

Einführung

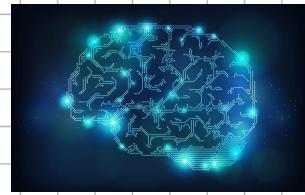
Zw. wichtigste Pfeiler: Daten & Modell

Entwicklung einer ML-Anwendung:

1. Problem modellieren \rightarrow 2. Lernen \rightarrow 3. Evaluation (\rightarrow evtl. iterieren)

No free lunch theorem:

für jedes Modell / jedem Algorithmus gibt es Probleme für welche die Resultate sehr schlecht sind



Supervised Learning (SL)

basiert auf Eingabewort $\vec{x} \in \mathcal{X}$ möchten wir einen Ausgabewort $\vec{z} \in \mathcal{T}$ vorhersagen wobei wir für den Lernprozess eine Daten-Menge $D = \{(t_k^i, \vec{x}_k)\mid 1 \leq k \leq N_D\}$ von beobachtbaren Wertepaaren zur Verfügung haben

\mathcal{X} ist oft \mathbb{R}^d , kann aber Menge mit anders strukturierten Elementen sein

\mathcal{T} kann \mathbb{R}^d sein (Regressionsproblem) oder endliche Menge $\{c_1, \dots, c_m\}$ (Klassifizierungsproblem)

Vorgehen: gesuchte Funktion $f: \mathcal{X} \rightarrow \mathcal{T}$ spezifizieren, die an Trainingsdaten D angepasst wird und deren Funktionswerte $f(\vec{x})$ zur Vorhersage von $\vec{z} \in \mathcal{T}$ verwendet werden

als Qualitätsmaß benutzt man oft den Verallgemeinerungsfehler

$$\overline{L}_{f,T} = \frac{\sum_{k=1}^{N_T} L(t_k^i, f(\vec{x}_k))}{N_T} \quad \text{über alle Testdaten gemittelter Verlust}$$

wobei L ein für die Anwendung geeignetes Verlustmaß und $T = \{(t_k^i, \vec{x}_k)\mid 1 \leq k \leq N_T\}$ unabhängige Testdaten bezeichnen

\rightarrow es ist niemals zulässig $\overline{L}_{f,D}$ als Qualitätsmaß zu verwenden!

man möchte ein Modell p_m der bedingten Verteilung (Verteilung im Sinne von pmf / pdf) $p(\vec{z} \mid \vec{x})$ erkennen

unter Verwendung eines geeigneten Verlustmaßes können wir aus der gegebenen bedingten Verteilung $p_m(\vec{z} \mid \vec{x})$ dann eine optimale Vorhersagefunktion $f: \mathcal{X} \rightarrow \mathcal{T}$ bestimmen

Unsupervised Learning (UL)

(im Gegensatz zu) überwachtem Lernproblem gibt es keinen ausgewählten Ausgabewort und die Lernbasis ist schlicht eine Menge von beobachteten Daten $D = \{\vec{x}_k \mid 1 \leq k \leq N_D\}$

Ziel ist es, die zugrundeliegende Verteilung zu erfassen / repräsentieren und dadurch im weiteren Sinne nutzbar zu machen

Vorteil: keine (aufgrund) Labeling von Daten nötig

Wahrscheinlichkeitstheoretisch kann man unbefristetes Lernproblem als Anpassung eines probabilistischen Modells p_m an die den Trainings-Daten D zugrundeliegende Verteilung $p(\vec{x})$ anpassen (density estimation)

Häufige Fragestellung:

für Eingabe \vec{x} kompakte, komprimierte sowie für weitere Verarbeitung effiziente und informative Repräsentation \vec{z} (den verborgenen Zustand) finden

indem wir prob. Modell $p_m(\vec{z}, \vec{x})$ an Trainingsdaten D anpassen, können wir danach für neue Eingabe \vec{x} z.B. das Maximum oder Mittelwert \vec{z} der mithilfe trainierten Modells berechneten bedingten Verteilung $p_m(\vec{z}, \vec{x})$ als Repräsentanten für \vec{x}

Beispiele:

- Clustering: Unterteilung von Daten in Gruppen von unterschiedlichen inhärenten Daten
- Dimensionsreduktion: (hochdimensionale) Eingaben \vec{x} niedrigdimensionale Repräsentanten \vec{z} zuordnen so, dass die für die nachfolgende Verarbeitung relevanten Informationen erhalten bleibt und einfache Zugänglichkeit ist
→ der verborgene Zustand wird vielfach als Projektion der Eingabe auf eine Untermannigfaltigkeit (im einfachsten Fall ein offener Unterraum) konstruiert
- Zustandsraummodell

als Qualitätsmaß eignet sich die Likelihood:

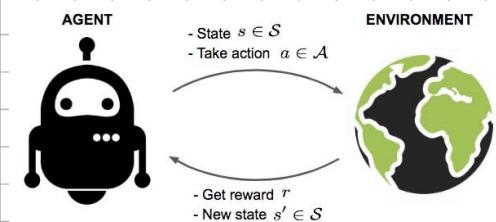
$$L(T|M) = p_M(T) = \prod_{k=1}^{n_T} p_M(\vec{x}_k)$$

es beschreibt die Wahrscheinlichkeit der Testdaten $T = \{\vec{x}_k \mid 1 \leq k \leq n_T\}$ unter dem auf den Daten D trainierten Modells p_M

reinforcement learning (RL)

Agent kann mittels Aktionen mit Umgebung wechselwirken

auf Aktion $a_t \in A$ zum Zeitpunkt t reagiert Umgebung mit
mit Veränderung ihres Zustandes $s_t \in S$ ins $s_{t+1} \in S$ und der Aus-
sendung einer Belohnung r_{t+1} , welche beide vom Agenten beobachtet werden



Agent ist definiert durch seine Handlungsstrategie (policy) $\pi: S \rightarrow A$ mit welcher er basierend auf
aktuellerem Zustand $s \in S$ der Umgebung die Aktion $\pi(s) \in A$ wählt

Ziel ist Finden eines policy $\pi: S \rightarrow A$ mit mit welcher die über sehr viele Episoden gemittelte
(erwartete) akkumulierte Belohnung $g_t = r_1 + \dots + r_t$ maximal gross wird

Herausforderungen:

- Wirkung einer einzelnen Aktion auf akkumulierte Belohnung nicht unmittelbar erfahrbar
- bei einigen Umgebungen sind Belohnungen bei meisten Schritten ≤ 0 und nur gelegentlich (z.B. bei Er-
reichen eines technischer Zustände) Aktionen mit positiver Belohnung verbunden sind

als Abgrenzung zu UU gibt es (in Form der Belohnung) Rückmeldung mittels welcher policy als gesamtes
evaluiert werden muss, im Gegensatz zu SL ist Rückmeldung nicht in Form von (Ausgabe-, Eingabe-) Paaren (a, s)

welche Aktionen in einem Zustand lohnend sind, muss Agent aus seinem akkumulierten
Wissen über Umgebung extrahieren ⇒ hierzu ist einerseits Erforschung (Exploration) der Wirkung
der Aktionen auf akkumulierte Belohnung nötig, aber andererseits sollen auch gute Aktionen
gewählt werden (exploitation), weil ja akkumulierte Belohnung gross sein soll

Modellierung des Lernproblems meist als Form von sequentiellen, markovschen Entscheidungsprozessen
(MDP), dabei ist Umgebung definiert durch bedingte Verteilungen $p(\tilde{s}, \tilde{r} \mid s, a)$ über den neuen
Zustand \tilde{s} der Umgebung und ausgesandte Belohnung \tilde{r} , gegeben dass im aktuellen Zustand s die
Aktion a gewählt wird. policy ist dabei durch bedingte Verteilung $\pi(a \mid s)$ über Aktionen a im
Zustand s der Umgebung definiert
→ beides zusammen erlaubt die Verteilung über Episoden zu berechnen und bildet das Funda-
mentum von Methoden für RL zu unterscheiden

zentrale Begriffe

anhand eines Regressionsproblems mit Eingaben $x \in \mathbb{R}$ und Ausgaben $t \in \mathbb{R}$

Daten: zugrundeliegende Funktion $g(x) = \sin(2\pi x)$ überlagert mit normalverteiltem ($\mu=0, \sigma=0.3$) Rauschen r : $t = g(x) + r$

Datensatz generieren mit auf $[0, 1]$ gleichm. vert. Eingaben Ausgabe gemäß obiger Formel ergibt $D = \{(t_k, x_k) \mid 1 \leq k \leq n_D\}$ und $T = \{(t_k, x_k) \mid 1 \leq k \leq n_T\}$

Verlustfunktion: $L(t, p) = (t - p)^2$

$$\text{Vergleichbarer Fehler: } \overline{L}_{T,T} = \frac{\sum_{k=1}^{n_T} (t_k - f(x_k))^2}{n_T} \text{ MSE oder } \sqrt{\frac{\sum_{k=1}^{n_T} (t_k - f(x_k))^2}{n_T}} \text{ RMSE}$$

für ein genügend großes n_T erwarten wir für die wahre Funktion g einen RMSE in der Größenordnung der Standardabweichung des Rauschens also $\text{Errs}(g, T) \approx \sigma$

Modell für zugrundeliegende Regelmäßigkeit:

$$y(x, \vec{w}) = \sum_{j=0}^m w_j x^j = w_0 + w_1 x + \dots + w_m x^m$$

Modell für Abweichung: Normalverteilung

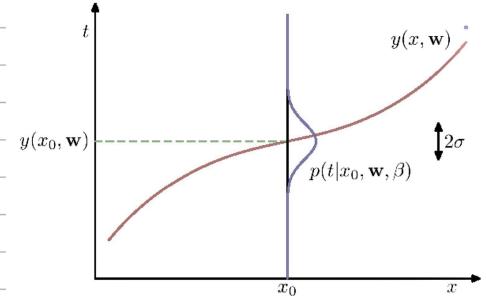
Insgesamt erhalten wir das Wahrscheinlichkeitsfkt.

Modell → Präzision (inverse Varianz)

$$p(t|x, \vec{w}, \beta) = N(t | y(x, \vec{w}), \beta^{-1}) \quad \text{polynomiale Regression mit Gaussianem Rauschen}$$

mit adaptiven Parametern $\vec{\theta} = (\vec{w}, \beta)$ mit $\vec{w} \in \mathbb{R}^{M+1}$, $\beta > 0$ wobei $\beta = V^{-1} = \sigma^{-2}$

es folgen 3 allgemeine Methoden für die Anpassung von polyn. Modellen an Daten



Maximum-Likelihood-Methode Karte 1

$$\text{Likelihood der Daten } D = L(\vec{\theta}^* | D) = p(D | \vec{\theta}^*) = \prod_{k=1}^{n_0} p(t_k | \vec{\theta}^*)$$

→ pmf (diskrete Daten) oder pdf (kont. Daten)

Ziel: Parameter $\vec{\theta}_{\text{ML}} = \underset{\vec{\theta}}{\operatorname{argmax}} p(D | \vec{\theta})$ finden welcher die Likelihood maximiert

mithilfe der relativen Entropie (Kullback-Leibler Distanz resp. KL Divergenz) kann ML-Methode auch als die im Sinne der Entropie beste Näherung an die empirische Verteilung der Daten D erkannt werden

Optimierungsproblem, welches nur manchmal analytisch (geschlossen) lösbar ist, in vielen Fällen jedoch numerisch angegangen werden muss

für das polynomiale Regressionsproblem erhalten wir eine geschlossene Lösung:

$$L(\vec{\theta} | D) = \prod_{k=1}^{n_0} N(t_k | y(x_k, \vec{w}), \beta^{-1}) = \prod_{k=1}^{n_0} \left((2\pi)^{-\frac{1}{2}} / \beta^{\frac{1}{2}} \exp \left(-\frac{1}{2\beta} (t_k - y(x_k, \vec{w}))^2 \right) \right) \quad (1)$$

$$\ln L(\vec{\theta} | D) = -\frac{\beta}{2} \sum_{k=1}^{n_0} (t_k - y(x_k, \vec{w}))^2 + \frac{n_0}{2} (\ln(\beta) - \frac{n_0}{2} \ln(2\pi)) \quad (2)$$

$$\frac{1}{2} \sum_{k=1}^{n_0} (t_k - y(x_k, \vec{w}))^2 = \frac{1}{2} \sum_{k=1}^{n_0} \left(t_k - \sum_{j=0}^m w_j x_k^j \right)^2 = \frac{1}{2} \|\vec{t} - \underline{X} \vec{w}\|_{\mathbb{R}^{n_0}}^2 \text{ mit } X_{k,j} = x_k^{j-1} \quad (3)$$

Maximierung von (1) bz. \vec{w} ist äquivalent zu Minimierung von (3) bz. \vec{w}

ML-Methode liefert für das polynomiale Regressionsproblem die Punktschätzungen

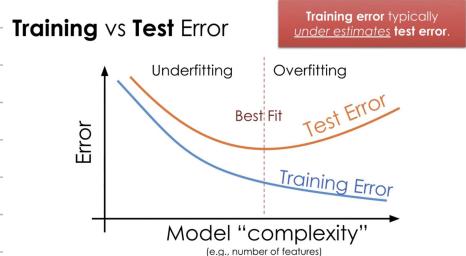
$$\vec{w}_{\text{ML}} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \vec{t} \quad \text{und} \quad \beta_{\text{ML}}^{-1} = \frac{1}{n_0} \|\vec{t} - \underline{X} \vec{w}_{\text{ML}}\|_{\mathbb{R}^{n_0}}^2$$

mit Systemmatrix $\underline{X} \in \mathbb{R}^{n_0 \times (m+1)}$ und Daten $\vec{t}, \underline{X} \in \mathbb{R}^{n_0}$

bei Implementierung \vec{w}_{ML} nicht gemäss Formel sondern direkt als Lösung des Kleinstfehler-Quadrat-Problems berechnen (mit QR- oder Singularitätenzerlegung)

Evaluation der Qualität

beim polynomialem Regressionsmodell tritt im Falle $M \approx 10$ das Phänomen overfitting auf, das dann die Datenpunkte exakt wiederzugeben vermögen kann ($E_{\text{MS}}(y_i, D) = 0$), die Verallgemeinerungsqualität jedoch sehr schlecht ist.



ML-Methode tendiert dazu ein Modell sehr fein an die Daten (höchste Rauschen) anzupassen, falls das Modell im Bezug auf die Menge an verfügbaren Daten zu flexibel ist.

bei konstanter gehaltener Menge Trainingsdaten nimmt bei ML-Methode Testfehlerfehler in Abhängigkeit von zunehmender Flexibilität typischerweise U-Form an während Trainingsfehler monoton abfällt

i. A. reduziert stets Überanpassungsproblem mit zunehmender Trainingsdatenmenge bei gleichbleibender Modellflexibilität

Modellselektion

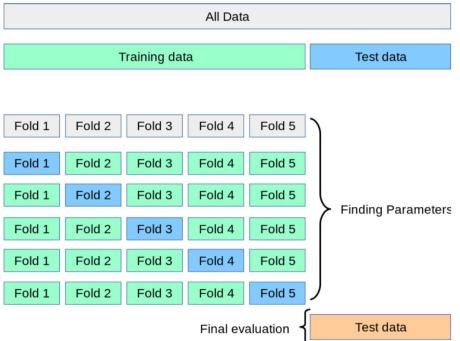
Ziel ist unter verschiedenen Modellen mit gleichbleibender Flexibilität das Modell mit der besten Verallgemeinerungsqualität auszuwählen. Dafür müsste man einen weiteren Teil der zur Verfügung stehenden Daten für die Validierung des Modells reservieren (zusätzlich zu den Testdaten für die finale Evaluation)

die Varianz der Schätzung des Verallgemeinerungsfehler ist umso kleiner je mehr Daten wir für das Evaluieren verwenden (Gesetz der grossen Zahlen). andererseits möchten wir einen möglichst grossen Anteil der Daten zum Trainieren benutzen, um möglichst für alle Modelle die Überanpassung zu verhindern

mit Hilfe der Kreuzvalidierung versucht man diesen Konflikt zu entschärfen und die verfügbaren Daten effizienter zu nutzen, was mit Redundanz bezahlt werden muss, weil alle Modelle darauf mehrmals trainiert werden müssen

bei der k-fach Kreuz-Validierung werden verfügbare Daten D in k Teile D_1, \dots, D_k zerlegt. anschließend werden für alle Teile $j \in \{1, \dots, k\}$ alle Modelle M_1, \dots, M_m auf $D \setminus D_j$ trainiert und dann auf D_j der Fehler $E(M_i, D_j)$ evaluiert. danach wird für jedes Modell der gemittelte Fehler

$$E^{k-\text{cv}}(M_i, D) = \frac{1}{k} \left(\sum_{j=1}^k E(M_i, D_j) \right)$$



als Verallgemeinerungsfehler verwendet. kann eingesetzt werden um ein einzelnes Modell zu evaluieren oder um unter verschiedenen Modellen zu selektieren. für die Validierung zum Selektion muss für die die Evaluation des schliesslich selektierten Modells eben in für die UV verwendete Testdatenset abgesondert werden dabei dass CV-selektierte Modell auf gesamten Trainingsdaten nur trainiert wird, bevor die Qualität auf den Testdaten evaluiert wird

Parameter, welche die Flexibilität/Komplexität kontrollieren, können nicht mit ML-Methode auf Trainingsdaten alleine adaptiert werden, da die ML-M. immer die flexibelsten Modelle bevorzugt und damit Überanpassung produziert wird

