

MA333 Introduction to Big Data Science Mathematical Preliminary

Zhen Zhang

Southern University of Science and Technology

Basics

Three Sources of Uncertainty

1. Inherent randomness in the system being modeled.
(e.g. quantum mechanics, particles are probabilistic.)
2. Incomplete observability
(e.g. Monty Hall problem. 3 dogs, one of which has a bonus behind it.)
3. Incomplete modeling
(e.g. linear regression.)

Random Variables (r.V.)

X : r.V., its values x

$X = x$ happens with a certain probability $p(X = x)$. Range of X may be discrete or continuous.

Discrete variables & Probability mass functions (PMF)

$$P(X = x) = p(x) \text{ or } X \sim p(x), 0 \leq p(x) \leq 1.$$

Joint probability distribution :

$$p(X = x, Y = y) = p(x, y)$$

Properties :

(1) $Dom(P) = Range(X) =$ set of all possible states of X .

(2) $\forall x \in Range(X), 0 \leq p(x) \leq 1.$

(3) $\sum_{x \in Range(X)} p(x) = 1$ (Normalized)

e.g. uniform distribution : $P(X = x_i) = \frac{1}{k}, \quad i = 1, \dots, k.$

Continuous variables & Probability density function (p.d.f)

Properties :

- (1) $Dom(P) = Range(X) =$ set of all possible states of X .
- (2) $\forall x \in Range(X), p(x) \geq 0$. (No need for $p(x) \leq 1$)
- (3) $\int p(x) dx = 1$.

$p(x)$ does not give the probability of a specific state x , but the prob of X landing inside the interval $[x, x + \delta x)$ with an infinitesimal δx . ($prob = p(x)\delta x$)

e.g. uniform distribution function on an interval $[a, b]$:

$$u(x; a, b) = \begin{cases} 0 & x \notin [a, b] \\ \frac{1}{b-a} & x \in [a, b] \end{cases} \quad x \sim U(a, b)$$

Marginal Probability & Conditional Probability

Marginal Probability :

prob over a subset of all variables.

discrete : $p(x) = P(X = x) = \sum_y P(X = x, Y = y) = \sum_y p(x, y)$

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

Conditional Probability :

prob of some event given some other events have happened.

$$P(Y = y | X = x) = \frac{P(X = x, Y = y)}{P(X = x)}$$

(well-defined if $P(X = x) > 0$)

Chain Rules :

$$P(X^{(1)}, \dots, X^{(n)}) = P(X^{(1)})P(X^{(2)} | X^{(1)})P(X^{(3)} | X^{(1)}, X^{(2)}) \dots P(X^{(n)} | X^{(1)}, \dots, X^{(n-1)})$$

Independence & Conditional Independence

X and Y are independent, if

$$\forall x, y, P(X = x, Y = y) = P(X = x)P(Y = y)$$

X, Y are conditionally independent given Z if $\forall x, y, z,$

$$P(X = x, Y = y | Z = z) = P(X = x | Z = z)P(Y = y | Z = z).$$

short-hand notation : $X \perp Y, X \perp Y | Z$.

Expectation Variance and Covariance

$$E_{x \sim p}[f(x)] = \begin{cases} \sum_x p(x)f(x) \\ \int p(x)f(x)dx \end{cases},$$

$$E[\alpha f(x) + \beta g(x)] = \alpha E f(x) + \beta E g(x)$$

$$\text{Var}[f(x)] = E[(f(x) - E f(x))^2], \text{std}[f(x)] = \sqrt{\text{Var}[f(x)]}$$

$$\text{Cov}[f(X), g(Y)] = E[(f(x) - E f(x))(g(Y) - E g(Y))]$$

measures how much two variable are linearly related to each other as well as their scales.

$$\text{Cov}[f(X), g(Y)] = \frac{\text{Cov}[f(X), g(Y)]}{\text{std}[f(X), g(Y)]} \in [-1, 1] \text{ normalized}$$

$$X \perp Y \Rightarrow \text{Cov}[X, Y] = 0 \text{ or } \text{Corr}[X, Y] = 0$$

" \Leftarrow " is only true when $X, Y \sim$ normal distribution

$$\text{Covariance Matrix : } \text{Cov}(\vec{X}) = (\text{Cov}(X_i, Y_j))_{1 \leq i, j \leq n}$$

$$\text{Diagonal } \text{Cov}(X_i, X_i) = \text{Var}(X_i)$$

Common Prob Distributions

- 1) Bernoulli : binary r.v. one parameter $\phi \in [0, 1]$

$$p(X = 1) = \phi, \quad p(X = 0) = 1 - \phi \text{ or}$$

$$p(X = x) = \phi^x (1 - \phi)^{1-x} \quad x \in \{0, 1\},$$

$$E_X[X] = \phi \quad \text{Var}_X[X] = \phi(1 - \phi).$$

- 2) Multinoulli (Categorical)

$$\text{Range}(X) = \{0, 1, \dots, k\}$$

$$p(X = i) = p_i \quad \sum_{i=1}^k p_i = 1 \quad \text{only } k - 1 \text{ parameters.}$$

- 3) Gaussian (Normal)

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

$$E(X) = \mu, \quad \text{Var}[X] = \sigma^2 > 0$$

$$\beta = \sigma^{-2} = \text{precision.}$$

- 4) ..., continue after the importance of Gaussian.

