

Wissensentdeckung in Datenbanken

Optimierung, Überanpassung

Nico Piatkowski und Uwe Ligges

Informatik—Künstliche Intelligenz
Computergestützte Statistik
Technische Universität Dortmund

04.05.2017

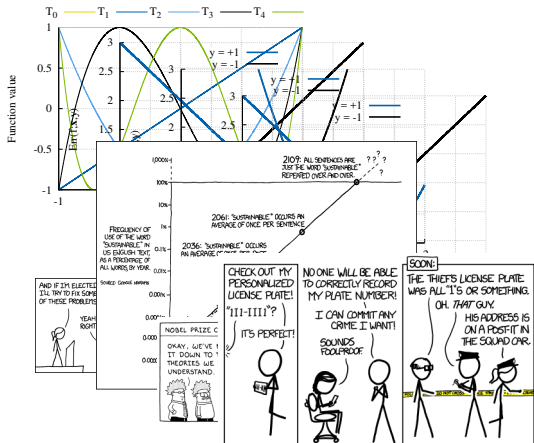
Überblick

Was bisher geschah...

- Modellklassen
- Verlustfunktionen

Heute

- Optimierung
- Overfitting



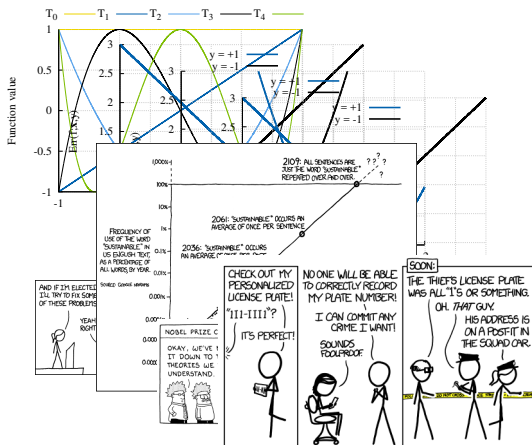
Überblick

Was bisher geschah...

- Modellklassen
- Verlustfunktionen

Heute

- Optimierung
- Overfitting





Differenzierbarkeit

- Die Funktion $f : \mathbb{R}^n \rightarrow \mathbb{R}$ ist *partiell differenzierbar* an β_i , falls

$$\forall \beta \in \mathbb{R}^n : \frac{\partial f(\beta)}{\partial \beta_i} = \lim_{h \rightarrow 0} \frac{f(\beta + h e_i) - f(\beta)}{h}$$

- Der Vektor

$$\nabla f(\beta) = \begin{pmatrix} \frac{\partial f(\beta)}{\partial \beta_1} \\ \frac{\partial f(\beta)}{\partial \beta_2} \\ \dots \\ \frac{\partial f(\beta)}{\partial \beta_n} \end{pmatrix}$$

heißt *Gradient* von f an der Stelle β .

Differenzierbarkeit

- Die Funktion $f : \mathbb{R}^n \rightarrow \mathbb{R}$ ist *partiell differenzierbar* an β_i , falls

$$\forall \beta \in \mathbb{R}^n : \frac{\partial f(\beta)}{\partial \beta_i} = \lim_{h \rightarrow 0} \frac{f(\beta + h e_i) - f(\beta)}{h}$$

- Der Vektor

$$\nabla f(\beta) = \begin{pmatrix} \frac{\partial f(\beta)}{\partial \beta_1} \\ \frac{\partial f(\beta)}{\partial \beta_2} \\ \dots \\ \frac{\partial f(\beta)}{\partial \beta_n} \end{pmatrix}$$

heißt *Gradient* von f an der Stelle β .



