

Databases

Normalization

Danilo Montesi

danilo.montesi@unibo.it

Normal Forms

- “Normal form” is a property of a relational database, that guarantees its quality, *i.e.*, the absence of certain flaws
- When a relation is not in normal form:
 - it has redundancies
 - it can have undesired behaviours during updates
- Normal forms are usually defined on the relational model, but they make sense also in other contexts, for example in the E-R model

Normalization

- Such task allows to transform non-normalized schema into schema that satisfy the normal form
- Normalization has to be used as a **verification technique** to test the result of the database design
- It is not a methodology for database design



A Relation with Anomalies

<u>Employee</u>	Wage	<u>Project</u>	Budget	Role
Jones	20	Mars	2	Technician
Smith	35	Jupiter	15	Designer
Smith	35	Venus	15	Designer
Williams	55	Venus	15	Chief
Williams	55	Jupiter	15	Consultant
Williams	55	Mars	2	Consultant
Brown	48	Mars	2	Chief
Brown	48	Venus	15	Designer
White	48	Venus	15	Designer
White	48	Jupiter	15	Director

Anomalies

- **Redundancy**: each employee's wage is repeated throughout the relation
- **Update Anomaly**: when an employee's wage changes, then we have to change all its occurrences
- **Deletion Anomaly**: when an employee stops participating in all its project, it is completely removed from the database
- **Insertion Anomaly**: we cannot create an employee without an associated project



Why is this Situation Undesirable?

- Because different pieces of information are represented within the same relation
 - Employees and their wages
 - Projects and their budgets
 - The role of each employee within a project they're working on

Functional Dependencies

- In order to study in a systematic way the concept we just introduced, we need to use the **functional dependency**
- Functional dependency is a specific integrity constraint for the relational model that describes functional bounds among the attributes of a relation

The Idea

- Each employee always has the same wage (even if he participates in several projects)
- Each project always has the same budget
- Each employee in each project always has the same role (even though they may have different functions in different projects)

Functional Dependencies: Definition

- Given a relation r having a schema $R(X)$
- Given two non-empty subsets of attributes Y and Z of X
- A **functional dependency** (FD) exists between Y and Z if:
 - for each couple of tuples t_1 and t_2 of r having the same values on the attributes Y , it results that t_1 and t_2 have the same values also on the attributes Z



Functional Dependencies: Notation

$$Y \rightarrow Z$$

Examples:

Employee \rightarrow Wage

Project \rightarrow Budget

Employee Project \rightarrow Role

Functional Dependencies: an Example

<u>Employee</u>		<u>Wage</u>		<u>Project</u>		<u>Budget</u>		<u>Role</u>
Jones		20		Mars	●	2		Technician
Smith	●	35		Jupiter	●	15		Designer
Smith	●	35		Venus	●	15		Designer
Williams	●	55		Venus	●	15		Chief
Williams	●	55		Jupiter	●	15		Consultant
Williams	●	55		Mars	●	2		Consultant
Brown	●	48		Mars	●	2		Chief
Brown	●	48		Venus	●	15		Designer
White	●	48		Venus	●	15		Designer
White	●	48		Jupiter	●	15		Director

- Employee → Wage
- Project → Budget
- Employee Project → Role

FD: Some Properties

Employee Project \rightarrow Project

- Such FD is “trivial” (it is always true)
- A FD $Y \rightarrow A$ is **nontrivial** if one of the following conditions are met:
 - A is an attribute, and doesn't belong to Y
 - A is a set of attributes and none of the attributes in A belong to Y



Anomalies Depend on some FDs

- Employees must have only one Wage

Employee \rightarrow Wage

- Projects must have only one Budget

Project \rightarrow Budget



Not all the FD Provoke Anomalies

- In each Project, an Employee has only one Role

Employee Project \rightarrow Role

- The constraint is “trivially” satisfied because **Employee Project** is a key

FD and Anomalies

- The first two FD **are not keys** and **cause anomalies**
- The third FD **is a key** (Employee Project) and **does not cause anomalies**
- The relation contains **some informations linked to the key** and **other informations linked to attributes** that do **not** compose a **key**
- Hence, **anomalies are caused by heterogeneous informations:**
 - Employee's properties (the Wage)
 - Projects' properties (the Budget)
 - Properties for the key **Employee Project** (the Role)



