

# Introduction

2017-12-06 14:23

## An example

- Hand writing recognition: machine learning is a kind of method of
- Curve Fitting:

- $$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2$$

- Root-mean-square error (RMS-Error)

- $$E_{\text{RMS}} = \sqrt{2E(\mathbf{w}^*)/N}$$

- To prevent over fitting

- More data
- By regularization (正则化)

- $$\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

- Suggest: taking the available data and partitioning it into a training set, used to determine the coefficients  $\mathbf{w}$ , and a separate validation set, also called a hold-out set, used to optimize the model complexity (either  $M$  or  $\lambda$ ). However it is a kind of wasting training data
- 这一部分作为导入很贴切，基本说明了是在用什么思路解决什么问题、解决问题中可能遇到的瓶颈以及常用的解决策略

## Probability Theory

### 0.基本概念

- Sum rule of probability (加法法则 边缘概率)

- $$p(X = x_i) = \sum_{j=1}^L p(X = x_i, Y = y_j) \quad i = \frac{c_i}{N}.$$

$$\therefore p(X, Y) = p(X|Y)p(Y) = p(Y|X)p(X) = p(Y, X)$$

- conditional probability (条件概率)

- $$p(Y = y_j | X = x_i) = \frac{n_{ij}}{c_i}.$$

$$\therefore p(Y|X) = p(X|Y)p(Y)/p(X)$$

- Joint probability (联合概率)

- $$p(X = x_i, Y = y_j) = \frac{n_{ij}}{N} = \frac{n_{ij}}{c_i} \cdot \frac{c_i}{N} = p(Y = y_j | X = x_i) p(X = x_i)$$

- Productive rule (乘法法则)

- $$p(X, Y) = p(Y|X)p(X).$$

- Bayes' theorem

- $$p(Y|X) = \frac{p(X|Y)p(Y)}{p(X)}$$

- Prior probability (先验概率): 在随机变量的值确定前获得的概率。如从蓝色盒子中抽取的概率是40%
- Posterior probability (后验概率): 在随机变量的值确定后推算出的概率。如已知抽取到的是苹果，推算出之前是从蓝盒子中抽取的概率。

### 1.概率密度 (当随机变量的取值是连续型时)

- The probability that  $x$  will lie in an interval  $(a, b)$  is then given by

- $$p(x \in (a, b)) = \int_a^b p(x) dx.$$

- 原本要求 $p(y)$ 需要用 $p_y$ ，在知道随机变量 $x$ 和 $y$ 之间的关系后，能将原问题转化成用 $p_x$ 来求

Under a nonlinear change of variable, a probability density transforms differently from a simple function, due to the Jacobian factor. For instance, if we consider a change of variables  $x = g(y)$ , then a function  $f(x)$  becomes  $\tilde{f}(y) = f(g(y))$ . Now consider a probability density  $p_x(x)$  that corresponds to a density  $p_y(y)$  with respect to the new variable  $y$ , where the suffices denote the fact that  $p_x(x)$  and  $p_y(y)$  are different densities. Observations falling in the range  $(x, x + \delta x)$  will, for small values of  $\delta x$ , be transformed into the range  $(y, y + \delta y)$  where  $p_x(x)\delta x \simeq p_y(y)\delta y$ , and hence

? 这里的约等于关系是怎么来的

- are different densities. Observations falling in the range  $(x, x + \delta x)$  will, for small values of  $\delta x$ , be transformed into the range  $(y, y + \delta y)$  where  $p_x(x)\delta x \simeq p_y(y)\delta y$ , and hence

$$p_y(y) = p_x(x) \left| \frac{dx}{dy} \right| = p_x(g(y)) |g'(y)|. \quad (1.27)$$

? 这里的约等于关系是怎么来的

- 如果随机变量的取值是离散的，那么将  $p(x)$  称为概率质量函数 probability mass function
- 对于连续型变量，其概率的加法公式和乘法公式为

$$p(x) = \int p(x, y) dy$$

$$p(x, y) = p(y|x)p(x).$$

## 2. Expectations and covariance 期望和协方差

- 在知道一个函数自变量取值的概率分布情况下，可以求该函数的期望

$$\mathbb{E}[f] = \sum_x p(x)f(x)$$

$$\mathbb{E}[f] = \int p(x)f(x) dx.$$

- 方差

$$\text{var}[f] = \mathbb{E}[(f(x) - \mathbb{E}[f(x)])^2]$$

$$\mathbb{E}[(f(x) - \mathbb{E}[f(x)])^2] = \mathbb{E}[f(x)^2 - 2f(x)\mathbb{E}[f(x)] + \mathbb{E}[f(x)]^2]$$

$$= \mathbb{E}[f(x)^2] - 2\mathbb{E}[f(x)]^2 + \mathbb{E}[f(x)]^2 = \mathbb{E}[f(x)^2] - \mathbb{E}[f(x)]^2$$

