## 7.5.2 LOGISTIC REGRESSION FOR $K > 2$ CLASSES

☛ Following the case where $K = 2$, in the logistic case, it may seem that we could assume that, for $k = 1, \ldots, K$,

$$m_k(\mathbf{x}) = E(Y_k | \mathbf{X} = \mathbf{x}) = P(G = k | \mathbf{X} = \mathbf{x}) = \frac{\exp(\beta_{k0} + \beta_k^T \mathbf{x})}{1 + \exp(\beta_{k0} + \beta_k^T \mathbf{x})}.$$

☛ However as seen earlier, we must have $m_1(\mathbf{x}) + \ldots + m_K(\mathbf{x}) = 1$.

☛ This can be satisfied if we use the following slightly different model: for $k = 1, \ldots, K - 1$:

$$m_k(\mathbf{x}) = \frac{\exp(\beta_{k0} + \beta_k^T \mathbf{x})}{1 + \sum_{\ell=1}^{K-1} \exp(\beta_{\ell 0} + \beta_\ell^T \mathbf{x})}.$$

☛ Then we take

$$m_K(\mathbf{x}) = 1 - \sum_{k=1}^{K-1} m_k(\mathbf{x}) = 1 - \frac{\sum_{k=1}^{K-1} \exp(\beta_{k0} + \beta_k^T \mathbf{x})}{1 + \sum_{\ell=1}^{K-1} \exp(\beta_{\ell 0} + \beta_\ell^T \mathbf{x})}$$

$$= \frac{1}{1 + \sum_{\ell=1}^{K-1} \exp(\beta_{\ell 0} + \beta_\ell^T \mathbf{x})}.$$

☛ For $k = 1, \ldots, K - 1$, estimate the $\beta_{k0}$'s and $\beta_k$'s by their maximum likelihood estimators $\hat{\beta}_{k0}$, $\hat{\beta}_k$ (again can be applied to PLS components instead) and deduce estimators $\hat{m}_k$ of $m_k$. Then estimate $m_K(\mathbf{x})$ by

$$\hat{m}_K(\mathbf{x}) = 1 - \sum_{k=1}^{K-1} \hat{m}_k(\mathbf{x}).$$

☛ Classify the new individual in group $\hat{G}$ which gives the max value among $\hat{m}_1(\mathbf{x}), \ldots, \hat{m}_K(\mathbf{x})$:

$$\widehat{G} = \arg \max_{k=1,\ldots,K} \widehat{m}_k(\mathbf{x})$$

☛ Note that this is again a linear method. To see why, note that to find the max it suffices to compute, for $k = 1, \ldots, K - 1$

$$\log \frac{\hat{P}(G = k | \mathbf{X} = \mathbf{x})}{\hat{P}(G = K | \mathbf{X} = \mathbf{x})} = \log \frac{\hat{m}_k(\mathbf{x})}{\hat{m}_K(\mathbf{x})} = \hat{\beta}_{k0} + \hat{\beta}_1^T \mathbf{x} \, ,$$

and compare the groups 2 by 2 based on those log ratios.

☞ Ex: suppose $K > 5$ and $\hat{m}_5(\mathbf{x})$ is the max. Then

$$\hat{m}_5(\mathbf{x}) > \hat{m}_K(\mathbf{x}) \iff \log\{\hat{m}_5(\mathbf{x})/\hat{m}_K(\mathbf{x})\} > 0$$

and for all $k \neq 5$ and $\neq K$,

$$\hat{m}_k(\mathbf{x}) < \hat{m}_5(\mathbf{x}) \iff \frac{\hat{m}_k(\mathbf{x})}{\hat{m}_K(\mathbf{x})} < \frac{\hat{m}_5(\mathbf{x})}{\hat{m}_K(\mathbf{x})} \iff \log\frac{\hat{m}_k(\mathbf{x})}{\hat{m}_K(\mathbf{x})} < \log\frac{\hat{m}_5(\mathbf{x})}{\hat{m}_K(\mathbf{x})}$$

so we are able to make our decision based only on these log ratios which are linear in $\mathbf{x}$.

*LD, QD based on normality assumption.*

☛ Using the same ideas, LD and QD methods can be generalised to the case where data come from $K > 2$ groups.

☛ Let $\pi_1 = P(G = 1), \ldots, \pi_K = P(G = K)$.

