# Machine Learning

# Week 2: Pattern Classification

Mahesan Niranjan

School of Electronics and Computer Science
University of Southampton

Autumn Semester 2016/17

# Overview (Week 2)

- Review of what we learnt in Lab One
  - Multivariate Gaussian
  - Drawing samples from $\mathcal{N}(\boldsymbol{m}, \boldsymbol{C})$
  - Principal directions
- Introduction to Bayesian Decision Theory
- Bayes' Classifier for Simple Gaussian Distributions
- Simple Classifiers
  - Distance to mean classifier
  - Nearest Neighbour classifier
  - Linear classifier (more on this later)
  - Perceptron (formal setting later)
- What will we learn in Lab Two?

# Bayesian Decision Theory

- Classes: $\omega_i$, $i = 1, ..., K$
- Prior Probabilities: $P[\omega_1], ..., P[\omega_K]$;
  $P[\omega_i] \geq 0$, $\sum_{i=1}^{C} P[\omega_i] = 1$
- Likelihoods (class conditional probabilities): $p(\mathbf{x}|\omega_i)$, $i = 1, .., K$
- Posterior Probability: $P[\omega_j | \mathbf{x}]$

$$P[\omega_j | \mathbf{x}] = \frac{p(\mathbf{x}|\omega_j) \; P[\omega_j]}{\sum_{i=1}^{K} p(\mathbf{x}|\omega_i) \; P[\omega_i]}$$

- From prior knowledge: $P[\omega_i]$; From traing data: $p(\mathbf{x}|\omega_i)$
- Decision rule: Assign $\mathbf{x}$ to the class that maximizes posterior probability.
- The denominator is a constant; *i.e.* does not depend on $\omega_j$
- Hence the decision rule becomes:

$$\mathbf{x} \in \max_{j} p(\mathbf{x}|\omega_j) \; P[\omega_j]$$

# Bayes' Classifier for Gaussian Densities
Make assumptions, cancel common terms when making comparisons...

- Decision rule from: $p(\mathbf{x}|\omega_j) \; P[\omega_j]$
- Assume the two classes are Gaussian distributed with distinct means and identical covariance matrices
  $p(\mathbf{x}|\omega_j) = \mathcal{N}(\mathbf{m}_j, \mathbf{C})$
- Substitute into Bayes' classifier decision rule

$$P[\omega_1|\mathbf{x}] \lessgtr P[\omega_2|\mathbf{x}]$$
$$p(\mathbf{x}|\omega_1) \; P[\omega_1] \lessgtr p(\mathbf{x}|\omega_2) \; P[\omega_2]$$

$$\frac{1}{(2\pi)^{p/2}(\det(C))^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \mathbf{m}_1)^t \mathbf{C}^{-1}(\mathbf{x} - \mathbf{m}_1)\right\} P[\omega_1] \lessgtr$$
$$\frac{1}{(2\pi)^{p/2}(\det(C))^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \mathbf{m}_2)^t \mathbf{C}^{-1}(\mathbf{x} - \mathbf{m}_2)\right\} P[\omega_2]$$

# Bayes' classifier for simple densities (cont'd)
Distinct Means; Equal, isotropic covariance matrix

- Suppose the densities are isotropic and priors are equal
  *i.e.* $C = \sigma^2 I$ and $P[\omega_1] = P[\omega_2]$
- The comparison simplifies to (see algebra on board):

$$(x - m_1)^t(x - m_1) \lesseqgtr (x - m_2)^t(x - m_2)$$
$$|x - m_1| \lesseqgtr |x - m_2|$$

- The above is a simple *distance to mean* classifier
- Under the above simplistic assumptions, we only need to store one template per class (the means)!

# Bayes' classifier for simple densities (cont'd)
Distinct Means; Common covariance matrix (but not isotropic)

- Cancel common terms and take log

$$(x - m_1)^t C^{-1} (x - m_1) \lesseqgtr (x - m_2)^t C^{-1} (x - m_2) - \log \left\{ \frac{P[\omega_1]}{P[\omega_2]} \right\}$$

- Also simplifies to a $\boxed{\text{linear classifier}}$ !

$$w^t x + b \lesseqgtr 0$$

$$w = 2C^{-1} (m_2 - m_1)$$
$$b = \left( m_1^t C^{-1} m_1 - m_2^t C^{-1} m_2 \right) - \log \left\{ \frac{P[\omega_1]}{P[\omega_2]} \right\}$$

- Also a distance to template classifier, where the distance is

$$(x - m_1)^t C^{-1} (x - m_1)$$

Known as Mahalanobis distance

Linear classifier

$$\boldsymbol{w}^t \boldsymbol{x} + b \lesseqgtr 0$$

Expand dimensions: $\boldsymbol{a} = [\boldsymbol{w}^t \, b]^t$ and $\boldsymbol{y} = [\boldsymbol{x}^t \, 1]^t$

$$\boldsymbol{a}^t \boldsymbol{y} \lesseqgtr 0$$

```
random guess of the weights
repeat
    select data at random
    if not correctly classified
        update weights
until (all data correctly classified)
```

Update:

$$\boldsymbol{a}^{(k+1)} = \boldsymbol{a}^{(k)} + \eta \, \boldsymbol{y}^{(k)}$$

# Lab 2

1. Some plotting
2. Bayes' optimal class boundary
3. Implement your own perceptron algorithm