

Business Intelligence & Process Modelling

Frank Takes

Universiteit Leiden

Lecture 5 — Reinforcement Learning & Python

Recap

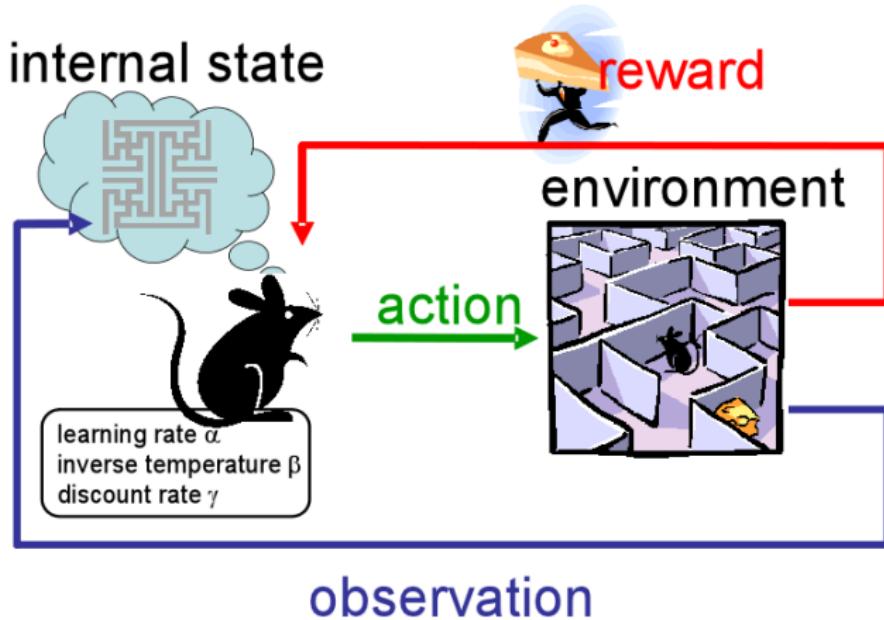
- **Business Intelligence:** anything that aims at providing actionable information that can be used to support business decision making
 - Business Intelligence
 - Visual Analytics
 - **Descriptive Analytics**
 - **Predictive Analytics**
- Process Modelling (April and May)

Categories of techniques

- Machine learning
 - Supervised learning: learning on labeled data
 - Semi-supervised learning: partially labeled data
 - Unsupervised learning: leaning/mining on unlabeled data
 - **Reinforcement learning:** agents learning to act in an environment

Reinforcement learning

Reinforcement learning



https://www.cs.utexas.edu/~eladlieb/rl_interaction.png

Google Deepmind

- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves, Martin Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharshan Kumaran, Daan Wierstra, Shane Legg & Demis Hassabis,
Human-level control through deep reinforcement learning,
Nature 518, 529–533, 2015.



<http://dx.doi.org/10.1038/nature14236>

Google Deepmind

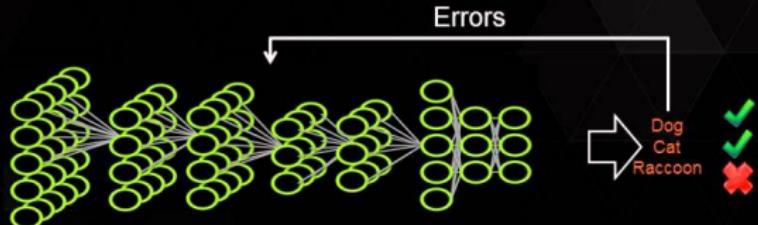
- Silver et al.. **Mastering the game of Go with deep neural networks and tree search**, *Nature* 529, 484–489, 2016.



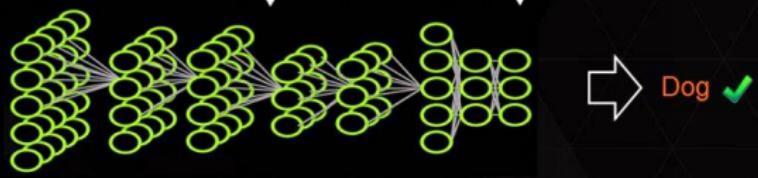
Deep learning

DEEP LEARNING APPROACH

Train:



Deploy:

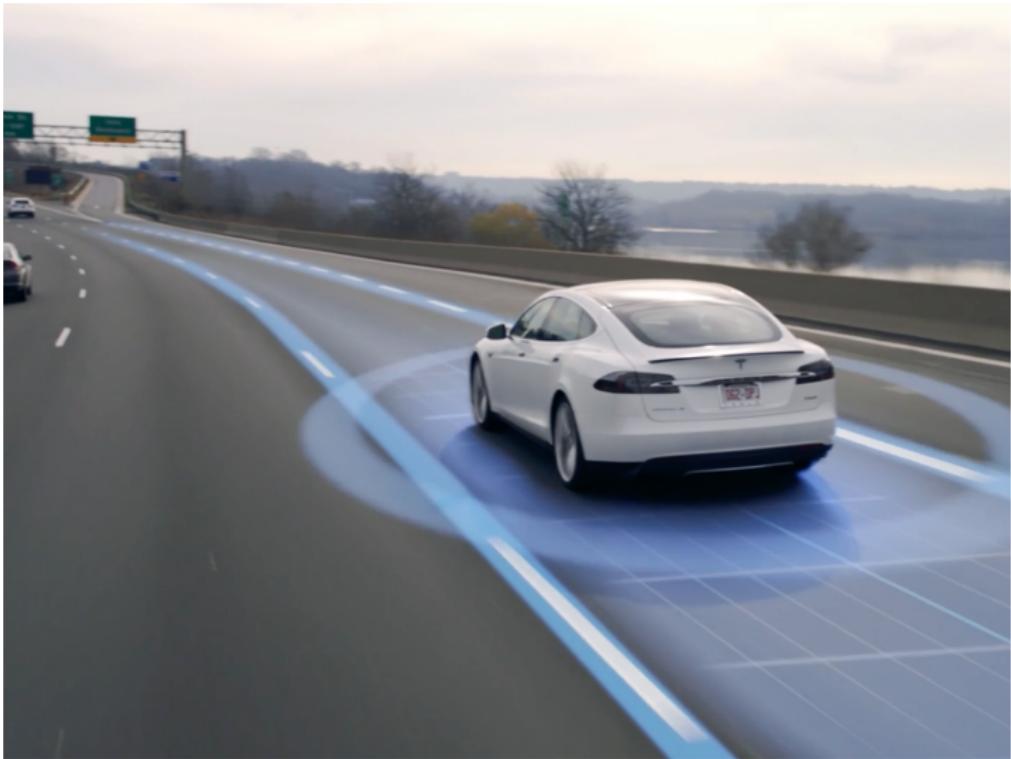


Watson wins Jeopardy



<https://www.youtube.com/watch?v=YgYSv2KSyWg>

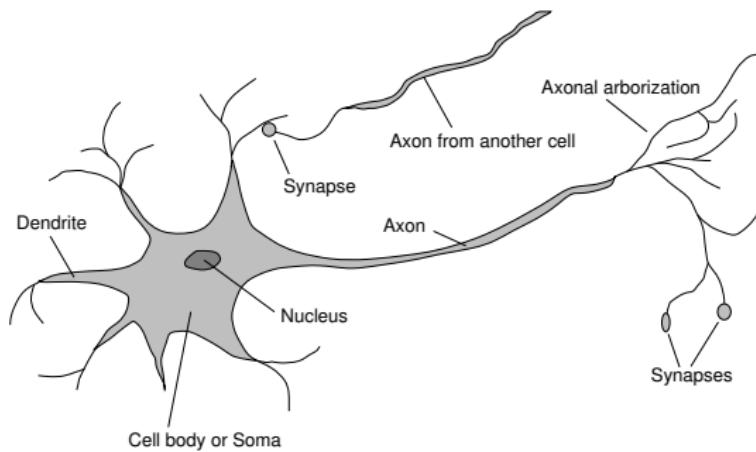
Self-driving cars



Neurale netwerken

Hersen

De menselijke hersenen bestaan uit 10^{11} **neuronen** (grootte ≈ 0.1 mm; meer dan 20 types), die onderling met **axonen** (lengte ≈ 1 cm) en **synapsen** verbonden zijn:

