

Sistemas Recomendadores

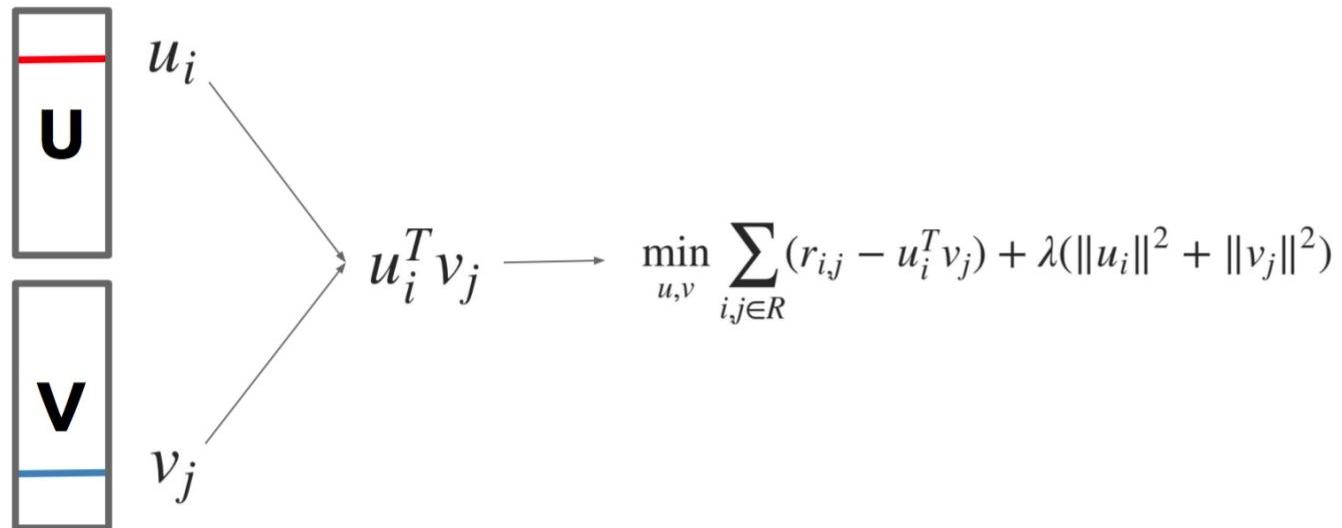
IIC-3633

Recomendación basada en contenido
Parte 3

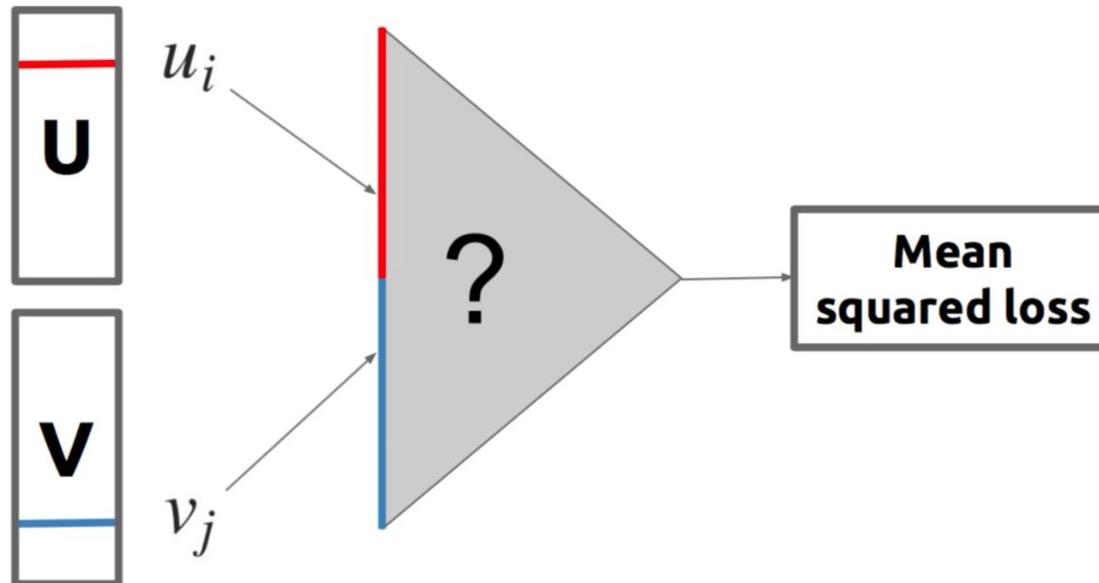
Esta clase

1. Recomendación basada en contenido imágenes y audio.

Factorización Matricial con Redes Neuronales



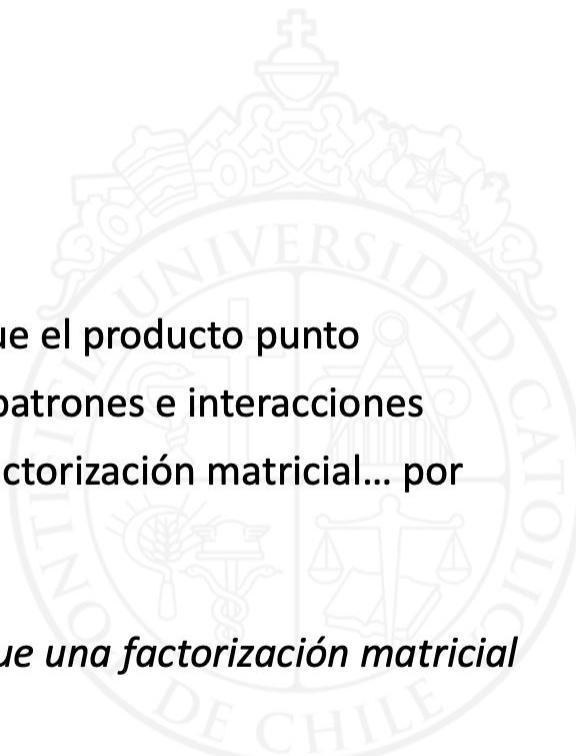
Factorización Matricial con Redes Neuronales



Factorización Matricial con Redes Neuronales

- Redes neuronales y factorización matricial son similares
 - Uso de embeddings
 - Pérdida mínimos cuadrados
 - Óptimo con descenso de gradiente
- Una red neuronal puede aprender más combinaciones que el producto punto
- Una red neuronal requiere muchos datos para aprender patrones e interacciones
- El uso de redes neuronales no es mejor que una buena factorización matricial... por ahora...

Las redes neuronales permiten capturar más información que una factorización matricial como información de texto, imágenes, música, etc.



