $\Theta_m = [R_{jm}, \theta_{jm}]^{J_m}$
of the next tree.

Given R_{jm} , finding optimal θ_{jm} is easy

$$\hat{\theta}_{jm} = \arg \min_{\theta_{jm}} \sum_{x_i \in R_{jm}} L(y_i, f_{m-1}(x_i) + \theta_{jm})$$

Two-class classification + exp. loss \rightarrow AdaBoost

Otherwise,

$$\hat{\Theta}_m = \arg \min \sum_{i=1}^N w_i^{(n)} e^{-y_i f_{m-1}(x_i)}$$

$$\hat{\theta}_m = \arg \min_{\theta_m} \sum_{i=1}^N w_i^{(n)} \exp(-y_i T(x_i, \theta_m))$$

A decision tree algorithm can use exp loss as splitting criterion

$$\Rightarrow \hat{\theta}_{jm} = \ln \left(\frac{\sum_{x_i \in R_{jm}} w_i^{(n)} \mathbb{1}_{\{y_i=1\}}}{\sum_{x_i \in R_{jm}} w_i^{(n)} \mathbb{1}_{\{y_i=-1\}}} \right)$$

This requires a specialized tree-growing algorithm, but we prefer weighted least squares regression tree.

Fast approximate algorithms for solving

$$\hat{\theta}_m = \arg \min_{\theta_m} \sum_{i=1}^N L(y_i, f_{m-1}(x_i) + T(x_i, \theta_m))$$

With any differentiable loss can be derived similar to numerical optimization by minimizing

$$L(f) = \sum_{i=1}^N L(y_i, f(x_i))$$

(with respect to f of the form $f(x) = \sum_{m=1}^M T(x, \theta_m)$, e.g.

$$\hat{f} = \arg \min_f L(f)$$

where the "parameters" $f \in \mathbb{R}^N$ are $f = \{f(x_1), \dots, f(x_N)\}$.

Numerical optimization solves as

$$f_m = \sum_{n=0}^M h_m, \quad h_m \in \mathbb{R}^N,$$

where $f_0 = h_0$ is an initial guess and each f_m is induced based on current vector f_{m-1} , the sum of previous updates. Different procedures compute "steps" h_m differently.

Different procedures compute "step" "in" "one" /

10.10.1 Steepest Descent

$h_m = -\gamma_m g_m$ where $\gamma_m \in \mathbb{R}$, $g_m \in \mathbb{R}^N$ is the gradient of $L(f)$ evaluated at $f = f_{m-1}$

$$\gamma_m = \left[\frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right]_{f(x_i) = f_{m-1}(x_i)}$$

and step length is $\gamma_m = \arg \min_s L(f_{m-1} - s g_m)$, then

$$f_m = f_{m-1} - \gamma_m g_m$$

very greedy since $L(f)$ is most rapidly decreasing at $f = f_{m-1}$.

10.10.3 Implementations of Gradient Boosting

Algorithm 10.7

generic gradient tree-boosting algorithm
for regression

- ① Initialize $f_0(x) = \arg \min_f \sum_{i=1}^n L(y_i, f)$ ← initialize a single terminal node tree

- specific algorithms
specify the loss
- terminal node tree
- ② For $m=1$ to M
- (a) For $i=1, \dots, N$,
- $$r_{im} = - \left[\frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right]_{f=f_{m-1}}$$
- (b) Fit a regression tree to targets r_{im} giving $R_{jm}, j=1, \dots, J_m$
- (c) For $j=1, \dots, J_m$
- $$\gamma_{jm} = \arg \min_{\gamma} \sum_{x_i \in R_{jm}} L(y_i, f_{m-1}(x_i) + \gamma)$$
- (d) Update $f_m(x) = f_{m-1}(x) + \sum_{j=1}^{J_m} \gamma_{jm} \mathbb{1}_{x \in R_{jm}}$
- ③ Output $\hat{f}(x) = f_M(x)$

