

图像形成

齐次坐标表示: $\tilde{x} = (\tilde{x}, \tilde{y}, \tilde{w}) \in \mathcal{P}^2$, 其中 $\mathcal{P}^2 = \mathcal{R}^3 - (0,0,0)$ 为二维射影空间. 齐次坐标 (x, y, w) 对应笛卡尔坐标 $(\frac{x}{w}, \frac{y}{w})$. $\tilde{x} = (\tilde{x}, \tilde{y}, \tilde{w}) = \tilde{w}(x, y, 1) = \tilde{w}\tilde{x}$. 其中 $\tilde{x} = (x, y, 1)$ 被称为增广向量.当齐次向量 \tilde{x} 的 $\tilde{w} = 0$ 时, 也可称为理想点或**无穷远点**. 它没有对应的非齐次表示.

Translation 平移 2D translations can be written as $\tilde{x}' = x + t$ or $\tilde{x}' = [I \quad t] \tilde{x}$, I is the (2×2) identity matrix or $\tilde{x}' = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \tilde{x}$. z 2 自由度, orientation 不变.

Rotation-Translation 平移+旋转 也称为 $2D$ rigid body motion or the $2D$ Euclidean transformation. It can be written as $\tilde{x}' = R\tilde{x} + t$ or $\tilde{x}' = [R \quad t] \tilde{x}$ where $R = \begin{bmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{bmatrix}$ is an orthonormal rotation matrix with $RR^T = I$ and $|R| = 1, 3$ 自由度, lengths 不变. 注意这个是逆时针旋转.

Scaled Rotation 缩放旋转; 相似变换 similarity transform, this transformation can be expressed as $\tilde{x}' = sR\tilde{x} + t$, where s is an arbitrary scale factor. also as $\tilde{x}' = [sR \quad t] \tilde{x} = \begin{bmatrix} a & -b & t_x \\ b & a & t_y \end{bmatrix} \tilde{x}$, we no longer require that $a^2 + b^2 = 1$. The similarity transform preserves angles between lines. 4 自由度, angles 不变.

Affine 仿射 The affine transformation is written as $\tilde{x}' = A\tilde{x}$, where A is an arbitrary 2×3 matrix, i.e., $\tilde{x}' = \begin{bmatrix} a_{00} & a_{01} & a_{02} \\ a_{10} & a_{11} & a_{12} \end{bmatrix} \tilde{x}$. Parallel lines remain parallel under affine transformations. 6 自由由

Projective 投影 This transformation, also known as a perspective transform (透视变换) or homography, operates on homogeneous coordinates, $\tilde{x}' = \tilde{H}\tilde{x}$, where \tilde{H} is an arbitrary 3×3 matrix. Note that \tilde{H} is homogeneous, i.e., it is only defined up to a scale, and that two \tilde{H} matrices that differ only by scale are equivalent. The resulting homogeneous coordinate \tilde{x}' must be normalized in order to obtain an inhomogeneous result x , i.e., $x' = \frac{h_{00}x+h_{01}y+h_{02}}{h_{20}x+h_{21}y+h_{22}}$, $y' = \frac{h_{10}x+h_{11}y+h_{12}}{h_{20}x+h_{21}y+h_{22}}$. Perspective transformations preserve straight lines (i.e., they remain straight after the transformation).8 自由由

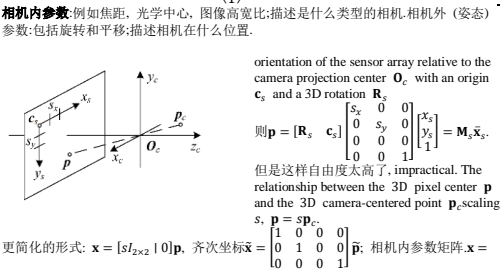
三维变换: Translation: $[I \quad t]_{3 \times 4}$, 3 自由由; Rigid(Euclidean): $[R \quad t]_{3 \times 4}$, 6 自由由; Similarity: $[sR \quad t]_{3 \times 4}$, 7 自由由; Affine: $[A]_{3 \times 4}$, 12 自由由; Projective: $[H]_{3 \times 4}$, 15 自由由.