## 3. Bayesian probabilities

- 先验概率 (incorporating with observed data) 后验概率

$$p(\mathbf{w}|\mathcal{D}) = \frac{p(\mathcal{D}|\mathbf{w})p(\mathbf{w})}{p(\mathcal{D})} \quad (1.43)$$

Likelihood function(似然函数): 描述了数据集中对于不同的  $\mathbf{w}$  观测的情况。



integrating both sides of (1.43) with respect to  $\mathbf{w}$ , we can express the denominator in Bayes' theorem in terms of the prior distribution and the likelihood function

$$p(\mathcal{D}) = \int p(\mathcal{D}|\mathbf{w})p(\mathbf{w}) d\mathbf{w}. \quad (1.45)$$

Note: 这里的似然函数并非观测在  $\mathbf{w}$  上的概率分布，所以似然函数累加求和也并一定等于1。

- Maximum likelihood(极大似然法): in which  $\mathbf{w}$  is set to the value that maximizes the likelihood function  $p(\mathcal{D}|\mathbf{w})$

## 4. Gaussian distribution (高斯分布 正态分布)

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left\{-\frac{1}{2\sigma^2}(x - \mu)^2\right\} \quad (1.46)$$

$\sigma$  标准差

- $\beta = 1/\sigma^2$ , is called the precision

$$\mathbb{E}[x] = \int_{-\infty}^{\infty} \mathcal{N}(x|\mu, \sigma^2) x dx = \mu. \quad (1.49)$$

$$\mathbb{E}[x^2] = \int_{-\infty}^{\infty} \mathcal{N}(x|\mu, \sigma^2) x^2 dx = \mu^2 + \sigma^2. \quad (1.50)$$

- 如果随机变量  $\mathbf{x}$  是  $D$  维的:

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

$$\begin{aligned} \mathbb{E}[(f(x) - \mathbb{E}[f(x)])^2] &= \mathbb{E}[f(x)^2 - 2f(x)\mathbb{E}[f(x)] + \mathbb{E}[f(x)]^2] \\ &= \mathbb{E}[f(x)^2] - 2\mathbb{E}[f(x)]^2 + \mathbb{E}[f(x)]^2 \\ &= \mathbb{E}[f(x)^2] - \mathbb{E}[f(x)]^2 \end{aligned} \quad (1.52)$$

- 最大化似然函数，已知  $\mathbf{x}$  服从正态分布

$$p(\mathbf{x}|\mu, \sigma^2) = \prod_{n=1}^N \mathcal{N}(x_n|\mu, \sigma^2). \quad (1.53)$$

$$\ln p(\mathbf{x}|\mu, \sigma^2) = -\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 - \frac{N}{2} \ln \sigma^2 - \frac{N}{2} \ln(2\pi). \quad (1.54)$$

直接求偏导，令偏导数为零即可得

$$\mu_{\text{ML}} = \frac{1}{N} \sum_{n=1}^N x_n \quad (1.55)$$

$$\sigma_{\text{ML}}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{\text{ML}})^2$$

$$\mathbb{E}[\mu_{\text{ML}}] = \mu \quad (1.57)$$

ML: max likelihood



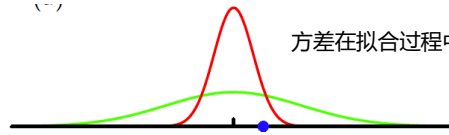
方差在拟合过程中缩小了

$$\sigma_{\text{ML}}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{\text{ML}})^2$$

$$\mathbb{E}[\mu_{\text{ML}}] = \mu \quad (1.57)$$

$$\mathbb{E}[\sigma_{\text{ML}}^2] = \left( \frac{N-1}{N} \right) \sigma^2 \quad (1.58)$$

这里并不straightforward



方差在拟合过程中缩小了

## 5. Curve fitting re-visited (重新审视曲线拟合) P47