

Chapter 1: Introduction

- People want to make models to explain phenomena. But relevant factors may be left out, and measurements to validate come with some error. Statistical models are used here because they incorporate uncertainty.
- Linear model
 - How can an observed quality y be explained by a number of other quantities x_1, \dots, x_p ?
 - Explain the response in terms of the explanatory variables
 - Betas are called parameters and x s are called regressors, predictors...
 - Y is called the response variable

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \epsilon, \quad (1.1.1)$$

where $\beta_0, \beta_1, \dots, \beta_p$ are constants and ϵ is an error term that accounts for uncertainties. We shall refer to y as the *response variable*. It is also

- Error term accounts for unaccounted factors, missing considerations, etc.
- Why Linear models?
 - Simple and easy to interpret - often used as the first step to other models
 - Useful even for certain non-linear models
 - Probability models with multivariate normal distributions
- Description of linear models and notations
 - If you have n observations, p predictors, and x_{ij} is the i th observation of the j th predictor variable and E_i is the unobservable error

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}. \quad (1.3.2)$$

In this model,

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} 1 & x_{11} & \cdots & x_{1p} \\ 1 & x_{21} & \cdots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \cdots & x_{np} \end{pmatrix}, \quad \boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix}, \quad \boldsymbol{\epsilon} = \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{pmatrix}.$$

- It's kinda like a linear transformation $\mathbf{XB} = \mathbf{y}$ where \mathbf{X} is an $n \times p$ matrix, \mathbf{B} is a $p \times 1$ matrix, and \mathbf{y} is an $n \times 1$ matrix
- Errors are assumed to have zero mean with variances known to a scale factor

$$E(\boldsymbol{\epsilon}) = \mathbf{0}, \quad D(\boldsymbol{\epsilon}) = \sigma^2 \mathbf{V}, \quad (1.3.3)$$

where the notation E stands for expected value and D represents the dispersion (or variance-covariance) matrix. The vector $\mathbf{0}$ denotes a vector with zero elements (in this case n elements) and \mathbf{V} is a known matrix of order $n \times n$. The parameter σ^2 is unspecified, along with the vector parameter $\boldsymbol{\beta}$. The elements of $\boldsymbol{\beta}$ are real-valued, while σ^2 is nonnegative.

- $(y, XB, \sigma^2 V)$ is a shorthand triplet for a linear model (errors are uncorrelated and each has a variance σ^2)
 - where $V = I$ (the identity matrix), $(y, XB, \sigma^2 I)$ = homoscedastic linear model
 - More general structure with V is the general linear model
- If the explanatory variables are observed quantities, the x_{ij} s are assumed to be random, so the model is a conditional model of y given X

$$E(y|X) = X\beta, \quad D(y|X) = \sigma^2 V. \quad (1.3.4)$$

Thus, the error term in (1.3.2) is the difference $y - E(y|X)$. The mean and dispersion of the error given in (1.3.3) should be interpreted as conditional on X . The representation (1.3.4) is called the *linear regression model*. In this context, β is called the vector of regression parameters or regression coefficients (see Section 3.4 for a brief discussion on regression). An important aspect of the linear regression model is that the error $y - E(y|X)$ must be uncorrelated with X (see Exercise 1.6).

Suppose that the observations $(y_i, x_{i1}, \dots, x_{ip})$ for $i = 1, 2, \dots, n$ are statistically independent (in which case $V = I$). Then the conditional model (1.3.4) can be written in the simpler form

$$\begin{aligned} E(y|x_1, x_2, \dots, x_p) &= \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p, \\ \text{Var}(y|x_1, \dots, x_p) &= \sigma^2, \end{aligned} \quad (1.3.5)$$

- Scope of the linear model

- Polynomial regression

- When there is a single explanatory variable, the mean response is sought to be explained as a polynomial function of the explanatory variable

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \beta_p x^p + \epsilon,$$

where y and x are the response and the explanatory variable, respectively. This is known as the *polynomial regression model*. This model can be viewed as a special case of (1.1.1) with $x_j = x^j$, $j = 1, 2, \dots, p$. Therefore, the methodology to be developed for the model (1.1.1) will be applicable to the polynomial regression model. \square

- Linearity in the parameters makes the model linear (as long as the right-hand side is linear in the parameters even if it is not linear in the explanatory variables. (parameters are your beta coefficients)
 - The nonlinearity in the explanatory variables can be removed with transformations (linearization)

