

Deep learning notes

Tototitata

六大挂科王

Lecture 1

Deep learning VS Machine Learning.

- ① DL: 用神经元训练。
- ② ML: 选 feature, 训练。

Tensor Internal: 内核?

Tensor 在储存的时候，都会被组织
为向量。

用 .storage() 直接。

储存的映射也有 stride 和 offset 参数。

Lecture 2

2-1

机器学习3大类：

↑
(按照目的)
来区分)

1. Classification
2. Regression → 预测值
3. Density Estimation

→ 判断属于哪类

Capture the struct
of data

(e.g. outlier detection/
image generation)

Expected Risk VS Empirical Risk

E.R.: 基础风险，模型于关于输入输出

↓ 不知道真实 (不知道 $P(X, Y)$)
输入.

Emp.R.: 经验风险：模型于关于训练集的平均
用 $\hat{P}(X, Y)$ 代替，损失.

Training: Minimize the empirical risk.

2-2.

KNN: k-nearest-neighbors

对于新输入，在训练集中与它最近的 k 个点，
属于 k 中的多数类。

k 太小，overfitting. 学习 noise.

k 太大，underfitting. 没学习

Capacity: 拟合函数的能力.
参数较少

Regularization: 让模型泛化.

↓
保留所有 features, 减少特征数量级。
① ↓. 越光滑.

2-3 Bias - Variance.

Generalization error (泛化误差),

: 训练集的 loss 和一般化数据集的 loss
之间的差异.

Bias. Variance. 用来描述学习的模型和真实模
型之间的差距.



所有训练集的预测输出的 mean output 和

real data for output 之间的关系 \downarrow trade off

不同训练集训练出的模型之间的差异.

2-5 Clustering and Embedding.

Embedding: 高维集合的元素映射到一个向量。

$$\begin{aligned}x_1 &\rightarrow \vec{v}_1 & \text{if } x_1 \text{ is similar with } x_2 \\x_2 &\rightarrow \vec{v}_2 & \rightarrow \vec{v}_1 \text{ is close to } \vec{v}_2\end{aligned}$$

PCA Algo.

Clustering: 美化元素在一起...
K-means.

Lecture 3

3-1 Perceptron

The basic mathematical model
for a neuron (神经元) was

Threshold Logic Unit.
(阈值逻辑单元)

$$f(x) = 1_{\{w \sum_i x_i + b \geq 0\}}$$

可以演变成例如

$$\text{or}(u, v) = 1_{\{u + v - 0.5 \geq 0\}}$$

$$w=1, b=-0.5$$

类似的基本逻辑门

调整参数



o/w: otherwise.

Perception:

$$f(x) = \begin{cases} 1 & \text{if } \sum_i w_i x_i + b \geq 0 \\ 0 & \text{o/w} \end{cases}$$

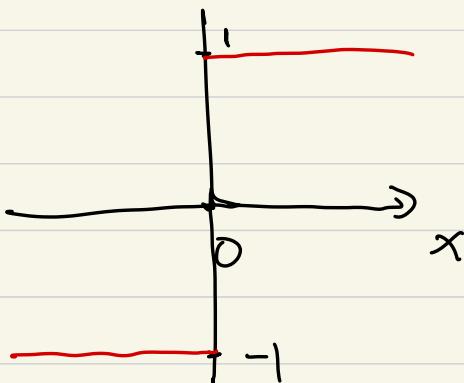
Rosenblatt

w: synaptic weights

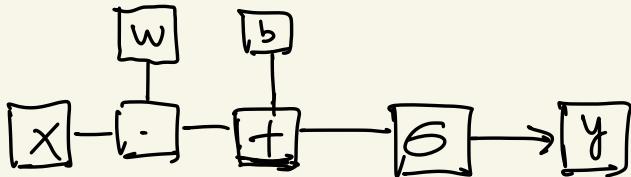
⇒ activation function σ

$$f(x) = \sigma(w \cdot x + b)$$

$$\sigma(x) = \begin{cases} 1 & \text{if } x \geq 0 \\ -1 & \text{o/w} \end{cases}$$



$$f(x) = \sigma(W \cdot x + b)$$



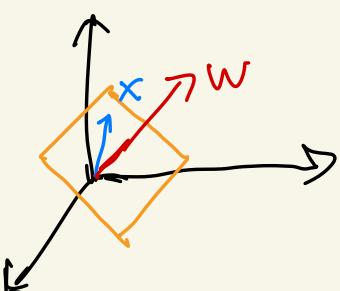
超平面: hyperplane

n维空间中 n-1 维的子平面，能把线性空间分割成不相交的两部分。

Ex: 2维空间中的一条直线

3维空间中的一个平面。

法向量 w: 垂直于超平面的向量



□ 为过原点的超平面。

w 为法向量

x 为超平面上的任意一点

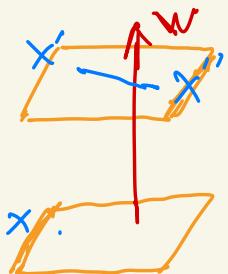
显然有 $w^T x = 0$ (法向量与平面上的向量积为0)

$$w = (w_1, w_2, w_3)$$

$$x = (x_1, x_2, x_3)$$

观察超平面沿着 w 方向移动。
此时 x 不再是超平面内的点，

$$w^T x \neq 0 .$$



在新的超平面内有2点, x' 和 x''

