

Advanced Regression Modelling ARM – ZF

- Inferenz:** Gibt es signifikanten Zusammenhang? Wie sicher sind wir mit geschätzten Koeffizienten?
- Modellwahl:** Welche Variablen haben grossen Anteil an Erklärung der Zielgrösse?
- Vorhersagen:** Vorhersagen mit Genauigkeit?
- Goodness-of-Fit:** Beschreibt Modell Daten adäquat?
- Residuenanalyse:** Stimmen Modellannahmen?

Zielgrösse Y (ihre Verteilung)	Modell
Stetig mit konstanter Varianz	Multiple lineare Regression
Binär	Logistische Regression
Anzahl	Poisson Regression
Stetig mit proportionaler Varianz	Gamma Regression
Zensierte Wartezeit stetig	Weibull-/Cox-Regression
Zensierte Wartezeit diskret	Grupp. prop. Hazardmodell

Generalisierte lineare und additive Modelle (GLM, GAM) → `par(mfrow = c(2, 3)); plot(Y ~., data = data)`

Logistische Regressionsmodell

- Zielvariable Y ist binär (0/1) und $\pi =$ Wahrscheinlichkeit, dass der Wert 1 angenommen wird ($\pi = P(Y = 1)$) ⇔ D.h.: Modellieren Y mittels Bernoulli-Verteilung mit unbekanntem Parameter π . π wird aus Daten geschätzt.
- Idee: bedingte Erwartungswert Zielgrösse, $E(Y_i|x_i) = \pi(x_i)$ soll erkl. Var. x_i abhängen. $E(Y_i|x_i) = P(Y_i = 1|x_i)$. Ansatz $\pi(x_i) = \beta_0 + \beta_1 x_i$ funktioniert nicht, da $\pi(x_i)$ Werte [0, 1] muss. ⇒ brauchen adäquate Transformation auf [0, 1]
- Suche nicht-linearen Funktion G , welche Wertebereich berücksichtigt: $\pi(x_i) = G(\beta_0 + \beta_1 x_i)$ bzw. $G^{-1}(\pi(x_i)) = \beta_0 + \beta_1 x_i$
- Wie G aussehen: Funktion G meistens S-Form. Beliebte G sind unten ersichtlich. $\eta = g(\pi)$ ist der lineare Prädiktor.

	$G(\eta)$	$g(\pi)$	Name
Logistische Verteilung	$\exp(\eta)/(1 + \exp(\eta))$	$\log(\pi/(1 - \pi))$	Logit-Modell; log Regression
Normal-Verteilung	$\Phi(\eta)$	$\Phi^{-1}(\pi)$	Probit-Modell
Extremwert-Verteilung	$1 - \exp(-\exp(\eta))$	$\log(-\log(1 - \pi))$	Komplementäres Log-Log-Modell

- Link-Funktion $g(\cdot)$ ist identisch zur inversen kumulativen Verteilungsfunktion $G^{-1}(\cdot)$.
- Beiden ersten Modelle sind theoretisch sehr schwer auseinander zu halten ausser in ihren äussersten Extremen
- Dritte Link-Funktion **nicht symmetrisch** um $\pi = 0.5$ ⇒ **entscheidend**, was 1 (= **Erfolg**) und 0 (= **Misserfolg**) in **ZV** ist.
- Logistische** Verteilung (Logit-Modell) wird i. A. **bevorzugt**, weil es eine **einfache Interpretation** der Parameter erlaubt.
- Kann anstelle Ereigniseintretenswahr' **Wettverhältnis (odds)** für (Nicht-)Erfolg: $\log(odds) = \log(\frac{\pi}{1-\pi}) = \beta_0 + \sum_{j=1}^m \beta_j x_j^{(j)}$
- Somit haben Koeffizienten β_j gleiche Bedeutung bei Beschreibung der log-odds wie bei der multiplen lin. Regression.

Die logistische Regression (**Logit-Modell**) ist durch folgende drei Elemente bestimmt:

- ZV Y_i Bernoulli mit Erfolgswahrscheinlichkeit π_i .
- Lineare Prädiktor $\eta_i = \beta_0 + \beta_1 x_i^{(1)} + \dots + \beta_m x_i^{(m)}$
- Lineare Prädiktor η_i ist mit Erfolgsparameter über logistische Funktion verlinkt: $\eta_i = \log(\frac{\pi_i}{1-\pi_i}) \Leftrightarrow \pi = G(\eta) = \frac{\exp(\eta)}{1 + \exp(\eta)}$

Parameterschätzung: Um geeignete Schätzung herzuleiten ⇒ Maximum Likelihood. Verfahren wird für Fall erläutert, dass Zielvariable unabhängig und binomial verteilt, verwenden logistische Link und nur eine erklärende Variable.

- Um Maximum-Likelihood-Schätzung von β_0, β_1 zu erhalten, nehmen wir partiellen Ableitungen von $l(\beta|y)$ bezüglich β_0 und β_1 durch Anwendung der Kettenregel. ML-Schätzung erhalten wir durch Gleichung $\sum_k (y_k^* - m_k \hat{\pi}_k) x_k^{(j)} = 0, j = 0, 1$
- Da $\pi_k = G(x_k^T \hat{\beta})$ eine nichtlineare Funktion der unbekannt Parameter $\hat{\beta}$ ist, ist das **Gleichungssystem nicht linear** → iterativ Algorithmen (Gradientenverfahren) wie z.B. der IRLS-Algorithmus müssen verwendet werden.
- Kommt **nicht darauf an, ob wir binären Daten oder gruppierten** für Anpassung verwenden. Schätzwerte sind gleich
- Anpassung mit binären Zielvariablen: `MAC.glm1 <- glm(Y ~ ac, family = binomial, MAC)` #Bei Default wird Logit-Modell
- Anpassung mit gruppierten Daten: `MAC.glmG1 <- glm(cbind(Y, m-Y) ~ ac, family = binomial, data = MAC1)` #angepasst
- Achtung:** Anpassung ist Resultat nichtlinearen Gleichungssystems ⇒ kann vorkommen, dass folgende erscheint: **Warning message: glm.fit: algorithm did not converge.** Algorithmus hat nicht konvergiert ⇒ **kann Schätzwerten nicht trauen**
- Leider nicht möglich Konvergenz Algorithmus sicherzustellen, ausser max. Anzahl Schritte erhöhen `glm(maxit = 1000)`.

Binäre Regression – Besonderheiten

- Punkte liegen auf zwei Kurven wegen binären Werten. Hat nichts mit Verletzungen Modellannahmen zu tun
- fast unmöglich, Ausreisser zu erkennen • Situation kann durch Datenaggregation erheblich verbessert werden
- Bei binären ZV ist Normal-QQ-Plot oft bedeutungslos. Selbst wenn Residuen nicht normalverteilt, sieht Plot gut aus.
- Leverage-Werte oft sehr klein. Können nicht Faustregel beurteilen, da zu viele Beob. verglichen Anz. unbek. Parameter
- Beispiel:** Geschätzte Modellgleichung für Erwartungswert binomial Logit-Modell mit einer erkl. Variable **weight**:
- $\hat{\eta}_i = -4.561 + 0.0051 \cdot \text{weight}_i \rightarrow \hat{\pi}_i = \hat{\mu}_i = \frac{\exp(\hat{\eta}_i)}{1 + \exp(\hat{\eta}_i)}$, Werte $-4.561, 0.0051$ sind β_0, β_1 Koeffizient aus Summary Output

Koeffizient log. Regression = **Änderung log. Wettverhältnis** → $\exp(0.0051) - 1 \approx 0.005$ ← Wenn Erhöhung **weight** um 1 Einheit, führt dazu dass das Log. Odds Wettverhältnis für zutreffen bspw. Infekt (Wert 1) sich um 0.5% ($0.005 \cdot 100$) erhöht.

```
sunflowerplot(x = birth$weight, y = birth$Y/birth$m, number = birth$m) | Verteilungsannahme: (Y_i|x_i) → (infect_i|wcc_i)
x = erkl. Variable, y = Y Anz. zutreffende Argumente, m = Total Anz. | ZV unabh. bin/bern: Y ~ Bin(m=6, pi = P(infect_i = 1))
x <- seq(min(birth$weight), max(birth$weight), length=50) #Prognose | E(Y_i) = m_i pi_i → E(pi_i) = pi_i = mu_i | eta_i = log(pi_i/(1-pi_i)) = beta_0 + beta_1 * wcc_i
mu.p <- predict(birth.glm1, newdata = data.frame(weight = x), | Ab welche Wert erkl. Var. wcc_i haben 80% Erfolg
type = "response"); | lines(x, mu.p, col = 'blue', lwd = 4, lty = 6) | wcc_i = 1/beta_1 * [log(pi_i/(1-pi_i)) - beta_0], pi_i = 0.8, beta_0 = summary
```

Asymptotisch Normalverteilt (Wald Asymptotik): Theorie ML-Schätzer: Jeder ML-Schätzer asymptotisch normalverteilt.
• asymptotisch im Sinne Zentralen Grenzwertsatzes. D.h. Approximation umso besser, je grösser Anzahl Beobachtungen.
• Können deshalb Summary-Output gleich interpretieren wie bei Kleinsten-Quadrate. Testresultate, VI gelten **nur approx**.

Das umfassende Modell – Die einfache Exponentialfamilie

Sowohl das lineare wie auch das logistische Regressionsmodell haben folgende Gemeinsamkeiten:

- Erklärende Variablen $x^{(1)}, x^{(2)}, \dots, x^{(m)}$ • eine Link-Funktion g , sodass $g(\mu) = x^T \underline{\beta}$
- eine Zielvariable Y mit Erwartungswert $E(Y) = \mu$ • eine Verteilung für Variabilität in Zielvariablen Y .
- Klassische lin. Reg. Identität $g(\mu) = \mu$ als Link und Normalverteilung mit $E(Y_i) = \mu_i$ und konstanter Varianz σ^2 verwendet.
- Zufallsvariable Y mit Erwartungswert $\mu = E(Y)$ und Dispersion ϕ hat Verteilung in einfachen Exponentialfamilie, falls
- Dichte oder Wahrscheinlichkeitsfunktion von Y geschrieben werden kann als $f(y_i; \mu_i, \phi) = \exp(\frac{y_i b(\mu_i) - c(\mu_i)}{\phi/w_i} + d(y_i; \phi, w_i))$
- $b(\cdot), c(\cdot)$ spez. Verteilung aus Exp.familie • $d(\cdot)$ Normalisierung Dichte/Wahrscheinlichkeitsfunktion auf 1
- Dispersionsparameter ϕ Teil der Varianz von Y
- «Variable» w_i ist für jede Beobachtung i eine bekannte Zahl. w_i kann unter den versch. Beobachtungen unterschiedlich sein. Interpretation als Gewicht. Wird in diesem Modulteil verwendet, um Binomialverteilung mit $m_k > 1$ abzudecken.
- Kann zeigen, dass für Erwartungswert von Y gilt $\mu_i = E(Y_i) = \frac{c'(\mu_i)}{b'(\mu_i)}$, wobei $c'(\mu_i), b'(\mu_i)$ Ableitung nach μ sind
- Varianz ist Produkt positiven Funktion $V(\mu_i) = \frac{1}{b''(\mu_i)}$ von Erwartungswert, Dispersion ϕ und Gewicht w_i : $var(Y_i) = \frac{\phi}{w_i} V(\mu_i)$

Gamma-Verteilung (Verallgemeinerung Exponentialverteilung)

- Zeit zwischen zwei konsekutiven Ereignissen, die Poisson verteilt sind, ist exponentialverteilt.
- beide sind geeignet, um Betragsgrössen wie (Überlebens-)zeiten, Schadenshöhen, Verluste oder Kosten zu modellieren
- Dichte gammadverteilten ZV Y_i ist $f(y, \alpha, \beta) = e^{-\beta y} \cdot y^{\alpha-1} \cdot \frac{\beta^\alpha}{\Gamma(\alpha)}$ mit $\alpha, \beta > 0$ auf Definitionsbereich $y \in [0, \infty)$
- Der Parameter α ist **Formparameter** und β ist die **Rate** ($1/\beta$ ein **Skalenparameter**).
- $\Gamma(\cdot)$ ist die Gamma-Funktion, eine Verallgemeinerung der Fakultätsfunktion auf die reellen Zahlen.
- $\alpha = 1, \beta = \lambda \rightarrow$ **Exp.verteilung** mit $\lambda: T \sim \text{Exp}(\lambda)$ • $\alpha = \frac{m}{2}, \beta = \frac{1}{2} \rightarrow$ **χ^2 -Verteilung** mit m Freiheitsgraden.
- m positive ganze Zahl, dann wird Gamma-Verteilung mit $\alpha = m$ und $\beta = \lambda$ auch **Erlang-Verteilung** genannt.
- Summe von zwei exponentialverteilten Z.V. ist gammadverteilt.** Allgemein: Falls $T_i, i = 1, 2, \dots, m$ unabhängig $\sim \text{Exp}(\lambda)$ dann ist Summe $T_1 + T_2 + \dots + T_m \sim \Gamma(\alpha = m, \beta = \lambda)$. Wenn wir annehmen, dass Wartezeit in einer Warteschlange unabhängig exponentialverteilt mit Parameter λ ist, dann ist Zeit zur Bedienung von m Kunden gammadverteilt.

