

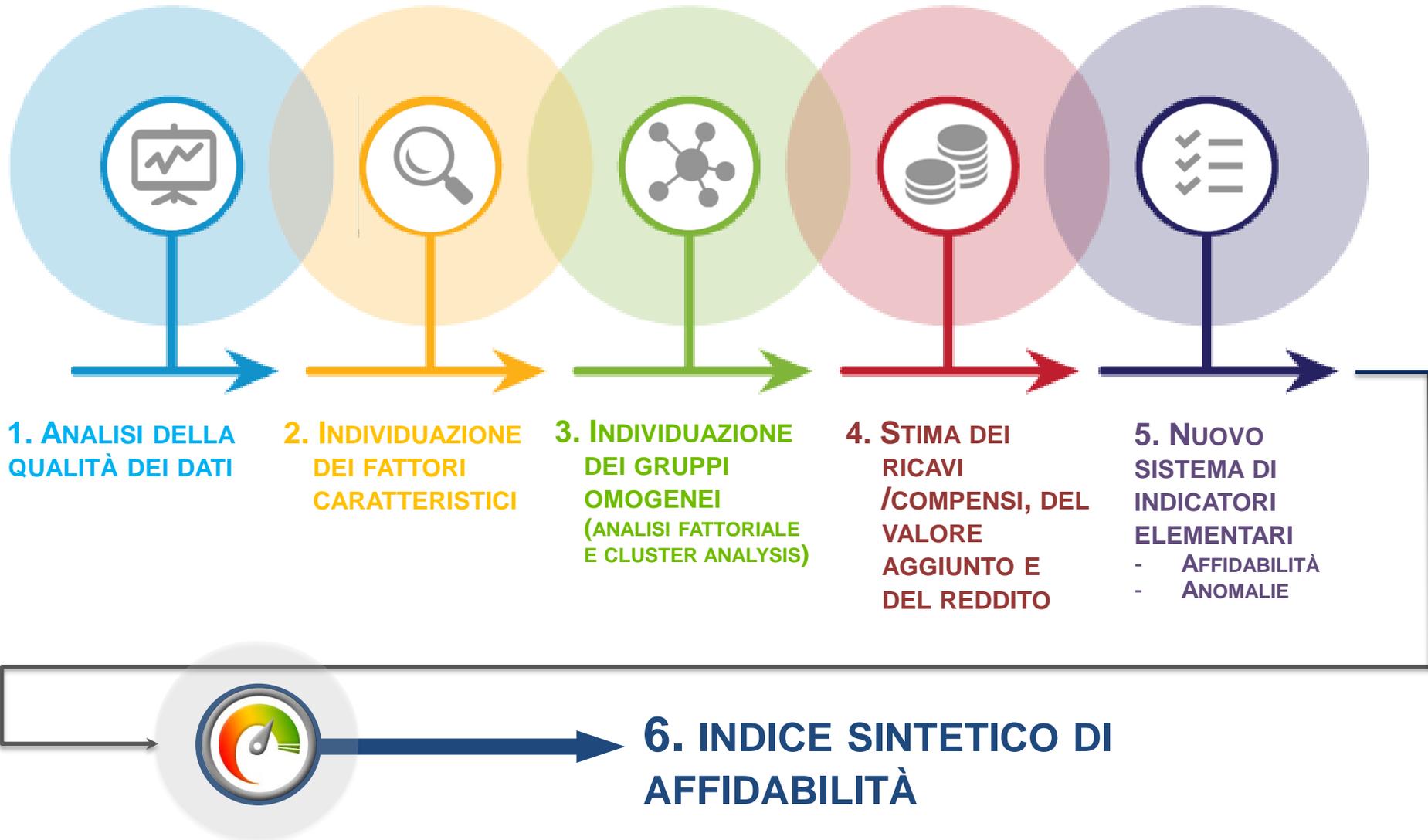
DAGLI STUDI DI SETTORE ALL' INDICE SINTETICO DI AFFIDABILITÀ

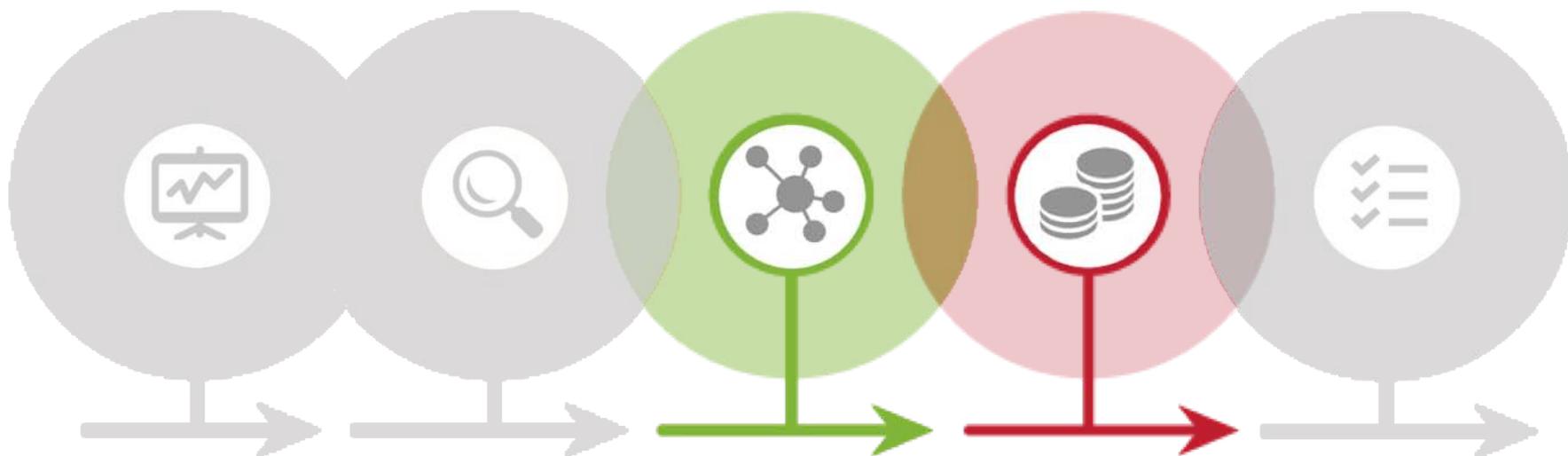
EVOLUZIONE E INNOVAZIONE DELLA METODOLOGIA

