

Unit 10

Ensemble Methods and Decision Trees

Prof. Phil Schniter



THE OHIO STATE UNIVERSITY

ECE 5307: Introduction to Machine Learning, Sp23

Learning objectives

- Understand intuition behind ensemble methods: “the wisdom of the crowd”
- Understand parallel ensemble methods
 - bagging, pasting
 - random feature selection
- Understand decision trees
 - feature thresholding and decision regions
 - training via top-down tree induction
 - homogeneity metrics: variance reduction, gini impurity
 - ensemble extension: random forests
- Understand boosting, or sequentially trained ensemble methods
 - Adaboost
 - gradient boosting
 - XGBoost

Outline

- Parallel Ensemble Methods: Bagging and Pasting
- Decision Trees and Random Forests
- Boosting: Sequentially Trained Ensemble Methods

The wisdom of the crowd

- Suppose that you want to determine the heavy side of a biased coin
 - Suppose “heads” occurs 70% and “tails” occurs 30%
- To predict the heavy side, you flip it N times and choose the **majority vote**
 - You'll be correct with probability

N	1	11	21	31
prob	70%	92%	97%	99%
 - These values follow from the cdf of the binomial distribution
 - Interpretation: As you use more independent trials, the variance decreases
- Implications for classification:
 - Say you have 31 binary classifiers, each correct only 70% of the time
 - Individual classifiers (“**weak learners**”) do only a bit better than random guessing
 - A majority vote (“**strong learner**”) will classify with 99% accuracy!
 - Caveat: We've assumed the classifiers generate **statistically independent** outcomes
 - As they become more dependent, there will be less improvement from **ensembling**

Voting classifiers

- Given an **ensemble** (i.e., a collection) of base classifiers ...
 - the **hard voting classifier** takes the majority vote of hard decisions
 - the **soft voting classifier** averages their soft outputs (i.e., pmfs) and then chooses the maximizing class ... which usually works better
 - Both are implemented in `sklearn.VotingClassifier`
- The voting classifier will be better than the base classifiers if the base classifiers are sufficiently **diverse** (i.e., sufficiently independent)
- There are different ways to generate **diverse** ensembles of classifiers:
 - 1 Train a common prediction model using **different datasets**
 - 2 Train **different prediction models** using a common dataset:
 - Use different types of model (e.g., logistic regression, SVM, neural net)
 - Use same type but inject randomness (e.g., use random subsets of features)
- Similar ideas can be used to generate an ensemble of base regressors

Bagging and pasting

- Two main ways to introduce data diversity:
 - 1 **pasting**: draw training samples i **without replacement**
 - each sample is used exactly once by one predictor
 - 2 **bootstrap aggregating** (or **bagging**): draw training samples i **with replacement**
 - each sample may be used several times, and some samples never used
- bagging usually works better than pasting
- Bagging:
 - Usually draw n samples per predictor (where n is total # training samples)
 - When n is large, this implies $e^{-1} \approx 37\%$ samples go unused (per predictor)
 - Can use these “**out of bag**” samples for cross-validation!
- Implemented in `sklearn.BaggingClassifier` & `sklearn.BaggingRegressor`:
 - Choose `bootstrap=True` for bagging (or `=False` for pasting)
 - Choose `oob_score=True` to report out-of-bag cross-validation

Random patches and random subspaces

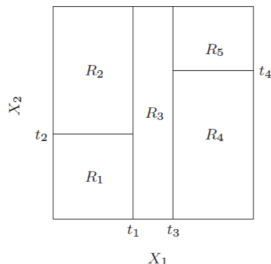
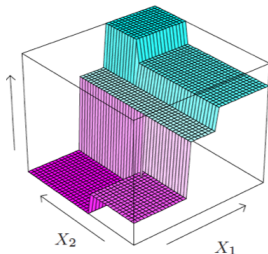
- Recall that using **random subsets of features** is another way to add diversity
 - Useful when features are high-dimensional and redundant
 - Can sample features j with replacement (i.e., “**bootstrap**”) or without
- Terminology:
 - Using both random data i and features j is known as “**random patches**”
 - Using random features j but full data is known as “**random subspaces**”
- Implemented in `sklearn.BaggingClassifier` & `sklearn.BaggingRegressor`:
 - Set `max_features < d` for feature randomization
 - Set `bootstrap_features=True` for bootstrapping

