

Linear classification methods

Lecture 2a



Overview

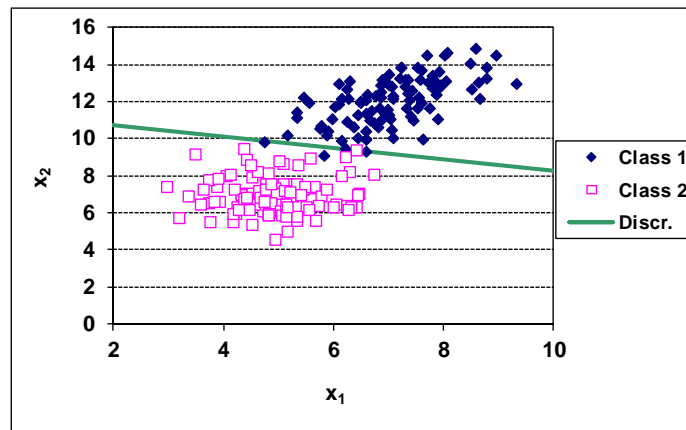
- Discriminant Analysis models
- Logistic regression
- Generalized Linear Model

Classification

- Given data $D = ((X_i, Y_i), i = 1 \dots N)$
 - $Y_i = Y(X_i) = C_j \in \mathcal{C}$
 - Class set $\mathcal{C} = (C_1, \dots, C_K)$

Classification problem:

- Decide $\hat{Y}(x)$ that maps **any** x into some class C_K
 - Decision boundary



Classifiers

- **Deterministic**: decide a rule that directly maps X into \hat{Y}
- **Probabilistic**: define a model for $P(Y = C_i|X), i = 1 \dots K$

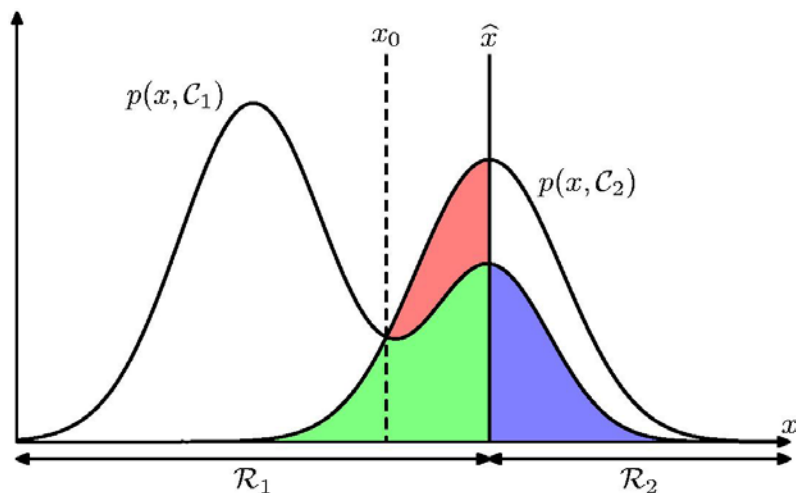
Disadvantages of deterministic classifiers:

- Sometimes simple mapping is not enough (risk of cancer)
- Difficult to embed loss \rightarrow rerun of optimizer is often needed
- Combining several classifiers into one is more problematic
 - Algorithm A classifies as spam, Algorithm B classifies as not spam \rightarrow ???
 - $P(\text{Spam} | A) = 0.99, P(\text{Spam} | B) = 0.45 \rightarrow$ better decision can be made

Probabilities into decision

- Loss minimization

$$\min_{\hat{f}} EL(y, \hat{f}) = \min_{\hat{f}} \int L(y, \hat{f}) p(y, x|w) dx dy$$



When loss is

$\begin{cases} 1, \text{wrongly classified} \\ 0, \text{correctly classified} \end{cases}$