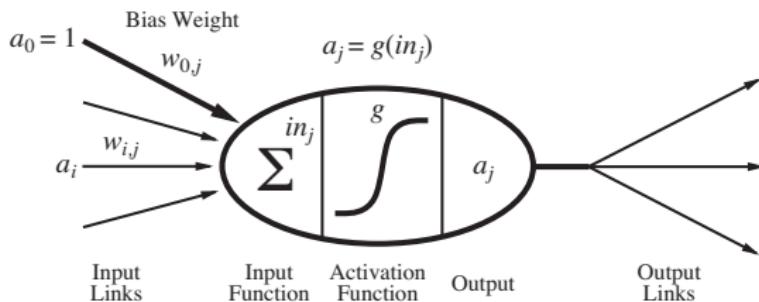


Werking hersenen

- Signalen worden doorgegeven via een nogal gecompliceerde electro-chemische reactie.
- Als de electrische potentiaal van het cellichaam een zekere **drempelwaarde** haalt, wordt een puls/actie-potentiaal/spike-train op het axon gezet.
- Er zijn **excitatory** (verhogen potentiaal) en **inhibitory** (verlagen potentiaal) synapsen.

Units

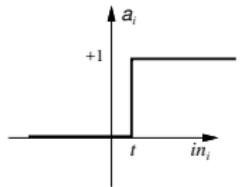
We modelleren de neuronen door middel van eenheden (units) die we ook weer **neuronen** noemen:



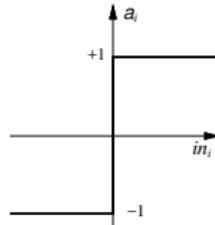
Er geldt: de input van neuron i is $in_i = \sum_j W_{j,i} a_j$ (de gewogen som van de inputs), waarbij de a_j 's de **activaties** van de inkomende verbindingen (inputs) zijn, en $W_{j,i}$ hun gewichten ($W_{j,i}$ weegt de verbinding tussen de twee neuronen j en i). De activatie (output) van neuron i is $a_i = g(in_i)$, met g de **activatie-functie**.

Activatie-functies

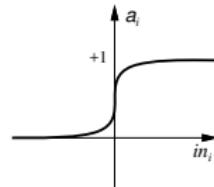
Veelgebruikte activatie-functies zijn:



(a) Step function



(b) Sign function



(c) Sigmoid function

$$\text{step}_t(x) = 1 \text{ als } x \geq t; 0 \text{ als } x < t \quad (\text{drempelwaarde } t)$$

$$\text{sign}(x) = 1 \text{ als } x \geq 0; -1 \text{ als } x < 0$$

$$\text{sigmoid}(x) = 1/(1 + e^{-\beta x}) \quad (\text{vaak kiest men } \beta = 1)$$

Bias-knopen

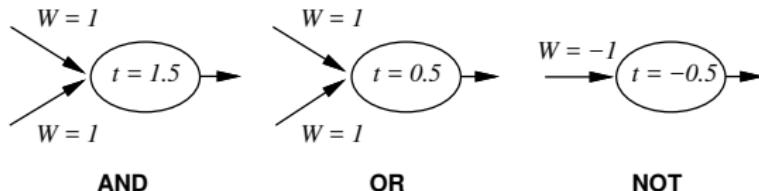
De drempelwaarde t binnen de neuronen kan “gesimuleerd” worden door een extra (constante) activatie -1 in een zogeheten **bias-knoop** 0 met gewicht $W_{0,i} = t$ op de verbinding van 0 naar i :

$$\text{step}_{\color{red}t\color{black}} \left(\sum_{j=1}^n W_{j,i} a_j \right) = \text{step}_{\color{red}0\color{black}} \left(\sum_{j=0}^n W_{j,i} a_j \right)$$

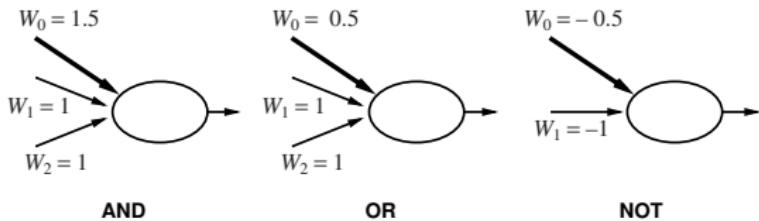
Zo worden drempelwaardes en gewichten uniform behandeld.
De functie step_0 heet ook wel de Heaviside-functie, en lijkt op de functie sign .

Booleaanse functies

Alle Booleaanse functies kunnen gerepresenteerd worden door netwerken met geschikte neuronen:

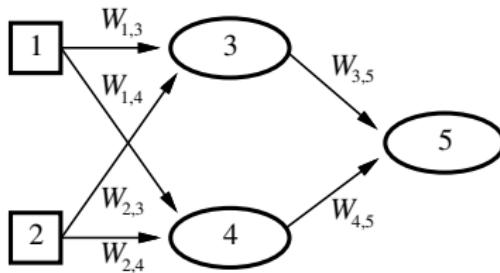


En mét bias-knopen:



Soorten netwerken

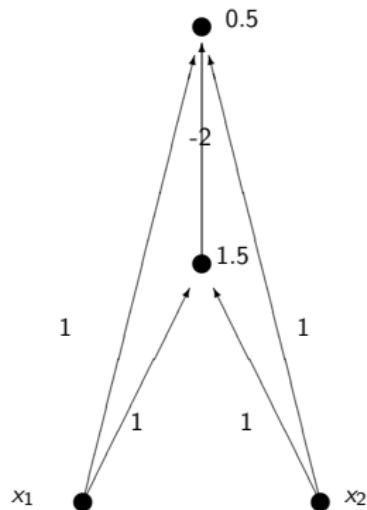
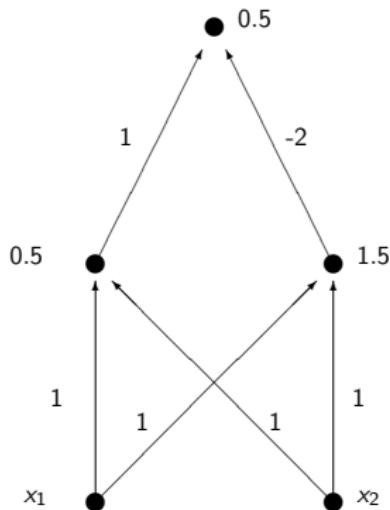
Een **feed-forward netwerk** heeft gerichte takken en geen cykels (i.t.t. een **recurrent netwerk**). Vaak is zo'n netwerk in **lagen** georganiseerd:



Dit netwerk heeft twee input-knopen 1 en 2, twee **verborgen** (= hidden) knopen 3 en 4, en één uitvoer-knoop 5. Het represeneert de functie:

$$\begin{aligned}
 a_5 &= g(W_{3,5} a_3 + W_{4,5} a_4) \\
 &= g(W_{3,5} g(W_{1,3} a_1 + W_{2,3} a_2) \\
 &\quad + W_{4,5} g(W_{1,4} a_1 + W_{2,4} a_2))
 \end{aligned}$$

Voorbeelden



De drempels staan naast de knopen.

Het linker feed-forward netwerk representeert de XOR-functie, de verborgen units zijn in feite een OR en een AND.

Het rechter netwerk representeert dezelfde functie, maar is niet “conventioneel” (geen lagen).