Outline

- Parallel Ensemble Methods: Bagging and Pasting
- Decision Trees and Random Forests
- Boosting: Sequentially Trained Ensemble Methods

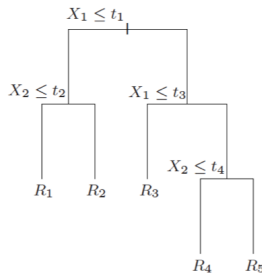
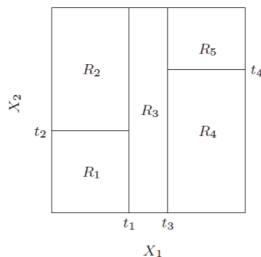
Decision trees

- Consider a supervised learning task (e.g., regression or classification)
 - Suppose data is $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ with $\mathbf{x}_i \in \mathbb{R}^d$ and $y_i \in \mathbb{R}$
- Approach:
 - 1 Partition feature domain \mathbb{R}^d into L regions $\{R_\ell\}_{\ell=1}^L$
 - 2 Output $z_\ell \in \mathbb{R}$ whenever $\mathbf{x} \in R_\ell$
 - For regression, could set z_ℓ at sample mean of $\{y_i\}_{i \in S_\ell}$, where $S_\ell = \{i : \mathbf{x}_i \in R_\ell\}$
 - For classification, could set z_ℓ at sample mode, i.e., most common value in $\{y_i\}_{i \in S_\ell}$
- Regression example with dimension $d = 2$:



Why are they “trees”?

- We normally construct the regions by thresholding one feature at a time!
 - Thus, the decision boundaries are always parallel to the coordinate axes
- Then can view prediction as a **decision tree**:
 - The domain of \mathbf{x} (i.e., \mathbb{R}^d) forms the **root** of the tree
 - The decision regions $\{R_\ell\}$ are the L **leaves** of the tree
- Example with dimension $d = 2$ and $L = 5$ leaves:



Top-down induction of decision trees

Decision trees are trained in a **top-down** manner:

- Start with entire data set: $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ and define index set $S \triangleq \{1, \dots, n\}$
- First split S into subsets (S_1, S_2) , so that labels $\{y_i\}$ are most **homogeneous** (i.e., similar) within each subset
 - Splitting is performed by thresholding some feature j :
 - Ordinal feature: $S_1 = \{i : x_{ij} \leq t\}$ and $S_2 = \{i : x_{ij} > t\}$
 - Categorical feature: $S_1 = \{i : x_{ij} \in \{A, C\}\}$ and $S_2 = \{i : x_{ij} \in \{B, D\}\}$
 - To find “best” split, must search jointly over feature j and threshold t
 - Many ways to quantify homogeneity (will discuss later)
- Then split each subset S_1 and S_2 further, using the same procedure
- Repeat until ...
 - labels are perfectly homogeneous within a subset (e.g., all y_i are same), or
 - stopping condition: e.g., subsets have min # samples, tree has max depth, etc

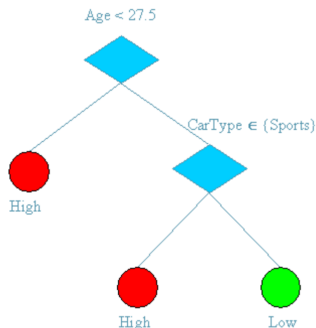
Decision tree example

Example: Risk prediction with $n = 6$ samples, $d = 2$ features, $L = 3$ leaves.
 Feature $j = 1$ (age) is ordinal and feature $j = 2$ (car type) is categorical:

Tid	Age	Car Type	Class
0	23	Family	High
1	17	Sports	High
2	43	Sports	High
3	68	Family	Low
4	32	Truck	Low
5	20	Family	High

Numeric

Categorical



Age=40, CarType=Family \Rightarrow Class=Low

<https://web.fhnw.ch/personenseiten/taoufik.nouri/Data%20Mining/Course/Course3/DM-Part%203.htm>

Homogeneity metrics for regression

- Recall that the goal is to *maximize homogeneity* within S_1 and S_2
 - Equivalently, we want to *minimize inhomogeneity* within S_1 and S_2
- Let us first consider **regression**, where labels $y_i \in \mathbb{R}$
- Inhomogeneity could be measured by the **variance** after splitting $S \rightarrow (S_1, S_2)$:
 - The mean in subset S_ℓ is $\mu_\ell \triangleq \frac{1}{|S_\ell|} \sum_{i \in S_\ell} y_i$, for $\ell \in \{1, 2\}$
 - The variance in subset S_ℓ is $v_\ell \triangleq \frac{1}{|S_\ell|} \sum_{i \in S_\ell} (y_i - \mu_\ell)^2$, for $\ell \in \{1, 2\}$
 - Thus the (average) variance after splitting is $v \triangleq \frac{|S_1|}{|S|} v_1 + \frac{|S_2|}{|S|} v_2$
- Another option is to use **absolute error** $\frac{1}{|S_\ell|} \sum_{i \in S_\ell} |y_i - \mu_\ell|$ instead of v_ℓ , since it is more robust to outliers
- We'll see a more sophisticated metric when discussing XGBoost on page 31

Homogeneity metrics for classification

- Now consider K -ary **classification**, where $y_i \in \{1, \dots, K\} \forall i$
 - For subset S_ℓ , the **empirical label pmf** is $\{p_{\ell k}\}_{k=1}^K$, where $p_{\ell k} = \frac{\sum_{i \in S_\ell} \mathbb{1}_{y_i=k}}{|S_\ell|}$
- Again, want to minimize inhomogeneity within S_1 and S_2
- One popular measure of inhomogeneity is **Gini impurity**:

$$I_G = \frac{|S_1|}{|S|} I_{G,1} + \frac{|S_2|}{|S|} I_{G,2}, \quad \text{where } I_{G,\ell} \triangleq 1 - \sum_{k=1}^K p_{\ell k}^2$$

- Note that $I_{G,\ell} \in [0, 1 - \frac{1}{K}]$
 - $I_{G,\ell} = 0$ means perfectly pure/homogeneous, i.e., constant y_i for $i \in S_\ell$
- Other criteria include **entropy**, **misclassification error**, and **chi-square**
 - See **Hastie, Tibshirani, Friedman, *The Elements of Statistical Learning*, 2nd Ed., 2009**

Top-down induction for generic loss

- Say we want a tree $f(\cdot)$ that minimizes a generic loss $\sum_{i=1}^n L(y_i, f(\mathbf{x}_i))$
- To design the tree, it helps to write $f(\mathbf{x}) = z_{q(\mathbf{x})}$, where
 - z_ℓ is the value assigned to leaf ℓ
 - $q(\cdot)$ is the decision function that maps \mathbf{x} to ℓ
- Thus tree design becomes “ $\arg \min_{\mathbf{z}, q} \sum_{\ell} \sum_{i: q(\mathbf{x}_i) = \ell} L(y_i, z_\ell)$ ”
- Remember that we design the tree one split at a time. Consider splitting $S \rightarrow (S_1, S_2)$. This would cause the loss to go from

$$\min_{z_1} \sum_{i \in S} L(y_i, z_1) \rightarrow \min_{z_1} \sum_{i \in S_1} L(y_i, z_1) + \min_{z_2} \sum_{i \in S_2} L(y_i, z_2)$$

- So, we try all feature/threshold pairs (j, t) , each of which leads to a split (S_1, S_2) . We evaluate the loss for each split and choose the best one.
- If the best split does not improve the loss, we do no splitting.

