# 1   Kernel Machines

Consider a linear model

$$Y_i = x_i^T \beta^0 + \varepsilon_i, \ i = 1, \ldots, n, \ x_i \in \mathbb{R}^p \text{ fixed}$$

where $\mathbb{E}\varepsilon = 0$, $\mathrm{Var}(\varepsilon) = \sigma^2 I_n$. We have

$$\begin{aligned}
\hat{\beta}^{\mathrm{ols}} &= \mathrm{argmin}_{\beta \in \mathbb{R}^p} \sum_{i=1}^n (Y_i - x_i^T \beta)^2 \\
&= \mathrm{argmin}_{\beta \in \mathbb{R}^p} ||Y - X\beta||^2 \\
&= (X^T X)^{-1} X^T Y.
\end{aligned}$$

Classical theory:

- $\hat{\beta}^{\mathrm{ols}}$ unbiased,

$$\mathrm{Var}(\hat{\beta}^{\mathrm{ols}}) = \sigma^2 (X^T X)^{-1} = i^{-1}(\beta^0)$$

  Where $i$ is the Fisher information.

- Cramér-Rao lower bound: if an estimator $\tilde{\beta}$ is unbiased then

$$\mathrm{Var}(\tilde{\beta}) - i^{-1}(\beta^0) \underbrace{\geq}_{\text{positive semi-definite}} 0.$$

- If $\varepsilon \sim \mathcal{N}(0, \sigma^2 I_n)$, then $\hat{\beta}^{\mathrm{ols}}$ is the MLE of $\beta^0$. Furthermore $\sqrt{n}(\hat{\beta}^{\mathrm{ols}} - \beta^0) \sim \mathcal{N}(0, n\sigma^2 (X^T X)^{-1})$. From this we can derive confidence intervals, hypothesis test, etc.

In a general model with parameter $\theta \in \mathbb{R}^p$, $n$ independent observations, under regularity, we have asymptotic normality, i.e $\sqrt{n}(\hat{\theta}^{\mathrm{MLE}} - \theta^0) \xrightarrow{d} \mathcal{N}(0, I^{-1}(\theta^0))$ (with $p$ fixed).

Question: what happens when $p$ is large relative to $n$?

- If $p > n$, $\hat{\beta}^{\mathrm{ols}}$ is not even defined.

- If $p \approx n$, $\mathrm{Var}(\hat{\beta}^{\mathrm{ols}})$ explodes since $X^T X$ is near singular.

- More generally, if $p, n \to \infty$ then asymptotic normality can break down.

Recall the bias-variance decomposition:

$$\begin{aligned}
\mathrm{mse}(\tilde{\beta}) &= \mathbb{E}_{\beta^0, \sigma^2} \left[ (\tilde{\beta} - \beta^0)(\tilde{\beta} - \beta^0) \right] \\
&= \mathbb{E}_{\beta^0, \sigma^2} \left\| \tilde{\beta} - \mathbb{E}\tilde{\beta} + \mathbb{E}\tilde{\beta} - \beta^0 \right\| \\
&= \mathrm{Var}(\tilde{\beta}) + \left\| \mathbb{E}(\tilde{\beta}) - \beta^0 \right\|^2.
\end{aligned}$$

We introduce bias to reduce the variance.

## 1.1    Ridge regression

Define

$$(\hat{\mu}_\lambda^R, \hat{\beta}_\lambda^R) = \text{argmin}_{(\mu,\beta) \in \mathbb{R} \times \mathbb{R}^p} \left[ ||Y - \mu\mathbf{1} - X\beta||^2 + \underbrace{\lambda||\beta||^2}_{\text{penalty for large } \beta} \right].$$

$\lambda$ is called a *regularisation* or *tuning* parameter. We shall assume the columns of $X$ have been standardised (mean 0, variance 1).

After standardisation, we can show that

$$\hat{\mu}_\lambda^R = \frac{1}{n} \sum_{i=1}^n Y_i = \bar{Y}.$$

Hence, if we replace $Y$ with $Y - \mathbf{1}\bar{Y}$ we can write

$$\hat{\beta}_\lambda^R = \text{argmin}_{\beta \in \mathbb{R}^p} \left[ ||Y - X\beta||^2 + \lambda||\beta||^2 \right]$$
$$= \underbrace{(X^T X + \lambda I_p)^{-1}}_{\text{always invertible}} X^T Y.$$

**Theorem 1.1.** *For $\lambda > 0$ sufficiently small,*

$$\mathbb{E}||\hat{\beta}^{ols} - \beta^0||^2 - \mathbb{E}||\hat{\beta}_\lambda^R - \beta^0||^2 > 0. \qquad (*)$$

*Proof.* We have
$$Y = X\beta^0 + \varepsilon.$$

The bias of $\hat{\beta}_\lambda^R$ is

$$\mathbb{E}(\hat{\beta}_\lambda^R - \beta^0) = (X^T X + \lambda I)^{-1} X^T X \beta^0 - \beta^0$$
$$= (X^T X + \lambda I)^{-1} (X^T X + \lambda I - \lambda I)\beta^0 - \beta^0$$
$$= -\lambda(X^T X + \lambda I)^{-1}\beta^0.$$

While we have variance

$$\text{Var}(\hat{\beta}_\lambda^R) = \mathbb{E} \left\| (X^T X + \lambda I)^{-1} X^T \varepsilon \right\|^2$$
$$= \sigma^2 \left[ (X^T X + \lambda I)^{-1} X^T X (X^T X + \lambda I)^{-1} \right].$$

Then $(*)$ becomes

$$\mathbb{E}||\hat{\beta}^{\text{ols}} - \beta^0||^2 - \mathbb{E}||\hat{\beta}_\lambda^R - \beta^0||^2$$
$$= \sigma^2 (X^T X)^{-1} - \sigma^2 (X^T X + \lambda I) X^T X (X^T X + \lambda I)^{-1}$$
$$\quad - \lambda^2 (X^T X + \lambda I)^{-1} \beta^0 (\beta^0)^T (X^T X + \lambda I)^{-1}$$
$$= \quad \vdots \qquad\qquad\qquad\qquad (\text{use SVD } X = UDU^T)$$
$$= \lambda(X^T X + \lambda I)^{-1} \left[ \sigma^2 \left\{ 2I_p + \lambda(X^T X)^{-1} \right\} - \lambda\beta^0(\beta^0)^T \right] (X^T X + \lambda I)^{-1}.$$

We want to show this is positive definite. This is equivalent to

$$\sigma^2 \left[ 2I + \lambda(X^TX)^{-1} \right] - \lambda\beta^0(\beta^0)^T > 0$$
$$\iff 2\sigma^2 I - \lambda\beta^0(\beta^0)^T > 0$$
$$\iff 2\sigma^2 ||z||^2 - \lambda(z^T\beta^0)^2 > 0 \quad \forall z \in \mathbb{R}^p. \tag{$\dagger$}$$

We also have $(z^T\beta^0)^2 \leq ||z||^2||\beta^0||^2$ by Cauchy-Schwartz. Hence ($\dagger$) holds for all $\lambda < \frac{2\sigma^2}{||\beta^0||^2}$. $\qquad\square$