

Widzenie Komputerowe - Klasyfikacja

Wykład 11.

Magdalena Mazur-Milecka

Katedra Inżynierii Biomedycznej, WETI, PG

11 kwietnia 2019

- 1 Wprowadzenie do klasyfikacji
- 2 kNN - k Nearest Neighbours
- 3 Regresja logistyczna
- 4 Naive Bayes
- 5 SVM
- 6 Drzewa Decyzyjne
- 7 Random Forest

Klasyfikacja - realny cel rozpoznania, widzenia komputerowego.
Oszacowanie prawdopodobieństwa tego, że obraz analizowany jest obiektem, o którym posiadamy informacje w systemie.

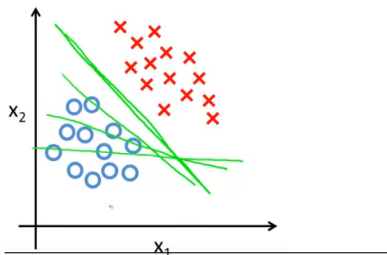
Dąży się do osiągnięcia 100% pewności co do zaklasyfikowania analizowanego obiektu.

Czym jest klasyfikacja w odniesieniu do umiejętności człowieka?

Podział klasyfikacji obrazów:

Podział klasyfikacji obrazów:

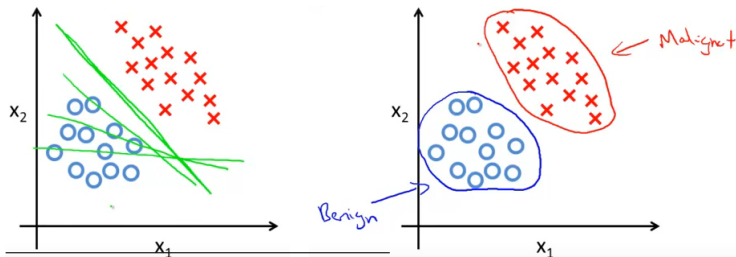
- podejście całościowe - wzorce i obiekt reprezentowane są przez wektor cech. Wektor cech umieszcza obiekt w przestrzeni (o wymiarze takim, jak ilość cech). Przestrzeń cech jest podzielona na obszary decyzyjne. Przykładem jest regresja logistyczna.



Stanford OpenClassroom, Andrew Ng

Podział klasyfikacji obrazów:

- podejście całościowe - wzorce i obiekt reprezentowane są przez wektor cech. Wektor cech umieszcza obiekt w przestrzeni (o wymiarze takim, jak ilość cech). Przestrzeń cech jest podzielona na obszary decyzyjne. Przykładem jest regresja logistyczna.
- podejście strukturalne (Generative) - wzorce reprezentowane są przez elementy bazowe oraz ich relacje, budowany jest model klas. Przykładem jest Naive Bayes.



Podział klasyfikacji obrazów:

- podejście całościowe - wzorce i obiekt reprezentowane są przez wektor cech. Wektor cech umieszcza obiekt w przestrzeni (o wymiarze takim, jak ilość cech). Przestrzeń cech jest podzielona na obszary decyzyjne. Przykładem jest regresja logistyczna.
- podejście strukturalne (Generative) - wzorce reprezentowane są przez elementy bazowe oraz ich relacje, budowany jest model klas. Przykładem jest Naive Bayes.

Podejście strukturalne jest bardziej złożone, ale pozwala na analizę bardziej złożonych obrazów

Miary oceny klasyfikatorów

- Czułość (sensitivity) = $\frac{TP}{TP+FN}$

Miary oceny klasyfikatorów

- Czułość (sensitivity) = $\frac{TP}{TP+FN}$
- Specyficzność (specificity, precision) = $\frac{TN}{TN+FP}$

Miary oceny klasyfikatorów

- Czułość (sensitivity) = $\frac{TP}{TP+FN}$
- Specyficzność (specificity, precision) = $\frac{TN}{TN+FP}$



- Czułość (sensitivity) = $\frac{TP}{TP+FN}$
- Specyficzność (specificity, precision) = $\frac{TN}{TN+FP}$
- Ilość pozytywnej predykcji = $\frac{TP}{TP+FP}$

- Czułość (sensitivity) = $\frac{TP}{TP+FN}$
- Specyficzność (specificity, precision) = $\frac{TN}{TN+FP}$
- Ilość pozytywnej predykcji = $\frac{TP}{TP+FP}$
- Ilość negatywnej predykcji = $\frac{TN}{TN+FN}$

- Czułość (sensitivity) = $\frac{TP}{TP+FN}$
- Specyficzność (specificity, precision) = $\frac{TN}{TN+FP}$
- Ilość pozytywnej predykcji = $\frac{TP}{TP+FP}$
- Ilość negatywnej predykcji = $\frac{TN}{TN+FN}$
- Dokładność (accuracy) = $\frac{TP+TN}{TP+TN+FP+FN}$

- ① **Leve-one-out** - w przypadkach małej ilości przypadków. Jeżeli liczba przypadków = N , trenowanie odbywa się na $N-1$ przypadkach, testowanie na 1, następnie wybierany jest kolejny pojedynczy przypadek i uczenie jest powtarzane. Na końcu wynik zostaje uśredniony.

- 1 **Leve-one-out** - w przypadkach małej ilości przypadków. Jeżeli liczba przypadków = N , trenowanie odbywa się na $N-1$ przypadkach, testowanie na 1, następnie wybierany jest kolejny pojedynczy przypadek i uczenie jest powtarzane. Na końcu wynik zostaje uśredniony.
- 2 **Powtarzalne losowanie podpróby** - w przypadkach większej ilości przypadków. Losowanie zbioru testowego i uczącego z zadany parametrem jest powtarzane n -krotnie. Na końcu wynik zostaje uśredniony.

- 1 Wprowadzenie do klasyfikacji
- 2 **kNN - k Nearest Neighbours**
- 3 Regresja logistyczna
- 4 Naive Bayes
- 5 SVM
- 6 Drzewa Decyzyjne
- 7 Random Forest

kNN - k najbliższych sąsiadów

Założenia:

- Dany jest zbiór uczący(?) zawierający obserwacje, które składają się z wektora cech (X_1, \dots, X_n) oraz wartość zmiennej objaśnianej Y ,

kNN - k najbliższych sąsiadów

Założenia:

- Dany jest zbiór uczący(?) zawierający obserwacje, które składają się z wektora cech (X_1, \dots, X_n) oraz wartość zmiennej objaśnianej Y ,
- Dana jest obserwacja C opisana wektorem cech, dla której chcemy prognozować wartość Y

kNN - k najbliższych sąsiadów

Założenia:

- Dany jest zbiór uczący(?) zawierający obserwacje, które składają się z wektora cech (X_1, \dots, X_n) oraz wartość zmiennej objaśnianej Y ,
- Dana jest obserwacja C opisana wektorem cech, dla której chcemy prognozować wartość Y

Algorytm polega na:

kNN - k najbliższych sąsiadów

Założenia:

- Dany jest zbiór uczący(?) zawierający obserwacje, które składają się z wektora cech (X_1, \dots, X_n) oraz wartość zmiennej objaśnianej Y ,
- Dana jest obserwacja C opisana wektorem cech, dla której chcemy prognozować wartość Y

Algorytm polega na:

- 1 porównaniu wartości cech dla obserwacji C z wartościami cech dla obserwacji ze zbioru uczącego (usytuowanie obserwacji C w przestrzeni cech uczących),

Założenia:

- Dany jest zbiór uczący(?) zawierający obserwacje, które składają się z wektora cech (X_1, \dots, X_n) oraz wartość zmiennej objaśnianej Y ,
- Dana jest obserwacja C opisana wektorem cech, dla której chcemy prognozować wartość Y

Algorytm polega na:

- 1 porównaniu wartości cech dla obserwacji C z wartościami cech dla obserwacji ze zbioru uczącego (usytuowanie obserwacji C w przestrzeni cech uczących),
- 2 wyborze parametru k , który powinien być dobierany odpowiednio do ilości klas,

Założenia:

- Dany jest zbiór uczący(?) zawierający obserwacje, które składają się z wektora cech (X_1, \dots, X_n) oraz wartość zmiennej objaśnianej Y ,
- Dana jest obserwacja C opisana wektorem cech, dla której chcemy prognozować wartość Y

Algorytm polega na:

- 1 porównaniu wartości cech dla obserwacji C z wartościami cech dla obserwacji ze zbioru uczącego (usytuowanie obserwacji C w przestrzeni cech uczących),
- 2 wyborze parametru k , który powinien być dobierany odpowiednio do ilości klas,
- 3 sprawdzeniu wartości Y dla k najbliższych obserwacji ze zbioru uczącego,

kNN - k najbliższych sąsiadów

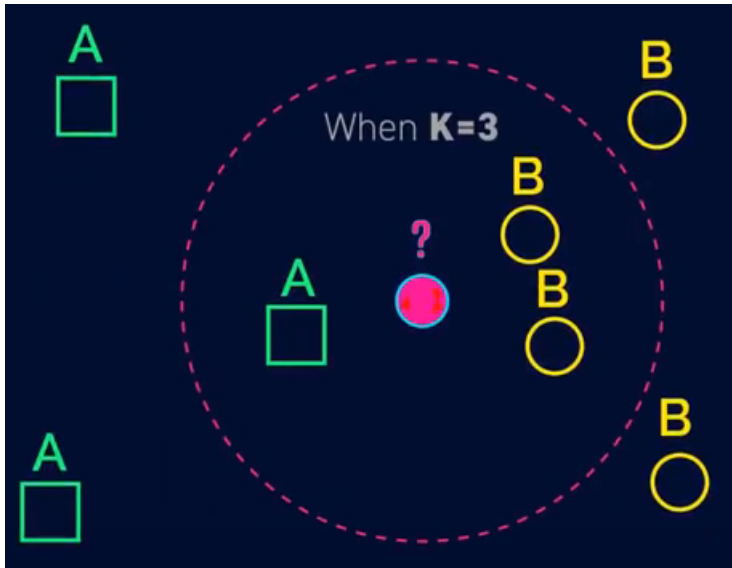
Założenia:

- Dany jest zbiór uczący(?) zawierający obserwacje, które składają się z wektora cech (X_1, \dots, X_n) oraz wartość zmiennej objaśnianej Y ,
- Dana jest obserwacja C opisana wektorem cech, dla której chcemy prognozować wartość Y

Algorytm polega na:

- 1 porównaniu wartości cech dla obserwacji C z wartościami cech dla obserwacji ze zbioru uczącego (usytuowanie obserwacji C w przestrzeni cech uczących),
- 2 wyborze parametru k , który powinien być dobierany odpowiednio do ilości klas,
- 3 sprawdzeniu wartości Y dla k najbliższych obserwacji ze zbioru uczącego,
- 4 uśrednieniu (lub wybranie większości) wartości Y dla wybranych obserwacji, co daje prognozę.

Wikipedia



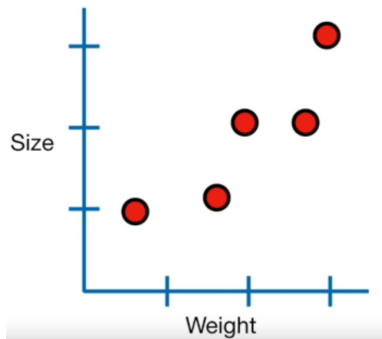
- Uczenie nie jest skomplikowane,
- Predykcja jest czasochłonna,
- Algorytm jest prosty,
- Zwiększenie k powoduje zmniejszenie wydajności,
- Zmniejszenie k powoduje zwiększenie podatności na błędy (zmniejszenie dokładności),
- Algorytm nie zawiera modelu danych.

- 1 Wprowadzenie do klasyfikacji
- 2 kNN - k Nearest Neighbours
- 3 **Regresja logistyczna**
- 4 Naive Bayes
- 5 SVM
- 6 Drzewa Decyzyjne
- 7 Random Forest

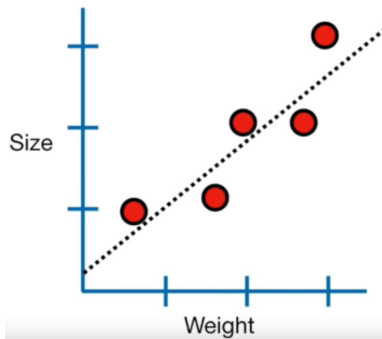
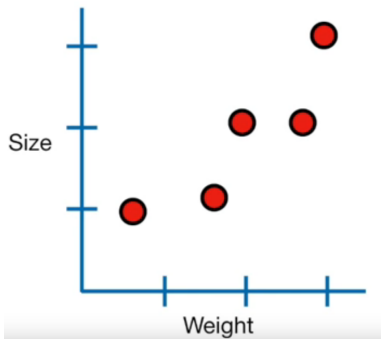
Jest to szczególny przypadek uogólnionej regresji liniowej, gdy zmienna zależna jest dychotomiczna - przyjmuje tylko dwie wartości (sukces, porażka)

$$y \in 0, 1$$

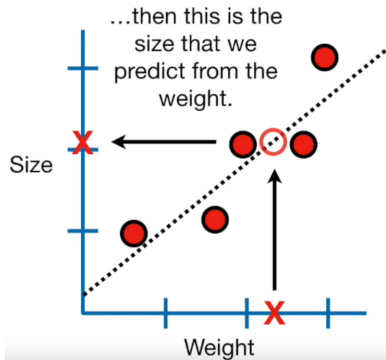
Regresja liniowa



Regresja liniowa



Regresja liniowa

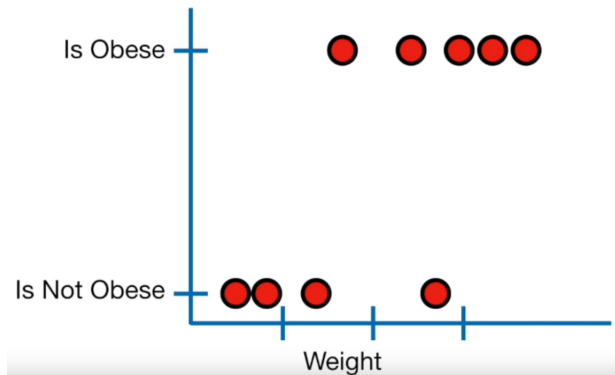


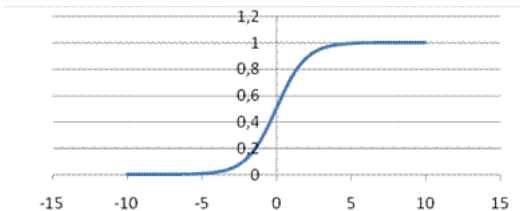
$$h_{\theta}(x) = \theta^T x$$

h - hipoteza θ - parametr

Regresja logistyczna

Jest to szczególny przypadek uogólnionej regresji liniowej, gdy **zmienna zależna jest dychotomiczna** - przyjmuje tylko dwie wartości (sukces, porażka)

