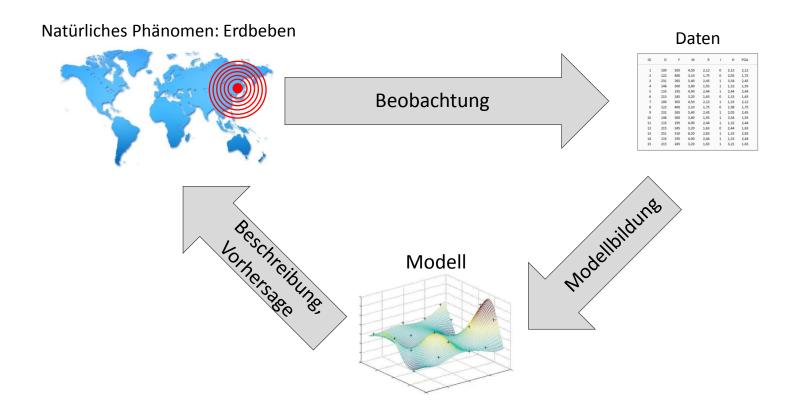


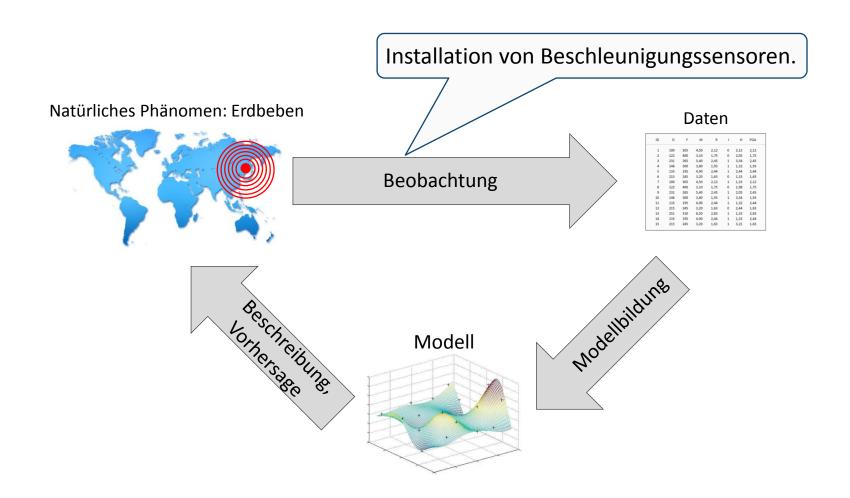
Maschinelles Lernen zur Modellbildung in den Naturwissenschaften

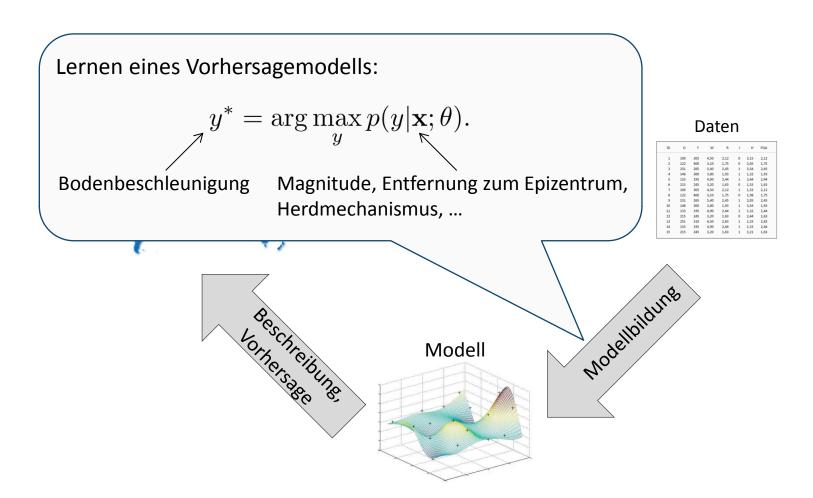
Niels Landwehr
Institut für Informatik, Universität Potsdam

Vorstellung Forschungsgruppe

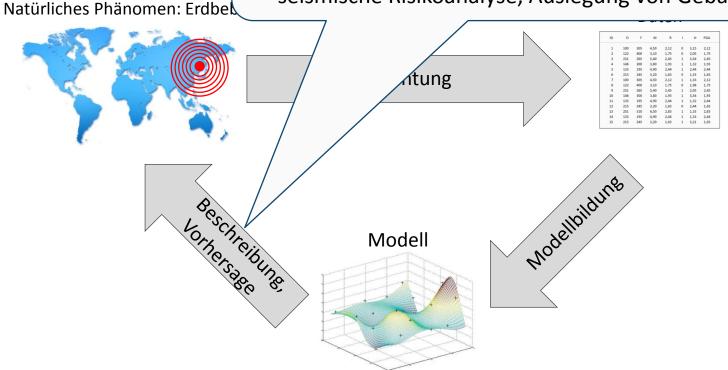
- Nachwuchsgruppe "Machine Learning and Scientific Data Analysis"
 - Emmy Noether-Programm der DFG, seit Februar 2013
 - Forschungsprogramm: Fragestellungen des maschinellen Lernens im Zusammenhang mit naturwissenschaftlicher Modellbildung
 - Kooperation mit Wissenschaftlern aus der Geophysik (AG Frank Scherbaum) und der kognitiven Psychologie (AG Reinhold Kliegl)
 - Angegliedert an die Arbeitsgruppe "Maschinelles Lernen" (Tobias Scheffer), Institut für Informatik, Universität Potsdam
- Heute: Überblick über Forschungsprogramm und Vorarbeiten