$$x' = (x'_1, x'_2, x'_3)$$

$$x'' = (x''_1, x''_2, x''_3)$$

$$\text{可知 } (\vec{x}' - \vec{x}'') \cdot \vec{w} = 0$$

$$(x'_1 - x''_1, x'_2 - x''_2, x'_3 - x''_3) \cdot (w_1, w_2, w_3) = 0$$

$$x'_1 w_1 + x'_2 w_2 + x'_3 w_3 = x''_1 w_1 + x''_2 w_2 + x''_3 w_3$$

$$w^T x' = w^T x''$$

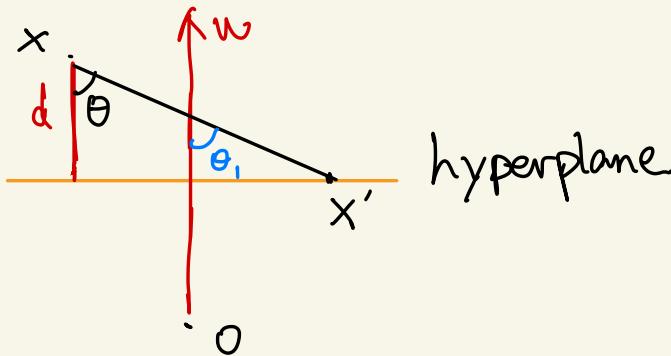
$$\sum b = -w^T x'' . x' \rightarrow x$$

由以得证

$$w^T x + b = 0$$

超平面
公式

点到超平面的公式



x 是平面外一点, 到平面距离为 d .

$$\cos \theta = \frac{d}{\|x - x'\|}$$

w 有向方向
→ $\|w\|$

$$|(x - x') \cdot w| = \|x - x'\| \cdot \|w\| \cdot \cos \theta$$

因为 $d \parallel w$, 有 $\theta = \theta_1$, 可利用点乘求 θ

$$d = \frac{|(x - x') \cdot w|}{\|w\|} = \frac{|wx - wx'|}{\|w\|}$$

在 hyperplane 内, $b = -w^T x'$

$$d = \frac{|wx + b|}{\|w\|}$$

④ 3. M. Perceptron.

Perceptron 本质上也是一个二分类的线性模型。

$$f(x) = \begin{cases} 1 & \text{if } \sum_i w_i x_i + b \geq 0 \\ 0 & \text{o/w} \end{cases}$$

\downarrow

$w^T x + b$

超平面.

也就是说有一个超平面 $w^T x + b = 0$ 存在。
才可以二分类， ≥ 0 和 < 0 . 即超平面的左右。

在 Perceptron algo 中.

- ① initialize w with 0. ($w^0 = 0$)
 - ② While $\exists n_k$ s.t. $y_{n_k}(w^k \cdot x_{n_k}) \leq 0$.
 update $w^{k+1} = w^k + y_{n_k} x_{n_k}$
- y_{n_k} to target output.
 $x_{n_k}(w^k \cdot x_{n_k}) \leq 0$
 说明 target 与 pred 不一样！ $-1 \neq 1$. $w^k \cdot x_{n_k} \neq 0$
 prediction.

所以要更新 w.

① 关于 $y(w^T x + b)$ 用乘判定分类是否正确.

y 为观察到的值.

$w^T x + b$ 映射到 y 的关系是

$$\begin{cases} 1 & w^T x + b \geq 0 \\ -1 & w^T x + b < 0 \end{cases}$$

所以, $w^T x + b$ 与 y 应该是同号

若出现异号说明当前这个 x 被分错了.

② 同时, $w^T x + b$ 也代表 x 到超平面的距离.

依据点到超平面的公式.

$$d = \frac{|w^T x + b|}{\|w\|}$$

当对 w 归一化, $\|w^*\| = 1$ 时.

$$d = |w^T x + b|.$$

$(w^T x + b)$ 则表示到超平面的
距离和方向。

在更新 w 时，Loss function =

M 为误差点

$$-\frac{1}{\|w\|} \sum_{x_i \in M} y_i (w^T x + b)$$

保证 L.F. > 0 距离 y_i : 错误的分类结果.

比如 $w^T x + b = 2$. y 应该 = 1.

观察到的是 -1. $\|w\| = 1$.

$$-\frac{1}{\|w\|} \sum y_i (w^T x + b) = 2$$

算法优化的目标是. $\min \underline{L(w, b)}$
↑
Loss. Fun.

用 L 分别对 w 和 b 求偏导.

$$\frac{\partial L(w, b)}{\partial w} = \frac{\partial -\sum_{x_i \in M} y_i (w^T x_i + b)}{\partial w} = -\sum_{x_i \in M} y_i x_i$$

$$\frac{\partial L(w, b)}{\partial b} = -\sum y_i$$

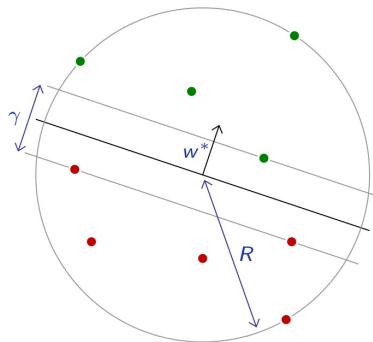
如果把 b 放进 $w \times$ 矩阵中。

$$w = (b, w_1, w_2, \dots, w_n)$$

$$x = (1, x_1, x_2, \dots, x_n)$$

这样写
 $w^{k+1} = w^k + y_i x_i$
 $(-(-y_i x_i))$

We can get a convergence result under two assumptions:



1. The x_n are in a sphere of radius R :

$$\exists R > 0, \forall n, \|x_n\| \leq R.$$

2. The two populations can be separated with a margin γ :

$$\exists w^*, \|w^*\| = 1, \exists \gamma > 0, \forall n, y_n (x_n \cdot w^*) \geq \gamma/2.$$

1. 有半径为 R 的球，能把所有样本包含进去。

2. 有一个超平面满足 $\|\hat{w}^*\| = 1$.

$$\hat{w}^T \hat{x} = 0$$

能将数据点完全分开，且最近的点到 hyperplane 的距离为 $\frac{\gamma}{2}$.

1. + 2. \Rightarrow Algo convergence

Proof: Suppose there still is a misclassified sample at iteration k .