Maximum-a-posteriori-Methode Karte 2

besitzt Parameter um sich mithilfe einer Prior-Verteilung mit den Prior-Verteilung $p(\vec{\theta})$ repräsentiert unser aktuelles Wissen / unsere Annahmen über Parameter des Modells, bevor wir (neue) Daten beobachten

bei polynomialen Regression können wir sagen dass grosse Koeffizienten a posteriori weniger wahrscheinlich sind als kleine (da diese zu starken Oszillationen führen) indem eine Verteilung mit von $\vec{\theta} \sim \mathbb{R}^{n+1}$ konzentrierter Wahrscheinlichkeitsmasse, zur Verstärkung von spitzen Beobachtungen in Form einer Normalverteilung

$$p(\vec{w} | \alpha) = N(\vec{w} | \vec{\theta}, \alpha^{-1} \underline{E}_{m+1}) = \left(\frac{\alpha}{2\pi} \right)^{\frac{n+1}{2}} \exp\left(-\frac{\alpha}{2} \vec{w}^T \vec{c}_w\right)$$

mit Mittelwert $\vec{\theta}$ und Kovarianz $\alpha^{-1} \underline{E}_{m+1}$ wählen

um nun sowohl die Prior-Verteilung $p(\vec{\theta})$ als auch Daten $D = \{(t_k, x_k) | 1 \leq k \leq n_D\}$ miteinbeziehen können wir mit der Bayes-Formel die Posterior-Verteilung des Parameters $\vec{\theta}$ berechnen:

$$p(\vec{\theta} | D) = \frac{p(D|\vec{\theta}) p(\vec{\theta})}{p(D)} \quad \text{posterior} \propto \text{likelihood} \times \text{prior}$$

bei MAP-Methode berechnet man

$$\vec{\theta}_{MAP} = \operatorname{argmax}_{\vec{\theta}} p(\vec{\theta} | D) = \operatorname{argmax}_{\vec{\theta}} p(D|\vec{\theta}) p(\vec{\theta}) = \operatorname{argmax}_{\vec{\theta}} \prod_{k=1}^{n_D} p(t_k | \vec{\theta}) p(\vec{\theta})$$

wodurch Posterior-Verteilung $p(\vec{\theta} | D)$ des Parameters $\vec{\theta}$ gegeben, die Daten D enthalten

Evidenz oder marginale Likelihood $p(D) = \int p(D|\vec{\theta}) p(\vec{\theta}) d\vec{\theta}$, welche analytisch meistens nicht geschlossen berechenbar und deren numerische Berechnung nicht immer einfach und manchmal sehr CPU-intensiv ist, muss für MAP-Methode nicht berechnet werden

für das polynomiale Regressionsproblem erhalten wir:

$$p(\vec{t} | \vec{x}, \vec{w}, \beta) p(\vec{w} | \alpha) = \prod_{k=1}^{n_D} N(t_k | y(x_k, \vec{w}), \beta^{-1}) N(\vec{w} | \vec{\theta}, \alpha^{-1} \underline{E}_{m+1}) \quad (1)$$

$$\frac{1}{2} \sum_{k=1}^{n_D} (t_k - y(x_k, \vec{w}))^2 + \frac{\alpha}{\beta} \frac{\vec{w}^T \vec{c}_w}{2} = \frac{1}{2} \left\| \begin{pmatrix} \vec{t} \\ \vec{\theta} \end{pmatrix} - \begin{pmatrix} \vec{x} \\ \sqrt{\frac{\alpha}{\beta}} \underline{E}_{m+1} \end{pmatrix} \vec{w} \right\|^2 \quad (2) \quad \text{neg, log, etc}$$

(2) wird augmentierter (regularisierter) Summe-der-Quadrat-Fehler genannt

Maximierung von (1) bez. \vec{w} ist äquivalent zu Minimierung von (2) bez. \vec{w}

MAP-Methode liefert die Punktschätzung

$$\vec{w}_{MAP} = (\vec{x}^T \vec{x} + \lambda \underline{E}_{m+1})^{-1} \vec{x}^T \vec{t} \quad \text{mit } X_{k,j} = x_k^{j-1}$$

wobei $\lambda = \frac{\alpha}{\beta}$, $\vec{x} \in \mathbb{R}^{n_D \times (m+1)}$, $\underline{E}_{m+1} \in \mathbb{R}^{(m+1) \times (m+1)}$

wir erhalten eine durch $\alpha > 0$ parametrisierte Familie von Modellen. α (oder äquivalent λ) kontrolliert die Flexibilität des Modells wobei kleinere Werte von λ zu flexibleren Modellen führt. Wir können α (oder auch Grad d. Polyn. M) allerdings noch immer nicht auf Basis der Trainingsdaten alleine bestimmen!

Plug-In Modell: $\vec{\theta}_{PI} = \int p(\vec{\theta} | D) \vec{\theta} d\vec{\theta}$ Mittelwert ausfertig Maximieren

Bayes-Methode Karte 3

Führt auf prinzipiellen Ansatz zur Selektion der Flexibilität von Modellen (basierend auf Trainingsdaten, ohne Kreuzvalidierung)

Ursache der Überanpassung bei ML und MAP: Vernachlässigung der nach dem Lernen verbleibende Unsicherheit in den Parametern $\vec{\theta}$

BM bestimmt Posterior-Voraussage-Verteilung mit Marginalisierung über Modellparameter $\vec{\theta}$

$$p(t|D) = \int p(t|\vec{\theta}) p(\vec{\theta}|D) d\vec{\theta}$$

unter Einbezug des aktuellen Wissens (in Form der Posterior-Verteilung)

«Mischung der Modelle $p(t|\vec{\theta})$ zu allen möglichen Parametern wobei jeder Parameter $\vec{\theta}$ entsprechend der Posterior-Verteilung gewichtet wird»

$p(t|D)$ berücksichtigt auch die nach Beobachtung der Daten verbleibende Unsicherheit über die Parameter $\vec{\theta}$

für das polynomiale Regressionsproblem erhalten wir:

$$p(t|x, \vec{t}, \vec{x}, \alpha, \beta) = \int p(t|x, \vec{w}, \beta) p(\vec{w}|\vec{t}, \vec{x}, \alpha, \beta) d\vec{w}$$

wobei $p(t|x, \vec{w}, \beta) = N(t|y(x, \vec{w}), \beta^{-1})$, $p(\vec{w}|\vec{t}, \vec{x}, \alpha, \beta) = \frac{p(\vec{t}|\vec{x}, \vec{w}, \beta) p(\vec{w}|\alpha)}{\int p(\vec{t}|\vec{x}, \vec{w}, \beta) p(\vec{w}|\alpha) d\vec{w}}$, $p(\vec{w}|\alpha) = N(\vec{w}|\vec{0}, \alpha^{-1}E)$

für dieses einfaches Modell sind die Integrale analytisch geschlossen barthalbar und wir erhalten:

$$p(t|x, \vec{t}, \vec{x}, \alpha, \beta) = N(t|m(x), v(x))$$

mit $m(x) = \vec{b}(x)^T \vec{m}$, $v(x) = \beta^{-1} + \vec{b}(x)^T \vec{b}(x)$, $\vec{b}(x) = (1, x, x^2, \dots, x^m)^T$, $\vec{m} = \alpha E + \beta \vec{x}^T \vec{x}$, $\vec{v} = \beta E - \vec{x}^T \vec{x}$

diese Voraussage-Verteilung hängt noch immer von den Hyper-Parametern α, β, M ab. nun unter dieser Familie das beste zu wählen, gibt es im Rahmen der Bayes-Methode des Ansatz die Unsicherheit in den (Hyper-)Parametern wieder explizit mittels einer Prior-Verteilung zu bewerten.

bspw. kann bei einer Familie von parametrischen Modellen $p(t|\vec{\theta}, M)$ mit dazugehörigen Prior-Verteilungen $p(\vec{\theta}|M)$ die auf den Trainingsdaten $D = \{t_k | 1 \leq k \leq n_D\}$ berechnete Evidenz oder marginale Likelihood

$$p(D|M) = \int p(D|\vec{\theta}, M) p(\vec{\theta}|M) d\vec{\theta} = \int_{k=1}^{n_D} \prod_{k=1}^{n_D} p(t_k|\vec{\theta}, M) p(\vec{\theta}|M) d\vec{\theta}$$

zur Schätzfunktion von Modellen verwendet werden (nur mithilfe der Trainingsdaten alleine!)

Modelle für diskrete Daten und die Bayes-Methode

im Fall von diskreten Daten ist Bayes-Methode berechenbar

Endliche Familien diskreter Verteilungen

Annahme: esse Θ^* von endlich vielen mögl. den Hypothesen $\Theta \in \{1, \dots, r\}$ generiert diskrete Daten $x \in \mathcal{X}$ mit Verteilung $p(x|\Theta^*)$ und wir möchten basierend auf Testdaten $D = \{x_i | 1 \leq i \leq n\}$ eine von zukünftigen Werten abhängende Größe voraussagen
→ gelingt am besten falls wir $p(x|\Theta^*)$ rekonstruieren können

falls Werte nur endlich viele Werte $| \mathcal{X} | = m$ aufweisen können wir diese notwendigen $\mathcal{X} = \{1, \dots, m\}$ und Verteilungen aller Hypothesen in einer $m \times r$ -Matrix $M = p(x|\Theta)$ mit den Einträgen $M_{k,l} = p(x=k|\Theta=l)$ kodieren

Likelihood und ML-Methode

für die Likelihood der Daten D unter der Hypothese Θ gilt $p(D|\Theta) = \prod_{i=1}^n p(x_i|\Theta) = \prod_{k=1}^m p(k|\Theta)^{n_k}$

wobei $n_k = |\{i \in \{1, \dots, n\} | x_i = k\}|$ die Anzahl der Datenspunkte mit Wert $k \in \{1, \dots, m\}$ bezeichnet
die Anzahlen $\vec{n} = (n_1, \dots, n_m)$ werden erschöpfende Statistik genannt

mit der ML-Methode können die Hypothese Θ_M wählen welche das Produkt maximiert und es ist nicht immer bei der ML-Methode das Risiko einer Überanpassung bestellt

Prior/Posterior - Verteilung und MAP- und Bayes-Methode

falls Kenntnis über Hypothesen $\Theta = \{1, \dots, r\}$ vorhanden ist, können wir dies auch hier in Form einer diskreten Prior-Verteilung $p(\Theta) = (\pi_1, \dots, \pi_r)$ über die Parameter Θ des parametrisierten Modells $p(x|\Theta)$ vernutzen und die Voraussage-Verteilung mittels Marginalisierung berechnen:

$$\text{Voraussage-Verteilung: } p(x) = \sum_{\Theta} p(x|\Theta)p(\Theta)$$

falls Prior-V. $p(\Theta)$ und Daten D vorliegen kann man Posterior-V. berechnen:

$$\text{Posterior-Verteilung: } p(\Theta|D) = \frac{p(D|\Theta)p(\Theta)}{p(D)} = \frac{\prod_{k=1}^m p(k|\Theta)^{n_k} p(\Theta)}{\sum_{l=1}^r (\prod_{k=1}^m p(k|l)^{n_k} p(l))}$$

MAP-Methode verwendet Maximum a posteriori der Posterior-V. als Basis für Voraussagen: $p(x|\Theta_{MAP})$

Bayes-Methode benutzt induzierte Posterior-Voraussage-Verteilung: $p(x|D) = \sum_{\Theta} p(x|\Theta)p(\Theta|D)$

Sequentielles Lernen

mittels Induktion kann man zeigen, dass wenn bei der Bayes-Methode auf das gleiche Ergebnis kommt wenn man folgendermassen vorgeht: D in n Pakete D_1, \dots, D_n zerlegen, mit Prior-Verteilung $p(\Theta)$ und D_1 Posterior-Verteilung $p(\Theta|D_1)$ berechnen, mit $p(\Theta|D_1)$ als neuer Prior-V. und D_2 wieder dazugehörige Posterior-Verteilung berechnen welche genau $p(\Theta|D_1, D_2)$ entspricht, usw.

im Grenzzustand von unendlich viel Daten liegt die Posteriorverteilung auf den MAP-Parameter konzentriert, welches dann zusammenfällt mit dem ML-Parameter zusammenfällt, weil der Einfluss der Prior-Verteilung vom Gewicht der Daten (über Likelihood) vollständig überdeckt wird
→ bei wenigen Daten kann Unterschied der 3 Lernmethoden jedoch sehr gross sein!

Bernoulli-Modelle

→ Basisfond für flexible & konsolidierte Modelle

Bernoulli-Verteilung: $p(x|\theta) = \text{Bern}_\theta(x) = \theta^x (1-\theta)^{1-x}$, $x \in \{0, 1\}$, $\theta \in [0, 1]$

erschöpfende Statistik der Datensetze $D = \{x_1, \dots, x_n\}$: $n_1 = |\{x \in D | x=1\}|$, $n_0 = |\{x \in D | x=0\}|$

$$L(\theta|D) = p(D|\theta) = \theta^{n_1} (1-\theta)^{n_0} \Rightarrow \ln(L(\cdot)) \rightarrow \text{maximieren} \Rightarrow \hat{\theta}_{ML} = \frac{n_1}{n_1 + n_0}$$

$$\text{Posterior-V.}: p(\theta|D) = \frac{p(D|\theta)p(\theta)}{p(D)} = \frac{\theta^{n_1} (1-\theta)^{n_0} p(\theta)}{\int_0^1 \theta^{n_1} (1-\theta)^{n_0} p(\theta) d\theta}$$

Definition konjugierte Prior-V.-Verteilung:

Familie K von Verteilungen $p(\theta)$ über die Parameter $\vec{\theta}$ eines prob. Modells $p(x|\theta)$ heißt konjugiert zu $p(D|\theta)$ falls für alle Daten D und alle Verteilungen $p \in K$ auch induzierte Posterior-V. $p(\theta|D) \in K$. Eine Verteilung der Familie ist, also $p(\theta|D) \in K$.

Likelihood der Bern.-V.:

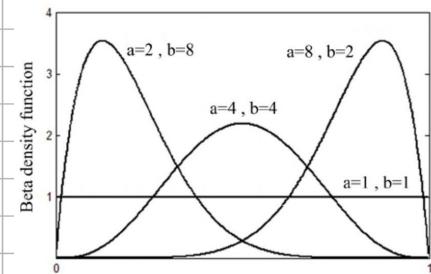
Beta-Prior-V. $p(\theta) = \text{Beta}_{a,b}(\theta)$ ist konjugiert zur Likelihood der Bernoulli-V. $p(x|\theta) = \text{Bern}_\theta(x)$ und für Daten D mit Statistik (n_1, n_0) erhalten wir die Posterior-V. $p(\theta|D) = \text{Beta}_{a+n_1, b+n_0}(\theta)$

Beta-Verteilung:

$$\text{Beta}_{a,b}(\theta) = \frac{\theta^{a-1} (1-\theta)^{b-1}}{B(a,b)} \quad \text{wobei } \theta \in [0, 1], \quad a > 0, b > 0$$

Normalisierung im Nenner heißt Beta-Funktion:

$$B(a,b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)} \quad \text{mit } \Gamma(x) = \int_0^\infty t^{x-1} e^{-t} dt$$



Kenngrößen:

$$\text{mean} = \frac{a}{a+b}, \quad \text{variance} = \frac{a \cdot b}{(a+b)^2 (a+b+1)}, \quad \text{mode} = \frac{a-1}{a+b-2} \quad (\text{exists only if } a>1, b>1), \quad \text{Stärke} = a+b$$

mit $p(\theta) = \text{Beta}_{a,b}(\theta)$ als Prior-V. erhält man für die Posterior-V.:

$$p(\theta|D) = \frac{1}{p(D)} p(D|\theta) p(\theta) = \frac{1}{p(D)} \theta^{n_1} (1-\theta)^{n_0} \frac{1}{B(a,b)} \theta^{a-1} (1-\theta)^{b-1} \propto \text{Beta}_{a+n_1, b+n_0}$$

Zwei Wahrscheinlichkeitsdichten, die proportional zueinander sind müssen jedoch normalisiert identisch sein (wegen Integral $\equiv 1$) ■

prior $\xrightarrow{\text{data}}$ posterior dieser Übergang legt es nahe die Parameter a, b als Pseudo-Beschreibungen $\text{Beta}_{a,b} \xrightarrow{n_1, n_0} \text{Beta}_{a+n_1, b+n_0}$ von 1-Werten respektive 0-Werten zu interpretieren

Stärke $a+b$ der Prior-V. ist Mass für relativen Einfluss der Prior-V. auf Posterior-V., je stärker der Prior (bei gleich viel Daten) desto mehr gleicht sich Posterior-V. der Prior-V. an

Variance nimmt bei Übergang prior \xrightarrow{D} posterior i.A. ab, was als Zunahme des Wissens über Parameter θ interpretiert werden kann

