

$$\begin{aligned}
 & \text{Entropy of Gaussian} \\
 H(x) &= -\int_{-\infty}^{\infty} p(x) \log[p(x)] dx \\
 p(x) &= \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \\
 \log p(x) &= \log\left(\frac{1}{\sqrt{2\pi}\sigma}\right) - \frac{(x-\mu)^2}{2\sigma^2} \\
 \Rightarrow H(x) &= -\int_{-\infty}^{\infty} p(x) \left[\log\left(\frac{1}{\sqrt{2\pi}\sigma}\right) - \frac{(x-\mu)^2}{2\sigma^2} \right] dx \\
 &= -\log\left(\frac{1}{\sqrt{2\pi}\sigma}\right) \int_{-\infty}^{\infty} p(x) dx + \frac{1}{2\sigma^2} \int_{-\infty}^{\infty} (x-\mu)^2 p(x) dx \\
 &= -\left[\frac{1}{2} \log(2\pi) + (-1) \log \sigma\right] + \frac{1}{2\sigma^2} \int_{-\infty}^{\infty} (x-\mu)^2 p(x) dx \\
 &= \frac{1}{2} \log(2\pi) + \log \sigma + \frac{1}{2\sigma^2} \int_{-\infty}^{\infty} (x-\mu)^2 p(x) dx \\
 &= \frac{1}{2} \log(2\pi) + \log \sigma + \frac{1}{2\sigma^2} \\
 &= \frac{1}{2} \log(2\pi) + \log \sigma + \frac{1}{2} \log e
 \end{aligned}$$

$$H(x) = \frac{1}{2} (\log(2\pi) + 2\log \sigma + \log e)$$

$$H(x) = \frac{1}{2} \log(2\pi e^{\sigma^2})$$

REGRESSION LOSS

① Mean Squared Error

$$\text{MSE} = \left(\frac{1}{n}\right) \sum_{i=0}^n (\hat{y}_i - y_i)^2$$

sensitive to outliers

② Mean Absolute Error

$$\text{MAE} = \frac{1}{n} \sum_{i=0}^n |\hat{y}_i - y_i|$$

Not sensitive to outliers, not differentiable at zero

③ Huber Loss \rightarrow MLE + MAE.

δ : Huber's threshold. Linear for large values, quadratic for smaller ones

$$L = \begin{cases} \frac{1}{2} (y - \hat{y})^2 & \text{if } |y - \hat{y}| \leq \delta \\ \delta |y - \hat{y}| - \frac{1}{2} \delta^2 & \text{if } |y - \hat{y}| > \delta \end{cases}$$

Small $\delta \rightarrow$ MAE

Large $\delta \rightarrow$ MSE

Problem: Threshold effect; not smooth

④ Log-Cosh Loss

Smaller version of Huber's loss

$$L = \sum_{i=1}^n \log(\cosh(\hat{y}_i - y_i))$$

\rightarrow No critical threshold

\rightarrow Smooth diff function that handles outliers better than MSE

$$\cosh(x) = \frac{e^x + e^{-x}}{2}$$

BINARY CLASSIFICATION LOSS

① Avg. Zero-One Loss

$$L_{0-1}(y_i, \hat{y}_i) = \begin{cases} 0 & \hat{y}_i = y_i \\ 1 & \hat{y}_i \neq y_i \end{cases}$$

$$L_{0-1} = \frac{1}{n} \sum_{i=1}^n L_{0-1}(\hat{y}_i, y_i)$$

\rightarrow non-differentiable, non-continuous, no gradients
 \rightarrow no margin info \rightarrow treats all misclassifications equally

\rightarrow can help determine optimum decision threshold

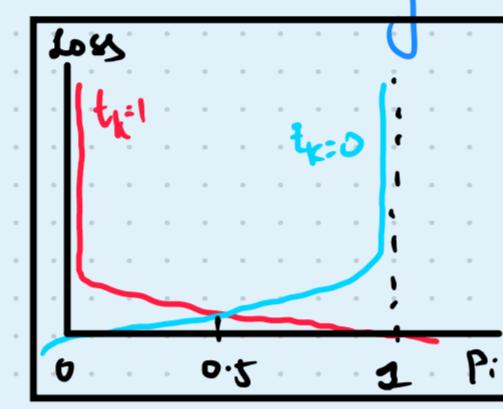
② Log Loss / Binary cross entropy / Logistic

$$L_i = -y_i \log p_i - (1-y_i) \log(1-p_i)$$

\rightarrow lower = better
 \rightarrow penalises harshly when predicted probability strays

