Mathematic Foundation and Basic Concept

Hao Dong

Peking University

Recap: Prerequisites

- Basic knowledge of probabilities
 - Bayes rule, chain rule, probability distribution ...
- Basic knowledge of information theory
 - Self-information, Shannon entropy, differential entropy ...
 - Kullback-Leibler (KL) divergence
- Basic knowledge of machine learning/deep learning
 - "Machine Learning", "Pattern Recognition and Machine Learning"
 - "Computer Vision", "Natural Language Processing" ...
- Basic programming language
 - Python

Mathematic Foundation and Basic Concept

- Example: Regression Polynomial Curve Fitting
- Probability Theory
- Decision Theory
- Information Theory

• Example: Regression - Polynomial Curve Fitting

- Probability Theory
- Decision Theory
- Information Theory

• Problem Definition



• A training dataset of N = 10 points

1

Predict the target value *t* given the input *x* ٠

$$\mathbf{x} \equiv (x_1, \dots, x_N)^{\mathrm{T}}$$
 A row vector with N elements $\mathbf{t} \equiv (t_1, \dots, t_N)^{\mathrm{T}}$

• Simple Model: Polynomial Function



M order polynomial:

$$y(x, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \ldots + w_M x^M = \sum_{j=0}^M w_j x^j$$

 $\mathbf{w} = (w_1, \ldots, w_M)^{\mathrm{T}}$ Weights: a row vector with M elements

The appropriate value \hat{t} :

 $\hat{t} = \mathbf{x}^{\mathrm{T}} \mathbf{w}$

• Error Function



• Model Capacity and Overfitting



M order polynomial

Root mean square error:

$$E_{\rm RMS} = \sqrt{2E(\mathbf{w}^{\star})/N}$$
 Averaged/mean error
$$E(\mathbf{w}^{\star}) = 0$$
 Optimised weights



• Solution1 : More Training Data Point



More data points = Better generalization

• Solution2 : Weight Regularisation



• Solution2 : Weight Regularisation

$$\widetilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^{N} \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$
$$\|\mathbf{w}\|^2 \equiv \mathbf{w}^{\mathrm{T}} \mathbf{w} = w_0^2 + w_1^2 + \ldots + w_M^2$$



• Solution2 : Weight Regularisation



- Example: Regression Polynomial Curve Fitting
- Probability Theory
- Decision Theory
- Information Theory

- The probability of an event is the fraction of times that event occurs out of the total number of trials
- The probability must lie in the interval [0, 1]
- A box contains red and blue balls, we randomly pick 100 balls from a box, and 90 balls are red

$$p(ball = red) = \frac{90}{100} = 0.9$$





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

Given ..

$$p(X = x_i) = \frac{c_i}{N}$$

So.. We have the sum rule of probability

$$p(X = x_i) = \sum_{j=1}^{L} p(X = x_i, Y = y_j)$$

Conditional probability: $Y = y_i$ given $X = x_i$



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

As we know:

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

We can have the product rule of probability

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

Variable Value

• The Rules of Probability

$$\label{eq:posterior} \begin{array}{ll} \mbox{sum rule} & p(X) = \sum_Y p(X,Y) \\ \mbox{product rule} & p(X,Y) = p(Y|X)p(X) \end{array}$$

• Bayes' theorem

Given the **product rule**:

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

 $p(Y,X) = p(X|Y)p(Y)$
symmetry property $p(X,Y) = p(Y,X)$

We can have the **Bayes' theorem:**

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

Using the sum rule ($p(X) = \sum_{Y} p(X, Y)$), we can have:

$$p(X) = \sum_{Y} p(X|Y)p(Y)$$

• Bayes' theorem

The prior probability: it is the probability available before we observe the X The likelihood : It expresses how probable the 先验 observed data set is for different settings of the parameter y 似然 $p(Y|X) = \frac{p(X|Y)p(Y)}{p(X)}$ $p(X) = \frac{p(X|Y)p(Y)}{p(X)}$ $p(X) = \int_{Y} p(X|Y)p(Y)dY$ The posterior probability: it is the probability obtained after we have observed X 后验

posterior \propto likelihood \times prior

• The Rules of Probability

sum rule
$$p(X) = \sum_{Y} p(X, Y)$$
product rule $p(X, Y) = p(Y|X)p(X)$ Bayes' theorem $p(Y|X) = \frac{p(X|Y)p(Y)}{p(X)}$ where $p(X) = \sum_{Y} p(X|Y)p(Y)$

