

①

The Delta Method

Random variable $Y = f(X)$, and we know $E(X)$ and $V(X)$.

What can we say about $E(Y)$ and $V(Y)$?

\Rightarrow setup CI for Y based
on CI for X

- Linearize Y in X around $E(X)$

$$Y \approx f(E(X)) + f'(X)|_{E(X)} (X - E(X)) + \text{higher order terms}$$

↓
 derivative
 evaluated @ $E(X)$
Ignore these

$\Rightarrow E(Y) \approx f(E(X))$

$$V(Y) \approx \left(f'(X)|_{E(X)}\right)^2 V(X)$$

- Example of application

Half-life $T = \frac{\ln 2}{\delta}$

δ ← decay rate

$\hat{T} = \frac{\ln 2}{\hat{\delta}}$ plug-in estimator

$$V(\hat{T}) \approx \left(-\frac{\ln 2}{\hat{\delta}^2}\right)^2 V(\hat{\delta}) \Rightarrow CI(T) = \left[\frac{\ln 2}{\hat{\delta}} \pm q_{\alpha/2} \left(\frac{\ln 2}{\hat{\delta}^2} \right) SE(\hat{\delta}) \right]$$

CART

Classification & Regression Trees

(2)

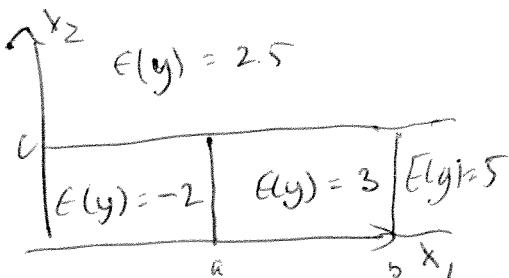
Idea - partition X -space into rectangular regions

whose $E(y)$ takes on only one value

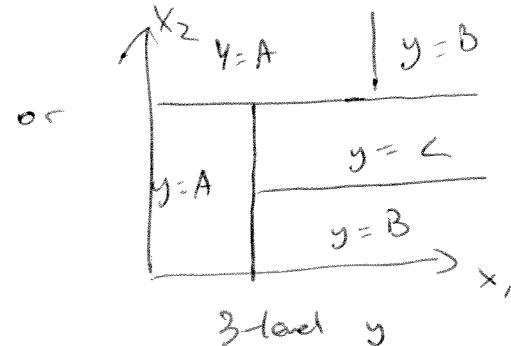
y regression tree; $E(Y | \text{region } k) = \mu_k$

y classification tree; Y in region k = category l

Example



continuous/numerical y.



The model looks like this

$$E(y | X \in R_k) = \mu_k$$

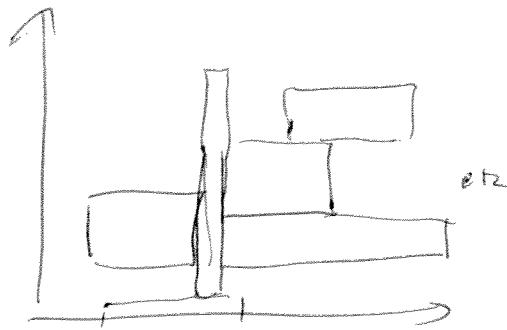
$$R_k = \{x_j > \tau_{kj}\} \times \{x_j \leq \tau_{kj}\} \times \{x_{j''} \leq \tau_{kj''}\} \dots \text{etc}$$

Note can have same x-variable appear multiple times above

E.g. in example

$$E(y) = 3 \text{ region is } \{x_2 \leq c\} \times \{x_1 > a\} \times \{x_1 \leq b\}$$

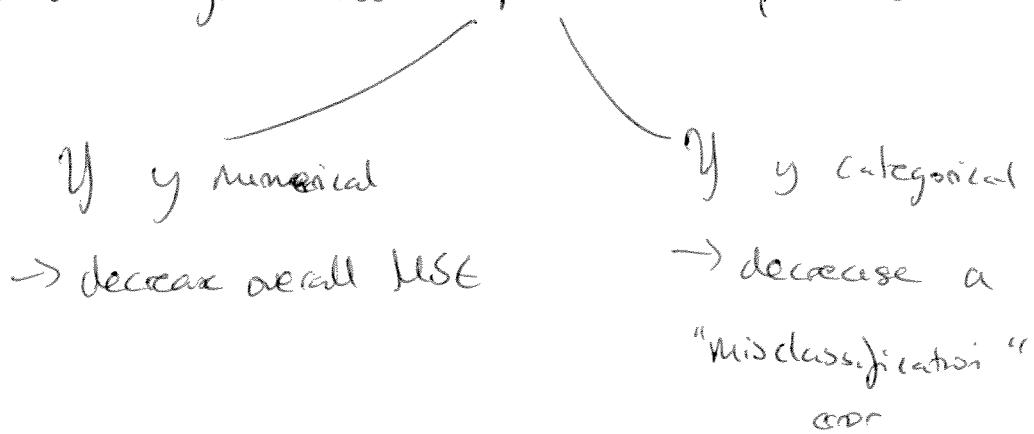
We can't build any type of rectangular regions... ③