三维到二维投影: 正交投影是最简单的投影模型. 去掉三维坐标 \mathbf{p} 的 z 分量, 得到对应的二维点 \mathbf{x} . 记作 $\mathbf{x} = [I_{2 \times 2} \quad 0] \mathbf{p}$. 齐次坐标的表示形式为: $\tilde{\mathbf{x}} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \tilde{\mathbf{p}}$, 去掉了 z 分量, 但保留了 w 分量. 缩放正交投影: $\mathbf{x} = [I_{2 \times 2} \quad 0] \mathbf{p}$.

透视投影: $\tilde{\mathbf{x}} = \mathcal{P}_c(\mathbf{p}) = \begin{bmatrix} x/z \\ y/z \\ 1 \end{bmatrix}$, 齐次坐标表示形式: $\tilde{\mathbf{x}} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \tilde{\mathbf{p}}$. 去掉

了 w 分量, 因此无法从图像中恢复距离信息. 如针孔相机的投影模型: $\begin{pmatrix} x \\ y \\ 1 \end{pmatrix} \sim \begin{pmatrix} fX \\ fY \\ Z \end{pmatrix}$

相机内参数:例如焦距, 光学中心, 图像高宽比;描述的是什么类型的相机(相机外 (姿态) 参数:包括旋转和平移-描述相机在什么位置.



更简化的形式: $\mathbf{x} = [I_{2 \times 2} \quad 0] \mathbf{p}$. 齐次坐标表示形式: $\tilde{\mathbf{x}} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \tilde{\mathbf{p}}$. 相机内参数矩阵 $\mathbf{K} = \begin{bmatrix} f & 0 & 0 \\ 0 & f & 0 \\ 0 & 0 & 1 \end{bmatrix}$. 相机外参数矩阵: 标定校准, \mathbf{u}, \mathbf{v} 代表二维

分量, X, Y, Z 代表三维分量, 取代 x, y, z . 使用已知的三维空间点及其像素位置, 求解矩阵中的 11 个未知数.

汇总表示: $\mathbf{x} = \mathbf{K}[R \quad t] \tilde{\mathbf{x}}$. 其中 \mathbf{x} : Image Coordinates: $(u, v, 1)$; \mathbf{K} : Intrinsics Matrix (3×3) ; R : Rotation (3×3) ; t : Translation (3×1) ; $\tilde{\mathbf{x}}$: World Coordinates: $(X, Y, Z, 1)$

内参矩阵就是将图像坐标系转化为像素坐标系. $\begin{bmatrix} f/dx & 0 & u_0 \\ 0 & f/dy & v_0 \\ 0 & 0 & 1 \end{bmatrix}$, 其中 u_0, v_0 主点的实际位置, 单位也是像素. 外参矩阵就是真实物理空间坐标到相机坐标系的旋转加平移变换. 本质就是一个变换矩阵 $[R \quad t]$. 内外参数矩阵形式: $\begin{bmatrix} X \\ Y \\ 1 \end{bmatrix} = \begin{bmatrix} f/dx & 0 & u_0 \\ 0 & f/dy & v_0 \\ 0 & 0 & 1 \end{bmatrix} [R \quad t] \begin{bmatrix} x_w \\ y_w \\ z_w \end{bmatrix} = M \begin{bmatrix} x_w \\ y_w \\ z_w \end{bmatrix}$ 优点: 相机所有特性都集中在一个矩阵中; 能够计算任意三维点在图像中的成像位置; 缺点: 不能单独计算特定参数; 内部参数与外部参数混合, 导致相机位置不能变, 一旦变了就失效. 此时可以使用运动校准, 也就是利用运动标定板来实现, 不需要知道位置和朝向;

将中间的 image plane 坐标系的坐标移动到左上角, 我们就需要增加一些映射关系 $\begin{pmatrix} x_s \\ y_s \end{pmatrix} = \begin{pmatrix} f_x x_c / z_c + s_y / z_c + c_x \\ f_y y_c / z_c + c_y \end{pmatrix} \Leftrightarrow \tilde{\mathbf{x}}_s = \begin{bmatrix} f_x & s_c & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix} \tilde{\mathbf{x}}_c$

