

Machine Learning

Support Vector Machines

Kernels for (Splice) Site Recognition

Soeren.Sonnenburg@first.fraunhofer.de



Fraunhofer Institut
Rechnerarchitektur
und Softwaretechnik

MASCHINELLES LERNEN ... ?

• Unüberwachtes Lernen

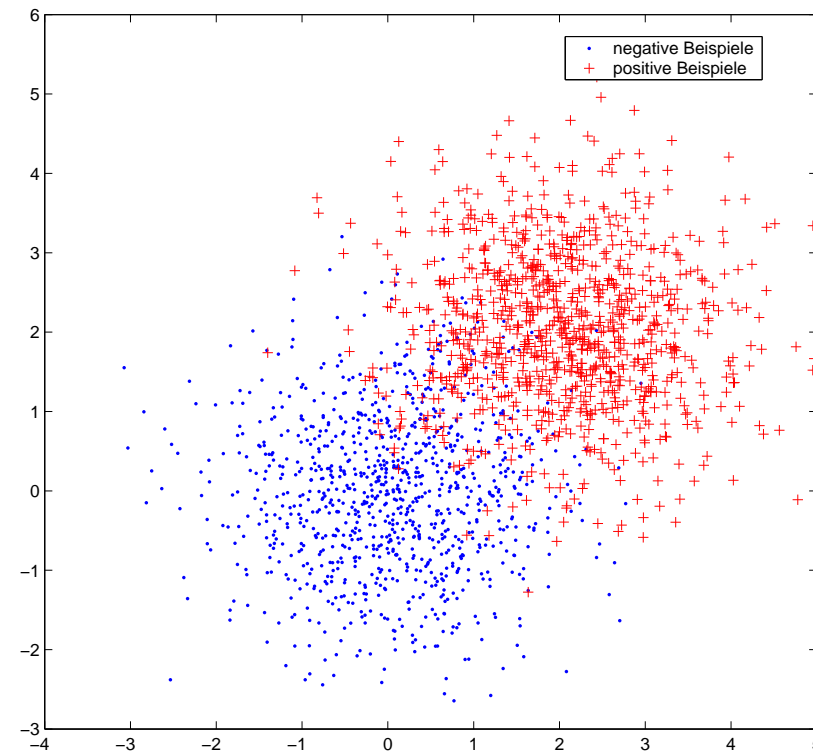
- Visualisierung
- Dimensionsreduktion
- Quellentrennung
- Clustering
- Merkmals Extraktion

• Überwachtes Lernen

- K-Nächster Nachbar
- Neuronale Netze, Perceptron
- Support Vector Machine
- Regression

KLASSIFIKATION ?

- Trainingsdaten $(\mathbf{x}_i, y_i) \in \mathcal{X} \times \{-1, +1\}$, $i = 1 \dots N$ (\mathcal{X} zB. \mathbb{R}^2)
- Funktion $f(\mathbf{x}, \alpha) : \mathcal{X} \times \{-1, +1\}$ mit $\mathbf{x} \mapsto y$



KLASSIFIKATION ?

- Trainingsdaten $(\mathbf{x}_i, y_i) \in \mathcal{X} \times \{-1, +1\}$, $i = 1 \dots N$ (\mathcal{X} zB. \mathbb{R}^2)
- Funktion $f(\mathbf{x}, \boldsymbol{\alpha}) : \mathcal{X} \times \{-1, +1\}$ mit $\mathbf{x} \mapsto y$

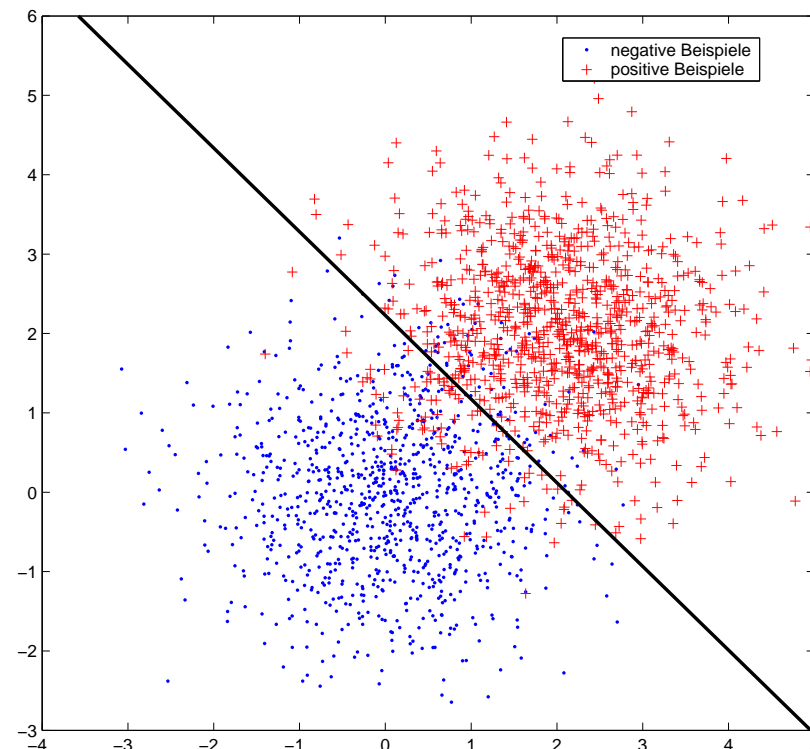
zB. Lineare Trennfunktion:

Parameter:

$$\boldsymbol{\alpha} = (\mathbf{w}, b)$$

Klassifikator:

$$f(\mathbf{x}) = \text{sign}(\mathbf{w} \cdot \mathbf{x} + b)$$



KLASSIFIKATOR

Klassifikator

- Trainingsdaten $(\mathbf{x}_i, y_i) \in \mathcal{X} \times \{-1, +1\}$, $i = 1 \dots N$ unabhängig voneinander und gleichmäßig von fester Wahrscheinlichkeitsverteilung aufgenommen (i.i.d.)
- Klassifikator: Funktion $f(\mathbf{x}, \boldsymbol{\alpha}) : \mathcal{X} \times \{-1, +1\}$ mit $\mathbf{x} \mapsto y$
- Vorhersage auf **ungesehenen** Daten (von gleicher Verteilung wie Trainingsdaten)

KLASSIFIKATOR ALA FISHER (1936)

Annahme 2 normalverteilte Klassen,

$$i \in \{1, 2\} : \mathcal{N}_{\mathbf{m}_i, \Sigma_i}(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} |\Sigma_i|} e^{\frac{1}{2}(\mathbf{x} - \mathbf{m}_i)^\top \Sigma_i^{-1} (\mathbf{x} - \mathbf{m}_i)}$$

Klassifikator:

$$\begin{aligned} f(\mathbf{x}) &= \text{sign}(\mathcal{N}_{\mathbf{m}_1, \Sigma_1}(\mathbf{x}) - \mathcal{N}_{\mathbf{m}_2, \Sigma_2}(\mathbf{x})) \\ &= \text{sign}\left(\ln \frac{\mathcal{N}_{\mathbf{m}_1, \Sigma_1}(\mathbf{x})}{\mathcal{N}_{\mathbf{m}_2, \Sigma_2}(\mathbf{x})}\right) \\ &= \text{sign}\left(\frac{1}{2}(\mathbf{x} - \mathbf{m}_1)^\top \Sigma_1^{-1} (\mathbf{x} - \mathbf{m}_1) - \frac{1}{2}(\mathbf{x} - \mathbf{m}_2)^\top \Sigma_2^{-1} (\mathbf{x} - \mathbf{m}_2) + \ln \frac{|\Sigma_2|}{|\Sigma_1|}\right) \end{aligned}$$

KLASSIFIKATOR ALA FISHER (1936)

