

## 1 Linear Regression

$y_i = \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \epsilon_i$ ,  $x_{i1} \equiv 1$ , so  $\beta_1$  is intercept  $Y = X\beta + \epsilon$   
 $\epsilon_1, \dots, \epsilon_n$  are indep.,  $E(\epsilon_i) = 0$ ,  $\text{Var}(\epsilon_i) = \sigma^2$  (homoscedastic)  
**Least square solution**  $\hat{\beta} = \arg\min_{\beta} \|Y - X\beta\|_2^2 = (X^T X)^{-1} X^T Y \sim N_p(\beta, \sigma^2 (X^T X)^{-1})$

$RSS = \|Y - X\hat{\beta}\|_2^2$ ,  $\sigma^2 \approx \frac{1}{n-p} RSS$ , if  $\epsilon \sim N$  then  $\hat{Y} \sim N_n(X\hat{\beta}, \sigma^2 P)$ ,  
 $r \sim N_n(0, \sigma^2(I - P))$ ,  $\hat{\sigma}^2 \approx \frac{1}{n-p} \|r\|_2^2$  where  $P = X(X^T X)^{-1} X^T$

**Interpreting R-output:** Residuals  $y - \text{model}$  \$fitted.values;  
Residual Standard Error:  $\hat{\sigma}$ ; Coeff. Estimate =  $\hat{\beta}$ ;  
Coeff. Std Error =  $\hat{\sigma}(X^T X)^{-1}$ ; Coeff. t value; Estimate / Std. Error.

### Linear Regression

```
fit <- lm(y~x1+x2) # Fit only x1 and x2 (so p=3)
predict(fit, pred.data.frame)
# Manual fit
X <- cbind(1, x1, x2) # n = 3
X1X_inv <- solve(X) %>% X
beta.hat <- X1X_inv %>% t(X_inv) %>% y
res <- y - X_inv %>% beta.hat # Residuals
RSE <- sqrt(sum(res^2)/(n-p)) # Residual std. error.
# Est. of the std of the noise in the linear model
se.x1 <- RSE * sqrt(X1X_inv[2,2]) # Std. error of x1
val.x1 <- beta.hat[2] / se.x1 # t value of x1
p.val.x1 <- 2 * pt(abs(val.x1), df=n-p, lower=F)
# Alternative: car::AIC, AIC, AICc, AICc, AICc, AICc
coef <- summary(fit)$coef # Coefficients
t1 <- coef[1, 1], Estimate / coef[1, 1], "Std. Error"
# Poly regression
lm(wage~poly(age, 4)) # Orthogonal polynomials
lm(wage~poly(age, 4, raw=1)) # same as below
lm(wage~age+I(age^2)+I(age^3)+I(age^4)) # Alternative
```

### 1.1 Tests and model selection

#### Entity-wise test

$H_0: y = X\beta + \epsilon$  with  $\beta_j = 0$   $H_A: y = X\beta + \epsilon$  with  $\beta_j \neq 0$

Under  $H_0$ :  $\frac{\beta_j - E(\beta_j) = 0}{\sqrt{\sigma^2(X^T X)^{-1}}} \sim N(0,1)$  t-statistic:  $\frac{\beta_j}{\sqrt{\sigma^2(X^T X)^{-1}}} \sim t_{n-p}$

**P-Value:** P(robs a value of the test stat. that is as rare or rarer than the one we saw if  $H_0$  is true). Rarest (if unique) would have p-value of 0. If  $\alpha$  then reject  $H_0$ . **ANOVA (Anlysis of variance)**  
 $\|Y - \hat{Y}\|_2^2 = \|Y - \hat{Y}\|_2^2 + \|Y - \hat{Y}\|_2^2$

sum of squares, degrees of freedom, mean square, E[mean square]

regression	$\ Y - \hat{Y}\ _2^2$	$p-1$	$\ Y - \hat{Y}\ _2^2 / (p-1)$	$\sigma^2 + \frac{1}{p-1} \frac{\ Y - \hat{Y}\ _2^2}{p-1}$
error	$\ Y - \hat{Y}\ _2^2$	$n-p$	$\ Y - \hat{Y}\ _2^2 / (n-p)$	$\sigma^2$
total sum of squares	$\ Y\ _2^2$	$n-1$	-	-