Konvexität

- f konvex, gdw.

$$\forall \mathbf{a}, \mathbf{b} \in \mathbb{R}^n : f(\mathbf{a}) \geq f(\mathbf{b}) + \langle \nabla f(\mathbf{b}), \mathbf{a} - \mathbf{b} \rangle$$

- f konvex, gdw.

$$\forall \mathbf{a}, \mathbf{b} \in \mathbb{R}^n : \langle \nabla f(\mathbf{a}) - \nabla f(\mathbf{b}), \mathbf{a} - \mathbf{b} \rangle \geq 0$$

≡ Gradient ist monotoner Operator

- f konvex \Rightarrow Jedes lokale Minimum ist ein globales Minimum

(Neue Literatur im Moodle: Stephen Boyd und Lieven Vandenberghe: Convex Optimization)

Konvexität

- f konvex, gdw.

$$\forall \mathbf{a}, \mathbf{b} \in \mathbb{R}^n : f(\mathbf{a}) \geq f(\mathbf{b}) + \langle \nabla f(\mathbf{b}), \mathbf{a} - \mathbf{b} \rangle$$

- f konvex, gdw.

$$\forall \mathbf{a}, \mathbf{b} \in \mathbb{R}^n : \langle \nabla f(\mathbf{a}) - \nabla f(\mathbf{b}), \mathbf{a} - \mathbf{b} \rangle \geq 0$$

≡ Gradient ist monotoner Operator

- f konvex \Rightarrow Jedes lokale Minimum ist ein globales Minimum

(Neue Literatur im Moodle: Stephen Boyd und Lieven Vandenberghe: Convex Optimization)

Konvexität

- f konvex, gdw.

$$\forall \mathbf{a}, \mathbf{b} \in \mathbb{R}^n : f(\mathbf{a}) \geq f(\mathbf{b}) + \langle \nabla f(\mathbf{b}), \mathbf{a} - \mathbf{b} \rangle$$

- f konvex, gdw.

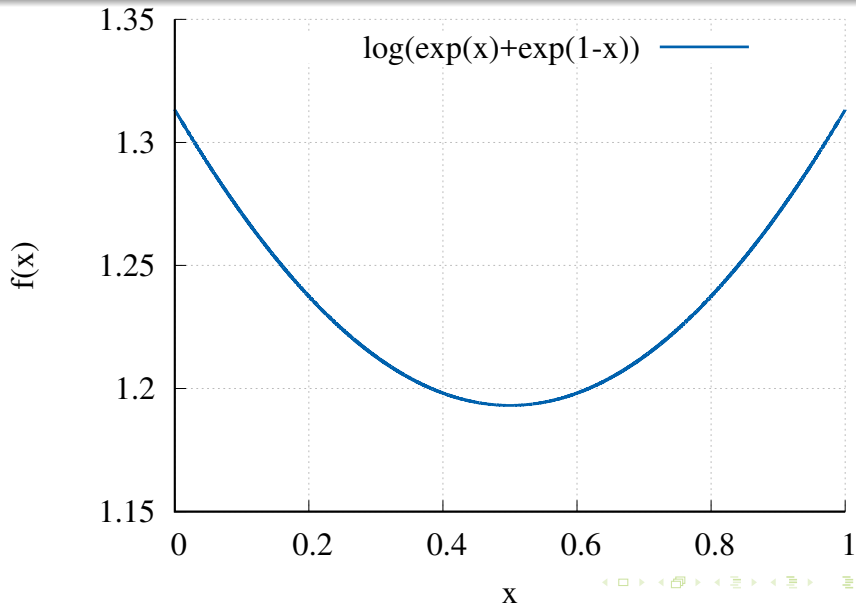
$$\forall \mathbf{a}, \mathbf{b} \in \mathbb{R}^n : \langle \nabla f(\mathbf{a}) - \nabla f(\mathbf{b}), \mathbf{a} - \mathbf{b} \rangle \geq 0$$

≡ Gradient ist monotoner Operator

- f konvex \Rightarrow Jedes lokale Minimum ist ein globales Minimum

(Neue Literatur im Moodle: Stephen Boyd und Lieven Vandenberghe: Convex Optimization)

Konvexität (II)



Lipschitz-Stetigkeit für $\mathbb{R}^n \rightarrow \mathbb{R}^m$

- f heißt *Lipschitz-Stetig mit Konstante* $L > 0$, falls

$$\forall \mathbf{a}, \mathbf{b} \in \mathbb{R}^n : |f(\mathbf{a}) - f(\mathbf{b})| \leq L \|\mathbf{a} - \mathbf{b}\|_2$$

Hilfreich bei $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$:

- Falls

$$K = \sup_{\mathbf{c} \in \mathbb{R}^n} \|\nabla f(\mathbf{c})\|_2 < \infty ,$$

dann ist f Lipschitz stetig mit Konstante $K > 0$.

