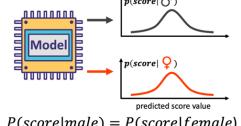


## Fairness definition: score-based statistical parity

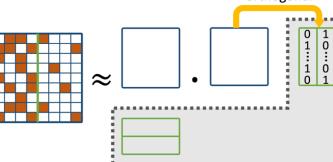
## Fairness-Aware Recommendation (FAR) algorithm

## Similarity-based approach

Statistical parity requires the **prediction** (predicted scores) from the model to be **independent** to the **group labels**.



$$P(score|male) = P(score|female)$$



Mainstream users have high similarity to other users.

$$MS_u^{\text{sim}} = \frac{\sum_{v \in \mathcal{U} \setminus u} J_{u,v} / (N - 1)}{\text{Average similarity to other users}}$$

Jaccard similarity between  $u$  and  $v$  w.r.t. historical interactions  
Number of users

## Global method – Data Augmentation algorithm (DA)

- For a niche user  $u$ , generate synthetic users based on the distribution (following the **Distribution Calibration** method<sup>1)</sup>:

$$\mathbf{p}_u = \alpha \mathbf{o}_u + (1 - \alpha) \frac{1}{|\mathcal{N}_u|} \sum_{v \in \mathcal{N}_u} \mathbf{o}_v$$

Historical interaction record of  $u$   
Orthogonal  
Remove sensitive dimensions  
Neighbor users similar to  $u$ .

Data from users of different preferences may **not be helpful** (even play **negative roles**) for learning preference of a target group of users.

Global methods have to **trade-off** between niche and mainstream users.

	ML1M						
	NDCG @20	low	Subgroups of mainstream levels	med-low	medium	med-high	high
VAE	0.315	0.209	0.264	0.283	0.337	0.483	
DA	0.317	0.222	0.266	0.282	0.335	0.480	
WL	0.319	0.232	0.272	0.284	0.332	0.476	
LOCA	0.323	0.242	0.276	0.286	0.335	0.476	
LFT	<b>0.337</b>	<b>0.255</b>	<b>0.288</b>	<b>0.298</b>	<b>0.349</b>	<b>0.496</b>	
LFT vs. VAE	6.98%	22.01%	9.09%	5.30%	3.54%	2.69%	
LFT vs. LOCA	4.40%	5.55%	4.09%	4.23%	4.24%	4.22%	

Proposed LFT outperforms other alternatives significantly.

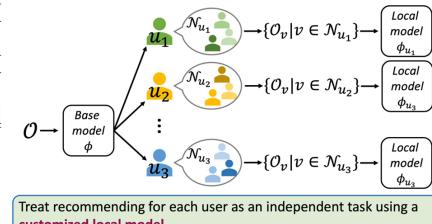
## Global method – Weighted Loss algorithm (WL)

- Add **weights** to users in loss function.
- Users with lower mainstream levels get higher weights:

$$\mathcal{L}_{WL} = \sum_{u \in \mathcal{U}} w_u \cdot \mathcal{L}_{VAE}(u), \quad w_u \propto (\frac{1}{MS_u})^\beta$$

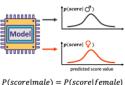
Original loss for user  $u$

## Local method – Local Fine Tuning algorithm



Treat recommending for each user as an independent task using a **customized local model**.

## Takeaways



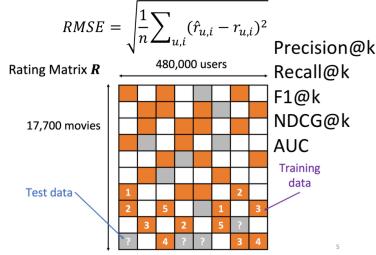
$$P(score|male) = P(score|female)$$

- Define recommendation fairness as fairness. **score-based statistical parity**.

- Propose an **embedding intervention** method to enhance

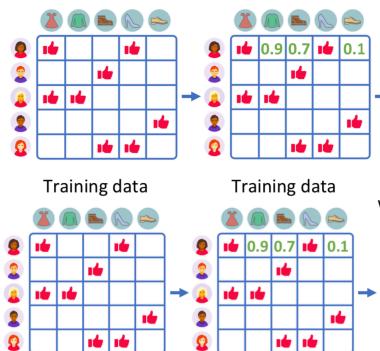
- Effective performance for the challenging **multi-attribute** and **multi-category** scenario.

