

Course ID 80250993 Dates: 9/16-12/23/2021 @ I-205
Tencent/VooV Meeting 521 4678 6257 (passwd 1205)



Part I

Deterministic Machines for Supervised Learning

Xuegong Zhang
September 23, 2021



Xuegong Zhang

1

Course ID 80250993 Dates: 9/16-12/23/2021 @ I-205
Tencent/VooV Meeting 521 4678 6257 (passwd 1205)



Chapter 3

Linear Learning Machines

Xuegong Zhang
September 23, 2021



Xuegong Zhang

2



3.1 Fisher's Linear Discriminant

Xuegong Zhang

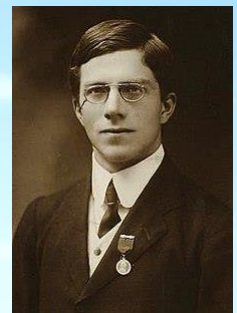
3

Sir Ronald Aylmer Fisher (R.A. Fisher)

(17 February 1890 – 29 July 1962)



- British statistician and geneticist
 - “a genius who almost single-handedly created the foundations for modern statistical science”
 - “the single most important figure in 20th century statistics”
 - “the greatest of Darwin's successors”.
- Some of the stuff he invented or popularized
 - ANOVA (analysis of variance)
 - Maximum likelihood
 - Fisher's z-distribution (F distribution)
 - Fisher's method for data fusion (meta-analysis)
 - The 0.05 cutoff of p-value, the notion of null hypothesis
 - Fisher's exact test
 - [Fisher's Discriminant Analysis \(in 1936\)](#)
 -
 - *The Genetical Theory of Natural Selection* (1930)
 - *The Design of Experiments* (1935)



From Wikipedia

"Natural selection is a mechanism for generating an exceedingly high degree of improbability."

Xuegong Zhang

4

The controversial side of R.A. Fisher



Judea Pearl, *The Book of Why*, 2018

FIGURE 4.3. R. A. Fisher with one of his many innovations: a Latin square experimental design, intended to ensure that one plot of each plant type appears in each row (fertilizer type) and column (soil type). Such designs are still used in practice, but Fisher would later argue convincingly that a randomized design is even more effective. (Source: Drawing by Dakota Harr.)



Homework (optional): Why proving smoking harmful to human health was so difficult?

Xuegong Zhang

5

Linear Classifier / Linear Discriminant

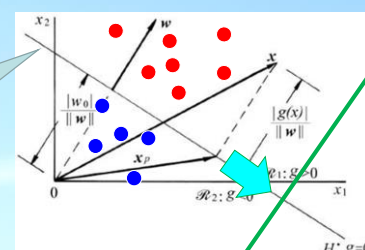


$$g(\mathbf{x}) = \sum_{i=1}^d w_i x_i + w_0$$

$$y = \text{sgn} \left(\sum_{i=1}^d w_i x_i + w_0 \right)$$

$$y = \begin{cases} +1 & \Rightarrow \text{class A or } \mathbf{x} \in \omega_1 \\ -1 & \Rightarrow \text{class B or } \mathbf{x} \in \omega_2 \end{cases}$$

Decision line
Classification boundary
Discriminant

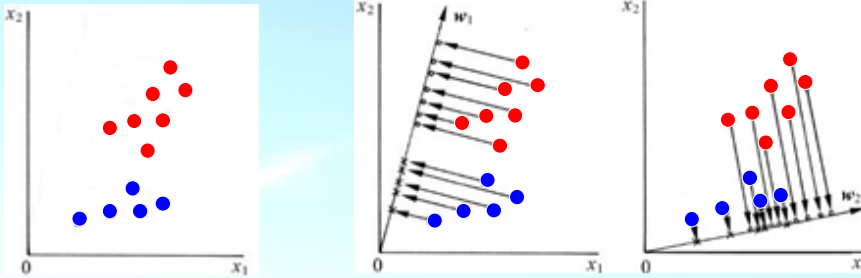


Projecting to one dimension and make decision based on the projection

Xuegong Zhang

6

Fisher's criterion for “the best” discriminant



- Which projection is better? Why and How?
- Fisher: The best projection: samples of the two classes “*best separated*”
 - Within-class: as close as possible
 - Between-classes: as far as possible

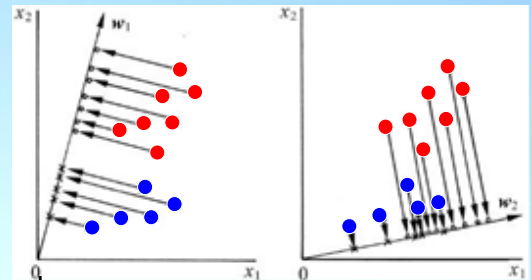
Xuegong Zhang

7

Fisher's criterion for “the best” discriminant



- Sample set: $\mathcal{X} = \{x_1, \dots, x_N\}$,
 where, Class 1 (ω_1): $\mathcal{X}_1 = \{x_1^1, \dots, x_{N_1}^1\}$,
 Class 2 (ω_2): $\mathcal{X}_2 = \{x_1^2, \dots, x_{N_2}^2\}$
- Projection:
 $\mathcal{X} \rightarrow \mathcal{Y}$: $y_i = w^T x_i$, $i = 1, \dots, N$
- Task: find the w that best separates the two classes



Xuegong Zhang

8



The “closeness” of samples before projection

- In \mathcal{X} space:

Class mean $\mathbf{m}_i = \frac{1}{N_i} \sum_{x_j \in \mathcal{X}_i} x_j, i = 1, 2$

Within-class scatter matrix

$$\mathbf{S}_i = \sum_{x_j \in \mathcal{X}_i} (x_j - \mathbf{m}_i)(x_j - \mathbf{m}_i)^T, \quad i = 1, 2$$

Total within-class scatter matrix $\mathbf{S}_w = \mathbf{S}_1 + \mathbf{S}_2$

Between-class scatter matrix

$$\mathbf{S}_b = (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^T$$

Xuegong Zhang

9



The “closeness” of samples after projection

- In \mathcal{Y} space:

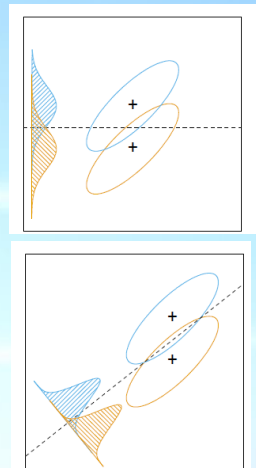
Class mean $\tilde{\mathbf{m}}_i = \frac{1}{N_i} \sum_{y_j \in \mathcal{Y}_i} y_j, i = 1, 2$

Within-class scatter

$$\tilde{\mathbf{S}}_i = \sum_{y_j \in \mathcal{Y}_i} (y_j - \tilde{\mathbf{m}}_i)(y_j - \tilde{\mathbf{m}}_i)^T, \quad i = 1, 2$$

Total within-class scatter $\tilde{\mathbf{S}}_w = \tilde{\mathbf{S}}_1 + \tilde{\mathbf{S}}_2$

Between-class scatter $\tilde{\mathbf{S}}_b = (\tilde{\mathbf{m}}_1 - \tilde{\mathbf{m}}_2)(\tilde{\mathbf{m}}_1 - \tilde{\mathbf{m}}_2)^T$



Xuegong Zhang

10



Fisher's Criterion

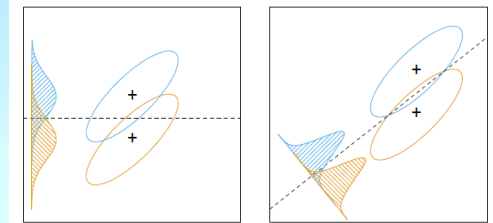
- Fisher's Criterion (of best separation):

$$\max_{\mathbf{w}} J_F(\mathbf{w}) = \frac{(\tilde{m}_1 - \tilde{m}_2)^2}{\tilde{S}_1 + \tilde{S}_2}$$

$$y_i = \mathbf{w}^T \mathbf{x}_i$$

- i.e., $\mathbf{w}^* = \operatorname{argmax}_{\mathbf{w}} J_F(\mathbf{w})$

$$J_F(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_b \mathbf{w}}{\mathbf{w}^T \mathbf{S}_w \mathbf{w}}$$



Xuegong Zhang

11



Solution for $\max_{\mathbf{w}} J_F(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_b \mathbf{w}}{\mathbf{w}^T \mathbf{S}_w \mathbf{w}}$?

