

Module 5: Classification

Linear Methods: Logistic Regression

STAT/BIOSTAT 527, University of Washington

Emily Fox

May 28th, 2013

©Emily Fox 2013

1

Very convenient!

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

implies

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

1 - P(y=0|x)

Examine ratio:

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

implies *log odds*

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

linear

linear classification rule!

©Emily Fox 2013

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 ∈ ℝ^d

fixed in training data

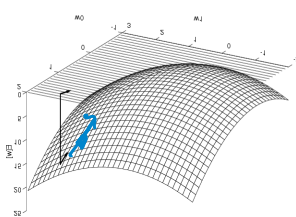
Good news: $l(\beta)$ is concave function of β , no local optima problems

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

Good news: concave functions easy to optimize

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 can be much better

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

Gradient Ascent for LR

revisit soon

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

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

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

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

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

repeat

©Emily Fox 2013

5

Linear Separability

$\exists \beta$ s.t. all pos. examples have $\beta_0 + \sum \beta_j x_j > 0$

+ neg. examples have $\beta_0 + \sum \beta_j x_j < 0$

$$\beta_0 + \sum \beta_j x_j > 0$$

$$2\beta_0 + \sum 2\beta_j x_j > 0$$

"more sure" for no reason ... same separation

$$\beta_0 + \sum \beta_j x_j = 0$$

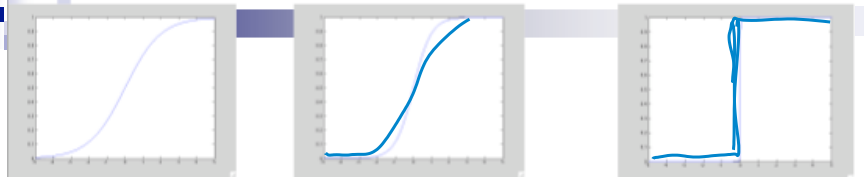
$$2\beta_0 + \sum 2\beta_j x_j = 0$$

$\Rightarrow \alpha \beta$ is also a separating hyperplane $\forall \alpha > 0$

©Emily Fox 2013

6

Large Parameters → Overfitting



$$\frac{1}{1 + e^{-x}}$$

$$\frac{1}{1 + e^{-2x}}$$

$$\frac{1}{1 + e^{-100x}}$$

- If data is linearly separable, weights go to infinity

$$p(y=0 | \beta, x) = \frac{1}{1 + e^{\beta_0 + \sum_j \beta_j x_j}} \quad \left. \vphantom{p(y=0 | \beta, x)} \right\} \text{ increases as } \|\beta\| \rightarrow \infty$$

- In general, leads to overfitting:

- Penalizing high weights can prevent overfitting...

regularization → $\|\beta\|_2^2$ ★
 → $\|\beta\|_1$
 ↓

©Emily Fox 2013

7

Standard v. Regularized Updates

- Maximum conditional likelihood estimate

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

w/o regularization

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

- Regularized maximum conditional likelihood estimate

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

pushes toward 0

$$\beta_j^{(t+1)} \leftarrow \beta_j^{(t)} + \eta \left\{ \underbrace{-\lambda \beta_j^{(t)}}_{\text{pushes toward 0}} + \sum_i x_{ij} (y_i - \hat{p}(y=1 | x_i, \beta^{(t)})) \right\}$$

©Emily Fox 2013

8

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

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

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

©Emily Fox 2013

9

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:

- “True” gradient ascent rule:

- How do we estimate expected gradient?

©Emily Fox 2013

10

SGD: Stochastic Gradient Ascent (or Descent)

- “True” gradient: $\nabla l(\beta) = E_x[\nabla l(\beta, x)]$
- Sample based approximation:
- What if we estimate gradient with just one sample???
 - Unbiased estimate of gradient
 - Very noisy!
 - Called stochastic gradient ascent (or descent)
 - Among many other names
 - VERY useful in practice!!!

