

实验报告

预处理与词频统计

首先是文本的预处理，在这一阶段我并没有使用 `textblob` 进行分词，而是参照网上的方法重新将文本打印了一遍，同时打印的时候每行只打印一个单词，程序中使用 `createFiles(oriPath, path)` 完成这一过程，在这一方法中调用了 `createProcessFile` 将单词分行打印，这样在后续读取时也可以达到分词的效果。

之后调用 `lineProcess` 方法对文本单词进行处理，使用 `nlTK` 中现成的方法去停用词、提取词干，将非英文单词的部分去除，使用 `countWords(path)` 进行词频统计。

计算 IDF-TF 值，建立向量空间

统计所有 Train 文档中出现的词干，多滤掉低频词，其他单词使用如下格式保存

```
Doc = {} # <word, (doc1,...,docN)>
```

```
IDF = {} # <word, IDF>
```

使用 `CreateVSM(path)` 方法建立空间向量表示，每个训练数据都对应一个向量，向量形式以 `<category, document, (word1, para1), (word2, para2),...>` 存入文件。

KNN 实现

将测试集中每个文档的 `{word, TF}` 与训练集 `<类_文件名, <word, TFIDF>>` 进行比较，通过余弦函数找到最相近的 20 个文档，将其按照距离远近加权，通过 20 个文档所属类别计算出与测试文档向量距离和最小的类。

最终准确率在 0.78 左右

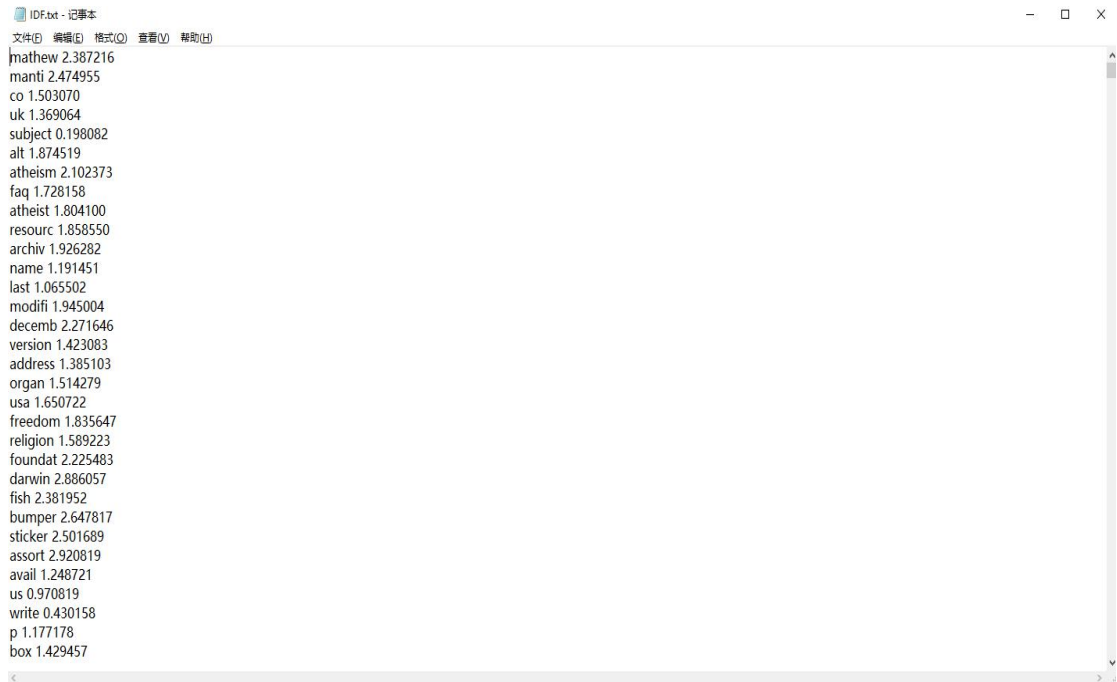
词频统计



词频统计结果 (部分示例):

