

1 Multiple Linear Regression

1.1 Introduction

First it is important to differentiate between inference and prediction: in Inference we are interested in the relation between the response Y and the covariates X_1, \dots, X_n . More specifically

- Which predictors are associated with the response ?
- What is the relationship between the response and each predictor ?
- Can the relationship between response and covariates be adequately summarized using a linear equation, or is the relationship more complicated ?

Hence for interpretation purpose, there is a trade-off between the facility of interpretation and the complexity of the model e.g parametrizing the relation as linear (yielding linear regression) give a clear interpretation of the parameter i.e the relation between Y and X_1, \dots, X_n whereas parametrizing f in $Y = f(X)$ as a neural network gives a hard time to interpret the parameter "think black box model". The latter is in contrast useful for prediction. There is also a trade-off in quality of prediction with respect to the model complexity and the number of data we see. As the goal is to have a low error on prediction we need to avoid overfitting on the data we see. In the extreme case of fitting the function f perfectly to the data, we have fitted noise that is not taken into account in the relation between Y and X_1, \dots, X_n .

1.2 Statistical Learning Framework

Definition 1. A statistical learning problem is a tuple $(\mathcal{H}, \mathcal{X}, \mathcal{Y}, \mathcal{D}, l)$ where

- \mathcal{H} is the class of functions $h : \mathcal{X} \rightarrow \mathcal{Y}$ and is called the hypotheses class (and h is called a prediction rule, hypothesis or classifier)
- $\mathcal{X} := \mathcal{X} \times \mathcal{Y}$ is the domain, where
 - \mathcal{X} is the state space of observation
 - \mathcal{Y} is the label space of observations
- \mathcal{D} is a probability distribution on \mathcal{X}
- l is a measurable loss function $l : \mathcal{H} \times \mathcal{X} \rightarrow \mathbb{R}$

The learner has to find a predictor $h \in \mathcal{H}$ which minimizes the true loss (Risk)

$$R_{\mathcal{D}}(h) := \mathbb{E}_{(x,y) \sim \mathcal{D}} l(h(x), y) \quad (1)$$

Note that $\mathcal{D} = \mathcal{D}^x \times \underbrace{\mathcal{D}^y}_{\text{Unknown}}^x$ is unknown. Hence we cannot compute the true risk.

For this the learner has access to data sample $\mathcal{S} := \{(x_1, y_1), \dots, (x_n, y_n)\} \in (\mathcal{X} \times \mathcal{Y})^n$ draw *iid* with respect to \mathcal{D} . Given that sample the learner has to find an algorithm

$\mathcal{A} : \cup_{n \in \mathbb{N}} (\mathcal{X} \times \mathcal{Y})^n \rightarrow \mathcal{H}$ that returns an estimator $h_{\mathcal{S}} := \mathcal{A}(\mathcal{S})$. Note that the underlying distribution is still unknown, hence a common method is to use a proxy i.e the empirical distribution (or measure) yield the ERM methods

$$R(h_{\mathcal{S}}) = \mathbb{E}_{\mathcal{S} \sim \underbrace{\otimes_{i \in [n]} \mathcal{D}}_{\text{unknown}}} [l(h_{\mathcal{S}}, \mathcal{S})] \quad (2)$$

$$\mathbb{P}_{\mathcal{S}} \rightarrow \hat{\mathbb{P}}_{\mathcal{S}} \implies \quad (3)$$

$$\hat{R}(h_{\mathcal{S}}) := \int_{(\mathcal{X} \times \mathcal{Y})^{\otimes n}} l(h_{\mathcal{S}}, \mathcal{S}) d\hat{\mathbb{P}}_{\mathcal{S}} = \frac{1}{n} \sum_{i \in [n]} l(h(x_i), y_i) \quad (4)$$

$$\text{ERM: } h_{\mathcal{S}}^{\text{ERM}} = \arg \min_{h_{\mathcal{S}} \in \mathcal{H}} \hat{R}(h_{\mathcal{S}}) \quad (5)$$

Note: This is equivalent to the decision theoretic framework, and notation might change in this document in order to match the literature.

