

## Module 5: Classification

# Linear Methods: Logistic Regression

STAT/BIOSTAT 527, University of Washington

Emily Fox

May 22<sup>nd</sup>, 2014

©Emily Fox 2014

1

Very convenient!

$$p(y = 0 | x, \beta) = \frac{1}{1 + \exp(\beta_0 + \sum_j \beta_j x_j)}$$

implies

$$p(y = 1 | x, \beta) = \frac{\exp(\beta_0 + \sum_j \beta_j x_j)}{1 + \exp(\beta_0 + \sum_j \beta_j x_j)}$$

Handwritten notes:  $g = \text{logit} = \log\left(\frac{p(y=1|x)}{p(y=0|x)}\right)$ , link fn.  $\rightarrow g(\eta) = \beta_0 + \beta_1 x_1 + \dots$ ,  $1 - P(y=0|x)$

Examine ratio:

$$\frac{p(y = 1 | x, \beta)}{p(y = 0 | x, \beta)} = \exp(\beta_0 + \sum_j \beta_j x_j) > 1 \Rightarrow \begin{matrix} \text{class 1 wins,} \\ \text{else class 0} \\ \text{(under 0-1 loss)} \end{matrix}$$

implies

$$\log \frac{p(y = 1 | x, \beta)}{p(y = 0 | x, \beta)} = \beta_0 + \sum_j \beta_j x_j > 0 \Rightarrow \text{class 1 wins, as before}$$

Handwritten notes:  $\log \text{ odds}$ ,  $g(\eta)$ , linear, linear classification rule!

©Emily Fox 2014

2

# Maximizing Conditional Log Likelihood

$$p(y=0 | x, \beta) = \frac{1}{1 + \exp(\beta_0 + \sum_j \beta_j x_j)}$$

$$p(y=1 | x, \beta) = \frac{\exp(\beta_0 + \sum_j \beta_j x_j)}{1 + \exp(\beta_0 + \sum_j \beta_j x_j)}$$

$$l(\beta) = \sum_i \log p(y_i | x_i, \beta)$$

$$= \sum_i y_i (\beta_0 + \sum_j \beta_j x_{ij}) - \log(1 + \exp(\beta_0 + \sum_j \beta_j x_{ij}))$$

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

fixed in training data

**Good news:**  $l(\beta)$  is concave function of  $\beta$ , no local optima problems

**Bad news:** no closed-form solution to maximize  $l(\beta)$

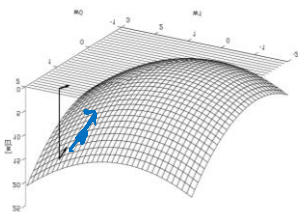
**Good news:** concave functions easy to optimize

©Emily Fox 2014

3

# Optimizing Concave Function – Gradient Ascent

- Conditional likelihood for logistic regression is concave
- Find optimum with gradient ascent



**Gradient:**  $\nabla_{\beta} l(\beta) = \left[ \frac{\partial l(\beta)}{\partial \beta_0}, \dots, \frac{\partial l(\beta)}{\partial \beta_d} \right]'$

Step size,  $\eta > 0$

**Update rule:**  $\Delta \beta = \eta \nabla_{\beta} l(\beta)$

$$\beta_j^{(t+1)} \leftarrow \beta_j^{(t)} + \eta \frac{\partial l(\beta)}{\partial \beta_j}$$

- Gradient ascent is simplest of optimization approaches

- e.g., Conjugate gradient ascent can be much better

Often, esp. proofs,  $\eta$  gets smaller w/ iterations  
e.g.  $\eta_t = \frac{1}{t}$  const.

©Emily Fox 2014

4

# Gradient Ascent for LR

start w/  $\beta^{(0)}$  (e.g. 0)

revisit soon

Gradient ascent algorithm: iterate until change  $< \epsilon$

can do in //  $\forall j$  repeat

$$\beta_0^{(t+1)} \leftarrow \beta_0^{(t)} + \eta \sum_i \left( y_i - \hat{p}(y=1 | x_i, \beta^{(t)}) \right) \frac{\partial L(\beta)}{\partial \beta_0}$$

For  $j=1, \dots, d$ ,

$$\beta_j^{(t+1)} \leftarrow \beta_j^{(t)} + \eta \sum_i x_{ij} \left( y_i - \hat{p}(y=1 | x_i, \beta^{(t)}) \right) \frac{\partial L(\beta)}{\partial \beta_j} \quad j=1, \dots, d$$

