



SOI1010

Machine Learning II

Lecture 3: Linear Regression

Department of Data Science
Hanyang University

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



cat

- Task: determine if the image is cat or not
- Performance measure: probability of misclassification (error rate)

Example 1: Image classification



cats

Experience/Data:
Images with labels



dogs

Example 1: image classification

- A few terminologies
 - Training data: the images given for learning
 - Validation data: for model selection
 - Test data: the images to be classified (for evaluation)
 - Binary classification: classify into two classes
 - E.g., dogs vs cats
 - Multiclass classification: classify into three or more classes

Performance measure

- Specific to task T
- Evaluated using datasets
 - Training/validation/test sets
- Often challenging to choose
 1. Difficult to decide what to measure
 2. Know ideal measure but measurement could be impractical

Central challenge in ML

- Generalization
 - Ability to perform well on previously unseen examples
- Generalization error
 - Expected error on new examples \Rightarrow implausible to calculate
- Training error
 - Measured on a training set \Rightarrow bad proxy for generalization error
- Test error
 - Measured on a test set (not used in training)
 - \Rightarrow better proxy for generalization error

Example 2: clustering images



- Task: partition the images into 2 groups
- Performance: similarities within groups
- Data: a set of images

Example 2: clustering images

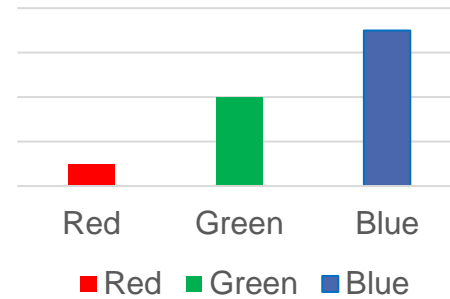
- A few terminologies
 - Unlabeled data vs labeled data
 - Unsupervised learning vs supervised learning
- Examples of unsupervised learning
 - Principal component analysis / dimensionality reduction (e.g., PCA)
 - Cluster analysis (e.g., k-means)
- Examples of supervised learning
 - SVM, Linear regression, Logistic regression, k-NN, neural networks
 - Naïve Bayes, Linear Discriminant Analysis, Decision trees

Math formulation



Extract
features

Color Histogram



Feature vector: $\mathbf{x}^{(i)}$

Outdoor

1

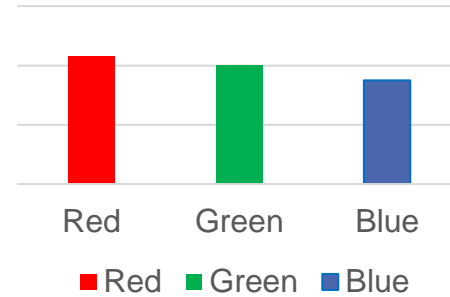
Label: $y^{(i)}$

Math formulation



Extract
features

Color Histogram



Feature vector: $\mathbf{x}^{(j)}$

Indoor

0

Label: $y^{(j)}$

Math formulation

- Given training data $\{(\mathbf{x}^{(i)}, y^{(i)}) : 1 \leq i \leq n\}$

But.. What kind of functions?

- Find $y = f(\mathbf{x})$ using training data
- s.t. f produces correct results on new/unseen test data

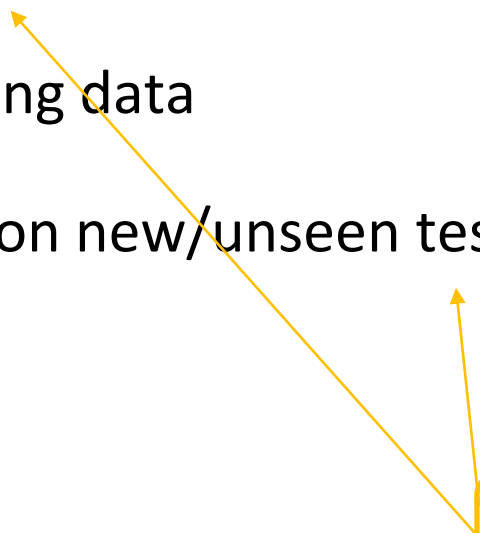
Math formulation

- Given training data $\{(\mathbf{x}^{(i)}, y^{(i)}) : 1 \leq i \leq n\}$
- Find $y = f(\mathbf{x}) \in \mathcal{H}$ using training data
- s.t. f produces correct results on new/unseen test data

