

Lecture 10 - Feb 12

Assessing Classifiers

Fisher's Linear Discriminant Analysis (LDA)

Reading

Week 6 Notes in GitHub

[*Data Mining and Machine Learning*](#)

1.3.3: Orthogonal Projection

20.1: Linear Discriminant Analysis

[Elements of Statistical Learning](#)

4.3 Discriminant Analysis

Upcoming Deadlines

Homework 2 (Feb 18)

Suppose $M: \mathbb{R}^d \rightarrow \{c_1, \dots, c_k\}$ is a classifier

To build a classifier requires a training set.
We assess its performance on a testing set.

22.1 Classification Performance Measures

D - testing set of n points $x_i \in \mathbb{R}^d$ with labels y_i

$\hat{y}_i = M(x_i)$ - predicted label in $\{c_1, \dots, c_k\}$

$$\mathbb{1}_A = \begin{cases} 0, & \text{if } A \text{ is false} \\ 1, & \text{if } A \text{ is true} \end{cases}$$

Error Rate: $\frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{y_i \neq \hat{y}_i\}}$ → fraction that is incorrectly classified

Accuracy: $\frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\{y_i = \hat{y}_i\}}$ → fraction correctly classified

These metrics are global and do not consider how different classes contribute to the error.

Let $D_j = \{x_i \mid y_i = c_j\}$ points in class j

* $n_j = |D_j| = \# \text{ of points in class } j$

$R_j = \{x_i \mid \hat{y}_i = c_j\}$ points with predicted class j

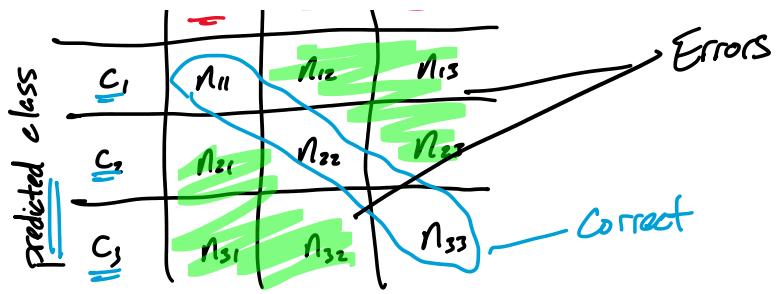
* $M_j = |R_j| = \# \text{ of points with predicted class } j$

A Confusion matrix is a $k \times k$ table N with

$N_{ij} = n_{ij}$ - number of points with predicted class c_i with true label c_j

$n_{ij} = |R_i \cap D_j| = \{x_a \in D \mid \hat{y}_a = c_i \text{ and } y_a = c_j\}$

		<u>true class</u>		
		<u>c_1</u>	<u>c_2</u>	<u>c_3</u>
<u>c_1</u>	<u>c_1</u>	n_{11}	n_{12}	n_{13}
	<u>c_2</u>			
		Errors		



Precision: the class-specific precision of M for class C_i is

$$\text{prec}_i = \frac{N_{ii}}{M_i} = \frac{\text{true class } i \text{ predictions}^*}{\text{class } i \text{ predictions}^*} \quad (I = \max)$$

i.e. the fraction of our class i predictions that are correct.

Global precision is

$$\text{precision} = \sum_{i=1}^k \left(\frac{m_i}{n} \right) \text{prec}_i = \sum_{i=1}^k \frac{N_{ii}}{n} \quad (= \text{accuracy})$$

Recall:

$$\text{recall}_i = \frac{N_{ii}}{N_i} \quad (I = \max)$$

the fraction of class i points we correctly identify

There is a trade-off between precision and recall...

for example, if $M(x_i) = y_2$ for all x_i , then recall_2 is

$$\frac{N_2}{N_2} = 1$$

but precision is $\frac{N_2}{n}$, which is likely low.