Redes Neuronales para Factorización Matricial

- <https://arxiv.org/pdf/1907.06902.pdf>

Are We Really Making Much Progress? A Worrying Analysis of Recent Neural Recommendation Approaches

Maurizio Ferrari Dacrema
Politecnico di Milano, Italy
maurizio.ferrari@polimi.it

Paolo Cremonesi
Politecnico di Milano, Italy
paolo.cremonesi@polimi.it

Dietmar Jannach
University of Klagenfurt, Austria
dietmar.jannach@aau.at

- Si solo usamos datos de interacciones, puede que las redes neuronales no nos ofrezcan tanto beneficio.

Ventajas de las Redes Neuronales

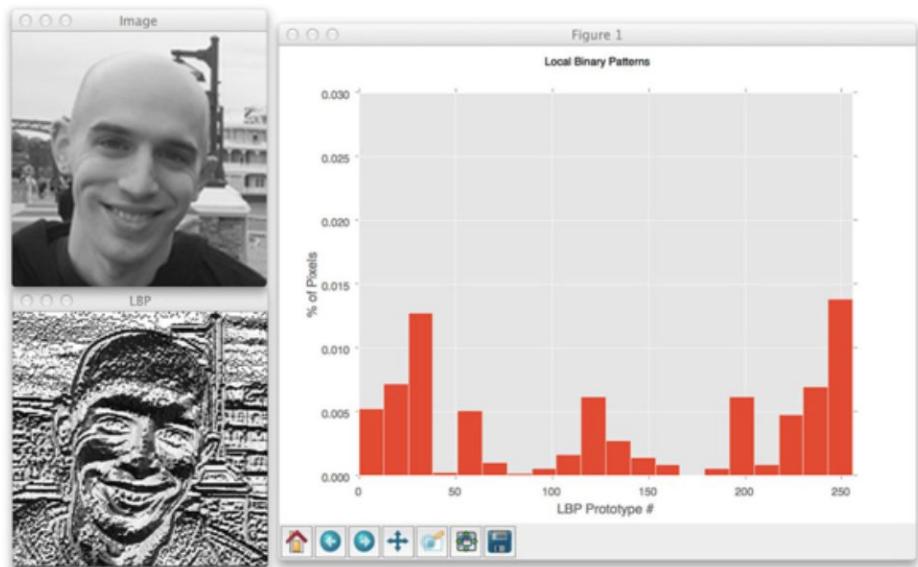
- Información de texto
 - Ejemplo: descripción de productos, comentarios de usuarios
 - Extracción: RNN
 - Aplicaciones: noticias, libros, e-commerce
- Imágenes
 - Ejemplo: foto de productos, thumbnail de videos
 - Extracción: CNN
 - Aplicaciones: moda, video
- Música / audio
 - Ejemplo: spotify
 - Extracción: CNN y RNN
 - Aplicación: música



Extracción de features manuales de imágenes

LBP

- Finalmente se calcula un histograma que tabula el número de ocasiones en que cada patrón LBP ocurrió.
- Podemos pensar en este histograma como un vector de features.



Fuente: <https://www.pyimagesearch.com/2015/12/07/local-binary-patterns-with-python-opencv/>

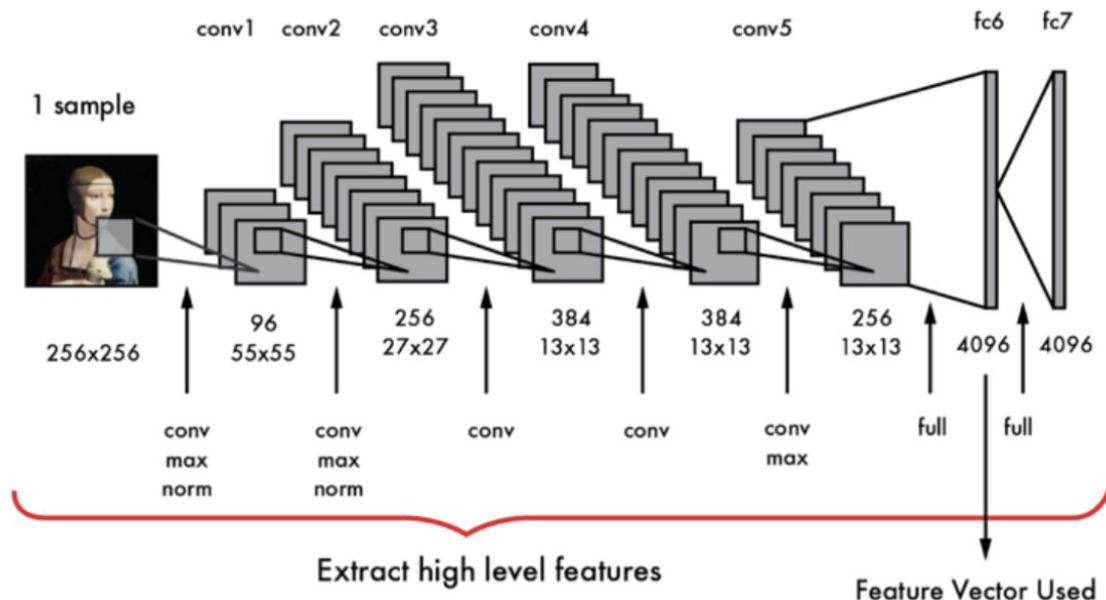
Atractivas



No Atractivas

Features manuales versus Deep Learning

- Con DL podemos usar features aprendidas automáticamente con una red neuronal pre-entrenada para otra tarea: clasificación de objetos del dataset Imagenet.



Krizhevsky, A., Sutskever, I., & Hinton, G. E. (2012). Imagenet classification with deep convolutional neural networks. In *Advances in neural information pro*