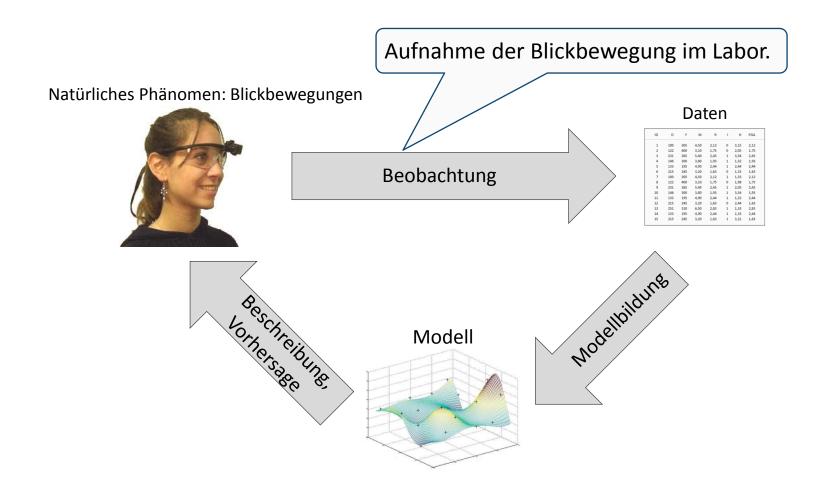




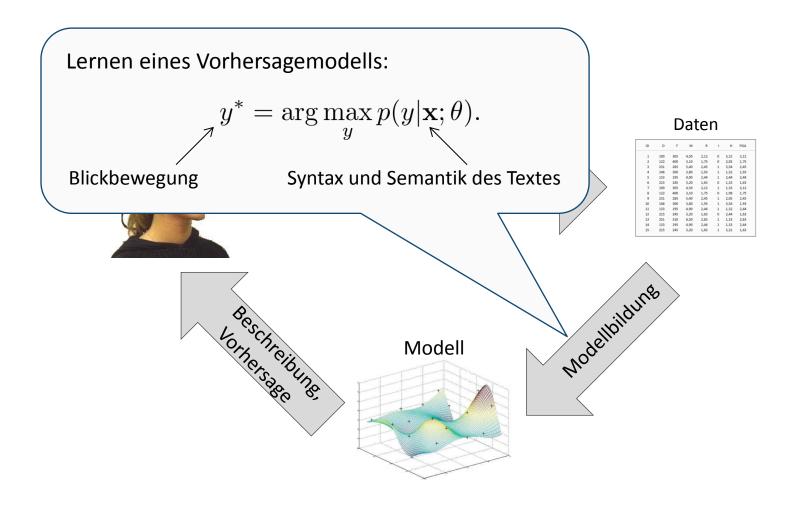
- Einsicht in natürliche Prozesse.
- Vorhersagen lokaler Bodenbeschleunigung: seismische Risikoanalyse, Auslegung von Gebäuden.



