

Introduction to Data Mining

Deng Cai (蔡登)

College of Computer Science
Zhejiang University

dengcai@gmail.com



Introduction to Machine Learning

Deng Cai (蔡登)

College of Computer Science
Zhejiang University

dengcai@gmail.com





Short Bio

- ▶ Dr. Deng Cai (蔡登)
 - dengcai@gmail.com, dengcai@cad.zju.edu.cn
- ▶ Professor at CS college (the state key lab of CAD&CG).
 - 紫金港校区蒙民伟楼508
- ▶ Research interests:
 - Machine learning
 - Data mining
 - Computer vision
 - ...
- ▶ <http://dengcai.zjulearning.org:8081/>



Course Information

- ▶ Web: <http://dengcai.zjulearning.org:8081/Courses/DM/>
- ▶ Homework: <http://assignment.zjulearning.org:8081/>
 - 缺省用户名和密码：学号，登陆之后修改密码
- ▶ Time:
 - **Monday, 14:05 – 15:35**
 - **Thursday, 14:05 – 15:35**
- ▶ Place: Room 504, 7th teaching building, Yuquan Campus
- ▶ QQ group: 397340601(DM_ZJU) (**Apply with name and student ID**)
- ▶ TA: 张永辉、胡津铭



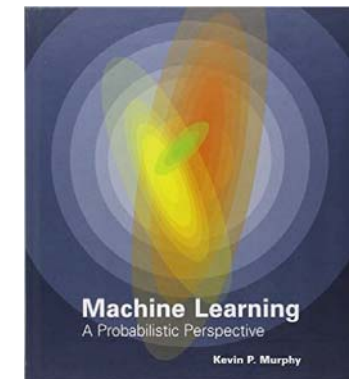
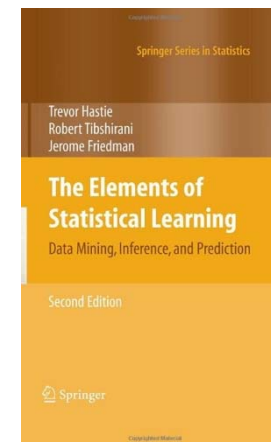
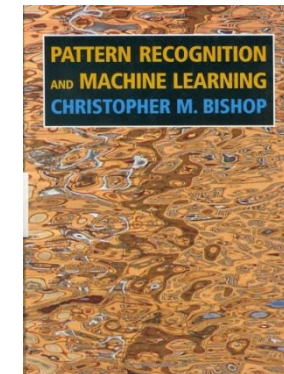
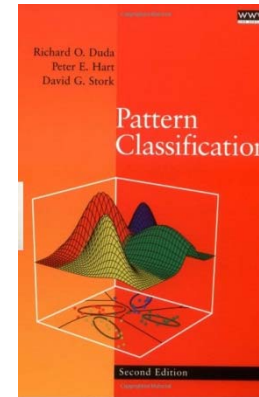
Course information (Cont'd)

- ▶ Prerequisite:
 - Linear algebra, analysis, probability theory
 - Basic programming skills
- ▶ Course textbook: No textbook is required. (Papers and other materials are available at the class web page)
- ▶ Objective:
 - Basic understandings of some of the important machine learning methods.
 - Basic ability to use some machine learning techniques to solve real world problems.



Reference Books

- ▶ R. Duda, P. Hart & D. Stork, *Pattern Classification* (2nd ed.), Wiley, 2000
- ▶ C. M. Bishop, *Pattern Recognition and Machine Learning*, Springer, 2006
- ▶ T. Hastie, R. Tibshirani & J. Friedman, *The Elements of Statistical Learning: Data Mining, Inference, and Prediction* (2nd ed.), Springer, 2009
- ▶ Kevin Murphy, *Machine Learning: A Probabilistic Perspective*, The MIT Press, 2012





Reference Books

- ▶ You can download all the books from the QQ group



Evaluation

- ▶ Quizzes (15%)
- ▶ Four assignments (10% each)
 - Everyone do it by himself
- ▶ Final exam (45%)
- ▶ Programming language:
 - Matlab
 - Tutorials
 - <http://www.math.ufl.edu/help/matlab-tutorial/>
 - <http://www.math.mtu.edu/~msgocken/intro/node1.html>
 - Python



Course Policies

- ▶ Class
 - No laptop, no cellphone.
- ▶ Cheating
 - No.
- ▶ Homework:
 - You have to write you own solution/program.
- ▶ Late Policy:
 - 0~24 hours: 90%
 - 24~48 hours: 50%
 - 48 hours ~: 25%
- ▶ Questions?



Why Take This Course?

- ▶ It is NOT
 - Easy course with high scores
 - Recommendation letter for US school application
 - Rank 1st
- ▶ You should
 - Work hard
 - Be honest



What is machine learning?