Contenido Visual y Musical

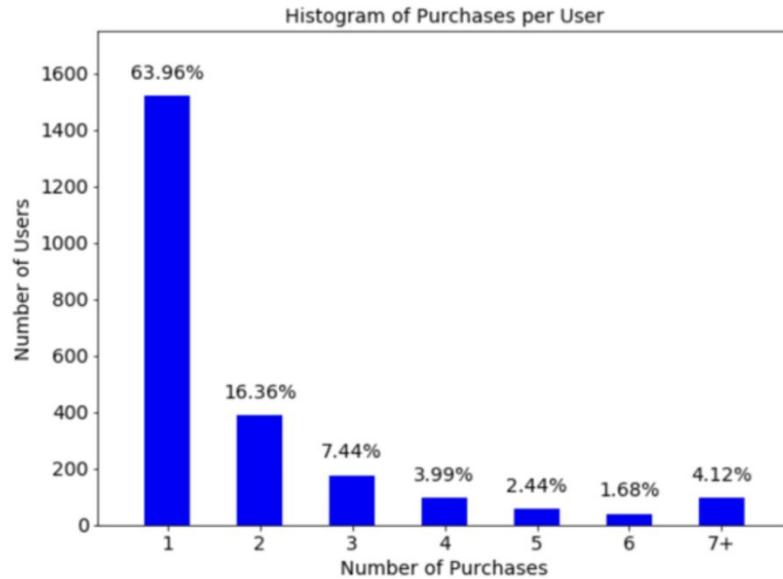
- La representación del contenido visual y musical no es tan intuitiva como en el caso de texto.
- Si bien hay investigación madura en como representar música y texto, los modelos de Deep Learning de los años recientes han modificado profundamente esta área:
 - Modelos anteriores hacían feature (características) engineering
 - Modelos modernos usan Deep Learning (DL) para aprender las características.

Ejemplo

- Messina, P., Dominguez, V., Parra, D., Trattner, C., & Soto, A. (2019). Content-based artwork recommendation: integrating painting metadata with neural and manually-engineered visual features. *User Modeling and User-Adapted Interaction*, 29(2), 251-290.

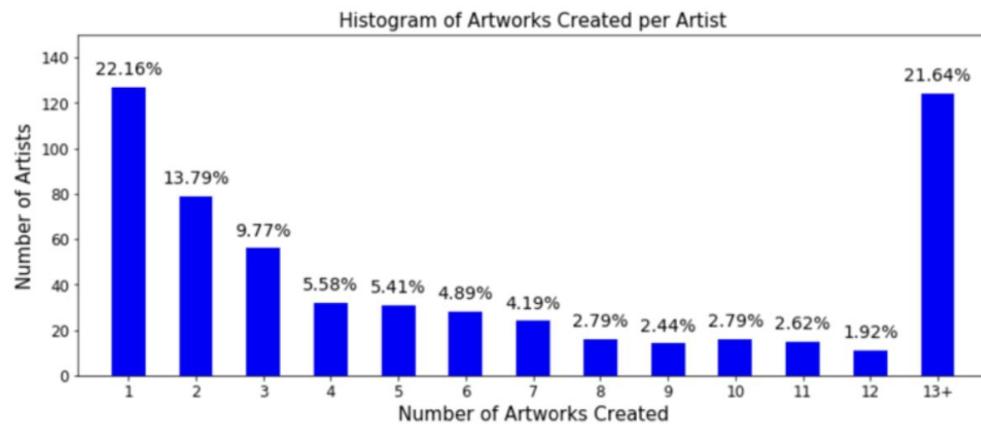
Dataset

- 5336 transacciones (compras)
- 2378 usuarios
- 6040 obras de arte



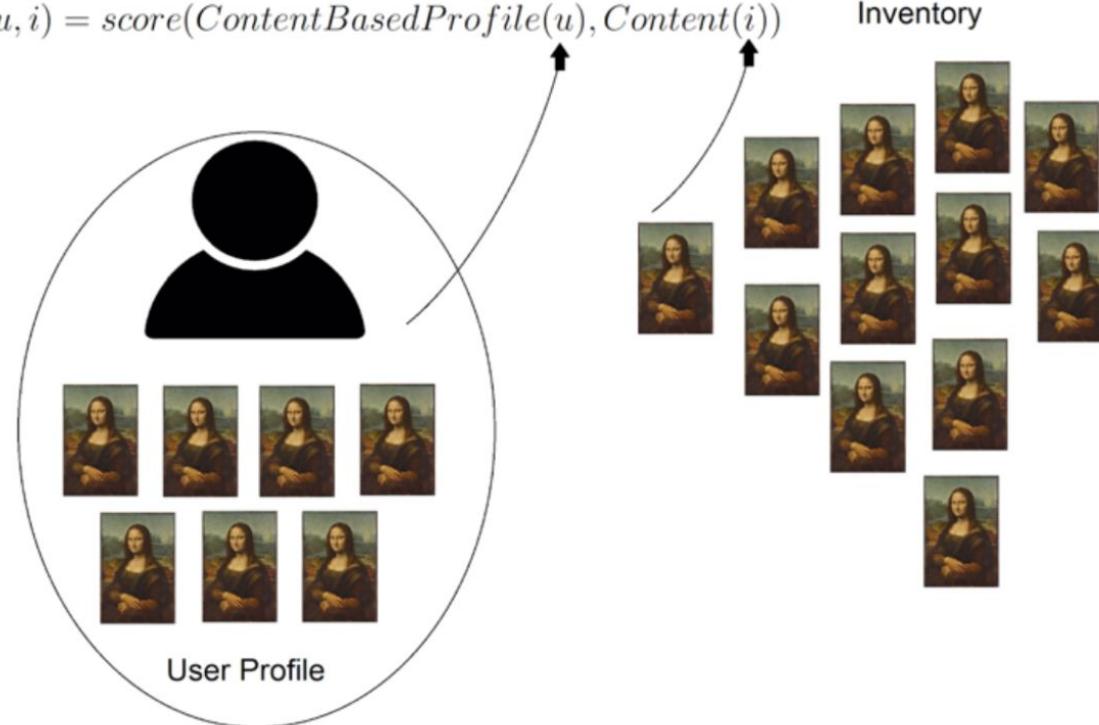
Dataset

- 573 artistas en total
- Un artista por obra
- 10,54 obras por artista en promedio



Recomendación basada en contenido

$$s(u, i) = \text{score}(\text{ContentBasedProfile}(u), \text{Content}(i))$$



Métodos utilizados

Attribute values: color, subject, style, mood.