So far, focus on explicit recommendation



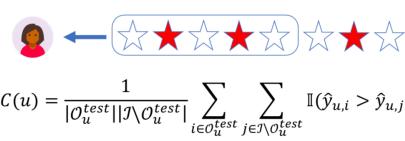
## Implicit Recommendation

A user set  $\mathcal{U} = \{1, 2, \dots, N\}$ ; an item set  $\mathcal{I} = \{1, 2, \dots, M\}$ ; a binary variable  $y_{u,i}$ ; an implicit feedback dataset  $\mathcal{O} = \{(u, i)\}$ ; and the goal is to predict  $\hat{y}_{u,i}$  when  $y_{u,i} = 0$ , and return a ranked list of  $k$  items based on  $\hat{y}_{u,i}$ .



$$\begin{aligned} \text{Precision}@k &= \frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \frac{|\mathcal{O}_u^{\text{test}} \cap \mathcal{O}_u^{\text{rec}}|}{k} \\ \text{Recall}@k &= \frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \frac{|\mathcal{O}_u^{\text{test}} \cap \mathcal{O}_u^{\text{rec}}|}{|\mathcal{O}_u^{\text{test}}|} \\ F1@k &= \frac{2 \cdot \text{precision}@k \cdot \text{recall}@k}{\text{precision}@k + \text{recall}@k} \end{aligned}$$

Do not consider different importance of different ranking positions



## Area Under the Curve

Probability of a positive item being ranked higher than a negative item.

$$L_{BPR} = - \sum_{(u, i, j) \in \mathcal{O}^{\text{train}}} \log(\sigma(\mathbf{P}_u^T \mathbf{Q}_i - \mathbf{P}_u^T \mathbf{Q}_j))$$

- Gradient Descent to solve the minimization problem
- Add regularization to counteract overfitting

$$DCG@k = \sum_{i=1}^k \mathbb{I}(y_{u,i} = 1) \frac{1}{\log(1+i)}$$

## Discounted Cumulative Gain

$$IDCG@k = \sum_{i=1}^{\min(|\mathcal{O}_u^{\text{test}}|, k)} \frac{1}{\log(1+i)}$$

## Ideal Discounted Cumulative Gain

$$NDCG@k = \frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \frac{DCG@k}{IDCG@k}$$

## MF for Implicit Recommendation

### Normalized Discounted Cumulative Gain

$$AUC(u) = \frac{1}{|\mathcal{O}_u^{\text{test}}| |\mathcal{I} \setminus \mathcal{O}_u^{\text{test}}|} \sum_{i \in \mathcal{O}_u^{\text{test}}} \sum_{j \in \mathcal{I} \setminus \mathcal{O}_u^{\text{test}}} \mathbb{I}(\hat{y}_{u,i} > \hat{y}_{u,j})$$

$$AUC = \frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} AUC(u)$$

### Matrix Factorization

$$\begin{aligned} \mathbf{P}, \mathbf{Q} &= \arg \min_{\mathbf{P}, \mathbf{Q}} \mathcal{L} \\ &= \sum_{(u, i) \in \mathcal{O}^{\text{train}}} (\mathbf{P}_u^T \mathbf{Q}_i - y_{u,i})^2 + \lambda (\|\mathbf{P}\| + \|\mathbf{Q}\|) \end{aligned}$$

where  $\mathcal{O}^{\text{train}} = \mathcal{O} \cup \mathcal{O}^{\text{neg}}$

$$P_{\text{train}} = \mathcal{O} \cup \mathcal{O}^{\text{neg}}$$

$$\begin{aligned} \text{BPR: Bayesian Personalized Ranking} \\ AUC(u) &= \frac{1}{|\mathcal{O}_u^{\text{test}}| |\mathcal{I} \setminus \mathcal{O}_u^{\text{test}}|} \sum_{i \in \mathcal{O}_u^{\text{test}}} \sum_{j \in \mathcal{I} \setminus \mathcal{O}_u^{\text{test}}} \mathbb{I}(\hat{y}_{u,i} > \hat{y}_{u,j}) \end{aligned}$$

Main idea: to maximize AUC by making all positive items ranked higher than negative items.