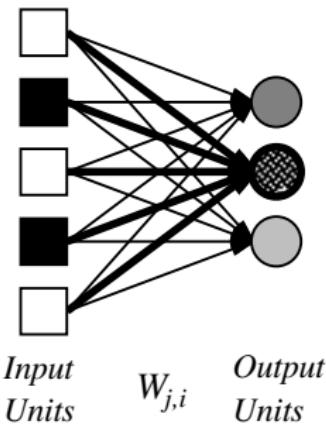
Feed-forward netwerken

Een feed-forward netwerk zonder verborgen neuronen heet een **perceptron**. Een **meerlaags** (= multi-layer) netwerk heeft één of meer verborgen lagen, en alle pijlen van laag ℓ gaan naar laag $(\ell + 1)$. Met één (voldoende grote) laag met verborgen units kunt je elke continue functie benaderen, en met twee lagen zelfs elke discontinue functie (Cybenko).

De output van een netwerk hangt af van de parameters, de gewichten. Bij *te veel* parameters is er gevaar voor **overfitting**: het netwerk generaliseert dan niet goed.

Perceptron

In een perceptron (geen verborgen units!) zijn de uitvoerknopen onafhankelijk:

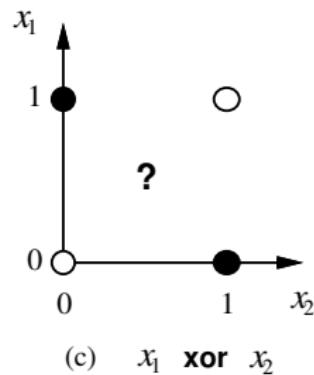
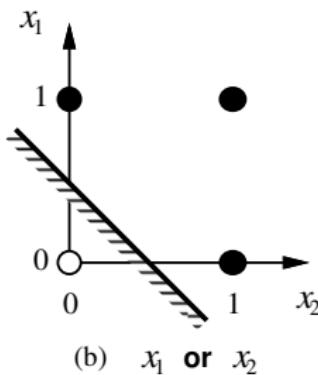
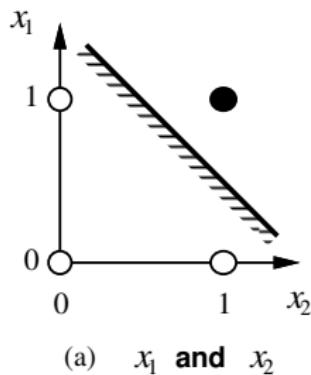


Voor een enkele output-unit geldt dat de uitvoer 1 is indien $\sum_j W_j x_j = W_0 x_0 + W_1 x_1 + \dots + W_n x_n \geq 0$ en anders 0. Hierbij is (x_1, \dots, x_n) de invoer en $x_0 = -1$ de bias-knoop.

Wat kan een perceptron?

De **majority functie** (uitvoer 1 \Leftrightarrow meer dan de helft van de n inputs is 1) kan eenvoudig gemaakt worden: kies $W_0 = n/2$ en $W_j = 1$ voor $j = 1, 2, \dots, n$.

De vergelijking $-W_0 + W_1x_1 + \dots + W_nx_n \geq 0$ laat zien dat je precies Booleaanse functies kunt maken die **lineair te scheiden** zijn, de XOR-functie dus niet:



Hoe leert een perceptron?

Gegeven genoeg trainings-voorbeelden, kan een perceptron elke Booleaanse lineair te scheiden functie leren. Een (hyper)vlak scheidt positieve en negatieve voorbeelden.

Rosenblatt's algoritme uit 1957/60 werkt als volgt:

$$W_j \leftarrow W_j + \alpha \cdot x_j \cdot \text{Error}$$

met Error = correcte uitvoer – net-uitvoer en $\alpha > 0$ de **leersnelheid** (= learning rate). De correcte uitvoer heet wel de **target**, het doel.

Als Error = 1 en $x_j = 1$, wordt W_j ietsje opgehoogd, in de hoop dat de net-uitvoer hoger wordt.

Perceptron — voorbeeld

We willen een perceptron x_1 AND x_2 leren, met $\alpha = 0.1$:

	x0	x1	x2	w_0	w_1	w_2	uitvoer	target	error
1	-1	1	1	-0.160	-0.606	-0.217	0.000	1.000	1.000
2	-1	1	0	-0.260	-0.506	-0.117	0.000	0.000	0.000
3	-1	0	1	-0.260	-0.506	-0.117	1.000	0.000	-1.000
4	-1	1	0	-0.160	-0.506	-0.217	0.000	0.000	0.000
5	-1	1	0	-0.160	-0.506	-0.217	0.000	0.000	0.000
6	-1	1	1	-0.160	-0.506	-0.217	0.000	1.000	1.000
7	-1	0	0	-0.260	-0.406	-0.117	1.000	0.000	-1.000
8	-1	0	0	-0.160	-0.406	-0.117	1.000	0.000	-1.000
...									
70	-1	1	0	0.140	0.194	0.183	1.000	0.000	-1.000
71	-1	1	1	0.240	0.094	0.183	1.000	1.000	0.000
...									

Probleem: wanneer stop je?

Perceptron — anders leren

We kunnen in plaats van de discontinue stapfunctie ook een gladde sigmoide g gebruiken in de knopen. Een vergelijkbaar leeralgoritme gaat dan als volgt.

Zij $E = \frac{1}{2}\text{Error}^2 = \frac{1}{2}\left(y - g\left(\sum_{j=0}^n W_j x_j\right)\right)^2$ met y de target. Met behulp van **gradient descent** bepalen we in welke richting de fout het snelst *stijgt*:

$$\frac{\partial E}{\partial W_j} = \text{Error} \cdot \frac{\partial \text{Error}}{\partial W_j} = -\text{Error} \cdot g'(\text{in}) \cdot x_j$$

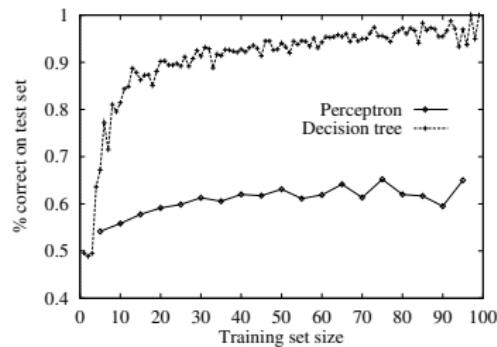
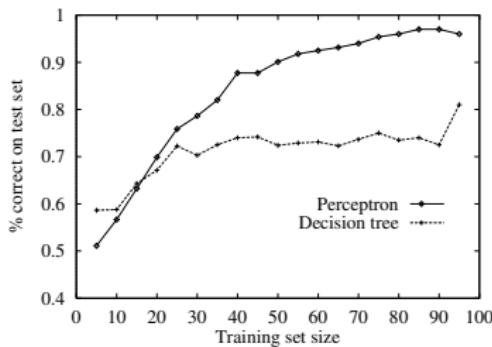
met $\text{in} = \sum_{j=0}^n W_j x_j$. Dus leerregel (let op de extra $-$, we willen de fout laten *dalen*):

$$W_j \leftarrow W_j + \alpha \cdot \text{Error} \cdot g'(\text{in}) \cdot x_j .$$

Perceptron — vergelijking

Voor een lineair te scheiden majority probleem (11 inputs) is het perceptron veel beter dan een decision tree (links)

Maar voor een meer complex probleem is het omgekeerd (rechts)



Codering

Er zijn allerlei keuzes om invoer en uitvoer te coderen. Je kunt bijvoorbeeld **locaal coderen**: stel dat een variabele 3 waarden kan hebben: Geen, Gemiddeld en Veel; dit kun je dan in één knoop coderen als 0.0, 0.5 en 1.0 respectievelijk. Maar ook met drie aparte knopen, en dan als 1–0–0, 0–1–0 en 0–0–1 respectievelijk: **gedistribueerd coderen**. Je kunt zelfs bij de foutmaat ervoor zorgen dat je waarden als 1–1–0.9 niet fijn vindt, en daar van weg trainen (“**softmax**”).