Under the null hypothesis  $H_0: \beta = 0$  We have that:  
 $F = \frac{\|Y - \hat{Y}\|_2^2 / (p-1)}{\|Y - \hat{Y}\|_2^2 / (n-p)} \sim F_{p-1, n-p}$  Bigger F = better fit = lower P-value

### P-Values & ANOVA

```
fit.smaller <- lm(y ~ x1)
# Partial F test. Anova uses RSS and DoF
# of largest (last) model, so use ascending order!
anova(fit.smaller, fit, fit.all)
# Overall F-test
fit.empty <- lm(y ~ 1, data=...) # Empty model
anova(fit.empty, fit) # Compare models
# Alternative F-test
fTest <- summary(fit)$statistic
pval <- 1 - pf(fTest[1], df1=fTest[2], df2=fTest[3])
```

### R<sup>2</sup> (Coefficient of determination):

the proportion of the total variation of the response  $y$  around its mean  $\bar{y}$  that is explained by the regression  $\hat{y}$  via the Analysis Of Variance decomposition  
 $R^2 = \frac{\|Y - \bar{y}\|_2^2}{\|Y\|_2^2} = 1 - \frac{\|Y - \hat{Y}\|_2^2}{\|Y\|_2^2} = \frac{RSS}{TSS} = 1 - \frac{RSS}{TSS}$  Bigger  $R^2$  = better fit

### R squared

```
RSE <- sqrt(sum(residuals(fit)^2)/(n-p))
RSS <- sum(res^2) # Residual sum of squares
TSS <- sum((y - mean(y))^2) # Total sum of squares
Rsq <- 1 - RSS / TSS
AdjR2 <- 1 - (RSS/(n-p)) / (TSS/(n-1))
```

## 1.2 Model selection

Penalized RSS:  $\|Y - X\hat{\beta}\|_2^2 + \lambda \|\hat{\beta}\|_2^2$

**Forward Selection:** Greedily add feature that brings best improvement (CV) or penalized RSS, iterate.  
**Backward Selection:** Same but start from full feature set and eliminate features.

### Stepwise methods

```
# Backward selection
mortal.full <- lm(mortality ~ ., data=mortality)
mortal.lw <- stepAIC(mortal.full, direction = "backward")
# Forward selection
mortal.empty <- lm(mortality ~ 1, data = mortality)
mortal.fw <- stepAIC(mortal.empty, direction = "forward", scope = list(lower=mortal.empty, upper=mortal.full))
lower=mortal.empty)
```

### 1.3 R Diagnostic plots (plotting, which)

#1 **Takey-Anscombe** plots the points follow the line, else  $E(\epsilon) = 0$  violated (if so, try data transforms, log if variance grows linearly, sqrt if variance grows as sqrt, or weighted regression) **#2-Q-Q Plot** should follow a straight line, else error not Gaussian (still all fine). **#3 Scale-Location:** Checks for heteroscedasticity (similar to Turkey, but here we have standardized residuals on the y axis). **#4 Cook distance:** Shows if some data points have a larger impact on the fit than others (outliers) **#5 Residuals vs Leverage** Leverage = measure of points being outliers in terms of the feature distribution. Points with high leverage normally also have high residuals, but not always  
**Categorical Variables:** For two levels:  $y_i = \beta_1 x_{i1} + \dots + \beta_p x_{ip} + \lambda_i \epsilon_i$  so if  $i$  is in category, then  $d_{i1} = 1$  else  $d_{i1} = 0$ . This acts as a different intercept  $E(y_i - E(y)) = \lambda_i$ . If more categories, add more dummy variables. This can be done for any feature too, not just for intercept

### Categorical variables

```
# overwrite club as a factor, otherwise will be
# considered a continuous variable
# If done correctly, for club in {1,2,3}, the summary
# will show club2, club3 (first is baseline)
clubs.dset[club] = factor(clubs.dset[club])
reg2 = lm(cbind = factor(clubs.dset[club])
```