- Generalized linear model

- A nonlinear function of expected response is model as a linear function of the regression parameters

$$\eta(E(y|x_1, \dots, x_p)) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p. \quad (1.4.4)$$

This model is known as the *generalized linear model*, and the function η is called the *link function*, which is assumed to be known. One can

- You can redefine $z = \eta(y)$ where z includes the E error term

- Error due to measurement of ignored factors and linearization
- Related Models
 - Term linear model can be more general (any model which connects variables or their transformed versions through a linear relationship)
 - Autoregressive model
 - Where the explanatory variables are collections of past values of the response itself AR(p)
 - Serves as a vehicle for linear prediction (predicting future values using a linear combination of its past values)

A more general linear model for time series data is the autoregressive moving average model of order (p, q) , also known as ARMA(p, q) and given by

$$x_t = \phi_1 x_{t-1} + \phi_2 x_{t-2} + \dots + \phi_p x_{t-p} + \theta_0 \epsilon_t + \theta_1 \epsilon_{t-1} + \dots + \theta_q \epsilon_{t-q}, \quad (1.5.2)$$
 - Mixed-effects linear model
 - A combination of fixed and random parameters are involved
 - The simplest form of linear modeling is usually used as a validation benchmark
- Uses of the linear model
 - Regression analysis: average response is explained via regressors
 - Or used as a vehicle of analysis:
 - Relationship with a particular regressor (using your beta coefficient)
 - Examine statistically certain empirical beliefs: x_2 effects y 3 times more than x_1 ($B_2 \geq 3B_1$) - this can be tested using available data
 - Area of prediction: unobserved values of response may be predicted on the basis of the fitted model and the values of the explanatory variables corresponding to the unobserved response -> calibration
 - Designed experiments: effects of certain variables or testing
 - Control: if one of the explanatory variables can be controlled, which value produces the desired level of response?
 - Optimization: choosing the right explanatory variable
 - Maximum or minimum of the estimated response surface within a certain range of the variables
 - Fill missing data using related variables like a prediction or detect bad

Review of Linear Algebra

- Matrices and Vectors
 - Matrix = rows and columns of numbers
 - M rows and n columns = order $m \times n$
 - Product: number of rows in A = number of columns in B
 - Order matters
 - Diagonal matrix: all off-diagonal elements are 0
 - Sum of diagonals = trace = $\text{tr}(A)$

$$\text{tr}(\mathbf{A}_{m \times n} \mathbf{B}_{n \times m}) = \text{tr}(\mathbf{B}_{n \times m} \mathbf{A}_{m \times n}) = \sum_{i=1}^m \sum_{j=1}^n a_{i,j} b_{j,i}.$$

■ Transpose of $\mathbf{A} = \mathbf{A}' = ((a_{ij}))$

- $(\mathbf{AB})' = \mathbf{B}'\mathbf{A}'$
- Symmetric if $\mathbf{A}' = \mathbf{A}$

$\mathbf{A} + \mathbf{B}$ is defined as

$$\mathbf{A} + \mathbf{B} = ((a_{i,j} + b_{i,j})).$$

The scalar product of a matrix \mathbf{A} with a real number c is defined as

$$c\mathbf{A} = ((ca_{i,j})).$$

The difference of $\mathbf{A}_{m \times n}$ and $\mathbf{B}_{m \times n}$ is defined as $\mathbf{A} + (-1)\mathbf{B}$. Thus,

$$\mathbf{A} - \mathbf{B} = ((a_{i,j} - b_{i,j})).$$

The product of the matrix $\mathbf{A}_{m \times n}$ with the matrix $\mathbf{B}_{n \times k}$, denoted by \mathbf{AB} , is defined as

$$\mathbf{AB} = \left(\left(\sum_{l=1}^n a_{i,l} b_{l,j} \right) \right).$$

- Matrix with single column = column vector
 - Synonymous to a *vector* = column vector
 - Order of a vector = $n \times 1$ = order n for brevity
- Matrix with single row = row vector
- Non-trivial vector or matrix = not identically equal to 0
- $\mathbf{1}$ = vector of 1s
 - Identity matrix = all diagonal elements 1
 - $\mathbf{AI} = \mathbf{IA} = \mathbf{A}$
- Kronecker product

The *Kronecker product* of two matrices $\mathbf{A}_{m \times n}$ and $\mathbf{B}_{p \times q}$, denoted by

$$\mathbf{A} \otimes \mathbf{B} = ((a_{ij}\mathbf{B})),$$

is a partitioned $mp \times nq$ matrix with $a_{ij}\mathbf{B}$ as its (i,j) th block. This product is found to be very useful in the manipulation of matrices with special block structure. It follows from the definition of the Kronecker product that