- Not unique: If we change the norm of \mathbf{w} , $J_F(\mathbf{w})$ won't change.
- Fix the denominator $\mathbf{w}^T \mathbf{S}_w \mathbf{w} = c \neq 0$ and maximize the numerator $\mathbf{w}^T \mathbf{S}_b \mathbf{w}$, i.e.

$$\begin{aligned} \max \quad & \mathbf{w}^T \mathbf{S}_b \mathbf{w} \\ \text{s. t.} \quad & \mathbf{w}^T \mathbf{S}_w \mathbf{w} = c \end{aligned}$$

- Define the Lagrange function (*Lagrangian*)

$$L(\mathbf{w}, \lambda) = \mathbf{w}^T \mathbf{S}_b \mathbf{w} - \lambda(\mathbf{w}^T \mathbf{S}_w \mathbf{w} - c)$$



Xuegong Zhang

12



Let $\frac{\partial L}{\partial \mathbf{w}} = 0$, we get

$$\mathbf{S}_w^{-1} \mathbf{S}_b \mathbf{w}^* = \lambda \mathbf{w}^*$$

本征向量 (特征向量)

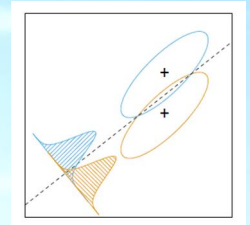
i.e., \mathbf{w}^* is the *eigenvector* of matrix $\mathbf{S}_w^{-1} \mathbf{S}_b$

Substitute $\mathbf{S}_b = (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^T$ in, we have

$$\begin{aligned} \lambda \mathbf{w}^* &= \mathbf{S}_w^{-1} (\mathbf{m}_1 - \mathbf{m}_2) (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w}^* \triangleq \mathbf{S}_w^{-1} (\mathbf{m}_1 - \mathbf{m}_2) R \\ R &\triangleq (\mathbf{m}_1 - \mathbf{m}_2)^T \mathbf{w}^* \end{aligned}$$

We only care about the direction, so

$$\mathbf{w}^* = \mathbf{S}_w^{-1} (\mathbf{m}_1 - \mathbf{m}_2)$$



Xuegong Zhang

13

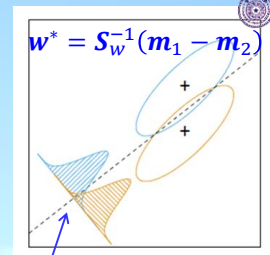
- Are we done for the solution?
- What was our original goal?
 - a “best” linear classifier, not a direction of projection



$$y = \text{sgn} \left(\sum_{i=1}^n w_i x_i + w_0 \right) = \text{sgn}(\mathbf{w}^T \mathbf{x} + w_0),$$

$$y = \begin{cases} +1 & \Rightarrow \mathbf{x} \in \omega_1 \\ -1 & \Rightarrow \mathbf{x} \in \omega_2 \end{cases}$$

- How shall we choose w_0 ?



Xuegong Zhang

14



- Commonly used thresholds:

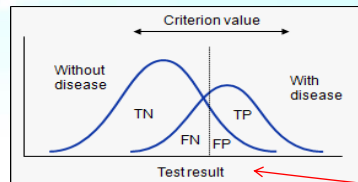
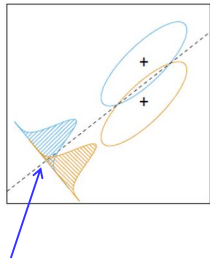
$$w_0 = -\frac{1}{2}(\tilde{m}_1 + \tilde{m}_2)$$

$$w_0 = -\tilde{m}$$

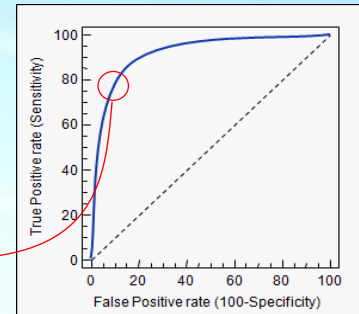
$$w_0 = -\frac{1}{2}(\tilde{m}_1 + \tilde{m}_2) + \frac{1}{N_1 + N_2 - 2} \ln \frac{P(\omega_1)}{P(\omega_2)}$$



- Choose the threshold with ROC curve



adjusting the
threshold

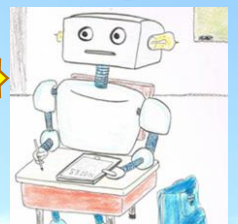


Xuegong Zhang

15

- “Wait, was that Machine Learning, not mathematics?”

$$w^* = S_w^{-1}(m_1 - m_2)$$



- Learning is the process of acquiring new, or modifying existing, knowledge, behaviors, **skills**, values, or preferences.
- Many popular learning algorithms are iterative, but ML does not have to be iterative.

Xuegong Zhang

16



1-minute stretch break



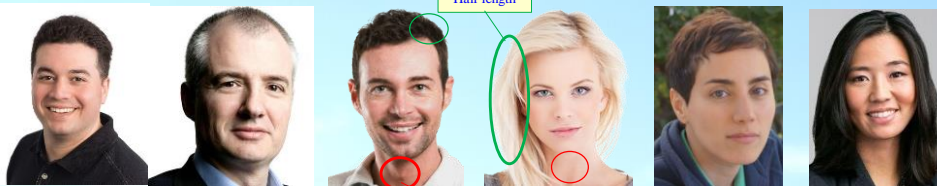
Xuegong Zhang

17



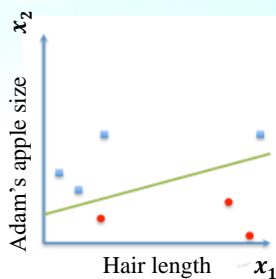
A toy example

- How to teach a machine to recognize man or woman?



The Adam's apple, or laryngeal prominence, is a feature of the human neck, and is the lump or protrusion that is formed by the angle of the thyroid cartilage surrounding the larynx seen especially in males.

An iterative
learning
procedure



For the green line:
 $w_1x_1 + w_2x_2 + w_0 = 0$
 For each blue dots:
 $w_1x_1 + w_2x_2 + w_0 > 0$
 For each red dots:
 $w_1x_1 + w_2x_2 + w_0 < 0$



Xuegong Zhang

18



3.2

Perceptron: “the first machine that learns to classify”

Xuegong Zhang

19

Perceptron (感知器)



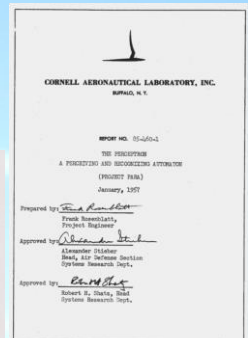
Frank Rosenblatt, *The Perceptron – a perceiving and recognizing automaton*,
Report 85-460-1, Cornell Aeronautical Laboratory, Jan. 1957

- The first learning machine: perceptron

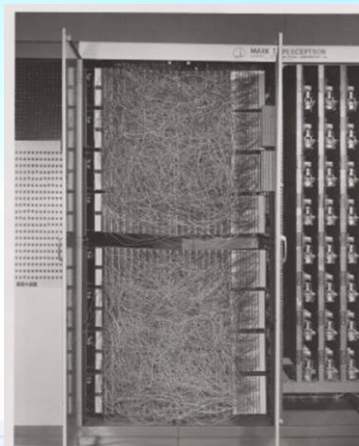
$$y = \text{sgn}\left(\sum_{i=1}^d w_i x_i + w_0\right)$$

- Why was it called a “learning machine”?

① It was a machine!



<https://en.wikipedia.org/wiki/Perceptron>



20



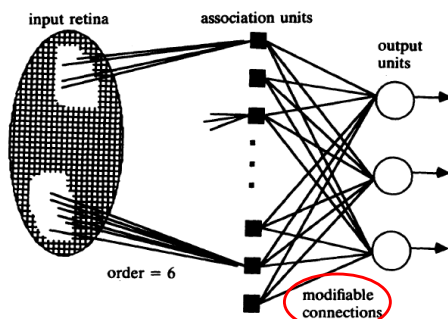
Perceptron

- Why is it a “learning machine”?

② **Because it can learn!**

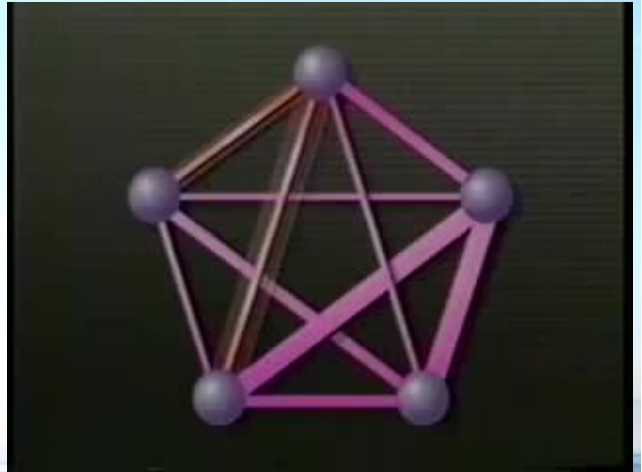
- It is not hard-coded in a Von-Neumann computing program, but rather, it is modifiable based on training data.

$$y = \text{sgn}\left(\sum_{i=1}^n w_i x_i + w_0\right)$$



M. Olazaran, A sociological study of the official history of the perceptrons controversy, *Social Studies of Science*, 1996

Xuegong Zhang



How Perceptron learns?



- To train the machine $y = \text{sgn}(\sum_{i=1}^d w_i x_i + w_0)$ with data $\{(x_1, y_1), \dots, (x_N, y_N)\}, x_i \in \mathbb{R}^d$

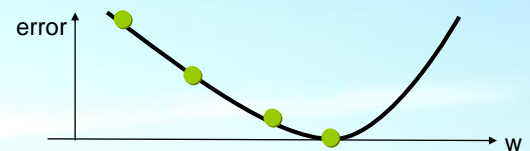
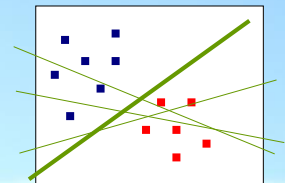
- Goal: to optimize an objective function $J(w)$

- The training error (i.e., # of mistakes)

- Learning algorithm

- do $w(t+1) = w(t) - \eta \nabla J$
- until $(\nabla J < \text{threshold})$

—— Basic Gradient Descent



Xuegong Zhang

22

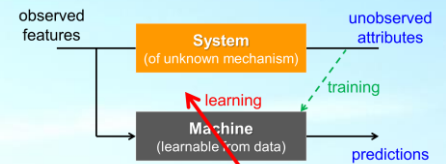
Basic concepts of (supervised) Machine Learning



- How can we make a learning machine?