Notice the following two main interesting loss

- 0-1 Loss: assuming $\mathcal{Y} = \{1, \dots, k\}$ this corresponds to a classification problem (binary if $k = 2$) then an intuitive cost function is the 0-1 loss i.e $l_{0-1}(h, z) := \mathbb{I}_{\{h(x) \neq y\}}$ which yield

$$R_{\mathcal{D}}(h) = P_{(x,y) \sim \mathcal{D}}(h(x) \neq y) \quad \text{and} \quad \hat{R}_{\mathcal{S}}(h) = \frac{1}{n} \sum_{i \in [n]} \mathbb{I}_{\{h(x_i) \neq y_i\}} \quad (6)$$

- Regression and square Loss: in case $\mathcal{Y} = \mathbb{R}^q$ the task is called a regression task (see later chapter) where in this case the squared loss is used $l(h, z) := (h(x) - y)^2$ which yield

$$R_{\mathcal{D}}(h) = \mathbb{E}_{(x,y) \sim \mathcal{D}} (h(x) - y)^2 \quad \text{and} \quad \hat{R}_{\mathcal{S}}(h) = \frac{1}{n} \sum_{i \in [n]} (h(x_i) - y_i)^2 \quad (7)$$

Also note that in Guarantees for machine learning we define the Excess Risk as follow

$$\mathcal{E}(h_{\mathcal{S}}, h^*) = R_{\mathcal{D}}(h_{\mathcal{S}}) - \inf_{h \in \mathcal{H}} R_{\mathcal{D}}(h) \quad (8)$$

where h^* defines the true underlying rule $h^* : \mathcal{X} \rightarrow \mathcal{Y}$. where we developed methods to control this excess risk.

1.3 Basics of Regression Analysis and prediction task

Given the setup in the last section we aim at predicting $\underline{\mathbf{Y}}$ given $\underline{\mathbf{X}}$ where $\underline{\mathbf{Y}}$ depends non trivially to $\underline{\mathbf{X}}$. As in the last section this boils down to finding a measurable map $g : \mathcal{X} \rightarrow \mathcal{Y}$ such that $g(\underline{\mathbf{X}})$ is close to $\underline{\mathbf{Y}}$ i.e $\underline{\mathbf{Z}} = |g(\underline{\mathbf{X}}) - \underline{\mathbf{Y}}|$ is small. But as both quantity are random variable, it is not clear what is the right notion of closeness. A somewhat arbitrary idea of small is to define a random variable $\underline{\mathbf{Z}}$ small if $\mathbb{E}\underline{\mathbf{Z}}^2 = (\underbrace{\mathbb{E}\underline{\mathbf{Z}}}_{\text{mean}})^2 + \underbrace{\text{var}(\underline{\mathbf{Z}})}_{\text{fluctuation around the squared mean}}$ is small, meaning that the mean of the

random variable is small as well as the fluctuations around the mean.

Fitting this idea with last paragraph we notice that it corresponds to the notion of risk i.e for any measurable function $g : \mathcal{X} \rightarrow \mathbb{R}$ the l_2 risk is

$$R_{\mathcal{D}(x,y)}(g) = \mathbb{E}(\underline{\mathbf{Y}} - g(\underline{\mathbf{X}}))^2 \quad (9)$$

where we assumed $q = 1$ or equivalently $\mathcal{Y} = \mathbb{R}$. Now a natural question is what is the best measurable function g i.e the one minimizing the L2 risk ? the answer is

$$f_{best}(x) = \mathbb{E}_{\mathcal{Y}|x} [\mathbf{Y}|\mathbf{X} = x] \quad \forall x \in \mathcal{X} \quad (10)$$

Proof:

Assume $\underline{\mathbf{Z}}$ with $\mathbb{E}\underline{\mathbf{Z}} < \infty$ and $var(\underline{\mathbf{Z}}) < \infty$ then

$$\arg \min_{a \in \mathbb{R}} \mathbb{E}(\underline{\mathbf{Z}} - a)^2 = \mathbb{E}\underline{\mathbf{Z}}^2 - 2a\mathbb{E}\underline{\mathbf{Z}} + a^2 \quad (11)$$

$$\Leftrightarrow \quad (12)$$

$$\frac{\partial}{\partial a} (\mathbb{E}\underline{\mathbf{Z}}^2 - 2a\mathbb{E}\underline{\mathbf{Z}} + a^2)_{|a=a_{min}} = 0 \quad (13)$$