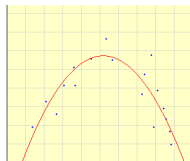
©Emily Fox 2014

5

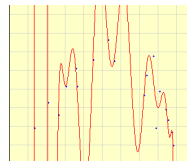
# Regularization in Linear Regression

- Overfitting usually leads to very large parameter choices, e.g.:

$$-2.2 + 3.1 X - 0.30 X^2$$



$$-1.1 + 4,700,910.7 X - 8,585,638.4 X^2 + \dots$$



even for  $n \gg p$ ,  $p$  large

- Regularized** or **penalized** regression aims to impose a “complexity” penalty by penalizing large weights
  - “Shrinkage” method

$$\hat{\beta} = \arg \min_{\beta} \sum_{i=1}^n \underbrace{(y_i - (\beta_0 + \beta^T x_i))^2}_{\text{RSS}} + \underbrace{\lambda \|\beta\|}_{\text{penalty}}$$

©Emily Fox 2014

6





## Regularized Conditional Log Likelihood

- Add regularization penalty, e.g.,  $L_2$ :

$$l(\beta) = \log \prod_{i=1}^n p(y_i | x_i, \beta) - \frac{\lambda}{2} \|\beta\|_2^2$$

log-likelihood

- Practical note about  $\beta_0$ :
- Gradient of regularized likelihood:

©Emily Fox 2014

9

## Standard v. Regularized Updates

- Maximum conditional likelihood estimate

$$\hat{\beta} = \arg \max_{\beta} \log \prod_{i=1}^n p(y_i | x_i, \beta)$$

normal  
logistic

$$\beta_j^{(t+1)} \leftarrow \beta_j^{(t)} + \eta \sum_i x_{ij} (y_i - \hat{p}(y = 1 | x_i, \beta^{(t)}))$$

- Regularized maximum conditional likelihood estimate

$$\hat{\beta} = \arg \max_{\beta} \log \prod_{i=1}^n p(y_i | x_i, \beta) - \frac{\lambda}{2} \sum_{j=1}^d \beta_j^2$$

regularized

$$\beta_j^{(t+1)} \leftarrow \beta_j^{(t)} + \eta \left\{ -\lambda \beta_j^{(t)} + \sum_i x_{ij} (y_i - \hat{p}(y = 1 | x_i, \beta^{(t)})) \right\}$$

pushes coeff. towards 0

©Emily Fox 2014

10

# Stopping Criterion

$$l(\beta) = \log \prod_{i=1}^n p(y_i | x_i, \beta) - \frac{\lambda}{2} \|\beta\|_2^2$$

- When do we stop doing gradient ascent?

$$\|l(\beta^*) - l(\beta^{(t)})\| \leq \epsilon$$

stopping criterion

- Because  $l(\mathbf{w})$  is strongly concave:

- i.e., because of some technical condition

$$l(\beta^*) - l(\beta) \leq \frac{1}{2\lambda} \|\nabla l(\beta)\|_2^2$$

- Thus, stop when:

$$\|l(\beta^*) - l(\beta)\| < \frac{1}{2\lambda} \|\nabla l(\beta)\|_2^2 < \epsilon$$

©Emily Fox 2014

11

## Digression:

## Logistic Regression for $K > 2$

- Logistic regression in more general case ( $K$  classes), where  $Y$  in  $\{1, \dots, K\}$

©Emily Fox 2014

12

## Digression: Logistic Regression for $K > 2$

- Logistic regression in more general case, where  $Y \in \{1, \dots, K\}$

for  $k < K$

$$p(y = k | \mathbf{x}, \beta) = \frac{\exp(\beta_{k0} + \sum_{j=1}^d \beta_{kj} x_j)}{1 + \sum_{k'=1}^{K-1} \exp(\beta_{k'0} + \sum_{j=1}^d \beta_{k'j} x_j)}$$

for  $k=K$  (normalization, so no weights for this class)

$$p(y = K | \mathbf{x}, \beta) = \frac{1}{1 + \sum_{k'=1}^{K-1} \exp(\beta_{k'0} + \sum_{j=1}^d \beta_{k'j} x_j)}$$

**Estimation procedure is basically the same  
as what we derived!**

©Emily Fox 2014

13

## The Cost, The Cost!!! Think about the cost...

- What's the cost of a gradient update step for LR???

$$\beta_j^{(t+1)} \leftarrow \beta_j^{(t)} + \eta \left\{ -\lambda \beta_j^{(t)} + \sum_i x_{ij} (y_i - \hat{p}(y = 1 | x_i, \beta^{(t)})) \right\}$$

$\forall j$

$O(d)$

$O(nd)$

Naively,  $O(nd \times d = nd^2)$   
 but if you "cache"  $\hat{p}$  (same  $\forall i$ )  $\rightarrow O(nd)$ .  
 However, if "n" is huge (or online streaming), this  
 is slow.