Lipschitz-Stetigkeit für $\mathbb{R}^n \rightarrow \mathbb{R}^m$

- f heißt *Lipschitz-Stetig mit Konstante* $L > 0$, falls

$$\forall \mathbf{a}, \mathbf{b} \in \mathbb{R}^n : |f(\mathbf{a}) - f(\mathbf{b})| \leq L \|\mathbf{a} - \mathbf{b}\|_2$$

Hilfreich bei $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$:

- Falls

$$K = \sup_{\mathbf{c} \in \mathbb{R}^n} \|\nabla f(\mathbf{c})\|_2 < \infty ,$$

dann ist f Lipschitz stetig mit Konstante $K > 0$.

Konvexität + Lipschitz stetige Gradienten

$$\forall \mathbf{a}, \mathbf{b} \in \mathbb{R}^n : \|\nabla f(\mathbf{a}) - \nabla f(\mathbf{b})\|_2 \leq L\|\mathbf{a} - \mathbf{b}\|_2$$

Multiplikation mit $\|\mathbf{a} - \mathbf{b}\|_2$ und Cauchy-Schwarz Ungl.:

$$\forall \mathbf{a}, \mathbf{b} \in \mathbb{R}^n : \langle \nabla f(\mathbf{a}) - \nabla f(\mathbf{b}), \mathbf{a} - \mathbf{b} \rangle \leq L\|\mathbf{a} - \mathbf{b}\|_2^2$$

Substitution von $g(\mathbf{x}) = (L/2)\|\mathbf{x}\|_2^2 - f(\mathbf{x})$ mit $\nabla g(\mathbf{x}) = L\mathbf{x} - \nabla f(\mathbf{x})$ zeigt g ist konvex. Substitution von $g(\mathbf{x})$ in

$$\forall \mathbf{a}, \mathbf{b} \in \mathbb{R}^n : g(\mathbf{a}) \geq g(\mathbf{b}) + \langle \nabla g(\mathbf{b}), \mathbf{a} - \mathbf{b} \rangle$$

führt zu

$$f(\mathbf{a}) \leq f(\mathbf{b}) + \langle \nabla f(\mathbf{b}), \mathbf{a} - \mathbf{b} \rangle + (L/2)\|\mathbf{a} - \mathbf{b}\|_2^2 \quad (1)$$

Konvexität + Lipschitz stetige Gradienten

$$\forall \mathbf{a}, \mathbf{b} \in \mathbb{R}^n : \|\nabla f(\mathbf{a}) - \nabla f(\mathbf{b})\|_2 \leq L\|\mathbf{a} - \mathbf{b}\|_2$$

Multiplikation mit $\|\mathbf{a} - \mathbf{b}\|_2$ und Cauchy-Schwarz Ungl.:

$$\forall \mathbf{a}, \mathbf{b} \in \mathbb{R}^n : \langle \nabla f(\mathbf{a}) - \nabla f(\mathbf{b}), \mathbf{a} - \mathbf{b} \rangle \leq L\|\mathbf{a} - \mathbf{b}\|_2^2$$

Substitution von $g(\mathbf{x}) = (L/2)\|\mathbf{x}\|_2^2 - f(\mathbf{x})$ mit $\nabla g(\mathbf{x}) = L\mathbf{x} - \nabla f(\mathbf{x})$ zeigt g ist konvex. Substitution von $g(\mathbf{x})$ in

$$\forall \mathbf{a}, \mathbf{b} \in \mathbb{R}^n : g(\mathbf{a}) \geq g(\mathbf{b}) + \langle \nabla g(\mathbf{b}), \mathbf{a} - \mathbf{b} \rangle$$

führt zu

$$f(\mathbf{a}) \leq f(\mathbf{b}) + \langle \nabla f(\mathbf{b}), \mathbf{a} - \mathbf{b} \rangle + (L/2)\|\mathbf{a} - \mathbf{b}\|_2^2 \quad (1)$$