Hypothesis class

Math formulation

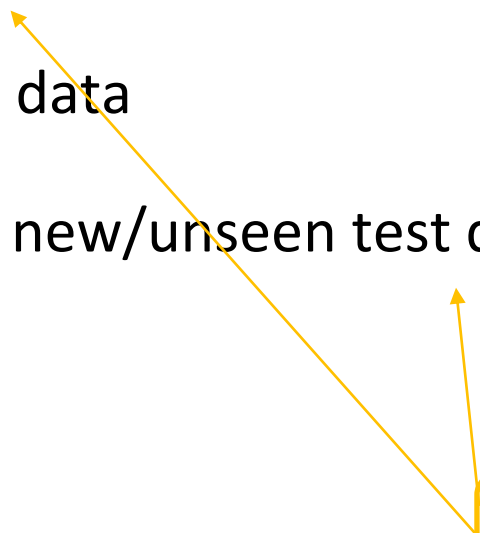
- Given training data $\{(\mathbf{x}^{(i)}, y^{(i)}) : 1 \leq i \leq n\}$
- Find $y = f(\mathbf{x}) \in \mathcal{H}$ using training data
- s.t. f produces correct results on new/unseen test data



Connection between training and test data?

Math formulation

- Given training data $\{(\mathbf{x}^{(i)}, y^{(i)}) : 1 \leq i \leq n\}$ i.i.d. from distribution \mathcal{D}
- Find $y = f(\mathbf{x}) \in \mathcal{H}$ using training data
- s.t. f produces correct results on new/unseen test data i.i.d. from distribution \mathcal{D}

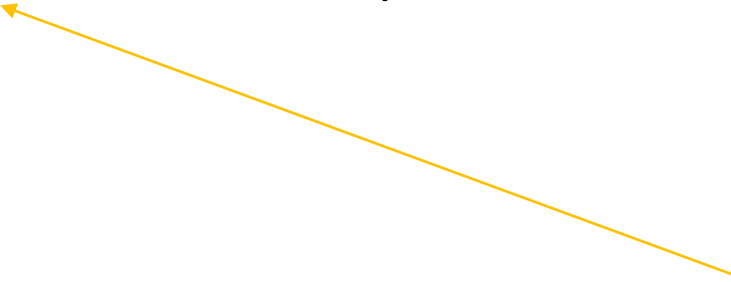


Connection between training and test data?

i.i.d.: independently identically distributed

Math formulation

- Given training data $\{(\mathbf{x}^{(i)}, y^{(i)}) : 1 \leq i \leq n\}$ i.i.d. from distribution \mathcal{D}
- Find $y = f(\mathbf{x}) \in \mathcal{H}$ using training data
- s.t. f produces correct results on new/unseen test data i.i.d. from distribution \mathcal{D}



What kind of performance measure?

Math formulation

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

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

Various loss functions



Math formulation

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

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

- Examples of loss functions:

➤ 0-1 loss

- $l(f, \mathbf{x}, y) = \mathbb{I}[f(\mathbf{x}) \neq y]$
- $\mathcal{L}(f) = P[f(\mathbf{x}) \neq y]$

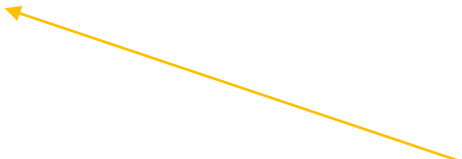
➤ l_2 loss

- $l(f, \mathbf{x}, y) = [f(\mathbf{x}) - y]^2$
- $\mathcal{L}(f) = \mathbb{E}[f(\mathbf{x}) - y]^2$

Math formulation

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

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



How to use?