- ▶ Machine learning is the study of computer systems that improve their performance through experience.
 - Learn existing and known structures and rules.
 - Discover new findings and structures.
 - Face recognition
 - News summarization
- ▶ In machine learning, we study two types of problems



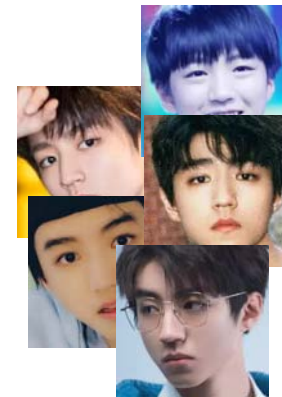
The first kind of problems



刘德华



章子怡



王俊凯

.....



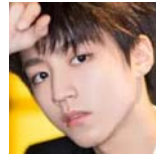
章子怡



The first kind of problems



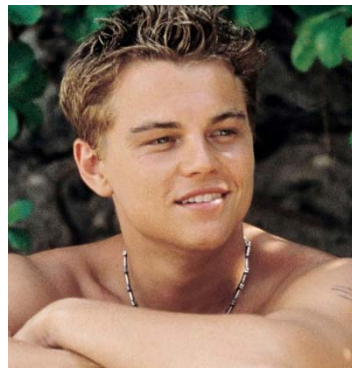
同一个人



不同人



同一个人





The first kind of problems



30岁



28岁



18岁



14岁



57岁

... ..



33岁



The second kind of problems

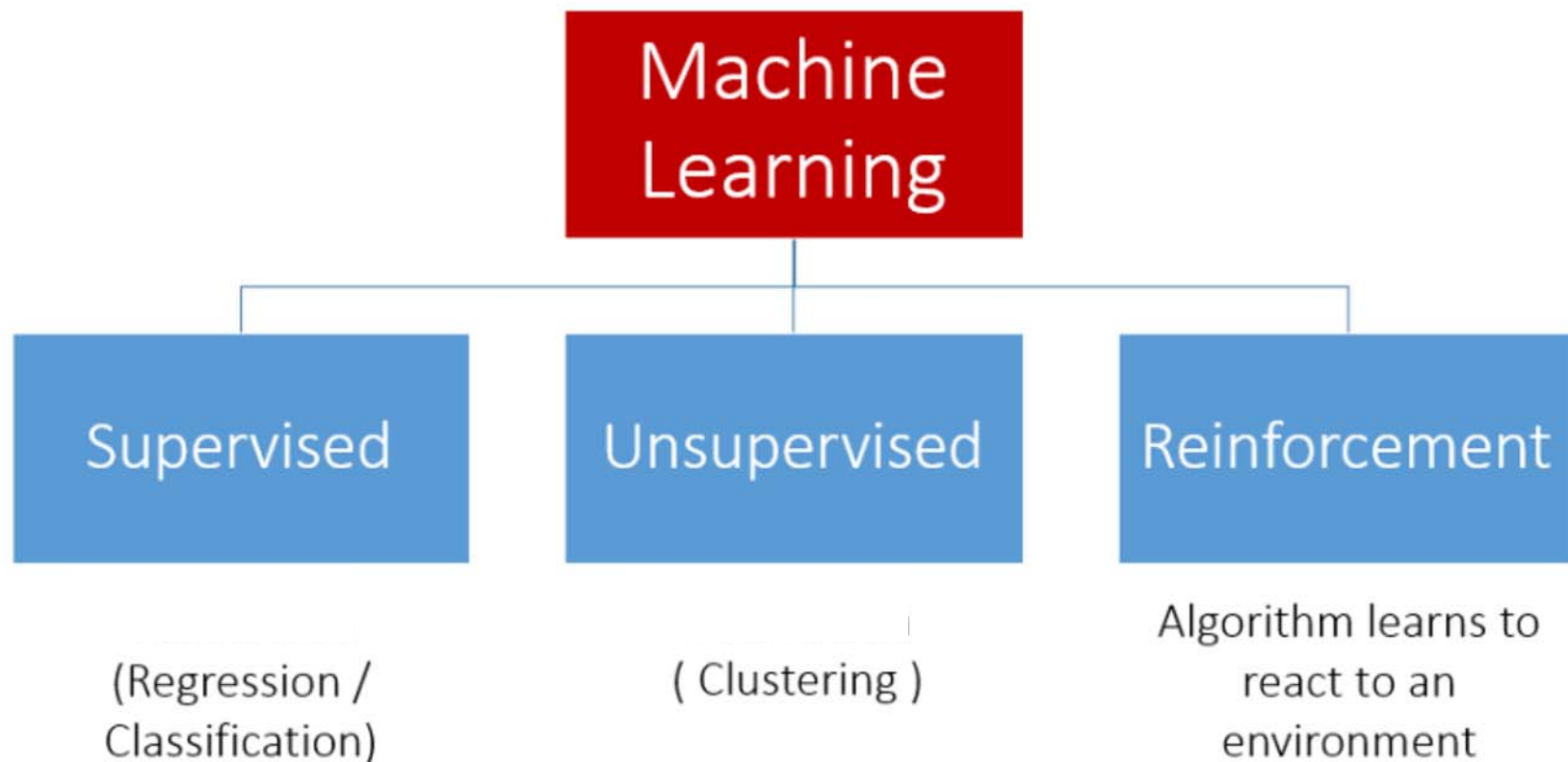




Two kinds of problems

- ▶ What are the differences?
- ▶ Supervised learning **vs.** Unsupervised learning

