

Referências bibliográficas

BINGHAM, D. R.; SITTER, R. R. **Design Issues in fractional factorial split-plot experiments.** Journal of Quality Technology; Jan 2001; v.33, 1, ABI/INFORM Global pag. 39.

BISGAARD, S. **The Design and Analysis of $2^{k-p} \times 2^{q-r}$ Split Plot Experiments.** Journal of Quality Technology; Jan 2000; 32, 1, ABI/INFORM Global pag. 39.

BISGAARD, S.; FULLER, H. T.; BARRIOS, E. **Quality quandaries: two-level factorials run as split plot experiments.** Quality Engineering; 1996; v.8, p. 705-708.

BOX, G. E. P. **Must we randomize our experiment?** Quality Engineering; 1990; v.2, p. 497-502.

CALADO, V.; MONTGOMERY, D. C. **Planejamento de experimentos usando *Statística*.** E-papers Serviços Editoriais: Rio de Janeiro, 2003.

GALDÁMEZ, Edwin Vladimir Cardoza. **Aplicação das técnicas de planejamento e análise de experimentos na melhoria da qualidade de um processo de fabricação de produtos plásticos,** 2002. Disponível em: <<http://www.teses.usp.br/teses>>. Acesso em 28 out 2007.

GANJU, J.; LUCAS, J. M. **Detecting randomization restrictions caused by factors.** Journal of Statistical Planning and Inference; 1999; v.81, p. 129-140.

GANJU, J.; LUCAS, J. M. **Randomized and random run order experiments** Journal of Statistical Planning and Inference; 2004; v.133, pag. 199-210.

GIL, Antonio C. **Como elaborar projetos de pesquisa.** Atlas: São Paulo, 1991.

GOMES, U. R.. **Otimização do processo de laminação a frio através de planejamentos de experimentos.** Dissertação (Mestrado em Engenharia Industrial) – Pontifícia Universidade Católica do Rio de Janeiro, Rio de Janeiro, 2007.

JU, H. L.; LUCAS, J. M. **L^k Factorial experiments with hard-to-change and easy-to-change factors.** Journal of Quality Technology; Oct 2002; vol. 34, no 4, ABI/INFORM Global pag. 411.

LETSINGER, Jennifer D.; MYERS, Raymond H.; LENTNER, Marvin. **Response surface methods for bi-randomization structures.** Journal of Quality Technology; Oct 1996; vol. 28, no 4, ABI/INFORM Global pag. 381.

LOEPPKY, J.L.; SITTER, R. R. **Analyzing unreplicated blocked or split-plot fractional factorial designs.** Journal of Quality Technology; Jul 2002; v.34, pag. 229.

MARTINS, G. A. **Manual para elaboração de monografias e dissertações.** Atlas: São Paulo, 1994.

MONTGOMERY, Douglas C. **Design and Analysis of Experiments.** Wiley: New York, 2001.

MONTGOMERY, Douglas C.; RUNGER, George C. **Estatística Aplicada e Probabilidade para Engenheiros.** LTC: Rio de Janeiro, 2003.

MONTGOMERY, Douglas C.; MYERS, Raymond H. **Response Surface Methodology: Process and Product Optimization Using Designed Experiments.** Wiley-Interscience: 2002.

SKF/CSN. **Relatório de estudo sobre chatter nos laminadores LTF 1 e LTF 2.** SKF do Brasil Ltda: Rio de Janeiro, 2007.

VERGARA, Sylvia C. **Projetos e Relatórios de Pesquisa em Administração.** São Paulo: Atlas, 1997.

VIEIRA, Antonio F. C. **Análise da média e dispersão em experimentos fatoriais não replicados para otimização de processos industriais.** Tese (Doutorado em Engenharia Industrial) – Pontifícia Universidade Católica do Rio de Janeiro, Rio de Janeiro, 2004.

VINING, G. Geoffrey; KOWALSKI, Scott M.; MONTGOMERY, Douglas C. **Response Surface Designs within a Spli-Plot structure.** Journal of Quality Technology; Apr 2005; v.37, p. 115.

WEBB, Derek F.; LUCAS, James M.; BORKOWSKI, John J. **Factorial experiments when factor levels are not necessarily reset.** Journal of Quality Technology; Jan 2004; vol. 36, pag. 1.

Apêndices

APÊNDICE A – Etapas de execução de um projeto *split-plot* para o exemplo de força de tensão do papel.

Elaborando o projeto

Utilizando um projeto fatorial geral, por se tratar de fatores com mais de dois níveis, conforme a Figura A.1, introduz-se o número de fatores a serem analisados no experimento. Sabe-se que o experimento apresenta dois fatores, entretanto, por sugestão do tutorial e para auxiliar na análise considera-se a replicação como um fator.

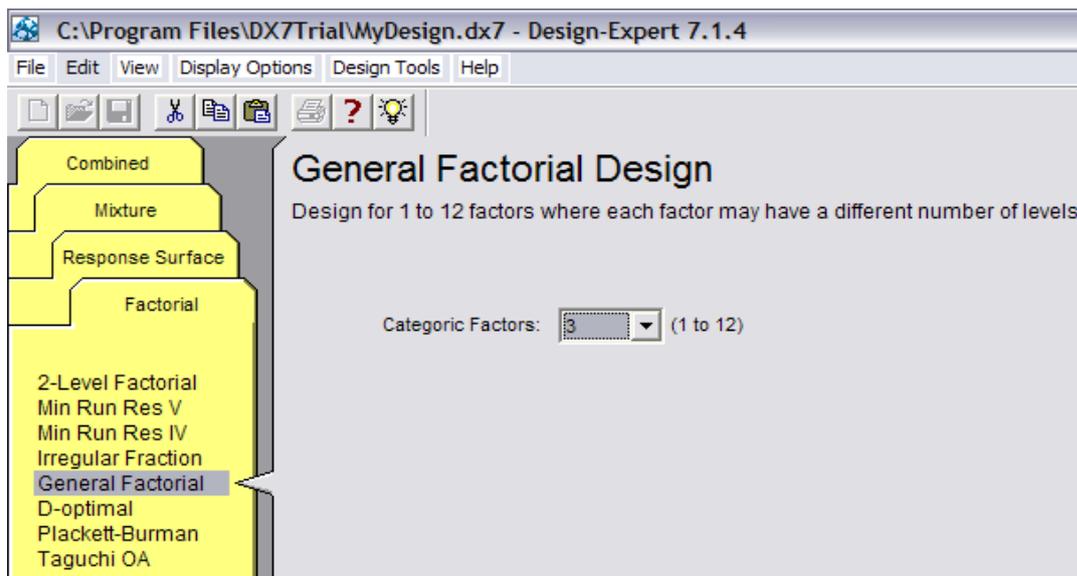


Figura A.1 – Seleção do tipo de projeto para a execução do experimento *split-plot*.

Dando continuidade ao projeto, na próxima tela introduz-se **réplica** para o nome do fator, com **3** níveis, denominados **dia 1**, **dia 2** e **dia 3**, conforme apresentado na Figura A.2.

General Factorial Design

Factor A: Name: Current number of Rows: 12

Units: Maximum number of Rows: 32766

Levels: (2 to 999) Categoric contrasts:

Nominal Ordinal

Treatments
Dia 1
Dia 2
Dia 3

Figura A.2 – Inserção do nome e das características do fator A (Bloco).

Clicando novamente no botão *Continue*, no canto direito inferior da tela , a próxima visualização refere-se ao fator *B*. Então se insere o nome do fator **Preparação da polpa (Prep polpa)**, com os **3** níveis **B1**, **B2** e **B3**, como na Figura A.3.

General Factorial Design

Factor B: Name: Current number of Rows: 18

Units: Maximum number of Rows: 32766

Levels: (2 to 999) Categoric contrasts:

Nominal Ordinal

Treatments
B1
B2
B3

Figura A.3 – Inserção do nome e das características do fator B (whole plot).

