

Assumption: $E(y|x)$ is linear in the inputs x_1, \dots, x_p .

* best for : small no. of training cases

low signal-to-noise ratio

sparse.

3.2. Linear Regression models and least squares.

$x^T = (x_1, x_2, \dots, x_p)$: input $\rightarrow y$: output.

linear regression model has the form:

$$f(x) = \beta_0 + \sum_{j=1}^p x_j \beta_j$$

β_j : unknown parameters or coefficients

x_j : come from different sources.

(i) quantitative inputs.

(ii) transforms of quantitative inputs (log, $\sqrt{\cdot}$, squareroot)

(iii) basic expansions such as: $x_2 = x_1^2$, $x_3 = x_1^3$

leading to polynomial representation.

(iv) numeric or "dummy" coding of the levels of qualitative inputs.

(v) interactions between variables ($x_3 = x_1 \cdot x_2$)

- Model is linear in parameters, irrespective of source.

data: $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$

parameters: $\beta = (\beta_0, \beta_1, \dots, \beta_p)^T$

features: $x_i = (x_{i1}, x_{i2}, \dots, x_{ip})^T$

Residual Sum of Squares: $RSS(\beta) = \sum_{i=1}^n (y_i - f(x_i))^2$

$$= \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p x_{ij} \beta_j \right)^2$$

- From statistical p.o.v., this criterion is reasonable if the training observations (x_i, y_i) represent independent random draws from their population.

- Even if x_i 's are not drawn randomly, criterion is still valid if $y_i | x_i$ are independent (i.e. y_i 's are conditionally independent)

- RSS: measures the average error of fit.

$$RSS(\beta) = (y - X\beta)^T (y - X\beta) \rightarrow \text{quadratic function in } p+1 \text{ parameters}$$

Differentiate wrt β :

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

$$\Rightarrow \frac{\partial^2 RSS}{\partial \beta \partial \beta^T} = 2X^T X$$

Assuming X has full column rank,

$\Rightarrow X^T X$ is positive definite.

We set first derivative to zero.

$$X^T(y - X\beta) = 0$$

To obtain unique solution:

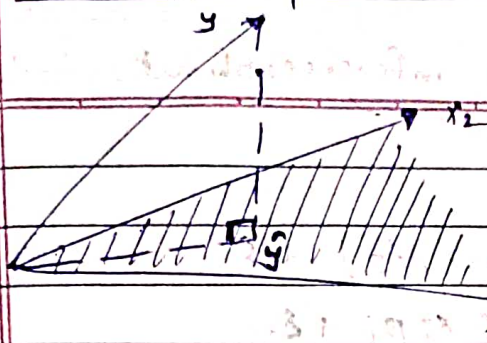
$$\Rightarrow \hat{\beta} = (X^T X)^{-1} X^T y$$

Fitted values of training input:

$$\hat{y} = X\hat{\beta} = X(X^T X)^{-1} X^T y$$

$$\hat{y} = Hy \quad \text{where } H = X(X^T X)^{-1} X^T \rightarrow \text{as it puts hat on } y.$$

* Geometrical representation:



We minimize $RSS(\beta) = \|y - X\beta\|^2$ by choose $\hat{\beta}$ so that, $y - \hat{y}$ is orthogonal to \mathbb{R}^n .

$\Rightarrow \hat{y}$ is orthogonal projection of y onto ~~the~~ subspace (column space of X)

As H computes \hat{y} , thus H is orthogonal projection matrix!

- If ~~not~~ all columns of X are linearly independent, then, then $X^T X$ is singular $\Rightarrow \hat{\beta}$ are not uniquely defined.

but, $\hat{y} = X\hat{\beta}$ ~~are~~ are still projection of y onto the column space of X , but there is more than one way to express that projection in terms of column vectors of X .

- To (make $X^T X$ non-singular: recoding) dropping redundant columns in X .

- If no of features $>$ no of observations.

\hookrightarrow Filtering
 \hookrightarrow regularization.

* To get sampling properties of $\hat{\beta}$:

Assumption: (1) y_i uncorrelated

(2) constant variance: $\text{var}(y_i) = \sigma^2$

(3) x_i are fixed (non-random).