- (a) $(\mathbf{A}_1 + \mathbf{A}_2) \otimes \mathbf{B} = \mathbf{A}_1 \otimes \mathbf{B} + \mathbf{A}_2 \otimes \mathbf{B}$;
- (b) $\mathbf{A} \otimes (\mathbf{B}_1 + \mathbf{B}_2) = \mathbf{A} \otimes \mathbf{B}_1 + \mathbf{A} \otimes \mathbf{B}_2$;
- (c) $(\mathbf{A}_1 \mathbf{A}_2) \otimes (\mathbf{B}_1 \mathbf{B}_2) = (\mathbf{A}_1 \otimes \mathbf{B}_1)(\mathbf{A}_2 \otimes \mathbf{B}_2)$;
- (d) $(\mathbf{A} \otimes \mathbf{B})' = \mathbf{A}' \otimes \mathbf{B}'$;
- (e) $\mathbf{A}_{m \times n} \mathbf{b} = (\mathbf{b}' \otimes \mathbf{I}_{m \times m}) \text{vec}(\mathbf{A})$, where $\text{vec}(\mathbf{A})$ is the vector obtained by successively concatenating the columns of \mathbf{A} .

- Column rank = maximum number of linearly independent columns of A
 - Row rank is the same for rows
 - Column rank of $A_{m \times n}$ is n = full rank (full row rank same)
 - Rank deficient = neither full row rank nor full column rank
 - An important result of matrix theory: row rank of any matrix = column rank
 - $p(A) \leq \min\{m, n\}$
 - A square matrix $B_{n \times n}$ is called full rank or non-singular if $p(B) = n$
 - $p(B) < n$ = singular matrix \rightarrow rank deficient
- Inner product: a generalization of the dot product
 - $a'b = b'a$
- Norm of a vector: magnitude (length of the vector from origin to point)
 - $\|a\| = (a'a)^{1/2} = \sqrt{x_1^2 + \dots + x_n^2}$
 - If $\|a\| = 1$, vector is unit norm or unit vector
- Inverse and generalized Inverses
 - If $AB = I$, then B is the right-inverse of A and A is the left inverse of B
 - Right inverse of $A = A^+R$ (only exists when A is full row rank)
 - Left inverse of $B = B^+L$ (only when B is full column rank)
 - Both inverses only exist if A is a square matrix and full rank (RREF pivots = I , so all vector lin indep). Then A^+L and A^+R are unique and equal
 - The inverse of a non-singular matrix
by A^{-1} . By definition, the inverse exists and is unique if and only if A is nonsingular, and $AA^{-1} = A^{-1}A = I$. If A and B are both nonsingular with the same order, then $(AB)^{-1} = B^{-1}A^{-1}$.
 - Non_singular matrix is an invertible matrix
 - If A is invertible and $A^{-1} = A'$, then A is an orthogonal matrix
- Vector space and projection
 - A vector u is orthogonal to another vector v if $u'v = 0$.
 - If a vector is orthogonal to all vectors in a vector space, orthogonal to S
 - Two vector spaces can be orthogonal to each other
 - Number of bases is a uniquely defined attribute of vector spaces
 - The number is called the dimension of the vector space: $\dim(S)$
 - If S consists of n -component vectors, $\dim(S) \leq n$
 - Basis
 - Orthogonal bases = vectors in basis set are orthogonal to each other
 - Unit norm = basis is called orthonormal
 - You can always construct these special bases for a vector space using Gram-Schmidt Orthogonalization
 - Projection matrices = confusing (pg. 35)
- Column Space
 - Spanned by the columns of $A = C(A)$
 - Row space = $C(A')$

Let the matrix \mathbf{A} have the columns $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n$. If \mathbf{x} is a vector having components x_1, x_2, \dots, x_n , then the matrix-vector product

$$\mathbf{A}\mathbf{x} = x_1\mathbf{a}_1 + x_2\mathbf{a}_2 + \dots + x_n\mathbf{a}_n,$$

represents a linear combination of the columns of \mathbf{A} . The set of all vectors that may be expressed as linear combinations of the columns of \mathbf{A} is a vector space and is called the *column space* of \mathbf{A} . We denote it

