

Խնդրի դրվածք

Վարկային կազմակերպությունները հանդիպում են հիմնական մարտահրավերի. գնահատել հաճախորդների վճարունակությունն ու վարկային պորտֆելի ընդհանուր ռիսկայնությունը՝ համապատասխանելով ՖՀՄՄ 9 ստանդարտի պահանջներին: ՖՀՄՄ 9-ը պարտադրում է գնահատել ակնկալվող վարկային վնասները (Expected Credit Loss, ECL)՝ հիմնված ոչ միայն պատմական տվյալների, այլև ընթացիկ տնտեսական ցուցանիշների և կանխատեսվող մակրոտնտեսական սցենարների վրա: Այս գործընթացը պահանջում է մեծածավալ ֆինանսական տվյալների մշակում, բազմազործոն վերլուծություն և մոդելների կիրառություն, որոնք ապահովում են կանխատեսումների ճշգրտություն և վստահելիություն:

Խնդրն այն է, որ ավանդական տեղային (on-premises) համակարգերը դժվարանում են ապահովել անհրաժեշտ հաշվարկային հզորությունը, տվյալների անվտանգ պահպանումը, մոդելների ձկուն թարմացումը, ծավալային մասշտաբայնացնելը, մասնավորապես, եթե գործ ունենք հարյուր հազարավոր հաճախորդների, վարկային պատմությունների, ֆինանսական հաշվետվությունների և մակրոտնտեսական ժամանակաշարերի հետ:

Այս ավարտական աշխատանքի նպատակն է մշակել վարկային ռիսկի մոդելավորման և ակնկալվող վարկային վնասների գնահատման ամպային համակարգ օգտագործելով Microsoft Azure հարթակը, որը հնարավորություն է տալիս մեծ տվյալների արդյունավետ մշակում:

Ֆինանսական տվյալների հավաքում և նախապատրաստում

- Վարկառուների պատմական տվյալներ
- Վարկերի մարում/չմարում
- Վարկային ոեյթինգային կարգավիճակներ (Stage 1/2/3՝ ըստ ՖՀՄՄ 9)
- Մակրոտնտեսական ցուցանիշներ
- Տվյալների մաքրում, «մեղալիոն» ձարտարապետությամբ

Գնահատել IFRS 9 ECL-ի համար պահանջվող երեք հիմնական բաղադրիչները.

- PD – Probability of Default (վերաբեռնման հավանականություն)
 - LGD – Loss Given Default (վնասի տեսակարար չափ)
 - EAD – Exposure at Default (մնացորդ պարտավորություն)
- Սրանց հիման վրա՝ $ECL = PD \times LGD \times EAD$

Տվյալների խողովակաշարերի կառուցում և պահպանում Azure-ում

- Azure Data Factory՝ տվյալների կորզում, մշակում և արտահանում
- Azure SQL / Data Lake՝ արդյունքների պահպանում

Մասշտարավորվող ավարտական համակարգի նախագծում

- Ապահովել ապահովություն, մատչելիություն, տվյալների գաղտնիություն

Այսպիսով աշխատանքի հիմնական նպատակն է մշակել ամպային հարթակով աշխատող, ավտոմատացված և մասշտաբավորվող վարկային ռիսկի գնահատման համակարգ, որը.

- Իրականացնում է վարկային ռիսկի մոդելավորում՝ հիմնված պատմական և կանխատեսվող տվյալների վրա
- Վերահաշվարկում է Expected Credit Loss ըստ IFRS 9-ի
- Թույլ է տալիս ֆինանսական կազմակերպություններին ընդունել տվյալահեն որոշումներ
- Նվազեցնում է ձեռքով կատարվող հաշվարկների սխալները և ռիսկերը

Մատակարարում է վարկառուների ռիսկային պրոֆիլի օպերատիվ գնահատում՝ Azure-ի հաշվարկային ռեսուրսներով

Microsoft Azure-ի և ամպային տեխնոլոգիաներին վերաբերվող հիմնական հասկացությունները