这个映射矩阵就成为相机内参矩阵 (camera intrinsic), 表示把三维坐标映射到图像平面. **图像的颜色空间.** additive colors red, green, and blue; subtractive colors cyan, magenta, and yellow. 加色混合: 红色、绿色和蓝色可以混合产生青色、品红色、黄色和白色. 减色混合: 蓝、洋红和黄色可以混合产生红色、绿色、蓝色和黑色 RGB: $R = \int L(\lambda) S_R(\lambda) d\lambda$, $G = \int L(\lambda) S_G(\lambda) d\lambda$, where $L(\lambda)$ is the incoming spectrum of light at a given pixel and $S_R(\lambda), S_G(\lambda), S_B(\lambda)$ are the red, green, and blue spectral sensitivities of the corresponding sensors.

负值问题. 以及很难区分, 故 CIE 定义了 XYZ 空间. The transformation from RGB to XYZ is given by $\begin{bmatrix} X \\ Y \\ Z \end{bmatrix} = \begin{bmatrix} 0.49 & 0.31 & 0.20 \\ 0.17697 & 0.81240 & 0.10663 \\ 0.00 & 0.01 & 0.99 \end{bmatrix} \begin{bmatrix} R \\ G \\ B \end{bmatrix}$

Because the response of the human visual system is roughly logarithmic (we can perceive relative luminance differences of about 1%), the CIE defined a non-linear re-mapping of the XYZ space called 'L*a*b*' (also sometimes called CIELAB), where differences in luminance or chrominance are more perceptually uniform, as shown in Figure 2.30 b. The 'L' component of lightness is defined as $L^* = 116(f(\frac{Y}{Y_n}) - \frac{16}{116})$, where Y_n is the luminance value for nominal white (Fairchild 2013) and $f(t) = \begin{cases} t^{1/3} & t > \delta^3 \\ \frac{1}{29}(\delta^3 + 16t) & \text{else} \end{cases}$ is a finite-slope approximation to the cube root with $\delta = 6/29$. The resulting 0...100 scale roughly measures equal amounts of lightness perceptibility. In a similar fashion, the a^* and b^* components are defined as $a^* = 500[f(\frac{X}{X_n}) - f(\frac{Y}{Y_n})]$ and $b^* = 200[f(\frac{X}{X_n}) - f(\frac{Z}{Z_n})]$ where again, (X_n, Y_n, Z_n) is the measured white point.

If we divide the Y^* values by the sum of $X + Y + Z$, we obtain the chromaticity coordinates $x = \frac{X}{X+Y+Z}$, which sum to 1. 此外注意 $\mathbf{F} = \frac{\mathbf{E}}{R+G+B}$

其他BRDF: Bidirectional Reflectance Distribution Function. 表面光反射分布函数. BRDF: $f(\theta_i, \phi_i; \theta_r, \phi_r) = \frac{\text{Luminance}(\theta_r, \phi_r)}{\text{Irradiance}(\theta_i, \phi_i)}$. 像像素 = 体反射 + 表面反射. 枕状畸变: 内凹型; 桶状畸变: 外凸型; 薄透镜成像: $\frac{1}{d} = \frac{1}{d_o} + \frac{1}{d_i}$, 曝光 = 光圈 + 快门速度; 光圈的直径 D 限制光的入射范围 (光圈可以在透镜的任一侧); 快门速度决定从光源射入的光的量; 更改光圈尺寸影响景深; 小光圈能增加成像基本清晰的范围

图像处理 **数字图像:** 数字图像可以视为一种函数: $f: R^2 \rightarrow R$. 通常情况下, 我们期望图像 $\gamma \hat{s}_i + \beta_i$

包含有限数量的像素, 且幅值在有限范围内. 例如 $f: [a, b] \times [c, d] \rightarrow [0, 255]$. 彩色图像可以视为三个函数的组合. 数字图像经过采样和量化处理: 像素坐标和幅值均离散化. 二维采样: 图像空间坐标的离散化. 幅值量化: 图像幅度值的离散化.