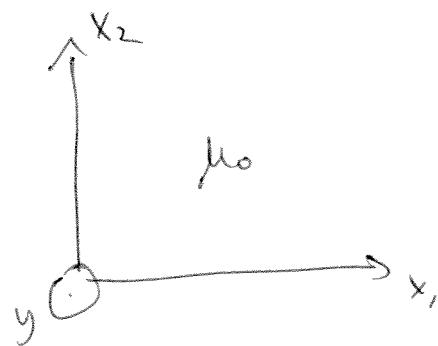
- because the model space to search over
 - which variables & which thresholds
- is too large.

What CART does is use a forward search method.

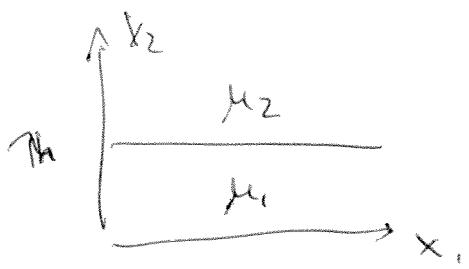
The idea is, to build the regions one variable ~~&~~ threshold at a time. Each step tries to make a newly created region as "pure" as possible



(4)

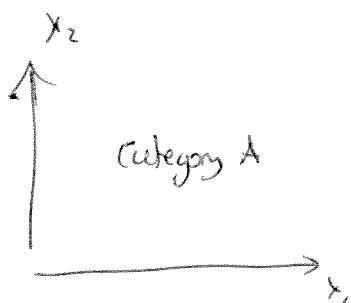
ExampleBefore split MSE₀

$$= \frac{1}{n} \sum_{i=1}^n (y_i - \mu_0)^2$$



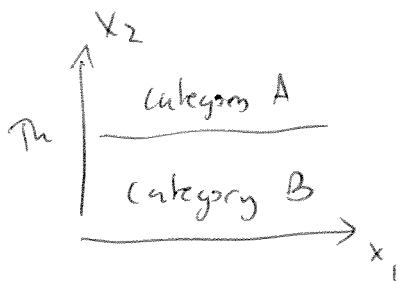
$$MSE_1 = \frac{1}{n} \sum_{i=1}^n (y_i - \mu_1 \mathbb{1}\{x_{i2} \leq Th\} - \mu_2 \mathbb{1}\{x_{i2} > Th\})^2$$

Find the x-variable and value for Th that minimizes the MSE

Example 2

Before split

$$\text{Misclass. error rate} = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{y_i \neq A\}$$

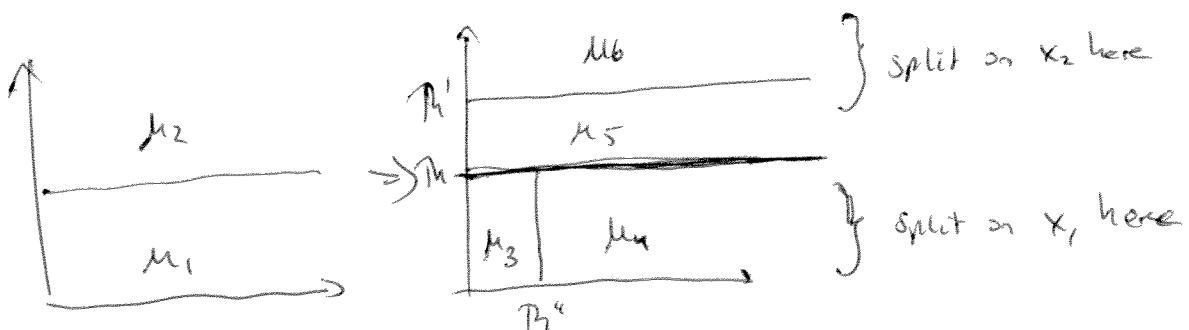


After

$$\text{Misclass. error rate} = \frac{1}{n} \sum_{i=1}^n \left(\mathbb{1}\{x_{i2} \geq Th_1\} \cdot \mathbb{1}\{y_i \neq A\} + \mathbb{1}\{x_{i2} \leq Th_2\} \cdot \mathbb{1}\{y_i \neq B\} \right)$$

