# WEBIST 2014

## 10th International Conference on Web Information Systems and Technologies
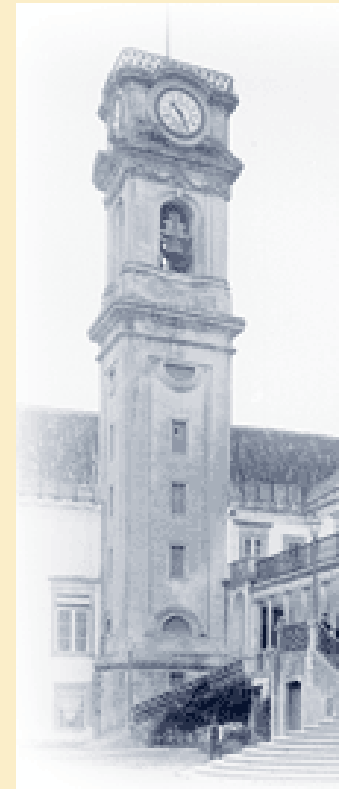
### Barcelona, Spain | 3 - 5 April, 2014
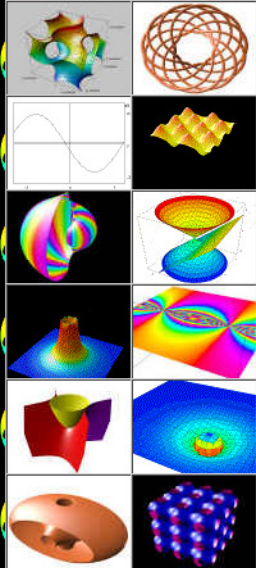
# AUTOMATIC WEB PAGE CLASSIFICATION USING VISUAL CONTENT

António Videira and Nuno Gonçalves

Institute for Systems and Robotics – University of Coimbra, Portugal
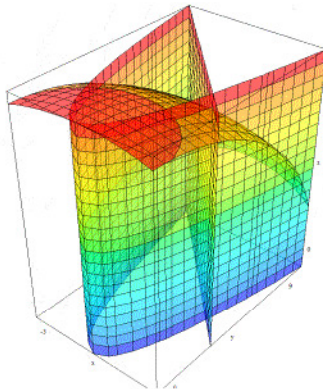
nunogon@isr.uc.pt

# MOTIVATION



"Quality: Excellent. Value: Excellent." ... "DPGraph is one of the most exciting Windows-PC programs I've ever seen for creating beautiful, even stunning, mathematical graphics.", Dr. Michael W. Ecker, Recreational & Educational Computing, *and DPGraph runs under Wine on Linux, or under SoftWindows or Virtual PC on the Mac.*

**DPGraph: Dynamic Photorealistic 3D Graphing Software for Math and Physics Visualization**

Subscribe | List of Site Subscribers | Free Viewer
Latest news: 4 Jul 2008 | Update to newest version: 4 Jul 2008
Math Art Gallery | Documentation | Links | Privacy | Contact

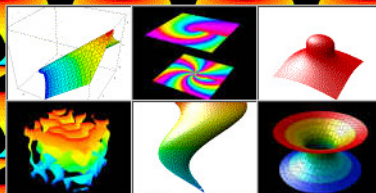★★★★★ "You'll be dazzled ...", ZDNet

The world's most powerful software for math and physics visualization. Create beautiful, interactive, dynamic, photorealistic 2D, 3D, 4D, 5D, 6D, 7D and 8D graphs. So easy to use that even junior high and senior high students have had their graphs published. Includes hundreds of examples contributed by users from around the world.

Over two million mathematicians, physicists, teachers and students at over 1,000 colleges, universities and K-12 schools worldwide are already licensed. Comes with a free subscription to the Flaming Thunder programming language.
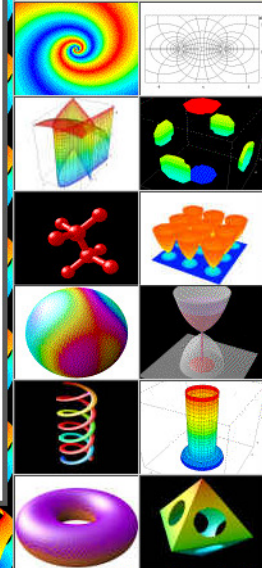
Optimized for the internet -- search for DPGraph using your favorite search engine to find ways that people are using DPGraph for both classroom and distance learning. Used for pre-algebra, geometry, trigonometry and general physics, through multivariable calculus, field theory, quantum mechanics and gravitation.

Use time and color as extra dimensions (to create motion or encode momentum, for example). Use the scrollbar to vary parameters in realtime, to slice through graphs, or to vary transparency. Programmed entirely in assembly language for maximum speed.

Graph functions, equations, conic sections, planes, spheres, toruses, parametric curves and surfaces, implicit equalities and inequalities, volume intersections, volumes of integration, vector fields, surfaces of revolution, equipotential surfaces, and much more, in rectangular, polar, cylindrical, or spherical coordinates.