**Interaction Definition** if there is no interaction, then the difference in the predicted values between samples with type=0 and type=1 should be the same if the other predictors coincide, no matter whether partners=1 or partners=8. Mathematically, this means that  $\gamma_{1,0} - \gamma_{1,1} = \gamma_{8,0} - \gamma_{8,1}$ .  
There is an interaction effect between  $X_1$  and  $X_2$  in relation to  $Y$  if the slope of  $X_1$  in the linear regression model depends on  $X_2$ .  
**Correlation Definition** Two variables  $X_1$  and  $X_2$  are correlated if the product of their expectations does not equal the expectation of their product. **Bias Definition:** Estimate - TrueValue

### 2 Confidence intervals

We find some CI of level  $\alpha$ ,  $P(\text{CI}) = 1 - \alpha$   
For coefficients:  $\beta_j \in [\hat{\beta}_j \pm \hat{\sigma}(\hat{\beta}_j) \cdot t_{n-p-1, \alpha/2}]$  (meaning the  $1 - \alpha/2$  quantile of a  $t_{n-p}$  distribution)  
For predictions: there are 2 types: confidence (where the true model should lie) and prediction (where individual predictions might end up), the latter is always wider.

### Confidence intervals

```
confInt(fit, level=0.95) # CI for coeffs
# CI for new samples, see predict.lm
predict(reg, sample, interval="confidence", level)
predict(reg, sample, interval="prediction", level)
```

### Bias Variance Trade-Off

Expected Test MSE at  $x_0$ :  $E[(\hat{f}(x_0) - y_0)^2] = \text{Bias}^2(\hat{f}(x_0)) + \text{Var}(\hat{f}(x_0)) + \sigma^2$ , where  $\text{Bias}^2(\hat{f}(x_0)) = [E(\hat{f}(x_0)) - f(x_0)]^2$ .

### Bias Variance Trade-Off of a Method

```
Bias <- mean(Estimate.bingCV) - TrueValueSimulated
MSE <- Bias^2 + var(Estimate.bingCV)
```

## 3 Non-parametric Density Estimation

Estimate true density  $f(x)$  with  $\hat{f}(x)$ .  
ie. minimize  $MS E(x) = E[(\hat{f}(x) - f(x))^2]$   
**Kernel density est.**  $\hat{f}(x) = \frac{1}{n} \sum_{i=1}^n K\left(\frac{x - x_i}{h}\right)$   
Conditions for  $K$ :  $K(x) \geq 0 \forall x$ ,  $\int K(x) dx = 1$ ,  $K(x) = K(-x) \forall x$ . Common choices for  $K$ :  
Naive/histogram-like approach:  $K(x) = \frac{1}{2} 1_{|x| \leq 1}$   
Gaussian:  $K(x) = (2\pi)^{-\frac{1}{2}} \exp(-\frac{1}{2} x^2)$

### Non Parametric Density Estimation

```
denshat <- density(data, bw=0.1, kernel="gaussian")
# construct approximate pdf from density estimate
# rule says how to pick outside [min(x), max(x)]
pdf <- approxfun(denshat$y, denshat$x, rule=2)
# taking the log likelihood of some data
logLikelihood = mean(log(pdf(eval.data)))
```

### 4 Non-parametric Regression

We want to estimate  $m(x) = E[Y | X = x]$   
**4.1 Nadaraya-Watson kernel method**  
 $\hat{m}(x) = \frac{\sum_{i=1}^n K((x - x_i)/h)}{\sum_{i=1}^n K((x - x_i)/h)}$ ,  $h$  is the bandwidth.  $K$  is the kernel function  
**Local bandwidth** to improve one idea is to have  $h(x)$  depend on  $x$  and not be the same everywhere.  
**Local polynomial:** Extension of Watson method that fits a local polynomial instead of a local constant.  
**Smoothing Splines:**  
 $G = \{g : [a, b] \rightarrow \mathbb{R} : g'' \text{ exists and } \int_a^b g''(x)^2 dx < \infty\}$  is the class of functions to consider.  $\hat{g} = \arg\min_{g \in G} \sum_{i=1}^n (y_i - g(x_i))^2 + \lambda \int_a^b g''(x)^2 dx$ . Note: if  $\lambda = 0$ ,  $\hat{g}$  is any function in  $G$  that passes through all data points. If  $\lambda = \infty$ , then least squares estimate. Shrink version of natural spline with knots at  $x_1, \dots, x_n$ .  
Closed form solution:  $\hat{\beta} = (B^T B + \lambda I)^{-1} B^T Y$ , where  $O_{jk} := \int B_j^T(x) B_k(x) dx$  with spline basis functions  $B_j$ ,  $j \in \{1, \dots, n\}$ . Smoother matrix  $S = B(B^T B + \lambda I)^{-1} B^T Y$ .