☛ LD: Assume that for $k = 1, \ldots, K$, $\mathbf{X}|G = k \sim N_p(\mu_k, \Sigma)$ and classify new obs $\mathbf{x}$ in group $k$ that maximises $\hat{m}_k$, or equivalently,

$$\delta_k(\mathbf{x}) = \log \hat{\pi}_k + \mathbf{x}^T \hat{\Sigma}^{-1} \hat{\mu}_k - \frac{1}{2} \hat{\mu}_k^T \hat{\Sigma}^{-1} \hat{\mu}_k.$$

☛ QD: Assume that for $k = 1, \ldots, K$, $\mathbf{X}|G = k \sim N_p(\mu_k, \Sigma_k)$ and classify new obs $\mathbf{x}$ in group $k$ that maximises

$$\delta_k(\mathbf{x}) = \log(\hat{\pi}_k) - \frac{1}{2} \log\{\det(\hat{\Sigma}_k)\} - \frac{1}{2}(\mathbf{x} - \hat{\mu}_k)^T \hat{\Sigma}_k^{-1}(\mathbf{x} - \hat{\mu}_k).$$

☛ As in case $K = 2$ groups, for simplicity, relabel the data $\mathbf{X}_i$ making the distinction of which group they come from: $\mathbf{X}_{i,k}$, $i = 1, \ldots, n_k$ denote all the $\mathbf{X}_i$'s that are from group $k$;

☛ $\hat{\pi}_k = n_k/n$ is the proportion of training data that are from group $k$ or $\hat{\pi}_k = 1/K$ depending on our beliefs,

$$\hat{\mu}_k = \overline{\mathbf{X}}_k = \frac{1}{n_k} \sum_{i=1}^{n_k} \mathbf{X}_{i,k}$$

$$\widehat{\Sigma}_k = \frac{1}{n_k - 1} \sum_{i=1}^{n_k} (\mathbf{X}_{i,k} - \widehat{\mu}_k)(\mathbf{X}_{i,k} - \widehat{\mu}_k)^T$$

and

$$\widehat{\Sigma} = \frac{1}{n_1 + \ldots + n_K - K} \sum_{k=1}^{K} \sum_{i=1}^{n_k} (\mathbf{X}_{i,k} - \widehat{\mu}_k)(\mathbf{X}_{i,k} - \widehat{\mu}_k)^T.$$

then $p$

*logistic、LDQD、linear regression 都是参数模型、要估计参数,还有很多假设(正态)*

# 8 CLASSIFICATION AND REGRESSION TREES AND RELATED METHODS

## 8.1 CLASSIFICATION AND REGRESSION TREES (CART)

*确实有功能 是种 approximation.*

Hastie et al. (2017), section 9.2 (second edition, 12th printing).

### 8.1.1 INTRODUCTION

☞ Methods from previous chapter rely on strong parametric assumptions (linear model, logistic model or normality assumptions).

☞ When these assumptions are too far from the truth, the performance of classifiers can be poor; e.g. recall the example at page 253.

☞ We need more flexible models that are less driven by strong parametric assumptions.

☞ Instead of using linear or logistic regression assumptions, one possibility is to use regression trees.

*based on non-parametric assumption*

## 8.1.2 REGRESSION TREES

☞ Suppose we observe an i.i.d. sample $(\mathbf{X}_1, Y_1), \ldots, (\mathbf{X}_n, Y_n)$ coming from the regression model

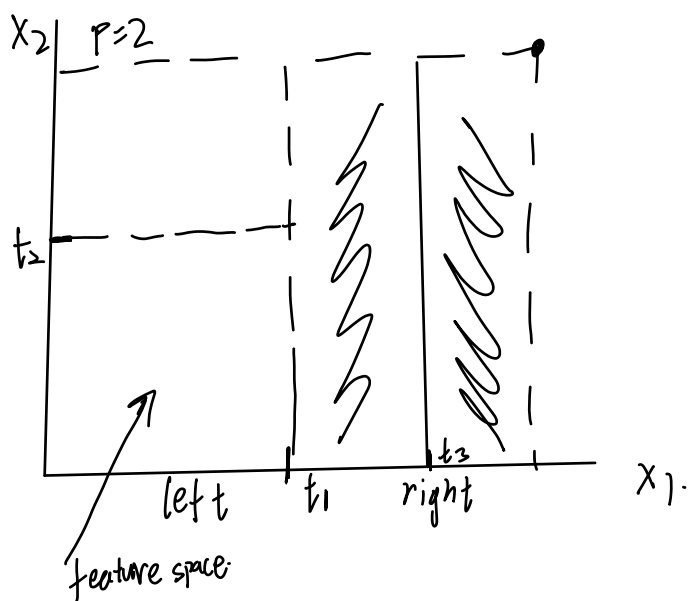$$E(Y_i | \mathbf{X}_i = \mathbf{x}) = m(\mathbf{x}),$$

where $Y_i \in \mathbb{R}$ and $\mathbf{X}_i = (X_{i1}, \ldots, X_{ip})^T \in \mathbb{R}^p$ is continuous.