Math formulation

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

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



Empirical loss

Machine Learning 1-2-3

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

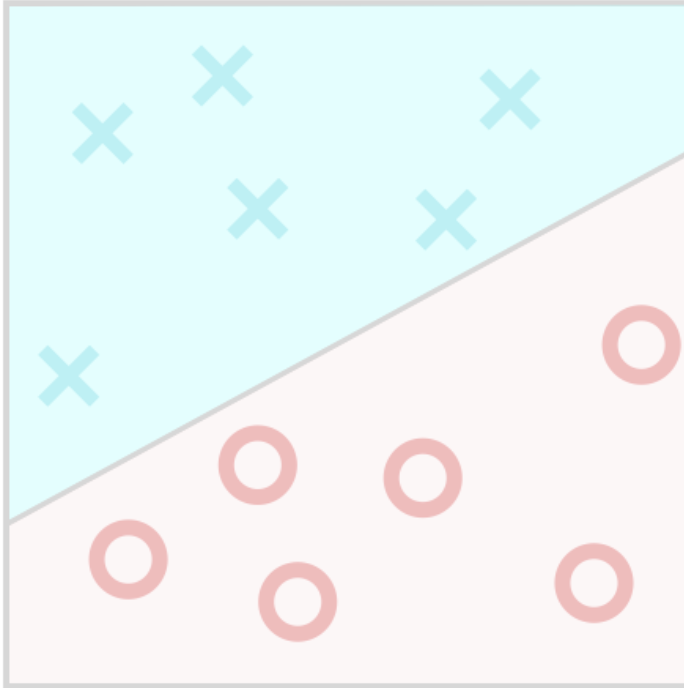
Wait...

- Does Machine Learning-1-2-3 include all approaches?
 - Include many but not all
 - For now, we will focus on Machine Learning-1-2-3

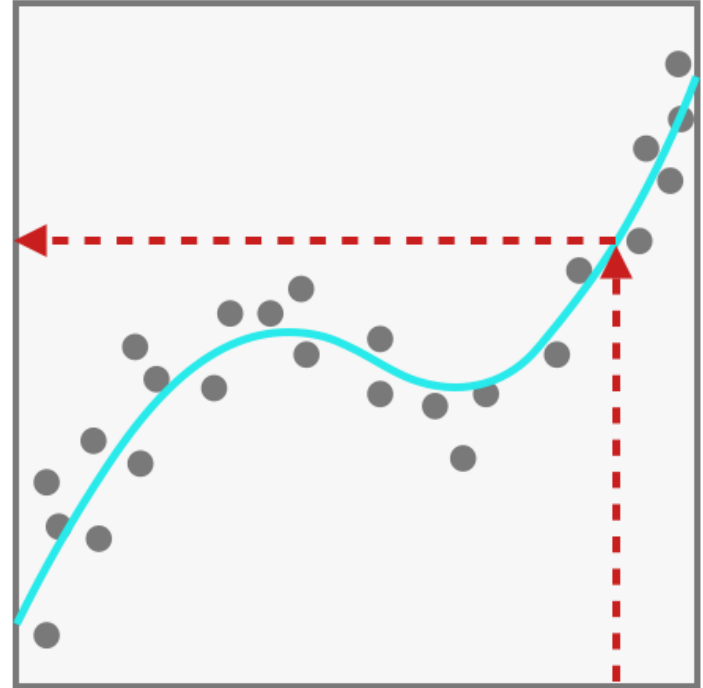
Wait...

