# Statistical learning for biological data

David Causeur
*Department of Statistics and Computer Science*
*Agrocampus Ouest*
*IRMAR CNRS UMR 6625*
*http://www.agrocampus-ouest.fr/math/causeur/*

# Outline

# Outline

# Impact of correlation on the least-squares fit

Case $p = 1$: $\mathbb{E}_x(Y) = \beta_0 + \beta x$

$$\text{Var}(\hat{\beta}) = \frac{\sigma^2}{n} \frac{1}{S_x^2}, \text{ where } S_x^2 = \frac{\sum_{i=1}^{n}(x_i - \bar{x})^2}{n}$$

Estimation accuracy depends on the spread of the $x-$profile.

## Impact of correlation on the least-squares fit

Case $p = 2$: $\mathbb{E}_x(Y) = \beta_0 + \beta_1 x_1 + \beta_2 x_2$

$$\text{Var}(\hat{\beta}) = \frac{\sigma^2}{n} \left[ \begin{array}{cc} S_1^2 & S_{12} \\ S_{12} & S_2^2 \end{array} \right]^{-1} = \left[ \begin{array}{cc} \text{Var}(\hat{\beta}_1) & \text{Cov}(\hat{\beta}_1, \hat{\beta}_2) \\ \text{Cov}(\hat{\beta}_1, \hat{\beta}_2) & \text{Var}(\hat{\beta}_2) \end{array} \right]$$

Hence,

$$\text{Var}(\hat{\beta}_1) = \frac{\sigma^2}{n} \frac{1}{S_1^2} \frac{1}{1 - r_{12}^2}, \ \text{Var}(\hat{\beta}_2) = \frac{\sigma^2}{n} \frac{1}{S_2^2} \frac{1}{1 - r_{12}^2}$$

Estimation accuracy also depends on the correlation (here $r_{12}$) within the $x-$profile.

▶ Illustration using high-dimensional data in R

# Impact of correlation on the least-squares fit

Ordinary Least Squares fitting

$$\hat{\beta} \text{ minimizes } \sum_{i=1}^{n} \Big( Y_i - \bar{Y} - \big[ \beta_1(x_{i1} - \bar{x}_1) + \ldots + \beta_p(x_{ip} - \bar{x}_p) \big] \Big)^2$$

A closed-form solution ... provided $S_x^{-1}$ exists

$$\hat{\beta} = S_x^{-1} s_{xy}.$$

$\hat{\beta}$ is unbiased with variance

$$\text{Var}_x(\hat{\beta}) = \frac{\sigma^2}{n} S_x^{-1}.$$

In high-dimension, $\text{Var}_x(\hat{\beta}_j)$ can be very large ...

# A biased alternative: penalized regression

Least-squares optimization under constraint

$$
\begin{aligned}
\text{RSS}(\beta) \;=\; & \sum_{i=1}^{n} \big( Y_i - \bar{Y} - \beta_1 [x_i^{(1)} - \bar{x}^{(1)}] - \ldots - \beta_p [x_i^{(p)} - \bar{x}^{(p)}] \big)^2, \\
& \text{subject to } \sum_{j=1}^{p} \beta_j^2 \leq \kappa
\end{aligned}
$$

# A biased alternative: penalized regression

Least-squares optimization under constraint (equivalent form)

$$
\begin{aligned}
\mathrm{RSS}(\beta) \;=\; & \sum_{i=1}^{n} \big( Y_i - \bar{Y} - \beta_1 [x_i^{(1)} - \bar{x}^{(1)}] - \ldots - \beta_p [x_i^{(p)} - \bar{x}^{(p)}] \big)^2, \\
& \text{subject to } \sum_{j=1}^{p} \beta_j^2 = \kappa
\end{aligned}
$$

