

deep-learning 笔记

徐世桐

1 基本定义

深度学习: 得到数据深度的特征, 不一定使用神经网络

label 标签: 输出结果, \hat{y} 为计算得到的结果, y 为实际测量结果

feature 特征: 用于预测标签的输入变量, $x_j^{(i)}$ 为第 i 组 sample 第 j 号特征

sample 样本: 一组特征的取值和对应的标签输出

batch: batch size 个 sample 被分为一组, 进行向量化的计算, 称 B

hyperparameter 超参数: 人为设定的参数。如样本个数 (批量大小 batch size) $|B|$, 学习率 η 。少数情况下通过学习得到

W 一层 layer 的权重矩阵: 行数 = 前层节点数, 列数 = 当前层节点数

全连接层 fully-connected layer/稠密层 dense layer: 此层所有节点都分别和前一层所有节点连接

sigmoid 函数:

$\sigma(x) = \frac{1}{1+e^{-x}}$ 可能造成 gradient vanishment

为特殊的 softmax: $\frac{1}{1+e^{-x}} = \frac{e^0}{e^x + e^0} = \frac{e^{x_i}}{\sum_{i \in \{0,1\}} e^{x_i}}$

tanh 函数:

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

平均值在 0, 允许输出为负值, 帮助 converge。计算开销较大

多层 tanh 或 sigmoid 都将降低学习效率

softmax: $\text{softmax}(X) = \frac{\exp(-x)}{\sum_{x' \in X} \exp(-x')}$

softmax 函数: $\text{softmax}(X) = \frac{\exp(x)}{\sum_{x' \in X} \exp(x')}$, 将数值输出转化为概率值, 1. 值为正 2. 值总和为 1

ReLU 激活函数

仅有一半 neuron activate, 前向计算更快速。dying relu problem: relu 值为 0 导致权重不被更新

Parapetric ReLU:

$$PReLU(x) = \begin{cases} \alpha x & x < 0 \\ x & x \geq 0 \end{cases}$$

α 训练时更新

leaky ReLU:

同 PReLU, α 为超参数

Exponential LU

$$ELU(x) = \begin{cases} \alpha(e^x - 1) & x < 0 \\ x & x \geq 0 \end{cases}$$

形状从-1按 $\exp(x)$ 增长到0,随后线性。 α 从0增大时 $x \leq 0$ 包含更大部分的 \exp 图像,值 $> -\alpha$ 仅在 $\alpha = 1$ 时 continuously differentiable

CELU

$$ELU(x) = \begin{cases} \alpha(e^{\frac{x}{\alpha}} - 1) & x < 0 \\ x & x \geq 0 \end{cases}$$

相对ELU对 x 除超参数 α 。

α 从0增大时图像 $x \leq 0$ 近似linear部分逐渐增大 增长变平缓,值 $> -\alpha$

当 α 不为1时求导为 continuous function 即 continuously differentiable

SELU

$$SELU(x) = scale * ELU, scale > 1$$

self normalizing: 当每层激活函数使用 internal normalizing 函数时,最终的特征 mean variance converge 到一定值

激活函数需要经过(0,0),有正负值。否则重复 apply 导致 mean 值为无限

有区域输入绝对值较大时斜率较小。当较浅网络层 variance 较高时用于 damp variance

临近0位置斜率较高。增加较浅网络层 variance

为一 continuous curve。保证有 mean variance 的 fix point

GELU

$$GELU(x) = x\Phi(x)$$

Φ 为 Gaussian 的 CDF。图像形状为在0附近有小幅降低后 linear 升高

SoftPlus 激活函数:

$$SoftPlus(x) = \frac{1}{\beta} \log(1 + \exp(\beta x))$$

形如 \exp ,拐点连续可导。当 β 增大时,接近 ReLU

为保证 Numerical stability,当 $\beta x >$ 阈值时,改用线性激活函数

LogSigmoid 代价函数: $\log(\text{sigmoid}(x))$ 形状约到 $x=2$ 时线性上升到0,随后水平

LogSoftmax 代价函数: $\log(\text{softmax}(x))$ linear 形状,不经过(0,0)

cross entropy: 代价函数

定义:分部 p 和分部 q 间的 cross entropy $H(p, q) = -E_p(\log(q))$ 。为 expected value of $\log(q)$ with respect to distribution p

$$\text{针对单类别 classification: } CELoss(x, target_class) = -\log\left(\frac{\exp(x_{target_class})}{\sum \exp(x)}\right)$$

即-LogSoftmax

categorical cross entropy

$$\text{公式: 对一样本的输出 } \hat{y}^{(i)}, H(y^{(i)}, \hat{y}^{(i)}) = -\frac{1}{n} \sum_{j < n} y_j^{(i)} \log(\hat{y}_j^{(i)})$$

n 为输出标签数

对一批量的输出使用 cross entropy: 对每一样本 i 的 $H(y^{(i)}, \hat{y}^{(i)})$ 求和

使用: 联系两个值概率分部间的差异,即可将数值输出 \hat{y} 和分类结果 y 直接做对比

不适用于 multi-class classification, 仅有唯一类别作为输出

原因：multi-class 单一节点输出 binary classification 即代表 2 节点 softmax 输出，

而 categorical cross entropy 假设 其余类别输出 随 期望类别输出升高 而降低，则忽略其余类别的输出是否为 0

仍可和 softmax 同时使用，softmax 将可能性先转换为正数并和为 1，随后使用 cross entropy

binary cross entropy

公式： $H(y, \hat{y}) = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y})$

\hat{y}, y 为单一标量

使用：

使用于单个输出的 **logistic regression**，或 **multi-label regression**

multi label 时将每一 label 分别进行 binary cross entropy 后求和

multi-label regression 代价值为每一 label 预测值的 binary cross entropy 求和

Smooth L1 代价函数

$$loss(x, y) = \frac{1}{n} \sum_i z_i$$

$$z_i = \begin{cases} \alpha(x_i - y_i)^2 & |x_i - y_i| < 1 \\ |x_i - y_i| - \alpha & otherwise \end{cases}$$

α 为超参数，目标为在 $(-\alpha, \alpha)$ 范围内用二次函数替代 L1 拐点，拐点更平缓。

Neg log likelihood loss 代价函数

每一样本输出为对 m 个类别的 neg log likelihood，即 $\hat{y} \in \mathbb{R}^m$

对一样本 n ，期望输出类别为 c

定义 $w_n = [0, \dots, weight[c], 0, \dots]$ 。即仅有 c 位置为 $weight[c]$ 的向量， $weight$ 为 c 类别的权重

定义 $l_n = -w_n \hat{y}_n$ ，即取 c 类别的预测输出，乘以 c 类别权重 取负

$$NLLLoss(x, y) = \begin{cases} \sum_{n=1}^N l_n & reduction = sum \\ \frac{\sum_{n=1}^N l_n}{\sum_{n=1}^N w_n} & reduction = mean \end{cases}$$

Margin Ranking loss

$$loss(x, y) = \max(0, -y \cdot (x_1 - x_2) + margin)$$

x_1, x_2 为两样本到类别 centroid 的距离非两样本的 **embedding**， y 为正代表 x_1 应当大于 x_2 常取 ± 1

Triplet Margin Loss

$$loss(x_a, x_p, x_n) = \max(0, m + |f(x_a) - f(x_b)| - |f(x_a) - f(x_n)|)$$

即使正样本对 x_a, x_p 距离降为 0 负样本对 x_a, x_n 距离高于 threshold m

Cosine Embedding loss

$$loss(x, y) = \begin{cases} 1 - cosine(x_1, x_2) & y = 1 \\ \max(0, cosine(x_1, x_2) - m) & y = -1 \end{cases}$$

所有 x_i 特征为 normalized。即使正样本对 x_1, x_2 点乘接近 1 负样本对 x_a, x_n 点乘小于 threshold m

指数加权移动平均

$$y_t = \gamma y_{t-1} + (1 - \gamma)x_t$$

小批量乘法

对 n 个形状为 (a, b) 矩阵 X_1, X_2, \dots, X_n , n 个形状为 (b, c) 矩阵 Y_1, Y_2, \dots, Y_n , 小批量乘法结果为 $X_1 Y_1, X_2 Y_2, \dots, X_n Y_n$

* 特指被 word2vec 使用的乘法, 其他模型仍根据矩阵乘法, 非所有批量乘法都使用 小批量乘法

Sigmoid 二元交叉熵损失函数

得到向量 $P = [p_1, \dots, p_n]$, 向量 $L = [l_1, \dots, l_n]$, 掩码向量 $M = [m_1, \dots, m_n]$

p_i 标记 i 位置事件的可能性, l_i 标记期望 i 位置的事件发生 ($l_i = 1$) 或不发生 ($l_i = 0$)

计算:

1. 遍历 L , 若 $l_i = 0$, $p'_i = -p_i$

2. 对 P' 中每一 p'_i 取 sigmoid 值, 对所有 $m_i = 1$ 的 p'_i 求平均值

整体计算为: $f(P, L, M) = \frac{\sigma(P) \odot (-1)^L \odot M}{\text{sum}(M)}$

Laplace smoothing

令 $n(x)$ 为 x 的出现个数, m 为词的个数, n 为所有词出现次数总和

目标: 解决在特殊词组 $[x, x']$ 出现次数过小时 $n(x, x')$ 较小。导致原公式 $P(x'|x) = \frac{n(x, x')}{n(x)}$ 值过小

$$P(x) = \frac{n(x) + \epsilon_1 / m}{n + \epsilon_1}$$

$$P(x'|x) = \frac{n(x, x') + \epsilon_2 P(x')}{n(x) + \epsilon_2}$$

$$P(x''|x, x') = \frac{n(x, x', x'') + \epsilon_3 P(x'|x)}{n(x, x') + \epsilon_3}$$

Laplace operator Δ

对函数 $f(x, y)$: $\Delta f = \frac{d^2 f}{dx^2} + \frac{d^2 f}{dy^2}$

函数卷积

对函数 $f, g: \mathbb{R}^n \mapsto \mathbb{C}$, 定义函数卷积 $(f * g)(x) = \int_{\mathbb{R}^n} f(\tau)g(x - \tau) d\tau$

f, g 取离散值时, $(f * g)(x) = \sum_{\tau} f(\tau)g(x - \tau)$

$f(m), g(m)$ 都为标量。 $f(x)g(x - \tau)$ 意为将卷积核左移 n 后和 $f(x)$ 按元素相乘求和, 求和结果

为输出中 x 位置的值

性质:

Commutativity: $f * g = g * f$

Associativity: $f * (g * h) = (f * g) * h$, $a(f * g) = (af) * g$

Distributivity: $f * (g + h) = (f * g) + (f * h)$

求导: $\frac{d}{dx}(f * g) = \frac{df}{dx} * g = f * \frac{dg}{dx}$

$\int_{\mathbb{R}^n} f(\tau)g(x - \tau) d\tau = \int_{\mathbb{R}^n} f(x - \tau)g(\tau) d\tau$ 由交换律

$$= \int_{\mathbb{R}^n} f(x + \tau)g(-\tau) d\tau$$

即计算时先将卷积核倒序

shift invarient

对于一操作 f , $f(x) = f(x + \epsilon)$ 称 shift invarient。即当输入图片目标移动 仍能输出正确分类

如: 全连接层, 最大池化层。不论目标位置, 输出标签/张量值相同

shift equivariance

输出应随输入目标移动而移动

如: 卷积层。目标移动后输出矩阵值改变

Cosine distance

$$d(x, y) = \frac{x \cdot y}{\|x\| \|y\|}$$

即向量 x, y 夹角

Euler-Lagrange equation

定义: u 为一函数, $u'(t) = \frac{du(t)}{dt}$, $S(u) = \int_t L(t, u(t), u'(t)) dt$ 。

Theorem: u 为一 S 的 stationary point

$$\text{iff } \frac{L(t, u(t), u'(t))}{du} - \frac{L(t, u(t), u'(t))}{d(u'(t))dt} = 0$$

bias variance

定义:

实际数据模型为 f , 样本集合为 $\{x_i, y_i\}$, 训练得到模型为 \hat{f} 。

有 $y_i = f(x_i) + \epsilon$, ϵ 为噪声, 有 variance σ^2 均值为 0

令 $f(x), \hat{f}(x)$ 简写为 f, \hat{f}

\hat{f} 和 f 间 MSE $E[(y - \hat{f})^2] = E[((f - E[\hat{f}]) + \epsilon - (\hat{f} - E[\hat{f}]))^2]$

$$= (f - E[\hat{f}])^2 + E[\epsilon^2] + E[(\hat{f} - E[\hat{f}])^2]$$

$$= Bias^2 + Var^2 + \sigma^2$$

假设 ϵ , Bias, Variance 间 independent。交叉项乘积全部为 0

Bagging

从 Bootstrap Aggregating 得名

算法:

同时训练 n 个模型, 最终输出为 n 个模型的平均值

每一模型训练集为总训练集上 bootstrapping 采样结果

bootstrapping 采样: 从 m 大小的数据集上随机选取 m 大小数据集, 样本可重复

验证集 out-off-bag oob 样本: 当使用 bagging 选取时, 平均只有 $1 - e^{-1}$ 样本被选择, 余下样

本称 oob 样本, 为验证集

bagging 降低方差, 不降低偏差。方差较大时 bagging 降方差更加明显, 方差较大模型不稳定

证明降低方差:

令每一模型函数为 h , 合并结果 $\hat{f}(x) = E[h(x)]$ 。由 $E[x]^2 \leq E(x^2)$

真实模型函数为 f , 为固定函数, 非变量

$$(E[f - \hat{f}])^2 \leq E[(f - h)^2]$$

$$(E[f - E[h]])^2 \leq E[(f - h)^2]$$

$$f^2 - 2fE[h] + E[h]^2 \leq f^2 - 2fE[h] + E[h^2]$$

$$E[h]^2 \leq E[h^2]$$

针对决策树变化:

random patches 随机贴片: 对特征和训练集同时取子集进行训练

random subspace 随机子空间: 对特征取子集, 对整个总训练集进行训练

extra-trees 极度随机森林: '使用随机 t_k 而不使用最小化数据混杂度的 t_k '

kfeature importance 特征重要性: 对所有取 k 为判断条件的节点 N_i , 计算加权平均值 $\sum_i (S_i \text{ 熵增百分比})$

boosting

将多个较弱模型合并为较强模型, 每一模型使用之前模型预测效果较差样本作为训练集。降低偏差

AdaBoost: 对上一预测机制遗漏的样本加更高权重, 进行训练

gradient boosting:

第 t 时间合并的模型称 H_t

第 t 模型 h_t 使用训练集 $\{x_i, y_i - H_{t-1}(x_i)\}$ ，随后更新总模型 $H_t = h_t + \eta h_t$

模型应使用 regularization，否则较强模型导致拟合整个训练集

决策树使用 boosting：合并多模型即不同决策树输出取加权平均值

训练 sequential，耗时长，使用 XGBoost lightGBM 减少训练时间

stacking

同一样本经过不同模型，输出连接后经过全连接层。降低方差

multi stacking:

多层 stacking，每层输入为样本/前层连接结果

易于 overfitting，由于单层特征 stacking 已经为强模型

每层训练使用不同数据集。不同层通常使用相同模型，相同超参数，不共享权重

训练第 i 层时使用第 i 训练集，输入为 i 训练集经过前层输出（不经过全连接层），输出连接后经过全连接层进行预测

可对不同层输出预测取平均，降低偏差

可使用 k-fold 创建不同层样本

Neural Architecture Search 搜索神经网络结构

使用 RNN 生成定义网络结构的文本。类似强化学习，环境即数据集，决策即网络结构

One shot approach：训练极大模型，包含所有想要验证的模型

Differentiable Architecture Search：多层 categorical search，每层包含多个不同超参数模型 + softmax

softmax 对每一模型输出加权平均，输入的权重为可学习参数

训练时输出为层中模型结果加权平均，预测时仅选择权重最大的唯一模型计算输出

由于权重可学习，称 differentiable

DARTS

Scaling CNN：efficientNet 调整 CNN 参数时满足固定关系

Hypernetworks/Neural Processes：预测其他网络超参数的网络

network with memory：kanerva model，非 case by case 模型

2 反向传播公式推导

向量求导定义

$$\text{标量对向量求导: } \frac{ds}{dv} = \begin{bmatrix} \frac{ds}{dv_0} \\ \frac{ds}{dv_1} \\ \dots \\ \frac{ds}{dv_n} \end{bmatrix}$$

$$\text{向量对向量求导: } \frac{du}{dv} = \begin{bmatrix} \frac{du_0}{dv_0} & \dots & \frac{du_0}{dv_n} \\ \dots & \dots & \dots \\ \frac{du_m}{dv_0} & \dots & \frac{du_m}{dv_n} \end{bmatrix}$$

对 $z = xv = vx$ ， v 为向量， z 为标量 有：

$$\frac{dz}{dx} = v$$

对 $z = Wx$ 有：

$$\frac{dz}{dx} = W$$

$$\frac{dL}{dW} = \frac{dL}{dz} x$$

对 $z = xW$ 有:

$$\frac{dz}{dx} = W^T$$

$$\frac{dL}{dW} = x^T \frac{dL}{dz}$$

对 $z = x^T W x$ 有:

$$\frac{dz}{dx} = (W^T + W)x$$

对一层前向传播 $g(Z)$, $Z = XW + \vec{b}$

X 为前一层输入, 即前一层 $g(X)$ 矩阵

g 此处为广播操作, 对矩阵 Z 中每一元素求 activation 值

+ 此处为广播操作, 对每行 XW 加偏差

求权重斜率

单一权重值求导: $\frac{dL}{dw_{ij}} = \sum_k \frac{dL}{dZ_{kj}} X_{ki}$

对权重矩阵求导: $\frac{dJ}{dW} = X^T \frac{dL}{dZ}$

对前一层输出 X 求导

$$\frac{dL}{dX_{ij}} = \sum_k W_{jk} \frac{dL}{dZ_{ik}}$$

$$\frac{dL}{dX} = \frac{dL}{dZ} W^T$$

对偏差求导

$$\frac{dL}{db_i} = \sum_k \frac{dL}{dZ_{ki}}$$

$$\frac{dL}{db} = \vec{1}^T \frac{dL}{dZ}$$

激活函数求导

反向传播: $\frac{dL}{dx} = \frac{dL}{dg(x)} \cdot g'(x)$

g', \cdot 为按元素操作

sigmoid $g(z) = \frac{1}{1+e^{-z}}$ 求导: $g'(z) = g(z)(1 - g(z))$

tanh $g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$ 求导: $g'(z) = 1 - g(z)^2$

softmax $\hat{y}_i = \text{softmax}(z_i)$ $J = -\frac{1}{N} \sum_{i=1}^N \sum_{c=1}^C y_{ic} \log(\hat{y}_{ic})$,

求导: $\frac{dL}{dz} = \frac{1}{N}(\hat{y} - y)$

验证 gradient 计算

1. 得到斜率 $\frac{dL}{dW}$
2. 微调 $W + \epsilon$, 计算 $L(W)$, $\frac{dL}{dW} \approx \frac{L(W+\epsilon) - L(W-\epsilon)}{2\epsilon}$
3. 比较两 $\frac{dL}{dW}$ 值是否相近

