

Historia Rachunku Prawdopodobieństwa i  
Statystyki  
WYKŁAD V: Quetelet. Legendre i Gauss. CTG  
Laplace'a

MiNI PW

# Adolf Quetelet (1796-1874)

Adolf Qutelet: dobrze wykształcony matematyk, pierwszy doktorat z nauk ścisłych Uniwersytetu w Gandawie.

Członek Królewskiej Akademii Nauk i Literatury od 1820 roku.

Założyciel Królewskiego Obserwatorium w Brukseli (1828).

Bibliografia jego prac zawiera ponad 300 pozycji.

Główne dzieło *O człowieku i rozwoju jego cech. Rozprawa o fizyce społecznej* (dwa tomy opublikowane w 1835 roku). Drugie wydanie *Rozprawa o fizyce społecznej* (1869)

Fizyka społeczna: analiza ilościowa danych dotyczących społeczeństw.

# Adolf Quetelet (1796-1874), statystyk, socjolog, astronom



# Adolf Quetelet (1796-1874)

## Główne osiągnięcia:

- ▶ Próba wprowadzenia metody ilorazowej Laplace'a i dyskusja metodologiczna z tym związana;
- ▶ koncepcja człowieka przeciętnego (*l'homme moyen*)
- ▶ wkład do analizy danych (ilościowa analiza danych społecznych)
- ▶ dopasowanie rozkładu normalnego do danych
- ▶ prawo przyczyn przypadkowych

## Metoda ilorazowa Laplace'a

Laplace: metoda oceny wielkości populacji na podstawie przeskalowania danych z jednego departamentu.

$$\frac{N}{N_D} = \frac{U}{U_D},$$

gdzie:

$N$ ,  $N_D$  - liczby mieszkańców kraju i ustalonego departamentu odpowiednio;

$U$ ,  $U_D$  - liczby urodzin w kraju i departamencie odpowiednio.

Zatem

$$N = N_D \times \frac{U}{U_D},$$

czyli znając liczbę urodzin w departamencie i kraju oraz liczbę mieszkańców w departamencie byliśmy w stanie oszacować liczbę  $N$ .

Podobnie

$$N = N_D \times \frac{Z}{Z_D},$$

gdzie  $Z$  i  $Z_D$  odpowiednie liczby zgonów.

Krytyka de Kevenberga: Zmienność liczby zgonów i urodzin w departamentach jest bardzo duża i zależy od wielu czynników (gęstość zaludnienia, wysokość nad poziomem morza..) więc dane z jednego departamentu nie są *reprezentatywne*. Można wprowadzić zmodyfikować wzór

$$N = \frac{1}{K} \sum_{k=1}^K N_{D,k} \times \frac{Z}{Z_{D,k}},$$

ale wdg Kevenberga prowadziłyby to do wybrania tylu departamentów, że metoda straciłaby swoją użyteczność w porównaniu z pełnym spisem.

AQ zgodził się z Kevenbergiem i odstąpił od propagowania metody ilorazowej.

Kwestia złych metodologii i błędnych badań ...

# Koncepcja człowieka przeciętnego (*l'homme moyen*)

Przeciętność dla AQ oznaczała bycie blisko wartości średniej charakterystyki w rozpatrywanej populacji.

Koncepcja człowieka przeciętnego ewoluowała w czasie.

# Koncepcja człowieka przeciętnego (*l'homme moyen*)

Koncepcja ilościowego spojrzenia na społeczeństwo i prawa społeczne rozwijana w tamtym czasie.

Jedna z nagród Królewskiego Towarzystwa w Brukseli za rok 1828:  
*Przedstawić teorię matematyczną ludzi i zwierząt traktowanych jako motory i maszyny.*

Zaproponował koncepcję człowieka przeciętnego jako posiadacza średnich wartości cech policzonych dla pewnej społeczności. Miało to dotyczyć nie tylko cech antropometrycznych ale również kategorii moralnych i obyczajowych.

