de-7. SQLによるデータ分析: GROUP BYを用いたグループ化 と集約

(データベース演習)

URL: https://www.kkaneko.jp/de/de/index.html

金子邦彦







① データベースのグループ化と集約の概念 を理解し、基本的な分析手法を習得

② Access の SQLビューで,集計のための SQL を実行. 結果を確認

7-1. イントロダクション

リレーショナルデータベースの仕組み

- データをテーブルと呼ばれる表形式で保存
- テーブル間は関連で結ばれる
- 複雑な構造を持ったデータを効率的に管理することを可能

商品

ID	商品名	単価
1	みかん	50
2	りんご	100
3	メロン	500

関連

購入

購入者	商品番号
X	1
X	3
Υ	2

商品テーブルと購入テーブル

商品

ID	商品名	単価
1	みかん	50
2	りんご	100
3	メロン	500

購入

商品番号

関連

購入 有	冏品番亏
X	1
X	3
Υ	2

Xさんは、**1**の**みかん**と,

3のメロンを買った

Yさんは、2のりんごを買った

購入テーブルの情報 商品

商品テーブルの情報

集約の方法

AVG, MAX, MIN, SUM:平均值,最大值,最小值,合計值を算出

COUNT:行数を計算

記録テー ブル

名前	得点	居室
徳川家康	85	1階
源義経	78	2階
西郷隆盛	90	3階
豊臣秀吉	82	1階
織田信長	75	2階

SELECT AVG(得点) FROM 記録;

結果:82

SELECT MAX(得点) FROM 記録;

結果:90

SELECT MIN(得点) FROM 記録;

結果: 75

SELECT SUM(得点) FROM 記録;

結果: 410

SELECT COUNT(*) FROM 記録;

結果:5

グループ化

グループ化は、同じ属性値を共有するデータを集めるプロセス.

例:科目の「国語」、「算数」、「理科」でグループ化

科目	受講者	得点
国語	Α	85
国語	В	90
算数	Α	90
算数	В	96
理科	Α	95



科目	受講者	得点
理科	Α	95

科目	受講者	得点
算数	Α	90
算数	В	96

例:受講者の「AI、「BIでグループ化

科目	受講者	得点
国語	Α	85
国語	В	90
算数	Α	90
算数	В	96
理科	A	95



科目	受講者	得点
国語	Α	85
算数	Α	90
理科	Α	95

科目	受講者	得点
国語	В	90
算数	В	96

それぞれの値ごとにグループに分けることで、データの分析が 容易になる

成績テーブルの「科目」によるグループ化

① 成績テーブルには科目、受講者、得点が記載されている

科目	受講者	得点
国語	Α	85
国語	В	90
算数	Α	90
算数	В	96
理科	Α	95

②科目の「国語」、「算数」、「理科」のグループを形成

科目	受講者	得点
国語	Α	85
国語	В	90

科目	受講者	得点
算数	Α	90
算数	В	96

科目	受講者	得点
理科	Α	95

データのグループ化を用いた集約

③ グループで、集約(行数、平均、合計など)を実施

科目	受講者	得点
国語	А	85
国語	В	90
算数	А	90
算数	В	96
理科	Α	95



科目	受講者	得点
国語	Α	85
国語	В	90

科目	受講者	得点
理科	Α	95

科目	受講者	得点
算数	Α	90
算数	В	96



国語

→ 行数: 2

平均:87.5 93 95

合計:175

算数

186

理科

SQL における GROUP BY の基礎

GROUP BY は、特定の属性(例:科目,受講者)を基準 として、グループ化を行う

集約をGROUP BYと組み合わせることで、グループごとの 集約結果を得ることができる

これにより、データの傾向を効率的に分析することが可能である。

GROUP BY の役割と書き方

- SQL 問い合わせ「SELECT ...」の中で、GROUP BY を使用してデータをグループ化する
- 1つ以上の属性を GROUP BY に指定してグループ化の基準とする。