3 linear regression 线性回归

平方代价函数: $J(\theta) = \frac{1}{|B|} \sum_{i=1}^{|B|} J^{(i)}(\theta) = \frac{1}{2|B|} \sum_{i=1}^{|B|} (\hat{y}^{(i)} - y^{(i)})^2$, 为所有样本误差的平均值

迭代: $\theta_i = \theta_i - \frac{\eta}{|B|} \sum_{j \in B} \frac{dJ^{(j)}(\theta)}{d\theta_i}$, 即对所有 sample 训练一次, 得到 label 差值, 对每一参数减 斜率 * 学习率 的平均值

当使用平方代价函数:

$$\theta_i = \theta_i - \frac{\eta}{|B|} \sum_{i \in B} x_i^{(j)} (x_1^{(j)} \theta_1 + x_2^{(j)} \theta_2 + \dots + \text{const} - y^{(j)}) = \theta_i - \frac{\eta}{|B|} \sum_{i \in B} x_i^{(j)} (\hat{y}^{(j)} - y^{(j)})$$

$$\text{const} = \text{const} - \frac{\eta}{|B|} \sum_{i \in B} (x_1^{(j)} \theta_1 + x_2^{(j)} \theta_2 + \dots + \text{const} - y^{(i)}) = \text{const} - \frac{\eta}{|B|} \sum_{i \in B} (\hat{y}^{(j)} - y^{(j)})$$

$$\text{对样本 } i \text{ 的偏导数向量为 } \nabla_{\theta} J^{(i)}(\theta) = \begin{bmatrix} x_1^{(i)} \\ x_2^{(i)} \\ \dots \\ 1 \end{bmatrix} (\hat{y}^{(i)} - y^{(i)})$$

交叉熵代价函数: $J(\theta) = \frac{1}{|B|} \sum_{i \in B} H(y^{(i)}, \hat{y}^{(i)})$

softmax 线性回归: 单层神经网络, 使用 softmax 函数得到分类, 使用 cross entropy 计算代价
过拟合问题

1. **权重衰减 regularisation**: 在代价函数中惩罚高权重的值, 尽可能使所有权重值减小

每一权重的正则化值都会影响其余所有权重的斜率, 由于正则化项为标量加入代价函数

L2 正则化代价函数 $= J(\theta) + \frac{\lambda}{2} \sum_{w \in W} w^2$, 即 $J(\theta) + \frac{\lambda}{2} * \text{所有权重的平方和}$ 。 λ 为超参数, 决定权重衰减的程度

$$\text{求导} = \frac{dL}{dw} + \lambda w$$

L1 正则化代价函数 $= J(\theta) + \lambda \sum_{w \in W} |w|$, 即 $J(\theta) + \lambda * \text{所有权重绝对值和}$ 。

$$\text{求导} = \frac{dL}{dw} + \lambda \text{sign}(w)$$

2. **丢弃法 dropout**

每一权重 (不包括 const) 有 p 的几率 $\theta' = 0$, 有 1-p 的几率 $\theta' = \frac{\theta}{1-p}$

常用 p = 0.5

仅将输出值部分设为 0, 不删除原权重

反向传播时被隐藏的神经元向上一层传斜率为 0

为了得到确切的值, 在测试模型时较少使用

取消 dropout 层后需将输出按元素全部乘 p 值, 否则测试时输出参数值为训练时 $\frac{1}{p}$ 倍

初始化参数

1. **MXNet 默认随机初始化**: 所有权重 $\sim N(0, 1)$ 的 normal distribution, 所有 const 取 0

2. **Xavier 随机初始化**: 对一全连接层, 输入个数 a, 输出个数 b, 则所有参数 $\sim U(-\sqrt{\frac{6}{a+b}}, \sqrt{\frac{6}{a+b}})$

预处理数据集

1. **特征标准化 standardization/z-normalization**: $x' = \frac{x - \mu}{\sigma}$, 即统计中 z 值

2. **离散值转换成指示特征**: 对于一个可取值为 A, B, C 的离散输入值, 转换成 3 个数值输入。即如果原输入为 A, 转换后 3 个数值输入为 1, 0, 0。原离散值为 B 则转换后为 0, 1, 0

3. **Min-max normalisation**: 将每一特征值拉伸为 [a, b], 取一特征的最值 X_{min}, X_{max} , 则拉伸后每一特征值 $X' = a + \frac{(X - X_{min})(b - a)}{X_{max} - X_{min}}$

存储 a, b, X_{min}, X_{max} 值可将数据逆操作得到原数据集

结构

- 将训练集分组, 每组 batch_size 个 sample。

- 对这个 batch 的数据进行向量化计算, 计算 loss, **斜率清零**, 计算斜率, 调用优化函数。

- 即每一 batch 使用相同的权重 偏差。一次训练一共历多次所有 sample, 一次遍历进行 $\frac{\text{sample_size}}{|B|}$

次向量化计算

隐藏层必定使用激活函数, 输出层可选使用激活函数

全连接神经网络常用神经元数同输入特征个数, 每层神经元个数相同

4 restricted Boltzman machine RBM 受限玻尔兹曼机

组成

单层可见全连接层：有可见层向量 v ，元素数 = 节点数 N_v 。对应输入数据

单层隐藏全连接层：有隐藏层向量 h ，元素数 = 节点数 N_h 。简化输入数据的特征， $N_h < N_v$

$$\begin{aligned} \text{定义可能性 } P(v, h) &= \frac{\exp(-E(v, h))}{\sum_{v, h} \exp(-E(v, h))} \\ &= \frac{\exp(-E(v, h))}{Z} \quad \text{令 } Z = \sum_{v, h} \exp(-E(v, h)) \end{aligned}$$

bernoulli - bernoulli RBM: v, h 每一元素 $\in \{0, 1\}$

$$\text{能量值 } E(v, h) = -a^T v - b^T h - h^T W v$$

$W \in \mathbb{R}^{N_h \times N_v}$ 为隐藏层权重

a, b 分别为两层 bias

$$\begin{aligned} \text{Theorem 1: } P(h|v) &= \prod_i \frac{\exp(b_i h_i + h_i W_i v)}{\sum_{h_i} \exp(b_i h_i + h_i W_i v)} \\ &= \prod_i P(h_i|v) \end{aligned}$$

W_i 为 W 第 i 行向量，即隐藏层第 i 节点参数

隐藏层间没有连接，则可分别计算每一 $P(h_i|v)$ 后相乘

$$\begin{aligned} \text{Theorem 2: } P(h_i = 1|v) &= \frac{\exp(b_i + W_i v)}{\exp(b_i * 0 + 0 * W_i v) + \exp(b_i * 1 + 1 * W_i v)} \\ &= \sigma(b_i + W_i v), \quad \sigma \text{ 即激活函数 sigma} \end{aligned}$$

$$\text{可得 } P(h = \vec{1}|v) = \sigma(Wv + b)$$

$$\text{同 theorem 2, } P(v = \vec{1}|h) = \sigma(W^T h + a)$$

Gaussian - bernoulli RBM: $v \in \mathbb{R}^{N_v}, h \in \{0, 1\}^{N_h}$

$$E(v, h) = \frac{(v-a)^T (v-a)}{2} - b^T h - h^T W v$$

Theorem: $P(v|h)$ 为 $v \sim \mathcal{N}(W^T h + a, I)$ 的 pdf

I 为 covariance matrix, 为 identity matrix

性质:

$$\text{free energy } F(v) = -\log(\sum_h \exp(-E(v, h)))$$

$$\text{则 } P(v) = \sum_h P(h, v)$$

$$\begin{aligned} &= \sum_h \frac{\exp(-E(v, h))}{Z} \\ &= \frac{\sum_h \exp(-E(v, h))}{Z} \\ &= \frac{\exp(-F(v))}{\sum_v \exp(-F(v))} \end{aligned}$$

当 RBM 不包含隐藏层时 $p_n(v) = \frac{\sum_h \exp(-E_n(v, h))}{Z_n}$, n 为可见层节点数

$$\begin{aligned} &= \frac{\prod_{i=1}^n \sum_{h_i} \exp(b_i h_i + h_i W_i v)}{Z_n} \\ &= \frac{\prod_{i=1}^{n-1} \sum_{h_i} \exp(b_i h_i + h_i W_i v) (1 + \exp(b_n + W_n v))}{Z_n} \\ &= p_{n-1}(v) \frac{Z_{n-1}}{Z_n} (1 + \exp(b_n + W_n v)) \end{aligned}$$

$$\text{代价函数 } J(v) = -\log(P(v)) = F(v) + \log(\sum_v \exp(-F(v)))$$

训练:

一次迭代中，输入样本为可见层向量 v_0 ，根据 $\text{pdf}P(h|v_0)$ 取 h_0

1. 由 h_t ，根据 $\text{pdf}P(v_{t+1}|h_t)$ 产生 v_{t+1}
2. 由 v_{t+1} ，根据 $\text{pdf}P(h_{t+1}|v_{t+1})$ 产生 h_{t+1}
3. 重复 1.2. 得到 h_T, v_T

使用近似学习，则求导为

$$\begin{aligned}\frac{dJ}{dW} &= [h_T v_T^T - h_0 v_0^T] \\ \frac{dJ}{da} &= [v_T - v_0] \\ \frac{dJ}{db} &= [h_T - h_0]\end{aligned}$$

5 convolutional neural network 卷积神经网络

互相关运算：

输入一个二维数组，和二维核 **kernel** 进行互相关运算，得到二维数组

二维核/卷积核/filter 过滤器：在输入数组上滑动，每次和二维数组矩阵一部分按元素相乘求和，作为输出矩阵的元素

二维卷积层：

将输入和卷积核做互相关运算，结果加上 **const** 作为输出

特征图：输出矩阵可看做是输入矩阵的表征，称特征图

感受野 receptive field：

对输出矩阵一元素 x ，所有可能影响其值的输入矩阵元素称感受野

感受野可能大于实际输入的矩阵边界

填充 padding：

在输入矩阵外侧添加全零元素，使得输出矩阵的维度增加，由于可用的感受野增加。

常使用奇数 **kernel**，添加 $\lfloor \frac{\text{kernel}}{2} \rfloor$ 的填充，使得输出矩阵和输入矩阵纬度一样

步幅 stride：

定义每次感受野向左/向下移动的纬度

dilation

在较大卷积核中仅选择某些矩阵元素赋权重进行卷积计算，其余权重为 0

多通道输入输出：

当输入的数据包含多个矩阵，即多通道输入，例：RGB 图像有 3 个输入通道

对 c_i 输入， c_o 输出的卷积层，**kernel shape** 为 $(c_o, c_i, \text{行数}, \text{列数})$

每一输入通道有唯一的 **kernel** (c_i , 行数, 列数) 对应，进行互相关运算后结果矩阵相加，作为一条输出通道的结果

多组 (c_i , 行数, 列数) 分别产生输出通道的结果矩阵，则有 c_o 条输出通道

池化层：

作用：1 为了防止当输入变化时，输出立即随之更改。2 减少计算量

池化窗口，同卷积层的感受野。限定某块输入被同时考虑，同样有 **stride**，可对输入 **padding**

1. 最大池化层：取池化窗口内最大的输入

2. 平均池化层：取池化窗口平均值

多输入通道间池化结果不相加，即 **输入通道数 = 输出通道数**

LeNet 卷积神经网络 (针对 (28, 28, 1) 形状的输出图片)

1.2 组 卷积计算层 激活函数层 池化层

6 输出通道卷积层 **kernel** 为 (5, 5), 2 填充

sigmoid 激活函数层 对每一元素做 **sigmoid**

最大池化层 窗口 (2, 2), 步幅为 2

16 输出通道卷积层 **kernel** 为 (5, 5)

sigmoid 激活函数层 对每一元素做 sigmoid

最大池化层 窗口 (2, 2), 步幅为 2

2.3 组全连接层

节点数 120, 84, 输出节点数。除输出层使用 sigmoid 激活函数, 即 120 84 节点层

将 (批量大小, 通道数, height, width) 看做 (批量大小, 通道数 * height * width) 处理

AlexNet 深度卷积神经网络:

除输出层和丢弃层, 全部使用 relu 做激活函数, 针对 (224, 224, 3) 形状图片

卷积部分

- 2 组 卷积层 + 最大池化层

```
nn.Conv2D(96, kernel_size=11, strides=4, activation='relu')
```

可不使用大卷积核

```
nn.MaxPool2D(pool_size=3, strides=2)
```

```
nn.Conv2D(256, kernel_size=5, padding=2, activation='relu')
```

```
nn.MaxPool2D(pool_size=3, strides=2)
```

- 3 卷积层 + 1 最大池化层, 高输出通道, 低卷积窗口

```
nn.Conv2D(384, kernel_size=3, padding=1, activation='relu')
```

```
nn.Conv2D(384, kernel_size=3, padding=1, activation='relu')
```

```
nn.Conv2D(256, kernel_size=3, padding=1, activation='relu')
```

```
nn.MaxPool2D(pool_size=3, strides=2)
```

全连接层部分

- 两 hidden layer 全连接层 使用丢弃法

```
nn.Dense(4096, activation="relu"), nn.Dropout(0.5)
```

```
nn.Dense(4096, activation="relu"), nn.Dropout(0.5)
```

```
nn.Dense(10) // 根据需求改变输出层节点, 原论文为 1000
```

除最后一全连接层, 所有参数都被分配到 2GPU 上, 特征通道数被平分

第三 conv 层和所有全连接层得到前层的所有输入进行计算

VGG 使用重复元素网络

VGG 基础块

数个 (3, 3)kernel 1 填充卷积层 + 1 个 (2, 2) 窗口 2 步幅最大池化层

卷积层 层数 通道数为超参数, 一 VGG 块中每一卷积层有相同通道数

VGG 神经网络由 数个 VGG 块 + 数个全连接层 组成

例: **VGG-11**

1. (1, 64) (1, 128) (2, 256) (2, 512) (2, 512) 5 层 VGG 块

(n, m) 代表此 VGG 块使用 n 层卷积层, 各有 m 通道

2. 3 层全连接层, 实现同 AlexNet 的全连接层部分

共 8 层卷积层 + 3 层全连接层, 所以称 VGG-11

NiN 神经网络

NiN 块

1 个自定义卷积层 + 2 层 (1, 1)kernel 卷积层, 3 层卷积层都不包含池化层

自定义卷积层可设置 kernel, 步幅, 填充。

(1, 1) 卷积层可设置通道数 (= 自定义层通道数), 其余固定为默认值

NiN 神经网络有多组 (NiN 块 + 池化层)

例: NiN 模型

- NiN 块部分

```
nin_block(96, kernel_size=11, strides=4, padding=0)
```

```
nn.MaxPool2D(pool_size=3, strides=2)
```

```
nin_block(256, kernel_size=5, strides=1, padding=2)
```

```
nn.MaxPool2D(pool_size=3, strides=2)
```

```
nin_block(384, kernel_size=3, strides=1, padding=1)
```

```
nn.MaxPool2D(pool_size=3, strides=2)
```

- 在 NiN 块部分结束后加入丢弃层

```
nn.Dropout(0.5)
```

- 转化为对应分类个数的输出

```
nin_block(10, kernel_size=3, strides=1, padding=1)
```

```
nn.GlobalAvgPool2D() // 全局平均池化层, 每一通道取矩阵所有元素的平均值
```

```
nn.Flatten() // 平均池化层结果即分类结果, flatten 仅用于改变 shape
```

GoogLeNet 含并行结构神经网络

Inception 块

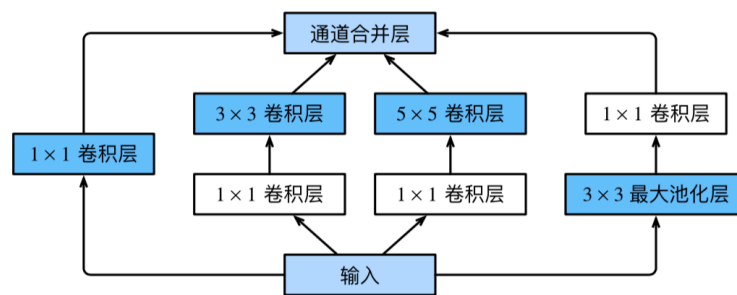


图 5.8: Inception 块的结构

inception 块结构表示: $(n_1, (n_{21}, n_{22}), (n_{31}, n_{32}), n_4)$

第一线路使用 n_1 通道

第二线路第一卷积层使用 n_{21} 通道, 第二层卷积层使用 n_{22} 通道 1 填充

第三线路第一卷积层使用 n_{31} 通道, 第二层卷积层使用 n_{32} 通道 2 填充

第四线路第一最大池化层使用 $(3, 3)$ 窗口 1 填充, 第二层卷积层使用 n_4 通道

每一卷积层都使用 relu 激活函数

所有层输出作为不同通道结果, 即最终有 $n_1 + n_{22} + n_{32} + n_4$ 通道

使用 inception 块使得参数量比使用卷积层参数量更小

GoogLeNet 结构:

5 个串联模块, 每一卷积层使用 relu 激活函数, 每一模块间使用 $(3, 3)$ 窗口 步幅 2 1 填充 池化层

1. 64 通道 $(7, 7)$ kernel 2 步幅 3 填充卷积层 + 模块间池化层

2. 64 通道 $(1, 1)$ kernel 卷积层 + 64 * 3 通道 $(3, 3)$ kernel 1 填充卷积层 + 模块间池化层

3. 串联 2 inception 块 + 模块间池化层, 分别有结构

$(64, (96, 128), (16, 32), 32)$, 通道比 2:4:1:1

(128, (128, 192), (32, 96), 64), 通道比 4:6:3:2

4. 串联 5 inception 块 + 模块间池化层,

(192, (96, 208), (16, 48), 64)

(160, (112, 224), (24, 64), 64)

(128, (128, 256), (24, 64), 64)

(112, (144, 288), (32, 64), 64)

(256, (160, 320), (32, 128), 128)

5. 串联 2 inception 块 + 全局平均池化层

(256, (160, 320), (32, 128), 128)

(384, (192, 384), (48, 128), 128)

6. 全连接层, 节点数和分类类别数相同

Inception-V2 (-BN): 使用 batch normalization

Inception-V3: 修改 inception 块

3. 中的所有 inception 块 (5, 5) 卷积核换为 2 个串联 (3, 3) 卷积核

4. 中 (3, 3) 换为 (7, 1)+(1, 7) 卷积核, (5, 5) 换为 2 组 (7, 1)+(1, 7) 卷积核

5. 中 (3, 3) 换为并行的 (3, 1) 和 (1, 3) 卷积核, (5, 5) 换为 (3, 3)+ 并行 (3, 1)+(1, 3) 卷积核

Inception-V4: 增加残差连接

6 CNN 优化方法

批量归一化 batch normalization

每一通道有拉伸 γ , 偏移 β 参数参与反向传播。初始化所有 $\gamma = \vec{1}, \beta = \vec{0}$

使用时无需和 dropout 合用, 由于使用 batch 的 mean/var 导致每次训练有来自样本的 noise, 和 dropout 一样对特征加入 noise。batch size 64-256 间, 保证计算的 mean variance 有 noise

