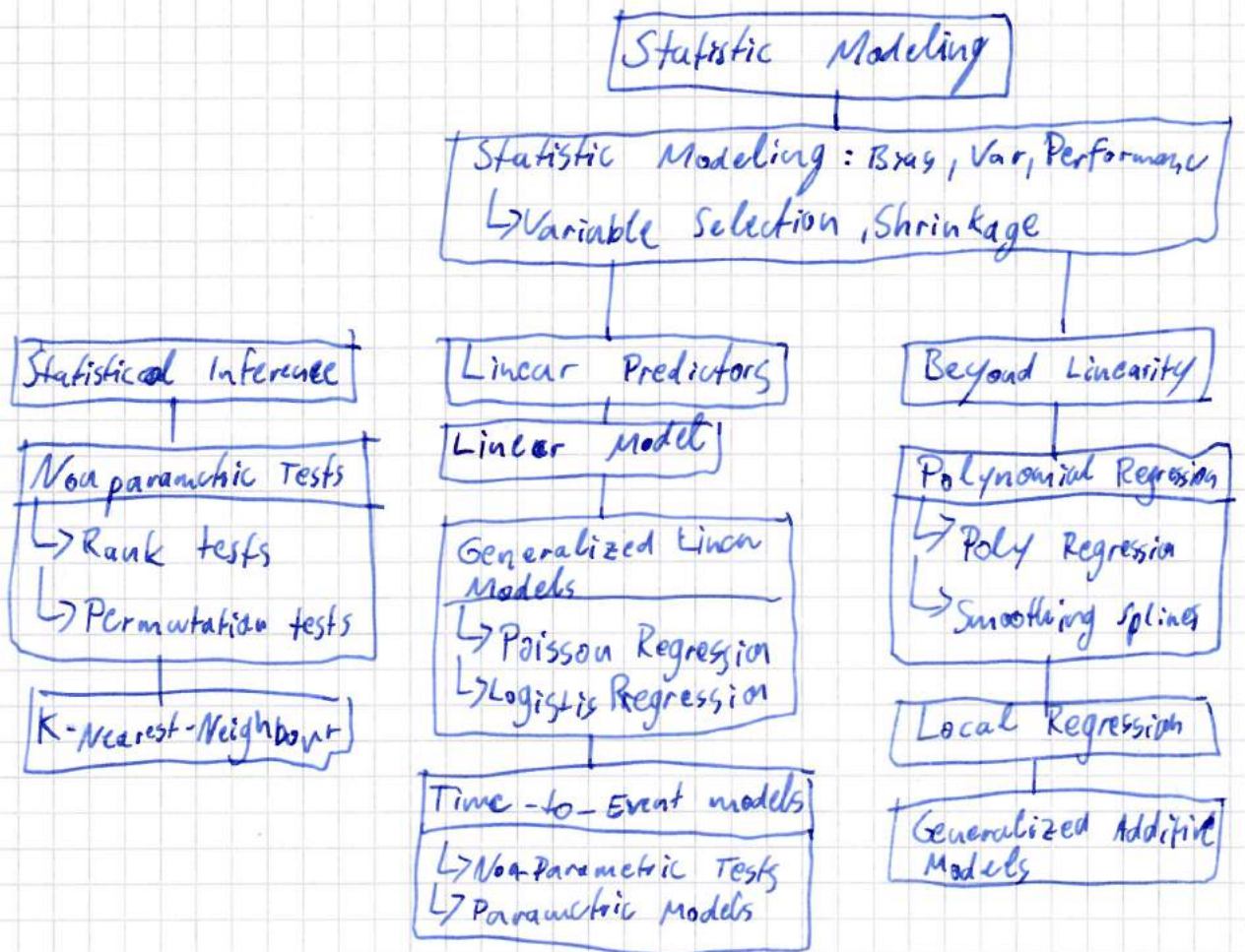


Übersicht



Vorhersagefehler

□ Bias - Variance - Tradeoff

↳ Expected Test Error = Bias² + Var + Irreducible Error

↳ $MSE = \text{Bias}[\hat{f}(x)] + \text{Var}[\hat{f}(x)] + \sigma^2$

$E[(y - \hat{f}(x))^2]$

Modellfehler

Datenfehler

deswegen ist der auch irreducible

↳ Varianz + \Rightarrow overfitting

↳ Bias + \Rightarrow Underfitting

□ Anpassungsgüte

↳ $AIC = -2LL + 2(p+1)$

↳ $BIC = -2LL + p \log(n)$

Erwartungswert

$L(b) = \prod_{i=1}^n P(y=y_i | b, x_i)$

$LL(b) = \sum_{i=1}^n y_i (b^T x_i) - \exp(b^T x_i) - \log(y_i!)$

↳ Der AIC ist immer nur auf den gleichen Daten vergleichbar, da dieser ein Vergleichswert und kein Richtwert ist (wächst mit der Anzahl an Daten)

