

Lecture 3

- Classification specific concepts
- Kinds of predictors
 - Nearest Neighbors

- Next Lecture : in person
- HW1 - t.b. posted TODAY
- 12-oct4
13-oct6 } t.b. posted --"
 ↳ on web : handouts
- Q1 → Tue oct 18
- Steve WR : Tutorial Refresher

Lecture Notes I – Examples of Predictors

Marina Meilă
`mmp@stat.washington.edu`

Department of Statistics
University of Washington

September 29, 2022

Prediction problems by the type of output ✓

The “learning” paradigm and vocabulary ✓

→ Classification concepts

The Nearest-Neighbor and kernel predictors ←

?

Linear predictors

- Least squares regression
- Linear Discriminant Analysis (LDA)
- QDA (Quadratic Discriminant Analysis)
- Logistic Regression

NOT Linear

X The PERCEPTRON algorithm

• Classification and regression tree(s) (CART)

• The Naive Bayes classifier

Hastie & al

Reading HTF Ch.: 2.3.1 Linear regression, 2.3.2 Nearest neighbor, 4.1–4 Linear classification, 6.1–3. Kernel regression, 6.6.2 kernel classifiers, 6.6.3 Naive Bayes, 9.2 CART, 11.3 Neural networks, Murphy Ch.: 1.4.2 nearest neighbors, 1.4.4 linear regression, 1.4.5 logistic regression, 3.5 and 10.2.1 Naive Bayes, 4.2.1–3 linear and quadratic discriminant, 14.7.3– kernel regression, locally weighted regression, 16.2.1–4 CART, (16.5 neural nets), Bach Ch.: —

N/A

The “sign trick” for transforming a regressor into a classifier

The *sign* function $\text{sgn}(y) = y/|y|$ if $y \neq 0$ and 0 iff $y = 0$ turns a real valued variable Y into a discrete-valued one. This function is used to allow one to construct *real-valued classifiers*. In these classifiers, the model $f(x)$ is a real-valued function, and the prediction \hat{y} is given by $\text{sgn}(f(x))$.

Note that in a vanishingly small fraction of cases, when the value of $f(x)$ is exactly 0, no label will be assigned to the input x .

Classifiers with Continuous output

Binary

$$y \in \{\pm 1\}$$

Classifier $f: \mathcal{X} \rightarrow \{\pm 1\}$

↑ space of inputs \mathcal{X}
e.g. \mathbb{R}^d

But sometimes

$f: \mathcal{X} \rightarrow \mathbb{R}$, ($f \in \mathcal{F}$ ^{learned} model class)

$$\hat{y} = \text{sgn } f(x)$$

↑ classifier output
↑ label assigned by classifier to input x

$$\text{sgn } z = \begin{cases} +1, & z > 0 \\ -1, & z < 0 \\ 0, & z = 0 \end{cases}$$

Data $\mathcal{D} = \{(x^i, y^i), i=1:n\}$ $y^i \in \{\pm 1\}$ in our model

\mathcal{E}_x : Naive Bayes $\rightarrow f(x) \in (0,1)$ $f(x) = \mathbb{P}_i[Y=1|X=x] \Rightarrow$

logistic regression

$$\Rightarrow \hat{y}(x) = \text{sgn}\left(f(x) - \frac{1}{2}\right)$$

[Support Vector Machines]

[kernel classifiers] later

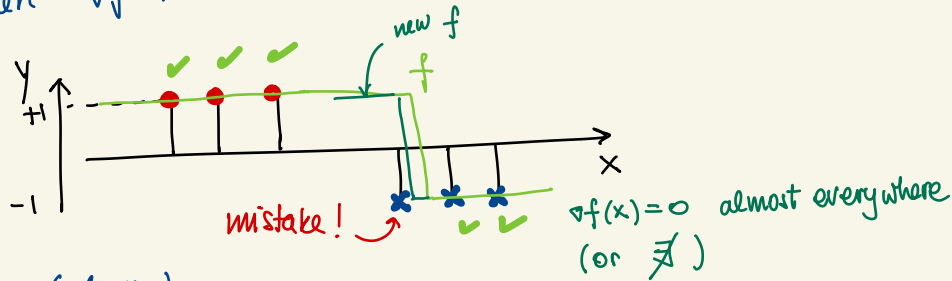
[Neural Net]

$$\Leftrightarrow \hat{y}(x) = 1 \text{ iff } \mathbb{P}[Y=1|X=x] > \frac{1}{2}$$

Why (continuous) \mathbb{R} -valued f ?