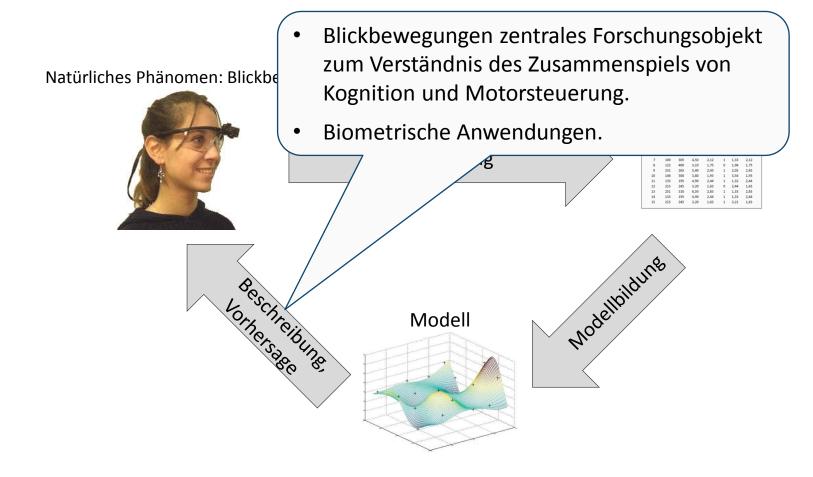
Modellbildung aus Experimentaldaten

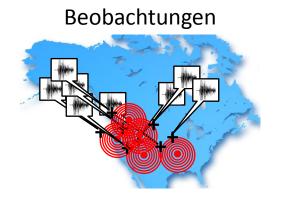


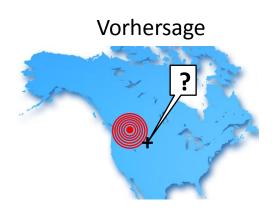
Modellbildung aus Experimentaldaten

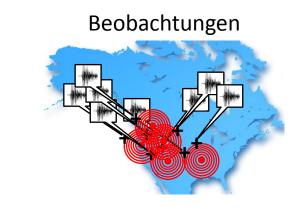


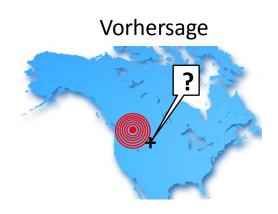
Modellbildung aus Experimentaldaten

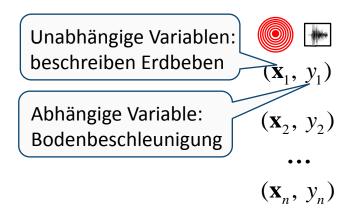


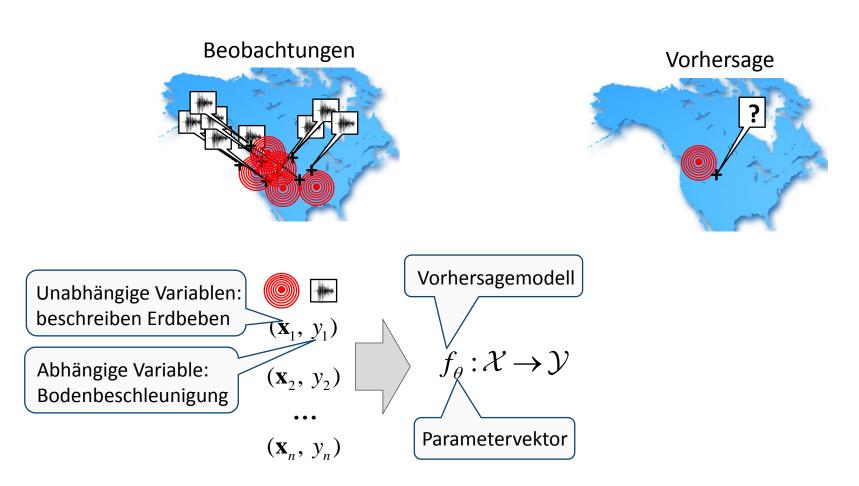


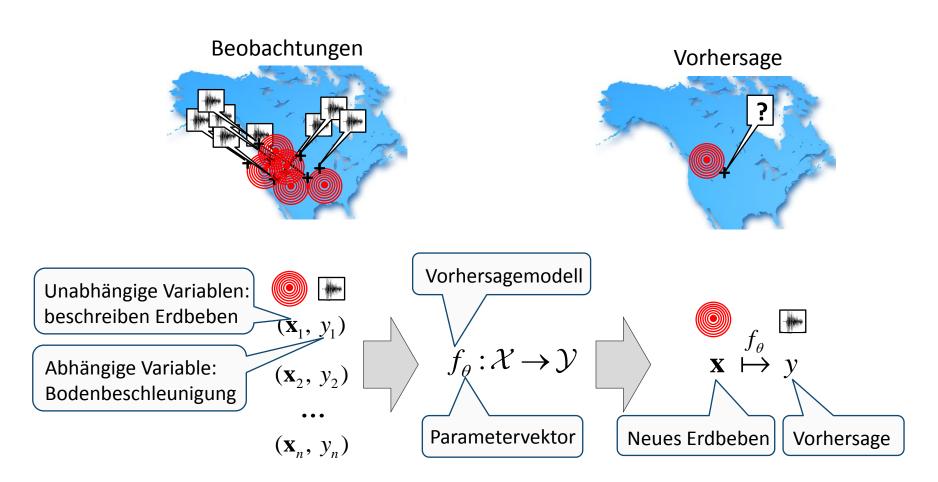












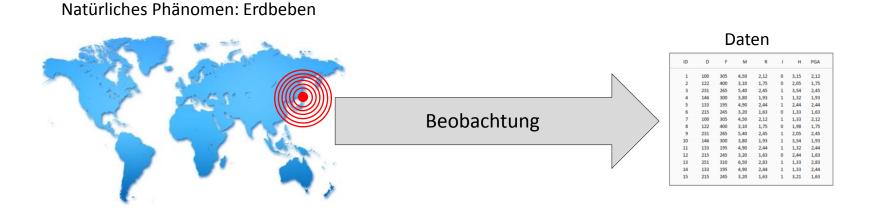
Maschinelles Lernen: Modellinferenz

- Zentraler Schritt: Inferenz eines geeigneten Modells gegeben Daten
- Viele Modellklassen und Verfahren
 - Probabilistische Modelle mit max likelihood /max aposteriori / Bayes'scher Inferenz
 - parametrische / nichtparametrische (kernelisierte) Modelle
- Lösen eines Optimierungsproblems
 - Suchproblem über komplexen Raum möglicher Modelle
 - Algorithmisch aufwendig
 - Approximative numerische Verfahren

Modellbildung: statistische Annahmen

I.I.D. Annahme:

- Alle Datenpunkte folgen einer einheitlichen Verteilung $p(\mathbf{x}, y)$.
- Datenpunkte sind unabhängige Stichproben aus $p(\mathbf{x}, y)$.

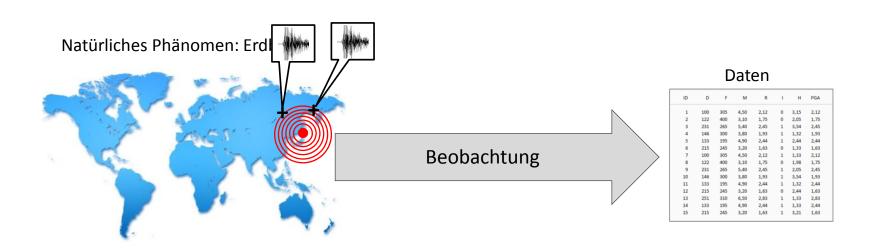


Abhängige Datenpunkte

• **Abhängige Datenpunkte**: Verschiedene Sensoren registrieren dasselbe Ereignis (beispielsweise verschiedene Messstationen eine seismische Welle).

