



Wrocław University
of Science and Technology

Modele generujące GDA, Naive Bayes

Szymon Zaręba

`szymon.zareba@pwr.edu.pl`

Modelowanie generujące

Rozkład warunkowy $p(y|\mathbf{x})$ można wyznaczyć korzystając ze wzoru Bayesa i twierdzenia o prawdopodobieństwie całkowitym

$$\begin{aligned} p(y|\mathbf{x}) &= \frac{p(\mathbf{x}|y)p(y)}{p(\mathbf{x})} \\ &= \frac{p(\mathbf{x}|y)p(y)}{\sum_{y'} p(\mathbf{x}|y')p(y')} \end{aligned}$$

Wielkości, które będą modelowane to $p(\mathbf{x}|y, \Theta)$ oraz $p(y|\Theta)$.

Podójście generujące: Naive Bayes

Rozkład $p(y|\Theta)$ będzie rozkładem wielopunktowym:

$$p(y|\Theta) = M(y|\pi) = \prod_{k=1}^K \pi_k^{y^{(k)}}$$

Model Naive Bayes zakłada niezaleźność cech:

$$p(\mathbf{x}|y, \Theta) = \prod_{d=1}^D p(x^{(d)}|y, \Theta)$$

Podejście generujące: Naive Bayes

Przypadek dyskretny (1)

Każda z cech modelowana jest rozkładem dwupunktowym:

$$p(x^{(d)}|y = k, \Theta) = B(x^{(d)}|\theta_{k,d}) = \theta_{k,d}^{x^{(d)}} (1 - \theta_{k,d})^{1-x^{(d)}}$$

Rozkład $p(\mathbf{x}|y = k, \Theta)$ dla wektora x przyjmuje postać:

$$p(\mathbf{x}|y = k, \Theta) = \prod_{d=1}^D B(x^{(d)}|\theta_{k,d})$$

Podejście generujące: Naive Bayes

Przypadek dyskretny (2)

W ogólnym przypadku rozkład $p(\mathbf{x}|y, \Theta)$ przyjmuje postać:

$$p(\mathbf{x}|y, \Theta) = \prod_{k=1}^K \left[\prod_{d=1}^D \text{B}(x^{(d)}|\theta_{k,d}) \right]^{y^{(k)}}$$

Zatem rozkład łączny $p(\mathbf{x}, y|\Theta)$ można zapisać jako:

$$p(\mathbf{x}, y|\Theta) = \prod_{k=1}^K \left[\pi_k \prod_{d=1}^D \text{B}(x^{(d)}|\theta_{k,d}) \right]^{y^{(k)}}$$

Podejście generujące: Naive Bayes

Przypadek ciągły (1)

Każda z cech modelowana jest rozkładem normalnym:

$$p(x^{(d)}|y = k, \Theta) = \mathcal{N}(x^{(d)}|\mu_{k,d}, \sigma_{k,d}^2)$$

Rozkład $p(\mathbf{x}|y = k, \Theta)$ dla wektora x przyjmuje postać:

$$p(\mathbf{x}|y = k, \Theta) = \prod_{d=1}^D \mathcal{N}(x^{(d)}|\mu_{k,d}, \sigma_{k,d}^2) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

$$\boldsymbol{\mu}_k = \begin{bmatrix} \mu_{k,1} \\ \cdots \\ \mu_{k,D} \end{bmatrix}$$

$$\boldsymbol{\Sigma}_k = \text{diag}(\sigma_k^2) = \begin{bmatrix} \sigma_{k,1}^2 & \cdots & 0 \\ \cdots & \cdots & \cdots \\ 0 & \cdots & \sigma_{k,D}^2 \end{bmatrix}$$

Podejście generujące: Naive Bayes

Przypadek ciągły (2)

W ogólnym przypadku rozkład $p(\mathbf{x}|y, \Theta)$ przyjmuje postać:

$$p(\mathbf{x}|y, \Theta) = \prod_{k=1}^K \left[\prod_{d=1}^D \mathcal{N}(x^{(d)} | \mu_{k,d}, \sigma_{k,d}^2) \right]^{y^{(k)}}$$

