

# Modelos de Regressão Linear Simples - parte I

Erica Castilho Rodrigues

27 de Setembro de 2017

Introdução

Coeficiente de Correlação Linear

## Objetivos

Ao final deste capítulo você deve ser capaz de:

- ▶ Usar modelos de regressão para construir modelos para dados coletados.
- ▶ Entender como método de mínimos é usado para estimar parâmetros desconhecidos.

# Introdução

---

- ▶ Podemos estar interessados em explorar a relação entre duas ou mais variáveis.
- ▶ Essa técnica é chamada **Análise de Regressão**.
- ▶ Exemplo: qual relação entre nível de escolaridade e renda?

- ▶ A ferramenta inicial para sabermos se existe relação entre as variáveis é o **gráfico de dispersão**.
- ▶ A correlação poderá:
  - ▶ não existir;
  - ▶ ser uma correlação linear (ao longo de uma reta);
  - ▶ ser uma correlação não linear (ao longo de uma curva).

## Correlação

Há um relacionamento entre as variáveis?

- ▶ Elas aumentam juntas?
- ▶ Aumentando uma variável a outra aumenta ou diminui?
- ▶ Exemplo: nota na prova e horas de estudo.
- ▶ Variam juntas.
- ▶ Se uma aumenta, a outra também aumenta.

**Exemplo:**

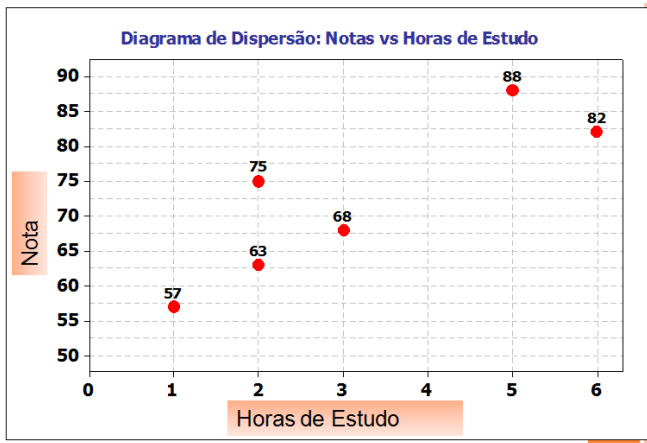
- ▶ Vamos considerar as variáveis nota na prova e horas de estudo.
- ▶  $Y$  - nota na prova (variável resposta).
- ▶  $x$  - horas de estudo (variável explicativa ou preditora).
- ▶ Os dados são os seguintes

Aluno	Horas de Estudo	Nota na Prova
A	6	82
B	2	63
C	1	57
D	5	88
E	3	68
F	2	75



## Exemplo: (continuação)

- ▶ O gráfico de dispersão é mostrado a seguir



- ▶ Parece existir uma correlação positiva (uma aumenta a outra aumenta).

## Exemplo:

- ▶ 50 municípios de um estado são analisados.
- ▶ Deseja-se verificar se existe relação entre duas variáveis:
  - ▶ nível de pobreza da população;
  - ▶ taxa de roubos e furtos do município.
- ▶ Os dados coletados são apresentados a seguir.

Município	Índice de pobreza	Taxa de roubos e furtos	Município	Índice de pobreza	Taxa de roubos e furtos
1	0,520243	4,6796	26	0,861767	8,0104
2	0,573892	5,8444	27	0,863136	7,644
3	0,588446	6,5371	28	0,881571	8,3115
4	0,610826	5,8645	29	0,886649	8,9649
5	0,626467	6,5505	30	0,895995	9,6845
6	0,634336	7,3872	31	0,909950	8,3012
7	0,640390	6,0816	32	0,921256	9,7942
8	0,644050	6,569	33	0,924994	10,3435
9	0,663735	5,6804	34	0,927112	9,202
10	0,676159	7,2103	35	0,934594	8,5556
11	0,692946	5,7382	36	0,944849	10,3021
12	0,719850	6,3769	37	0,952311	9,909
13	0,721679	7,7359	38	0,954732	10,7073
14	0,764735	8,6539	39	0,956666	9,2998
15	0,781271	8,0238	40	0,958970	8,8892
16	0,786335	7,3774	41	0,958975	9,7285
17	0,797273	8,5611	42	0,965260	9,0648
18	0,806268	7,9516	43	0,968608	10,4605
19	0,811569	7,0226	44	0,969550	9,9478
20	0,821166	8,7187	45	0,980406	8,8263
21	0,821884	7,8224	46	0,982279	9,5566
22	0,828157	8,2076	47	0,982572	10,8394
23	0,845534	8,3666	48	0,983857	9,2011
24	0,857904	9,4827	49	0,994508	9,2753
25	0,858685	8,8655	50	0,995053	10,0008

## Exemplo: (continuação)

### Pergunta

Existe associação entre níveis de pobreza e taxas de roubos e furtos do município?

### Ou seja...

O nível de pobreza de um município determina (explica, prediz, interfere no) nível de criminalidade (roubo/furto) do município?

### Hipótese

Quanto maior os níveis de pobreza de um município, maior a taxa de criminalidade.

**Exemplo:** (continuação)

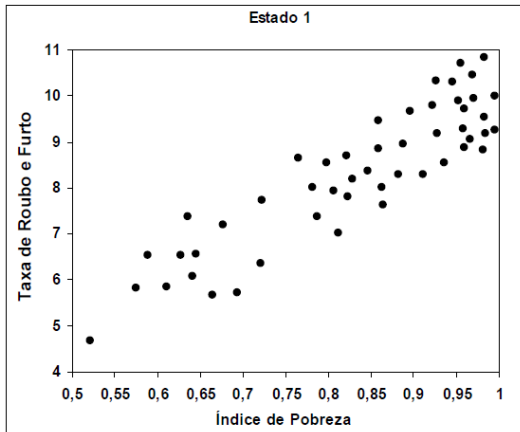
- ▶ **Variável Resposta ou Dependente:**

**Exemplo:** (continuação)

- ▶ **Variável Resposta ou Dependente:** criminalidade.
- ▶ **Variável Explicativa ou Independente:**

## Exemplo: (continuação)

- ▶ **Variável Resposta ou Dependente:** criminalidade.
- ▶ **Variável Explicativa ou Independente:** pobreza.
- ▶ Gráfico de dispersão é mostrado a seguir:



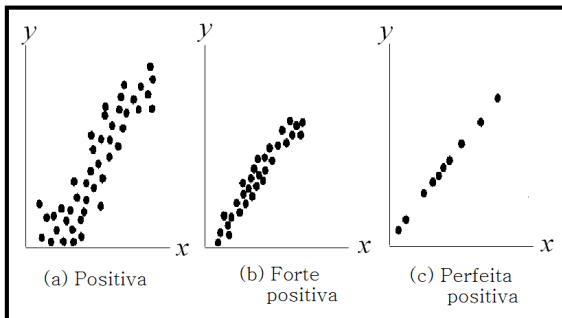
- ▶ Qual a conclusão que você tira?