☞ When we don't want to highlight the dependence on the sample, we use the notation $(\mathbf{X}, Y)$, where

$$E(Y | \mathbf{X} = \mathbf{x}) = m(\mathbf{x}),$$

$Y \in \mathbb{R}$ and $\mathbf{X} = (X_1, \ldots, X_p)^T \in \mathbb{R}^p$ is continuous.

☞ Note the difference between the notation $\mathbf{X}_i$ for the $i$th training vector, and $X_i$ for the $i$th component of the vector $\mathbf{X}$. We will use both notations a lot in this chapter.

$X_2$ | P=2

$t_2$

left    $t_1$    right    $t_3$    $X_1$.

feature space.

1st split: devide    $(X_1 < t_1$    $X_1 > t_1)$ into 2 regions.

2st split    $(X_2 < t_2$    $X_2 > t_2)$

$t_1 < X_1 < t_3$        $X_1 > t_3$    · · · · · · , 不断地划分区域来找X的范围
                                            来创造空间.
                            same.

gradually → divide big rectangles into smaller ~~rectangles~~
rectangles each time. and each time
work on one rectangle. (either vertically or horizontally).

$(X_1$ less/greater to $t_2)$
/$X_2$

☛ The main idea of regression trees is to partition the "feature space", i.e. the domain of $\mathbf{X}$, i.e. the subset of $\mathbb{R}^p$ of all possible values of $\mathbf{X}$.

☛ On each partition we approximate $m(\mathbf{x})$ by a constant.

☛ To understand regression trees, consider the case where $p = 2$. Here we need to estimate the regression curve

$$m(\mathbf{x}) = E(Y|\mathbf{X} = \mathbf{x})$$

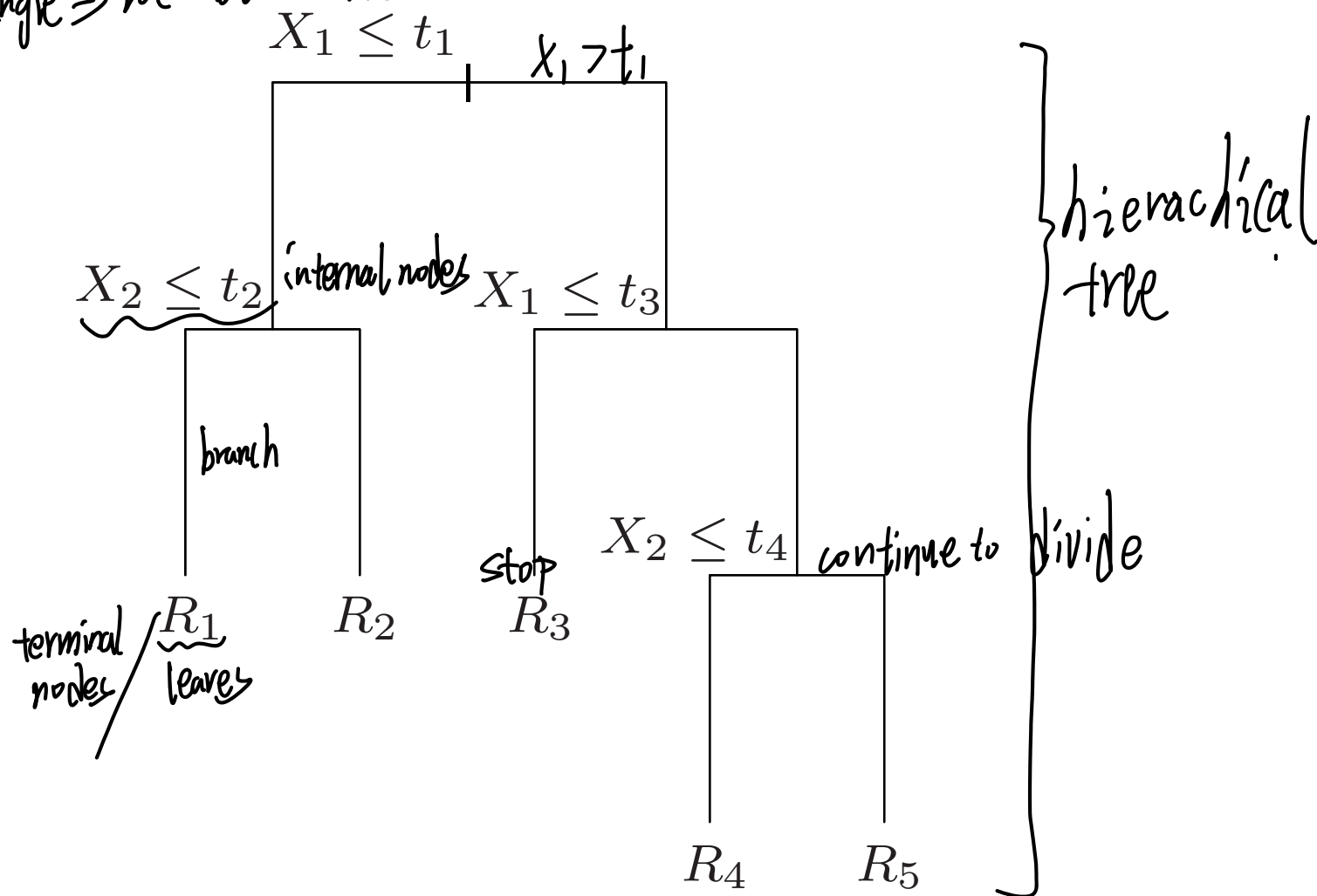where $Y \in \mathbb{R}$ and $\mathbf{X} = (X_1, X_2)^T \in \mathbb{R}^2$ is continuous.