# MOTIVATION

# MOTIVATION

# MOTIVATION

# MOTIVATION

# MOTIVATION

| Newspaper | Hotel | Celebrities |
|-----------|-------|-------------|

# MOTIVATION

# MOTIVATION

## Newspaper



## Hotel

## Celebrities

# MOTIVATION

News ...brities

# MOTIVATION

## Newspaper

## Hotel

## Celebrities

# MOTIVATION

Newspaper ...rities

# MOTIVATION

## Newspaper





## Hotel



## Celebrities

# OUTLINE

- Related work
- Problem statement
- Our approach:
  - Feature extraction
  - Feature selection
  - Classifiers: training and testing
- Experiments
- Conclusions
- Future work

# RELATED WORK

- Classification of web pages has been traditionally done using <u>only</u> text.
- Text is rich in semantic content.
- Classification is focused on the <u>topic</u>.
- Text sources:
  - html code (tags: keywords, title, description, ...)
  - html code (text in paragraphs or anchors)
  - url (self url or page hyperlinks)

# RELATED WORK

- Text only:
  - Neural networks + PCA [Selamat and Omatu, 2004]
  - SVM + semantic [Chen and Hsieh, 2006]
  - Web crawlers (Google, Yahoo, …)

- Text + images:
  - Structure of page + images [Asirvatham and Ravi, 2001]
  - Visual Adjcency Multigraph (graphs) [Kovacevic et al., 2004].

- Visual only:
  - Several features + classifiers [Boer et al., 2010].

# RELATED WORK

- **Positives** of using text:
  - very fast
  - easy to implement
  - rich content to classify the topic of the page.

- **Negatives** of using text:
  - it has a lot of irrelevant data
  - absolutely no information on some subjective variables

# PROBLEM STATEMENT

- Classification of web pages
- Using <u>only</u> the visual content
- Subjective variables
  - whether a page is beautiful or ugly – aesthetic value
  - whether a page is old fashioned or has newer design – design recency
- Topic of the web page

# VISUAL CONTENT

* <u>Positives</u> of using visual content:
    + rich content that we can not find in text
    + can ben applied to subjective variables
    + there is a growing trend to use images and banners in the design/layout of the page – these images contain visual text that is not caught by crawlers

* <u>Negatives</u> of using visual content:
    + slow

# RICH CONTENT IN IMAGES OF TEXT

# OUR APPROACH

# FEATURE VECTORS

- **1 – Low- level descriptor** (166 dimensions)
  - Color Histogram
  - Edge Histogram
  - Gabor
  - Tamura
- **2 – Mid- level descriptor** (100, 200 or 500 dimensions)
  - SIFT
  - Bag Of Words (BOW)

# COLOR HISTOGRAM

# COLOR HISTOGRAM

# EDGE HISTOGRAM

# EDGE HISTOGRAM

# GABOR FEATURES

# GABOR FEATURES

# TAMURA FEATURES

# TAMURA FEATURES

# SIFT + BoW

- SIFT features – David Lowe [1999, 2003]
- They are:
  - scale and orientation invariant
  - local descriptors
  - slow to compute
  - keypoint descriptor is an orientation histogram
- Bag of Words (BoW) – we use Jialu Liu [2013]
  - All keypoint descriptors are used to build a dictionary by grouping SIFT features in "visual words".
  - A good trade-off between size and accuracy.

# SIFT + BoW



(i) SIFT Keypoint Detection

(ii) SIFT Descriptors

(iii) Quantizing descriptors
k-means clustering

w1

w2

w3

w4

(iv) Bag-of-Words

# SELECTION OF FEATURES

- Features are selected by their discriminative power.

- Chi Square criterion ($\chi^2$)
  - it relies on maximizing the matching between observed and expected frequencies.

- Principal Component Analysis (PCA)
  - it relies on explaining the most variability of the sample

# CLASSIFIERS

- Most representative classifiers were used in a <u>supervised learning</u> scheme:

    - Naïve Bayes
    - SVM
    - Decision Tree
    - AdaBoost

# EXPERIMENTS

- 90 images for each class
- Images are renderings of landing web pages

- 2 binary classifications
  - Aesthetic value (beautiful/ugly)
  - Design recency (new/old fashioned)
- 1 multi-class classification
  - Web page topic (4 classes and 8 classes)

# AESTHETIC VALUE

- This is an inherently <u>subjective variable</u>
- Aesthetic notion is a valuable concept for marketing and psychology.

- Ugly pages are retrieved from blogs and public listings: [Andrade, 2009] and [Shuey, 2013].
- Beautiful pages are retrieved from designer blogs [Crazyleafdesign.com, 2013].

# DESIGN RECENCY

- This is also an inherently <u>subjective variable</u>, however, design recency can be more rationalized.

- Design notion of recency is a highly valuable concept for marketing.

- Pages were retrieved by consulting the most popular pages from today (2013) and for 1999, using internet archives and alexa.com.

# DESIGN RECENCY

# WEB PAGE TOPIC

- Topic is the main focus of web page classification.
- We started with 4 classes (inspired in [Boer et al., 2010]):
  - Newspapers
  - Hotels
  - Celebrities
  - Conferences

# WEB PAGE TOPIC

- We then added 4 additional ones:
  - Classified advertisements
  - Social networks
  - Gaming
  - Video-sharing

- These 8 classes represent a big majority of all web pages. They are, however, not a systematic study on the topic.

