

LINEAR MODEL

Bùi Tiến Lên

01/09/2019



KHOA CÔNG NGHỆ THÔNG TIN
TRƯỜNG ĐẠI HỌC KHOA HỌC TỰ NHIÊN



1. Linear Classification

2. Linear Regression

Simple Linear Model

Linear Basis Function Models

3. Capacity, Overfitting and Underfitting

Model Comparision

4. Logistic Regression

Notation



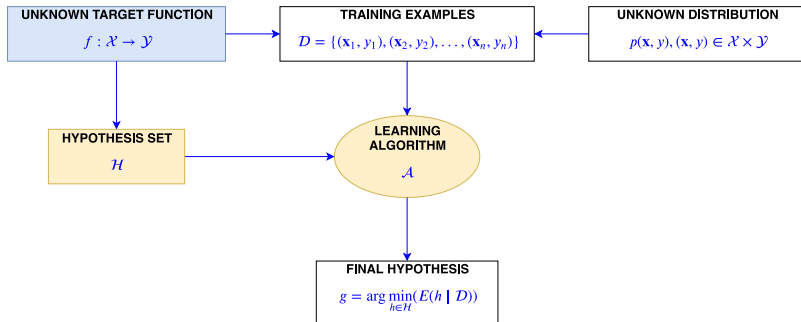
symbol	meaning
$a, b, c, N \dots$	scalar number
$\mathbf{w}, \mathbf{v}, \mathbf{x}, \mathbf{y} \dots$	column vector
$\mathbf{X}, \mathbf{Y} \dots$	matrix
\mathbb{R}	set of real numbers
\mathbb{Z}	set of integer numbers
\mathbb{N}	set of natural numbers
\mathbb{R}^D	set of vectors
$\mathcal{X}, \mathcal{Y}, \dots$	set
\mathcal{A}	algorithm

operator	meaning
\mathbf{w}^T	transpose
$\mathbf{X}\mathbf{Y}$	matrix multiplication
\mathbf{X}^{-1}	inverse

Learning Goal



- Learning diagram revisited



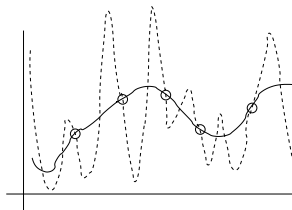


Inductive Bias

Theorem 1 (No Free Lunch Theorems)

An unbiased learner can never generalize.

- The converse of the Inductive Learning Hypothesis is that generalization only possible if we make some assumptions, or introduce some priors. We need an Inductive Bias.
- **Consider:** arbitrarily wiggly functions or random truth tables or non-smooth distributions.



0 0 0	0
0 0 1	?
0 1 0	1
0 1 1	1
1 0 0	0
1 0 1	?
1 1 0	1
1 1 1	?

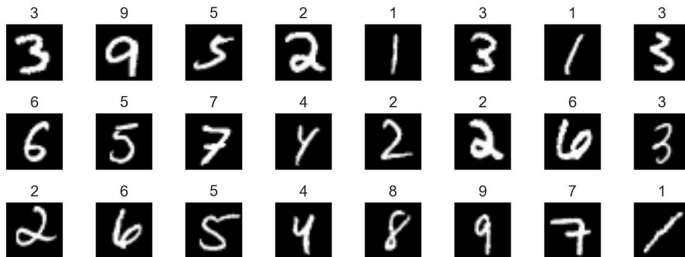


Linear Classification

A real data set



- Some 16-by-16 pixel grayscale image from the MNIST database





Input representation

Input representation or feature extraction

- “raw” input

$$\begin{array}{ll} \text{pixels} & \mathbf{x}^T = (x_0 \quad x_1 \quad \dots \quad x_{256}) \\ \text{linear model} & \mathbf{w}^T = (\mathbf{w}_0 \quad \mathbf{w}_1 \quad \dots \quad \mathbf{w}_{256}) \end{array}$$

- **Feature extraction:** extract useful information

$$\begin{array}{ll} \text{intensity and symmetry} & \mathbf{x}^T = (x_1 \quad x_2) \\ \text{linear model} & \mathbf{w}^T = (\mathbf{w}_1 \quad \mathbf{w}_2) \end{array}$$

