

## Università degli Studi di Milano-Bicocca

Progetto Finale di *Data Visualization*

### “Just Repeat-Hit”

Misurare la ripetitività della musica  
tramite la compressione testuale

### Valutazione della qualità delle infografiche

---

**Borroni Alessandro** / mat. **800069** / [a.borroni2@campus.unimib.it](mailto:a.borroni2@campus.unimib.it)

**Giugliano Mirko** / mat. **800226** / [m.giugliano@campus.unimib.it](mailto:m.giugliano@campus.unimib.it)

**Prade Angela** / mat. **838540** / [a.prade@campus.unimib.it](mailto:a.prade@campus.unimib.it)

---

## Introduzione

Il progetto nasce dall'idea di verificare se la *ripetitività dei testi delle canzoni sia in aumento* con il passare del tempo. Dopo aver calcolato per ogni testo un indice di ripetitività, si è cercato di verificare dapprima se ci fosse un trend di ripetitività in crescita negli ultimi anni e successivamente si è cercato di verificare quale sia il genere o i generi che trainano questo trend. Una volta completato il lavoro di raccolta dati, di pre-processing ed elaborazione e analisi dei dati, sono state create delle infografiche e, successivamente, ne è stata valutata la qualità sottoponendole al parere degli utenti. Il software utilizzato per implementare le visualizzazioni è stato *Tableau*, mentre per i grafici riguardanti la valutazione delle infografiche è stato utilizzato *R*.

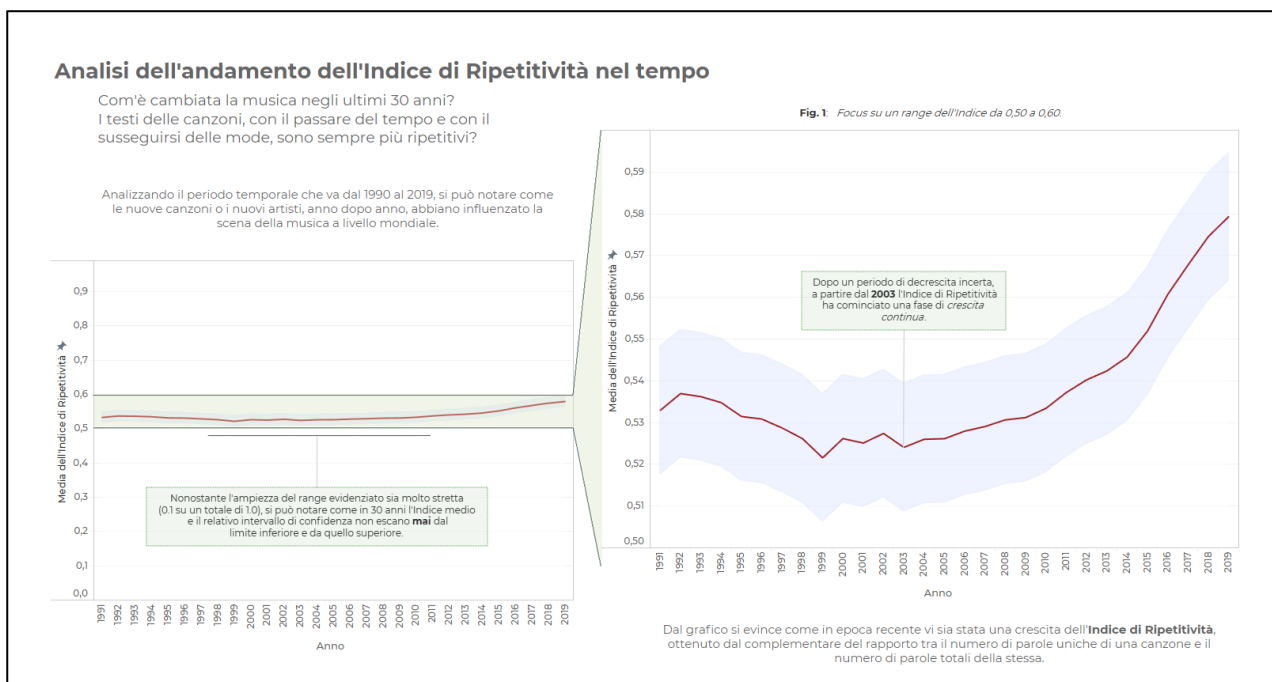
Sono state sviluppate, in totale, 5 infografiche:

- **Infografica 0** Introduzione e Descrizione dell'Indice di Ripetitività
- **Infografica 1** Analisi dell'andamento dell'Indice di Ripetitività nel tempo
- **Infografica 2** Indice di Ripetitività per Genere
- **Infografica 3** Distribuzione dell'Indice di Ripetitività
- **Infografica 4** Analisi della ripetitività delle singole canzoni per artista

## Think aloud

Prima di procedere con i questionari psicometrici e con i task da sottoporre agli utenti, è stato chiesto a *tre* persone di esprimere a voce alta il loro pensiero mentre interagivano con l'infografica in modo da poter comprendere gli aspetti della User Experience che potessero creare delle difficoltà.

### ◦ Infografica 1



Il problema che è emerso mostrando questa visualizzazione è il fatto che agli utenti non risultasse immediatamente chiaro che il grafico a destra rappresentasse le stesse informazioni del grafico a sinistra ma su una scala diversa. Per ovviare a questo problema, si è deciso di specificare in modo più visibile cosa rappresentasse il grafico a destra e di evidenziare tramite un segmento colorato la porzione di asse y che assume i valori da 0.5 a 0.6 del line chart a sinistra e l'asse del line chart a destra. Inoltre, è stata aggiunta nel line chart di sinistra una linea indicante l'Indice medio di Ripetitività misurato su 280 articoli di giornale e nel line chart l'indicazione dei valori di massimo e di minimo assunti dall'indice stesso. Il risultato ottenuto dopo le modifiche è mostrato nella figura successiva.