Zusammenstellung	Verteilung	Wertebereich Y	$E(Y) = \mu$	$var(Y)$	$b(\mu)$	$V(\mu)$	ϕ	w
wichtigsten Grössen für einige der populärsten Verteilungen aus der einfachen Exponential-Familie.	Gauss (μ, σ^2)	$(-\infty, +\infty)$	μ	σ^2	μ	1	σ^2	1
	Binomial (m, π)	$(0, 1, \dots, m)$	π	$\frac{\pi(1-\pi)}{m}$	$\log(\frac{\mu}{1-\mu})$	$\mu(1-\mu)$	1	m
	Poisson (λ)	0, 1, 2, ...	λ	λ	$\log(\mu)$	μ	1	1
	Gamma (α, β)	$(0, +\infty)$	$\frac{\alpha}{\beta}$	$\frac{\alpha}{\beta^2}$	$-\frac{1}{\mu}$	μ^2	$\frac{1}{\alpha}$	1
→ →	Invers Gauss	$(0, +\infty)$	μ	$\frac{\mu^3}{\lambda}$	$-\frac{1}{\mu^2}$	μ^3	$\frac{1}{\lambda}$	1

Das generalisierte lineare Modell – GLM

- 1 Verteilungselement:** Verteilung Zielvariable Y_i , gegeben erklärenden Variablen x_i , ist Mitglied der einfachen Exponentialfamilie mit Erwartungswert $E(Y_i|x_i) = \mu_i$ und Dispersion ϕ . Zielvariable $Y_i, i = 1, \dots, n$, ist unabhängig verteilt.
- 2 Strukturelement:** Erwartungswert μ_i wird zum linearen Prädiktor $\eta_i := x_i^T \underline{\beta}$ der erklärenden Variablen mittels einer oft nichtlinearen Funktion $g(\cdot)$ in Beziehung gebracht: $g(\mu_i) = \eta_i = x_i^T \underline{\beta}$. Funktion $g(\cdot)$ wird Link-Funktion genannt.

Link-Funktion

- Ihre **Wahl** ist zu einem **gewissen Teil recht willkürlich**, sofern sie folgende grundlegende Eigenschaft besitzt:
 - Inverse Link-Funktion $g^{-1}(\cdot)$ soll Werte linearen Prädiktor $\eta = x^T \underline{\beta}$ auf Raum der möglichen Werte abbilden, die Erwartungswert von Y annehmen kann. Funktion $g^{-1}(\cdot)$ **muss unmögliche Werte vermeiden**. Link-Funktionen:
 - $g(\mu) = \mu$, wenn $\mu = E(Y)$ keinen Einschränkungen unterliegt, • $g(\mu) = \log(\mu)$, wenn $\mu = E(Y) > 0$ sein muss,
 - $g(\mu) = \text{logit}(\mu) := \log(\mu/(1-\mu))$, falls $\mu = E(Y)$ zwischen 0 und 1 liegen muss.

Kanonische Link-Funktion

- Für jede Verteilung aus einfachen Exponentialfamilie gibt spezielle **Link-Funktion**, die im gewissen Sinn **natürlich** ist.
- Führen Erwartungswert μ in kanonischen Parameter $\theta := b(\mu)$ über, bzw. wird angenommen, dass Effekte linearen Prädiktors linear auf kanonischen Parameter wirken, $\theta = \eta$. Diese Link-Funktionen nennt man kanonische Link-Funktionen.
- Also entspricht $b(\mu)$ gerade dem kanonischen Link für die jeweilige Verteilung. Siehe Tabelle oben Spalte $b(\mu)$
- ACHTUNG: kanonische Link für Gammadverteilung erfüllt Anforderungen nicht.** Genau Anpassungsergebnis überprüfen!

Link-Funktionen, die in R implementiert sind	binomial	gaussian	Gamma	inverse.gauss	poisson	quasi
Kanonische Link wird mit D bezeichnet (Spaltenname entspricht Verteilung (= family) in R, Zeilenname dem Namen der Link-Funktion	D					
• z.B. <code>glm(Y~X, family=poisson(link = identity))</code>	●				●	●
Poisson-Verteilung/Regression		D	●	●	●	●
• Bei Poisson-Verteilung ist Erwartungswert gleich Varianz → Normalverteilung kann nicht angenommen werden, da mit höheren $E(X) = 1/\mu^2$ dann auch eine höhere Varianz vorhanden ist. <code>sqr</code>	●	●	●	●	●	●
• Bei Schätzungen sollte man ganze natürliche Zahlen nehmen, da die Zielgrösse immer ganze Zahlen (Anzahlen) hat.			D		●	●

Prüfen Modellgüte von GLMs: Sollten Modellannahmen überprüfen bevor GLM als aussagekräftig akzeptieren und es verwenden!

- Überprüfung beruhen in erster Linie auf angepassten Werten und Residuen. **Ziel Residuenanalyse? Aufdecken systematische Abweichung Daten vom Modell und Identifikation ungenügend beschriebener Bereiche / Ausreisser**
- 4 Definition für GLM Residuum, die alle im Fall Normalverteilung identisch mit uns bekannten Definition Residuums sind.
- Einfachste und naheliegendste Definition Residuums ist **Response Residuen**: $R_i := Y_i - \hat{\mu}_i$, R_i sind i.A. nicht normalverteilt und ihre Varianz ist in erster Näherung $\phi V(\mu_i)(1 - H_{ii})/W_{ii}$, Varianzfunktion $V(\mu_i)$ ist i.A. nicht konstant
- Somit Definition sinnvoll **Pearson-Residuen**: $R_i^{(p)} := R_i \cdot w_i / \sqrt{V(\hat{\mu}_i)} \rightarrow$ ungefähr konstant Varianz, i.A. nicht normalverteilt.
- Basierend auf IRLS: **Working Residuen**: $R_i^{(w)} := R_i g'(\hat{\mu}_i)$. \rightarrow gewichtete Residuen; $\sqrt{w_i} R_i^{(w)} = \frac{w_i}{\sqrt{\phi(\hat{\mu}_i)g'(\hat{\mu}_i)^2}} (Y_i - \hat{\mu}_i) g'(\hat{\mu}_i) \equiv R_i^{(p)}$
- Eigenschaft klassisch Residuen: **Unterschiede zwischen beobachteten und angepassten Responsewerten messen**
- Summe quadrierten Residuen zur Schätzung σ^2 verwenden und Unterschiede Summen quadrierten Residuen versch. Modellen zu Modellvergleich einsetzen $\rightarrow R_i^{(w)} := \text{sign}(y_i - \hat{\mu}_i) \sqrt{d_i}$, $d_i =$ Summanden Residuen-Devianz.
- Wie Erfahrung zeigt, sind Unterschiede zwischen Devianz- und Pearson-Residuen üblicherweise sehr klein.

Tukey-Anscombe-Diagramm (Residuen-Plot): Pearson-Residuen $R_i^{(p)}$ gegen angepassten lineare Prädiktorwerte $\hat{\eta}_i$

- Wird verwendet, um auf **nicht erwartete Nichtlinearitäten** im linearen Prädiktor (oder allenfalls auf **unpassende Link-Funktionen**) und nebenbei auf **Ausreisser** aufmerksam zu werden.
- Abweichungen sichtbar machen, Glätter hineinlegen, der bei passendem Modell um horizontale Null-Achse schwankt
- Beurteilung kann herausfordern, da bei nicht normalverteilten Zielfunktionen irritierende Artefakte auftreten können.

Scale-Location-Diagramm: Wurzel absoluten Werte standardisierten Residuen gegen angepassten linearen Prädiktor.

- Untersucht **konstante Varianz** entlang angepassten Werte und relevanten Strukturen mit Glätter herausfiltern.

Half-Normal QQ-Plot: Indirekt zu prüfen, ob Verteilungsannahme zutrifft

- Falls **Verteilungsannahme** für Daten zutrifft, müssen Devianz-Residuen genähert normalverteilt \rightarrow um Gerade streuen
- Fokus Abweichungen (Ausreisser rechts im Plot). Deshalb wird half-normal QQ-Plots verwendet.
- Half-normal QQ-Plot werden geordneten absoluten Devianz-Residuen gegen entsprechenden Quantile der Standardnormalverteilung aufgetragen, d.h. gegen die $\frac{n+i-1/2}{2 \cdot n+1/2}$ Quantile der Standardnormalverteilung ($i = 1, 2, \dots, n$)

- Devianz-Residuen nur genähert normverteilt. Möglich, trotz Verteilungsannahme «korrekt», keine lin. Struktur sichtbar ist
- Darstellung nützlich um festzustellen, ob allfällige Overdispersion auf kleine Anzahl Ausreisser zurückzuführen ist.
- Bei binären ZV ist Normal-QQ-Plot oft bedeutungslos. Selbst wenn Residuen nicht normalverteilt, sieht Plot gut aus.

• `source("ARM/RFn_Plot-glmSim.R"); par(mfrow = c(2, 4)); plot(turb.glm); plot.glmSim(turb.glm, SEED = 184)`

Bootstap-Simulation: Idee: Erzeugen Beobachtungen y_i^* gem. Modell unter Verwendung neuer Zufallszahlen. Erzeugt n Zufallszahlen y_i^* entsprechend der verwendeten Verteilungsfamilie mit Erwartungswert $\hat{\mu}_i$ und Streuung $\hat{\phi}$.

- Berechnen GLM-Anpassung mit erkl. Var. aus Datensatz und neu generierten Responsewerten y_i^* . Anschliessend wird glatte Linie TA-Diagramm bestimmt und graue Linie eingezeichnet. Wiederholen diese beiden Schritte $n_{rep} = 19$ mal
- Gibt Fälle, wo rote Glätter starke nichtlineare Strukturen zeigt, aber noch innerhalb grauen Bootstrap-Simulationen liegt.
- Andererseits überdecken Bootstrap-Simulationen jedoch die horizontale Nulllinie nicht! Was heisst das jetzt?
- Können daraus schliessen, dass es **keine** Evidenz gibt, dass Missspezifikation im linearen Prädiktor vorliegt.

Residual gegen Leverage — zu einflussreiche Beobachtungen → Ausreisser und/oder Hebelpunkt

- Messen Hebelwirkung einer Beobachtung mittels der Diagonalelemente H_{ii} der GLM-Hutmatrix
- Beobachtung wird als Leverage Punkt bezeichnet, wenn $H_{ii} > 2p/n$ ist, p Anzahl der Koeffizienten ist.
- Cook's distance misst wieder Einfluss einzelner Beobachtungen und ist definiert als $d_i^{\odot} := \frac{R_i^{\odot}}{p \phi V(\hat{\mu}_i)(1-H_{ii})} \frac{H_{ii}}{1-H_{ii}}$, $1 \leq i \leq n$
- Faustregel**: Beobachtungen mit einer Cook's Distance d_i^{\odot} **grösser als 1** gelten als zu einflussreich

(Zeitliche) Korrelation in den Residuen: Residuen z. B. gegen Reihenfolge oder räumliche Anordnung Versuchseinheiten
Alternative: Streudiagramm Paare $(R_i^{(p)}, R_{i+1}^{(p)})$ `plot(fit, type = "h")` #Falls Beob. unab. streuen Punkte wild (weiss Rauschen)

Partielle Residuen-Plots (term plot). Ziel: Entdecken Nichtlinearitäten in erklärenden Grössen

- Auf x-Achse erkl. Variable $x_i^{(k)}$, auf y-Achse Pseudo-Zielvariable minus «Effekt» lin. Prädiktors ohne Komponente $x_i^{(k)}$:
- $Z_i - (\hat{x}_i^T \hat{\beta} - \beta_k x_i^{(k)}) = R_i^{(w)} + \hat{\beta}_k x_i^{(k)} \rightarrow$ Falls Modell richtig, sollten dargestellten Punkte um Gerade streuen.