Types of Machine Learning





Two kinds of problems

- ▶ What are the differences?
- ▶ Supervised learning **vs.** Unsupervised learning
- ▶ Supervised learning
 - Goal: learn a mapping from inputs \mathbf{x} to outputs y
 - Training data: a labeled set of input-output pairs
 - Classification (Categorization, Decision making...)
 - y is a categorical variable
 - Regression
 - y is real-valued



Two kinds of problems

- ▶ What are the differences?
- ▶ Supervised learning **vs.** Unsupervised learning
- ▶ Unsupervised learning
 - We are only given inputs
 - Goal: find “interesting patterns”
 - Much less well-defined problem
 - Discovering clusters, Clustering
 - Discovering latent factors
 - Dimensionality reduction, Matrix factorization, Topic modeling



Two kinds of problems

- ▶ What are the differences?
- ▶ Supervised learning **vs.** Unsupervised learning
- ▶ Reinforcement learning
 - It is a supervised learning scenario
 - No desired category signal is given
 - The only teaching feedback is that the tentative category is right or wrong.
 - This is useful for learning how to act or behave when given occasional reward or punishment signals.



Focus of This Course

- ▶ What are the typical machine learning **problems**?
 - Supervised Learning
 - Classification (decision making)
 - Regression
 - Unsupervised Learning
 - Cluster analysis
 - Latent factor analysis
- ▶ What are the basic machine learning **tools (methods, algorithms)**?
- ▶ Matlab/Python programming



Basic Concepts of Supervised Learning

- ▶ Sample, example, pattern



- ▶ Features, predictors, independent variables

- $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n$

- ▶ State of the nature, labels, pattern class, class, responses, dependent variables

- $\omega_1, \omega_2, \dots, \omega_c$ or y_1, y_2, \dots, y_c or z_1, z_2, \dots, z_c

- ▶ Training data

- $(\mathbf{x}_1, \omega_1), (\mathbf{x}_2, \omega_2), \dots, (\mathbf{x}_n, \omega_n)$