O próximo passo é inserir o terceiro fator no projeto e nomeá-lo como **Temperatura (Temp)**, com **graus Fahrenheit (graus F)** como unidades, **4** níveis, especificados como **200**, **225**, **250** e **275** para os tratamentos (Figura A.4). Para reconhecer a natureza numérica desse fator, deve-se mudar a categoria dos contrastes para **ordinal**.

General Factorial Design

Factor C: Name: Current number of Rows: 36

Units: Maximum number of Rows: 32766

Levels: (2 to 999) Categorical contrasts:
 Nominal Ordinal

Treatments
200
225
250
275

Figura A.4 – Inserção do nome e das características do fator C (Subplot).

A tela seguinte solicita o número de replicações para o experimento. Como, para auxiliar na análise, réplica foi considerada como um fator, sugere-se utilizar o valor “1” para produzir 36 corridas (3x3x4 níveis dos três fatores) e prosseguir. Nota-se, no próximo passo, a solicitação do número de respostas para o experimento. Apesar de se medir apenas uma resposta no experimento, para efeito de análise, esta resposta será particionada em 4 opções com as respectivas unidades referentes ao problema, conforme a Figura A.5. As demais informações solicitadas na tela não são requisitos para a execução do experimento.

General Factorial Design

Optional Power Wizard: For each response, you may enter the minimum change the design should detect as statistically significant and also the estimated standard deviation of each response (generally obtained from historical data). The ratio will then be calculated in the Delta/Sigma field. Press Continue to see the calculated power for each response. A probability of 80% or higher is recommended. If power is low, consider adding runs by choosing a larger design or replication, or reconcile yourself to not detecting a signal this small.

Leave Sigma and Delta fields blank to skip power calculation.

Responses: (1 to 999)

Name	Units	Diff. to detect Delta("Signal")	Est. Std. Dev. Sigma("Noise")	Delta/Sigma (Signal/Noise Ratio)
Whole plot	tensão			
Sub plot	tensão			
Interação	tensão			
Todos os efeitos	tensão			

Figura A.5 – Definição das quatro respostas a serem analisadas.

A tela do projeto se apresenta agora em ordem completamente aleatória. Como não é o caso, deve-se configurar a ordem de execução das corridas conforme a alocação das observações em relação aos fatores. As três figuras a seguir ilustram como deve ser o arranjo no *software* para que o projeto reflita a realidade. Primeiramente, ordena-se as corridas com base no fator “preparação da polpa” (Figura A.6). Em seguida, faz-se o mesmo procedimento, só que desta vez ordenando as corridas com base no fator “réplica” (Figura A.7). Por fim, na coluna “run”, atribui-se a reordenação das corridas com base na atual configuração (Figura A.8).

Select	Std	Run	Factor 1 A: Réplica	Factor 2 B: P...	Factor 3	Response
	2	2	Dia 2			
	20	3	Dia 2			
	28	5	Dia 1			
	29	11	Dia 2	B1		275
	3	13	Dia 3	B1		200
	1	15	Dia 1	B1		200

Figura A.6 – Re-organização dos dados com base no fator B.

Select	Std	Run	Factor 1 A: Réplica	Factor 2 B: P...	Factor 3 C: Temp (aus F)	Response
	28	5				275
	1	15				200
	10	26				225
	19	36	Dia 1	B1		250
	4	20	Dia 1	B2		200
	22	22	Dia 1	B2		250

Figura A.7 – Re-organização dos dados com base nos blocos.

Select	Std	R	Factor 1	Factor 2	Factor 3
	28				Temp aus F
	1				275
	10	3	Dia 1	B1	225
	19	4	Dia 1	B1	250
	4	5	Dia 1	B2	200
	22	6	Dia 1	B2	250
	13	7	Dia 1	B2	225
	31	8	Dia 1	B2	275

Figura A.8 – Re-organização da ordem das corridas experimentais.

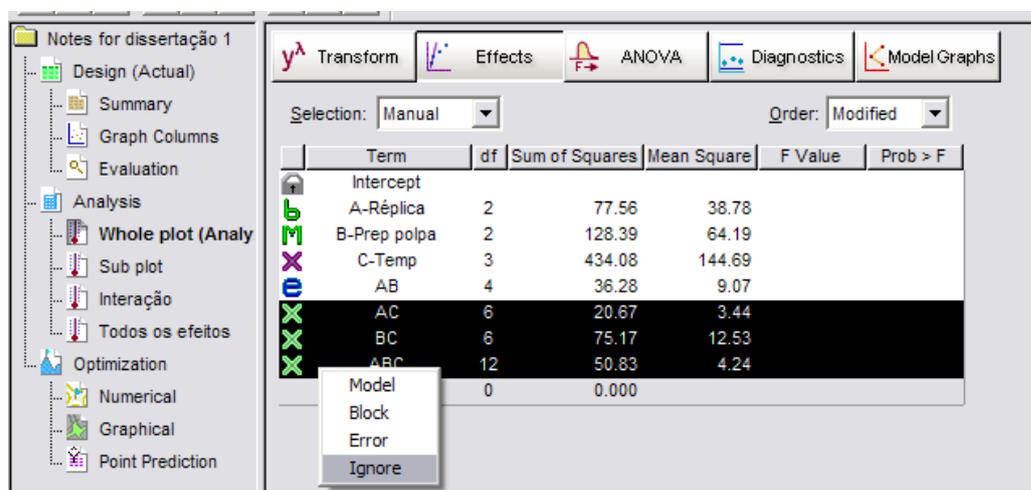
Analizando os dados

Para que seja possível analisar um projeto *split-plot* no software *Design Expert 7.1.4*, é necessário criar ANOVA's separadas, e isto é feito manualmente. Cria-se uma ANOVA para o tratamento *whole plot* (preparação da polpa), outra para o tratamento *subplot* (temperatura) e, finalmente a interação *whole plot* e *subplot*, para que individualmente cada uma seja testada corretamente. Ainda ajusta-se um modelo completo para ter o diagnóstico de significância e os gráficos do modelo, mas ignora-se a ANOVA do modelo completo.

Select	Std	Run	Factor 1 A: Réplica	Factor 2 B: Prep polpa	Factor 3 C: Temp graus F	Response 1 Whole plot tensão	Response 2 Sub plot tensão	Response 3 Interação tensão	Response 4 Todos os efeit tensão
	28	1	Dia 1	B1	275	36	36	36	36
	1	2	Dia 1	B1	200	30	30	30	30
	10	3	Dia 1	B1	225	35	35	35	35
	19	4	Dia 1	B1	250	37	37	37	37
	4	5	Dia 1	B2	200	34	34	34	34
	22	6	Dia 1	B2	250	38	38	38	38
	13	7	Dia 1	B2	225	41	41	41	41
	31	8	Dia 1	B2	275	42	42	42	42
	34	9	Dia 1	B3	275	36	36	36	36
	16	10	Dia 1	B3	225	26	26	26	26
	7	11	Dia 1	B3	200	29	29	29	29
	25	12	Dia 1	B3	250	33	33	33	33
	2	13	Dia 2	B1	200	28	28	28	28
	20	14	Dia 2	B1	250	40	40	40	40
	29	15	Dia 2	B1	275	41	41	41	41
	11	16	Dia 2	B1	225	32	32	32	32
	32	17	Dia 2	B2	275	40	40	40	40
	23	18	Dia 2	B2	250	42	42	42	42
	14	19	Dia 2	B2	225	36	36	36	36
	5	20	Dia 2	B2	200	31	31	31	31
	35	21	Dia 2	B3	275	40	40	40	40
	17	22	Dia 2	B3	225	30	30	30	30
	8	23	Dia 2	B3	200	31	31	31	31
	26	24	Dia 2	B3	250	32	32	32	32
	3	25	Dia 3	B1	200	31	31	31	31
	30	26	Dia 3	B1	275	40	40	40	40
	21	27	Dia 3	B1	250	41	41	41	41
	12	28	Dia 3	B1	225	37	37	37	37

Figura A.9 – Arranjo do projeto *split-plot*.