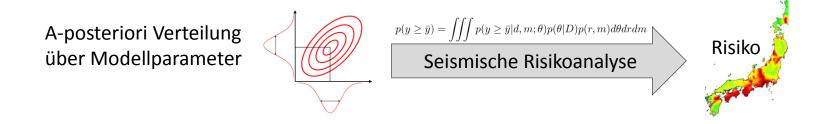


Widerspruch zur **Unabhängigkeitsannahme**.



Abhängige Datenpunkte

 Nicht berücksichtigte Datenabhängigkeiten führen zu einer Unterschätzung der nach der Modellbildung verbleibenden Unsicherheit.



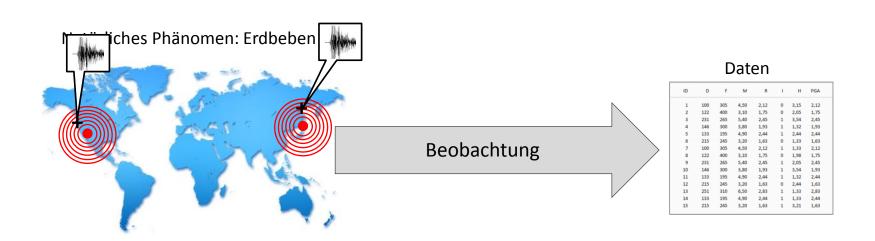
• Falsche Schlussfolgerungen: beispielsweise Unterschätzung der Wahrscheinlichkeit extremer Ereignisse.

Verteilungsverschiebungen

• Verteilungsverschiebungen: Verteilung der Daten abhängig vom Ort der Messung (Aufnahme seismischer Daten in unterschiedlichen Regionen).

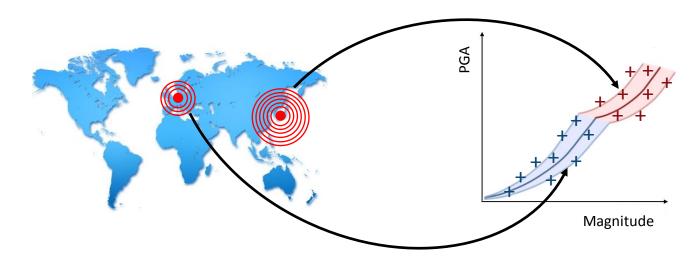


Widerspruch zur Annahme einer einheitlichen Datenverteilung.



Verteilungsverschiebungen

- Bildung von Modellen für seismisches Risiko in Deutschland:
 - In Deutschland stehen nur Beobachtungen niedriger Intensität zur Verfügung.
 - Extrapolation zu Erdbeben h\u00f6herer Intensit\u00e4t mit Daten aus anderen Regionen.



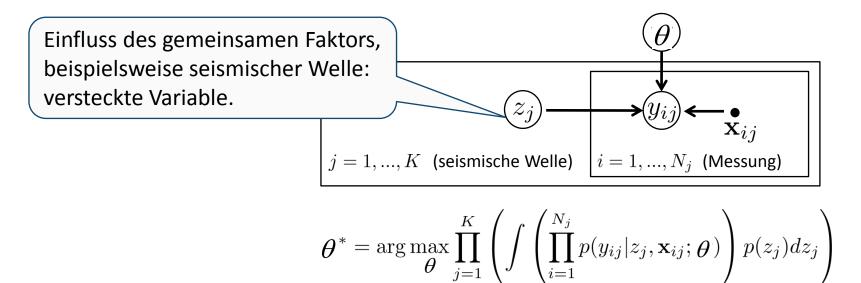
Verteilungsunterschiede können zu falschen Vorhersagen führen.

Übersicht

- Einleitung & Motivation
- Lernen unter Datenabhängigkeiten
- Lernen unter Verteilungsverschiebungen
- Modellevaluierung

Abhängigkeiten: Stand der Forschung

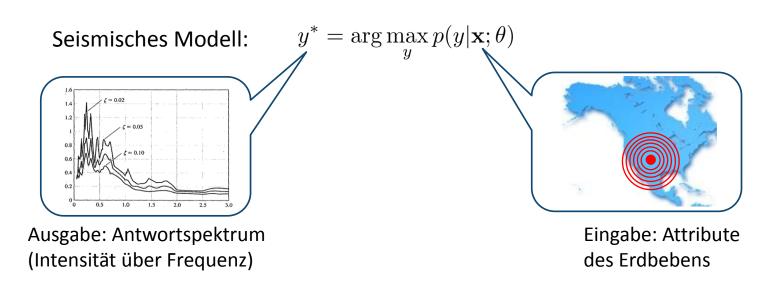
• Stand der Forschung zur Modellierung abhängiger Daten: Mixed Models.



