

Caso d'uso

è necessaria l'AI per realizzare questo obiettivo?

Canvas: AI per la PA

Una guida per l'adozione etica di strumenti AI a sostegno della PA

Canvas: AI per la PA

1 Dati

Questa sezione è per una pianificazione generale dell'algoritmo/modelo. Per specifiche di maggior dettaglio utilizzare le AI product cards

| | | | |
|---|---|--|--|
| ATTIVO CONSENTITI NELLA RACCOLTA DATI Quali sono i dati raccolti? Chi sono i soggetti della raccolta dati? (se presenti) Chi raccoglie i dati? Si corrisponde alle finalità raccolte? Con quali obiettivi sono stati raccolti i dati da analizzare? (Identificare anche scenari di riuso) | Possono esserci bias nei dati raccolti? Etica [x] [] Sicurezza [x] [] Classe sociale [x] [] Altro: [x] [] AI training data set (se presenti) [x] [] Misurare i bias nei dati Qual è la base giuridica del trattamento dei dati personali? (vedi tabella, confermare con DPO) | SONO STATE TECNICHE DI ANONIMIZZAZIONE E PSEUDONIMIZZAZIONE PER LA RACCOLTA DEI DATI PERSONALI RACCOLTI? Esistono ulteriori dataset che possono generare conseguenze indesiderate (es. re-identificazione di soggetti anonimizzati)? Quali precauzioni sono state prese? | Protezione dei dati (raccolti e terzi) Anonimizzazione Effacement Classificazione Altro Finalizzazione dati e processi (Data stewardship) Chi è il data steward? I dati sono accessibili a tutti i livelli di legge? I dati sono ben documentati? (metadata, lineage...) I dati sono ben definiti? (scopo, scopo...) Non c'è la possibilità di certificato? Chi si occupa dell'addestramento del modello? Separazione dei dati [x] [] Pubblicazione [x] [] Segregazione dei dati [x] [] Con quali modalità sono aggiornati i dati? Si effettua un backup crittografato? |
|---|---|--|--|

Canvas AI Ethics © 2024 by Provincia Autonoma di Trento e Fondazione Bruno Kessler. Is licensed under CC BY-SA 4.0

Canvas: AI per la PA

2 Algoritmi

Questa sezione è per una pianificazione generale dell'algoritmo/modelo. Per specifiche di maggior dettaglio utilizzare le AI product cards

| | |
|---|---|
| Qual è la principale funzione dello strumento AI? [] Predizione [] Raccomandazione [] Analisi di impatto ex ante (Indicare una o più opzioni con una crocetta) [] Analisi di impatto ex post [] Altro | Definisci il tipo di strumento AI [] Chatbot [] Strumento predittivo [] Altro |
| Quali sono i principali rischi presenti nell'algoritmo? Sicurezza: quale categoria di rischio si classifica lo strumento AI (vedi AI Act)? [] Tecnico [] Legale [] Etico [] Inaccettabile C'è la possibilità di certificato? [] AI Act [] GDPR [] Altro Non c'è la possibilità di certificato | Vanno attribuiti i guai per la categoria sottorappresentata nel training dataset? Etica [x] [] Sicurezza [x] [] Classe sociale [x] [] Altro [x] [] Il risultato finale dell'algoritmo è originale (copyright) e verificato/corretto/modificato da figure esperte umane? In caso di riuso di un algoritmo, assicurarsi che il contesto di riuso e di prima implementazione siano equivalenti IL RISULTATO E IL PROCESSO DELL'ALGORITMO SONO SPERABILI? L'ALGORITMO PRESENTA I PROCESSI DI BIAS/ASSOCIAZIONE DI DATI TALI DA RENDERE I SOGGETTI VULNERABILI? |
| Il training dataset contiene dei bias? I soggetti dei dati del training dataset sono reidentificabili? | |

Canvas AI Ethics © 2024 by Provincia Autonoma di Trento e Fondazione Bruno Kessler. Is licensed under CC BY-SA 4.0

