

LINEAR REGRESSION (M2)

$Y = \beta_0 + \beta_1 X + \epsilon$
classical normal regression model
 $\epsilon \sim N(0, \sigma^2)$
OLS estimator: $\hat{\beta} = (X^T X)^{-1} X^T Y$
 $\hat{\sigma}^2 = \frac{1}{n-2} RSS$

Confidence interval for β_j :
 $\hat{\beta}_j \pm t_{\alpha/2, n-2} \sqrt{\hat{\sigma}^2 (X^T X)^{-1}_{jj}}$

Hypothesis test:
 $H_0: \beta_j = 0$ vs $H_1: \beta_j \neq 0$
Test statistic: $F = \frac{R^2 / (k-1)}{RSS / (n-k)}$

MODEL SELECTION (M6)

Model selection metrics:
AIC = $-2 \log \text{likelihood} + 2 \text{ (number of parameters)}$
BIC = $-2 \log \text{likelihood} + \log(n)$

When model selection is done on test set or with CV, then AIC, BIC, MSE can be used.

Ridge regression: $\text{minimize } RSS + \lambda \sum \beta_j^2$
Lasso regression: $\text{minimize } RSS + \lambda \sum |\beta_j|$

MODEL REGULARIZATION (M6)

Ridge regression: $\text{minimize } RSS + \lambda \sum \beta_j^2$
Lasso regression: $\text{minimize } RSS + \lambda \sum |\beta_j|$

Partial least squares (PLS): $\text{minimize } \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \int \hat{g}(x)^2 dx$

MOVING BEYOND LINEARITY (M7)

Polynomial regression: $b(x) = I(x \in X < C_{k+1})$

Regression splines: combine polynomials & steps at knots

Smoothing splines: $\text{minimize } \sum_{i=1}^n (y_i - g(x_i))^2 + \lambda \int g'(x)^2 dx$

BIAS-VARIANCE trade-off in regression settings

Expected test mean squared error at x_0 :
 $E[(Y - \hat{f}(x_0))^2] = \text{Bias}^2 + \text{Variance}$

Plot of bias and variance vs model complexity:
- Bias decreases as model complexity increases.
- Variance increases as model complexity increases.

RESAMPLING METHODS (M5)

DATA rich situation (sometimes, e.g. when we generate data ourselves):

TRAIN	VALIDATE	TEST
-------	----------	------

Bootstrap: B_1, \dots, B_m samples of size n from X .

Conduct: $Z = \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{\{f(B_i) > \hat{f}\}}$

RESAMPLING METHODS (M5)

Partial least squares: Competitor to PCR, but now Z is found also taking Y into account. For details in M870L.

Local regression: smoothed version of KNN regression

Additive model: put it all together

RESAMPLING METHODS (M5)

How to represent CV: All tests of model fitting must be on test data unless \rightarrow selection bias. Eg: cross predicted perfectly from Sigmas example.

How to represent CV: All tests of model fitting must be on test data unless \rightarrow selection bias. Eg: cross predicted perfectly from Sigmas example.

CLASSIFICATION (M4)

Bayes classifier: classify to the class with the highest probability
 $P(Y = k | X = x)$

Sampling paradigm: LDA & QDA
LDA: $f(x) = \mu_k$ if $x \in C_k$
QDA: $f(x) = \mu_k$ if $x \in C_k$

CLASSIFICATION (M4)

Bayesian classification: $P(Y = k | X = x) = \frac{P(x | Y = k) \pi_k}{\sum_{k=1}^K P(x | Y = k) \pi_k}$

Sampling paradigm: estimate π_k and $f_k(x)$ and classify to largest $\pi_k f_k(x)$

Diagnostic paradigm: directly estimate $P(Y = k | X = x)$

TREE-BASED METHODS (M8)

Both for regression & classification - and we have nonlinear hypothesis between covariates!

