



UNIVERSITÀ DEGLI STUDI DI SALERNO
DIPARTIMENTO DI INFORMATICA

Laurea triennale in Informatica

Fondamenti di Intelligenza Artificiale

Lezione 14 - Ingegneria del Machine Learning



Ingegneria del Machine Learning

Ingegneria del Machine Learning?

Che vi piaccia o no, l'ingegneria del software è la chiave di tutto il software che c'è al mondo... senza questa, semplicemente, non potremmo sviluppare alcun sistema software affidabile o nessun sistema che consenta di fornire funzionalità agli utenti.

I sistemi software non adeguatamente ingegnerizzati falliscono nel 94% dei casi, rischiando di danneggiare cose, persone e ambiente circostante.

Recenti statistiche hanno dimostrato che la mancanza di metodi di ingegneria del software hanno causato (1) oltre tre trilioni di dollari in perdite finanziarie nel 2020; (2) oltre quattro miliardi di persone a fronteggiare problemi software; e (3) oltre 50 mila decessi all'anno.

Tutto questo perché l'ingegneria del software non è (solo) documentazione, ma anche e soprattutto gestione dei processi di sviluppo, misurazione delle attività condotte durante lo sviluppo, e monitoraggio della qualità del software sviluppato.

Non è un caso che il 97% dei neo-laureati in Italia inizia con un lavoro da software engineer. Non è un caso che molti di questi finiranno la propria carriera come project manager —> altre discipline sono importanti, ma è l'ingegneria del software ad abilitarle!

Il machine learning e l'intelligenza artificiale sono, di fatto, un particolari esempi di software. Cosa accade se questi sistemi non sono ingegnerizzati?

Ingegneria del Machine Learning?

Arizona. March 18, 2018 - 9:58 p.m

'Uber should be shut down': friends of self-driving car crash victim seek justice

Loved ones are in shock over the death of Elaine Herzberg in Arizona, but questions remain as to whether Uber will be held accountable

● **Self-driving Uber kills woman in first fatal pedestrian crash**

Elaine Herzberg è stata la prima vittima di un'auto a guida autonoma.

Sfortunatamente, non è stata l'ultima.

Tutto questo è successo per due ragioni, stando al report fornito dalla US National Transportation Safety Board:

- (1) Il modello di machine learning utilizzato dalla macchina a guida autonoma ha erroneamente classificato il pedone come un altro oggetto. L'aggravante è che la classificazione errata è durata per almeno 15 secondi, in cui l'auto ha ripetutamente ravvisato oggetti diversi sulla strada, senza mai accorgersi del pedone.
- (2) Anche se il sistema di guida autonoma avesse optato per una frenata di emergenza in una situazione del genere, questa manovra fu disabilitata "per ridurre il potenziale comportamento irregolare del veicolo".

In altri termini, il modello di machine learning non fu solo non testato abbastanza in termini di errore, ma fu anche lasciato libero di "apprendere" sul campo.

Ingegneria del Machine Learning?

Qualcuno di voi potrebbe pensare che questi siano solo casi rari o eccezionali... sfortunatamente, non è così.

Why Amazon's Automated Hiring Tool Discriminated Against Women



By [Rachel Goodman](#), Staff Attorney, ACLU Racial Justice Program
OCTOBER 12, 2018 | 1:00 PM

Nel 2018, Amazon è stata al centro di polemiche enormi per via dei suoi modelli automatici di reclutamento.

Il modello di machine learning utilizzato da Amazon dal 2014 al 2018 utilizzava i curriculum degli attuali dipendenti come base per classificare se, un nuovo curriculum, sarebbe dovuto essere considerato o meno.

Il problema nacque dal fatto che il modello considerava dati sensibili, come il genere di un dipendente. Al tempo, molti dei dipendenti Amazon erano di genere maschile e, per questa ragione, i curriculum provenienti da persone di altro genere venivano sistematicamente scartati.

Nel 2018, Amazon è stata costretta a ritirare il modello di machine learning.

Problemi di questo tipo hanno riguardato negli anni non solo il genere, ma anche altri attributi personali, come l'invalidità civile, l'età, la razza, e altri attributi che non dovrebbero influenzare in nessun modo le scelte operative di un'azienda.

Ingegneria del Machine Learning?

Ma perché tutti questi problemi e perché tutto così frequentemente?

L'intelligenza artificiale impara ma non capisce

Sebbene si pensi che le macchine stiano diventando intelligenti, c'è sempre un fattore poco considerato: gli algoritmi di machine learning vengono progettati ed implementati da *uomini*. In altri termini, se un machine learner è progettato male, allora imparerà a fare cose stupide!

Come abbiamo già iniziato a comprendere, un machine learner “non è niente di più” che un modello matematico capace di fare inferenza sulla base dei dati che ha a disposizione. Chiaramente, questi dati rappresenteranno il modello di apprendimento.

Could 'fake text' be the next global political threat?

An AI fake text generator that can write paragraphs in a style based on just a sentence has raised concerns about its potential to spread false information

Alcuni politici italiani - ogni riferimento a cose o persone è puramente (non) casuale - fanno costante uso di questi tool, influenzando volontariamente l'opinione pubblica!

Ingegneria del Machine Learning

Ingegneria del Machine Learning!

Ma fortunatamente, voi avete opzionato un esame a scelta come FIA. E, fortunatamente (spero), non avete e non utilizzerete mai un fake news generator e non sarete mai dei politici disonesti - in caso contrario, vi pregherei di lasciare il corso! :)

A valle di quanto detto finora, sono quindi due le cose principali a cui pensare quando progettare una soluzione basata su machine learning: data & software engineering!

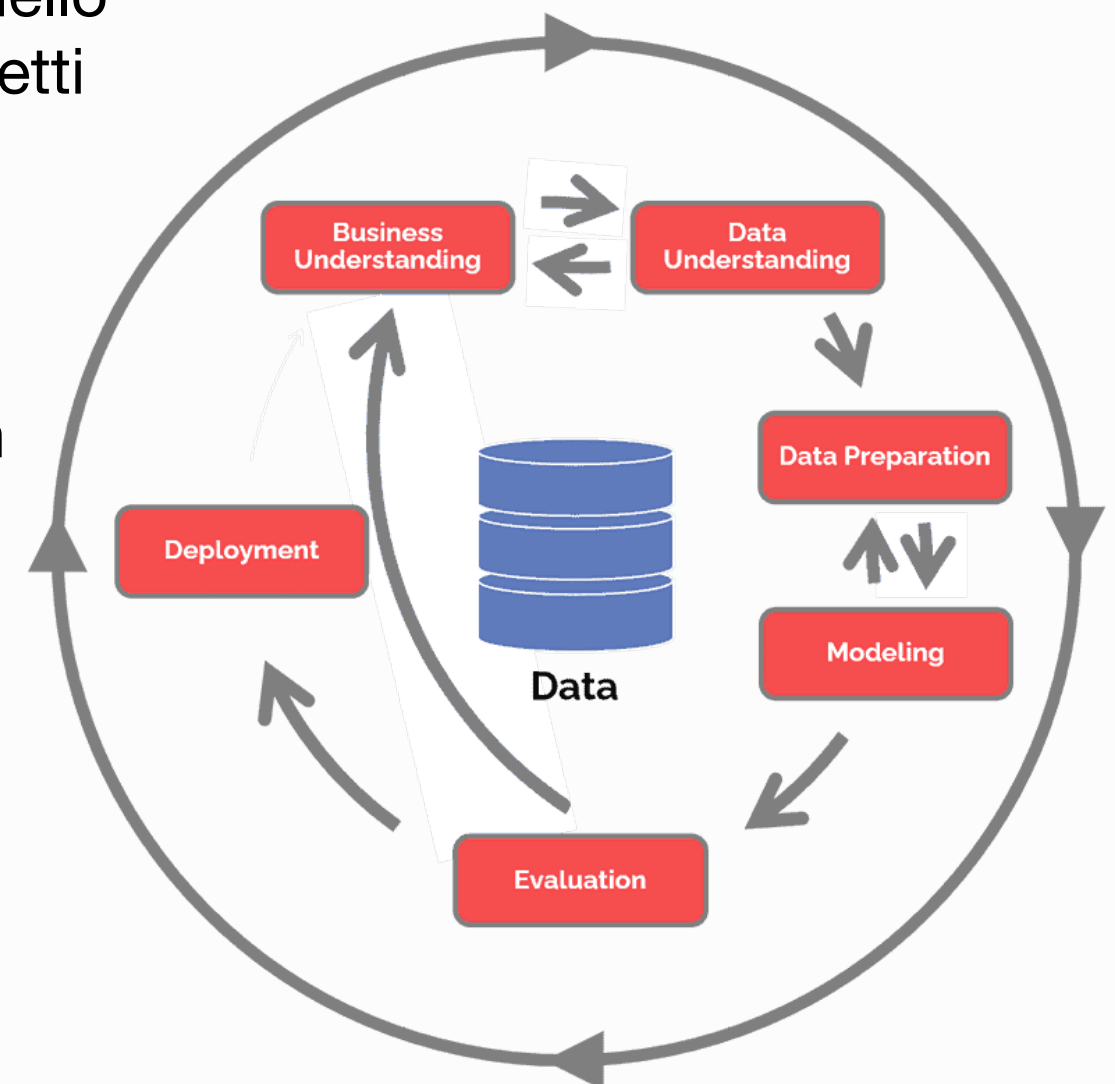
Tutto ciò può essere riassunto dal cosiddetto modello CRISP-DM, che rappresenta il ciclo di vita di progetti basati su intelligenza artificiale e data science.