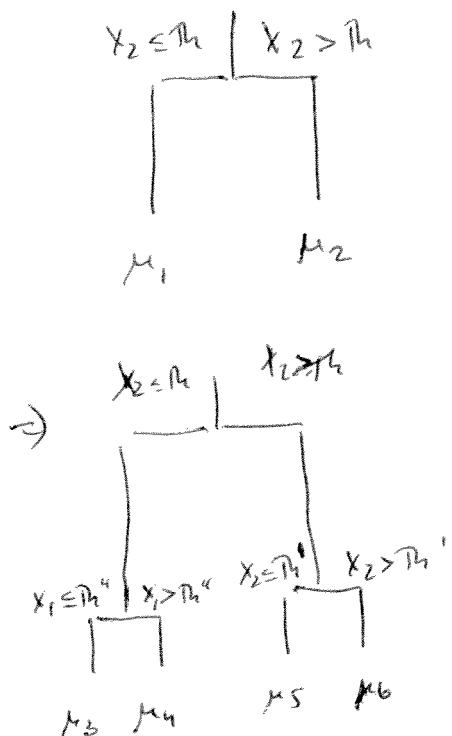
Next step — look at data in each region separately

and find the ultimate split here



We can write these models in a tree format

(5)



The length of the branches
are proportional to the
reduction in MSE
due to the split.

Here, the first split ($X_2 > \bar{x}_h$
 $X_2 \leq \bar{x}_h$)
reduced the overall MSE
much more than subsequent splits.

How many splits should we use?

- This is like choosing the number of variables in a regression model.

The most commonly used method is Cross-validation

Build a big tree, lots of splits. Then try "pruning" it by cutting branches. Evaluate the CV MSE or error rate after each pruning step and stop pruning when the CV error increases.

DEMO.

(6)

(Autonomy remark.)

(CART is notoriously unstable)

That is, if you run CART on a slightly different data set, the tree can look very very different - especially if x 's are correlated.

- One of the best methods for

- prediction is to run CART on

- several bootstrap data sets (resample observations)

 (x_i, y_i)

- and use an average model, average

- prediction from all trees built.

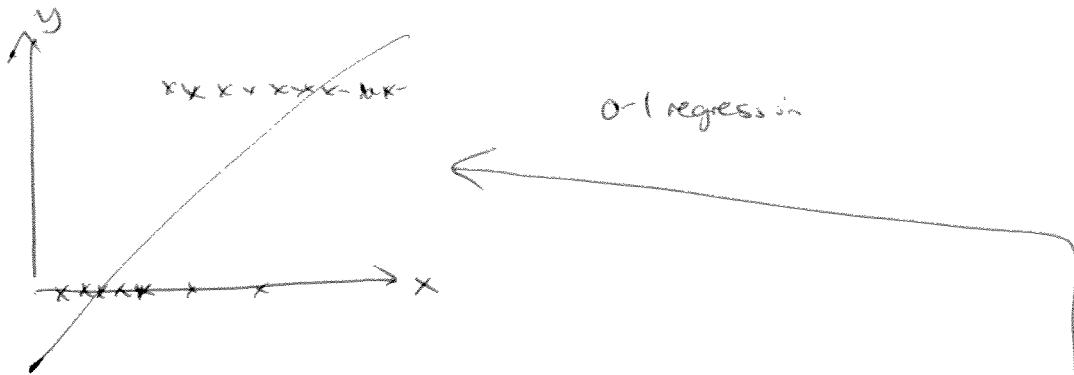
- This is called Bagging (but similar

- procedures are also used for non-tree models).

(Generalized) Linear Models (GLM)

①

What if $y \in \{0,1\}$, can we still run regression?



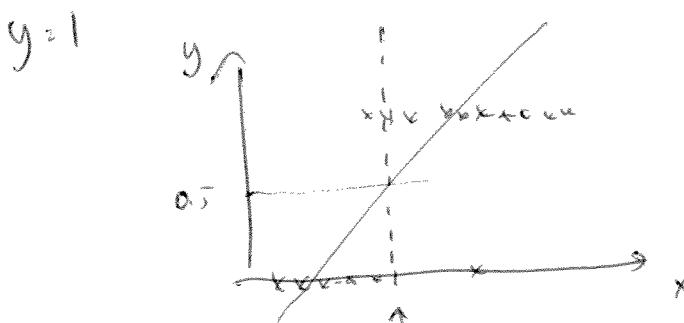
$$E(y) = x\beta$$

When $y \in \{0,1\}$ $E(y) = P(y=1)$

do interpretation of
regression line $\$$

$$P(y=1 | X=x) = x\beta$$

Natural to use $P > 0.5$ as a cut-off for predicting



$X: \{x\beta = 0.5\}$ is called the decision boundary.