ROMA, 4 MARZO

RELATORI: ARIANNA CAMPAGNA E GIANCARLO FERRARA

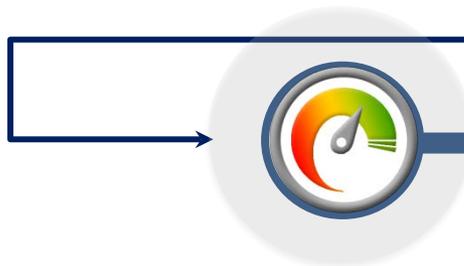
Il Processo metodologico





**3. INDIVIDUAZIONE
DEI GRUPPI
OMOGENEI
(ANALISI FATTORIALE
E CLUSTER ANALYSIS)**

**4. STIMA DEI
RICAVI
/COMPENSI, DEL
VALORE
AGGIUNTO E
DEL REDDITO**



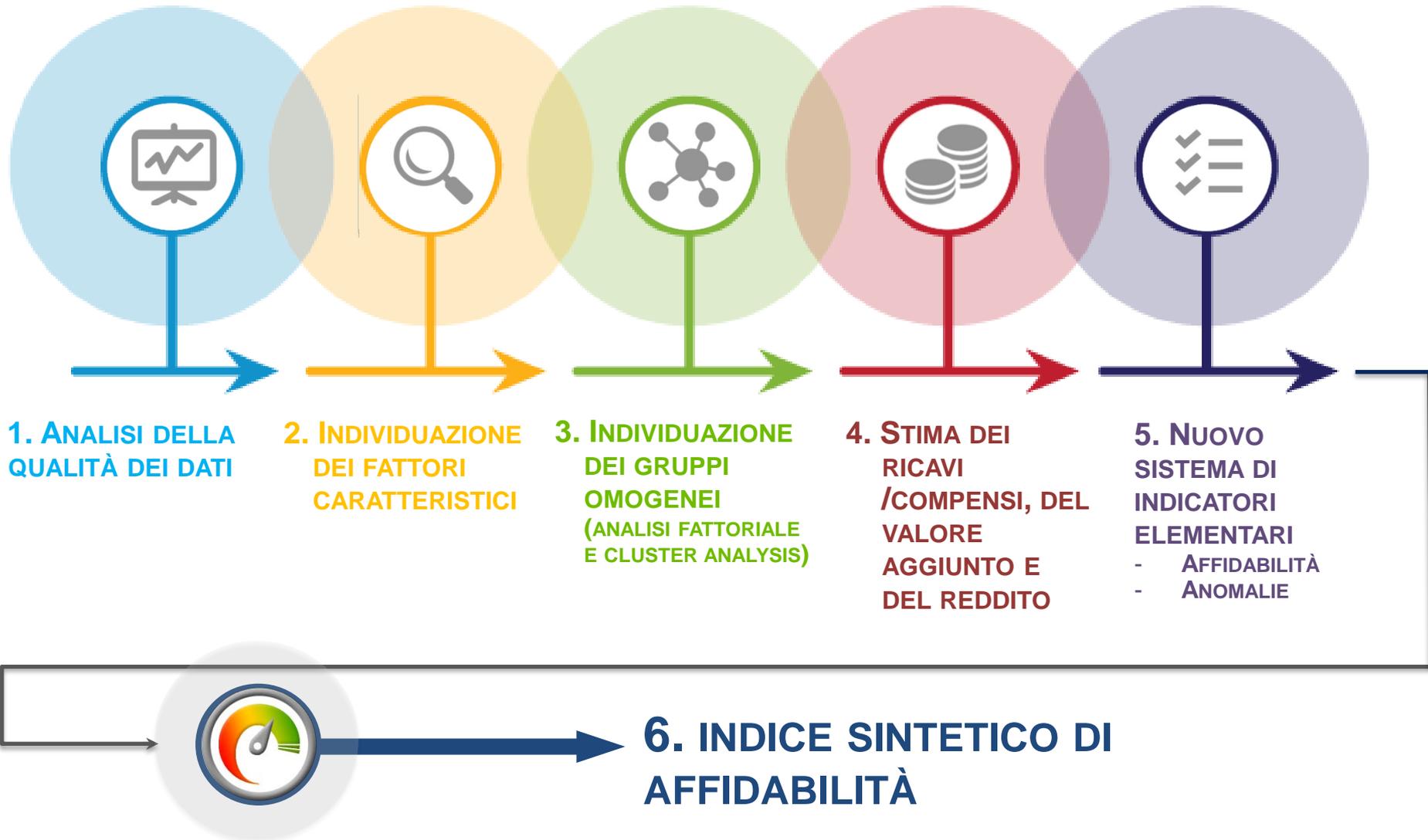
**6. INDICE SINTETICO DI
AFFIDABILITÀ**

Obiettivo

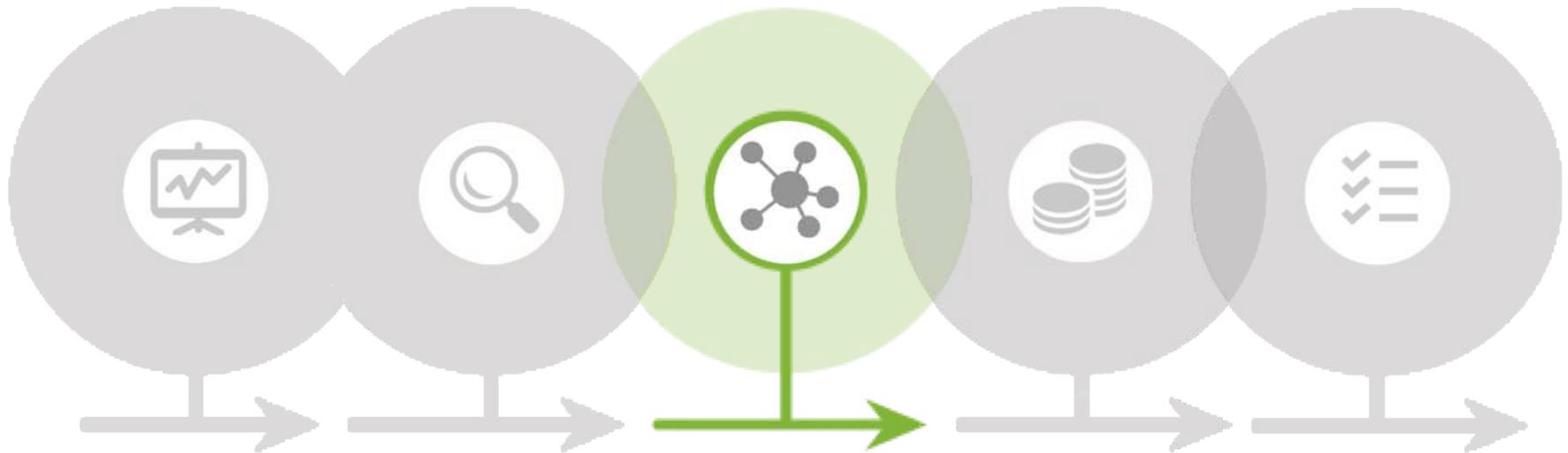
- ✓ **Misurare l'affidabilità del contribuente**
- 
- ✓ **Definizione di un valore di riferimento «normale» delle basi imponibili**
- 
- ✓ **Stima delle basi imponibili**

- ✓ Approccio economico alla base della **analisi dei gruppi omogenei**
- ✓ Approccio statistico nella stima dei gruppi
- ✓ Approccio econometrico nella definizione delle **funzioni di stima**

Il Processo metodologico



Cluster analysis



**3. INDIVIDUAZIONE
DEI GRUPPI
OMOGENEI
(ANALISI FATTORIALE
E CLUSTER ANALYSIS)**



**6. INDICE SINTETICO DI
AFFIDABILITÀ**

Analisi dei gruppi: perché?

Eterogeneità indotta dal contesto



**Contesto misurato attraverso lo studio delle
modalità organizzative**



Gruppi omogenei (cluster)



Variabilità non spiegata sintetizzata dalla variabile cluster

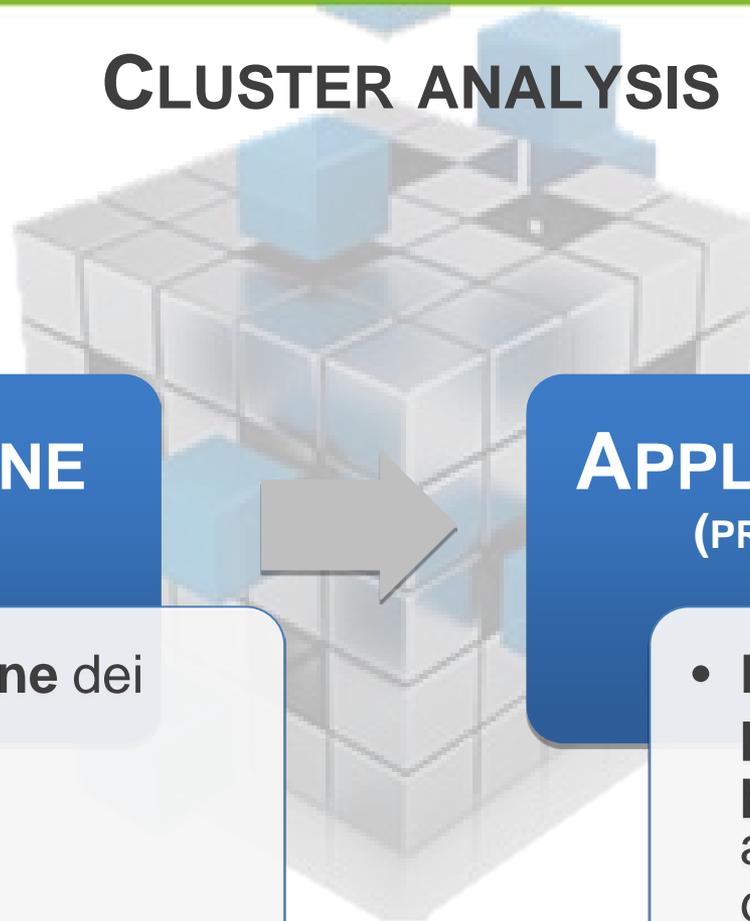
CLUSTER ANALYSIS

COSTRUZIONE
(STIMA)

- **Definizione dei gruppi**

APPLICAZIONE
(PREDIZIONE)

- **Definizione delle probabilità a posteriori di appartenenza ai gruppi**



IDENTIFICAZIONE



Con riferimento agli aspetti statistici associati alla teoria economica sottostante, i diversi blocchi logici alla base del metamodello vengono analizzati attraverso

single Analisi fattoriali,

in modo da individuare coerentemente le diverse dimensioni di analisi utilizzate per la definizione dei gruppi.

“BLOCCHI COSTITUTIVI”