Advantages and disadvantages of decision trees

Advantages:

- Very **general**
 - Very few assumptions made on data; can work with any dataset
 - Don't need to standardize the features.
- Very **interpretable**: “white box model”
 - Easy to understand how a tree arrives at its prediction
- Prediction is **fast**: $O(\log_2 n)$ decisions to process a test sample x

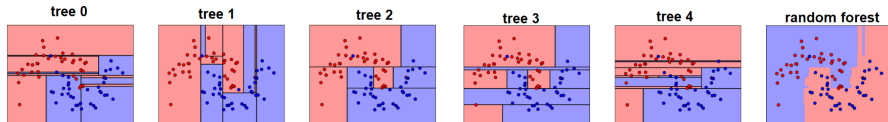
Disadvantages:

- Training can be **expensive** if all features j are considered at every split
 - Solution: Restrict to a few randomly chosen features j at each split
- Prone to **overfitting**
 - Highly dependent on training data: changing one sample can change entire tree!
 - Can regularize by enforcing min # samples per subset, max depth, etc

~> There is a **good overview of decision trees in the sklearn documentation!**

Random forests

- Decision trees tend to **overfit**: low bias & high variance
- Idea: Use a **random forest**: An ensemble of trees generated by ...
 - using a random subset of data to train each tree (i.e., bagging or pasting)
 - considering a random subset of features j at each *split*
- As with other ensemble methods, results are averaged to make a prediction
- Implemented in `sklearn.RandomForestClassifier` & `Regressor`



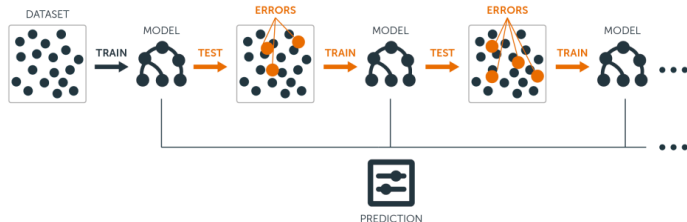
from [machine-learning-algorithms-ensemble-methods-bagging-boosting-and-random-forests](#)

Outline

- Parallel Ensemble Methods: Bagging and Pasting
- Decision Trees and Random Forests
- Boosting: Sequentially Trained Ensemble Methods

Boosting

- Previously we described ensemble methods that train **in parallel**
- Now we discuss ensemble methods that train **sequentially**
 - This is called “**boosting**”
- Some of the best known boosting methods are:
 - adaptive boosting, or **Adaboost** (1996)
 - **gradient boosting** (2001)
 - extreme gradient boosting or **XGBoost** (2016)



Adaboost

- **Adaboost** sequentially trains an ensemble of predictors $\{f_m(\cdot)\}_{m=1}^M$ as follows:

For $m = 1 \dots M$,

- 1 Train predictor $f_m(\cdot)$ to minimize some **weighted loss** $\sum_{i=1}^n w_{mi} L(y_i, f_m(\mathbf{x}_i))$
- 2 For next round, assign larger weights $w_{m+1,i}$ to samples i with higher loss

- The final prediction is done using a **weighted average**

$$F_M(\mathbf{x}) = \sum_{m=1}^M \alpha_m f_m(\mathbf{x}), \quad \text{for some learned weights } \{\alpha_m\}$$

- For simplicity, we will focus on the design of binary classifiers $f_m(\cdot) \in \{-1, 1\}$
 - Other versions of Adaboost exist for regression and non-binary classification