Proposition 2.4.1

- (a) $\mathcal{C}(\mathbf{A} : \mathbf{B}) = \mathcal{C}(\mathbf{A}) + \mathcal{C}(\mathbf{B})$.
- (b) $\mathcal{C}(\mathbf{AB}) \subseteq \mathcal{C}(\mathbf{A})$.
- (c) $\mathcal{C}(\mathbf{AA}') = \mathcal{C}(\mathbf{A})$. Consequently, $\rho(\mathbf{AA}') = \rho(\mathbf{A})$.
- (d) $\mathcal{C}(\mathbf{C}) \subseteq \mathcal{C}(\mathbf{A})$ only if \mathbf{C} is of the form \mathbf{AB} for a suitable matrix \mathbf{B} .
- (e) If $\mathcal{C}(\mathbf{B}) \subseteq \mathcal{C}(\mathbf{A})$, then $\mathbf{AA}^-\mathbf{B} = \mathbf{B}$, irrespective of the choice of the g -inverse. Similarly, $\mathcal{C}(\mathbf{B}') \subseteq \mathcal{C}(\mathbf{A}')$ implies $\mathbf{BA}^-\mathbf{A} = \mathbf{B}$.
- (f) $\mathcal{C}(\mathbf{B}') \subseteq \mathcal{C}(\mathbf{A}')$ and $\mathcal{C}(\mathbf{C}) \subseteq \mathcal{C}(\mathbf{A})$ if and only if $\mathbf{BA}^-\mathbf{C}$ is invariant under the choice of the g -inverse.
- (g) $\mathbf{B}'\mathbf{A} = \mathbf{0}$ if and only if $\mathcal{C}(\mathbf{B}) \subseteq \mathcal{C}(\mathbf{A})^\perp$.
- (h) $\dim(\mathcal{C}(\mathbf{A})) = \rho(\mathbf{A})$.
- (i) If \mathbf{A} has n rows, then $\dim(\mathcal{C}(\mathbf{A})^\perp) = n - \rho(\mathbf{A})$.
- (j) If $\mathcal{C}(\mathbf{A}) \subseteq \mathcal{C}(\mathbf{B})$ and $\rho(\mathbf{A}) = \rho(\mathbf{B})$, then $\mathcal{C}(\mathbf{A}) = \mathcal{C}(\mathbf{B})$. In particular, $\mathcal{C}(\mathbf{I}_{n \times n}) = \mathbb{R}^n$.
- (k) $\rho(\mathbf{AB}) \leq \min\{\rho(\mathbf{A}), \rho(\mathbf{B})\}$.
- (l) $\rho(\mathbf{A} + \mathbf{B}) \leq \rho(\mathbf{A}) + \rho(\mathbf{B})$.

○ ^ confusing

- Matrix Decompositions

- Rank-factorization

Any non-null matrix $\mathbf{A}_{m \times n}$ of rank r can be written as $\mathbf{B}_{m \times r}\mathbf{C}_{r \times n}$, where \mathbf{B} has full column rank and \mathbf{C} has full row rank. This is called a *rank-factorization*.

- Singular value decomposition

- Non-zero elements of \mathbf{D} are singular values

- Columns of \mathbf{U} and \mathbf{V} corresponding to the singular values are called the left and right singular vectors of \mathbf{A}

Any matrix $\mathbf{A}_{m \times n}$ can be written as \mathbf{UDV}' , where $\mathbf{U}_{m \times m}$ and $\mathbf{V}_{n \times n}$ are orthogonal matrices and $\mathbf{D}_{m \times n}$ is a diagonal matrix with nonnegative diagonal elements. This is called a *singular value decomposition*

- $p(\mathbf{A}) = p(\mathbf{D})$, so number of positive singular values = rank

Example 2.5.1 Consider the matrix

$$\mathbf{A} = \begin{pmatrix} 4 & 5 & 2 \\ 0 & 3 & 6 \\ 4 & 5 & 2 \\ 0 & 3 & 6 \end{pmatrix}.$$

The rank of \mathbf{A} is 2. An SVD of \mathbf{A} is \mathbf{UDV}' , where

$$\mathbf{U} = \begin{pmatrix} \frac{1}{2} & \frac{1}{2} & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & -\frac{1}{2} & \frac{1}{2} & -\frac{1}{2} \\ \frac{1}{2} & \frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} \\ \frac{1}{2} & -\frac{1}{2} & -\frac{1}{2} & \frac{1}{2} \end{pmatrix}, \mathbf{D} = \begin{pmatrix} 12 & 0 & 0 \\ 0 & 6 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}, \mathbf{V} = \begin{pmatrix} \frac{1}{3} & \frac{2}{3} & \frac{2}{3} \\ \frac{2}{3} & \frac{1}{3} & -\frac{2}{3} \\ \frac{2}{3} & -\frac{2}{3} & \frac{1}{3} \end{pmatrix}.$$