During training, we do not have access to  $\mathcal{O}_u^{\text{test}}$ , but we have  $\mathcal{O}_u$  and  $\mathcal{O}_u^{\text{neg}}$ .

$$\hat{y}_{u,i} = \bar{r}_u + \frac{\sum_{u' \in \mathcal{N}} s(u, u') (r_{u',i} - \bar{r}_{u'})}{\sum_{u' \in \mathcal{N}} s(u, u')}$$

RMSE loss can be replaced by **cross-entropy loss**

$$\mathbf{P}, \mathbf{Q} = \arg \min_{\mathbf{P}, \mathbf{Q}} \mathcal{L}$$

$$= - \sum_{(u, i) \in \mathcal{O}^{\text{train}}} (y_{u,i} \log(\sigma(\mathbf{P}_u^T \mathbf{Q}_i)) + (1 - y_{u,i}) \log(1 - \sigma(\mathbf{P}_u^T \mathbf{Q}_i))) + \lambda (\|\mathbf{P}\| + \|\mathbf{Q}\|)$$

## Item-item CF

$$\hat{y}_{u,i} = \frac{\sum_{i' \in \mathcal{N}} s(i, i') y_{u,i'}}{\sum_{i' \in \mathcal{N}} s(i, i')}$$

## User-user CF

$$\hat{y}_{u,i} = \frac{\sum_{u' \in \mathcal{N}} s(u, u') y_{u',i}}{\sum_{u' \in \mathcal{N}} s(u, u')}$$

Class 1	Truth	Micro Ave. Table	
		True	False
Classifier: yes	100	20	1860
Classifier: no	10	890	900
Total	110	910	1860

Macrovaved precision:  $(0.5 + 0.9)/2 = 0.7$   
Macrovaved recall:  $100/120 = .83$   
Macrovaved score is dominated by score on frequent classes

```

mean = train_df["Rating"].mean()
unique_users_in_training = train_df["UserID"].unique()
unique_movies_in_training = train_df["MovieID"].unique()

mean_user_ratings = {}
for each in unique_users_in_training:
    looped_df = train_df.loc[train_df["UserID"] == each]
    mean_user_ratings[each] = looped_df["Rating"].mean()

mean_item_rating = {}
for each in unique_movies_in_training:
    looped_df = train_df.loc[train_df["MovieID"] == each]
    mean_item_rating[each] = looped_df["Rating"].mean()

```

```

mean_item_rating = {}
for each in unique_movies_in_training:
    looped_df = train_df.loc[train_df["MovieID"] == each]
    mean_item_rating[each] = looped_df["Rating"].mean()

pred_mat = []

for user_index in range(0, UNIQUE_USERS):
    movie_matrix = []
    for movie_index in range(0, UNIQUE_MOVIES):
        if ((movie_index not in unique_movies_in_training) or (mean_item_rating[movie_index] == 0.0)):
            b_i = 0
        else:
            b_i = mean_item_rating[movie_index] - miu

        if ((user_index not in unique_users_in_training) or (mean_user_ratings[user_index] == 0.0)):
            b_u = 0
        else:
            b_u = mean_user_ratings[user_index] - mu

        movie_matrix.append(miu + b_i + b_u)

    pred_mat.append(movie_matrix)

## Shuffle the dataset
ds = list(zip(self.sample_user, self.sample_movie))
np.random.shuffle(ds)

```

```

batch_gradient_updates = False

## Batching updates Part 1 of 3
if batch_gradient_updates:
    (glp_updates, glq_updates) = ([], [])

```

```

for u, i in ds:
    if self.train_indicator_mat[u, i] == 0.0:
        continue
    r_u_i = self.train_mat[u, i]
    (Pu, Qi) = (self.P[u], self.Q[i])
    ls = 2 * (Pu @ Q[i] - r_u_i)
    gL_Pu = ((ls * Qi) + (self.reg * Pu))
    gL_Qi = ((ls * Pu) + (self.reg * Qi))

    self.P[u] += self.lr * gL_Pu
    self.Q[i] += self.lr * gL_Qi

