

드롭아웃과 Yolo

연준모

1. 옵티나이저와 드롭아웃
2. Yolo 아키텍처 알아보기
3. 수식으로 Yolo 열어보기

1. 옵티마이저와 드롭아웃 옵티마이저?

ONECLICK AI

옵티마이저란?

손실함수 값을 최소화하는 최적의 파라미터를 찾는 알고리즘

일반적으로, 경사하강법을 기반으로 한다

4. 기본 개념 알고가기 경사하강법

ONECLICK AI



오차를 최소화 하는 방향으로
가중치, 편향을 업데이트 한다

$$w_{\text{new}} = w_{\text{old}} - \eta \nabla J(w)$$

새로운 위치 = 현재 위치 + 움직인 방향과 거리

w_{new} : 업데이트 된 가중치. 다음 학습에서 사용된다.

w_{old} : 현재 가중치. 출발하는 값.

η : 학습률. 경사를 얼마나 많이 이동할 지 정한다.

$\nabla J(w)$: 현재 위치에서 기울기가 가장 가파른 방향.

여기에 사용되는 기울기를 역전파가 알려준다

1. 옵티마이저와 드롭아웃 확률적 경사 하강법 SGD?

ONECLICK AI

Stochastic Gradient Descent SGD?

확률적 경사하강법

기본적인 경사 하강법은 전체 학습 데이터를 사용해서 한 번에 기울기를 계산

데이터가 매우 클 경우 계산량이 엄청나게 많아져서 느려 터진다.

=> 전체 데이터가 아닌, 무작위로 선택된 일부 데이터(미니배치 'mini batch')를 사용해 계산, 파라미터 업데이트
속도 증가, 노이즈가 있을 수 도 있는 데이터를 거름으로써 지역 최솟값에 빠질 위험을 줄여준다

$$W \leftarrow W - \eta \frac{\partial L_i(W)}{\partial W}$$

미니배치에 속하는 데이터 $x^{(i)}$ 와 정답 $y^{(i)}$ 에 대한 손실 함수를 $L_i(W)$ 라 할 때 식이 위와 같이 나온다

딥러닝은 전체 노드로, 손실 계산 + 가중치 업데이트만 미니 배치에 있는 노드만으로

1. 옵티마이저와 드롭아웃 모멘텀 Momentum?

ONECLICK AI

Momentum?

‘관성’의 개념 도입

언덕 내려오는 공 점점 빨라지는 것 처럼

모멘텀은 이전의 기울기 방향을 현재 기울기 계산에 반영.

기울기가 같은 방향으로 계속 이동할 때 더 빠르게 학습, 기울기 방향이 자주 바뀌는 경우 진동을 줄여줌

모멘텀을 나타내는 변수 v 와 모멘텀 계수 γ (보통 0.9와 같은 값을 사용)를 도입

1. 모멘텀 누적 v_t : 이전 스텝의 모멘텀 γv_{t-1} 과 현재 스텝의 기울기 $\eta \nabla_W L(W)$ 를 합산

$$v_t = \gamma v_{t-1} + \eta \nabla_W L(W)$$

2. 가중치 업데이트: 누적된 모멘텀 값으로 가중치를 업데이트

$$W \leftarrow W - v_t$$

1. 옵티마이저와 드롭아웃 AdaGrad?

ONECLICK AI

Adaptive Gradient?

파라미터마다 서로 다른 학습률을 적용

언덕 내려오는 공 점점 빨라지는 것 처럼

학습 과정에서 변화가 적었던 파라미터는 더 큰 학습률로 업데이트하고,
변화가 많았던 파라미터는 작은 학습률로 업데이트

과거의 모든 기울기 제곱값을 합산하는 변수 G 를 사용.

1. 기울기 제곱 누적 G_t : 현재 스텝의 기울기 제곱($\nabla_W L(W)$)²을 이전 값에 더한다

$$G_t \leftarrow G_{t-1} + (\nabla_W L(W))^2$$

2. 가중치 업데이트: 학습률 η 를 G_t 의 제곱근으로 나누어 업데이트 강도를 조절. ϵ 은 분모가 0이 되는 것을 방지하는 작은 값(e.g., 1e-8)

$$W \leftarrow W - \frac{\eta}{\sqrt{G_t + \epsilon}} \nabla_W L(W)$$

학습률이 급격히 감소하는 문제가 있다 해결 위해 RMSProp 가 있다

1. 옵티마이저와 드롭아웃 Adam?

ONECLICK AI

Adaptive Moment Estimation?

모멘텀과 RMSProp의 장점을 결합

기울기의 지수 이동 평균(1차 모멘텀)과 기울기 제곱의 지수 이동 평균(2차 모멘텀)을 함께 사용하여 파라미터를 업데이트

모멘텀을 위한 1차 모멘텀 추정값 m_t 와 RMSProp을 위한 2차 모멘텀 추정값 v_t 를 사용하며, 각각의 감쇠율로 β_1, β_2 를 사용

1. 1차 모멘텀 계산 m_t :

$$m_t \leftarrow \beta_1 m_{t-1} + (1 - \beta_1) \nabla_W L(W)$$

2. 2차 모멘텀 계산 v_t :

$$v_t \leftarrow \beta_2 v_{t-1} + (1 - \beta_2) (\nabla_W L(W))^2$$

1. 옵티마이저와 드롭아웃 Adam?

