Machine Learning Basics Linear Regression

HKUST MSBD 6000B

Instructor: Yu Zhang

Machine learning basics

What is machine learning?

 "A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P, if its performance at tasks in T as measured by P, improves with experience E."

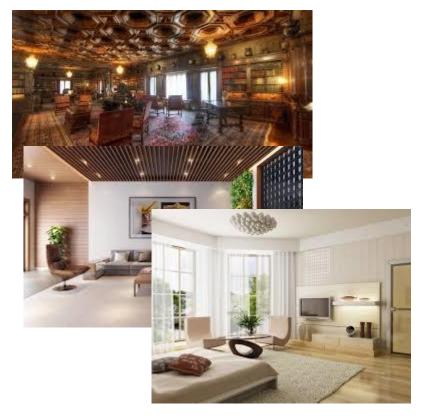
----- Machine Learning, Tom Mitchell, 1997

Example 1: image classification



Task: determine if the image is indoor or outdoor Performance measure: probability of misclassification

Example 1: image classification



Experience/Data: images with labels



Indoor outdoor

Example 1: image classification

- A few terminologies
 - Training data: the images given for learning
 - Test data: the images to be classified
 - Binary classification: classify into two classes

Example 1: image classification (multi-class)



Example 2: clustering images

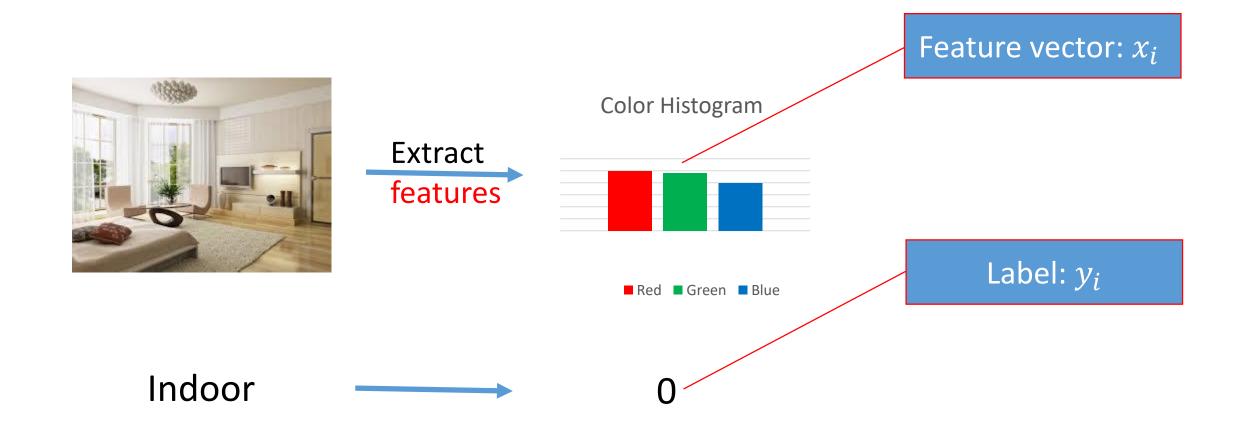


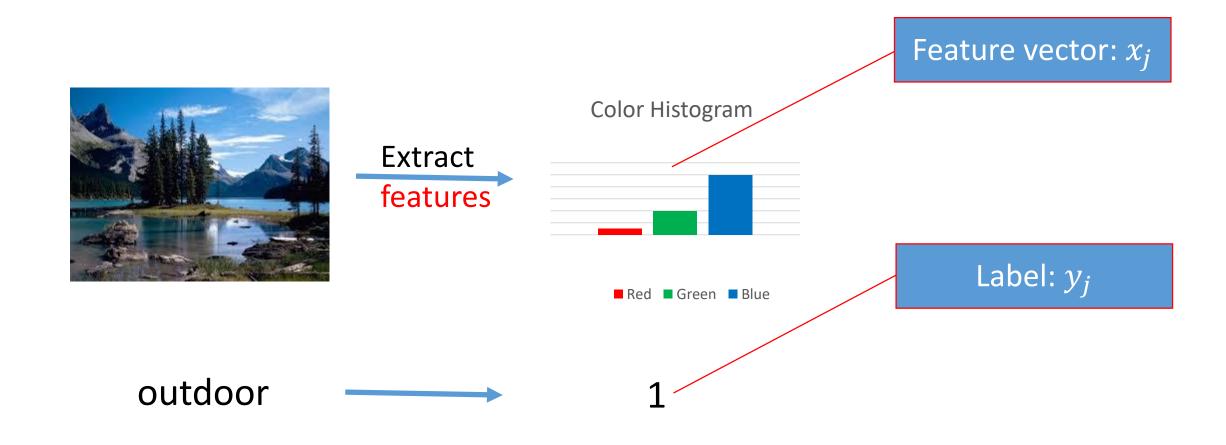
Task: partition the images into 2 groups Performance: similarities within groups