- ▶ Qual a conclusão que você tira?
- ▶ Parece existir uma associação positiva entre as variáveis.
- ▶ Vamos ver como podemos medir se essa associação é fraca ou forte.



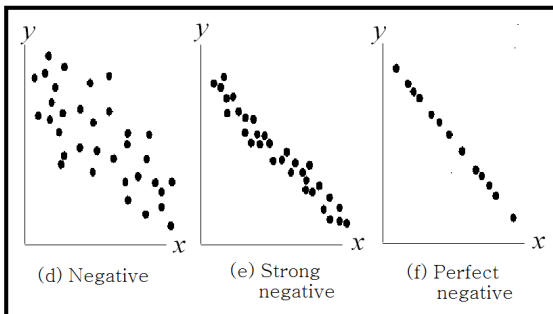
## Correlação linear positiva

Uma variável aumenta a outra também aumenta.



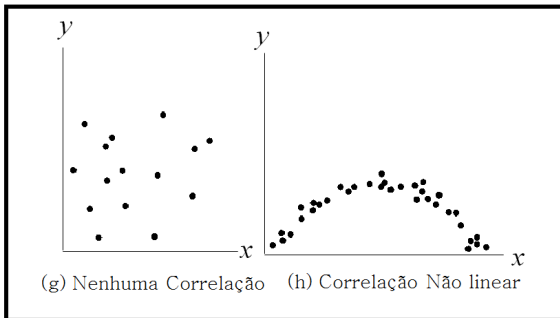
## Correlação linear negativa

Uma variável aumenta a outra diminui.



## Correlação não linear

Não existe relação linear entre as variáveis.



## Coeficiente de correlação linear

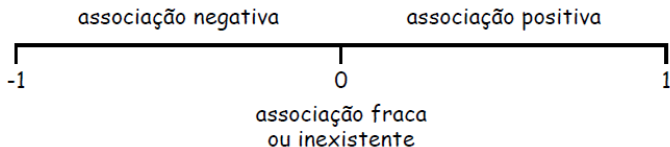
Mede a força da **relação linear** entre duas variáveis.

- ▶ Denotado por  $r$  se for amostral e  $\rho$  se for populacional.
- ▶ Mede o grau de associação linear entre duas variáveis.
- ▶ Indica também se essa associação é positiva ou negativa.
- ▶ É dado por

$$r = \frac{n \sum_i x_i y_i - (\sum x_i)(\sum y_i)}{\sqrt{n(\sum_i x_i^2) - (\sum_i x_i)^2} \sqrt{n(\sum_i y_i^2) - (\sum_i y_i)^2}}$$

## Propriedades do coeficiente de correlação linear

- ▶  $-1 \leq r \leq 1$
- ▶ Se  $r \approx 1 \Rightarrow$  correlação forte positiva.
- ▶ Se  $r \approx -1 \Rightarrow$  correlação forte negativa.
- ▶ Se  $r \approx 0 \Rightarrow$  não existe correlação **linear**.
- ▶ O valor de  $r$  não é influenciado pelas escalas de  $x$  e  $y$ .



- ▶ A significância dos valores de  $r$  depende muito da área em que estamos trabalhando.
- ▶ Em algumas áreas, não se espera que uma variável explique bem a outra.
- ▶ Assim não se espera valores muito altos para  $r$ .
- ▶ A tabela a seguir apresenta um guia geral:

Valor de $r$	Intensidade da associação	Valor de $r$
$-0,2 < r \leq 0$	inexistente	$0 \leq r < 0,2$
$-0,4 < r \leq -0,2$	muito fraca	$0,2 \leq r < 0,4$
$-0,7 < r \leq -0,4$	fraca/moderada	$0,4 \leq r < 0,7$
$-0,9 < r \leq -0,7$	forte	$0,7 \leq r < 0,9$
$-1 < r \leq -0,9$	muito forte	$0,9 \leq r < 1$
$r = -1$	perfeita	$r = 1$

## Exemplo

- ▶ Vamos retomar o exemplo dos níveis de pobreza e criminalidade.

## Pergunta

Existe associação entre níveis de pobreza e taxas de roubos e furtos do município?

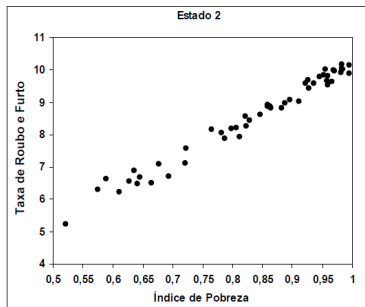
## Ou seja...

O nível de pobreza de um município determina (explica, prediz, interfere no) nível de criminalidade (roubo/furto) do município?

## Hipótese

Quanto maior os níveis de pobreza de um município, maior a taxa de criminalidade.

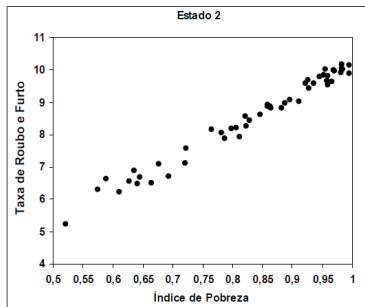
## Exemplo (continuação)



- ▶ Existe associação linear entre as variáveis?

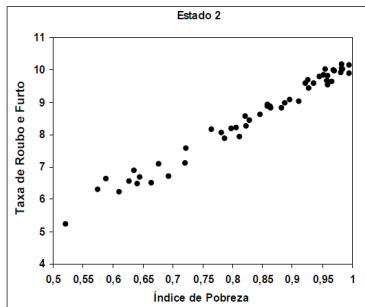


## Exemplo (continuação)



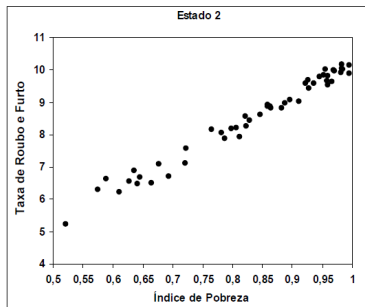
- ▶ Existe associação linear entre as variáveis? Sim.
- ▶ Ela é forte ou fraca?

## Exemplo (continuação)



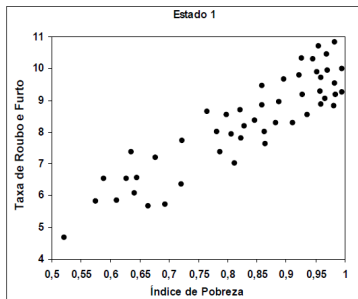
- ▶ Existe associação linear entre as variáveis? Sim.
- ▶ Ela é forte ou fraca? Forte.
- ▶ O valor de  $r$  deve estar em torno de quanto?