• **Residual Plot**: Da der rote Glätter klar inner-/ausserhalb stochastischen Fluktuation (graue Spaghetti) liegt, gibt es keine/starke Evidenz, dass Einfluss der erklärenden Variablen auf die Erwartungswerte falsch spezifiziert wurde.

• **Half-Normal QQ-Plot**: Da schwarzen Punkte inner-/ausserhalb stochastischen Fluktuation (graue Punkte) liegen, gibt es keine/klare Evidenz, dass die Verteilung inadäquat spezifiziert ist und auch noch Hinweise auf zwei Ausreisser Wegen Overdispersion haben grauen Punkte eine andere Steigung als schwarzen.

• **Scale-Location-Plot**: Da rote Glätter vollständig inner-/ausserhalb der stochastischen Fluk. `a <- 2*length(coef(glm)) / nrow(data)`; `abline(v=a, col=3)` tuation liegt (graue Spaghetti), gibt es keine/klare Evidenz, dass Varianz inadäquat spezifiziert wurde. Wegen Overdispersion liegen grauen Spaghettis auch viel tiefer als rote Glätter! Y: Ausreisser, X: Hebelpunkte

- Residuen vs. Leverage**: Keine/Mehrere Beobachtungen haben Cook's Distance grösser als 1 und gibt/sind deshalb (keine) zu einflussreich(en) Beob. Bei den einen ist es bedingt durch Ausreisser, bei den anderen durch Hebelarme.
- Fazit**: Alle Plots zeigen (keine)klare Evidenz, dass Modell Daten (nicht) adäquat beschreiben kann (Inferenz wertlos)

Behandlung von Unzulänglichkeiten: J.W. Tukey nannte sie First Aid Transformations

- Logarithmus-Transformation für Konzentrationen/Beträge, • Arcus-Sinus-Wurzel $\hat{y} = \text{arcsin}(\frac{\bar{y}}{\sqrt{y}})$
- Wurzeltransformation für Zähldaten • **Nicht anwenden Zeitvariablen/zu wenig Streuung**
- Logit für Anteile (Prozentzahlen/100): $\hat{y} = \log(\frac{y+0.005}{1.01-y})$ • First-Aid/Log Faktor 10

Generalisierte additive Modelle (GAM)

- GAM erweitert GLM indem einige/alle linearen Terme linearen Prädiktor durch **glatte Funktionen** ersetzen: $\eta_i = \beta_0 + \sum_{k=1}^m x_i^{(k)} \beta_k \rightarrow \eta_i = \sum_{k=1}^m f_k(x_i^{(k)})$. Analyse grafisch mit geschätzten Funktionen: Glatte Kurve partiellen Residuen-Plot
- Schätzen unbekannte Funktionen **Backfitting-Algorithmus** \rightarrow Zielvariable Y_i durch Pseudo-ZV ersetzen und gewichtete Glättung wie IRLS-Algorithmus durchführen. GAM gut geeignet, da Transformation **nicht-parametrisch** schätzen.
- Hauptzweck für Einsatz von GAM ist **Aufdeckung geeigneter Transformationen** bei erklärenden Variablen im GLM.

`library(gam); baby.gam <- gam(Survival ~ lo(Weight) + lo(Age) + lo(pH), family = binomial, bf.maxit=100, data = baby)`
`par(mfrow=c(2, 3)) #2*3=Anz.Erkl.Var.; plot(baby.gam, se = T, residuals = T)` ↓ Bessere Darstellung für Variable EXTRP
`lo()` geht für Faktorvariablen nicht, diese regulär ohne `lo()` Modell nehmen. `plot(b.gam, residuals = T, terms = "lo(EXTRP)")`
Faustregel 1: Gerade passt (nicht) zwischen gestrichelte Vertrauensband, gem. Faustregel keine Transformation nötig.

Faustregel 2: summary(baby.gam) #Anova bei Nonparametric: wenn p significant, dann transformieren
`baby$Weight <- ifelse(baby$Weight < 1000, 0, baby$Weight-1000)` #bis zu einem gewissen Schwellenwert Prädiktor
`baby$tpH <- ifelse(baby$PH < 7.27, 0, baby$PH-7.27)` #beibehalten und danach diesen transformieren, Hockey Stick
`glm4 <- glm(Survival ~ Weight + tWeight + Age + Appar1 + Appar5 + pH + tpH, family = binomial, data = baby)`
`DaR$HRtr <- ifelse(DaR$HR < 12.5, DaR$HR, 12)` #Mehrfacher Hockeystick mit HR \rightarrow HRtr und HRhigh
`DaR$HRhigh <- as.integer(DaR$HR > 19.5); DaR$IF <- as.factor(cut(DaR$F, breaks=c(0,0.1,0.4,0.7,1)))` #kat. Var.
`DaR.QP <- glm(RDR ~ sVH + sPOP + IAR + HRtr + HRhigh + ff + IND, data = DaR, family = quasipoisson)`

Modelle für Raten (Rate Models – Poisson-Raten-Regression) \rightarrow **Term `log(ex_i)` heisst offset und ist bekannt.**

- Wichtige Annahme, dass **sowohl die Intensität** (oder Rate) des **Auftretens** von Ereignissen **und die Gelegenheit** (oder Exposition) zur Zählung **konstant** sind bei den gleichen erklärenden Variablen für alle entsprechenden Beobachtungen.
- Wenn jedoch z.B. Anzahl Autounfälle pro Jahr modelliert werden, kann sinnvoll sein, dass Unfallrisiko und damit Rate von gefahrenen Kilometern pro Jahr abhängt \rightarrow **Unfallrate pro gefahrener Kilometer und ZEIT modellieren**
- Rate ρ definiert als erwartete Zählungen pro Expositionseinheit. z.B. Unfälle pro 1'000 gefahrene Km, Festplattenausfälle pro Betriebsstunde. D.h., $\rho := \mu / ex$, ex ist **Exposition** und μ erwartete Anzahl innerhalb einer Exposition
- $\log(\rho_i) := \log(\frac{\mu_i}{ex_i}) = \beta_0 + \beta_1 x^{(1)} + \dots + \beta_m x^{(m)}$ | $\log(\mu_i) := \log(ex_i) + \beta_0 + \dots + \beta_m x^{(m)} \rightarrow \mu_i = ex_i \cdot \exp(\beta_0 + \beta_1 x^{(1)} + \dots)$

Ähnlichkeit zwischen Raten und Proportionen im Binomialmodell: (Anzahl Ereignisse) / (Gelegenheiten für Ereignisse)

- Die zwei Modellansätze sind jedoch unterschiedlich. Kann sie **unterscheiden**: Ist es denkbar, dass jede **Expositionseinheit mehr als ein Ereignis** erzeugt, oder muss **jede Expositionseinheit genau 0 oder 1 Ereignis hervorrufen**?
- Falls mehr als ein Ereignis in einer Expositionseinheit möglich \rightarrow Poisson-Raten besser, weil Annahmen Binom. verletzt
- Z.B. gibt es nichts, was verhindert, dass in 1'000 gefahrenen Kilometern mehrere Unfälle passieren können
- Andererseits wird jedes Baby eine Frühgeburt überleben oder nicht. Dafür sollte Binomialmodell verwendet werden.
- Parameter $\beta_0, \beta_1, \dots, \beta_m$ werden wie zuvor mit **Maximum-Likelihood-Methode** geschätzt.
- Ihre Interpretation erfolgt jedoch in Bezug auf die Rate und nicht auf den Erwartungswert.**
- β_0 ist Ereignisrate, wenn wir alle erkl. Var. auf 0 setzen, β_1 ist Änderung der Ereignisrate pro Einheit Anstieg von $x_i^{(1)}$
- Die Masseinheit für die Rate ist Zählungen pro Expositionseinheit (z.B. Unfälle pro 1'000 gefahrene Kilometer).

`glm1 <- glm(Y ~ offset(IService) + type + period + year, family = poisson, data = Ships)` #offset() für Raten-Modelle

Quasi-Modelle – Quasi-Poissonmodell

- Bei Zähldaten kann sein, dass Poisson-Modell die Daten nicht adäquat beschreibt, weil Dispersion zu gross ist.
- Alternative wo Erwartungswert-Varianz-Beziehung Form $\text{var}(Y_i) = \phi \cdot \mu_i$ für unbekannte Konstante ϕ . \leftarrow Fall $\phi > 1$
- Keine ML-Schätzung, da zu Schätzgleichung keine Likelihood definiert werden kann. Schätzungen asymptotisch normalverteilt. Können keine likelihood Methoden anwenden, wie Devianz-Test, Profiling, AIC. Bootstrap/GAM funktioniert
- Kovarianzmatrix zu Koeffizienten $\hat{\beta}$ wird geschätzt als $\hat{\phi}$ -fache der Kovarianzmatrix beim Poisson-Regressionsmodell.

`qPois <- glm(RDR ~ sVH + sPOP + IAR + HR + F + IND, family = quasipoisson, data = DaR1)`

GLM – Polytome Zielvariable \rightarrow für kategoriale Variablen

Kreuz-/Kontingenztafel und Mosaic-Plot \rightarrow `mosaicplot(table(wafer$method, wafer$quality))`

- Mosaic-Plots**: Fläche der Rechtecke proportional zur Anzahl Beobachtungen für Merkmalskombination in Stichprobe. Säulenbreite proportional zur relativen Häufigkeit der Stufen in ersten Variable. Höhe Rechtecke in jeder Säule proportional zur relativen Häufigkeit der Stufen aus zweiten Variablen. Anzahl Beobachtungen ist nicht ablesbar.

Multinomiale Logit-Modell für nominale Zielvariablen \rightarrow Zielvariable kategoriale Variable mit mehr als zwei Stufen

- Grösse von Interesse: **Wahrscheinlichkeit**, dass Y_i eine bestimmte **Ausprägung** ℓ annimmt, d.h. $\pi_{i\ell} = P(Y_i = \ell)$
- Ziel Beziehung zwischen $\pi_{i\ell}$ und erklärenden Variablen $x^{(k)}$, $k = 1, \dots, m$ finden, wobei garantiert sein soll, dass **Wahrscheinlichkeiten** zwischen **0 und 1** liegen und **Summe** $\pi_{i\ell}$ gleich 1 ist. Modellierung **ähnliche Idee wie Logit-Modell**.

• **Schritt 1**: Betrachten logarithmierte Wettverhältnis gegenüber Referenzkategorie (üblich = 1): $\log(\frac{\pi_{i\ell}^{(1)}}{\pi_{i1}^{(1)}}) = \log(\frac{\pi_{i\ell}^{(1)}}{\pi_{i1}^{(1)}})$

• **Schritt 2**: Setzen für jedes log. Wettverhältnis $\ell = 2, \dots, L$ sep. lin. Modell: $\log(\frac{\pi_{i\ell}^{(1)}}{\pi_{i1}^{(1)}}) = \eta_{i\ell}^{(1)} = \beta_0^{(\ell)} + \beta_1^{(\ell)} x_i^{(1)} + \dots + x_i^{(m)} \beta_m^{(\ell)}$

• Wahr modelliert: $\pi_{i\ell}^{(\ell)} = \exp(\eta_{i\ell}^{(\ell)}) \cdot \pi_{i1}^{(1)}$, $\ell > 1$ und $\pi_{i1}^{(1)} = \frac{1}{1 + \sum_{\ell=2}^L \exp(\eta_{i\ell}^{(\ell)})}$. Referenzkategorie sorgt **Summe** $\pi_{i\ell} = 1$

• Für jede Kategorie ℓ der Zielgrösse Y_i wird eigene Abhängigkeit der Wahrscheinlichkeit $\pi_{i\ell}^{(\ell)}$ von erkl. Variable geschätzt.

