

Exam Formula Sheet

Epi 202: Probability

$$\begin{aligned}\text{Var}(\tilde{a} \cdot \tilde{X}) &= \text{Var}\left(\sum_{i=1}^n a_i X_i\right) \\ &= \tilde{a}^\top \text{Var}(\tilde{X}) \tilde{a} \\ &= \sum_{i=1}^n \sum_{j=1}^n a_i a_j \text{Cov}(X_i, X_j) \\ \text{E}[Y] &= \text{E}[\text{E}[Y | X]]\end{aligned}$$

$$\text{E}[Y | Z] = \text{E}[\text{E}[Y | X, Z] | Z]$$

$$\text{Var}(Y) = \text{E}[\text{Var}(Y | X)] + \text{Var}(\text{E}[Y | X])$$

$$\text{Cov}(Y, Z) = \text{E}[\text{Cov}(Y, Z | X)] + \text{Cov}(\text{E}[Y | X], \text{E}[Z | X])$$

Epi 203: Statistical inference

$$\mathcal{L}(\theta) \stackrel{\text{def}}{=} \text{p}(\tilde{X} = \tilde{x} | \Theta = \theta)$$

$$\ell \stackrel{\text{def}}{=} \log\{\mathcal{L}(\tilde{x} | \theta)\}$$

$$\ell' \stackrel{\text{def}}{=} \frac{\partial}{\partial \theta} \ell(\tilde{x} | \theta)$$

$$\ell'' \stackrel{\text{def}}{=} \frac{\partial}{\partial \tilde{\theta}} \frac{\partial}{\partial \tilde{\theta}^\top} \ell(\tilde{x} | \tilde{\theta})$$

$$\ell''_{ij} = \frac{\partial}{\partial \theta_i} \frac{\partial}{\partial \theta_j} \ell(\tilde{X} = \tilde{x} | \tilde{\theta})$$

$$I \stackrel{\text{def}}{=} -\ell''(\tilde{x} | \tilde{\theta})$$

$$\mathcal{J} \stackrel{\text{def}}{=} \text{E}[I(\tilde{x} | \theta)]$$

$$\hat{\theta}_{ML} \sim \text{N}(\theta, [\mathcal{J}(\tilde{\theta})]^{-1})$$

For one parameter θ_k :

$$\text{CI}_{1-\alpha}(\theta_k) = \left[\hat{\theta}_k \pm z_{1-\frac{\alpha}{2}} \widehat{\text{SE}}(\hat{\theta}_k) \right]$$

$$Z_k \stackrel{\text{def}}{=} \frac{\hat{\theta}_k - \theta_{k,0}}{\widehat{\text{SE}}(\hat{\theta}_k)} \sim \text{N}(0, 1) \quad \text{under } H_0 : \theta_k = \theta_{k,0} \quad z_k = \text{observed } Z_k$$

$$\begin{aligned}p\text{-value} &= 2 \Pr(|Z| \geq |z_k|), \quad Z \sim \text{N}(0, 1) \\ &= 2(1 - \Phi(|z_k|))\end{aligned}$$

Sta 108: Linear regression

$$t_k \stackrel{\text{def}}{=} \frac{\hat{\beta}_k - \beta_{k,0}}{\widehat{\text{SE}}(\hat{\beta}_k)} \sim t_{n-p} \quad \text{under } H_0 : \beta_k = \beta_{k,0} \quad t_k^{\text{obs}} = \text{observed } t_k$$

$$\text{CI}_{1-\alpha}(\beta_k) = \left[\hat{\beta}_k \pm t_{n-p} \left(1 - \frac{\alpha}{2}\right) \widehat{\text{SE}}(\hat{\beta}_k) \right]$$

Let $T_{n-p} \sim t_{n-p}$.

$$\begin{aligned} p\text{-value} &= 2 \Pr(|T_{n-p}| \geq |t_k^{\text{obs}}|) \\ &= 2(1 - F_{t_{n-p}}(|t_k^{\text{obs}}|)) \end{aligned}$$

$$\text{CI}_{1-\alpha}(\mu(\tilde{x}^*)) = \left[\hat{\mu}(\tilde{x}^*) \pm t_{n-p} \left(1 - \frac{\alpha}{2}\right) \widehat{\text{SE}}(\hat{\mu}(\tilde{x}^*)) \right]$$

