

欧氏平面上平方距离最小的直线拟合

BY JIN

2019-06-10

摘要

从平面上的一族点拟合一条直线是一个常见问题。然而统计上常用的针对 $y = kx + b$ 的最小二乘优化 $\{k = (n\sum xy - \sum x\sum y) / (n\sum x^2 - (\sum x)^2) = \text{cov}(x, y) / \text{cov}(x, x), b = (\sum y - k\sum x) / n = \bar{y} - k\bar{x}\}$, 最小化的是 y 方向上的误差 $\sum (y_i - kx_i - b)^2$, 并不是欧几里得平面上的距离最小。本文将针对点到直线的垂直距离, 使用完全最小二乘法求最优, 得出为更优雅的结果, 并自多种角度对这一结果进行解释和验证。

1 问题描述

待拟合的直线用法向式表示为

$$x \cos\alpha + y \sin\alpha - r = 0 \quad (1)$$

则任意一点 (x, y) 到该直线的距离为

$$d = |x \cos\alpha + y \sin\alpha - r| \quad (2)$$

则问题是求已知点集 $\{(x, y)\}_n$, 以最小化 $F(\alpha, r) = \sum d^2$ 为目标, 求最优解

$$\underset{\alpha, r}{\text{ArgMin}} \sum_i (x_i \cos\alpha + y_i \sin\alpha - r)^2 \quad (3)$$

2 求解最优

记 $\cos\alpha = c, \sin\alpha = s$, 待优化的和式有如下展开

$$\begin{aligned} F(\alpha, r) &= \sum (x \cos\alpha + y \sin\alpha - r)^2 \\ &= \sum (c^2 x^2 + s^2 y^2 + 2csxy - 2rcx - 2rsy + r^2) \\ &= c^2 \sum x^2 + s^2 \sum y^2 + 2cs \sum xy - 2rc \sum x - 2rs \sum y + nr^2 \end{aligned} \quad (4)$$

求偏导, 有

$$\begin{aligned} 0 = \frac{\partial F}{\partial r} &= -2(c \sum x + s \sum y - nr) \\ \text{即 } r &= (c \sum x + s \sum y) / n = c\bar{x} + s\bar{y} \end{aligned} \quad (5)$$

以及

$$\begin{aligned} 0 = \frac{\partial F}{\partial \alpha} &= -2(cs(\sum x^2 - \sum y^2) - (c^2 - s^2)\sum xy - r(s\sum x - c\sum y)) \\ \text{代入 } r, \quad 0 &= cs(\bar{x}^2 - \bar{y}^2) - (c^2 - s^2)\bar{x}\bar{y} - (c\bar{x} + s\bar{y})(s\bar{x} - c\bar{y}) \\ &= cs(\bar{x}^2 - \bar{y}^2) - (c^2 - s^2)\bar{x}\bar{y} - cs(\bar{x}^2 - \bar{y}^2) + (c^2 - s^2)\bar{x}\bar{y} \\ &= cs(\text{cov}(x, x) - \text{cov}(y, y)) - (c^2 - s^2)\text{cov}(x, y) \end{aligned} \quad (6)$$

整理, 得

$$\frac{2 \text{cov}(x, y)}{\text{cov}(x, x) - \text{cov}(y, y)} = \frac{2 \cos\alpha \sin\alpha}{\cos^2\alpha - \sin^2\alpha} = \tan 2\alpha \quad (7)$$

由此解出

$$\alpha = \frac{1}{2} \arctan \frac{2 \operatorname{cov}(x, y)}{\operatorname{cov}(x, x) - \operatorname{cov}(y, y)} + \frac{k\pi}{2} \quad (8)$$

包含相垂直的两个解，正好对应直线的方向向量 θ 和法向量方向 α ，也是误差函数的最大和最小值。

3 分析验根

统计角度上看， $\operatorname{cov}(x, x) - \operatorname{cov}(y, y)$ 的符号表示点集在x, y方向上的分散程度大小，而 $\operatorname{cov}(x, y)$ 的符号表示了点集中x, y是正相关还是负相关：据此可以确定 α 的正确取值。

