Pattern Information Processing: ¹ Linear Models and Least-Squares

> Masashi Sugiyama (Department of Computer Science)

Contact: W8E-505 <u>sugi@cs.titech.ac.jp</u> http://sugiyama-www.cs.titech.ac.jp/~sugi/

Supervised Learning: Learning from Examples

There exists a function y = f(x).
We do not know f(x), but we are given its samples {(x_i, y_i)}ⁿ_{i=1}.
The goal of supervised learning is to obtain an approximation f(x) to f(x).

from $\{(x_i, y_i)\}_{i=1}^n$.

Regression: Real Outputs

For real output $y \in \mathbb{R}$, the supervised learning problem is called regression.



Classification: Categorical Outputs

For categorical output $y \in \{1, ..., m\}$, the problem is called classification.

$$y = 1 \xrightarrow{\times \times} \\ y = 1 \xrightarrow{\times} \\ y =$$

For the moment, let us focus on regression problems.

Notations

- $\square \mathcal{D} \subset \mathbb{R}^d$: Input domain
- $f(\boldsymbol{x})$:Learning target function ($\mathcal{D} \to \mathbb{R}$)
- lacksquare $x_i \in \mathcal{D}$: Training input point
- $y_i \in \mathbb{R}$: Training output value
- $\{(\boldsymbol{x}_i, y_i)\}_{i=1}^n$: Training examples
- $\hat{f} \in \mathcal{M}$:Learned function
- \blacksquare *M* :Model (a set of functions)

Today's Plan

Model:

- Linear model
- Kernel model
- Learning method:
 - Least-squares learning

Linear/Non-Linear Models

7

Model is a set of functions from which learning result functions are searched.

We use a parameterized family of functions

$$\{f_{\boldsymbol{\alpha}}(\boldsymbol{x}) \mid \boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_b)^\top\}$$

Linear model: f_{\alpha}(x) is linear with respect to \alpha
 (Note: not necessarily linear with respect to x)
 Non-linear model: Otherwise

Linear Model

$$f_{\boldsymbol{\alpha}}(\boldsymbol{x}) = \sum_{i=1}^{b} \alpha_i \varphi_i(\boldsymbol{x})$$

 $\{\varphi_i(\boldsymbol{x})\}_{i=1}^b: Linearly independent basis functions \\ \mbox{For example, when } d=1: \label{eq:product}$

Polynomial bases

$$1, x, x^2, \dots, x^{b-1}$$

Sinusoidal bases

$$1, \sin x, \cos x, \dots, \sin kx, \cos kx$$

$$b = 2k + 1$$

Multi-Dimensional Linear Model ⁹

For multidimensional input (d > 1), a product model could be used.

$$f_{\alpha}(\boldsymbol{x}) = \sum_{i_1=1}^{c} \cdots \sum_{i_d=1}^{c} \alpha_{i_1,\dots,i_d} \varphi_{i_1}(x^{(1)}) \cdots \varphi_{i_d}(x^{(d)})$$

$$\boldsymbol{x} = (x^{(1)}, \dots, x^{(d)})^{\top}$$

The number of parameters is b = c^d, growing exponentially with respect to d.
Infeasible for large d !

Additive Model

For large d, we have to reduce the number of parameters.

Additive model: $\boldsymbol{x} = (x^{(1)}, \dots, x^{(d)})^{\top}$ $f_{\boldsymbol{\alpha}}(\boldsymbol{x}) = \sum_{j=1}^{d} \sum_{i=1}^{c} \alpha_{i,j} \varphi_i(x^{(j)})$

The number of parameters is only b = cd.

However, additive model is too simple so its representation capability may not be rich enough in some application.

Kernel Model

Linear model:

 $\{\varphi_i(\boldsymbol{x})\}_{i=1}^b$ do not depend on $\{(\boldsymbol{x}_i, y_i)\}_{i=1}^n$

Kernel model:

$$f_{\boldsymbol{\alpha}}(\boldsymbol{x}) = \sum_{i=1}^{n} \alpha_i K(\boldsymbol{x}, \boldsymbol{x}_i)$$

• Example: Gaussian kernel

$$K(\boldsymbol{x}, \boldsymbol{x}') = \exp\left(-\frac{\|\boldsymbol{x} - \boldsymbol{x}'\|^2}{2h^2}\right)$$

 $h(\geq 0)$: Bandwidth

Gaussian Kernel Model

12

$$f_{\boldsymbol{\alpha}}(\boldsymbol{x}) = \sum_{i=1}^{n} \alpha_i \exp\left(-\frac{\|\boldsymbol{x} - \boldsymbol{x}_i\|^2}{2h^2}\right)$$

Put kernel functions at training input points.