Para se analisar inicialmente o efeito *whole plot* (*B*), deve-se clicar no rótulo referente ao nome no diretório *Analysis* na janela à esquerda. Logo em seguida, deve-se clicar no botão *Effects*, com o qual aparecerá a lista de efeitos do experimento. O *software* oferece quatro alternativas para designar os efeitos, dentre as quais: *Model* (“*M*”); *Block* (“*b*”); *Error* (“*e*”); *Ignore* (“*X*”). O tratamento *whole plot* deverá ser testado em relação a **Réplica** através da interação *AB*. Dessa forma, será selecionado para o modelo o fator *B*; como bloco, a Réplica; como erro, a interação *AB*; e serão ignorados do modelo os demais efeitos (Figura A.10).



Term	df	Sum of Squares	Mean Square	F Value	Prob > F
Intercept					
A-Réplica	2	77.56	38.78		
B-Prep polpa	2	128.39	64.19		
C-Temp	3	434.08	144.69		
AB	4	36.28	9.07		
AC	6	20.67	3.44		
BC	6	75.17	12.53		
ABC	12	50.83	4.24		
Model	0	0.000			

Figura A.10 – Lista de efeitos para o projeto *whole-plot*.

A seguir, seleciona-se o botão da ANOVA para analisar a significância dos tratamentos. Com base na Figura A.11, nota-se que o efeito *whole plot* (preparação da polpa) é significativo, uma vez que o *p-value* < 0,05. Aconselha-se, neste ponto, a não analisar o diagnóstico (*diagnostics*) e os gráficos do modelo, visto que os modelos ainda estão incompletos nesta etapa.

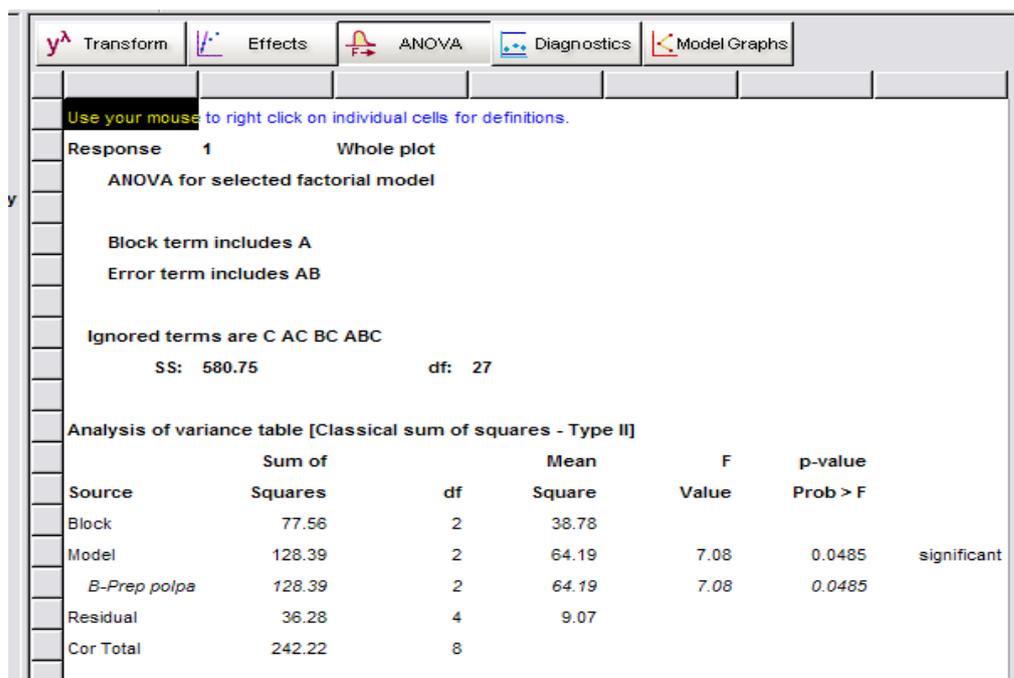


Figura A.11 – ANOVA para o projeto *whole plot*.

Para analisar o tratamento *subplot* (temperatura), seleciona-se a guia *subplot* no diretório de *Analysis* e depois, clica-se no botão *Effects*. O tratamento *subplot* deve ser testado em relação à interação Réplica x Temperatura, isto é, a interação *AC*. Desse modo, a lista de efeitos deve estar configurada conforme a Figura A.12 (os fatores *A* e *B* definidos como bloco; o fator *C* selecionado para o modelo; a interação *AC* listada como erro; e as interações *AB*, *BC* e *ABC* ignoradas). Prosseguindo com a visualização da ANOVA (Figura A.13), é possível notar que o fator *subplot* (Temperatura) é significativo ($p\text{-value} < 0,05$). Novamente não se deve olhar o diagnóstico (*diagnostics*) e os gráficos do modelo, visto que, nesta etapa, os modelos ainda estão incompletos.

Term	df	Sum of Squares	Mean Square	F Value	Prob > F
Intercept					
A-Réplica	2	77.56	38.78		
B-Prep polpa	2	128.39	64.19		
C-Temp	3	434.08	144.69		
AB	4	36.28	9.07		
AC	6	20.67	3.44		
BC	6	75.17	12.53		
ABC	12	50.83	4.24		
Residuals	0	0.000			

Figura A.12 – Lista de efeitos para o projeto *subplot*.

Use your mouse to right click on individual cells for definitions.

Response 2 Sub plot

ANOVA for selected factorial model

Block term includes A, B
Error term includes AC

Ignored terms are AB BC ABC

SS: 162.278 df: 22

Analysis of variance table [Classical sum of squares - Type II]

Source	Sum of Squares	df	Mean Square	F Value	p-value	Prob > F
Block	205.94	4	51.49			
Model	434.08	3	144.69	42.01	0.0002	significant
C-Temp	434.08	3	144.69	42.01	0.0002	
Residual	20.67	6	3.44			
Cor Total	660.69	13				

Figura A.13 – ANOVA para o projeto *subplot*.

No caso da interação *whole plot* x *subplot*, clica-se na opção interação no diretório de *Analysis*. Semelhante aos demais, no botão *Effects*, a interação *BC* deve ser testada em relação à interação *ABC*, logo esta é selecionada como erro. A réplica e os fatores *B* e *C* são considerados blocos. A interação *BC* é selecionada para o modelo, e os demais são ignorados, como mostrado na Figura A.14. Pela ANOVA (Figura A.15), observa-se que a interação *BC* é considerada não significativa pelo programa, visto que $p\text{-value} > 0,05$ em um valor limítrofe.

Dessa forma, pela proximidade do valor cabe ao experimentalista a decisão de ignorar a recomendação do software e considerar *BC* significativo.

Term	df	Sum of Squares	Mean Square	F Value	Prob > F
Intercept					
A-Réplica	2	77.56	38.78		
B-Prep polpa	2	128.39	64.19		
C-Temp	3	434.08	144.69		
AB	4	36.28	9.07		
AC	6	20.67	3.44		
BC	6	75.17	12.53		
ABC	12	50.83	4.24		
Residuals	0	0.000			

Figura A.14 – Lista de efeitos para o projeto de interação *whole plot x subplot*.

Use your mouse to right click on individual cells for definitions.

