

Multiple Linear Regression

1. The Linear Model

Assume $Y_i = x_i^\top \beta + \epsilon_i$ or $Y = X \times \beta + \epsilon$ with $X \in \mathcal{R}^{(n \times p)}$; ($n > p$) and $\mathbb{E}[\epsilon_i] = 0, Var(\epsilon_i) = \sigma^2$. X is often augmented with $(1_{N \times 1})$ to use β_1 as bias.

2. Least Squares Method

LS estimator is $\hat{\beta} = \arg \min_{\beta} ||Y - X\beta||_2^2 = (X^\top X)^{-1} X^\top Y =$
 PY (orth. proj. of Y onto $span(X)$). Estimate $\hat{\sigma}^2 =$
 $\frac{1}{n-p} \sum_{i=1}^n r_i^2$ with $\mathbb{E}[\hat{\sigma}^2] = \sigma^2$.

Assumptions for Linear Model

- i Linear regression equation is correct, i.e. $\mathbb{E}[\epsilon_i] = 0 \forall i$.
- ii We measure x_i 's exactly. Else, need correction (?).
- iii Error is homoscedastic, i.e. $Var(\epsilon_i) = \sigma^2 \forall i$. Else, use "Weighted LS".
- iv Errors are uncorrelated, i.e. $Cov(\epsilon_i, \epsilon_j) = 0 \forall i \neq j$. Else "Generalized LS".
- v Errors are jointly normally distributed. Else "Robust Methods".

Moments of least squares estimates

Assume $Y = X\beta + \epsilon, \mathbb{E}[\epsilon] = 0, Cov(\epsilon\epsilon^\top) = \sigma^2 I$ (all assumptions satisfied). Then

- i $\mathbb{E}[\hat{\beta}] = \beta$ ($\hat{\beta}$ is unbiased).
- ii $\mathbb{E}[\hat{Y}] = \mathbb{E}[Y] = X\beta$ and $\mathbb{E}[r] = 0$.
- iii $Cov(\hat{\beta}) = \sigma^2 (X^\top X)^{-1}$.
- iv $Cov(\hat{Y}) = \sigma^2 P, Cov(r) = \sigma^2 (I - P)$.

If additionally $\epsilon_1, \dots, \epsilon_n$ i.i.d. $\sim \mathcal{N}(0, \sigma^2)$, then

- i $\hat{\beta} \sim \mathcal{N}_p(\beta, \sigma^2 (X^\top X)^{-1})$
- ii $\hat{Y} \sim \mathcal{N}_n(X\beta, \sigma^2), r \sim \mathcal{N}_n(0, \sigma^2 (I - P))$
- iii $\hat{\sigma}^2 \sim \frac{\sigma^2}{n-p} \chi^2_{n-p}$.

Even when normality assumption doesn't hold, central limit theorem is a justification.

3. Tests and Confidence Regions

T-test

Assume linear model with Gaussian errors (or "large enough" sample size), s.t. $\hat{\beta} \sim \mathcal{N}_p(\beta, \sigma^2 (X^\top X)^{-1})$ is normally distributed. Then we can test the null-hypothesis $H_{0,j} : \beta_j = 0$ against $H_{A,j} : \beta_j \neq 0$:

$$\frac{\hat{\beta}_j}{\sqrt{\sigma^2 (X^\top X)^{-1}_{jj}}} \sim \mathcal{N}(0, 1) \Rightarrow T_j = \frac{\hat{\beta}_j}{\sqrt{\hat{\sigma}^2 (X^\top X)^{-1}_{jj}}} \sim t_{n-p}$$

under the null-hypothesis $H_{0,j}$. Unknown σ^2 is replaced by $\hat{\sigma}^2$. Note that $t_{n-p} \approx \mathcal{N}$. An individual t-test for $H_{0,j}$ gives the effect of β_j after subtracting the linear effect of all $\beta_i \neq j$. Note that in `summary.lm`, the term *Std. Error* is $\sqrt{\hat{\sigma}^2 (X^\top X)^{-1}_{jj}} = \sqrt{Var(\hat{\beta}_j)}$.

Global null hypothesis and ANOVA

We can also check the global null-hypothesis $H_0 : \beta_2 = \dots = \beta_p = 0$ using the *analysis of variance* (ANOVA), which decomposes

$$||Y - \hat{Y}||_2^2 = ||\hat{Y} - \bar{Y}||_2^2 + ||Y - \hat{Y}||_2^2.$$

Under the global null-hypothesis $\mathbb{E}[Y] = \mathbb{E}[\hat{Y}] = const.$ (no effect of predictor variables). $\sigma^2 / \hat{\sigma}^2$ yields F-statistic:

$$F = \frac{||\hat{Y} - \bar{Y}||^2 / (p - 1)}{||Y - \hat{Y}||^2 / (n - p)} \sim F_{p-1, n-p}$$

under the global null-hypothesis H_0 . ANOVA also yields *goodness of fit* $R^2 = \frac{||\hat{Y} - \bar{Y}||_2^2}{||Y - \bar{Y}||_2^2}$, which should be around 1. Finally, we can also build a confidence interval using $\hat{\beta}_j \pm \sqrt{\hat{\sigma}^2 (X^\top X)^{-1}_{jj}} \cdot t_{n-p; 1-\alpha/2}$.

```
anova(fit) # global F test
# partial F test - sig. of predictors in .full but not .part
anova(fit.part, fit.full)
```