A

B

C

D

ANALISI FATTORIALE

A

$X_{1A} X_{2A} X_{3A} X_{4A} X_{5A}$

B

$X_{1B} X_{2B}$

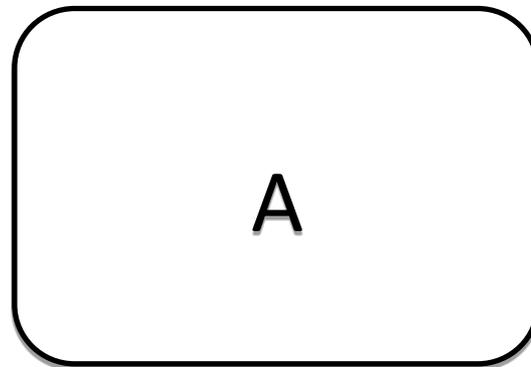
C

$X_{1C} X_{2C} X_{3C}$

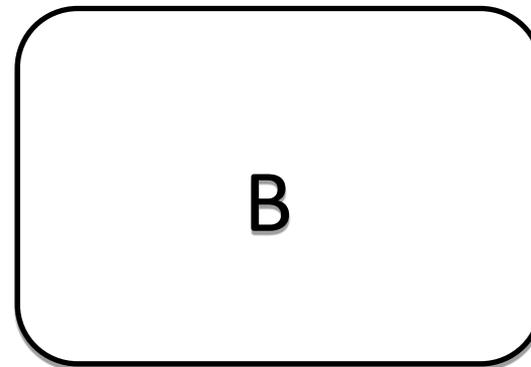
D

$X_{1D} X_{2D} X_{3D} X_{4D}$

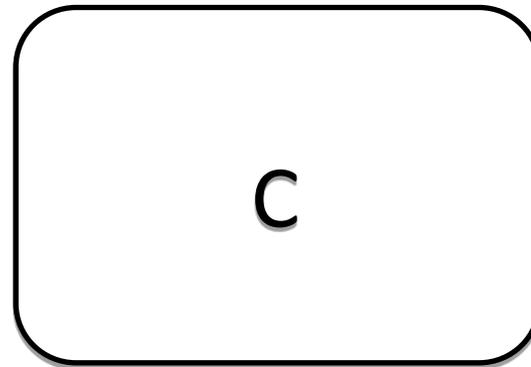
ANALISI FATTORIALE



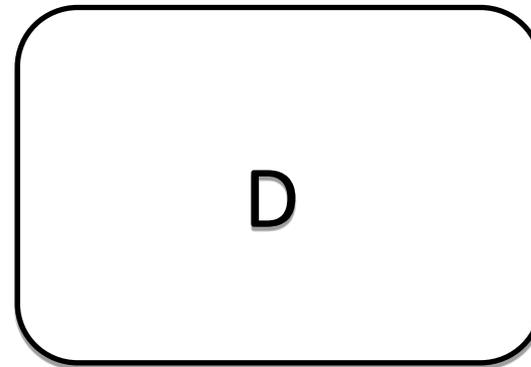
$f_{1A} f_{2A} f_{3A}$



f_{1B}

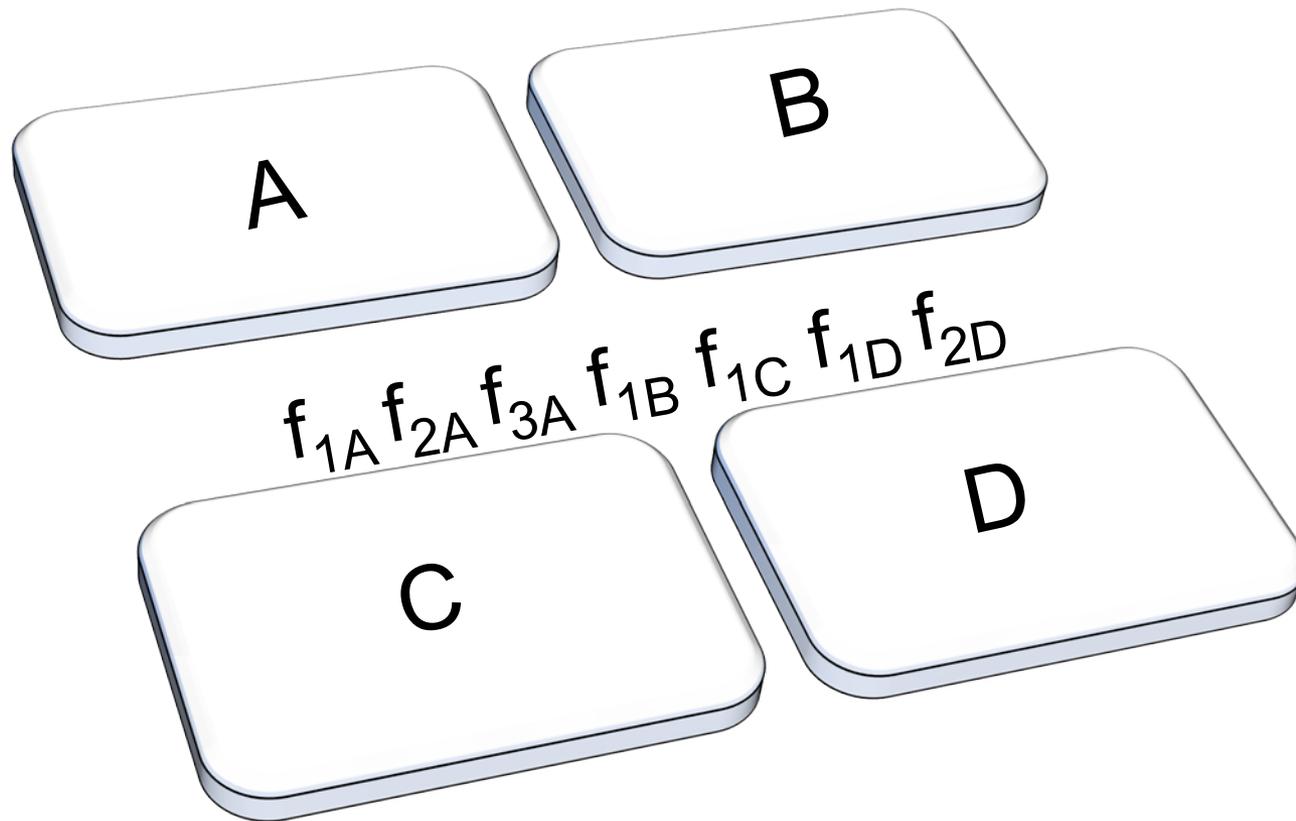


f_{1C}



$f_{1D} f_{2D}$

ANALISI FATTORIALE



$$f(\mathbf{x}) = \sum_{g=1}^G p_g f_g(\mathbf{x})$$

dove:

- $f_g(\mathbf{x})$ sono funzioni (di densità) di probabilità (componenti \equiv sottopopolazioni) che possono assumere differenti forme analitiche (normale, Poisson, esponenziale, etc.);
- p_g (pesi o probabilità *a priori*).

$f(\mathbf{x})$ è una funzione di densità, ovvero non negativa e con integrale pari a 1 (Everitt, 1981).

Una osservazione \mathbf{x} viene classificata nelle G sottopopolazioni attraverso il calcolo delle probabilità a posteriori (McLachlan and Peel, 2000):

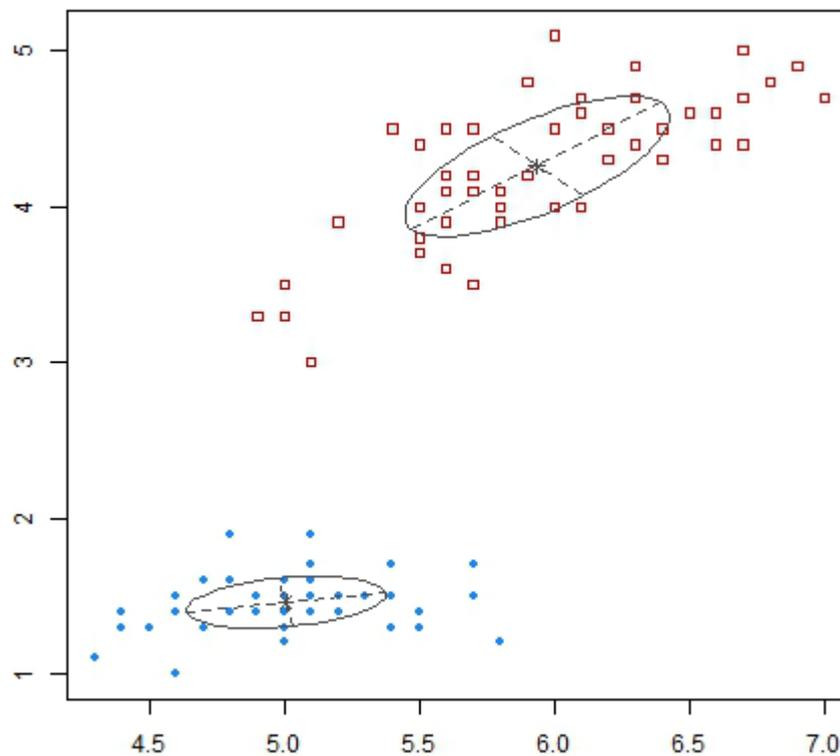
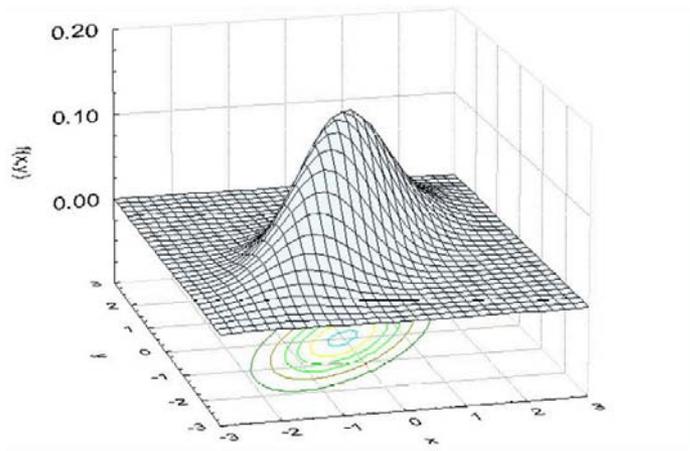
$$p(g|\mathbf{x}) = \frac{p_g f_g(\mathbf{x})}{\sum_{h=1}^G p_h f_h(\mathbf{x})}$$

$$f(\mathbf{x}) = \sum_{g=1}^G p_g \phi(\mathbf{x}; \boldsymbol{\mu}_g, \boldsymbol{\Sigma}_g)$$

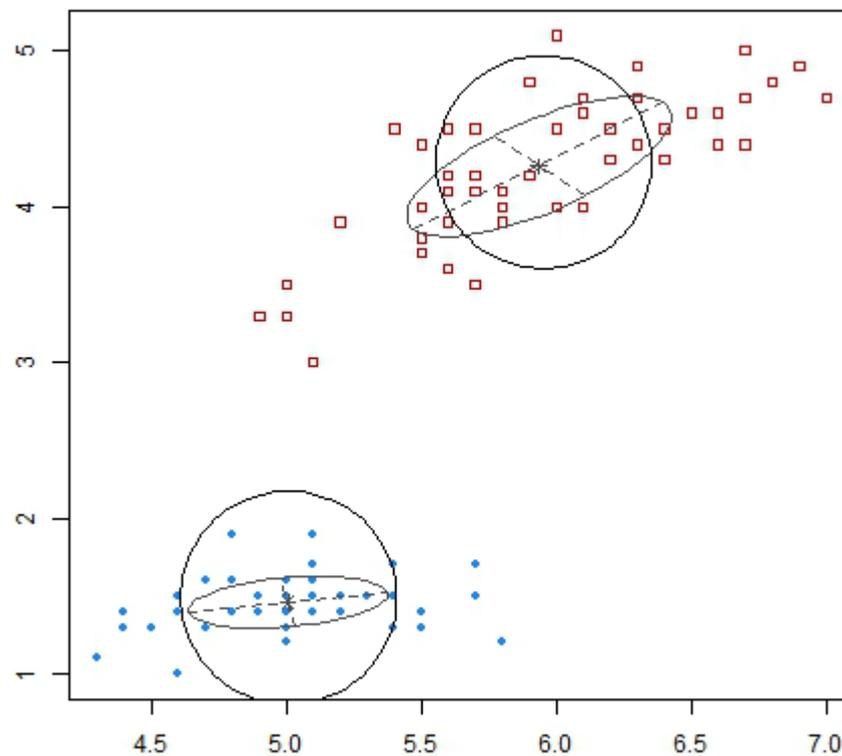
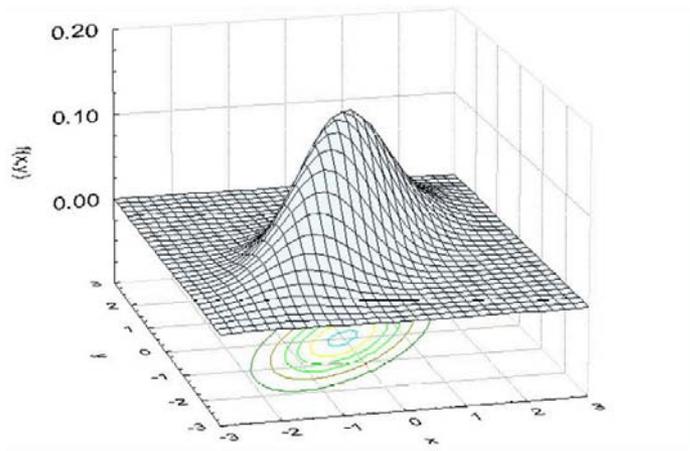
dove:

- $\phi(\mathbf{x})$ è la funzione di densità normale
- p_g sono le probabilità *a priori*
- $\boldsymbol{\mu}_g$ medie delle componenti
- $\boldsymbol{\Sigma}_g$ varianze/covarianze delle componenti