*If one individual at any given epoch of society possessed all the qualities of the average man, he would represent all that is great, good or beautiful*



Koncepcja człowieka przeciętnego zyskała w tamtych czasach dużą popularność, podobnie jak teraz.

Jednak:

Recenzenci mieli wątpliwości co do rozciągnięcia idei 'normalności' na sferę moralną:

*Jeśli chodzi o sferę moralności taka przeciętność byłaby nie do zaakceptowania*

A. Cournot: Idea określania typowości poprzez liczenie średniej nie ma sensu, gdyż obiekt średni nie musi być elementem badanej klasy.

Argument: Trójkąt mający boki będące średnimi boków trójkątów prostokątnych nie musi być trójkątem prostokątnym.

# Koncepcja normalności

Sam AQ widział, że żeby porównywać wartość ze średnią, konieczna jest miara zmienności wokół średniej.

Wprowadził, na podstawie odstępstwa od średniej, wskaźnik względnej skłonności do popełnienia przestępstwa w zależności od wieku ...

Koncepcja normalności: jeśli odchylenie jakiejś jednostki od średniej było nieistotne, AQ traktował ją jako 'normalną'...

Co jest normalne ?

American Psychiatric Association diagnostic and statistical manual

DSM-V: znaczne rozszerzenie liczby zaburzeń traktowanych jako choroby

..

Problem wielokrotnego testowania w przypadku wielu cech !!!

# Wkład do analizy danych (ilościowa analiza danych społecznych)

AQ był zafascynowany metodami syntetycznego spojrzenia na dane i możliwościami wykresów.

Nie wyszedł ponad reprezentację dwóch zmiennych jednocześnie, ale pomijając to jego ciekawość była nieograniczona:

zgony i urodziny versus miesiąc, miasto, temperatura, pora dnia;

zgony versus wiek, zawód, miejsce zamieszkania;

zgony w więzieniach i szpitalach;

statystyki dotyczące pijaństwa, przestępczości i chorób psychicznych.

Zadziwiając regularności: liczba małżeństw

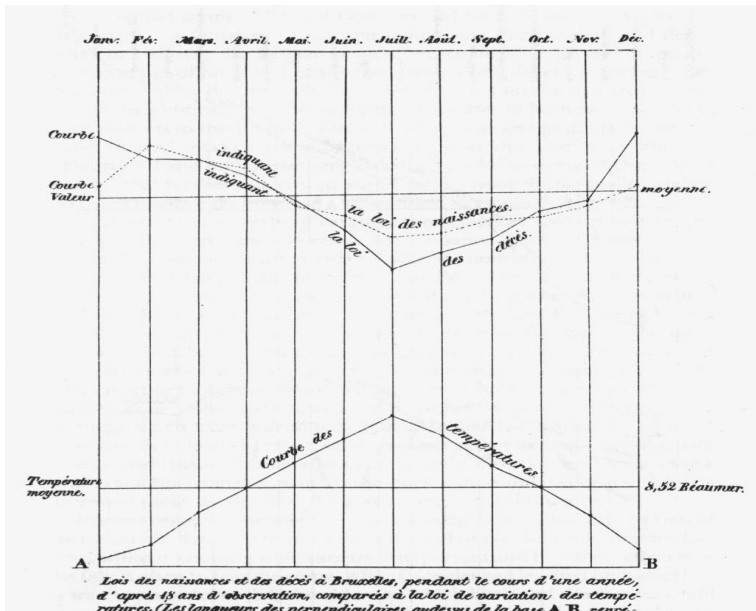
(kobieta  $\geq 60$ , mężczyzna  $\leq 30$  w Belgii w latach 1841 - 1845):

7, 6, 8, 5, 5

odpowiednio.

To doprowadziło go do prawa przyczyn przypadkowych.