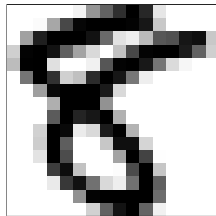
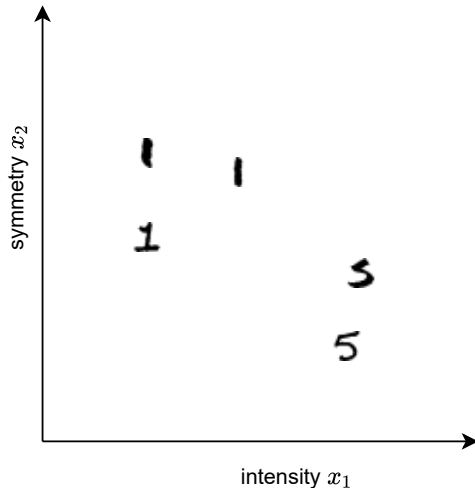


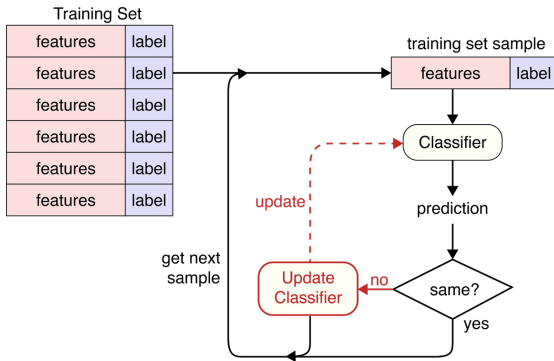
Illustration of Features



Classifier Training



- **Select** the learning model for **classifier**, e.g., Perceptron
- **Train** the classifier



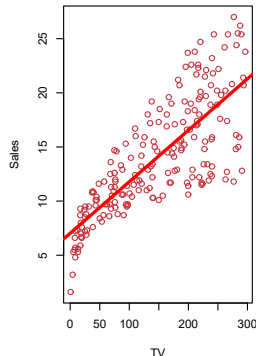
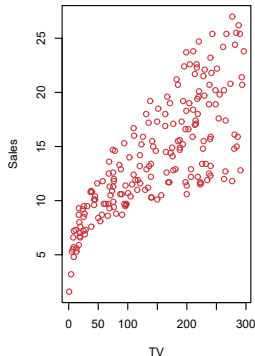


Linear Regression



Simple Problem

- Consider the Advertising data set \mathcal{D}_{train} consists of the **sales** of that product in 200 different markets, along with advertising budgets for the product in each of those markets for the media **TV**. Find the *relationship* between **TV** and **sales**





Solving Problem by Learning

- **The requirement** is to build a system that can take a vector $\mathbf{x} \in \mathbb{R}^{D+1}$ as **input** and **predict** the value of a scalar $y \in \mathbb{R}$ as its **output**
- The hypothesis set \mathcal{H}

$$y \approx \hat{y} = h_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} \quad (1)$$

where \hat{y} be the value that our model (function) predicts y and $\mathbf{w} \in \mathbb{R}^{D+1}$ is a *vector of parameters* of the model



Solving Problem by Learning (cont.)

- **Task T :** to predict y from \mathbf{x} by outputting $\hat{y} = h_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$
- **Performance measure P :**

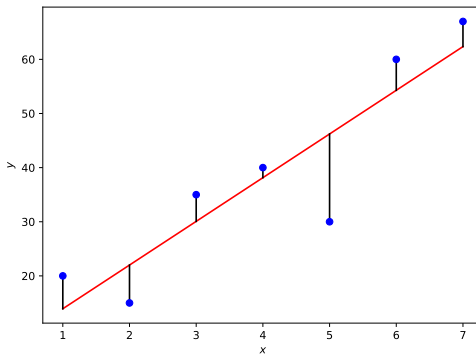
The *mean squared error* MSE_{train} of the model on the train set \mathcal{D}_{train} denoted as (\mathbf{X}, \mathbf{y}) including N samples $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2) \dots (\mathbf{x}_N, y_N)\}$

$$MSE_{train} = \frac{1}{N} \|\hat{\mathbf{y}} - \mathbf{y}\|^2 \quad (2)$$

Construct the matrix \mathbf{X} and the vectors \mathbf{y} and $\hat{\mathbf{y}}$

$$\underbrace{\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_N^T \end{bmatrix}}_{\text{input data matrix}}, \quad \underbrace{\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}}_{\text{target vector}}, \quad \underbrace{\hat{\mathbf{y}} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_N \end{bmatrix}}_{\text{output vector}} \quad (3)$$

Solving Problem by Learning (cont.)