Classify Y as

$$\hat{Y} = \arg \max_c p(Y = c|X)$$

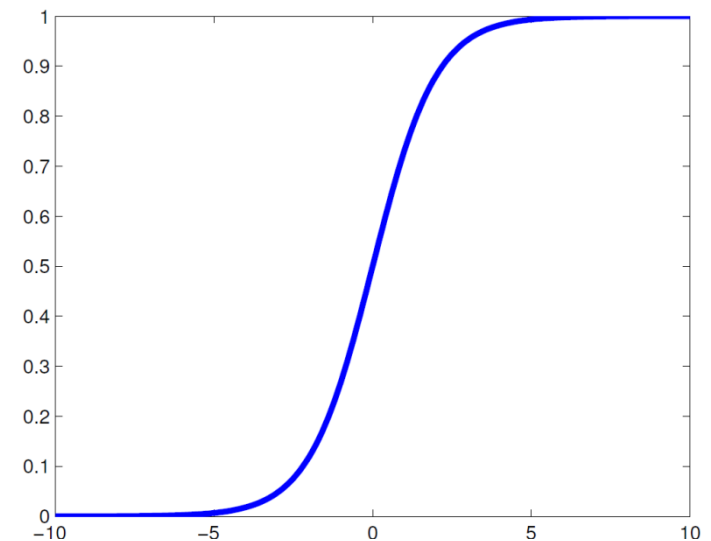
Logistic regression

- Discriminative model
- Model for binary output
 - $C = \{C_1 = 1, C_2 = 0\}$
 $p(Y = C_1|X) = \text{sigm}(\mathbf{w}^T \mathbf{x})$

$$\text{sigm}(a) = \frac{1}{1 + e^{-a}}$$

- Alternatively
 $Y \sim \text{Bernoulli}(\text{sigm}(a)), a = \mathbf{w}^T \mathbf{x}$
$$\text{sigm}(a) = \frac{1}{1 + e^{-a}}$$

What is $P(Y = C_2|X)$?



Logistic regression

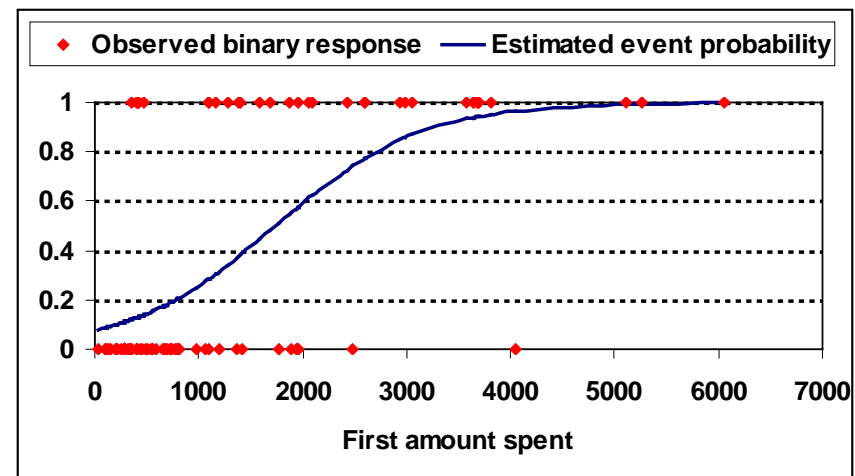
- Logistic model- yet another form

$$\ln \frac{p(Y = 1|X = x)}{P(Y = 0|X = x)} = \ln \frac{p(Y = 1|X = x)}{1 - P(Y = 1|X = x)} = \text{logit}(p(Y = 1|X = x)) = \mathbf{w}^T \mathbf{x}$$

**The log of the odds
is linear in \mathbf{x}**

- Here $\text{logit}(t) = \ln \left(\frac{t}{1-t} \right)$
- Note $p(Y|X)$ is connected to $\mathbf{w}^T \mathbf{x}$ via logit link