すべての科目ごとに、受講者の数を計算

SELECT 科目, COUNT(*) FROM 成績 GROUP BY 科目;

科目	COUNT(*)
国語	2
理科	1
算数	2

データ分析とビジネスインテリジェンス

主要な分析手法

グループ化(データをまとめること)と集約(まとめた データの計算)

行数,平均,合計などの生成

・データ分析

カテゴリ別,時系列別(時間の経過に沿った変化)のデータ分析

・ビジネスインテリジェンス(BI)

売上のトレンド分析(傾向や変化を調べること)

顧客セグメント分析(顧客を年齢や購買傾向などで分類して分析すること)

7-2. 演習

いまから演習で行うこと、注意点

・次のテーブルを作成



【Access での注意点】

・SQLビューでは、<u>SQL文を1つずつ</u>実行

(複数まとめての一括実行ができない)

- CREATE TABLE では、「実行」の後、画面が変化しない が実行できている
- INSERT INTO では、「実行」の後、確認表示が出る。その後、画面が変化しないが実行できている

SQL 理解のための前提知識

〇 テーブル

データをテーブルと呼ばれる表形式で保存

ID	商品名	単価	
1	みかん	50	
2	りんご	100	
3	メロン	500	

購入者	商品番号
X	1
X	3
Υ	2

- 問い合わせ (クエリ)
- 問い合わせ(クエリ)は、データベースから必要なデータ を検索、加工するための指令
- SELECT, FROM, WHERE など、**多様**なコマンドが存在。
- 結合、集計、ソート、副問い合わせなど、高度な操作も可能

SQL によるテーブル定義

- ・テーブル名:成績
- 属性名:科目、受講者、得点
- ・属性のデータ型:テキスト、テキスト、数値
- データの整合性を保つための制約:なし

```
CREATE TABLE 成績 (
科目 TEXT,
受講者 TEXT,
得点 INTEGER);
```

データ追加のSQL

成績

作日	文神有	行品
国語	Α	85
国語	В	90
算数	Α	90
算数	В	96
理科	Α	95

```
INSERT INTO 成績 VALUES('国語', 'A', 85);
INSERT INTO 成績 VALUES('国語', 'B', 90);
INSERT INTO 成績 VALUES('算数', 'A', 90);
INSERT INTO 成績 VALUES('算数', 'B', 96);
INSERT INTO 成績 VALUES('理科', 'A', 95);
```



演習 1. Access の SQL ビューを用いたテーブル定義 とデータの追加

【トピックス】

- ・SQLビューを開く
- ・SQL文の編集
- create table
- insert into
- ・SQL文の実行

演習

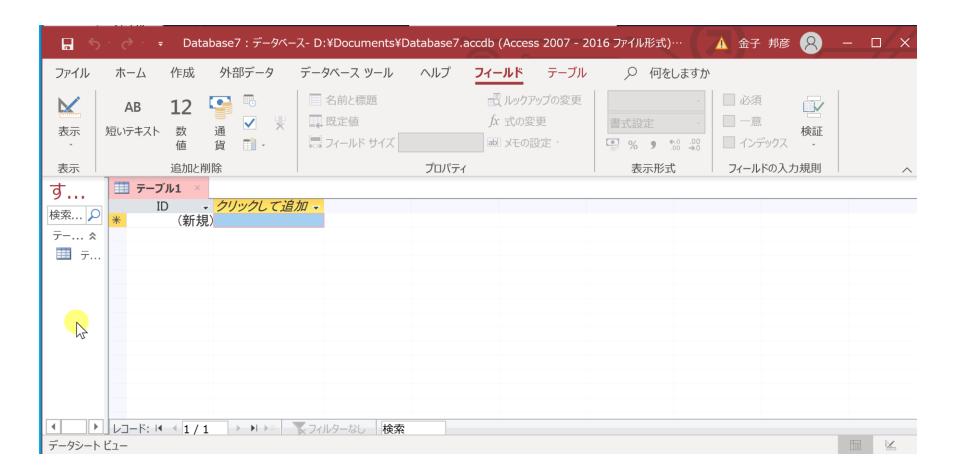
1. パソコンを使用する **前もって Access をインストールしておくこと**