Zatem rozkład łączny $p(\mathbf{x}, y | \Theta)$ można zapisać jako:

$$p(\mathbf{x}, y | \Theta) = \prod_{k=1}^K \left[\pi_k \prod_{d=1}^D \mathcal{N}(x^{(d)} | \mu_{k,d}, \sigma_{k,d}^2) \right]^{y^{(k)}}$$

Podójście generujące: GDA

Można zadać następujące rozkłady $p(y|\Theta)$ i $p(\mathbf{x}|y = k, \Theta)$:

$$p(y|\Theta) = M(y|\boldsymbol{\pi}) = \prod_{k=1}^K \pi_k^{y^{(k)}}$$

$$p(\mathbf{x}|y = k, \Theta) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

W ogólnym przypadku rozkład $p(\mathbf{x}|y, \Theta)$ przyjmuje postać:

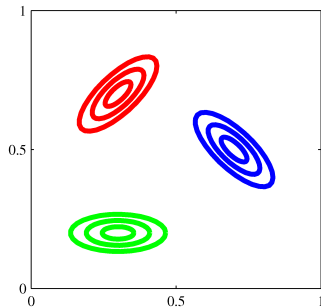
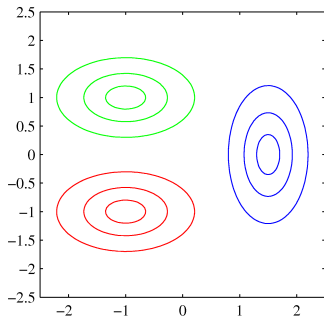
$$p(\mathbf{x}|y, \Theta) = \prod_{k=1}^K [\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]^{y^{(k)}}$$

Zatem rozkład łączny $p(\mathbf{x}, y|\Theta)$ można zapisać jako:

$$p(\mathbf{x}, y|\Theta) = \prod_{k=1}^K [\pi_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]^{y^{(k)}}$$

Zestawienie (1)

Model	$p(x, y)$	Θ	$\text{card}(\Theta)$
NB-d	$\prod_{k=1}^K \left[\pi_k \prod_{d=1}^D \text{B}(x^{(d)} \theta_{d,k}) \right] y^{(k)}$	$\pi_1, \dots, \pi_K,$ $\theta_{1,1}, \dots, \theta_{D,K}$	K KD
NB-c	$\prod_{k=1}^K \left[\pi_k \prod_{d=1}^D \mathcal{N}(x^{(d)} \mu_{k,d}, \sigma_{k,d}^2) \right] y^{(k)}$	$\pi_1, \dots, \pi_K,$ $\mu_{1,1}, \dots, \mu_{D,K},$ $\sigma_{1,1}^2, \dots, \sigma_{D,K}^2$	K KD KD
GDA	$\prod_{k=1}^K [\pi_k \mathcal{N}(\mathbf{x} \mu_k, \Sigma_k)] y^{(k)}$	$\pi_1, \dots, \pi_K,$ $\mu_1, \dots, \mu_K,$ $\Sigma_1, \dots, \Sigma_K$	K KD $K \frac{D(D+1)}{2}$



Uczenie

GDA

Dany jest rozkład $p(\mathbf{x}|y, \Theta)$:

$$p(\mathbf{x}|y, \Theta) = \prod_{k=1}^K [\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]^{y^{(k)}}$$

Można wyznaczyć funkcję wiarygodności:

Uczenie

GDA

Dany jest rozkład $p(\mathbf{x}|y, \Theta)$:

$$p(\mathbf{x}|y, \Theta) = \prod_{k=1}^K [\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]^{y^{(k)}}$$

Można wyznaczyć funkcję wiarygodności:

$$p(\mathcal{D}|\Theta) = \prod_{n=1}^N \prod_{k=1}^K [\mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]^{y_n^{(k)}}$$

