

De-Novo-Sequencing using Spectrum-Graphs, enabling Open Searches

19. Juni 2023

Dominik Habermann

Ruhr Universität Bochum

Hintergrund

AS Sequenzierung

De-Novo-Sequenzierung

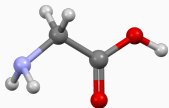
pNovo+ Algorithmus

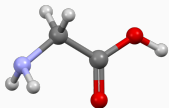
Open-pNovo Algorithmus

Zusammenfassung

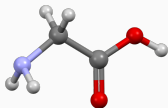
Hintergrund

- Peptide: Kette von Aminosäuren (AS)

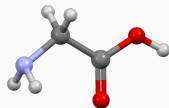




- Peptide: Kette von Aminosäuren (AS)
- 20 relevante AS



- Peptide: Kette von Aminosäuren (AS)
- 20 relevante AS
- Reihenfolge von AS ist weitestgehend beliebig

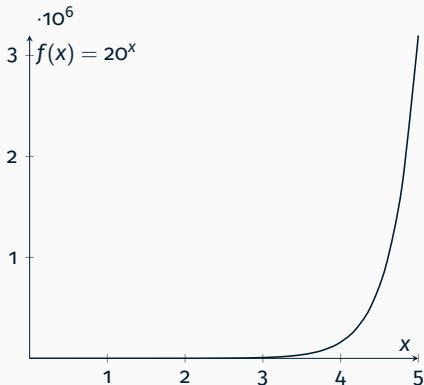


- Peptide: Kette von Aminosäuren (AS)
- 20 relevante AS
- Reihenfolge von AS ist weitestgehend beliebig
- $f(x) = 20^x$ x : Anzahl an AS

- Bereits bei wenigen AS: Kaum händelbare Anzahl an Kombinationen

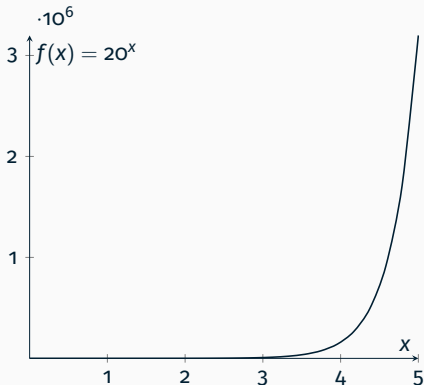
Anzahl an Kombinationen

- Bereits bei wenigen AS: Kaum händelbare Anzahl an Kombinationen



Anzahl an Kombinationen

- Bereits bei wenigen AS: Kaum händelbare Anzahl an Kombinationen



- Zum Vergleich: Proteine bis zu mehreren zehntausend AS



- Zuverlässige Bestimmung *kurzkettiger* Peptide möglich?



- Zuverlässige Bestimmung *kurzkettiger* Peptide möglich?
- Biomedizinisch relevant:



- Zuverlässige Bestimmung *kurzkettiger* Peptide möglich?
- Biomedizinisch relevant:
 - Katalogisierung von Proteinen



- Zuverlässige Bestimmung *kurzkettiger* Peptide möglich?
- Biomedizinisch relevant:
 - Katalogisierung von Proteinen
 - Wechselwirkungen von Proteinen



- Zuverlässige Bestimmung *kurzkettiger* Peptide möglich?
- Biomedizinisch relevant:
 - Katalogisierung von Proteinen
 - Wechselwirkungen von Proteinen
 - Analyse von Enzymen

AS Sequenzierung

- AS Sequenzierung: Bestimmung der AS-Sequenz

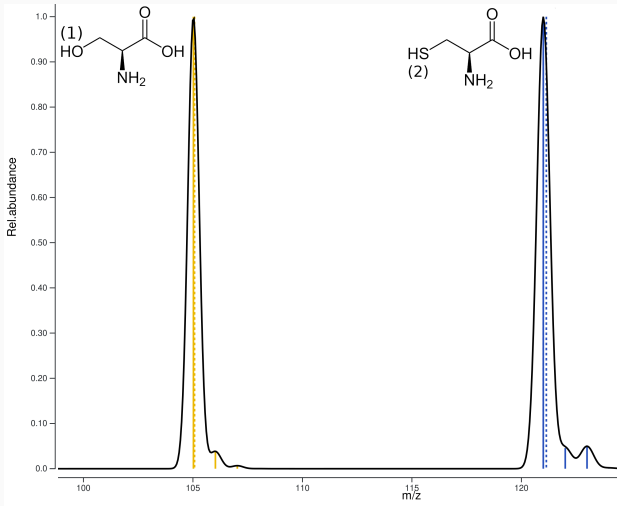
- AS Sequenzierung: Bestimmung der AS-Sequenz
- Hilfsmittel: Massenspektrometrie (MS)

