

Descrizione del servizio

1/1/2025

Introduzione	4
I dati sintetici	5
Confronto con anonimizzazione classica	5
Allineamento con il GDPR e l'Al Act	6
Perché scegliere Aindo	7
Moduli	10
Dati Sintetici Anonimi	11
Data Curation	11
Predizione	12
Valutazione Dati Sintetici	12
Funzionalità	13
Sintesi	13
Generazione di Dati Sintetici Relazionali	13
Generazione Parzialmente Sintetica	14
Utilizzo di Generatori Pre-Addestrati	14
Valutazione dei dati sintetici	14
Anonimizzazione	15
Anonimizzazione o pseudonimizzazione classica	15
Redazione delle PII nei testi e nei documenti	15
Modellistica predittiva	15
Predizione su Dati Relazionali	16
Imputazione intelligente dei Dati Mancanti	16
Estrapolazione di Tendenze e Scenari What-If	17
Funzionalità in Interfaccia Utente, REST API, SDK	17
Caratteristiche tecniche	19
Design di alto livello	19
Connettori	20
Design di basso livello	20
Archiviazione	21
Flusso di Dati	24
Scalabilità e Benchmarks	28
CPU benchmark	28
Dataset: Adult	28
Dataset: Basket	29
Dataset: Airbnb	29
GPU Benchmarks	30
Dataset: Berka	30
Dataset: Porto	30
Sicurezza, Certificazioni e Compliance	32
Conformità alle Normative (GDPR, Al Act)	32
Certificazioni (ISO, EuroPrivacy)	32

TOM's (Technical and Organizational Measures)	32
Sicurezza e Crittografia	33
Servizi	34
Setup	34
Supporto Cliente	34
Richieste di Supporto	34
Lingua	35
Livelli di Severità	35
Obiettivi del Livello di Servizio	35
Ambito	36
Richieste di Modifica	36
Gestione Operativa della Piattaforma	37
Guida all'uso	38
Formazione	38
Casi d'Uso	38
Benefici per le Pubbliche Amministrazioni e gli Enti Regionali	38
Applicazioni in Al/ML e Analisi Dati	39
Protezione delle Informazioni Personali nei Testi e Documenti	40
Licenza	41
Sottoscrizione Sperimentale	41
Sottoscrizione Commerciale Enterprise	41
Monte Ore per Sviluppo e Consulenza (Opzionale)	41

Introduzione

Nel panorama digitale odierno, i dati sono il motore dell'innovazione, ma la loro gestione pone sfide complesse, soprattutto quando si tratta di proteggere la privacy e rispettare normative sempre più stringenti come il GDPR e l'Al Act. Aindo nasce proprio per risolvere questo dilemma, offrendo una soluzione all'avanguardia che permette alle organizzazioni di sfruttare al massimo il potenziale dei propri dati, senza compromettere la sicurezza delle informazioni sensibili.

La piattaforma di Aindo si basa su una tecnologia innovativa: i dati sintetici. A differenza delle tradizionali tecniche di anonimizzazione, che spesso riducono l'utilità dei dati o lasciano margini per la re-identificazione, i dati sintetici di Aindo sono generati attraverso modelli Al avanzati. Questi dati mantengono le proprietà statistiche, le correlazioni e le dipendenze temporali dei dati originali, ma non contengono alcun collegamento diretto con individui reali. In altre parole, offrono un migliore compromesso tra l'utilità per analisi e modelli predittivi, e la sicurezza per la privacy.

Con Aindo, le organizzazioni possono finalmente superare i limiti imposti dalla scarsità di dati, dalla loro qualità variabile o dai vincoli normativi. La piattaforma consente di generare dataset sintetici di alta qualità, ideali per una vasta gamma di applicazioni, dalla ricerca medica alla finanza, dalla customer analytics alla pianificazione strategica. Grazie a strumenti integrati per il ribilanciamento dei dati, l'imputazione intelligente dei valori mancanti e la creazione di scenari "what-if", Aindo non solo preserva l'integrità dei dati, ma li migliora, rendendoli più rappresentativi e utili per il processo decisionale.

La nostra piattaforma offre soluzioni a sfide tecniche chiave, come l'accessibilità limitata ai dati, la scarsità e la variabilità della qualità dei dati, garantendo al contempo la conformità con rigorosi requisiti normativi come il GDPR e l'Al Act. Grazie a soluzioni avanzate di dati sintetici, le organizzazioni possono sfruttare con fiducia i propri dati per l'innovazione e il processo decisionale, trasformando gli ostacoli in opportunità.

Uno dei punti di forza di Aindo è la sua flessibilità. La piattaforma è progettata per adattarsi alle esigenze di ogni organizzazione, sia che si tratti di una piccola startup o di una grande azienda. Disponibile come servizio SaaS, in cloud privato o on-premise, Aindo supporta scalabilità orizzontale e verticale, con ottimizzazione per CPU, GPU e multi-GPU. Questo significa che, indipendentemente dalla dimensione o dalla complessità dei dataset, Aindo è in grado di gestirli in modo efficiente e sicuro.

Ma Aindo non è solo tecnologia: è anche conformità e sicurezza. La piattaforma è certificata EuroPrivacy per il GDPR, ISO 9001 per la qualità e ISO 27001 per la sicurezza delle informazioni. Con crittografia end-to-end, penetration test regolari e strumenti avanzati come la redazione automatica delle PII (Informazioni Personali Identificabili) in testi e documenti, Aindo garantisce che i dati siano protetti a ogni livello, sia a riposo che in transito.

Infine, Aindo si distingue per la sua semplicità d'uso. Grazie a un'interfaccia utente intuitiva, una REST API robusta e un SDK per sviluppatori, la piattaforma è accessibile a utenti con diversi livelli di competenza tecnica.

Aindo non è solo una piattaforma per la generazione di dati sintetici: è un ecosistema completo che permette alle organizzazioni di innovare in modo sicuro, etico e conforme alle normative, permettendo di:

- Sbloccare l'accesso e il riutilizzo di dataset soggetti a vincoli di privacy attraverso la sintesi basata su modelli di Al Generativa e tecniche di anonimizzazione e pseudonimizzazione"tradizionali";
- Migliorare la qualità dei dati tramite ribilanciamento, augmentation e imputazione dei valori mancanti;
- Creare modelli predittivi direttamente sui dati relazionali, senza necessità di feature engineering;
- Generare scenari di What-If, per simulazioni e analisi previsionali;
- Identificare e redarre le informazioni personali identificabili (PII) all'interno di dati testuali e documenti¹

Con Aindo, i dati non sono più un limite, ma un'opportunità per trasformare le sfide in soluzioni e guidare il futuro dell'innovazione.

I dati sintetici

I dati sintetici rappresentano un'innovazione rivoluzionaria nella gestione e nell'utilizzo delle informazioni, offrendo una soluzione più sicura ed efficace rispetto alle tradizionali tecniche di anonimizzazione. Grazie all'uso di modelli avanzati di intelligenza artificiale, i dati sintetici preservano le proprietà statistiche e strutturali dei dati reali senza mantenere alcun legame con individui specifici. Questo approccio garantisce un'elevata protezione della privacy, mantenendo al contempo l'utilità dei dati per analisi, ricerca e sviluppo di modelli Al/ML. In questa sezione, esploreremo i principali vantaggi dei dati sintetici rispetto all'anonimizzazione classica e il loro allineamento con normative come il GDPR e l'Al Act.

Confronto con anonimizzazione classica

La tecnologia di sintesi che Aindo offre è intrinsecamente più sicura dell'anonimizzazione dei dati per la protezione dei dati personali:

Rischio di re-identificazione minimizzato

L'anonimizzazione classica (es. mascheramento, generalizzazione) mantiene un collegamento uno-a-uno tra i dati originali e quelli anonimizzati, rendendo possibile la re-identificazione attraverso correlazioni o inferenze. I dati sintetici, invece, vengono generati da zero utilizzando modelli di intelligenza artificiale addestrati sui pattern dei dati originali. Non esiste un collegamento diretto con individui reali, minimizzando il rischio di re-identificazione.

_

¹ In fase di sviluppo

Mantenimento dell'utilità dei dati

L'anonimizzazione classica spesso riduce la qualità dei dati rimuovendo o distorcendo informazioni per proteggere la privacy. I dati sintetici, invece, preservano le proprietà statistiche, le correlazioni e le dipendenze temporali dei dati originali, garantendo un'elevata utilità per l'addestramento Al/ML, le analisi e il processo decisionale.

Garanzie avanzate di privacy

La nostra piattaforma utilizza tecniche come l'addestramento differenzialmente privato e la soppressione di valori rari per limitare matematicamente la fuga di informazioni. Questi metodi forniscono garanzie di privacy più forti rispetto all'anonimizzazione classica, che manca di protezioni formali di questo tipo.

Scalabilità

A differenza del processo di anonimizzazione, che è un processo manuale che deve essere effettuato da un esperto del settore al fine di minimizzare i rischi, la sintesi dei dati risulta assai più automatizzabile non richiedendo competenze specifiche in termini di privacy, permettendo una migliore scalabilità del processo.

In sintesi, i dati sintetici offrono un livello più elevato di protezione della privacy perché rompono il legame tra i record reali e quelli sintetici, mantenendo al contempo le proprietà statistiche dei dati originali.

