

Bayes decision rule

Victor Kitov

v.v.kitov@yandex.ru

Table of Contents

- 1 Minimum cost and maximum probability solutions
- 2 Gaussian classifier
- 3 Naive Bayes assumption
- 4 Model examples with naive Bayes assumption

Costs

Classification

- supervised learning
- $y \in \{1, 2, \dots, C\}$ takes finite discrete set of values
- λ_{yf} is the cost of predicting true class y with forecasted class f .
- Examples with costs: diagnosis prediction, fraud detection, spam filtering, intrusion detection.

Costs

- Matrix of outcomes:

	$f = 1$	$f = 2$	\dots	$f = C$
$y = 1$	λ_{11}	λ_{12}	\dots	λ_{1C}
$y = 2$	λ_{21}	λ_{22}	\dots	λ_{2C}
\dots	\dots	\dots	\dots	\dots
$y = C$	λ_{C1}	λ_{C2}	\dots	λ_{CC}

- Expected cost of solution $\hat{y}(x) = f$:

$$\mathcal{L}(f) = \sum_y p(y|x) \lambda_{yf}$$

Decision rule

- Which best prediction $\hat{y}(x)$ for object x to select?

Decision rule

- Which best prediction $\hat{y}(x)$ for object x to select?

Bayes minimum risk decision rule

Assign class, yielding minimum expected cost:

$$\hat{y}(x) = \arg \min_f \mathcal{L}(f) \quad (1)$$

Decision rule

- Which best prediction $\hat{y}(x)$ for object x to select?

Bayes minimum risk decision rule

Assign class, yielding minimum expected cost:

$$\hat{y}(x) = \arg \min_f \mathcal{L}(f) \quad (1)$$

- This rule minimizes expected cost among all rules (if $p(y|x)$ are correct).

Simplifications

- $\lambda_{yf} \equiv \lambda_y \mathbb{I}[y \neq f]$: constant within class cost of misclassification.

Simplifications

- $\lambda_{yf} \equiv \lambda_y \mathbb{I}[y \neq f]$: constant within class cost of misclassification.

Matrix of outcomes:

	$f = 1$	$f = 2$	\dots	$f = C$
$y = 1$	0	λ_1	\dots	λ_1
$y = 2$	λ_2	0	\dots	λ_2
\dots	\dots	\dots	\dots	\dots
$y = C$	λ_C	λ_C	\dots	0

Simplifications

- $\lambda_{yf} \equiv \lambda_y \mathbb{I}[y \neq f]$: constant within class cost of misclassification.

Matrix of outcomes:

	$f = 1$	$f = 2$	\dots	$f = C$
$y = 1$	0	λ_1	\dots	λ_1
$y = 2$	λ_2	0	\dots	λ_2
\dots	\dots	\dots	\dots	\dots
$y = C$	λ_C	λ_C	\dots	0

- Expected cost of solution $\hat{y}(x) = f$:

$$\mathcal{L}(f) = \sum_y p(y|x) \lambda_y \mathbb{I}[f \neq y]$$

Equal misclassification costs

- Then cost of prediction equals:

$$\mathcal{L}(f) = \sum_y p(y|x)\lambda_y \mathbb{I}[f \neq y] = \overbrace{\sum_y p(y|x)\lambda_y}^{\text{const}(f)} - p(f|x)\lambda_f$$

- So (1) becomes:

$$\hat{y}(x) = \arg \min_f \mathcal{L}(f) = \arg \max_f \lambda_f p(f|x) \quad (2)$$

- Suppose further $\lambda_y \equiv \lambda \forall y$, then

$$\hat{y}(x) = \arg \max_f p(f|x)$$

- This is termed **maximum posterior probability rule** or **Bayes minimum error rule** because it yields minimum probability of misclassification among all decision rules (given that $p(f|x)$ is correct)

Equal misclassification costs

- This rule minimizes expected error rate.
 - if $p(y|x)$ are known

Equal misclassification costs

- This rule minimizes expected error rate.
 - if $p(y|x)$ are known
- If x and y are independent, then (2) reduces to

$$\hat{y}(x) = \arg \max_f p(f|x) = \arg \max_f p(f)$$

Table of Contents

- 1 Minimum cost and maximum probability solutions
- 2 Gaussian classifier
- 3 Naive Bayes assumption
- 4 Model examples with naive Bayes assumption