- We can use prior knowledge to design suitable features
- But.. why handcraft the feature vectors?
- Can computers learn the features on the raw images?
 - Learn features directly on the raw images: Representation Learning
 - Deep Learning \subseteq Representation Learning \subseteq Machine Learning \subseteq Artificial Intelligence
- We could use prior knowledge to facilitate representation learning

Classification vs Regression

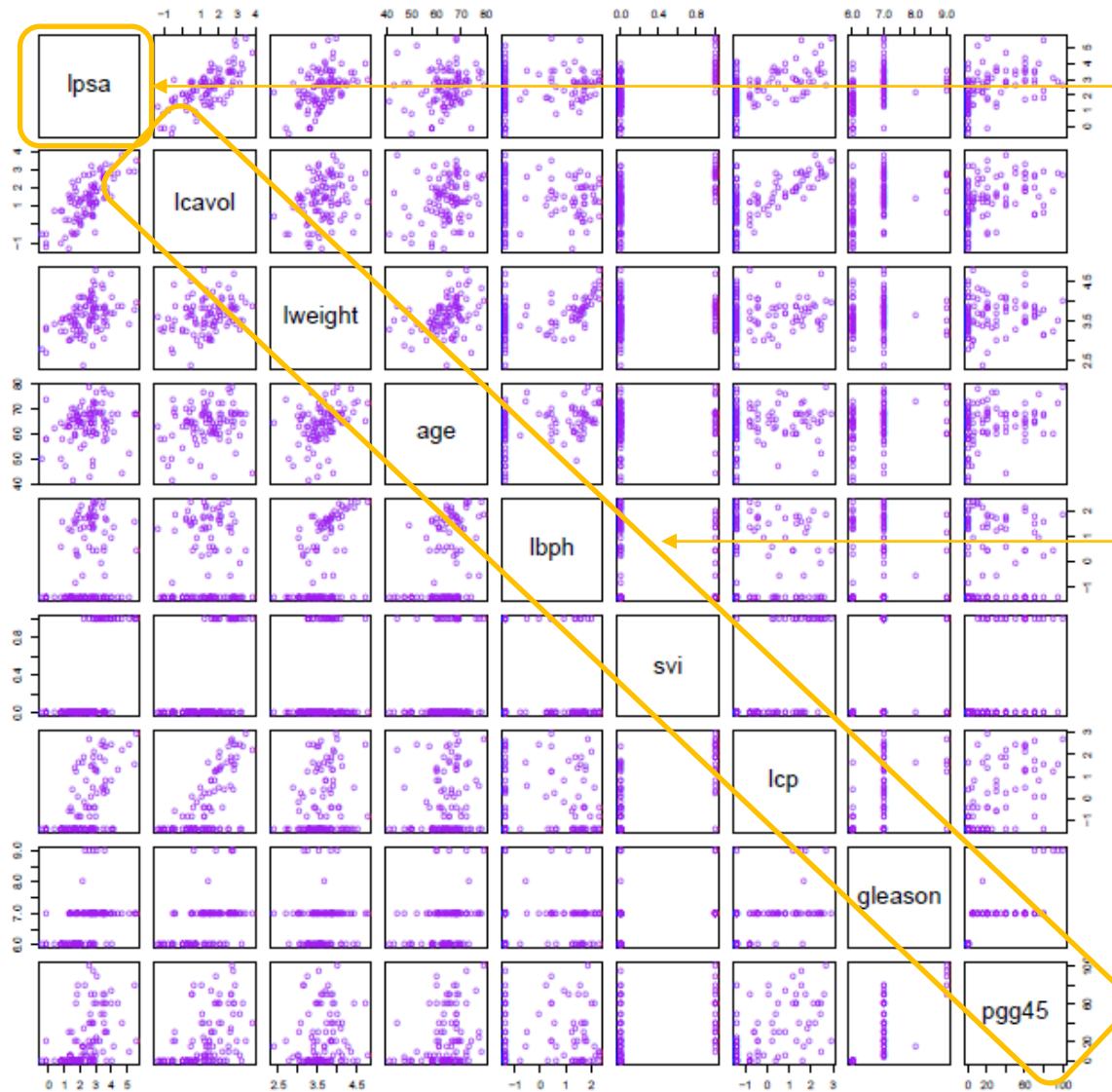


Classification



Regression

Linear regression



y : size/extent of cancer

$x = (x_1, \dots, x_8)$: clinical measures

Linear regression: Math formulation

- Given training data $\{(\mathbf{x}^{(i)}, y^{(i)}) : 1 \leq i \leq n\}$ i.i.d. from distribution \mathcal{D}
- Find $f_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ that minimizes $\hat{\mathcal{L}}(f_{\mathbf{w}}) = \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}^{(i)} - y^{(i)})^2$