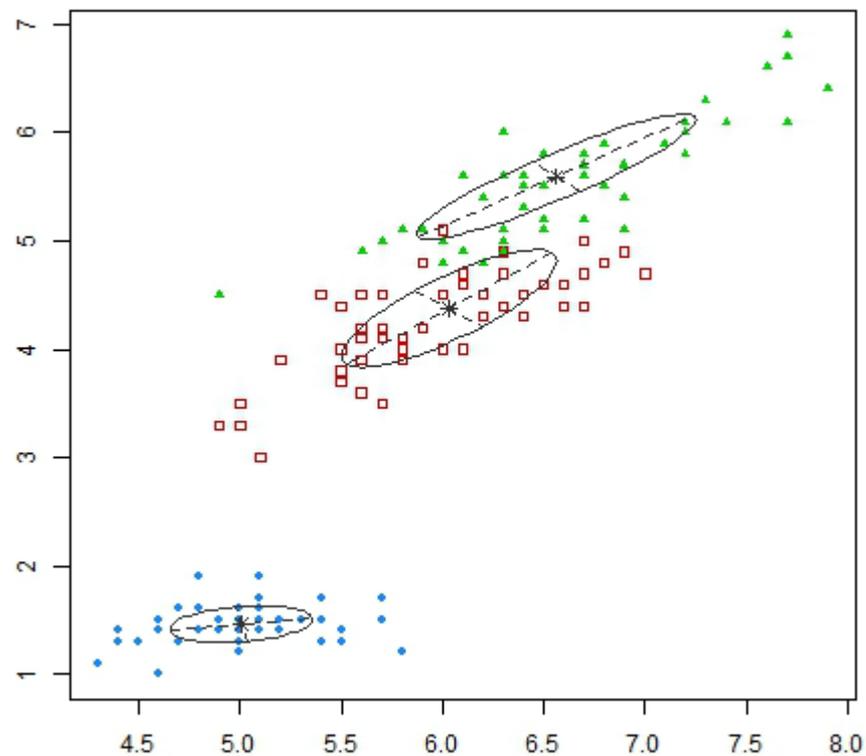
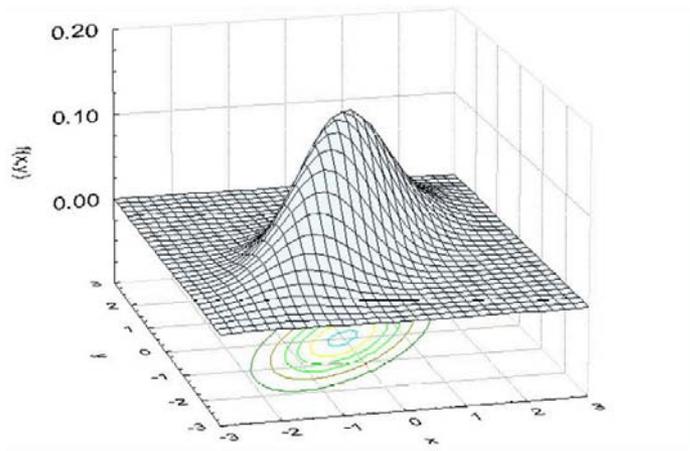
INTUITIVAMENTE ...



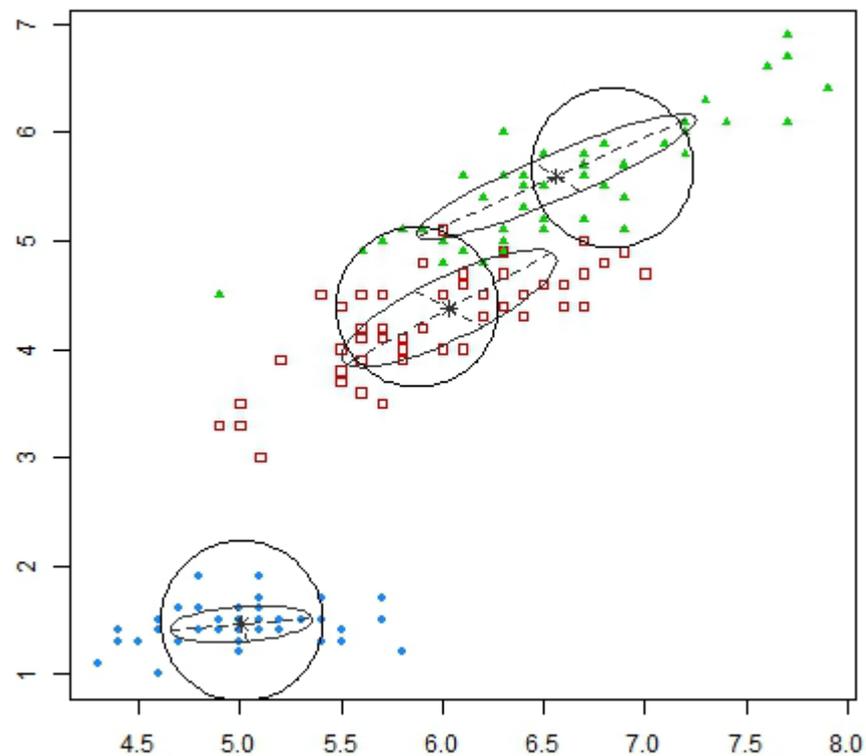
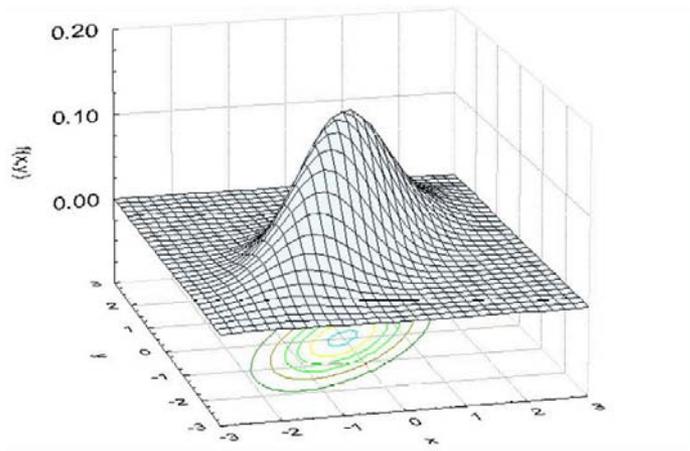
... INTUITIVAMENTE ...



... INTUITIVAMENTE ...



... INTUITIVAMENTE ...



Per tenere sotto controllo la eventuale presenza di noise e/o outlier può essere aggiunta nella mistura finita una ulteriore componente che rappresenta il noise.

Fraley *et al.* (2012) specificano tale componente attraverso un processo di *Poisson* di primo ordine:

$$f(\mathbf{x}) = p_0 \frac{1}{V} + \sum_{g=1}^G p_g \phi(\mathbf{x}; \boldsymbol{\mu}_g, \boldsymbol{\Sigma}_g)$$

dove:

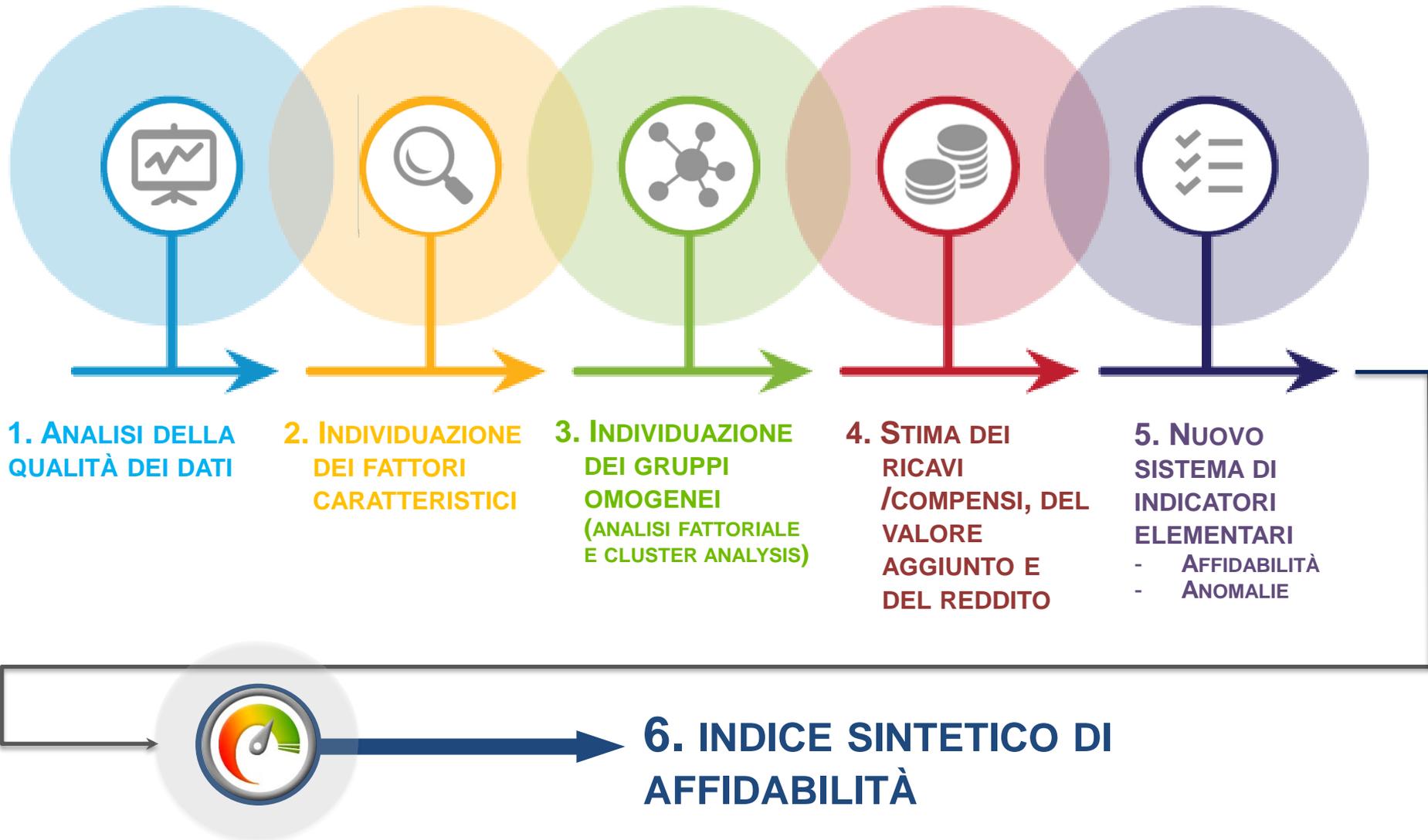
- $\phi(\mathbf{x})$ è la funzione di densità normale
- p_g sono le probabilità *a priori*
- V è un threshold tipicamente pari all'ipervolume dei dati