# Reprezentacje graficzne: narodziny i zgony vs czas



# Reprezentacje graficzne Minard

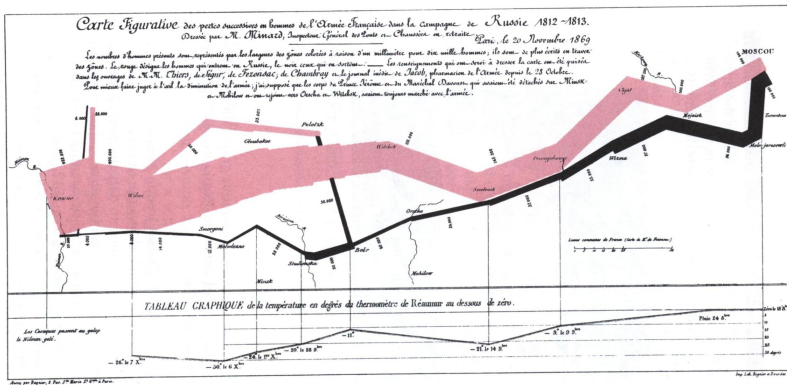


Figure 3.12 Napoleon's march to and retreat from Moscow. (Source: Edward R. Tufte, The Visual Display of Quantitative Information, Cheshire, CT: Graphic Press, 1983.)

Autor: Karol Minard

# Berezyna, listopad 1812



# Modelowanie zależności

Próba modelowania związku między wagą i wzrostem

$$y + \frac{y}{1000(T - y)} = ax + \frac{t + x}{1 + (4/3)x}$$

$t, T$  - wzrost w momencie urodzin i zgonu;

$x$ - waga,  $y$ - wzrost,  $a$  - stała zależna od populacji.

# Stabilność charakterystyk populacji

Teza AQ:

Jeśli związki przyczynowe nie ulegają zmianie to średnie charakterystyki w populacji (nawet niejednorodnej) się nie zmieniają.

Wprowadził tzw indeks Queteleta otyłości (od 1972 roku zwany BMI).



# Prawo przyczyn przypadkowych

Ogólne prawo rządzące wszechświatem: prawo przyczyn przypadkowych (*loi des causes accidentales*).

Zależności są wynikiem działania dwóch sił: sił stałych (*naturelles*) i sił zakłócających (*perturbatrices*). Siły stałe określają stan przeciętny (*etat moyen*), a siły zakłócające (losowe) odpowiedzialne są za odstępstwa. W przypadku wzrostu: stan przeciętny: wzrost średni, siły zakłócające: klimat, odżywianie itp.

W przypadku zjawisk społecznych siła zakłócającą może być wolna wola. Wpływ wolnej woli na przebieg zjawisk społecznych ogranicza się do roli przyczyny przypadkowej i jej wypadkowa maleje w przypadku dużych populacji.

Buckle *History of Civilisation in England* 1857

*Quteletismus*

# Dopasowanie rozkładu normalnego do danych: Quetelet

Quetelet zrozumiał znaczenie CTG Laplace'a dla analizy danych społecznych. Przykład analizy: obwód klatki piersiowej 5738 żołnierzy szkockich (najczęściej używane dane w XIX wieku) (moda 40 cali=102 cm).

AQ wykorzystał swoje, sprytnie policzone, tablice rozkładu normalnego oparte na losowaniu 999 kul z 'dużej' urny zawierającej jednakową liczbę kul białych i czarnych (de facto  $\text{Bin}(999, 1/2)$ ).

Metoda opierała się na policzeniu skumulowanych prawdopodobieństw normalnych

$$S(r) = P(500 \leq X < 500 + r)$$

policzeniu analogicznego skumulowanego prawdopodobieństwa  $\hat{p}$  i

$$\hat{r} = S^{-1}(\hat{p})$$

Natępnie interpolacja liniowa  $\hat{r} \rightarrow \tilde{r}$  i wyliczenie odpowiednich prawdopodobieństw obwodu klatki na podstawie  $S(\tilde{r})$ .