- for statistical models, $f = \Pr[Y=1|X]$ (probabilistic) (e.g. logistic regression)
- for non-statistical, $|f(x)|$ large \approx confidence in \hat{y} (e.g. SVM, Neural Net, kernel)
 \uparrow qualitative
- f differentiable ($\Leftrightarrow \nabla f$ exists) and ∇f not 0 almost everywhere)

then ∇f is used in learning algorithm



Margin of classifier (at x)

- $z = y f(x) \Rightarrow \underline{\text{P1}} \quad \begin{array}{ll} z > 0 & \text{iff } f(x) \text{ correct at } x \\ z < 0 & \text{iff } \text{mistake at } x \end{array}$

\uparrow
true

P2 $|z| = |f(x)|$ is classifier's confidence

$\hookrightarrow z > 0$ and large \rightarrow classifier "knows" it's correct
 Robustness!!

$z \geq 0 \rightarrow$ Not robust

Continuous valued Multiway $y \in \{1, 2, \dots, m\}$

Train $f_{1:m}$ "classifiers"

$$\hat{y}(x) = \operatorname{argmax}_{c=1:m} f_c(x)$$

Margin

$$z_{(x)} = \underbrace{f_{\hat{y}}(x)}_{\text{true}} - \max_{c \neq \hat{y}} f_c(x)$$

$$\begin{aligned} \overset{\text{correct}}{\hat{y}} = y &\Rightarrow z = f_y - f_{\text{next largest}} > 0 \\ \overset{\text{error}}{\hat{y}} \neq y &\Rightarrow z = f_y - f_{\hat{y}} < 0 \end{aligned}$$

Decision regions, decision boundary of a classifier

Let $f(x)$ be a classifier (not necessarily binary)

- ▶ $f(x)$ takes only a finite set of values
- ▶ The **decision region** associated with class y = the region in X space where f takes value y , i.e. $D_y = \{x \in \mathbb{R}^d, f(x) = y\} = f^{-1}(y)$.
- ▶ The boundaries separating the decision regions are called **decision boundaries**.

Some label

Binary

$$D_+ = \{x \in \mathbb{R}^d, f(x) > 0\}$$

$$D_- = \{x \in \mathbb{R}^d, f(x) < 0\}$$

$$\text{decision boundary} \left\{ \begin{array}{l} D_0 = \{x \in \mathbb{R}^d, f(x) = 0\} \\ f \text{ real valued} \end{array} \right.$$

