

STA 302: Regression Analysis

1.1

Greeting, Syllabus, Rules

Statistics

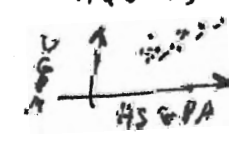
Goal: To obtain knowledge from noisy numerical data.

One approach: Study relationships between variables.

← Case: Sample of n cases (people, rats, cars, hospitals)

← Variable: Piece of information recorded for each case (age, sex, mpg, # nurses etc)

Data file: Rectangular array (Spreadsheet) $\begin{cases} \text{Rows} = \text{Cases} \\ \text{Cols} = \text{Variables} \end{cases}$

Scatter plot, least squares line 

Related = Not independent

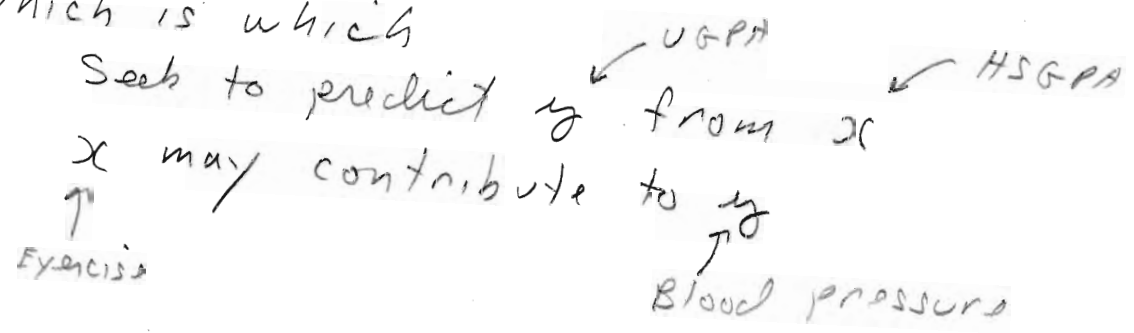
Independent: Conditional distribution of Y given $X=x$ does not depend on x

$$f(y|x) = f(y) \Leftrightarrow \frac{f(x,y)}{f(x)} = f(y)$$

$$\Leftrightarrow f(x,y) = f(x)f(y)$$

X = Independent (Explanatory) Var
 Y = Dependent (Response) Var

Which is which



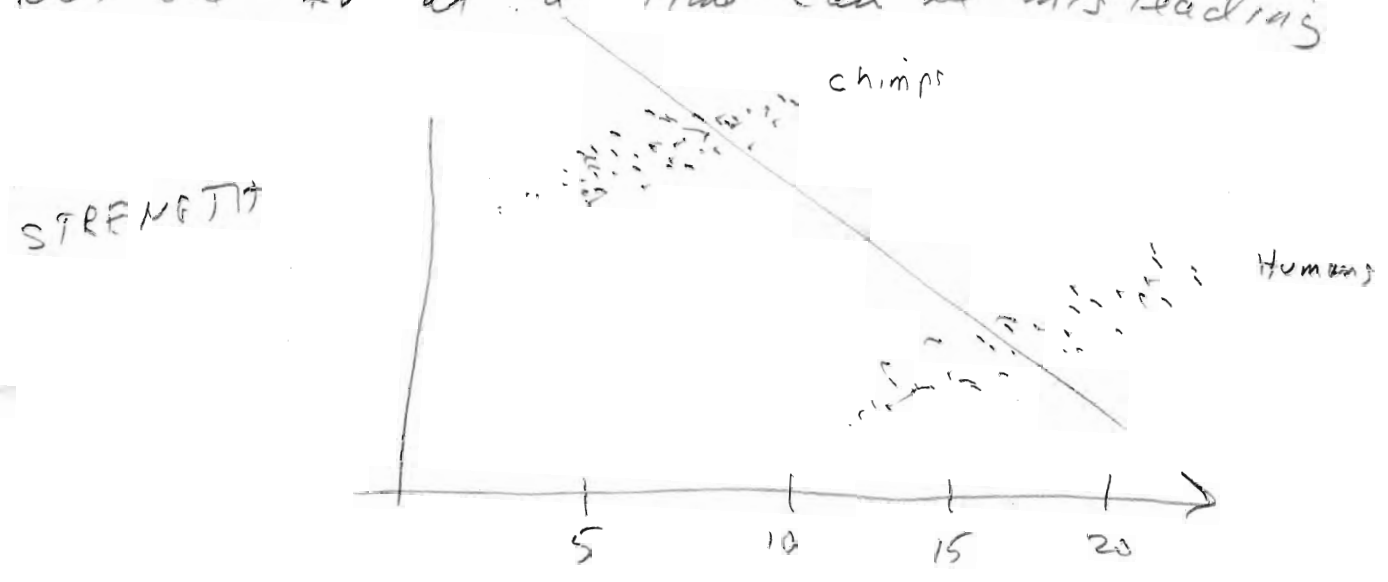
Regression (term based on a misunderstanding) is one way the distribution of Y might depend on X . Incl for $i = 1, \dots, n$

