

ETAPY BUDOWY MODELU EKONOMETRYCZNEGO

1. Sformułowanie modelu

a. wybór zmiennych: y, x_1, x_2, \dots

b. wybór postaci matematycznej modelu: liniowa, potęgowa,...

2. Zebranie danych statystycznych (różne źródła)

3. Selekcja zmiennych objaśniających

4. Estymacja parametrów modelu:

a. parametrów strukturalnych: a_0, a_1, a_2, \dots

b. parametrów stochastycznych: $s(a_i), s(y), R^2, R$

5. Weryfikacja modelu

MODEL BEZ WERYFIKACJI NIE MA ŻADNEJ WARTOŚCI

NIE NALEŻY KORZYSTAĆ Z PROGRAMÓW KOMPUTEROWYCH NIE DAJĄCYCH MOŻLIWOŚCI WERYFIKACJI

6. Interpretacja modelu

ETAP 1a. WYBÓR ZMIENNYCH

- zmienna objaśniana Y:
- zmienne objaśniające X_i (jak najwięcej dla modelu przyczynowo-skutkowego) z następujących źródeł (w kolejności):
 - teoria danej dziedziny wiedzy
 - doświadczenie zleceniodawcy i statystyka
 - metoda prób i błędów (intuicyjnie)
- wybrane zmienne muszą mieć dużą zmienność ($V > 30\%$)
- najczystszy błąd – „masa małania” prowadzące do związku funkcyjnego i nie dające żadnej informacji o zmiennej objaśnianej

ETAP 1b. WYBÓR POSTACI MATEMATYCZNEJ

- modele przyczynowo-skutkowe – najbardziej zalecane jest równoczesne prowadzenie obliczeń dla dwu postaci:
 - liniowej
$$y = \sum a_i x_i + \xi$$
 - potęgowej
$$y = \prod x_i^{a_i} \varepsilon \quad \ln y = \sum a_i \ln x_i + \xi$$
 - stosuje się te modele nieliniowe o narzuconej postaci nieliniowej, których parametry ustala się przez programowanie liniowe lub innymi metodami
- modele tendencji rozwojowej:
 - funkcja liniowa
 - proste funkcje nieliniowe
 - wielomiany
 - modele kombinowane: trend + wahania okresowe

ETAP 2. GROMADZENIE DANYCH STATYSTYCZNYCH

- rodzaje danych: dane przekrojowe i szeregi czasowe
- Źródła danych: roczniki statystyczne, różne działy przedsiębiorstwa, badania marketingowe, wywiady itd.
- wiarygodność danych: do jakiego celu zostały one przygotowane?
- porównywalność danych: inflacja (ceny bieżące a ceny stałe), zmiany procesów technicznych
- zmienność zjawisk: trzeba sprawdzić czy wybrana w etapie 1a zmienna jest rzeczywiście zmienną losową

$$V_x = \frac{s(x)}{\bar{x}} 100\% \quad V_x \text{ musi wynosić co najmniej 30-40\%}$$

ETAP 3. SELEKCJA ZMIENNYCH OBJAŚNIAJĄCYCH

KADK ZMIENNA X WYTYPOWANA W ETAPIE 1a TRAKTUJEMY JAKO KANDYDATKA NA ZMIENNA OBJAŚNIAJĄCĄ

- w modelu nie może być zbyt wielu zmiennych (nieczytelny)
- kandydatka może nie mieć wpływu na zmienną Y
- kandydatka może wnosić prawie tę samą informację o Y co inna kandydatka
- dwie bardzo podobne kandydatki mogą sobie nawzajem przeszkadzać (efekt katalityczny)

Kryteria, jakie musi spełniać kandydatka X_i , aby nadawać się do modelu:

- musi być silnie powiązana ze zmienną Y
- nie może być powiązana z inną kandydatką X_j

Metody selekcji zmiennych objaŝniajŝcych:

- badanie istotnoŝci korelacji
- grafowa
- Hellwiga (pojemnoŝci informacji)
- taksonomiczne (clustering)

$$r_{x_i x_j}$$

Przykŝad Macierz powiŝzania zmiennych ze sobŝ przedstawia tabela ($n=20$). Jak jŝmoŝna zinterpretowaŝ? Ktŝre zmienne sŝ powiŝzane ze sobŝ **w sposób istotny?**

	Y	X ₁	X ₂	X ₃
Y		0,52	0,64	-0,21
X ₁	0,52		0,82	-0,18
X ₂	0,64	0,82		0,08
X ₃	-0,21	-0,18	0,08	

Macierz współczynników korelacji

• **Testowanie istotnoŝci współczynnika korelacji**

$H_0: \rho = 0$

$H_1: \rho \neq 0$

Moŝna przeprowadziŝ testem **Studenta (t)**;

Wallace'a-Snedecora (R)

TEST Wallece'a-Snedecora

Fragment tablicy rozkładu Wallece'a-Snedecora

Stopnie swobody	Liczba zmiennych					
	2		3		4	
	0,05	0,01	0,05	0,01	0,05	0,01
8	0,632	0,765	0,726	0,827	0,777	0,860
18	0,444	0,561	0,532	0,633	0,587	0,678
28	0,361	0,463	0,439	0,530	0,490	0,573

Reguła decyzyjna: — jeżeli $|r_{obl}| > R_{tabl}$, odrzucamy H_0 (korelacja istotna)
 — jeżeli $|r_{obl}| < R_{tabl}$, przyjmujemy H_0 (brak korelacji)

W przykładzie, jeżeli przyjmujemy $\alpha = 0,05$, to $R_{tabl} = 0,444$.

	Y	X ₁	X ₂	X ₃
Y		0,52	0,64	-0,21
X ₁	0,52		0,82	-0,18
X ₂	0,64	0,82		0,08
X ₃	-0,21	-0,18	0,08	

Pozostanę zatem tylko trzy istotne powiązania: y-x₁, y-x₂, x₁-x₂

Test Studenta (t)

$H_0: \rho = 0$

$H_1: \rho \neq 0$

$$r_{kr} = \sqrt{\frac{t_\alpha^2}{n-2 + t_\alpha^2}} \quad r_{x_i x_j}, \quad r_{y x_i}$$

Reguła decyzyjna: — jeżeli $|r_{ij}| > r_{kr}$, odrzucamy H_0 (korelacja istotna)
 — jeżeli $|r_{ij}| < r_{kr}$, przyjmujemy H_0 (brak korelacji)

W przykładzie, jeżeli przyjmujemy $\alpha = 0,05$, to $r_{kr} = 0,3778$.