```

```

def masked_rmse(x,y, mask=1):
    inner = ((x*mask) - (y*mask)) ** 2
    o = (np.sum(inner) / np.sum(mask))
    return np.sqrt(o)

```

```

pred = np.matmul(self.P, self.Q.T)
regularization = (self.reg * (np.sum(self.P**2) + np.sum(self.Q**2)))
epoch_loss = (np.sum((pred+self.train_indicator_mat - self.train_mat)**2)) + regularization

```

```

def get_support(arr, items):
    # Returns the mean support of an iterable containing items that are valid indices of our ohe array
    # Equivalent to (sumation(N_i*sumation_k{[(user0_iitem0 ... U ... user_k_itemN)] / K}) / N,
    # where K is the number of users and N is the number of items in the iterable
    return np.product([arr[:,item] for item in items], axis=0).mean()

```

Suppose we have  $N$  data samples,  $C$  different class labels,  $K$  clusters and  $a_{ck}$  number of data points belonging to the class  $c$  and cluster  $k$ .

$$\text{Homogeneity: } h = \frac{H(C|K)}{H(C)}$$

$$\text{Completeness: } c = \frac{H(K|C)}{H(K)}$$

$$\text{V-measure: } V_B = \frac{(1+\beta)hc}{\beta h + c}$$

### Cross-entropy error

$$E(\mathbf{w}_1, \dots, \mathbf{w}_C) = - \sum_{n=1}^N \sum_{c=1}^C y_{nc} \log p(y_{nc} = 1 | \mathbf{w}_1, \dots, \mathbf{w}_C)$$

- Optimization with gradient descent, convex function
- Computational details are out of scope
- But the gradient vector w.r.t. each weight  $\mathbf{w}_c$  looks like this

$$\nabla E_{\mathbf{w}_c} = \sum_{n=1}^N [p(y_{nc} = 1 | \mathbf{w}_1, \dots, \mathbf{w}_C) - y_{nc}] \mathbf{x}_n$$

```

running_sum = 0
count_of_ratings = 0
for user_index in range(0, UNIQUE_USERS):
    movie_matrix = []
    for movie_index in range(0, UNIQUE_MOVIES):
        if test_mat[user_index][movie_index] != 0.0:
            running_sum = running_sum + (test_mat[user_index][movie_index])
            count_of_ratings += 1
RMSE_ERROR = (running_sum / count_of_ratings)**(1/2)

```

## BPR Loss to Optimize AUC

$$AUC(u) = \frac{1}{|\mathcal{O}_u||\mathcal{O}_u^{neg}|} \sum_{i \in \mathcal{O}_u} \sum_{j \in \mathcal{O}_u^{neg}} \mathbb{I}(\hat{y}_{u,i} > \hat{y}_{u,j})$$

$$\mathcal{L}_{BPR}(u) = - \sum_{i \in \mathcal{O}_u} \sum_{j \in \mathcal{O}_u^{neg}} \log(\sigma(\hat{y}_{u,i} - \hat{y}_{u,j}))$$

$$\mathcal{L}_{BPR} = - \sum_{(u,i,j) \in \mathcal{O}_{train}} \log(\sigma(\hat{y}_{u,i} - \hat{y}_{u,j}))$$

Completeness:

$$c = 1 - \frac{H(K|C)}{H(K)}$$

$$H(K|C) = - \sum_{c=1}^C \sum_{k=1}^K \frac{a_{ck}}{N} \log \frac{a_{ck}}{\sum_{k=1}^K a_{ck}}$$

$$H(K) = - \sum_{k=1}^K \frac{\sum_{c=1}^C a_{ck}}{N} \log \frac{\sum_{c=1}^C a_{ck}}{N}$$

Entropy of the whole dataset w.r.t cluster labels

42

$$h = 1 - \frac{H(C|K)}{H(C)}$$

$$H(C|K) = - \sum_{k=1}^K \sum_{c=1}^C \frac{a_{ck}}{N} \log \frac{a_{ck}}{\sum_{k=1}^C a_{ck}}$$

$$H(C) = - \sum_{c=1}^C \frac{\sum_{k=1}^K a_{ck}}{N} \log \frac{\sum_{k=1}^K a_{ck}}{N}$$