- It needs a teacher.
 - We (people) *design* it (features, the model) and *train* it.
- We need materials to train it. / It needs materials to learn from.
 - Training data
- We need to tell what is the goal of the learning.
 - Objective function/criterion
- We need to tell it how to learn.
 - Learning/training algorithm

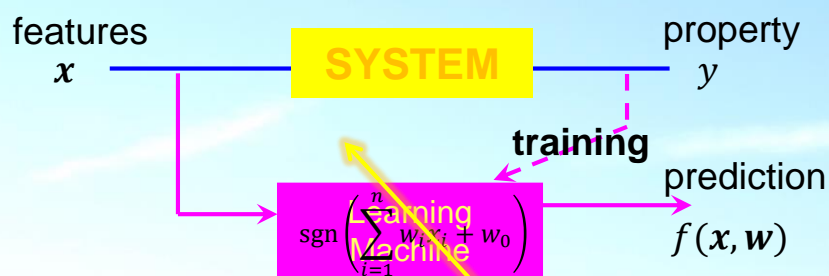


---- The basic elements of *Supervised Learning*

Xuegong Zhang

23

- Linear Machines



Xuegong Zhang

24



Linear discriminant function

$$g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = \boldsymbol{\alpha}^T \mathbf{y}$$

$\mathbf{y} = \begin{bmatrix} 1 \\ \mathbf{x} \end{bmatrix}$ augmented feature vector

$\boldsymbol{\alpha} = \begin{bmatrix} w_0 \\ \mathbf{w} \end{bmatrix}$ augmented weight vector

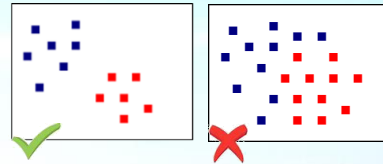
To simplify the treatment, we normalize the function by the following replacement

$$\mathbf{y}'_i = \begin{cases} \mathbf{y}_i, & \text{if } \mathbf{y}_i \in \omega_1, i = 1, \dots, N \\ -\mathbf{y}_i, & \text{if } \mathbf{y}_i \in \omega_2, i = 1, \dots, N \end{cases}$$

then we have $\boldsymbol{\alpha}^T \mathbf{y}'_i > 0, i = 1, \dots, N$ if all samples are correctly classified. We still use \mathbf{y}_i to represent a “*normalized augmented feature vector*” for convenience.

- Linearly separable:

$$\exists \boldsymbol{\alpha}, \quad \boldsymbol{\alpha}^T \mathbf{y}_i > 0, i = 1, \dots, N$$



Xuegong Zhang

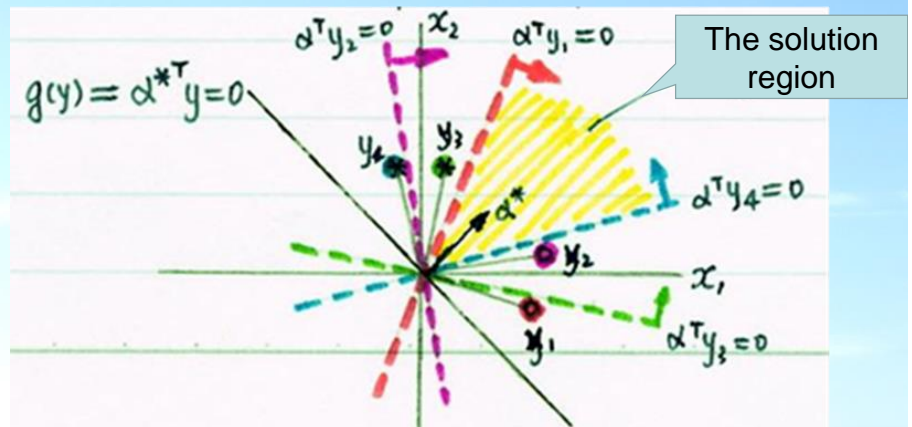
25

Solution vectors and solution region



- A solution $\boldsymbol{\alpha}^*$: a weight vector that satisfies

$$\boldsymbol{\alpha}^{*T} \mathbf{y}_i > 0, \\ i = 1, \dots, N$$
- Solution region: the region of all solution vectors in the weight space

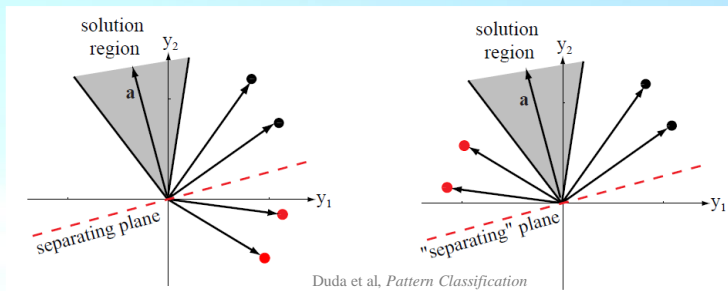


Xuegong Zhang

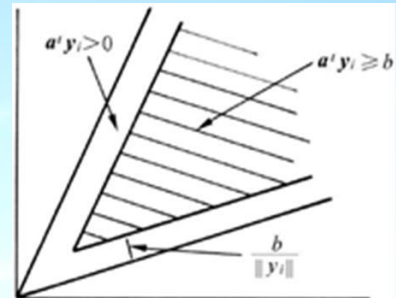
26



Solution vectors and solution region



- Introducing margins
 $\alpha^T \mathbf{y}_i \geq b > 0$



Xuegong Zhang

张学工《模式识别》

27

The Perceptron Criterion

$$J_P(\alpha) = \sum_{\mathbf{y}_j \in Y^k} (-\alpha^T \mathbf{y}_j)$$

- The Perceptron solution for linearly separable cases

$$J_P(\alpha^*) = \min J_P(\alpha) = 0$$

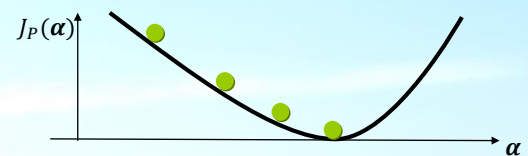
- Gradient descent:**

$$\alpha(k+1) = \alpha(k) - \rho_k \nabla J$$

$$\nabla J = \frac{\partial J_P(\alpha)}{\partial \alpha} = \sum_{\mathbf{y}_j \in Y^k} (-\mathbf{y}_j)$$

$$\therefore \alpha(k+1) = \alpha(k) + \rho_k \sum_{\mathbf{y}_j \in Y^k} \mathbf{y}_j$$

Y^k : set of samples misclassified by α at step k



Xuegong Zhang

28

The Perceptron Algorithm



- Fixed increment rule

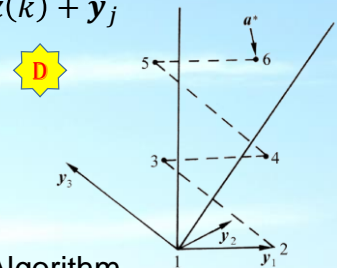
$$\alpha(k+1) = \alpha(k) + \rho_k \sum_{y_j \in Y^k} y_j$$

- ① Initialization (set initial weights to anything, say, zeros)
- ② For sample y_j , if $\alpha(k)^T y_j \leq 0$ (or b), then $\alpha(k+1) = \alpha(k) + y_j$
- ③ Repeat (2) for all samples till $J_p = 0$

- Variable increment rule $\rho_k = \frac{|\alpha(k)^T y_j|}{\|y_j\|^2}$

- Perceptron Convergence Theorem:

If training samples are linearly separable, the Perceptron Algorithm will **converge** to a solution vector in a **finite number of updates**.



Xuegong Zhang

29

Discussion



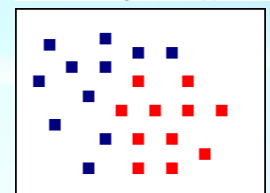
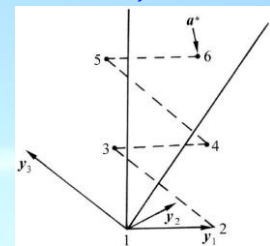
- Questions:

- What if the data are not linearly separable?
- How can we know the data are linearly separable?
- Can the machine still find a “useful” solution?
- Explore the answer with your homework.

- Solutions:

- Force the learning to stop, e.g., by shrinking the learning rate (step size)
- Use other linear methods that allow for errors
- Turn to nonlinear methods

$$\alpha(k+1) = \alpha(k) + \rho_k \sum_{y_j \in Y^k} y_j$$



Xuegong Zhang

30

单选题 1分

⚙ 设置

4-min break

- ☒ A Yes, I'm back to my computer.
- ☐ B Sorry, not yet.



3.3 Linear Regression



Machines to predict a real value?



e.g., Can it predict the score a student will get based on some observations?

• How shall we make such a machine?

- What features may predict the score?
 - Study hours, background, IQ, EQ/looking/shape/size, ...?
- What do we need for making such a machine?
 - A function set (the machine)
 - Training data
 - Objective function
 - Algorithm

$$y = f(x)$$

↑ prediction ↑ Real function ↑ sample

• How can we make a learning machine?