Canvas: AI per la PA

3 Metodi di analisi

| | | |
|---|--|---|
| CHI PRIME NECESSITÀ DELLA BASE DELL'ALGORITMO NE SA COMPRENDERE IL PROCESSO LOGICO? Comprensione del funzionamento logico dell'algoritmo [x] [] Se l'algoritmo non è spiegabile, fornire strumenti di explainability (e relative istruzioni) all'utente Valutazione della generalità delle conseguenze [x] [] Processo di valutazione del feedback loop [x] [] Processo di valutazione del feedback loop [x] [] | Quali attività di formazione bisogna svolgere per riconoscere i bias nell'algoritmo e nei dati? Chi utilizza lo strumento AI lo fa nell'ambito delle proprie mansioni pre-esistenti o acquisisce nuove mansioni con l'introduzione dello strumento? (Questo può influenzare il tempo a disposizione per l'apprendimento dello strumento) Pianificare le seguenti attività: Post processing dell'algoritmo per verificare bias di dati Verifica della qualità dei dati pre-processing | CHI SI ASSUME LA RESPONSABILITÀ DELLA DECISIONE FINALE PRESA DALLA BASE DELL'ALGORITMO NEI SOGGETTI CASI? Identificare le figure responsabili in caso di violazioni nei seguenti ambiti: Non discriminazione Privacy Tutela del copyright Altro |
|---|--|---|

Canvas AI Ethics © 2024 by Provincia Autonoma di Trento e Fondazione Bruno Kessler. Is licensed under CC BY-SA 4.0

Canvas: AI per la PA

4 Elementi culturali e sociali

Questa sezione è per pianificare la comunicazione verso la cittadinanza sull'uso dell'AI e di dati personali da parte della PA

| | |
|---|---|
| SONO PRESENTI DATI DA TERZE PARTI CHE POSSONO CREARE CONSEGUENZE ETICHE INDESIDERATE (ES. REIDENTIFICAZIONE DI SOGGETTI DI CUI SI SAIA RACCOLTI DATI PERSONALI)? [] Sì, Quali? [] No QUALI PRECAUZIONI SONO PRESE PER LIMITARE LE POTENZIALI CONSEGUENZE INDESIDERATE? | SONO INCLUSI I SOGGETTI A RISCHI NELL'INFORMATICA DEL TRATTAMENTO DEI DATI PERSONALI? Quali categorie di dati sono raccolti? Per quali finalità sono raccolti i dati? Quali sono le basi giuridiche su cui si fonda il trattamento dei dati? Con quali modalità sono raccolti i dati? |
| I DATI RACCOLTI POSSONO ESSERE ELASCATI COME OPEN DATA IN FORMA AGGIORNATA E? | NEL RACCOLTARE I DATI PERSONALI OCCORRE ELABORARE: 1. Un'informazione contenente tutte le info necessarie ai fini del GDPR 2. Una comunicazione dal linguaggio semplice e comprensibile a chiunque per favorire la piena comprensione del trattamento dati da parte dei soggetti interessati |

Canvas AI Ethics © 2024 by Provincia Autonoma di Trento e Fondazione Bruno Kessler. Is licensed under CC BY-SA 4.0

Canvas: AI per la PA

5 Requisiti funzionali

Sulla base delle risposte alle domande precedenti, elencare i requisiti funzionali da implementare nella piattaforma AI

| | | |
|--|--|---|
| Implementazione dei principi tecnici Notifiche sulla completezza e sulla correttezza in anticipo dei dati Inclusione di strumenti di explainability Notifiche sulla sostituzione degli output generati Notifiche sulle modifiche ai dati operanti dati Altro | Implementazione dei principi sociali e culturali Notifiche sul trattamento per utenti generici Notifiche sulla correttezza tra i dati della struttura e l'utente per cui è stato creato Notifiche sui potenziali rischi privacy (es. re-identificabilità) Altro | Quali dei seguenti principi sono prioritari nell'implementazione dello strumento di elaborazione? Indicare un ordine di priorità Privacy Equa distribuzione di benefici e oneri Non discriminazione su base etica e di genere Altro (es. sostenibilità ambientale) |
|--|--|---|

Considerare i seguenti strumenti per la documentazione (e le istruzioni per l'utente): Data cards, Model cards, AI Product cards. Utilizzare AI Product cards per specificare le caratteristiche tecniche del modello abbozzate in questa pagina del canvas. Utilizzare l'Assessment List for Trustworthy Artificial Intelligence (ALTAI) per una seconda valutazione della coerenza dell'algoritmo con le norme UE.