Voraussage-Verteilung liefert $p(x|a,b) = \int \text{Bern}(x|\theta) \text{Beta}_\theta(\theta|a,b) d\theta = \text{Bern}(x|\theta_{\text{prior}})$ wobei:

$$\theta_{\text{prior}} = p(x=1|a,b) = \int \text{Bern}(x=1|\theta) \text{Beta}_\theta(\theta|a,b) d\theta = \int \theta \text{Beta}_\theta(\theta|a,b) d\theta = \text{Mean}(\text{Beta}_{a,b}) = \frac{a}{a+b}$$

$$\text{MAP-Methode wählt } \hat{\theta}_{\text{MAP}} = \underset{\theta}{\operatorname{argmax}} p(x|D) = \text{Mode}(\text{Beta}_{a+n_1, b+n_0}) = \frac{a+n_1-1}{a+n_1+b+n_0-2}$$

$$\text{Bayessche Voraussage-Verteilung liefert } \text{Bern}(x|\theta_{\text{mean}}) \text{ mit } \theta_{\text{mean}} = \frac{a+n_1}{a+n_1+b+n_0}$$

Bayessche Methode interpoliert zwischen Prior-Voransage (nur Prior, keine Daten) und der ML-Voransage (kein Prior, nur Daten):

$$\theta_{\text{mean}} = \lambda \theta_{\text{prior}} + (1-\lambda) \theta_{\text{ML}} \quad \text{mit Mischungsparameter } \lambda = \frac{1}{1 + \frac{n_0 + n_1}{a + b}}$$

Voransage einer Sequenz

falls mehrere, unklare Werte voransgesagt werden sollen \rightarrow Binomialverteilung

- θ_{ML} und θ_{MAP} können direkt in $Bin_{m,\theta}(y) = \binom{m}{y} \theta^y (1-\theta)^{m-y}$ eingesetzt werden
- mit Bayes ergibt sich $\int Bin(y|m,\theta) \text{Beta}(\theta|c,d) d\theta = \binom{m}{y} \frac{B(y+c, m-y+d)}{B(c,d)} =: \text{BetaBin}(y|c,d,m)$

Multitask Bernoulli:

Ziel: r Parameter θ_k von r ähnlich strukturierten, nicht notwendigerweise gleich großen Populationen mit Datensätzen D_k mit statistischen ($n_{k,1}, n_{k,0}$) bestimmen

$$p(D_k, \theta_k | a, b) = p(D_k | \theta_k) p(\theta_k | a, b) = \theta_k^{n_{k,1}} (1-\theta_k)^{n_{k,0}} \text{Beta}(\theta_k | a, b) \quad \text{Modell}$$

$$p(\theta_k | D_k, a, b) = \text{Beta}(\theta_k | a + n_{k,1}, b + n_{k,0}) \quad \text{Posterior-V.}$$

hierarchischer Ansatz:

$\xrightarrow{\text{Prior}}$ $\xrightarrow{\text{Hyperprior: Prior über Prior}}$

$$p(D, \vec{\theta}, a, b) = p(D | \vec{\theta}) p(\vec{\theta} | a, b) p(a, b) = \left(\prod_{k=1}^r \theta_k^{n_{k,1}} (1-\theta_k)^{n_{k,0}} \text{Beta}(\theta_k | a, b) \right) p(a, b)$$

$$\text{mit } D = \bigcup_{k=1}^r D_k, \quad \vec{\theta} = (\theta_1, \dots, \theta_r)$$

gemeinsame Verteilung

nach Bayes-Methode können wir mittels Marginalisierung die Posterior-V. $p(\theta_k | D)$ berechnen, welche analytisch leider schwierig zu berechnen ist

Vorteil ist, dass statistische Stärke von Populationen mit vielen Daten auf Populationen mit wenigen Daten transferiert werden kann, da die θ_k nun abhängig sind

Multinomelli-Modelle

Kategorische / Multinomelli-Verteilung: $p(k|\vec{\theta}) = \text{Cat}_{\vec{\theta}}(k) = \theta_k$ Verallgemeinerung Bernoulli-V.

mit $k \in \{1, \dots, K\}$, $\vec{\theta} = (\theta_1, \dots, \theta_K) \in \mathbb{R}_+^K$ (welcher $\sum \theta_k = 1$ erfüllt)

erschöpfende Statistik der Datensammlung $D = \{x_1, \dots, x_n\}$: $\vec{n} = (n_1, \dots, n_K)$
mit $n_k = |\{x \in D \mid x=k\}|$

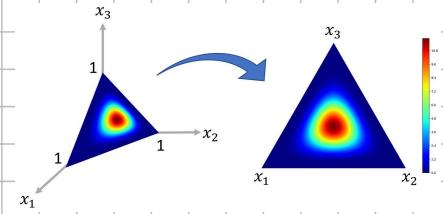
$$L(\vec{\theta}|D) = p(D|\vec{\theta}) = \prod_{k=1}^K \theta_k^{n_k} \Rightarrow \ln(\cdot) \Rightarrow \text{maximieren mit Lagrange-Multiplikatoren}$$

$$\Rightarrow \vec{\theta}_{ML} = \frac{\vec{n}}{\|\vec{n}\|_1} = \frac{(n_1, \dots, n_K)}{n_1 + \dots + n_K} \quad \text{relative Häufigkeit}$$

Dirichlet-Verteilung:

$$\text{Dir}_{\vec{\alpha}}(\vec{\theta}) = \frac{\prod_{k=1}^K \theta_k^{\alpha_k-1}}{B(\vec{\alpha})} = \frac{\prod_{k=1}^K \theta_k^{\alpha_k-1} (1 - \sum_{k=1}^K \theta_k)^{\alpha_K-1}}{B(\vec{\alpha})} \quad \text{Verallgemeinerung Beta-Verteilung}$$

mit $B(\vec{\alpha}) = \prod_{k=1}^K \Gamma(\alpha_k) / \Gamma(\sum_{k=1}^K \alpha_k)$ Verallgemeinerung Beta-Funktion



Kenngrößen:

$$\text{mean}_k = \frac{\alpha_k}{\|\vec{\alpha}\|_1}, \quad \text{variance}_k = \frac{\text{mean}_k(1-\text{mean}_k)}{\|\vec{\alpha}\|_1 + 1}, \quad \text{mode}_k = \frac{\alpha_k - 1}{\|\vec{\alpha}\|_1 - K} \quad (\text{existiert nur falls } \alpha_k > 1 \forall k), \quad \text{Stärke} = \|\vec{\alpha}\|_1$$

Dirichlet-Prior-V. $p(\vec{\theta}) = \text{Dir}_{\vec{\alpha}}(\vec{\theta})$ ist konjugiert zur Likelihood der Multinomelli-V. $p(x|\vec{\theta}) = \text{Cat}_{\vec{\theta}}(x)$ und für Daten D mit Statistik $\vec{n} = (n_1, \dots, n_K)$ erhalten wir die Posterior-V. $p(\vec{\theta}|D) = \text{Dir}_{\vec{\alpha}+\vec{n}}(\vec{\theta})$

$$\vec{\theta}_{MAP} = \frac{\vec{\alpha} + \vec{n} + \vec{l}_K}{\|\vec{\alpha} + \vec{n}\|_1 - K}$$

Voraussage-Verteilung: $\text{Cat}(x | \vec{\theta}_{\text{prior}})$ wobei $\vec{\theta}_{\text{prior}} = \frac{\vec{\alpha}}{\|\vec{\alpha}\|_1}$

$$\vec{\theta}_{\text{mean}} = \frac{\vec{\alpha} + \vec{n}}{\|\vec{\alpha} + \vec{n}\|_1}$$

Naive-Bayes-Klassifikator



Annahmen: $t \in \{1, \dots, C\}$, $\vec{x} \in \mathcal{X} = \{x_1, \dots, x_m\} \times \dots \times \{x_1, \dots, x_m\}$, $D = \{(t_k, \vec{x}_k) \mid 1 \leq k \leq n_D\}$

Anzahl verschiedener möglicher \vec{x} -Werte: $|\mathcal{X}| = \prod_{m=1}^M K_m$

Klassen-Prior-Verteilung: $p(t | \vec{\theta}_0) = \text{Cat}_{\vec{\theta}_0}(t)$ modelliert Wissen über Verteilung der Klassenzugehörigkeit

auf Klasse t bedingte kategoriale Verteilung: $p(\vec{x} | t, \vec{\theta}_t)$

gemeinsame Verteilung der (Ausgabe, Eingabe)-Paare: $p(t, \vec{x} | \vec{\theta} = (\vec{\theta}_0, \vec{\theta}_1, \dots, \vec{\theta}_C)) = p(\vec{x} | t, \vec{\theta}_t) p(t | \vec{\theta}_0)$

trainiertes Modell: $p(t | \vec{x}_n) = \frac{p(t, \vec{x}_n | D)}{p(\vec{x}_n | D)} = \frac{p(t, \vec{x}_n | D)}{\sum_{t \in C} p(t, \vec{x}_n | D)}$

für die $p(\vec{x} | t, \vec{\theta}_t)$ mit $t = 1, \dots, C$ brauchen wir $C \cdot (M \cdot K_m - 1)$ unabhängige Komponenten in $[0, 1]$
→ dies übersteigt schon bei einfacheren Problemen die Eddington-Number von $\approx 10^{80}$ ($\#\rho^+$ im Universum)

daher behandelt man die einzelnen Komponenten x_m als unabhängig ("naiv"):

$$p(\vec{x} | t, \vec{\theta}_t) = \prod_{m=1}^M p(x_m | t, \vec{\theta}_{t,m}) = \prod_{m=1}^M \text{Cat}_{\vec{\theta}_{t,m}}(x_m)$$

Anzahl $\sum_{m=1}^M (K_m - 1)$ unabh. Parametern $\vec{\theta}_t = (\vec{\theta}_{t,1}, \dots, \vec{\theta}_{t,M})$ ist linear in Dimension M des Eingabekontexts

Naive-Bayes-Klassifikator:

$$p(t, \vec{x} | \vec{\theta}) = p(\vec{x} | t, \vec{\theta}_t) p(t | \vec{\theta}_0) = \prod_{m=1}^M p(x_m | t, \vec{\theta}_{t,m}) p(t | \vec{\theta}_0) = \prod_{m=1}^M \text{Cat}_{\vec{\theta}_{t,m}}(x_m) \cdot \text{Cat}_{\vec{\theta}_0}(t)$$

mit Parametern $\vec{\theta} = (\vec{\theta}_0, \vec{\theta}_1, \dots, \vec{\theta}_C) = (\vec{\theta}_{0,1}, \vec{\theta}_{0,2}, \dots, \vec{\theta}_{0,M}, \dots, \vec{\theta}_{C,1}, \dots, \vec{\theta}_{C,M})$

ML-Methode:

$$\begin{aligned} p(D | \vec{\theta}) &= \prod_{(t, \vec{x}) \in D} (p(\vec{x} | t, \vec{\theta}_t) p(t | \vec{\theta}_0)) = \prod_{(t, \vec{x}) \in D} p(\vec{x} | t, \vec{\theta}_t) \prod_{(t, \vec{x}) \in D} p(t | \vec{\theta}_0) = \prod_{t=1}^C \prod_{m=1}^M \prod_{\vec{x} \in D_t} \text{Cat}_{\vec{\theta}_{t,m}}(x_m) \cdot \prod_{(t, \vec{x}) \in D} \text{Cat}_{\vec{\theta}_0}(t) \\ \Rightarrow \vec{\theta}_{ML,0} &= \vec{n}_0 / \| \vec{n}_0 \|_1, \quad \vec{\theta}_{ML,t,m} = \vec{n}_{t,m} / \| \vec{n}_{t,m} \|_1 \end{aligned}$$

Bayessche Methode:

$$\begin{aligned} \text{Prior } p(\vec{\theta} | \vec{\alpha}) &= \prod_{t=1}^C \prod_{m=1}^M \text{Dir}_{\vec{\alpha}_{t,m}}(\vec{\theta}_{t,m}) \cdot \text{Dir}_{\vec{\alpha}_0}(\vec{\theta}_0) \text{ ist proportional zu Likelihood d. NBK} \\ \Rightarrow \text{Posterior: } p(\vec{\theta} | \vec{\alpha} + \vec{n}) &= \prod_{t=1}^C \prod_{m=1}^M \text{Dir}_{\vec{\alpha}_{t,m} + \vec{n}_{t,m}}(\vec{\theta}_{t,m}) \cdot \text{Dir}_{\vec{\alpha}_0 + \vec{n}_0}(\vec{\theta}_0) \end{aligned}$$

Modellierung beamföhliger Lernprobleme

Verlustmasse und end-zu-end Lernen

MSE nicht tolerant gegen Ausreißer resp. nicht optimal für bedingte Ausgabeverteilungen mit heavy tails
→ falls man grosse Abweichungen weniger stark sanktionieren möchte, kann man auch mittlere absolute Abweichung als Verallgemeinerungsfehler wählen

end-zu-end: direkt Voraussagefunktion lernen (nicht Verteilung der Ausgabe gegeben Ergebnis)

Verallgemeinerungsproblem

Konstruktion eines prob. Modells $p_m(\vec{t}^* | \vec{x}, \vec{\theta})$ für die den Daten zugrundeliegende Verteilung $p(\vec{t}^* | \vec{x})$

Unterschiede zum end-zu-end-Arbeitsweise:

- i.A. einfacher vorhergesagtes Wissen / Hypothesen über Problem kontrolliert und direkt ins Modell einfließen zu lassen
- erlaubt speziell entdeckte Algorithmen für Dichte-Schätzung anstelle von allg. Optimierung ins Spiel zu bringen
- Verlustmasse L geht nicht ins Verallgemeinerungsproblem ein

Entscheidungsproblem

basiert auf gelerntem prob. Modell p_m die bez. Verlustmasse beste Voraussage-Funktion finden:

$$f_{m,\text{opt}}(\vec{x}) = \underset{\vec{p}}{\operatorname{argmin}} \int L(\vec{t}, \vec{p}) p_m(\vec{t} | \vec{x}) d\vec{t}$$

unter Annahme, dass Modell p_m exakt ist, ist dies die beste Voraussage-Funktion

weitere Konsequenzen der Aufteilung

Flexibilität, Modularität

Lineare Modelle für Regressionsprobleme

häufigster Bestandteil komplexerer Modelle → „Arbeitspferd des m.L.“

Familie von diskritiven, parametrischen, generalistischen Modellen

(lineares) Regressionsmodell:

> hier als unabhängig von \vec{x} angenommen (homoskedastische Regression)

$$p(t|\vec{x}, \vec{w}, \beta) = N(t | \vec{\phi}(\vec{x})^T \vec{w}, \beta^{-1}) \text{ mit } \vec{\phi}(\vec{x})^T \vec{w} = \sum_{k=0}^m w_k \phi_k(\vec{x})$$

mit adaptiven Parametern $\vec{w} \in \mathbb{R}^{m+1}$, $\beta > 0$

vielfach verwendete generische Familien von Basisfunktionen: Monome, Spline-Funktionen, Gaußsche Funktionen (meist exakt kodiziert wie Splines), Fourier-Basen (idealisiert im Frequenzbereich), Wavelet-Basen (idealisiert in zeit- und Frequenzbereich)

klassische Regression: $\phi_0(\vec{x}) = 1$, $\phi_k(\vec{x}) = x_k \Rightarrow \vec{\phi}(\vec{x})^T \vec{w} = \sum_{k=1}^m w_k x_k + w_0$ Ausgabe hängt affin von Eingabe ab

eine Familie von Basisfunktionen induziert eine Transformation $\vec{x} \rightarrow \vec{\phi}(\vec{x}) = [\phi_0(\vec{x}), \phi_1(\vec{x}), \dots, \phi_m(\vec{x})]^T$ des Eingaberaums ⇒ jedes lineare Modell kann als klassische Regression auf dem Raum der transformierten Eingaben $\vec{\phi}(\vec{x})$ aufgefasst werden resp. Transformation gefolgt von klassischer Regression

transformierte Eingabegrößen $\vec{\phi}(\vec{x})$ nennt man auch **Features**

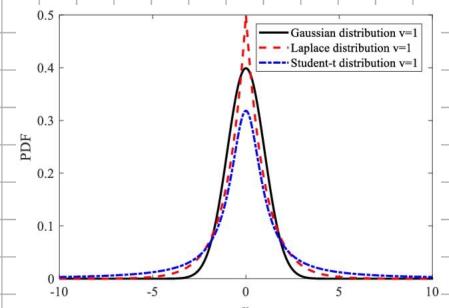
Ziel der Feature Extraction ist lineare Beziehung der Ausgabe aus Features zu ermöglichen/herzuführen

alternative Ausgabeverteilungen

viele der geschlossenen Formeln hängen von Annahme der Normalverteilung als Ausgabeverteilung ab

zur Rechtfertigung kann man die Form des zentralen Grenzwertsatzes oder auch Prinzip der Entropie-Maximierung unter Vorgabe der ersten zwei Momente herbeiziehen

N ist jedoch nicht robust gegen Ausreißer, in solchen Fällen sollte mit Vorteil eine andere Verteilung mit einer Wahrscheinlichkeitsmasse in Einklang, z.B. Laplace- oder Student-t-Verteilung, gewählt werden



Standardisierung der Daten

Zentrierung der Ausgabe und Features in Null zusammen mit Skalierung der Varianz der Features auf Eins

→ begünstigt Vergleichbarkeit der Sensitivität der Parameter \vec{w}

→ ermöglicht Wahl einer uniformen Stärke der Prior-Verteilung in allen Dimensionen

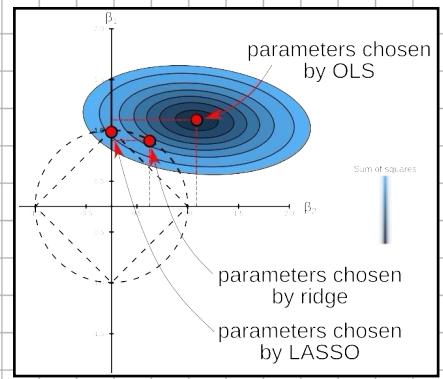
Maximum Likelihood Lernen

$$\text{Log-Likelihood} = \ln p(\vec{t} | \vec{x}, \vec{w}, \beta) = \sum_{i=1}^N \ln \mathcal{N}(t_i | \vec{\phi}(\vec{x}_i)^T \vec{w}, \beta^{-1}) \\ = \frac{N}{2} \ln(\beta) - \frac{N}{2} \ln(2\pi) - \frac{\beta}{2} \sum_{i=1}^N (t_i - \vec{\phi}(\vec{x}_i)^T \vec{w})^2$$

Maximierung der Log-Likelihood ergibt

$$\vec{w}_{ML} = (\underline{\Phi}(\vec{x})^T \underline{\Phi}(\vec{x}))^{-1} \underline{\Phi}(\vec{x})^T \vec{t}$$

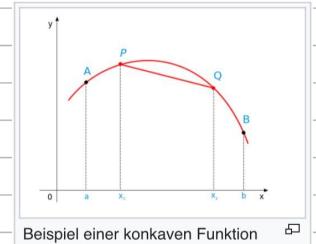
$$\beta^{-1} = \frac{1}{N} \| \vec{t} - \underline{\Phi}(\vec{x}) \vec{w}_{ML} \|^2 \quad \text{mittlere quadratische Abweichung der Lösung des Least-Squares-Problems}$$



wobei für die Systemmatrix (Designmatrix) $\Phi_{j,k}(\vec{x}) = \Phi_{k,1}(\vec{x}_j)$ gilt, also

$$\underline{\Phi}(\vec{x}) = \begin{bmatrix} \Phi_0(\vec{x}_1) & \Phi_1(\vec{x}_1) & \cdots & \Phi_m(\vec{x}_1) \\ \vdots & \vdots & & \vdots \\ \Phi_0(\vec{x}_N) & \Phi_1(\vec{x}_N) & \cdots & \Phi_m(\vec{x}_N) \end{bmatrix} \in \mathbb{R}^{N \times (m+1)}$$

Likelihood bei Linearer Modellen konkav
 → Vorteil gegenüber komplexeren Modellen!