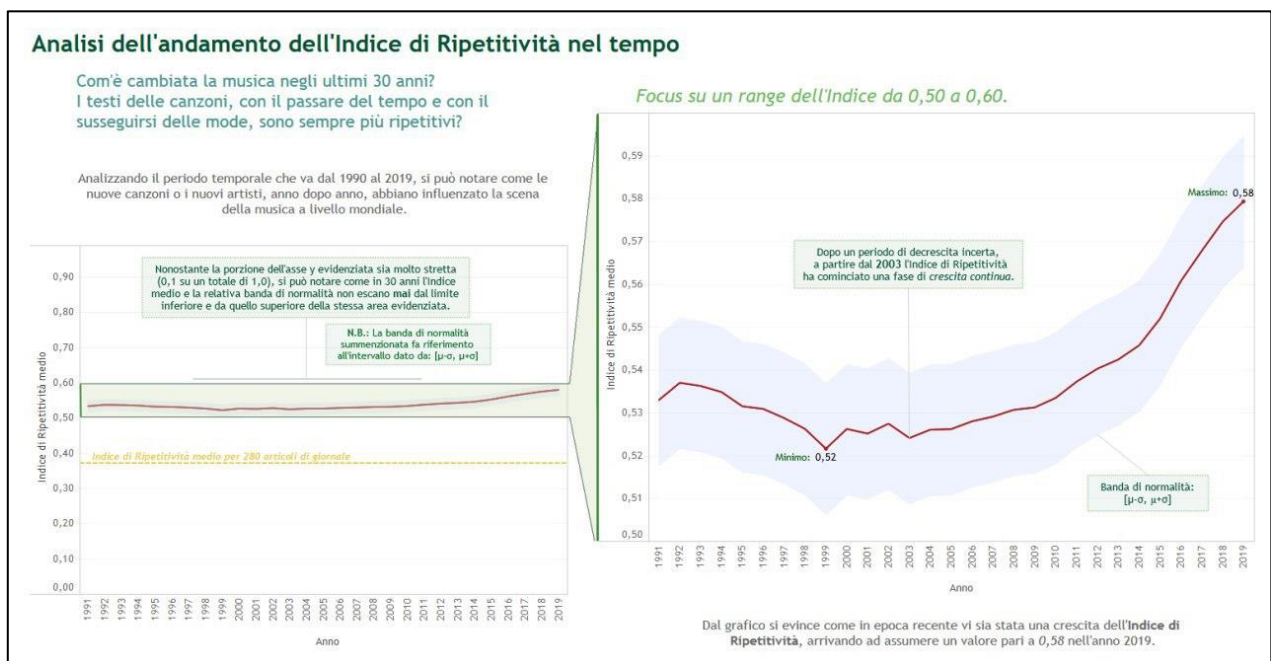
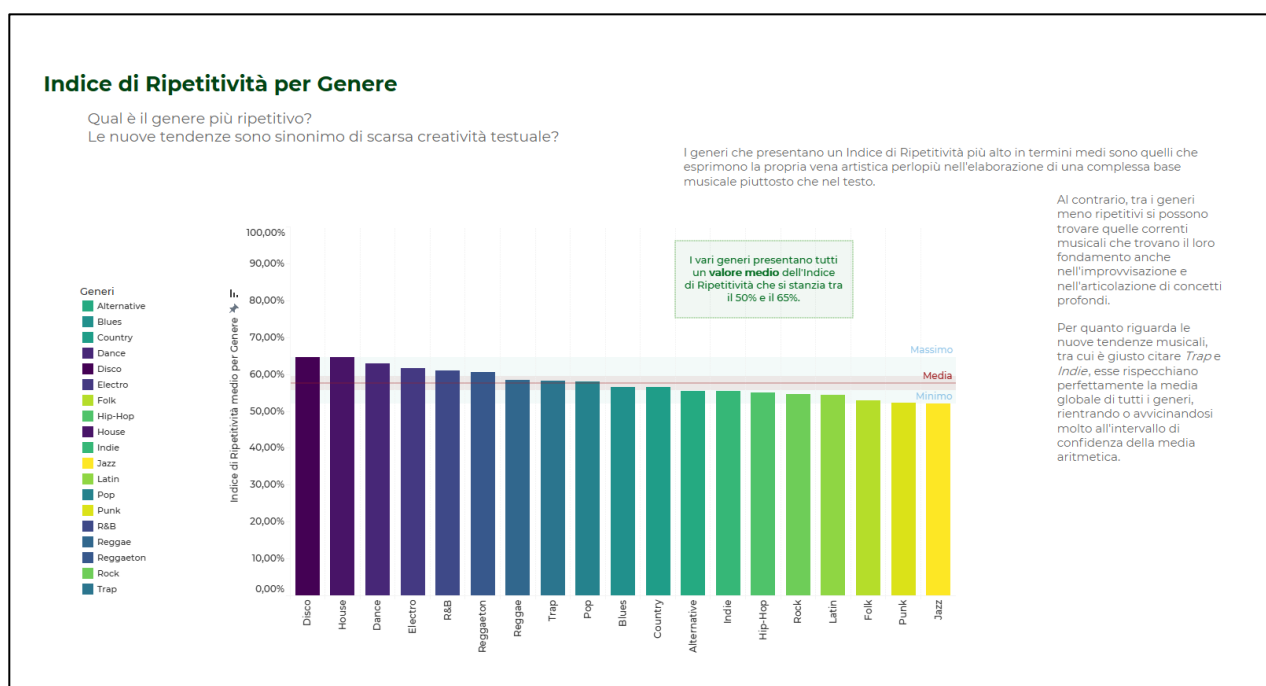
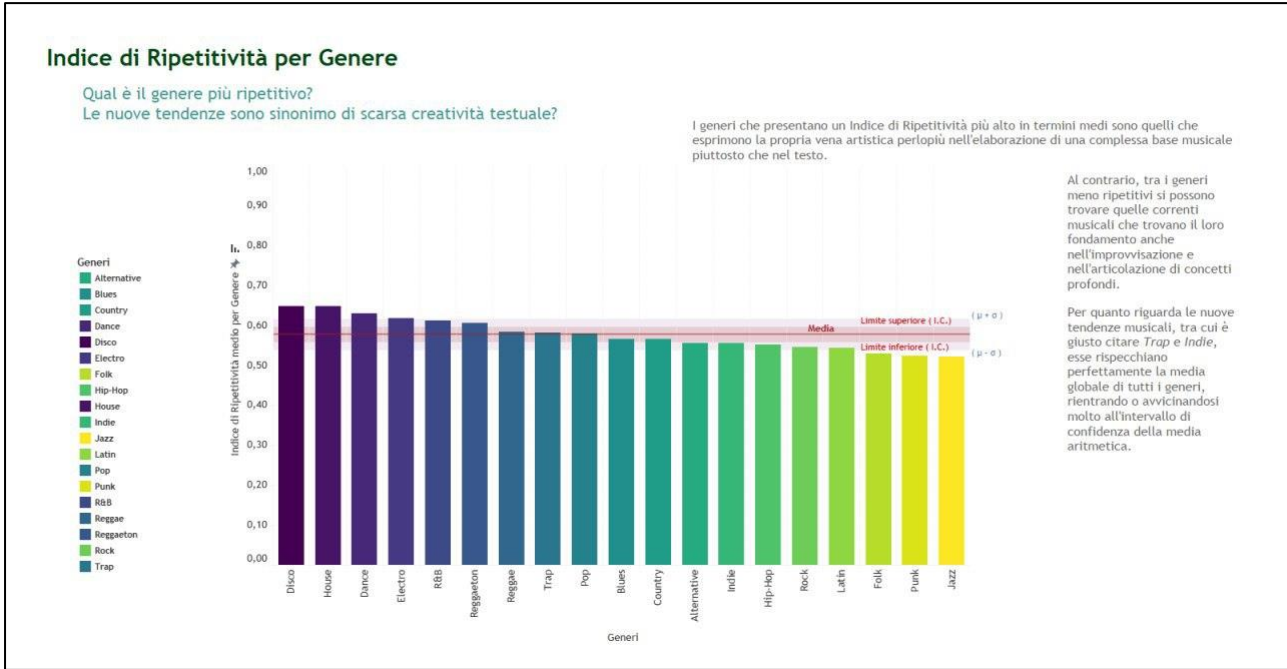


Figura 1

## Infografica 2



L'approccio degli utenti con questa Infografica ha evidenziato come l'indicazione di "Massimo" e "Minimo" sia stata ritenuta non rilevante. Inoltre, la banda di colore rosa intorno alla media ha creato confusione. Si è preferito quindi eliminare l'indicazione di "Massimo" e "Minimo", indicare l'utilità della banda rosa e aggiungere la banda di normalità data da  $[\mu - \sigma, \mu + \sigma]$ .



Successivamente sono stati aggiunti gli intervalli di confidenza per ogni singolo genere e, al fine di renderli più evidenti agli utenti, si è deciso di troncare l'asse y evidenziando solo i valori appartenenti al range [0.30, 0.80].