1. Most Popular Curated Attribute Value (MPCAV)
2. Personalized Most Popular Curated Attribute Value (PMPCAV)
3. Personalized Favorite Artist (FA)
4. **Learned Visual Features: Deep Convolutional Neural Networks (CNN)**
5. **Handcrafted Visual Features (HVF)**
6. Hybrid Recommendations (Hybrid)

Score con HVF (manuals)

- Análogo a las CNNs: similaridad coseno + agregaciones (max, average, average-top-k)

$$\text{sim}(V_i^{\text{Attract}}, V_j^{\text{Attract}}) = \cos(V_i^{\text{Attract}}, V_j^{\text{Attract}}) \quad \text{sim}(V_i^{\text{LBP}}, V_j^{\text{LBP}}) = \cos(V_i^{\text{LBP}}, V_j^{\text{LBP}})$$

- También probamos un híbrido Attractiveness + LBP:

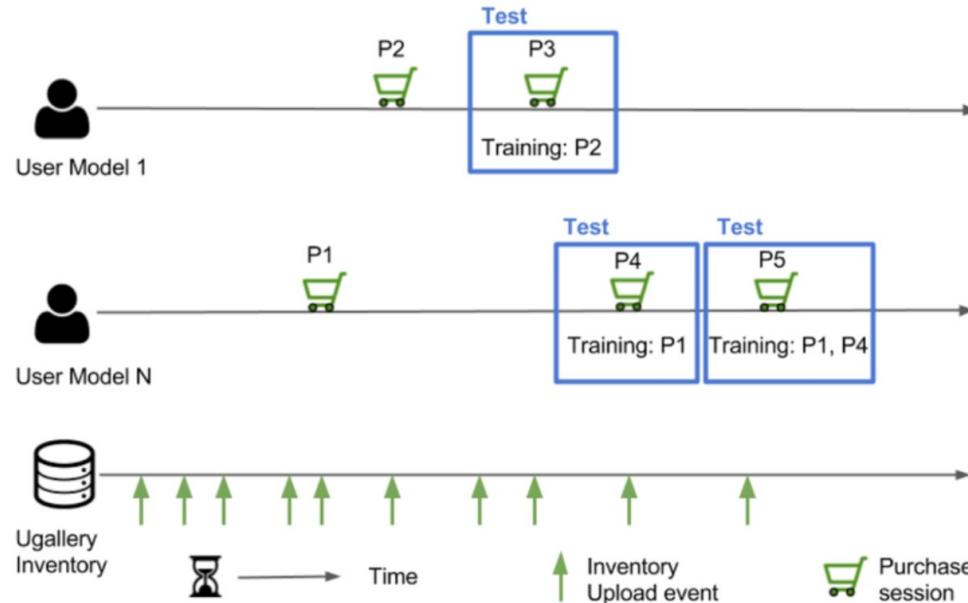
$$\begin{aligned} \text{score}(u, i)_{\text{HVF}} &= \alpha_1 \cdot \text{score}(u, i)_{\text{Attractiveness}} \\ &\quad + \alpha_2 \cdot \text{score}(u, i)_{\text{LBP}} \end{aligned}$$

Score con features de red neuronal CNN

$$score(u, i)_X = \begin{cases} \max_{j \in P_u} \{sim(V_i^X, V_j^X)\} & (maximum) \\ \frac{\sum_{j \in P_u} sim(V_i^X, V_j^X)}{|P_u|} & (average) \\ \frac{\sum_{r=1}^{\min\{K, |P_u|\}} \max_{j \in P_u} {}^{(r)}\{sim(V_i^X, V_j^X)\}}{\min\{K, |P_u|\}} & (average top K) \end{cases}$$

$$sim(V_i, V_j) = cos(V_i, V_j) = \frac{V_i \cdot V_j}{\|V_i\| \|V_j\|}$$

Evaluación offline



Resultados

| ID | Method | nD@20 | R@20 | P@20 | F1@20 |
|----|-------------------------|---------------------|---------------------|---------------------|---------------------|
| 1 | CNN (All) | .1295 ³ | .1702 ⁴ | .0151 ² | .0248 ² |
| 2 | CNN (ResNet50) | .1247 ³ | .1628 ⁴ | .0145 ⁵ | .0236 ⁴ |
| 3 | CNN (AlexNet) | .1081 ⁴ | .1461 ⁴ | .0135 ⁵ | .0216 ⁴ |
| 4 | CNN (VGG19) | .1008 ⁵ | .1398 ⁸ | .0124 ⁶ | .0205 ⁵ |
| 5 | CNN (InceptionV3) | .1007 ⁶ | .1332 ⁸ | .0125 ⁴ | .0201 ⁶ |
| 6 | CNN (NASNet Large) | .0998 ⁸ | .1379 ⁸ | .0120 ⁸ | .0197 ⁷ |
| 7 | CNN (InceptionResNetV2) | .0932 ⁸ | .1300 ⁸ | .0119 ⁸ | .0192 ⁸ |
| 8 | HVF (LBP) | .0507 ⁹ | .0736 ¹¹ | .0068 ⁹ | .0107 ⁹ |
| 9 | HVF (LBP + Attr.) | .0493 ¹¹ | .0728 ¹¹ | .0064 ¹⁰ | .0103 ¹¹ |
| 10 | HVF (Attractiveness) | .0407 ¹¹ | .0628 ¹¹ | .0059 ¹¹ | .0095 ¹¹ |
| 11 | Random | .0097 | .0200 | .0015 | .0025 |

Stat. significance by multiple t-tests, Bonferroni corr.