$$\Leftrightarrow \quad (14)$$

$$a_{min} = \mathbb{E}\underline{\mathbf{Z}} \quad (15)$$

Now we can rewrite the l2 Risk as follow

$$\arg \min_{g \in \mathcal{G}} R_{\mathcal{D}}(g) = \arg \min_{g \in \mathcal{G}} \mathbb{E}_{\underline{\mathbf{X}}, \underline{\mathbf{Y}}} (\underline{\mathbf{Y}} - g(\underline{\mathbf{X}}))^2 \quad (16)$$

$$= \arg \min_{g \in \mathcal{G}} \mathbb{E}_{\underline{\mathbf{X}}} \mathbb{E}_{\underline{\mathbf{Y}}|\underline{\mathbf{X}}} (\underline{\mathbf{Y}} - \underbrace{g(\underline{\mathbf{X}}|\underline{\mathbf{X}})}_{=: a \in \mathbb{R}})^2 \quad (17)$$

$$\Leftrightarrow \quad (18)$$

$$g_{min}(\underline{\mathbf{X}}) = f_{best}(\underline{\mathbf{X}}) = \mathbb{E}_{\mathcal{Y}|x} [\mathbf{Y}|\mathbf{X}] =: \text{Regression function} \quad (19)$$

$$\square. \quad (20)$$

Hence our goal is to have an estimate of the regression function. Taking any estimator or algorithm $f_{\mathcal{S}} := \mathcal{A} : \cup_{n \in \mathbb{N}} (\mathcal{X} \times \mathcal{Y})^n \rightarrow \mathcal{H}$ given $\mathcal{S} := \{(x_1, y_1), \dots, (x_n, y_n)\} \in (\mathcal{X} \times \mathcal{Y})^n$ where now we need to choose \mathcal{H} with objective to approximate $\mathbb{E}_{\mathcal{Y}|x} [\mathbf{Y}|\mathbf{X}]$ we can decompose the l2 risk as follow

$$R_{\mathcal{D}}(\hat{\mathbf{f}}_{\mathcal{S}}) = \mathbb{E}_{\mathcal{D}} [\underline{\mathbf{Y}} - f_{best}(\underline{\mathbf{X}}) + f(\underline{\mathbf{X}}) - \hat{\mathbf{f}}_{\mathcal{S}}(\underline{\mathbf{X}})]^2 \quad (21)$$

$$= \underbrace{\mathbb{E}_{\mathcal{D}} [\underline{\mathbf{Y}} - f_{best}(\underline{\mathbf{X}})]^2}_{\text{deterministic and irreducible}} + \underbrace{\|\hat{\mathbf{f}}_{\mathcal{S}} - f_{best}\|_2^2}_{\text{random}} \quad (22)$$

hence as a measure of quality we are interesting in the quantity

$$\|\hat{\mathbf{f}}_{\mathcal{S}} - f_{best}\|_2^2 = \int_{\mathcal{X}} (\hat{\mathbf{f}}_{\mathcal{S}} - f_{best})^2 dP_{\underline{\mathbf{X}}} = \|\hat{\mathbf{f}}_{\mathcal{S}} - f_{best}\|_{L_2(P_{\underline{\mathbf{X}}})}^2 \quad (23)$$

Note that one of the biggest advantage of using the square loss is that our function space is a Hilbert space ! Hence as we will show later in linear regression, we will make extensive use of orthogonality.

Note that we aim that the quantity of interest goes to Zero as $n \rightarrow \infty$

Hence defining a sequence of positive numbers $\{\phi_{n \rightarrow \infty}\}_n \rightarrow 0$ we bound the random quantity of interest as follow

- Bounds in Expectation: $\mathbb{E}_{\mathcal{S}_n} \|\hat{\mathbf{f}}_{\mathcal{S}_n} - f_{best}\|_2^2 \leq \phi_n$ **represent the average behaviour of the estimator for several sample of size n** But notice that this bound does say anything about the deviation of the random quantity hence follows the above bounds

- Bounds with high probability: $\mathbb{P}[\|\hat{\mathbf{f}}_{\mathcal{S}_n} - f_{best}\|_2^2 > \phi_n(\delta)] \leq \delta$ in this setting the bounds control the tail of the quantity if interest. Usually favored in learning theory or PAC learning
- Bounds with high probability: the above bounds usually follows from concentration around the mean of the quantity i.e $\mathbb{P}[\|\hat{\mathbf{f}}_{\mathcal{S}_n} - f_{best}\|_2 - \mathbb{E}\|\hat{\mathbf{f}}_{\mathcal{S}_n} - f_{best}\|_2 > t]$