Data: a set of images

Example 2: clustering images

- A few terminologies
 - Unlabeled data vs labeled data
 - Supervised learning vs unsupervised learning





- Given training data $\{(x_i, y_i): 1 \le i \le n\}$
- Find y = f(x) using training data
- s.t. f correct on test data

What kind of functions?

- Given training data $\{(x_i, y_i): 1 \le i \le n\}$
- Find $y = f(x) \in \mathcal{H}$ using training data
- s.t. *f* correct on test data

Hypothesis class

- Given training data $\{(x_i,y_i): 1 \le i \le n\}$ Find $y=f(x) \in \mathcal{H}$ using training data
- s.t. f correct on test data

Connection between training data and test data?

- Given training data $\{(x_i, y_i): 1 \le i \le n\}$ i.i.d. from distribution D
- Find $y = f(x) \in \mathcal{H}$ using training data
- s.t. f correct on test data i.i.d. from distribution D

They have the same distribution

i.i.d.: independently identically distributed

- Given training data $\{(x_i, y_i): 1 \le i \le n\}$ i.i.d. from distribution D
- Find $y = f(x) \in \mathcal{H}$ using training data
- s.t. f correct on test data i.i.d. from distribution D

What kind of performance measure?

- Given training data $\{(x_i, y_i): 1 \le i \le n\}$ i.i.d. from distribution D
- Find $y = f(x) \in \mathcal{H}$ using training data
- s.t. the expected loss is small

$$L(f) = \mathbb{E}_{(x,y)\sim D}[l(f,x,y)] -$$

Various loss functions

- Given training data $\{(x_i, y_i): 1 \le i \le n\}$ i.i.d. from distribution D
- Find $y = f(x) \in \mathcal{H}$ using training data
- s.t. the expected loss is small

$$L(f) = \mathbb{E}_{(x,y)\sim D}[l(f,x,y)]$$

- Examples of loss functions:
 - 0-1 loss: $l(f, x, y) = \mathbb{I}[f(x) \neq y]$ and $L(f) = \Pr[f(x) \neq y]$
 - l_2 loss: $l(f, x, y) = [f(x) y]^2$ and $L(f) = \mathbb{E}[f(x) y]^2$

- Given training data $\{(x_i, y_i): 1 \le i \le n\}$ i.i.d. from distribution D
- Find $y = f(x) \in \mathcal{H}$ using training data
- s.t. the expected loss is small

$$L(f) = \mathbb{E}_{(x,y)\sim D}[l(f,x,y)]$$

How to use?

- Given training data $\{(x_i, y_i): 1 \le i \le n\}$ i.i.d. from distribution D
- Find $y = f(x) \in \mathcal{H}$ that minimizes $\hat{L}(f) = \frac{1}{n} \sum_{i=1}^{n} l(f, x_i, y_i)$
- s.t. the expected loss is small

