



日期: /

1. 线性回归 linear regression -

$$y = w * x + b + \epsilon$$

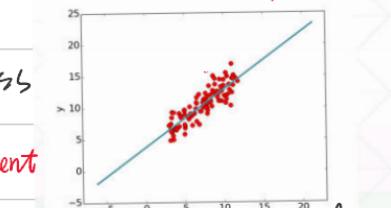
$\epsilon \in N(0,1)$ 噪声 正态分布

$$\text{loss} = \sum_i (w * x_i + b - y_i)^2$$

(误差). minimize loss

→ 梯度下降. Gradient Descent

$$w' = w - lr * \frac{dy}{dw}$$



$$\text{exp. } x = x - 0.05 + \frac{dy}{dx}$$

w 的值向最低点 (导数为0) 方向移动.

learning rate ↗ 步长. 衰减因子.

$$w' = w - lr * \boxed{\frac{\partial \text{loss}}{\partial w}}$$

$$b' = b - lr * \frac{\partial \text{loss}}{\partial b}$$

$$w' * x + b' \rightarrow y$$

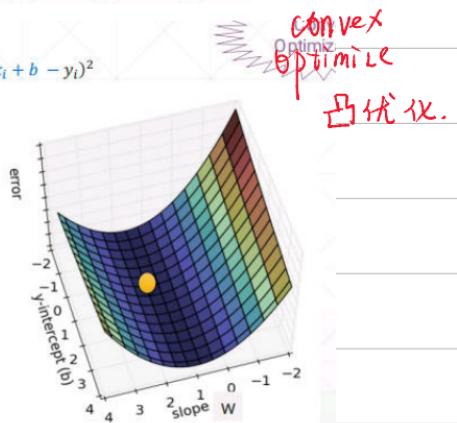
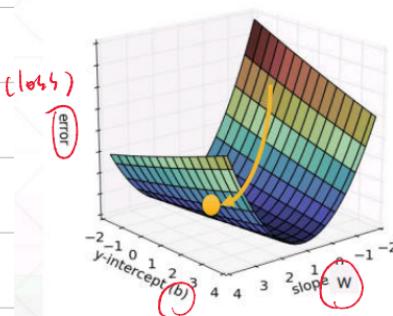
梯度. } 向梯度为负的方向行进.

梯度误差最小
(最优)

w, b.

surface

$$\text{loss} = \sum_i (w * x_i + b - y_i)^2$$



日期: /

Linear Regression.

Logistic Regression 二分类
classification 离散,

点集: $[(x_0, y_0) \dots (x_{99}, y_{99})]$

随机梯度下降. w_0, b_0 $\text{且 } w_0 = b_0 = 0$

Step 1. compute Loss.

$$\text{loss} = \frac{1}{N} \sum_i (w * x_i + b - y_i)^2$$

Step 2. Compute Gradient and update.

$$w' = w - lr * \frac{\partial \text{loss}}{\partial w}$$

$$b' = b - lr * \frac{\partial \text{loss}}{\partial b}$$

计算梯度.

$$\sum_i 2(w * x_i + b - y_i) * x_i$$

$$\frac{\partial \text{loss}}{\partial w} = 2(w * x_0 + b - y_0) * x_0 + 2(w * x_1 + b - y_1) * x_1 \dots$$

$$\frac{\partial \text{loss}}{\partial b} = 2(w * x_0 + b - y_0) * 1 + \dots$$

$$\sum_i 2(w * x_i + b - y_i)$$

日期: /

Step3. set $w = w'$ and loop

循环更新 w, b.

$$w \leftarrow w' \quad b \leftarrow b'$$

离散预测. Discrete Prediction.

计算机视觉, 分类. Image Classification.

手写数字识别. Hand-written Digits Recognition.

Mnist \rightarrow 数据集.

Image. 大小: $[28, 28, 1]$ 灰度图.

打平: 每行相连. \downarrow $[784] \quad 28 \times 28$.

↑ b38图.

Input $x = [b, 784]$

prediction.

- 1. dog = 0, cat = 1, fish = 2, ...
- 2. dog = [1, 0, 0, ...]
- cat = [0, 1, 0, ...]
- fish = [0, 0, 1, ...]

日期: /

regression vs. classification:

$$y = w \cdot x + b \quad y \in \mathbb{R}^d$$

out = $x @ w + b$. 祛除, 预测向量

out = [0.1, 0.8, 0.02, 0.08] 置信度.

pred = argmax(out) 预测 \rightarrow 置信度最大

$$\text{pred} = 1$$

$$\text{label} = 2$$

computation Graph

$$\text{out} = x @ w + b$$

$x = [b, 784] \rightarrow 1 \uparrow \text{batch}$.

$$w = [784, 10] \quad b = 0 \sim 9,$$

$$b = [10] \quad [b, 784] @ [784, 10] + [10]$$

$$\text{out} = [b, 10] \quad [b, 10] + [10] \rightarrow [b, 10]$$

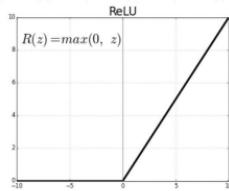
out

日期: /

It's linear!

添加非线性因子.

$$\text{out} = \text{relu}(x @ w + b)$$

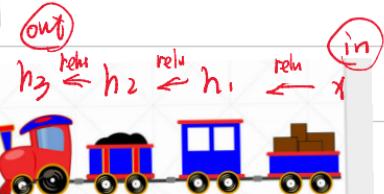


激活函数 = ReLU, sigmoid

h : hidden 隐藏层

经过多道工序

$$\text{out} = \text{relu}(x @ w + b)$$



$$h_1 = \text{relu}(x @ W_1 + b_1)$$

(relu 3次)

$$h_2 = \text{relu}(h_1 @ W_2 + b_2)$$

$$\text{out} = \text{relu}(h_2 @ W_3 + b_3)$$

particularly

每一道工序

- $X = [v1, v2, \dots, v784]$

- $X: [1, 784]$

- $h_1 = \text{relu}(X @ W_1 + b_1)$

- $W_1: [784, 512]$

- $b_1: [1, 512]$

$$\rightarrow [1, 512] \quad (h_1)$$



- $h_2 = \text{relu}(h_1 @ W_2 + b_2)$

- $W_2: [512, 256]$

- $b_2: [256]$

$$\rightarrow [1, 256] \quad (h_2)$$



- $\text{out} = \text{relu}(h_2 @ W_3 + b_3)$

- $W_3: [256, 10]$

- $b_3: [10]$

$$\rightarrow [1, 10] \quad (\text{out})$$



$$[0, 0, 0, 0.01, 0.1, \textcolor{red}{0.8}, 0, \dots]$$

日期: /

- out: [1, 10]

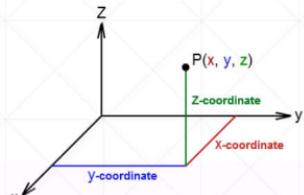
输出 → 确定值, 打标签.

- Y/label: 0~9

one-hotting

eg.: 1 → [0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0]

eg.: 3 → [0, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0]



- 将标签与输出层计算结果
对比

MSE 欧式距离

(矩阵 → 多维度)

$$\sum_i (y_i - \text{out})^2$$

$$\text{out} = \text{relu} \circ \text{relu} \circ \text{relu} [x @ w_1 + b_1] @ w_2 + b_2 @ w_3 + b_3$$

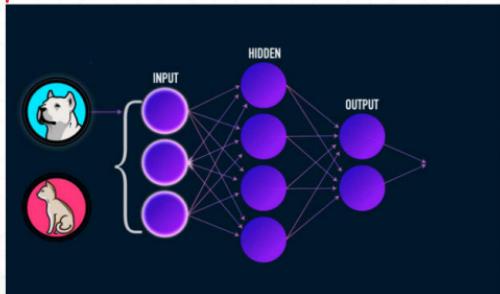
$$\text{pred} = \arg \max(\text{out})$$

$$\text{loss} = \text{MSE}(\text{out}, \text{label})$$

$$\text{minimize } \text{loss}$$

$$\{w_1, b_1; w_2, b_2; w_3, b_3\}$$

与训练集符合度最高.



日期: /

阅读摘要

classification procedure.

↓
将输入映射

Step 1. compute $[h_1, h_2, \text{out}]$

Step 2. compute loss.

Step 3. compute gradient and update.

$[w_1, b_1, w_2, b_2, w_3, b_3]$

Step 4. loop.

优化后
参数.

学习过程:

数据集 \rightarrow 打好标签 (label) 的图像.

计算 y_i \downarrow $\text{label} = [1, 10]$ 矩阵.

通过神经网络. (输入 图像. $[6, 784]$) \rightarrow (输出 $[1, 10]$ 矩阵. \downarrow loop)
 \downarrow (relu) \uparrow 梯度. \uparrow 更新参数

* 每次 loop 都计算 out \downarrow (loop) \downarrow (loss 最小)
~~平均损失~~ (batch size)
每次 loop 都计算 $\sum_r (y_i - \text{out}_i)$ \downarrow 1个 batch (使训练结果与标签相差最小)

即得到训练好的模型.

日期: /

Data container.

list 整数, 字母, 数组, 类.

np.array.

tf.Tensor. 黑体字 数组类型.

Tensor?

scalar = 1.1 标量. dim=0

vector: [1.1], [1.1, 2.2, ...] 向量. dim=1.

matrix: [[1.1, 2.2], [3.3, 4.4]]

tensor = rank > 2 {向量} > 2

tensorflow 极限 $\xrightarrow{\text{运算}}$ 3. imp.

TF is a computing lib.

int, float, double.

bool 布尔型. equal (T/F).

String

日期: /

Create:

整数 = tf.constant(1, dtype=)

(True, False).

(Hello world)

Tensor Property.

① tf.device. 没有值. ("cpu") ("gpu")

结果: a.gpu / a.cpu.

随机数.

② b.numpy. tensor \rightarrow numpy

b.shape. 查看一个 tensor 的 shape.

rank / ndim numpy=1 \rightarrow 3维

name

check Tensor Type.

① isinstance(a, tf.Tensor).

② .tf.is_tensor(b) \rightarrow 返回 true / false.

③ a.dtype.

日期: /

convert

numpy → tf. int64.

① aa = tf. convert_to_tensor(a) / (a, dtype=int32).

② tf.cast(aa, dtype=tf.float32).

tf.double.

tf.int32.

bool返回的是[0,1]→tf.bool

③ bool → int

tf.cast(aa, dtype=tf.int32).

false true
↓ ↓
0 1

b = tf.Variable(a) ↑ tensor
自动包装 b.

↓
b 只有 { b.name
b.trainable. 求导. 可训练. 梯度信息.

To numpy:

a.numpy.

a.tf.one([]). → int(a) / float(a)

↑ 1x1 是 scalar.

日期: /

(6) 建 Tensor:

1. from numpy . list

→ data.

([1. 2]) ↑ size.

tf . convert_to_tensor . (np . ones ([2, 3])).

① tf . zeros ([]) → shape . () ↓
[2, 2] [2, 3, 3]

tf . zeros_like (a) 把 a 传入 shape - 按 a 的形状分配。

② tf . ones . 全为 1 . 同样传入 shape .

③ tf . fill ([2, 2], 10)
↓ shape ↓ 值 .

Normal 正态

正态分布

tf . random . normal ([2, 2], mean=0 , stddev=1)

tf . random . truncated_normal ([2, 2], mean=0 , stddev=1)

裁取正态分布的一部分

Gradient vanish.

用于 sigmoid

日期: /

Uniform . 均匀 $\rightarrow \text{shape}$

`tf.random.uniform([2,2], minval, maxval)`

Random permutation. 随机打乱.

对数分布 (有对应关系的)

`idx = tf.range(10) [0, 9]`

`idx = tf.random.shuffle(idx)`

随机
方法

In [24]: `idx=tf.range(10)` 得到一个放数据的顺序
In [26]: `idx=tf.random.shuffle(idx)`

Out[27]: <tf.Tensor: id=67, shape=(10,), dtype=int32, numpy=array([2, 1, 9, 3, 8, 7, 0, 5, 4, 6])>

In [29]: `a=tf.random.normal([10,784])`

In [31]: `b=tf.random.uniform([10], maxval=10, dtype=tf.int32)`

Out[32]: <tf.Tensor: id=78, shape=(10,), dtype=int32, numpy=array([1, 4, 6, 9, 4, 9, 3, 2, 0, 5])>

In [33]: `a=tf.gather(a, idx)`

In [34]: `b=tf.gather(b, idx)`

Out[35]: <tf.Tensor: id=83, shape=(10,), dtype=int32, numpy=array([6, 4, 5, 9, 0, 2, 1, 9, 4, 5])>

`tf.constant` = `tf.convert_to_tensor`.

每行的维度一致

日期: /

Typical Dim Data.

1. Scalar : 标量、

[]

loss = mse (out, y)

accuracy

LOSS

out

y:

```
In [28]: out=tf.random.uniform([4,10])
<tf.Tensor: id=78, shape=(4, 10), dtype=float32, numpy=
array([[0.11917067, 0.499925 , 0.5349126 , 0.12400055, 0.5961182 , ...
       0.48180282, 0.47066116, 0.96908057, 0.6354896 , 0.44619405]],

In [29]: y=tf.one_hot(y,depth=10) 输出：归一化。
<tf.Tensor: id=95, shape=(4, 10), dtype=float32, numpy=
array([[1., 0., 0., 0., 0., 0., 0., 0., 0., 0.], ...
       [0., 0., 0., 1., 0., 0., 0., 0., 0., 0.]], dtype=float32)>

In [31]: loss=tf.keras.losses.mse(y, out) Σ(y - out)
Out[32]: <tf.Tensor: id=99, shape=(4,), dtype=float32, numpy=array([0.37688 , 0.27573627,
0.5074741 , 0.3111775 ], dtype=float32)>

In [33]: loss=tf.reduce_mean(loss)
Out[34]: <tf.Tensor: id=102, shape=(), dtype=float32, numpy=0.36781698>
```

Dim = 1.

scalar

Vector

Bias: $x @ w + b \rightarrow$ bias

[out_dim] 5 out 维度一致。

矩阵。

```
In [20]: net=layers.Dense(10) → 7行，8→10
In [21]: net.build((4,8))
In [22]: net.kernel W
<tf.Variable 'kernel:0' shape=(8, 10) dtype=float32, numpy=
array([[ 0.19603091, -0.2846143 ,  0.14747244, -0.03926206, -0.02218038,
       0.02417278, -0.31929022, -0.07324755,  0.17448658,  0.06289428]],

In [23]: net.bias b: dim: 10
Out[23]: <tf.Variable 'bias:0' shape=(10,) dtype=float32, numpy=array([0., 0., 0., 0., 0., 0., 0., 0., 0., 0.], dtype=float32)>
```

日期:

Dim=r

Matrix 矩阵.

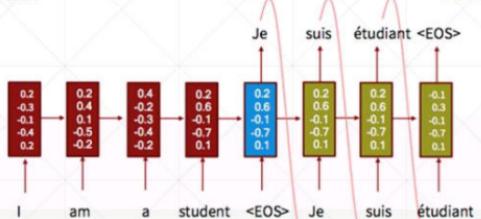
input $x: [b, \text{vec_dim}]$

weight

$x: [b, \text{seq_len}, \text{word_dim}]$

Dim=3 Tensor^r

Encoder



Decoder

```
In [35]: x = tf.random.normal([4, 784])
Out[36]: TensorShape([4, 784])
```

```
In [38]: net = layers.Dense(10) 784 → 10
In [39]: net.build((4, 784)) ✓
```

```
In [40]: net(x).shape [4, 10]
Out[40]: TensorShape([4, 10])
```

```
In [41]: net.kernel.shape
Out[41]: TensorShape([784, 10])
```

```
In [42]: net.bias.shape
Out[42]: TensorShape([10])
```

```
In [4]: (X_train, y_train), (X_test, y_test) =
keras.datasets.imdb.load_data(num_words=10000)
In [5]: x_train = keras.preprocessing.sequence.pad_sequences(X_train, maxlen=80)
In [14]: x_train.shape
```

Out[14]: TensorShape([25000, 80]) ↗ words.

run ✓

In [12]: emb = embedding(x_train) ↗ 100维.

In [13]: emb.shape

Out[13]: TensorShape([25000, 80, 100]) ↗ emb

→ 100维.

In [19]: out=rnn(emb[:4])

In [20]: out.shape

Out[20]: TensorShape([4, 256])

日期: /

Dim=4 Tensor.

image. [b, h, w, 3] \rightarrow R.G.B.

卷积神经网络图.

feature maps. [b, h, w, c]

卷积.

```
In [4]: x=tf.random.normal((4,32,32,3))
```

```
In [5]: net=layers.Conv2D(16,kernel_size=3)
```

```
In [6]: net(x)  
# [4, 32, 32, 16]
```

卷积.

Dim=5 Tensor.

single task = [b, h, w, 3]

meta-learning. (多任务)

任务 [4, 64, 28, 28, 1]

训练合集 \rightarrow task.

日期: /

索引与切片.

basic indexing. 给定每一个维度.

```
 In [4]: a=tf.ones([1,5,5,3])
 In [5]: a[0][0] →
<tf.Tensor: id=16, shape=(5, 3), dtype=float32, numpy=
array([[1., 1., 1.],
       [1., 1., 1.],
       [1., 1., 1.],
       [1., 1., 1.]], dtype=float32)>
→ 3行 3列 [d×o][d×1]

 In [6]: a[0][0][0] → 第一行
Out[6]: <tf.Tensor: id=29, shape=(3,), dtype=float32, numpy=array([1., 1., 1.]),
dtype=float32>

 In [7]: a[0][0][0][2] → 第一个
Out[7]: <tf.Tensor: id=46, shape=(), dtype=float32, numpy=1.0>
```

可读性差、取样方式单一.

Numpy-style indexing.

```
 In [8]: a=tf.random.normal([4,28,28,3])
4张照片. 28×28, 彩色.
 In [9]: a[1].shape → 2张.
Out[9]: TensorShape([28, 28, 3])

 In [10]: a[1,2].shape → 2张的第3行.
Out[10]: TensorShape([28, 3])

 In [11]: a[1,2,3].shape RGB数值.
Out[11]: TensorShape([3])

 In [12]: a[1,2,3,2].shape B通道.
Out[12]: TensorShape([]) 色彩图.(0~255)
(0~1)
```

日期: /

start = end

起
末
切片 背景 $A = B$

```
In [8]: a=tf.range(10)      对前索引进行偏移。  
Out[9]: <tf.Tensor: numpy=array([0, 1, 2, 3, 4, 5, 6, 7, 8, 9])>  
  
In [14]: a[-1:] [L-1=0] → 9 取末  
Out[14]: <tf.Tensor: id=48, shape=(1,), dtype=int32, numpy=array([9])> Vector  
In [15]: a[-2:]  
Out[15]: <tf.Tensor: id=53, shape=(2,), dtype=int32, numpy=array([8, 9])>  
In [16]: a[:-2]  
Out[16]: <tf.Tensor: id=58, shape=(2,), dtype=int32, numpy=array([0, 1])>  
In [17]: a[-1]  
Out[17]: <tf.Tensor: id=63, shape=(9,), dtype=int32, numpy=array([0, 1, 2, 3, 4,  
5, 6, 7, 8])>
```

Indexing by:

* 这是简单开000

多种取样方式

步长 : Step.

▪ start:end:step

▪ ::step

```
In [14]: a.shape #TensorShape([4, 28, 28, 3])  
In [15]: a[0].shape #TensorShape([28, 28, 3])  
  
In [16]: a[0,:,:,:].shape  
Out[16]: TensorShape([28, 28, 3])  
  
In [17]: a[0,1,:,:].shape  
Out[17]: TensorShape([28, 3])  
  
In [18]: a[:, :, :, 0].shape  
Out[18]: TensorShape([4, 28, 28])  
  
In [19]: a[:, :, :, 2].shape  
Out[19]: TensorShape([4, 28, 28])  
  
In [20]: a[:, 0, :, :].shape  
Out[20]: TensorShape([4, 28, 3])
```

```
In [21]: a.shape  
Out[21]: TensorShape([4, 28, 28, 3])  
  
In [22]: a[0:2,:,:,:].shape  
Out[22]: TensorShape([2, 28, 28, 3])  
  
In [23]: a[:, 0:28:2, 0:28:2, 0].shape  
Out[23]: TensorShape([4, 14, 14, 3])  
        ↓↓↓  
In [24]: a[:, 14, :, 14, :].shape  
Out[24]: TensorShape([4, 14, 14, 3])  
        ↓↓↓  
In [25]: a[:, 14, 14, :, :].shape  
Out[25]: TensorShape([4, 14, 14, 3])  
  
In [26]: a[:, :, 0:28:2, 0:28:2, :].shape  
Out[26]: TensorShape([4, 14, 14, 3])
```

日期: /

221 例序.

```
In [27]: a=tf.range(4)
Out[28]: <tf.Tensor: id=118, shape=(4,), dtype=int32, numpy=array([0, 1, 2, 3], dtype=int32)>
In [29]: a[::-1]
Out[29]: <tf.Tensor: id=123, shape=(4,), dtype=int32, numpy=array([3, 2, 1, 0], dtype=int32)>
In [30]: a[::2]
Out[30]: <tf.Tensor: id=128, shape=(2,), dtype=int32, numpy=array([3, 1], dtype=int32)>
In [31]: a[2::-2]
Out[31]: <tf.Tensor: id=133, shape=(2,), dtype=int32, numpy=array([2, 0], dtype=int32)>
```

B ← A

Selective Indexing

[4, 2, 8, 2, 8, 3]
▪ tf.gather

[3, 2, 7, 9, 13] 收集 (无规律).

▪ tf.gather_nd

▪ tf.boolean_mask

[4, 35, 8] 第一个值/接收
将升序、降序、白值索引 (第 x, 3 个值)
轴数
In [46]: tf.gather(a, axis=0, indices=[2, 3]).shape
Out[46]: TensorShape([2, 35, 8]) ↓ [2, 35, 8]
In [47]: a[2:4].shape
Out[47]: TensorShape([2, 35, 8])
4个值, 顺序不同.
In [48]: tf.gather(a, axis=0, indices=[2, 1, 4, 0]).shape
Out[48]: TensorShape([4, 35, 8])
顺序固定.
In [49]: tf.gather(a, axis=1, indices=[2, 3, 7, 9, 16]).shape
Out[49]: TensorShape([4, 5, 8])
输出乱序.
In [50]: tf.gather(a, axis=2, indices=[2, 3, 7]).shape
Out[50]: TensorShape([4, 35, 3])

日期:

tf.gather_nd

- data: [classes, students, subjects]

- What if sample several students and their several subjects?

$$[4, 3] \rightarrow [4, 3] \rightarrow [4, 3]$$

▪ aa = tf.gather(a, axis, [several students])

▪ aaa = tf.gather(aa, axis, [several subjects])

tf.gather_nd

- data: [classes, students, subjects]

- What if sample several (classes and students)?

多个维度指定 index.

看例題

- for instance:

• [class1_student1, class2_student2, class3_student3, class4_student4]
• $\rightarrow [4, 8]$

四个样本，8特征

indices [0] $\rightarrow [3, 8]$
indices [0, 1] $\rightarrow [8]$

In [55]: a.shape
Out[55]: TensorShape([4, 35, 8])

In [60]: tf.gather_nd(a, [0]).shape
Out[60]: TensorShape([35, 8])

In [61]: tf.gather_nd(a, [0, 1]).shape
Out[61]: TensorShape([8])

In [62]: tf.gather_nd(a, [0, 1, 2]).shape
Out[62]: TensorShape([])

In [63]: tf.gather_nd(a, [[0, 1, 2]]).shape
Out[63]: TensorShape([1])

日期:

```
●●●  
In [55]: a.shape  
Out[55]: TensorShape([4, 35, 8])  
  
In [56]: tf.gather_nd(a, [[0,0],[1,1]]).shape 2行 -> 8 ] 指定  
Out[56]: TensorShape([2, 8])  
  
In [57]: tf.gather_nd(a, [[0,0],[1,1],[2,2]]).shape 3行指针 [~.~.~]  
Out[57]: TensorShape([3, 8])  
  
In [58]: tf.gather_nd(a, [[0,0,0],[1,1,1],[2,2,2]]).shape 3行指针 [~.~.~] 指定值  
Out[58]: TensorShape([3])  
  
In [59]: tf.gather_nd(a, [[[0,0,0],[1,1,1],[2,2,2]]]).shape  
Out[59]: TensorShape([1, 3])  
↓ [~.~.~] (一个指针) .
```

tf.boolean_mask

```
●●●  
In [75]: a.shape  
Out[75]: TensorShape([4, 28, 28, 3])  
  
In [76]: tf.boolean_mask(a, mask=[True,True,False,False]).shape  
Out[76]: TensorShape([2, 28, 28, 3])  
  
In [77]: tf.boolean_mask(a, mask=[True,True,False],axis=3).shape  
Out[77]: TensorShape([4, 28, 28, 2])  
↑ 指针  
  
In [78]: a=tf.ones([2,3,4])  
  
In [79]: tf.boolean_mask(a,mask=[[True,False,False],[False,True,True]])  
<tf.Tensor: id=354, shape=(3, 4), dtype=float32, numpy:  
array([[1., 1., 1., 1.],  
       [1., 1., 1., 1.],  
       [1., 1., 1., 1.]], dtype=float32)>
```

日期: /

维度变换

Outline

- shape, ndim
- reshape
- expand_dims/squeeze
- transpose
- broadcast_to

batch row column channel.
 $[b, 28, 28, 1]$
 ↑ height
 ↗ width.
 ↘ axis

View

- 每张 [b, 28, 28] content, [b, h, w] view1 $\{h \times w\}$ view2 $\{b \times h \times w\}$ 数量元被坏
- $\rightarrow [b, 28*28]$ view1 - 行 - (无行) view2 -
 - $\rightarrow [b, 2, 14*28]$ view1 - 列 -
 - $\rightarrow [b, 28, 28, 1]$ view1 - 章节 - view2 -

5

5

5

Reshape

```

In [80]: a=tf.random.normal([4,28,28,3])
Out[81]: TensorShape([4, 28, 28, 3])

In [82]: a.shape, a.ndim
Out[82]: (TensorShape([4, 28, 28, 3]), 4)

In [83]: tf.reshape(a,[4,784,3]).shape → [b, pixel, c]
Out[83]: TensorShape([4, 784, 3])

In [84]: tf.reshape(a,[4,(-1,3)]).shape
Out[84]: TensorShape([4, 784, 3]) View1 - 自由计算.

In [85]: tf.reshape(a,[4,784*3]).shape
Out[85]: TensorShape([4, 2352]) View1 - 没有指定 channel / 什么
Out[85]: TensorShape([4, 2352]) data point. 784*3. 元素数量.

In [86]: tf.reshape(a,[4,-1]).shape
Out[86]: TensorShape([4, 2352])

```

没指定 channel / 什么
 784*3. 元素数量.



Reshape is flexible

```
...  

In [80]: a=tf.random.normal([4,28,28,3])  

Out[81]: TensorShape([4, 28, 28, 3])  

In [87]: tf.reshape(tf.reshape(a,[4,-1]),[4,28,28,3]).shape  

Out[87]: TensorShape([4, 28, 28, 3])  

In [88]: tf.reshape(tf.reshape(a,[4,-1]),[4,14,56,3]).shape  

Out[88]: TensorShape([4, 14, 56, 3])  

In [89]: tf.reshape(tf.reshape(a,[4,-1]),[4,1,784,3]).shape  

Out[89]: TensorShape([4, 1, 784, 3])
```

size 括号即可。
输出相同。
意义？

Reshape could lead to potential bugs!

- images: [4, 28, 28, 3]
 - [b, h, w, s]
- reshape to: [4, 784, 3]
 - [b, pixel, s]
- [4, 784, 3] $\xrightarrow{\text{height: 28, width: 28}}$ [4, 28, 28, 3]
- [4, 784, 3] $\xrightarrow{\text{height: 14, width: 56}}$ [4, 14, 56, 3]
- [4, 784, 3] $\xrightarrow{\text{width: 28, height: 28}}$ [4, 28, 28, 3]

恢复？ 需要 content 信息。
即：图片的长宽，顺序。

tf.transpose 改变 content. 转置.

$[h,w] \rightarrow [w,h]$

```
In [93]: a=tf.random.normal((4,3,2,1))
```

```
In [94]: a.shape  
Out[94]: TensorShape([4, 3, 2, 1])
```

```
In [95]: tf.transpose(a).shape  
Out[95]: TensorShape([1, 2, 3, 4])  

          全部转置。  

          ↓  

          Now 放到 1, 3, 1, 2 → [1, 3, 2, 4]
```

```
In [97]: tf.transpose(a,perm=[0,1,3,2]).shape  
Out[97]: TensorShape([4, 3, 1, 2])  

          ↓  

          重新放回 0, 1, 2, 3 → [4, 3, 1, 2]
```

日期: /

→ [b, 3, h, w] [b, h, w, 3] 行在前, 列在后

```
In [98]: a=tf.random.normal([4,28,28,3])  
In [99]: tf.transpose(a,[0,2,1,3]).shape  
Out[99]: TensorShape([4, 28, 28, 3])  
  
In [101]: tf.transpose(a,[0,3,2,1]).shape  
Out[101]: TensorShape([4, 3, 28, 28])  
  
In [102]: tf.transpose(a,[0,3,1,2]).shape  
Out[102]: TensorShape([4, 3, 28, 28])
```

Expand dim

- a: [classes, students, classes]
 - [4, 35, 8]
- add school dim *增加一个轴* (维数)
 - [1, 4, 35, 8] + [1, 4, 35, 8] *axis dim*
 - [2, 4, 35, 8]

```
In [103]: a=tf.random.normal([4,35,8])  
In [105]: tf.expand_dims(a,axis=0).shape  
Out[105]: TensorShape([1, 4, 35, 8])  
  
In [106]: tf.expand_dims(a,axis=3).shape  
Out[106]: TensorShape([4, 35, 8, 1])
```

axis

0 1 2

```
In [103]: a=tf.random.normal([4,35,8])  
In [105]: tf.expand_dims(a,axis=0).shape  
Out[105]: TensorShape([1, 4, 35, 8])  
          轴 = 轴前 00  
In [106]: tf.expand_dims(a,axis=3).shape  
Out[106]: TensorShape([4, 35, 8, 1])  
          轴 = 轴后 00  
In [107]: tf.expand_dims(a,axis=-1).shape  
Out[107]: TensorShape([4, 35, 8, 1])  
          轴 = 轴后 00  
In [108]: tf.expand_dims(a,axis=-4).shape  
Out[108]: TensorShape([1, 4, 35, 8])
```

日期: /

Squeeze dim

- Only squeeze for shape=1 dim P.能用 shape = 1

- [4, 35, 8, 1]
- [1, 4, 35, 8]
- [1, 4, 35, 1]

Broadcasting

- expand
- without copying data
- tf.broadcast_to

key idea. Insert 1 dim if needed

插值，自动补充。

① 小维度先对齐。

② 插入 1. 维度 → 扩张。

a [4, 32, 32, 3]

b. [3] \rightarrow [1, 1, 1, 32] \rightarrow [4, 32, 32, 3]

日期: /

[4,3]	[4,3]	=	[10,10,10,20,20,20,30,30,30]	+ [0,1,2,0,1,2,0,1,2]	=	[10,10,10,20,20,20,30,30,30]	+ [0,1,2,0,1,2,0,1,2]	=	[10,10,10,20,20,20,30,30,30]	+ [0,1,2,0,1,2,0,1,2]	=	[10,10,10,20,20,20,30,30,30]	+ [0,1,2,0,1,2,0,1,2]	=	[10,10,10,20,20,20,30,30,30]
[4,3]	[1,3]	=	[0,0,0,10,20,30,10,20,30]	+ [0,1,2,0,1,2,0,1,2]	=	[0,0,0,10,20,30,10,20,30]	+ [0,1,2,0,1,2,0,1,2]	=	[0,0,0,10,20,30,10,20,30]	+ [0,1,2,0,1,2,0,1,2]	=	[0,0,0,10,20,30,10,20,30]	+ [0,1,2,0,1,2,0,1,2]	=	[0,0,0,10,20,30,10,20,30]
[4,1]	[1,3]	=	[0,10,20,30]	+ [0,1,2,0,1,2,0,1,2]	=	[0,10,20,30]	+ [0,1,2,0,1,2,0,1,2]	=	[0,10,20,30]	+ [0,1,2,0,1,2,0,1,2]	=	[0,10,20,30]	+ [0,1,2,0,1,2,0,1,2]	=	[0,10,20,30]

没有意义的数学!!!

意义?

▪ When it has no axis

- Create a new concept

• [classes, students, scores] + [scores]

大轴数 4 35 8 小轴数.

默认对所有概念适用

无序数 / 学生概念

▪ When it has dim of size 1

- Treat it shared by all

• [classes, students, scores] + [students, 1]

加偏置.

对某些学生加分、

[35, 1] 前多少学生成绩

日期: /

Why broadcasting?

tf.tile 扩展内存.

1. for real demanding

- [classes, students, scores]
- Add bias for every student: +5 score
- [4, 32, 8] + [4, 32, 8]
- [4, 32, 8] + [5.0] [1] → [4, 32, 8]

偏移

每个学生加分.

2. memory consumption

- [4, 32, 8] → 1024 字节数. 1024×4
- bias=[8]: [5.0, 5.0, 5.0, ...] → 8 偏移

8×4 节省内存.

Broadcastable:

[1] ✓

a [4] ✗

b [1, 3]

Match from Last dim!

底层维度.

- If current dim=1, expand to same
- If either has no dim, insert one dim and expand to same
- otherwise, NOT broadcastable

对齐.

广播维度

相等.

广播维

度.

```
In [25]: x=tf.random.normal([4,32,32,3])  
In [27]: a+tf.random.normal([3]).shape  
Out[27]: TensorShape([4, 32, 32, 3])  
  
In [28]: (x+tf.random.normal([32,32,1])).shape  
Out[28]: TensorShape([4, 32, 32, 3])  
  
In [29]: (x+tf.random.normal([4,1,1,1])).shape  
Out[29]: TensorShape([4, 32, 32, 3])  
  
In [31]: (x+tf.random.normal([1,1,1,1])).shape  
InvalidArgumentError: Incompatible shapes: [4, 32, 32, 3] vs. [1,4,1,1] [Op:Add]  
name: add/  
    ↪ 对应，两个不对应
```

```
In [35]: x.shape  
Out[35]: TensorShape([4, 32, 32, 3])  
  
In [36]: (x+tf.random.normal([4,1,1,1])).shape  
Out[36]: TensorShape([4, 32, 32, 3])  
    ↪ 对应  
  
In [37]: b=tf.broadcast_to(tf.random.normal([4,1,1,1]), [4,32,32,3])  
        ↪ 4维的shape.  
  
In [38]: b.shape  
Out[38]: TensorShape([4, 32, 32, 3])
```

```
In [4]: a=tf.ones([3,4])  
In [5]: a1=tf.broadcast_to(a, [2,3,4])  
<tf.Tensor: id=7, shape=(2, 3, 4), dtype=float32, numpy=  
array([[[1., 1., 1., 1.],  
       [1., 1., 1., 1.],  
       [1., 1., 1., 1.]],  
      [[1., 1., 1., 1.],  
       [1., 1., 1., 1.],  
       [1., 1., 1., 1.]]], dtype=float32)>  
  
In [7]: a2=tf.expand_dims(a, axis=0)  
Out[8]: TensorShape([1, 3, 4])  
In [10]: a2=tf.tile(a2, [2,1,1])  
<tf.Tensor: id=12, shape=(2, 3, 4), dtype=float32,  
numpy=[[[[1., 1., 1., 1.],  
       [1., 1., 1., 1.],  
       [1., 1., 1., 1.]],  
      [[1., 1., 1., 1.],  
       [1., 1., 1., 1.],  
       [1., 1., 1., 1.]]], dtype=float32)>
```

日期: /

数学运算

- $+-*/$ 加减乘除
 - $\text{平方: } \underline{\text{**}}, \text{pow}, \text{square}$
 - sqrt
 - $\text{/}, \underline{\%}$ 取余
 - exp, log 底数
 - $@, \text{matmul}$ 矩阵
 - linear layer
- element-wise 对应元素
 - $+-*/$
- matrix-wise 矩阵相乘。
▪ $@, \text{matmul}$ $[b, 3, 4] @ [b, 4, 5]$
- dim-wise 对某一个维度。
▪ $\text{reduce_mean/max/min/sum}$

```
          + - * / 分是 x.  
In [134]: b=tf.fill([2,2],2) ~ 什么一样。  
In [135]: a=tf.ones([2,2]) ~ 什么一样。  
  
In [136]: a+b, a-b, a*b, a/b  
(<tf.Tensor: id=462, shape=(2, 2), dtype=float32, numpy=  
array([[3., 3.],  
       [3., 3.]], dtype=float32),  
<tf.Tensor: id=463, shape=(2, 2), dtype=float32, numpy=  
array([[-1., -1.],  
       [-1., -1.]], dtype=float32),  
<tf.Tensor: id=464, shape=(2, 2), dtype=float32, numpy=  
array([[2., 2.],  
       [2., 2.]], dtype=float32),  
<tf.Tensor: id=465, shape=(2, 2), dtype=float32, numpy=  
array([[0.5, 0.5],  
       [0.5, 0.5]], dtype=float32))  
  
In [137]: b//a, b%a  
(<tf.Tensor: id=470, shape=(2, 2), dtype=float32, numpy=  
array([[2., 2.],  
       [2., 2.]], dtype=float32),  
<tf.Tensor: id=471, shape=(2, 2), dtype=float32, numpy=  
array([[0., 0.],  
       [0., 0.]], dtype=float32))
```

```
In [138]: a  
<tf.Tensor: id=461, shape=(2, 2), dtype=float32, numpy=  
array([[1., 1.],  
       [1., 1.]], dtype=float32))  
  
In [140]: tf.math.log(a) log e. 无法使用 log e. log 10.  
<tf.Tensor: id=475, shape=(2, 2), dtype=float32, numpy=  
array([[0., 0.],  
       [0., 0.]], dtype=float32))  
  
In [141]: tf.exp(a)  
<tf.Tensor: id=477, shape=(2, 2), dtype=float32, numpy=  
array([2.7182817, 2.7182817],  
     [2.7182817, 2.7182817]), dtype=float32)>
```

日期:

使用 \log_2 , \log_{10} ?

```
In [22]: tf.math.log(8.)/tf.math.log(2.)  
Out[22]: <tf.Tensor: id=54, shape=(), dtype=float32,  
numpy=3.0>  
In [23]: tf.math.log(100.)/tf.math.log(10.)  
Out[23]: <tf.Tensor: id=60, shape=(), dtype=float32,  
numpy=2.0>
```

$\frac{\log_b}{\log_2} \cdot \log_b$.

pow sqrt

```
In [142]: b  
<tf.Tensor: id=458, shape=(2, 2), dtype=float32, numpy=  
array([[2., 2.],  
       [2., 2.]], dtype=float32)>  
  
In [143]: tf.pow(b,3)  
<tf.Tensor: id=481, shape=(2, 2), dtype=float32, numpy=  
array([[8., 8.],  
       [8., 8.]], dtype=float32)>  
  
In [144]: b**3  
<tf.Tensor: id=484, shape=(2, 2), dtype=float32, numpy=  
array([[8., 8.],  
       [8., 8.]], dtype=float32)>  
  
In [145]: tf.sqrt(b)  
<tf.Tensor: id=486, shape=(2, 2), dtype=float32, numpy=  
array([[1.4142135, 1.4142135],  
       [1.4142135, 1.4142135]], dtype=float32)>
```

矩阵

② MatMul /

```
In [146]: a,b  
(<tf.Tensor: id=461, shape=(2, 2), dtype=float32, numpy=  
array([[1., 1.],  
       [1., 1.]], dtype=float32),  
<tf.Tensor: id=458, shape=(2, 2), dtype=float32, numpy=  
array([[2., 2.],  
       [2., 2.]], dtype=float32))  
  
In [147]: a@b  
<tf.Tensor: id=490, shape=(2, 2), dtype=float32, numpy=  
array([[4., 4.],  
       [4., 4.]], dtype=float32)>  
  
In [148]: tf.matmul(a,b) // 如果是 2D tensor  
<tf.Tensor: id=492, shape=(2, 2), dtype=float32, numpy=  
array([[4., 4.],  
       [4., 4.]], dtype=float32)>
```

日期: 11/11

```
...  
In [150]: a=tf.ones([4,2,3])          batch → 重複4次.  
In [151]: b=tf.fill([4,3,5], 2.)      [2,3]  
In [152]: a@b  
<tf.Tensor: id=503, shape=(4, 2, 5), dtype=float32, numpy=  
array([[6., 6., 6., 6., 6.],  
       ...  
       [6., 6., 6., 6., 6.]])], dtype=float32>  
  
In [153]: tf.matmul(a,b)  
<tf.Tensor: id=505, shape=(4, 2, 5), dtype=float32, numpy=  
array([[6., 6., 6., 6., 6.],  
       ...  
       [6., 6., 6., 6., 6.]])], dtype=float32>
```

```
...  
In [164]: a.shape # TensorShape([4, 2, 3])  
In [165]: b.shape # TensorShape([3, 5])  
In [166]: bb=tf.broadcast_to(b, [4,3,5])  
In [167]: a@bb  
<tf.Tensor: id=516, shape=(4, 2, 5), dtype=float32, numpy=  
array([[6., 6., 6., 6., 6.],  
       [6., 6., 6., 6., 6.], ...  
       [6., 6., 6., 6., 6.],  
       [6., 6., 6., 6., 6.]])], dtype=float32>
```

Recap

概念: relu.

