

ch1 Basis Functions and Piecewise Regression

* Polynomial Regression 多项式回归

Linear regression: $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$ usually don't use degree > 4: overly flexible
 Polynomial regression: $y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^3 + \dots + \beta_d x_i^d + \epsilon_i$

* Step Function: define regions which different models will be fit

→ creating thresholds: $c_1, c_2, \dots, c_k \rightarrow$ intervals / C_k as indicator variables (1 or 0)

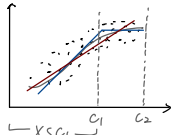
$$y_i = \beta_0 + \beta_1 C_1 x_i + \beta_2 C_2 x_i + \dots + \beta_k C_k x_i + \epsilon_i$$

$$C_1 = 1 \text{ if } c_0 < x < c_1$$

$$C_2 = 1 \text{ if } c_1 < x < c_2$$

$$\vdots$$

$$C_k = 1 \text{ if } c_{k-1} < x < c_k$$



* Regression Splines

① knot selection (K knots)

② piecewise polynomial fitting → cubic splines

$$y_i = \beta_0 + \beta_1 b_1(x_i) + \beta_2 b_2(x_i) + \dots + \beta_{k+1} b_{k+1}(x_i) + \epsilon_i$$

- Add truncated power basis to knot

$$y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \beta_3 x_i^3 + \sum_{j=1}^k \beta_{k+j} (x_i - \xi_j)_+^3 + \epsilon_i$$

$$h(x, \xi) \equiv (x - \xi)_+^3 = \begin{cases} (x - \xi)^3 & \text{if } x > \xi \\ 0 & \text{otherwise} \end{cases}$$

③ add constraints → natural cubic splines (smoother, boundary constraint to be linear)

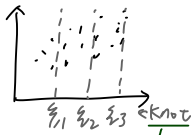
* Smoothing splines: want RSS ↓ + smooth function \hat{f}

$$RSS = \sum_{i=1}^n (y_i - g(x_i))^2$$

$$\rightarrow RSS + \lambda \int g''(x)^2 dx$$

Penalty
 tuning parameter (select: CV, AIC, BIC)

λ small: fit data closely
 λ large: smoother curve



if number of knots $\geq n$, model will overfit → Use Cross Validation

ch2 Tree-based method

* Recursive split (partition) data to minimize ① sum of RSS ② sum of misclassification in new regions

→ branch, leaves, root node

* Steps

① Split data

② Fit full tree → overfitting (down to 1 observ.)

③ Prune tree

* Pruning - reduce overfitting by simplifying model

- Usually post-pruning: explore data thoroughly before decide which branch are redundant

④ Validate pruned tree: for each trees, use k-fold CV to estimate MSE
 ⑤ Average MSE for each k-folds
 ⑥ Select best tree: lowest MSE

Advantage	Disadvantage
Interpretable Close to human reasoning Nice Graph No need to create dummy variables	Not good at predicting (regression) high variance (sensitive to small change in data)

* Regression Tree (greedy algo) Quantitative

→ Select region R_k , predictor X_j , splitting point $s \rightarrow$ splitting R_k with criterion $X_j < s$ (max RSS ↓)

- stop splitting Rule: each region with $5 \downarrow$ observations

- Overfitting problem if 1 point / region.

1. weakest link pruning $\frac{RSS(T_i) - RSS(T_0)}{|T_i| - |T_0|} \downarrow$ (T_0 : original tree)

2. cost complexity pruning $\left(\sum_{i=1}^n (y_i - \hat{y}_{R_m})^2 + \alpha |T| \right) \downarrow$
 choose by CV

* Bagging, Random Forest & Boosting

- Tree → High Variance: small changes / partitions in data can have a big effect on fitted model

① Bagging (Bootstrap Aggregating): variance ↓ (subset with replacement)

→ Bootstrapping: results in similar datasets → highly correlated trees $\hat{f}(x) = \frac{1}{B} \sum_{b=1}^B \hat{f}_b(x)$