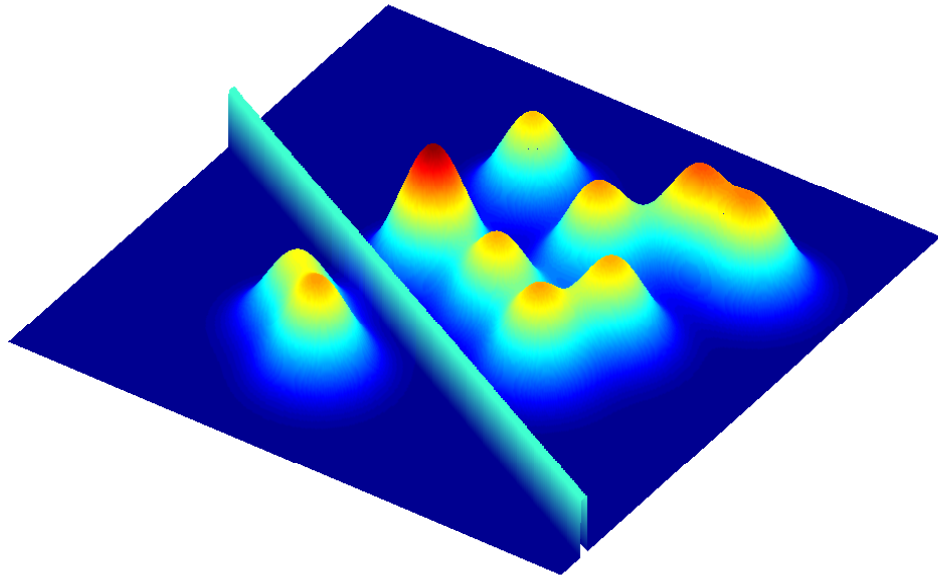
$$f_{sq}(\mathbf{x}) = \text{sign} \left(\frac{1}{2}(\mathbf{x} - \mathbf{m}_1)^\top \Sigma_1^{-1}(\mathbf{x} - \mathbf{m}_1) - \frac{1}{2}(\mathbf{x} - \mathbf{m}_2)^\top \Sigma_2^{-1}(\mathbf{x} - \mathbf{m}_2) + \ln \frac{|\Sigma_2|}{|\Sigma_1|} \right)$$

$$f_{lin}(\mathbf{x}) = \text{sign} \left((\mathbf{m}_2 - \mathbf{m}_1)^\top \Sigma^{-1} \mathbf{x} - \frac{1}{2} (\mathbf{m}_2^\top \Sigma^{-1} \mathbf{m}_2 + \mathbf{m}_1^\top \Sigma^{-1} \mathbf{m}_1) \right)$$

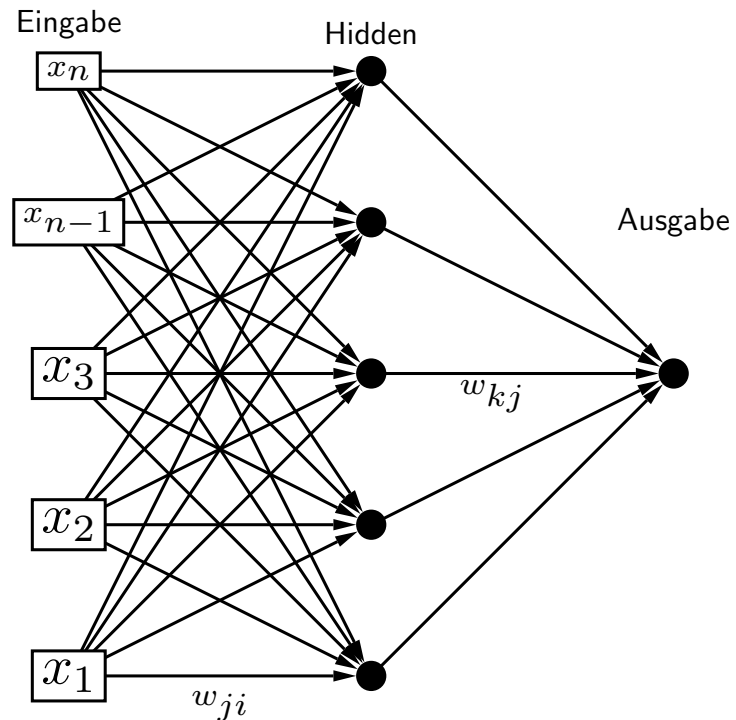
- $O(n^2)$ freie Parameter im quadratischen Fall (n Anzahl der Dimensionen)
- $O(n)$ Parameter im linearen Fall
- Fisher schlägt vor: $\Sigma = \tau \Sigma_1 + (1 - \tau) \Sigma_2$ wenn $N < 10 * n^2$
- SVM nur linear

PROBLEME

- Normalerweise Klassenverteilungen unbekannt
- „Curse of Dimensionality“
- Nur Trennfläche für Klassifikation interessant

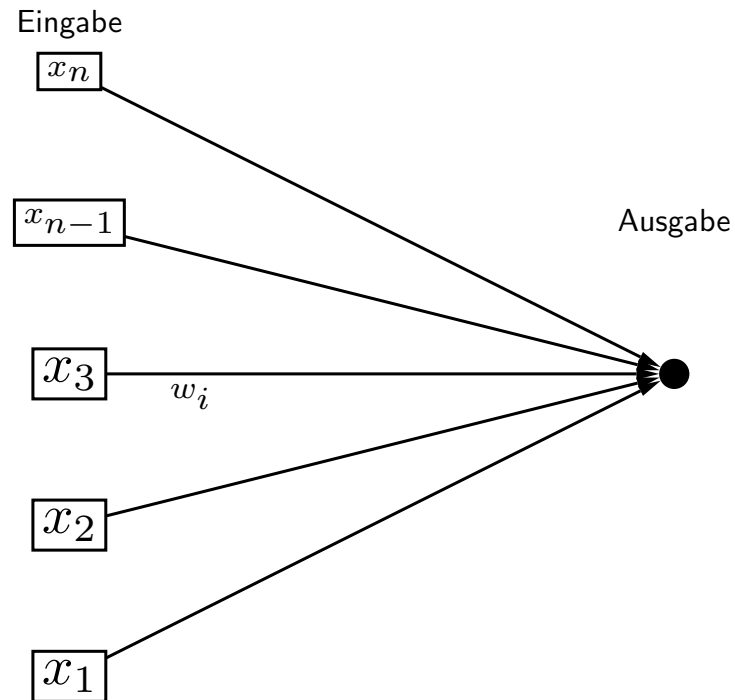


NEURONALES NETZ



$$f(\mathbf{x}) = \Phi \left(\sum_{j=1}^{n_H} w_{kj} \Phi \left(\sum_{i=1}^N w_{ji} x_i + w_{j0} \right) + w_{k0} \right)$$

PERCEPTRON



Linearer Klassifikator !

$$f(\mathbf{x}) = \sum_{i=1}^N w_i x_i + w_0 = \mathbf{w} \cdot \mathbf{x} + b$$

PROBLEME DES LERNENS

- Dichteschätzungen → „Curse of Dimensionality“
- Nur Trennfläche für Klassifikation interessant
- Lokale Minima vermeiden

Problem: Generalisierung vs. auswendig lernen

Beispiel Botaniker

FEHLER

- Empirischer Fehler

$$R_{emp}(\boldsymbol{\alpha}) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(y_i, f(\mathbf{x}_i, \boldsymbol{\alpha}))$$

$$\mathcal{L}(y, y') = \begin{cases} 0, & \text{wenn } y = y' \\ 1, & \text{sonst} \end{cases}$$

- Erwartungswert des wahren Fehlers bei bekannter Verteilung $P(\mathbf{x}, y)$

$$R(\boldsymbol{\alpha}) = \int \mathcal{L}(y, f(\mathbf{x}, \boldsymbol{\alpha})) dP(\mathbf{x}, y)$$