- $y = w * x + b$ 這等於
- $Y = X @ W + b$

$$\begin{bmatrix} x_0^0 & x_0^1 \\ x_1^0 & x_1^1 \end{bmatrix} \begin{bmatrix} w_{00} & w_{01} & w_{02} \\ w_{10} & w_{11} & w_{12} \end{bmatrix} + [b_0, b_1, b_2] \rightarrow \begin{bmatrix} y_0^0 & y_0^1 & y_0^2 \\ y_1^0 & y_1^1 & y_1^2 \end{bmatrix}$$

- $[b, 2] \rightarrow [b, 3]$

$$y_0^0 = x_0^0 W_{00} + x_0^1 W_{10} + b_0$$

$$Y = X@W + b$$

日期:

```
In [168]: x=tf.ones([4,2]) > [4,1]
In [169]: W=tf.ones([2,1])
In [170]: b=tf.constant(0.1) → After broadcast
In [171]: x@W+b
<tf.Tensor: id=526, shape=(4, 1), dtype=float32, numpy=
array([[2.1],
       [2.1],
       [2.1],
       [2.1]], dtype=float32)>
```

$$out = \text{relu}(X@W + b)$$

```
In [171]: x@W+b
<tf.Tensor: id=526, shape=(4, 1), dtype=float32, numpy=
array([[2.1],
       [2.1],
       [2.1],
       [2.1]], dtype=float32)>
In [172]: out=x@W+b    非线性因子，去负数。
In [173]: out=tf.nn.relu(out)
<tf.Tensor: id=530, shape=(4, 1), dtype=float32, numpy=
array([[2.1],
       [2.1],
       [2.1],
       [2.1]], dtype=float32)>
```

日期: /

前向传播 (张量) 实战

增加复杂度

Recap:

- out = $\text{relu} \{ \text{relu} \{ \text{relu} \{ X @ W_1 + b_1 \} @ W_2 + b_2 \} @ W_3 + b_3 \}$

非线性

层

- pred = $\text{argmax}(\text{out})$

- loss = $MSE(\text{out}, \text{label})$

误差

- minimize loss

使误差最小

- $[W'_1, b'_1, W'_2, b'_2, W'_3, b'_3]$

张量的合并与分割 . Merge and split

- tf.concat

连接

- tf.split

分割

- tf.stack

堆叠

- tf.unstack

拆分

关键在维度!!!

concat

- Statistics about scores

- [class1-4, students, scores]
- [class5-6, students, scores]

$$\begin{bmatrix} 4, 35, 8 \\ 2, 35, 8 \end{bmatrix}$$

4/2+2=4维

$$[6, 35, 8]$$

注: 其他维度一致, 小会增加维数.

```
In [3]: a=tf.ones([4,35,8])
In [4]: b=tf.ones([2,35,8])
In [6]: c=tf.concat([a,b],axis=0)
In [7]: c.shape
Out[7]: TensorShape([6, 35, 8])
In [8]: a=tf.ones([4,32,8])
In [9]: b=tf.ones([4,3,8])
In [10]: tf.concat([a,b],axis=1).shape
Out[10]: TensorShape([4, 35, 8])
```

拼接维度(4+2)

插入维度(32+3)

日期: /

Stake.

▪ Statistics about scores

- + (• School1:[classes, students, scores]
- School2:[classes, students, scores]
- [schools, classes, students, scores]

注: 所有已存储数据都相等.

```
In [19]: a.shape  
Out[19]: TensorShape([4, 35, 8])  
In [23]: b.shape  
Out[23]: TensorShape([4, 35, 8])  
  
In [20]: tf.concat([a,b],axis=-1).shape  
Out[20]: TensorShape([4, 35, 16])  
  
In [21]: tf.stack([a,b],axis=0).shape  
Out[21]: TensorShape([2, 4, 35, 8])  
In [22]: tf.stack([a,b],axis=1).shape  
Out[22]: TensorShape([4, 35, 8, 2])
```

在最后一步
在第 2 步
在第 3 步
在第 4 步

Un Stack.

```
In [30]: a.shape # TensorShape([4, 35, 8])  
In [32]: b=tf.ones([4,35,8])  
  
In [33]: c=tf.stack([a,b])  
In [34]: c.shape  
Out[34]: TensorShape([2, 4, 35, 8])  
  
In [35]: aa,bb=tf.unstack(c,axis=0)  
In [36]: aa.shape,bb.shape  
Out[36]: (TensorShape([4, 35, 8]), TensorShape([4, 35, 8]))  
  
# [2, 4, 35, 8]  
In [41]: res=tf.unstack(c,axis=3)  
  
In [42]: res[0].shape, res[7].shape  
Out[42]: (TensorShape([2, 4, 35]), TensorShape([2, 4, 35]))
```

Split

假定行数的数据

```
# [2, 4, 35, 8]  
In [43]: res=tf.unstack(c,axis=3)  
In [44]: len(res)  
Out[44]: 8  
In [45]: res=tf.split(c,axis=3, num_or_size_splits=2)  
In [46]: len(res)  
Out[46]: 2  
In [47]: res[0].shape  
Out[47]: TensorShape([2, 4, 35, 4])  
In [48]: res[1].shape  
Out[48]: TensorShape([2, 4, 35, 4])  
In [49]: res[0].shape, res[2].shape  
Out[49]: (TensorShape([2, 4, 35, 2]), TensorShape([2, 4, 35, 4]))
```

日期: /

数据统计.

▪ tf.norm 范数.

▪ tf.reduce_min/max

▪ tf.argmax/argmin

▪ tf.equal

▪ tf.unique

Vector Norm

Eukl. Norm $\|x\|_2 = [\sum_k x_k^2]^{1/2}$

二范数.

Max.norm $\|x\|_\infty = \max_k |x_k|$

L_1 -Norm $\|x\|_1 = \sum_k |x_k|$

-范数.

▪ Here talks about Vector Norm

tf. norm ().

```
●●●  
In [50]: a=tf.ones([2,2]) [ 1 1 ]  
In [52]: tf.norm(a)  
Out[52]: <tf.Tensor: id=192, shape=(), dtype=float32, numpy=2.0>  
  
In [54]: tf.sqrt(tf.reduce_sum(tf.square(a)))  
Out[54]: <tf.Tensor: id=197, shape=(), dtype=float32, numpy=2.0>  
  
In [55]: a=tf.ones([4,28,28,3])  
  
In [56]: tf.norm(a)  
Out[56]: <tf.Tensor: id=206, shape=(), dtype=float32, numpy=96.99484>  
  
In [57]: tf.sqrt(tf.reduce_sum(tf.square(a)))  
Out[57]: <tf.Tensor: id=211, shape=(), dtype=float32, numpy=96.99484>
```

```
●●●  
In [68]: b = tf.ones([2,2])  
  
In [69]: tf.norm(b)  
Out[69]: <tf.Tensor: id=250, shape=(), dtype=float32, numpy=2.0>  
  
In [73]: tf.norm(b,ord=2, axis=1) J+I  
Out[73]: <tf.Tensor: id=271, shape=(2,), dtype=float32, numpy=array([1.4142135, 1.4142135], dtype=float32)>  
  
In [78]: tf.norm(b,ord=1)  
Out[78]: <tf.Tensor: id=255, shape=(), dtype=float32, numpy=4.0>  
  
In [79]: tf.norm(b,ord=1, axis=0)  
Out[79]: <tf.Tensor: id=260, shape=(2,), dtype=float32, numpy=array([2., 2.], dtype=float32)>  
  
In [72]: tf.norm(b,ord=1, axis=1)  
Out[72]: <tf.Tensor: id=265, shape=(2,), dtype=float32, numpy=array([2., 2.], dtype=float32)>
```

日期:

reduce_min/max/mean

```
In [76]: a=tf.random.normal([4,10])
In [78]: tf.reduce_min(a),tf.reduce_max(a),tf.reduce_mean(a)
Out[78]:
<tf.Tensor: id=283, shape=(), dtype=float32, numpy=-1.1872448>,
<tf.Tensor: id=285, shape=(), dtype=float32, numpy=2.1355827>,
<tf.Tensor: id=287, shape=(), dtype=float32, numpy=0.3523524>
In [79]: tf.reduce_min(a, axis=1),tf.reduce_max(a, axis=1),tf.reduce_mean(a, axis=1)
Out[79]:
<tf.Tensor: id=292, shape=(4,), dtype=float32, numpy=array([-0.3937837,
-1.1872448, -1.0798895, -1.1366792], dtype=float32)>,
<tf.Tensor: id=294, shape=(4,), dtype=float32, numpy=array([1.9718986, 1.1612172,
2.1353827, 2.0984378], dtype=float32)>,
<tf.Tensor: id=296, shape=(4,), dtype=float32, numpy=array([ 0.61504304,
-0.01389184, 0.606747 , 0.20151145], dtype=float32)>
```

減小

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \\ 2 & 2 \end{bmatrix} \xrightarrow{\text{max}} \begin{bmatrix} 2 \\ 4 \\ 2 \end{bmatrix}$$

每一行找最大值

max
mean.

argmax/argmin

得到位置。

```
axis=0
In [80]: a.shape
Out[80]: TensorShape([4, 10])
In [81]: tf.argmax(a).shape
Out[81]: TensorShape([10])
In [83]: tf.argmax(a)
Out[83]: <tf.Tensor: id=305, shape=(10,), dtype=int64, numpy=array([0, 0, 2, 3, 1, 3, 0, 1, 2, 0])>
In [82]: tf.argmax(a).shape
Out[82]: TensorShape([10]) → 位置
```

tf.equal

比较。

```
In [44]: a=tf.constant([1,2,3,2,5])
In [45]: b=tf.range(5) → [0,1,2,3,4] false → 0
In [46]: tf.equal(a,b)
Out[46]: <tf.Tensor: id=170, shape=(5,), dtype=bool, numpy=array([False, False, True, False, False])>
In [47]: res=tf.equal(a,b) → 真偽値
In [48]: tf.reduce_sum(tf.cast(res, dtype=tf.int32))
Out[48]: <tf.Tensor: id=175, shape=(), dtype=int32, numpy=8>
```

```
In [99]: a
<tf.Tensor: id=308, shape=(2, 3), dtype=float32, numpy=
array([[0.1 , 0.2 , 0.7 ],
       [0.9 , 0.05, 0.05]], dtype=float32)>
In [100]: pred=tf.cast(tf.argmax(a, axis=1), dtype=tf.int32)
Out[101]: <tf.Tensor: id=324, shape=(2,), dtype=int32, numpy=array([2, 0])>
In [110]: y
Out[110]: <tf.Tensor: id=328, shape=(2,), dtype=int32, numpy=array([2, 1], dtype=int32)>
In [112]: tf.equal(y,pred)
Out[112]: <tf.Tensor: id=335, shape=(2,), dtype=bool, numpy=array([ True, False])>
In [113]: correct=tf.reduce_sum(tf.cast(tf.equal(y,pred), dtype=tf.int32))
In [114]: correct_
Out[114]: <tf.Tensor: id=340, shape=(), dtype=int32, numpy=1>
In [115]: correct/2
Out[115]: <tf.Tensor: id=345, shape=(), dtype=float64, numpy=0.5>
```

两个概率和

[0,1] → [0,1]

错误的 0.1 二值化。

日期:
tf.unique

去重.

去掉重复
元素

```
In [116]: a=tf.range(5)
In [117]: tf.unique(a)
Out[117]: Unique(y<tf.Tensor: id=351, shape=(5,), dtype=int32, numpy=array([0, 1, 2, 3, 4], dtype=int32>, idx<tf.Tensor: id=352, shape=(5,), dtype=int32, numpy=array([0, 1, 2, 3, 4], dtype=int32>))
In [118]: a=tf.constant([4,2,2,4,3]) 可以重
          0 1 0 4. ↓ [0, 1, 2] 去重.
In [119]: tf.unique(a)
Out[119]: Unique(y<tf.Tensor: id=356, shape=(3,), dtype=int32, numpy=array([4, 2, 3], dtype=int32>, idx<tf.Tensor: id=357, shape=(5,), dtype=int32, numpy=array([0, 1, 1, 0, 2], dtype=int32>))
```

在不重中的索引号.

排序 序号

▪ Sort.argsort

▪ Topk 前--↑, 后--↓

▪ Top-5 Acc.

Sort.argsort

```
In [86]: a=tf.random.shuffle(tf.range(5)) #numpy=array([2, 0, 3, 4, 1])
In [87]: tf.argsort(a,direction='DESCENDING') 降序.
Out[87]: <tf.Tensor: id=397, shape=(5,), dtype=int32, numpy=array([4, 3, 2, 1, 0])>
In [88]: tf.argsort(a,direction='DESCENDING')
Out[88]: <tf.Tensor: id=400, shape=(5,), dtype=int32, numpy=array([0, 1, 2, 3, 4])>
In [89]: idx=tf.argsort(a,direction='DESCENDING')
In [90]: tf.gather(a, idx) 从大到小的位置.
Out[90]: <tf.Tensor: id=422, shape=(5,), dtype=int32, numpy=array([4, 3, 2, 1, 0])>
```

高维 → 对某一个维度.

0~9

```
In [95]: a=tf.random.uniform([3,3],maxval=10, dtype=tf.int32)
array([[4, 6, 8],
       [9, 4, 7],
       [4, 5, 11]])
In [97]: tf.sort(a) 每行
array([[4, 6, 8],
       [4, 7, 9],
       [1, 4, 5]])
In [98]: tf.sort(a, direction='DESCENDING')
array([[8, 6, 4],
       [9, 7, 4],
       [5, 4, 11]])
In [99]: idx=tf.argsort(a)
array([[0, 1, 2],
       [1, 2, 0],
       [2, 0, 1]])
```

日期:

Top-K

indices

values.

accuracy

Top-5

位置对错

• Prob:[0.1, 0.2, 0.3, 0.4]

• Label:[2] 真实值.

• Only consider top-1 prediction: [3]

$\frac{1}{3} = 33\% \rightarrow$ 正确.

• Only consider top-2 prediction: [3, 2]

$\frac{1}{2} = 50\%$

正确/错误.

• Only consider top-3 prediction: [3, 2, 1]

$\frac{1}{3} = 33\%$

正确/错误.

● ● ●

看有错没

● ● ●

In [59]: prob=tf.constant([[0.1, 0.2, 0.7], [0.2, 0.7, 0.1]])

In [60]: target=tf.constant([2, 0]) 真实[2, 0]

In [61]: k_b=tf.math.top_k(prob, 3).indices

array([[2, 1, 0], [1, 0, 2]]) 位置-

In [63]: k_b=tf.transpose(k_b, [1, 0]) Top3>(0%)

array([[2, 1], [1, 0], [0, 2]]) Top1 → $\frac{1}{3} = 33\%$, Top2 → $\frac{1}{2} = 50\%$, Top3 → $\frac{1}{1} = 100\%$

In [65]: target=tf.broadcast_to(target, [3, 2]) 行广播

<tf.Tensor: id=214, shape=(3, 2), dtype=int32, numpy=

array([[2, 0], [2, 0], [2, 0]])>

日期

```
def accuracy(output, target, topk=(1,)):
    maxk = max(topk)
    batch_size = target.shape[0]
    pred = tf.math.top_k(output, maxk).indices
    pred = tf.transpose(pred, perm=[1, 0])
    target_ = tf.broadcast_to(target, pred.shape)
    correct = tf.equal(pred, target_) # [k, b]
    res = []
    for k in topk:
        correct_k = tf.cast(tf.reshape(correct[:k], [-1]), dtype=tf.float32)
        correct_k = tf.reduce_sum(correct_k)
        acc = float(correct_k / batch_size) → 精确度.
        res.append(acc)
    return res
```

target:
[670 671
670 671
670 671]

解題

Top1: [b, b,]
Top2: [b, b,]
Top3: [b, b,]
pred

此段為：為何？ 轉化成 [0,1]
↓ 為什麼。

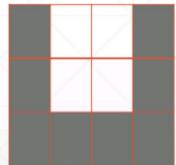
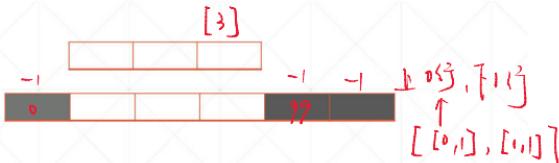
日期: /

填充与复制

pad

tile

broadcast_to



上行
下行
左列 右列

```
●●●  
In [120]: a=tf.reshape(tf.range(9), [3,3])  
<tf.Tensor: id=365, shape=(3, 3), dtype=int32, numpy=  
array([[0, 1, 2],  
       [3, 4, 5],  
       [6, 7, 8]], dtype=int32)>
```

In [123]: tf.pad(a, [[0,0],[0,0]]) %} 没有 padding .

```
Out[123]:  
<tf.Tensor: id=369, shape=(3, 3), dtype=int32, numpy=  
array([[0, 1, 2],  
       [3, 4, 5],  
       [6, 7, 8]], dtype=int32)>
```

In [124]: tf.pad(a, [[1,0],[0,0]]) 上行 .

```
Out[124]:  
<tf.Tensor: id=372, shape=(4, 3), dtype=int32, numpy=  
array([[0, 0, 0],  
       [0, 1, 2],  
       [3, 4, 5],  
       [6, 7, 8]], dtype=int32)>
```

●●●
In [125]: tf.pad(a, [[1,1],[0,0]])
<tf.Tensor: id=375, shape=(5, 3), dtype=int32, numpy=
array([[0, 0, 0],
 [0, 1, 2],
 [3, 4, 5],
 [6, 7, 8],
 [0, 0, 0]], dtype=int32)>

In [126]: tf.pad(a, [[1,1],[1,0]]) 上下行、左-右 .

```
<tf.Tensor: id=378, shape=(5, 4), dtype=int32, numpy=  
array([[0, 0, 0, 0],  
       [0, 0, 1, 2],  
       [0, 3, 4, 5],  
       [0, 6, 7, 8],  
       [0, 0, 0, 0]], dtype=int32)>
```

In [127]: tf.pad(a, [[1,1],[1,1]]) + 下行 .

```
<tf.Tensor: id=381, shape=(5, 5), dtype=int32, numpy=  
array([[0, 0, 0, 0, 0],  
       [0, 0, 1, 2, 0],  
       [0, 3, 4, 5, 0],  
       [0, 6, 7, 8, 0],  
       [0, 0, 0, 0, 0]], dtype=int32)>
```

日期: /

image padding.

```

●●●
In [128]: a=tf.random.normal([4,28,28,3])
          ↗ 28x28
In [129]: b=tf.pad(a, [[0,0], [2,2], [2,2], [0,0]])
          ↗ batch
          ↓ (1) 列
In [130]: b.shape
Out[130]: TensorShape([4, 32, 32, 3])

```

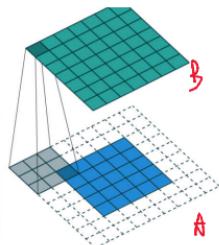
pad batch.

$B \rightarrow A$

$A \rightarrow A'$

$A' \rightarrow A$

使 $A' = B$



$[1, 5, 5, 1]$

$[0, 0], [2, 2], [2, 2], [0, 0]$

↓

$[1, 9, 9, 1]$

tile.

repeat data along dim n. times.

$[a, b, c], 2 \rightarrow [a, b, c, a, b, c]$

选择绕着小同的维度

```

In [132]: a
<tf.Tensor: id=396, shape=(3, 3), dtype=int32, numpy=
array([[0, 1, 2],
       [3, 4, 5],
       [6, 7, 8]], dtype=int32)
          分组
          ↗ 部分是原本的2倍.
In [133]: tf.tile(a,[1,2])
<tf.Tensor: id=399, shape=(3, 6), dtype=int32, numpy=
array([[0, 1, 2, 0, 1, 2],
       [3, 4, 5, 3, 4, 5],
       [6, 7, 8, 6, 7, 8]], dtype=int32)
          ↗ 分组
          ↗ 小麦.
In [134]: tf.tile(a,[2,1]) → 3|小麦.
<tf.Tensor: id=402, shape=(6, 3), dtype=int32, numpy=
array([[0, 1, 2],
       [3, 4, 5],
       [6, 7, 8],
       [0, 1, 2],
       [3, 4, 5],
       [6, 7, 8]], dtype=int32)

```

排列到原数据

W1[1,2]

都复制成2倍,

```

In [135]: tf.tile(a,[2,2])
<tf.Tensor: id=405, shape=(6, 6), dtype=int32, numpy=
array([[0, 1, 2, 0, 1, 2],
       [3, 4, 5, 3, 4, 5],
       [6, 7, 8, 6, 7, 8],
       [0, 1, 2, 0, 1, 2],
       [3, 4, 5, 3, 4, 5],
       [6, 7, 8, 6, 7, 8]], dtype=int32)

```

日期: /

张量操作.

▪ clip_by_value

根据值来裁剪

▪ relu

▪ clip_by_norm

▪ gradient clipping

```
●●●  
In [149]: a  
Out[149]: <tf.Tensor: id=422, shape=(10,), dtype=int32, numpy=array([0, 1, 2, 3,  
4, 5, 6, 7, 8, 9], dtype=int32)>  
  
In [150]: tf.maximum(a,2) max(x,y) if x<y, res=y. else x,  
Out[150]: <tf.Tensor: id=426, shape=(10,), dtype=int32, numpy=array([2, 2, 3,  
4, 5, 6, 7, 8, 9], dtype=int32)>  
  
In [151]: tf.minimum(a,8) min(x,y) if x>y, res=y. else x,  
Out[151]: <tf.Tensor: id=429, shape=(10,), dtype=int32, numpy=array([0, 1, 2, 3,  
4, 5, 6, 7, 8, 8], dtype=int32)>  
  
In [152]: tf.clip_by_value(a,2,8) → [2,8] 2~8. [min,max]  
Out[152]: <tf.Tensor: id=434, shape=(10,), dtype=int32, numpy=array([2, 2, 2, 3,  
4, 5, 6, 7, 8, 8], dtype=int32)>
```

relu.

```
●●●  
In [153]: a=a-5  
Out[154]: <tf.Tensor: id=437, shape=(10,), dtype=int32, numpy=array([-5, -4, -3,  
-2, -1, 0, 1, 2, 3, 4], dtype=int32)>  
  
In [155]: tf.nn.relu(a)✓  
Out[155]: <tf.Tensor: id=439, shape=(10,), dtype=int32, numpy=array([0, 0, 0, 0,  
0, 0, 1, 2, 3, 4], dtype=int32)>  
  
In [156]: tf.maximum(a,0)  
Out[156]: <tf.Tensor: id=442, shape=(10,), dtype=int32, numpy=array([0, 0, 0, 0,  
0, 0, 1, 2, 3, 4], dtype=int32)>
```

日期: / $\sqrt{2}x^2$ 等于放缩.

clip_by_norm 范数. 改变 gradient 方向.

```
In [157]: a=tf.random.normal([2,2],mean=10)
<tf.Tensor: id=449, shape=(2, 2), dtype=float32, numpy=
array([[12.217459 , 10.1498375],
       [10.84643 , 10.972536 ]], dtype=float32)>
 $\sqrt{2}x^2$ .  
)即 - 到 1. 再乘以 x.  
↓  
和大于范数.
```

```
In [159]: tf.norm(a)
Out[159]: <tf.Tensor: id=455, shape=(), dtype=float32, numpy=22.14333>
```

```
In [161]: aa=tf.clip_by_norm(a,15)  $\downarrow$  new norm.
<tf.Tensor: id=473, shape=(2, 2), dtype=float32, numpy=
array([[8.276167 , 6.8755493],
       [7.3474245, 7.43285 ]], dtype=float32)>
```

```
In [162]: tf.norm(aa)
Out[162]: <tf.Tensor: id=496, shape=(), dtype=float32, numpy=15.000001>
```



Gradient clipping

- Gradient Exploding or vanishing

过大 过小

- set $\|r\|_1$

learning rate.

对所有的进行缩放. \rightarrow 整体 norm.

- new_grads, total_norm = tf.clip_by_global_norm(grads, 25)
 \rightarrow 缩放后 grads. 没有 clipping in norm.

$$\sum \|W_i\| + \|W_j\| \dots$$

一个网络包含了很多维度方向 $[w_1, w_2, w_3 \dots]$

每个方向都存在梯度 $gradient - w_1$
 $gradient - w_r$

保持整体缩放的 方向不变.

```
(x, y), _ = datasets.mnist.load_data()
x = tf.convert_to_tensor(x, dtype=tf.float32) / 50.

i@z68:~/TutorialsCN/code_TensorFlow2.0/lesson18-数据限制$ python main.py
2.0.0-dev20190225
x: (60000, 28, 28) y: (60000, 10)
sample: (128, 28, 28) (128, 10) norm.建议ovo.
```

Step1

gw
y-w

```
==before==
tf.Tensor(89.03711, shape=(), dtype=float32)
tf.Tensor(2.6175494, shape=(), dtype=float32)
tf.Tensor(118.17449, shape=(), dtype=float32)
tf.Tensor(2.1617627, shape=(), dtype=float32)
tf.Tensor(134.27968, shape=(), dtype=float32)
tf.Tensor(2.5254946, shape=(), dtype=float32)
0 loss: 28.8848876953125
```

Step2

过大 gradient exploring.

```
==before==
tf.Tensor(1143.292, shape=(), dtype=float32)
tf.Tensor(35.847225, shape=(), dtype=float32)
tf.Tensor(1279.236, shape=(), dtype=float32)
tf.Tensor(24.312374, shape=(), dtype=float32)
tf.Tensor(1185.6311, shape=(), dtype=float32)
tf.Tensor(17.80448, shape=(), dtype=float32)
```

裁剪

```
...
print('==before==')
for g in grads:
    print(tf.norm(g))
grads, _ = tf.clip_by_global_norm(grads, 15)

print('==after==')
for g in grads:
    print(tf.norm(g))
```

```
i@z68:~/TutorialsCN/code_TensorFlow2.0/lesson18-数据限制$ python main.py
2.0.0-dev20190225
x: (60000, 28, 28) y: (60000, 10)
sample: (128, 28, 28) (128, 10)
==before==
tf.Tensor(118.00854, shape=(), dtype=float32)
tf.Tensor(3.5821552, shape=(), dtype=float32)
tf.Tensor(146.76697, shape=(), dtype=float32)
tf.Tensor(2.830059, shape=(), dtype=float32)
tf.Tensor(183.28879, shape=(), dtype=float32)
tf.Tensor(3.4088597, shape=(), dtype=float32)
==after==
tf.Tensor(6.734187, shape=(), dtype=float32)
tf.Tensor(0.20441659, shape=(), dtype=float32)
tf.Tensor(8.375294, shape=(), dtype=float32)
tf.Tensor(0.16149803, shape=(), dtype=float32)
tf.Tensor(10.45942, shape=(), dtype=float32)
tf.Tensor(0.19452743, shape=(), dtype=float32)
0 loss: 41.25679016113281
```

gradient 出现问题会导致 loss:nan 小稳定

日期: /

高阶操作.

where.

Scatter

meshgrid

布尔型.

[3,3]

Where.tensor:

True	False	False
False	True	False
False	False	True

返回 true 的坐标.

[3,2] → [0,0]
shape.
[1,1]
[2,2]]

```
In [3]: a=tf.random.normal([3,3])
<tf.Tensor: id=11, shape=(3, 3), dtype=float32, numpy=
array([[ 1.6420907 ,  0.43938753, -0.31872085],
       [ 1.144599 , -0.02425919, -0.9576591 ],
       [ 1.5931814 ,  0.1182256 , -0.39948994]], dtype=float32)>

In [5]: mask=a>0  这是布尔值.
<tf.Tensor: id=14, shape=(3, 3), dtype=bool, numpy=
array([[ True,  True,  False],
       [ True,  False,  False],
       [ True,  True,  False]])>

In [7]: tf.boolean_mask(a,mask) 得到为true的具体值.
<tf.Tensor: id=42, shape=(5,), dtype=float32, numpy=
array([1.6420907 , 0.43938753, 1.144599 , 1.5931814 , 0.1182256 ],
      dtype=float32)> (具体值)

In [8]: indices=tf.where(mask)
<tf.Tensor: id=44, shape=(5, 2), dtype=int64, numpy=
array([[0, 0],
       [0, 1],
       [1, 0],
       [2, 0],
       [2, 1]])>

In [10]: tf.gather_nd(a,indices) 将原值取出来. 之后起来实现 boolean.
```

```
In [11]: mask
<tf.Tensor: id=14, shape=(3, 3), dtype=bool, numpy=
array([[ True,  True, False],
       [ True, False, False],
       [ True,  True, False]])>

In [12]: A=tf.ones([3,3]) [1 1 1  
1 1 1  
1 1 1]

In [13]: B=tf.zeros([3,3])

In [14]: tf.where(mask, A, B)
<tf.Tensor: id=55, shape=(3, 3), dtype=float32, numpy=
array([[1., 1., 0.],
       [1., 0., 0.],
       [1., 1., 0.]], dtype=float32)>
```

where (cond, A, B).
↓
(mask) T M. AB 追.
F MAB E 追.

Scatter-nd

- `tf.scatter_nd`
 - `indices`,
 - `updates`,
 - `shape`

buffer indices. 将 update 位置存入
依次更新到 shape 中。返回 output。
底档是全局。

更新现有 tensor = A

$$A - A' \rightarrow \text{clear}$$

updates

$A^{(1)} \rightarrow \text{scatter_nd}$ (保留小度的数据)
将两个 output 相累加

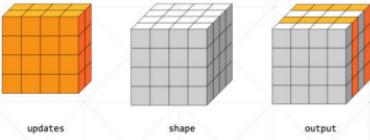
将两个 output 相累加

```
In [17]: indices = tf.constant([[4], [3], [1], [7]])
In [18]: updates = tf.constant([9, 10, 11, 12])
In [19]: shape = tf.constant([8])

In [20]: tf.scatter_nd(indices, updates, shape)
Out[20]: <tf.Tensor: id=60, shape=(8,), dtype=int32, numpy=array([ 0, 11,  0, 10,
   9,  0,  0, 12], dtype=int32)>
```

日期:

二/三/



updates

shape

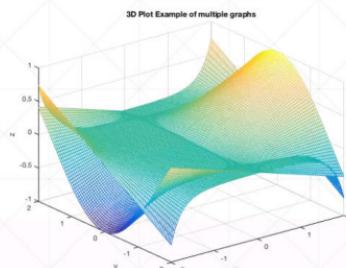
output

```
•••  
In [21]: indices = tf.constant([[0], [2]]) 低層、  
In [22]: updates = tf.constant([[[5, 5, 5, 5], [6, 6, 6, 6],  
...:            [7, 7, 7, 7], [8, 8, 8, 8]],  
...:  
...:  
...:  
In [24]: updates.shape  
Out[24]: TensorShape([2, 4, 4])  
In [23]: shape = tf.constant([4, 4, 4])  
In [25]: tf.scatter_nd(indices, updates, shape)  
<tf.Tensor: id=65, shape=(4, 4, 4), dtype=int32, numpy=  
array([[[5, 5, 5, 5],  
       [6, 6, 6, 6],  
       [7, 7, 7, 7],  
       [8, 8, 8, 8]],  
      [[0, 0, 0, 0],  
       [0, 0, 0, 0],  
       [0, 0, 0, 0],  
       [0, 0, 0, 0]],  
      [[5, 5, 5, 5],  
       [6, 6, 6, 6],  
       [7, 7, 7, 7],  
       [8, 8, 8, 8]],  
      [[0, 0, 0, 0],  
       [0, 0, 0, 0],  
       [0, 0, 0, 0],  
       [0, 0, 0, 0]]], dtype=int32)>
```

→ 及新
→ 底板.

meshgrid.

- [-2, -2]
- [-1, -2]
- [0, -2]
- [-2, -1]
- [-1, -1]
- ...
- [2, 2]



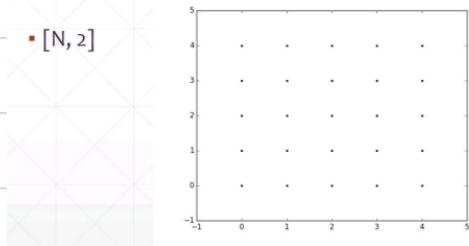
日期:

1

- $[y, x, z]$
 - $[5, 5, 2]$

Points.

这样
那样



```
points = []
for y in np.linspace(-2,2,5):
    for x in np.linspace(-2,2,5):
        points.append([x,y])
return np.array(points)
```

$$\begin{array}{l} x[-2, 2] \\ y[-2, 2] \end{array}$$

如何通过 GPU acceleration 实现

points
[n, 2]

```
In [39]: y=tf.linspace(-2.,2,5)
In [40]: y
Out[40]: <tf.Tensor: id=136, shape=(5,), dtype=float32, numpy=array([-2., -1., 0., 1., 2.], dtype=float32)>

In [41]: x=tf.linspace(-2.,2,5)

In [42]: points_x, points_y=tf.meshgrid(x,y)

In [43]: points_x.shape
Out[43]: TensorShape([5, 5])
```

•••

堵在一起 [5, 5, 2] 坐标部分。

In [44]: points_x

<tf.Tensor: id=162, shape=(5, 5), dtype=float32, numpy=
array([[-2., -1., 0., 1., 2.],
 [-2., -1., 0., 1., 2.],
 [-2., -1., 0., 1., 2.],
 [-2., -1., 0., 1., 2.],
 [-2., -1., 0., 1., 2.]], dtype=float32)>

*坐标点。

In [45]: points_y

<tf.Tensor: id=163, shape=(5, 5), dtype=float32, numpy=
array([[-2., -2., -2., -2., -2.],
 [-1., -1., -1., -1., -1.],
 [0., 0., 0., 0., 0.],
 [1., 1., 1., 1., 1.],
 [2., 2., 2., 2., 2.]], dtype=float32)>

坐标点。

```
[51]: points_x.shape  
Out[51]: TensorShape([5, 5])  
  
In [52]: points=tf.stack([points_x,points_y], axis=2)  
<tf.Tensor: id=176, shape=(5, 5, 2), dtype=float32, numpy=  
array([[[[-2., -2.],  
       [-1., -2.],  
       [ 0., -2.],  
       [ 1., -2.],  
       [ 2., -2.]],  
  
      [[-2., -1.],  
       [-1., -1.],  
       [ 0., -1.],  
       [ 1., -1.],  
       [ 2., -1.]]])
```

日期: /

数据集加载

小型数据集.

keras.datasets

- keras.datasets
- tf.data.Dataset.from_tensor_slices
 - shuffle
 - map
 - batch
 - repeat
- we will talk **Input Pipeline** later

变为 tensor

boston housing

- Boston housing price regression dataset.

mnist/fashion mnist

- MNIST/Fashion-MNIST dataset.