② Random forest

→ Bagging + restrict number of predictors, the algo can see when it splits (randomly choose msp predictors)

→ there are 1~2 predictors dominate splitting

→ Allowing models to be fit without dominant variables: the set of B trees is decorrelated

averaging B trees (predictions), $\hat{f}_{me}(x)$ will have lower variance

③ Boosting: Regularization, bias ↓

→ Adding shrinkage tree to every trees sequentially, correcting previous error

$$\hat{f}(x) = \hat{f}(x) + \lambda \hat{f}^b(x) \leftarrow \text{shrinkage tree}$$

$$r_i = r_i - \lambda \hat{f}^b(x) \text{ for all } i$$

$$\hat{f}(x) = \sum_{b=1}^B \lambda \hat{f}_b(x)$$

④ Bayesian Additive Regression Tree (BART)

→ grow trees like Bagging and Random Forest also capture new info like boosting

$\hat{f}_k^b(x)$ is the prediction at k^{th} regression tree (of k) and b^{th} iteration (of B)

$$\hat{f}_k^b(x) = \sum_{i=1}^n \hat{f}_{k,i}^b(x)$$

ch10 Deep Learning "structure" of the model

* Multiple Layer Neural Networks

- Input layer

- Hidden layer(s)

$$A_k^{(1)} = h^{(1)}(x) = g(w_{k0}^{(1)} + \sum_{j=1}^n w_{kj}^{(1)} x_j) \text{ for } k=1, \dots, K_1$$

$$A_k^{(2)} = h^{(2)}(x) = g(w_{k0}^{(2)} + \sum_{j=1}^{K_1} w_{jk}^{(2)} A_j^{(1)}) \text{ for } k=1, \dots, K_2$$

- Output layer $E_m = \beta_{m0} + \sum_{k=1}^{K_2} \beta_{mk} A_k^{(2)}$ last hidden layer

$$= \beta_{m0} + \sum_{k=1}^{K_2} \beta_{mk} A_k^{(2)}$$

$$\text{Softmax } f_m(x) = p_r(Y=m|X) = \frac{e^{z_m}}{\sum_{l=1}^L e^{z_l}}$$

→ we need to have a decision rule. Usually, assign to the class with highest probability

To train, we want to find coefficient estimates such that we minimize cross entropy $-\sum_{i=1}^n \sum_{m=0}^M \log(f_m(x_i))$

* Convolutional Neural Networks (CNN)

① Convolutional Layer: Feature contraction

② Pooling Layer: classification, max pooling = max{each 2x2 subset of image}

$$\begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \\ j & k & l \end{bmatrix} \begin{bmatrix} \alpha & \beta \\ \gamma & \delta \end{bmatrix} \begin{bmatrix} a\alpha + b\gamma + d\delta + e\beta \\ \dots \end{bmatrix}$$

original → convolutional → convolved Image

ReLU

$$g(z) = (z)_+ = \begin{cases} z & \text{if } z > 0 \\ 0 & \text{otherwise} \end{cases}$$

① Flexible: useful when not linear

② Simple and Interpretability

③ Control over complexity: high degree → fit complex pattern

overfitting → small dataset

underfitting → simple model

* Classification Tree Qualitative

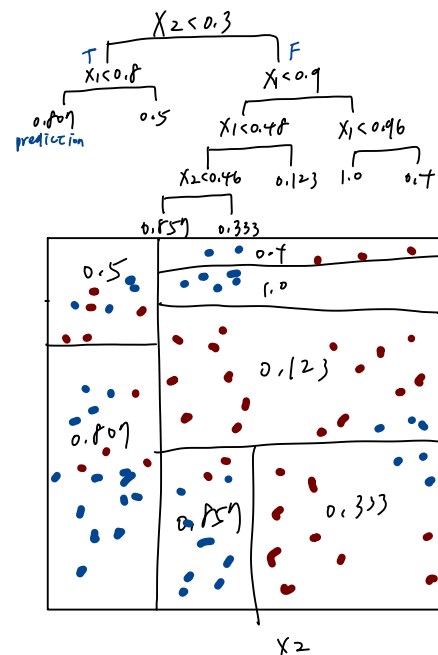
→ minimize Loss Function

For growing tree
 1. 0-1 loss / misclassification rate $\sum_{m=1}^M \sum_{i \in R_m} 1(y_i \neq \hat{y}_{R_m})$