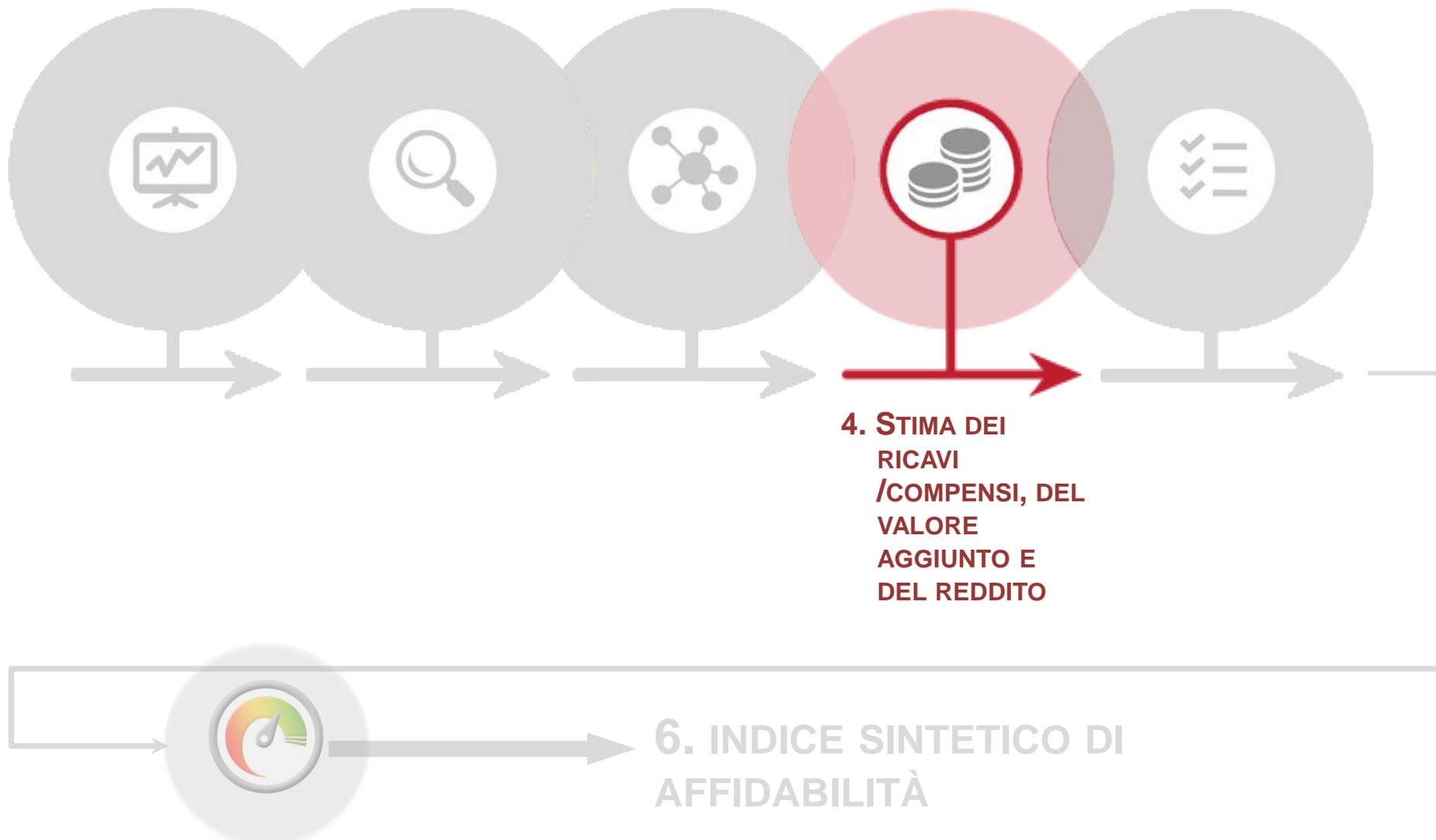
Metodologie a confronto

Obiettivo	Individuazione dei Modelli Organizzativi di Business (MOB) e stima delle probabilità di appartenenza del contribuente ai suddetti	
Ambito	Studi di Settore	Indice Sintetico di Affidabilità
Individuazione cluster (costruzione)	<ul style="list-style-type: none"> • Analisi fattoriale (ACP) • K-medie sui fattori 	<ul style="list-style-type: none"> • Analisi fattoriale (ACP) • Misture finite sui fattori
Stima delle probabilità (applicazione)	<ul style="list-style-type: none"> • Analisi discriminante su variabili costituenti i fattori 	
Caratteristica principale	L'assegnazione del contribuente ai gruppi realizzata in costruzione NON COINCIDE NECESSARIAMENTE con la probabilità massima di appartenere al medesimo in applicazione	La probabilità di appartenenza ai gruppi realizzata in costruzione COINCIDE ESATTAMENTE con quella in applicazione

Il Processo metodologico



Modelli predittivi



Considerazioni preliminari

- ✓ *trasformazione logaritmica*
- ✓ *fondamento microeconomico*
- ✓ *introduzione della costante*
- ✓ *da valori assoluti a valori pro-capite*
- ✓ *disegno longitudinale*

Modello di stima

$$y = \alpha + x\beta + \varepsilon$$

Specificazione lineare nei logaritmi.

Implicitamente si considera una funzione di produzione Cobb-Douglas:

$$y = \alpha x^\beta$$

Uno dei principali vantaggi derivanti da tale specificazione risiede nella relativa interpretazione dei parametri:

- α indica in che misura l'attività produttiva è efficiente, misura la scala di produzione;
- β misura l'elasticità dell'output prodotto rispetto alla corrispondente variabile

Specificazione pro-capite

Si è ritenuto opportuno adottare un modello che faccia riferimento all'input fondamentale utilizzato per la creazione del valore, ovvero il **fattore lavoro**.

In tal senso, il modello proposto prevede che la variabile risposta ed i corrispondenti regressori siano espressi in termini di rapporto rispetto al numero di addetti (i.e. *pro capite*).

I regressori costituiscono degli indicatori che forniscono una immediata lettura dal punto di vista economico.



L'analisi viene condotta a partire dalla banca dati degli studi di settore su un panel in modo da tener conto nella definizione della stima anche del **comportamento individuale** del contribuente nel tempo e dell'**andamento congiunturale**.



Struttura dei dati

Dati di Input

- variabili osservate per impresa/professionista :
 - costo del venduto per addetto
 - spese per lavoro dipendente per addetto
 - valore dei beni strumentali per addetto
 - ...
- gruppi omogenei
- misure economiche territoriali (OMI, reddito disponibile, ...)
- misure di ciclo economico

Osservazioni ripetute
per impresa/professionista

Dati di Output

- Predizione per impresa/professionista
 - ricavi
 - valore aggiunto

Pooled OLS

$$y_{it} = \alpha + x_{it}\beta + \varepsilon_{it}$$

Fixed effects

$$y_{it} = \alpha_i + x_{it}\beta + \varepsilon_{it}$$

Random effects

$$y_{it} = \alpha_i + x_{it}\beta + \varepsilon_{it}$$

$$\alpha_i = \alpha + w_i\vartheta + \gamma_i$$

Fixed vs Random

- Nel caso di FE l'inferenza è limitata al comportamento di un particolare insieme di unità.
- Nel caso di modelli ad effetti random, gli stimatori OLS sono ancora non distorti e consistenti, ma non più efficienti.
- Quando è vero il modello ad effetti fissi, gli stimatori OLS non sono più non distorti e consistenti. Tale risultato è conseguenza di un problema di omissioni di variabili dovute al fatto che nello stimatore OLS non sono presenti le dummy per ciascuna unità di analisi.
- La specificazione RE è parsimoniosa, consente di modellizzare l'effetto legato alle singole unità attraverso un solo parametro.
- FE least squares, anche noti come least squares dummy variables (LSDV), prevedono la stima di $(n - 1)$ extra parameteri, pertanto un numero elevato di dummy può aggravare il problema della multicollinearità fra regressori.
- Lo stimatore FE non permette di stimare l'effetto di variabili time-invariant.

Fixed vs Random

Differenza principale

RE assume che la distribuzione degli effetti random non dipende dai regressori mentre FE non richiede tale ipotesi.

In molte applicazioni l'ipotesi di indipendenza può essere violata perché gli effetti random sono legati alla presenza di variabili time-invariant che sono state omesse dal modello. L'ipotesi di indipendenza è equivalentemente ad assumere che le variabili omesse time-invariant sono indipendenti dai regressori.

È possibile verificare se sia presente una correlazione, ed eventualmente continuare a mantenere una specificazione RE, includendo le medie delle variabili esplicative per gruppo/livello (Mundlak, 1978; Snijders and Bosker, 1999; Hsiao, 2014).

$$\left. \begin{array}{l} \mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\boldsymbol{\gamma} + \boldsymbol{\varepsilon} \\ \boldsymbol{\gamma} \sim N(\mathbf{0}, \mathbf{G}) \\ \boldsymbol{\varepsilon} \sim N(\mathbf{0}, \mathbf{R}) \\ \boldsymbol{\gamma} \perp \boldsymbol{\varepsilon} \end{array} \right\} \mathbf{y} | \mathbf{X} \sim N(\mathbf{X}\boldsymbol{\beta}, \mathbf{Z}\mathbf{G}\mathbf{Z}' + \mathbf{R})$$

dove:

- \mathbf{y} il vettore delle osservazioni della variabile risposta;
- $\boldsymbol{\beta}$ vettore degli effetti fissi;
- $\boldsymbol{\gamma}$ vettore degli effetti casuali;
- \mathbf{X} matrice del disegno degli effetti fissi;
- \mathbf{Z} matrice del disegno degli effetti casuali.

Specificazione degli effetti

$$y_{it} = \alpha + \mathbf{x}_{it}\boldsymbol{\beta} + \bar{\mathbf{x}}_i \cdot \boldsymbol{\beta}^* + \mathbf{w}_t\boldsymbol{\eta} + \\ + \delta_1 clu_{it}^{(1)} + \dots + \delta_{G-1} clu_{it}^{(G-1)} + \gamma_i + \varepsilon_{it}$$

Si è tenuto conto con un'unica funzione delle possibili differenze di risultati economici riconducibili agli aspetti territoriali congiuntamente alle diverse caratteristiche strutturali ed organizzative definite dalla *Cluster Analysis*.

Effetto **casuale** definito dall'impresa stessa, analizzando in tal modo il comportamento individuale.

