

# MATH 462

---

- problem defn

False Verification

- Important Fns.  
similarity  $S^n, \mathbb{R}^n$

- big project

---

Lecture 16

---

3.11.2024

---



## Face Verification Problem (FV)

Given two images of { same object (person's face)  
different object

Decide (Binary class)

if they represent same +1  
diff -1.

[ Related classification K people  
Given image  $x$ , which person is it?  
K-classification ]

## Two types of generalization

(G1)  $m_1$  photos of  $K$  different faces (balanced)  
Training Set Train class-model.  
 $m_2$  photos of same  $K$  faces. acc.  $p_{\text{train}}$   
New photos what is the accuracy?  
 $p_{\text{test}}$

Gap  $p_{\text{train}} - p_{\text{test}}$  generalization gap

Ex  $.97 - .85 = .12$

model overfits if  $\text{gap} > 0$ .

Better model overfits less

---

$$.91 - .89 = .02$$

problem

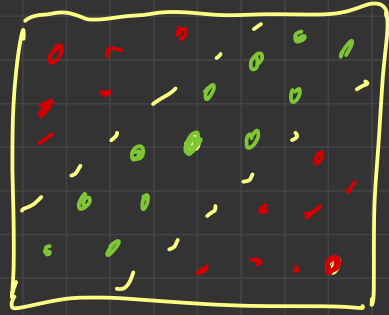
F.V.

$x$  image of a face

$p = (x^1, x^2)$  pair of images  $x$

$$y \in Y_{\pm} = \{-1, +1\}$$

$$y(p) = \begin{cases} +1 & \text{similar pair (same face)} \\ -1 & \text{otherwise} \end{cases}$$



$x$

- diff picture of same face
- picture of diff. faces

$$S_m = \{(p_1, y_1), \dots, (p_m, y_m)\}$$

$$= \{(x_1^1, x_1^2, y_1), \dots, (x_m^1, x_m^2, y_m)\}$$

GOAL Given  $p = (x^1, x^2)$  classify  $y$  as similar or not

G1  $S_{\text{train}}$  most of pairs

$S_{\text{test}}$ .

train model using  $S_{\text{train}}$  accuracy  $P_{\text{train}}$

test model using  $S_{\text{test}}$  " "  $P_{\text{test}}$

if  $P_{\text{train}} - P_{\text{test}}$  too large. model overfits.

"Regularization" Deep Learning Book  
see chapter

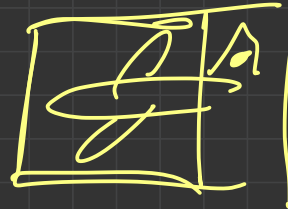
→ Data Augmentation

$x$  replace with  $\tilde{x} = T_i(x)$  for some  
transformation.

$T_i$  = add noise, crop, smooth,

change  $x$ , same image.

EG "happy features"



G-2 trained on  $P_{celeb}$ . out of dist gen.  
test on reserved images from  
 $P_{celeb}$ .

model

FV on other faces not in  $P_{celeb}$

- New objects. to verify  
never been seen before.
  - still fails.
  - not e.g. recognize dog breeds.
- what do the objects have in common?

# In distribution generalization

G1  $p_{\text{data}}(x)$   $S_m$  samples i.i.d.  
independent, identically distributed.

generalization  
accuracy  $p_{\text{train}}$  v.s.  $p_{\text{test}}$

when  $S_{\text{train}}$  &  $S_{\text{test}}$   
each consisting of samples drawn i.i.d.  
from same  $p_{\text{data}}$ .

G2 out of distribution generalization.

partially defined

$h(x)$

trained on  $p_{\text{data}}(x)$

$\sum_m$  i.i.d  $x$  from  $p_{\text{data}}(x)$

$\Rightarrow \{(x_i, y_i)\}$

$y_i = y(x_i)$   
true label

same starting point.

E.g. classification.

Keep unknown the label for  $y(x)$ .

Change  $p_{\text{data}}(x)$  to  $p_{\text{new}}(x)$ .

this part is not defined.!!!



G2 for  $\$V$ :

new facts! people facts not in  $P_{train}$

Works! Great!

So ood generalization is possible.

Research problems define it in such a way  
cover cases works, not too broad.

→ Hard ood gen.