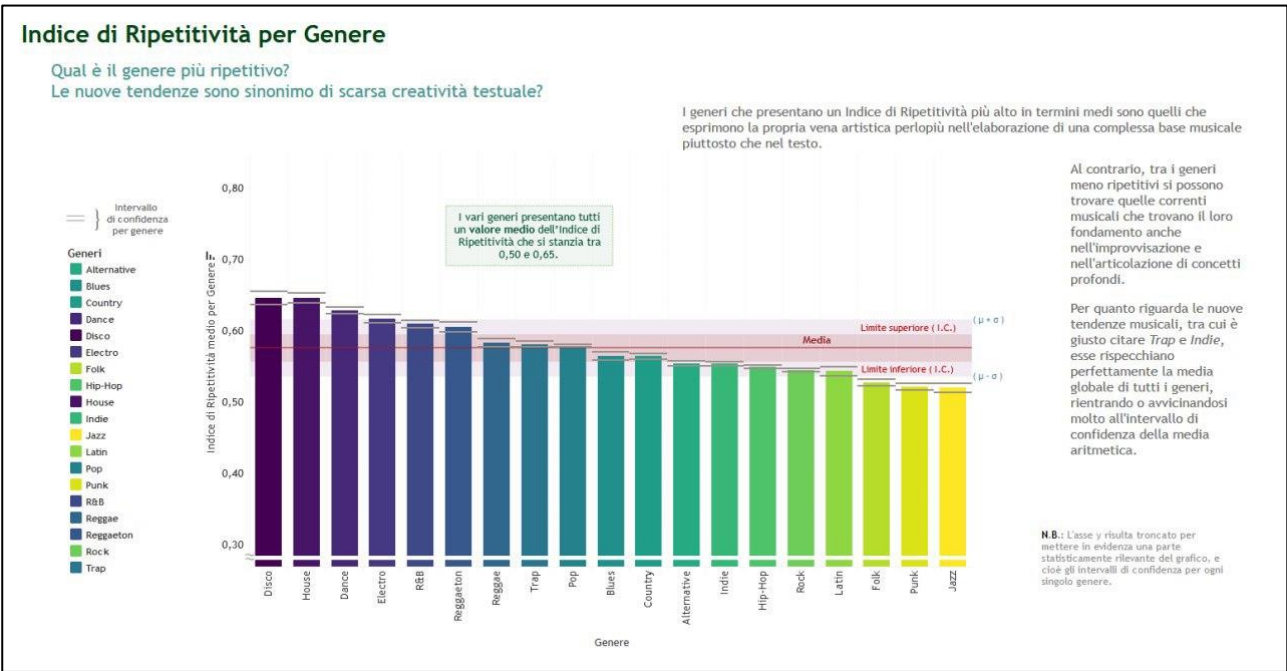
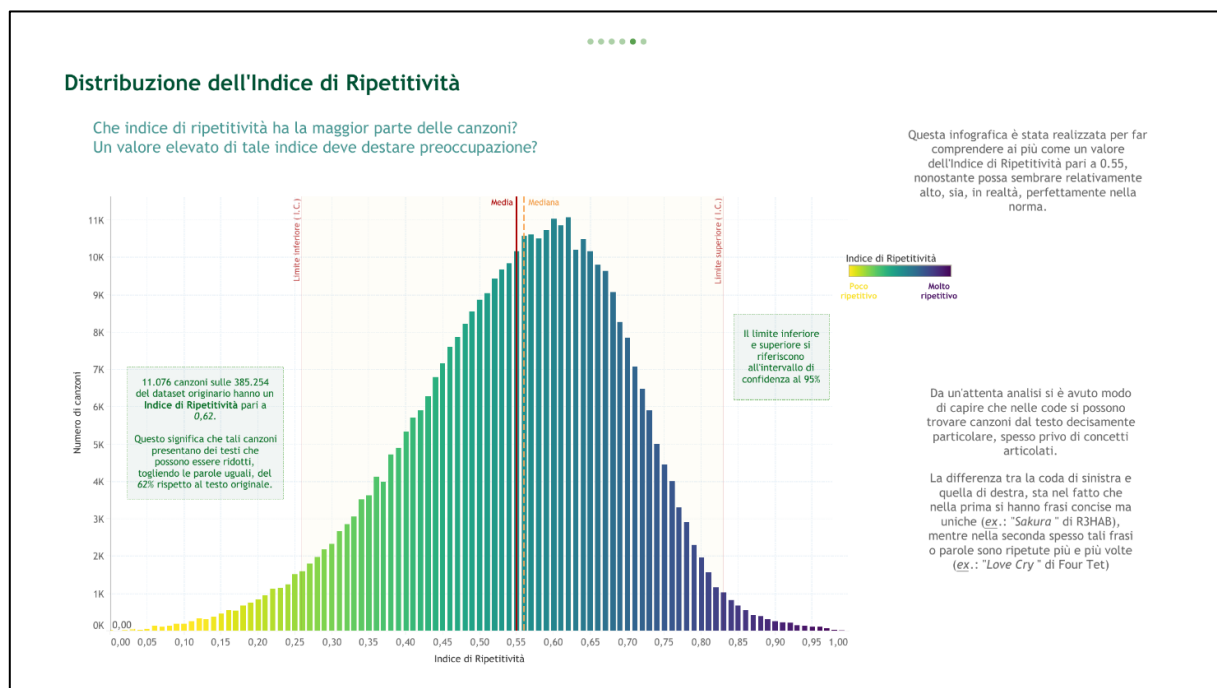


Figura 2

## Infografica 3



Inizialmente si è deciso di calcolare l'intervallo di confidenza al 95% intorno alla media e di riportarlo sul grafico. Utenti più esperti che hanno interagito con l'infografica ci hanno fatto notare la stranezza del valore di tale intervallo portandoci a individuare un errore nel calcolo.

Visto che l'intervallo corretto sarebbe stato eccessivamente stretto e non avrebbe portato ad alcuna informazione aggiuntiva, si è scelto di indicare la banda di normalità  $[\mu - \sigma, \mu + \sigma]$ . Il risultato è quello riportato nella figura successiva.