**Code for non-parametric regression methods**  
# Nadaraya-Watson, x.points are predicted on!  
# CAREFUL! This reorders things internally!  
watson = ksmooth(x, y, bandwidth=0.2, x.points=x)  
plot(watson\$fit, watsom\$y)  
# Local polynomial, span or emp.target for smoothing  
poly = loess(y ~ x)  
# Splines ~ df, spar, lambda for deg free  
spline <- smooth.spline(x, y)  
y.hat <- predict(spline, x = x)\$y # or spline\$y  
# Local optimum bandwidth  
# global optimum bandwidth  
# global optimum bandwidth  
# global optimum bandwidth  
plot(global.opt\$y.out, global.opt\$y.out) # pred on x.out  
# Local optimum bandwidth  
local.opt = loess(x, y, x.out=x)  
plot(local.opt\$y.out, local.opt\$y.out)

### 4.2 Degrees of freedom

**Hat Matrix**  $\hat{Y} = SY$ . Every linear method in  $Y$  has such a matrix  
**Degrees of freedom**  $df = \text{tr}(S)$ .

**Compute hat matrix and df**  
Id <- diag(nrow(data))  
for (i in 1:nrow(data))  
for (j in 1:nrow(data))  
S[i, j] <- ksmooth(x[i], Id[j], x.point=x)\$y  
df <- sum(diag(S))

### 5 Cross Validation

Can be used for model assessment (estimate test MSE) and model selection (choose tuning parameters, variable selection). But not both at the same time (use double CV instead).  
**Validation set:** split data into two halves, train on one, test on the other (most bias). **k-Fold:** same, but with many folds. Try all folds for test and average metrics over the folds (in between).  
 $\text{Var}(\hat{\theta}_k) = 1/K \cdot \text{Var}(MS E_k)$  **LOOCV:** extreme version where each

data point is a fold (least bias, high variance since the datasets are correlated).  
**LOOCV efficient computation:** Assuming a linear model ( $\hat{Y} = SY$ ) we have a fast way to find LOOCV.

$\text{LOOCV} = \frac{1}{n} \sum_{i=1}^n \left( \frac{Y_i - m(X_i)}{1 - S_{ii}} \right)^2$   
A historical approximation of this is GCV.

$\text{GCV} = \frac{1}{n} \sum_{i=1}^n \left( \frac{Y_i - m(X_i)}{1 - \frac{1}{n} \text{tr}(S)} \right)^2$

### Cross validation

```
Indices = sample(0:n.elements, replace=F)
# or, if ordering is needed
indices = 1:n.elements
fold.indicator = cut(indices, breaks=n_folds, labels=F)
for (i in 1:n_folds) {
  test.indices <- which(fold.indicator == i)
  train <- data[test.indices,]
  test <- data[test.indices,]
```

### 6 Bootstrap

Sample uniform from data points with replacement, compute bootstrapped estimator. For a large dataset  $x_1, \dots, x_n$ , the probability that  $x_1$  is contained in a random bootstrap dataset is:  $1 - (1 - 1/n)^n \approx 2/3$  (for large  $n$ , limit goes to  $1 - 1/e$ )

#### Bootstrap CI

**Reversed Quantile CI:** (intelligent)  $(\hat{\theta}_n - q_{\alpha}^*, \hat{\theta}_n - q_{1-\alpha}^*)$  (type="basic")

**Quantile Bootstrap CI:** (naive) not theoret. justified unless  $\hat{\theta}_n$  is symm.:  $[q_{\alpha}^*(\hat{\theta}_n), q_{1-\alpha}^*(\hat{\theta}_n)]$  (type="perc"). Same as reversed.