$$\alpha_{bonf} = \alpha/n = 0.05/55 = .00091.$$

Diversidad

| ID | Method | F1@20 | D@10 visual cluster | D@10 visual pairwise | D@10 artist | D@10 jaccard pairwise | D@10 color | D@10 medium |
|----|-------------------------------------|---------------------|-----------------------------|---------------------------|-----------------------------|---------------------------|-----------------------------|----------------------------|
| 1 | Hybrid ₁ (FA+CNN+PMPCAV) | .0333 ² | 10.0697 ⁴ | .3952 ² | 8.4375 ³ | .7433 ² | 11.7362¹² | 2.2719 ³ |
| 2 | Hybrid ₂ (FA+CNN) | .0325 ⁵ | <u>9.1883</u> | <u>.3803</u> | <u>7.6165⁴</u> | .7730 ¹ | 12.0959 ³ | 2.7902 ⁴ |
| 3 | Hybrid ₃ (FA+PMPCAV) | .0312 ⁴ | 11.8327 ⁹ | .4297 ⁹ | <u>7.8472²</u> | .7214 ⁴ | <u>11.8309¹²</u> | <u>2.0459¹²</u> |
| 4 | FA | .0295 ⁵ | 9.7124 ² | .4092 ⁸ | 2.8809 | .7068 | 11.9983 ³ | 2.3864 ¹ |
| 5 | CNN (All) | .0248 ⁶ | <u>9.6688²</u> | <u>.3913²</u> | 12.6822 ¹ | .8488 ¹⁶ | 12.6514 ⁷ | 3.3951 ² |
| 6 | CNN (ResNet50) | .0236 ⁸ | 10.1429 ⁴ | .3968 ⁷ | 12.6804 ¹ | .8524 ⁵ | 12.7164 ⁷ | 3.4399 ² |
| 7 | CNN (AlexNet) | .0216 ⁸ | 10.1732 ⁴ | <u>.3923²</u> | 13.0314 ⁵ | .8502 ¹⁶ | 12.4317 ² | 3.5119 ⁵ |
| 8 | CNN (VGG19) | .0205 ⁹ | 10.6845 ⁷ | .4016 ⁶ | 14.3341 ¹¹ | .8648 ⁶ | 13.0546 ¹⁵ | 3.5386 ⁶ |
| 9 | CNN (InceptionV3) | .0201 ¹⁰ | 11.2208 ⁸ | .4195 ¹¹ | 13.8768 ⁷ | .8712 ⁸ | 13.1360 ¹⁵ | 3.6926 ¹¹ |
| 10 | CNN (NASNet Large) | .0197 ¹² | 11.0767 ⁸ | .4144 ⁸ | 14.0180 ¹⁶ | .8697 ⁸ | 13.1435 ¹⁵ | 3.6827 ⁸ |
| 11 | CNN (InceptionResNetV2) | .0192 ¹³ | 11.1313 ⁸ | .4151 ⁴ | 14.0232 ¹⁶ | .8703 ⁸ | 13.1871 ¹⁵ | 3.6072 ⁶ |
| 12 | PMPCAV(All) | .0156 ¹³ | 13.6607 ³ | .4498 ³ | 14.4608 ¹¹ | .7429 ³ | 11.0691 | 1.8303 ¹⁶ |
| 13 | HVF (LBP) | .0107 ¹⁴ | 14.6874 ¹² | .4667 ¹² | 15.8733 ¹² | .8949¹⁵ | 13.9820 ¹¹ | 4.1296⁹ |
| 14 | HVF (LBP + Attr.) | .0103 ¹⁶ | 15.3969 ¹³ | .4732 ¹³ | 16.3359¹³ | .8961¹⁵ | 14.0628¹¹ | 4.0633⁹ |
| 15 | HVF (Attractiveness) | .0095 ¹⁷ | 15.4358¹³ | .4743¹³ | 16.5584¹³ | .8850 ⁹ | 12.8210 ⁷ | 4.0569 ⁹ |
| 16 | MPCAV(Medium) | .0081 ¹⁷ | 15.4375¹³ | .4829¹⁵ | 13.7440 ⁷ | .7844 ² | 14.3841¹⁴ | 1.0017 |
| 17 | Random | .0025 | 17.4006¹⁶ | .4972¹⁶ | 18.4069¹⁵ | .9123¹⁴ | 14.2869¹⁴ | 4.5804¹³ |

Statistical significance was obtained using multiple pairwise t-tests with Bonferroni correction,
 $\alpha_{bonf} = \alpha/n = 0.05/136 = .00037$.

Evaluación online (8 curadores de UGallery)

| Madeleine's profile | | method 1 | method 2 | method 3 | method 4 | method 5 |
|---------------------|--|---------------------|---------------------|---------------------|---------------------|---------------------|
| Liked Artworks | | | | | | |
| | | | | | | |
| | | | | | | |
| | | | | | | |
| | | | | | | |
| | | | | | | |
| Rating | | Successfully rated! |
| Rating | | Successfully rated! |
| Rating | | Successfully rated! |
| Rating | | Successfully rated! |

Evaluación online (8 curadores de UGallery)

| Name | nD@5 | nD@10 | P@5 | P@10 |
|--------------------|---------------|---------------|---------------|---------------|
| Hybrid(FA+CNN+HVF) | 0.9042 | 0.8913 | 0.7500 | 0.6750 |
| Hybrid(CNN+HVF) | 0.6747 | 0.6638 | 0.5000 | 0.4250 |
| CNN | 0.7176 | 0.6947 | 0.5000 | 0.4000 |
| FA | 0.4276 | 0.5662 | 0.3000 | 0.4000 |
| HVF | 0.5498 | 0.5314 | 0.3500 | 0.2625 |

Otro método: Visual BPR

- VBPR = Visual Bayesian Personalized Ranking (R. He & McAuley, 2016)

$$\hat{x}_{u,i} = \beta_i + \gamma_u^T \gamma_i + \theta_u^T (Ef_i) + \beta'^T f_i$$

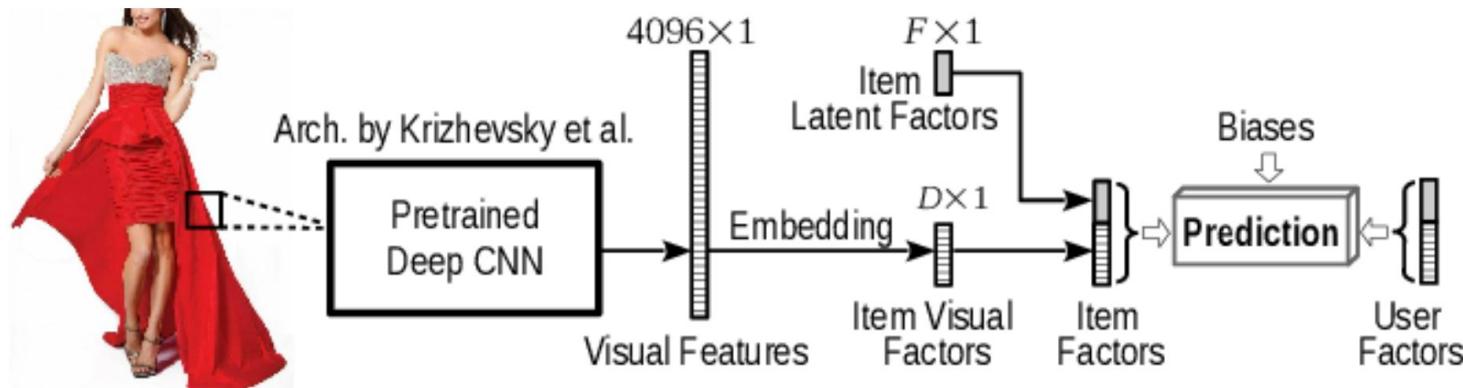
Vector de features desde
AlexNet CNN

- Variables se aprenden con BPR-OPT (Rendle et al., 2009)

$$D_S = \{(u, i, j) | u \in U \wedge i \in I_u^+ \wedge j \in I \setminus I_u^+\}$$

$$\sum_{(u,i,j) \in D_S} \ln(\sigma(\hat{x}_{u(i)}(\Theta))) - \lambda_\Theta \|\Theta\|^2 \quad \hat{x}_{u(i)}(\Theta) = \hat{x}_{u,i} - \hat{x}_{u,j}$$

VBPR



He, R., & McAuley, J. (2016). VBPR: Visual Bayesian Personalized Ranking from Implicit Feedback. AAAI.