Example: Probability to buy
more than once as function of
First Amount Spend



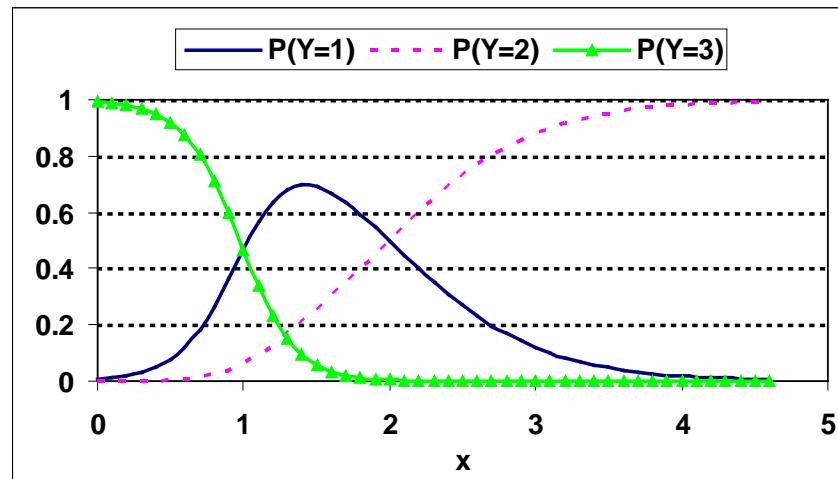
Logistic regression

- When Y is categorical,

$$p(Y = C_i | x) = \frac{e^{w_i^T x}}{\sum_{j=1}^K e^{w_j^T x}} = \text{softmax}(w_i^T x)$$

- Alternatively

$$Y \sim \text{Multinoulli}(\text{softmax}(w_1^T x), \dots, \text{softmax}(w_K^T x))$$



Logistic regression

Fitting logistic regression

- In binary case,

$$\log P(D|w) = \sum_{i=1}^N y_i \log(\text{sigm}(w^T x_i)) + (1 - y_i) \log(1 - \text{sigm}(w^T x_i))$$

- Can not be maximized analytically, but unique maximizer exists
- To maximize loglikelihood, optimization used
 - Newton's method traditionally used (Iterative Reweighted Least Squares)
 - Steepest descent, Quasi-newton methods...

Estimation:

For new x , estimate $p(y) = [p_1, \dots, p_c]$ and classify as $\arg \max_i p_i$

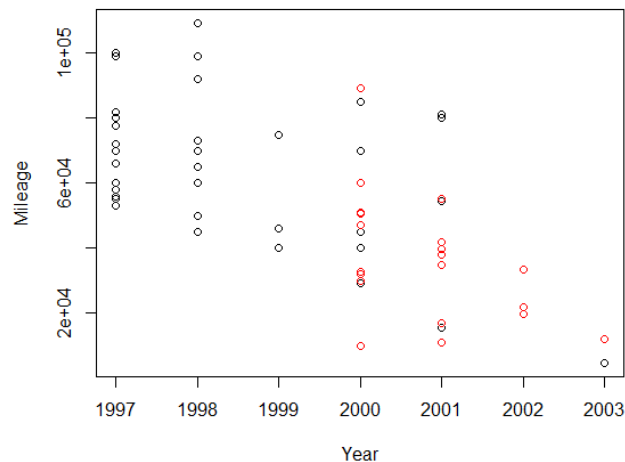
Decision boundaries of logistic regression are linear

Logistic regression

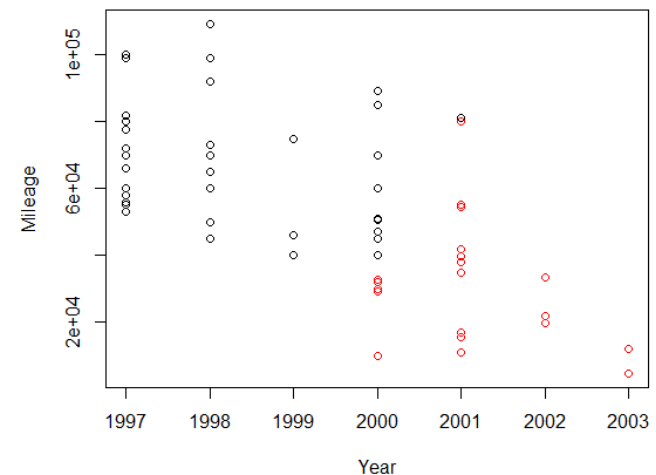
- In R, use `glm()` with `family="binomial"`
 - Predicted probabilities: `predict(fit,newdata, type="response")`