Response 3 Interação

ANOVA for selected factorial model

Block term includes A, B, C
Error term includes ABC

Ignored terms are AB AC

SS: 56.9444 df: 10

Analysis of variance table [Classical sum of squares - Type II]

Source	Sum of Squares	df	Mean Square	F Value	p-value	Prob > F
Block	640.03	7	91.43			
Model	75.17	6	12.53	2.96	0.0520	not significant
BC	75.17	6	12.53	2.96	0.0520	
Residual	50.83	12	4.24			
Cor Total	766.03	25				

Figura A.15 – ANOVA para a interação *whole-plot x sub-plot*.

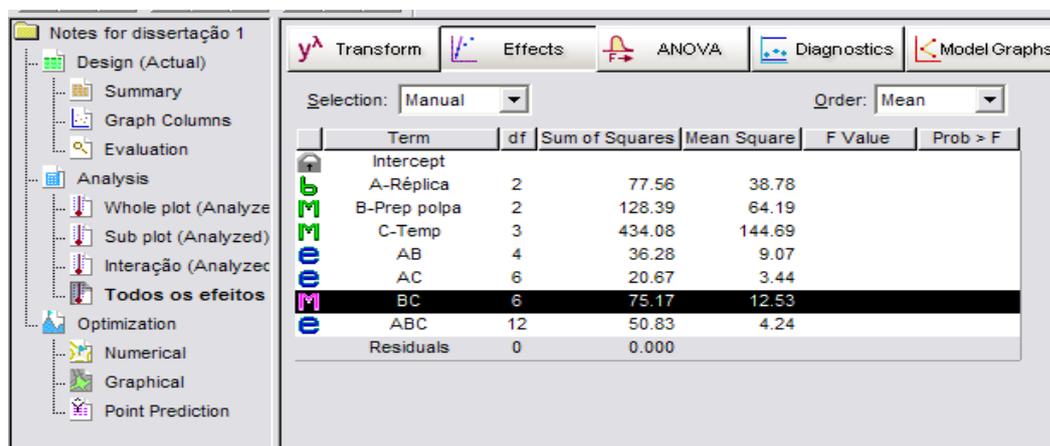
Antes que se avance para a fase final do projeto, devem ser levados em consideração alguns aspectos importantes, que justificam a análise estatística feita até o momento:

- Em um experimento bloqueado as interações bloco *versus* tratamento são usadas para estimar o erro.

- Em um projeto bloqueado completamente aleatorizado, todas as interações bloco *versus* tratamento são combinadas em uma única estimativa do erro.
- Em um projeto *split-plot*, devido às restrições quanto à aleatorização, determinadas interações bloco *versus* tratamento devem ser designadas para estimar o erro de um tratamento particular que esteja sendo testado.

Com base nestas justificativas, as análises de variância executadas até o momento neste projeto experimental foram: o tratamento *whole plot* (preparação da polpa) em relação à interação *AB* (bloco *versus* preparação da polpa); o tratamento *subplot* (temperatura) em relação à interação *AC* (bloco *versus* temperatura); e, por fim, a interação *BC* (polpa *versus* temperatura) em relação à interação *ABC* (bloco *versus* polpa *versus* temperatura).

Para que seja possível obter diagnósticos significativos e a análise gráfica do modelo, deve-se ajustar o modelo completo (Figura A.16), selecionando-se os fatores *B*, *C* e a interação *BC* para o modelo; *A* como bloco; e os demais fatores como erro. Por conseguinte, faz-se a análise do diagnóstico e dos gráficos do modelo, desconsiderando-se a ANOVA.



Term	df	Sum of Squares	Mean Square	F Value	Prob > F
Intercept					
A-Réplica	2	77.56	38.78		
B-Prep polpa	2	128.39	64.19		
C-Temp	3	434.08	144.69		
AB	4	36.28	9.07		
AC	6	20.67	3.44		
BC	6	75.17	12.53		
ABC	12	50.83	4.24		
Residuals	0	0.000			

Figura A.16 – Lista de efeitos para o projeto completo.

Diagnóstico e influência dos resíduos

- Gráfico da Normal x Resíduos: este gráfico tem a finalidade de identificar alguma alteração no que tange a normalidade e a presença de observações

atípicas. Pode-se analisar, para o experimento em questão, segundo a Figura A.17, que não há indícios de observações atípicas e de perda de normalidade uma vez que os valores dos resíduos estão uniformemente dispostos sobre a reta.

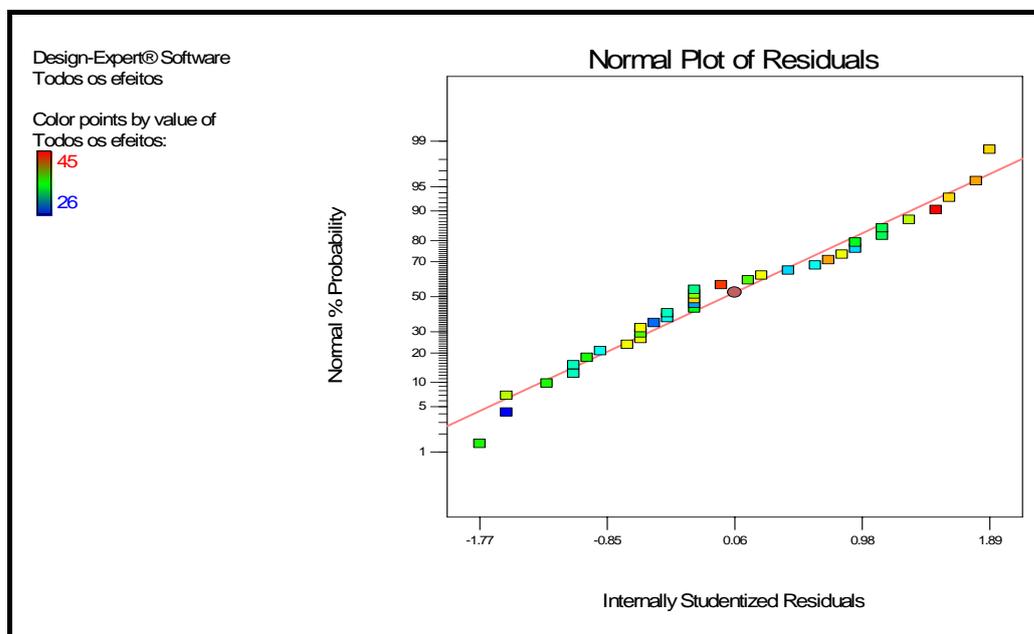


Figura A.17 – Gráfico de probabilidade normal para o experimento da polpa de papel.

- Gráfico de Resíduos x Diagnóstico (Valores previstos): As características importantes que este gráfico permite avaliar são: a adequação do modelo quanto à variância constante e a aditividade. Dessa forma, analisando o gráfico da Figura A.18 obtido para o experimento em questão, é possível identificar que não há indícios de padrão de comportamento estabelecido pelas observações, o que permite concluir que a variância é constante. No que tange a aditividade, pode-se observar que os valores estão distribuídos de modo uniforme em torno de zero, mas não se deve assegurar a aditividade dos efeitos.