Pozostanę również tylko trzy istotne powiązania: y-x₁, y-x₂, x₁-x₂

Metoda grafowa

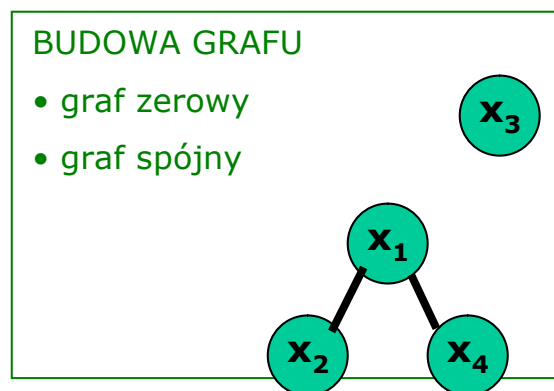
- obliczenie macierzy współczynników:
- wyodrębnienie z macierzy powiązań istotnych
- budowa grafu z powiązań istotnych
- wybranie zmiennych na podstawie grafu

$$r_{x_i x_j}, r_{y x_i}$$

DO MODELU WYBIERA SIĘ ZMIENNE:

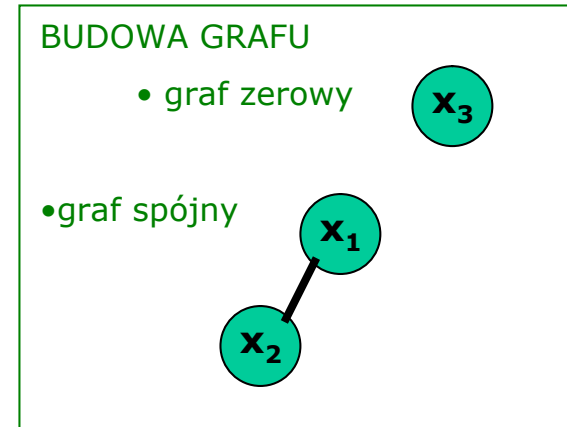
- każdy graf zerowy (jeżeli $|r_{yx}| > 0,1$)
- po jednej reprezentantce grafu spójnego;
 - reprezentantką grafu jest zmienna, która ma najwięcej powiązań z innymi kandydatkami
 - jeżeli kilka zmiennych ma tę samą maksymalną liczbę powiązań wybiera się ta która jest najsilniej powiązana ze zmienną Y ($\max r_{yx}$)
 - jeżeli graf jest rozległy, może mieć dwie reprezentantki, ale muszą one leżeć na przeciwległych stronach grafu

Przykład grafu



Przykład cd. Wybierzemy zmienną x_2 (reprezentantka grafu spójnego) i zmienną x_3 (graf zerowy)

	Y	X_1	X_2	X_3
Y		0,52	0,64	-0,21
X_1	0,52		0,82	-0,18
X_2	0,64	0,82		0,08
X_3	-0,21	-0,18	0,08	



Metoda Hellwiga

- wypisujemy wszystkie możliwe kombinacje kandydatek;

jest ich $l=2^m-1$

- obliczamy pojemność indywidualną i kłosa informacji (dla każdej zmiennej w każdej kombinacji)

$$h_{lj} = f(r_{yx_j}, r_{x_i x_j})$$

$$h_{lj} = \frac{r_{y,j}^2}{1 + \sum_{i \neq j} |r_{ij}|}$$

$i, j = 1, 2, \dots, m$

$l = 2^m - 1$ m - ilość kandydatek

r_j - wsp. korelacji j -tej kandydatki ze zmienną objaśnianą

r_{ij} - wsp. korelacji i -tej i j -tej zmiennej

- obliczamy pojemność całkowitą dla każdej kombinacji

$$H_l = \sum_j h_{lj}$$

- wybieramy kombinację o największej pojemności $H_l = \max$

Metody taksonomiczne

Taksonomia wrocławska (cluster analysis) to metoda grupowania obiektów (zmiennych) w grupy jednorodne pod względem n cech (wymiarów) łącznie. Podstawą grupowania jest odległość euklidesowa, która w przypadku zmiennych

$$d_{ij} = f(r_{x_i x_j})$$

WSZYSTKIE METODY TO SELEKCJA WSTĘPNA

ETAP 4. ESTYMACJA PARAMETRÓW MODELU

Cel etapu: wyznaczenie parametrów strukturalnych i stochastycznych

Estymacja: szacowanie parametrów populacji na podstawie próbki

Metody estymacji: analiza regresji i korelacji (KMNK i inne)

Założenia dla KMNK

1. Zmienne losowe są zmiennymi nie powiązаныmi ze sobą (nie występuje współliniowość)
2. Składnik losowy ξ jest zmienną losową $E(\xi)=0$; $V^2=\text{const}$ (stała wariancja, niezależna od zmiennej x lub t)
3. Składnik losowy ξ nie jest powiązany ze zmiennymi objaśnającymi

4. Wartości reszt u_i są niezależne od siebie

$$u_i = y_i - \hat{y}_i$$

5. m - liczba zmiennych objaśnających;

n - liczność próby:

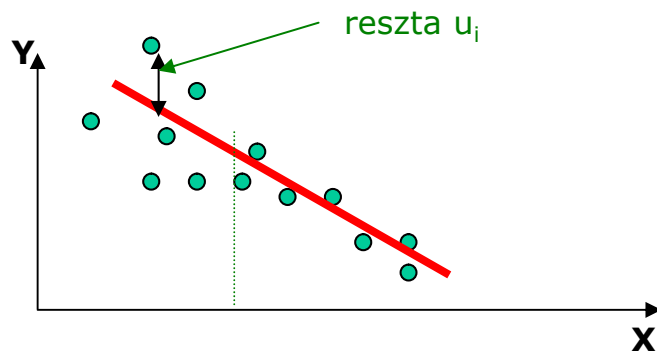
$$m < n$$

Skutki niedotrzymania założenia o błędach niezależnych

1. Model nieprzydatny; niekiedy absurdalny (nie uwarunkowane dane)

Stosuje się do próbek ($n > 100$), regresję grzbietową (ridge regression)

2. Lewa część zbioru ma dużą wariancję a prawa – małą. Stosuje się specjalny wariant MNK z korektą na różne wariancje



$u_t = y_t - \hat{y}_t$	u_{t-1}
u_1	—
u_2	u_1
u_3	u_2
u_4	u_3
u_5	u_4

3. Jeśli reszty u_i są ze sobą powiązane (skorelowane) tzn. że występuje *autokorelacja składnika losowego* (najczęściej zjawisko występuje przy szeregach czasowych). Oznacza to, że istnieje istotna zależność

$$u_t = f(u_{t-k}) \quad t = 1, 2, \dots$$

Przyczyny autokorelacji:

zakłócenia (dodatnie lub ujemne) w jednym okresie wpływają na poziom zjawiska w następnym okresie

Występowanie autokorelacji powoduje nieprzydatność modelu

4. Składnik losowy jest skorelowany ze zmienną objaśnianą wtedy gdy została pominięta jakaś ważna zmienna – przyczyna. Model taki nie ma żadnej wartości; trzeba dbać o jak najwyższy współczynnik determinacji ($R^2 > 0,9$)

ETAP 5. WERYFIKACJA MODELU

- Cele:**
1. opis rzeczywistości (populacji generalnej)
 2. dokładna (ostateczna) selekcja zmiennych objaśniających
 3. poznanie składowika losowego (spełnienie założeń KMNK)

Narzędzia: hipotezy i testy statystyczne

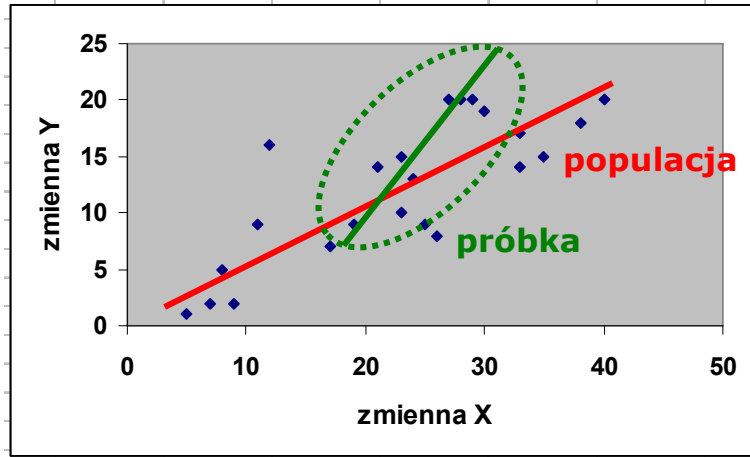
Metodyka: każde równanie oddzielnie; weryfikacja obejmuje 9 etapów
(od najmniej do najbardziej pracochłonnego)

