

STAT 391

5/25/23

Lecture 18

logistic regression

Regression examples

Clustering 2 lectures

HW 6 optional

Last HW deadline
Friday 5 pm ← Solution
5, 6

T.B Posted

- notes on linear
- z logistic reg.
- slides on
- Clustering
- + b. edited

logistic Regression

used for: Classification

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

$$y \in \{0, 1\} \text{ "Label"}$$

$$\mathcal{D} = \{(x^i, y^i)_{i=1:n}\}$$

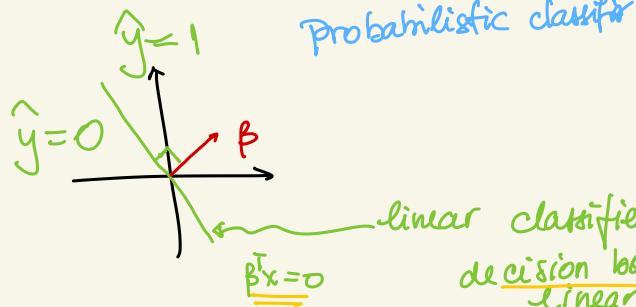
Model $f(x) = \ln$

$$\frac{P_{y|x}(y=1|x)}{P_{y|x}(y=0|x)} = \underline{\beta^T x} \quad \beta \in \mathbb{R}^d$$

meaning
 $\underbrace{\frac{P_{y|x}(y=1|x)}{P_{y|x}(y=0|x)}}_{\text{odds } \in (0, \infty)}$
 $\underbrace{\ln \frac{P_{y|x}(y=1|x)}{P_{y|x}(y=0|x)}}_{\text{log odds } \in (-\infty, \infty)}$

Classification with $f = \beta^T x \quad (\beta \text{ known})$

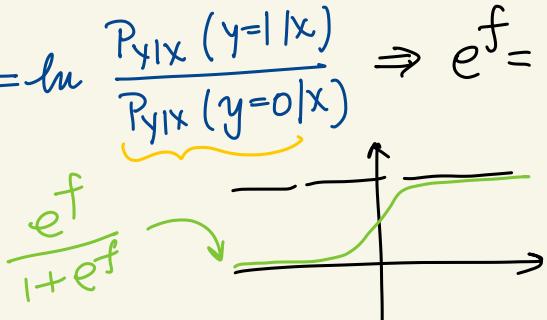
$$\hat{y} = \frac{\text{sign } f(x) + 1}{2} \iff f(x) > 0 \Leftrightarrow \frac{P(y=1|x)}{P(y=0|x)} > \frac{P(y=0|x)}{P(y=1|x)}$$



- $P_{Y|X}[y=1|x] = p$

$$f = \ln \frac{p}{1-p} = \ln \frac{P_{Y|X}(y=1|x)}{P_{Y|X}(y=0|x)}$$

$$\hat{\beta}^T x$$



$$P = \frac{e^f}{1+e^f}$$

$$1-p = \frac{1}{1+e^f}$$

$$p+(1-p)=1$$

- Estimating β by Max likelihood

Likelihood

$$L(\beta) = P_{Y|X}[y^{1:n}|x^{1:n}, \beta]$$

$$= \prod_{i=1}^n P_{Y|X}[y^i|x^i, \beta] = \prod_{i=1}^n \frac{e^{y^i \beta^T x^i}}{1+e^{\beta^T x^i}}$$

P.
for $y^i=1$

1-P.
for $y^i=0$

confidence
of
classification

log likelihood

$$\ell(\beta) = \sum_{i=1}^n \left[y^i \beta^T x^i - \ln(1+e^{\beta^T x^i}) \right]$$

maximize over β
by Gradient Ascent

$$l(\beta) = \sum_{i=1}^n \left[y_i \beta^T x_i - \ln(1 + e^{\beta^T x_i}) \right]$$

Gradient $\frac{\partial l}{\partial \beta} \in \mathbb{R}^d$

$$\frac{\partial l}{\partial \beta_j} = \frac{\partial l}{\partial f} \cdot \frac{\partial f}{\partial \beta_j} = \sum_{i=1}^n \left[y_i - \frac{e^{f(x_i)}}{1 + e^{f(x_i)}} \right] x_i$$

$$\frac{\partial l}{\partial \beta} = \left[\frac{\partial l}{\partial \beta_j} \right]_{j=1:d} = \frac{\partial l}{\partial f} \left[\frac{\partial f}{\partial \beta_j} \right]_{j=1:d}$$

$$\frac{\partial f}{\partial \beta} = \frac{\partial}{\partial \beta} (\beta^T x) = X$$

