

MATHEMATICAL STATISTICS

MINSEOK SONG

Mathematical Statistics

1 Sufficient statistic

Definition 1. *Statistic T is called sufficient if $X | T$ does not depend on θ .*

- (*) We need **measure theoretic** definition of sufficiency. What exactly is meant that samples from $X | T$ follow the same distribution as X ?
- (*) What I am getting at is that $P(X \in A) = P(X \in A | \sigma(w : T(X(w)) \in T(A)))(w)$ for $w \in \{w : T(X(w)) \in T(A)\}$. To be continued (check!)
- The setting here is that we are trying to infer θ based on data X . Bayesian definition would be phrased as $X | T \perp \theta$
- The intuition here is that having X given T does not give further information about θ . Just by definition, X (n samples) has the same distribution as a sample from $X|T$ (also n samples).

Example 1.

For example $X|T$ has a distribution depending on $T(X)$ such as

$$N(1^T \bar{X}, \begin{pmatrix} 1 - 1/n & -1/n & \dots \\ -1/n & 1 - 1/n & \dots \\ \dots & \dots & 1 - 1/n \end{pmatrix})$$

and if we sample \tilde{X}_i 's from this distribution, it does follow $N(\theta, 1)$ (the conditional distribution follows that above distribution).

Another example would be sum for Bernoulli, and maximum value for uniform distribution. By the way, maximum and minimum of samples from uniform distribution are both sufficient, illustrating that 1d sufficient statistic is not necessarily minimal.

- Any one to one function of minimal sufficient statistic is also minimal sufficient statistic. For example, $(T(x), T(x))$ would do.
- Lehmann-Scheffe theorem establishes the relationship between unbiasedness and sufficiency of statistic.

Theorem 1. (*factorization theorem*) *Sufficient statistic T satisfies $P(x) = f_\theta(T(x))g(x)$, that is, θ depends only through $T(x)$.*

Definition 2. *A sufficient statistic T is called minimal if any other sufficient statistic M is a function of T .*

Lemma 2. *Suppose $H_0 \subset H$. Suppose S is minimal sufficient for $P_\theta, \theta \in H_0$ and sufficient for $P_\theta, \theta \in H$. Then it is minimal sufficient for $P_\theta, \theta \in H$.*

- This is because any sufficient statistic in a larger parameter space is sufficient in a smaller set of parameter space by definition.

Theorem 3. For $P_\theta : \theta \in \{\theta_1, \dots, \theta_d\}$, $T(x) = (\frac{P_{\theta_1}(x)}{P_{\theta_0}(x)}, \frac{P_{\theta_2}(x)}{P_{\theta_0}(x)}, \dots, \frac{P_{\theta_d}(x)}{P_{\theta_0}(x)})$ is a minimal sufficient statistic.

Proof. Sufficiency by definition. Minimality due to factorization theorem. \square

•

DS procedure

The dichotomy of Bayesian and Frequentist approach is mostly a matter of division based on foundational philosophical differences, and in some situation is not necessarily natural; each instance can be formally simply seen as a choice of analysis. The dichotomy has mainly developed based on two philosophical school of thoughts: subjective belief (via prior distribution) vs. long-run frequency (via sampling distribution).

(1) Frequentist

- They avoid making probabilistic claims about model parameters and hypothesis. Instead, they describe the behavior of statistics and procedures over many hypothetical repeated samples.
- They use Statistics derived from data, and use deterministic approach for inference, with methods like hypothesis testing and confidence intervals to draw conclusions about population parameters based on sample data.
- examples: T-test, linear regression, etc

(2) Bayesian

- They assign prior beliefs (distribution function) on parameters of model.
- They use probabilistic argument on specific hypothesis or parameter values.
- examples: MCMC, Bayesian hierarchical modeling, etc

2 Data Science Procedure

[list of topics that will be discussed]

(1) Preparation of data

- handling missing data (imputation, etc)
- transformation
- Box-Cox transformation
- Outliers

(2) Interpret descriptive statistics about data

- plots
- t-distribution
- F-distribution

(3) Modeling (parametric, nonparametric)

- ANOVA
- Ensemble methods
- Boosting
- AIC-BIC
- Regularization

(4) Predict based on model

- Confidence/Prediction intervals
- Real-time prediction

(5) Inference on statistics

- Hypothesis testing
- p-value

- Type I/II error
- Causal inference
- (6) Evaluation of models
 - Cross-validation
- (7) Communication of models

2.1 Types of Missing Data

- (1) MCAR
- (2) MAR
- (3) MNAR

3 Outliers

4 Inference on statistics

5 p-value

Definition 3. *Given some hypothesis, the p-value is the probability of obtaining test results at least as extreme as the result actually observed, under the assumption that the null hypothesis is correct.*

- If p-value is less than significance level, we reject the null hypothesis. Significance level is a threshold we set to decide when we have enough evidence to reject the null hypothesis in favor of the alternative hypothesis.
- In statistics, we do not quite "prove" that the hypothesis is true. We can only reject the false statements or assertions (the theory of falsification, proposed by Karl Popper).