Aspetto	Anonimizzazione classica	Dati sintetici
Rischio di re-identificazione	Alto (tramite attacchi di correlazione)	Minimo (nessuna mappatura 1:1 con individui reali)
Utilità dei dati	Spesso degradata	Preservata (statistiche, correlazioni, formati)
Garanzie di privacy	Non disponibili	Differential privacy
Scalabilità	Limitata da processi manuali	Generazione automatizzata e scalabile

Allineamento con il GDPR e l'Al Act

La piattaforma di Aindo è pensata esplicitamente per soddisfare i requisiti deil GDPR e dell'Al act:

Conformità al GDPR

- I dati sintetici sono intrinsecamente anonimi secondo il GDPR (Considerando il recital 26: "The principles of data protection should therefore not apply to anonymous information, namely information which does not relate to an identified or identifiable natural person or to personal data rendered anonymous in such a manner that the data subject is not or no longer identifiable"), poiché non contengono informazioni identificabili.
- La piattaforma di Aindo è certificata EuroPrivacy, confermando la conformità al GDPR per la generazione di dati sintetici, in particolare in settori sensibili come la sanità.
- Funzionalità come la crittografia (dei dati a riposo/in transito) e la redazione delle PII garantiscono un'ulteriore aderenza ai principi di sicurezza e responsabilità del GDPR.

Conformità all'Al Act

Grazie alla generazione dei dati sintetici la piattaforma di Aindo favorisce uno sviluppo sicuro, affidabile ed etico degli strumenti di intelligenza artificiale, e in conformità con quanto previsto dagli articoli 10 e 59 dell'Al ACT.

Secondo l'Al ACT, i dati sintetici sono considerati anonimi e non personali, e il loro utilizzo è preferibile, ove possibile, per l'addestramento dei modelli di intelligenza artificiale. Questo approccio consente di evitare l'impiego di dati reali, oltre a mitigare bias e pregiudizi tramite funzionalità di rebalancing e augmentation, offerte nativamente dalla piattaforma Aindo. Anche il DDL n. 1146 del 20 maggio 2024 sull'intelligenza artificiale, in ambito sanitario, disciplina l'uso di tali strumenti, sottolineando l'importanza del cosiddetto "unbiasing" e suggerendo l'adozione di soluzioni adeguate di riequilibrio. Le tecniche di generazione controllata, come il rebalancing e l'augmentation, consentono di produrre set di dati equilibrati e rappresentativi, contribuendo a ridurre le disparità e a migliorare l'equità dei sistemi di intelligenza artificiale.

Perché scegliere Aindo

Aindo si distingue come leader nella generazione di dati sintetici, offrendo una tecnologia avanzata, certificata e testata da organismi internazionali. La sua piattaforma è progettata per garantire il massimo livello di protezione della privacy, qualità dei dati e conformità normativa, consentendo alle organizzazioni di innovare in modo sicuro ed efficace.

Tecnologia all'avanguardia riconosciuta internazionalmente

La tecnologia di Aindo è riconosciuta come **best in class** e ha superato test condotti da organismi internazionali di riferimento, come il **NIST (National Institute of Standards and**

Technology), che valuta le soluzioni più avanzate nel campo della privacy e della sicurezza dei dati. Maggiori dettagli sono disponibili <u>qui</u>.

Integrità del dato e delle relazioni

Aindo genera dati sintetici multi-tabella mantenendo chiavi esterne, correlazioni cross-table e dipendenze temporali, senza richiedere l'appiattimento dei dati in strutture monodimensionali.

Flessibilità d'uso per ogni livello di competenza

Aindo offre modalità di utilizzo adatte sia agli utenti esperti che a chi non ha competenze tecniche avanzate:

- **Interfaccia grafica intuitiva**: consente a chiunque di generare e gestire dati sintetici in pochi clic.
- **SDK per sviluppatori**: offre pieno controllo e integrazione con workflow esistenti, permettendo personalizzazioni avanzate.
- REST API: garantisce la compatibilità con sistemi e infrastrutture già in uso.

Combinazione flessibile delle tecniche di anonimizzazione

Possibilità di sintetizzare solo colonne critiche (es.: dati biometrici) e applicare tecniche tradizionali (hashing, generalizzazione) alle restanti, massimizzando il controllo sulla privacy.

Gestione avanzata dei vincoli nei dati sintetici

La piattaforma permette di **definire regole matematiche e di business** per garantire che i dati sintetici rispettino i vincoli presenti nei dataset originali. Questo include relazioni tra colonne, vincoli di uguaglianza/diseguaglianza e la conservazione della struttura relazionale dei dati.

Creazione di dati sintetici da descrizioni testuali

Generazione di dati tabulari partendo da prompt in linguaggio naturale (es.: "Crea un dataset di clienti con età, reddito e storico acquisti"), senza necessità di dati di training iniziali.

Privacy garantita con tecniche avanzate

Aindo integra tecniche di **privacy differenziale (Differential Privacy)**, che forniscono **garanzie matematiche** sulla protezione delle informazioni sensibili. Inoltre, per evitare rischi di re-identificazione, supporta la **rimozione di valori rari** nei dataset sintetici, garantendo un anonimato effettivo.

Supporto per molteplici casi d'uso

La piattaforma di Aindo è progettata per affrontare diverse sfide legate ai dati, tra cui:

- Sviluppo e addestramento di modelli Al/ML senza l'utilizzo di dati reali.
- Data sharing sicuro tra aziende, ricercatori e partner, mantenendo la conformità al GDPR
- Analisi avanzate e data augmentation per bilanciare dataset sbilanciati e migliorare la qualità delle previsioni.

- Redazione automatica delle informazioni personali identificabili (PII)² nei documenti di testo.
- Costruzione di modelli predittivi direttamente sul dato strutturato, senza necessità di feature engineering.
- Simulazione di scenari ipotetici ("what-if analysis") per prendere decisioni basate su dati realistici ma privi di rischi per la privacy.

Certificazioni e conformità normativa

Aindo è la **prima azienda di dati sintetici certificata EuroPrivacy**, confermando la sua conformità al GDPR, in particolare per il settore sanitario. Inoltre, dispone di certificazioni **ISO 9001** (gestione della qualità) e **ISO 27001** (sicurezza delle informazioni), assicurando i più elevati standard di protezione e affidabilità.

Massima sicurezza e protezione dei dati

Aindo implementa **crittografia dei dati a riposo e in transito**, prevenendo accessi non autorizzati. Inoltre, la piattaforma è sottoposta a **penetration test regolari**, garantendo la massima sicurezza contro attacchi informatici.

Connettività e Integrazione con stack legacy

Supporto nativo per database relazionali (Oracle, SQL Server) e sistemi on-premise, oltre a cloud pubblici (AWS, GCP).

Scalabilità e supporto per diverse architetture

La soluzione di Aindo è altamente scalabile e supporta differenti architetture, tra cui:

- CPU e GPU per un'elaborazione efficiente.
- Multi-GPU training per la gestione di dataset complessi e di grandi dimensioni.
- Opzioni di installazione flessibili: può essere usata come SaaS, su cloud privato o on-premises.

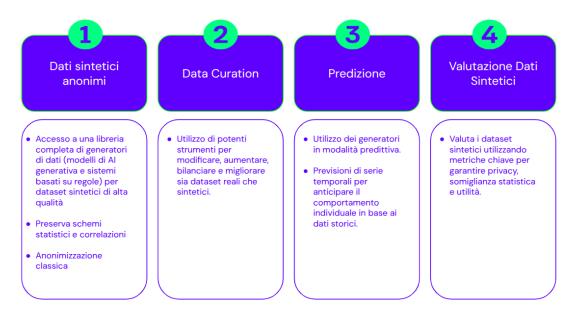
Grazie a queste caratteristiche, Aindo consente alle aziende di sfruttare al meglio i propri dati, garantendo innovazione, conformità e protezione della privacy.

² Svi	lub	og	in	CO	rsc
OVI	up	PΟ		OO.	50

-

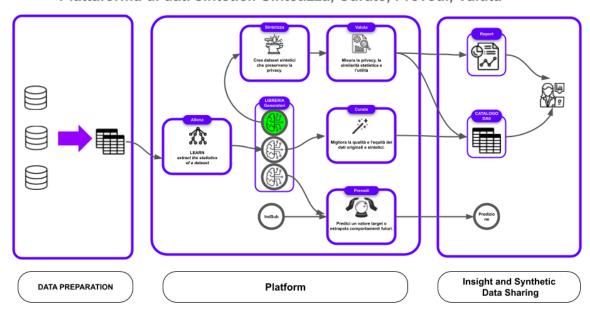
Moduli

La piattaforma Aindo si articola in quattro moduli principali: Dati sintetici anonimi, data curation, predizione e valutazione dei dati sintetici. Ciascun modulo è progettato per rispondere a specifiche esigenze di gestione, generazione e analisi dei dati.



Uno schema del flusso dati di alto livello e dell'organizzazione dei moduli è riprodotta nell'illustrazione seguente.

Piattaforma di dati sintetici: Sintetizza, Curate, Prevedi, Valuta



Nelle seguenti sezioni viene fornita una descrizione dei moduli e delle relative funzionalità.

Dati Sintetici Anonimi

Questo modulo consente di generare dati sintetici che riproducono le proprietà statistiche dei dati originali senza conservarne alcuna informazione identificabile. Rispetto ai metodi di anonimizzazione classica, i dati sintetici offrono una maggiore protezione dalla re-identificazione e mantengono l'utilità analitica. Il modulo supporta dataset relazionali, garantendo la coerenza tra le tabelle e rispettando eventuali vincoli presenti nei dati originali.