$$D_0 = \text{boundary between } D_+, D_-$$

any x, y, f

$$f: X \rightarrow Y$$

$$A \subseteq Y$$

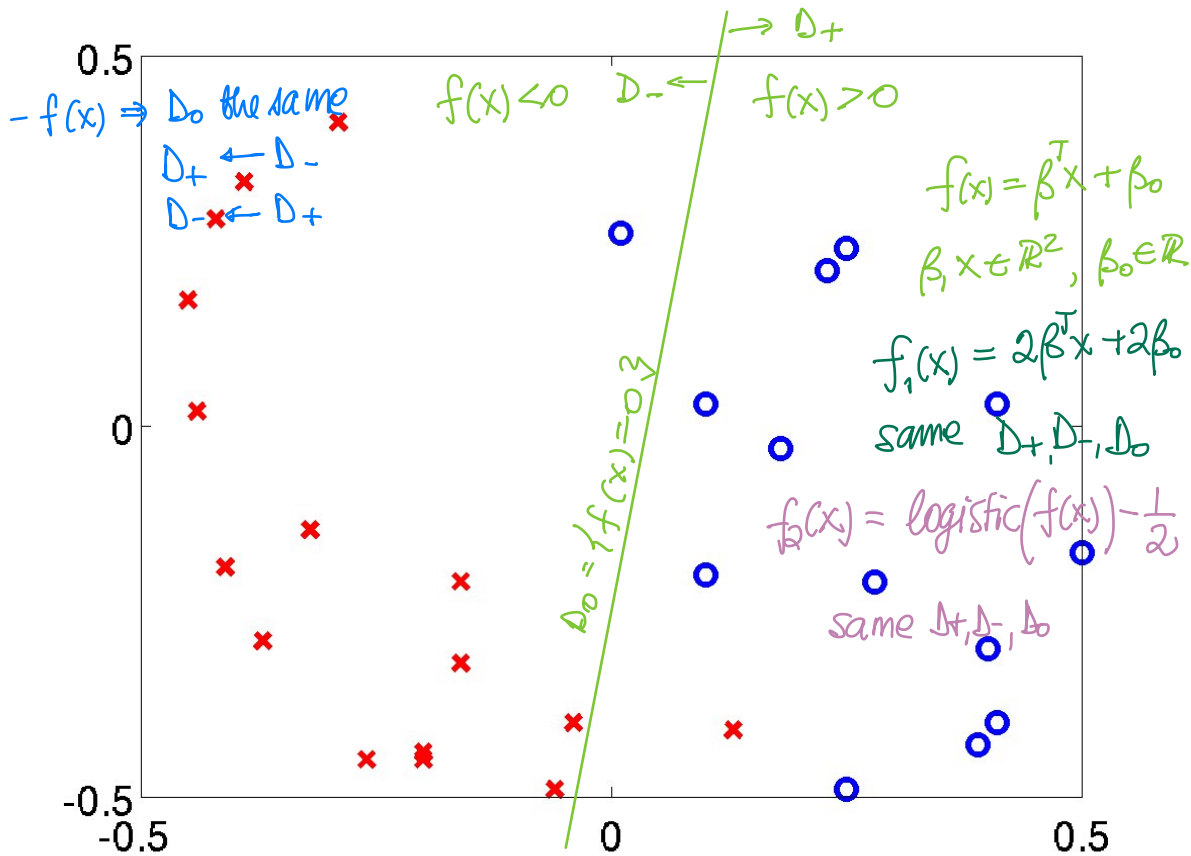
$$f^{-1}(A) = \{x \in X \mid f(x) \in A\}$$

Decision regions, decision boundary of a classifier

Let $f(x)$ be a classifier (not necessarily binary)

- ▶ $f(x)$ takes only a finite set of values
- ▶ The **decision region** associated with class y = the region in X space where f takes value y , i.e. $D_y = \{x \in \mathbb{R}^d, f(x) = y\} = f^{-1}(y)$.
- ▶ The boundaries separating the decision regions are called **decision boundaries**.

Multiway Δ_0 = boundary between any two $D_y, D_{y'}$
 $y \neq y'$



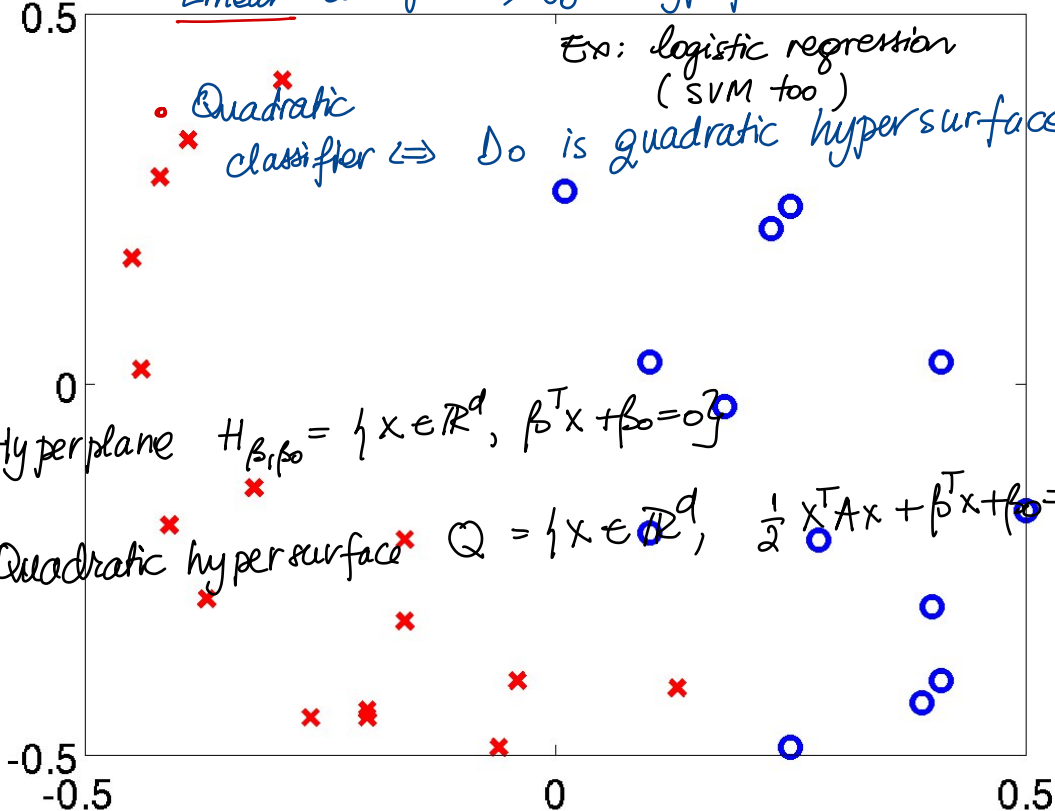
• Linear classifier $\Leftrightarrow D_0$ is hyperplane in \mathbb{R}^d "linear"

