**Jakub Gałka**

jgalka@agh.edu.pl, C2-419
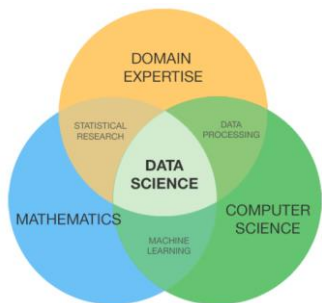office hours 2021: Monday 11:30-14:00

---

# Data Analysis and Pattern Recognition
Analiza Danych i Rozpoznawanie Wzorców

» Grades
» Laboratory classes - 50%
» Brief lecture test - 50%

---

» Data Science, Data Engineering
» Machine Learning, Deep Learning
» Artificial Intelligence
» Big Data, IoT
» Industry 4.0

---

# Jobs

Deep / Machine Learning Engineer
Data Analyst / Scientist / Engineer
Big Data Engineer

---

---

# Data Science Job

» Python, R, libraries
» Machine learning
» Statistical methods, inference, data visualisation
» Domain knowledge (eg. financial)
» SQL, noSQL
» Web, Cloud, API

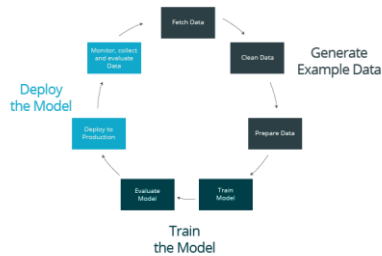## Machine Learning Job

» Python, libraries, C++
» ML: regression, classification, cluster analysis, model evaluation
» Deep learning (DNN, GAN, RL ...)
» Feature engineering, DSP
» Linear algebra, applied statistics
» Serving, Docker, API, TFX, Git
» General algorithms, data structures
» Domain knowledge (NLP, computer vision, financial)
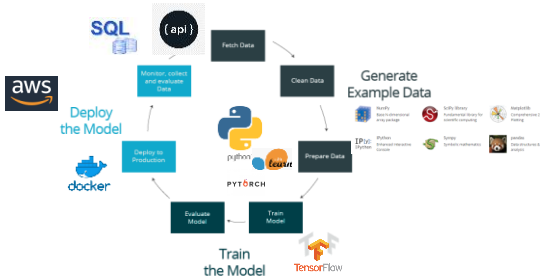
www.agh.edu.pl


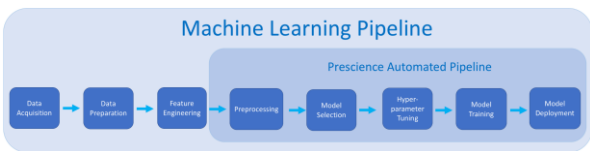Machine Learning Pipeline

www.agh.edu.pl


Machine Learning Pipeline

www.agh.edu.pl

## Tools



www.agh.edu.pl


Machine Learning Pipeline
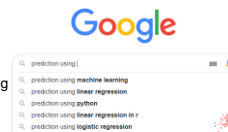
www.agh.edu.pl

## Types of ML predictions

Classification
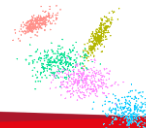  Predicts category (class, label)
  Known labelled input  x : {c1, c2, ….}

Regression
  Predicts value
  Known I/O value mapping y=f(x)
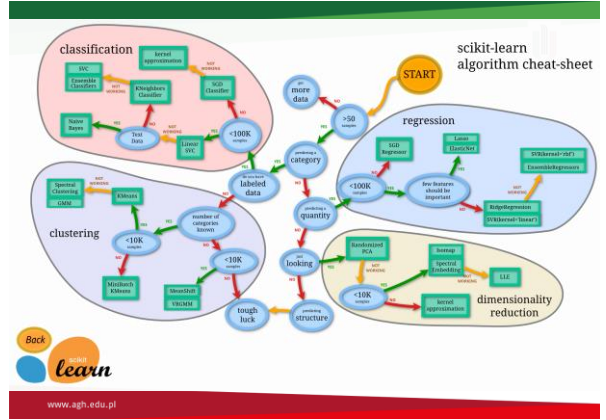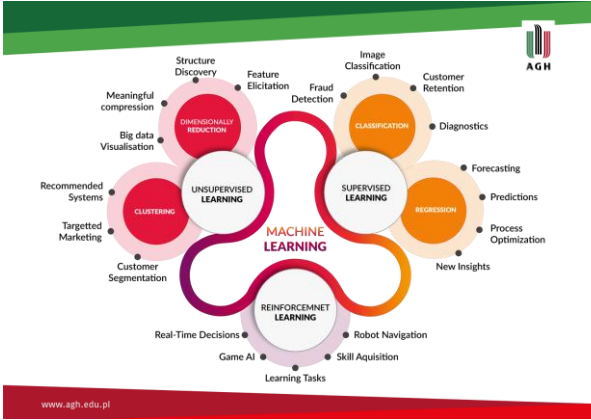
Clustering
  Groups similar elements (no labels)

www.agh.edu.pl

12/14/2021

scikit-learn algorithm cheat-sheet

www.agh.edu.pl

## Classification

koty or psy?

Predict known label for unseen data

$x$ - *feature vector*

$$x = [x_1, \ldots, x_L]^T$$

$c_i$ – *class, label* from the set of $M$ classes, $i=1,\ldots,M$

$$x \xrightarrow{?} c_i, \quad \bigcup_{i=1}^{M} c_i = \Omega$$



www.agh.edu.pl

## Bayes Classifier (MAP)

MAP (*Maximum a Posteriori prob.*): we are looking for a specific label $c_i^*$

$$c_i^* = \arg\max_{c_i, i=1,\ldots,M} P(c_i \mid x)$$

That would maximize the observed probability of the class, given the input $x$.

Binary calssification: $M=2$

Class priors: $P(c_1)$, $P(c_2)$

If we knew likelihood $p(x|c_i)$:

$$P(c_i \mid x) = \frac{p(x \mid c_i) P(c_i)}{p(x)}$$

Thomas Bayes
(1702-1761)

www.agh.edu.pl

## Binary Naive Bayes

$$p(x \mid c_1) P(c_1) \overset{?}{><} p(x \mid c_2) P(c_2)$$



www.agh.edu.pl

## Naive Bayes – 2D
## Dyscrimination Hyperplane

$$P(c_i \mid x) - P(c_j \mid x) = 0$$



www.agh.edu.pl

## Slide 1

### Decision function (Fischer discriminator)

Monotonic function of MAP

$$g_i(\boldsymbol{x}) \equiv f\big(P(c_i \mid \boldsymbol{x})\big) \qquad \forall j \neq i \quad g_i(\boldsymbol{x}) > g_j(\boldsymbol{x}) \Rightarrow \boldsymbol{x} \in c_i$$

Decision plane

$$g_{ij}(\boldsymbol{x}) \equiv g_i(\boldsymbol{x}) - g_j(\boldsymbol{x}) = 0, \quad i, j = 1, \ldots, M, \ i \neq j$$
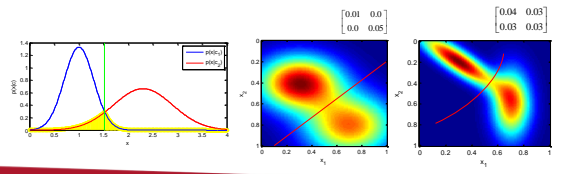


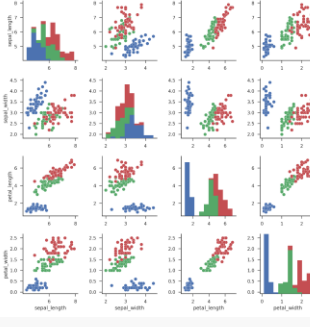How to tackle an *N*-class problem using binary classifiers only ?