Funzionalità principali

- **Generazione di dati sintetici relazionali**: creazione di dataset multi-tabella con mantenimento di chiavi esterne, correlazioni cross-table e dipendenze temporali.
- Privacy differenziale: addestramento dei modelli con garanzie matematiche per limitare la fuga di informazioni sensibili.
- **Generazione parziale**: sintetizzazione selettiva di colonne critiche (es. dati biometrici) mentre si preservano altre informazioni.
- **Anonimizzazione classica:** creazione di dataset anonimi tramite l'applicazione di tecniche di mascheramento, mocking, generalizzazione, etc...
- Redazione di PII nei documenti: rimozione di informazione di identificazione personale (PII) nei documenti.³
- Generatori pre-addestrati: creazione di dati tabulari partendo da descrizioni in linguaggio naturale, senza necessità di dati di training iniziali.
- Gestione vincoli: applicazione di regole di business (es. relazioni matematiche tra colonne) per garantire coerenza con i dati originali.

Data Curation

Questo modulo offre strumenti per il miglioramento della qualità dei dati, inclusi il ribilanciamento delle classi, la gestione dei valori mancanti e l'armonizzazione dei dataset. L'obiettivo è garantire che i dati utilizzati per l'analisi e l'addestramento dei modelli siano rappresentativi e privi di distorsioni. Le tecniche impiegate comprendono l'imputazione intelligente dei dati mancanti e l'applicazione di tecniche di campionamento per correggere squilibri nelle distribuzioni dei dati.

Funzionalità principali

- Ribilanciamento dei dati: correzione di squilibri nelle distribuzioni (es. dataset con categorie sottorappresentate).
- **Imputazione intelligente**: sostituzione di valori mancanti basata su relazioni contestuali tra colonne e tabelle.
- Augmentation: espansione di dataset esistenti con dati sintetici per aumentarne la rappresentatività.

_

³ Sviluppo in corso

Predizione

Il modulo di predizione permette di eseguire analisi predittive direttamente su dati relazionali, senza richiedere operazioni di trasformazione manuale (feature engineering). Utilizza modelli generativi per analizzare le relazioni tra le variabili, consentendo previsioni più accurate e contestualizzate. Inoltre, il modulo rende possibile l'estrapolazione del contenuto di strutture dati complesse nel futuro. Questa funzionalità permette anche la simulazione di scenari (what-if analysis), utili per valutare l'impatto di diverse condizioni iniziali sulle estrapolazioni stesse.

Funzionalità principali

- Modellistica su dati relazionali: previsioni contestuali senza necessità di feature engineering manuale o appiattimento dei dati.
- Estrapolazione di tendenze: analisi di dati storici multi tabella per proiezioni temporali accurate.
- **Scenari what-if**: L'estrapolazione di tendenze può essere effettuata anche in modalità "what-if", permettendo la simulazione di eventi ipotetici.

Valutazione Dati Sintetici

Questo modulo fornisce metriche per misurare la qualità e la privacy dei dati sintetici generati. Le valutazioni includono il confronto con i dati originali per verificare la preservazione delle distribuzioni statistiche e delle correlazioni, oltre a test per garantire la protezione contro il rischio di re-identificazione. Il modulo genera report automatici per supportare la conformità normativa e il controllo qualità dei dataset sintetici.

Funzionalità principali

- Metriche di privacy: test che permettono di valutare i rischi di re-identificazione.
- Valutazione statistica: confronto di distribuzioni, correlazioni e dipendenze tra dati sintetici e originali.
- Report automatici: generazione di documenti PDF con risultati dettagliati.

Funzionalità

La nostra piattaforma è progettata per offrire strumenti all'avanguardia che garantiscono un equilibrio tra privacy, utilità e scalabilità, supportando diversi scenari d'uso.

In questa sezione esploreremo le principali funzionalità della nostra tecnologia di sintesi, dalla generazione di dati relazionali alla valutazione della qualità dei dataset sintetici, fino all'implementazione di tecniche predittive basate su dati strutturati.

Sintesi

La nostra soluzione offre un'ampia gamma di generatori di dati sintetici, in grado di soddisfare la maggior parte delle esigenze di mercato.

Generazione di Dati Sintetici Relazionali

Una robusta libreria di generatori di dati per creare dataset sintetici di alta qualità, composti da una o più tabelle collegate tra loro tramite chiavi esterne. I generatori addestrano direttamente sui dati del cliente, senza uscire dall'ambiente di installazione. I nostri generatori preservano le statistiche delle tabelle originali, comprese le correlazioni tra tabelle diverse o tra righe all'interno della stessa tabella, supportando di *default* le serie temporali. Allo stesso tempo, i dati sintetici generati non possono essere collegati a nessuna persona reale, risultando quindi anonimi secondo il GDPR. Inoltre, preservano anche l'integrità delle chiavi e supportano diversi tipi di dati, tra cui categorici, booleano, data, ora, data e ora, interi, numero in virgola mobile, coordinate, testo libero, codice fiscale etc... In particolare, anche la colonna di testo libero viene generata mantenendo le correlazioni con le altre colonne della stessa tabella.

I nostri generatori offrono un'eccellente protezione della privacy: supportano l'addestramento con privacy differenziale per garantire matematicamente il corretto livello di protezione delle informazioni. Inoltre, per proteggersi da fonti secondarie di perdita di privacy, l'utente ha la possibilità di generare dati sintetici privi di valori rari presenti nei dati originali.

Infine, la nostra soluzione include un modulo che consente all'utente di definire *business rules* o vincoli presenti nei dati, che devono essere rispettati sistematicamente dai dati sintetici. Attualmente gestiamo relazioni matematiche generiche tra colonne, espresse come equazioni di uguaglianza o disuguaglianza (espressioni come: "Colonna A / Colonna C + 3 > Colonna D") e mappature tra colonne.

Casi d'uso:

- **Sanità**: Generazione di cartelle cliniche sintetiche interconnesse (pazienti, diagnosi, terapie) per ricerca medica senza esporre dati reali.
- **Finanza**: Replica di database transazionali (conti, movimenti, filiali) per testare algoritmi anti-frode in ambienti sicuri.

Generazione Parzialmente Sintetica

La nostra soluzione permette di sintetizzare solo un sottoinsieme di colonne di un dataset, preservando le altre colonne dai dati originali. Questa opzione consente un controllo dettagliato sulle informazioni originali rivelate e sul livello di protezione della privacy.

Casi d'uso:

- **Sanitario**: Preservare valori di diagnosi ed esami, sintetizzando dati sensibili come informazioni personali ed altri potenziali identificativi come le date.
- **E-commerce**: Preservare gli ID ordine e i prodotti acquistati, sintetizzando dati sensibili come indirizzi o numeri di carta di credito.
- **Telecomunicazioni**: Mantenere metriche di utilizzo (minuti, dati consumati) mentre si sostituiscono numeri di telefono e dati biometrici.

Utilizzo di Generatori Pre-Addestrati

La nostra soluzione consente di generare dati tabulari a partire da una descrizione in linguaggio naturale, utilizzando un generatore pre-addestrato su un vasto set di dati strutturati. Ciò consente la generazione di dati anche in assenza di esempi, per casi d'uso come il test del software in cui è necessario avere dei dati che non si avrebbero normalmente a disposizione. A differenza di altre soluzioni basate su LLM disponibili sul mercato, la nostra soluzione garantisce il formato delle righe e delle colonne, evitando errori di formattazione nei processamenti successivi.

Se ritenuto utile, la soluzione pre-addestrata può essere ulteriormente *fine-tuned*, un'opzione particolarmente utile quando l'utente dispone di un numero limitato di esempi. Questo fine-tuning può essere eseguito su più dataset di formati diversi, permettendo di apprendere da più sorgenti di dati contemporaneamente.

Il generatore pre-addestrato è anche in grado di ampliare un dataset aggiungendo colonne non originariamente presenti in esso.

Casi d'uso:

- Sintesi di dati di piccole dimensioni: I modelli pre-addestrati sono particolarmente utili nel caso in cui il dato di partenza da sintetizzare si di numerosità scarsa.
- Sviluppo software: Creare dati realistici per testare applicazioni CRM senza accedere a database reali.

Valutazione dei dati sintetici

La nostra soluzione offre un set completo di metriche per garantire il rispetto degli standard essenziali di privacy e fedeltà statistica. Questo framework di valutazione è fondamentale per mantenere la qualità dei dati e la conformità normativa nei diversi casi d'uso. Le valutazioni sono disponibili in formato PDF, generato automaticamente ad ogni creazione di dato. Per maggiori dettagli fare riferimento alla documentazione https://docs.aindo.com/evaluation/.

Casi d'uso:

- **Assicurazioni**: Validare che i dati sintetici utilizzati per modelli di rischio mantengano le distribuzioni originali di sinistri e premi.
- Pubbliche amministrazioni: Verificare la conformità GDPR prima di pubblicare dataset per open data.

Anonimizzazione

Oltre alle opzioni di sintesi, la nostra piattaforma offre ulteriori opzioni di protezione della privacy grazie a strumenti di anonimizzazione classica.

Anonimizzazione o pseudonimizzazione classica

La piattaforma offre una libreria di strumenti per l'anonimizzazione o la pseudonimizzazione dei dati tabulari, che combina diverse tecniche: masking, generalizzazione, randomizzazione, permutazione, hashing, mock data. Questa funzione può essere combinata con la sintesi parziale per una grande flessibilità di configurazione. Per maggiori dettagli fare riferimento alla documentazione https://docs.anonymize.aindo.com/.

Casi d'uso:

- **Ricerca di mercato**: Condividere dataset con partner esterni sostituendo nomi e indirizzi con valori mock, ma preservando pattern di acquisto.
- **Logistica**: Pseudonimizzare ID veicoli in dati di telemetria per analisi prestazioni senza esporre flotte specifiche.