Why are Gaussian important

(1) Central Limit Theorem (many versions)

$\{X_i\}_{i=1}^n$ i.i.d. (independent and identically distributed),

$S_n = \frac{1}{n}(X_1 + \cdots + X_n)$. (Law of Large numbers

$\Rightarrow S_n \rightarrow EX_1 = \mu$ as $n \rightarrow \infty$)

stronger result : $\sqrt{n}(S_n - \mu) \Rightarrow N(0, \sigma^2)$ as $n \rightarrow \infty$), in distribution $\sigma^2 = \text{Var}[X_1]$.

$\lim_{n \rightarrow \infty} P_r[\sqrt{n}(S_n - \mu) \leq z] = \Phi(\frac{z}{\sigma})$, where $\Phi(x)$ is c.d.f of $N(0, 1)$.

Why are Gaussian important

- (2) Information theory : Gaussian encodes the maximum amount of uncertainty

$$H[X] = E_{x \sim p}[-\log p(x)] = - \int p(x) \log p(x) dx.$$

$$\max_p E[X] \text{ s.t.}$$

$$\lambda_1 \quad \int p(x) dx = 1 \Rightarrow p^*(x) = \sqrt{\frac{1}{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right)$$

$$\lambda_2 \quad \int xp(x) dx = \mu$$

$$\lambda_3 \quad \int (x - \mu)^2 p(x) dx = \sigma^2,$$

$$p(x) = \exp\{\lambda_1 - 1 + \lambda_2(x - \mu) + \lambda_3(x - \mu)^2\}, \quad \lambda_2 = 0,$$

$$\lambda_1 : \text{normalization constant}, \quad \lambda_3 = -1/(2\sigma^2).$$

Multivariate normal :

$$N(\vec{x}; \vec{\mu}, \vec{\Sigma}) = \sqrt{\frac{1}{(2\pi)^n \det(\vec{\Sigma})}} \exp\left(-\frac{1}{2}(\vec{x} - \vec{\mu})^T \vec{\Sigma}^{-1}(\vec{x} - \vec{\mu})\right)$$

$\vec{\Sigma}^{-1} = \Omega$ precision matrix.

$$\text{Cov}(\vec{x}) = \vec{\Sigma}, \quad E\vec{x} = \vec{\mu}$$

Common Prob Distributions (Continued)

4) Exponential and Laplace (Doubly exponential)

exponential : $p(x; \lambda) = \lambda \mathbb{1}_{x \geq 0} \exp(-\lambda x)$

Laplace $(x; \mu, r) = \frac{1}{2r} \exp(-\frac{|x-\mu|}{r})$.

5) Dirac and Empirical

Dirac : $p(x) = \delta(x - \mu)$ which puts prob $\frac{1}{m}$ on each of $x^{(i)}$.
(samples), i.e. $P(X = x^{(i)}) = \frac{1}{m}$, $i = 1, \dots, m$.

Common Prob Distributions (Continued)

6) Mixture :

$P(x) = \sum P(c = i)P(x|c = i)$, where $P(c)$ is multimoulli.

Empirical is mixture with Dirac.

Latent variable c , $P(X|c)$ relates the latent variable c to the visible variables X .

e.g. Gaussian Mixtures $P(x|c = i)$ is Gaussian for $\forall i$,
 $\sim N(x; \mu^{(i)}, \Sigma^{(i)})$.

$P(c = i)$ is prior prob. $P(c|x)$ is posterior prob.

Gaussian mixture model is universal approximator of density in the sense that any smooth density can be approximated with any specific non-zero amount of error by a Gaussian mixture model with enough components.

Bayes' Rule

$$P(x|y) = \frac{P(x)P(y|x)}{P(y)} \text{ i.e. } P(x|y)P(y) = P(x, y) = P(x)P(y|x)$$

$$P(y) = \sum_x P(x)P(y|x)$$

$P(x)$: prior ; $P(x|y)$: posterior

Transformations

X, Y two r.v.s. $Y = g(X)$, g is invertible

p.d.f. of X is P_x , p.d.f. of Y is P_y , then $P_k(\vec{y}) = P_x(g^{-1}(\vec{y})) \left| \frac{\partial \vec{x}}{\partial \vec{y}} \right|$

or $P_k(\vec{x}) = P_y(g^{-1}(\vec{x})) \left| \frac{\partial \vec{y}}{\partial \vec{x}} \right| = P_y(g(\vec{x})) |\det J|$.

Let $J = \left(\frac{\partial y_i}{\partial x_j} \right)_{i,j}$.

Information Theory

information measures the amount of uncertainty :

- (1) Likely events have low information
- (2) Less likely events have higher information
- (3) Information events have additive information. (toss a coin twice)

Self-information at event $X = x$, $I(x) \triangleq -\log P(x)$, other logarithms (base-2) is called bits or shannons.

Shannon entropy : $H(X) = E_{x \sim p}[I(X)] = -E_{x \sim p}[\log P(x)]$.

It gives a lower bound on the number of bits need on average to encode symbols drawn from P .

If X is continuous, $H(P)$ is differential entropy.

e.g. Discrete uniform distribution maximite discrete entropy within the distributions having the same number of states

$$\max_{\{p_i\}} \sum_{i=1}^k -p_i \log p_i = E \log p, \text{ s.t. } \sum_{i=1}^k p_i = 1. \Rightarrow p_i = \frac{1}{k}.$$

Examples of Entropy

If the sun rises in the East there is no information content, the sun always rises in the East.

If you toss an unbiased coin then there is information in whether it lands heads or tails up. If the coin is biased towards heads then there is more information if it lands tails.

Surprisal associated with an event is the negative of the logarithm of the probability of the event $-\log_2(p)$.

People use different bases for the logarithm, but it doesn't make much difference, it only makes a scaling difference. But if you use base 2 then the units are the familiar bits. If the event is certain, so that $p = 1$, the information associated with it is zero. The lower the probability of an event the higher the surprise, becoming infinity when the event is impossible.

Examples of Entropy

But why logarithms? The logarithm function occurs naturally in information theory. Consider for example the tossing of four coins. There are 16 possible states for the coins, HHHH, HHHT, ..., TTTT. But only four bits of information are needed to describe the state. HTHH could be represented by 0100.

$$4 = \log_2(16) = -\log_2(1/16).$$

Going back to the biased coin, suppose that the probability of tossing heads is $3/4$ and $1/4$ of tossing tails. If I toss heads then that was almost expected, there's not that much information. Technically it's $-\log_2(0.75) = 0.415$ bits. But if I toss tails then it is $-\log_2(0.25) = 2$ bits.

Examples of Entropy

This leads naturally to looking at the average information, this is our entropy :

$$-\sum p \log_2(p),$$

where the sum is taken over all possible outcomes.

(Note that when there are only two possible outcomes the formula for entropy must be the same when p is replaced by $1 - p$. And this is true here.)

Cross Entropy

Suppose you have a *model* for the probability of discrete events, call this p_k^M where the index k just means one of K possibilities. The sum of these probabilities must obviously be one.

And suppose that you have some empirical data for the probabilities of those events, p_k^E . With the sum again being one.

The cross entropy is defined as

$$-\sum_k p_k^E \ln(p_k^M).$$

It is a measure of how far apart the two distributions are.

Example of Cross Entropy

Suppose that you have a machine-learning algorithm that is meant to tell you whether a fruit is a passion fruit, orange or guava.

As a test you input the features for an orange. And the algorithm is going to output three numbers, perhaps thanks to the softmax function, which can be interpreted as the probabilities of the fruit in question (the orange) being one of P , O or G . Will it correctly identify it as an orange?

The model probabilities come out of the algorithm as

$$p_P^M = 0.13, \quad p_O^M = 0.69, \quad \text{and} \quad p_G^M = 0.18.$$

Ok, it's done quite well. It thinks the fruit is most likely to be an orange. But it wasn't 100% sure.

Example of Cross Entropy

Empirically we know that

$$p_P^E = 0, \quad p_O^E = 1, \quad \text{and} \quad p_G^E = 0,$$

because it definitely is an orange. The cross entropy is thus

$$-(0 \times \ln(0.13) + 1 \times \ln(0.69) + 0 \times \ln(0.18)) = 0.371.$$

The cross entropy is minimized when the model probabilities are the same as the empirical probabilities.

Example of Cross Entropy

To see this we can use Lagrange multipliers. Write

$$L = - \sum_k p_k^E \ln(p_k^M) - \lambda \left(\sum_k p_k^M - 1 \right).$$

The second term on the right is needed because the sum of the model probabilities is constrained to be one.

