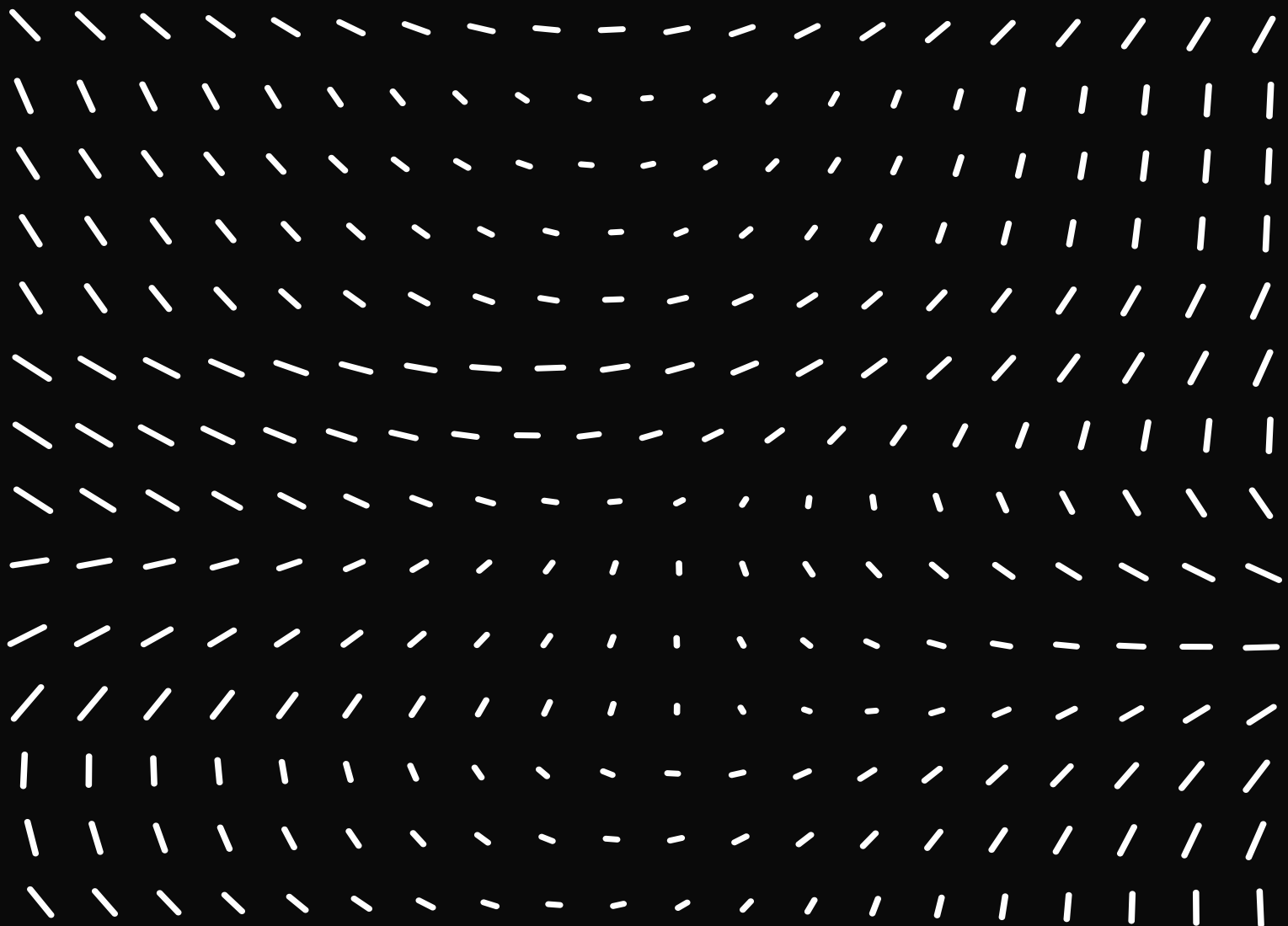# Additive Model & Decision Tree
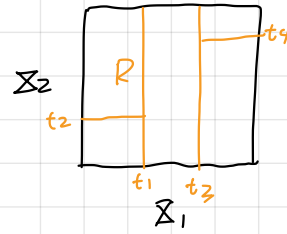
# 1. Generalized Additive Model

