



M12: FINAL
08.04.2019

UNSUPERVISED (PCA, clustering $\begin{cases} \nearrow \text{hierarchical} \\ \searrow \text{K-means} \end{cases}$)

REGRESSION

Continuous

$Y_i = f(x_i) + \epsilon_i$ where $E(\epsilon_i) = 0$ and $\text{Var}(\epsilon_i) = \sigma^2$
 and $(x_i, Y_i) \quad i=1, \dots, n$
 independent pairs

Labels: $f(x_i)$ is systematic, ϵ_i is random additive errors.

AIM: estimate $f(x)$

Cost: $(Y - f(x))^2$
squared error

Bias-variance trade-off (test mse at a new x_0)

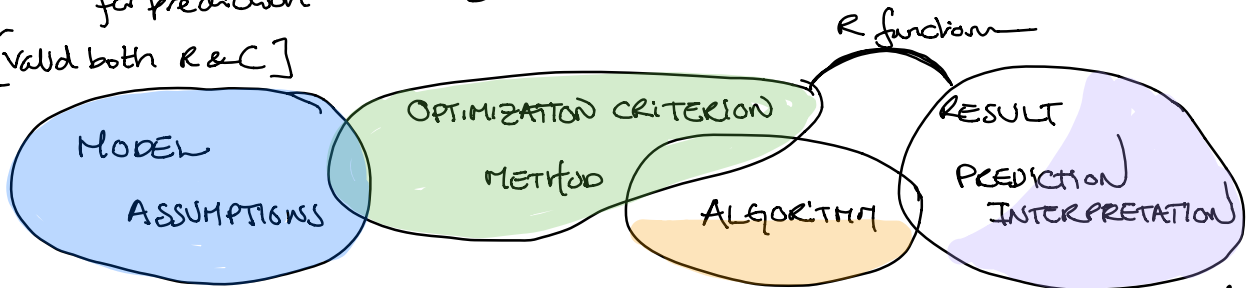
$$E[(Y - \hat{f}(x_0))^2] = \dots = [E(\hat{f}(x_0)) - f(x_0)]^2 + \text{Var}(\hat{f}(x_0)) + \text{Var}(\epsilon)$$

Labels: Y is new obs, $\hat{f}(x_0)$ is at x_0 , $E(\hat{f}(x_0)) - f(x_0)$ is bias², $\text{Var}(\hat{f}(x_0)) + \text{Var}(\epsilon)$ is irreducible error.

Q: how state this orally

- there is random error that you never can fit (deterministic world)
- some times an unbiased answer is the best, but if this has high variance, maybe a biased answer is better for prediction?

[Valid both R & C]



$f(x)$: how complex? and how much data do we have?
 → guide our choice of solution

1) ● $Y_i = \overbrace{X_i^T \beta}^{f(x_i)} + \epsilon_i$ and often $\epsilon_i \sim N$ M3

$$Y = \sum \beta + \epsilon$$

$n \times 1$ $n \times (p+1)$ $(p+1) \times 1$ $n \times 1$

$$\hat{y}_i = X_i^T \hat{\beta}$$

"core model"

linear in parameters

& linear in parameters x

↑
 $f(x)$ is a hyperplane

a) ● $\underset{\beta}{\operatorname{argmin}} \text{RSS}$ train MSE

$\sum_{i=1}^n (y_i - \hat{y}_i)^2 \leftarrow$ least squares estimation M3

● $\underset{\beta}{\operatorname{argmax}} \ell(\beta, \sigma^2)$ maximum likelihood $\epsilon_i \sim N$
 ↑
 $\propto \text{RSS} + C$

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$

$$\hat{Y} = X \hat{\beta} = \underbrace{X (X^T X)^{-1} X^T}_{H} Y = HY$$

b) ● $\underset{\beta}{\operatorname{argmin}} \text{RSS} + \lambda \cdot \sum_{j=1}^p \beta_j^2$

ridge

$$\hat{\beta} = (X^T X + \lambda I)^{-1} X^T Y$$

c) ● $\underset{\beta}{\operatorname{argmin}} \text{RSS} + \lambda \cdot \sum_{j=1}^p |\beta_j|$ lasso