$\begin{pmatrix} \beta_0^{(t)} \\ \vdots \\ \beta_d^{(t)} \end{pmatrix}$

©Emily Fox 2014

14

## Gradient ascent in Terms of Expectations

- “True” objective function:

$$l(\beta) = E_x[l(\beta, x)] = \int p(x)l(\beta, x)dx$$

- Taking the gradient:

$$\nabla_{\beta} l(\beta) = E_x[\nabla_{\beta} l(\beta, x)]$$

- “True” gradient ascent rule:

$$\beta^{(t+1)} \leftarrow \beta^{(t)} + \eta E_x[\nabla_{\beta} l(\beta, x)]$$

- How do we estimate expected gradient?

We have been min.  
training loss s.t.  
complexity penalty  
but we really  
want to min.

Can't compute it  
So approx. from  
sample

©Emily Fox 2014

15

## SGD: Stochastic Gradient Ascent (or Descent)

- “True” gradient:  $\nabla l(\beta) = E_x[\nabla l(\beta, x)]$

- Sample based approximation: take  $x_i$  iid

$$E_x[\nabla l(\beta, x)] \approx \frac{1}{n} \sum_{i=1}^n \nabla_{\beta} l(\beta, x_i)$$

use this for  
normal  
logistic  
reg.

- What if we estimate gradient with just one sample???

- ☐ Unbiased estimate of gradient
- ☐ Very noisy! *high var.*
- ☐ Called stochastic gradient ascent (or descent)
  - Among many other names
- ☐ VERY useful in practice!!!

©Emily Fox 2014

16

# Stochastic Gradient Ascent for Logistic Regression

- Logistic loss as a stochastic function:

$$E_x[l(\beta, x)] = E_x \left[ \log p(y | x, \beta) - \frac{\lambda}{2} \|\beta\|_2^2 \right]$$

*Handwritten note:  $\frac{1}{n} \sum_{i=1}^n \log p(y_i | x_i, \beta)$*

- Batch** gradient ascent updates:

$$\beta_j^{(t+1)} \leftarrow \beta_j^{(t)} + \eta \left\{ -\lambda \beta_j^{(t)} + \frac{1}{n} \sum_{i=1}^n x_{ij} \left( y_i - \hat{p}(y = 1 | x_i, \beta^{(t)}) \right) \right\}$$

*Handwritten notes: "normal LR" with a circle around "Batch"; "new" with a circle around the fraction  $\frac{1}{n}$*

- Stochastic gradient ascent updates:

- Online setting:

$$\beta_j^{(t+1)} \leftarrow \beta_j^{(t)} + \eta \left\{ -\lambda \beta_j^{(t)} + x_{i(t),j} \left( y_{i(t)} - \hat{p}(y = 1 | x_{i(t)}, \beta^{(t)}) \right) \right\}$$

*Handwritten notes: "take 1 data pt. at iter. t: x\_i(t)"; "use this same data pt. x\_i(t), \forall j"*

©Emily Fox 2014

17

## What you should know...

- Classification: predict discrete classes rather than real values
- Logistic regression model: Linear model
  - Logistic function maps real values to [0,1]
- Optimize conditional likelihood
- Gradient computation
- Overfitting
- Regularization
- Regularized optimization
- Cost of gradient step is high, use stochastic gradient descent

©Emily Fox 2014

18

## Module 5: Classification

### Linear Methods: LDA and QDA

STAT/BIOSTAT 527, University of Washington

Emily Fox

May 22<sup>nd</sup>, 2014

©Emily Fox 2014

19

## Discriminative vs. Generative

- So far, we have considered modeling/fitting

$$p(Y | X)$$

input predictor

"discriminative" method  
for fixed  $X$ , max  $P(Y|X)$   
for fitting

- There are also a large set of **generative** methods

- Model:

- Class-conditional densities  $f_k(X) = P(X|Y=k)$
- Class prior probabilities  $\pi_k = P(Y=k) \rightarrow \sum \pi_k = 1$

- Via Bayes' rule:

$$P(Y=k | X=x) = \frac{\pi_k f_k(x)}{\sum \pi_l f_l(x)} = P(Y=k | x)$$

©Emily Fox 2014

20

# Generative Classifiers

$$p(Y = k | X = x) = \frac{\pi_k f_k(x)}{\sum_{\ell} \pi_{\ell} f_{\ell}(x)}$$

## ■ Examples include:

- Linear and quadratic discriminative analysis (LDA and QDA)