Gaussian Kernel Model (cont.)¹³

When training inputs are unevenly distributed, Gaussian kernel model automatically focuses on the region where training inputs exist.



Kernel Model (cont.)

$$f_{\boldsymbol{\alpha}}(\boldsymbol{x}) = \sum_{i=1}^{n} \alpha_i K(\boldsymbol{x}, \boldsymbol{x}_i)$$

- The number of parameters is n, which is independent of the input dimensionality d.
- Although kernel model is linear w.r.t. α , the number of parameters grows as the number of training samples increases.
- Mathematical treatment could be different from ordinary linear models (called a "nonparametric model" in statistics).

Summary of Linear Models

- Linear model (product): High flexibility, high complexity
- Linear model (additive): Low flexibility, low complexity
- Kernel model:
 - Moderate flexibility, moderate complexity
- Good model depends on applications.
- Later, we discuss how to choose an appropriate model ("model selection").

Learning Methods

Linear learning method:

Parameter vector $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_b)^\top$ is estimated linearly with respect to

$$\boldsymbol{y} = (y_1, \dots, y_n)^{\top}$$

Non-linear learning methods: Otherwise

Linear Learning for ¹⁷ Linear Models / Kernel Models

$$f_{\boldsymbol{lpha}}(\boldsymbol{x}) = \sum_{i=1}^{b} \alpha_i \varphi_i(\boldsymbol{x})$$

In linear learning methods, a learned parameter vector is given by

$$\hat{\boldsymbol{\alpha}} = \boldsymbol{L} \boldsymbol{y}$$
 L : Learning matrix

Least-Squares Learning

Learn α such that the squared error at training input points is minimized:

$$\hat{oldsymbol{lpha}}_{LS} = \operatorname*{argmin}_{oldsymbol{lpha} \in \mathbb{R}^b} J_{LS}(oldsymbol{lpha})$$

$$egin{aligned} J_{LS}(oldsymbollpha) &= \sum_{i=1}^n \left(f_{oldsymbollpha}(oldsymbol x_i) - y_i
ight)^2 \ &= \|oldsymbol Xoldsymbollpha - oldsymbol y\|^2 \end{aligned}$$

 $\boldsymbol{X}_{i,j} = \varphi_j(\boldsymbol{x}_i)$:Design matrix $(n \times b)$

How to Obtain A Solution

Extreme-value condition:

$$abla J_{LS}(\hat{oldsymbol{lpha}}_{LS}) = 2 oldsymbol{X}^{ op} (oldsymbol{X} \hat{oldsymbol{lpha}}_{LS} - oldsymbol{y}) = oldsymbol{0}$$

$$\hat{oldsymbol{lpha}}_{LS} = (oldsymbol{X} \ ^{-1}oldsymbol{X} \ ^{-1}oldsymbol{X} \ ^{-1}oldsymbol{Y}$$

(We assume $(\mathbf{X}^{\top}\mathbf{X})^{-1}$ exists.)

Therefore, LS is linear learning.

$$\hat{oldsymbol{lpha}}_{LS} = oldsymbol{L}_{LS}oldsymbol{y}$$
 $oldsymbol{L}_{LS} = (oldsymbol{X}^ opoldsymbol{X})^{-1}oldsymbol{X}^ op$

If you are not familiar with vector-derivatives, see e.g, "Matrix Cookbook" (http://matrixcookbook.com)

Example of LS

$$f_{\boldsymbol{lpha}}(\boldsymbol{x}) = \sum_{i=1}^{b} lpha_i \varphi_i(\boldsymbol{x})$$

Trigonometric polynomial model:

 $1, \sin x, \cos x, \dots, \sin 15x, \cos 15x \quad (b = 31)$



Homework

$$f_{oldsymbol{lpha}}(oldsymbol{x}) = \sum_{i=1}^b lpha_i arphi_i(oldsymbol{x})$$

1. Prove that the LS solution in kernel models is given by

$$\hat{oldsymbol{x}}_{LS} = oldsymbol{L}_{LS} oldsymbol{y}$$
 $oldsymbol{L}_{LS} = oldsymbol{K}^{-1}$ $oldsymbol{K}_{i,j} = K(oldsymbol{x}_i,oldsymbol{x}_j)$ (Kernel matrix)

Homework (cont.)

22

2. For your own toy 1-dimensional data, perform simulations using

 Gaussian kernel models and least-squares learning

and analyze the results when, e.g.,

 target functions, number of samples, noise level, and width of Gaussian kernel

are changed.

Tips: If matrix K is unstable to invert, you may add a small positive constant to diagonals.

Deadline: Next class