

# intro to machine learning

zhangyifeng

October 21, 2018

## Contents

<b>1</b>	<b>Intro</b>	<b>1</b>
1.1	Course Requirements and Grading . . . . .	1
1.2	Course . . . . .	2
<b>2</b>	<b>Logistic Regression</b>	<b>3</b>

## 1 Intro

### 1.1 Course Requirements and Grading

#### Lab(30%)

- Python
- Synthetic data
- 2 deliverables, distributed over moodle

#### Theory exercises(0/20)

close to the end(early December)

#### Final exam(70%)

- Theory questions(judgement-oriented)
- Simulate running algorithms by hand

#### Meeting hours

- Office: 104B, 68-72 Gower street
- Meeting hours: Tuesday, 14:00-15:00

#### Prerequisites

Linear Algebra; Calculus; Probability; Programming

## 1.2 Course

### Machine Learning

data -> model -> prediction

#### Least squares model

- least squares solution for linear regression
- Least squares solution for generalized linear regression
- Least squares solution for ridge regression

#### notation

$D$ : problem dimension, e.g. 1D, 2D( can visualize)

$N$ : training set size

Training set: input-output pairs  $S = \{\mathbf{x}_i, y_i\}, i = 1, \dots, N$  where  $\mathbf{x}_i = \{x_{i1}, \dots, x_{iD}\}^T \in \mathbb{R}^D, y_i \in \mathbb{R}$ , generally can be  $\mathbf{x}$

$\mathbf{w}$ : weight,  $\mathbf{w} = \{w_1, \dots, w_D\}^T \in \mathbb{R}^D$

$\epsilon_i$ : noise

$X = \{x_1, x_2, \dots, x_N\}^T = \{x_1^T; x_2^T; \dots; x_N^T\}$

Remark: "," represent column vector

$\mathbf{y} = \{y_1, \dots, y_N\}^T$

$\boldsymbol{\epsilon} = \{\epsilon_1, \dots, \epsilon_N\}^T$ : residuals vectors

#### linear regression model

$$\mathbf{y} = X\mathbf{w} + \boldsymbol{\epsilon} \quad \text{or} \quad \mathbf{y}^T = \mathbf{w}^T X^T + \boldsymbol{\epsilon}^T$$

that is  $y_i = \mathbf{x}_i^T \mathbf{w} + \epsilon_i$ , or  $y_i = \mathbf{w}^T \mathbf{x}_i + \epsilon_i, i = 1, \dots, N$

Loss function:  $L(\mathbf{w}) = \sum_{i=1}^N (y_i - \mathbf{w}^T \mathbf{x}_i)^2$

goal:  $\min_w L(\mathbf{w})$  Least squares solution for linear regression:  $\mathbf{w}^* = (X^T X)^{-1} X^T \mathbf{y}$

#### Generalized linear regression model

$\mathbf{x} \rightarrow [\boldsymbol{\phi}(\mathbf{x})] = [\phi_1(\mathbf{x}), \dots, \phi_M(\mathbf{x})]^T$ , where  $\phi_i(\mathbf{x}), i = 1, \dots, M$ ,  $\boldsymbol{\phi}(\mathbf{x})$  can be other form besides  $x_i$  ( if  $x_i$ , and  $M = D$ , it is just the linear regression model)

If  $D = 1$ , and  $\phi_i(x) = x^{i-1}$ , then it is **k-th degree polynomial fitting**

If the highest order of  $\phi_i(\mathbf{x})$  is 2, then it is **second-order polynomials fitting**

set  $\Phi = [\boldsymbol{\phi}(\mathbf{x}_1)^T; \boldsymbol{\phi}(\mathbf{x}_2)^T; \dots; \boldsymbol{\phi}(\mathbf{x}_N)^T]$

then the model is:

$$\mathbf{y} = \Phi \mathbf{w} + \boldsymbol{\epsilon}$$

Least squares solution for generalized linear regression:  $\mathbf{w}^* = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y}$

### approximations

If  $N > D$  (e.g. 30 points, 2 dimensions): overdetermined system

If  $N < D$  (e.g. 30 points, 3000 dimensions): underdetermined system (overfitting)

### How to control complexity (Regularized linear regression)

1. use vector norm (L2, L1, Lp norm) to measure residual vector

Remark: different norm represent different regularized linear regression, here we use L2 norm

2. rewrite loss function:  $L(\mathbf{w}) = \|\mathbf{y} - X\mathbf{w}\|^2 + \lambda\|\mathbf{w}\|^2$

this is **Ridge regression**, a.k.a, L2-regularized linear regression

Remark:  $\lambda$  is "hyperparameter", select  $\lambda$  with cross-validation (use cross-validation for diff values of  $\lambda$  - pick value minimizes cross-validation error)

Cross-validation: least glorious, most effective of all methods (teacher said)

3. Least squares solution for ridge regression:  $\mathbf{w}^* = (X^T X + \lambda \mathbf{I})^{-1} X^T \mathbf{y}$

## 2 Logistic Regression

### Machine Learning variants-Supervised learning

- Classification
- Regression

### Gaussian ( or Normal ) distribution

one-dimensional case

- Mean  $\mu$
- Variance  $\sigma^2$
- Pdf:  $\mathcal{N}(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$

Multi-dimensional case

- Mean  $\boldsymbol{\mu}$
- Covariance  $\boldsymbol{\Sigma}$
- Pdf:  $\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right\}$

Remark:  $|\boldsymbol{\Sigma}|$  represent matrix norm, e.g. Frobenius norm of  $\boldsymbol{\Sigma}$

### Parameter estimation

Given:  $X = \{x_1, x_2, \dots, x_N\}$ , parametric form of distribution, parameters  $\theta$

Learning goal: estimate  $\theta$

Likelihood of  $\theta$ :  $L(\theta) = p(X; \theta) = \prod_{n=1}^N p(x_n; \theta)$

$$\begin{aligned} \text{Log-likelihood: } \ln L(\theta) &= \sum_{n=1}^N \ln p(x_n; \theta) \\ \max \ln L(\theta) &\Rightarrow \frac{\partial \ln L(\theta)}{\partial \theta} = 0 \Rightarrow \hat{\theta} \end{aligned}$$