→ linear (+quad) decision boundaries

- Mixture of Gaussians (saw in BNP module)

→ non-linear boundary

- Nonparametric density estimation for  $f_k(x)$

KDE, very flexible

- Naïve Bayes

assumes a simple form for  $f_k(x)$

©Emily Fox 2014

21

# Linear Discriminative Analysis

- Assume Gaussian class-conditional densities

$$f_k(X) = \frac{1}{(2\pi)^{d/2} |\Sigma_k|^{1/2}} e^{-\frac{1}{2}(x-\mu_k)^T \Sigma_k^{-1} (x-\mu_k)}$$

↑  $\mathbb{R}^d$

- Furthermore, consider equal covariances

$$\Sigma_k = \Sigma, \forall k$$

- Log odds

$$\begin{aligned} \log \frac{p(Y = k | X = x)}{p(Y = \ell | X = x)} &= \log \frac{\pi_k}{\pi_{\ell}} + \log \frac{f_k(x)}{f_{\ell}(x)} \\ &= \log \frac{\pi_k}{\pi_{\ell}} - \frac{1}{2} (\mu_k + \mu_{\ell})^T \Sigma^{-1} (\mu_k - \mu_{\ell}) + x^T \Sigma^{-1} (\mu_k - \mu_{\ell}) \end{aligned}$$

(Σ's same)      ↑ linear in x

©Emily Fox 2014

22

# Linear Discriminative Analysis



$$\log \frac{p(Y = k | X = x)}{p(Y = \ell | X = x)} = \log \frac{\pi_k}{\pi_\ell} - \frac{1}{2}(\mu_k + \mu_\ell)^T \Sigma^{-1}(\mu_k - \mu_\ell) + x^T \Sigma^{-1}(\mu_k - \mu_\ell)$$

- Equivalently,

$$\log \frac{p(Y = k | X = x)}{p(Y = \ell | X = x)} = \delta_k(x) - \delta_\ell(x) = 0$$

where

$$\delta_k(x) = x^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \log \pi_k$$

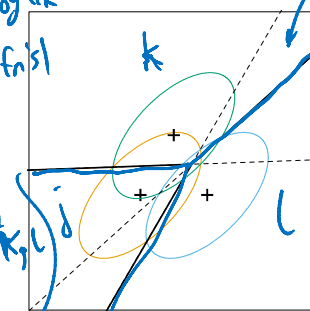
- Decision rule:

$$\hat{g}(x) = \underset{k}{\operatorname{argmax}} \delta_k(x)$$

- Linear decision boundaries

$$\{x : \delta_k(x) = \delta_\ell(x)\}, \mu_k, \mu_\ell$$

if  $\Sigma = \sigma^2 I$ ,



From Hastie, Tibshirani, Friedman book

©Emily Fox 2014

23

## LDA Parameter Estimation

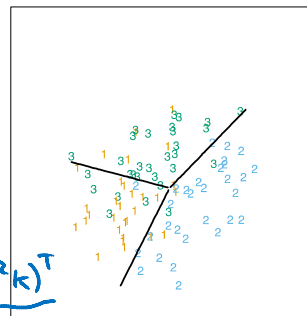
$$\log \frac{p(Y = k | X = x)}{p(Y = \ell | X = x)} = \log \frac{\pi_k}{\pi_\ell} - \frac{1}{2}(\mu_k + \mu_\ell)^T \Sigma^{-1}(\mu_k - \mu_\ell) + x^T \Sigma^{-1}(\mu_k - \mu_\ell)$$

- Based on the training class labels,  $\{y_i\}_{i=1}^n$ ,  $\Sigma$ ,  $\mu_k$ ,  $\pi_k$  estimate parameters:

$$\hat{\pi}_k = \frac{n_k}{n}$$

$$\hat{\mu}_k = \frac{\sum_{y_i=k} x_i}{n_k}$$

$$\hat{\Sigma} = \frac{1}{n-K} \sum_{k=1}^K \sum_{y_i=k} (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T$$