• Independent Events

If X and Y are independent

Bayes' theorem
$$p(Y|X) = \frac{p(X|Y)p(Y)}{p(X)}$$
 $p(Y|X) = p(Y)$
product rule $p(X,Y) = p(Y|X)p(X)$ $p(X,Y) = p(X)p(Y)$

Probability Densities •



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

$$p(x \in (a,b)) = \int_{a}^{b} p(x) \, \mathrm{d}x$$

$$p(x) \ge 0$$

$$\int_{-\infty}^{\infty} p(x) \, \mathrm{d}x = 1$$

• Probability Densities



cumulative distribution function (CDF)

$$P(z) = \int_{-\infty}^{z} p(x) \, \mathrm{d}x$$

which satisfies P'(x) = p(x)

• Expectations and Covariances

The average value of some function f(x) under a probability distribution p(x) is called the **expectation** of f(x) and will be denoted by E[f]:

For a continuous distribution:

$$\mathbb{E}[f] = \int p(x)f(x) \,\mathrm{d}x$$

For a discrete distribution:

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

$$\mathbb{E}[f] \simeq \frac{1}{N} \sum_{n=1}^{N} f(x_n)$$

It becomes exact in the limit $N \to \infty$

• Expectations and Covariances

In the case of **continuous variables**, expectations are expressed in terms of an integration with respect to the corresponding probability density

$$\mathbb{E}[f] = \int p(x)f(x) \, \mathrm{d}x$$

The conditional expectation with respect to a conditional distribution

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

• Expectations and Covariances

The variance of f(x):

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

For two random variables x and y, the **covariance** is:

$$cov[x, y] = \mathbb{E}_{x,y} \left[\left\{ x - \mathbb{E}[x] \right\} \left\{ y - \mathbb{E}[y] \right\} \right] \\ = \mathbb{E}_{x,y} [xy] - \mathbb{E}[x] \mathbb{E}[y]$$

For two vectors of random variables **x** and **y**, the **covariance matrix**:

The **variance** of the variable *x*:

$$\operatorname{var}[x] = \mathbb{E}[x^2] - \mathbb{E}[x]^2$$

$$\begin{aligned} \operatorname{cov}[\mathbf{x}, \mathbf{y}] &= & \mathbb{E}_{\mathbf{x}, \mathbf{y}} \left[\{ \mathbf{x} - \mathbb{E}[\mathbf{x}] \} \{ \mathbf{y}^{\mathrm{T}} - \mathbb{E}[\mathbf{y}^{\mathrm{T}}] \} \right] \\ &= & \mathbb{E}_{\mathbf{x}, \mathbf{y}}[\mathbf{x}\mathbf{y}^{\mathrm{T}}] - \mathbb{E}[\mathbf{x}]\mathbb{E}[\mathbf{y}^{\mathrm{T}}]. \end{aligned}$$

The covariance of the components of a vector **x** with each other:

$$\operatorname{cov}[\mathbf{x}] \equiv \operatorname{cov}[\mathbf{x}, \mathbf{x}]$$

- Bayesian Probabilities
 - Frequentist: Probability in terms of frequencies of random, repeatable events.
 - Do not work if the event is not repeatable, e.g., the probability of the Arctic ice cap disappear.
 - **Bayesian**: Degree of belief.

- Frequentist Probability
 - Maximum Likelihood Estimation, MLE

A widely used frequentist estimator is maximum likelihood, in which θ is set to the value that maximises the likelihood function $p(X | \theta)$.