Konvexität + Lipschitz stetige Gradienten

$$\forall \mathbf{a}, \mathbf{b} \in \mathbb{R}^n : \|\nabla f(\mathbf{a}) - \nabla f(\mathbf{b})\|_2 \leq L\|\mathbf{a} - \mathbf{b}\|_2$$

Multiplikation mit $\|\mathbf{a} - \mathbf{b}\|_2$ und Cauchy-Schwarz Ungl.:

$$\forall \mathbf{a}, \mathbf{b} \in \mathbb{R}^n : \langle \nabla f(\mathbf{a}) - \nabla f(\mathbf{b}), \mathbf{a} - \mathbf{b} \rangle \leq L\|\mathbf{a} - \mathbf{b}\|_2^2$$

Substitution von $g(\mathbf{x}) = (L/2)\|\mathbf{x}\|_2^2 - f(\mathbf{x})$ mit $\nabla g(\mathbf{x}) = L\mathbf{x} - \nabla f(\mathbf{x})$ zeigt g ist konvex. Substitution von $g(\mathbf{x})$ in

$$\forall \mathbf{a}, \mathbf{b} \in \mathbb{R}^n : g(\mathbf{a}) \geq g(\mathbf{b}) + \langle \nabla g(\mathbf{b}), \mathbf{a} - \mathbf{b} \rangle$$

führt zu

$$f(\mathbf{a}) \leq f(\mathbf{b}) + \langle \nabla f(\mathbf{b}), \mathbf{a} - \mathbf{b} \rangle + (L/2)\|\mathbf{a} - \mathbf{b}\|_2^2 \quad (1)$$

Konvexität + Lipschitz stetige Gradienten

$$\forall \mathbf{a}, \mathbf{b} \in \mathbb{R}^n : \|\nabla f(\mathbf{a}) - \nabla f(\mathbf{b})\|_2 \leq L\|\mathbf{a} - \mathbf{b}\|_2$$

Multiplikation mit $\|\mathbf{a} - \mathbf{b}\|_2$ und Cauchy-Schwarz Ungl.:

$$\forall \mathbf{a}, \mathbf{b} \in \mathbb{R}^n : \langle \nabla f(\mathbf{a}) - \nabla f(\mathbf{b}), \mathbf{a} - \mathbf{b} \rangle \leq L\|\mathbf{a} - \mathbf{b}\|_2^2$$

Substitution von $g(\mathbf{x}) = (L/2)\|\mathbf{x}\|_2^2 - f(\mathbf{x})$ mit $\nabla g(\mathbf{x}) = L\mathbf{x} - \nabla f(\mathbf{x})$ zeigt g ist konvex. Substitution von $g(\mathbf{x})$ in

$$\forall \mathbf{a}, \mathbf{b} \in \mathbb{R}^n : g(\mathbf{a}) \geq g(\mathbf{b}) + \langle \nabla g(\mathbf{b}), \mathbf{a} - \mathbf{b} \rangle$$

führt zu

$$f(\mathbf{a}) \leq f(\mathbf{b}) + \langle \nabla f(\mathbf{b}), \mathbf{a} - \mathbf{b} \rangle + (L/2)\|\mathbf{a} - \mathbf{b}\|_2^2 \quad (1)$$



Methoden

- **Spezialisiert für bestimmte Modellklassen**
 - Stützvektormethode und Sequential Minimal Optimization
 - Markov Random Fields und Iterative Proportional Fitting
- **Generisch**
 - Erster Ordnung, mit Gradient
 - Einfache Implementierung
 - Geringer Ressourcenverbrauch
 - Vergleichsweise langsame Konvergenz
 - Zweiter Ordnung, mit Hesse-Matrix
 - Hoher Ressourcenverbrauch
 - Schnelle Konvergenz
 - Proximal-Point Methoden Erster Ordnung
 - Unterstützung für Nebenbedingungen und nicht-überall differenzierbare Funktionen
 - Heuristiken, mit Glück (Evolutionäre Algorithmen, Randomisierte Suche)