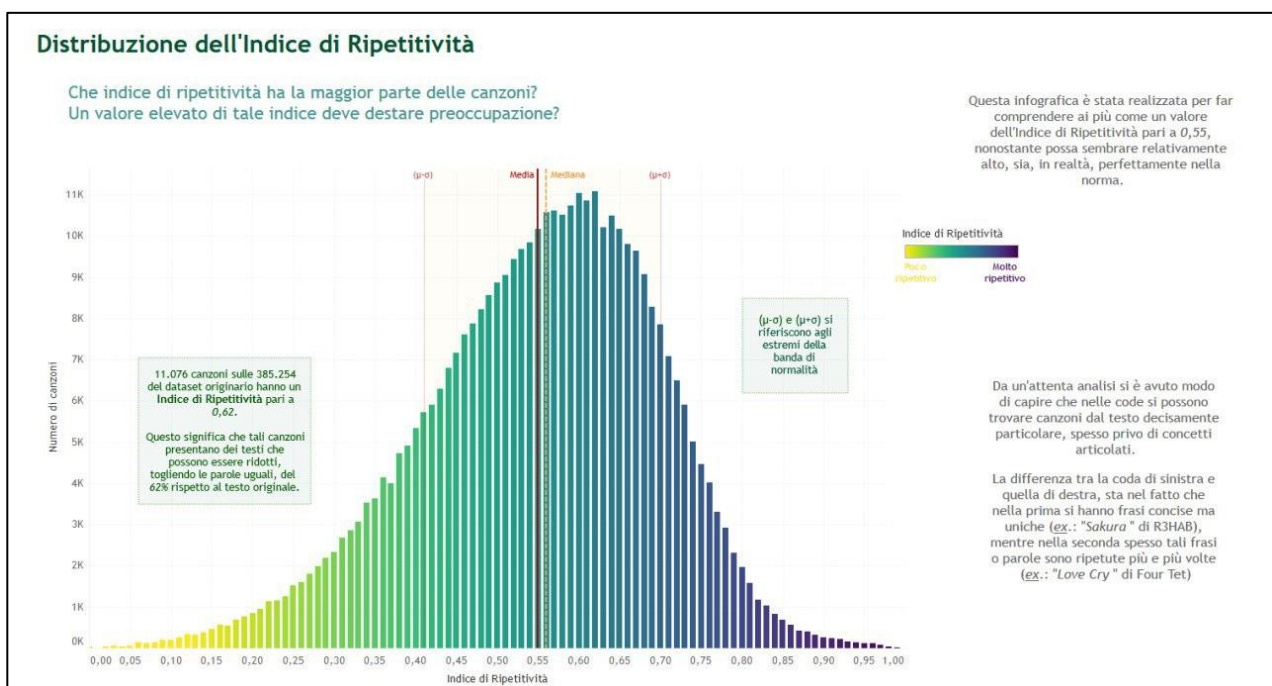


Figura 3



## Infografica 4

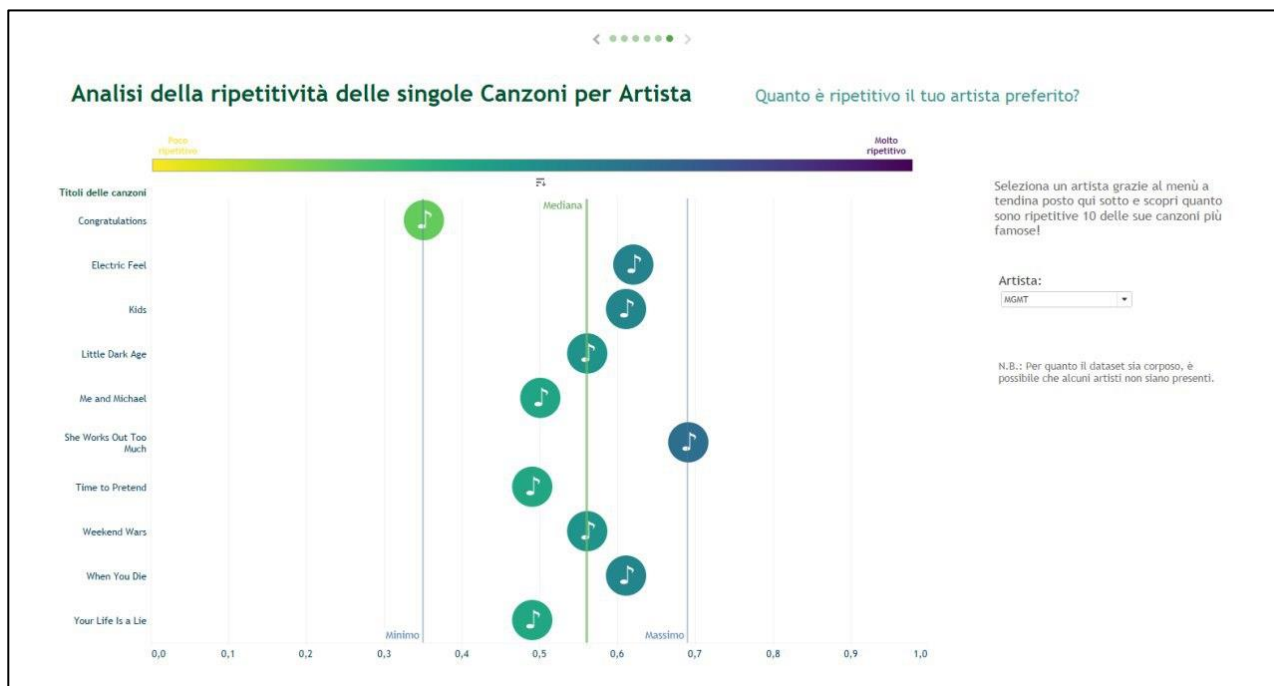


Figura 4

L'infografica 4 non ha presentato particolari difficoltà di comprensione da parte degli utenti.

Inoltre, vista la curiosità dei soggetti intervistati nel comprendere come sia stato calcolato l'indice di ripetitività, si è scelto di inserire un'infografica iniziale nella quale venisse approfondito questo aspetto.

### Indice di Ripetitività: definizione

L'indice di Ripetitività è un indicatore in grado di spiegare quanto il testo di una canzone risulti ripetitivo.

Si ottiene dal *complementare* del rapporto tra il numero di parole uniche nel testo (ottenuto mantenendo solo una parola di quelle ripetute più di una volta) e il numero totale di parole presenti prima di procedere a ogni sorta di modifica.

In formula:

$$1 - \frac{\text{parole\_uniche}}{\text{parole\_totali}}$$

Il valore ottenuto dal rapporto ci indica la percentuale di porzione di testo della canzone residuale dopo la manipolazione. Il suo complementare, invece, ci indicherà la percentuale di testo rimossa dal totale, e cioè l'insieme di tutte le parole della canzone che si ripetono.

Per efficientare la riduzione del testo, abbiamo deciso di eliminare ogni segno di punteggiatura quali apostrofi, punti, virgole e similari, e di rendere tutte le lettere minuscole (così da eliminare la distinzione tra "Cavolfiore" e "cavolfiore").

Volendo fare un esempio, analizziamo qui di fianco la canzone "Il pescatore" di Fabrizio De André.

N.B.: Non abbiamo applicato al testo le modifiche di cui sopra, al fine di facilitare la comprensione del procedimento.

1) Qui di seguito viene riportata una parte del testo originario della canzone:

"All'ombra dell'ultimo sole  
S'era assopito un pescatore  
E aveva un solco lungo il viso  
Come una specie di sorriso  
Veniva alla spiaggia un assassino  
Due occhi grandi da bambino  
Due occhi enormi di paura  
Eran gli specchi di un'avventura  
E chiese al vecchio "dammi il pane  
Ho poco tempo e troppa fame"  
E chiese al vecchio "dammi il vino  
Ho sete e sono un assassino"  
Gli occhi dischiuse il vecchio al giorno  
Non si guardò neppure intorno  
Ma versò il vino e spezzò il pane  
Per chi diceva "ho sete, ho fame"

2) Volendone evidenziare, in rosso, le parole ripetute nel testo appena mostrato, si avrebbe:

"All'ombra dell'ultimo sole  
S'era assopito un pescatore  
E aveva **un** solco lungo il viso  
Come una specie di sorriso  
Veniva alla spiaggia **un** assassino  
Due occhi grandi da bambino  
**Due** occhi **enormi** di paura  
Eran gli specchi di un'avventura  
**E** chiese al vecchio "dammi il pane  
Ho poco tempo e troppa fame"  
**E** chiese al vecchio "dammi il vino  
**Ho** sete e sono **un** assassino"  
Gli **occhi** dischiuse il vecchio al giorno  
Non si guardò neppure intorno  
Ma versò **il** **vino** e spezzò **il** **pane**  
Per chi diceva "ho **sete**, **ho** fame"

3) Procedendo, quindi, all'eliminazione delle parole appena evidenziato, avremmo un testo (privo di senso logico), che recita più o meno così:

"All'ombra dell'ultimo sole  
S'era assopito un pescatore  
E aveva solco lungo il viso  
Come una specie di sorriso  
Veniva alla spiaggia assassino  
Due occhi grandi da bambino  
enormi paura  
Eran gli specchi 'avventura  
chiese al vecchio "dammi pane  
Ho poco tempo e troppa fame"  
"vino  
sete sono "  
Gli dischiuse giorno  
Non si guardò neppure intorno  
Ma versò spezzò  
Per chi diceva "ho , , "

Quindi avremo:

- parole\_uniche = 65
- parole\_totali = 96
- Indice di Ripetitività =  $1 - (65/96) = 0,32$

Applicando la formula vista all'inizio alla porzione di testo qui riportata, il testo si riduce del **32%**.

Se invece applicassimo tale procedimento di riduzione all'intero testo della canzone, qui non riportato per non tediarlo ulteriormente il lettore, si arriva a ottenere un Indice di Ripetitività pari a 0,52, il quale indica che la canzone, nella sua completezza, presenta il 52% di parole che si ripetono.

N.B.: Le Figure 1, 2, 3 e 4, e cioè le versioni finali delle infografiche, sono state successivamente mostrate agli utenti per completare la sezione “Task”.

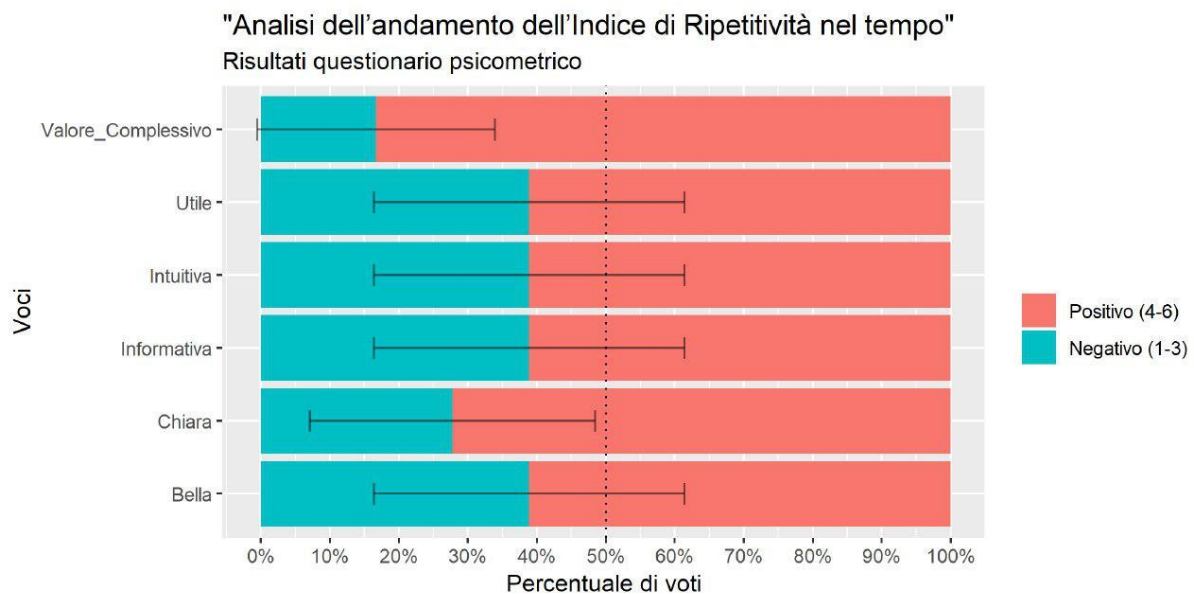
## Questionario psicometrico

Sono inoltre stati sottoposti ad ulteriori 78 utenti i questionari psicometrici. Una volta visualizzata l'infografica, agli utenti venivano sottoposti 5 quesiti:

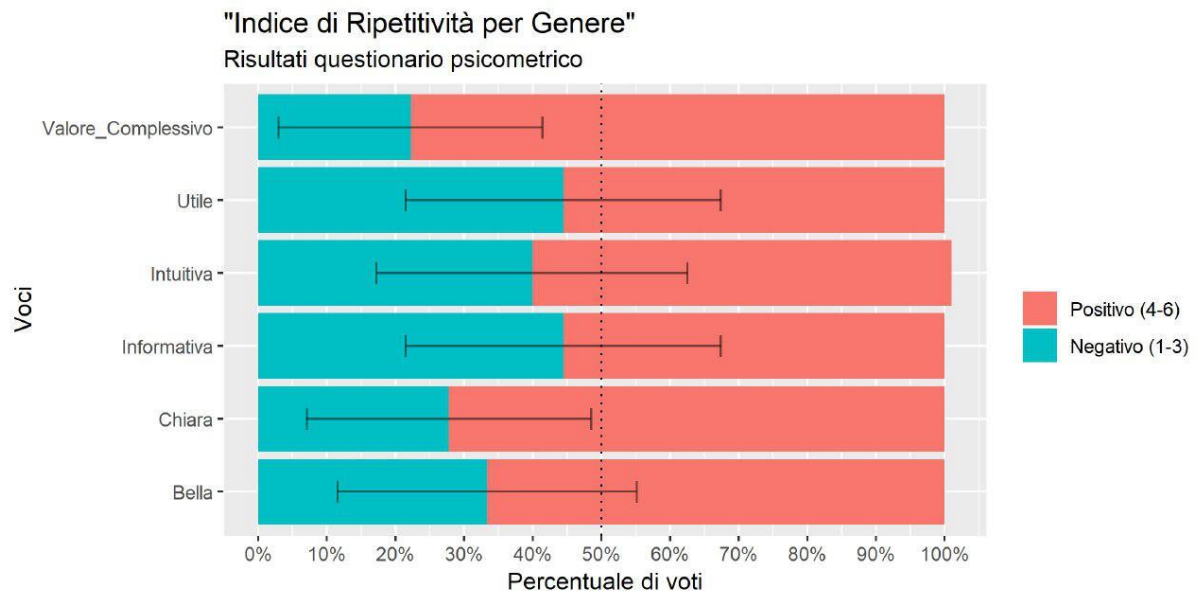
- Quanto ritieni l'infografica Utile?
- Quanto ritieni l'infografica Intuitiva?
- Quanto ritieni l'infografica Chiara?
- Quanto ritieni l'infografica Informativa?
- Quanto ritieni l'infografica Bella?
- Qual è il valore complessivo?

Di seguito vengono riportati i risultati ottenuti attraverso gli stacked barchart.

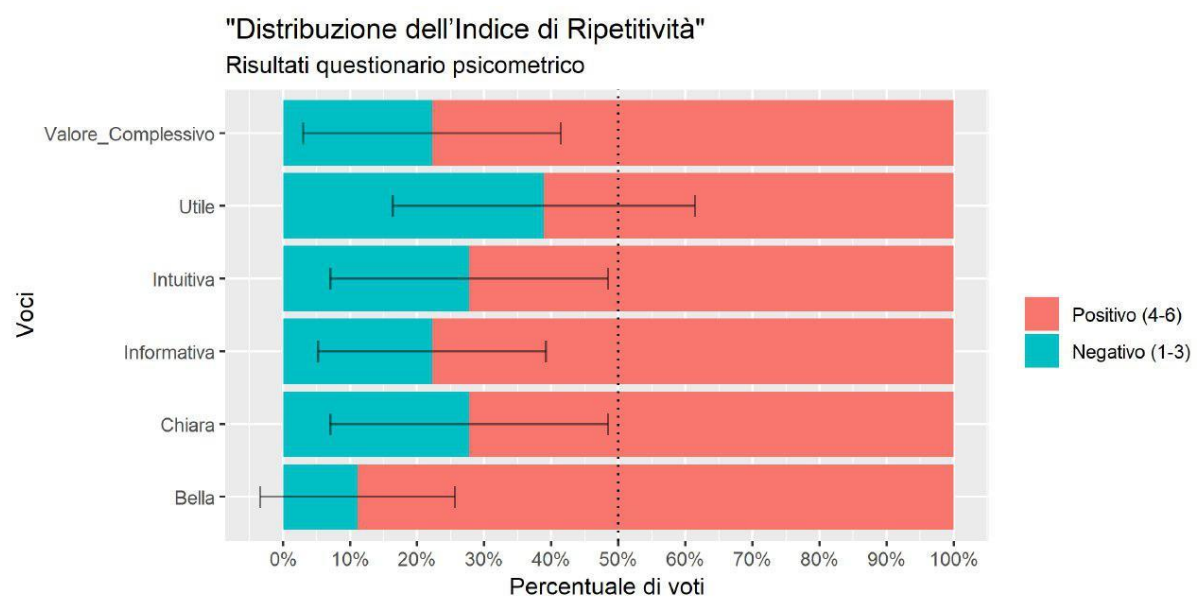
### ◦ Infografica 1



## ◦ Infografica 2

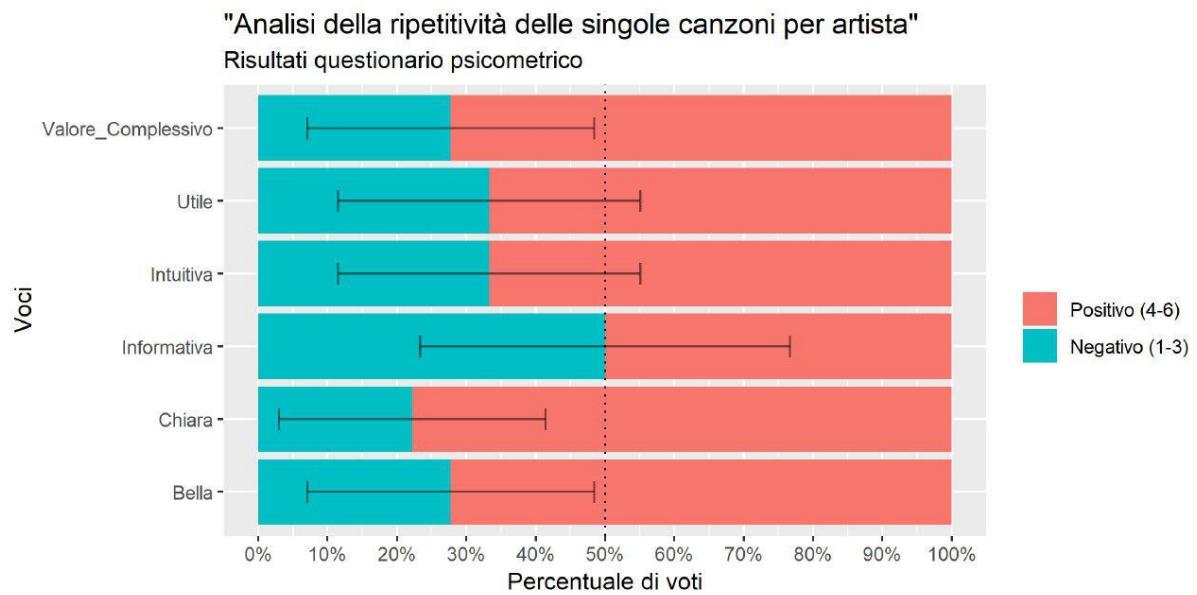


## ◦ Infografica 3

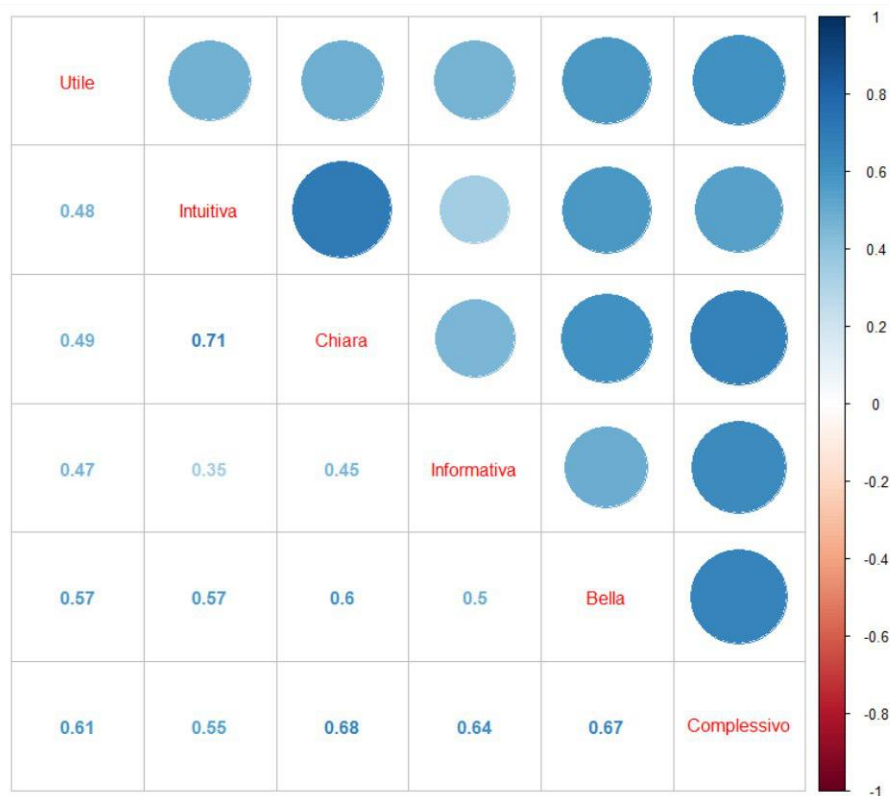




#### ◦ Infografica 4



Si è inoltre verificato se esistesse correlazione tra le risposte date ai questionari.