点运算符: 常见点运算符有控制对比度、gamma 变换、直方图均衡化. 图像合成与蒙版. **直方图均衡化:** histogram: $h(I)$; $c(I) = \frac{1}{N} \sum_{i=0}^L h(i) = c(I-1) + \frac{1}{b} h(I)$, where N is the number of pixels in the image. For any given grade or intensity, we can look up its corresponding percentile $c(I)$ and determine the final value that the pixel should take. When working with eight-bit pixel values, the I and c axes are rescaled from $[0, 255]$. 完全直方图均衡化: $f(I) = c(I)$; 部分直方图均衡化: $f(I) = ac(I) + (1-a)I$ 蒙版 $C = (1-a)B + aF$.

邻域运算: 线性滤波使用邻域的**线性组合**代替原来的像素值 $g(i, j) = \sum_k l(f(i+k, j)+h(k, l))$, 其中 $h(k, l)$ 为滤波核. 等价的**卷积**形式: $g(i, j) = \sum_k l(f(i-k, j)-D)h(k, l)$, 或记为 $g = f * h$.

箱式滤波: 滑动平均 - 通过滑动窗口平均得到新的像素值. 以局部均值代替原始像素值; 也称为平方平均滤波; 具有平滑和模糊的效果. 滤波核是 $\frac{1}{b \times b}$ 的全是 1 的 3×3 矩阵. **高斯滤波:** 加权滑动平均 - 使用滑动窗口的加权平均得到新的像素值. 从二维连续的高斯函数采样得到. 滤波器的值从中心向外逐渐减小. $G_\sigma = \frac{1}{2\pi\sigma^2} e^{-\frac{(x^2+y^2)}{2\sigma^2}}$ 高斯核使得滤波后的图像更加平滑. σ 决定平滑程度. 性质: 去除图像中的“高频”成分 (低通滤波器); 高斯核与自身的卷积仍然是高斯核. 如何选择滤波核的尺寸? 尽量使边缘的值接近 0; 经验法: 滤波核宽度为 6σ .

若图像尺寸为 $M \times M$, 卷积核尺寸为 $N \times N$. 采用不可分滤波器需要进行多少次**乘法运算**? 每像素 $N \times N$ 次, 共 $N^2 \times M^2$ 次. 采用可分滤波器需要进行多少次乘法运算? 每像素 $2 \times N$ 次, 共 $2 \times N \times M^2$ 次. 可分卷积核: 一些二维卷积核可写成行向量和列向量的乘积. 比如前边箱式滤波那里.

傅里叶变换: 连续域 $H(\omega_x, \omega_y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} h(x, y) e^{-j(\omega_x x + \omega_y y)} dx dy$, 离散域 $H(k_x, k_y) = \frac{1}{MN} \sum_{n=0}^{M-1} \sum_{m=0}^{N-1} h(n, m) e^{-j2\pi(k_x n + k_y m)/N}$. 其中 M, N 分别表示图像的长和宽. $F[g * h] = F[g] F[h]$ 空域卷积等价于频域乘积! 为什么需要傅里叶变换? 使用快速傅里叶变换允许我们在采样上与内核大小无关的大型内核卷积; 此外, 转换到频域可以帮助我们理解滤波器的特性. 用傅里叶变换可以理解为多少高通滤波核式滤波更平滑.

全局运算: 高斯滤波后降采样: 原始图片, 先高斯滤波, 再删去偶数行, 偶数列, 获得图片 1. 对图片 1 先高斯, 再删去偶数行, 偶数列. 高斯金字塔的构建: 1) 滤波 2) 降采样 3) 重复 1,2 直到达到目标分辨率. 几何变换: 在 Point Process 是对 range of the image 做变换. $g(\mathbf{x}) = h(f(\mathbf{x}))$, 而这里我们要转换 domain, $g(\mathbf{x}) = f(h(\mathbf{x}))$. 参数变换是由少量参数控制的几何变换. 比如前边那些二维三维变换.

前向算法: 输入 f, h , 输出 g . For every pixel x in $f(x)$, 1) Compute the destination location $x' = h(x)$, 2) Copy the pixel $f(x)$ to $g(x')$. 前向算法存在的问题: x' 具有非整数数值时, 将像素 $f(x)$ 复制到 g 中的位置 x' 的过程没有得到很好的定义; 若将 x' 四舍五入到整数, 则将导致锯齿. 一些 x' 没有赋值, 导致裂缝和孔洞; 用邻近填补洞会产生模糊和混叠. **后向算法:** 输入 f, h , 输出 g . For every pixel x' in $g(x')$, 1) Compute the source location $x = h^{-1}(x')$, 2) Resample $f(x)$ at location x and copy to $g(x')$.

