# Chapter 1

- Stat learning – vast set of tools for understanding data
- Supervised learning – involves building a model for predicting an output based on 1+ inputs
- unsupervised learning – there are inputs but no supervising output → we try to learn patterns from the data
- regression problem – continuous / quantitative output
- classification problem – categorical or qualitative output
- clustering problem – find groups based on observed characteristics of inputs
- & PICK METHODS based on response!

# Chapter 5

- Resampling – involve repeatedly drawing samples from training set
  → model assessment: evaluating a model's performance
  → model selection: selecting proper level of flexibility
  # Bootstrap – accuracy of parameters (non parametric)
  # CV = estimate test error associated or level of flexibility

- Validation set approach: randomly split into 2 sets
  → different results depend on different splits
  → validation estimate of test error highly variable
  → only trained w/ subset = over estimates test error

- Leave-One-Out-CV (LOOCV): single observation used for validation
  → unbiased estimated of test error
  → $CV_n = \frac{1}{n}\sum_{i=1}^{n} MSE$  tends not to overestimate test error
  → advantage = less bias than validation, always yields same results
  → disadvantage = time consuming, $CV_n = \frac{1}{n}\sum \left(\frac{y_i - \hat{y}}{1 - h_i}\right)$ least squares

- k fold = randomly divide into k-folds, fit k-1 folds
  $CV_k = \frac{1}{k}\sum MSE$, LOOCV = special case k=n, adv: computations

Validation Set      k-fold      LOOCV
                    bias        variance

# Chapter 3 = Regression, supervised learning

SLR = $Y \approx B_0 + B_1 X \Rightarrow \hat{y} = \hat{B}_0 + \hat{B}_1 X$ (least squares line)
→ least squares criterion: $e_i = y_i - \hat{y}_i$ (measuring closeness)
→ RSS: $e_1^2 + e_2^2 + ... + e_n^2 = (y_1 - \hat{B}_0 - \hat{B}_1 x) + ... + (y_n - B_0 - B_1 \hat{x}_n)$
→ least square approach, minimize RSS: $\hat{B}_1 = \frac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{(x_i - \bar{x})^2}$
→ Standard Error:
$SE(\hat{B}_0)^2 = \sigma^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{\sum(x_i - \bar{x})^2}\right]$, $Var(\hat{u}) = SE(\hat{u})^2 = \frac{\sigma^2}{n}$  standard deviation
$SE(\hat{B}_1)^2 = \frac{\sigma^2}{\sum_{i=1}^{n}(x_0 - \bar{x})^2}$ | $\sigma^2 = var(\epsilon)$, $\sigma$, $RSE = \sqrt{\frac{RSS}{n-2}}$

Confidence Interval = 95% CI = range of values such 95% probabilities that the range will contain true unknown value
→ in LR: $\hat{B} \pm 2 \cdot SE(B) | H_0: B_1 = 0$ | t-stat: $t = \frac{\hat{B}_1 - 0}{SE(\hat{B})}$
Hypothesis Testing | $Ha = B_1 \neq 0$ | p-value: probability of $x \geq |t|$
→ small p value = reject null hypothesis
Model Accuracy = model fit → $RSE = \sqrt{\frac{RSS}{n-2}} = \sqrt{\frac{1}{n-2}\sum(y_i - \hat{y}_i)^2}$
RSE (lack of fit): deviation from regression line.
$R^2$: independent from scale of Y, measured in proportion (explained)
$R^2 = \frac{TSS - RSS}{TSS} = \frac{ESS}{TSS}$ | $TSS = \sum(y_i - \bar{y})^2$  0: does not explain / 1: has explained | explained