Redazione delle PII nei testi e nei documenti⁴

Offriamo un avanzato strumento di elaborazione testuale per identificare e rimuovere le informazioni personali identificabili (PII) nei testi e nei documenti. L'utente può scegliere di sostituire le informazioni personali con dati fittizi, preservando la struttura complessiva del documento.

Casi d'uso:

- Legale: Redigere contratti o verbali prima della condivisione con parti terze.
- **Healthcare**: Rimuovere PII da referti medici per archiviazione in database di ricerca.

Modellistica predittiva

La piattaforma di Aindo supera l'approccio tradizionale alla modellistica predittiva, integrando in modo nativo la capacità di generare dati sintetici con strumenti avanzati di analisi predittiva. Questo connubio unico consente alle organizzazioni non solo di preservare la

⁴ Sviluppo in corso

privacy, ma anche di accelerare lo sviluppo di modelli di machine learning e di ottimizzare processi decisionali complessi, tutto direttamente su dati strutturati e relazionali, senza i colli di bottiglia tipici dei metodi convenzionali.

Predizione su Dati Relazionali

Aindo supera i limiti dei framework tradizionali consentendo previsioni direttamente sul dato relazionale grezzo, senza richiedere operazioni di *feature engineering* manuale o ristrutturazione del database. Mentre i metodi classici obbligano a:

- Appiattire tabelle di database in strutture monodimensionali, perdendo relazioni critiche e contesto:
- **Dedicare tempo di sviluppo** alla creazione manuale di feature rilevanti;
- Affidarsi a team di data scientist per iterazioni complesse;

Aindo automatizza tutto ciò grazie a modelli generativi che comprendono automaticamente le relazioni tra tabelle, colonne e serie temporali. Esempio concreto: in un database sanitario con tabelle separate per pazienti, diagnosi e terapie, la piattaforma può prevedere il rischio di riammissione ospedaliera analizzando in modo contestuale anamnesi, trattamenti e outcome, senza richiedere l'unificazione forzata dei dati. Ciò offre dei vantaggi chiave:

- Riduzione dei tempi di sviluppo: bypassando feature engineering e preprocessing, i modelli sono pronti per la produzione in tempi ridotti
- Maggior accuratezza: preservando la struttura relazionale, i modelli catturano pattern nascosti che i metodi tradizionali tralasciano.

Imputazione intelligente dei Dati Mancanti

La piattaforma di Aindo permette una funzionalità avanzata per l'imputazione di valori mancanti, basata sulla capacità predittiva dei suoi modelli generativi. A differenza dei metodi tradizionali—come l'interpolazione lineare o la sostituzione con valori medi—Aindo sfrutta le relazioni contestuali presenti nei dati strutturati per inferire informazioni mancanti in modo dinamico e coerente. I modelli generativi di Aindo, addestrati sui dati originali, identificano pattern e correlazioni tra colonne e tabelle. Quando incontrano un valore mancante:

- Analizzano il contesto: considerano tutte le variabili correlate (es.: età, transazioni precedenti, dati geografici) anche se distribuite in tabelle relazionali separate.
- Generano un valore probabilistico: il modello propone un'imputazione che rispetta le dipendenze statistiche del dataset (es.: inferire il reddito mancante di un cliente basandosi sul suo codice postale, professione e storico di acquisti).

Rispetto alle Tecniche Tradizionali, l'imputazione fatta con Aindo permette minore distorsione: l'imputazione contestuale evita errori sistematici introdotti da approcci semplificati (es.: sostituire con la mediana). Inoltre le relazioni tra tabelle sono sfruttate per imputazioni più accurate, senza richiedere l'appiattimento dei dati. Anche questo processo viene eseguito senza preprocessing manuale: il sistema opera direttamente sui dati grezzi,

eliminando la necessità di trasformazioni preliminari. Esempi di casi d'uso per questa funzionalità sono:

- Dati sanitari: completare campi mancanti nelle cartelle cliniche (es.: pressione sanguigna non registrata) basandosi su diagnosi, terapie e dati demografici correlati.
- Analisi finanziaria: stimare valori mancanti in report transazionali (es.: importi di vendita) utilizzando informazioni contestuali come periodo, filiale e categoria prodotto.
- **Customer analytics**: ricostruire dati anagrafici parziali (es.: età o genere) partendo da comportamenti d'acquisto e interazioni con il servizio clienti.

Estrapolazione di Tendenze e Scenari What-If

La nostra soluzione include un modulo in grado di effettuare estrapolazioni sui dati relazionali. In questa modalità, il sistema permette di prevedere tendenze future basandosi su dati storici. Questo modulo utilizza un generatore specifico addestrato su dati passati per fornire previsioni accurate e contestualmente rilevanti. Poiché la previsione viene eseguita su dati strutturati complessi, è possibile includere informazioni in formato variegato nelle estrapolazioni. Il modulo è anche in grado di generare scenari di *what-if*, consentendo di valutare le differenze nei risultati nel caso in cui si verifichi un determinato evento. Il funzionamento dello strumento è basato sui seguenti passaggi:

- Apprendimento contestuale: Il sistema analizza dati storici multitable (es.: vendite, comportamenti clienti, dati macroeconomici) identificando relazioni non lineari e dipendenze temporali.
- Generazione di scenari: Basandosi su modelli probabilistici, produce proiezioni realistiche (es.: "Come cambierebbe il tasso di abbandono clienti se aumentassero i prezzi del 10%?").
- Analisi di sensitività: Valuta l'impatto di singole variabili (es.: l'effetto di una campagna marketing su segmenti demografici specifici).

Esempi di casi d'uso per questa funzionalità sono:

- **Finanza**: Prevedere l'andamento di portafogli d'investimento sotto diversi shock di mercato.
- Retail: Simulare l'effetto di promozioni stagionali su inventario e logistica.
- Sanità: Fornire la probabilità di sviluppare una certa patologia in una determinata popolazione di pazienti e quindi destinare i soggetti maggiormente a rischio a campagne di screening.

Funzionalità in Interfaccia Utente, REST API, SDK

La piattaforma è progettata per supportare sia progetti di sviluppo che di produzione, offrendo diverse opzioni di interazione per soddisfare le esigenze dei clienti. La piattaforma dispone di un'interfaccia utente (UI) intuitiva che consente anche agli utenti meno esperti di sfruttare la maggior parte delle funzionalità disponibili. Le stesse funzionalità presenti nell'UI

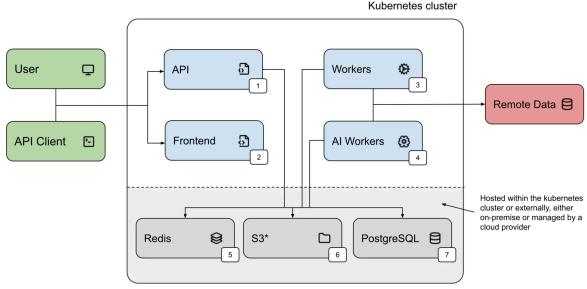
sono accessibili anche tramite una REST API per interazioni macchina-macchina. Inoltre, offriamo un SDK adatto agli sviluppatori. Le funzionalità disponibili nelle diverse opzioni di implementazione sono elencate nella seguente tabella.

ITEM	SDK	UI / API
Generazione di dati sintetici relazionali		
Allenamento con privacy differenziale		
Protezione valori rari		
Gestione vincoli		
Generazione dati parzialmente sintetici		
Generatori pre-addestrati		
Valutazione dei dati sintetici		
Ribilanciamento/curation dei dati		
Anonimizzazione o pseudonimizzazione classica		
Redazione delle PII nei testi e nei documenti	in fase di sviluppo	
Predizione su dati relazionali		
Imputazione intelligente dei Dati Mancanti	in fase di sviluppo	
Estrapolazione su dati relazionali	•	

Caratteristiche tecniche

Design di alto livello

La Piattaforma di Dati Sintetici Aindo è alimentata da container Docker in esecuzione in un cluster Kubernetes, garantendo scalabilità e affidabilità per una varietà di casi d'uso. L'architettura della piattaforma può essere rappresentata schematicamente come segue:



* or compatible object storage provider

Come illustrato nello schema, la piattaforma è composta da più servizi che lavorano insieme per raggiungere modularità, flessibilità e affidabilità. Le componenti principali sono le seguenti:

#	Nome	Descrizione
1	API	Componente che gestisce le richieste API provenienti dai browser degli utenti o dai client API
2	Frontend	Componente che serve le risorse frontend
3	Workers	Pool di worker che gestiscono compiti leggeri o medi
4	Al Workers	Pool di worker che gestiscono compiti di intelligenza artificiale
5	Redis	Cache della piattaforma e message broker
6	S3	Archiviazione della piattaforma
7	PostgreSQL	Database della piattaforma

L'applicazione è progettata per scalare sia orizzontalmente che verticalmente in base alle esigenze. Ecco alcuni esempi di come questa scalabilità può essere sfruttata:

Scalabilità orizzontale per aumentare il numero di lavori di sintesi concorrenti (come l'addestramento e la generazione di dataset sintetici), è possibile aumentare il numero di repliche del componente #4;

Scalabilità orizzontale per supportare un numero maggiore di utenti concorrenti sulla piattaforma, è possibile aumentare le repliche dei componenti #1 e #3;

Scalabilità verticale per sintetizzare dataset più grandi, è possibile potenziare la capacità computazionale dei nodi che ospitano il componente #4.

Grazie alla flessibilità dell'approccio Kubernetes, la piattaforma può essere utilizzata come SaaS, ospitata su cloud privati o installata *on-premise*.