批量归一化层有可变 μ', σ'^2 用于预测时进行批量归一。每一迭代由动量 η 更新

初始化 $\mu' = \vec{0}, \sigma'^2 = \vec{1}$

训练时 μ, σ^2 直接由输入 x 求得, 每一迭代更新 $\mu' = \eta\mu' + (1 - \eta)\mu, \sigma'^2 = \eta\sigma'^2 + (1 - \eta)\sigma^2$

variance 计算使用所有同一通道内的参数总体平均值, 不可使用 torch.var

$\mu'\sigma'^2$ 仅在预测时使用, 训练时使用 batch 本身的 mean var

1. 对全连接层做批量归一

处于 输入的仿射变换 和 激活函数 间, 即 输出 = $\phi(BN(x))$ BN 为批量归一计算

1. 对于批量仿射 $x = Wu + b$, 求标准化 $\hat{x}_i = \frac{x_i - \mu_i}{\sqrt{\sigma_i^2 + \epsilon}}$ 。 μ 和 σ^2 为 x mean var

即对不同样本同一特征向量元素求 mean 和 var

2. $BN(x) = \gamma * \hat{x} + \beta$, γ 拉伸 β 偏移, 和 x 形状相同。* + 为按元素加法乘法

2. 对卷积层做批量归一

处于 卷积计算 和 激活函数 间, 卷积计算 + 批量归一 + 激活函数 + 池化层

各通道独立计算, 各有独立拉伸 γ 偏移 β 值。

σ^2, μ 为此通道 所有样本内 的所有元素 的总体 var, mean

得到 σ^2, μ 后对此通道此批量内所有元素求标准化

ResNet 残差网络

使用原因:

解释 1: 由于网络输入常为 (0, 1) 区间的值, 使用乘法导致前向传输值不断减小。残差网络将两层有助于维持输出值大小

存在前向传输路径跨过多层网络, 有助于维持输出值大小

解释 2:

generic function class: 增加一网络层后不一定能够实现原先网络的预测 即增加网络层导致离正确网络结构更远。

大多数没有残差连接网络为此类

nested function class: 增加网络层后仍能做原先网络预测

残差块

训练时期望输出为 $f(x) - x$, 而非直接使用 $f(x)$ 期望输出。得到 $f(x) - x$ 后 $+x$ 得到 $f(x)$

1. 卷积层 (批量归一) + relu + 卷积层 (批量归一) 得到 $f(x) - x$

2. $[f(x) - x] + [(1, 1) \text{ 卷积层对 } (x) \text{ 卷积结果}] + \text{relu 激活函数}$

第一卷积层: (3, 3)kernel 1 填充 (通道数 步幅自定义)

第 234 组残差组第一残差块 第一卷积层步幅为 2, 否则为 1

第二卷积层: (3, 3)kernel 1 填充 1 步幅 (通道数自定义)

$+x$ 步骤的 (1, 1) 卷积层: (通道数 步幅自定义)

第 234 组残差组 第一残差块使用 (1, 1) 卷积层, 步幅为 2, 否则直接 $+x$

3 层卷积层通道数共享同一自定义值, **要求 2 层卷积层输入输出通道数一致**

ResNet-18 模型: 共 18 卷积层

1. 64 通道 (7, 7)kernel 2 步幅 3 填充 批量归一卷积层 + (3, 3) 窗口 2 步幅 1 填充最大池化层

2. 4 组 残差块, 每组包含多个残差块

第一组 2 个残差块 输出通道数和 1 中输出通道数一致

第二三四组 各 2 个残差块 输出通道数为前一层通道数 *2

3. 全局平均池化层 + 对应输出结果数全连接层

ResNext: 将 resnet 卷积层通道平分为并行通道, 减少计算量。

DenseNet 稠密连接网络

类似 ResNet 残差网络, $+x$ 步变为 concat x 连在输出结果后, 即 x 直接传向下一层

稠密块

多组 (批量归一 + relu + (3, 3)kernel 1 填充卷积层 + concat x) 卷积层通道数相同

concat 操作为在通道纬度的 concat, 即输入 x 作为额外输出通道。

增长率 = 输出通道 - 输入通道 = 卷积层通道数

过渡层

批量归一 + relu + (1, 1) 卷积层 + (2, 2) 窗口 2 步幅平均池化层

使用 (1, 1) 卷积层减小通道数, 2 步幅平均池化层减小矩阵大小

卷积层通道数 = 输出通道数 / 2

DenseNet 模型

1. 64 通道 (7, 7)kernel 2 步幅 3 填充 批量归一卷积层 + (3, 3) 窗口 2 步幅 1 填充最大池化层

2. 4 组稠密块, 由 3 个过渡层分隔

4 层稠密块卷积层数可以不相同

3. 批量归一 + relu + 全局平均池化层 + 对应输出结果数全连接层

Squeeze-Excite Net

每一通道计算权重，取 softmax。结果的 (1, 1, C) 和 C 个原通道参数按元素相乘。即 channel wise 的 attention 操作

ShuffleNet

1.group convolution: 使得卷积操作仅对部分 channel 进行操作。2. 交换 channel 顺序, 使得 group conv 得到其他 channel 信息

7 RNN 循环神经网络

记录数据状态, 根据以往状态和当前输入决定输出

仅仅使用前 n 个词进行预测会将误差叠加, 导致产生的输出最终趋向一固定值

n 阶马尔科夫链: 一个词的出现仅和前 n 个词有关

语言模型: 词序 (w_1, w_2, \dots, w_T) 的出现可能性为

$$P(w_1, w_2, \dots, w_T) = \prod_{t=1}^T P(w_t | w_{t-(n-1)}, \dots, w_{t-1})$$

称 n 元语法, 每一 w_t 为一时间步中出现的词

处理语言模型: 将每一文字转化为索引, 使用索引做训练参数集

one_hot 表示: 索引为 i 的词对应 one_hot 向量 $[v_0 = 0, v_1 = 0, \dots, v_i = 1, v_{i+1} = 0, \dots]$

采样方式:

B batch size

T 每个样本包含的时间步数

S 总时间步数

1 随机采样:

在 $[0, T)$ 范围内选择一起始位置 s , 每 T 长度分为一样本。共得到 $\lfloor \frac{S-s}{T} \rfloor$ 样本

随机选择 B 个样本作为一批量数据。不同批量间不重复使用样本

训练不同样本时不能将前一次隐藏层结果纳入计算

2 相邻取样:

同随机采样得到 $\lfloor \frac{S-s}{T} \rfloor$ 样本

选择批量中样本时, 不同批量的同一位置样本连续。即第 2 批量 0 位置的样本为第 1 批量 0 位置样本的下一样本

仅需在一 epoch 开始时初始化隐藏层结果, 而非在每一批量开始初始化

裁剪梯度

clip by value: 所有斜率限定在 $[v_{min}, v_{max}]$ 区间

clip by norm: 将所有参数斜率拼接成向量 g , 进行裁剪: $g' = \min(\frac{\theta}{\|g\|}, 1)g$

即所有斜率一次迭代中最多更改 θ

困惑度 perplexity

T 时间步困惑度 $= \exp(-\frac{1}{T} \sum_{i \in T} \log(P(n_i | n_1, \dots, n_{i-1})))$

n_i 为期望的输出词, $P(n_i | n_1, \dots, n_{i-1})$ 即 softmax 后期望输出词位置的 \hat{g}

即对每一样本, T 时长里每一次预测的 Softmax Caragorical 交叉熵损失值求和, 取 exp

求多迭代总困惑度: 交叉熵按时长 T 加权求和后取 exp 值

强制学习 teaching forcing: 对自然语言人工智能通用

当预测第 i 时间步 y_i 时使 y_{i-1} 为样本实际上一 label, 而不是上一次预测的计算结果 \hat{y}_{i-1}

auto regressive model

模型输出基于自身前一时点的输出, RNN 为一 auto regressive model

RNN 实现

模型:

输入 X_t , 输出 O : 时间步数 个 (批量大小, 词典大小) 矩阵

1. 隐藏层 $H_t = \phi(X_t W_{xh} + H_{t-1} W_{hh} + b_h)$

X_t 为上一时间步的词

H_{t-1} 项为隐藏层前一次输出, 称隐藏状态, 第一个时间步使用全零 H_{t-1} 矩阵

简化计算: $X_t W_{xh} + H_{t-1} W_{hh} = [X_t, H_{t-1}] [W_{xh}, W_{hh}]^T$

激活函数 ϕ 为 \tanh

2. 输出层 $O = H W_{hq} + b_q$

实现:

训练时输入一段词序 $[w_0, w_l]$, 输出即为预测 $[w'_1, w'_N]$ 的词

批量输入: (批量大小, 时间步数 l) 每一元素为单个词, 即一批量的句子前缀

输入转换为 时间步数 个 (批量大小, 词典大小) 矩阵, 第 i 矩阵第 j 行对应 批量中第 j 样本 第 i 时间步的词的 one_hot 向量

计算: 预测 T 长度的词序

时间步 t 从 1 开始

1. $t = 0$ 的隐藏状态 H_0 设为全 0 向量

2. $t \leq T$: 使用 (prefix 在 t 位置的词, 隐藏状态 H_{t-1}) 计算 (t 位置预测词, H_t)

输出词 $[w'_1, w'_T]$ 不一定和输入 $[w_1, w_T]$ 一致, 但在 $[1, T]$ 时间段内使用输入词 w_{t-1} , 而非上一预测词 w'_{t-1} 进行预测

即 forced learning。或称 warm-up, 过程中计算了隐藏状态

3. $t > T$: 使用 (上一预测词, H_{t-1}) 预测

训练时标签长度和 prefix 长度相同, 即到达 $t = T$ 时预测停止

代价函数对每组 预测词和目标词的 one_hot 使用 softmax 交叉熵, 求和/求平均

反向传播仍使用此批量 T 时长的总代价值, 不使用困惑度反向传播

通过时间反向传播

有关时间步的损失函数: $L = \frac{1}{T} \sum_{t=1}^T l(o_t, y_t)$ 此处为单一批量的 t 时间步输出

反向传播公式

针对单一样本, 假设隐藏层不使用激活函数

$$h_t = W_{xh} x_t + W_{hh} h_{t-1} + b_h$$

$$o = W_{hq} h_t$$

$$\frac{dL}{do_t} = \frac{1}{T} \frac{dl(o_t, y_t)}{do_t}$$

$$\frac{dL}{dW_{qh}} = \sum_{t=1}^T \frac{dL}{do_t} h_t^T$$

$$\frac{dL}{dh_T} = W_{qh}^T \frac{dL}{do_T}$$

$$\text{当 } t < T: \frac{dL}{dh_t} = W_{hh}^T \frac{dL}{dh_{t+1}} + W_{qh}^T \frac{dL}{do_t}$$

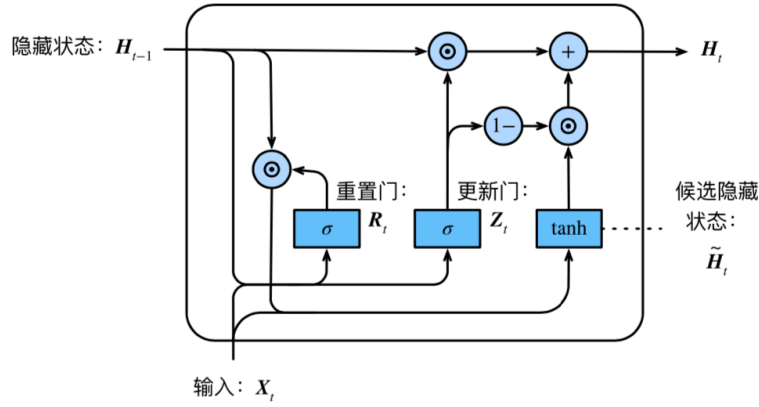
$$= \sum_{i=t}^T (W_{hh}^T)^{T-i} W_{qh}^T \frac{dL}{do_{T-i+1}}$$

$$\frac{dL}{dW_{hx}} = \sum_{t=1}^T \frac{dL}{dh_t} x_t^T$$

$$\frac{dL}{dW_{hh}} = \sum_{t=1}^T \frac{dL}{dh_t} h_{t-1}^T$$

GRU 门控循环单元

替代原计算隐藏状态方法, 应对梯度衰减



reset gate 重置门 update gate 更新门:

得到上一层隐藏层结果 H_{t-1} 当前时间步输入 X_t

重置门 $R_t = \sigma(X_t W_{xr} + H_{t-1} W_{hr} + b_r)$

更新门 $Z_t = \sigma(X_t W_{xz} + H_{t-1} W_{hz} + b_z)$

W 为权重参数 b 为偏差参数 σ 为 sigmoid 函数

候选隐藏状态

候选隐藏状态 $\tilde{H}_t = \tanh(X_t W_{xh} + (R_t \odot H_{t-1}) W_{hh} + b_h)$

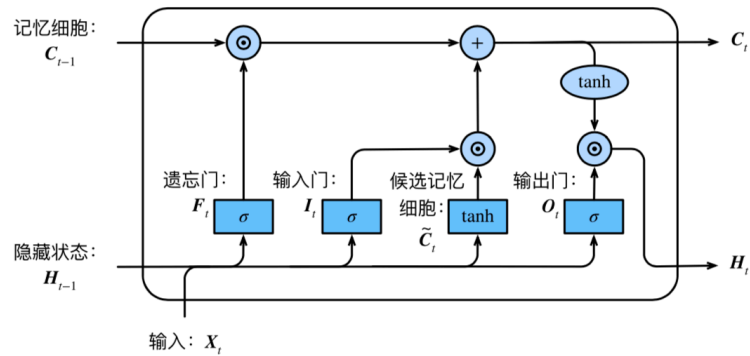
\odot 为按元素相乘, 使得重置门中对应位置值为 0 的元素被丢弃

隐藏状态

$H_t = Z_t \odot H_{t-1} + (1 - Z_t) \odot \tilde{H}_t$

输出公式不变: $O_t = H_t W_{hq} + b_q$

LSTM 长短期记忆门控循环网络



input gate 输入门 forget gate 遗忘门 output gate 输出门候选记忆细胞

输入门 $I_t = \sigma(X_t W_{xi} + H_{t-1} W_{hi} + b_i)$

遗忘门 $F_t = \sigma(X_t W_{xf} + H_{t-1} W_{hf} + b_f)$

输出门 $O_t = \sigma(X_t W_{xo} + H_{t-1} W_{ho} + b_o)$

候选记忆细胞 $\tilde{C}_t = \tanh(X_t W_{xc} + H_{t-1} W_{hc} + b_c)$

记忆细胞

$C_t = F_t \odot C_{t-1} + I_t \odot \tilde{C}_t$

隐藏状态

$H_t = O_t \odot \tanh(C_t)$

输出公式同 RNN/GRU 不变

深度循环神经网络

包含多个隐藏层的 RNN

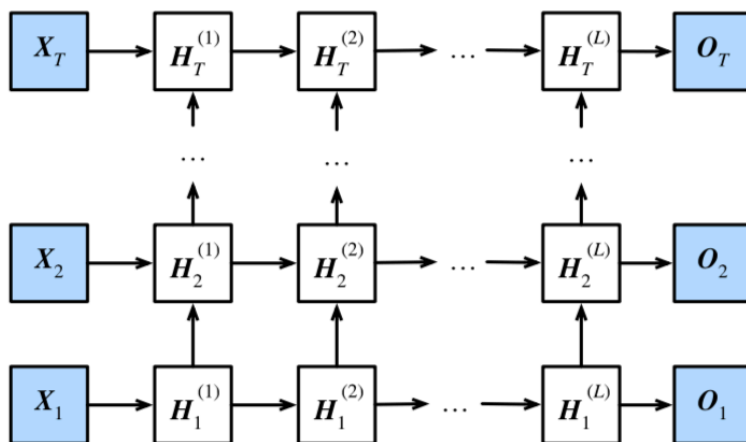


图 6.11: 深度循环神经网络的架构

结构:

隐藏层输出为 $H_i^{(1)} H_i^{(2)} \dots H_i^{(n)}$, 第 i 次 forward 对应图中一行, 即 $X_i \rightarrow H_i^{(1)} \rightarrow \dots \rightarrow H_i^{(n)} \rightarrow O_i$

第 i 次计算中:

$$l = 1 \text{ 隐藏层输出 } H_i^{(1)} = \phi(X_i W_{xh}^{(1)} + H_{i-1}^{(l)} W_{hh}^{(1)} + b_h^{(1)})$$

和单层隐藏层 RNN 的隐藏层输出公式一致

$$l = 2, 3 \dots n \text{ 隐藏层输出 } H_i^{(l)} = \phi(H_i^{(l-1)} W_{xh}^{(l)} + H_{i-1}^{(l)} W_{hh}^{(l)} + b_h^{(l)})$$

即 $l = 1$ 公式中将输入变为前一隐藏层输出

$$\text{输出层 } O_i = H_i^{(n)} W_{hq} + b_q$$

双向循环神经网络

允许神经网络根据前后的词序决定当前时间步的词

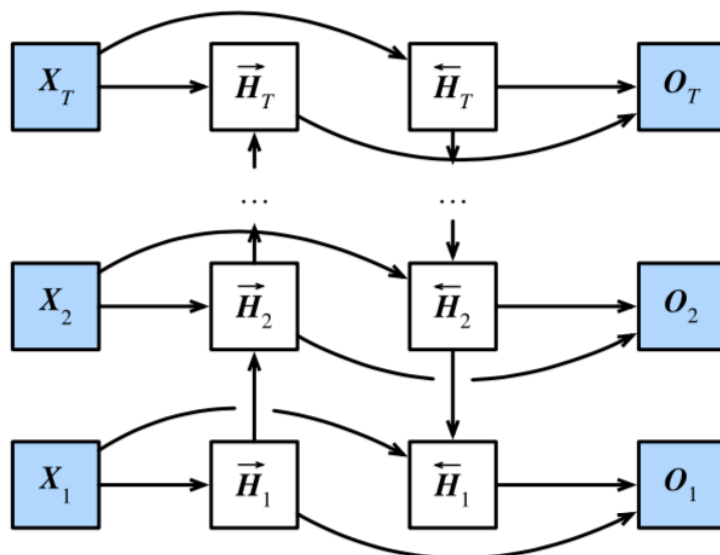


图 6.12: 双向循环神经网络的架构

结构:

仅有 2 隐藏层, 分为 正向隐藏层 反向隐藏层, 分别输出 $H^{(f)}$ $H^{(b)}$

第 i 次计算中

$$H_i^{(f)} = \phi(X_i W_{xq}^{(f)} + H_{i-1}^{(f)} W_{hh}^{(f)} + b_h^{(f)})$$

$$H_i^{(b)} = \phi(X_i W_{xq}^{(b)} + H_{i+1}^{(b)} W_{hh}^{(b)} + b_h^{(b)})$$

$$O_i = H_i W_{hq} + b_q$$

$$H_i = (H_i^{(f)}, H_i^{(b)}), \text{ 为正向 反向隐藏层的输出 concat}$$

双向 + deep RNN: 应 stack 多个双向 RNN 组成 deep bi-RNN, 非每一 direction 为 deep RNN

经验总结

1. 增加一批量中的时间步数有助于提高周期性函数的预测效果
2. ReLU 做激活函数在周期性函数上预测效果较好

8 代码算法优化

Stochastic Gradient Decent 随机梯度下降: 采样方法

每次迭代随机选择一个**样本**, 得到权重针对此样本的斜率进行梯度下降 而不是求参数针对所有样本的代价函数斜率。开销从 $O(n)$ 变为 $O(1)$

$$\text{即任取的样本 } j, \text{ SGD 迭代为 } \theta_i = \theta_i - \eta \frac{dJ^{(j)}(\theta)}{d\theta_i}$$