www.agh.edu.pl

## Slide 2

### Naive Bayes – Normal PDF

Multivariate N-PDF
$$p(\boldsymbol{x} \mid c_i) = \frac{1}{(2\pi)^{l/2}|\Sigma_i|^{1/2}} \exp\left(-\frac{1}{2}(\boldsymbol{x} - \boldsymbol{\mu}_i)^T \Sigma_i^{-1}(\boldsymbol{x} - \boldsymbol{\mu}_i)\right)$$

$$\boldsymbol{\mu}_i = E[\boldsymbol{x}] \qquad \Sigma_i = E\big[(\boldsymbol{x} - \boldsymbol{\mu}_i)(\boldsymbol{x} - \boldsymbol{\mu}_i)^T\big] \qquad \boldsymbol{\Theta}_i = \{\boldsymbol{\mu}_i, \Sigma_i\}$$

$$\begin{bmatrix} 0.01 & 0.0 \\ 0.0 & 0.05 \end{bmatrix} \qquad \begin{bmatrix} 0.04 & 0.03 \\ 0.03 & 0.03 \end{bmatrix}$$
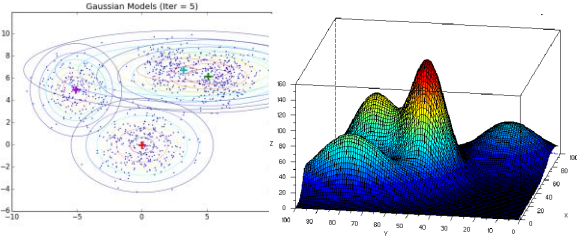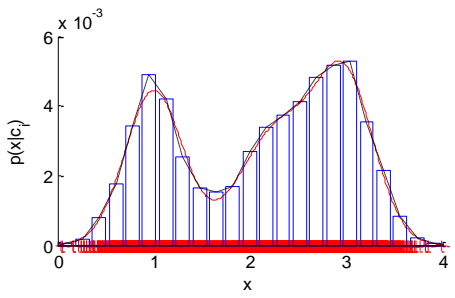


www.agh.edu.pl

## Slide 3



```
>>> from sklearn import datasets
>>> iris = datasets.load_iris()
>>> from sklearn.naive_bayes import GaussianNB
>>> gnb = GaussianNB()
>>> y_pred = gnb.fit(iris.data, iris.target).predict(iris.data)
>>> print("Number of mislabeled points out of a total %d points : %d"
... % (iris.data.shape[0],(iris.target != y_pred).sum()))
Number of mislabeled points out of a total 150 points : 6
```

## Slide 4



Gaussian Models (Iter = 5)

www.agh.edu.pl

## Slide 5



www.agh.edu.pl

## Slide 6

### Gaussian Mixture Model

Carl Friedrich Gauss
(1777-1855)

GMM – statistical model of complex distributions

$$p(\boldsymbol{x} \mid c_i) = \sum_{j=1}^{J} p(\boldsymbol{x} \mid \boldsymbol{\Theta}_{ij}) P_j$$

$$\sum_{j=1}^{J} P_j = 1, \quad \int_{x \in X} p(\boldsymbol{x} \mid \boldsymbol{\Theta}_{ij}) d\boldsymbol{x} = 1$$

$$\max\left\{ p^{ML} = \prod_k p(\boldsymbol{x}_k, \boldsymbol{\Theta}, P_1, \ldots, P_J) \right\}$$

mixture.GaussianMixture

www.agh.edu.pl

4

**Slide 1:**

```
import numpy as np
from sklearn import mixture
np.random.seed(1)
g = mixture.GaussianMixture(n_components=3)
# Generate random observations with two modes centered on 0
# and 10 to use for training.
obs = np.concatenate((np.random.randn(100, 1), 10 + np.random.randn(300, 1), 5 + np.random.randn(200, 1)))
g.fit(obs)
l=g.predict([[0], [2], [6], [9], [10]])
print(l)
p=g.predict_proba([[0], [2], [6], [9], [10]])
print(p)
```

```
[1 1 2 0 0]
[[1.74450697e-23 9.99985504e-01 1.44962481e-05]
 [2.78709519e-14 8.18856673e-01 1.81143327e-01]
 [3.62624405e-04 3.43844650e-10 9.99637375e-01]
 [9.98599654e-01 3.64032460e-22 1.40034569e-03]
 [9.99984542e-01 1.86890101e-27 1.54577593e-05]]
```
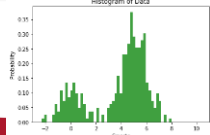
www.agh.edu.pl

**Slide 2:**

```
import numpy as np
from sklearn import mixture
np.random.seed(1)
g1 = mixture.GaussianMixture(n_components=3)
g2 = mixture.GaussianMixture(n_components=3)
# Generate random observations with two modes centered on 0
# and 10 to use for training.
obs1 = np.concatenate((np.random.randn(100, 1), 5 + np.random.randn(300, 1))) # 0 & 5
obs2 = np.concatenate((2.5 + np.random.randn(100, 1), 6 + np.random.randn(200, 1))) # 2.5 & 6
g1.fit(obs1)
g2.fit(obs2)
p1=g1.score([[4]])
print("log p(x=4|g1)=",p1)
p2=g2.score([[4]])
print("log p(x=4|g2)=",p2)
p1=g1.score([[2]])
print("log p(x=2|g1)=",p1)
p2=g2.score([[2]])
print("log p(x=2|g2)=",p2)
print("What is P(g|x)=? , Bayes")
print("P(g1|x=2)=", np.exp(p1)*0.5/(np.exp(p1)*0.5+np.exp(p2)*0.5))
print("P(g2|x=2)=", np.exp(p2)*0.5/(np.exp(p1)*0.5+np.exp(p2)*0.5))
```
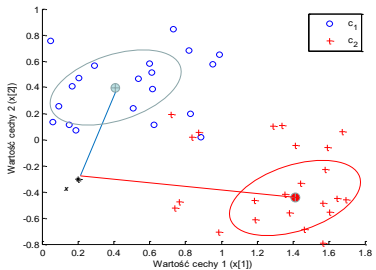
```
log p(x=4|g1)= -1.7072020753278263
log p(x=4|g2)= -2.2636967356701074
log p(x=2|g1)= -4.153287400767302
log p(x=2|g2)= -2.212724708571532
What is P(g|x)=? , Bayes
P(g1|x=2)= 0.1255860595922101
P(g2|x=2)= 0.8744139404077899
```


Histogram of Data

www.agh.edu.pl

**Slide 3:**

# Nearest neighbor



www.agh.edu.pl

**Slide 4:**

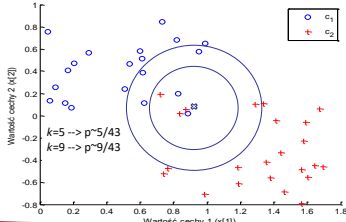# $k$NN p(x) estimator (*k-Nearest Neighbor*)

» $k$NN: what whoul be the radius (volume) of a hypersphere with k samples inside ?
» We measure the radius or volume using some metric (e.g. Euclidean)
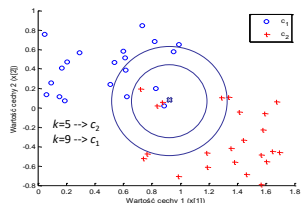
$$\tilde{p}(x) = \frac{k}{KV(x)}$$

$k$=5 --> p~5/43
$k$=9 --> p~9/43



www.agh.edu.pl

**Slide 5:**

# kNN Classifier

$$c*(x) = \arg\max_{c_i, i=1,...,M} \left( \left| \{ x_{c_i} \in V_k(x) \} \right| \right)$$