**Normal Bootstrap CI:** Assume  $\hat{\theta}_n$  is to be asympt. normal:  $\hat{\theta}_n \pm q_{\alpha}^*(1 - \alpha/2) \hat{\text{sd}}(\hat{\theta}_n)$  where  $z \sim N(0,1)$  and  $\hat{\text{sd}}(\hat{\theta}_n) = \sqrt{\text{Var}(\hat{\theta}_n)}$  (type="norm")

**Double bootstrap CI:** For a regular bootstrap CI we have some error:  $P(\hat{\theta} \in (1 - \alpha) = 1 - \alpha + \Delta_n)$ . The goal is to adjust the interval such that we actually get the coverage probability that we originally aimed for:  $P(\hat{\theta} \in (1 - \alpha)) = 1 - \alpha$   
**Parametric Bootstrap:** We assume  $X_1, \dots, X_n \sim P_0$ . We estimate  $\hat{\theta}$  then we draw bootstraps  $X_1^*, \dots, X_n^* \sim P_0$   
**Regression bootstrap:** We do  $Y_i = m(X_i) + \epsilon_i$  **parametric:**  $\epsilon_i \sim N(0, \hat{\sigma}^2)$  **model-based:**  $\epsilon_i \sim \hat{P}$ , the empirical dist. of residuals

```
# Nonparametric
statistic <- function(orig.data, ind)
{
  {mean(orig.data[ind,])}
  non_para.boot = boot.data, statistic, R=100)
  non_para.boot$0 # final bootstrap estimate
  non_para.boot$1 # each individual estimate
  boot.ci(non_para.boot, conf, type=c("basic", "norm", "perc"))
}
```

**Parametric**  
statistic <- function(data) {mean(data)}  
fitted.distr = fiddistr(data, "gamma")  
ran.gen <- function(data, me) {

```
rgamma(length(data), shape = mle[1], rate = mle[2])
}
para.boot = boot(data, statistic, R, sim="parametric",
  ran.gen=ran.gen, me=fitted.distr)
# Intervals by hand
ci.quantile =
  quantile(boot.estimates, probs=c(0.025, 0.975))
center = 2 * theta.hat - mean(boot.estimates)
ci.normal = c(center - spread, center + spread)
ci.reversed <- theta.hat - quantile(boot.estimates - theta.hat, probs=c(0.975, 0.025))
```

### 7 Classification

for any class we estimate the probability of being in such class.  
 $\pi_j(x) = P(Y = j | X = x)$

The Bayes estimator is  $\hat{c}(x) = \arg \max_j \pi_j(x)$ .

**Linear & quadratic discriminant analysis:** Assumption:  $X \mid Y \sim N(\mu_j, \Sigma)$ . To approximate

$$\hat{\mu}_j = \sum_{i=1}^n X_i 1_{Y_i=j} / \sum_{i=1}^n 1_{Y_i=j} = \frac{1}{n_j} \sum_{i:Y_i=j} X_i$$

$$\hat{\Sigma}_j = \frac{1}{n-j} \sum_{i=1}^n (X_i - \hat{\mu}_j)(X_i - \hat{\mu}_j)^T 1_{Y_i=j}$$

**Quadratic discriminant analysis:** Compute  $\Sigma$  group-wise by only looking at elements in the group and by correcting with  $\frac{1}{n-j-1}$  instead of  $\frac{1}{n-j}$

**Logistic regression:** Let  $\pi(x) = P(Y = 1 \mid X = x)$ . Estimate with a linear model  $g(x) = \logit(\pi(x)) = \log \frac{\pi(x)}{1-\pi(x)}$ . Assuming  $g(x) = X\beta$ , for the log-likelihood  $\ell(\beta; Y, X) = \sum_{i=1}^n \left( Y_i \beta^T X_i - \log(e^{\beta^T X_i} + 1) \right)$

**Comment1: Null deviance vs Residual deviance:** how well the response variable can be predicted by a model with only an intercept term vs how well the response is predicted by the model when the predictors are included.

**Comment2:** LDA is a logistic regression model, since log-odds ratio is linear. But not every logistic regression model is LDA due to normal feature assumption in LDA.