- AS Sequenzierung: Bestimmung der AS-Sequenz
- Hilfsmittel: Massenspektrometrie (MS)
- MS kann chemische Strukturen bestimmen

- AS Sequenzierung: Bestimmung der AS-Sequenz
- Hilfsmittel: Massenspektrometrie (MS)
- MS kann chemische Strukturen bestimmen
- Rückschluss auf die AS-Sequenz

- Ergebnisse einer MS: Spektrogramm

■ Ergebnisse einer MS: Spektrogramm



De-Novo-Sequenzierung

- Beide Verfahren bestimmen die AS-Sequenz

- Beide Verfahren bestimmen die AS-Sequenz

| AS Sequenzierung | De-Novo-Sequenzierung |
|-----------------------------|-----------------------|
| Datenbanken als Hilfsmittel | |

- Beide Verfahren bestimmen die AS-Sequenz

| AS Sequenzierung | De-Novo-Sequenzierung |
|--|-----------------------|
| Datenbanken als Hilfsmittel | |
| Identifizierung von <i>bekannten</i> Sequenzen | |

- Beide Verfahren bestimmen die AS-Sequenz

| AS Sequenzierung | De-Novo-Sequenzierung |
|--|--------------------------|
| Datenbanken als Hilfsmittel | Ohne weitere Hilfsmittel |
| Identifizierung von <i>bekannten</i> Sequenzen | |

- Beide Verfahren bestimmen die AS-Sequenz

| AS Sequenzierung | De-Novo-Sequenzierung |
|--|---|
| Datenbanken als Hilfsmittel | Ohne weitere Hilfsmittel |
| Identifizierung von <i>bekannten</i> Sequenzen | Bestimmung <i>unbekannter</i> Sequenzen |

- Beide Verfahren bestimmen die AS-Sequenz

| AS Sequenzierung | De-Novo-Sequenzierung |
|--|---|
| Datenbanken als Hilfsmittel | Ohne weitere Hilfsmittel |
| Identifizierung von <i>bekannten</i> Sequenzen | Bestimmung <i>unbekannter</i> Sequenzen |

- De novo: lat. „Von neuem“

- Zusätzliche Informationen notwendig

- Zusätzliche Informationen notwendig
- Verwendung einer 2. MS

- Zusätzliche Informationen notwendig
- Verwendung einer 2. MS
- Verfahren: Tandem-Massenspektrometrie MS2

- Ionen aus m/z Bereich auswählbar machen

- Ionen aus m/z Bereich auswählbar machen
- Quasi eine Filterung

- Ionen aus m/z Bereich auswählbar machen
- Quasi eine Filterung
- Ausgewählte Ionen werden für 2. MS verwendet

- Ionen aus 1. MS „fragmentieren“:

- Ionen aus 1. MS „fragmentieren“:
 - Energiezuführung

- Ionen aus 1. MS „fragmentieren“:
 - Energiezuführung
 - Ionen zerfallen

- Ionen aus 1. MS „fragmentieren“:
 - Energiezuführung
 - Ionen zerfallen
 - Ergebnis: „Fragment-Ionen“

- Ionen aus 1. MS „fragmentieren“:
 - Energiezuführung
 - Ionen zerfallen
 - Ergebnis: „Fragment-Ionen“
- Verschiedene Fragmentierungsmethoden mit spezifischen Fragmenten

- Ionen aus 1. MS „fragmentieren“:
 - Energiezuführung
 - Ionen zerfallen
 - Ergebnis: „Fragment-Ionen“
- Verschiedene Fragmentierungsmethoden mit spezifischen Fragmenten
- 2. MS wird auf Fragment-Ionen angewendet

- Höhere Genauigkeit durch Filterung nach 1. MS

- Höhere Genauigkeit durch Filterung nach 1. MS
- Bessere Selektivität beim 2. MS

- Höhere Genauigkeit durch Filterung nach 1. MS
- Bessere Selektivität beim 2. MS
- Ionen zerfallen spezifisch