- ax: 62551, edu: 29228, subject: 20657, com: 17770, x: 16554, one: 15386, use: 15133, write: 14920, would: 14805, c: 12306, article: 11887, ani: 11544, r: 11147, w: 11147, s: 11147, t: 11147, e: 11147, n: 11147, d: 11147, o: 11147, a: 11147, i: 11147, u: 11147, y: 11147, p: 11147, h: 11147, m: 11147, l: 11147, k: 11147, j: 11147, q: 11147, z: 11147, v: 11147, b: 11147, f: 11147, g: 11147, x: 11147, c: 11147, d: 11147, e: 11147, f: 11147, g: 11147, h: 11147, i: 11147, j: 11147, k: 11147, l: 11147, m: 11147, n: 11147, o: 11147, p: 11147, q: 11147, r: 11147, s: 11147, t: 11147, u: 11147, v: 11147, w: 11147, x: 11147, y: 11147, z: 11147, aa: 11147, ab: 11147, ac: 11147, ad: 11147, ae: 11147, af: 11147, ag: 11147, ah: 11147, ai: 11147, aj: 11147, ak: 11147, al: 11147, am: 11147, an: 11147, ao: 11147, ap: 11147, aq: 11147, ar: 11147, as: 11147, at: 11147, au: 11147, av: 11147, aw: 11147, ax: 11147, ay: 11147, az: 11147, ba: 11147, bb: 11147, bc: 11147, bd: 11147, be: 11147, bf: 11147, bg: 11147, bh: 11147, bi: 11147, bj: 11147, bk: 11147, bl: 11147, bm: 11147, bn: 11147, bo: 11147, bp: 11147, bq: 11147, br: 11147, bs: 11147, bt: 11147, bu: 11147, bv: 11147, bw: 11147, bx: 11147, by: 11147, bz: 11147, ca: 11147, cb: 11147, cc: 11147, cd: 11147, ce: 11147, cf: 11147, cg: 11147, ch: 11147, ci: 11147, cj: 11147, ck: 11147, cl: 11147, cm: 11147, cn: 11147, co: 11147, cp: 11147, cq: 11147, cr: 11147, cs: 11147, ct: 11147, cu: 11147, cv: 11147, cw: 11147, cx: 11147, cy: 11147, cz: 11147, da: 11147, db: 11147, dc: 11147, dd: 11147, de: 11147, df: 11147, dg: 11147, dh: 11147, di: 11147, dj: 11147, dk: 11147, dl: 11147, dm: 11147, dn: 11147, do: 11147, dp: 11147, dq: 11147, dr: 11147, ds: 11147, dt: 11147, du: 11147, dv: 11147, dw: 11147, dx: 11147, dy: 11147, dz: 11147, ea: 11147, eb: 11147, ec: 11147, ed: 11147, ee: 11147, ef: 11147, eg: 11147, eh: 11147, ei: 11147, ej: 11147, ek: 11147, el: 11147, em: 11147, en: 11147, eo: 11147, ep: 11147, eq: 11147, er: 11147, es: 11147, et: 11147, eu: 11147, ev: 11147, ew: 11147, ex: 11147, ey: 11147, ez: 11147, fa: 11147, fb: 11147, fc: 11147, fd: 11147, fe: 11147, ff: 11147, fg: 11147, fh: 11147, fi: 11147, fj: 11147, fk: 11147, fl: 11147, fm: 11147, fn: 11147, fo: 11147, fp: 11147, fq: 11147, fr: 11147, fs: 11147, ft: 11147, fu: 11147, fv: 11147, fw: 11147, fx: 11147, fy: 11147, fz: 11147, ga: 11147, gb: 11147, gc: 11147, gd: 11147, ge: 11147, gf: 11147, gg: 11147, gh: 11147, gi: 11147, gj: 11147, gk: 11147, gl: 11147, gm: 11147, gn: 11147, go: 11147, gp: 11147, gq: 11147, gr: 11147, gs: 11147, gt: 11147, gu: 11147, gv: 11147, gw: 11147, gx: 11147, gy: 11147, gz: 11147, ha: 11147, hb: 11147, hc: 11147, hd: 11147, he: 11147, hf: 11147, hg: 11147, hh: 11147, hi: 11147, hj: 11147, hk: 11147, hl: 