Such techniques were studied in the courses

- Algorithmic foundations of Data Science
- Guarantees for Machine Learning

1.3.1 Random Design

The random design corresponds to the **statistical learning** setup as described in the last sections. Assume the model

$$\underline{\mathbf{y}}_i = f(\underline{\mathbf{x}}_i) + \epsilon_i \quad | \quad \epsilon_i \sim \mathcal{N}(0, \sigma^2) \quad (24)$$

then let $\mathcal{S}_{n+1} := \underbrace{\{(\underline{\mathbf{y}}_i, \underline{\mathbf{x}}_i)\}_{i \in [n]}}_{=: \underline{\mathcal{S}}_n} \cup (\underline{\mathbf{y}}_{n+1}, \underline{\mathbf{x}}_{n+1})$ of iid samples drawn from \mathcal{D} . Using $\underline{\mathcal{S}}_n$

the goal is to construct a function $\hat{\mathbf{f}}_n$ such that $\hat{\mathbf{f}}_n(\underline{\mathbf{x}}_{n+1})$ is a good predictor of $\underline{\mathbf{y}}_{n+1}$. Keep in mind that to construct our estimator only a realisation of the sample up to n is know/used i.e $\underline{\mathcal{S}}_n = \mathcal{S}_n$. A good measure of performance for a given $\underline{\mathcal{S}}_n = \mathcal{S}_n$ is L2 risk as described in the first section:

$$R(\hat{f}_n) = \mathbb{E}[\underline{\mathbf{y}}_{n+1} - \hat{f}_n(\underline{\mathbf{x}}_{n+1})]^2 = \underbrace{\mathbb{E}[\underline{\mathbf{y}}_{n+1} - f(\underline{\mathbf{x}}_{n+1})]^2}_{=\sigma^2} + \|\hat{f}_n(\underline{\mathbf{x}}_{n+1}) - f(\underline{\mathbf{x}}_{n+1})\|_{L^2(P_{x_{n+1}})}^2 \quad (25)$$

where $P_{x_{n+1}}$ is the marginal distribution of $\underline{\mathbf{x}}_{n+1}$.

Interpretation: the squared L2 norm $\|\hat{f}_n(\underline{\mathbf{x}}_{n+1}) - f(\underline{\mathbf{x}}_{n+1})\|_{L^2(P_{x_{n+1}})}^2$ measures how close \hat{f}_n is to f in average over the realizations of $\underline{\mathbf{x}}_{n+1}$ i.e how good is the prediction of $\underline{\mathbf{y}}_{n+1}$ in average over the realisations of $\underline{\mathbf{x}}_{n+1}$

1.3.2 Fix Design

In the fix design we assume that the vectors x_1, \dots, x_n are deterministic. Of course one could see x_1, \dots, x_n as realisations of random variables. But their is a fundamental difference in performance measure, here we don't have such thing as the marginal over $\underline{\mathbf{x}}_{n+1}$. Since the design points are considered deterministic our goal is to estimate f only at these point. This problem is sometime call denoising since we aim at recovering the $f(x_1), \dots, f(x_n)$ from noisy observations. We can define the fix model as follow

$$\underline{\mathbf{y}} = \mu^* + \epsilon \quad | \quad \mu_i^* = f(x_i) \quad \epsilon_i \sim \mathcal{N}(0, \sigma^2) \quad (26)$$

$$\underline{\mathbf{y}} \in \mathbb{R}^d \quad \mu^* \in \mathbb{R}^d \quad (27)$$

we define the quality of measure known as the **Mean Squared Error**

$$\text{MSE}(\hat{f}_n) = \frac{1}{n} \sum_{i \in [n]} (\hat{f}_n(x_i) - f(x_i))^2 = \frac{1}{n} \|\hat{\mu} - \mu^*\|_2^2 \quad (28)$$