- Höhere Genauigkeit durch Filterung nach 1. MS
- Bessere Selektivität beim 2. MS
- Ionen zerfallen spezifisch
- → Rekonstruktion der ursprünglichen Ionen möglich

- Höhere Genauigkeit durch Filterung nach 1. MS
- Bessere Selektivität beim 2. MS
- Ionen zerfallen spezifisch
- → Rekonstruktion der ursprünglichen Ionen möglich
- MS2 Ergebnisse haben eine höhere Qualität

pNovo+ Algorithmus

- Algorithmus für die De-Novo-Sequenzierung

- Algorithmus für die De-Novo-Sequenzierung
- Auswertung von MS2 Spektrogrammen

- Algorithmus für die De-Novo-Sequenzierung
- Auswertung von MS2 Spektrogrammen
- Rekonstruktion der AS-Sequenz

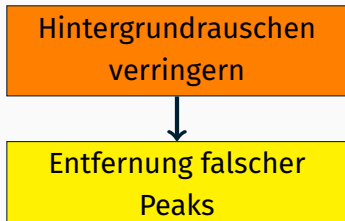
- Algorithmus für die De-Novo-Sequenzierung
- Auswertung von MS2 Spektrogrammen
- Rekonstruktion der AS-Sequenz
- Hilfsmittel: Spektrums-Graph

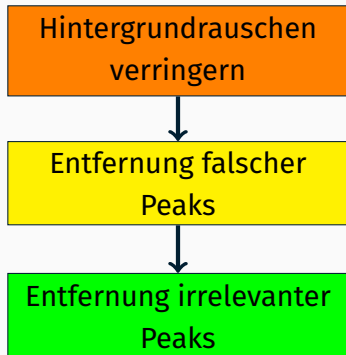
- Zwei MS2 Spektren verwenden

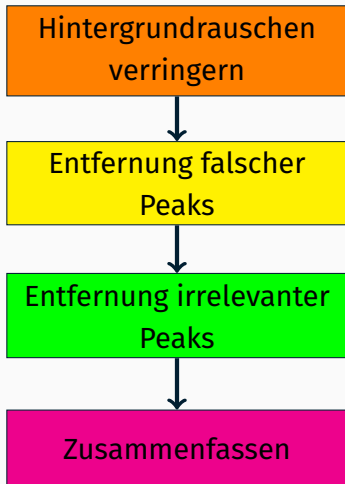
- Zwei MS2 Spektren verwenden
- Unterschiedliche Fragmentierungsmethoden pro Spektrum

- Zwei MS2 Spektren verwenden
- Unterschiedliche Fragmentierungsmethoden pro Spektrum
- Ziel: bessere Sequenzierungsergebnisse

Hintergrundrauschen
verringern



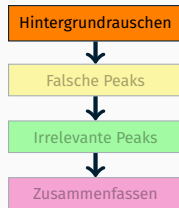




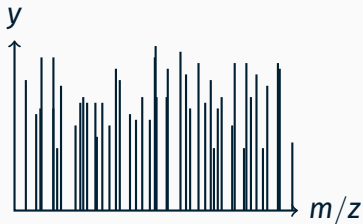
- Überpriorisierung fehlerhafter Daten vermeiden



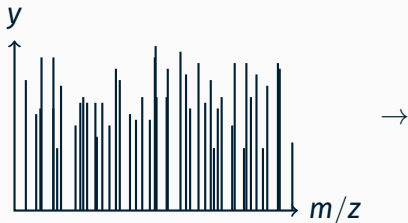
- Überpriorisierung fehlerhafter Daten vermeiden
- Tool: $\ln()$



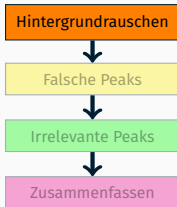
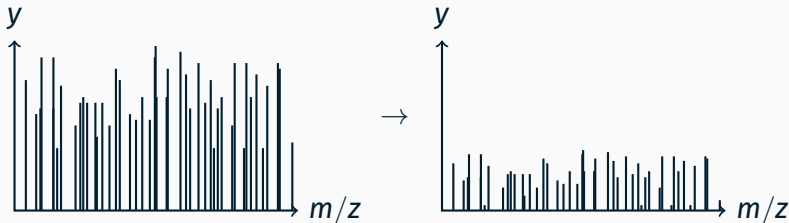
- Überpriorisierung fehlerhafter Daten vermeiden
- Tool: $\ln()$