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\hat{\mathbf{V}}^{-1}\mathbf{X})^{-1}\mathbf{X}'\hat{\mathbf{V}}^{-1}\mathbf{y}$$

Henderson (1984)

$$\hat{\boldsymbol{\gamma}} = \hat{\mathbf{G}}\mathbf{Z}'\hat{\mathbf{V}}^{-1}(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})$$

$$\hat{\mathbf{V}} = \mathbf{Z}\hat{\mathbf{G}}\mathbf{Z}' + \hat{\mathbf{R}}$$

$$\hat{\mathbf{G}} = \hat{\sigma}_G^2 \mathbf{I}_n$$

$$\hat{\mathbf{R}} = \hat{\sigma}_R^2 \mathbf{I}_{nT}$$

Predizione effetto individuale

$$\hat{\beta} = (\mathbf{X}'\hat{\mathbf{V}}^{-1}\mathbf{X})^{-1}\mathbf{X}'\hat{\mathbf{V}}^{-1}\mathbf{y}$$

Henderson (1984)

$$\hat{\mathbf{y}} = \hat{\mathbf{G}}\mathbf{Z}'\hat{\mathbf{V}}^{-1}(\mathbf{y} - \mathbf{X}\hat{\beta})$$

$$\hat{y}_i = \hat{\mathbf{G}}\mathbf{Z}'_i \hat{\mathbf{V}}_i^{-1}(\mathbf{y}_i - \mathbf{X}_i\hat{\beta})$$

$i = 1, \dots, n$

Caso
particolare
T=1

$$\left\{ \begin{array}{l} \hat{\mathbf{V}}_i = \mathbf{Z}_i\hat{\mathbf{G}}\mathbf{Z}'_i + \hat{\mathbf{R}}_i = \hat{\sigma}_G^2 + \hat{\sigma}_R^2 \\ \hat{y}_i = \frac{\hat{\sigma}_G^2}{\hat{\sigma}_G^2 + \hat{\sigma}_R^2} (\mathbf{y}_i - \mathbf{x}_i\hat{\beta}) \end{array} \right.$$

Definizione del campione di stima

L'analisi viene condotta a partire dalla banca dati degli studi di settore in relazione ai periodi d'imposta 2007-2014 in modo da tener conto nella definizione della stima anche del comportamento individuale del contribuente nel tempo e dell'andamento congiunturale.

- Panel a partire dalle imprese presenti nell'ultimo periodo di imposta, utilizzato nella fase di costruzione dei gruppi omogenei, ed in maniera **retrospettiva** viene valutata la presenza dei contribuenti fino all'anno di imposta 2007
- Panel «completo» costruito a partire da tutte le imprese presenti nei periodi di imposta osservati



SAS/HPA Procedure Highlights

- **PROC HPLMIXED**
 - High-performance version of PROC MIXED
 - Not to be confused with HPMIXED procedure in SAS/STAT
 - Supports
 - » RANDOM statements
 - » REPEATED statement
 - » Covariance structures from PROC MIXED
 - Sparse MMEQs with > 40,000 unknowns
 - » Impossible in MIXED
 - » 12 hours in HPMIXED
 - » 3 minutes in HPLMIXED

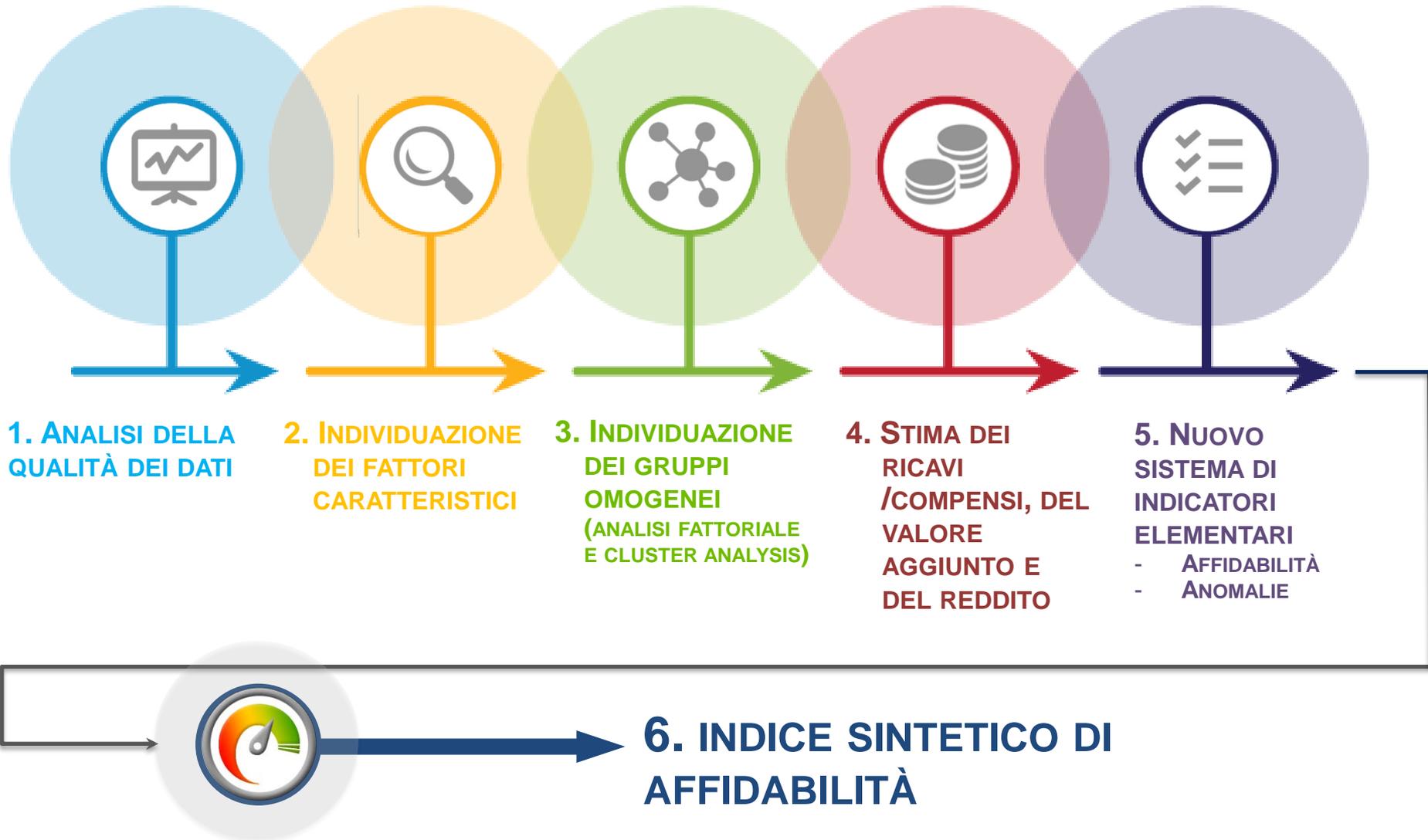
#SASGF11

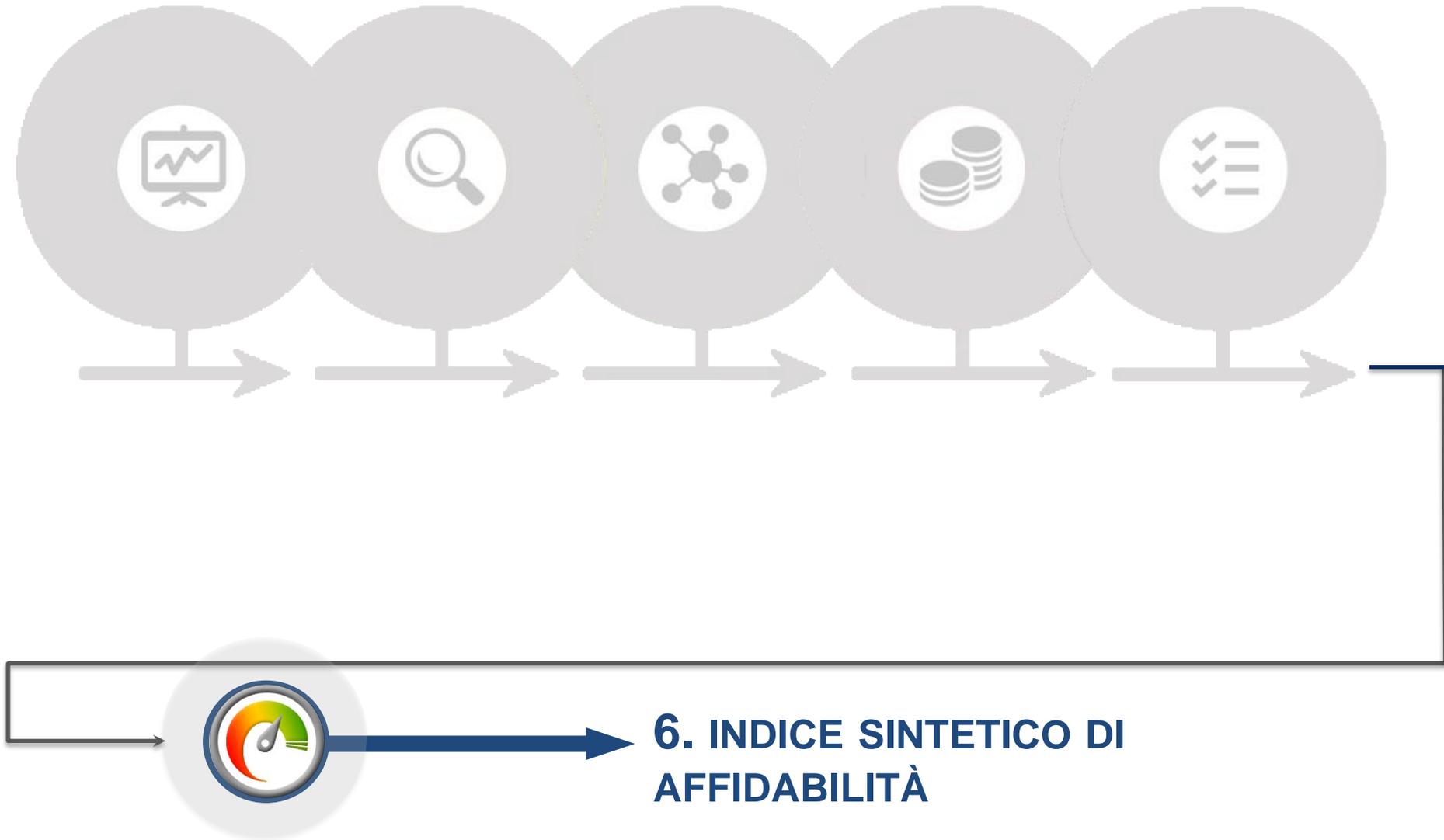


Metodologie a confronto

Obiettivo	Previsione delle basi imponibili in uno specifico periodo di imposta	
Ambito	Studi di Settore	Indice Sintetico di Affidabilità
Variabili obiettivo	Ricavo	Valore aggiunto/Ricavo per addetto (funzione di produzione)
Metodi di stima	<ul style="list-style-type: none"> • WLS (per gruppo) • indicatori di normalità economica • correttivi crisi 	<ul style="list-style-type: none"> • mixed models (multilevel approach)
Dati	cross-section	panel

Il Processo metodologico





L'INDICE SINTETICO DI AFFIDABILITÀ



RAPPRESENTA IL POSIZIONAMENTO DI OGNI CONTRIBUENTE RISPETTO ALL'AFFIDABILITÀ DEI SUOI COMPORTAMENTI FISCALI. E' UNA MEDIA SEMPLICE DI INDICATORI ELEMENTARI.