$$L(f) = \mathbb{E}_{(x,y)\sim D}[l(f,x,y)]$$

Empirical loss

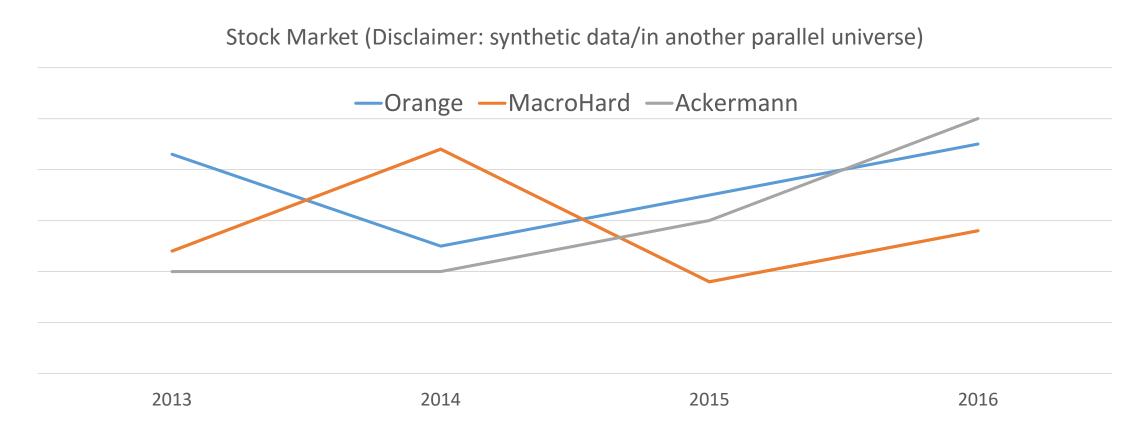
Machine learning

- Collect data and extract features
- Build model: choose hypothesis class $m{\mathcal{H}}$ and loss function l
- Optimization: minimize the empirical loss

Representation Learning

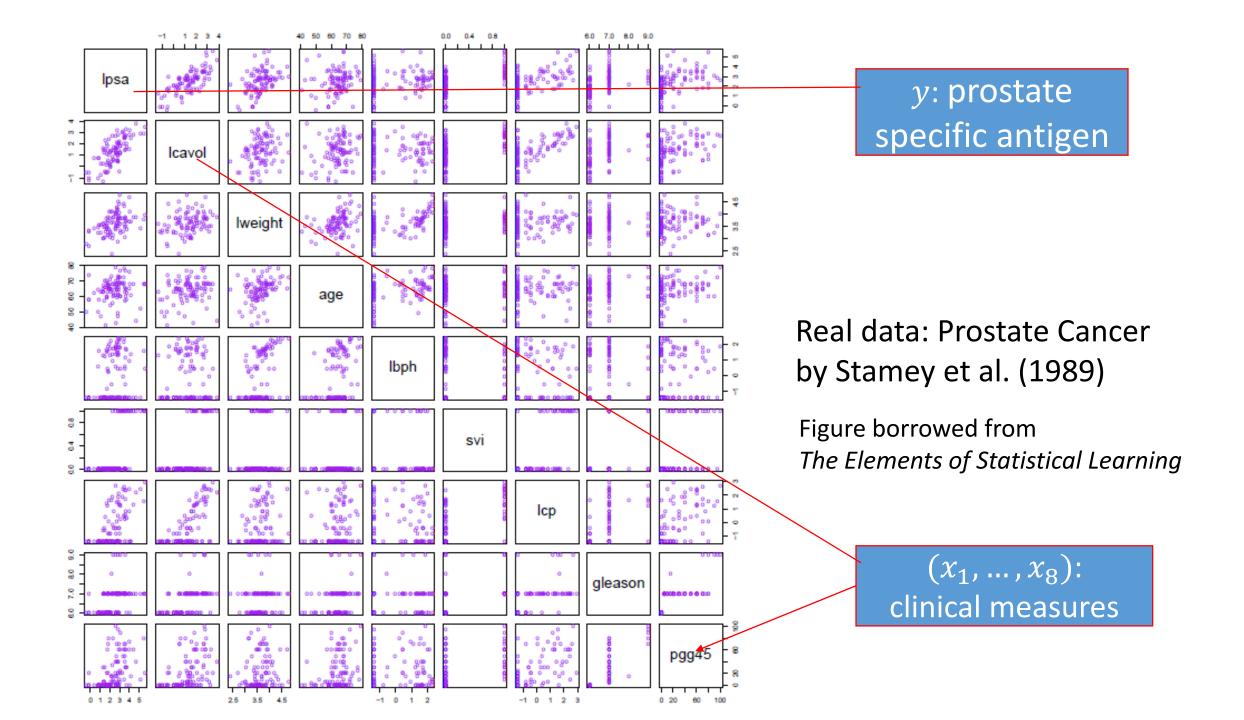
- Why handcraft the feature vectors x, y?
 - Can use prior knowledge to design suitable features
- Can computer learn the features on the raw images?
 - Learn features directly on the raw images: Representation Learning
 - Deep Learning ⊆ Representation Learning ⊆ Machine Learning ⊆ Artificial Intelligence