# A biased alternative: penalized regression

Least-squares optimization under constraint (equivalent form)

$$\text{RSS}(\beta) = \sum_{i=1}^{n} \big(Y_i - \bar{Y} - \beta_1[x_i^{(1)} - \bar{x}^{(1)}] - \ldots - \beta_p[x_i^{(p)} - \bar{x}^{(p)}]\big)^2,$$

$$\text{subject to } \sum_{j=1}^{p} \beta_j^2 = \kappa$$

The *Lagrange multiplier* trick:

$$\text{RSS}(\beta; \lambda) = \sum_{i=1}^{n} \big(Y_i - \bar{Y} - \beta_1[x_i^{(1)} - \bar{x}^{(1)}] - \ldots - \beta_p[x_i^{(p)} - \bar{x}^{(p)}]\big)^2$$

$$+ \ \lambda \sum_{j=1}^{p} \beta_j^2$$

# A biased alternative: penalized regression

Let's do it when $p = 1$:

$$
\begin{aligned}
\text{RSS}(\beta; \lambda) &= \sum_{i=1}^{n} \big(Y_i - \bar{Y} - \beta[x_i - \bar{x}]\big)^2 + \lambda\beta^2, \\
&= \sum_{i=1}^{n}(Y_i - \bar{Y})^2 - 2\beta\sum_{i=1}^{n}(Y_i - \bar{Y})(x_i - \bar{x}) \\
&\quad + \beta^2\big[\lambda + \sum_{i=1}^{n}(x_i - \bar{x})^2\big], \\
&= n\big(s_Y^2 - 2\beta s_{xY} + \beta^2\big[s_x^2 + \frac{\lambda}{n}\big]\big).
\end{aligned}
$$

# A biased alternative: penalized regression

Let's do it when $p = 1$:

$$\text{RSS}(\beta; \lambda) \;\;=\;\; n\big(s_Y^2 - 2\beta s_{xY} + \beta^2 \big[s_x^2 + \frac{\lambda}{n}\big]\big)$$

Let's differentiate w.r.t $\beta$:

$$\frac{\partial \text{RSS}}{\partial \beta}(\beta; \lambda) \;\;=\;\; -2n\big(s_{xY} - \beta \big[s_x^2 + \frac{\lambda}{n}\big]\big).$$

By equating to zero:

$$\hat{\beta}_\lambda = \frac{s_{xY}}{s_x^2 + \frac{\lambda}{n}}$$

# A biased alternative: penalized regression

Ridge estimation of $\beta$:

$$\hat{\beta}_\lambda = \left[ S_{xx} + \frac{\lambda}{n} I_p \right]^{-1} s_{xy}$$

The ridge estimator

- is biased: $b_\lambda = \mathbb{E}_x \left[ \hat{\beta}_\lambda - \beta \right] \neq 0$      [shrinkage estimation]
- but has a smaller variance than $\hat{\beta}$:
  - $\lambda \approx 0$: small bias, large variance
  - $\lambda \gg 0$: large bias, small variance

$\lambda$ is chosen so that MSEP$(\lambda)$ is as small as possible

▶ Ridge regression using R

# Shrinkage and selection: *LASSO* regression

Least-squares optimization under constraint

$$
\begin{aligned}
\text{RSS}(\beta; \lambda) \;=\; & \sum_{i=1}^{n} \big( Y_i - \bar{Y} - \beta_1 [x_i^{(1)} - \bar{x}^{(1)}] - \ldots - \beta_p [x_i^{(p)} - \bar{x}^{(p)}] \big)^2 \\
& + \; \lambda \sum_{j=1}^{p} |\beta_j|
\end{aligned}
$$

# Shrinkage and selection: *LASSO* regression

Let's try to do it when $p = 1$:

$$
\begin{aligned}
\text{RSS}(\beta; \lambda) &= \sum_{i=1}^{n} \left( Y_i - \bar{Y} - \beta[x_i - \bar{x}] \right)^2 + \lambda|\beta|, \\
&= n\left( s_Y^2 - 2\beta s_{xY} + \beta^2 s_x^2 + \frac{\lambda}{n}|\beta| \right), \\
&= n\left\{ \begin{array}{ll} s_Y^2 - 2\beta[s_{xY} + \frac{\lambda}{2n}] + \beta^2 s_x^2 & \text{for} \quad \beta < 0, \\ s_Y^2 - 2\beta[s_{xY} - \frac{\lambda}{2n}] + \beta^2 s_x^2 & \text{for} \quad \beta \geq 0 \end{array} \right.
\end{aligned}
$$