Adaboost: Derivation

- Consider binary classification, with training data $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ for $y_i = \pm 1$
- Adaboost's base classifiers $f_m(\cdot) = \pm 1$, for $m = 1 \dots M$, are trained as follows:
 - Define the step- m **boosted classifier** $F_m(\mathbf{x}) \triangleq F_{m-1}(\mathbf{x}) + \alpha_m f_m(\mathbf{x})$
 - Initialize $F_0(\mathbf{x}) = 0$. For $m = 1, \dots, M$, design $\alpha_m > 0$ and f_m so that F_m minimizes the **exponential loss** $\mathcal{L}_m = \sum_{i=1}^n e^{-y_i F_m(\mathbf{x}_i)}$
- Plugging $F_m(\cdot)$ into \mathcal{L}_m , we find that

$$\begin{aligned}
 \mathcal{L}_m &= \sum_{i=1}^n e^{-y_i (F_{m-1}(\mathbf{x}_i) + \alpha_m f_m(\mathbf{x}_i))} = \sum_{i=1}^n \overbrace{e^{-y_i F_{m-1}(\mathbf{x}_i)}}^{\triangleq w_{mi}} e^{-y_i \alpha_m f_m(\mathbf{x}_i)} \\
 &= e^{-\alpha_m} \sum_{i: y_i = f_m(\mathbf{x}_i)} w_{mi} + e^{\alpha_m} \sum_{i: y_i \neq f_m(\mathbf{x}_i)} w_{mi} \quad \text{since } f_m \text{ and } y_i \text{ are } \pm 1 \\
 &= e^{-\alpha_m} \sum_{i=1}^n w_{mi} + (e^{\alpha_m} - e^{-\alpha_m}) \sum_{i: y_i \neq f_m(\mathbf{x}_i)} w_{mi}
 \end{aligned}$$

Adaboost: Derivation (cont.)

- From the previous expression, we see that the \mathcal{L}_m -minimizing classifier $f_m(\cdot)$ is that which minimizes the **weighted misclassification loss**

$$\sum_{i: y_i \neq f_m(\mathbf{x}_i)} w_{mi} = \sum_{i=1}^n w_{mi} \mathbb{1}_{y_i \neq f_m(\mathbf{x}_i)}$$

- At step $m=1$, use $w_{1i} = 1 \ \forall i$, and thus train $f_1(\cdot)$ on $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ as usual
- At step $m > 1$, the **weight** w_{mi} is large if F_{m-1} made a wrong decision on y_i
- The \mathcal{L}_m -minimizing α_m can be found by solving $\partial \mathcal{L}_m / \partial \alpha_m = 0$, which gives

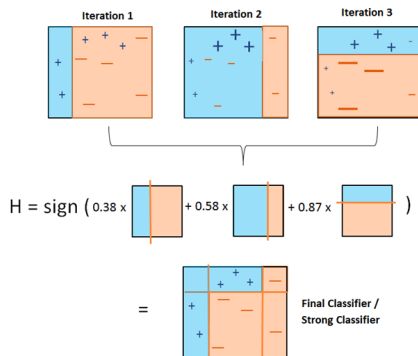
$$\alpha_m = \frac{1}{2} \left(\frac{1 - \epsilon_m}{\epsilon_m} \right) \text{ for } \text{weighted error rate } \epsilon_m = \frac{\sum_{i: y_i \neq f_m(\mathbf{x}_i)} w_{mi}}{\sum_i w_{mi}}$$

- In summary, to train Adaboost, initialize $w_{1i} = 1 \ \forall i$ and then, for $m = 1 \dots M$, compute $f_m(\cdot)$, ϵ_m , α_m , $F_m(\cdot)$, and $\{w_{m+1,i}\}_{i=1}^n$
- Once trained, predict the label of a test sample \mathbf{x} using $F_M(\mathbf{x})$

Adaboost: Implementation

- Adaboost classification implemented via `sklearn.AdaBoostClassifier`
 - By default, it uses 50 decision tree classifiers of depth 1 ("stumps")
- Adaboost regression implemented via `sklearn.AdaBoostRegressor`
 - By default, it uses 50 decision tree regressors of depth 3

AdaBoost Classifier Working Principle with Decision Stump as a Base Classifier



Gradient boosting

- Goal: Design predictor $F_M(\cdot)$ to minimize a training loss of the form

$$\mathcal{L}(F_M) = \sum_{i=1}^n L(y_i, F_M(\mathbf{x}_i)) \quad \text{for some } L(\cdot, \cdot)$$

where $F_M(\cdot) = \sum_{m=1}^M \alpha_m f_m(\cdot)$ with $f_m(\cdot) \in \mathcal{F}$ is sequentially learned