4. Checking Model Assumptions

Tukey-Anscombe Plot

Error should fluctuate randomly. If error increases linearly, do log-transform $Y \mapsto \log Y$. If error increases with \sqrt{Y} , do a square-root-transform $Y \mapsto \sqrt{Y}$.

QQ-Plot/Normal-Plot

Plot empirical quantiles of residuals on y versus the theoretical quantiles of $\mathcal{N}(0, 1)$ on x. If assumption holds, get straight line with intercept μ and slope σ . Z-shape: long-tailed distr.; Curved: skewed distr.

5. Model Selection

Assume again $\mathbb{E}[\epsilon_i] = 0, Var(\epsilon_i = \sigma^2)$. We need to address *bias-variance trade-off*. Bias is defined as $\mathbb{E}[f(x)] - f(x)$, variance as $q/n \cdot \sigma^2$ with $q \leq p$.

Mallows C_p statistic

Let $SSE(\mathcal{M})$ the residual sum of squares. Then $n^{-1} \sum_{i=1}^n \mathbb{E} \left[(f(x) - \hat{f}_x(x))^2 \right] \approx n^{-1} SSE(\mathcal{M}) - \hat{\sigma}^2 + 2\hat{\sigma}^2 |M|/n$, with \mathcal{I} the indices of selected predictors and $|\mathcal{I}| = q$. Thus, we search for the model that minimizes the C_p -statistic with $C_p(\mathcal{M}) = \frac{SSE(\mathcal{M})}{\hat{\sigma}^2} - n + 2|M|$. Otherwise Akaike's information criterion (AIC) or Bayesian information criterion (BIC). AIC is equivalent to C_p for linear Gaussian models.

```
require(leaps); fit.all <- regsubsets(y~., data=data)
p.regsubsets(fit.all)
```

Forwards and backwards selection

Forward selection: (i) Start with empty model. (ii) (Greedy) Keep adding variable that reduces the residual sum of squares the most. (iii) When done, pick submodel which minimizes C_p . **Backward selection:** (i) Start with full model. (ii) (Greedy) Keep excluding predictor that increases the residual sum of squares the least. (iii) When done, pick submodel which minimizes C_p . Backwards selection typically better but more expensive. When $p \geq n$, use forward selection. Both methods prone to overfitting — p-values (and similar values) are *not* valid anymore and effects look too significant.

```
fit.empty <- lm(y~1, data=data)
fit.full <- lm(y~., data=data)
fit.bw <- step(fit.full, direction="backward")
fit.fw <- step(fit.empty, direction="forward",
  ↪ scope=list(upper=fit.full, lower=fit.empty))
```

Nonparametric Density Estimation

Kernel estimator

Estimate density $\hat{f}(x) = \frac{1}{nh} \sum_{i=1}^n w((x - X_i)/h)$. Kernels include (i) rectangular ($w(x) = 0.5 \cdot 1_{|x| < 1}$), (ii) triangular, or (iii) Gaussian. We require $\int_{\mathbb{R}} K(x) dx = 1$. The bandwidth parameter h is crucial and determines the "smoothness" of the density estimate.

Choosing a bandwidth h

A simple approach is using k -nearest neighbors, i.e. $h(x) = \max_{x_i \in KNN_k(x)} ||x - x_i||_2$ with tuning parameter k . Note that $\int_{\mathbb{R}} K(x) dx = 1$ might be violated. Naturally, the bandwidth also induces a *bias-variance trade-off*. Note that $MSE(x) = \mathbb{E} \left[(f(x) - \hat{f}(x))^2 \right] = (\mathbb{E}[f(x)] - f(x))^2 +$

$Var(\hat{f}(x))$, so we can try to minimize the integrated *MSE* over all points to find the best bandwidth.

Density estimation in higher dimensions

Basically use $\hat{f}(x) = \frac{1}{nh^d} \sum_{i=1}^n K((x - X_i)/h)$ with a Kernel that supports vectors. The Gaussian kernel is the only one that is radially symmetric. Note that in higher dimensions, density estimation becomes very hard, due to data points becoming very sparse.

Nonparametric Regression

Nonparametric regression with *one* predictor variable, i.e. $Y_i = m(x_i) + \epsilon_i$ with $\epsilon_{1:n}$ i.i.d and $\mathbb{E}[\epsilon_i] = 0$. We want $m(x) = \mathbb{E}[Y|x]$ and "some" smoothness.