→ easy ood gen. (some confused)  
sequel

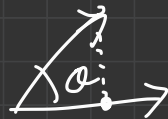
Background.

$$\mathbb{S}^d = \{v \in \mathbb{R}^d \mid \|v\|=1\}$$

similarity

① Given  $v_1, v_2 \in \mathbb{S}^d$   $v_1 \cdot v_2 = \cos \theta$

$$\text{sim}(v_1, v_2) = \cos \theta$$



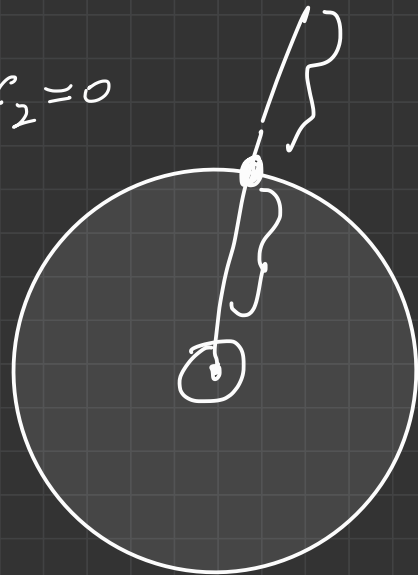
② Given  $f_1, f_2 \in \mathbb{R}^d$   
 $\cos \theta$  undefined if  $f_1 \cdot f_2 = 0$

Defn  $\text{sim}_\varepsilon : \mathbb{R}^d \times \mathbb{R}^d \rightarrow [-1, 1]$

$$N_\varepsilon : \mathbb{R}^d \rightarrow \mathbb{R}^d$$

$$N_\varepsilon(f) = \begin{cases} \frac{f}{\|f\|} & \text{if } \|f\| \geq \varepsilon \\ 0 & \text{o.w.} \end{cases}$$

$$\text{sim}_\varepsilon(f_1, f_2) = N_\varepsilon(f_1) \cdot N_\varepsilon(f_2)$$



$$\text{sim}_{\epsilon}(f_1, f_2) = \begin{cases} \frac{f_1 \cdot f_2}{\|f_1\| \|f_2\|} & \text{if } \|f_1\| \geq \epsilon \text{ and } \|f_2\| \geq \epsilon \\ 0 & \text{if either of } \|f_1\| < \epsilon \text{ or } \|f_2\| < \epsilon \end{cases}$$

Method feature representation Learning

$$\mathcal{H} = \{ h(x) = s \mid h(x) = \text{sim}_{\epsilon}(f(x^1, w), f(x^2, w)) - t \}$$

$f(x, w)$  DNN

$s \rightarrow s - t$  threshold score classifier.

$f: \underset{\mathbb{R}^d}{X} \rightarrow \mathbb{R}^k$  features

$$p = (x^1, x^2)$$

$$x^1 \rightarrow f_1 = f(x^1, w)$$

$$x^2 \rightarrow f_2 = f(x^2, w)$$

$$s = \text{sim}_v(f_1, f_2)$$

$$h = s - t.$$

$$c(h) = \text{sgn}(h).$$

Binary class same as before

$$l_{\text{class}}(h, y) = l_{\text{log}}(h, y) \text{ or } l_{\text{margin}}(h, y)$$

Loss

$$\hat{L}(w, t) = \frac{1}{m} \sum_{i=1}^m l_{\text{class}}(h(p_i), y_i)$$

# Generalization

<u>GI</u>	K-classification	FV
shallow	proof ✓	✓
CNN	X	X

Bound on  
Gen gap:  
with high prob.  
 $\text{gap} \leq O(\frac{1}{\sqrt{m}})$

G 2	K-classification	FV
end to end hlp)	X <del>proof</del> practical	X
$f(x, w)$ Prun. [new threshold.] $t^*$ $p_{\text{new}}$	✓ possible	✓

train  $f(x, w)$  on  $p_{\text{train}}$   
using  $\hat{L}(w, t)$

strip of the final layer (score classifier).

New dataset  $p_{\text{new}}$  different faces.

① fine tuning. where allow  $w$  to change  
no proof. it can overfit.

② just change threshold.  
i.e. only  $t$  will generalize.

$$f(x).$$

model  $h(w, x) = w \cdot f(x)$

$\hat{\mathcal{L}}(w)$  (train loss.  
 $\text{train acc}(h) \leq \hat{\mathcal{L}}(h(w, x))$ )

Thm

$$\mathcal{L}(w) = \mathbb{E}_{x \sim \rho} \mathcal{L}(h(w, x), y) \text{ in theory}$$

$$\hat{\mathcal{L}}_{\text{test}}(w) \text{ or } \frac{1}{m} \sum_{i=1}^m \mathcal{L}(h(w, x_i), y_i) \text{ test loss.}$$

Thm

$$|\hat{\mathcal{L}}_{\text{test}} - \hat{\mathcal{L}}_{\text{train}}| \leq \frac{C_\varepsilon}{\sqrt{m}} + \mathcal{R}_m(\mathcal{H})$$

with prob  $\geq 1 - \varepsilon$

$f(x, w_0)$  ImageNet feats.

Apply New data of vegetables.

New classifier

$$h(x, \tilde{w}) = f(x, w_0) \cdot w$$

freeze  $w_0$ , train  $w$ .

Acc Train, 93

Bound Test on veg.

$$93 - \frac{c_\epsilon}{\sqrt{m}}$$

with high prob.

Allow  $w_0$  to train.

No Bound Acc Train .99

≠ enough train sample