6 t-distribution

Definition 4. *t-distribution is defined as $\frac{Z}{\sqrt{\frac{V}{k}}}$ where $Z \perp V$, $Z \sim N(0, 1)$, and $V \sim \chi_k^2$.*

- The t-distribution has the thicker tails than the normal distribution.
- Note that the χ_k^2 distribution has mean k and variance $2k$. Hence as k gets larger the t-distribution looks like standard normal distribution due to CLT.
- In multiple regression setting, $\frac{\frac{\hat{\beta}_1 - \beta_1}{\sigma}}{\sqrt{\frac{\sum_i (X_i - \bar{X})^2}{(n-2) \cdot \frac{\hat{\sigma}^2}{n-2}}}} \sim t_{n-2}$ under normal assumption of residual.

7 Chi-squared distribution

- It is possible to have $f(X)$ and $g(X)$ be independent. The key example concerns chi-distribution: Say θ is uniformly distributed in $[0, 2\pi]$. Let us say R^2 follows *chi*-squared distribution with degree 2: this acts like a distance square in the 2d space. Then $X = r \cos(\theta)$ and $Y = r \sin(\theta)$ are independent.
- The definition of chi-square looks contrived, in particular compared to normal distribution, but we use it a lot in the context of sum of squares.
- Under moment assumption, it's approximately true for large n via CLT.
- Under moment assumption, it's approximately true for large n via CLT.

8 Prediction interval, Confidence interval

- In the context of linear regression, a **confidence interval** provides a range of values within which we expect the true regression coefficient (e.g., β_0 or β_1) to lie with a certain level of confidence. This reflects our uncertainty about the value of the coefficient based on the data we have observed.

- On the other hand, a **prediction interval** provides a range of values within which a new observation y (given a particular x) is expected to fall with a certain level of confidence. We have $y = \beta_0 + \beta_1 x + \epsilon$. Note the additional uncertainty introduced by the random error ϵ .

8.1 Fisher information

Definition 5. Fisher information is defined by $\mathcal{I}(\theta) = E[(\frac{\partial}{\partial \theta} \log f(X; \theta))^2 | \theta]$ for appropriately regular f .

- Note that $E[\frac{\partial}{\partial \theta} \log f(X; \theta) | \theta] = 0$.
- Think of this as how much (the logarithm of) likelihood varies in the vicinity of θ .
- Large variation of likelihood we have more "information." (for example, if a deviation causes a significant drop in likelihood, it means that the data is most probable under the original θ value.)

Theorem 4 (Cramer-Rao bound). Say $\hat{\theta}(X)$ is an unbiased estimator. We have $\text{Var}(\hat{\theta}) \geq \frac{1}{\mathcal{I}(\theta)}$.

- This says that as the fisher information gets larger, we have better precision.
- There is a similar result for biased estimator.

Linear Model

- Two assumptions: normality($Y_i \sim N(\beta_0 + \beta_1 X_i, \sigma^2)$) and moment assumption.

Fact 1. For Normal distribution, pair-wise independence implies mutual independence.

Proof. It comes from the fact that the mutual independence holds if and only if density function factorises. In particular, multivariate normal distribution with diagonal matrix covariate satisfies this property. \square

- For 1d case, we could calculate $\hat{\beta}$ based on $\partial/\partial \beta(RSS)$...
- After computing β_0 and β_1 , we can also compute their respective variance and covariance between them. I need to use the fact that \bar{Y} and $\hat{\beta}_1$, which could also be easily calculated.
- We could further define $\hat{\epsilon}$, and unbiased estimator $\hat{\sigma}$. Showing that it is unbiased and is independent to $\hat{\beta}_0$ and $\hat{\beta}_1$ is more involved (need to check).
- Knowing that $\frac{\hat{\beta}_1 - \beta_1}{\sigma/\sqrt{\sum_i (X_i - \bar{X})^2}} \sim N(0, 1) \perp$ (because $\hat{\beta}_1 \perp \hat{\sigma}^2$) $(n-2) \cdot \frac{\hat{\sigma}^2}{\sigma^2} \sim \chi_{n-2}^2$ we can do some more analysis involving confidence interval of predictive value and mean value at a given point.
- Some **intuition** about why the variance of β is inversely proportional to dispersion: well, as the data gets concentrated, the prediction is rather **unreliable**.
- The reason we used t-distribution is to cancel out the **unknown value** σ .
 - RSS = portion that is not explained by the model
 - TSS = total sum of squares = $\|Y - \bar{Y}1_n\|^2$. Note that $Y - \bar{Y}1_n \in \text{span}(X)$
 - ESS = explained sum of squares = $\|\hat{Y} - \bar{Y}1_n\|^2$
- $RSS + ESS = TSS$. The role of \bar{y} is like the "standardized line." Refer to here.
- R square is ESS/TSS .
- We would use F test... F statistic is given by $F = \frac{\|U_2^T Y\|^2/k}{\|U_3^T Y\|^2/(n-p)} = \frac{RSS_{\text{partial model}} - RSS_{\text{full model}}}{RSS_{\text{full model}}/(n-p)}$
- We want F to be as small as possible for the hypothesis $\beta = 0$ to be true, to accept partial model. In other words, $RSS_{\text{partial model}}$ tends to be larger if $\beta_S \neq 0$.