2. Access を起動する

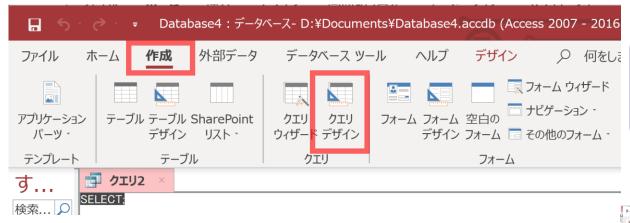
3. Access で、「**空のデータベース**」を選び、「**作成**」を クリック.



4. テーブルツール画面が表示されることを確認



5. 次の手順で、**SQLビュー**を開く.



①「作成」タブで、「クエリデザイン」をクリック





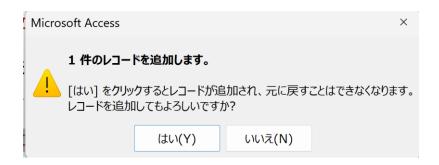
このような 表示が出た ときは 「**閉じる**」を クリック



②「**デザイン**」タブで、 「**表示**」を展開し「**SQL** ビュー」を選ぶ

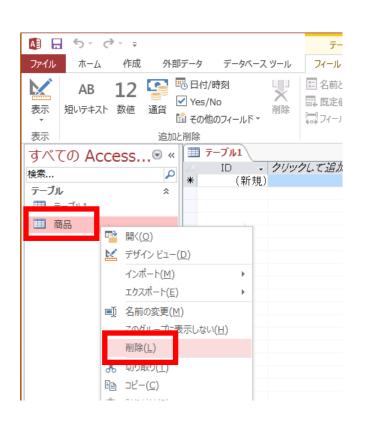
6. **SQL ビュー**に、次の SQL を 1 つずつ入れ、「**実行**」ボタンで、**SQL文**を実行.結果を確認

```
CREATE TABLE 成績 (
   科目 TEXT,
   受講者 TEXT,
   得点 INTEGER);
INSERT INTO 成績 VALUES('国語', 'A', 85);
INSERT INTO 成績 VALUES('国語', 'B', 90);
INSERT INTO 成績 VALUES('算数', 'A', 90);
INSERT INTO 成績 VALUES('算数', 'B', 96);
INSERT INTO 成績 VALUES('理科', 'A', 95);
```



INSERT INTOでは、「実行」の後、確認表示が出る。その後、**画面が変化しない**が実行できている

間違ってしまったときは、テーブルの削除 を行ってからやり直した方が早い場合がある



テーブルビューで、削除したいテーブルを**右クリック**して、 **削除**」

テーブルを削除するときは、 間違って必要な**テーブル**を削除しない ように、十分に注意する! (元に戻せない)



演習 2. SQL によるグループ 化と集約. Access の SQL ビューを使用.