The class-specific F-measure tries to balance them

$$F_i = \frac{2}{\frac{1}{\text{prec}_i} + \frac{1}{\text{recall}_i}} = \frac{2 \text{prec}_i \cdot \text{recall}_i}{\text{prec}_i + \text{recall}_i} = \frac{2 \cdot \frac{N_{ii}}{M_i} \cdot \frac{N_{ii}}{N_i}}{\frac{N_{ii}}{M_i} + \frac{N_{ii}}{N_i}}$$

$2N_{ii}^2$ $2N_{ii}$

$$= \frac{\sum_{i=1}^k m_i^2}{\sum_{i=1}^k m_i + \sum_{i=1}^k m_i} = \frac{\sum_{i=1}^k m_i^2}{2 \sum_{i=1}^k m_i}$$

m_i m_i
 for F_i to be near 1,
 we need both to be
 near 1.

Global F-score = $\frac{1}{k} \sum_{i=1}^k F_i$ AKA F1-score

Binary Classification: Suppose there are only 2 classes, we can call

		True Class	
		pos.	neg.
predicted class	pos	n_{11} = True Positives	n_{12} = False Positives
	neg	n_{21} = False Negatives	n_{22} = True Negatives

$$\text{prec}_p = \frac{TP}{TP+FP}$$

$\underbrace{TP}_{\% \text{ of pos. prediction that are correct}}$
 $+ FP$

$$\text{recall}_p = \frac{TP}{TP+FN}$$

$\underbrace{TP}_{\% \text{ of pos. points that are correct}}$
 $+ FN$

Receiver Operating Characteristic (ROC)

2 classes

popular strategy for assessing a binary classifier which outputs a score value for the positive class for each point in the test set, e.g. Bayes classifier outputs a posterior probability $P(c_1|x_i)$.

decision threshold

Typically, $M(x_i) = \begin{cases} c_1, & \text{if score} > \underline{s} \\ c_2, & \text{else} \end{cases}$ but \underline{s} is a somewhat arbitrary threshold

now like positive rate vs. true positive rate.

ROC plots false positive rate vs. true positive rate
 for each value of γ from γ^{\min} to γ^{\max}

$$\frac{FP}{\text{neg. points}} \quad \min S(x_i) \quad \frac{TP}{\text{pos. points}} \quad \max S(x_i)$$

γ -score for pos. class

At γ^{\min} ...

		truth			
		P	N		
pred.	P	O	O		
	N	FN	TN		

At γ^{\max} ...

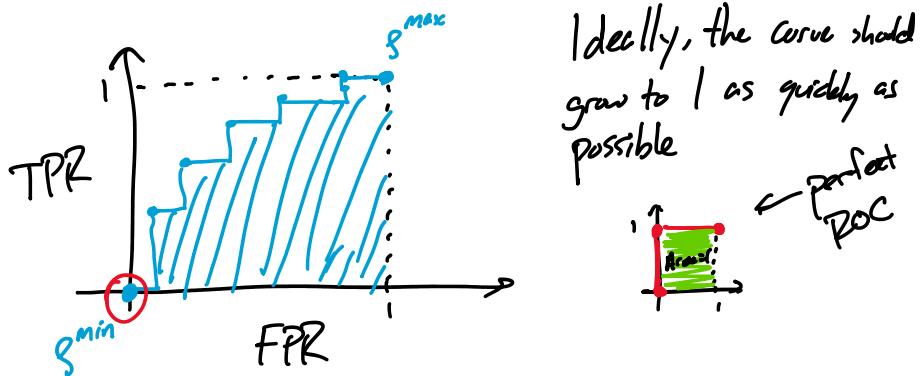
		truth			
		P	N		
pred.	P	TP	FP		
	N	O	O		

Ideal Classifier

		P	N
P	P	O	
	TP	O	
N	O	TN	

between...

$$R(\gamma) = \{x_i \in D \mid S(x_i) > \gamma\} - \text{points classified pos. for threshold } \gamma$$



Ideally, the curve should grow to 1 as quickly as possible



Area under ROC curve (AUC) is a one-number performance

metric in $[0,1]$ - the prob. classifier will rank a random pos. test point higher than a random neg. point.