- Spectral decomposition, eigenvalues, and eigenvectors are confusing (pg 42-45)
- Lower order
 - Come back to this
- Solution of linear equations
 - $Ax = b$ has a solution if and only if b is in $C(A)$
 - (b) The equations $Ax = b$ have a unique solution if and only if $b \in C(A)$ and $\rho(A) = n$.
 - (c) If $b \in C(A)$, every solution to the equations $Ax = b$ is of the form $A^-b + (I - A^-A)c$ where A^- is any fixed g-inverse of A and c is an arbitrary vector.
- Optimization of quadratic forms and functions
 - Consider $q(x) = x'Ax + b'x + c$, where A is a symmetric without loss of generality
 - To minimize or maximize $q(x)$ with respect to x , differentiate $q(x)$ w.r.t components of x one at a time and set derivatives equal to 0. Check where the gradient is equal to 0.

Let $x = (x_1, \dots, x_n)'$. The gradient of a function $f(x)$ is defined as

$$\frac{\partial f(x)}{\partial x} = \begin{pmatrix} \frac{\partial f(x)}{\partial x_1} \\ \frac{\partial f(x)}{\partial x_2} \\ \vdots \\ \frac{\partial f(x)}{\partial x_n} \end{pmatrix}.$$

- Come back to this section

Chapter 3: Review of Statistical Results

- Covariance Adjustment

$E(u)$ and $D(u)$, respectively. For an ordered pair of random vectors u and v , we denote the matrix of covariances by $Cov(u, v)$, that is,

$$Cov(u, v) = E[\{u - E(u)\}\{v - E(v)\}'].$$

According to this notation, $D(u) = Cov(u, u)$, and $Cov(v, u) = Cov(u, v)'$. It is easy to see that $Cov(Au, Bv) = ACov(u, v)B'$ and $D(Au) = AD(u)A'$.

 - Variance-covariance matrix = dispersion matrix $D(u)$
 - For a single dispersion matrix, column space does not contain all the vectors having the same order of its columns
 - So it's important to know which vectors are contained
- Basic distributions
 - Multivariate normal distribution: if for every fixed vector l with the same order as y , the RV $l'y$ has the univariate normal distribution

The multivariate normal distribution is completely characterized by the mean vector and the dispersion matrix. If $E(\mathbf{y}) = \boldsymbol{\mu}$ and $D(\mathbf{y}) = \mathbf{V}$, the notation $\mathbf{y} \sim N(\boldsymbol{\mu}, \mathbf{V})$ indicates that \mathbf{y} has the multivariate normal distribution with mean $\boldsymbol{\mu}$ and dispersion matrix \mathbf{V} . The joint probability density of such a random vector with n components is

$$f(\mathbf{y}) = (2\pi)^{-n/2} |\mathbf{V}|^{-1/2} \exp\left[-\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu})' \mathbf{V}^{-1} (\mathbf{y} - \boldsymbol{\mu})\right],$$

- Come back to this section
- Distribution of quadratic forms
 - Come back to this section
- Regression
 - Approximate a random vector \mathbf{y} by a suitable function of a random vector \mathbf{x} the vector function $\mathbf{g}(\mathbf{x})$ that minimizes $E[(\mathbf{y} - \mathbf{g}(\mathbf{x}))' \mathbf{W}(\mathbf{x}) (\mathbf{y} - \mathbf{g}(\mathbf{x}))]$, where \mathbf{W} is an arbitrary, positive definite matrix that may depend on \mathbf{x} , is $\mathbf{g}(\mathbf{x}) = E(\mathbf{y}|\mathbf{x})$ (Exercise 3.7). This conditional expectation is called the *regression of \mathbf{y} on \mathbf{x}* . It is sometimes referred to as the *regression function*, when viewed as a function of \mathbf{x} . If \mathbf{y} is a scalar (denoted as y)
 - Minimum mean squared error criterion is to minimize $E[y - g(x)]^2$ with respect to function g
 - This occurs when $g(x) = E(y|x)$, the regression of y on x
 - The approximation error $y - E(y|x)$ is uncorrelated through a *linear* function of \mathbf{x} , including a constant. The solution to the minimization problem