**Techniques for multi-class:**

- 1) One versus rest: Train a model per class against all other classes to receive  $\pi_j(x)$ . Normalization needed.
- 2) Everyone versus reference: Train all other  $j-1$  classes against reference class 0:  $g(x) = \log(\pi_j(x)/\pi_0(x))$
- 3) Multinomial Distribution: Use multinomial distribution instead of Binomial
- 4) One vs. one: Train a model to distinguish each class from each other. In total  $\binom{J}{2}$  models.
- 5) Ordered classes:  $\logit(P(Y \leq k \mid x)) = a_k + g(x)$  with  $a_0 \leq a_1 \leq \dots \leq a_{J-1}$

#### Classification in R

```
# train and test contain numeric features, while labels are a factor
knn_preds = knn(train = train, test = test,
  cl = train_labels, k=5, prob=f)

acc = mean(knn_preds == test_labels)

# LDA/QDA
c_ida=c(1da(x=trial, c["x1", "x2"]), grouping=data[, "y"])
c_qda=c(qda(x=trial, c["x1", "x2"]), grouping=data[, "y"])
accmean(c_ida, newdata) # class
# posterior if we need probabilities

# Logistic Regression
log reg
mod <- glm(y ~ ., data = data, family = gaussian)
predict(mod, test_data, type = "response") # return
# the probability of being 1
# for multinomial classification
class_multinom <- multinom(y ~ ., data = data)
```

#### ROC Receiver Operating Characteristic:

Predict  $\hat{Y}(x) = I(\hat{\pi}(x) \geq 0)$ . Basics

$$TP + FN = P, TN + FP = N$$

$$\text{sensitivity} = TPR = \frac{TP}{TP+FN}, \text{estimates } P(\hat{Y} = 1 | Y = 1)$$

$$\text{specificity} = TNR = \frac{TN}{TN+FP}, \text{estimates } P(\hat{Y} = 0 | Y = 0)$$

precision =  $\frac{TP}{TP+FP}$ , recall=synonym for specificity

Changing all values of  $\theta \in [0, 1]$  look at specificity vs. sensitivity. Integrating the resulting function yields the area under the curve. ROC AUC =  $P(\hat{\pi}_0 < \pi_1)$ , where  $I_0$  and  $I_1$  are sampled uniformly at random from indexes with classes  $Y = 0$  and  $Y = 1$ , respectively. default cost = misclass. rate =  $\frac{EP+FN}{n}$

#### Evaluation

```
# I is a vector of 0/1 (the truth)
# p is a vector of predicted probs to be 1
pred <- prediction(predictions = p, labels = 1)
perf <- performance(pred, "pr", "auc")
plot(perf) # plots the ROC curve
cost <- performance(pred, "cost", cost.fp = 1,
  cost.fn = 2)
plot(cost)
```

#### 8 Flexible regression

Making some structural assumption, we estimate

$$g(x) = \begin{cases} E[Y \mid x] & \text{in regression} \\ \logit(\pi(x)) & \text{in classification} \end{cases}$$

##### 8.1 Additive model

In additive model we assume that

$$g(x) = \mu + \sum_j g_j(x_j) \quad E[g_j(x_j)] = 0 \quad \forall j$$

To find  $g_j$ , use **backfitting**. Let  $s_j = (s_{j1}, \dots, s_{jn_j})^T \rightarrow (s_{j1}, \dots, s_{jn_j})^T$  be a smoother given any prev model that we used, eg. lin reg. kernel estimation (splines)

Initialize  $\hat{\mu} = \frac{1}{n} \sum_{i=1}^n y_i$ ,  $\hat{g}_j = 0 \quad \forall j$   
Do until convergence:  
for each  $j = 1, \dots, p$ :

$$\hat{g}_j \leftarrow s_j(y - \hat{\mu} - \sum_{k \neq j} \hat{g}_k)$$

$$\hat{g}_j \leftarrow \hat{g}_j - \frac{1}{n} \sum_{i=1}^n \hat{g}_i(x_{ij})$$

##### 8.2 MARS (Multivariate Adaptive Regression Splines)

We make a regression with functions in  $M$ . We start with  $M = \{h_0(\cdot) = 1\}$ . We then iteratively add functions in  $M$  in the form:

$$h_{new1-2}(x) = h_j(x) \pm x_j \mp d_j$$

Where  $h_j$  is in  $M$ . We then prune back useless stuff.