We have,

$$\begin{aligned}
 w^{k+1} \cdot w^* &= (w^k + y_{nk} x_{nk}) \cdot w^* \\
 &= w^k \cdot w^* + y_{nk} (x_{nk} \cdot w^*) \\
 &\geq \underbrace{w^k \cdot w^*}_{\text{: by induction.}} + \frac{\gamma}{2} \\
 &\geq (k+1) \frac{\gamma}{2}
 \end{aligned}$$

$$\therefore \|w^k\| \|w^*\| \geq w^k \cdot w^*$$

$$\begin{aligned}
 \therefore \|w^k\|^2 &\geq (w^k \cdot w^*)^2 / \|w^*\|^2 \\
 &\geq k^2 \gamma^2 / 4
 \end{aligned}$$

若还有misclassified.

$$\begin{aligned}
 \|w^{k+1}\|^2 &= w^{k+1} \cdot w^{k+1} \\
 &= (w^k + y_{nk} x_{nk})^2 \\
 &= w^k \cdot w^k + \underbrace{2y_{nk} w^k \cdot x_{nk}}_{\leq 0} + \underbrace{\|x_{nk}\|^2}_{\leq R^2}
 \end{aligned}$$

$$\leq \|\mathbf{w}^k\|^2 + R^2$$

$$\leq (k+1)R^2$$

$$\frac{k^2 \gamma^2}{4} \leq \|\mathbf{w}\|^2 \leq kR^2$$

$$k \leq \frac{4R^2}{\gamma^2}$$

在 $\frac{4R^2}{\gamma^2} \uparrow$ iteration 當分差以遞減

3-2 Probabilistic view of a linear classifier

方差 Variance: 度量单个随机变量的
离散程度.

$$\sigma_x^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

n 为样本量. \bar{x} 为样本均值

协方差 covariance:

用来刻画两个随机变量的
相似程度.

$$Cov(x, y) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

两个变量同反向(正)还是同方向变化?

$|Cov(x, y)| \uparrow$, 同向/逆向程度↑

方差可以被视作 x 关于自身的协方差 $Cov(x, x)$

方差是协方差的一种特殊形式!!!
⇒ 同向程度很大

方差→协方差矩阵

给 d 个 随机变量 X_k , $k=1, \dots, d$

$$X_k \text{ 的方差为 } \sigma^2(X_k) = \frac{1}{n-1} \sum_{i=1}^n (x_{ki} - \bar{x}_k)^2$$
$$k = 1, \dots, d$$

Where, x_{ki} : i-th sample

n : sample size

求这些随机变量两两之间的协方差

$$\sigma^2(X_m, X_k) = \frac{1}{n-1} \sum_{i=1}^n (x_{mi} - \bar{x}_m)(x_{ki} - \bar{x}_k)$$

得到协方差矩阵

$$\Sigma = \begin{bmatrix} \sigma^2(x_1, x_1) & \cdots & \sigma^2(x_1, x_d) \\ \vdots & \ddots & \vdots \\ \sigma^2(x_d, x_1) & \cdots & \sigma^2(x_d, x_d) \end{bmatrix} \in \mathbb{R}^{d \times d}$$

↑ Symmetric matrix ↓ 对角线上为每个随机变量的方差

多元正态分布与线性转换

Multi-variate Gaussian distribution

假设一个向量 X 服从均值 μ , 协方差矩阵为 Σ 的
多元正态分布, 则

$$P(x) = \frac{1}{\sqrt{|2\pi\Sigma|}} \exp\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

.....

Homoskedasticity 同方差性.

(因 \hat{y} 的殘差項是隨機的) variance 不變

↓
residual: (覈算值 - 樣本估計值)

$$f(x; w, b) = G(w \cdot x + b)$$

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

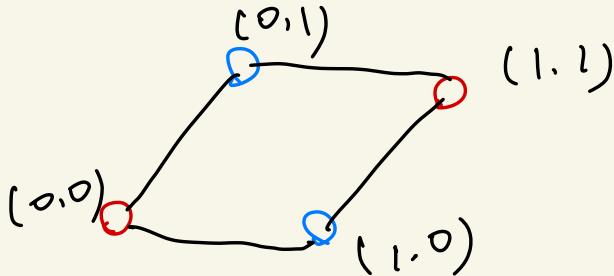
if $X \in \{-1, 1\}$

$$P(Y=y | X=x) = G(y(wx+b))$$

3.3.

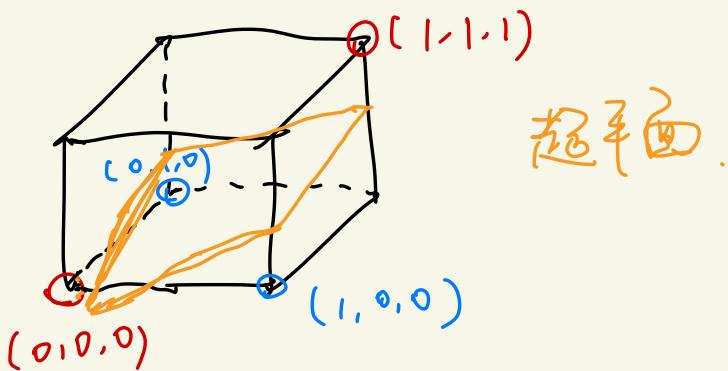
如果有时候数据不可线性分离
可以通过加一维的方式。

Ex.