- ▶ Model, statistical model, pattern class model, classifier

- f

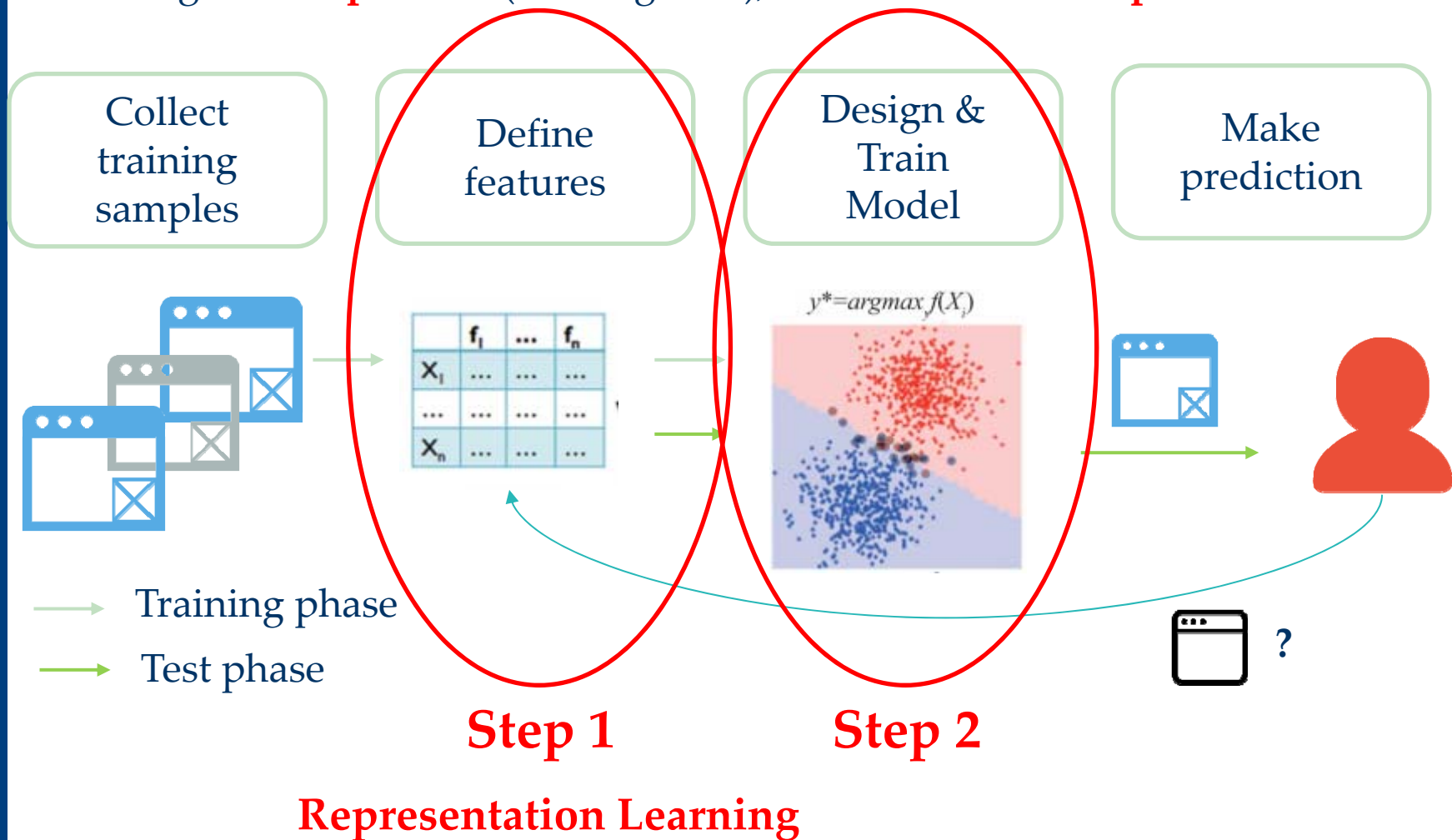
- ▶ Test data

- ▶ Training error & test error



Supervised Learning

Learning from **experience**(training data), and build **model** to **predict** the future





Supervised Learning

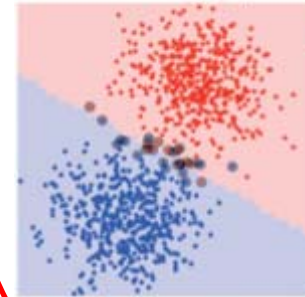
Define
features

	f_1	...	f_n
x_1
...
x_n

Step 1

Design &
Train
Model

$$y^* = \operatorname{argmax}_j f(X_j)$$



Step 2

- ▶ Which step is more important in building a successful system?
- ▶ Which one is the focus of this course?



Why general classification hard?

- Intra-class variability



The letter "T" in different typefaces

Define
features

	f_1	...	f_n
x_1
...
x_n

**Step 1 is not
good enough**



Same face under different expression, pose, illumination



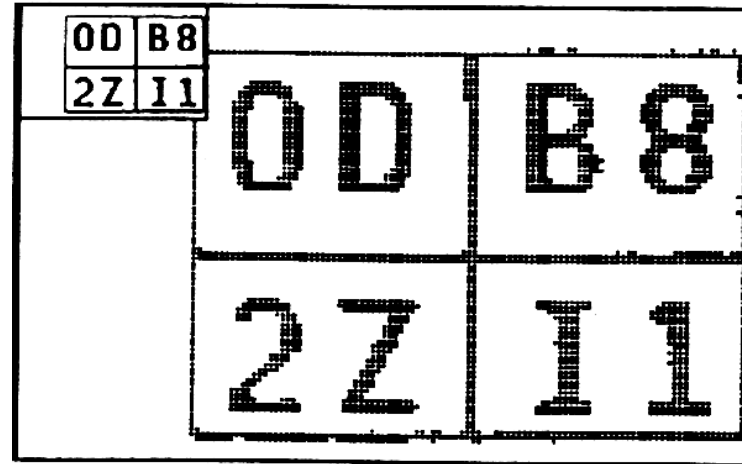
Why general classification hard?

- Inter-class similarity

Define
features

	f_1	...	f_n
x_1
...
x_n

**Step 1 is not
good enough**





Semantic Gap



Looks similar
But semantically
different



Looks different
But semantically
the same



Representation: Features

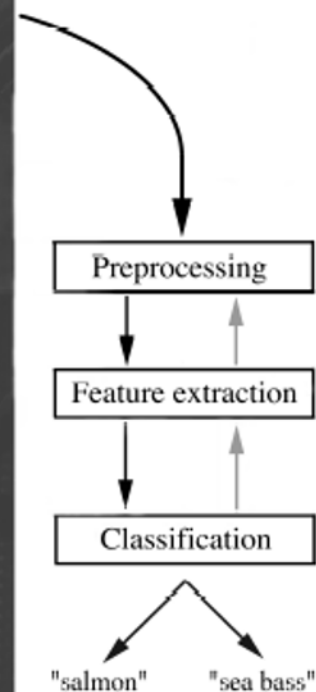
- ▶ Extract features to represent the samples
- ▶ Feature vector
- ▶ Good representation:
 - Low intra-class variability
 - Low inter-class similarity



Fish Classification: Salmon v. Sea Bass

Preprocessing involves
image enhancement
and segmentation;