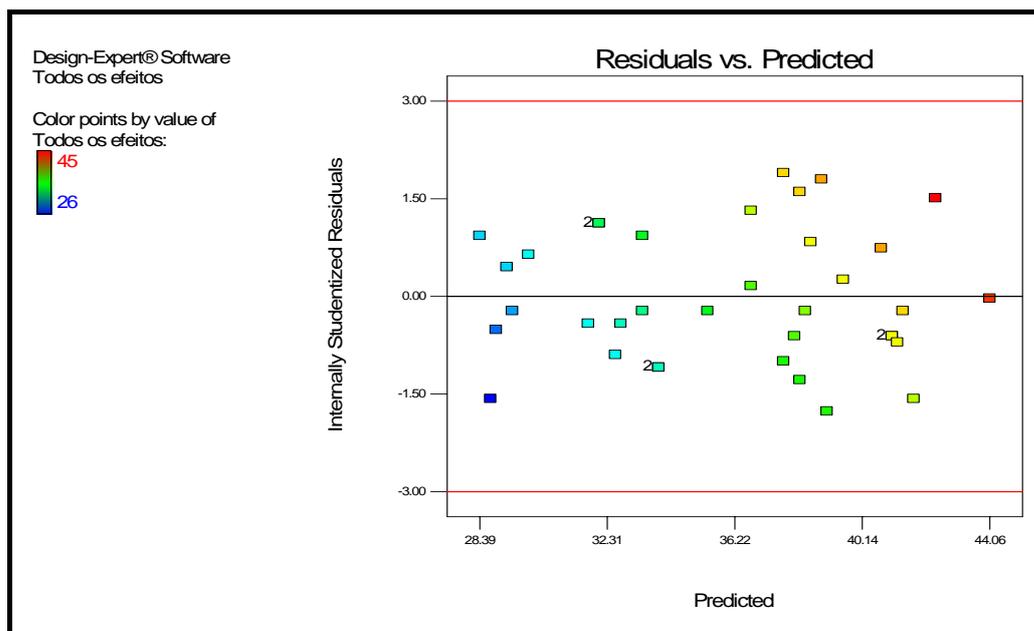


Figura A.18 – Gráfico dos resíduos *versus* os valores previstos.

- Gráfico da Distância de Cook (Cook's Distance - D_i): Para analisar a influência dos resíduos é possível destacar na observação do gráfico abaixo que não há valores de resíduos maiores ou iguais a 0,5 ($D_i \geq 0,5$), logo se confirma a ausência de observações influentes (*outlier*). Ressalta-se que quanto maior este resíduo (Valores de $D_i \geq 0,5$) mais influente a observação será para o modelo.

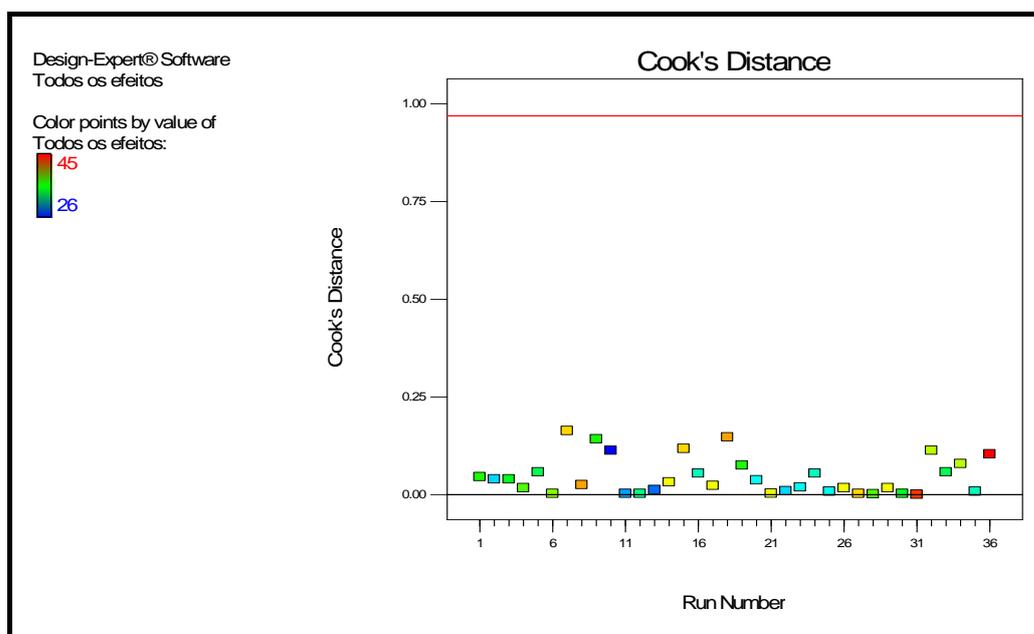


Figura A.19 – Gráfico da distância de Cook para o experimento *split-plot*.

APÊNDICE B – Lista dos efeitos *whole-plot* para análise do gráfico *half-normal*.

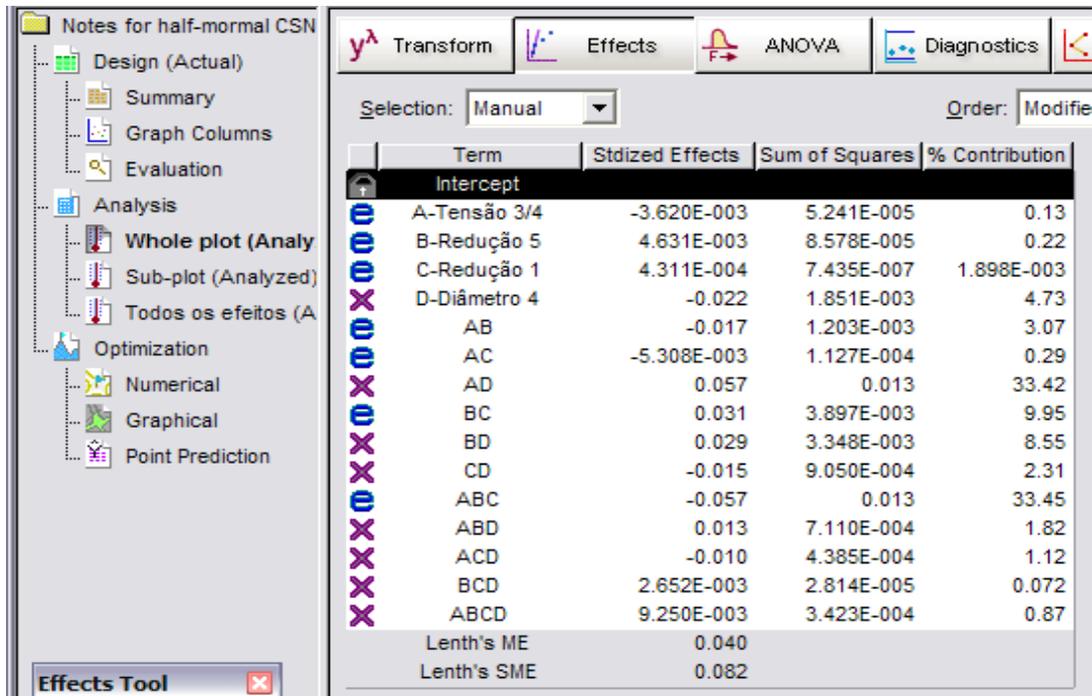


Figura B.1 – Lista dos efeitos *whole-plot* para análise do gráfico *half-normal*.

APÊNDICE C – Lista dos efeitos *whole-plot* para análise do gráfico *half-normal* após seleção dos efeitos significativos.

