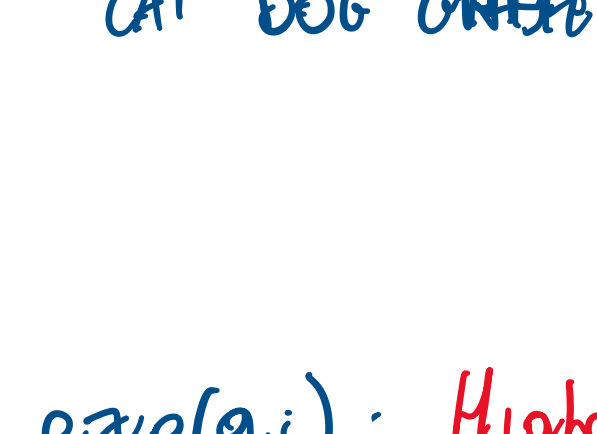


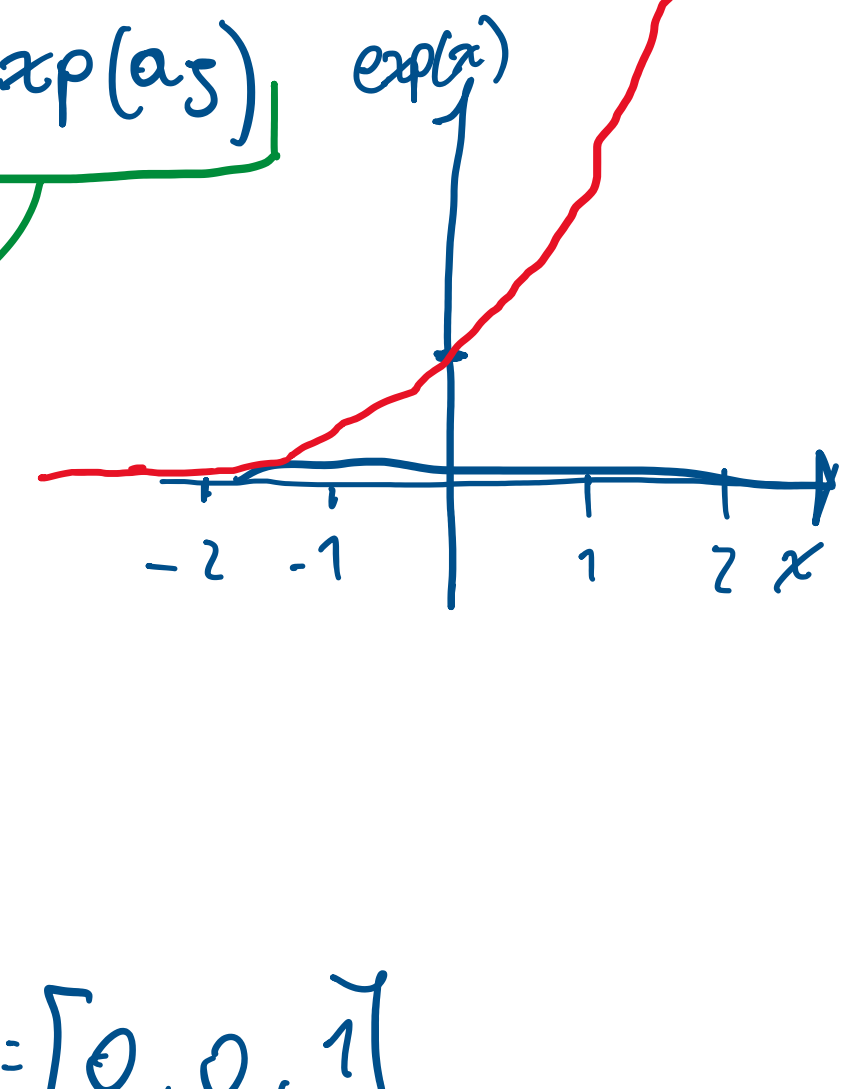
# Generalization for $C > 2$ Multiclass - Classification



Softmax maps vectors to probability simplex.

$$\text{softmax}(a_i) = \frac{\exp(a_i)}{\sum_j \exp(a_j)}$$

$\exp(a_i)$ : Higher  $a_i$  get large probability  
Lower  $a_i$  get scaled down



ensure everything sum up to 1

$$f(\mathbf{x}^i) = \text{softmax}(\Theta^T \mathbf{x}^i) \rightarrow \text{Logits}$$

cat = [1, 0, 0]    dog = [0, 1, 0]    giraffe = [0, 0, 1]

Probability distribution with all the mass on a single class.  
 $f(\mathbf{x})$  encodes the probability distribution over all the classes