CRISP-DM è l'acronimo di Cross-Industry Standard Process for Data Mining.

Possiamo paragonare il modello CRISP-DM ad un modello a cascata con feedback utilizzato per lo sviluppo di sistemi software tradizionali.

Il CRISP-DM è un modello *non* sequenziale in cui le diverse fasi possono essere eseguite un numero illimitato di volte.

Andiamo a vedere più nel dettaglio ciascuna di queste fasi.



Il Modello CRISP-DM

Business Understanding

La prima fase è chiaramente quella di raccolta dei requisiti e di *definizione degli obiettivi di business* che si intende raggiungere (ovvero, che cosa deve fare il machine learner che stiamo progettando).

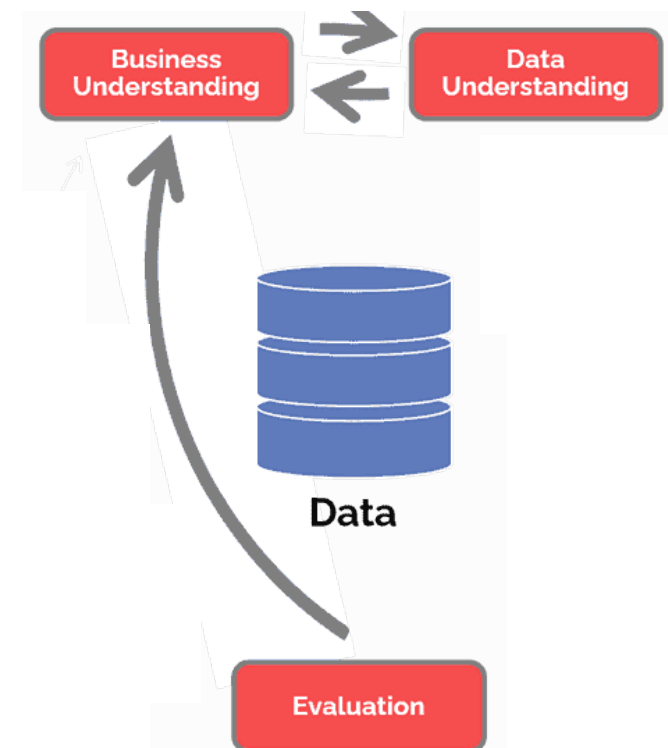
La fase di business understanding prevede la definizione dei cosiddetti *business success criteria*, ovvero i criteri secondo i quali potremo accertare che il sistema costruito è in linea con gli obiettivi di business.

In questa fase, bisogna inoltre determinare la *disponibilità delle risorse*, *stimare i rischi*, *definire i relativi piani di contingenza* e condurre una *analisi costi-benefici*.

Oltre che definire i criteri di successo da un punto di vista di business, è inoltre necessario definire gli *obiettivi tecnici* che si intendono raggiungere.

Infine, verranno selezionate le *tecnologie ed i tool* necessari agli obiettivi.

Output: *Piano di progetto*, il documento che spiega l'esecuzione del progetto da un punto di vista di gestione tecnica e socio-tecnica.



Il Modello CRISP-DM

Data Understanding

Sulla base degli obiettivi definiti nella fase precedente, il secondo passo consiste *nell'identificazione, collezione e analisi dei dataset* che possono portare al raggiungimento degli obiettivi.

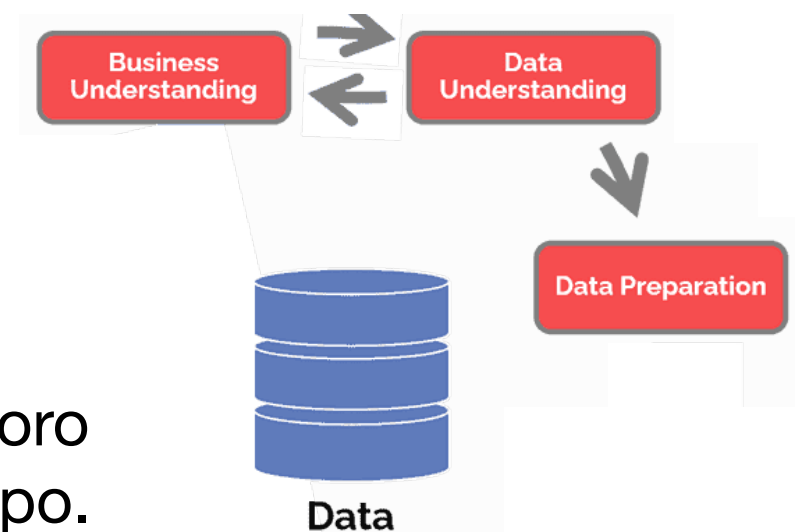
Innanzitutto, quindi, vengono *acquisiti i dati necessari* al raggiungimento degli obiettivi di business e tecnici. I dati verranno poi caricati in un tool di analisi dei dati.

I dati vengono quindi *esaminati e documentati* rispetto al loro formato, il numero di record e il significato di ciascun campo.

Il terzo task consiste nell'*esplorazione dei dati*: questi vengono visualizzati e, cosa più importante, eventuali relazioni vengono identificate.

Infine, il processo di *qualità dei dati*. Vengono identificati e documentati possibili problemi di qualità dei dati (ad esempio, dati mancanti).

Output: Documento di analisi dei dati, che riporta i metodi di estrazione ed analisi, oltre che le relazioni tra i dati ed eventuali problemi di qualità.



Il Modello CRISP-DM

Data Preparation

L'obiettivo di questa fase è quello di *preparare i dati* in maniera tale che possano essere utilizzati nei successivi passi del processo.