Gli indici di correlazione più alti si presentano nelle coppie *complessivo - bella*, *complessivo - utile*, e *complessivo - informativa*, a testimoniare l'apprezzamento sia estetico che di contenuto.

## Task

Successivamente a 12 utenti sono stati sottoposti dei task, ovvero delle domande le cui risposte potevano essere dedotte osservando l'infografica. Questo è stato fatto al fine di verificare quanto le informazioni che si volevano trasmettere attraverso l'infografica fossero effettivamente percepite.

Successivamente a questi utenti sono stati sottoposti i questionari psicometrici.

Si riportano di seguito dapprima i violin plot, i quali riportano i risultati dei task, e subito dopo vengono mostrati gli stacked barchart indicanti le risposte ai questionari psicometrici.

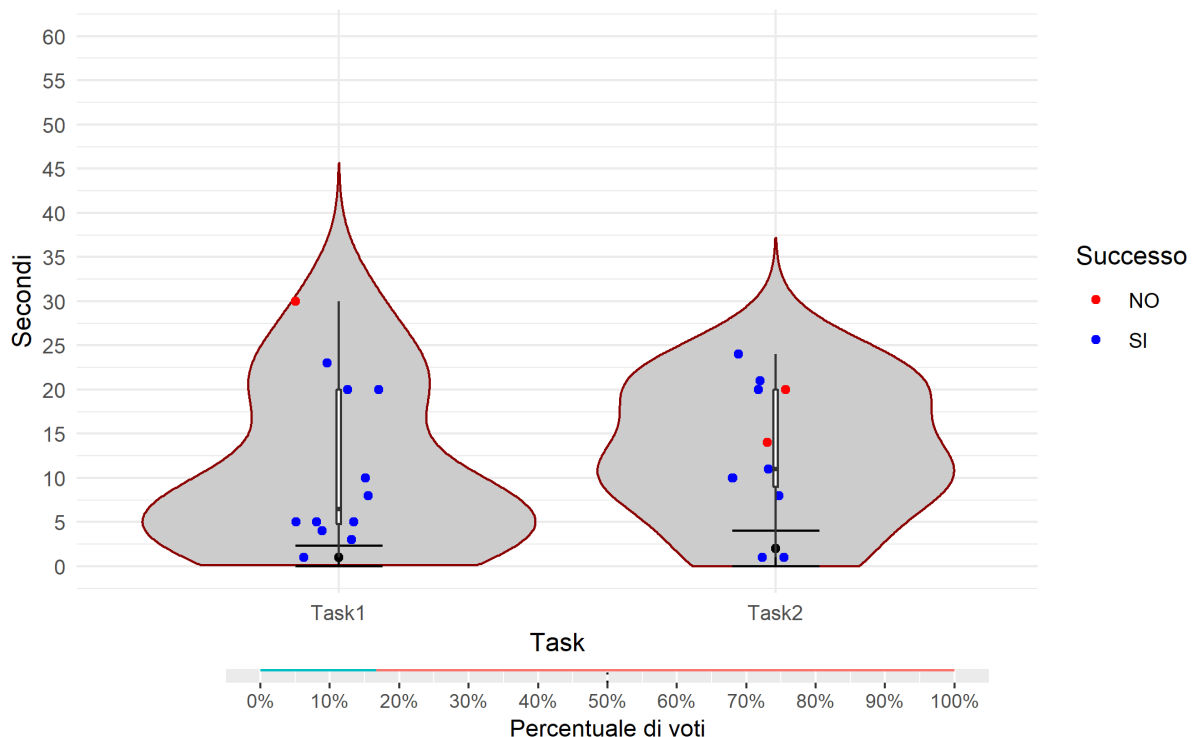
### ◦ **Infografica 1** - *Analisi dell'andamento dell'Indice di Ripetitività nel tempo*

Per la prima visualizzazione i task richiesti agli utenti sono stati:

1. Indicare se le tendenze degli ultimi anni presentano qualche trend visibile;
2. Indicare se grafico a destra e quello a sinistra rappresentano l'andamento delle stesse variabili.

#### Tempi di esecuzione per Task

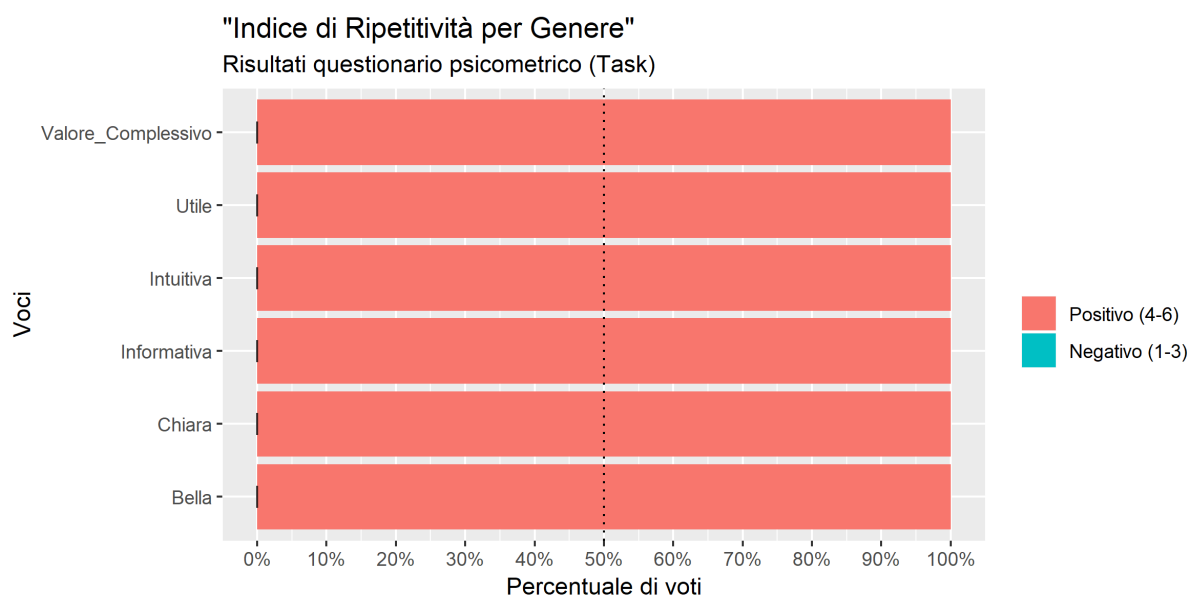
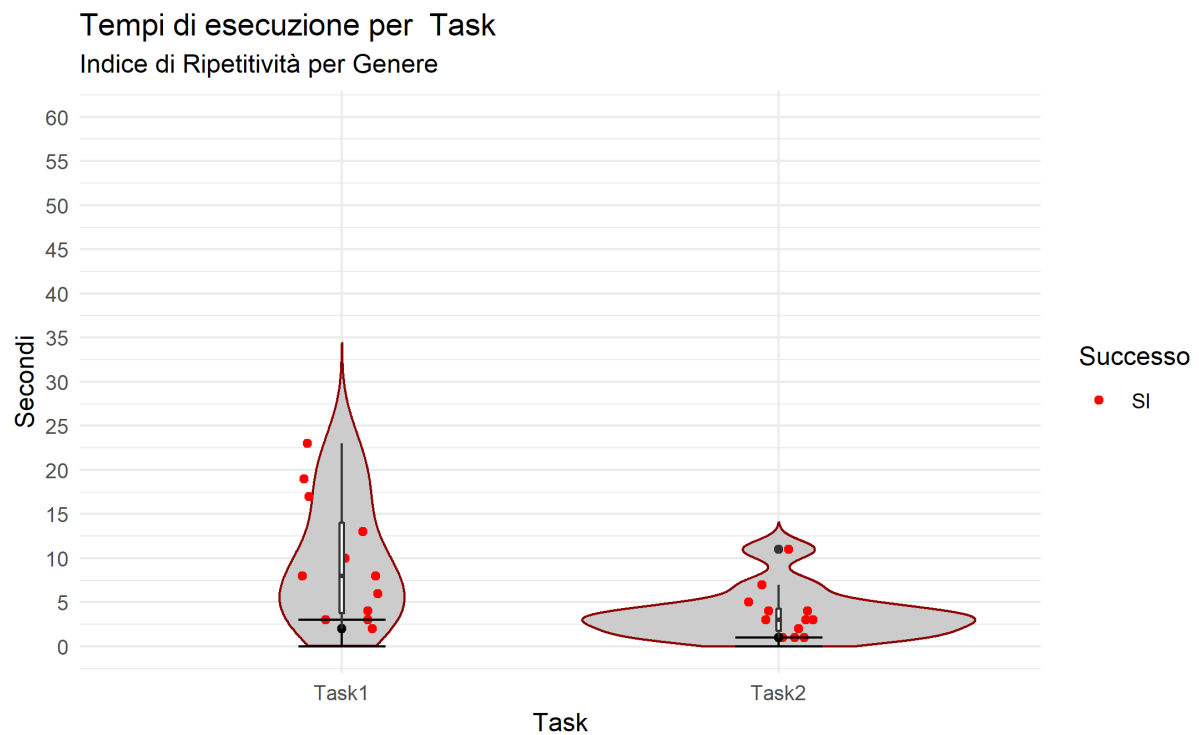
Analisi dell'andamento dell'Indice di Ripetitività nel tempo