This easily generalizes to more than one x , and we can run model selection etc.

Problem? $x\beta$ can be > 1 and < 0 so cannot really represent $P(y=1 | X=x)$.

Logistic regression

②

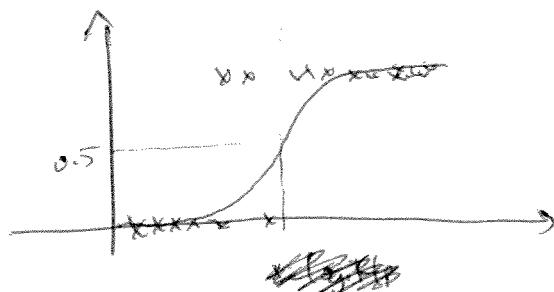
What if we don't assume $p(y=1|X=x)$ is linear in X
 but rather that some transformation of p is linear in X .

Example of a transformation that works, Logit

$$\text{Logit: } \log \frac{p(y=1|X=x)}{1-p(y=1|X=x)} = \log \frac{p(y=1|X=x)}{p(y=0|X=x)} = X\beta$$

"Log-odds"

$$\Rightarrow p(y=1|x=x) = \frac{e^{x\beta}}{1+e^{x\beta}} \in (0,1) \text{ always}$$



$$x: \frac{e^{x\beta}}{1+e^{x\beta}} = 0.5$$

General idea

$$\mathbb{E}(y_i) = \mu_i \Rightarrow g(\mu_i) = x_i'\beta$$

/

link function!

Recap of models we know how to work with

(3)

① LM $y = X\beta + \varepsilon$ $\begin{cases} E(\varepsilon) = 0 \\ V(\varepsilon) = \sigma^2 I \\ \varepsilon \sim N(0, \sigma^2) \end{cases}$ $\Rightarrow y \sim N(X\beta, \sigma^2 I)$
 $\Rightarrow \hat{\beta} = (X'X)^{-1}X'y$
 $\hat{\sigma}^2 = \frac{RSS(\hat{\beta})}{n-p}$

Closed form solutions

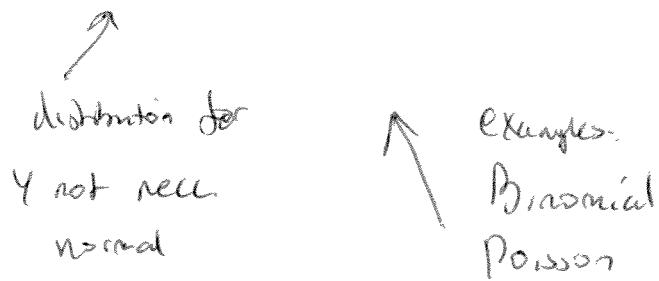
② NLIN $y = f(X; \beta) + \varepsilon$ $\begin{cases} E(\varepsilon) = 0 \\ V(\varepsilon) = \sigma^2 I \\ \varepsilon \sim N(0, \sigma^2) \end{cases}$ $\Rightarrow y \sim N(f(X; \beta), \sigma^2 I)$

\Rightarrow fit by iterative least squares
(ILS)

\rightarrow Careful w. starting values, convergence
& confidence intervals

Can use F-test for model selection still.

③ GLM $y \sim f_Y(x; \beta)$ - allow for non-normal errors



Binomial : $y_i \sim \text{Bin}(m, \pi_i)$, $\text{logit}(\pi_i) = x_i \beta$
(need $\pi_i \in (0, 1)$)

Poisson : $y_i \sim \text{Poi}(\mu_i)$, $\log(\mu_i) = x_i \beta$
(need $\mu_i > 0$)

Assumptions for GLM

(9)

- y_i independent observations (like uncorr. ϵ in LM)

- $E(y_i) = \mu_i$, $g(\mu_i) = \eta_i = x_i' \beta$



 η_i called the linear predictor.
 link
function

- $V(y_i)$ related to μ_i in a known fashion

- Example $y_i \sim \text{Poisson}(\mu_i) \Rightarrow \begin{cases} E(y_i) = \mu_i \\ V(y_i) = \mu_i \end{cases}$

$$y_i \sim \text{Bin}(n, \pi_i) \Rightarrow \begin{cases} E(y_i) = n\pi_i \\ V(y_i) = n\pi_i(1-\pi_i) \end{cases}$$