Boyce-Codd Normal Form (BCNF)

- **Normal forms** are (useful) properties that are satisfied only in absence of anomalies, by defining constraints on functional dependencies
- The most important normal form the the one named after Boyce and Codd (BCNF)

Definition:

A relation **r** is in **BCNF** if, for each functional dependency (non trivial) $X \rightarrow Y$ defined on **r**, **X** has a key **K** of **r**, that is, **X** is superkey for **r**



When a Relation does not Satisfy BCNF

- In most cases, we can replace it with two or more normalized relations satisfying the BCNF. Such process is called **normalization**
- This process is based on a simple criteria:
 - If a relation represent more than one dependent concept, then it must be decomposed in smaller relations, one for each concept



Decomposition Example (1)

<u>Employee</u>	<u>Wage</u>	<u>Project</u>	<u>Budget</u>	<u>Role</u>
Jones	20	Mars	2	Technician
Smith	35	Jupiter	15	Designer
Smith	35	Venus	15	Designer
Williams	55	Venus	15	Chief
Williams	55	Jupiter	15	Consultant
Williams	55	Mars	2	Consultant
Brown	48	Mars	2	Chief
Brown	48	Venus	15	Designer
White	48	Venus	15	Designer
White	48	Jupiter	15	Director



Decomposition Example (2)

<u>Employee</u>	Wage
Jones	20
Smith	35
Williams	55
Brown	48
White	48

<u>Project</u>	Budget
Mars	2
Jupiter	15
Venus	15

<u>Employee</u>	<u>Project</u>	Role
Jones	Mars	Technician
Smith	Jupiter	Designer
Smith	Venus	Designer
Williams	Venus	Chief
Williams	Jupiter	Consultant
Williams	Mars	Consultant
Brown	Mars	Chief
Brown	Venus	Designer
White	Venus	Designer
White	Jupiter	Director

Decomposition with Loss!

Employee	Project	Office
Jones	Mars	Rome
Smith	Jupiter	Milan
Smith	Venus	Milan
White	Saturn	Milan
White	Venus	Milan

Employee → Office

Employee	Office
Jones	Rome
Smith	Milan
White	Milan

Project → Office

Project	Office
Mars	Rome
Jupiter	Milan
Venus	Milan
Saturn	Milan

We Try to Rebuild

Employee	Office
Jones	Rome
Smith	Milan
White	Milan



Office	Project
Rome	Mars
Milan	Jupiter
Milan	Venus
Milan	Saturn



Employee	Office	Project
Jones	Rome	Mars
Smith	Milan	Jupiter
Smith	Milan	Venus
Smith	Milan	Saturn
White	Milan	Jupiter
White	Milan	Saturn
White	Milan	Venus



Employee	Office	Project
Jones	Rome	Mars
Smith	Milan	Jupiter
Smith	Milan	Venus
White	Milan	Saturn
White	Milan	Venus

DIFFERENT FROM THE ORIGINAL RELATION!

Memories from Algebra

- The result of the join between R_1 and R_2 has a number of tuples between zero and $|R_1| \times |R_2|$

$$0 \leq |R_1 \bowtie R_2| \leq |R_1| \times |R_2|$$



- If the join involves a key from R_2 , then the number of the resulting tuples is within 0 and $|R_1|$

$$0 \leq |R_1 \bowtie R_2| \leq |R_1|$$



- If the join involves a key from R_2 and a Referential Integrity Constraint, the number of the tuples is $|R_1|$

$$|R_1 \bowtie R_2| = |R_1|$$

Decomposition: Lossless Join Property

Definition:

A relation r can be **decomposed lossless** in two relations $q(X)$ and $s(Y)$ if the join of the projection of r on X and Y is the same as r :

$$\pi_X(r) \bowtie \pi_Y(r) = r$$

This property is verified if the common attributes contains a key for at least one of the decomposed relations