ONECLICK AI

3. 편향 보정 (Bias Correction): 학습 초반에 m_t 와 v_t 가 0에 가깝게 추정되는 것을 보정

$$\widehat{m}_t = \frac{m_t}{1 - \beta_1^t}$$

$$\widehat{v}_t = \frac{v_t}{1 - \beta_2^t}$$

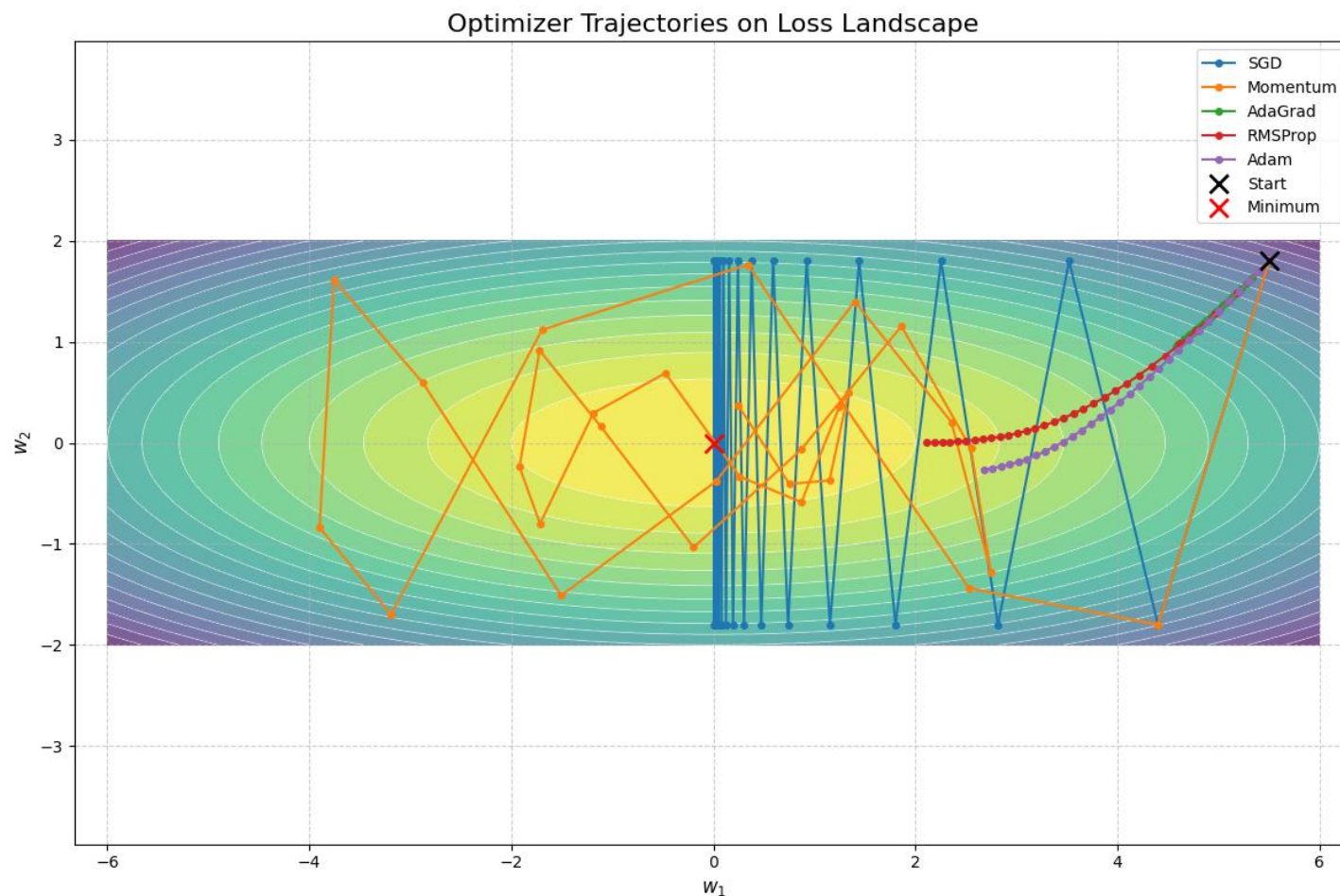
4. 가중치 업데이트: 보정된 모멘텀 값들을 사용하여 최종 업데이트를 수행

$$W \leftarrow W - \eta \frac{\widehat{m}_t}{\sqrt{\widehat{v}_t} + \epsilon}$$

Adam은 방향과 학습률을 모두 적응적으로 조절하여 대부분의 경우 좋은 성능을 보인다
그래서 요즘도 잘 쓰이는 그런거다 Adam, AdamW 이렇게 많이 쓰인다

1. 옵티마이저와 드롭아웃 비교?

ONECLICK AI



SGD와 Momentum은 발산해버림(너무 높은 학습율)

1. 옵티마이저와 드롭아웃 드롭아웃?

ONECLICK AI

Dropout?

신경망의 과적합(overfitting)을 방지하기 위한 강력한 정규화(regularization) 기법
한 줄로 설명하자면, 훈련 중에 뉴런을 무작위로 끈다

신경망의 한 계층(layer)이 수행하는 연산은 활성화 함수까지 해서 이렇게 나온다

$$A^{(l)} = f(W^{(l)}x^{(l-1)} + b^{(l)})$$

$W^{(l)}x^{(l-1)}$ 이게 행렬곱셈인데,

이 연산이 바로 신경망을 통해 데이터의 특징이 변환되는 과정의 본질
하드마르곱 + 마스킹으로 구현

입력 벡터: $x^{(l-1)}$

가중치 행렬: $W^{(1)}$

편향 벡터: $b^{(1)}$

활성화 함수: f

1. 옵티나이저와 드롭아웃 드롭아웃?