- Adaboost did this with exponential loss $L(y, F) = e^{-yF}$ and binary $f_m = \pm 1$
- Gradient boosting tries to do this for general L and \mathcal{F}
- Ideally, we would like to do the following:
 - Train $f_1 \in \mathcal{F}$ to minimize $\sum_{i=1}^n L(y_i, f_1(\mathbf{x}_i))$, set $F_1 = f_1$
 - For $m = 2 \dots M$: $F_m = F_{m-1} + \arg \min_{\alpha_m, f_m \in \mathcal{F}} \sum_{i=1}^n L(y_i, F_{m-1}(\mathbf{x}_i) + \alpha_m f_m(\mathbf{x}_i))$
 - But for many choices of loss L , this optimization is too difficult
- Idea: Instead of exact minimization, settle for a **gradient step**, i.e.,

For $m = 2 \dots M$: $F_m = F_{m-1} - \alpha_m \nabla \mathcal{L}(F_{m-1})$

Gradient boosting: Details

■ Problem:

- We are limited to updates of the form

$$F_m = F_{m-1} + \alpha_m f_m, \text{ for } f_m \in \mathcal{F}, \text{ where } \mathcal{F} \text{ is set of base predictors}$$

- Gradient descent does not constrain the update to \mathcal{F} :

$$F_m = F_{m-1} + \alpha_m [-\nabla \mathcal{L}(F_{m-1})], \text{ with } \nabla \mathcal{L}(F_{m-1}) = \sum_{i=1}^n \underbrace{\frac{\partial L(y_i, F_{m-1}(\mathbf{x}_i))}{\partial F}}_{\triangleq -r_{mi}}$$

- Solution: Train $f_m \in \mathcal{F}$ to be *close to* $-\nabla \mathcal{L}(F_{m-1})$

- In particular, for $m = 2, \dots, M$, we do the following:

- First train $f_m \in \mathcal{F}$ to minimize $\sum_{i=1}^n (r_{mi} - f_m(\mathbf{x}_i))^2$

- Then choose α_m via line-search:

$$\alpha_m = \arg \min_{\alpha} \sum_{i=1}^n L(y_i, F_{m-1}(\mathbf{x}_i) + \alpha f_m(\mathbf{x}_i))$$

- Finally, update boosted predictor: $F_m(\mathbf{x}) = F_{m-1}(\mathbf{x}) + \alpha_m f_m(\mathbf{x})$

Gradient boosting: Implementation

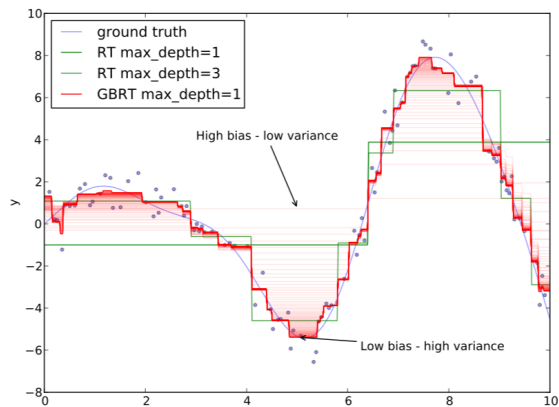
- In practice, it helps to slow down the updates using “learning rate” $\mu \in (0, 1]$:

$$F_m(\mathbf{x}) = F_{m-1}(\mathbf{x}) + \mu \alpha_m f_m(\mathbf{x})$$

- Classification implemented via `sklearn.GradientBoostingClassifier`
 - By default, it uses 100 decision tree classifiers of depth 3
 - By default, uses $\mu = 0.1$
- Regression implemented via `sklearn.GradientBoostingRegressor`
 - By default, it uses 100 decision tree regressors of depth 3
 - By default, uses $\mu = 0.1$