- **The learning goal:** find the vector of parameter \mathbf{w} such that

$$\mathbf{w} = \arg \min_{\mathbf{w}} (MSE_{train})$$



Solving Problem by Learning (cont.)

Solution

- Compute the gradient of MSE_{train}

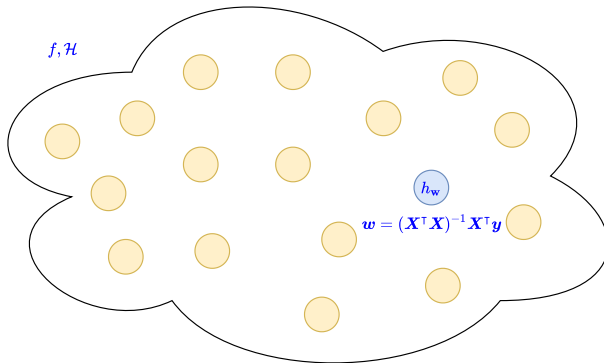
$$\begin{aligned}\nabla_{\mathbf{w}}(MSE_{train}) &= \nabla_{\mathbf{w}}(\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{w}^T \mathbf{X}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \mathbf{w} + \mathbf{y}^T \mathbf{y}) \\ &= 2\mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{X}^T \mathbf{y}\end{aligned}$$

- If MSE_{train} reach the min value then $\nabla_{\mathbf{w}}(MSE_{train}) = 0$

$$\begin{aligned}\nabla_{\mathbf{w}}(MSE_{train}) &= 0 \\ \mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y} &= 0 \\ \mathbf{X}^T \mathbf{X} \mathbf{w} &= \mathbf{X}^T \mathbf{y} \\ \mathbf{w} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}\end{aligned}\tag{4}$$



Solving Problem by Learning (cont.)





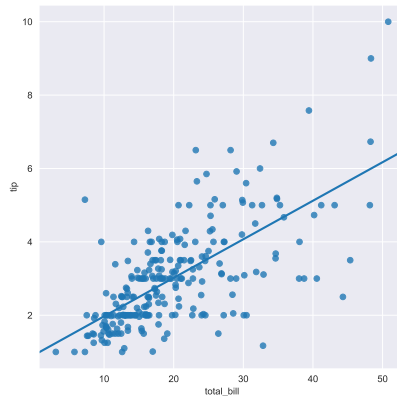
Programming Example

- Use seaborn to read tips dataset and find the linear relationship between total_bill and tip

```
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt

sns.set_style("darkgrid")
tips = sns.load_dataset("tips")
sns.regplot(x="total_bill", y="tip", data=tips, ci=None)
plt.show()
```

Programming Example (cont.)





Word Example

1. Find the linear regression function $y = f(x) = w_0 + w_1x$ given the following data set \mathcal{D}

input x	target y
1	2
2	3
3	3
4	5

2. Find the linear regression function $y = f(\mathbf{x}) = f(x_1, x_2) = w_0 + w_1x_1 + w_2x_2$ given the following data set \mathcal{D}

input \mathbf{x}	target y
(1, 1)	1
(2, 3)	3
(3, 4)	4
(4, 3)	5

Discussion



- D is a large number
- Online learning
- Limitations of **the model (hypothesis set)**



Linear in What?