Connettori

La piattaforma offre diversi connettori per l'input/output dei dati. Attualmente, sono disponibili i seguenti connettori per database:

- PostgreSQL
- MySQL
- MariaDB
- Google Big Query
- Microsoft SQL Server
- Oracle Database

Supportiamo inoltre le connessioni ai seguenti sistemi di storage a oggetti:

- Amazon S3
- Google Cloud Storage

Permettiamo inoltre il caricamento di file via browser o via URL

Design di basso livello

Questa sezione fornisce una panoramica tecnica approfondita e di basso livello del design della piattaforma di dati sintetici Aindo. In questa sezione, vengono approfonditi i dettagli circa:

Archiviazione: l'architettura di archiviazione, inclusi i diversi sistemi di archiviazione utilizzati e i tipi di dati gestiti da ciascuno

Flusso di Dati: come i dati si muovono all'interno dell'architettura e dei sistemi di archiviazione durante le interazioni degli utenti con le funzionalità chiave della piattaforma

Archiviazione

La piattaforma di dati sintetici Aindo utilizza tre distinti sistemi di archiviazione: PostgreSQL, S3 e Redis. Questa sezione descrive il tipo di informazioni memorizzate in ciascun sistema.

PostgreSQL

PostgreSQL funge da database relazionale primario, memorizzando metadati essenziali sugli utenti e sulle risorse create all'interno della piattaforma.

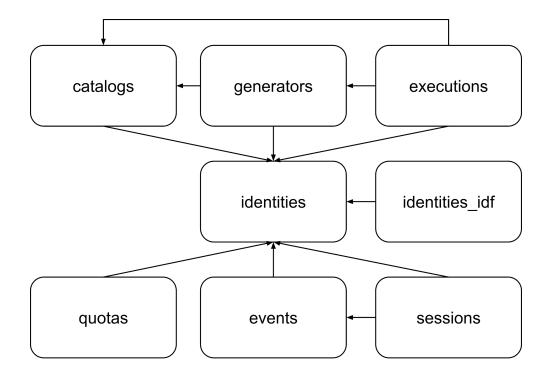
I dati più sensibili memorizzati in PostgreSQL sono le password degli utenti, che vengono hashate in modo sicuro utilizzando l'algoritmo di hashing bcrypt nel campo "identities.pwh".

Elenco delle tabelle del database e del loro schema:

- catalogs: memorizza metadati sulle origini e destinazioni dei dati configurate dagli utenti.
 - Colonne: id (stringa), ownedBy (stringa), createdBy (stringa), createdAt (timestamp), modifiedBy (stringa), modifiedAt (timestamp), name (stringa), description (stringa), type (enumerazione), deletedAt (timestamp), collectAfter (timestamp), hidden (booleano), sample (booleano), shared (booleano)
- generators: memorizza le configurazioni dei generatori di dati sintetici creati dagli utenti
 - Colonne: id (stringa), ownedBy (stringa), createdBy (stringa), createdAt (timestamp), modifiedBy (stringa), modifiedAt (timestamp), name (stringa), description (stringa), sourceId (stringa), deletedAt (timestamp), collectAfter (timestamp), executionCounter (intero), shared (booleano)
- **executions:** tiene traccia dei processi di addestramento e generazione di dati sintetici avviati dagli utenti.
 - Colonne: id (stringa), ownedBy (stringa), createdBy (stringa), createdAt (timestamp), status (enumerazione), generatorId (stringa), sourceId (stringa), destinationId (stringa), deletedAt (timestamp), collectAfter (timestamp), errorCode (enumerazione), modifiedBy (stringa), modifiedAt (datetime), name (stringa), description (stringa), index (intero), resources (lista di stringhe), checkpointVersion (stringhe), train (booleano), generate (booleano)
- **quotas:** tiene traccia del consumo delle quote degli utenti.
 Colonne: owner (stringa), metric (enumerazione), anchor (intero), used (float)
- identities: memorizza informazioni sugli utenti.
 Colonne: id (stringa), traits (json), createdAt (timestamp), modifiedAt (timestamp), pwh (bytes, password hashata con algoritmo bcrypt con sale casuale), email (stringa), avatarUrl (stringa), blockedAt (timestamp), wellcomeAt (timestamp), totp (stringa), isAdmin (booleano), username (stringa), lastActivityAt (timestamp), deletedAt (timestamp), collectAfter (timestamp)

- Identities_idf: Mappa più metodi di autenticazione a un singolo account utente (rilevante quando sono abilitati più metodi di accesso).
 Colonne: id (stringa), createdAt (timestamp), modifiedAt (timestamp), identityId (stringa), type (enumerazione), identifier (stringa)
- sessions: gestisce i dettagli delle sessioni degli utenti.
 Colonne: id (stringa), createdAt (timestamp), modifiedAt (timestamp), identityId (stringa), expiresAt (timestamp), expiresUpdatedAt (timestamp), revokedAt (timestamp), revokeReason (enumerazione), creationLog (stringa), revokeLog (stringa), entropy (stringa)
- events: Registra eventi relativi all'autenticazione (ad esempio, accessi, verifiche OTP).
 - Colonne: id (stringa), identityld (stringa), ipAddress (stringa), userAgent (stringa), type (enumerazione), createdAt (timestamp), location (jsonb)
- api_tokens: tiene traccia dei token API generati dagli utenti.
 Colonne: id (stringa), ownedBy (stringa), createdBy (stringa), createdAt (timestamp), modifiedBy (stringa), modifiedAt (timestamp), expiration (timestamp), name (stringa), description (stringa), enabled (booleano), lastUsed (timestamp), deletedAt (timestamp), collectAfter (timestamp)
- mail: Registra le notifiche email inviate dalla piattaforma (rilevante quando la configurazione della piattaforma richiede notifiche email).
 Colonne: id (stringa), createdAt (timestamp), identityId (stringa), email (stringa), type (enumerazione)
- **sqlrl**: registra eventi di limitazione della velocità. Colonne: id (stringa), at (timestamp), level (float)
- accepts: tiene traccia dell'accettazione dei termini da parte degli utenti (non applicabile per installazioni on-premise).
 Colonne: identityld (stringa), itemld (stringa), itemCode (stringa), choice (booleano), log (stringa)
- alembic_version: gestisce le migrazioni dello schema del database (ambito tecnico)

Le relazioni tra le tabelle chiave del database possono essere rappresentate come segue (per chiarezza, sono illustrate solo le relazioni più critiche):



S3

L'applicazione utilizza S3 per memorizzare file e dati di configurazione, inclusi le configurazioni di origine/destinazione e le configurazioni di sintesi.

I dati memorizzati in S3 che possono essere sensibili sono:

- Configurazioni di origine/destinazione: sono crittografate utilizzando l'algoritmo Fernet (https://github.com/fernet/spec/blob/master/Spec.md), con la chiave di crittografia memorizzata in modo sicuro in un segreto Kubernetes. Questi file sono memorizzati nel bucket "attachments".
- File creati quando un utente richiede un download dei dati (se la funzionalità è abilitata). Questi file sono memorizzati nel bucket "CD", con una politica di ciclo di vita automatica che elimina i file più vecchi di 1 giorno.
- File caricati dagli utenti per creare origini dati basate su file (se la funzionalità è abilitata). Questi file sono memorizzati nei bucket "temporary" (con regole di ciclo di vita per l'eliminazione automatica) e "catalog files".

L'applicazione utilizza quattro distinti bucket S3, ciascuno con uno scopo specifico:

- **temporary:** in questo bucket l'applicazione memorizza temporaneamente i file caricati dagli utenti (utilizzato solo se la funzionalità è abilitata). Una regola di ciclo di vita elimina automaticamente i file più vecchi di 1 giorno.
- **CD:** in questo bucket l'applicazione posiziona i file che possono essere scaricati dagli utenti. Gli utenti possono scaricare dati di origine o sintetici (se la funzionalità è abilitata) e report di sintesi. I file vengono eliminati automaticamente dopo 1 giorno.
- catalog files: in questo bucket l'applicazione salva le origini file (quando un utente crea un'origine caricando file) e i dataset sintetici (quando un utente seleziona l'archiviazione dell'applicazione come destinazione di sintesi)
- attachments: memorizza i file collegati ai record del database SQL
 - Configurazioni di origini e dataset sintetici

- Schema relazionale di origini e dataset sintetici (tabelle, colonne, tipi di dati, chiavi esterne, ecc.)
- Modelli di sintesi Al
- Report di sintesi
- Metadati del processo di addestramento e sintesi (passaggi, tempi, punteggi)
- - Metadati del dataset (numero di record, dimensione della memoria, ecc.)
- Metadati statistici del dataset (istogrammi, valori distinti, valori unici, ecc., se la funzionalità è abilitata)

Redis

Redis è utilizzato per memorizzare dati temporanei e tecnici, inclusi:

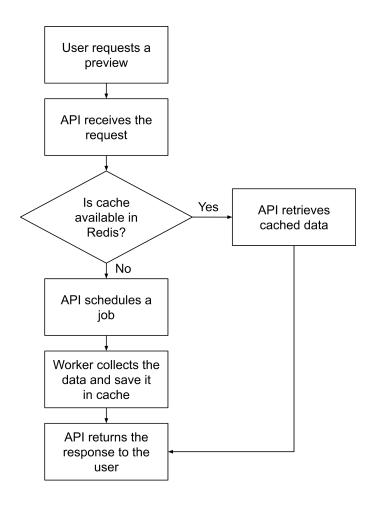
- Metadati dei job: dettagli tecnici come elenco di job, input e output dei job, ecc.
- Cache dell'anteprima dei dati: memorizza temporaneamente una pagina di anteprima con un TTL di 30 minuti (se la funzionalità di anteprima è abilitata).
- Informazioni PII: memorizza temporaneamente i flag PII con un TTL di 30 minuti (se la funzionalità di rilevamento PII è abilitata).