#### Flexible regression in R

```
# Additive model
gamf orm <- warp(formula(logprob ~ ., data=d.ozone.e)
  gam_fit <- gam(gamf orm, data = mydata)
  predicted_y <- predict(gam_fit, testdata)
# MARS model
mars_fit <- earth(y ~ ., data = mydata, degree = 3) #
  no. of interactions
predict(mars_fit, newdata = test_data)
```

##### 8.3 CART (Classification and regression trees)

The tree decides the domain in rectangles. In the end you have a piece-wise constant function

$$\hat{g}_{tree}(x) = \sum_{M_i} \beta_i R_{M_i}(x)$$

We can penalize larger tree to reduce variance. We choose the simplest model which is within standard error to the best one in cross-validation. To fit we are greedy. **Pros and cons**

- Highly flexible and interpretable
- Too simplistic, only piece-wise constant
- **One Std Error Rule:** First find the model with the lowest cross-validation error, then choose the simplest model, which is at most one standard-error worse than that model.

**Residual mean deviance:** a measure of the error remaining in the tree after construction

#### CART in R

```
tree = rpart(y ~ ., data = data, control =
  #cp = complexity factor, the R^2 must increase at least
  # of a factor cp after split
  preds = predict(tree, newdata) # type = "class" for
  # classification
  # for no. reg. at all: cp=0, minsplit=
  plotcp(tree) or printcp(tree) # estimate the best cp
  cp.opt = tree$cp[which.min(tree$cp)]
  # Same things with package "tree"
  single_tree = tree(y, train ~ ., data=Boston, train,
  # control = tree.control(nobs, ...))
  plot(single_tree, data)
  pruned_tree = prune.tree(single_tree, k=1)
```

#### 9 Autoencoder

Find  $E: \mathbb{R}^d \rightarrow \mathbb{R}^d, D: \mathbb{R}^d \rightarrow \mathbb{R}^d$ , such that  $\|\sum_{i=1}^n (E(x_i) - D(E(x_i)))\|_2^2$  is minimized.

For linear restriction of  $E$  and  $D$ , this corresponds to doing PCA. **PCA:** keeps the largest  $k$  singular values (dimensionally reduction).

##### 10 Bagging-Boosting

###### 10.1 Bootstrap Aggregating (Bagging)

Bagging is a variance reduction technique. We need a simple estimator  $\hat{g}$  (for instance a tree).

- 1) We take  $B$  bootstrap samples
- 2) For each bootstrap sample we train an estimator  $\hat{g}_1, \dots, \hat{g}_B$ .
- 3) Aggregation:  $\hat{g}_{bag}(\cdot) = \frac{1}{B} \sum_{b=1}^B \hat{g}_b(\cdot)$

**classification:** In this case we have majority voting (or we can average base estimator probs of being in each class)

**Subsample aggregating (Subbagging):** Like Bagging but we just take subset of the dataset (without repetition) instead of bootstrapping.  $m < n$  size of the subsets (argued to be computationally cheaper)

**Out-of-Bag Error:** Some bags have not trained on a particular sample. Can predict this only by the bags that have not been trained on it (should be  $\sim 1/3$  for all samples and average to get a valid estimate for the test error).

**Random Forests:** Essentially bagged trees. Have  $B$  bootstrap samples  $\rightarrow$  create trees. They reduce dependence between tree estimates by only allowing a random subset of predictors at each split. Default: regression  $p/3$ , classification  $\sqrt{p}$  (in R option mtry).

#### Random forest

```
library(randomForest)
rf <- randomForest(y ~ ., data = mydata, mtry=p-1,
  # importance = TRUE, ntree = 100)
# mtry: no. of features from which we chose a split
# if mtry = n, predictors then we have bagging trees
test_preds = predict(tree, newdata)
oob_preds = rfpredicted
importance(rf) # get feature importance
```

##### 10.2 Feature Importance Measures for RF

One idea is "screwing up" a variable at a time and checking how bad the accuracy drops (or mse increases). Second is to sort variables by how much they decrease the gini index/purity.

These should be taken with a pinch of salt because of predictor correlation (if we retain leaving the original most important predictor out, we may still get the same accuracy)

##### 10.3 Boosting

Boosting is a bias reduction technique. We have  $g$  to be very simple estimator (stump or small tree). We iteratively train a  $g$  predictor on the current model and then update the model:  $f \leftarrow f + \eta g$  ( $\eta$  usually is .1).