1.4 Theory of Linear Regression

1.4.1 Multiple Linear Regression in Compact Form

Model given $\underline{\mathbf{X}} = \mathbb{X} = [x_1, \dots, x_n]$:

$$\mathcal{X} = \mathbb{R}^r \quad \mathcal{Y} = \mathbb{R}^q \quad \beta \in \mathbb{R}^d \quad X_i \in \{A \in \mathbb{R}^{q \times p} | \text{rank}(A) = \min(q, d)\} \quad (29)$$

$$\underline{\epsilon}_i \sim_{iid} \mathcal{N}(0, \Sigma) \quad \Psi : \mathcal{X} = \mathbb{R}^r \rightarrow \mathbb{R}^{q \times p} (\text{feature map}) \quad r, q, d \in \mathbb{N} \quad (30)$$

then we assume the linear model

$$\underline{\mathbf{y}}_i = X_i \beta + \underline{\epsilon}_i = \Psi(x_i) \beta + \underline{\epsilon}_i \quad \forall i \in [n] \quad (31)$$

with hypotheses space

$$\mathcal{H} = \{x \rightarrow \Psi(x) \beta | \psi : \mathbb{R}^r \rightarrow \mathbb{R}^{q \times d}, \beta \in \mathbb{R}^d\} \quad (32)$$

MLE: Given iid sample $\underline{\mathcal{S}}_{\mathbf{n}} = \{(\underline{\mathbf{X}}_i, \underline{\mathbf{y}}_i)\}_{i \in [n]} \sim P_{\beta}^{\otimes n} = \mathbb{P}_{\beta}$. We notice that maximizing wrt to β mle of $\underline{\mathcal{S}}_{\mathbf{n}} | \underline{\mathbf{X}} = \mathbb{X}$ and $\underline{\mathcal{S}}_{\mathbf{n}}$ yield the same programm (random and fix design). Also recall that in the statistical learning framework in section 1 we assumed that $\mathcal{D}^{y|x}$ was the unknown part of the distribution. Hence now we have assumed some model class or hypotheses testing class parametrized by β . In the above model we assume gaussian distribution which yield

$$\tilde{\mathbf{y}}_i := \underline{\mathbf{y}}_i | \underline{\mathbf{X}}_i = x_i \sim \mathcal{N}(X_i \beta, \Sigma) \quad (33)$$

hence knowing it is continous wrt to the lebesgue measure we can write the Likelihood function as follow

$$\mathcal{L}_{\underline{\mathcal{S}}_{\mathbf{n}} | \mathbb{X}}(\beta) = \prod_{i \in [n]} f_{\mathcal{N}(X_i \beta, \Sigma)}(\tilde{y}_i) \quad (34)$$

$$= \prod_{i \in [n]} (2\pi)^{-q/2} \det(\Sigma)^{-1/2} \exp\left(-\frac{1}{2}(\tilde{y}_i - X_i \beta)^T \Sigma^{-1}(\tilde{y}_i - X_i \beta)\right) \quad (35)$$

Now taking the maximum yield

$$\beta_{MLE} = \arg \max_{\beta \in \mathbb{R}^d} \log \mathcal{L}_{\underline{\mathcal{S}}_{\mathbf{n}} | \mathbb{X}}(\beta) \quad (36)$$

$$= \arg \max_{\beta \in \mathbb{R}^d} \text{const} - \frac{1}{2} \sum_{i \in [n]} \|\Sigma^{-1/2}(\tilde{y}_i - X_i \beta)\|^2 \quad (37)$$

$$= \arg \min_{\beta \in \mathbb{R}^d} \frac{1}{2} \sum_{i \in [n]} \|\Sigma^{-1/2}(\tilde{y}_i - X_i \beta)\|^2 \quad (38)$$

we have a sum of convex functions which is again convex over a convex domain, meaning a global minimum exists and is attained at $\nabla(\cdot)|_{\beta=\beta_{MLE}} = 0$ i.e

$$\nabla_{\beta=\beta_{MLE}} \left\{ \frac{1}{2} \sum_{i \in [n]} \|\Sigma^{-1/2}(\tilde{y}_i - X_i \beta)\|^2 \right\} = \sum_{i \in [n]} (-2X_i^T \Sigma^{-1} \tilde{y}_i + 2X_i^T \Sigma^{-1} X_i \beta) = 0 \quad (39)$$

$$\Leftrightarrow \quad (40)$$

$$\beta_{MLE} = \left(\sum_{i \in [n]} X_i^T \Sigma^{-1} X_i \right)^{\dagger} \left(\sum_{i \in [n]} X_i^T \Sigma^{-1} \tilde{y}_i \right) \quad (41)$$

Note: Here we assumed X_i to have full rank, hence the pseudo inverse is an inverse as Σ is psd matrix