• regularization: add penalty term

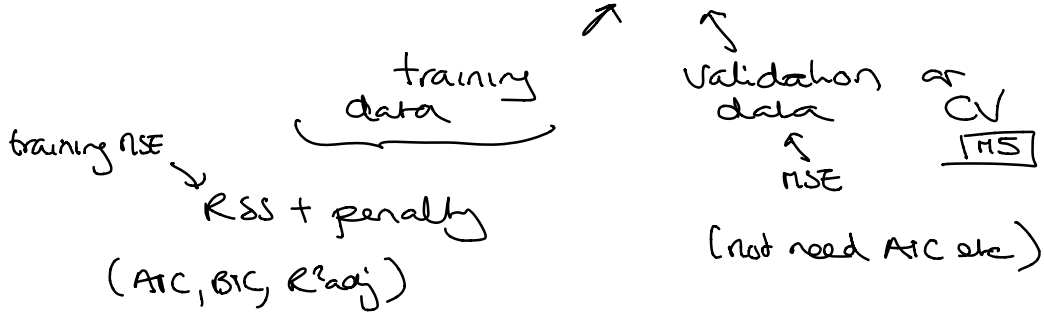
• penalty (hyper) parameter λ

↑
 cross-validation

FIGURES

• dual optimization problem

Alternative to regularization: subset selection MS



2) keep linearity in parameters - but ^{go} nonlinear in covariates.
original

use LS estimation

if [categorical ← dummy variable coding]

$$X_i^T = [1 \quad \textcircled{X_{i1}} \quad X_{i2} \quad \dots \quad X_{ip}]$$

might be an interaction term

focus on X_{i1} , but same for the others
numerical
replace X_{i1} with

polynomial regr.

$$X_{i1} \quad X_{i1}^2 \quad X_{i1}^3 \quad \dots \quad X_{i1}^d$$

step functions

$$I(c_1 \leq X_{i1} < c_2) \quad I(c_2 \leq X_{i1} \leq c_3) \quad \dots$$

$$(X_{i1} - c_k)_+^s = \begin{cases} (X_{i1} - c_k)^s & \text{if } X_{i1} \geq c_k \\ 0 & \text{else} \end{cases}$$

truncated power basis

order of spline with K knots
[$(M-1) + K + 1 \text{ int}$]

$$b_j(X_{i1}) = X_{i1}^j \quad j = 1, \dots, M-1$$

$$b_{M-1+K}(X_{i1}) = (X_{i1} - c_k)_+^{K-1}$$

$$x \quad x^2 \quad x^3$$

$$k = 1, \dots, K \quad \downarrow \text{one per knot}$$

natural cubic splines

$$b_1(x_{i1}) = x_{i1}, \quad b_{K+2}(x_{i1}) = d_K(x_{i1}) - d_{K+1}(x_{i1}) \quad k = 0, \dots, K-1$$

$$c_0, \dots, c_{K+1}$$

$$d_k(x_{i1}) = \frac{(x_{i1} - c_k)_+^3 - (x_{i1} - c_{k+1})_+^3}{c_{k+1} - c_k} \quad 3$$

$k \quad K$

3) Local regression and KNN

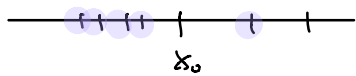
linear in parameters locally

$$\hat{f}(x_0) = \hat{\beta}_0 + \hat{\beta}_1 x_0 + \hat{\beta}_2 x_0^2$$

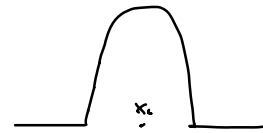
Use $\frac{k}{n}$ closest x_i to x_0 and weigh $K_{i0} = K(x_i, x_0)$ and find $\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2$ by running

$$\sum_{i=1}^n K_{i0} \cdot (y_i - \beta_0 - \beta_1 x_i - \beta_2 x_i^2)^2$$

weighted LS



k default



$$K(x_i, x_0) = \left(1 - \left|\frac{x_0 - x_i}{x_0 - x_{(k)}}\right|\right)_+^3 \begin{cases} (1 - |t|)^3 & |t| \leq 1 \\ 0 & |t| > 1 \end{cases}$$

tricube
(N, Epanechnikov)