From Hastie, Tibshirani, Friedman book

©Emily Fox 2014

24



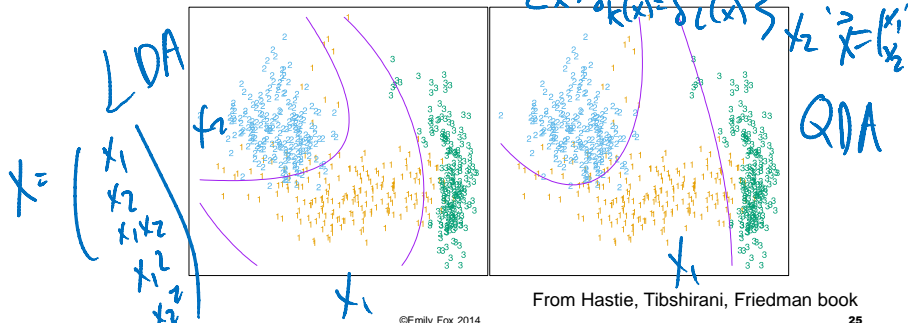
# Quadratic Discriminative Analysis

- Same setup as LDA, but allow class-specific covariances

- Quadratic discriminant functions:

$$\delta_k(x) = -\frac{1}{2} \log |\Sigma_k| - \frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) + \log \pi_k$$

- Quadratic decision boundaries



## QDA Parameter Estimation

- Based on the training class labels, estimate parameters:

$\hat{\pi}_k, \hat{\mu}_k$  same as before

$$\hat{\Sigma}_k = \sum_{i \in k} (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T / (n_k - 1)$$

- Number of parameters:

LDA	QDA
$(k-1)(d+1)$	$(k-1) \left( \frac{d(d+3)}{2} + 1 \right)$

many more parameters from  $\hat{\Sigma}_k$

- Can also consider shrinkage estimators

$$\hat{\Sigma}_k(\alpha) = \alpha \hat{\Sigma}_k + (1 - \alpha) \hat{\Sigma} \quad \hat{\Sigma}(\gamma) = \gamma \hat{\Sigma} + (1 - \gamma) \sigma^2 I$$

↑ distinct  $\hat{\Sigma}_k$

↑ common  $\hat{\Sigma}$

↑ common  $\hat{\Sigma}$

↑ spherical

# Notes on QDA and LDA

- LDA + QDA tend to perform very well in practice
- It is not true that data are Gaussian or, furthermore, that covariances are equal (LDA)
- Performance is likely attributed to the fact that the data can only support simple decision boundaries
  - Also, estimates for Gaussian models are stable

bias-variance tradeoff  
 ↑ linear boundary      ↑ lower than in more complex classifiers

©Emily Fox 2014

27

# LDA vs. Logistic Regression

- Both have linear log odds:

$$\log \frac{p(Y = k | X = x)}{p(Y = K | X = x)} = \alpha_{k0} + \alpha_k^T x$$

$$\log \frac{p(Y = k | X = x)}{p(Y = K | X = x)} = \beta_{k0} + \beta_k^T x$$

- Difference is in how the coefficients are estimated

$$p(X, Y = k) = p(X) p(Y = k | X)$$

↑  
~~h~~ h<sub>k</sub>?

↑ same form for both models

©Emily Fox 2014

28

# LDA vs. Logistic Regression

$$p(X, Y = k) = p(X)p(Y = k | X)$$

- Marginal likelihood term

- Logistic regression:

arbitrary... just maximize likelihood  
kind of like estimating  $p(x)$  nonparametrically  
empirically w/ mass  $\frac{1}{n}$  @ each  $x_i$ .

- LDA:

$$p(X) = \sum_k \pi_k N(x; \mu_k, \Sigma)$$

mixture density  
↑  
params.

©Emily Fox 2014

29

# LDA vs. Logistic Regression

- In LDA, the data inform the parameters more

- If data are indeed Gaussian, then asymptotically maximizing just conditional likelihood requires 30% more data to perform as well

- Data far from boundary affect  $\Sigma$  in LDA, but are ignored by logistic regression

→ LDA is not robust to outliers

- Observations without class labels can be used in mixture model case, but not in logistic regression

$x_i$ 's w/o  $y_i$ 's (LDA)

- Marginal likelihood  $p(X)$  acts as a regularizer

2-class, lin. separable + log. reg. → ML est. are undefined

→ LDA coeff. for some data are well-defined

- Logistic regression tends to be more robust than LDA and can handle qualitative  $X$  variables, but performance is often similar.

©Emily Fox 2014

30

## Module 5: Classification

# Nonparametric Methods: KDE and Naïve Bayes

STAT/BIOSTAT 527, University of Washington

Emily Fox

May 22<sup>nd</sup>, 2014

©Emily Fox 2014

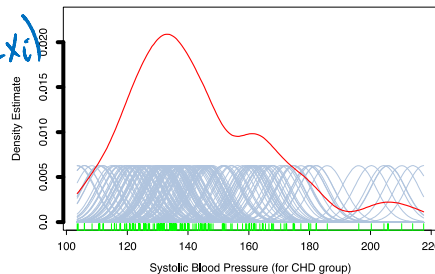
31

## KDE for Classification

$$p(Y = k \mid X = x) = \frac{\pi_k f_k(x)}{\sum_{\ell} \pi_{\ell} f_{\ell}(x)}$$