训练时 noisy

batch gradient descent 小批量随机梯度下降: 采样方法

每次迭代随机选择一组样本 B , 计算参数关于 B 的代价函数斜率

分组目的为减少 noisy, 并避免需要迭代过整个训练集才能计算斜率

重复采样: 允许 B 中重复出现同一样本。反之不重复采样,

随迭代次数增加, 学习率可减小, 使学习率和斜率的乘积方差减小

动量法: 迭代方法

解决固定一学习率值无法同时满足多个参数的学习率范围, 导致某些参数发生学习太慢, 某些参数不断越过最优解

定义:

向量 v_t 为第 t 次迭代每一参数的速度变量

向量 x_t 为第 t 次迭代参数向量

向量 g_t 为第 t 次迭代每一参数的斜率

迭代:

$$v_t = \gamma v_{t-1} + \eta_t g_t$$

$$x_t = x_{t-1} - v_t$$

adaptive learning rate

对每一参数使用不同 learning rate, 当参数更新较慢时增加学习率, 反之减少

以下迭代方法均为 adaptive learning rate 算法

AdaGrad 算法: 迭代方法

根据参数斜率调整学习率

定义向量 s_t 为第 t 次迭代累加变量, 累计每一参数的斜率平方和

迭代:

$$s_t = s_{t-1} + g_t \odot g_t$$

$$x_t = x_{t-1} - \frac{\eta}{\sqrt{s_t + \epsilon}} \odot g_t$$

即每一参数有变量记录所有斜率历史, 每次迭代时 学习率/斜率历史 l2 norm

由于累加变量斜率历史, 学习率始终降低, 造成学习缓慢

RMSProp 算法: 迭代方法

解决 AdaGrad 末期学习率过低问题

迭代:

$$s_t = \gamma s_{t-1} + (1 - \gamma) g_t \odot g_t \quad \text{即对 } s_t \text{ 按元素平方做指数加权}$$

x_t 同 AdaGrad 算法

AdaDelta 算法: 迭代方法

解决 AdaGrad 末期学习率过低问题, 不使用超参数

定义:

向量 g'_t 记录参数变化量

向量 Δx_t 对 g'_t 按平方做指数加权

迭代:

$$s_t = \gamma s_{t-1} + (1 - \gamma) g_t \odot g_t \quad \text{同 RMSProp}$$

$$\Delta x_t = \gamma \Delta x_{t-1} + (1 - p) g'_t \odot g'_t \quad \Delta x_t \text{ 初始化为全零向量}$$

$$g'_t = \sqrt{\frac{\Delta x_t + \epsilon}{s_t + \epsilon}} \odot g_t$$

$$x_t = x_{t-1} - g'_t \quad g'_t \text{ 即参数变化率}$$

Adam 算法: 迭代方法

定义:

v_t t 时间步的 动量变量

s_t t 时间步的 指数加权移动平均

g_t t 时间步 小批量随机梯度

$0 \leq \beta_1 < 1$ 所有 动量变量的权重超参数。常取 0.9

$1 \leq \beta_2 < 1$ 所有 指数加权移动平均变量超参数。常取 0.999

迭代:

$$v_t = \beta_1 v_{t-1} + (1 - \beta_1) g_t$$

$$s_t = \beta_2 s_{t-1} + (1 - \beta_2) g_t \odot g_t$$

$$\hat{v}_t = \frac{v_t}{1 - \beta_1^t}$$

$$\hat{s}_t = \frac{s_t}{1 - \beta_2^t}$$

使用 \hat{v}_t, \hat{s}_t 偏差修正

使得每一时间步 t 前的权重和为 1, 使得在时间步数较小时动量值不受权重影响

例: 不使用偏差修正: $v_1 = 0.1 g_1$

$$g'_t = \frac{\eta \hat{v}_t}{\sqrt{\hat{s}_t^2 + \epsilon}}$$

ϵ 常取 10^{-8}

$$x_t = x_{t-1} - g'_t$$

9 计算机视觉

提升模型泛化能力方法: 图像增广 图像微调

图像增广

对图像随机变换, 产生相似样本。扩大训练集

方法:

翻转: 有固定几率上下/左右翻转

裁剪: 随机裁剪原图 10%-100% 面积的图像, 宽高比 0.5-2, 并拉伸至固定像素大小

变化颜色: 调整亮度 色调 对比度 饱和度

微调

假设:

源模型包含的知识和目标模型紧密相连

源模型输出层不能直接用于目标模型

目标数据集远小于源数据集

方法:

1. 在源数据集上训练一神经网络, 称源模型
2. 创建一新神经网络, 称目标模型, 复制源模型除了输出层的所有结构和参数
3. 为目标模型添加输出层, 节点数对应输出种类, 随机初始化模型参数
4. 在目标数据集上训练目标模型, 微调 hidden layer 参数, 从头训练输出层参数

通过提高输出层学习率, 降低隐藏层学习率达到重新训练输出层, 微调隐藏层的目的。学习率

相差可达 1000 倍

目标检测 创建 匹配锚框算法

锚框: 标记一个物体的范围框

针对参数 $s = s_1, \dots, s_n, r = r_1, \dots, r_m$

为避免生成过多锚框, 将图片分为不同区域 grid, 以每一区域中心生成锚框。区域 grid 可重叠

每一 grid 生成 $(s_1, r_1), (s_1, r_2), \dots, (s_1, r_m), (s_2, r_1), \dots, (s_n, r_m)$ 共 $n + m - 1$ 个锚框

对 H, W 大小的图片, h 行 w 列个区域情况中, 每一区域中以 (s_i, r_j) 为参数的锚框有 高 = $s_i * \sqrt{r_j} * \frac{H}{h}$, 宽 = $\frac{s_i}{\sqrt{r_j}} * \frac{W}{w}$

`contrib.ndarray.MultiBoxPrior` 和 `d2l.torch.multibox_prior` 得到高 = $\frac{s_i}{\sqrt{r_j}}$ 宽 = $s_i * \sqrt{a_j}$

交并比 IoU: 两锚框 相交面积/相并面积

训练集: 每一锚框有对应的目标的标签 相对目标范围框的偏移

赋目标框

1. 对锚框组 A_1, A_2, \dots, A_n , 目标框组 B_1, B_2, \dots, B_m 定义矩阵 X ,

X 为 (n, m) , 包含每一锚框相对每一目标框的交并比

2. 找到 X 中值最大项 x_{ij} , 则 A_i 对应目标 B_j , 移除 X 中 i 行 j 列

3. 重复找到剩余矩阵中的最大项并移除, 最终有 m 锚框对应目标, 剩余 $n - m$ 锚框

4. 对每一未对应锚框 A_i , 寻找 X 中 i 行最大交并比, 若值大于阈值, 为 A_i 分配对应目标框, 若没有交并比大于阈值, 目标框为整个图片。

被赋予背景目标框的锚框称负类锚框, 否则为正类锚框

赋偏移量

对锚框 A_i 有坐标 + 长宽 (x_a, y_a, h_a, w_a) , 对应目标框 B_j 有 (x_b, y_b, h_b, w_b)

设置常数 $\mu_x = \mu_y = \mu_w = \mu_h = 0, \sigma_x = \sigma_y = 0.1, \sigma_w = \sigma_h = 0.2$

偏移量为 $(\frac{x_b - x_a}{\sigma_x}, \frac{y_b - y_a}{\sigma_y}, \frac{\log(\frac{w_b}{w_a}) - \mu_w}{\sigma_w}, \frac{\log(\frac{h_b}{h_a}) - \mu_h}{\sigma_h})$

非极大值抑制: non-maximum suppression

在显示结果阶段, 去除重复对一个物体分类的锚框。在训练结束后输出目标时使用, 训练时直接使用所有锚框和标签集对比

锚框置信度: 一个锚框 A_i 针对所有目标锚框 B 计算概率, A_i 最大概率符合的目标锚框对应 A_i 的预测类别, 此最大概率为 p 。概率不是锚框和目标锚框的交并比, 而是输出层输出的 每一锚框与每一类别 的可能性

算法:

1. 计算每一锚框置信度, 选取最高的锚框 A_i

2. 将所有和 A_i 交并比高于阈值的锚框删除

重复 1.2. 步, 直至没有锚框剩余

多尺度目标检测: 仅适用部分像素点作为锚框的中心, 减少锚框数量, 减少计算

SSD 单发多框检测: 一种目标检测算法

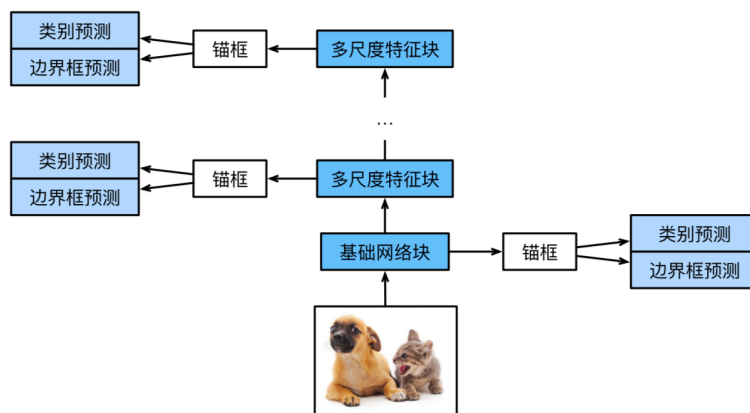


图 9.4: 单发多框检测模型主要由一个基础网络块和若干多尺度特征块串联而成

类别预测层:

得到卷积网络的张量, 形状 (批量大小, 通道数, 高, 宽), 处于同一高宽坐标的元素对应原图片中一区域 **grid** 的像素, 即原图区域为此元素的感受野。此区域中可包含多个锚框。

假设锚框中心像素有 (h, w) 个, 每个中心点生成 a 个锚框, 每个锚框需匹配进 $(q+1)$ 个分类, 多 1 分类对应负类分类

将锚框看做多层, 每层 $h * w$ 个, 共 a 层, 层编号为 c_1, c_2, \dots, c_a

方法 0: 对每一锚框使用全连接层分类, 即每一像素输入全连接网络, 最终有匹配类别个数的输出层, 判断锚框分类

造成参数过多, 无法训练

方法 1: 使用卷积层通道数输出类别

包含一个保持输入高 宽的卷积层, 如 1 填充 (3, 3)kernel 的卷积层

输出通道数为 $a * (q+1)$, 每一通道输出仍为矩阵, 第 $(i-1) * (q+1) + j$ 通道包含锚框层 c_i 中所有锚框对分类 j 的匹配可能性。由于输入张量每一元素对应一区域 **grid**, 卷积层输出 (h, w) 位置的值的的可能性即为此锚框层属于 (h, w) 区域的锚框对一类别的预测可能性

边界框预测层:

输入结构同类别预测层, 输出通道数为 $a * 4$, 4 对应一个锚框有 4 个偏移量, 每一通道输出仍为矩阵, 第 $(i-1) * (q+1) + j$ 通道包含锚框层 c_i 中所有锚框的第 i 号偏移量

CNN 预测偏移量, 而非锚框边界值

连接多尺度预测:

将不同多尺度特征块和基础网络块产生的 类别 边界预测层结果 合并
(批量大小, 通道数, 高, 宽) 转为 (批量大小, 通道数 * 高 * 宽)

高宽减半块:

1. 两组 [(3, 3)kernel 1 填充 卷积层 + 批量归一化 + relu]
2. (2, 2) 窗口 2 步幅 最大池化层

基础网络块

串联 3 块高宽减半块, 分别有 16, 32, 64 通道

整体模型

结构对应图例, 每一模块 l_i 对结果张量 (批量大小, h_i, w_i) 有 类别 边界预测层

1. 基础网络块
- 2-4. 高宽减半块, 通道数都为 128
5. 仅包含全局最大池化层 高宽降至 1

得到输出集

5 层 类别 边界预测结果分别通过 **concat** 得到总体类别预测 边界预测, 用于计算代价。

类别分类 **concat** 结果: (批量大小, $\sum_i h_i * w_i * \text{锚框数}$, 类别数 + 1)。

边界预测 **concat** 结果: (批量大小, $\sum_i h_i * w_i * \text{锚框数} * 4$), 4 对应偏移量个数

类别预测使用 **SoftmaxCrossEntropyLoss**, 边界预测使用 **L1Loss**。由于所有坐标通过百分比表示, 可直接和标签集一一对应, 不受卷积层对矩阵大小影响

得到标签集

对于 5 层中所有使用过的锚框, 使用**目标检测**章节的算法对每一锚框添加标签 偏移量。

结果得到 3 个张量:

每一锚框的类别, 此张量所有元素和输出的类别预测进行交叉熵代价计算

锚框数 *4 大小的 bitmap，代表哪些锚框的偏移量被纳入代价函数计算。即将此 bitmap 和第三张量按元素相乘 结果取平方和，为偏移量代价值

锚框数 *4 大小的张量，代表每一锚框的偏移量。负类锚框的偏移量被忽略

针对同一网络，不同迭代中产生的锚框位置 大小不变。由于目标框计算基于交并比，偏移量计算值不变。计算偏移量代价作用为：检查网络对偏移量的预估是否正确

Two-stage object detection

包含 region proposal 得到锚框位置，classifier 对每一锚框得到标签。相较 SSD 耗时长，精度高

selective search: region proposal 方法

将有类似颜色/材质/亮度的区域合并，合并后区域对应一锚框

R-CNN: selective search + SVM classification

Fast R-CNN: selective search + CNN classification

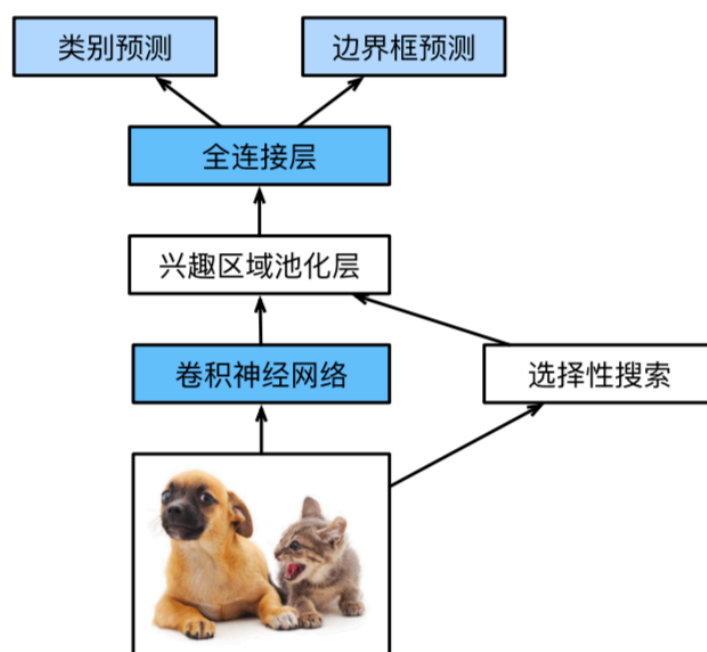


图 9.6: Fast R-CNN模型

1. CNN 卷积层，负责抽取特征。输出通道数为 c
2. 使用 selective search 得到 n 个锚框。输入仍为原图像，非 1. 中结果
3. ROI Pooling 兴趣区域池化层

根据 2. 中锚框 (h_1, w_1) ，从 1. 中选取锚框对应位置的特征作为 ROI Pooling 输入。

ROI Pooling 将不同锚框 i 的输出 (c, h_i, w_i) 统一为 (n, c, h, w) 形状输出

ROI Pooling 例：
$$\begin{bmatrix} x_{11} & x_{12} & x_{13} \\ x_{21} & x_{22} & x_{23} \\ x_{31} & x_{32} & x_{33} \end{bmatrix}$$
 可取 $\begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \end{bmatrix}$ 和 $\begin{bmatrix} x_{31} & x_{32} \end{bmatrix}$ 做输出 y_{11} 和 y_{21} 的来源

4. 将 3. 中结果 flatten，经过全连接层。输出 (n, d) ， d 为超参数。即每一锚框有特征向量，长度为 d

5. classification 层: 包含 2 个全连接层, 输入同为 4. 中结果

预测类别, 对 n 个锚框预测 q 种类别, 形状为 (n, q)

预测边界框, 对 n 个锚框分别计算 4 个偏移量, 形状为 $(n, 4)$

代价函数: 对每一锚框特征向量: $L_{cls}(\hat{y}, y) + \lambda(\vec{1} \cdot L_{loc}(t, t^*))$

L_{cls} 为分类代价值, 对锚框中特征分类

L_{loc} 为边界框预测代价函数, 进一步优化锚框位置。 t_i, t_i^* 为此锚框预测/期望的偏移量

Faster R-CNN: regional proposal network (CNN) + CNN classification

所有部分同 Fast R-CNN, 除了 2. 中 selective search 由 RPN 替代

Regional Proposal Network RPN: region proposal 方法

得到 1. 中结果 (), 得到是否为一可能的锚框位置 锚框的边界

代价函数: $\sum_i^n L_{cls}(p_i, p_i^*) + \lambda \sum_i^n \vec{1} \cdot L_{loc}(t_i, t_i^*)$

输出仅为锚框位置, 锚框大小由边界框定义, 类别由随后 classification 决定

L_{cls} 为 classification loss, 判断锚框 i 是否包含一物体

L_{loc} 为边界框预测代价函数, 算法同 5. 中 L_{loc}

图像分割

区分图片中一个物体的不同部分, 不将每一物体和标签对应

实例分割

区分一个像素属于哪一物体, 即使两物体为同一标签下的实例

semantic segmentation 语义分割

图像预处理: 不使用拉伸, 仅适用随机图片裁剪, 裁剪得到图片一块固定形状的小图, 作为训练集

FCN 全卷积网络

结构:

1. 使用卷积神经网络 抽取图像特征

2.(1, 1) 卷积层 通道数变为类别个数

3. 转置卷积层 将图像的高宽变为输入图像的尺寸, 作用仅为将 2. 中的图片分类反卷积得到图片表示, 所以输入通道数 = 输出通道数 = 类别个数

例:

1.ResNet-18 预先训练神经网络抽取图像特征 丢弃最后全局平均池化层 + 全连接层

2.(1, 1) 卷积层, 通道数 = 类别数 n

3. 转置卷积层将 (1, 1) 卷积层每一通道的输出转换回图片大小的矩阵, 每一元素为类别对应像素的 index

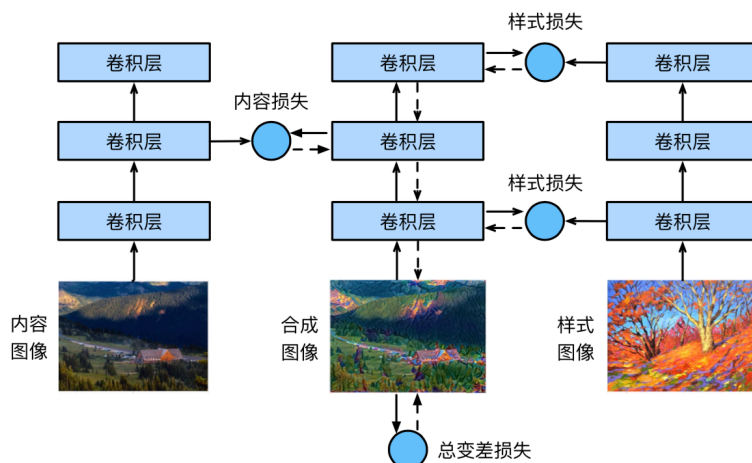
输出图片:

1. 对一个图片的 n 个通道中结果 合并为唯一矩阵, 新矩阵每一元素 = n 个矩阵中对应位置最大值

即 输出的 n 个矩阵代表图像符合每一类别的区域, 对一个位置的元素取最大值即判断此位置最优先属于哪一类别

2. 替换 1. 中输出的矩阵元素, 每一元素替换为此位置的 RGB 向量

样式迁移



得到 内容图像 样式图像，输出合成图像

图像预处理：将输入图片在 RGB 3 通道内分别做标准化

图像后处理：将输出矩阵标准化的值 转回 像素值

抽取特征：

使用 VGG-19 预训练参数抽取图像特征

选择 VGG 靠近输出的层，称内容层，避免保留过多内容图像的细节

选择多个中间层的输出匹配样式，称样式层

损失函数

损失函数为内容 样式 总变差损失函数的加权和

内容损失函数：

使用平均平方代价函数，计算 合成图像 和 内容图像 在 内容特征上的误差

输入张量为 内容图像和合成图像 通过抽取特征后内容层的输出

样式损失函数：

输入张量为样式层输入

对每一样式层：

1. 将 c 通道 (h, w) 的样式层输出转为 (c, hw) 的矩阵 X
2. 计算 X 的 Gram matrix, XX^T 。元素 $(XX^T)_{ij}$ 即 $x_i \cdot x_j$ ，代表通道 i 和 j 的相关性。可

对 Gram matrix 每一元素先除以 $|X|$ 元素值，避免样式损失值过高

即将 c 通道的矩阵两两按元素相乘 求和，得到 $c \times c$ 矩阵

3. 传入样式图像 合成图像，对每一样式层输出计算 Gram matrix 的平均平方代价函数，每层代价相加为总样式损失值

总变差损失：

输入仅为合成图像，对合成图像使用总变差降噪

$$J = \sum_{i,j} |x_{i,j} - x_{i+1,j}| + |x_{i,j} - x_{i,j+1}|$$

即 对每一合成图片的像素，求其和右方 下方像素的差值。最后求和

样式迁移模型：

使用每一 VGG 块的第一卷积层做样式层，第四卷积块的最后一卷积层做内容层

更改的参数仅有生成的图片，不对网络参数进行更改

10 自然语言识别

词编码

bag of word BOW: 矩阵, 每列对应一词, 每行对应一句话, 每一元素对应词是否出现在句中
矩阵稀疏

map to concept: 将词和实际语义对应, 如 WordNet。

无法处理一词多意问题

Feature extraction: 词嵌入

一. 词嵌入

词赋予两个向量 v_o, v_c , 分别为将词作为中心词和作为背景词使用的词向量

词间夹角 $\frac{u_o^T v_c}{\|u_o\| \|v_c\|} \in [-1, 1]$ 为词的相似度

创建方法:

Count-based: co-occurrence matrix

Prediction-based: word2vec, GLOVE, BERT

定义:

S : 词序集合, 包含多个词序, 包含重复的词

\mathcal{V} : 所有词的索引集合, 即 $nub(S)$

w_i : 一个词, 索引为 i

v_i : 词 w_i 对应的实数向量

向量夹角余弦值 $\frac{x^T y}{\|x\| \|y\|}$ 代表两个词的关联程度

背景窗口大小 n : 定义对一个中心词 w_i , 背景词取值范围 = [中心词 $-n$, 中心词 $+n$]。

范围为词序内位置关系, 不是索引序号关系。背景词和中心词必须在同一词序内

当中心词一侧背景词个数 $l < n$, 填充 $n - l$ 个无效词, 随后由掩码舍弃不进入代价函数计算。

给定中心词 w_c , 有单个背景词 w_o 的概率 $P(w_o|w_c) = \frac{\exp(u_o^T v_c)}{\sum_{i \in \mathcal{V}} \exp(u_i^T v_c)}$

$\log(P(w_o|w_c)) = u_o^T v_c - \log(\sum_{i \in \mathcal{V}} \exp(u_i^T v_c))$

co-occurrence matrix

矩阵行列数同词个数 $\|\mathcal{V}\|$

每一元素 (x, y) 标记是否存在一句话同时包含词 x y

或记录同时包含词 x y 的句子个数

矩阵为 symmetric matrix, 词向量即矩阵中词所在的行 (列)

TF-IDF(w, d): 词 w 在词序 d 中的权重, 避免常用词和其余词相似度过高

$TF - IDF(w, d) = TF(w, d) \times IDF(w, d)$

$TF(w, d)$: 词 w 在词序 d 中出现的次数

$IDF(w, d) = \log(\frac{1}{w})$

跳字模型:

基于一个中心词生成其周围的多个背景词

似然函数: 代表一段词序发生的可能性

对中心词 w_c , 产生所有窗口内的背景词的概率为 (假设背景词间 independent) $P(c) = \prod_{c-n < j < c+n, j \neq c} P(w_j|w_c)$

(假设每一中心词产生背景词的可能性 independent) 产生整个词序 S 的可能性:

$P(S) = \prod_{0 \leq c \leq |S|} P(c) = \prod_{c=0}^{|S|-1} \prod_{j=c-n, j \neq c}^{c+n} P(w_j|w_c)$

代价函数: 基于似然函数, 高似然代表低代价函数值

$$J = -\log(P(S)) = -\sum_{c=0}^{|S|-1} \sum_{j=c-n, j \neq c}^{c+n} \log(P(w_j|w_c))$$

$$\text{求导: } \frac{dJ}{dv_c} = -\sum_{c=0}^{|S|-1} \sum_{j=c-n, j \neq c}^{c+n} \frac{d\log(P(w_j|w_c))}{dv_c}$$

$$\text{对于每一中心词 } v_c \text{ 有 } \frac{dP(w_j|w_c)}{dv_c} = u_o^T - \sum_{j \in \mathcal{V}} P(w_j|w_c) u_j$$

连续词袋模型

与跳字模型不同处：中心词由背景词产生，和跳字模型相反

给定背景词 $W_o = w_{o_1}, w_{o_2}, \dots, w_{o_{2m}}$ ，中心词 w_c 有出现概率：

$$P(w_c|w_{o_1}, \dots, w_{o_{2m}}) = \frac{\exp(\frac{1}{2m} u_c^T (v_{o_1} + \dots + v_{o_{2m}}))}{\sum_{i \in \mathcal{V}} \exp(\frac{1}{2m} u_i^T (v_{o_1} + \dots + v_{o_{2m}}))}$$

$$= \frac{\exp(u_c^T \bar{v}_o)}{\sum_{i \in \mathcal{V}} \exp(u_i^T \bar{v}_o)}$$

1. 使用 **posterior estimate** 计算得到，忽视 $P(w_c)$ 项

2. $\frac{1}{2m}$ 项为额外添加，使 $\bar{v}_o = \frac{(v_{o_1} + \dots + v_{o_{2m}})}{2m}$

3. 背景词使用背景向量 w_o 而非 w_c

似然函数：

$$\prod_{c \in \mathcal{V}} P(w_c|W_o)$$

代价函数：基于似然函数，高似然代表低代价函数值

$$J = -\log(\prod_{c \in \mathcal{V}} P(w_c|W_o))$$

$$= -\sum_{c \in \mathcal{V}} (u_c^T \bar{v}_o - \log(\sum_{i \in \mathcal{V}} \exp(u_i^T \bar{v}_o)))$$

$$\text{求导: } \frac{dJ}{dv_{o_i}} = -\frac{1}{2m} (u_c - \sum_{j \in \mathcal{V}} P(w_j|W_o) u_j)$$

近似训练：对跳字模型和词袋模型的优化

避免每次梯度计算都包含词典大小的项数计算

负采样

定义：

背景词 w_o 出现在 w_c 背景窗的概率为 $P(D = 1|w_c, w_o) = \sigma(u_o^T v_c)$

σ 为 sigmoid 激活函数， $D = 1$ 代表事件发生， $D = 0$ 代表未发生

更改 $P(w_o|w_c)$ 定义： $P(w_o|w_c) = P(D = 1|w_c, w_o) \prod_{k=1, w_k \sim P(w)}^K P(D = 0|w_c, w_k)$

根据 $P(w)$ 分部取 K 个反样本，称噪声词，反样本不能为背景词

根据 word2vec 论文，选择每一噪声词 w 的几率为 w 出现几率的 0.75 次方

加入反样本的原因：若最大化 $P(D = 1|w_c, w_o) = \sigma(u_o^T v_c)$ ，则所有词向量方向相同且长度极大

代价函数： $-\log(P(w_o|w_c)) = -\log(\sigma(u_o^T v_c)) - \sum_{k=1, w_k \sim P(w)}^K \log(1 - \sigma(u_k^T v_c))$

$$= -\log(\sigma(u_o^T v_c)) - \sum_{k=1, w_k \sim P(w)}^K \log(\sigma(-u_k^T v_c)) \text{ 层序 softmax}$$

对整个词典有二叉树，每一分支节点 i 有背景词向量 u_i ，每一叶节点对应一词

定义：

$L(w)$: w 在二叉树中的深度，包括根节点和叶节点

$n(w, i)$: 从根节点到 w 叶节点的路径上第 j 个节点，有背景词向量 $u_{n(w, i)}$

w 在此为背景词，非中心词

一节点只有唯一背景词向量，当一节点出现在多个路径上时背景词向量被共享

更改 $P(w_o|w_c) = \prod_{j=1}^{L(w_o)-1} \sigma(isLeftChild(n(w_o, j+1))) \cdot u_{n(w_o, j)}^T v_c$

$$isLeftChild(x) = \begin{cases} 1 & x < 0 \\ -1 & otherwise \end{cases}$$

对固定中心词 w_c ，所有词的几率和为 $\sum_{o \in \mathcal{V}} P(w_o|w_c) = 1$

证明：选择任意节点 k 使得其左右子节点 i, j 都为叶节点

$$\begin{aligned} P(w_i|w_c) + P(w_j|w_c) &= \prod \dots \cdot (\sigma(x) + \sigma(-x)) \\ &= \prod \dots \cdot 1 \end{aligned}$$

即，任意仅有 2 个子节点的节点，子节点可能性和都为 1。循环合并子节点，最终得到跟节点可能性为 1

二次采样

作用：对于一个词 w_0 ，和低频词同时出现的情况 比 和低频词同时出现 对模型训练更加有用
实现：

取样的背景词中每个词有 $P(w_i) = \max(1 - \sqrt{\frac{t}{f(w_i)}}, 0)$ 几率被丢弃

$f(w_i)$ = 出现 w_i 个数 / 总词数，即词 w_i 在整个数据集中出现频率

t 为超参数

对每个 $w \in S$ ，使用二次采样随机丢弃 w ，创建新的 S' 作为训练集

word2vec 词嵌入模型实现

目标：得到每一词的词向量，使有关联的词间向量夹角最小

嵌入层

得到词嵌入的层，输入词 w_i 索引 i ，输出权重矩阵第 i 行作为词向量

嵌入层有权重矩阵，形状（词典大小，每个词向量纬度）

前向计算

输入：中心词索引矩阵 $\begin{bmatrix} [c_1] \\ \dots \\ [c_b] \end{bmatrix}$ + (背景词, 噪声词) 索引矩阵 $\begin{bmatrix} [o_{11}, \dots, o_{1n}, q_{11}, \dots, q_{1k}] \\ \dots \\ [o_{b1}, \dots, o_{bn}, q_{b1}, \dots, q_{bk}] \end{bmatrix}$

b 为批量大小

n 为窗口大小

k 为噪声词数量

两个输入矩阵元素皆为常数，非向量

1. 通过嵌入层变换为中心词向量张量 ($b, 1$, 词向量长) $\begin{bmatrix} [u_{c1}] \\ \dots \\ [u_{cb}] \end{bmatrix}$ 背景噪声词向量张量 ($b, n+k$,

词向量长) $\begin{bmatrix} [v_{o_{11}}, \dots, v_{o_{1n}}, v_{q_{11}}, \dots, v_{q_{1k}}] \\ \dots \\ [v_{o_{b1}}, \dots, v_{o_{bn}}, v_{q_{b1}}, \dots, v_{q_{bk}}] \end{bmatrix}$

2. 对两张量小批量相乘，即得到 ($b, 1, n+k$) 张量 $\begin{bmatrix} u_{c1}^T [o_{11}, \dots, o_{1n}, q_{11}, \dots, q_{1k}] \\ \dots \\ u_{cb}^T [o_{b1}, \dots, o_{bn}, q_{b1}, \dots, q_{bk}] \end{bmatrix}$

乘法为矩阵乘法，行向量 u_{c1}^T 乘 矩阵 $[o_{11}, \dots, o_{1n}, q_{11}, \dots, q_{1k}]$

输出张量 ($b, 1, k+n$)，每一元素为背景词 噪声词和中心词的点乘，对应计算词向量夹角

代价函数

使用 Sigmoid 二元交叉熵函数，传入前向传播结果 P ，label L 代表样本是否为正样本，掩码 M 代表样本是否有效

目标为使预测结果 P 每一有效位置和 L 对应位置相同

Feed-forward Language Model FFLM

为一种 n 元语法, 得到前 n 个词, 将每一词转为词向量相连。连接后向量经过全连接层得到输出

fastText 子词嵌入: 基于 word2vec 优化

产生子词: 在单词首尾添加 $< >$, 取所有长度为 s 的子字符串 + 单词本身

如, 对单词 'where' 有子词 $\{'<wh', 'whe', 'her', 'ere', 're>', '<where>'\}$

每一子词都存在词典中, 单个词对应的向量为子词向量和

GloVe 全局向量的词嵌入: 基于 word2vec 优化

对每一中心词 w_i , 合并所有 w_i 的背景词, 称多重集 \mathcal{C}_i 。多重集直接合并原集合, 不删除重复项

一个背景词 w_j 在 \mathcal{C}_i 中的出现次数称 重数, 记做 x_{ij}

定义 $x_i = |\mathcal{C}_i|$, 为多重集集合大小

代价函数 1: $J = -\sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{V}} x_{ij} \log(P(x_j|x_i))$

$$= -\sum_{i \in \mathcal{V}} x_i \sum_{j \in \mathcal{V}} \frac{x_{ij}}{x_i} \log(P(x_j|x_i))$$

第二层求和即 实际样本中 w_j 出现在 w_i 背景集合中的概率 $\frac{x_{ij}}{x_i}$ 和 模型预测的出现概率 $P(x_j|x_i)$ 的交叉熵

* 较少使用, 计算开销较大。并包含大量生僻词, 导致预测准确率较低

代价函数 2: 使用平方代价函数 $J = \sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{V}} h(x_{ij})(u_j^T v_i + b_i + c_j - \log(x_{ij}))^2$

$h(x_{ij})$ 为权重, 在 $[0, 1]$ 单调递增

例:

$$h(x) = \begin{cases} 0 & x = 0 \\ (x/c)^{0.75} & x < c \\ 1 & x > c \end{cases}$$

c 可取 100

b_i 为中心词偏差项, c_j 为背景词偏差项

$u_j^T v_i$ 即 x_{ij} , $\log(x_{ij})$ 即 $P(x_j|x_i)$ 除去分子

GloVe 每一词的背景词向量和中心词向量相同, 由于每对词互为背景词。最终每一词的背景词向量 = 中心词向量 = 所求的 $u_c + u_o$

预测 $P(w_i|w_j)$ 公式不变

证明代价函数

针对词向量, 定义函数 $f(u_i, u_j, u_k)$ 使得对一个词 w_k , 其余两个词相对出现的几率 $\frac{P(w_i|w_k)}{P(w_j|w_k)} \approx f(u_i, u_j, u_k)$

函数 $f(u_i, u_j, u_k)$ 可定义为标量函数 $g((u_i - u_j)^T u_k)$

$g(x)$ 需满足 $g(x)g(-x) = 1$, 则 g 可取 $g(x) = \exp(x)$

则需要 $\frac{\exp(u_i^T u_k)}{\exp(u_j^T u_k)} \approx \frac{P(w_i|w_k)}{P(w_j|w_k)}$

则 $P(w_i|w_k) = \frac{x_{ki}}{x_k} = \alpha \cdot \exp(u_i^T u_k)$

$$u_i^T u_k = \log(\alpha) + \log(x_{ki}) - \log(x_k)$$

中心词偏差项 b_k 和背景词偏差项 c_i 之和 $b_k + c_i$ 模拟 $-\log(\alpha) + \log(x_k)$

则最小化 $u_i^T u_k + b_k + c_i - \log(x_{ki})$, 使用平方代价函数

求近义词: word2vec 应用

使用 KNN, 在训练结束的词向量中寻找余弦值最小的 K 个向量

求类比词: word2vec 应用

给定词 w_a, w_b, w_c , 求 w_d 使得 $w_a : w_b$ 关系类似 $w_c : w_d$

算法：使用 KNN，寻找向量临近 $v_b - v_a + v_c$

Transformer

layer norm:

有可学习参数 $\gamma = \vec{1}, \beta = \vec{0}$ ，形状同一样本特征形状

nlp 中，两参数元素数同特征数 d_h 。初始化 layer norm 为 `layer_norm(特征数)`

CNN 中，两参数对 每一通道每一像素都有对应元素。初始化 layer norm 为 `layer_norm((通道数, 高, 宽))`

对一样本的所有参数进行 `normalize`, $= \frac{x - E(x)}{\sqrt{Var(x) + \epsilon}} * \gamma + \beta$

$E(x), Var(x)$ 为初始化 layer norm 时指定的维度内所有元素的均值/方差

例：CNN 中 $E(x)$ 即一样本所有通道所有像素的均值

为即时计算结果，非可学习参数。得到测试数据时同样直接使用测试数据的对应维均值/方差

使用 layer norm 原因：

当样本存在部分参数为无效参数，(如 nlp 中不满最大词序长度的位置)，使用批量归一结果将受无效参数影响

layer norm 可仅对部分参数求均值/方差，(如仅计算有效时序的均值/方差)

self attention：使用词序中其余词代表一词

定义：

W^Q : queries 矩阵, W^K : keys 矩阵, W^V : values 矩阵

三矩阵 $\in \mathbb{R}^{D \times d_h}$ ，参与反向传播

D 为词向量长度, d_h 为每一词 query key value 向量长度

$X \in \mathbb{R}^{T \times D}$ 为所有词向量 vstack, T 为时序长度。

$Q = XW^Q, K = XW^K, V = XW^V$

令词 i 有词向量 w_i ，则有 query $w_i W^Q$ ，key $w_i W^K$ ，value $w_i W^V$

由于全部使用同一向量 w_i ，为词向量 w_i 对自身的 attention，所以称 self attention。 w_i 取不同值时为广义 attention

self-attention 值: $Z \in \mathbb{R}^{T \times d_h} = \text{softmax}(\frac{QK^T}{\sqrt{d_h}})V$

对一词 i 计算中，将其 query 值与所有其余词的 key 值点乘 除以 $\sqrt{d_h}$ 经过 softmax，作为其余词和词 i 的相似度

除以 $\sqrt{d_h}$ ，用于减小 softmax 后的 peak

相似度做为权重，求 value 的加权平均值

softmax 后可加 dropout。

广义 attention：

K, V 时间维度长度相同，由于每一 key 对应一 value。 K, Q 特征维度长度相同，由于对每一 query 需和每一 key 点乘得到相似度

d_h 使用 Q 特征维度长度

Multi-head self-attention：

1. 输入 w_q, w_k, w_v 分别经过 3 全连接层，得到 3 个 (D, d_h) 形状结果。结果各分为 H 段。即结果为 $3 \times H$ 组 $(D, d_h / H)$ 结果