☛ Then we construct a sequence of partitions of the type:
- $\{X_1 \leq t_1\}, \{X_1 > t_1\}$
- $\{X_2 \leq t_2\}, \{X_2 > t_2\}$
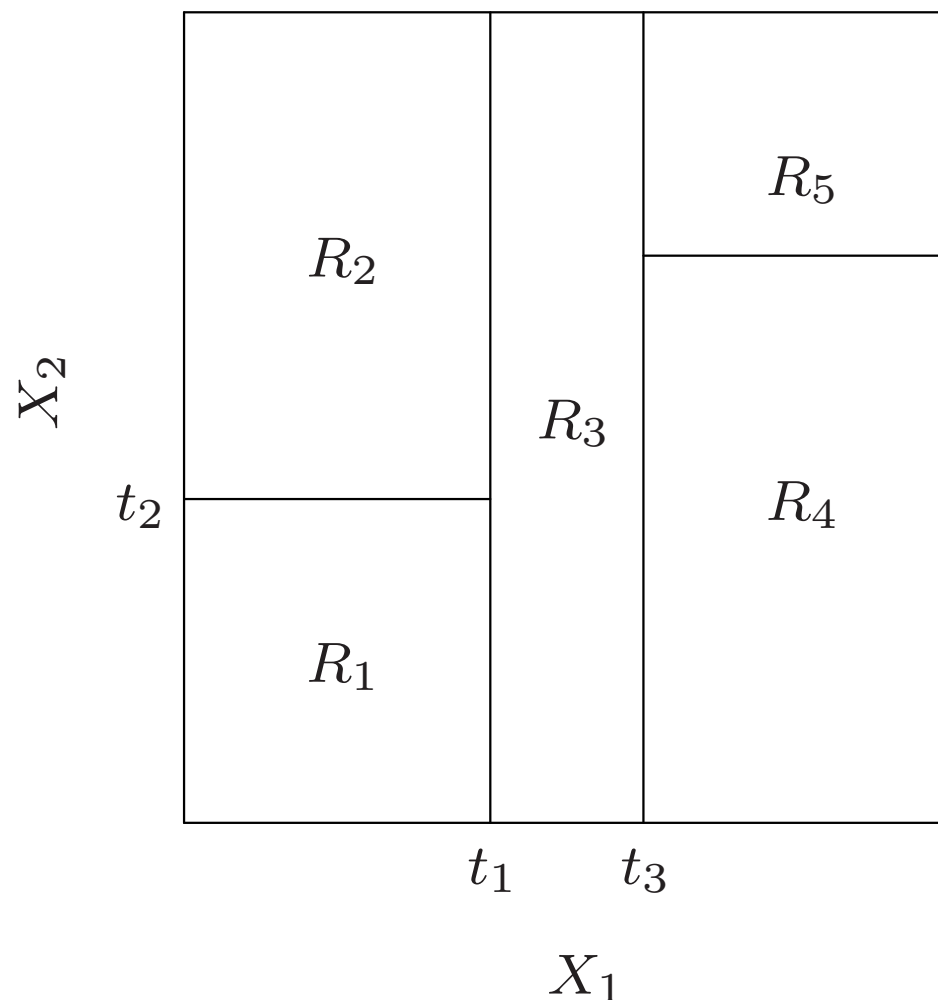- $\{X_1 \leq t_3\}, \{X_1 > t_3\}$