11147, hm: 11147, hn: 11147, ho: 11147, hp: 11147, hq: 11147, hr: 11147, hs: 11147, ht: 11147, hu: 11147, hv: 11147, hw: 11147, hx: 11147, hy: 11147, hz: 11147, ia: 11147, ib: 11147, ic: 11147, id: 11147, ie: 11147, if: 11147, ig: 11147, ih: 11147, ii: 11147, ij: 11147, ik: 11147, il: 11147, im: 11147, in: 11147, io: 11147, ip: 11147, iq: 11147, ir: 11147, is: 11147, it: 11147, iu: 11147, iv: 11147, iw: 11147, ix: 11147, iy: 11147, iz: 11147, ja: 11147, jb: 11147, jc: 11147, jd: 11147, je: 11147, jf: 11147, jg: 11147, jh: 11147, ji: 11147, jj: 11147, jk: 11147, jl: 11147, jm: 11147, jn: 11147, jo: 11147, jp: 11147, jq: 11147, jr: 11147, js: 11147, jt: 11147, ju: 11147, jv: 11147, jw: 11147, jx: 11147, jy: 11147, jz: 11147, ka: 11147, kb: 11147, kc: 11147, kd: 11147, ke: 11147, kf: 11147, kg: 11147, kh: 11147, ki: 11147, kj: 11147, kk: 11147, kl: 11147, km: 11147, kn: 11147, ko: 11147, kp: 11147, kq: 11147, kr: 11147, ks: 11147, kt: 11147, ku: 11147, kv: 11147, kw: 11147, kx: 11147, ky: 11147, kz: 11147, la: 11147, lb: 11147, lc: 11147, ld: 11147, le: 11147, lf: 11147, lg: 11147, lh: 11147, li: 11147, lj: 11147, lk: 11147, ll: 11147, lm: 11147, ln: 11147, lo: 11147, lp: 11147, lq: 11147, lr: 11147, ls: 11147, lt: 11147, lu: 11147, lv: 11147, lw: 11147, lx: 11147, ly: 11147, lz: 11147, ma: 11147, mb: 11147, mc: 11147, md: 11147, me: 11147, mf: 11147, mg: 11147, mh: 11147, mi: 11147, mj: 11147, mk: 11147, ml: 11147, mm: 11147, mn: 11147, mo: 11147, mp: 11147, mq: 11147, mr: 11147, ms: 11147, mt: 11147, mu: 11147, mv: 11147, mw: 11147, mx: 11147, my: 11147, mz: 11147, na: 11147, nb: 11147, nc: 11147, nd: 11147, ne: 11147, nf: 11147, ng: 11147, nh: 11147, ni: 11147, nj: 11147, nk: 11147, nl: 11147, nm: 11147, nn: 11147, no: 11147, np: 11147, nq: 11147, nr: 11147, ns: 11147, nt: 11147, nu: 11147, nv: 11147, nw: 11147, nx: 11147, ny: 11147, nz: 11147, oa: 11147, ob: 11147, oc: 11147, od: 11147, oe: 11147, of: 11147, og: 11147, oh: 11147, oi: 11147, oj: 11147, ok: 11147, ol: 11147, om: 11147, on: 11147, oo: 11147, op: 11147, oq: 11147, or: 11147, os: 11147, ot: 11147, ou: 11147, ov: 11147, ow: 11147, ox: 11147, oy: 11147, oz: 11147, pa: 11147, pb: 11147, pc: 11147, pd: 11147, pe: 11147, pf: 11147, pg: 11147, ph: 11147, pi: 11147, pj: 11147, pk: 11147, pl: 11147, pm: 11147, pn: 11147, po: 11147, pp: 11147, pq: 11147, pr: 11147, ps: 11147, pt: 11147, pu: 11147, pv: 11147, pw: 11147, px: 11147, py: 11147, pz: 11147, qa: 11147, qb: 11147, qc: 11147, qd: 11147, qe: 11147, qf: 11147, qg: 11147, qh: 11147, qi: 11147, qj: 11147, qk: 11147, ql: 11147, qm: 11147, qn: 11147, qo: 11147, qp: 11147, qq: 11147, qr: 11147, qs: 11147, qt: 11147, qu: 11147, qv: 11147, qw: 11147, qx: 11147, qy: 11147, qz: 11147, ra: 11147, rb: 11147, rc: 11147, rd: 