◦ **Infografica 2** - *Indice di Ripetitività per Genere*

Per la seconda visualizzazione i task effettuati sono stati:

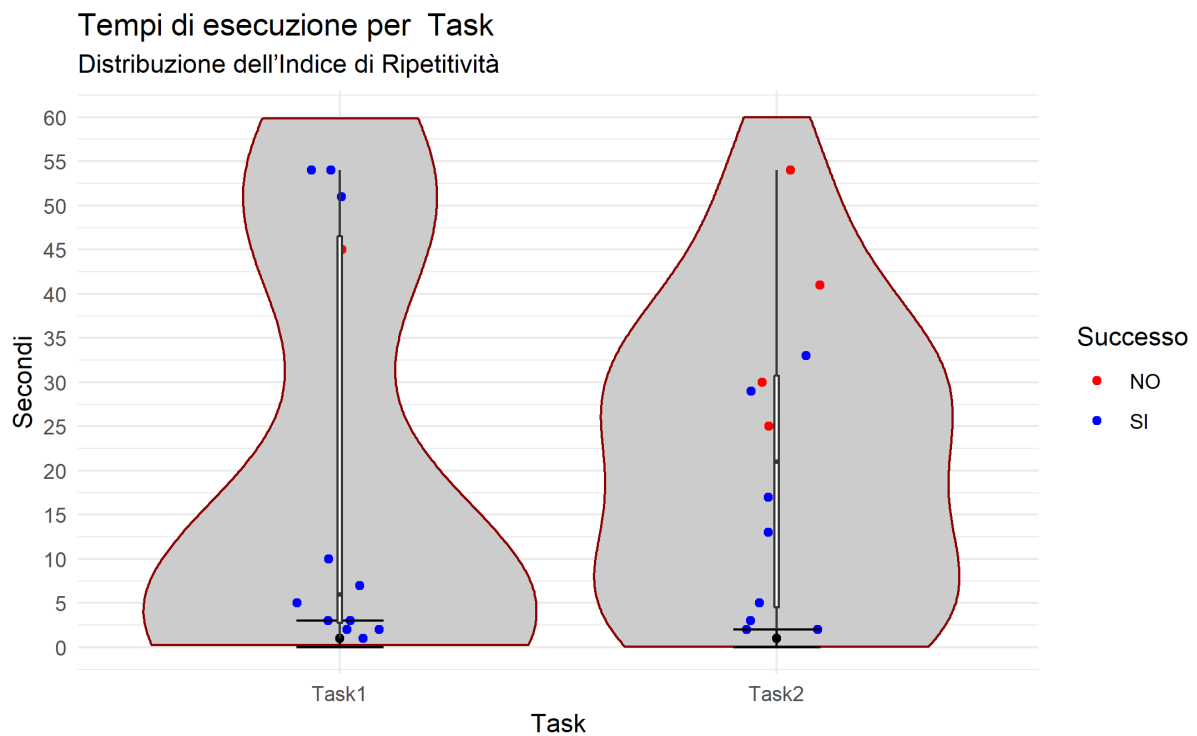
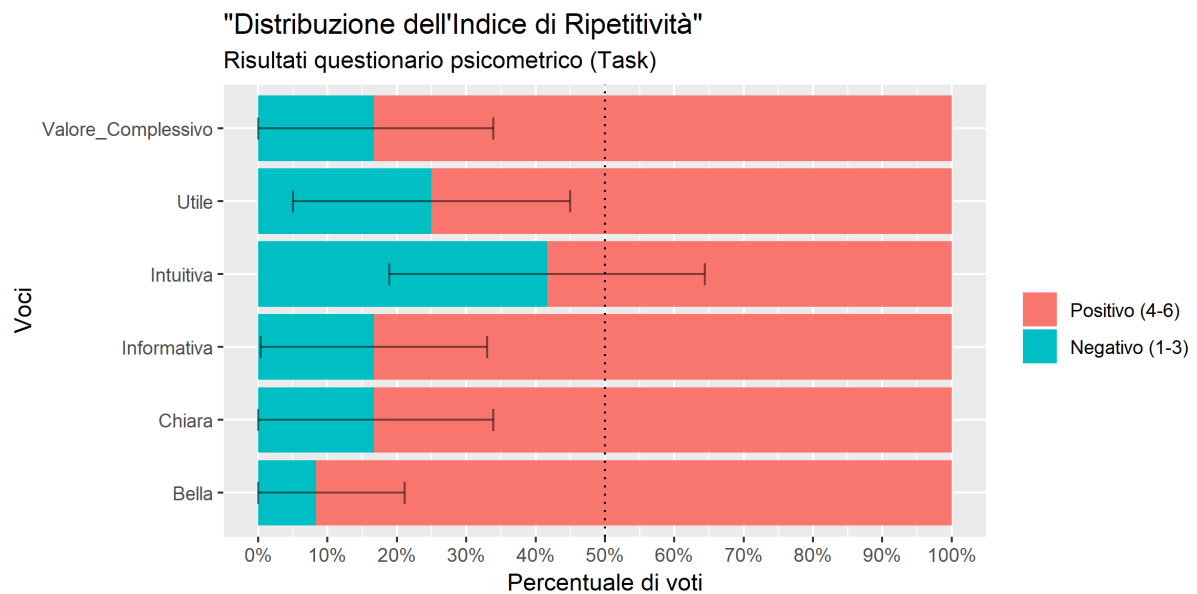
1. Indicare la colonna "Trap" e dire se è significativamente più ripetitiva rispetto alla media;
2. Indicare il genere più ripetitivo.