©Emily Fox 2013

11

Stochastic Gradient Ascent for Logistic Regression

- Logistic loss as a stochastic function:

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

- Batch gradient ascent updates:

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

- Stochastic gradient ascent updates:

- Online setting:

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

©Emily Fox 2013

12

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 2013

13

Module 5: Classification

Linear Methods: LDA and QDA

STAT/BIOSTAT 527, University of Washington

Emily Fox

May 28th, 2013

©Emily Fox 2013

14

Discriminative vs. Generative

- So far, we have considered modeling/fitting

$$p(Y | X)$$

- There are also a large set of **generative** methods
- Model:
 - Class-conditional densities $f_k(X) =$
 - Class prior probabilities π_k
- Via Bayes' rule:

©Emily Fox 2013

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)
 - Mixture of Gaussians (saw in BNP module)
 - Nonparametric density estimation for $f_k(x)$
 - Naïve Bayes

©Emily Fox 2013

Linear Discriminative Analysis

- Assume Gaussian class-conditional densities

$$f_k(X) =$$

- Furthermore, consider equal covariances

- Log odds

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

©Emily Fox 2013

Linear Discriminative Analysis

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

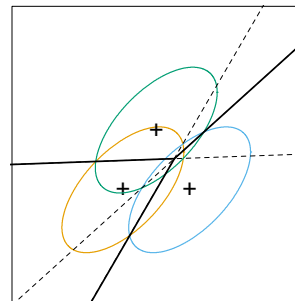
- Equivalently, $\log \frac{p(Y = k | X = x)}{p(Y = \ell | X = x)} = \delta_k(x) - \delta_\ell(x)$

where

$$\delta_k(x) =$$

- Decision rule:

- Linear decision boundaries



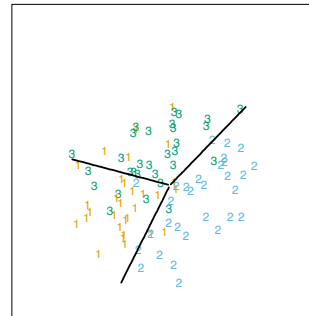
From Hastie, Tibshirani, Friedman book

©Emily Fox 2013

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, estimate parameters:



From Hastie, Tibshirani, Friedman book

©Emily Fox 2013

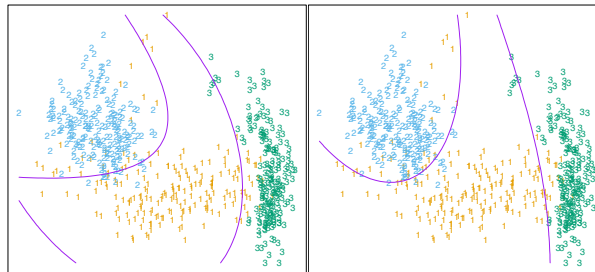
Quadratic Discriminative Analysis

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

- Quadratic discriminant functions:

$$\delta_k(x) =$$

- Quadratic decision boundaries



From Hastie, Tibshirani, Friedman book

©Emily Fox 2013

QDA Parameter Estimation

- Based on the training class labels, estimate parameters:

- Number of parameters:

- Can also consider shrinkage estimators

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

©Emily Fox 2013

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

©Emily Fox 2013

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

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

©Emily Fox 2013

LDA vs. Logistic Regression

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

- Marginal likelihood term

- Logistic regression:

- LDA:

©Emily Fox 2013

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 Σ in LDA, but are ignored by logistic regression
- Observations without class labels can be used in mixture model case, but not in logistic regression
- Marginal likelihood $p(X)$ acts as a regularizer