模型融合与优化

Regularization: $\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \|\mathbf{x}_n^T \mathbf{w} - t_n\|^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$. 选择合适的 λ 约束权重, 可以避免欠拟合或过拟合.

高斯分布 $N(\mathbf{x} | \mu, \Sigma) = \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1}(\mathbf{x} - \mu)\right)$, MLE $\mu_{ML} = \frac{1}{N} \sum_{n=1}^N x_n$, $\Sigma_{ML}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{ML})(x_n - \mu_{ML})^T$.

In more detail, the interpolated function f is a weighted sum (or superposition) of basis functions centered at each input data point $\mathbf{f}(\mathbf{x}) = \sum_k w_k \phi(\|\mathbf{x} - \mathbf{x}_k\|) = \mathbf{w}^T \phi(\mathbf{x})$ where the \mathbf{x}_k are the locations of the scattered data points, the ϕ are the radial basis functions (or kernels), and \mathbf{w}_k are the local weights associated with each kernel. 即 $\mathbf{w} = (w_1, w_2, \dots, w_N)^T$ 为基函数核. $\phi(\mathbf{x}) = (\phi(\|\mathbf{x} - \mathbf{x}_1\|), \phi(\|\mathbf{x} - \mathbf{x}_2\|), \dots, \phi(\|\mathbf{x} - \mathbf{x}_N\|))^T$ 为数据点附近的径向基函数. 拥有该特性的函数可作为径向基函数: $\phi(\|\mathbf{x}\|) = \phi(x)$. 即函数取值仅取决于径向距离. 如 Gaussian: $\phi(r) = \exp(-r^2/c^2)$. 多二次基函数 $\phi(r) = \sqrt{r^2 + c^2}$. 逆二次基函数: $\phi(r) = \frac{1}{\sqrt{r^2 + c^2}}$. 薄板样条函数: $\phi(r) = r^2 \ln(r)$.

If we want our function to exactly interpolate the data values, we solve the linear system of equations $\mathbf{f}(\mathbf{x}_k) = \sum_i \mathbf{w}_i \phi(\|\mathbf{x}_k - \mathbf{x}_i\|) = \mathbf{d}_k$ to obtain the desired set of weights \mathbf{w}_k . 但权值 \mathbf{w} 的求解条件可能是病态的: 数据位置或数据值的微小改变可能引起插值函数的巨大变化. 正则化: $E(\{\mathbf{w}_k\}) = E_D + \lambda E_W = \sum_k \|\mathbf{y}_k - \mathbf{f}(\mathbf{x}_k)\|^2 - \mathbf{d}_k^T \mathbf{I}^2 + \lambda \sum_k \|\mathbf{w}_k\|^2$. $p = 2$ 时, 得到纯最小二乘问题. 可以用 normal equations 求解: $\mathbf{w} = (\Phi^T \Phi + \lambda I)^{-1} \Phi^T \mathbf{D}$. $p = 1$ 时, 称为 Lasso 正则化. 权重项趋近于零, 使用较为稀疏的基函数集合.

稳健的数据拟合: 比起使用二次惩罚, 我们可以使用一个鲁棒损失函数 $\rho, E_R = \sum_k \rho(\|\mathbf{r}_k\|)$, 其中 $\mathbf{r}_k = \mathbf{y}_k - \mathbf{f}(\mathbf{x}_k) - \mathbf{d}_k$. 该损失函数给较大的数据拟合误差赋予较低的权重. 增加训练集数量; 降低模型复杂度; 增加正则化约束项; 集成学习; Dropout; BatchNorm; early stopping; 欠拟合: 增加特征数; 增加模型复杂度; 减小正则化系数.

