# Introduction to Machine Learning and Pattern Recognition
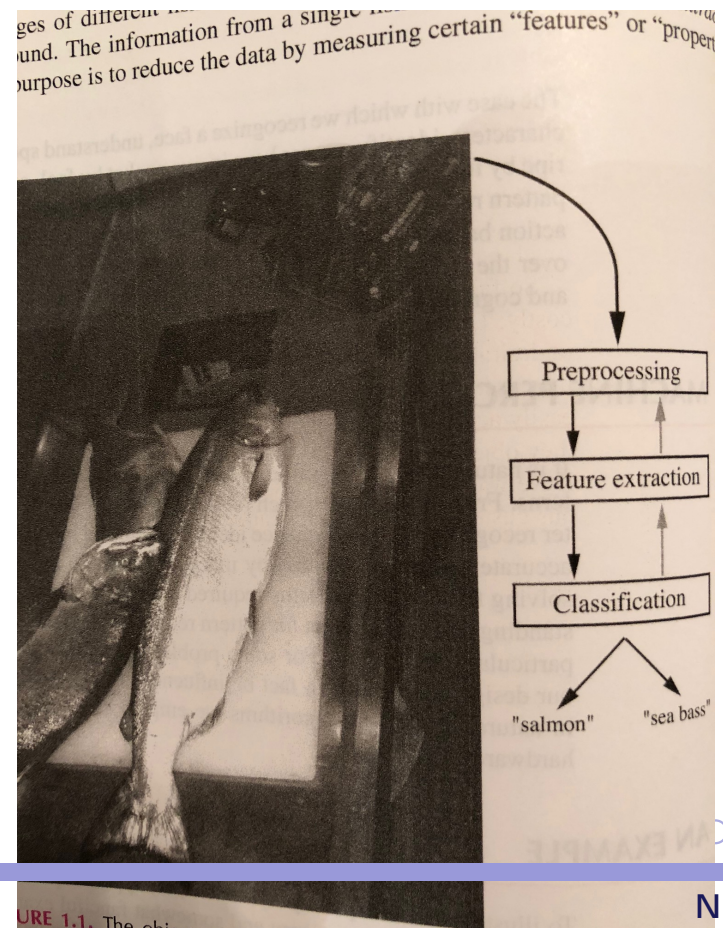
## David Brady[1]

[1]ECE Department
Northeastern University

Fall 2018

# The Classification Problem

Solutions to classification problems involve:

- ▶ data sensing

  - ▶ equipment, installation, data storage and retrieval

- ▶ preprocessing

  - ▶ segmentation (isolation of fish)
  - ▶ grouping (tail + torso + head)

- ▶ feature extraction

  - ▶ invariance to scale, 3D rotation, etc.
  - ▶ length, weight, lightness, etc.

- ▶ classification design: (this course)

  - ▶ bass or salmon?
  - ▶ use feature space ! (next)
  - ▶ missing data



ges of different
und. The information from a single
ourpose is to reduce the data by measuring certain "features" or "propert

Preprocessing

Feature extraction

Classification

"salmon"     "sea bass"

URE 1.1. The objects

# Decision Boundary in Feature Space

Placing the decision boundary depends on:

- ▶ cost of each type of error

- ▶ *prior* probabilities

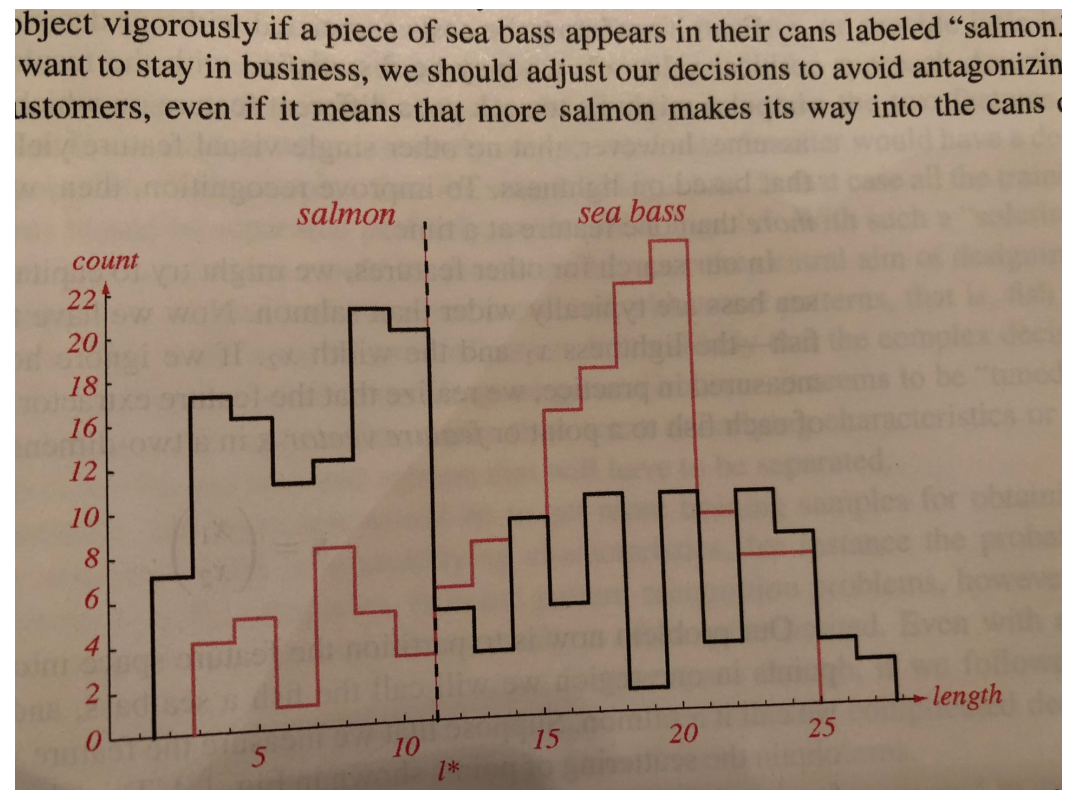Why is *length* a poor feature for distinguishing fish?



Figure: Histograms for length feature.

# Feature Selection

Feature selection depends on:

- ▶ sensing noise
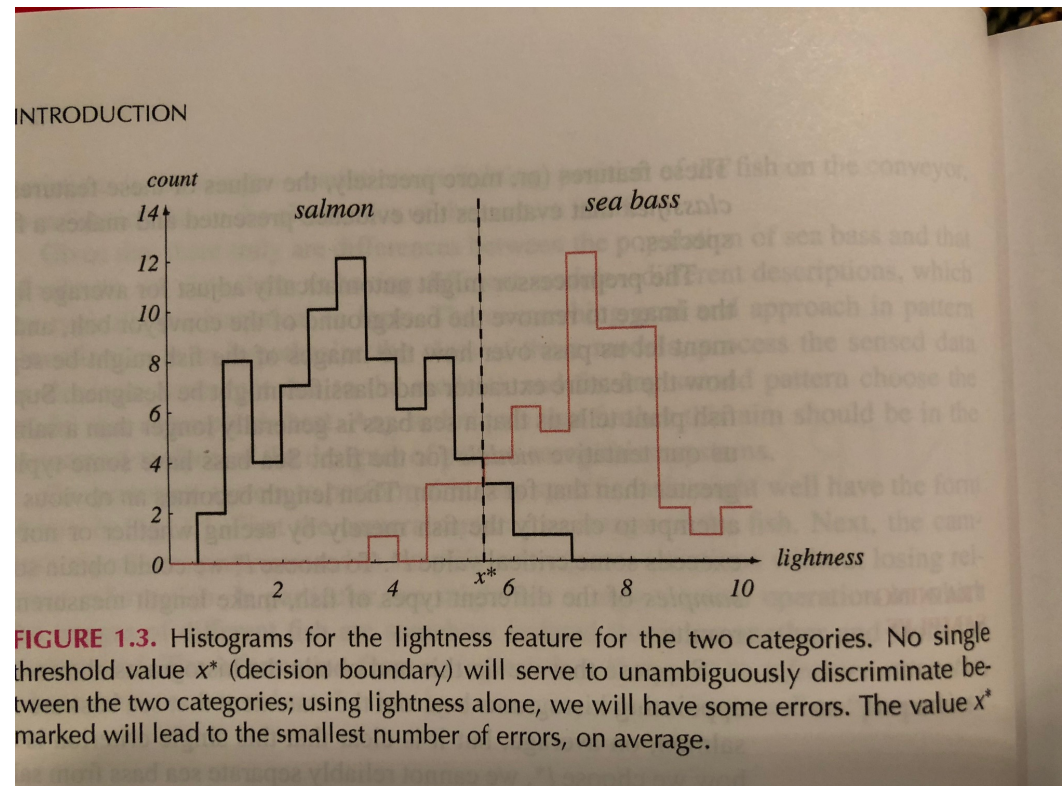- ▶ classification problem
- ▶ class invariance
- ▶ training samples



Figure: Histograms for lightness feature.