- Für zunehmendes $x^{(k)}$ bedeutet **positiver Koeffizient** $\beta_k^{(\ell)}$ **steigende Neigung** zur Kategorie ℓ im **Vergleich zur Referenzkategorie**. Für jede Kategorie **separates Modell** \rightarrow **viele Parameter**: Gleichungssystem $(L - 1)(m + 1)$ Parameter
- Sollten genügend Beobachtungen für Anzahl geschätzter Koeffizienten haben: (**min. 5-10 Datenpunkte pro Koeffizient!**)
- Wahl Referenzkategorie keine Auswirkung auf Modell** • Beobachtungen können gruppiert sein \rightarrow Kreuztabellen.

Nicht gruppierte Daten: library(nnet); library(MASS); waver.fitsm <- multinom(quality ~ method, data = wafer)
R nimmt jeweils erste Faktorstufe als Referenz. Diese ist dann im summary(waver.fitsm) auch nicht direkt ersichtlich.
Bei Gruppierung: wG.fitsm <- multinom(cbind(excellent, ok, bad) ~ method, data = waverG); summary(wG.fitsm)
Zielvariable als Matrix übergeben, Faktorstufen/Levels explizit angeben, hier bspw. excellent, etc. Outputs sind identisch!

Interpretation Koeffizient: Wert Koeffizient ist immer im Vergleich zur Referenzausprägung! Log-Odds von ok zu excellent
nehmen für Methode2 im Vergleich zu Methode1 um 0.734 ab. Objekte Methode2 haben im Vergleich zu Methode1 eine geringere Neigung ok als excellent zu sein
Inferenz: Ob erklärende Grösse Einfluss auf Zielgrösse hat, sollte nicht anhand Standardfehler bestimmen (da L Koeffizienten null sein müssten, wenn kein Einfluss da ist!) → hierarchischer Modellvergleich via Devianzen mit anova()
• testen, ob an allen Standorten gleiche Verteilung verwendet werden könnte (Homogenitätshypothese)
• drop1() geht für multinomiale Modelle nicht. Alternativ library(MASS); dropterm(waver.fitsm, test = "Chisq")
• Gleiche Resultat mit Devianztest: w1.fitsm <- multinom(quality ~ 1, data = wafer); anova(waver.fitsm, w1.fitsm)
• Modellvergleich zeigt, beide Modelle signifikant verschieden → Homogenitätshypothese wird auf 5% Niveau verworfen, da P-Wert 4.9e-13 kleiner als 5%. Saturierte Modell ist also notwendig, um Daten adäquat zu beschreiben. Folglich haben Methoden statistisch signifikanten Unterschied auf Qualität! (Methode / Qualität auf erkl. / ZV anpassen!)

Vorhersage (Prediction): Klassenzugehörigkeit vorhersagen: predict(waver.fitsm, newdata = w0, type = "class") ↓ besser
w0 <- data.frame(method = c("method1", "method2", "method3")); predict(waver.fitsm, newdata = w0, type = "probs")

Modell-Diagnostik: Gibt keine aussagekräftige Residuen Definition. Keine Prüfung, ob Modell Daten adäquat beschreibt
Modelle ordinale Zielvariable wafer\$quality <- factor(wafer\$quality, level = c("bad", "ok", "excellent"), ordered=T)
• Ordnung in Stufen der Zielvariable berücksichtigen. Ordnen Stufen in Zielvariable natürlich aufsteigende Ordnung.
• Modell, wie bei logistischen Regression mit Annahme latenten kontinuierlichen Variable T. In Realität kann nur Yi beobachtet werden, die sich als diskretisierte Version von Ti ergibt mit Schwellenwerten: alpha_0 = -inf < alpha_1 < ... < alpha_{L-1} = inf
• Schwellenwerte alpha müssen nicht äquidistant sein und sind i.d.R. nicht a priori bekannt → mit Daten schätzen
• Weniger Parameter zum Schätzen als bei multinom. Regr. Residuenanalyse: Keine, wie bei multinom Logit Modell
• Modell: tau_i = P(Yi <= l | xi_i) = P(Ti >= alpha_l | xi_i) = 1 - F(alpha_l - (beta_0 + beta_1 xi_i^(1) + ... + beta_m xi_i^(m))) / beta_0 ist unbestimmt, jedoch beta_0 = 0
library(MASS); w.polar <- polr(quality ~ method, data = wafer); summary(w.polar) #Faktoren mit ordered = T übergeben!

Statistische Wartezeitanalyse (Analysis of Time-to-Event Data, Event History Analysis, Failure Time Analysis)

- Werden zensierte Beobachtungen weggelassen/ignoriert, kann systematische Fehler in geschätzten Grössen geben. Was wollen wir mit der Analyse von Wartezeitdaten erreichen? Zwei wichtige Modellierungskonzepte:
o Wahrscheinlichkeit, dass Ereignis später als zum Zeitpunkt t eintritt → Überlebensfunktion
o Möglichkeit, dass ein Ereignis zu einem bestimmten Zeitpunkt t eintritt → Hazard-Rate
• (Warte-)zeit interessante Variable, modelliert T ≥ 0 reellwertige Zufallsvariable Zeit bis Ereigniseintreten. Zwei Fälle:
o T genau gemessen → T ist kontinuierlich o nur bestimmt Zeitpunkt erfasst (zweimal im Jahr) → T diskret
• Typ I Zensierung: Experimentstart gleichzeitig. Zeitpunkt t_x Abbruch. Komponenten ohne Ausfall sind zensiert.
• Typ II Zensierung: Gleich wie Typ I, jedoch Abbruch, sobald r der insgesamt n Objekte ausgefallen sind.
• Zufällige Zensierung: Oft gleichzeitiger Start unmöglich. Zensierung erfolgt, weil Patient weggezogen (Loss to Follow-up); Behandlung wegen Nebeneffekt oder sonst abgebrochen (Drop Out); oder Studie beendet (Study Termination).
• Beschreibung Wartezeitverteilung: stetige Verteilungsfun. F. Dazugeh. Dichtefun. f. → P(T ≤ t) = F(t) = ∫_0^t f(x) dx
• Survivorfunktion S (Zuverlässigkeitsfunktion): Wahrscheinlichkeit, dass Wartezeit mindestens t beträgt:
S(t) := P(T > t) = 1 - F(t) = ∫_t^inf f(x) dx S(t) monoton fallend, weil F monoton steigend ist; also S(0) = 1 und S(inf) = 0.
• h(t) = lim_{dt->0} (P(T ≤ t+dt | T ≥ t) - P(T ≤ t | T ≥ t)) / dt = f(t) / S(t) → Hazardrate (Ausfall-/Sterberate) ist keine Wahrscheinlichkeit. Tritt Ereignis nicht vor Zeitpunkt t ein, so informiert die Hazardrate über «weiteren Verlauf»:
o Ist konstant → Objekte «altern» & «verjüngen» nicht o Steigt Ausfallrate → Objekte «altern», Verschleiss
o Abfallende Ausfallrate (selten) → Ausfall, wenn er eintritt, dann meistens früh
• Hazardrate bed. Wahr' Ereigniseintritt in [a_{t-1}, a_t), Bedingung Erreichen Intervall: lambda_t := P(T ∈ [a_{t-1}, a_t) | T ≥ a_{t-1})

Schätzung Hazardrate und Survivorfunktion

Sterbetafel (life-time table) → Eintreten Ereignis beobachtet (z_i = 1) oder nicht beobachtet → zensiert (z_i = 0) ist
• Einfachste, gebräuchlichste Methode diskreten Wartezeiten. Einfluss erkl. Merkmale explizit nicht berücksichtigt.
• Zeitintervalle t o Anzahl Komponenten, die im Zeitintervall dem Risiko unterliegen, dass sie ausfallen können
o Anzahl Komponenten, die ausgefallen sind o Anzahl Komponenten, die zensierte Beobachtungszeit haben.
• n Gesamtzahl gefährdeten Objekte/Subjekte zu Beginn • d_t Anzahl Ereignisse, die t-ten Intervall [a_{t-1}, a_t) eintreffen
• w_t Anzahl Zensierung Intervall t [a_{t-1}, a_t) → erreichen t-te Intervall, aber weder Ereigniseintritt noch nächste Intervall
• n_t Anzahl gefährdeten Objekte, bei denen Ereignis im t-ten Intervall [a_{t-1}, a_t) eintreten könnte → Risikomenge «at risk»
• Schätzen aus Sterbetafeln: Ohne Zensierung lambda_hat_t = d_t / n_t, mit Zensierung (Wann gezählt?): Ende t-ten Intervall lambda_hat_t = d_t / n_t, Anfang t-ten Intervall lambda_hat_t = d_t / n_{t-wt/2} → übliche Schätzverfahren t lambda_hat_t = d_t / n_{t-wt/2}. Schätzung Survivorfunktion S_hat(t) = ∏_{s=1}^t (1 - lambda_hat_s)

LDT <- data.frame(Time = c(6.8, 17.6, 17.6, ...), status = c(1, 1, 0, ...)) h <- LDT\$Tdiscrete <- ceiling((LDT\$Time/24))
dt <- table(factor(h[LDT\$Tstatus == 1], levels = 1:6)) #6 = Anz. Intervalle; wt <- table(factor(h[LDT\$Tstatus == 0], levels = 1:6))
et <- c(0, cumsum(dt[-length(dt)] + wt[-length(wt)])); nt <- nrow(LDT) - et; hazard <- dt / (wt - wt/2); S <- cumprod(1 - hazard)
library(discSurv); RR <- read.table("Rossi.txt", ...); RR4 <- data.frame(week=(RR\$week-1)%/4 + 1, arrest=RR\$arrest)
h.lab1 <- paste("W", seq(1,49,by=4), sep=""); h.lab2 <- paste("W", seq(4,52,by=4), sep="") #auf je vier Wochen aggregieren
h.lab <- paste("I", h.lab1, "I", h.lab2, "I", sep="") # ↓ Spalte "dropouts" ≡ w_t; Spalte "atRisk" ≡ Nenner (n_t - w_t/2)
IT <- lifeTable(dataShort = RR4, timeColumn = "week", eventColumn = "arrest", intervalLimits = h.lab); IT\$Output
15.01.2025 4. Semester / 7. Semester Seite 7 von 12

Verteilung Survivorfunktions-Schätzung: Keine zensierte Beobachtung (w_t = 0) : S_hat_t = n_t / n = (n - d_1 - d_2 - ... - d_{t-1}) / n. Anz. Ereignis, die bis Zeitpunkt a_t noch nicht eingetroffen, ist binomialverteilt, N_t ~ B(n, S_t). Folglich gilt E(S_hat_t) = S_t, Var(S_hat_t) = S_t(1-S_t) / n

- Genähertes 95%-Vertrauensintervall: S_hat(t) ± 2 * sqrt(S_hat(t)(1-S_hat(t))/n), genähertes 95%-VI Hazardfunktion: lambda_hat_t ± 2 * sqrt(d_hat_t(1-d_hat_t)/n_t), lambda_hat_t = d_hat_t / n_t
• Falls zensierte Beobachtungen (w_t ≥ 0): Schätzer lambda_hat_t = d_hat_t / (n_t - w_t/2) ist nicht konsistenter Schätzer.
• Verzerrung ignoriert und Varianzschätzung analog Schätzer korrigiert: lambda_hat_t ± 2 * sqrt(d_hat_t(1-d_hat_t) / (n_t - w_t/2)), genähertes, konservativ 95%-VI
• Varianzschätzung Survivorfunktionssschätzung S_hat(t) oft Greenwoods Formel Var(S_hat(t)) = S_hat(t)^2 * sum_{k=1}^t (lambda_hat_k / ((1-lambda_hat_k)(n_k - w_k/2)))
• Geglättete Hazardrate: Schätzung Hazardfunktion nach Sterbetafel kann sehr volatil sein vor allem in einzelnen Zeitintervallen, in denen Risikomenge n_t klein wird. Dies ist vor allem bei grossen t der Fall.
• Da es sich bei gewissen Daten eigentlich um GLM handelt mit erkl. Var. «Nummer Zeitintervall», können GAM einsetzen, um geglättete Hazardrate zu bestimmen. Glätten mit Hilfe gam hilft bei vielen Zeitintervallen weiter.