## Exemplo (continuação)



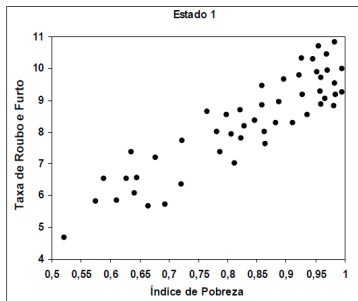
- ▶ Existe associação linear entre as variáveis? Sim.
- ▶ Ela é forte ou fraca? Forte.
- ▶ O valor de  $r$  deve estar em torno de quanto?  
 $r = 0,989$ .

## Exemplo (continuação)



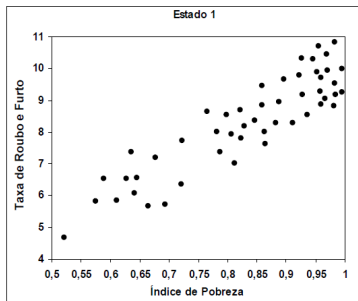
- ▶ Existe associação linear entre as variáveis?

## Exemplo (continuação)



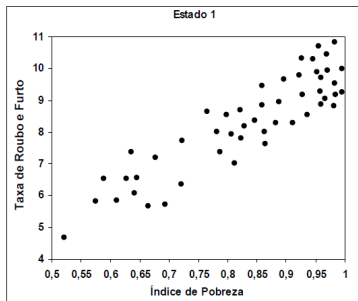
- ▶ Existe associação linear entre as variáveis? Sim.
- ▶ Ela é forte ou fraca?

## Exemplo (continuação)



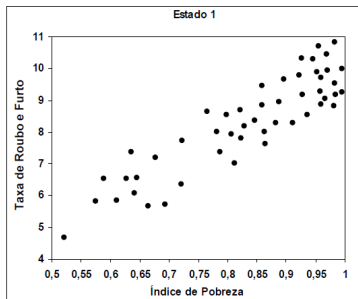
- ▶ Existe associação linear entre as variáveis? Sim.
- ▶ Ela é forte ou fraca? Forte.
- ▶ O valor de  $r$  é mais baixo ou mais alto que o anterior?

## Exemplo (continuação)



- ▶ Existe associação linear entre as variáveis? Sim.
- ▶ Ela é forte ou fraca? Forte.
- ▶ O valor de  $r$  é mais baixo ou mais alto que o anterior? Mais baixo.
- ▶ O valor de  $r$  deve estar em torno de quanto?

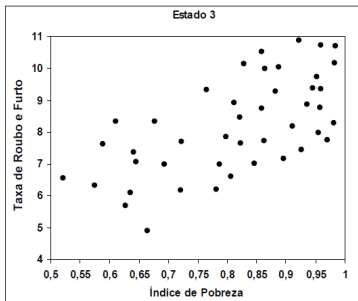
## Exemplo (continuação)



- ▶ Existe associação linear entre as variáveis? Sim.
- ▶ Ela é forte ou fraca? Forte.
- ▶ O valor de  $r$  é mais baixo ou mais alto que o anterior? Mais baixo.
- ▶ O valor de  $r$  deve estar em torno de quanto?  
 $r = 0,898$ .

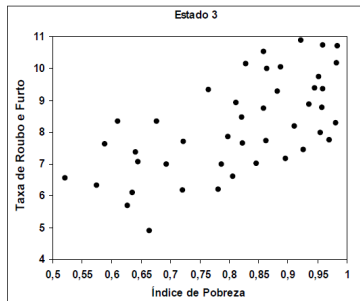


## Exemplo (continuação)



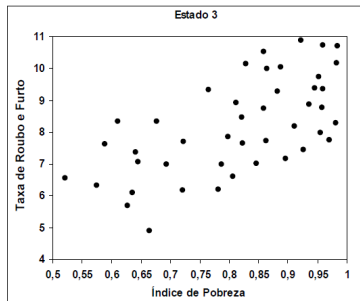
- ▶ Existe associação linear entre as variáveis?

## Exemplo (continuação)



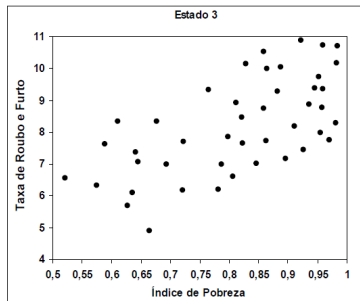
- ▶ Existe associação linear entre as variáveis? Sim.
- ▶ Ela é forte, moderada ou fraca?

## Exemplo (continuação)



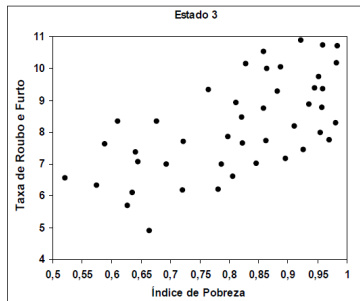
- ▶ Existe associação linear entre as variáveis? Sim.
- ▶ Ela é forte, moderada ou fraca? Moderada.
- ▶ O valor de  $r$  é mais baixo ou mais alto que o anterior?

## Exemplo (continuação)



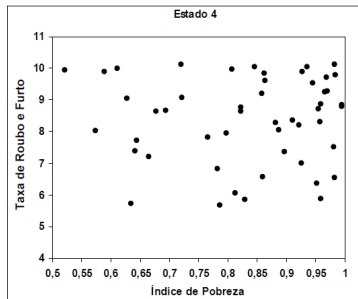
- ▶ Existe associação linear entre as variáveis? Sim.
- ▶ Ela é forte, moderada ou fraca? Moderada.
- ▶ O valor de  $r$  é mais baixo ou mais alto que o anterior? Mais baixo.
- ▶ O valor de  $r$  deve estar em torno de quanto?

## Exemplo (continuação)



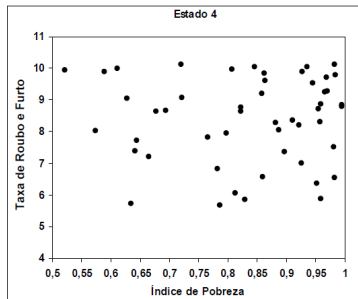
- ▶ Existe associação linear entre as variáveis? Sim.
- ▶ Ela é forte, moderada ou fraca? Moderada.
- ▶ O valor de  $r$  é mais baixo ou mais alto que o anterior? Mais baixo.
- ▶ O valor de  $r$  deve estar em torno de quanto?  
 $r = 0,692$ .

## Exemplo (continuação)



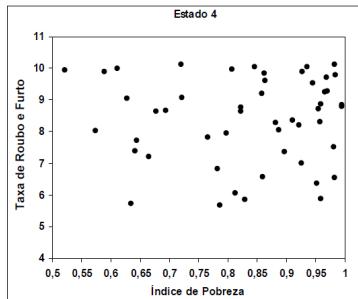
- ▶ Existe associação linear entre as variáveis?

## Exemplo (continuação)



- ▶ Existe associação linear entre as variáveis? Não.
- ▶ O valor de  $r$  é mais baixo ou mais alto que o anterior?