- Use KDE to estimate class-conditional densities
- Recall commonly used Gaussian KDE in 1D

$$\hat{f}_k(x) = \frac{1}{n_k} \sum_{i \in k} \phi_h(x - x_i)$$



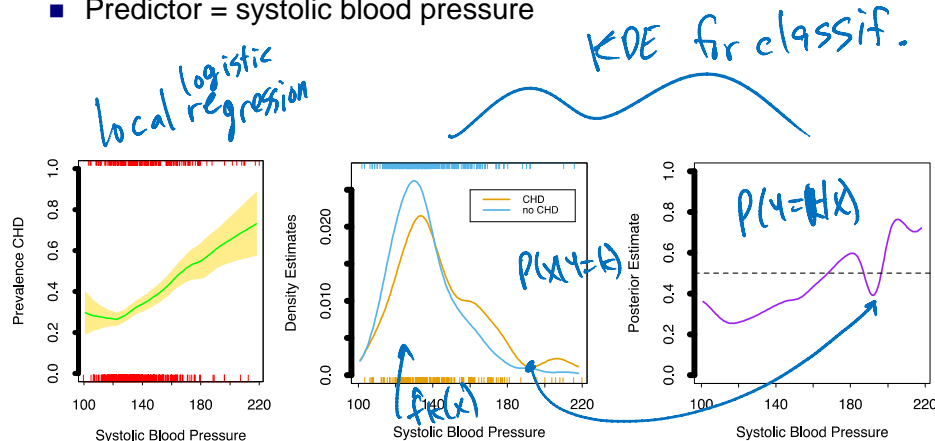
From Hastie, Tibshirani, Friedman book

©Emily Fox 2014

32

# Example: Heart Disease Data

- Binary response = CHD (coronary heart disease)
- Predictor = systolic blood pressure

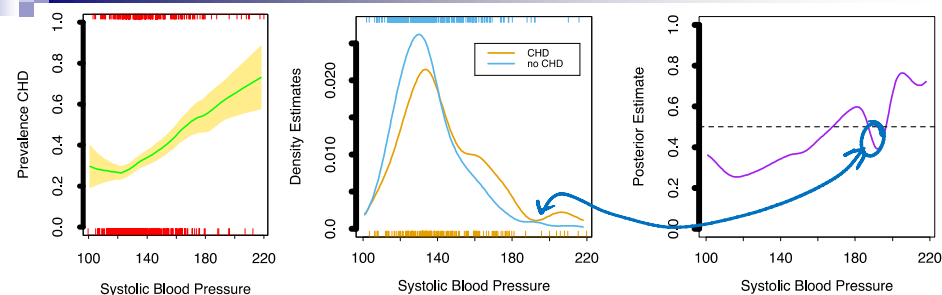


©Emily Fox 2014

33

# Example: Heart Disease Data

From Hastie, Tibshirani, Friedman book

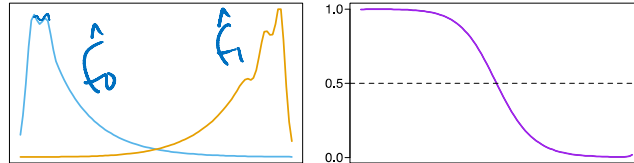


- KDE estimates are poor in regions with little data
- Local linear model uses variable bandwidth based on k-NN  
→ smooths out over these regions
- For classification tasks, do not need to estimate each class-conditional density well. Just need good estimates of the posterior near the decision boundary

©Emily Fox 2014

34

# Class-Conditionals vs. Posterior



## ■ Example:

- Both densities are multimodal
- Might opt for rougher, high-variance estimator to capture features
- However, posterior is quite smooth
- Fine-scale features are irrelevant for classification here

©Emily Fox 2014

35

# Multivariate KDE

■ In 1d

$$\hat{p}(x_0) = \frac{1}{n\lambda} \sum_{i=1}^n K_{\lambda}(x_0, x_i)$$

- In  $\mathbb{R}^d$ , assuming a product kernel,

$$\hat{p}(x_0) = \frac{1}{n\lambda_1 \cdots \lambda_d} \sum_{i=1}^n \left\{ \prod_{j=1}^d K_{\lambda_j}(x_{0j}, x_{ij}) \right\}$$

lots of params to choose

- Typical choice = Gaussian RBF

→ Gaussian KDE

$$e^{-\frac{\|x_0 - x_i\|^2}{\lambda}}$$

©Emily Fox 2014

36

# Naïve Bayes Classifier

$$p(Y = k | X = x) = \frac{\pi_k f_k(x)}{\sum_{\ell} \pi_{\ell} f_{\ell}(x)}$$

- Useful in high-dimensional settings ( $d$  large)
- Assumes factored form for class-conditional densities

$$f_k(X) = \prod_{j=1}^d f_{kj}(x_j)$$

↑  
1D  
KDE

generally, not true,  
but often performs well

- Benefits:
  - Estimate  $f_{kj}(X_j)$  separately for each  $j$  using only 1D KDE
  - If  $X_j$  of  $X$  is discrete, then can combine using a histogram estimate

