

5.2 Probabilistic Classifiers

Probabilistic = output is prob of a new object belonging to a particular class

train data in Matrix, vector = (X, t)

Probability for class c:

$$P(T_{\text{new}} = c | x_{\text{new}}, X, t)$$

5.2.1 ~~Naive Bayes~~ The Bayes classifier

5.2.2 Logistic Regression

5.3.2 Support vector machines

Bayes Classifier

- takes name from equation it is based on
- Aim is to compute predictive probs of classes

From Bayes Rule, expression for Predictive Prob:

$$P(T_{\text{new}} = c | X_{\text{new}}, X, f) =$$

Bayes Classifies

(2)

from Bayes Rules, Expression for pred prob:

$$P(T_{\text{new}} = c | X_{\text{new}}, X, t) = \frac{P(X_{\text{new}} | T_{\text{new}} = c, X, t) P(T_{\text{new}} = c | X, t)}{P(X_{\text{new}} | X, t)}$$

① $P(X_{\text{new}} | T_{\text{new}} = c, X, t)$

- the probability of observing the new feature vector X_{new} given that the new data point belongs to class c
 - + given training set X & parameters t
- the likelihood of observing the new data points feature if we assume it belongs to class c

② $P(T_{\text{new}} = c | X, t)$

- The Prior Probability that the new data point T_{new} belongs to class c
 - + given training set X & parameters t

③ $P(X_{\text{new}} | X, t) \leftarrow$ marginal likelihood

- Probability of observing the new feature vector X_{new} given *training set & parameters t

(3)

the marginal likelihood can be expanded to a sum over all possible C class

~~P(Xnew | X, t) = P(Tnew = c' | X, t) * P(Xnew | Tnew = c')~~

$$P(X_{\text{new}} | X, t) \rightarrow$$

$$\sum_{c'=1}^C P(X_{\text{new}} | T_{\text{new}} = c', X, t) P(T_{\text{new}} = c' | X, t)$$

where c' = All possible classes

Why this expansion work? law of total prob

this says the prob of an event occurring (X_{new}) can be calculated by summing the probs occurring under each condition, i.e., the classes

Data is assumed to be complete & mutually exclusive, i.e. all classes known and data only belongs to one class

(4)

This Expansion is the product of the likelihood & prior prob for each class

How to define/identify these new eqs?

$$p(X_{\text{new}} | T_{\text{new}} = c', X, \epsilon) = \text{likelihood of } X_{\text{new}} \text{ belong to Class } C$$

$$P(T_{\text{new}} = c | X, \epsilon) = \text{prior prob of class } C$$

Note the full Expanded looks like:

$$\frac{P(X_{\text{new}} | T_{\text{new}} = c, X, \epsilon) P(T_{\text{new}} = c | X, \epsilon)}{\sum_{c'=1}^C P(X_{\text{new}} | T_{\text{new}} = c', X, \epsilon) P(T_{\text{new}} = c' | X, \epsilon)}$$

This equation works by allowing us to calc the probability of a new data point belonging to a specific class by comparing the likelihood & prior for that class to the likelihood & prior for all classes

Same exp on top & bottom w/ addition of sum

(5)

How to calculate likelihood?

$$p(X_{\text{new}} \mid T_{\text{new}} = c, X, t)$$

This is a distribution specific to the c^{th} class ($T_{\text{new}} = c$) evaluated at X_{new}

To create a Bayes classifier, we need to

- Define C of these class condition distributions
- it is common to use same dist for ea class
- But no strict ruling to do this
- Dist should be based on data being modelled & knowledge we have about data
- once dist chose, select parameters
- e.g. gaussian, mean & co-variance

(6)

How to calc prior probability?

$$P(T_{\text{new}} = c \mid X, t)$$

This is the probability that an object belongs to class c on just the training set, (X, t)

This piece of info allows us to give prior beliefs + about the class of a new data instance X_{new} before we see it

if $P(T_n = c, X, t)$ is really low then there exists a prior bias against X_{new} being c

we would only classify it as c if the likelihood (update) is very high

the only restrictions on setting the probs is that they are positive & \sum to 1

(7)

two popular choices are:

