



Lecture 5

- · a's on predictors
- · Parametric vs Non-parametric
- . Generative classifiers

- · Q1 on 10/19
- · LT posted
- · Hwa't. b. posted
 - HW1 due [now, 10/19]

1. QDA intuition 2-LDA QDA parameters 3. Decision Region DT $\hat{y}(x) = +1$ iff $f_{+}(x) > f_{-}(x)$ Dt $\mathcal{N}^{4} = \mathcal{N}^{-} = \frac{3}{\mathcal{N}}$ $\Lambda_{\mu_{-}}$, $\tilde{\Sigma}_{-}^{2} = \sigma_{-}^{2}$ $\hat{M}_{+}, \hat{\Sigma}_{+} = \hat{\sigma}_{+}^{2}$





3. DT decision regions



Lecture II: Prediction – Basic concepts

Marina Meilă mmp@stat.washington.edu

> Department of Statistics University of Washington

October, 2023



Generative and discriminative models for classification



Generative classifiers Discriminative classifiers Generative vs discriminative classifiers

Loss functions Bayes loss

Variance, bias and complexity

 $^{^{1}}$ Neither textbook is close to these notes except in a few places; take them as alternative perspectives or related reading

The "learning" problem

- Given
 a problem (e.g. recognize digits from m × m gray-scale images)
 - a sample or (training set) of labeled data
 - $\mathcal{D} = \{(x^1, y^1), (x^2, y^2), \dots, (x^n, y^n)\}$

drawn i.i.d. from an unknown P_{XY}

- model class $\mathcal{F} = \{f\}$ = set of predictors to choose from
 - Wanted $f \in \mathcal{F}$ that performs well on future samples from the same P_{XY}
 - "choose a predictor $f \in \mathcal{F}$ " = training/learning \leftarrow training alg
 - "performs well on future samples" (i.e. f generalizes well) how do we measure this? how can we "guarantee" it?

ec-

(later)

choosing *F* is the model selection problem – about this later

A zoo of predictors

- October, 2023
- Linear regression
- Logistic regression
- Linear Discriminant (LDA)
- Quadratic Discriminant (QDA)
- CART (Decision Trees)
- K-Nearest Neighbors
- Nadaraya-Watson (Kernel regression)
- Naive Bayes
- Neural networks/Deep learning
- Support Vector Machines
- Monotonic Regression

Classes of predictors - by output (classification regression - by dimension < finite P - by decision bdary" generative discriminative

Parametric vs. non-parametric models

Example (Parametric and non-parametric predictors)



CART
$$L = \#$$
 leaves $= \#$ splits ± 1
 $\cdot L$ fixed \Rightarrow "p" fixed $\propto L \Rightarrow J_L = 1$ CART with L leaves?
 $\cdot L$ grows with $m \Rightarrow F = 1$ all cherts
 $e.g \ L = n$ non-parametric
Para metric $-(typically)$ interpretable $Ex: LDA \ \mu \leq 5 \leq 1$
 $-[easier to analyze] -1/2$
 $Ex: \ N \sim N(\mu, \sigma^2)$
wanted $\mu : \ \hat{\mu} = \frac{1}{n} \sum_{i=1}^{n} E[\|\hat{\mu} - \mu\||^2] = \frac{\sigma^2}{m} \propto \frac{1}{n} = n^1$ Parametric
 $-typically)$ not interpretable $[\pi + \mu]^2$
 $-hander to analyze + 1000 fast when $\sum m^{1/2}$ metric $metric$
 $-hander to analyze + 1000 fast when $\sum m^{1/2}$ metric $metric$
 $+ adapts to shape of data + 1000 fast when $p = \frac{1}{n} \sum_{i=1}^{n} \frac{1}{n} = 10$$$$



A mathematical definition

A model class F is parametric if it is finite-dimensional, otherwise it is non-parametric

> F = vector space of f's

In other words

- When we estimate a parametric model from data, there is a fixed number of parameters, (you can think of them as one for each dimension, although this is not always true), that we need to estimate to obtain an estimate $\hat{f} \in \mathcal{F}$.
- The parameters are meaningful. E.g. the β_i in logistic regression has a precise meaning: the component of the normal to the decision boundary along coordinate *i*.
- The dimension of β does not change if the sample size *n* increases.