Exercise prove this

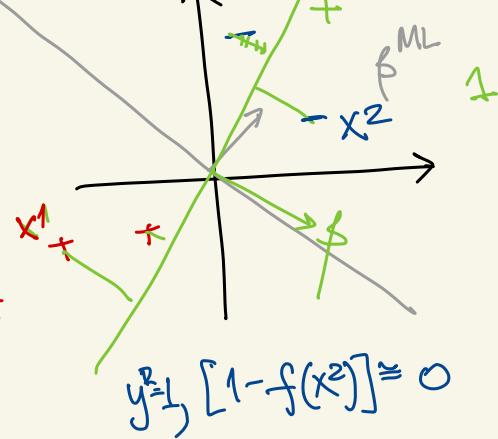
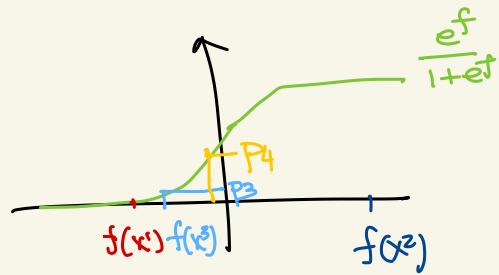
$$\frac{\partial f}{\partial \beta_j} = x_j$$

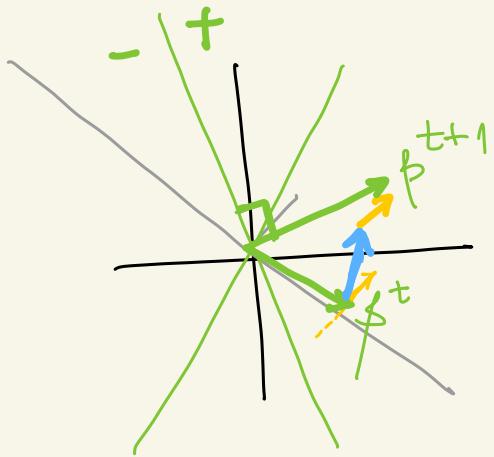
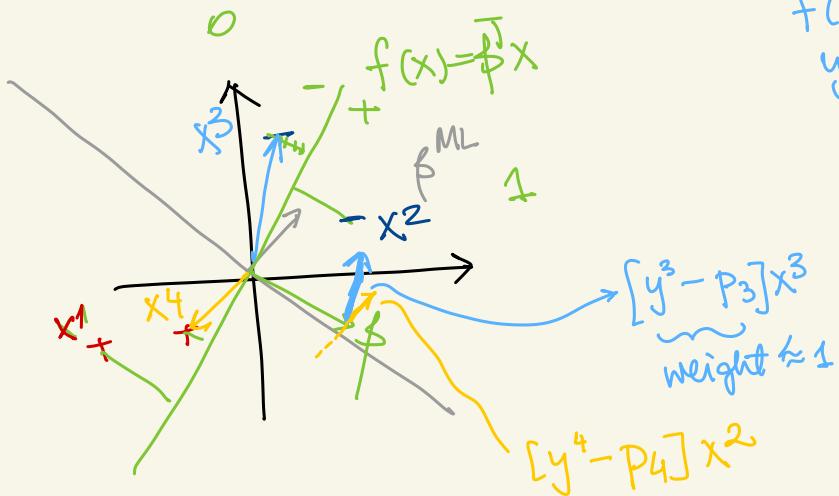
$$\frac{\partial l}{\partial f} = \sum_{i=1}^n \left[y_i - \underbrace{\frac{e^{f(x_i)}}{1 + e^{f(x_i)}}}_{P(y_i=1|x_i)} \right] \in \mathbb{R}$$

by the model

$$0 \approx \frac{e^{f(x)}}{1 + e^{f(x)}} = P_1 \quad \begin{cases} |f(x)| \text{ large} \\ f(x) < 0 \text{ correct} \end{cases}$$

$$0 > [y_i - P_1] \approx 0$$





$$\begin{aligned}\hat{\beta}^{t+1} &= \hat{\beta}^t + \eta \frac{\partial L}{\partial \beta} \\ &= \hat{\beta}^t + \eta \sum_{i=1}^n [y_i - \hat{\beta}^t] x^i\end{aligned}$$

linear combination
of x^i

Prediction \rightarrow regression $y \in \mathbb{R}$ (linear)
 \equiv Supervised learning
 classification $y \in \{0, 1\}$ or discrete
 $D = \{(x^i, y^i)\}_{i=1:n}$ (logistic)
 labels, output