Flusso di Dati

Questa sezione descrive come i dati si muovono attraverso l'applicazione quando gli utenti interagiscono con varie funzionalità della piattaforma. Alcuni flussi sono opzionali, il che significa che si verificano solo se la corrispondente funzionalità è abilitata; questi flussi sono contrassegnati come "OPZIONALI". Sebbene questo elenco non sia esaustivo, fornisce una panoramica schematica di alcuni importanti flussi di dati all'interno della piattaforma.

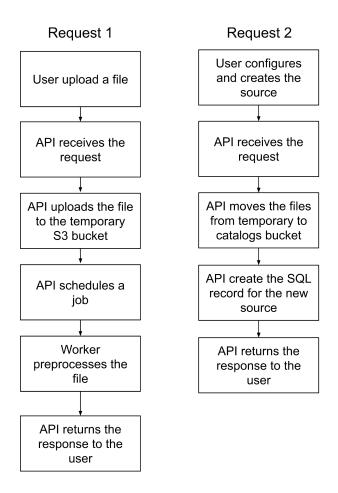
Flusso di Anteprima - OPZIONALE

La funzionalità di anteprima consente agli utenti di visualizzare i dati sia dai dataset di origine che dai dataset sintetici. Per una panoramica dettagliata sulla funzionalità, fare riferimento alla documentazione ufficiale https://docs.aindo.com/sources/view/#view-page. Il flusso per la funzionalità di anteprima è il seguente:



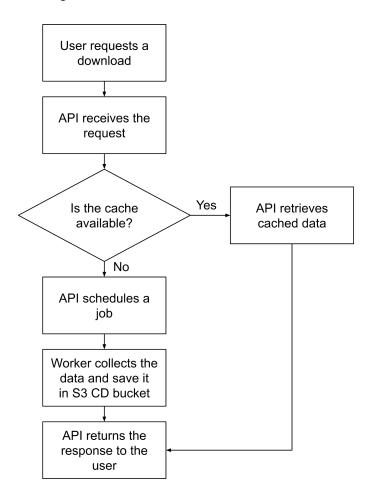
Flusso di Caricamento File - OPZIONALE

La funzionalità di caricamento file consente agli utenti di caricare file strutturati (ad esempio, CSV) e creare un nuovo dataset di origine. Maggiori dettagli possono essere trovati nella documentazione https://docs.aindo.com/sources/create/#file-source. Questo flusso consiste in due passaggi sequenziali: il primo gestisce il caricamento, convalida la sua struttura ed esegue il preprocessing, il secondo crea un nuovo dataset di origine basato sui file caricati e sulla configurazione definita dall'utente.



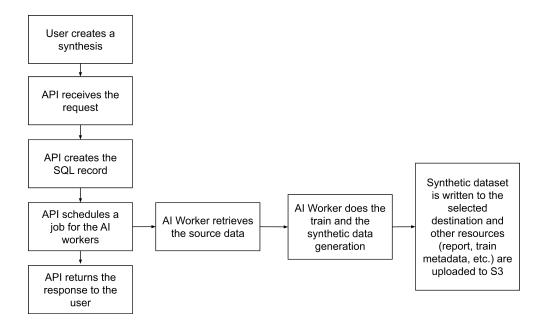
Flusso di Download Dati - OPZIONALE

La funzionalità di download consente agli utenti di esportare file strutturati contenenti dati da un dataset di origine o da un dataset sintetico generato. Ulteriori dettagli possono essere trovati nella documentazione https://docs.aindo.com/sources/manage/#download-source. In questo caso il flusso è il seguente:



Flusso di Addestramento e Sintesi

Il processo di addestramento e sintesi prevede l'addestramento di un modello AI e la generazione di un dataset sintetico. Maggiori informazioni possono essere trovate nella documentazione https://docs.aindo.com/synthesis_platform/create/. Il flusso può essere rappresentato schematicamente come segue:



Scalabilità e Benchmarks

La soluzione può essere eseguita su CPU, ma supporta l'addestramento multi-GPU per i generatori, consentendo la sintesi di dati altamente complessi e voluminosi.

CPU benchmark

In questa sezione, presentiamo i benchmark relativi ai tempi tipici di addestramento e sintesi del modello Aindo su vari dataset. Tutti i test sono stati eseguiti su un'istanza EC2 c6i.4xlarge di AWS, dotata di 16 core CPU e 32 GB di RAM.

Ad eccezione della dimensione del batch, tutti i parametri sono stati impostati ai valori predefiniti sulla piattaforma Aindo.

Dataset: Adult

Il dataset <u>Adult</u>, noto anche come Census Income dataset, è un dataset con un'unica tabella contenente 15 colonne e 29.305 record.

Model	batch	training time	training	time per	generation time
size	size	(sec)	steps	step (sec)	(sec)

Small	256	320	4200	0.076	5
Medium	256	425	3200	0.133	6
Large	256	608	2600	0.233	8

Dataset: Basket

Il dataset <u>Basket</u> è composto da tre tabelle: **"all_star"**, **"season"** e **"players"**. La tabella "players" funge da tabella principale, mentre "all_star" e "season" sono tabelle figlie collegate tramite chiavi esterne. Le dimensioni delle tre tabelle sono:

players: 9 colonne, 4,556 records;
all_star: 7 colonne, 1,487 records;
season: 11 colonne, 2,1415 records;

Model size	batch size	training time (sec)	training steps	time per step (sec)	generation time (sec)
Small	64	1923	5800	0.332	8
Medium	64	2339	3800	0.616	16
Large	64	3312	2800	1.183	18

Dataset: Airbnb

Questo <u>dataset</u> contiene informazioni sugli annunci di Airbnb a New York City, suddivise in due tabelle:

host: 2 colonne, 33712 rows;listings: 13 colonne, 43885 rows.

Model size	batch size	training time (sec)	training steps	time per step (sec)	generation time (sec)
Small	32	1892	18200	0.104	24
Medium	32	2933	16400	0.179	29
Large	32	3526	11200	0.315	39

GPU Benchmarks

In questa sezione, presentiamo i benchmark relativi ai tempi tipici di addestramento e sintesi del modello Aindo eseguito su una macchina virtuale con 4x GPU L40S e una CPU AMD EPYC 9354 a 32 core.

Dataset: Berka

Il dataset <u>Berka</u> è un insieme di informazioni finanziarie provenienti da una banca ceca. Il dataset contiene dati di oltre **5.300 clienti bancari** e circa **1.000.000 di transazioni**.

Questo dataset è stato utilizzato con la seguente configurazione:

• Transazioni massime per cliente: 100

I seguenti parametri sono stati utilizzati per tutte le esecuzioni:

• **Batch Size**: 512

Training steps: 20.000

Device		CPU		GPU		4x GPUs
Model size	time per step (sec)	RAM (GiB)	time per step (sec)	VRAM (GiB)	time per step (sec)	VRAM max per GPU (GiB)
Small	18,6	5,6	0,63	3,4	0,26	4,6
Medium	34,5	7,9	1,09	6,4	0,32	10,4
Large	63,7	14,2	1,91	15,5	0,50	16,0

Dataset: Porto

Un <u>dataset</u> contenente **1.710.671** tragitti in taxi registrati nel corso di un anno (dal **01/07/2013 al 30/06/2014**) nella città di Porto, in Portogallo.

Questo dataset è stato utilizzato con la seguente configurazione:

• Viaggi in taxi: 100.000

Lunghezza massima del tragitto: 100 record

I seguenti parametri sono stati utilizzati per tutte le esecuzioni:

Batch Size: 2048Training steps: 20.000

Device		CPU		GPU		4x GPUs
Model size	time per step (sec)	RAM (GiB)	time per step (sec)	VRAM (GiB)	time per step (sec)	VRAM max per GPU (GiB)
Small	12,0	6,7	0,37	3,2	0,26	4,4
Medium	24,6	9,5	0,70	5,9	0,22	7,1
Large	48,7	14,9	1,30	11,4	0,35	12,6

Sicurezza, Certificazioni e Compliance

Conformità alle Normative (GDPR, Al Act)

La piattaforma Aindo è progettata per aderire alle normative internazionali sulla protezione dei dati e l'intelligenza artificiale:

GDPR: I dati sintetici generati sono considerati anonimi ai sensi del Regolamento (Considerando il recital 26: "The principles of data protection should therefore not apply to anonymous information, namely information which does not relate to an identified or identifiable natural person or to personal data rendered anonymous in such a manner that the data subject is not or no longer identifiable"), eliminando il collegamento diretto con individui reali. I processi di sintesi integrano principi di *privacy by design* e *by default*, garantendo la minimizzazione dei dati e la protezione di informazioni sensibili (es.: dati sanitari).

Al Act: La piattaforma supporta la mitigazione dei rischi legati all'Al (bias, trasparenza) attraverso dataset bilanciati e meccanismi di tracciabilità delle decisioni algoritmiche.

Certificazioni (ISO, EuroPrivacy)

Aindo ha ottenuto certificazioni riconosciute a livello internazionale:

- EuroPrivacy: Prima azienda di dati sintetici certificata in Europa, con focus specifico sul rispetto del GDPR in ambito sanitario.
- ISO 27001: Certificazione per il Sistema di Gestione della Sicurezza delle Informazioni (SGSI), che copre controlli su accessi fisici/logici, gestione degli incidenti e continuità operativa.
- **ISO 9001**: Certificazione del Sistema di Gestione per la Qualità, con audit regolari sui processi di sviluppo e manutenzione della piattaforma.