- etc. See handwritten construction on video.

by dividing rectangle $\Rightarrow$ we obtain a tree.

$X_1 \leq t_1$    $X_1 > t_1$

$X_2 \leq t_2$   internal nodes   $X_1 \leq t_3$

branch

hierachical tree

stop   $X_2 \leq t_4$   continue to divide

terminal nodes / $R_1$    $R_2$    $R_3$
   leaves

$R_4$    $R_5$

☛ page 306 of Hastie et al. (2017): At the end of the sequence of binary partitions, we have partitioned the feature space in rectangles, say $R_1, R_2, \ldots, R_L$. Example when $p = 2$ and $L = 5$:

divided

$X_j$ greater/less than

☞ The regions $R_1, \ldots, R_L$ obtained at the end of the process are called terminal nodes or leaves of the tree.

終點　　　葉子

☞ The splits such as $\{X_1 \le t_1\}$, inside the tree, are called internal nodes.

☞ The segments of the tree that connect the nodes are called branches of the tree.

☞ Once we have partitioned the space into regions $R_1, \ldots, R_L$, on each region we approximate the regression curve $m$ by a constant:
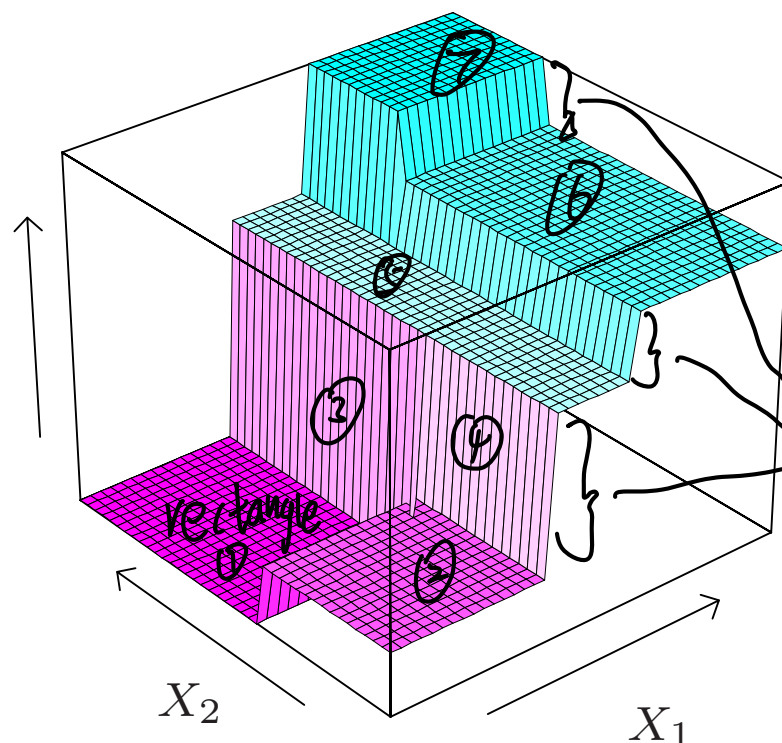
估計到的　(STD)

sum of all of my L regions

$$\text{For all } \mathbf{x} \text{ in feature space, } m(\mathbf{x}) \approx \sum_{\ell=1}^{L} c_\ell I\{\mathbf{x} \in R_\ell\}$$

Indication Function.

where

$$c_\ell I\{\mathbf{x} \in R_\ell\} = \begin{cases} c_\ell & \text{if } \mathbf{x} \in R_\ell \\ 0 & \text{otherwise.} \end{cases}$$

$c_\ell$ constant used in that particular region

$X$ indicator of $X$ in the particular region. 系統

page 306 of Hastie et al., 2017: piecewise constant approximations on regions $R_1, R_2, \ldots, R_5$ of $\mathbb{R}^2$ constructed earlier:

by each of region we approximate the regression curve by a constant

不同的 rectangle 跟不同的 constant

函数叫做的一个意思

height is difference on each rectangle

☛ Why is this flexible? As long as we partition a feature space in small enough pieces, we can always approximate well a regression curve by constants on each piece.

just use constant in each region instead of using sophiscated functions.