Ներածություն

Ամպային տեխնոլոգիաները հանդիսանում են ժամանակակից ծառայությունների մատուցման մոդել, որի հաշվարկային և ծրագրային ռեսուրսները հասանելի են համացանցի միջոցով՝ ըստ պահանջի: Այս մոտեցումը ազատում է կազմակերպություններին սեփական սերվերների, տվյալների պահպանման համակարգերի և ցանցային ենթակառուցվածքի ստեղծման ու սպասարկման դժվարություններից: Ամպային ծառայությունները ներառում են վիրտուալ սերվերներ և հաշվարկային հզորություն, տվյալների պահեստավորում և արխիվացում, տվյալների բազաներ ու վերլուծական ծառայություններ, ցանցային գործիքակազմ՝ firewall, VPN, ինչպես նաև ծրագրային ապահովում՝ DevOps միջավայրեր: Ամպային հաշվարկի հիմնական իմաստը կայանում է նրանում, որ բոլոր ռեսուրսները հասանելի են ցանկացած տեղ՝ համացանցի առկայության դեպքում, և վճարումն իրականացվում է օգտագործված ռեսուրսների չափով:

Ամպային մոդելի կարևորագույն առավելություններից է մասշտաբելիությունը, որը թույլ է տալիս արագ ավելացնել կամ նվազեցնել օգտագործվող ռեսուրսների քանակը ըստ անհրաժեշտության, ինչը հատկապես կարևոր է սեղոնային կամ անկանոն բեռնվածություն ունեցող բիզնեսների համար: Մյուս կարևոր առավելությունն է ցածր նախնական ծախսերը. ընկերությունները չեն գնում թանկարժեք սարքավորումներ, չեն կատարում սպասարկման վճարներ և նվազեցնում են IT անձնակազմի ծանրաբեռնվածությունը, քանի որ համակարգը գործում է «վճարիր օգտագործելուց հետո» սկզբունքով: Բացի այդ, ամպային համակարգերը ապահովում են բարձր հուսալիություն, քանի որ տվյալների կենտրոնները, որոնք բաշխված են աշխարհով մեկ, ապահովում են ծառայությունների անընդհատ աշխատանքը, և եթե մի կենտրոնում խնդիր առաջանա, բեռնվածությունն ավտոմատ կերպով տեղափոխվում է մյուս կենտրոնների վրա:

Անվտանգության առումով ամպային հարթակները ներառում են ներկառուցված մեխանիզմներ՝ մշտական մոնիթորինգ, տվյալների կոդավորում, վերահսկում և պաշտպանություն DDoS հարձակումներից, որոնք համապատասխանում են միջազգային ստանդարտներին, ինչպիսիք են ISO, GDPR և HIPAA: Ամպային ծառայությունները նաև ապահովում են արագ զարգացում և ավտոմատացում՝ տրամադրելով CI/CD գործընթացներ, ավտոմատացման գործիքներ և թեստավորուման համակարգեր, ինչը արագացնում է ծրագրային ապահովման թողարկումը:

Microsoft Azure

Microsoft Azure-ը Microsoft ընկերության ստեղծած լայնածավալ և բազմաֆունկցիոնալ ամպային հարթակն է, որը առաջարկում է ավելի քան 200 ծառայություն և համապատասխանում է ինչպես փոքր ընկերությունների, այնպես էլ մեծ ձեռնարկությունների պահանջներին: Azure-ը հնարավորություն է տալիս կառուցել, փորձարկել, տեղակայել և կառավարել ծրագրեր՝ առանց ֆիզիկական ենթակառուցվածքներ ունենալու անհրաժեշտության:

Azure-ի հիմնական ծառայությունները ներառում են հաշվողական հզորություն, ինչպիսիք են Virtual Machines (VMs), որոնք տրամադրում են տարբեր օպերացիոն համակարգերով վիրտուալ սերվերներ, Azure Functions՝ առանձ սերվերների հաշվարկ, որտեղ կողը աշխատում է իրադարձությունների հիման վրա, և App Services՝ վեր ծրագրերի հոստինգի համար: Տվյալների պահեստավորման ոլորտում Azure-ն առաջարկում է Blob Storage, որը նախատեսված է մեծ չստրոկուլար պահպանման համար, Azure Files՝ կիսվող ֆայլային համակարգ, ինչպես նաև Data Lake Storage՝ մեծ չափի պահեստ տարբեր տիպի տվյալների համար:

Տվյալների բազաների ոլորտում ծառայությունների շարքում են Azure SQL Database, որը հանդիսանում է կառավարվող տվյալների բազա, Cosmos DB՝ գլոբալ բաշխված և բարձր աշխատանքային կարողություններով NoSQL տիպի տվյալների բազա, և PostgreSQL/MySQL ծառայություններ: Մերենայական ուսուցման և արհեստական բանականության համար Azure-ը տրամադրում է Azure Machine Learning ծառայությունը՝ մերենայական ուսուցման մոդելների մշակման, ուսուցման և տեղակայման ամբողջական միջավայր, ինչպես նաև Cognitive Services՝ տեսողության, խոսքի, լեզվի, որոնման և այլ տիպերի API-ների համախումբ:

Տվյալների պրոցեսավորման և ինտեգրման ոլորտում Azure-ի ծառայությունները ներառում են Azure Data Factory-ն՝ տվյալների տեղափոխման և փոխակերպման համար, և Azure Synapse Analytics՝ մեծածավալ տվյալների վերլուծության և բիզնես ինտելեկտի ապահովման համար: Ցանցային ծառայությունների շարքում ընդգրկված են Virtual Network, Load Balancer, VPN Gateway, Firewall, ինչպես նաև CDN, DDoS հարձակումներից պաշտպանություն և այլ ցանցային գործիքներ:

Azure-ի հիմնական առավելությունն այն է, որ այն ապահովում է ձկունություն, գլոբալ հասանելիություն և անվտանգություն, միաժամանակ նվազեցնելով կազմակերպությունների տեխնիկական և ֆինանսական բեռնվածությունը:

«Մեղալիոն» ճարտարապետություն

Ամպային տվյալների կառավարման ժամանակակից մոտեցումներում «մեղալիոն» ճարտարապետությունը հանդիսանում է որպես հստակ կառուցվածքային մոդել, որն էլ իր հերթին թույլ է տալիս կազմակերպություններին քայլ առ քայլ ձևակերպել, մաքրել և կարգավորել իրենց տվյալները: Այս մոտեցումը հիմնված է այն գաղափարի վրա, որ տվյալները պետք է անցնեն մի քանի հաջորդական քայլերով՝ յուրաքանչյուրում ձեռք բերելով ավելի հստակ և կարգավորված ձև:

«Մեղալիոն» ճարտարապետության հիմքում ընկած է շերտային տվյալային Lakehouse մոդելը: Առաջին շերտում, որը սովորաբար անվանում են բրոնզե, պահվում են հում, դեռևս չմշակված տվյալները: Այստեղ տվյալները գրեթե չեն ենթարկվում փոփոխության, դրանք պահպանվում են այսպես, ինչպես փոխանցվել են աղբյուրներից՝ թույլ տալով հետագա փոփոխության վերադառնալ սկզբնական տեղեկատվությանը և անհրաժեշտության դեպքում վերամշակել այն:

Բրոնզե շերտից հետո տվյալները տեղափոխվում են դեպի արծաթե մակարդակ, որտեղ արդեն իրականացվում է մաքրում, սխալների շտկում, տիպերի և տվյալների կառուցվածքի ձևափոխում: Այս փուլը կարելի է դիտարկել որպես տվյալների հիմնական ձևափոխման քայլ, քանի որ այստեղ ձևավորվում է միջանկյալ շերտ, որը միաժամանակ ունակ է ապահովելու թե՛ հաշվետվությունների համար անհրաժեշտ որակը, թե՛ խորքային վերլուծություններ իրականացնելու հիմքը:

Վերջին՝ ոսկե շերտը ներկայացնում է առավել մշակված և բիզնես տրամաբանությամբ հարստացված տվյալները: Այստեղ արդեն ամփոփվում են հաշվարկված ցուցանիշները, պատրաստվում են վերլուծական աղյուսակները, ստեղծվում են տվյալների այնպիսի կառուցվածքներ, որոնք անմիջապես կարող են օգտագործվել dashboard-ների հաշվետվություններում կամ մեքենայական ուսուցման մոդելներում: Ուսկե՛ շերտը հանդիսանում է ամբողջ ճարտարապետության նպատակակետը, որտեղ տվյալները հասնում են իրենց վերջնական արժեքին:

Այս երեք փոփային մոտեցումը ոչ միայն ապահովում է տվյալների հստակ ձևապարհը հում աղբյուրից մինչև պատրաստ վերլուծական արդյունք, այլ նաև մեծացնում է տվյալների կառավարելիությունը: Շերտերը թույլ են տալիս հեշտությամբ հետևել տվյալների ծագմանը, վերահսկել որակը և ապահովել, որ յուրաքանչյուր քայլում տվյալները օգտագործվեն ճիշտ նպատակի համար:

«Մեղալիոն» ճարտարապետությունը դարձել է մեծ տվյալների հետ աշխատող կազմակերպությունների կարևոր գործիքներից մեկը: Այն միաժամանակ տրամադրում է ձկունություն, կառուցվածք և մասշտարելիություն՝ դարձնելով տվյալների աշխարհը ավելի կանխատեսելի ու հուսալի:

ETL/ELT տվյալների ինտեգրման եղանակներ

ETL-ը տվյալների ինտեգրման դասական եղանակ է, որի նպատակը տարբեր աղբյուրներից ստացված տվյալները մաքրելն, փոխակերպելն ու պահպանելն է տվյալների պահեստում: ETL-ի գործընթացը բաղկացած է երեք փուլից՝ Extract, Transform և Load:

Առաջին փուլը՝ Extract, որի ընթացքում տվյալները վերցվում են տարբեր աղբյուրներից, օրինակ՝ տվյալների բազաներից, ինչպիսիք են SQL database կամ Oracle, API-ներից և վեբ ծառայություններից, CSV, JSON, Excel ֆայլերից, ինչպես նաև օպերացիոն համակարգերից կամ լոգ ֆայլերից: Այս փուլում տվյալները հիմնականում վերցվում են իրենց սկզբնական վիճակում (raw):

Երկրորդ փուլը՝ Transform, որտեղ տվյալները ձևափոխվում են միջանկյալ պրոցեսավորման միջավայրերում, օրինակ՝ ETL սերվերի վրա կամ հատուկ ETL գործիքի միջոցով: Զևափոխումը ընդգրկում է տվյալների մաքրում, ֆորմատների փոխակերպում, աղյուսակների միավորում, ագրեգացիաների և հաշվարկների ընդգրկում:

Վերջնին փուլը՝ Load, որը ապահովում է ձևափոխված տվյալների պահպանումը նպատակային տվյալների պահեստում, սովորաբար Data Warehouse-ում, ինչպիսիք են Azure SQL Data Warehouse կամ Teradata: ETL ինտեգրման եղանակը արդյունավետ է, եթե անհրաժեշտ է ապահովել ծանր և բարդ հաշվարկներ ձևափոխման փուլում, և եթե նպատակային պահեստը չունի մեծ հաշվարկային հզորություն, եթե անհրաժեշտ է նախապես մաքրել և հետո միայն պահել վավերացված տվյալներ, ինչպես նաև եթե տվյալները անկանոն կառուցվածքով են:

ELT-ն (Extract - Load - Transform) մոտեցումը տարբերվում է ETL-ից նրանով, որ գործընթացի փուլերի հերթականությունը փոխված է, և վերափոխումները կատարվում են հենց նպատակային տվյալների պահեստում: Սա դարձել է հնարավոր ժամանակակից ամպային տվյալների պահեստների մեծ հաշվարկային հզորության շնորհիվ: ELT-ում առաջին փուլը Extract է, որի ընթացքում աղբյուրներից ստացվում են տվյալները՝ առանց մեծ վերափոխումների: Երկրորդ փուլը՝ Load, ենթադրում է raw ֆորմատի տվյալների անմիջական պահպանումը տվյալների պահեստում, օրինակ՝ Azure Synapse Analytics, Snowflake, BigQuery կամ Databricks Delta Lake: Երրորդ փուլը՝ Transform, իրականացվում է հենց պահեստի ներսում՝ օգտագործելով SQL-ի, Spark-ի կամ MPP (Massively Parallel Processing) հաշվարկային մեխանիզմները:

ELT մոտեցման առավելությունն այն է, որ ժամանակակից պահեստներն ունեն մեծ հաշվարկային հզորություն և կարող են արագ մշակել հսկայական ծավալի տվյալներ: ELT-ն հատկապես օգտակար է, եթե աշխատում ենք մեծ ծավալի տվյալների հետ, եթե պահեստը ունի բարձր հաշվարկային հզորություն, եթե վերափոխումները կարելի է իրականացնել SQL կամ Spark պահեստի ներսում, ինչպես նաև եթե անհրաժեշտ է արագ բեռնել տվյալները և միայն հետո մշակել դրանք: Դրանց համեմատության աղյուսակը բերված է ներքում:

Բնութագիր	ETL	ELT
Վերափոխման տեղը	Միջանկյալ համակարգ	Տվյալների պահեստում
Առավել հարմար է	Փոքր և միջին տվյալների համակարգերի համար	Մեծ տվյալների, ամպային պահեստների համար
Հաշվարկային ծանրաբեռնվածություն	ETL սերվերի վրա	Cloud/Data Warehouse -ի վրա
Տվյալների պահպանում	Պահպանվում են միայն վերափոխված տվյալները	Պահպանվում են raw տվյալները և վերափոխված տվյալները
Գործիքներ	SSIS, Informatica, Talend	Azure Synapse, Snowflake, BigQuery, Databricks
Ակզրնական ձկունություն	Ավելի կոշտ գործընթացներ	Շատ ձկուն՝ մոդելավորման և SQL transform-ի համար