Gradient descent to minimize $L(x)$

$$x_{i+1} = x_i - \alpha \nabla L(x_i)$$

"go in the steepest downhill direction"

where α is a fixed constant

Improving gradient descent

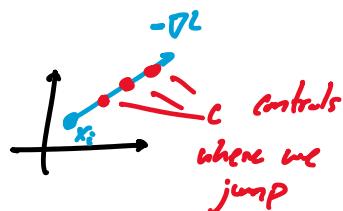
Line search: $x_{i+1} = x_i - c \nabla L(x_i)$

this is how
normal gradient
boosting works

choose c minimizing $L(x_i - c \nabla L(x_i))$

to pick a distance to jump

easy if f is convex and cheap to compute



If L is infinitely differentiable, the Taylor expansion of L at $x + \Delta x$ is

$$\sum_{i=1}^{\infty} \frac{\alpha^i}{i!} f^{(i)}(x)$$

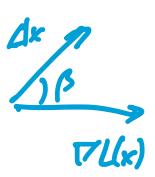
$$L(x + \alpha \Delta x) = L(x) + \underline{\alpha \Delta x^T \nabla L(x)} + \underline{\frac{\alpha^2}{2!} \|\Delta x\|_2^2 \nabla^2 L(x)} + \dots$$

Gradient descent: $L(x + \alpha \Delta x) \approx L(x) + \alpha \Delta x^T \nabla L(x)$

minimize w.r.t. Δx ... we need $\Delta x^T \nabla L(x)$ to be minimized,
 $\underline{\Delta x^T \nabla L(x)} \leq 1$

minimize w.r.t. Δx ... we need

$$-1 \leq \text{Cos} \beta = \frac{\Delta x^T \nabla L(x)}{\|\Delta x\| \|\nabla L(x)\|} \leq 1$$



$$-\|\Delta x\| \|\nabla L(x)\| = \Delta x^T \nabla L(x)$$

$$\beta = \pi \text{ here, so } \Delta x^T = -\nabla L(x)$$

move opposite direction to gradient

$$L(x + \Delta x) \approx L(x) + \Delta x^T \nabla L(x) + \frac{1}{2} \|\Delta x\|_2^2 \nabla^2 L(x)$$

Newton's method:

minimize right side w.r.t. Δx

$$0 = \frac{\partial}{\partial \Delta x} \left(L(x) + \Delta x^T \nabla L(x) + \frac{1}{2} \|\Delta x\|^2 \nabla^2 L(x) \right)$$

$$0 = \nabla L(x) + \Delta x^T \nabla^2 L(x)$$

$$-\nabla L(x) = \Delta x^T \nabla^2 L(x)$$

$$-\nabla L(x) [\nabla^2 L(x)]^{-1} = \Delta x^T \nabla^2 L(x) [\nabla^2 L(x)]^{-1}$$

$$\Rightarrow x_{i+1} = x_i - \nabla L(x) [\nabla^2 L(x)]^{-1}$$

$$Z = \{(x_1, y_1), \dots, (x_n, y_n)\}, \quad x_i \in \mathbb{R}^d, y_i \in \mathbb{R}$$

Tree ensemble

$$\hat{y}_i = f(x_i) = \sum_{k=1}^K T_k(x_i; \theta_k)$$

trees with weights w
(e.g. CART)

Let T = number of leaves on tree

Regularized objective

$$L(f) = \sum_{i=1}^n L(y_i, f(x_i)) + \sum_{k=1}^K \mathcal{S}(T_k)$$

differentiable + convex loss function

$$\text{where } \mathcal{S}(T) = \gamma T + \frac{1}{2} \lambda \|w\|^2$$

γ hyperparam
 λ penaltie
high trees

The model is trained additively, so we have iterations

where we minimize

$$\text{Loss}^t = \sum_{i=1}^n L(y_i, f_{t-1}(x_i) + T_t(x_i; \theta_t)) + \mathcal{S}(T_t)$$

Taylor expansion at $f + f - f = T$ as product

Second-order approximation

$$\text{Loss}^t \approx \sum_{i=1}^n \left(L(y_i, f_{t-1}(x_i)) + \underbrace{\frac{\partial L(y_i, f_{t-1}(x_i))}{\partial f_{t-1}(x_i)} T_t(x_i; \theta_t) + \frac{1}{2} \frac{\partial^2 L(y_i, f_{t-1}(x_i))}{\partial f_{t-1}(x_i)^2} T_t(x_i; \theta_t)^2}_{\text{Taylor's theorem... Newton's method}} \right) + \mathcal{S}(T_t)$$

$\underbrace{\text{regularization term}}$

If we do $\arg \min_{\theta_t} \text{Loss}^{(t)}$, this is independent of $L(y, f_{t-1}(x))$, so