ONECLICK AI

Dropout?

1. 마스크 벡터 생성

드롭아웃을 적용할 활성화 벡터 $x^{(l)}$ 와 동일한 차원의 마스크 벡터 $d^{(l)}$ 를 생성

이 벡터의 각 원소는 뉴런을 유지할 확률 p 를 따르는 베르누이 분포에서 샘플링 된다

즉, 각 원소는 확률 p 로 1이 되고, 확률 $1-p$ 로 0이 된다

베르누이 분포가 궁금하다면?

예시로, $p=0.8$ 이고 $x^{(l)}$ 가 4차원 벡터일 때, 마스크 $d^{(l)}$ 는 $\begin{pmatrix} 1 \\ 0 \\ 1 \\ 1 \end{pmatrix}$ 와 같이 생성될 수 있다

살아남을 확률이 80%

2. 마스크 적용

벡터 $x^{(l)}$ 에 마스크 벡터 $d^{(l)}$ 를 하드마르 곱 \odot 하여 드롭아웃이 적용된 새로운 활성화 벡터 $\widetilde{x^{(l)}}$ 를 만든다

$$\widetilde{a^{(l)}} = a^{(l)} \odot d^{(l)}$$

$$\begin{pmatrix} 0.8 \\ 0.2 \\ 0.5 \\ 0.9 \end{pmatrix} \odot \begin{pmatrix} 1 \\ 0 \\ 1 \\ 1 \end{pmatrix} = \begin{pmatrix} 0.8 \\ 0.0 \\ 0.5 \\ 0.9 \end{pmatrix}$$

마스크된 벡터 $\widetilde{x^{(l)}}$ 가 다음 계층의 입력으로 전달

일부러 죽은 노드를 만들어 낸다

1. 옵티마이저와 드롭아웃 드롭아웃?

ONECLICK AI

Dropout?

? 과정 늘어났는데 그러면 시간 더 걸리고 느려지는거 아님?

위에서 설명한 하드마르 곱은 대각 행렬을 이용한 행렬 곱셈으로 완벽하게 치환할 수 있다.

$$d^{(l)} = \begin{pmatrix} d_1 \\ d_2 \\ \vdots \\ d_n \end{pmatrix} \rightarrow D^{(l)} = \begin{pmatrix} d_1 & 0 & \cdots & 0 \\ 0 & d_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & d_n \end{pmatrix}$$

이제 마스크 연산은 다음과 같이 표현된다

$$\widetilde{x^{(l)}} = x^{(l)} \odot d^{(l)} = D^{(l)}x^{(l)}$$

$$x^{(l+1)} = f(W^{(l+1)}\widetilde{x^{(l)}} + b^{(l+1)}) = f(W^{(l+1)}(D^{(l)}x^{(l)}) + b^{(l+1)})$$

행렬곱셈은 결합법칙이 성립. 괄호의 위치를 바꿀 수 있다

$$x^{(l+1)} = f((W^{(l+1)}D^{(l)})x^{(l)} + b^{(l+1)})$$

1. 옵티마이저와 드롭아웃 드롭아웃?

ONECLICK AI

Dropout?

$W^{(l+1)}D^{(l)}$ 는 $W^{(l+1)}$ 의 열(column)들을 $D^{(l)}$ 의 대각 원소(0 또는 1)로 스케일링하는 것과 같다.

$D^{(l)}$ 의 대각 원소가 0인 위치에 해당하는 $W^{(l+1)}$ 의 열 벡터 전체가 0으로 바뀐다.

이는 드롭 아웃마냥 비활성 되는거다

결론적으로, 훈련 과정의 매 스텝마다 신경망은 무작위로 가중치 행렬의 일부 열들이 0으로 설정 더 얇은 부분 신경망을 학습하게 된다

? 근데 렐루는 0 되는게 문제라 하지 않았음?? 죽은 렐루라고??

근데 죽은 렐루는 컨트롤이 안됨. 드롭아웃은 컨트롤 된다.

그리고 죽은 렐루는 한번 죽으면 영원히 업데이트 없지만,

드롭아웃은 다음 스텝의 학습에선 학습될 가능성 있음

1. 옵티마이저와 드롭아웃 드롭아웃?

ONECLICK AI

Dropout?

드롭아웃의 문제점

훈련 시 뉴런의 p 비율만 사용했기 때문에, 활성화 값의 총 기대값이 p 배만큼 작아진다.
이를 보정하기 위해 추론(모델을 실사용) 시에는 모든 뉴런을 사용하되,
가중치 행렬 전체에 p 를 곱해줘야 $W_{\text{test}} = p \cdot W_{\text{train}}$ 훈련 때와 스케일을 맞출 수 있다

이러한 불편함을 없앤 것이 바로 **역전파 드롭아웃**
원하는 총량을 유지하기 위해 활성화 비율로 나누어 준다

1. 옵티마이저와 드롭아웃 역전파 드롭아웃?

ONECLICK AI

Dropout? 예시를 통해 이해해 보자

입력신호 x : [1, 5, 2, 8]

가중치 w : [0.1, 0.2, 0.3, 0.4]

출력 z : $0.1 + 1 + 0.6 + 3.2 = 4.9$

여기서, 4.9가 목표 출력이다

1. 일반 드롭아웃인 경우

1. 마스크 생성 : [1, 0, 1, 0]

2. 입력에 마스크 적용 : [1, 0, 2, 0]

3. 출력 계산 : $(1 * 0.1) + (0 * 0.2) + (2 * 0.3) + (0 * 0.4) = 0.1 + 0 + 0.6 + 0 = 0.7$

최종적으로, 0.7이라는 너무 많이 약해진 출력의 신호를 기준으로 학습한다

1. 옵티마이저와 드롭아웃 역전파 드롭아웃?