- It needs a teacher.
 - We (people) design it (features, the model) and train it.
- We need materials to train it. / It needs materials to learn from.
 - Training data
- We need to tell what is the goal of the learning.
 - Objective function/criterion
- We need to tell it how to learn.
 - Learning/training algorithm

---- The basic elements of Supervised Learning

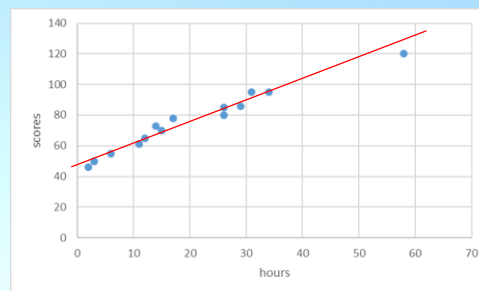
Xuegong Zhang

33

Toy data example



Student id	Final score	Study Hours per Week
1	50	3
2	95	34
3	78	17
4	55	6
5	65	12
6	70	15
7	80	26
8	86	29
9	73	14
10	120	58
11	46	2
12	95	31
13	85	26
14	61	11



Simple Linear Regression

$$y = w_0 + w_1x$$

Xuegong Zhang

34

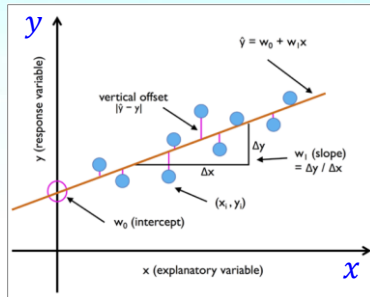


Linear Regression

y
dependent variable
response variable
regressand
output variable
target variable
predication
...

- Simple Linear Regression

$$y = w_0 + w_1 x$$



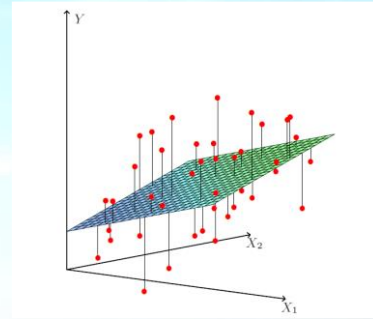
Raschka & Mirjalili, *Python Machine Learning*

x
independent variables
explanatory variables
predictor variables
regressors
input variable
features
...

- Multiple Linear Regression

$$y = w_0 + w_1 x + \dots + w_d x_d$$

$$= \sum_{i=0}^d w_i x_i = \mathbf{w}^T \mathbf{x}$$



Hastie, Tibshirani, Friedman, *Elements of Statistical Learning*

Xuegong Zhang

35

Linear Regression as a Learning Machine



- How can it *learn*?

– It needs a teacher.

→ We *design* it: $f(\mathbf{x}) = \sum_{i=0}^d w_i x_i = \mathbf{w}^T \mathbf{x}$

– We need materials to train it. / It needs materials to learn from.

→ Training data: $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$, $\mathbf{x}_j \in \mathbb{R}^{d+1}, y_j \in \mathbb{R}$

– We need to tell what is the goal of the learning.

→ Objective function: ?

– We need to tell it how to learn.

→ Learning algorithm: ?

Xuegong Zhang

36

Linear Regression as a Learning Machine



- How can it *learn*?

- It needs a teacher.

→ We design it: $f(\mathbf{x}) = \sum_{i=0}^d w_i x_i = \mathbf{w}^T \mathbf{x}$

- We need materials to train it. / It needs materials to learn from.

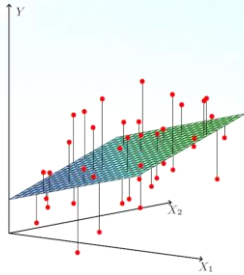
→ Training data: $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$, $\mathbf{x}_j \in \mathbb{R}^{d+1}, y_j \in \mathbb{R}$

- We need to tell what is the goal of the learning.

→ Objective function: $\min E = \frac{1}{N} \sum_{j=1}^N (f(\mathbf{x}_j) - y_j)^2$
 E : mean squared error, sum of squares, in-sample error, ...

- We need to tell it how to learn.

→ Learning algorithm: ?



Xuegong Zhang

37

Linear Regression Algorithm



$$\min_{\mathbf{w}} E(\mathbf{w}) = \frac{1}{N} \sum_{j=1}^N (f(\mathbf{x}_j) - y_j)^2 = \frac{1}{N} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|^2 = \frac{1}{N} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y})$$

where $\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_N^T \end{bmatrix}, \mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$.

Solution:

Let $\nabla E(\mathbf{w}) = \frac{\partial E(\mathbf{w})}{\partial \mathbf{w}} = \frac{2}{N} \mathbf{X}^T (\mathbf{X}\mathbf{w} - \mathbf{y}) = 0$,

we have $\mathbf{X}^T \mathbf{X} \mathbf{w} = \mathbf{X}^T \mathbf{y}$.

Therefore, $\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$, if $(\mathbf{X}^T \mathbf{X})$ is invertible.

Pseudo-inverse: $\mathbf{X}^+ = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$

$$\text{dim:} \quad \left[\begin{array}{c} \text{ } \\ \text{ } \end{array} \right] \left[\begin{array}{c} \text{ } \\ \text{ } \end{array} \right]^{-1} \left[\begin{array}{c} \text{ } \\ \text{ } \end{array} \right]$$

$(d+1) \times N \quad N \times (d+1) \quad (d+1) \times N$



Xuegong Zhang

38



Linear regression algorithm:

- 1: Construct the matrix X and the vector y from the data set $(x_1, y_1), \dots, (x_N, y_N)$, where each x includes the $x_0 = 1$ bias coordinate, as follows

$$X = \underbrace{\begin{bmatrix} -x_1^T \\ -x_2^T \\ \vdots \\ -x_N^T \end{bmatrix}}_{\text{input data matrix}}, \quad y = \underbrace{\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}}_{\text{target vector}}.$$

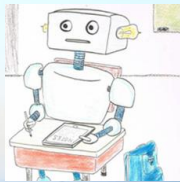
- 2: Compute the pseudo-inverse X^\dagger of the matrix X . If $X^T X$ is invertible,

$$X^\dagger = (X^T X)^{-1} X^T.$$

- 3: Return $w_{\text{lin}} = X^\dagger y$.

--- Ordinary least squares (OLS) algorithm

- Is this “learning” or just an analytic solution?
 - Well, LS methods was published by Legendre in 1805 and by Gauss in 1809, long before the idea of “machine learning” was born.
 - As long as it gets a model from the data, we are happy to call it “learning”.



Xuegong Zhang



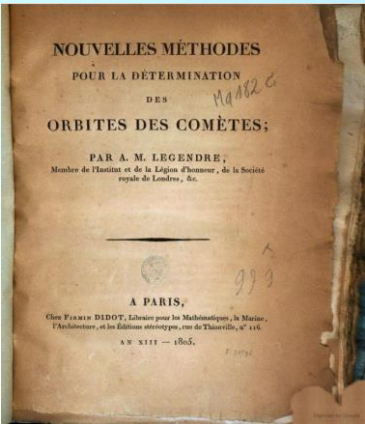
A.M. Legendre. Nouvelles méthodes pour la détermination des orbites des comètes, Firmin Didot, Paris, 1805. “Sur la Méthode des moindres carrés” appears as an appendix.

Adrien-Marie Legendre



1820 watercolor caricature of Adrien-Marie Legendre by French artist Julien-Leopold Boilly (see portrait debate), the only existing portrait known.^[1]

Born 18 September 1752
Paris, France
Died 10 January 1833 (aged 80)
Paris, France



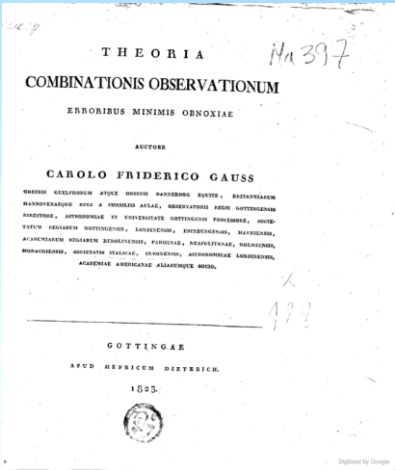
C.F. Gauss. Theoria Motus Corporum Coelestium in Sectionibus Conicis Solem Ambientum. (1809)

Carl Friedrich Gauss



Carl Friedrich Gauß (1777–1855), painted by Christian Albrecht Jensen

Born Johann Carl Friedrich Gauss
30 April 1777
Brunswick, Principality of Brunswick-Wolfenbüttel
Died 23 February 1855 (aged 77)
Göttingen, Kingdom of Hanover, German Confederation



Xuegong Zhang

Any question?



“ $\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$, if $(\mathbf{X}^T \mathbf{X})$ is invertible.”

- What is “invertible”? What happens if it is not?

- Invertible (nonsingular, nondegenerate, full rank)

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_N^T \end{bmatrix} = \begin{bmatrix} & & & \\ & & & \\ & & & \\ & & & \end{bmatrix}_{N \times (d+1)}$$

- $(\mathbf{X}^T \mathbf{X})$ invertible: \mathbf{X} full column rank : features linearly independent

- Usually true when $N \gg d + 1$

- When features not linearly independent, pseudo-inverse can still be defined, but solution is not unique.

- Solutions:

- Remove redundancy by feature selection or transformation
- Introduce other criteria (e.g., SVD or regularization) to restrict the solution



Xuegong Zhang

41

Recall: Similar assumption in FLD



- Fisher's criterion: $\max \quad J_F(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_b \mathbf{w}}{\mathbf{w}^T \mathbf{S}_w \mathbf{w}}$

Solution:

$$\begin{aligned} & \max \mathbf{w}^T \mathbf{S}_b \mathbf{w} \\ & \text{s. t. } \mathbf{w}^T \mathbf{S}_w \mathbf{w} = c \end{aligned}$$

Lagrange function

$$L(\mathbf{w}, \lambda) = \mathbf{w}^T \mathbf{S}_b \mathbf{w} - \lambda(\mathbf{w}^T \mathbf{S}_w \mathbf{w} - c)$$

Let $\frac{\partial L}{\partial \mathbf{w}} = 0$, we get $\mathbf{S}_w^{-1} \mathbf{S}_b \mathbf{w}^* = \lambda \mathbf{w}^*$

and

$$\mathbf{w}^* = \mathbf{S}_w^{-1}(\mathbf{m}_1 - \mathbf{m}_2)$$

Xuegong Zhang

42

How to evaluate the regression?

