

Friday, February 5, 2021 3:04 PM

Lecture 6Vanishing Gradient Problem

$$f(x) \equiv A(w_h - A(w_{h-1} - A(w_{h-2} - \dots - A(w_1 x) \dots))$$

$$\frac{\partial f}{\partial w_k} \equiv G_k x_k$$

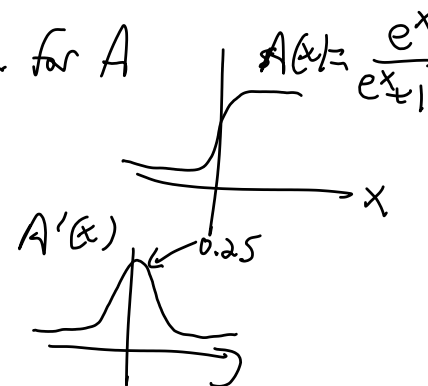
$$x_k \equiv A(w_{k-1} x_{k-1})$$

$$G_{k-1} \equiv G_k w_k A'(w_{k-1} x_{k-1})$$

In early days of NN's, common to use sigmoid fn. for A

weights getting stuck far from optimum

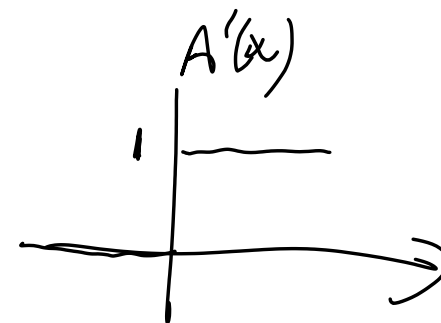
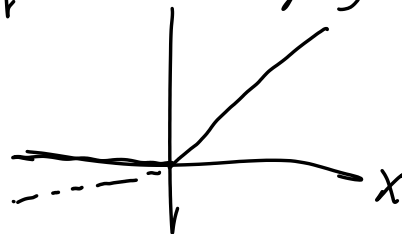
decreasing or suppressed gradients esp. for initial layers. as network gets deeper.



Friday, February 5, 2021 3:32 PM

Sol'n: use a better activation fn

e.g. ReLU: $\max(0, x)$



enabled much deeper networks

Friday, February 5, 2021 3:37 PM

Back to binary classification

Key result in binary classification: Neyman-Pearson Lemma: optimal binary classifier is the Likelihood Ratio

$$R(x) = \frac{P(x|S)}{P(x|B)}$$

Proof! consider BCE loss

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

$$= \int dx \left[P(x|S) \log f(x) + P(x|B) \log(1 - f(x)) \right]$$

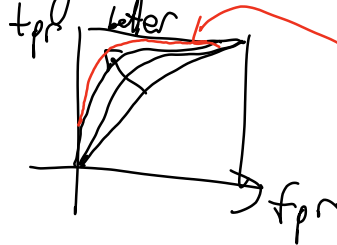
$f \rightarrow f + \delta f$, set 1st der. to zero

$$0 = \frac{P(x|S)}{f} - \frac{P(x|B)}{1-f}$$

$$\rightarrow f = \frac{P(x|S)}{P(x|S) + P(x|B)}$$

$$= \frac{R(x)}{R(x) + 1} \quad f > c \Leftrightarrow R > c'$$

uniformly best ROC curve.



"best" everywhere.

at fixed tpr $R(x)$ achieves lowest possible fpr for every tpr.

($f(x) = P(S|x)$?)
 Bayes thm: $P(S|x)P(x) = P(x|S)P(S)$

Friday, February 5, 2021 3:49 PM

f is monotonic in $R \rightarrow$ equivalent classifiers ✓
 LR minimizes BCE loss ✓

Proof #2: (directly showing that R gives best possible ROC curve)