WYKAZ ETAPÓW WERYFIKACJI MODELU

- 5.1. Badanie istotności korelacji
- 5.2. Badanie wyrazistości modelu
- 5.3. Badanie istotności parametrów
- 5.4. Badanie symetrii składowego
- 5.5. Badanie losowości składowego
- 5.6. Badanie stacjonarności składowego
- 5.7. Badanie wartości oczekiwanej składowego
- 5.8. Badanie autokorelacji składowego
- 5.9. Badanie normalności składowego

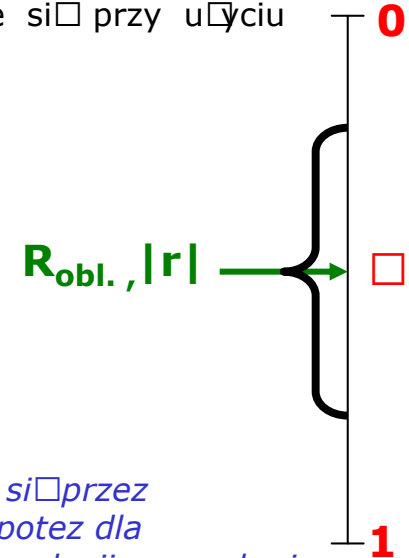
ETAP 5.1. Badanie istotności korelacji

Celem etapu jest sprawdzenie, czy istnieje w populacji generalnej powiązanie pomiędzy zmienną Y i wszystkimi zmiennymi objaśniającymi



Przedział ufności dla nieznanego współczynnika korelacji ρ dla populacji buduje się przy użyciu błędów

$$s_R = \sqrt{\frac{1 - R^2}{n - k}}$$



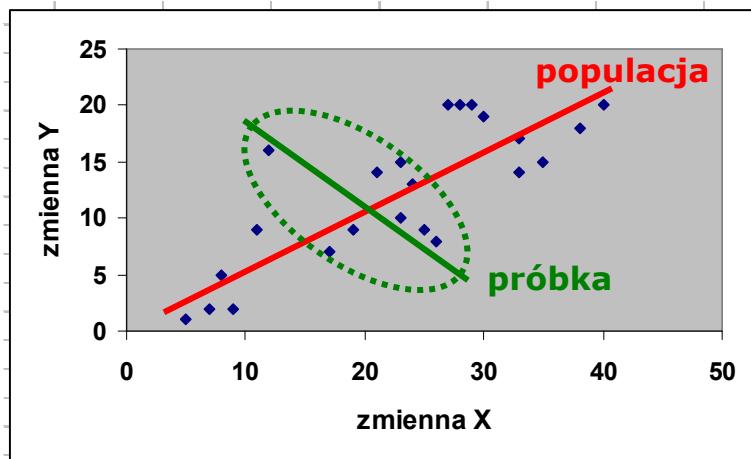
Istotność korelacji weryfikuje się przez postawienie następujących hipotez dla współczynnika korelacji dla populacji generalnej:

$H_0 : \rho = 0$

$H_1 : \rho \neq 0$

Brak korelacji, nie ma powiązania...

Korelacja istotna, jest powiązanie...



Hipotezy te można weryfikować trzema równoważnymi testami:

- testem t Studenta (tylko dla regresji dwu zmiennych)
- testem F Fishera
- testem R Wallace'a-Snedecora

UWAGA!

W przypadku regresji wielorakiej, gdy liczba zmiennych objaśnianych jest duża w porównaniu z liczbą obserwacji (n), współczynnik determinacji R^2 może dawać zawyżoną ocenę stopnia wyjaśnienia zmienności zmiennej objaśnianej; dlatego wprowadzono

skorygowany współczynnik determinacji R_a^2 (i korelacji):

$$R_a^2 = 1 - \left(\frac{n-1}{n-k} \right) \left(\frac{\text{SSE}}{\text{TOSS}} \right) = 1 - (n-1) \left(\frac{s^2}{\text{TOSS}} \right)$$

k – ilość parametrów w modelu regresji

$$R_a^2 = 1 - \frac{\frac{\sum(\hat{y}_i - \bar{y})^2}{n-k}}{\frac{\sum(y_i - \bar{y})^2}{n-1}}$$

- *adjusted coefficient of multiple determination* (wydruki komputerowe)
- jeżeli k jest małe, nie ma większej różnicy pomiędzy normalnym a skorygowanym R^2

TEST STUDENTA

$$t_{obl} = r \sqrt{\frac{n-2}{1-r^2}} = \frac{r}{s_r} \qquad t_{tabl} = t_{\alpha/2}\{n-2\}$$

TEST FISHERA

$$F_{obl} = \frac{MSTR}{MSE} = \frac{R^2}{1-R^2} \frac{n-k}{k-1} \qquad F_{tabl} = F_{\alpha}\{k-1, n-k\}$$

Źródło zmienności	Liczba stopni swobody	Suma kwadratów	Średni kwadrat	Statystyka F
Model (czynniki)	k-1	SSTR	MSTR	$F_{obl} = \frac{MSTR}{MSE}$
Błąd (reszta)	n-k	SSE	MSE	
Razem	n-1	SSTO		

TEST WALLACE'A-SNEDECORA

$$R_{obl} = \sqrt{R^2} \qquad R_{tabl} = R_{\alpha}\{k, n-k\}$$

Zmienna (czynnik)	Wartość oszacowana	Błąd oszacowania	Statystyka t_{obl}	Rzeczywisty poziom istotności P
Wyraz wolny	a_0	$s(a_0)$	$t(a_0)$	$P(a_0)$
Czynnik X_1	a_1	$s(a_1)$	$t(a_1)$	$P(a_1)$
Czynnik X_2	a_2	$s(a_2)$	$t(a_2)$	$P(a_2)$

Współczynniki: determinacji R^2 , zbłądności ϕ^2 , błąd resztowy $s(y)$ i inne

Odczyt R_{tabl} z tablicy testu R Wallace'a-Snedecora

Stopnie swobody	Liczba zmiennych					
	2		3		4	
	0,05	0,01	0,05	0,01	0,05	0,01
8	0,632	0,765	0,726	0,827	0,777	0,860
18	0,444	0,561	0,532	0,633	0,587	0,678
28	0,361	0,463	0,439	0,530	0,490	0,573

Wnioski rozkładu R Wallace'a-Snedecora:

- im wyższy poziom istotności, tym niższe R_{tabl}
- im większa liczba zmiennych w modelu, tym wyższe R_{tabl}
- im większa liczba stopni swobody (większa próbka), tym niższe R_{tabl}
- tablica R powstała z przeliczenia tablic t oraz F (odwrócenie wzorów)

Tablica testu R jest najszybszym i najwygodniejszym narzędziem do weryfikacji istotności korelacji

Reguła decyzyjna (podsumowanie etapu 5.1.):

jeżeli $R_{obl} > R_{tabl}$, model jest poprawny, można przejść do etapu 5.2

jeżeli $R_{obl} < R_{tabl}$, model jest niepoprawny, trzeba zmienić albo zestaw zmiennych objaśnianych albo jego postać matematyczną

