# Linear Methods for Regression (Chap 3 ESC)

## I. Introduction.

A linear regression model assumes that the regression function $\mathbb{E}[Y|X]$ is linear in the inputs $X_1, ..., X_p$. Reminder : $\mathbb{E}[Y|X] = \varphi(x)$.

## II. Linear Regression Models and Least Squares.

. We have an input vector $X \in \mathbb{R}^p$ and we want to predict an output $Y \in \mathbb{R}$.

. Linear regression model : $Y = f(X) + \varepsilon$ where $f(x) = \beta_0 + \sum_{j=1}^{c} X_j \beta_j^*$ and $\mathbb{E}[\varepsilon] = 0$

The coefficients $\beta_1, ..., \beta_p$ are unknown and the variables $(x_j)$ can come from different sources : quantitative inputs, transformation of quantitative inputs $(log, \sqrt{\cdot}, ...^e)$, ..

. The model is linear in the parameters.

. Data : Collect $(y_1, x_1), ..., (y_n, x_n)$ where $\forall i \in [n]$ $y_i \in \mathbb{R}$ and $x_i \in \mathbb{R}^p$.

. From those data we estimate $\beta^* := (\beta_0^*, ..., \beta_p^*)$ (in ML we "learn").

. How ? Minimize the residual sum of squares : $\hat{\beta}_{LS} \in \text{argmin} \|Y - X\beta\|_2^2$.

. RSS $(\beta) = \sum_{i=1}^{n} (y_i - \beta_0' - \sum_{j=1}^{c} x_{ij} \beta_j)^2$. How to minimize ? Matrix notation !

. $X \in \mathbb{R}^{n \times (p+1)}$, matrix with each row being an input vector $x_i$ (with 1 in first position).

. $y \in \mathbb{R}^n$ is the vector of outputs and $\beta \in \mathbb{R}^{p+1}$ is the parameter to "learn".

$\hookrightarrow$ RSS$(\beta) = (y - X\beta)^T (y - X\beta)$ and $\frac{\partial RSS}{\partial \beta}(\beta) = -2X^T(y - X\beta)$ and $\frac{\partial^2 RSS}{\partial \beta^2}(\beta) = 2X^TX$.

$\hookrightarrow$ If $X$ is full rank (ie $\text{rank}(X) = p+1$) then $\text{rank}(X^TX) = p+1$ and $X^TX$ invertible.

Moreover $X^TX$ invertible ensures $X^TX$ is positive definite $(\forall u \in \mathbb{R}^{p+1} \ u^T X^T X u = \|Xu\|^2 > 0)$.

Hessian positive-definite ensures convexity of the function. Hence RSS is convex in $\beta$.

This implies that any critical point is the global minimizer. Hence to

minimize RSS it suffices to find $\hat{\beta}$ s.t. RSS$(\hat{\beta}_{LS}) = 0$ ie $X^TX\hat{\beta}_{LS} = X^Ty$.

$\hookrightarrow \hat{\beta}_{LS} = (X^TX)^{-1} X^Ty$. where $Y = X\beta^* + \varepsilon$.

$\hookrightarrow$ the predicted value at an input vector $x_{n+1}$ is $(1, x_{n+1,1}, ..., x_{n+1,p})^T \hat{\beta}$.

$\hookrightarrow$ The fitted values at the training inputs are $\hat{y} = X\hat{\beta} = X(X^TX)^{-1}X^Ty$.

. What happens if $X^TX$ is not invertible ?

$\hookrightarrow$ perfect Multicolinearity : one predictor is a linear combination of the others $\rightarrow$ remove it.

$\hookrightarrow$ High-dimension : $p+1 > n$ $\rightarrow$ Regularization techniques.

. What happens when mild multicolinearity : predictors have close to exact linear relationship.

$\hookrightarrow$ LS estimates for $\beta_j$ is well defined but have large variance $\rightarrow$ Regularization.