«Մեղալիոն» Ճարտարապետություն

Վերջին տարիներին տվյալային համակարգերի զարգացման մեջ լայնորեն կիրառվող մեղալիոնային ճարտարապետությունը ներկայացվում է որպես հատակ կառուցվածք, որի նպատակն է ստեղծել հասկանալի, վերահսկելի և երկարաժամկետ վստահելի տվյալային միջավայր: Այդ մոտեցումը հիմնված է այն գաղափարի վրա, որ տվյալները պետք է անցնեն հաջորդական փուլերով, որտեղ յուրաքանչյուր փուլ ապահովում է ավելի բարձր որակ և մաքրություն: Այս պատճառով էլ ամբողջ գործընթացը իրենց անվանումը ստացել է «մեղալիոն»՝ շերտերով կառուցված համակարգի գաղափարով:

Տվյալների առաջին կանգառը համարվում է բրոնզե շերտը, որտեղ դրանք պահպում են հում տեսքով՝ այնպես, ինչպես ստացվում են տարբեր աղբյուրներից: Այստեղ կարող են գտնվել տարբեր ձևաչափերով տվյալներ՝ լոգեր, սենսորների հոսքեր, բազաների արտահանումներ կամ այլ չմշակված ֆայլեր: Այս փուլում տվյալների վրա որևէ փոփոխություն չի կատարվում, քանի որ նպատակն է պահպանել դրանց ամբողջականությունը և ապահովել հետագա քայլերի համար անհրաժեշտ հուսալի հիմք: Բրոնզե շերտը նաև պահպանում է տվյալների ծագումնաբանությունը, ինչը կարևոր է հետագա վերլուծությունների և առողջության համար:

Հաջորդ փուլը արծաթե շերտն է, որտեղ տվյալները ենթարկվում են մաքրումի և կառուցվածքային ընդհանրացման: Այս փուլում վերացվում են կրկնօրինակները, շտկվում են սխալ կամ անհամապատասխան ձևաչափերը, լրացվում են բացակայող արժեքները, և տարբեր աղբյուրներից ստացված տվյալները համարվում են միմյանց հետ: Արծաթե շերտը դարձնում է տվյալները կանոնավոր, միատեսակ և տրամաբանորեն կապված: Սա այն միջին մակարդակն է, որտեղ արդեն կարող է կիրառվել բիզնես տրամաբանության մի մասը, սակայն տվյալները դեռ չեն հասել իրենց վերջնական մոդելին:

Ճարտարապետության ամենավերջին և ամենահաստատուն փուլը ուսկե շերտն է: Այստեղ տվյալները դառնում են առավել որակյալ, ամբողջապես մաքրված և պատրաստ բիզնես վերլուծության համար: Այս փուլում դրանք խմբագրվում և ազրեգացվում են՝ համապատասխանեցված կազմակերպության պահանջներին: Ուսկե շերտից առաջացած տվյալները օգտագործվում են հաշվետվությունների, վերլուծական մոդելների, ցուցանիշների հաշվարկների և կառավարման որոշումների համար: Շատ կազմակերպություններ այս շերտը դիտարկում են որպես իրենց միավորային ճշմարտության աղբյուր, քանի որ այստեղ գտնվող տվյալները ունեն առավելագույն վստահելիություն:

Մեղալիոնային ճարտարապետությունը գնահատվում է իր պարզությամբ և բարձր արդյունավետությամբ: Այն թույլ է տալիս կազմակերպություններին ունենալ պատշաճ

տվյալային համակարգ, որտեղ տվյալների որակը բարձրանում է աստիճանաբար, իսկ ամբողջ գործընթացը դառնում է կանխատեսելի և թափանցիկ: Այս մոտեցումը դարձնում է տվյալների կառավարումն ավելի հեշտ և ապահովում է, որ տարբեր բաժիններ ու թիմեր օգտվեն նույն հուսալի տվյալներից՝ իրենց որոշումների և վերլուծությունների համար:

OLTP/OLAP համակարգեր

OLTP (Online Transaction Processing) համակարգերը նախատեսված են օպերացիոն աշխատանքների համար, որտեղ առաջնահերթ են արագությունն ու կայունությունը, ինչպես նաև տվյալների մանրամասն վերահսկումը: Դրանք զբաղվում են առօրյա գործարքներով, որոնք փոքր ծավալով են, բայց տեղի են ունենում հաճախակի: Այս համակարգերը կարող են միաժամանակ սպասարկել հազարավոր կամ միլիոնավոր կարճ տևողությամբ տրանզակցիաներ ապահովելով բարձր արագությունը և տվյալների ամբողջականությունը պահպանելով ACID սկզբունքները, որոնցում ընդգրկված են Atomicity, Consistency, Isolation և Durability: Տվյալները սովորաբար պահպում են նորմալիզացված բազաներում՝ բազմաթիվ աղյուսակների բաժանված, որպեսզի նվազեցնեն կրկնօրինակումները և բարձրացնեն ամբողջականությունը: Օրինակների շարքում կարելի է նշել բանկային համակարգերը, որոնք ապահովում են հաշվիների միջև փոխանցումները և քարտային վճարումները, առցանց խանութները, որոնք իրականացնում են պատվերների գրանցում, վճարումներ և պահեստային փոփոխություններ, ինչպես նաև CRM համակարգերը, որոնք թույլ են տալիս օպերացիոն կառավարում հաճախորդների տվյալների:

OLAP (Online Analytical Processing) համակարգերը նախատեսված են վերլուծական և բիզնես-վերլուծական խնդիրների համար, եթե անհրաժեշտ է աշխատել մեծ ծավալի պատմական տվյալների հետ, կատարել բարդ հարցումներ և ստանալ խորքային վերլուծություն: Այս համակարգերը ապահովում են բարդ վերլուծական հարցումներ, որոնք կարող են ներառել խմբավորումներ, ազրեգացիաներ և բազմաչափ վերլուծություններ: Տրանզակցիոն հաճախականությունը սովորաբար ցածր է, և հարցումները կարող են տևել վայրկյաններ կամ նույնիսկ րոպեներ: Տվյալները պահպում են դենորմալիզացված կառուցվածքով՝ օգտագործելով Star կամ Snowflake Schema, ինչը ապահովում է վերլուծական հաշվարկների արագագործությունը: Այս համակարգերը հաճախ օգտագործվում են BI և Dashboard-ների ստեղծման համար, օրինակ՝ Power BI, Tableau և reporting համակարգերում: Օրինակներ են հանդիսանում Azure Synapse Analytics՝ ամպային մեծ տվյալների և վերլուծության հարթակը, ինչպես նաև Data Warehouse համակարգերը, ինչպիսիք են Snowflake, Google BigQuery և Amazon Redshift:

Տվյալային խողովակաշարեր (Data Pipelines)

Տվյալային խողովակաշարերը տվյալների ավտոմատացված հոսքերի համակարգված հավաքածու է, որի միջոցով տվյալները վերցվում են տարբեր աղբյուրներից, մաքրվում, ձևափոխվում, մշակում և տեղափոխվում են նպատակային համակարգեր: Այն ապահովում է, որ տվյալները հոսեն շարունակաբար, կանխատեսելի և հուսալի ձևով՝ առանց մարդկային միջամտության: Data Pipeline-ը կարող է ներառել ինչպես պարզ գործողություններ, օրինակ՝ ֆայլի պատճենում, տեղափոխում, այնպես էլ բարդ գործընթացներ, ինչպիսիք են մեծ ծավալնով տվյալների վերլուծությունը, real-time տվյալների մշակումն կամ մեքենայական ուսուցման մոդելների թարմացումը:

Տվյալային խողովակաշարի հիմնական քայլերը սկսվում են տվյալների հավաքագրումով, որի ընթացքում տվյալները ստացվում են տարբեր աղբյուրներից, ինչպիսիք են տվյալների բազաները (SQL/NoSQL), API-ներն ու web services-ը, ֆայլային համակարգերը (CSV, JSON, Parquet) և IoT սարքերը: Այս փուլում տվյալները հիմնականում հավաքվում են իրենց raw տեսքով:

Մյուս փուլը ձևափոխումն է (ETL/ELT), որտեղ տվյալները անցնում են մաքրման, սխալների ուղղման, ֆորմատների ստանդարտացման, ֆիլտրավորման կամ բարդ վերափոխումների միջով: ETL-ի դեպքում վերափոխումները կատարվում են միջանկյալ համակարգում, իսկ ELT-ի դեպքում դրանք իրականացվում են տվյալների պահեստում՝ օգտագործելով SQL կամ Spark:

Վերջնական փուլը տվյալների տեղափոխումն է տարբեր պահեստներ. մշակված տվյալները տեղափոխվում են Data Warehouse՝ վերլուծության համար, Data Lake մեծ չմշակված տվյալների պահպանման համար, կամ Operational Store և analytics engine: Տեղափոխումն կարող է կատարվել batch, near-real-time կամ real-time ձևաչափով:

Մոնիթորինգն ու սխալների կառավարումը ապահովում են, որ pipeline-ը աշխատի առանց ընդհատման: Այս փուլը ներառում է workflow-ի վերահսկում, սխալների ավտոմատ հայտնաբերում և վերագործարկում, ինչպես նաև լոգավորում և ծանուցումների հավաքագրում (alerting):

Վերջնական օգտագործման փուլում pipeline-ի միջոցով ստացված տվյալները օգտագործվում են վերլուծական հաշվետվությունների (BI dashboards) ստեղծման, մեքենայական ուսուցման մոդելների ուսուցման և թարմացման, օպերացիոն համակարգերի որոշումների կայացման, ինչպես նաև real-time համակարգերի ավտոմատացման համար:

Data Pipeline-ների կարևորությունն այն է, որ դրանք ապահովում են տվյալների հուսալի շրջանառություն՝ անկախ ծավալից կամ աղբյուրների բազմազանությունից, կրկնության դիմացկուն են՝ կարող են ինքնաշխատ վերագործարկվել կամ շարունակել սխալի դեպքում, ավտոմատացված են և նվազեցնում են ձեռքով միջամտության անհրաժեշտությունը, ինչպես նաև ապահովում են, որ տվյալները հասանելի լինեն ճիշտ պահին, ինչը կարևոր է վերլուծության, որոշումների կայացման և AI/ML համակարգերի համար: Pipeline-ների միջոցով համակարգերը դառնում են ձկուն, քանի որ տվյալները կարող են հոսել մեծածավալ և արագ աճող համակարգերի միջև:

Microsoft Azure Data Factory

Azure Data Factory-ն Microsoft Azure-ի ամպային տվյալների ինտեգրման ծառայությունն է, որը թույլ է տալիս կառուցել ամբողջական տվյալային հոսքեր՝ սկսած տվյալների հավաքագրումից մինչև դրանց մշակում, փոխակերպում և պահպանում տարբեր նպատակային համակարգեր: Այն հանդիսանում է որպես կապող օղակ տարբեր տվյալների աղբյուրների միջև՝ ապահովելով ավտոմատացված, վերահսկվող և մասշտաբավորվող տվյալային գործընթացներ: Azure Data Factory-ն լայնորեն կիրառվում է ETL և ELT գործընթացների իրականացման, Data Pipeline-ների ստեղծման ու կառավարման, բազմաթիվ աղբյուրներից տվյալների տեղափոխելը դեպի պահեստներ, ինչպես նաև տարբեր քայլերի հաջորդական կառավարման, ժամանակացույցի և սխալների մշակման համար: Ծառայությունը հնարավորություն է տալիս աշխատել ինչպես Azure-ի ներքին, այնպես էլ արտաքին կամ տեղային տվյալների աղբյուրների հետ, ապահովելով ճկուն ինտեգրում:

Azure Data Factory-ի հիմնական բաղադրիչներն են Pipelines, Activities, Linked Services, Datasets և Integration Runtime: Pipelines-ը ADF-ի հիմնական կառուցվածքային միավորն է, որը ներկայացնում է գործողությունների տրամաբանական հավաքածու՝ միաժամանակ կատարելով գործընթացներ տվյալների հետ, օրինակ՝ տվյալների ստացում, մաքրում և պահեստավորում: Pipeline-ները կարող են աշխատել ժամանակացույցով, արձագանքել տրիգերներին, պարունակել կախվածություններ և լինել մասշտաբավորվող ու վերահսկվող:

Activities-ը Pipeline-ի մեջ գործող քայլերն են, որտեղ ամեն activity կատարում է կոնկրետ առաջադրանք, ինչպիսիք են տվյալների պատճենումը, SQL սքրիպտի գործարկումը կամ API կանչը: Հիմնական activity-ներն են Copy Activity, որը իրականացնում է տվյալների փոխանցում աղբյուրից դեպի նպատակային պահեստ, Data Flow, որը ապահովում է տվյալների գրաֆիկական ձևափոխում Spark-ին նման parallel engines-ի վրա, Stored Procedure Activity, որը գործարկում է SQL stored procedure տվյալների բազայում, Web/HTTP Activity, որը կատարում է վեր API-ների կանչ տվյալներ ստանալու կամ գործողություն անելու համար, և Execute Pipeline, որը թույլ է տալիս մեկ pipeline-ի կանչը մյուսից՝ բարդ workflows կառուցելու ժամանակ:

Linked Services-ը կապ է, որը սահմանում է, թե ինչպես է ստեղծվում կապը տվյալների աղբյուրի և նպատակային սերվիսի հետ: Այն պարունակում է connection string, authentication տվյալներ և endpoint հասցե: Օրինակներ են Azure Blob Storage, Azure SQL Database, On-premises SQL Server, Amazon S3, Salesforce և Oracle: Linked Services-ը ADF-ի օբյեկտների միջև կապերի հաստատման գործիքն է:

Dataset-ը ներկայացնում է այն տվյալային օբյեկտը, որի վրա աշխատում են activity-ները: Այն նկարագրում է տվյալների տեսակը, սիեման, ֆայլի ձևաչափը (CSV, JSON, Parquet): Dataset-ը տվյալներն ինքնին չեն, այլ դրանց ներկայացումն է:

Integration Runtime-ը ADF-ի հաշվարկային ենթակառուցվածքն է, որը իրականացնում է տվյալների փոխանցումն ու ձևափոխումը: Այն է engine-ը, որը ֆիզիկական կատարում է Copy Activity-ը, Data Flow-ը և այլ transform-ներ: IR-ի

տեսակներն են Azure IR, որը ամբողջությամբ կառավարվում է Azure-ի կողմից և օգտագործվում է cloud-to-cloud տվյալների շարժի համար, աշխատելով առանձ սերվեր, Self-hosted IR, որը տեղակայվում է օգտագործողի սեփական on-premises կամ VM միջավայրում և թույլ է տալիս աշխատել firewalled միջավայրերում On-premises ↔ Cloud տվյալների շարժի համար, և Azure SSIS IR, որը նախատեսված է SSIS (SQL Server Integration Services) packages-ները փոխադրել Azure և աշխատում է որպես լիարժեք SSIS կատարման միջավայր ամպում:

Գլխի ամփոփում

Azure-ը ապահովում է միասնական միջավայր, որտեղ տվյալների հավաքագրումը, փոխանցումը, պահպանումը և վերլուծությունը իրականացվում են նույն հարթակի վրա՝ բարձր արդյունավետության և ավտոմատացման մակարդակով: Այս մոդելը՝ թույլ է տալիս կազմակերպություններին արագ կառուցել տեխնոլոգիական լուծումներ՝ առանց բարդ ենթակառուցվածք կազմակերպելու:

Տվյալների շարժի և կառավարման միասնական մոտեցումը ապահովում է Azure Data Factory-ի միջոցով, որը կարգավորում է տվյալների հոսքերի ամբողջական ցիկլը՝ թույլ տալով կապել տարբեր բիզնես համակարգեր, ավտոմատացնել ժամանակացույցով գործընթացներ և ապահովել տվյալների մշտական թարմացում: Սա դարձնում է տվյալային ճարտարապետությունը կանխատեսելի, վերահսկվող և ճկուն:

Azure-ի տվյալային ծառայությունները համակցված աշխատելով ապահովում են տվյալների ճշգրտությունը, վերահսկվող շարժ, բռնկումների և սխալների ավտոմատ մշակելիություն, ինչպես նաև տվյալների հասանելիություն տարբեր թիմերի և գործիքների համար:

Վերլուծության և որոշումների արագացման համար Azure-ի վերլուծական ծառայությունները անմիջապես օգտագործում են Data Factory-ի միջոցով բեռնված և մշակված տվյալները, ինչը թույլ է տալիս ստեղծել իրական ժամանակի dashboard-ներ, կատարել պատմական տվյալների խորքային վերլուծություն, արագ պատրաստել մերենայական ուսուցման մոդելներ և ընդունել բիզնես որոշումներ:

Բիզնեսի պահանջներին համապատասխան Azure-ը ապահովում է ճկունություն՝ թույլ տալով տարբեր աշխատանքաներ տեղափոխել ամպ, ընդլայնել կամ նվազեցնել ռեսուրսները ըստ պահանջի: Սա հատկապես կարևոր է այն կազմակերպությունների համար, որոնք աշխատում են մեծ տվյալների ծավալներով կամ ունեն փոփոխական բեռնվածություն: