

16. Dimensionality reduction

- Recommender systems
- Representations
- PCA
- ICA

Recommender systems

Recommender Systems

The image is a screenshot of the Netflix homepage. At the top left is the Netflix logo. To its right are buttons for "Browse" and "MY LIST". A search bar with a magnifying glass icon is on the far right. Below the header, the page is divided into three main sections:

- Top Picks for Michael**: This section displays five show posters: "COMMUNITY", "NARCOS", "HAPPINESS", "That '70s Show", and "HORNS".
- Trending Now**: This section displays five show posters: "IT'S ALWAYS SUNNY IN PHILADELPHIA", "ORANGE IS THE NEW BLACK", "SUITS", "TRAILER PARK BOYS", and "FRIDAY NIGHT LIGHTS".
- Watch It Again**: This section displays five show posters: "ARMISTICE DEVELOPMENT", "BROOKLYN NINE-NINE", "the office", "HOUSE of CARDS", and "THE THICK OF IT".

Is this product any good?

- Up until now, we recommend a product based on its features
- e.g. Wizard of Oz any good?
 - cute dog +1
 - lots of magic +1
 - pretty shoes +1
 - cheesy messaging -1
- Predict $y_{\text{like}}(x_{\text{WizOfOz}}) = \text{sign}(\theta^T x_{\text{WizOfOz}})$
- Issue: How to collect so much data?



Recommendation without features

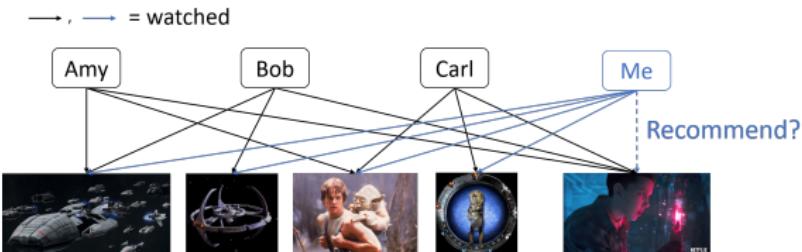
- I love SciFi.
- I've seen *Battlestar Galactica*, *Star Trek*, *Star Wars*, *Stargate*, ...



- *Stranger things?*
- It's not “typical SciFi” (no spaceships, aliens).



Recommendation by collaboration



Assumptions:

- There aren't that many types of people ("archetypes", "archetypical")
- There aren't that many types of movies (genres, styles, ...)
- "Does Bob like Star Wars?" hard to learn (I'm not a mind reader!)
- "Does Bob-like people like Star Wars-like movies?" (Tons of data!)

The Netflix Challenge

Netflix provided a *training* data set of 100,480,507 ratings that 480,189 users gave to 17,770 movies. Each training rating is a quadruplet of the form `<user, movie, date of grade, grade>`. The user and movie fields are integer IDs, while grades are from 1 to 5 (integral) stars.^[3]

The *qualifying* data set contains over 2,817,131 triplets of the form `<user, movie, date of grade>`, with grades known only to the jury. A participating team's algorithm must predict grades on the entire qualifying set, but they are only informed of the score for half of the data, the *quiz* set of 1,408,342 ratings. The other half is the *test* set of 1,408,789, and performance on this is used by the jury to determine potential prize winners. Only the judges know which ratings are in the quiz set, and which are in the test set —this arrangement is intended to make it difficult to *hill climb* on the test set. Submitted predictions are scored against the true grades in terms of *root mean squared error* (RMSE), and the goal is to reduce this error as much as possible. Note that while the actual grades are integers in the range 1 to 5, submitted predictions need not be. Netflix also identified a *probe* subset of 1,408,395 ratings within the *training* data set. The *probe*, *quiz*, and *test* data sets were chosen to have similar statistical properties.

(wikipedia)



The Netflix Challenge

In summary, the data used in the Netflix Prize looks as follows:

- Training set (99,072,112 ratings) not including the probe set, 100,480,607 including the probe set
 - Probe set (1,408,395 ratings)
- Qualifying set (2,817,131 ratings) consisting of:
 - Test set (1,408,789 ratings), used to determine winners
 - Quiz set (1,408,342 ratings), used to calculate leaderboard scores

For each movie, title and year of release are provided in a separate dataset. No information at all is provided about users. In order to protect the privacy of customers, "some of the rating data for some customers in the training and qualifying sets have been deliberately perturbed in one or more of the following ways: deleting ratings; inserting alternative ratings and dates; and modifying rating dates".^[2]

The training set is such that the average user rated over 200 movies, and the average movie was rated by over 5000 users. But there is wide **variance** in the data—some movies in the training set have as few as 3 ratings,^[4] while one user rated over 17,000 movies.^[5]

$$R = \begin{bmatrix} 1 & 5 \\ 2 & 3 \end{bmatrix} = U V^T$$

(wikipedia)

| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 | 32 | 33 | 34 | 35 | 36 | 37 | 38 | 39 | 40 | 41 | 42 | 43 | 44 | 45 | 46 | 47 | 48 | 49 | 50 | 51 | 52 | 53 | 54 | 55 | 56 | 57 | 58 | 59 | 60 | 61 | 62 | 63 | 64 | 65 | 66 | 67 | 68 | 69 | 70 | 71 | 72 | 73 | 74 | 75 | 76 | 77 | 78 | 79 | 80 | 81 | 82 | 83 | 84 | 85 | 86 | 87 | 88 | 89 | 90 | 91 | 92 | 93 | 94 | 95 | 96 | 97 | 98 | 99 | 100 | 101 | 102 | 103 | 104 | 105 | 106 | 107 | 108 | 109 | 110 | 111 | 112 | 113 | 114 | 115 | 116 | 117 | 118 | 119 | 120 | 121 | 122 | 123 | 124 | 125 | 126 | 127 | 128 | 129 | 130 | 131 | 132 | 133 | 134 | 135 | 136 | 137 | 138 | 139 | 140 | 141 | 142 | 143 | 144 | 145 | 146 | 147 | 148 | 149 | 150 | 151 | 152 | 153 | 154 | 155 | 156 | 157 | 158 | 159 | 160 | 161 | 162 | 163 | 164 | 165 | 166 | 167 | 168 | 169 | 170 | 171 | 172 | 173 | 174 | 175 | 176 | 177 | 178 | 179 | 180 | 181 | 182 | 183 | 184 | 185 | 186 | 187 | 188 | 189 | 190 | 191 | 192 | 193 | 194 | 195 | 196 | 197 | 198 | 199 | 200 | 201 | 202 | 203 | 204 | 205 | 206 | 207 | 208 | 209 | 210 | 211 | 212 | 213 | 214 | 215 | 216 | 217 | 218 | 219 | 220 | 221 | 222 | 223 | 224 | 225 | 226 | 227 | 228 | 229 | 230 | 231 | 232 | 233 | 234 | 235 | 236 | 237 | 238 | 239 | 240 | 241 | 242 | 243 | 244 | 245 | 246 | 247 | 248 | 249 | 250 | 251 | 252 | 253 | 254 | 255 | 256 | 257 | 258 | 259 | 260 | 261 | 262 | 263 | 264 | 265 | 266 | 267 | 268 | 269 | 270 | 271 | 272 | 273 | 274 | 275 | 276 | 277 | 278 | 279 | 280 | 281 | 282 | 283 | 284 | 285 | 286 | 287 | 288 | 289 | 290 | 291 | 292 | 293 | 294 | 295 | 296 | 297 | 298 | 299 | 300 | 301 | 302 | 303 | 304 | 305 | 306 | 307 | 308 | 309 | 310 | 311 | 312 | 313 | 314 | 315 | 316 | 317 | 318 | 319 | 320 | 321 | 322 | 323 | 324 | 325 | 326 | 327 | 328 | 329 | 330 | 331 | 332 | 333 | 334 | 335 | 336 | 337 | 338 | 339 | 340 | 341 | 342 | 343 | 344 | 345 | 346 | 347 | 348 | 349 | 350 | 351 | 352 | 353 | 354 | 355 | 356 | 357 | 358 | 359 | 360 | 361 | 362 | 363 | 364 | 365 | 366 | 367 | 368 | 369 | 370 | 371 | 372 | 373 | 374 | 375 | 376 | 377 | 378 | 379 | 380 | 381 | 382 | 383 | 384 | 385 | 386 | 387 | 388 | 389 | 390 | 391 | 392 | 393 | 394 | 395 | 396 | 397 | 398 | 399 | 400 | 401 | 402 | 403 | 404 | 405 | 406 | 407 | 408 | 409 | 410 | 411 | 412 | 413 | 414 | 415 | 416 | 417 | 418 | 419 | 420 | 421 | 422 | 423 | 424 | 425 | 426 | 427 | 428 | 429 | 430 | 431 | 432 | 433 | 434 | 435 | 436 | 437 | 438 | 439 | 440 | 441 | 442 | 443 | 444 | 445 | 446 | 447 | 448 | 449 | 450 | 451 | 452 | 453 | 454 | 455 | 456 | 457 | 458 | 459 | 460 | 461 | 462 | 463 | 464 | 465 | 466 | 467 | 468 | 469 | 470 | 471 | 472 | 473 | 474 | 475 | 476 | 477 | 478 | 479 | 480 | 481 | 482 | 483 | 484 | 485 | 486 | 487 | 488 | 489 | 490 | 491 | 492 | 493 | 494 | 495 | 496 | 497 | 498 | 499 | 500 | 501 | 502 | 503 | 504 | 505 | 506 | 507 | 508 | 509 | 510 | 511 | 512 | 513 | 514 | 515 | 516 | 517 | 518 | 519 | 520 | 521 | 522 | 523 | 524 | 525 | 526 | 527 | 528 | 529 | 530 | 531 | 532 | 533 | 534 | 535 | 536 | 537 | 538 | 539 | 540 | 541 | 542 | 543 | 544 | 545 | 546 | 547 | 548 | 549 | 550 | 551 | 552 | 553 | 554 | 555 | 556 | 557 | 558 | 559 | 560 | 561 | 562 | 563 | 564 | 565 | 566 | 567 | 568 | 569 | 570 | 571 | 572 | 573 | 574 | 575 | 576 | 577 | 578 | 579 | 580 | 581 | 582 | 583 | 584 | 585 | 586 | 587 | 588 | 589 | 590 | 591 | 592 | 593 | 594 | 595 | 596 | 597 | 598 | 599 | 600 | 601 | 602 | 603 | 604 | 605 | 606 | 607 | 608 | 609 | 610 | 611 | 612 | 613 | 614 | 615 | 616 | 617 | 618 | 619 | 620 | 621 | 622 | 623 | 624 | 625 | 626 | 627 | 628 | 629 | 630 | 631 | 632 | 633 | 634 | 635 | 636 | 637 | 638 | 639 | 640 | 641 | 642 | 643 | 644 | 645 | 646 | 647 | 648 | 649 | 650 | 651 | 652 | 653 | 654 | 655 | 656 | 657 | 658 | 659 | 660 | 661 | 662 | 663 | 664 | 665 | 666 | 667 | 668 | 669 | 670 | 671 | 672 | 673 | 674 | 675 | 676 | 677 | 678 | 679 | 680 | 681 | 682 | 683 | 684 | 685 | 686 | 687 | 688 | 689 | 690 | 691 | 692 | 693 | 694 | 695 | 696 | 697 | 698 | 699 | 700 | 701 | 702 | 703 | 704 | 705 | 706 | 707 | 708 | 709 | 710 | 711 | 712 | 713 | 714 | 715 | 716 | 717 | 718 | 719 | 720 | 721 | 722 | 723 | 724 | 725 | 726 | 727 | 728 | 729 | 730 | 731 | 732 | 733 | 734 | 735 | 736 | 737 | 738 | 739 | 740 | 741 | 742 | 743 | 744 | 745 | 746 | 747 | 748 | 749 | 750 | 751 | 752 | 753 | 754 | 755 | 756 | 757 | 758 | 759 | 760 | 761 | 762 | 763 | 764 | 765 | 766 | 767 | 768 | 769 | 770 | 