

Gaussian Naive Bayes

Evelyn Yosiana / 13522083

1. Konsep gaussian naive bayes pada intinya yaitu probabilitas (tepatnya likelihood probability). Model ini cocok digunakan untuk data numerik berdistribusi normal (sesuai dengan namanya, gaussian).

Proses train: untuk proses train, dicari mean, variansi, dan prior (probabilitas awal) untuk tiap label target

Proses predict: untuk proses ini, tiap data yang akan diprediksi (X_{test}) dihitung nilai pdf[label, fitur] (posterior probability). Misal terdapat 2 label (0 dan 1) dan 2 fitur (fitur_1 dan fitur_2), maka pdf yang dicari yaitu pdf0,fitur_1, pdf0,fitur_2, pdf1_fitur_1, dan pdf1,fitur_2. Persamaan untuk mencari pdf yaitu

$$pdf_{n,m} = \exp(- (x - mean_n)^2 / 2 var_y) / \sqrt{2 \cdot \pi \cdot var_y}$$

dengan n merupakan label yang diprediksi dan m merupakan banyak fitur. Setelah itu dicari log posterior (likelihood probability) dari tiap label dengan persamaan sebagai berikut.

$$\log posterior_n = \log(prior_n) + \log(pdf_{n,1}) + \dots + \log(pdf_{n,m})$$

dengan n merupakan label yang diprediksi dan m merupakan banyak fitur. Setelah didapatkan log posterior dari seluruh label, label dengan log posterior tertinggi akan menjadi label hasil. Nilai log posterior sendiri menunjukkan seberapa kuat suatu data berasal dari kelas tertentu. Proses tersebut diulangi untuk seluruh data di X_{test} .

Parameter:

- classes : seluruh nilai prediksi (distinct).
 - mean: rata-rata awal.
 - var: variansi awal.
 - priors: probabilitas awal.
4. Pada model ini, hasil perbandingan model yang saya buat dan model dari library kurang sama. Hal ini dapat disebabkan oleh beberapa faktor, antara lain:
 - **Penggunaan epsilon:** dalam kode yang saya buat, terdapat penambahan epsilon (angka yang sangat kecil) untuk menghindari error (misalnya log(0)) pada variansi awal, sedangkan kode dalam library sudah dioptimasi menggunakan laplace smoothing yang lebih komprehensif untuk perhitungan yang lebih presisi.

- **Perhitungan prior:** dalam kode saya, prior dihitung secara manual sedangkan dalam library, prior dihitung dengan cara yang sudah dioptimasi untuk data yang imbalanced.
4. Improvement yang dapat saya lakukan antara lain:
- **Hyperparameter tuning** untuk mendapatkan kombinasi parameter yang optimal (salah satunya dengan menggunakan metode grid search atau library optuna).
 - Menggunakan **confusion metrics** untuk model evaluation.

Lampiran

Contoh perhitungan manual:

Evelyn Yonana / 1522093

GAUSSIAN NAIVE BAYES

contoh data

$$x_{\text{train}} = \begin{bmatrix} 1 & 3 \\ 4 & 2 \\ 5 & 2 \end{bmatrix} \quad y_{\text{train}} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \quad \text{learning rate} = 0.01$$

$$\text{mean}_0 = \left[\frac{1+5}{2}, \frac{3+2}{2} \right] = [3, 2.5]$$

Rumus!

$$\text{var}_0 = \frac{(1-3)^2 + (5-3)^2}{2-1}, \frac{(3-2.5)^2 + (2-2.5)^2}{2-1} = [8, 0.5]$$

$$pdf_y(x) = \frac{\exp\left(\frac{-(x-\text{mean}_y)^2}{2 \cdot \text{var}_y}\right)}{\sqrt{2\pi \cdot \text{var}_y}}$$

$$\text{prior}_0 = \frac{2}{3} \approx 0.666$$

$$\text{mean}_1 = \left[\frac{4}{1}, \frac{2}{1} \right] = [4, 2]$$

$$\text{var}_1 = [0, 0] \quad \text{tapi di codenya ditambah } 1e^{-9} \text{ biar ga 0}$$

$$\text{prior}_1 = \frac{1}{3} \approx 0.333$$

prediksi [3,7] :

$$pdf_{0,1} = \frac{\exp\left(\frac{-(3-3)^2}{2 \cdot 8}\right)}{\sqrt{2\pi \cdot 8}} \approx 0.1413$$

$$pdf_{1,1} = \frac{\exp\left(\frac{-(3-4)^2}{2 \cdot 1e-9}\right)}{\sqrt{2\pi \cdot 1e-9}} \approx 2.3 \times 10^9$$

$$pdf_{0,2} = \frac{\exp\left(\frac{-(7-2.5)^2}{2 \cdot 0.5}\right)}{\sqrt{2\pi \cdot 0.5}} \approx 5.9$$

$$pdf_{1,2} = \frac{\exp\left(\frac{-(7-2)^2}{2 \cdot 1e-9}\right)}{\sqrt{2\pi \cdot 1e-9}} \approx 1.4 \times 10^{10}$$

$$\log \text{posterior}_0 = \log(\text{prior}_0) + \log(pdf_{0,1}) + \log(pdf_{0,2}) \\ \approx -23.8891$$

$$\log \text{posterior}_1 = \log(\text{prior}_1) + \log(pdf_{1,1}) + \log(pdf_{1,2}) \\ \approx 32.4609$$

$\log \text{posterior}_1 > \log \text{posterior}_0$, so prediction = 1 //