Methoden

- **Spezialisiert für bestimmte Modellklassen**
 - Stützvektormethode und Sequential Minimal Optimization
 - Markov Random Fields und Iterative Proportional Fitting
- **Generisch**
 - **Erster Ordnung, mit Gradient**
 - Einfache Implementierung
 - Geringer Ressourcenverbrauch
 - Vergleichsweise langsame Konvergenz
 - **Zweiter Ordnung, mit Hesse-Matrix**
 - Hoher Ressourcenverbrauch
 - Schnelle Konvergenz
 - **Proximal-Point Methoden Erster Ordnung**
 - Unterstützung für Nebenbedingungen und nicht-überall differenzierbare Funktionen
 - **Heuristiken, mit Glück (Evolutionäre Algorithmen, Randomisierte Suche)**



Methoden

- **Spezialisiert für bestimmte Modellklassen**
 - Stützvektormethode und Sequential Minimal Optimization
 - Markov Random Fields und Iterative Proportional Fitting
- **Generisch**
 - **Erster Ordnung, mit Gradient**
 - Einfache Implementierung
 - Geringer Ressourcenverbrauch
 - Vergleichsweise langsame Konvergenz
 - **Zweiter Ordnung, mit Hesse-Matrix**
 - Hoher Ressourcenverbrauch
 - Schnelle Konvergenz
 - **Proximal-Point Methoden Erster Ordnung**
 - Unterstützung für Nebenbedingungen und nicht-überall differenzierbare Funktionen
 - **Heuristiken, mit Glück (Evolutionäre Algorithmen, Randomisierte Suche)**



Methoden

- **Spezialisiert für bestimmte Modellklassen**
 - Stützvektormethode und Sequential Minimal Optimization
 - Markov Random Fields und Iterative Proportional Fitting
- **Generisch**
 - **Erster Ordnung, mit Gradient**
 - Einfache Implementierung
 - Geringer Ressourcenverbrauch
 - Vergleichsweise langsame Konvergenz
 - **Zweiter Ordnung, mit Hesse-Matrix**
 - Hoher Ressourcenverbrauch
 - Schnelle Konvergenz
 - **Proximal-Point Methoden Erster Ordnung**
 - Unterstützung für Nebenbedingungen und nicht-überall differenzierbare Funktionen
 - **Heuristiken, mit Glück (Evolutionäre Algorithmen, Randomisierte Suche)**



Methoden

- **Spezialisiert für bestimmte Modellklassen**
 - Stützvektormethode und Sequential Minimal Optimization
 - Markov Random Fields und Iterative Proportional Fitting
- **Generisch**
 - Erster Ordnung, mit Gradient
 - Einfache Implementierung
 - Geringer Ressourcenverbrauch
 - Vergleichsweise langsame Konvergenz
 - Zweiter Ordnung, mit Hesse-Matrix
 - Hoher Ressourcenverbrauch
 - Schnelle Konvergenz
 - Proximal-Point Methoden Erster Ordnung
 - Unterstützung für Nebenbedingungen und nicht-überall differenzierbare Funktionen
 - Heuristiken, mit Glück (Evolutionäre Algorithmen, Randomisierte Suche)

Gradientenabstieg

Jetzt: Modelle repräsentiert durch Parameter vektor $\beta \in \mathbb{R}^n$.
 Lernen von β mittels *Gradientenabstieg*:

- 1 Wähle beliebigen Startwert β^0 sowie $\eta_0, \eta_1, \eta_2, \dots$,
- 2 Erzeuge Sequenz von Modellen $\beta^0, \beta^1, \beta^2, \dots$ mittels

$$\beta^{t+1} = \beta^t - \eta_t \nabla \ell(\beta^t; \mathcal{D})$$

Sei $\ell(\beta; \mathcal{D})$ eine konvexe Funktion (in β), mit Lipschitz stetigem Gradienten (Konstante L) und $\eta_t = \eta = 1/L$. Dann gilt nach t Schritten: $\ell(\beta^t; \mathcal{D}) - \ell(\beta^*; \mathcal{D}) \leq \frac{L}{2t} \|\beta^0 - \beta^*\|_2^2$.