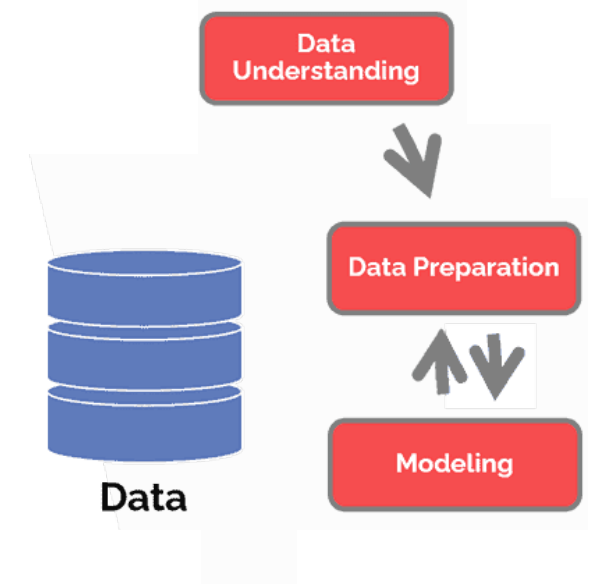
Questo include un processo fondamentale che è noto come *feature engineering*, ovvero la selezione delle caratteristiche del problema che hanno maggiore *potenza predittiva*.

Inoltre, questa fase include l'implementazione dei processi di *pulizia dei dati* sulla base dei problemi di qualità riscontrati nella fase precedente.

Sulla base della pulizia fatta così come dell'analisi della potenza predittiva delle caratteristiche considerate, il progettista può considerare di estendere le caratteristiche da considerare.

Infine, i dati vengono *formattati* in una maniera tale che possano essere prese in input da un modello di machine learning - questo potrebbe dipendere dai tool selezionati in fase di business understanding.

Output: Insieme di dati di input, ovvero l'insieme di dati che verranno considerati in fase di modellazione dell'algoritmo di machine learning.



Il Modello CRISP-DM

Data Modeling

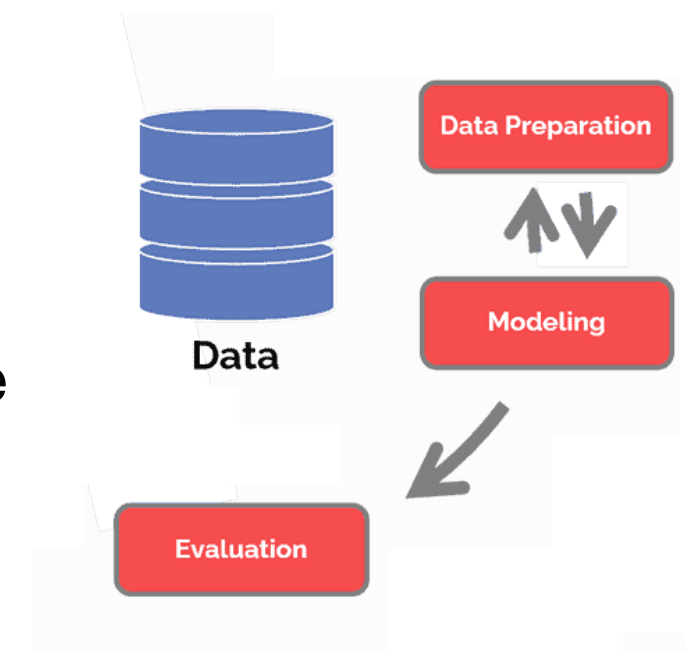
Una volta sistemati i dati, è ora di iniziare la *fase di modellazione*, che è sempre particolarmente complicata e “unica”, nel senso che la definizione di un algoritmo di machine learning dipende strettamente dal problema in esame e dai dati a disposizione.

In primo luogo, va *selezionata la tecnica o l'algoritmo da utilizzare*: ad esempio, conviene modellare il problema come un problema di classificazione o regressione? Quale soluzione sarà più adatta e utile al raggiungimento degli obiettivi?

Dopodiché, si passa alla *fase di addestramento*. In questo caso, si configurano i parametri del modello selezionato, si addestra il modello e si descrivono i risultati ottenuti in fase di addestramento.

Molto spesso, il progettista sarà costretto a tornare nella fase di data preparation per effettuare ulteriori operazioni sui dati.

Output: *Il modello di machine learning*, ovvero l'algoritmo addestrato e configurato sui dati a disposizione.



Il Modello CRISP-DM

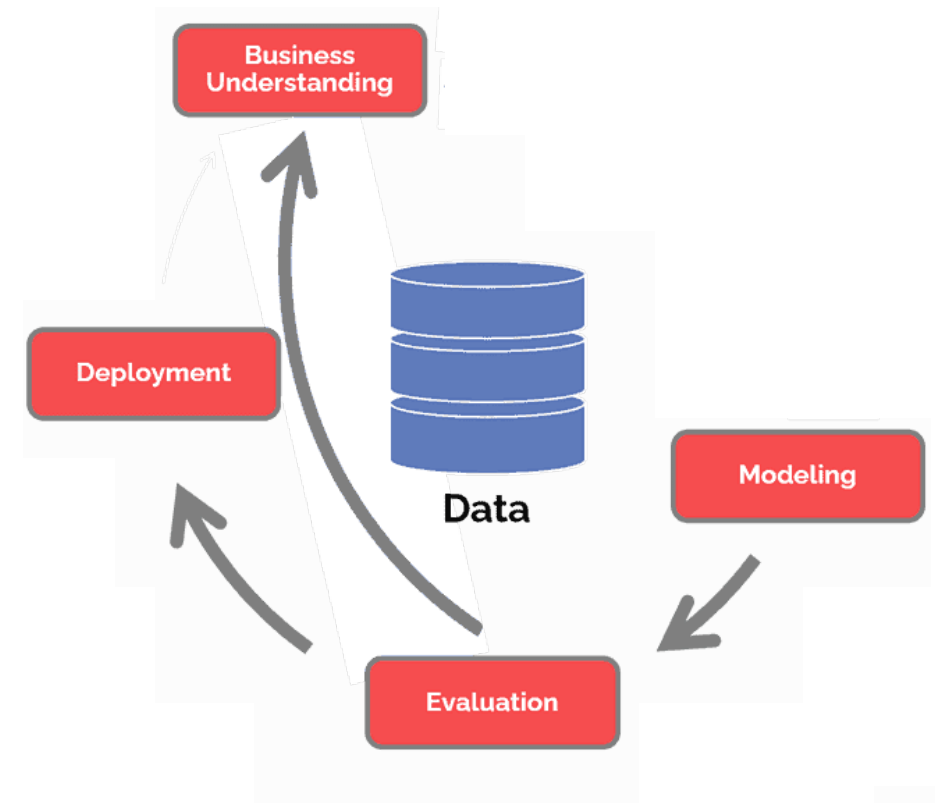
Evaluation

La *fase di validazione* ha l'obiettivo di valutare se i risultati sono chiari, se sono in linea con gli obiettivi di business e se rivelano delle prospettive aggiuntive alle quali il progettista non aveva pensato.

In questa fase, è inoltre importante verificare la consistenza e la solidità dell'intero processo. Ad esempio, ci sono degli aspetti che non sono convincenti? Ci sono delle alternative metodologiche che potrebbero portare a risultati diversi? E come queste alternative impattano i risultati ottenuti?

Una volta avute le risposte, si può procedere alla definizione dei prossimi passi da effettuare. Possiamo considerare la definizione dell'approccio completa? Oppure è necessario fare un passo indietro e valutare opzioni diverse?

Output: *Risultati della validazione del modello*, che rivelano il grado di attendibilità e conformità dell'algoritmo rispetto agli obiettivi di business.



