Lecture II

Neyman–Pearson Lemma

$$L = -\left( \sum_{i \in S} \log f(x_i) + \sum_{i \in B} \log (1 - f(x_i)) \right)$$

what $f$ minimizes $L$ in ideal case?
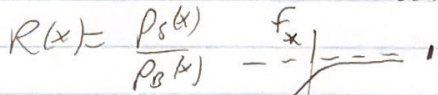
prob density of sg & bg

continuum limit ($\infty$ data)

$$L = -\int dx \left[ p_S(x) \log f(x) + p_B(x) \log (1 - f(x)) \right]$$

$\int dx\, p(x) \approx \sum_{x \sim p(x)}$    $f \to f + \delta f$, $\left. \dfrac{\partial L}{\partial \delta f} \right|_{f = f_*} = 0$    $\leftarrow$ optimal $f$

(assume $N_S = N_B$ WLOG)

$$0 = \frac{p_S(x)}{f_*} - \frac{p_B(x)}{1 - f_*} \implies \boxed{f_* = \frac{p_S(x)}{p_S(x) + p_B(x)}}$$

So optimal classifier ⊘ $f_* = \dfrac{R(x)}{R(x) + 1}$

$$R(x) = \frac{p_S(x)}{p_B(x)}$$

$f_*$ is monotonic w/ $R(x)$

$\implies \boxed{f_* > f_c \iff R(x) > R_c}$   cut on $f$ equiv to cut on $R$

So optimal classifier is likelihood ratio

"Neyman–Pearson Lemma"

- Powerful result! fundamental result in statistics
  (LR uniformly most powerful test
  for simple hypothesis)

later will see:
- Can turn around → learn $\overset{ML}{\wedge}$ classifier from samples
  → approach LR.
  "likelihood ratio trick"

- Another perspective: Bayes Thm

$$f_* = \frac{P_S(x)}{P_S(x) + P_B(x)} = \frac{p(x|s)}{p(x|s) + p(x|B)} = \frac{p(x|s)\, p(s)^{\frac{1}{2}}}{p(x)}$$

$$= p(s|x) \qquad \frac{= p(x|s) + p(x|0)}{2}$$

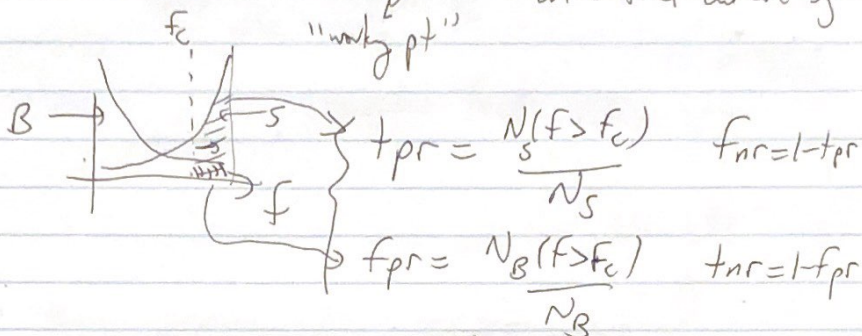So $f_*$ is the prob of signal given x which is where we started!

i.e. BCE is minimized when our model for this prob ✓
= true prob.

_____

Back to NP lemma: what does "most powerful" test mean?
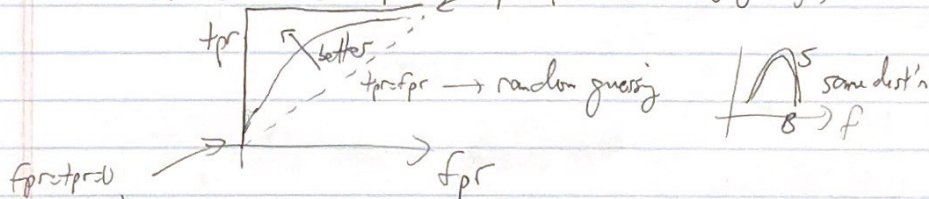→ metrics for binary classification

How do we use binary classifier?

Usually: cut $f(x) > f_c$    all x satisfy cut are "signal" <sub>classified as</sub>

all x fail cut are "bg"

$f_c$    "working pt"



$tpr = \dfrac{N_s(f > f_c)}{N_s}$    $fnr = 1 - tpr$

$fpr = \dfrac{N_B(f > f_c)}{N_B}$    $tnr = 1 - fpr$

- accuracy $= \max_{f_c} \dfrac{tpr + (1 - fpr)}{2}$

- ROC curve   $tpr = 1$  ← $fpr = tpr = 1$ (let everything through)



$tpr = fpr$ → random guessing

same dist'n

$fpr = tpr = 0$
reject everything    Vary $f_c$, sweep out "ROC curve" $(fpr(f_c), tpr(f_c))$

$\longrightarrow$ AUC "Area under the curve"