Bestimmung wöchentlichen Hazardfunktion mit Sterbetafelmethode und mit geglätteter Sterbetafelmethode GAM
IT <- lifeTable(dataShort = RR, timeColumn = "week", eventColumn = "arrest"); h <- IT\$Output; h\$fit <- 1:nrow(h)
library(mgcv); h.gam <- gam(cbind(events, round(atRisk) - events) ~ s(tl), family = binomial, gamma = 0.7, data = h)
h.gam.p <- predict(h.gam, type = "response", se = TRUE); h <- data.frame(h, sfit = h.gam.p\$sfit, se.sfit = h.gam.p\$se.sfit);
plot(h\$fit, h\$se.hazard, type = "l"); lines(h\$fit, h\$se.hazard, lty = 2); lines(h\$fit, h\$se.hazard, lty = 2); lines(h\$fit, h\$se.hazard, lty = 2);
lines(h\$fit, h\$se.sfit, col="red"); lines(h\$fit, h\$se.sfit, col="red", lty=2) #GAM-Fit, noch mit -2*
Vorhersage auf Link-Ebene (besser): gam.link <- predict(h.gam, type = "link", se = T); h1.fit <- gam.link\$sfit
h1.se <- gam.link\$se.fit; h1 <- data.frame(fit=h1.fit, lci=h1.se, uci=h1.fit + 2*h1.se, uci=h1.fit - 2*h1.se); h1.trans <- exp(h1)/(1+exp(h1)) #Rücktrans lines(h\$fit, h1.trans\$sfit); lines(h\$fit, h1.trans\$lci); lines(h\$fit, h1.trans\$uci)

Kaplan-Meier-Schätzung Survivorfunktion für stetige Wartezeit

- Keine Zensierung: Schätzung S(t) wie emp. Verteilungsfun. (stetigen) Messgrösse: S_hat_n(t) = #Beobachtungen > t / n = #t_i >= t / n
• S_hat_n(t) konsistent Schätzung unzensiert. n * S_hat_n(t) ist B(n, S_t)-verteilt. ZGWS approx. VI für S(t_0): S_hat_n(t_0) ± 2 * sqrt(S_hat_n(t_0)(1-S_hat_n(t_0))/n)
• Zensierten Beobachtung: (t_i^*, z_i) t_i^* enthalten Wartezeiten und z_i hält fest, ob «+» zensiert (z_i = 0) oder nicht (z_i = 1)
• Positive Halbchase mit disjunkten Intervallen I_k = [a_{k-1}, a_k), wobei a_0 = 0 ist. Survivorfunktion wird geschätzt mit S_hat(t) = ∏_{k=1}^t (1 - lambda_hat_k), lambda_hat_k = Wahr', Ereignis im Intervall zu haben. lambda_hat_k = 1 falls kein Ereignis im Zeitintervall beobachtet wurde.
• Stetigen Fall Intervall ganz klein gemacht, sodass höchstens eine Ereigniszeit vorkommt: S_hat_KM(t) = ∏_{i:t_i^* <= t} (1 - d_i/n_i) = ∏_{i:t_i^* <= t} (n_i - d_i/n_i)
• KM-Schätzer ist wie emp. Schätzer stückweise horizontal und hat nur an Stellen t_i^* Sprünge (von unzensierten Beob.)
• KM ist asymptotisch normalverteilt: S_hat_KM(t) ~ N(S(t), sigma_hat_KM^2(t)). Asymptotische Varianz wir mit Greenwood's Formel geschätzt: sigma_hat_KM^2(t) ≡ se_G^2(t) := S_hat_KM^2(t) * sum_{i:t_i^* <= t} (d_i / (n_i(n_i - d_i))), wobei t_i^*(k) ≤ t < t_i^*(k+1). Genähertes 95%-VI für S(t) bei festem t:
S_hat_KM(t) ± 2 * q_{0.975} * se_G(t). Ansatz kann dazu führen, dass Vertrauensintervall nicht innerhalb [0, 1] zu liegen kommt.
library(survival); KM <- survfit(Surv(Time, status) ~ 1, data = LDT); summary(KM) #status = zensiert mit 0, beobachtet = 1
plot(KM, conf.int = T, las = 1, mark.time = TRUE, xlab = "Lebensdauer", ylab = "Geschätzte Survivorfunktion S")

Parametrische Modelle für kontinuierliche Wartezeiten

Exponential-Verteilung E(lambda) → Modellierung Ausfallzeiten Objekten, die nicht altern → gedächtnislose Verteilung!
• Dichtefunktion f(t) = lambda * e^{-lambda*t} • Erwartungswert E(X) = 1/lambda • Varianz var(X) = 1/lambda^2
• Survivor-Funktion Verweildauer T: S(t) = P(T > t) = e^{-lambda*t} • Verteilungsfunktion F(t) = 1 - e^{-lambda*t} mit lambda > 0 für t ≥ 0
• Hazard- (Ausfall-) Rate Exponential-Verteilung ist zu jedem Zeitpunkt t konstant gleich Parameter lambda, da h(t) = f(t)/S(t) = lambda
• Vorzeitige Erneuerung noch nicht ausgefallener Objekte bei exponentialverteilter Lebensdauern ist sinnlos!

Weibull-Verteilung W(lambda, alpha) → Wurde 1939 von W. Weibull Beschreibung Materialermüdungserscheinungen erfunden
• Verteilungsfunktion F(t) = 1 - exp(-(lambda*t)^alpha) mit alpha, lambda > 0 • Falls alpha = 1, erhalten Exponential-Verteilung
• Dichte f(t) = lambda(alpha*t)^{alpha-1} * exp(-(lambda*t)^alpha) • Falls alpha = 2, erhalten Rayleigh-Verteilung
• lambda ist Skalenparameter und alpha Formparameter • Survivorfunktion S(t) = exp(-(lambda*t)^alpha)
• Hazardrate h(t) = lambda(alpha*lambda*t)^{alpha-1}, sie ist monoton steigend falls alpha > 1, monoton fallend falls alpha < 1 und konstant falls alpha = 1

Weibull-Diagramm: komplementäre log-log-transformierte Survivorfunktion einer Weibullverteilung ist lineare Funktion in log(t) : log(-log(S(t))) = alpha(log(lambda) + log(t)) → Plot log(-log(S(t))) vs. log(t) Klärung Zeiten Weibull
• S_hat(t) wird mit Kaplan-Meier-Methode geschätzt. • Aufgetragen werden nur unzensierten Beobachtungen.
• Falls Daten von Weibull-Verteilung stammen, müssen Punkte bis auf zufällige Abweichungen auf einer Geraden liegen.

library(survival); source("addFn4SurvAnal.R"); KM <- survfit(Surv(Time, status) ~ 1, type = "kaplan-meier", data = LDT)
plot(KM, xlab = "Zeit in Stunden bis zum Ausfall", ylab = expression(hat(S)(KM)(t)), conf.int = TRUE)
Überprüfen Daten exp. verteilt: plot.weibull(Surv(Time, status) ~ 1, data = LDT, line = F, xlab = "log(Std)"); h <- (log(-log(KM\$surv)) - log(KM\$time)); (h <- h[is.finite(h)]); abline(a = mean(h), b = 1, lty = 1, col = 4) #log(lambda) von Hand geschätzt.

Überprüfung Verteilungsannahme: plot.weibull(Surv(Time, status) ~ 1, data = diab) #Weibull, mit scale=1 Exp Verteil.
plot(KM, fun = "dloglog", col = 2:3, mark.time = T, xlab = "log(Std)(KM)(t)"), ylab = expression(log(-log(hat(S)(KM)(t))))

Gumbel-Verteilung G(mu, sigma) → Lage-Skalen-Verteilung

- Ist Zeit Weibull-verteilt, so ist logarithmierte Zeit Gumbel verteilt: T ~ W(alpha, lambda) → Y = log(T) ~ G(mu = -log(lambda), sigma = 1/alpha)
• Lageparameter mu, Skalenparameter sigma auf ganzen reellen Achse mit Verteilungsfunktion F(y) = 1 - exp(-exp(y - mu/sigma))

QQ-Plot für Gumbel-Verteilung

- Tragen geordneten logarithmierten Wartezeiten (nur jene die beobachtet wurden) gegen entsprechenden Quantile Standard-Gumbel-Verteilung auf. Quantile Standard-Gumbel-Verteilung berechnen mit $F_G^{-1}(F_n(y_k)) = \log(-\log(1 - F_n(y_k)))$
- $F_n(y_k)$ ist empirische Verteilungsfunktion der Wartezeiten. Tragen im Streudiagramm $\log(t)$ gegen $\log(-\log(\hat{S}_{KM}))$ auf.
- Um zensierten Beobachtungen zu berücksichtigen wird $1 - F_n(y_k) = S_n(y_k)$ mit Kaplan-Meier-Methode geschätzt.
- Folglich ist QQ-Plot für Gumbel-Verteilung gerade dem an der Winkelhalbierenden gespiegelten Weibull-Diagramm.

Verteilung	$\mathbb{E}(T)$	$var(Y)$	Median	$h(t)$	$S(t)$
Weibull (α, λ)	$\frac{1}{\lambda} \Gamma(1 + \frac{1}{\alpha})$	$\frac{1}{\lambda^2} (\Gamma(1 + \frac{2}{\alpha}) - \Gamma^2(1 + \frac{1}{\alpha}))$	$\frac{1}{\lambda} \log(2)^{\frac{1}{\alpha}}$	$\lambda \alpha (\lambda t)^{\alpha-1}$	$\exp(-(\lambda t)^\alpha)$
Gumbel (μ, σ)	$\mu - \gamma \sigma$	$\frac{\pi^2}{6} \sigma^2$	$\mu + \sigma \log(\log(2))$	kompiziert	$\exp(-\exp(\frac{y-\mu}{\sigma}))$

Loglogistische Verteilung – Wartezeit T ist loglogistisch verteilt, falls $Y := \log(T)$ logistisch verteilt ist

- Dichte $f(t) = \lambda \cdot \alpha \cdot \frac{(\lambda t)^{\alpha-1}}{(1+(\lambda t)^\alpha)^2}$
- Survivorfunktion $S(t) = \frac{1}{1+(\lambda t)^\alpha}$
- Hazardrate $h(t) = \frac{\lambda \cdot \alpha \cdot (\lambda t)^{\alpha-1}}{1+(\lambda t)^\alpha}$
- Wettverhältnisse mit Survivorfunktion: $\frac{S(t)}{1-S(t)} = \frac{1}{(\lambda t)^\alpha}$ somit $\log(\frac{S(t)}{1-S(t)}) = -\alpha \cdot \log(\lambda) - \alpha \cdot \log(t)$

Inferenz am Beispiel der Exponential-Verteilung – Maximum-Likelihood Schätzung bei zensierten Beobachtungen

- Sei t_i beobachteten Zeiten, n_u Anzahl unzensierten Beobachtungen und z_i gibt an, ob **Eintreten Ereignis beobachtet wurde** ($z_i = 1$) oder **zensiert** ($z_i = 0$) ist. Bei zensierten Beobachtungen wird Beitrag Log-Likelihood mit $\log(S(t_i))$ erfasst
- $\ell(\lambda) = \sum_{i=1}^n z_i \log(f(t_i)) + (1 - z_i) \log(S(t_i)) = n_u \log(\lambda) - \lambda \sum_{i=1}^n t_i \rightarrow 0 = \frac{\partial \ell(\lambda)}{\partial \lambda} = \frac{n_u}{\lambda} - \sum_{i=1}^n t_i$, ML-Schätzer $\hat{\lambda} = \frac{n_u}{\sum_{i=1}^n t_i}$
- Wie unzensierten Schätzung Kehrwert Mittelwert, ausser statt Anzahl Beob $n \rightarrow$ Anzahl unzensierten Beob n_u geteilt.
- Schätzung $\hat{\lambda}$ ist asymptotisch normalverteilt mit $\mathbb{E}(\hat{\lambda}) = \lambda$, Varianz $var(\hat{\lambda}) = \frac{\lambda^2}{\mathbb{E}(n_u)}$. Genäherte 95%-VI Wald: $\hat{\lambda} \pm q_{0.975}^N \cdot \frac{\hat{\lambda}}{\sqrt{n_u}}$
- Approximation bei **kleinen Stichproben schlecht**, u.a., da asympt. Varianz vom unbekanntem Parameter abhängt.