Leerproces: training

Hoe verloopt nu het leren bij Neurale Netwerken, het aanpassen van de gewichten, kortom: de **training**?

Je begint met een **random initialisatie** van de gewichten. Vervolgens biedt je één voor één voorbeelden aan uit een zogeheten **trainingsset**. Deze voorbeelden moeten in een willekeurige volgorde staan. Steeds geef je een invoer, en met behulp van de juiste uitvoer pas je volgens je trainings-algoritme de gewichten aan.

Je gaat net zolang door totdat (bijvoorbeeld) de fout op een apart gehouden **validatieset** niet meer daalt — of zelfs gaat stijgen (overtraining).

Je rapporteert tot slot over het gedrag op een (weer apart gehouden) **testset**.

Cross-validation

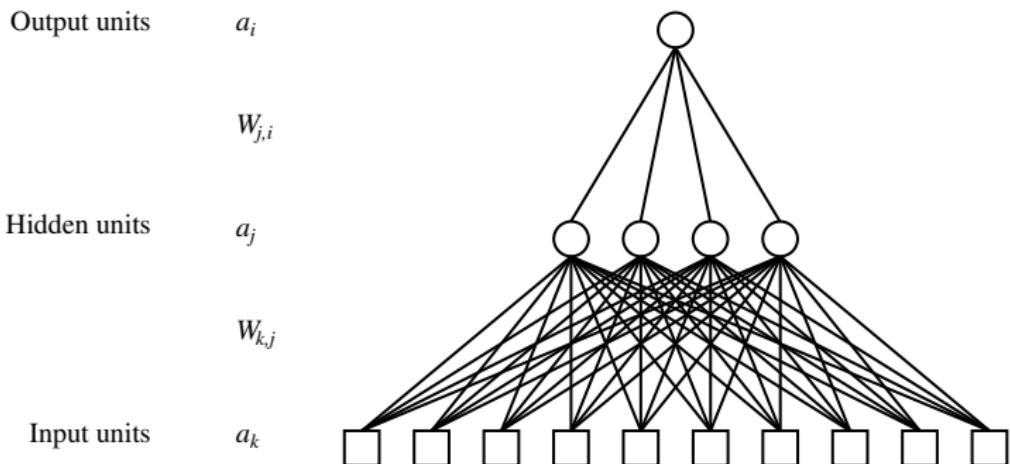
Bij elk leer-algoritme kun je **cross-validation** gebruiken om overfitting tegen te gaan.

Bij “ k -fold cross-validation” (vaak $k = 5$ of $k = 10$) draai je k experimenten, waarbij je steeds een $1/k$ -deel van de data apart zet om als testset te gebruiken — en de rest als trainingsset. De testset is dus steeds een ander random gekozen deel!

Als $k = n$ met n de grootte van de dataset, heet deze techniek wel “leave-one-out”. En dan heb je ook nog “ensemble leren”, AdaBoost, enzovoorts. Zie verder het college Data Mining.

Meerlaags netwerk

We kijken nu naar een **meerlaags neuraal netwerk**. Meestal zijn alle lagen onderling *volledig* verbonden.



De notatie is ietwat dubbelzinnig: zit $W_{1,2}$ op twee plaatsen? Je kunt of knopen doornummers of meerdere W 's hanteren. We gebruiken index i voor de uitvoerlaag, j voor de verborgen laag en k voor de invoerlaag. De a_k 's zijn de input(s) — wat we eerder x_k noemden.

Back-propagation 1

Het meest bekende leerschema voor Neurale Netwerken is **back-propagation** (1969, 198?), waarbij we — nadat een invoer een uitvoer heeft opgeleverd — de fout teruggeven (propageren) van uitvoerlaag richting invoerlaag.

Definieer Error; als de fout in de i -de uitvoer (dat wil zeggen: i -de target minus i -de uitvoer van het netwerk). De leerregel is dan net als eerder:

$$W_{j,i} \leftarrow W_{j,i} + \alpha \cdot a_j \cdot \Delta_i \quad \text{met} \quad \Delta_i = \text{Error}_i \cdot g'(\text{in}_i)$$

voor gewichten $W_{j,i}$ van knoop j in de verborgen laag naar knoop i in de uitvoerlaag.

Back-propagation 2

Voor gewicht $W_{k,j}$ van de k -de invoerknoop naar knoop j in de verborgen laag is de leerregel:

$$W_{k,j} \leftarrow W_{k,j} + \alpha \cdot a_k \cdot \Delta_j \quad \text{met} \quad \Delta_j = g'(\text{in}_j) \sum_i W_{j,i} \Delta_i$$

Hierbij geldt:

α is de leersnelheid

a_k is de activatie van de k -de invoerknoop

in_j is de gewogen invoer voor de j -de verborgen knoop

$W_{j,i}$ is het gewicht op de verbinding tussen

de j -de verborgen knoop en de i -de uitvoerknoop

Δ_i is de i -de "uitvoer-delta", zie de vorige sheet

Back-propagation algoritme

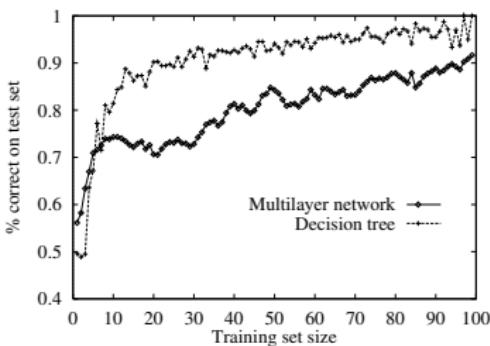
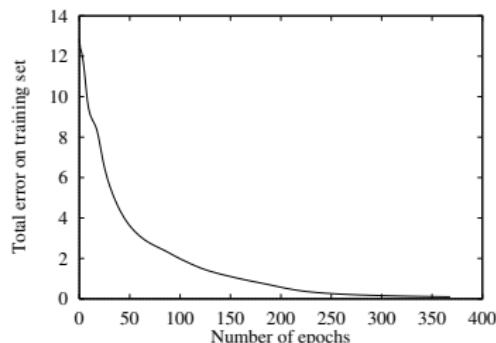
Het back-propagation-algoritme voor een netwerk met één verborgen laag gaat nu als volgt:

```
repeat
    for each  $e$  in trainingsset do
        maak de  $a_k$ 's gelijk aan de  $x_k$ 's
        bereken de  $a_j$ 's en de  $a_i$ 's (outputs)
        bereken de  $\Delta_i$ 's en de  $\Delta_j$ 's
        update de  $W_{j,i}$ 's en de  $W_{k,j}$ 's
    until netwerk "geconvergeerd"
```

Let erop dat de gewichten goed random (?) geinitialiseerd worden. En bied de voorbeelden in random volgorde aan.

Grafieken

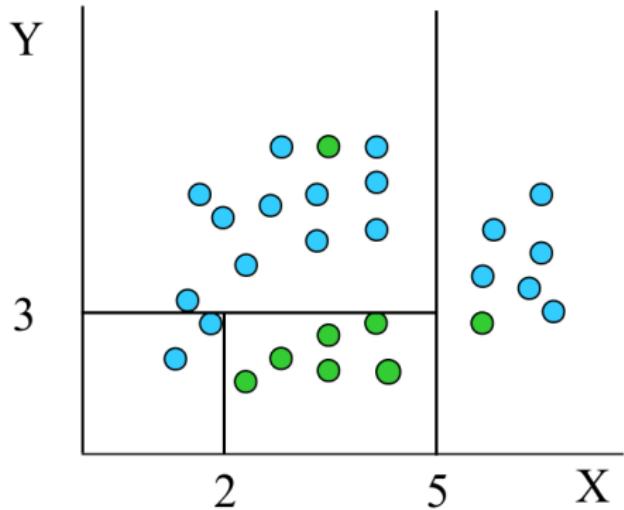
Respectievelijk een trainings-curve en een leercurve:



Een **epoch** is een ronde waarin alle voorbeelden uit de trainingsset één keer één voor één in random volgorde door het netwerk gegaan zijn. Soms heb je ∞ veel voorbeelden.