771 | 772 | 773 | 774 | 775 | 776 | 777 | 778 | 779 | 780 | 781 | 782 | 783 | 784 | 785 | 786 | 787 | 788 | 789 | 790 | 791 | 792 | 793 | 794 | 795 | 796 | 797 | 798 | 799 | 800 | 801 | 802 | 803 | 804 | 805 | 806 | 807 | 808 | 809 | 810 | 811 | 812 | 813 | 814 | 815 | 816 | 817 | 818 | 819 | 820 | 821 | 822 | 823 | 824 | 825 | 826 | 827 | 828 | 829 | 830 | 831 | 832 | 833 | 834 | 835 | 836 | 837 | 838 | 839 | 840 | 841 | 842 | 843 | 844 | 845 | 846 | 847 | 848 | 849 | 850 | 851 | 852 | 853 | 854 | 855 | 856 | 857 | 858 | 859 | 860 | 861 | 862 | 863 | 864 | 865 | 866 | 867 | 868 | 869 | 870 | 871 | 872 | 873 | 874 | 875 | 876 | 877 | 878 | 879 | 880 | 881 | 882 | 883 | 884 | 885 | 886 | 887 | 888 | 889 | 890 | 891 | 892 | 893 | 894 | 895 | 896 | 897 | 898 | 899 | 900 | 901 | 902 | 903 | 904 | 905 | 906 | 907 | 908 | 909 | 910 | 911 | 912 | 913 | 914 | 915 | 916 | 917 | 918 | 919 | 920 | 921 | 922 | 923 | 924 | 925 | 926 | 927 | 928 | 929 | 930 | 931 | 932 | 933 | 934 | 935 | 936 | 937 | 938 | 939 | 940 | 941 | 942 | 943 | 944 | 945 | 946 | 947 | 948 | 949 | 950 | 951 | 952 | 953 | 954 | 955 | 956 | 957 | 958 | 959 | 960 | 961 | 962 | 963 | 964 | 965 | 966 | 967 | 968 | 969 | 970 | 971 | 972 | 973 | 974 | 975 | 976 | 977 | 978 | 979 | 980 | 981 | 982 | 983 | 984 | 985 | 986 | 987 | 988 | 989 | 990 | 991 | 992 | 993 | 994 | 995 | 996 | 997 | 998 | 999 | 1000 | 1001 | 1002 | 1003 | 1004 | 1005 | 1006 | 1007 | 1008 | 1009 | 10010 | 10011 | 10012 | 10013 | 10014 | 10015 | 10016 | 10017 | 10018 | 10019 | 10020 | 10021 | 10022 | 10023 | 10024 | 10025 | 10026 | 10027 | 10028 | 10029 | 10030 | 10031 | 10032 | 10033 | 10034 | 10035 | 10036 | 10037 | 10038 | 10039 | 10040 | 10041 | 10042 | 10043 | 10044 | 10045 | 10046 | 10047 | 10048 | 10049 | 10050 | 10051 | 10052 | 10053 | 10054 | 10055 | 10056 | 10057 | 10058 | 10059 | 10060 | 10061 | 10062 | 10063 | 10064 | 10065 | 10066 | 10067 | 10068 | 10069 | 10070 | 10071 | 10072 | 10073 | 10074 | 10075 | 10076 | 10077 | 10078 | 10079 | 10080 | 10081 | 10082 | 10083 | 10084 | 10085 | 10086 | 10087 | 10088 | 10089 | 10090 | 10091 | 10092 | 10093 | 10094 | 10095 | 10096 | 10097 | 10098 | 10099 | 100100 | 100101 | 100102 | 100103 | 100104 | 100105 | 100106 | 100107 | 100108 | 100109 | 100110 | 100111 | 100112 | 100113 | 100114 | 100115 | 100116 | 100117 | 100118 | 100119 | 100120 | 100121 | 100122 | 100123 | 100124 | 100125 | 100126 | 100127 | 100128 | 100129 | 100130 | 100131 | 100132 | 100133 | 100134 | 100135 | 100136 | 100137 | 100138 | 100139 | 100140 | 100141 | 100142 | 100143 | 100144 | 100145 | 100146 | 100147 | 100148 | 100149 | 100150 | 100151 | 100152 | 100153 | 100154 | 100155 | 100156 | 100157 | 100158 | 100159 | 100160 | 100161 | 100162 | 100163 | 100164 | 100165 | 100166 | 100167 | 100168 | 100169 | 100170 | 100171 | 100172 | 100173 | 100174 | 100175 | 100176 | 100177 | 100178 | 100179 | 100180 | 100181 | 100182 | 100183 | 100184 | 100185 | 100186 | 100187 | 100188 | 100189 | 100190 | 100191 | 100192 | 100193 | 100194 | 100195 | 100196 | 100197 | 100198 | 100199 | 100200 | 100201 | 100202 | 100203 | 100204 | 100205 | 100206 | 100207 | 100208 | 100209 | 100210 | 100211 | 100212 | 100213 | 100214 | 100215 | 100216 | 100217 | 100218 | 100219 | 100220 | 100221 | 100222 | 100223 | 100224 | 100225 | 100226 | 100227 | 100228 | 100229 | 100230 | 100231 | 100232 | 100233 | 100234 | 100235 | 100236 | 100237 | 100238 | 100239 | 100240 | 100241 | 100242 | 100243 | 100244 | 100245 | 100246 | 100247 | 100248 | 100249 | 100250 | 100251 | 100252 | 100253 | 100254 | 100255 | 100256 | 100257 | 100258 | 100259 | 100260 | 100261 | 100262 | 100263 | 100264 | 100265 | 100266 | 100267 | 100268 | 100269 | 100270 | 100271 | 100272 | 100273 | 100274 | 100275 | 100276 | 100277 | 100278 | 100279 | 100280 | 100281 | 100282 | 100283 | 100284 | 100285 | 100286 | 100287 | 100288 | 100289 | 100290 | 100291 | 100292 | 100293 | 100294 | 100295 | 100296 | 100297 | 100298 | 100299 | 100300 | 100301 | 100302 | 100303 | 100304 | 100305 | 100306 | 100307 | 100308 | 100309 | 100310 | 100311 | 100312 | 100313 | 100314 | 100315 | 100316 | 100317 | 100318 | 100319 | 100320 | 100321 | 100322 | 100323 | 100324 | 100325 | 100326 | 100327 | 100328 | 100329 | 100330 | 100331 | 100332 | 100333 | 100334 | 100335 | 100336 | 100337 | 100338 | 100339 | 100340 | 100341 | 100342 | 100343 | 100344 | 100345 | 100346 | 100347 | 100348 | 100349 | 100350 | |