ONECLICK AI

Dropout? 예시를 통해 이해해 보자

1. 일반 드롭아웃인 경우 예측시

훈련시 신호가 50%로 약해졌기 때문에 이를 보정해주기 위해 가중치에 p 인 0.5를 곱한다

1. 가중치 스케일링 : $w_{test} = w * p = [0.1, 0.2, 0.3, 0.4] * 0.5 = [0.05, 0.1, 0.15, 0.2]$

2. 최종출력 계산 : $z_{test} = (1 * 0.05) + (5 * 0.1) + (2 * 0.15) + (8 * 0.2)$

$$z_{test} = 0.05 + 0.5 + 0.3 + 1.6 = 2.45$$

이러면, 목표치의 정확히 절반의 값이 나온다 이렇게, 강제로 스케일을 맞춰 일관성을 유지한다
강제로 맞춰도, 절반으로 나온다

1. 옵티마이저와 드롭아웃 역전파 드롭아웃?

ONECLICK AI

Dropout? 예시를 통해 이해해 보자

2. 역전파 드롭아웃인 경우

같은 뉴런이 꺼졌다고 가정해보자

1. 마스크 생성 : $[1, 0, 1, 0]$

2. 입력에 마스크 적용 : $[1, 0, 2, 0]$

3. 살아남은 신호를 p로 나누어서 증폭 : $x_{inverted} = x_{masked} / p = [1, 0, 2, 0] / 0.5 = [2, 0, 4, 0]$

4. 출력 계산 : $z_{inverted_{train}} = (2 * 0.1) + (0 * 0.2) + (4 * 0.3) + (0 * 0.4) = 0.2 + 0 + 1.2 + 0 = 1.4$

1.4는 4.9와는 너무 다른 값이다. 하지만, 다른 마스크를 적용하면 또 다른 값이 나오는 등, 평균적인 기댓값이 4.9에 맞춰지도록 훈련이 진행된다

1. 옵티마이저와 드롭아웃 역전파 드롭아웃?

ONECLICK AI

Dropout? 예시를 통해 이해해 보자

2. 일반 드롭아웃인 경우 예측시

훈련 과정에서 이미 보정을 해 주었기 때문에, 예측시에는 그냥 하면 된다

$$z_{test} = (1 * 0.1) + (5 * 0.2) + (2 * 0.3) + (8 * 0.4) = 4.9$$

이렇게 해서, 목표한 값 4.9가 그래도 출력되게 한다

1. 옵티나이저와 드롭아웃 비교

ONECLICK AI

일반 드롭아웃과 역전파 드롭아웃 비교

구분	일반 드롭아웃	역전파 드롭아웃 (Inverted Dropout)
핵심 동작	훈련 시 신호를 약하게 하고, 예측 시 가중치를 보정	훈련 시 약해진 신호를 미리 증폭
훈련 출력 (예시)	0.7 (약해진 신호)	1.4 (보정된 신호)
예측 출력 (예시)	2.45 (보정된 결과)	4.9 (원래 결과)
장점	개념이 단순함	훈련/예측 모델이 동일하여 편리함 (현재 표준)

2. Yolo 아키텍처 알아보기 어떻게 생겼을까?

ONECLICK AI

Yolo와 함께 볼 것들?

실시간 객체 검출을 위한 모델

많은 버전이 있으며, 그 중에도 가장 최신인 Yolo v8을 기준으로 알아보자

이미지 내 객체의 위치와 종류를 한 번에 예측하는 단일 단계 검출기

분류와 지역화를 매우 빠른 속도로 한 번에 처리한다

2. Yolo 아키텍처 알아보기 어떻게 생겼을까?

ONECLICK AI

Yolo와 함께 볼 것들?

1. Conv는 CNN때 많이 해서 잘 알고 있으리라 믿는다

2. 배치 정규화

컨볼루션 연산을 거치며 숫자들의 분포가 한쪽으로 치우치거나 너무 커지는 것을 막아준다

각 레이어의 작업 결과를 "표준 규격"에 맞게 보정하여 다음 레이어가 작업을 더 쉽고 안정적으로 할 수 있게 돕는 역할이다

2. Yolo 아키텍처 알아보기 어떻게 생겼을까?

ONECLICK AI

Yolo와 함께 볼 것들?

3. 활성화 함수 SiLU

컨볼루션이 찾아낸 패턴이 얼마나 중요한지를 판단하는 스위치 역할
탐지된 패턴이 중요하고 의미 있다면 스위치를 활짝 열어(강한 신호 전달),
별 의미 없다면 스위치를 닫아(약한 신호 전달) 불필요한
정보를 걸러낸다

$$SiLU(x) = x \cdot \sigma(x)$$

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

시그모이드를 통해서 입력값 x 를 0 ~ 1 사이로 정규화 (비선형성 부여)
값의 크기를 통해서 얼마나 값을 통과시킬지 스케일 값 역할을 하게 된다
입력값이 매우 크면 시그모이드의 출력이 커 연산되었을 때 값이 크게 나오고(문이 많이 열림)
입력값이 작으면 시그모이드의 연산값이 작게 된다(문이 덜 열림)