↳ AIC ist ein Maß für die Anpassungsgüte und bestraft Overfitting

↳ In der Praxis nimmt man das Modell mit dem geringsten AIC bei etwa gleichen MSE auf den Testdaten

Nichtparametrisch

Nicht Parametrisch

- ▣ Verteilung nicht bekannt
- ▣ Nullhypothese nicht über Verteilung definierbar
 - ↳ Ablehnbereich über kritische Werte

Parametrisch

- ▣ Für jeden Parameter θ ist die Verteilung der ZV bekannt
- ▣ Nullhypothese ist Punkt/Intervall und somit ist die Verteilung bekannt
 - ↳ Ablehnbereiche über Quantile der Verteilung

Nichtparametrische Tests

- ↳ Rangtests: Werte bekommen Ränge zugewiesen und über diese wird getestet
- ↳ Permutations tests: Beobachtete Werte werden als fest betrachtet und die Verteilung zufällig behandelt
- ↳ Monte-Carlo-Permutations test: Simulationstest über Permutationen

① Einstichprobentests / Verbundene Stichproben

$$P(X \leq x_{\text{med}}) = P(X \geq x_{\text{med}}) = \frac{1}{2}$$

▣ Vorzeichenstest

- ↳ X_i i.i.d. mit stetiger Verteilung
- ↳ $H_0: \{x_{\text{med}} = \theta_0\}$ vs $H_1: \{x_{\text{med}} \neq \theta_0\}$ zu α
- ↳ $Y_i = \begin{cases} 1, & X_i < x_{\text{med}} \\ 0, & X_i \geq x_{\text{med}} \end{cases} \Rightarrow Y_i \stackrel{H_0}{\sim} \mathcal{B}(1, \frac{1}{2}) \Rightarrow T = \sum_{i=1}^n Y_i \stackrel{H_0}{\sim} \mathcal{B}(n, \frac{1}{2})$
- ↳ $x_i = \theta_0$ können wegen stetig theoretisch nicht auftreten \Rightarrow weg reduzieren
- ↳ nominal und ordinal möglich
- ↳ $n \geq 30: \mathcal{B}(n, \frac{1}{2}) \sim \mathcal{N}(\frac{n}{2}, \frac{1}{4}n)$

▣ Wilcoxon-Vorzeichen-Rangtest

- ↳ X_i i.i.d. symmetrisch mit stetiger Verteilung
- ↳ $H_0: \{x_{\text{med}} = \theta_0\}$ vs $H_1: \{x_{\text{med}} \neq \theta_0\}$
- ↳ $D_i = X_i - \theta_0 \Rightarrow$ Bilde Ränge der $|D_i|$ (Ordinalskala ab 1 ↑)
- ↳ $W^+ = \sum_{D_i > 0} \text{rk}(|D_i|)$, $W^- = \sum_{D_i < 0} \text{rk}(|D_i|) \Rightarrow T = \min(W^+, W^-)$
- ↳ $W^+ \stackrel{H_0}{\approx} W^- \Rightarrow$ Ablehnbereich für T über n
- ↳ Für $X_i = X_j$ bekommen die Werte den durchschnittlichen Rang
- ↳ $n \geq 30$: Approximation über Normalverteilung möglich

▣ Permutationstest

- ↳ X_i i.i.d. mit $D_i = X_i - Y_i$ symmetrisch um θ verteilt $\leftarrow M \subset \{1, \dots, n\}$ beliebig
- ↳ $H_0: \{\theta = 0\}$ vs $H_1: \{\theta \neq 0\}$
- ↳ $(X_i - Y_i) \stackrel{H_0}{\sim} (Y_i - X_i) \Rightarrow T = \frac{1}{2} \sum D_i \sim T_M = \frac{1}{2} (\sum_{i \in M} X_i - Y_i + \sum_{i \in M^c} Y_i - X_i)$
- \Rightarrow Test über alle M -Wahlweise mit Monte-Carlo zufällig probieren.