手写字

cifar10/100

- small images classification dataset.

imdb

- sentiment classification dataset.

Mnist

bulk train test

numpy

```
In [5]: (x,y), (x_test,y_test)=keras.datasets.mnist.load_data()
In [6]: x.shape # (60000, 28, 28)
In [7]: y.shape # (60000,) 0~7 → one-hot (label)
In [8]: x.min(),x.max(),x.mean()
Out[8]: (0, 255, 33.318421449829934)           $\frac{0}{255} \sim 1/255$ 
In [9]: x_test.shape, y_test.shape
Out[9]: ((10000, 28, 28), (10000,))

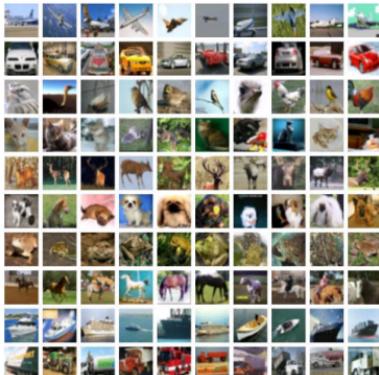
In [10]: y[:4] label
Out[10]: array([5, 0, 4, 1], dtype=uint8)

In [11]: y_onehot=tf.one_hot(y, depth=10)    独热向量 one-hot
In [12]: y_onehot[:2]                         10类
<tf.Tensor: id=8, shape=(2, 10), dtype=float32, numpy=
array([[0., 0., 0., 0., 0., 1., 0., 0., 0., 0.], $\rightarrow$  1位是 1.
[1., 0., 0., 0., 0., 0., 0., 0., 0., 0.]], dtype=float32)>
↓ 0位是 1.
```

日期:

CIFAR10 / 100.

airplane



大类 → 10小类

automobile



size = [32, 32, 3]

bird 大类



cat



deer



dog



frog



horse



ship



truck



```
... ● ● ●
In [14]: (x,y),(x_test,y_test)=keras.datasets.cifar10.load_data()
          (0K, 0K, cifar10.

In [15]: x.shape,y.shape,x_test.shape,y_test.shape
Out[15]: ((50000, 32, 32, 3), (50000, 1), (10000, 32, 32, 3), (10000, 1))

In [16]: x.min(),x.max()
Out[16]: (0, 255)

In [17]: y[:4]
Out[17]:
array([6,
       9,
       9,
       4], dtype=uint8)
```

tf.data.Dataset

from_tensor_slices()

x: → x-tensor

```
... ● ● ●
In [5]: (x,y),(x_test,y_test)=keras.datasets.cifar10.load_data()

In [6]: db=tf.data.Dataset.from_tensor_slices(x_test)
        it=iter(db).

In [7]: next(it).shape
Out[7]: TensorShape([32, 32, 3])           it=iter(db).
                                                next(it)  迭代器

# can not use [x_test,y_test]
In [9]: db=tf.data.Dataset.from_tensor_slices((x_test,y_test))
        image + label

In [11]: next(it)[0].shape
Out[11]: TensorShape([32, 32, 3])           image + label
```

label 顺序.

日期: / shuffle 将图片顺序打乱
(Dataset)

```
In [12]: db=tf.data.Dataset.from_tensor_slices((x_test,y_test))  
In [13]: db=db.shuffle(10000)
```

正确的

数据预处理.

$x \rightarrow \text{Tensor} \rightarrow \text{float32} \rightarrow \text{one-hot}$

Map.

```
In [16]: def preprocess(x,v): [32,32,3]  
...: x=tf.cast(x, dtype=tf.float32)/255. # -> (V.  
...: y=tf.cast(y, dtype=tf.int32)  
...: y=tf.one_hot(y, depth=10)  
...: return x,y  
  
In [17]: db2=db.map(preprocess)  
  
In [18]: res=next(iter(db2))  
In [19]: res[0].shape, res[1].shape  
Out[19]: (TensorShape([32, 32, 3]), TensorShape([1, 10]))  
In [20]: res[1][2]  
Out[20]: <tf.Tensor: id=58, shape=(1, 10), dtype=float32, numpy=array([[1., 0.,  
0., 0., 0., 0., 0., 0., 0., 0.]])>, dtype=float32>
```

.batch.

batch 处理

```
In [21]: db3=db2.batch(32) [0, 32, 32, 3] 批处理 [60, 1] 1  
In [25]: res=next(iter(db3))  
  
In [26]: res[0].shape, res[1].shape  
Out[26]: (TensorShape([32, 32, 32, 3]), TensorShape([32, 1, 10]))  
x
```

StopIteration.

```
In [28]: db_iter = iter(db3)  
  
In [28]: while True:  
...:     next(db_iter)  
...:  
  
OutOfRangeException Traceback (most recent call last)  
StopIteration:
```

for x,y in db4

.repeat().

停止.

```
In [29]: db4=db3.repeat()
```

```
In [30]: db4=db3.repeat(2) ② 迭代 次数.
```

日期: /

For example.

```
def prepare_mnist_features_and_labels(x, y):
    x = tf.cast(x, tf.float32) / 255.0
    y = tf.cast(y, tf.int64)
    return x, y

def mnist_dataset(): ↗ok.
    (x, y), (x_val, y_val) = datasets.fashion_mnist.load_data()
    y = tf.one_hot(y, depth=10) ↗ok, 10
    y_val = tf.one_hot(y_val, depth=10)

    ds = tf.data.Dataset.from_tensor_slices((x, y)) →dataset
    ds = ds.map(prepare_mnist_features_and_labels) →预处理
    ds = ds.shuffle(60000).batch(100) →打散+batch.
    ds_val = tf.data.Dataset.from_tensor_slices((x_val, y_val))
    ds_val = ds_val.map(prepare_mnist_features_and_labels)
    ds_val = ds_val.shuffle(10000).batch(100)
    return ds,ds_val
```

日期: /

测试张量，实践

- Matmul
- Neural Network
- Deep Learning
- Multi-Layer

全连接层

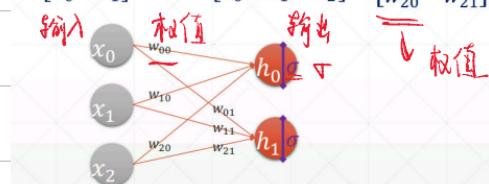
↑非线性。

$$\text{out} = f(X @ W + b)$$

$$\text{out} = \text{relu}(X @ W + b)$$

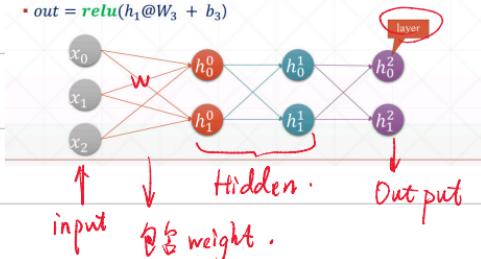
低维 \rightarrow 高维。

$$[h_0, h_1] = \begin{bmatrix} h_0^0 & h_1^0 \\ h_0^1 & h_1^1 \end{bmatrix} = \text{relu}\left(\begin{bmatrix} x_0^0 & x_1^0 & x_2^0 \\ x_0^1 & x_1^1 & x_2^1 \end{bmatrix} @ \begin{bmatrix} w_{00} & w_{01} \\ w_{10} & w_{11} \\ w_{20} & w_{21} \end{bmatrix} + [b_0, b_1]\right)$$



$X @ W + b$ 是神经网络的层。

- $h_0 = \text{relu}(X @ W_1 + b_1)$
- $h_1 = \text{relu}(h_0 @ W_2 + b_2)$
- $\text{out} = \text{relu}(h_1 @ W_3 + b_3)$

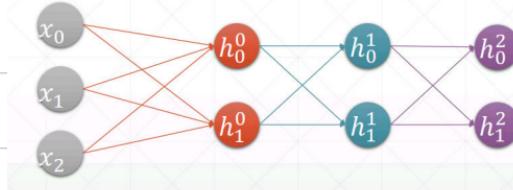


权值通过 deep learning 得到，loss 减少。→ 损失下降。

日期:

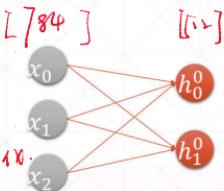
▪ Neural Network in the 1980s

- 3~5 layers



deeplearning now: $n \approx 1000$ layers.

- BigDATA
- ReLU *relu*.
- Dropout
- BatchNorm
- ResNet
- Xavier Initialization ~~和Dense~~
- Caffe/TensorFlow/PyTorch



tf.keras.layers.Dense.

```
In [49]: x=tf.random.normal([4,784])  
In [48]: net=tf.keras.layers.Dense(512)  
In [50]: out=net(x)  
In [51]: out.shape  
Out[51]: TensorShape([4, 512])  
In [52]: net.kernel.shape, net.bias.shape  
Out[52]: (TensorShape([784, 512]), TensorShape([512]))
```

① *自动生成 W*

② *使用 build 生成 W*

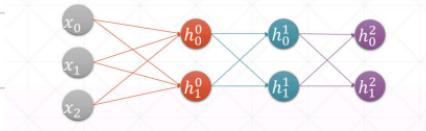
```
In [3]: net=tf.keras.layers.Dense(10)  
In [4]: net.bias  
# AttributeError: 'Dense' object has no attribute 'bias'  
In [5]: net.get_weights()  
Out[5]: []  
In [6]: net.weights  
Out[6]: []  
  
In [13]: net.build(input_shape=(None,4))  
In [14]: net.kernel.shape,net.bias.shape  
Out[14]: (TensorShape([4, 10]), TensorShape([10]))  
  
In [15]: net.build(input_shape=(None,20))  
In [16]: net.kernel.shape,net.bias.shape  
Out[16]: (TensorShape([20, 10]), TensorShape([10]))  
  
In [10]: net.build(input_shape=(2,4))  
In [11]: net.kernel  
<tf.Variable 'kernel:0' shape=(4, 10) dtype=float32, numpy=  
array([-0.28106192, -0.2522246 ,  0.16050524,  0.43587887, -0.50773597,
```

手动生成

```
In [15]: net.build(input_shape=(None,20))  
In [16]: net.kernel.shape,net.bias.shape  
Out[16]: (TensorShape([20, 10]), TensorShape([10]))  
  
In [17]: out=net(tf.random.rand((4,12)))  
InvalidArgumentError: Matrix size-incompatible: In[0]: [4,12], In[1]: [20,10]  
[Op:MatMul]  
In [19]: out=net(tf.random.normal((4,20)))  
  
In [20]: out.shape  
Out[20]: TensorShape([4, 10])
```

日期: / 設定好 Dense.

- keras.Sequential([layer1, layer2, layer3])



訓練网络结构

```
x = tf.random.normal([2, 3])  
model = keras.Sequential([  
    keras.layers.Dense(2, activation='relu'),  
    keras.layers.Dense(2, activation='relu'),  
    keras.layers.Dense(2)  
])  
model.build(input_shape=[None, 3])  
model.summary()  
for p in model.trainable_variables:  
    print(p.name, p.shape)
```

输出方式

- $y \in R^d$

实数

- $y_i \in [0, 1], i = 0, 1, \dots, y_d - 1$

概率。

- $y_i \in [0, 1], \sum_{i=0}^{y_d} y_i = 1, i = 0, 1, \dots, y_d - 1$

概率和 = 1

- $y_i \in [-1, 1], i = 0, 1, \dots, y_d - 1$

1. $y \in R^d$ 实数.

- linear regression

线性回归.

- naïve classification with MSE 损差.

- other general prediction

- $out = \text{relu}(X@W + b)$ (正数用 relu)

• logits

日期: /

2) $y_i \in [0, 1]$

- binary classification

- $y > 0.5, \rightarrow 1$
- $y < 0.5, \rightarrow 0$

二分类

- Image Generation

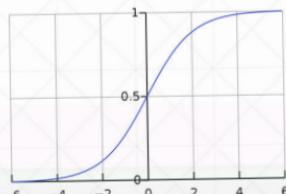
- $\text{rgb} \quad 0 \sim 255 \rightarrow [0, 1] \quad 128 \cdot 128$.

如何实现将数据压缩

- $\text{out} = \text{relu}(X @ W + b)$

▪ sigmoid

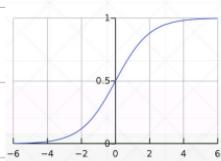
\downarrow
ReLU



- $\text{out}' = \text{sigmoid}(\text{out})$

Sigmoid

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



tf. sigmoid

```
In [10]: a=tf.linspace(-6,6,(10))  
<tf.Tensor: id=19, shape=(10,), dtype=float32, numpy=  
array([-6., -4.6666665, -3.333333, -2.  
0.6666667, 2., 3.333334, 4.666667, 6.  
])  
  
In [12]: tf.sigmoid(a)  
<tf.Tensor: id=21, shape=(10,), dtype=float32, numpy=  
array([0.00247264, 0.00931597, 0.03444517, 0.11920291, 0.33924365,  
0.6607564, 0.8807971, 0.96555483, 0.99068403, 0.9975274 ])  
  
In [13]: x=tf.random.normal([1,28,28])*5  
In [15]: tf.reduce_min(x), tf.reduce_max(x)  
(<tf.Tensor: id=32, shape=(), dtype=float32, numpy=-18.78872>,  
<tf.Tensor: id=34, shape=(), dtype=float32, numpy=15.466431>)  
  
In [16]: x=tf.sigmoid(x)  
In [17]: tf.reduce_min(x), tf.reduce_max(x)  
(<tf.Tensor: id=39, shape=(), dtype=float32, numpy=0.0>,  
<tf.Tensor: id=41, shape=(), dtype=float32, numpy=0.99999976>)
```

日期:

2019.1.

3. $y_i \in [0, 1] \quad \sum y_i = 1$

▪ sigmoid

```
In [21]: a=tf.linspace(-2., 2, 5)
In [22]: tf.sigmoid(a)
<tf.Tensor: id=54, shape=(5,), dtype=float32, numpy=
array([0.11920291, 0.26894143, 0.5       , 0.7310586 , 0.880797 ],
```

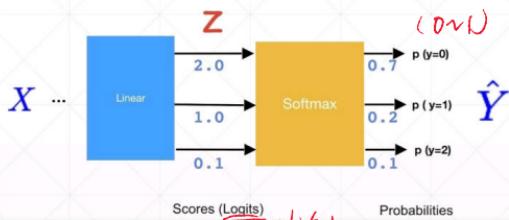
▪ softmax 概率和 = 1

```
In [23]: tf.nn.softmax(a)
<tf.Tensor: id=56, shape=(5,), dtype=float32, numpy=
array([0.01165623, 0.03168492, 0.08612854, 0.23412167, 0.6364086 ],
```

分类问题

Meet Softmax

$$\sigma(\mathbf{z})_j = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}} \quad \text{for } j = 1, \dots, K.$$



softmax
概率函数的输出

-1 ~ 1

```
In [24]: logits=tf.random.uniform([1,10],minval=-2,maxval=2)
<tf.Tensor: id=64, shape=(1, 10), dtype=float32, numpy=
array([-1.5842109 ,  1.6475668 , -0.71567106,  1.5819931 ,  0.35972595,
       0.12556812,  1.0662012 , -0.70131207, -1.5194197 ,  1.5553613 ],
```

dtype=float32>

[...]

```
In [27]: prob=tf.nn.softmax(logits, axis=1)
<tf.Tensor: id=67, shape=(1, 10), dtype=float32, numpy=
array([[0.00946281, 0.23964217, 0.02255392, 0.22443208, 0.06610908,
       0.05229748, 0.13399215, 0.02288011, 0.01009621, 0.21853393]],
```

dtype=float32>

求和 = 1

```
In [29]: tf.reduce_sum(prob, axis=1)
Out[29]: <tf.Tensor: id=70, shape=(1,), dtype=float32, numpy=array([0.99999994],
```

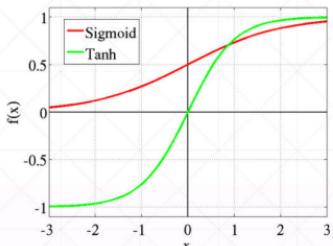
dtype=float32>

日期: /

4. $y_i \in [-1, 1]$

▪ tanh

$\sinh = 1$



$$\tanh(x) = \sinh(x)/\cosh(x) = (e^x - e^{-x})/(e^x + e^{-x})$$

Sigmoid 放大 2倍
↓ 下移

tanh.

```
●●●  
In [30]: a  
Out[30]: <tf.Tensor: id=53, shape=(5,), dtype=float32, numpy=array([-2., -1., 0., 1., 2.], dtype=float32)>  
  
In [33]: tf.tanh(a)  
<tf.Tensor: id=73, shape=(5,), dtype=float32, numpy=  
array([-0.9640276, -0.7615942, 0. , 0.7615942, 0.9640276],  
      dtype=float32)>
```

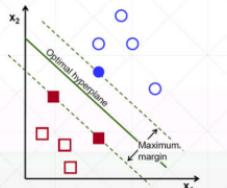
日期: /

损失函数 MSE 均方差.

MSE
▪ Cross Entropy Loss

Hinge Loss 支持向量机

$$\sum_i \max(0, 1 - y_i * h_\theta(x_i))$$



MSE

$$\text{loss} = \frac{1}{N} \sum (y - \text{out})^2 \quad N = \text{batch}.$$

二范数.

$$\underline{L_2\text{-norm}} = \sqrt{\sum (y - \text{out})^2}$$

```
y = tf.constant([1, 2, 3, 0, 2])
y = tf.one_hot(y, depth=4) 0, 1, 2, 3.
y = tf.cast(y, dtype=tf.float32)

out = tf.random.normal([5, 4]) logits.
    > MSE. 均值.
loss1 = tf.reduce_mean(tf.square(y-out))
    > L2
loss2 = tf.square(tf.norm(y-out))/(5*4)
    (b).
loss3 = tf.reduce_mean(tf.losses.MSE(y, out)) # VS MeanSquaredError is a class

tf.Tensor(1.0918376, shape=(), dtype=float32)
tf.Tensor(1.0918376, shape=(), dtype=float32)
tf.Tensor(1.0918376, shape=(), dtype=float32)
```

日期: /

Entropy. 熵.

- Uncertainty
- measure of surprise

低熵.

- lower entropy → more info. 不确定

离散.

$$Entropy = - \sum_i P(i) \log P(i)$$

分布.

信息量大



Claude Shannon

Lottery.

概率均等 → 高熵. (精细度小)

```

In [3]: a=tf.fill([4],0.25) 4类, [0.25, 0.25, 0.25, 0.25]
In [6]: a*tf.math.log(a)/tf.math.log(2.) a*log2 a.
Out[6]: <tf.Tensor: id=13, shape=(4,), dtype=float32, numpy=array([-0.5, -0.5,
-0.5, -0.5], dtype=float32)
In [7]: -tf.reduce_sum(a*tf.math.log(a)/tf.math.log(2.)) 手动.
Out[7]: <tf.Tensor: id=22, shape=(), dtype=float32, numpy=-2.0> entropy. 高熵

In [8]: a=tf.constant([0.1,0.1,0.1,0.7]) 概率不等.
In [9]: -tf.reduce_sum(a*tf.math.log(a)/tf.math.log(2.))
Out[9]: <tf.Tensor: id=32, shape=(), dtype=float32, numpy=1.3567796> 精细度变大

In [10]: a=tf.constant([0.01,0.01,0.01,0.97])
In [11]: -tf.reduce_sum(a*tf.math.log(a)/tf.math.log(2.))
Out[11]: <tf.Tensor: id=42, shape=(), dtype=float32, numpy=0.24194068> 精细度高

```

Cross Entropy.

$$H(p, q) = - \sum p(x) \log q(x)$$

两个分布

$$H(p, q) = H(p) + D_{KL}(p||q).$$

相似度 / 差异.

交叉熵

分类问题.

当 $p=q$ 时, $H(p, q)=0$

- for $p = q$

▪ Minima: $H(p, q) = H(p)$

p 预测
q: prob.

- for p : one-hot encoding

▪ $H([0,1,0]) = -1\log 1 = 0$

▪ $H([0,1,0], [p_0, p_1, p_2]) = 0 + D_{KL}(p||q) = -\log q_1$

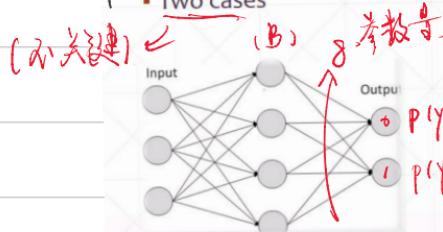
P.

D: one-hot.
label.

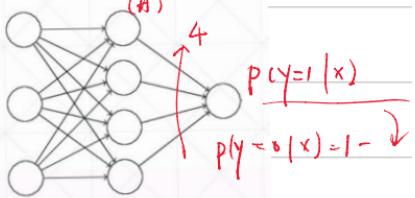
日期: /

Binary Classification

- Two cases



只有一个输出
(A)



Single Output

$$H(P, Q) = -P(cat) \log Q(cat) - (1 - P(cat)) \log(1 - Q(cat))$$

\therefore 2类 $\therefore P(dog) = (1 - P(cat))$

$$\text{loss} = H(p, q) \quad \text{One-hot}$$

$$H(P, Q) = - \sum_{i=(cat, dog)} P(i) \log Q(i)$$

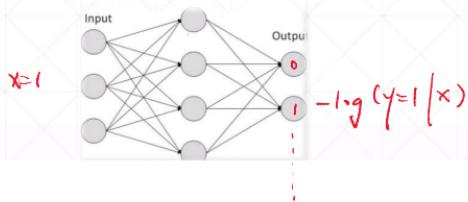
$$= -P(cat) \log Q(cat) - P(dog) \log Q(dog)$$

$$\text{loss} = -(y \log(p) + (1 - y) \log(1 - p))$$

$p(cat)$ 真实分布 p 概率值.

使节点概率
最大化。

$$\bullet H([0,1,0], [p_0, p_1, p_2]) = 0 + D_{KL}(p|q) = -\underline{\underline{1 \log q_1}}$$



Classification



label

$$P_1 = [1 \ 0 \ 0 \ 0 \ 0]$$

预测值、置信度

$$Q_1 = [0.4 \ 0.3 \ 0.05 \ 0.05 \ 0.2]$$

$$H(P_1, Q_1) = - \sum_i P_1(i) \log Q_1(i)$$

loss.

$$= -(1 \log 0.4 + 0 \log 0.3 + 0 \log 0.05 + 0 \log 0.05 + 0 \log 0.2)$$

$$= -\log 0.4$$

$$\approx 0.916$$

(label) log (prob.)

优化

$$Q_1 = [0.98 \ \underline{0.01} \ 0 \ 0 \ 0.01]$$

$$H(P_1, Q_1) = - \sum_i P_1(i) \log Q_1(i)$$

$$= -(1 \log 0.98 + 0 \log 0.01 + 0 \log 0 + 0 \log 0 + 0 \log 0.01)$$

$$= -\log 0.98$$

$$\approx 0.02$$

Categorical Cross Entropy. 什么情况。

```
••• (y_true) (prob)
In [15]: tf.losses.categorical_crossentropy([0,1,0,0],[0.25,0.25,0.25,0.25])
Out[15]: <tf.Tensor: id=98, shape=(), dtype=float32, numpy=1.3862944>

In [16]: tf.losses.categorical_crossentropy([0,1,0,0],[0.1,0.1,0.8,0.1])
Out[16]: <tf.Tensor: id=117, shape=(), dtype=float32, numpy=2.3978953>

In [17]: tf.losses.categorical_crossentropy([0,1,0,0],[0.1,0.7,0.1,0.1])
Out[17]: <tf.Tensor: id=136, shape=(), dtype=float32, numpy=0.35667497>

In [18]: tf.losses.categorical_crossentropy([0,1,0,0],[0.01,0.97,0.01,0.01])
Out[18]: <tf.Tensor: id=155, shape=(), dtype=float32, numpy=0.030459179>
```

••• ↓ = crossEntropy.

```
In [20]: criteron([0,1,0,0],[0.1,0.7,0.1,0.1])
Out[20]: <tf.Tensor: id=186, shape=(), dtype=float32, numpy=0.35667497>

In [21]: criteron([0,1],[0.9,0.1])
Out[21]: <tf.Tensor: id=216, shape=(), dtype=float32, numpy=2.3025851>
          ↗ 换入二分类 ↘ 换出二分类
In [22]: tf.losses.BinaryCrossentropy()([],[0,1])
Out[22]: <tf.Tensor: id=254, shape=(), dtype=float32, numpy=2.3025842>

In [23]: tf.losses.binary_crossentropy([1],[0.1])
Out[23]: <tf.Tensor: id=281, shape=(), dtype=float32, numpy=2.3025842>
```

日期:

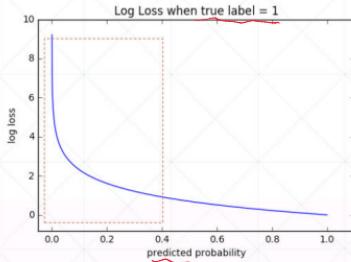
▪ sigmoid + MSE

▪ gradient vanish

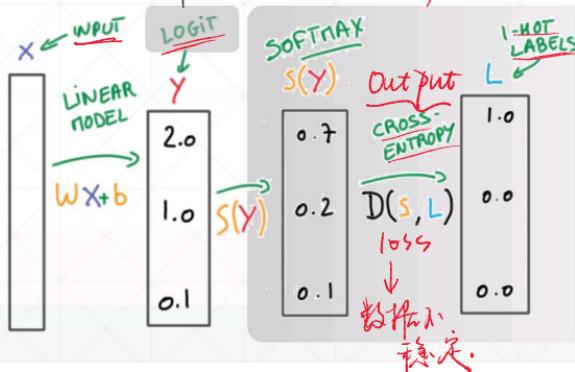
▪ converge slower

▪ However

▪ e.g. meta-learning



logits \rightarrow Cross Entropy \rightarrow loss. 合在一起。



```
...  
In [24]: x=tf.random.normal([1,784])  
In [25]: w=tf.random.normal([784,2])  
In [26]: b=tf.zeros([2])  
  
In [27]: logits=x@w+b  
Out[29]: <tf.Tensor: id=299, shape=(1, 2), dtype=float32, numpy=array([-26.27812, 28.63038]), dtype=float32>  
  
In [30]: prob=tf.math.softmax(logits, axis=1)  
Out[31]: <tf.Tensor: id=301, shape=(1, 2), dtype=float32,  
numpy=array([[1.4241021e-24, 1.0000000e+00]], dtype=float32>  
  
In [34]: tf.losses.categorical_crossentropy([0,1],logits, from_logits=True)  
Out[34]: <tf.Tensor: id=393, shape=(1,), dtype=float32, numpy=array([0.],  
dtype=float32)>  
  
In [35]: tf.losses.categorical_crossentropy([0,1],prob)  
Out[35]: <tf.Tensor: id=411, shape=(1,), dtype=float32, numpy=array([1.192093e-07],  
dtype=float32)>
```

用 logit.

直接计算 loss

日期: /

梯度下降. Gradient Descent.

What's Gradient?

通用

- 导数, derivative $\lim_{x \rightarrow 0} \frac{xy}{x}$ 变动量. (-维) 方向



(标量) 具体

- 偏微分, partial derivative 沿着某一方向. (x / y ?)

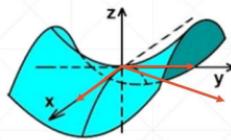
所有方向的偏微分.

- 梯度, gradient 合在一起

$$z = y^2 - x^2$$

$$\frac{\partial z}{\partial x} = -2x$$

$$\frac{\partial z}{\partial y} = 2y$$

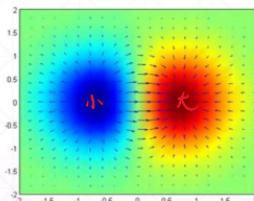
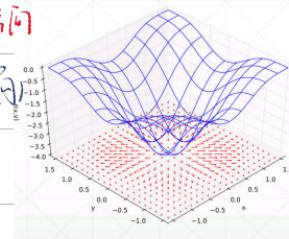


$\nabla f = \left(\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right)$ 形成一个向量. (有方向)

箭头: 梯度方向

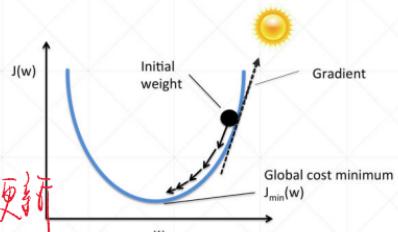
(增大方向)

坡度 (速率)



How to search?

- $\nabla f(\theta) \rightarrow \text{larger value}$



向梯度减小的方向更新

- Search for minima:

$$\text{lr} \propto \eta$$

$$\theta_{t+1} = \theta_t - \alpha_t \nabla f(\theta_t)$$

$$\downarrow \text{lr}$$

梯度.

For instance

$$\begin{array}{l} \text{U} \\ \theta_1 = 0 \\ \theta_2 = 0 \end{array}$$

$$\theta_{t+1} = \theta_t - \alpha_t \nabla f(\theta_t).$$

$$\theta_1 = 4 \quad \theta_2 = -4$$

$$\theta_1 = 4 - 0.01 \times 8$$

$$\theta_2 = -4 - 0.01 \times -8$$

Function:

$$J(\theta_1, \theta_2) = \theta_1^2 + \theta_2^2$$

Objective:

$$\min_{\theta_1, \theta_2} J(\theta_1, \theta_2)$$

Update rules:

$$\begin{aligned} \theta_1 &:= \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1, \theta_2) \\ \theta_2 &:= \theta_2 - \alpha \frac{d}{d\theta_2} J(\theta_1, \theta_2) \end{aligned}$$

Derivatives:

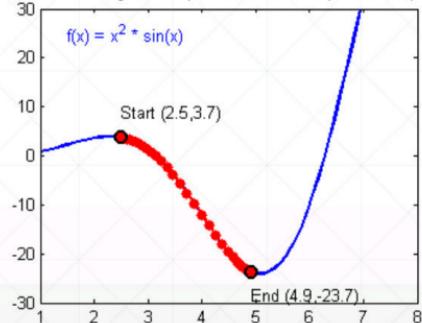
$$\frac{d}{d\theta_1} J(\theta_1, \theta_2) = \frac{d}{d\theta_1} \theta_1^2 + \frac{d}{d\theta_1} \theta_2^2 = 2\theta_1$$

$$\frac{d}{d\theta_2} J(\theta_1, \theta_2) = \frac{d}{d\theta_2} \theta_1^2 + \frac{d}{d\theta_2} \theta_2^2 = 2\theta_2$$

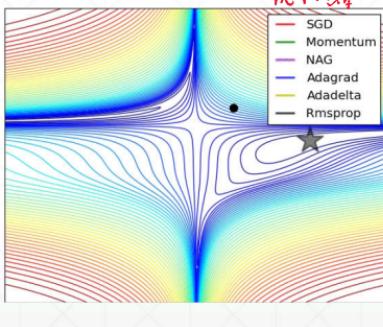
偏微分

单边

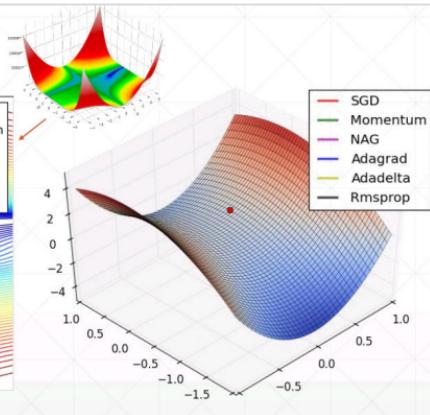
Descending with step coefficient 0.005 (iteration 50)



优化器



Learning Process-2



日期: /

AutoGrad

在一个环境中

用 tensorflow

With `Tf.GradientTape()` as tape:

- Build computation graph

$$\bullet \ loss = f_\theta(x)$$

对应的梯度.

$\bullet [w_grad] = tape.gradient(loss, [w])$

参数

```
...  
In [3]: w=tf.constant(1.) 定义参数.  
In [4]: x=tf.constant(2.)  
In [5]: y=x*w  
In [6]: loss  
In [8]: with tf.GradientTape() as tape: 创建环境.  
...:     tape.watch([w])  
...:     y2=x*w  
In [11]: grad1=tape.gradient(y,[w])  
Out[12]: [None] 自由 y 没有在环境中.  
In [18]: with tf.GradientTape() as tape:  
...:     tape.watch([w])  
...:     y2=x*w  
In [19]: grad2=tape.gradient(y2,[w])  
Out[20]: [tf.Tensor: id=8, shape=(), dtype=float32, numpy=2.0>]
```

耗时的资源.

```
...  
In [3]: w=tf.constant(1.)  
In [4]: x=tf.constant(2.)  
In [5]: y=x*w  
  
In [18]: with tf.GradientTape() as tape:  
...:     tape.watch([w])  
...:     y2=x*w  
In [19]: grad2=tape.gradient(y2,[w])  
Out[20]: [tf.Tensor: id=8, shape=(), dtype=float32, numpy=2.0>]  
  
In [19]: grad2=tape.gradient(y2,[w])  
RuntimeError: GradientTape.gradient can only be called once on non-persistent  
tapes.  
  
In [18]: with tf.GradientTape(persistent=True) as tape:  
...:     tape.watch([w]) 调用 2 次. 可以计算多个参数  
...:
```

可以计算多个参数

日期: /

二阶梯度

2nd-order

$$\bullet y = xw + b$$

$$\bullet \frac{\partial y}{\partial w} = x$$

x 与 w 无关

$$\bullet \frac{\partial^2 y}{\partial w^2} = \frac{\partial y'}{\partial w} = \frac{\partial x}{\partial w} = \underline{\text{None}}$$

```
with tf.GradientTape() as t1:  
    y = x * w + b  
    dy_dw, dy_db = t2.gradient(y, [w, b])  
    d2y_dw2 = t1.gradient(dy_dw, w)
```

原因
None. ↗ 梯度 1. ↘ 梯度 1.

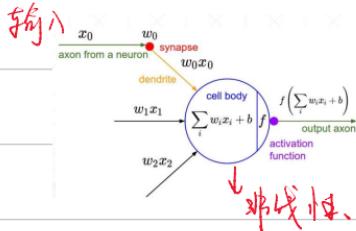
日期: /

激活函数及其梯度.

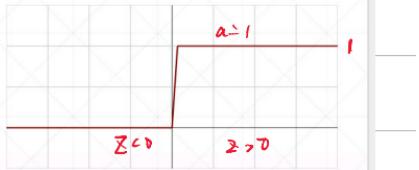
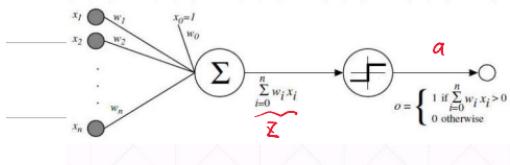
Sigmoid tanh. relu.



阈值响应机制



“激活” * 不可导

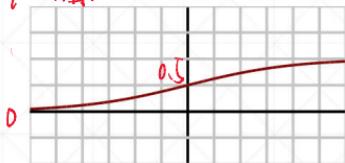


Sigmoid / logistic.

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

值域: [0, 1]

压缩 可导



prob, RGB...

缺省 $x \rightarrow \infty$ 时, \log
又更新

导数:

$$\begin{aligned} \frac{d}{dx} \sigma(x) &= \frac{d}{dx} \left(\frac{1}{1 + e^{-x}} \right) \\ &= \frac{e^{-x}}{(1 + e^{-x})^2} \\ &= \frac{(1 + e^{-x}) - 1}{(1 + e^{-x})^2} \\ &= \frac{1 + e^{-x}}{(1 + e^{-x})^2} - \left(\frac{1}{1 + e^{-x}} \right)^2 \\ &= \frac{\sigma(x) - \sigma(x)^2}{\sigma(x)^2} \\ \sigma' &= \sigma(1 - \sigma) \end{aligned}$$

$$\Gamma = \sigma(1 - \sigma)$$

```

a = tf.linspace(-10., 10., 10)
with tf.GradientTape() as tape:
    tape.watch(a)
    y = tf.sigmoid(a)

grads = tape.gradient(y, [a]) 梯度
x: [-10.          -7.77777777 -5.5555553  -3.333333  -1.1111107   1.1111116
      3.333334   5.5555563  7.7777786  10.          ]
y: [4.5388937e-05 4.1878223e-04 3.8510561e-03 3.4445226e-02 2.4766389e-01
    7.5233626e-01 9.6555448e-01 9.9614894e-01 9.9958128e-01 9.9995458e-01]
grad: [4.5386874e-05 4.186685e-04 3.8362255e-03 3.3258751e-02 1.8632649e-01
       1.8632641e-01 3.3258699e-02 3.8362255e-03 4.1854731e-04 4.5416677e-05]
  
```

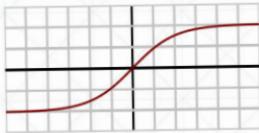
日期: /

Tanh 在 RNN 中使用

sigmoid x_2 变形 tanh.

$$f(x) = \tanh(x) = \frac{(e^x - e^{-x})}{(e^x + e^{-x})}$$

$$= 2 \text{sigmoid}(2x) - 1$$



$$\frac{d}{dx} \tanh(x) = \frac{(e^x + e^{-x})(e^x + e^{-x}) - (e^x - e^{-x})(e^x - e^{-x})}{(e^x + e^{-x})^2}$$

$$= 1 - \frac{(e^x - e^{-x})^2}{(e^x + e^{-x})^2} = 1 - \tanh^2(x)$$

```
In [5]: a=tf.linspace(-5.,5.,10)
```

```
In [6]: tf.tanh(a)
```

```
Out[6]:  
<tf.Tensor: id=10, shape=(10,), dtype=float32, numpy=  
array([-0.9999092 , -0.9991625 , -0.99229795, -0.9311096 , -0.50467217,  
       0.5046725 , 0.93110967, 0.99229795, 0.9991625 , 0.9999092 ],  
      dtype=float32)>
```

Rectified Linear Unit (Relu)

优势：梯度不变。

$$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$$

没有限制。



$$f'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$$

[-1, 1]

```
In [11]: a=tf.linspace(-1.,1.,10)
```

```
In [12]: tf.nn.relu(a)
```

```
<tf.Tensor: id=24, shape=(10,), dtype=float32, numpy=  
array([0.          , 0.          , 0.          , 0.          , 0.          ,  
       0.11111116, 0.33333337, 0.55555556, 0.77777778, 1.        ],  
      dtype=float32)>
```

$x < 0$ 时 $= k=0$. 通过函数平移。

```
In [13]: tf.nn.leaky_relu(a)
```

```
<tf.Tensor: id=26, shape=(10,), dtype=float32, numpy=  
array([-0.2          , -0.15555556, -0.11111112, -0.06666666, -0.02222222,  
       0.11111116, 0.33333337, 0.55555556, 0.77777778, 1.        ],  
      dtype=float32)>
```

$x > 0$ 时 $= k=1$.