Gradientenabstieg

Jetzt: Modelle repräsentiert durch Parameter vektor $\beta \in \mathbb{R}^n$.
Lernen von β mittels *Gradientenabstieg*:

- 1 Wähle beliebigen Startwert β^0 sowie $\eta_0, \eta_1, \eta_2, \dots$,
- 2 Erzeuge Sequenz von Modellen $\beta^0, \beta^1, \beta^2, \dots$ mittels

$$\beta^{t+1} = \beta^t - \eta_t \nabla \ell(\beta^t; \mathcal{D})$$

Sei $\ell(\beta; \mathcal{D})$ eine konvexe Funktion (in β), mit Lipschitz stetigem Gradienten (Konstante L) und $\eta_t = \eta = 1/L$. Dann gilt nach t Schritten: $\ell(\beta^t; \mathcal{D}) - \ell(\beta^*; \mathcal{D}) \leq \frac{L}{2t} \|\beta^0 - \beta^*\|_2^2$.



Gradientenabstieg—Konvergenz

Sei $\ell(\beta; \mathcal{D})$ eine konvexe Funktion (in β), mit Lipschitz stetigem Gradienten (Konstante L) und $\eta_t = \eta = 1/L$. Dann gilt nach t Schritten,

$$\ell(\beta^t; \mathcal{D}) - \ell(\beta^*; \mathcal{D}) \leq \frac{L}{2t} \|\beta^0 - \beta^*\|_2^2.$$

Also erfordert $\ell(\beta^t; \mathcal{D}) - \ell(\beta^*; \mathcal{D}) \leq \epsilon$ höchstens $\mathcal{O}(L/\epsilon)$ Schritte.



Gradientenabstieg—Konvergenz (II)

Beweis. Substitution von $\mathbf{a} = \beta^{t+1}$ und $\mathbf{b} = \beta^t$ in (1),

$$\begin{aligned}\ell(\beta^{t+1}; \mathcal{D}) &\leq \ell(\beta^t; \mathcal{D}) - (\eta - (L\eta^2)/2) \|\nabla \ell(\beta^t; \mathcal{D})\|_2^2 \\ &\leq \ell(\beta^*; \mathcal{D}) + \langle \nabla \ell(\beta^t; \mathcal{D}), \beta^t - \beta^* \rangle - (\eta/2) \|\nabla \ell(\beta^t; \mathcal{D})\|_2^2 \\ &= \ell(\beta^*; \mathcal{D}) + (1/(2\eta)) (\|\beta^t - \beta^*\|_2^2 - \|\beta^{t+1} - \beta^*\|_2^2) \quad (2)\end{aligned}$$

Die Sequenz der Funktionswerte $\ell(\beta^t; \mathcal{D})$ ist monoton fallend (erste Umformung).

Damit ist $\ell(\beta^t; \mathcal{D}) - \ell(\beta^*; \mathcal{D})$ kleiner als die übrigen Folgeglieder $\ell(\beta^j; \mathcal{D}) - \ell(\beta^*; \mathcal{D})$ mit $j < t$.



Gradientenabstieg—Konvergenz (II)

Beweis. Substitution von $\mathbf{a} = \beta^{t+1}$ und $\mathbf{b} = \beta^t$ in (1),

$$\begin{aligned}\ell(\beta^{t+1}; \mathcal{D}) &\leq \ell(\beta^t; \mathcal{D}) - (\eta - (L\eta^2)/2) \|\nabla \ell(\beta^t; \mathcal{D})\|_2^2 \\ &\leq \ell(\beta^*; \mathcal{D}) + \langle \nabla \ell(\beta^t; \mathcal{D}), \beta^t - \beta^* \rangle - (\eta/2) \|\nabla \ell(\beta^t; \mathcal{D})\|_2^2 \\ &= \ell(\beta^*; \mathcal{D}) + (1/(2\eta)) (\|\beta^t - \beta^*\|_2^2 - \|\beta^{t+1} - \beta^*\|_2^2) \quad (2)\end{aligned}$$

Die Sequenz der Funktionswerte $\ell(\beta^t; \mathcal{D})$ ist monoton fallend (erste Umformung).