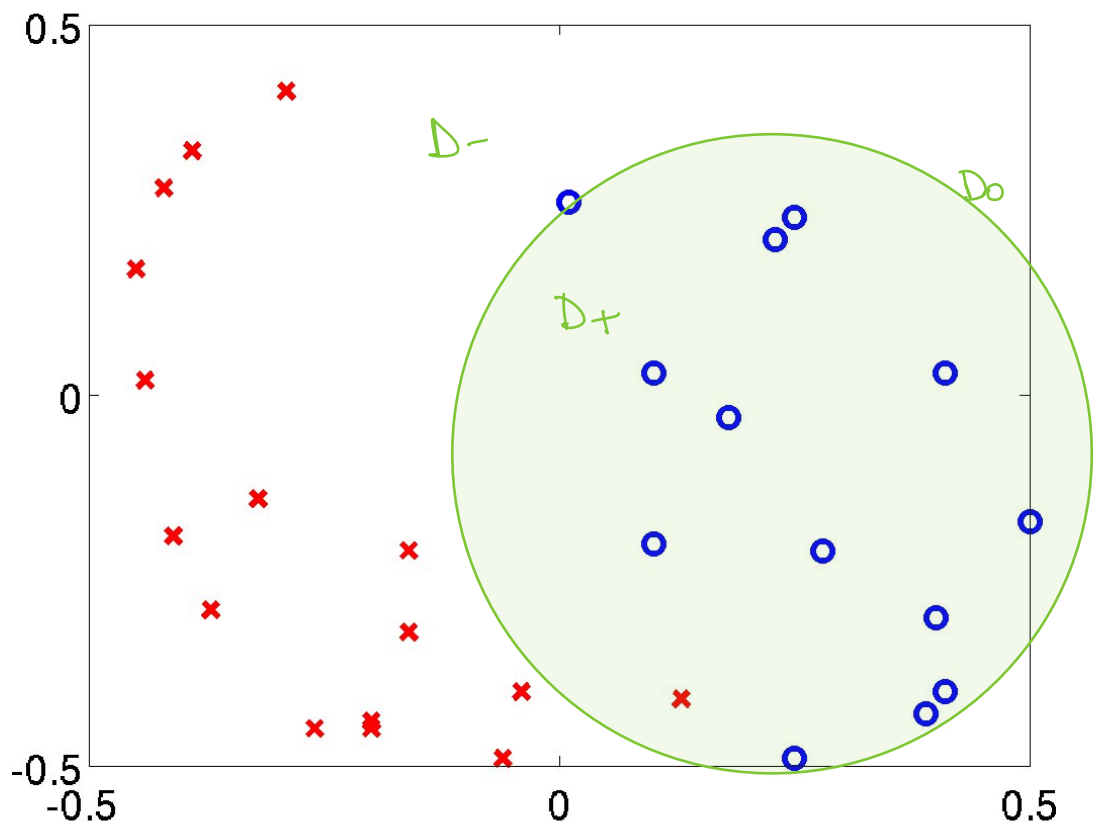
Ex: logistic regression
(SVM too)

• Quadratic classifier $\Leftrightarrow D_0$ is quadratic hypersurface

Hyperplane $H_{\beta, \beta_0} = \{x \in \mathbb{R}^d, \beta^T x + \beta_0 = 0\}$