Example Equipment=f(Year, mileage)

Original data



Classified data



Moving beyond typical distributions

- We know how to model
 - Normally distributed targets -> linear regression
 - Bernoulli and Multinomial targets → logistic regression
 - What if target distribution is more complex?

Example 1: Daily Stock prices NASDAQ

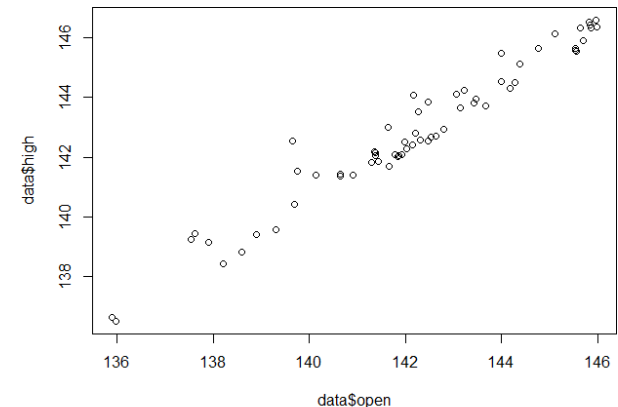
- Open
- High (within day)

Does it seem that the error is normal here?

Example 2: Number of calls to bank

- Y = Number of calls
- X = time

Endless amount of classes → multinomial does not work... (Poisson)



Exponential family

- More advanced error distributions are sometimes needed!
- Many distributions belong to **exponential** family:
 - Normal, Exponential, Gamma, Beta, Chi-squared..
 - Bernoulli, Multinoulli, Poisson...

$$p(\mathbf{x}|\boldsymbol{\eta}) = h(\mathbf{x})g(\boldsymbol{\eta})e^{(\boldsymbol{\eta}^T u(\mathbf{x}))}$$

- Easy to find MLE and MAP
- Non-exponential family distributions: uniform, Student t

Example: Bernoulli

Generalized linear models

- Assume Y from the exponential family
- **Model** is $Y \sim EF(\mu, \dots)$, $f(\mu) = \mathbf{w}^T \mathbf{x}$
 - Alt $\mu = f^{-1}(\mathbf{w}^T \mathbf{x})$
 - f^{-1} is activation function
 - f is link function (in principle, arbitrary)
- Arbitrary f will lead to (s – dispersion parameter)

$$p(y|w, s) = h(y, s)g(\mathbf{w}, \mathbf{x})e^{\frac{b(\mathbf{w}, \mathbf{x})y}{s}}$$

- If f is a canonical link, then

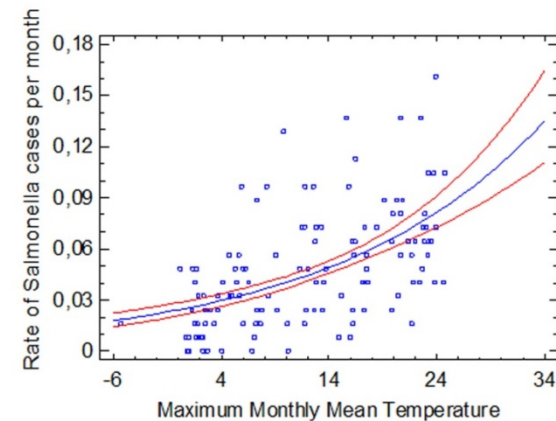
$$p(y|w, s) = h(y, s)g(\mathbf{w}, \mathbf{x})e^{\frac{(\mathbf{w}^T \mathbf{x})y}{s}}$$