日期: /

損失函數及其梯度.

Typical loss.

Mean Square Error

$$(y - \bar{y})^2$$

Cross Entropy Loss

$$-p \log q.$$

MSE

$$\text{loss} = \sum [y - f_\theta(x)]^2$$

$$\text{loss} = \sum \frac{(y - (x @ w + b))^2}{f(x; w, b)}.$$

$$\frac{\nabla \text{loss}}{\nabla \theta} = 2 \sum [y - f_\theta(x)] * \frac{\nabla f_\theta(x)}{\nabla \theta}$$

$$\text{L}_2\text{-norm} = \left| \left| y - (x @ w + b) \right| \right|_2$$

$$\sqrt{2(y_1 - y_2)^2}$$

$$\text{loss} = \text{norm}(y - (x @ w + b))$$

$x @ w + b \rightarrow \text{logit.}$

```
...  
In [3]: x=tf.random.normal([2,4])  
In [4]: w=tf.random.normal([4,3])  
In [5]: b=tf.zeros([3])  
In [6]: y=tf.constant([2,0]) [2]  
  
In [9]: with tf.GradientTape() as tape:  
...:     tape.watch([w,b]) if Variable.  
...:     prob = tf.nn.softmax(x@w+b, axis=1)  
...:     loss = tf.reduce_mean(tf.losses.MSE(tf.one_hot(y, depth=3), prob))
```

In [10]: grads = tape.gradient(loss, [w,b])

In [11]: grads[0]

[[-0.00967887, -0.00335512, 0.01303399],

[-0.04446869, 0.06194263, -0.01747394], [4,3] 5w-數

[-0.04530644, 0.01043231, 0.03487412],

[0.02006017, -0.03638988, 0.0163297]]

In [12]: grads[1] # [-0.02585024, 0.06217915, -0.03632889]

[3] 5b-數

$\frac{\partial L}{\partial w}$

$\frac{\partial L}{\partial b}$.

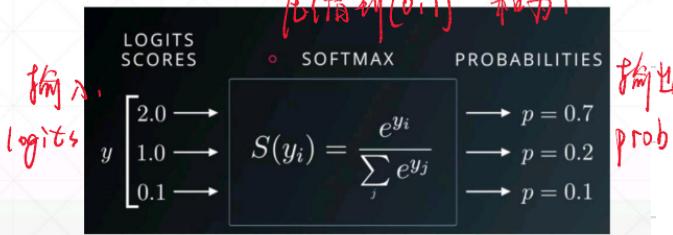
日期:

/

▪ soft version of max

Softmax.

$$\begin{array}{l} \text{logits} \\ \hline a_0 \rightarrow p_0 \\ a_1 \rightarrow p_1 \\ a_2 \rightarrow p_2 \\ a_3 \rightarrow p_3 \end{array}$$

Derivative

$$p_i = \frac{e^{a_i}}{\sum_{k=1}^N e^{a_k}}$$

when $i = j$

$$\begin{aligned} \frac{\partial \frac{e^{a_i}}{\sum_{k=1}^N e^{a_k}}}{\partial a_j} &= \frac{e^{a_i} \sum_{k=1}^N e^{a_k} - e^{a_j} e^{a_i}}{\left(\sum_{k=1}^N e^{a_k}\right)^2} \\ &= \frac{e^{a_i} \left(\sum_{k=1}^N e^{a_k} - e^{a_j}\right)}{\left(\sum_{k=1}^N e^{a_k}\right)^2} \\ &= \frac{e^{a_j}}{\sum_{k=1}^N e^{a_k}} \times \frac{\left(\sum_{k=1}^N e^{a_k} - e^{a_j}\right)}{\sum_{k=1}^N e^{a_k}} \\ &= p_i (1 - p_j) \end{aligned}$$

aj 对应 (j) 的值

Derivative

$$p_i = \frac{e^{a_i}}{\sum_{k=1}^N e^{a_k}}$$

$$\frac{\partial p_i}{\partial a_j} = \frac{\partial \frac{e^{a_i}}{\sum_{k=1}^N e^{a_k}}}{\partial a_j}$$

$$f(x) = \frac{g(x)}{h(x)}$$

$$f'(x) = \frac{g'(x)h(x) - h'(x)g(x)}{h(x)^2}$$

$$g(x) = e^{a_i}$$

$$h(x) = \sum_{k=1}^N e^{a_k}$$

when $i \neq j$

$$\begin{aligned} \frac{\partial \frac{e^{a_i}}{\sum_{k=1}^N e^{a_k}}}{\partial a_j} &= \frac{0 - e^{a_j} e^{a_i}}{\left(\sum_{k=1}^N e^{a_k}\right)^2} \\ &= \frac{-e^{a_j}}{\sum_{k=1}^N e^{a_k}} \times \frac{e^{a_i}}{\sum_{k=1}^N e^{a_k}} \\ &= -p_j \cdot p_i \end{aligned}$$

$$p_i / p_j \in [0, 1]$$

$$\frac{\partial p_i}{\partial a_j} = \begin{cases} p_i(1 - p_j) & \text{if } i = j \\ -p_j \cdot p_i & \text{if } i \neq j \end{cases}$$

正负

$$\text{Or using Kronecker delta } \delta_{ij} = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{if } i \neq j \end{cases}$$

$$\frac{\partial p_i}{\partial a_j} = p_i (\delta_{ij} - p_j)$$

Y I N

日期: /

```
In [3]: x=tf.random.normal([2,4])
In [4]: w=tf.random.normal([4,3])
In [5]: b=tf.zeros([3])
In [6]: y=tf.constant([2,0])

In [14]: with tf.GradientTape() as tape:
...:     tape.watch([w,b])
...:     logits = x@w+b
...:     loss = tf.reduce_mean(tf.losses.categorical_crossentropy(tf.one_hot(y,depth=3), logits,
...: from_logits=True)) prob / logits

In [15]: grads = tape.gradient(loss, [w,b])
<tf.Tensor: id=163, shape=(4, 3), dtype=float32, numpy=
array([[-0.08729011, -0.10937974,  0.19666985],
       [-0.22951077,  0.36995798, -0.14044718],
       [-0.3506433 , -0.2172048 ,  0.56784815],
       [ 0.08480322, -0.26216313,  0.17735994]], dtype=float32)>
In [17]: grads[1] # [-0.07538486,  0.51023775, -0.4348529 ]
```

日期: /

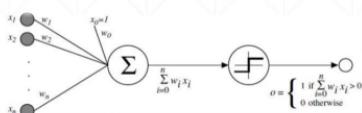
单输出感知机及其梯度

$$y = XW + b$$

单层
↓
输出层

$$y = \sum x_i * w_i + b$$

bias.



→ 离散 (输出层).

Derivative

$$E = \frac{1}{2} (O_0 - t)^2$$

$$\frac{\partial E}{\partial w_{j0}} = (O_0 - t) \frac{\partial O_0}{\partial w_{j0}}$$

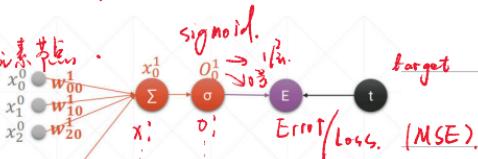
$$\frac{\partial E}{\partial w_{j0}} = (O_0 - t) \frac{\partial \sigma(x_0)}{\partial w_{j0}}$$

$$\frac{\partial E}{\partial w_{j0}} = (O_0 - t) \sigma(x_0)(1 - \sigma(x_0)) \frac{\partial x_0^1}{\partial w_{j0}}$$

$$\frac{\partial E}{\partial w_{j0}} = (O_0 - t) O_0 (1 - O_0) \frac{\partial x_0^1}{\partial w_{j0}}$$

$$\frac{\partial E}{\partial w_{j0}} = (O_0 - t) O_0 (1 - O_0) x_j^0$$

$$O_0 = \sigma(x_0)$$



对应权值
节点

只跟输出 O_0 , t 和输入 x_j^0 有关。

```

x=tf.random.normal([1,3]) 3维.
w=tf.ones([3,1])
b=tf.ones([1])
y = tf.constant([1]) label.

with tf.GradientTape() as tape:
    tape.watch([w, b])
    logits = tf.sigmoid(x@w+b) 1x3 -> 1x1 prob.
    loss = tf.reduce_mean(tf.losses.MSE(y, logits))

grads = tape.gradient(loss, [w, b])
w grad: tf.Tensor
[[ -0.00478814]
 [-0.00588211]
 [ 0.00186196]], shape=(3, 1), dtype=float32)
b grad: tf.Tensor([-0.00444918], shape=(1,), dtype=float32)

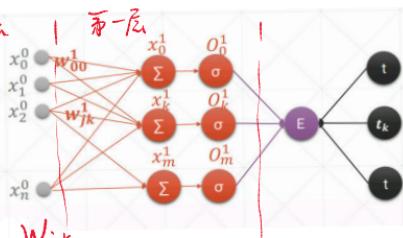
```

输出层

日期:

/

多输出感知机及其梯度.



$$E = \frac{1}{2} \sum (O_i^1 - t_i)^2$$

$$\frac{\partial E}{\partial w_{jk}} = (O_k - t_k) \frac{\partial O_k}{\partial w_{jk}}$$

$$\frac{\partial E}{\partial w_{jk}} = (O_k - t_k) \frac{\partial \sigma(x_k)}{\partial w_{jk}}$$

$$\frac{\partial E}{\partial w_{jk}} = (O_k - t_k) \sigma(x_k)(1 - \sigma(x_k)) \frac{\partial x_k^1}{\partial w_{jk}}$$

$$\frac{\partial E}{\partial w_{jk}} = (O_k - t_k) O_k (1 - O_k) \frac{\partial x_k^1}{\partial w_{jk}}$$

$$\frac{\partial E}{\partial w_{jk}} = (O_k - t_k) O_k (1 - O_k) x_j^0 \quad \frac{\partial E}{\partial w_{jk}} = (O_k - t_k) O_k (1 - O_k) x_j^0$$

w_{jk}

$j \rightarrow k$

x_j^0

$O_k = \sigma(x_k)$

只和第 j 层 x_j^0 以及 x_k^1 对应的

t_k, O_k 有关

可以构成单层

```

In [3]: x=tf.random.normal([2,4])          # 2维输入 x2
In [4]: w=tf.random.normal([4,3])          # 4维降维到3维
In [5]: b=tf.zeros([3])                   # 3个偏置
In [6]: y=tf.constant([2,0])               # 2个目标值
In [7]: y[0] = 1                          # 将 y[0] 改为 1
In [8]: y[1] = 0                          # 将 y[1] 改为 0
In [9]: with tf.GradientTape() as tape:    # 使用 GradientTape
      ...:     tape.watch([w,b])
      ...:     prob = tf.nn.softmax(x@w+b, axis=1)
      ...:     loss = tf.reduce_mean(tf.losses.MSE(tf.one_hot(y, depth=3), prob))
In [10]: grads = tape.gradient(loss, [w,b])
In [11]: grads[0]
[[[-0.00967887, -0.00335512, 0.01303399],
  [-0.04446869, 0.06194263, -0.01747394],
  [-0.04530644, 0.01043231, 0.03487412],
  [0.02006017, -0.03638988, 0.0163297]]]
In [12]: grads[1] # [-0.02585024, 0.06217915, -0.03632889]

```

$\log 1 + \log 2$

日期: /

会话式法则

3) 中间变量

Rules	Function	Derivative
Multiplication by constant	cf	cf'
Power Rule	x^n	nx^{n-1}
Sum Rule	$f + g$	$f' + g'$
Difference Rule	$f - g$	$f' - g'$
Product Rule	fg	$f g' + f' g$
Quotient Rule	f/g	$(f' g - g' f)/g^2$
Reciprocal Rule	$1/f$	$-f'/f^2$
Chain Rule (as "Composition of Functions")	$f \circ g$	$(f' \circ g) \times g'$
Chain Rule (using ')	$f(g(x))$	$f'(g(x))g'(x)$
Chain Rule (using $\frac{dy}{dx}$)	$\frac{dy}{dx} = \frac{dy}{du} \frac{du}{dx}$	

Basic Rule.

$f + g$.

$f - g$.

Product rule.

$$\bullet (fg)' = f'g + fg'$$

乘积。

$$\bullet x^4' = (x^2 * x^2)' = 2x * x^2 + x^2 * 2x = 4x^3$$

Quotient Rule

$$\bullet \left(\frac{f}{g}\right)' = \frac{f'g - fg'}{g^2}$$

▪ e.g. Softmax

$$p_i = \frac{e^{a_i}}{\sum_{k=1}^N e^{a_k}}$$

$$\frac{\partial p_i}{\partial a_j} = \frac{\partial \frac{e^{a_i}}{\sum_{k=1}^N e^{a_k}}}{\partial a_j}$$

$$\begin{aligned} \frac{\partial \frac{e^{a_i}}{\sum_{k=1}^N e^{a_k}}}{\partial a_j} &= \frac{e^{a_i} \sum_{k=1}^N e^{a_k} - e^{a_j} e^{a_i}}{\left(\sum_{k=1}^N e^{a_k}\right)^2} \\ &= \frac{e^{a_i} \left(\sum_{k=1}^N e^{a_k} - e^{a_j}\right)}{\left(\sum_{k=1}^N e^{a_k}\right)^2} \\ &= \frac{e^{a_j}}{\sum_{k=1}^N e^{a_k}} \times \frac{\left(\sum_{k=1}^N e^{a_k} - e^{a_j}\right)}{\sum_{k=1}^N e^{a_k}} \\ &= p_i (1 - p_j) \end{aligned}$$

日期:

Chain rule

$$\frac{\partial y}{\partial x} = \frac{\partial y}{\partial u} \frac{\partial u}{\partial x}$$

$x \xrightarrow[w_1]{b_1} u \xrightarrow[w_2]{b_2} y$
(hidden layer).

$$y_1 = x @ w_1 + b_1$$

$$y_2 = y_1 @ w_2 + b_2$$

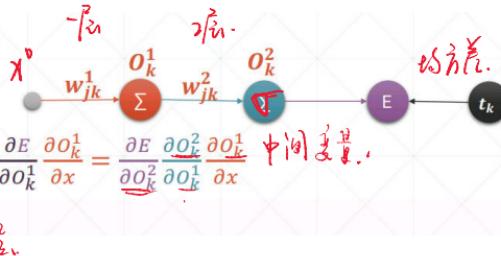
▪ $y_2 = y_1 w_2 + b_2$ $f(y_1) = y_2 = y_1 \cdot w_2 + b_2 \dots$

▪ $y_1 = x w_1 + b_1$ 中間变量.

▪ $\frac{\partial y_2}{\partial w_1} = \frac{\partial f(y_1)}{\partial w_1} = \frac{\partial f(y_1)}{\partial y_1} \frac{\partial y_1}{\partial w_1} = w_2 * x$ $\frac{\partial f(y_1)}{\partial y_1} \rightarrow w_2$

▪ $y_2 = (x w_1 + b_1) * w_2 + b_2$

$$\frac{\partial y_1}{\partial w_1} = x$$



▪ $\frac{\partial E}{\partial w_{jk}} = \frac{\partial E}{\partial o_k^1} \frac{\partial o_k^1}{\partial x} = \frac{\partial E}{\partial o_k^1} \frac{\partial o_k^2}{\partial o_k^1} \frac{\partial o_k^1}{\partial x}$ 中间变量.

```
● ● ●
x = tf.constant(1.)
w1 = tf.constant(2.)
b1 = tf.constant(1.)
w2 = tf.constant(2.)
b2 = tf.constant(1.)
```

```
with tf.GradientTape(persistent=True) as tape:
    tape.watch([w1, b1, w2, b2])
```

```
y1 = x * w1 + b1
y2 = y1 * w2 + b2
```

```
dy2_dy1 = tape.gradient(y2, [y1])[0]
dy1_dw1 = tape.gradient(y1, [w1])[0]
dy2_dw1 = tape.gradient(y2, [w1])[0]
tf.Tensor(2.0, shape=(), dtype=float32)
tf.Tensor(2.0, shape=(), dtype=float32)
```

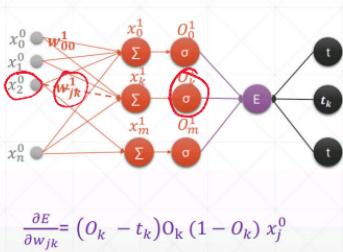
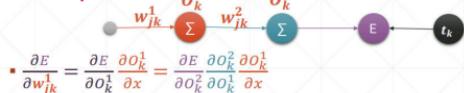
更新:

$$w_1' = w_1 - \frac{\delta L}{\delta w_1} \cdot \text{lr.}$$

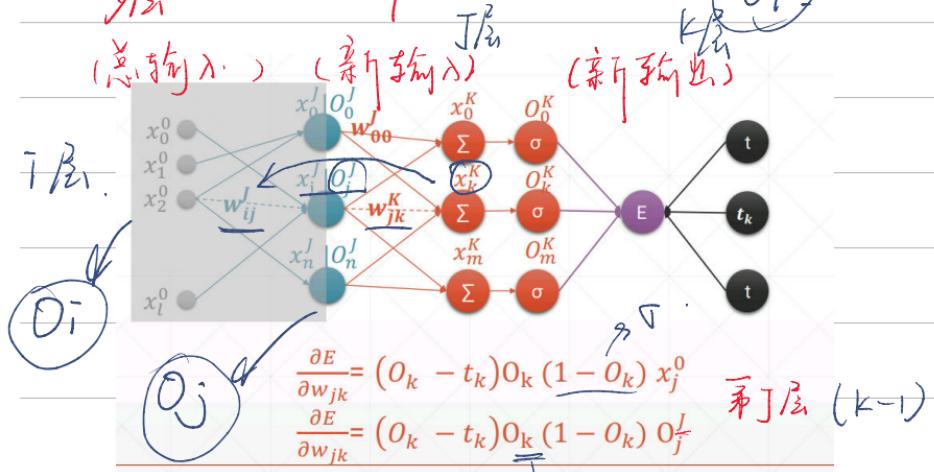
日期: /

多层感知机梯度

链式



Multi-Layer



$$\frac{\partial E}{\partial w_{jk}} = (O_k - t_k) O_k (1 - O_k) O_j^K$$

$$\frac{\partial E}{\partial w_{jk}} = \underbrace{\delta_k^K}_{\downarrow} \underbrace{O_j^K}_{\text{输入.}}$$

跟 O_k 节点有关
(下一层)

O_j^K = 中间变量

$w_{ij} \rightarrow O_j \rightarrow x_k$

日期: $j \rightarrow j - k$

$$\frac{\partial E}{\partial W_{ij}} = \mathcal{O}_j(1 - \mathcal{O}_j)\mathcal{O}_i \sum_{k \in K} (\mathcal{O}_k - t_k)\mathcal{O}_k(1 - \mathcal{O}_k)W_{jk}$$

$$\frac{\partial E}{\partial W_{ij}} = \mathcal{O}_i\mathcal{O}_j(1 - \mathcal{O}_j) \sum_{k \in K} \delta_k W_{jk} \rightarrow j \text{ 与 } k \text{ 无关}$$

$$\delta_k = (\mathcal{O}_k - t_k)\mathcal{O}_k(1 - \mathcal{O}_k).$$

\therefore 与 k 无关

k 层 \rightarrow 最终输出层

For an output layer node $k \in K$

$$\frac{\partial E}{\partial W_{jk}} = \mathcal{O}_j \delta_k$$

上一层节点输出 \mathcal{O}_j

where

$$\delta_k = \mathcal{O}_k(1 - \mathcal{O}_k)(\mathcal{O}_k - t_k)$$

当前层信息

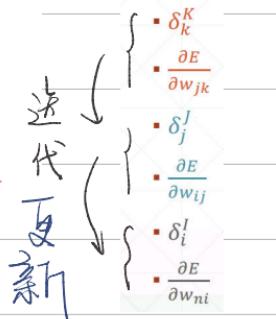
For a hidden layer node $j \in J$

$$\frac{\partial E}{\partial W_{ij}} = \mathcal{O}_i \delta_j$$

where

$$\delta_j = \mathcal{O}_j(1 - \mathcal{O}_j) \sum_{k \in K} \delta_k W_{jk}$$

上一层节点输出



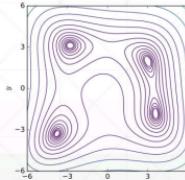
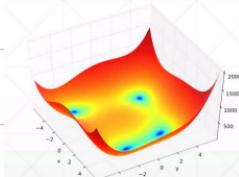
意义: 任何一层的输出都可以由上一层输出、和这一层的节点信息 \mathcal{O} 相乘得到

日期: /

Himmelblau 函数优化

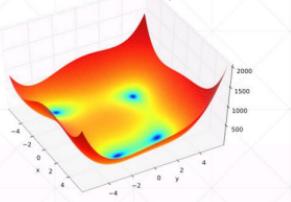
Himmelblau function

$$f(x, y) = (x^2 + y - 11)^2 + (x + y^2 - 7)^2.$$



有四个角 四个局部最小值 一个全局最小值。

- $f(3.0, 2.0) = 0.0,$
- $f(-2.805118, 3.131312) = 0.0,$
- $f(-3.779310, -3.283186) = 0.0$
- $f(3.584428, -1.848126) = 0.0.$



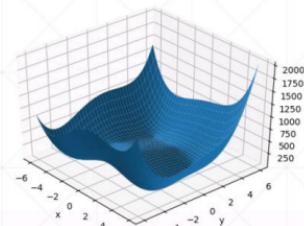
初略是否影响结果：

构造函数。

```
def himmelblau(x):
    return (x[0] ** 2 + x[1] - 11) ** 2 + (x[0] + x[1] ** 2 - 7) ** 2
```

x = np.arange(-6, 6, 0.1) 步进
y = np.arange(-6, 6, 0.1) 生成一系列点
print('x,y range:', x.shape, y.shape)
X, Y = np.meshgrid(x, y) | if meshgrid 生成点
print('X,Y maps:', X.shape, Y.shape)
Z = himmelblau([X, Y])
loss = y

fig = plt.figure('himmelblau')
ax = fig.gca(projection='3d')
ax.plot_surface(X, Y, Z)
ax.view_init(60, -30)
ax.set_xlabel('x')
ax.set_ylabel('y')
plt.show()



部署状态. constant → Tensor 需要 watch
 dense layer 构建 → variable.

```

  ...
  # [1., 0.], [-4, 0.], [4, 0.]
  x = tf.constant([-4., 0.])
  for step in range(200):
    with tf.GradientTape() as tape:
      tape.watch([x])
      y = himmelblau(x)
    grads = tape.gradient(y, [x])[0] 梯度 得到 gradients
    x -= 0.01*grads
  
```

if $x = x - lr * \text{grad}$ 更新参数

Mnist 手写问题实战

Epoch (时期)

当一个完整的数据集通过了神经网络一次。
 (所有样本在神经网络中, 进行一次正向, 一次反向)
 并返回了一次. 这个过程称为一次 epoch.
 另一个 epoch 就是将所有训练样本训练一次
 然而一个 epoch 的样本数量过大, 所以需要
 将分成 batch 进行训练.

日期: /

Batch (批, 样本)

将数据集分成若干个 batch.

Batch-Size 每批样本的大小

Iteration (-次迭代)

训练一个 Batch 就是一次迭代 β - 一个 step.

① 下载数据. 预处理 \rightarrow preprocess.

转化为 Tensor: 定义 batch shuffle 打乱

训练集. 测试集. 迭代器 * reshape.

model $g_{\text{in}} \rightarrow g_{\text{out}} = [b, 10]$

② 构建神经网络
 $[b, 784] \xrightarrow{\text{Dense 1}} [b, 156] \xrightarrow{\text{Dense 2}} [b, 128] \xrightarrow{\text{Dense 3}} [b, 64] \xrightarrow{\text{Dense 4}} [b, 32]$

使用 layer. Dense.

参数 layer. Dense (256, activation=tf.nn.relu).

每个连接. 中间层输出可以随便取.

权值 param: $330 = 32 \times 10 + 10 \rightarrow \text{bias}$.

上一层 \downarrow 上一层 \downarrow

日期： /

③ 完成前向传播

定义多少个 epoch. 对每个 epoch 对 batch 处理
X reshape $28 \times 28 = 784$

将 reshape 之后的 X 放入到网络中

* logit = model(x) 将 x] logit

y → one-hot (形状与 x 一样) -> y)

计算 loss, 计算 grads

④ 完成更新.

keras 自带.

* 定义优化器 optimizer. 自动完成更新逻辑.

使用 optimizer 完成更新

⑤ 测试：使用 test 集.

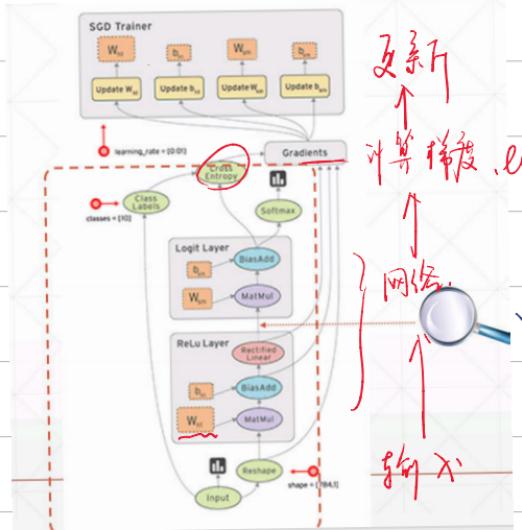
logit $\xrightarrow{\text{softmax}}$ prob $(y) \text{ 为 } \text{one-hot}$
[b, 10] $\xrightarrow{\text{int64} \rightarrow \text{int32}}$ [b] 和 label 进行比较.
True: +1

correct / total_num.

日期: /

Tensor Board 可視化.

如何查看數值?



可視化:

TensorBoard / Visdom (pytorch)

- Installation 安裝

- Listen logdir 監聽日志

- Curves 曲線

- build summary instance 新建日志

- Image Visualization

- fed data into summary instance

Step1.run listener

- open URL: <http://localhost:6006>

cmd → 選入路徑 → dir → tensorboard

```
i@z68:~/TutorialsCN/code_TensorFlow2.0/lesson28-可视化$ tensorboard --logdir logs
TensorBoard 1.13.0 at http://z68:6006 (Press CTRL+C to quit)
[::1], total_correct/total)
```

日期: Step2. build summary

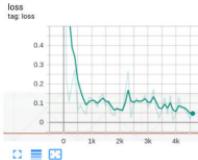
在程序中，创建 writer.

```
...  
current_time = datetime.datetime.now().strftime("%Y%m%d-%H%M%S")  
log_dir = 'logs/' + current_time  
summary_writer = tf.summary.create_file_writer(log_dir) 备注.
```

Step3.fed scalar

标签。

```
...  
with summary_writer.as_default():  
    tf.summary.scalar('loss', float(loss), step=epoch)  
    tf.summary.scalar('accuracy', float(train_accuracy), step=epoch)
```



最近 acc in loss.

Step3.fed single Image

```
...  
# get x from (x,y)  
sample_img = next(iter(db))[0] -张图.  
# get first image instance  
sample_img = sample_img[0] b  
sample_img = tf.reshape(sample_img, [1, 28, 28, 1])  
with summary_writer.as_default():  
    tf.summary.image("Training sample:", sample_img, step=0) (通过报)
```

Step3.fed multi-images

显示时会分单张。

```
...  
val_images = x[:25]  
val_images = tf.reshape(val_images, [-1, 28, 28, 1])  
with summary_writer.as_default():  
    tf.summary.scalar('test-acc', float(test_acc), step=step)  
    tf.summary.image("val-onebyone-images:", val_images, max_outputs=25,  
                    step=step)
```

将多张图片组合成一张。

```
...  
val_images = tf.reshape(val_images, [-1, 28, 28])  
figure = image_grid(val_images)  
tf.summary.image('val-images:', plot_to_image(figure), step=step)
```

日期: /

Keras 高层 API.

Keras Metrics.

默认是 tf.keras.

- Metrics *(继承于, 新建)*
- update_state *(计算)*
- result().numpy() *(输出)*
- reset_states

optimizers.

▪ datasets

▪ layers

▪ losses

▪ metrics

Step1. Build a meter

```
acc_meter = metrics.Accuracy()  
loss_meter = metrics.Mean()
```

Step2. Update data

输入数据.

```
loss_meter.update_state(loss)  
Accuracy. label  
acc_meter.update_state(y, pred)
```

Clear buffer 清理缓存 reset

Step3. Get Average data

```
print(step, 'loss:', loss_meter.result().numpy())  
...  
print(step, 'Evaluate Acc:', total_correct/total, acc_meter.result().numpy())
```

```
if step % 100 == 0:  
    print(step, 'loss:', loss_meter.result().numpy())  
    loss_meter.reset_states()  
  
    # evaluate  
    if step % 500 == 0:  
        total, total_correct = 0., 0  
        acc_meter.reset_states()
```

日期: /

Compile & Fit.

- Compile

裝載 loss 和評估器的邏輯

- Fit

- Evaluate

評估測試

- Predict

預測

固定流程.

```
for i, y in enumerate(db):
    with tf.GradientTape() as tape:
        # [b, 28, 28] => [b, 784]
        x = tf.reshape(x, (-1, 28*28))
        # [b, 784] => [b, 10]
        out = network(x)
        # [b] => [b, 10]
        y_onehot = tf.one_hot(y, depth=10)
        # [b]
        loss = tf.reduce_mean(tf.losses.categorical_crossentropy(y_onehot, out,
from_logits=True))

    grads = tape.gradient(loss, network.trainable_variables)
    optimizer.apply_gradients(zip(grads, network.trainable_variables))
```

↓ 打包成

Compile - Fit

使用后:

```
network.compile(optimizer=optimizers.Adam(lr=0.01),
                 loss=tf.losses.CategoricalCrossentropy(from_logits=True),
                 metrics=['accuracy'])
```

指定: lr
loss.

Accuracy
(準確率)

```
for epoch in range(epochs):
    for step, (x, y) in enumerate(db):
        ...
```



```
network.compile(optimizer=optimizers.Adam(lr=0.01),
                 loss=tf.losses.CategoricalCrossentropy(from_logits=True),
                 metrics=['accuracy'])

network.fit(db, epochs=10)
```

指定數量
for epoch 數量

```
4590/4699 [=====] - 41s 9ms/step - loss: 0.1085
Epoch 2/10
4590/4699 [=====] - 41s 9ms/step - loss: 0.0597
Epoch 3/10
4590/4699 [=====] - 41s 9ms/step - loss: 0.0448
Epoch 4/10
4590/4699 [=====] - 41s 9ms/step - loss: 0.0417
Epoch 5/10
4590/4699 [=====] - 40s 8ms/step - loss: 0.0513
Epoch 6/10
4590/4699 [=====] - 40s 9ms/step - loss: 0.1354
Epoch 7/10
4590/4699 [=====] - 41s 9ms/step - loss: 0.1738
Epoch 8/10
4590/4699 [=====] - 41s 9ms/step - loss: 0.2567
Epoch 9/10
4590/4699 [=====] - 39s 8ms/step - loss: 0.3887
```

日期: /

评估

```
# evaluate
if step % 500 == 0:
    total, total_correct = 0., 0
    for step, (x, y) in enumerate(ds_val):
        # [b, 28, 28] => [b, 784]
        x = tf.reshape(x, (-1, 28*28))
        # [b, 784] => [b, 10]
        out = network(x)
        # [b, 10] => [b]
        pred = tf.argmax(out, axis=1)
        pred = tf.cast(pred, dtype=tf.int32)
        # bool type
        correct = tf.equal(pred, y)
        # bool tensor => int tensor => numpy
        total_correct += tf.reduce_sum(tf.cast(correct,
                                                dtype=tf.int32)).numpy()
        total += x.shape[0]

    print(step, 'Evaluate Acc:', total_correct/total)
```

epoch = 1

test = val

```
network.compile(optimizer=optimizers.Adam(lr=0.01),
                 loss=tf.losses.CategoricalCrossentropy(from_logits=True),
                 metrics=['accuracy'])

network.fit(db, epochs=10, validation_data=ds_val,
            validation_steps=2)
```

固定训练步数

每多个 epoch (步长),

在训练过程中评估样本来一次 val train.

```
169/469 [=====] - 5s 10ms/step - loss: 0.2794 - accuracy: 0.8408 - val_loss:
0.0864 - val_accuracy: 0.9805
Epoch 2/10
169/469 [=====] - 4s 8ms/step - loss: 0.1378 - accuracy: 0.9590 - val_loss:
0.0637 - val_accuracy: 0.9805 relu'
Epoch 3/10 [=====] - 4s 8ms/step - loss: 0.1077 - accuracy: 0.9693 - val_loss:
0.0801 - val_accuracy: 0.9805 relu'
Epoch 4/10 [=====]
169/469 [=====] - 4s 8ms/step - loss: 0.0941 - accuracy: 0.9729 - val_loss:
0.0915 - val_accuracy: 0.9688
Epoch 5/10
169/469 [=====] - 4s 8ms/step - loss: 0.0833 - accuracy: 0.9775 - val_loss:
0.1232 - val_accuracy: 0.9609
Epoch 6/10
169/469 [=====] - 4s 8ms/step - loss: 0.0793 - accuracy: 0.9789 - val_loss:
0.0678 - val_accuracy: 0.9766 logits=True)
Epoch 7/10
169/469 [=====] - 4s 8ms/step - loss: 0.0741 - accuracy: 0.9810 - val_loss:
0.0360 - val_accuracy: 0.9883
Epoch 8/10
```

日期: / 在训练结束后进行.

Test

```
● ● ●  
network.compile(optimizer=optimizers.Adam(lr=0.01),  
                 loss=tf.losses.CategoricalCrossentropy(from_logits=True),  
                 metrics=['accuracy'])  
                )  
  
if test_acc > 0.97  
    SAVE  
    break  
  
network.fit(db, epochs=10, validation_data=ds_val,  
            validation_steps=2)  
  
test ← network.evaluate(ds_val) ← ds_test
```

```
Epoch 8/10  
469/469 [=====] - 4s 8ms/step - loss: 0.0710 - accuracy: 0.9807 - val_loss:  
0.1060 - val_accuracy: 0.9805  
Epoch 9/10  
469/469 [=====] - 4s 8ms/step - loss: 0.0680 - accuracy: 0.9821 - val_loss:  
0.0927 - val_accuracy: 0.9766  
Epoch 10/10  
469/469 [=====] - 4s 9ms/step - loss: 0.0594 - accuracy: 0.9845 - val_loss:  
0.0428 - val_accuracy: 0.9922  
79/79 [=====] - 0s 6ms/step - loss: 0.1205 - accuracy: 0.9757
```

test ←

predict

预测一组数据 .

```
● ● ●  
sample = next(iter(ds_val))  
x = sample[0]  
y = sample[1] # one-hot  
pred = network.predict(x) # [b, 10]  
# convert back to number  
y = tf.argmax(y, axis=1)  
pred = tf.argmax(pred, axis=1)  
  
print(pred)  
print(y)
```

日期: /

自定义网络

▪ keras.Sequential

容器 (层) 串联在一起.

▪ keras.layers.Layer

层 生成层.

▪ keras.Model

```
network = Sequential([layers.Dense(256, activation='relu'),  
                      layers.Dense(128, activation='relu'),  
                      layers.Dense(64, activation='relu'),  
                      layers.Dense(32, activation='relu'),  
                      layers.Dense(10)])  
  
参数  
创建 w. b.  
network.build(input_shape=(None, 28*28))  
network.summary()  
打印信息.  
[ ] 生成 list 生成 5 个.  
network (X).  
激活函数
```

参数管理

▪ model.trainable_variables

[w. b. .w. b. ...] 传给
Dense1 Dense2
optimizer

前向传播

▪ model.call() forward

model (x) 调用网络

Layer/Model (类)

(sequential).

* Inherit from keras.layers.Layer keras.Model

抽象类:

① • __init__

自定义层

② • call

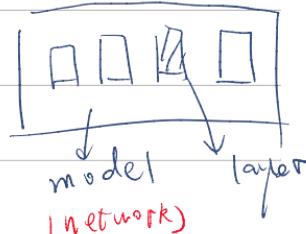
串联.
对子类的调用

model (x).

* Model: compile/fit/evaluate

predict.