ETAP 5.1. OBOWIĄZUJE DLA WSZYSTKICH MODELI

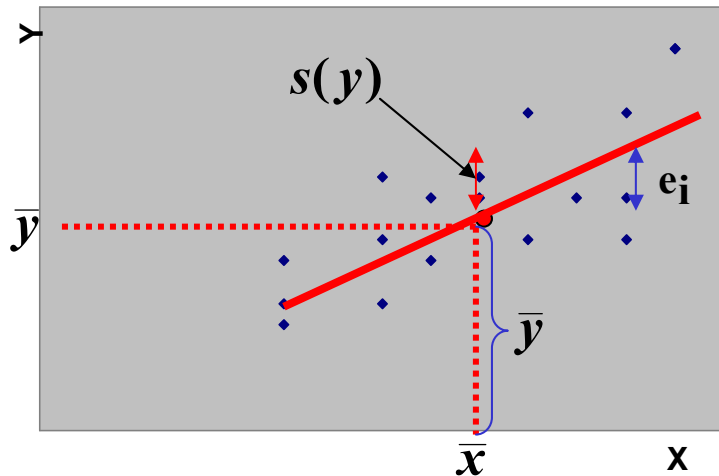
ETAP 5.1. PRZEPROWADZA SIĘ DLA TRANSFORMATY (dla postaci pierwotnej można uzyskać $R > 1$)

Rola współczynnika determinacji R^2

- korelacja może być istotna przy małym R i bardzo małym R^2
- małe R^2 oznacza niski stopień wyjaśnienia rzeczywistości i stanowi zagrożenie dla modelu
- należy dążyć (poprzez odpowiedni dobór zmiennych-przyczyn i postaci matematycznej modelu) do jak największego R^2 (dla postaci pierwotnej)
- wysoka wartość R^2 świadczy o dobrym poznaniu badanego zjawiska
- wysoka wartość R^2 bardzo często wynika jednak ze złego dobrania zmiennych objaśnianych

ETAP 5.2. Badanie wyrazistości modelu

Celem etapu jest kontrola rozrzutu danych



Wyrazistość modelu dana jest wzorem

$$V_{obl} = \frac{s(y)}{\bar{y}} 100 \%$$

Współczynnik zmienności losowej $V_{obl} < 30\%$

(w przeciwnym przypadku rozrzut danych jest zbyt duży)

Uwaga: gdy \bar{y} jest bliskie 0 trudno jest w ustaleniu czy model poprawny czy niepoprawny

ETAP 5.2. OBOWIĄZUJE DLA WSZYSTKICH MODELI (ale nie ma on charakteru statystycznego)

ETAP 5.2. PRZEPROWADZA SIĘ DLA POSTACI PIERWOTNEJ

ETAP 5.3. Badanie istotności parametrów (współczynników) modelu

Celem etapu jest sprawdzenie:

- czy poszczególne zmienne objaśnialne mają istotny wpływ na zmienną objaśnianą?
- czy zmienne objaśnialne są wybrane prawidłowo?
- czy wyraz wolny różni się istotnie od zera?

W etapie 5.3. następuje ostateczna selekcja zmiennych objaśniających:

- jeżeli wszystkie a_i okazały się istotne, model jest poprawny:
 - model przyczynowo-skutkowy: do interpretacji
 - model tendencji rozwojowej: do etapu 5.4
- jeżeli choć jedno a_i okazało się nieistotne, model jest niepoprawny i wymaga poprawy przez usunięcie nieistotnych zmiennych:
 - zmienne należy usuwać po jednej (ze względu na efekt katalityczny)
 - usuwa się zawsze zmienną o najniższej wartości $|t(a_i)|$ [$\max P(a_i)$]
 - usunięcie ostatniej zmiennej nieistotnej kończy proces selekcji kandydatek na zmienne,
 - selekcja nie jest ostateczna, gdyż zawsze istnieje możliwość zamiany zmiennych, które są powiązane ze sobą

ETAP 5.3. OBOWIĄZUJE DLA WSZYSTKICH MODELI

ETAP 5.3 PRZEPROWADZA SIĘ DLA TRANSFORMATY

ETAP 5.4. Badanie symetrii składownika losowego

Badanie symetrii: dla $n > 30$ test z (r-d normalny); dla $n < 30$ test t -Studenta

$$H_0: \frac{m}{n} = \frac{1}{2}$$

$$H_1: \frac{m}{n} \neq \frac{1}{2}$$

$$t_{obl} = \frac{\left| \frac{m}{n} - \frac{1}{2} \right|}{\sqrt{\frac{\frac{m}{n} \left(1 - \frac{m}{n} \right)}{n-1}}} \quad t_{\alpha, v=n-1}$$

m – liczba reszt dodatnich (lub ujemnych)

n – licznok próby

Brak symetrii wymaga zmiany matematycznej postaci modelu

ETAP 5.4. OBOWIĄZUJE DLA MODELI TENDENCJI ROZWOJOWEJ

ETAP 5.4. PRZEPROWADZA SIĘ DLA POSTACI PIERWOTNEJ

Zmienna (czynnik)	Wartość oszacowana	Błąd oszacowania	Statystyka t_{obl}	Rzeczywisty poziom istotności P
Wyraz wolny	a_0	$s(a_0)$	$t(a_0)$	$P(a_0)$
Czynnik X_1	a_1	$s(a_1)$	$t(a_1)$	$P(a_1)$
Czynnik X_2	a_2	$s(a_2)$	$t(a_2)$	$P(a_2)$
Czynnik X_3	a_3	$s(a_3)$	$t(a_3)$	$P(a_3)$

Współczynniki: determinacji R^2 , zbliżności φ^2 , błąd resztowy $s(y)$ i inne

Istotność parametrów a_i można sprawdzać (tak jak problem średniej dla populacji) na dwa sposoby:

- konstruujemy przedziały ufności dla nieznannej wartości α_i
- weryfikujemy hipotezy — $H_0: \alpha_i=0$ wobec $H_1: \alpha_i \neq 0$

$$P(a_i - t_{\alpha/2}\{n - k\}s(a_i) < \alpha_i < a_i + t_{\alpha/2}\{n - k\}s(a_i)) = 1 - \alpha$$

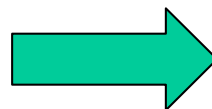
W praktyce stosuje się hipotezy:

Badanie istotności przeprowadza się dla każdego α_i oddzielnie: $t_{obl}(a_i) = \frac{a_i - 0}{s(a_i)}$ $t_{tabl} = t_{\alpha/2}\{n - k\}$

- jeżeli $|t_{obl}(a_i)| > t_{tabl}(a_i)$, odrzucamy hipotezę zerową parametr jest istotny z błędem równym co najwyżej α
- jeżeli $|t_{obl}(a_i)| < t_{tabl}(a_i)$, nie ma podstaw do odrzucenia hipotezy zerowej; parametr jest nieistotny lub
- porównujemy rzeczywisty poziom istotności $P(a_i)$ z przyjętym *a priori* α

Jeżeli $P(a_i)$ jest mniejsze od α odrzucamy H_0

Odrzucając H_0



ZMIENNA X_i
MA WPŁYW NA
ZMIENNĄ Y

ETAP 5.5. Badanie losowości składnika losowego

Badanie losowości przeprowadza się testem t -Studenta lub testem serii

Test serii:

$$H_0 : \xi_t \text{ jest składnikiem losowym} \quad H_1 : \xi_t \text{ nie jest losowy}$$

Wartościom $u_t > 0$ nadajemy symbol a (liczba symboli a : n_1), wartościom $u_t < 0$ symbol b (liczba symboli b : n_2). Otrzymujemy podciąg czyli serie z kolejnych symboli a lub b . Liczbę wszystkich serii (podciągów) oznaczmy jako K . W tablicy liczby serii (dla n_1 ; n_2 ; α) odczytujemy liczbę krytyczną K_α .