Kernel regression estimator

A "locally weighted" approach yields the NW kernel estimator $\hat{m}(x) = \frac{\sum_{i=1}^n \omega_i Y_i}{\sum_{i \in \mathbb{R}} \omega_i} = \arg \min \sum_{i \in \mathbb{R}} \omega_i (Y_i - m_x)^2$ (3.1)

with $\omega_i = K\left(\frac{x_i - x}{h}\right)$ a kernel centered at x_i and bandwidth h . As h small \rightarrow large then (high variance) \rightarrow (high bias). For x_i equidistant there exists $h_{opt} = f(\sigma_\epsilon^2, m''(x))$ which can be iteratively found.

```
ksmooth(x, y, kernel="normal", bandwidth=0.2,
  ↪ x.points=x)$y$
# automatic bandwidth
fit.lo<-lokerns(X, Y, x.out=X, hetero=TRUE, is.rand=TRUE)
fit.gl<-gkerns(X, Y, x.out=X, hetero=TRUE, is.rand=TRUE)
```

Local polynomial regression estimator

Instead of finding a local constant m_x we can also find a *local polynomial*, i.e. we replace m_x with $\beta_1 + \sum_{i=2}^p \beta_i (x_i - x)^{i-1}$ (usually $p = 2$ or $p = 4$). Often better at edges and yields first derivative.

```
fit.loess <- loess(y ~ x, data=data.frame(x=x, y=y_pert),
  ↪ span=0.2971339, surface='direct')
fit.loess.pred <- predict(fit.loess, newdata=x)
```

The hat matrix S

We want to construct S with $\hat{Y} = SY$, i.e. the linear operator mapping the labels to the predictions. Given the regression (smoothing) function s , we compute $S_{.j} = s(x, e_j, h)$ with e_j the j -th unit vector. Then $Cov(\hat{m}(x)) = Cov(SY) = SCov(Y)S^\top = \sigma_\epsilon^2 SS^\top$, i.e. $Cov(\hat{m}(x_i), \hat{m}(x_j)) = \sigma_\epsilon^2 (SS)^\top_{ij}$, and $Var(\hat{m}(x_i)) = \sigma_\epsilon^2 (SS^\top)_{ii}$. Estimate $\hat{\sigma}_\epsilon^2 \approx \sum_{i=1}^n (Y_i - \hat{m}(x_i))^2 / (n - df)$. Then

- $\widehat{s.e.}(\hat{m}(x_i)) = \sqrt{\widehat{Var}(\hat{m}(x_i))} = \hat{\sigma}_\epsilon \sqrt{(SS^\top)_{ii}}$
- $\hat{m}(x_i) \approx \mathcal{N}(\mathbb{E}[\hat{m}(x_i)], Var(\hat{m}(x_i)))$
- $I = \hat{m}(x_i) \pm 1.96 \cdot \widehat{s.e.}(\hat{m}(x_i)) \rightarrow$ (pointwise) CI

Additionally we can compute the *degrees of freedom* for regression estimators with $df = \text{tr}(S)$.

```
# Construct S matrix
N <- length(x); Eye <- diag(N)
S.nw <- S.lp <- S.ss <- matrix(0, nrow=N, ncol=N)
for (j in 1:N) {
  y_~ <- Eye[, j]
  S.nw[, j] <- ksmooth(x, y_~, kernel="normal", bandwidth=0.2,
    ↪ x.points=x)$y; $}
est.nw<-est.lp<-est.ss<-matrix(0,nrow=length(x),ncol=nrep)
for (i in 1:nrep) {
  # generate y with disturbance
  y_pert <- y + rnorm(length(x), mean=0, sd=1)
  # try to fit NW
  est.nw[, i] <- ksmooth(x=x, y=y_pert, kernel="normal",
    ↪ bandwidth=0.2, x.points=x)$y; $
  sig_sq.nw <- sum((y_pert - est.nw[, i])^2) / (length(y) -
    ↪ sum(diag(S.nw)))
se.nw[, i] <- sqrt(sig_sq.nw * diag(S.nw \%*\% t(S.nw)))}
```

Smoothing splines and penalized regression

High-order polynomials do not work, so splines are used. We discuss splines *without* having to specify the knots. Find $\arg \min_{m \in C^0(\mathbb{R})} \sum_{i=1}^n (Y_i - m(x_i))^2 + \lambda \int_{\mathbb{R}} m''(z)^2 dz$. Note that the minimizer is *finite dimensional* — it is a cubic spline that can be computed using a set of basis functions $m_\lambda(x) = \sum_{j=1}^n \beta_j B_j(x)$ or $||Y - B\beta||^2 + \lambda \beta^\top \Omega \beta \Rightarrow \hat{\beta} = (B^\top B + \lambda \Omega)^{-1} B^\top Y$. Choose λ on the scale of $df = \text{tr}(S_\lambda)$. Note that this is Ridge-type regression, which saves us from being overparametrized (n points, n parameters). In the exam, this is *not* considered "standard" least squares.

```
fit.ss<-smooth.spline(x, y_pert, spar=0.6) # attr. cvcrit =
  ↪ loocv
fit.ss.pred[, i]<-predict(fit.ss, newdata=x)$y; $
```