Note that when $q = 1$ we have the more standard form

$$\underline{\mathbf{y}}_i = \langle X_i, \beta \rangle + \underline{\epsilon} = \langle \psi(x_i), \beta \rangle + \underline{\epsilon} \quad (42)$$

$$\underbrace{\underline{\mathbf{Y}}}_{\in \mathbb{R}^n} = \underbrace{X}_{\in \mathbb{R}^{n \times d}} \beta + \underline{\epsilon} \quad | \quad X = [\psi(x_1), \dots, \psi(x_n)]^T \quad (43)$$

$$(44)$$

1.4.2 Multivariate Linear Regression

Model:

1.5 Theory of Hypotheses Testing

The goal of hypotheses testing in a statistical testing problem is to decide wether a hypotheese that has been formulated is correct or not. The choice lie between two decisions, accepting or rejecting the hypotheses. A decision procedure aiming at deciding wether or not a hypotheese has to be rejected is called a test of the hypotheses.

1.5.1 Modeling the Problem

Let $\underline{\mathbf{x}} : (\Omega, \mathcal{A}, \mathbb{P}) \rightarrow (\mathcal{X}, \mathcal{B})$ be a random variable, $(\Omega, \mathcal{A}, \mathbb{P})$ a probability space and $(\mathcal{X}, \mathcal{B})$ a measurable space. The random variable $\underline{\mathbf{x}}$ induces a probability measure P_x on $(\mathcal{X}, \mathcal{B})$ given by $P_x(B) = \mathbb{P}(\underline{\mathbf{x}}(\omega) \in B) \quad \forall B \in \mathcal{B} \quad \forall \omega \in \Omega$. Assume the distribution of $\underline{\mathbf{x}}$ to be parametrized by some unknown θ and belonging to some parametrized family of distribution i.e

$$P_x = P_{\theta} \in \mathcal{P} := \{P_{\theta} : \theta \in \Theta\}$$

The decision procedure to accept or reject a hyptoheese has to be based on realisations of $\underline{\mathbf{x}}$. Defining the hypotheses $H_0 : \theta \in \Theta_0$ and the alternative $H_1 : \theta \in \Theta_1$ with $\Theta = \Theta_0 \cup \Theta_1; \Theta_0 \cap \Theta_1 = \emptyset$. Note that we could also define an hypotheses as follow

$$g : \Theta \rightarrow \Gamma \quad ; H_0 = \gamma_0 \in \Gamma; \quad \theta \in \{\vartheta \in \Theta_0 : g(\vartheta) = \gamma_0\} \quad (45)$$

for sake of simplicity we model g as the identity in this document. As the procedure has to be made on realisations of $\underline{\mathbf{x}}$ we can divide the space \mathcal{X} in the following regions

$$\mathcal{X} = \mathcal{S}_0 \cup \mathcal{S}_1 \quad (46)$$

$$\mathcal{S}_0 := \{x \in \mathcal{X} \quad \text{s.t} \quad H_0 \quad \text{is true}\} = \text{acceptance region} \quad (47)$$

$$\mathcal{S}_1 := \{x \in \mathcal{X} \quad \text{s.t} \quad H_0 \quad \text{is rejected}\} = \text{rejection region} \quad (48)$$

Definition 2 (critical function). A critical function ϕ is any function of the form $\phi(x) \in [0, 1] \quad \forall x \in \mathcal{X}$

Definition 3 (test function). A test function is a critical function $\phi(x)$ such that $\forall x \in \mathcal{X}$ we accept H_0 w.p $1 - \phi(x)$ and reject H_1 w.p $\phi(x)$

A test can made two type of errors: accepting H_0 when it is actually wrong or rejecting H_0 when it is actually true. This can be summarised as follow

Definition 4 (Error Type-I). *"Rejecting when True": for $\theta \in \Theta_0$ the function $\theta \rightarrow \mathbb{E}_\theta \phi(\underline{\mathbf{x}})$ is called type-I error*

Definition 5 (Power). *The power of a test procedure is defined as : for $\theta \in \Theta_1$ the function $\beta\theta = \theta \rightarrow \mathbb{E}_\theta \phi(\underline{\mathbf{x}})$ is called the power of the testing procedure. "Think: how powerful the test is for rejecting Null Hypotheses when actually wrong"*

Definition 6 (Error Type-II). *"Accepting when False": the function $1 - \beta(\theta)$ is called error of type-II*