Er bestaat naast deze **incrementele** benadering ook een **batch**-benadering.

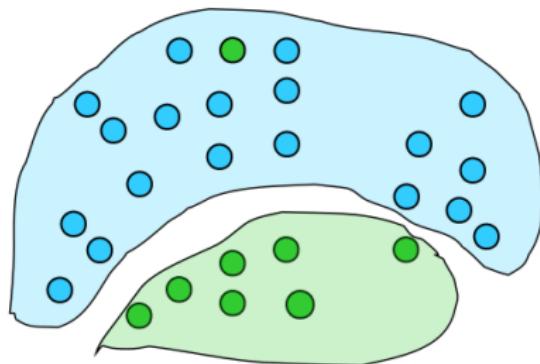
Classification: Decision trees



Decision Tree ($d = 3$)

```
if(X > 5) return BLUE;  
elseif(Y > 3) return BLUE;  
elseif(X > 2) return GREEN;  
else return BLUE;
```

Classification: Neural networks



Neural Networks

- Perceptron
- Multi-level (meerlaags) neural network
- Convolutional neural network

Nog meer soorten

Er zijn nog vele andere soorten netwerken, zoals

- Kohonen's Self Organizing Maps (SOM's)
- Support Vector Machines (SVM's), kernel machines
- Convolutional Neural Networks
(→ Deep Learning), met activatiefunctie
 $g(x) = \max(0, x)$, de "rectifier"
- ...

Python

Python

- Guido van Rossem (1991)
- Directly interpreted (no compilation)
- Multi-paradigm: imperative, functional, procedural
- Cross-platform
- Version 2.7 and version 3.0



tweakers.net/video/11095/polderpioniers-guido-van-rossum-ontwerper-van-de-programmeertaal-python.html

Zen of Python

- Beautiful is better than ugly
- Explicit is better than implicit
- Simple is better than complex
- Complex is better than complicated
- Readability counts
- ...

Zen of Python

- Beautiful is better than ugly
- Explicit is better than implicit
- Simple is better than complex
- Complex is better than complicated
- Readability counts
- ...
- “Pythonic” code

PEP 20 – The Zen of Python: <https://www.python.org/dev/peps/pep-0020/>

Zen of Python

```
import this
"""The Zen of Python, by Tim Peters. (poster by Joachim Jablon)"""

1 Beautiful is better than ugly.
2 Explicit is better than implicit.
3 Simple is better than complex.
4 Complex is better than complicated.
5 Flat is better than nested.
6 Sparse is better than dense.
7 Readability counts.
8 Special cases aren't special enough to break the rules.
9 Although practicality beats purity.
10 raise PythonicError("Errors should never pass silently.")
11 # Unless explicitly silenced.
12 In the face of ambiguity, refuse the temptation to guess.
13 There should be one-- and preferably only one --obvious way to do it.
14 # Although that way may not be obvious at first unless you're Dutch.
15 Now is better than ... never.
16 Although never is often better than rightnow.
17 If the implementation is hard to explain, it's a bad idea.
18 If the implementation is easy to explain, it may be a good idea.
19 Namespaces are one honking great idea -- let's do more of those!
```

Control flow

```
1 # this is a comment
2 def foo(x):
3     if x == 0:
4         bar()
5         baz()
6     else:
7         qux(x)
8         foo(x - 1)
```

Control flow

```
1 # this is a comment
2 def foo(x):
3     if x == 0:
4         bar()
5         baz()
6     else:
7         qux(x)
8         foo(x - 1)
```

- Indentation determines control flow
- Colon : starts indentation
- No semicolon ;
- Parentheses (and) often optional

Types

- Integers int
- Floats float
- True or False bool
- Null object: None
- `type()` function returns type

Boolean operators

Operation	Result
<code>x or y</code>	if <code>x</code> is false, then <code>y</code> , else <code>x</code>
<code>x and y</code>	if <code>x</code> is false, then <code>x</code> , else <code>y</code>
<code>not x</code>	if <code>x</code> is false, then <code>True</code> , else <code>False</code>

<https://docs.python.org/2/library/stdtypes.html>

Comparison operators

Operation	Meaning
<	strictly less than
<=	less than or equal
>	strictly greater than
>=	greater than or equal
==	equal
!=	not equal
is	object identity
is not	negated object identity

<https://docs.python.org/2/library/stdtypes.html>

Control flow

```
1 if statement1 and statement2:  
2     print("1 and 2 are True")  
3 elif statement2 or statement3:  
4     print("2 or 3 or both are True")  
5 else:  
6     print("then just do this")  
7  
8 while i < 5:  
9     print(i)  
10    i += 1 # ++ does not exist  
11  
12 for x in [1,2,3]:  
13     print(x)
```

Lists

```
1 emptylist = []
2 alist = [ 'a' , 'b' , 'c' ]
3 print (alist)
4 >>> [ 'a' , 'b' , 'c' , 'x' ]
5 alist.append(8)
6 print (alist)
7 >>> [ 'a' , 'b' , 'c' , 'x' , 8]
8 print (len(alist))
9 >>> 5
10 fastlist = [x**2 for x in
               range(0,5)]
```

Lists

```
1 emptylist = []
2 alist = [ 'a' , 'b' , 'c' ]
3 print (alist)
4 >>> [ 'a' , 'b' , 'c' , 'x' ]
5 alist.append(8)
6 print (alist)
7 >>> [ 'a' , 'b' , 'c' , 'x' , 8]
8 print (len(alist))
9 >>> 5
10 fastlist = [x**2 for x in
               range(0,5)]
```

- Ordered sequences
- Indexed from $L[0]$ to $\text{len}(L)-1$
- Can be joined and sliced

List operators

Operation	Result
<code>x in s</code>	True if an item of <code>s</code> is equal to <code>x</code> , else False
<code>x not in s</code>	False if an item of <code>s</code> is equal to <code>x</code> , else True
<code>s + t</code>	the concatenation of <code>s</code> and <code>t</code>
<code>s * n, n * s</code>	equivalent to adding <code>s</code> to itself <code>n</code> times
<code>s[i]</code>	<i>i</i> th item of <code>s</code> , origin 0
<code>s[i:j]</code>	slice of <code>s</code> from <i>i</i> to <i>j</i>
<code>s[i:j:k]</code>	slice of <code>s</code> from <i>i</i> to <i>j</i> with step <i>k</i>
<code>len(s)</code>	length of <code>s</code>
<code>min(s)</code>	smallest item of <code>s</code>
<code>max(s)</code>	largest item of <code>s</code>
<code>s.index(x)</code>	index of the first occurrence of <code>x</code> in <code>s</code>
<code>s.count(x)</code>	total number of occurrences of <code>x</code> in <code>s</code>

<https://docs.python.org/2/library/stdtypes.html>

Functions

```
1 def printthis(wellthis):  
2     print (wellthis)  
3     return;  
4  
5 def multi(val1, val2):  
6     return (val1 * val2)  
7  
8 def returnmeanddouble(me):  
9     return me, me*2;  
10  
11 print (multi(val2=8, val1=4))  
12 >>> 32  
13 x, y = returnmeanddouble(8)  
14 print (x, y)  
15 >>> 8 16
```

Functions

```
1 def printthis(wellthis):  
2     print (wellthis)  
3     return;  
4  
5 def multi(val1, val2):  
6     return (val1 * val2)  
7  
8 def returnmeanddouble(me):  
9     return me, me*2;  
10  
11 print (multi(val2=8, val1=4))  
12 >>> 32  
13 x, y = returnmeanddouble(8)  
14 print (x, y)  
15 >>> 8 16
```

- May or may not return a value
- Optional arguments
- Optionally named arguments
- May return tuple (immutable list) of values