# WEB PAGE TOPIC EXAMPLES



**Newspapers** | **Conferences** | **Hotels** | **Celebrities**

# WEB PAGE TOPIC EXAMPLES



**Classifieds**

**Gaming**

**Social Networks**

**Video Sharing**

# RESULTS – AESTHETIC VALUE

| | Naïve Bayes | SVM | Decision Tree | AdaBoost |
|---|---|---|---|---|
| Color Histogram | 65% | 85% | 70% | 85% |
| Low-level (50%) | 75% | 65% | 80% | 80% |
| SIFT+BoW | 80% | 80% | 75% | 95% |

✖ We observed that most of the features with higher discriminative power are collected from Color Histogram.

# RESULTS – DESIGN RECENCY

|  | Naïve Bayes | SVM | Decision Tree | AdaBoost |
| --- | --- | --- | --- | --- |
| Gabor only | 85% | 100% | 95% | 100% |
| Low-level (5%) | 100% | 85% | 90% | 90% |
| SIFT+BoW | 90% | 90% | 90% | 90% |

- We observed that for design recency, Gabor features have the higher discriminative power. Gabor features are intrinsically related to texture and spatial frequency.

# RESULTS – WEB PAGE TOPIC

## 4 CLASSES

|  | Naïve Bayes | SVM | Decision Tree | AdaBoost |
|---|---|---|---|---|
| Low-level | 62,5% | 72,5% | 72,5% | 70% |
| SIFT+BoW | 70% | 75% | 82,5% | 72,5% |

- ✖ SIFT+BoW always improved results.
- ✖ Higher discriminative power is related with Tamura and Gabor features (texture, coarseness, frequency, …).

# RESULTS – WEB PAGE TOPIC

## 4 CLASSES

| Predict\Actual | Newspapers | Conferences | Celebrities | Hotel |
|---|---|---|---|---|
| Newspapers | 10 | 0 | 1 | 0 |
| Conferences | 0 | 9 | 0 | 1 |
| Celebrities | 0 | 1 | 7 | 2 |
| Hotel | 0 | 0 | 2 | 7 |

# RESULTS WEB PAGE TOPIC

## 8 CLASSES

|  | Naïve Bayes | SVM | Decision Tree | AdaBoost |
|---|---|---|---|---|
| SIFT+BoW | 64% | 59% | 49% | 39% |

| Prd\Act | Newsp. | Conf. | Celeb. | Hotel | Classif. | Gaming | Social | Video |
|---|---|---|---|---|---|---|---|---|
| Newsp. | 9 | 1 | 1 | 1 | 4 | 0 | 1 | 2 |
| Conf. | 1 | 8 | 0 | 0 | 0 | 0 | 0 | 0 |
| Celeb. | 0 | 0 | 4 | 2 | 0 | 3 | 2 | 1 |
| Hotel | 0 | 0 | 0 | 7 | 0 | 1 | 1 | 1 |
| Classif. | 0 | 0 | 1 | 0 | 6 | 1 | 0 | 0 |
| Gaming | 0 | 0 | 4 | 0 | 0 | 5 | 2 | 0 |
| Social | 0 | 1 | 0 | 0 | 0 | 0 | 2 | 0 |
| Video | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 6 |

# MAIN CONCLUSIONS

- Subjective variables are simpler to classify and features like color and edges have high discriminative power.
  - Beautiful web pages tend to have more soft colors
  - New fashioned web pages tend to have lower visual frequency and much lower density on information

- SIFT+BOW features are slower to compute and add only a small amount of discrimination.

# MAIN CONCLUSIONS

- Web page topic is harder to classify using only visual content

- However, accuracies of 80% are generally achieved.

- SIFT+BoW features are able to enhance significantly the accuracy.

- A bigger database is, however, needed to achieve better results in classification of the topic.

# CONCLUSION

- Visual content of web pages effectively has rich content to the task of classification in several variables: aesthetic value, design recency and web page topic.

- Classification using Visual content can be added to the traditional classification relying only on text to achieve a much <u>powerful crawler</u>, boosting its accuracy.

# FUTURE DIRECTIONS

- Integrate visual content with text (html, url, ...) to achieve that powerful classifier.

- Explore the semantic content of the visual content, by segmenting the page and analyzing items separately (images, banners, flash, advertisement, ...).

- Analyze group differences of subjective variables and their implications on marketing (gender, age, culture, income rank, ...).

# APPLICATIONS

✖ Building of an <u>advice system</u> to assist in the design of web pages.

✖ <u>Marketing</u>: targeting of consumers, campaign design, ...

✖ Content filtering and <u>content scoring</u>, including ranking of web sites based on <u>user profile</u>.

✖ Recommender systems for product web pages.

✖ Sentiment and emotional implications of web pages.

# THANK YOU

- Some questions?

- I am searching for partners for <u>H2020</u>!

- Contact: nunogon@deec.uc.pt