$$k \% M \neq 0$$

$k$=5 --> $c_2$
$k$=9 --> $c_1$



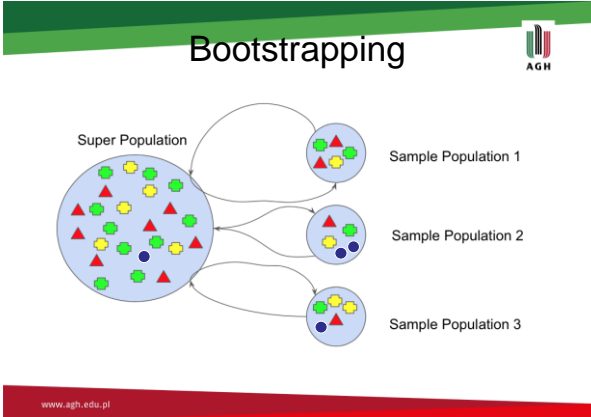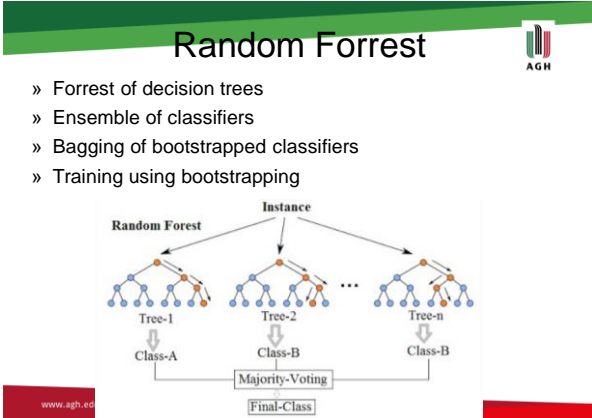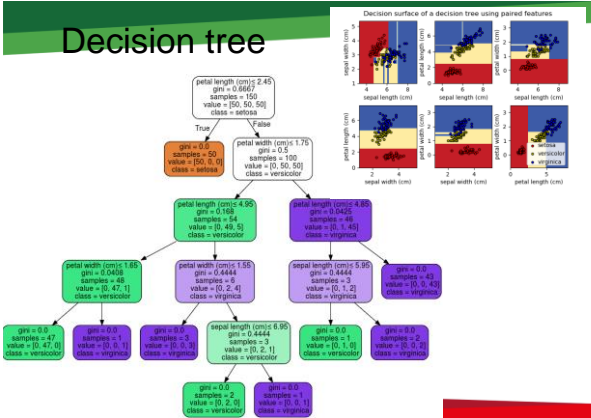www.agh.edu.pl

**Slide 6:**

# kNN – ScikitLearn

```
X = [[0], [1], [2], [3]]
y = [0, 0, 1, 1]
from sklearn.neighbors import KNeighborsClassifier
neigh = KNeighborsClassifier(n_neighbors=3)
neigh.fit(X, y)
print(neigh.predict([[1.1]]))
print(neigh.predict_proba([[1.1]]))
print(neigh.predict([[1.9]]))
print(neigh.predict_proba([[1.9]]))
print(neigh.predict([[2.1]]))
print(neigh.predict_proba([[2.1]]))
```

```
[0]
[[0.66666667 0.33333333]]
[1]
[[0.33333333 0.66666667]]
[1]
[[0.33333333 0.66666667]]
```
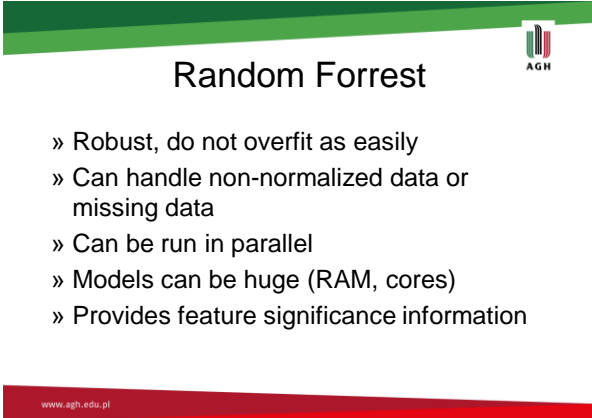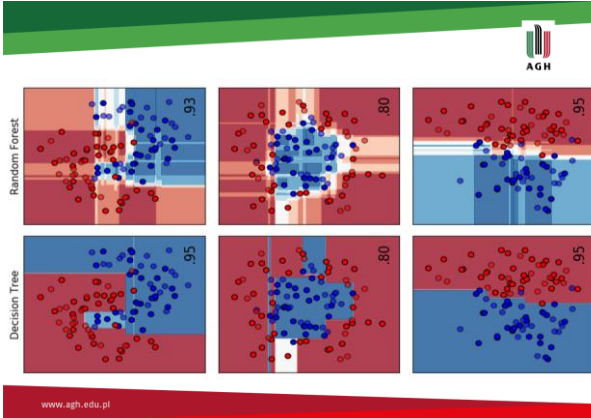
www.agh.edu.pl

## Decision tree



# Random Forrest

» Forrest of decision trees
» Ensemble of classifiers
» Bagging of bootstrapped classifiers
» Training using bootstrapping



## Bootstrapping



Super Population

Sample Population 1

Sample Population 2

Sample Population 3

```
>>> from sklearn.ensemble import RandomForestClassifier
>>> from sklearn.datasets import make_classification

>>> X, y = make_classification(n_samples=1000, n_features=4,
...                            n_informative=2, n_redundant=0,
...                            random_state=0, shuffle=False)
>>> clf = RandomForestClassifier(n_estimators=100, max_depth=2,
...                            random_state=0)
>>> clf.fit(X, y)
RandomForestClassifier(bootstrap=True, class_weight=None, criterion='gini',
          max_depth=2, max_features='auto', max_leaf_nodes=None,
          min_impurity_decrease=0.0, min_impurity_split=None,
          min_samples_leaf=1, min_samples_split=2,
          min_weight_fraction_leaf=0.0, n_estimators=100, n_jobs=None,
          oob_score=False, random_state=0, verbose=0, warm_start=False)
>>> print(clf.feature_importances_)
[0.14205973 0.76664038 0.0282433  0.06305659]
>>> print(clf.predict([[0, 0, 0, 0]]))
[1]
```



# Random Forrest

» Robust, do not overfit as easily
» Can handle non-normalized data or missing data
» Can be run in parallel
» Models can be huge (RAM, cores)
» Provides feature significance information

www.agh.edu.pl

6

## Boosting

## Boosting algorithm

## XG Boost
### eXtreme Gradient Boosting

» XGBoost library: Implementation for efficient boosting
    huge parallelization, efficiency
» Usually used with Trees or other weak (simple) classifiers
» Sparse aware (can deal with missing data)
» Block structure (for parallelization)
» Continued training (for boosting already trained models)

$ pip install xgboost

## Artificial Neural Networks

## Rosenblatt perceptron

Algebraic function which defines a discriminating hyperplane

Frank Rosenblatt
(1928-1971)

$$g(\boldsymbol{x}) = \boldsymbol{w}^T \boldsymbol{x} + w_0, \quad \boldsymbol{w} = [w(1), \ldots, w(L)]$$

$$g_{ij} = \boldsymbol{w}^T (\boldsymbol{x}_i - \boldsymbol{x}_j) = 0$$

$\boldsymbol{w}$ - hyperplane normal vector



$$g(\boldsymbol{x}) = f(\boldsymbol{w}^T \boldsymbol{x} + \boldsymbol{w}_0)$$

## Perceptron inference example

x=[2, 3], y* = 1     w=[-1, -3, 2],
Linear:
Y(x) =
Relu:

Loss linear:
L1 = | y − y* |   =
MSE, (L2) =
Loss Relu:
       L1:                    L2:

---

.

## Gradient descent example

$x=[2, 3]$, $y^* = 1$   $w=[-1, -3, 2]$, $L = 0.1$     $F_{MSE} = (y-y^*)^2$     $y(x) = ?$

$Y(x) = -1 + (-3)*2 + 2*3 = -1$     $-> F(w,x) = (-1-1)^2 = 4$

$w' = w - L * dF(w, X)/dw$

$dF/dw = d( w0 + w1*2 + w2*3 - 1 )^2/dw = ...$

/w0:

## FFN/MLP, one-hot classifier

```
>>> from sklearn.neural_network import MLPClassifier
>>> X = [[0., 0.], [1., 1.]]
>>> y = [0, 1]
>>> clf = MLPClassifier(solver='lbfgs', alpha=1e-5,
...                     hidden_layer_sizes=(5, 2), random_state=1)
...
>>> clf.fit(X, y)
MLPClassifier(activation='relu', alpha=1e-05, batch_size='auto',
        beta_1=0.9, beta_2=0.999, early_stopping=False,
        epsilon=1e-08, hidden_layer_sizes=(5, 2),
        learning_rate='constant', learning_rate_init=0.001,
        max_iter=200, momentum=0.9, n_iter_no_change=10,
        nesterovs_momentum=True, power_t=0.5, random_state=1,
        shuffle=True, solver='lbfgs', tol=0.0001,
        validation_fraction=0.1, verbose=False, warm_start=False)
>>> clf.predict([[2., 2.], [-1., -2.]])
array([1, 0])
```

```
import keras
from keras.models import Sequential
from keras.layers import Dense, Dropout, Activation
from keras.optimizers import SGD

# Generate dummy data
import numpy as np
x_train = np.random.random((1000, 20))
y_train = keras.utils.to_categorical(np.random.randint(10, size=(1000, 1)), num_classes=10)
x_test = np.random.random((100, 20))
y_test = keras.utils.to_categorical(np.random.randint(10, size=(100, 1)), num_classes=10)

model = Sequential()
# Dense(64) is a fully-connected layer with 64 hidden units.
# in the first layer, you must specify the expected input data shape:
# here, 20-dimensional vectors.
model.add(Dense(64, activation='relu', input_dim=20))
model.add(Dropout(0.5))
model.add(Dense(64, activation='relu'))
model.add(Dropout(0.5))
model.add(Dense(10, activation='softmax'))

sgd = SGD(lr=0.01, decay=1e-6, momentum=0.9, nesterov=True)
model.compile(loss='categorical_crossentropy',
              optimizer=sgd,
              metrics=['accuracy'])

model.fit(x_train, y_train,
          epochs=20,
          batch_size=128, verbose=0)
score = model.evaluate(x_test, y_test, batch_size=128)
print(score)
```
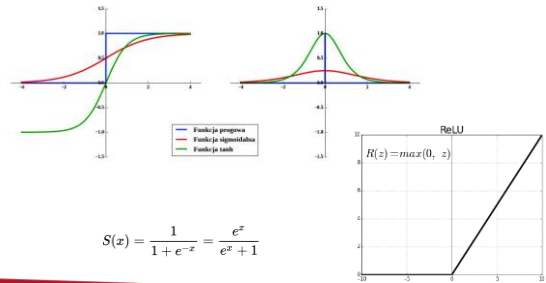```
100/100 [==============================] - 0s 704us/step
[2.2926833629608154, 0.12999999523162842]
```

## Activation functions and derivatives



$$S(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{e^x + 1}$$

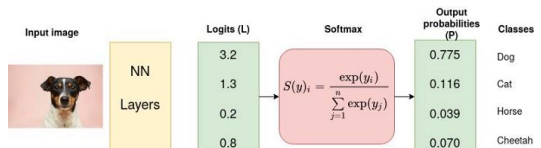$R(z) = max(0, z)$

## Softmax network outputs

- » Usually used as the network output activations
- » 1 example – one class only (no multilabel problems, one-hot approach)
- » Trained to produce probability estimation
- » Often trained under a log-loss / cross-entropy loss function

$$\sigma(\mathbf{z})_i = \frac{e^{z_i}}{\sum_{j=1}^{k} e^{z_j}}$$

## Softmax example



| Input image | Logits (L) | Softmax | Output probabilities (P) | Classes |
|---|---|---|---|---|
| NN Layers | 3.2 | $S(y)_i = \frac{\exp(y_i)}{\sum_{j=1}^{n} \exp(y_j)}$ | 0.775 | Dog |
| | 1.3 | | 0.116 | Cat |
| | 0.2 | | 0.039 | Horse |
| | 0.8 | | 0.070 | Cheetah |

9

## Slide 1 — Softmax example

$$\sigma(\mathbf{z})_i = \frac{e^{z_i}}{\sum_{j=1}^{k} e^{z_j}}$$

Softmax example

» Y1 = 0.3     y2 = 0.1

» P(y1) = e^0.3  / (e^0.3 + e^0.1) = a
» P(y2) = e^0.1 / (  ……  )  = 1 – a
» ==1

www.agh.edu.pl

## Slide 2 — Common Loss Functions

# Common Loss Functions

» L1:              |y – y^|
» L2, MSE:        (y – y^)$^2$

» Binary Cross Entropy:
    –( y^ * log(p) + (1-y^)*log(1-p) )

  Multi-Class Cross Entropy:
    L(x$_i$, y$_j$) = - sum$_{\_j}$ { y$_{ij}$ * log(p$_{ij}$) }
    One-hot encoding:    y^ = [0, 0, 1, 0, …, 0]

» Divergence (for ratio-measures):
        log ( y^/y + y/y^ )        1:1 -> log(2)

www.agh.edu.pl

## Slide 3 — Cross-Entropy example

# Cross-Entropy example

» x: y = [0.1, 0.6, 0.3]     y^ = [0, 0, 1]
» L = – (y^ * log(p) + (1-y^)*log(1-p))
» L = - ( 0*log(0.1) + 1 *log(0.9) +
        0* log(0.6) + 1*log(0.4) +
        1*log(0.3) + 0*log(0.7) )

www.agh.edu.pl

## Slide 4 — Backpropagation - assumptions

*Backpropagation* - assumptions



$$w_{ij}^{t+1} = w_{ij}^{t} + \Delta w_{ij}$$

w – weights
i – layer (i=1,…,l)
j,m – neuron, synapse (j,m=1,…,N$_i$)
k – sample (k=1,…,K)
μ – learning rate
J – loss function
y – expected outputs
δ – gradient

$$\Delta w_{ij} = -\mu \frac{\delta J}{\delta w_{ij}}$$

$$J = \frac{1}{2}\sum_{k=1}^{K}\sum_{m=1}^{M}(y_k(m) - \hat{y}_k(m))^2$$

$$\Delta w_{ij} = -\mu \sum_{k=1}^{K} \delta_{ij,k} y_{i-1,k}$$

$$\delta_{ij,k} \equiv \frac{\delta\left(\sum_{m=1}^{M}(y_k(m)-\hat{y}_k(m))^2\right)}{\delta(w_{ij,k} \cdot y_{i-1,j,k})}$$

www.agh.edu.pl

## Slide 5 — Backpropagation - solution

*Backpropagation* - solution

$$w_{ij}^{t+1} = w_{ij}^{t} - \mu \sum_{k=1}^{K} \delta_{ij,k} y_{i-1,k}$$

$$f(x) = \frac{1}{1+\exp(-\alpha x)}$$
$$f'(x) = \alpha f(x)(1-f(x))$$

Solution for J=L*MSE(y-y^)

$$i = I : \quad \delta_{i=I,j,k} = \left[f(w_{ij} y_{i-1,k}) - \hat{y}_k(j)\right] f'(w_{ij} y_{i-1,k})$$

$$i < I : \quad \delta_{i-1,j,k} = \left(\sum_{m=1}^{N_i} \delta_{im} w_{im}(j)\right) f'(w_{i-1,j} y_{i-2,k})$$

w – weights
i – layer (i=1,…,l)
j,m – neuron, synapse (j,m=1,…,N$_i$)
k – sample (k=1,…,K)
μ – learning rate
J – loss function
y – expected outputs
δ – gradient



www.agh.edu.pl

## Slide 6 — Backpropagation - implementation

*Backpropagation* - implementation

1. Inicjalizacja – wylosuj wszystkie wagi **w**$^{t=0}$
2. Wyznacz odpowiedzi y$_{i,k}$=f(**w**y$_{i-1,k}$) wszystkich neuronów sieci dla każdego wektora treningowego (*forward step*)
3. Oblicz funkcję kosztu *J* (krok pośredni)
4. *i*=*l*, Dla wszystkich wzorców (*k*) wyznacz wartość δ$_{l,jk}$ każdego neuronu wyjściowego (*backward step*)
5. Dla wszystkich warstw (kolejno wstecz) *i*=*l,l-1*,…,2 wyznacz δ$_{i-1, jk}$ wszystkich neuronów dla każdego wzorca (*k*)
6. Ustaw wszystkie wagi

$$w_{ij}^{t+1} = w_{ij}^{t} - \mu \sum_{k=1}^{K} \delta_{ij,k} y_{i-1,k}$$

7. t=t+1, *goto* 2, *until*   J<J$_{tr}$  OR  dJ/d**w**<e  OR  *t*>T

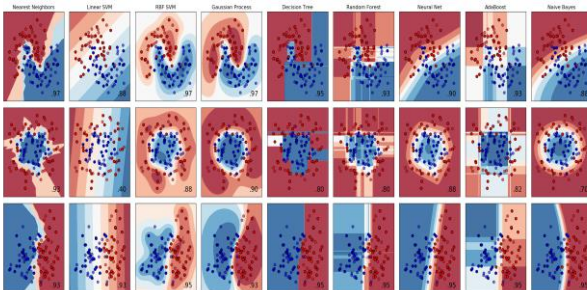www.agh.edu.pl

## Stochastic gradient descent

» GD Problems
– RAM
– Computational power
» SGD Solution
– Batch
– Batch-size
– Epoch

www.agh.edu.pl

---

» https://losslandscape.com/gallery/
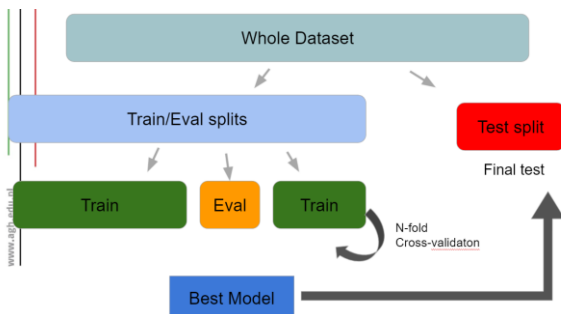


---

## Comparison of classifiers



www.agh.edu.pl

---

## Model evaluation

» How good is our model?
» Generalization, Robustness

» Loss Function
» Performance: Accuracy, Error rate (%)
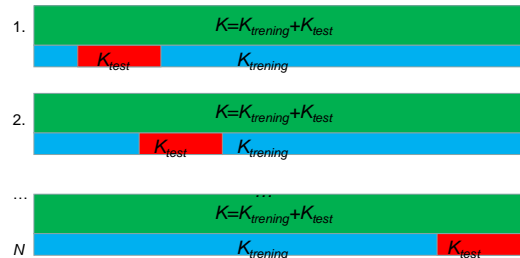» Domain-specific metric
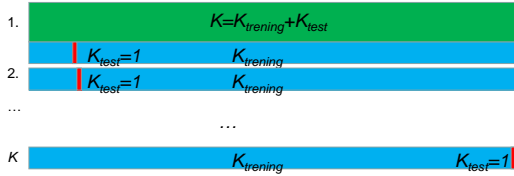
www.agh.edu.pl

---

## Model evaluation schema



---

## *N*-fold Cross-validation $ERR = \frac{1}{N}\sum_{n=1}^{N} ERR_n$



www.agh.edu.pl

---

## Leave-one-out



1. $K = K_{trening} + K_{test}$
2. $K_{test}=1$ $K_{trening}$
   $K_{test}=1$ $K_{trening}$
...
...
$K$ $K_{trening}$ $K_{test}=1$

Maximizing the count of both training and validation sets.

www.agh.edu.pl

## Model training plots

Loss, Error (%)

Iterations, epochs

www.agh.edu.pl

## Overfitting

» Train/Eval mismatch
» Data problem
» Model problem
» Optimizer problem

www.agh.edu.pl

## How to avoid overfitting

» Better Data (quantity, quality, representativeness)
» Data augumentation
» Feature engineering
» Hyper-parameter optimization (x-val)
» Model (architecture and size change)
» Regularization L1, L2
» Batch-Normalization
» Dropout
» Shorter training

www.agh.edu.pl

## Binary classification

» Two-class problem
» Often stated as thresholding problem
» How to assess detection performance and confidence
» Is the decision true or false ?
» How good is it ?



www.agh.edu.pl

## Types of binary predictions



Positive → True Positive / False Negative

Negative → True Negative / False Positive

Model Prediction

www.agh.edu.pl

## MAP probability



$P_A > P_B$

## Performance of Binary Classifiers (detectors)

» Accuracy (%) not the best choice (class quantity bias)
» TP, TN, FP, FN
» Recall (TPR, Sensitivity), Precision
» AUC, F-score

| TP | FP |
|----|----|
| FN | TN |
| 1  | 1  |

## Binary classification evaluation



|  |  | Condition (as determined by "Gold standard") | | |
|---|---|---|---|---|
|  |  | Condition Positive | Condition Negative | |
| Test Outcome | Test Outcome Positive | True Positive | False Positive (Type I error) | Positive predictive value = Σ True Positive / Σ Test Outcome Positive |
|  | Test Outcome Negative | False Negative (Type II error) | True Negative | Negative predictive value = Σ True Negative / Σ Test Outcome Negative |
|  |  | Sensitivity = Σ True Positive / Σ Condition Positive | Specificity = Σ True Negative / Σ Condition Negative | |

## Binary classification performance

» Recall, Sensitivity, True Positive Rate

Recall $= TP / (TP + FP)$

» Precision $= TP / (TP + FN)$
» Accuracy $= (TP + TN) / All$
» F-measure $= 2*Precision*Recall/(Precision+Recall)$

## ROC - *Receiver Operating Characteristic* (TPR vs FPR) AUC – Area Under Curve

## Precision-Recall Plot
for inbalanced data



www.agh.edu.pl

## Why going deep ?



www.agh.edu.pl



## Why Deep Learning ?



www.agh.edu.pl

## Why Now ? (2010+)

» More data
» More computational power
» More interest -> more people
» Better algorithms
» Better results
» More applications
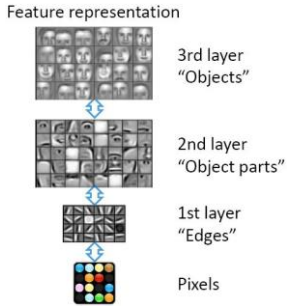
and loop…

www.agh.edu.pl

## What is Deep Learning

WiKi:
» Deep learning is a class of machine learning algorithms
» use a cascade of multiple layers of nonlinear processing units for feature extraction and transformation. Each successive layer uses the output from the previous layer as input.
» learn in supervised (e.g., classification) and/or unsupervised (e.g., pattern analysis) manners.
» learn multiple levels of representations that correspond to different levels of abstraction; the levels form a hierarchy of concepts.

www.agh.edu.pl

## Layered abstraction

Feature representation



3rd layer "Objects"

2nd layer "Object parts"

1st layer "Edges"

Pixels

www.agh.edu.pl

## How To „Deep Learning" ?



Caffe — Caffe2 — Chainer — mxnet — PaddlePaddle — PYTORCH — TensorFlow — theano — torch — WOLFRAM

www.agh.edu.pl

## Normal vs Deep Feed-Forward Network



www.agh.edu.pl

## Deep network training

» Network pretraining (1-3) and Fine-tunning (4)
» Unsupervised and supervised training (loss function)



1)  2)  3)  4)

» Regularization (L2, L1, Dropout)
» Data augumentation
» Meta-parameters optimization (Batch size, net size)
» Evaluation (Metric)

www.agh.edu.pl