- (i) separate touching
or occluding fishes
and
- (ii) extract fish
contour





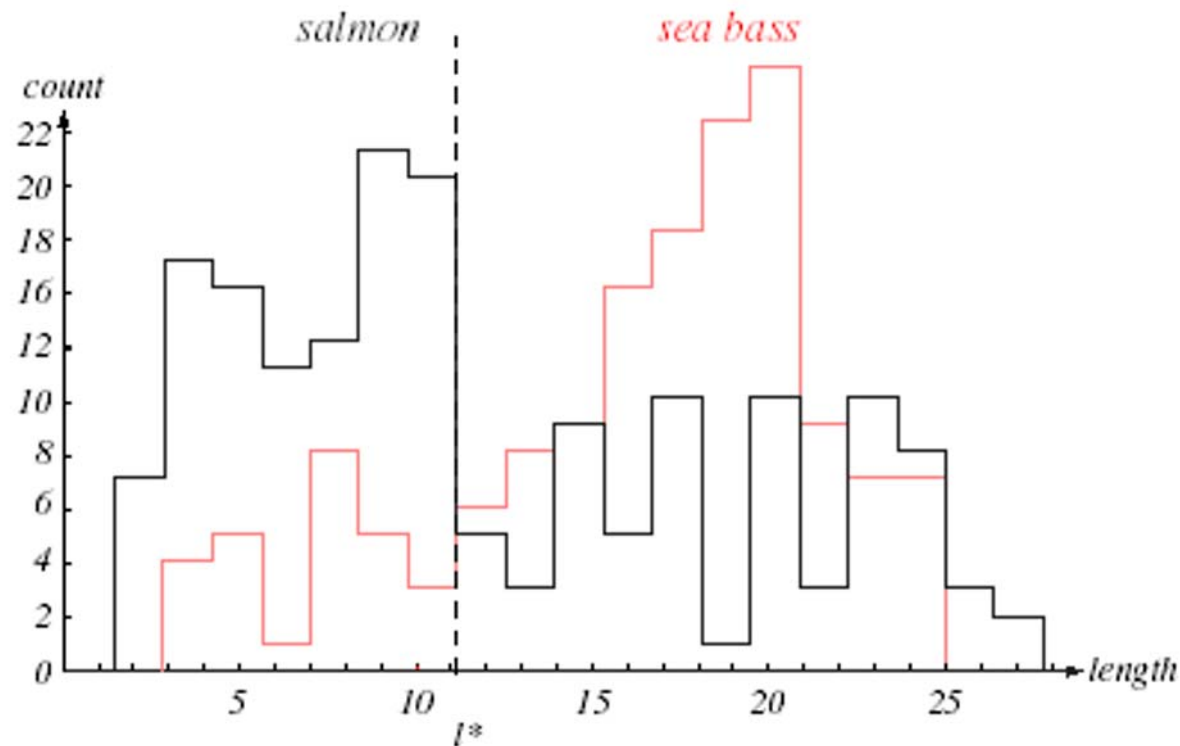
Representation: Fish Length As Feature

- ▶ How to design a classifier?