$\hookrightarrow$ Ridge regression for example.

### 1. Gauss-Markov theorem (why $\hat{\beta}^{LS}$ and not minimizing another criterion ?)

. The least squares estimate $\hat{\beta}^{LS}$ has the smallest variance among all linear unbiased estimates.

. Estimation of any linear combination of the parameters : $\theta = a^T\beta^*$.

$\hookrightarrow \hat{\theta}_{LS} = a^T\hat{\beta}_{LS} = a^T(X^TX)^{-1}X^Ty$. $\rightarrow \mathbb{E}[a^T\hat{\beta}_{LS}] = a^T(X^TX)^{-1}X^T\mathbb{E}[y] = a^T(X^TX)^{-1}X^TX\beta = a^T\beta^*$.

$\hookrightarrow$ If we have any other linear estimator $\tilde{\theta} = c^Ty$ that is unbiased $(\mathbb{E}[\tilde{\theta}] = a^T\beta^*)$ :

$$\mathbb{V}(\hat{\theta}_{LS}) \leq \mathbb{V}(\tilde{\theta})$$

For any estimator $\tilde{\theta}$ of $\theta^*$ : $MSE(\tilde{\theta}) = \mathbb{E}_{\theta^*}[(\tilde{\theta}-\theta^*)^2] = \mathbb{V}_{\theta^*}(\tilde{\theta}) + (\mathbb{E}[\tilde{\theta}]-\theta^*)^2$ .

$\quad$└ Variance + Squared bias. Gauss - Markov → The LS estimator has the smallest
$\qquad$ MSE among all unbiased linear estimators.

$\quad$└ However we may find a biased estimator with smaller MSE.

$\quad$└ Add a little bias for a huge reduction in variance.

$\quad$└ Any estimator that shrinks the coefficients of the LS estimator is biased.


## III. Shrinkage Methods

### 1. Ridge Regression.

$$\hat{\beta}_\lambda^R \in \underset{\beta \in \mathbb{R}^p}{\text{argmin}} \; \|Y - X\beta\|_2^2 + \lambda \|\beta\|_2^2 . \quad \text{(We omit } \beta_0 \text{ in the penalty).}$$

Idea : When there is multicolinearity, the coefficients of the parameter are poorly
determined by the OLS estimator. A wildly large positive coefficient on one variable
can be canceled by a similarly large negative coefficient on its correlated cousin.
Ridge imposes a size constraint on the coefficients → reduces this problem.

How to find $\hat{\beta}_\lambda^R$ ? → Objective is differentiable.

$\quad$└ $RSS(\beta) = (y - x\beta)^T (y - x\beta) + \lambda \beta^T \beta \longrightarrow \dfrac{\partial RSS}{\partial \beta}(\beta) = -2x^T(y - x\beta) + 2\lambda\beta$ .

$\quad$└ $\dfrac{\partial^2 RSS}{\partial \beta^2}(\beta) = 2(x^Tx + \lambda I_p)$ is positive definite because $\forall u \in \mathbb{R}^p : 2u^Tx^Txu + 2\lambda u^Tu \geqslant 0$ .

$\quad$└ $\beta \longmapsto RSS(\beta)$ is convex and thus $\hat{\beta}_\lambda^R$ satisfies $\dfrac{\partial RSS}{\partial \beta}(\hat{\beta}_\lambda^R) = 0 \longrightarrow \hat{\beta}_\lambda^R = (x^Tx + \lambda I_p)^{-1}x^Ty$ .

### 2. Lasso.

$$\hat{\beta}_\lambda^L \in \underset{\beta \in \mathbb{R}^p}{\text{argmin}} \; \|Y - X\beta\|_2^2 + \lambda \|\beta\|_1 . \quad \text{(We omit } \beta_0 \text{ in the penalty).}$$

$\quad$└ No closed form. Quadratic programming problem. Efficient algorithms with same
computational cost as for ridge.