- Given $x \sim iid \sim p(x|\theta)$, $p(X|\theta) = \prod_{i=1}^{N} p(x_i|\theta)$
- $\theta_{\text{MLE}} = \arg \max_{\theta} \log p(X|\theta) = \arg \max_{\theta} \sum_{i=1}^{N} \log p(x_i|\theta)$

- Bayesian Probabilities
 - θ : random variables, $\theta \sim p(\theta) \leftarrow$ prior
 - MAP: Maximum A Posterior

 $p(\theta|X) = \frac{p(X|\theta)p(\theta)}{p(X)} \propto p(X|\theta)p(\theta)$ posterior
normalisation constant

•
$$\theta_{MAP} = \arg \max_{\theta} p(\theta|X) = \arg \max_{\theta} p(X|\theta)p(\theta)$$

Bayes can evaluate the **uncertainty** in θ after we have observed X in the form of the posterior probability $p(\theta|X)$.

• The Gaussian Distribution



• The Gaussian Distribution



$$\mathbb{E}[x] = \int_{-\infty}^{\infty} \mathcal{N}(x|\mu, \sigma^2) x \, \mathrm{d}x = \mu$$

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

$$\operatorname{var}[x] = \mathbb{E}[x^2] - \mathbb{E}[x]^2 = \sigma^2$$

• The Gaussian Distribution

D-dimensional vector **x** of continuous variables

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

 μ is called the mean Σ is called the covariance $|\Sigma|$ denotes the determinant of Σ

The Gaussian Distribution



Now suppose that we have a dataset of observations $\mathbf{x} = (x_1, ..., x_N)^T$, representing N observations of the scalar variable x.

How to determine the mean and variance according to the dataset??

All data points are sampled independently from the same distribution, they are **independent and identically distributed (i.i.d)**.

For i.i.d variables, we can have the joint probability as follow

$$p(\mathbf{x}) = p(x_1) p(x_2) p(x_3) \dots p(x_N)$$

$$p(\mathbf{x}|\boldsymbol{\mu}, \sigma^2) = \prod_{n=1}^N \mathcal{N}\left(x_n | \boldsymbol{\mu}, \sigma^2\right)$$

The Gaussian Distribution



Maximum Likelihood Estimation (MLE)

To find the "best" mean and variance, we can maximise the probability of the parameters given the data

$$p(\mathbf{X}|\boldsymbol{\mu}, \sigma^2) = \prod_{n=1}^{N} \mathcal{N}\left(x_n | \boldsymbol{\mu}, \sigma^2\right)$$

In practice, it is more convenient to maximise the log of the likelihood function. The logarithm is a monotonically increasing function, maximisation of the log of a function is equivalent to maximisation of the function itself

$$\ln p\left(\mathbf{x}|\mu,\sigma^{2}\right) = -\frac{1}{2\sigma^{2}}\sum_{n=1}^{N}(x_{n}-\mu)^{2} - \frac{N}{2}\ln\sigma^{2} - \frac{N}{2}\ln(2\pi)$$

So.. the maximum likelihood solution:

$$\mu_{\rm ML} = \frac{1}{N} \sum_{n=1}^{N} x_n$$

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

- Example: Regression Polynomial Curve Fitting
- Probability Theory
- Decision Theory
- Information Theory

- An example
- Basic intuition
- Minimising the misclassification rate
- Minimising the expected loss
- The reject option
- Loss functions for regression

• An example

A medical diagnosis problem in which we have taken an X-ray image of a patient, and we wish to determine whether the patient has cancer or not.

We are interested in **the probabilities of the two classes** given the **image**, which are given by p(Ck|x). Using Bayes' theorem, these probabilities can be expressed in the form :

$$p(\mathcal{C}_k|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)}{p(\mathbf{x})}.$$

Note that any of the quantities appearing in Bayes' theorem can be obtained from the joint distribution p(x, Ck) by either marginalizing or conditioning with respect to the appropriate variables.

We can now interpret p(Ck) as the prior probability for the class Ck, and p(Ck|x) as the corresponding posterior probability. Thus p(C1) represents the probability that a person has cancer, before we take the X-ray measurement.

Similarly, p(C1|x) is the corresponding probability, revised using Bayes' theorem in light of the information contained in the X-ray.

- Basic intuition
 - If our aim is to minimise the chance of assigning x to the wrong class, then intuitively we would choose the class having the higher posterior probability.
 - We now show that this intuition is correct, and we also discuss more general criteria for making decisions.



