## Neural Network Console クラウド版 ネットワーク解説 -物体検出編-

ソニーネットワークコミュニケーションズ株式会社

#### 概要

本ドキュメントではNeural Network Console(NNC)にある物体検出のサンプルプロジェクト (tutorial.object detection.synthetic image object detection)のネットワークを解説します。 サンプルプロジェクトをベースに自らネットワークを変更して精度改善を図るために、現状のサンプル プロジェクトのネットワーク構造を理解したい方を読者と想定しております。 これから物体検出を始めようという方は、物体検出の一連の流れを解説した<u>スターターガイド-物体解</u> <u>説編-</u>をまずはご確認ください。

本ドキュメントでは各レイヤー<sup>※1</sup>がネットワーク全体でどのような役割を担っているかに焦点を当て説 明していますので、ネットワークの中で用いられているレイヤーの具体的な機能は<u>レイヤーリファレン</u> <u>ス</u>をご確認ください。

※1 レイヤーとはDeep Learningでネットワークを作成するための関数で、NNCに限らず一般的なもののため、それぞれの 詳細な仕組みなどは入門書などでも確認することができます。 物体検出のサンプルプロジェクトにはネットワークとデータセットが含まれています。

ネットワークの全体像を理解するうえで、データセットとネットワークの入出力の理解が必須であるため、 まずはこれらを解説します。そのあとにネットワークの詳細な構造を解説していきます。

データセットとネットワークの入出力 1

ネットワークの構造 2

### -般的な物体検出の入出力

- 一般的な物体検出では画像に映る検出対象のラベルと、それを囲う四角形(バウンディングボックス)を検出 します。
- 本サンプルプロジェクトには、正方形の背景の上に色、個数、形状がランダムな物体(楕円、三角形、四角形、 五角形)を配置したデータセットがすでに登録されています。



### データセットの説明

データセットには入力として画像データ(image)が、出力としてアンカーボックスとグリッドごとのラベル (label)とバウンディングボックス(region)が含まれています。

labelとregionはバウンディングボックスの中心が含まれるグリッドに値が記載されています。中心が含まれない箇所はlabelは-1、regionは0で穴埋めがされています。

アンカーボックスやグリッドは学習効率化のための手法で、詳細は<u>スターガーガイド</u>で解説していますので、 ご参照ください。

データセットの入出力

ラベルと物体の対応



※相対位置はグリッドサイズで規格化をし、バウンディングボックスのサイズはグリッドサイズで規格化をした後に対数変換をしています

### ネットワークの出力情報

実際のネットワークではデータセットの出力のlabelをグリッド内に物体中心を含む確率(score)とグリッド内 に含まれる物体中心の分類ラベルの指数<sup>※1</sup>(category)の2つに分けて予測をします。 検出物体の分類ラベル予測にあわせて、検出物体の有無(score)を予測することで効率的な学習ができます。 したがってネットワークの出力はscore、category、regionの3つになります。

#### ネットワークの出力情報



アンカーボックス数(5)

出力情報の説明

<b>^</b>	
出力情報	説明
score	ラベルによらずグリッド内に物 体中心を含む確率
category	グリッド内に含まれる物体中心 の分類ラベルの指数 <sup>※1</sup>
region	グリッド内に物体中心があるバ ウンディングボックスのグリッ ド内での相対位置(横、縦)とサイ ズ(幅、高さ) <sup>※2</sup>

#### ※1 softmax変換により確率値に変換できます

※2 グリッドサイズにより規格化されており、バウンディングボックスのサイズはさらに対数変換した値です

ネットワークの出力ファイル

scoreは確率値を表すモノクロ画像がアンカーボックス数分出力されます。

categoryとregionは1つのCSVファイルとして出力するために、x方向に分類クラスの確率やバウンディングボックスの情報を、y方向にアンカーボックス数分の情報をそれぞれ並べ、2次元の情報に変換します。



ここまでネットワークの全体像を理解いただくために、データセットとネットワークの入出力を解説をして きました。ここからはネットワークの構造を解説していきます。

## 1 データセットとネットワークの入出力

### 2 ネットワークの構造

### ネットワークの全体構成

本サンプルプロジェクトでは学習用、検証用、推論用のネットワークの入出力や誤差評価指数などを個別に 設定するため、Trainingタブ、Validationタブ、Runtimeタブと別々に記載しています<sup>※</sup>。 各ネットワークでは共通のNetworkタブを引用しており、NetworkタブはConvUnitタブを引用しています。

ネットワーク画面の説明

ネットワークタブの引用関係



※それぞれのネットワークはCONFIGタブの設定画面で学習用、検証用、推論用に対応づけられています

### NetworkタブとConvUnitタブの概要

Networkタブでは画像処理で一般的なCNN(<u>C</u>onvolutional <u>N</u>eural <u>N</u>etwork)型のネットワークが記述されています。 Unitレイヤーで引用しているConvUnitタブがCNNの基本構造を示しています。