Gaussian classifier

- In Gaussian classifier

$$p(x|y) = \frac{1}{(2\pi)^{D/2} |\Sigma_y|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu_y)^T \Sigma_y^{-1} (x - \mu_y) \right\}$$

Gaussian classifier

- In Gaussian classifier

$$p(x|y) = \frac{1}{(2\pi)^{D/2} |\Sigma_y|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu_y)^T \Sigma_y^{-1} (x - \mu_y) \right\}$$

- It follows that

$$\begin{aligned} \log p(y|x) &= \log p(x|y) + \log p(y) - \log p(x) \\ &= -\frac{1}{2} (x - \mu_y^T) \Sigma_y^{-1} (x - \mu_y) - \frac{1}{2} \log |\Sigma_y| \\ &\quad - \frac{D}{2} \log(2\pi) + \log p(y) - \log p(x) \end{aligned}$$

Gaussian classifier

- In Gaussian classifier

$$p(x|y) = \frac{1}{(2\pi)^{D/2} |\Sigma_y|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu_y)^T \Sigma_y^{-1} (x - \mu_y) \right\}$$

- It follows that

$$\begin{aligned} \log p(y|x) &= \log p(x|y) + \log p(y) - \log p(x) \\ &= -\frac{1}{2} (x - \mu_y)^T \Sigma_y^{-1} (x - \mu_y) - \frac{1}{2} \log |\Sigma_y| \\ &\quad - \frac{D}{2} \log(2\pi) + \log p(y) - \log p(x) \end{aligned}$$

- Removing common additive terms, we obtain discriminant functions:

$$g_y(x) = \log p(y) - \frac{1}{2} \log |\Sigma_y| - \frac{1}{2} (x - \mu_y)^T \Sigma_y^{-1} (x - \mu_y) \quad (3)$$

Practical application

- In practice we replace theoretical terms μ_y , Σ_y with their sample estimates $\hat{\mu}_y$, $\hat{\Sigma}_y$.
- $\hat{p}(y) = \frac{N_y}{N}$.

$$g_y(x) = \log \hat{p}(y) - \frac{1}{2} \log |\hat{\Sigma}_y| - \frac{1}{2} (x - \hat{\mu}_y)^T \hat{\Sigma}_y^{-1} (x - \hat{\mu}_y)$$

- Analysis:
 - depends on normality assumptions (in particular - on unimodality)
 - needs to specify:
 - CD parameters to estimate $\hat{\mu}_y$, $y = 1, 2, \dots, C$.
 - $CD(D + 1)/2$ parameters to estimate $\hat{\Sigma}_y$, $j = 1, 2, \dots, C$.

Simplifying assumptions

- $CD(D + 3)/2$ may be too large for multidimensional tasks with small training sets.
- Simplifying assumptions:
 - **Naive Bayes**: assume that $\Sigma_1, \Sigma_2, \dots, \Sigma_C$ are diagonal.
 - **Project data onto a subspace**: for example on first few principal components.
 - **Proportional covariance matrices**: assume that $\Sigma_1 = \alpha_1 \Sigma, \Sigma_2 = \alpha_2 \Sigma, \dots, \Sigma_C = \alpha_C \Sigma$.
 - **Fisher's linear discriminant analysis**: assume that $\Sigma_1 = \Sigma_2 = \dots = \Sigma_C$.

QDA vs. LDA

Gaussian classifier is called:

- *Quadratic discriminant analysis* (QDA) when $\Sigma_1, \Sigma_2, \dots, \Sigma_C$ are arbitrary.
 - class boundaries are quadratic¹
- *Linear discriminant analysis* (LDA) when $\Sigma_1 = \Sigma_2 = \dots = \Sigma_C$
 - class boundaries are linear²

¹prove this

²prove this

Table of Contents

- 1 Minimum cost and maximum probability solutions
- 2 Gaussian classifier
- 3 Naive Bayes assumption**
- 4 Model examples with naive Bayes assumption

High dimensional problem

$$p(x^1, x^2, \dots, x^D) = p(x^1)p(x^2|x^1)\dots p(x^D|x^1, x^2, \dots, x^{D-1})$$

Problem: exponential to D number of observations needed for estimation.

High dimensional problem

$$p(x^1, x^2, \dots, x^D) = p(x^1)p(x^2|x^1)\dots p(x^D|x^1, x^2, \dots, x^{D-1})$$

Problem: exponential to D number of observations needed for estimation.

Solution: make simplifying assumptions.