Damit ist $\ell(\beta^t; \mathcal{D}) - \ell(\beta^*; \mathcal{D})$ kleiner als die übrigen Folgeglieder $\ell(\beta^j; \mathcal{D}) - \ell(\beta^*; \mathcal{D})$ mit $j < t$.



Gradientenabstieg—Konvergenz (III)

Das Minimum $\ell(\beta^t; \mathcal{D}) - \ell(\beta^*; \mathcal{D})$ ist kleiner als der Mittelwert der vorherigen Folgeglieder $\ell(\beta^j; \mathcal{D}) - \ell(\beta^*; \mathcal{D})$. Darum

$$\ell(\beta^t; \mathcal{D}) - \ell(\beta^*; \mathcal{D}) \leq \frac{1}{t} \sum_{j=1}^{t-1} \ell(\beta^j; \mathcal{D}) - \ell(\beta^*; \mathcal{D}) \leq \frac{L \|\beta^0 - \beta^*\|_2^2}{2t}.$$

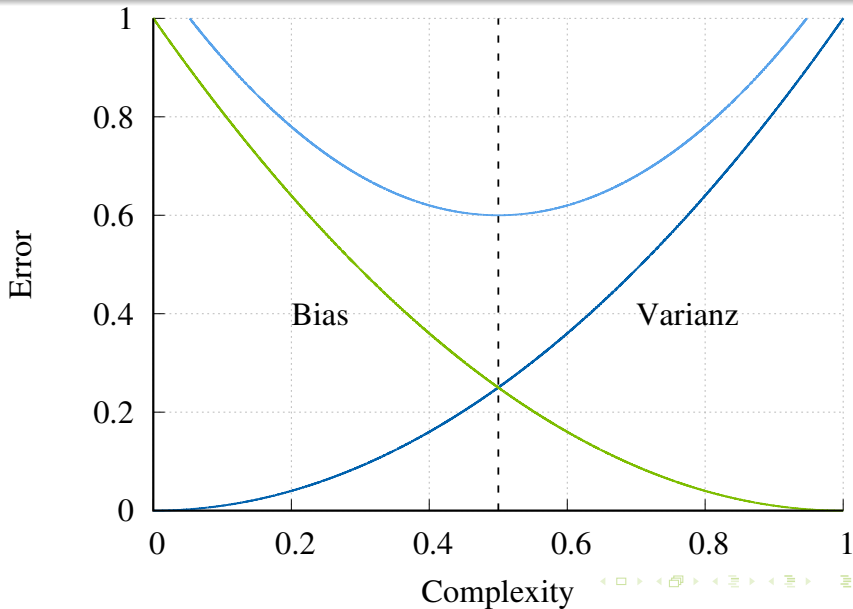
Die letzte Ungleichung folgt aus $\eta_t = \eta = 1/L$ und (2). ■



Überanpassung

Überanpassung

Bias und Varianz





Regularisierung

Intuition:

- Modell passt sich ggf. an das Rauschen ϵ in den Daten an
- Optimum der Verlustfunktion $\ell(f; \mathcal{D})$ liefert möglicher suboptimale Vorhersagen

Idee: bestrafe Überanpassung mittels *Regularisierung*

$$R: \mathcal{M} \rightarrow \mathbb{R}$$

$$f^* = \arg \min_{f \in \mathcal{M}} \ell(f; \mathcal{D}) + \lambda R(f)$$

Der Parameter $\lambda > 0$ bestimmt den Einfluss der Regularisierung.

Oft: $R(f) = \|f\|_q$