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i \quad E(\epsilon_i) = 0, V(\epsilon_i) = \sigma^2$$

$$E(Y_i) = \beta_0 + \beta_1 X_i, V(Y_i) = \sigma^2 \quad \text{or } \epsilon_i \sim N(0, \sigma^2)$$

A MODEL

But one IV at a time can be misleading



Multiple Regression

1.3

In scalar form

$$Y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + \epsilon_i$$

$\epsilon_1, \dots, \epsilon_n$ independent $E(\epsilon_i) = 0$, $V(\epsilon_i) = \sigma^2$
or $\epsilon_i \sim N(0, \sigma^2)$

In matrix form $Y = X\beta + \epsilon$

$$\begin{matrix} Y \\ \sim \\ \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} \end{matrix} = \begin{matrix} X \\ \begin{bmatrix} 1 & 14.2 & \dots & 1 \\ \vdots & 11.9 & & 0 \\ \vdots & 3.7 & & 0 \\ \vdots & \vdots & & \vdots \\ 1 & 6.75 & & 1 \end{bmatrix} \end{matrix} \begin{matrix} \beta \\ \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix} \end{matrix} + \begin{matrix} \epsilon \\ \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{bmatrix} \end{matrix}$$

$$\epsilon \sim MVN(0, \sigma^2 I_n)$$

Tells us we need

- Matrix Algebra
- Random vectors, esp multivariate normal
- Software to compute all this

Tests & CI all depend on the normal distribution

Lead to χ^2 , t , F

We'll do the distribution theory use change of vars $E(g(x)) = \int g(x) f(x) dx$ a lot.

Moment-generating functions

X is a RV

$$M_X(t) = E(e^{xt}) = \begin{cases} \sum_x e^{xt} f(x) \\ \int_{-\infty}^{\infty} e^{xt} f(x) dx \end{cases}$$

Correspond uniquely to probability distributions. we will use (without proof) MGF of

$$X \sim N(\mu, \sigma^2) \quad M_X(t) = e^{\mu t + \frac{1}{2} \sigma^2 t^2}$$

$$X \sim \chi^2(\nu) \quad M_X(t) = (1 - 2t)^{-\nu/2}$$

No derivation required here

$$e^{\mu t + \frac{1}{2} \sigma^2 t^2}$$

$$(1 - 2t)^{-1/2}$$

(15)

MGFs correspond uniquely to probability distributions (no proof: Advanced)

Ex $X \sim N(\mu, \sigma^2)$, $Z = \frac{X - \mu}{\sigma}$

$$\begin{aligned} M_Z(t) &= E e^{\frac{X - \mu}{\sigma} t} = E \left(e^{X(t/\sigma)} e^{-\frac{\mu t}{\sigma}} \right) \\ &= e^{-\frac{\mu t}{\sigma}} M_X(t/\sigma) = e^{-\frac{\mu t}{\sigma}} e^{\mu(t/\sigma) + \frac{1}{2} \sigma^2 (t/\sigma)^2} \\ &= e^{-\frac{\mu t}{\sigma}} e^{\frac{\mu t}{\sigma}} e^{\frac{1}{2} t^2} = e^{\frac{1}{2} t^2} = e^{0t + \frac{1}{2} 1 t^2} \quad N(0, 1) \end{aligned}$$

Properties (well known)

$$M_{ax}(t) = M_x(at)$$

$$\begin{aligned} M_{(x+c)}(t) &= e^{ct} M_x(t) & E(e^{(x+c)t}) &= E(e^{xt+ct}) \\ & & &= E(e^{xt} e^{ct}) \\ & & &= e^{ct} E(e^{xt}) \end{aligned}$$

If X_1, X_2 are ind.

$$M_{X_1 + X_2}(t) = M_{X_1}(t) M_{X_2}(t) \quad \text{same } t$$

Proof

$M_{X_1 + X_2}$ (Next page)

$$M_{X_1 + X_2}(t) = \iint e^{(x_1 + x_2)t} f_{X_1, X_2}(x_1, x_2) dx_1 dx_2$$

$$= \iint e^{x_1 t} e^{x_2 t} f_{X_1, X_2}(x_1, x_2) dx_1 dx_2$$

ind
↓
=

$$\iint e^{x_1 t} e^{x_2 t} f_{X_1}(x_1) f_{X_2}(x_2) dx_1 dx_2$$

$$\int e^{x_2 t} f_{X_2}(x_2) \left[\int e^{x_1 t} f_{X_1}(x_1) dx_1 \right] dx_2$$