The Netflix winner

The BellKor 2008 Solution to the Netflix Prize

Robert M. Bell
~~AT&T Labs~~ - Research
Florham Park, NJ

Yehuda Koren
~~Yahoo!~~ Research
Haifa, Israel

Chris Volinsky
AT&T Labs - Research
Florham Park, NJ

BellKor@research.att.com

1. Introduction

Our RMSE=0.8643² solution is a linear blend of over 100 results. Some of them are new to this year, whereas many others belong to the set that was reported a year ago in our 2007 Progress Prize report [3]. This report is structured accordingly. In Section 2 we detail methods new to this year. In general, our view is that those newer methods deliver a superior performance compared to the methods we used a year ago. Throughout the description of the methods, we highlight the specific predictors that participated in the final blended solution. Nonetheless, the older methods still play a role in the blend, and thus in Section 3 we list those methods repeated from a year ago. Finally, we conclude with general thoughts in Section 4.

Solution spurred study in collaborative filtering

Matrix factorization

- m users, user i represent as $\underline{u_i} \in \mathbb{R}^d$ *ratings* $d \ll \{m, n\}$
- n movies, movie j represent as $\underline{v_j} \in \mathbb{R}^d$
- Predicted rating of user i on movie j : $R_{ij} \approx g(u_i^\top v_j)$
- Form matrices

$$U = [u_1, u_2, \dots, u_m], \quad V = [v_1, v_2, \dots, v_n], \quad R \in \mathbb{R}^{m \times n}$$

- Assumption: $R_{ij} = f(u_i, v_j)$
- Matrix factorization: $R_{ij} = u_i^\top v_j$
 $R = UV^T$
 U looks like $u_i^\top v_j$
 V looks like $u_i^\top v_j$
- If $d \ll \min\{m, n\}$, then at most d directions of freedom
- (meta) \rightarrow constrains “how many types” of people, movies

$$d=1$$

Basic implementation

$$u_1 = 1$$

$$v_2 = 1$$

- Collect data R_{ij} from observed user i , movie j
- Sparse mask: $S_{ij} = 1$ if user i saw movie j , $S_{ij} = 0$ otherwise
- Find u_i, v_j to minimize

$$R = \begin{bmatrix} \cdot & \cdot \\ \cdot & \cdot \end{bmatrix}$$

$$\underset{u_i \in \mathbb{R}^d, v_j \in \mathbb{R}^d}{\text{minimize}} \sum_{i=1}^m \sum_{j=1}^n S_{ij} (R_{ij} - u_i^T v_j)^2$$

what is u_2, v

- Predict unseen $R_{ij} = u_i^T v_j$
- Idea: by enforcing $d \ll \min\{m, n\}$, we “learn” the unseen R_{ij} via “archetype” information sharing

$$u_2 = 1 \quad v_2 = 1$$

Matrix factorization is a nonconvex problem

$$\underset{u_i \in \mathbb{R}^d, v_j \in \mathbb{R}^d}{\text{minimize}} \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^n S_{ij} (R_{ij} - u_i^T v_j)^2$$

Proof: Take $d = m = n = 1$, $S = 1$. Then objective is

$$f(u, v) = \frac{1}{2}(R - uv)^2, \nabla f(u, v) = (uv - R) \begin{bmatrix} u \\ v \end{bmatrix}, \nabla^2 f(u, v) = \begin{bmatrix} u^2 & 2uv \\ 2uv & v^2 \end{bmatrix}$$

Take $u = v = 1$, then Hessian is $\begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix}$ which is not positive semidefinite.

$$\begin{bmatrix} 1 \\ -1 \end{bmatrix}^T \begin{bmatrix} 1 & 2 \\ 2 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \end{bmatrix} = \begin{bmatrix} -1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \end{bmatrix} = -2$$

Better-than-basic: De-biasing

$$\begin{bmatrix} 5 \\ 5 \\ 5 \end{bmatrix}$$

- Scalar R_0 = average rating for everything
- Some people are just critics: Vector $\underline{u_0} \in \mathbb{R}^m$ = average user rating (after $-R_0$)
- Some movies just suck: Vector $v_0 \in \mathbb{R}^n$ = average movie rating (after $-R_0$)