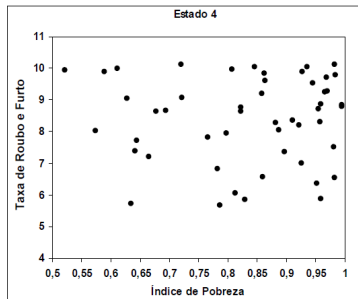
## Exemplo (continuação)



- ▶ Existe associação linear entre as variáveis? Não.
- ▶ O valor de  $r$  é mais baixo ou mais alto que o anterior? Mais baixo.
- ▶ O valor de  $r$  deve estar em torno de quanto?



## Exemplo (continuação)



- ▶ Existe associação linear entre as variáveis? Não.
- ▶ O valor de  $r$  é mais baixo ou mais alto que o anterior? Mais baixo.
- ▶ O valor de  $r$  deve estar em torno de quanto?  
 $r = 0,019$ .

## Exemplo

- ▶ Queremos verificar se existe associação entre idade e pressão sanguínea.
- ▶ Os dados são mostrados a seguir:

Aluno	Age	Blood Pressure	Age* BP	age <sup>2</sup>	BP <sup>2</sup>
A	43	128	5504	1849	16384
B	48	120	5760	2304	14400
C	56	135	7560	3136	18225
D	61	143	8723	3721	20449
E	67	141	9447	4489	19881
F	70	152	10640	4900	23104
Soma	345	819	47640	20399	112443

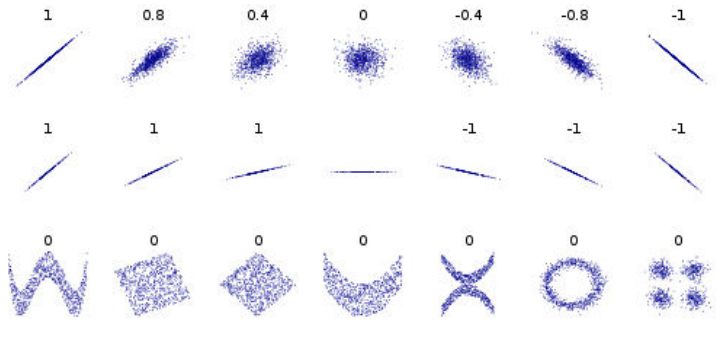
## Exemplo (continuação)

- ▶ Temos então que

$$r = \frac{(6)(47634) - (345)(819)}{\sqrt{(6)(20399) - 345^2} \sqrt{(6)(112443) - 819^2}}$$
$$= 0,897$$

- ▶ As variáveis estão fortemente associadas.
- ▶ A associação é positiva:
  - ▶ quanto maior a idade  $\Rightarrow$  maior a pressão sanguínea.

- ▶ A figura abaixo mostra valores de  $r$  para vários conjuntos de dados distintos.



## Coeficiente de correlação populacional $\rho$

- ▶ É uma característica da população e não da amostra.
- ▶ Não sabemos seu valor verdadeiro.
- ▶ Mas podemos estimá-lo pelo coeficiente amostral.
- ▶ Pode-se fazer testes e intervalos de confiança para esse parâmetro.
- ▶ Pode ser que na população  $\rho = 0$ ,
  - ▶ mas na amostra  $r \neq 0$  por mero acaso.
- ▶ Vamos ver a seguir como podemos testar se o coeficiente é significativo.

## Teste do Coeficiente de Correlação

- ▶ O coeficiente de correlação  $r$  é apenas uma estimativa amostral.
- ▶ Ele é calculado com base em uma amostra de tamanho  $n$ .
- ▶ Os valores amostrais podem apresentar uma correlação, mas a população não.
- ▶ Se  $r \neq 0$  não garante que  $\rho \neq 0$ .
- ▶ Podemos fazer um teste de hipótese para verificar se de fato  $\rho \neq 0$ .

- ▶ As hipóteses a serem testadas são as seguintes:

$$H_0 : \rho = 0 \quad \text{vs} \quad H_1 : \rho \neq 0 .$$

- ▶ Sob  $H_0$  temos que  $\rho = 0$ .
- ▶ Além disso, pode-se mostrar que

$$\text{Var}(r) = \frac{1 - \rho^2}{n - 2}$$

que é estimada por

$$\text{Var}(r) = \frac{1 - r^2}{n - 2} .$$

- ▶ A estatística de teste é dada por:

$$t = \frac{r - \text{valor sob } H_0}{\sqrt{\text{Var}(r)}}$$

- ▶ A estatística de teste fica

$$t = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}} .$$

- ▶ Como a variância está sendo estimada, essa estatística tem uma distribuição t-student com  $n - 2$  graus de liberdade.



## Exemplo

- ▶ Considere o exemplo de índices de pobreza e violência.
- ▶ Temos que

$$r = \frac{n \sum_i x_i y_i - (\sum_i x_i)(\sum_i y_i)}{\sqrt{n(\sum_i x_i^2) - (\sum_i x_i)^2} \sqrt{n(\sum_i y_i^2) - (\sum_i y_i)^2}} = 0,898.$$

- ▶ Queremos testar

$$H_0 : \rho = 0 \quad \text{vs} \quad H_1 : \rho \neq 0.$$

- ▶ A estatística de teste é dada por

$$t = \frac{r\sqrt{n-2}}{\sqrt{1-r^2}} = \frac{0,898\sqrt{50-2}}{\sqrt{1-0,898^2}} = 14,25.$$

## Exemplo (continuação)

- ▶ Vamos considerar  $\alpha = 0,05$ .
- ▶ Temos  $t_{48;0,025} \approx z_{0,025}$  pois  $n$  é grande.
- ▶  $z_{0,025} = 1,96$
- ▶ Qual a região crítica do teste?

## Exemplo (continuação)

- ▶ Vamos considerar  $\alpha = 0,05$ .
- ▶ Temos  $t_{48;0,025} \approx z_{0,025}$  pois  $n$  é grande.
- ▶  $z_{0,025} = 1,96$
- ▶ Qual a região crítica do teste?

$$t > 1,96 \text{ ou } t < -1,96 .$$

- ▶ Qual a conclusão?

## Exemplo (continuação)

- ▶ Vamos considerar  $\alpha = 0,05$ .
- ▶ Temos  $t_{48;0,025} \approx z_{0,025}$  pois  $n$  é grande.
- ▶  $z_{0,025} = 1,96$
- ▶ Qual a região crítica do teste?

$$t > 1,96 \text{ ou } t < -1,96 .$$

- ▶ Qual a conclusão?
- ▶ Como  $t = 14,25 > 1,96$ , rejeitamos  $H_0$ .
- ▶ Com 5% de significância há evidências de que  $\rho \neq 0$ , ou seja, de que existe associação entre as variáveis.