2. 所有 1. 中同一组输出的 $Q, K, V \in \mathbb{R}^{T \times d_h / H}$ 计算 attention，得到 H 个 $(T, d_h / H)$ 张量结果。

此处 Q, K, V 为已经和 W^Q, W^K, W^V 相乘后矩阵

3. 2. 中输出连接后经过另一全连接层，做最终 multi head 输出。张量形状 (T, d_h) 不变

目标为使得 attention 能类似多通道学习不同类型特征

Masked multi-head self-attention:

基于 multi-head self-attention

若 mask 一位置为否, 每一 head 权重部分 $\frac{QK^T}{\sqrt{d_h}}$ 中对应位置设为负极大值, 常取 -10^9 。

position-wise feed-forward network PWFFN:

$FNN(X) = ReLU(XW_1 + b_1)W_2 + b_2$, 其中 $W_1 \in \mathbb{R}^{D \times d_{ff}}, W_2 \in \mathbb{R}^{d_{ff} \times D}$

即两层全连接层, 由 ReLU + dropout 连接。层间使用 d_{ff} 长度向量代表信息, d_{ff} 常取 2048

称 position-wise, 由于每一词 $\in \mathbb{R}^D$ 使用相同权重

Positional encoding: 对词在句中位置进行编码, 使 transformer 类似 RNN 有时序性

PE matrix $\in \mathbb{R}^{maxT, D}$

$$PE_{pos,i} = \begin{cases} \sin(\frac{pos}{10000^{\frac{i}{D}}}) & even(i) \\ \cos(\frac{pos}{10000^{\frac{i}{D}}}) & odd(i) \end{cases}$$

maxT 为最大可能的词序长度, D 为词向量长度。

encoding 为一矩阵, 对每一时间步每一词向量元素有取值, 较低词向量位置取值周期变化较快。类似二进制值 **encoding**

前向传播时对 T 长度的词序 (Batch_size, T, D), 按元素相加 PE[:T] 矩阵, 输出张量形状仍为 (Batch_size, T, D)

性质: 对 pos 时间词 $2i, 2i+1$ 位置 encoding $PE_{pos,2i}, PE_{pos,2i+1}$, 在 δ 时间步以外的 encoding 为:

$$\begin{bmatrix} PE_{pos,2i} \\ PE_{pos,2i+1} \end{bmatrix} = \begin{bmatrix} \cos(\frac{\delta}{10000^{\frac{2i}{D}}}) & \sin(\frac{\delta}{10000^{\frac{2i}{D}}}) \\ -\sin(\frac{\delta}{10000^{\frac{2i}{D}}}) & \cos(\frac{\delta}{10000^{\frac{2i}{D}}}) \end{bmatrix} \begin{bmatrix} PE_{pos+\delta,2i} \\ PE_{pos+\delta,2i+1} \end{bmatrix}$$

即词位置变化后 encoding 变化 (式中转换矩阵) 仅和移动距离 δ 有关, 和初始位置无关

encoder 部分:

定义 source mask: 形状 (批量大小, 1, 时序长度)

所有 '<pad>' 词对应位置权重为极低值, 保证 softmax 后其余词和此词 attention 值近似为 0。

非必须, 一下仅使用 source mask 的 attention 仅称为 multi-head attention

encoder layer: 包含两残差块

输入为 (批量大小, 时序长度, 词向量长度), 在以下操作间形状不变

1. 第一残差块 $f(x) - x$ 为 multi-head self-attention 输出, 使用 source mask。 x 部分为原输入不变

2. 残差块输出通过 layer norm 层

3. 第二残差块 $f(x) - x$ 为 PWFFN 输出

PWFFN 对每一批量中每一时序的词使用相同全连接层, 所以称 position wise

4. layer norm 层

encode 结构:

1. embedding, 将词序转为词向量。embedding 参数形状 (词典大小, 词向量长度)

2. position encoding 层

3. 多层 encoder layer

最终 encoder 输出结果为最后一层 encoder layer 输出

decoder 部分:

定义 target mask: 形状 (批量大小, 时序长度, 时序长度)

所有 '<pad>' 位置权重设为极小值

第 i 号词序从 i 位置以后权重设为极低值, 模拟逐词预测整个 decoder 输入词序。即整个 mask 为三角矩阵, 上三角全为负极大值

对单一 T 长度词序 mask 时, 使用广播机制, 产生 T 个 T 长度词序。

decoder layer: 三层 (残差块, layer norm), 每一层残差块 $f(x) - x$ 分别为以下部分

1. masked multi-head self-attention:

使用 target mask, QKV 计算中使用同一 X, 即上一 decoder layer 输出 或词向量 + Position embedding 结果

2. multi-head self-attention:

使用 source mask

$Q = XW^Q$ 中 X 为 1. 的输出

$K = XW^K$, $V = XW^K$ 中 X 为 encoder 输出, 不同 decoder layer 使用相同 唯一的 encoder 输出

3. PWFFN

decoder 结构

1. 得到 (1..T-1) 词序。经过 embedding, position encoding 层。

2. 多层 decoder layer。层数无需同 encoder 层数

3. 全连接层, 对每一 decoder 特征进行全连接。即能够输出所有词序的词, 并非只能输出单一 词。

在翻译任务中由于只能基于前次预测结果预测下一词, 仍为 auto regressive

训练:

source 词序经过 encoder 得到输出。target 词序取 (1..T-1) 长度经过 decoder, 过程中使用 encoder 输出。

decoder 输出经过全连接层得到预测的下一词编号, 和期望的 target 词序使用 softmax cross entropy 计算代价

BERT

输入为两词序连接后的一整词序, 当输入仅有一词序时 segment embedding 为全 0

输入第一词为 <cls>, 第一句结尾加入 <sep>。若有第二句, 第二句句末加入 <sep>。添加完特殊词后进行 pre-training encoding

pre-training: 由 Input Embedding 得到输入, 输出为 NSP 和 MLM 连接结果。填充 '<pad>' 使得所有词序长度相同

Input Embedding: 为以下 3 项按元素相加

token embedding: 转为词向量。

segment embedding: 区分两输入的词序, 前句 segment embedding 为全 0 向量, 后句为全 1 向量。

position embedding: 随机初始化, 训练时更新

Masked Language Model MLM

更改文本中 15% 的词: 其中 80% 的情况词被掩盖, 替换为 [MASK]。10% 的情况替换为任意词, 10% 的情况不更改词

目标为改正所有被更改的词。由于不知道哪些词被更改, 每一输入词都为可更改的词。

计算代价时仅取被 mask 的词求代价。输出为词 one_hot, 非词向量

输出预测时每一词序位置输出通过同一 MLP(全连接 + ReLU + layer norm + 全连接) 得到输出。

输出为词典大小向量, 代价函数为输出向量 softmax 和目标词的 onehot 交叉熵

Next Sentence Prediction NSP:

一 binary 输出, 判断输入两句是否为上下文关系

直接将对应 <cls> 位置的输出传入全连接层, 得到 binary 输出

BERT 为 semi supervised learning, 在较大无标签数据集上训练, 在有标签数据集 fine tune 实现具体任务

在小数据集上微调 BERT 效果不稳定, 由于 1. 原 BERT 使用的 adam 删除了 bias correction 步骤, 实际使用完整 BERT。2. 增加微调 epoch 数

evaluate 时使用 BLEU 判断翻译准确度

二. 文本情感分类

分析文本作者的情感

输入: 多组 词序长度 n + 标记的情感类型

使用循环神经网络:

1. 使用预先训练的嵌入层得到词向量

2. 双向循环神经网络

每一隐藏状态分别有形状 (批量大小, $2 * \text{隐藏单元个数}$), *2 由于为双向网络

3. 全连接层, 得到分类的情感

输入为双向循环神经网络 第一和最后一隐藏状态 连接后的张量, 有形状 (批量大小, $4 * \text{隐藏单元个数}$)

使用卷积神经网络 textCNN

1. 使用预先训练的嵌入层得到词向量, 即输入形状为 (词向量长度 d , 词数 L)

2. 卷积计算: 多个一维卷积核 k_i , 每一卷积核为 (输出通道个数 c_i , 词向量长度 d , 卷积核宽 w_i) 形状的矩阵

对一个卷积核一通道的计算, 将每一卷积核窗口内的输入按元素相乘求和。输出形状 (1 , 词数 - 核宽 + 1)

总输出为一列表矩阵, 有形状 $[(c_0, L - w_0 + 1), (c_1, L - w_1 + 1), \dots]$

同一维卷积核有唯一核宽长度相同, 不同一维卷积核可使用不同核宽

3. 时序最大池化层: 类似一维全局最大池化层

针对单一一维卷积核的输出 $(c_i, L - w_i + 1)$, 对 c_i 个输入通道各得出 $L - w_i + 1$ 时序内最大值。即输出一向量, 长度为 c_i

4. 将所有时序最大池化层结果 concat

5. 全连接层, 将 concat 结果分类为情感

seq2seq 编码器-解码器

对不定长输入 $[x_1, \dots, x_T]$ 允许输出为不定长序列 $[y_1, \dots, y_{T'}]$, 使用定长的背景向量 c 做连接 解码器-编码器 的中间向量

定义: 输出词集合为 \mathcal{Y} , 其中包含一个 < eol > 特殊词代表词序末尾

编码器: 将不定长输入序列变为定长背景变量 c

对每一时间步 i , 字符 w_i , 使用循环神经网络隐藏层计算隐藏状态 h_i 。即输出 $[h_0, \dots, h_T]$

长度为 T 序列 s_i 有背景向量 $c = q(h_1, \dots, h_T)$ 。q 为自定义函数

解码器：根据前一输出序列和 c ，输出结果序列

根据解码器隐藏状态 $s_{t'} = g(y_{t'-1}, c, s_{t'-1})$ 计算 $P(y_i|y_1, \dots, y_{i-1}, c)$

对一个序列的输出，总可能性为 $P(y_1, \dots, y_{T'}|x_1, \dots, x_T)$ 。此值目标为此值最大化

代价函数

$$\begin{aligned} -\log(P(y_1, \dots, y_{T'}|x_1, \dots, x_T)) &= -\log(\prod_{i=1}^{T'} P(y_i|x_1, \dots, x_T)) \\ &= -\sum_{i=1}^{T'} \log(P(y_i|x_1, \dots, x_T)) \end{aligned}$$

解码器算法：

贪婪搜索：每一 $y_i = \operatorname{argmax}_{y \in \mathcal{Y}} P(y|y_1, \dots, y_{i-1})$

即每次取 y_i 使得 $[1, i]$ 范围内 P 值最高，直至 $i = T'$

y_{i-1} 对 y_i 的取值造成影响，若最优解 y_i 位置的词非第 i 时间步的 argmax 则贪婪算法无法取到最优解

beam search 束搜索

定义：beam size 束宽 k

计算：

1. 时间步 $i = 1$ ，选取条件概率最大 k 个值，做 k 个可选序列的首词

2. 时间步 $i = 2, 3, \dots, L$ ，对每一可选输出序列考虑整个 \mathcal{Y} ，将 $k * \mathcal{Y}$ 个新时间序列看做一整体，选出 P 值最高的 k 个作为此序列的下一词。最终每一时间步的输出序列都作为可选序列进入第三步，共 $L * K$ 个时间序列。

即，永远保持 k 个序列，而一个时间序列可能添加不同的 \hat{y} 而在下一时间步预测中有多个分支。

3. 将 $L * k$ 个序列筛选，仅保留包含 $\langle eol \rangle$ 的词序，舍弃 $\langle eol \rangle$ 后的词

4. 在 3. 的序列集合中选择值 $\frac{1}{l^a} \log(P(y_1, \dots, l)) = \frac{1}{l^a} \sum_{i=1}^l \log(P(y_i|y_1, \dots, y_{i-1}, c))$ 最大序列作为输出

l 为每一序列长度， a 为参数，常取 0.75。 $\frac{1}{l^a}$ 惩罚长度较大项

注意力机制

编码器对每一时间步产生背景向量 c'_i 而非对整个词序产生唯一背景向量 c

令 s_t 为解码器时间 t 的隐藏状态， h_t 为编码器时间 t 的隐藏状态

计算背景变量 $c_{t'}$ ：

计算每一编码器隐藏状态 h_t 的权重 $\alpha_{t't} = \frac{\exp(e_{t't})}{\sum_{t=1}^T \exp(t't)}$ ，即对不同 t 的 $e_{t't}$ 求 softmax

$e_{t't} = a(s_{t'-1}, h_t)$ ，同时基于编码器时间步 t 和解码器时间步 t'

a 函数代表编码器解码器隐藏状态相似度。例：当 $s_{t'-1}, h_t$ 长度相等， $a(s_{t'-1}, h_t) = s_{t'-1}^T h_t$

注意力机制论文： $a(s_{t'-1}, h_t) = v^T \tanh(W_s s + W_h h)$ ， v, W_s, W_h 为可学习参数

对所有时间步内的编码器隐藏状态 h_t 求加权平均，即背景变量 $c_{t'} = \sum_{t=1}^T \alpha_{t't} h_t$

计算解码器隐藏状态：(类似 GRU 算法，加入背景变量项)

$$s_{t'} = z_{t'} \odot s_{t'-1} + (1 - z_{t'}) \odot \tilde{s}_{t'} g(y_{t'-1}, c_{t'}, s_{t'-1})$$

$$\text{重置门 } r_{t'} = \sigma(W_{yr} y_{t'-1} + W_{sr} s_{t'-1} + W_{cr} c_{t'} + b_r)$$

$$\text{更新门 } z_{t'} = \sigma(W_{yz} y_{t'-1} + W_{sz} s_{t'-1} + W_{cz} c_{t'} + b_z)$$

$$\text{候选隐藏状态 } \tilde{s}_{t'} = \tanh(W_{ys} y_{t'-1} + W_{ss} (s_{t'-1} \odot r_{t'}) + W_{cs} c_{t'} + b_s)$$

三. NLP 中 CNN

将每一词转为词向量，得到 (词序长度, 词向量长度 v) 形状矩阵

一层卷积操作包含 n 通道，每个通道内卷积核个数 m 相同，形状 (h_c, v) 相同。所有卷积核宽同词向量长度，通道间卷积核高 h_c 可不同

一通道计算：

在第一维 (高) 方向上移动卷积核做卷积计算，一个卷积核的输出通过全局最大池化层

即所有通道输出形状为 $(m,)$

将所有通道输出连接，得到一层卷积操作结果

即期望每一卷积核理解不同情感的词，全局池化层保证词在词序中 spatial invariant

11 推荐算法

collaborative filter CF

分类：

memory-based CF：如 nearest neighbour-based CF (包括 user-based CF, item-based CF),

model-based CF：matrix factorization

hybrid CF

pointwise approach：从 implicit feedback 得到推荐

针对每一组 (用户，产品) 关系得到推荐

pairwise approach：从 implicit feedback 得到推荐

对每一用户，比较任意 2 产品的推荐程度

定义：

用户集合 U ，产品集合 I ，用户 u 评分过的产品 I_u^+

interaction matrix $R \in \mathbb{R}^{m \times n}$ ：对 m 个用户 n 个产品，创建 $m \times n$ 矩阵，元素 (i, j) 为用户 i 对 j 产品的评分

sparsity: $1 - \text{interaction matrix 中非空项个数} / (n \times m)$

seq-aware mode：分离测试集方法，将时间戳最近的数据留作测试集

训练集 $D = (u, i, j) | \forall u \in U, i \in I_u^+, j \in I \setminus I_u^+$

matrix factorization model

定义：

参数使用梯度下降学习，如 SGD 或 Adam

k : latent factor size。远小于 m, n

代表每一产品将有 k 个特征，用户对某一 k 元素向量的特征感兴趣

$P \in \mathbb{R}^{m \times k}$: user latent matrix

第 i 行代表用户 i 对最感兴趣的产品特征

$Q \in \mathbb{R}^{n \times k}$: item latent matrix

第 i 行代表产品 i 的特征

b_u : 用户 u 的 bias。 b_i : 产品 i 的 bias

$$\hat{R}_{ui} = PQ_{ui}^T + b_u + b_i$$

代价函数: $\sum_{(u,i)} \|R_{ui} - \hat{R}_{ui}\|^2 + \lambda(\|P\|_F^2 + \|Q\|_F^2 + b_u^2 + b_i^2)$

即，对 \hat{R} 使用平方代价函数，对 P, Q, b_u, b_i 使用 L2 正则

item based AutoRec

得到一个产品 i 的评价向量 $R_{*i} \in \mathbb{R}^m$ ，每一元素为一用户对此产品的评价。

输出一向量 $\hat{R}_{*i} \in \mathbb{R}^m$, 每一元素 R_{ji} 对应此产品对用户 j 的吸引力。

$$\hat{R}_{*i} = f(W \cdot g(VR_{*i} + \mu) + b)$$

f, g 为激活函数。可取 f 为 linear activation, g 为 sigmoid + 0.05 dropout 层

代价函数: $\sum_{i=1}^m \|R_{*i} - \hat{R}_{*i}\|_0^2 + \lambda(\|W\|_F^2 + \|V\|_F^2)$

即对每一用户, 计算对其偏好向量的预测, 和实际的偏好向量求平方代价函数。

此模型假设所有产品的数据使用同一转换得到实际用户对产品的需求。

Bayesian Personalized Ranking Loss BPRLoss: 一种 pairwise ranking 代价函数

预测: 对用户 u , 在产品 i, j 中用户偏向产品 i 的可能性为: $p(u, i, j) = \sigma(\hat{y}_{ui} - \hat{y}_{uj})$ 。

\hat{y} 为前向计算得到的用户 u 对产品 i 的偏向程度

目标: 最大化 posterior probability $p(\theta | >_u) \propto p(>_u | \theta)p(\theta)$

θ 代表参数, $>_u$ 代表实际用户对产品的排名。

$$\begin{aligned} \text{则最大化 posterior estimate } \ln(p(\theta | >_u)) &\propto \ln(\prod_{(u,i,j) \in D} \sigma(\hat{y}_{ui} - \hat{y}_{uj})p(\theta)) \\ &= \sum_{(u,i,j) \in D} \ln(\sigma(\hat{y}_{ui} - \hat{y}_{uj})) + \ln(p(\theta)) \end{aligned}$$

Hinge Loss: 一种 pairwise ranking 代价函数

$$\sum_{(u,i,j) \in D} \max(m - \hat{y}_{ui} + \hat{y}_{ij}, 0)$$

m 为超参数

即: 需要前向计算结果对两产品预测差值至少达到 m , 否则增加代价值

NeuMF: 一种 neural CF

GMF 部分:

针对 user, item latent matrix $P \in \mathbb{R}^{m \times k}$, $Q \in \mathbb{R}^{n \times k}$

$\hat{y}_{ui}^{MLP} = P_u \odot V_i$ 。即对 P, Q 的 u 行和 i 行按元素相乘

反向传播更新 P, Q

MLP 部分:

有 $U \in \mathbb{R}^{m \times k}$, $V \in \mathbb{R}^{n \times k}$

连接 U_u 和 V_i 行, 使用全连接神经网络, 输出 \hat{y}_{ui}^{MLP}

反向传播更新 U, V , 所有神经网络层参数

前向传播

连接 $[\hat{y}_{ui}^{GMF}, \hat{y}_{ui}^{MLP}]$ 。

使用单一全连接层, 输出标量结果 $\hat{y}_{ui} = f(h^T[\hat{y}_{ui}^{GMF}, \hat{y}_{ui}^{MLP}])$

代价函数使用 BPRLoss

** 分析结果 **:

$$\text{hit rate at cut off } l: \text{Hit}@l = \frac{1}{m} \sum_{u \in U} (\text{rank} \leq l)$$

取用户 u 倾向的产品, 求其中在推荐列表中排名前 l 名的个数。对所有用户重复操作, 求个数和。

AUC: ROC 曲线下方的面积

$$AUC = \frac{1}{m} \sum_{u \in U} \frac{1}{|T|}$$

12 Information Theory

self information

对长度为 n 的二进制数组 X , self information $I(X) = -\log_2(n)$

information 定义

0. 针对一 random variable

1. 两个 random variable 的总体 information \leq 分别得到两 variable information 后求和值。当两 random variable independent, 总体 information = 两者 information 之和

2. 越固定的 random variable, information 值越接近 0 **Shannon Entropy**

针对一分部

对 random variable $X \sim P$, pdf 为 $p(X)$ 。

Shannon entropy 为 expected information $H(X) = -E_{x \sim P}(\log(p(x)))$

X 离散时 $H(X) = -\sum_x p(x) \log(p(x))$

X 连续时 $H(X) = -\int_x p(x) \log(p(x)) dx$

取-log 原因

取 $X_1, X_2 \sim P$, 由 $p(X_1, X_2) = p(X_1)p(X_2)$, $\log p(X_1, X_2) = \log(p(X_1))\log(p(X_2))$ 。满足 information 第 1 条定义

取负, 由于当 $p(x)$ 升高时希望 information 值降低

Entropy 性质

令 $X \sim P$, pdf 为 $p(x)$ 。当使用新分部 Q , pdf $q(x)$, estimate X 时:

$$H(X) = -E_{x \sim P}(\log(x)) \leq -E_{x \sim P}(\log(q(x)))。$$

仅当 $Q = P$ 时等式成立

令 $X \sim P$ 。当 X 均匀分部时 entropy 最大。

$H(X) \leq \log(\frac{1}{k})$, 离散 P 中 k 为类别数, 连续 P 中 k 为值域

即均匀分部时 random variable 在不同类别/不同取值中的总体 information 最大化

joint entropy

$$H(X, Y) = -E_{(x, y) \sim P}(\log(p_{X, Y}(x, y)))$$

X, Y 为离散 random variable, $H(X, Y) = -\sum_x \sum_y p_{X, Y}(x, y) \log(p_{X, Y}(x, y))$

X, Y 为连续 random variable, $H(X, Y) = -\int_{x, y} p_{X, Y}(x, y) \log(p_{X, Y}(x, y)) dx dy$

$$H(X), H(Y) \leq H(X, Y) \leq H(X) + H(Y)$$

证: 当 X, Y 完全 dependent, 则 $H(X) = H(Y) = H(X, Y)$ 。当 X, Y 完全 independent, 则 $H(x) + H(Y) = H(X, Y)$

conditional entropy

$$H(Y|X) = -E_{(x, y) \sim P}(\log(p(y|x)))$$

X, Y 离散: $H(Y|X) = \sum_x \sum_y p_{X, Y}(x, y) \log(\frac{p_{X, Y}(x, y)}{p_X(x)})$

theorem: $H(Y|X) = H(X, Y) - H(X)$

mutual information

定义: $I(X, Y) = H(X, Y) - H(X|Y) - H(Y|X)$

' 被 X, Y 共享的 information。除去给定 X 或 Y 就能决定另一 random variable 的 entropy'

$$= E_x E_y (p_{X, Y}(x, y) \log(\frac{p_{X, Y}(x, y)}{p_X(x)p_Y(y)}))$$

$$= H(X) - H(X|Y) = H(Y) - H(Y|X)$$

$$= H(X) + H(Y) - H(X, Y)$$

性质:

$$1. \text{symmetric: } I(X, Y) = I(Y, X)$$

2.non-negative: $I(X, Y) \geq 0$

3. $I(X, Y)$ iff X, Y independent

4. 与 3. 相反, $I(X, Y) = I(X) = I(Y)$ iff X, Y 有 bijection

Pointwise mutual information

$$pmi(x, y) = \log\left(\frac{p_{X,Y}(x,y)}{p_X(x)p_Y(y)}\right)$$

(x, y) 和假设 X, Y independent 时的近似程度

Kullback-Leibler Divergence/relative entropy

计算两分部之间的距离

令 $X \sim P$, 有 pdf $p(x)$ 。使用另一分部 Q 估计 P , 有 pdf $q(x)$

relative entropy $D_{KL}(P\|Q) = E_{x \sim P}(\log(\frac{p(x)}{q(x)}))$

当 $p(x) \ll q(x)$, $\log(\frac{p(x)}{q(x)})$ 为负值, 绝对值较大

当 $p(x) \gg q(x)$, $\log(\frac{p(x)}{q(x)})$ 为正值, 绝对值较大

性质:

non-symmetric: $D_{KL}(P\|Q) \neq D_{KL}(Q\|P)$

non-negative: $D_{KL}(P\|Q) \geq 0$ 只有 $P = Q$ 时取 0

当存在 x 使得 $p(x) > 0$, $q(x) = 0$ 。定义 $D_{KL}(P\|Q) = \infty$

$$I(X, Y) = D_{KL}(P(X, Y)\|P(X)P(Y))$$

$$= E_Y(D_{KL}(P(X|Y)\|P(X)))$$

$$= E_X(D_{KL}(P(Y|X)\|P(Y)))$$

13 图像处理 (CV 课程笔记)

CMOS: RGB sensor

Bayer filter: 将输入光过滤为仅有 R/G/B 强度的光线, 每一像素被分配为 R/G/B 中一种颜色

在 (2, 2) 形状区域内对角线为 R 像素, 剩余 2 像素分别为 RB

令一像素被分配颜色 A, 显示时一 RGB 除 A 外的颜色值取相邻像素的平均值

聚焦公式:

令 f 为凸透镜焦距, u 为物距, v 为相距。则 $\frac{1}{f} = \frac{1}{u} + \frac{1}{v}$

quantisation: 相机得到颜色值为连续值, 存储时为离散值。连续值转离散值称 quantisation

RBGA:

RGB + alpha channel 代表透明度

Subtractive colour space CMYK

由白色开始, 减去 RGB 得到颜色

包含 Magenta, Cyan, Yellow。用于打印

image filter:

卷积操作, 步幅 1 填充 (卷积核边长 / 2)。

Average filter: 卷积核内所有值相同, 相加为 1。

Gaussian filter: 令卷积核中心为 (0, 0)

$$h(i, j) = \frac{1}{2\pi\sigma^2} e^{-\frac{i^2+j^2}{2\sigma^2}}$$

$$= h_x(i)h_y(j)$$

$$h_x(i) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{i^2}{2\sigma^2}}, h_y(j) \text{ 同理}$$

High-pass filter:

identity 卷积核: 仅有中心值为 1, 其余值都为 0

high-frequency 部分:

identity 卷积核 - average filter

或 卷积核中心为 1, 其余 $kernel_size^2 - 1$ 位置值为 $\frac{1}{kernel_size^2 - 1}$

high-pass 卷积核 = identity + high-frequency 卷积核。即卷积核中心以外的值相等, 为负。中心值 > 1

median filter: 对卷积核覆盖的像素取中位数输出

impulse response: 当输入为一 impulse, filter 输出称 impulse response

impulse 函数 $\delta(x)$ 定义为 $\int_x \delta(x) dx = 1$, 并仅有唯一 x 使 $\delta(x) \neq 0$

time-invariant filter: 当输入移动一 offset 时输出移动相同 offset, 并输出值不变

linear filter:

令 filter 输出为 h , 输入为 $f(x)$, linear filter 满足 $h(\alpha f(x_1) + \beta f(x_2)) = \alpha g(x_1) + \beta g(x_2)$

edge detection

离散值中求导:

forward difference: $f'(x) = f(x+1) - f(x)$

backward difference: $f'(x) = f(x+1) - f(x)$

central difference: $f'(x) = \frac{f(x+1) - f(x-1)}{2}$

卷积核中求导: forward: $[1, -1, 0]$, backward: $[[0, 1, -1]]$, central: $[1, 0, -1]$

此卷积核为 h 函数定义, 在矩阵上卷积时先取对称向量

Prewitt filter:

$$\text{竖直卷积核 } h_v = \begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix} \quad \text{水平卷积核 } h_h = \begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{bmatrix}$$

Sobel filter:

$$\text{竖直卷积核 } h_v = \begin{bmatrix} 1 & 0 & -1 \\ 2 & 0 & -2 \\ 1 & 0 & -1 \end{bmatrix} \quad \text{水平卷积核 } h_h = \begin{bmatrix} 1 & 2 & 1 \\ 0 & 0 & 0 \\ -1 & -2 & -1 \end{bmatrix}$$

Prewitt, Sobel filter 不同行作用为 smoothing。从图片多行中取值, 减少斜率受紧邻的同行像素极值影响

Prewitt filter, Sobel filter 2 卷积核分别计算图像在水平/竖直方向的斜率

$g_h = f * h_h$ 每一像素水平方向斜率, $g_v = f * h_v$ 像素竖直方向斜率

$g = \sqrt{g_h^2 + g_v^2}$ 此像素总体斜率

$\theta = \tan^{-1}(g_v, g_h)$ 每一像素仅有唯一角度值

Canny edge detection:

使边界输出 1, 非边界为 0。边界仅用一像素表示

1. 使用 Gaussian filter 减少 noise
2. 使用 Prewitt 或 Sobel 计算像素斜率
3. 每一像素进行非极大值抑制得到单一标量, 代表其为边界的权重

非极大值抑制 (不同于 SSD):

将像素周围角度分为 8 个区间, 每一区间中心为一临近像素

选择 2 区间使得 此像素角度值穿过两区间

当两区间对应的像素 总体斜率 都小于此像素总体斜率, 则此像素有权重为总体斜率, 否则权重为 0

4. 边界权重在 high threshold 以上的像素有像素值 1, 在 low high threshold 间的像素有值 0.5。否则像素值为 0

5. Hysteresis thresholding: 将 0.5 权重的像素设为 0 或 1

将对所有值为 0.5 的像素, 若临近 8 像素有值为 1 像素, 则此像素值调为 1

迭代将 0.5 权重设为 1。直至没有额外像素值可改变

Hough transformation:

将边界检测图像结果转换为边界的函数表达式, 允许有多个边界表达式

算法:

1. 对每一边界像素 (x, y) 值带入边界函数作为系数。

2. 在函数的 parameter space 中, 将参数 (m, b, r, ...) 取值范围分为网格。

3. 将 1. 中所有经过的 grid 权重 + 1。多函数经过同一 grid 时权重叠加

4. 最终权重大于一 threshold 并为 8 临近 grid local maximum 的 (m, b) 点为一边界函数参数。

slope intercept form: $y = mx + b$

求解: 参数为 m, b。在 1. 中使用的边界函数为 $b = y - mx$

double intercept form: $\frac{x}{a} + \frac{y}{b} = 1$

normal form: $x \cos(\theta) + y \sin(\theta) = \rho$

参数为 θ, ρ

相较 slope intercept form 取值范围有限, 易于定义 grid

圆弧边界:

方法 1: $(x - a)^2 + (y - b)^2 = r^2$

方法 2: $x = a + r \cos(\theta), y = b + r \sin(\theta)$

对每一可能的 r 值求 $a = x - r \cos(\theta), b = y - r \sin(\theta)$, θ 取像素 (x, y) 位置的角度

对 (a, b, r) 位置的权重 +1, 最终选择一组 (a, b, r) 使得权重高于一 threshold 并为 8 临近

grid local maximum

Interest Point Detection

image matching: 将两张图中像素对应, 如图片连接。为一种 interest point detection

Harris detector:

对图片一区域 W 内的像素

定义 $E(u, v) = \sum_{(x, y) \in W} w(x, y) [I(x + u, y + v) - I(x, y)]^2$

$I(x, y)$ 即 (x, y) 位置像素亮度

u, v 为在水平/竖直方向的移动距离。 $I(x + u, y + v)$ 可在 W 外

$w(x, y)$ 为 window function,

idle function: (x, y) 在 W 内时为 1, 否则为 0。

Gaussian function: 不论 (x, y) 是否在 W 内, 根据 (x, y) 距离 W 中心的距离取 Gaussian

分部值

Taylor expansion:

$I(x + u, y + v) = I(x, y) + uI_x(x, y) + vI_y(x, y)$

I_x, I_y 为图像在 x, y 位置的斜率

带入得 $E(u, v) = [u, v] \sum_{(x,y) \in W} w(x, y) \begin{bmatrix} I_x(x, y)^2 & I_x(x, y)I_y(x, y) \\ I_x(x, y)I_y(x, y) & I_y(x, y)^2 \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix}$

令 $M = \sum_{(x,y) \in W} w(x, y) \begin{bmatrix} I_x(x, y)^2 & I_x(x, y)I_y(x, y) \\ I_x(x, y)I_y(x, y) & I_y(x, y)^2 \end{bmatrix}$

则 M symmetric, 有 $M = P \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix} P^T$

当 M 为全零矩阵/ $\lambda_1, \lambda_2 \sim 0$ 时, (x, y) 不在边界

当 $M_{00} \gg M_{11}/\lambda_1 \gg \lambda_2$, 在竖直 edge 上

当 $M_{00} \ll M_{11}/\lambda_1 \ll \lambda_2$, 在水平 edge 上

当 M_{00}, M_{11} 较大/ λ_1, λ_2 较大时, 在边界拐点

定义一点在拐点几率: M 矩阵仅和斜率 x, y 相关, 和 u, v 无关

Harris and Stephens: $R = \lambda_1 \lambda_2 - k(\lambda_1 + \lambda_2)^2$, k 常取 0.05

Kanade and Tomasi: $R = \min(\lambda_1, \lambda_2)$

Noble: $R = \frac{\lambda_1 \lambda_2}{\lambda_1 + \lambda_2 + \epsilon}$

简化以上计算: $\det(M) = \lambda_1 \lambda_2$, 对角线求和 $\text{trace}(M) = \lambda_1 + \lambda_2$

Harris detector 为 rotation invariant, 非 scale invariant

由于 scale 后 corner 判断值偏向 edge

scaling 后图像进行 interest point detection:

Theorem: 对函数 $f(x), g(sx) = f(x)$, 有 $\frac{df(x)}{dx} = s \frac{dg(sx)}{dx}$

$M = \sum_{(x,y) \in W} w(x, y) \sigma^2 \begin{bmatrix} I_x(x, y, \sigma)^2 & I_x(x, y, \sigma)I_y(x, y, \sigma) \\ I_x(x, y, \sigma)I_y(x, y, \sigma) & I_y(x, y, \sigma)^2 \end{bmatrix}$

其中 $I_x(x, y, \sigma)$ 为图像在 (x, y) 位置 放大 σ 在水平方向的斜率

scale space extrema:

对多个 scale 的图像输入, 选择 key point 满足: 此像素为不同 scale 层输出中 (x, y) 位置最大值, 并为第 i 结果 (x, y) 位置临近 8 像素中最大值

scale adapted harris detector:

对图片使用 n 个不同 σ 的 gaussian 卷积, 对每一卷积结果每一像素计算 harris detector 值

harris detector 结果做 scale space extrema 输入, 得到可选 key point

若可选 key point harris detect 值高于阈值。则此点做最终 key point

Laplacian of Gaussian LoG:

在使用 Gaussian filter 后求二阶导

求二阶导即对亮度 I 使用 Laplace operator ΔI , 或使用卷积核 $\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$

合并 Gaussian filter 和求导操作:

f 为输入图像, g 为二维 Gaussian 卷积函数, Δ 为 laplace operator

$$\Delta(f * g) = f * \Delta g$$

$$= f * \left(-\frac{1}{\pi \sigma^4} \left(1 - \frac{x^2 + y^2}{2\sigma^2} \right) e^{-\frac{x^2 + y^2}{2\sigma^2}} \right)$$

其中 $\Delta g(x, y)$ 部分型如草帽, 称 Marr's wavelet。即为使用的卷积核 LoG

$$LoG_{norm}(x, y, \sigma) = \sigma^2 (I_{xx}(x, y, \sigma) + I_{yy}(x, y, \sigma))$$

其中 I_{xx}, I_{yy} 为 gaussian 的二次求导

Difference of Gaussian DoG

$$DoG(x, y, \sigma) = I * G(k\sigma) - I * G(\sigma)$$

$G(\sigma)$ 为使用 σ 的二维 Gaussian 卷积函数

即使用两不同 σ Gaussian 卷积核计算输出相减

卷积核为 $G(k\sigma) - G(\sigma)$, 由于卷积有分配率

$$DoG(x, y, \sigma) \approx (k-1)\sigma^2 \Delta^2(G(x, y, \sigma))$$

$$\begin{aligned} \text{证明: 给定 } \frac{dG}{d\sigma} &\approx \frac{G(k\sigma) - G(\sigma)}{k\sigma - \sigma}, DoG(x, y, \sigma) = G(k\sigma) - G(\sigma) \\ &= \frac{DoG(x, y, \sigma)}{(k-1)\sigma} \end{aligned}$$

$$\text{则 } DoG(x, y, \sigma) \approx (k-1)\sigma^2 \Delta G(x, y, \sigma)$$

Feature Description

根据 descriptor 将不同图像上的像素对应

pixel intensity: 相同 intensity 的像素一一对应

patch intensity: 使一区域以内的像素 intensity 对应

gradient orientation: 相同斜率和角度的像素对应

histogram: 收集一区域以内的像素 intensity 分部, 相同分部的区域对应

Scale-invariant Feature Transform SIFT

目标: 将图像转换为 interest point, 每一点有 feature vector

计算 N 个 Gaussian filter 结果, 其中 $k\sigma$ 值取 $\sigma, \sqrt{2}\sigma, 2\sigma, 2\sqrt{2}\sigma, 4\sigma, \dots$

相邻 Gaussian filter 间计算差值, 得到 N-1 个 DoG 结果

对 DoG 结果使用 scale space extrema 得到 key point

key point refine:

假设 key point 周围像素值以二次函数分部, 得到 refine 的 key point $v = (x, y, \sigma)$ 使得 $DoG(v)$ 值为极值。此时 v 不一定在一整数像素/ σ 上

$$\text{通过 Taylor 展开: } DoG(v + \delta) = DoG(v) + \frac{dDoG^T}{dv} \delta + \frac{1}{2} \delta^T \frac{d^2 DoG}{dv^2} \delta$$

$$\text{估计一阶导值 } \frac{DoG(x)}{dx} = \frac{DoG(x+1, y, \sigma) + DoG(x-1, y, \sigma)}{2}$$

二阶导值使用一阶导差值

$$\text{则 } \frac{dDoG(v+\delta)}{d\delta} = \frac{dDoG}{dv} + \frac{d^2 DoG}{dv^2} \delta$$

$$\text{取 } \frac{dDoG(v+\delta)}{d\delta} = 0, \text{ 则 } \delta = -\frac{d^2 DoG}{dv^2}^{-1} \frac{dDoG}{dv}。 DoG(v + \delta) \text{ 有最值}$$