1. 当 $\operatorname{cov}(x, x) - \operatorname{cov}(y, y) > 0, \operatorname{cov}(x, y) > 0$ ，对应 $0 < \theta < \frac{\pi}{4}, -\frac{\pi}{2} < \alpha < -\frac{\pi}{4}$ ，因此

$$\begin{aligned}\theta &= \frac{1}{2} \arctan 2(2 \operatorname{cov}(x, y), \operatorname{cov}(x, x) - \operatorname{cov}(y, y)) \\ \alpha = \theta \pm \frac{\pi}{2} &= \frac{1}{2} \arctan 2(-2 \operatorname{cov}(x, y), \operatorname{cov}(y, y) - \operatorname{cov}(x, x))\end{aligned}$$

2. 当 $\operatorname{cov}(x, x) - \operatorname{cov}(y, y) < 0, \operatorname{cov}(x, y) > 0$ ，对应 $\frac{\pi}{4} < \theta < \frac{\pi}{2}, -\frac{\pi}{4} < \alpha < 0$ ，因此仍有

$$\begin{aligned}\theta &= \frac{1}{2} \arctan 2(2 \operatorname{cov}(x, y), \operatorname{cov}(x, x) - \operatorname{cov}(y, y)) \\ \alpha = \theta \pm \frac{\pi}{2} &= \frac{1}{2} \arctan 2(-2 \operatorname{cov}(x, y), \operatorname{cov}(y, y) - \operatorname{cov}(x, x))\end{aligned}$$

3. 同理（或由连续性）可证 $\operatorname{cov}(x, x) - \operatorname{cov}(y, y) < 0, \operatorname{cov}(x, y) < 0, -\frac{\pi}{2} < \theta < -\frac{\pi}{4}, 0 < \alpha < \frac{\pi}{4}$ 的情况，以及 $\operatorname{cov}(x, x) - \operatorname{cov}(y, y) > 0, \operatorname{cov}(x, y) < 0, -\frac{\pi}{4} < \theta < 0, \frac{\pi}{4} < \alpha < \frac{\pi}{2}$ 的情况也有同样的结论

因此统一地，有结论

$$\alpha = \frac{1}{2} \arctan 2(-2 \operatorname{cov}(x, y), \operatorname{cov}(y, y) - \operatorname{cov}(x, x)) \quad (9)$$

4 统计角度

按 $r = \bar{x} \cos \alpha + \bar{y} \sin \alpha$ ，若对点集 $X = \begin{bmatrix} \vdots & \vdots \\ x_i & y_i \\ \vdots & \vdots \end{bmatrix}_{n \times 2}$ 作平移变换 $X' = X - \bar{X}$ ，可将直线方程写成齐次形式：

$$x' \cos \alpha + y' \sin \alpha = 0 \quad (10)$$

再进行旋转变换 $X'' = X' R$, $R R^T = I$ ，此时新的协方差矩阵

$$\begin{aligned}\Sigma(X'') &= E((X' R - 0)^T (X' R - 0)) = R^T E(X'^T X') R \\ &= R^{-1} E((X - \bar{X})^T (X - \bar{X})) R \\ &= R^{-1} \Sigma(X) R\end{aligned} \quad (11)$$

为原协方差矩阵的一个相似变换。

由 X'' 为对称正定矩阵，故必存在恰当的正交变换 R 使其对角化，即

$$C(X'') = \begin{bmatrix} \sigma_{x''}^2 & 0 \\ 0 & \sigma_{y''}^2 \end{bmatrix} \quad (12)$$

如此可将点集的误差分解到彼此独立且正交的两个方向上。设 $\sigma_{x''}^2 > \sigma_{y''}^2$, 则 $\sigma_{x''}^2$ 对应了直线方向, $\sigma_{y''}^2$ 对应了直线的法线方向。由此引出算法:

1. 求点集 $\{(x, y)\}_n$ 的协方差矩阵 $C = \begin{bmatrix} \text{cov}(x, x) & \text{cov}(x, y) \\ \text{cov}(x, y) & \text{cov}(y, y) \end{bmatrix}$
2. 对 C 作特征分解 $C = T\Sigma'T^{-1}$, 其中 $C' = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}$, $\sigma_1^2 > \sigma_2^2$, $T = \begin{bmatrix} u & v \\ -v & u \end{bmatrix} = \begin{bmatrix} \cos\beta & \sin\beta \\ -\sin\beta & \cos\beta \end{bmatrix}$
3. 该特征分解对应了点集(围绕其均值)的一个旋转变换 $[x'', y''] = [x', y']T^{-1}$, 该变换将直线的方向向量旋转到了0度。而反过来 $\begin{bmatrix} x' \\ y' \end{bmatrix} = \begin{bmatrix} \cos\beta & -\sin\beta \\ \sin\beta & \cos\beta \end{bmatrix} \begin{bmatrix} x'' \\ y'' \end{bmatrix}$ 意味着 X'' 旋转 $-\beta$ 即可到达 X' 的角度, 因此直线方向为 $\theta = -\beta$, 法向 $\alpha = \pi - \beta$, 即 $\begin{bmatrix} \cos\alpha \\ \sin\alpha \end{bmatrix} = \begin{bmatrix} -\sin\theta \\ \cos\theta \end{bmatrix} = \begin{bmatrix} \sin\beta \\ \cos\beta \end{bmatrix} = \begin{bmatrix} v \\ u \end{bmatrix}$.

实践中, 没有必要对 $X^T X$ 做特征分解, 完全可以等效地对 X 做奇异值分解:

5 完全最小二乘(SVD分解)

同样按(10)的方式齐次化直线方程。令 $X' = \begin{bmatrix} \vdots & \vdots \\ x'_i & y'_i \\ \vdots & \vdots \end{bmatrix}_{n \times 2}$, $A = \begin{bmatrix} \cos\alpha \\ \sin\alpha \end{bmatrix}$, 此时优化目标可以写成

$$E(X) = (XA)^T(XA) = A^T X^T X A, \quad (13)$$

求 $A = \text{Argmin}_A E(X)$ 。对 X 作SVD分解: $X = U_{n \times 2} \Sigma_{2 \times 2} V_{2 \times 2}^T$

$$E(X) = A^T V \Sigma^T U^T U \Sigma V^T A = A^T V \Sigma^T \Sigma V^T A \quad (14)$$

设 $V = \begin{pmatrix} \cos\theta & -\sin\theta \\ \sin\theta & \cos\theta \end{pmatrix}$, $\Sigma = \begin{pmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{pmatrix}$, $\sigma_1 \geq \sigma_2 \geq 0$, 可直接求出

$$E(X) = \frac{1}{2}(\sigma_1^2 + \sigma_2^2 + (\sigma_1^2 - \sigma_2^2)\cos(2\alpha - 2\theta)) \quad (15)$$

当 $2(\alpha - \theta) = \pi + 2k\pi$ 即 $\alpha = \theta + \frac{\pi}{2} + k\pi$ 时取最小值。考虑到 $\alpha + \pi$ 仅仅是 α 的 180° 反方向, 故只取一个解 $\alpha = \theta + \frac{\pi}{2}$ 即可。此时

$$A = \begin{bmatrix} \cos\alpha \\ \sin\alpha \end{bmatrix} = \begin{bmatrix} -\sin\theta \\ \cos\theta \end{bmatrix} = V[2] \quad (16)$$

正好与特征分解算法的结果一致。

另一方面, 按Total Least Squares求解, 也能有:

$$A = \text{Argmin}_A \|\tilde{X}\|_F, \quad (X + \tilde{X})A = 0 \quad (17)$$

SVD分解

$$\begin{aligned} X &= [U_1 \ U_2] \begin{bmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{bmatrix} [V_1 \ V_2]^T \\ X + \tilde{X} &= [U_1 \ U_2] \begin{bmatrix} \sigma_1 & 0 \\ 0 & 0 \end{bmatrix} [V_1 \ V_2]^T = U_1 \sigma_1 V_1^T \\ \text{两式相减 } \tilde{X} &= -U_2 \sigma_2 V_2^T \end{aligned} \quad (18)$$

由此也有, 在 $X + \tilde{X}$ 的补空间中取单位向量

$$A = \pm V_2$$

可满足 $(X + \tilde{X})A = U_1 \sigma_1 V_1^T \cdot V_2 = U_1 \sigma_1 (V_1^T V_2) = 0$ 。

Q.E.D.