$$P(K \leq K_\alpha) = \alpha \quad \text{Gdy } K > K_\alpha \text{ nie ma podstaw do odrzucenia } H_0$$

$$P(K > K_\alpha) = 1 - \alpha \quad \text{Gdy } K \leq K_\alpha \text{ odrzucamy } H_0$$

Brak losowości wymaga zmiany matematycznej postaci modelu

ETAP 5.5. OBOWIĄZUJE DLA MODELI TENDENCJI ROZWOJOWEJ

ETAP 5.5. PRZEPROWADZA SIĘ DLA POSTACI PIERWOTNEJ

ETAP 5.6. Badanie stacjonarności składnika losowego

W etapie tym sprawdzamy niezależność wariancji składnika losowego od zmiennej objaśniającej t (warunek stosowalności KMNK)

Bada się istotność współczynnika korelacji reszt i zmiennej czasowej t

$$r_{u_t, t}$$

Stosuje się test istotności t (slajd: badanie współczynnika korelacji)

Przyczyny braku stacjonarności:

- niewłaściwa postać analityczna modelu
- niewłaściwa metoda szacowania parametrów strukturalnych modelu

ETAP 5.6. OBOWIĄZUJE DLA MODELI TENDENCJI ROZWOJOWEJ

ETAP 5.6 PRZEPROWADZA SIĘ DLA TRANSFORMATY

ETAP 5.7. Badanie wartości oczekiwanej składowika losowego

Z zał. KMNK wynika, że parametry składowika losowego:

$$E(\xi) = 0 \quad \left\{ \sum u_i = 0 \right\} \quad V(\xi) = \text{const}$$

Po retransformacji do postaci pierwotnej mamy nowe reszty u' , dla których

$$\sum u'_i \neq 0 \quad \bar{u}'_i \neq 0$$

Celem etapu jest sprawdzenie, czy odchylenie od „0” nie jest zbyt duże (sprawdza się to testem Studenta)

ETAP 5.7. OBOWIĄZUJE DLA MODELI TENDENCJI ROZWOJOWEJ

ETAP 5.7. PRZEPROWADZA SIĘ DLA POSTACI PIERWOTNEJ

ETAP 5.8. Badanie autokorelacji składowika losowego

Składowik losowy nie jest czysto losowy, lecz zależy od wskaźnika i , czyli zmienne losowe są zależne od poprzednich wartości t .

Autokorelacja to korelacja wartości zmiennej z jej wartościami z okresów wcześniejszych o jeden lub więcej okresów.

Na ogół autokorelację można wyrazić w postaci relacji: $\xi_i = f(\xi_{i-1}, \xi_{i-2}, \dots, \xi_{i-\tau}) \quad u_t = f(u_{t-k}) \quad t = 1, 2, \dots$

W praktyce przyjmuje się, że funkcja f jest funkcją liniową a maksymalne opóźnienie czasowe wynosi jeden lub dwa.

Estymator współczynnika autokorelacji r_1 (rzędu pierwszego):

$$r_1 = \frac{\sum_{i=2}^n (u_i - \bar{u}_i)(u_{i-1} - \bar{u}_{i-1})}{\sqrt{\sum_{i=2}^n (u_i - \bar{u}_i)^2 \sum_{i=2}^n (u_{i-1} - \bar{u}_{i-1})^2}}$$

Przyczyny autokorelacji:

- zakłócenia (dodatnie lub ujemne) w jednym okresie wpływają na poziom zjawiska w następujących okresach: skutki niektórych zdarzeń rozciągają się na wiele okresów (natura procesów gospodarczych, społecznych);
- psychologia i sposób podejmowania decyzji, na które duży wpływ mają zdarzenia z najbliższej przeszłości;
- niepoprawna postać funkcyjna modelu (model nie uwzględnia cykliczności zjawiska, aproksymacja zależności nieliniowej przez funkcję liniową);
- wadliwa struktura dynamiczna modelu: w roli zmiennej objaśniającej nie występuje – a powinna – opóźniona zmienna objaśniana; brak opóźnionych zmiennych niezależnych lub zmiennej czasowej;
- pominięcie w modelu istotnej zmiennej objaśniającej (reszty mogą układać się seriami, nawet mieć tendencję do stałego zwiększania swej wartości bezwzględnej);
- interpolacja, wygładzanie oraz agregacja (np. przekształcanie danych miesięcznych w kwartalne).

Badanie autokorelacji można przeprowadzić

- testem R istotności korelacji $r_{u_t, u_{t-k}}$
- testem Durbina-Watsona

Test Durbina-Watsona służy do sprawdzenia hipotezy: $H_0 : \rho_1 = 0$ $H_1 : \rho_1 < 0$ lub $H_1 : \rho_1 > 0$

Statystyka d:

$$d = \frac{\sum_{i=2}^n (u_i - u_{i-1})^2}{\sum_{i=2}^n u_i^2}$$

$$d = 2(1-r_1)$$

Z relacji wynika, że $d \in [0,4]$:

- jeżeli $r_1=0$ to $d=2$ (brak autokorelacji)
- jeżeli $r_1=1$ to $d=0$ (silna autokorelacja dodatnia)
- jeżeli $r_1=-1$ to $d=4$ (silna, ujemna autokorelacja)

Rozkład statystyki d przy założeniu, że H_0 jest prawdziwa i składniki losowe mają rozkład normalny $N(0; \sigma^2)$ zależy od liczby obserwacji n oraz liczby zmiennych objaśnianych k i $d \in \langle d_L; d_U \rangle$.

Wartości krytyczne d_L i d_U zawiera tablica testu Durbina-Watsona dla poziomu istotności α .

Reguła decyzyjna:

$$H_0 : \rho_1 = 0 \quad H_1 : \rho_1 < 0$$

- jeżeli $d < d_L$ to H_0 odrzucamy
- jeżeli $d_L < d < d_U$?
- jeżeli $d > d_U$ nie ma podstaw do odrzucenia H_0

Reguła decyzyjna:

$$H_0 : \rho_1 = 0 \quad H_1 : \rho_1 > 0$$

- jeżeli $d > 4 - d_L$ to H_0 odrzucamy
- jeżeli $4 - d_U < d < 4 - d_L$?
- jeżeli $d < 4 - d_U$ nie ma podstaw do odrzucenia H_0

W przypadku stwierdzenia autokorelacji mamy trzy możliwości:

- usunąć przyczyny autokorelacji;
- zastosować procedury estymacji w warunkach autokorelacji, aby zapewnić wysoką efektywność estymatorów (usuwanie autokorelacji);
- pozostać przy KMNK godząc się na mniejszą efektywność estymatorów.

ETAP 5.8. OBOWIĄZUJE DLA MODELI TENDENCJI ROZWOJOWEJ

ETAP 5.8. PRZEPROWADZA SIĘ DLA TRANSFORMATY

Etap 5.9. Badanie normalności składnika losowego

Celem etapu jest stwierdzenie, czy reszty mają rozkład normalny

Stosuje się znane ze statystyki testy nieparametryczne:

- λ - Kolmogorowa

lub

- test χ^2

Analiza reszt — oddzielny dział analizy regresji i korelacji

ETAP 6. INTERPRETACJA MODELU

Celem etapu jest wydobyć z modelu całą nową wiedzę, której nie widać „gołym okiem”

Interpretacja modelu przyczynowo-skutkowego

Polega na określeniu wpływu poszczególnych czynników na badane zjawisko:

- podział wszystkich czynników na trzy grupy

A — czynniki nie mające wpływu na Y

B — czynniki mające wpływ na Y i wprowadzone do modelu

C — czynniki mające wpływ na Y, ale nie występujące w modelu

- ocena jakościowa wpływu czynników B i C
- ocena ilościowa wpływu czynników B i C

Ocena jakościowa

Na podstawie znaków stojących przy współczynnikach r oraz a_i

Ocena ilościowa

• **Model potęgowy – FUNKCJA PRODUKCJI COBBA-DOUGLASA** $P = \alpha_0 X_1^{\alpha_1} X_2^{\alpha_2} e^{\gamma \xi}$

α_1

Współczynniki elastyczności produkcji względem X_1 i X_2

α_2

$\alpha_1, \alpha_2 > 0$

$K = \alpha_1 + \alpha_2$ **Efekt skali produkcji**

Funkcja produkcji to specjalny model, określający zależność pomiędzy produkcją (P), a czynnikami produkcji: majątkiem produkcyjnym (X_1) i nakładami pracy żywej (X_2).