Maximum Posterior Lernen

wir betrachten β als externen Parameter und wählen mithilfe indirekt durch Konjuguität

$$p(\vec{w} | \vec{\alpha}) = \prod_{i=0}^m N(w_i | 0, \alpha_i^{-1}) \quad \text{Prior}$$

als Prior-Verteilung auf dem Raum der Gewichte $\vec{w} \in \mathbb{R}^{m+1}$

Kovarianzmatrix ist diagonal gewählt, weil Gewichte der einzelnen Basisfunktionen nicht a priori korreliert sein sollen und der Präzisions-Hyperparameter $\alpha \in \mathbb{R}_{+}^{m+1}$ kontrolliert Stärke, mit welcher Gewichte im Nähe von $\vec{\beta}$ gezwungen werden

$$\text{Log von Likelihood} \times \text{Prior} = \ln(p(t_i | \vec{x}, \vec{w}, \beta) p(\vec{w} | \vec{\alpha})) = -\frac{\beta}{2} \sum_{j=1}^N (t_j - \phi(\vec{x}_j)^T \vec{w})^2 - \frac{1}{2} \sum_{k=0}^m \alpha_k w_k^2$$

Maximierung davon in \vec{w} liefert

$$\vec{w}_{\text{MAP}} = (\Phi(\vec{x})^T \Phi(\vec{x}) + \underline{\alpha}^{-1})^{-1} \Phi(\vec{x})^T \vec{t}$$

mit $\underline{\alpha} = \vec{\alpha} / \beta$

Klassische Ridge-Regression: MAP-Methode mit eindimensionalen Prior-V. mit Kovarianz $\alpha^{-1} E_{m+1}$

bei gegebener Präzision β kontrolliert der Hyper-Parameter $\vec{\alpha}$ (oder äquivalent $\underline{\alpha}$) die Flexibilität des Modells und dieser Wert ist oft numerisch optimiert (z.B. mit Kruz.-Validation)

für $\vec{\alpha} \rightarrow \vec{0}$ konvergiert $\vec{w}_{\text{MAP}} \rightarrow \vec{w}_{\text{ML}}$

Lasso-Regression: faktorierte Laplace-V. als Prior

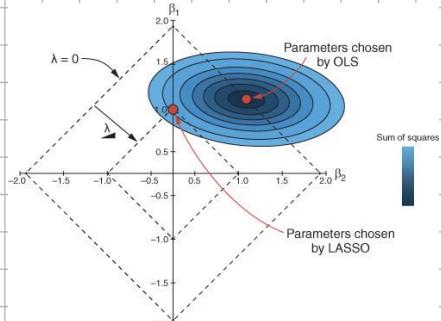
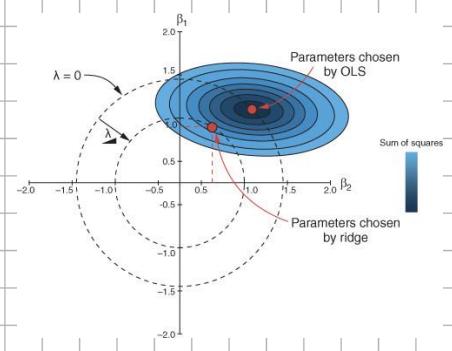
→ besser als Ridge falls aus grosser Familie von Basisfunktionen die besten ausgewählt werden sollen, da sie die Eigenschaft hat dass das damit gebaute Modell viele Koeffizienten exakt bei Null hat (sparsity erfordert)

→ man kann den die Stärke der Prior-V. kontrollierenden Faktor λ als Lagrange-Multiplikator eines Minimierungsproblems mit der Nebenbedingung $\sum_{i=0}^m w_i = C$ auffassen und dann aufgrund der Geometrie der Nebenbedingung argumentieren, dass in der Lösung \vec{w} viele Nullstellen vorhanden sollten

falls man Gauss-Prior zur Selektion von Features verwenden möchte, wird die allgemeine Form mit gleich vielen Hyperparametern $\vec{\alpha}$ wie Gewichten \vec{w} und nicht die Ridge-Regression benötigt

→ bedingt, dass auch die Hyperparameter $\vec{\alpha}$ adaptiert werden

→ automatische-Relevant-Bestimmung (ARD)



Bayesianes Lernen \rightarrow für sequentielles Lernen sowie permanente Adaptivität geeignet

Likelihoodhood: $p(\vec{t} | \vec{w}, \vec{x}, \beta) = \mathcal{N}(\vec{t} | \underline{\Phi}(\vec{x})\vec{w}, \beta^{-1}\mathbb{E}) = \prod_{j=1}^n \mathcal{N}(t_j | \vec{\Phi}(\vec{x}_j)^T \vec{w}, \beta^{-1})$

Prior: $p(\vec{w} | \vec{m}_0, \Sigma_0) = \mathcal{N}(\vec{w} | \vec{m}_0, \Sigma_0)$

Posterior:

$$p(\vec{w} | \vec{t}, \vec{x}, \eta) = \mathcal{N}(\vec{w} | \vec{m}_n, \Sigma_n)$$

mit $\vec{m}_n = \Sigma_n (\beta \underline{\Phi}(\vec{x})^T \vec{t} + \Sigma_0^{-1} \vec{m}_0)$, $\Sigma_n = (\Sigma_0^{-1} + \beta \underline{\Phi}(\vec{x})^T \underline{\Phi}(\vec{x}))^{-1}$

marginal Likelihood:

$$p(\vec{t} | \vec{x}, \beta, \vec{m}_0, \Sigma_0) = \mathcal{N}(\vec{t} | \underline{\Phi}(\vec{x})\vec{m}_0, \beta^{-1}\mathbb{E} + \underline{\Phi}(\vec{x})\Sigma_0 \underline{\Phi}(\vec{x})^T)$$

Posterior-Vorhersage-Verteilung für $p(t | \vec{x}, \vec{w}, \beta) = \mathcal{N}(t | \vec{\Phi}(\vec{x})^T \vec{w}, \beta^{-1})$ und $p(\vec{w} | \vec{m}, \Sigma) = \mathcal{N}(\vec{w} | \vec{m}, \Sigma)$:

$$p(t | \vec{x}, \beta, \vec{m}, \Sigma) = \int p(t | \vec{x}, \vec{w}, \beta) p(\vec{w} | \vec{m}, \Sigma) d\vec{w} = \mathcal{N}(t | m(\vec{x}), s(\vec{x}))$$

mit $m(\vec{x}) = \vec{\Phi}(\vec{x})^T \vec{m}$ und $s(\vec{x}) = \beta^{-1} + \vec{\Phi}(\vec{x})^T \Sigma \vec{\Phi}(\vec{x})$

die Bayessche Vorhersage-Varianz $s(\vec{x}) = \beta^{-1} + \vec{\Phi}(\vec{x})^T \Sigma \vec{\Phi}(\vec{x})$ wird kleiner, wenn Eingabe \vec{x} nahe bei Eingabevektoren der Trainingsdaten liegt und wird allgemein kleiner wird, nach Anzahl N der Trainingsdaten größer wird

falls \vec{x} an Rändern des Supports aller Basisfunktionen liegt \rightarrow Kovarianz bleibt weil alle Features fast null sind, auch wenn \vec{x} weit weg von Trainingsdaten (schwer zu beobachtendes Problem von diskriminativen Modellen)

Grenzen Linearer Modelle

Basisfunktionen lassen sich für hochdimensionale Eingaben \vec{x} nicht einfach generisch wählen, weil (vereinfacht gesagt) der Träger jeder in der euklidischen Norm lokalisierter Basisfunktion nicht in der Nähe der Daten liegt

motiviert durch unsere Erfahrung/Intuition könnte man versuchen ein allgemeines, generisches Regressionsmodell zu konstruieren, indem man eine geeignete flexible Familie von Basisfunktionen so wählt, dass deren Träger den Raum der Eingaben \vec{x} gleichmäßig überdecken

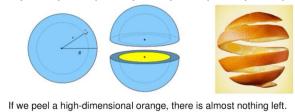
Falls Eingaben \vec{x} in zentrierten, d-dimensionalen Würfel enthalten sind wie dazu jede Achse gleichmäßig ϵ Stücke strecken, ist Anzahl Basis-Funktionen $\propto h^d$ und somit hängt auch die benötigte Menge an Daten exponentiell von d ab (\rightarrow keine gute Idee)

Curse of Dimensionality:

wir betrachten Gleichverteilung auf d-dimensionalen Einheitsball $B = \{\vec{x} \in \mathbb{R}^d \mid \|\vec{x}\| \leq 1\}$ und berechnen Anteil des Volumens (oder Äquivalent Wkeit der Daten) welche in der äußersten Schale $S_\epsilon = \{\vec{x} \in \mathbb{R}^d \mid 1 - \epsilon \leq \|\vec{x}\| \leq 1\}$ der Dicke ϵ enthalten ist:

$$V(\epsilon) = \frac{C(d) 1^d - C(d)(1-\epsilon)^d}{C(d) 1^d} = 1 - (1-\epsilon)^d$$

\rightarrow bei höheren Dimensionen liegt der Anteil $V(\epsilon)$ auch schon für kleine ϵ nahe bei 1
 \rightarrow fast alle Daten liegen in dünner Schale am Rand der Kugel



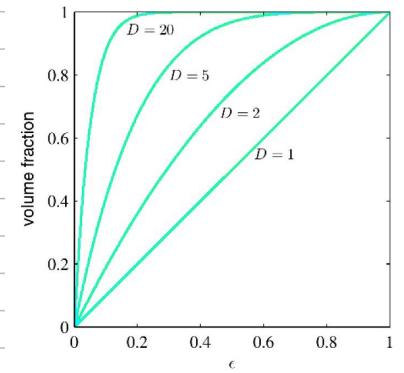
If we peel a high-dimensional orange, there is almost nothing left.

wenn wir folglich den Raum mithilfe der euklidischen Metrik in lokale Träger unterteilen, liegen gleichmäßige Daten mit hoher Wkeit am Rand eines Trägers und diese Lokalität ist v. gernet eine verhinderte Glättigkeit der Ausgabe t bezüglich der Eingabe \vec{x} aufgrund davon

in realen Regressionsproblemen ist andererseits Verteilung der Eingaben $\vec{x} \in \mathbb{R}^d$ meistens sehr weit von einer Gleichverteilung entfernt, die \vec{x} liegen meist sogar konzentriert in der Nähe einer (i. A. nichtlinearen) Untermannigfaltigkeit $M \subset \mathbb{R}^d$ viel kleinerer Dimension und zusätzlich hängt t häufig nur von einigen möglichen Veränderungsrichtungen der Eingabe \vec{x} sensitiv ab

diese Untermannigfaltigkeit und daran die lokalen sensiblen Richtungen (für die Abhängigkeit der Ausgabe von der Eingabe) definieren ein Distanzmaß (Metrik) das i. A. von der euklidischen Distanz des Eingraumes \mathbb{R}^d erheblich verschieden ist

Familie von Basisfunktionen deren Träger auf der Untermannigfaltigkeit M liegen und lokal bezüglich dieses Distanzmaßes gewählt sind kann prinzipiell die Glättigkeit der Ausgabe t viel effizienter ausnutzen und bis zu einem gewissen Grade den Fluss der Dimensionalität verhindern
 \rightarrow leider kann eine solche Familie von Basisfunktionen nur in seltenen Fällen einfach von Hand konstruiert werden!



Lineare Modelle für Klassifikationsprobleme

Zwei grundlegende lineare Modelle für bearbeitigtes Klassifikationsproblem:

- generativ $\rightarrow p(\vec{t}, \vec{x}) (= p(\vec{t}|\vec{x})p(\vec{x}) = p(\vec{x}|\vec{t})p(\vec{t}))$
- diskriminativ $\rightarrow p(\vec{t}|\vec{x})$

bei generativ manchmal natürlicher $p(\vec{t}|\vec{x})p(\vec{x})$ oder $p(\vec{x}|\vec{t})p(\vec{t})$ zu verwenden

da bedingte Verteilung $p(\vec{t}|\vec{x})$ aus jedem generativen Modell $p(\vec{t}, \vec{x})$ berechnet werden kann, sind generative Modelle universeller und damit potentiell schwieriger zu lernen als diskriminative Modelle

für Entscheidungsproblem nur $p(\vec{t}|\vec{x})$ benötigt, welche mit beiden Modellarten verfügbar ist

Kenntnis von $p(\vec{x})$ hat insbesondere den Vorteil, dass damit ungewöhnliche Ergebnisse \vec{x} , für welche das Modell nicht trainiert wurde, detektiert werden können (novelty detection)

generativ

generatives Modell $p(\vec{x}, \vec{t})$ für bearbeitigtes Klassifikationsproblem:

Klassen-Prior - V.:

$p(\vec{t})$, oft als diskrete $p(\vec{t}) = \text{Cat}_{\theta_0}(\vec{t})$ modelliert

Klassen-konditionierte - V.:

$p(\vec{x}|\vec{t})$

Klassen-Posterior - V.:

$$p(\vec{t}|\vec{x}) = \frac{p(\vec{x}|\vec{t})p(\vec{t})}{p(\vec{x})} = \frac{p(\vec{x}|\vec{t})p(\vec{t})}{\sum_{j=1}^C p(\vec{x}|j)p(j)}$$

Klassifikation:

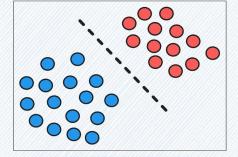
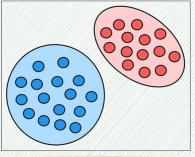
$$\vec{x} = \operatorname{argmax}_{\vec{t}} p(\vec{t}|\vec{x}) = \operatorname{argmax}_{\vec{t}} p(\vec{x}|\vec{t})p(\vec{t})$$

Entscheidungsbereiche:

$$R_k = \{\vec{x} \mid k = \operatorname{argmax}_{\vec{t}} p(\vec{t}|\vec{x})\}$$

Achtung: Bezeichnungen Prior / Posterior beziehen sich in diesem Zusammenhang nicht auf die Parameter des Modells!

verschiedene generative Klassifikationsmodelle unterscheiden sich in der Form der Klassen-Verteilung $p(\vec{x}|\vec{t})$

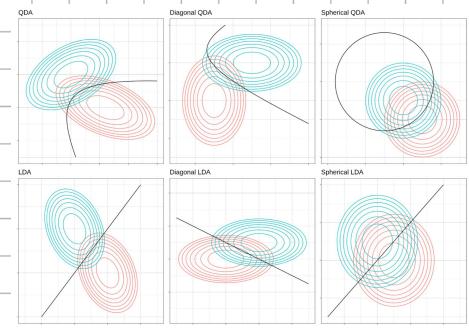
Discriminative Models	Generative Models
	
Learns the decision boundary between classes	Learns the input distribution
Maximizes the conditional probability: $P(Y X)$	Maximizes the joint probability: $P(X, Y)$
Directly estimates $P(Y X)$	Estimates $P(X Y)$ to find $P(Y X)$ using Bayes' rule
Cannot generate new data	Can generate new data
Specifically meant for classification tasks	Typically, they are NOT used to solve classification tasks
Discriminative models don't possess generative properties	Generative models possess discriminative properties
Logistic Regression	Hidden Markov Models
Random Forests	Naive Bayes
SVMs	Gaussian Mixture Models
Neural Networks	Gaussian Discriminant Analysis
Decision Tree	LDA
kNN	Bayesian Networks