5 of 10



layer.Dense (512)

build ([None, 784])

W:[784, 512]

调用类 b [512]

自定义层

类

class MyDense(layers.Layer):

def __init__(self, inp_dim, outp_dim):
super(MyDense, self).__init__()

{ w: self.kernel = self.add_variable('w', [inp_dim, outp_dim])
b: self.bias = self.add_variable('b', [outp_dim])

def call(self, inputs, training=None):

--call-- out = inputs @ self.kernel + self.bias
w b
return out

在 layers 中实现。

784 512

在 model 作用。
model

构建网络。

class MyModel(keras.Model):
def __init__(self):
super(MyModel, self).__init__()
self.fc1 = MyDense(28*28, 256)
self.fc2 = MyDense(256, 128)
self.fc3 = MyDense(128, 64)
self.fc4 = MyDense(64, 32)
self.fc5 = MyDense(32, 10)

前向传播的逻辑。
def call(self, inputs, training=None):

(x = self.fc1(inputs)

x = tf.nn.relu(x)

x = self.fc2(x)

x = tf.nn.relu(x)

x = self.fc3(x)

x = tf.nn.relu(x)

x = self.fc4(x)

x = tf.nn.relu(x)

x = self.fc5(x)

logits:

return x

Dense. Layer !

五层全连接层 fn: x = x - 1

五层输出数

五层全连接层

日期: /

模型保存与加载

Save / load weights 参数

Save / load entire model

saved-model ONNX

管理参数

```
# Save the weights  
model.save_weights('./checkpoints/my_checkpoint')  
  
# Restore the weights  
model = create_model() 加载  
model.load_weights('./checkpoints/my_checkpoint')  
  
loss, acc = model.evaluate(test_images, test_labels)  
print("Restored model, accuracy: {:.2f}%".format(100*acc))
```

fair

```
network.save_weights('weights.ckpt')  
print('saved weights: ')  
del network 释放内存  
  
network = Sequential([layers.Dense(256, activation='relu'),  
                      layers.Dense(128, activation='relu'),  
                      layers.Dense(64, activation='relu'),  
                      layers.Dense(32, activation='relu'),  
                      layers.Dense(10)])  
network.compile(optimizer=optimizers.Adam(lr=0.01),  
                loss=tf.losses.CategoricalCrossentropy(from_logits=True),  
                metrics=['accuracy'])  
network.load_weights('weights.ckpt')  
network.evaluate(ds_val)
```

日期:

/

所有状态 + 结构 + 参数

```
network.save('model.h5')
print('saved total model.')
del network

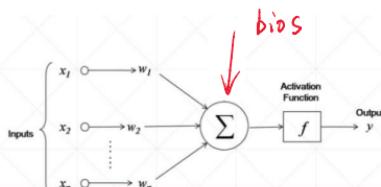
print('load model from file')
network = tf.keras.models.load_model('model.h5')

network.evaluate(x_val, y_val)
```

```
● ● ●
tf.saved_model.save(m, '/tmp/saved_model/')

imported = tf.saved_model.load(path)
f = imported.signatures["serving_default"]
print(f(x=tf.ones([1, 28, 28, 3])))
```

keras 完成 CIFAR10
使用 my dense 层.



保存: 优先保存 权值. 后缀名随便取.

日期: /

过拟合和欠拟合.

factor: noise 噪声、误差. $loss = (wx + b - y)^2$.
 $y = w \cdot x + b + \epsilon$ (eps).

$\epsilon \sim N(0, 0.1)$ 正态分布.

model capacity 容量. (参数量) ↗ 多层.

简单 复杂度不足. Estimated Ground Truth
under-fitting 模型复杂度 < 真实复杂度.
train acc & loss are bad

复杂
Over-fitting Estimated > Ground Truth
train acc & loss are very good
预测值集中在曲线上, but test acc is bad.

Generalization Performance.

日期: /

交叉验证.
Splitting. (train (train)
test x = [60k] train | val
shuffle? 50k (0k.)

tf.split(x, num_or_size_splits=[50000, 10000])

network.fit(train, epochs=5, validation_data=val, validation_freq=1)
训练. 评估. 验证. 测试.

network.evaluate(test) → 测试

Kaggle. train / val / Test.
k-fold cross-validation.
↓
train val | test
+ 随机部分 val.
↓
50k份. 每次50-100.

重新切新
index
索引

```
for epoch in range(500):
    idx = tf.range(60000)
    idx = tf.random.shuffle(idx)
    x_train, y_train = tf.gather(x, idx[:50000]), tf.gather(y, idx[:50000])
    x_val, y_val = tf.gather(x, idx[-10000:]), tf.gather(y, idx[-10000:])

    db_train = tf.data.Dataset.from_tensor_slices((x_train,y_train))
    db_train = db_train.map(preprocess).shuffle(50000).batch(batchsz)

    db_val = tf.data.Dataset.from_tensor_slices((x_val,y_val))
    db_val = db_val.map(preprocess).shuffle(10000).batch(batchsz)

    # training...
    # evaluation...
```

每个 epoch 翻转

日期:

6月12

随机分割

network.fit(db_train_val, epochs=6, validation_split=0.1, validation_freq=2)

Regularization.

Occam's Razor

模型必要，勿增实体。

Reduce Overfitting

More data

Constraint model complexity.

Dropout

① Shallow

Data augmentation.

② regularization 去弊？ 模拟相似？ Early Stopping
(weight Decay) Acc.

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m [y_i \ln \hat{y}_i + (1-y_i) \ln (1-\hat{y}_i)] + \lambda \sum_{i=1}^n (\theta_i)^2$$
$$|\theta| \rightarrow J(W; X, y) + \frac{1}{2} \lambda \|W\|^2$$

通过参数的范数为。

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \underbrace{\beta_n x^n}_{\text{高维参数}} + \varepsilon.$$

高维参数 \rightarrow

但是 regularization \rightarrow 过拟合 \rightarrow 低值。

日期: /

```
l2_model = keras.models.Sequential([
    keras.layers.Dense(16, kernel_regularizer=keras.regularizers.l2(0.001),
                      activation=tf.nn.relu, input_shape=(NUM_WORDS,)),
    keras.layers.Dense(16, kernel_regularizer=keras.regularizers.l2(0.001),
                      activation=tf.nn.relu),
    keras.layers.Dense(1, activation=tf.nn.sigmoid)
])
```

约束器

Flexible regularization

网易云课堂

```
for step, (x, y) in enumerate(db):
    with tf.GradientTape() as tape:
        # ...
        loss = tf.reduce_mean(tf.losses.categorical_crossentropy(y_onehot, out,
                                                               from_logits=True))

        loss_regularization = []
        for p in network.trainable_variables:
            loss_regularization.append(tf.nn.l2_loss(p))
        loss_regularization = tf.reduce_sum(tf.stack(loss_regularization))

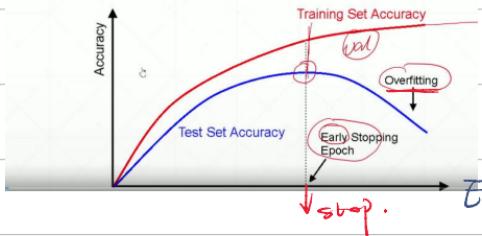
        loss = loss + 0.0001 * loss_regularization
        # grads = tape.gradient(loss, network.trainable_variables)
        optimizer.apply_gradients(zip(grads, network.trainable_variables))
```

梯度
[=+]
[=>]

日期: /

Early stop. Drop out. Stochastic Gradient Decent.
early stop 没定阈值. (val).

▪ Regularization



Monitor

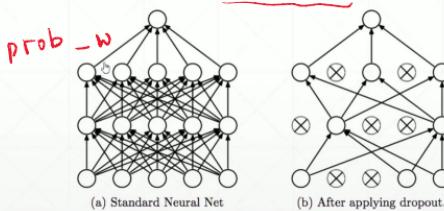
Stop at the
brightest val prof.

Epoch.

Dropout

使 $|W|$ 变小
connection.

- Learning less to learn better
- Each connection has $p = [0, 1]$ to lose



当 $\text{prob_w} \rightarrow 0$ 时。
将连接断开。

随机
阻滞.

```
network = Sequential([layers.Dense(256, activation='relu'),  
                     layers.Dropout(0.5), # 0.5 rate to drop  
                     layers.Dense(128, activation='relu'),  
                     layers.Dropout(0.5), # 0.5 rate to drop  
                     layers.Dense(64, activation='relu'),  
                     layers.Dense(32, activation='relu'),  
                     layers.Dense(10)])
```

随机.

→ 随机 50% drop

$$\frac{\partial L}{\partial w} = \sum_{i=0}^{batch} \frac{\partial L}{\partial w}$$

大概率
一定 100%
随机 -> batch 样本的梯度.

Stochastic Gradient Descent

Deterministic

-- 对应

$x \rightarrow f(x)$.

SGD

所有梯度

信息平均

日期: Convolution / 卷积关注, 滑动, 共享权值.

卷积神经网络. 什么是卷积? channel

Feature map / image

[b, h, w, c]

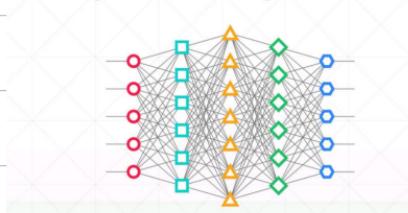
位置信息 特征 {RGB}

[1, 28, 28, 3] \rightarrow [3, 3, 512] \rightarrow [b, 10]

为什么用全连接层?

参数量: 335K or 1.3MB.

▪ 4 Layers: [784, 256, 256, 256, 10]



Layer (type)	Output Shape	Param #
dense (Dense)	multiple	200960
dense_1 (Dense)	multiple	65792
dense_2 (Dense)	multiple	65792
dense_3 (Dense)	multiple	2570
Total params: 335,114		
Trainable params: 335,114		
Non-trainable params: 0		

内存占用巨大. Batch X, Gradient Cache.

只能使用 Dense 层.

整张图

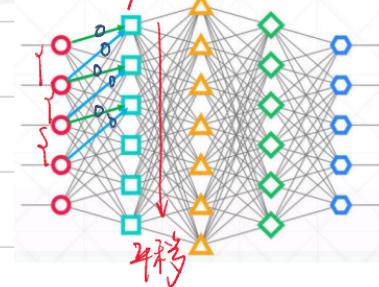
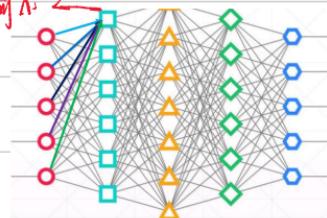
观察窗口

Receptive Field 感知域 (局部)



部分输入 观察方式一样.

所有输入

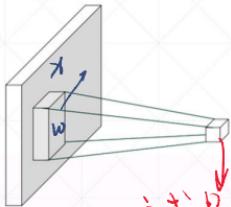


输出相同
一样

Sharing

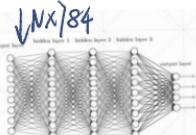
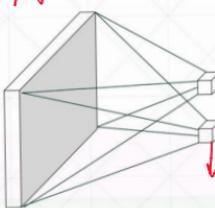
使用相同view滑动

部分：



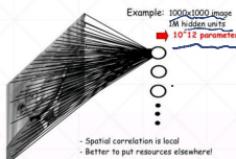
形成一幅新的信息

全局：

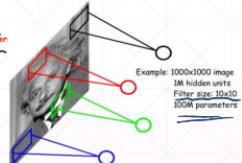


参数减少

FULLY CONNECTED NEURAL NET

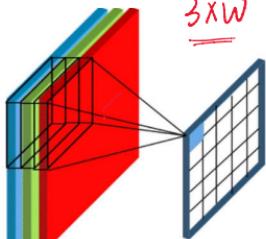
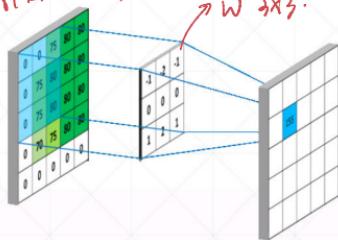


LOCALLY CONNECTED NEURAL NET



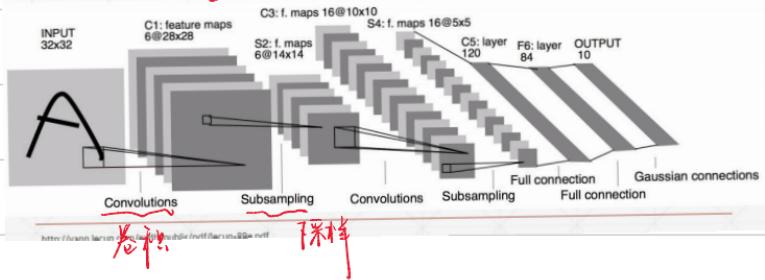
$100 \times M$
→ 100个连接成

X(窗口) 4x4



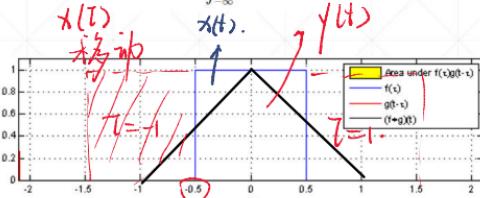
- 6 Layers
- ~60k parameters
- 4 layers, 335K 有连接, 335K

$\sum_{ij} W_{ij} X_{ij}$. 对应相乘两相加



日期:

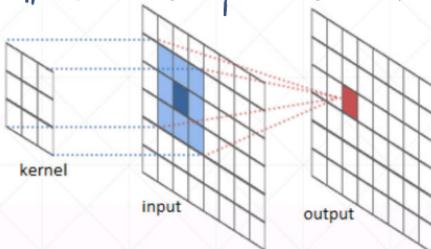
$$y(t) = x(t) * h(t) = \int_{-\infty}^{\infty} x(\tau)h(t-\tau)d\tau$$



2D Convolution $\times(t)$

$$y(t) = x(t) * h(t) = \int_{-\infty}^{\infty} x(\tau)h(t-\tau)d\tau$$

偏置 = x 方向, y 方向. $(\Delta x, \Delta y)$



$g(x, y)$

$f(x, y)$

$o(\Delta x, \Delta y)$

$\xrightarrow{\text{新}} \text{image.}$

当前 $\Delta x, \Delta y$ 偏置
的积. 分.

$\Delta x \in [0, h]$

$\Delta y \in [0, w]$

Computer Vision-

Sharpen:

固定.

0	0	0	0	0	0
0	0	-1	0	0	0
0	-1	5	-1	0	0
0	0	-1	0	0	0
0	0	0	0	0	0



Blur: 打散, 模糊用

0	0	0	0	0
0	1	1	1	0
0	1	3	1	0
0	1	3	1	0
0	0	0	0	0



Sobel 第 8. 样板

Edge Detect:

w	1	0
1	-4	1
0	1	0

↓

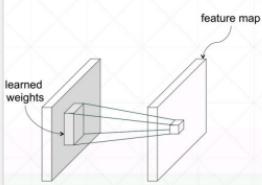
Sobel



改变 kernel \rightarrow 改变观察方式

日期: /

卷积神经网络.



1 channel

1x1	1x0	1x1	0	0
0x0	1x1	1x0	1	0
0x1	0x0	1x1	1	1
0	0	1	1	0
0	1	1	0	0

1x1	1x0	1x1	0	0
0x0	1x1	1x0	1	0
0x1	0x0	1x1	1	1
0	0	1	1	0
0	1	1	0	0

1x1	1x0	1x1	0	0
0x0	1x1	1x0	1	0
0x1	0x0	1x1	1	1
0	0	1	1	0
0	1	1	0	0

1x1	1x0	1x1	0	0
0x0	1x1	1x0	1	0
0x1	0x0	1x1	1	1
0	0	1	1	0
0	1	1	0	0

1x1	1x0	1x1	0	0
0x0	1x1	1x0	1	0
0x1	0x0	1x1	1	1
0	0	1	1	0
0	1	1	0	0

1x1	1x0	1x1	0	0
0x0	1x1	1x0	1	0
0x1	0x0	1x1	1	1
0	0	1	1	0
0	1	1	0	0

1x1	1x0	1x1	0	0
0x0	1x1	1x0	1	0
0x1	0x0	1x1	1	1
0	0	1	1	0
0	1	1	0	0

1x1	1x0	1x1	0	0
0x0	1x1	1x0	1	0
0x1	0x0	1x1	1	1
0	0	1	1	0
0	1	1	0	0

1x1	1x0	1x1	0	0
0x0	1x1	1x0	1	0
0x1	0x0	1x1	1	1
0	0	1	1	0
0	1	1	0	0

1x1	1x0	1x1	0	0
0x0	1x1	1x0	1	0
0x1	0x0	1x1	1	1
0	0	1	1	0
0	1	1	0	0

1x1	1x0	1x1	0	0
0x0	1x1	1x0	1	0
0x1	0x0	1x1	1	1
0	0	1	1	0
0	1	1	0	0

1x1	1x0	1x1	0	0
0x0	1x1	1x0	1	0
0x1	0x0	1x1	1	1
0	0	1	1	0
0	1	1	0	0

1x1	1x0	1x1	0	0
0x0	1x1	1x0	1	0
0x1	0x0	1x1	1	1
0	0	1	1	0
0	1	1	0	0

1x1	1x0	1x1	0	0
0x0	1x1	1x0	1	0
0x1	0x0	1x1	1	1
0	0	1	1	0
0	1	1	0	0

1x1	1x0	1x1	0	0
0x0	1x1	1x0	1	0
0x1	0x0	1x1	1	1
0	0	1	1	0
0	1	1	0	0

1x1	1x0	1x1	0	0
0x0	1x1	1x0	1	0
0x1	0x0	1x1	1	1
0	0	1	1	0
0	1	1	0	0

1x1	1x0	1x1	0	0
0x0	1x1	1x0	1	0
0x1	0x0	1x1	1	1
0	0	1	1	0
0	1	1	0	0

1x1	1x0	1x1	0	0
0x0	1x1	1x0	1	0
0x1	0x0	1x1	1	1
0	0	1	1	0
0	1	1	0	0

1x1	1x0	1x1	0	0
0x0	1x1	1x0	1	0
0x1	0x0	1x1	1	1
0	0	1	1	0
0	1	1	0	0

1x1	1x0	1x1	0	0
0x0	1x1	1x0	1	0
0x1	0x0	1x1	1	1
0	0	1	1	0
0	1	1	0	0

1x1	1x0	1x1	0	0
0x0	1x1	1x0	1	0
0x1	0x0	1x1	1	1
0	0	1	1	0
0	1	1	0	0

1x1	1x0	1x1	0	0
0x0	1x1	1x0	1	0
0x1	0x0	1x1	1	1
0	0	1	1	0
0	1	1	0	0

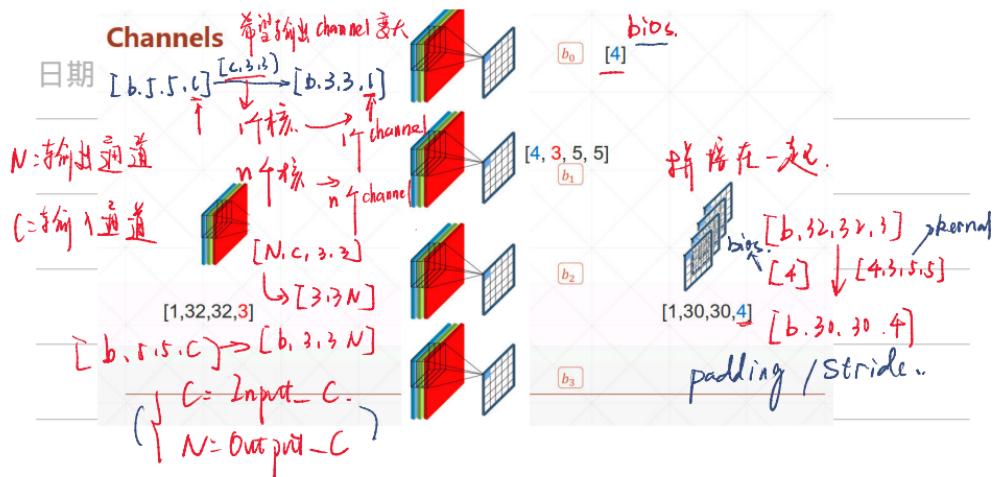
1x1	1x0	1x1	0	0
0x0	1x1	1x0	1	0
0x1	0x0	1x1	1	1
0	0	1	1	0
0	1	1	0	0

1x1	1x0	1x1	0	0
0x0	1x1	1x0	1	0
0x1	0x0	1x1	1	1
0	0	1	1	0
0	1	1	0	0

1x1	1x0	1x1	0	0
0x0	1x1	1x0	1	0
0x1	0x0	1x1	1	1
0	0	1	1	0
0	1	1	0	0

1x1	1x0	1x1	0	0
0x0	1x1	1x0	1	0
0x1	0x0	1x1	1	1
0	0	1	1	0
0	1	1	0	0

1x1	1x0	1x1	0	0
0x0	1x1	1x0	1	0
0x1	0x0	1x1	1	1
0	0	1	1	0
0	1	1	0	0</td



For instance

D

G

B

Input Volume (+pad 1) (7x7x3)
 $x(t, i, j, 0)$

0 0 0 0 0 0 0	-1 0 1
0 0 0 1 0 2 0	0 0 1
0 1 0 2 0 1 0	1 -1 1
0 1 2 0 2 2 0 0	w0(t, i, 1, 1)
0 2 0 0 2 0 0	1 0 1
0 2 1 2 2 0 0	1 -1 1
0 0 0 0 0 0 0	0 1 0

Filter W0 (3x3x3)
 $w0(t, i, j, 0)$

0 1 -1	0 + 0x + 0x3
0 -1 0	0 + 0x + 0x3
0 -1 1	0 + 0x + 0x3

Output Volume (3x3x2)
 $o(t, i, 0)$

2 3 3	[b, 3, 3, 2]
3 7 3	
8 10 3	

filter movement

Bias b0 (1x1x1)
 $b0(t, i, 0)$

Bias b1 (1x1x1)
 $b1(t, i, 0)$

t, i, r

0 + 0x + 0x3

Output Volume (3x3x2)
 $o(t, i, 0)$

2 3 3	[b, 3, 3, 2]
3 7 3	
8 10 3	

filter movement

filter W1 (3x3x3)
 $w1(t, i, j, 0)$

0 1 -1	0 + 0x + 0x3
0 -1 0	0 + 0x + 0x3
0 -1 1	0 + 0x + 0x3

filter W1 (3x3x2)
 $w1(t, i, j, 1)$

1 0 0	0 + 0x + 0x3
1 -1 0	0 + 0x + 0x3
1 0 0	0 + 0x + 0x3

filter W1 (3x3x2)
 $w1(t, i, j, 2)$

1 1 -1	0 + 0x + 0x3
0 -1 -1	0 + 0x + 0x3
1 0 0	0 + 0x + 0x3

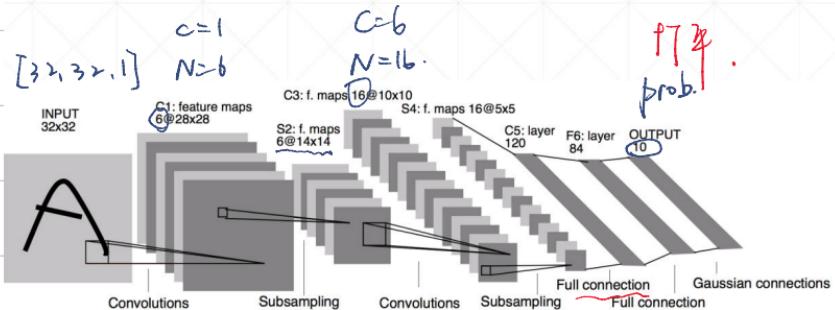
filter W1 (3x3x2)
 $w1(t, i, j, 3)$

1 1 1	0 + 0x + 0x3
0 -1 1	0 + 0x + 0x3
1 0 0	0 + 0x + 0x3

x: [b, 28, 28, 1]
one k: (3, 3, 3)
multi-k: (16, 3, 3, 3)
stride: 1
padding: [1, 1, 1, 1]
bias: (16)
out: [b, 28, 28, 16]

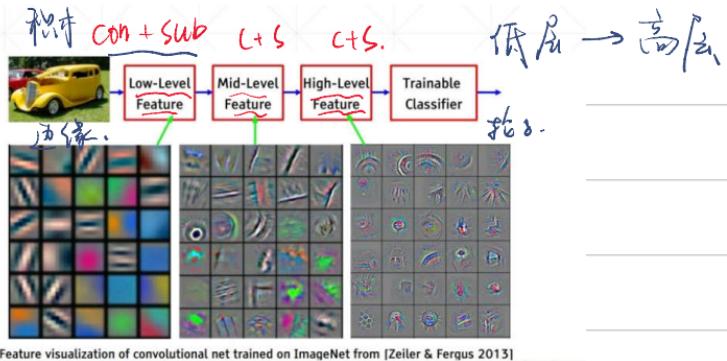
<https://skymind.ai/wiki/convolutional-network>

LeNet-5



日期:

作业量.



In [1]: import tensorflow as tf
In [2]: from tensorflow.keras import layers
In [6]: layer = layers.Conv2D(4, kernel_size=5, strides=1, padding='valid')
In [8]: out = layer(x)
Out[9]: TensorShape([1, 28, 28, 4])

In [10]: layer = layers.Conv2D(4, kernel_size=5, strides=1, padding='same')
In [11]: out = layer(x)
Out[12]: TensorShape([1, 32, 32, 4])

In [13]: layer = layers.Conv2D(4, kernel_size=5, strides=2, padding='same')
In [14]: out = layer(x)
Out[15]: TensorShape([1, 16, 16, 4])

In [16]: layer.call(x).shape
Out[16]: TensorShape([1, 16, 16, 4])

逐層分析：

逐層分析：

In [13]: layer = layers.Conv2D(4, kernel_size=5, strides=2, padding='same')
In [14]: out = layer(x)
Out[15]: TensorShape([1, 16, 16, 4])

In [17]: layer.kernel
Out[17]: $[5, 5, 3, 4]$

In [18]: layer.bias
Out[18]: $\text{tf.Variable } 'conv2d_3/bias:0' \text{ shape}=(4,) \text{ dtype}=\text{float32}, \text{ numpy}=\text{array}([0., 0., 0., 0.], \text{ dtype}=\text{float32})$

日期:

底稿

```

In [21]: w=tf.random.normal([5,5,3,4])
In [22]: b=tf.zeros([4])
In [23]: x.shape
Out[23]: TensorShape([1, 32, 32, 3])

In [29]: out=tf.nn.conv2d(x, w, strides=1, padding='VALID')
Out[30]: TensorShape([1, 28, 28, 4])

In [31]: out = out + b
Out[32]: TensorShape([1, 28, 28, 4])

In [33]: out=tf.nn.conv2d(x,w,strides=2,padding='VALID')
Out[34]: TensorShape([1, 14, 14, 4])

```

自己建 w, b

写 kernel

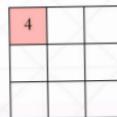
conv2d

kernel

Gradient ?

$\frac{\partial \text{Loss}}{\partial w}$

1x1	1x0	1x1	0	0
0x0	1x1	1x0	1	0
0x1	0x0	1x1	1	1
0	0	1	1	0
0	1	1	0	0



Input

$X [b, 3, 3, 1]$

x_{00}	x_{01}	x_{02}
x_{10}	x_{11}	x_{12}
x_{20}	x_{21}	x_{22}

weight
kernel
 $[2, r.c.N]$

w_{00}	w_{01}
w_{10}	w_{11}

Output

O_{00}	O_{01}
O_{10}	O_{11}

$\xrightarrow{\text{FC/CNN}}$

prob.

0
1
0
0

$[b, 2, 2, 1]$

相乘、相加.

$$O_{00} = x_{00} * w_{00} + x_{01} * w_{01} + x_{10} * w_{10} + x_{11} * w_{11} + b$$

$$O_{01} = x_{01} * w_{00} + x_{02} * w_{01} + x_{11} * w_{10} + x_{12} * w_{11} + b$$

$$O_{10} = x_{10} * w_{00} + x_{11} * w_{01} + x_{20} * w_{10} + x_{21} * w_{11} + b$$

$$O_{11} = x_{11} * w_{00} + x_{12} * w_{01} + x_{21} * w_{10} + x_{22} * w_{11} + b$$

DI 中间变量.

$$\frac{\partial \text{Loss}}{\partial w_{00}} = \sum_{i \in [0,0,1,10,11]} \frac{\partial \text{Loss}}{\partial O_i} \frac{\partial O_i}{\partial w_{00}}$$

$$\frac{\partial O_{00}}{\partial w_{00}} = \frac{\partial (x_{00} * w_{00} + x_{01} * w_{01} + x_{10} * w_{10} + x_{11} * w_{11} + b)}{\partial w_{00}} = x_{00}$$

日期: /

池化与采样.

Max/Avg pooling.
• stride=2
 $\frac{1}{2}$
 $\frac{1}{2}$
 $\frac{1}{2}$

• stride=1

4x4

1	3	2	4
1	3	2	4
5	7	6	8
5	7	6	8

3x3

3	3	4
7	7	8
7	7	8

stride=2

• stride=2

1	3	2	4
1	3	2	4
5	7	6	8
5	7	6	8

$\frac{1+3+3+1}{4}$

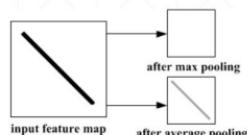
2	3
6	7

Average pooling

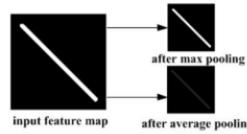
4x4

2x2

Max pooling



(a) Illustration of max pooling drawback



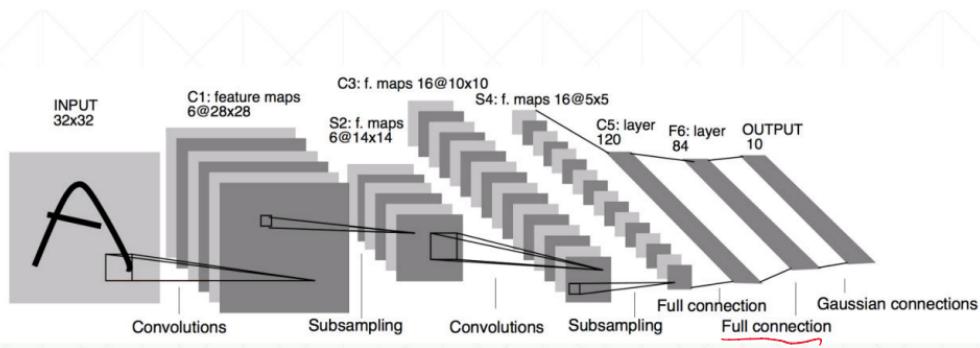
0	255	255	255
255	0	255	255
255	255	0	255
255	255	255	0

digit express of the pooling process

255	0
0	255
0	0
0	255

digit express of the pooling process

(b) Illustration of average pooling drawback



日期:

无 channel
无关

```
In [36]: x # TensorShape([1, 14, 14, 4])
In [37]: pool=layer.MaxPool2D(2, strides=2)
In [38]: out=pool(x)
Out[39]: TensorShape([1, 7, 7, 4])
```

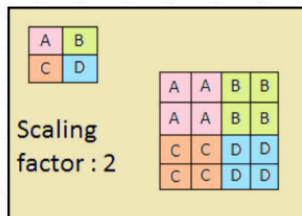
✓ ✓ ✓
kernel size → 2x2

```
In [40]: pool=layer.MaxPool2D(3, strides=2)
In [41]: out=pool(x)
Out[42]: TensorShape([1, 6, 6, 4])
```

```
In [44]: out=tf.nn.max_pool2d(x, 2, strides=2, padding='VALID')
Out[45]: TensorShape([1, 7, 7, 4])
```

UpSampling

- nearest
- bilinear



```
In [47]: x=tf.random.normal([1,7,7,4])
In [48]: layer=layer.UpSampling2D(size=3)
In [49]: out=layer(x)
Out[50]: TensorShape([1, 21, 21, 4])
```

→ 3x3

```
In [51]: layer=layer.UpSampling2D(size=2)
In [52]: out=layer(x)
Out[53]: TensorShape([1, 14, 14, 4])
```

→ 2x2.

ReLU.

↑ → ↗

负值 → 0

```
In [55]: x=tf.random.normal([2,3])
<tf.Tensor: id=154, shape=(2, 3), dtype=float32, numpy=
array([-1.533682 , -2.705333 ,  0.36354962],
      [ 0.00713745,  0.69756126,  0.8053344 ], dtype=float32)>
```

```
In [57]: tf.nn.relu(x)
<tf.Tensor: id=156, shape=(2, 3), dtype=float32, numpy=
array([[0.          , 0.          , 0.36354962],
      [ 0.00713745,  0.69756126,  0.8053344 ]], dtype=float32)>
```

```
In [59]: layers.ReLU()(x)
<tf.Tensor: id=158, shape=(2, 3), dtype=float32, numpy=
array([[0.          , 0.          , 0.36354962],
      [ 0.00713745,  0.69756126,  0.8053344 ]], dtype=float32)>
```

日期: /

CIFAR 100 与 VGG 实战

▪ Load datasets

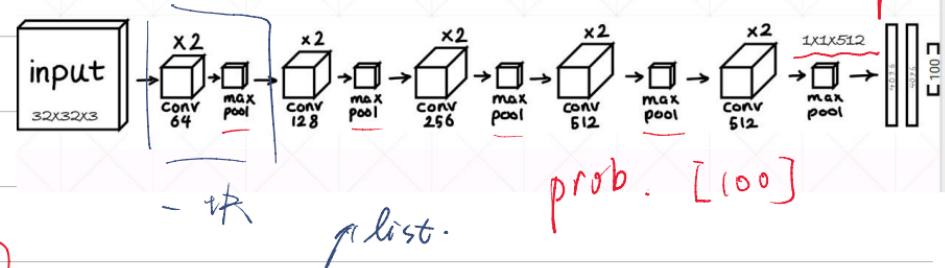
▪ Build Network

▪ Train

▪ Test

自连接.

13 Layers. 5 层



①

conv_layer []

生成特征层. N 直定 (输出)

strides

layer.Conv2D(64, kernel_size=[3,3], padding='same',
channel ↗ activation='tf.nn.relu'), 取活小数

生成池化层.

size

layer.MaxPool2D(pool_size=[2,2],

步长. strides=2, padding='same')

- 每 channel (N) 进后增大.

每张尺寸 (w, h) 进后减少

日期: /

②

main 函数:

conv_net = Sequential ([conv_layers])

[b, 32, 32, 3] \Rightarrow [b, ..., 512]

fc_net = Sequential ([...])

conv_net.build (input_shape=[None, 32, 32, 3]) \rightarrow 权值.

fc_net.build (input_shape=[None, 512])

③

preprocess. 预处理. ④ 构建 fw 批裁 ⑤ 创建数据集

out = conv_net(x)

[b, 32, 32, 3] \Rightarrow [b, 1, 1, 512]

⑥ out = tf.reshape(out, [-1, 512]) flatten \Rightarrow [b, 512] *

logits = fc_net(out) \rightarrow prob. [b, 512] \Rightarrow [b, 100]

one-hot \rightarrow loss \rightarrow gradient.

参数.

variables = conv_net.trainable_variables + fc_net.trainable_variables

Optimizer = Optimizer.Adam(learning_rate=1e-4)

test.

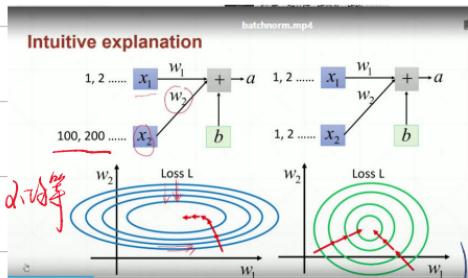
Logit $\xrightarrow{\text{softmax}}$ prob. $\xrightarrow{\arg\max}$ pred $\xrightarrow{\text{inverted}}$ correct.

日期: /

Batch Norm (BN)

Sigmoid \rightarrow 梯度离散.

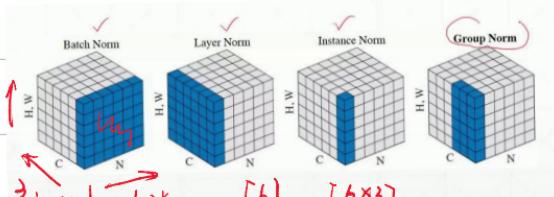
Normalization. 均匀化. 小公极大的分布.



Output 在 0 周围分布.

Batch Normalization Dynamic mean

batch $[N, C, H \times W]$ $[6, 3, 784]$



$[6, 3, 784]$ $[18]$

34 channel } $\text{cho} \rightarrow \text{mean} \rightarrow 1\text{数. batch:}$
 $\text{chi} \rightarrow \text{mean}$
 $\text{ch}_2 \rightarrow \text{mean}$ $\text{输出值} = C.$

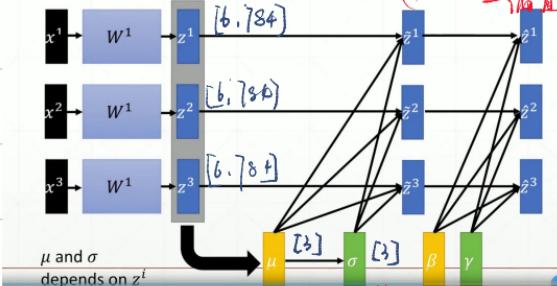
每个 channel feature 值. [3]

日期:

Batch normalization

$$\bar{z}^i = \frac{(z^i) - \mu}{\sigma} \sim N(0, 1)$$

$$\hat{z}^i = \gamma \odot \bar{z}^i + \beta \rightarrow N(\beta, \gamma)$$

running β running γ

(历史信息)

Input: Values of x over a mini-batch: $B = \{x_{1..m}\}$;Parameters to be learned: γ, β **Output:** $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$ μ, γ

均值

$$\mu_B \leftarrow \frac{1}{m} \sum_{i=1}^m x_i$$

// mini-batch mean

自动更新 方差

$$\sigma_B^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2$$

// mini-batch variance

 x_i

$$\hat{x}_i \leftarrow \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

 $N(\beta, \gamma)$ // normalize

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i)$$

缩放

// scale and shift

Algorithm 1: Batch Normalizing Transform, applied to activation x over a mini-batch.