## Exemplo (continuação)

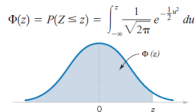
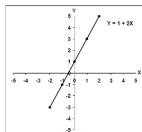


Table II Cumulative Standard Normal Distribution (continued)

z	0.00	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09
0.0	0.50000	0.50398	0.50797	0.51196	0.51595	0.51993	0.52392	0.52790	0.53188	0.53585
0.1	0.53982	0.54379	0.54778	0.55177	0.55576	0.55974	0.56373	0.56771	0.57169	0.57568
0.2	0.57926	0.58316	0.58706	0.59095	0.59483	0.59870	0.60258	0.60644	0.61026	0.61409
0.3	0.61791	0.62171	0.62551	0.62930	0.63307	0.63683	0.64057	0.64429	0.64802	0.65173
0.4	0.65542	0.65907	0.66275	0.66640	0.67003	0.67364	0.67724	0.68082	0.68438	0.68793
0.5	0.69146	0.69497	0.69846	0.70194	0.70540	0.70884	0.71226	0.71566	0.71904	0.72240
0.6	0.72574	0.72909	0.73237	0.73563	0.73891	0.74215	0.74537	0.74857	0.75174	0.75490
0.7	0.75803	0.76114	0.76423	0.76730	0.77035	0.77337	0.77637	0.77935	0.78230	0.78526
0.8	0.78814	0.79103	0.79389	0.79673	0.79954	0.80233	0.80510	0.80785	0.81057	0.81327
0.9	0.81594	0.81858	0.82121	0.82381	0.82639	0.82894	0.83147	0.83397	0.83645	0.83891
1.0	0.84134	0.84375	0.84613	0.84849	0.85083	0.85314	0.85542	0.85769	0.85992	0.86214
1.1	0.86434	0.86650	0.86864	0.87076	0.87285	0.87492	0.87697	0.87899	0.88100	0.88297
1.2	0.88493	0.88680	0.88876	0.89065	0.89251	0.89435	0.89616	0.89795	0.89972	0.90147
1.3	0.90319	0.90490	0.90658	0.90824	0.90987	0.91149	0.91308	0.91465	0.91620	0.91776
1.4	0.91924	0.92073	0.92219	0.92361	0.92506	0.92647	0.92785	0.92921	0.93056	0.93188
1.5	0.93319	0.93447	0.93574	0.93699	0.93820	0.93942	0.94062	0.94179	0.94294	0.94408
1.6	0.94520	0.94630	0.94738	0.94844	0.94949	0.95052	0.95154	0.95254	0.95351	0.95448
1.7	0.95543	0.95637	0.95728	0.95818	0.95907	0.95994	0.96079	0.96163	0.96246	0.96327
1.8	0.96407	0.96485	0.96562	0.96637	0.96711	0.96784	0.96855	0.96925	0.96994	0.97062
1.9	0.97128	0.97193	0.97257	0.97319	0.97381	0.97441	0.97500	0.97558	0.97614	0.97670

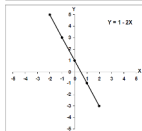
# Revisão de matemática: equação da reta

$$y = a + bx$$



$$y = 1 + 2x$$

X	Y
-2	-3
-1	-1
0	1
1	3
2	5



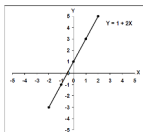
$$y = 1 - 2x$$

X	Y
-2	5
-1	3
0	1
1	-1
2	-3

- ▶ Coeficiente Linear ou Intercepto (a):

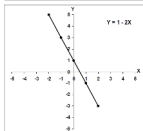
## Revisão de matemática: equação da reta

$$y = a + bx$$



$$y = 1 + 2x$$

X	Y
-2	-3
-1	-1
0	1
1	3
2	5



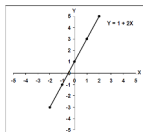
$$y = 1 - 2x$$

X	Y
-2	5
-1	3
0	1
1	-1
2	-3

- ▶ Coeficiente Linear ou Intercepto (a): valor de  $y$  quando  $x = 0$ .
- ▶ Coeficiente Angular ou Inclinação da reta (b):
  - ▶  $b > 0$  a reta é

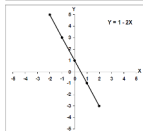
## Revisão de matemática: equação da reta

$$y = a + bx$$



$$y = 1 + 2x$$

X	Y
-2	-3
-1	-1
0	1
1	3
2	5



$$y = 1 - 2x$$

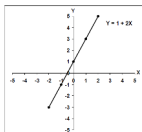
X	Y
-2	5
-1	3
0	1
1	-1
2	-3

- ▶ Coeficiente Linear ou Intercepto (a): valor de  $y$  quando  $x = 0$ .
- ▶ Coeficiente Angular ou Inclinação da reta (b):
  - ▶  $b > 0$  a reta é crescente ( $x$  cresce,  $y$  cresce)
  - ▶  $b < 0$  a reta é



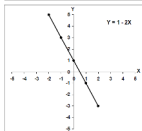
## Revisão de matemática: equação da reta

$$y = a + bx$$



$$y = 1 + 2x$$

X	Y
-2	-3
-1	-1
0	1
1	3
2	5



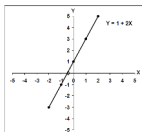
$$y = 1 - 2x$$

X	Y
-2	5
-1	3
0	1
1	-1
2	-3

- ▶ Coeficiente Linear ou Intercepto (a): valor de  $y$  quando  $x = 0$ .
- ▶ Coeficiente Angular ou Inclinação da reta (b):
  - ▶  $b > 0$  a reta é crescente ( $x$  cresce,  $y$  cresce)
  - ▶  $b < 0$  a reta é decrescente ( $x$  cresce,  $y$  decresce)
  - ▶  $b = 0$  reta é

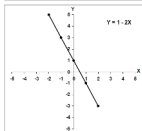
## Revisão de matemática: equação da reta

$$y = a + bx$$



$$y = 1 + 2x$$

X	Y
-2	-3
-1	-1
0	1
1	3
2	5



$$y = 1 - 2x$$

X	Y
-2	5
-1	3
0	1
1	-1
2	-3

- ▶ Coeficiente Linear ou Intercepto (a): valor de  $y$  quando  $x = 0$ .
- ▶ Coeficiente Angular ou Inclinação da reta (b):
  - ▶  $b > 0$  a reta é crescente ( $x$  cresce,  $y$  cresce)
  - ▶  $b < 0$  a reta é decrescente ( $x$  cresce,  $y$  decresce)
  - ▶  $b = 0$  reta é paralela ao eixo  $x$  ( $x$  cresce,  $y$  não muda).

- ▶ Qual a derivada de  $a + bx$  em relação a  $x$ ?

- ▶ Qual a derivada de  $a + bx$  em relação a  $x$ ?  $b$ .
- ▶ Qual interpretação da derivada?

- ▶ Qual a derivada de  $a + bx$  em relação a  $x$ ?  $b$ .
- ▶ Qual interpretação da derivada?
- ▶ Quanto  $y$  varia quando  $x$  varia.
- ▶ O  $b$  representa o número de unidades que  $y$  aumenta ou diminui quando  $x$  aumenta em uma unidade.