Gaußsche Diskriminanten Analyse (GDA) generativ

Klassen-Bedingte-V.: $p(\vec{x}|t, \vec{M}_t, \Sigma_t) = \mathcal{N}(\vec{x}|\vec{M}_t, \Sigma_t)$

Klassen-Prior-V.: $p(t|\vec{\pi}) = \text{Cat}_{\vec{\pi}}(t)$

generatives Modell: $p(t, \vec{x} | \vec{\theta}) = p(\vec{x}|t, \vec{M}_t, \Sigma_t) p(t|\vec{\pi})$



Aktivierung der Klasse k bei Eingabe \vec{x} :

$$\begin{aligned} a_k(\vec{x}) &= \ln(p(\vec{x}|k)p(k)) = \ln p(\vec{x}|k) + \ln p(k) \\ &= -\frac{1}{2} \left(d \ln(2\pi) + \ln(\det \Sigma_k) + (\vec{x} - \vec{M}_k)^T \Sigma_k^{-1} (\vec{x} - \vec{M}_k) \right) + \ln p(k) \\ &= -\frac{1}{2} \left(\vec{x}^T \Sigma_k^{-1} \vec{x} - 2 \vec{x}^T \Sigma_k^{-1} \vec{M}_k + c_k \right) \end{aligned}$$

Softmax-Funktion:

$$\begin{aligned} \vec{s}: \mathbb{R}^c &\rightarrow \{ \vec{p} \in \mathbb{R}_+^c \mid p_1 + \dots + p_c = 1 \} \\ \vec{a} &\mapsto \vec{s}(\vec{a}) = (s_1(\vec{a}), \dots, s_c(\vec{a}))^T \quad \text{wobei } s_k(\vec{a}) = e^{a_k} / \sum_{j=1}^c e^{a_j} \end{aligned}$$

transformiert Aktivierungen in Wahrscheinlichkeit

mit Softmax kann man für jedes gen. M. die Klassen-Posterior-V. $p(t|\vec{x})$ universell in der Form $p(t|\vec{x}) = S_t(\vec{a}(\vec{x}))$ schreiben wobei $S_t(\vec{a}(\vec{x})) = \text{Cat}_{\vec{\pi}(\vec{a}(\vec{x}))}(t)$ gilt

Entscheidungsbereiche: $R_k = \{ \vec{x} \mid k = \arg \max_i a_i(\vec{x}) \}$

Grenze zwischen R_j und R_k ist in der Menge $\{ \vec{x} \mid a_j(\vec{x}) = a_k(\vec{x}) \}$ enthalten, welche durch die Quadrik

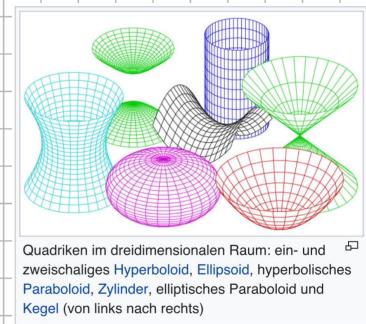
$$\vec{x}^T \Sigma_k^{-1} \vec{x} - 2 \vec{x}^T \Sigma_k^{-1} \vec{M}_k + c_k = \vec{x}^T \Sigma_j^{-1} \vec{x} - 2 \vec{x}^T \Sigma_j^{-1} \vec{M}_j + c_j$$

gegeben ist (Quadratische Diskriminanten Analyse (QDA))

falls $\Sigma_k = \Sigma_j = \Sigma$ gilt, erhält man eine Hyperschneide: $\vec{x}^T \Sigma^{-1} (\vec{M}_k - \vec{M}_j) = \frac{c_k - c_j}{2}$

falls dies gefordert wird, spricht man von Linearer Diskriminanten Analyse (LDA)

↳ Parameter-Sharing / -Sharing



Quadriken im dreidimensionalen Raum: ein- und zweischaliges Hyperboloid, Ellipsoid, hyperbolisches Paraboloid, Zylinder, elliptisches Paraboloid und Kegel (von links nach rechts)

ML-Methode

$$\text{QDA: } \vec{\pi}_m = \frac{\vec{n}_m}{\|\vec{n}_m\|}, \quad \vec{M}_{m,k} = \frac{\sum_{t_i=k} \vec{x}_i}{n_k}, \quad \Sigma_{m,k} = \frac{\sum_{t_i=k} (\vec{x}_i - \vec{M}_{m,k})(\vec{x}_i - \vec{M}_{m,k})^T}{n_k}$$

$$\text{LDA: } \Sigma_m = \sum_{k=1}^C \frac{n_k}{N} \sum_{m,k}$$

MAP- und Bayes-Methode

LDA als Reduktion d. Flexibilität einer Kunde \rightarrow mit geeigneter Prior-V. in Kombination mit MAP kann Flexibilität auch feiner kontrolliert werden

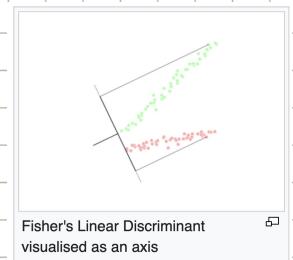
MAP und Bayes sind insbesondere für konj. Prior-V. berechenbar

Feature Extraktion und Dimensionsreduktion

anstatt Eingaben \vec{x} direkt zu verwenden, kann es auch sinnvoll sein, diese zuerst geednet zu transformieren und die für die Klassifikation relevanten Features herauszuarbeiten

bei sehr hochdimensionalen Eingaben \vec{x} kann Flexibilität der GDA durch Dimensionsreduktion begrenzt werden, wozu häufig linear auf Unterräum projiziert wird

→ z.B. Fisher's - Lineare - Diskriminanten - Analyse



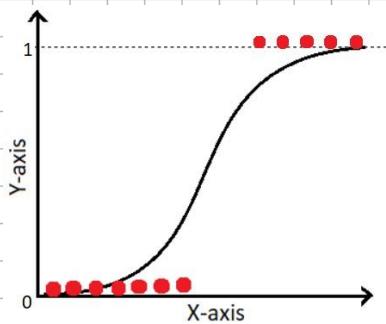
Logistische Regression diskriminativ

diskr. Modell: $p(t|\vec{x}, \underline{w}) = \text{Cot}_{\vec{\phi}(\underline{w}\vec{\phi}(\vec{x}))}(t) = S_t(\underline{w}\vec{\phi}(\vec{x}))$

Basisfunktionen: $\vec{\phi}(\vec{x}) = (\phi_0(\vec{x}), \dots, \phi_n(\vec{x}))^\top$

Parameter:

$$\underline{w} = (\vec{w}_1, \dots, \vec{w}_c)^\top = \begin{bmatrix} w_{1,0} & \dots & w_{1,n} \\ \vdots & & \vdots \\ w_{c,0} & \dots & w_{c,n} \end{bmatrix}$$



für $\phi_0(\vec{x}) = 1$, $\phi_k(\vec{x}) = x_k$ resultiert affine Form der log. Regression

obwohl Symmetrie in Parametrisierung angesprochen ist, hat sie den Nachteil, dass \underline{w} nicht identifizierbar sind; man kann einen beliebigen Vektor $\vec{\theta}^\top$ zu allen Zeilen addieren und man verändert nichts an der Verteilung der Ausgabe
 → diese Überparametrisierung kann man eliminieren, indem man \underline{w} geeignet normalisiert
 → häufig fordert man, dass letzte Zeile $= \vec{\theta}^\top \vec{1}$ ist

mit Fehlklassifikations-Kostenmaß erhält man die Entscheidungsbereiche

$$R_t = \{ \vec{x} \mid t = \arg \max_k p(k|\vec{x}, \underline{w}) = \arg \max_k \vec{w}_k^\top \vec{\phi}(\vec{x}) \}$$

wobei die Grenze zw. R_u und R_v in $R_u \cap R_v \subseteq \{ \vec{x} \mid (\vec{w}_u - \vec{w}_v)^\top \vec{\phi}(\vec{x}) = 0 \}$ enthalten ist

Entscheidungsbereiche sind im Feature-Raum $\vec{y} = \vec{\phi}(\vec{x})$ konvex und Entscheidungsbereiche in Hyperebenen enthalten
 → bei affiner Basisfunktion gilt dies auch im Edgetraum

ML- und MAP-Parameter

Likelihood nicht mehr analytisch maximierbar (kann gesagt weil $\ln()$ einer Summe nicht vereinfacht werden kann), jedoch immer noch konkav, was numerische Maximierung einfacher und robust macht

falls Klassen separierbar sind, existiert jedoch kein Maximum und ML-Methode konvergiert gegen Funktion im Stile der Heaviside resp. einem numerischen Overflow
 → kann mit MAP verhindern lassen

falls Prior $p(\vec{w}|\alpha) = N(\vec{w} | \vec{\theta}, \alpha^{-1} E)$ tritt zusätzlicher Summand $-\frac{\alpha}{2} \|\vec{w}\|_2^2$ auf wenn man logarithmiert (MAP)

→ dies wird L_2 -regularisierte Maximum Likelihood Methode genannt

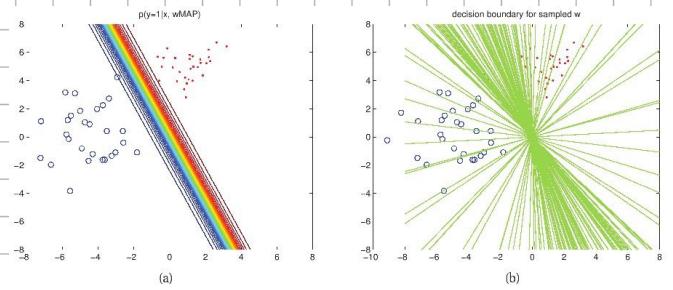
numerisch wird Optimierungsproblem der ML- oder MAP-Methode mit verschiedenen Standardmethoden gelöst (historisch relevant aber heute nicht mehr verwendet: IRLS (basiert auf Newton-Methode))

Bayessches Lernen

$$\text{Posterior - V.: } p(\underline{w} | \vec{t}, \underline{x}) = \frac{p(\vec{t} | \underline{x}, \underline{w}) p(\underline{w})}{p(\vec{t} | \underline{x})}$$

$$= \frac{p(\vec{t} | \underline{x}, \underline{w}) p(\underline{w})}{\int p(\vec{t} | \underline{x}, \underline{w}) p(\underline{w}) d\underline{w}}$$

hochdimensionales Integral
numerisch leicht!



für keine Prior - V. geschlossen berechenbar
(im Gegensatz zu linearer Regression)
→ müssen approximieren

Monte-Carlo-Approximation einer Verteilung:

$$p(\vec{x}) \approx q(\vec{x}) = \frac{1}{|S|} \sum_{s \in S} \delta_s(\vec{x})$$

S: Menge von Samples

→ viel Rechnzeit nötig!

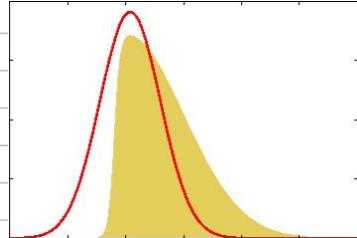
Variationsapproximation:

aus Familie geeigneter parametrisierter Verteilungen Q und diejenige Verteilung $q(\vec{x})$ gewählt, die ein Distanzmaß (häufig KL-Divergenz) zur exakten Verteilung $p(\vec{x})$ minimiert
→ basiert auf num. Optimierung
→ weniger rechenaufwändig als Monte-Carlo-Appr., aber schwieriger zu spezifizieren

Laplace-Approximation:

$$p(\vec{z}) = \frac{f(\vec{z})}{\int f(\vec{z}) d\vec{z}} = \frac{f(\vec{z})}{\vec{z}} \quad \begin{matrix} \text{bemerkbar} \\ \text{nicht bemerkbar} \end{matrix} \quad \left. \begin{matrix} \text{Anwendung} \\ \text{Anwendung} \end{matrix} \right\}$$

$$\simeq q(\vec{z}) = \mathcal{N}(\vec{z} | \vec{z}_0, \Lambda^{-1}) = \mathcal{N}(\vec{z} | \arg \max_{\vec{z}} p(\vec{z}), (-H \ln f(\vec{z}))^{-1})$$



→ Normalverteilung mit gleichem Modus und Krümmung bei diesem

$$\vec{z} \stackrel{\text{def}}{=} \frac{f(\vec{z})}{q(\vec{z})} \stackrel{\text{ln}(\cdot)}{\rightarrow} \ln \vec{z} \stackrel{\text{def}}{=} (\ln f(\vec{z}) - \ln q(\vec{z})) = \ln f(\vec{z}) - \frac{1}{2} \ln \det \Lambda + \frac{d}{2} \ln(2\pi)$$

bei Posterior - V. $p(\vec{\theta}|D) = p(D|\vec{\theta})p(\vec{\theta})/p(D)$ ist vielfach $p(D|\vec{\theta})p(\vec{\theta})$ analytisch verfügbar, aber Normalisationskonstante $p(D) = \int p(D|\vec{\theta})p(\vec{\theta})d\vec{\theta}$ nicht berechenbar

Laplace-Posterior-Approximation (allgemein):

$$q(\vec{\theta}) = \mathcal{N}(\vec{\theta} | \vec{\theta}_{\text{MAP}}, \Lambda^{-1}) \text{ mit } \Lambda = -H \ln p(D|\vec{\theta}_{\text{MAP}}) - H \ln p(\vec{\theta}_{\text{MAP}})$$

$$\Rightarrow \ln p(D) \stackrel{\text{def}}{=} \ln p(D|\vec{\theta}_{\text{MAP}}) + \ln p(\vec{\theta}_{\text{MAP}}) - \frac{1}{2} \ln \det \Lambda + \frac{d}{2} \ln(2\pi) \quad (\text{Logarithmierte marginale Likelihood})$$

Occams Factor → bestrafen flexible/komplexe Modelle

log. Regression mit Bayesscher Methode:

Prior: $p(\underline{w}) = \mathcal{N}(\underline{w} | \vec{w}_0, \Lambda_0^{-1})$

Lapl.-Post.-Appr.: $p(\underline{w} | \vec{w}_0, \Lambda_0, \underline{x}, \vec{t})$

Post.-Voraussage: $p(t | \underline{x}, \underline{x}, \vec{t}) = \int p(t | \underline{x}, \underline{w}) p(\underline{w} | \vec{w}_0, \Lambda_0, \underline{x}, \vec{t}) d\underline{w}$

nicht geschlossen berechenbar
→ müssen numerisch approximieren

$$\simeq \frac{1}{|S|} \sum_{s \in S} p(t | \underline{x}, \underline{w}) \quad \text{Monte-Carlo-Approximation}$$

zuerst numerisch mit \underline{x}, \vec{t} $\underline{w}_{\text{MAP}}$ bestimmen und damit anschließend $\Lambda = -H \ln p(\vec{t} | \underline{x}, \underline{w}_{\text{MAP}}) + \Lambda_0$

Erweiterungen und Schwächen der logistischen Regression

generelle lineare Modelle:

- Verallgemeinerung d. log. Regressions auf andere Ausgabeverteilungen
- anstatt Softmax wird eine zur Verteilung passende Aktivierungsfunktion verwendet
- lineare Abhängigkeit von Gewichten und Features

Schwächen:

- Fluss der Dimensionen in Kombination mit fixer Wahl von nicht adaptiven Features
- Probleme bei hochdimensionalen Eingaben

Komplexität und Modellselektion

erwarteter Verlust und Flexibilität von Modellen

der Einfachheit halber betrachten wir ein-dimensional Regressionsprobleme und wählen quadratisches Verlustmaß zur Bewertung der Vorhersagegenauigkeit

Verlustmaß: $L(t_{\text{true}}, t_{\text{pred}}) = (t_{\text{true}} - t_{\text{pred}})^2$

exakte Verteilung: $p(t, \vec{x}) \rightarrow \text{nicht bekannt}$

Voransagefunktion: $f(\vec{x})$

optimale Voransagef.: $m(\vec{x}) = E_{p(\cdot | \vec{x})}[t]$

erwarteter Verlust:

$$\begin{aligned} E[L_f] &= \int (t - f(\vec{x}))^2 p(t, \vec{x}) dt d\vec{x} = \int (t - m(\vec{x}) + m(\vec{x}) - f(\vec{x}))^2 p(t, \vec{x}) dt d\vec{x} \\ &= \left[((t - m(\vec{x}))^2 + 2(t - m(\vec{x}))(m(\vec{x}) - f(\vec{x})) + (m(\vec{x}) - f(\vec{x}))^2 \right] p(t, \vec{x}) dt d\vec{x} \\ &= \int (t - m(\vec{x}))^2 p(t, \vec{x}) dt d\vec{x} + \int 2(t - m(\vec{x}))(m(\vec{x}) - f(\vec{x})) p(t, \vec{x}) dt d\vec{x} + \int (m(\vec{x}) - f(\vec{x}))^2 p(t, \vec{x}) dt d\vec{x} \end{aligned}$$

$$\int (t - m(\vec{x}))^2 p(t, \vec{x}) dt d\vec{x} = E[L_m] \quad \text{unvermeidbarer erwarteter Verlust}$$

$$\begin{aligned} \int 2(t - m(\vec{x}))(m(\vec{x}) - f(\vec{x})) p(t, \vec{x}) dt d\vec{x} &= \left[\left(\int 2(t - m(\vec{x}))(m(\vec{x}) - f(\vec{x})) p(t | \vec{x}) dt \right) p(\vec{x}) d\vec{x} \right] = \left[(2(m(\vec{x}) - f(\vec{x})) \int (t - m(\vec{x})) p(t | \vec{x}) dt) p(\vec{x}) d\vec{x} \right] \\ &= \left[(2(m(\vec{x}) - f(\vec{x})) \left(\int t p(t | \vec{x}) dt - \int m(\vec{x}) p(t | \vec{x}) dt \right)) p(\vec{x}) d\vec{x} \right] = \left[(2(m(\vec{x}) - f(\vec{x})) \left(\int t p(t | \vec{x}) dt - m(\vec{x}) \int p(t | \vec{x}) dt \right)) p(\vec{x}) d\vec{x} \right] \\ &= \left[(2(m(\vec{x}) - f(\vec{x})) \left(\int t p(t | \vec{x}) dt - m(\vec{x}) \cdot 1 \right)) p(\vec{x}) d\vec{x} \right] = \left[(2(m(\vec{x}) - f(\vec{x})) (m(\vec{x}) - m(\vec{x}))) \right] p(\vec{x}) d\vec{x} = 0 \end{aligned}$$