☛ The finer the partition of the feature space, the better is the approximation of $m$ by constants $c_1, c_2, c_3, \ldots$ on the regions $R_1, R_2, R_3, \ldots$.

☛ To understand this, here is an example in the case where $p = 1$.

分区分得 越细, regression 就越精准.
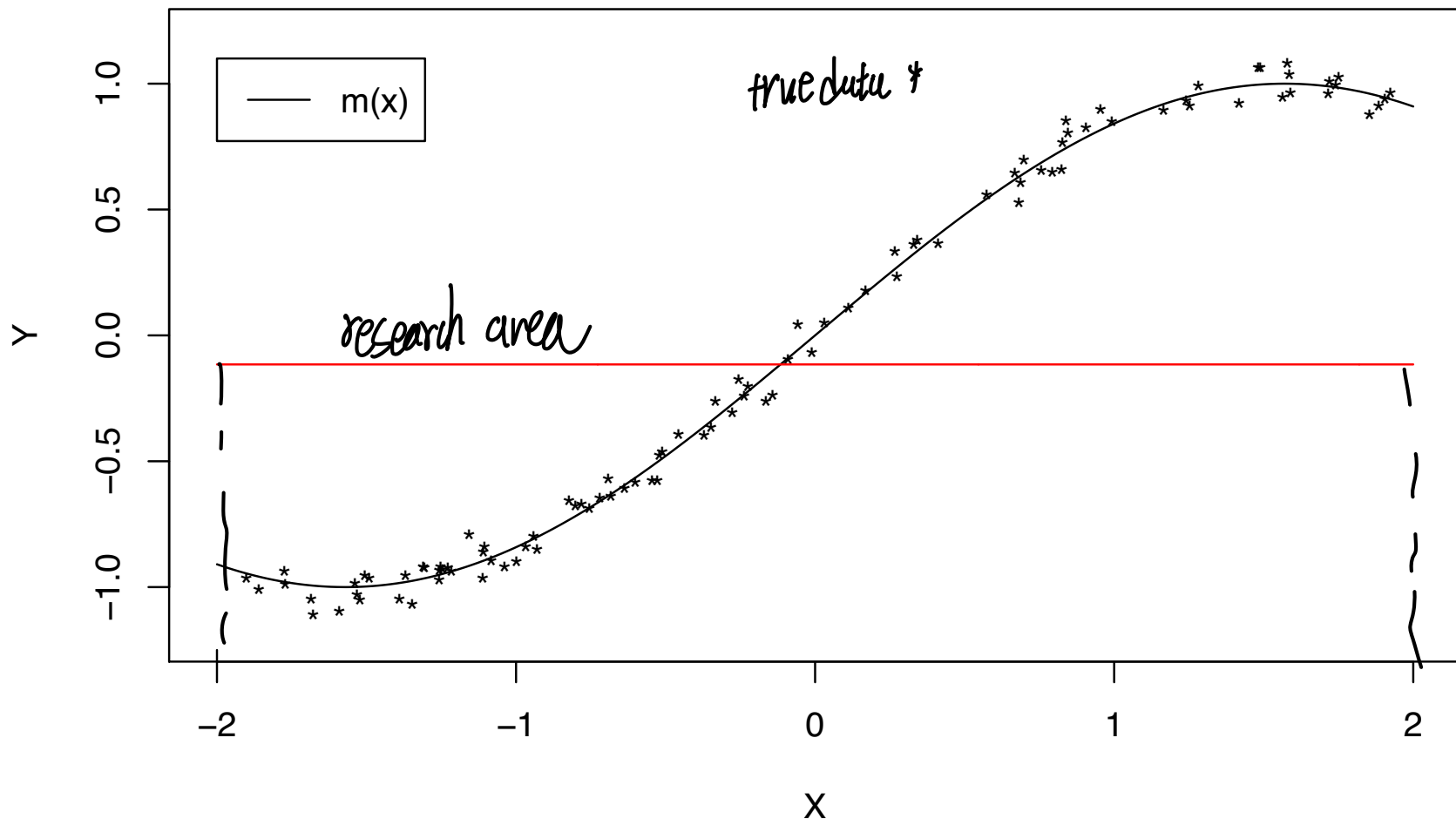
Example with $p = 1$. Here, partitioning the feature space means splitting the range of values of the variable $X$ into intervals.