• Bei der Parameterschätzung wird über die versteckten Variablen integriert

Abhängigkeiten: Ziele

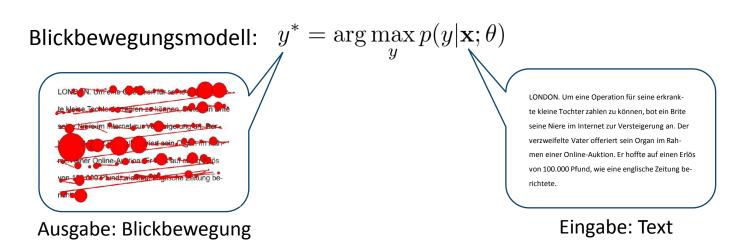
 Mixed Models nur bekannt für einfache Ausgabevariablen, in Anwendungen ist Ausgabevariable aber oft komplex strukturiert.



- Ziel: Modelle für komplexe Ausgaben unter Datenabhängigkeiten.
- Aufwendige Optimierungsprobleme bei Modellbildung und Vorhersage.

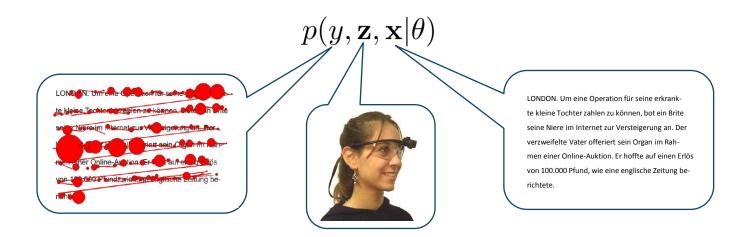
Abhängigkeiten: Ziele

 Mixed Models nur bekannt für einfache Ausgabevariablen, in Anwendungen ist Ausgabevariable aber oft komplex strukturiert.

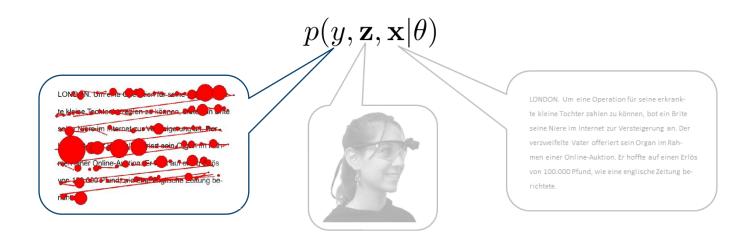


- Ziel: Modelle für komplexe Ausgaben unter Datenabhängigkeiten.
- Aufwendige Optimierungsprobleme bei Modellbildung und Vorhersage.

• **Ziel**: gemeinsame Modelle von Blickbewegungen, Textinhalt und Attributen des Lesers.

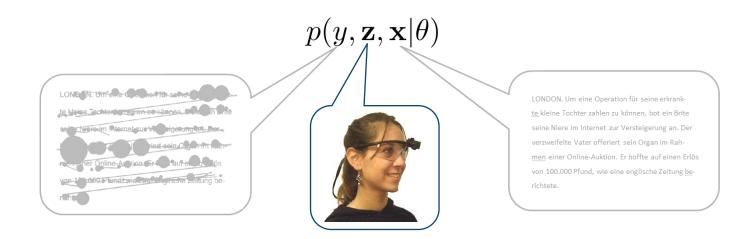


• **Ziel**: gemeinsame Modelle von Blickbewegungen, Textinhalt und Attributen des Lesers.



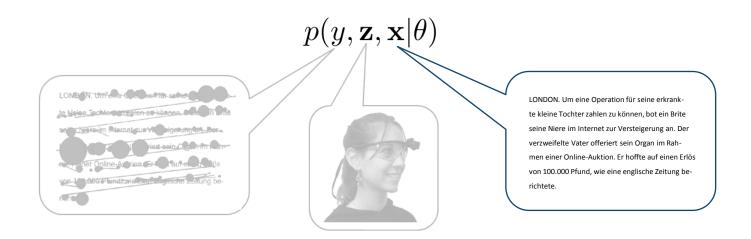
Vorhersage der Blickbewegung (Psychologie).

 Ziel: gemeinsame Modelle von Blickbewegungen, Textinhalt und Attributen des Lesers.



- Rückschlüsse auf die Identität des Lesers (Biometrie).
- Inferenz von Eigenschaften des Lesers: Muttersprache, Alter, Kompetenz (Textpersonalisierung, Kompetenztests).

• **Ziel**: gemeinsame Modelle von Blickbewegungen, Textinhalt und Attributen des Lesers.



 Rückschluss auf die Lesbarkeit von Texten, durch Vorhersage der zu erwartenden Blickbewegungen (Werkzeuge zur Textanalyse).

Übersicht

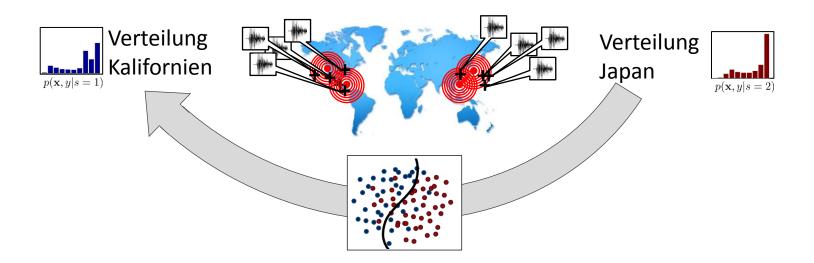
- Einleitung & Motivation
- Lernen unter Datenabhängigkeiten
- Lernen unter Verteilungsverschiebungen
- Modellevaluierung

Transferlernen



- Messungen an unterschiedlichen Orten: Verteilungsverschiebungen.
- Standardansatz: Verteilungsunterschiede ignorieren und Daten poolen.
- Transferlernen: Verteilungsunterschiede zwischen Teilmengen der Daten korrigieren.