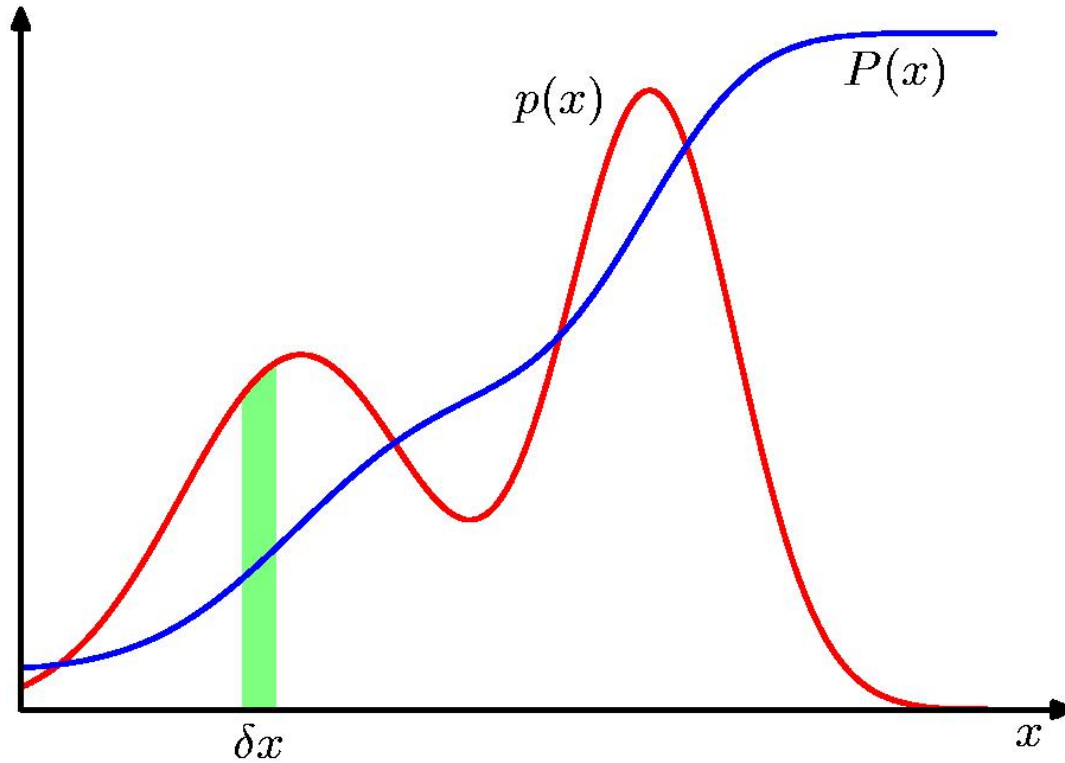
Representation: Fish Length As Feature

Training (design or learning) Samples





Probability Densities



$$p(x \in (a, b)) = \int_a^b p(x) dx$$

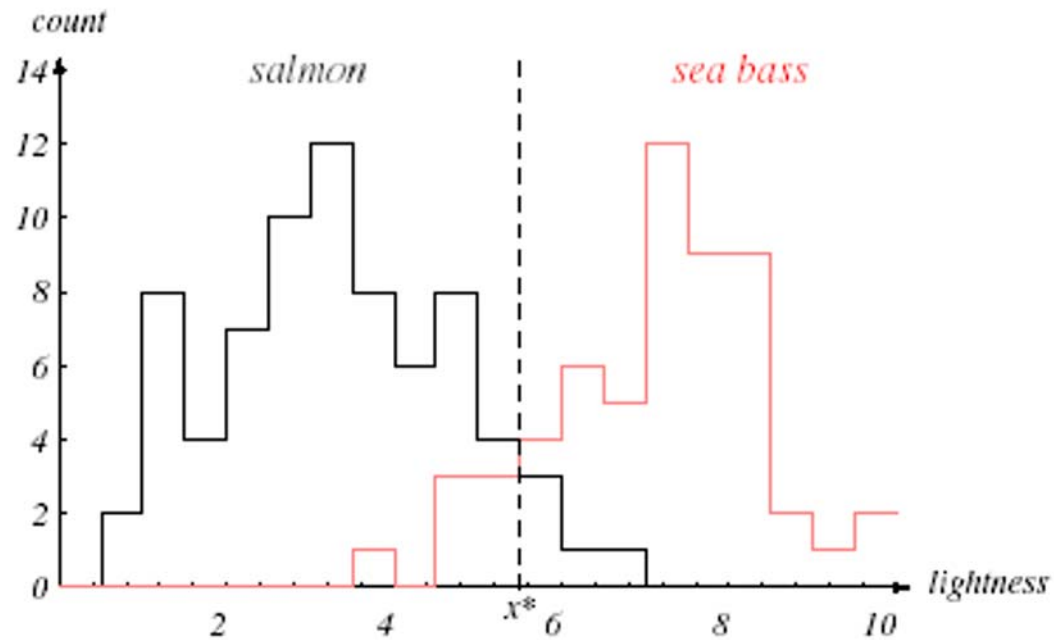
$$P(z) = \int_{-\infty}^z p(x) dx$$

$$p(x) \geq 0$$

$$\int_{-\infty}^{\infty} p(x) dx = 1$$



Fish Lightness As Feature

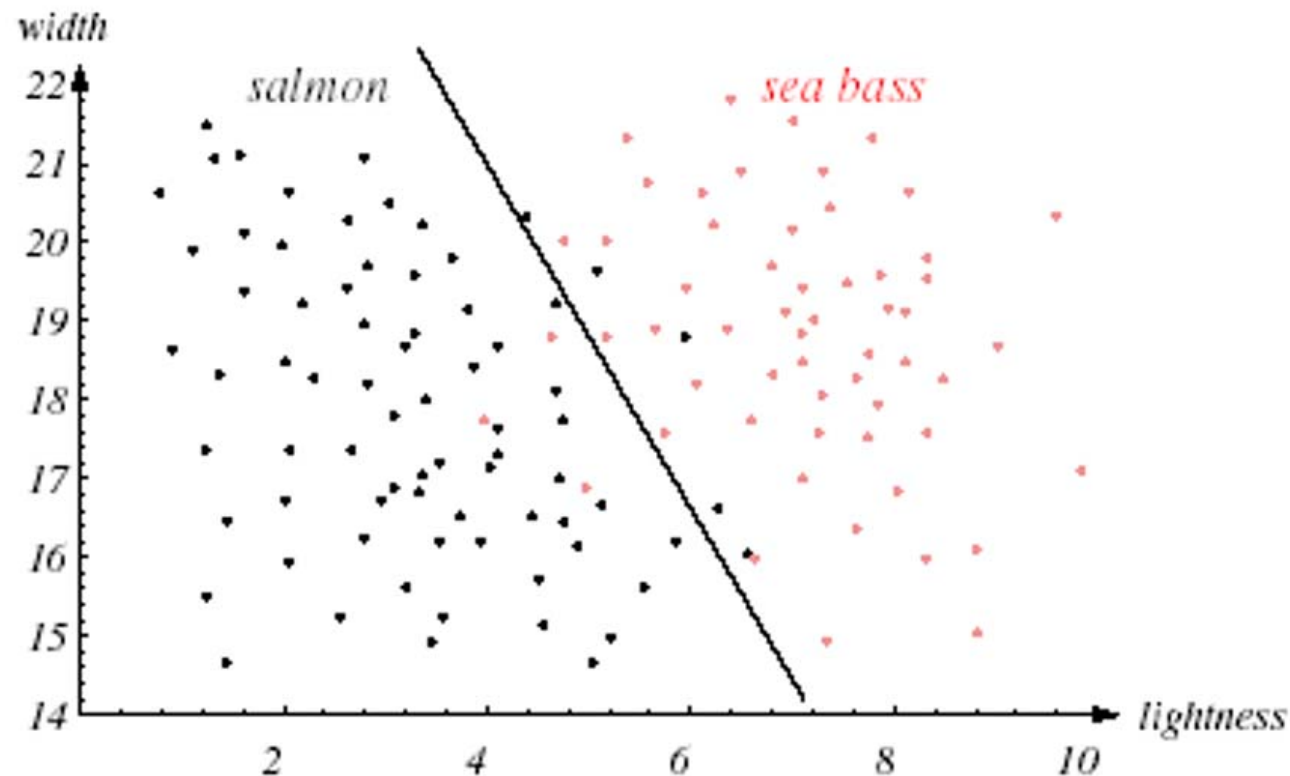


Overlap of these histograms is small compared to length feature