## Regressão Linear

- ▶ Verificamos até agora se existe correlação ou não entre as variáveis.
- ▶ Se existe, podemos querer descobrir a forma dessa associação.
- ▶ Queremos estimar a função que determina a relação entre as variáveis.
- ▶ Podemos usar a equação ajustada para prever valores da variável resposta.

## Exemplo:

- ▶ Estamos analisando um processo químico.
- ▶ O rendimento do produto está relacionado com a temperatura do processo.
- ▶ Podemos construir um modelo que seja capaz de:
  - ▶ prever o rendimento para uma dada temperatura.
- ▶ Esse modelo pode ser usado na otimização do processo:
  - ▶ encontrar a temperatura que maximiza o rendimento.

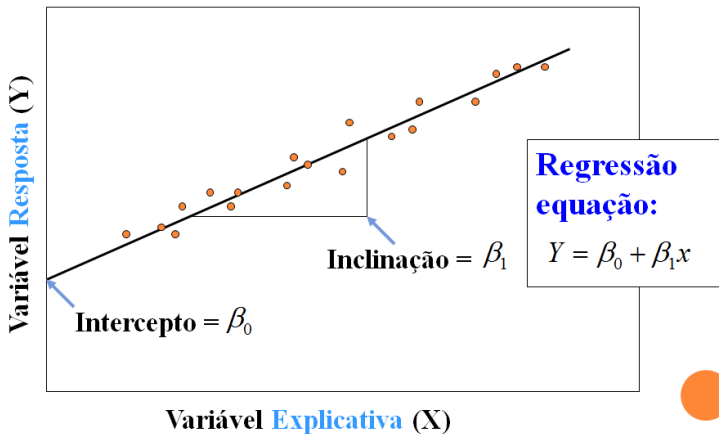
- ▶ É estabelecida uma equação

$$\underbrace{Y}_{\text{resposta}} = \beta_0 + \beta_1 \underbrace{X}_{\text{explicativa}}$$

onde

- ▶  $\beta_0$  é o intercepto em  $Y$  ( $x=0$ );
- ▶  $\beta_1$  é inclinação (taxa de mudança).
- ▶ Veremos que essa equação não é exata.
- ▶ Precisamos incluir um erro aleatório.
- ▶ Vamos considerar assim por enquanto.



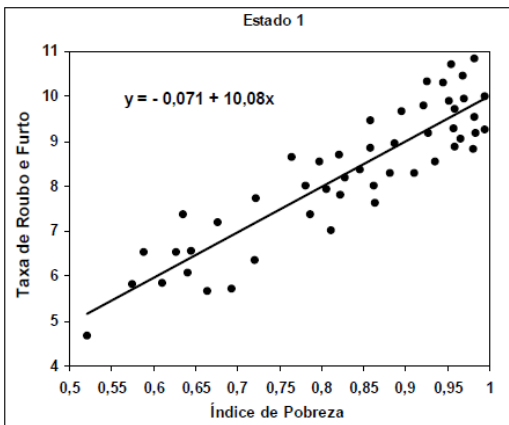


## Exemplo:

- ▶ Considere novamente o exemplo de associação entre pobreza e criminalidade.
- ▶ Vimos que existe uma forte associação entre as variáveis.
- ▶ Podemos escrever a variável taxa de criminalidade em função da variável pobreza:

$$\text{Taxa Criminalidade} = -0,7 + 10,08\text{Pobreza}$$

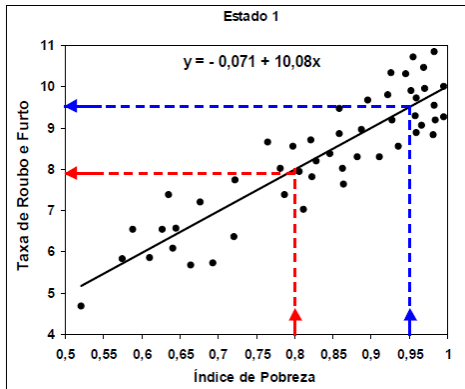
## Exemplo: (continuação)



$$r = 0,898$$

$$R^2 = 80,7\%$$

- ▶ Considere um município com índice de pobreza  $X = 0,8$ .
- ▶ O valor esperado da taxa de furto é  
 $-0,7 + 10,08 * (0,8) = 7,99$  casos por mil habitantes



## Exemplo:

- ▶ Vamos olhar a relação entre:
  - ▶ y - pureza do oxigênio produzido em um processo de destilação;
  - ▶ x - porcentagem de hidrocarbonetos presentes no condensador.

Table 11-1 Oxygen and Hydrocarbon Levels

Observation Number	Hydrocarbon Level x(%)	Purity y(%)
1	0.99	90.01
2	1.02	89.05
3	1.15	91.43
4	1.29	93.74
5	1.46	96.73
6	1.36	94.45
7	0.87	87.59
8	1.23	91.77
9	1.55	99.42
10	1.40	93.65
11	1.19	93.54
12	1.15	92.52
13	0.98	90.56
14	1.01	89.54
15	1.11	89.85
16	1.20	90.39
17	1.26	93.25
18	1.32	93.41
19	1.43	94.98
20	0.95	87.33

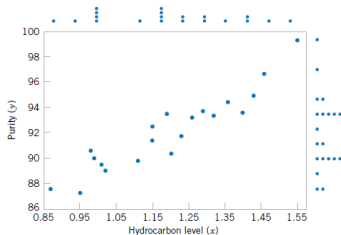


Figure 11-1 Scatter diagram of oxygen purity versus hydrocarbon level from Table 11-1.

**Exemplo: (continuação)**

- ▶ O **gráfico de dispersão** que representa cada par  $(x_i, y_i)$  como um ponto.
- ▶ Nenhuma curva simples passa exatamente por todos pontos.
- ▶ Os pontos parecem estar dispersos aleatoriamente em torno de uma reta.
- ▶ É razoável considerar que a média de  $Y$  esteja relacionada linearmente com  $x$

$$E(Y|x) = \mu_{Y|x} = \beta_0 + \beta_1 x .$$

- ▶  $\beta_0$  e  $\beta_1$  são chamados **coeficientes de regressão**.
- ▶ O valor de  $y$  não cai exatamente sobre a reta.

- ▶ O modelo

$$E(Y|x) = \mu_{Y|x} = \beta_0 + \beta_1 x .$$

descreve a média de  $Y$  e não seu valor observado.

- ▶ Podemos generalizar para um modelo probabilístico.
- ▶ Consideramos que o valor esperado de  $x$  é função linear de  $Y$ .
- ▶ Para um valor fixo de  $x$ , o valor de  $Y$  é dado pela função do valor médio mais um erro aleatório

$$Y = \underbrace{\beta_0 + \beta_1 x}_{\text{Valor médio}} + \epsilon$$

onde  $\epsilon$  é um erro aleatório.

- ▶ Esse modelo é chamado **modelo de regressão linear simples**.

- ▶ Se tivéssemos várias variáveis ( $x_1, x_2, \dots, x_p$ ) no modelo

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \epsilon$$

é chamado **modelo de regressão linear múltipla**.