[3, 3, N, C]

net = layers.BatchNormalization()

axis=-1,

center=True,

scale=True

trainable=True

• 通用，指定维度 (轴, axis)

偏移 缩放

net(x, training=None)

true. 打开 training / test

日期

```
In [3]: net=layers.BatchNormalization()
In [5]: x=tf.random.normal([2,3])
In [6]: out=net(x) 金鐘秀.
```

In [7]: net.trainable_variables
[<tf.Variable 'batch_normalization/gamma:0' shape=(3,) dtype=float32, numpy=array([1., 1., 1.],
dtype=float32),
<tf.Variable 'batch_normalization/beta:0' shape=(3,) dtype=float32, numpy=array([0., 0., 0.],
dtype=float32)>]

In [8]: net.variables
[<tf.Variable 'batch_normalization/gamma:0' shape=(3,) dtype=float32, numpy=array([1., 1., 1.],
dtype=float32),
<tf.Variable 'batch_normalization/beta:0' shape=(3,) dtype=float32, numpy=array([0., 0., 0.],
dtype=float32),
<tf.Variable 'batch_normalization/moving_mean:0' shape=(3,) dtype=float32, numpy=array([0., 0., 0.],
dtype=float32),
<tf.Variable 'batch_normalization/moving_variance:0' shape=(3,) dtype=float32, numpy=array([1., 1., 1.],
dtype=float32)]

隨便寫
東西
更新

Backward update.

```
for i in range(10):
    with tf.GradientTape() as tape:
        out = net(x, training=True)
        loss = tf.reduce_mean(tf.pow(out,2)) - 1
    grads = tape.gradient(loss, net.trainable_variables)
    optimizer.apply_gradients(zip(grads, net.trainable_variables))

backward(10 steps):
[<tf.Variable 'batch_normalization/gamma:0' shape=(3,) dtype=float32, numpy=array([0.93549937,
0.9356556 , 0.9355564 ], dtype=float32),
<tf.Variable 'batch_normalization/beta:0' shape=(3,) dtype=float32, numpy=array([
1.3411044e-09, 1.3411045e-08, -4.0978188e-10], dtype=float32)>...]
```

偏差變小，收斂變快。

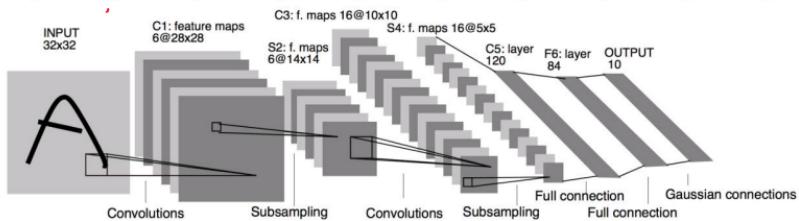
Better performance Converge faster
Robust (Stable, larger learning rate)

日期: /

经典卷积神经网络. VGG, Google net, Inception.

LeNet-5 手写数字识别, Mnist

- 99.2% acc.
- 5/6 layers



AlexNet

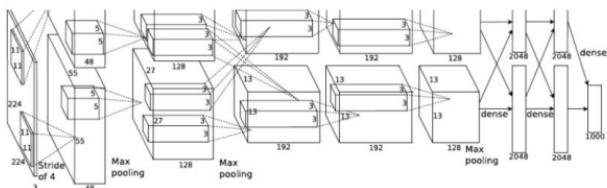
- GTX 580
- 3GBx2
- 11X11
- 8 layers

5+3.

[224, 224, 3]

Kernel: [11x11]

AlexNet: ILSVRC 2012 winner



- Similar framework to LeNet but:
 - Max pooling, ReLU nonlinearity
 - More data and bigger model (7 hidden layers, 650K units, 60M params)
 - GPU implementation (50x speedup over CPU)
 - Trained on two GPUs for a week
 - Dropout regularization

A. Krizhevsky, I. Sutskever, and G. Hinton,
[ImageNet Classification with Deep Convolutional Neural Networks](#), NIPS 2012

VGG

11 x 11

參數大

VGGNet: ILSVRC 2014 2nd place

- 3x3 ✓
 - 1x1 ✓
 - 11-19 layers

b个版本

V6G 11 / 16 / 19...

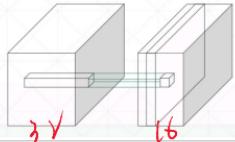
Conv + pooling

Table 2: Number of parameters (in millions).

- Sequence of deeper networks trained progressively
 - Large receptive fields replaced by successive layers of 3×3 convolutions (with ReLU in between)
 - One 7×7 conv layer with C feature maps needs $49C^2$ weights, three 3×3 conv layers need only $27C^2$ weights
 - Experimented with 1x1 convolutions

1x1 Convolution

- less computation
 - $c_{in} \Rightarrow c_{out}$



[16, 32, 1, 1]

17

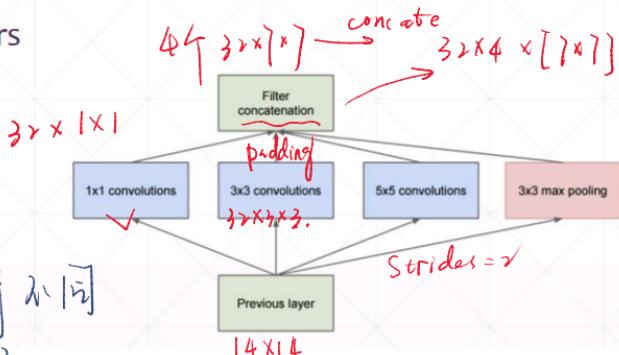
改變 channel 數.!!

维度改变

GoogLeNet

- 1st in 2014 ILSVRC
 - 22 layers

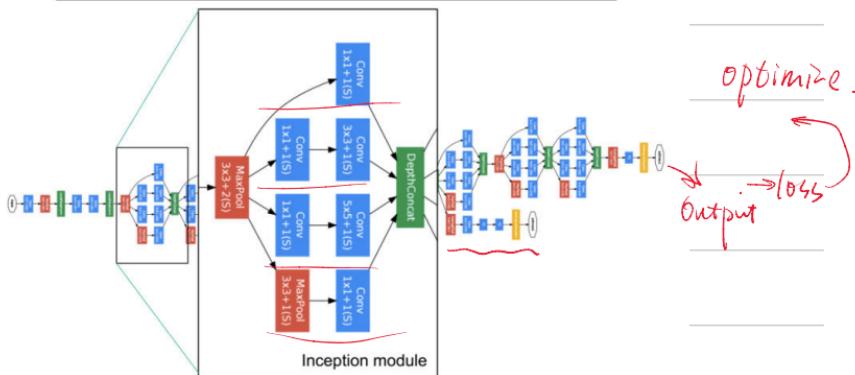
答：[12]-12，用小12表示差异数。



感受視野不同
肩 + 肩部。

日期：

GoogLeNet



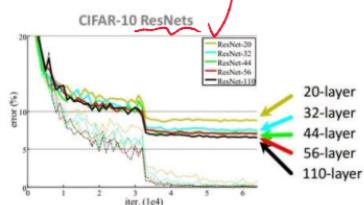
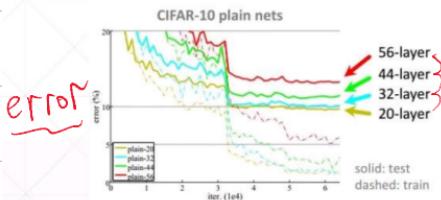
C. Szegedy et al., [Going deeper with convolutions](#), CVPR 2015

Stack more layers? No!!!

- 1000 layers?

But ResNet

CIFAR-10 experiments



日期:

深度残差网络

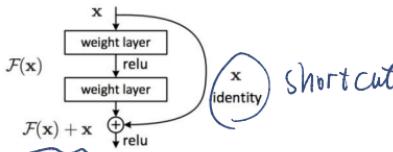
ResNet & DenseNet

$30 \rightarrow 2V$

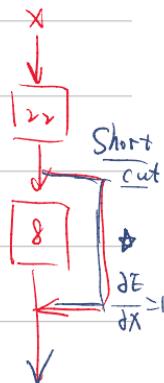
ResNet

The residual module

- Introduce skip or shortcut connections (existing before in various forms in literature)
- Make it easy for network layers to represent the identity mapping
- For some reason, need to skip at least two layers



Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun,
[Deep Residual Learning for Image Recognition](#), CVPR 2016 (Best Paper)



$\dots + \text{conv} + \text{relu} + \text{conv} + \text{relu} + \dots$

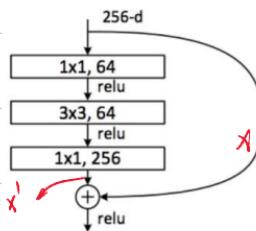
$256 \rightarrow 64 \rightarrow 64 \rightarrow 256$

ResNet

每个 unit 的输入输出 size 一致

Deeper residual module (bottleneck)

Unit



Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun,
[Deep Residual Learning for Image Recognition](#), CVPR 2016 (Best Paper)

- Directly performing 3×3 convolutions with 256 feature maps at input and output:
 $256 \times 256 \times 3 \times 3 \sim 600K$ operations

直接使用

[256, 256, 3, 3]

- Using 1×1 convolutions to reduce 256 to 64 feature maps, followed by 3×3 convolutions, followed by 1×1 convolutions to expand back to 256 maps:

$256 \times 64 \times 1 \times 1 \sim 16K$

$64 \times 64 \times 3 \times 3 \sim 36K$

$64 \times 256 \times 1 \times 1 \sim 16K$

Total: ~70K

参数较少。

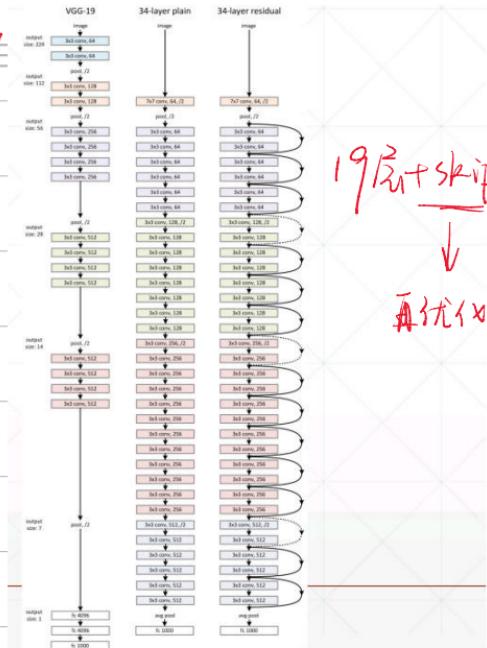
[64, 64, 1, 1]

[64, 64, 3, 3]

↓

[256, 256, 3, 3]

19

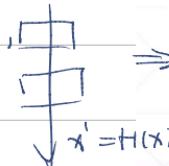


19层+skip

优点



Why call Residual?

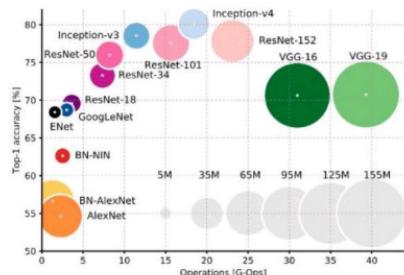
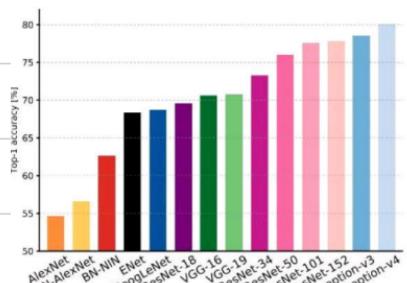


$\rightarrow x' = H(x)$

$$F(x) := H(x) - x$$

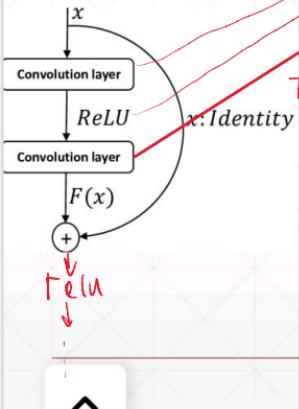
→ 简化梯度流

→ VGG-19



Basic Block

$\text{conv} - \text{BN} - \text{ReLU}$



```

class BasicBlock(layers.Layer):
    def __init__(self, filter_num, stride=1):
        super(BasicBlock, self).__init__()

        self.conv1 = layers.Conv2D(filter_num, (3, 3), strides=stride, padding='same')
        self.bn1 = layers.BatchNormalization()
        self.relu = layers.Activation('relu')
        self.conv2 = layers.Conv2D(filter_num, (3, 3), strides=1, padding='same')
        self.bn2 = layers.BatchNormalization()

        if stride == 1:
            self.downsample = Sequential()
            self.downsample.add(layers.Conv2D(filter_num, (1, 1), strides=stride))
            self.downsample.add(layers.BatchNormalization())
        else:
            self.downsample = lambda x: x

        self.stride = stride

    def call(self, inputs, training=None):
        residual = self.downsample(inputs)

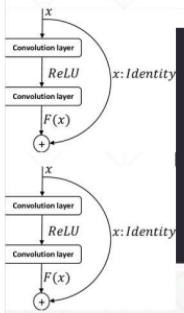
        conv1 = self.conv1(inputs)
        bn1 = self.bn1(conv1)
        relu1 = self.relu(bn1)
        conv2 = self.conv2(relu1)
        bn2 = self.bn2(conv2)

        add = layers.add([bn2, residual])
        out = self.relu(add)
        return out
  
```

Diagram annotations in red:

- Red arrows point from the input x to the first convolution layer and from the output of the first convolution layer to the second convolution layer.
- A red circle highlights the $x: \text{Identity}$ label.
- A red circle highlights the ReLU label at the bottom right.
- A red circle highlights the $\text{bn}_2 + \text{res} \Rightarrow \text{out}$ label.

Res Block

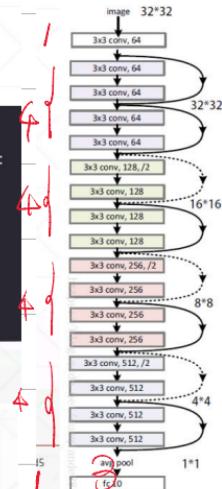


```

def _build_resblock(self, block, filter_num, blocks, stride=1):
    res_blocks = keras.Sequential()
    res_blocks.add(block(filter_num, stride))

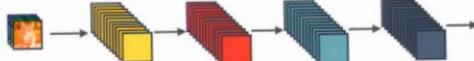
    for _ in range(1, blocks):
        res_blocks.add(block(filter_num, stride=1))

    return res_blocks
  
```

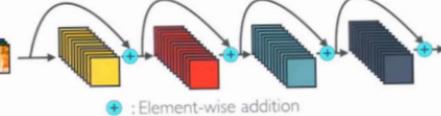


DenseNet

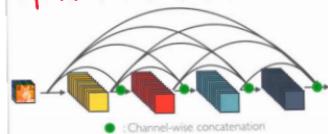
Flat



Res



双流?



日期:

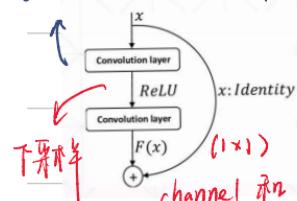
class Basic Block

ResNet 実装

固定变量

① Basic Block

CONV + BN + ReLU.



结构定义

self.conv1=layer.Conv2D(stride=Stride)

Self.bn1=layer.BatchNormalization()

Self.relu=layer.Activation('relu')

正向. self.conv1=layer.Conv2D(stride=Stride)

Self.bn1=layer.BatchNormalization()

if: stride != 1 (输出维度发生改变,

Self.downsample=sequential()

Self.downsample.add(layers.Conv2D((1,1))

传播过程

def call().

out = self.conv1(input)

out = self.bn1(out)

out = self.relu(out)

out = self.bn2(out)

out = self.conv2(input)

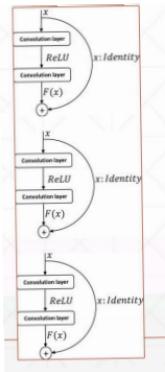
identity = self.downsample(inputs)

output = layers.add([out, identity]) 跳层相加

output = self.relu(output)

relu 层没有参数 当作函数

return out



Res Block (\nrightarrow 4 basic block)

class Res Block

def build_resblock () .

res_block = sequential ()

res_blocks.add (BasicBlock (N, stride))

for _ in range (1, blocks) = ↓ 不一样

res_blocks.add (BasicBlock (stride=))

return res_block.

Res Net, 由 3 个 Res Block 构成

→ Res Block

结构: def __init__(self, layer_dim): [2, 2, 2, 2]

方法:

示例: self.stem = sequential ([layer.conv]D[64, 3, 3])

Stride=1, layers, Batch Normalization.

layers, Activation ('relu'), layers, MaxPool2D(...)

Res Block

由 3 个 Res Block.

元素

self.layer1 = self.build_resblock(64, layer_dim[0])

self.layer2 = self.build_resblock(128, layer_dim[1], stride=2)

self.layer3 = self.build_resblock(256, layer_dim[2], stride=2)

self.layer4 = self.build_resblock(512, layer_dim[3], stride=2)

Fc.

10

日期: / out put = [b, 512, h, w] $[b, 512]$
512 vector
 $[512 \times 1 \times 1]$

分类: self.avgpool = layers.Global Average Pooling 2D[1].

self.fc = layers.Dense (num_classes). ↓

前向传播. def call ↓ 输出维度 100

x = self.stem(input)

x = self.layer1(x)

x = self.layer2(x)

x = self.layer3(x)

x = self.layer4(x)

x = self.avgpool(x) [b, c]

x = self.fc(x) [b, 100]

return x.

定义 ResNet 18

ResNet 34.

[3, 4, 6, 3]

def ResNet18():

return ResNet([2, 2, 2, 2])

日期: /

ResNet - train

```
from resnet import resnet18 导入 resnet
```

preprocess 预处理, 把 b.

$$\downarrow [b, 32, 32, 3] \rightarrow [b, 1, 1, 512]$$

main model = resnet18() *(参数)*

```
model.build(input_shape=[None, 32, 32, 3])
```

weight_decay. model.summary() *(参数)*

```
optimizer = optimizer.Adam(lr=1e-3)
```

for epoch in range(...):

for step in enumerate(train_db):

```
logit  $\rightarrow$  loss  $\rightarrow$  grad  $\rightarrow$  optimize.
```

解决办法

convolution

- 1. decrease batch size

- 2. tune resnet [2, 2, 2, 2]

- 3. try Google Colab \rightarrow 1<80 VGP 显存

- 4. buy new NVIDIA GPU Card

日期： /

预训练网络（迁移学习）

↓
保存在 Keras 库中包含
VGG16、VGG19、
ResNet50、
Inception v3、
Xception 等经典的模型架构。

可移植的、

在 ImageNet 上

进行训练。

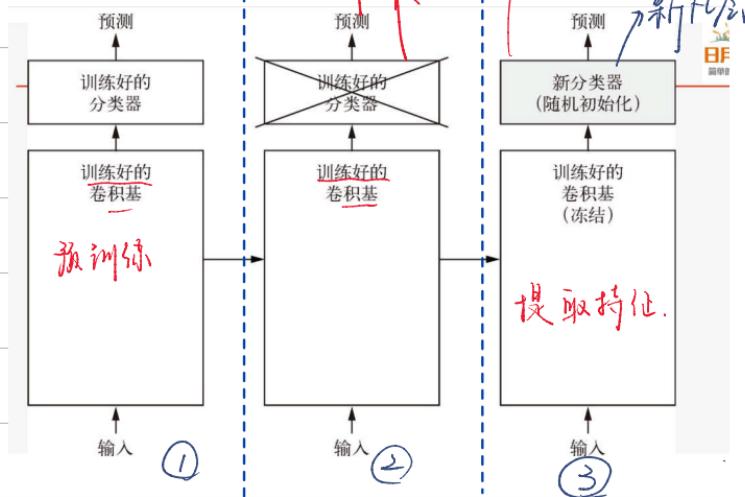
直接用 tf、keras。

Fc。
通过提取特征。

特征进行分类。

网络架构：

3步。



卷积：提取特征

全连接：综合(组合所有特征), 进行分类, 映射输出

日期： /

预训练微调

卷积基 $\begin{cases} \text{底层} & \rightarrow \text{调节, 特征} \\ \text{顶层} & \rightarrow \text{偏向分类任务} \end{cases}$

- 步骤
- ① 在预训练卷积基上添加自定义层 (FC)
 - ② 添加卷积基所有层 (conv)
 - ③ 训练分类层 (FC)
 - ④ 冻结顶部卷积基 (conv)
 - ⑤ 融合训练冻结层和分类层 (FC+conv)

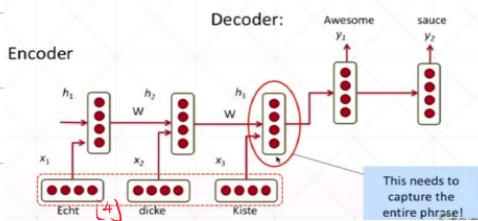
日期:

位置 及 时间 (语言, 文字)

时间序列表示方法

Sequence 序列

文本 → 数值



Sequence embedding 嵌入

b句, 10个词, 4个维度.

[b, seq_len, feature_len]

[b, 10, 4]

e.g. I like it.

[1000] 37词.

```
<tf.Tensor: id=5, shape=(3, 3), dtype=float32, numpy=
array([[[1., 0., 0.],
       [0., 1., 0.],
       [0., 0., 1.]], dtype=float32)>
```

时间箭头 (采样) ✓

[b, 28, 28] 像素 ✗

↑ feature 描述 ✓

6



b个曲线, t

[b, word num, word vec]

[word num, b, word vec] → 特征向量,

每-时间)采样, 得到b个点.

主要维度 = 时间,

[words, word vec]

日期:

- How to represent a word

- [Rome, Italy, ...]

- one-hot

label

语意识别

on: 80%

V 一个包点, word V

$$\begin{aligned} \text{Rome} &= [1, 0, 0, 0, 0, 0, 0, \dots, 0] \\ \text{Paris} &= [0, 1, 0, 0, 0, 0, 0, \dots, 0] \end{aligned}$$

$$\text{Italy} = [0, 0, 1, 0, 0, 0, 0, \dots, 0]$$

$$\text{France} = [0, 0, 0, 1, 0, 0, 0, \dots, 0]$$

label 故事.

high-dim 词汇量 10^4

(10^4, 1)

semantic similarity

trainable

语意相关性 "similarity."

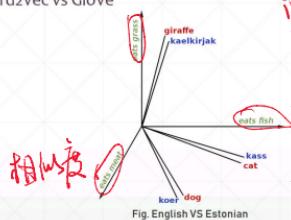
model.most_similar('king', topn=10)	model.most_similar('queen', topn=10)
(word, similarity with 'king')	(word, similarity with 'queen')
('Kings', 0.897245)	('Queens', 0.942618)
('baratheon', 0.809678)	('Joffrey', 0.933756)
('son', 0.763614)	('Margery', 0.931099)
('Robert', 0.708522)	('sister', 0.928902)
('lady', 0.698844)	('prince', 0.927364)
('Joffrey', 0.696455)	('Daenerys', 0.925207)
('prince', 0.695669)	('Varys', 0.918421)
('brother', 0.685239)	('need', 0.917492)
('Arya', 0.684527)	('melliferous', 0.915403)
('stannis', 0.682932)	('robb', 0.915272)

index	Embedding
0	[1.2, 3.1]
1	[0.1, 4.2]
2	[1.0, 3.1]
3	[0.3, 2.1]
4	[2.2, 1.4]
5	[0.7, 1.7]
6	[4.1, 2.0]

1000 个语意 [1000]

形成的语意表达方法

▪ Word2Vec vs GloVe



- Random initialized embedding

```

input
In [8]: from tensorflow.keras import layers
In [9]: x=tf.range(5)
In [10]: x=tf.random.shuffle(x)
Out[10]: <tf.Tensor: id=9, shape=(5,), dtype=int32, numpy=array([0, 4, 1, 3, 2])>

In [9]: net=layers.Embedding(10, 4) -> 嵌入数
In [10]: net(x)
<tf.Tensor: id=25, shape=(5, 4), dtype=float32, numpy=
array([-0.03535435,  0.01710499,  0.024379 , -0.010867 ],
[-0.03977622,  0.01753286,  0.00805125,  0.00836002],
[-0.0119989,  0.01030685, -0.01133521,  0.02052242],
[ 0.02230989, -0.02186236,  0.01720804, -0.03888531],
[-0.04997355,  0.00911248,  0.01886252,  0.01570504]), dtype=float32)>

```

input [4] $\xrightarrow{\text{嵌入数}} [4, 4]$ \downarrow embedding

```

In [11]: net.trainable
Out[11]: True
In [12]: net.trainable_variables
[<tf.Variable 'embedding/embeddings:0' shape=(10, 4) dtype=float32, numpy=
array([-0.03535435,  0.01710499,  0.024379 , -0.010867 ],
[-0.03977622,  0.01753286,  0.00805125,  0.00836002],
[-0.0119989,  0.01030685, -0.01133521,  0.02052242],
[ 0.02230989, -0.02186236,  0.01720804, -0.03888531],
[-0.04997355,  0.00911248,  0.01886252,  0.01570504]),
dtype=float32)>]

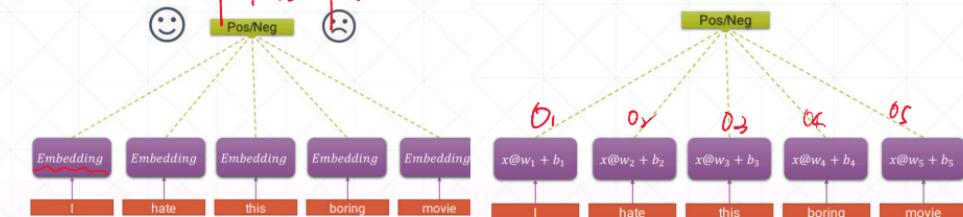
```

日期: /

循环神经网络 (RNN)

多个词，多个全连接层

时间 / 空间.



① Long sentence 的

- 100+ words
- too much parameters $[w_N, b_N]$

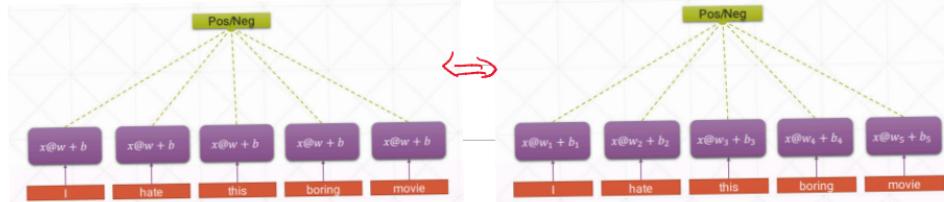
参数量过大
没有背景/语义相关性. 变孤独

Weight sharing.
利用同一张矩阵

② No context information

- consistent tensor

单独操作.



Consistent memory.

历史信息

加上之前的信息

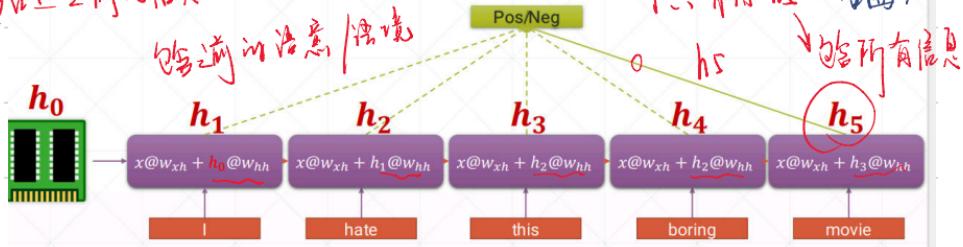


包含前的语义/语境

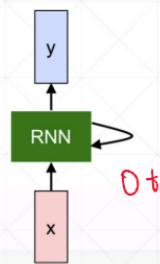
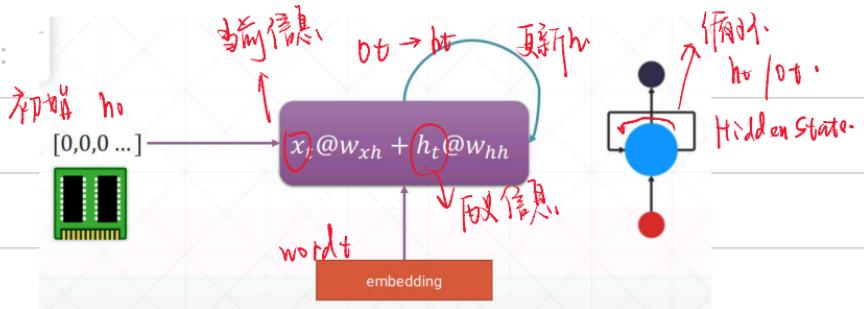
} 所有 h 相关.

} 只需用 h 1 层面.

→ 给所有信息.



日期:



$$h_t = f_W(h_{t-1}, x_t)$$

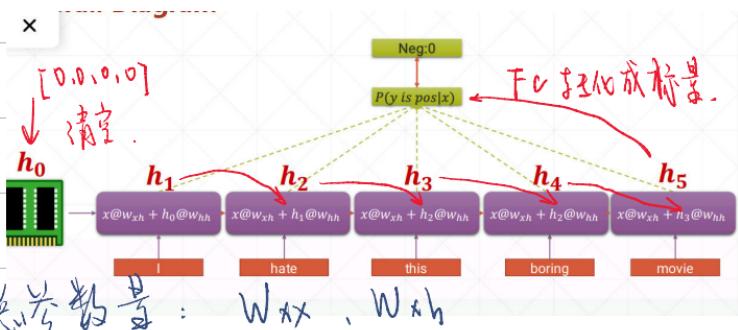
activation function

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

$$y_t = W_{hy}h_t + b.$$

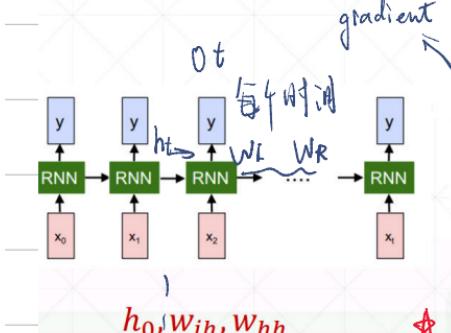
input

只用2个参数



参数数量: W_{xh} , W_{hh}

Now to train:



$$h_t = \tanh(W_I x_t + W_R h_{t-1})$$

$$y_t = \overline{W_O} h_t + b$$

$$\text{Loss: } \text{MSE}(y_{\tau, t}, t)$$

$$\frac{\partial E_t}{\partial W_R} = \sum_{i=0}^t \frac{\partial E_t}{\partial y_t} \frac{\partial y_t}{\partial h_t} \frac{\partial h_t}{\partial h_i} \frac{\partial h_i}{\partial W_R}$$

$$\frac{\partial h_t}{\partial h_i} = \underbrace{\frac{\partial h_t}{\partial h_{t-1}}}_{\frac{\partial h_{t-1}}{\partial h_{t-2}}} \frac{\partial h_{t-1}}{\partial h_{t-2}} \cdots \frac{\partial h_{i+1}}{\partial h_i} = \prod_{k=i}^{t-1} \frac{\partial h_{k+1}}{\partial h_k}$$

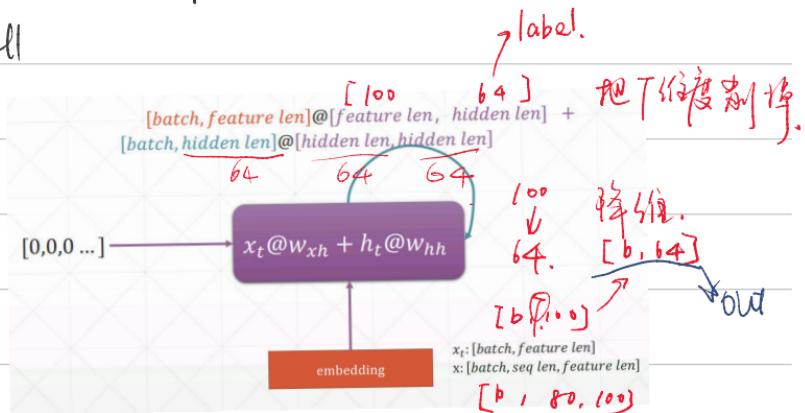
$$\frac{\partial h_{k+1}}{\partial h_k} = \text{diag}(f'(W_I x_i + W_R h_{i-1})) W_R$$

$$\frac{\partial h_k}{\partial h_1} = \prod_{i=1}^k \text{diag}(f'(W_I x_i + W_R h_{i-1})) W_R$$

累乘

日期: / ANN layer 使用.

RNN cell

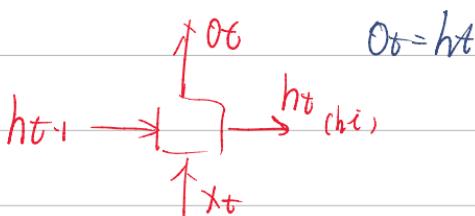


... $xW_{xh} + hW_{hh} \rightarrow$ Sample. RNN.

In [17]: $\text{cell} = \text{layers.SimpleRNNCell}(3)$ \rightarrow hidden
In [18]: $\text{cell}.\text{build}(\text{input_shape}=(\text{None}, 4))$ \rightarrow 80N(4維).

In [19]: $\text{cell}.\text{trainable_variables}$
 $\leftarrow \text{tf.Variable 'kernel:0' shape: (4, 3) dtype=float32, numpy=$
 $\text{array}([[0.23682523, 0.24167228, 0.19834113},$
 $[0.58464265, -0.44347632, -0.23693317],$
 $[0.5130104, -0.86219984, 0.16108215],$
 $[-0.37421566, -0.6311711, 0.46914995]], \text{dtype}=\text{float32})$,
 $\text{<tf.Variable 'recurrent_kernel:0' shape: (3, 3) dtype=float32, numpy=}$
 $\text{array}([[0.7126031, 0.39248434, 0.5815092],$
 $[-0.34029973, 0.9182056, -0.2027184],$
 $[0.61350864, 0.05342963, -0.7878784]], \text{dtype}=\text{float32})$,
 $\text{<tf.Variable 'bias:0' shape: (3,) dtype=float32, numpy=array([0., 0., 0.], \text{dtype}=\text{float32})}]$

- 前向: input $\uparrow t-1$
- $\text{out}, h_1 = \text{call}(x, h_0)$
- $x: [b, \text{seq len}, \text{word vec}]$
 $\uparrow 80 \quad 100$
- $h_0/h_1: [b, h \text{ dim}]$
- $\text{out}: [b, h \text{ dim}]$



Out [b, units]

日期: /

hidden_len.

```

state0 = [tf.zeros([batchsz, units])] → Memory.
state1 = [tf.zeros([batchsz, units])]
for word in tf.unstack(x, axis=1): # word: [b, 100]
    # ht = xtwxh+0*wth
    # out0: [b, 64] ← 更新 ↓ xt ↓ ht-1
    out0, state0 = self.rnn_cell0(word, state0, training)
    # out1: [b, 64] ↓ 2个cell
    out1, state1 = self.rnn_cell1(out0, state1, training)

```

固定维度数量.

```

self.rnn = keras.Sequential([
    layers.SimpleRNN(units, dropout=0.5, return_sequences=True, unroll=True),
    layers.SimpleRNN(units, dropout=0.5, unroll=True)
])
# x: [b, 80, 100] => [b, 64] t=80 且 no 断句   ht0=80
x = self.rnn(x)

```

返回状态

情感分类实战

定义: unit=64, epoch=4

加载 IMDB (电影评论) 数据集.

▪ SimpleRNNCell

- single layer
- multi-layers

▪ RNNCell

构造网络 class My RNN (keras.Model)

def __init__(self), [b, 80) → [b, 80, 100]

self.embedding = layers.embedding(total_words,

embedding_len, input_length=max_review_len.

100

80

日期: / [b, 80, 100] h-dim=64

定義 cell

self.rnn_cell0 = layers.SimpleRNNCell(units, dropout=0.5).

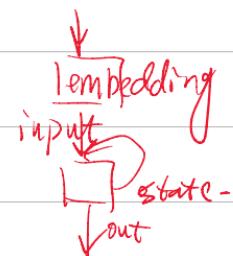
self.rnn_cell1 = layers.SimpleRNNCell(units, dropout=0.5).

fc: [b, 80, 100] → [b, 64] → [b, 1] 每一時刻 FC. 防止過拟合

self.outlayer = layers.Dense(1).

def.call(self.inputs, training=None):

net(x), net(x, training=True.)



input [b, 80]

x=input [b, 80]

x=self.embedding(x). [b, 64]

for word in tf.unstack(x, axis=1): (每个word)

[b, 64] 在时间上展开

initial state.

outstate = self.rnn_cell0(word, state0) outD=[b, 64]

state0 = state1

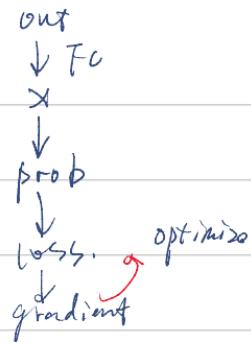
out1=[b, 64]

outstate = self.rnn_cell1(word, state1)

日期: /

$x = \text{soft_outlayer}(\text{out})$

$\text{prob} = f \cdot \text{sigmoid}(x)$



Main:

unit = 14 (hidden_dim) epoch = 4

model = My RNN (units) → 调用类(模型)

model.compile(optimizer = ...)

交叉验证: train, test loss =

epoch. loss - metrics = ...)

model.fit(db_train, epochs = epochs, validation_data=

model.evaluate(db_train))

日期: /

梯度弥散与梯度爆炸.

无限等价? 无限展开?