Y^* denotes a new observation (not in the training data), with corresponding covariate pattern \tilde{x}^* . Let $\hat{Y}^* \stackrel{\text{def}}{=} \hat{\mu}(\tilde{x}^*)$.

$$\text{PI}_{1-\alpha}(Y^*|\tilde{x}^*) = \left[\hat{Y}^* \pm t_{n-p} \left(1 - \frac{\alpha}{2}\right) \widehat{\text{SE}}(Y^* - \hat{Y}^*) \right]$$

$$\widehat{\text{SE}}(Y^* - \hat{Y}^*) = \hat{\sigma} \sqrt{1 + (\tilde{x}^*)^\top (\mathbf{X}'\mathbf{X})^{-1} \tilde{x}^*}$$

$$\text{Var}(Y^* - \hat{Y}^*) = \sigma^2 (1 + (\tilde{x}^*)^\top (\mathbf{X}'\mathbf{X})^{-1} \tilde{x}^*)$$

Let $\hat{\Sigma} \stackrel{\text{def}}{=} \widehat{\text{Var}}(\hat{\beta}) = \hat{\sigma}^2 (\mathbf{X}'\mathbf{X})^{-1}$.

$$\widehat{\text{Var}}(\hat{\mu}(\tilde{x})) = \tilde{x}^\top \hat{\Sigma} \tilde{x}$$

Let $\Delta\mu(\tilde{x}, \tilde{x}^*) = \mu(\tilde{x}) - \mu(\tilde{x}^*)$, and let $\Delta\tilde{x} = \tilde{x} - \tilde{x}^*$; then:

$$\widehat{\text{Var}}(\widehat{\Delta\mu}(\tilde{x}, \tilde{x}^*)) = \Delta\tilde{x}^\top \hat{\Sigma} \Delta\tilde{x}$$

Epi 204: Generalized linear models

Generalized linear models have three components:

1. The **outcome distribution** family: $p(Y|\mu(\tilde{x}))$
2. The **link function**: $g(\mu(\tilde{x})) = \eta(\tilde{x})$
3. The **linear component**: $\eta(\tilde{x}) = \tilde{x} \cdot \beta$

$$\theta_\omega(\tilde{x}, \tilde{x}^*) = \exp\{(\Delta\tilde{x}) \cdot \tilde{\beta}\}$$

Estimates of odds ratios from 2x2 contingency tables

$$\hat{\theta} = \frac{ad}{bc}$$

Summary of logistic regression definitions and results

Odds and log-odds

Odds

$$\omega \stackrel{\text{def}}{=} \frac{\Pr(A)}{\Pr(\neg A)}$$

Conditional odds

$$\omega(A|B) \stackrel{\text{def}}{=} \frac{\Pr(A|B)}{\Pr(\neg A|B)}$$

Odds function

$$\text{odds}\{\pi\} \stackrel{\text{def}}{=} \frac{\pi}{1 - \pi}$$

Probability to odds

$$\omega = \frac{\pi}{1 - \pi}$$

Odds function equals odds

$$\omega = \text{odds}\{\pi\}$$

Simplified odds expressions

$$\text{odds}\{\pi\} = \frac{1}{\pi^{-1} - 1} = (\pi^{-1} - 1)^{-1}$$

Odds of a non-event

$$\omega(\neg A) = \frac{1 - \pi}{\pi} = \pi^{-1} - 1$$

Odds ratio

$$\theta(\omega_1, \omega_2) \stackrel{\text{def}}{=} \frac{\omega_1}{\omega_2}$$

OR as ratio of probability ratios

$$\begin{aligned}\theta(\omega_1, \omega_2) &= \frac{\omega_1}{\omega_2} \\ &= \frac{\left(\frac{\pi_1}{1 - \pi_1}\right)}{\left(\frac{\pi_2}{1 - \pi_2}\right)}\end{aligned}$$