• Minimising the expected loss

Loss function, also called a cost function, which is a single, overall measure of loss incurred in taking any of the available decisions or actions.

Our goal is then to minimise the total loss incurred.

We seek to minimise the average loss, where the average is computed with respect to this distribution, which is given by

$$\mathbb{E}[L] = \sum_{k} \sum_{j} \int_{\mathcal{R}_{j}} L_{kj} p(\mathbf{x}, \mathcal{C}_{k}) \, \mathrm{d}\mathbf{x}$$

• The reject option



- Example: Regression Polynomial Curve Fitting
- Probability Theory
- Decision Theory
- Information Theory

Basic intuition

- Learning that an unlikely event has occurred is more informative than learning that a likely event has occurred.
 - E.g.
 - "the sun rose this morning" : uninformative
 - "there was a solar eclipse this morning": informative.

Basic intuition

- We would like to quantify information in a way that formalises this intuition. Specifically:
 - Likely events should have low information content, and in the extreme case, events that are guaranteed to happen should have no information content whatsoever.
 - Less likely events should have higher information content.
 - Independent events should have additive information.
 - E.g. Finding out that a tossed coin has come up as heads twice should convey twice as much information as finding out that a tossed coin has come up as heads once.

Self-information

• In order to satisfy all three of these properties, we define the self information of an event x to be

$$I(x) = -\log P(x).$$

(use log to mean the natural logarithm, with base e. I(x) is therefore written in units of nats)



Shannon entropy

• Self-information deals only with a single outcome. We can quantify the amount of uncertainty in an entire probability distribution using the Shannon entropy:

$$H(\mathbf{x}) = \mathbb{E}_{\mathbf{x} \sim P}[I(x)] = -\mathbb{E}_{\mathbf{x} \sim P}[\log P(x)].$$

• The Shannon entropy of a distribution is the expected amount of information in an event drawn from that distribution. It gives a lower bound on the number of bits (if the logarithm is base 2, otherwise the units are different) needed on average to encode symbols drawn from a distribution P.

Example for Shannon entropy



KL divergence

• If we have two separate probability distributions P(x) and Q(x) over the same random variable x, we can measure how different these two distributions are using the Kullback-Leibler (KL) divergence:

$$D_{\mathrm{KL}}(P||Q) = \mathbb{E}_{\mathbf{x}\sim P}\left[\log\frac{P(x)}{Q(x)}\right] = \mathbb{E}_{\mathbf{x}\sim P}\left[\log P(x) - \log Q(x)\right].$$

 In the case of discrete variables, it is the extra amount of information needed to send a message containing symbols drawn from probability distribution P, when we use a code that was designed to minimize the length of messages drawn from probability distribution.



KL divergence

- Because the KL divergence is non-negative and measures the difference between two distributions, it is often conceptualised as measuring some sort of distance between these distributions.
 - Not a true distance measure because it is not symmetric.
 - It is nonnegative.

KL divergence



- Which direction of the KL divergence to use?
 - Some applications require an approximation that usually places high probability anywhere that the true distribution places high probability: left one
 - Other applications require an approximation that rarely places high probability anywhere that the true distribution places low probability: right one

$$D_{\mathrm{KL}}(P \| Q) = \mathbb{E}_{\mathbf{x} \sim P} \left[\log \frac{P(x)}{Q(x)} \right] = \mathbb{E}_{\mathbf{x} \sim P} \left[\log P(x) - \log Q(x) \right].$$

Cross-entropy

• A quantity that is closely related to the KL divergence is the crossentropy $H(P,Q) = H(P) + D_{KL}(P||Q)$, which is similar to the KL divergence but lacking the term on the left:

$$H(P,Q) = -\mathbb{E}_{\mathbf{x}\sim P} \log Q(x).$$

 Minimising the cross-entropy with respect to Q is equivalent to minimising the KL divergence, because does not participate in the omitted term. Mathematic Foundation and Basic Concept

- Example: Regression Polynomial Curve Fitting
- Probability Theory
- Decision Theory
- Information Theory

Reference





Free Download

https://github.com/zsdonghao/deep-learningbook/blob/master/All-in-one-pdf 《Deep Learning》 all-in-one.pdf

Thanks