$P_{y|x}$ wanted

Clustering: wanted Clusters = groups in data
 Unsupervised Wanted some feature of P_x distribution of x

$D = \{x^i\}_{i=1:n}$
 no output

Ex: estimate - a density
 - a discrete distribution

Ex: dimension reduction PCA

Ex: causal inference

$$x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad x_1 \xrightarrow{\text{causes}} x_2$$

x_1 = treatment, x_2 = effect, health
 smoking cancer

Lecture Notes IX – Clustering

Marina Meilă
`mmp@stat.washington.edu`

Department of Statistics
University of Washington

May, 2023

Paradigms for clustering 

+ what is clustering

Parametric clustering algorithms (K given)

Cost based / hard clustering

K-means clustering and the quadratic distortion

Model based / soft clustering

Issues in parametric clustering

Selecting K

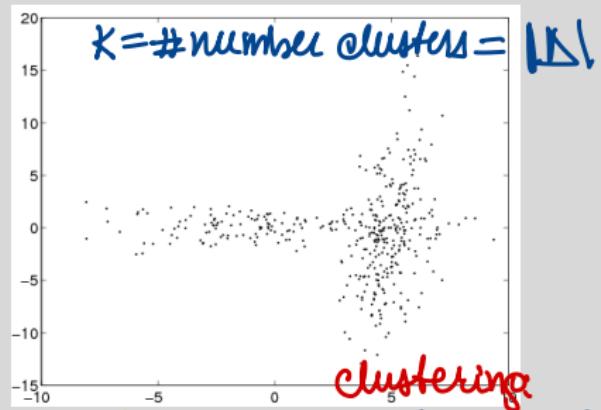
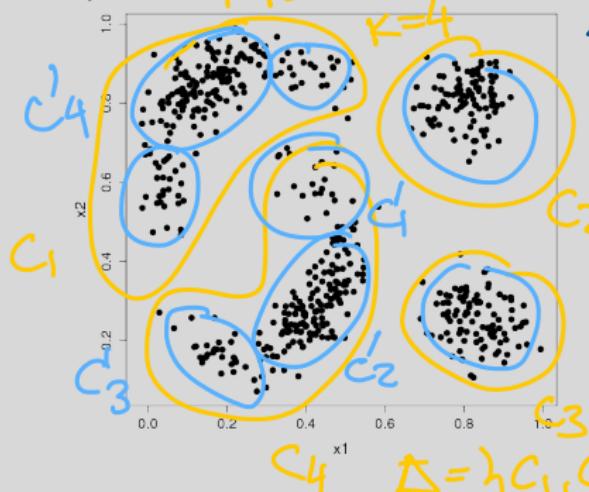
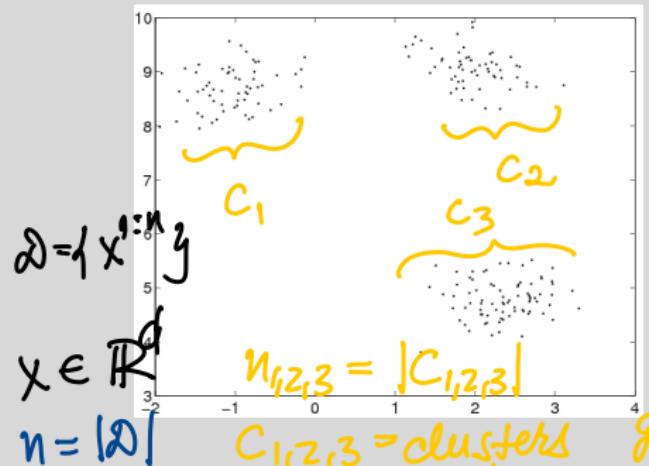
Reading: Ch. 18