For pruning tree
 2. Gini Index $\sum_{m=1}^M q_m \sum_{k=1}^K \hat{p}_{mk} (1 - \hat{p}_{mk})$
 \hat{p}_{mk} : proportion of class k with R_m
 q_m : proportion of samples in R_m

3. Cross-entropy $-\sum_{m=1}^M q_m \sum_{k=1}^K \hat{p}_{mk} \log(\hat{p}_{mk})$

Bagging	Random Forest	Boosting	BART
variance ↓ high correlated bias ↓ x	variance ↓ correlated ↓ bias ↑ (random)	variance ↓ overfitting bias ↓ slow to train	capture complex non-linear relation slow to train



* Graphics
 data type [Quantitative: continuous, discrete
 Qualitative: ordered, nominal (categories)]

Other:
 - Classification: Confusion Matrix, ROC Curve, Precision-recall curve
 - Access model fitting: Residual plot

Histogram	proportions distribution	hist()	Quant
Density plot	smoothing histogram	plot(density())	Quant
Box plot	statistical info	boxplot()	Qual + Quant
Rug plot	additional distribution	plot(density()) rug()	Quant
Dot chart	comparison values	dotchart()	Qual + Quant
Mosaic plot	proportions across 2 variables	mosaicplot()	Qual + Qual
Scatter plot	relationship between 2 var.	plot()	Quant + Quant
Bar plot	frequency of category	barplot()	Qual

abline() add line
 points() add point
 lines() add line segments
 text() add text
 plot(type="x", p: points, l: lines, n: nothing
 xlim=c(0,1), range for scale of axis
 xlab="...", x axis label
 main="...", plot title
 col=c(), vector of colors
 log="xy", use log scale on y/x/xy
 par(mfrow=c(x,y))
 grid(nx=NULL, ny=NULL, lwd=2) add grid line
 legend("topright", fill=colors, legend=c(A,B,C,D), title="")
 lwd=2, thickness of line
 pch=19, dots type
 cex=0.5, 50% smaller
 lty=2, line type
 las=1 style of tick mark labels

* R code

- Bootstrapping

```
library(boot) # mean & median estimation
set.seed(8436)
results_mean <- boot(data, bootstrap=mean, R=1000)
results_median <- boot(data, bootstrap=median, R=1000)
# confidence interval
ci_mean <- boot.ci(results_mean, type="perc")
ci_median <- boot.ci(results_median, type="perc")
```

- Basic Neural Network

```
input_shape <- c(32, 32, 3)
model_basic <- keras_model_sequential() %>%
  layer_flatten(input_shape)
  layer_dense(units=512, activation='relu')
  layer_dense(units=10, activation='softmax')
model_basic %>% compile(
  loss='categorical_crossentropy',
  optimizer=optimizer_adam(),
  metrics=c('accuracy')
)
history_basic <- model_basic %>% fit(
  x_train, y_train,
  epochs=10,
  batch_size=64,
  validation_split=0.2
)
evaluation_results_basic <- model_basic %>% evaluate(x_test, y_test)
```

```
model_cnn <- keras_model_sequential() %>%
  layer_conv_2d(filters=32, kernel_size=c(3,3), activation='relu', input_shape)
  layer_max_pooling_2d(pool_size=c(2,2))
  layer_conv_2d(filters=64, kernel_size=c(3,3), activation='relu')
  layer_max_pooling_2d(pool_size=c(2,2))
  layer_conv_2d(filters=64, kernel_size=c(3,3), activation='relu')
  layer_max_pooling_2d(pool_size=c(2,2))
  layer_flatten()
  layer_dense(64, 'relu')
  layer_dense(10, 'softmax')
```