Generalized linear models

- Canonical links are normally used
 - MLE computations simplify
 - MLE $\hat{w} = F(X^T Y) \rightarrow$ computations do not depend on all data but rather a summary (sufficient statistics) \rightarrow computations speed up

Example: Poisson regression

$$f^{-1}(\mu) = e^{\mu}, Y \sim \text{Poisson}(e^{w^T x})$$



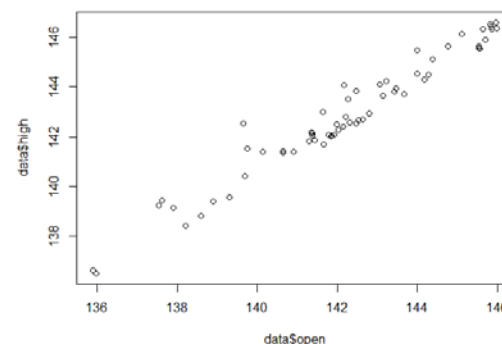
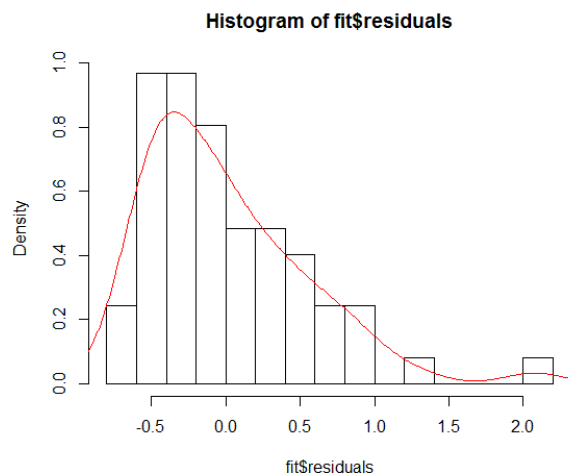
Generalized linear model: software

- Use `glm(formula, family, data)` in R

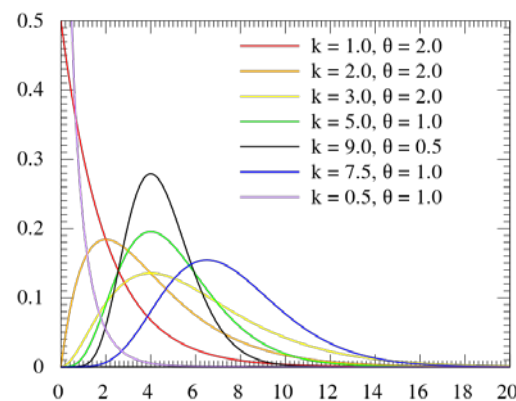
Example: Daily Stock prices NASDAQ

- Open
- High (within day)

