

ここで、 i はチャンネルのインデックスを表し、 u は空間位置のインデックス（簡素化のため高さおよび幅を意味する）です。異常局所化は異常領域を特定し、各ピクセル u に対して異常スコアを割り当てる異常スコアマップ $s(u)$ を生成します。 $s(u)$ は特微量差分ベクトル $d(:, u)$ の $L2$ ノルムとして計算されます。

$$s(u) = \|d(:, u)\|_2. \quad (3)$$

異常検出は、画像に異常な領域が含まれているかどうかを検出することを目的とします。直感的に、平均プールされた $s(u)$ の最大値を画像全体の異常スコアとして採用します。

3.2 トランスフォーマーによる「同一マッピング」の防止

CNNと比較して、注意層のクエリ埋め込みがTransformerが「同一マッピング」を学習しにくくすると推測されます。正常な画像の特性を $x^+ \in \mathbb{R}^{K \times C}$ と表し、ここで K は特性数、 C はチャンネル次元です。異常画像の特性を $x^- \in \mathbb{R}^{K \times C}$ と表します。再構築ネットワークとして1層のネットワークを採用し、 x^+ に対してMSE損失で訓練し、 x^- の異常領域を検出するためにテストします。CNNの畳み込み層。まず、全接続層を訪れます。その重みとバイアスはそれぞれ $w \in \mathbb{R}^{C \times C}$ 、 $b \in \mathbb{R}^C$ で表されます。この層を正常サンプルの再構築モデルとして使用する場合、次のように表せます。

$$\hat{x} = x^+ w + b \in \mathbb{R}^{K \times C}. \quad (4)$$

MSE損失関数が \hat{x} を x^+ に押し上げるため、モデルは $w \rightarrow I$ （単位行列）、 $b \rightarrow 0$ への近道を取る可能性があります。最終的に、このモデルは x^- を適切に再構築できるため、異常検出に失敗する可能性があります。1×1カーネルを持つ畳み込み層は、全接続層と等価です。さらに、 $n \times n$ ($n > 1$) のカーネルはより多くのパラメーターと大きな容量を持ち、1×1カーネルが可能なことはすべて実行できます。したがって、畳み込み層もショートカットを学習する可能性があります。

クエリ埋め込みを含むトランスフォーマーは、学習可能なクエリ埋め込み $q \in \mathbb{R}^{K \times C}$ を持つ注意層を含みます。この注意層は次のように表せます：

$$\hat{x} = \text{softmax}(q(x^+)^T / \sqrt{C}) x^+ \in \mathbb{R}^{K \times C}. \quad (5)$$

\hat{x} を x^+ に押し上げるため、アテンションマップ ($\text{softmax}(q(x^+)^T / \sqrt{C})$ は I （単位行列）に近似する必要があり、 q は x^+ と密接に関連している必要があります。訓練されたモデルで q が正常なサンプルに関連しているため、モデルは x^- を適切に再構築できません。第4.4節の消去実験では、注意層またはクエリ埋め込みを削除すると、トランスフォーマーの性能がそれぞれ2.4%または3%低下し、これはCNNとほぼ同じです。これは、注意層のクエリ埋め込みがトランスフォーマーが「同一のショートカット」を学習するのを防ぐのに役立つことを示しています。

3.3 異常データ利用可能ケースへの適応

実践では、生産ラインの稼働に伴い異常が徐々に増加するため、これらの増加する異常に対応する互換性が求められます。そのため、ADTRをADTR+に拡張し、異常データが利用可能なケースに対応できるように適応させます。



where i represents the index of channel, u is the index of spatial position (height together with width for simplicity). *Anomaly localization* aims to localize anomalous regions, producing an anomaly score map, $\mathbf{s}(u)$, which assigns an anomaly score for each pixel, u . $\mathbf{s}(u)$ is calculated as the $L2$ norm of the feature difference vector, $\mathbf{d}(:, u)$.

$$\mathbf{s}(u) = \|\mathbf{d}(:, u)\|_2. \quad (3)$$

Anomaly detection aims to detect whether an image contains anomalous regions. We intuitively take the maximum value of the averagely pooled $\mathbf{s}(u)$ as the anomaly score of the whole image.

3.2 Preventing “Identical Mapping” with Transformer

We suspect that, compared with CNN, the query embedding in attention layer makes transformer difficult to learn an “identical mapping”. We denote the features in a normal image as $\mathbf{x}^+ \in \mathbb{R}^{K \times C}$, where K is the feature number, C is the channel dimension. The features in an anomalous image are denoted as $\mathbf{x}^- \in \mathbb{R}^{K \times C}$. We take a 1-layer network as the reconstruction net, which is trained on \mathbf{x}^+ with the MSE loss and tested to detect anomalous regions in \mathbf{x}^- .

Convolutional layer in CNN. We first visit a fully-connected layer, whose weights and bias are denoted as $\mathbf{w} \in \mathbb{R}^{C \times C}$, $\mathbf{b} \in \mathbb{R}^C$, respectively. When using this layer as the reconstruction model of normal samples, it can be written as,

$$\hat{\mathbf{x}} = \mathbf{x}^+ \mathbf{w} + \mathbf{b} \in \mathbb{R}^{K \times C}. \quad (4)$$

With the MSE loss pushing $\hat{\mathbf{x}}$ to \mathbf{x}^+ , the model may take shortcut to regress $\mathbf{w} \rightarrow \mathbf{I}$ (identity matrix), $\mathbf{b} \rightarrow \mathbf{0}$. Ultimately, this model could also reconstruct \mathbf{x}^- well, failing in anomaly detection. A convolutional layer with 1×1 kernel is equivalent to a fully-connected layer. Besides, An $n \times n$ ($n > 1$) kernel has more parameters and larger capacity, and can complete whatever 1×1 kernel can. Thus, the convolutional layer also has the chance to learn a shortcut.

Transformer with query embedding contains an attention layer with a learnable query embedding, $\mathbf{q} \in \mathbb{R}^{K \times C}$. This attention layer can be denoted as,

$$\hat{\mathbf{x}} = \text{softmax}(\mathbf{q}(\mathbf{x}^+)^T / \sqrt{C}) \mathbf{x}^+ \in \mathbb{R}^{K \times C}. \quad (5)$$

To push $\hat{\mathbf{x}}$ to \mathbf{x}^+ , the attention map, $\text{softmax}(\mathbf{q}(\mathbf{x}^+)^T / \sqrt{C})$, should approximate \mathbf{I} (identity matrix), so \mathbf{q} must be highly related to \mathbf{x}^+ . Considering that \mathbf{q} in the trained model is relevant to normal samples, the model could not reconstruct \mathbf{x}^- well. The ablation study in Sec. 4.4 shows that without the attention layer or the query embedding, the performance of transformer respectively drops by 2.4% or 3%, which is almost the same as CNN. This reflects that the query embedding in attention layer helps prevent transformer from learning an “identical shortcut”.

3.3 Adaptation with Anomaly-available Case

In practice, anomalies gradually increase with the runs of production lines, which brings the demands of compatibility with these increasing anomalies. Thus we adapt ADTR to ADTR+ for compatibility with the anomaly-available case.