: Linear predictor Parametric parametric $\beta_0, \beta \in \mathbb{R}^d \implies p = d+1$ independent of n, Dp = # parameters Non para metric K-NN p= # params to discribe Do - grows with n - growth depends on a (HURDINGH PXY)

Non-parametric models - Some intuition

- \blacktriangleright When the model is non-parametric, the model class ${\cal F}$ is a function space.
- The f̂ that we estimate will depend on some numerical values (and we could call them parameters), but these values have little meaning taken individually.
- The number of values needed to describe \hat{f} generally grows with *n*. Examples In the Nearest neighbor and kernel predictors, we have to store all the data points, thus the number of values describing the predictor *f* grows (linearly) with the sample size. Exercise Does the number of values describing *f* always grow linearly with the sample size? Does it have to always grow to infinity? Does it have to always grow in the same way for a given \mathcal{F} ?
- Non-parametric models often have a smoothness parameter.

Examples of smoothness parameters K in K-nearest neighbor, h the kernel bandwidth in kernel regression.

To make matters worse, a smoothness parameter is not a parameter! More precisely it is not a parameter of an $f \in \mathcal{F}$, because it is not estimated from the data, but a descriptor of the model class \mathcal{F} .

• We will return to smoothness parameters later in this lecture.

Generative classifiers

One way to define a classifier is to assume that each class is generated by a distribution $g_y(X) = P(X|Y = y)$. If we know the distributions g_y and the class probabilities P(Y = y), we can derive the *posterior probability* distribution of Y for a given x. This is

$$P(Y = y|X) = \frac{P(Y = y)g_y(X)}{\sum_{y'} P(Y = y')g_{y'}(X)} = \frac{P(Y = y)g_y(X)}{P(X)}$$
(1)

The "best guess" for Y(X) (i.e. the decision rule) is

$$f(X) = \operatorname{argmax}_{y} P(Y = y | x) = \operatorname{argmax}_{y} P(Y = y) g_{y}(x)$$
(2)

(1) amounts to a likelihood ratio test for Y.

The functions g_y(x) are known as generative models for the classes y. Therefore, the resulting classifier is called a generative classifier. Examples: LDA, QDA, Naive Bayes.

- In contrast, a classifier defined directly in terms of f(x) (or P_{Y|X}), like the linear, quadratic, decision tree is called a discriminative classifier.
- ▶ In practice, we may not know the functions $g_y(x)$, in which case we estimate them from the sample \mathcal{D} .

Generative classifiers (Binary)
given
$$\vartheta = J(x_1^i y_1^i)$$
, $i=1:n_3^i$ $y_1^i eft 13$
1. fit $g_{+} \stackrel{e}{=} P_{X|Y=+}$ a generative models
 $g_{-} \stackrel{e}{=} P_{X|Y=-}$
t deuxity or PMF
Ex: Classify
objects by images $g_{car} = f(image | car)$
 $g_{person} = f(-u-| person)$
 $\vartheta = person = f(-u-| person)$

Generative classifier and the likelihood ratio

$$P(Y = y|X) = \frac{P(Y = y)g_y(X)}{\sum_{y'} P(Y = y')g_{y'}(X)} = \frac{P(Y = y)g_y(X)}{P(X)}$$

 $f(x) = \operatorname{argmax}_{y} P(Y = y | x) = \operatorname{argmax}_{y} g_{y}(x) P(Y = y)$

Likelihood Ratio test (for $y \in \{\pm 1\}$)

 $\frac{g_+(x)P(Y=+)}{g_-(x)P(Y=-)}$

October, 2023

Example (Fisher's LDA in one dimension)

Assume $Y = \pm 1$, $g_y(x) = N(x, \pm \mu, \sigma^2 I)$, i.e each class is generated by a Normal distribution with the same spherical covariance matrix, but with a different mean. Let $P(Y = 1) = p \in (0, 1)$. Then, the posterior probability of Y is

$$P(Y = 1|x) \propto p e^{-||x-\mu||^2/(2\sigma^2)} \quad P(Y = -1|x) \propto (1-p) e^{-||x+\mu||^2/(2\sigma^2)}$$
(3)

and f(x) = 1 iff $\ln P(Y = 1|x) / P(Y = -1|x) \ge 0$, i.e iff

$$\ln \frac{p}{1-p} - \frac{1}{2\sigma^2} [||x^2|| - 2\mu^T x + ||\mu||^2 - ||x^2|| - (2\mu)^T x - ||\mu||^2] = \left(\frac{2\mu}{\sigma^2}\right)^T x + \ln \frac{p}{1-p} \ge 0$$
(4)

Hence, the classifier f(x) turns out to be a linear classifier. The decision boundary is perpendicular to the segment connecting the centers μ , $-\mu$. This classifier is known as **Fisher's Linear Discriminant**. [Exercises Show that if the generative models are normal with different variances, then we obtain a quadratic classifier. What happens if the models g_y have the same variance, but it is a full covariance matrix Σ ?]