Now differentiate with respect to each model probability, and set the results to zero :

$$\frac{\partial L}{\partial p_k^M} = -\frac{p_k^E}{p_k^M} - \lambda = 0.$$

But since the sums of the two probabilities must be one we find that $\lambda = -1$ and $p_k^M = p_k^E$.

Example of Cross Entropy

Because the cross entropy is minimized when the model probabilities are the same as the empirical probabilities we can see that cross entropy is a candidate for a useful cost function when you have a classification problem.

If you take another look at the sections on MLE and on cost functions, and compare with the above on entropy you'll find a great deal of overlap and similarities in the mathematics. The same ideas keep coming back in different guises and with different justifications and uses.

Kullback-Leibler (KL) divergence

Two distributions $P(x)$, $Q(x)$

$H(P, Q) = -E_{x \sim p} \log Q(x)$ cross-entropy

$D_{KL}(P||Q) = E_{x \sim p} [\log \frac{P(x)}{Q(x)}] = -H(P) + H(P, Q)$ Extra information “distance”.

Properties : $D_{KL}(P||Q) = 0$ iff $P = Q$ a.e. $D_{KL}(P||Q) \geq 0$

But $D_{KL}(P||Q) \neq D_{KL}(Q||P)$, $0 \log 0 = \lim_{x \rightarrow 0} x \log x = 0$.

Outlines

Probability and Information Theory

Statistics and Machine Learning

References

Tasks of Machine Learning

- (1) Classification : find $f : \mathbb{R}^n \rightarrow \{1, \dots, k\}$ s.t. $y \approx f(x)$.
- (2) Regression : find $f : \mathbb{R}^n \rightarrow \mathbb{R}$, s.t. $y \approx f(x)$.
- (3) Anomaly Detection : find $P(x)$ for normal samples, predict abnormal samples \hat{x} with small prob $P(\tilde{x}) < \varepsilon$.
- (4) Imputation of missing values : predict $P(x)$ for $x = (x_1, \dots, x_n)$ then insert the value of x_i for a new sample \vec{x} with x_i missing.
- (5) Denoising : Given a corrupted example $\tilde{x} \in \mathbb{R}^n$, find a clean example $\tilde{x} \in \mathbb{R}^n$ s.t. $x \approx \tilde{x}$ with some noise, i.e. find $P(x|\tilde{x})$.

Tasks of Machine Learning

(6) Density Estimation or PMF estimation :

find $P_{model} : \mathbb{R}^n \rightarrow \mathbb{R}$, $P_{model}(x)$ for given samples \vec{x} . All previous problems, and clustering dimensionality reduction etc, could fall into this category.

This is rather difficult, computationally intractable.

supertrained learning : $P(y|x) = \frac{P(x,y)}{\sum_{y'} P(x,y')}$.

Learning conditional statistics (e.g. expectation). given a measure of diviation (loss function). $L(y, f)$ want to find the best f^* that minimize it i.e. $\min_f E_{x,y \sim P_{data}} L(y, f(x))$.

If $L = \|y - f\|_2^2$, then $f^*(x) = E_{y \sim P_{data}(y|x)}[y]$.

If $L = \|y - f\|_1$, then $f^*(x) = \text{conditional median}$.

Linear Regression

$$P_{data}(x, y) = \frac{1}{m} \sum_{i=1}^m \delta(x - x^{(i)}, y - y^{(i)})$$

$$f(x) = w^T x,$$

$$\min_{f=w^T x} E_{x,y \sim P_{data}(x,y)} \|y - f(x)\|_2^2 \approx \min_w \frac{1}{m} \sum_{i=1}^m \|y^{(i)} - w^T x^{(i)}\|_2^2$$

$$P_{data}(x, y) \approx P_{train}(x, y)$$

$$\nabla_w MSE_{train} = 0 \Leftrightarrow w = (X^{(train)T} X^{(train)})^{-1} X^{(train)T} y^{(train)}$$

Goals :

- 1) Make MSE_{train} small (under fitting if this is not achieved)
- 2) Make $MSE_{train} - MSE_{test}$ small (over fitting if this is not achieved)

Linear Regression

Performance : $MSE_{test} = \frac{1}{m^{(test)}} \|X^{(test)}_w - y^{(test)}\|_2^2$.

If $(X^{(train)}, y^{(train)})$ and $(X^{(test)}, y^{(test)}) \sim P_{data}(x, y)$, then $MSE_{train} = MSE_{test}$ for \hat{w} computed from 1).

However, in general, $MSE_{train} \leq MSE_{test}$, since \hat{w} is computed from 2).