Canvas AI Ethics © 2024 by Provincia Autonoma di Trento e Fondazione Bruno Kessler. Is licensed under CC BY-SA 4.0

Legenda



AI scientist

AI engineer

AI user

Come si compila il canvas

Compilare il canvas per ogni strumento AI in elaborazione all'inizio della fase di sviluppo.
Utilizzare le info inserite nel canvas per compilare le AI Product Cards (vedere allegato) a prodotto sviluppato.

Legenda

- **AI scientist**: ricerca e sviluppa i sistemi AI (può essere esterno all'ente)
- **AI engineer**: allena, implementa, mantiene i sistemi AI
- ▲ **AI user**: usa i sistemi AI sul lavoro (es. assessore, funzionario)

Indicazioni per la lettura

1. Le domande sono in riquadri di colori diversi, i quali corrispondono ai colori nella legenda per indicare quale profilo debba affrontare la domanda in questione.
2. Ogni colore è accompagnato da una forma geometrica. Se chi utilizza il canvas è daltonico, può fare riferimento alle forme invece che ai colori.
3. I riquadri di colore misto (e senza forme di riferimento) indicano domande che devono essere affrontate da tutti e tre i profili insieme.

Consiglio d'uso

Consigliamo di indicare una figura moderatrice a inizio lavori. Quando non si conosce la risposta ad una domanda, invitare al "tavolo" nuove figure professionali.

SEZIONI del CANVAS

Dati

Per ottimizzare i dati, ridurne i bias e tutelare la privacy

Algoritmi

Per definire le regole per ottimizzare l'automazione delle decisioni, riducendo i bias.

Metodi di analisi

Per interpretare in modo critico il risultato del processo algoritmico e decidere se e come applicarlo

Elementi sociali e culturali

Per assicurarsi di aver comunicato correttamente al pubblico la natura del trattamento dei dati

Requisiti funzionali

Per definire le caratteristiche dello strumento AI che si sta sviluppando

Glossario

Un glossario del gergo tecnico (spesso ricco di termini inglese e anglicismi) per l'utente che non ha competenze specifiche in campo AI. Sono omessi per brevità i termini inglesi di uso comune (es. privacy o dataset)

AI Product Cards

Documentazione di uno strumento AI che ne illustra le caratteristiche dei dati per l'addestramento dell'algoritmo, dei modelli e degli usi previsti.

Simile: "data cards"/"model cards" per dataset e modelli.

Bias

Pregiudizio che il creatore di uno strumento AI può inavvertitamente trasferire all'AI nel suo funzionamento.

Clusterizzazione

Raggruppamento dei dati in insiemi prestabiliti in base alle caratteristiche dei dati stessi (metriche di vicinanza).

Data steward

Figura incaricata di aggiornare uno o più dataset all'interno di un'organizzazione, aggiornandone anche i metadati e assicurandone l'usabilità.

Explainability

Spiegabilità del processo logico e del risultato di un algoritmo.

Feedback loop

Quando i risultati del processo algoritmico sono immessi nella successiva iterazione del processo algoritmico.

Metadati

Informazioni riguardo ai dati contenuti in un dataset (esempio: periodicità dell'aggiornamento, formato, granularità).

Multifactor Authentication

Autenticazione a più fattori per accedere a una piattaforma, a un file o a un dataset.

Post-processing

Affinamento del dettaglio di un insieme di dati.

Training dataset

Insieme di dati utilizzato per allenare (training) un algoritmo di AI.

Segregazione del dato

Separazione logica o fisica dei dati, limitandone l'accesso solo a utenti, sistemi o processi autorizzati.