11147, re: 11147, rf: 11147, rg: 11147, rh: 11147, ri: 11147, rj: 11147, rk: 11147, rl: 11147, rm: 11147, rn: 11147, ro: 11147, rp: 11147, rq: 11147, rr: 11147, rs: 11147, rt: 11147, ru: 11147, rv: 11147, rw: 11147, rx: 11147, ry: 11147, rz: 11147, sa: 11147, sb: 11147, sc: 11147, sd: 11147, se: 11147, sf: 11147, sg: 11147, sh: 11147, si: 11147, sj: 11147, sk: 11147, sl: 11147, sm: 11147, sn: 11147, so: 11147, sp: 11147, sq: 11147, sr: 11147, ss: 11147, st: 11147, su: 11147, sv: 11147, sw: 11147, sx: 11147, sy: 11147, sz: 11147, ta: 11147, tb: 11147, tc: 11147, td: 11147, te: 11147, tf: 11147, tg: 11147, th: 11147, ti: 11147, tj: 11147, tk: 11147, tl: 11147, tm: 11147, tn: 11147, to: 11147, tp: 11147, tq: 11147, tr: 11147, ts: 11147, tt: 11147, tu: 11147, tv: 11147, tw: 11147, tx: 11147, ty: 11147, tz: 11147, ua: 11147, ub: 11147, uc: 11147, ud: 11147, ue: 11147, uf: 11147, ug: 11147, uh: 11147, ui: 11147, uj: 11147, uk: 11147, ul: 11147, um: 11147, un: 11147, uo: 11147, up: 11147, uq: 11147, ur: 11147, us: 11147, ut: 11147, uu: 11147, uv: 11147, uw: 11147, ux: 11147, uy: 11147, uz: 11147, va: 11147, vb: 11147, vc: 11147, vd: 11147, ve: 11147, vf: 11147, vg: 11147, vh: 11147, vi: 11147, vj: 11147, vk: 11147, vl: 11147, vm: 11147, vn: 11147, vo: 11147, vp: 11147, vq: 11147, vr: 11147, vs: 11147, vt: 11147, vu: 11147, vv: 11147, vw: 11147, vx: 11147, vy: 11147, vz: 11147, wa: 11147, wb: 11147, wc: 11147, wd: 11147, we: 11147, wf: 11147, wg: 11147, wh: 11147, wi: 11147, wj: 11147, wk: 11147, wl: 11147, wm: 11147, wn: 11147, wo: 11147, wp: 11147, wq: 11147, wr: 11147, ws: 11147, wt: 11147, wu: 11147, wv: 11147, ww: 11147, wx: 11147, wy: 11147, wz: 11147, xa: 11147, xb: 11147, xc: 11147, xd: 11147, xe: 11147, xf: 11147, xg: 11147, xh: 11147, xi: 11147, xj: 11147, xk: 11147, xl: 11147, xm: 11147, xn: 11147, xo: 11147, xp: 11147, xq: 11147, xr: 11147, xs: 11147, xt: 11147, xu: 11147, xv: 11147, xw: 11147, xx: 11147, xy: 11147, xz: 11147, ya: 11147, yb: 11147, yc: 11147, yd: 11147, ye: 11147, yf: 11147, yg: 11147, yh: 11147, yi: 11147, yj: 11147, yk: 11147, yl: 11147, ym: 11147, yn: 11147, yo: 11147, yp: 11147, yq: 11147, yr: 11147, ys: 11147, yt: 11147, yu: 11147, yv: 11147, yw: 11147, yx: 11147, yy: 11147, yz: 11147, za: 11147, zb: 11147, zc: 11147, zd: 11147, ze: 11147, zf: 11147, zg: 11147, zh: 11147, zi: 11147, zj: 11147, zk: 11147, zl: 11147, zm: 11147, zn: 11147, zo: 11147, zp: 11147, zq: 11147, zr: 11147, zs: 11147, zt: 11147, zu: 11147, zv: 11147, zw: 11147, zx: 11147, zy: 11147, zz: 11147, aa: 11147, ab: 11147, ac: 11147, ad: 11147, ae: 11147, af: 11147, ag: 11147, ah: 11147, ai: 11147, aj: 11147, ak: 11147, al: 11147, am: 11147, an: 11147, ao: 11147, ap: 11147, aq: 11147, ar: 11147, as: 11147, at: 11147, au: 11147, av: 