Odds ratios are reversible

$$\theta_\omega(A|B) = \theta_\omega(B|A)$$

Conditional ORs are reversible

$$\theta_\omega(A|B, C) = \theta_\omega(B|A, C)$$

Inverse-odds and probability recovery

Odds to probability

$$\pi = \frac{\omega}{1 + \omega}$$

Inverse-odds function

$$\text{invodds}\{\omega\} \stackrel{\text{def}}{=} \frac{\omega}{1 + \omega}$$

Probability as inverse-odds

$$\pi = \text{invodds}\{\omega\}$$

Simplified inverse-odds

$$\text{invodds}\{\omega\} = \frac{1}{1 + \omega^{-1}} = (1 + \omega^{-1})^{-1}$$

One minus inverse-odds

$$1 - \pi = \frac{1}{1 + \omega}$$

Complement of inverse-odds

$$1 + \omega = \frac{1}{1 - \pi}$$

Log-odds (logit) and expit

Log-odds

$$\eta \stackrel{\text{def}}{=} \log\{\omega\}$$

Log-odds from probability

$$\eta = \log\left\{\frac{\pi}{1 - \pi}\right\}$$

Logit function

$$\text{logit}(\pi) \stackrel{\text{def}}{=} \log\{\text{odds}\{\pi\}\}$$

Logit expanded

$$\text{logit}(\pi) = \log\left\{\frac{\pi}{1 - \pi}\right\}$$

Log-odds equals logit

$$\eta = \text{logit}\{\pi\}$$

Odds from log-odds

$$\omega = \exp\{\eta\}$$

Probability from log-odds

$$\pi = \frac{\exp\{\eta\}}{1 + \exp\{\eta\}}$$

Expit / inverse-logit function

$$\text{expit}(\eta) \stackrel{\text{def}}{=} \text{invodds}\{\exp\{\eta\}\}$$

Expit expressions

$$\text{expit}(\eta) = \frac{\exp\{\eta\}}{1 + \exp\{\eta\}} = (1 + \exp\{-\eta\})^{-1}$$

Probability as expit

$$\pi = \text{expit}\{\eta\} \tag{1}$$

Logit and expit are inverses

$$\text{logit}\{\text{expit}\{\eta\}\} = \eta \quad \text{expit}\{\text{logit}\{\pi\}\} = \pi$$

$$\left[\pi \stackrel{\text{def}}{=} \Pr(Y = 1 | \tilde{X} = \tilde{x}) \right] \xrightarrow[\frac{\omega}{1+\omega}]{\frac{\pi}{1-\pi}} \underbrace{\left[\omega \stackrel{\text{def}}{=} \text{odds}(Y = 1 | \tilde{X} = \tilde{x}) \right]}_{\text{expit}(\eta)} \xrightarrow[\exp\{\eta\}]{\log\{\omega\}} \left[\eta(\tilde{x}) \stackrel{\text{def}}{=} \text{log-odds}(Y = 1 | \tilde{X} = \tilde{x}) \right]$$

logit(π)

Figure 1: Diagram of logistic regression link and inverse link functions

Rare events

Odds minus probability

$$\omega - \pi = \frac{\pi^2}{1 - \pi}, \quad \text{where } \omega = \frac{\pi}{1 - \pi}$$

Derivatives

	π	ω	η	\tilde{x}	$\tilde{\beta}$
π	1	$(1 + \omega)^2$	$\frac{(1 + \omega)^2}{\omega}$	undef	undef
ω	$(1 - \pi)^2$	1	$\frac{1}{\omega}$	undef	undef
η	$\pi(1 - \pi)$	ω	1	undef	undef
\tilde{x}	$\tilde{\beta}\pi(1 - \pi)$	$\tilde{\beta}\omega$	$\tilde{\beta}$	\mathbb{I}	$\mathbf{0}$
$\tilde{\beta}$	$\tilde{x}\pi(1 - \pi)$	$\tilde{x}\omega$	\tilde{x}	$\mathbf{0}$	\mathbb{I}

Column labels indicate the numerators of the derivatives; row labels indicate the denominators.