Entropy of the whole dataset w.r.t class labels

## Maximum Likelihood Estimation (MLE)

If we don't know the parameter  $\theta$  of the biased coin but we observe  $y_1, \dots, y_n$  from  $n$  independent experiments, then we can estimate  $\theta$  by maximizing the log likelihood:

$$\theta^* = \arg \max_{\theta} \sum_{i=1}^n y_i \log \theta + (1 - y_i) \log(1 - \theta)$$

$$\theta^* = \frac{\sum_{i=1}^n y_i}{n}$$

### $F_1$ or F-score

$$\bullet F = \frac{2 \times P \times R}{P + R} = \frac{2 \cdot 0.9440 \cdot 0.9934}{0.9440 + 0.9934} = 0.9681$$

### $F_{\beta}$ : weighted combination

$$\bullet F = \frac{(1+\beta^2) \times P \times R}{\beta^2 \times P + R}$$

### Logistic Regression

- Linear combination of input features  $\mathbf{w}^T \mathbf{x}$
- Transform through sigmoid function  $\sigma(\mathbf{w}^T \mathbf{x}) \rightarrow$  interpretable as probability
- Decision rule based on whether  $\sigma(\mathbf{w}^T \mathbf{x}) \leq 0.5$
- Evaluation through data likelihood, or cross-entropy error

$$E(\mathbf{w}) = - \sum_{n=1}^N \{y_n \log [\sigma(\mathbf{w}^T \mathbf{x}_n)] + (1 - y_n) \log [1 - \sigma(\mathbf{w}^T \mathbf{x}_n)]\}$$

### Optimization through gradient descent

### Multinomial Regression

- Conditional logit model:  $p(y = c | \mathbf{x}, \mathbf{w}_c) = \frac{e^{\mathbf{w}_c^T \mathbf{x}}}{\sum_{c=1}^C e^{\mathbf{w}_c^T \mathbf{x}}}$
- Similar to 2-class logistic regression
  - compute negative cross-entropy and perform gradient descent

```

def kneans(x, number_of_centroids:int = 10, maxIt:int = 5, centroid_init_fn=None, debug=False, l1=None, labels=None):
    # Initialize centroids
    if centroid_init_fn is None:
        centroids = np.random.uniform(min(x), max(x), size=(number_of_centroids, len(x)))
    else:
        centroids = centroid_init_fn(x, number_of_centroids)

    # Compute initial error
    minum_loc = np.argmin(x, axis=0)
    errors = []
    for i in range(maxIt):
        # Compute working v0
        for i in range(number_of_centroids):
            working_v0 = np.zeros(len(x))

        # Compute distances
        def get_distances():
            working_v0_by = working_v0
            for i in range(number_of_centroids):
                working_v0_by = working_v0_by + zip(rep(x), centroids))

        # Compute distances to centroids
        points_distances_to_centroids = get_distances()
        minimum_locs = np.argmin(points_distances_to_centroids, axis=0)

        # Compute errors
        errors.append(np.sum(group_distances**2))

        # Compute mean error
        mean_error = np.sum(errors, )
        homogeneity, completeness, v_measure = sklearn.metrics.homogeneity_completeness_v_measure(labels, minimum_locs)

    def prep_return():
        points_distances_to_centroids = pdc = get_distances()
        loc = np.argmax(pdcc, axis=0)

        (accuracy, preds, _labels) = get_accuracy(loc, labels)

        rd = dict(
            points=points_distances_to_centroids,
            group_ids=group_ids,
            group_members=group_members,
            centroids=centroids,
            mean_error=mean_error,
            points_distances_to_centroids=points_distances_to_centroids,
            homogeneity=homogeneity,
            completeness=completeness,
            v_measure=v_measure,
            loc=loc,
            labels=labels,
            preds=preds,
            accuracy=accuracy,
            step=1,
        )
        rd[(knp.array(v) for k, v in rd.items())]
        return rd

    # If l1 is not None or i == maxIt - 1:
    if (l1 == mean_error) or (i == maxIt - 1):
        return prep_return()
    loc = mean_error

```