Il Modello CRISP-DM

Deployment

La *fase di deployment* ha l'obiettivo di mettere in funzione l'approccio definito e, quindi, renderlo usabile.

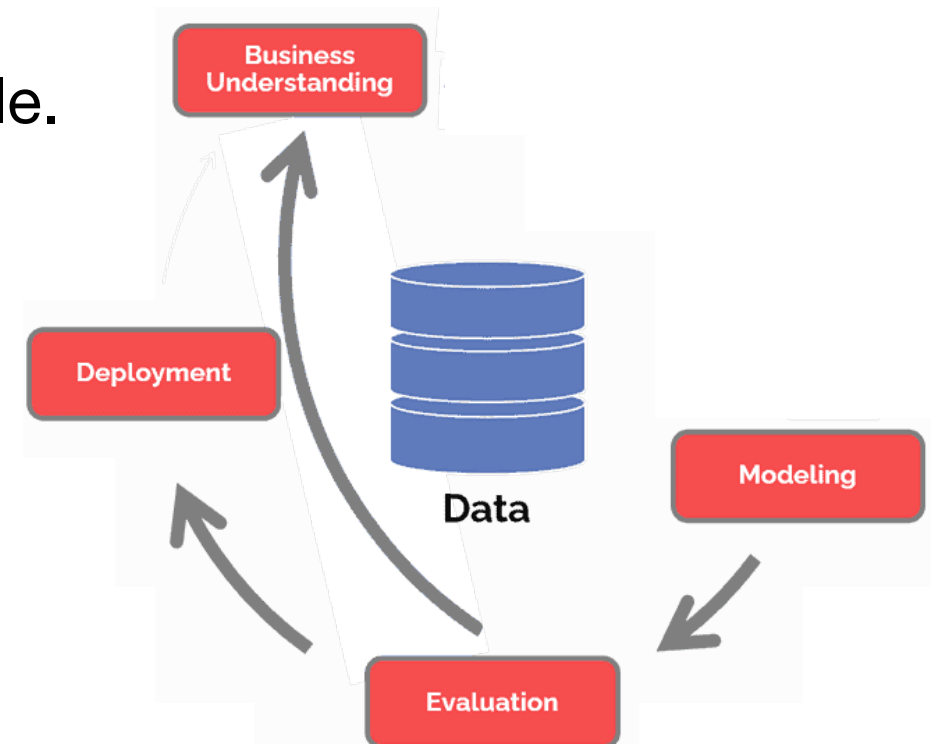
Il modello generato, precisamente, come dovrà essere reso disponibile? Quale sarà il grado di interazione con gli utenti? Ed in che modo gli utenti utilizzeranno il modello?

In altri termini, questa fase vede il passaggio dall'ingegneria del machine learning all'ingegneria del software e all'ingegneria dell'usabilità!

Questo è ancora più evidente se consideriamo che questo modello non resterà nel suo stato in eterno —> avrà bisogno di essere costantemente monitorato e mantenuto!

Occhio alle prime due leggi di Lehman sull'evoluzione di sistemi software: (1) un sistema che non cambia è un sistema che diventerà presto inutile; (2) la complessità del sistema crescerà inesorabilmente nel tempo.

Output: *Report finale di progetto*, che descrive tutte le fasi condotte, oltre che il piano di manutenzione e monitoraggio.



Il Modello CRISP-DM: Un modello a cascata o un modello agile?

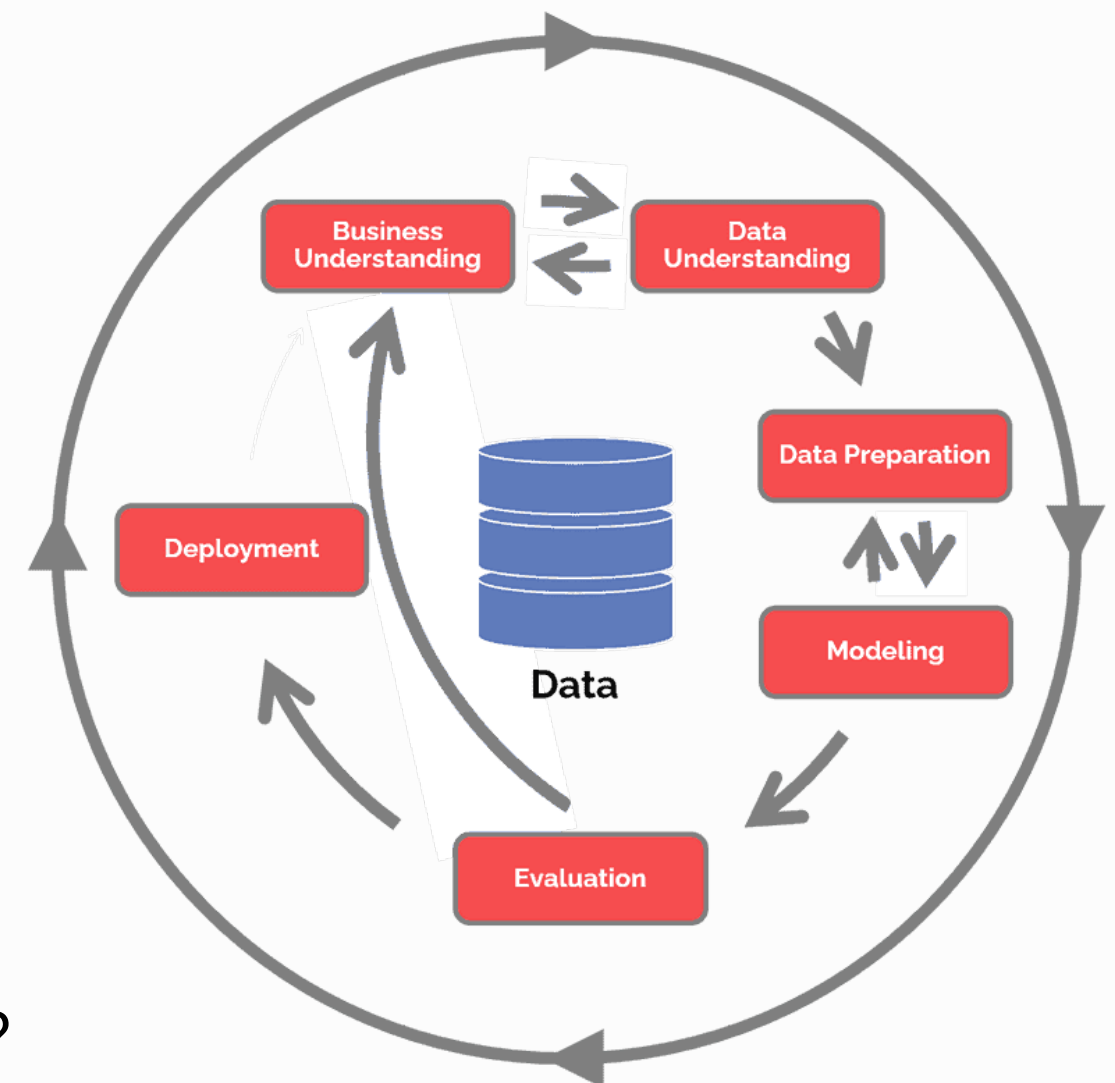
Come detto in precedenza, il modello CRISP-DM è simile ad un ciclo di vita a cascata. Di fatti, seguendo pedissequamente le regole definite, è facile simulare un modello di ingegneria del machine learning del genere.

Sebbene questo sia vero, è anche bene notare che il modello CRISP-DM non è, per definizione, rigido. Possiamo facilmente notare che sono molte le fasi in cui un progettista potrebbe trovarsi a tornare sui suoi passi e/o effettuare azioni più e più volte, finendo quindi in una situazione agile.