Xuegong Zhang

43

Evaluation of regression models

R^2 : the **goodness-of-fit**, the **coefficient of determination**, ...



$$R^2 = 1 - \frac{\sum_{i=1}^N (y_i - \hat{y}_i)^2}{\sum_{i=1}^N (y_i - \bar{y})^2}$$

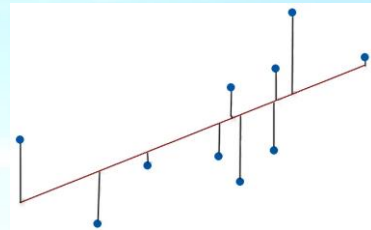
Unexplained variation

Total variation

For OLS,

$$R^2 = \frac{\sum_{i=1}^N (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^N (y_i - \bar{y})^2}$$

and $0 \leq R^2 \leq 1$.



Xuegong Zhang

44

Evaluation of regression models



R^2 : the **goodness-of-fit**, the **coefficient of determination**, ...

- For OLS, $0 \leq R^2 \leq 1$.

$R^2 = 1$: Perfect regression.

$R^2 = 0$: Baseline model. Predictions are the average.

- For other types of regression as a general measure of goodness-of-fit

$$R^2 = 1 - \frac{\sum_{i=1}^N (y_i - \hat{y}_i)^2}{\sum_{i=1}^N (y_i - \bar{y})^2} \in [-\infty, 1]$$

and should no longer be called R^2 .

Xuegong Zhang

45

R^2 is not enough for evaluating regression

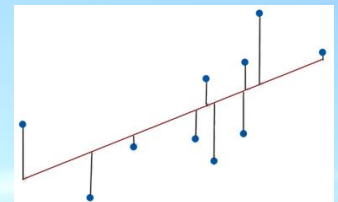


$$y_i = w_0 + \mathbf{w}^T \mathbf{x}_i + \epsilon_i, \quad i = 1, \dots, N$$

Dependent
variable

Fitted values
(deterministic)

Error, residual, noise
(stochastic)



- R^2 : percentage of dependent variable variations that the linear model explains
- R^2 does not indicate if the regression model provides an adequate fit to the data

Xuegong Zhang

46

R^2 is not enough for evaluating regression



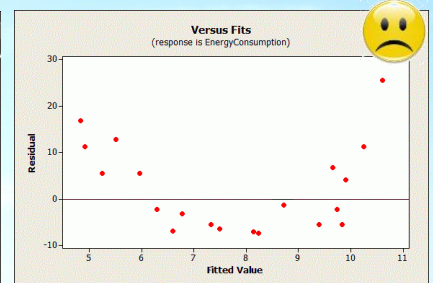
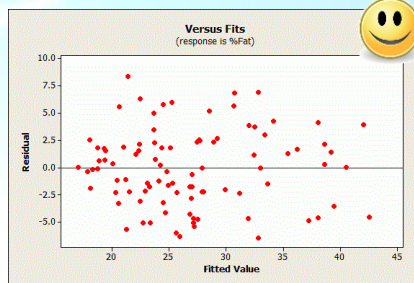
$$y_i = w_0 + \mathbf{w}^T \mathbf{x}_i + \epsilon_i, \quad i = 1, \dots, N$$



- R^2 : percentage of dependent variable variations that the linear model explains
- R^2 does not indicate if the regression model provides an adequate fit to the data

- **Residual plots**

- To check if model is adequate
- Poor fitting if error/residual not random



Xuegong Zhang

47

Evaluating each coefficient



$$y_i = w_0 + \mathbf{w}^T \mathbf{x}_i + \epsilon_i = w_0 + \sum_{j=1}^d w_j x_{ij} + \epsilon_i, \quad \epsilon_i \sim N(0, \sigma^2), i = 1, \dots, N,$$

Does each w_j contribute?

Test for statistical significance of regression coefficients

$$\frac{\hat{w}_j - w_j}{s_{\hat{w}_j}} \sim t_{N-d-1}, \quad j = 0, 1, \dots, d.$$



Xuegong Zhang

48



1-minute break



Xuegong Zhang

49



3.4 Minimum Squared Error (MSE) Method

Xuegong Zhang

50

Recap on Perceptron

- Objective function: classification error

$$J_P(\alpha) = \sum_{y_j \in Y^k} (-\alpha^T y_j)$$

- The Perceptron Criterion:

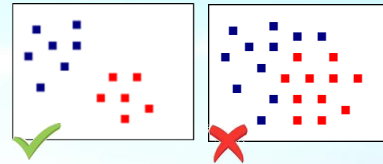
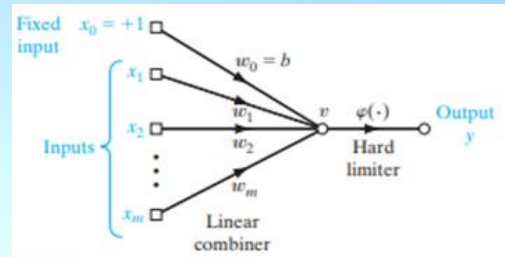
$$J_P(\alpha^*) = \min J_P(\alpha) = 0$$

- The learning algorithm:

- Gradient descent: $\alpha(k+1) = \alpha(k) - \rho_k \nabla J$

- Major limitation:

- Designed for linearly separable cases



Xuegong Zhang

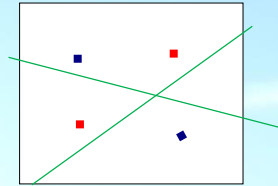
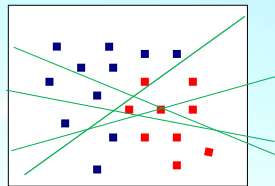
51

Linearly non-separable cases

- How about the following cases?



Not perfectly separable, but can achieve reasonable separation with a linear classifier.



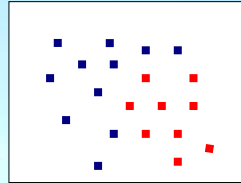
No reasonable linear solution.

Xuegong Zhang

52

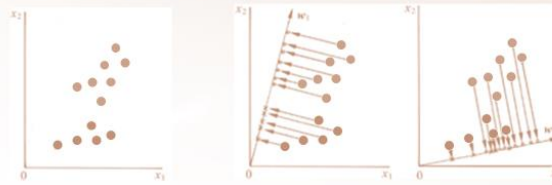


- Can we get a “best” linear classifier when there is no perfect one?



- What is “the best”?
 - A basic principle that differentiates many different methods.

Fisher's criterion for “the best” discriminant



- Which projection is better? Why and How?
- Fisher: The best projection: samples of the two classes “best separated”
 - Within-class: as close as possible
 - Between-classes: as far as possible



Xuegong Zhang

53

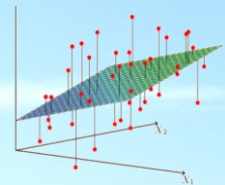
Minimum Squared Error (MSE) method



- How do we calculate and minimize “errors”?
- We knew how to minimize mean squared error for linear regression

$$\min E = \frac{1}{N} \sum_{j=1}^N (f(x_j) - y_j)^2$$

- Our goal was to find the \mathbf{w} that $y = \mathbf{w}^T \mathbf{x}$ for all samples
- Now our goal is to find the \mathbf{w} that $\alpha^T \mathbf{y} > 0$ for all samples
 - well, for as many samples as possible.



- How about to define a $b_j > 0$ for each sample so that

$$\alpha^T \mathbf{y}_j = b_j, \quad j = 1, \dots, N$$

for as many sample as possible?

Xuegong Zhang

54



Resuming the setting of perceptron

$$g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = \boldsymbol{\alpha}^T \mathbf{y}$$

$\mathbf{y} = \begin{bmatrix} 1 \\ \mathbf{x} \end{bmatrix}$ augmented feature vector

$\boldsymbol{\alpha} = \begin{bmatrix} w_0 \\ \mathbf{w} \end{bmatrix}$ augmented weight vector

To simplify the treatment, we normalize the function by the following replacement

$$\mathbf{y}'_i = \begin{cases} \mathbf{y}_i, & \text{if } \mathbf{y}_i \in \omega_1, \\ -\mathbf{y}_i, & \text{if } \mathbf{y}_i \in \omega_2, \end{cases} \quad i = 1, \dots, N$$

then we have $\boldsymbol{\alpha}^T \mathbf{y}'_i > 0, i = 1, \dots, N$ if all samples are correctly classified. We still use \mathbf{y}_i to represent a “*normalized augmented feature vector*” for convenience.