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

©Emily Fox 2013

Module 5: Classification

Nonparametric Methods: KDE and Naïve Bayes

STAT/BIOSTAT 527, University of Washington

Emily Fox

May 28th, 2013

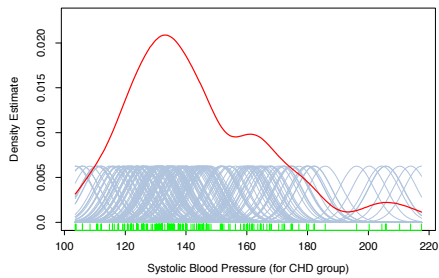
©Emily Fox 2013

26

KDE for Classification

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

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

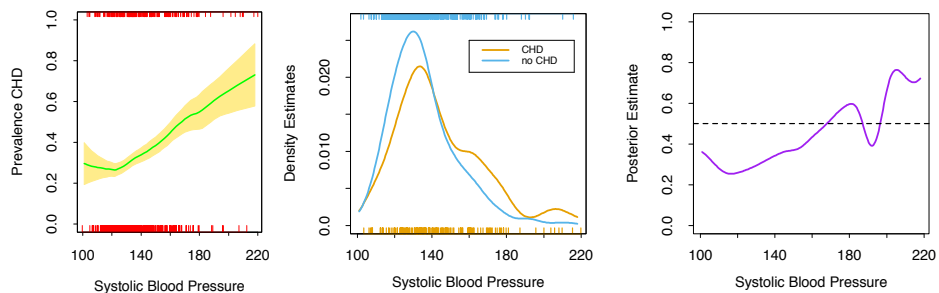


From Hastie, Tibshirani, Friedman book

©Emily Fox 2013

Example: Heart Disease Data

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

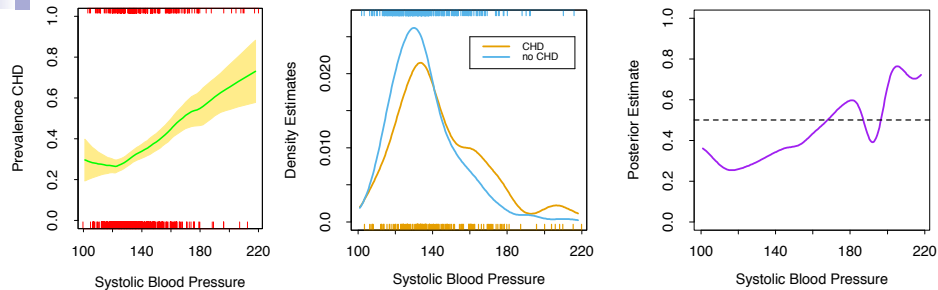


From Hastie, Tibshirani, Friedman book

©Emily Fox 2013

Example: Heart Disease Data

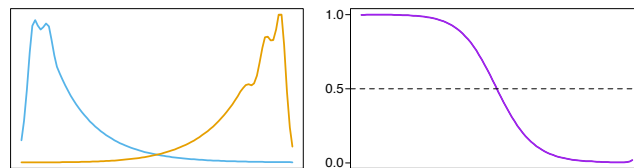
From Hastie, Tibshirani, Friedman book



- 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 2013

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 2013

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

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

- Typical choice = Gaussian RBF \rightarrow Gaussian KDE

lots of params to choose

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

©Emily Fox 2013

31

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

$$f_k(X) =$$

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

©Emily Fox 2013

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

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

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

©Emily Fox 2013

Module 5: Classification

Mixture Models for Classification

STAT/BIOSTAT 527, University of Washington

Emily Fox

May 28th, 2013

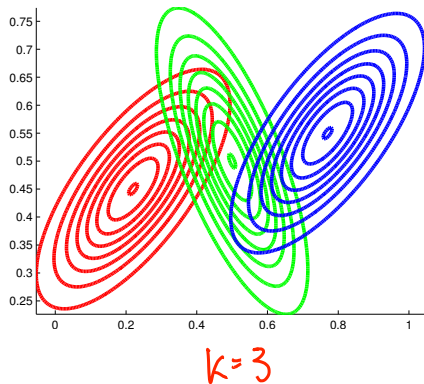
©Emily Fox 2013

34

Density as Mixture of Gaussians

- Approximate density with a mixture of Gaussians

Mixture of 3 Gaussians