Quindi, il modello è di fatto flessibile e può essere considerato sia tradizionale che agile, dipendentemente dal livello di flessibilità che si intende mettere in atto.

In tutti i casi, c'è qualcosa che CRISP-DM non considera esplicitamente. I progetti software sono, per natura, **socio-tecnici**. Il codice ed i modelli di machine learning sono solo parte della storia: lo sviluppo è fatto dagli uomini!

Il modello CRISP-DM non chiarisce chi è responsabile di fare cosa. C'è quindi un modo migliore di sviluppare modelli di machine learning?



CRISP-DM, ti presento SCRUM...

Dovreste già essere “pratici” di SCRUM, ma in poche parole possiamo dire che SCRUM è un modello di ciclo di vita che prevede la divisione dei blocchi di lavoro in *sprint*, ovvero dei brevi periodi di tempo in cui determinati requisiti vengono definiti, analizzati, progettati e sviluppati.

A differenza dei cicli di vita tradizionali, SCRUM prevede lo sviluppo incrementale di porzioni del sistema. L'incrementalità serve da un lato ad avere una migliore interazione con gli stakeholder (dovuta al rilascio di piccole parti del sistema) e dall'altro a minimizzare i rischi (che nel processo tradizionale sono dovuti al possibile misunderstanding dei requisiti).

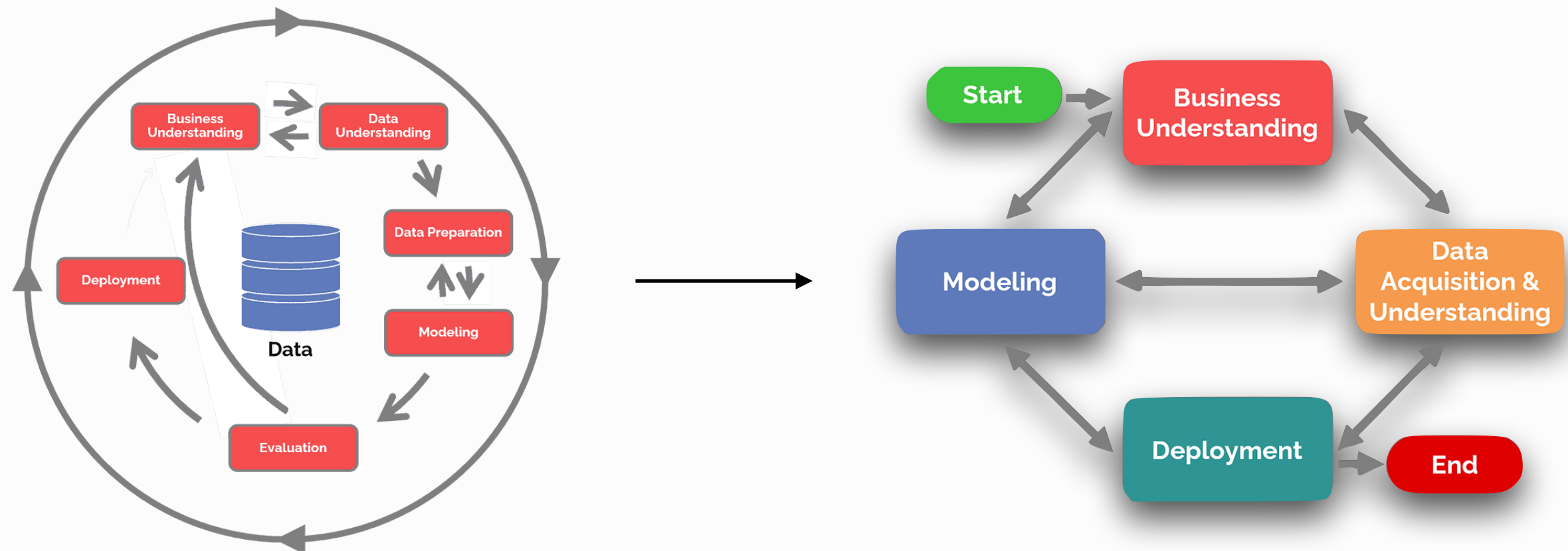
SCRUM prevede, tra gli altri, tre ruoli fondamentali:

- (1) Il *product owner*, colui il quale rappresenta gli stakeholder, definisce e prioritizza i requisiti, per aggiungerli al backlog;
- (2) Lo *SCRUM Master*, colui che agevola il lavoro e coordina le attività di sviluppo. Sebbene sia un ruolo manageriale, non è da confondere con un project manager, poiché non ha alcuna responsabilità di gestione del personale.
- (3) Il *team di sviluppo*, che è responsabile della consegna del prodotto. Il team dovrebbe essere composto da 3-9 persone che competenze trasversali, le quali si occupano della documentazione e dello sviluppo del prodotto in ogni sprint.

CRISP-DM, ti presento SCRUM...

In SCRUM, ogni sprint inizia con uno *sprint planning*, in cui le attività da compiere nel prossimo periodo temporale vengono definite. Il lavoro è monitorato tramite i *daily scrum*, una riunione quotidiana in cui i progressi vengono riportati. Alla fine dello sprint, viene effettuata una *sprint review*, dove il team mostra i progressi effettuati. Infine, lo sprint viene chiusa con una *sprint retrospective*, dove il team riflette sui problemi riscontrati e definisce i miglioramenti da fare nel prossimo sprint.

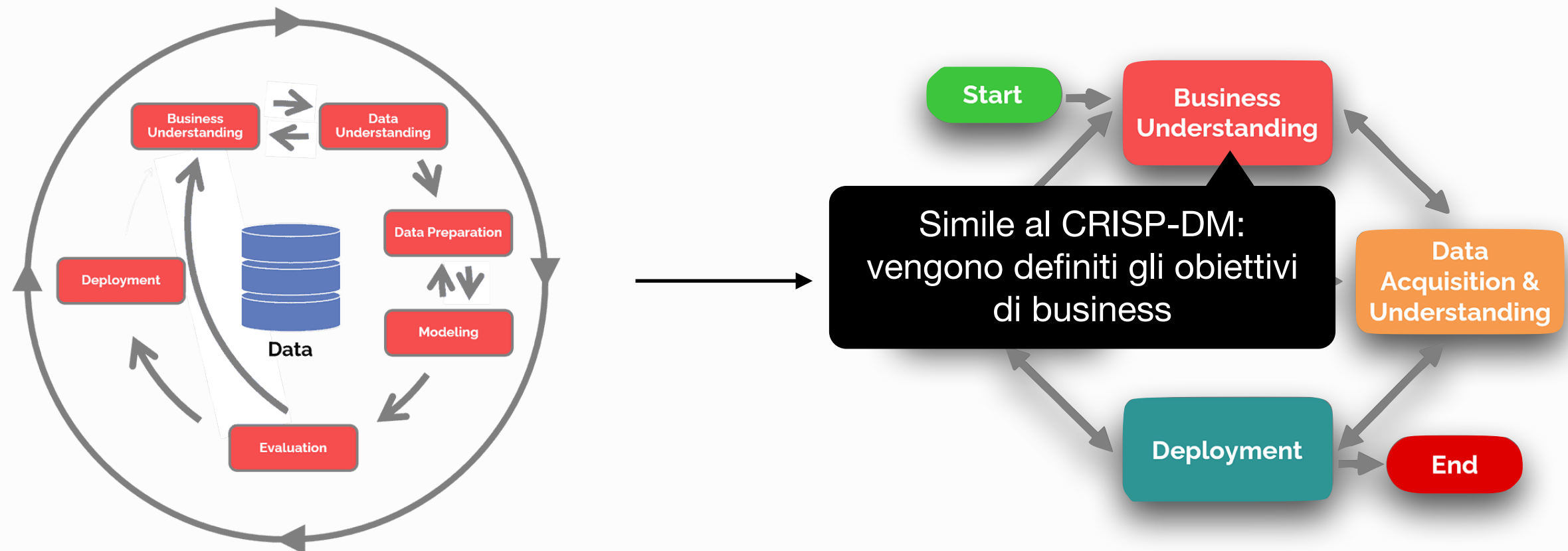
Combinando SCRUM e CRISP-DM, si ottiene il **TDSP**, acronimo di Team Data Science Process. Sviluppato da Microsoft nel 2020 (e ancora in via di affinamento), TDSP incorpora aspetti socio-tecnici nell'esecuzione dei progetti.