Metoda estymacji parametrów strukturalnych funkcji Cobba-Douglasa (jak dla modelu potęgowego) to MNK

Graficznym obrazem funkcji jest krzywa wypukła do początku układu współrzędnych X_1, X_2 . Poruszając się po krzywej P otrzymujemy ten sam wolumen produkcji P przy różnych kombinacjach czynników-nakładów X_1 i X_2 .

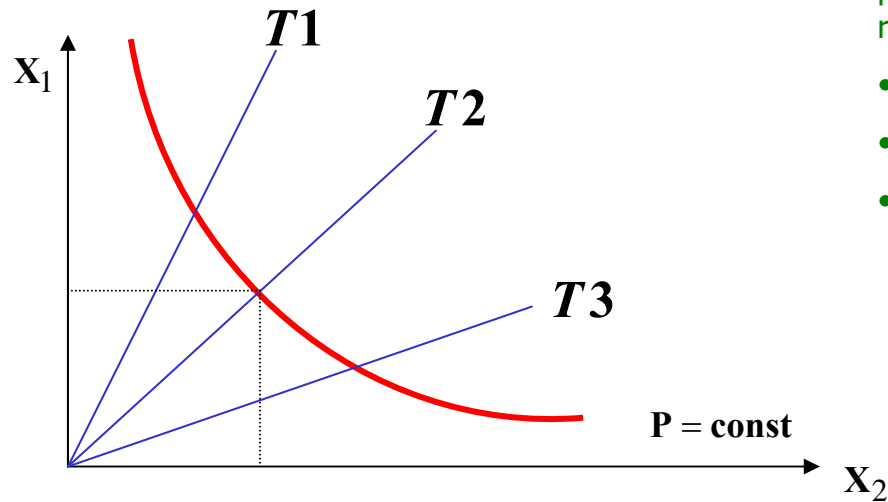
Ustalanie różnych proporcji czynników X_1 i X_2 w celu otrzymania tej samej wielkości produkcji P jest możliwe ze względu na wzajemnie substytucyjne czynniki.

Współczynnik elastyczności substytucji: $\varepsilon = \frac{\Delta X_1}{X_1} / \frac{\Delta X_2}{X_2}$

Elastyczność ε to współczynnik określający reakcję zmiany jednej zmiennej na inną zmienną – na przykład elastyczność cenowa popytu (jak popyt zmienia się wraz z ceną towaru).

Z ekonomicznego punktu widzenia elastyczność produkcji ε określa o ile należy zwiększyć nakłady na rodki trwałe X_1 , aby nakłady pracy X_2 zmniejszyło o jednostkę (przy stałym poziomie produkcji P).

Izolinie funkcji produkcji



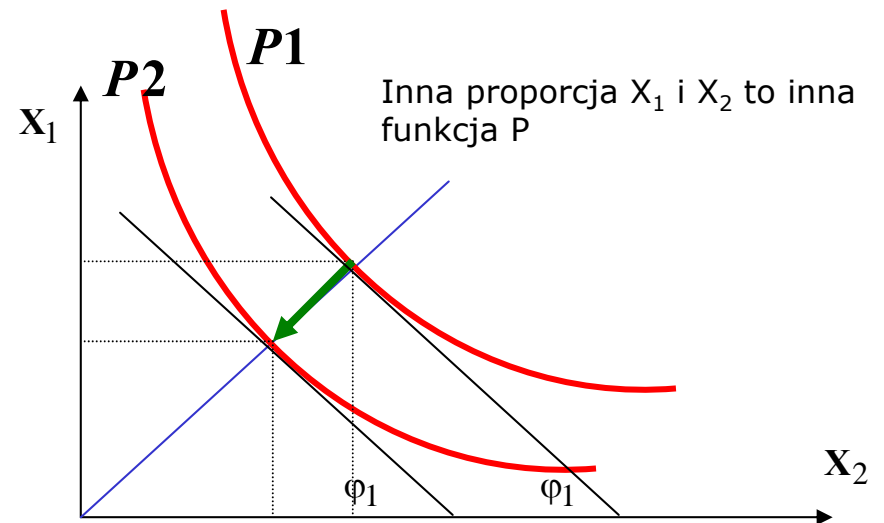
Wartość produkcji P można osiągnąć przy różnych proporcjach pomiędzy majątkiem X_1 i nakładami siły roboczej X_2 :

- T1 – technika kapitałochłonna, np. pełna automatyzacja
- T2 – pośredni poziom techniki
- T3 – niski poziom techniki, produkcja pracochłonna

Efekt poprawy zarządzania bez zmiany techniki (lepsze wykorzystanie ludzi i sprzętu):

- zmniejszenie kapitałochłonności
- zmniejszenie pracochłonności
- zwiększenie produkcji (przejście z funkcji P_1 na P_2)

Krańcowa stopa substytucji $\text{tg}\varphi_1 = R$



Przykład. Zinterpretuj zweryfikowaną funkcję produkcji Cobba-Douglasa

$$P_i = 45,2M_i^{0,65}Z_i^{0,73}e^{0,04\xi}$$

- wartość 45,2 nie interpretuje się
- wartość 0,65 to współczynnik elastyczności majątkowej (kapitałowej) produkcji: zmiana majątku o 10% daje średnio zmianę produkcji o 6,5% (przy stałym zatrudnieniu = pozostałe czynniki ceteris paribus)
- wartość 0,73 to współczynnik elastyczności zatrudnieniowej produkcji: zmiana zatrudnienia o 10% daje zmianę produkcji średnio o 7,3% (przy stałym majątku)
- suma obu współczynników $K=0,65+0,73=1,38$ to współczynnik skali produkcji:
 - jeżeli $K < 1$, to firma rozwija się ekstensywnie (szybsze tempo wzrostu czynników produkcji aniżeli przyrostu produkcji)
 - jeżeli $K > 1$, to firma rozwija się intensywnie (szybsze tempo przyrostu produkcji aniżeli wzrostu czynników produkcji)
- wartość $\gamma = 0,04$ to współczynnik postępu organizacyjnego:
 - jeżeli $\gamma > 0$, to miało miejsce postępowanie organizacyjne
 - jeżeli $\gamma < 0$, to miało miejsce regres organizacyjny

MODEL TENDENCJI ROZWOJOWEJ

Model tendencji rozwojowej to konstrukcja teoretyczna (równanie lub układ równań) opisująca kształtowanie się określonego zjawiska jako funkcji:

- zmiennej czasowej t
- wahań okresowych (periodycznych)
- wahań przypadkowych (nieregularnych).

Czyli na zmienność zjawiska w czasie ma wpływ:

- tendencja rozwojowa (trend)
- wahania typu okresowego
- wahania typu przypadkowego (losowego).