深度学习

感知器: Minsky 证明感知器本质上是一种线性模型. 只能处理线性分类问题. 就连最简单的 XOR(异或)问题都无法正确分类. 可以使用二层感知器来解决. 一种思路: $\mathbf{y}_1 = \text{ReLU}(\mathbf{x}_1 - x_2)$, $\mathbf{y}_2 = \text{ReLU}(\mathbf{x}_2 - x_1)$, $\mathbf{y}_3 = \text{ReLU}(\mathbf{y}_1 + \mathbf{y}_2)$. 或者 $A \cdot B = (A \cdot \mathbf{A} - B) \cdot V = (A - A \cdot B)$. 则用两层感知器处理这个问题: $A \wedge B = \text{sgn}_1(f_1(A, B)), f_1(A, B) = w_0 + w_1 A + w_2 B, A \cdot B \cdot \text{sgn}_2 \text{ xor } (f_2(A, B)), f_2(A, B) = w_0 + w_1 \text{sgn}_1(f_1(A, B)) + w_2 \text{sgn}_1(f_1(A, B))$. 没有隐藏层的网络在它们可以学习建模的输入输出映射方面非常有限. 更多的线性单元层并没有帮助, 它仍然是线性的. 固定输出的非线性是不够的.

激活函数: ReLU: 线性整流函数. $\max(0, x)$. Sigmoid: $\sigma(x) = \frac{1}{1+e^{-x}}$. tanh: $\tanh(x)$. Leaky ReLU: $\max(0, 1.5x)$.

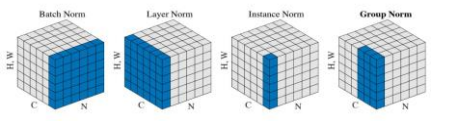
Softmax function: $p(C_k | \mathbf{x}) = \frac{p(\mathbf{x}|C_k)p(C_k)}{\sum_{j=1}^J p(\mathbf{x}|C_j)p(C_j)} = \frac{\exp(a_k)}{\sum_{j=1}^J \exp(a_j)}$. 求导: $\nabla_{\mathbf{x}} \mathbf{A} \mathbf{x} = \mathbf{A}^T$, $\nabla_{\mathbf{x}} \mathbf{x}^T \mathbf{A} = \mathbf{A}$, $\nabla_{\mathbf{x}} \mathbf{x}^T \mathbf{A} \mathbf{x} = (\mathbf{A} + \mathbf{A}^T) \mathbf{x}$. 若 $\mathbf{y} = \mathbf{A} \mathbf{x}$, 则 $\frac{\partial \mathbf{y}}{\partial \mathbf{A}} = \mathbf{x}^T$. $\tanh'(x) = 1 - \tanh^2(x)$. sigmoid $(x) = \text{sigmoid}(x) \cdot (1 - \text{sigmoid}(x))$

$f(x) = \frac{1}{1+e^{-x}}$, $f'(x) = \frac{e^{-x}}{(1+e^{-x})^2}$, $g(x) = \frac{e^{-x}-e^x}{(1+e^{-x})^2}$, $g'(x) = \frac{4}{(1+e^{-x})^3}$

梯度下降: BGD: $L(W) = \frac{1}{N} \sum_{i=1}^N L_i(x_i, y_i, W) + \lambda R(W)$; $\nabla_W L(W) = \frac{1}{N} \sum_{i=1}^N \nabla_W L_i(x_i, y_i, W) + \lambda \nabla_W R(W)$. 影响因素: 权值初始化方式、步数、学习率、batch 尺寸、data sampling. 当 N 很大时, 开始计算大 SGD 可能存在的问题: 如果损失在一个方向上快速变化, 在另一个方向上缓慢变化会怎样? 如果损失函数存在局部最小值或鞍点, 则梯度为 0, 梯度下降卡住. 此外, 我们的梯度来至于小批量数据. 因此可能会有噪声. 随机梯度下降 (SGD) 比批量梯度下降 (BGD) 更快地接近最小值, 但它可能永远无法“收敛”到最小值.

优化器: SGD: $x_{t+1} = x_t - \alpha \nabla f(x_t)$; SGD+Momentum: $v_{t+1} = \rho v_t + \nabla f(x_t)$, $x_{t+1} = x_t - \alpha v_{t+1}$; Nesterov: $v_{t+1} = \rho v_t - \alpha \nabla f(x_t + \rho v_t)$, $x_{t+1} = x_t + v_{t+1}$; AdaGrad; RMSProp; Adam 约等于 RMSProp + Momentum. Adam 满足 momentum, adaptive learning rates; Leaky second moments; bias correction for moment estimates.