- ▶ Exemplo: renda explicada pelo sexo, faixa etária, escolaridade, etc.



- ▶ O modelo pode aparecer a parte de uma relação teórica.
- ▶ Exemplo:

$$p = x_0 + v \times t$$

onde

- ▶  $p$  é posição ( $y$ )
  - ▶  $t$  tempo ( $x$ )
  - ▶  $v$  velocidade ( $\beta_1$ )
  - ▶  $x_0$  posição inicial ( $\beta_0$ ).
- ▶ Em outros casos, não sabemos qual relação entre  $y$  e  $x$ .
  - ▶ Então o modelo é escolhido a partir do diagrama de dispersão.
  - ▶ Como foi feito para os dados do oxigênio.

- ▶ Quando não conhecemos uma relação teórica chamamos de **modelo empírico**.
- ▶ Escolhemos a forma mais adequada a partir de uma análise empírica dos dados.
- ▶ Essa forma não precisa ser necessariamente uma reta.
- ▶ Só iremos tratar aqui esse caso mais simples.

- ▶ Considere novamente o modelo

$$Y = \beta_0 + \beta_1 X + \epsilon.$$

- ▶ Consideramos que  $\epsilon \sim N(0, \sigma^2)$ .
- ▶ Se fixarmos  $x$  temos que

$$E(Y|x) = E(\beta_0 + \beta_1 x + \epsilon) =$$

- ▶ Considere novamente o modelo

$$Y = \beta_0 + \beta_1 X + \epsilon.$$

- ▶ Consideramos que  $\epsilon \sim N(0, \sigma^2)$ .
- ▶ Se fixarmos  $x$  temos que

$$E(Y|x) = E(\beta_0 + \beta_1 x + \epsilon) = E(\beta_0) + E(\beta_1 x) + E(\epsilon)$$

=

- ▶ Considere novamente o modelo

$$Y = \beta_0 + \beta_1 x + \epsilon.$$

- ▶ Consideramos que  $\epsilon \sim N(0, \sigma^2)$ .
- ▶ Se fixarmos  $x$  temos que

$$E(Y|x) = E(\beta_0 + \beta_1 x + \epsilon) = E(\beta_0) + E(\beta_1 x) + E(\epsilon)$$

$$= \beta_0 + \beta_1 x + 0 = \beta_0 + \beta_1 x$$

(como na definição anterior).

- ▶ A variância é dada por

$$\text{Var}(Y|x) = \text{Var}(\beta_0 + \beta_1 x + \epsilon) =$$

- ▶ Considere novamente o modelo

$$Y = \beta_0 + \beta_1 x + \epsilon.$$

- ▶ Consideramos que  $\epsilon \sim N(0, \sigma^2)$ .
- ▶ Se fixarmos  $x$  temos que

$$E(Y|x) = E(\beta_0 + \beta_1 x + \epsilon) = E(\beta_0) + E(\beta_1 x) + E(\epsilon)$$

$$= \beta_0 + \beta_1 x + 0 = \beta_0 + \beta_1 x$$

(como na definição anterior).

- ▶ A variância é dada por

$$\text{Var}(Y|x) = \text{Var}(\beta_0 + \beta_1 x + \epsilon) = \text{Var}(\beta_0) + \text{Var}(\beta_1 x) + \text{Var}(\epsilon)$$

=

- ▶ Considere novamente o modelo

$$Y = \beta_0 + \beta_1 x + \epsilon .$$

- ▶ Consideramos que  $\epsilon \sim N(0, \sigma^2)$ .
- ▶ Se fixarmos  $x$  temos que

$$E(Y|x) = E(\beta_0 + \beta_1 x + \epsilon) = E(\beta_0) + E(\beta_1 x) + E(\epsilon)$$

$$= \beta_0 + \beta_1 x + 0 = \beta_0 + \beta_1 x$$

(como na definição anterior).

- ▶ A variância é dada por

$$\text{Var}(Y|x) = \text{Var}(\beta_0 + \beta_1 x + \epsilon) = \text{Var}(\beta_0) + \text{Var}(\beta_1 x) + \text{Var}(\epsilon)$$

$$= 0 + 0 + \sigma^2 = \sigma^2 .$$

## Interpretação do $\beta_1$

- ▶ Temos que  $\beta_1$  é a inclinação da reta.
- ▶ Então  $\beta_1$  representa:
  - ▶ o aumento **esperado** em  $Y$  quando  $x$  aumenta uma unidade.
- ▶ Exemplo:  $Y$  (renda em mil reais) e  $x$  (escolaridade em anos)

$$Y = 0,5 + 1,5x + \epsilon.$$

- ▶ Espera-se um aumento de R\$ 1500,00 no salário para cada ano a mais de estudo.



## Interpretação do $\beta_0$

- ▶  $\beta_0$  é o intercepto da reta.
- ▶ Então  $\beta_0$  representa:
  - ▶ o valor **esperado** de  $Y$  quando  $x = 0$ .
- ▶ Exemplo:  $Y$  (renda em mil reais) e  $x$  (escolaridade em anos)

$$Y = 0,5 + 1,5x + \epsilon.$$

- ▶ A renda esperada de uma pessoa sem estudo algum é de R\$ 500,00.

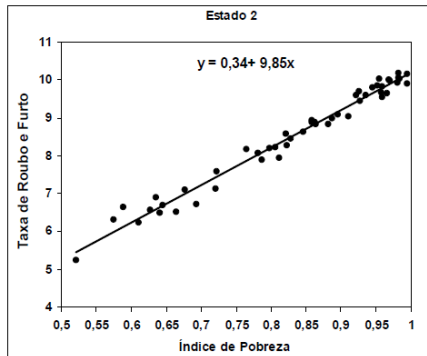
## Exemplo:

- ▶ Considere novamente o exemplo de associação entre pobreza e criminalidade.
- ▶ Vimos que existe uma forte associação entre as variáveis.
- ▶ Temos que  $\beta_0 = -0,7$ .
- ▶ Esse coeficiente não tem sentido prático.
- ▶ Não existe taxa negativa.
- ▶  $\beta_1 = 10,08$ .
- ▶ Interpretação: O aumento em uma unidade nos índices de pobreza aumenta em 10,08 o número de casos esperados de casos de roubos/furtos por mil habitantes.

## Exemplo: (continuação)

- ▶ Para o Estado 2 temos que a reta de regressão é dada por:

$$\text{Taxa Criminalidade} = 0,34 + 9,85\text{Pobreza}$$



**Exemplo: (continuação)**

- ▶ Nesse caso a reta é bastante informativa.
- ▶ Há pouca dispersão dos pontos em torno dela.
- ▶  $\beta_0$  não tem interpretação.
- ▶  $\beta_1 = 9,85$ .
- ▶ Qual interpretação?