# Dane: obwody klatki piersiowej

MESSURES de la poitrine.	NOMBRE d'individus.	NOMBRE proportionnel.	PROBABILITÉ d'être l'un d'eux.	RANG dans la TABLE.	RANG d'être le cas.	PROBABILITÉ d'être la TABLE.	NOMBRE d'individus total.
Pouces.							
55	5	5	0,5000			0,5000	7
54	18	31	0,4995	52	50	0,4995	29
53	81	141	0,4964	42,5	42,5	0,4964	119
52	185	322	0,4883	33,5	34,5	0,4854	325
51	420	732	0,4591	20,0	20,5	0,4551	733
50	749	1505	0,3760	18,0	18,5	0,3799	1533
49	1075	1807	0,2464	10,5	10,5	0,2465	1838
			0,0387	2,5	2,5	0,0488	
40	1679	1882	0,1285	3,5	3,5	0,1330	1987
41	854	1628	0,2913	13	13,5	0,2954	1675
42	638	1148	0,4061	91	91,5	0,4130	1096
45	370	645	0,4706	30	29,5	0,4690	560
44	92	160	0,4916	35	37,5	0,4911	221
45	50	87	0,4053	41	43,5	0,4080	69
46	21	38	0,4091	46,5	51,5	0,4090	16
47	4	7	0,4098	56	61,8	0,4099	3
48	1	2	0,5000			0,5000	1
	5738	1,0000					1,0000

Figure 5.3. Quetelet's analysis fitting a normal distribution to data on the chest circumferences of Scottish soldiers. Column 1 gives the chest circumference in inches; columns 2 and 3 give the frequency and relative frequency distributions for 5738 individuals, columns 4–7 give the details of Quetelet's calculations (see text), and column 8 gives the fitted relative frequency distribution. (From Quetelet, 1846, p. 400.)

# Henry Buckle 'Historia cywilizacji w Anglii'



# Henry Buckle 'Historia cywilizacji w Anglii'

Henry Buckle (1821-1862) : historia rządzi się materialistycznymi prawami, wydarzenia następują w konsekwencji wydarzeń poprzednich, a nie w wyniku arbitralnych decyzji jednostek czy działania sił nadprzyrodzonych. Klimat i dostępność żywności są podstawowymi czynnikami wpływającymi na rozwój cywilizacji.

Istnieje 'konieczność dziejowa', ale nie należy jej przyspieszać (inaczej sądził Marks)

Koncepcja społeczeństwa jako zbioru dużej liczby jednostek rządzącego się pewnymi ogólnymi prawami propagowana przez Queteleta i Buckle'a została wykorzystana przez fizyków do stworzenia podstaw **fizyki statystycznej**.

# Główne postaci

- ▶ Rudolph Clausius (1822-1866)
- ▶ James Maxwell (1831-1879)
- ▶ Ludwig Boltzmann (1844 -1906)
- ▶ Marian Smoluchowski
- ▶ Josiah Gibbs (1839-1903)



# Legendre i metoda najmniejszych kwadratów

Adrien Marie Legendre (1752-1833)

Pierwsza wzmianka o MNK w pracy o wyliczaniu orbit komet.

Później wykorzystana przez niego przy liczeniu długości południka (od równika do bieguna północnego, *meridian quadrant*) przechodzącego przez Paryż (1805).

Na niej był oparty wzorec metra (metr =  $10^{-7}$  długości). Dopasowanie równania parametrycznego do danych empirycznych: szukamy minimum funkcji kryterialnej.

MNK traktowana jako metoda analizodanowa, bez modelu probabilistycznego.

W szczególności zdawał sobie sprawę, że reguła brania średniej obserwacji to przypadek szczególny metody MNK.

## Legendre i metoda najmniejszych kwadratów cd.

Spór Legendre'a z Gaussem o odkrycie MNK.

AMK był niewątpliwie pierwszy, który tę metodę wykorzystał w opublikowanej pracy (1805), Gauss miał zwyczaj pisania do szuflady ... Podobnie jak I. Newton, który toczył analogiczny spór z G. Leibnizem o odkrycie rachunku różniczkowego i całkowego.