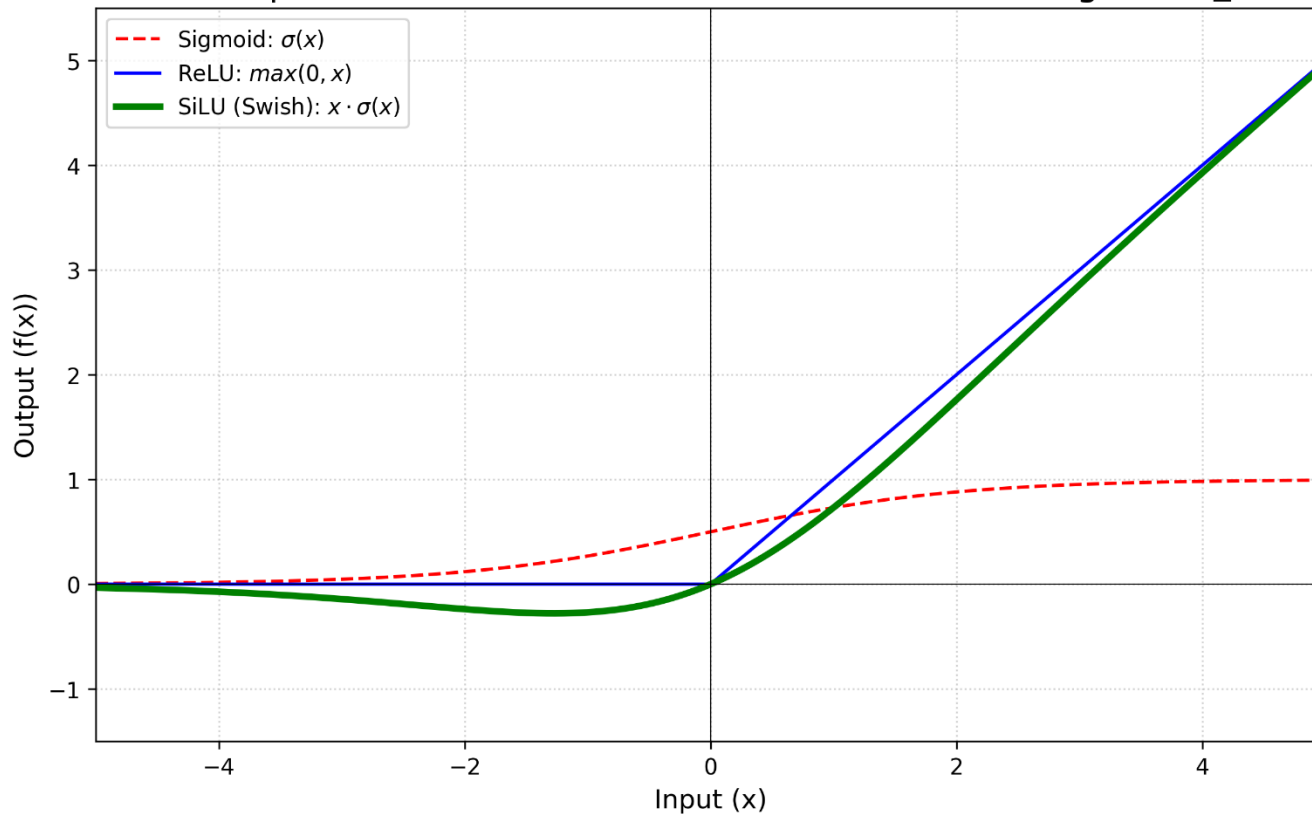
입력 벡터를 그대로 유지 (선형성)

2. Yolo 아키텍처 알아보기 어떻게 생겼을까?

ONECLICK AI

Yolo와 함께 볼 것들?

Comparison of Activation Functions (SiLU, ReLU, Sigmoid) □



? 렐루랑 생긴거 비슷한데 이거 왜 쓰는거임?

동적 연산:

단순히 0이냐 그냥 출력이냐를 비교하는 렐루와 다르게 동적 연산을 진행

죽은 노드 값 완화:

ReLU는 0이하 값이 들어오면 역전파 때 죽어버리는 문제가 있지만, SiLU는 0에 수렴하지 0이 아니기 때문에 그럴일이 없다

2. Yolo 아키텍처 알아보기 어떻게 생겼을까?

ONECLICK AI

Yolo?

다음과 같은 3개의 구조를 가진다

1. 백본 : 픽셀을 의미로 바꾸는 특징 추출기
간단하게 표현하면, CNN이다
2. 넥 : 흠어진 정보를 종합하는 융합기
여러 정보를 전부 종합한다
3. 헤드 : 최종 예측기
여러 정보를 전부 종합한다

2. Yolo 아키텍처 알아보기 어떻게 생겼을까?

ONECLICK AI

Yolo?

1. 백본(Backbone)

CNN에서 했던, 특징 추출 과정을 그대로 따라간다
여기서는, C2f 모듈을 수십번 반복해서 이미지의 특징을 점진적으로 추출해 낸다

C2f 모듈 구성

1. Conv

2. 배치 정규화

3. 활성화 함수 SiLU

2. Yolo 아키텍처 알아보기 어떻게 생겼을까?

ONECLICK AI

Yolo?

