

Problemy uczenia sieci neuronowych

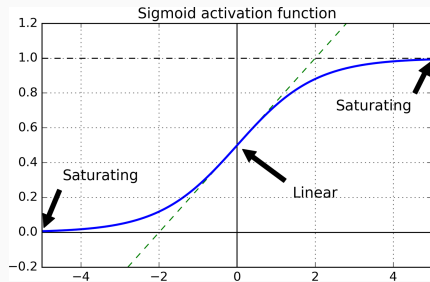
dr inż. Sebastian Ernst

Przedmiot: Uczenie Maszynowe

Uczenie sieci neuronowych

Problem znikających/eksplodujących gradientów

- Algorytm propagacji wstecznej przenosi gradienty błędów idąc od warstwy wyjściowej do wejściowej (czyli wstecz)
- Często wartości gradientów spadają przy przechodzeniu do coraz niższych warstw – problem *znikających gradientów*
- Czasami jest wręcz przeciwnie – problem *eksplodujących gradientów*
- W 2010 Glorot i Bengio odkryli związek między niestabilnością gradientów a używaniem funkcji sigmoidalnej oraz popularnego wówczas sposobu inicjalizacji wag (rozkład normalny $\mu = 0$, $\sigma = 1$)



Inicjalizacja Glorota i He

- Teoretycznie: wariancja wejść każdej warstwy musi być równa wariancji jej wyjść
- W praktyce: wagi połączeń warstwy powinny być równe:
 - rozkładowi normalnemu o $\mu = 0$, $\sigma^2 = \frac{1}{fan_{avg}}$ lub
 - rozkładowi jednostajnemu pomiędzy $-r$ a $+r$ przy $r = \sqrt{\frac{3}{fan_{avg}}}$

Metoda	Funkcje aktywacji	σ^2
Glorot	brak, tanh, sigmoid, softmax	$1/fan_{avg}$
He	ReLU & co.	$2/fan_{in}$
LeCun	SELU	$1/fan_{in}$

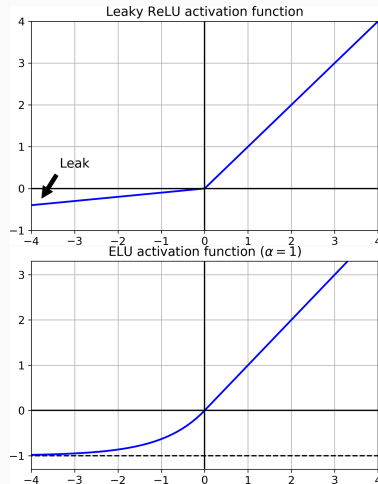
Nienasycające funkcje aktywacji

- Przy ReLU: problem „umierających” neuronów
 - suma ważona wejść zawsze ujemna
- Dwa rozwiązania:
 - leaky ReLU:

$$\text{LeakyReLU}_{\alpha}(z) = \max(\alpha z, z)$$

- Exponential Linear Unit:

$$\text{ELU}_{\alpha}(z) = \begin{cases} \alpha(e^z - 1) & \text{dla } z < 0 \\ z & \text{dla } z \geq 0 \end{cases}$$



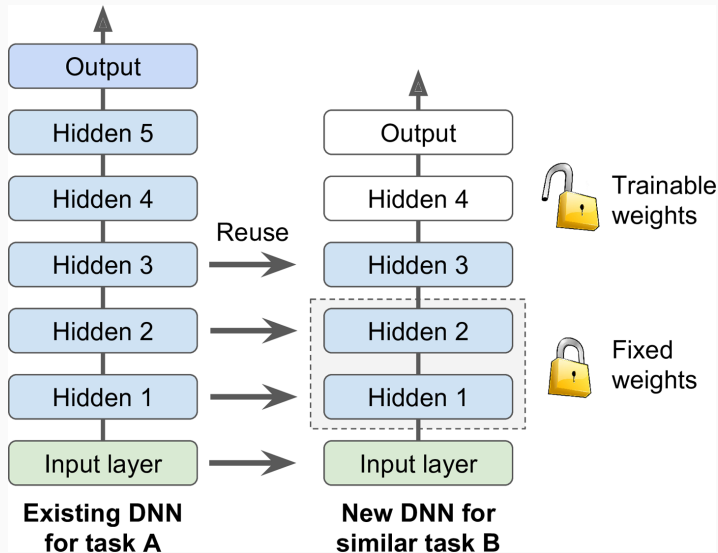
Normalizacja wsadów

- Stosowanie inicjalizacji He i ELU/ReLU & co. zmniejsza ryzyko niestabilności gradientów na początku uczenia, ale problem może pojawić się później.
- Technika *Batch Normalization* polega na dodaniu operacji tuż przed lub po funkcji aktywacji każdej warstwy ukrytej – wycentrowanie i normalizacja wejścia + skalowanie i przesunięcie wyniku.

```
model = keras.models.Sequential([
    keras.layers.Flatten(input_shape=[28, 28]),
    keras.layers.BatchNormalization(),
    keras.layers.Dense(300, activation="relu"),
    keras.layers.BatchNormalization(),
    keras.layers.Dense(100, activation="relu"),
    keras.layers.BatchNormalization(),
    keras.layers.Dense(10, activation="softmax")
])
```

Uczenie transferowe

Wykorzystanie wytrenowanych warstw



Uczenie transferowe w Keras

```
model_A = keras.models.load_model("my_model_A.h5")
model_B_on_A = keras.models.Sequential(model_A.layers[:-1])
model_B_on_A.add(keras.layers.Dense(1, activation="sigmoid"))
```

Uwaga: trening nowego modelu będzie modyfikował wagi warstw również w model_A – aby tego uniknąć, klonujemy model:

```
model_A_clone = keras.models.clone_model(model_A)
model_A_clone.set_weights(model_A.get_weights())
model_B_on_A = keras.models.Sequential(model_A_clone.layers[:-1])
model_B_on_A.add(keras.layers.Dense(1, activation="sigmoid"))
```

Wstępne przyuczenie warstwy wyjściowej

Początkowo, warstwa wyjściowa może zwracać wartości dalekie od ideału, a propagacja zmian może zniszczyć wagi przeniesione z pierwotnego modelu. W tym celu blokujemy możliwość modyfikacji wag przeniesionych warstw.

```
for layer in model_B_on_A.layers[:-1]:  
    layer.trainable = False
```

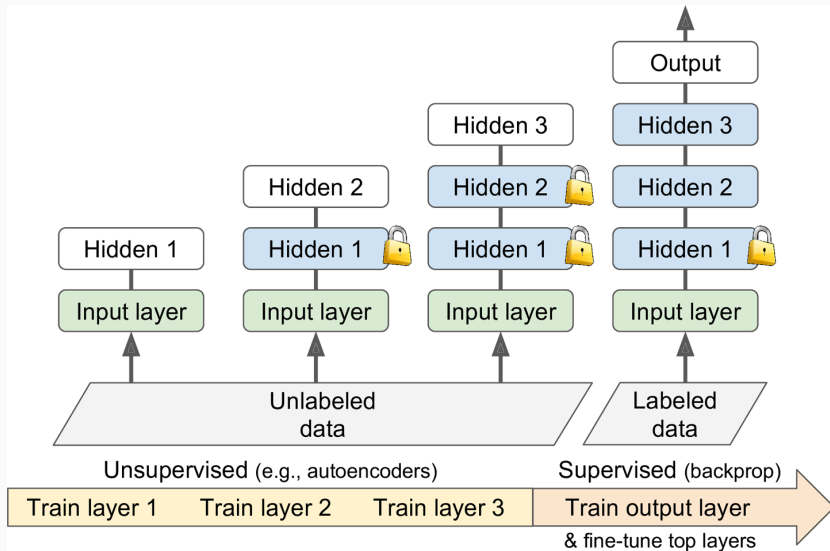
```
model_B_on_A.compile(loss="binary_crossentropy",  
                      optimizer=keras.optimizers.SGD(learning_rate=1e-3),  
                      metrics=["accuracy"])
```

```
history = model_B_on_A.fit(X_train_B, y_train_B, epochs=4,  
                           validation_data=(X_valid_B, y_valid_B))
```

Uczenie nowego modelu

```
for layer in model_B_on_A.layers[:-1]:  
    layer.trainable = True  
  
model_B_on_A.compile(loss="binary_crossentropy",  
                     optimizer=keras.optimizers.SGD(learning_rate=1e-3),  
                     metrics=["accuracy"])  
history = model_B_on_A.fit(X_train_B, y_train_B, epochs=16,  
                          validation_data=(X_valid_B, y_valid_B))
```

Nienadzorowane uczenie wstępne



Inne (szybsze) algorytmy optymalizacji

Przyspieszanie procesu uczenia

Dotychczas znamy **cztery sposoby** na przyspieszenie uczenia:

1. dobra strategia inicjalizacji wag połączeń
2. dobra funkcja aktywacji
3. korzystanie z normalizacji wsadów (BN)
4. wykorzystanie części wstępnie przyuczonej sieci

Piąty sposób: wykorzystanie algorytmu optymalizacji innego niż gradientowy (*gradient descent*).

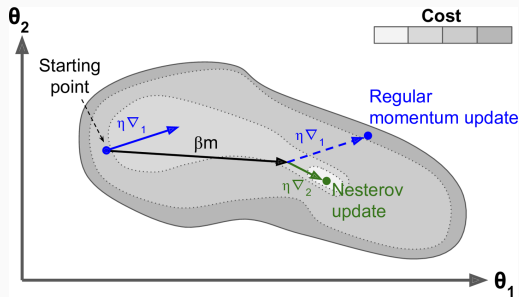
Algorytm gradientowy z pędem

- Klasyczny algorytm gradientowy nie uwzględnia wcześniejszych wartości gradientów – bierze pod uwagę tylko wartość chwilową.
- Algorytm z pędem dodaje *wektor pędu*, a więc wartość gradientu przekłada się na *przyspieszenie* a nie *prędkość*:
 1. $\mathbf{m} \leftarrow \beta \mathbf{m} - \eta \nabla_{\Theta} J(\Theta)$
 2. $\Theta \leftarrow \Theta + \mathbf{m}$
- Np. dla $\beta = 0.9$ prędkość może wzrosnąć 10x w stosunku do zwykłego algorytmu gradientowego.
- Przyspieszenie widoczne szczególnie przy różnych skalach wejść (efekt wydłużonej misy).
- W Keras wystarczy ustawić parametr momentum:

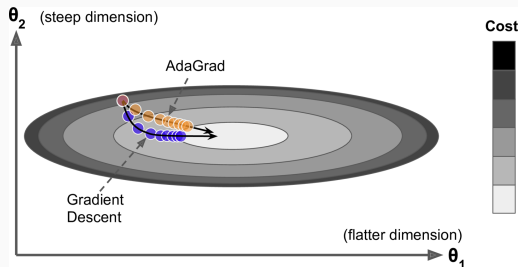
```
optimizer = keras.optimizers.SGD(lr=0.001, momentum=0.9)
```

Nesterov Accelerated Gradient (NAG)

- Zaproponowany w 1983 przez Nesterowa, prawie zawsze szybszy niż „zwykły” gradient z pędem.
- Patrzy „z wyprzedzeniem”:
 1. $\mathbf{m} \leftarrow \beta \mathbf{m} - \eta \nabla_{\Theta} J(\Theta + \beta \mathbf{m})$
 2. $\Theta \leftarrow \Theta + \mathbf{m}$
- W Keras dodajemy argument `nesterov=True`.



- Koryguje kierunek gradientu w stronę globalnego minimum, skalując wektor gradientów wedle „stromości”:
 1. $\mathbf{s} \leftarrow \mathbf{s} + \nabla_{\Theta} J(\Theta) \otimes \nabla_{\Theta} J(\Theta)$
 2. $\Theta \leftarrow \Theta - \eta \nabla_{\Theta} J(\Theta) \oslash \sqrt{\mathbf{S} + \epsilon}$
- Problem przy sieciach głębokich: zatrzymuje się za wcześnie.



- Rozwinięcie AdaGrad – bierze pod uwagę tylko gradienty z ostatnich iteracji (a nie od początku uczenia)
- Patrzy „z wyprzedzeniem”:

1. $\mathbf{s} \leftarrow \beta \mathbf{s} + (1 - \beta) \nabla_{\Theta} J(\Theta) \otimes \nabla_{\Theta} J(\Theta)$
2. $\Theta \leftarrow \Theta - \eta \nabla_{\Theta} J(\Theta) \oslash \sqrt{\mathbf{s} + \varepsilon}$

- Typowa wartość β : 0.9.
- W Keras:

```
keras.optimizers.RMSprop(lr=0.001, rho=0.9)
```

Adam

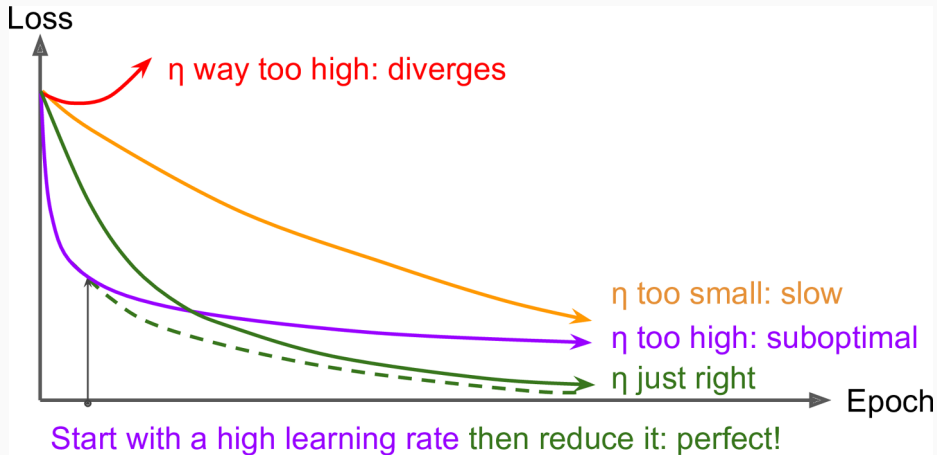
- Adam = *adaptive moment estimation*
- Łączy cechy gradientu z pędem oraz RMSProp:
 - tak jak gradient z pędem, pamięta wykładniczo zanikającą średnią poprzednich gradientów,
 - tak jak RMSProp, pamięta wykładniczo zanikającą średnią kwadratów poprzednich gradientów.

1. $\mathbf{m} \leftarrow \beta_1 \mathbf{m} - (1 - \beta_1) \nabla_{\Theta} J(\Theta)$
2. $\mathbf{s} \leftarrow \beta_2 \mathbf{s} + (1 - \beta_2) \nabla_{\Theta} J(\Theta) \otimes \nabla_{\Theta} J(\Theta)$
3. $\hat{\mathbf{m}} \leftarrow \frac{\mathbf{m}}{1 - \beta_1^t}$
4. $\hat{\mathbf{s}} \leftarrow \frac{\mathbf{s}}{1 - \beta_2^t}$
5. $\Theta \leftarrow \Theta + \eta \hat{\mathbf{m}} \oslash \sqrt{\hat{\mathbf{s}} + \varepsilon}$

```
keras.optimizers.Adam(lr=0.001  
                        beta_1=0.9,  
                        beta_2=0.999)
```

Harmonogramowanie kroku uczenia

Długość kroku a efekty uczenia



Harmonogramowanie z potęgowaniem (*power scheduling*)

- Krok w funkcji numeru iteracji t : $\eta(t) = \eta_0 / (1 + t/s)^c$
- Hiperparametry:
 - η_0 – krok początkowy
 - c – wykładnik potęgi
 - s – liczba kroków
- Po wykonaniu s kroków, η spada do $\eta_0/2$, potem do $\eta_0/3$, $\eta_0/4$, itd.

Harmonogramowanie wykładnicze (*exponential scheduling*)

- $\eta(t) = \eta_0 \cdot 0.1^{t/s}$
- η maleje 10-krotnie co s kroków
- nie wyhamowuje tak jak harmonogramowanie z potęgowaniem

Harmonogramowanie ze stałymi wartościami (*piecewise constant scheduling*)

- Sekwencja par (długość kroku, liczba iteracji)
- Wymaga ręcznego strojenia całej sekwencji

Harmonogramowanie oparte o wydajność (*performance scheduling*)

- Zmierz błąd walidacyjny co N kroków (tak jak *early stopping*)
- Zmniejsz długość kroku o λ gdy nie ma poprawy

- Zaproponowane w 2018 przez Liesliego Smitha
- Zaczynamy od η_0 , zwiększamy liniowo do η_1
- Potem obniżamy z powrotem do η_0 w kolejnej części procesu uczenia
- Pod koniec uczenia zmniejszamy krok, nadal liniowo, ale o kilku rzędów wielkości
- η_1 ustawiamy tak jak zwykle statyczny krok; $n_0 \approx \eta_1/10$
- Jeżeli korzystamy z pędu, zaczynamy od wysokiej wartości (np. 0.95) i w pierwszej fazie nieco ją obniżamy (np. 0.85), potem podnosimy

Harmonogramowanie w Keras

- Harmonogramowanie z potęgowaniem:

```
optimizer = keras.optimizers.SGD(lr=0.01, decay=1e-4)
```

- Harmonogramowanie wykładnicze i ze stałymi wartościami możemy uzyskać definiując funkcję zaniku i dołączając jako callback do procesu uczenia:

```
def exponential_decay_fn(epoch):  
    return 0.01 * 0.1**(epoch / 20)
```

```
lr_scheduler =  
    keras.callbacks.LearningRateScheduler(exponential_decay_fn)  
history = model.fit(X_train_scaled, y_train, epochs=n_epochs,  
                    validation_data=(X_valid_scaled, y_valid),  
                    callbacks=[lr_scheduler])
```

Regularyzacja: unikanie przeuczenia

Regularyzacja: co już znamy

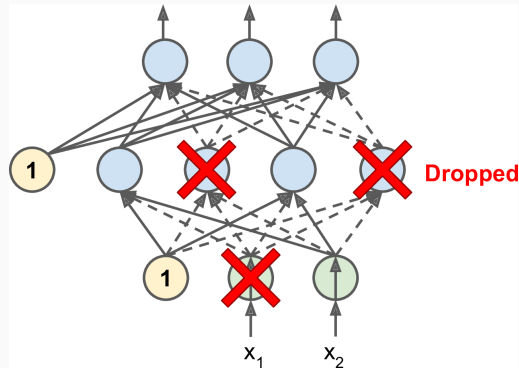
- *Early stopping* – przerywanie uczenia przy braku poprawy
- Normalizacja wsadów – ma takie działanie, mimo że powstała aby unikać niestabilności gradientów

Regularyzacja ℓ_1 i ℓ_2

- Regularyzacja ℓ_1 – ograniczenie wag połączeń w sieci
- Regularyzacja ℓ_2 – model rzadki (wiele wag ustawionych na 0)
- W Keras dodajemy do warstwy argument `kernel_regularizer`, przekazując jeden z obiektów:
 - `keras.regularizers.l1`
 - `keras.regularizers.l2`
 - `keras.regularizers.l1_l2`

Dropout

- Metoda zaproponowana i rozwinięta w 2012 i 2014 roku
- W każdym kroku, każdy neuron z prawdopodobieństwem p zostanie opuszczony (ang. *dropped out*)
- Opuszczony neuron jest ignorowany, ale tylko w bieżącym kroku uczenia
- Hiperparametr p (współczynnik opuszczenia – *dropout rate*) ustawiamy na 10–50% (zakres węższy w sieciach rekurencyjnych i konwolucyjnych)



Dropout w Keras

```
model = keras.models.Sequential([
    keras.layers.Flatten(input_shape=[28, 28]),
    keras.layers.Dropout(rate=0.2),
    keras.layers.Dense(300, activation="elu",
                        kernel_initializer="he_normal"),
    keras.layers.Dropout(rate=0.2),
    keras.layers.Dense(100, activation="elu",
                        kernel_initializer="he_normal"),
    keras.layers.Dropout(rate=0.2),
    keras.layers.Dense(10, activation="softmax")
])
```