
Data Mining & Knowledge Discovery

Lesson: Data Normalization

Lan Man

Department of Computer Science and Technology

East China Normal University

©2017 All rights reserved.

Data Normalization

- A function that maps the entire set of values of a given attribute to a new set of replacement values s.t. each old value can be identified with one of the new values.
- **Normalization**: scaled to fall within a small, specified range
 - min-max normalization
 - z-score normalization
 - normalization by decimal scaling

Normalization

- **Min-max normalization:** to $[new_min_A, new_max_A]$

$$v' = \frac{v - min_A}{max_A - min_A} (new_max_A - new_min_A) + new_min_A$$

- Ex. Let *income* range \$12,000 to \$98,000 normalized to $[0.0, 1.0]$.

Then \$73,000 is mapped to

$$\frac{73,600 - 12,000}{98,000 - 12,000} (1.0 - 0) + 0 = 0.716$$

- **Z-score normalization (zero-mean)**

(μ : mean, σ : standard deviation):
$$v' = \frac{v - \mu_A}{\sigma_A}$$

- Ex. Let $\mu = 54,000$, $\sigma = 16,000$. Then
$$\frac{73,600 - 54,000}{16,000} = 1.225$$

- **Normalization by decimal scaling**

$$v' = \frac{v}{10^j} \quad \text{Where } j \text{ is the smallest integer such that } \text{Max}(|v'|) < 1$$