找不到 linear separator

$$\phi: (x_u, x_v) \rightarrow (x_u, x_v, x_u x_v)$$



MLP multi-layer Perceptrons

① A linear classifier :

$$\mathbb{R}^D \rightarrow \mathbb{R}$$

$$x \mapsto g(w \cdot x + b)$$

with $w \in \mathbb{R}^D$, $b \in \mathbb{R}$, $g: \mathbb{R} \rightarrow \mathbb{R}$

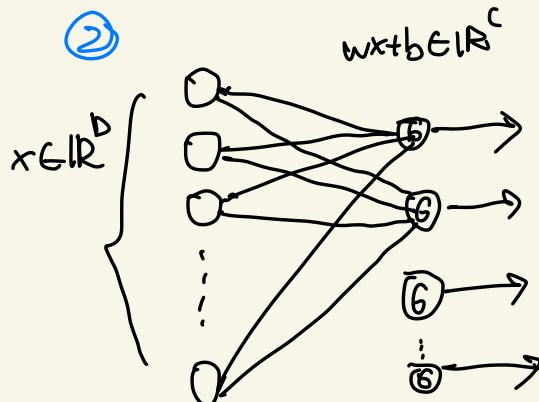
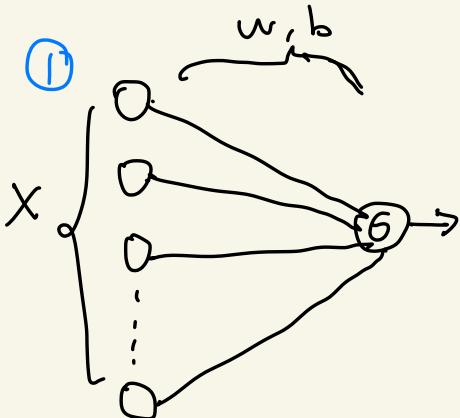
② Component-wise: piece-wise 线性

$$\mathbb{R}^D \rightarrow \mathbb{R}^C$$

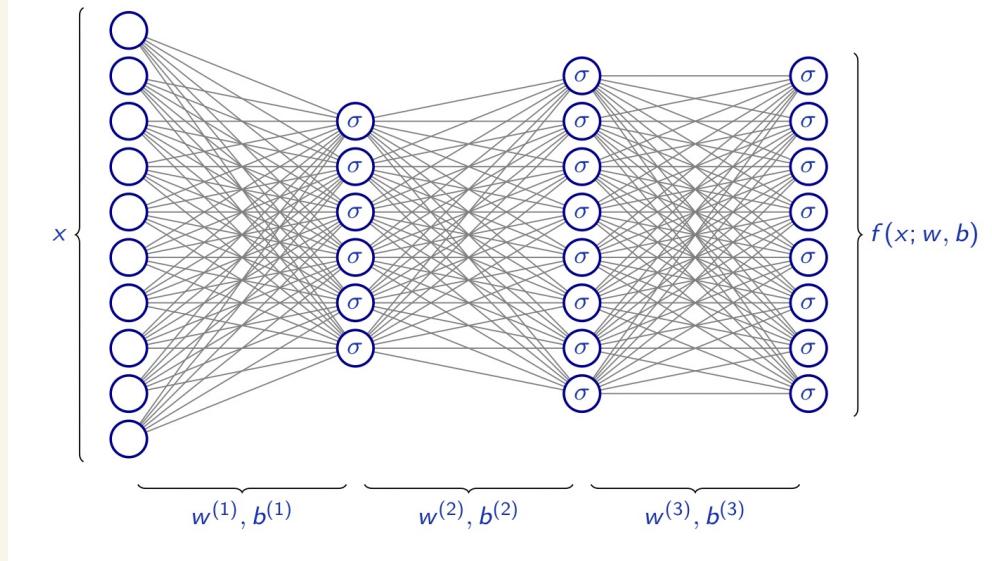
$$x \mapsto g(wx + b)$$

with $w \in \mathbb{R}^{C \times D}$, $w_c, b \in \mathbb{R}^C$

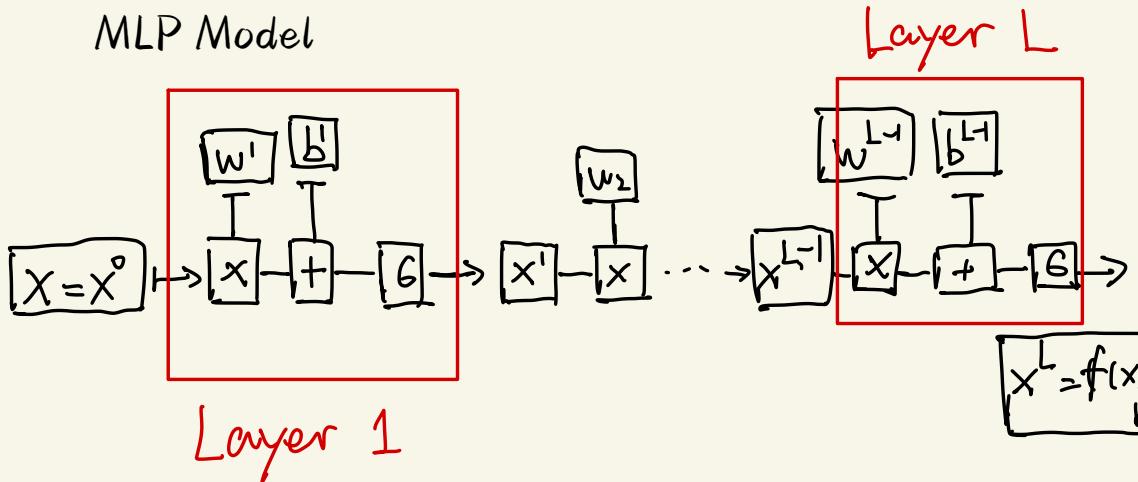
extended to multi-dimension output



Even though it has no practical value implementation-wise, we can represent such a model as a combination of units. More importantly, we can extend it.