Regularization

Instead of minimizing MSE_{train} , we take into account the model complexity, in the case of linear regression, $\lambda \|w\|_2^2$,
 $J(w) = MSE_{train} + \lambda \|w\|_e^2$, $\lambda \rightarrow 0$ over fitting.

$\lambda \rightarrow +\infty$, underfitting, optimal $\lambda \in (0, +\infty)$.

Statistical Estimators

Point Estimation : estimate Q in $f(x; Q)$ by \hat{Q} .

$$\begin{aligned}\hat{Q} &= \arg \min_Q E_{x,y \sim P_{data}(x,y)} L(y, f(x, Q)) \\ &= \arg \min_Q \frac{1}{m} \sum_{i=1}^m L(y^{(i)}, f(x^{(i)}; Q)) \\ \Rightarrow \hat{Q}_m &= g((x^{(i)}, y^{(n)}) \sim (x^{(m)}, y^{(m)}); Q)\end{aligned}$$

Function Estimation : nonparameterized $f(x)$ in some functional space \mathfrak{H} , then the least square rule gives \hat{f} as the best approximation of f among \mathfrak{H} .—point estimation in \mathfrak{H} .

$$\text{bais}(\hat{Q}_m) = E[\hat{Q}_m] - Q$$

\hat{Q} is unbiased if $\text{bais}(\hat{Q}_m) = 0$.

asymptotically unbiased if $\lim_{m \rightarrow \infty} \text{bais}(\hat{Q}_m) = 0$.

Statistical Estimators

e.g. Bernoulli, $\{x^{(1)}, \dots, x^{(m)}\}$, $P(x^{(i)}; Q) = Q^{x^{(i)}}(1 - Q)^{1-x^{(i)}}$.

Let $\hat{Q}_m = \frac{1}{m} \sum_{i=1}^m x^{(i)}$, $E[\hat{Q}_m] = \frac{1}{m} \sum_{i=1}^m E[x^{(i)}] = Q$, unbiased.

e.g. Gaussian, $\{x^{(1)}, \dots, x^{(m)}\}$,

$P(x^{(i)}) = N(x^{(i)}; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(x^{(i)} - \mu)^2}{2\sigma^2})$

$\hat{\mu}_m = \frac{1}{m} \sum_{i=1}^m x^{(i)}$, $E[\hat{\mu}_m] = \frac{1}{m} \sum_{i=1}^m E[x^{(i)}] = \mu$. unbiased

$\hat{\sigma}_m^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \hat{\mu}_m)^2$ sample variances

$E[\hat{\sigma}_m^2] = \frac{m-1}{m} \sigma^2 \neq \sigma^2$ biased

but $\hat{\sigma}_m^2 = \frac{1}{m-1} \sum_{i=1}^m (x^{(i)} - \hat{\mu}_m)^2 = \frac{m-1}{m} \hat{\sigma}_m^2$ is unbiased.

Variance and Standard Error

\hat{Q}_m : estimator. $\sqrt{\text{Var}(\hat{Q}_m)}$ = standard error = $SE(\hat{Q}_m)$

e.g. $SE(\hat{\mu}_m) = \sqrt{\text{Var}[\frac{1}{m} \sum_{i=1}^m x^{(i)}]} = \frac{\sigma}{\sqrt{m}}$

$CLT \Rightarrow \hat{\mu}_m \sim N(\mu, SE^2(\hat{\mu}_m))$

Confidence interval : $(\hat{\mu}_m - 1.96SE(\hat{\mu}_m), \hat{\mu}_m + 1.96SE(\hat{\mu}_m))$.

e.g. Bernoulli : $\text{Var}(\hat{Q}_m) = \frac{1}{m^2} \sum_{i=1}^m \text{Var}(x^{(i)}) = \frac{1}{m} Q(1 - Q),$

$SE(\hat{Q}_m) = \frac{\sqrt{Q(1-Q)}}{\sqrt{m}}$ decreasing function of m .

Bias-variance Tradeoff

$$\begin{aligned} MSE &= E[(\hat{Q}_m - Q)^2] = \text{Bias}(\hat{Q}_m)^2 + \text{Var}(\hat{Q}_m) \\ &= E[(\hat{Q}_m - E\hat{Q}_m)^2 + 2(\hat{Q}_m - E\hat{Q}_m)(E\hat{Q}_m - Q) + (E\hat{Q}_m - Q)^2] \end{aligned}$$

Variance and Standard Error

Underfitting : High bias, low variance

Overfitting : Low bias, high variance

Consistency :

\hat{Q}_m is a consistent estimator if $\hat{Q}_m \xrightarrow{P} Q \ (m \rightarrow \infty)$.

i.e. $P(\hat{Q}_m - Q|_{>\varepsilon}) \rightarrow 0 \ (m \rightarrow \infty)$ for $\forall \varepsilon > 0$.