Sets

```
1 my_set = {1, 2, 3, 4, 8,
2   102835643}
3 print(my_set)
4 someSet = {0, (), False}
5 emptyset = {}
```

Sets

```
1 my_set = {1, 2, 3, 4, 8,
2   102835643}
3 print(my_set)
4 someSet = {0, (), False}
5 emptyset = {}
```

- Unordered collection
- Not indexed
- Apply common set operators

Set operators

Operation	Equivalent	Result
<code>len(s)</code>		number of elements in set s (cardinality)
<code>x in s</code>		test x for membership in s
<code>x not in s</code>		test x for non-membership in s
<code>s.issubset(t)</code>	$s \leq t$	test whether every element in s is in t
<code>s.issuperset(t)</code>	$s \geq t$	test whether every element in t is in s
<code>s.union(t)</code>	$s \mid t$	new set with elements from both s and t
<code>s.intersection(t)</code>	$s \& t$	new set with elements common to s and t
<code>s.difference(t)</code>	$s - t$	new set with elements in s but not in t
<code>s.symmetric_difference(t)</code>	$s \wedge t$	new set with elements in either s or t but not both
<code>s.copy()</code>		new set with a shallow copy of s

<https://docs.python.org/2/library/stdtypes.html>

Dictionaries

```
1 stuff = { 'name': 'Frank', 'age':30, 'height': 6*30+14}
2 print (stuff['name'])
3 >>> Frank
4 print (stuff['age'], stuff['height'])
5 >>> 30 194
6 stuff['city'] = "Leiden"
7 print stuff['city']
8 >>> Leiden
```

Dictionaries

```
1 stuff = { 'name': 'Frank', 'age':30, 'height': 6*30+14}
2 print (stuff['name'])
3 >>> Frank
4 print (stuff['age'], stuff['height'])
5 >>> 30 194
6 stuff['city'] = "Leiden"
7 print stuff['city']
8 >>> Leiden
```

- Unordered indexed collection
- Freely combine datatypes
- Hashing, so 'fast'

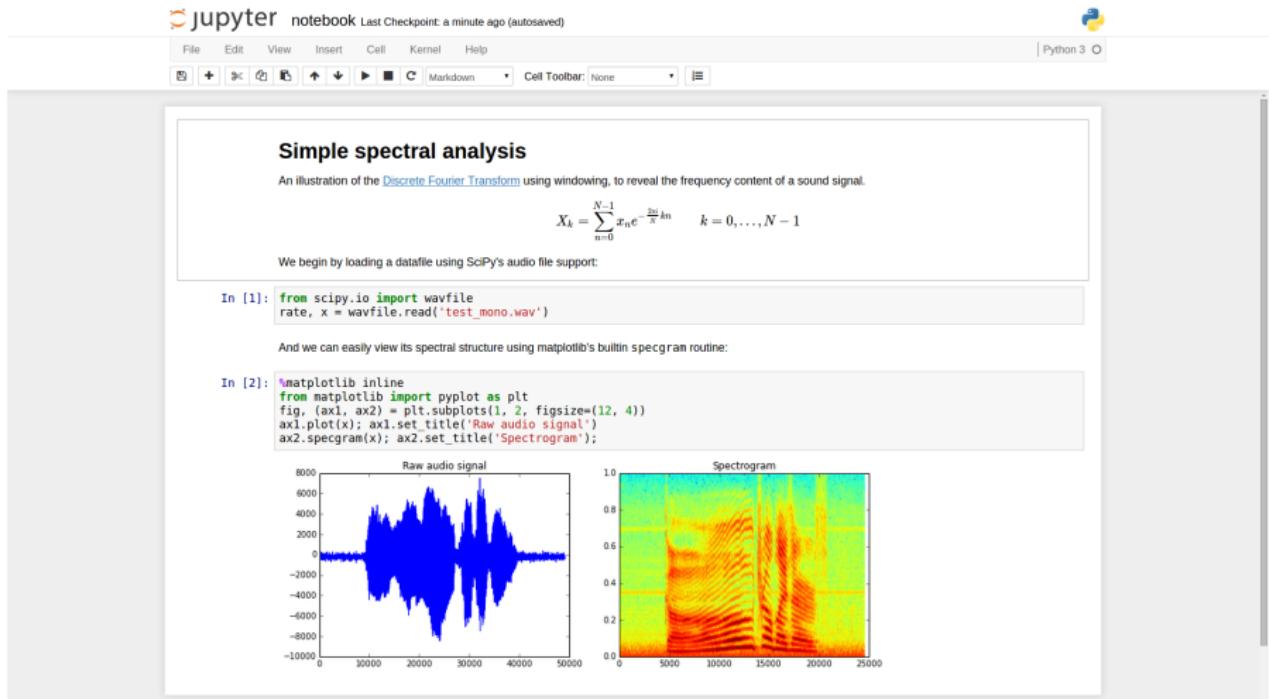
Modules/packages

- File `fibo.py` with Fibonacci functions `fib1()` and `fib2()`
- `import fibo`
Call using `fibo.fib1(...)`
- `from fibo import fib, fib2`
Call using `fib1(...)`
- `from fibo import *`
Call using `fib1(...)`
- Standard modules: `str`, `sys`, etc.

Console or script?

- Execute interactively from `python` console
- Execute one script using `python filename.py` command
- Use Jupyter notebooks
- Spyder

... notebook



The screenshot shows a Jupyter Notebook interface with the title "jupyter notebook Last Checkpoint: a minute ago (autosaved)". The toolbar includes File, Edit, View, Insert, Cell, Kernel, Help, and a Cell Toolbar dropdown set to "None". A Python 3 logo is also present.

Simple spectral analysis

An illustration of the [Discrete Fourier Transform](#) using windowing, to reveal the frequency content of a sound signal.

$$X_k = \sum_{n=0}^{N-1} x_n e^{-\frac{2\pi}{N} kn} \quad k = 0, \dots, N - 1$$

We begin by loading a datafile using SciPy's audio file support:

```
In [1]: from scipy.io import wavfile
rate, x = wavfile.read('test_mono.wav')
```

And we can easily view its spectral structure using matplotlib's builtin specgram routine:

```
In [2]: %matplotlib inline
from matplotlib import pyplot as plt
fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(12, 4))
ax1.plot(x); ax1.set_title('Raw audio signal')
ax2.specgram(x); ax2.set_title('Spectrogram')
```

Two plots are displayed:

- The first plot, titled "Raw audio signal", shows a blue line graph of the audio waveform over time, ranging from 0 to 50,000 samples. The y-axis ranges from -10,000 to 8,000.
- The second plot, titled "Spectrogram", is a heatmap showing the frequency spectrum. The x-axis represents time from 0 to 25,000 samples, and the y-axis represents frequency from 0.0 to 1.0. The plot shows several vertical bands of energy, indicating discrete frequency components.