- Bestmöglicher Klassifikator minimiert $R(\boldsymbol{\alpha})$

BAYES KLASSIFIKATOR

$$\begin{aligned}R(\boldsymbol{\alpha}) &= \int \mathcal{L}(y, f(\mathbf{x}, \boldsymbol{\alpha})) dP(\mathbf{x}, y) \\ &= \int \mathcal{L}(-1, f(\mathbf{x}, \boldsymbol{\alpha})) P(y = -1, \mathbf{x}) d\mathbf{x} \\ &\quad + \int \mathcal{L}(+1, f(\mathbf{x}, \boldsymbol{\alpha})) P(y = +1, \mathbf{x}) d\mathbf{x}\end{aligned}$$

minimal für

$$\mathcal{L}(-1, f(\mathbf{x}, \boldsymbol{\alpha})) = \begin{cases} 0, & P(y = -1, \mathbf{x}) \geq P(y = +1, \mathbf{x}) \\ 1, & P(y = -1, \mathbf{x}) < P(y = +1, \mathbf{x}) \end{cases}$$
$$\mathcal{L}(+1, f(\mathbf{x}, \boldsymbol{\alpha})) = \begin{cases} 0, & P(y = +1, \mathbf{x}) \geq P(y = -1, \mathbf{x}) \\ 1, & P(y = +1, \mathbf{x}) < P(y = -1, \mathbf{x}) \end{cases}$$

Bestmöglicher Klassifikator: Bayes Klassifikator

$$\begin{aligned}f(\mathbf{x}) &= \text{sign}(P(y = +1, \mathbf{x}) - P(y = -1, \mathbf{x})) \\ &= \text{sign}(P(y = +1|\mathbf{x})P(\mathbf{x}) - P(y = -1|\mathbf{x})P(\mathbf{x})) \\ &= \text{sign}(P(y = +1|\mathbf{x}) - P(y = -1|\mathbf{x}))\end{aligned}$$

OBERE SCHRANKE FÜR FEHLER

- Abschätzung des Erwartungswertes des Fehlers $0 \leq \eta \leq 1$ beliebig, mit Wahrscheinlichkeit $1 - \eta$ hält diese Ungleichung (Vapnik,1995)

$$R(\alpha) \leq R_{emp}(\alpha) + \underbrace{\sqrt{\left(\frac{VC(\log \frac{2N}{VC} + 1) - \log \frac{\eta}{4}}{N} \right)}}_{\text{VC-Sicherheit}}$$

VAPNIK CHERVONENKIS (VC) DIMENSION

Definition

- VC ist die Vapnik Chervonenkis Dimension, ein Maß für die Komplexität einer Klasse von Funktionen $\{f(\mathbf{x}, \boldsymbol{\alpha})\}$
- Maximale Anzahl von Punkten beliebiger Färbung, die günstig angeordnet durch mindestens eine Funktion aus dieser Klasse auseinandergelassen werden kann.
Hier: Funktionen $\{f(\mathbf{x}, \boldsymbol{\alpha})\} \in \{-1, 1\}$, daher 2^N verschiedene Färbungen bei N Punkten.
- Separierung durch Hyperebenen mit Ausrichtung; $\mathbb{R}^n \rightarrow n + 1$ Punkte, Beispiel \mathbb{R}^2
- Beispiel: k -Nächster Nachbarschafts-klassifikator für $k = 1$ hat $R_{emp} = 0$ und VC Dimension unendlich. Abschätzung gibt keine Information, dennoch Klassifikator gut.

STRUKTURELLE RISIKO MINIMIERUNG (SRM)

- VC-Sicherheit nur von Klasse von Funktionen abhängig, **aber** empirisches/tatsächliches Risiko von der genauen Funktion abhängig.
- Minimierung durch Unterteilung der Klasse von Funktionen in verschachtelte Unterklassen für die wiederum VC bestimmt/abgeschätzt wird.
- SRM sucht die Unterklasse, die kleinste obere Schranke zum tatsächlichen Risiko ist, indem in jeder Unterklasse ein Klassifikator mit möglichst geringem empirischen Risiko bestimmt wird und letztendlich die mit kleinster Summe aus empirischem Risiko und VC Sicherheit genommen wird.

Occams Razor, KiSS

SVM

Der separierbare Fall

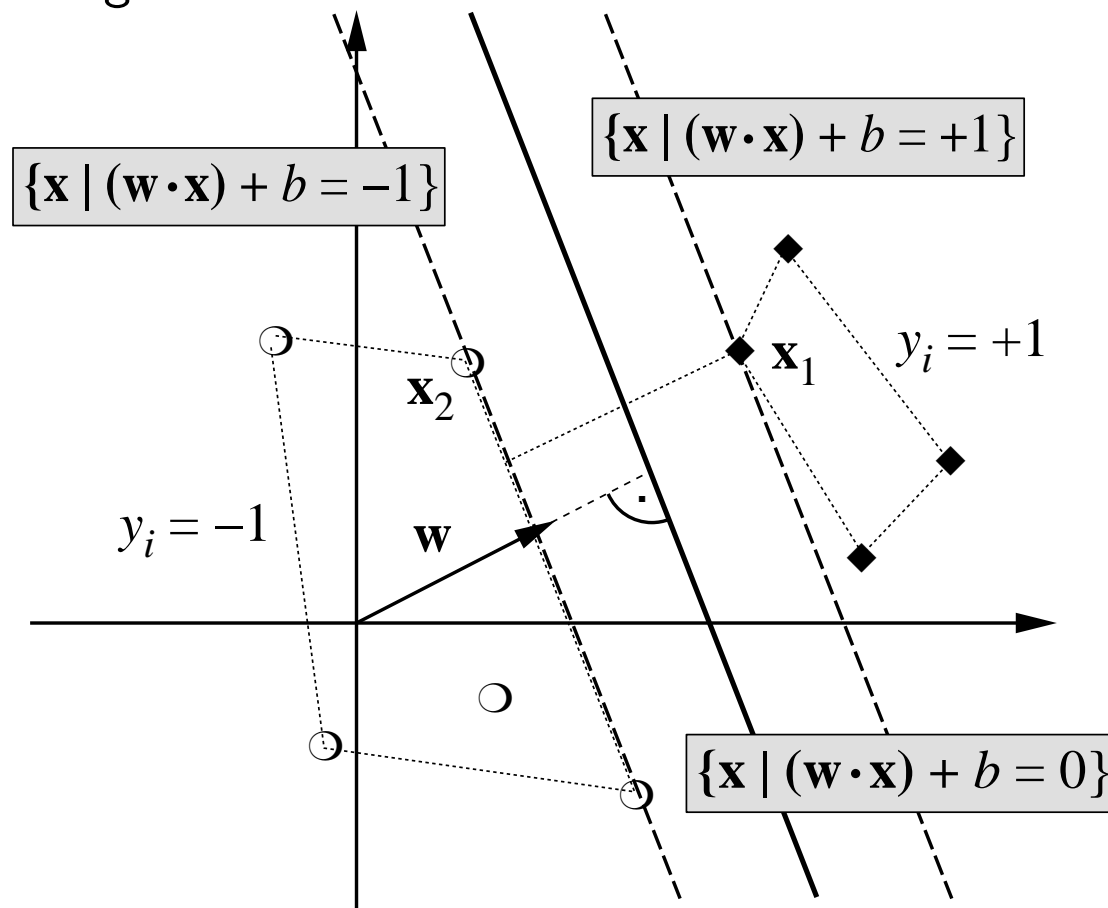
- Trainings Daten $(\mathbf{x}_i, y_i) \in \mathcal{X} \times \{-1, +1\}$ i.i.d.
- Support Vector Machine lernt *linearen* Klassifikator $f(\mathbf{x}, \boldsymbol{\alpha})$, $\boldsymbol{\alpha} = (\mathbf{w}, b)$.