- **Linearity in the weights**

$$h_{\mathbf{w}}(\mathbf{x}) = w_0 + w_1x_1 + \dots + w_Dx_D \quad (5)$$

- We extend **the model** by introducing **linear combinations** of **fixed nonlinear functions** of the input variables

$$h_{\mathbf{w}}(\mathbf{x}) = w_0 + w_1\phi_1(\mathbf{x}) + \dots + w_M\phi_M(\mathbf{x}) \quad (6)$$

where $\phi_i(\mathbf{x})$ are basis functions



Some types of basis functions

- Power

$$\phi_j(x) = x^j \quad (7)$$

- Gaussian

$$\phi_j(x) = \exp\left(-\frac{(x - \mu_j)^2}{s^2}\right) \quad (8)$$

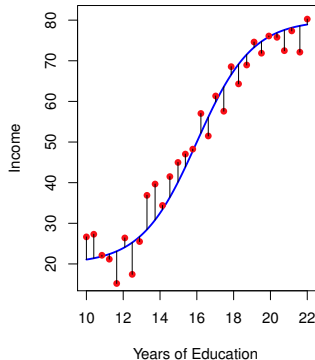
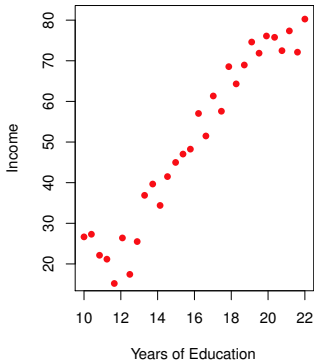
- Sigmoid

$$\phi_j(x) = \frac{1}{1 + e^{-\frac{x - \mu_j}{s}}} \quad (9)$$



Another Problem

- Find the relationship between Years of Education and Income based on the given data





Solving Problem by Learning

- **The requirement** is to build a system that can take a vector $\mathbf{x} \in \mathbb{R}^D$ as **input** and **predict** the value of a scalar $y \in \mathbb{R}$ as its **output**
- The hypothesis set \mathcal{H}

$$y \approx \hat{y} = h_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) \quad (10)$$

where \hat{y} be the value that our model (function) predicts y , $\mathbf{w} \in \mathbb{R}^{M+1}$ is a *vector of parameters* of the model and $\boldsymbol{\phi}$ is a set of $M + 1$ basis functions

$$\boldsymbol{\phi}(\mathbf{x}) = \begin{bmatrix} \phi_0(\mathbf{x}) \\ \phi_1(\mathbf{x}) \\ \vdots \\ \phi_M(\mathbf{x}) \end{bmatrix} \quad (11)$$



Solving Problem by Learning (cont.)

- **Task T :** to predict y from \mathbf{x} by outputting $\hat{y} = h_{\mathbf{w}}(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$
- **Performance measure P :**
The *mean squared error* MSE_{train} of the model on the train set \mathcal{D}_{train} including N samples $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2) \dots (\mathbf{x}_N, y_N)\}$
- **The learning goal:** find the vector of parameter \mathbf{w} such that

$$\mathbf{w} = \arg \min_{\mathbf{w}} (MSE_{train})$$

Solving Problem by Learning (cont.)



1. Construct the matrix Φ and the vectors \mathbf{y}

$$\underbrace{\Phi = \begin{bmatrix} \phi(\mathbf{x}_1)^\top \\ \phi(\mathbf{x}_2)^\top \\ \vdots \\ \phi(\mathbf{x}_N)^\top \end{bmatrix}}_{\text{design matrix}}, \quad \underbrace{\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}}_{\text{target vector}} \quad (12)$$

2. Calculate the vector of parameters

$$\mathbf{w} = (\Phi^\top \Phi)^{-1} \Phi^\top \mathbf{y} \quad (13)$$

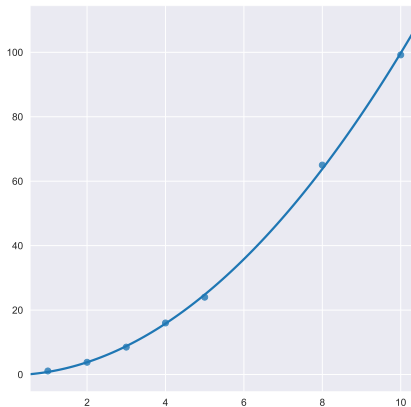
Programming Example



```
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt

sns.set_style("darkgrid")
x = [1, 2, 3, 4, 5, 8, 10]
y = [1.1, 3.8, 8.5, 16, 24, 65, 99.2]
sns.regplot(x, y, order=2, ci=None)
plt.show()
```

Programming Example (cont.)