1

Canvas: AI per la PA

Dati

ATTORI COINVOLTI NELLA RACCOLTA DATI

| | |
|--|--|
| Quali sono i dati raccolti? | |
| Chi sono i soggetti della raccolta dati? (se presenti) | |
| Chi raccoglie i dati? | |
| Si sovrappongono competenze nella raccolta dati? | |

Con quali obiettivi sono stati raccolti i dati da analizzare? (Identificare anche scenari di riuso)

Possono esserci bias nei dati raccolti?

| | |
|-----------------------------|------------------------------------|
| Etnia [si] [no] | |
| Genere [si] [no] | |
| Classe sociale [si] [no] | |
| Altro: _____ | Ad esempio, attori economici (PMI) |

Misurare i bias nel dataset

Qual è la base giuridica del trattamento dei dati personali?
(Nel dubbio, confrontarsi col DPO)

SONO USATE TECNICHE DI ANONIMIZZAZIONE O PSEUDONIMIZZAZIONE PRIMA DELL'USO DEI DATI PERSONALI RACCOLTI?

Esistono ulteriori dataset che possono generare conseguenze indesiderate (es. re-identificazione di soggetti anonimizzati)?

Quali precauzioni sono state prese?

Protezione dei dati (raccolti e terzi)

| | |
|------------------|--|
| Anonimizzazione | |
| Offuscamento | |
| Clusterizzazione | |
| Altro | |

L'accesso ai dati è protetto da:

| | |
|---|--|
| Crittografia [si] [no] | |
| Segregazione del dato [si] [no] | |
| Multifactor authentication [si] [no] | |

Monitoraggio dati e processi (Data stewardship)

| | |
|--|--|
| Chi è il data steward? | |
| I dati sono accessibili a tutti (salvo limiti di legge)? | |
| I dati sono ben documentati? (metadati, ontologie...) | |
| Con quale periodicità sono aggiornati i dati? | |
| si effettua un backup crittografato? | |

2

Una guida per l'adozione etica di strumenti AI nella PA

Algoritmi

Questa sezione è per una pianificazione generale dell'algoritmo/modello. Per specifiche di maggior dettaglio utilizzare le **AI product cards**

Qual è la principale funzione dello strumento AI?

(Indicare una o più opzioni con una crocetta)

[] Predizione [] Raccomandazione [] Analisi di impatto ex ante
[] Analisi di impatto ex post [] Altro _____

Definisci il tipo di strumento AI

[] Chatbot [] Strumento predittivo [] Altro_____

Secondo quale categoria di rischio è classificabile lo strumento AI (vedi AI Act)?

[] minimo
[] limitato
[] elevato
[] inaccettabile

QUALI BIAS POSSONO ESSERE PRESENTI NELL'ALGORITMO?

| | |
|-----------------------------|--|
| Etnia [si] [no] | |
| Genere [si] [no] | |
| Classe sociale [si] [no] | |
| Altro: _____ [si] [no] | |

Vanno attribuiti i giusti pesi alle categorie sottorappresentate nel training dataset?

| | |
|-----------------------------|--|
| Etnia [si] [no] | |
| Genere [si] [no] | |
| Classe sociale [si] [no] | |
| Altro: _____ [si] [no] | |

Il risultato finale dell'algoritmo è originale (copyright) e verificato/curato/modificato da figure esperte umane?

In caso di riuso di un algoritmo, assicurarsi che il contesto di riuso e di prima implementazione siano equivalenti

IL RISULTATO E IL PROCESSO DELL'ALGORITMO SONO SPIEGABILI?

L'ALGORITMO PREVIENE I PROCESSI DI RI-ASSOCIAZIONE DI DATI TALI DA RENDERE I SOGGETTI RICONOSCIBILI?

CHI SI OCCUPA DELL'ADDESTRAMENTO DEL MODELLO?

Definisci gli obiettivi dell'addestramento del modello

Il training dataset contiene dei bias?

I soggetti dei dati del training dataset sono reidentificabili?

Metodi di analisi

CHI PRENDE DECISIONI SULLA BASE DELL'ALGORITMO NE SA COMPRENDERE IL PROCESSO LOGICO?

| | |
|--|---|
| Comprensione del funzionamento logico dell'algoritmo | [si][no] Se l'algoritmo non è spiegabile, fornire strumenti di explainability (e relative istruzioni) all'utente |
| Valutazione delle potenziali conseguenze | [si][no] |
| Processo di valutazione dei feedback loop | [si][no] Fornire strumenti di prevenzione |

Quali attività di formazione bisogna svolgere per riconoscere i bias nell'algoritmo e nei dati?