```
1  from keras.models import Sequential
2  from keras.layers import Dense
3  import numpy
4  # fix random seed for reproducibility
5  numpy.random.seed(7)
6
7  dataset = numpy.loadtxt("my_dataset.csv", delimiter=",")
8  # input (X) and output (Y)
9  X = dataset[:,0:39]
10 Y = dataset[:,39]
11
12 # create model
13 # 12, 8, 4, and 1 neuron in consecutive layers
14 model = Sequential()
15 # input layer
16 model.add(Dense(12, input_dim=39, activation='relu'))
17 model.add(layers.Dropout(0.2, noise_shape=None, seed=None))
18 model.add(Dense(8, activation='relu'))
19 model.add(layers.Dropout(0.1, noise_shape=None, seed=None))
20 model.add(Dense(4, activation='relu'))
21 model.add(Dense(1, activation='sigmoid'))
22 # Output layer have one leuron - '1' for Speech, '0' for non-speech frame
23
24 # Compile model
25 model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
26
27 # Fit the model
28 model.fit(X, Y, epochs=100, batch_size=20)
29
30 # evaluate the model
31 scores = model.evaluate(X, Y)
32 print("\n%s: %.2f%%" % (model.metrics_names[1], scores[1]*100))
33
34 # predict output
35 predictions = model.predict(X)
```

## Recurrent Neural Networks



Unfold

www.agh.edu.pl

15

## Recurrent Neural Networks

BAIDU (12.2014) Deep Speech: Scaling up end-to-end ASR



Softmax output

Recurrent
Recurrent

3 clipped-ReLU FF

FFT spectrogram input (5000 hrs)

## RNN vs LSTM

» Vanishing gradient problem
» RNN acts as a very deep FF



www.agh.edu.pl

## Standard RNN unwrapped



www.agh.edu.pl

## Long Short-Term Memory Networks (1997)



www.agh.edu.pl

## Steps of LSTM operation
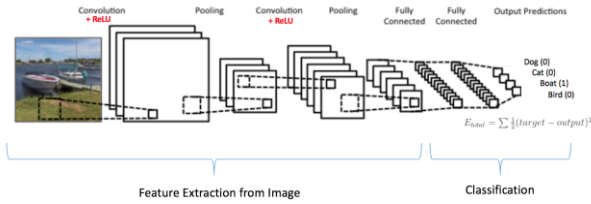
A Propagate $C_{t-1}$?
B Update $C_t$?
C Update $C_t$!
D Prepare $h_t$



## Convolutional Neural Networks

» In fully conected FF networks number of weights is to big for huge inputs.
» Same features can be observed in different places of the same signal (as in images)
» We can **group neurons** to analyse different parts of data –> lower number of weights

## CNN workflow



## Convolutional Neural Networks



## Max-pooling, Average-pooling



## Elementary CNN Subsampling



## Exemple of typical CNN-FC model



## CNN classification

## CNN concepts

AGH

» Kernel / filter size,
» Number of kernels
» Stride – kernel shift (usually 1)
» Padding
» Pooling – reduces dimensionality
     Max-pooling, Average-pooling
» Flattening

www.agh.edu.pl

---

## CNN in Keras

AGH