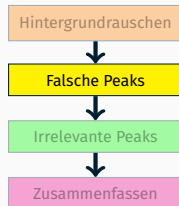
- Überpriorisierung fehlerhafter Daten vermeiden
- Tool: $\ln()$



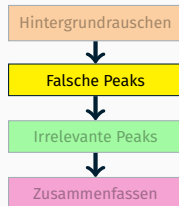
- Überpriorisierung fehlerhafter Daten vermeiden
- Tool: $\ln()$



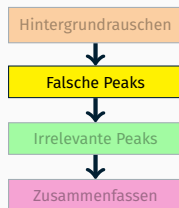
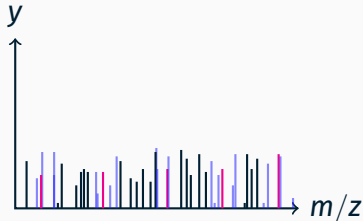
- Hintergrundrauschen könnte auch AS mit Isotop sein



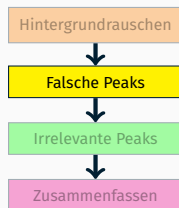
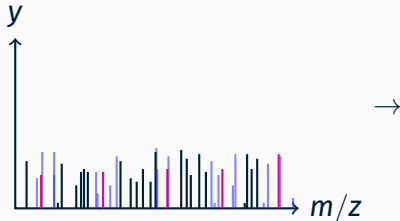
- Hintergrundrauschen könnte auch AS mit Isotop sein
- m/z wählen, die garantiert von AS stammen



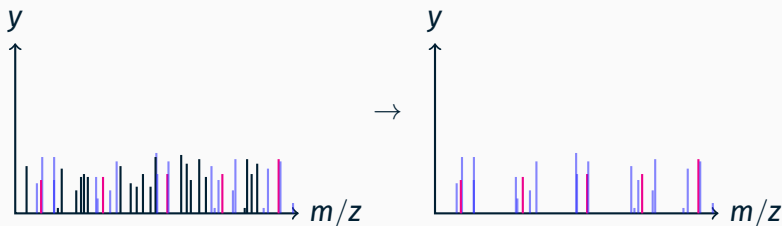
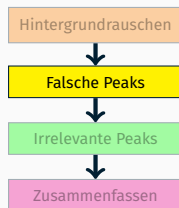
- Hintergrundrauschen könnte auch AS mit Isotop sein
- m/z wählen, die garantiert von AS stammen
- Peaks mit definierten Abstand auswählen



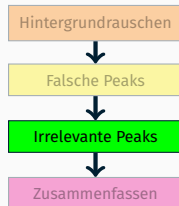
- Hintergrundrauschen könnte auch AS mit Isotop sein
- m/z wählen, die garantiert von AS stammen
- Peaks mit definierten Abstand auswählen



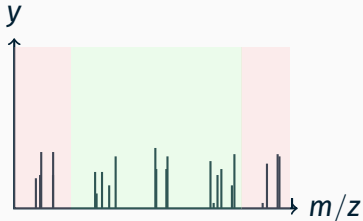
- Hintergrundrauschen könnte auch AS mit Isotop sein
- m/z wählen, die garantiert von AS stammen
- Peaks mit definierten Abstand auswählen



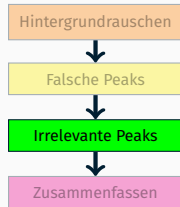
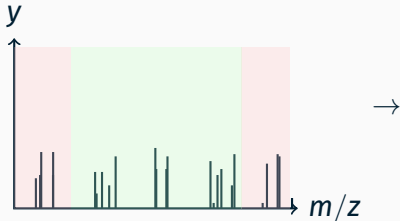
- Peaks aus irrelevantem Intervall entfernen