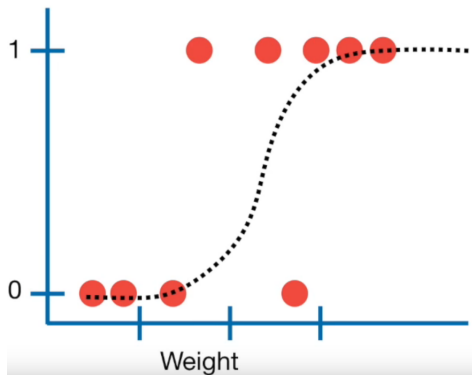




funkcja sigmoidalna, logistyczna

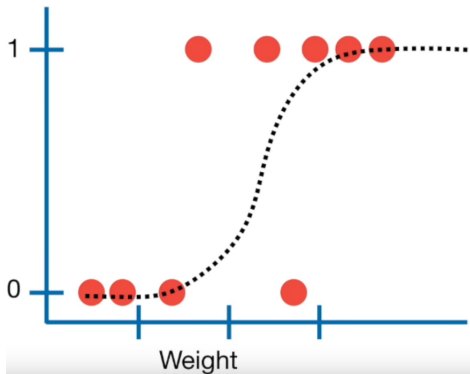
$$g(x) = \frac{1}{1+e^{-x}}$$

Regresja logistyczna



$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

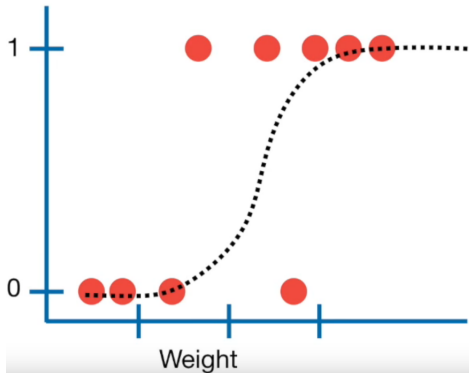
Regresja logistyczna



$$h_{\theta}(x) = \frac{1}{1+e^{-\theta^T x}}$$

$h_{\theta}(x)$ - estymowane
prawdopodobieństwo tego,
że $y=1$ dla danego x

Regresja logistyczna

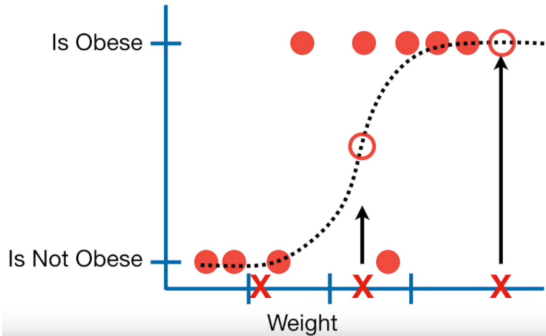


$$h_{\theta}(x) = \frac{1}{1+e^{-\theta^T x}}$$

$$0 \leq h_{\theta}(x) \leq 1$$

$h_{\theta}(x)$ - estymowane
prawdopodobieństwo tego,
że $y=1$ dla danego x

Regresja logistyczna

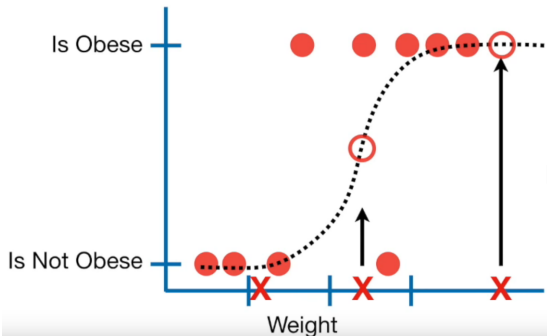


$$h_{\theta}(x) = \frac{1}{1+e^{-\theta^T x}}$$

$$0 \leq h_{\theta}(x) \leq 1$$

$h_{\theta}(x)$ - estymowane
prawdopodobieństwo tego,
że $y=1$ dla danego x

Regresja logistyczna



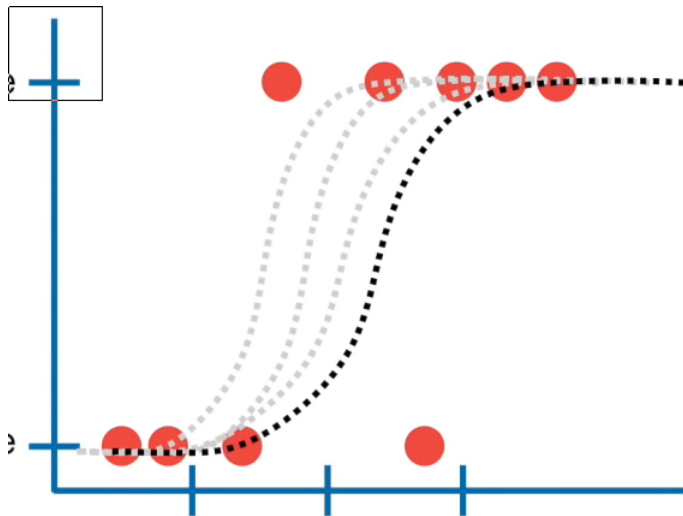
$$h_{\theta}(x) = \frac{1}{1+e^{-\theta^T x}}$$

$$0 \leq h_{\theta}(x) \leq 1$$

$h_{\theta}(x)$ - estymowane prawdopodobieństwo tego, że $y=1$ dla danego x

$h_{\theta}(x) = P(y = 1|x; \theta)$
"prawdopodobieństwo, że $y=1$, mając x , przy określonym parametrze θ "

Regresja logistyczna



Dopasowanie parametru θ

Zbiór uczący składa się z m par x, y :

$$(x^1, y^1), (x^2, y^2), \dots, (x^m, y^m),$$

gdzie każda para posiada n cech wejściowych $x \in \begin{array}{|l} x_0 \\ x_1 \\ \dots \\ x_n \end{array}$ oraz

wyście: $y \in 0, 1$

$$h_{\theta}(x) = \frac{1}{1+e^{-\theta^T x}}$$

Funkcja kosztu dla regresji liniowej:

$$\text{Cost}(h_{\theta}(x), y) = (h_{\theta}(x) - y)^2$$

Dopasowanie parametru θ - Funkcja kosztu

Funkcja kosztu dla regresji liniowej:

$$\text{Cost}(h_{\theta}(x), y) = (h_{\theta}(x) - y)^2$$

Funkcja kosztu dla regresji logistycznej:

$$\text{Cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{gdy } y = 1 \\ -\log(1 - h_{\theta}(x)) & \text{gdy } y = 0 \end{cases}$$

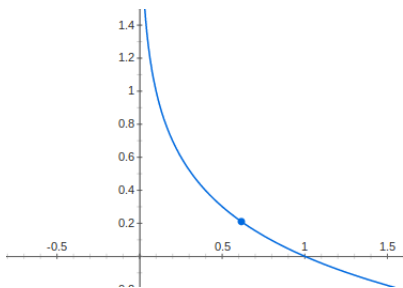
Funkcja kosztu dla regresji logistycznej:

$$\text{Cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{gdy } y = 1 \\ -\log(1 - h_{\theta}(x)) & \text{gdy } y = 0 \end{cases}$$

Dopasowanie parametru θ - Funkcja kosztu

Funkcja kosztu dla regresji logistycznej:

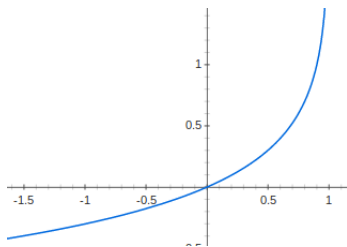
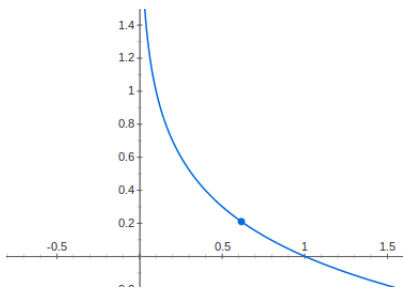
$$\text{Cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{gdy } y = 1 \\ -\log(1 - h_{\theta}(x)) & \text{gdy } y = 0 \end{cases}$$



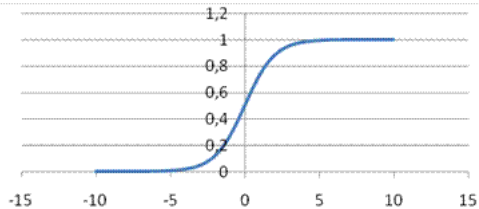
Dopasowanie parametru θ - Funkcja kosztu

Funkcja kosztu dla regresji logistycznej:

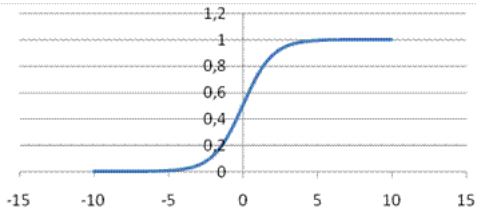
$$\text{Cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{gd}y \ y = 1 \\ -\log(1 - h_{\theta}(x)) & \text{gd}y \ y = 0 \end{cases}$$



Decision Boundary

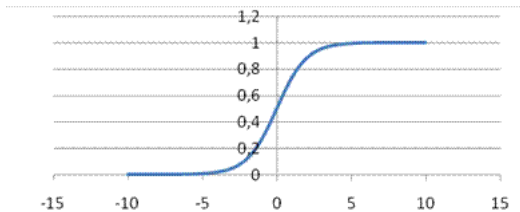


Decision Boundary



$$g(z) = \frac{1}{1+e^{-z}}$$
$$h_{\theta}(x) = g(\theta^T x)$$

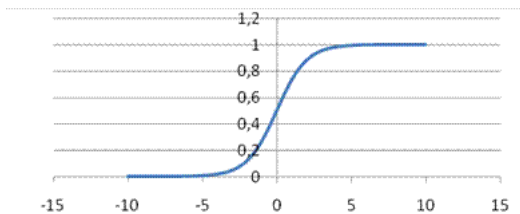
Decision Boundary



$$g(z) = \frac{1}{1+e^{-z}}$$
$$h_{\theta}(x) = g(\theta^T x)$$

Założmy predykcję $y=1$ jeśli $h_{\theta}(x) \geq 0.5$

Decision Boundary

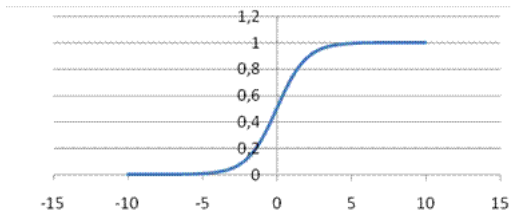


Założmy predykcję $y=1$ jeśli $h_{\theta}(x) \geq 0.5$

$$g(z) = \frac{1}{1+e^{-z}}$$
$$h_{\theta}(x) = g(\theta^T x)$$

$$g(z) \geq 0.5 \text{ gdy } z \geq 0$$

Decision Boundary

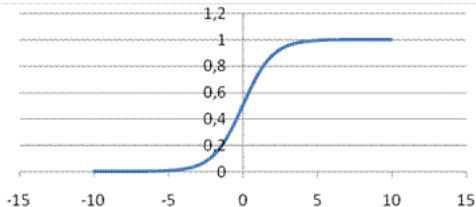


Założmy predykcyję $y=1$ jeśli $h_{\theta}(x) \geq 0.5$

$$g(z) = \frac{1}{1+e^{-z}}$$
$$h_{\theta}(x) = g(\theta^T x)$$

$$g(z) \geq 0.5 \text{ gdy } z \geq 0$$
$$h_{\theta}(x) \geq 0.5 \text{ gdy } \theta^T x \geq 0$$

Decision Boundary



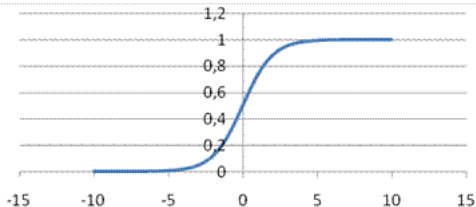
Założmy predykcję $y=1$ jeśli $h_{\theta}(x) \geq 0.5$

predykcję $y=0$ jeśli $h_{\theta}(x) < 0.5$

$$g(z) = \frac{1}{1+e^{-z}}$$
$$h_{\theta}(x) = g(\theta^T x)$$

$$g(z) \geq 0.5 \text{ gdy } z \geq 0$$
$$h_{\theta}(x) \geq 0.5 \text{ gdy } \theta^T x \geq 0$$

Decision Boundary



Założmy predykcję $y=1$ jeśli $h_{\theta}(x) \geq 0.5$

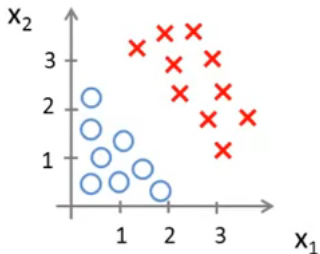
predykcję $y=0$ jeśli $h_{\theta}(x) < 0.5$

$$g(z) = \frac{1}{1+e^{-z}}$$
$$h_{\theta}(x) = g(\theta^T x)$$

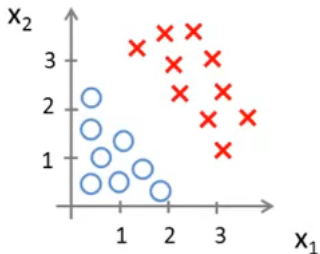
$$g(z) \geq 0.5 \text{ gdy } z \geq 0$$
$$h_{\theta}(x) \geq 0.5 \text{ gdy } \theta^T x \geq 0$$

$$h_{\theta}(x) < 0.5 \text{ gdy } \theta^T x < 0$$

Decision Boundary - przykład

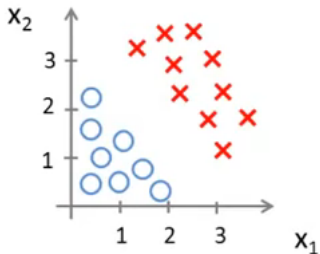


Decision Boundary - przykład



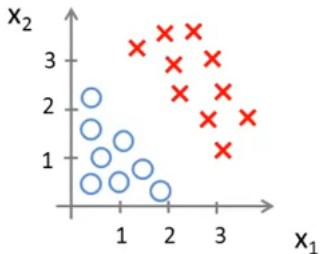
$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

Decision Boundary - przykład



$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$
$$\theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix}$$

Decision Boundary - przykład

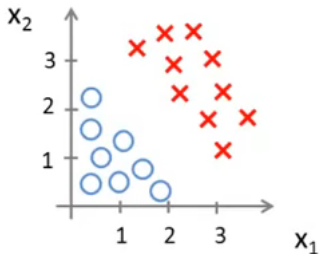


$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

$$\theta = \begin{vmatrix} -3 \\ 1 \\ 1 \end{vmatrix}$$

$$-3 + x_1 + x_2 \geq 0$$

Decision Boundary - przykład



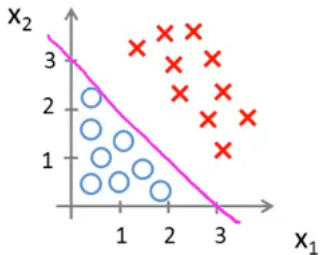
$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

$$\theta = \begin{vmatrix} -3 \\ 1 \\ 1 \end{vmatrix}$$

$$-3 + x_1 + x_2 \geq 0$$

$$x_1 + x_2 = 3 \Rightarrow h_{\theta}(x) = 0.5$$

Decision Boundary - przykład



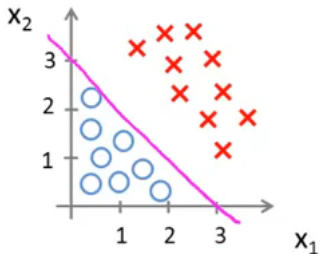
$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

$$\theta = \begin{vmatrix} -3 \\ 1 \\ 1 \end{vmatrix}$$

$$-3 + x_1 + x_2 \geq 0$$

$$x_1 + x_2 = 3 \Rightarrow h_{\theta}(x) = 0.5$$

Decision Boundary - przykład



$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

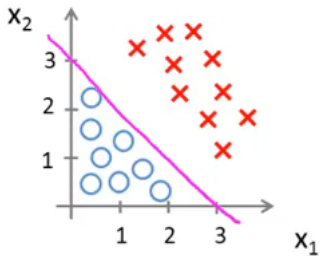
$$\theta = \begin{vmatrix} -3 \\ 1 \\ 1 \end{vmatrix}$$

$$-3 + x_1 + x_2 \geq 0$$

$$x_1 + x_2 = 3 \Rightarrow h_{\theta}(x) = 0.5$$

Predykcja $y=1$ dla $x_1 + x_2 \geq 3$

Decision Boundary - przykład



$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$$

$$\theta = \begin{vmatrix} -3 \\ 1 \\ 1 \end{vmatrix}$$

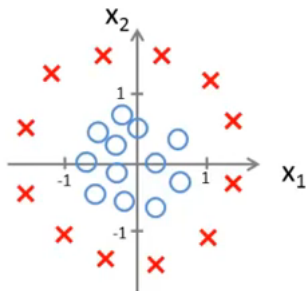
$$-3 + x_1 + x_2 \geq 0$$

$$x_1 + x_2 = 3 \Rightarrow h_{\theta}(x) = 0.5$$

Predykcja $y=1$ dla $x_1 + x_2 \geq 3$

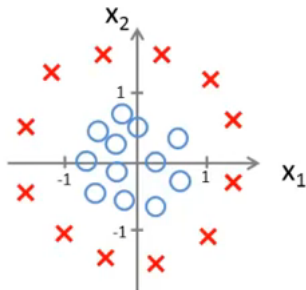
Predykcja $y=0$ dla $x_1 + x_2 < 3$

Decision Boundary - przykład nieliniowy



$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$$

Decision Boundary - przykład nieliniowy



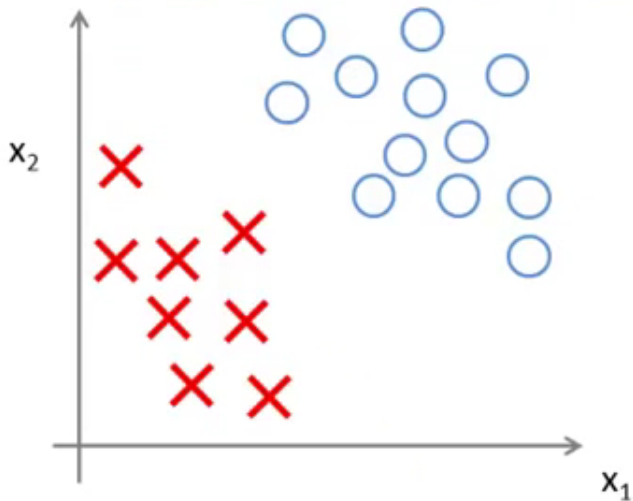
$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$$

Uwaga: Granice decyzyjne nie są obliczane na podstawie zbioru uczącego, zbiór uczący optymalizuje parametr θ .

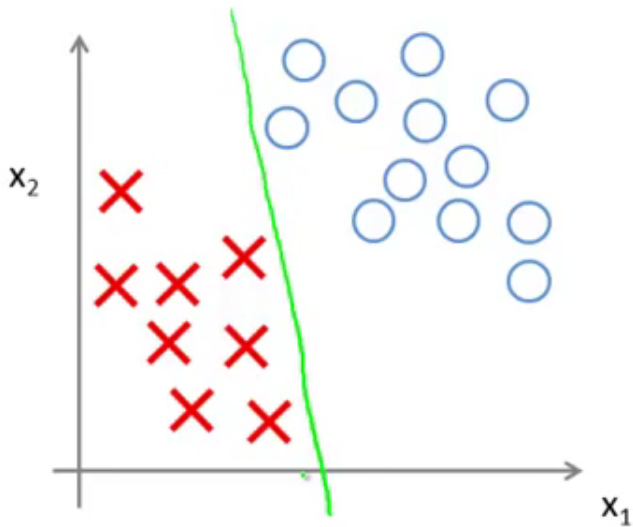
- 1 Wprowadzenie do klasyfikacji
- 2 kNN - k Nearest Neighbours
- 3 Regresja logistyczna
- 4 **SVM**
- 5 Naive Bayes
- 6 Drzewa Decyzyjne
- 7 Random Forest

Problem regresji logistycznej

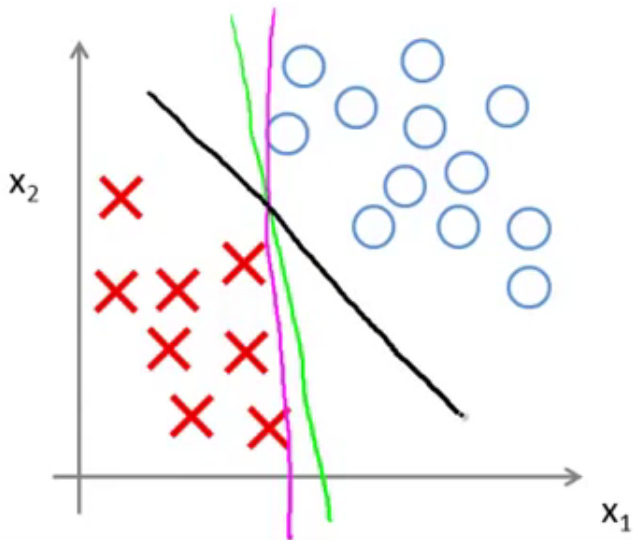
Problem regresji logistycznej



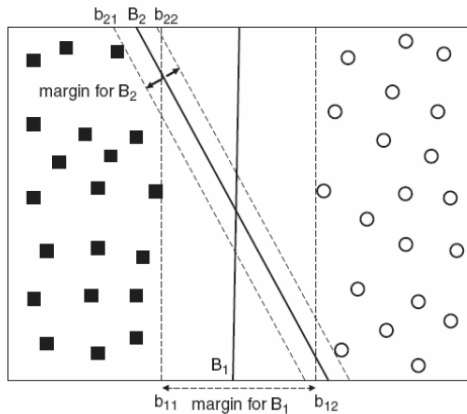
Problem regresji logistycznej



Problem regresji logistycznej



Metoda Wektorów nośnych - Support Vector Machine



Marginesy są to proste (hiperpłaszczyzny) otrzymane przez równoległe przesuwanie granicy aż do pierwszych punktów z obu klas.

Odległość między nimi to **marginesy** klasyfikatora.

Szerszy margines to lepsze własności generalizacji oraz mniejsza podatność na przeuczenie.

Support Vectors

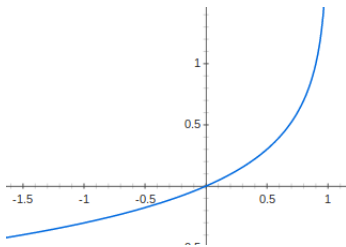
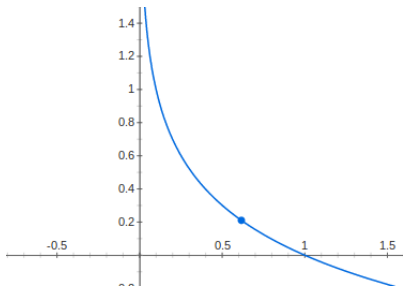
Alice Zhao

Magdalena Mazur-Milecka

Funkcja kosztu

Funkcja kosztu dla regresji logistycznej:

$$\text{Cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{gd}y\ y = 1 \\ -\log(1 - h_{\theta}(x)) & \text{gd}y\ y = 0 \end{cases}$$



Funkcja kosztu dla regresji logistycznej:

$$\text{Cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{gd}y \ y = 1 \\ -\log(1 - h_{\theta}(x)) & \text{gd}y \ y = 0 \end{cases}$$

$$\text{Cost} = -(y \log h_{\theta}(x) + (1 - y) \log(1 - h_{\theta}(x)))$$

Funkcja kosztu dla regresji logistycznej:

$$\text{Cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{gdym } y = 1 \\ -\log(1 - h_{\theta}(x)) & \text{gdym } y = 0 \end{cases}$$

$$\begin{aligned} \text{Cost} &= -(y \log h_{\theta}(x) + (1 - y) \log(1 - h_{\theta}(x))) \\ &= -y \log \frac{1}{1 + e^{-\theta^T x}} - (1 - y) \log \left(1 - \frac{1}{1 + e^{-\theta^T x}}\right) \end{aligned}$$

Funkcja kosztu dla regresji logistycznej:

$$\text{Cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{gd}y \ y = 1 \\ -\log(1 - h_{\theta}(x)) & \text{gd}y \ y = 0 \end{cases}$$

$$\begin{aligned} \text{Cost} &= -(y \log h_{\theta}(x) + (1 - y) \log(1 - h_{\theta}(x))) \\ &= -y \log \frac{1}{1 + e^{-\theta^T x}} - (1 - y) \log\left(1 - \frac{1}{1 + e^{-\theta^T x}}\right) \end{aligned}$$

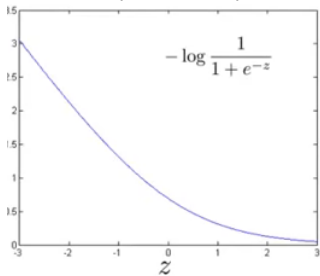
$$\min_{\theta} \frac{1}{m} \underbrace{\left[\sum_{i=1}^m y^{(i)} (-\log h_{\theta}(x^{(i)})) + (1 - y^{(i)}) (-\log(1 - h_{\theta}(x^{(i)}))) \right]}_{J(\theta)}$$

$$\begin{aligned} \text{Cost} &= -(y \log h_{\theta}(x) + (1 - y) \log(1 - h_{\theta}(x))) \\ &= -y \log \frac{1}{1+e^{-\theta^T x}} - (1 - y) \log\left(1 - \frac{1}{1+e^{-\theta^T x}}\right) \end{aligned}$$

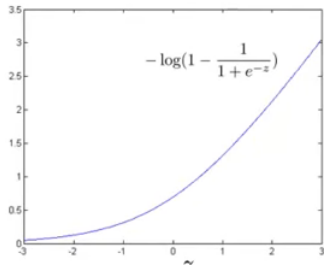
$$\begin{aligned} \min_{\theta} \frac{1}{m} \left[\sum_{i=1}^m y^{(i)} (-\log h_{\theta}(x^{(i)})) + (1 - y^{(i)}) (-\log(1 - h_{\theta}(x^{(i)}))) \right] + \\ \underbrace{\hspace{15em}}_{J(\theta)} \\ + \lambda \underbrace{\frac{\partial J(\theta)}{\partial \theta_j}}_B \end{aligned}$$

Funkcja kosztu

dla $y=1$ ($\theta^T x \gg 0$)

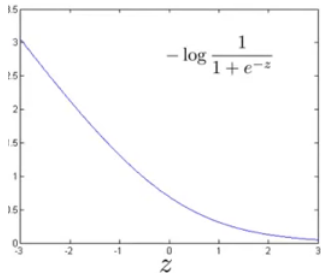


dla $y=0$ ($\theta^T x \ll 0$)

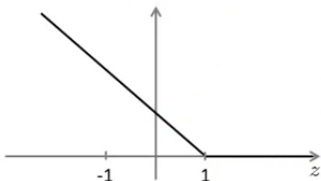
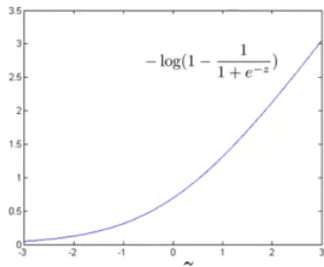


Funkcja kosztu

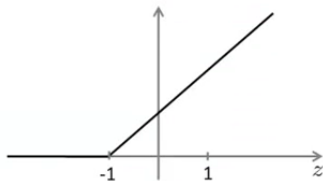
dla $y=1$ ($\theta^T x \gg 0$)



dla $y=0$ ($\theta^T x \ll 0$)



$\theta^T x \geq 1$



$\theta^T x \leq -1$

Regresja logistyczna:

$$\min_{\theta} \frac{1}{m} \left[\sum_{i=1}^m y^{(i)} (-\log h_{\theta}(x^{(i)})) + (1 - y^{(i)}) (-\log(1 - h_{\theta}(x^{(i)}))) \right] + \lambda B$$

Regresja logistyczna:

$$\min_{\theta} \frac{1}{m} \left[\sum_{i=1}^m y^{(i)} (-\log h_{\theta}(x^{(i)})) + (1 - y^{(i)}) (-\log(1 - h_{\theta}(x^{(i)}))) \right] + \lambda B$$
$$J(\theta) + \lambda B$$

gdzie λ - parametr regularyzacji

Regresja logistyczna:

$$\min_{\theta} \frac{1}{m} \left[\sum_{i=1}^m y^{(i)} (-\log h_{\theta}(x^{(i)})) + (1 - y^{(i)}) (-\log(1 - h_{\theta}(x^{(i)}))) \right] + \lambda B$$

$$J(\theta) + \lambda B$$

SVM:

$$CJ(\theta) + B$$

Regresja logistyczna:

$$\min_{\theta} \frac{1}{m} \left[\sum_{i=1}^m y^{(i)} (-\log h_{\theta}(x^{(i)})) + (1 - y^{(i)}) (-\log(1 - h_{\theta}(x^{(i)}))) \right] + \lambda B$$

$$J(\theta) + \lambda B$$

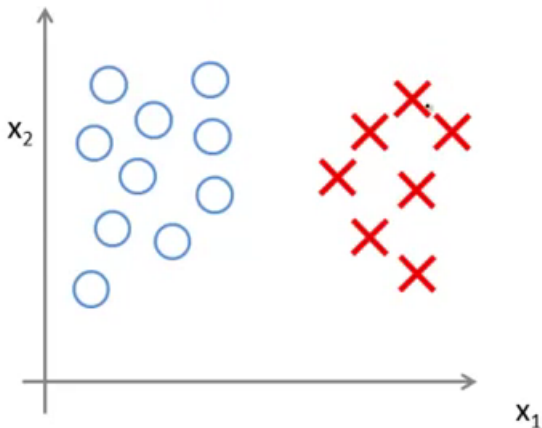
SVM:

$$CJ(\theta) + B$$

$$\min_{\theta} C \left[\sum_{i=1}^m y^{(i)} (\text{cost}_1(\theta^T x^{(i)})) + (1 - y^{(i)}) (\text{cost}_0(\theta^T x^{(i)})) \right] + B$$

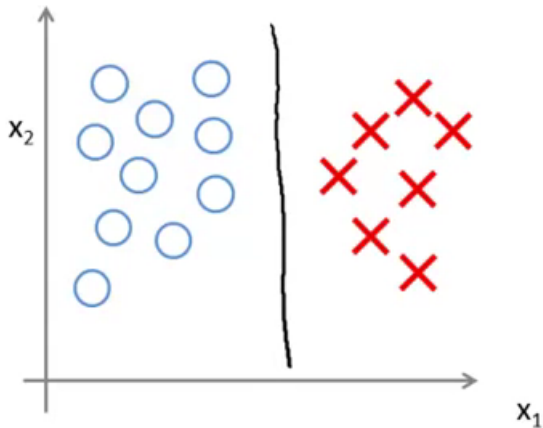
Funkcja kosztu

$$\min_{\theta} CJ(\theta) + B$$



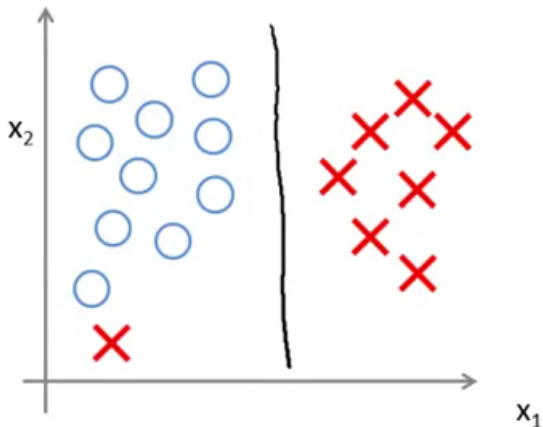
Funkcja kosztu

$$\min_{\theta} CJ(\theta) + B$$



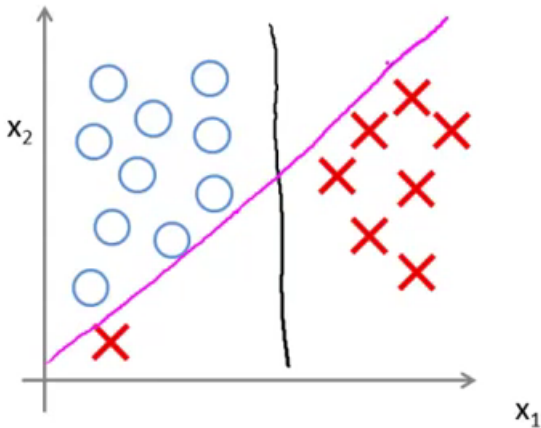
Funkcja kosztu

$$\min_{\theta} CJ(\theta) + B$$

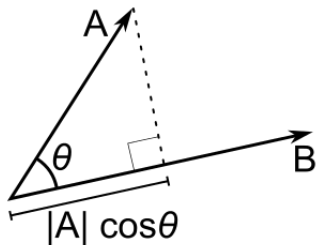


Funkcja kosztu

$$\min_{\theta} CJ(\theta) + B$$

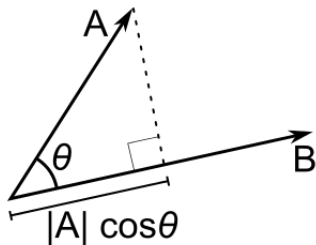


Wektory nośne - Support Vectors



$$\mathbf{A} = \begin{bmatrix} A_1 \\ A_2 \end{bmatrix} \quad \mathbf{B} = \begin{bmatrix} B_1 \\ B_2 \end{bmatrix}$$

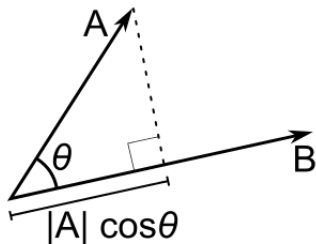
Wektory nośne - Support Vectors



$$\mathbf{A} = \begin{bmatrix} A_1 \\ A_2 \end{bmatrix} \quad \mathbf{B} = \begin{bmatrix} B_1 \\ B_2 \end{bmatrix}$$

$$\mathbf{B}^T \mathbf{A} = B_1 A_1 + B_2 A_2 = |\mathbf{B}| |\mathbf{A}| \cos \theta$$

Wektory nośne - Support Vectors

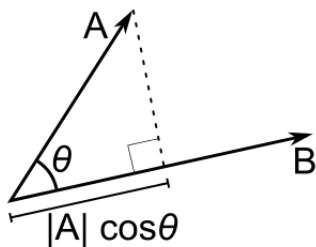


$$\mathbf{A} = \begin{bmatrix} A_1 \\ A_2 \end{bmatrix} \quad \mathbf{B} = \begin{bmatrix} B_1 \\ B_2 \end{bmatrix}$$

$$\mathbf{B}^T \mathbf{A} = B_1 A_1 + B_2 A_2 = |\mathbf{B}| |\mathbf{A}| \cos \theta = p |\mathbf{B}|$$

$$|\mathbf{A}| \cos \theta = p$$

Wektory nośne - Support Vectors

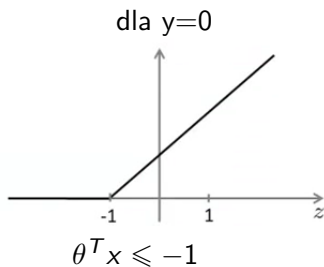
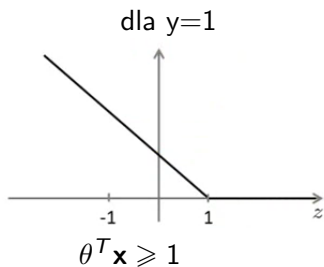


$$\mathbf{A} = \begin{bmatrix} A_1 \\ A_2 \end{bmatrix} \quad \mathbf{B} = \begin{bmatrix} B_1 \\ B_2 \end{bmatrix}$$

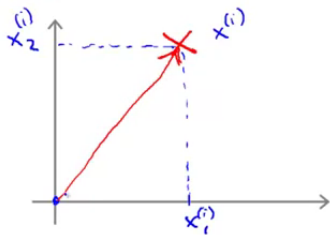
$$\mathbf{B}^T \mathbf{A} = B_1 A_1 + B_2 A_2 = |\mathbf{B}| |\mathbf{A}| \cos \theta = p |\mathbf{B}|$$

$$|\mathbf{A}| \cos \theta = p$$

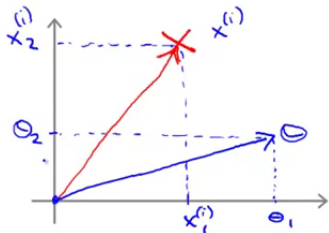
p jest < 0 dla $\theta > 90^\circ$



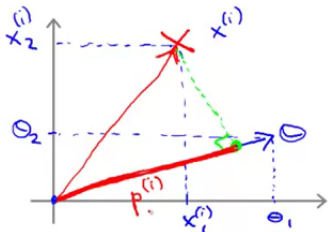
$$\theta^T \mathbf{x}^{(i)} \geq 1 \text{ dla } y^{(i)} = 1$$
$$\theta^T \mathbf{x}^{(i)} \leq -1 \text{ dla } y^{(i)} = 0$$



$$\theta^T \mathbf{x}^{(i)} \geq 1 \text{ dla } y^{(i)} = 1$$
$$\theta^T \mathbf{x}^{(i)} \leq -1 \text{ dla } y^{(i)} = 0$$



$$\theta^T \mathbf{x}^{(i)} \geq 1 \text{ dla } y^{(i)} = 1$$
$$\theta^T \mathbf{x}^{(i)} \leq -1 \text{ dla } y^{(i)} = 0$$

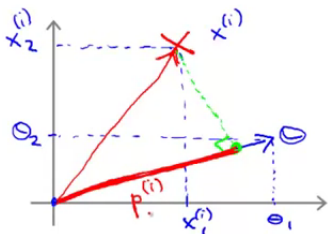


Wektory nośne

$$\theta^T \mathbf{x}^{(i)} \geq 1 \text{ dla } y^{(i)} = 1$$
$$\theta^T \mathbf{x}^{(i)} \leq -1 \text{ dla } y^{(i)} = 0$$

$$\theta^T \mathbf{x}^{(i)} = \rho^{(i)} |\theta|$$

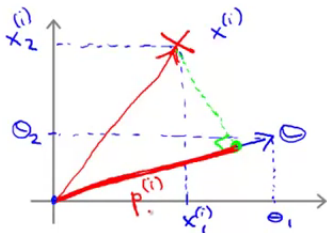
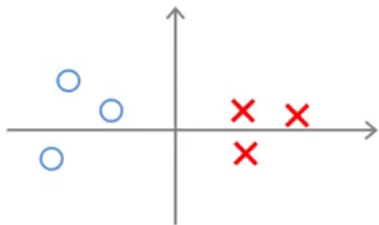
$\rho^{(i)}$ - rzut wektora $\mathbf{x}^{(i)}$ na θ



Wektory nośne

$$\theta^T \mathbf{x}^{(i)} \geq 1 \text{ dla } y^{(i)} = 1$$
$$\theta^T \mathbf{x}^{(i)} \leq -1 \text{ dla } y^{(i)} = 0$$

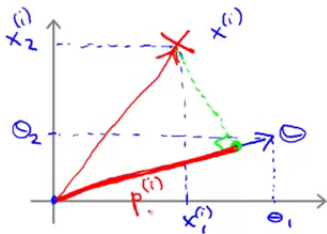
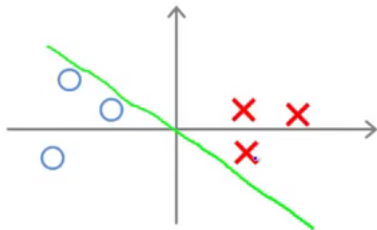
$$\theta^T \mathbf{x}^{(i)} = \rho^{(i)} |\theta|$$



Wektory nośne

$$\theta^T \mathbf{x}^{(i)} \geq 1 \text{ dla } y^{(i)} = 1$$
$$\theta^T \mathbf{x}^{(i)} \leq -1 \text{ dla } y^{(i)} = 0$$

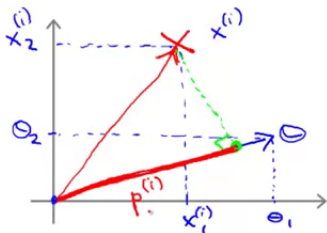
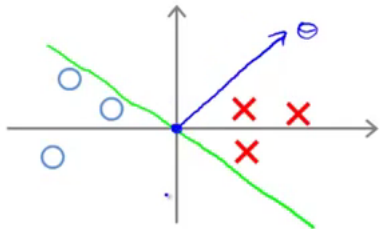
$$\theta^T \mathbf{x}^{(i)} = p^{(i)} |\theta|$$



Wektory nośne

$$\theta^T \mathbf{x}^{(i)} \geq 1 \text{ dla } y^{(i)} = 1$$
$$\theta^T \mathbf{x}^{(i)} \leq -1 \text{ dla } y^{(i)} = 0$$

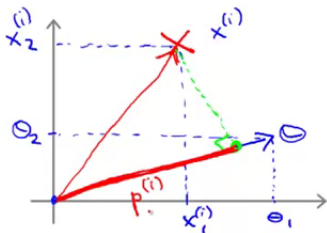
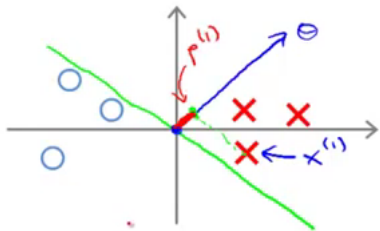
$$\theta^T \mathbf{x}^{(i)} = p^{(i)} |\theta|$$



Wektory nośne

$$\theta^T \mathbf{x}^{(i)} \geq 1 \text{ dla } y^{(i)} = 1$$
$$\theta^T \mathbf{x}^{(i)} \leq -1 \text{ dla } y^{(i)} = 0$$

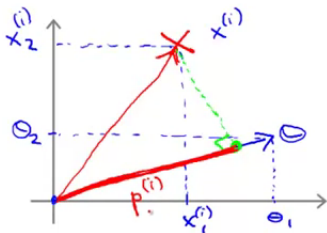
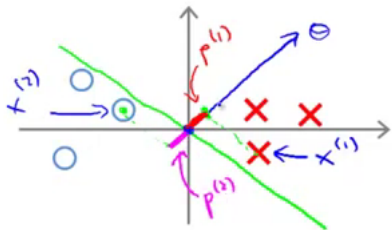
$$\theta^T \mathbf{x}^{(i)} = p^{(i)} |\theta|$$



Wektory nośne

$$\theta^T \mathbf{x}^{(i)} \geq 1 \text{ dla } y^{(i)} = 1$$
$$\theta^T \mathbf{x}^{(i)} \leq -1 \text{ dla } y^{(i)} = 0$$

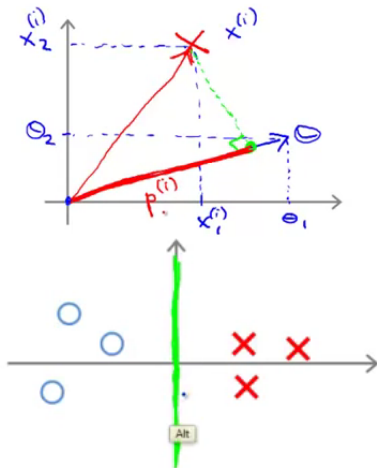
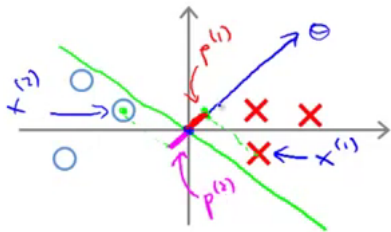
$$\theta^T \mathbf{x}^{(i)} = p^{(i)} |\theta|$$



Wektory nośne

$$\theta^T \mathbf{x}^{(i)} \geq 1 \text{ dla } y^{(i)} = 1$$
$$\theta^T \mathbf{x}^{(i)} \leq -1 \text{ dla } y^{(i)} = 0$$

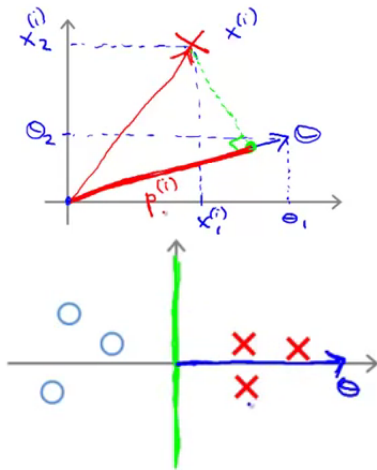
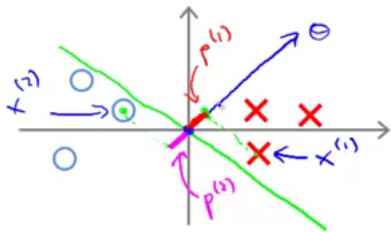
$$\theta^T \mathbf{x}^{(i)} = \rho^{(i)} |\theta|$$



Wektory nośne

$$\theta^T \mathbf{x}^{(i)} \geq 1 \text{ dla } y^{(i)} = 1$$
$$\theta^T \mathbf{x}^{(i)} \leq -1 \text{ dla } y^{(i)} = 0$$

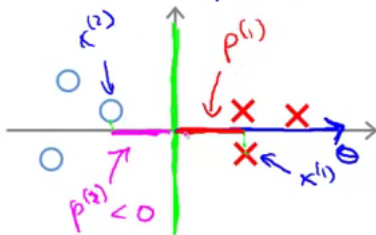
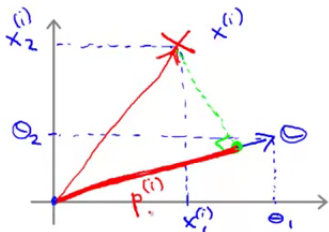
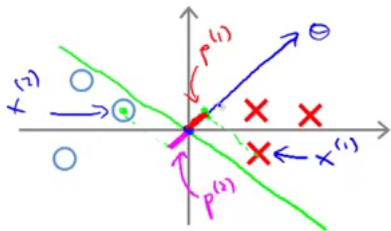
$$\theta^T \mathbf{x}^{(i)} = \rho^{(i)} |\theta|$$



Wektory nośne

$$\theta^T \mathbf{x}^{(i)} \geq 1 \text{ dla } y^{(i)} = 1$$
$$\theta^T \mathbf{x}^{(i)} \leq -1 \text{ dla } y^{(i)} = 0$$

$$\theta^T \mathbf{x}^{(i)} = \rho^{(i)} |\theta|$$



Używane podczas poszukiwania nieliniowych granic decyzyjnych.

$$\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 + \theta_5 x_2^2 + \dots$$

Używane podczas poszukiwania nieliniowych granic decyzyjnych.

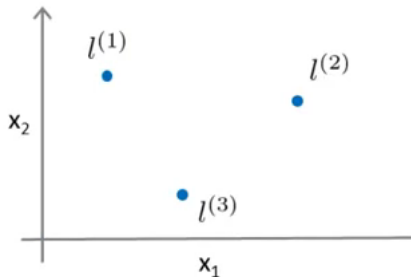
$$\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 + \theta_5 x_2^2 + \dots$$

$$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 + \theta_4 f_4 + \theta_5 f_5 + \dots$$

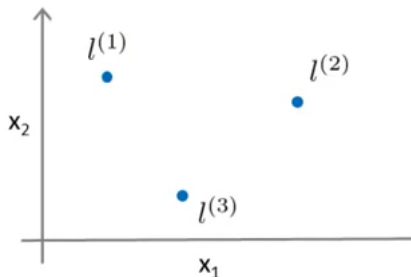
gdzie f -funkcja podobieństwa danej do znaczników (landmarks) używająca jądra

$$f = \text{similarity}(x, l^{(i)}) = k(x, l^{(i)})$$

$$f = \text{similarity}(x, l^{(i)}) = k(x, l^{(i)})$$



$$f = \text{similarity}(x, l^{(i)}) = k(x, l^{(i)})$$



Jądro przekształca wektor przypadku x w nowy wektor f m -wymiarowy, gdzie m -wielkość zbioru treningowego.

Przykład jądra: gaussowskie

$$f_i = \textit{similarity}(x, l^{(i)}) = \exp\left(-\frac{|x - l^{(i)}|^2}{2\sigma^2}\right)$$

Przykład jądra: gaussowskie

$$f_i = \text{similarity}(x, l^{(i)}) = \exp\left(-\frac{|x - l^{(i)}|^2}{2\sigma^2}\right)$$

dla $x \approx l^{(i)}$

$f_i \approx 1$

Przykład jądra: gaussowskie

$$f_i = \text{similarity}(x, l^{(i)}) = \exp\left(-\frac{|x - l^{(i)}|^2}{2\sigma^2}\right)$$

dla $x \approx l^{(i)}$

$$f_i \approx 1$$

dla x dalekiego od $l^{(i)}$

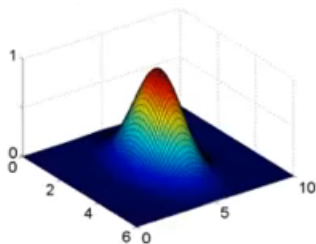
$$f_i \approx 0$$

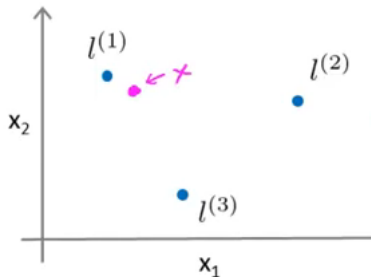
Przykład jądra: gaussowskie

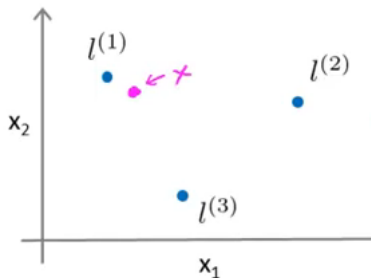
$$f_i = \text{similarity}(x, l^{(i)}) = \exp\left(-\frac{|x - l^{(i)}|^2}{2\sigma^2}\right)$$

dla $x \approx l^{(i)}$
 $f_i \approx 1$

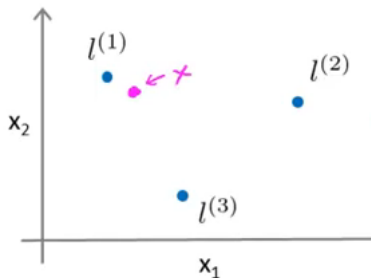
dla x dalekiego od $l^{(i)}$
 $f_i \approx 0$







$$x^{(i)} = \begin{bmatrix} f_1^{(i)} \\ f_2^{(i)} \\ \dots \\ f_m^{(i)} \end{bmatrix}$$



$$x^{(i)} = \begin{bmatrix} f_1^{(i)} \\ f_2^{(i)} \\ \dots \\ f_m^{(i)} \end{bmatrix}$$

Uwaga: wszystkie założenia zmieniają się z $\theta^T x$ na: $\theta^T f$

- C - duże C =overfitting, małe C =niedopasowanie,

SVM- Wybór parametrów

- C - duże C =overfitting, małe C =niedopasowanie,
- jądro (lub jego brak)

- C - duże C =overfitting, małe C =niedopasowanie,
- jądro (lub jego brak)
- σ w jądrze gaussowskim,

Mocne strony SVM

Słabe strony SVM

Mocne strony SVM

- Potrafi poradzić sobie z outliersami i danymi nieseparowalnymi liniowo

Słabe strony SVM

Mocne strony SVM

- Potrafi poradzić sobie z outliersami i danymi nieseparowalnymi liniowo
- Minimalizacja funkcji kosztu wymaga mniej obliczeń niż w przypadku regresji logistycznej

Słabe strony SVM

Mocne strony SVM

- Potrafi poradzić sobie z outliersami i danymi nieseparowalnymi liniowo
- Minimalizacja funkcji kosztu wymaga mniej obliczeń niż w przypadku regresji logistycznej
- Stopień skomplikowania jest niezależny od liczby wymiarów

Słabe strony SVM

Mocne strony SVM

- Potrafi poradzić sobie z outliersami i danymi nieseparowalnymi liniowo
- Minimalizacja funkcji kosztu wymaga mniej obliczeń niż w przypadku regresji logistycznej
- Stopień skomplikowania jest niezależny od liczby wymiarów
- Nie jest czuły na przetrenowanie

Słabe strony SVM

Mocne strony SVM

- Potrafi poradzić sobie z outliersami i danymi nieseparowalnymi liniowo
- Minimalizacja funkcji kosztu wymaga mniej obliczeń niż w przypadku regresji logistycznej
- Stopień skomplikowania jest niezależny od liczby wymiarów
- Nie jest czuły na przetrenowanie
- Dzięki dopasowaniu jądra, algorytm osiąga dużą skuteczność w praktyce

Słabe strony SVM

Mocne strony SVM

- Potrafi poradzić sobie z outliersami i danymi nieseparowalnymi liniowo
- Minimalizacja funkcji kosztu wymaga mniej obliczeń niż w przypadku regresji logistycznej
- Stopień skomplikowania jest niezależny od liczby wymiarów
- Nie jest czuły na przetrenowanie
- Dzięki dopasowaniu jądra, algorytm osiąga dużą skuteczność w praktyce

Słabe strony SVM

- Długotrwały trening przez minimalizację funkcji

Mocne strony SVM

- Potrafi poradzić sobie z outliersami i danymi nieseparowalnymi liniowo
- Minimalizacja funkcji kosztu wymaga mniej obliczeń niż w przypadku regresji logistycznej
- Stopień skomplikowania jest niezależny od liczby wymiarów
- Nie jest czuły na przetrenowanie
- Dzięki dopasowaniu jądra, algorytm osiąga dużą skuteczność w praktyce

Słabe strony SVM

- Długotrwały trening przez minimalizację funkcji
- Brak możliwości wprowadzenia prior knowledge

- 1 Wprowadzenie do klasyfikacji
- 2 kNN - k Nearest Neighbours
- 3 Regresja logistyczna
- 4 SVM
- 5 **Naive Bayes**
- 6 Drzewa Decyzyjne
- 7 Random Forest

Naiwny klasyfikator Bayesa

Prosty klasyfikator probabilistyczny. Model prawdopodobieństwa można wprowadzić korzystając z **twierdzenia Bayesa**.

Prosty klasyfikator probabilistyczny. Model prawdopodobieństwa można wprowadzić korzystając z **twierdzenia Bayesa**.

$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$$

gdzie: $P(A|B)$ - prawdopodobieństwo warunkowe,
prawdopodobieństwo zajścia zdarzenia A jeśli zajdzie zdarzenie B .

Naiwny klasyfikator Bayesa

Prosty klasyfikator probabilistyczny. Model prawdopodobieństwa można wprowadzić korzystając z **twierdzenia Bayesa**.

Naiwność - założenie o *wzajemnej niezależności predyktorów* (cech) - często błędne.

Naiwny klasyfikator Bayesa

Prosty klasyfikator probabilistyczny. Model prawdopodobieństwa można wprowadzić korzystając z **twierdzenia Bayesa**.

Naiwność - założenie o *wzajemnej niezależności predyktorów* (cech) - często błędne.

Naiwne klasyfikatory bayesowskie można uczyć z nadzorem. Estymacja parametru używa często **metody maksymalnego prawdopodobieństwa a posteriori**.

Naiwny klasyfikator Bayesa

Prosty klasyfikator probabilistyczny. Model prawdopodobieństwa można wprowadzić korzystając z **twierdzenia Bayesa**.

Naiwność - założenie o *wzajemnej niezależności predyktorów* (cech) - często błędne.

Naiwne klasyfikatory bayesowskie można uczyć z nadzorem. Estymacja parametru używa często **metody maksymalnego prawdopodobieństwa a posteriori**.

Pomimo uproszczonych założeń, klasyfikatory te często dają bardzo dobre wyniki.

$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$$

$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$$

$$P(C|F_1, \dots, F_n) = \frac{P(F_1, \dots, F_n|C)P(C)}{P(F_1, \dots, F_n)}$$

gdzie $P(C|F_1, \dots, F_n)$ - prawdopodobieństwo przynależności do klasy C pod warunkiem cech F_1, \dots, F_n

$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$$

$$P(C|F_1, \dots, F_n) = \frac{P(F_1, \dots, F_n|C)P(C)}{P(F_1, \dots, F_n)}$$

gdzie $P(C|F_1, \dots, F_n)$ - prawdopodobieństwo przynależności do klasy C pod warunkiem cech F_1, \dots, F_n

$$P(C|F_1, \dots, F_n) = P(F_1|C)P(F_2|C)\dots P(F_n|C)P(C) = P(C) \prod_{i=1}^n P(F_i|C)$$

$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$$

$$P(C|F_1, \dots, F_n) = \frac{P(F_1, \dots, F_n|C)P(C)}{P(F_1, \dots, F_n)}$$

gdzie $P(C|F_1, \dots, F_n)$ - prawdopodobieństwo przynależności do klasy C pod warunkiem cech F_1, \dots, F_n

$$P(C|F_1, \dots, F_n) = P(F_1|C)P(F_2|C)\dots P(F_n|C)P(C) = P(C) \prod_{i=1}^n P(F_i|C)$$

$$P(C|F_1, \dots, F_n) = \frac{1}{Z} P(C) \prod_{i=1}^n P(F_i|C)$$

gdzie Z - współczynnik skalowania zależny od F_1, \dots, F_n

Model uczy się:

$P(F|C)$ oraz $P(C)$ - inf. a priori

Na podstawie nauczonych prawdopodobieństw można obliczyć prawdopodobieństwo przynależności do klasy $C = j$ w danym przypadku (dla danych cech $F = F_1, \dots, F_n$):

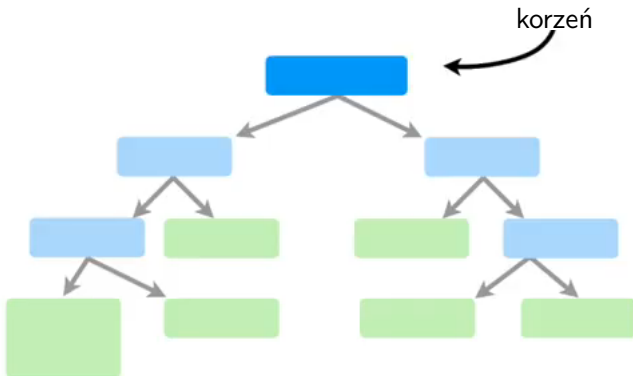
$$P(C = j|F) = \frac{P(F|C = j)P(C = j)}{P(F)}$$

Naiwny model probabilistyczny Bayesa - Przykład

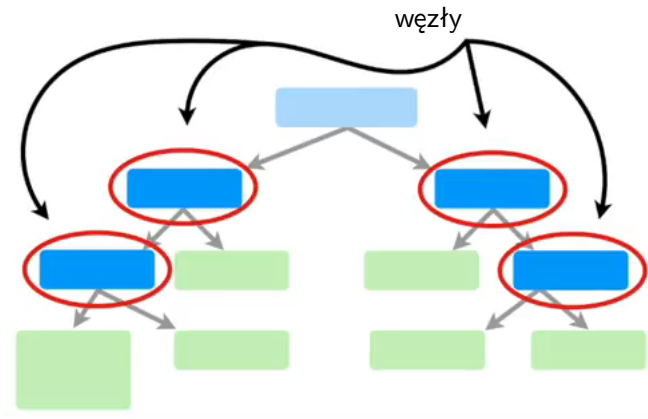
Outlook	Temp.	Humidity	Windy	Play Golf
Rainy	Hot	High	False	No
Rainy	Hot	High	True	No
Overcast	Hot	High	False	Yes
Sunny	Mild	High	False	Yes
Sunny	Cool	Normal	False	Yes
Sunny	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Rainy	Mild	High	False	No
Rainy	Cool	Normal	False	Yes
Sunny	Mild	Normal	False	Yes
Rainy	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Sunny	Mild	High	True	No

- 1 Wprowadzenie do klasyfikacji
- 2 kNN - k Nearest Neighbours
- 3 Regresja logistyczna
- 4 SVM
- 5 Naive Bayes
- 6 **Drzewa decyzyjne**
- 7 Random Forest

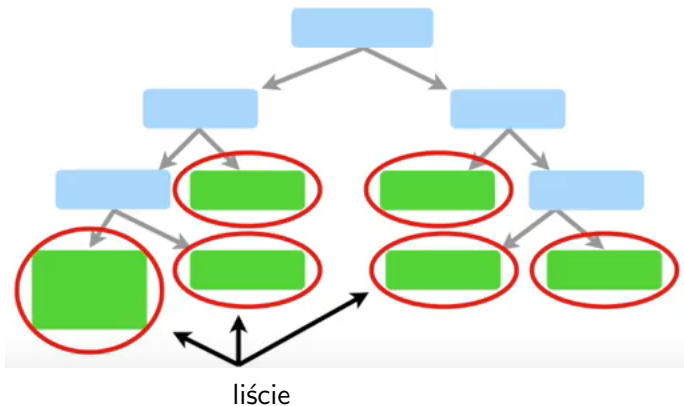
Algorytm zaczerpnięty z teorii decyzji, gdzie wspomaga procesy decyzyjne. W uczeniu maszynowym służy do pozyskiwania wiedzy ze zbioru uczącego.



Budowa drzewa



Budowa drzewa



- Drzewo nie ma pętli (cykli),
- Istnieje tylko jedna ścieżka między dwoma różnymi węzłami,
- Reprezentuje podział zbioru na klasy:
 - węzły opisują sposób podziału,
 - liście to klasy
 - krawędzie reprezentują wartości cech, na podstawie których dokonano podziału

Rodzaje cech możliwych do klasyfikacji przez drzewa:

- cechy binarne (tak, nie)
- cechy ilościowe
- cechy jakościowe (ranking)
- cechy wielokrotnego wyboru.

- 1 utwórz korzeń drzewa

Drzewa decyzyjne - Algorytm

- 1 utwórz korzeń drzewa
- 2 utwórz element drzewa, który:

Drzewa decyzyjne - Algorytm

- 1 utwórz korzeń drzewa
- 2 utwórz element drzewa, który:
- 3 jest liściem, jeżeli wszystkie przykłady ze zbioru uczącego są pozytywne lub negatywne

Drzewa decyzyjne - Algorytm

- 1 utwórz korzeń drzewa
- 2 utwórz element drzewa, który:
- 3 jest liściem, jeżeli wszystkie przykłady ze zbioru uczącego są pozytywne lub negatywne
- 4 jest węzłem w przeciwnym wypadku

- 1 utwórz korzeń drzewa
- 2 utwórz element drzewa, który:
- 3 jest liściem, jeżeli wszystkie przykłady ze zbioru uczącego są pozytywne lub negatywne
- 4 jest węzłem w przeciwnym wypadku
 - na węzeł wybierz cechę, który najlepiej klasyfikuje zbiór uczący

Drzewa decyzyjne - Algorytm

- 1 utwórz korzeń drzewa
- 2 utwórz element drzewa, który:
- 3 jest liściem, jeżeli wszystkie przykłady ze zbioru uczącego są pozytywne lub negatywne
- 4 jest węzłem w przeciwnym wypadku
 - na węzeł wybierz cechę, który najlepiej klasyfikuje zbiór uczący
 - dla każdej wartości cechy (węzła) stwórz gałąź i rozdziel zbiór uczący według tej cechy

Drzewa decyzyjne - Algorytm

- 1 utwórz korzeń drzewa
- 2 utwórz element drzewa, który:
- 3 jest liściem, jeżeli wszystkie przykłady ze zbioru uczącego są pozytywne lub negatywne
- 4 jest węzłem w przeciwnym wypadku
 - na węzeł wybierz cechę, który najlepiej klasyfikuje zbiór uczący
 - dla każdej wartości cechy (węzła) stwórz gałąź i rozdziel zbiór uczący według tej cechy
 - jeśli gałąź jest pusta, utwórz liść z wartością, jaką najczęściej przyjmuje cecha

Drzewa decyzyjne - Algorytm

- 1 utwórz korzeń drzewa
- 2 utwórz element drzewa, który:
- 3 jest liściem, jeżeli wszystkie przykłady ze zbioru uczącego są pozytywne lub negatywne
- 4 jest węzłem w przeciwnym wypadku
 - na węzeł wybierz cechę, który najlepiej klasyfikuje zbiór uczący
 - dla każdej wartości cechy (węzła) stwórz gałąź i rozdziel zbiór uczący według tej cechy
 - jeśli gałąź jest pusta, utwórz liść z wartością, jaką najczęściej przyjmuje cecha
 - w przeciwnym wypadku przejdź do kroku 4.

Miara wyboru cechy, która najlepiej klasyfikuje zbiór uczący:

Miara wyboru cechy, która najlepiej klasyfikuje zbiór uczący:

- **Entropia (Z)**

charakteryzuje zanieczyszczenie zbioru uczącego

$$Entropia(Z) = -p_+ \log_2 p_+ - p_- \log_2 p_-$$

Miara wyboru cechy, która najlepiej klasyfikuje zbiór uczący:

- **Entropia (Z)**

charakteryzuje zanieczyszczenie zbioru uczącego

$$Entropia(Z) = -p_+ \log_2 p_+ - p_- \log_2 p_-$$

Przykład: dla 14 przykładów (9 pozytywnych i 5 negatywnych)

$$Entropia([9+, 5-]) = -\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right) = 0,94$$

Miara wyboru cechy, która najlepiej klasyfikuje zbiór uczący:

- **Entropia (Z)**

charakteryzuje zanieczyszczenie zbioru uczącego

$$Entropia(Z) = -p_+ \log_2 p_+ - p_- \log_2 p_-$$

Przykład: dla 14 przykładów (9 pozytywnych i 5 negatywnych)

$$Entropia([9+, 5-]) = -\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right) = 0,94$$

Własności entropii:

- Entropia=0 - Wszystkie rozpatrywane przykłady należą do tej samej grupy
- Entropia=1 - przykładów pozytywnych jest tyle co negatywnych
- Entropia<1 - liczby pozytywnych i negatywnych przykładów nie są sobie równe

Miara wyboru cechy, która najlepiej klasyfikuje zbiór uczący:

- **Entropia (Z)**
- **Współczynnik Czerwińskiego**

Miara wyboru cechy, która najlepiej klasyfikuje zbiór uczący:

- **Entropia (Z)**
- **Współczynnik Czerwińskiego**

Obliczany z tablice wielodzielczych

Wsp. Czerwińskiego przyjmuje wartości ze zbioru 0,1.

Miara wyboru cechy, która najlepiej klasyfikuje zbiór uczący:

- **Entropia (Z)**
- **Współczynnik Czerwińskiego**

Obliczany z tablice wielodzielczych

Wsp. Czerwińskiego przyjmuje wartości ze zbioru 0,1.

Własności:

- wartość 1 - pełna zależność między cechami
- wartość 0 - brak zależności między cechami

Miara wyboru cechy, która najlepiej klasyfikuje zbiór uczący:

- **Entropia (Z)**
- **Współczynnik Czerwińskiego**
- **Współczynnik Ginięgo**

Jakość podziału (zanieczyszczenie) obliczana dla każdego z liści ze wzoru:

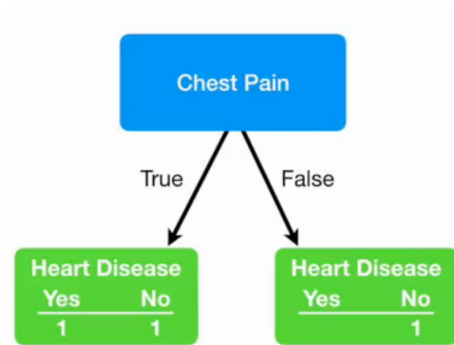
$$Wspczynnik_{Ginięgo} = 1 - \sum_{j=1}^n (p_j)^2$$

Drzewa decyzyjne - Przykład

Chest Pain	Good Blood Circulation	Blocked Arteries	Heart Disease
No	No	No	No
Yes	Yes	Yes	Yes
Yes	Yes	No	No
Yes	No	???	Yes
etc...	etc...	etc...	etc...

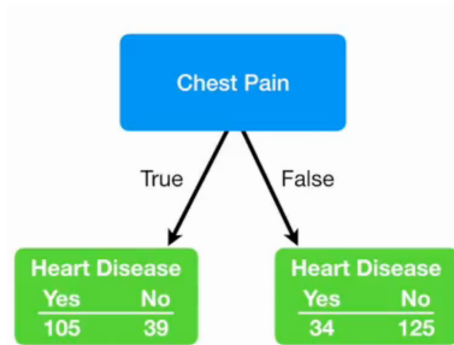
Drzewa decyzyjne - Przykład

Chest Pain	Good Blood Circulation	Blocked Arteries	Heart Disease
No	No	No	No
Yes	Yes	Yes	Yes
Yes	Yes	No	No
Yes	No	???	Yes
etc...	etc...	etc...	etc...



Drzewa decyzyjne - Przykład

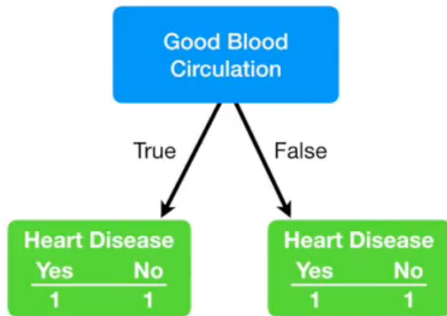
Chest Pain	Good Blood Circulation	Blocked Arteries	Heart Disease
No	No	No	No
Yes	Yes	Yes	Yes
Yes	Yes	No	No
Yes	No	???	Yes
etc...	etc...	etc...	etc...



Drzewa decyzyjne - Przykład

Chest Pain	Good Blood Circulation	Blocked Arteries	Heart Disease
No	No	No	No
Yes	Yes	Yes	Yes
Yes	Yes	No	No
Yes	No	???	Yes
etc...	etc...	etc...	etc...

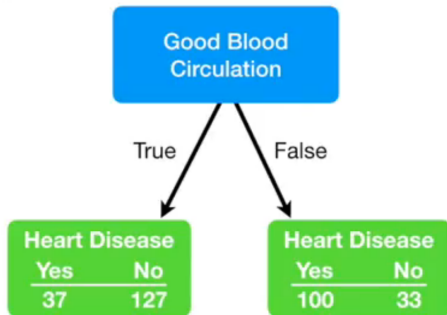
1.



Drzewa decyzyjne - Przykład

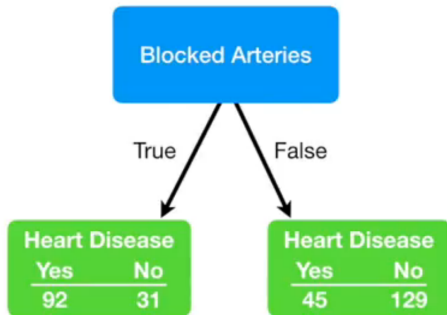
Chest Pain	Good Blood Circulation	Blocked Arteries	Heart Disease
No	No	No	No
Yes	Yes	Yes	Yes
Yes	Yes	No	No
Yes	No	???	Yes
etc...	etc...	etc...	etc...

1.

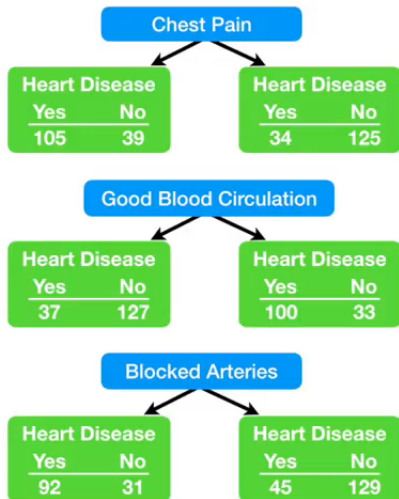


Drzewa decyzyjne - Przykład

Chest Pain	Good Blood Circulation	Blocked Arteries	Heart Disease
No	No	No	No
Yes	Yes	Yes	Yes
Yes	Yes	No	No
Yes	No	???	Yes
etc...	etc...	etc...	etc...



Drzewa decyzyjne - Przykład

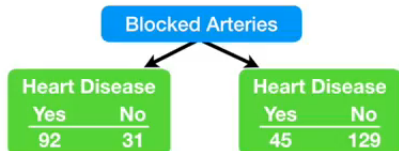
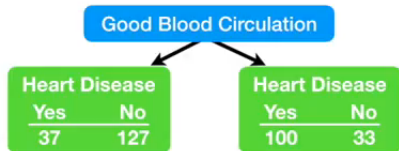
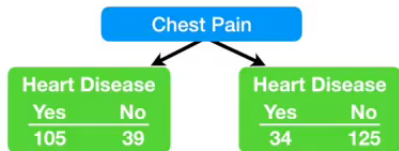


Drzewa decyzyjne - Przykład

$Wspczynnik_{Ginięgo} =$

$$1 - (\text{prawdopodobieństwo "Tak"})^2 - (\text{Prawdopodobieństwo "Nie"})^2$$

$$1 - \left(\frac{105}{105+39}\right)^2 - \left(\frac{39}{105+39}\right)^2 = 0,395$$



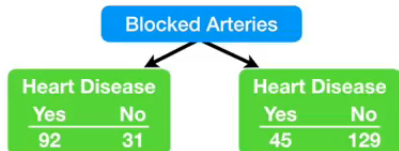
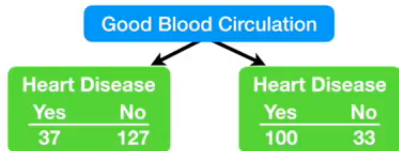
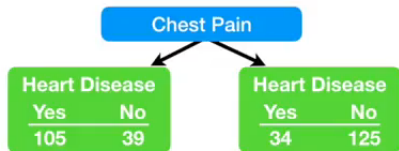
Drzewa decyzyjne - Przykład

Wspczynnik $Gini$ ego =

$$1 - (\text{prawdopodobieństwo "Tak"})^2 - (\text{Prawdopodobieństwo "Nie"})^2$$

$$1 - \left(\frac{105}{105+39}\right)^2 - \left(\frac{39}{105+39}\right)^2 = 0,395$$

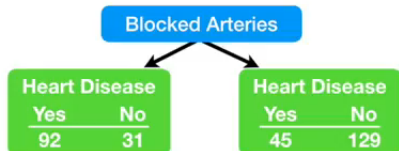
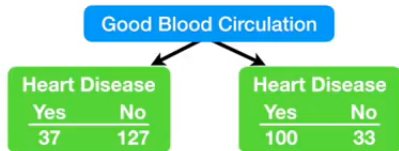
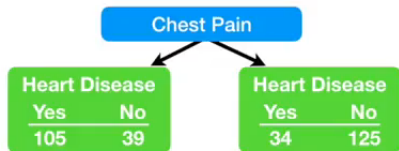
$$1 - \left(\frac{34}{34+125}\right)^2 - \left(\frac{125}{34+125}\right)^2 = 0,336$$



Drzewa decyzyjne - Przykład

Wspczynnik $Gini$ ego =

$$1 - (\text{prawdopodobieństwo "Tak"})^2 - (\text{Prawdopodobieństwo "Nie"})^2$$



$$1 - \left(\frac{105}{105+39}\right)^2 - \left(\frac{39}{105+39}\right)^2 = 0,395$$

$$1 - \left(\frac{34}{34+125}\right)^2 - \left(\frac{125}{34+125}\right)^2 = 0,336$$

Całkowita jakość podziału to
ważona średnia jakości podziału
obu liści:

$$\left(\frac{144}{144+159}\right)0,395 +$$

$$\left(\frac{159}{144+159}\right)0,336 = 0,364$$

Drzewa decyzyjne - Przykład

$Wspczynnik_{Ginięgo} =$

$$1 - (\text{prawdopodobieństwo "Tak"})^2 - (\text{Prawdopodobieństwo "Nie"})^2$$

Chest Pain	
Heart Disease	
Yes	No
105	39

Chest Pain	
Heart Disease	
Yes	No
34	125

0,395

0,336

0,364

Good Blood Circulation	
Heart Disease	
Yes	No
37	127

Good Blood Circulation	
Heart Disease	
Yes	No
100	33

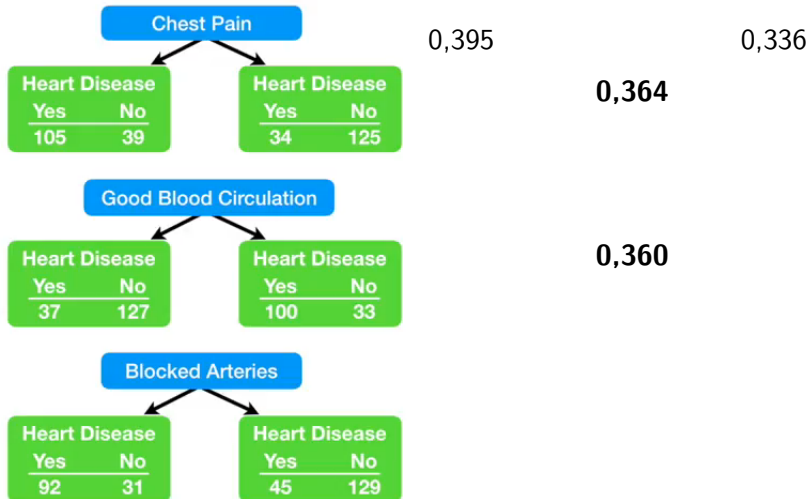
Blocked Arteries	
Heart Disease	
Yes	No
92	31

Blocked Arteries	
Heart Disease	
Yes	No
45	129

Drzewa decyzyjne - Przykład

Wspczynnik $Gini$ ego =

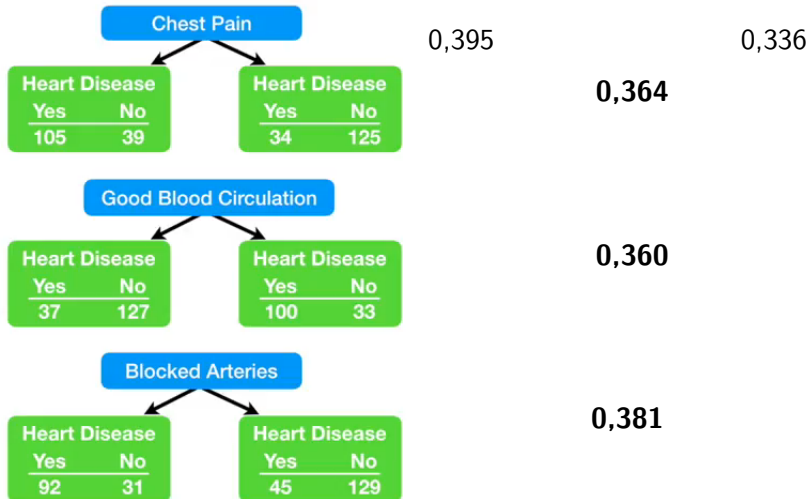
$$1 - (\text{prawdopodobieństwo "Tak"})^2 - (\text{Prawdopodobieństwo "Nie"})^2$$



Drzewa decyzyjne - Przykład

$Wspczynnik_{Ginięgo} =$

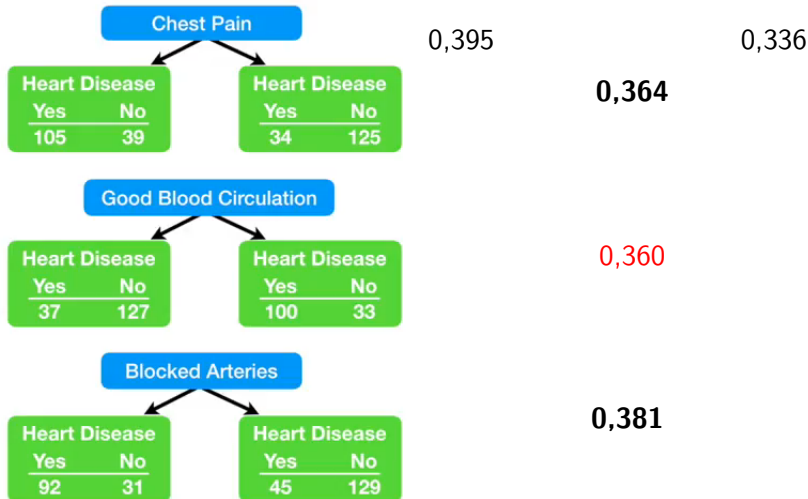
$$1 - (\text{prawdopodobieństwo "Tak"})^2 - (\text{Prawdopodobieństwo "Nie"})^2$$



Drzewa decyzyjne - Przykład

$Wspczynnik_{Ginięgo} =$

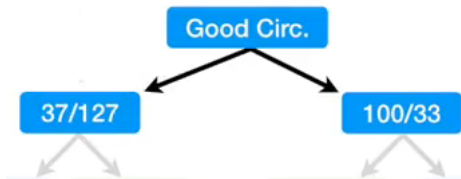
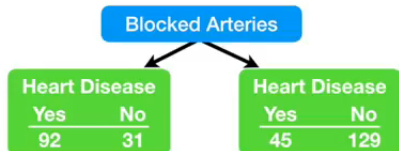
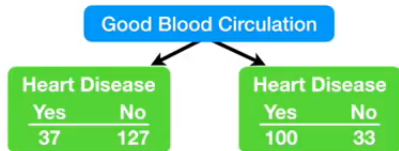
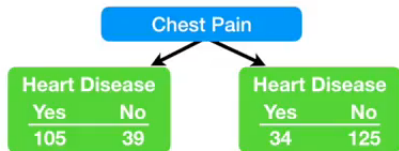
$$1 - (\text{prawdopodobieństwo "Tak"})^2 - (\text{Prawdopodobieństwo "Nie"})^2$$



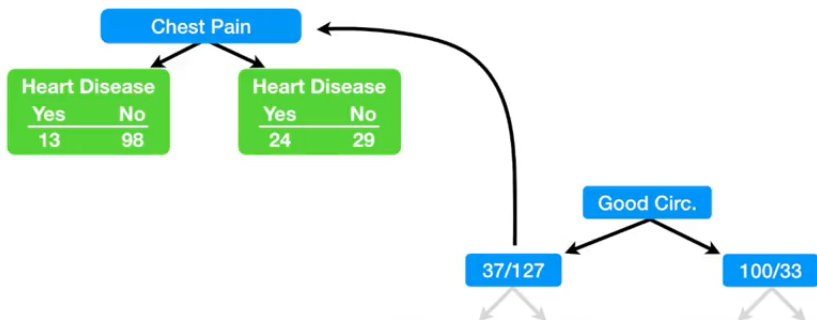
Drzewa decyzyjne - Przykład

$Wspczynnik_{Ginięgo} =$

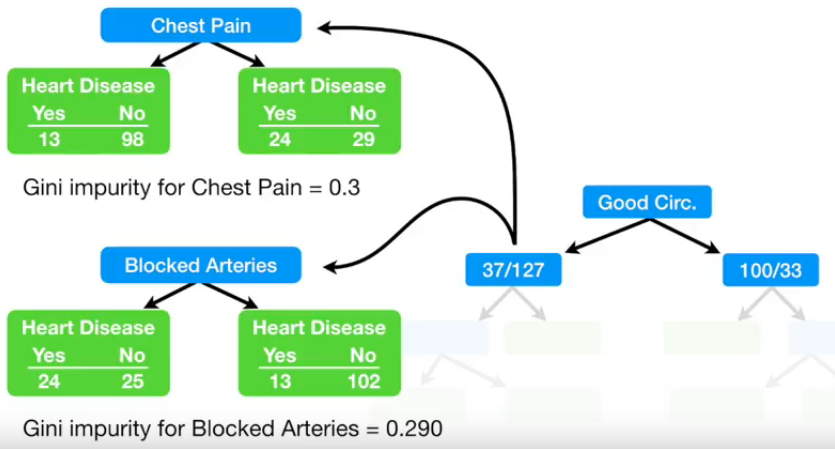
$$1 - (\text{prawdopodobieństwo "Tak"})^2 - (\text{Prawdopodobieństwo "Nie"})^2$$