©Emily Fox 2014

37

# Naïve Bayes Classifier

$$p(Y = k | X = x) = \frac{\pi_k \prod_j f_{kj}(x_j)}{\sum_{\ell} \pi_{\ell} \prod_j f_{\ell j}(x_j)}$$

- Log odds

$$\log \frac{p(Y = k | X = x)}{p(Y = \ell | X = x)} = \log \frac{\pi_k}{\pi_{\ell}} + \sum_i \log \frac{f_{ki}(x_i)}{f_{\ell i}(x_i)}$$

"discriminative"  
↓

$$= \alpha_{k_0} + \sum_{j=1}^d g_{kj}(x_j)$$

- Has form of GAM, but fit very differently
  - Analogous to difference between LDA and logistic regression

NB: generative

©Emily Fox 2014

38

## Module 5: Classification

# Mixture Models for Classification

STAT/BIOSTAT 527, University of Washington

Emily Fox

May 22<sup>nd</sup>, 2014

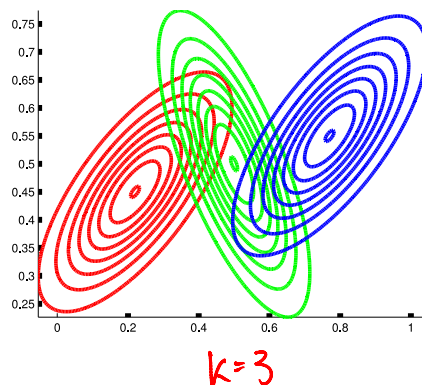
©Emily Fox 2014

39

## Density as Mixture of Gaussians

- Approximate density with a mixture of Gaussians

Mixture of 3 Gaussians



$$p = p(x_i | \pi, \mu, \Sigma) = \sum_{k=1}^K \pi_k N(x_i | \mu_k, \Sigma_k)$$

*Handwritten notes:*

- $K$ : # of mix comp.
- $\pi_k$ : mix. weights
- $\mu_k, \Sigma_k$ : shape params
- Gauss. kernel, just like in KDE, but not centered at obs
- In 1D:  $P = \text{target density}$
- $\sum \pi_k = 1$

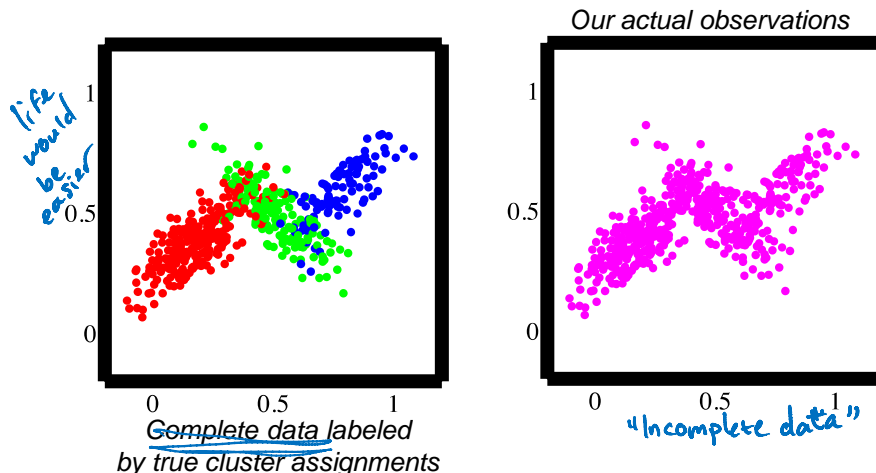
©Emily Fox 2014

40



# Clustering our Observations

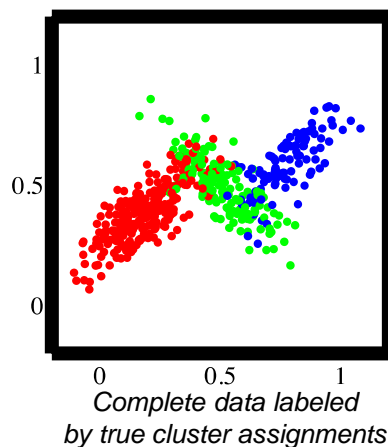
- Imagine we have an assignment of each  $x_i$  to a Gaussian