$$E(Y \mid X_1, X_2 \cdots X_P) = \alpha + f_1(X_1) + \cdots + f_P(X_P)$$

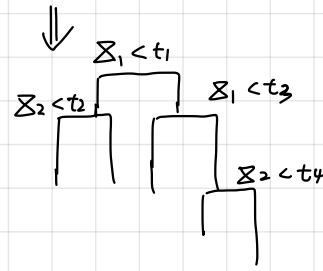# 2. Decision Tree

$$f(z) = \sum_{m=1}^{M} c_m I\{(z_1, z_2) \in R_m\}$$



## 2.1 given data $(x_i, y_i)$

$$x_i = (x_{i1}, x_{i2}, \cdots, x_{ip})$$

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

$$\Downarrow$$

$$\hat{c}_m = avg(y_i | x_i \in R_m)$$

(circled) greedy algorithm

$z_1 < t_1$

$z_2 < t_2 \qquad z_1 < t_3 \qquad$ binary partition

$z_2 < t_4$

proof: given the data $(x_i, y_i)$ and $\begin{cases} \text{splitting var. } j \\ \text{splitting point } s \end{cases}$

$$R_1(j,s) = \{z | z_j \le s\} \qquad R_2(j,s) = \{z | z_j > s\}$$

object: $\min_{j,s} \left[ \min_{c_1} \sum_{x_i \in R_1(j,s)} (y_i - c_1)^2 + \min_{c_2} \sum_{x_i \in R_2(j,s)} (y_i - c_2)^2 \right]$

$$\Downarrow$$

$$c_1 = avg(y_i | x_i \in R_1(j,s)) \qquad c_2 = avg(y_i | x_i \in R_2(j,s))$$

repeat on $R_1$ and $R_2$.

## 2.2 Tree size.

strategy: grow large tree $T_0$. stop when minimum node size reached.

==cost-complexity pruning:== collapse any number of its internal (not terminal) nodes.

$|T|$: # of terminal nodes $\quad (R_1 R_2 \cdots R_m, \text{ index by } m)$

$N_m$: # $\{x_i \in R_m\}$

$$\hat{c}_m = \frac{1}{N_m} \cdot \sum_{x_i \in R_m} y_i$$

$$Q_m(T) = \frac{1}{N_m} \cdot \sum_{x_i \in R_m} (y_i - \hat{c}_m)^2$$

> $T \subseteq T_0$ means $T$ can be obtained by pruning $T_0$

complexity criterion: $C_\alpha(T) = \sum_{m=1}^{|T|} N_m Q_m(T) + \alpha |T|$

objective: for each $\alpha$, find $T_\alpha \subseteq T_0$ that minimize $C_\alpha(T)$

### 2.2.1 * for each $\alpha$, there is a unique smallest $T_\alpha$ that minimize $C_\alpha(T)$

==weakest link pruning:== ☐ successively collapse internal nodes that produces the smallest per-node increase

in $\sum N_m \, Q_m(T)$

until produce a single node tree

☐ collect a sequence of subtrees, must contain $T_\alpha$

☐ use CV to find the best $\hat{\alpha}$ minimize MSE

$\Downarrow$

$T_{\hat{\alpha}}$

2.3. CART: leaf only contain decision values

Tree for classification: ($1, 2, \cdots, K$ categories)

$$\begin{cases} \text{in regression}: & Q_m(T) = \frac{1}{N_m} \sum_{x_i \in R_m} (y_i - \hat{C}_m)^2 \quad \leftarrow \text{called node impurity} \\ \text{in classification}: & \hat{P}_{mk} = \frac{1}{N_m} \sum_{x_i \in R_m} I(y_i = k) \end{cases}$$