Decision tree: $\text{minimize } \sum_{i=1}^n \min_{k \in \{1, \dots, K\}} |y_i - k|$

Random Forest: $\text{minimize } \sum_{i=1}^n \min_{k \in \{1, \dots, K\}} |y_i - k|$

TREE-BASED METHODS (M8)

Decision tree: $\text{minimize } \sum_{i=1}^n \min_{k \in \{1, \dots, K\}} |y_i - k|$

Random Forest: $\text{minimize } \sum_{i=1}^n \min_{k \in \{1, \dots, K\}} |y_i - k|$

DIAGNOSTIC PARADIGM: KNN & LOGISTIC REGRESSION (trees & SVM too)

KNN: $\hat{f}(x) = \frac{1}{K} \sum_{j \in \text{neighbors}(x)} y_j$

Logistic regression (k=2): $Y_i \in \{0, 1\}$

SVM: $\text{minimize } \frac{1}{2} \sum_{i=1}^n \max(0, 1 - y_i f(x_i)) + \lambda \sum_{i=1}^n \beta_i^2$

SUPPORT VECTOR MACHINES (M9)

A method both for regression and classification, but we only consider classification (high dimensional). X 's are standardized and two classes preferred!

Goal: find hyperplane that separates (perfectly) the two classes coded as 0 and 1.

Margin: $M = \frac{1}{\|\beta\|} \min_{i=1, \dots, n} |y_i - \beta^T x_i|$

TREE-BASED METHODS (M8)

Variable importance plots: G_{ini}

Boosting: $f(x) = \sum_{t=1}^T f_t(x)$

Random Forest: $\text{minimize } \sum_{i=1}^n \min_{k \in \{1, \dots, K\}} |y_i - k|$

TREE-BASED METHODS (M8)

Boosting: $f(x) = \sum_{t=1}^T f_t(x)$

Random Forest: $\text{minimize } \sum_{i=1}^n \min_{k \in \{1, \dots, K\}} |y_i - k|$

SUPPORT VECTOR MACHINES (M9)

Support vector classifier: non-separable case - ϵ_1, ϵ_2 slack variables

Classification rule: $f(x) = \beta^T x + \beta_0$ and if $f(x) > 0$ set $y = +1$, if $f(x) < 0$ set $y = -1$

Support vector machines: $f(x) = \sum_{i \in S} K(x_i, x)$

SUPPORT VECTOR MACHINES (M9)

Support vector classifier: non-separable case - ϵ_1, ϵ_2 slack variables

Classification rule: $f(x) = \beta^T x + \beta_0$ and if $f(x) > 0$ set $y = +1$, if $f(x) < 0$ set $y = -1$

Support vector machines: $f(x) = \sum_{i \in S} K(x_i, x)$

NEURAL NETWORKS (M11)

Feedforward network: $Y = \sum_{i=1}^n \lambda_i f(x_i)$

Hidden layer: $Y = \sum_{i=1}^n \lambda_i f(x_i)$

Output layer: $Y = \sum_{i=1}^n \lambda_i f(x_i)$

NEURAL NETWORKS (M11)

Hidden layer: $Y = \sum_{i=1}^n \lambda_i f(x_i)$

Output layer: $Y = \sum_{i=1}^n \lambda_i f(x_i)$

UNSUPERVISED LEARNING (M10)

Principal component analysis (PCA): $\text{minimize } \sum_{i=1}^n \sum_{j=1}^k (x_{ij} - \hat{x}_{ij})^2$

K-means clustering: $\text{minimize } \sum_{i=1}^n \sum_{k=1}^K \|x_i - \mu_k\|^2$

Hierarchical clustering: $\text{minimize } \sum_{i=1}^n \sum_{j=1}^n \|x_i - x_j\|^2$

UNSUPERVISED LEARNING (M10)

K-means clustering: $\text{minimize } \sum_{i=1}^n \sum_{k=1}^K \|x_i - \mu_k\|^2$

Hierarchical clustering: $\text{minimize } \sum_{i=1}^n \sum_{j=1}^n \|x_i - x_j\|^2$

UNSUPERVISED LEARNING (M10)

K-means clustering: $\text{minimize } \sum_{i=1}^n \sum_{k=1}^K \|x_i - \mu_k\|^2$

Hierarchical clustering: $\text{minimize } \sum_{i=1}^n \sum_{j=1}^n \|x_i - x_j\|^2$

UNSUPERVISED LEARNING (M10)

K-means clustering: $\text{minimize } \sum_{i=1}^n \sum_{k=1}^K \|x_i - \mu_k\|^2$

Hierarchical clustering: $\text{minimize } \sum_{i=1}^n \sum_{j=1}^n \|x_i - x_j\|^2$