1. Try to fit usual linear regression, study histogram of residuals



Gamma distribution: Wikipedia



Generalized linear model

Assume

$$High \sim \text{Gamma}\left(1, \frac{1}{w_0 + w_1 \text{Open}}\right)$$

What is link function
here?

```
call:
glm(formula = high ~ open, family = Gamma(link = "inverse"),
    data = data)
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-0.0052879	-0.0028896	-0.0006678	0.0016598	0.0148083

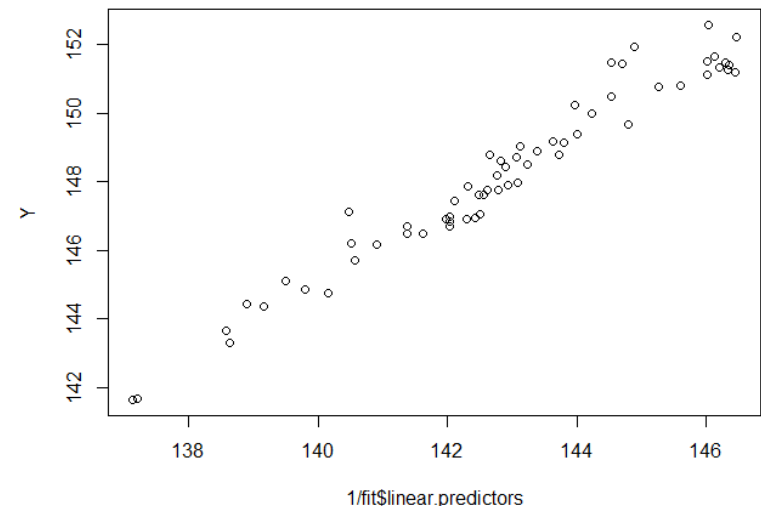
Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	1.355e-02	1.962e-04	69.06	<2e-16 ***
open	-4.604e-05	1.379e-06	-33.39	<2e-16 ***

signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

New generated data

- Has similar pattern as original data!

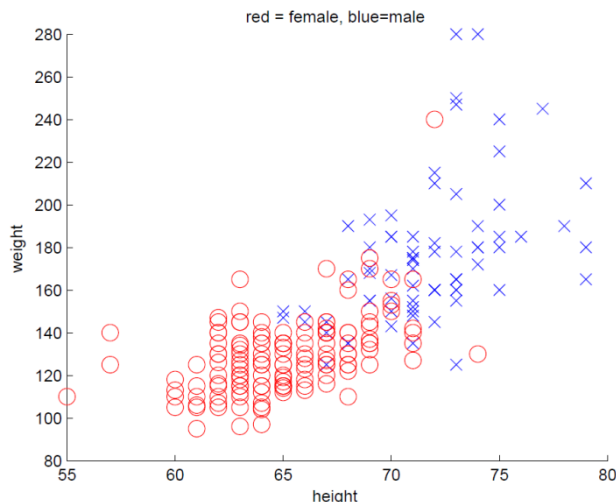
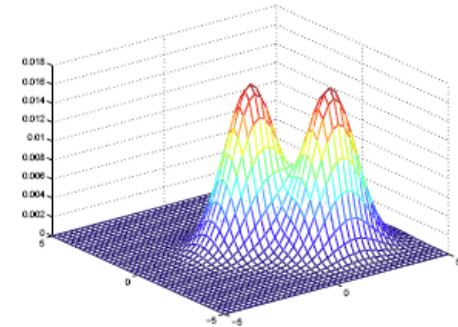


Quadratic discriminant analysis

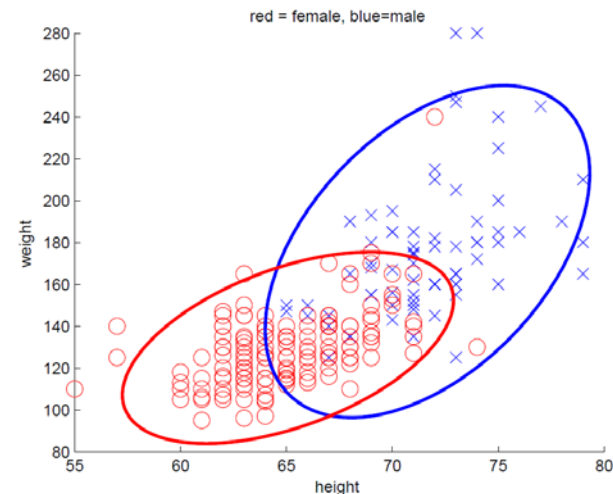
- Generative classifier
- Main assumptions:
 - \mathbf{x} is now **random** as well as y

$$p(\mathbf{x}|y = C_i, \theta) = N(\mathbf{x}|\mu_i, \Sigma_i)$$

Unknown parameters $\theta = \{\mu_i, \Sigma_i\}$



Source: Probabilistic Machine Learning by Murphy

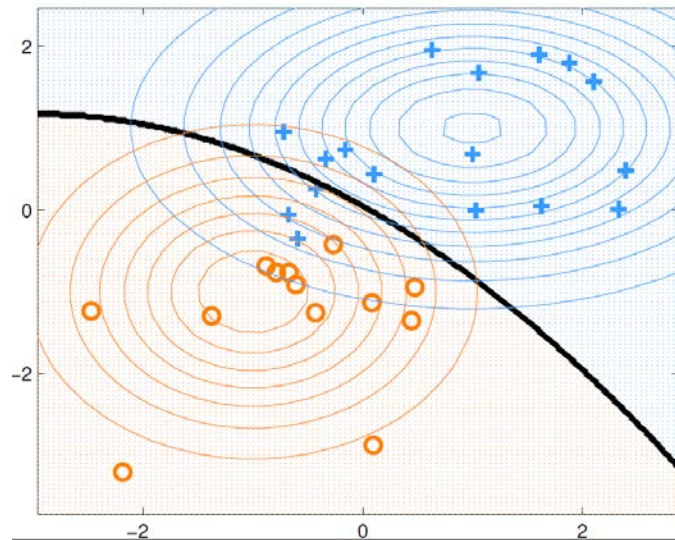


Source: Probabilistic Machine Learning by Murphy

Quadratic discriminant analysis

- If parameters are estimated, classify:

$$\hat{y}(\mathbf{x}) = \arg \max_c p(y = c | \mathbf{x}, \theta)$$