11147, aw: 11147, ax: 11147, ay: 11147, az: 11147, ba: 11147, bb: 11147, bc: 11147, bd: 11147, be: 11147, bf: 11147, bg: 11147, bh: 11147, bi: 11147, bj: 11147, bk: 11147, bl: 11147, bm: 11147, bn: 11147, bo: 11147, bp: 11147, bq: 11147, br: 11147, bs: 11147, bt: 11147, bu: 11147, bv: 11147, bw: 11147, bx: 11147, by: 11147, bz: 11147, ca: 11147, cb: 11147, cc: 11147, cd: 11147, ce: 11147, cf: 11147, cg: 11147, ch: 11147, ci: 11147, cj: 11147, ck: 11147, cl: 11147, cm: 11147, cn: 11147, co: 11147, cp: 11147, cq: 11147, cr: 11147, cs: 11147, ct: 11147, cu: 11147, cv: 11147, cw: 11147, cx: 11147, cy: 11147, cz: 11147, da: 11147, db: 11147, dc: 11147, dd: 11147, de: 11147, df: 11147, dg: 11147, dh: 11147, di: 11147, dj: 11147, dk: 11147, dl: 11147, dm: 11147, dn: 11147, do: 11147, dp: 11147, dq: 11147, dr: 11147, ds: 11147, dt: 11147, du: 11147, dv: 11147, dw: 11147, dx: 11147, dy: 11147, dz: 11147, ea: 11147, eb: 11147, ec: 11147, ed: 11147, ee: 11147, ef: 11147, eg: 11147, eh: 11147, ei: 11147, ej: 11147, ek: 11147, el: 11147, em: 11147, en: 11147, eo: 11147, ep: 11147, eq: 11147, er: 11147, es: 11147, et: 11147, eu: 11147, ev: 11147, ew: 11147, ex: 11147, ey: 11147, ez: 11147, fa: 11147, fb: 11147, fc: 11147, fd: 11147, fe: 11147, ff: 11147, fg: 11147, fh: 11147, fi: 11147, fj: 11147, fk: 11147, fl: 11147, fm: 11147, fn: 11147, fo: 11147, fp: 11147, fq: 11147, fr: 11147, fs: 11147, ft: 11147, fu: 11147, fv: 11147, fw: 11147, fx: 11147, fy: 11147, fz: 11147, ga: 11147, gb: 11147, gc: 11147, gd: 11147, ge: 11147, gf: 11147, gg: 11147, gh: 11147, gi: 11147, gj: 11147, gk: 11147, gl: 11147, gm: 11147, gn: 11147, go: 11147, gp: 11147, gq: 11147, gr: 11147, gs: 11147, gt: 11147, gu: 11147, gv: 11147, gw: 11147, gx: 11147, gy: 11147, gz: 11147, ha: 11147, hb: 11147, hc: 11147, hd: 11147, he: 11147, hf: 11147, hg: 11147, hh: 11147, hi: 11147, hj: 11147, hk: 11147, hl: 11147, hm: 11147, hn: 11147, ho: 11147, hp: 11147, hq: 11147, hr: 11147, hs: 11147, ht: 11147, hu: 11147, hv: 11147, hw: 11147, hx: 11147, hy: 11147, hz: 11147, ia: 11147, ib: 11147, ic: 11147, id: 11147, ie: 11147, if: 11147, ig: 11147, ih: 11147, ii: 11147, ij: 11147, ik: 11147, il: 11147, im: 11147, in: 11147, io: 11147, ip: 11147, iq: 11147, ir: 11147, is: 11147, it: 11147, iu: 11147, iv: 11147, iw: 11147, ix: 11147, iy: 11147, iz: 11147, ja: 11147, jb: 11147, jc: 11147, jd: 11147, je: 11147, jf: 11147, jg: 11147, jh: 11147, ji: 11147, jj: 11147, jk: 11147, jl: 11147, jm: 11147, jn: 11147, jo: 11147, jp: 11147, jq: 11147, jr: 11147, js: 11147, jt: 11147, ju: 11147, jv: 11147, jw: 11147, jx: 11147, jy: 11147, jz: 11147, ka: 11147, kb: 11147, kc: 11147, kd: 11147, ke: 11147, kf: 11147, kg: 11147, kh: 11147, ki: 11147, kj: 11147, kk: 11147, kl: 11147, km: 