CRISP-DM, ti presento SCRUM...

In SCRUM, ogni sprint inizia con uno *sprint planning*, in cui le attività da compiere nel prossimo periodo temporale vengono definite. Il lavoro è monitorato tramite i *daily scrum*, una riunione quotidiana in cui i progressi vengono riportati. Alla fine dello sprint, viene effettuata una *sprint review*, dove il team mostra i progressi effettuati. Infine, lo sprint viene chiusa con una *sprint retrospective*, dove il team riflette sui problemi riscontrati e definisce i miglioramenti da fare nel prossimo sprint.

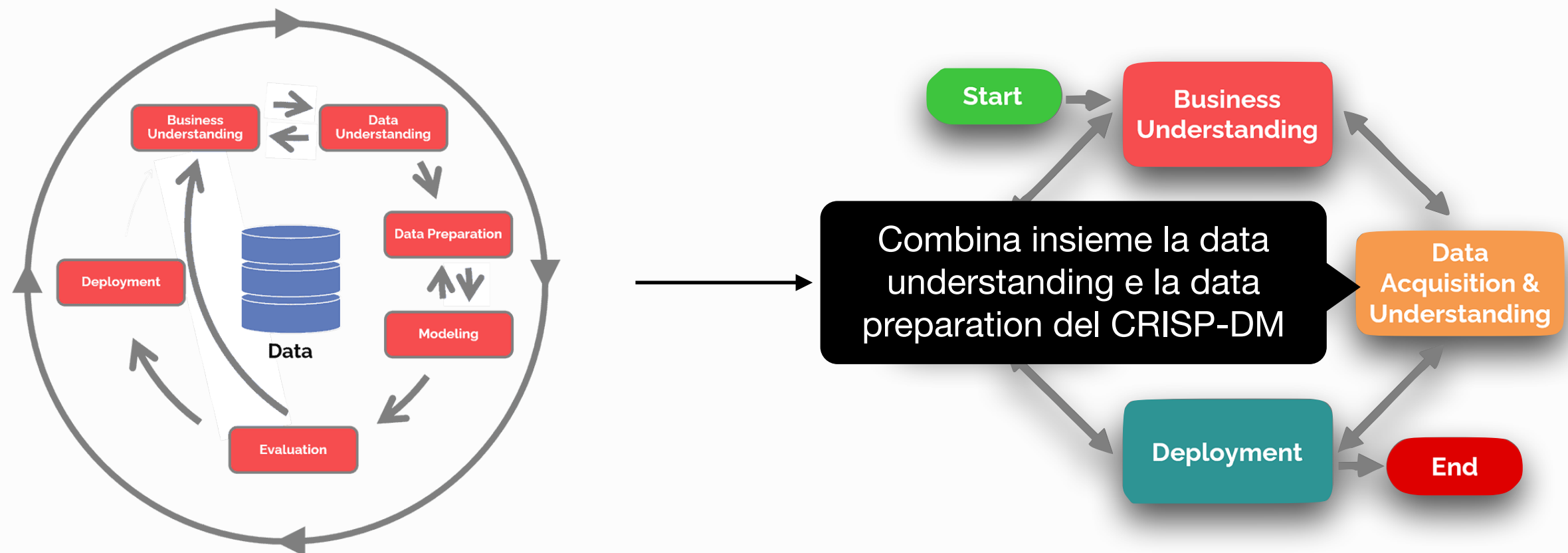
Combinando SCRUM e CRISP-DM, si ottiene il **TDSP**, acronimo di Team Data Science Process. Sviluppato da Microsoft nel 2020 (e ancora in via di affinamento), TDSP incorpora aspetti socio-tecnici nell'esecuzione dei progetti.



CRISP-DM, ti presento SCRUM...

In SCRUM, ogni sprint inizia con uno *sprint planning*, in cui le attività da compiere nel prossimo periodo temporale vengono definite. Il lavoro è monitorato tramite i *daily scrum*, una riunione quotidiana in cui i progressi vengono riportati. Alla fine dello sprint, viene effettuata una *sprint review*, dove il team mostra i progressi effettuati. Infine, lo sprint viene chiusa con una *sprint retrospective*, dove il team riflette sui problemi riscontrati e definisce i miglioramenti da fare nel prossimo sprint.

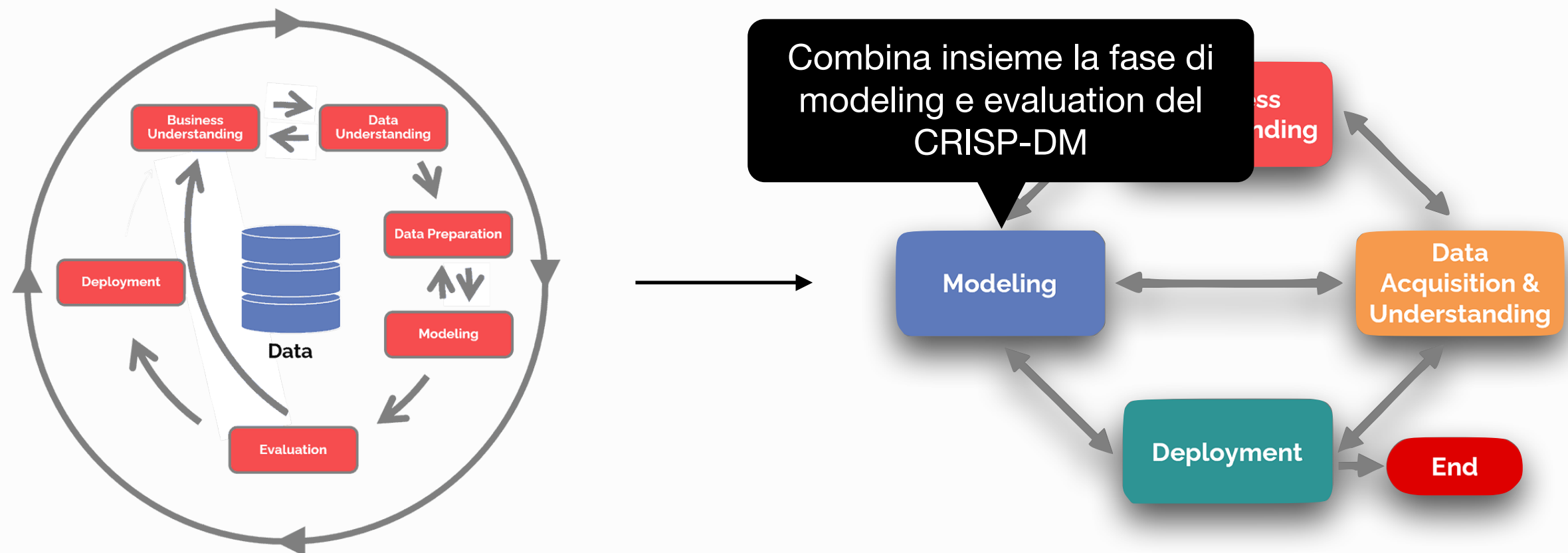
Combinando SCRUM e CRISP-DM, si ottiene il **TDSP**, acronimo di Team Data Science Process. Sviluppato da Microsoft nel 2020 (e ancora in via di affinamento), TDSP incorpora aspetti socio-tecnici nell'esecuzione dei progetti.