【トピックス】

- 1. グループ化
- 2. 集約
- 3. GROUP BY
- 4. AVG
- 5. COUNT(*)
- 6. SUM

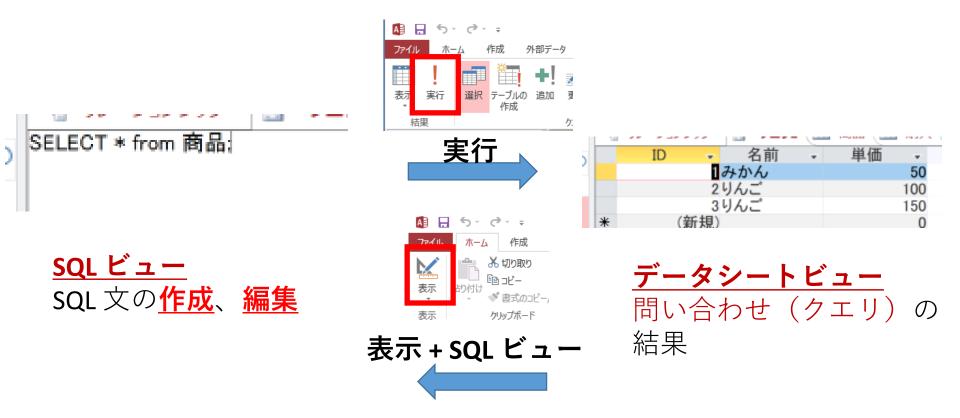
Access の SQL ビューを用いた問い合わせ

- ① Access の **SQLビュー**開く
- ② **SQL 文の編集**。**select, from, where** を使用例: select * from テーブル名 where 列1 = 値1;
- ③ SQL 文の実行

実行の結果、**データシートビュー**に画面が変わり、そこに**問い合わせの結果**が表示される

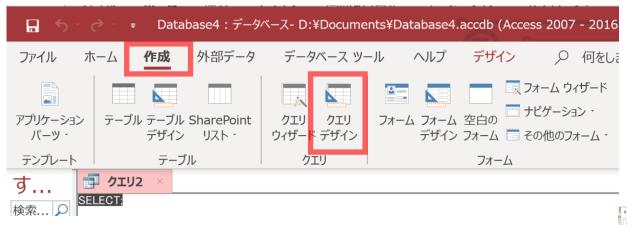
④ さらにSQL 文の編集、実行を続ける場合には、<u>画面を SQL</u>ビューに切り替える

SQL 問い合わせ(クエリ)で使用する2つのビュー



マウス操作でビューを切り替え

1. 次の手順で、**SQLビュー**を開く.



①「作成」タブで、「クエリデザイン」をクリック





このような 表示が出た ときは 「**閉じる**」を クリック



②「**デザイン**」タブで、 「**表示**」を展開し「**SQL ビュー**」を選ぶ

2. **SQL ビュー**に、次の SQL を 1 つずつ入れ、「**実 行**」ボタンで、**SQL文**を実行. 結果を確認

1. 単純な表示

SELECT * FROM 成績;

科目	v	受講者	-	得点	•
国語	Α				85
国語	В				90
算数	Α				90
算数	В				96
理科	Α				95

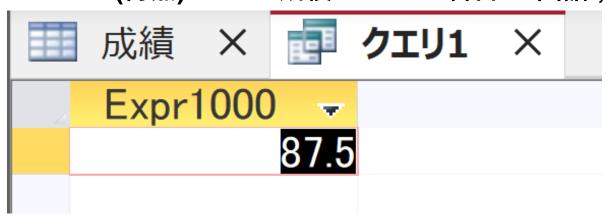
2. 得点の平均

SELECT AVG(得点) FROM 成績;



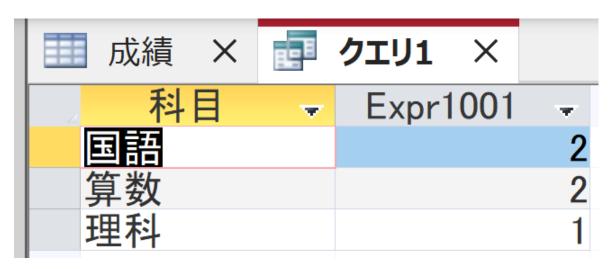
3. 国語の得点の平均

SELECT AVG(得点) FROM 成績 WHERE 科目 = '国語';