Pierwszy z publikacją był Leibniz *Nova methodus pro maximis et minimis*. GL w sporze odwołał się do Royal Society (której prezesem w tym czasie był Newton ..).

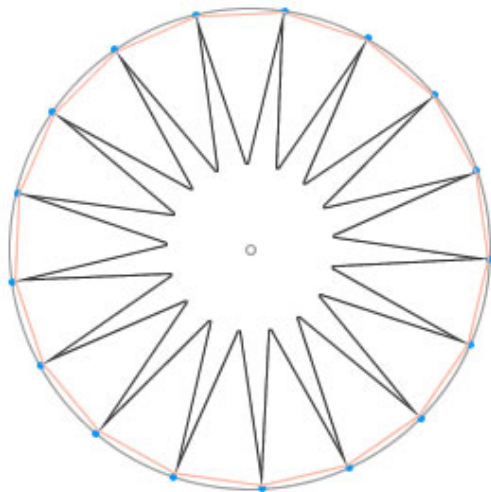


# Carl Friedrich Gauss



# Carl Friedrich Gauss

Konstrukcja geometryczna 17-kąta foremnego:



Gauss' Regular "17-gon"

# Carl Friedrich Gauss

W 1809 roku Gauss publikuje 'Teorię ruchu ciał niebieskich wokół słońca'. Pojawia się w niej MNK oparta na modelu probabilistycznym. Co zrobił Gauss ?

$$V = ap + bq + cr + ds + etc.$$

$$V' = a'p + b'q + c'r + d's + etc.$$

$$V = a''p + b''q + c''r + d''s + etc.$$

$V, V', V''$  - wartości oczekiwane odpowiedzi,  $a, b, c$  - zmienne niezależne,  $p, q, r..$  -współczynniki.

Równania strukturalne dopasowywane do danych.  $M, M', M''$  - obserwowane wartości zmiennej zależnej ( $EM = V$  etc).

$$\Delta = V - M, \Delta' = V' - M', \Delta'' = V'' - M'' \text{ etc.}$$

błędy pomiaru

Prawdopodobieństwa błędów pomiaru

$$\phi(\Delta), \phi(\Delta'), \phi(\Delta'') \text{ etc.}$$

Wartości błędów znajdziemy maksymalizując

$$\Omega = \phi(\Delta) \times \phi(\Delta') \times \phi(\Delta'') \text{ etc.}$$

O  $\phi$  przyjęt postulatory:

- ▶  $\operatorname{argmax}\phi(s) = 0$ ;
- ▶  $\phi(s) = \phi(-s)$ ;
- ▶ Gdyby  $V = V' = V'' = \dots = p$  to wartośćia maksymalizującą  $\Omega$  byłoby  $\hat{p} = k^{-1}(M + M' + M'' + \text{etc})$ , gdzie  $k$  jest liczbą obserwacji.

Końcówść 3-ciego postulatu: zakładamy coś, co mamy pokazać !

Przy tych założeniach Gauss pokazał, że

$$\phi(\Delta) = \frac{h}{\sqrt{\pi}} e^{-h^2 \Delta^2}$$

i że przyjęcie takiej krzywej błędów prowadzi do estymacji metodą najmniejszych kwadratów.

Estymator NW dla modelu regresji liniowej z błędami normalnymi jest estymatorem MNK.

$$\mathcal{L} = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(Y_i - x_i'\beta)^2}{2\sigma^2}\right) = \frac{1}{\sqrt{2\pi\sigma^2}^n} \exp\left(-\frac{\sum_{i=1}^n (Y_i - x_i'\beta)^2}{2\sigma^2}\right)$$

# Postać estymatora MNK metodą Gaussa

FG: inne rozwiązanie nie wykorzystujące równań normalnych. Napisał model liniowy

$$y = X\beta + \varepsilon$$

w postaci

$$X'X\beta = X'y - X'\varepsilon = X'y + z,$$

gdzie  $z = -X'\varepsilon$ .