We can solve

$$\underset{\theta_t}{\operatorname{argmin}} \sum_{i=1}^n \left[\frac{\partial L(y_i, f_{t-1}(x_i))}{\partial f_{t-1}(x_i)} T_t(x_i; \theta_t) + \frac{1}{2} \frac{\partial^2 L(y_i, f_{t-1}(x_i))}{\partial f_{t-1}(x_i)^2} T_t(x_i; \theta_t)^2 \right] + \underline{S(T)}$$

$$S(T) + \frac{1}{2} \sum_{j=1}^T w_j^2$$

$$\underset{\theta_t}{\operatorname{argmin}} \sum_{j=1}^T \left[\sum_{i \in R_j} \frac{\partial L(y_i, f_{t-1}(x_i))}{\partial f_{t-1}(x_i)} \right] w_j + \frac{1}{2} \left(\sum_{i \in R_j} \frac{\partial^2 L(y_i, f_{t-1}(x_i))}{\partial f_{t-1}(x_i)^2} + \lambda \right) w_j^2 \right] + S(T)$$

$$\frac{\partial g}{\partial w_j} = \sum_{j=1}^T \left[\left(\sum_{i \in R_j} \frac{\partial L(y_i, f_{t-1}(x_i))}{\partial f_{t-1}(x_i)} \right) + \left(\sum_{i \in R_j} \frac{\partial^2 L(y_i, f_{t-1}(x_i))}{\partial f_{t-1}(x_i)^2} + \lambda \right) w_j \right] = 0$$

$$w_j^* = - \frac{\sum_{i \in R_j} \frac{\partial L(y_i, f_{t-1}(x_i))}{\partial f_{t-1}(x_i)}}{\sum_{i \in R_j} \frac{\partial^2 L(y_i, f_{t-1}(x_i))}{\partial f_{t-1}(x_i)^2} + \lambda}$$

For a fixed tree structure,

$$\Rightarrow \text{loss}^t = - \sum_{j=1}^T \frac{\left(\sum_{i \in R_j} \frac{\partial L(y_i, f_{t-1}(x_i))}{\partial f_{t-1}(x_i)} \right)^2}{\sum_{i \in R_j} \frac{\partial^2 L(y_i, f_{t-1}(x_i))}{\partial f_{t-1}(x_i)^2} + \lambda} + \frac{1}{2} \left(\sum_{i \in R_j} \frac{\partial^2 L(y_i, f_{t-1}(x_i))}{\partial f_{t-1}(x_i)^2} + \lambda \right) \left(- \frac{\left(\sum_{i \in R_j} \frac{\partial L(y_i, f_{t-1}(x_i))}{\partial f_{t-1}(x_i)} \right)^2}{\sum_{i \in R_j} \frac{\partial^2 L(y_i, f_{t-1}(x_i))}{\partial f_{t-1}(x_i)^2} + \lambda} \right)^2 + S(T)$$

$$\text{loss}^m = \frac{-1}{2} \sum_{j=1}^{|T_m|} \frac{\left(\sum_{i \in R_j} \frac{\partial L}{\partial f_{t-1}} \right)^2}{\left(\sum_{i \in R_j} \frac{\partial^2 L}{\partial f_{t-1}^2} + \lambda \right)} + \beta |T_m|$$

Scoring function
for quality of
a tree
structure T_t

like impurity but derived from
a wider class of objective
functions

$$r = \gamma_1 r^2 \quad 1 - \frac{\partial L}{\partial f} r^2 \rightarrow$$

a tree
structure T_t

$$Loss^{split} = \frac{1}{2} \left[\frac{\left(\sum_{i \in R_L} \frac{2^L}{2f_{t-1}} \right)^2}{\sum_{i \in R_L} \frac{2^{2L}}{2f_{t-1}^2} + \lambda} + \frac{\left(\sum_{i \in R_R} \frac{2^L}{2f_{t-1}} \right)^2}{\sum_{i \in R_R} \frac{2^{2L}}{2f_{t-1}^2} + \lambda} - \frac{\left(\sum_{i \in L} \frac{2^L}{2f_{t-1}} \right)^2}{\sum_{i \in L} \frac{2^{2L}}{2f_{t-1}^2} + \lambda} \right] - \beta$$