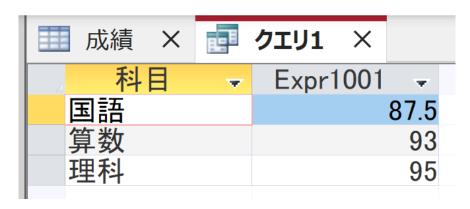
4. それぞれの科目の受講者数

SELECT 科目, COUNT(*) FROM 成績 GROUP BY 科目;

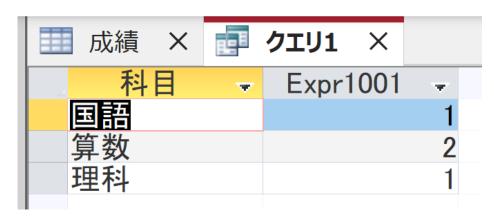


5. それぞれの科目の平均得点

SELECT 科目, AVG(得点) FROM 成績 GROUP BY 科目;

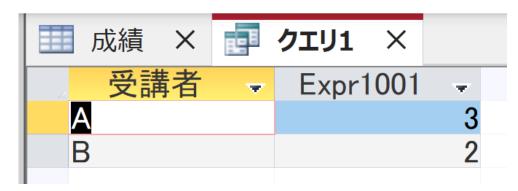


6. それぞれの科目について、得点が90点以上である受講者数 SELECT 科目, COUNT(*) FROM 成績 WHERE 得点 >= 90 GROUP BY 科目;



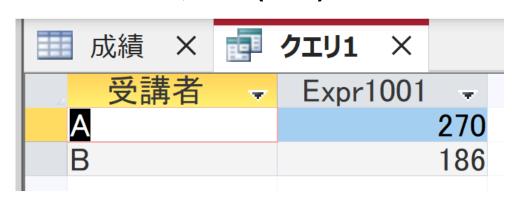
7. それぞれの受講者が受講している科目数

SELECT 受講者, COUNT(*) FROM 成績 GROUP BY 受講者;



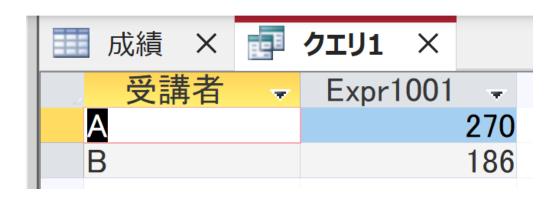
8. それぞれの受講者の得点合計

SELECT 受講者, SUM(得点) FROM 成績 GROUP BY 受講者;



9. それぞれの受講者の得点平均

SELECT 受講者, AVG(得点) FROM 成績 GROUP BY 受講者;



7-3. 実データを用いた演習

演習の目的と形式

•目的:実データを使い、グループ化と集約の有用性を確認する。SQLのスキルアップも行う

・形式:発展演習形式(**資料を見ながら各自実施し** てください)

演習の内容

SQL を用いたグループ化と集約、そのバリエーションと有用 性を知る

・米国成人調査データを利用

調査に協力した人たちの<u>年</u> 齢分布は?



教育と年収の関係を見る

教育	- 年収5万ドル -	Expr1002 -
10th	<=50K	871
10th	>50K	62
11th	<=50K	1115
11th	>50K	60
12th	<=50K	400
12th	>50K	33
1st-4th	<=50K	162
1st-4th	>50K	6
4年制大学	<=50K	3134
4年制大学	>50K	2221
5th-6th	<=50K	317
5th-6th	>50K	16
7th-8th	<=50K	606
7th-8th	>50K	40
9th	<=50K	487
9th	>50K	27
Preschool	<=50K	51
何らかの大学	<=50K	5904
何らかの大学	>50K	1387
高校	<=50K	8826
高校	>50K	1675
職業技術訓練校	<=50K	1021
職業技術訓練校	>50K	361
専門職大学院	<=50K	153
専門職大学院	>50K	423
大学院修士	<=50K	764
大学院修士	>50K	959
大学院博士	<=50K	107
大学院博士	>50K	306
短大、コミュニティカレッジ	<=50K	802
短大、コミュニティカレッジ	>50K	265