- Also - assume $g(\cdot)$ is the "correct" link (fits the data)

- and no outliers

- Example, radioactive decay data

$$\text{Nth } y_i = d + \beta e^{-\gamma t_i} + \epsilon_i$$

$$\text{Now } y_i \sim \text{Poi}(\mu_i), \mu_i = \beta e^{-\gamma t_i} = e^{\log \beta - \gamma t_i} = e^{(\beta' + \gamma' t_i)} = e^{x_i' \beta}$$

linear predictor $\beta' + \gamma' t_i$, (β', γ') reparametrization

of problem into what is called "canonical parameters"

Poisson is an example of an exponential family distribution. (5)

- properties of Gb & well understood here.

Dehnison

Exp. Family distribution $f(y | \theta, \phi) = \exp \left[\frac{y\theta - b(\theta)}{a(\phi)} + l(y, \phi) \right]$

\rightarrow can generalize to $f(y_i | \theta_i, \varphi)$

Example : normal error

$$f(y) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{1}{2}(y-\mu)^2\right\} \quad \mu = x\beta$$

$$-\log f(y) = -\frac{1}{2\sigma^2} (y - \mu)^2 - \frac{1}{2} \log(2\pi\sigma^2)$$

$$= \frac{1}{6^2} \left(y\mu - \frac{1}{2}\mu^2 - \frac{1}{2}y^2 \right) - \frac{1}{2} \log 2\pi e^2$$

$$\Rightarrow \begin{cases} \theta = \mu = x\beta \\ \psi = 6^2 \end{cases} \quad \text{that is } g(\mu) = \theta = \text{identity link}$$

$$\Rightarrow b(\theta) = \frac{1}{2} \mu^2$$

Poisson

$$f(y) = \frac{e^{-\mu} \mu^y}{y!}$$

$$\log(f(y)) = y \log(\mu) - \mu - \log(y!)$$

$$\Rightarrow \begin{cases} \theta = \log(\mu) \\ (\approx x, \beta) \\ \phi = 1 \end{cases} \Rightarrow \mu = e^\theta \quad \text{log-link}$$

$$\Rightarrow b(\theta) = \mu = e^\theta$$

Binomial

$$f(y) = \binom{m}{y} \pi^y (1-\pi)^{1-y}$$

$$\log(f(y)) = \log(m) + y \log(\pi) + (1-y) \log(1-\pi)$$

$$= y \log\left(\frac{\pi}{1-\pi}\right) + \log(1-\pi) + \log\binom{m}{y}$$

$$\Rightarrow \begin{cases} \theta = \log\left(\frac{\pi}{1-\pi}\right) \\ \phi = 1 \end{cases}$$

(9)

WLS

let's say $V(y_i) = \sigma^2$ not constant σ^2

$\Rightarrow V(y) = \sigma^2$ ~~Y~~ varying component. y ~~Y~~ non-diagonal — y 's are correlated as well.

\Rightarrow Use LS $\Rightarrow \hat{\beta} = (X'X)^{-1}X'y$

$$E(\hat{\beta}) = \beta, V(\hat{\beta}) = \sigma^2 (X'X)^{-1} (X'VX)(X'X)^{-1}$$

~~transformed variables~~

WLS

$$\min_{\beta} \sum_{i=1}^n w_i (y_i - x_i' \beta)^2 = \sum_{i=1}^n (w_i^{-1/2} (y_i - x_i' \beta))^2 = \| W^{1/2} y - W^{1/2} X \beta \|_2^2$$

$$= \| \tilde{y} - \tilde{X} \beta \|_2^2 \quad \text{transformed variables}$$

$$\begin{cases} \tilde{y} = W^{1/2} y \\ \tilde{X} = W^{1/2} X \end{cases}$$

$$\Rightarrow \hat{\beta} = (\tilde{X}' \tilde{X})^{-1} \tilde{X}' \tilde{y}, \hat{y} = Hy = \tilde{X}' \hat{\beta}$$

$$= (X' W X)^{-1} X' W y = \underbrace{(W^{1/2} X (X' W X)^{-1} X' W^{1/2})}_H y$$

Now if $V(y) = \sigma^2 I$ and we use $W = V^{-1}$

$$\Rightarrow V(\hat{\beta}) = \sigma^2 (X' W X)^{-1} \quad \text{best possible}$$

That is, use weights $w_i \sim \frac{1}{V(y_i)}$