Cross Validation

Let $(X_1, Y_1), \dots, (X_n, Y_n)$ i.i.d $\sim P$. We would like to compute $\mathbb{E}_{(X_{new}, Y_{new})} [\rho(Y_{new}, m_{train}(X_{new}))]$.

Constructing cross-validation datasets

Approaches include
Validation set: —
Leave-one-out CV: $n^{-1} \sum_{i=1}^n \rho(Y_i, \hat{m}_{n-1}^{(-i)}(X_i))$ ca. unbiased.
K-fold CV: $K^{-1} \sum_{i=1}^K |B_k|^{-1} \sum_{i \in B_k} \rho(Y_i, \hat{m}_{n-|B_k|}^{(-B_k)}(X_i))$.
Smaller variance than 1-CV.
Random division: Like K-fold, but build B_k by sampling without replacement ($\approx 10\%$). Usually fastest.

Tricks using hat matrix

For linear fitting operators and the loss $\rho(y, x) = (y - x)^2$ we can exploit the hat matrix and get the full 1CV result in a single step using

$$n^{-1} \sum_{i=1}^n \left(Y_i - \hat{m}_{n-1}^{(-i)}(X_i) \right)^2 = n^{-1} \sum_{i=1}^n \left(\frac{Y_i - \hat{m}(X_i)}{1 - S_{ii}} \right)^2.$$

It can be cheaper to just compute $\text{tr}(S)$ (instead of all S_{ii}), which leads to the *generalized cross-validation*

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

The two equations coincide if $S_{ii} = c \forall i$.

Bootstrap

Efron's *parametric* and *nonparametric bootstrap* can be described as "simulating from an estimated model" and can be used for *statistical inference* (*confidence intervals and testing*) and *estimating the predictive power of a model or algorithm*.

Nonparametric Bootstrap

Let $Z_{1:n}$ i.i.d $\sim P$ with $Z_i = (X_i, Y_i), X_i \in \mathbb{R}^p, Y_i \in \mathbb{R}$, and let $\theta_n = g(Z_{1:n})$ be an estimator. We would like to know the *distribution* of θ_n . We approximate P by the *empirical distribution* \hat{P}_n that assigns $\mathbb{P}[X_i] = 1/n \forall i$. Then we can repeatedly sample $Z_{1:n}^*$ i.i.d. $\sim \hat{P}_n$ and compute $\hat{\theta}_n^* = g(Z_{1:n}^*)$. The histogram (or any density estimator) then describes the distribution of $\hat{\theta}_n^*$. The algorithm reads

- a) Sample (with replacement) $Z_{1:n}^*$ i.i.d $\sim \hat{P}_n$.
- b) Compute the bootstrapped estimator $\hat{\theta}_n^* = g(Z_{1:n}^*)$.
- c) Repeat B times to obtain $\hat{\theta}_n^{*1:B}$.
- d) Approximate $\mathbb{E}^*[\hat{\theta}_n^*] \approx B^{-1} \sum_{i=1}^B \hat{\theta}_n^{*i}$ and $Var^*(\hat{\theta}_n^*) \approx (B - 1)^{-1} \sum_{i=1}^B \left(\hat{\theta}_n^{*i} - B^{-1} \sum_{j=1}^B \hat{\theta}_n^{*j} \right)^2$. Then α -quantile of $\hat{\theta}_n^* \approx$ empirical α -quantile of $\hat{\theta}_n^{*1:B}$.

Central limit theorem

Let X_i be a random variable with $\mathbb{E}[X_i] = 0$ and $Var(X_i) = \sigma^2$. Then $n^{-1} \sum_{i=1}^n X_i \xrightarrow{n \rightarrow \infty} \mathcal{N}(\mu, \sigma^2/n)$.

Bootstrap consistence

Consistency of the bootstrap typically holds if the limiting distribution of $\hat{\theta}_n$ is Normal and if $Z_{1:n}$ are i.i.d. Mathematically, for an increasing sequence a_n and $\forall x$, $\mathbb{P}[a_n(\hat{\theta}_n - \theta) \leq x] - \mathbb{P}^*[a_n(\hat{\theta}_n^* - \hat{\theta}_n) \leq x] \xrightarrow{P} 0$ as $n \rightarrow \infty$. Then $Op^*(\hat{\theta}_n^*)/Op(\hat{\theta}_n) \xrightarrow{P} 1$ with $Op \in \{Var, \mathbb{E}\}$.