$$f(\mathbf{x}) = \text{sign}(\mathbf{w} \cdot \mathbf{x} + b)$$

- welcher Lineare Klassifikator „optimal“ ?
- der mit größtem „Margin“ (störungsunempfindlich in Daten und Parametern)
- **globales Optimum**, Komplexität nur von Anzahl der Supportvektoren nicht der Dimension des Problems abhängig

SVM

Hyperebene beschrieben durch $w \cdot x + b = 0$ separiert unsere Trainingsdaten in positive/negative Exemplare mit $f(x) = \text{sign}(w \cdot x + b)$ als Entscheidungsfunktion. Problem Bestimmung der optimalen Hyperebene mit maximalem Margin:



Note:

$$(w \cdot x_1) + b = +1$$

$$(w \cdot x_2) + b = -1$$

$$\Rightarrow (w \cdot (x_1 - x_2)) = 2$$

$$\Rightarrow \left(\frac{w}{\|w\|} \cdot (x_1 - x_2) \right) = \frac{2}{\|w\|}$$

SVM

- Senkrechter Abstand der Hyperebenen: $\frac{2}{\|\mathbf{w}\|}$ **Margin**
- \Rightarrow maximal für minimales $\|\mathbf{w}\|$
- Bedingungen: Hyperebenen parrallel zur Entscheidungsfläche und \mathbf{x}_i auf der „richtigen“ Seite

$$\mathbf{x}_i \cdot \mathbf{w} + b \geq +1 \text{ for } y_i = +1$$

$$\mathbf{x}_i \cdot \mathbf{w} + b \leq -1 \text{ for } y_i = -1$$

$$\Rightarrow y_i(\mathbf{x}_i \cdot \mathbf{w} + b) \geq +1 \text{ for } y_i = \pm 1.$$

SVM

Zu lösen:

$$\text{minimiere } \frac{1}{2} \|\mathbf{w}\|^2$$

bei Einhaltung von $y_i(\mathbf{x}_i \cdot \mathbf{w} + b) - 1 \geq 0$ for $y_i = \pm 1, i = 1 \dots N$

⇒ *Standard* quadratisches Optimierungsproblem

- positive Lagrange Multiplikatoren: $\alpha_i, i = 1, \dots, N$. für jede Ungleichung einen. werden mit der jeweiligen Ungleichung multipliziert und von der gegenläufigen Bedingung subtrahiert.

$$L_p \equiv \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^N \alpha_i y_i (\mathbf{x}_i \cdot \mathbf{w} + b) + \sum_{i=1}^N \alpha_i$$

SVM

- \mathcal{C}_1 : Minimierung von L_p bzgl. w, b unter Berücksichtigung des Verschwindens aller Ableitungen nach α_i unter Einhaltung von $\alpha_i \geq 0$. Das ist ein konvexes quadratisches Optimierungsproblem daher äquivalente Lösung des „Dualen Problems.“
- Das Duale Optimierungsproblem
 \mathcal{C}_2 : Maximierung von L_p unter Berücksichtigung des Verschwindens der Ableitungen von L_p nach w, b und unter Einhaltung von $\alpha_i \geq 0$. „Wolf Dual“ mit Eigenschaft Maximum von L_p mit Bedingungen \mathcal{C}_2 hat selbe Lösungen w, b, α wie das Minimum unter den Bedingungen \mathcal{C}_1 .

SVM

$$L_p \equiv \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^N \alpha_i y_i (\mathbf{x}_i \cdot \mathbf{w} + b) + \sum_{i=1}^N \alpha_i$$

$$\frac{dL_p}{d\mathbf{w}} = 0 = \mathbf{w} - \sum_i \alpha_i y_i \mathbf{x}_i$$

$$\frac{dL_p}{db} = 0 = \sum_i \alpha_i y_i$$

einsetzen:

$$\begin{aligned} L_D &= \frac{1}{2} \left| \sum_i \alpha_i y_i \mathbf{x}_i \right|^2 - \sum_i \alpha_i y_i \mathbf{x}_i \cdot \sum_j \alpha_j y_j \mathbf{x}_j + \sum_i \alpha_i \\ &= \frac{1}{2} \sum_i \alpha_i y_i \mathbf{x}_i \sum_j \alpha_j y_j \mathbf{x}_j - \sum_i \alpha_i y_i \mathbf{x}_i \sum_j \alpha_j y_j \mathbf{x}_j + \sum_i \alpha_i \\ &= \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j \end{aligned}$$

SVM

The Karush-Kuhn-Tucker komplementär Bedingung sagt:

$$\alpha_i(y_i(\mathbf{x}_i \cdot \mathbf{w} + b) - 1) = 0 \quad i = 1, \dots, N$$

Die Lösung ist nun:

$$\mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i$$

$$b = \frac{1}{y_k} - \mathbf{x}_k \cdot \mathbf{w} = y_k - \mathbf{x}_k \cdot \left(\sum_{i=1}^N \alpha_i y_i \mathbf{x}_i \right)$$

Entscheidungsfunktion:

$$f(\mathbf{x}) = \text{sign} \left(\sum_{i=1}^N \alpha_i y_i \mathbf{x}_i \cdot \mathbf{x} + b \right) = \text{sign} \left(\left(\sum_{i=1}^N \alpha_i y_i \mathbf{x}_i \cdot (\mathbf{x} - \mathbf{x}_k) \right) + \frac{1}{y_k} \right).$$

DER NICHTSEPARIERBARE FALL

- Einführung von „Schlupfvariablen“ ξ_i

$$\mathbf{x}_i \cdot \mathbf{w} + b \geq +1 - \xi_i, \text{ for } y_i = +1$$

$$\mathbf{x}_i \cdot \mathbf{w} + b \leq -1 + \xi_i, \text{ for } y_i = -1$$

$$\xi_i \geq 0 \quad \forall i.$$

($\sum_i \xi_i$ obere Schranke für Trainingsfehler)

- Zielfunktion nun:

Minimiere $\frac{\|\mathbf{w}\|^2}{2} + C (\sum_i \xi_i)$, mit C Straffaktor (je größer desto mehr wird bestraft)

- „Wolf Dual“ Maximiere:

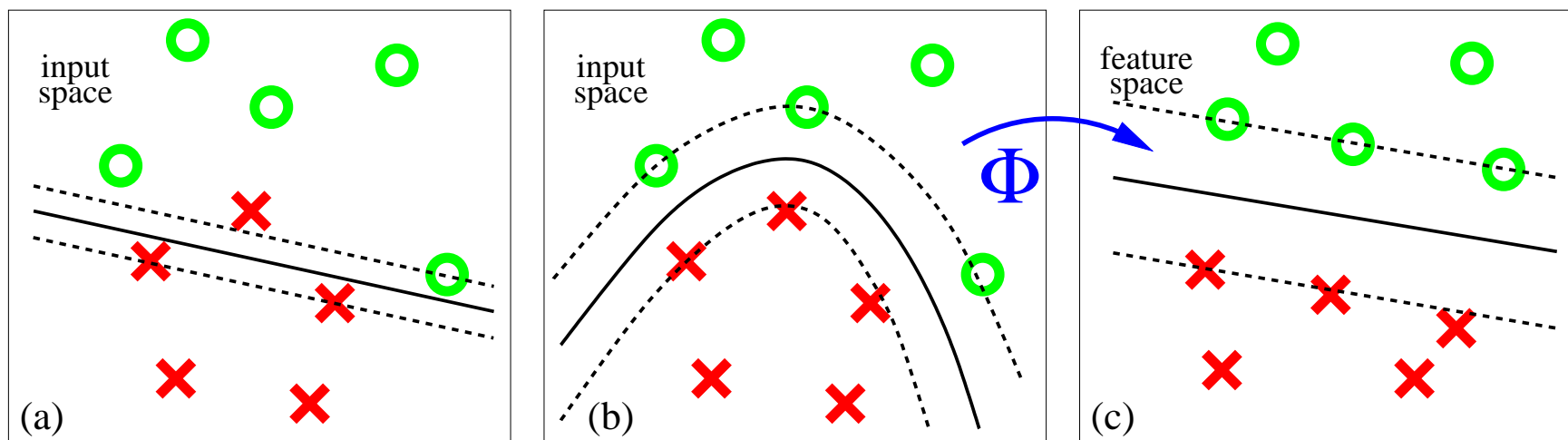
$$L_D = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j$$

unter Beachtung von $0 \leq \alpha_i \leq C$ und $\sum_i \alpha_i y_i = 0$.