FG szukał estymatora  $b$  parametru  $\beta$  w postaci

$$\beta = b + Qz.$$

Zatem

$$QX'X\beta = QX'y + Qz = QX'y + \beta - b$$

ma zachodzić dla dowolnego  $\beta$  zatem

$$QX'X = I_p \quad Q = (X'X)^{-1} \quad b = QX'y$$

# Analiza rezyduów

FG udowodnił i stwierdził, że jest to ciekawe, iż

$$\text{Var}(\varepsilon'\varepsilon/n) = 2\sigma^4/n$$

$$\text{Var}(e'e/n) = 2\sigma^4/(n-p),$$

gdzie  $p$  jest liczbą predyktorów. Wniosek:  $e'e$  może być traktowana jako suma kwadratów  $n-p$  niezależnych zmiennych losowych.

Komentarz: niestety nie zauważył, że

$$\varepsilon'\varepsilon = e'e + (b - \beta)'(X'X)(b - \beta)$$

i że oba składniki są niezależne, bo to by dało mu postać testu  $F$  odkrytego przez Fishera 100 lat później ( $H_0 : \beta = 0$ . Przy  $H_0$

$$\frac{p^{-1}b'(X'X)b}{(n-p)^{-1}e'e} \sim F_{p,n-p}$$



# CTG Laplace'a

Równoległe do prac nad podejściem bayesowskim (inverse probability)  
Laplace zajmował się CTG.

1776: CTG dla zmiennych losowych o rozkładzie prostokątnym

1781: dla zmiennych losowych o gęstości kawałkami ciągłej.

1810: przypadek ogólny. ....

Laplace po zobaczeniu wyniku Gaussa zrozumiał, że w oparciu o CTG ma znacznie lepsze uzasadnienie krzywej błędów, czego przedtem nie widział. Dodał to jako uzupełnienie do CTG w swojej pracy z 1810. Spór o pierwszeństwo między Legendre'm i Gaussem. Legendre mówił o MNK jako o naszej (tj Gaussa i jego) metodzie, Gauss twierdził, że używał tej metody od 1795 roku i mówił swoim współpracownikom o niej przed 1805 rokiem.

# Twierdzenie Gaussa-Markowa : Laplace

L. rozpatrywał sytuację

$$\varepsilon^{(i)} = p^{(i)}z - \alpha^{(i)}, \quad i = 1, \dots, n$$

$\varepsilon^{(i)}$  błędy aproksymacji  $y_i = \alpha^{(i)}$  funkcją liniową  $p^{(i)}z$  ( $p^{(i)} = x_i$ ).

Rozumowanie było oparte na prostej obserwacji: gdyby dla pewnych wag  $m^{(i)}$  mielibyśmy  $\sum m^{(i)}\varepsilon^{(i)} = 0$  wtedy

$$\sum m^{(i)}\varepsilon^{(i)} = \sum m^{(i)}p^{(i)}z - \sum m^{(i)}\alpha^{(i)}$$

i

$$z = \frac{\sum m^{(i)}\alpha^{(i)}}{\sum m^{(i)}p^{(i)}}$$

Oczywiście z reguły  $\sum m^{(i)}\varepsilon^{(i)} \neq 0$ . Wtedy modyfikacja rozumowania daje

$$z = \frac{\sum m^{(i)}\alpha^{(i)}}{\sum m^{(i)}p^{(i)}} + u,$$

gdzie

$$u = \frac{\sum m^{(i)}\varepsilon^{(i)}}{\sum m^{(i)}p^{(i)}}.$$

Uogólniając swoje CTG pokazał w 1811 roku, że  $u$  ma rozkład asymptotycznie normalny oraz estymator minimalizujący oczekiwany błąd to

$$z = \frac{\sum p^{(i)}\alpha^{(i)}}{\sum p^{(i)2}},$$

czyli estymator najmniejszych kwadratów.

W 12 lat później Gauss udowodnił to samo dla błędu średniokwadratowego i zauważył, że wystarczy tylko skończony drugi moment błędów. Twierdzenie Gaussa-Markowa.