◦ **Infografica 3** - *Distribuzione dell'Indice di Ripetitività*

Per la terza visualizzazione i task effettuati sono stati:

1. Indicare qual è il valore mediano dell'indice di ripetitività;
2. Indicare i limiti della banda di normalità.



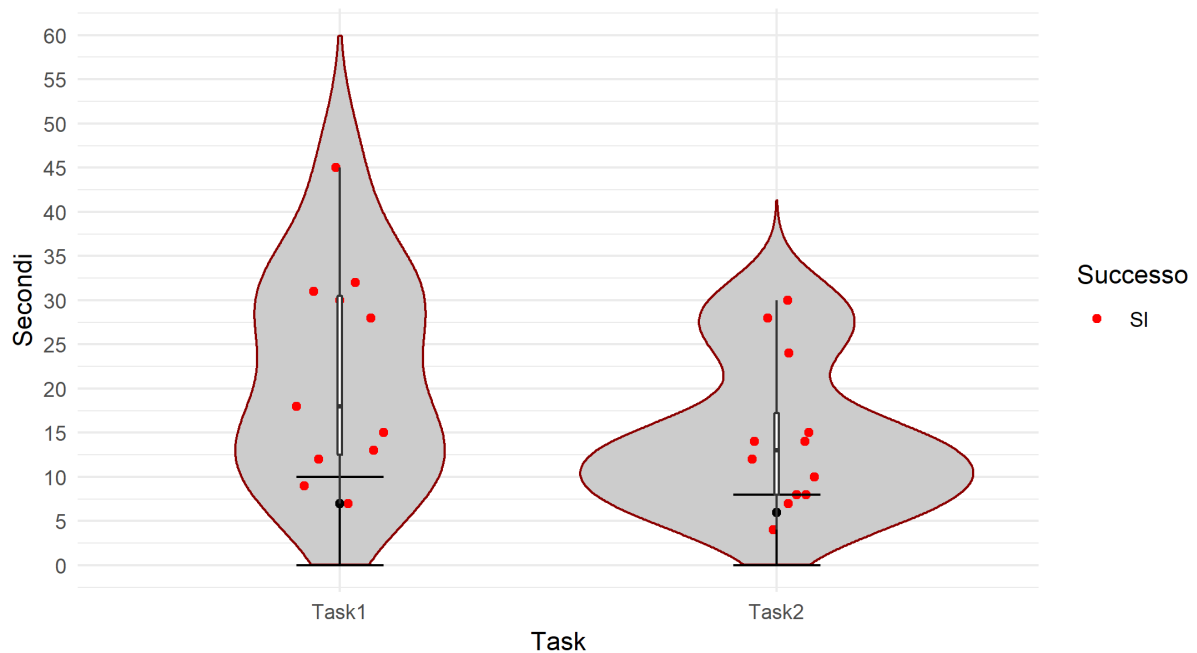
◦ **Infografica 4** - *Analisi della ripetitività delle singole canzoni per artista*

Per la quarta visualizzazione i task effettuati sono stati:

1. Indicare qual è il valore dell'indice di ripetitività della canzone "A casa de Sandro?" di Achille Lauro
2. Qual è la canzone meno ripetitiva di Capo Plaza

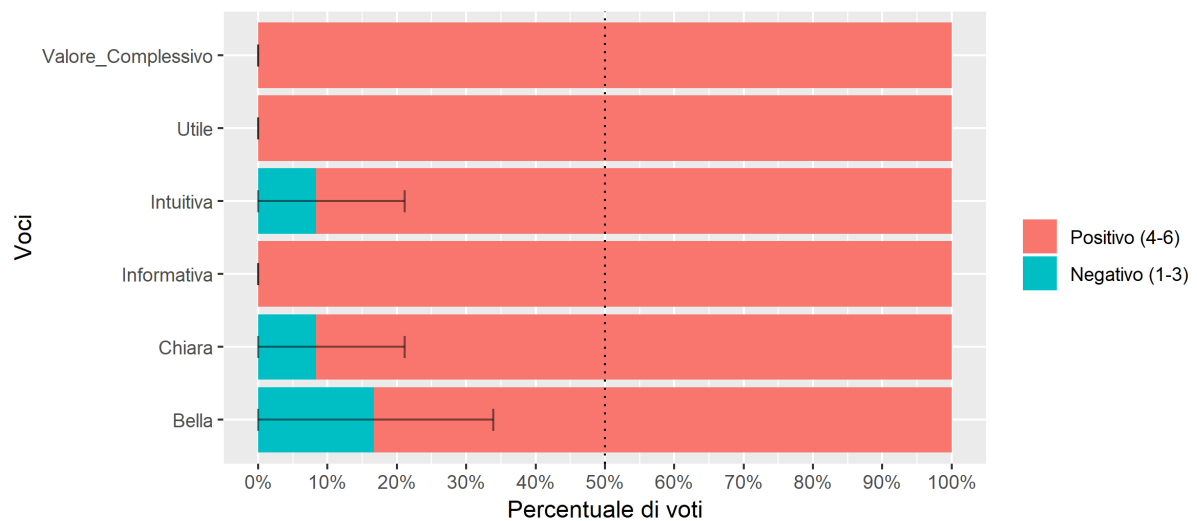
**Tempi di esecuzione per Task**

Analisi della ripetitività delle singole canzoni per artista



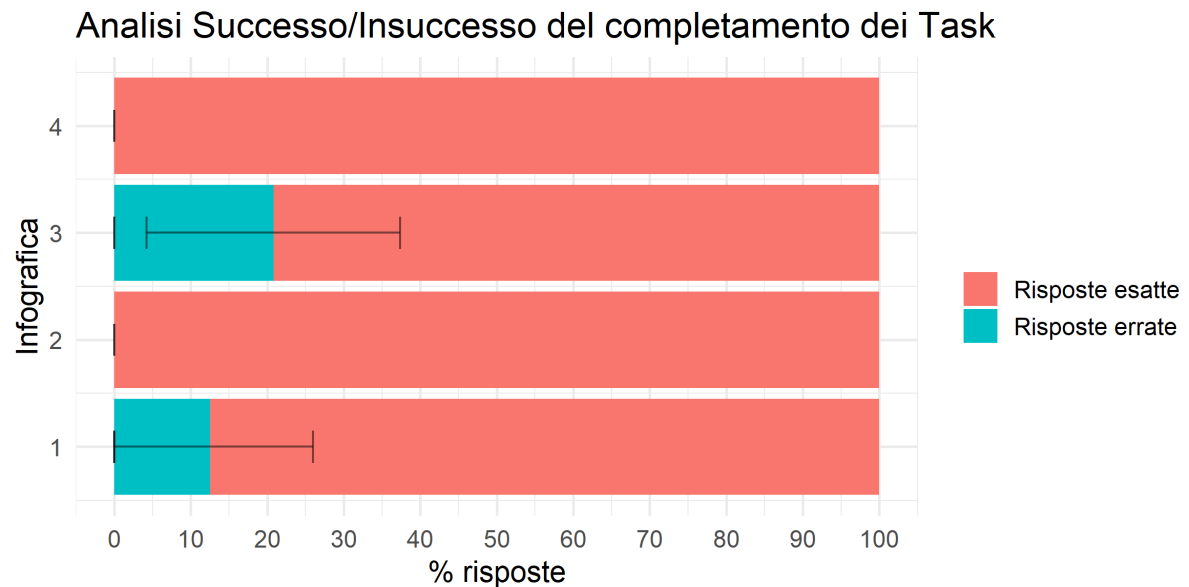
**"Analisi della ripetitività delle singole canzoni per artista"**

Risultati questionario psicometrico (Task)



È interessante notare come gli utenti a cui è stato richiesto di svolgere i task, abbiano apprezzato maggiormente le infografiche.

Più precisamente, le visualizzazioni più apprezzate sono state la seconda (*Figura 2*) e la quarta (*Figura 4*).



Si può notare come i task relativi alle Infografiche 2 e 4 siano stati tutti completati con successo, mentre per quanto riguarda le Infografiche 1 e 3 vi sono state delle risposte errate. Inoltre, anche la densità di distribuzione dei tempi di risposta emersa dai violin plot si può differenziare tra il gruppo di infografiche 1 e 3 e il gruppo di infografiche 2 e 4. Questa differenza si suppone possa essere dovuta al diverso background degli intervistati, il quale non è stato tuttavia registrato. Si è però notato che utenti con background umanistico facessero fatica a comprendere concetti prettamente statistici/matematici (quali intervallo di confidenza, variabili e trend), nonostante questi fossero riportati in modo chiaro e ben esposti all'interno delle visualizzazioni stesse, e nonostante la risposta ai task non richiedesse una vera e propria comprensione del concetto ma una semplice comprensione della "legenda" (e cioè saper leggere le indicazioni poste nelle infografiche).

### Link alle Infografiche

Qui seguito si riporta il link alla pagina di Tableau in cui è possibile visualizzare una versione interattiva delle Infografiche qui analizzate:

- [https://public.tableau.com/views/ProgettodiDataVisualization-/DataVisualization?:embed=y&:display\\_count=yes&:toolbar=no&:origin=viz\\_share\\_link](https://public.tableau.com/views/ProgettodiDataVisualization-/DataVisualization?:embed=y&:display_count=yes&:toolbar=no&:origin=viz_share_link)