Xuegong Zhang

55

The MSE Criterion



$$\boldsymbol{\alpha}^T \mathbf{y}_i > 0 \quad \Longleftrightarrow \quad \boldsymbol{\alpha}^T \mathbf{y}_i = b_i > 0, i = 1, \dots, N$$

- From inequation to equation

$$\mathbf{Y}\boldsymbol{\alpha} = \mathbf{b}, \quad \mathbf{b} = [b_1, b_2, \dots, b_N]^T$$

- The MSE Criterion

$$\boldsymbol{\alpha}^*: \min_{\boldsymbol{\alpha}} J_S(\boldsymbol{\alpha})$$

$$J_S(\boldsymbol{\alpha}) = \|\mathbf{Y}\boldsymbol{\alpha} - \mathbf{b}\|^2 = \sum_{i=1}^N (\boldsymbol{\alpha}^T \mathbf{y}_i - b_i)^2$$

Xuegong Zhang

56



The MSE Criterion

$$\min J_S(\alpha) = \|\mathbf{Y}\alpha - \mathbf{b}\|^2 = \sum_{i=1}^N (\alpha^T \mathbf{y}_i - b_i)^2$$

- Solutions

- Pseudo-Inverse

$$\alpha^* = (\mathbf{Y}^T \mathbf{Y})^{-1} \mathbf{Y}^T \mathbf{b} = \mathbf{Y}^+ \mathbf{b},$$

$$\mathbf{Y}^+ = (\mathbf{Y}^T \mathbf{Y})^{-1} \mathbf{Y}^T$$

- Gradient Descent

$$\nabla J_S(\alpha) = 2\mathbf{Y}^T (\mathbf{Y}\alpha - \mathbf{b})$$

$$\alpha(k+1) = \alpha(k) + \rho_k (b_k - \alpha(k)^T \mathbf{y}^k) \mathbf{y}^k$$

--- Widrow-Hoff Algorithm
or the **ADALINE**



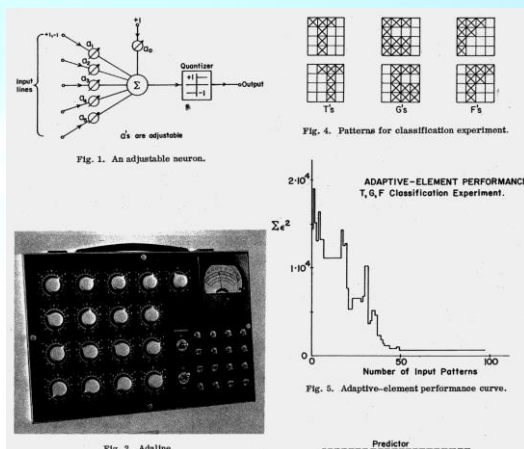
Xuegong Zhang

57

ADALINE



Widrow & Hoff, Adaptive switching circuits, 1960 IRE Western Electric
Show and Convention Record, Part 4, pp.96-104, Aug, 1960



Widrow & Hoff, Adaptive switching circuits, 1960 IRE Western
Electric Show and Convention Record, Part 4, pp.96-104, Aug, 1960

Xuegong Zhang

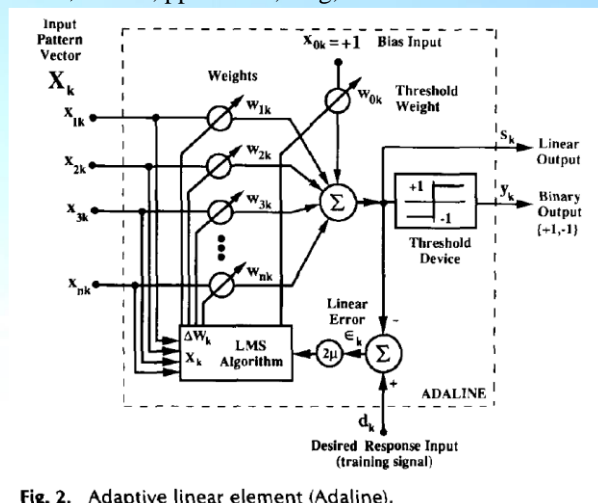


Fig. 2. Adaptive linear element (Adaline).

Widrow & Lehr, 30 years of adaptive neural networks: Perceptron, Madaline, and
Backpropagation, *Proceedings of the IEEE*, 78(9): 1415-1442, 1990

58

ADALINE / Widrow-Hoff Algorithm / LMS Algorithm



Algorithm: Widrow-Hoff (ADALINE)

- ① Normalize the augmented feature vectors \mathbf{z}_j of all training samples by

$$\mathbf{z}'_i = \begin{cases} \mathbf{z}_i, & \text{if } \mathbf{z}_i \in \omega_1, \\ -\mathbf{z}_i, & \text{if } \mathbf{z}_i \in \omega_2, \end{cases} \quad i = 1, \dots, N$$

- ② Initialization: Set $k = 0$; Set initial weights to all zeros $\alpha(0) = \mathbf{0}$; Set proper target values b_i for all samples;

- ③ Pick up a sample \mathbf{z}_j from the training set, compute the gradient and update the weight

$$\alpha(k+1) = \alpha(k) + \rho_k (b_j - \alpha(k)^T \mathbf{z}_j) \mathbf{z}_j.$$

- ④ Let $k = k + 1$. Repeat ② for all samples till the stopping criterion is met.

Xuegong Zhang

59

Any questions?

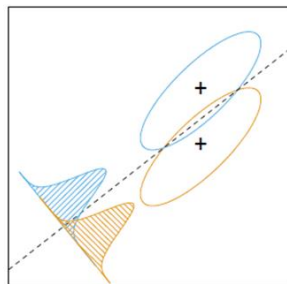


How do we assign the \mathbf{b} ?

- It can be proven that if we assign the \mathbf{b} as

$$b_i = \begin{cases} N/N_1, & \text{if } \mathbf{y}_i \in \omega_1 \\ N/N_2, & \text{if } \mathbf{y}_i \in \omega_2 \end{cases},$$

then the MSE solution is equal to the FLD solution with $w_0 = \hat{m}$.



Xuegong Zhang

60



The “best”?

- It can also be proven that if we assign the \mathbf{b} as

$$b_i = 1, i = 1, \dots, N,$$

then when $N \rightarrow \infty$, the MSE solution best approximates the Bayesian Discriminant

$$g_0(\mathbf{x}) = P(\omega_1|\mathbf{x}) - P(\omega_2|\mathbf{x})$$

with minimal square error $e^2 = \int [\boldsymbol{\alpha}^T \mathbf{y} - g_0(\mathbf{x})]^2 p(\mathbf{x}) d(\mathbf{x})$.

the optimal classifier when class distribution densities are known

Let's leave it for future classes.

Xuegong Zhang

61

单选题 1分

设置

4-min break

- ☒ A Yes, I'm back to my computer.
- ☐ B Sorry, not yet.



Xuegong Zhang

提交

62



3.5 Logistic Regression

Xuegong Zhang

63



Regression for binary outcome?

$y = f(x)$

prediction ↑

real function ↑

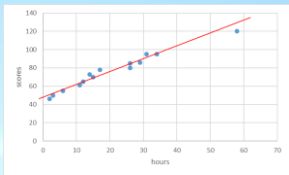
sample ↑

$X = [x_1, \dots, x_m] = \begin{bmatrix} x_{11} & \dots & x_{m1} \\ \vdots & \ddots & \vdots \\ x_{1n} & \dots & x_{mn} \end{bmatrix}$

← features

↑ A sample

Student id	Final score	Study Hours per Week
1	50	3
2	95	34
3	78	17
4	55	6
5	65	12
6	70	15
7	80	26
8	86	29
9	73	14
10	120	58
11	46	2
12	95	31
13	85	26
14	61	11



Simple Learn Regression
 $y = w_0 + w_1x$

Xuegong Zhang

64



Regression for binary outcome?

$$y = f(x)$$

prediction ↑

binary function ↑

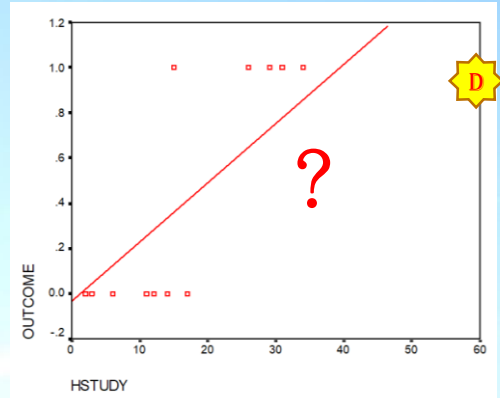
sample ↑

$$X = [x_1, \dots, x_m] = \begin{bmatrix} x_{11} & \dots & x_{m1} \\ \vdots & \ddots & \vdots \\ x_{1n} & \dots & x_{mn} \end{bmatrix}$$

← features

A sample →

Student id	Final score	Study Hours per Week
1	0	3
2	1	34
3	0	17
4	0	6
5	0	12
6	1	15
7	1	26
8	1	29
9	0	14
10	1	58
11	0	2
12	1	31
13	1	26
14	0	11



Xuegong Zhang

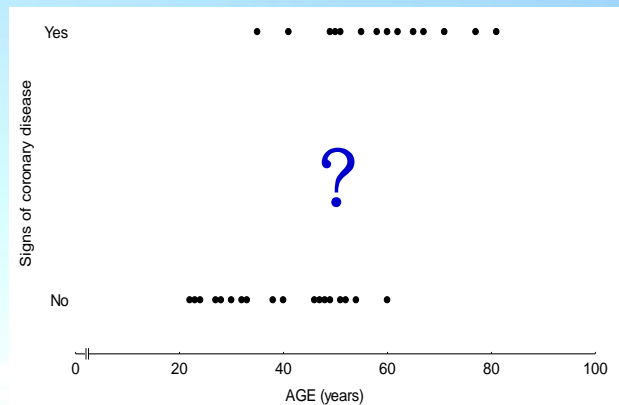
65



Another example

Age and signs of coronary heart disease (CD)

Age	CD	Age	CD	Age	CD
22	0	40	0	54	0
23	0	41	1	55	1
24	0	46	0	58	1
27	0	47	0	60	1
28	0	48	0	60	0
30	0	49	1	62	1
30	0	49	0	65	1
32	0	50	1	67	1
33	0	51	0	71	1
35	1	51	1	77	1
38	0	52	0	81	1