GLI INDICATORI ELEMENTARI PRENDONO IN CONSIDERAZIONE:

- ✓ la plausibilità dei ricavi/compensi, del valore aggiunto e del reddito
- ✓ l'affidabilità dei dati dichiarati
- ✓ le anomalie economiche



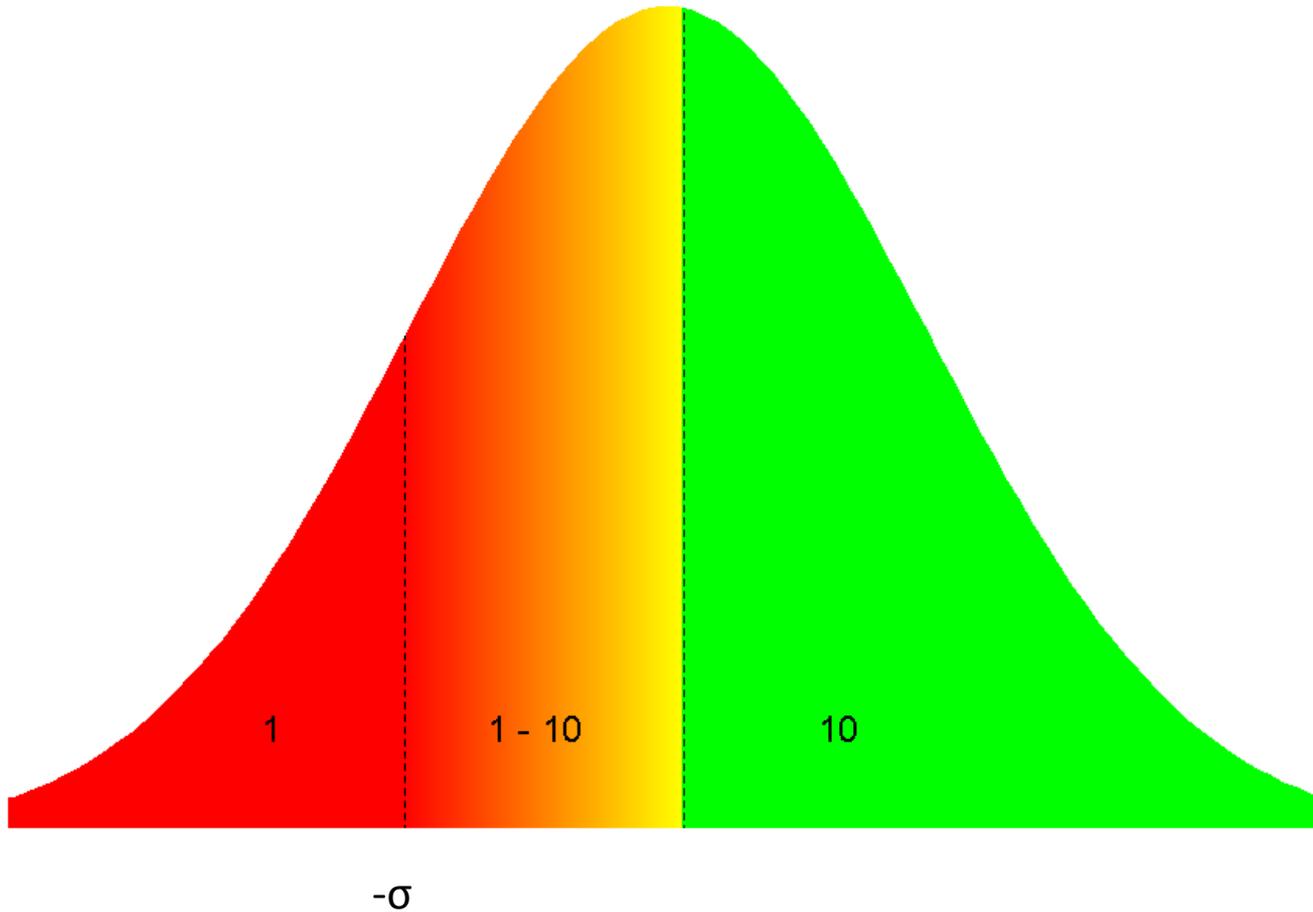
Il valore dell'indice sintetico assegnato al contribuente è compreso tra 1 e 10