Lossless Condition

- Given a relation $r(X)$ and X_1 and X_2 subset of X , such that $X_1 \cup X_2 = X$
- Given $X_0 = X_1 \cap X_2$
- $r(X)$ can be **decomposed lossless** in in two relations $q(X_1)$ and $s(X_2)$ if:
 - $X_0 \rightarrow X_1$ is satisfied or
 - $X_0 \rightarrow X_2$ is satisfied



Decomposition without Loss

Employee	Project	Office
Jones	Mars	Rome
Smith	Jupiter	Milan
Smith	Venus	Milan
White	Saturn	Milan
White	Venus	Milan

Employee	Office
Jones	Rome
Smith	Milan
White	Milan

Employee	Project
Jones	Mars
Smith	Jupiter
Smith	Venus
White	Saturn
White	Venus

Example: Insert a Tuple (1)

- Suppose that a new tuple (White, Mars, Milan) is inserted

Employee	Project	Office
Jones	Mars	Rome
Smith	Jupiter	Milan
Smith	Venus	Milan
White	Saturn	Milan
White	Venus	Milan
White	Mars	Milan



Example: Insert a Tuple (2)

Employee	Project	Office
Jones	Mars	Rome
Smith	Jupiter	Milan
Smith	Venus	Milan
Smith	Mars	Milan
White	Saturn	Milan
White	Venus	Milan
White	Mars	Milan

Employee	Office
Jones	Rome
Smith	Milan
White	Milan

Project	Office
Mars	Rome
Jupiter	Milan
Venus	Milan
Saturn	Milan
Mars	Milan

Example: Insert a Tuple (3)

Employee	Project	Office
Jones	Mars	Rome
Smith	Jupiter	Milan
Smith	Venus	Milan
White	Saturn	Milan
White	Venus	Milan
White	Mars	Milan

Employee → Office

Project → Office

Employee	Office
Jones	Rome
Smith	Milan
White	Milan

Employee	Project
Jones	Mars
Smith	Jupiter
Smith	Venus
White	Saturn
White	Venus
White	Mars

Dependency-preserving Decompositions

We say that a decomposition **preserves the dependencies** if each functional dependency of the original schema involves attributes that appear all together in one of the decomposed schemas

■ Project \rightarrow Office is not preserved

Decompositions Quality

- The process of normalization through decomposition must also confirm the existence of additional properties that the relational schemas, taken together, should have:
 - the **lossless join property**, that assures the rebuilding of the original information
 - the **dependency preservation**, that assures the keeping of the original integrity constraints

A Relation not in Normal Form

Chief	<u>Project</u>	<u>Office</u>
Smith	Mars	Rome
Johnson	Jupiter	Milan
Johnson	Mars	Milan
White	Saturn	Milan
White	Venus	Milan

Project Office → Chief

Chief → Office



Decomposition has some Problems

Project Office → Chief

- Involves all the attributes, so no decomposition could preserve the dependency
- In some cases BCNF “cannot be reached”

Third Normal Form (3NF)

Definition:

A relation r is in **third normal form** if, for each non-trivial FD $X \rightarrow Y$ on r , at least one of the following conditions is met:

- $K \subset X$, X is superkey in r
- Each attribute Y is in at least one key of r

BCNF and 3NF

- BCNF is stricter than 3NF (3NF admits relations with some anomalies)
- 3NF can always be reached (there is a theorem)
- If a relation has only one key, it is in BCNF if and only if it is in 3NF

A Relation not in Normal Form

Chief	<u>Project</u>	<u>Office</u>
Smith	Mars	Rome
Johnson	Jupiter	Milan
Johnson	Mars	Milan
White	Saturn	Milan
White	Venus	Milan

Project Office \rightarrow Chief \Leftarrow **X is key**

Chief \rightarrow Office \Leftarrow **Y \subset key**

3NF Decomposition

- Create a relation for each set of attributes involved in a functional dependency
- Check that in the end at least a relation has a key of the original relation
- It depends on the found dependencies

A Possible Solution

- If the relation is not normalized, decompose it to 3NF
- Then, check if the resulting schema is in BCNF
- In most cases, decomposing to reach the 3NF also allows the BCNF to be reached

A Relation that cannot be put in BCNF

Chief	<u>Project</u>	<u>Office</u>
Smith	Mars	Rome
Johnson	Jupiter	Milan
Johnson	Mars	Milan
White	Saturn	Milan
White	Venus	Milan

Project Office → Chief

Chief → Office



A Possible Reorganization

Chief	<u>Project</u>	<u>Office</u>	Dept.
Smith	Mars	Rome	1
Johnson	Jupiter	Milan	1
Johnson	Mars	Milan	1
White	Saturn	Milan	2
White	Venus	Milan	2

Dept. Office → Chief

Chief → Office Dept.