Pandas

Pandas

- Python package
- Good for tabular data
- Based on DataFrames
- Extensions for statistics
- Extensions for visualization



<http://pandas.pydata.org>

Pandas essentials

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4
5 df = pd.Function()
```

Loading and saving data

```
1 import pandas as pd
2
3 # tab separated
4 df = pd.read_csv('filename.csv', sep='\t')
5
6 # no index column
7 df2 = pd.read_csv('filename.csv', index_col=None)
8
9 # limit number of rows
10 df2 = pd.read_csv('filename.csv', nrows=20)
11 df.to_csv('out.csv')
```

Viewing data

```
1 import pandas as pd
2
3 df = pd.read_csv('filename.csv', sep='\t')
4 df.head(10)
5 df.tail(30)
6 df.describe
7 df.shape
8 df.columns
9 df.values
```

Selecting data

```
1 import pandas as pd
2
3 df = pd.read_csv('filename.csv', sep='\t')
4
5 # get one column
6 df['columnX']
7
8 # get slice of data
9 df.iloc[3:5,0:2] # rows,columns
10
11 # select based on value
12 df[df.columnX > 0]
13
14 # get copy of data where column value is in list
15 df2 = df[df['columnX'].isin(['something','else'])]
```

Data types and Missing data

```
1 import pandas as pd
2
3 df = pd.read_csv('filename.csv', sep='\t')
4
5 # drop any row with missing data
6 df.dropna(how='any')
7
8 # fill any missing value with integer 5
9 df.fillna(value=5)
10
11 # interpret column as certain type
12 df["col1"] = df["col1"].astype("float")
```

http://pandas.pydata.org/pandas-docs/stable/missing_data.html

Adding stuff to data

```
1 # add a row
2 s = df.iloc[2] # get row 2
3 df.append(s, ignore_index=True)
4
5 # add multiple columns
6 df1 = pd.read_csv(sys.argv[1], sep=';')
7 df2 = pd.read_csv(sys.argv[2], sep="\t")
8 df3 = pd.merge(df1, df2, left_on=["idcol1"], right_on=[ "idcol2"], how="left")
```

Grouping data

```
1 import pandas as pd
2
3 df = pd.read_csv('filename.csv', sep='\t')
4
5 # group the data by somefield, sum numvalue into
6 # mysum
7 df['mysum'] = df.groupby('somefield')['numvalue'].transform('sum')
8
9 # remove duplicate rows
10 df = df.drop_duplicates(cols = 'somefield')
```

Matplotlib

Matplotlib

- Library for plotting
- 2D plotting
- Within Python: PYPLOT
- Large collection of example code snippets

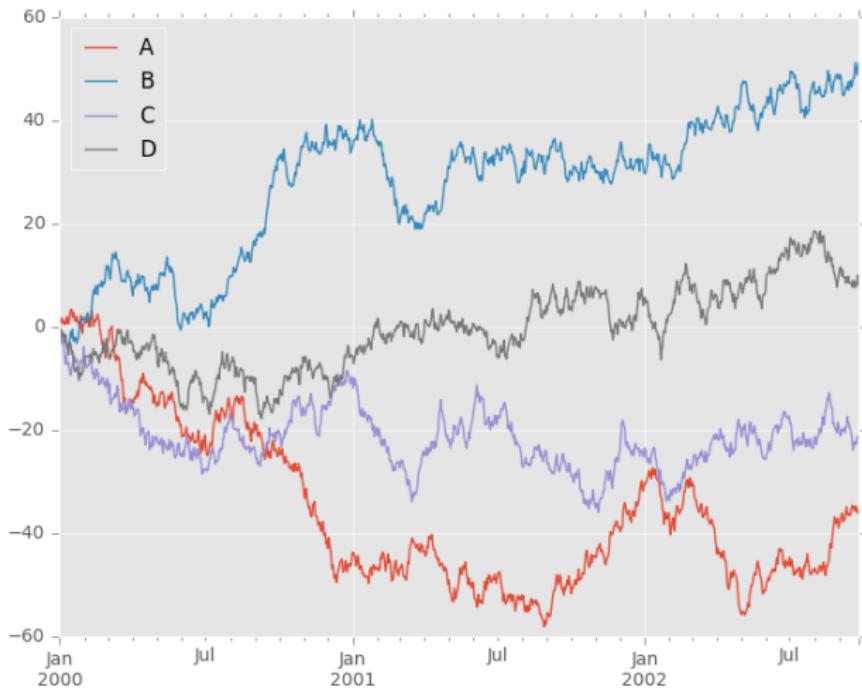
matplotlib



Matplotlib

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 df = pd.read_csv('filename.csv', sep='\t')
5
6 # plot all columns against index to window
7 df.plot(kind='line')
8 df.plot.bar() # alternative
9
10 # plot two columns to window
11 df.plot(x='colA', y='colB')
12
13 # save last plot to file
14 plt.savefig('resultfile.pdf', dpi=600)
```

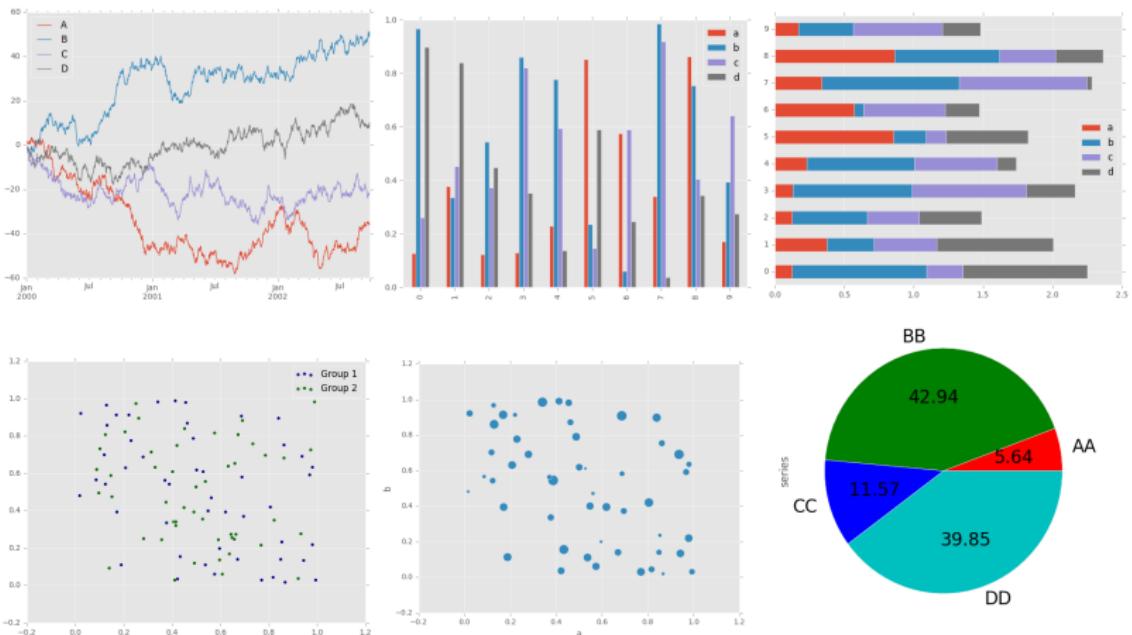
Matplotlib 'line'



Plot types

- 'bar' or 'barh' for bar plots
- 'hist' for histogram
- 'box' for boxplot
- 'kde' or 'density' for density plots
- 'area' for area plots
- 'scatter' for scatter plots
- 'hexbin' for hexagonal bin plots
- 'pie' for pie plots