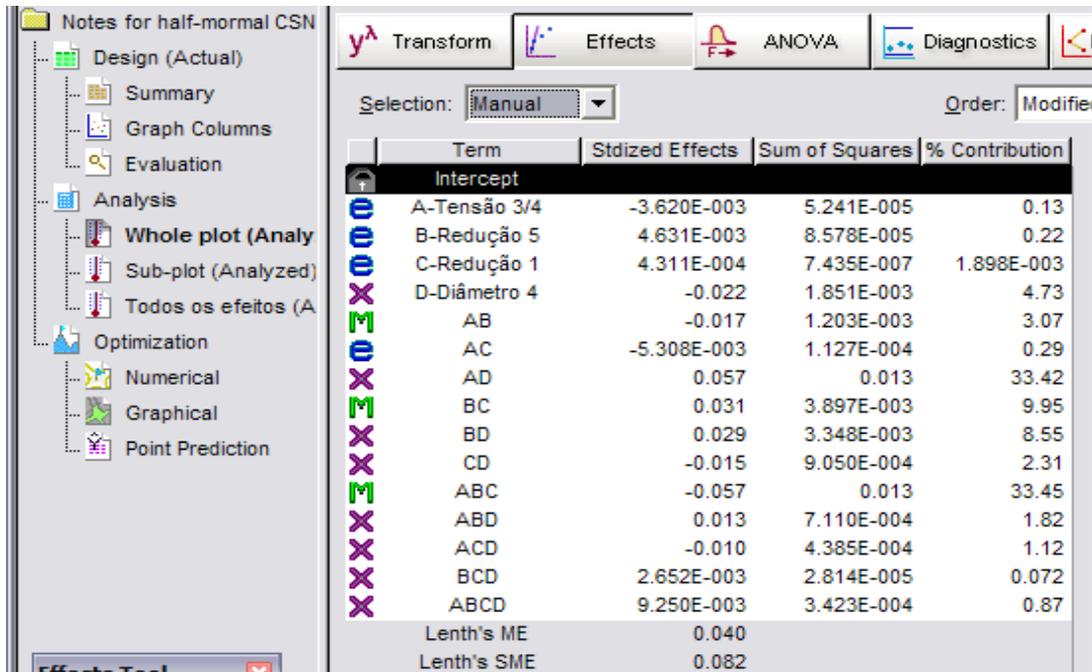


Figura C.1 – Lista dos efeitos *whole-plot* para análise do gráfico *half-normal* após seleção dos efeitos significativos

APÊNDICE D - Lista dos efeitos *sub-plot* para análise do gráfico *half-normal*

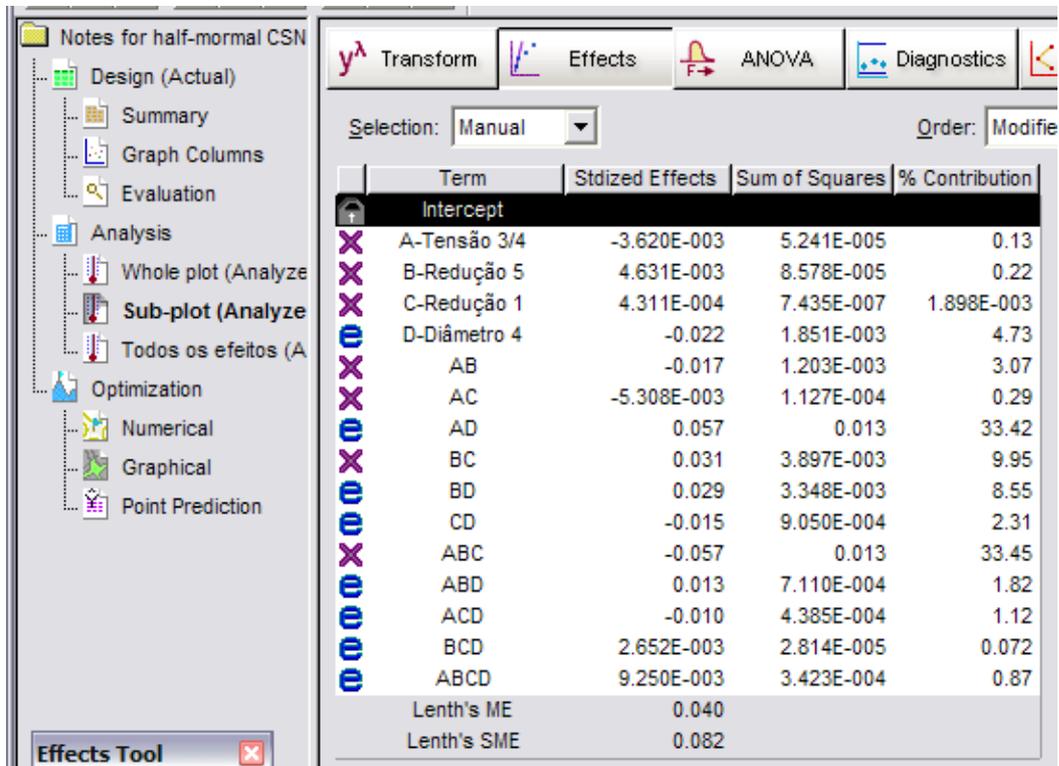


Figura D.1 – Lista dos efeitos *sub-plot* para análise do gráfico *half-normal*.

APÊNDICE E - Lista dos efeitos *sub-plot* para análise do gráfico *half-normal* após seleção dos efeitos significativos.

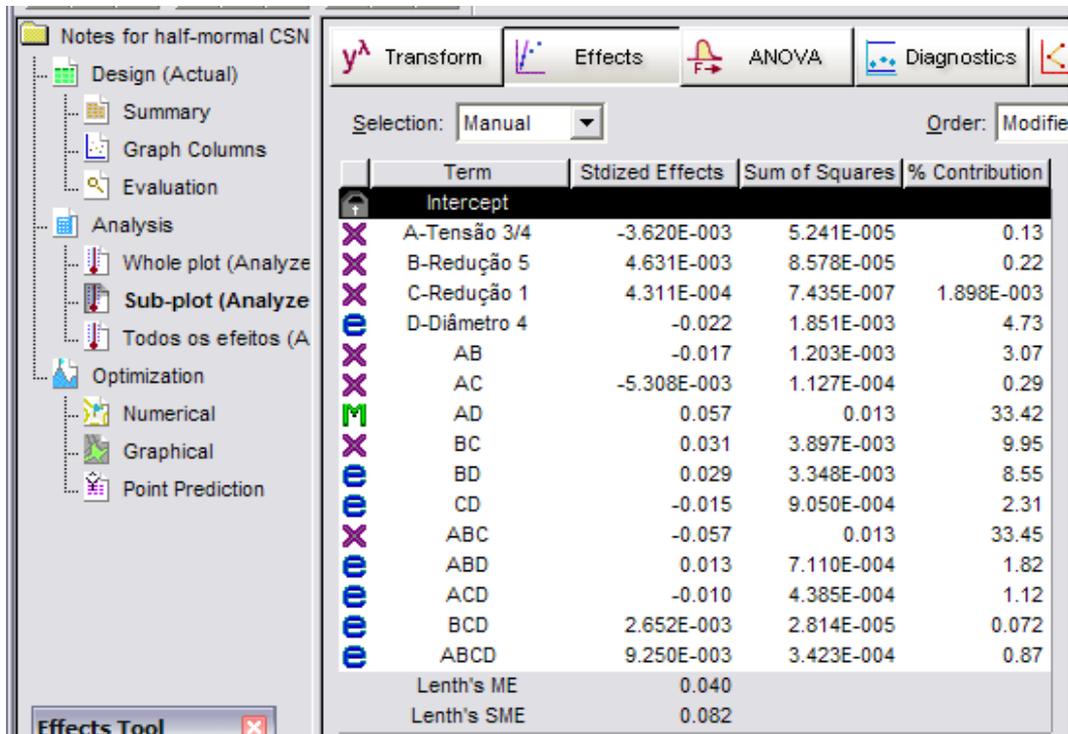


Figura E.1 – Lista dos efeitos *sub-plot* para análise do gráfico *half-normal* após seleção dos efeitos significativos.

Anexos

ANEXO A – Revisão de regressão linear e procedimento para cálculo de intervalo de confiança para a média da resposta e de intervalo de previsão para uma futura resposta (Vieira, 2004).

A.1 Regressão Linear

A análise de regressão linear procura estabelecer a relação entre uma variável de resposta y e um conjunto de variáveis independentes x_1, x_2, \dots, x_k . Entende-se por variável de resposta uma medida de desempenho ou uma característica da qualidade de um produto ou um processo produtivo, os quais são influenciados pelas variáveis independentes, que também podem ser chamadas de variáveis de regressão. O termo “linear” é utilizado pelo fato de a equação da variável de resposta y representar uma função linear dos parâmetros desconhecidos $\beta_0, \beta_1, \beta_2, \dots, \beta_k$, como será visto a seguir.