TOM's (Technical and Organizational Measures)

Aindo implementa varie misure di sicurezza tecniche e organizzative per garantire la riservatezza, l'integrità, l'accuratezza e la disponibilità dei propri dati e di quelli dei suoi clienti. Queste misure sono allineate agli standard **ISO/IEC 27001**, al **GDPR** e ad altre normative sulla protezione dei dati. Tra le più rilevanti elenchiamo

Sviluppo Sicuro del Software

 Le pratiche di sviluppo sicuro del software sono integrate nell'architettura di sicurezza di Aindo, aderendo a linee guida riconosciute come OWASP, CWE, MISRA, AUTOSAR e CERT.

- Ogni commit è soggetto a revisione indipendente da parte di un altro sviluppatore.
 Sia test automatizzati (smoke test, unit test) che test manuali vengono utilizzati per garantire la qualità del codice.
- L'ambiente di produzione contiene solo codice che è stato accuratamente analizzato e testato.

Valutazione delle Vulnerabilità e Test di Penetrazione

 Aindo definisce linee guida per la pianificazione e l'implementazione di attività di valutazione delle vulnerabilità e test di penetrazione, con l'obiettivo di identificare, valutare e mitigare le vulnerabilità presenti nei sistemi informativi e nelle reti aziendali.

Queste misure vengono monitorate e migliorate costantemente per adattarsi alle minacce emergenti e garantire la protezione continua degli asset informativi di Aindo.

Gestione delle terze parti

- Viene seguito un processo sistematico di valutazione del rischio per ogni potenziale fornitore per garantire la protezione adeguata delle informazioni e dei dati.
- Il Processore e gli eventuali Sub-Processori si impegnano a implementare tecniche di crittografia a tutti i livelli, utilizzando standard crittografici non deprecati.

Gestione degli Incidenti

- Gli incidenti relativi alla sicurezza delle informazioni vengono segnalati immediatamente attraverso i canali di comunicazione stabiliti.
- Aindo fornisce una risposta tempestiva, efficace e coordinata agli incidenti che potrebbero compromettere la sicurezza delle reti e dei sistemi informativi.
- Il processo di gestione degli incidenti include rilevamento, presa in carico, analisi, contenimento, ripristino, chiusura e registrazione dell'incidente.
- Sono previste procedure cicliche per l'analisi e la verifica del comportamento dannoso/presenza di malware.
- Aindo protegge la propria rete dalle minacce online adottando misure di sicurezza strutturate per garantire la protezione dell'infrastruttura tecnologica.

Sicurezza e Crittografia

La nostra piattaforma cripta tutti i dati sia a riposo che in transito, garantendo alti livelli di sicurezza e impedendo l'accesso non autorizzato ai dati. Offre inoltre diversi tipi di autenticazione, includendo sistemi di Oauth e include l'opzione per l'autenticazione a **due fattori (2FA)**. Tutti gli accessi e le operazioni effettuate sulla piattaforma sono tracciate e conservati sui log della piattaforma.

Servizi

Aindo offre una gamma completa di servizi professionali a supporto dell'adozione, dell'integrazione e dell'operatività della nostra piattaforma di dati sintetici. I nostri servizi sono progettati per garantire un'implementazione senza problemi, un'assistenza continua e un supporto strategico nell'uso della tecnologia.

Setup

Per facilitare l'implementazione della nostra soluzione, offriamo servizi di installazione professionale personalizzati in base alle esigenze del cliente. Il setup può avvenire sia **on-premise** (all'interno dell'infrastruttura IT del cliente) sia in **cloud** (su ambienti AWS, Azure, GCP o altri provider).

I nostri servizi di setup includono:

- Installazione e configurazione della piattaforma nell'ambiente scelto.
- **Supporto all'integrazione** con i sistemi esistenti, inclusi database, strumenti di analisi e infrastrutture AI/ML.
- **Verifica delle performance** e ottimizzazione per garantire il corretto funzionamento e l'efficienza della piattaforma.
- Assistenza per la conformità normativa, garantendo che la soluzione sia allineata alle politiche aziendali di sicurezza e protezione dei dati.

Supporto Cliente

Questa sezione descrive il supporto tecnico fornito da Aindo al Cliente. Definisce le procedure per le richieste di supporto al fine di garantire la collaborazione tra le parti e la tempestiva risoluzione dei problemi.

Richieste di Supporto

Alla scoperta di un problema, il Cliente deve informare tempestivamente Aindo tramite l'indirizzo email dedicato al supporto: support@aindo.com.

Ogni richiesta di supporto deve includere tutte le informazioni necessarie affinché Aindo possa identificare, mitigare e risolvere correttamente il problema. Ogni richiesta di supporto deve contenere almeno:

- Livello di severità (definito di seguito)
- Descrizione dettagliata del problema (descrizione, screenshot, log)
- Informazioni di supporto, come screenshot e/o log
- Riferimento al sistema che presenta il problema
- Informazioni di contatto dell'utente che richiede il supporto

Durante la risoluzione della richiesta di supporto, in base al livello di severità, il Cliente deve garantire la disponibilità ragionevole dei propri rappresentanti per fornire tutte le informazioni richieste da Aindo per indagare adeguatamente il problema.

Il team di supporto di Aindo farà il possibile per fornire una prima risposta in conformità con gli Obiettivi del Livello di Servizio (definiti di seguito). Aindo si impegnerà a diagnosticare il problema e a fornire una soluzione, che potrebbe includere l'eliminazione del difetto, il rilascio di aggiornamenti o la dimostrazione al Cliente di come evitare il problema. Tuttavia, alcuni problemi potrebbero non essere risolvibili a livello di supporto. In questi casi, il team di supporto avvierà un'escalation interna per i ticket Business Critical e High Priority, assicurando che la correzione del difetto abbia la massima priorità. Il Cliente sarà aggiornato sui progressi attraverso il canale email di supporto.

Lingua

La lingua predefinita per il supporto è l'inglese. Su richiesta, è disponibile anche il supporto in italiano.

Livelli di Severità

Il supporto fornito da Aindo si basa sul livello di severità della richiesta di supporto:

Livello	Descrizione	
Critico / Business Critical	Le funzionalità critiche sono compromesse in modo irreversibile, rendendo il servizio completamente inutilizzabile e influenzando l'uso regolare.	
Maggiore / Servizio Degradato	Alcune funzionalità critiche sono compromesse, causando un'uso parziale del servizio e influenzando aspetti minori dell'uso regolare.	
Minore / Problema Generale	Funzionalità non critiche sono compromesse o il problema può essere temporaneamente evitato dal Cliente, causando una degradazione minore del servizio.	
Miglioramento	Le funzionalità non sono compromesse, ma la risoluzione del problema migliorerà il servizio.	

Il Cliente è invitato a proporre un livello di severità per ogni richiesta di supporto. In ogni caso, Aindo valuterà il livello di severità e si riserva il diritto di modificarlo.

Obiettivi del Livello di Servizio

La casella di posta elettronica **support@aindo.com** è monitorata durante l'orario lavorativo italiano (dalle 9:00 alle 17:00, dal lunedì al venerdì). Per le richieste di supporto ricevute al di

fuori dell'orario lavorativo, non può essere garantita alcuna azione fino al giorno lavorativo successivo.

Le richieste di supporto sono elaborate secondo i seguenti obiettivi:

Priorità	Tempo di presa in carico
Critico	4 ore
Maggiore	24 ore
Minore	2 giorni

Il tempo di presa in carico si riferisce al tempo necessario per riconoscere la ricezione della richiesta.

Ambito

In Aindo, la soddisfazione del cliente è la nostra massima priorità e ci impegniamo a fornire un supporto solido per il nostro software. In quanto sviluppatore e fornitore di software, il supporto tecnico copre esclusivamente le funzionalità del software di Aindo. La tabella seguente illustra quali componenti della soluzione sono gestiti e supportati da Aindo e quali aspetti sono responsabilità del Cliente o di terzi. Il nostro obiettivo è garantire una collaborazione fluida ed efficiente per soddisfare le esigenze dei nostri clienti.

Parte	Ambito di Supporto
Funzionalità del software AINDO	Aindo
Aggiornamenti del software AINDO	Aindo
Patch di sicurezza del software AINDO	Aindo
Autenticazione del software	Condiviso
HTTPS / SSL / TLS del software	Condiviso

Le parti con un Ambito di Supporto Condiviso coinvolgono sistemi e servizi di entrambe le parti. Pertanto, le parti devono concordare la collaborazione e la manutenzione regolare dei sistemi per garantire il corretto funzionamento della parte coinvolta.

Richieste di Modifica

Le richieste per modificare una funzionalità esistente o per introdurre nuove funzionalità sono ben accette e il Cliente è incoraggiato a fornire tali feedback ad Aindo. Le richieste di

modifica devono essere inviate tramite lo stesso indirizzo email dedicato al supporto: support@aindo.com e saranno opportunamente segnalate.

Le richieste di modifica saranno valutate dal team di prodotto e, se ritenute appropriate, saranno incluse in un prossimo aggiornamento del prodotto.

Le richieste di modifica non seguono la procedura di supporto e, pertanto, non sono soggette agli Obiettivi del Livello di Servizio, né il team di supporto è obbligato a tenere aggiornato il Cliente sul rilascio della modifica. Tuttavia, il team di prodotto potrebbe informare il Cliente se la modifica è stata presa in carico e, quando possibile, fornire una data stimata di rilascio.