CRISP-DM, ti presento SCRUM...

In SCRUM, ogni sprint inizia con uno *sprint planning*, in cui le attività da compiere nel prossimo periodo temporale vengono definite. Il lavoro è monitorato tramite i *daily scrum*, una riunione quotidiana in cui i progressi vengono riportati. Alla fine dello sprint, viene effettuata una *sprint review*, dove il team mostra i progressi effettuati. Infine, lo sprint viene chiusa con una *sprint retrospective*, dove il team riflette sui problemi riscontrati e definisce i miglioramenti da fare nel prossimo sprint.

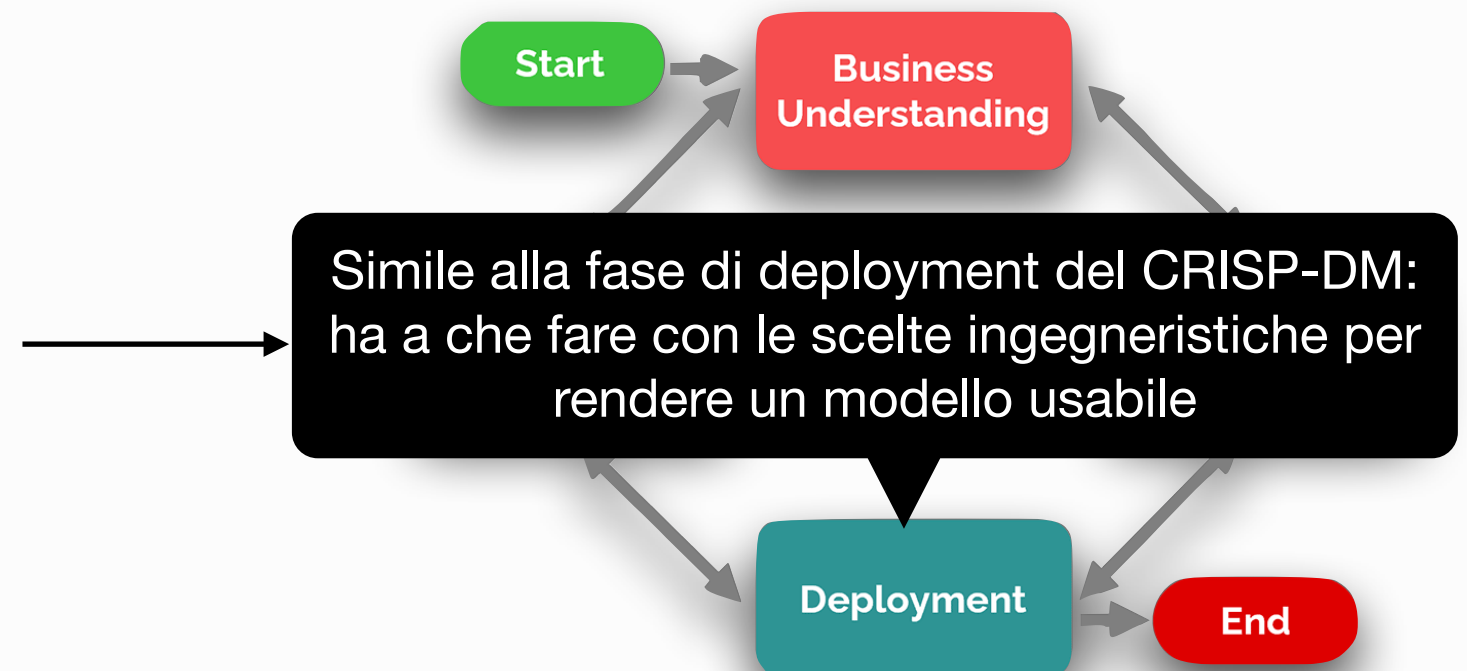
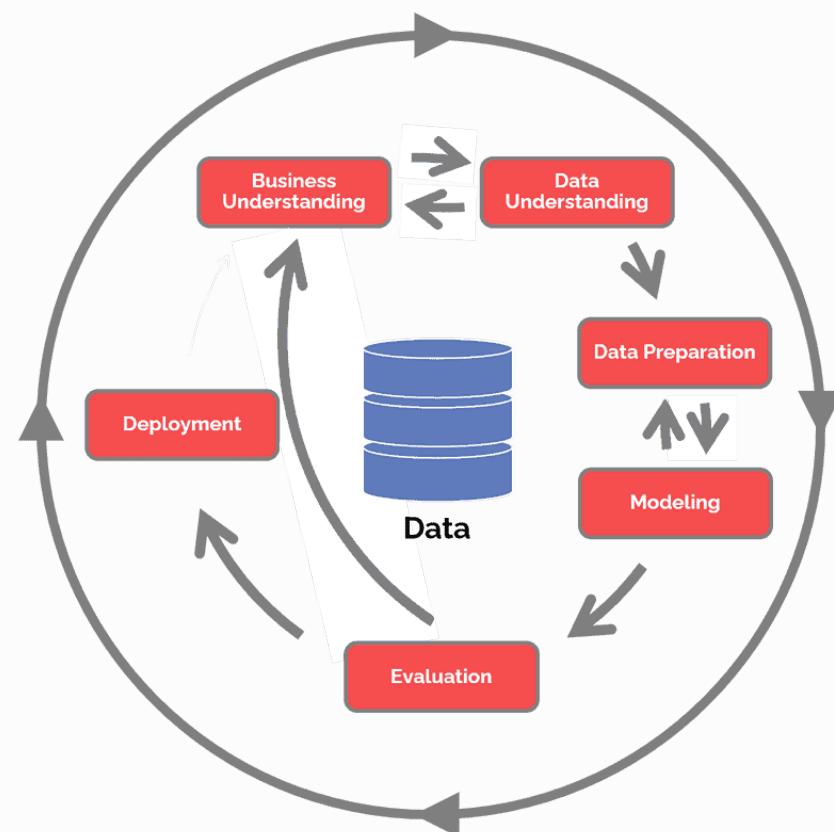
Combinando SCRUM e CRISP-DM, si ottiene il **TDSP**, acronimo di Team Data Science Process. Sviluppato da Microsoft nel 2020 (e ancora in via di affinamento), TDSP incorpora aspetti socio-tecnici nell'esecuzione dei progetti.



CRISP-DM, ti presento SCRUM...

In SCRUM, ogni sprint inizia con uno *sprint planning*, in cui le attività da compiere nel prossimo periodo temporale vengono definite. Il lavoro è monitorato tramite i *daily scrum*, una riunione quotidiana in cui i progressi vengono riportati. Alla fine dello sprint, viene effettuata una *sprint review*, dove il team mostra i progressi effettuati. Infine, lo sprint viene chiusa con una *sprint retrospective*, dove il team riflette sui problemi riscontrati e definisce i miglioramenti da fare nel prossimo sprint.

Combinando SCRUM e CRISP-DM, si ottiene il **TDSP**, acronimo di Team Data Science Process. Sviluppato da Microsoft nel 2020 (e ancora in via di affinamento), TDSP incorpora aspetti socio-tecnici nell'esecuzione dei progetti.



CRISP-DM, ti presento SCRUM...