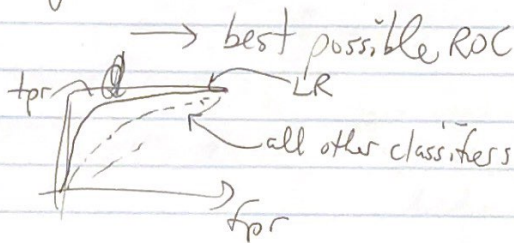$\begin{cases} AUC = 1 \text{ perfect} \\ AUC = 0.5 \text{ random} \end{cases}$    (same for acc but AUC $\neq$ ACC)

- AUC, ACC mostly sensitive to $tpr, fpr \sim O(1)$ part of ROC curve

When $N_{bg} \gg N_{sig}$ (as in many cases, eg LHC)

care more about fpr << 1 ~~mostly~~ part of ROC curve

→ often report fpr @ fixed tpr eg tpr=50%

or 30%

or rejection factor $\begin{cases} R50 = \dfrac{1}{fpr \, @ \, tpr=50\%} \\ \\ R30 \quad etc \end{cases}$

• Neyman-Pearson LR classifier is optimal

→ best possible ROC curve



tpr — LR

all other classifiers

fpr

• In practice cannot achieve NP optimality

— exact likelihoods unknown

— finite training data → Bera of big data

— limited model capacity (expressivity)

↳ NNs very ~~expressly~~ expressive

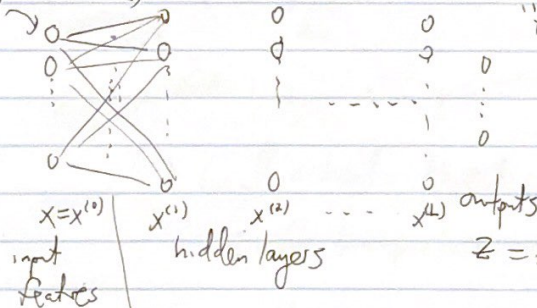⟹ can ~~maybe~~ be set very close to optimal!

# Neural Networks

- So far have not specified family of fit fns $f(x; \theta)$

- Many choices (BDTs, SVMs, Boltzmann machines, ...)

  $\longrightarrow$ deep learning / "modern ML" : $f(x; \theta) = NN$
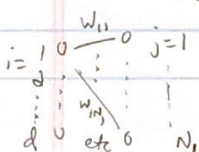
  why? — universal approx thm : roughly, "NN can approximate any fn"

  — differentiable (extremely large models can be trained on a lot of data using backprop & SGD)

  not well understood, but surprisingly resistant to overfitting — "generalization puzzle" "inductive bias"

- Basic structure of NN    "feed-forward"
  "MLP" "DNN"

"nodes" or "neurons"      "fully connected"



$x = x^{(0)}$    $x^{(1)}$    $x^{(2)}$   $\cdots$   $x^{(L)}$ outputs

input features     hidden layers        $z = F(x; \theta)$

$i = 1$   $w_{11}$   $j = 1$

$\vdots$    $w_{1N_1}$

$d$   etc   $N_1$

each connection: "weight"
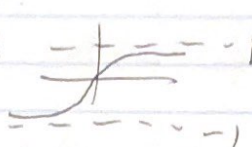all together: weight matrix $w_{ij}$

"biases"

$x^{(1)} = A^{(1)}\left(w^{(1)} x^{(0)} + b^{(1)}\right)$
     "activation"

- activation: only source of non linearity.

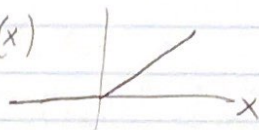    examples: $A(x) = \dfrac{e^x}{1+e^x}$   "sigmoid"



$$= \tanh x$$



best choice!
others lead
to "vanishy gradient
problem" → more later

$$= ReLU(x)$$



$$= \begin{cases} x & \text{if } x > 0 \\ 0 & \text{if } x \le 0 \end{cases}$$

- Note: activation acts element wise

$$A^{(i)}\begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{pmatrix} = \begin{pmatrix} A^{(i)}(x_1) \\ \vdots \\ A^{(i)}(x_n) \end{pmatrix}$$

- Structure of NN is recursive

$$x^{(n)} = A^{(n)}\left( w^{(n)} x^{(n-1)} + b^{(n)} \right)$$

$$\underset{z = \, \overset{x^{(L+1)}}{\overbrace{\phantom{xx}}}}{} A^{(L+1)}\left( w^{(L+1)} x^{(L)} + b^{(L+1)} \right)$$

    ↳ final activation depends on problem

    eg for binary classification $z = \begin{pmatrix} p(5|x) \\ p(0|x) = 1 - p(5|x) \end{pmatrix}$

    want outputs to sum to 1

    popular choice: soft max $\left( \dfrac{e^{x_1}}{\sum_{j=1} e^{x_j}}, \dfrac{e^{x_2}}{\sum e^{x_j}} \right)$