imp: ROC/AUC is not sensitive to class imbalance

Projections

The angle between vectors a and b is

$$\cos \theta = \frac{a^T b}{\|a\| \|b\|} = \left(\frac{a}{\|a\|} \right)^T \left(\frac{b}{\|b\|} \right)$$

unit vectors

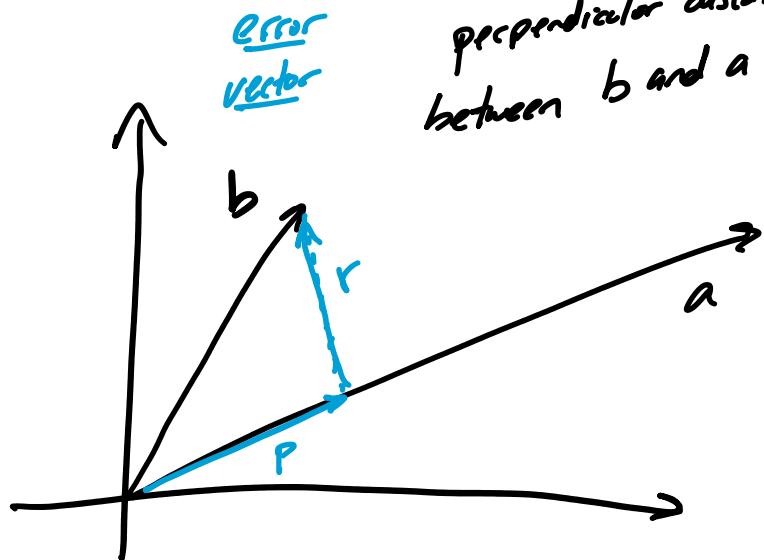
$$\Rightarrow \theta = 0 \Leftrightarrow \frac{a^T}{\|a\|} \frac{b}{\|b\|} = 0 \Leftrightarrow a^T b = 0$$

Any vector b can be written

$$b = b_{||} + b_{\perp} = p + r$$

where $b_{||} = p$ is parallel to a vector a and $b_{\perp} = r$ is orthogonal to a .

Orthogonal projection of b onto a



Note $p = ca$ for some constant c .

note $P = Ca$ for some constnt C .

Since P and r are orthogonal,

$$0 = (Ca)^T r = Ca^T(b - ca) = \cancel{ca^T b} - C^2 a^T a$$

$$ca^T a = a^T b$$

$$C = \frac{a^T b}{a^T a}$$

scalar of b on a *projection* $= \text{proj}_a(b)$ = "offset" along a

$$\Rightarrow P = Ca = \left(\frac{a^T b}{a^T a} \right) a$$

Eigenvalues and Eigenvectors

Let $A \in \mathbb{R}^{n \times n}$, $u \in \mathbb{R}^n$, $\lambda \in \mathbb{R}$

If $Au = \lambda u$, u is an eigenvector of A
 λ is an eigenvalue of A

a vector that
points in the
same direction
when multiplied by A

How do we find eigenvectors?

$$Au = \lambda u \Rightarrow Au - \lambda u = 0 \\ (A - \lambda I)u = 0$$

By the invertible matrix theorem, this system has a nonzero solution if and only if

computed at
 $O(n^{2.3n})$

Fundamental
theorem
of algebra

$$\det(A - \lambda I) = 0$$

$$(\lambda - \lambda_1) \cdot \dots \cdot (\lambda - \lambda_n) = 0$$

where $\lambda_1, \dots, \lambda_n$ are solutions (i.e. the eigenvalues)

$(x_i, y_i) \in \mathbb{R}^d \times \{c_1, \dots, c_k\}$ dataset

LDA seeks a vector w that maximizes the separation between classes after projection onto w "discriminate"

20.1 Optimal Linear Discriminant

D - dataset of n points $x_i \in \mathbb{R}^d$ + labels $y_i \in \{c_1, \dots, c_k\}$

D_i - class c_i points $\{x_j \in D \mid y_j = c_i\}$

$$|D_i| = n_i$$

Assume $k=2$, so $D = D_1 + D_2$

Let w be a unit vector. i.e. $\|w\|^2 = w^T w = 1$

The projection of x_i onto w is

$$x'_i = \left(\frac{w^T x_i}{w^T w} \right) w = \underbrace{(w^T x_i)}_{a_i} w$$