Hypothesis class \mathcal{H}

l_2 loss; a.k.a.
mean squared error (MSE)

Linear regression: Math formulation

- Given training data $\{(\mathbf{x}^{(i)}, y^{(i)}) : 1 \leq i \leq n\}$ i.i.d. from distribution \mathcal{D}
- Find $f_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ that minimizes $\hat{\mathcal{L}}(f_{\mathbf{w}}) = \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}^{(i)} - y^{(i)})^2$
- Let \mathbf{X} be a matrix whose i -th row is \mathbf{x}_i^T .
- And let \mathbf{y} be the vector $(y_1, \dots, y_n)^T$

$$\hat{\mathcal{L}}(f_{\mathbf{w}}) = \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}^{(i)} - y^{(i)})^2 = \frac{1}{n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2$$

Linear regression: Math formulation & optimization

- Set the gradient to 0 to get the minimizer

$$\nabla_{\mathbf{w}} \hat{\mathcal{L}}(f_{\mathbf{w}}) = 0$$

$$\nabla_{\mathbf{w}} \frac{1}{n} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 = 0$$

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

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

$$\frac{\partial \mathbf{x}^T \mathbf{A}}{\partial \mathbf{x}} = \mathbf{A}$$

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

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

$$\frac{\partial \mathbf{x}^T \mathbf{A} \mathbf{x}}{\partial \mathbf{x}} = (\mathbf{A} + \mathbf{A}^T) \mathbf{x}$$

$\mathbf{X}^T \mathbf{X}$ is symmetric

Linear regression: Math formulation & optimization

- Algebraic view of the minimizer

- If \mathbf{X} is invertible, just solve $\mathbf{X}\mathbf{w} = \mathbf{y}$ and get $\mathbf{w} = \mathbf{X}^{-1}\mathbf{y}$

- But typically \mathbf{X} is a tall matrix

$$\begin{bmatrix} \mathbf{X} \end{bmatrix} \begin{bmatrix} \mathbf{w} \end{bmatrix} = \begin{bmatrix} \mathbf{y} \end{bmatrix} \quad \longrightarrow \quad \begin{bmatrix} \mathbf{X}^T \mathbf{X} \end{bmatrix} \begin{bmatrix} \mathbf{w} \end{bmatrix} = \begin{bmatrix} \mathbf{X}^T \mathbf{y} \end{bmatrix}$$
$$\downarrow$$
$$\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Linear regression

- Given training data $\{(\mathbf{x}^{(i)}, y^{(i)}) : 1 \leq i \leq n\}$ i.i.d. from distribution \mathcal{D}
- Find $f_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ that minimizes the loss $\hat{\mathcal{L}}(f_{\mathbf{w}}) = \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^T \mathbf{x}^{(i)} - y^{(i)})^2$

Linear regression with bias

- Given training data $\{(\mathbf{x}^{(i)}, y^{(i)}) : 1 \leq i \leq n\}$ i.i.d. from distribution \mathcal{D}
- Find $f_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$ that minimizes the loss



Bias term

- Reduce to the case without bias:
 - Let $\mathbf{w}' = [\mathbf{w}; b]$, $\mathbf{x}' = [\mathbf{x}, 1]$
 - Then, $f_{\mathbf{w}, b}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b = (\mathbf{w}')^T (\mathbf{x}')$
- Thus, whatever we did in previous slides still holds