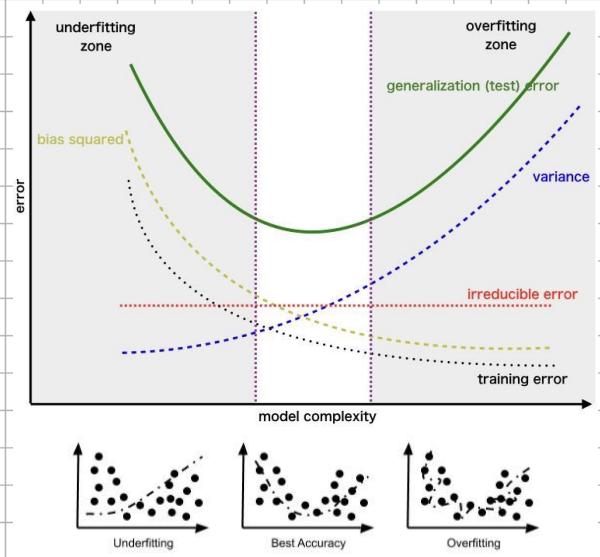
$$\begin{aligned} \int (m(\vec{x}) - f(\vec{x}))^2 p(t, \vec{x}) dt d\vec{x} &= \left[\left(\int (t - m(\vec{x}))^2 p(t | \vec{x}) dt \right) p(\vec{x}) d\vec{x} \right] = \left[((m(\vec{x}) - f(\vec{x}))^2 \int p(t | \vec{x}) dt) p(\vec{x}) d\vec{x} \right] \\ &= \left[((m(\vec{x}) - f(\vec{x}))^2 \cdot 1 \right] p(\vec{x}) d\vec{x} = \int (m(\vec{x}) - f(\vec{x}))^2 p(\vec{x}) d\vec{x} = \int (m(\vec{x}) - f(\vec{x}))^2 p(\vec{x}) d\vec{x} \end{aligned}$$

$$E[L_f] = \int (m(\vec{x}) - f(\vec{x}))^2 p(\vec{x}) d\vec{x} + E[L_m]$$

unser Modell (und damit auch erw. Verlust) hängt von Daten ab:

$$\begin{aligned} E[L_{f_D}] &= E_D [E[L_{f_D(\cdot | D)}]] = \int E_D [(m(\vec{x}) - f_D(\vec{x} | D))^2] p(\vec{x}) d\vec{x} + E[L_m] \\ &= \int ((m(\vec{x}) - E_D[f_D(\vec{x} | D)])^2 + \text{Var}_D[f_D(\vec{x} | D)]) p(\vec{x}) d\vec{x} + E[L_m] \\ &= \int (m(\vec{x}) - E_D[f_D(\vec{x} | D)])^2 p(\vec{x}) d\vec{x} + \int \text{Var}_D[f_D(\vec{x} | D)] p(\vec{x}) d\vec{x} + E[L_m] \\ &= \text{bias}^2 + \text{variance} + \text{noise} \end{aligned}$$

und $\mathbb{D} = \{D \mid \text{mit } |D| = N\}$ alle möglichen Trainingsdaten vom Umfang N bezeichnet



Modelle selektieren mit der Bayes-Methode

Posterior auf Raum der Modelle:

$$p(M|D) = \frac{p(D|M) p(M)}{p(D)}$$

wichtigste
Formel/
dieses Abschnitts

Method	Definition
Maximum likelihood	$\hat{\theta} = \operatorname{argmax}_{\theta} p(\mathcal{D} \theta)$
MAP estimation	$\hat{\theta} = \operatorname{argmax}_{\theta} p(\mathcal{D} \theta)p(\theta \eta)$
ML-II (Empirical Bayes)	$\hat{\eta} = \operatorname{argmax}_{\eta} \int p(\mathcal{D} \theta)p(\theta \eta)d\theta = \operatorname{argmax}_{\eta} p(\mathcal{D} \eta)$
MAP-II	$\hat{\eta} = \operatorname{argmax}_{\eta} \int p(\mathcal{D} \theta)p(\theta \eta)p(\eta)d\theta = \operatorname{argmax}_{\eta} p(\mathcal{D} \eta)p(\eta)$
Full Bayes	$p(\theta, \eta \mathcal{D}) \propto p(\mathcal{D} \theta)p(\theta \eta)p(\eta)$

Evidenz / marginale Likelihood:

$$p(D|M) = \int p(D|\vec{\theta}, M) p(\vec{\theta}|M) d\vec{\theta}$$

für viele Modelle nicht geschlossen berechenbar
für lineare Modelle berechenbar ✓

Posterior über Parameter:

$$p(\vec{\theta}|D, M) = \frac{p(D|\vec{\theta}, M) p(\vec{\theta}|M)}{p(D|M)}$$

falls man nur 1 Modell betrachtet, kann man Bezeichnung M weglassen

Normalisierungsfaktor in $p(M|D)$:

$$p(D) = \int p(D|M) p(M) dM$$

Hyperprior $p(M)$ hat auch wieder Parameter → mit Bayes führt das konsequenterweise auf hierarchische bayessche Modelle (oder zu graphischen probabilistischen Modellen)
→ prinzipiell sinnvoll aber selten berechenbar
→ Annäherungen

Bayessche Posterior-Voransage-Verteilung (bayessche Modell-Mittelung):

$$p(\vec{E}|\vec{x}, D) = \int p(\vec{E}|\vec{x}, M, D) p(M|D) dM \quad \text{Mittelung aller Modelle}$$

$$\text{wobei } p(\vec{E}|\vec{x}, M, D) = \int p(\vec{E}|\vec{x}, \vec{\theta}, M) p(\vec{\theta}|D, M) d\vec{\theta}$$

Bayessche Modell-Selektion / Typ II MAP:

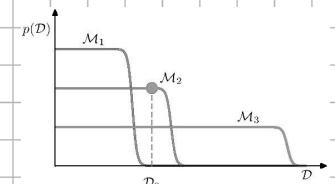
$$p(\vec{E}|\vec{x}, \hat{M}, D) \quad \text{wobei } \hat{M} = \operatorname{argmax}_M p(M|D) = \operatorname{argmax}_M p(D|M) p(M)$$

empirische Bayessche Modell-Selektion / Typ II ML:

$$p(\vec{E}|\vec{x}, \hat{M}, D) \quad \text{wobei } \hat{M} = \operatorname{argmax}_M p(D|M)$$

↓ Vereinfachung

↓ Vereinfachung



Selektion eines Modells mittels empirischer bayesscher Modellselektion, also basierend auf der höchsten marginalen Likelihood $p(D|M)$ der Trainingsdaten (bei welches über Parameter $\vec{\theta}$ mit Hilfe einer Prior-Verteilung integriert und nicht maximiert wird) neigt viel weniger zu Überanpassung als Selektion des Modells mit maximaler Likelihood

Bayessche Hyperparameterselektion für lineare Regressionsprobleme

Modell:

$$p(t|\vec{x}, \vec{w}, \beta) = \mathcal{N}(t | \Phi(\vec{x})^T \vec{w}, \beta^{-1})$$

gaussscher Prior:

$$p(\vec{w}|\alpha) = \mathcal{N}(\vec{w} | \vec{0}, \alpha^{-1} \mathbb{E})$$

(logarithmierte) marginale Likelihood:

$$\ln p(\vec{t}|\vec{X}, \alpha, \beta) = -\frac{1}{2}(N \ln(2\pi) + \ln \det(\Sigma) + \vec{t}^T \Sigma^{-1} \vec{t}) \quad \text{wobei } \Sigma = \beta^{-1} \mathbb{E} + \alpha^{-1} \Phi(\vec{X}) \Phi(\vec{X})^T$$

man kann nun logarithmierte marginale Likelihood (Typ II ML) mit Optimierung maximieren

man kann auch Prior über α, β wählen und Typ II MAP oder sogar volle bayessche Modellierung mittels Monte-Carlo-Methoden oder Variationsmethoden approximieren

Automatische - Relevanz - Bestimmung (ARD)

kann man als Feature Selection für lineare Regressionsmodelle anfassen

Modell:

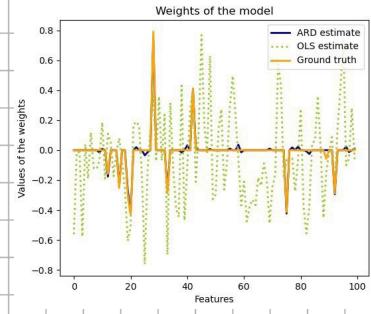
$$p(t|\vec{x}, \vec{w}, \beta) = \mathcal{N}(t | \Phi(\vec{x})^T \vec{w}, \beta^{-1})$$

gaussscher Prior:

$$p(\vec{w}|\vec{x}) = \mathcal{N}(\vec{w} | \vec{0}, D_{\vec{x}}^{-1}) = \prod_{k=0}^m \mathcal{N}(w_k | 0, \alpha_k^{-1})$$

(logarithmierte) marginale Likelihood:

$$\ln p(\vec{t}|\vec{X}, \vec{\alpha}, \beta) = -\frac{1}{2}(N \ln(2\pi) + \ln \det(\Sigma) + \vec{t}^T \Sigma^{-1} \vec{t}) \quad \text{wobei } \Sigma = \beta^{-1} \mathbb{E} + \Phi(\vec{X}) D_{\vec{x}}^{-1} \Phi(\vec{X})^T$$



um maximiert nun logarithmierte marginale Likelihood (Typ II ML) mit Optimierung

dabei streben viele der Präzisions-Hyper-Parameter $\alpha_k \rightarrow \infty$ und folglich werden die dazugehörigen Parameter $w_k \rightarrow 0$ gezerrt

bei ARD wird vielfach durch noch die Präzision β der Ausgabeverteilung optimiert und mittels Gamma-Prior-Verteilung über α und β kann regularisiert werden (Typ II MAP), was erlaubt die Stärke der Selektion zu kontrollieren

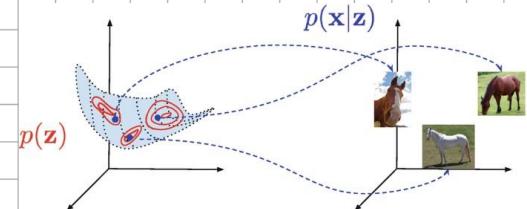
mit ARD können auch bei (logistischen) Regressien Features selektiert werden, wobei die marginale Likelihood jedoch nicht exakt berechenbar ist und somit approximiert werden muss

Modelle mit verborgenen Variablen

Definition:

$$p(\vec{x}, \vec{z} | \vec{\theta}) \quad \vec{x}: \text{beobachtete Variablen}$$

\vec{z} : verborgene Variablen



induzierte Verteilung über beobachteten Variablen:

$$p(\vec{x} | \vec{\theta}) = \int p(\vec{x}, \vec{z} | \vec{\theta}) d\vec{z}$$

manchmal aber nicht immer haben verborgene Variablen ohne Bedeutung

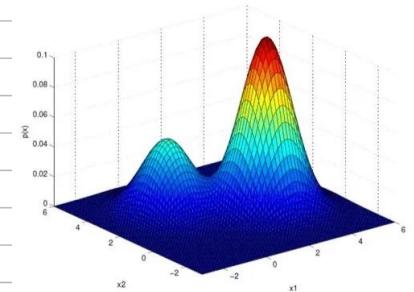
Mischungsmodelle

$$p(\vec{x}, z | \vec{\theta}, \vec{\pi}) = p(\vec{x} | z, \vec{\theta}) p(z | \vec{\pi})$$

mit $p(z=k | \vec{\pi}) = \pi_k$ und $p(\vec{x} | z=k, \vec{\theta}) = p_k(\vec{x} | \vec{\theta}_k)$

marginal Verteilung:

$$p(\vec{x} | \vec{\theta}, \vec{\pi}) = \sum_{k=1}^K p(\vec{x}, z=k | \vec{\theta}, \vec{\pi}) = \sum_{k=1}^K \pi_k p_k(\vec{x} | \vec{\theta}_k)$$



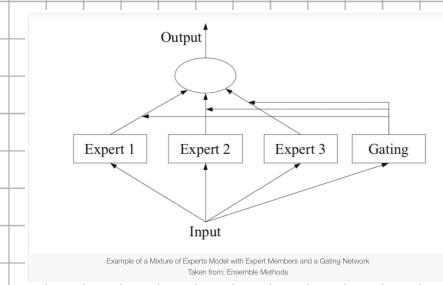
Gaußsche Mischungsmodelle:

$$p(\vec{x} | \vec{\theta}, \vec{\pi}) = \sum_{k=1}^K \pi_k \mathcal{N}(\vec{x} | \vec{\mu}_k, \Sigma_k)$$

Mischung von Experten:

$$\begin{aligned} p(t, z | x, \vec{\theta}, \vec{\pi}) &= p(t | x, z=k, \vec{\theta}) p(z | x, \vec{\pi} = w) \\ &= \mathcal{N}(t | \vec{\phi}(x)^T \vec{w}_k, \beta_k^{-1}) \text{Cat}_{\vec{s}(\vec{w} \cdot \vec{\phi}(x))}(z) \end{aligned}$$

- z hängt (über Logistisches Regressionsmodell) von x ab
- lineare Regressionsmodelle für die einzelnen z bedingten Komponenten
- einzelnen Komponenten können sich auf Regionen der Eingaben spezialisieren («Experten»)



Inferenz in Mischungsmodellen:

$$p(z=j | \vec{x}, \vec{\theta}, \vec{\pi}) = \frac{p(\vec{x} | z=j, \vec{\theta}) p(z=j | \vec{\pi})}{p(\vec{x} | \vec{\theta}, \vec{\pi})} = \frac{\pi_j \cdot p_j(\vec{x} | \vec{\theta}_j)}{\sum_{k=1}^K \pi_k \cdot p_k(\vec{x} | \vec{\theta}_k)}$$

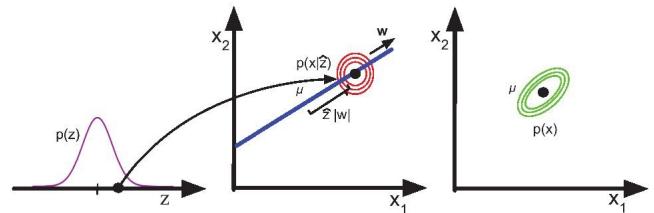
«responsibility»

die Bezeichnungen $\{1, \dots, K\}$ der Mischungskomponenten (genauer die Wkeit π_k zusammen mit Parametern $\vec{\theta}_k$) permutiert werden können ohne Verteilung $p(\vec{x} | \vec{\theta}, \vec{\pi})$ über beobachtbarem Variablen zu verändern, sind Parameter $\vec{\pi}, \vec{\theta}$ höchstens modulo dieser Permutationssymmetrie bestimbar
→ indem zusätzliche Bedingungen formuliert werden kann diese Symmetrie bei Berücksichtigung machen

Faktorenanalysemodelle

können zur (linearen) Dimensionsreduktion verwendet werden

zur Hauptkomponentenanalyse analoges, generatives Modell mit kontinuierlichem \vec{z}



$$p(\vec{x}, \vec{z}) = p(\vec{x} | \vec{z}) p(\vec{z}) \text{ mit } \vec{x} \in \mathbb{R}^d, \vec{z} \in \mathbb{R}^l \text{ wobei typischerweise } l \ll d \text{ ist}$$

$$p(\vec{z}) = \mathcal{N}(\vec{z} | \vec{0}, E) \quad \text{Standardnormalverteilte verborgene Variable}$$

$$p(\vec{x} | \vec{z}) = \mathcal{N}(\vec{x} | \underline{W}\vec{z} + \vec{\mu}, \text{diag}(\vec{C})) \quad W = \text{Faktormatrix}, \underline{W} \in \mathbb{R}^{d \times l}$$

da konditionale Verteilung $p(\vec{x} | \vec{z})$ unkorreliert ist, muss die affine Transformation $\underline{W}\vec{z} + \vec{\mu}$ der niederdimensionalen, auch unkorrelierten, verborgenen Variable \vec{z} die Korrelationen zu den beobachteten Variablen \vec{x} erklären. In diesem Sinne ist verborgene Variable ohne komprimierte Darstellung der beobachteten Variable in verschiedene unkorrelierte Faktoren

diese Auffassung wird durch die induzierte Verteilung

$$p(\vec{x}) = \mathcal{N}(\vec{x} | \vec{\mu}, \text{diag}(\vec{C}) + \underline{W}\underline{W}^T)$$