Bootstrap confidence intervals

Given bootstrap consistence, we can compute confidence intervals:

- i quantile: $[q_{\hat{\alpha}^*}(\alpha/2), q_{\hat{\alpha}^*}(1 - \alpha/2)]$
- ii rev. quantile: $[\hat{\theta} - q_{\hat{\alpha}^*-\hat{\alpha}}(1 - \alpha/2), \hat{\theta} - q_{\hat{\alpha}^*-\hat{\alpha}}(\alpha/2)]$
- iii normal: $2\hat{\theta} - \bar{\hat{\theta}}^* \pm q_X(1 - \alpha/2) \cdot s.d(\hat{\theta})$ - corrects for bias $\hat{\theta} - \hat{\alpha}^*, X \sim \mathcal{N}(0, 1)$

Note $\hat{q}_\alpha = q_\alpha^* - \hat{\theta}_n$ with $q_\alpha^* = \alpha$ -bootstrap quantile of $\hat{\theta}_n^*$. Thus $[\hat{\theta}_n - \hat{q}_{1-\alpha/2}, \hat{\theta}_n - \hat{q}_{\alpha/2}] = [2\hat{\theta}_n - q_{1-\alpha/2}^*, 2\hat{\theta}_n - q_{\alpha/2}^*]$.

```
require("boot")
tm <- function(x, ind) {mean(x[ind], trim = 0.1)}
res.boot <- boot(data=sample, statistic=tm, R=10000,
  sim="ordinary")
# 'basic'='rev.quant.', 'norm'='normal', 'perc'='quant.'
boot.ci(res.boot, conf=0.95, type=c("basic", "norm", "perc"))
?quantile ?qnorm
```

```
require(MASS); fit.gamma <- fitdistr(boogg, "gamma")
par.est <- fit.gamma$estimate $ # for parametric gen
boot.est <- matrix(NA, nrow=R, ncol=1)
for (i in 1:R) {
  boogg.s <- rgamma(N, shape=par.est[1], rate=par.est[2])
  # boogg.s <- sample(boogg, N, replace=T) # NP
  boot.est[i] <- quantile(boogg.s, probs=0.75)
}; a <- 0.05
# QUANTILE
quantile(boot.est, probs=c(a/2, 1-a/2))
# NORMAL APPROXIMATION
mean.est <- mean(boot.est)
sd.hat <- sqrt(1/(R-1)*sum((boot.est-mean.est)^2))
2*est-mean.est + c(-1, +1)* qnorm(1-a/2)*sd.hat
# REVERSED QUANTILE
est - quantile(boot.est-est, probs=c(1-a/2, a/2))
```

Double bootstrap

Idea: Find α' s.t. actual coverage of bootstrap CI $I^*(1 - \alpha')$ is equal to α .

- i Draw BS sample Z^* . Sample from Z^* to obtain Z^{**} . Compute CI $I^{**}(1 - \alpha)$ for $\hat{\theta}^*$ based on B draws Z^{**} . Compute coverage of $\hat{\theta}$ by I^{**} (1 or 0).
- ii Repeat i) M times to obtain M coverage values. Compute mean to obtain actual coverage of I^{**} .
- iii Adjust α and repeat previous steps until coverage $(I^{**}(1 - \alpha')) = 1 - \alpha$. Use CI $I^*(1 - \alpha')$

Parametric Bootstrap

Assume $Z = (Z_1, \dots, Z_n)$ i.i.d. $\sim P_\theta$. Fit $\hat{\theta} = MLE(Z, P)$ and generate samples $Z^* \text{ i.i.d. } \sim P_{\hat{\theta}}$. Usually better than nonparametric version when $P_{\hat{\theta}}$ is a good fit (e.g. known model structure P) and few data points available.

```
require(MASS); mle <- fitdistr(boogg, "gamma")
fun.theta <- function(x) {quantile(x, probs = 0.75)}
fun.gen <- function(x, mle)
  <- {rgamma(length(x), shape=mle[1], rate=mle[2])}
res.boot <- boot(data, fun.theta, R=1000, sim="parametric",
  <- ran.gen=fun.gen, mle=fit.gamma$estimate); $
```

Bootstrap error estimate

Generalization error (loss $\rho(y, m(x))$) of model m (fitted to full data set) can be estimated by fitting models $m^{*,i}$ to bootstrap samples and computing

- errors on full data set $e^{*,i} = n^{-1} \sum_{i=1}^n \rho(y, m^{*,i}(x_i))$
- OOB errors $e_{ob}^{*,i} = n^{-1} \sum_{i=1}^{n_{ob,i}} \rho(y_{ob,i}, m^{*,i}(x_{ob,i}))$

The error of m is then approximated by $R^{-1} \sum_{i=1}^R e^{*,i}$.