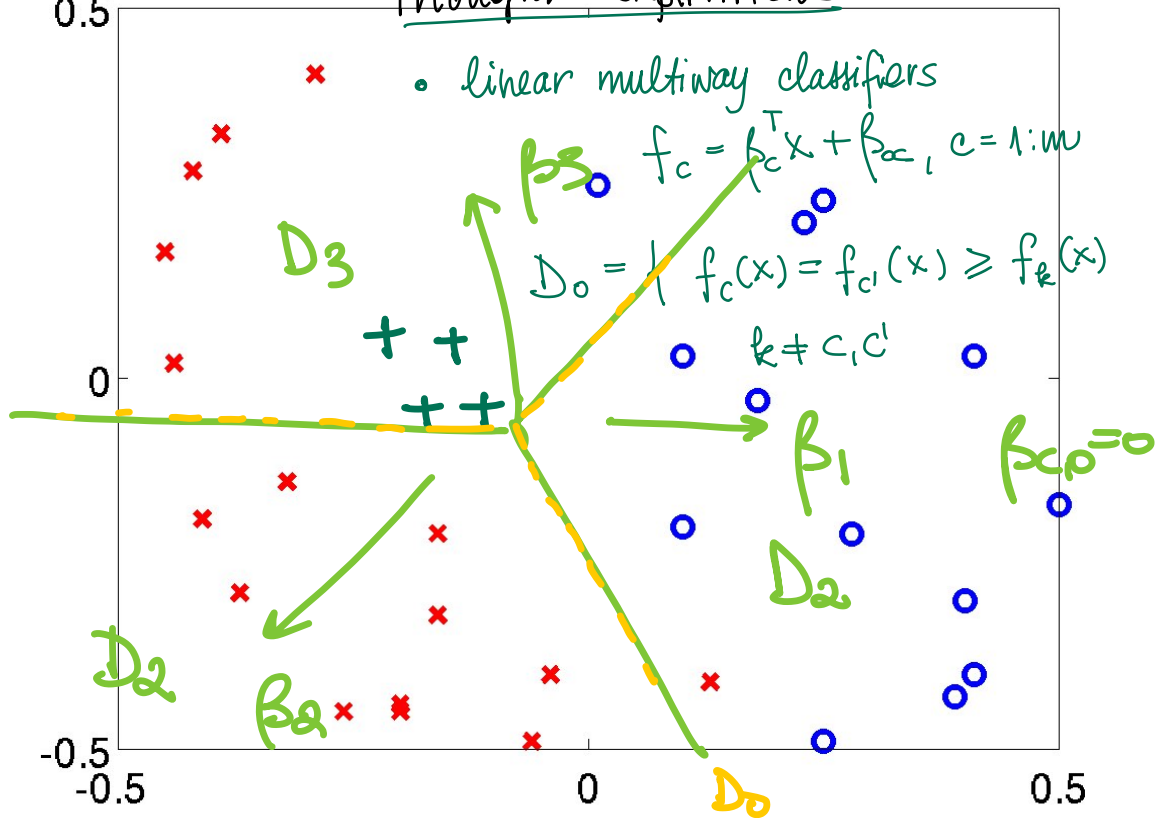
Quadratic hypersurface $Q = \{x \in \mathbb{R}^d, \frac{1}{2} x^T A x + \beta^T x + \beta_0 = 0\}$





Thought experiments

- linear multiway classifiers



Decision regions, decision boundary of a classifier

Let $f(x)$ be a classifier (not necessarily binary)

- ▶ $f(x)$ takes only a finite set of values
- ▶ The **decision region** associated with class y = the region in X space where f takes value y , i.e. $D_y = \{x \in \mathbb{R}^d, f(x) = y\} = f^{-1}(y)$.
- ▶ The boundaries separating the decision regions are called **decision boundaries**.
- ▶ For a binary classifier, we have two decision regions, D_+ and D_- . By convention $f(x) = 0$ on the decision boundary.
- ▶ For binary classifier with real valued $f(x)$ (i.e. $\hat{y} = \text{sgn}f(x)$) we define $D_+ = \{x \in \mathbb{R}^d, f(x) > 0\}$, $D_- = \{x \in \mathbb{R}^d, f(x) < 0\}$ and the decision boundary $\{x \in \mathbb{R}^d, f(x) = 0\}$

$$\ln\left(\frac{P[Y=1|X]}{P[Y=-1|X]}\right) = \beta_0 + \beta_1 X + \frac{1}{2}\beta_2 X^2$$

$$D_0 = \{x, P[Y=1|X] = P[Y=-1|X]\}$$

$$\Leftrightarrow \beta_0 + \beta_1 X + \frac{1}{2}\beta_2 X^2 = 0 \quad \text{quadratic}$$

dim $X > 1$ (parabola, hyperb, ellipse)