Transferlernen: Vorstudie

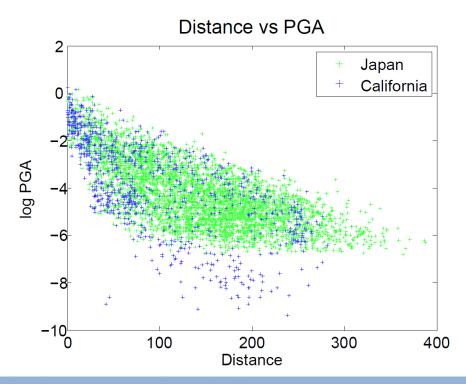


- Vorstudie zum Transferlernen, mit AG Frank Scherbaum, Geophysik.
- Ansatz der Verteilungskorrektur [Bickel, Sawade, Scheffer; 2008]
 - Annahme: jeweils einheitliche Verteilung für Daten aus Kalifornien, Japan.
 - Explizites Modell des Verteilungsunterschieds aus Daten lernen.

Erdbebendaten Japan/Kalifornien

- Bodenbeschleunigungsdaten aus Kalifornien (841) und aus Japan (3542)
- Ziel: Modell für Kalifornien ("Zielverteilung"/"Zieldaten")
- Zusätzlich können Daten aus Japan verwendet werden ("Hilfsdaten")

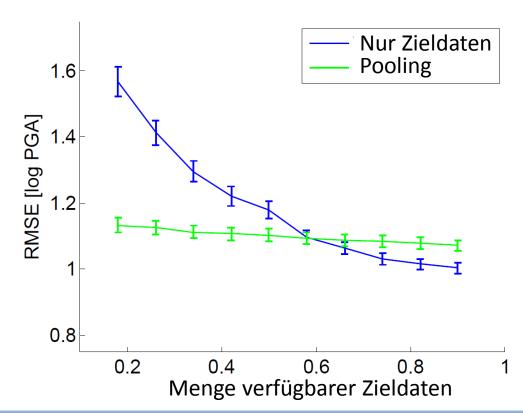
Leicht unterschiedliche Verteilung über Daten aus Kalifornien und Japan



Vergleich "Nur Zieldaten", "Pooling"

- Baseline 1: "Nur Zieldaten" (ignoriert Hilfsdaten)
- Baseline 2: "Pooling" (ignoriert Verteilungsunterschied)

Fehler in Abhängigkeit der zum Training verfügbaren Zieldaten: Bias vs Varianz



Verteilungskorrektur

Transferlernen (Verteilungskorrektur)

Zieldaten (Kalifornien):
$$(\mathbf{x}, y) \sim p(\mathbf{x}, y | t = 1)$$

Hilfsdaten (Japan):
$$(\mathbf{x}, y) \sim p(\mathbf{x}, y | t = 0)$$

• Idee: Umgewichten der Daten, so dass $p(\mathbf{x}, y | t = 0)$ zu $p(\mathbf{x}, y | t = 1)$ korrigiert wird.

Verteilungskorrektur

- Verteilungskorrektur formal
 - Pool folgt Mischverteilung

$$p(\mathbf{x}, y) = \underbrace{p(\mathbf{x}, y \mid t = 1)}_{\text{Zieldaten}} \underbrace{p(t = 1)}_{\text{Anteil Zieldaten}} + \underbrace{p(\mathbf{x}, y \mid t = 0)}_{\text{Hilfsdaten}} \underbrace{p(t = 0)}_{\text{Anteil Hilfsdaten}}$$

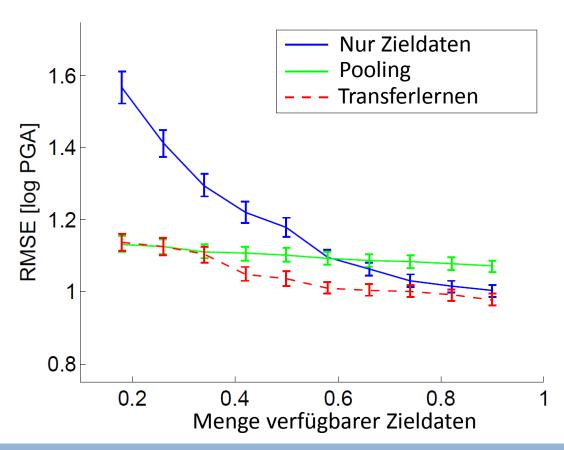
- Wir hätten gerne alle Daten aus Zielverteilung $p(\mathbf{x}, y | t = 1)$
- Daten umgewichten mit $r(\mathbf{x}, y) = \frac{p(\mathbf{x}, y \mid t = 1)}{p(\mathbf{x}, y)}$
- Umgewichtungsfaktoren mit Modell schätzen

$$\mathbf{r}(\mathbf{x}, y) = \frac{p(\mathbf{x}, y \mid t = 1)}{p(\mathbf{x}, y)} = \frac{p(t = 1 \mid \mathbf{x}, y) p(\mathbf{x}, y)}{p(\mathbf{x}, y) p(t = 1)} \propto p(t = 1 \mid \mathbf{x}, y)$$

Diskriminatives Modell: Zieldaten gegen Hilfsdaten (hier: logistische Regression, Poly-Kernel)