Recomendación de música

Música

- Ejemplo de recomendación de Spotify
 - Muchos sistemas, incluso a la fecha, representan música con diferentes features manuales, siendo MFCC (Mel Frequency Cepstral Coefficients), los más populares. Se obtienen así:
 - Separar la señal en pequeños tramos.
 - A cada tramo aplicarle la Transformada de Fourier discreta y obtener la potencia espectral de la señal.
 - Aplicar el banco de filtros correspondientes a la Escala Mel al espectro obtenido en el paso anterior y sumar las energías en cada uno de ellos.
 - Tomar el logaritmo de todas las energías de cada frecuencia mel
 - Aplicarle la transformada de coseno discreta a estos logaritmos.

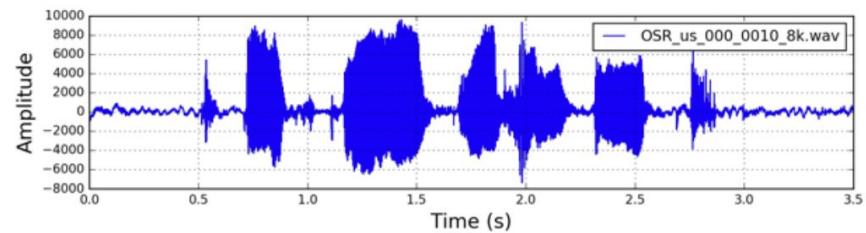
Escala Mel

- La escala Mel es una escala que relaciona la frecuencia percibida de un tono con la frecuencia medida real. Escala la frecuencia para que coincida más con lo que el oído humano puede escuchar (los humanos son mejores para identificar pequeños cambios en el habla a frecuencias más bajas).
- Una frecuencia en Hertz se convierte a escala Mel con:

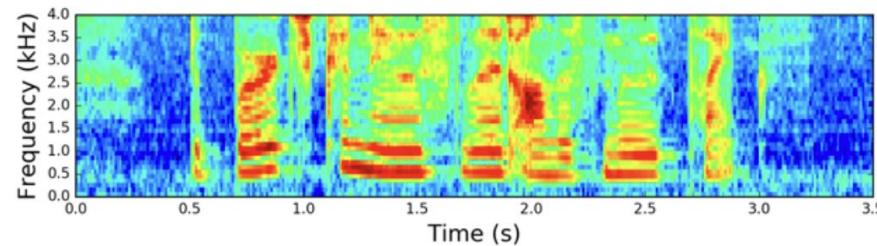
$$\text{Mel}(f) = 2595 \log\left(1 + \frac{f}{700}\right)$$

<https://medium.com/prathena/the-dummys-guide-to-mfcc-aceab2450fd>

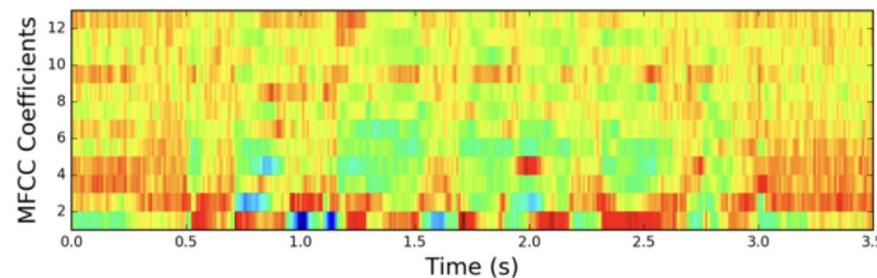
MFCC



Señal de audio



Espectrograma



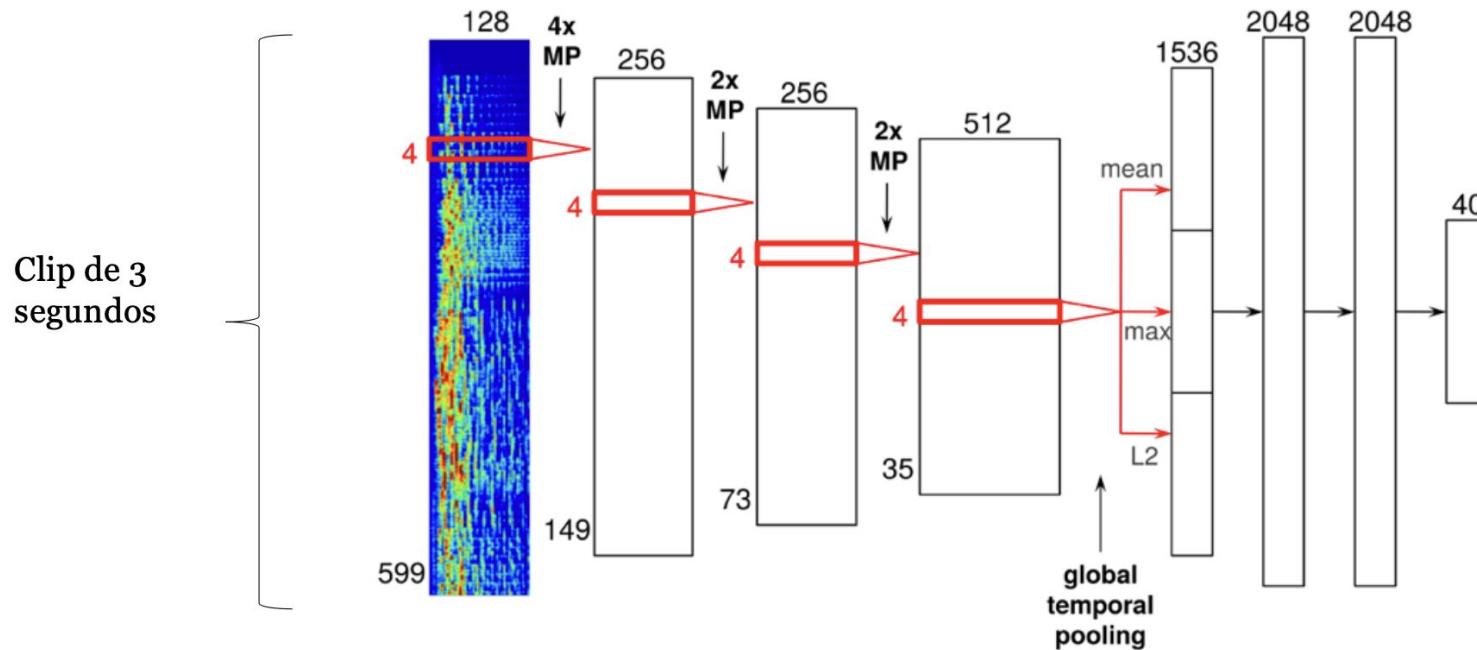
MFCC

Approach con redes neuronales

- En Van den Oord et al (2013) comparan un approach tradicional basado en MFCC con DL features.
- El approach tradicional:
 - **Extract MFCCs from the audio signals.** We computed 13 MFCCs from windows of 1024 audio frames, corresponding to 23 ms at a sampling rate of 22050 Hz, and a hop size of 512 samples. We also computed first and second order differences, yielding 39 coefficients in total.
 - **Vector quantize the MFCCs.** We learned a dictionary of 4000 elements with the K-means algorithm and assigned all MFCC vectors to the closest mean.
 - **Aggregate them into a bag-of-words representation.** For every song, we counted how many times each mean was selected. The resulting vector of counts is a bag-of-words feature representation of the song.