$k(m) = \underset{k}{\arg\max} \; \hat{P}_{mk}$  (the majority in node $m$)

$$Q_m(T) = \begin{cases} \text{misclassification error}: \quad \text{(non differentiable)} \\ \quad \frac{1}{N_m} \sum_{i \in R_m} I(y_i \neq k(m)) = 1 - \hat{P}_{mk(m)} \\ \text{Gini index}: \\ \quad \sum_{k \neq k'} \hat{P}_{mk} \hat{P}_{mk'} = \sum_{k=1}^{K} \hat{P}_{mk} (1 - \hat{P}_{mk}) \\ \text{Cross-entropy or deviance}: \\ \quad -\sum_{k=1}^{K} \hat{P}_{mk} \log \hat{P}_{mk} \end{cases}$$

☐ Gini index / cross-entropy used for $T_0$ tree growing

☐ misclassification rate use to guide cost-complexity pruning

# 3. Random Forest. Modification of bagging

trees {
- can capture complex interaction structure
- if grown sufficiently deep, can have low bias
- noisy, so will benefit great from average
- each tree in RF is identical, so $E(\text{average}) = E(\text{itself})$

## 3.1 Bagging for variance reduction

$B$ i.d. trees with correlation $\rho$ and variance $\sigma^2$: $\rho\sigma^2 + \frac{1-\rho}{B}\sigma^2$

when $B \to \infty$, $\frac{1-\rho}{B}\sigma^2 \to 0$, only $\rho\sigma^2$ left

$\rho$ will affect averaging $B$ trees (bagging)

RF: reduce $\rho$ while control $\sigma^2$

## 3.2. from bagging to RF in classification & regression

also called feature bagging
↗ reduce correlation between those strong predicting features

input at each split has $p$ variables,
choose $m \leq p$ at random as candidates for splitting

{
- classification: obtain a class vote from each tree, then use majority

  ($m = \sqrt{p}$ and minimum node size is $1$ )

- regression: average

  ($m = p/3$ and minimum node size is $5$ )

## 3.3 Out of bag samples

for each observation $z_i = (x_i, y_i)$, construct its RF predictor by averaging only those trees corresponding to bootstrap samples in which $z_i = (x_i, y_i)$ did not included.

(kind like CV, and it's built in so we only have to run it in one sequence)

## 3.4 variable importance:

Fit RF to data set, calculate out-of-bag error. Then for each feature $j$, permute it, calculate new out-of-bag error.

importance score = average dif of out-of-bag error → normalize

# 4. XGBoost : decision tree ensembles

consists of a set of CART : $f_k(x)$

Model formula:
$$\hat{y_i} = \sum_{k=1}^{K} f_k(x_i), \quad f_k \in F$$

<span style="color:orange">$K$ : number of trees<br>$F$ : set of all possible CARTs<br>same as RF, dif is how we train them.</span>

Objective function:
$$obj(\theta) = \sum_{i=1}^{n} l(y_i, \hat{y_i}) + \sum_{k=1}^{K} \Omega(f_k)$$

## 4.1  Additive Training

$$\hat{y_i}^{(0)} = 0$$

$$\hat{y_i}^{(1)} = f_1(x_i) = \hat{y_i}^{(0)} + f_1(x_i)$$

$$\hat{y_i}^{(2)} = f_1(x_i) + f_2(x_i) = \hat{y_i}^{(1)} + f_1(x_i)$$

$$\vdots$$

$$\hat{y_i}^{(t)} = \sum_{k=1}^{t} f_k(x_i) = \hat{y_i}^{(t-1)} + f_t(x_i)$$

### 4.1.1  objective at step $t$ :

$$obj^{(t)} = \sum_{i=1}^{n} l(y_i, \hat{y_i}^{(t)}) + \sum_{k=1}^{t} \Omega(f_k) \qquad \textcolor{orange}{constant}$$

$$= \sum_{i=1}^{n} l(y_i, \hat{y_i}^{(t-1)} + f_t(x_i)) + \Omega(f_t) + \textcolor{orange}{\boxed{\sum_{k=1}^{t-1} \Omega(f_k)}}$$

if $l$ is MSE :

$$obj^{(t)} = \sum_{i=1}^{n} [(y_i - \hat{y_i}^{(t-1)}) - f_t(x_i)]^2 + \Omega(f_t) + constant$$