```
# convert class vectors to binary class matrices
y_train = keras.utils.to_categorical(y_train, num_classes)
y_test = keras.utils.to_categorical(y_test, num_classes)

model = Sequential()
model.add(Conv2D(32, kernel_size=(3, 3),
                 activation='relu',
                 input_shape=input_shape))
model.add(Conv2D(64, (3, 3), activation='relu'))
model.add(MaxPooling2D(pool_size=(2, 2)))
model.add(Dropout(0.25))
model.add(Flatten())
model.add(Dense(128, activation='relu'))
model.add(Dropout(0.5))
model.add(Dense(num_classes, activation='softmax'))

model.compile(loss=keras.losses.categorical_crossentropy,
              optimizer=keras.optimizers.Adadelta(),
              metrics=['accuracy'])

model.fit(x_train, y_train,
          batch_size=batch_size,
          epochs=epochs,
          verbose=1,
          validation_data=(x_test, y_test))
score = model.evaluate(x_test, y_test, verbose=0)
print('Test loss:', score[0])
print('Test accuracy:', score[1])
```
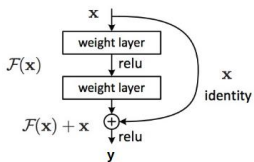
```
from __future__ import print_function
import keras
from keras.datasets import mnist
from keras.models import Sequential
from keras.layers import Dense, Dropout, Flatten
from keras.layers import Conv2D, MaxPooling2D
from keras import backend as K

batch_size = 128
num_classes = 10
epochs = 12

# input image dimensions
img_rows, img_cols = 28, 28

# the data, split between train and test sets
(x_train, y_train), (x_test, y_test) = mnist.load_data()

if K.image_data_format() == 'channels_first':
    x_train = x_train.reshape(x_train.shape[0], 1, img_rows, img_cols)
    x_test = x_test.reshape(x_test.shape[0], 1, img_rows, img_cols)
    input_shape = (1, img_rows, img_cols)
else:
    x_train = x_train.reshape(x_train.shape[0], img_rows, img_cols, 1)
    x_test = x_test.reshape(x_test.shape[0], img_rows, img_cols, 1)
    input_shape = (img_rows, img_cols, 1)