Data points:∗. Approximate $m$ by a constant (red line) without partitioning feature space does not give a good approximation of $m$:
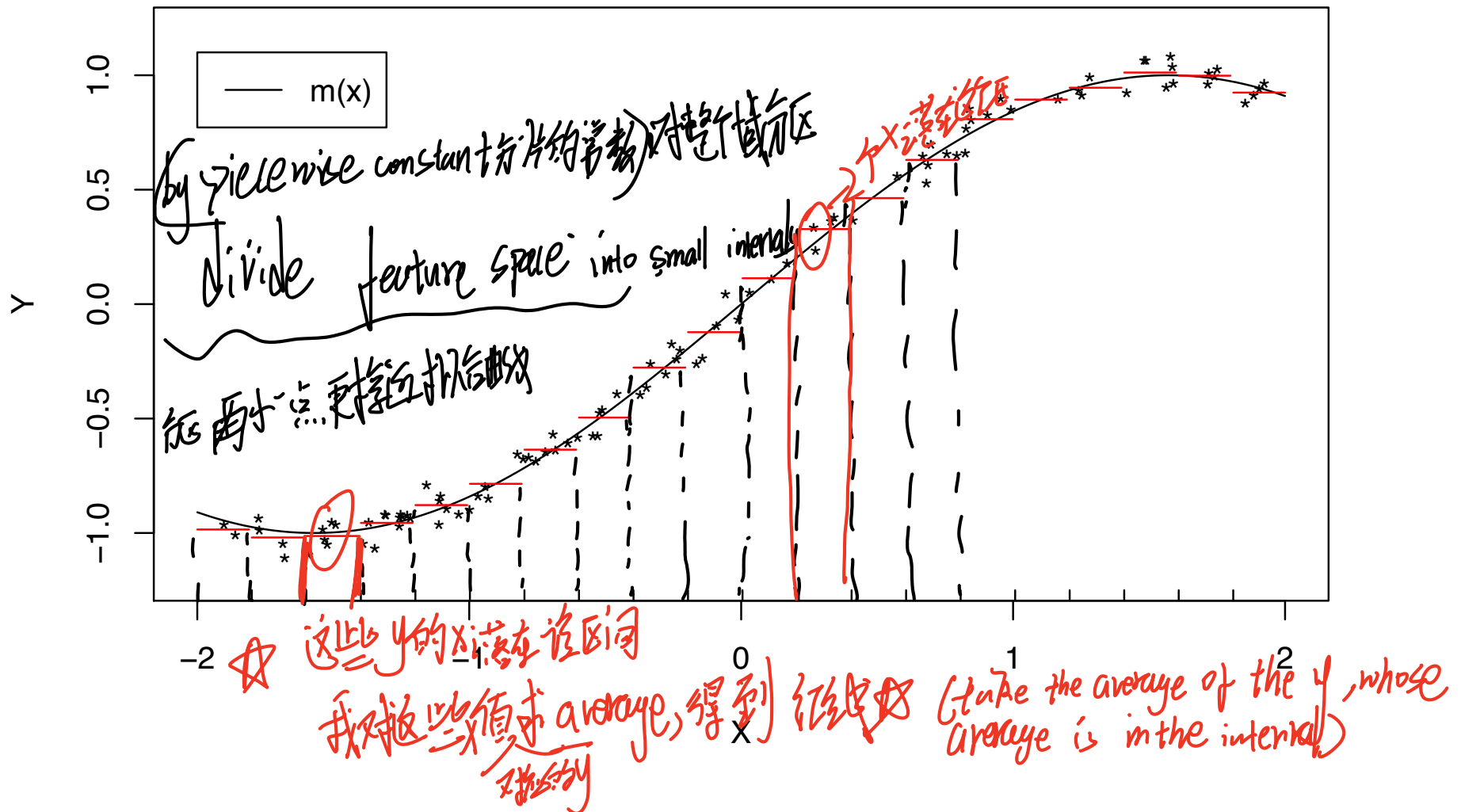
Partition the real line into small intervals and approximate $m$ by a well chosen constant on each interval gives a much better approximation of $m$

**reasonable partition**

常数值 *(in each piece is (best) possible approximation to my curve)*

# FITTING THE TREE (CHOOSING THE CONSTANTS) 分两个

☞ In practice <u>don't know $m$</u> and <u>can't choose $c_j$ that approximates $m$</u>
"the best" on each $R_j$.

