

Notation

\mathcal{D} : data $\{(x_i, y_i)\}_{i=1}^n$

\mathcal{D}_i : data sample (x_i, y_i)

\cancel{x} dimension d

θ : parameter of Bernoulli, $Ber(y|\theta)$

θ : hypothesis for MLE $h \in H$ in Bernoulli: $\theta \in [0, 1]$

X, x_i : design matrix, features of example 'i'

y, y_i : target vector, label of example 'i'

Naive Bayes

Issues ✓ model makes no sense

- div by 0

- MLE?

Advantages

- simple

- fast train

- parallel

- small # parameters

- not much data needed

Disadvantages

- strong mutual dependence assumption

↳ limited modelling power

only certain data sets

Project idea

- explore ways to relax NB assumption
(TAN Bayes, Bayesian Network classifiers)
free augmented

Today

What if $x_i \in \mathbb{R}^p$ real valued features instead of binary

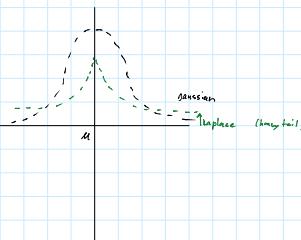
$$\overline{x}_i \in \mathbb{R}^d$$

Gaussian Distribution

distributed
normal

$$x \sim N(\mu, \sigma^2)$$

$$p(x | \mu, \sigma^2) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

Motivation

- central limit theorem (converge to mean)
- analytic properties (smooth, differentiable, symmetric)
- defined $(-\infty, \infty)$
- easy to sample
- data is Gaussian

MLE

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N x_i \quad \sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \hat{\mu})^2 \quad \Rightarrow p(x^i | y_i)$$

$$\bar{x} \sim N(\mu, \Sigma) \quad p(\bar{x} | \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(\bar{x} - \mu)^T \Sigma^{-1} (\bar{x} - \mu)\right)$$

$\mu \in \mathbb{R}^d$ (anything)
symmetric

if $\Sigma = I$, you get ellipses
it Σ is I_d similar to regular naive bayes,
(however it does not have to be.)

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N \bar{x}_i$$

$$\Sigma = \frac{1}{N} \sum_{i=1}^N (\bar{x}_i - \mu)(\bar{x}_i - \mu)^T + I$$

outer product of 2 vectors is a rank 1 matrix

$\Sigma \in \mathbb{S}^d_+$
positive definite

$\Sigma \geq 0$
strictly greater than zero

$\int_{-\infty}^{\infty} p(x | \mu, \Sigma) = 1$

$\Sigma \in \mathbb{S}^d_+$
positive semi-definite

(some reason why no 2+0)

Problems

1. $\hat{\Sigma}$ might not be valid

2. unimodal



3. Not robust (far away from mean)



Gaussian Discriminant Analysis

$$p(\bar{x}_i, y_i) = p(\bar{x}_i | y_i) p(y_i)$$

$\bar{x}_i \in \mathbb{R}^d$

$y_i \in \{1, 2, 3, \dots\}$

$$(\bar{x}_i | y_i) \sim N(\mu_c, \Sigma_c)$$

$$\begin{matrix} \mu_1 & \Sigma_1 \\ \mu_2 & \Sigma_2 \\ \mu_3 & \Sigma_3 \\ \vdots & \vdots \end{matrix}$$

$$\begin{matrix} \bar{x}_1 \\ \bar{x}_2 \\ \bar{x}_3 \\ \vdots \end{matrix} \sim \Sigma_0$$

$$\begin{matrix} \bar{x}_1 \\ \bar{x}_2 \\ \bar{x}_3 \\ \vdots \end{matrix} \sim N(\mu_0, \Sigma_0)$$

know class, we are in some gaussian distribution

$$p(y_i | \bar{x}_i) = \alpha p(\bar{x}_i | y_i)$$

we don't care about $p(\bar{x}_i)$, does not depend on y , so we can already make a decision.

$$p(\bar{x}_i) = \sum_{c=1}^C p(\bar{x}_i | y_i = c) p(y_i = c)$$

$$p(\bar{x}_i, y_i = c)$$

parameters Σ

$$C + \frac{cd}{n} + \frac{C(C+1)d}{2} = O(Cd^2)$$

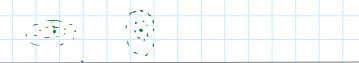
This is the sum of a matrix

that is generic: $\frac{d(d+1)}{2}$. There are C of the Σ terms.

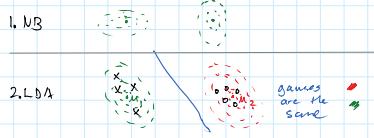
Here is also C of the n terms. That gives the big O given, upper bounded by sigma ($d(d+1) = d^2$)

1. Naive Bayes assumption: Σ diagonal isotropic (circles of the same size)
2. Linear discriminant analysis: $\Sigma_c = \Sigma$, $\forall c$ each Σ similar to same common sigma
3. Quad. Discr. anal.: general Σ_c

1. NB



2. LDA



James are the same

3. Quad



one class envelope another
(right class is a subset)

Gaussian Regression

both real

$$y \quad \quad \quad x = \begin{bmatrix} \dots \end{bmatrix} \quad y = \begin{bmatrix} \dots \end{bmatrix}$$

"linear model"

"Linear regression" $y_i = w^\top x_i + \varepsilon_i$

linear

$$y_i | \bar{x}_i, w \sim N(w^\top \bar{x}_i, 1)$$

$$w_{MLE} = \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} -\sum_{i=1}^n \log p(y_i | \bar{x}_i, w)$$

$$= \underset{w \in \mathbb{R}^d}{\operatorname{argmin}} \frac{1}{2} \sum_{i=1}^n (y_i - w^\top \bar{x}_i)^2 + \text{constant (doesn't depend, doesn't change argmax)}$$

"Least Squares"

convex problem (so set $\frac{\partial L}{\partial w}$ to 0 and solve)

$$X_w = \begin{bmatrix} \bar{w}^\top \bar{x}_1 \\ \bar{w}^\top \bar{x}_2 \\ \vdots \\ \bar{w}^\top \bar{x}_n \end{bmatrix}$$

(inner products, sum over i in the space \mathbb{R}^d)

negative log likelihood

NLL(w) =

$\underset{w \in \mathbb{R}^d}{\operatorname{argmin}}$

$\frac{1}{2} (Y - Xw)^\top (Y - Xw)$

scalar

related to $\frac{1}{2} \|Y - Xw\|^2$

$\frac{1}{2} \sum_i (Y_i - w^\top \bar{x}_i)^2$

$\frac{1}{2} \sum_i (Y_i - w^\top \bar{x}_i)^2$