Log-basiertes Vertrauensintervall:

- Abhilfe: Transformieren $\hat{\lambda}$ so, dass asympt. Varianz unabhängig vom unbekanntem Parameter ist: $\log(\hat{\lambda})$ betrachten
- Wird approximatives VI gerechnet und Endpunkte zurücktransformiert: $\hat{\lambda} \cdot \exp(\pm q_{0.975}^N \cdot \frac{1}{\sqrt{n_u}}) \rightarrow$ gute Näherung 95%-VI
- Aufgrund asympt. Resultat $\log(\hat{\lambda}) \sim N(\log(\lambda), \frac{1}{\mathbb{E}(n_u)})$ können genähertes 95 %-VI für $\log(\lambda)$ durch $\log(\hat{\lambda}) \pm q_{0.975}^N \cdot \frac{1}{\sqrt{n_u}}$

Parameter schätzen der Exponentialverteilung mit ML-Schätzer. Wie gross ist erwartete Überlebenszeit?

```
library(survival); LDT.E <- survreg(Surv(Time, status) ~ 1, data = LDT, dist = "exponential"); exp(-coef(LDT.E))
#(oben) Rücktransformation Schätzwert für λ; exp(-confint(LDT.E))[2:1] #Rücktrans. VI für λ zu Reihenfolge Intervall
Schätzwert erwartete Lebenszeit (= 1/λ) → exp(coef(LDT.E)); exp(confint(LDT.E)) #gleich, einfach ohne das Minus
Wahrscheinlichkeit, dass mehr als 26 Wochen vergehen: exp(-26*exp(-coef(LDT.E))) #Transformation aufgrund Gumbel
LDT.W <- survreg(Surv(Time, status) ~ 1, data = LDT, dist = "weibull") #Anpassung in survreg(...) erfolgen via Gumbel-
Modell! → Schätzungen für Exp./Weibull müssen rücktransformiert werden! summary(LDT.W) # ↓ P-Wert «log(scale)»
Log(scale) ist auf 5% Niveau (nicht) signifikant von 0 verschieden; d.h. scale ist auf 5% Niveau (nicht) signifikant von 1
verschieden. Da α = 1/scale folgt, dass α auf 5% Niveau (nicht) signifikant von 1 verschieden ist. Somit wird Nullhypothese,
dass Daten mit Exponential-Verteilung modelliert werden können, (nicht) verworfen. D.h. keinen Sinn frühzeitig Austausch
plot(LDT.KM, conf.int = F, col = 2, lty = 2, las = 1); h.tp <- seq(0, 3, 0.5); lines(h.tp, pexp(h.tp, rate = exp(-coef(LDT.E))),
lower.tail = F); lines(h.tp, pweibull(h.tp, scale = exp(coef(LDT.W))), shape = 1/LDT.W$scale, lower.tail = F), col = 3)
```

- KM: Weil grössten Lebensdauerdaten zensiert sind, können keine mittlere Lebenszeit schätzen.
- Mit Exponential-/Weibullmodell ist möglich erwartete Lebensdauer zu schätzen (siehe Tabelle oben Median Gumbel):
- **Bsp** (Werte Summary: Intercept 4.503, Scale fixed at 1) Median = $\exp(\mu + \sigma \log(\log(2))) = \exp(4.503 + 1 \cdot \log(\log(2)))$

Beschleunigte Ausfallzeit-Modelle (accelerated failure time models – AFT) → Einfluss erklärenden Variablen berücksichtigen

- Lineare Prädiktor (Linearkombination erkl. Var.) soll auf Lage logarithmierten Zielvariable einwirken: $\log(T_i) := Y_i = \beta_0 + x_i^T \beta + \sigma \cdot W_i$, wobei W entweder standard Gumbel, standard logistisch oder standard normalverteilt ist.
- Sei nun $W_i^* := \beta_0 + \sigma \cdot W_i$ (Referenzverteilung), kann obiges Modell geschrieben werden als $\log(T_i) := Y_i = x_i^T \beta + W_i^*$, wobei W_i^* Gumbel, loglogistisch oder normalverteilt mit jeweils Lageparameter $\mu = \beta_0$ und Skalensparameter σ ist.
- Rücktrans. W -, loglogistischen, Lognormalvert. $T_i^* := \exp(W_i^*) \rightarrow T_i = \exp(Y_i) = \exp(x_i^T \beta) \cdot \exp(W_i^*) = \exp(x_i^T \beta) \cdot T_i^*$
- Der exponentierte lineare Prädiktor $x_i^T \beta$ wirkt also multiplikativ auf eine «Referenz»-Wartezeitverteilung T_i^* ein.

Survivorfunktion von T: Beschleunigungsfaktor $\gamma := \exp(-x^T \beta)$, Survivorfunktion von T gegeben x lässt mit Survivorfunktion von T^* ausdrücken: $S(t|x) = S_0(t \cdot \gamma|x) \rightarrow$ Der exponentierte lineare Prädiktor ändert also Skala von der Zeit t.

- Erklärenden Variablen beschleunigen/verlangsamten Ausfall-/Wartezeit; d.h. Zeit bis Ausfall, abhängig $\gamma < 1, \gamma > 1$
- Gibt für AFT-Modelle verschiedene Schreibweisen. $\log(T_i) := Y_i = \beta_0 + x_i^T \beta + \sigma \cdot W_i$ heisst **log-lineare Darstellung**.

Hazardrate: Hazardrate $h_0(t)$ Zufallsvariablen T^* ist unbeeinflusst von erkl. Var. x und Parametervektor β

- Für Hazardrate $h(t|x)$ von $T = e^{x^T \beta} \cdot T^*$ gilt dann $h(t|x) = -\frac{d \log(S(t|x))}{dt} = h_0(t \cdot \gamma|x) \cdot \gamma$, mit $\gamma := e^{-x^T \beta}$
- Darstellung zeigt, dass **Beschleunigungsfaktor γ sowohl multiplikativ auf Hazardrate als multiplikativ auf Zeit einwirkt**.

Weibull-Regressionsmodell (gleiches Modell kann auch mit Gumbel-Verteilung beschrieben werden)

- Weibull beliebt, weil Survivorfunktion $S(t|x) = e^{-(\lambda t e^{x^T \beta})^\alpha} = e^{-(\lambda t \gamma)^\alpha}$ explizit darstellen und mathematisch «einfach» sowie es nicht nur ein beschleunigtes Ausfallzeiten-Modell ist, sondern es hat auch **proportionale Hazard-Eigenschaft**.

- Hazardfunktion Weibullverteilung ist $h(t|x) = \lambda \alpha (\lambda t)^{\alpha-1} \cdot e^{-x^T \beta} = \lambda \alpha (\lambda t \gamma)^\alpha$ und damit ist Verhältnis Weibull-Hazardraten bei zwei verschiedenen Zuständen der erkl. Var., gekennzeichnet durch x_i und x_k , $\frac{h(t|x_i)}{h(t|x_k)} = \left(\frac{e^{-x_i^T \beta}}{e^{-x_k^T \beta}}\right)^\alpha = \left(\frac{\gamma_{x_i}}{\gamma_{x_k}}\right)^\alpha$

Eigenschaft **Verhältnis Hazardraten unabhängig von Wartezeit t** ist, heisst **Proportionale-Hazard-Eigenschaft**.

- Unterscheiden erkl. Vektoren x_i, x_k nur darin, dass beim einen Behandlung A vorgenommen $x_i^{(1)} = 1$ und anderen B $x_k^{(1)} = 0$, dann ist bei positiven Koeffizienten β_1 Hazardrate Gruppe Behandlung A immer um Faktor $(e^{\beta_1})^\alpha = e^{\beta_1 \cdot \alpha}$ kleiner (wegen Minuszeichen, das nicht mehr Faktor ist) als jene Gruppe B. Folglich Wartezeit Eintritt Ereignis Behandlung A länger als B. Weil **Faktor unabhängig Wartezeit, Behandlung A immer bessere unabhängig Wartezeit**.

proportionale Hazard-Eigenschaft führt weiter dazu, dass Hazardrate konzeptionell auch als $h(t|x) = h_0(t) \cdot (e^{-x^T \beta})^\alpha = h_0(t) \cdot e^{-x^T \beta \alpha}$ gilt, $h_0(t) = \lambda \alpha (\lambda t)^{\alpha-1}$ ist **Referenz-Hazardrate** (baseline). Hazardrate Referenzzustandes in erkl. Var.

$x_i = 0$. $S(t|x) = S_0(t) \cdot e^{-x^T \beta \alpha}$ Survivorfunktion, $S_0(t)$ Survivorfunktion der Referenz-Hazardrate $\exp(-(\lambda \cdot t)^\alpha)$

- **Weibull-Regressionmodell** in log-lineare Darstellung mit **Gumbelparameter** $\log(T_i) = Y_i = \beta_0 + x_i^T \beta + \sigma \cdot W_i$, W_i Standard-Gumbel $G(\mu = 0, \sigma = 1)$ verteilt. Weibull-Regressionsmodell Weibull-Parameter $h(t|x, \alpha, \lambda, \theta) = \lambda \alpha (\lambda t)^\alpha \cdot e^{x^T \theta}$
- Weibull- $(\alpha, \lambda, \theta)$ Gumbel (β_0, β, σ) Beziehungen $\beta_0 = -\log(\lambda)$ oder $\alpha = \frac{1}{\sigma} = e^{-\log(\sigma)}$, $|\beta = -\frac{\theta}{\sigma}$ oder $\lambda = e^{-\beta \theta} \lambda$ von Baseline $|\sigma = \frac{1}{\alpha}$ oder $\theta = -\frac{\beta}{\sigma}$. In Verteilungen: $\log(T_i) \sim G(\mu_i = \beta_0 + x_i^T \beta, \sigma = \sigma)$ und $T_i \sim W(\lambda_i = \exp(-(\beta_0 + x_i^T \beta)^\alpha), \alpha = \frac{1}{\sigma})$

Maximum-Likelihood-Schätzung: Unbekannten Parameter β_0, β, σ Gumbelparameter werden maximalen Likelihood geschätzt. Falls $r_i := y_i - (\beta_0 + x_i^T \beta)$ ist, dann lautet entsprechende Log-Likelihood $\ell(\beta_0, \beta, \sigma, y, x, z) = \sum_{i=1}^n \{z_i \cdot \log(f(r_i)) + (1 - z_i) \cdot \log(S(r_i))\}$ wobei y_i gleich $\log(t_i)$ ist und z_i angibt, ob die **Beobachtung zensiert ($z_i = 0$) ist oder nicht ($z_i = 1$)**.

- (Allgemein gilt:) Die ML-Schätzungen $\hat{\beta}_0, \hat{\beta}, \log(\hat{\sigma})$ sind asympt. normalvert.
- mit Erwartungswerten $\beta_0, \beta, \log(\sigma)$, Respektive σ Varianzen, gegeben durch inverse Informationsmatrix

```
Mo.W <- survreg(Surv(time, status) ~ x, data = Mo, dist = "weibull"); summary(Mo.W); h.Mo <- data.frame(x = 1000)
#Prognose Weibull-Parameter; (h.a <- 1/Mo.W$scale) #â; (h.li <- predict(Mo.W, newdata = h.Mo, type = "response", se = T))
#Prognose 1/λ(1000) und StdFhrl; h.li$fit * gamma(1+1/h.a) #Schätzung E(X): E(Weibullvert) = 1/α · Γ(1 + 1/α)
1/h.li$fit #λ(x); 1/(h.li$fit + 1.96*c(1,-1)*h.li$se.fit) #VI für λ(x) | Scale im Summary, hier Scale = 0.711 → (α = 1/scale = 1.41)
```

Intercept ist $\hat{\beta}_0 = 5.3$. Folglich ist $\exp(-\hat{\beta}_0)$ geschätzte Parameter λ in Referenzweibullverteilung $\rightarrow \exp(-5.3) = 0.0049$

Baselin-Hazard-R $h_0(t) = \lambda \alpha (\lambda t)^{\alpha-1} = \alpha \cdot \lambda^\alpha \cdot t^{\alpha-1} = 1.41 \cdot 0.0049^{1.41} \cdot t^{1.41-1} = 0.0008 \cdot t^{0.41} \rightarrow$ Exponent $1/Mo.W$scale-1$

- Weil Exponent zu t positiv und kleiner 1 ist Baseline-Hazardrate monoton steigend, jedoch abflachend wachsendem t.
- Weil Exponent zu t positiv und grösser 1 ist, ist Baseline-Hazardrate monoton steigend und wird steiler mit wachsendem t.
- Weil Exponent zu t negativ ist, ist Baseline-Hazardrate monoton fallend und wird flacher mit wachsendem t.
- Weil Exponent zu t gleich 1 ist, ist Baseline-Hazardrate konstant, also eine Gerade \rightarrow Exponential-Verteilung
- Stufe «yes» bei fit mit 0.3042 signifikant positiv, verschiebt Lage nach oben und Wartezeit wird grösser, wenn Variable Beobachtung «yes» ist.
- Koeffizient wPrio (Wurzel Anzahl Verurteilungen) signifikant negativ, verschiebt sich Lage log. Wartezeit nach unten, also wird Wartezeit kleiner je mehr Verurteilungen vorausgegangen sind.
- Einfluss erkl. Var. auf Hazardrate $h(t|x): h(t|x) = h_0(t) \cdot \exp(-\eta \cdot \alpha)$, $\eta = \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_m x_m$ wobei $x =$ erkl. Var.