归一化: $\mu_i = \frac{1}{|S|} \sum_{n \in B} s_i^{(n)}$, $\sigma_i^2 = \frac{1}{|S|} \sum_{n \in B} (s_i^{(n)} - \mu_i)^2$, $s_i^{(n)} = \frac{s_i^{(n)} - \mu_i}{\sigma_i^2 + \epsilon}$, $y_i = \gamma \hat{s}_i + \beta_i$



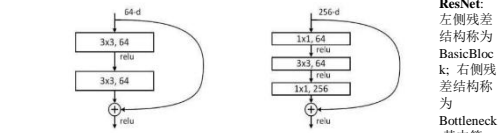
数据增强: Distorted; Cutout; Mixup.

CNN 特点一: 稀疏连接; **特点二:** 参数共享; **特点三:** 等变表示. Pooling 可以分为 Max-pooling(返回池化窗口内最大值); Average pooling; Stochastic pooling.

多通道卷积: 图像有多个颜色通道, 各通道分别进行卷积, 再把卷积的结果相加. 一个卷积核提取一个特征, 多个卷积核产生多个特征. 每个特征对应于一个通道. 原始图 $N \times N$, 卷积核 $m \times m$, 卷积核 $(N-m+1) \times (N-m+1)$. 池化: 最大池化, 平均池化, 随机池化.

令 $O =$ 输出图像的尺寸, $I =$ 输入图像的尺寸, $K =$ 卷积核的核尺寸, $S =$ 移动步长, $P =$ 填充数. 则输出图像尺寸 $O = \frac{I - K + 2P}{S} + 1$. 严格来说这个尺寸是要向下取整的

参数个数: (核长 x 核宽 x 输入通道数) x 输出通道数. AlexNet: VGGNet 大的感受野被 (之间带有 ReLU 的) 连续的 3×3 卷积层所取代. 相似的感受野, 更少的参数, 更多的非线性; GoogleNet: 1×1 3×3 5×5 convolutions 并联. 具有不同感受野大小和操作的并行路径旨在捕获特征图重要性中稀疏相关模式. Inception v2, v3 使用批量归一化规范训练, 降低辅助分类器的权重; ResNet 引入 skip or shortcut connections.



层的 1×1 的卷积核的作用是对特征矩阵进行降维操作, 将特征矩阵的深度由 256 降为 64. 第三层的 1×1 的卷积核是对特征矩阵进行升维操作, 将特征矩阵的深度由 64 升为 256. 降低特征矩阵的深度主要是为了减少参数的个数. 如果采用 BasicBlock, 参数的个数应该是: $256 \times 256 \times 3 \times 3 \times 2 = 1179648$. 采用 Bottleneck, 参数的个数是: $1 \times 1 \times 256 \times 64 + 3 \times 3 \times 64 \times 64 + 1 \times 1 \times 256 \times 64 = 69632$

识别

图像分类: 输入图片, 输出猫, 单个物体, 没有空间范围. 语义分割, 没有物体, 只有像素, 输入图片, 分割草, 树, 天空. 目标检测: 输入图片, 用框把猫和猫框住. 实例分割: 输入图片, 把猫和猫用不同颜色覆盖. 图像分析: 基于特征的方法: 词袋(Bag of words). 在关键点处提取特征; 量化: 获取视图词 (特征聚类中心) 上的直方图; 用分类算法如支持向量机对特征分布直方图上进行决策. 基于部件的模型 (Part-based model). 原理: 通过寻找物体的组成部分并测量它们之间的几何关系来识别物体.

目标检测: 输入: 单张 RGB 图像; 输出: 检测物体的集合 (每个物体的预测内窗: 类别标签 (属于已知类别集合); 检测框 (Bounding Box) 的位置 (x, y) + 形状 (width, height). 挑战: 多输出, 多种类输出 (需要同时预测 what and where), 图像尺寸大. **人脸识别** 算法在鲁棒性上需要考虑: 光照、表情、尺度、姿态.