In SCRUM, ogni sprint inizia con uno *sprint planning*, in cui le attività da compiere nel prossimo periodo temporale vengono definite. Il lavoro è monitorato tramite i *daily scrum*, una riunione quotidiana in cui i progressi vengono riportati. Alla fine dello sprint, viene effettuata una *sprint review*, dove il team mostra i progressi effettuati. Infine, lo sprint viene chiusa con una *sprint retrospective*, dove il team riflette sui problemi riscontrati e definisce i miglioramenti da fare nel prossimo sprint.

Combinando SCRUM e CRISP-DM, si ottiene il **TDSP**, acronimo di Team Data Science Process. Sviluppato da Microsoft nel 2020 (e ancora in via di affinamento), TDSP incorpora aspetti socio-tecnici nell'esecuzione dei progetti.

La vera, grande differenza sta nel fatto che, ad intervalli regolari, viene effettuata una fase di *customer acceptance*, ovvero una validazione di come il sistema implementa i requisiti di business.

Inoltre, TDSP definisce sei ruoli espliciti:

Project manager: Il responsabile dell'intero progetto;

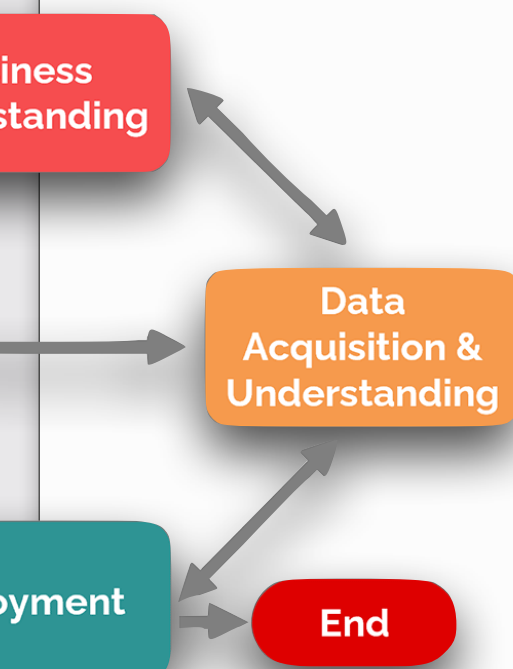
Project lead: Simile allo SCRUM master, è il team leader.

Data engineer: Il responsabile della data acquisition;

Data scientist: Il responsabile della data understanding;

Application developer: Il responsabile dell'implementazione dell'applicazione;

Solution architect: Progetta e manutene le architetture delle applicazioni;



TDSP, pro e contro

By design, TDSP enfatizza la necessità di rilasci incrementali, il che minimizza i rischi connessi a potenziali misunderstanding dei requisiti.

Più importante, TDSP riconosce esplicitamente la *complementarietà* tra ingegneria del machine learning e ingegneria del software. I ruoli previsti da TDSP includono esperti di project management e lead, esperti di ingegneria del software, oltre che esperti di data engineering e science.

Altro aspetto da non sottovalutare: il TDSP utilizza una terminologia e dei tool simili a quelli di SCRUM, il che semplifica la comprensione delle responsabilità.

Di contro, la necessità di definire degli sprint potrebbe essere controproducente: a differenza dei progetti software tradizionali, la definizione di sistemi di machine learning potrebbe incappare in problemi tecnici non di poco conto, come ad esempio la mancanza o la bassa qualità dei dati!

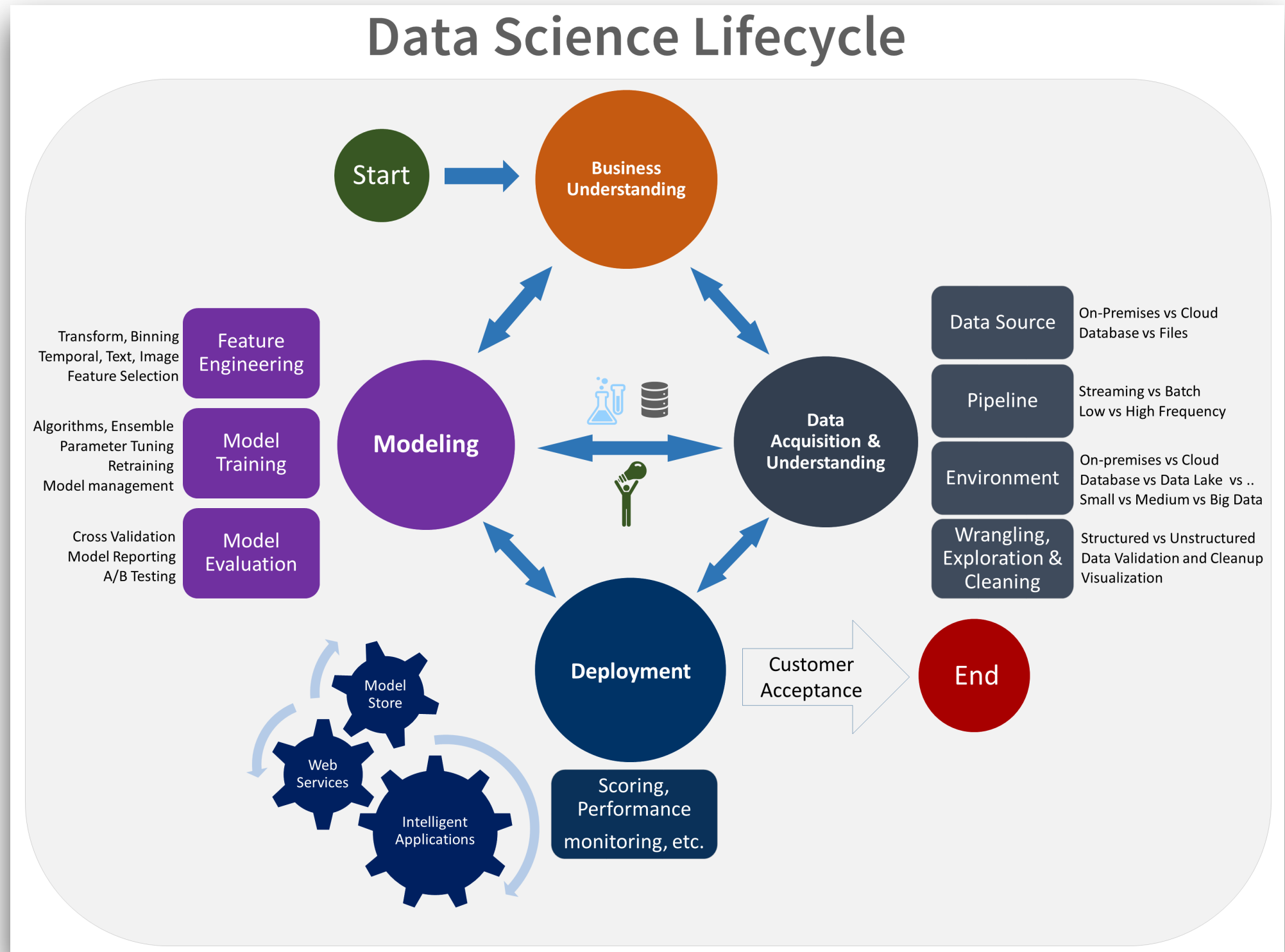
TDSP, pro e contro

a i rischi

gneria
lono
e che

I simili a

nte: a
ne
esempio





UNIVERSITÀ DEGLI STUDI DI SALERNO
DIPARTIMENTO DI INFORMATICA

Laurea triennale in Informatica

Fondamenti di Intelligenza Artificiale

Lezione 14 - Ingegneria del Machine Learning