11147, kn: 11147, ko: 11147, kp: 11147, kq: 11147, kr: 11147, ks: 11147, kt: 11147, ku: 11147, kv: 11147, kw: 11147, kx: 11147, ky: 11147, kz: 11147, la: 11147, lb: 11147, lc: 11147, ld: 11147, le: 11147, lf: 11147, lg: 11147, lh: 11147, li: 11147, lj: 11147, lk: 11147, ll: 11147, lm: 11147, ln: 11147, lo: 11147, lp: 11147, lq: 11147, lr: 11147, ls: 11147, lt: 11147, lu: 11147, lv: 11147, lw: 11147, lx: 11147, ly: 11147, lz: 11147, ma: 11147, mb: 11147, mc: 11147, md: 11147, me: 11147, mf: 11147, mg: 11147, mh: 11147, mi: 11147, mj: 11147, mk: 11147, ml: 11147, mm: 11147, mn: 11147, mo: 11147, mp: 11147, mq: 11147, mr: 11147, ms: 11147, mt: 11147, mu: 11147, mv: 11147, mw: 11147, mx: 11147, my: 11147, mz: 11147, na: 11147, nb: 11147, nc: 11147, nd: 11147, ne: 11147, nf: 11147, ng: 11147, nh: 11147, ni: 11147, nj: 11147, nk: 11147, nl: 11147, nm: 11147, nn: 11147, no: 11147, np: 11147, nq: 11147, nr: 11147, ns: 11147, nt: 11147, nu: 11147, nv: 11147, nw: 11147, nx: 11147, ny: 11147, nz: 11147, oa: 11147, ob: 11147, oc: 11147, od: 11147, oe: 11147, of: 11147, og: 11147, oh: 11147, oi: 11147, oj: 11147, ok: 11147, ol: 11147, om: 11147, on: 11147, oo: 11147, op: 11147, oq: 11147, or: 11147, os: 11147, ot: 11147, ou: 11147, ov: 11147, ow: 11147, ox: 11147, oy: 11147, oz: 11147, pa: 11147, pb: 11147, pc: 11147, pd: 11147, pe: 11147, pf: 11147, pg: 11147, ph: 11147, pi: 11147, pj: 11147, pk: 11147, pl: 11147, pm: 11147, pn: 11147, po: 11147, pp: 11147, pq: 11147, pr: 11147, ps: 11147, pt: 11147, pu: 11147, pv: 11147, pw: 11147, px: 11147, py: 11147, pz: 11147, qa: 11147, qb: 11147, qc: 11147, qd: 11147, qe: 11147, qf: 11147, qg: 11147, qh: 11147, qi: 11147, qj: 11147, qk: 11147, ql: 11147, qm: 11147, qn: 11147, qo: 11147, qp: 11147, qq: 11147, qr: 11147, qs: 11147, qt: 11147, qu: 11147, qv: 11147, qw: 11147, qx: 11147, qy: 11147, qz: 11147, ra: 11147, rb: 11147, rc: 11147, rd: 11147, re: 11147, rf: 11147, rg: 11147, rh: 11147, ri: 11147, rj: 11147, rk: 11147, rl: 11147, rm: 11147, rn: 11147, ro: 11147, rp: 11147, rq: 11147, rr: 11147, rs: 11147, rt: 11147, ru: 11147, rv: 11147, rw: 11147, rx: 11147, ry: 11147, rz: 11147, sa: 11147, sb: 11147, sc: 11147, sd: 11147, se: 11147, sf: 11147, sg: 11147, sh: 11147, si: 11147, sj: 11147, sk: 11147, sl: 11147, sm: 11147, sn: 11147, so: 11147, sp: 11147, sq: 11147, sr: 11147, ss: 11147, st: 11147, su: 11147, sv: 11147, sw: 11147, sx: 11147, sy: 11147, sz: 11147, ta: 11147, tb: 11147, tc: 11147, td: 11147, te: 11147, tf: 11147, tg: 11147, th: 11147, ti: 11147, tj: 11147, tk: 11147, tl: 11147, tm: 11147, tn: 11147, to: 11147, tp: 11147, tq: 11147, tr: 11147, ts: 11147, tt: 11147, tu: 11147, tv: 11147, tw: 11147, tx: 11147, ty: 11147, tz: 11147, ua: 11147, ub: 11147, uc: 1

IDF



TF-IDF



Knn