Chi utilizza lo strumento AI lo fa nell'ambito delle proprie mansioni pre-esistenti o acquisisce nuove mansioni con l'introduzione dello strumento?
(Questo può influenzare il tempo a disposizione per l'apprendimento dello strumento)

Pianificare le seguenti attività

| | |
|--|--|
| Post-processing dell'algoritmo per attribuire pesi ai dati | |
| Verifica della legalità del post-processing | |

CHI SI ASSUME LA RESPONSABILITÀ DELLA DECISIONE FINALE PRESA SULLA BASE DELL'ALGORITMO NEI SEGUENTI CASI?

Identificare le figure responsabili in caso di violazioni nei seguenti ambiti

| | |
|----------------------|--|
| Non discriminazione | |
| Privacy | |
| Tutela del copyright | |
| Altro: _____ | |

4

Una guida per l'adozione etica di
strumenti AI nella PA

Elementi culturali e sociali

Questa sezione è per pianificare la comunicazione verso la cittadinanza sull'uso dell'AI e di dati personali da parte della PA

SONO PRESENTI DATI DA TERZE PARTI CHE POSSONO CREARE CONSEGUENZE ETICHE INDESIDERATE (ES. REIDENTIFICAZIONE DI SOGGETTI DI CUI SI SIANO RACCOLTI DATI PERSONALI)?

☐ Sì. Quali?_____

☐ No

QUALI PRECAUZIONI SONO PRESE PER LIMITARE LE POTENZIALI CONSEGUENZE INDESIDERATE?

SONO INCLUSI I SEGUENTI ASPETTI NELL'INFORMATIVA SUL TRATTAMENTO DEI DATI PERSONALI?

Quali categorie
di dati sono
raccolte?

Per quale
finalità sono
raccolti i dati?

Quali sono le
basi giuridiche
su cui si fonda
il trattamento
dei dati?

Con quali
metodi sono
raccolti i dati?

**I DATI RACCOLTI POSSONO ESSERE RILASCIATI
COME OPEN DATA IN FORMA AGGREGATA?**

NEL RACCOLGERE I DATI PERSONALI OCCORRE ELABORARE:

1. Un'informativa contenente tutte le info necessarie ai fini del GDPR
2. Una comunicazione dal linguaggio semplice e comprensibile a chiunque per favorire la piena comprensione del trattamento dati da parte dei soggetti interessati

5

Canvas: AI per la PA

Requisiti funzionali

Sulla base delle risposte alle domande precedenti, elencare i requisiti funzionali da implementare nella piattaforma AI

Implementazione dei principi tecnici

| | |
|--|--|
| Notifiche sulla completezza e sulla compatibilità tra ontologie dei dati | |
| Inclusione di strumenti di explainability | |
| Notifiche sulla soddisfazione degli obblighi giuridici | |
| Notifiche sulle modifiche ai dati apportate dall'AI | |
| Altro | |

Implementazione dei principi sociali e culturali

| | |
|---|--|
| Istruzioni sul funzionamento per utenti generalisti | |
| Notifiche sulla coerenza tra l'uso dello strumento e l'obiettivo per cui è stato creato | |
| Notifiche sui potenziali rischi privacy (es. re-identificabilità) | |
| Altro | |

Quali dei seguenti principi sono prioritari nell'implementazione dello strumento AI elaborato? Indicare un ordine di priorità

| | |
|--|--|
| Privacy | |
| Equa distribuzione di benefit economici | |
| Non discriminazione su base etnica e di genere | |
| Altro (es. sostenibilità ambientale) | |

Considerare i seguenti strumenti per la documentazione (e le istruzioni per l'utente): Data cards, Model cards, AI Product cards.

Utilizzare **AI Product cards** per specificare le caratteristiche tecniche del modello abbozzate in questa pagina del canvas.

Utilizzare la **Assessment List for Trustworthy Artificial Intelligence (ALTAI)** per una seconda valutazione della coerenza dell'algoritmo con le norme UE.