演習で使うデータベース

米国成人調査データ

(1994年、米国における統計調査データのうち 32561人分)

米国成人	調査データ	(
ID →	年齢 →	職業の分類・	教育 ▼	教育年数 →	職業	*	性別 🔻	週当たり労働時間。	母国	→ 年収5万ドノ→
1	39	州政府	4年制大学	13	管理、事務		男性	40	米国	<=50K
2	50	法人でない自営業	4年制大学	13	執行、経営		男性	13	3 米国	<=50K
3		民間	高校		各種取扱者、		男性	40	米国	<=50K
4		民間	11th		各種取扱者、		男性	40	米国	<=50K
5		民間	4年制大学	13	専門職		女性	40	キューバ	<=50K
6		民間	大学院修士	14	執行、経営		女性	40	米国	<=50K
7		民間	9th		その他のサー		女性		ジャマイカ	<=50K
8		法人でない自営業	高校		執行、経営		男性		米国	>50K
9		民間	大学院修士		専門職		女性		米国	>50K
10		民間	4年制大学	13	執行、経営		男性	40	米国	>50K
11	37	民間	何らかの大学	10	執行、経営		男性	80	米国	>50K
12		州政府	4年制大学	13	専門職		男性	40	コインド	>50K
13		民間	4年制大学		管理、事務		女性	30	米国	<=50K
14		民間	短大、コミュニティカレッジ	12	販売		男性	50	米国	<=50K
4 =	*^	P88	THY 35 4 + 4 4 4 5 1 1 6 未 4 六	4.4	— <i>I</i> ⊬ <i>I</i> .⊘⊤m		FFI .W4	A 2	30	NEAR.

※ このデータを使います

(演習では、特定の職業、学歴、性別、母国を差別的に見ないようにしてください)

データの出典: Lichman, M. (2013).

UCI Machine Learning Repository [http://archive.ics.uci.edu/ml].

Irvine, CA: University of California, School of Information and Computer Science (米国)

演習用のデータベースファイル

・演習用の Access データベースファイル

セレッソの利用者は, セレッソからもダウンロード可能 ファイル名: **db4-4.accdb**

「コンテンツの有効化」のメッセージが出たときは、確認のうえ、次にすすむ。

! **セキュリティの警告** 一部のアクティブ コンテンツが無効にされました。クリックすると詳細が表示されます。

コンテンツの有効化

つぎのような表示が出たときは、確認のうえ、「はい」



米国成人調査データ



SELECT 年齢, count(*) FROM 米国成人調査データ GROUP BY 年齢;

調査に協力した人たちの年齢分布は?

■ 米国成人訓	音データ 🗗 クエリ1
年齢 🕶	Expr1001 -
17	395
18	550
19	712
20	753
21	720
22	765
23	877
24	798
25	841
26	785
27	835
28	867

SELECT 教育, count(*) FROM 米国成人調査データ GROUP BY 教育;

調査に協力した人たちの教育の分布は?

5	教育	Expr1001 -
1	10th	933
	11th	1175
	12th	433
	1st-4th	168
	4年制大学	5355
	5th-6th	333
	7th-8th	646
	9th	514
	Preschool	51
	何らかの大学	7291
	高校	10501
	職業技術訓練校	1382
	専門職大学院	576
	大学院修士	1723
	大学院博士	413
	短大、コミュニティカレッジ	1067

SELECT 週当たり労働時間, count(*)

FROM 米国成人調査データ

GROUP BY 週当たり労働時間;

調査に協力した人たちの週当たり<u>労働時間の分布</u>は?