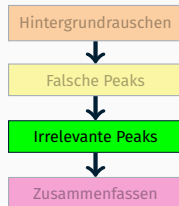
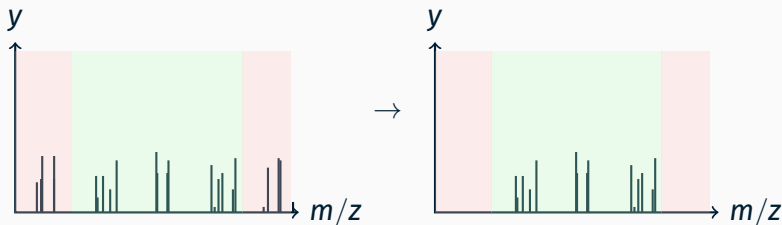
- Peaks aus irrelevantem Intervall entfernen



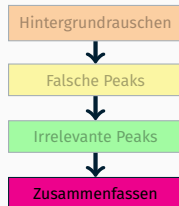
- Peaks aus irrelevantem Intervall entfernen



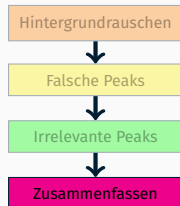
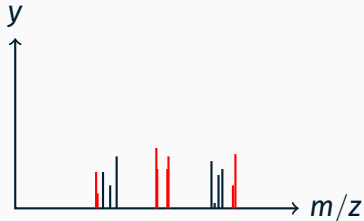
- Peaks aus irrelevantem Intervall entfernen



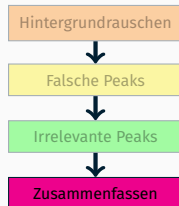
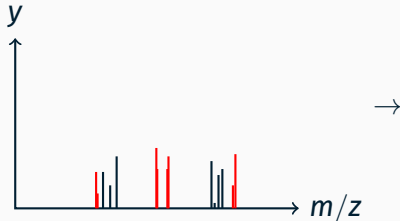
- Zusammenfassen von Peaks



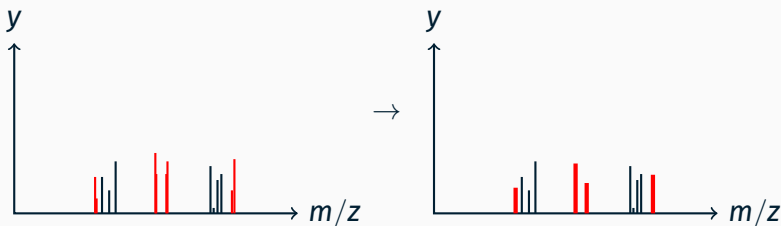
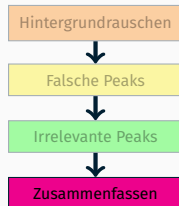
- Zusammenfassen von Peaks
- Abstand einen Schwellwert unterschreitet

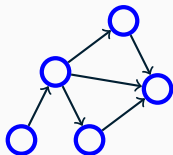
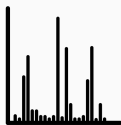


- Zusammenfassen von Peaks
- Abstand einen Schwellwert unterschreitet
- $y = \text{Median}$ aus zusammengefassten Peaks

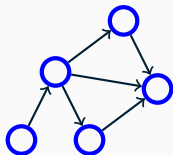
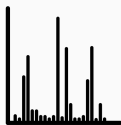


- Zusammenfassen von Peaks
- Abstand einen Schwellwert unterschreitet
- $y = \text{Median}$ aus zusammengefassten Peaks



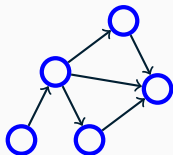
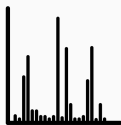


- Verwendung vorverarbeiteter MS2 Spektren

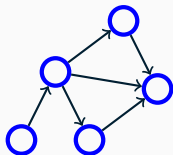
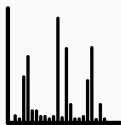


- Verwendung vorverarbeiteter MS2 Spektren
- Peaks $\hat{=}$ Knoten

pNovo+ Algorithmus – Bildung eines Spektrumsgraphen

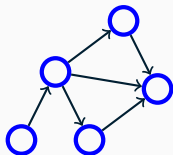
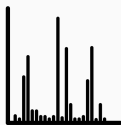


- Verwendung vorverarbeiteter MS2 Spektren
- Peaks $\hat{=}$ Knoten
- Knoten bekommen eine „Masse“