über die beobachtete Variable bestätigt, weil \underline{W} die Korrelation daran berichtet

die Flexibilität der Korrelation wird durch die Dimension l des Raums der verborgenen Variablen (Bottleneck), festgelegt und liegt zwischen unkorreliert (diagonale Kovarianz) und voller Kopplung (allgemeine symmetrische Kovarianz)

indem wir $\text{diag}(\vec{C}) = \sigma E$ wählen erhalten wir im Wesentlichen die generative Version der Hauptkomponentenanalyse (probabilistische PCA, PPCA)

ein Vorteil des Faktorenanalysemodells gegenüber PPCA ist, dass die nicht durch die affine Transformation $\underline{W}\vec{z} + \vec{\mu}$ erklärbaren Anteile von \vec{x} nicht in allen Komponenten dieselbe Varianz haben müssen

Inferenzproblem für das Faktorenanalysemodell:

$$\vec{m}_{\vec{x}} = E[\vec{z} | \vec{x}] \quad \text{Latenz-Faktor der Beobachtung } \vec{x}$$

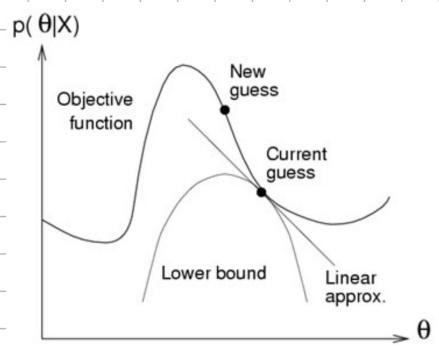
$$p(\vec{z} | \vec{x}) = \mathcal{N}(\vec{z} | \vec{m}_{\vec{x}}, \Sigma) \text{ mit } \Sigma = (E + \underline{W}^T \text{diag}(\vec{C})^{-1} \underline{W})^{-1} \text{ und } \vec{m}_{\vec{x}} = \Sigma \underline{W}^T \text{diag}(\vec{C})^{-1} (\vec{x} - \vec{\mu})$$

Faktormatrix \underline{W} nicht identifizierbar da mit einer beliebigen orthogonalen Matrix $R \in \mathbb{R}^{l \times l}$ von rechts multipliziert eine andere Matrix $\tilde{W} = \underline{W}R$ dieselbe Verteilung $p(\vec{x})$ resultiert
 → dies muss man beachten, wenn man z.B. die Spalten von \underline{W} interpretieren möchte
 → um den Grad der Nicht-Identifizierbarkeit zu verringern, kann man zusätzliche Bedingungen an \underline{W} formulieren, z.B. orthogonale Spaltenvektoren fordern (→ Singulärwertzerlegung!)

Lernproblem für Modelle m. v. Variablen: EM - Algorithmus

«Anleitung um Algorithmus herzuleiten, mit dem man eine Annäherung an die ML- oder MAP-Parameter erhält wenn verborgene Variablen vorhanden sind»

um die Parameter $\vec{\theta}$ eines Modells $p(\vec{x}, \vec{z} | \vec{\theta})$ mit verborgenen Variablen \vec{z} an beobachtbare Daten D anzupassen, könnte man versuchen direkt ML auf induzierte Verteilung $p(\vec{x} | \vec{\theta}) = \int p(\vec{x}, \vec{z} | \vec{\theta}) d\vec{z}$ anzuwenden (maximieren in $\vec{\theta}$)



Problem: analytisch selten möglich und numerisch schnell schwierig (Funktion kompliziert und nicht konkav, benötigte Marginalisierung über verborgene Variable ziemlich verhältnistisch)

(eine robuste und häufig verwendete) Lösung: EM-Algorithmus (iterativ)

Konvergiert garantiert gegen lokales Maximum der Likelihood $\ln(p(D | \vec{\theta})) = \sum_{k=1}^{N_0} \ln(p(\vec{x}_k | \vec{\theta})) = \sum_{k=1}^{N_0} (\ln(\int p(\vec{x}_k, \vec{z} | \vec{\theta}) d\vec{z}))$

2 Phasen:

Expectation - Phase:

1. bedingte Verteilung $p(\vec{z} | D, \vec{\theta}_t) = \prod_{j=1}^{N_0} p(\vec{z}_j | \vec{x}_j, \vec{\theta}_t)$ der unbeobachteten Variablen $\vec{z} = \{\vec{z}_k | 1 \leq k \leq N_0\}$ gegeben die Daten $D = \{\vec{x}_k | 1 \leq k \leq N_0\}$ unter aktuellen Parametern $\vec{\theta}_t$ (analytisch oder numerisch) bestimmen

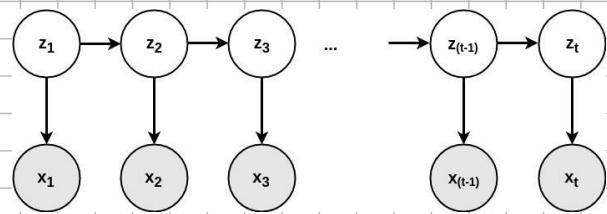
2. erwartete totale logarithmierte Likelihood $Q(\vec{\theta} | D, \vec{\theta}_t) = E_{p(\vec{z} | D, \vec{\theta}_t)} [\ln p(D, \vec{z} | \vec{\theta})] = \sum_{j=1}^{N_0} E_{p(\vec{z}_j | \vec{x}_j, \vec{\theta}_t)} [\ln(p(\vec{x}_j, \vec{z}_j | \vec{\theta}))]$ der Daten D in Abhängigkeit von $\vec{\theta}$ berechnen

Maximization - Phase:

$$\vec{\theta}_{t+1} = \operatorname{argmax}_{\vec{\theta}} Q(\vec{\theta} | D, \vec{\theta}_t)$$

Modelle für sequentielle Daten

einzelne Beobachtungen \rightarrow nicht unabhängig!
 \rightarrow hängt von vorangegangenem Zustand ab



Annahme dass hier Lernzykluszyklus besteht Einfluss auf Bildung hat

Annahmen:

- $p(\vec{x}_{t+1} | \vec{x}_1, \dots, \vec{x}_t) = p(\vec{x}_{t+1} | \vec{x}_t) = p_{\text{trans}}(\vec{x}_{t+1} | \vec{x}_t)$
 - $p(\vec{x}_1) = p_{\text{init}}(\vec{x}_1)$

Markov - Modell:

$$p(\vec{x}_1, \dots, \vec{x}_T) = p_{\text{init}}(\vec{x}_1) \prod_{t=2}^T p_{\text{trans}}(\vec{x}_t | \vec{x}_{t-1})$$

→ oft zu restitutiv, da nur unmittelbarer Vorgänger → MN höherer Ordnung

Markov-Modell höherer Ordnung $j > 1$:

$$p(\vec{x}_1, \dots, \vec{x}_T) = p_{\text{init}}(\vec{x}_1, \dots, \vec{x}_j) \prod_{t=2}^T p_{\text{trans}}(\vec{x}_t | \vec{x}_{t-1}, \dots, \vec{x}_{t-j})$$

Zustandsgrammmodelle

Annahmen:

- Transitionsmodell: $p(\vec{z}_{t+1} | \vec{z}_t) = p_{\text{trans}}(\vec{z}_{t+1} | \vec{z}_t)$ für alle Zeiten t stationäres MM 1ter Order
 - Beobachtungsmodell: $p(x_t | \vec{z}_t) = p_{\text{obs}}(x_t | \vec{z}_t)$

gemeinsame Verteilung einer Sequenz von beob. und verb. Variablen:

$$p(\vec{x}_{1:T}, \vec{z}_{1:T}) = p_{\text{init}}(\vec{z}_1) \prod_{t=2}^T p_{\text{trans}}(\vec{z}_t | \vec{z}_{t-1}) \prod_{t=1}^T p_{\text{obs}}(\vec{x}_t | \vec{z}_t)$$

→ mittels Marginalisierung über alle verborgenen Variablen $\vec{z}_{1:T}$ erhält man die induktive Verteilung $p(\vec{x}_{1:T})$ über die beobachteten Werte

i.A. werden für θ_{ans} & θ_{obs} parametrisierte Vert. verwendet $\rightarrow \vec{\theta} = (\vec{\theta}_{\text{mit}}, \vec{\theta}_{\text{ans}}, \vec{\theta}_{\text{obs}})$

in einigen Anwendungen hängen diese Verteilungen auch noch von (immer beobachteten) Ertragbelasten ab:

$$p(\vec{x}_{1:T}, \vec{z}_{1:T} | \vec{u}_{1:T}, \vec{\theta}) = p_{\text{init}}(\vec{z}_1 | \vec{u}_1, \vec{\theta}_{\text{init}}) \prod_{t=2}^T p_{\text{trans}}(\vec{z}_t | \vec{z}_{t-1}, \vec{u}_t, \vec{\theta}_{\text{trans}}) \prod_{t=1}^T p_{\text{obs}}(\vec{x}_t | \vec{z}_t, \vec{u}_t, \vec{\theta}_{\text{obs}})$$

Inferenzproblem: $p(\vec{z}_{1:T} | \vec{x}_{1:T}, \theta)$, $p(\vec{z}_t | \vec{x}_{1:T}, \theta)$ oder auch $p(\vec{x}_{1:T} | \theta)$

es folgen 2 Spezialisierungen von Zustandstrukturmödellen

- HMMs $\rightarrow z_t \in \{1, \dots, K\}$ diskret, \vec{x}_t diskret oder kontinuierlich
 - LDSs $\rightarrow \vec{z}_t \in \mathbb{R}^d$ kontinuierlich, p_{obs} & p_{dis} linear, alles Gauss-Verteilungen

Hidden Markov Modelle (HMMs)

$$p_{\text{init}}(z_1 = k) = \pi_k \rightarrow \vec{\pi} \in \mathbb{R}^K$$

$$p_{\text{trans}}(z_t = k | z_{t-1} = i) = A_{i,k} \rightarrow A \in \mathbb{R}^{K \times K} \quad \text{Transitionsmatrix}$$

Bemerkungen:

- $\vec{\pi} \in \mathbb{R}^K$ mit $K-1$ Freiheitsgraden da $\sum_k \pi_k = 1$
- $A \in \mathbb{R}^{K \times K}$ mit $K(K-1)$ Freiheitsgraden da $\sum_i A_{i,k} = 1$ für alle i

für Bedeutungsmodell müssen die K bedingten V. $p_{\text{obs}}(x_t | z_t = k)$ definiert werden

falls Beobachtungen $x_t \in \{1, \dots, L\}$ auch diskret sind, liefert allgemeinsten Parameterfamilie ohne Wahrscheinlichkeitsmatrix $O \in \mathbb{R}^{K \times L}$ mit $K(L-1)$ Freiheitsgraden mit Elementen $p_{\text{obs}}(x_t = l | z_t = i) = O_{i,l}$

bei kontinuierlichen Beobachtungen $\vec{x}_t \in \mathbb{R}^d$ werden hinsichtlich gaußische Beobachtungsmodell $p_{\text{obs}}(\vec{x}_t | z_t = k) = \mathcal{N}(\vec{x}_t | \mu_k, \Sigma_k)$ eingesetzt
 → können als gaußische Mischungsmodelle interpretiert werden, in welchen die Mischungskomponenten zeitlich korreliert sind

natürlich können auch anders parametrisierte Beobachtungsmodelle eingesetzt werden wie z.B. Neuronale Netze

Algorithmen dieses Abschnitts sind Spezialfälle von Message-Passing-Algorithmen (dynamische Programmierung) für das Inferenzproblem bei probabilistischen Graphischen Modellen

$$p(\vec{x}_{1:T}) = \sum_{z_{1:T} \in \{1, \dots, K\}^T} p(\vec{x}_{1:T}, z_{1:T}) = \sum_{z_{1:T} \in \{1, \dots, K\}^T} p_{\text{init}}(z_1) \prod_{t=2}^T p_{\text{trans}}(z_t | z_{t-1}) \prod_{t=1}^T p_{\text{obs}}(\vec{x}_t | z_t)$$

dieser Ausdruck besteht zwar implementiert aus K^T Summanden und lässt sich geschickt organisieren
 Filterung

Def. Filterung: rekursive Berechnung von $p(z_t | \vec{x}_{1:t})$ ausgehend von $p(z_{t-1} | \vec{x}_{1:t-1})$ resp. $\vec{\alpha}_t$ aus $\vec{\alpha}_{t-1}$

Th. Filterung: für ein HMM mit der Transitionsmatrix A und der Initialverteilung $\vec{\pi}$ kann man $\vec{\alpha}_t$ wie folgt rekursiv berechnen:

$$(1) \vec{\alpha}_1 = \frac{\vec{\lambda}_1}{l_1} \text{ mit } \vec{\lambda}_1 = \vec{\pi} \odot \vec{\pi} \text{ und } l_1 = p(\vec{x}_1) = \|\vec{\lambda}_1\|_1 \text{ für } t=1$$

$$(2) \vec{\alpha}_t = \frac{\vec{\lambda}_t}{l_t} \text{ mit } \vec{\lambda}_t = \vec{\lambda}_{t-1} \odot A^T \vec{\alpha}_{t-1} \text{ und } l_t = p(\vec{x}_t | \vec{x}_{1:t-1}) = \|\vec{\lambda}_t\|_1 \text{ für } t \geq 2$$

wobei $\vec{\lambda}_t$ die lokale Likelihood mit Komponenten $\lambda_{t,k} = p_{\text{obs}}(\vec{x}_t | z_t = k)$ bezeichnet

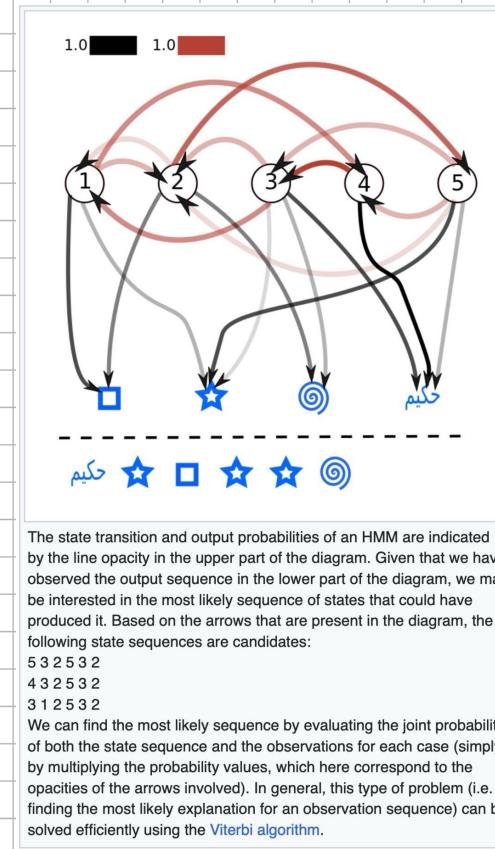
Bew. Filterung:

$$(1) p(z_t = k | \vec{x}_1) \stackrel{\text{Bayes}}{=} \frac{p(\vec{x}_1 | z_t = k) p(z_t = k)}{p(\vec{x}_1)} \stackrel{\text{t.t. Wklt}}{=} \frac{p(\vec{x}_1 | z_t = k) p(z_t = k)}{\sum_{k=1}^K p(\vec{x}_1 | z_t = k) p(z_t = k)} = \frac{(\vec{\lambda}_1 \odot \vec{\pi})_k}{\|\vec{\lambda}_1 \odot \vec{\pi}\|_1} = \frac{\vec{\alpha}_{1,k}}{\|\vec{\alpha}_1\|_1}$$

$$(2) \text{ mit } p(\vec{x}_t | z_t = k, \vec{x}_{1:t-1}) = p_{\text{obs}}(\vec{x}_t | z_t = k), p(\vec{x}_t | \vec{x}_{1:t-1}) = \sum_{k=1}^K p_{\text{obs}}(\vec{x}_t | z_t = k) p(z_t = k | \vec{x}_{1:t-1}), p(z_t = k | \vec{x}_{1:t-1}) = \sum_{k=1}^K p(z_t = k, z_{t-1} = k' | \vec{x}_{1:t-1}) \\ = \sum_{k=1}^K p(z_t = k | z_{t-1} = k') p(z_{t-1} = k' | \vec{x}_{1:t-1}) \stackrel{\text{Markov}}{=} \sum_{k=1}^K p_{\text{trans}}(z_t = k | z_{t-1} = k') p(z_{t-1} = k' | \vec{x}_{1:t-1}) = \sum_{k=1}^K p_{\text{trans}}(z_t = k, z_{t-1} = k' | \vec{x}_{1:t-1}) \text{ folgt:}$$

$$p(z_t = k | \vec{x}_{1:t}) = p(z_t = k | \vec{x}_t, \vec{x}_{1:t-1}) \stackrel{\text{Bayes}}{=} \frac{p(\vec{x}_t | z_t = k, \vec{x}_{1:t-1}) p(z_t = k | \vec{x}_{1:t-1})}{p(\vec{x}_t | \vec{x}_{1:t-1})} = \frac{p_{\text{obs}}(\vec{x}_t | z_t = k) \sum_{k=1}^K p_{\text{trans}}(z_t = k, z_{t-1} = k' | \vec{x}_{1:t-1})}{\sum_{k=1}^K p_{\text{obs}}(\vec{x}_t | z_t = k') p(z_t = k' | \vec{x}_{1:t-1})} \\ = \frac{p_{\text{obs}}(\vec{x}_t | z_t = k) \sum_{k=1}^K A_{k,k'} \alpha_{t-1,k'}}{\sum_{k=1}^K p_{\text{obs}}(\vec{x}_t | z_t = k') p(z_t = k' | \vec{x}_{1:t-1})} = \frac{(\vec{\lambda}_t \odot A^T \vec{\alpha}_{t-1})_k}{\|\vec{\lambda}_t \odot A^T \vec{\alpha}_{t-1}\|_1} = \frac{\vec{\alpha}_{t,k}}{\|\vec{\alpha}_t\|_1}$$

aus gefilterten Verteilung kann beispielsweise Voraussageverteilung $p(\vec{x}_{t+k} | \vec{x}_{1:t})$ oder auch die marginale Verteilung $p(\vec{x}_{1:T})$ berechnen