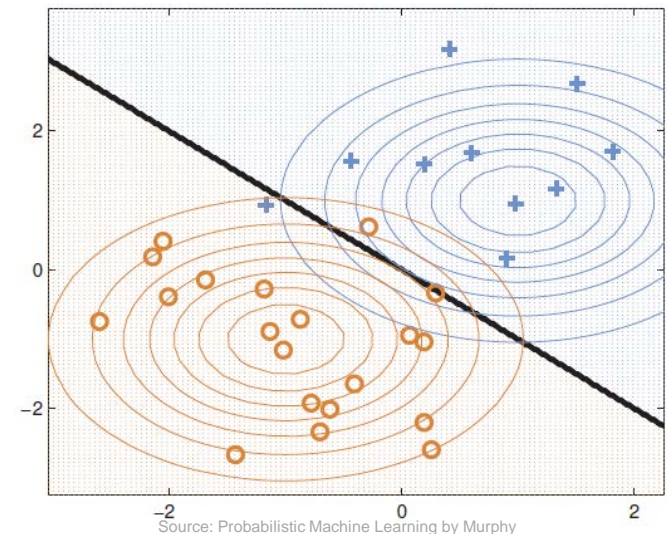


Source: Probabilistic Machine Learning by Murphy

Linear discriminant analysis (LDA)

- Assumption $\Sigma_i = \Sigma, i = 1, \dots, K$
- Then $p(y = c_i | x) = \text{softmax}(w_i^T x + w_{0i}) \rightarrow$ exactly the same form as the logistic regression
 - $w_{0i} = -\frac{1}{2} \mu_i^T \Sigma^{-1} \mu_i + \log \pi_i$
 - $w_i = \Sigma^{-1} \mu_i$
- Decision boundaries are linear
 - **Discriminant function:**

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



Linear discriminant analysis (LDA)

- Difference LDA vs logistic regression??
 - Coefficients will be estimated differently! (models are different)
- How to estimate coefficients
 - find MLE.

$$\hat{\mu}_c = \frac{1}{N_c} \sum_{i:y_i=c} \mathbf{x}_i, \quad \hat{\Sigma}_c = \frac{1}{N_c} \sum_{i:y_i=c} (\mathbf{x}_i - \hat{\mu}_c)(\mathbf{x}_i - \hat{\mu}_c)^T$$

$$\hat{\Sigma} = \frac{1}{N} \sum_{c=1}^k N_c \hat{\Sigma}_c$$

- Sample mean and sample covariance are MLE!
- If class priors are parameters (**proportional priors**),

$$\hat{\pi}_c = \frac{N_c}{N}$$

LDA and QDA: code

- Syntax in R, library **MASS**

`lda(formula, data, ..., subset, na.action)`

- Prior – class probabilities
- Subset – indices, if training data should be used

`qda(formula, data, ..., subset, na.action)`

`predict(..)`

LDA: output