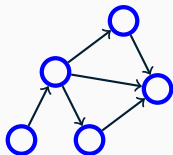
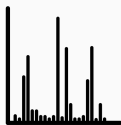
- Verwendung vorverarbeiteter MS2 Spektren
- Peaks $\hat{=}$ Knoten
- Knoten bekommen eine „Masse“
- Masse $\hat{=}$ m/z Wert

pNovo+ Algorithmus – Bildung eines Spektrumsgraphen

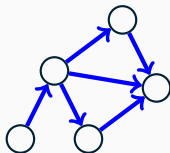
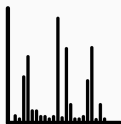


- Verwendung vorverarbeiteter MS2 Spektren
- Peaks $\hat{=}$ Knoten
- Knoten bekommen eine „Masse“
- Masse $\hat{=}$ m/z Wert
- Startknoten (Masse = 0)

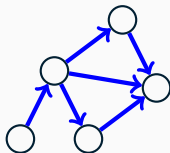
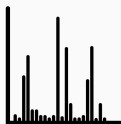
pNovo+ Algorithmus – Bildung eines Spektrumsgraphen



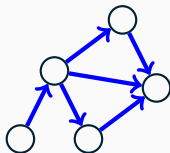
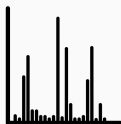
- Verwendung vorverarbeiteter MS2 Spektren
- Peaks $\hat{=}$ Knoten
- Knoten bekommen eine „Masse“
- Masse $\hat{=}$ m/z Wert
- Startknoten (Masse = 0)
- Endknoten (Masse = vorheriger Knoten - 18, 02)



- Gerichtete Kanten zwischen Knotenpaar, wenn:

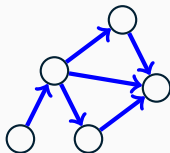
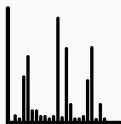


- Gerichtete Kanten zwischen Knotenpaar, wenn:
 - Massendifferenz genau Masse einer AS entspricht



■ Gerichtete Kanten zwischen Knotenpaar, wenn:

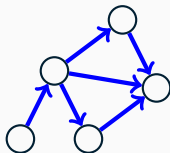
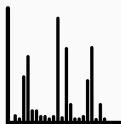
- Massendifferenz genau Masse einer AS entspricht
- Massendifferenz genau Masse zwei AS entsprechen



- Gerichtete Kanten zwischen Knotenpaar, wenn:

- Massendifferenz genau Masse einer AS entspricht
- Massendifferenz genau Masse zwei AS entsprechen

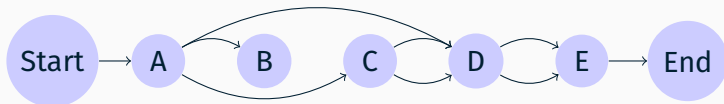
- $N + \binom{n+N-1}{N-1}$ Differenzen $n = 2$ $N = 20$



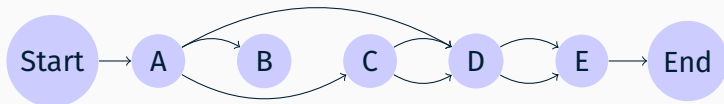
- Gerichtete Kanten zwischen Knotenpaar, wenn:
 - Massendifferenz genau Masse einer AS entspricht
 - Massendifferenz genau Masse zwei AS entsprechen
- $N + \binom{n+N-1}{N-1}$ Differenzen $n = 2 \quad N = 20$
- 230 Differenzen

- Ergebnis: Directed acyclic graph (DAG)

- Ergebnis: Directed acyclic graph (DAG)

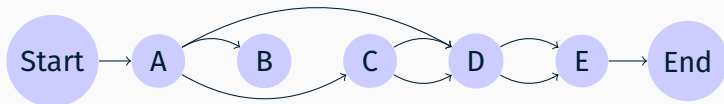


- Ergebnis: Directed acyclic graph (DAG)



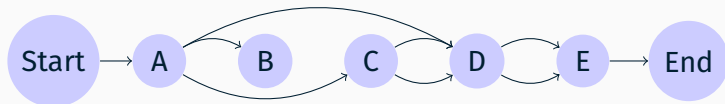
- Alle möglichen Pfade von Start nach End ermitteln

- Ergebnis: Directed acyclic graph (DAG)



- Alle möglichen Pfade von Start nach End ermitteln
- Scoring Funktion „bewertet“ jeden Pfad

- Ergebnis: Directed acyclic graph (DAG)



- Alle möglichen Pfade von Start nach End ermitteln
- Scoring Funktion „bewertet“ jeden Pfad
- Pfad mit dem höchsten Scoring Wert ist das Ergebnis

- 8677 Datensätze





- 8677 Datensätze
- Erfolgreiche Sequenzierungen: 81,2%



- 8677 Datensätze
- Erfolgreiche Sequenzierungen: 81,2%
- Konkurrenzalgorithmus: 71.8%



- 8677 Datensätze
- Erfolgreiche Sequenzierungen: 81,2%
- Konkurrenzalgorithmus: 71.8%
- pNovo+ besser als die Konkurrenz!