Modellvergleich: Statistisch Vergleich von zwei ineinander geschachtelten Modellen erfolgt Likelihood-Quotienten-Test.

- Diese Teststatistik ist äquivalent zum Devianztest. Unter Null-Hypothese ist Teststatistik D asymptotisch χ_q^2 -verteilt mit q gleich Anzahl in Null-Hypothese festgelegten Parameter. Umsetzung in R gleich wie bei anderen Modellen, anova()

Bsp: loglogistisches vs. weibull AFT: Welches Modell aufgrund max Log-Likelihood wählen?

Mo.LL <- survreg(Surv(time, status) ~ x, data = mot, dist = "loglogistic"); summary(Mo.LL) #gleiches Modell mit dist weibull Grösseres Max Log-Likelihood bessers Modell. Vergleich nur, wenn beide Modelle gleich viele Parameter besitzen.

Variablenselektion: Wie bei GLM AIC zur Variablenselektion benutzt. Man wählt Modell, dass das AIC minimiert. step()

Modellierung: Weil AFT-Modelle Lage-Skalen-Verteilungsfamilien beruhen, wirken erklärenden Variablen durch linearen Prädiktor auf Lage der Verteilung der Zielvariablen ein. Aber Achtung! **Lage ≠ Erwartungswert; E(log(T)) = μ - γσ**

- Deshalb ergibt sich folgende naheliegende Definition für Residuen: $r_i = y_i - (\beta_0 + x_i^T \beta)$
- Falls Parameter bekannt wären, dann Residuen Gumbel verteilt mit $\mu = 0$ und σ unbekannt. \rightarrow Gumbel-QQ-Plot
- Wenn Beobachtungen rechts-zensiert sind, dann Schätzung $x_i^T \beta$ üblicherweise grösser als $y_i = \log(t_i)$ sein.
- Folglich sind entsprechenden Residuen tendenziell negativ und können als negative Ausreisser in Erscheinung treten.

T-A-/Gumbel-QQ-Plot: source("addFn4survAnal.R"); par(mfrow = c(1,2), las=1); plotGumbelRes(Mo.LL, resType = "raw")

Cox-Snell-Residuen: Sind Stat. Wartezeit-Analyse beliebt: $r_i^{(CS)} = -\log(\hat{S}_{KM}(r_i|x_i))$, r_i Residuen aus Regressionsmodell

- CS-Res Weibull-Regr.: $r_i^{(CS)} = -\log(S(r_i|x_i)) = \exp(\frac{\gamma_i - \beta_0 + x_i^T \beta}{\sigma}) = e^{\frac{r_i}{\sigma}} \rightarrow$ Originalskala zurücktransformierten $\hat{\sigma}$ skalierten Gumbelresiduen. Können $\log(r_i^{(CS)})$ mit Gumbel-QQ-Plot überprüfen, ob (approximativ) exponentialverteilt mit $\lambda = 1$.
- Da $\hat{S}_{KM}(r_i|x_i)$ immer im Intervall [0, 1] liegt, sind alle $r_i^{(CS)}$ nicht negativ und entsprechend schwierig zu interpretieren.
- Sind keine Residuen, da sie aus Transformation der Differenz zwischen beobachteten und angepassten Werten sind.

Devianz-Residuen: Martingal-Residuum $r_i^{(M)} := z_i - r_i^{(CS)}$, Wert Intervall $(-\infty, 1]$. **Zensierte Beobachtungen immer negativ!**

- Verbesserte, symmetrisierte Martingale Residuen: Devianz-Residuen $r_i^{(D)} := \text{sign}(r_i^{(M)}) \cdot \sqrt{-2 \cdot (r_i^{(M)})^2 + z_i \cdot \log(z_i - r_i^{(M)})}$
- R multipliziert $r_i^{(D)}$ mit (-1). Devianz-Residuen sollten symmetrisch um 0 verteilt sein. Summieren sich nicht auf 0 wie in klass. lin. Reg. Wenn nur schwach bis moderat Zensierung, dann sollten $r_i^{(D)}$ aussehen wie iid normalverteilt qqnorm()

Tukey-Anscombe Devianzresiduen/Gumbel-QQ-Diagramm: plot(GumbelRes(Mo.LL, smooth = TRUE))

- **Tukey-Anscombe-Plot** streuen Punkte (nicht) gut um horizontale Null-Linie, wie Glätter zeigt. Auch Streuung Punkte bleibt über linearen Prädiktor-Werte gleich. Ob dies statistisch signifikant ist oder relevant kann nicht beurteilt werden.
- **Gumbel QQ-Plot** streuen Punkte gut um Linie "0 + 1*x", welche Exponentialverteilung mit $\lambda = 1$ repräsentiert (log Cox-Snell Residuen müssen so verteilt sein). Allenfalls haben wir einen "leichten" Ausreisser links unten. → Gesamtfazit...

Cox-Regressionsmodell → bei kontinuierlichen Wartezeiten

- Hazardrate $h_c(t|x)$ ist Produkt Baseline-Hazardrate $h_0(t)$ und erkl. Var.: $h_c(t|x) = h_0(t) \cdot \exp(x^T \theta)$
- Baseline-Hazardrate $h_0(t)$ ist unabhängig von erkl. Var. x , jedoch Form nicht wie Weibull-Reg.modell explizit spezifiziert.
- Funktion $h_0(t)$ kann als Hazardrate für Beobachtung betrachtet werden, bei der **alle erkl. Var. die Werte 0 haben**.
- **Effekte konstant über Zeit** auf Hazardrate einwirken: Unterscheiden sich erkl. Vektoren $x_1^{(1)}, x_2^{(1)}$ nur durch Behandlung A und B, dann ist bei positiven Koeffizienten θ_1 Hazardrate für verbleibende Dauer B immer Faktor e^{θ_1} grösser als A.

library(survival); D.cox <- coxph(Surv(zeit, tod) ~ sex + diab, data = D); length(unique(D\$zeit)); which(table(D\$zeit) > 1)

Vergleich Koeffizientenschätzungen **Cox- mit Weibull-Regressionsmodell** → transformieren Gumbelparametrisierung (-coef(D.weib)[-1]/D.weib\$scale; coef(D.cox) #können nun verglichen werden mit: summary(D.cox\$coefficients[,3]) $\hat{h}(t) = h_0(t) \cdot \exp(0.14 \cdot \text{sex}W) \cdot \exp(0.53 \cdot \text{diab}Yes)$ #Werte aus 1. Block, sig. positiver Koeffizient = kürzer Wartedauer

- **2 Block e^{θ_k} , Faktor mit Hazardrate multipliziert, wenn andere Variablen gleich und k-te Variable um 1 Einheit ändert.**
- **Harrell's concordance index (C-Index)** Masszahl **Anpassungsgüte** (Wartezeit-) Modell: Konkordanz betrachtet alle möglichen Paare von Personen, bei dem eine gestorben und andere mind. so lange lebt. Falls vorhergesagte Lebensdauer grösser für Person, die tatsächlich länger lebte, so ist Vorhersage für Paar konkordant mit tatsächlichen Ausgang. C-Index Anteil konkordanten Paare unter allen Paaren. Werte [0, 1]. Je näher 1, desto besser beschreibt Modell Daten.

Variablenlektion und Modellvergleich step(); anova() → basieren auf «max. partial Log-Likelihoods».

- **Maximum Partial Likelihood:** Um Parametervektor θ zu schätzen, neue Schätzmethodik → maximum partial likelihood
- Schätzansatz ignoriert Baseline-Hazardrate. Berücksichtigt nur Reihenfolge der beobachteten (**unzensierten**) Wartezeiten. **Damit Schätzer zulässig, dürfen keine beobachteten Wartezeitwerte mehrmals vorkommen!**
- Da sonst Wartezeiten nicht rangiert werden können. Trotz Modifikation Likelihoodansatz, können Resultate aus allgemeinen Max-Likelihood-Theorie übernommen werden. **Nicht möglich, damit Vorhersagen zu machen.**

Wertgleiche Wartezeiten: Schätzverfahren funktioniert i.A. nur, wenn **keine wertgleichen Wartezeiten** (tied event times)

- Häufigste Ausweg **Breslow's Approx.** verwenden. Funktioniert gut, sofern **Ties nur selten**. Hat **viele Ties** kann Approximation zu **unzuverlässigen Resultaten** führen. Bessere Approximation ist **Efron** → rechnerisch **effizienter**

Exakter Cox-Regressionsansatz I → nicht in R implementiert, da in Praxis selten genutzt

- Annahme: Ties durch **grobe Rundung Wartezeit**. Folglich gäbe wahre Reihenfolge Wartezeiten. Weil diese unbekannt, alle Reihenfolgen berücksichtigt und Likelihood-Terme hinzufügen. 5 Zeiten gleich → 5! = 120 mögliche Reihenfolge **Möglichkeiten** bei mehr Ties **explodieren** förmlich und **Berechnung** wird auch mit heutigen Mitteln zu **zeitintensiv**.

Exakter Cox-Regressionsansatz II – diskreter Cox-Regressionsansatz → ties = "exact"

- Annahme: Ties durch Diskretheit Wartezeiten verursacht. Folglich kann keine wahre Reihenfolge geben. Ansatz beruht auf anderen Modell, welches kein proportional Hazardmodell ist. Schätzmethodik beruht «partiellen Likelihood-Ansatz».

coxph aus Package **survival** hat **Efron's Ansatz** als StandardEinstellung. Argument **ties** kann **"breslow"** oder **"exact"** wählen.