$$\min_{\mathbf{l}, c} E[y - \mathbf{l}'\mathbf{x} - c]^2$$
 may be called the *linear regression of y on \mathbf{x}* . It is known more commonly as the *best linear predictor* (BLP) of y given \mathbf{x} . We denote it by $\hat{E}(y|\mathbf{x})$.
 - Very complicated
- Basic concepts of inference
 - Sufficiency: a statistic $t(\mathbf{y})$ is called sufficient for the parameter θ if the conditional distribution of \mathbf{y} given $t(\mathbf{y}) = t_0$
 - You are summarizing information to make the summary brief
 - $t(\mathbf{y})$ is called minimal if it is almost surely equal to a function of any other sufficient statistic almost surely for all θ - you have no more than what you need to know from \mathbf{y} about θ
 - Ancillary: $z(\mathbf{y})$ for all θ if marginal distribution of $z(\mathbf{y})$ does not depend on θ
 - Confusing
- Point estimation
 - Come back to this section
- Bayesian estimation
 - Making use of prior knowledge that is often expressed in terms of a prior distribution θ , denoted by $\pi(\theta)$ -> called the prior

- So you get the conditional distribution of y given θ
- Average risk

$$r(\mathbf{t}, \pi) = \int R(\boldsymbol{\theta}, \mathbf{t}) d\pi(\boldsymbol{\theta}),$$

- Come back to this section

Chapter 4: Estimation in the Linear Model

Consider the homoscedastic linear model $(\mathbf{y}, \mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I})$. This model is a special case of (1.3.2)–(1.3.3) where the model errors have the same variance and are uncorrelated. The unknown parameters of this model are the coefficient vector $\boldsymbol{\beta}$ and the error variance σ^2 . In this chapter we deal with the problem of estimation of these parameters from the observables \mathbf{y} and \mathbf{X} . We assume that \mathbf{y} is a vector of n elements, \mathbf{X} is an $n \times k$ matrix and $\boldsymbol{\beta}$ is a vector of k elements.

- Zero linear functions
 - An important tool in the theory of best linear unbiased estimation
 - Linear functions of the response which have zero expectation
- Linear estimation: some basic facts

Much of the classical inference problems related to the linear model $(\mathbf{y}, \mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I})$ concern a *linear parametric function* (LPF), $\mathbf{p}'\boldsymbol{\beta}$. We often estimate it by a linear function of the response, $\mathbf{l}'\mathbf{y}$. Since \mathbf{y} itself is modelled as a linear function of the parameter $\boldsymbol{\beta}$ plus error, it is reasonable to expect that one may be able to estimate $\boldsymbol{\beta}$ by some kind of a linear transformation in the reverse direction. This is why we try to estimate LPFs by linear estimators, that is, as linear functions of \mathbf{y} .

- Linear unbiased estimator and linear zero function
 - LZFs are important because they contain information about the error in the model (useful for estimating σ^2), and the mean and variance of LZFs do not depend on $\boldsymbol{\beta}$ (decoupled from $\mathbf{X}\boldsymbol{\beta}$)
 - Can be used to isolate noise
 - Ex:

$$\begin{pmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 0 \\ 0 & 1 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} \alpha_1 \\ \alpha_2 \end{pmatrix} + \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \epsilon_3 \\ \epsilon_4 \end{pmatrix}$$

The observations y_2 and y_4 , being direct measurements of error, may be used to estimate the error variance. These are LZFs. The other two observations carry information about the two parameters. There are several unbiased estimators of α_1 , such as y_1 , $y_1 + y_2$ and $y_1 + y_4$. It appears that y_1 would be a natural estimator of α_1 since it is free from the baggage of any LZF. We shall formalize this heuristic argument later. \square

- But we usually don't get this nice of a split of errors in any real linear model—our goal is to achieve it
- LZFs are sometimes referred to as linear error functions

Proposition 4.1.10 *A necessary and sufficient condition for the estimability of an LPF $(\mathbf{p}'\boldsymbol{\beta})$ is that $\mathbf{p} \in \mathcal{C}(\mathbf{X}')$.* \square

- Come back to this section
- Least squares estimation
 - Error vector can be written as $(y - XB)$ so you estimate B to minimize the sum of squared elements from this error vector

$$\hat{\beta}_{LS} = \arg \min_{\beta} (y - X\beta)'(y - X\beta). \quad \hat{\beta}_{LS} = (X'X)^{-1}X'y.$$

- $X'X$ has to be singular (X has full column rank so B is estimable)
- Best linear unbiased estimation