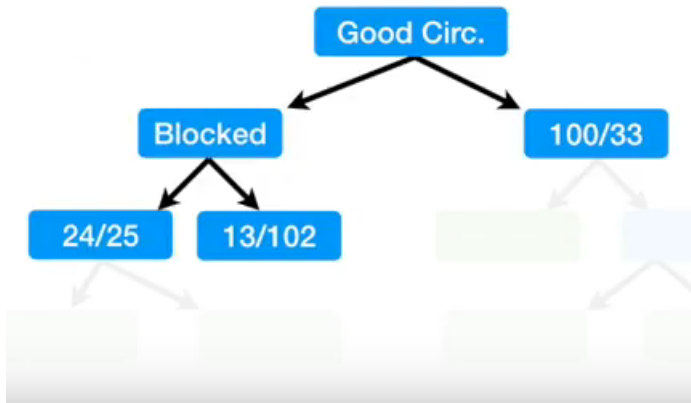
Drzewa decyzyjne - Przykład



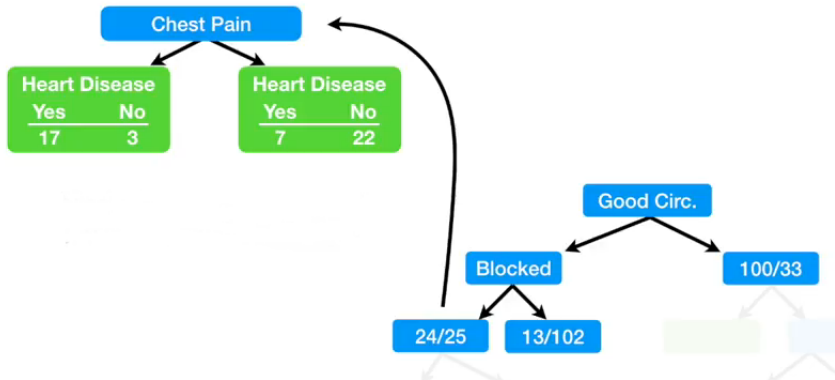
Drzewa decyzyjne - Przykład



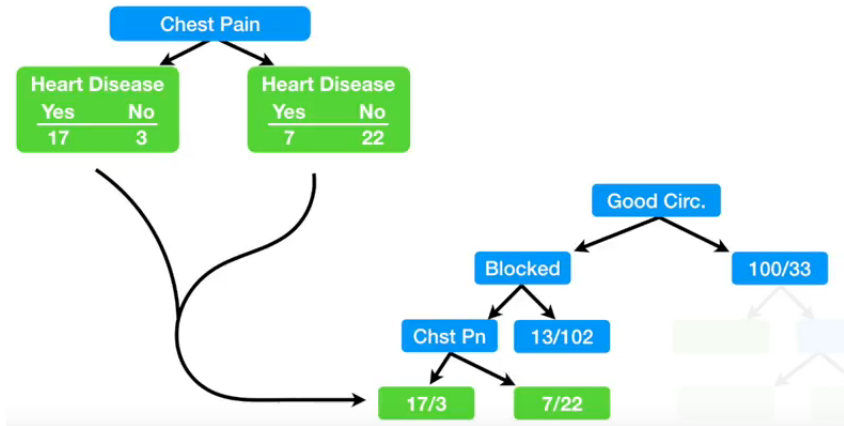
Drzewa decyzyjne - Przykład



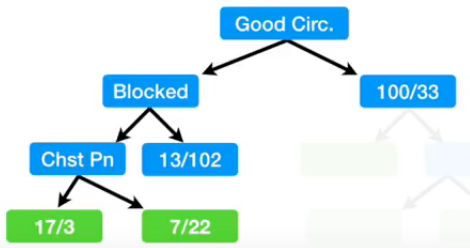
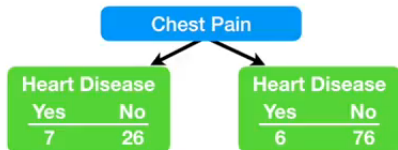
Drzewa decyzyjne - Przykład



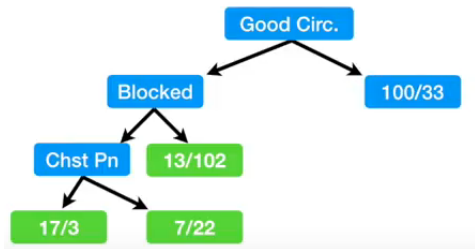
Drzewa decyzyjne - Przykład



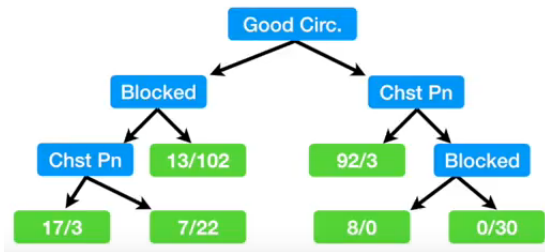
Drzewa decyzyjne - Przykład



Drzewa decyzyjne - Przykład



Drzewa decyzyjne - Przykład



Przykład - Dane numeryczne

Kroki algorytmu w przypadku danych numerycznych:

Weight	Heart Disease
220	Yes
180	Yes
225	Yes
190	No
155	No

Przykład - Dane numeryczne

Kroki algorytmu w przypadku danych numerycznych:

- 1 sortowanie

	Weight	Heart Disease
Lowest	155	No
	180	Yes
	190	No
	220	Yes
Highest	225	Yes

Przykład - Dane numeryczne

Kroki algorytmu w przypadku danych numerycznych:

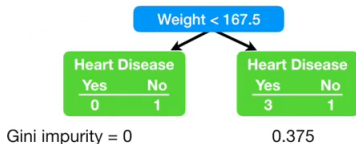
- 1 sortowanie
- 2 obliczenie średniej sąsiednich wartości

Weight	Heart Disease
155	No
167.5	
180	Yes
185	
190	No
205	
220	Yes
222.5	
225	Yes

Przykład - Dane numeryczne

Kroki algorytmu w przypadku danych numerycznych:

- 1 sortowanie
- 2 obliczenie średniej sąsiednich wartości
- 3 obliczenie jakości podziału dla każdej średniej wartości



Weight	Heart Disease
155	No
167.5	
180	Yes
185	
190	No
205	
220	Yes
222.5	
225	Yes

Przykład - Dane numeryczne

Kroki algorytmu w przypadku danych numerycznych:

- 1 sortowanie
- 2 obliczenie średniej sąsiednich wartości
- 3 obliczenie jakości podziału dla każdej średniej wartości

Weight	Heart Disease	
155	No	
167.5		→ Gini impurity = 0.3
180	Yes	
185		→ Gini impurity = 0.47
190	No	
205		→ Gini impurity = 0.27
220	Yes	
222.5		→ Gini impurity = 0.4
225	Yes	

Zalety:

- Efektywny pamięciowo zapis wiedzy,
- Cechy ilościowe, binarne, rankingowe, wielokrotnego wyboru,
- Nie są czułe na błędy i braki danych,
- Szybki proces klasyfikacji,
- Zapis intuicyjny dla człowieka.

Wady:

- Brak możliwości uczenia adaptacyjnego,
- Możliwe duże rozmiary i poziom skomplikowania,
- Niewielka zmiana wymaga uruchomienia algorytmu budowy od początku,
- Zbyt duże dopasowanie do zbioru uczącego, duża niedokładność.

- 1 Wprowadzenie do klasyfikacji
- 2 kNN - k Nearest Neighbours
- 3 Regresja logistyczna
- 4 SVM
- 5 Naive Bayes
- 6 Drzewa decyzyjne
- 7 **Random Forest**

Idea:

- opary na zbiorze drzew decyzyjnych
- losowanie (ze zwracaniem) podzbioru z całego zbioru uczącego
- losowanie (ze zwracaniem) cech

Random Forest - Algorytm uczenia

- 1 Losowanie z n-elementowej próby n wektorów obserwacji,

Original Dataset

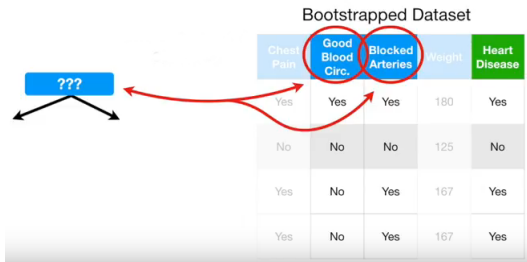
Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
No	No	No	125	No
Yes	Yes	Yes	180	Yes
Yes	Yes	No	210	No
Yes	No	Yes	167	Yes

Bootstrapped Dataset

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	Yes	Yes	180	Yes
No	No	No	125	No
Yes	No	Yes	167	Yes
Yes	No	Yes	167	Yes

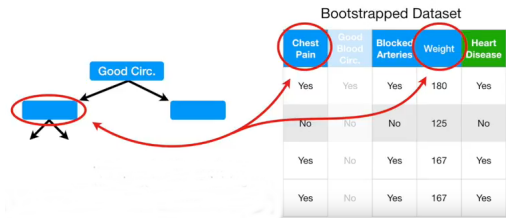
Random Forest - Algorytm uczenia

- 1 Losowanie z n -elementowej próby n wektorów obserwacji,
- 2 Tworzenie drzewa decyzyjnego, gdzie dla każdego węzła losowane jest m z p cech (potem k z m , gdzie $p > m > k$)



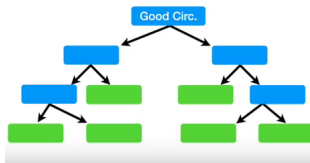
Random Forest - Algorytm uczenia

- 1 Losowanie z n -elementowej próby n wektorów obserwacji,
- 2 Tworzenie drzewa decyzyjnego, gdzie dla każdego węzła losowane jest m z p cech (potem k z m , gdzie $p > m > k$)



Random Forest - Algorytm uczenia

- 1 Losowanie z n -elementowej próby n wektorów obserwacji,
- 2 Tworzenie drzewa decyzyjnego, gdzie dla każdego węzła losowane jest m z p cech (potem k z m , gdzie $p > m > k$)

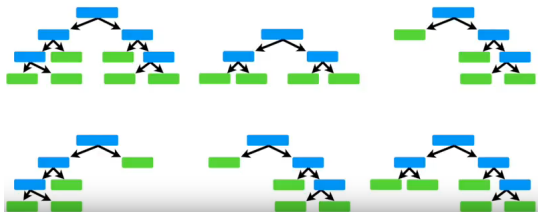


Bootstrapped Dataset

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	Yes	Yes	180	Yes
No	No	No	125	No
Yes	No	Yes	167	Yes
Yes	No	Yes	167	Yes

Random Forest - Algorytm uczenia

- 1 Losowanie z n -elementowej próby n wektorów obserwacji,
- 2 Tworzenie drzewa decyzyjnego, gdzie dla każdego węzła losowane jest m z p cech (potem k z m , gdzie $p > m > k$)
- 3 Powrót do kroku 1. i stworzenie kolejnych losowych drzew.

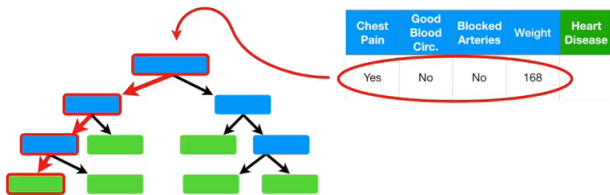


Random Forest - Klasyfikacja

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	No	No	168	

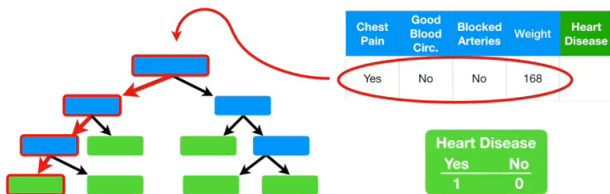
Random Forest - Klasyfikacja

- 1 Klasyfikacja obserwacji przez jedno drzewo,



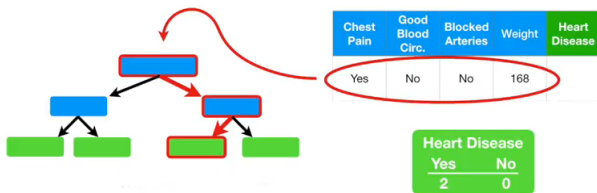
Random Forest - Klasyfikacja

- 1 Klasyfikacja obserwacji przez jedno drzewo,
- 2 Zapisanie wyniku,



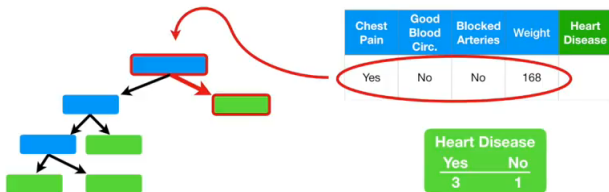
Random Forest - Klasyfikacja

- 1 Klasyfikacja obserwacji przez jedno drzewo,
- 2 Zapisanie wyniku,
- 3 Powtórzenie kroku 1. dla wszystkich drzew,



Random Forest - Klasyfikacja

- 1 Klasyfikacja obserwacji przez jedno drzewo,
- 2 Zapisanie wyniku,
- 3 Powtórzenie kroku 1. dla wszystkich drzew,



Random Forest - Klasyfikacja

- 1 Klasyfikacja obserwacji przez jedno drzewo,
- 2 Zapisanie wyniku,
- 3 Powtórzenie kroku 1. dla wszystkich drzew,
- 4 Wybór klasy, która wystąpiła najczęściej.

Chest Pain	Good Blood Circ.	Blocked Arteries	Weight	Heart Disease
Yes	No	No	168	YES

Heart Disease	
Yes	No
5	1

Cechy:

- jest dokładny
- jest skuteczny na dużych bazach danych
- szacuje dane w przypadku ich braku
- nie wymaga wiedzy eksperckiej
- nie jest podatny na overfitting
- określa przydatność cech
- umożliwia klasteryzację
- umożliwia oszacowanie błędu klasyfikacji