Xuegong Zhang

66



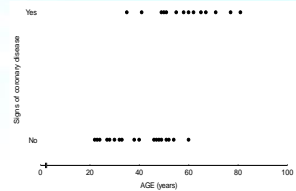
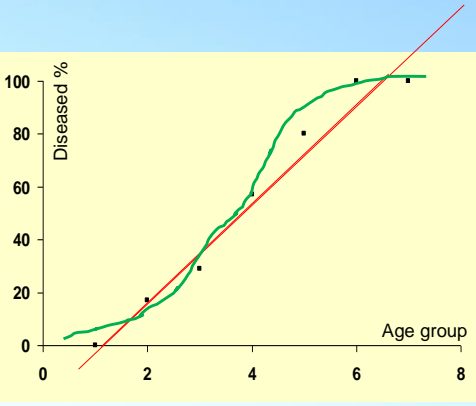
Let's count the percentage

Age and signs of coronary heart disease (CD)

Age	CD	Age	CD	Age	CD
22	0	40	0	54	0
23	0	41	1	55	1
24	0	46	0	58	1
27	0	47	0	60	1
28	0	48	0	60	0
30	0	49	1	62	1
30	0	49	0	65	1
32	0	50	1	67	1
33	0	51	0	71	1
35	1	51	1	77	1
38	0	52	0	81	1

Prevalence (%) of signs of CD according to age group

Age group	# in group	Diseased	
		#	%
20 - 29	5	0	0
30 - 39	6	1	17
40 - 49	7	2	29
50 - 59	7	4	57
60 - 69	5	4	80
70 - 79	2	2	100
80 - 89	1	1	100

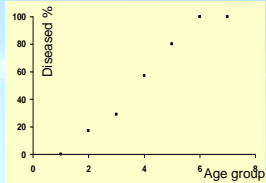
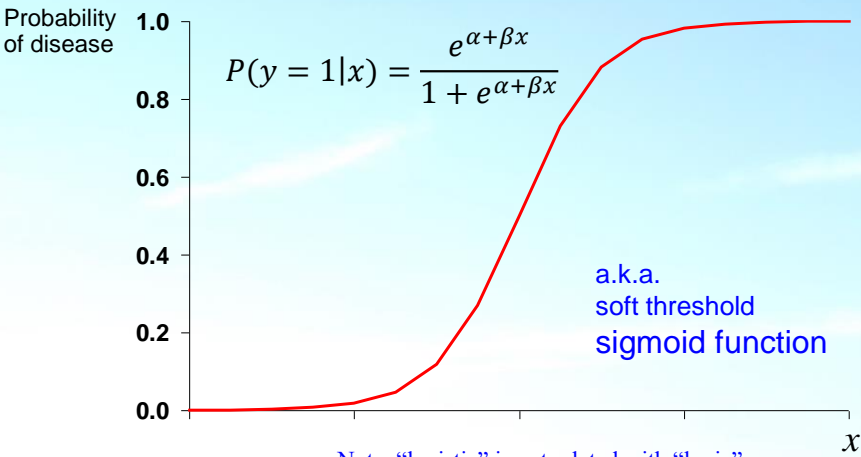


Xuegong Zhang

67



Logistic function



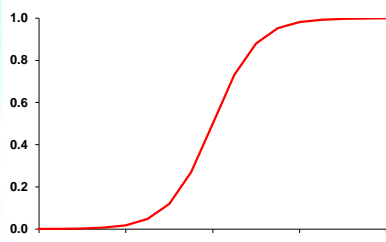
Note: “logistic” is not related with “logic”.
So I prefer the translation 罗杰斯特 or 罗杰斯蒂 .

Xuegong Zhang

68



Logistic regression



$$P(y|x) = \frac{e^{\alpha + \beta x}}{1 + e^{\alpha + \beta x}}$$

↑ Inverse function

$$\ln \left[\frac{P(y|x)}{1 - P(y|x)} \right] = \alpha + \beta x$$

logit of $P(y|x)$

---- log odds 对数几率

Xuegong Zhang

69



Multiple logistic regression

$$P = \frac{e^{w_0 + w_1 x_1 + w_2 x_2 + \dots + w_p x_p}}{1 + e^{w_0 + w_1 x_1 + w_2 x_2 + \dots + w_p x_p}} = \frac{e^{w^T x}}{1 + e^{w^T x}} \triangleq \theta(w^T x)$$

$$odds = \frac{P}{1 - P} = e^{w_0 + w_1 x_1 + w_2 x_2 + \dots + w_p x_p} = e^{w^T x}$$

$$\ln \left(\frac{P}{1 - P} \right) = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_p x_p = w^T x$$

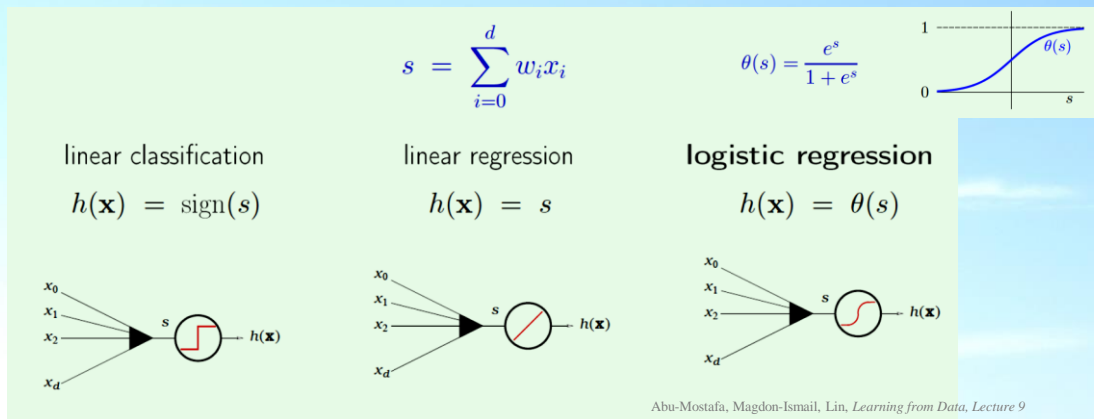
- Interpretation of w_i

- Increase in log-odds for a one unit increase in x_i with all other x_i 's constant
- Measures association between x_i and log-odds adjusted for all other x_i

Xuegong Zhang

70

Three types of linear machines



$h(\mathbf{x}) = \theta(s)$ is the estimate of the **probability** of $y = 1$.
 $s = \mathbf{w}^T \mathbf{x}$ is the signal of the event ("risk score")

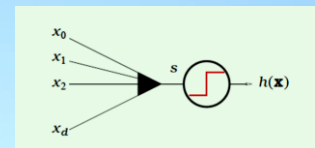
Xuegong Zhang

71

Basic concepts of ML: Perceptron



- How can we make a learning machine?
 - It needs a teacher.
 - The model: $y = \text{sgn}(\sum_{i=1}^d w_i x_i + w_0)$
 - We need materials to train it. / It needs materials to learn from.
 - Training data: $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}, \mathbf{x}_j \in \mathbb{R}^{d+1}, y_j \in \{-1, 1\}$
 - We need to tell what is the goal of the learning.
 - Objective function: $\min J_P(\boldsymbol{\alpha}) = \sum_{\mathbf{y}_j \in \mathcal{Y}^k} (-\boldsymbol{\alpha}^T \mathbf{y}_j)$
 - We need to tell it how to learn.
 - Learning algorithm: $\boldsymbol{\alpha}(k+1) = \boldsymbol{\alpha}(k) - \rho_k \nabla J = \boldsymbol{\alpha}(k) + \rho_k \sum_{\mathbf{y}_j \in \mathcal{Y}^k} \mathbf{y}_j$



Xuegong Zhang

72



Basic concepts of ML: Linear Regression

- How can we make a learning machine?

- It needs a teacher.

- The model: $f(\mathbf{x}) = \sum_{i=0}^d w_i x_i = \mathbf{w}^T \mathbf{x}$

- We need materials to train it. / It needs materials to learn from.

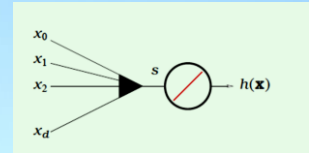
- Training data: $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}, \mathbf{x}_j \in R^{d+1}, y_j \in R$

- We need to tell what is the goal of the learning.

- Objective function: $\min E = \frac{1}{N} \sum_{j=1}^N (f(\mathbf{x}_j) - y_j)^2$

- We need to tell it how to learn.

- Learning algorithm: $\mathbf{w}(k+1) = \mathbf{w}(k) - \rho_k \nabla E$



Xuegong Zhang

73



Basic concepts of ML: Logistic Regression

- How can we make a learning machine?

- It needs a teacher.

- The model: $h(\mathbf{x}) = \theta(\mathbf{w}^T \mathbf{x})$

- We need materials to train it. / It needs materials to learn from.

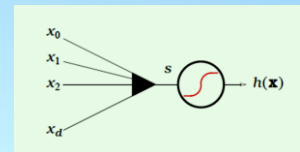
- Training data: $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}, \mathbf{x}_j \in R^{d+1}, y_j \in \{-1, 1\}$

- We need to tell what is the goal of the learning.

- Objective function: ?

- We need to tell it how to learn.

- Learning algorithm: ?