Classification

Given $(X_1, Y_1), \dots, (X_n, Y_n)$ i.i.d. with $Y_i \in \{0, \dots, J-1\}$, determine $\pi_j(x) = \mathbb{P}[Y = j | X = x] \forall j = 0, 1, \dots, J-1$. The optimal classifier is the *Bayes classifier*, which is simply $C_{Bayes}(x) = \arg \max_{0 \leq j \leq J-1} \pi_j(x)$. Then, the zero-one test set error is called *Bayes risk*, i.e. $\mathbb{P}[C_{Bayes}(X_{new}) \neq Y_{new}]$.

Discriminant analysis

Linear case: Assume $(X|Y = j) \sim \mathcal{N}_p(\mu_j, \Sigma)$, $\mathbb{P}[Y = j] = p_j$, and $\sum_{j=0}^{J-1} p_j = 1$. Then by Bayes formula

$$\pi_j(x) = \frac{f_X|Y=j(x) \cdot p_j}{\sum_{k=0}^{J-1} f_X|Y=k(x) \cdot p_k}$$

with each $f_X|Y=j$ a Gaussian $\mathcal{N}(\mu_j, \Sigma_j)$. We can estimate μ_j and even Σ/Σ_j using closed formulas, but we also need priors for Y_i , which often is picked as $p_j = n_j/n$. This results in $\hat{\delta}_j(x) = (x - \hat{\mu}_j/2)^\top \Sigma^{-1} \hat{\mu}_j + \log(\hat{p}_j)$ with (linear in x) decision boundaries $\hat{\delta}_j(x) = \hat{\delta}_j, \hat{p}_j(x) \geq 0$ and $C(x) = \arg \max_j \hat{\delta}_j(x)$.

Quadratic case: Now we assume different Σ_j for each class and obtain quadratic decision boundaries $\hat{\delta}_j(x) = -\log(\det(\hat{\Sigma}_j))/2 - (x - \hat{\mu}_j)^\top \hat{\Sigma}_j^{-1} (x - \hat{\mu}_j)/2 + \log(\hat{p}_j)$. The price: $J \cdot p(p+1)/2$ parameters (for all Σ s) vs. $p(p+1)/2$ for a single Σ .

Logistic regression for binary classification

Given some model $g: \mathbb{R}^p \rightarrow \mathbb{R}$ (e.g. a linear model) we can use the logistic transform $\pi \mapsto \log(\pi/(1 - \pi))$ to get probabilities: $\log(\pi(x)/(1 - \pi(x))) = g(x)$ and $\pi(x) = 1/(1 + \exp(-g(x)))$. This implies $Y_i \sim \text{Bernoulli}(\pi(x_i))$ (e.g. weighted coin flip). The likelihood is $L(\beta; (X_i, Y_i)_{i=1:n}) = \prod_{i=1}^n \pi(x_i)^{Y_i} (1 - \pi(x_i))^{1-Y_i}$. We typically estimate β using e.g. (Newton's) gradient descent (due to a non-linear problem). As $n \rightarrow \infty$ we can asymptotically compute the standard errors $s.e.(\hat{\beta}_j)$ and t-test statistics $\hat{\beta}_j/s.e.(\hat{\beta}_j) \sim \mathcal{N}(0, 1)$ (under H_0 , $j: \beta_j \neq 0$).

```
fit <- glm(Y~., data=data, family="binomial")
mean((predict(fit, type="response") > 0.5) == data$Y)$
```

Linear predictors

Note that both *LDA* and *Logistic regression* are linear in the prediction variables. For LDA that comes from the Gaussian assumption (i.e. "linearization" of the true distribution), for Logistic regression it comes from the linear log-odds function.

Multiclass case ($J > 2$)

- a) J classes $\rightarrow J$ binary variables: $\hat{\pi}_j(x) = \frac{\hat{p}_{i,j}(x)}{\sum_{j=0}^{J-1} \hat{\pi}_j(x)}$
- b) Using *multinomial distribution* (parametric linear logistic) (see multinom)
- c) "Reference class" $\log(\pi_j(x)/\pi_0(x)) = g_j(x)$
- d) Pairwise 1-vs-1, fitting $\binom{J}{2} \cdot p$ parameters
- e) Exploiting "ordered" classes with proportional odds

Flexible regression and classification methods

We fight the *curse of dimensionality* by making some structural assumptions (although staying with methods $g(\cdot): \mathbb{R}^p \rightarrow \mathbb{R}$ of nonparametric nature).

1. Additive models

Decompose multivariate function in bias plus sum of univariate functions, i.e. $g_{add}: \mathbb{R}^p \rightarrow \mathbb{R}, x \mapsto g_{add}(x) = \mu + \sum_{j=1}^p g_j(x_j)$ with $g_j(\cdot): \mathbb{R} \rightarrow \mathbb{R}, \mathbb{E}[g_j(X_j)] = 0$. Note that the zero-mean requirement for each $g_j(\cdot)$ makes the problem well posed. This approach is a generalization of linear models, and similarly can not model interaction terms $g_{j,k}(x_j, x_k)$. Due to the way they are constructed, additive linear models avoid the *curse of dimensionality*!