$\text{Var}(\hat{\beta}) = (X^T X)^{-1} \sigma^2 \rightarrow$ var cov. matrix of $\hat{\beta}$.

$$\sigma^2 = E \left[\frac{1}{N-p-1} \sum_{i=1}^N (y_i - \hat{y}_i)^2 \right]$$

$$\sigma^2 = E(\hat{\sigma}^2)$$

* Additional assumptions to draw inferences about parameters of model

(1) $\varepsilon \sim N(0, \sigma^2)$

(2) $E(Y|X) = Y = E(Y|X_1, \dots, X_p) + \varepsilon$

$$Y = \beta_0 + \sum_{j=1}^p X_j \beta_j + \varepsilon$$

$$\Rightarrow \hat{\beta} \sim N(\beta, (X^T X)^{-1} \sigma^2)$$

$$\Rightarrow (N-p-1) \hat{\sigma}^2 \sim \sigma^2 \chi_{N-p-1}^2$$

$\Rightarrow \hat{\beta}$ & $\hat{\sigma}^2$ are statistically independent.

* Testing of hypothesis & confidence intervals of β_j

To test: $H_0: \beta_j = 0$ vs $H_1: \beta_j \neq 0$

$$t_j = \frac{\hat{\beta}_j}{\hat{\sigma} \sqrt{v_j}}$$

Where v_j : j th diagonal element of $(X^T X)^{-1}$

Under H_0 : $t_j \sim t_{N-p-1}$

→ rejected for large absolute value of t

→ if $\hat{\sigma}$ is replaced by σ , t_j will have standard normal distribution.

→ as sample size T , tail quantiles of t & Z are almost same, so Z is used.

* Test for significance of groups of coefficients simultaneously.

$$F = \frac{(RSS_0 - RSS_1) / (p_1 - p_0)}{RSS_1 / (N - p_1 - 1)}$$

where :



RSS_1 : residual sum of squares for the least squares fit of the bigger model with p_1+1 parameters.

RSS_0 : -- nested smaller model with p_0+1 parameters.

$p_1 - p_0$ parameters constrained to be 0.

- F statistics: measures the change in residual sum of squares per additional parameter in the bigger model, normalized by estimate of σ^2 .

Ho: smaller model is correct.

$$F \sim F_{p_1 - p_0, N - p_1 - 1}.$$

For large N ,

$$F \sim \chi^2_{p_1 - p_0}$$

* confidence interval for β_j

$$\left[\hat{\beta}_j - z_{(1-\alpha)/2} \sqrt{\hat{\sigma}^2 v_j}, \hat{\beta}_j + z_{(1-\alpha)/2} \sqrt{\hat{\sigma}^2 v_j} \right]$$

Where $z_{(1-\alpha)/2}$ is the $(1-\alpha)$ percentile of normal distribution

* confidence set for β

$$C_\beta = \{ \beta \mid (\hat{\beta} - \beta)^T X^T X (\hat{\beta} - \beta) \leq \hat{\sigma}^2 \chi^2_{p_0+1, (1-\alpha)} \}$$

where: $\chi^2_{p_0+1, (1-\alpha)}$ is the $(1-\alpha)$ th percentile of the chi-squared distribution on p_0+1 degrees of freedom.

3.2.2 The Gauss - Markov Theorem:

The least squares estimates of the parameters β have the smallest variance among all linear unbiased estimates.

~~Estimation:~~

Let $\theta = a^T \beta$ for

- focus on estimation of any linear combination of the parameters $\theta = a^T \beta$, e.g. predictions $f(x_0) = x_0^T \beta$ are of this form.

$$\hat{\theta} = a^T \hat{\beta} = a^T (X^T X)^{-1} X^T y$$

→ $a^T \hat{\beta}$ is unbiased for $a^T \beta$.

$$\begin{aligned} E(a^T \hat{\beta}) &= E(a^T (X^T X)^{-1} X^T y) \\ &= a^T (X^T X)^{-1} X^T X \beta \\ &= a^T \beta \end{aligned}$$

* Gauss Markov Thm:

If we have any other linear estimator $\tilde{\theta} = c^T y$ that is unbiased for $a^T \beta$ i.e. $E(c^T y) = a^T \beta$ then:

$$\text{Var}(a^T \hat{\beta}) \leq \text{Var}(c^T y)$$

⇒ Least squares estimator has the smallest MSE of all linear estimators with no bias.

~~But~~ BUT: there might exist a biased estimator with smaller MSE.

→ Least MSE estimator is preferred.



consider the prediction of the new response at input x_0 :

$$y_0 = f(x_0) + \varepsilon_0$$

Expected prediction error of estimate:

$$\hat{f}(x_0) = x_0^T \hat{\beta}$$

is:

$$\begin{aligned} E(y_0 - \hat{f}(x_0))^2 &= \sigma^2 + E(x_0^T \hat{\beta} - f(x_0))^2 \\ &= \sigma^2 + \text{MSE}(\hat{f}(x_0)) \end{aligned}$$

(Variance of y_0).