- 8677 Datensätze
- Erfolgreiche Sequenzierungen: 81,2%
- Konkurrenzalgorithmus: 71.8%
- pNovo+ besser als die Konkurrenz!
- Side Note: pNovo+ ist frei verfügbar

Open-pNovo Algorithmus

- Peptide sind nicht zwingend stabil

- Peptide sind nicht zwingend stabil
- Wechselwirkungen können die Sequenz abändern

- Peptide sind nicht zwingend stabil
- Wechselwirkungen können die Sequenz abändern
- Posttranslationale Proteinmodifikationen (PTM)

- Peptide sind nicht zwingend stabil
- Wechselwirkungen können die Sequenz abändern
- Posttranslationale Proteinmodifikationen (PTM)
- Mit De-Novo-Algorithmen an sich kein Problem

- Bildung von nicht proteinogenen AS möglich

- Bildung von nicht proteinogenen AS möglich
- AS, die normalerweise nicht in Peptiden vorkommen

- Bildung von nicht proteinogenen AS möglich
- AS, die normalerweise nicht in Peptiden vorkommen
- Spektrogramm zeigt solche AS

- Bildung von nicht proteinogenen AS möglich
- AS, die normalerweise nicht in Peptiden vorkommen
- Spektrogramm zeigt solche AS
- Änderungen können von pNovo+ nicht erkannt werden

- Bildung von nicht proteinogenen AS möglich
- AS, die normalerweise nicht in Peptiden vorkommen
- Spektrogramm zeigt solche AS
- Änderungen können von pNovo+ nicht erkannt werden
- pNovo+ erzeugt zwangsweise Fehler

- Neue Scoring Funktion: RankBoost

- Neue Scoring Funktion: RankBoost
- Machine Learning Algorithmus aus 2003

- Neue Scoring Funktion: RankBoost
- Machine Learning Algorithmus aus 2003
- Erweiterung des AdaBoost Algorithmus

- Neue Scoring Funktion: RankBoost
- Machine Learning Algorithmus aus 2003
- Erweiterung des AdaBoost Algorithmus
- Präferenzen in Datensätzen zu erkennen

- Neue Scoring Funktion: RankBoost
- Machine Learning Algorithmus aus 2003
- Erweiterung des AdaBoost Algorithmus
- Präferenzen in Datensätzen zu erkennen
- Filterung der nicht gültigen AS

■ 20259 Datensätze



- 20259 Datensätze
- Erfolgreiche Sequenzierungen: 76,3%





- 20259 Datensätze
- Erfolgreiche Sequenzierungen: 76,3%
- pNovo+: 68,5%



- 20259 Datensätze
- Erfolgreiche Sequenzierungen: 76,3%
- pNovo+: 68,5%
- Zwei Konkurrenzalgorithmen: 65,8% sowie 39,9%



- 20259 Datensätze
- Erfolgreiche Sequenzierungen: 76,3%
- pNovo+: 68,5%
- Zwei Konkurrenzalgorithmen: 65,8% sowie 39,9%
- Open-pNovo ebenfalls besser als die Konkurrenz!

Zusammenfassung

- De-Novo-Sequenzierung leichter durchführbar

- De-Novo-Sequenzierung leichter durchführbar
- Beide Algorithmen liefern **erstklassige** Ergebnisse

- De-Novo-Sequenzierung leichter durchführbar
- Beide Algorithmen liefern **erstklassige** Ergebnisse
- pNovo+ Ansatz mit Spektrumsgraphen ist wirkungsvoll

- De-Novo-Sequenzierung leichter durchführbar
- Beide Algorithmen liefern **erstklassige** Ergebnisse
- pNovo+ Ansatz mit Spektrumsgraphen ist wirkungsvoll
- Open-pNovo erkennt zuverlässig Proben mit PTMs