# Shrinkage and selection: *LASSO* regression

Let's try to do it when $p = 1$:

$$\text{RSS}(\beta; \lambda) = n \left\{ \begin{array}{ll} s_Y^2 - 2\beta[s_{xY} + \frac{\lambda}{2n}] + \beta^2 s_x^2 & \text{for} \quad \beta < 0, \\ s_Y^2 - 2\beta[s_{xY} - \frac{\lambda}{2n}] + \beta^2 s_x^2 & \text{for} \quad \beta \geq 0 \end{array} \right.$$

Not differentiable w.r.t $\beta$ for $\beta = 0$...

$$\hat{\beta}_\lambda = \left\{ \begin{array}{ll} \frac{s_{xY} + \frac{\lambda}{2n}}{s_x^2} & \text{if} \quad s_{xY} \leq -\frac{\lambda}{2n}, \\ 0 & \text{if} \quad -\frac{\lambda}{2n} \leq s_{xY} \leq \frac{\lambda}{2n}, \\ \frac{s_{xY} - \frac{\lambda}{2n}}{s_x^2} & \text{if} \quad s_{xY} \geq \frac{\lambda}{2n} \end{array} \right.$$

▶ Implementation using R

# Shrinkage and selection: *LASSO* regression

Extension to a multivariate profile of *x*s: no closed-form expression for $\hat{\beta}_\lambda$

... estimation achieved using an iterative algorithm (e.g. cyclic coordinate descent)

The LASSO estimator is also biased

- $b_\lambda = \mathbb{E}_x \left[ \hat{\beta}_\lambda - \beta \right] \neq 0$         [shrinkage estimation]

- ... but $\hat{\beta}_\lambda$ has a smaller variance than $\hat{\beta}$

- ... and increasing $\lambda$ kills the $\beta$s         [selection]

▶ LASSO regression using R

# LASSO or Ridge?

Both can be very performant

LASSO regression models are sparse

Warning!! When the *x*s are highly correlated, the selection is unstable ...

Mixing ridge and lasso penalties (elastic net) may bring stability:

$$\hat{\beta}(\lambda, \alpha) \text{ minimizes } \mathsf{RSS}(\beta) + \lambda \Big[ \alpha \sum_{j=1}^{p} |\beta_j| + (1 - \alpha) \sum_{j=1}^{p} \beta_j^2 \Big]$$

Warning!! Two hyper-parameters to be tuned ...

# Ridge/LASSO logistic regression modeling

Aim: guessing production site of coffees using NIRS

- $Y$: production site named $y_1, \ldots, y_6$
- $x$ is a NIRS:

  $x = \big(x(\omega_1), \ldots, x(\omega_p)\big)'$, where $\omega_i$ is the $i$th wave number.

▶ Import data in R session

# Ridge/LASSO logistic regression modeling

Minimization of the penalized deviance

$$
\begin{aligned}
\mathcal{D}(\beta; \lambda) &= \mathcal{D}(\beta) + \lambda ||\beta||_2^2, \quad \text{[Ridge]} \\
\mathcal{D}(\beta; \lambda) &= \mathcal{D}(\beta) + \lambda ||\beta||_1, \quad \text{[Lasso]}
\end{aligned}
$$

where $||\beta||_2^2 = \sum_{k=2}^{6} \sum_{j=1}^{p} [\beta_j^{(k)}]^2$ and $||\beta||_1 = \sum_{k=2}^{6} \sum_{j=1}^{p} |\beta_j^{(k)}|$.

▶ Lasso estimation of a multinomial logistic regression model using R