Gradient boosting: Example

```
from sklearn.ensemble import GradientBoostingRegressor
est = GradientBoostingRegressor(n_estimators=2000, max_depth=1).fit(X, y)
for pred in est.staged_predict(X):
    plt.plot(X[:, 0], pred, color='r', alpha=0.1)
```



from Gradient Boosted Regression Trees by Prettenhofer & Louppe

Extreme gradient boosting (XGBoost)

- **XGBoost** is an evolution of gradient boosting that has won many recent machine-learning contests
 - On low-dimensional tasks, XGBoost often performs as well as neural nets, but it is much easier to design/train!
- Compared to gradient boosting, XGBoost uses...
 - second-order Taylor series approximation
 - regularization
 - sophisticated tree-search
 - many parallelization and hardware acceleration tricks
- The **XGBoost package** is not part of **sklearn**, but compatible with it
 - Once you download the XGBoost package, you can **use it just like sklearn**

XGBoost: Derivation

- **XGBoost** sequentially designs a boosted predictor $F_M(\mathbf{x}) = \sum_{m=1}^M f_m(\mathbf{x})$

- Similar to gradient boosting, we'd ideally like to do the following:

- Train $f_1 \in \mathcal{F}$ to minimize $\sum_{i=1}^n L(y_i, f_1(\mathbf{x}_i))$, set $F_1 = f_1$

- For

$$m = 2 \dots M: \quad F_m = F_{m-1} + \arg \min_{f_m \in \mathcal{F}} \left[\sum_{i=1}^n L(y_i, F_{m-1}(\mathbf{x}_i) + f_m(\mathbf{x}_i)) + \phi(f_m) \right]$$

- Note that we added **regularization** ϕ to the cost function

- Unfortunately, this optimization problem is too difficult to solve exactly

- Idea: Simplify using a **2nd-order Taylor series approximation**

- Use $L(y_i, F_{m-1}(\mathbf{x}_i) + f_m(\mathbf{x}_i)) \approx L(y_i, F_{m-1}(\mathbf{x}_i)) + g_{mi}f_m(\mathbf{x}_i) + \frac{1}{2}h_{mi}f_m^2(\mathbf{x}_i)$

$$\text{where } g_{mi} \triangleq \frac{\partial L(y_i, F_{m-1}(\mathbf{x}_i))}{\partial F} \text{ and } h_{mi} \triangleq \frac{\partial^2 L(y_i, F_{m-1}(\mathbf{x}_i))}{\partial^2 F}$$

- Then the optimization problem becomes

$$\arg \min_{f_m \in \mathcal{F}} \left[\sum_{i=1}^n (g_{mi}f_m(\mathbf{x}_i) + \frac{1}{2}h_{mi}f_m^2(\mathbf{x}_i)) + \phi(f_m) \right] \quad (1)$$

XGBoost: Derivation 2

- Now suppose that $f(\cdot)$ is a **decision tree with max L leaves**
 - As on page 15, write $f(\mathbf{x}) = z_{q(\mathbf{x})}$ where ...
 - $z_\ell \in \mathbb{R}$ is the output assigned to leaf ℓ
 - $q(\cdot)$ is the function that assigns $\mathbf{x} \in \mathbb{R}^d$ to a leaf $\ell \in \{1, \dots, L\}$
- XGBoost uses the regularization

$$\phi(f) = \frac{\lambda}{2} \sum_{\ell=1}^L z_\ell^2 + \gamma L, \quad \text{for tunable } \lambda > 0 \text{ and } \gamma > 0$$