- Lösung genauso (einziger Unterschied obere Schranke C)

SUPPORT VEKTOR KERNE

- bisher nur lineare Lösung in Merkmalsräumen, aber Lineare Trennung für viele Probleme nicht möglich
- → nichtlineare Transformation $\Phi : \mathcal{X} \rightarrow \mathcal{F}$, dort wieder linear Lösung suchen
- nur Skalarprodukte $\Phi(\mathbf{x}) \cdot \Phi(\mathbf{x}')$ benötigt und nicht explizit $\Phi(\mathbf{x})$
- Kernel Trick: $k(\mathbf{x}, \mathbf{x}') = (\Phi(\mathbf{x}) \cdot \Phi(\mathbf{x}'))$



SUPPORT VEKTOR KERNE

- $K(\mathbf{x}, \mathbf{x}') = (\Phi(\mathbf{x}) \cdot \Phi(\mathbf{x}'))$ in Duales Optimierungsproblem einsetzen:

$$L_D = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j)$$

- Entscheidungsfunktion wird:

$$\begin{aligned} f(\mathbf{x}) &= \text{sign} \left(\sum_{i=1}^N \alpha_i y_i (\Phi(\mathbf{x})_i \cdot \Phi(\mathbf{x})) + b \right) \\ &= \text{sign} \left(\sum_{i=1}^N \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) + b \right) \end{aligned}$$

KERNE - BEISPIELE

- Polynomieller Kernel

$$\begin{aligned}
 k(\mathbf{x}, \mathbf{x}') &:= (\mathbf{x} \cdot \mathbf{x}')^d \\
 &= \left(\sum_{i=1}^N x_i x'_i \right)^d \\
 &= \sum_{i_1, \dots, i_d \in \{1 \dots N\}} x_{i_1} \cdots x_{i_d} x'_{i_1} \cdots x'_{i_d} \\
 &= (C_d(\mathbf{x}) \cdot C_d(\mathbf{x}'))
 \end{aligned}$$

$C_d(\mathbf{x})$ ordnet \mathbf{x} alle bis zum Grad d möglichen Monome (nach Grad geordnet) zu Bsp: $C_2((x_1, x_2)) = (x_1^2, x_2^2, x_1x_2, x_2x_1)$

- Gauss RBF kernel:

$$k(\mathbf{x}, \mathbf{x}') := e^{-\left(\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{\sigma^2}\right)}$$

KLASSIFIKATOR (SVM) TRAINING

Standard Fall

- Aufteilung der Daten in Training, Validierung, Test
- Auswahl einer Menge geeigneter Kerne (oder neu designen)
- SVM auf Trainingsdatensatz trainieren und auf Validierungsdatensatz evaluieren. C und Kern (und Kernparameter) so einstellen, dass Validierungsfehler minimal.
- tatsächlichen Fehler auf den Testdaten schätzen

Wenig Daten \Rightarrow Kreuzvalidierung

- Aufteilung der Daten in s Teile



- $s - 1$ Teile für Training, den anderen Teil aufsplitten in Validierung und Test
- s mal trainieren (Extremfall $s = N - 1$)
- Auswahl einer Menge geeigneter Kerne (oder neu designen)
- SVM auf Trainingsdatensatz trainieren und auf Validierungsdatensatz evaluieren. C und Kern (und Kernparameter) so einstellen, dass Validierungsfehler minimal.
- tatsächlichen Fehler auf den Testdaten schätzen und mit Standardabweichung angeben

VERGLEICH VON KLASSIFIKATOREN

- Test Fehler
- Fehler für bestimmten Wert von falsch negativen auf dem Receiver Operating Characteristic (ROC) Graphen
- ROC Score (Fläche unter der ROC)

KONFUSIONS MATRIX

	Klassifikator Vorhersagen		
wahre Label	tp	fn	+
	fp	tn	-
	+	-	

- Sensitivität

$$\text{sensitivity} = \frac{tp}{tp + fn}$$

- Spezifität (Prezision)

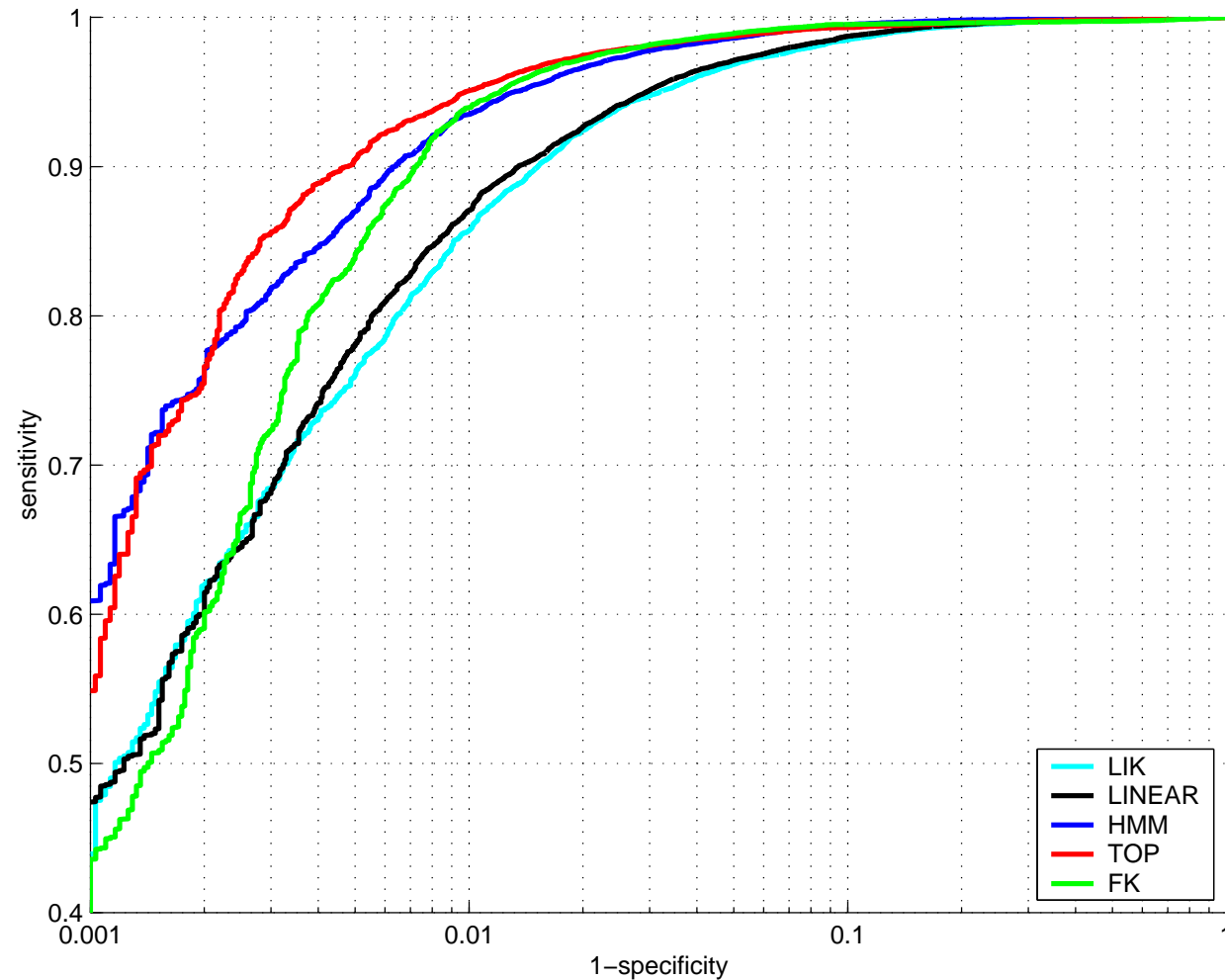
$$\text{specificity} = \frac{tn}{fp + tn}$$