↑
for evaluating split
candidates

Example

```
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sn

from sklearn.datasets import load_breast_cancer
from sklearn.datasets import load_wine
from sklearn.inspection import permutation_importance
from sklearn.metrics import classification_report
from sklearn.metrics import confusion_matrix
from sklearn.model_selection import GridSearchCV
from sklearn.model_selection import train_test_split

from tensorflow.keras.datasets import mnist
from tensorflow.keras.datasets import cifar10

from xgboost import XGBClassifier
```

Breast Cancer Classification

```
cancer = load_breast_cancer()

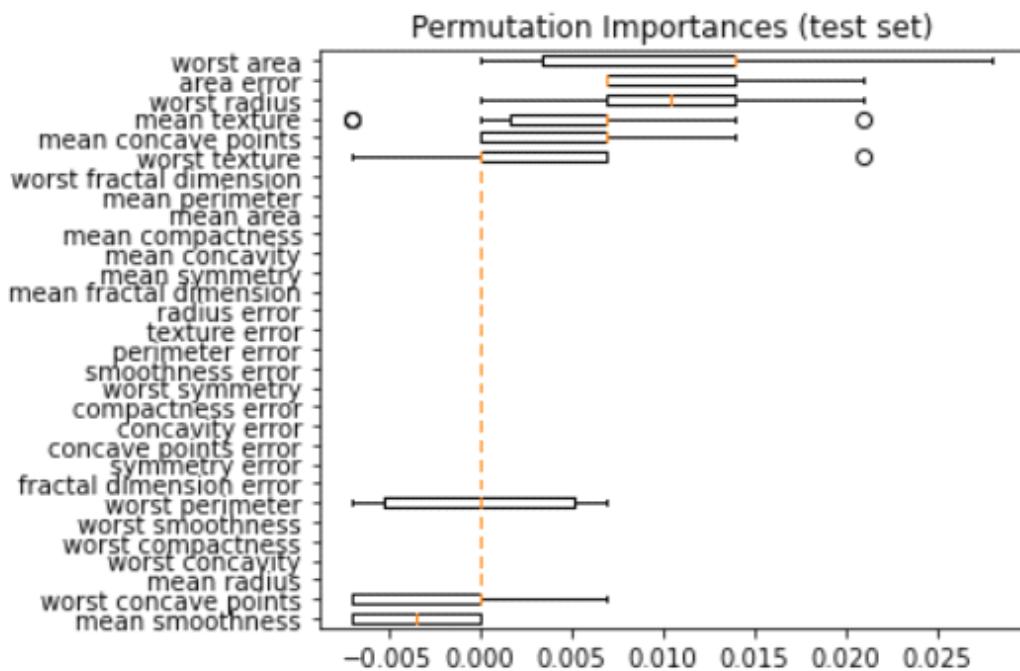
X = cancer['data']
Y = cancer['target']

# split the data into train and test sets
trainX, testX, trainY, testY = train_test_split(X, Y, test_size = 0.25)

# build the classifier
model = XGBClassifier(n_jobs = -1) A
:
result = permutation_importance(model, testX, testY, n_repeats = 10, n_jobs = -1)
sorted_idx = result.importances_mean.argsort()

fig, ax = plt.subplots()
ax.boxplot(result.importances[sorted_idx].T, vert = False,
           labels = cancer.feature_names[sorted_idx])
ax.set_title("Permutation Importances (test set)")
fig.tight_layout()
plt.show()
```

Permutation Importances (test set)



Example

Wine Example

```
: wine = load_wine()

X = wine['data']
Y = wine['target']

# split the data into train and test sets
trainX, testX, trainY, testY = train_test_split(X, Y, test_size = 0.25)

# build the classifier
model = XGBClassifier(n_jobs = -1)
```