- Schätzverfahren bei Cox-Regression geht davon aus, dass **keine wertgleichen Wartezeiten** beobachtet werden
- Wird stark gerundet oder stammen Wartezeiten von diskreten Zeitskala, dann versagt eigentliche Schätzverfahren
- **Bei wenigen wertgleichen Wartezeiten kann man Problem mit Approximationen umgehen. Bei vielen nicht!**

Regressionsmodelle bei diskreten Wartezeiten

- Wenn Wartezeit **diskrete Zeitskala** gibt **zwei exakte Vorschläge**. Sind **rechenintensiv** → **nur kleinen Stichproben**
- Diskret Wartezeit mit **GLM** → ist einfacher, Datenaufbereitung aufwändiger. Möglich zeitveränderlich erkl. Var. einbinde
- Einfluss erkl. Var. auf Hazardfunktion $\lambda(t)$ wie in binomialen Regression (GLM): $g(\lambda(t|x)) = \beta_0 + \beta_1 x^{(1)} + \dots + \beta_m x^{(m)} = \beta_0 + x^T \beta$ wobei $g(\cdot)$ Linkfunktion ist. Modellierung $\lambda(t|x)$ nicht nur bei einem t , sondern für alle Intervalle $t = 1, 2, \dots, q$.
- Einfacher Ansatz: **nur Achsenabschnitt β_0 von Wartezeit t abhängig** machen (mit $\beta_0 \rightarrow \gamma_t$): $g(\lambda(t|x)) = \gamma_t + x^T \beta$
- **Interpretation dieses Ansatzes hängt von Wahl Linkfunktion ab!** (Logit-Modell oder Komplementäres Log-Log-Modell)
- **Modellansatz modelliert Effekte der erkl. Variablen auf Hazardfunktion** → es wird ganze **Survivor-Funktion** modelliert

Logit-Modell für diskrete Wartezeiten: Logit-Modell modelliert Hazardfunktion wie folgt: $\log\left(\frac{\lambda(t|x)}{1-\lambda(t|x)}\right) = \gamma_t + x^T \beta$

- variieren nur Achsenabschnitte γ_t über Wartezeit t . • Weiteren Koeffizienten β sind unabhängig von t .
- Linke Seite = log. Wettverhältnis (log-odds) Hazardfunktion. Mit Wahr' lautet (bedingte) Wettverhältnis: $\frac{P(T=t|T \geq t, x)}{1-P(T=t|T \geq t, x)}$
- $x = 0$, dann Wettverhältnis Hazardwahrscheinlichkeit = exponenzierten Parameter γ_t : $\frac{\lambda(t|x=0)}{1-\lambda(t|x=0)} = \exp(\gamma_t)$
- **Parameter γ_t , $t = 1, \dots, q$ legen über Hazard-Wettverhältnis, Baseline-Hazardfunktion fest. Erkl. Var. skalieren** (transformierte) **Baseline** einheitlich über Zeitintervalle $t = 1, \dots, q$ mit Faktor: $\frac{\lambda(t|x)}{1-\lambda(t|x)} = \exp(\gamma_t) \cdot \exp(x^T \beta)$, $t = 1, \dots, q$
- **Odds-Ratio Hazardwahrscheinlichkeiten ändert um Faktor $\exp(\beta_j)$, wenn sich erkl. Variable x^j um eine Einheit ändert.**
- Baseline-Hazard hängt von Parametrisierung in erkl. Var. ab, weil sie der Hazardfunktion entspricht, bei der $x = 0$ ist!
- Quotient Wettverhältnis Hazardwahr' bei **versch. Wartezeiten** t und s : $\exp(\gamma_t - \gamma_s)$ unabhängig von x ; hängt also nur von Baseline-Hazards ab. Quotienten Wettverhältnis Hazardwahr' bei **verschiedenen Werten** x_1 und x_2 der erkl. Var.: $\exp((x_1 - x_2)^T \beta)$ unabhängig von Wartezeit t ; hängt nur von $(x_1 - x_2)$ ab.
- 2 Survivorfun. Versch. Prädiktoren kreuzen nie: $\lambda(t|x_1) > \lambda(t|x_2) \rightarrow$ Survivorfun x_1 alle Intervalle t kleiner als jene bei x_2

Gruppiertes proportionales Hazardmodell → link = cloglog

- Proportionales Hazardmodell, jedoch (stark) gerundete Wartezeiten → Zeitachse unterteilt in disjunkte Intervalle
- Äquivalent: $\log(-\log(1 - \lambda(t|x))) = \gamma_t + x^T \theta$ der mit komplement. log-log Link aus binären Regression übereinstimmt.
- Was Einfluss erkl. Var. x betrifft, wird **diskrete Modell gleiche Analyse erlauben wie proportionale Hazard-Modell**.
- $\lambda(t|x)$ des diskretisierten proportionalen Hazardmodells kann **nicht als Produkt** einer Baseline-Hazard-Wahrscheinlichkeit (\neq Hazardrate) und des Einflusses der erklärenden Variablen dargestellt werden.
- Quotient logarithmierten Survivorfunktionen ist unabhängig von Wartezeit: $\frac{\log(S(t|x_1))}{\log(S(t|x_2))} = \exp((x_1 - x_2)^T \theta)$
- Obwohl Ansatz nicht proportionale Hazeideigenschaft besitzt, heisst es **«gruppiertes proportionale Hazardmodell»**, weil aus proportionalen Hazardmodell von kontinuierlichen Wartezeiten bei Annahme (starken) Rundungen ableitet.

Schätzung bei diskreten Wartezeiten | proportional continuation ratio model → link = logit

- ML-Schätzung Regressionsmodell diskrete Wartezeit führt zu äquivalent Lösung wie GLM. Müssen Daten aufbereiten
- Jede **Wartezeitbeobachtung** ist die Realisierung **zweier Zufallsvariablen**: Eine Zufallsvariable ist **Wartezeit T_i** , und die andere Zufallsvariable ist **Indikator Z_i** , der angibt, ob **Ereignis beobachtet wurde ($Z_i = 1$) oder nicht ($Z_i = 0$)**.
- Zensierung Ende jeweiligen Intervall. Annahme: **Zensierung nicht informativ** (Zensierung hängt nicht von Parametern Wartezeitprozess ab) → Zensierungselemente können in Likelihood zu **Konstanten c_i** zusammengefasst werden.
- Mit Hazardwahrscheinlichkeiten $\lambda(t_i|x)$ gilt dann für Likelihood $L_i = c_i \cdot \lambda(t_i|x)^{z_i} \cdot (1 - \lambda(t_i|x))^{1-z_i} \cdot \prod_{j=1}^{t_i-1} (1 - \lambda(t_j|x))$
- ML-Schätzung Parameter im diskreten Wartezeit-Regressionsmodell hat gleiche asymptotische Verteilung wie jene, die auf binären Repräsentation der Übergänge zwischen Wartezeitkategorien (d.h. GLM-Anpassung) beruht.

library(discSurv); long <- dataLong(data, timeColumn = "Dauer", eventColumn = "Beob", timeAsFactor = T) gphm <- glm(y ~ 1 + timeInt + Kundengruppe + Region, data = long, family = binomial(link="cloglog")); summary(gphm) drop1(gphm, test = "Chisq") Vorteil von drop1() gegenüber summary(): → erkl. Faktorvariablen

- informiert, welche Variablen signifikant sind. Können damit klären, ob timeInt (und damit Hazardwahr) konstant sind.
- P-Werte zuverlässiger, da beruhen Teststatistiken/Deviantests → Asymptotik wirkt schneller (bzgl. Stichprobengrösse).
- gibt für Variable an, ob Einfluss auf ZV signifikant ist unter Berücksichtigung, dass andere Variablen auch Einfluss haben

Darstellung Hazardwahr' mit KI: h.q <- length(dummy.coef(gphm)\$timeInt); h.hc <- dummy.coef(gphm)\$timeInt h.hcse <- summary(gphm)\$coefficients[1:h.q, 2]; plot(1-exp(-exp(h.hc)), type="l", panel.first=abline(v=(1.4)^1/3/2), ylab = "Hazard-Wahr", xlab = "Dauer in ..."); lines(1-exp(-exp(h.hc+2*h.hcse)), lty=2, col=2); lines(1-exp(-exp(h.hc+2*h.hcse))) Alternative mit predict: newdata = Referenz-Hazardwahr'. R findet heraus, welche Levels Variablen sonst noch haben h.new <- data.frame(timeInt=as.factor(1:52), fin="no", race="black", wexp="no", paro="no", mar="married", educ="2", fAge = "(0,21]", wPriO = 0); h.p <- predict(gphm, newdata=h.new, se.fit=T, type="link") #h.hc=h.p\$fit und h.hcse=h.p\$se.fit

- **Variablenlektion:** wie GLM AIC: step(g.glm, scope = list(lower = ~ 1 + timeInt, upper = formula(g.glm)))
- Bei Wechselwirkungen/Interaktionen: Sobald Variable **timeInt** mit einer erkl. Var. (z.B. V1) in Wechselwirkung tritt, bedeutet dies, dass unterschiedliche Baseline-Hazards haben für Population, die durch V1 beschrieben wird.
- **Additive Modelle:** wenn Hazardwahr' aus vielen Intervallen zusammensetzen (Faktorvariable **timeInt** hat viele Stufen).
- Wenn viele Stufen gibt, ist dies einzige Möglichkeit Hazardwahr' zu schätzen, weil sonst Design-Matrix zu gross wird.
- **gam()** aus **library(mgcv)**: + numerisch geschickter implementiert, beruht auf Splines-Schätzungen der unbekanntem Transformation. +/- wählt automatisch optimalen Glättungsparameter, führt oft zu einer zu glatten Funktion
- **gam()** aus **library(gam)**: + Schätzungen unbekanntem Transformation mit loess. Kann zu lokalen Minima konvergieren. (-) Glättungsparameter von Hand wählen. Ziehe leicht Unterglättung vor, weil von Auge etwas «weiter glätten» kann. long\$Ntr <- as.integer(as.character(long\$timeInt)); f.gam <- gam(y ~ lo(Ntr) + fin + mar + lo(wPriO) + lo(AGE), data = long, family = binomial(link = "cloglog")); par(mfrow=c(2, 3)); plot(f.gam, se = T) #Residuen wie unten analysieren.

- **Devianz-Residuen** sind nicht für Überprüfung Modelleignung geeignet, weil Pseudo-Beobachtungen y_i aus n Beobachtungen abgeleitet wurden. Ursprünglichen Daten bestehen aus n unabhängigen Beobachtungen der Form (t_i, z_i, x_i) .
 - o Bei GLM wurde als Gütemass Devianz D eingeführt als (-2) mal Log-Likelihood, $D = -2 \cdot \ell$. $D = \sum_{i=1}^n d_i$
 - o «Elemente» d_i der Devianz als Grundlagedevianz Residuen für disk. Zeiten: $r_i^{(D)} = \sqrt{d_i}$ ohne Vorzeichen (GLM mit)
 - o Zensierung ist in $r_i^{(D)}$ nicht vollständig berücksichtigt, da Zensierungsmechanismus, der in Zensierungskonstanten $\log(c_i)$ festgehalten ist, weggelassen wurde. Definition für Residuen ungewöhnlich, da alle **Residuen positive Werte**.

library(discSurv); haz <- predict(glm, type = "response"); dR <- devResid(dataLong = DL, hazards = haz) qqnorm(dR\$Output\$DevResid) #Darstellung der berechneten Devianz Residuen. Werden Praxis nicht weiter einsetzen

- **Adjustierte Devianz-Residuen:** mit einer Transformation adjustierten Residuen.
 - o genähert **normalverteilt**, werden Normal-QQ-Plot zur Beurteilung, ob Modell Daten adäquat beschreibt, eingesetzt.
 - o Beachten: Beobachtungen, bei denen eine **Hazardwahrscheinlichkeiten 0** geschätzt ist, zu ganz **groben Ausreissern** in den adjustierten Devianz-Residuen führen. adR <- adjDevResid(dataLong = DL, hazards = haz); qqnorm(adR\$Output\$AdjDevResid); qqline(adR\$Output\$AdjDevResid) #Ist Normalverteilung ersichtlich? → i.O.

- **Martingal-Residuen:** sind bei disk. Wartezeiten analog kontinuierlichen Wartezeiten definiert: $r_i^{(M)} = z_i - \sum_{s=1}^{t_i} \hat{\lambda}(s|x_i)$
 - o Mit Hilfsgrössen y_{is} kann Martingal-Residuum schreiben: $r_i^{(M)} = \sum_{s=1}^{t_i} y_{is} - \hat{\lambda}(s|x_i)$. Folglich entspricht Martingal-Residuum Differenz zwischen Übergangsindikator y_i und geschätzten Eintrittswahrscheinlichkeit.
 - o Martingal-Residuen nur Logit-Modell Summe 0, $\sum_{i=1}^n r_i^{(M)} = 0$, Wertebereich Residuen [-1, 1], asymmetrisch um 0.
 - o Falls Modell Daten adäquat beschreibt und entsprechend alle erkl. Var. enthält, dann sind Martingal-Residuen zufällig und unkorreliert mit erkl. Var.. Deshalb werden sie auch eingesetzt, um Effekte der erkl. Var. zu überprüfen.

mR <- martingaleResid(dataSetLong=DL, hazards=haz); plot(mR, covariates=c("fin", "erkl. Var X"), dataSetLong = DL) Faktorvar. = Boxplot, sonst Streudiagramm + Glätter. Boxplot mean rechnen (da Median in Boxplot). aggregate(as.vector(mR), by = list(fin = data\$fin, FUN = mean) #data=data normal, nicht lang, sum(as.vector(mR)) #Werte=Null → gut