① uniform prior: $P(T_{\text{new}} = c | X, t) = \frac{1}{C}$

② Class size prior: " $= \frac{N_c}{N}$

N = train set N_c = class num in train

Note neither dist does not need to be conditioned
on X & t . ① uses neither ② uses t via C

(8)

Example Gaussian Class-Conditional

Each train object consists of a two-Dim attribute vector $x_n = [x_{n1}, x_{n2}]^T$ & an assoc label $t_n = \{1, 2, 3\}$

We can create a dataset from these 3 classes using distributions

This example uses gaussian as attributes are real valued

$$P(x_n | t_n = c, X, t) = N(\mu_c, \Sigma_c)$$

μ & Σ need to be chosen based on train pairs of class c , know X^c

This is nooo a machine learning class.
we have X^c & we want to infer the parameters of the model: μ_c, Σ_c

find paras that Max likelihood of obs X^c

(9)

Alternatively could use a Bayesian approach

Define a prior density for these parameters, $p(\mu_c, \Sigma_c)$ to compute a posterior from Bayes rules

$$p(\mu_c, \Sigma_c | X^c) = \frac{p(X^c | \mu_c, \Sigma_c) p(\mu_c, \Sigma_c)}{p(X^c)}$$

calc dist of paras given
obsd X^c

then compute the likelihood of X_{new} by taking the expectation:

How likely our values of μ, Σ are after obs data

$$p(X_{\text{new}} | T_{\text{new}} = c, X, t) = E_p(\mu_c, \Sigma_c | X^c)$$

calc average value $\rightarrow \sum p(X_{\text{new}} | \mu_c, \Sigma_c)^2$ likelihood of X_{new} given μ, Σ
by their distrib weights

This equation calcs the likelihood of a new data point X_{new} belonging to class C by taking the Expectation of $p(X_n | \mu_c, \Sigma_c)$ with respect to the posterior distribution of μ_c, Σ_c

Expectation means calc the average of a function w/
respect to a distribution function

(10)

 μ Σ

Max likelihood estimates for mean μ & var Σ of a gaussian given N data points is obtained by Diff the log-likelihood w/ respect to each parameter, setting to 0 & solving

- ① likelihood = Prob of obs data given μ & Σ
- ② instead of work w/ likelihood, work w/ log-likelihood $\ln(L)$. Log is monotonic so preserves location of the max. Log is easier to work for summing probs
- ③ Differentiation - Derivative tells you how log-likelihood changes as you vary parameters
- ④ Set to 0 & Solve = Values that Max log-likelihood

notes:

- have function that tells you how well your data fits a gaussian dist for diff paras
- take derive to find paras that best fit the data

(11)

(MEE)

the Max likelihood estimates for
the mean & var become

$$\mu_c = \frac{1}{N_c} \sum_{n=1}^{N_c} X_n$$

$$\Sigma_c = \frac{1}{N_c} \sum_{n=1}^{N_c} (X_n - \mu_c)(X_n - \mu_c)^T$$

(12)

How to make predictions?

Once you have class-conditional distributions & the prior you can make predictions

Example: compute class probs for

$$X_{\text{new}} = [2, 0]^\top$$

$$\frac{P(T_{\text{new}}=c | X_{\text{new}}, X, t)}{\sum_{c'=1}^C P(T_{\text{new}}=c' | X_{\text{new}}, X, t)} = \frac{P(X_{\text{new}} | T_{\text{new}}=c, X, t) P(T_{\text{new}}=c | X, t)}{\sum_{c'=1}^C "c'"} \quad "c"$$

Various calcs needed:

	Likelihood	Prior	Product
c	$P(X_n T_n=c, \mu_c, \Sigma_c)$	$P(T_n=c X, t)$	
1	0.198	$\frac{1}{3}$	0.0046
2	0.0661	$\frac{1}{3}$	0.0020
3	0.0002	$\frac{1}{3}$	0.0001

use $\frac{N_c}{N}$

0.0067

$X_{\text{new}} [2, 0]$

this col summed
the denominator

$P(T=1) = 0.6890$

$P(T=2) = 0.3024$

$P(T=3) = 0.0087$