To construct the models, let S_j be a smoothing technique (e.g. *Nadaraya-Watson Gaussian kernel estimators*). Then, the **backfitting** algorithm works as follows:

- Compute $\hat{\mu} = n^{-1} \sum_{i=1}^n Y_i$ and initialize $\hat{g}_j(\cdot) := 0$.
- Cycle through the indices $j = 1, 2, \dots, p, 1, 2, \dots, p, 1, 2, \dots$ and update $\hat{g}_j = S_j(Y - \hat{\mu} - \sum_{k \neq j} \hat{g}_k)$. Stop each function at convergence.
- Normalize the functions: $\tilde{g}_j(\cdot) = \hat{g}_j(\cdot) - n^{-1} \sum_{i=1}^n \hat{g}_j(X_{ij})$.

This basically makes the algorithm repeatedly solve the 1-dimensional fitting problem. The algorithm may be slow but often works and can use any 1-dimensional fitting technique.

When fitting Additive models in R with the function `gam`, the smoothers S_j penalized regression spline, and the degrees of freedom for each spline (i.e. each variable) will be determined through cross-validation.

```
fit <- gam(Y ~ s(x1) + s(x2) + ..., data=data)
plot(fit, pages=1, shade=TRUE)
sfsmsic::TA.plot(fit, labels="o")
```

2. Multivariate adaptive regression splines

MARS

$g(x) = \mu + \sum_{m=1}^M \beta_m h_m(x) = \sum_{m=0}^M \beta_m h_m$ Find $h \in \mathcal{M}$ functions by forward selection and pruning:

- i Initialize $\mathcal{M} = \{h_0 = 1\}, \beta_0 = \bar{Y}$
- ii For $r = 1, 2, \dots$: Find best pair (most reduction of RSS) $h_{2r-1} = h_l \circ x(x_j - x_{i,j}) + h_{2r} = h_l(\cdot) \times (x_{i,j} - x_j) +$, where $h_l \in \mathcal{M}$ does not already depend on x_j . Estimate β_{2r-1}, β_{2r} by LS. Add h_{2r-1}, h_{2r} to \mathcal{M} .
- iii Repeat until M large enough. Prune by repeatedly removing one function from pairs h_{2r-1}, h_{2r} (least increase in RSS). Stop when GCV score is optimized. $(x_j - d)_+ = x_j - d$, if $x_j \geq d$, zero otherwise.

```
require("earth");
fit <- earth(formula=y~., data=data, degree=2)
plotmo(fit, degree2=FALSE, caption="main effects")
```

Neural networks

$g(x)_k = f_0(\alpha_k + \sum_{h=1}^q w_{hk} \sigma(\tilde{\alpha}_h + \sum_{j=1}^p \tilde{w}_{hj} x_j)) \forall k = 1..J$, q hidden nodes, J output dim. Activation $\sigma(t) = \frac{\exp t}{1 + \exp t}$. For regression, f_0 identity, for classification $f_0 = \sigma$ and $C_{NN} = \arg \max_j g_j(x)$. Many other architectures possible, e.g. including a component directly connecting input to output by linear regression.

```
library(nnet); ?nnet ?ppr
```

3. Trees

Classification and Regression Trees

Let $g_{tree}(x) = \sum_{r=1}^M \beta_r 1_{[x \in \mathcal{R}_r]}$, where $\mathcal{P}\{\mathcal{R}_1, \dots, \mathcal{R}_M\}$ is a partition of \mathbb{R}^p . The function is piecewise constant. When partition is given, estimate $\hat{\beta}_r = \sum_{i=1}^n Y_i 1_{[x \in \mathcal{R}_r]} / \sum_{i=1}^n 1_{[x \in \mathcal{R}_r]}$. For multiclass classification $\hat{\pi}_j(x) = \sum_{i=1}^n 1_{[Y_i=j]} 1_{[x \in \mathcal{R}_r]} / \sum_{i=1}^n 1_{[x \in \mathcal{R}_r]}$ for $x \in \mathcal{R}_r$. Greddy algorithm to find axes parallel partition:

- i Initialize $M = 1$ subset $\mathcal{P} = \{\mathcal{R} = \mathbb{R}^p\}$
- ii Split \mathcal{R} at d in dimension j , where d is from the set of midpoints of observed values. Select j, k s.t. neg. log-likelihood decrease is maximized by refinement.
- iii Apply ii) to one cell of the current partition (select like above). Add the resulting two cells and remove the refined one.
- iv Iterate iii) until until specified max. partition size is achieved.
- v Prune tree by removing leaves resulting in smalles increase in some (CV) metric.