Ergebnisse: Testfehler

Vergleich "Nur Zieldaten", "Pooling", und "Transferlernen"

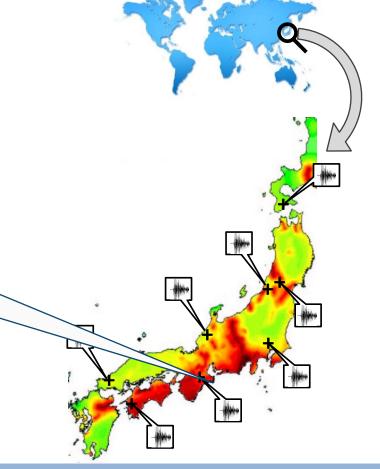


Transferlernen: Ziele

 Aufteilung in Teilmengen Vereinfachung: Verteilung ändert sich kontinuierlich mit dem Ort der Messung.

 Ziel: Modelle, die eine kontinuierliche Verteilungsänderung abbilden.

> Vorhersage für Standort: Aufwendige Inferenzprobleme, approximative Inferenzalgorithmen.

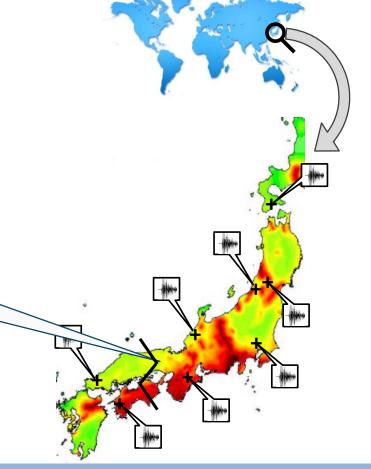


Transferlernen: Ziele

 Aufteilung in Teilmengen Vereinfachung: Verteilung ändert sich kontinuierlich mit dem Ort der Messung.

 Ziel: Modelle, die geophysikalisches Vorwissen berücksichtigen.

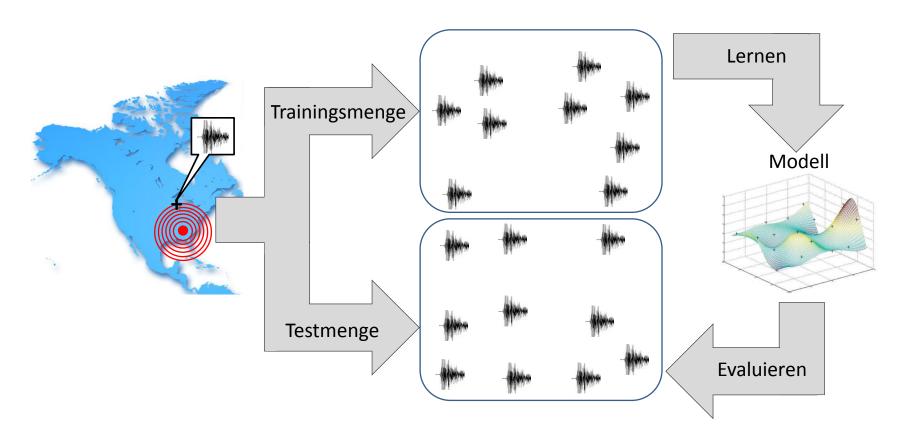
Bekannte seismische Bruchlinie: stärkere Verteilungsänderung erwartet.



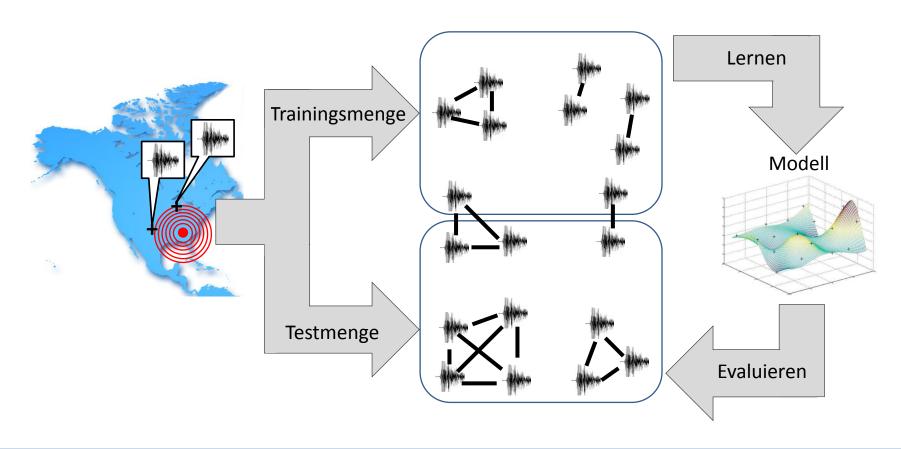
Übersicht

- Einleitung & Motivation
- Lernen unter Datenabhängigkeiten
- Lernen unter Verteilungsverschiebungen
- Modellevaluierung