x_train = x_train.astype('float32')
x_test = x_test.astype('float32')
x_train /= 255
x_test /= 255
print('x_train shape:', x_train.shape)
print(x_train.shape[0], 'train samples')
print(x_test.shape[0], 'test samples')
```

www.agh.edu.pl

---

## Deep CNN

AGH



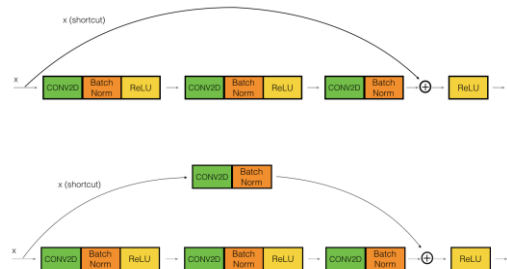www.agh.edu.pl

---

## ResNet - Residual Networks

AGH



www.agh.edu.pl

---

## ResNet - Residual Networks

AGH



$$y = x + F(x)$$

$$\frac{\delta E}{\delta x} = \frac{\delta E}{\delta y} * \frac{\delta y}{\delta x} = \frac{\delta E}{\delta y} * (1 + F'(x))$$

$$= \frac{\delta E}{\delta y} + \frac{\delta E}{\delta y} * F'(x)$$

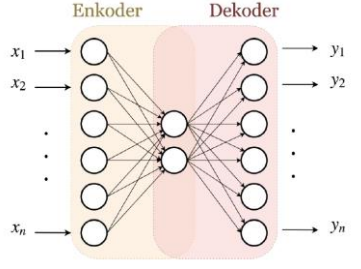www.agh.edu.pl

---

AGH

X_shortcut = X     # Store the initial value of X in a variable

## Here perform convolution + batch norm operations on X

X = Add()([X, X_shortcut])   # SKIP Connection
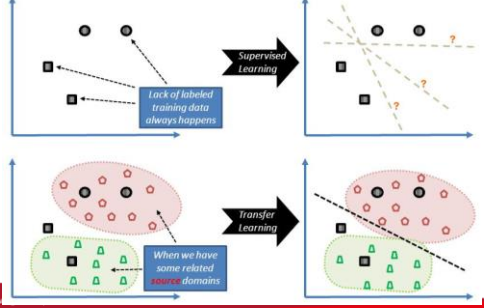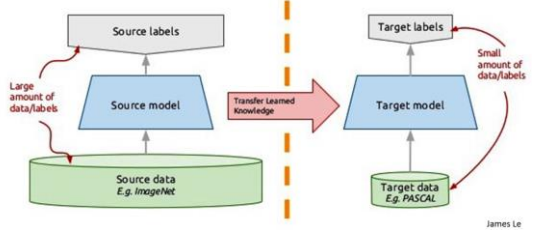


www.agh.edu.pl

## Autoencoders



## Transfer Learning



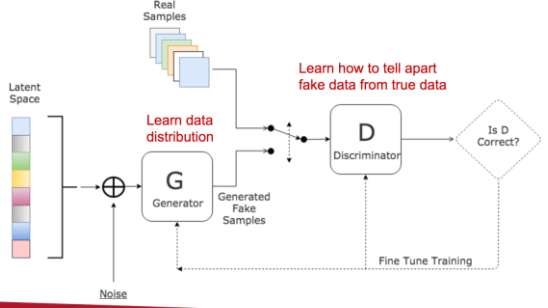## Transfer learning: idea



## Generative Adversarial Networks



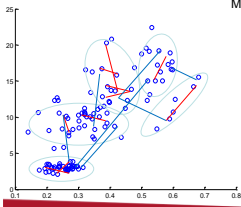## Cluster analysis



## Cluster analysis

## Analiza skupień

» Przyporządkowanie obiektów (wektorów cech) do $K$ zbiorów $X_{k=1...K}$, w których obiekty z danego zbioru $X_k$ są do siebie podobne, a między zbiorami różne.
» Kryterium podobieństwa: miara odległości (np. metryka euklidesowa)
» Kryterium jakości przyporządkowania

Minimalizacja zmienności w zbiorze:

$$q = \min_{\{X_k\}} \sum_{k=1}^{K} \sum_{\substack{x_i, x_j \in X_k \\ x_i \neq x_j}} \delta(x_i, x_j)$$

Maksymalizacja zmienności między zbiorami:

$$r = \max_{\{X_k\}} \sum_{k=1}^{K} \sum_{\substack{x_i \in X_k \\ x_j \notin X_k}} \delta(x_i, x_j)$$

www.agh.edu.pl

---

## Analiza skupień - zastosowania

» *Data-mining*, wykrywanie prawidłowości w zbiorach danych
» Generowanie hipotez na podstawie danych
» Weryfikacja hipotez na podstawie danych
» Redukcja wymiarowości, kwantyzacja
» Kompresja, kodowanie
» Modelowanie

www.agh.edu.pl

---

## Metody hierarchiczne
### *Linkage*

» Metody iteracyjne
» Proste w implementacji
» Dość szybkie
» Algorytm *Linkage*:
  – Łączenie (aglomeracja) skupień otrzymanych w poprzednim kroku działania algorytmu.
  – Łączenie na podstawie podobieństwa minimalnego między obiektami oraz zbiorami. (*x~X*, *x~x*, *X~X*)
  – Iteracyjne powtarzanie łączenia
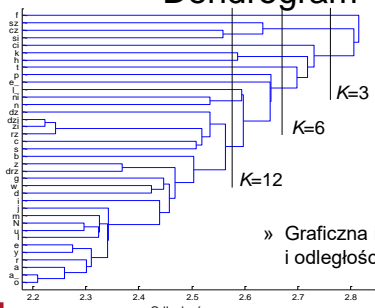
www.agh.edu.pl

---

## Algorytm *Linkage*

1. Inicjalizacja: każdy wektor $x_i$ jest 1-elementowym zbiorem $X_i$ (skupieniem)
2. Znajdź najmniejszą **odległość między** wszystkimi **zbiorami**
   $$d_{min} = \min(\ d(X_i, X_j)\ ),\ i \neq j$$
3. Połącz te zbiory, dla których znaleziono $d_{min}$
4. Powtarzaj 2) i 3) aż otrzymasz 1 zbiór (skupienie) lub oczekiwaną liczbę skupień (*K*)

K=6       ...       K=2       K=1

www.agh.edu.pl

---

## Dendrogram

*K*=3

*K*=6

*K*=12

» Graficzna ilustracja skupień i odległości między zbiorami

Odległość

www.agh.edu.pl
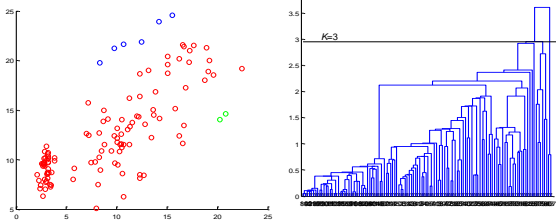
---

## Definicje odległości dla zbiorów

» *Single Linkage*
$$d_{\text{single}} = \min_{\substack{x \in X_i \\ y \in X_j}} \delta(x, y)$$

» *Complete Linkage*
$$d_{\text{complete}} = \max_{\substack{x \in X_i \\ y \in X_j}} \delta(x, y)$$

» *Average Linkage*
$$d_{\text{average}} = \frac{1}{|X_i \| X_j|} \sum_{x \in X_i} \sum_{y \in X_j} \delta(x, y)$$
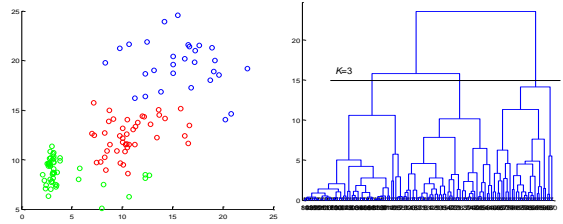
www.agh.edu.pl

## Single linkage
*przykład*

$$d_{single} = \min_{\substack{x \in X_i \\ y \in X_j}} \delta(x, y)$$
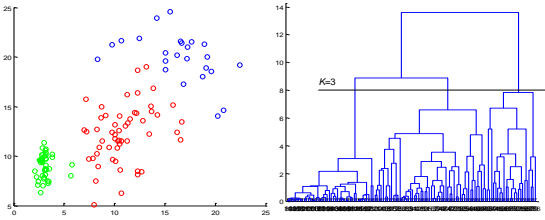


www.agh.edu.pl

## Complete linkage
przykład

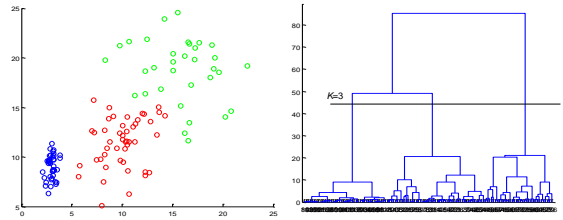$$d_{complete} = \max_{\substack{x \in X_i \\ y \in X_j}} \delta(x, y)$$



www.agh.edu.pl

## Average linkage
przykład

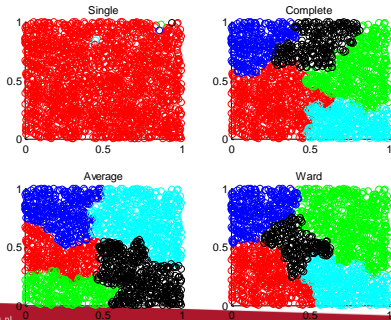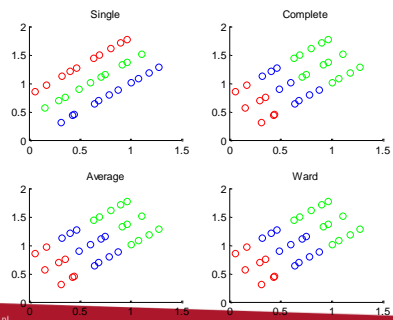$$d_{average} = \frac{1}{|X_i||X_j|} \sum \sum \delta(x, y)$$



www.agh.edu.pl

## Ward's linkage
przykład


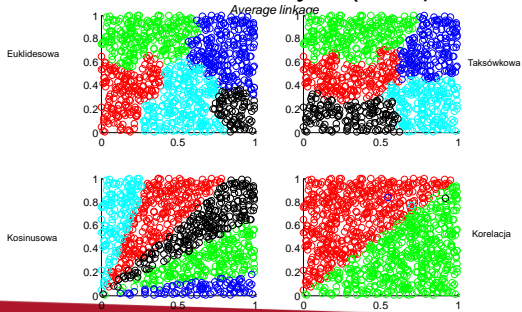
www.agh.edu.pl

## Analiza wektorów losowych
*Linkage K*=5



www.agh.edu.pl

## Algorytm *linkage*
przypadek szczególny, *K*=3



www.agh.edu.pl

## Różne metryki (*K*=5)



## *K-Means* założenia

» Bardzo popularny
» Występuje w wielu odmianach
» Nie hierarchiczny
» Zakładamy z góry *K* - oczekiwaną liczbę skupisk
» Funkcja celu to minimalizacja rozmiarów skupień:

$$q = \min_{\{X_k\}} \sum_{k=1}^{K} \sum_{\substack{x_i, x_j \in X_k \\ x_i \neq x_j}} \delta(x_i, x_j)$$

» Zawsze zbieżny,
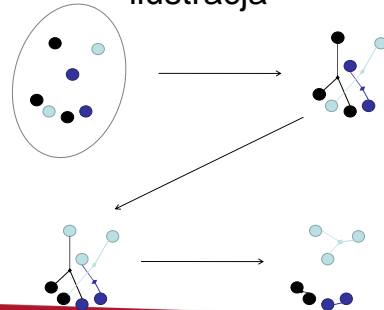» ale nie zawsze do globalnego optimum.

## *K-Means*

1. Wybierz wartość *K*.