MLP Model



$$\forall I=1, \dots, L \quad x^{(I)} = G(w^I x^{I-1} + b^I)$$

$$\text{and } f(x; w, b) = x^L$$

仿射变换 = 线性变换 + 平移！

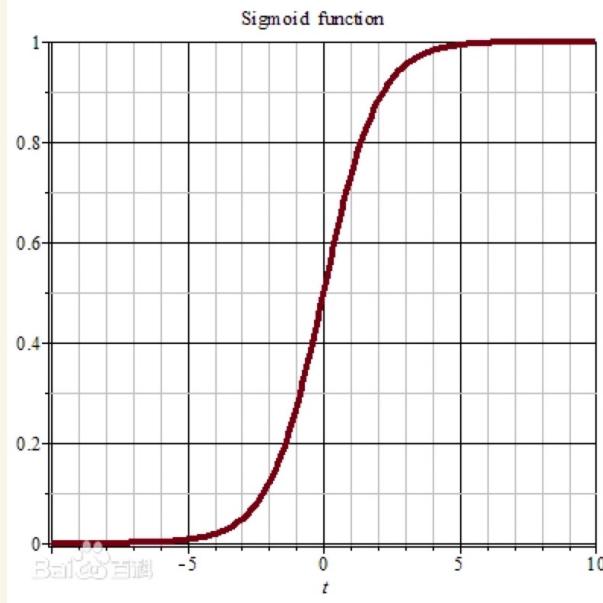
直线经过变换后依旧是直线且平行
线段的长度比例不变

激活函数不应该是仿射的（通常是非线性的）。

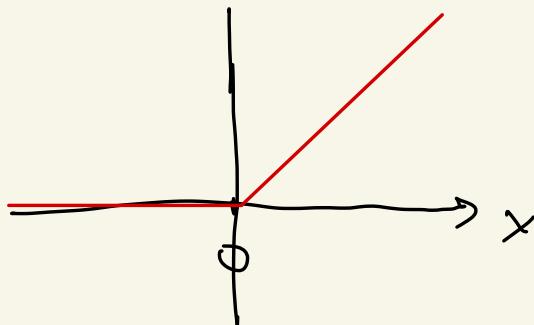
两个常见的激活函数：

1. Sigmoid
2. ReLU

Sigmoid: $S(x) = \frac{1}{1+e^{-x}}$



ReLU. $x \mapsto \max(0, x)$



Universal Approximation

U.A Theorem: 人工神经网络能逼近
任何函数.

So we can approximate any continuous function

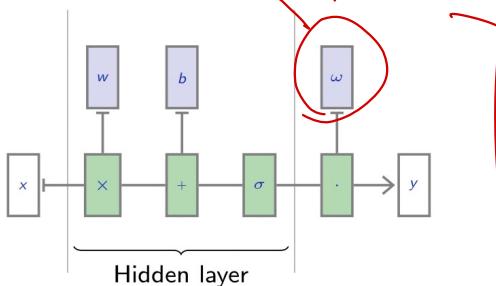
✓ 这已是连续的

$$\psi: [0, 1]^D \rightarrow \mathbb{R}$$

with a one hidden layer perceptron

$$x \mapsto \omega \cdot \sigma(wx + b),$$

where $b \in \mathbb{R}^K$, $w \in \mathbb{R}^{K \times D}$, and $\omega \in \mathbb{R}^K$. 那参数?



This is the **universal approximation theorem**.

Universal approximation theorem

通用逼近定理

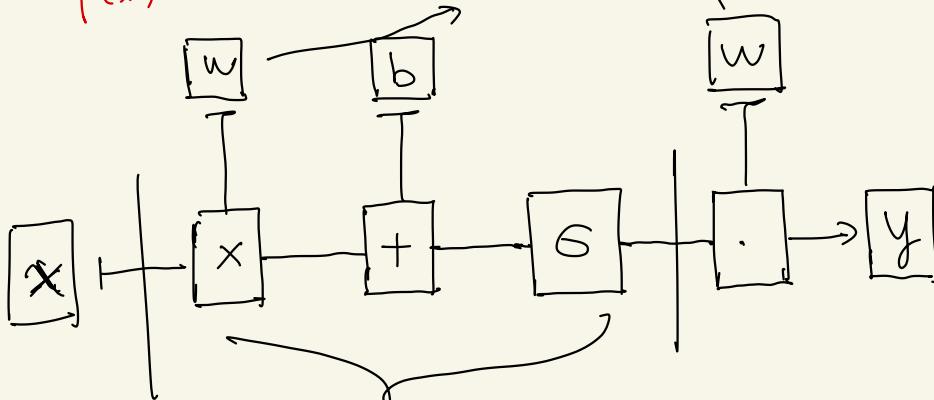
We approximate any $\psi \in C([a,b], \mathbb{R})$ with [a linear combination] of translated / scaled ReLU functions.

$$f(x) = G(w_1x + b_1) + G(w_2x + b_2) + \dots + G(w_nx + b_n)$$

重复尝试 ReLU 单元的不同组合，

直到 $f(x)$ 才以含目标函数。

$$f(x) = w \cdot G(wx + b)$$



Hidden layer.

Hidden layer larger ↑

Approximation precision ↑

training error ↓

No business with

test error

U.A.T states ↑

Gradient Descent

MSE:

$$J(w, b) = \frac{1}{N} \sum_n (f(x_n; w, b) - y_n)^2$$

通常常是 $\frac{1}{2N}$

← Loss

↓
min

↑
For classification, certain regressions ...

For logistic regression:

$$P_w = \underbrace{(Y=1 | X=x)}_{G(x)} = G(wx+b)$$

$$G(x) = \frac{1}{1+e^{-x}} \leftrightarrow \text{Sigmoid}$$

$$J(w, b) = - \sum_n \log G(y_n(w \cdot x_n + b))$$

逻辑回归试图找到分类概率 $P(Y=1)$ 与 x 的关系，比
较概率值来判断。

与实际值。

Given a function:

$$f: \mathbb{R}^D \rightarrow \mathbb{R}$$
$$x \mapsto f(x_1, x_2, \dots, x_D)$$

its gradient:

$$\nabla f: \mathbb{R}^D \rightarrow \mathbb{R}$$
$$x \mapsto \left(\frac{\partial f}{\partial x_1}(x), \dots, \frac{\partial f}{\partial x_D}(x) \right)$$

Back-propagation 反向传播.

