

LINEAR REGRESSION (M2) MODEL SELECTION (M6) MODEL REGULARIZATION (M6) MOVING BEYOND LINEARITY (M7)

LINEAR REGRESSION (M2)
 $Y = X\beta + \epsilon$
 $\beta = (X^T X)^{-1} X^T Y$
 $\sigma^2 = \frac{1}{n-1} RSS$
 Hypothesis test: $H_0: \beta_j = 0$ vs $H_1: \beta_j \neq 0$
 Residuals: $\epsilon_i = y_i - \hat{y}_i$

MODEL SELECTION (M6)
 can be used in model selection using 'only' the training data
 $R^2 = 1 - \frac{RSS}{TSS}$
 $AIC = -2 \log likelihood + 2 \text{ (number of param)}$
 $BIC = -2 \log likelihood + \log(n)$
 When model selection is done on test set or with CV, then R^2, AIC, BIC can be used.

MODEL REGULARIZATION (M6)
 Ridge (L2) regression: minimize $RSS + \lambda \sum \beta_j^2$
 Lasso (L1) regression: minimize $RSS + \lambda \sum |\beta_j|$
 Interpretation of plots: ridge vs lasso

MOVING BEYOND LINEARITY (M7)
 $Y = f(x) + \epsilon$
 Polynomial regression: $f(x) = \sum_{k=0}^d \beta_k x^k$
 Regression splines: combine polynomials & steps at knots
 Natural cubic spline = cubic spline that is linear at ends

BIAS-VARIANCE trade-off in regression settings:
 Expected test mean squared error at x_0
 $E[(Y - \hat{f}(x_0))^2] = \text{Var}(\epsilon) + \text{Bias}(\hat{f}(x_0))^2 + \text{Var}(\hat{f}(x_0))$
 Plot of this: model complexity vs error

RESAMPLING METHODS (M5)
 DATA rich situation (sometimes, e.g. when we generate data ourselves)
 TRAIN: fit model, VALIDATE: model selection, TEST: model assessment
 Bootstrap: data set of n observations. Estimate f by \hat{f} .
 Draw random sample from f → draw with replacement from data → one bootstrap sample.
 The probability that x is in bootstrap sample is $1 - (1 - p)^n \approx np$

CLASSIFICATION (M4)
 (Y_i, X_i) with $Y_i \in \{1, \dots, K\}$ or $Y_i \in \{0, 1\}$
 $P(Y = k | X = x)$ posterior class probability
 $P(Y = k)$ prior class probability
 $f(x) = P(Y = k | X = x)$ is pdf of X in class k conditional distribution

CLASSIFICATION (M4)
 Bayesian classifier: classify to the class with the highest probability
 $P(Y = k | X = x)$ → end in simulation experiments this is known.
 Bayesian-class boundary: boundary of Bayes classifier.
 Bayes error = error rate of Bayes classifier.

Tree-based Methods (M8)
 Both for regression & classification - and we have nonlinear hedges & interactions between covariates!
 Want to minimize an error criterion on the whole tree construction - but since computationally infeasible → greedy approach "recursive binary splitting" is used.
 $R_L(y) = 1(x < c)$, $R_R(y) = 1(x \geq c)$: find predictor f_j and splitting point c that minimize "ess"
 Regression: $\sum_{i: x_i < c} (y_i - \hat{y}_j)^2 + \sum_{i: x_i \geq c} (y_i - \hat{y}_j)^2$
 Classification: similar, but using impurity measure
 Gini index: $1 - \sum_{k=1}^K p_k^2$ with N_j obs.
 Cross entropy: $D = -\sum_{k=1}^K p_k \log(p_k) = -\sum_{k=1}^K \left(\frac{N_k}{N}\right) \log\left(\frac{N_k}{N}\right)$
 Deviance: $-2 \sum_{k=1}^K n_{jk} \log(\hat{p}_{jk}) = -2 \sum_{k=1}^K n_{jk} \log\left(\frac{N_{jk}}{N_j}\right)$

Support Vector Machines (M9)
 A method both for regression and classification, but we only consider classification, standardized and two classes preferred!
 Aim: find hyperplane that separates (perfectly) the two classes
 $p_0 + x^T p = 0$ with normalized p
 $p_0 + x^T p > 0$ one side of hyperplane, < 0 the other side
 Maximal margin classifier: $\max_{p, p_0} \frac{1}{\|p\|} \text{ subject to } y_i (p_0 + x_i^T p) \geq 1$
 $y_i (p_0 + x_i^T p) \geq 1$ if correct classified x_i

Diagnosis paradigm: KNN & logistic regression (trees & SVM too)
 KNN: $\hat{A}(x_j | X = x_0) = \frac{1}{K} \sum_{i \in N(x_0)} I(y_i = j)$
 Logistic regression ($k=2$): $Y_i \in \{0, 1\}$ with prob π_i
 Prediction at x_0 : $\hat{p}(x_0) = \frac{\exp(x_0^T \beta)}{1 + \exp(x_0^T \beta)}$

Support Vector Machines (M9)
 Connection SVM & logistic regression: hinge loss vs ridge penalty
 Support vector classifier: non-separable case - ϵ_1, ϵ_2 slack variables
 $\max_{\beta, p_0, \epsilon_1, \epsilon_2} \frac{1}{\|p\|} \text{ subject to } y_i (p_0 + x_i^T p) \geq 1 - \epsilon_i, \epsilon_i \geq 0, \sum \epsilon_i \leq C$
 C is tuning parameter chosen by CV

Neural Networks (M10)
 Look for underlying structure or groupings in data - no Y only X
 Principal component analysis (see M6)
 PC loadings: interpret effect of each covariate on each component
 PC score: $(\text{plot observations in } PC_1 \text{ & } PC_2)$ can be used for quality control
 PC score → MLR = PCR (M6)

Neural Networks (M10)
 Possible to represent MLR & logistic regression as graph (one input and one output layer)
 $Y = p_0 + p_1 x_1 + p_2 x_2 + \dots + p_n x_n + \epsilon$
 $\log\left(\frac{p}{1-p}\right) = p_0 + p_1 x_1 + \dots + p_n x_n$
 $p = \frac{\exp(p_0 + p_1 x_1 + \dots + p_n x_n)}{1 + \exp(p_0 + p_1 x_1 + \dots + p_n x_n)}$
 Instead of M7 where we'd add nonlinear function of each covariate, we instead look at nonlinear function of sums of covariates pr. layer - and add many layers. Popular non-linear activation function is $\text{relu}(x) = \max(0, x)$.
 Choice of activation function = feedforward network.
 Tuning parameters: λ = shrinkage parameter, d = number of tree splits

UNSUPERVISED LEARNING (M10)
 Look for underlying structure or groupings in data - no Y only X
 Hierarchical clustering: work in a sequential way by connecting observations that are similar
 Similarity measures: Euclidean (distance), Correlation
 Linkage: how to calculate dissimilarity between groups of observations?
 - single: minimum
 - average: average
 - complete: maximum
 Presented in a dendrogram - choose where to cut the dendrogram to get a number of clusters.