# >2D Feature Space

Additional features provide:

- ► discrimination
- ► information
- ► *complexity*
- ► *increased training size*
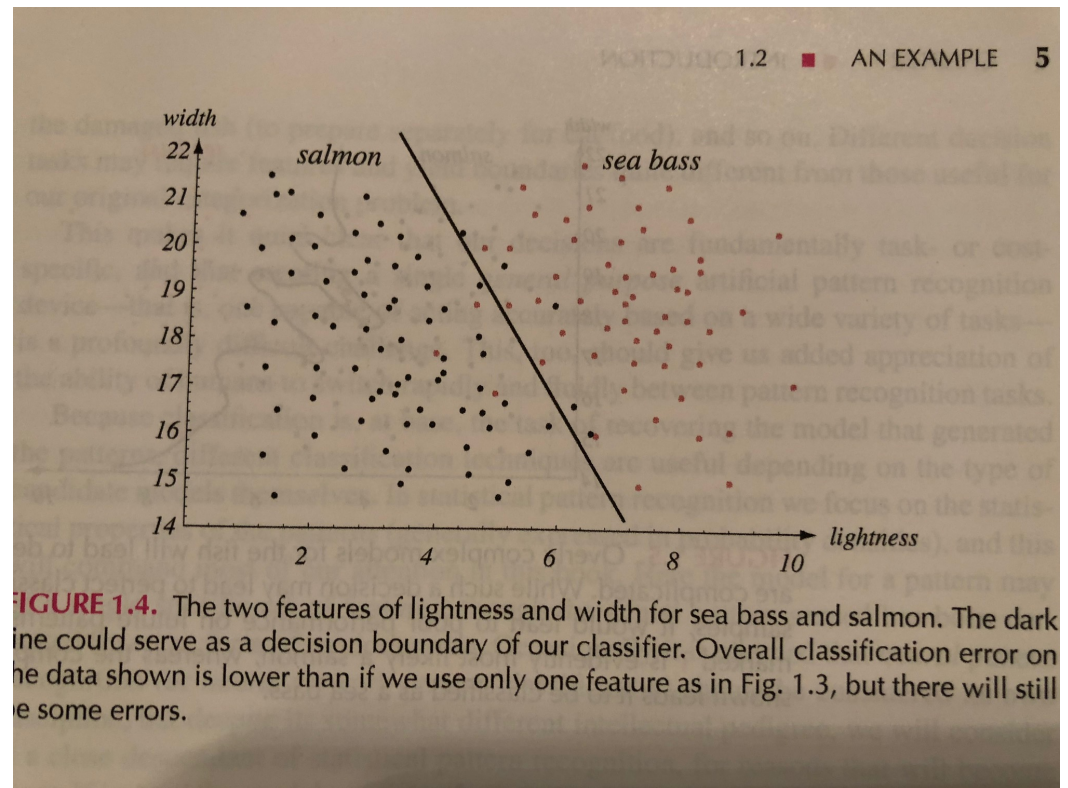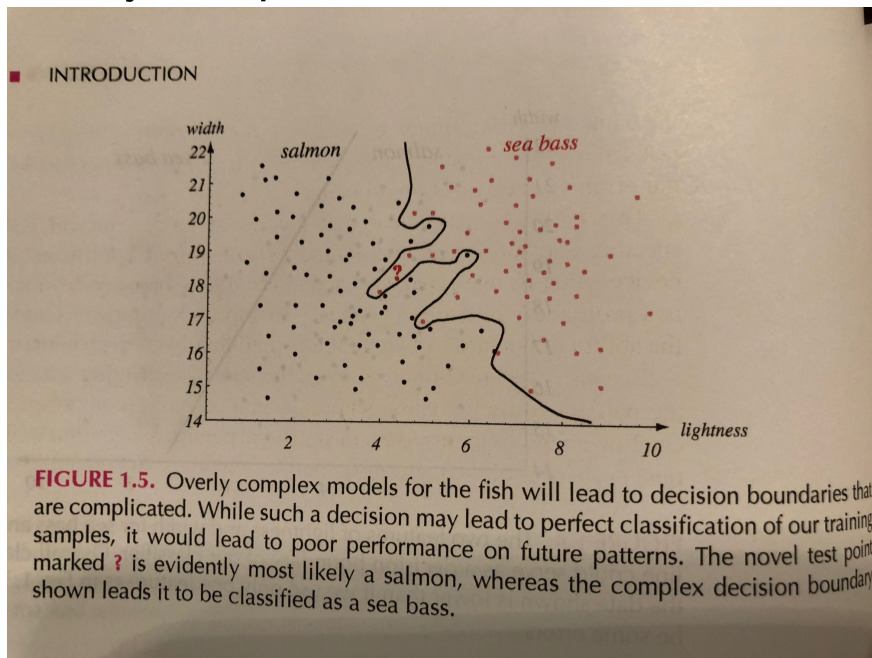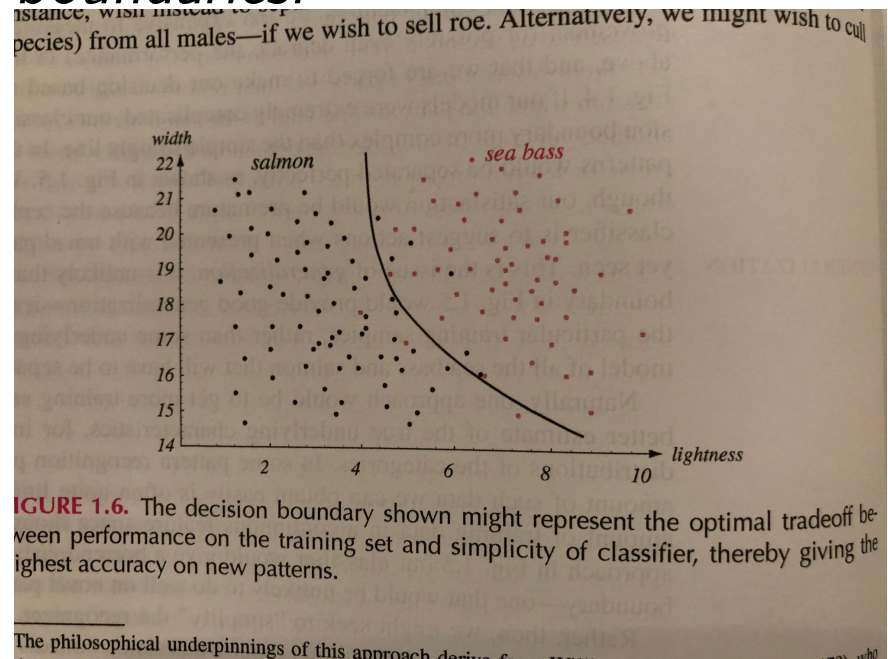
How should you select features?