③ Brier Score

$$BS = \frac{1}{N} \sum_{i=1}^N (t_{y_i} - \hat{y}_i)^2$$



Much gentler penalty

MULTICLASS CLASSIFICATION LOSS

① Gross Entropy Loss

$$L_{ce} = -\sum_{i=1}^n y_i \log(\hat{p}_i)$$

② Softmax

$$\hat{p}_i = \frac{e^{v_i}}{\sum_{j=1}^C e^{v_j}}$$

Converts model output into probabilities

③ Categorical Cross Entropy (many classes)

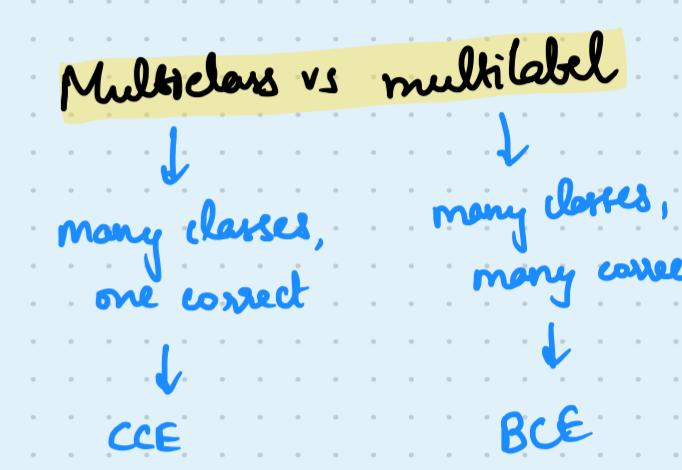
$$L_{cat} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^C y_{ij} \log(\hat{p}_{ij})$$

\downarrow cost func: iterates over samples

\downarrow loss func: iterates over all classes

\downarrow sparse \rightarrow converts integer to one-hot

\downarrow N \rightarrow number of samples



HINGE LOSS

Hinge opens in one direction \rightarrow loss increases only if certain criterion not satisfied

\downarrow
does not distinguish level of satisfaction

$$L = \sum_{i=0, i \neq c}^{N-1} \max(0, y_i - y_c + m)$$

\downarrow if $y_i - y_c > m$:
 \downarrow no change to loss, 0
 \downarrow then loss changes

typically $m=1$

Squared hinge loss

$$L = \frac{1}{N} \sum_{i=0}^{N-1} \max(0, 1 - y_i \cdot \hat{y}_i)^2$$

Focal Loss: IMBALANCED DATASETS

$$L = -(1-y)^{\alpha} \log y$$

where $y_t = \begin{cases} y & \text{if GT is class 1} \\ 1-y & \text{if GT is class 0} \end{cases}$

CONTRASTIVE LOSS

open-set task \rightarrow model needs to pull similar instances together, generalize classes, push dissimilar away

self-supervised / pretext task \rightarrow supervised loss function

$$L(w, y, \vec{x}_1, \vec{x}_2) = (1-y) L_s(D_w) + y L_d(D_w)$$

L_s : loss function for similar pair \vec{x}_1, \vec{x}_2
 \downarrow as $D_w \uparrow \rightarrow L_s \uparrow$
 L_d : loss function for dissimilar pair \vec{x}_1, \vec{x}_2
 \downarrow as $D_w \uparrow, L_d \uparrow$ but only until margin m
 \downarrow after certain distance m , no point pushing further apart

$$D_w(\vec{x}_1, \vec{x}_2) = \|G(\vec{x}_1) - G(\vec{x}_2)\|_2$$

where $G(\cdot)$ \rightarrow mapping function to transformed input space

where distance has meaning

Say we use Euclidean distance

$$L_s(D_w) = \frac{1}{2} D_w^2$$

$$L_d(D_w) = \frac{1}{2} [\max(0, m - D_w)]^2$$

{ once D_w becomes greater than m , no penalty }

$$L_s(D_w) = (1-y) \left[\frac{1}{2} D_w^2 \right] + y \left[\frac{1}{2} [\max(0, m - D_w)]^2 \right]$$

Focal Loss \rightarrow imbalanced learning

for i^{th} sample:

$$L(y_i) = \begin{cases} -(\bar{y})^\alpha \log y_i & \text{if } \text{gt} = 1 \\ -\bar{y}^\alpha \log(1-y_i) & \text{if } \text{gt} = 0 \end{cases}$$