Word Example

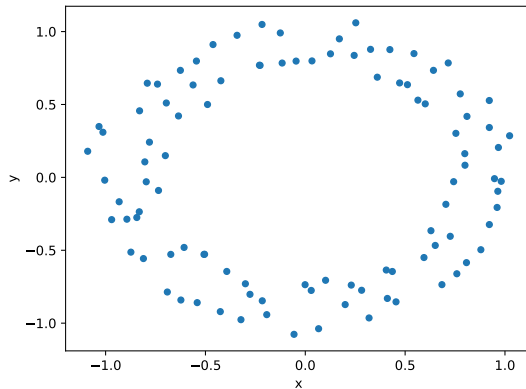
- Find a polynomial regression function $y = f(x) = w_0 + w_1x + w_2x^2$ given the data set \mathcal{D}

input x	target y
1	2
2	3
3	3
4	5

Puzzle



- What basis functions?



Capacity, Overfitting and Underfitting



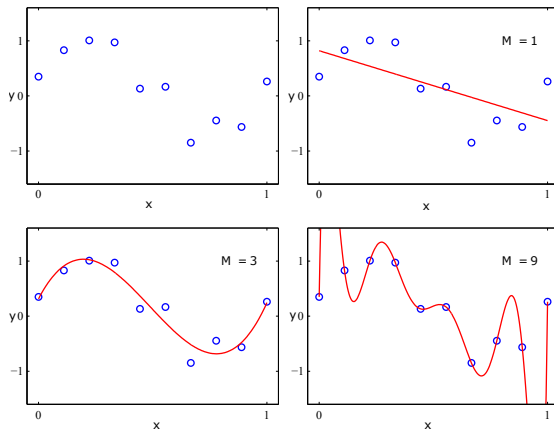


Model Training

We consider three hypothesis sets (polynomial functions) $\mathcal{H}_1, \mathcal{H}_3$ and \mathcal{H}_9 and the results of fitting the models to the data set \mathcal{D}

Which one

- Under-fitting
- Over-fitting
- **Appropriate fitting**



Model Training (cont.)



	$M = 1$	$M = 3$	$M = 9$
w_0	0.82	0.31	0.35
w_1	-1.27	7.99	232.37
w_2		-25.43	-5321.83
w_3		17.37	48568.31
w_4			-231639.30
w_5			640042.26
w_6			-1061800.52
w_7			1042400.18
w_8			-557682.99
w_9			125201.43



Model Performance

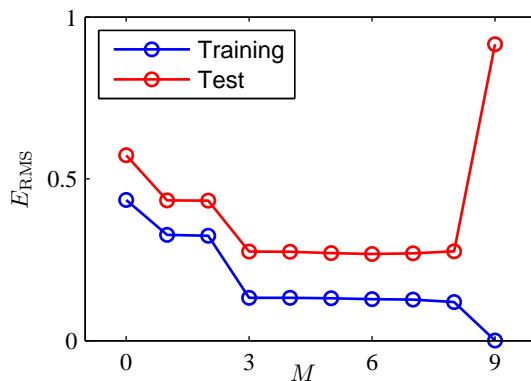


Figure 1: Graphs of the root-mean-square error evaluated on the **training set** and on an **independent test set** for various values of M