假设现在要训练一个MLP模型.

目标: Minimize losses over the training set

Loss function: $L(w, b) = \sum_n l(f(x_n; w, b), y_n)$

如果用 Gradient descent, 那就对 $L(w, b)$.

求关于 w 和 b 的偏导.

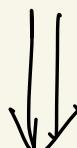
$$\frac{\partial L_n}{\partial w_{i,j}} \quad \text{and} \quad \frac{\partial L_n}{\partial b_i}$$

Forward Pass 正向传播

首先引入 $s^{(1)}, s^{(2)}, \dots, s^{(L)}$

记录在 activation function \rightarrow FA 值.

$$x^{(0)} = x \xrightarrow{w^{(1)} b^{(1)}} s^{(1)} \xrightarrow{g} x^{(1)} \rightarrow \dots \xrightarrow{w^{(L)} b^{(L)}} s^{(L)} \xrightarrow{g} x^{(L)} \rightarrow f(x, w, b)$$



$$\forall I = 1, \dots, L. \quad \left\{ \begin{array}{l} s^{(I)} = w^{(I)} \cdot x^{(I-1)} + b^{(I)} \\ x^{(I)} = \sigma(s^{(I)}) \end{array} \right.$$

Output: $x^{(L)} = f(x; w, b)$

Backpass:

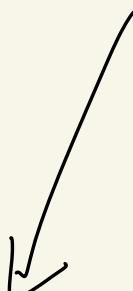
B-P 算法基于 "Chain Rule"

$$(g \circ f)' = (g' \circ f) \cdot f'$$

$$x^{(I-1)} \xrightarrow{w^{(I)}, b^{(I)}} s^{(I)} \xrightarrow{\sigma} x^I$$

Since $s_i^{(I)}$ influences f only through $x_i^{(I)}$ with

$$x_i^{(I)} = \sigma(s_i^{(I)})$$



We have

$$\frac{\partial l}{\partial s_i^{(I)}} = \frac{\partial l}{\partial x_i^{(I)}} \cdot \frac{\partial x_i^{(I)}}{\partial s_i^{(I)}} = \frac{\partial l}{\partial x_i^{(I)}} g'(s_i^{(I)})$$

And since $x_j^{(I-1)}$ influences l only through the $s_i^{(I)}$ with

$$s_i^{(I)} = \sum_j w_{i,j}^{(I)} x_j^{(I-1)} + b_i^{(I)}$$

We have

$$\frac{\partial l}{\partial x_j^{(I-1)}} = \sum_i \frac{\partial l}{\partial s_i^{(I)}} \cdot \frac{\partial s_i^{(I)}}{\partial x_j^{(I-1)}} = \sum_i \frac{\partial l}{\partial s_i^{(I)}} w_{i,j}^{(I)}$$

Since $w_{i,j}^{(I)}$ and $b_i^{(I)}$ influences l only through $s_i^{(I)}$ with

$$s_i^{(I)} = \sum_j w_{i,j}^{(I)} x_j^{(I-1)} + b_i^{(I)}$$

We have

$$\frac{\partial l}{\partial w_{i,j}^{(I)}} = \frac{\partial l}{\partial s_i^{(I)}} \cdot \frac{\partial s_i^{(I)}}{\partial w_{i,j}^{(I)}} = \frac{\partial l}{\partial s_i^{(I)}} x_j^{(I-1)},$$

$$\frac{\partial l}{\partial b_i^{(I)}} = \frac{\partial l}{\partial s_i^{(I)}}$$

Summary: We can compute $\frac{\partial l}{\partial x_i^{(I)}}$ from the def. of l .
and recursively propagate backward the derivatives

of the loss $\frac{\partial l}{\partial s_i^{(I)}} = \frac{\partial l}{\partial x_i^{(I)}} g'(s_i^{(I)})$

$$\frac{\partial l}{\partial x_j^{(I-1)}} = \sum_i \frac{\partial l}{\partial s_i^{(I)}} w_{i,j}^{(I)}$$

And then compute the derivate

$$\frac{\partial \ell}{\partial w_{i,j}^{(l)}} = \frac{\partial \ell}{\partial s_i^{(l)}} \times_j^{(l-1)}$$

$$\frac{\partial \ell}{\partial b_i^{(l)}} = \frac{\partial \ell}{\partial s_i^{(l)}}$$

Forward pass

Compute the activations.

$$x^{(0)} = x, \quad \forall l = 1, \dots, L, \quad \begin{cases} s^{(l)} = w^{(l)}x^{(l-1)} + b^{(l)} \\ x^{(l)} = \sigma(s^{(l)}) \end{cases}$$

Backward pass

Compute the derivatives of the loss w.r.t. the activations.

$$\left\{ \begin{array}{l} \left[\frac{\partial \ell}{\partial x^{(l)}} \right] \text{ from the definition of } \ell \\ \text{if } l < L, \left[\frac{\partial \ell}{\partial x^{(l)}} \right] = (w^{(l+1)})^\top \left[\frac{\partial \ell}{\partial s^{(l+1)}} \right] \end{array} \right. \quad \left[\frac{\partial \ell}{\partial s^{(l)}} \right] = \left[\frac{\partial \ell}{\partial x^{(l)}} \right] \odot \sigma'(s^{(l)})$$

Compute the derivatives of the loss w.r.t. the parameters.

$$\left[\frac{\partial \ell}{\partial w^{(l)}} \right] = \left[\frac{\partial \ell}{\partial s^{(l)}} \right] (x^{(l-1)})^\top \quad \left[\frac{\partial \ell}{\partial b^{(l)}} \right] = \left[\frac{\partial \ell}{\partial s^{(l)}} \right].$$

Gradient step

Update the parameters.

$$w^{(l)} \leftarrow w^{(l)} - \eta \left[\frac{\partial \ell}{\partial w^{(l)}} \right] \quad b^{(l)} \leftarrow b^{(l)} - \eta \left[\frac{\partial \ell}{\partial b^{(l)}} \right]$$

Pytorch 的 Tensor 有 `requires_grad`

的属性，默認為 `False`

⇒ 可以用来计算 gradient.