Log-likelihood and score function

Log-likelihood component

$$\ell_i(\pi_i) = y_i \eta_i - \log\{1 + \omega_i\}$$

Score function as sum

$$\tilde{\ell}'(\tilde{\beta}) = \sum_{i=1}^n \tilde{\ell}'_i(\tilde{\beta})$$

Score component

$$\ell'_i(\tilde{\beta}) = \tilde{x}_i e_i$$

Score function

$$\tilde{\ell}'(\tilde{\beta}) = \sum_{i=1}^n \tilde{x}_i e_i = \mathbf{X}^\top \tilde{e}$$

MLE

One-sample MLE for odds

$$\hat{\omega} = \frac{x}{n - x}$$

Odds ratios in logistic regression

General OR formula

$$\theta_{\omega}(\tilde{x}, \tilde{x}^*) = \exp\{\eta(\tilde{x}) - \eta(\tilde{x}^*)\}$$

Difference in log-odds

$$\Delta\eta \stackrel{\text{def}}{=} \eta(\tilde{x}) - \eta(\tilde{x}^*)$$

OR in terms of $\Delta\eta$

$$\theta_{\omega}(\tilde{x}, \tilde{x}^*) = \exp\{\Delta\eta\}$$

$\Delta\eta$ from covariates

$$\Delta\eta = (\tilde{x} - \tilde{x}^*) \cdot \tilde{\beta}$$

Difference in covariate patterns

$$\Delta\tilde{x} \stackrel{\text{def}}{=} \tilde{x} - \tilde{x}^*$$

$\Delta\eta$ from $\Delta\tilde{x}$

$$\Delta\eta = \Delta\tilde{x} \cdot \tilde{\beta}$$

OR in terms of $\Delta\tilde{x}$

$$\theta_{\omega}(\tilde{x}, \tilde{x}^*) = \exp\{(\Delta\tilde{x}) \cdot \tilde{\beta}\}$$

Log OR equals $\Delta\eta$

$$\log\{\theta_{\omega}(\tilde{x}, \tilde{x}^*)\} = \Delta\eta$$

Inference for log-odds and odds ratios

Estimated SE of log-odds

$$\begin{aligned}\widehat{\text{Var}}(\hat{\eta}(\tilde{x})) &= \tilde{x}^{\top} \hat{\Sigma} \tilde{x} \\ \widehat{\text{SE}}(\hat{\eta}(\tilde{x})) &= \sqrt{\tilde{x}^{\top} \hat{\Sigma} \tilde{x}}\end{aligned}$$

Estimated SE of $\Delta\hat{\eta}$

$$\begin{aligned}\widehat{\text{Var}}(\Delta\hat{\eta}) &= \Delta\tilde{x}^{\top} \hat{\Sigma}(\Delta\tilde{x}) \\ \widehat{\text{SE}}(\Delta\hat{\eta}) &= \sqrt{\Delta\tilde{x}^{\top} \hat{\Sigma}(\Delta\tilde{x})}\end{aligned}$$

Comparing probabilities

Risk difference

$$\delta(\pi_1, \pi_2) \stackrel{\text{def}}{=} \pi_1 - \pi_2$$

Risk ratio

$$\rho(\pi_1, \pi_2) \stackrel{\text{def}}{=} \frac{\pi_1}{\pi_2}$$

Relative risk difference

$$\xi(\pi_1, \pi_2) \stackrel{\text{def}}{=} \frac{\delta(\pi_1, \pi_2)}{\pi_2}$$

RRD equals RR minus 1

$$\xi(\pi_1, \pi_2) = \rho(\pi_1, \pi_2) - 1$$

Logistic regression model

Logistic regression

$$Y_i | \tilde{X}_i \sim_{\perp} \text{Ber}(\pi(\tilde{X}_i)), \quad \text{logit}\{\pi(\tilde{x})\} = \tilde{x}^{\top} \tilde{\beta}$$