New objective function

$$\underset{\substack{u_i \in \mathbb{R}^d, v_j \in \mathbb{R}^d \\ i, j : S_{ij} = 1}}{\text{minimize}} \sum \underbrace{((R_{ij} - R_0 - u_0[i] - \cancel{v_0[j]}) - u_i^T v_j)^2}_{\text{de-biased ratings}}$$

pred : $(u_i^T v_j + u_0[i] + v_0[j]) + R_0$

Implementation of better-than-basic method

$$\underset{u_i \in \mathbb{R}^d, v_j \in \mathbb{R}^d}{\text{minimize}} \quad f(U, V) := \sum_{i,j: S_{ij}=0} \underbrace{((R_{ij} - R_0 - u_0[i] - v_j[0]) - u_i^T v_j)^2}_{\text{de-biased ratings}}$$

- Alternating gradient descent: Start with average ratings

$$\begin{aligned} U^{(t+1)} &= U^{(t)} - \alpha \nabla_U f(U^{(t)}, V^{(t)}) \\ V^{(t+1)} &= V^{(t)} - \alpha \nabla_V f(U^{(t+1)}, V^{(t)}) \end{aligned}$$



- Trial and error to pick best α
- Personal experience: about 5-10 passes through data is enough

Further details: [https://datajobs.com/data-science-repo/Recommender-Systems-\[Netflix\].pdf](https://datajobs.com/data-science-repo/Recommender-Systems-[Netflix].pdf)

RMSE scores on MovieLens 100K

Baseline Methods

| Software | Method | 5-fold CV all-but-10 References | |
|------------------|------------------|---------------------------------|--------|
| MyMediaLite 3.07 | GlobalAverage | 1.1256 | 1.1238 |
| MyMediaLite 3.07 | UserAverage | 1.0437 | 1.0518 |
| MyMediaLite 3.07 | ItemAverage | 1.0246 | 1.0453 |
| MyMediaLite 3.07 | UserItemBaseline | 0.9413 | 0.9656 |

kNN-based Collaborative Filtering

| Software | Method | 5-fold CV all-but-10 | References |
|------------------|---------|----------------------|------------|
| MyMediaLite 3.07 | UserKNN | 0.9283 | 0.9572 |
| MyMediaLite 3.07 | ItemKNN | 0.9182 | 0.9445 |

Matrix Factorization

| Software | Method | 5-fold CV all-but-10 References | |
|------------------|----------------------------------|---------------------------------|--------|
| MyMediaLite 3.07 | BiasedMatrixFactorization | 0.9220 | 0.9475 |
| MyMediaLite 3.07 | SVDPlusPlus | 0.9112 | 0.9409 |
| MyMediaLite 3.07 | SigmoidUserAsymmetricFactorModel | 0.8939 | 0.9232 |

RMSE scores on MovieLens 100K

Browse > Miscellaneous > Recommendation Systems > MovieLens 100K dataset

Recommendation Systems on MovieLens 100K



<https://paperswithcode.com/sota/collaborative-filtering-on-movielens-100k>

RMSE scores on MovieLens 100K

| RANK | MODEL | RMSE (U1 SPLITS) | RMSE (RANDOM 90/10 SPLITS) | RMSE@80%TRAIN | EXTRA TRAINING DATA | PAPER | CODE | RESULT | YEAR |
|------|--|------------------------|-------------------------------------|---------------|---------------------------|---|-------------------|-------------------|------|
| 1 | Bayesian timeSVD++ flipped + Feat w/ Ordered Probit Regression | 0.882 | | | ✓ | On the Difficulty of Evaluating Baselines: A Study on Recommender Systems | 🔗 | 📄 | 2019 |
| 2 | Bayesian timeSVD++ flipped + Feat | 0.884 | | | ✓ | On the Difficulty of Evaluating Baselines: A Study on Recommender Systems | 🔗 | 📄 | 2019 |
| 3 | Bayesian timeSVD++ flipped | 0.886 | | | ✓ | On the Difficulty of Evaluating Baselines: A Study on Recommender Systems | 🔗 | 📄 | 2019 |
| 4 | GraphRec + Feat | 0.897 | 0.883 | | ✓ | Attribute-aware non-linear co-embeddings of graph features | 🔗 | 📄 | 2019 |
| 5 | GraphRec | 0.904 | 0.887 | | ✗ | Attribute-aware non-linear co-embeddings of graph features | 🔗 | 📄 | 2019 |
| 6 | IGMC | 0.905 | | | ✗ | Inductive Matrix Completion Based on Graph Neural Networks | 🔗 | 📄 | 2019 |
| 7 | GC-MC + feat | 0.905 | | | ✓ | Graph Convolutional Matrix Completion | 🔗 | 📄 | 2017 |
| 8 | GC-MC | 0.910 | | | ✗ | Graph Convolutional Matrix Completion | 🔗 | 📄 | 2017 |
| 9 | Self-Supervised Exchangeable Model | 0.91 | | | ✗ | Deep Models of Interactions Across Sets | 🔗 | 📄 | 2018 |
| 10 | GRAEM | 0.9174 | | | ✓ | Scalable Probabilistic Matrix Factorization with Graph-Based Priors | 🔗 | 📄 | 2019 |

<https://paperswithcode.com/sota/collaborative-filtering-on-movielens-100k>

Ongoing research!

- If $R_{ij} \in \{\text{blue thumbs up}, \text{blue thumbs down}\}$, replace mean squared with logistic loss
- If $R_{ij} = 1$ if watched, 0 if not watched, need something more creative
- Side information: we shouldn't just ignore features, if we do have them!
- Sharing accounts, how to disambiguate?
- Incorporate movie reviews? NLP?