We define the size of a tree \mathcal{T} as the number of leaves (1 + +cuts). For some goodness-of-fit measure $\mathcal{R}(\mathcal{T})$ (e.g. SSE, NLL), the cost-complexity measure is $\mathcal{R}_\alpha(\mathcal{T}) = \mathcal{R}(\mathcal{T}) + \alpha \text{size}(\mathcal{T})$. For some α , we thus choose $\mathcal{T}(\alpha) = \arg \min_{\mathcal{T} \subset \mathcal{T}_M} \mathcal{R}_\alpha$ is then chosen by CV. 1 s.e. rule: Choose smalles tree s.t. its performance is at most one standard error larger than the minimal one.

```
library(rpart); require(rpart.plot);
# cp = alpha
rp <- rpart(y~., data = data, control=rpart.control(cp=0.0,
  <- minsplit=30))
plotcp(rp); cps = printcp(tree)
nid <- 10 # select from plot by ise rule
cp.opt <- cps[which(cps[, 'nsplit']==nid), 'CP']
pruned.tree = prune.rpart(tree, cp=cp.opt)

rf <- randomForest(Boston.train, y.train, xtest=Boston.test,
  <- ytest.y.test, ntree=500, mtry=ncol(Boston.train))
mean(rf$rmse)
mean(rf$test$rmse) $
```

4. Ridge and Lasso

Trade-off between $\|\beta\|_1$ and $\|\beta\|_2$ regularization, i.e. $L(\lambda_1, \lambda_2, \beta) = \|Y - X\beta\|_2^2 + \lambda_2 \|\beta\|_2^2 + \lambda_1 \|\beta\|_1$ with $\alpha = \lambda_2/(\lambda_1 + \lambda_2)$. $\alpha = 1 \Leftrightarrow \lambda_2 = 1$ means full L_2 -regularization, i.e. *Lasso*, otherwise *Ridge*.

```
fit <- glmnet(x.train, y.train, lambda=100, alpha=1)
fit$beta
fit_2 <- cv.glmnet(x.train, y.train, lambda=grid.lambda,
  <- alpha=1, type.measure="mse")
mean((predict(fit_2, newx=x.test, s="lambda.min") -
  <- y.test)^2)
fit_3 <- cv.glmnet(x.train, y.train, lambda=grid.lambda,
  <- alpha=0, type.measure="mse")
mean((predict(fit_3, newx=x.test, s=fit_3$lambda.min) -
  <- y.test)^2)
```

Bagging and Boosting

Bagging and Subbagging

Bootstrap aggregating (bagging) (mostly on trees), uses $\hat{g}(\cdot): \mathbb{R}^p \rightarrow \mathbb{R}$ and ensembles them (which comes at the loss of interpretability).

- i Generate bootstrap sample $(X_1^*, Y_1^*), \dots, (X_n^*, Y_n^*)$ and compute $\hat{g}_i^*(\cdot)$. Repeat B times.
- ii Aggregate bootstrap estimates with $\hat{g}_{Bag}(\cdot) = B^{-1} \sum_{i=1}^B \hat{g}_i^*(\cdot) \approx \mathbb{E}^*[\hat{g}^*(\cdot)]$.

Note that $\hat{g}_{Bag}(\cdot) = \hat{g}(\cdot) + \frac{\mathbb{E}^*[\hat{g}^*(\cdot)] - \hat{g}(\cdot)}{\text{bootstrap bias estimate}}$. We can reduce variance at price of higher bias (at least for trees). In fact, for many x , $\text{Var}(\hat{g}_{Bag}(x)) < \text{Var}(\hat{g}(x))$. We can use larger trees (higher variance) to balance the bias-variance trade-off. For **Subsample aggregating** (Subbagging), we draw $(X_1^*, Y_1^*), \dots, (X_m^*, Y_m^*)$ without replacement (e.g. with $m = \lfloor n/2 \rfloor$), which can be cheaper overall and is equivalent to Bagging in some simple settings.

L2 Boosting

Similar to Bagging, iterates on a "base-learner" by continually adding a fit on the residuals.

- i Get first fit $\hat{g}_1(\cdot)$ by fitting on the full data. Compute residuals $U_i = Y_i - \hat{g}_1(X_i)$ and let $f_1(\cdot) = \nu \hat{g}_1(\cdot)$ with $0 < \nu \leq 1$ (typically $\nu = 0.1$).
- ii For $m = 2, 3, \dots, M$ fit $\hat{g}_m(\cdot)$ on residuals U_i and set $\hat{f}_m(\cdot) = \hat{f}_{m-1}(\cdot) + \nu \hat{g}_m(\cdot)$ (and update residuals using $\hat{f}_m(\cdot)$).

The main tuning parameter is the stopping point M . Boosting increases the bias and can be used to ensemble trees to fit more complex data. See e.g.

```
?mboost; ?xgboost; ?gbm;
```