# CS395T - DEEP LEARNING SEMINAR

Philipp Krähenbühl

# OVERVIEW

Philipp Krähenbühl
office GDC 4.824
office hours by
appointment (send email)

**TA** Huihuang Zheng
office GDC 6.802
office hours We 10-11am

Try piazza first!

# OVERVIEW

https://philkr.github.io/CS395T/

# OVERVIEW

before class

Read and review
two papers

is covariate shift
important?

how about this baseline?

why not …?

in class

10 min intro (Philipp)

10-15 min paper 1
(one of you)

25 min paper 1
(all of you)

10-15 min paper 2
(one of you)

25 min paper 2
(all of you)

# PRESENTATIONS

no incredibly long walls of text or math than nobody can understand, follow or otherwise parse.

visual

show presentation
to Philipp 1 week ahead

no scrolling
through paper

# OVERVIEW

- two deep learning projects

  - train a deep network

  - write a latex report

  - present your work

# PROJECTS

- two deep learning projects

  - Teams up to 3

  - **GPU access**

  - Python preferred

    - Caffe, TF or Theano

# PROJECT 1

1960s

# PROJECT 2

- Open ended

  - can be your research

  - **not** published by Dec 31

- level of a top tier workshop publication

  - CVPR, ICCV, ICML, NIPS, ACL, SIGGRAPH

  - SIGBOVIK



SIGBOVIK, APRIL 2015

## Visually Identifying Rank

David F. Fouhey, *Mathematicians Hate Him!*
Daniel Maturana, *Random Forester* Rufus von Woofles, *Good Boy*

**Abstract**—The visual estimation of the rank of a matrix has eluded researchers across a myriad of disciplines many years. In this paper, we demonstrate the successful visual estimation of a matrix's rank by treating it as a classification problem. When tested on a dataset of tens-of-thousands of colormapped matrices of varying ranks, we not only achieve state-of-the-art performance, but also distressingly high performance on an absolute basis.

**Index Terms**—perceptual organization; vitamin and rank deficiencies; egalitarianism in the positive-semi-definite cone; PAC bounds for SVDs; class-conscious norms

Fig. 1. What are the ranks of these matrices? Which ones are rank-deficient? In this paper, we investigate how one can guesstimate the rank of a matrix from visual features alone. See footnote on page 2 for answer.

that require access to the matrix, our work gives guarantee-free solutions that can operate on only an colormapped version of a matrix. By treating matrix rank as an image classification problem, we are able to consistently achieve distressingly high performance – ≈ 40% accuracy on 10-way classification; ≈ 80% accuracy on rank-deficient/not-rank-deficient binary classification. In subsequent experiments we show the following: 1) Our method can identify what matrices seem low rank, and why; 2) Our method is easily extended to structured prediction; 3) That activations of our network can be even used as a feature for semantic image classification with non-embarrassing performance (20.9% on Caltech 101 with 15 samples).

## 1 INTRODUCTION

Consider Figure 1(b): what is the rank of the matrix?

# PREREQUISITES

- 391L - Intro Machine learning (or equivalent)

- 311 or 311H - Discrete math for computer science (or equivalent)



- Proficiency in Python

- Basic **deep learning background**

# GOALS

- Review a deep learning paper

- Give an interesting DL presentation

- Devise and execute a DL project

# GRADES

- 30% paper presentation

- 30% project 1 (10% presentation, 20% project)

- 40% project 2 (10% presentation, 30% project)

- (optional) 12.5% volunteering for second presentation

# GRADES

```python
def grade(p):
  from math import floor
  if p < 50: return 'F'
  v = (100-p) * 4 / (50 + 1e-5)
  return chr(ord('A')+floor(v)) +
         ['+','','','-'][floor((v-floor(v))*4)]
```

we might train a deep network to grade instead

# AUDITING

- Allowed

  - No homework or presentation required

  - Paper review and discussion required

# TOPICS

week 2

week 3

week 4

# TOPICS



week 5

week 6

week 7     Project 1 presentations

# TOPICS

week 8



Single RGB Image → Depth Map

week 9



week 10

# TOPICS

week 11



week 12



week 13

# TOPICS

week 14                      TPD


week 15          Project 2 presentations

# THE N-WORD



- Neural

  - try to keep Neuroscience out of this class

  - try to motivate through optimization and ML

    - instead of biology

# REVIEW

Data

Model

Output

# TRAINING

Data

Data

Data

Data

Model →

?

Output

Output

Output

Output

# TESTING / INFERENCE

Data

Model ➡

?

# DATA

feature

Data

$f_1 = \{a,b,c,\ldots\}$

$f_2 = \{d,e,f,\ldots\}$

…

# EXAMPLE: GRADING

Data

| feature | score Project 1 | score Project 2 | age | height |
|---|---|---|---|---|
| $f_1$ | 51 | 44 | 23 | 182 |
| $f_2$ | 25 | 80 | 26 | 172 |

# EXAMPLE: GRADING

f₁   51   44   23   182       A+

**Model** →

f₂   25   80   26   172       A-

continuous and (differentiable)

# EXAMPLE: GRADING

f₁   51   44   23   182   → Model →   |

A   B   C   D   F

**regression**:
predict continuous value and round

linear regression
g = A f + b

# EXAMPLE: GRADING

$f_1$    51    44    23    182    Model

A    B    C    D    F
40%  30%  10%  20%  0%

**classification**:
predict continuous distribution over grades

logistic regression
$P_g = \text{softmax}(A\,f + b)$

$\text{softmax}(x) = \exp(x) / (\Sigma \exp(x))$

# TRAINING

$f_1$     51   44   23   182                                 A+

**Model** ➡

$f_2$     25   80   26   172                                 A-

?

# EXAMPLE: GRADING

**regression**
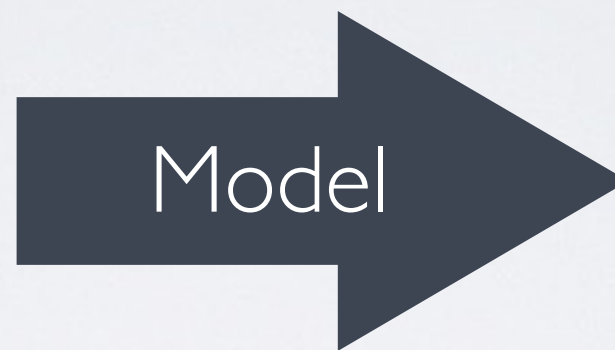
minimize $\Sigma_i(g_i - A\,f_i - b)^2$
  A,b

**classification**

maximize $\Sigma_i \log P(g_i)$
  A,b

# FAIRNESS

$$f_1 > f_2 \quad \Rightarrow \quad g_1 \geq g_2$$

# EXAMPLE: GRADING

**regression**

$$\text{minimize } \sum_i (g_i - A f_i - b)^2$$
$$A,b \quad + \sum_{i,j:f_i<f_j} \max(Af_i - Af_j, 0)$$

**classification**

$$\text{maximize } \sum_i \log P(g_i)$$
$$A,b \quad + \sum_{i,j:f_i<f_j} \max(Af_i - Af_j, 0)$$

# EXAMPLE: GRADING

**regression**

minimize $(g_i - A f_i - b)^2$
        A,b


A $\geq$ 0

**classification**

maximize log $P(g_i)$
        A,b


$A_0 \geq A_1$

# NEXT CLASS

- Look at list of papers

  - Send Huihuang your top picks per email

    - instructions on Piazza