KNN: Exam 2018 Prob 1

4) Smoothing spline

$$Y_i = f(x_i) + \epsilon_i$$

nonlinear

$$\min \sum_{i=1}^n (y_i - f(x_i))^2 + \lambda \int f''(t)^2 dt$$

$$\hat{y} = S y \quad \text{where } S = X(X^T X + \lambda I)^{-1} X^T$$

$$\neq X(X^T X)^{-1} X^T Y$$

lot of theory here - not in this course

$$X(X^T X + \lambda I)^{-1} X^T$$

5) Trees

R_1, \dots, R_J non-overlapping regions in predictor space

\hat{y}_{kj} = mean response of all training samples $x_i \in R_j$

Prediction: $\hat{y}_0 = \hat{y}_{kj}$ when $x_0 \in R_j$ "step function"

Find R_j 's to

$$\text{minimize } \sum_{j=1}^J \sum_{i \in R_j} (y_i - \hat{y}_{kj})^2$$

$$\left(\text{alt: } \sum_{i=1}^n (y_i - \hat{y}_{R_j}(x_i))^2 \right)$$

but do this by greedy recursive binary splitting

- Pruning \leftarrow aka lasso penalty
 - Bagging & random forest
 - Boosting
- } algorithms

6) Add linear transformation $\xrightarrow{\text{to MLR}} \text{PCR}$
principal comp. regr

Rename: $Y_i = Z_i^T \beta + \epsilon_i$

$$Z_i = \begin{bmatrix} e_1^T x_i \\ e_2^T x_i \\ \vdots \\ e_m^T x_i \end{bmatrix} \text{ where } S = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T$$

and $S e_i = \lambda_i e_i$

(x_i may be scaled also, so correlation and not covariance matrix used)

7) Neural networks

$$\hat{y}_1(x_i) = \beta_{01} + \sum_{m=1}^m \beta_{m1} \cdot z_m(x_i)$$

$$z_m(x_i) = \phi_n \left(\alpha_{0m} + \sum_{j=1}^p \alpha_{jm} \cdot x_{ij} \right)$$

$$\phi_n(a) = \max(a, 0) \quad \text{or} \quad \phi_n(a) = \frac{1}{1 + \exp(-a)}$$

And, may add more ϕ_n 's inside the ϕ_n 's =
more layers of the NN.

As before:
$$\operatorname{argmin}_{\alpha, \beta} \sum_{i=1}^n (y_i - \hat{y}_1(x_i))^2$$

and may add regularization of various types

CLASSIFICATION

$$Y = \{1, 2, \dots, J\} \text{ or } Y = \{0, 1\}$$

Observe (x_i, Y_i) jointly, $i = 1, \dots, n$
independent pairs

1) Bayes classifier (Utopia)

assign a new obs x_0 to the most likely class

$$P(Y=j | X=x_0) \leftarrow \text{this will minimize the } 0/1 \text{-loss}$$

Error rate on obs x_0

$$1 - \max_j P(Y=j | X=x_0)$$

Bayes error rate = optimal error rate

"we know $P(Y=j | X)$ "

$$= 1 - E_X \left(\max_j P(Y=j | X) \right)$$

↑
compare to irreducible error



if you go lower \rightarrow you prob. have overfitted

2) Two paradigms:

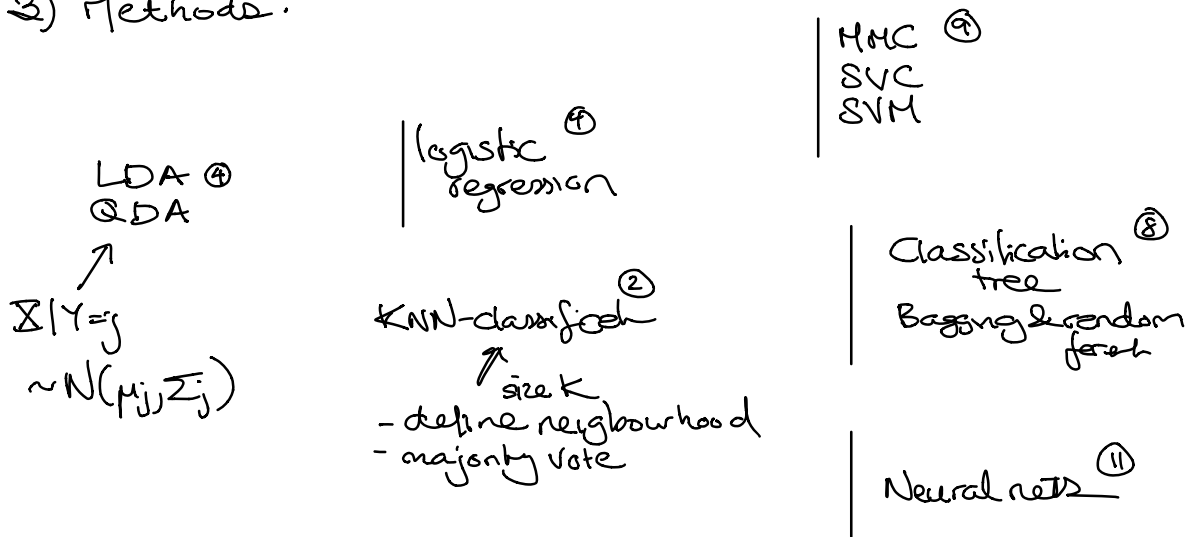
Diagnostic: model $P(Y=j | X)$ directly

Sampling: model $P(X | Y=j)$ and $P(Y=j)$ and use

Bayes theorem to get

$$P(Y=j | X=x) = \frac{P(X=x | Y=j) P(Y=j)}{P(X=x)}$$

3) Methods:



4)

Logistic regression [two-class problem]

$$P(Y=1 | x_i) \rightarrow p_i = \frac{1}{1 + \exp(-x_i^T \beta)}$$

$$\prod_{i=1}^n p_i^{y_i} (1-p_i)^{1-y_i}$$

$$\underset{\beta}{\operatorname{argmax}} \ell(\beta) = \underset{\beta}{\operatorname{argmax}} \sum_{i=1}^n (y_i \log p_i + (1-y_i) \log(1-p_i))$$

odds \rightarrow FRASE

5) Classification tree :

R_1, \dots, R_J non-overlapping regions in predictor space

$\hat{y}_j =$ majority vote in R_j or $\hat{p}_{jk} = n_{jk} / N_j$
when $x_0 \in R_j$
↑ train in region
↑ train of class k in region j

Find R_j 's to minimize

$$G = \sum_{k=1}^K \hat{p}_{jk} \cdot (1 - \hat{p}_{jk})$$

Gini node impurity

$$D = - \sum_{k=1}^K \hat{p}_{jk} \cdot \log \hat{p}_{jk}$$

cross entropy (closely to deviance)

but do this by greedy recursive binary splitting

- Pruning ← aka lasso penalty
- Bagging & random forest } algorithms

6) SVM support vector machine $y_i \in \{-1, 1\}$

$$\max_{\beta, \alpha, \epsilon} \gamma \quad \text{subject to} \quad \sum \beta_j^2 = 1, \quad \sum \epsilon_i \leq C, \quad \epsilon_i \geq 0$$

\swarrow slack \nwarrow budget
 \swarrow \nwarrow

$$y_i \cdot f(x_i) \geq \gamma(1 - \epsilon_i)$$

$$f(x) = \beta_0 + \sum_{i \in S} \alpha_i \cdot K(x, x_i)$$

\nwarrow kernel

alternative way: when $f(x_i) = \beta_0 + x_i^T \beta$ (SVC)

$$\min_{\beta} \left[\sum_{i=1}^n \underbrace{\max(0, 1 - y_i(\beta_0 + x_i^T \beta))}_{\text{hinge loss}} + \lambda \sum \beta_j^2 \right]$$

compared to logistic regression with $(-1, 1)$

$$l(\beta) \propto \sum_{i=1}^n \log(1 + e^{-y_i(\beta_0 + x_i^T \beta)})$$

\uparrow add ridge penalty
if needed

F) Neural networks

$$\hat{y}_c(x) = \phi_0(\beta_{0c} + \sum_{m=1}^M \beta_{mc} z_m)$$

$$z_m = \phi_h(\alpha_{0m} + \sum_{j=1}^r \alpha_{jm} x_j)$$

one hidden layer
 may add more

$$\phi_0(a) = \frac{1}{1 + \exp(-a)}$$

c=2: binary cross-entropy

$$J(\theta) = -\frac{1}{n} \sum [y_i \ln \hat{y}_i(x_i) + (1-y_i) \ln (1-\hat{y}_i(x_i))]$$

$$\phi_0(a_j) = \frac{\exp(a_j)}{\sum_{j=1}^c \exp(a_j)} \quad \text{softmax} \quad \left(\text{out } \frac{1}{2}\right)$$

c > 2: categ. c-2.

$$J(\theta) = -\frac{1}{n} \sum \frac{1}{c} \sum (y_{ic} \cdot \ln \hat{y}_{ic}(x_i))$$

[log vs ln]