C. Bishop, Pattern Recognition & Machine Learning

# Clustering our Observations

- Imagine we have an assignment of each  $x_i$  to a Gaussian



- Introduce latent cluster indicator variable  $z_i$

$$z_i \in \{1, \dots, K\}$$

$$\Pr(z_i = k) = \pi_k$$

- Then we have

$$p(x_i | z_i, \pi, \mu, \Sigma) =$$

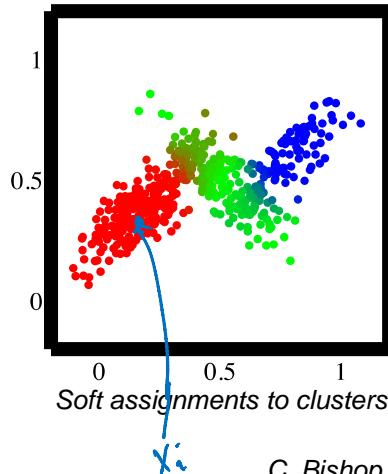
$$N(x_i | \mu_{z_i}, \Sigma_{z_i})$$

param. est. is easy if we have  $\pi, \mu, \Sigma$   
 $\Rightarrow$  decoupled into  $K$  Gauss. est.

C. Bishop, Pattern Recognition & Machine Learning

# Clustering our Observations

- We must infer the cluster assignments from the observations



- Posterior probabilities of assignments to each cluster \*given\* model parameters:

$$r_{ik} = p(z_i = k \mid x_i, \pi, \theta) =$$

$$= \frac{\pi_k N(x_i \mid \mu_k, \Sigma_k)}{\sum_j \pi_j N(x_i \mid \mu_j, \Sigma_j)}$$

*motivates an iterative alg.*

C. Bishop, Pattern Recognition & Machine Learning

## Mixture Models for Classification

- Can use mixture models as a generative classifier in the unsupervised setting

- EM algorithm = iteratively:

- Estimate responsibilities given parameter estimates

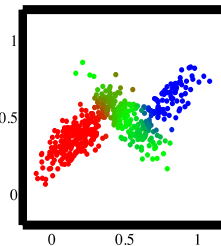
$$\hat{r}_{ik} = \frac{\hat{\pi}_k N(x_i, \hat{\mu}_k, \hat{\Sigma}_k)}{\sum_{\ell} \hat{\pi}_{\ell} N(x_i, \hat{\mu}_{\ell}, \hat{\Sigma}_{\ell})}$$

- Maximize parameters given responsibilities

- For classification, threshold the estimated responsibilities

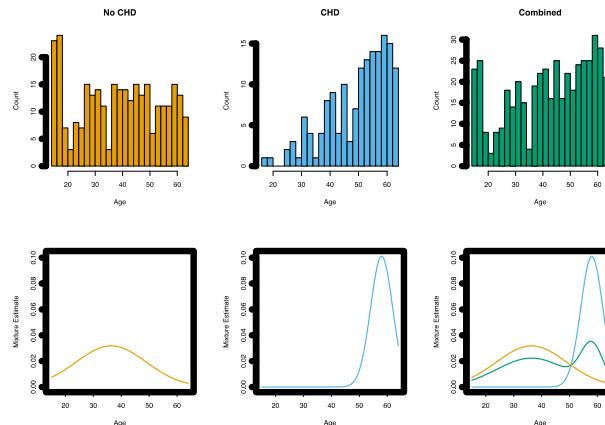
- E.g.,  $\hat{g}(x_i) = \arg \max_k \hat{r}_{ik}$

- Note: allows non-linear boundaries as in QDA



# Example: Heart Disease Data

- Binary response = CHD (coronary heart disease)
- Predictor = systolic blood pressure



From Hastie, Tibshirani, Friedman book

©Emily Fox 2014

45

## What you need to know

- Discriminative vs. Generative classifiers
- LDA and QDA assume Gaussian class-conditional densities
  - Results in linear and quadratic decision boundaries, respectively
- KDE for classification
  - Challenging in areas with little data or in high dimensions
  - Estimating class-conditionals is not optimizing classification objective
- Naïve Bayes assumes factored form
  - Results in log odds that have GAM form
- Mixture models allow for unsupervised generative approach

©Emily Fox 2014

46

# Readings



- Hastie, Tibshirani, Friedman – 4.3, 4.4.5, 6.6.2-6.6.3, 6.8