Representations

| Sentiment | Tweets |
|-----------------|---|
| Negative | @united is the worst. Nonrefundable First class tickets? Oh because when you select Global/FC their system auto selects economy w/upgrade. |
| Neutral | @united I will not be flying you again @VirginAmerica my drivers license is expired by a little over a month. Can I fly Friday morning using my expired license? @VirginAmerica any plans to start flying direct from DAL to LAS? |
| Positive | @VirginAmerica done! Thank you for the quick response, apparently faster than sitting on hold ;) @united I appreciate your efforts getting me home! |

ipullrank.com/step-step-twitter-sentiment-analysis-visualizing-united-airlines-pr-crisis/

Feature extraction?

- Count # good / bad words (works well, easy to trick)
- Count # good / bad known phrases (too many combinations!)
- How to turn a non-numeric feature into a number-like feature?

$$u_i^T u_j \gg 0$$

$$u_{\text{son}}^T u_{\text{brown}} \gg 0$$

Co-occurrence = similarity

- ID words
- the quick brown fox jumped over the brown fence
1 T 2 5, 4 5, 6 T 3 7

1 = the, 2 = quick, 3 = brown, 4 = fox, 5 = jumped, 6 = over, 7 = fence

- Frequency vector:

$$f = (2, 1, 2, 1, 1, 1, 1)$$

- Co-occurrence matrix, sliding window = 3

$$C = \begin{bmatrix} 2 & 1 & 2 & 0 & 0 & 1 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 & 0 \\ 2 & 1 & 2 & 1 & 0 & 0 & 1 \\ 0 & 0 & 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 & 0 \\ 1 & 0 & 0 & 0 & 1 & 1 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 1 \end{bmatrix}$$

- Word embedding: $C_{ij} \approx u_i^T u_j$

Word embeddings

Example: vector = histogram of cooccurring words (Lund, Burgess '96)

I have a big fluffy gray
cat named Fluffy. We saw
a red balloon fly by.

$$x_{\text{cat}} = [0, 0, 2, 1]$$

apple balloon fluffy gray

$$x_{\text{red}} = [0, 1, 0, 0]$$

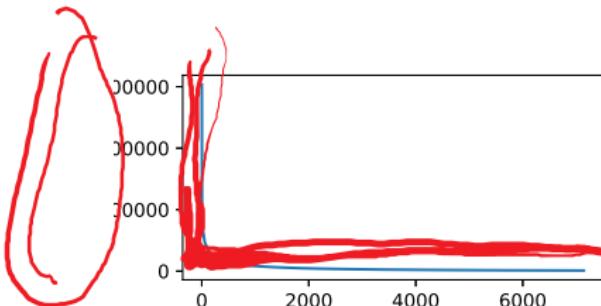
apple balloon fluffy gray

2-D projection of output →



Figure 2. Multidimensional scaling of co-occurrence vectors.

Linguistics: Very few common words, a lot of rare words



Disclaimer: This histogram is after pre-processing, which removes common and rare words—real histogram is more pronounced

Pairwise mutual information (PMI) matrix

$$\text{PMI}_{i,j} = \log \left(\frac{\text{corpus length} \times \# \text{ co-occurrence of word } i \text{ with word } j}{\text{freq. of word } i \times \text{freq. of word } j} \right)$$

Positive pairwise mutual information (PPMI) matrix

$$\text{PPMI}_{i,j} = \max \{ \text{PMI}_{i,j}, 0 \}$$

Word embeddings

Other notable methods

- Deerwester, et al '90 (Latent semantic analysis)
- Bengio et al 03, Collobert and Weston 08 (Deep neural nets)
- Pennington, Socher, Manning '14 (GloVe), Mikolov, Chen, Corrado, Dean '13 (word2vec)

~~Levy and Goldberg '14~~: many modern methods can be interpreted as implicit matrix factorization of a specific matrix: $C = UU^T$

- $C \in \mathbb{R}^{n \times n}$ encodes similarity between all pairs of n words (n is large)
- $U \in \mathbb{R}^{n \times d}$, each row contains d -dim. word vector representation ($d \ll n$)

Principle component analysis

Principle component analysis

Given $X \in \mathbb{R}^{m \times n}$, find $U \in \mathbb{R}^{m \times r}$, $V \in \mathbb{R}^{n \times r}$ where

$$UV^T = \min_{U \in \mathbb{R}^{m \times r}, V \in \mathbb{R}^{n \times r}} \|UV^T - X\|_F^2$$

$$X \approx \hat{U}_r \hat{\Sigma}_r \hat{V}_r^T$$
$$U = \hat{U}_r \hat{\Sigma}_r$$
$$V = \hat{V}_r \hat{\Sigma}_r$$

Example: compress 2D to 1D

The diagram illustrates the dimensionality reduction process. It starts with a blue 2D ellipse labeled d (representing the original dimension). An arrow points to a red line segment, representing the compressed 1D representation. This visualizes how the data is mapped from a higher-dimensional space (d) to a lower-dimensional space ($r < d$).

$$X = \sum_{i=1}^d \sigma_i u_i v_i^T \rightarrow \sum_{i=1}^r \sigma_i u_i v_i^T$$

$r < d = m \times n$ ($r \ll n$)

Eckart-Young-Mirsky theorem

- Given some $X \in \mathbb{R}^{m,n}$
- Take the singular value decomposition of X

$$X = \sum_{i=1}^{\min\{m,n\}} \sigma_i u_i v_i^T$$

$WV^T \approx X$


where

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_{\min\{m,n\}}$$

are the singular values and correspond to singular vectors u_i, v_i .