Consistency \Rightarrow Asymptotic unbiasedness.

A counter-example for the reverse statement :

$x^{(i)} \sim N(x; \mu, \sigma^2)$, $\hat{\mu} = x^{(1)}$, then $E[\hat{\mu}] = E[x^{(1)}] = \mu$.

But $\hat{\mu}$ does not trends to μ as $m \rightarrow \infty$.

Chi-Squared Test

It tests a null hypothesis stating that two variables is statistically independent, and usually requires $R * C$ contingency table as following :

$G \setminus Y$	Y_1	Y_2	\dots	Y_C	sum
G_1	$A_{11} (T_{11})$	$A_{12} (T_{12})$		$A_{1C} (T_{1C})$	$\sum_{c=1}^C A_{1c}$
G_2	$A_{21} (T_{21})$	$A_{22} (T_{22})$		$A_{2C} (T_{2C})$	$\sum_{c=1}^C A_{2c}$
\dots	\dots	\dots			
G_R	$A_{R1} (T_{R1})$	$A_{R2} (T_{R2})$		$A_{RC} (T_{RC})$	$\sum_{c=1}^C A_{Rc}$
sum	$\sum_{r=1}^R A_{r1}$	$\sum_{r=1}^R A_{r2}$		$\sum_{r=1}^R A_{rC}$	N

Given the hypothesis of independence, define the theoretical frequencies $T_{ij} = \frac{\sum_{c=i}^C A_{ic} \sum_{r=1}^R A_{rj}}{N}$.

The chi-squared test statistic χ^2 , which resembles a normalized sum of squared deviations between observed and theoretical frequencies, is valid to perform when the test statistic is chi-squared distributed under the null hypothesis.

Chi-Squared Test

Sustain or reject the null hypothesis based on whether the test statistic exceeds the critical value t of χ^2 . If the test statistic exceeds the critical value of χ^2 , the null hypothesis (H_0 = the row variable is independent of the column variable) can be rejected, and the alternative hypothesis (H_1 = there is an association) can be accepted, both with the selected level of confidence.

$$\chi^2 = \sum_{i=1}^R \sum_{j=1}^C \frac{(A_{ij} - T_{ij})^2}{T_{ij}} \sim \chi^2((R-1)(C-1))$$

In the chi merge example in class, interval and class frequency are selected as two variables. If the calculated χ is smaller than the selected threshold t , the null hypothesis is accepted and the two intervals are merged.

Maximum Likelihood Estimation (MLE)

Usually used in parametric density estimation.

$P_{data}(x) \approx P_{model}(x; Q)$, $\{x^{(i)}\}_{i=1}^m$ drawn i.i.d from $P_{data}(x)$.

Assume $\exists Q$, as if $x^{(i)} \sim P_{model}(x; Q)$.

MLE for Q is

$$Q_{ML} = \arg \max_Q P_{model}(X; Q) = \arg \max_Q \prod_{i=1}^m P_{model}(x^{(i)}; Q)$$

or

$$\begin{aligned} Q_{ML} &= \arg \max_Q \prod_{i=1}^m \log P_{model}(x^{(i)}; Q) \\ &= \arg \max_Q E_{X \sim \hat{P}_{data}} \log P_{model}(X; Q) \end{aligned}$$

Property : Q_{ML} minimizes KL divergence (dissimilarity) between \hat{P}_{data} and P_{model}

$$D_{KL}(\hat{P}_{data} || P_{model}) = E_{X \sim \hat{P}_{data}} [\log \hat{P}_{data}(X) - \log P_{model}(X; Q)]$$

$$\min_Q D_{KL} \Leftrightarrow \min_Q E_{X \sim \hat{P}_{data}} [-\log P_{model}(X; Q)]$$

Examples of MLE

Maximum Likelihood Estimation (MLE) is a common method for estimating parameters in a statistical/probabilistic model.

In words, you simply find the parameter (or parameters) that maximizes the likelihood of observing what actually happened.

Let's see this in a few classical examples.

Examples of MLE

Example : Taxi numbers

You arrive at the train station in a city you've never been to before. You go to the taxi rank so as to get to your final destination. There is one taxi, you take it. While discussing European politics with the taxi driver you notice that the cab's number is 1234. How many taxis are in that city?

To answer this we need some assumptions. Taxi numbers are positive integers, starting at 1, no gaps and no repeats. We'll need to assume that we are equally likely to get into any cab. And then we introduce the parameter N as the number of taxis.

What is the MLE for N ?

Examples of MLE

Example : Taxi numbers

What is the MLE for N ?