MODEL ADDYTYWNY $Y_t = F(t) + G(t) + \xi(t)$

MODEL MULTIPLIKATYWNY $Y_t = F(t) \cdot G(t) \cdot \xi(t)$

gdzie:

Y_t - poziom badanego zjawiska

$F(t)$ - funkcja trendu

$G(t)$ - funkcja wahań okresowych

$\xi(t)$ - składnik losowy o rozkładzie normalnym, $E(\xi)=0$, $V(\xi)=const$

Te trzy części trzeba zidentyfikować, a potem złożyć razem w model:

- addytywny (jeśli amplituda wahań jest stała)
- multiplikatywny (jeśli amplituda wahań rośnie lub maleje regularnie)

WYZNACZANIE TRENDU

Metody mechaniczne: $y_t = f(t) + \xi_t$

- wskaźniki bezwzględne (stały przyrost = postęp arytmetyczny = funkcja liniowa)

$$\bar{x}_{chr} = \frac{\frac{1}{2}x_1 + \dots + \frac{1}{2}x_n}{n-1}$$

- wskaźniki względne (przyrost o stały procent = postęp geometryczny = funkcja potęgowa)

- Średnia ruchoma (wygładzanie danych liniogramem) $\bar{x}_g = \sqrt[n]{\prod x_i}$

Średnia ruchoma: mając szereg czasowy y_1, y_2, \dots, y_n przyjmujemy długość kroku $k = 3$, lub 5 , lub 7 itd. i liczymy

$$\bar{y}_{k/2+0,5} = \frac{y_1 + \dots + y_k}{k} \quad \text{a następnie} \quad \bar{y}_{k/2+1,5} = \frac{y_2 + \dots + y_{k+1}}{k} \quad \text{itd.}$$

dla $k=3$: $\bar{y}_2 = \frac{y_1 + y_2 + y_3}{3} \quad \bar{y}_3 = \frac{y_2 + y_3 + y_4}{3} \quad \text{itd.}$

Dla $k=5$ $\bar{y}_3 = \frac{y_1 + y_2 + y_3 + y_4 + y_5}{5} \quad \bar{y}_4 = \frac{y_2 + y_3 + y_4 + y_5 + y_6}{5} \quad \text{itd.}$

Gdy k jest liczbą parzystą to uzyskujemy tzw. średnie scentrowane

Np. $k=4$ $\bar{y}_3 = \frac{\frac{1}{2}y_1 + y_2 + y_3 + y_4 + \frac{1}{2}y_5}{4}$

CECHY ŚREDNIEJ RUCHOMEJ:

- TRACI SIĘ $k-1$ DANYCH
- IM WIĘKSZE k , TYM BARDZIEJ SZTYWNY TREND

Przykład. Wyznacz trend metodą średniej ruchomej dla następujących danych:

okres	t	1	2	3	4	5	6	7	8
wartość y_t		5	4	7	6	8	10	9	8

k=5

$$\bar{y}_3 = \frac{30}{5} = 6$$

$$\bar{y}_4 = \frac{35}{5} = 7$$

$$\bar{y}_5 = \frac{39}{5} = 6,6$$

$$\bar{y}_6 = \frac{41}{5} = 8,2$$

k=3

$$\bar{y}_2 = \frac{5+4+7}{3} = 5,33$$

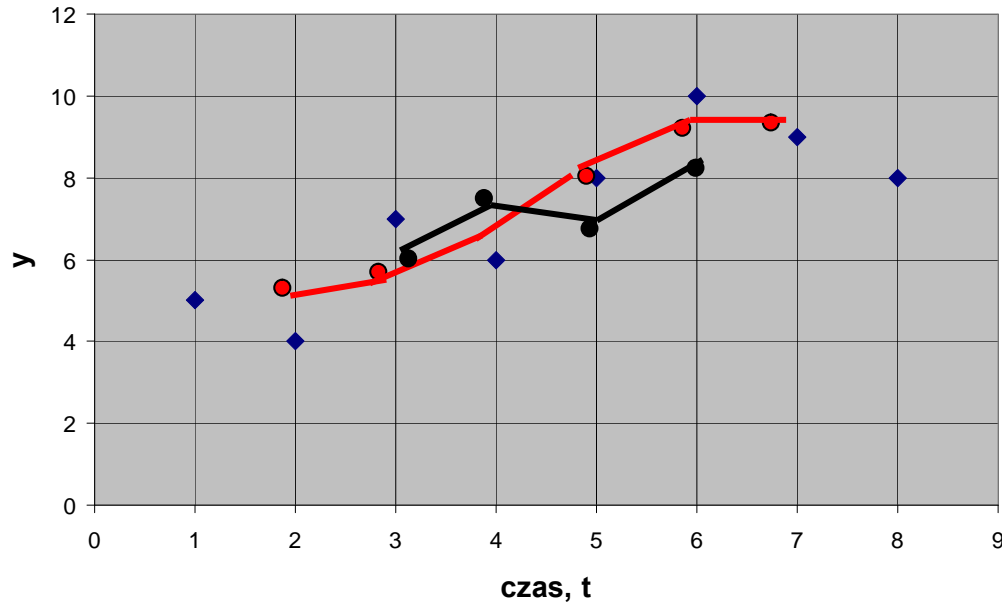
$$\bar{y}_3 = \frac{4+7+6}{3} = 5,67$$

$$\bar{y}_4 = \frac{7+6+8}{3} = 7,00$$

$$\bar{y}_5 = \frac{6+8+10}{3} = 8,00$$

$$\bar{y}_6 = \frac{8+10+9}{3} = 9,00$$

$$\bar{y}_7 = \frac{10+9+8}{3} = 9,00$$



Wygładzanie wykładnicze

Dla dowolnego momentu t operatorem wyrównania rzędu pierwszego dla szeregu y_t jest wyrażenie:

$$S_t = \alpha \cdot y_t + (1 - \alpha) \cdot S_{t-1}$$

$$0 < \alpha < 1$$

α - stała wygładzania

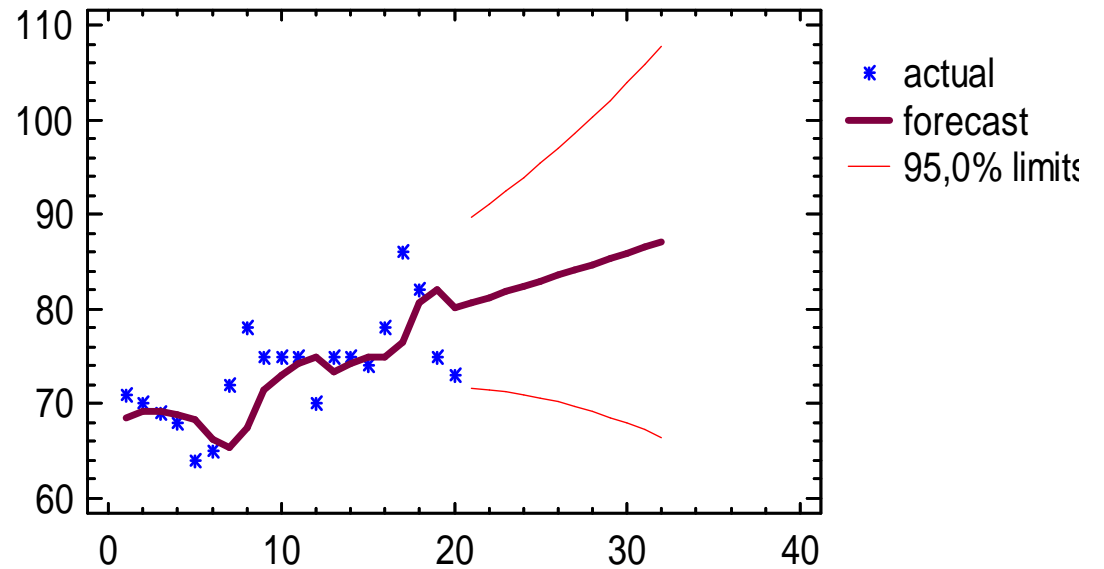
$$S_1 = y_1$$

Ogólnie operator wygładzania można zapisać wyrażeniem: $S_t = \alpha \sum_{i=0}^{t-1} (1 - \alpha)^i y_{t-i} + (1 - \alpha)^{t-1} y_1$