2. 넥(Neck)

백본을 통과하면서 이런 정보가 만들어진다

1. 깊은 층의 정보 : 해상도는 낮지만, 이건 사람이다 같은 의미 정보가 많이 있다
2. 얇은 층의 정보 : 해상도는 높지만, 여기에 모서리가 있다 같은 위치 정보만 있다

넥은 이 두 종류의 정보를 효과적으로 융합, 이 위치에는 사람이 있다 라는 종합적 정보를 만들어 낸다

여기서는 PANet을 사용한다

1. 하향식 경로 : 먼저 깊은 층의 의미 정보를 얇은 층으로 전달한다.
사람을 찾아봐 같은 지시를 내리는 것과 같다
2. 상향식 경로 : 그 후, 지시를 받은 얇은 층이 가진 정확한 위치 정보를 다시 깊은 층으로 전달.
사람은 이 픽셀 주변에 있다고 알려준다

이렇게 양방향 정보 교환을 통해서 넥은 크고 작은 모든 객체를 정확히 탐지한다

2. Yolo 아키텍처 알아보기 어떻게 생겼을까?

ONECLICK AI

Yolo?

3. 헤드(Head)

넥에서 완벽하게 만들어진 2가지의 특징을 맵을 박아서 우리가 원하는 형태의 출력을 한다(분류)

2가지의 특징을 가지는데,

1. 앵커 프리

과거는 미리 여러 크기의 예상 사각형(앵커)를 정해두고,

3번 예상 사각형을 약간 늘리면 객체에 맞겠군 이라고 예측하였다

Yolo v8은 예상 사각형 없이 객체의 중심은 여기고,

너비는 이만큼, 높이는 이만큼이라고 위치와 크기를 직접 예측한다. 훨씬 직관적이고 효율적이다

2. Yolo 아키텍처 알아보기 어떻게 생겼을까?

ONECLICK AI

Yolo?

3. 헤드(Head)

2. 분리형 헤드

헤드 안에는 두 개의 역할이 있다

1. 위치 역할 : 오직 객체의 위치(사각형 좌표)를 예측하는 것에만 특화되어 있다

2. 종류 역할 : 오직 객체의 종류를 예측하는 데만 특화되어 있다

이렇게, 역할을 분리하면 더 잘 예측해서 위치와 종류 모두에 대한 예측 정확도가 크게 향상된다

2. Yolo 아키텍처 알아보기 어떻게 생겼을까?

ONECLICK AI

다음과 같은 아키텍처라 하고 해 보자

입력 데이터

1. 백본

Conv1` -> `SiLU1` -> `Conv2` -> `SiLU2`

정규화는 이번엔 생략. 정규화 할 정도로 복잡하지 않은 구성

넥은 생략하겠다. 하행, 상행이 미준 적분 계속 반복하는거라 인간이 손으로 할 이유가 없다

2. 헤드

3. 수식으로 Yolo 열어보기 어떻게 생겼을까?

ONECLICK AI

다음과 같은 파라미터

입력 이미지 $I = \begin{pmatrix} 1 & 2 & 3 & 0 \\ 4 & 5 & 6 & 1 \\ 7 & 8 & 9 & 2 \\ 1 & 0 & 1 & 3 \end{pmatrix}$

가중치

백본 Conv 1 (3x3 필터) $W_{\text{conv1}} = \begin{pmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \end{pmatrix}$

백본 Conv 2 (2x2 필터) $W_{\text{conv2}} = \begin{pmatrix} 0.5 & 0.6 \\ 0.7 & 0.8 \end{pmatrix}$

넥 (1x1, 스칼라 값) $W_{\text{neck}} = 0.5$

헤드 (분리형)

Bounding Box 예측용 $W_{\text{head_bbox}} = \begin{pmatrix} 0.1 \\ 0.2 \end{pmatrix}$

Class 예측용 $W_{\text{head_cls}} = 0.3$

정답 레이블 : Bbox 좌표, 클래스

$$Y_{\text{true}} = \begin{pmatrix} 12 \\ 1520 \end{pmatrix}$$

학습율 : 0.01

3. 수식으로 Yolo 열어보기 1. 순전파

ONECLICK AI

순전파 1. c2f Conv 1 계산

$$I = \begin{pmatrix} 1 & 2 & 3 & 0 \\ 4 & 5 & 6 & 1 \\ 7 & 8 & 9 & 2 \\ 1 & 0 & 1 & 3 \end{pmatrix} \quad W_{\text{conv1}} = \begin{pmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 1 \end{pmatrix} \quad \text{최종 출력 } \text{Out}_{\text{conv1}} = \begin{pmatrix} 25 & 18 \\ 20 & 18 \end{pmatrix}$$

출력 0, 0

$$(1 \cdot 1 + 2 \cdot 0 + 3 \cdot 1) + (4 \cdot 0 + 5 \cdot 1 + 6 \cdot 0) + (7 \cdot 1 + 8 \cdot 0 + 9 \cdot 1) = 1 + 3 + 5 + 7 + 9 = 25$$

출력 0, 1

$$(2 \cdot 1 + 3 \cdot 0 + 0 \cdot 1) + (5 \cdot 0 + 6 \cdot 1 + 1 \cdot 0) + (8 \cdot 1 + 9 \cdot 0 + 2 \cdot 1) = 2 + 6 + 8 + 2 = 18$$

출력 1, 0

$$(4 \cdot 1 + 5 \cdot 0 + 6 \cdot 1) + (7 \cdot 0 + 8 \cdot 1 + 9 \cdot 0) + (1 \cdot 1 + 0 \cdot 0 + 1 \cdot 1) = 4 + 6 + 8 + 1 + 1 = 20$$

출력 1, 1

$$(5 \cdot 1 + 6 \cdot 0 + 1 \cdot 1) + (8 \cdot 0 + 9 \cdot 1 + 2 \cdot 0) + (0 \cdot 1 + 1 \cdot 0 + 3 \cdot 1) = 5 + 1 + 9 + 3 = 18$$