Nothing is straight forward.

$$h_t = \tanh(W_I x_t + W_R h_{t-1})$$

$$y_t = W_O h_t$$

$$\frac{\partial E_t}{\partial W_R} = \sum_{i=0}^t \frac{\partial E_t}{\partial y_i} \frac{\partial y_i}{\partial h_t} \frac{\partial h_t}{\partial h_i} \frac{\partial h_i}{\partial W_R}$$

$$h_t \rightarrow h_i \quad \frac{\partial h_t}{\partial h_i} = \frac{\partial h_t}{\partial h_{t-1}} \frac{\partial h_{t-1}}{\partial h_{t-2}} \dots \frac{\partial h_{i+1}}{\partial h_i} = \prod_{k=i}^{t-1} \frac{\partial h_{k+1}}{\partial h_k}$$

$$\frac{\partial h_{k+1}}{\partial h_k} = \text{diag}(f'(W_I x_t + W_R h_{i-1})) W_R \quad W_R \cdot N \rightarrow \text{gradient}$$

$$\frac{\partial h_k}{\partial h_1} = \prod_i^k \text{diag}(f'(W_I x_t + W_R h_{i-1})) W_R \quad W_h^k$$

当 W_h 越大时, 越容易 $W_h^k \rightarrow \infty$

当 W_h 越小时, 越不容易 $W_h^k \rightarrow 0$

查看 W_h gradient

dipping (对 gradient), 大于某一个阈值时.

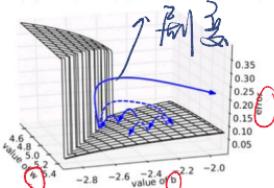
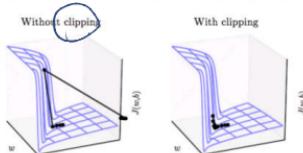


Figure 6. We plot the error surface of a single hidden unit recurrent network, highlighting the existence of high curvature walls. The solid lines depicts standard trajectories that gradient descent might follow. Using dashed arrow the diagram shows what would happen if the gradients is rescaled to a fixed size when its norm is above a threshold.

Algorithm 1 Pseudo-code for norm clipping

```
g ← ∂E/∂θ
if ||g|| ≥ threshold then
    g ← threshold * g / ||g||
end if
```

$$\hat{g} = \frac{g}{\|g\|} \cdot \text{threshold}$$



— Goodfellow et al., Deep Learning

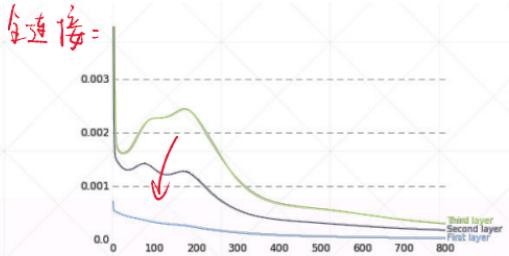
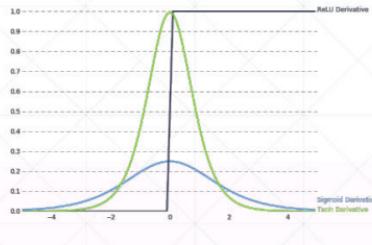
日期:

```
with tf.GradientTape() as tape:  
    logits = model(x)  
    loss = criteon(y, logits)  
  
grads = tape.gradient(loss, model.trainable_variables)  
# MUST clip gradient here or it will diverge!  
  
optimizer.apply_gradients(zip(grads, model.trainable_variables))
```

$\text{grads} = [\text{tf.clip_by_norm}(g, 5) \text{ for } g \text{ in grads}]$
clipping threshold 修剪值 -

反向传播，累计误差 $8^k \rightarrow \frac{\partial E}{\partial w}$

Gradient Vanishing



LSTM.

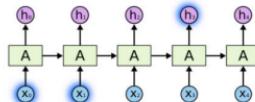
解决 gradient vanishing

长时记忆

The problem of long-term dependencies

(Vanilla) RNNs connect previous information to present task:

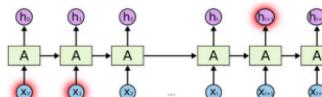
- enough for predicting the next word for “the clouds are in the *sky*”



short - term - memory.

- may not be enough when more context is needed

“I grew up in France... I speak fluent *French*. ”



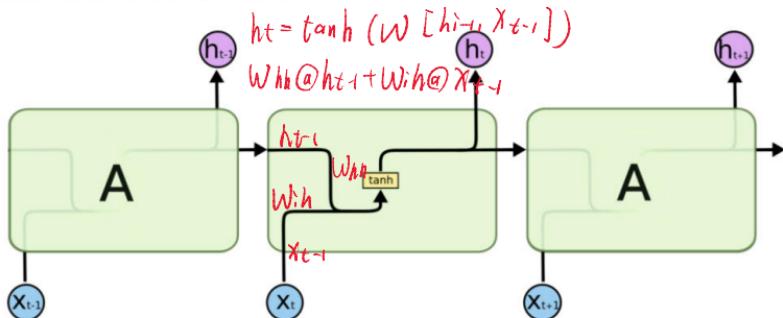
Long - term

Adapted from Christopher Olah

68

RNN:

All recurrent neural networks have the form of a chain of repeating modules of neural network



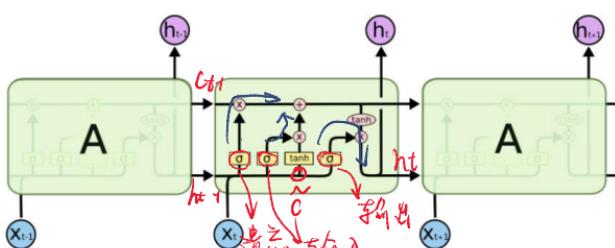
The repeating module in a standard RNN contains a single layer.

LSTM =

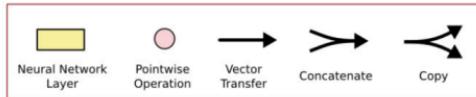
日期: /

有用的过滤器, “j” : σ (sigmoid.)
(遗忘)

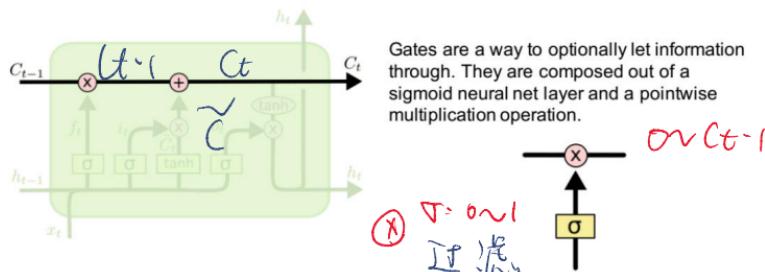
LSTMs also have this chain like structure, but the repeating module has a different structure. Instead of having a single neural network layer there are four, interacting in a very special way.



$$\sigma(w[h, x]) \downarrow \stackrel{b \sim 1}{\longrightarrow} h \rightarrow h'$$



- The Core Idea Behind LSTMs : Cell State



An LSTM has three of these gates, to protect and control the cell state.

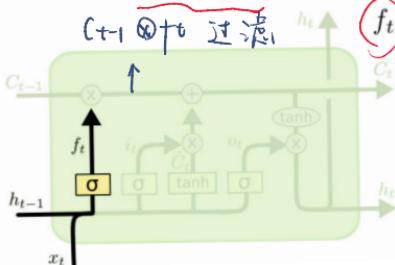
日期:

/

$t=1$ 全部记住 $t=0$ 忘记

Remember

LSTM : Forget gate



f_t 值:

控制量

输出值

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

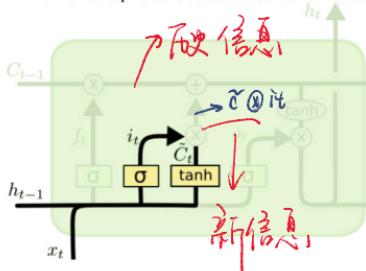
C_t It looks at h_{t-1} and x_t and outputs a number between 0 and 1 for each number in the cell state C_{t-1} .

A 1 represents "completely keep this" while a 0 represents "completely get rid of this".

W_c, W_i, W_o 自由参数
 b_c, b_i, b_o 偏差

LSTM : Input gate and Cell State

The next step is to decide what new information we're going to store in the cell state.



一个 sigmoid 层叫“**input gate layer**”决定哪些值我们将来会更新。

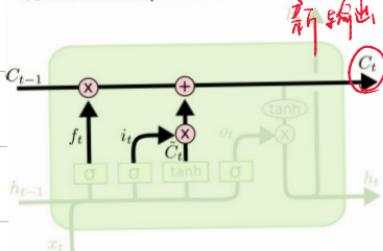
$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

一个 tanh 层创建一个新候选值向量，可能被添加到状态中。

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

LSTM : Input gate and Cell State

It's now time to update the old cell state into the new cell state.



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

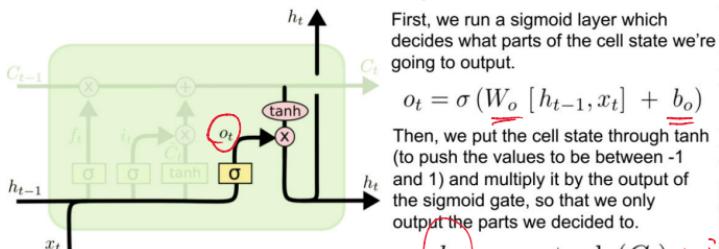
We multiply the old state by f_t forgetting the things we decided to forget earlier.

Then, we add the new candidate values, scaled by how much we decided to update each state value.

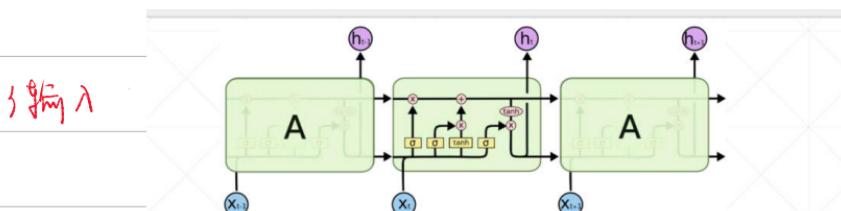
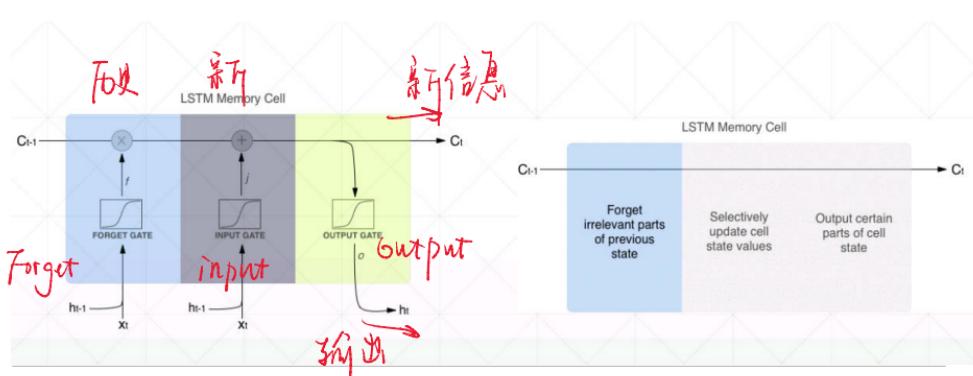
日期: /

LSTM : Output

Finally, we need to decide what we're going to output.



$C_t \rightarrow h_t$ (↑ 表示输出) 处理后是输出。



$$\begin{pmatrix} \mathbf{i}^{(t)} \\ \mathbf{f}^{(t)} \\ \mathbf{o}^{(t)} \\ \tilde{\mathbf{C}} \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} \mathbf{W} \begin{pmatrix} \mathbf{x}^{(t)} \\ \mathbf{h}^{(t-1)} \end{pmatrix} \quad (6)$$

$$\underline{\mathbf{c}^{(t)}} = \underline{\mathbf{f}^{(t)}} \circ \underline{\mathbf{c}^{(t-1)}} + \underline{\mathbf{i}^{(t)}} \circ \underline{\tilde{\mathbf{C}}} \quad (7)$$

$$\underline{\mathbf{h}^{(t)}} = \underline{\mathbf{o}^{(t)}} \circ \tanh(\underline{\mathbf{c}^{(t)}}) \quad (8)$$

过失 + 现在。

Out $\otimes \tanh(c^{(t)})$ (8)

记忆 + 短期输入

真值表:

input gate	forget gate	behavior
0	1	remember the previous value
1	1	add to the previous value
0	0	erase the value
1	0	overwrite the value

忽略当前信息

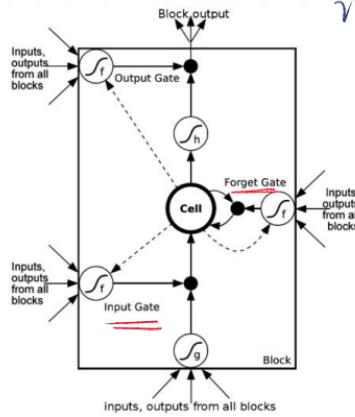
重写

删除之前信息

新信息

添加到

记忆中



How to solve Gradient Vanishing?

$$\frac{\partial h^k}{\partial h^i} = N \cdot \underbrace{W_h}_{\text{梯度向量}} \quad \text{从 } h \rightarrow C$$

$$\frac{\partial C_t}{\partial C_{t-1}} = \frac{\partial C_t}{\partial f_t} \frac{\partial f_t}{\partial h_{t-1}} \frac{\partial h_{t-1}}{\partial C_{t-1}} + \frac{\partial C_t}{\partial i_t} \frac{\partial i_t}{\partial h_{t-1}} \frac{\partial h_{t-1}}{\partial C_{t-1}}$$

$$+ \frac{\partial C_t}{\partial \tilde{C}_t} \frac{\partial \tilde{C}_t}{\partial h_{t-1}} \frac{\partial h_{t-1}}{\partial C_{t-1}} + \frac{\partial C_t}{\partial C_{t-1}}$$

是吧!!!

$$\frac{\partial C_t}{\partial C_{t-1}} = C_{t-1} \sigma'(\cdot) W_f * o_{t-1} \tanh'(C_{t-1})$$

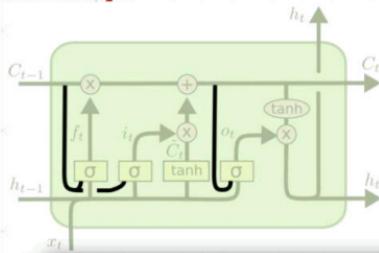
$$+ \tilde{C}_t \sigma'(\cdot) W_i * o_{t-1} \tanh'(C_{t-1})$$

$$+ i_t \tanh'(\cdot) W_C * o_{t-1} \tanh'(C_{t-1})$$

$$+ f_t$$

不可微梯度过大。

<https://zhuanlan.zhihu.com/p/30171115/Memory-Gradient.html>



$$\frac{\partial h_{k+1}}{\partial h_k} = \text{diag}(f'(W_I x_i + W_R h_{i-1})) W_R$$

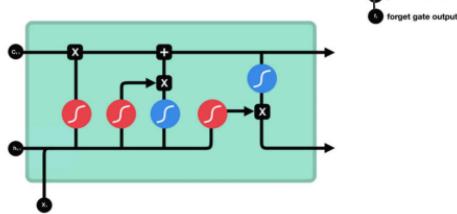
$$\frac{\partial h_k}{\partial h_1} = \prod_i^k \text{diag}(f'(W_I x_i + W_R h_{i-1})) W_R$$

日期: /

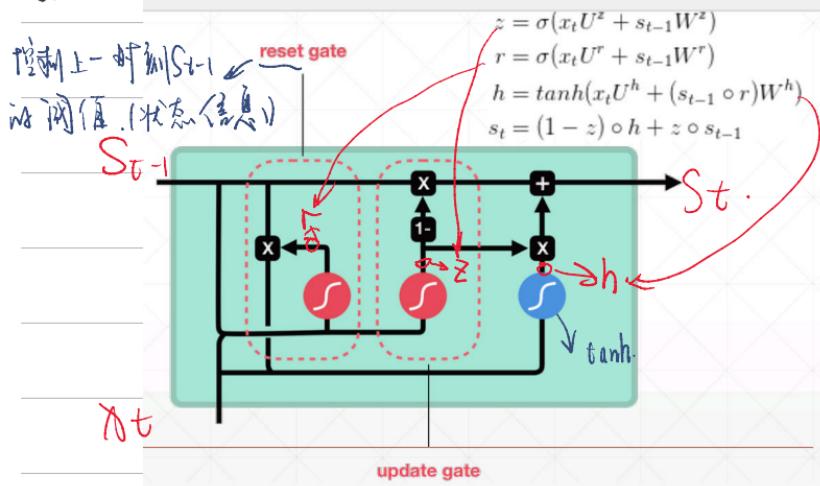
LSTM 实战

GRU: simpler, lower computation cost

LSTM,



GRU:



标注标志

RNN → GRU

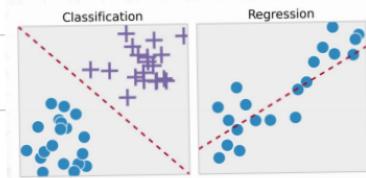
unroll = False.

self.rnn_cell = layer.GRUCell(units, dropout=0.5)

状态只用一个 self.state = [tf.zeros([batchsz, units])]
layers.LSTM layers.GRU.

日期：无监督学习.

Auto-Encoders.



"Label" 监督 通过人工标注完成
{
 源数据
 标注 } 进行 labeling

Massive Unlabel Data. Unsupervised Learning.

- Dimension reduction 降维.
- Preprocessing: Huge dimension, say 224x224, is hard to process

可视化.

- Taking advantages of unsupervised data

压缩 去噪 超分辨率.

- Compression, denoising, super-resolution ...

强化学习.

Pure Reinforcement Learning (cherry)

- The machine predicts a scalar reward given once in a while.
- A few bits for some samples



Supervised Learning (icing)

- The machine predicts a category or a few numbers for each input
- Predicting human-supplied data
- 10 → 10,000 bits per sample

Unsupervised/Predictive Learning (cake)

- The machine predicts anything its input for any observed part.
- Predicts future frames in videos

Millions of bits per sample

(Yes, I know, this picture is slightly offensive to RL folks. But I'll make it up)

特征提取

* 特征的生成层

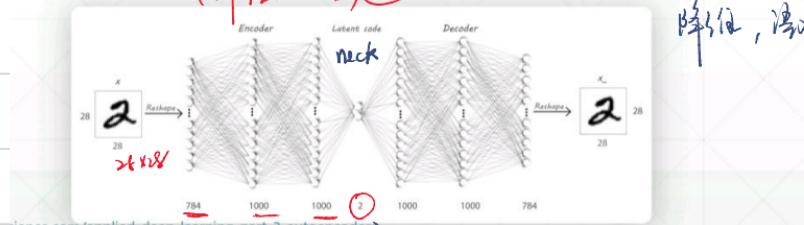
① 输入输出维度一致

② 有一个 code.

Auto-Encoder



reconstruct "Input" (神经结构调整) X



→ 24位 ~

降低，意义相等。

日期:

- Loss function for binary inputs $l(f(\mathbf{x})) = -\sum_k (x_k \log(\hat{x}_k) + (1 - x_k) \log(1 - \hat{x}_k))$
- Cross-entropy error function (reconstruction loss) $f(\mathbf{x}) \equiv \hat{\mathbf{x}}$

$$\begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

分类型

- Loss function for real-valued inputs

$$l(f(\mathbf{x})) = \frac{1}{2} \sum_k (\hat{x}_k - x_k)^2$$

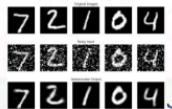
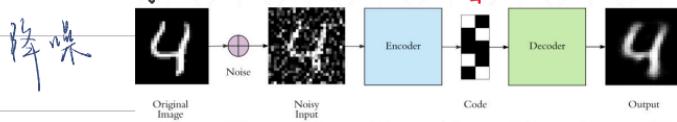
求重建损失的误差

- sum of squared differences (reconstruction loss)
- we use a linear activation function at the output

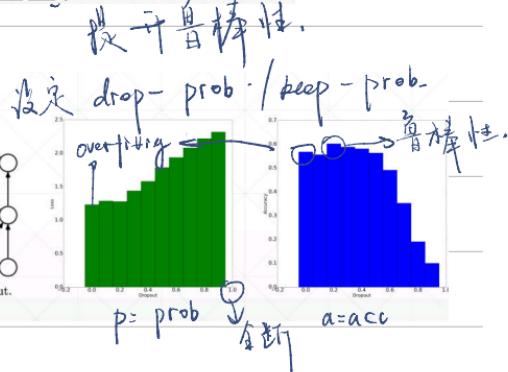
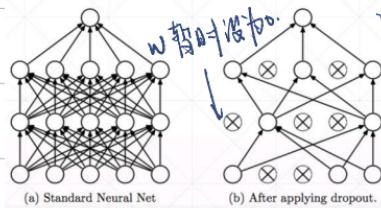
PCA vs. Auto-Encoder

- PCA, which finds the directions of maximal variance in high-dimensional data, select only those axes that have the largest variance.
- The linearity of PCA, however, places significant limitations on the kinds of feature dimensions that can be extracted.

Denoising Auto-Encoders.



Dropout Auto-Encoders.



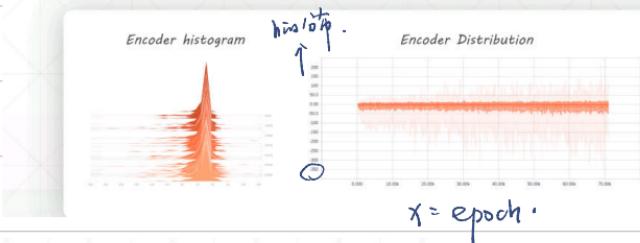
日期: /

Adversarial

- Distribution of hidden code

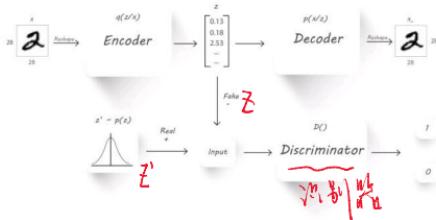
$x \rightarrow$ hidden vector $\rightarrow \bar{x}$

Auto Encoders.



- Give more details after GAN

希望 \bar{x} 属于预设的分布。



使 p, q 的分布相近
即 p, q 的分布相同。

目标: $x \rightarrow \bar{x} \rightarrow \bar{\bar{x}}$

Another Approach: $q(z) \rightarrow p(z)$

- Explicitly enforce

(AE)
重建误差 + 分布监督. 直接通过
 $I_i(\theta, \phi) = -E_{z \sim q_\theta(z|x_i)} [\log p_\phi(x_i|z)] + KL(q_\theta(z|x_i)||p(z))$

使逆向图

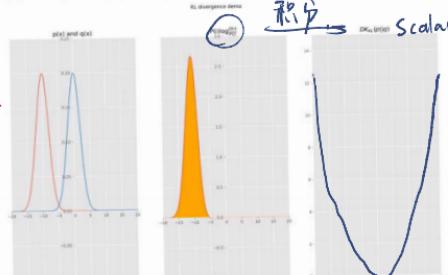
与后向图相似。

KL 效度.

两个分布 P, Q .

当 P, Q 的分布相近时
 $KL(P||Q)$ 变小

$[0, +\infty)$

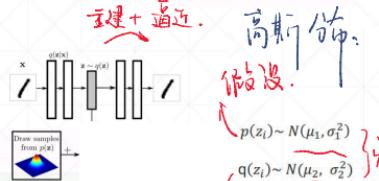


完全重叠时

$DKL(P||Q) = 0$

Minimize KL Divergence

- Evidence Lower BOund



$$KL(q_\theta(z|x_i) || p(z))$$

KL(p|q) 和 KL(q|p)

是大小相反的.

$$KL(p, q) = - \int p(x) \log q(x) dx + \int p(x) \log p(x) dx$$

$$= \frac{1}{2} \log(2\pi\sigma_2^2) + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2}(1 + \log 2\pi\sigma_1^2)$$

$$= \log \frac{\sigma_2}{\sigma_1} + \frac{\sigma_1^2 + (\mu_1 - \mu_2)^2}{2\sigma_2^2} - \frac{1}{2}$$

差异间
差异和均值

$x \rightarrow h \rightarrow \bar{x}$
 $q(h)$ simple \rightarrow 异样

假設

h 属于 g 分布.

$z \sim N(\mu_1, \sigma_1^2)$

Deterministic node
Random node

重建时进行
Sample.



decoder model

$$z = \mu + \sigma \odot \varepsilon$$

$$\varepsilon \sim N(0, 1)$$

$$z \sim N(\mu, \sigma^2)$$

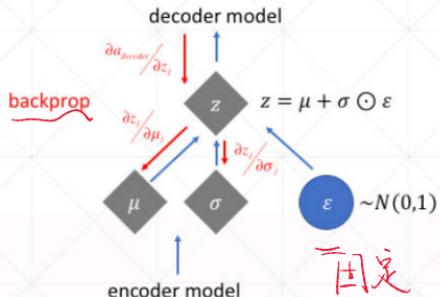
回溯进行
back prop

Reparameterization trick

将 $z \sim N(\mu, \sigma)$

表达为 $z = \mu + \sigma \cdot \varepsilon$.

$$\varepsilon \sim N(0, 1)$$



RL + $\beta \times KL(q||p)$

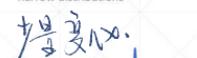
Penalizing reconstruction loss
encourages the distribution to
describe the input

Without regularization, our
network can "cheat" by learning
narrow distributions

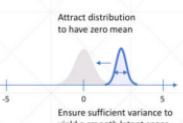
Penalizing KL divergence
acts as a regularizing force



Our distribution
deviates from the
prior to describe
some characteristic
of the data



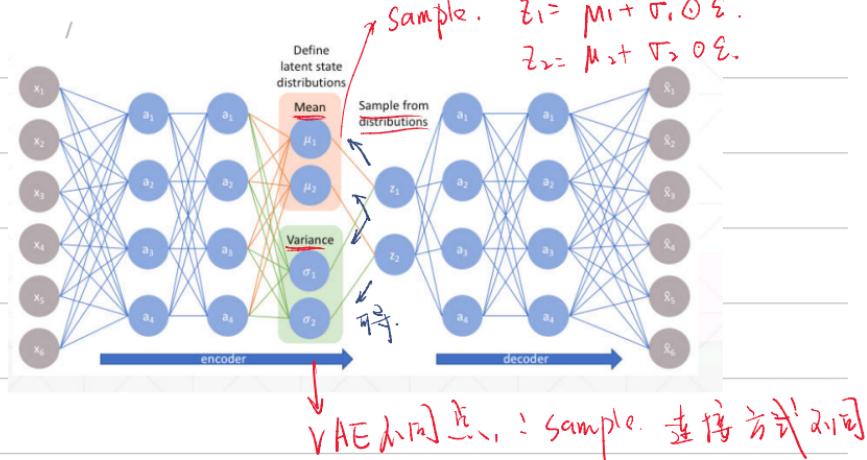
With a small
enough variance,
this distribution is
effectively only
representing a
single value



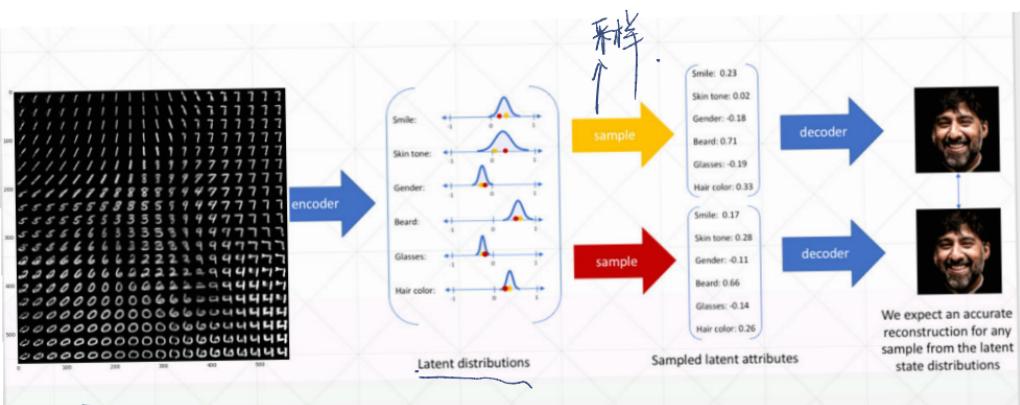
Attract distribution
to have zero mean

Ensure sufficient variance to
yield a smooth latent space

日期:



Variational autoencoders — VAE



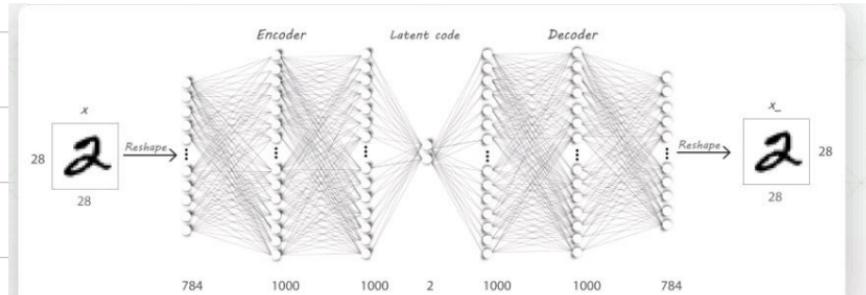
$$[h_0, h_1] \rightarrow$$

$$\mathcal{N}(\mu_1, \sigma_1^2)$$

$$\mathcal{N}(\mu_0, \sigma_0^2)$$

日期: /

VAE/AE 实战



- 定义参数: h-dim ($784 \rightarrow n$) batchsz, lr.
- 划分数据并拆分 (32) 训练集, 不用 y (label)
- 新建类 (模型) class AE(keras, model)
Encoder : $\text{se}[\text{t}].\text{encoder} = \text{sequential}([\text{layers.Dense}(256, \text{acti}=\text{tf.nn.relu}), \text{layers.Dense}(128, \text{acti}=\text{tf.nn.relu})])$
 $784 \rightarrow \text{h-dim}$ $784 \rightarrow 256$ $256 \rightarrow 128$ $128 \rightarrow \text{h-dim}$ $[\text{layers.Dense}(\text{h-dim})]$
- * $\text{h-dim} \rightarrow 784$ $128 \rightarrow \text{h-dim}$ $[\text{layers.Dense}(\text{h-dim})]$
- # Decoder : $\text{se}[\text{t}].\text{decoder} = \text{sequential}([\text{layers.Dense}(128, \text{acti}=\text{tf.nn.relu}), \text{layers.Dense}(256, \text{acti}=\text{tf.nn.relu}), \text{layers.Dense}(784, \text{acti}=\text{tf.nn.relu})])$
 $\text{h-dim} \rightarrow 128$ $128 \rightarrow 256$ $256 \rightarrow 784$ $[\text{layers.Dense}(784)]$

日期: /

前向传播. def call (self, inputs, training = None):

$[b, 784] \rightarrow [b, 10]$ $h = \text{self.encoder}(\text{inputs})$

$[b, 10] \rightarrow [b, 784]$ $x_{\hat{}} = \text{self.encoder}(h)$

return $x_{\hat{}}$

参数:

model, AE()

model.build (input_shape = [None, 784])

model.summary() optimizers.Adam (lr=lr)

for epoch in range (1):

for step, x in enumerate (train_db):

$x = \text{tf.reshape}$

with ... as tape:

$x - \text{rec_logits} = \text{model}(x)$

$\text{rec_loss} = \text{tf.losses.binary_crossentropy}(...)$

#求均值. $\text{rec_loss} = \text{tf.reduce_mean}(\text{rec_loss})$

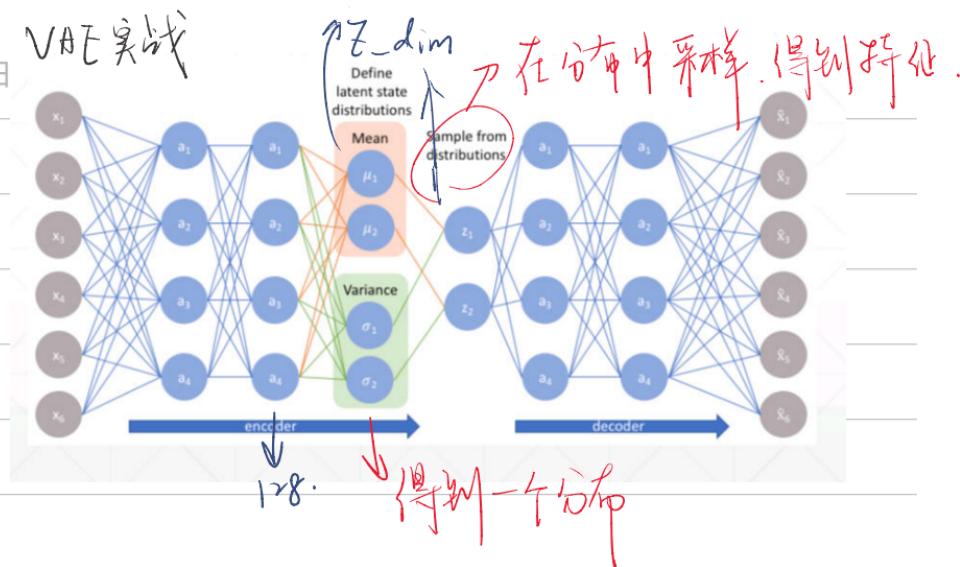
#计算梯度

$\text{ads} = \text{tape.gradient}(\text{loss}, \text{variable})$

#更新

$\text{optimizers.apply_gradient}(\text{ads}, \text{variable})$

VAE 实战



$$z_dim = 10 \quad \text{中间层维度}$$

```
class VAE(keras.Model):
```

~~def~~ ~~__init__(self):~~ # Encoder:

~~super(VAE, self).__init__()~~

$f_{c1} \rightarrow f_{c2} \rightarrow f_{c3}$ self.fc1 = layers.Dense(128)
M: self.fc2 = layers.Dense(z_dim)

α : self.fc3 = layers.Dense(z_dim)

Decoder: self.fc4 = layers.Dense(128)

[b, z] self.fc5 = layers.Dense(784.)

$\Rightarrow [b, 784]$

【】sequential?

日期:

/

$$x \rightarrow f_{c1} \leftarrow \begin{cases} f_{c2} \rightarrow mu \\ f_{c3} \rightarrow var \end{cases}$$

def Encoder (self, x):

mean: $h = tf.nn.relu(self.fc1(x))$
mu = self.fc2(h)

variance $\log-var = self.fc3(h)$.
return mu, logvar

def Decoder (self, x):

out = tf.nn.relu(fc4(z))

out = self.fc1(out).

return out

def reparameterize (self, mu, log_var)

z = ~~tf~~.^{s.}eps = tf.random.normal(inputs)

std = tf.exp(log_var) * 0.5

$z = mu + std * eps$.

$z = \mu + \sigma \odot \epsilon$.

return z.

日期: /

def call (self, input, training = None)

① mu, log-var = self.encoder (inputs)

reparameterize trick

② $\hat{z} = \text{self.reparameterize} (\mu, \log\text{-var})$

③ $\hat{x} = \text{self.decoder} (\hat{z})$

return, \hat{x} , mu, log-var

↓ 还原图.

\hat{x} , mu, log-var = model (*)

AEMN loss_rec_loss = tf.loss.binary_cross...

loss_rec_loss = tf.reduce_mean (rec_loss)

简化为 KL 损差. 符合正态分布

kl-div = $-0.5 * (\log\text{-var} + 1 - \mu^2) - \text{tf.exp}(\log\text{-var})$

kl-div = tf.reduce_sum (kl_dim) / x.shape[0].

loss = rec_loss + 1. * kl-div

两部分 直接误差. + KL ($x \star \star = x^2$)

样本 \hat{x} = | sample \hat{x} .
 \hat{x} (reconstructed) | ↗ 改好!

日期:

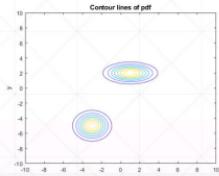
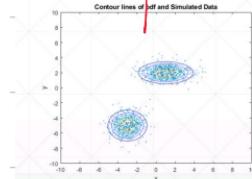
What I cannot create, I do not understand

对抗生成网络.

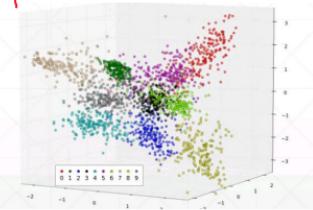
—Richard Feynman.

Our Goal: $p(x)$ 反观数据
猜.

```
13 mu = [1 2; 3 -5]
14 sigma = cat(3,[2 0; 0 .5],[1 0;0 1])
15 p = ones(1,2)/2
16 gm = gmdistribution(mu,sigma,p)
```



16步,
sample



Mnist. 降低. 量化.

画家的成长 → 学习优秀作品分布的过程. / 神经 / 鉴别.

同时成长 → 标准提高 → 优化 → 神经鉴别.
效果提高 ↗

Nash Equilibrium.

50% 混过关 · 50% 被鉴别

收获: 如何画一幅名作 / 成为一名画家.

Painter or Generator:

画家



Discriminator

鉴别师



→ 不被鉴别出

↓ 目标: 达到均衡

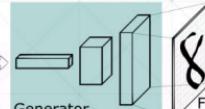
Training set

$P_r = P_g$.

$r = \text{real}$

$P_r(x)$

鉴别技术



Discriminator

Real 1

Fake 0

Random noise

Generator

Fake image

鉴别技术

真实 / 假造

生成 / 鉴别

生成技术

鉴别技术

真实 / 假造

$P_g(x) \rightarrow \text{模仿技巧.}$

日期: Falses. $\min_G \max_D L(D, G) = \mathbb{E}_{x \sim p_r(x)} [\log D(x)] + \mathbb{E}_{z \sim p_g(z)} [\log(1 - D(G(z)))]$ → 画出的产物
 Gene factor $\downarrow V(D, G) = \mathbb{E}_{x \sim p_r(x)} [\log D(x)] + \mathbb{E}_{x \sim p_g(x)} [\log(1 - D(x))]$ 等于 1.
 Discriminator 画出: 趋于 1.