↓

手写改为 `True`.

只有 floating point type 有梯度

$x = [1., 10.]$

```
[3]: t = torch.tensor([1., 2., 4.]).requires_grad_()
u = torch.tensor([10., 20.]).requires_grad_()
a = t.pow(2).sum() + u.log().sum()
```

```
[4]: torch.autograd.grad(a, (t, u))
```

$$\frac{\partial a}{\partial t}, \frac{\partial a}{\partial u}$$

/usr/local/lib/python3.8/dist-packages/torch/autograd/__init__
on your system. Please check that you have an NVIDIA GPU and i
(Triggered internally at /pytorch/c10/cuda/CUDAFunctions.cpp:
return Variable._execution_engine.run_backward()

```
[4]: (tensor([2., 4., 8.]), tensor([0.1000, 0.0500]))
```

$$t = [1., 2., 4.]$$

$$u = [10., 20.]$$

$$a = \sum t_i^2 + \sum \log(u)$$

$$(nx' = \frac{1}{x})$$

$$\frac{\partial a}{\partial t} = \left[\frac{\partial a}{\partial t_1}, \frac{\partial a}{\partial t_2}, \frac{\partial a}{\partial t_3} \right]$$

$$= [2t_1, 2t_2, 2t_3]$$

$$= [2., 4., 8.]$$

$x = \text{torch.tensor}([-3., 2., 5.]).\text{requires_grad}_()$

$u = x.\text{pow}(3).\text{sum}()$

$x.\text{grad}$

$u.\text{backward}()$ → ? 相当于关于 x 求导, $\frac{\partial u}{\partial x^2}$

$x.\text{grad}$ → $x = [-3, 2, 5]$ 时的 值

⇒ $\text{tensor}([27., 12., 75.])$

$$x = [-3, 2, 5]$$

$$u = \sum x_i^3$$

$$\frac{\partial u}{\partial x} = \left[\frac{\partial u}{\partial x_1}, \frac{\partial u}{\partial x_2}, \frac{\partial u}{\partial x_3} \right]$$

$$= [3x_1^2, 3x_2^2, 3x_3^2]$$

$$= [27, 12, 75]$$

好久没学习了！
从现在下学期开始

4.4 Convolution 卷积神经网络

① 来源：有时要处理数据量非常大（多维）的信息，

例如 image of RGB 256×256 . \rightarrow input

would require $(256 \times 256 \times 3)^2$ parameters.

发现信息在转换(transformation)时有一些共性(invariance).

于是将局部共性推广到全局。

为什么呢？

A transformation meaningful at a certain location
can be applied everywhere!

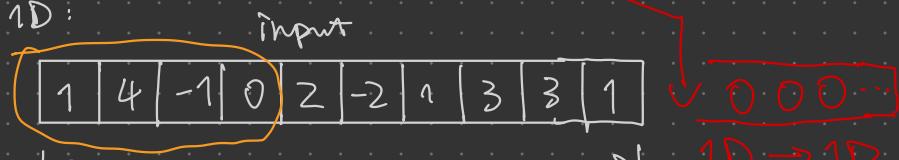
P.S. convolution 保留输入结构的一致性。

(the dimensions are same).

$$\begin{array}{c} 2d \\ 3d \end{array} \rightarrow \begin{array}{c} 2d \\ 3d \end{array}$$

+ 指明位置对应关系

1D:



$1D \rightarrow 1D$

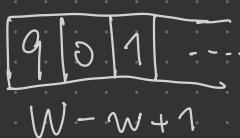
但输出比输入短
(+填充可以一样)

$$q = 1 \times 1 + 4 \times 2 + 0 \times (-1) + 0 \times (-1)$$

kernel



Output:



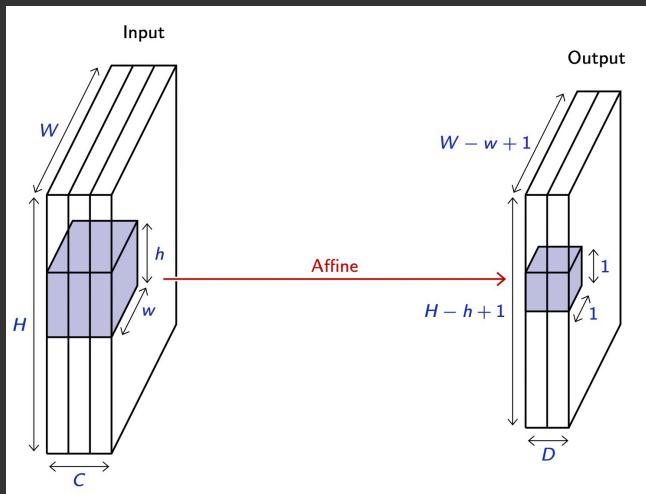
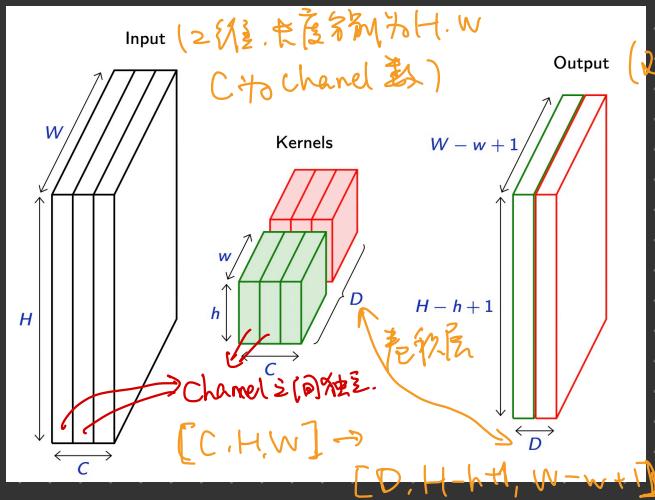
上圖的數字表述

input: $x = (x_1, x_2, \dots, x_w)$

kernel: $u = (u_1, u_2, \dots, u_w)$

$$(x * u) = (x_i, \dots, x_{i+w-1}) \cdot u$$

② 因素:



③ 亂數:

1. F. Conv2d

Convolution layer
↓ channel.

參數 weight : $D \times C \times \underbrace{h \times w}_{\text{size of filters}}$

bias : D

input : $\underline{N} \times C \times H \times W$

result : $\underline{N} \times D \times (H-h+1) \times (W-w+1)$
number of samples

2. torch.nn.Conv2d

參數}

① padding: input 在 0.

(2, 1): 上下各在 2 行, 左右各 1 行。

② stride: how coarsely the filter should be sniped across the signal.

(2, 2): 每隔 2 個距離，映射一次

③ dilation: filter 的擴張，填 0

Dilation, $1 \rightarrow 2$

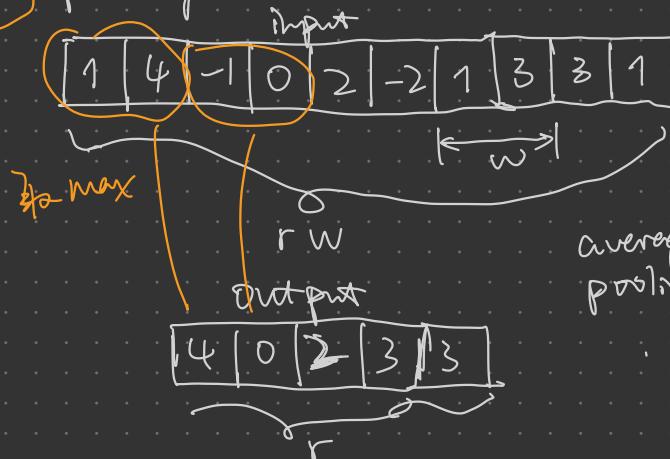


4.5 Pooling

簡化 Signal: 从高維信号中得到低維信号

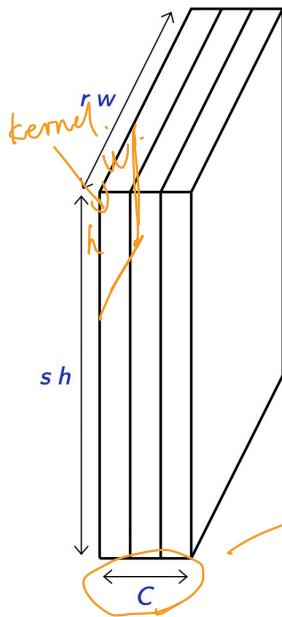
: Grouping → Pooling

① Max-pooling:

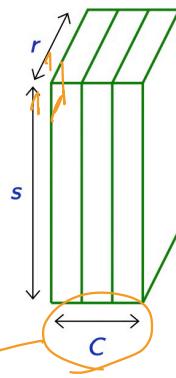


average: 平均
pooling

Input



Output



Convolution from
 $C \rightarrow D$

② 逐數:

1. $\text{F} \cdot \text{max_pool2d}$:

input: $N \times C \times H \times W$

kernel: (h, w)

Output: $N \times C \times \lfloor H/h \rfloor \times \lfloor W/w \rfloor$

2. `torch.nn.MaxPool2d`

卷积、神经网络 Convolution

Inspiration:

全连接神经网络 \Leftrightarrow 卷积

full connected

\hookrightarrow 参与过多

C.
解决
方案.

卷积、神经元: $\left\{ \begin{array}{l} \text{width} \\ \text{height} \\ \underline{\text{depth}} \end{array} \right.$ 不是网络深度.
激活数据样本的
第三维度.

卷积网络 {
 输入层
 卷积层 } 卷积层
 ReLU 层
 池化层
 全连接层

卷积层：产生网络中大部分的计算量！

- 作用：
- ① 卷积、偏置参数 = 可学习 不是参数。
即滤波器集合
保证(宽、高较小) 平直 ->
 - ② 做作神经元的一个输出。
与左右神经元共享参数
 - ③ 降低参数数量

感受野 Receptive field

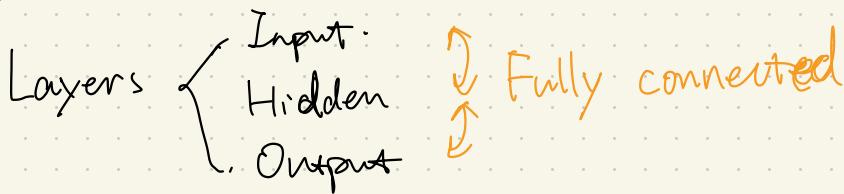
每个神经元和输入数据局部区域连接。
这个连接空间的大小 \Rightarrow 感受野

Receptive field 在深度方向和输入保持一致。



假设 input 为 RGB 图，
大小 $[32 \times 32 \times 3]$ 。
R-field: $[5 \times 5]$
权重 $5 \times 5 \times 3$ 个 + 1. bias

MLP



5 - 1

One-hot Coding: 将类别量转化为二进制特征。

Cross-entropy loss: 又叫熵损失,

描述了两个概率分布之间的距离。