Project Office → Dept.



Decomposition in BCNF

<u>Chief</u>	Office	Dept.
Smith	Rome	1
Johnson	Milan	1
White	Milan	2

<u>Project</u>	Office	Dept.
Mars	Rome	1
Jupiter	Milan	1
Mars	Milan	1
Saturn	Milan	2
Venus	Milan	2

A Theory for Dependency

The previous concepts could be automated in an algorithmic process. That is:

- Given a **relation** and a set of **functional dependencies** over it
- Generate a decomposition of such relation containing only **relations in normal form** that also satisfy the aforementioned decomposition properties:
 - lossless decomposition
 - dependencies preservation



Functional Dependencies: Implications

From the valid functional dependencies, we can determine other dependencies, we say that the first imply the seconds:

- A set of functional dependencies F implies f if each relation satisfying all the dependencies in F satisfies also f

An Example

Employee	Type	Wage
Smith	3	30.000
Johnson	3	30.000
White	4	50.000
Williams	4	50.000
Brown	5	72.000

Employee \rightarrow Type *and* Type \rightarrow Wage

IMPLY

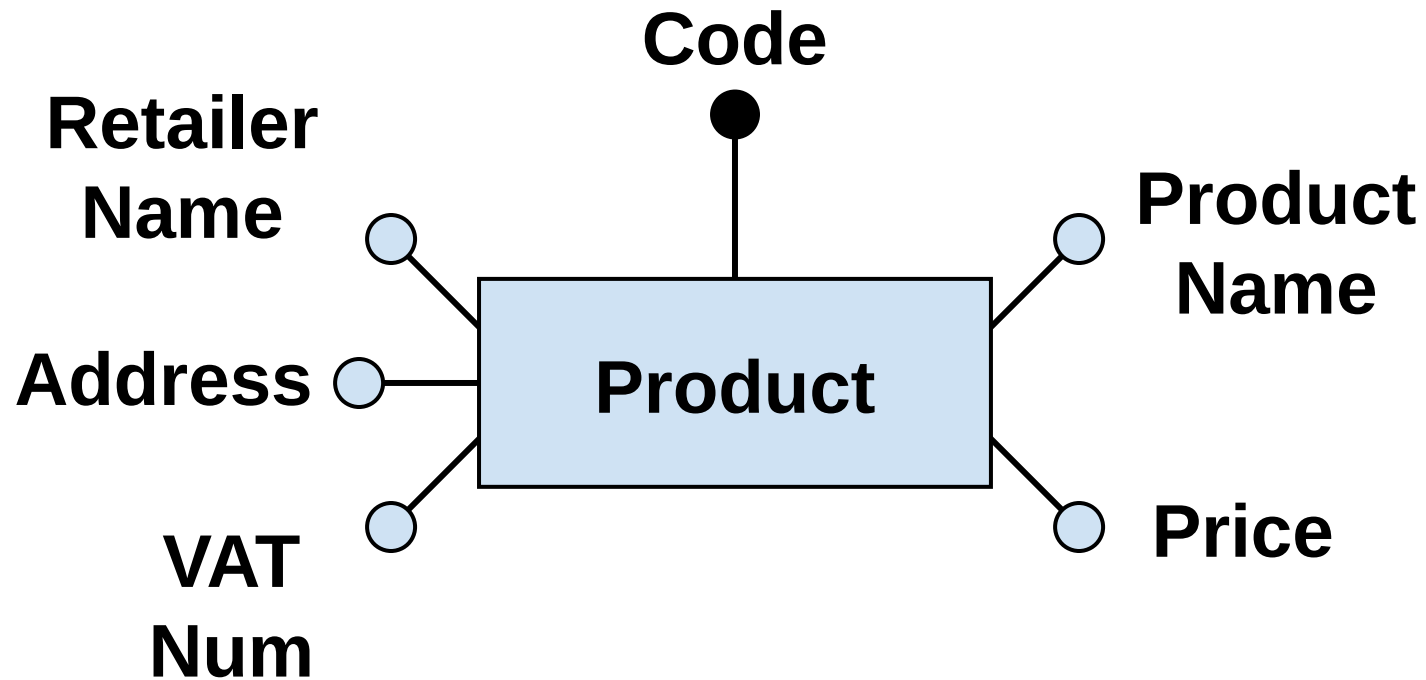
Employee \rightarrow Wage

Each relation satisfying the first two dependencies, satisfies also the third one

Design and Normalization

- Normalization theory can be used within the logical design to check the schema of the final relation
- It could be also used during the conceptual design phase to verify the quality of the conceptual schema

Normalization over Entities



VAT Num → RetailerName Address

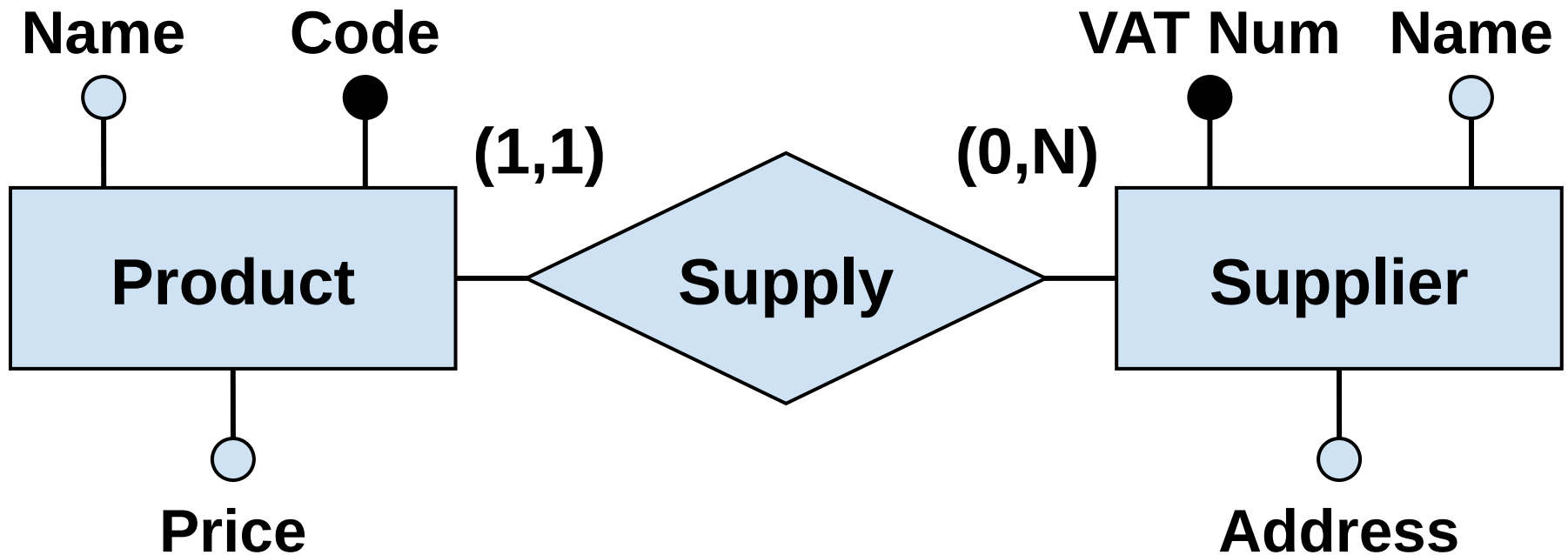
Check

- It violates the normal form due to the dependency:

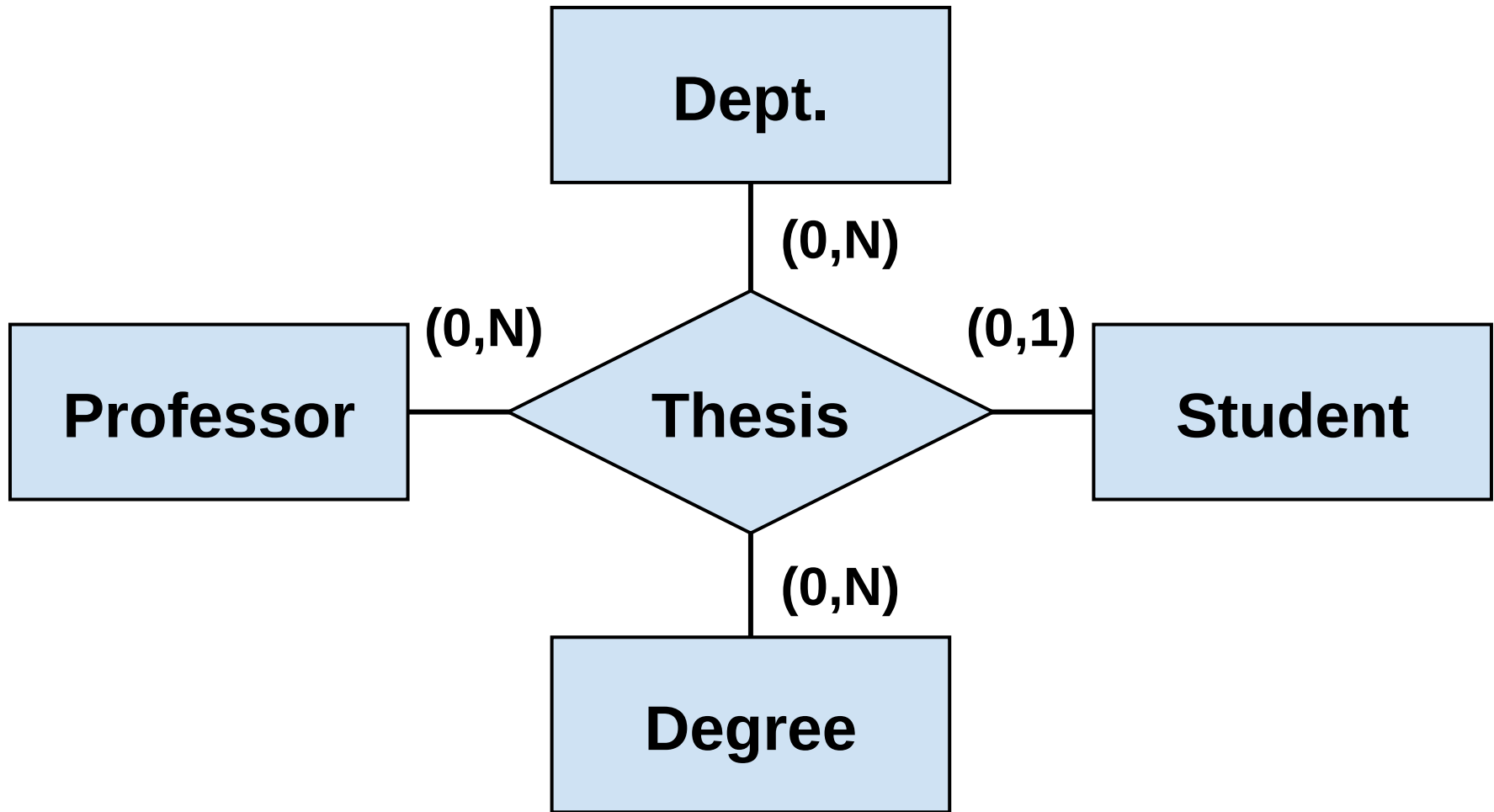
VAT Num → RetailerName Address

- We can decompose the entity using this dependency

Entity Decomposition



Normalization over Relationships



Student \rightarrow Degree

Student \rightarrow Professor

Professor \rightarrow Dept.