```
resLDA=lda(Equipment~Mileage+Year, data=mydata)
print(resLDA)
```

```
> print(resLDA)
Call:
lda(Equipment ~ Mileage + Year, data = mydata)

Prior probabilities of groups:
      0      1
0.6440678 0.3559322

Group means:
      Mileage      Year
0 63539.21 1998.447
1 36857.62 2000.762

Coefficients of linear discriminants:
              LD1
Mileage -1.500069e-05
Year      5.745893e-01
```

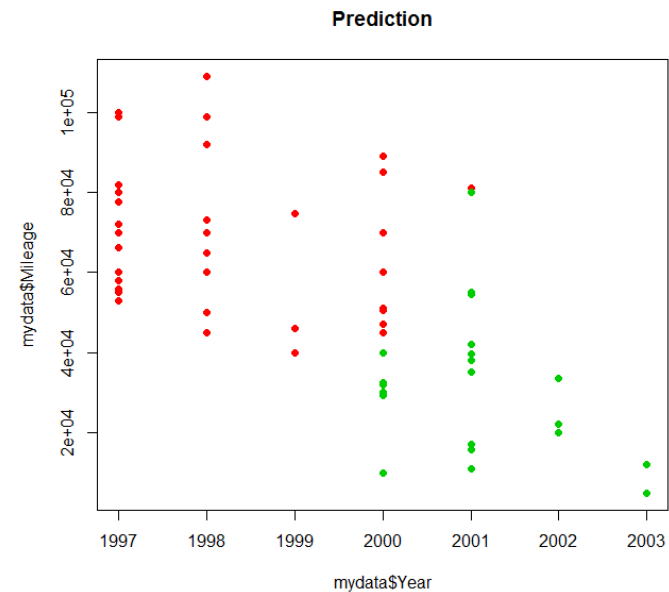
LDA: output

- Misclassified items

```
plot(mydata$Year, mydata$Mileage,  
col=as.double(Pred$class)+1, pch=21,  
bg=as.double(Pred$class)+1,  
main="Prediction")
```

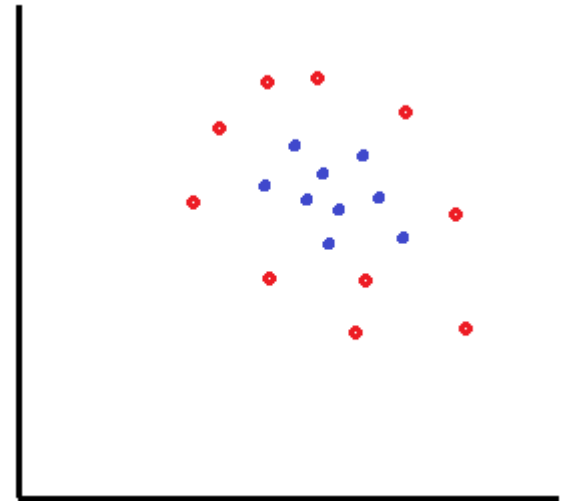
```
> table(Pred$class, mydata$Equipment)
```

```
      0  1  
0 31  6  
1  7 15
```



LDA versus Logistic regression

- Generative classifiers are easier to fit, discriminative involve numeric optimization
- LDA and Logistic have same model form but are fit differently
- LDA has stronger assumptions than Logistic, some other generative classifiers lead also to logistic expression
- New class in the data?
 - Logistic: fit model again
 - LDA: estimate new parameters from the new data
- Logistic and LDA: complex data fits badly unless interactions are included



LDA versus Logistic regression

- LDA (and other generative classifiers) handle missing data easier
- Standardization and generated inputs:
 - Not a problem for Logistic
 - May affect the performance of the LDA in a complex way
- Outliers affect $\Sigma \rightarrow$ LDA is not robust to gross outliers
- LDA is often a good classification method even if the assumption of normality and common covariance matrix are not satisfied.