Randomised Test: given $\underline{\mathbf{x}} = x$ and $\phi(x) \in [0, 1]$ a random test is defined as:

$$\underline{\delta}(x) = \begin{cases} d_0 & \text{w.p } 1 - \phi(x) \\ d_1 & \text{w.p } \phi(x) \end{cases} \equiv \underline{\delta}(x) \sim \text{Ber}(1 - \phi(x)) \quad (49)$$

$$d_0 = \text{decision to accept } H_0 \quad (50)$$

$$d_1 = \text{decision to reject } H_0 \quad (51)$$

The term random means that the experiment producing the outcomes decision d_0 and d_1 is random, explaining why the decision are accompanied by probabilities.

We notice that the test is completely characterized by the test function ϕ . Let us show the probability of error type - I for a randomized test "assume $\theta \in \Theta_0$ and reject H_0 "

$$\text{Pr}\{\underline{\delta}(\underline{\mathbf{x}}) = d_1\} = \mathbb{E}_{\underline{\delta}, \underline{\mathbf{x}}} \mathbb{I}\{\underline{\delta}(\underline{\mathbf{x}}) = d_1\} \quad (52)$$

$$= \mathbb{E}_{\underline{\mathbf{x}}} \mathbb{E}_{\underline{\delta}|\underline{\mathbf{x}}=x} \mathbb{I}\{\underline{\delta}(x) = d_1\} \quad (53)$$

$$= \mathbb{E}_{\underline{\mathbf{x}}} \text{Pr}\{\underline{\delta}(x) = d_1\} \quad (54)$$

$$= \mathbb{E}_{\underline{\mathbf{x}}} \phi(\underline{\mathbf{x}}) \quad (55)$$

$$= \mathbb{E}_\theta \phi(\underline{\mathbf{x}}) \quad (56)$$

Non-Randomised Test: given $\underline{\mathbf{x}} = x$ and $\phi(x) \in \{0, 1\}$ the same setup as before yield a non randomised procedure:

$$\delta(x) = \begin{cases} d_0 & \text{w.p } 1 - \phi(x) \\ d_1 & \text{w.p } \phi(x) \end{cases} \quad (57)$$

The test is no longer "random" as $\phi(x)$ takes 0 or 1 meaning that the underlying experiment outputs decisions to accept or reject the hypothesis (d_0, d_1) with probabilities either 1 or 0. The error of type-I is as follow

$$\text{Pr}\{\delta(\underline{\mathbf{x}}) = d_1\} = \mathbb{E}_{\underline{\mathbf{x}}} \mathbb{I}\{\delta(\underline{\mathbf{x}}) = d_1\} \quad (58)$$

$$= \mathbb{E}_\theta \phi(\underline{\mathbf{x}}) \quad (59)$$

we also notice

$$\mathbb{E}_\theta \phi(\underline{\mathbf{x}}) = \int_x \phi(x) dP_\theta(x) \quad (60)$$

$$= \int_{\mathcal{S}_0} 0 dP_\theta(x) + \int_{\mathcal{S}_1} 1 dP_\theta(x) \quad (61)$$

$$= \mathbb{I}\{x \in \mathcal{S}_1\} \quad (62)$$

were the error of type - I corresponds then to the indicator of the rejection region.

Goal: given the above setup we wish to control for error of type-I and type-II i.e find a function ϕ such that

$$\sup_{\theta \in \Theta_0} \mathbb{E}_\theta \phi(\underline{\mathbf{x}}) \leq \alpha \quad | \quad \alpha \in (0, 1) \quad (63)$$

$$\text{and} \quad (64)$$

$$\beta(\theta) \text{ is maximal } \forall \theta \in \Theta_1 \quad (65)$$

1.6 Statistical Modelling for Linear Regression

1.6.1 Gramm Schmidt and Multiple Regression from single Regression

- ESLII
- stat modelling slides and maelm summary
-

2 Non-parametric Density Estimation

- Good notes
- script compt stat 2023 (very summarized)
- Elements of computaitonal statistics

3 Non-parametric Regression

- Theory of statistics
- ESL
- Introduction to non parametric estimation

4 Classification

- Book: All of statistics (math compact)
- Book: Intro to Statistical Learning (well written applied)
- mathematics tool for ML (theory)
- Wikipedia for proof of bayes classifier very good
- (https://en.wikipedia.org/wiki/Bayes_classifier)