Uczenie

GDA

Dany jest rozkład $p(\mathbf{x}|y, \Theta)$:

$$p(\mathbf{x}|y, \Theta) = \prod_{k=1}^K [\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]^{y^{(k)}}$$

Można wyznaczyć funkcję wiarygodności:

$$p(\mathcal{D}|\Theta) = \prod_{n=1}^N \prod_{k=1}^K [\mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]^{y_n^{(k)}}$$

po wstawieniu rozkładu normalnego:

Uczenie

GDA

Dany jest rozkład $p(\mathbf{x}|y, \Theta)$:

$$p(\mathbf{x}|y, \Theta) = \prod_{k=1}^K [\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]^{y^{(k)}}$$

Można wyznaczyć funkcję wiarygodności:

$$p(\mathcal{D}|\Theta) = \prod_{n=1}^N \prod_{k=1}^K [\mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)]^{y_n^{(k)}}$$

po wstawieniu rozkładu normalnego:

$$p(\mathcal{D}|\Theta) = \prod_{n=1}^N \prod_{k=1}^K \left[\frac{1}{(2\pi)^{D/2} |\boldsymbol{\Sigma}_k|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) \right\} \right]^{y_n^{(k)}}$$

Uczenie

GDA

Można wyznaczyć logarytm funkcji wiarygodności:

Uczenie

GDA

Można wyznaczyć logarytm funkcji wiarygodności:

$$\ln p(\mathcal{D}|\Theta) = \sum_{n=1}^N \sum_{k=1}^K y_n^{(k)} \left[-\frac{D}{2} \ln(2\pi) - \frac{1}{2} \ln |\Sigma_k| - \frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) \right]$$

Uczenie

GDA

Można wyznaczyć logarytm funkcji wiarygodności:

$$\ln p(\mathcal{D}|\Theta) = \sum_{n=1}^N \sum_{k=1}^K y_n^{(k)} \left[-\frac{D}{2} \ln(2\pi) - \frac{1}{2} \ln |\Sigma_k| - \frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) \right]$$

wyznaczyć jego gradient ze względu na parametr $\boldsymbol{\mu}_j$ i Σ_j :

Uczenie

GDA

Można wyznaczyć logarytm funkcji wiarygodności:

$$\ln p(\mathcal{D}|\Theta) = \sum_{n=1}^N \sum_{k=1}^K y_n^{(k)} \left[-\frac{D}{2} \ln(2\pi) - \frac{1}{2} \ln |\Sigma_k| - \frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) \right]$$

wyznaczyć jego gradient ze względu na parametr $\boldsymbol{\mu}_j$ i Σ_j :

$$\nabla_{\boldsymbol{\mu}_j} \ln p(\mathcal{D}|\Theta) = \sum_{n=1}^N y_n^{(j)} \Sigma_j^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_j)$$

$$\nabla_{\Sigma_j} \ln p(\mathcal{D}|\Theta) = \sum_{n=1}^N y_n^{(j)} \left[-\frac{1}{2} \Sigma_j^{-1} + \frac{1}{2} \Sigma_j^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_j) (\mathbf{x}_n - \boldsymbol{\mu}_j)^T \Sigma_j^{-1} \right]$$

Uczenie

GDA

Dla parametru μ_j :

$$\nabla_{\mu_j} \ln p(\mathcal{D}|\Theta) = \sum_{n=1}^N y_n^{(j)} \Sigma_j^{-1} (\mathbf{x}_n - \mu_j)$$