$M_{X_1}(t)$

$$= M_{X_1}(t) M_{X_2}(t)$$

~~And by the uniqueness of MGFs, this is an if and only if. That is X_1 & X_2 are independent if and only if M_{X_1}~~

Extends to X_1, \dots, X_n ind, then

$$M_{\sum X_i}(t) = \prod_{i=1}^n M_{X_i}(t)$$

So, let $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$,

$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$. Find dist of \bar{X}_n .

$$\begin{aligned}
M_{\bar{X}}(t) &= M_{\frac{1}{n} \sum X_i}(t) = M_{\sum X_i}(t/n) \\
&= \prod_{i=1}^n M_{X_i}(t/n) = \prod_{i=1}^n e^{\mu t/n + \frac{1}{2} \sigma^2 (t/n)^2} \\
&= e^{\sum_{i=1}^n (\mu t/n) + \sum_{i=1}^n \frac{1}{2} \sigma^2 t^2/n^2} \\
&= e^{n(\mu t/n) + n(\frac{1}{2} \sigma^2 t^2/n^2)} \\
&= e^{\mu t + \frac{1}{2} \frac{\sigma^2}{n} t^2}
\end{aligned}$$

MGF of $N(\mu, \frac{\sigma^2}{n})$

So $\bar{X}_n \sim N(\mu, \frac{\sigma^2}{n})$

$X_1, \dots, X_n \stackrel{iid}{\sim} \chi^2(r_i)$ $Y = \sum X_i$

$$M_Y(t) = \prod_{i=1}^n M_{X_i}(t) = \prod_{i=1}^n (1-2t)^{-r_i/2}$$

$$= (1-2t)^{-\frac{\sum r_i}{2}} \quad \chi^2(\sum r_i)$$

~~That are some other problems~~

A KEY TOOL

for deriving the F
and t distributions

Let $X_1 \neq X_2$ be independent, $Y = X_1 + X_2$

$$X_1 \sim \chi^2(\nu_1), \quad Y \sim \chi^2(\nu_1 + \nu_2) \quad \begin{matrix} \nu_1 > 0 \\ \nu_2 > 0 \end{matrix}$$

Then $X_2 \sim \chi^2(\nu_2)$

Proof By independence, $M_Y(t) = M_{X_1}(t)M_{X_2}(t)$

$$\Rightarrow (1-2t)^{-\frac{\nu_1 + \nu_2}{2}} = (1-2t)^{-\frac{\nu_1}{2}} M_{X_2}(t)$$

multiply both sides by $(1-2t)^{\nu_1/2}$

$$\Rightarrow (1-2t)^{-\nu_2/2} = M_{X_2}(t)$$

MGF of $\chi^2(\nu_2)$, done

One more ...

1.9

$Z \sim N(0, 1)$, $Y = Z^2$. Show $Y \sim \chi^2(1)$

$$M_Y(t) = M_{Z^2}(t) = \int_{-\infty}^{\infty} e^{t z^2} \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}} dz$$

$$= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-\frac{1}{2}(z^2 - 2tz^2)} dz$$

$$= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-\frac{(1-2t)}{2} z^2} dz$$

← as $\frac{1}{2}$

$$= \frac{1}{(1-2t)^{1/2}} \int_{-\infty}^{\infty} \frac{(1-2t)^{1/2}}{\sqrt{2\pi}} e^{-\frac{(1-2t)}{2} z^2} dz$$

$$= (1-2t)^{-1/2} \quad \text{MGF of } \chi^2(1)$$

There are some homework problems

THIS will save technical headaches later (1.10)

Because MGFs correspond uniquely to probability distributions

(densities do not), distributions can be defined in terms of their MGFs.

So we will define $X \sim N(\mu, \sigma^2)$ as meaning $M_X(t) = e^{\mu t + \frac{1}{2}\sigma^2 t^2}$

where $\sigma^2 \geq 0$.

Now observe. Let Y be a degenerate RV with $P(Y = \mu) = 1$,

$$\begin{aligned} \text{So } M_Y(t) &= E(e^{tY}) = \sum_{\omega: P_Y(\omega) > 0} e^{t\omega} P_Y(\omega) \\ &= e^{\mu t} \cdot 1 = e^{\mu t} \end{aligned}$$

MGF of $N(\mu, 0)$

So in this sense, degenerate RVs are normal

Think of $\lim_{\sigma^2 \downarrow 0} f_X(x)$, $X \sim N(\mu, \sigma^2)$