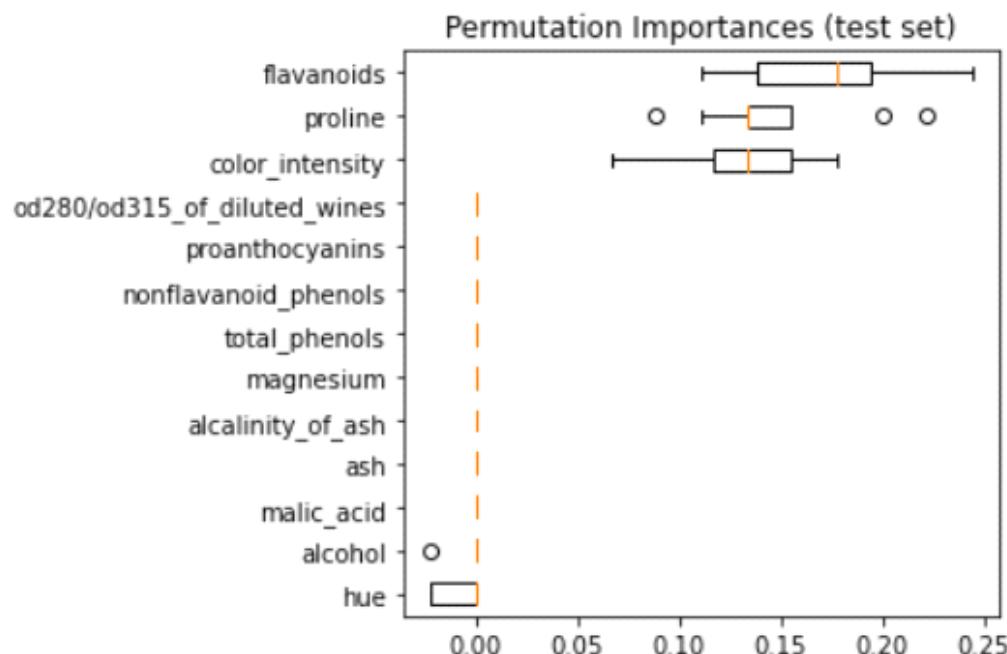
...

```

result = permutation_importance(model, testX, testY, n_repeats = 10, n_jobs = -1)
sorted_idx = result.importances_mean.argsort()

fig, ax = plt.subplots()
ax.boxplot(result.importances[sorted_idx].T, vert = False,
           labels = np.array(wine.feature_names)[sorted_idx])
ax.set_title("Permutation Importances (test set)")
fig.tight_layout()
plt.show()

```



Example

MNIST Example

```
(trainX, trainY), (testX, testY) = mnist.load_data()

# preprocess the data
trainX = trainX.reshape(trainX.shape[0], trainX.shape[1] * trainX.shape[2])
trainX = trainX.astype('float')/255.0

testX = testX.reshape(testX.shape[0], testX.shape[1] * testX.shape[2])
testX = testX.astype('float')/255.0

# build the classifier
model = XGBClassifier(n_jobs = -1)
```

... 100% train, 98% test

Let's try to tune XGBoost with 5-fold cross-validation to see if we can improve the results.

```
parameters = {'max_depth': [1, 5, 10, 100, None], 'reg_lambda': [0.001, 0.01, 0.1, 1, 10]}
model = GridSearchCV(XGBClassifier(n_jobs = -1, eval_metric='mlogloss'), parameters, n_jobs = -1)

# fit the classifier to the training data
model.fit(trainX, trainY)
```

No improvements...

Example

CIFAR10 Example

```
(trainX, trainY), (testX, testY) = cifar10.load_data()

# preprocess the data
trainX = trainX.reshape(trainX.shape[0], trainX.shape[1] * trainX.shape[2] * 3)
trainX = trainX.astype('float')/255.0

testX = testX.reshape(testX.shape[0], testX.shape[1] * testX.shape[2] * 3)
testX = testX.astype('float')/255.0

# build the classifier
model = XGBClassifier(n_jobs = -1)
```

Training Classification Report:

	precision	recall	f1-score	support
0	0.98	0.99	0.98	5000
1	1.00	1.00	1.00	5000
2	0.97	0.98	0.97	5000
3	0.99	0.97	0.98	5000
4	0.96	0.98	0.97	5000
5	0.99	0.98	0.99	5000
6	0.98	0.99	0.99	5000
7	1.00	0.99	0.99	5000
8	0.99	0.99	0.99	5000
9	1.00	0.99	1.00	5000
accuracy			0.99	50000
macro avg	0.99	0.99	0.99	50000
weighted avg	0.99	0.99	0.99	50000

Testing Classification Report:

	precision	recall	f1-score	support
0	0.61	0.62	0.62	1000

0	0.61	0.62	0.62	1000
1	0.64	0.63	0.64	1000
2	0.43	0.39	0.41	1000
3	0.39	0.38	0.39	1000
4	0.48	0.45	0.47	1000
5	0.48	0.48	0.48	1000
6	0.54	0.63	0.58	1000
7	0.60	0.57	0.59	1000
8	0.64	0.67	0.66	1000
9	0.58	0.60	0.59	1000
accuracy			0.54	10000