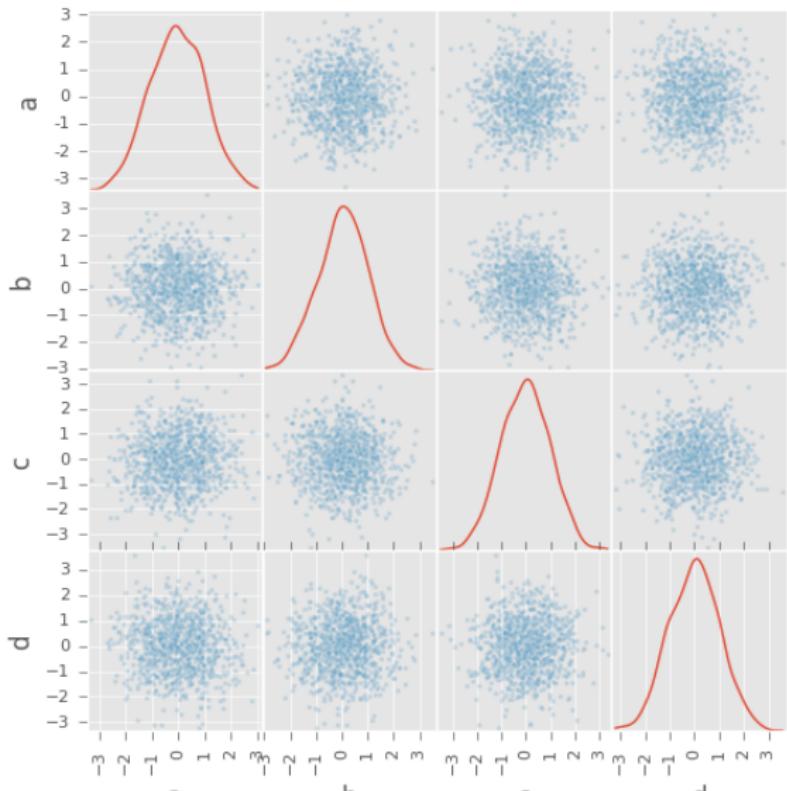
Matplotlib types



Scatter Matrix

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 from pandas.tools.plotting import scatter_matrix
4
5 # assume numeric input
6 df = pd.read_csv('filename.csv', sep='\t')
7
8 # plot all columns against index to window
9 scatter_matrix(df, alpha=0.2, figsize=(6, 6),
10                 diagonal='kde')
11
12 # save last plot to file
13 plt.savefig('resultfile.pdf', dpi=600)
```

Scatter Matrix



Control everything

- Axis scale
- Secondary axis
- Axis values and labels
- Custom legend
- Colors and gradients
- ...

Control everything

- Axis scale
- Secondary axis
- Axis values and labels
- Custom legend
- Colors and gradients
- ...
- <http://pandas.pydata.org/pandas-docs/stable/visualization.html>

Scikit-learn

Scikit-learn

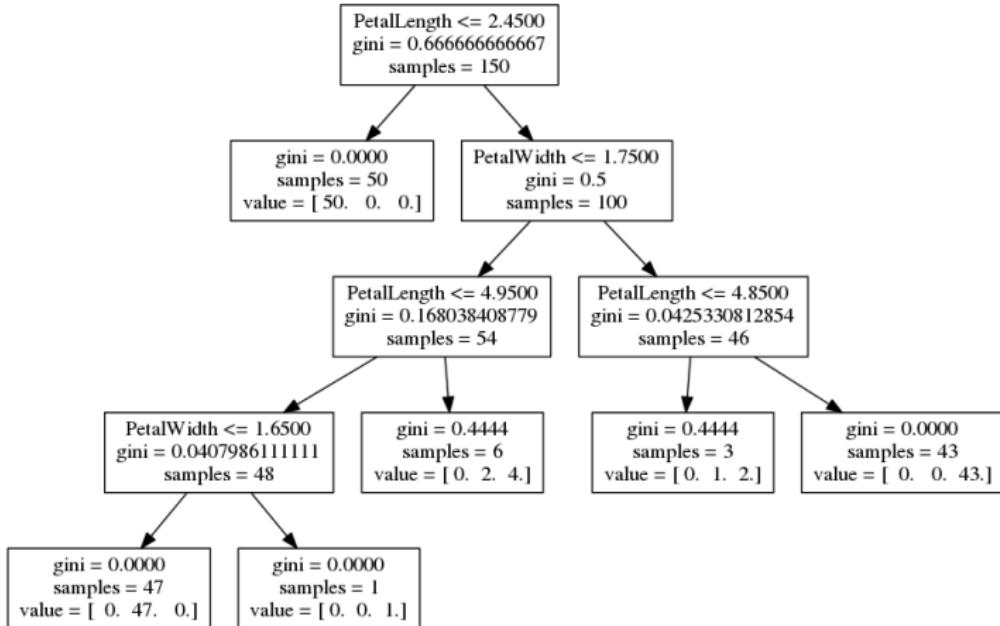
- Machine learning in Python
- Good integration with Matplotlib
- Run a learner in one line of code
- Large collection of example code snippets



Scikit-learn

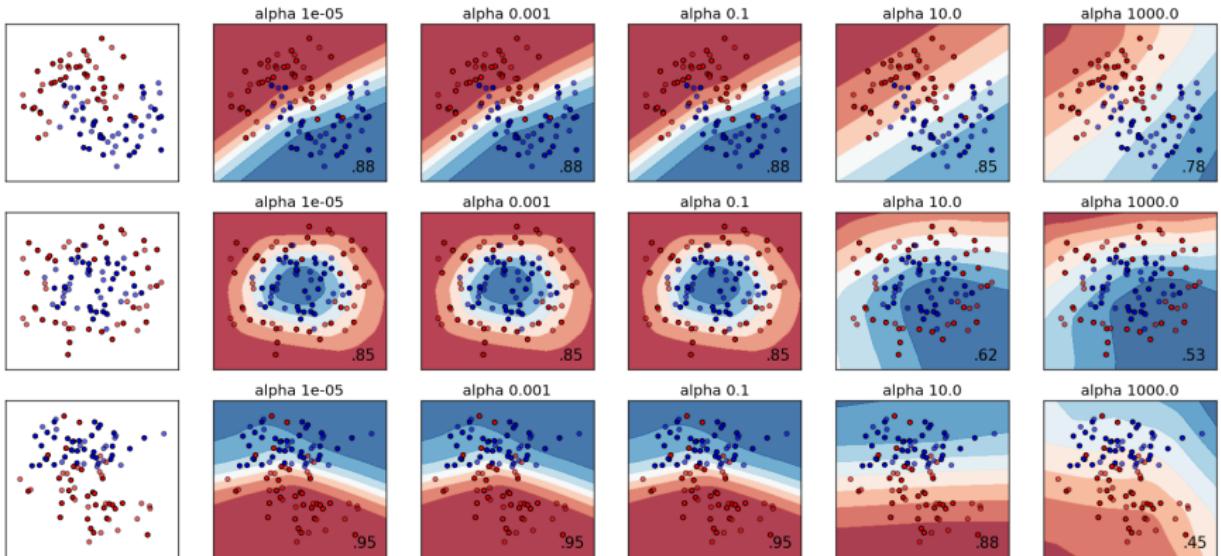
```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import numpy as np
4 from sklearn.tree import DecisionTreeClassifier
5
6 # assume proper input
7 df = pd.read_csv('filename.csv', sep='\t')
8
9 # prepare input
10 numfeatures = 8
11 features = list(df.columns[:numfeatures])
12 X = df[features]
13 y = df["target"]
14
15 # fit the model to the data
16 dt = DecisionTreeClassifier(min_samples_split=20,
17     random_state=99)
17 dt.fit(X, y)
```

Decision tree



http://chrissstrelieff.ws/sandbox/2015/06/08/decision_trees_in_python_with_scikit_learn_and_pandas.html

Neural networks α



Data mining disclaimers

- Data mining may find regularities from history, but history is inherently not the same as the future
- “Patterns are not people”
- Association (correlation) does not dictate trend nor causality
- Some abnormal data (outlier) could be caused by human error
- Interpretation is crucial
- The P-word is screaming for attention ...

Privacy



Solutions for privacy issues

- Give slightly randomized outputs
- Replace sensitive data with some anonymous ID or hash
- Multi-party hashing
- Multi-party computation
- Roberto J. Bayardo and Ramakrishnan Srikant, Technological Solutions for Protecting Privacy, *IEEE Computer* 36(9): 115–118, 2003.

Lab session March 9

- Set up your environment, see instructions in
`/vol/share/software/datascience/2016-README`
- Get familiar (again) with Python:
<https://github.com/jrjohansson/scientific-python-lectures>
- Do the tutorial 10 minutes to Pandas:
<http://pandas.pydata.org/pandas-docs/stable/10min.html>
- Replicate metrics from Assignment 1 using this toolstack
- Learn Matplotlib, for example via Lecture 4 of
 - <https://github.com/jrjohansson/scientific-python-lectures> or
 - <http://pandas.pydata.org/pandas-docs/stable/visualization.html>
- Replicate visualizations from Assignment 1 using this toolstack

Credits

Slides on neural networks based on slides from
“Kunstmatige Intelligentie” by Walter Kosters
(<https://liacs.leidenuniv.nl/~kosterswa/AI/neuraal.pdf>)