

Statistical methods, Formula sheet for final exam

Combinatorics

- Number of ways to choose k out of n objects:

$$\binom{n}{k} = \frac{n(n-1)\cdots(n-k+1)}{k!} = \frac{n!}{(n-k)!k!}$$

Basic probability

Always true:

- $P(A^c) = 1 - P(A)$
- $P(A \cup B) = P(A) + P(B) - P(A \cap B)$
- $P(A \cap B) = P(A)P(B|A)$
- $P(A^c) = 1 - P(A)$
- A, B mutually exclusive: $P(A \cup B) = P(A) + P(B)$
- A, B independent: $P(A \cap B) = P(A)P(B)$
- Conditional probability of A given B :

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

LTP and Bayes' theorem

$$P(A) = P(A|B)P(B) + P(A|B^c)P(B^c)$$

$$P(B|A) = \frac{P(A|B)P(B)}{P(A|B)P(B) + P(A|B^c)P(B^c)}$$

Discrete random variables

- Pmf: $p(x) = P(X = x)$

Special distributions:

- $X \sim \text{bin}(n, p)$: $p(k) = \binom{n}{k} p^k (1-p)^{n-k}$, $k = 0, 1, \dots, n$ (# of successes)

$$E[X] = np$$

- $X \sim \text{geom}(p)$: $p(k) = (1-p)^{k-1} p$, $k = 1, 2, \dots$ (wait for first success)

$$E[X] = \frac{1}{p}$$

Continuous random variables

- Pdf: $f(x) = F'(x)$, $x \in R$

- Cdf: $F(x) = \int_{-\infty}^x f(t) dt$, $x \in R$

Special distributions:

- $X \sim \text{unif}[a, b]$: $f(x) = \frac{1}{b-a}$, $a \leq x \leq b$ (choose “randomly”)

$$E[X] = \frac{a+b}{2}$$

- $X \sim \text{exp}(\lambda)$: $f(x) = \lambda e^{-\lambda x}$, $x \geq 0$ (memoryless)

$$E[X] = \frac{1}{\lambda}$$

- $X \sim N(0, 1)$: $\varphi(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$, $x \in R$

- $X \sim N(\mu, \sigma^2)$: $Z = \frac{X - \mu}{\sigma} \sim N(0, 1)$

Expected value

- $E[X] = \sum x_k p(x_k)$ if X is discrete with range $\{x_1, x_2, \dots\}$
- $E[X] = \int_{-\infty}^{\infty} x f(x) dx$ if X is continuous
- $E[g(X)] = \sum g(x_k) p_X(x_k)$
- $E[g(X)] = \int_{-\infty}^{\infty} g(x) f_X(x) dx$

Variance

- $\text{Var}[X] = E[(X - \mu)^2] = E[X^2] - (E[X])^2$
- Standard deviation: $\sigma = \sqrt{\text{Var}[X]}$

Sums of random variables

- X and Y independent, a and b constants:

$$E[aX + bY] = aE[X] + bE[Y] \quad \text{Var}[aX + bY] = a^2\text{Var}[X] + b^2\text{Var}[Y]$$

- X_1, \dots, X_n independent random variables with the same distributions, mean μ and variance σ^2 , $S_n = X_1 + \dots + X_n$.
- $E[S_n] = n\mu$ and $\text{Var}[S_n] = n\sigma^2$
- Central Limit Theorem: S_n is approximately $N(n\mu, n\sigma^2)$ and \bar{X} is approximately $N(\mu, \sigma^2/n)$.

Estimators

- Unbiased: $E[\hat{\theta}] = \theta$

- We want $\text{Var}[\hat{\theta}]$ to be as small as possible
- Sample mean \bar{X} , unbiased for mean μ , $\text{Var}[\bar{X}] = \sigma^2/n$
- Sample variance $s^2 = \frac{1}{n-1} \left(\sum_{k=1}^n X_k^2 - n\bar{X}^2 \right)$, unbiased for variance σ^2
- Standard error: $\sqrt{\text{Var}[\hat{\theta}]}$

Confidence intervals

- For μ in $N(\mu, \sigma^2)$ where σ^2 is known:

$$\mu = \bar{X} \pm z \frac{\sigma}{\sqrt{n}} \quad (q)$$

Use standard normal distribution, $\Phi(z) = (1+q)/2$.

- For μ in $N(\mu, \sigma^2)$ where σ^2 is unknown:

$$\mu = \bar{X} \pm t \frac{s}{\sqrt{n}} \quad (q)$$

Use t distribution, $\alpha = 1 - (1+q)/2$, $\nu = n - 1$.

- For unknown proportion p :

$$p = \hat{p} \pm z \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \quad (q)$$

- For two unknown proportions p_1 and p_2 :

$$p_1 - p_2 = \hat{p}_1 - \hat{p}_2 \pm z \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}} \quad (q)$$

In both cases, z is such that $\Phi(z) = (1+q)/2$.

Estimation methods

1. The maximum likelihood estimator (MLE) $\hat{\theta}$ maximizes the likelihood function:

$$L(\theta) = \prod_{k=1}^n f_{\theta}(X_k)$$

To find maximum, (i) take logarithm, (ii) differentiate w.r.t. θ and set = 0.

2. The r th moment and r th sample moment are:

$$\mu_r = E[X^r] \quad \text{and} \quad \hat{\mu}_r = \frac{1}{n} \sum_{k=1}^n X_k^r$$

The method of moments estimator (MOME) expresses the parameter as a function of moments, $\theta = g(\mu_1, \dots, \mu_r)$ and estimates it with the same function of the sample moments, $\hat{\theta} = g(\hat{\mu}_1, \dots, \hat{\mu}_r)$. Start with the first moment, if it is not enough, go on to the second, and so on.

Linear regression

- Model: $Y = a + bx + \varepsilon$, $\varepsilon \sim N(0, \sigma^2)$
- Observations: $(x_1, Y_1), \dots, (x_n, Y_n)$ where $Y_k \sim N(a + bx_k, \sigma^2)$
- Notation:

$$S_x = \sum_{k=1}^n x_k, \quad S_Y = \sum_{k=1}^n Y_k$$
$$S_{xx} = \sum_{k=1}^n x_k^2, \quad S_{xY} = \sum_{k=1}^n x_k Y_k$$

- Estimators:

$$\hat{b} = \frac{nS_{xY} - S_x S_Y}{nS_{xx} - S_x^2}$$
$$\hat{a} = \bar{Y} - \hat{b}\bar{x}$$

Estimated regression line: $y = \hat{a} + \hat{b}x$