MLR: why | not run separate SLR: ① single predictors? ② predictors ignore others!
minimize $RSS = \sum(y_i - \hat{y}_i)^2$
$\hat{Y} = \hat{B}_0 + \hat{B}_1 X_1 + \hat{B}_2 X_2 + ... + \hat{B}_p X_p$:
$H_0: B_1 = B_2 = B_p = 0$, $Ha:$ at least 1 $B_i$ is non zero | p value of
F stat: $\frac{TSS - RSS/P}{RSS - (n-p-1)}$ } $H_0 \Rightarrow f \approx 1$ } large n 1+ F stat important!
→ $R^2$ always increases when P increases, $RSE = \sqrt{\frac{RSS}{n-p-1}}$ } $Ha \Rightarrow f > 1$ } small n 1++
models can have in RSS | predictions: uncertainty ② model bias larger
with a higher is small relative to ① ↑ estimate! coefficient → assuming than
RSS, it increase in P in accuracy shape of
more values decrease → confidence intervals function

# Chapter 2

$Y = f(X) + \epsilon$  random error term $E(Y - \hat{y})^2 = E[f(x) + \epsilon - \hat{f}(x)]^2$
/independent of X $= [f(x) - \hat{f}(x)]^2 + var(\epsilon)$
systematic information (goal!)  mean (0)  reducible error | irreducible
→ linear → better for inference
→ non-linear → better for prediction

- Parametric Methods: 2 step model based approach ① make assump flexibility  increasing flexibility
  → reduces problem of estimating f down to estimating set of parameters. ② estimate parameters ↑ increases complexity → overfit
  advantage: simplifying problem, disadvantage: will not match the true the function  estimate can be poor! Learning the noise
- Non-Parametric Method: do not make explicit assumptions about functional forms, instead they choose to estimate f that gets as close to data points as possible.
  advantage: potential to accurately fit wider range of fs. disadvantage = large # of n needed to obtain estimate
  lower level of smoothness = better model fit = overfitting

Model Accuracy:
→ regression setting – quantify the extent to which predicted response is close to the response
$MSE = \frac{1}{n}\sum_{i=1}^{n}(y_i - \hat{f}(x_i))^2$  | Bias-Variance Trade-off:  test MSE cannot be lower than
→ test --- irreducible error | $E(y_0 - \hat{f}(x_0))^2 = Var(\hat{f}(x_0)) + (Bias(\hat{f}(x_0)))^2 + Var(\epsilon)$
→ train  min test MSE  → lowest Test MSE = low variance + low bias (ideal)
flexibility  high variance → small changes in training data lead to large in f, higher = more flex
- variance = amount by which $\hat{f}$ changes if we use different training data
- bias = error that is introduced by approx. real world problem w/ simpler model, lower = more flex
→ more flexibility = higher variance, lower bias  increases but bias decreases faster  MSE
test error/variance  as flexibility increases, initially variance increases but bias has little impact & variance increases faster, increases!
bias  so MSE falls but after a while bias has little impact & variance increases faster, increases!
flexibility  $\frac{1}{n}\sum_{i=1}^{n} I(y_i \neq \hat{y}_i)$, $I = 0$ or 1 $x$ {function of correct/incorrect classification} → test error rate: $Ave(I(y_0 \neq \hat{y}_0))$ is smallest!
→ classification setting ⇒ assign observation to the most likely class, given its predictor values:
Bayes Classifier: test error minimized when $Pr(Y = j | X = 0)$ is largest  lowest possible test error rate:
→ Bayes Decision Boundary = probability is exactly 50%  Bayes error rate = $1 - E(\max_j Pr(Y = j | X = x_0))$ = irreducible error
ideal scenario because conditional probabilities  ① identifies neighbors
$KNN = Pr(Y = j | X = x_0) = \frac{1}{K}\sum_{i \in N} I(y_i = j)$ → ② estimate conditional probability for class j
① MSE over f & k is small  ③ applies Bayes rule  CV in classification
④ classifies w/ highest prob  $CV_{(n)} = \frac{1}{n}\sum(I(y_i \neq \hat{y}_i))$
Bootstrap: unknown, SE, w/ a given estimate
$SE_B(\hat{a}) = \sqrt{\frac{1}{B-1}\sum_{r=1}^{B}\left(\hat{a}^r - \frac{1}{B}\sum_{r=1}^{B}\hat{a}^{r'}\right)^2}$  ⑦ binary setting better but can lead to unrealistic p(x)

# Chapter 4: Classification (supervised) ⇒ why not LR: ① no natural order Y belongs to a class rather than class directly.