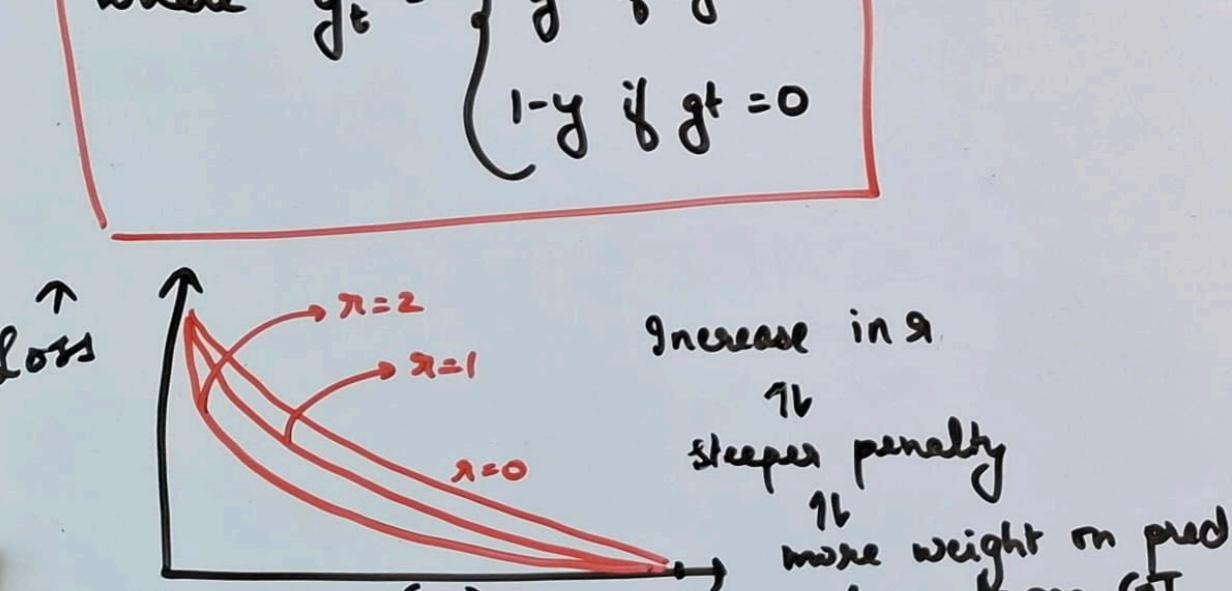
weight results by difference from gt
 $\text{GT}=1 \rightarrow 1-y$
 $\text{GT}=0 \rightarrow 0-y$

$$L(y) = -(\bar{y})^\alpha \log y$$

where $y_t = \begin{cases} y & \text{if } \text{gt} = 1 \\ 1-y & \text{if } \text{gt} = 0 \end{cases}$

$$L(y) = -(\bar{y})^\alpha \log y$$

where $y_t = \begin{cases} y & \text{if } \text{gt} = 1 \\ 1-y & \text{if } \text{gt} = 0 \end{cases}$



Contrastive Loss \rightarrow self-supervised learning

pull similar samples together, push dissimilar ones away

$$L(w, y, \vec{x}_1, \vec{x}_2) = (1-y) L_s(D_w) + y L_d(D_w)$$

$y=0$ gives similar loss

$y=1$ gives dissimilar loss

L_s : loss function for similar pair \vec{x}_1, \vec{x}_2

\downarrow as $D_w \uparrow \rightarrow L_s \uparrow$

L_d : loss function for dissimilar pair \vec{x}_1, \vec{x}_2

\downarrow as $D_w \uparrow, L_d \uparrow$ but only until margin m

after certain distance m , no point pushing further apart

$D_w(\vec{x}_1, \vec{x}_2) = \|G(\vec{x}_1) - G(\vec{x}_2)\|_2$

where $G(\cdot)$ \rightarrow mapping function to transformed input space

where distance has meaning

Say we use Euclidean distance

$$L_s(D_w) = \frac{1}{2} D_w^2$$

$$L_d(D_w) = \frac{1}{2} [\max(0, m - D_w)]^2$$

{ once D_w becomes greater than m , no penalty }

$$L_s(D_w) = (1-y) \left[\frac{1}{2} D_w^2 \right] + y \left[\frac{1}{2} [\max(0, m - D_w)]^2 \right]$$

$$d[w, y, \vec{x}_1, \vec{x}_2] = (1-y) \left[\frac{1}{2} D_w^2 \right] + y \left[\frac{1}{2} [\max(0, m - D_w)]^2 \right]$$

Triplet Loss

Similar points closer together than dissimilars only by some margin

$$L = \max(0, d(a,p) - d(a,n) + m)$$

\downarrow same label as anchor

\rightarrow uses three points: anchor, positive & negative

\downarrow diff label

\rightarrow from above:

Once distance from a to p differs from distance from a to n by the margin m or more,

NO PENALTY.

\rightarrow allows it to maintain some intra-class variance

\rightarrow considers margin for both +ve and -ve; contrastive considers only for -ve

contrastive for +ve: L_s

$$L_s = \frac{1}{2} D_w^2$$

minimizing L_s pushes D_w to zero

triplet only pushes to margin m

\rightarrow less greedy; more robust and allows it to learn more discriminative features