3. 수식으로 Yolo 열어보기 1. 순전파

ONECLICK AI

순전파 1. Conv 1에 SiLU 적용

이전 레이어 최종 출력 $\text{Out}_{\text{conv1}} = \begin{pmatrix} 25 & 18 \\ 20 & 18 \end{pmatrix}$

$$\text{SiLU}(x) = x \cdot \sigma(x) \quad \sigma(x) = \frac{1}{1 + e^{-x}}$$

$$\text{Out}_{\text{SiLU1}} = \begin{pmatrix} \text{SiLU}(25) & \text{SiLU}(18) \\ \text{SiLU}(20) & \text{SiLU}(18) \end{pmatrix} \approx \begin{pmatrix} 25 & 18 \\ 20 & 18 \end{pmatrix}$$

3. 수식으로 Yolo 열어보기 1. 순전파

ONECLICK AI

순전파 2. c2f Conv 2 계산

이전 레이어 최종 출력

Out_{SiLU1} 에 W_{conv2} 적용 후 값을 하나로 압축

$$Out_{conv2} = (25 \cdot 0.5) + (18 \cdot 0.6) + (20 \cdot 0.7) + (18 \cdot 0.8) = 12.5 + 10.8 + 14.0 + 14.4 = 51.7$$

왜 여기서 압축을 하는 걸까?? 이거 Flatten Layer 아닌가요?

정규화 대신 FlattenLayer 사용

3. 수식으로 Yolo 열어보기 1. 순전파

ONECLICK AI

순전파 2. c2f Conv 2 에 SiLU 계산

이전 레이어 최종 출력

$$\text{Out}_{\text{conv2}} = (25 \cdot 0.5) + (18 \cdot 0.6) + (20 \cdot 0.7) + (18 \cdot 0.8) = 12.5 + 10.8 + 14.0 + 14.4 = 51.7$$

여기에 활성화 함수 사용

$$\text{Out}_{\text{backbone}} = \text{SiLU}(51.7) \approx 51.7$$

3. 수식으로 Yolo 열어보기 1. 순전파

ONECLICK AI

순전파 3. 넥 스킵 헤더 계산

Bbox 예측

$$Y_{\text{bbox}} = W_{\text{bbox}} \times \text{Out}_{\text{backbone}} = (0.10.2) \times 51.7 = (5.1710.34)$$

클래스 예측

$$Y_{\text{cls}} = W_{\text{cls}} \times \text{Out}_{\text{backbone}} = 0.3 \times 51.7 =$$

3. 수식으로 Yolo 열어보기 1. 순전파

ONECLICK AI

순전파 4. 예측값, 손실값 계산

최종 예측. 두 클래스 값 결합

$$Y_{\text{pred}} = \begin{pmatrix} 5.17 \\ 10.3415.51 \end{pmatrix}$$

손실값 계산

예측과 정답의 평균 제곱 오차(MSE)를 계산합니다.

$$L = \frac{1}{3}((12 - 5.17)^2 + (15 - 10.34)^2 + (20 - 15.51)^2) \approx 29.51$$

3. 수식으로 Yolo 열어보기 2. 역전파

ONECLICK AI

역전파 1. 최종 출력 오차 계산

이전 레이어 최종 출력

순전파 과정 복기 :

최종 예측값 Y_{pred} 와 정답 Y_{true} 로 손실 함수 $L = \frac{1}{N} \sum \frac{1}{2} (Y_{true} - Y_{pred})^2$ 를 계산하였다.
(계산 편의상 계수는 $\frac{2}{3}$ 로 통일)

미분을 통한 공식 유도 :

손실 L 을 Y_{pred} 로 미분하면 $\frac{\partial L}{\partial Y_{pred}} = \frac{2}{N} (Y_{pred} - Y_{true})$ 가 된다.

이 값이 모든 역전파의 시작점이 되는 '오차의 기울기'입니다.

$$\frac{\partial L}{\partial Y_{pred}} = \frac{2}{3} \begin{pmatrix} 5.17 - 12 \\ 10.34 - 1515.51 - 20 \end{pmatrix} = \frac{2}{3} \begin{pmatrix} -6.83 \\ -4.66 - 4.49 \end{pmatrix} \approx \begin{pmatrix} -4.55 \\ -3.11 - 2.99 \end{pmatrix}$$

3. 수식으로 Yolo 열어보기 2. 역전파

ONECLICK AI

역전파 2. 헤더 기울기 계산

순전파 과정 복기 :

헤드는 입력 $X = \text{Out}_{\text{backbone}}$ 에 가중치 W_{head} 를 곱해 $Y_{\text{pred}} = W_{\text{head}} \cdot X$ 를 계산했습니다.

미분을 통한 공식 유도 :

위 식을 가중치 W_{head} 에 대해 미분하면 $\frac{\partial Y_{\text{pred}}}{\partial W_{\text{head}}} = X$ 입니다.

따라서, 연쇄 법칙에 의해 최종 손실에 대한 가중치의 기울기는 $\frac{\partial L}{\partial W_{\text{head}}} = \frac{\partial L}{\partial Y_{\text{pred}}} \cdot X$ 가 됩니다.