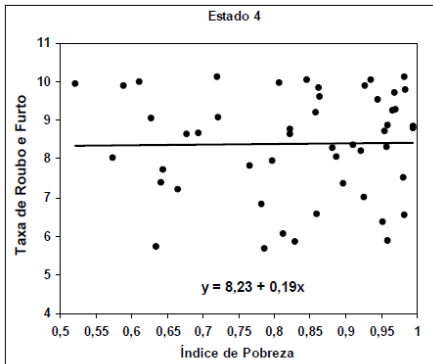
## Exemplo: (continuação)

- ▶ Nesse caso a reta é bastante informativa.
- ▶ Há pouca dispersão dos pontos em torno dela.
- ▶  $\beta_0$  não tem interpretação.
- ▶  $\beta_1 = 9,85$ .
- ▶ Qual interpretação?
- ▶ O aumento em uma unidade nos índices de pobreza aumenta em 9,85 o número de casos esperados de casos de roubos/furtos por mil habitantes.

**Exemplo: (continuação)**

- ▶ Para o Estado 4 temos que a reta de regressão é dada por:

$$\text{Taxa Criminalidade} = 8,23 + 0,19\text{Pobreza}$$



**Exemplo: (continuação)**

- ▶ A reta de regressão não é útil nesse caso.
- ▶ Não deve ser usada.
- ▶ Para esse estado, os níveis de pobreza não dizem nada sobre a criminalidade.
- ▶ Não podemos interpretar os coeficientes.

**Exemplo:**

- ▶ Considere o exemplo de destilação do oxigênio:
  - ▶  $y$  - pureza do oxigênio produzido em um processo de destilação;
  - ▶  $x$  - porcentagem de hidrocarbonetos presentes no condensador.
- ▶ Suponha que o verdadeiro modelo é dado por

$$Y = 75 + 15x + \epsilon$$

onde  $\epsilon \sim N(0, 2)$ .

- ▶ Então

$$Y \sim$$



**Exemplo:**

- ▶ Considere o exemplo de destilação do oxigênio:
  - ▶  $y$  - pureza do oxigênio produzido em um processo de destilação;
  - ▶  $x$  - porcentagem de hidrocarbonetos presentes no condensador.
- ▶ Suponha que o verdadeiro modelo é dado por

$$Y = 75 + 15x + \epsilon$$

onde  $\epsilon \sim N(0, 2)$ .

- ▶ Então

$$Y \sim N(75 + 15x, 2).$$

- ▶ A variância do efeito aleatório  $\sigma^2$ 
  - ▶ determina a variabilidade do  $Y$  em torno da reta.
- ▶ Se  $\sigma^2$  é grande  $\Rightarrow$

**Exemplo:**

- ▶ Considere o exemplo de destilação do oxigênio:
  - ▶  $y$  - pureza do oxigênio produzido em um processo de destilação;
  - ▶  $x$  - porcentagem de hidrocarbonetos presentes no condensador.
- ▶ Suponha que o verdadeiro modelo é dado por

$$Y = 75 + 15x + \epsilon$$

onde  $\epsilon \sim N(0, 2)$ .

- ▶ Então

$$Y \sim N(75 + 15x, 2).$$

- ▶ A variância do efeito aleatório  $\sigma^2$ 
  - ▶ determina a variabilidade do  $Y$  em torno da reta.
- ▶ Se  $\sigma^2$  é grande  $\Rightarrow$  as observações ficam longe da reta.
- ▶ Se  $\sigma^2$  é pequeno  $\Rightarrow$

**Exemplo:**

- ▶ Considere o exemplo de destilação do oxigênio:
  - ▶  $y$  - pureza do oxigênio produzido em um processo de destilação;
  - ▶  $x$  - porcentagem de hidrocarbonetos presentes no condensador.
- ▶ Suponha que o verdadeiro modelo é dado por

$$Y = 75 + 15x + \epsilon$$

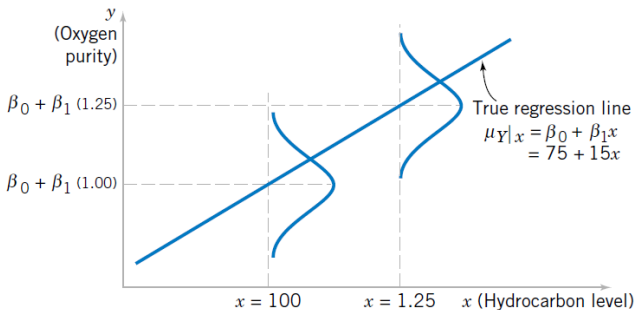
onde  $\epsilon \sim N(0, 2)$ .

- ▶ Então

$$Y \sim N(75 + 15x, 2).$$

- ▶ A variância do efeito aleatório  $\sigma^2$ 
  - ▶ determina a variabilidade do  $Y$  em torno da reta.
- ▶ Se  $\sigma^2$  é grande  $\Rightarrow$  as observações ficam longe da reta.
- ▶ Se  $\sigma^2$  é pequeno  $\Rightarrow$  as observações ficam perto da reta.

**Exemplo:** O modelo pode ser representado graficamente da seguinte forma:



**Figure 11-2** The distribution of  $Y$  for a given value of  $x$  for the oxygen purity-hydrocarbon data.

**Exemplo:**

- ▶ Podemos usar o modelo para respondermos :
  - ▶ qual a pureza esperada do oxigênio para uma determinada porcentagem de hidrocarbonetos?
- ▶ Considere que a porcentagem de hidrocarbonetos é 1,25% ( $x = 1,25$ ).
- ▶ Então a pureza esperada do oxigênio é de:

$$E(Y|x) = 75 + 15(1,25) = 93,75 .$$

- ▶ Esse é um exemplo hipotético.
- ▶ Geralmente não saberemos o valor real de  $(\beta_0, \beta_1)$  e  $\sigma^2$ .
- ▶ São estimados a partir de dados da amostra.
- ▶ Veremos a seguir o método mais usado.
- ▶ O método de mínimos quadrados.

## Abusos sobre a regressão

- ▶ Associação entre variáveis não implica relação causal.
- ▶ Planejamento de experimentos é a única forma de determinar relações causais.
- ▶ Relações de regressão são válidas apenas dentro da faixa dos dados coletados.
- ▶ Modelos de regressão podem não ser válidos para extrapolação.

## Uso da equação de regressão

- ▶ Podem ser úteis para predizer o valor de uma variável dado o valor de outra.
- ▶ Só podemos usá-la se o valor de  $r$  indica uma associação **linear** entre as variáveis.
- ▶ A reta de regressão precisa se ajustar bem aos dados.
- ▶ Se não existe uma correlação linear, a nossa melhor estimativa para  $Y$  é sua média.



- ▶ Se queremos prever o valor de  $Y$  usando  $x$ :
  - ▶ se não existe relação linear entre  $x$  e  $y$ , o melhor valor é a média;
  - ▶ se existe relação linear, substitui o valor de  $x$  na reta de regressão.
- ▶ Uma reta ajustada no passado pode não ser útil hoje.
- ▶ Não devemos fazer previsões para populações diferentes daquela de onde provem os dados amostrais.