Uczenie

GDA

Dla parametru μ_j :

$$\nabla_{\mu_j} \ln p(\mathcal{D}|\Theta) = \sum_{n=1}^N y_n^{(j)} \Sigma_j^{-1} (\mathbf{x}_n - \mu_j)$$

$$\mu_j^{\text{ML}} = \frac{1}{N_j} \sum_{n=1}^N y_n^{(j)} \mathbf{x}_n$$

Uczenie

GDA

Dla parametru μ_j :

$$\nabla_{\mu_j} \ln p(\mathcal{D}|\Theta) = \sum_{n=1}^N y_n^{(j)} \Sigma_j^{-1} (\mathbf{x}_n - \mu_j)$$

$$\mu_j^{\text{ML}} = \frac{1}{N_j} \sum_{n=1}^N y_n^{(j)} \mathbf{x}_n$$

Dla parametru Σ_j :

$$\nabla_{\Sigma_j} \ln p(\mathcal{D}|\Theta) = \sum_{n=1}^N y_n^{(j)} \left[-\frac{1}{2} \Sigma_j^{-1} + \frac{1}{2} \Sigma_j^{-1} (\mathbf{x}_n - \mu_j) (\mathbf{x}_n - \mu_j)^T \Sigma_j^{-1} \right]$$

Uczenie

GDA

Dla parametru μ_j :

$$\nabla_{\mu_j} \ln p(\mathcal{D}|\Theta) = \sum_{n=1}^N y_n^{(j)} \Sigma_j^{-1} (\mathbf{x}_n - \mu_j)$$

$$\mu_j^{\text{ML}} = \frac{1}{N_j} \sum_{n=1}^N y_n^{(j)} \mathbf{x}_n$$

Dla parametru Σ_j :

$$\nabla_{\Sigma_j} \ln p(\mathcal{D}|\Theta) = \sum_{n=1}^N y_n^{(j)} \left[-\frac{1}{2} \Sigma_j^{-1} + \frac{1}{2} \Sigma_j^{-1} (\mathbf{x}_n - \mu_j) (\mathbf{x}_n - \mu_j)^T \Sigma_j^{-1} \right]$$

$$\Sigma_j^{\text{ML}} = \frac{1}{N_j} \sum_{n=1}^N (\mathbf{x}_n - \mu_j) (\mathbf{x}_n - \mu_j)^T$$

Zestawienie (2)

Model	Estymatory
NB-d	$\pi_j = \frac{N_j}{N}$ $\theta_{j,d} = \frac{1}{N_j} \sum_{n=1}^N y_n^{(j)} x_n^{(d)}$
NB-c	$\pi_j = \frac{N_j}{N}$ $\mu_{j,d} = \frac{1}{N_j} \sum_{n=1}^N y_n^{(j)} x_n^{(d)}$ $\sigma_{j,d} = \frac{1}{N_j} \sum_{n=1}^N \left(x_n^{(d)} - \mu_{j,d} \right)^2$
GDA	$\pi_j = \frac{N_j}{N}$ $\mu_j = \frac{1}{N_j} \sum_{n=1}^N y_n^{(j)} \mathbf{x}_n$ $\Sigma_j = \frac{1}{N_j} \sum_{n=1}^N \left(\mathbf{x}_n - \mu_j \right) \left(\mathbf{x}_n - \mu_j \right)^T$

Pre dykcja

GDA

Wykorzystując wzór Bayesa i twierdzenie o prawdopodobieństwie całkowitym:

$$\begin{aligned} p(y = j|\mathbf{x}) &= \frac{p(\mathbf{x}, y = j)}{p(\mathbf{x})} = \frac{p(\mathbf{x}|y = j)p(y = j)}{\sum_{k=1}^K p(\mathbf{x}, y = k)} \\ &= \frac{p(\mathbf{x}|y = j)p(y = j)}{\sum_{k=1}^K p(\mathbf{x}|y = k)p(y = k)} \end{aligned}$$

Dla modelu GDA:

$$p(y = j|\mathbf{x}) = \frac{\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)\pi_j}{\sum_{k=1}^K \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)\pi_k}$$

Zestawienie (3)