2. Przypisz (losowo) wszystkie wektory $x_i$ do *K* skupień.
3. Wyznacz średnie $\mu_k$ wszystkich *K* skupień (*K-Means*).
4. Dla wszystkich wektorów $x_i$ znajdź najbliższe średnie $\mu_k$.

$$\mu^*(x_i) = \arg\min_{\mu_k} \delta(x_i, \mu_k)$$

5. Przypisz wszystkie wektory $x_i$ do skupień o najbliższej średniej.
6. Powtarzaj kroki 3–5 aż do ustabilizowania się pozycji średnich lub braku zmian w obsadzeniu skupień.

## Algorytm *K-Means* ilustracja



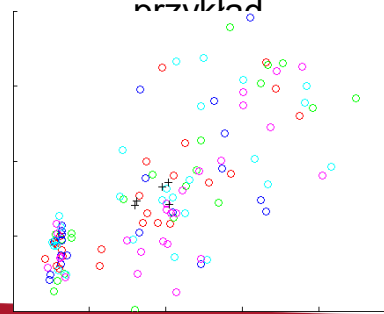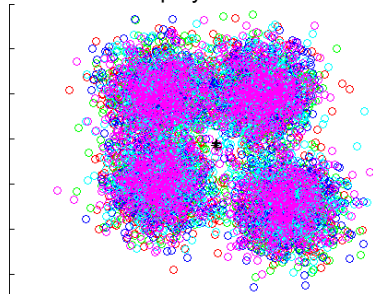## *K-Means* – przykład *K*=5



## *K-Means* przykład

*K-Means* przykład – 5 klas
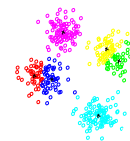


*K-Means* przykład



*K-Means* - problem
Co tu zrobić ?

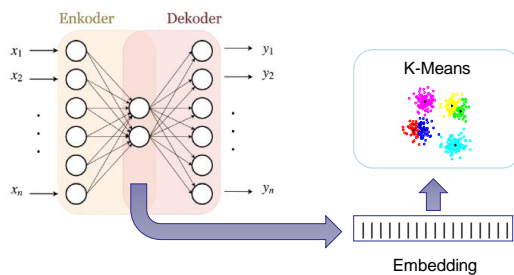## Podsumowanie

» Konieczność wyboru miary podobieństwa
» Ustalone a-prioti *K* - liczba skupień
» Zawsze zbieżne, ale
» Nie zawsze zbieżne do optimum globalnego
» Inicjowanie rozkładów początkowych



## Neural clustering w/embeddings



## Other methods

» K-Means, MiniBatch K-Means
» MeanShift, VBGMM
» Spectral: PCA, k-PCA
» Agglomerative, hierarchical
» Statistical, GMM
» Neural, embeddings

## Clustering evaluation

» Ground truth (%)
» Loss
» Adjusted Rand Index
» Mutual information (NMI, AMI)
» Completnes, Homogenity, v-measure
» Silhouette cefficients
» Contingency matrix
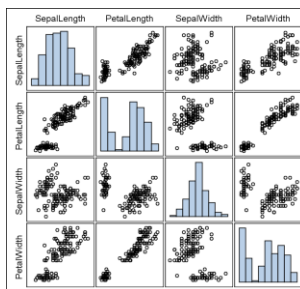» Business metrics

www.agh.edu.pl

## Feature Engineering

» Using domain knowledge of the data to create features that make machine learning algorithms work.
» Why ?
    Beter features = better results
    ML requires numerical inputs
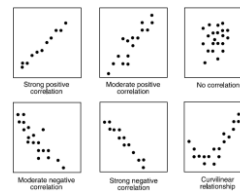» How ?

www.agh.edu.pl

## Data Exploration

» Scatter plots
» Histograms
» Distribution type
    Modality
» Ranges
» Correlation
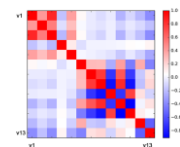» Quantity
» Understanding



www.agh.edu.pl

## Feature correlation/covariance



» Pearson correlation coefficient

$$r_{xy} = \frac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^{n}(x_i - \bar{x})^2}\sqrt{\sum_{i=1}^{n}(y_i - \bar{y})^2}}$$
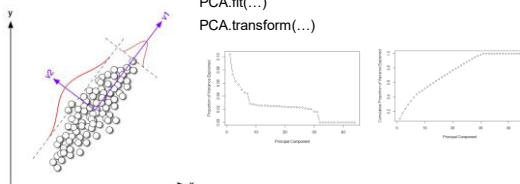
» Correlation/Covariance matrix



www.agh.edu.pl

## Correlation reduction

» PCA – Principal Component Analysis

>>> from sklearn.decomposition import PCA
PCA.fit(…)
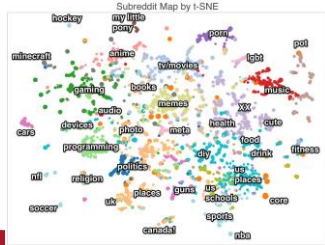PCA.transform(…)



www.agh.edu.pl

## Dimensionality reduction

» Feature Selection (quality, variance)
» Missing Values
» Feature Removal
(statistical tests – eg.Chi2, low variance, high x-corr)
» Feature space transformations
    PCA dimensionality reduction

www.agh.edu.pl

## t-Distributed Stochastic Neighbor Embedding (t-SNE )

» Great 2D/3D visualisation tool
» Dimensionality reduction
» Local similarities
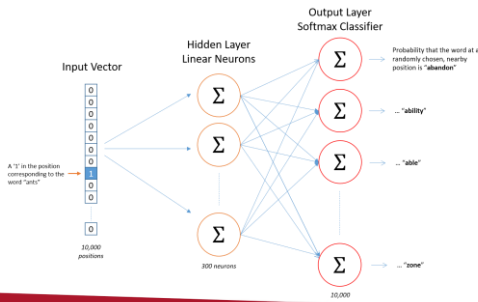» No global relations.



Subreddit Map by t-SNE

www.agh.edu.pl

## One-hot encoding

» Categorical data encoding
» Non-numerical data encoding
» No order of data categories
» High dimensionality

>>> sklearn.preprocessing.OneHotEncoder

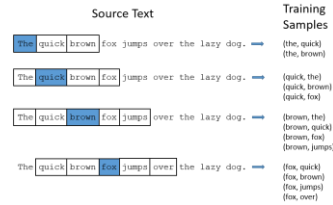| Cat | Dog | Cow | Frog | Fish | Bird | Bee |
|-----|-----|-----|------|------|------|-----|
| 0 | 0 | 0 | 0 | 1 | 0 | 0 |

www.agh.edu.pl

## One-Hot as NN input



www.agh.edu.pl

## Word2vec algorithm

» Train your network using hot-word representation
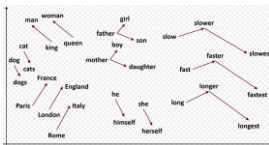» Use skip-gram method, or continuous bag of words
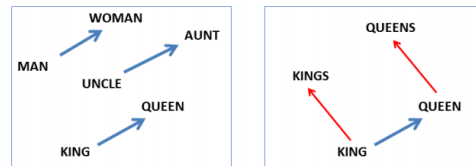


» Use bottleneck feature as embeddings

www.agh.edu.pl

## Word2Vec – word embeddings (Google, Mikolov)

» Converting text into high-dimensional vector of numbers (latent space embedding, 300-500 dim)
» Preserves linguistic and pragmatic information
» Embeddings are easy to manipulate and use in ML algorithms



www.agh.edu.pl
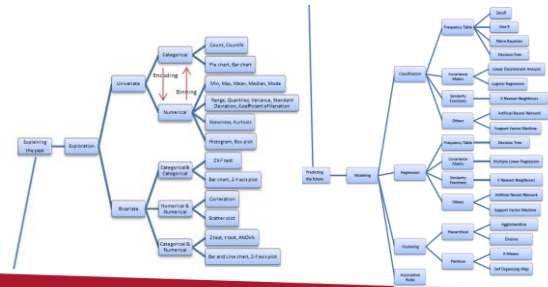
## Word embeddings - latent space



www.agh.edu.pl

## Typical data features

» Sound
» Video
» Image
» NLP
» General time series
» Frames, regresion models, histograms
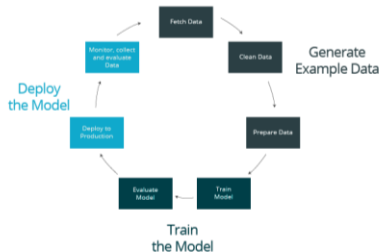
www.agh.edu.pl

## Data science cheatsheet

http://chem-eng.utoronto.ca/~datamining/dmc/data_mining_map.htm
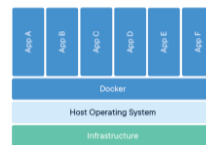


www.agh.edu.pl

## Deploying ML models



www.agh.edu.pl

## Model Deployment

» API for predictor (I/O)

» REST API (JSON)
» WebService (HTTP)
» Docker containers
» Cloud
  AWS, Azure, Google
» Big Data



www.agh.edu.pl

## TF Serving

» Scalability,
» Maintanance
» CI/CD
» TF Serving
» TF Lite
    Mobile, Embedded

www.agh.edu.pl

## TF Lite

» Embedded, mobile – runtime library with APIs and hardware support

1. Pick model
2. Optimize (quantization, etc.)
3. Convert to *Lite* format (FlatBuffer)
4. Deploy using TF library

www.agh.edu.pl