Two-dimensional Feature Space

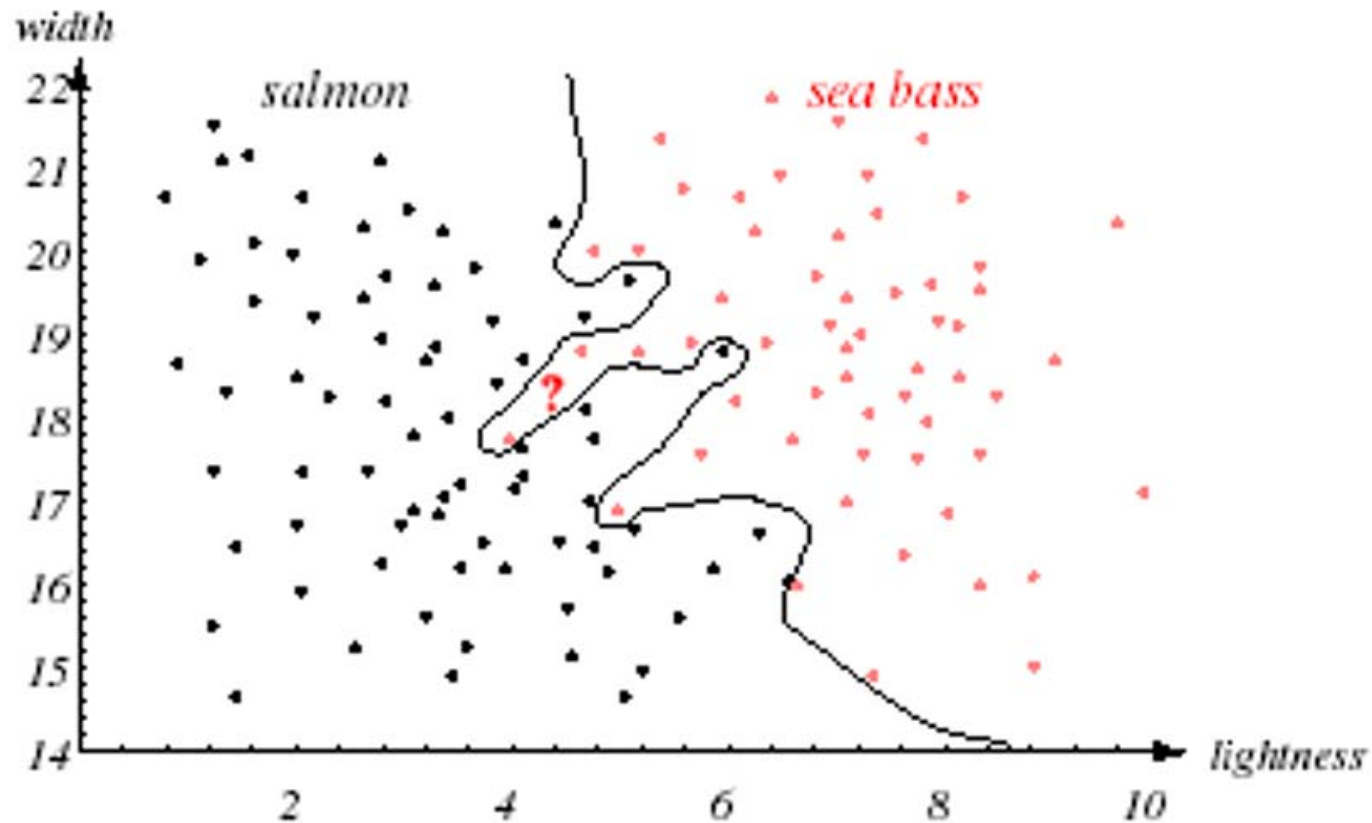
Linear (simple) decision boundary

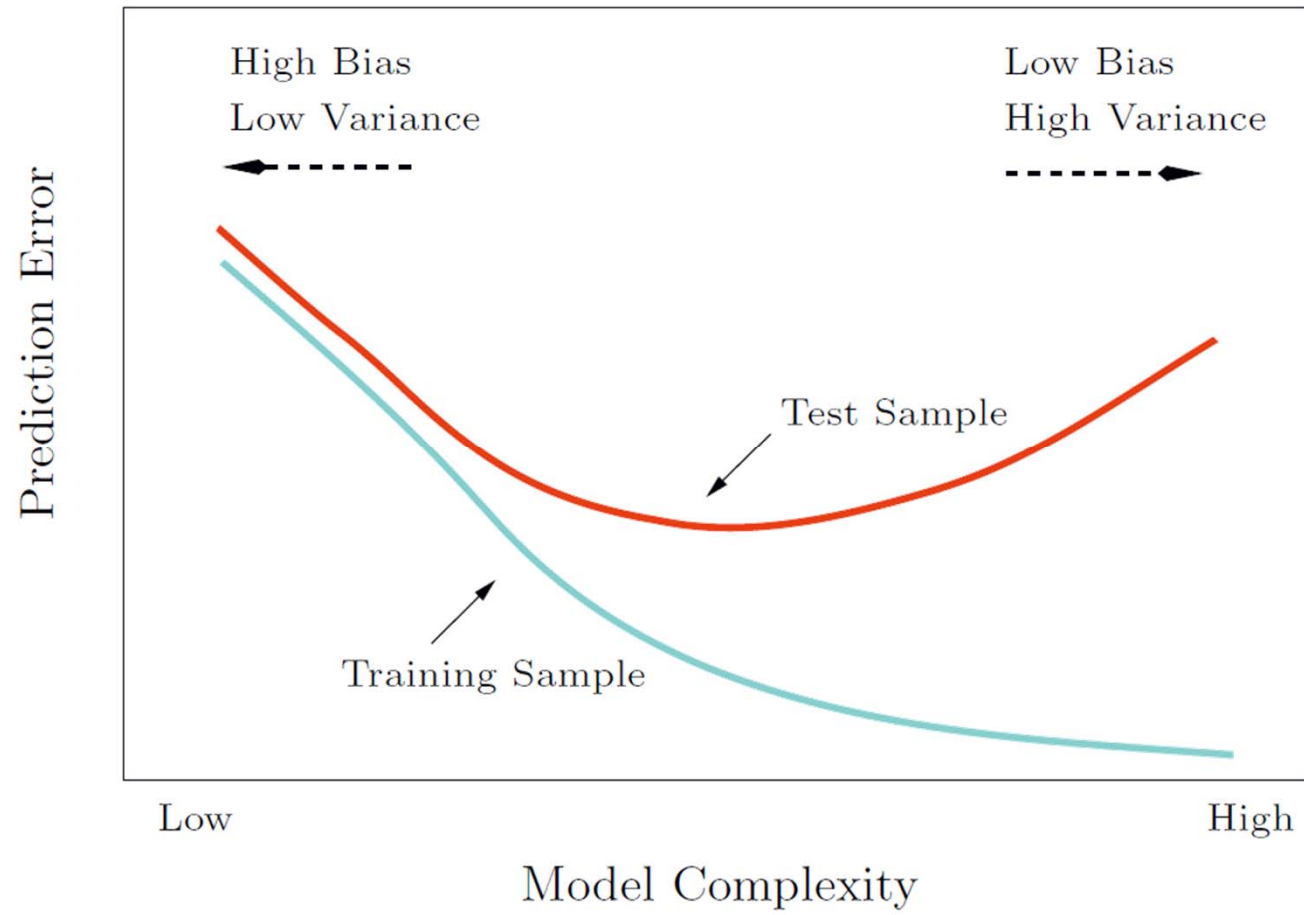


Two features together are better than individual features



Complex Decision Boundary







Generalization

- ▶ A generalization of a concept is an extension of the concept to less-specific criteria.
- ▶ Generalization of the classifier (model)
 - The performance of the classifier on **test** data.
- ▶ Training error:
 - ▶ Simple model \rightarrow large training error
 - ▶ Complex model \rightarrow less training error
- ▶ Test error:
 - ▶ Simple model \rightarrow ?
 - ▶ Complex model \rightarrow ?



Prerequisite Knowledge

- ▶ Probability:
 - Bayes theorem
- ▶ Analysis:
 - Gradient descent
- ▶ Linear Algebra
 - Linear space,
 - Matrix
 - Rank...
 - Positive definite matrix...
 - Eigenvector, eigenvalue
 - Singular vector, singular value