Xuegong Zhang

74



Reasoning behind logistic regression

- Data $\{(x_1, y_1), \dots, (x_N, y_N)\}$, $x_j \in R^{d+1}$, $y_j \in \{-1, 1\}$ were generated i.i.d. by the probability $f(x)$:

$$P(y|x) = \begin{cases} f(x) & \text{for } y = +1 \\ 1 - f(x) & \text{for } y = -1 \end{cases}$$

A generative model

- We use $h(x) = \theta(w^T x)$ to estimate $f(x)$.
- How to estimate the parameters of $h(x) = \theta(w^T x)$?



Xuegong Zhang

75



Reasoning behind logistic regression

- Likelihood** (n. 似然度) :

– For instance (x_j, y_j) of data, if $h = f$, how likely we get y_j from x_j ?

$$P(y_j|x_j) = \begin{cases} h(x_j) & \text{for } y_j = +1 \\ 1 - h(x_j) & \text{for } y_j = -1 \end{cases}$$

– In other words, now we have this instance of data, how likely h is the “model” that generated the data?



- Noting $\theta(-s) = 1 - \theta(s)$, we have the likelihood on x_j :

$$P(y_j|x_j) = \theta(y_j w^T x_j)$$

Xuegong Zhang

76

单选题 10分

设置

The likelihood $P(y_j|x_j) = \theta(y_j \mathbf{w}^T \mathbf{x}_j)$ is a function of whom?

- ☐ A of the output y_i
- ☒ B of the weights \mathbf{w}
- ☐ C of the features \mathbf{x}_i
- ☐ D of the sample (\mathbf{x}_j, y_j)

提交

77

Maximizing the likelihood



- Likelihood of getting the i.i.d. data $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$, $\mathbf{x}_j \in R^{d+1}, y_j \in \{-1, 1\}$ from the model (likelihood of the model on the data):

$$L(\mathbf{w}) = \prod_{j=1}^N P(y_j|x_j) = \prod_{j=1}^N \theta(y_j \mathbf{w}^T \mathbf{x}_j)$$

$$P(y_j|x_j) = \theta(y_j \mathbf{w}^T \mathbf{x}_j)$$



- Likelihood Maximization:

$$\min \quad E(\mathbf{w}) = -\frac{1}{N} \ln(L(\mathbf{w})) = -\frac{1}{N} \ln \left(\prod_{j=1}^N \theta(y_j \mathbf{w}^T \mathbf{x}_j) \right)$$

Xuegong Zhang

78



Maximizing the likelihood

- Likelihood Maximization:

$$\begin{aligned} \min \quad E(\mathbf{w}) &= -\frac{1}{N} \ln(L(\mathbf{w})) = -\frac{1}{N} \ln \left(\prod_{j=1}^N \theta(y_j \mathbf{w}^T \mathbf{x}_j) \right) \\ &= \frac{1}{N} \sum_{j=1}^N \ln \left(\frac{1}{\theta(y_j \mathbf{w}^T \mathbf{x}_j)} \right) = \frac{1}{N} \sum_{j=1}^N \ln \left(1 + e^{-y_j \mathbf{w}^T \mathbf{x}_j} \right) \end{aligned}$$



$$\left[\theta(s) = \frac{1}{1 + e^{-s}} \right]$$

- Gradient:

$$\nabla E = -\frac{1}{N} \sum_{j=1}^N \frac{y_j \mathbf{x}_j}{1 + e^{y_j \mathbf{w}^{(k)T} \mathbf{x}_j}}$$

Xuegong Zhang

79



Logistic Regression Algorithm

- Set $k = 0$, initialize $\mathbf{w}(0)$
- Do
 - Compute the gradient $\nabla E = -\frac{1}{N} \sum_{j=1}^N \frac{y_j \mathbf{x}_j}{1 + e^{y_j \mathbf{w}^{(k)T} \mathbf{x}_j}}$
 - Update the weights $\mathbf{w}(k+1) = \mathbf{w}(k) - \eta \nabla E$, set $k = k + 1$
 Until the stopping criterion met
- Return the final weights \mathbf{w}

- Initialization:**

All zeros work, but safer to initialize weights randomly, say, normal distribution with 0-mean and small variance

- Termination:**

An upper bound on iterations, or a threshold of the gradient



Xuegong Zhang

Abu-Mostafa, Magdon-Ismael, Lin, *Learning from Data, Lecture 9*

80

Basic concepts of ML: Logistic Regression



- How can we make a learning machine?
 - It needs a teacher.
 - The model: $h(\mathbf{x}) = \theta(\mathbf{w}^T \mathbf{x})$
 - We need materials to train it. / It needs materials to learn from.
 - Training data: $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}, \mathbf{x}_j \in R^{d+1}, y_j \in \{-1, 1\}$
 - We need to tell what is the goal of the learning.
 - Objective function: $\min E(\mathbf{w}) = \frac{1}{N} \sum_{j=1}^N \ln(1 + e^{-y_j \mathbf{w}^T \mathbf{x}_j})$
 - We need to tell it how to learn.
 - Learning algorithm: $\mathbf{w}(k+1) = \mathbf{w}(k) - \rho_k \nabla E$

Xuegong Zhang

81

3.6 Discussion



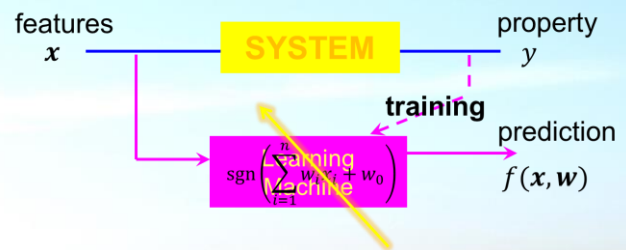
Xuegong Zhang

82

The basic elements of supervised learning



- It needs a teacher.
 - We (people) **design** it (features, the model) and **train** it.
- We need materials to train it. / It needs materials to learn from.
 - Training data
- We need to tell what is the goal of the learning.
 - Objective function/criterion
- We need to tell it how to learn.
 - Learning/training algorithm



Xuegong Zhang

83

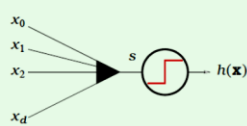
Three major types of linear machines



$$s = \sum_{i=0}^d w_i x_i$$

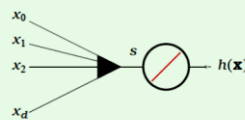
linear classification

$$h(\mathbf{x}) = \text{sign}(s)$$



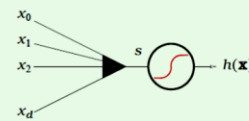
linear regression

$$h(\mathbf{x}) = s$$



logistic regression

$$h(\mathbf{x}) = \theta(s)$$

Abu-Mostafa, Magdon-Ismael, Lin, *Learning from Data*, Lecture 9

Xuegong Zhang

84

Optimization Problems



- For perceptron

$$\min J_P(\alpha) = \sum_{y_j \in Y^k} (-\alpha^T y_j)$$

Iterative learning

- For linear regression

$$\min E(\mathbf{w}) = \frac{1}{N} \sum_{j=1}^N (\mathbf{w}^T \mathbf{x}_j - y_j)^2$$

Closed-form solution
or iterative learning

- For logistic regression

$$\min E(\mathbf{w}) = \frac{1}{N} \sum_{j=1}^N \ln(1 + e^{-y_j \mathbf{w}^T \mathbf{x}_j})$$

Iterative learning

Xuegong Zhang

85

Optimization via Gradient Descent

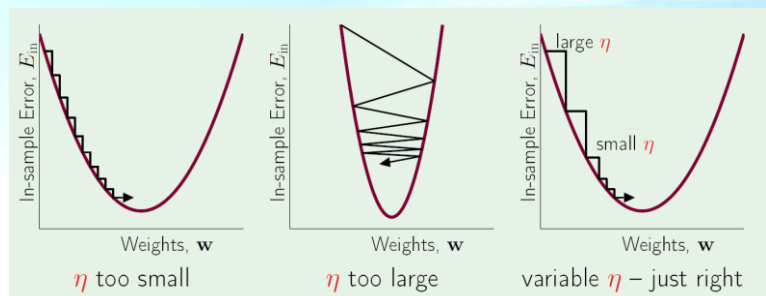
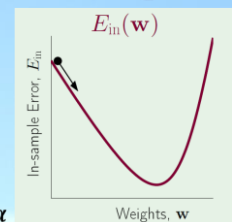
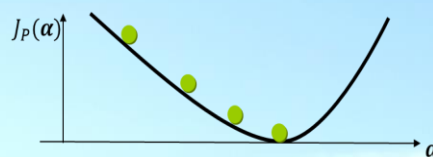


- General method of gradient descent:

$$\mathbf{w}(k+1) = \mathbf{w}(k) + \eta \hat{\mathbf{v}}$$

$$\hat{\mathbf{v}} = -\nabla E(\mathbf{w}(k))$$

- η : learning rate (step size)



Abu-Mostafa, Magdon-Ismail, Lin, *Learning from Data*, Lecture 9

Xuegong Zhang

86



Discussions

- Machines that can learn
 - Different types of machines for different tasks
 - Different assumptions for the same task
 - Shared general principles and techniques
 - One-size-fit-all solution?
 - Common theory for all machines?

Xuegong Zhang

87



Homework

- Problems (Pr2)
 1. (Op) Inference on smoking-health
 2. FLD and MSE
 3. Perceptron convergence
 4. Logistic function properties
- Deadline:
 - Sept. 29, 23:00 Beijing Time
- Computer exercises (Ex1)
 - Coding FLD, Perceptron and Logistic Regression in Python
 - Experimenting on a medical dataset
- Deadline:
 - Oct. 6, 23:00 Beijing Time



Xuegong Zhang

88

单选题 1分

设置

4-min break

- ☒ A Yes, I'm back to my computer.
- ☐ B Sorry, not yet.