※ 米国成人調査データ	ge F	クエリ1	
週当たり労働時間	w	Expr1001 -	
	0	20	
t l	2	32	
	3	39	
	4	54	
•	5	60	
	6	64	
	7	26	
	8	145	
	9	18	
	10	278	
	11	11	
	12	173	
	13	23	
	14	34	
	15	404	
	16	205	
	17	20	

SELECT 年収5万ドル以上か, count(*) FROM 米国成人調査データ GROUP BY 年収5万ドル以上か;

年収5万ドル以上の人とそうでない人の人数



SELECT 教育, 年収5万ドル以上か, count(*)

FROM 米国成人調査データ

GROUP BY 教育, 年収5万ドル以上か;

教育と年収の関係を見る

教育	- 年収5万ドル -	Expr1002 -
10th	<=50K	871
10th	>50K	62
11th	<=50K	1115
11th	>50K	60
12th	<=50K	400
12th	>50K	33
1st-4th	<=50K	162
1st-4th	>50K	6
4年制大学	<=50K	3134
4年制大学	>50K	2221
5th-6th	<=50K	317
5th-6th	>50K	16
7th-8th	<=50K	606
7th-8th	>50K	40
9th	<=50K	487
9th	>50K	27
Preschool	<=50K	51
何らかの大学	<=50K	5904
何らかの大学	>50K	1387
高校	<=50K	8826
高校	>50K	1675
職業技術訓練校	<=50K	1021
職業技術訓練校	>50K	361
専門職大学院	<=50K	153
専門職大学院	>50K	423
大学院修士	<=50K	764
大学院修士	>50K	959
大学院博士	<=50K	107
大学院博士	>50K	306
短大、コミュニティカレッジ	<=50K	802

全体まとめ

- グループ化:同じ属性値を共有するデータをまとめる. 科目や受講者といった特定の属性値でデータを分類し、分析を容易にする.
- GROUP BY: SQLにおいてグループ化を実行するためのもの、それぞれのグループに対して集計処理を行うことができる.
- 集約関数:グループ化されたデータに対して計算を行う. 代表的なものとしてAVG(平均値), COUNT(行数),
 SUM(合計値), MAX(最大), MIN(最小)がある.



① データベースのグループ化と集約の概念を理解し、 基本的な分析手法を習得

リレーショナルデータベースでは、テーブル形式でデータを管理する. GROUP BYを用いて同じ属性値(科目や受講者など)を持つデータをグループ化し、そのグループに対してAVG、COUNT、SUMなどの集約関数を適用することで、平均値や合計値などの統計値を算出できる. これにより、データの傾向や特徴を効率的に分析することが可能となる.

② Access の SQLビューで,集計のための SQL を 実行. 結果を確認

AccessのSQLビューでは、SELECT文にGROUP BY 句と集約関数を組み合わせたSQL文を作成・実行できる。実行結果はデータシートビューに表示され、即座に確認が可能である。WHERE句と組み合わせることで、特定の条件を満たすデータに対する集計も実現できる。SQLビューとデータシートビューを行き来しながら、段階的にSQLの理解を深めることができる。

発展演習 1. テーマ: 科目別の平均得点の計算

目的: GROUP BY を使用して、科目ごとに平均得点を計算する方法を学ぶ。

成績テーブルから、科目ごとに平均得点を計算するSQ文を書いてください。

ヒント: AVG と GROUP BY を組み合わせて使用し、科目 でグループ化します。

発展演習 2. テーマ: 得点が90点以上の受講者数の計算

特定の得点基準を満たす受講者の数を科目ごとに調べる。

各科目について、得点が90点以上である受講者数をカウントするSQL文を書いてください。

ヒント: WHERE を使って得点が90点以上のものを選択。 GROUP BY で科目ごとにグループ化します。 ・発展演習1の正解例

SELECT 科目, AVG(得点) FROM 成績 GROUP BY 科目;

・発展演習2の正解例

SELECT 科目, COUNT(*) FROM 成績 WHERE 得点 >= 90 GROUP BY 科目;