



黑客派

# 为什么分类问题用 cross entropy, 而回归问题用 MSE

作者: [moloee](#)

原文链接: <https://hacpai.com/article/1516696419994>

来源网站: [黑客派](#)

许可协议: [署名-相同方式共享 4.0 国际 \(CC BY-SA 4.0\)](#)

```

<h3 id="先说为什么分类问题用交叉熵">先说为什么分类问题用交叉熵</h3>
<script async src="https://pagead2.googlesyndication.com/pagead/js/adsbygoogle.js"></scr
pt>
<!-- 黑客派PC帖子内嵌-展示 -->
<ins class="adsbygoogle" style="display:block" data-ad-client="ca-pub-5357405790190342"
data-ad-slot="8316640078" data-ad-format="auto" data-full-width-responsive="true"></in
>
<script>
  (adsbygoogle = window.adsbygoogle || []).push({});
</script>
<h4 id="为啥不用-classification-error--">为啥不用 classification error ?</h4>
<p>来两个模型</p>
<p><strong>模型 1</strong></p>
<table>
<thead>
<tr>
<th>COMPUTED</th>
<th>TARGETS</th>
<th>CORRECT?</th>
</tr>
</thead>
<tbody>
<tr>
<td>0.3 0.3 0.4</td>
<td>0 0 1 (democrat)</td>
<td>yes</td>
</tr>
<tr>
<td>0.3 0.4 0.3</td>
<td>0 1 0 (republican)</td>
<td>yes</td>
</tr>
<tr>
<td>0.1 0.2 0.7</td>
<td>1 0 0 (other)</td>
<td>no</td>
</tr>
</tbody>
</table>
<p><strong>模型 2</strong></p>
<table>
<thead>
<tr>
<th>COMPUTED</th>
<th>TARGETS</th>
<th>CORRECT?</th>
</tr>
</thead>
<tbody>
<tr>
<td>0.1 0.2 0.7</td>
<td>0 0 1 (democrat)</td>
<td>yes</td>
</tr>

```

0.1	0.7	0.2
0	1	0 (republican)
yes		

  

0.3	0.4	0.3
1	0	0 (other)
no		

- 2 个模型的 classification error 相等，但模型 2 要明显优于模型 1。classification error 很难精确描述模型与理想模型之间的距离。
- 计算交叉熵，用交叉熵则可以更准确的衡量模型之间的差距

#### 为啥不用 MSE ?

总的来说，分类问题需要用 one hot 的形式计算个 label 的概率，然后用 argmax 来决定分类。计算概率的时候通常用 softmax。参考流程：计算 loss > 计算 softmax > argmax。

用 MSE 计算 loss 的问题在于，通过 Softmax 输出的曲线是波动的，有很多局部的极值点。即非凸优化问题 (non-convex)，如下图



而 cross entropy 计算 loss，则依旧是一个凸优化问题，用梯度下降求解时，凸优化问题有很好的收敛特性。

[为啥不用 cross entropy 计算 loss 的问题](https://link.hacpai.com/forward?goto=http%3A%2F%2Fjackon.me%2Fposts%2Fwhy-use-cross-entropy-error-for-loss-function%2F)

### 为什么回归问题用 MSE

#### 为啥不能用交叉熵

拿二项式的交叉熵定义来看



对于神经网络的分类问题等可以很好的使用，但是对于回归问题，任意取一个值比如 -1.5，就没法计算 log(-1.5)，所以一般不用交叉熵来优化回归问题。

#### 为什么用 MSE

最小二乘是在欧氏距离为误差度量的情况下，由系数矩阵所张成的向量空间内对于观测向量的最佳逼近点。

为什么用欧式距离作为误差度量（即 MSE），09 年 IEEE Signal Processing Magazine 的《Mean squared error: Love it or leave it?》这篇文章做了很好的讨论。链接：[文章](https://link.hacpai.com/forward?goto=https%3A%2F%2Flink.zhihu.com%2F%3Ftarget%3Dhttp%253A%2F%2Fwww2.units.it%2Fframponi%2Fteaching%2FDIP%2Fmateriale%2Fmse_bovik09.pdf)

这篇文章在“WHY DO WE LOVE THE MSE?”中说，MSE:

- 1. 它简单。
- 2. 它提供了具有很好性质的相似度的度量。例如：
  - 它是非负的;
  - 唯一确定性。只有  $x=y$  的时候， $d(x,y)=0$ ;
  - 它是对称的，即  $d(x,y)=d(y,x)$ ;
  - 符合三角性质。即  $d(x,z) \leq d(x,y)+d(y,z)$ 。
- 3. 物理性质明确，在不同的表示域变换后特性不变，例如帕萨瓦尔等式。

<li> <p>4. 便于计算。通常所推导得到的问题是凸问题，具有对称性，可导性。通常具有解析解，外便于通过迭代的方式求解。</p> </li>

<li> <p>5. 和统计和估计理论具有关联。在某些假设下，统计意义上是最优的。</p> </li></ul>

<p>然而，<strong>MSE 并非没有缺点</strong>。并不是所有的问题都可以套用该准则，在“IMPLICIT ASSUMPTIONS WHEN USING THE MSE”说，它基于了以下几点对于信号的假设：</p><script async src="https://pagead2.googlesyndication.com/pagead/js/adsbygoogle.js"></script>

<!-- 黑客派PC帖子内嵌-展示 -->

<ins class="adsbygoogle" style="display:block" data-ad-client="ca-pub-5357405790190342" data-ad-slot="8316640078" data-ad-format="auto" data-full-width-responsive="true"></ins>

<script>

(adsbygoogle = window.adsbygoogle || []).push({});

</script>

<ul>

<li> <p>1. 信号的保真度和该信号的空间和时间顺序无关。即，以同样的方法，改变两个待比较的信号本身的空间或时间排列，它们之间的误差不变。例如，[1 2 3], [3 4 5]两组信号的 MSE 和[3 2 1],[5 3]的 MSE 一样。</p> </li>

<li> <p>2. 误差信号和原信号无关。只要误差信号不变，无论原信号如何，MSE 均不变。例如，对固定误差[1 1 1]，无论加在[1 2 3]产生[2 3 4]还是加在[0 0 0]产生[1 1 1]，MSE 的计算结果不变。</p> </li>

<li> <p>3. 信号的保真度和误差的符号无关。即对于信号[0 0 0]，与之相比较的两个信号[1 2 3]和[-2 -3]被认为和[0 0 0]具有同样的差别。</p> </li>

<li> <p>4. 信号的不同采样点对于信号的保真度具有同样的重要性。</p> </li></ul>

<p><a href="https://link.hacpai.com/forward?goto=https%3A%2F%2Fwww.zhihu.com%2Fquestion%2F24095027%2Fanswer%2F30762001" target="\_blank" rel="nofollow ugc">参考</a></p>