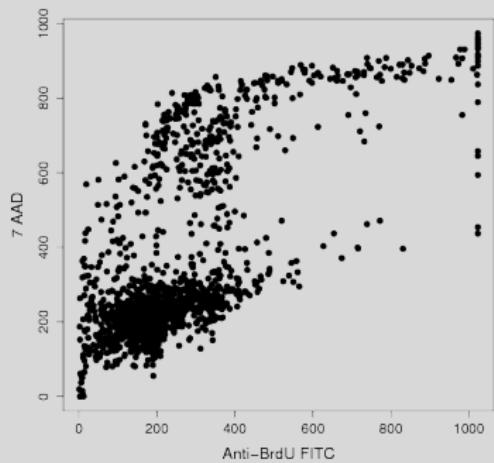
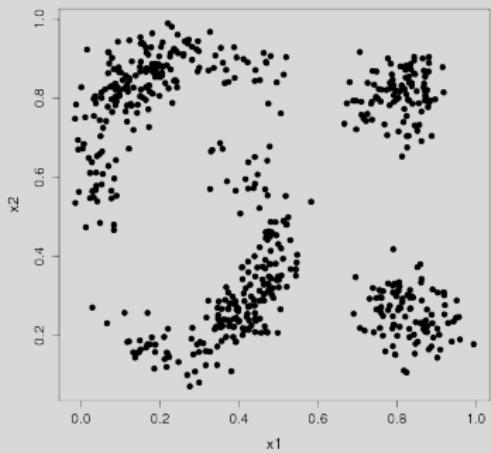
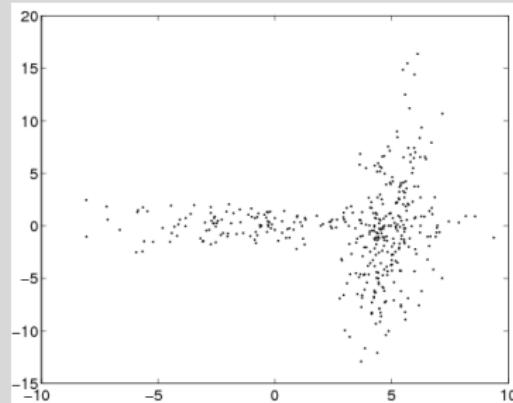
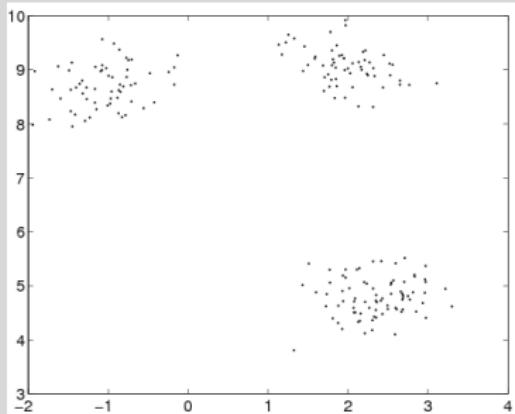
What is clustering? Problem and Notation

- ▶ **Informal definition** **Clustering** = Finding groups in data
- ▶ **Notation**
 - \mathcal{D} = $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ a **data set**
 - n = number of **data points**
 - K = number of **clusters** ($K \ll n$)
 - Δ = $\{C_1, C_2, \dots, C_K\}$ a partition of \mathcal{D} into disjoint subsets
 - $k(i)$ = the **label** of point i
 - $\mathcal{L}(\Delta)$ = cost (loss) of Δ (to be minimized)
- ▶ **Second informal definition** **Clustering** = given n **data points**, separate them into K **clusters**
- ▶ Hard vs. soft clusterings
 - ▶ **Hard** clustering Δ : an item belongs to only 1 cluster
 - ▶ **Soft** clustering $\gamma = \{\gamma_{ki}\}_{k=1:K}^{i=1:n}$
 γ_{ki} = the **degree of membership** of point i to cluster k

$$\sum_k \gamma_{ki} = 1 \text{ for all } i$$

(usually associated with a probabilistic model)





(from [Nugent and Meila, 2010])

Paradigms

← each method
has own definition of cluster

Depend on type of data, type of clustering, type of cost (probabilistic or not), and constraints (about K , shape of clusters)

- ▶ Data = vectors $\{x_i\}$ in \mathbb{R}^d

Parametric (K known)	Cost based [hard]
	Model based [soft]

Non-parametric	Dirichlet process mixtures [soft]
----------------	-----------------------------------

(K determined by algorithm)	Information bottleneck [soft]
	Modes of distribution [hard]

	Gaussian blurring mean shift[Carreira-Perpinan, 2007] [hard]
--	--

- ▶ Data = similarities between pairs of points $[S_{ij}]_{i,j=1:n}$, $S_{ij} = S_{ji} \geq 0$ **Similarity based clustering**

Graph partitioning	spectral clustering [hard, K fixed, cost based]
--------------------	---

	typical cuts [hard non-parametric, cost based]
--	--

Affinity propagation	[hard/soft non-parametric]
----------------------	----------------------------

Classification vs Clustering

	Classification	Clustering
Cost (or Loss) \mathcal{L}	Expected error	many! (probabilistic or not)
	Supervised	Unsupervised
Generalization	Performance on new data is what matters	Performance on current data is what matters
K	Known	Unknown
"Goal"	Prediction	Exploration <i>Lots of data to explore!</i>
Stage of field	Mature	Still young