A.1.1 Função de Resposta

Quando um produto ou um processo possui uma resposta y que depende das variáveis de entrada controláveis x_1, x_2, \dots, x_k , diz-se que há uma relação funcional entre y e x_1, x_2, \dots, x_k , que pode ser representada da seguinte forma:

$$y = f(\beta_1, \beta_2, \dots, \beta_k, x_1, x_2, \dots, x_k) + \varepsilon \quad (\text{A.1})$$

sendo o termo ε o erro, que representa outras fontes de variabilidade não consideradas em f , como erros de medição da resposta ou outras variações inerentes ao processo ou sistema.

Geralmente, a relação funcional apresentada na Equação (A.1) não é conhecida, o que propicia a utilização de modelos lineares de regressão, que podem ter a seguinte representação:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + \varepsilon \quad (\text{A.2})$$

onde os parâmetros desconhecidos, $\beta_0, \beta_1, \beta_2, \dots, \beta_k$, são chamados de coeficientes de regressão.

Os modelos aparentemente mais complexos também podem ser representados pelo modelo apresentado na Equação (A.2). Como exemplo,

considere a adição de um termo de interação a um modelo de primeira ordem com duas variáveis, ou seja,

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_{12} x_1 x_2 + \varepsilon \quad (\text{A.3})$$

Se for feita uma substituição do tipo $x_3 = x_1 x_2$ e $\beta_3 = \beta_{12}$, o modelo apresentado na equação (A.3) pode ser escrito da seguinte forma:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \varepsilon \quad (\text{A.4})$$

que tem a mesma forma da Equação (A.2).

Da mesma forma, considerando um modelo de segunda ordem com duas variáveis, tem-se

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_{12} x_1 x_2 + \beta_{11} x_1^2 + \beta_{22} x_2^2 + \varepsilon \quad (\text{A.5})$$

Fazendo uma substituição do tipo $x_3 = x_1^2$, $x_4 = x_2^2$, $x_5 = x_1 x_2$, $\beta_3 = \beta_{11}$, $\beta_4 = \beta_{22}$ e $\beta_5 = \beta_{12}$, o modelo da Equação (I.5) fica da seguinte forma:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 + \beta_5 x_5 + \varepsilon \quad (\text{A.6})$$

que também é um modelo de regressão linear.

A seguir, serão apresentados métodos para estimar os parâmetros dos modelos de regressão linear e para testar a significância dos coeficientes, o que é freqüentemente chamado de ajuste do modelo.

A.1.2 Estimação dos Parâmetros nos Modelos de Regressão Linear

O método dos mínimos quadrados é o método clássico de estimação dos parâmetros dos modelos de regressão linear.

Considere $n > k$ observações da variável de resposta, ou seja, y_1, y_2, \dots, y_n , sendo que, para cada resposta, tem-se as observações das variáveis de regressão, como mostrado na Tabela A.1.

Escrevendo a Equação (A.2) utilizando os dados apresentados na Tabela A.1, tem-se

$$y_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_k x_{ki} + \varepsilon_i, \quad i = 1, 2, \dots, n \quad (\text{A.7})$$

Assume-se que o erro ε_i são variáveis aleatórias independentes, com média zero, ou seja, $E(\varepsilon_i) = 0$, e variância constante igual a σ^2 .

Tabela A.1 – Dados para o Modelo de Regressão Linear

x_1	x_2	...	x_k	y
x_{11}	x_{12}	...	x_{1k}	y_1
x_{21}	x_{22}	...	x_{2k}	y_2
.
.
.
x_{n1}	x_{n2}	...	x_{nk}	y_n

Na forma matricial, a Equação (A.7) é representada da seguinte forma:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \quad (\text{A.8})$$

onde

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}, \quad \mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1k} \\ 1 & x_{21} & x_{22} & \cdots & x_{2k} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{nk} \end{bmatrix}, \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix} \quad e \quad \boldsymbol{\varepsilon} = \begin{bmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{bmatrix}$$

Nesse caso, \mathbf{y} é um vetor $n \times 1$ das observações, \mathbf{X} é uma matriz $n \times p$ ($p = k + 1$) dos níveis das variáveis independentes, $\boldsymbol{\beta}$ é um vetor $p \times 1$ dos coeficientes de regressão e $\boldsymbol{\varepsilon}$ é um vetor $n \times 1$ dos erros aleatórios.

O método dos mínimos quadrados escolhe os β 's na Equação (A.7) de tal forma que a soma dos quadrados dos erros ε_i sejam minimizados. Para isso, define-se primeiramente a função de mínimos quadrados L .

$$L = \sum_{i=1}^n \varepsilon_i^2 = \boldsymbol{\varepsilon}'\boldsymbol{\varepsilon} = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = \mathbf{y}'\mathbf{y} - \boldsymbol{\beta}'\mathbf{X}'\mathbf{y} - \mathbf{y}'\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\beta}'\mathbf{X}'\mathbf{X}\boldsymbol{\beta}$$

Sabendo que $\boldsymbol{\beta}'\mathbf{X}'\mathbf{y}$ é uma matriz 1×1 , a sua transposta $\mathbf{y}'\mathbf{X}\boldsymbol{\beta}$ é ela própria. Com isso, a função L pode ser expressa da seguinte forma:

$$L = \mathbf{y}'\mathbf{y} - 2\boldsymbol{\beta}'\mathbf{X}'\mathbf{y} + \boldsymbol{\beta}'\mathbf{X}'\mathbf{X}\boldsymbol{\beta} \quad (\text{A.9})$$

Como a função L deve ser minimizada em relação a $\boldsymbol{\beta}$, os estimadores de mínimos quadrados, ou seja, $\hat{\boldsymbol{\beta}}$, deve satisfazer a

$$\left. \frac{\partial L}{\partial \boldsymbol{\beta}} \right|_{\hat{\boldsymbol{\beta}}} = -2\mathbf{X}'\mathbf{y} + 2\mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} = 0$$

Simplificando,

$$\mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}'\mathbf{y} \quad (\text{A.10})$$

A Equação (A.10) é o conjunto de equações normais de mínimos quadrados representados na forma matricial. Desde que $\mathbf{X}'\mathbf{X}$ seja positiva definida, pode-se resolver a Equação (A.10) multiplicando ambos os seus membros por $(\mathbf{X}'\mathbf{X})^{-1}$.

Dessa forma, os estimadores de mínimos quadrados de $\boldsymbol{\beta}$ são

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y} \quad (\text{A.11})$$

e o modelo de regressão ajustado é

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}} \quad (\text{A.12})$$

Na forma escalar, o modelo ajustado é

$$\hat{y}_i = \hat{\beta}_0 + \sum_{j=1}^k \hat{\beta}_j x_{ij}, \quad i = 1, 2, \dots, n \quad (\text{A.13})$$

A diferença entre a observação y_i e o valor ajustado \hat{y}_i é o resíduo. Com isso, o vetor $n \times 1$ dos resíduos é

$$\mathbf{e} = \mathbf{y} - \hat{\mathbf{y}} \quad (\text{A.14})$$

Nos modelos de regressão linear, o método dos mínimos quadrados produz estimadores não enviesados dos parâmetros $\boldsymbol{\beta}$. Portanto, $E(\hat{\boldsymbol{\beta}}) = \boldsymbol{\beta}$ (Myers e Montgomery, 2002, pág. 25).