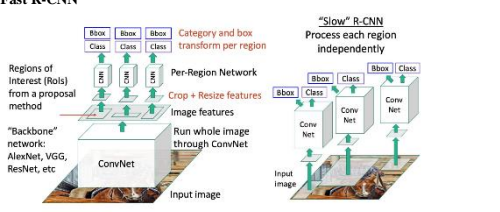
检测单个物体: 将定位看成是一个回归问题. 网络在 ImageNet 上预训练 (迁移学习), 之后一部分解决 “What” 的问题, 输出一个 class score, and correct label 计算出 softmax loss, 另一部分解决 “where” 的问题, 计算出 Box Coordinates, and Correct Box 计算出一个 L2 Loss. 两个 Loss 进行 Weighted Sum 之后计算出一个 MultiTask Loss.

检测多个物体: 滑动 (Sliding Window). 将不同的图像块输入卷积神经网络中. 网络将每一个图像块分类为物体或背景. 问题: 在一张尺寸为 $H \times W$ 的图像中有多个可能的检测框? 假设检测框的大小为 $h \times w$, 可能的位置: $W - w + 1$, 可能的位置: $H - h + 1$, 可能的位置 $(W - w + 1) \times (H - h + 1)$. 则所有可能的检测框数: $\sum_{h=1}^H \sum_{w=1}^W (W - w + 1) (H - h + 1) = \frac{H(H+1)}{2} \frac{W(W+1)}{2} \approx 800 \times 600$ 大小的图像有约 58M 个检测框! 全部计算是不现实的!

候选区域 (Region Proposals): 找到几乎包括所有物体的检测框的集合. 通常基于启发式 (heuristics): 如搜索“斑点状”的图像区域. 运行速度相对较快. 选择性搜索 (selective search) 在 CPU 上 几秒钟内可以给出 2000 个候选框.

R-CNN: Region-based CNN. 输入: 单张 RGB 图像. 1) 运行候选区域算法, 计算得到约 2000 个候选区域; 2) 将每个区域大小调整为 224×224 , 并分别输入 CNN 中预测其类别得分和检测框变换参数 (Forward each region through ConvNet). 3) 根据分类得分选择候选区域子集进行输出 (有多种选择: 背景图例, 类别图例, 或者输出每张图像得分前 K 个). 4) 跟真实的检测框标注进行比较. 问题: 非常慢! 对于每张图像需要约 2000 次的向前计算 (forward pass). 解决方案: 在 warping 前经过 CNN

Fast R-CNN



Per-Region 网络是相对轻量级的. 大多数计算在骨干网络 (Backbone Network) 中进行; 为重叠候选区域 (region proposals) 减少了工作量. AlexNet 作检测, 5 层卷积层作为 backbone, 2 层全连接层用于 per-region network; ResNet 作检测, 最后三个阶段用于 per-region network. 其余的部分用于 backbone. Fast R-CNN 运行时间被候选区域主导.

Fast R-CNN: 加入 Region Proposal Network (RPN) 来从特征中预测候选区域; 其他部分跟 Fast R-CNN 相同. Faster R-CNN 是一个 Two-stage object detector. 第一阶段, 每张图运行一次: Backbone Network; Region Proposal Network; 第二阶段: 每个区域运行一次. 表明特征: RoI Pool/Align; 预测物体类别; 预测检测框的偏移量. Fast R-CNN 计算时间: $\text{PropTime} + 1 \times \text{ConvTime} + \text{NumProp} \times \text{fTime}$

Feature Pyramid Networks (FPN). 如何提高多尺度目标的检测效果? Featureized image pyramid. 生成图片金字塔, 对不同尺寸的图片分别提取特征: Single feature map: 在网络的最后一层的特征图上进行提取 (SPFPNet) (Fast R-CNN, Faster R-CNN 等采用); Pyramid feature hierarchy: 从网络不同层抽取不同尺度的特征. 在不同尺度的特征上分别进行检测 (SSD 采用的多尺度融合策略). FPN: 为了使不同尺度的特征都包含丰富的语义信息, 同时又不复杂计算成本, FPN 采用 bottom-up, top-down and lateral connection 的方式, 让底层高分辨率 低语义的特征和高层低分辨率高语义的特征融合在一起, 使得最终得到不同尺度的特征图都有丰富的语义信息.