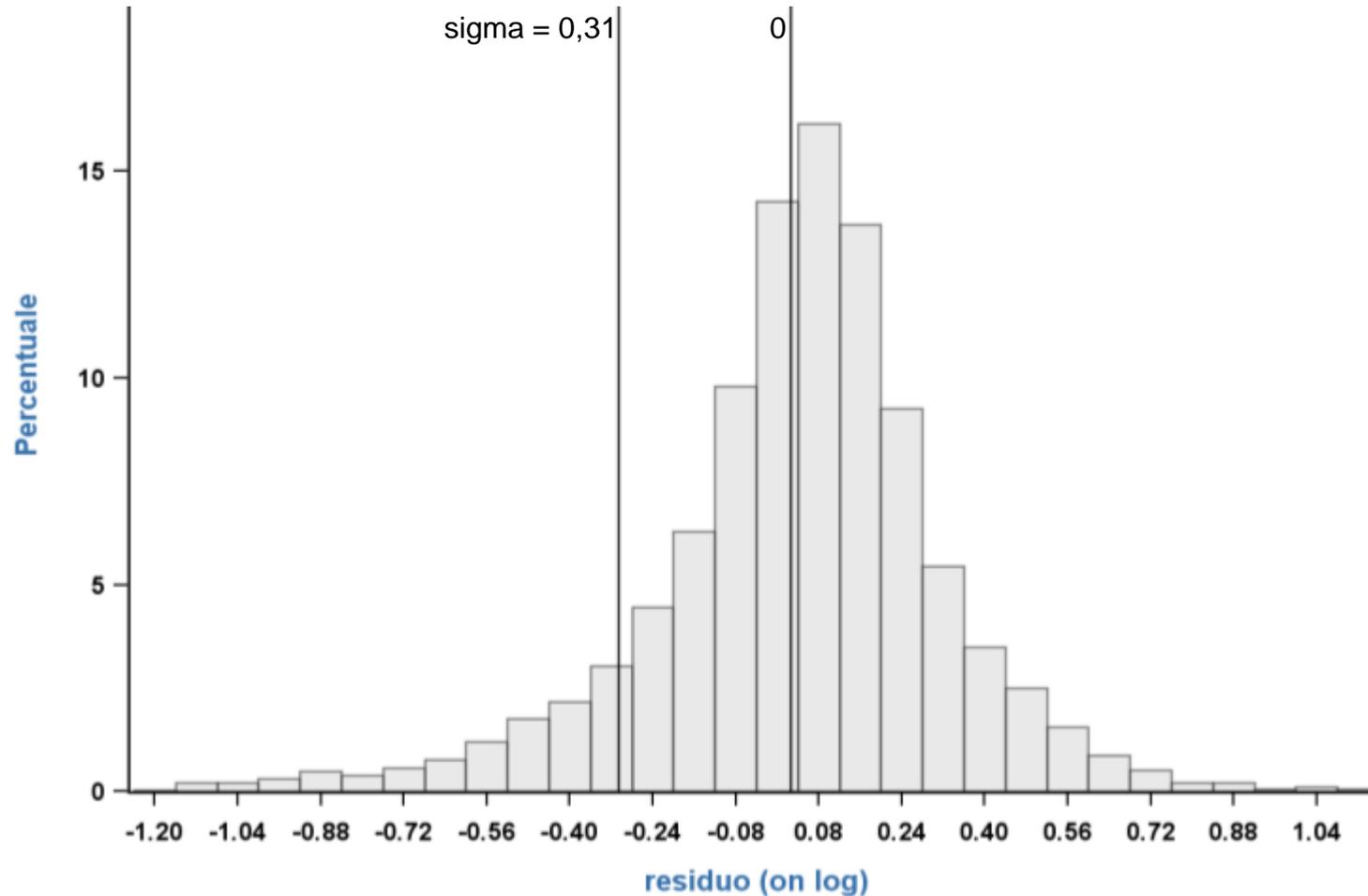


INDICATORI BASATI SU STIME

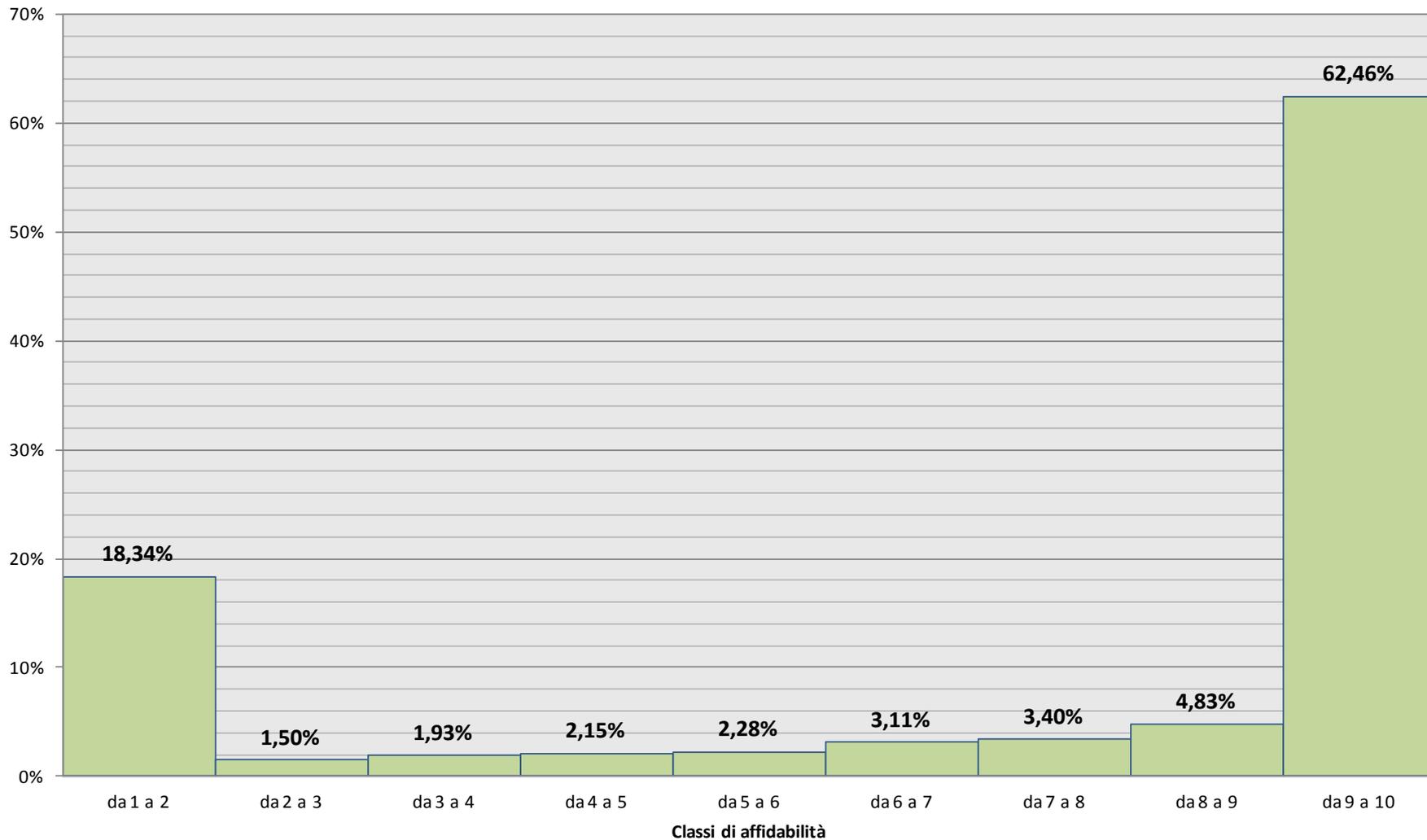
Calcolo punteggio teorico

$$\log\left(\frac{VA}{Add}\right) - \log\left(\frac{VA}{Add}\right)^*$$



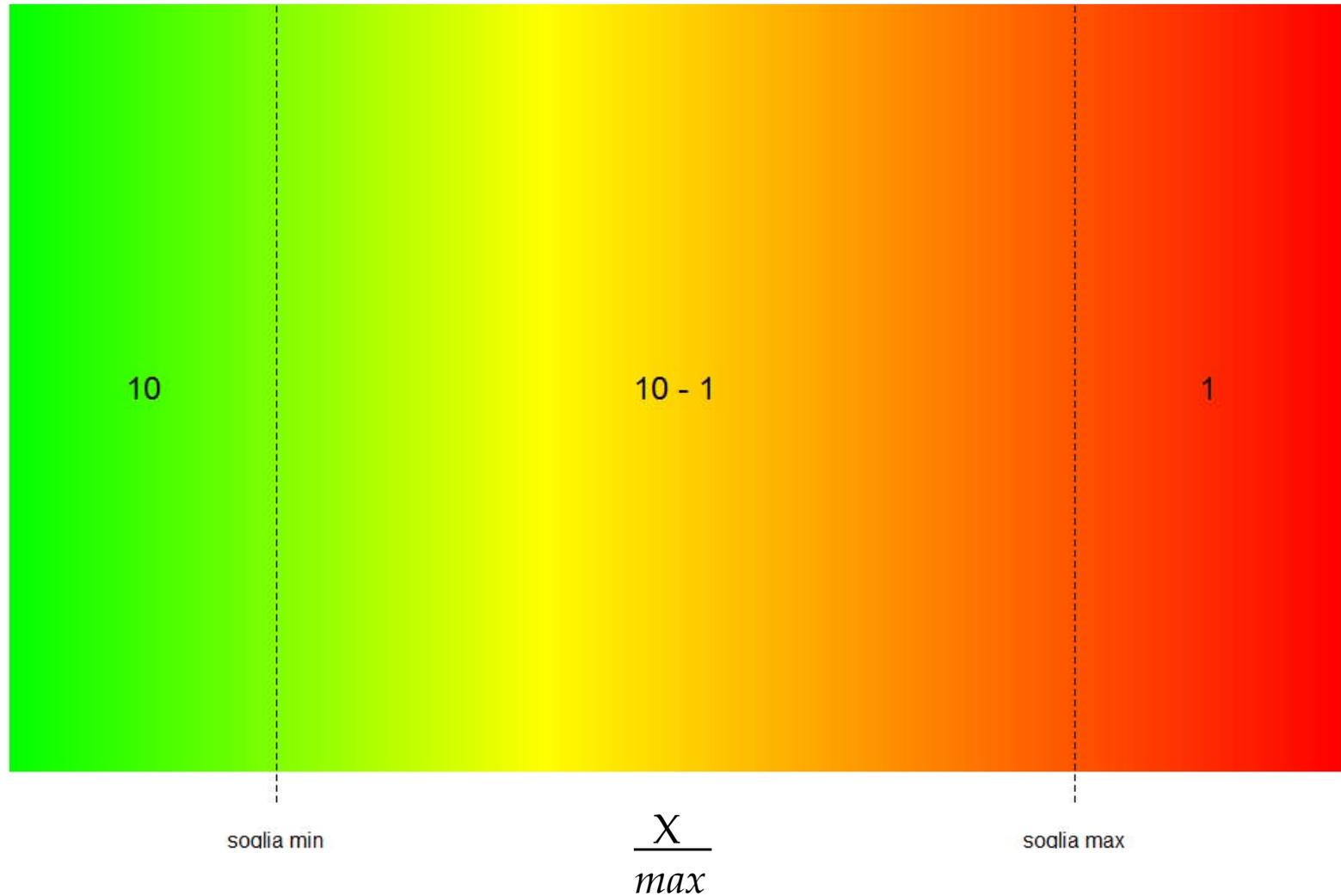


Affidabilità Valore aggiunto per addetto

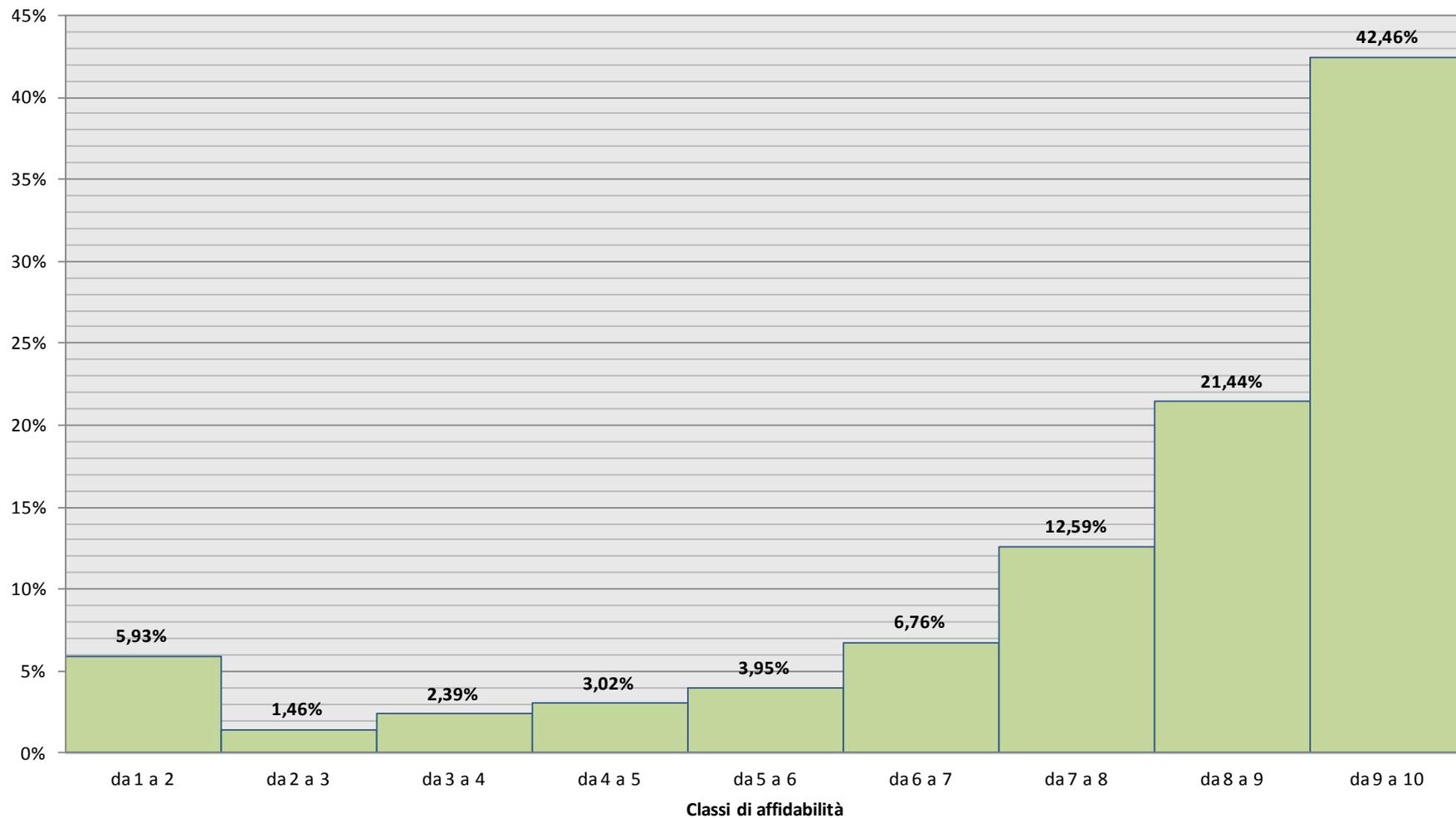


INDICATORI NON BASATI SU STIME

Calcolo punteggio teorico



Affidabilità Durata delle scorte



Giorni	526	468	409	351	292	234	175	117	58	50
Affidabilità	1	2	3	4	5	6	7	8	9	10

GRAZIE PER L'ATTENZIONE!

Arianna Campagna
acampagna@sose.it

Giancarlo Ferrara
gferrara@sose.it