What happen if increasing N

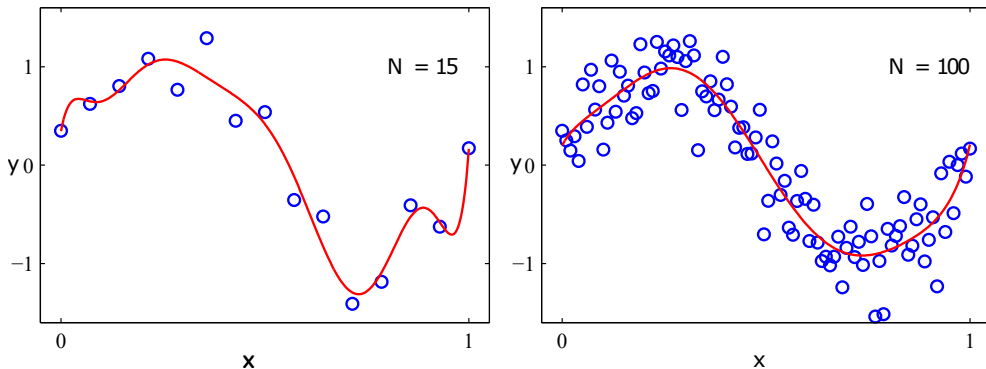


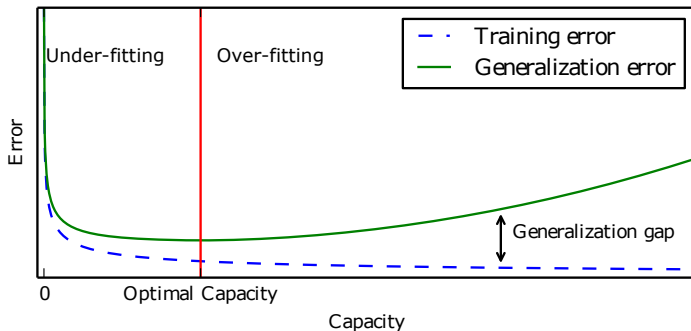
Figure 2: Using the $M = 9$ polynomial for $N = 15$ data points (left plot) and $N = 100$ data points (right plot). We see that increasing the size of the data set reduces the over-fitting problem



Generalization and Capacity

The criteria determining how well a machine learning model will perform:

1. Make the training error small.
2. Make the gap between training and test (generalization) error small.





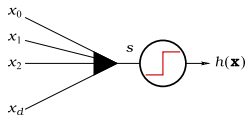
Logistic Regression



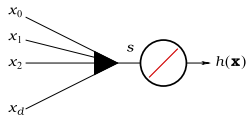
A Third Linear Model

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

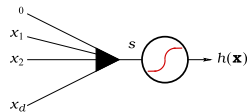
linear classification
 $h(\mathbf{x}) = \text{sign}(s)$



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



logistic regression
 $h(\mathbf{x}) = \sigma(s)$





The logistic function

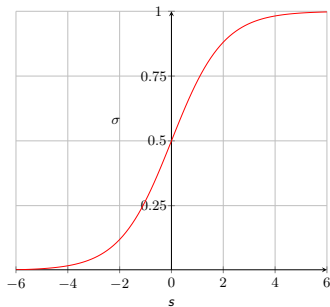
The formula

$$\sigma(s) = \frac{1}{1 + e^{-s}}$$

Some properties

$$\sigma(-s) = 1 - \sigma(s)$$

$$\sigma'(s) = \sigma(s)(1 - \sigma(s))$$





Probability Interpretation

- $h(\mathbf{x}) = \sigma(s)$ can be interpreted as a probability
- For example, prediction of heart attacks
 - Input \mathbf{x} : cholesterol level, age, weight, etc.
 - The signal $s = \mathbf{w}^T \mathbf{x}$: risk score
 - $\sigma(s)$: probability of a heart attack



Problem Statement

- The target function f is the probability

$$f : \mathbb{R}^D \rightarrow [0, 1]$$

- Data $\mathcal{D} = \{(\mathbf{x}_1, y_1) \dots (\mathbf{x}_N, y_N)\}$ with binary $y_i \in \{-1, 1\}$ generated by noisy target

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

- Hypothesis set $h_{\mathbf{w}}(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x})$



Error measure

- For each (\mathbf{x}, y) , y is generated by probability $h_{\mathbf{w}}(\mathbf{x})$.
- Plausible error measure based on **likelihood**
- Likelihood of \mathbf{x} is

$$P(y \mid \mathbf{x}, \mathbf{w}) = z^y (1 - z)^{1-y}$$

where $z = h_{\mathbf{w}}(\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x})$

- Likelihood of $\mathcal{D} = \{(\mathbf{x}_1, y_1) \dots (\mathbf{x}_N, y_N)\}$ is

$$\prod_{n=1}^N P(y_n \mid \mathbf{x}_n, \mathbf{w})$$



Error measure (cont.)