② Zweistichproben tests / Unverbundene Tests

□ Wilcoxon-Rangsummen-Test

- ↳ X_i, Y_i unabhängig. X_i stetig mit F_X Verteilung. Y_i stetig mit F_Y Verteilung.
- ↳ Es gelte: $F_X(x) = F_Y(x+\theta)$
- ↳ $H_0: \theta = 0$ vs $H_1: \theta \neq 0$ \Rightarrow Wahlweise \leftarrow / \rightarrow , um die Lage zu testen
- ↳ Bilde Ränge über $X_i, Y_i \Rightarrow T = \sum \text{Rk}(X_i)$
- ↳ $x = x_j \Rightarrow$ zufälliger Rang; $x_i = y_j$ Durchschnitt
- ↳ $n \geq 30$: Normalverteilungsapproximation

□ Permutationstest

- ↳ $Z = (X_1, \dots, X_n, Y_1, \dots, Y_m) \Rightarrow Z_i$
- ↳ wie Wilcoxon ohne Stetigkeitsannahme
- ↳ $Z \sim_{H_0} Z_\pi$ für alle Permutationen $Z_\pi = Z_{\pi(1)}, \dots, Z_{\pi(n+m)}$
- ↳ zweiseitige Testentscheidung mit $\frac{\alpha}{2}$

Durchschnittsrang

D_i	1.5	2.5	3.5	3.5	3.5	4.5
Rang ohne Durchschnitt	1	2	3	4	5	6
Rang mit Durchschnitt	1	2	4 (3+4+5)/3			6

Verallgemeinerte Lineare Modelle (GLM)

Idee: Die Fehlerverteilung der abhängigen Variablen (Response variable) bekommt eine andere Verteilung (anstatt einer Normalverteilung)

Bisher: $E(y_i | x_i) = b^T x_i$ ← Transformation der Erwarteten Verteilung

↳ Nun: mit Link Funktion $\eta_i = g(E(y_i | x_i))$
diskret!

□ Poisson-Regression

↳ y_1, \dots, y_n (bedingt) unabhängig aus $\mathcal{N} \Rightarrow y_i | x_i = x_i \sim \text{Poi}(\lambda_i)$

mit $\lambda_i = \exp(b^T x_i) \Rightarrow E(y_i | x_i) = \lambda_i = \exp(b^T x_i)$

↳ Optimierungsansatz

↳ $L(b) = \prod P(y_i = y_i | b, x_i) = \prod \exp(-\lambda_i) \frac{\lambda_i^{y_i}}{y_i!}$ ↗ ersetzen

$LL(b) = \sum y_i (b^T x_i) - \exp(b^T x_i) - \log(y_i!)$

$[LL(b)]' = \Delta l = \sum x_i (y_i - \exp(b^T x_i)) = s(b) = 0 \Rightarrow \text{GLS}$

↳ Modellannahmen: Pearson-Residuen streuen mit $\text{Var}=1$ um $E=0$

3 Kernelemente der GLM

① Exponentielle Familie von W-Verteilungen

② Linear Predictor $\eta = b^T x$

③ Link Funktion $g(E(y_i | x_i)) = \eta_i \Rightarrow E(y_i | x_i) = g^{-1}(\eta_i)$

□ Weitere GLM

↳ Änderung der Link Funktion g . Notwendig: g invertierbar.

□ Devianz \hat{b}_M ML-Schätzer

↳ $L_M(\hat{b}_M)$, $LL_M(\hat{b}_M) = \log(L_M)$

↳ $D_M = 2(LL_{\text{opt}} - LL_M)$ ← Devianz

⇓

□ Likelihood-Ratio-Test

↳ $H_0: \{b_{p+1} = \dots = b_{p+q} = 0\}$ ← $x_1, \dots, x_p, x_{p+1}, \dots, x_{p+q}$ Kovariaten

↳ $LR = D_M - D_m = 2(LL_M - LL_m) \underset{H_0}{\sim} \chi^2_q$

Basiserweiterung in X

Modell: $Y = f(X) + \varepsilon$ mit $E(\varepsilon) = 0$, $V(\varepsilon) = \sigma^2$

Idee: $f(X) = b_0 + b_1 h_1(X) + \dots + b_m h_m(X) = b_0 + \sum_{i=1}^m b_i h_i(X)$

$\rightarrow Z = \begin{pmatrix} 1 & h_1(x_1) & \dots & h_m(x_1) \\ \vdots & \vdots & & \vdots \\ 1 & h_1(x_n) & \dots & h_m(x_n) \end{pmatrix}$ Designmatrix

\rightarrow Nichtlineare Modelle mit linearen Methoden verarbeiten

□ Polynomregression: $Y = b_0 + b_1 X + b_2 X^2 + \dots + b_d X^d + \varepsilon$

□ Stückweise Polynomial:

$\rightarrow (-\infty, \xi_1) \cup [\xi_1, \xi_2) \cup \dots \cup [\xi_{K-1}, \xi_K) \cup [\xi_K, \infty)$ Intervalle
 $\leftarrow K$ Knoten ξ_k
 $\leftarrow K+1$ Intervalle