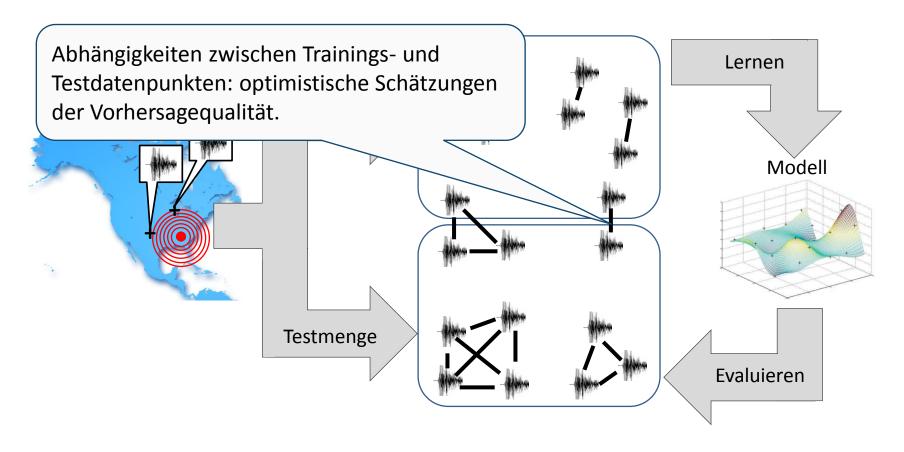
• Existierende Evaluierungsverfahren basieren auf der Annahme unabhängiger und einheitlich verteilter Beobachtungen.



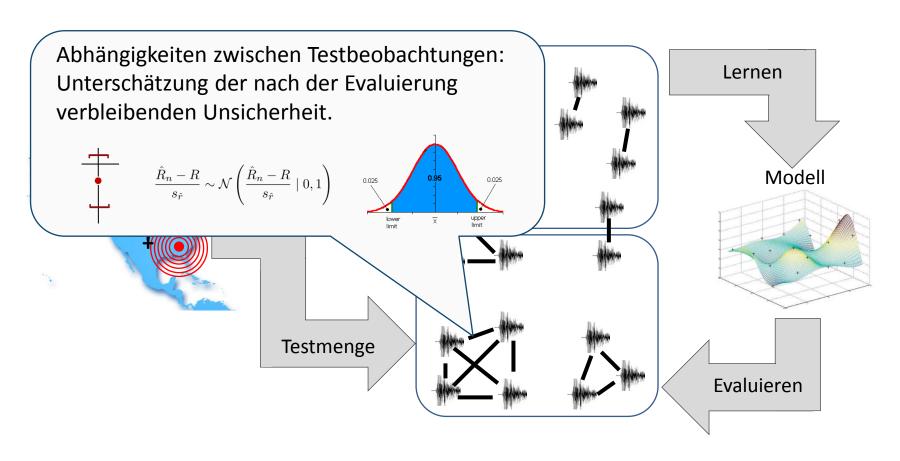
 Abhängige Datenpunkte verfälschen das Ergebnis bekannter Evaluierungsverfahren in verschiedener Weise.



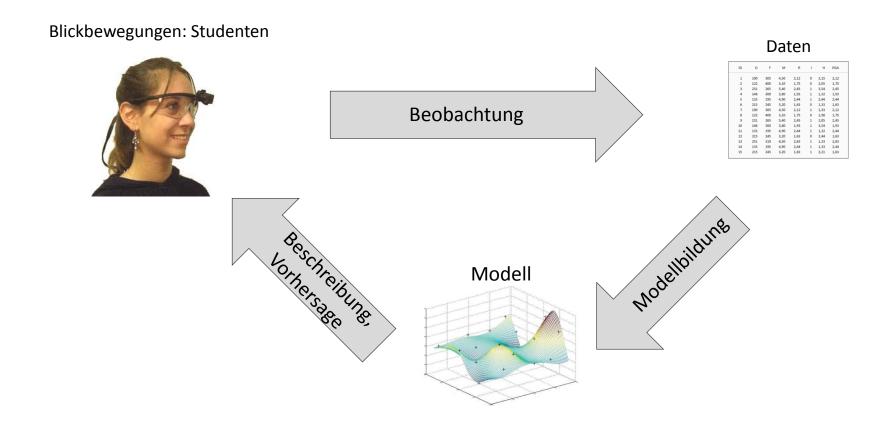
 Abhängige Datenpunkte verfälschen das Ergebnis bekannter Evaluierungsverfahren in verschiedener Weise.



 Abhängige Datenpunkte verfälschen das Ergebnis bekannter Evaluierungsverfahren in verschiedener Weise.



Evaluierung: Verteilungsverschiebungen



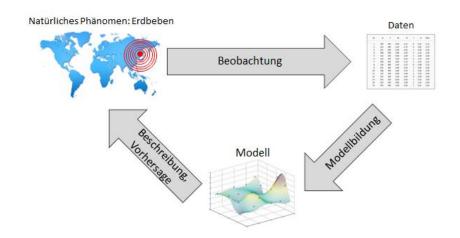
Evaluierung: Verteilungsverschiebungen

Wie gut beschreibt das mit studentischen Probanden gelernte Modell die Blickbewegungen der Gesamtbevölkerung? Blickbewegungen: Studenten Daten Beobachtung Modellaidune Blickbewegungen: Gesamtbevölkerung Modell Beschreibung, Vorhersage

Übersicht

- Einleitung & Motivation
- Lernen unter Datenabhängigkeiten
- Lernen unter Verteilungsverschiebungen
- Modellevaluierung

Zusammenfassung



- Maschinelles Lernen zur Modellbildung aus Beobachtungsdaten in den Naturwissenschaften.
- Besondere Verteilungseigenschaften, insbesondere Abhängigkeiten und Verteilungsverschiebungen.
- Zusammenarbeit mit Arbeitsgruppen aus der Geophysik und der kognitiven Psychologie.