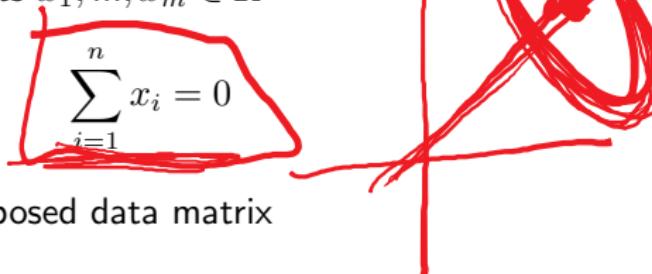
- Then $\hat{X} := \sum_{i=1}^r \sigma_i u_i v_i^T$ is the best rank- r approximation of X ; that is,

$$\|\hat{X} - X\|_* = \min_{X' \text{ rank } r} \|X' - X\|_*^2$$

where $\|\cdot\|_*$ is the Frobenius or spectral norm

Principle component analysis (PCA)

- Given m centered data points $x_1, \dots, x_m \in \mathbb{R}^n$



- Perform a SVD on the composed data matrix

$$X = \begin{bmatrix} x_1^T \\ \vdots \\ x_m^T \end{bmatrix} = \sum_{i=1}^{\min\{m,n\}} \sigma_i u_i v_i^T, \quad \sigma_1 \geq \sigma_2 \geq \dots$$

$$C = \hat{a} \hat{a}^T$$

- Reconstruct with some $r \ll n$

$$\hat{X} = \begin{bmatrix} \hat{x}_1^T \\ \vdots \\ \hat{x}_m^T \end{bmatrix} = \sum_{i=1}^r \sigma_i u_i v_i^T$$

Reconstruction error

EYMTum

\hat{X} = rank r approx of X

$$\sum_{i=1}^n \|x_i - \hat{x}_i^T\|_2^2 = \|\underbrace{X - \hat{X}}_F\|_F^2$$

$$= \left\| \sum_{i=r+1}^{\min\{m,n\}} \sigma_i u_i v_i^T \right\|_F^2$$

$$= \sum_{i=r+1}^{\min\{m,n\}} u_i^T \underbrace{XX^T}_{\text{gram matrix, } m \times m} u_i$$

$$= \sum_{i=r+1}^{\min\{m,n\}} v_i^T \underbrace{X^T X}_{\text{covariance matrix, } n \times n} v_i$$

[Turk, Pentland '91] Input images



#PK
#face

A red arrow points from the text "#PK" to the right edge of the grid. A red curved arrow points from the text "#face" down to the bottom center of the grid.

[Turk, Pentland '91] Principle components



[Turk, Pentland '91] Reconstruction

Each image corresponds to adding 8 principle components



$$u = \hat{u} Q^T$$

Independent component analysis

recommendation $\rightarrow R \approx U V^T$ $V^T V - I$

word embeddings $\rightarrow C \propto U U^T$

solver $\rightarrow X \propto U \tilde{\Sigma} V^T$

Blind source separation (BSS) problem



- d speakers, d sensors
- $x(t) = As(t)$, $A \in \mathbb{R}^{d \times d}$ is mixing matrix
- Given $x(t)$, can we recover $s(t)$?

Aside: chain rule for PDF of linear transformations

- Suppose that $s = Wx$ and W is invertible.
- If $p_X(x)$ is the probability distribution (pdf) of x , what is the pdf of s ?

Aside: chain rule for PDF of linear transformations

- Suppose that $s = Wx$ and W is invertible.
- If $p_X(x)$ is the probability distribution (pdf) of x , what is the pdf of s ?

Ans:

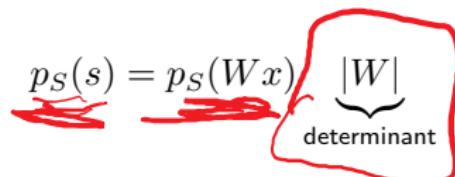
$$p_S(s) = p_S(Wx) \underbrace{|W|}_{\text{determinant}}$$

- Why?

Aside: chain rule for PDF of linear transformations

- Suppose that $s = \underline{Wx}$ and W is invertible.
- If $\underline{p_X(x)}$ is the probability distribution (pdf) of x , what is the pdf of s ?

Ans:

$$p_S(s) = p_S(Wx) \frac{|W|}{\text{determinant}}$$


- Why? Chain rule

1-D example: CDFs

$$F_S(s) = \Pr(S < s) = \Pr(wX < ws) = F_X(ws)$$

Now differentiate

$$p_S(s) = \frac{d}{ds} F_S(s) = \frac{d}{ds} F_X(ws) \stackrel{x=ws}{=} w p_X(x)$$



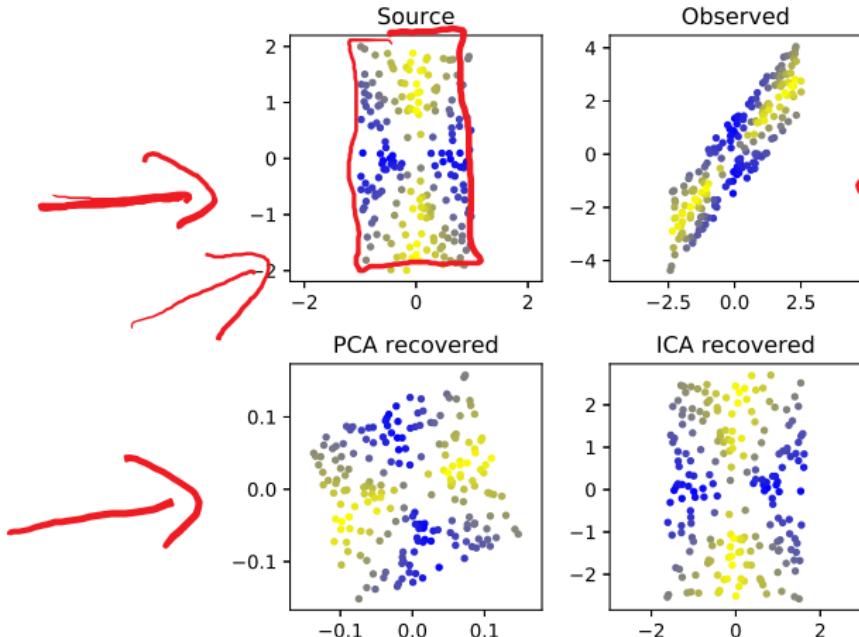
Formal ICA setup

- Given $x_1, \dots, x_T \in \mathbb{R}^d$, find $A \in \mathbb{R}^{d \times d}$ such that
 - $x_i = As_i$ for all i
 - s_1, \dots, s_T are all independent
- The MLE is

$$\begin{aligned}\frac{1}{T} \log(\Pr(W|x)) &= \frac{1}{T} \sum_{i=1}^T \log(\Pr(W|x_i)) \\ &= \frac{1}{m} \sum_{i=1}^T \sum_{j=1}^d \underbrace{\log(p_S(w_j^T x_i))}_{\text{independent sources}} + \underbrace{\log(|W|)}_{\text{chain rule}}\end{aligned}$$