Van den Oord, A., Dieleman, S., & Schrauwen, B. (2013). Deep content-based music recommendation. In Advances in neural information processing systems (pp. 2643-2651).

DL en Van den Oord et al (2013)



Latent factor prediction

- Como baseline se obtienen factores latentes de usuarios e items usando WMF (ALS)
- La tarea es predecir los factores latentes directamente desde el audio, con los siguientes métodos:
 - Linear regression trained on the bag-of-words representation
 - A multi-layer perceptron (MLP) trained on the same bag-of-words representation.
 - A convolutional neural network trained on log-scaled mel-spectrograms to minimize the mean squared error (MSE) of the predictions.
 - The same convolutional neural network, trained to minimize the weighted prediction error (WPE) from the WMF objective instead.

Latent factor prediction

- Dataset: Million Song Dataset (MSD)
- <http://millionsongdataset.com/>



Latent factor prediction: resultados

MFCC

| Model | mAP | AUC |
|-------------------|---------|---------|
| MLR | 0.01801 | 0.60608 |
| linear regression | 0.02389 | 0.63518 |
| MLP | 0.02536 | 0.64611 |
| CNN with MSE | 0.05016 | 0.70987 |
| CNN with WPE | 0.04323 | 0.70101 |

Table 2: Results for all considered models on a subset of the dataset containing only the 9,330 most popular songs, and listening data for 20,000 users.

| Model | mAP | AUC |
|-------------------|---------|---------|
| random | 0.00015 | 0.49935 |
| linear regression | 0.00101 | 0.64522 |
| CNN with MSE | 0.00672 | 0.77192 |
| upper bound | 0.23278 | 0.96070 |

Table 3: Results for linear regression on a bag-of-words representation of the audio signals, and a convolutional neural network trained with the MSE objective, on the full dataset (382,410 songs and 1 million users). Also shown are the scores achieved when the latent factor vectors are randomized, and when they are learned from usage data using WMF (upper bound).

Evaluación por muestras

| Query | Most similar tracks (WMF) | Most similar tracks (predicted) |
|--------------------------|--|--|
| Jonas Brothers - Hold On | Jonas Brothers - Games Miley Cyrus - G.N.O. (Girl's Night Out) Miley Cyrus - Girls Just Wanna Have Fun Jonas Brothers - Year 3000 Jonas Brothers - BB Good | Jonas Brothers - Video Girl Jonas Brothers - Games New Found Glory - My Friends Over You My Chemical Romance - Thank You For The Venom My Chemical Romance - Teenagers |
| Beyoncé - Speechless | Beyoncé - Gift From Virgo Beyoncé - Daddy Rihanna / J-Status - Crazy Little Thing Called Love Beyoncé - Dangerously In Love Rihanna - Haunted | Daniel Bedingfield - If You're Not The One Rihanna - Haunted Alejandro Sanz - Siempre Es De Noche Madonna - Miles Away Lil Wayne / Shanell - American Star |
| Coldplay - I Ran Away | Coldplay - Careful Where You Stand Coldplay - The Goldrush Coldplay - X & Y Coldplay - Square One Jonas Brothers - BB Good | Arcade Fire - Keep The Car Running M83 - You Appearing Angus & Julia Stone - Hollywood Bon Iver - Creature Fear Coldplay - The Goldrush |
| Daft Punk - Rock'n Roll | Daft Punk - Short Circuit Daft Punk - Nightvision Daft Punk - Too Long (Gonzales Version) Daft Punk - Aerodynamite Daft Punk - One More Time / Aerodynamic | Boys Noize - Shine Shine Boys Noize - Lava Lava Flying Lotus - Pet Monster Shotglass LCD Soundsystem - One Touch Justice - One Minute To Midnight |

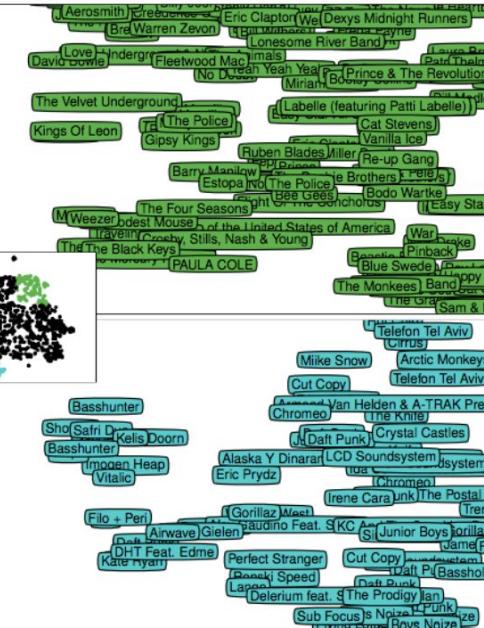
Table 4: A few songs and their closest matches in terms of usage patterns, using latent factors obtained with WMF and using latent factors predicted by a convolutional neural network.

Factores latentes predichos por CNN

HIP-HOP



POP



ALTERNATIVE
ROCK



ELECTRONIC

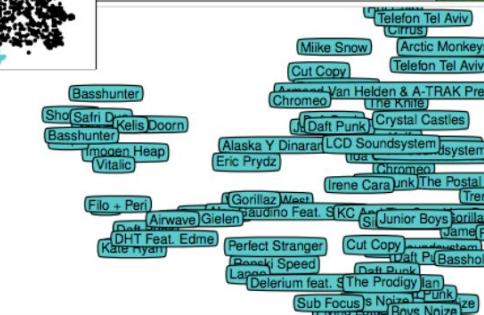


Figure 1: t-SNE visualization of the distribution of predicted usage patterns, using latent factors predicted from audio. A few close-ups show artists whose songs are projected in specific areas. We can discern hip-hop (red), rock (green), pop (yellow) and electronic music (blue). This figure is best viewed in color.