High dimensional problem

$$p(x^1, x^2, \dots, x^D) = p(x^1)p(x^2|x^1)\dots p(x^D|x^1, x^2, \dots, x^{D-1})$$

Problem: exponential to D number of observations needed for estimation.

Solution: make simplifying assumptions.

Independence assumption

Individual features are independent: $p(x) = p(x^1)p(x^2)\dots p(x^D)$

High dimensional problem

$$p(x^1, x^2, \dots, x^D) = p(x^1)p(x^2|x^1)\dots p(x^D|x^1, x^2, \dots, x^{D-1})$$

Problem: exponential to D number of observations needed for estimation.

Solution: make simplifying assumptions.

Independence assumption

Individual features are independent: $p(x) = p(x^1)p(x^2)\dots p(x^D)$

Naive Bayes assumption in classification

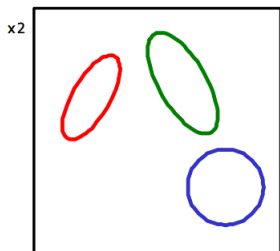
Individual features are **class conditionally** independent:

$$p(x|y) = p(x^1|y)p(x^2|y)\dots p(x^D|y)$$

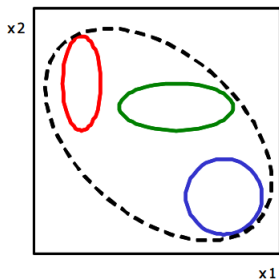
Under Naive Bayes assumption max-posterior probability rule becomes:

$$\hat{y}(x) = \arg \max_y p(y)p(x^1|y)p(x^2|y)\dots p(x^D|y)$$

Conditional independence visualization

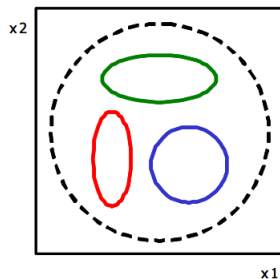


$$p(x|\omega_i) \neq \prod_{d=1}^D p(x_d|\omega_i)$$



$$p(x|\omega_i) = \prod_{d=1}^D p(x_d|\omega_i)$$

$$p(x) \neq \prod_{d=1}^D p(x_d)$$



$$p(x|\omega_i) = \prod_{d=1}^D p(x_d|\omega_i)$$

$$p(x) \cong \prod_{d=1}^D p(x_d)$$

Table of Contents

- 1 Minimum cost and maximum probability solutions
- 2 Gaussian classifier
- 3 Naive Bayes assumption
- 4 Model examples with naive Bayes assumption**

Text models

- Restrict attention to D words w_1, w_2, \dots, w_D
 - all unique words
 - possibly with stop words removal
 - possibly only most frequent words
 - or only words relevant to the topic of study
- Two major models:
 - Bernoulli
 - Multinomial

Bernoulli model⁵

- Document is represented with feature vector $x \in \mathbb{R}^D$
- $x^i = \mathbb{I}[w_i \text{ appeared in the document}]$
- $\theta_y^d = p(x^d = 1|y)$
- $p(x|y) = \prod_{d=1}^D (\theta_y^d)^{x^d} (1 - \theta_y^d)^{1-x^d}$
- $p(y) = \frac{N_y}{N}$
- $\theta_y^d = \frac{N_{yx^d}}{N_y}$
- Smoothed variant^{3,4}: $\theta_y^d = \frac{N_{yx^d} + \alpha}{N_y + 2\alpha}$

³interpret this in terms of adding artificial observations

⁴modify for smoothing towards unconditional word distribution

⁵is it linear classifier?

Multinomial model⁸

- Document is represented with feature vector $x \in \mathbb{R}^D$
- x^d = number of times w_d appeared in the document
- θ_y^d = probability of w_d on word position
- $p(x|y) = \frac{(\sum_d x^d)!}{\prod_d (x^d)!} \prod_{d=1}^D (\theta_y^d)^{x^d}$
- $p(y) = \frac{N_y}{N}$
- $\theta_y^d = \frac{n_{yd}}{n_y}$ where
 - n_{yd} - number of times word w_d appeared in documents $\in y$
 - n_y - number of words in documents $\in y$
- Smoothed version^{6,7}: $\theta_y^d = \frac{n_{yd} + \alpha}{n_y + \alpha D}$

⁶interpret this in terms of adding artificial observations

⁷modify for smoothing towards unconditional word distribution

⁸is it linear classifier?