t	y_t	\hat{y}_t
1	71,0	71,0
2	70,0	69,1
3	69,0	69,2
4	68,0	68,9
5	64,0	68,3
6	65,0	66,3
7	72,0	65,3
8	78,0	67,5
9	75,0	71,5
10	75,0	73,1
11	75,0	74,2
12	70,0	74,9
13	75,0	73,4
14	75,0	74,3
15	74,0	74,9
16	78,0	74,9
17	86,0	76,4
18	82,0	80,7
19	75,0	82,0
20	73,0	80,1

Model: Brown's linear exp. smoothing with $\alpha = 0,2$

Wyk\adzanie wyk\adnicze trendu dla $\alpha=0,2$



B. Metody analityczne

B1

Funkcja liniowa

$$y_t = \alpha + \beta \cdot t$$

Funkcja potęgowa

$$y_t = \alpha \cdot t^\beta$$

Funkcja wykładnicza

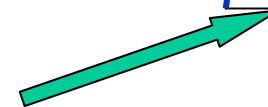
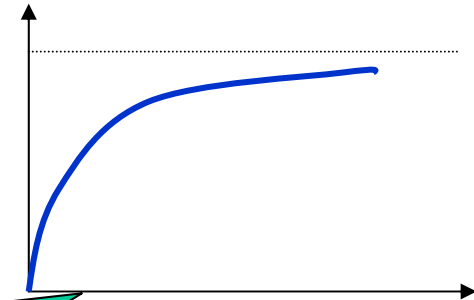
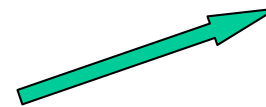
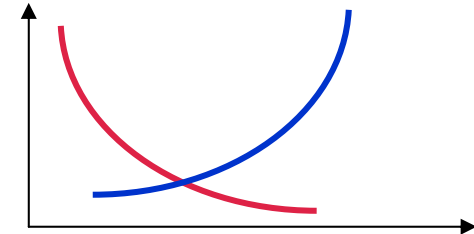
$$y_t = \alpha \cdot \beta^t$$

Funkcja logistyczna

$$y_t = \frac{\alpha}{1 + \beta \cdot e^{-\gamma \cdot t}}$$

Wielomian I rzędu

$$y_t = \alpha + \beta \cdot t + \gamma \cdot t^2$$



Parametry strukturalne funkcji oszacowywane MNK

Mierniki stopnia dopasowania modelu do danych empirycznych:

- Błąd przeciętny (Mean Error)

$$\text{M.E.} = \frac{1}{n} \sum (y_t - \hat{y}_t)$$

- Średni błąd kwadratowy (Mean Square Error)

$$\text{M.S.E.} = \frac{1}{n} \sum (y_t - \hat{y}_t)^2$$

- Średni błąd absolutny (Mean Absolute Error)

$$\text{M.A.E.} = \frac{1}{n} \sum |y_t - \hat{y}_t|$$

n – liczba reszt

B2. Wyrównywanie szeregów czasowych za pomocą wielomianów Lagrange'a

Zakładamy, że w szeregu czasowym nie występują silne wahania regularne.

$Y_t = f(t)$ zastępujemy wielomianem stopnia p .

Wtedy interpolacyjny wielomian Lagrange'a:

$$y_t = f(t_0) \frac{(t-t_1)(t-t_2)\dots(t-t_n)}{(t_0-t_1)(t_0-t_2)\dots(t_0-t_n)} + f(t_1) \frac{(t-t_0)(t-t_2)\dots(t-t_n)}{(t_1-t_0)(t_1-t_2)\dots(t_1-t_n)} + \dots$$

$$+ f(t_n) \frac{(t-t_0)(t-t_1)\dots(t-t_{n-1})}{(t_n-t_0)(t_n-t_1)\dots(t_n-t_{n-1})}$$

Liczba członów wielomianu i stopień zależy od liczby wyrazów szeregu.

Metoda bardzo pracochłonna

Wyznaczanie wahań okresowych

(dobowych, tygodniowych, miesięcznych, kwartalnych, rocznych, wieloletnich)

Przez wahania okresowe (sezonowe) należy rozumieć powtarzające się z roku na rok w tych samych jednostkach kalendarzowych doświadczenia regularne zmiany ilościowe.

Cechy wahań okresowych:

- roczny cykl z podokresami miesięcznymi, kwartalnymi, półrocznymi
- systematyczne powtarzanie się wahań w każdym roku
- określona regularność (stałe cykle zmian powtarzających się)

DWIE METODY: — metoda wskaźników sezonowości
 — metoda harmoniczna

Wszystkie metody mają jeden cel: uzyskanie przeciwnego obrazu jednego cyklu

METODA WSKAŹNIKÓW SEZONOWOŚCI

- obliczenie surowych wskaźników sezonowości (wylimowanie trendu):
- obliczenie oczyszczonych wskaźników sezonowości (wylimowanie wahań losowych):

Najprostszym sposobem wyodrębnienia wahań sezonowych jest metoda oparta na średnich okresach jednoimiennych.

Wskaźniki sezonowości oblicza się wg wzoru:

$$S_i = \frac{\bar{y}_i \cdot d}{\sum_{i=1}^d \bar{y}_i}$$

gdzie:

S_i – wskaźnik sezonowości dla i -tego podokresu (zwykle w %)

\bar{y}_i – średnia arytmetyczna dla jednoimiennych podokresów

d – liczba podokresów (podokresy miesięczne $d=12$; kwartalne $d=4$; półroczne $d=2$)

$$\sum_{i=1}^{12} S_i = 1200 \qquad \sum_{i=1}^4 S_i = 400 \qquad \sum_{i=1}^2 S_i = 200$$

Wskaźniki spełniające powyższe relacje to **oczyszczone wskaźniki sezonowości**.

Wskaźniki nie spełniające tych relacji to **surowe wskaźniki sezonowości**.

Współczynnik korygujący k : $k = \frac{d}{\sum_{i=1}^d S_i}$ pozwala sprowadzić surowe wskaźniki do oczyszczonych wg reguły:

Suma skorygowanych wskaźników: $S_i^k = k \cdot S_i$

Metoda harmoniczna (szeregi Fouriera)

Wahania okresowe przedstawia się jako sumę określonej liczby drgań harmonicznyc (sinusoid i cosinusoid) przesuniętych w fazie, lecz o jednakowym okresie

$$y_t = a_0 + a_1 t + \sum_{i=1}^{m/2} b_{1i} \sin \frac{2\pi}{m} i t + \sum_{i=1}^{m/2} b_{2i} \cos \frac{2\pi}{m} i t + \xi$$

Ogólnie, w przypadku m obserwacji liczba harmonik nie przekracza $m/2$.

Zmiany przebiegu funkcji okresowej dobrze daje się opisać za pomocą kilku początkowych harmonik (i – numer harmoniki).