对于 $V(DG)$ 磨削师试图使其变大
画家试图减小它。 \rightarrow 磨削为真。

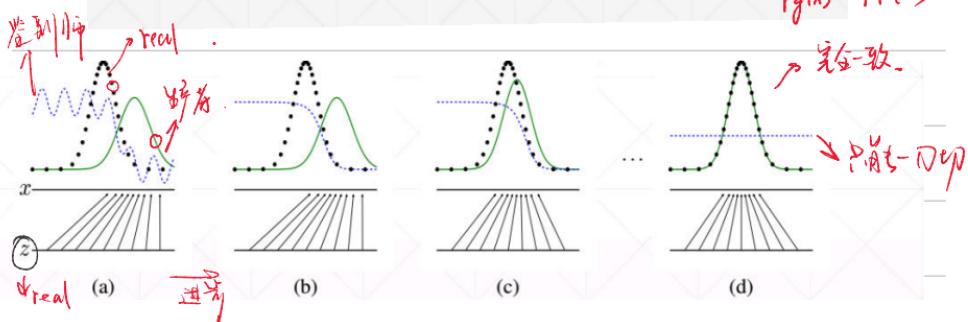
$$x \rightarrow D \rightarrow D(x) \qquad z \rightarrow G \rightarrow x_g$$

- <https://relinakano.github.io/gan-playground/>
 - <https://affinelayer.com/pixsrv/>
 - <https://www.youtube.com/watch?v=gReHvktowLY&feature=youtu.be>
 - <https://github.com/ajbrock/Neural-Photo-Editor>
 - <https://github.com/nashory/gans-awesome-applications>

- Q1. Where will D converge, given fixed G

- Q2. Where will G converge, after optimal D

$$\min_G \max_D L(D, G) = \mathbb{E}_{x \sim p_r(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))] \\ = \mathbb{E}_{x \sim p_r(x)} [\log D(x)] + \mathbb{E}_{x \sim p_g(x)} [\log(1 - D(x))]$$



日期: / Distributor 固定 Generator
 On here will D go (fixed G). * $D^* = \frac{\Pr}{\Pr + Pg}$

Proposition 1. For G fixed, the optimal discriminator D is

$$D_G^*(x) = \frac{p_{data}(x)}{p_{data}(x) + p_g(x)} \quad (2)$$

Proof. The training criterion for the discriminator D, given any generator G, is to maximize the quantity $V(G, D)$

$$\begin{aligned} V(G, D) &= \int_x p_{data}(x) \log(D(x)) dx + \int_z p_z(z) \log(1 - D(g(z))) dz \\ &= \int_x p_{data}(x) \log(D(x)) + p_g(x) \log(1 - D(x)) dx \end{aligned}$$

constant $\int_z p_z(z) dz = 1$ $P(z) = P(g)$.

$p_{data}(x)$
真实值.

Thus, set $\frac{df(\tilde{x})}{d\tilde{x}} = 0$, we get the best value of the discriminator:

$$D^*(x) = \tilde{x}^* = \frac{A}{A+B} = \frac{p_r(x)}{p_r(x) + p_g(x)} \in [0, 1]$$

最优解. G 固定 $\rightarrow D$ 成为 $\frac{P_r}{P_r + Pg}$.

$$\begin{aligned} f(\tilde{x}) &= A \log \tilde{x} + B \log(1 - \tilde{x}) \\ \frac{df(\tilde{x})}{d\tilde{x}} &= A \frac{1}{\ln 10} \frac{1}{\tilde{x}} - B \frac{1}{\ln 10} \frac{1}{1 - \tilde{x}} \\ &= \frac{1}{\ln 10} \left(\frac{A}{\tilde{x}} - \frac{B}{1 - \tilde{x}} \right) \\ &= \frac{1}{\ln 10} \frac{A - (A + B)\tilde{x}}{\tilde{x}(1 - \tilde{x})} = 0 \end{aligned}$$

$$A = (A + B) D(\tilde{x}),$$

$$D(\tilde{x}) = \frac{A}{A + B}.$$

下图
从左到右更接近

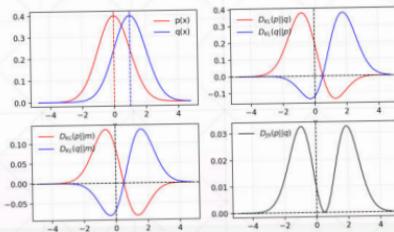
实机对称 KL divergence

差异
概率

$$D_{KL}(p \parallel q) = \int_x p(x) \log \frac{p(x)}{q(x)} dx$$

$$JS = KL(-\frac{1}{2} + \frac{1}{2})$$

$$D_{JS}(p \parallel q) = \frac{1}{2} D_{KL}(p \parallel \frac{p+q}{2}) + \frac{1}{2} D_{KL}(q \parallel \frac{p+q}{2})$$



定义 KL Where will G go (After D^*).

$$\begin{aligned} D_{JS}(p_r \parallel p_g) &= \frac{1}{2} D_{KL}(p_r \parallel \frac{p_r + p_g}{2}) + \frac{1}{2} D_{KL}(p_g \parallel \frac{p_r + p_g}{2}) \\ &= \frac{1}{2} \left(\log 2 + \int_x p_r(x) \log \frac{p_r(x)}{p_r(x) + p_g(x)} dx \right) + \\ &\quad \frac{1}{2} \left(\log 2 + \int_x p_g(x) \log \frac{p_g(x)}{p_r(x) + p_g(x)} dx \right) \\ &= \frac{1}{2} \left(\log 4 + L(G, D^*) \right) \end{aligned}$$

展开

$$D_{JS}(p_r \parallel p_g) \geq 0$$

$$p_r = p_g$$

固定训练样本在 D^*

$$\min_{G, D} \max_{D^*} V(G, D^*)$$

$$L(G, D^*) = 2D_{JS}(p_r \parallel p_g) - \log 2 \stackrel{!}{=} \log 4$$

固定

固定 D^* . $\min V(G, D^*) = L = 2 \log 4$.

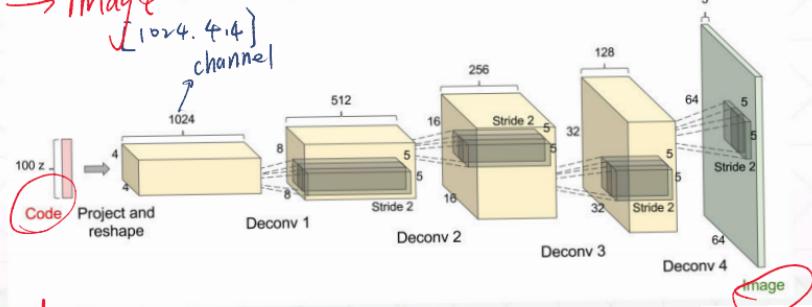
前提: $P_r = P_g$

$G, D \xrightarrow{\text{固定}} D^* \xrightarrow{\text{固定}} L(\min) = \log 4$.
 最优

日期: /

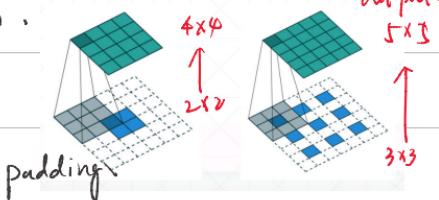
DC GAN (Deconvolution GAN).

Code → Image



Transposed

Deconvolution .



升维，
维度变大。

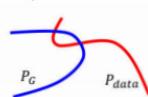
* Training Stability.

In most cases, P_G and P_{data} are not overlapped.

1. The nature of data

Both P_{data} and P_G are low-dim manifold in high-dim space.

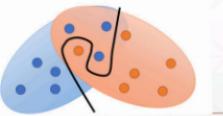
The overlap can be ignored.



2. Sampling

Even though P_{data} and P_G have overlap.

If you do not have enough sampling



不稳定.

GAN 稳定性 . Pretrain $P_{generator}$

小可 能有 overlap 问题 .

当 P_G 和 P_{data} 小重合时 .

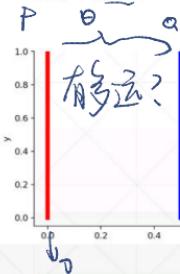
KL/JSD 无法正常运行 .

Toy example

极端情况：分布在线上

$$\forall (x, y) \in P, x = 0 \text{ and } y \sim U(0, 1)$$

$$\forall (x, y) \in Q, x = \theta, 0 \leq \theta \leq 1 \text{ and } y \sim U(0, 1)$$



+
KL不转

很好的衡量
两个分布的距离

$$D_{KL}(p\|q) = \int_x p(x) \log \frac{p(x)}{q(x)} dx$$

$$D_{JS}(p\|q) = \frac{1}{2} D_{KL}(p\|\frac{p+q}{2}) + \frac{1}{2} D_{KL}(q\|\frac{p+q}{2})$$

只要 $\theta \neq 0$:

小重合， $\int p \log \frac{p}{q}$

对于 $P \sim \delta_0$

$y \sim U(0, 1)$

$$D_{KL}(P\|Q) = \sum_{x=0, y \sim U(0,1)} 1 \cdot \log \frac{1}{0} = +\infty$$

$$D_{KL}(Q\|P) = \sum_{x=\theta, y \sim U(0,1)} 1 \cdot \log \frac{1}{\theta} = +\infty$$

$$D_{JS}(P, Q) = \frac{1}{2} \left(\sum_{x=0, y \sim U(0,1)} 1 \cdot \log \frac{1}{1/2} + \sum_{x=\theta, y \sim U(0,1)} 1 \cdot \log \frac{1}{1/2} \right) = \log 2$$

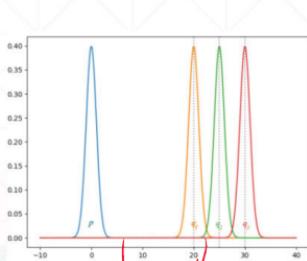
$$W(P, Q) = |\theta|$$

But when $\theta = 0$, two distributions are fully overlapped:

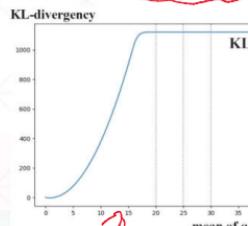
$$D_{KL}(P\|Q) = D_{KL}(Q\|P) = D_{JS}(P, Q) = 0$$

$$W(P, Q) = 0 = |\theta|$$

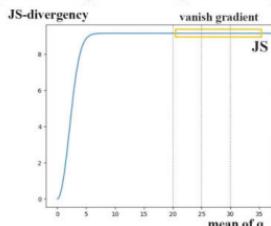
KL有一定的局限性。



MLE is kind of minimize KLD



GAN

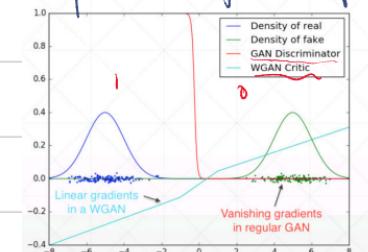


大面积都相似, KL差不多, 都很大.

梯度没了!!! ↓ gradient vanish!

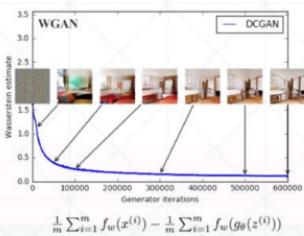
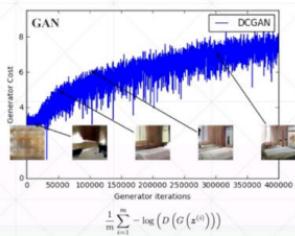
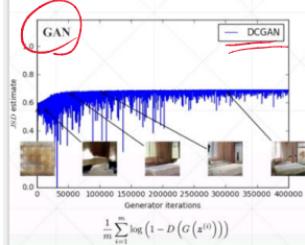
没有更新

所以开始时, P 和 Q 距离远, 找不到训练方向.

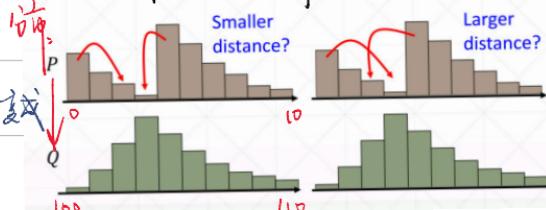


在不相关区域 GAN 梯度无

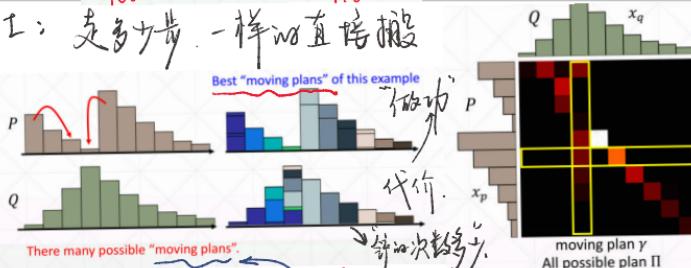
WGAN 还有,



JS / KL 都入能引導



钟工：走多步一样没直接搬



怎么钟最省力，时间关系

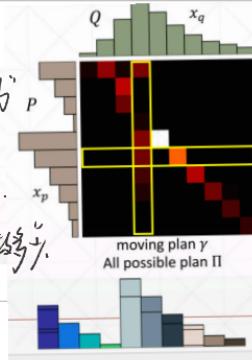
A "moving plan" is a matrix
The value of the element is the amount of earth from one position to another.

$$\text{Average distance of a plan } \gamma: B(\gamma) = \sum_{x_p, x_q} \gamma(x_p, x_q) \|x_p - x_q\|$$

Earth Mover's Distance:

$$W(P, Q) = \min_{\gamma \in \Pi} B(\gamma)$$

The best plan



Wasserstein Distance.

inf overlap.

$$W(\mathbb{P}_r, \mathbb{P}_g) = \inf_{\gamma \in \Pi(\mathbb{P}_r, \mathbb{P}_g)} \mathbb{E}_{(x, y) \sim \gamma} [\|x - y\|],$$

$$W \rightarrow [f(x^{(i)}) - f(G(z^{(i)}))]$$

新映射 \downarrow Discriminator/Critic

Generator

$$\text{GAN} \quad D \quad \nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m [\log D(x^{(i)}) + \log (1 - D(G(z^{(i)})))]$$

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m -\log (D(G(z^{(i)})))$$

$$\text{WGAN} \quad f \quad \nabla_{\theta_f} \frac{1}{m} \sum_{i=1}^m [f(x^{(i)}) - f(G(z^{(i)}))]$$

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m -f(G(z^{(i)}))$$

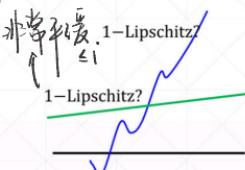
用神经网络法逼近
由 Wasserstein Distance
约束。在上往取两点 x_1, x_2 , 梯度小于 1.
学习得来

threshold.

Weight Clipping [Martin Arjovsky, et al., arXiv, 2017]

Force the parameters w between c and $-c$.
After parameter update, if $w > c$, $w = c$;
if $w < -c$, $w = -c$

How to fulfill this constraint?

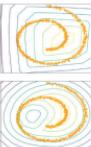
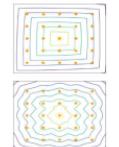
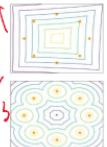


Sort of regularization

8 Gaussians

25 Gaussians

Swiss Roll



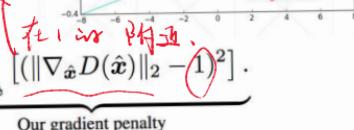
WGAN-Gradient Penalty

$$\hat{x} \in [x_r, x_g]$$

$$|f(x_1) - f(x_2)| \leq |x_1 - x_2|.$$

$$1 - L_F$$

$$\frac{\partial D}{\partial x}$$



$$L = \underbrace{\mathbb{E}_{\tilde{x} \sim \mathbb{P}_g} [D(\tilde{x})] - \mathbb{E}_{x \sim \mathbb{P}_r} [D(x)]}_{\text{Original critic loss}} + \lambda \underbrace{\mathbb{E}_{\hat{x} \sim \mathbb{P}_{\hat{x}}} [(\|\nabla_{\hat{x}} D(\hat{x})\|_2 - 1)^2]}_{\text{Our gradient penalty}}.$$

隨便 sampling

做线性差值。

where \hat{x} sampled from \tilde{x} and x with t uniformly sampled between 0 and 1

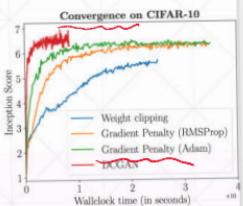
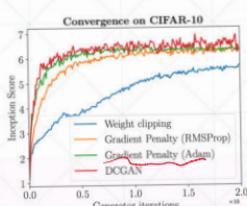
$$\hat{x} = t\tilde{x} + (1-t)x \text{ with } 0 \leq t \leq 1$$

$$\sum \frac{\partial D}{\partial x}, \text{ avg}$$

DCGAN: 没计参数。

WGAN + 增加复数。 $\frac{\partial P}{\partial x}$

DCGAN	LSGAN	WGAN (clipping)	WGAN-GP (ours)
Baseline (G : DCGAN, D : DCGAN)			
G : No BN and a constant number of filters, D : DCGAN			
G : 4-layer 512-dim ReLU MLP, D : DCGAN			
No normalization in either G or D			
Gated multiplicative nonlinearities everywhere in G and D			
tanh nonlinearities everywhere in G and D			
101-layer ResNet G and D			



日期:

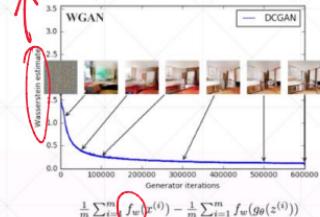
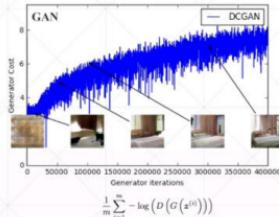
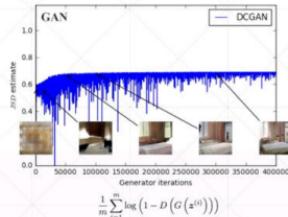
Unsupervised

Method	Score
ALI [8] (in [27])	5.34 ± .05
BEGAN [4]	5.62
DCGAN [22] (in [11])	6.16 ± .07
Improved GAN (-L+HA) [23]	6.86 ± .06
EGAN-Ent-VI [7]	7.07 ± .10
DFM [27]	7.72 ± .13
WGAN-GP ResNet (ours)	7.86 ± .07

Supervised

Method	Score
SteinGAN [26]	6.35
DCGAN (with labels, in [26])	6.58
Improved GAN [23]	8.09 ± .07
AC-GAN [20]	8.25 ± .07
SGAN-no-joint [11]	8.37 ± .08
WGAN-GP ResNet (ours)	8.42 ± .10
SGAN [11]	8.59 ± .12

可以反映距离高。



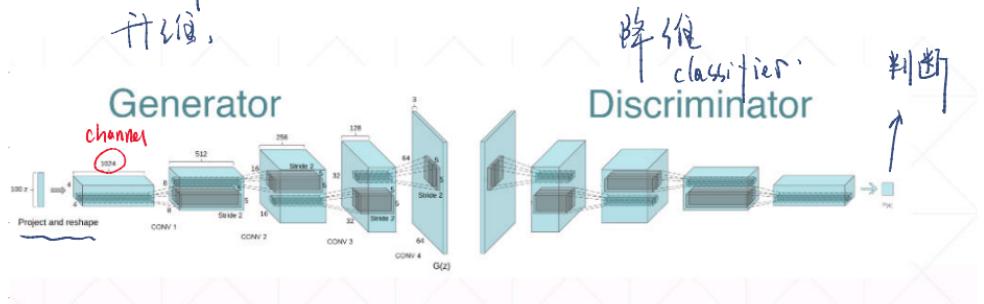
日期: /

GAN 实战

数据集: Anime [b 96x96. 3]

上采样

升维,



采用卷积 convolution 升维。

定义 Generator:

class Generator

定义类的 __init__

channel 降维, size +

$z = [b, 100] \rightarrow [b, 3 \times 3 \times 512] \rightarrow [b, 3 \times 3 \times 512] \rightarrow [b, 64, 64, 3]$

反向卷积, 升维、

self.fc = layers.Dense(3 * 3 * 512)

strides,

self.conv1 = layers.Conv2DTranspose(16, 3, 2)

self.bn1

self.conv2 = (16, 5, 2)

channel

conv3 = (3, 4, 3)

日期: / 參數 channel = $512 \rightarrow 256 \rightarrow 128 \rightarrow 3$
 size : 調節 strides! 使 $8 \times 8 \times 3$ 变成 $64 \times 64 \times 3$.

```

def call:
    x = self.fc(input) [2, 100] → [2, 3 * 3 * 512]
    x = tf.reshape(x, [-1, 3, 3, 512]) 3 * 3 * 512 → [3, 3, 512]
    x = tf.nn.relu(bn1(conv1(x)))
    x = tf.nn.tanh(bn2(conv2(x)))
    x = self.conv3.
    x = tf.tanh(x) [-1, 1] 激活函数.
    return x. size [b, 64, 64, 3]

```

識別器

Discriminator, class Discriminator:

#import ix def __init__:

[b, 64, 64, 3] → [b, 1] → 0 / 1

3个(conv + bn) + 1 FC.

channel 增加, size 降低.

conv1 = (64, 5, 3) conv2 = (128, 5, 3)

conv3 = (256, 5, 3)

[b, h, w, 3] → [b, -1] tf.layers.Flatten()

self.fc = layers.Dense(1)

日期： /

甘前向传播， def call.

$x = \text{tf}.nn.\text{leaky_relu}(\text{self}.conv1(inputs))$

$x = \text{tf}...(\text{self}.bn2(\text{self}.conv2(x)))$

$x = \text{tf}...(\text{bn3}(\text{conv3}(x)))$

$x = \text{self}.flatten.$

$\text{logits} = \text{self}.fc1(x) [b, -1] \rightarrow [b, 1]$

return logits. size: [b, 1]

测试网络。

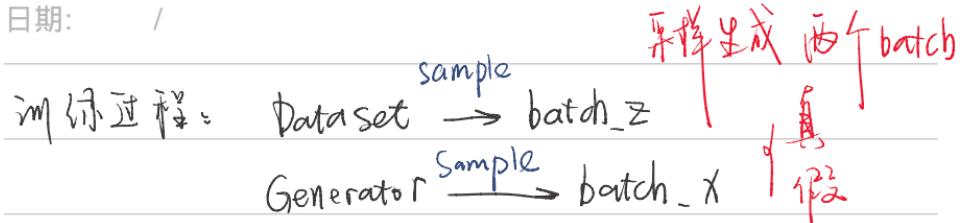
def main:

$d = \text{Discriminator}()$ $x = [[2, 64, 64, 3]]$

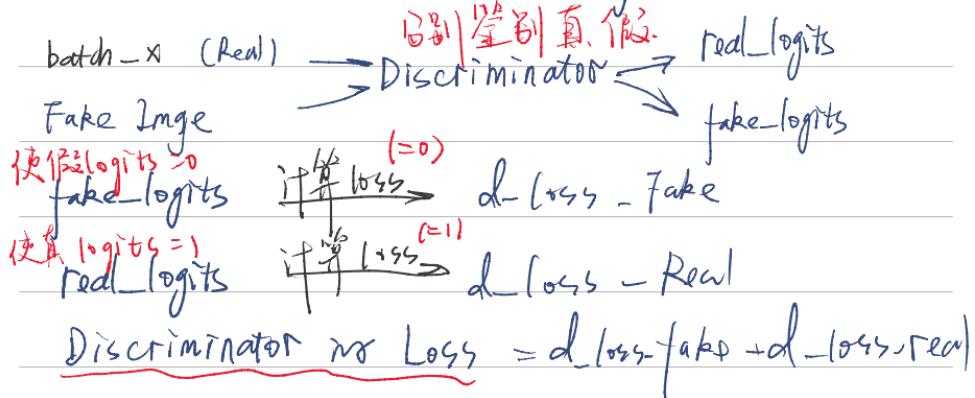
$g = \text{Generator}()$ $z = ([2, 100])$

$\text{prob} = g(x)$ $\hat{x} = g(z)$

日期:



真采样送入 $\xrightarrow{\text{Generator}}$ 得到一张假图
batch_Z $\xrightarrow{\text{Generator}}$ Fake Image.



G! 请将 Fake 希望、Fake 过 D 后得 1. (对)
fake_logits 使 $\text{loss} (=1)$ $\xrightarrow{\text{鉴别器真假值}} g_loss_fake$ 马过

Generator loss = g-loss-fake

d loss 计算梯度 $\xrightarrow{\text{梯度}}$ optimizer 更新
g-loss

日期: /

WGAN-GP

$$L = \underbrace{\mathbb{E}_{\tilde{x} \sim \mathbb{P}_g} [D(\tilde{x})] - \mathbb{E}_{x \sim \mathbb{P}_r} [D(x)]}_{\text{Original critic loss}} + \lambda \underbrace{\mathbb{E}_{\hat{x} \sim \mathbb{P}_{\hat{x}}} [(\|\nabla_{\hat{x}} D(\hat{x})\|_2 - 1)^2]}_{\text{Our gradient penalty}}$$

超参数

假 \rightarrow 微分逼近于 1



Our gradient penalty

where \hat{x} sampled from \tilde{x} and x with t uniformly sampled between 0 and 1

新图片、
 $\hat{x} = t\tilde{x} + (1-t)x$ with $0 \leq t \leq 1$
假 直 随机 sample.

使用

只用改一个地方。Loss 不变。gradient penalty

网络模型小修改

def gradient_penalty.

随机 batchsz = batch_x.shape.

$t = tf.random.uniform([batchsz, 1, 1, 1])$

[b, 1, 1, 1] \rightarrow [b, h, w, c]

t = tf.broadcast_to(t, (t, batch_x.shape))

interpolate = t * batch_x + (1-t) * fake_image.

tape.watch

d_interpolate_logits = discriminator(interpolate)

grads = tape.gradient(d_interpolate_logits, interpolate)

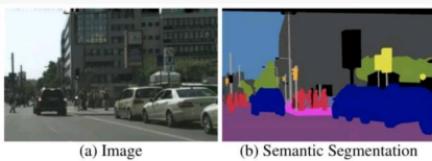
$\text{grads} = \text{tf.gradients}([\text{batch}_1, \dots, \text{batch}_n], [\text{gen}_1, \dots, \text{gen}_n])$
 $\text{gp} = \text{tf.norm}(\text{grads}, \text{axis}=1)$ [b] = 范数
 $\text{gp} = \text{tf.reduce_mean}((\text{gp}-1)^2) \text{ MSE}$
 return gp

def d_loss_fn

$\text{gp} = \text{gradient_penalty}(\text{discriminator}, \text{batchX}, \text{fake_image})$
 $\text{loss} = \text{d_loss_fake} + \text{d_loss_real} + \lambda * \text{gp}$
 $\downarrow \lambda \text{ 正则化系数}$

日期: /

图像语义分割



像素级，对每一个像素点进行分类

这里需要和实例分割区分开来。它没有分离同一类的实例；我们关心的只是每个像素的类别，如果输入对象中有两个相同类别的对象，则分割本身不将他们区分为单独的对象。

只隔

不知道
有几个



(c) Semantic segmentation



(d) This work segment individual object instances

语义分割的目标

$[h.w.3] / [h.w.1]$

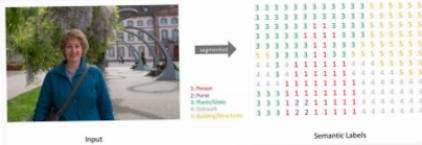
一般是将一张RGB图像 ($height*width*3$) 或是灰度图 ($height*width*1$) 作为输入，

输出的是分割图，其中每一个像素包含了其类别的标签 ($height*width*1$)

$[h.w.1]$

大小一样，没有像素值

每个像素都有 label

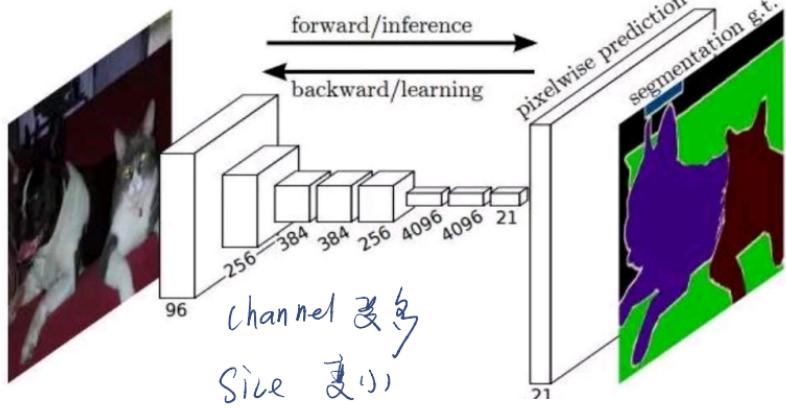


FCN

目前在图像分割领域比较成功的算法，有很大一部分都来自于同一个先驱：Long等人提出的Fully Convolutional Network (FCN)，或者叫全卷积网络。

FCN将分类网络转换成用于分割任务的网络结构，并证明了在分割问题上，可以实现端到端的网络训练。

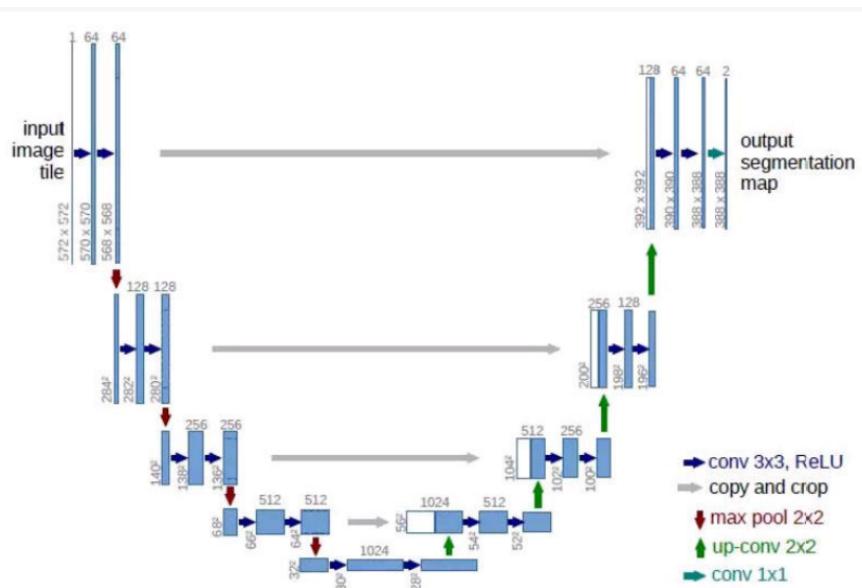
FCN成为了深度学习解决分割问题的奠基石。



UNET.

Unet是2015年诞生的模型，它几乎是当前segmentation项目中应用最广的模型。

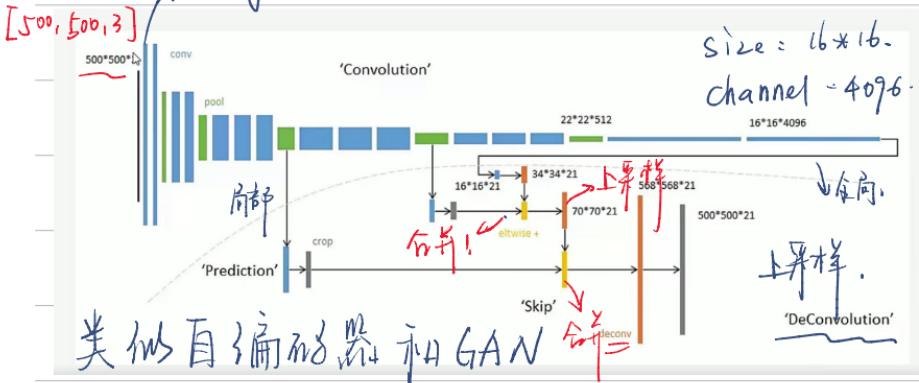
Unet能从更少的训练图像中进行学习。当它在少于 40 张图的生物医学数据集上训练时，IOU 值仍能达到 92%。



日期: /

图像语义分割的网络结构

Padding



1. 分类网络 (conv) → 上采样

2. 信息之前的信息

3. 输入任意尺寸

input_size = output_size

Output_Channel = $N + 1$
↓ 目标类数 背景

局部使用 conv, 无 FC. 先 strides 再卷积

上采样: Upsampling.

通过训练放大图像

差值法

反池化

反卷积

$$\boxed{9 \ 1} \rightarrow \boxed{\begin{matrix} 9 & 0 \\ 0 & 0 \end{matrix}}$$

最大

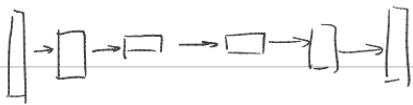
平均

$$\boxed{5} \rightarrow \boxed{\begin{matrix} 5 & 5 \\ 5 & 5 \end{matrix}}$$

$$\boxed{17}$$

$$\begin{matrix} w_1, w_2 \\ w_3, w_4 \\ w_5, w_6 \\ w_7, w_8 \end{matrix} \rightarrow \begin{matrix} a_{11}, a_{12} \\ a_{21}, a_{22} \\ a_{31}, a_{32} \\ a_{41}, a_{42} \end{matrix}$$

日期:



类自编码器结构

Skip. 跳级结构

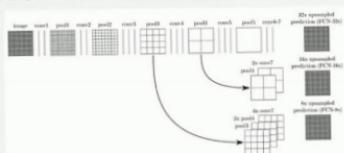
细节 → 抽象. 局部 → 全局.

1
32

- 将底层 (stride 32) 的预测 (FCN-32s) 进行 2 倍的上采样得到的图像，并与从 pool4 层 (stride 16) 进行的预测融合起来 (相加)，这一部分的网络被称为 FCN-16s。随后将这一部分的预测再进行一次 2 倍的上采样并与从 pool3 层得到的预测融合起来，这一部分的网络被称为 FCN-8s。

通过上采样使其 size 相同，可以运算。

增加 Skips 结构将最后一层的预测 (有富的全局信息) 和更浅层 (有更多的局部细节) 的预测结合起来，这样可以在遵守全局预测的同时进行局部预测。

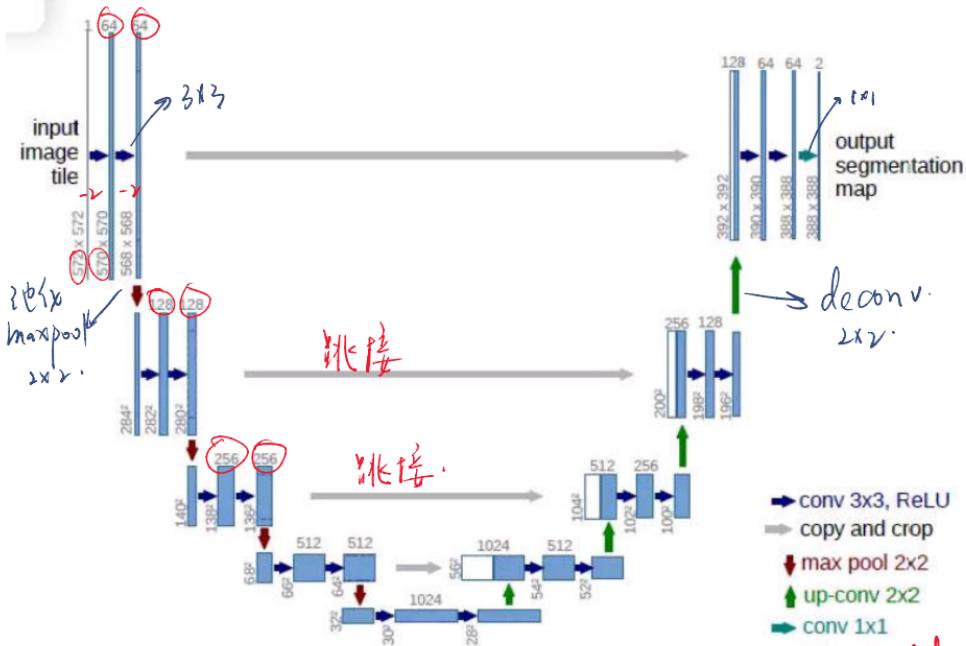


缺点：对细节不够敏感，效果不好。

日期: /

UNet 网络结构

前半部分 特征提取，后半部分是上采样。



采用的特征融合方式 = 跳接 (tf.concat).

对边像采用更高的权重。

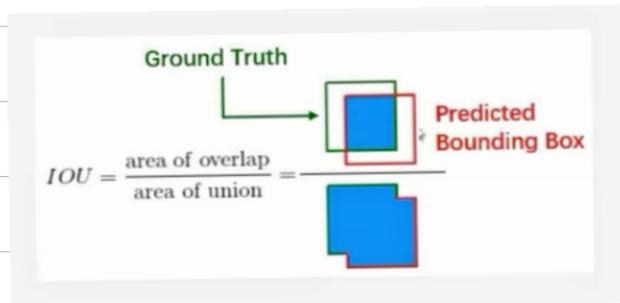
输入: 任意张图

输出: n(目标类别数) + 1 (背景)

日期: 2023.10.18.

IOU 评估指标.

表示: IoU 表示候选框 (candidate bound) 与原标记框 (ground truth bound)



日期： /

tf.data.

输入 + 处理文件。

tf.data, Dataset 表示一系列元素
包含多个 Tensor 对象。单个训练样本
-> 表示数据 / 框签的信号。

可以通过两种不同的方式来创建 tf.data.Dataset

1. 直接从 Tensor 创建 Dataset

例如 Dataset.from_tensor_slices() ;

当然 Numpy 也是可以的，TensorFlow 会自动将其转换为 Tensor。

2. 通过对一个或多个 tf.data.Dataset 对象来使用变换

(例如 Dataset.zip) 来创建 Dataset。

一个 Dataset 对象包含多个元素，每个元素的结构都相同。每个元素包含一个或多个 tf.Tensor 对象，这些对象被称为组件。

张量，元组。张量元组，
可以直接迭代？

日期： /