\Rightarrow projected points $\{a_1, \dots, a_n\}$ map \mathbb{R}^d to \mathbb{R} from d-dim points to 1-dim offsets

offset or scalar projection of x_i on the line w
AKA projected point

↑
 p_{c_i}' pts
○ - class 1
△ - class 2

Projected means:

$$M_1 = \frac{1}{n_1} \sum_{x_i \in D_1} a_i = \frac{1}{n_1} \sum w^T x_i$$

$$= w^T \left(\frac{1}{n_1} \sum x_i \right) = w^T M_1$$

$$M_2 = w^T M_2$$

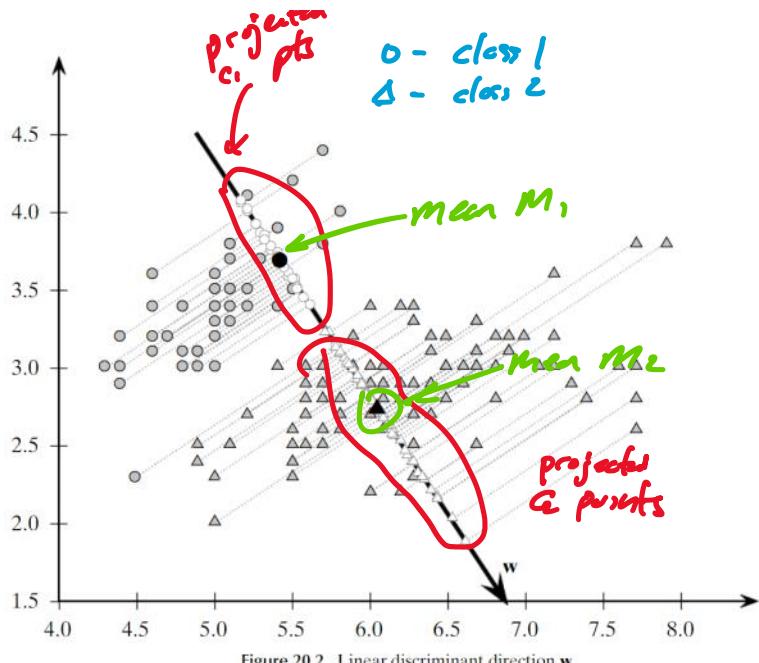


Figure 20.2. Linear discriminant direction w .

To maximize separation, maximizing $|M_1 - M_2|$ sounds reasonable, but we also need variance to not be too large so that the classes do not overlap too much

\Rightarrow LDA ensures the scatter s_i in each class is small

$$S_i = \sum_{x_j \in D_i} (a_j - M_i)^2 = n_i \sigma_i^2$$

variance for class c_i

Fisher LDA objective: $\max_w J(w) = \frac{(M_1 - M_2)^2}{S_1^2 + S_2^2}$

Solving this finds a direction w st. projected means are far apart without allowing large scatter

Let's solve it... first, express it in terms of input data:

$$(1^T - 1^T \dots 1^T) \Sigma (1^T / n_1 - 1^T / n_2) = (1^T / n_1 - 1^T / n_2) \Sigma (1^T / n_1 - 1^T / n_2)^T$$

$\omega^T \omega \sim$

$$(m_1 - m_2)^2 = (\omega^T(m_1 - m_2))^2 = \omega^T(m_1 - m_2) [\omega^T(m_1 - m_2)]^T$$

$$= \omega^T(m_1 - m_2)(m_1 - m_2)^T \omega$$

$$= \omega^T \mathcal{B} \omega$$

between-class scatter matrix

$$S_i^2 = \sum_{x_j \in D_i} (c_j - m_i)^2 = \sum_{x_j \in D_i} (\omega^T x_j - \omega^T m_i)^2$$

$$= \sum_{x_j \in D_i} (\omega^T(x_j - m_i))^2 =$$

$$= \omega^T \left(\sum_{x_j \in D_i} (x_j - m_i)(x_j - m_i)^T \right) \omega$$

scatter matrix for class $C_i = N_i \sum_i$

$$= \omega^T S_i \omega$$

Thus, $J(\omega) = \frac{\omega^T \mathcal{B} \omega}{\omega^T S_1 \omega + \omega^T S_2 \omega} = \frac{\omega^T \mathcal{B} \omega}{\omega^T S \omega}$ $S = S_1 + S_2$