Well, what is the probability of getting into taxi number 1234 when there are N taxis?

It is $\frac{1}{N}$ for $N \geq 1234$ and zero otherwise.

What value of N maximizes this expression? Obviously it is $N = 1234$. That is the MLE for the parameter N . It looks a bit disturbing because it seems a coincidence that you happened to get into the cab with the highest number. But then the probability of getting into any cab is equally likely. It is also disturbing that if there are N taxis then the average cab number is $(N + 1)/2$, and we somehow feel that should play a role.

Examples of MLE

Example : Coin tossing

Suppose you toss a coin n times and get h heads. What is the probability, p , of tossing a head next time?

The probability of getting h heads from n tosses is, assuming that the tosses are independent,

$$\frac{n!}{h!(n-h)!} p^h (1-p)^{n-h} = \binom{n}{h} p^h (1-p)^{n-h}.$$

Applying MLE is the same as maximizing this expression with respect to p .

Examples of MLE

Example : Coin tossing

Often with MLE when multiplying probabilities, as here, you will take the logarithm of the likelihood and maximize that.

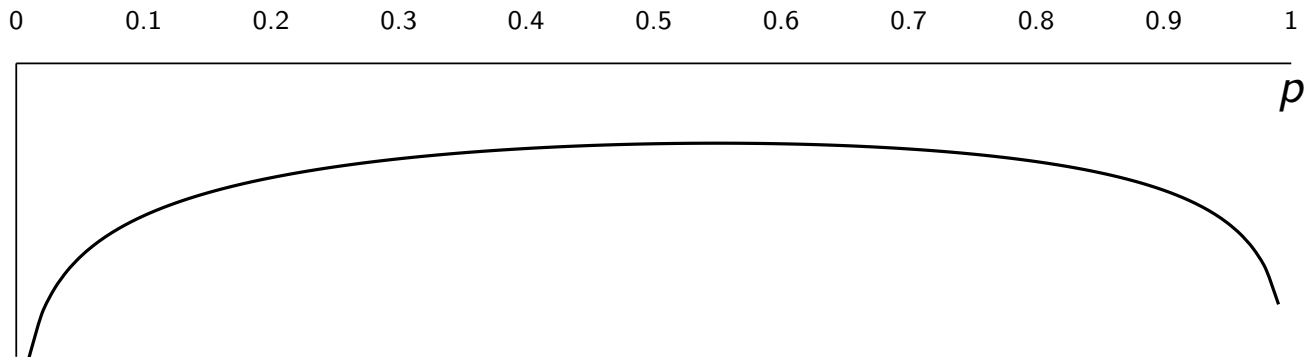
This doesn't change the maximizing value but it does stop you from having to multiply many small numbers, which is going to be problematic with finite precision.

(Look at the scale of the numbers on the vertical axis in the figure.)

Examples of MLE

Example : Coin tossing

Since the first part of this expression is independent of p we maximize $h \ln p + (n - h) \ln(1 - p)$ with respect to p . See below.



This just means differentiating with respect to p and setting the derivative equal to zero. This results in $p = \frac{h}{n}$, which seems eminently reasonable.

Conditional Log-likelihood

MLE for Q in $P(Y|X; Q)$

$$Q_{ML} = \arg \max_Q P(Y|X; Q) = \arg \max_Q \prod_{i=1}^m P(y^{(i)}|x^{(i)}; Q)$$

e.g. linear Regression : $y = w^T x + \varepsilon$, $\varepsilon \sim N(0, \sigma^2)$.

Then $P(y|x; w) \sim N(y; w^T x, \sigma^2)$, $\frac{1}{\sqrt{2\pi}\sigma} \exp(-\frac{(y-w^T x)^2}{2\sigma^2})$.

Conditional Log-likelihood

$$\begin{aligned}MLE &\Leftrightarrow \max_w \sum_{i=1}^m \log P(y^{(i)} | x^{(i)}; w) \\&= \max_w \sum_{i=1}^m \left(-\log \sigma - \frac{1}{2} \log(2\pi) - \frac{(y^{(i)} - w^T x^{(i)})^2}{2\sigma^2} \right) \\&= -m \log \sigma - \frac{m}{2} \log(2\pi) - \min_w \sum_{i=1}^m (y^{(i)} - w^T x^{(i)})^2 / \sigma^2 \\&\Leftrightarrow \min_w \frac{1}{m} \sum_{i=1}^m (y^{(i)} - w^T x^{(i)})^2 = \min_w MSE_{train}\end{aligned}$$

Properties of MLE

Under certain conditions (given below). MLE is a consistent estimator of the truth.