기울기 계산 :

$$\text{Bbox 계산 : } \frac{\partial L}{\partial W_{\text{bbox}}} = (-4.55 - 3.11) \times 51.7 \approx (-235.24 - 160.79)$$

$$\text{Class 계산 : } \frac{\partial L}{\partial W_{\text{cls}}} = -2.99 \times 51.7 \approx -154.58$$

3. 수식으로 Yolo 열어보기 2. 역전파

ONECLICK AI

역전파 2. backbone으로 가중치 전달

순전파 과정 복기 : $Y_{pred} = W_{head} \cdot X$ 에서 x 는 $Out_{backbone}$ 이었다

미분을 통한 공식 유도 : 이번에는 Y_{pred} 를 입력 x 에 대해 미분한다

$$\frac{\partial Y_{pred}}{\partial X} = W_{head} \quad \text{이므로, } x \text{에 대해 미분하면, } \frac{\partial L}{\partial X} = \frac{\partial L}{\partial Y_{pred}} \cdot W_{head}$$

분리형 헤드이므로, 각 헤드로부터 전파된 오차 신호를 합산해야한다

$$\begin{aligned} \text{기울기 계산 : } \frac{\partial L}{\partial Out_{backbone}} &= \left(W_{bbox}^T \cdot \frac{\partial L}{\partial Y_{bbox}} \right) + \left(W_{cls} \cdot \frac{\partial L}{\partial Y_{cls}} \right) \\ &= ((0.1 \quad 0.2)(-4.55 - 3.11)) + (0.3 \times -2.99) \\ &= (-0.455 - 0.622) - 0.897 = -1.974 \end{aligned}$$

3. 수식으로 Yolo 열어보기 2. 역전파

ONECLICK AI

역전파 3. c2f로 가중치 전달

SiLU 2 → Conv 2 오차 전파

순전파 과정 복기 : $\text{Out}_{\text{backbone}} = \text{SiLU}(\text{Out}_{\text{conv2}})$ 였으므로, 오차신호 -1.974 에 $\text{SiLU}'(51.7) \approx 1$ 를 곱해서 구한다 $\frac{\partial L}{\partial \text{Out}_{\text{conv2}}} \approx -1.974$

Conv2 기울기 계산 : $\text{Out}_{\text{conv2}} = \sum (W_{\text{conv2}} \cdot \text{Out}_{\text{SiLU1}})$ 였으므로, 가중치에 대한 기울기는 입력값인 $\text{Out}_{\text{SiLU1}}$ 가 된다

$$\frac{\partial L}{\partial W_{\text{conv2}}} = \left(\frac{\partial L}{\partial \text{Out}_{\text{conv2}}} \right) \times (\text{Out}_{\text{SiLU1}}) = -1.974 \times \begin{pmatrix} 25 & 18 \\ 20 & 18 \end{pmatrix} = \begin{pmatrix} -49.35 & -35.53 \\ -39.48 & -35.53 \end{pmatrix}$$

3. 수식으로 Yolo 열어보기 2. 역전파

ONECLICK AI

역전파 4. Conv2 -> Conv1 으로 가중치 전달

Conv2 -> SiLU 1 오차 전파

Conv2의 입력이었던 Out_{SiLU1} 에 대한 오차는 $\frac{\partial L}{\partial \text{Out}_{\text{conv2}}} \cdot W_{\text{conv2}}$ 로 계산된다.

$$\frac{\partial L}{\partial \text{Out}_{\text{SiLU1}}} = -1.974 \times \begin{pmatrix} 0.5 & 0.6 \\ 0.7 & 0.8 \end{pmatrix} = \begin{pmatrix} -0.987 & -1.184 \\ -1.382 & -1.579 \end{pmatrix}$$

SiLU 1 → Conv 1 오차 전파

SiLU의 미분값을 곱하여 오차를 전달

$$\frac{\partial L}{\partial \text{Out}_{\text{conv1}}} \approx \begin{pmatrix} -0.987 & -1.184 \\ -1.382 & -1.579 \end{pmatrix}$$

3. 수식으로 Yolo 열어보기 2. 역전파

ONECLICK AI

역전파 4. Conv2 -> Conv1 으로 가중치 전달

Conv 1 기울기 계산

Conv1의 기울기는 입력 이미지 I와 out_conv1에 대한 오차 행렬의 합성곱으로 계산

$$\frac{\partial L}{\partial W_{\text{conv1}}} \approx \begin{pmatrix} -12.44 & -16.03 & -8.99 \\ -21.49 & -27.56 & -14.93 \\ -10.03 & -12.42 & -7.56 \end{pmatrix}$$

3. 수식으로 Yolo 열어보기 3. 가중치 업데이트

ONECLICK AI

가중치 업데이트. 공식은 이미 알고 있으리라 믿는다 경사하강법

Bbox 헤드 $W'_{bbox} = (0.10.2) - 0.01 \times (-235.24 - 160.79) = (2.451.81)$

Class 헤드 $W'_{cls} = 0.3 - (0.01 \times -154.58) = 1.85$

Conv2 $W'_{conv2} = \begin{pmatrix} 0.5 & 0.6 \\ 0.7 & 0.8 \end{pmatrix} - 0.01 \times \begin{pmatrix} -49.35 & -35.53 \\ -39.48 & -35.53 \end{pmatrix} \approx \begin{pmatrix} 0.99 & 0.96 \\ 1.09 & 1.16 \end{pmatrix}$

Conv1 $W'_{conv1} = W_{conv1} - 0.01 \times \frac{\partial L}{\partial W_{conv1}} \approx \begin{pmatrix} 1.12 & 0.16 & 1.09 \\ 0.21 & 1.28 & 0.15 \\ 1.10 & 0.12 & 1.08 \end{pmatrix}$

감사합니다