A variância de $\hat{\boldsymbol{\beta}}$ pode ser obtida a partir da matriz de variância-covariância:

$$\text{var}(\hat{\boldsymbol{\beta}}) = E \left[(\hat{\boldsymbol{\beta}} - E(\hat{\boldsymbol{\beta}})) (\hat{\boldsymbol{\beta}} - E(\hat{\boldsymbol{\beta}}))' \right] \quad (\text{A.15})$$

que é uma matriz simétrica $p \times p$, cujo i -ésimo elemento da diagonal principal é a variância de β_i e o elemento (ij) é a covariância entre β_i e β_j . A matriz de covariância de $\hat{\boldsymbol{\beta}}$ é (Myers e Montgomery, 2002, pág. 27):

$$\text{var}(\hat{\boldsymbol{\beta}}) = \sigma^2 (\mathbf{X}'\mathbf{X})^{-1} \quad (\text{A.16})$$

O estimador dos mínimos quadrados de $\boldsymbol{\beta}$ é um estimador linear não enviesado e de variância mínima, o que lhe confere o título de melhor estimador linear não enviesado.

Pode-se demonstrar (Myers e Montgomery, 2002, pág. 27) que a estimativa da variância σ^2 do erro ε é relação entre a soma dos quadrados dos resíduos SS_E pelo número de graus de liberdade $(n - p)$, ou seja:

$$\hat{\sigma}^2 = \frac{SS_E}{n - p} \quad (\text{A.17})$$

sendo a soma dos quadrados dos resíduos obtida da seguinte forma:

$$SS_E = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (\text{A.18})$$

A.2 Intervalos para a Média e para a Previsão da Resposta

Dado um ponto $x_{01}, x_{02}, \dots, x_{0k}$, no espaço das variáveis regressoras, tem-se o vetor

$$\mathbf{x}_0 = \begin{bmatrix} 1 \\ x_{01} \\ x_{02} \\ \vdots \\ x_{0k} \end{bmatrix}$$

pode-se obter um intervalo de confiança para a média da resposta e o intervalo de previsão para uma futura observação de um valor individual da resposta.

A.2.1 Intervalo de Confiança para a Média da Resposta

A média da resposta no ponto \mathbf{x}_0 é

$$\mu_{y|\mathbf{x}_0} = \beta_0 + \beta_1 x_{01} + \beta_2 x_{02} + \dots + \beta_k x_{0k}$$

O estimador da média da resposta neste ponto é

$$\hat{\mu}_{y|\mathbf{x}_0} = \hat{y}(\mathbf{x}_0) = \mathbf{x}'_0 \hat{\boldsymbol{\beta}} \quad (\text{A.24})$$

O estimador não é enviesado, pois

$$E[\hat{\mu}_{y|\mathbf{x}_0}] = E(\mathbf{x}'_0 \hat{\boldsymbol{\beta}}) = \mathbf{x}'_0 \boldsymbol{\beta} = \mu_{y|\mathbf{x}_0}$$

A variância do estimador da média da resposta é

$$\text{var}[\hat{\mu}_{y|\mathbf{x}_0}] = \text{var}[\hat{y}(\mathbf{x}_0)] = \text{var}(\mathbf{x}'_0 \hat{\boldsymbol{\beta}}) \quad (\text{A.25})$$

Na Equação (2.16) tem-se que:

$$\text{var}(\hat{\boldsymbol{\beta}}) = \sigma^2 (\mathbf{X}'\mathbf{X})^{-1}$$

Então,

$$\text{var}[\hat{\mu}_{y|\mathbf{x}_0}] = \sigma^2 \mathbf{x}'_0 (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_0$$

Como a distribuição de y é normal, o quociente

$$t = \frac{\hat{\mu}_{y|\mathbf{x}_0} - \mu_{y|\mathbf{x}_0}}{\text{var}(\hat{\mu}_{y|\mathbf{x}_0})} \quad (\text{A.26})$$

tem distribuição t com $(n - p)$ graus de liberdade.

Portanto, para um intervalo de confiança de $100(1 - \alpha)\%$ tem-se que:

$$-t_{\alpha/2, n-p} \leq \frac{\hat{\mu}_{y|\mathbf{x}_0} - \mu_{y|\mathbf{x}_0}}{\sqrt{\text{var}(\hat{\mu}_{y|\mathbf{x}_0})}} \leq t_{\alpha/2, n-p} \quad (\text{A.27})$$

Nas Equações (A.24) e (A.25) tem-se que $\hat{\mu}_{y|\mathbf{x}_0} = \hat{y}(\mathbf{x}_0)$ e $\text{var}[\hat{\mu}_{y|\mathbf{x}_0}] = \sigma^2 \mathbf{x}'_0 (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_0$.

Substituindo em (A.27) tem-se que:

$$-t_{\alpha/2, n-p} \leq \frac{\hat{y}(\mathbf{x}_0) - \mu_{y|\mathbf{x}_0}}{\sqrt{\sigma^2 \mathbf{x}'_0 (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_0}} \leq t_{\alpha/2, n-p}.$$

O que é equivalente a

$$\begin{aligned} \hat{y}(\mathbf{x}_0) - t_{\alpha/2, n-p} \sqrt{\sigma^2 \mathbf{x}'_0 (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_0} &\leq \mu_{y|\mathbf{x}_0} \\ &\leq \hat{y}(\mathbf{x}_0) + t_{\alpha/2, n-p} \sqrt{\sigma^2 \mathbf{x}'_0 (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_0} \end{aligned} \quad (\text{A.28})$$

que é o intervalo de confiança de $100(1 - \alpha)\%$ para a média da resposta no ponto $\mathbf{x}_{01}, \mathbf{x}_{02}, \dots, \mathbf{x}_{0k}$.

A.2.2 Intervalo de Previsão para uma Futura Resposta

O modelo no ponto \mathbf{x}_0 é

$$y(\mathbf{x}_0) = \hat{\beta}_0 + \hat{\beta}_1 x_{01} + \hat{\beta}_2 x_{02} + \dots + \hat{\beta}_k x_{0k} + \varepsilon \text{ ou}$$

$$y(\mathbf{x}_0) = \mathbf{x}'_0 \hat{\boldsymbol{\beta}} + \varepsilon$$

A estimativa de uma nova resposta neste ponto é a mesma estimativa da média:

$$E[y(\mathbf{x}_0)] = \mathbf{x}'_0 \hat{\boldsymbol{\beta}}$$

A variância de uma nova resposta neste ponto é

$$\text{var}[y(\mathbf{x}_0)] = \text{var}(\mathbf{x}'_0 \hat{\boldsymbol{\beta}} + \boldsymbol{\varepsilon}) = \text{var}(\mathbf{x}'_0 \hat{\boldsymbol{\beta}}) + \text{var}(\boldsymbol{\varepsilon}) = \sigma^2 \mathbf{x}'_0 (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_0 + \sigma^2$$

ou

$$\text{var}[y(\mathbf{x}_0)] = \sigma^2 (\mathbf{x}'_0 (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_0 + 1) \quad (\text{A.29})$$

e o intervalo de $100(1 - \alpha)\%$ de probabilidade para uma nova resposta no ponto

$\mathbf{x}_{01}, \mathbf{x}_{02}, \dots, \mathbf{x}_{0k}$ é

$$\begin{aligned} \hat{y}(\mathbf{x}_0) - t_{\alpha/2, n-p} \sqrt{\text{var}[y(\mathbf{x}_0)]} &\leq y(\mathbf{x}_0) \leq \hat{y}(\mathbf{x}_0) + t_{\alpha/2, n-p} \sqrt{\text{var}[y(\mathbf{x}_0)]} \\ \hat{y}(\mathbf{x}_0) - t_{\alpha/2, n-p} \sqrt{\sigma^2 [\mathbf{x}'_0 (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_0 + 1]} &\leq y(\mathbf{x}_0) \\ &\leq \hat{y}(\mathbf{x}_0) + t_{\alpha/2, n-p} \sqrt{\sigma^2 [\mathbf{x}'_0 (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_0 + 1]} \end{aligned} \quad (\text{A.30})$$