Resumen

- Para imágenes y música, usar features aprendidas por modelos de deep learning puede ser muy útil y evita el costo de la “ingeniería de características manual”
- Una debilidad de este approach es que las características aprendidas son difícilmente explicables.

Ideas proyectos

The Effect of Explanations and Algorithmic Accuracy on Visual Recommender Systems of Artistic Images

Vicente Dominguez

Pontificia Universidad Católica de Chile & IMFD
Santiago, Chile
vidominguez@uc.cl

Ivania Donoso-Guzmán

Pontificia Universidad Católica de Chile & Conversica
Santiago, Chile
indonoso@uc.cl

Pablo Messina

Pontificia Universidad Católica de Chile & IMFD
Santiago, Chile
pamessina@uc.cl

Denis Parra

Pontificia Universidad Católica de Chile & IMFD
Santiago, Chile
dparra@ing.puc.cl

Setting

Prueban 3 interfaces para recomendación de arte:

- Baseline (rating predicho y top N recomendaciones)
- Explicación con ítems similares, te recomiendo esto porque consumiste esto.
- Explicación con features visuales (attractiveness)



Figure 3: Interface 1: Baseline recommendation interface without explanations.

| Recommended Artwork | Explanation |
|---|---|
|  <p>Successfully rated!</p> <p>★★★★★</p> | <p>Recommended because:</p> <p>it's 85.31% similar to this artwork it's 71.48% similar to this artwork it's 64.0% similar to this artwork that you like</p>    <p>With an average of 73.62%</p> |
| Recommended Artwork | Explanation |
|  | <p>Recommended because:</p> <p>it's 81.96% similar to this artwork it's 70.10% similar to this artwork it's 68.5% similar to this artwork that you like</p>    |

Figure 4: Interface 2: Explainable recommendation interface with textual explanations and top-3 similar images.

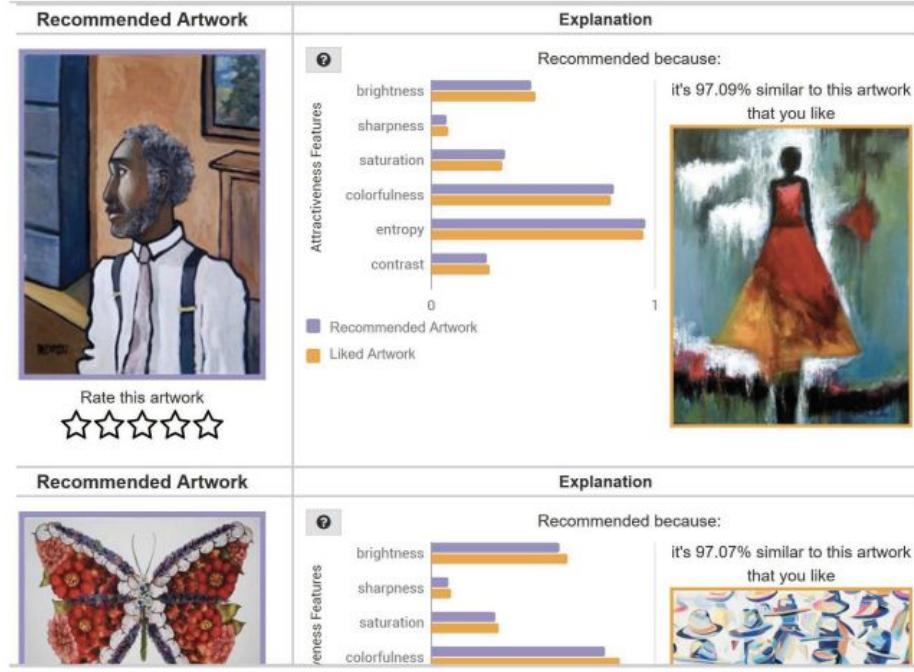


Figure 5: Interface 3: Explainable and transparent recommendation interface with features' bar chart and top-1 similar image.

Resultados de evaluación experta (N=121)

Columnas DNN y AVF son los algoritmos de recomendación basada en contenido para representar imágenes

| Condition | Explainable | | Relevance | | Diverse | | Interface Satisfaction | | Use Again | | Trust | | Average Rating | |
|---|---------------------|--------------------|---------------------|------|--------------------|-------|------------------------|------|-----------|------|-------|------|----------------|------|
| | DNN | AVF | DNN | AVF | DNN | AVF | DNN | AVF | DNN | AVF | DNN | AVF | DNN | AVF |
| Interface 1 (No Explanations) | 66.2* | 51.4 | 69.0* | 53.6 | 46.1 | 69.4* | 69.9 | 62.1 | 65.8 | 59.7 | 69.3 | 63.7 | 3.55* | 3.23 |
| Interface 2 (DNN & AVF: Top-3 similar images) | 83.5*↑ ¹ | 74.0↑ ¹ | 80.0* | 61.7 | 58.8 | 69.9* | 76.6* | 61.7 | 76.1* | 65.9 | 75.9* | 62.7 | 3.67* | 3.00 |
| Interface 3 (DNN: Top-3 similar, AVF: chart) | 84.2*↑ ¹ | 70.4↑ ¹ | 82.3*↑ ¹ | 56.2 | 65.3↑ ¹ | 71.2 | 69.9* | 63.3 | 78.2* | 58.7 | 77.7* | 55.4 | 3.90* | 2.99 |

Stat. sign. between interfaces by multiple t-tests, Bonferroni corr. $\alpha_{bonf} = \alpha/n = 0.05/3 = 0.0017$. Stat. sign. between algorithms using pairwise t-test, $\alpha = 0.05$.

Table 3: NASA TLX Results. The symbol \uparrow^2 indicates interface-wise significant difference (differences between interfaces using the same algorithms).

| Condition | Mental | | Hurry | | Insecure | |
|---|--------|-------|-------|-------|----------|--------------------|
| | DNN | AVF | DNN | AVF | DNN | AVF |
| Interface 1 (No Explanations) | 19.90 | 23.24 | 10.78 | 13.41 | 12.22 | 12.88 |
| Interface 2 (DNN & AVF: Top3 images) | 20.05 | 18.46 | 11.54 | 12.08 | 7.62 | 6.59 |
| Interface 3 (DNN: Top3 imag., AVF: chart) | 23.41 | 26.37 | 14.29 | 15.73 | 13.32 | 16.37 \uparrow^2 |