- Procedure: find $\underline{W_{\text{ICA}}}$ the maximizer of this loss function
- Find $(s_{\text{ICA}})_i = W_{\text{ICA}}^{-1}x_i, i = 1, \dots, T$

Example: Uniform 2-D distribution



- PCA recovers linear subspace, but rotationally ambiguous
- ICA can recover subspace and rotation

ICA vs PCA on faces

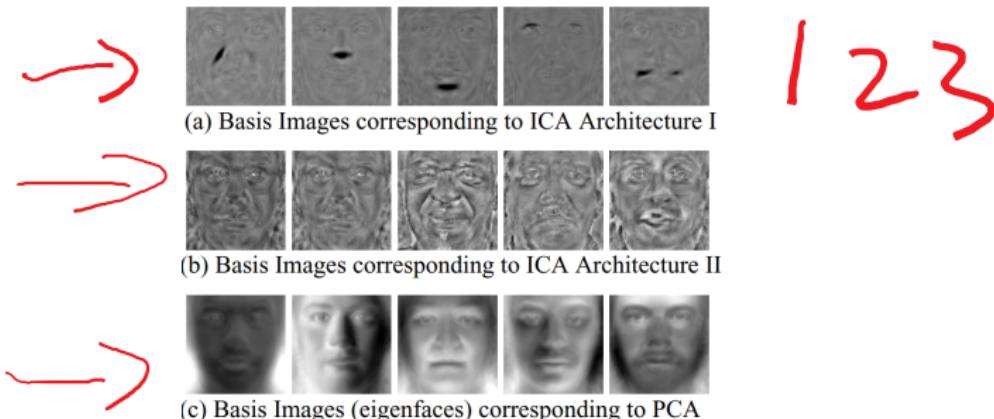


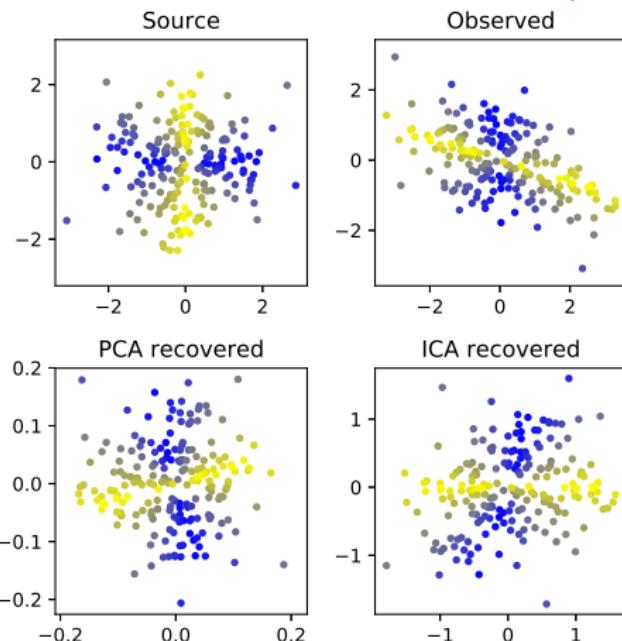
Figure 2. Basis images corresponding to ICA Architecture I, ICA Architecture II and PCA

Top: ICA, Middle: ICA with whiten preprocessing, Bottom: PCA

Yang, Zhang, Yang "Is ICA significantly better than PCA for face recognition?" ICCV
05

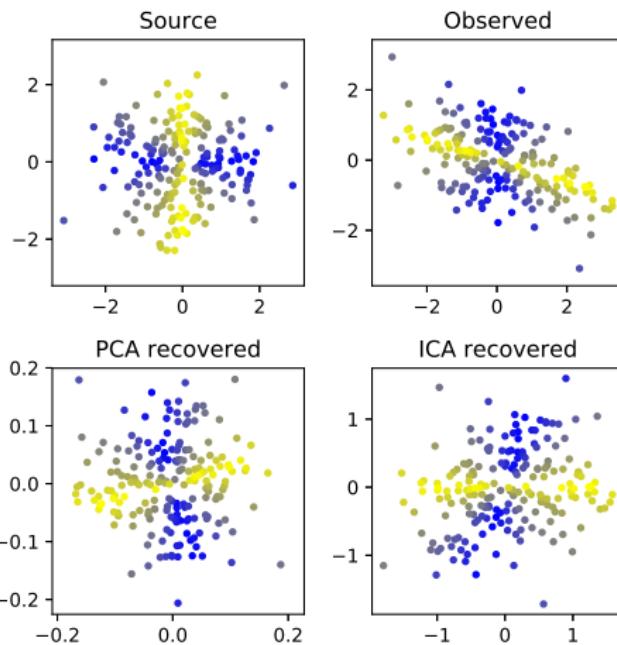
iid

ICA cannot do better than PCA on Gaussian data



Why? (Just answer based on what you see in the picture)

ICA cannot do better than PCA on Gaussian data



Why? (Just answer based on what you see in the picture)

Gaussian i.i.d. = rotationally invariant. (Only distribution with this property.)