ROC

Receiver Operating Characteristic

$$sens = \frac{tp}{tp + fn}$$

$$spec = \frac{tn}{fp + tn}$$



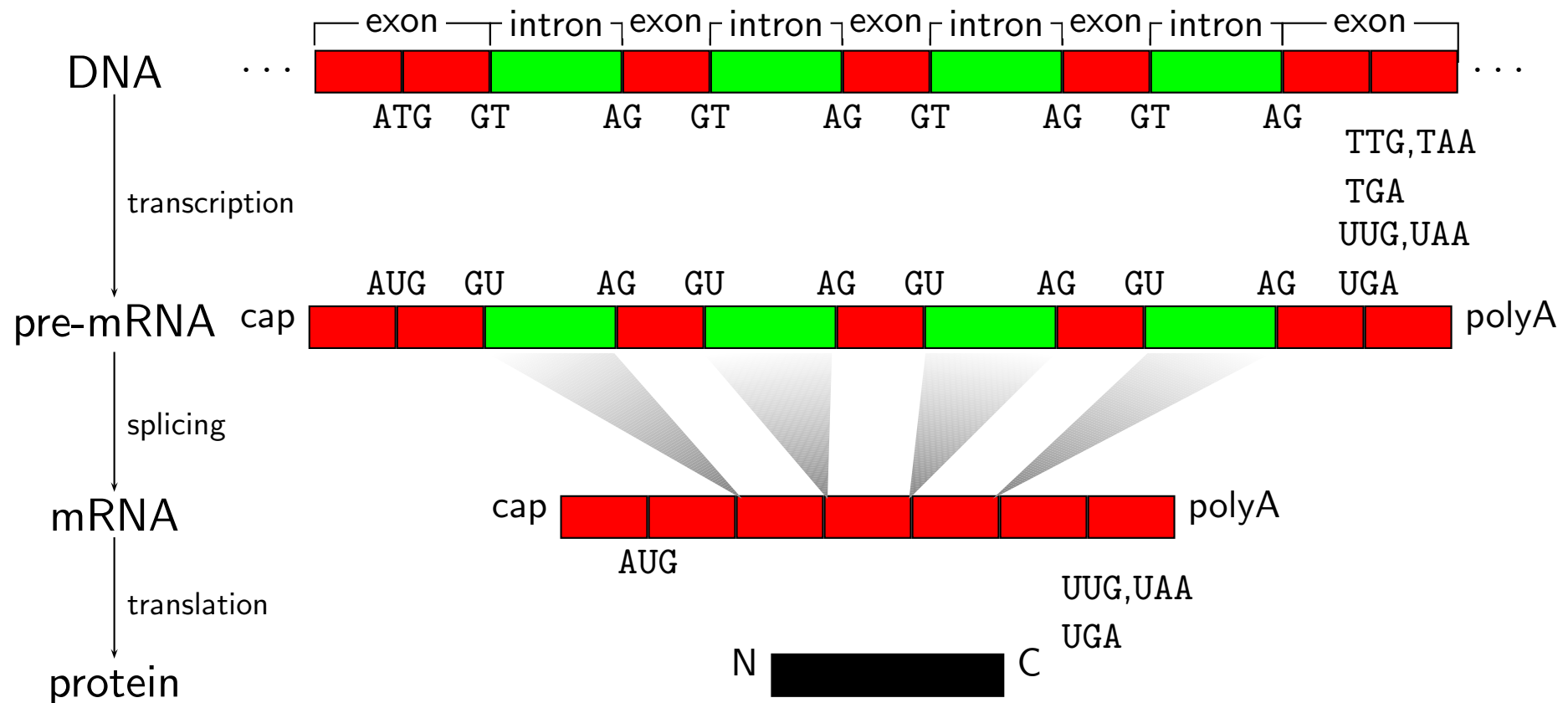
SPLEISSSTELLEN- ERKENNUNG

- SVM zur kanonischen Spleißstellenerkennung als Beispiel
- Vorwissen in Kern einbauen
- Was sind nützliche Eigenschaften
- Daten Extraktion
- Formulierung als 2 Klassen Klassifikationsproblem
- Welche Kerne machen Sinn ? Neudesign ?

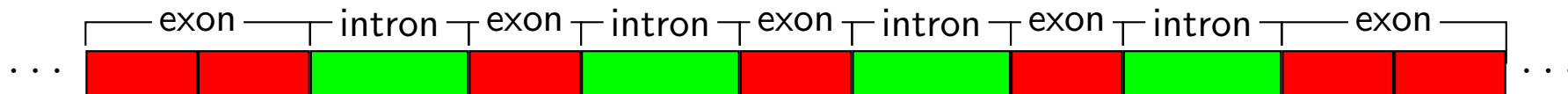
BIOLOGISCHER HINTERGRUND

Spleißstellen sind Übergangsstellen auf der DNA zwischen

- **exons** (Protein kodierend)
- **introns** (nicht kodierend)



FAKTEN ÜBER SPLEISSSTELLEN



- **Exons** sehr kurz (100 – 200 bp); **Introns** können sehr lang sein (> 1 kbp)
- Spleißen findet im *Zellkern* statt.
- Der Spleißapparat (“**Spliceosome**”) ist nicht Gewebespezifisch
- Spleißmechanismus sehr ähnlich in Eukaryoten
- Experimente zeigen, daß jede 5' Stelle mit jeder 3' Stelle verknüpft werden kann

Spleißmechanismus verwendet *locale* Information

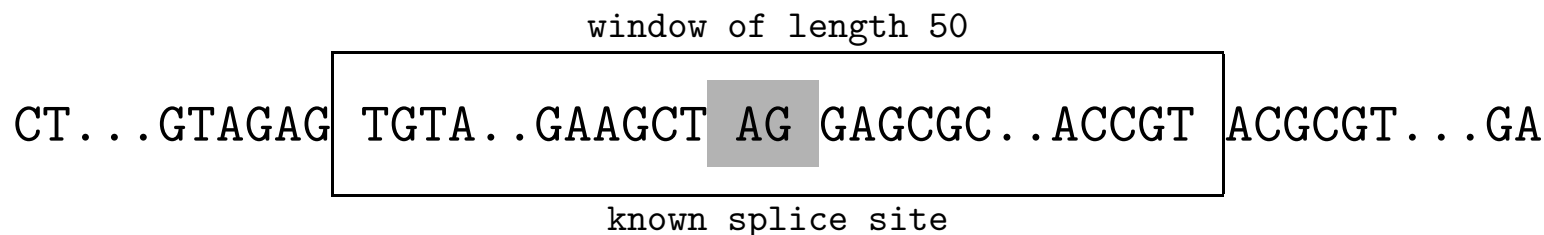
WIESO SPLEISSSTELLEN ERKENNUNG WICHTIG?