Consider any other classifier $f(x)$. $f(x) > c \rightarrow$ defines a region C in x space

$$\downarrow$$

$$tpr = \epsilon_S(f) = \int_C P(x|S)$$



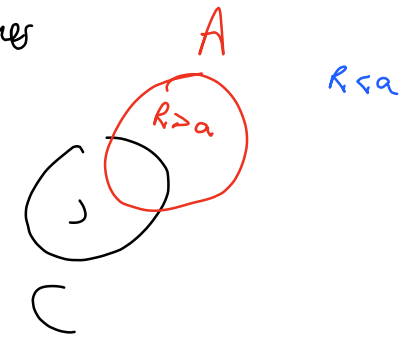
Now consider a cut on $R(x) > a$ w/ same tpr. This defines another region A

$$\epsilon_S(f) = \epsilon_S(R) = \int_A P(x|S)$$

Friday, February 5, 2021 3:56 PM

We want to compare $\epsilon_B(f)$ ✓ $\epsilon_B(R)$. (background efficiencies or fprs)

$$\int_C P(x|B) \quad \int_C P(x|B)$$



$$\epsilon_B(R) = \int_A P(x|B) = \int_{A \cap C} P(x|B) + \int_{A \cap C^c} P(x|B)$$

$$= \int_C P(x|B) - \int_{A^c \cap C} P(x|B) + \int_{A \cap C^c} P(x|B)$$

$$= \epsilon_B(f) - \int_{A^c \cap C} \frac{1}{R(x)} P(x|B) + \int_{A \cap C^c} \frac{1}{R(x)} P(x|B)$$

$$\begin{aligned}
 &< \epsilon_B(f) + \frac{1}{a} \left[\int_{A \cap C^c} P(x|B) - \int_{A^c \cap C} P(x|B) \right] \\
 &\quad + \frac{1}{a} \left[\int_{A \cap C} P(x|B) - \int_{A \cap C} P(x|B) \right] \\
 &= \epsilon_B(f) + \frac{1}{a} \int_{A \cap C^c} P(x|B) - \frac{1}{a} \int_{A^c \cap C} P(x|B) \\
 &= \epsilon_B(f) \quad \checkmark
 \end{aligned}$$

$$(R(x) = \frac{P(x|B)}{P(x|A)})$$

$\epsilon_B!$

Friday, February 5, 2021 4:02 PM

NP lemma is extremely useful!

- optimal classifier exists, given LR \rightarrow provides strategy for hypothesis testing, ^{e.g.} search for new physics
construct likelihood for signal & background.
- modern ML: "likelihood ratio trick"

Assume: NN can learn the optimal classifier $\rightarrow \frac{P(x|S)}{P(x|B)}$
 "likelihood free inference" \leftarrow gives access to likelihood ratio w/out knowing indiv. likelihoods!

Applications

- generation: GANs
- anomaly detection: "weak supervision", CWoLa
- phase space reweighting e.g. simulation vs data
- ⋮



Friday, February 5, 2021 3:56 PM

Multi class Classification

- Important & straightforward generalization of binary classification

- Objective: learn class probabilities $p(C_k | x)$ $k=1, \dots, N_{\text{class}}$

$$\sum_k p(C_k | x) = 1.$$

Ex: MNIST
classify 0-9
from each other.

- loss fn? MLE binary: $\sum_{i \in S} \log p(S | x_i) + \sum_{i \in B} \log p(B | x_i)$

multiclass:
$$L = \sum_k \sum_{i \in C_k} \log p(C_k | x_i)$$

"categorical cross entropy"

Friday, February 5, 2021 4:27 PM

Think of $p(C_k|x)$ $\xrightarrow{\text{vector}}$ $\begin{pmatrix} p(C_1|x) \\ \vdots \\ p(C_N|x) \end{pmatrix} = \vec{p}(x)$

introduce truth labels $k=1 \begin{pmatrix} 1 \\ \vdots \\ 0 \end{pmatrix} \leftarrow y \text{ for } k=1$, $k=2 \begin{pmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{pmatrix} \leftarrow y \text{ for } k=2$, ... $k=N \begin{pmatrix} 0 \\ \vdots \\ 1 \end{pmatrix} \leftarrow y \text{ for } k=N$
 "one-hot encoding"

$$L = \sum_{i \in \text{data}} \vec{y}_i \cdot \log \vec{p}(x_i)$$

\vec{y}_i truth label for data pt. i

\log acts elementwise over vector \vec{p} .

Friday, February 5, 2021 4:33 PM

Metrics:

- Accuracy: fraction of $\arg \max \vec{p}(x) = \text{true class}$ "top-1 accuracy"
- "top-n accuracy" \rightarrow correct label in n highest \vec{p} 's.
- one-vs-rest binary metrics
- confusion matrix

	true			
pred	C_1	C_2	C_3	...
C_1				
C_2				
C_3				
...				
...				