Example: Stock Market Prediction



Sliding window over time: serve as input x; non-i.i.d.

Linear regression



Linear regression

- Given training data $\{(x_i, y_i): 1 \le i \le n\}$ i.i.d. from distribution D
- Find $f_{w}(x) = w^{T}x$ that minimizes $\hat{L}(f_{w}) = \frac{1}{n}\sum_{i=1}^{n}(w^{T}x_{i} y_{i})^{2}$

 l_2 loss; also called mean square error

Hypothesis class ${m {\mathcal H}}$

Linear regression: optimization

- Given training data $\{(x_i, y_i): 1 \le i \le n\}$ i.i.d. from distribution D
- Find $f_w(x) = w^T x$ that minimizes $\hat{L}(f_w) = \frac{1}{n} \sum_{i=1}^n (w^T x_i y_i)^2$
- Let X be a matrix whose i-th row is x_i^T , y be the vector $(y_1, ..., y_n)^T$

$$\widehat{L}(f_w) = \frac{1}{n} \sum_{i=1}^{n} (w^T x_i - y_i)^2 = \frac{1}{n} ||Xw - y||_2^2$$

Linear regression: optimization

Set the gradient to 0 to get the minimizer

$$\nabla_{w} \hat{L}(f_{w}) = \nabla_{w} \frac{1}{n} ||Xw - y||_{2}^{2} = 0$$

$$\nabla_{w} [(Xw - y)^{T} (Xw - y)] = 0$$

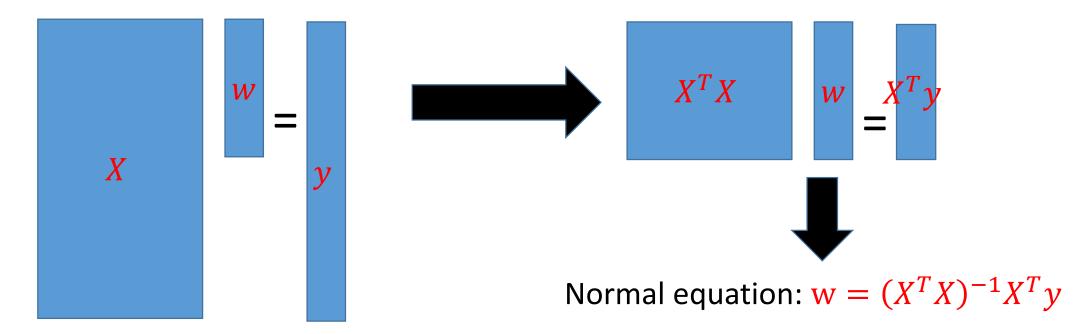
$$\nabla_{w} [w^{T} X^{T} X w - 2w^{T} X^{T} y + y^{T} y] = 0$$

$$2X^{T} X w - 2X^{T} y = 0$$

$$w = (X^{T} X)^{-1} X^{T} y$$

Linear regression: optimization

- Algebraic view of the minimizer
 - If X is invertible, just solve Xw = y and get $w = X^{-1}y$
 - But typically *X* is a tall matrix



Linear regression with bias

Bias term

- Given training data $\{(x_i, y_i): 1 \le i \le n\}$ i.i.d. from distribution D
- Find $f_{w,b}(x) = w^T x + b$ to minimize the loss
- Reduce to the case without bias:
 - Let w' = [w; b], x' = [x; 1]
 - Then $f_{w,b}(x) = w^T x + b = (w')^T (x')$