- erlaubt *akkurate* Vorhersage von mRNA and somit von Proteinen aus DNA
 - ⇒ wichtiger Schritt bei der Analyse des Genoms
- Spleißstellen können sehr gut erkannt werden
 - ⇒ wichtiger und exakter 'Marker' um Gene zu finden

Standard: Alignment to Datenbank Einträgen

Ziel: Spleißen mit Hilfe von ML verbessern

2-KLASSEN KLASSIFIKATIONS PROBLEM



- nur kanonische Spleißstellen (consensus AG,GT, 98%)
- **wahre Spleißstellen:** festes Fenster um die Spleißstelle
- **falsche:** erzeugt, durch verschieben des Fensters ± 25 bp

```

AAACAAATAAGTAACTAATCTTTTAGGAAGAACGTTTCAACCATTTTGAG
AAGATTAAAAAAAACAAATTTTTCATTACAGATATAATAATCTAATT
CACTCCCAAATCAACGATATTTTAGTTCACTAACACATCCGTCTGTGCC
TTAATTTCACTTCCACATACTTCCAGATCATCAATCTCCAAAACCAACAC
TTGTTTTAATATTCAATTTTTTACAGTAAGTTGCCAATTCAATGTTCCAC
TACCTAATTATGAAATTAATAATTCAGTGTGCTGATGGAAACGGAGAAGTC
  
```

Download at <http://mlg.anu.edu.au/~raetsch/splice>



NOCHMAL: EIGENSCHAFTEN VON SPLEISSSTELLEN

- Konservierte Sequenz
- Positionsabhängig
- Mutationen
- Einfügungen
- Löschungen
- Verschoben

SUPPORT VECTOR MACHINES

- Sequenzen $\mathbf{x}_i \in \mathcal{X}$ ($i = 1, \dots, N$) with respective labels y_i
- SVM Klassifikator (im Prinzip: Perceptron im Kernel Feature Space):

$$f(\mathbf{x}) = \text{sign} \left(\sum_{i=1}^N y_i \alpha_i \mathbf{k}(\mathbf{x}, \mathbf{x}_i) + b \right)$$

- lernt Parameter α durch Lösung des quadratischen Optimierungsproblems:

$$\max_{\alpha} \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{k}(\mathbf{x}_i, \mathbf{x}_j)$$

subject to $\alpha_i \in [0, C]$, $i = 1, \dots, N$, $\sum_{i=1}^N \alpha_i y_i = 0$.

Lösung hat keine lokalen Minima

KERNE - LINEAR/POLYNOMIELLER KERN

Polynomieller Kern vom Grad d :

$$k_{\text{POLY}}(\mathbf{x}, \mathbf{x}') = (\mathbf{x}, \mathbf{x}')^d$$

- nicht direkt anwendbar
- Ansatz, Eingabe Raum in binär Raum transformieren oder equivalent matches zählen

$$k_{\text{POLY}}(\mathbf{x}, \mathbf{x}') = \left(\sum_{p=1}^N l_p(\mathbf{x}, \mathbf{x}') \right)^d$$

⇒ Merkmalsraum alle Monome d -ter Ordnung, d.h.

$$x_{i_1} \cdot x_{i_2} \cdot \dots \cdot x_{i_d}, \quad i_{1\dots d} \in 1 \dots N$$

⇒ sehr hochdimensional, *nicht robust*

KERNELS - LOCALITY IMPROVED KERNEL

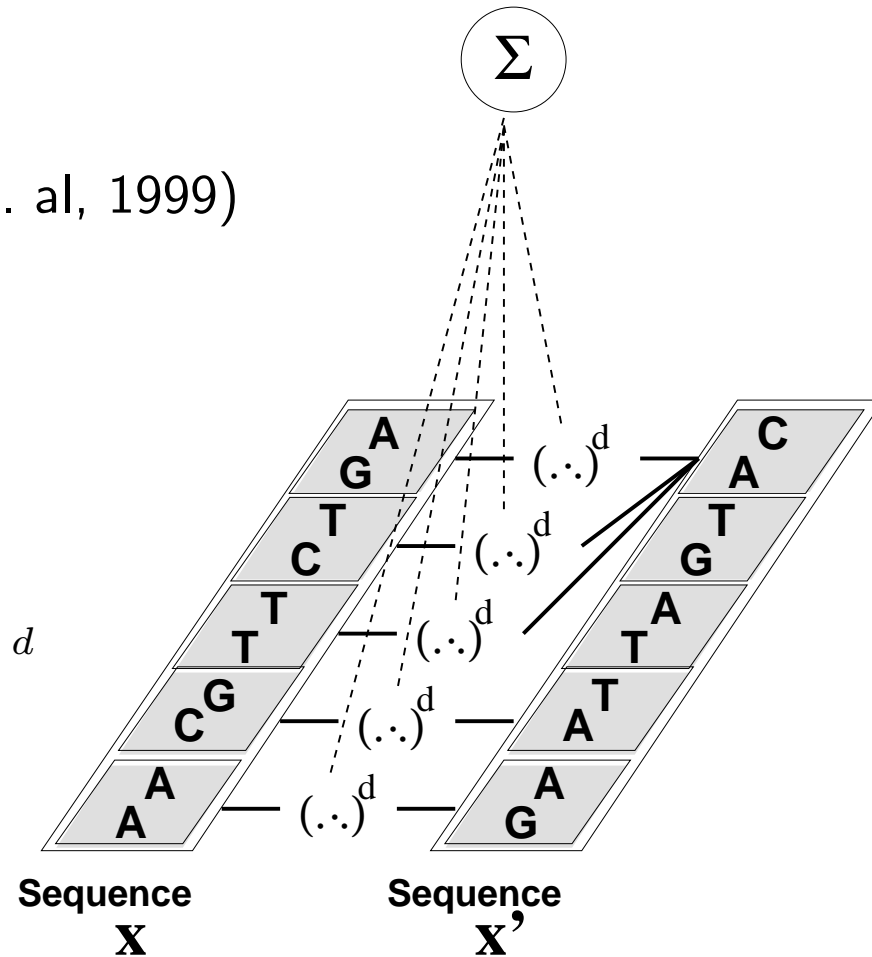
Locality Improved Kernel (Zien et. al, 1999)

- im Prinzip polynomieller Kern der locale Information betont

$$k_{LI}(\mathbf{x}, \mathbf{x}') = \sum_{p=1}^N \text{win}_p(\mathbf{x}, \mathbf{x}')$$

$$\text{win}_p(\mathbf{x}, \mathbf{x}') = \left(\sum_{j=-l}^{+l} p_j I_{p+j}(\mathbf{x}, \mathbf{x}') \right)^d$$

$$I_i(\mathbf{x}, \mathbf{x}') = \begin{cases} 1, & x_i = x'_i \\ 0, & \text{otherwise} \end{cases}$$



⇒ zählt letztendlich Matches von $2l + 1$ - Tupeln

WEIGHTED DEGREE KERNEL

- zählt Matches zwischen zwei Sequenzen x and x' zwischen den Worten $u_{\omega,i}(x)$ und $u_{\omega,i}(x')$
- $u_{\omega,i}(x) = x_i x_{i+1} \dots x_{i+\omega-1}$ für alle i und $1 \leq \omega \leq d$.
- ω ist die Ordnung (oder Länge des zu vergleichenden Wortes).