Model	Predykcja
NB-d	$\frac{\pi_j \prod_{d=1}^D \theta_{j,d}^{x^{(d)}} (1 - \theta_{j,d})^{1-x^{(d)}}}{\sum_{k=1}^K \pi_k \prod_{d=1}^D \theta_{k,d}^{x^{(d)}} (1 - \theta_{k,d})^{1-x^{(d)}}}$
NB-c	$\frac{\pi_j \prod_{d=1}^D \mathcal{N}(x^{(d)} \mu_{j,d}, \sigma_{j,d}^2)}{\sum_{k=1}^K \pi_k \prod_{d=1}^D \mathcal{N}(x^{(d)} \mu_{k,d}, \sigma_{k,d}^2)}$
GDA	$\frac{\pi_j \mathcal{N}(\mathbf{x} \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}$

Predykcja

Teoria decyzji mówi, że optymalnym wyborem jest y , który minimalizuje ryzyko (średnią stratę) przy zadanej funkcji straty.

Mając:

- $p(y|x, \Theta)$ - model
- $y \in \{1, \dots, K\}$ - zbiór możliwych decyzji
- $L(y, \hat{y}) = \begin{cases} 1, & y \neq \hat{y} \\ 0, & y = \hat{y} \end{cases}$ - funkcja straty

Predykcja

Teoria decyzji mówi, że optymalnym wyborem jest y , który minimalizuje ryzyko (średnią stratę) przy zadanej funkcji straty.

Mając:

- $p(y|x, \Theta)$ - model
- $y \in \{1, \dots, K\}$ - zbiór możliwych decyzji
- $L(y, \hat{y}) = \begin{cases} 1, & y \neq \hat{y} \\ 0, & y = \hat{y} \end{cases}$ - funkcja straty

$$R(\hat{y}) = \mathbb{E}[L(y, \hat{y})]$$

Predykcja

Teoria decyzji mówi, że optymalnym wyborem jest y , który minimalizuje ryzyko (średnią stratę) przy zadanej funkcji straty. Mając:

- $p(y|x, \Theta)$ - model
- $y \in \{1, \dots, K\}$ - zbiór możliwych decyzji
- $L(y, \hat{y}) = \begin{cases} 1, & y \neq \hat{y} \\ 0, & y = \hat{y} \end{cases}$ - funkcja straty

$$R(\hat{y}) = \mathbb{E}[L(y, \hat{y})] = \sum_{k=1}^K L(y = k, \hat{y})p(y = k|\mathbf{x})$$

Predykcja

Teoria decyzji mówi, że optymalnym wyborem jest y , który minimalizuje ryzyko (średnią stratę) przy zadanej funkcji straty. Mając:

- $p(y|x, \Theta)$ - model
- $y \in \{1, \dots, K\}$ - zbiór możliwych decyzji
- $L(y, \hat{y}) = \begin{cases} 1, & y \neq \hat{y} \\ 0, & y = \hat{y} \end{cases}$ - funkcja straty

$$\begin{aligned} R(\hat{y}) &= \mathbb{E}[L(y, \hat{y})] = \sum_{k=1}^K L(y = k, \hat{y})p(y = k|\mathbf{x}) \\ &= L(y = 1, \hat{y})p(y = 1|\mathbf{x}) + \dots + L(y = K, \hat{y})p(y = K|\mathbf{x}) \end{aligned}$$

Predykcja

Teoria decyzji mówi, że optymalnym wyborem jest y , który minimalizuje ryzyko (średnią stratę) przy zadanej funkcji straty. Mając:

- $p(y|x, \Theta)$ - model
- $y \in \{1, \dots, K\}$ - zbiór możliwych decyzji
- $L(y, \hat{y}) = \begin{cases} 1, & y \neq \hat{y} \\ 0, & y = \hat{y} \end{cases}$ - funkcja straty

$$\begin{aligned} R(\hat{y}) &= \mathbb{E}[L(y, \hat{y})] = \sum_{k=1}^K L(y = k, \hat{y})p(y = k|\mathbf{x}) \\ &= L(y = 1, \hat{y})p(y = 1|\mathbf{x}) + \dots + L(y = K, \hat{y})p(y = K|\mathbf{x}) \end{aligned}$$

$$\hat{y} = \operatorname{argmax}_k p(y = k|\mathbf{x})$$