Networkタブ



SONY

ConvUnitタブ

### 【参考】精度改善のためのNetworkタブの変更方法

NetworkタブのUnitレイヤーの増減やUnitレイヤーのパラメータ変更で、モデルの精度改善が期待できます。 変更の際には最終層のサイズをscore、category、regionのサイズに合わせることに注意してください。 各層のサイズはUnitレイヤーの詳細設定で決まります。

## Networkタブの設計ポイント



#### Unitレイヤーの詳細設定の例

ConvMapsで1次元目のサイズが決定し、ConvStrideで2,3次元目のサイズ が決定します。ConvStrideが1の場合には入力と同サイズ、2の場合には 入力の1/2のサイズになります。 詳細設定では以下のようにパラメータを用いた入力も可能です。

Unit		Unit_2		Unit_3	
Name	Unit	Name	Unit_2	Name	Unit_3
Input	1,112,112	Input	32,56,56	Input	32,56,56
Network	ConvUnit	Network	ConvUnit	Network	ConvUnit
ParameterScope	Unit	ParameterScope	Unit_2	ParameterScope	Unit_3
ConvRS	False	ConvRS	False	ConvBS	False
ConvMaps	32	ConvMaps	Input[0]	ConvMaps	Input[0]*2
ConvStride	2	ConvStride	1	ConvStride	2
Convolution.W.File		Convolution.W.File		Convolution.W.File	
Convolution.W.Initial 1		Convolution.W.Initial 1		Convolution.W.Initial 1	

※1 サンプルプロジェクトの場合、1次元目は(score1枚, category4枚, region4枚)がアン カーボックス数5セットで45、2,3次元目は格子状のグリッドで7,7にサイズ調整

### Trainingタブ、Validationタブの概要<sup>※1</sup>

CNN(Networkタブの引用)からの出力をscore、category、regionに分割し、データセットから作成したそれぞれの値との誤差を計算をしています。誤差は合計され、学習時や学習曲線の描画で利用されます。



※1 TrainingタブとValidationタブのネットワーク構造は同様のため、合わせて説明をしています。

### Trainingタブ、Validationタブのscore部分の解説

入力画像からCNN(Networkタブの引用)を通して得られたscoreと、データセットのlabelから作り出したscoreとの誤差を計算しています。scoreは物体を含まないことが多いため、学習のテクニックとして最後に物体を含む場合と含まない場合で誤差に重みづけをしています。



### Trainingタブ、Validationタブのcategory部分の解説

入力画像からCNN(Networkタブの引用)を通して得られたcategoryと、データセットのlabelとの誤差を計算しています。最初に全てのグリッドの誤差を求め、最後に物体中心を含まない箇所を対象外にしています。



※1 Reshapeはレイヤーの形状を調整しています ※2 グレーの部分はscore, regionに関連するため、他頁で説明をしています ※3 誤差評価時に-1ではエラーとなることの回避策です。最後に中心を含まない箇所を誤差評価の対象外にしているため、学習への影響はありません

### Trainingタブ、Validationタブのregion部分の解説

入力画像からCNN(Networkタブの引用)を通して得られたregionと、データセットのregionとの誤差を計算しています。最初に全てのグリッドの誤差を求め、最後に物体中心を含まない部分を対象外にしています。



### Runtimeタブの概要

CNN(Networkタブの引用)からの出力をscore、category、regionに分割し、それぞれ出力しています。 同一物体の複数のアンカーボックスでの検出を避けるため、1グリッドの出力を1つのアンカーボックスに 制限しています。また、確率値の閾値判定を行い、確度が高いもののみを出力しています。



### Runtimeタブのscore部分の解説

入力画像からCNN(Networkタブの引用)を通して得られたscoreを出力しています。 各グリッドで確率値の最も高いアンカーボックスの確率値だけを出力しています。



### Runtimeタブのcategory部分の解説

入力画像からCNN(Networkタブの引用)を通して得られたcategoryを出力しています。 各グリッドでアンカーボックス間のscoreの最大値が閾値を超えた箇所のcategoryだけを出力しています。



### Runtimeタブのregion部分の解説

入力画像からCNN(Networkタブの引用)を通して得られたregionを出力しています。 各グリッドでアンカーボックス間のscoreの最大値が閾値を超えた箇所のcategoryだけを出力しています。



# SONY

SONYはソニー株式会社の登録商標または商標です。

各ソニー製品の商品名・サービス名はソニー株式会社またはグループ各社の登録商標または商標です。その他の製品および会社名は、各社の商号、登録商標または商標です。