$\rightarrow Y = f_k(X) + \varepsilon$ Polynom in jedem Intervall

\rightarrow Anzahl freie Parameter: $(d+1)(K+1)$
 \leftarrow Anzahl der Polynome
 \uparrow
Komplexität pro Polynom

□ Regression Splines

$\rightarrow d-1$ mal stetig diff \Rightarrow Ableitungen und Punkte stimmen in Knoten überein

$\rightarrow d+K+1$ Freiheitsgrade

\rightarrow Bsp: $h_i(x) = x^{d-1}$, $i=1, \dots, d+1$; $h_i(x) = (x - \xi_{i-d-1})_+^d$, $i = d+2, \dots, d+K+1$

$$f(x) = \sum_{m=1}^{d+1} b_m x^{m-1} + \sum_{k=1}^K b_{d+1+k} (x - \xi_k)_+^d = \sum_{m=1}^{d+K+1} b_m h_m(x)$$

Truncated power basis

\rightarrow Knotenwahl z.B. äquidistante Quantile und CV für die Anzahl

\rightarrow Natürlicher Spline: $s_0''(x_0) = s_{n-1}''(x_n) = 0$ oder $(-\infty, \xi_K] \cup [\xi_K, \infty)$ linear verläuft

\rightarrow Anzahl freie Parameter: $K-d+3$

□ Smoothing-Regression-spline

$$\rightarrow \text{RSS}(f, X) = \sum_{i=1}^n (y_i - f(x_i))^2 + \lambda \int (f''(t))^2 dt$$

$\rightarrow f$ minimiert: natural cubic spline mit Knoten in x_i

$\rightarrow \lambda$ mittels CV

□ Effective Degrees of Freedom

\rightarrow Linear: $\hat{y}_i = (Hy)_i$, $H = X(X^T X)^{-1} X^T \Rightarrow p = \text{spur}(H)$

\rightarrow RSS: $\hat{y}_i = (S_\lambda y)_i \Rightarrow p = \text{spur}(S_\lambda)$

Hängt nur
von λ, X ab

Anzahl zu
schätzende Parameter
 \leftarrow
 $p = \text{spur}(H)$

Anzahl zu schätzende
Parameter
 $=$
Effective Degrees of
Freedom

Kernel Smoother

□ KNN (K-Nearest-Neighbor)

↳ Idee: Moving average

↳ $N_K(x) = \{x_{(1)}, \dots, x_{(K)}\}$ mit $|x - x_i|$ K nächste Nachbarn (kleinste Δ)

↳ $\hat{f}(x) = \frac{1}{K} \sum y_i$ für i von $x_i \in N_K(x)$

$$= \underset{b_0}{\operatorname{argmin}} \sum_{i=1}^n \underbrace{K(x_i, x_i)}_{\text{Kernel}} (y_i - b_0)^2 \quad \text{mit } K(x_i, x_i) = I(x_i \in N_K(x))$$

⇒ Lokale Polynomiale Regression von Grad 0

⇒ Nicht glatt sondern konstant auf Intervallen

↙ Verbesserung

□ Kernel Average Smoother

↳ Idee: $K(x_i, x_i)$ durch eine Gewichtungsfunktion ersetzen

↳ $K_\lambda: \mathbb{R}^2 \rightarrow \mathbb{R}^+$ mit Fensterbreite $\lambda \Rightarrow \int_{-\infty}^{\infty} K_\lambda(x_i, x_i) dx = 1$

↳ K_λ stetig $\Rightarrow \lambda$ als smoothing parameter

□ Lokale Lineare Regression

↳ $\hat{f}(x) = \hat{b}_0(x) + \hat{b}_1(x)x$

↳ $(\hat{b}_0(x), \hat{b}_1(x)) = \underset{(b_0, b_1) \in \mathbb{R}^2}{\operatorname{argmin}} \sum_{i=1}^n K_\lambda(x_i, x_i) (y_i - b_0 - b_1 x_i)^2$

↳ Anstatt quadratische Differenz wird nun über die optimale Verbindungsgerade geblättert

□ Lokale Polynomiale Regression

↳ $\hat{f}(x) = \hat{b}_0(x) + \hat{b}_1(x)x + \hat{b}_2(x)x^2$

↳ Wie Linear nur Erweiterung um Polynome höherer Ordnung

↳ d>2 möglich, aber unüblich.

↳ Grund: Quadrat erweitert das Modell und ermöglicht kurvige Glättung. Höhere Ordnungen würden Wendepunkte ermöglichen, nur diese sind eher durch die Fensterbreite λ abzufangen als durch den Kernel