- Omitting the “ m ” notation, the XGBoost optimization problem (1) becomes

$$\begin{aligned} & \arg \min_{q, \mathbf{z}} \sum_{\ell=1}^L \sum_{i: q(\mathbf{x}_i) = \ell} \left(g_i z_\ell + \frac{1}{2} h_i z_\ell^2 \right) + \frac{\lambda}{2} \sum_{\ell=1}^L z_\ell^2 + \gamma L \\ &= \arg \min_{q, \mathbf{z}} \sum_{\ell=1}^L \left[\left(\sum_{i: q(\mathbf{x}_i) = \ell} g_i \right) z_\ell + \frac{1}{2} \left(\lambda + \sum_{i: q(\mathbf{x}_i) = \ell} h_i \right) z_\ell^2 + \gamma \right] \end{aligned} \quad (2)$$

XGBoost: Derivation 3

- We can optimize over $z \in \mathbb{R}^L$ to yield

$$z_\ell = -\frac{\sum_{i:q(\mathbf{x}_i)=\ell} g_i}{\lambda + \sum_{i:q(\mathbf{x}_i)=\ell} h_i}, \quad \ell = 1, \dots, L$$

and plug this back into (2) to obtain an optimization problem over $q(\cdot)$:

$$\arg \min_q \sum_{\ell=1}^L \left[2\gamma - \frac{(\sum_{i:q(\mathbf{x}_i)=\ell} g_i)^2}{\lambda + \sum_{i:q(\mathbf{x}_i)=\ell} h_i} \right]$$

- Using the above loss, the tree $q(\cdot)$ can be similarly to page 15: Choose a feature j and search for the threshold t that maximally reduces loss after a split
- For example, consider the split $S \rightarrow (S_1, S_2)$. The loss would change as

$$2\gamma - \frac{(\sum_{i \in S} g_i)^2}{\lambda + \sum_{i \in S} h_i} \rightarrow \left[2\gamma - \frac{(\sum_{i \in S_1} g_i)^2}{\lambda + \sum_{i \in S_1} h_i} + 2\gamma - \frac{(\sum_{i \in S_2} g_i)^2}{\lambda + \sum_{i \in S_2} h_i} \right]$$

If the loss doesn't decrease, then it's better *not* to split!

- Further details can be found in [the original paper](#)

XGBoost: Example

```

1 # First XGBoost model for Pima Indians dataset
2 from numpy import loadtxt
3 from xgboost import XGBClassifier
4 from sklearn.model_selection import train_test_split
5 from sklearn.metrics import accuracy_score
6 # load data
7 dataset = loadtxt('pima-indians-diabetes.csv', delimiter=",")
8 # split data into X and y
9 X = dataset[:,0:8]
10 Y = dataset[:,8]
11 # split data into train and test sets
12 seed = 7
13 test_size = 0.33
14 X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=test_size, ran
15 # fit model no training data
16 model = XGBClassifier()
17 model.fit(X_train, y_train)
18 # make predictions for test data
19 y_pred = model.predict(X_test)
20 predictions = [round(value) for value in y_pred]
21 # evaluate predictions
22 accuracy = accuracy_score(y_test, predictions)
23 print("Accuracy: %.2f%%" % (accuracy * 100.0))

```

Running this example produces the following output.

```

1 Accuracy: 77.95%

```

This is a [good accuracy score on this problem](#), which we would expect, given the capabilities of the model and the modest complexity of the problem.

from [develop-first-xgboost-model-python-scikit-learn](#)

Tuning trees, random forests, and XGBoost

To get the best performance on test data, several hyperparameters must be tuned. Here are some good resources on how to do that:

- Importance of decision tree hyperparameters on generalization
- A beginner's guide to random forest hyperparameter tuning
- XGBoost: Notes on parameter tuning
- Complete guide to parameter tuning in XGBoost with code
- Hyperparameter tuning for hyperaccurate XGBoost model

Learning objectives

- Understand intuition behind ensemble methods: “the wisdom of the crowd”
- Understand parallel ensemble methods
 - bagging, pasting
 - random feature selection
- Understand decision trees
 - feature thresholding and decision regions
 - training via top-down tree induction
 - homogeneity metrics: variance reduction, gini impurity
 - ensemble extension: random forests
- Understand boosting, or sequentially trained ensemble methods
 - Adaboost
 - gradient boosting
 - XGBoost