Note: $S_1, S_2, S \in \mathbb{R}^{d \times d}$ are symmetric, positive semidefinite

$$\max_{\omega} J(\omega) = \max_{\omega} \frac{\omega^T \mathcal{B} \omega}{\omega^T S \omega}$$

Let's seek critical values...

$$Z \mathcal{B} \omega (\omega^T S \omega) - Z S \omega (\omega^T \mathcal{B} \omega)$$

..

$$\nabla J(\omega) = \frac{2B\omega(\omega^T S\omega) - 2S\omega(\omega^T B\omega)}{(\omega^T S\omega)^2} = 0$$

$$B\omega(\omega^T S\omega) = S\omega(\omega^T B\omega)$$

$$B\omega = S\omega \left(\frac{\omega^T B\omega}{\omega^T S\omega} \right)$$

$$B\omega = J(\omega) S\omega \quad \text{constant}$$

$$B\omega = \lambda S\omega$$

$\underbrace{}$

generalized eigenvalue problem

If S is nonsingular (S^{-1} exists), then

$$S^{-1}B\omega = \lambda S^{-1}S\omega = \lambda\omega$$

$$(S^{-1}B)\omega = \lambda\omega \Rightarrow \lambda \text{ is an eigenvalue}$$

Solving $\det(S^{-1}B - \lambda I) = 0$ gives several solutions, so we choose the largest since $\lambda = J(\omega)$ is to be maximized

(use linear algebra to compute eigenvalues)

Fisher's

Algorithm 20.1: Linear Discriminant Analysis

LINEARDISCRIMINANT (D):

Algorithm 20.1: Linear Discriminant Analysis

LINEARDISCRIMINANT (\mathbf{D}):

- 1 $\mathbf{D}_i \leftarrow \{\mathbf{x}_j^T \mid y_j = c_i, j = 1, \dots, n\}, i = 1, 2 // \text{ class-specific subsets}$
 - 2 $\boldsymbol{\mu}_i \leftarrow \text{mean}(\mathbf{D}_i), i = 1, 2 // \text{ class means}$
 - 3 ~~4~~ $\mathbf{B} \leftarrow (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T // \text{ between-class scatter matrix}$
 - 4 $\bar{\mathbf{D}}_i \leftarrow \mathbf{D}_i - \mathbf{1}_{n,i}\boldsymbol{\mu}_i^T, i = 1, 2 // \text{ center class matrices}$
 - 5 ~~6~~ $\mathbf{S}_i \leftarrow \bar{\mathbf{D}}_i^T \bar{\mathbf{D}}_i, i = 1, 2 // \text{ class scatter matrices}$
 - 6 ~~7~~ $\mathbf{S} \leftarrow \mathbf{S}_1 + \mathbf{S}_2 // \text{ within-class scatter matrix}$
 - 7 ~~8~~ $\lambda_1, \mathbf{w} \leftarrow \text{eigen}(\mathbf{S}^{-1} \mathbf{B}) // \text{ compute dominant eigenvector}$
-

$$(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T$$

In the 2-class case...

$$\begin{bmatrix} \boldsymbol{\mu}_{11} - \boldsymbol{\mu}_{21} \\ \vdots \\ \boldsymbol{\mu}_{1d} - \boldsymbol{\mu}_{2d} \end{bmatrix}_{d \times 1} \begin{bmatrix} \boldsymbol{\mu}_{11} - \boldsymbol{\mu}_{21} & \cdots & \boldsymbol{\mu}_{1d} - \boldsymbol{\mu}_{2d} \end{bmatrix}_{1 \times d}$$

$$\mathbf{B}\mathbf{w} = \lambda \mathbf{S}\mathbf{w}$$

$$c(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) = -\lambda \mathbf{S}\mathbf{w}$$

$$\frac{c}{\lambda} \mathbf{S}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) = \mathbf{w} \Rightarrow \text{optional solution } \mathbf{w} = \mathbf{S}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$$

constant