- (1) P_{data} lies in $\{P_{\text{model}}(; Q); Q\}$; otherwise, no estimator can recover P_{data} .
- (2) $\exists Q$, s.t. $P_{data} = P_{\text{model}}(; Q)$; otherwise, MLE can recover P_{data} , but not be able to determine Q .

$E_{X \sim P_{data}}[\hat{Q}_{ML} - Q]^2 \searrow$ Cramer-Rao bound as $m \rightarrow \infty$.

But MLE is not always unbiased, e.g. $\hat{\sigma}_{ML}^2 = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \hat{\mu})^2$

Bayesian Statistics

Previously, Q is fixed but unknown. \hat{Q} is a random variable as a function of data $\{x^{(i)}\}_{i=1}^m$ ($x^{(i)}$ is random).

Bayesian : $\{x^{(i)}\}_{i=1}^m$ is observed and non-random. Q is unknown and uncertain (random).

Prior prob distribution $P(Q)$, before observing the data. Given Q , $\{x^{(i)}\}_{i=1}^m$ is generated from $P(x^{(1)}, \dots, x^{(m)} | Q)$.

Bayes' rule $\Rightarrow P(Q | x^{(1)}, \dots, x^{(m)}) = \frac{P(x^{(1)}, \dots, x^{(m)} | Q) P(Q)}{P(x^{(1)}, \dots, x^{(m)})}$.

$P(Q)$ is usually given e.g. uniform or Gaussian with high entropy, observation of data causes the posterior to loose entropy and concentrate around a few highly likely values of parameters. Bayesian estimates the distribution of Q instead of point estimate.

$$P \in X^{(m+1)} | x^{(1)}, \dots, x^{(m)} = \int P(x^{(m+1)} | Q) P(Q | x^{(1)}, \dots, x^{(m)}) dQ$$

As more observations are given, knowledge about Q becomes different (more).

Bayesian Statistics

e.g. Bayesian Linear Regression $\hat{y} = w^T x$.

Given $(X^{(train)}, y^{(train)})$, $\hat{y}^{(train)} = X^{(train)} w$

$$P(y^{(train)} | X^{(train)}, w) = N(y^{(train)}; X^{(train)} w, I) \exp(-\frac{1}{2} (y^{(train)} - X^{(train)} w)^T (y^{(train)} - X^{(train)} w))$$

Prior of w : $P(w) = N(w; \mu_0, \Lambda_0) \exp(-\frac{1}{2} (w - \mu_0)^T \Lambda_0^{-1} (w - \mu_0))$

Posterior : $P(w | X^{(train)}, y^{(train)}) P(y^{(*)} | X^{(train)}, w) P(w)$

$$\exp(-\frac{1}{2} (y - Xw)^T (y - Xw)) \exp(-\frac{1}{2} (w - \mu_0)^T \Lambda_0^{-1} (w - \mu_0))$$

$$\exp(-\frac{1}{2} (-2y^T Xw + w^T X^T Xw + w^T \Lambda_0^{-1} w - 2\mu_0^T \Lambda_0^{-1} w))$$

$$\exp(-\frac{1}{2} (w - \mu_m)^T \Lambda_m^{-1} (w - \mu_m) + \frac{1}{2} \mu_m^T \Lambda_m^{-1} \mu_m)$$

$$\exp(-\frac{1}{2} (w - \mu_m)^T \Lambda_m^{-1} (w - \mu_m))$$

terms without w are normalizing constant.

If $\mu_0 = 0$, $\Lambda_0 = \frac{1}{\alpha} I$, $\mu_m = (X^T X + \alpha I)^{-1} X^T y$ is the ridge regression estimator of w .

also gives the variance $\Lambda_m = (X^T X + \alpha I)^{-1}$

Maximum A Posterior (MAP) similar to MLE

e.g. μ_m is MAP in the previous example.

MAP choose the point of maximum posterior prob.

$$Q_{MAP} = \arg \max_Q P(Q|X) = \arg \max_Q \{\log P(X|Q) + \log P(Q)\}$$

Bayesian linear regression, $\log -prior \|w\|_2^2$ ridge penalties.

Additional information in prior helps to reduce the variance in the MAP, but does not increase bias.

Different regularizations (penalties) corresponds to different log-prior. (but not all, smoe penalty may not be a logarithm of a prob distribution, some others depend on data)

Lsso Penalty \rightarrow Laplace distribution $\text{Laplace}(x; 0, \lambda^{-1})$.

Outlines

Probability and Information Theory

Statistics and Machine Learning

References

References

- Numerical Analysis, 9th Edition, by Richard L. Burden, J. Douglas Faires, Brooks/Cole, 2011.
- Machine learning : an applied mathematics introduction, by Wilmott, Paul. Panda Ohana Publishing, 2019.