- Maximizing Likelihood

$$\begin{aligned} & \text{Maximize} && \prod_{n=1}^N P(y_n \mid \mathbf{x}_n, \mathbf{w}) \\ \Leftrightarrow & \text{Minimize} && -\log \prod_{n=1}^N P(y_n \mid \mathbf{x}_n, \mathbf{w}) \end{aligned}$$

- We define error measurement (*cross-entropy error*)

$$E(h_{\mathbf{w}}) = - \sum_{n=1}^N (y_n \log z_n + (1 - y_n) \log(1 - z_n)) \quad (14)$$



Learning Algorithm

Linear
Classification

Linear
Regression

Simple Linear Model

Linear Basis Function
Models

Capacity,
Overfitting and
Underfitting

Model Comparison

Logistic
Regression

1. Initialize the weights (parameters) at $t = 0$ \mathbf{w}_0
2. For $t = 1, 2, 3, \dots$ do
 - 2.1 Compute the gradient

$$\nabla_{\mathbf{w}} E = - \sum_{n=1}^N \frac{y_n \mathbf{x}_n}{1 + e^{y_n \mathbf{w}_t^T \mathbf{x}_n}}$$

- 2.2 Update the weights

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta \nabla_{\mathbf{w}} E$$

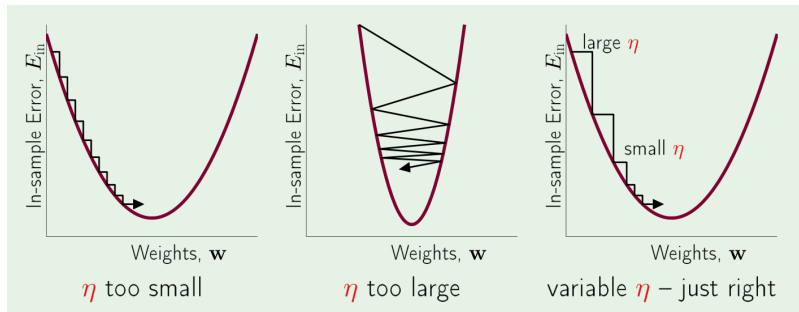
where η is a learning rate (*hyper-parameter*)

- 2.3 Iterate the next step until \mathbf{w} is not change
3. Return the final weights \mathbf{w}



Learning Rate

- How η affects the algorithm?



Programming Example



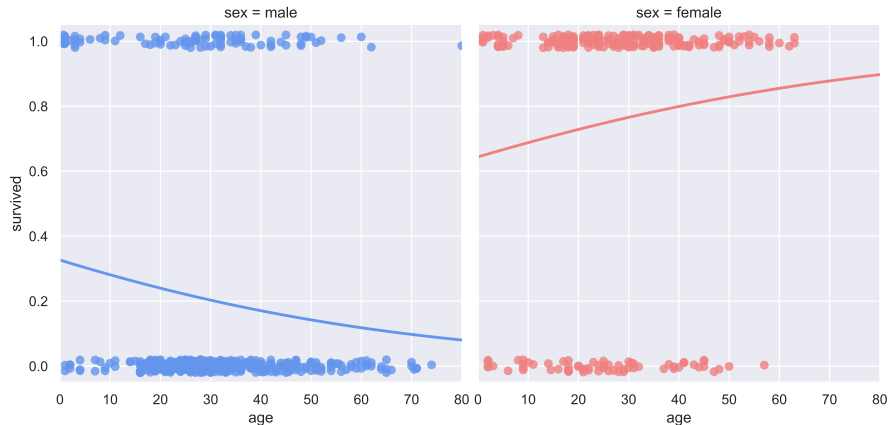
```
import matplotlib.pyplot as plt
import seaborn as sns
sns.set(style="darkgrid")

# Load the example titanic dataset
df = sns.load_dataset("titanic")

# Make a custom palette with gendered colors
pal = dict(male="#6495ED", female="#F08080")

# Show the survival probability as a function of age and sex
g = sns.lmplot(x="age", y="survived", col="sex", hue="sex", data=df,
               palette=pal, y_jitter=.02, logistic=True, ci=None)
g.set(xlim=(0, 80), ylim=(-.05, 1.05))
plt.show()
```

Programming Example (cont.)



References



Goodfellow, I., Bengio, Y., and Courville, A. (2016).

Deep learning.

MIT press.



Lê, B. and Tô, V. (2014).

Cở sở trí tuệ nhân tạo.

Nhà xuất bản Khoa học và Kỹ thuật.



Nguyen, T. (2018).

Artificial intelligence slides.

Technical report, HCMC University of Sciences.



Russell, S. and Norvig, P. (2016).

Artificial intelligence: a modern approach.

Pearson Education Limited.