Survival analysis

Probability distribution functions

Table 1: Probability distribution functions

Name	Symbols	Definition
Probability density function (PDF)	$f(t), p(t)$	$p(T = t)$
Cumulative distribution function (CDF)	$F(t), P(t)$	$P(T \leq t)$
Survival function	$S(t), \bar{F}(t)$	$P(T > t)$
Hazard function	$\lambda(t), h(t)$	$p(T = t T \geq t)$
Cumulative hazard function	$\Lambda(t), H(t)$	$\int_{u=-\infty}^t \lambda(u) du$
Log-hazard function	$\eta(t)$	$\log\{\lambda(t)\}$

Diagram of survival distribution function relationships

$$f(t) \xleftarrow[\frac{-S'(t)}{S(t)\lambda(t)}]{} S(t) \xleftarrow[\frac{\exp\{-\Lambda(t)\}}{\Lambda(t)}]{} \Lambda(t) \xleftarrow[\frac{\int_{u=0}^t \lambda(u) du}{\lambda(t)}]{} \lambda(t) \xleftarrow[\frac{\exp\{\eta(t)\}}{\eta(t)}]{} \eta(t)$$

$$f(t) \xrightarrow[\frac{\int_{u=t}^{\infty} f(u) du}{f(t)/\lambda(t)}]{} S(t) \xrightarrow[-\log S(t)]{} \Lambda(t) \xrightarrow[\Lambda'(t)]{} \lambda(t) \xrightarrow[\log\{\lambda(t)\}]{} \eta(t)$$

Survival likelihood contributions, assuming non-informative censoring

$$p(Y = y, D = d) = [f_T(y)]^d [S_T(y)]^{1-d}$$

$$= [\lambda_T(y)]^d [S_T(y)]$$

Nonparametric time-to-event distribution estimators

$$\hat{\lambda}_i = \frac{d_i}{r_i}$$

$$\hat{\kappa}_i = 1 - \hat{\lambda}_i = \frac{r_i - d_i}{r_i}$$

$$\hat{S}_{KM}(t) \stackrel{\text{def}}{=} \prod_{\{i: t_i \leq t\}} \hat{\kappa}_i$$

$$\hat{\Lambda}_{NA}(t) \stackrel{\text{def}}{=} \sum_{\{i: t_i \leq t\}} \hat{\lambda}_i$$

Proportional hazards model structure

Formula	Description
$\lambda(t \tilde{x}) = \lambda_0(t) \cdot \theta_\lambda(\tilde{x})$	Proportional hazards assumption
$\Lambda(t \tilde{x}) = \Lambda_0(t) \cdot \theta_\lambda(\tilde{x})$	Cumulative hazard factorization
$\eta(t \tilde{x}) = \eta_0(t) + \Delta\eta(\tilde{x})$	Log-hazard decomposition
$\Delta\eta(\tilde{x}) = \tilde{x} \cdot \tilde{\beta} = \beta_1 x_1 + \dots + \beta_p x_p$	Linear predictor
$\theta_\lambda(\tilde{x}) = \exp\{\Delta\eta(\tilde{x})\}$	Hazard multiplier
$\theta_\lambda(t \tilde{x} : \tilde{x}^*) \stackrel{\text{def}}{=} \frac{\lambda(t \tilde{x})}{\lambda(t \tilde{x}^*)}$	Hazard ratio definition
$\theta_\lambda(t \tilde{x} : \tilde{x}^*) = \exp\{\Delta\eta(\tilde{x}) - \Delta\eta(\tilde{x}^*)\} = \exp\{(\tilde{x} - \tilde{x}^*) \cdot \tilde{\beta}\}$	Hazard ratio formula

Proportional hazards model partial likelihood formula:

$$\mathcal{L}_i^* = \frac{\theta_\lambda(\tilde{x}_i)}{\sum_{k \in R(t_i)} \theta_\lambda(\tilde{x}_k)}$$

$$\mathcal{L}^* = \prod_{\{i: d_i=1\}} \mathcal{L}_i^*$$

Proportional hazards model baseline cumulative hazard estimator:

$$\hat{\Lambda}_0(t) = \sum_{t_i < t} \frac{d_i}{\sum_{k \in R(t_i)} \theta_\lambda(\tilde{x}_k)}$$