Glättung

Def. Glättung. Berechnung von $p(z_t | \vec{x}_{1:T})$ resp. $\vec{\gamma}_t$ gegeben alle verfügbaren Beobachtungen $\vec{x}_{1:T}$

Th. Glättung. für ein HMM mit der Transitionsmatrix A und der Initialverteilung $\vec{\pi}$ kann man $\vec{\gamma}_t$ wie folgt rekursiv berechnen:

(0) $\vec{\alpha}_t$ gemäß Filterung vorwärts in der Zeit für $t = 1, \dots, T$ berechnen

(1) $\vec{\beta}_t = A(\vec{\lambda}_{t+1} \odot \vec{\beta}_{t+1})$ mit Startwert $\vec{\beta}_T = \vec{1}_k$ rückwärts in der Zeit für $t = T, \dots, 1$ berechnen → sicheres Ergebnis

(2) $\vec{\gamma}_t = \frac{\vec{\gamma}_t}{n_t}$ mit $\vec{\gamma}_t = \vec{\alpha}_t \odot \vec{\beta}_t$ und $n_t = \|\vec{\gamma}_t\|_1$ für alle t berechnen

wobei $\vec{\lambda}_t$ die lokale Likelihood mit Komponenten $\lambda_{t,k} = p_{\text{obs}}(\vec{x}_t | z_t = k)$ bezeichnet

Bew. Glättung.

$$p(z_t = k | \vec{x}_{1:T}) = p(z_t = k | \vec{x}_{t+1:T}, \vec{x}_{1:t}) \stackrel{\text{Bayes}}{=} \frac{p(\vec{x}_{t+1:T} | z_t = k, \vec{x}_{1:t}) p(z_t = k | \vec{x}_{1:t})}{p(\vec{x}_{t+1:T} | \vec{x}_{1:t})} = \frac{p(\vec{x}_{t+1:T} | z_t = k) p(z_t = k | \vec{x}_{1:t})}{\sum_{k=1}^K p(\vec{x}_{t+1:T} | z_t = k') p(z_t = k' | \vec{x}_{1:t})}$$

$$p(z_t = k | \vec{x}_{1:t}) =: \alpha_{t,k} \quad (\text{Filterung})$$

$$p(\vec{x}_{t+1:T} | z_t = k) =: \beta_{t,k}$$

$$\begin{aligned} \beta_{t,k} &= p(\vec{x}_{t+2:T}, \vec{x}_{t+1}, z_t = k) \stackrel{\text{def. Wkst}}{=} \sum_{k''=1}^K p(\vec{x}_{t+2:T}, \vec{x}_{t+1}, z_{t+1}, z_{t+2} = k'' | z_t = k) \stackrel{\text{Multipl.}}{=} \sum_{k''=1}^K p(\vec{x}_{t+2:T} | \vec{x}_{t+1}, z_{t+1} = k'') p(\vec{x}_{t+1} | z_{t+1} = k'', z_t = k) p(z_{t+1} = k'' | z_t = k) \\ &= \sum_{k''=1}^K p(\vec{x}_{t+2:T} | z_{t+1} = k'') p_{\text{obs}}(\vec{x}_{t+1} | z_{t+1} = k'') p_{\text{trans}}(z_{t+1} = k'' | z_t = k) = \sum_{k''=1}^K \beta_{t+1,k''} \lambda_{t+1,k''} A_{k,k''} = \sum_{k''=1}^K A_{k,k''} \lambda_{t+1,k''} \beta_{t+1,k''} = (A(\vec{\lambda}_{t+1} \odot \vec{\beta}_{t+1}))_k \\ \Rightarrow p(z_t = k | \vec{x}_{1:T}) &= \frac{p(\vec{x}_{t+1:T} | z_t = k) p(z_t = k | \vec{x}_{1:t})}{\sum_{k=1}^K p(\vec{x}_{t+1:T} | z_t = k') p(z_t = k' | \vec{x}_{1:t})} = \frac{\beta_{t,k} \alpha_{t,k}}{\sum_{k=1}^K \beta_{t,k} \alpha_{t,k}} = \frac{(\vec{\alpha}_t \odot \vec{\beta}_t)_k}{\|\vec{\alpha}_t \odot \vec{\beta}_t\|_1} = \frac{\vec{\gamma}_{t,k}}{\|\vec{\gamma}_t\|_1} \end{aligned}$$

Glättung wird insbesondere für das Lernen benötigt

sowohl die Kosten Filter-Rekursion als auch die der Rückwärtsrekursion hängen linear von T und quadratisch von K ab und somit haben auch die Kosten der Glättung, also die Berechnung aller Verteilungen $p(z_t | \vec{x}_{1:T})$ für $1 \leq t \leq T$, insgesamt diese Abhängigkeit

mithilfe dieser beiden Rekursionen können beispielsweise die geglättete Nachfolgerverteilung $p(z_t, z_{t+1} | \vec{x}_{1:T})$ (resp. $\underline{\xi}_{t,t+1} \in \mathbb{R}^{k \times k}$) oder durch die wahrscheinlichste Abfolge $\vec{z}_{1:T}$ gegeben $\vec{x}_{1:T}$ (Viterbi-Algorithmus, linearer Aufwand in T) berechnet werden

Möglichkeiten für das Lernen der Parameter $\vec{\theta} = (A, \vec{\pi}, \vec{\theta}_{\text{obs}})$ gegeben die Sequenzen $D = \{\vec{x}_{j,1:T_j} \mid 1 \leq j \leq N\}$:

- EM-Algorithmus, welcher in diesem Kontext auch Baum-Welch-Algorithmus genannt wird
- SGD (v.a. bei grossen Differenzen)
- Bayesische Ansätze, welche auf Variationsapproximationen (Variational Bayes-EM) oder Monte-Carlo-Methoden basieren (da nicht exakt berechenbar)

Lineare dynamische Systeme (LDSs)

Transitionsmodell:

$$p_{\text{trans}}(\vec{z}_{t+1} | \vec{z}_t, \vec{u}_t) = \mathcal{N}(\vec{z}_{t+1} | E\vec{z}_t + B\vec{u}_t, Q) \quad \text{mit } E \in \mathbb{R}^{dxd}, B \in \mathbb{R}^{dxc}, Q \in \mathbb{R}^{dxd}$$

resp. $\vec{z}_t = E\vec{z}_{t-1} + B\vec{u}_{t-1} + \vec{q}_t$ mit Prozessrauschen $\vec{q}_t \sim \mathcal{N}(0, Q)$

Beobachtungsmodell:

$$p_{\text{obs}}(\vec{x}_t | \vec{z}_t) = \mathcal{N}(\vec{x}_t | H\vec{z}_t + D\vec{u}_t, R) \quad \text{mit } H \in \mathbb{R}^{nxl}, D \in \mathbb{R}^{nxc}, R \in \mathbb{R}^{nxn}$$

resp. $\vec{x}_t = H\vec{z}_t + D\vec{u}_t + \vec{v}_t$ mit Messrauschen $\vec{v}_t \sim \mathcal{N}(0, R)$

manchmal hängen einige der Parameter E, B, Q, H, D, R von der Zeit ab, was eine nicht-stationäre Verarbeitung liefert

ohne Ergebnissen ist kann die Darstellung natürlich unvollständig werden

Kalman Filter und Glätter

wie bei HMMs können auch bei LDSs explizite (rekursive) Formeln

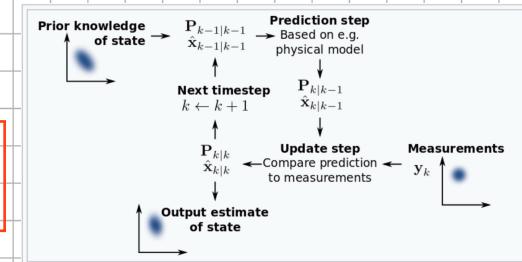
$$p(\vec{z}_t | \vec{x}_{1:t}, \vec{u}_t) = \mathcal{N}(\vec{z}_t | \vec{M}_{t|t}, \Sigma_{t|t})$$

$$p(\vec{z}_t | \vec{x}_{1:T}, \vec{u}_T) = \mathcal{N}(\vec{z}_t | \vec{M}_{t|T}, \Sigma_{t|T})$$

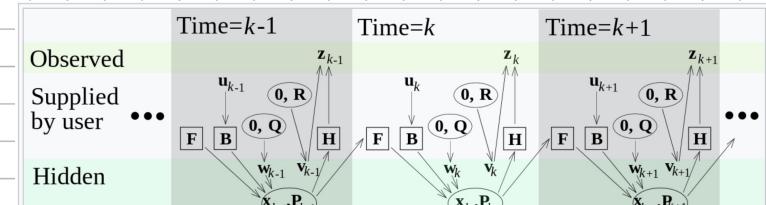
für die Filter- und Glätter-Verteilungen berechnet werden

für das Lernen gibt es mehrere Möglichkeiten:

- EM, wobei die Verteilungen über die verborgenen Zustände mithilfe der Kalman-Filter und Kalman-Glättung rekursiv berechnet werden
- Bayesische Ansätze, welche auf Variationsapproximationen (Variational Bayes-EM) oder Monte-Carlo-Methoden basieren



The Kalman filter keeps track of the estimated state of the system and the variance or uncertainty of the estimate. The estimate is updated using a state transition model and measurements. $\hat{x}_{k|k-1}$ denotes the estimate of the system's state at time step k before the k -th measurement y_k has been taken into account; $P_{k|k-1}$ is the corresponding uncertainty.



Model underlying the Kalman filter. Squares represent matrices. Ellipses represent multivariate normal distributions (with the mean and covariance matrix enclosed). Unenclosed values are vectors. For the simple case, the various matrices are constant with time, and thus the subscripts are not used, but Kalman filtering allows any of them to change each time step.

Neuronale Netze

Basisfunktionen adaptiv, Parameter der Basisfamilie werden im Lernprozess erhoerzt

Feedforward Neural Net (FNN)

eine verborgene Schicht:

$$\vec{f}_{2,\vec{\theta}_2} \circ \vec{f}_{1,\vec{\theta}_1} : \mathbb{R}^{n_{\text{in}}} \rightarrow \mathbb{R}^{n_{\text{hidden}}} \rightarrow \mathbb{R}^{n_{\text{out}}}$$

$$\vec{x} \mapsto \vec{f}_{2,\vec{\theta}_2}(\vec{f}_{1,\vec{\theta}_1}(\vec{x})) = \vec{\varphi}_2(\sum_{i=1}^n w_i (\vec{\varphi}_1(w_i \vec{x} + \vec{b}_1)) + \vec{b}_2)$$

ein solches Modell (1 verborgene Schicht resp 1 Basisexpansion) kann schon mit einer einfachen Aktivierungsfunktion (z.B. ReLU) und der linearen Ausgabeaktivierung jede stetige Funktion auf einer kompakten Menge gleichmaessig approximieren, wenn n_{hidden} genügend gross gewählt wird

es hat sich jedoch als lohnend erweisen iterierte Basisexpansion zu verwenden anstatt breiter zu wählen:

$$\vec{f}_{\vec{\theta}} : \mathbb{R}^{n_{\text{in}}} \rightarrow \dots \rightarrow \mathbb{R}^{n_{\text{out}}} \quad \text{mit } \vec{\theta} = (\vec{\theta}_1, \dots, \vec{\theta}_k) \in \mathbb{R}^{p_1} \times \dots \times \mathbb{R}^{p_k}$$

$$\vec{x} \mapsto \vec{f}_{k,\vec{\theta}_k} \circ \dots \circ \vec{f}_{1,\vec{\theta}_1}(\vec{x}) \quad \text{wobei: } \vec{f}_{j,\vec{\theta}_j}(\vec{a}) = \vec{\varphi}_j(w_j \vec{a} + \vec{b}_j)$$

die Ausgabe-Aktivierungsfunktion $\vec{\varphi}_k$ wird durch die Anwendung bestimmt

dense layer: volle Matrix und Verbindungsleiter sind adaptive Parameter

tiefe FNNs mit dichten layers nicht unmittelbar breite sind sehr flexibel (no bias, high variance)
→ Regularisierung und viel Trainingsdaten wichtig

obwohl durch Verwendung von partiellierten und strukturierten Matrizen auch komplexere Architekturen mit Faltungen, Skip-Connections durch das Schichtmodell dargestellt werden können ist es irgendwann sowohl für die Modellierung als auch die Implementierung natuerlicher das Netz als gerichteten, zyklischen Graphen von Basisexpansionen anzufassen

Backpropagation

$$\vec{f}(\vec{x}, \vec{\theta}_1, \dots, \vec{\theta}_k) = \vec{f}_k \left(\underbrace{\vec{f}_{k-1} \left(\dots \vec{f}_2 \left(\underbrace{\vec{f}_1(\vec{x}, \vec{\theta}_1), \vec{\theta}_2}_{\vec{a}_1} \right) \dots, \vec{\theta}_{k-1} \right)}_{\vec{a}_{k-1}}, \vec{\theta}_k \right)$$

Ableitung nach $\vec{\theta}_i$:

$$\frac{\partial}{\partial \vec{\theta}_i} \vec{f}(\vec{x}, \vec{\theta}_1, \dots, \vec{\theta}_k) = \frac{\partial}{\partial \vec{a}} \vec{f}_k(\vec{a}_{k-1}, \vec{\theta}_k) \dots \frac{\partial}{\partial \vec{a}} \vec{f}_{i+1}(\vec{a}_i, \vec{\theta}_{i+1}) \frac{\partial}{\partial \vec{\theta}} \vec{f}_i(\vec{a}_{i-1}, \vec{\theta}_i) \quad \text{Kettenregel}$$

diese Matrizen sollte man von links nach rechts berechnen (effizienter)
→ Ableitungsinformationstransfer rückwärts bis Netz

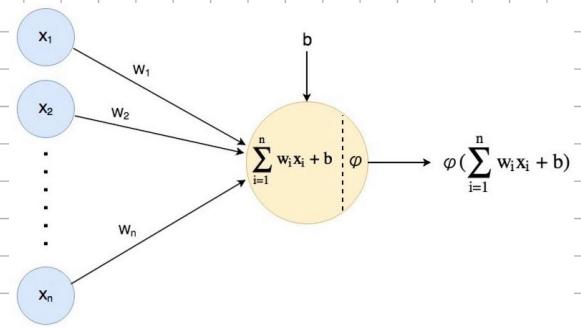
Trainieren

diskriminatives, probabilistisches Modell mit ML:

$$NLL(\vec{\theta}) = -\ln p(t|X, \vec{\theta}) = -\sum_{k=1}^{n_p} \ln p(t_k | \vec{x}_k, \vec{\theta})$$

allfällige Regularisierungsterme (MAP) können darüber addiert werden

Minimierung davon nicht analytisch berechenbar außerdem ist NLL nicht konkav, hat i.d.R. viele lokale nicht optimale Minima



WGL-Bereich - Gradient:

$$\frac{d\text{ML}}{d\vec{\theta}}(\vec{\theta}) \cong \frac{n_0}{n_0} \frac{d\text{ML}_{\text{B}}}{d\vec{\theta}}(\vec{\theta}) = -\frac{n_0}{n_0} \sum_{i=1}^{n_0} \frac{\partial}{\partial \vec{\theta}} \ln p(t_{k_i} | \vec{x}_{k_i}, \vec{\theta})$$

SGD - Iteration:

$$\vec{\theta}_{t+1} = \vec{\theta}_t - \eta_t \nabla_{\vec{\theta}} f(\vec{\theta}_t)$$

Robbins-Monro:

$$\sum_{t=1}^{\infty} \eta_t = \infty \quad \text{und} \quad \sum_{t=1}^{\infty} \eta_t^2 < \infty \quad \text{bekannte hinreichende Bedingungen für lokale Konvergenz}$$

Möglichkeiten, inhärente Flexibilität zu zähmen (Regularisierung):

- Prior-V. über Parameter, häufig unkorrelierte und zentrierte Normalv.
- Early Stopping: Minimieren steigen sobald Verlust auf Evaluationsdaten ansteigt
- Faltungsschichten (Form von Parameter-Sharing)
- Bayes
- viele weitere...