☞ Instead, for each $j$, using only the data $(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n)$, approximate "the best" $c_j$ by an estimator $\hat{c}_j$. How?

*把特征空间我成L个*

☞ Suppose we have partitioned the feature space of $\mathbf{X}$ into $R_1, \dots, R_L$.
For $\ell = 1, \dots, L$, for $\mathbf{x} \in R_\ell$, we approximate $m(\mathbf{x})$ by

*(的最好成: estimator of regression curve (常数) ⇒ 在$R_\ell$内的$x_i$对应的$Y$求平均 $Y_i$)*

$$\hat{m}(\mathbf{x}) = \hat{c}_\ell = \text{average of } Y_i's \text{ whose } \mathbf{X}_i \in R_\ell$$

$$= \sum_{i=1}^{n} Y_i \cdot I(\mathbf{X}_i \in R_\ell) / \sum_{i=1}^{n} I(\mathbf{X}_i \in R_\ell).$$ *Sum of indicators*

*(indicator of whether $X_i$ are in the region or not)*

*(分类选等)* ☞ Same as choosing $\hat{c}_j's$ that minimise residual sum of squares (RSS).

*Sum 每个区间*

$$\sum_{i=1}^{n} \left\{ Y_i - \sum_{\ell=1}^{L} c_\ell\, I(\mathbf{X}_i \in R_\ell) \right\}^2 = \sum_{\ell=1}^{L} \sum_{\mathbf{X}_i \in R_\ell} (Y_i - c_\ell)^2.$$

wrt $c_1, \dots, c_L$

*$\hat{m}(X_i)$*

*∑(每个区间的误差作一起计算)*

# BUILDING THE TREE (CHOOSING CONSECUTIVE SPLITS)

☞ Recall that for $\mathbf{X} = (X_1, \ldots, X_p)^T$, a split is of the type $\{X_j \leq t\}, \{X_j > t\}$ for some $1 \leq j \leq p$ and some value $t$.

☞ Ideally, choose splits so that final tree minimises RSS but not computationally possible.

☞ Why? Would involve constructing all possible sequences of splits $\{X_j \leq t\}, \{X_j > t\}$ for all $j = 1, \ldots, p$ and all values of $t$.

☞ Note: in fact, for split on $X_j$, only need to consider $t \in \{X_{1j}, \ldots, X_{nj}\}$.

Why? To find $\hat{c}_\ell$ on $R_\ell$, only need to know $I(\mathbf{X}_i \in R_\ell)$. Changing $t$ for split $\{X_j \leq t\}, \{X_j > t\}$ used in $R_\ell$ affects $I(\mathbf{X}_i \in R_\ell)$ through $I(X_{ij} \leq t)$, whose value only changes when $t$ changes from $t < X_{ij}$ to $t \geq X_{ij}$ (or vice versa).

☞ Even with this simplification, too time consuming to consider all partitions.

about tree: at the end, the entire feature space was divided into a number of regions.

数学情况下，我们只有树没有图 ⇒ 有纪录率 对 ⇒ 一 状树

做 regression tree 时: we approximate a regression curve by piecewise constant.
(ㄴ一区域一常数)

在新区域 ⇒ 图 we approximate regression curve by a constant (one constant per region).

In every small region

对于 每个区域的 constant: (the fit that we use in each region)
= average of $y_i$ (whose $x_i$ belongs to that region ($x_i \in R_j$)) [we need to check which data falls in a particular region. ⇒ we take the average of these $y_i$
of these particular individuals whose $x_i$ was in the region]

选择残差平方和最新的 C (RSS)

$$\sum_j \sum_{(i)\, x_i \in R_j} (y_i - c_j)^2$$

How do we choose the splitting point.        two regions
$$y = (x_1, \cdots, x_p) \Rightarrow \{x_j \le t\}\ \{x_j > t\}$$    (how do we choose $j$ & $t$?)

split is based on residue sum of squares

Choose the tree that can minimize the RSS → 但我们

所以 无法计算大的树 (计算资源)

we use the   Indicator   of variable to fit the tree (x在的区域就分别 不在题)