Figure: 2D feature space with a decision boundary

# Overfitting

Overly-complex boundaries follow noise!



FIGURE 1.5. Overly complex models for the fish will lead to decision boundaries that are complicated. While such a decision may lead to perfect classification of our training samples, it would lead to poor performance on future patterns. The novel test point marked ? is evidently most likely a salmon, whereas the complex decision boundary shown leads it to be classified as a sea bass.

*Use the simplest possible boundaries.*



FIGURE 1.6. The decision boundary shown might represent the optimal tradeoff between performance on the training set and simplicity of classifier, thereby giving the highest accuracy on new patterns.

The philosophical underpinnings of this approach derive from Willi...
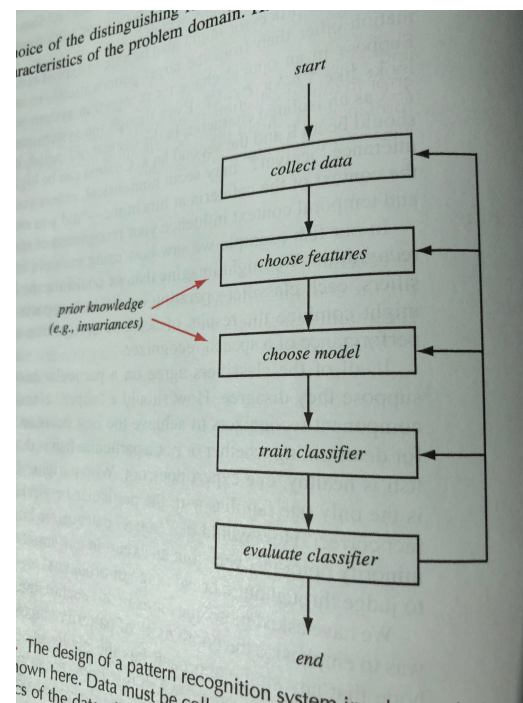
# Classifier Design Cycle

Design Tweaking Issues: (improve performance, reduce complexity)

- ▶ data collection
  - ▶ relevant data for features
- ▶ feature choice
  - ▶ retain fewest
- ▶ model choice
  - ▶ connects features & hypotheses
  - ▶ prefer simpler models
- ▶ training
  - ▶ balance *overfitting* with *typicality*



The design of a pattern recognition system shown here. Data must be coll... cs of the data...

- ▶ evaluation
  - ▶ criteria (error rates, costs)
  - ▶ confidence bounds on criteria (training size)
- ▶ computational complexity — as much as your budget can handle!