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

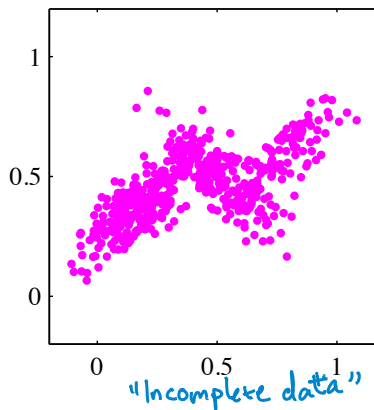
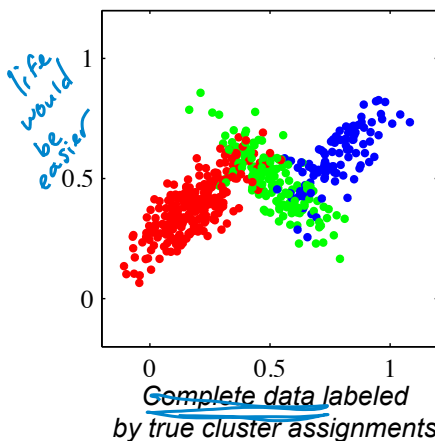
Handwritten notes:
 Gauss. kernel, just like in KDE, but not centered at obs.
 $\pi = [\pi_1, \dots, \pi_K]$
 $\mu = \{\mu_1, \dots, \mu_K\}$
 $\Sigma = \{\Sigma_1, \dots, \Sigma_K\}$
 K : # of mix comp.
 π_k : mix. weights
 μ_k, Σ_k : shape params
 In 1D: $P = \text{target density}$
 $\sum \pi_k = 1$

©Emily Fox 2013

Clustering our Observations

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

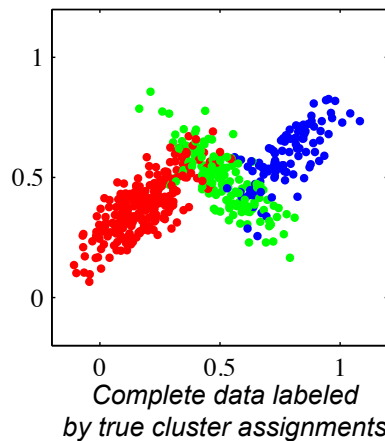
Our actual observations



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

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

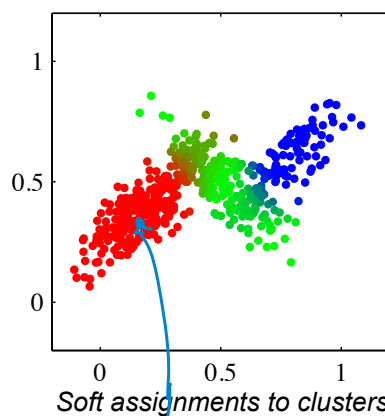
$$N(x_i | \mu_{z_i}, \Sigma_{z_i})$$

param. est. is easy if we have $\{z_i\}$
 \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:

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

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

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

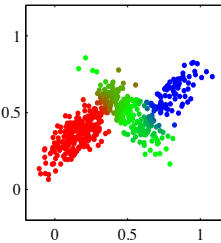
$$\hat{r}_{ik} = \frac{\hat{\pi}_k N(x_i, \hat{\mu}_k, \hat{\Sigma}_k)}{\sum_{\ell} \hat{\pi}_{\ell} N(x_i, \hat{\mu}_{\ell}, \hat{\Sigma}_{\ell})}$$

- Maximize parameters given responsibilities

- For classification, threshold the estimated responsibilities

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

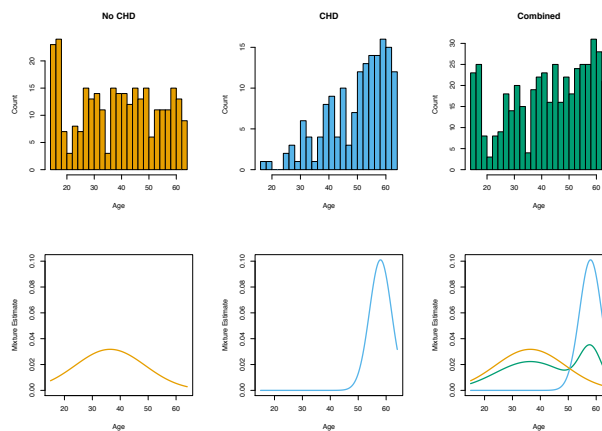
- Note: allows non-linear boundaries as in QDA



©Emily Fox 2013

Example: Heart Disease Data

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



From Hastie, Tibshirani, Friedman book

©Emily Fox 2013

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 2013

41

Readings

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

©Emily Fox 2013

42