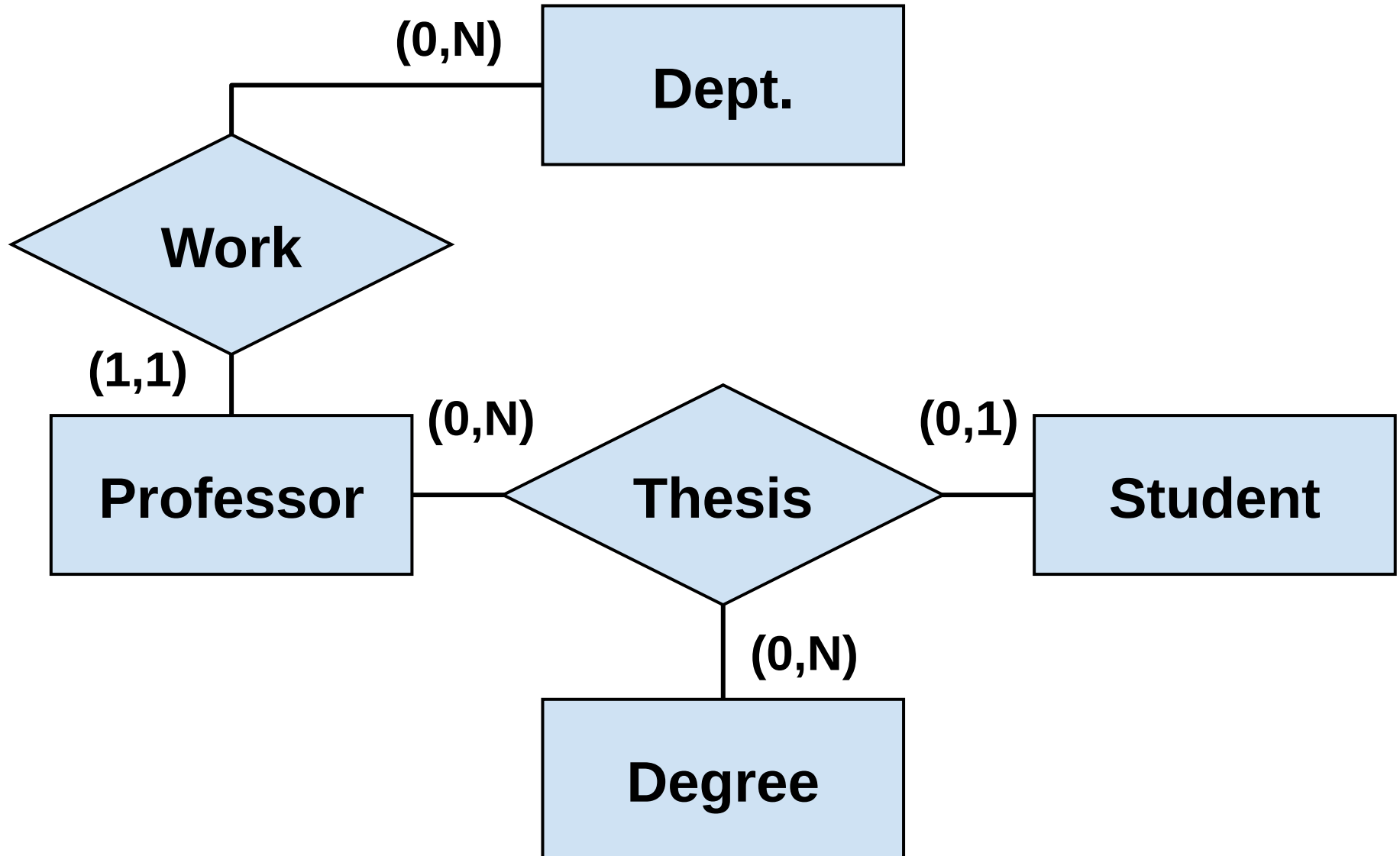
Checking the Relationship

- It has no 3NF due to the following dependency:

Professor \rightarrow Dept.

- We can decompose using this dependency

Relationship Decomposition





Yet another Dependency Analysis

- Thesis is in BCNF based on the dependencies:

Student \rightarrow Degree

Student \rightarrow Professor

- The two properties are independent
- We can perform a further decomposition

Relationship Decomposition (2)