In presence of high-dimensional predictors, boosting is also very useful as a regularization technique for additive or interaction modeling

**Boosting for linear models:** Update a single feature on each iteration (the one that has the highest correlation with the residuals). The outcome of this is the same as Lasso under certain assumptions.

##### 11 Ridge/Lasso:

We are in the situation where  $p \gg n$ . We center and scale all the  $x_i$  so that we can penalize all  $\beta_i$  equally. In Ridge we penalize with  $\|\beta\|_2^2$  (convex, rotation invariant). In Lasso  $\|\beta\|_1$  (convex, yields sparse solution)

**Elastic net:** Mix of Ridge and Lasso, the regularization is  $\lambda_2 \|\beta\|_2^2 + \lambda_1 \|\beta\|_1$ . We call  $\alpha = \frac{\lambda_1}{\lambda_1 + \lambda_2}$ .

**Adaptive Lasso:** First you get an estimate of  $\beta$ . Then we penalize  $\beta_i$  the more we think they are 0, we set  $w_i = |\beta_i|^{-\gamma}$ . The regularization term will be  $\lambda \sum w_i |\beta_i|$ .

**Relaxed Lasso:** First get an estimate of  $\beta$  via Lasso. Then do something similar to lasso but we use  $\phi(\lambda)$  instead of  $\lambda(\phi \in [0, 1])$ .

**Group Lasso:** Predictors are divided into  $J$  groups of size  $p_1, \dots, p_J$ .  $\sum \beta_j = p_j \beta_j^{group} = \arg \min_{\beta_j} RSS(\beta_j) + \lambda \sum_{i=1}^{p_j} |\beta_{ji}|$  (if  $L = p$ , we get Lasso). Acts like Lasso on a group level. Useful if there are

categorical variables with  $> 2$  categories (put all corresponding dummy variables in a group).

#### Ridge & Lasso

```
library(glmnet)
# creating a polynomial formula
ff <- warp(formula(y ~ ., data,
  # wrap string="poly(*, degree=3)")
  # have it take all 3-ws
  ff <- update(ff, y ~ as3)
  # data for glmnet must be a MATRIX
  # returns as many models as elements are in the grid
  grid <- 10^seq(log10=10, to=2, length=100)
  ridge <- glmnet(features, y, alpha=0, lambda.grid)
  lasso <- glmnet(features, y, alpha=1, lambda.grid)
  # coefficients for specific lambda_val
  coef(lasso, s=lambda_val) # or reFit with specific
  cv = lambda
  cv.glmnet(lmu, y, alpha=0.5, nfolds=10) # for input CV
```

#### R-Help

```
# Arguments per row, simpler alternative is type="class"
# probs <- predict(tree, newdata = data)
# pred <- colnames(probs)[max.col(probs)]

# L1ST OPS -----
mylist <- numeric(n) # Initialize empty list size n
# need to access with [[]] when storing objects!
# Test if an element is in a list
if ("X1" %in% names(coef(m, mod)))

# MATRIX OPS -----
%*% # matrix multiplication
solve(X) # matrix inversion

# DATAFRAME -----
# indexing always is df$row(s), col(s)
dfOrdered <- df[order(df[,x], decreasing = F), ]
col_induces = which(names(df) %in% c("col1", "col2"))
df <- df[, col_induces] # remove columns named above
df <- subset(df, select=c("col1", "col2")) # or this

# DATA PROCESSING -----
complete.cases # remove rows
which(array==target_value) # get idx of target val

# DISTRIBUTION OPS -----
runif(10, min=0, max=1), rnorm(10, mean=0, sd=1)
# constructing approximate pdf from density estimates
pdf = approxfun(denshat$xx, denshat$y, rule=2)
log(pdf(eval.data)) # get log likelihood of some data

# HIGH LEVEL OPS -----
apply(X, d, func) # apply func to all rows(d = 1) or
# col(d = 2) of matrix X
sapply(X, func, param1, param2) # apply func to elems
# in the list X, extra params go to func
replicate(n, times, fun(data, sample(1:nrow(data)))) #
# apply fn to the sampled data n times

# WEIRD ERROR MESSAGES
# If predict complains that num elements are
# different than expected -> check that the naming
# in the new data is the same as in the train data
```