Gestione Operativa della Piattaforma

Aindo offre un servizio di gestione operativa per la propria piattaforma di dati sintetici, consentendo ai clienti di delegare l'esecuzione di attività chiave legate all'utilizzo dello strumento. Questo servizio garantisce un uso efficiente e continuativo della piattaforma, senza necessità di risorse interne dedicate da parte del cliente. Il servizio, basato su di un monte ore, può prevedere le seguenti attività:

Caricamento Dati

- Ricezione e validazione dei dati forniti dal cliente.
- Verifica dell'integrità, coerenza e conformità ai formati richiesti dalla piattaforma.
- Pre-processing e trasformazione dei dati per la compatibilità con la piattaforma.
- Ingestione dei dati nella piattaforma, garantendo conformità agli standard richiesti.

Integrazione

 Supporto nella definizione e implementazioni di progetti di integrazione con i sistemi del cliente o di terze parti.

Generazione di Dati Sintetici

- Collaborazione con il cliente per identificare la modalità di sintesi più adatta agli obiettivi del particolare caso d'uso.
- Configurazione dei parametri di sintesi in base ai requisiti del cliente (dimensione del modello, gestione di valori rari, opzioni di ribilanciamento etc...).
- Applicazione di vincoli di business (es.: relazioni matematiche tra colonne, regole di dominio specifico).
- Esecuzione del processo di sintesi dei dati.
- Analisi statistica per verificare la fedeltà ai dati originali (distribuzioni, correlazioni).
- Test di privacy per assicurare l'assenza di rischi di re-identificazione.
- Produzione di report dettagliati (formato PDF/CSV) con metriche di valutazione.
- Formattazione dei dati sintetici secondo le esigenze del cliente.

Esecuzione di Analisi e Scenari Predittivi

- Collaborazione con il cliente per identificare scenari di analisi (es.: simulazioni what-if, previsioni di tendenza).
- Configurazione ed esecuzione di modelli predittivi e simulazioni "what-if".
- Verifica di performance dei modelli predittivi sviluppati.
- Elaborazione di report e dashboard per la visualizzazione dei risultati.
- Produzione di report analitici con insights quantitativi e raccomandazioni operative.

Guida all'uso

Aindo mette a disposizione una documentazione completa per facilitare l'uso della piattaforma, disponibile pubblicamente all'indirizzo https://docs.aindo.com/.

La documentazione include:

- Guida utente dell'interfaccia per l'uso della piattaforma tramite UI.
- Documentazione dell'SDK per l'integrazione con flussi di lavoro esistenti.
- Esempi pratici e best practice per la generazione e valutazione di dati sintetici.
- FAQ e troubleshooting, con soluzioni ai problemi più comuni.

Formazione

Per accelerare l'adozione della piattaforma e rendere i clienti autonomi nella gestione dei dati sintetici, Aindo offre servizi di formazione su richiesta che possono coprire:

- **Funzionalità base**: introduzione alla piattaforma, sintesi dei dati e utilizzo delle funzionalità principali.
- **Funzionalità avanzato**: configurazione avanzata, gestione di dataset complessi e valutazione dei dati sintetici.
- Corsi per sviluppatori: utilizzo dell'SDK e API per integrare i dati sintetici nei flussi aziendali.
- Workshop personalizzati, basati sulle esigenze specifiche del cliente.

I training possono essere erogati in modalità **remota o in presenza**.

Casi d'Uso

Benefici per le Pubbliche Amministrazioni e gli Enti Regionali

La piattaforma Aindo offre un sistema avanzato per la gestione dei dati sintetici, garantendo **privacy-by-design**, conformità normativa e un ecosistema tecnologico unificato per l'anonimizzazione e la predizione. Questo approccio consente agli enti pubblici di:

- Ottimizzare l'uso dei dati senza violare la privacy dei cittadini.
- Ridurre i costi operativi eliminando la necessità di strumenti separati per l'anonimizzazione e l'analisi predittiva.
- Incrementare la sicurezza e la scalabilità grazie a un'infrastruttura integrata e adattabile alle esigenze delle amministrazioni regionali.

L'adozione di Aindo consente di migliorare l'efficienza e la qualità della gestione pubblica in diversi ambiti. Ecco alcuni esempi pratici:

1. Studi di Real World Evidence (RWE)

L'uso di dati sintetici permette di condurre studi clinici e ricerche epidemiologiche senza compromettere la privacy dei pazienti.

 Esempio: Un'agenzia sanitaria regionale utilizza dati sintetici per analizzare l'efficacia di un nuovo trattamento oncologico, accelerando la ricerca senza necessità di accedere a dati sensibili.

2. Piattaforme di Sperimentazione

Creazione di ambienti digitali sicuri per testare soluzioni innovative in ambito healthcare, riducendo tempi e costi di sviluppo.

• **Esempio**: Un ospedale regionale utilizza Aindo per testare algoritmi di diagnosi assistita senza compromettere la privacy dei pazienti reali.

3. Ottimizzazione della Spesa Sanitaria

Previsione della domanda di farmaci e risorse sanitarie per migliorare la sostenibilità del sistema.

• **Esempio**: Analisi predittiva per stimare il fabbisogno di vaccini anti-influenzali, evitando sprechi e carenze.

4. Pianificazione delle Risorse Sanitarie

Utilizzo dell'Al per gestire al meglio letti ospedalieri, turni del personale e campagne sanitarie.

• **Esempio**: Previsione dei picchi di ricoveri durante l'inverno per ottimizzare la distribuzione del personale medico.

Applicazioni in Al/ML e Analisi Dati

1. Generazione di Dataset Anonimi

Creazione di versioni sintetiche di dataset che preservano i pattern originali, garantendo allo stesso tempo la privacy dei soggetti originali.

• **Esempio**: Anonimizzazione dei dati dei pazienti per consentire studi medici senza rischi di re-identificazione.

2. Condivisione Sicura dei Dati per la Ricerca

Possibilità di condividere dataset realistici senza esporre informazioni personali, favorendo la collaborazione tra istituzioni.

 Esempio: Un'università condivide dati sintetici sui pazienti con un'azienda farmaceutica per lo sviluppo di nuovi farmaci.

3. Aumento e Ribilanciamento dei Dati

Riequilibrio delle categorie di un dataset per ridurre i bias e migliorare la qualità delle analisi.

• **Esempio**: Un dataset medico con un numero sbilanciato di uomini e donne viene riequilibrato per addestrare un algoritmo di diagnosi equo e rappresentativo.

4. Modellazione Predittiva e Simulazione di Scenari

Previsione di tendenze future basate su dati storici, utile per il monitoraggio sanitario e la gestione delle risorse.

 Esempio: Un ente sanitario utilizza Aindo per stimare il rischio di malattie cardiovascolari in diverse fasce di popolazione, migliorando le campagne di prevenzione.

Protezione delle Informazioni Personali nei Testi e Documenti

Rilevamento e anonimizzazione automatica delle informazioni personali (PII) all'interno di documenti testuali.

Esempi:

- Rimozione di dati identificativi dalle cartelle cliniche elettroniche per consentire analisi senza rischi per la privacy.
- Mascheramento di nomi e indirizzi nei documenti legali prima della loro pubblicazione o condivisione con terze parti.

Licenza

La piattaforma di Aindo è disponibile attraverso un modello di sottoscrizione annuale con durata minima di 36 mesi, strutturato su due livelli principali: sottoscrizione sperimentale e sottoscrizione commerciale Enterprise. Questa suddivisione consente di rispondere a esigenze differenti, bilanciando flessibilità nella fase di testing e valutazione con una proposta commerciale solida e scalabile per l'adozione su larga scala.

Sottoscrizione Sperimentale

La sottoscrizione a scopo sperimentale è pensata per le aziende che desiderano testare le funzionalità della piattaforma prima di un'eventuale adozione in produzione. Questa opzione è fornita con finalità di testing e valutazione e non include elementi avanzati come supporto tecnico prioritario o aggiornamenti continuativi. È quindi indicata per realtà che necessitano di un primo approccio alla tecnologia, senza un immediato impegno a lungo termine.

Sottoscrizione Commerciale Enterprise

Per le organizzazioni che intendono integrare stabilmente la piattaforma nei propri flussi di lavoro, è disponibile la sottoscrizione commerciale Enterprise della durata di 36 mesi. Questo modello di sottoscrizione garantisce:

- Accesso completo a tutte le funzionalità della piattaforma, incluse quelle avanzate per la generazione e gestione di dati sintetici ed estrapolazione.
- Manutenzione continuativa per garantire l'aggiornamento costante della piattaforma con le ultime innovazioni tecnologiche e normative.
- Supporto tecnico standard incluso nella sottoscrizione, con possibilità di escalation in caso di problematiche critiche.

Servizio di Gestione Operativa (Opzionale)

Per supportare il cliente nelle attività di gestione quotidiana della piattaforma, Aindo offre una subscription opzionale per la di gestione operativa della piattaforma. Questo servizio consente ai clienti di concentrarsi sulle proprie attività senza doversi occupare direttamente degli aspetti tecnici e infrastrutturali, garantendo una gestione ottimale della soluzione Aindo.

Monte Ore per Sviluppo e Consulenza (Opzionale)

Oltre alla sottoscrizione della piattaforma e al servizio di gestione operativa, è possibile accedere a un monte ore opzionale, che consente di usufruire di figure professionali dedicate per esigenze specifiche di:

- Sviluppo personalizzato, per adattare la piattaforma a specifiche necessità tecniche e di integrazione con altri sistemi aziendali.
- Consulenza specialistica, per supporto strategico e operativo nell'uso dei dati sintetici e nell'ottimizzazione dei flussi di lavoro.

Il monte ore viene definito su richiesta del cliente e garantisce accesso a risorse qualificate in grado di rispondere a esigenze avanzate che vanno oltre i servizi inclusi nella sottoscrizione standard.