Logistic Regression: models the probabilities that
$p(X) = \frac{e^{B_0 + B_1 X}}{1 + e^{B_0 + B_1 X}}$ (logistic function) $\Rightarrow \left[\frac{p(X)}{1 - p(X)}\right] = odds = e^{B_0 + B_1 X}$ $B_1$ gives
0 = low prob → change in
$\Rightarrow \log\left(\frac{p(x)}{1 - p(x)}\right) = \log odds = B_0 + B_1 X$ ∞ = high prob  log odds.
maximum likelihood function:  $p(X)$ depends on current value
we seek $B_0 + B_1$ so $\hat{p}(x_i)$ for classes correspond as close as possible to observed!
$l(B_0, B_1) = \prod_{i: y_i = 1} p(x_i) \prod_{i': y_{i'} = 0} (1 - p(x_{i'}))$ ⇒ accuracy of coefficient  $Z = Stat = Z(B_0) = \frac{B_1}{SE(B_1)}$
⇒ confounding: correlation between predictors.  → large z, small p, reject $H_0$
$H_0: \frac{e^{B_0}}{1 + B_0}$

LDA: multiple classes → less direct approach: models distributions of predictors in each class, Bayes Theorem to flip into $Pr(Y = k | X = x)$
→ why not log reg: classes not well separated, n is small, more than 2 response classes.
→ Prior Probability = $\pi_k$ (randomly chosen n cases from class k
→ Density function = $f_k(x) = Pr(X = x | Y = k)$ → high prob if n from k has $X = x$  QDA
→ Posterior Prob = $P_k(x) = Pr(Y = k | X = x) = \frac{\pi_k f_k(x)}{\sum \pi_l f_l(x)} = \frac{\pi_k(P_k(x))}{\pi_k f_k(x) \sum \pi_l f_l(x)}$ – assume normal  no common $\Sigma$ $\delta_k \neq \Sigma_3$
LDA for p=1
assume $f_k(x)$ is normal: $f_k(x) = \frac{1}{\sqrt{2\pi}\sigma_k} \exp\left(\frac{-1}{2\sigma_k^2}(x - u_k)^2\right)$  mean / exclude training set range
$\sigma_1^2 = \sigma_2^2 = \sigma_k^2$ (shared variance)  LDA for p>1
⇒ assign to largest:
$\delta_k(x) = x \cdot \frac{u_k}{\sigma^2} - \frac{u_k^2}{2\sigma^2} + \log(\pi_k)$  $u_k = \frac{1}{n_k}\sum_{i:y_i=k} x_i$ (p, k)
when $k = 2, \pi_1 = \pi_2$:  $\sigma^2 = \frac{1}{n-k}\sum(x_i - u_k)^2$ $\left(\frac{p(p-1)}{2}\right)$
$x = \frac{u_1^2 + u_2^2}{2}$  $\hat{\sigma}^2$  $u_k = $ mean $(k)$  $\frac{\sum(x_i - u_k)}{n-k}$ $(k)$  $\delta_k(x) = x^T \Sigma^{-1} - \frac{1}{2}u_k^T \Sigma^{-1} u_k + \log(\pi_k)$

| Qualitative Predictors | Extensions | Potential Problems |
|---|---|---|
| # dummy = levels - 1 | - additive: predictors = independent | ① non-linearity of R-P relationships |
| $Y_i = B_0 + B_1 X_1 + B_2 X_2 + ...$ | → add interaction, hierarchical principle | ② correlation error terms |
| ③ irreducible error | reduce main effect | |
| → predictive intervals | -linear assumptions | ③ non-constant variance error |
| Y from $\hat{Y}$ | $B_0 + B_1 + f(class A)$ | |
| ② model bias larger | → transform X | ④ outliers: extreme Y |
| → assuming than | → more $R^2$ (better fit) | |
| shape of CI. | $B_0 + E(class C)$ | ⑤ high leverage: extreme X |
| | → might overfit. | |
| | $Y = B_0 + B_1 X + B_2 X^2 + \epsilon$ | ⑥ collinearity - decrease |