$$= \sum_{i=1}^{n} [\underbrace{\textcolor{orange}{(y_i - \hat{y_i}^{(t-1)})^2}}_{\textcolor{orange}{constant}} + f_t(x_i)^2 - 2(y_i - \hat{y_i}^{(t-1)}) f_t(x_i)] + \Omega(f_t) + \cdots$$

$$= \sum_{i=1}^{n} [f_t(x_i)^2 - 2(y_i - \hat{y_i}^{(t-1)}) f_t(x_i)] + \Omega(f_t) + constant$$

In general :

$$obj^{(t)} = \sum_{i=1}^{n} [l(y_i, \hat{y_i}^{(t-1)}) + g_i f_t(x_i) + \tfrac{1}{2} h_i f_t^2(x_i)] + \Omega(f_t) + constant$$

$$\begin{cases} g_i = \dfrac{\partial}{\partial \hat{y_i}^{(t-1)}} l(y_i, \hat{y_i}^{(t-1)}) & \textcolor{orange}{first\ derivative} \\[2mm] h_i = \dfrac{\partial^2}{\partial \hat{y_i}^{(t-1)2}} l(y_i, \hat{y_i}^{(t-1)}) & \textcolor{orange}{second\ derivative} \end{cases}$$

$$\vdots$$

$$= \sum_{i=1}^{n} [g_i f_t(x_i) + \tfrac{1}{2} h_i f_t^2(x_i)] + \Omega(f_t)$$

<span style="color:orange">this made XGBoost possible to support custom loss function</span>

### 4.1.2. regularization term $\Omega(f_t)$

$$f_t(x) = w_{q(x)} \quad , \quad w \in R^T \quad , \quad q \in R^d \rightarrow \{1, 2, 3, \cdots, T\}$$

$$\begin{cases} w: \text{ vector of scores on leaves} \\ T: \text{ number of leaves} \\ q: \text{ function that assign each data point } (R^d) \rightarrow \text{a leaf} \end{cases}$$

$$\Omega(f_t) = \gamma T + \tfrac{1}{2} \lambda \sum_{j=1}^{T} w_j^2$$

### 4.1.3. structure score

$$obj^{(t)} = \sum_{i=1}^{n} \left[ g_i w_{q(x_i)} + \tfrac{1}{2} h_i w_{q(x_i)}^2 \right] + \gamma T + \tfrac{1}{2} \lambda \sum_{j=1}^{T} w_j^2$$

$$= \sum_{j=1}^{T} \left[ \left( \sum_{i \in I_j} g_i \right) w_j + \tfrac{1}{2} \left( \sum_{i \in I_j} h_i + \lambda \right) w_j^2 \right] + \gamma T$$

$$I_j = I\{ i \mid q(x_i) = j \} \quad \text{set of indices of data point assign to leaf } j$$

$$= \sum_{j=1}^{T} \left[ G_j w_j + \tfrac{1}{2} (H_j + \lambda) w_j^2 \right] + \gamma T$$

$$\begin{cases} G_j = \sum_{i \in I_j} g_i \\[2mm] H_j = \sum_{i \in I_j} h_i \end{cases}$$

$w_j$ are independent to each other, by solve the $\min\{obj^{(t)}\}$

we get $\begin{cases} w_j^* = -\dfrac{G_j}{H_j + \lambda} \\[4mm] obj^{(t)*} = -\tfrac{1}{2} \sum_{j=1}^{T} \dfrac{G_j^2}{H_j + \lambda} + \gamma T \quad \leftarrow \text{ measures how good} \\ \qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\text{a tree is} \end{cases}$

when trying to creat a split of a leaf into L and R

$$\triangle obj = \gamma - \tfrac{1}{2} \left[ \frac{G_L^2}{H_L + \lambda} + \frac{G_R^2}{H_R + \lambda} - \frac{(G_L + G_R)^2}{H_L + H_R + \lambda} \right]$$

if it's smaller than $\gamma$, then obj $\nearrow$, no need to add this split (aka pruning)