Weighted Degree Kernel

$$k_{WD}(\mathbf{x}, \mathbf{x}') = \sum_{\omega=1}^d w_{\omega} \sum_{i=1}^{N-d} \mathbb{I}(u_{\omega,i}(\mathbf{x}) = u_{\omega,i}(\mathbf{x}'))$$

- w_{ω} Gewicht für Match der Länge ω
- funktioniert auch mit Mismatches können (geringere Gewichte für nicht perfekte Matches)

Weighted Degree Position Kernel

$$k_{WDpos}(\mathbf{x}, \mathbf{x}') = \sum_{\omega=1}^d w_{\omega} \sum_{i=1}^{N-d-\text{shift}_i} \sum_{j=1}^{\text{shift}_i} \mathbb{I}(\mathbf{u}_{\omega, i+j}(\mathbf{x}) = \mathbf{u}_{\omega, i+j}(\mathbf{x}'))$$

Kern kann folgende Eigenschaften ausnutzen

- Konservierte Sequenz
- Positionsabhängig
- Mutationen
- Einfügungen
- Löschungen
- Verschieben

SPECTRUM KERNEL

Spectrum Kernel (Leslie et. al, 2002):

$$k_{spec} = s(\mathbf{x})^\top s(\mathbf{x})$$

Skalarprodukt aus

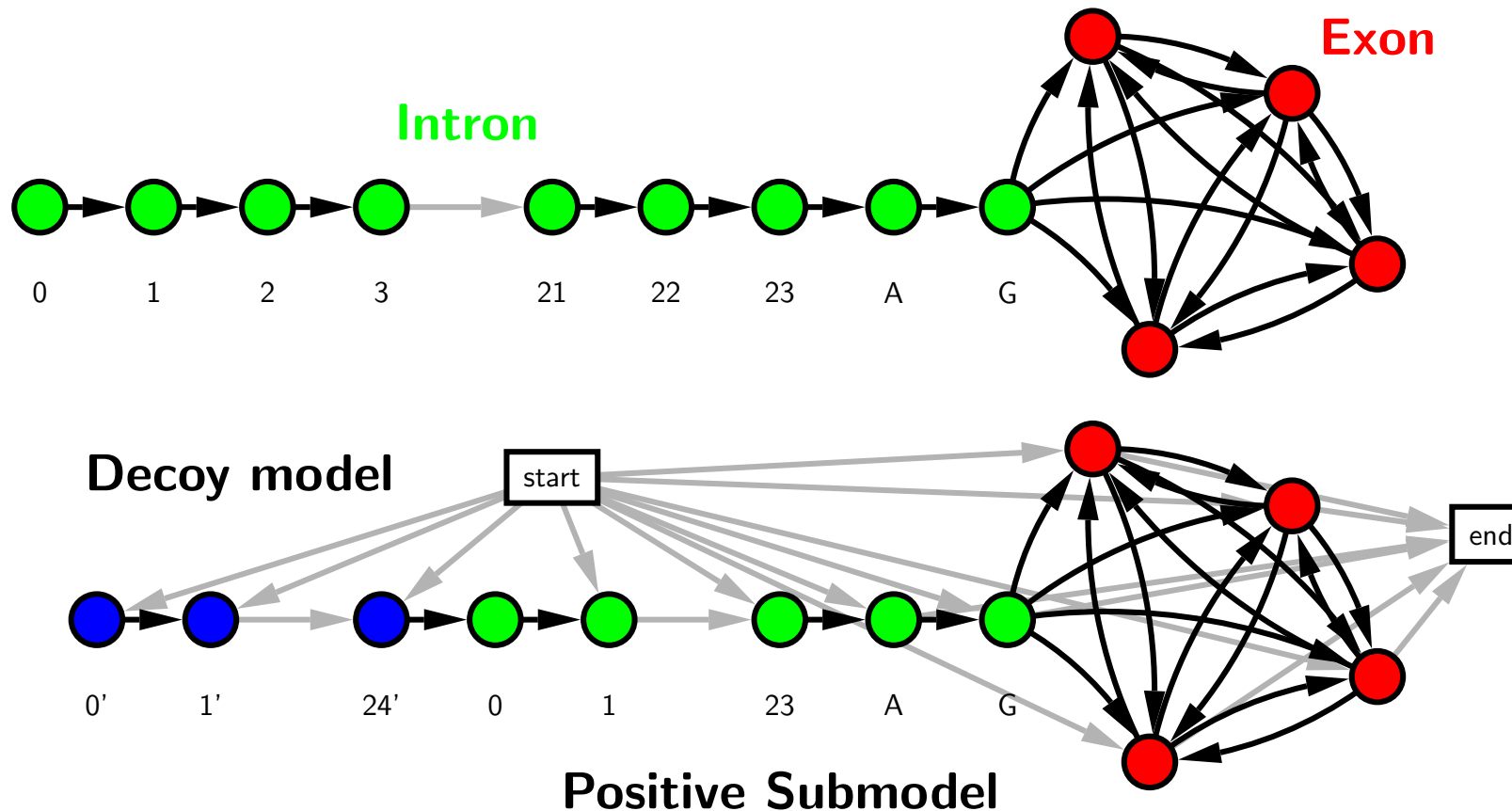
$$s(\mathbf{x}) = (\#A(\mathbf{x}))_{A \in \Sigma^k}$$

wobei $\#A(\mathbf{x}) =$ Häufigkeit des Auftretens von A in \mathbf{x} .

- zählt Matches von k -Tupeln
- keine positionsabhängige Information

GENERATIVE MODELLE

Verwenden generatives Model, z.B. designen ein HMM



(oben) positive Acceptor Model, (unten) negative Acceptor model

KERNE

- Kernel von generativen Modellen
 - vergleicht Objekte durch ein generative Model $\Pr(\mathbf{x}|\Theta)$
 - verwendet probabilistisches Model zum **discriminativen** Training
- Fisher Kernel (Jaakkola and Haussler, 1998)

$$\mathbf{k}_{\text{FK}}(\mathbf{x}, \mathbf{x}') = \mathbf{s}(\mathbf{x}, \hat{\boldsymbol{\theta}})^\top \mathbf{Z}^{-1}(\hat{\boldsymbol{\theta}}) \mathbf{s}(\mathbf{x}', \hat{\boldsymbol{\theta}})$$

$$\mathbf{s}(\mathbf{x}, \hat{\boldsymbol{\theta}}) = \nabla_{\boldsymbol{\theta}} \log \Pr(\mathbf{x}|\hat{\boldsymbol{\theta}})$$

Fisher score vector

$$\mathbf{Z}(\hat{\boldsymbol{\theta}}) = \mathbb{E}_{\mathbf{x}} \left[\mathbf{s}(\mathbf{x}, \hat{\boldsymbol{\theta}}) \mathbf{s}(\mathbf{x}, \hat{\boldsymbol{\theta}})^\top \mid \hat{\boldsymbol{\theta}} \right]$$

Fisher information matrix

- TOP Kernel (Tsuda et. al, 2002)

$$\mathbf{k}_{\text{TOP}}(\mathbf{x}, \mathbf{x}') = \mathbf{f}_{\hat{\boldsymbol{\theta}}}(\mathbf{x})^\top \mathbf{f}_{\hat{\boldsymbol{\theta}}}(\mathbf{x}')$$

$$\mathbf{f}_{\hat{\boldsymbol{\theta}}}(\mathbf{x}) = (v(\mathbf{x}, \hat{\boldsymbol{\theta}}), \nabla_{\boldsymbol{\theta}} v(\mathbf{x}, \hat{\boldsymbol{\theta}}))^\top$$

$$v(\mathbf{x}, \hat{\boldsymbol{\theta}}) = \log(\Pr(y = +1|\mathbf{x}, \hat{\boldsymbol{\theta}})) - \log(\Pr(y = -1|\mathbf{x}, \hat{\boldsymbol{\theta}}))$$

ERGEBNISSE

System	Neither	Donor	Acceptor	Total
LI-SVM	2.0±0.3%	0.8±0.2%	0.9±0.3%	3.7%
FK-SVM	2.1±0.4%	1.6±0.5%	1.6±0.4%	5.3%
TOP-SVM	2.2±0.4%	1.5±0.4%	1.7±0.3%	5.4%
HMM	2.6±0.5%	1.0±0.4%	2.4±0.7%	6.0%
RBF-SVM	n.d.	n.d.	n.d.	>10%
NN-BRAIN	n.d.	2.6%	4.3%	n.d.
BRAIN	4.0%	5.0%	4.0%	13.0%
KBANN	4.6%	7.6%	8.5%	20.7%
BackProp	5.3%	5.7%	10.7%	21.7%
PEBLS	6.9%	8.2%	7.6%	22.7%
Perceptron	4.0%	16.3%	17.4%	37.7%
ID3	8.8%	10.6%	14.0%	33.4%
COBWEB	11.8%	15.0%	9.5%	36.3%
Near. Neigh.	31.1%	11.7%	9.1%	51.9%

ERGEBNISSE