则 $v + \delta$ 中 xy 项为 key point 坐标, σ 为 refined scale

orientation assignment:

在 keypoint 周围采样像素, 采样时使用固定大小窗口。目标为得到此 key point 的方向, 用于生成 descriptor 时从旋转后的窗口取样

将 360 度角度分为 36 bin, 每一 bin 占 10 度。采集像素的角度被分入 36 个 bin 中, 斜率总和最大的 bin 作为此 key point 的角度

key point descriptor:

采样窗口大小 $\propto \sigma$, 即采样的像素个数不变, 间隔 $\propto \sigma$

将采样后的像素分入网格, 每一网格内方向分为 8 bin, 每一 bin 为一方向的斜率总和

descriptor 为将每一 grid 的 8 bin 连接

keypoint matching

1. 使用 euclidian distance 对应两图 key point

2. 对 1. 中对应上的 key point $(u, v)(x, y)$ 得到 key point 间转换表达式:

$$\begin{bmatrix} u \\ v \end{bmatrix}^T = \begin{bmatrix} m_1 & m_2 \\ m_3 & m_4 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}^T + \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}^T$$

对所有 key point pair $(u_i, v_i), (x_i, y_i)$:

$$\begin{bmatrix} x_1 & y_1 & 0 & 0 & 1 & 0 \\ 0 & 0 & x_1 & y_1 & 0 & 1 \\ x_2 & y_2 & 0 & 0 & 1 & 0 \\ 0 & 0 & x_2 & y_2 & 0 & 1 \\ \dots & & & & & \end{bmatrix} \begin{bmatrix} m_1 \\ m_2 \\ m_3 \\ m_4 \\ b_1 \\ b_2 \end{bmatrix} = \begin{bmatrix} u_1 \\ v_1 \\ u_2 \\ v_2 \\ \dots \end{bmatrix}$$

简写为 $Am = b$

求解 $m = (A^T A)^{-1} A^T b$

Random Sample Consensus RANSAC: 另一得到 key point 间转换方法

1. 任取两点, 得到两点定义的线性函数 $f(x)$
2. 计算满足 $f(x') - \epsilon \leq y' \leq f(x') + \epsilon$ 的样本 (x', y') 个数
3. 当满足的样本个数到达 threshold 则停止, 否则重复 1. 取点

求两图像的 key point matching 即求 key point 的线性函数。xy 权重即旋转矩阵, 偏差即位

移

Speeded-Up Robust Feature SRUF: feature descriptor

Haar wavelets: 包含 $2 * (2,)$ 向量, 分别为 $d_h = [-1, 1], d_v = [-1, 1]^T$ 。

将采样点周围像素值和 Haar wavelets 按元素相乘求和, 得到仅在水平/竖直方向斜率 d_h, d_v ,

一采样点 descriptor: $(\sum d_h, \sum d_v, \sum |d_h|, \sum |d_v|)$

颜色不变区域: descriptor 4 值都较小

纵向像素值不变 (纵向条纹): $\sum |d_h|$ 较大

纵向像素值不变 (横向渐变): $\sum d_h, \sum |d_h|$ 较大

一 key point descriptor 为将每一采样点 descriptor 向量连接

Binary Robust Independent Elementary Feature BRIEF:

$$\tau(p, q) = \begin{cases} 1 & h(p) > h(q) \\ 0 & otherwise \end{cases}$$

在一 key point 周围取 n_d 对像素 (p_i, q_i) , descriptor 为 $\sum_{i \leq n_d} \tau(p_i, q_i)$ 左移 i

两图中所有 key point 必须共用同一随机像素对分部, 由于需要相同分部才能保证描述相同像素变化趋势

descriptor 间距离为两 descriptor XOR 值取 XOR 值为 1 个数 (即 hamming distance)。

假设图像仅有移动, 不支持旋转, 缩放后 scale match

Histogram of Oriented Gradient HOG

对一区域的像素 或整张图片计算 feature

将图像整体分为网格 (非在某一像素周围分网格), 每一 grid 计算 gradient orientation histogram 向量 v

多个 grid 组成一 block, 有合并的 locally normalized descriptor $v_{norm} = \frac{[v_1, v_2, \dots]}{\sqrt{\| [v_1, v_2, \dots] \|_2^2 + \epsilon^2}}$

将 block 位置按卷积操作移动, HOG descriptor 为所有 block 的 descriptor 连接

Optic Flow

估计视频中的亮度在 x, y 方向分别移动 (u, v)

假设:

brightness constancy: 两帧间物体上同一点亮度相同

small motion: 两帧间物体上同一点移动距离较小

spatial coherence: 相邻像素的移动方向相近

由假设得 $I(x+u, y+v, t+1) = I(x, y, t) \approx I(x, y, t) + \frac{dI}{dx}u + \frac{dI}{dy}v + \frac{dI}{dt}$

则对每一像素 (x, y) 都有 $I_x u + I_y v + I_t = 0$

I_x, I_y 仍为一图片中亮度斜率, 由于取关于 x, y 偏导。 I_t 为固定位置亮度斜率

Lucas-Kanade Method: 计算 (u, v)

对每一像素 p_i : $[I_x(p_i), I_y(p_i)][u, v]^T = -I_t(p_i)$

令 $A = \begin{bmatrix} I_x(p_1) & I_y(p_1) \\ I_x(p_2) & I_y(p_2) \\ \dots & \dots \end{bmatrix}$, $x = [u, v]^T$, $b = -\begin{bmatrix} I_t(p_1) \\ I_t(p_2) \\ \dots \end{bmatrix}$ 则 $Ax = b$

根据假设 3, A 中点 p_i 取一较小区域内像素

取 $x = \operatorname{argmin}_x \|Ax - b\|^2$ 为最优 $[u, v]^T$

则有 $x = (A^T A)^{-1} A^T b$

$$A^T A = \sum_{p_i} \begin{bmatrix} I_x(p_i)^2 & I_x(p_i)I_y(p_i) \\ I_x(p_i)I_y(p_i) & I_y(p_i)^2 \end{bmatrix}$$

$$A^T b = -\sum_{p_i} \begin{bmatrix} I_x(p_i)I_t(p_i) \\ I_y(p_i)I_t(p_i) \end{bmatrix}$$

multi-scale Lucas-Kanade method: 得到较大范围的移动

对图片 up sampling: 一层 up sample 对 2×2 区间取一像素, scale 编号由 sampling 最少 scale 至原图 scale

1. 从上一 sampling 层 $l-1$ 中计算此层 (u, v) : $u^l = 2u^{l-1}, v^l = 2v^{l-1}$

2. 计算 warped image: $I_{warped}^l = I^l(x+u, y+v, t+1)$

即根据上层 u, v 预测移动此 scale 的 t 时刻图像

3. 计算 $I_t = I_{warped}^l(x+u^{l-1}, y+v^{l-1}, t+1) - I^l(x, y)$ 。 I_x, I_y 由 I_l 计算得

使用 2. 中移动后的图像计算精确的 u, v 量, refine scale。直接对此 scale 计算则无法得到较大 u, v

值

4. 计算 incremental flow: $[u^\delta, v^\delta]^T = (A^T A)^{-1} A^T b$

5. 更新 $u^l = u^{l-1} + u^\delta, v^l = v^{l-1} + v^\delta$

Horn-Schunck method:

定义函数 u, v , $u(x, y), v(x, y)$ 为在像素 x, y 位置的水平/竖直方向移动

最小化 $E = \int_y \int_x (I_x u + I_y v + I_t)^2 + \alpha (\|\nabla u^2\| + \|\nabla v^2\|) dx dy$

即最小化临近像素的方向向量改变量, 积分中所有函数针对 x, y 位置进行求导/取值

由 Euler Lagrange equation:

$$(I_x u + I_y v + I_t)I_x - \alpha \Delta u = 0$$

$$(I_x u + I_y v + I_t)I_y - \alpha \Delta v = 0$$

Δ 为 Laplace operator

带入 $\Delta u = \bar{u} - u$, \bar{u} 为一小像素区域内的 u, v 平均值

$$(I_x^2 + \alpha)u + I_x I_y v = \alpha \bar{u} - I_x I_t$$

$$(I_y^2 + \alpha)v + I_x I_y u = \alpha \bar{v} - I_y I_t$$

得到 u, v 关于 \bar{u}, \bar{v} 的公式:

$$u = \bar{u} - \frac{I_x(I_x \bar{u} + I_y \bar{v} + I_t)}{I_x^2 + I_y^2 + \alpha}$$

$$v = \bar{v} - \frac{I_y(I_x \bar{u} + I_y \bar{v} + I_t)}{I_x^2 + I_y^2 + \alpha}$$

初始化 u, v 对所有 x, y 像素取值为 0。根据以上公式进行迭代得到 motion 结果

Lucas-Kanade Tracker

针对原图像 I, 移动后图像 J

取一物体的锚框内像素 (u, v) , 使得将图像 I 按 (u, v) 方向移动后和 J 像素平方差之和最小

$$\text{即 } \min_{u,v} \sum_{(x,y) \in \text{pixels}} (I(x, y) - J(x + u, y + v))^2$$

求解: 对 u, v 微分

$$\frac{dE}{du} = 2 \sum_{(x,y)} (I(x, y) - J(x + u, y + v)) \frac{dJ}{du} \Big|_{(x+u)}$$

$$= -2 \sum_{(x,y)} (I(x, y) - J(x + u, y + v)) \frac{dJ}{dx} \Big|_{(x+u)}$$

$$= -2 \sum_{(x,y)} (I(x, y) - J(x, y) - \frac{dJ}{dx} u - \frac{dJ}{dy} v) \frac{dJ}{dx}$$

$$= -2 \sum_{(x,y)} (-I_t - I_x u - I_y v) I_x, \quad (\text{假设 } \frac{dJ}{dx} \approx \frac{dI}{dx})$$

$$v \text{ 公式同理: } = -2 \sum_{(x,y)} (-I_t - I_x u - I_y v) I_y$$

令两微分值为 0:

$$-\sum_{(x,y)} \begin{bmatrix} I_x I_t \\ I_y I_t \end{bmatrix} - \sum_{(x,y)} \begin{bmatrix} I_x^2 & I_x I_y \\ I_x I_y & I_y^2 \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = 0$$

$$\begin{bmatrix} u \\ v \end{bmatrix} = -(\sum_{(x,y)} \begin{bmatrix} I_x^2 & I_x I_y \\ I_x I_y & I_y^2 \end{bmatrix})^{-1} \sum_{(x,y)} \begin{bmatrix} I_x I_t \\ I_y I_t \end{bmatrix}$$

correlation filter method:

令原物体位置为 template, 移动后物体可能在的区域称 search window

template 经过 CNN 得到特征, 其中不通过全连接层。特征仍为多通道矩阵

对 search window 使用同一 CNN 得到特征。

将 template 的特征作为互相关计算的卷积核, 和 search window 的特征得到卷积结果。

由于两特征通道数相同, 结果仅有 1 通道, 为一矩阵。

结果称 score map, 其中值较大的区域即为可能的移动后物体位置

可对特征进行卷积操作后做互相关计算, 加入的卷积层参与学习。迁移学习时改变此卷积层
创建数据集:

人为标注每一帧物体位置。

得到物体在 search window 相对位置, 按比例缩放后将期望的 score map 中对应位置设为 1, 其余位置为 -1

tracking by detection:

将目标检测的 CNN 中标签个数设为 2, 仅区分物体和背景

迁移学习时仅改变最后全连接层

Camera Model

已知物体函数表达式和 2d 图像, 预测观测点位置

$$\text{令 } X = \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix}, \begin{bmatrix} x \\ y \\ z \end{bmatrix} = \begin{bmatrix} p_1^T \\ p_2^T \\ p_3^T \end{bmatrix} \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix}$$

prospective 投影后 $x = \frac{p_1^T X}{p_3^T X}, y = \frac{p_2^T X}{p_3^T X}$

可得 $X^T p_1 - X^T p_3 x = 0, X^T p_2 - X^T p_3 y = 0$

$$\begin{bmatrix} X^T & 0 & -X^T x \\ 0 & X^T & -X^T y \end{bmatrix} \begin{bmatrix} p_1 \\ p_2 \\ p_3 \end{bmatrix} = 0$$

令第一矩阵为 A, 有形状为 (2, 12)。第二矩阵 \vec{p} 形状 (12, 1)

以上针对一像素 (x, y, z) 位置。对 n 个像素, 第一矩阵形状为 (2*n, 12), 解 $\vec{p} = \operatorname{argmin}_p \|A\vec{p}\|^2$

加入条件 $\|\vec{p}\| = 1$, 使 $\vec{p} = 0$ 不为有效解

求解: 求 SVD: $A = USV^T$, \vec{p} 为 AV 中对应最小 singular value 的列

14 Deep Feedforward Network (Deep Learning 第 6 章笔记)

SVM 支持向量机

仍通过 $w^T x + b$ 得到输出, 输出仅表示 identity, 正值说明有 identity, 负值说明没有

依据: 一个平面的公式为 $\beta_0 + \beta_1 x_1 + \beta_2 x_2 = 0$, 则当计算 $w^T x + b$ 得到值后, >0 则为平面上方的数据点, <0 为下方数据点

kernel trick

kernel method 将数据集表示成相近的两个数据点一组的集合 (x_i, x_j) , kernel method 将一对数据变为单一数据点 $x = k(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle$

kernel method 使用 ϕ 转换数据的纬度, 而点乘化简后无需先计算 $\phi(x_i), \phi(x_j)$ 即可得到新数据点 x

manifold hypothesis:

当训练数据集包含大量无规律的数据, 则将其大部分视为无效数据, 并只关心落在一个 manifold 上的数据。

例: 生成图像文字声音时数据大多很集中, 当像素文字随机分布时生成图像大多无意义

deep feedforward network/feedforward neural network/multilayer perceptrons MLP:

找到 θ 使得 $f(x; \theta)$ 最接近数据 y 值。 f^* 为最理想的 f, 即 $f^*(x) = y$ 。 θ 可为多个参数, 如 $f(x; w, b) = x^T w + b$

$f^*(x) = f^{(3)}(f^{(2)}(f^{(1)}(x)))$, $f^{(1)}$ 为 network 第一层。每一 $f^{(i)}(x) = \phi(x; \theta)^T w$

神经网络

1. 结构:

输入层没有 weight, 第一 hidden layer 得到所有输入层的值。

hidden layer 和输出层所有输出都为 0/1, 非连续的值

2. 一层 hidden layer 计算方法: $f^{(i)}(x; W, c) = \sigma(W^T x + c)$

x 为前一层的输出向量, 输入层 x 即为训练参数向量。

c 为此层常数向量

$z = W^T x$ 为一层 hidden layer 对输入取得的中间值向量, 称 logit。 $a = \sigma(z + c)$ 为对 $z + c$ 每一元素取 σ 的结果向量, a 即此层的输出。

W 为此层参数矩阵, 行数 = 前层节点数, 列数 = 当前层节点数

X 为多个参数点的训练集中前一层的输出矩阵, 行数为数据点个数, 列数为前一层节点数

XW 当 W 对参数集矩阵操作时, 每行向量 z_i^T 此时为一层 hidden layer 各节点对第 i 参数点的中间值向量。对每行 $+c^T$ 并分别取 σ 得到输出矩阵, a_{ij} 为当使用第 i 个参数点时此层第 j 节点的输出

cross entropy

分部 p 和分部 q 间的 cross entropy $H(p, q) = -E_p(\log(q))$ 。为 expected value of $\log(q)$ with respect to distribution p

cost function

当使用 maximum likelihood 估计参数时, cost function $J(\theta)$ 为训练输入参数的分部和训练结果参数的分部间 cross-entropy: $J(\theta) = -E_{x, y \sim \text{training_dataset}}(\log(p_{\text{model}}(y|x)))$

对于每一在训练集内的 (x, y), 求 $\log(p_{\text{model}}(y|x))$, 并求 expected value。 $p_{\text{model}}(y|x)$ 即训练得到的 y 关于 x 的分部

例: 当 model 为 $y = N(f(x; \theta), 1)$ 正则分部时, $J(\theta) = -E_{x, y \sim \text{data}}(y - f(x; \theta))^2 + \text{const}$

output layer

当输出层的结果和不为 1 时, 代表数据没有被准确分到某一类中, 使用 exponentiation and normalisation

normalisation 后结果 $p = \frac{\tilde{p}}{\sum \tilde{p}'}$, 为 \tilde{p} 在所有结果中占的比例。 \tilde{p} 为未 normalise 值

假设输出层结果 $\tilde{P}(y|x)$ 有 $\log(\tilde{P}(y|x)) = yz$

$$\tilde{P}(y|x) = \exp(yz)$$

$$P(y|x) = \frac{\exp(yz)}{\sum_{y'=0}^1 y'z}, \text{ 称 softmax function}$$

$$P(y|x) = \sigma((2y - 1)z), \text{ y, y' 为训练目标结果, 所以 } \sum_{y'=0}^1 \text{ 包含所有 y'}$$

对 softmax function 使用 log likelihood 原因: $\log \text{softmax}(z)_i = z_i - \log \sum_j \exp(z_j)$ 。

当 z_i 为 dominant, 并对应期望的输出项。 $\log \text{softmax}(z)_i = 0$ 。则此项不产生高 cost, 否则产生 cost。

hidden unit

代表一个 hidden layer 节点的激发函数。

1. rectified linear unit: $g(x) = \max(0, x)$

无法用于 gradient based learning, 由于一阶导为 0

基于 rectified linear unit 的优化: $g(x) = \max(0, x) + a * \min(0, x)$

a = -1: absolute value rectifier

a 为极小值: leaky ReLU

a 为可学习值: Parametric ReLU, PReLU

2. Maxout units

将 x 分为多组, 每组 $h(x)$ 为组内最高值

backward propagation

一种计算 gradient 的方法, 区别于使用 gradient 进行学习的 stochastic gradient descent 算法:

After the forward computation, compute the gradient on the output layer:
 $\mathbf{g} \leftarrow \nabla_{\mathbf{y}} J = \nabla_{\mathbf{y}} L(\hat{\mathbf{y}}, \mathbf{y})$
for $k = l, l-1, \dots, 1$ **do**
 Convert the gradient on the layer's output into a gradient into the pre-
 nonlinearity activation (element-wise multiplication if f is element-wise):
 $\mathbf{g} \leftarrow \nabla_{\mathbf{a}^{(k)}} J = \mathbf{g} \odot f'(\mathbf{a}^{(k)})$
 Compute gradients on weights and biases (including the regularization term,
 where needed):
 $\nabla_{\mathbf{b}^{(k)}} J = \mathbf{g} + \lambda \nabla_{\mathbf{b}^{(k)}} \Omega(\theta)$
 $\nabla_{\mathbf{W}^{(k)}} J = \mathbf{g} \mathbf{h}^{(k-1)\top} + \lambda \nabla_{\mathbf{W}^{(k)}} \Omega(\theta)$
 Propagate the gradients w.r.t. the next lower-level hidden layer's activations:
 $\mathbf{g} \leftarrow \nabla_{\mathbf{h}^{(k-1)}} J = \mathbf{W}^{(k)\top} \mathbf{g}$
end for