$f(\mathbf{x})_i$  represent the probability for the single class  $i$

$$P(Y | f(\mathbf{x}^i)) = \prod_i [f(\mathbf{x}^i)]^{y_i}$$

since  $\mathbf{y}$  is one-hot vector only  $f(\mathbf{x})_i$  is going to contribute in the product  
 $\mathbf{y} = [0, 1, 0]$

$$P(Y^i | f(\mathbf{x}^i)) = f(\mathbf{x}^i)_{y^i}$$

$$f(\mathbf{x}^i)_1^0 \cdot f(\mathbf{x}^i)_2^1 \cdot f(\mathbf{x}^i)_3^0$$

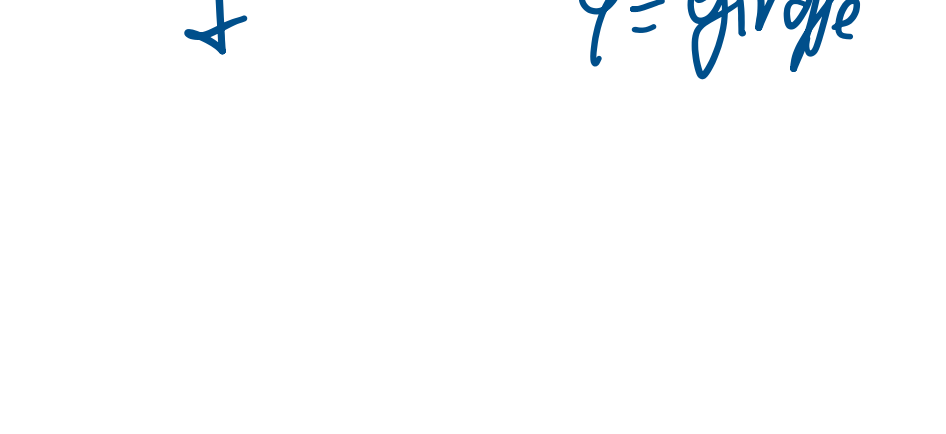
$$L = \prod_i P(Y^i | f(\mathbf{x}^i))$$

$$\log L = \sum^n \log P(Y^i | f(\mathbf{x}^i))$$

$$= \sum^n \log f(\mathbf{x}^i)_{y^i}$$

$$NLL = - \sum^n \log f(\mathbf{x}^i)_{y^i}$$

Cross Entropy



$$y = \text{cat} = - \left( 1(-1) + 0(-0.698) + 0(-0.15) \right)$$

$$= 1$$

$$y = \text{giraffe} = - \left( 0(-1) + 0(-0.69) + 1(-0.15) \right)$$

$$= 0.15$$

$$\mathbf{x}_1 \quad y_1 = \text{cat} \quad [1, 0, 0] \quad f(\mathbf{x}_1) = [0.1, 0.2, 0.7]$$

$$\mathbf{x}_2 \quad y_2 = \text{dog} \quad [0, 1, 0] \quad f(\mathbf{x}_2) = [0.01, 0.98, 0.01]$$

$$CE(\theta) = - \sum_i \sum_j y_j^i \log f(\mathbf{x}^i)_j$$

$$= \mathbf{x}_1 \cdot \left( y_1^1 \log(f(\mathbf{x}_1)_1) + y_1^2 \log(f(\mathbf{x}_1)_2) + y_1^3 \log(f(\mathbf{x}_1)_3) \right) +$$

$$\mathbf{x}_2 \cdot \left( y_2^1 \log(f(\mathbf{x}_2)_1) + y_2^2 \log(f(\mathbf{x}_2)_2) + y_2^3 \log(f(\mathbf{x}_2)_3) \right)$$

$$= - \left( 1(-0.91) + 0(-1.6) + 0(-0.91) + 0(-4.6) + 1(-0.02) + 0(-4.6) \right)$$

$$\log(f(\mathbf{x}_1)) = [-0.91, -1.6, -0.91]$$

$$\log(f(\mathbf{x}_2)) = [-4.6, -0.02, -4.6]$$

$$= -(-0.91 - 0.02) = 0.93$$

Maximise the probability of the true class at the expense of the other outputs.

Softmax: unstable :  $-\log\left(\frac{\exp p_i}{\sum_j \exp p_j}\right)$

$$-p_i + \log \sum \exp(p)$$

$$\log(\sum \exp(p_i - c)) + c$$

avoid 0s messing with our optimization by adding a small constant  $c$

if you set  $c = \max(p)$  ensure numerical instabilities will never occur

$$-\frac{1}{2n} (\text{Loss})$$

used for differentiation of square differences  
Normalize by number of samples in the dataset

$f(\mathbf{x}) = \text{softmax}(\Theta^T \mathbf{x}^i) \rightarrow$  Multiclass Logistic Regression  
best set of parameters obtained by minimizing  $CE(\theta)$

Entropy  $H(p) = - \sum_i p_i \log p_i$     measure uncertainty of distribution  $p$

Cross Entropy  $CE(p, q) = - \sum_i p_i \log q_i$     measure uncertainty of distribution  $p$  when we encode it with  $q$

KL Divergence  $KL(p||q) = \sum_i p_i \log \frac{p_i}{q_i}$     how much information is lost when you use  $q$  instead of  $p$

$$CE(p, q) = H(p) + KL(p||q)$$

uncertainty of true distribution    additional cost of using  $q$  instead of  $p$

$f(\mathbf{x})_i$ : logits cannot be considered as real probabilities.

$$P(Y=i | \mathbf{x}) = f(\mathbf{x})_i \text{ unless this happens!}$$

How to measure if a model is calibrated.

- 1. pick a validation set
- 2. split  $[0, 1]$  into  $m$  bins  $[1/m]$

$B_m$ : number of sample whose confidence falls into bin  $m$

$p_m$ : average confidence for each bin

$a_m$ : average accuracy for each bin

$$\text{Expected Calibration Error (ECE)} = \sum_m \frac{B_m}{n} |a_m - p_m|$$

## Hyperplanes

$N$ -dimensional space  $\rightarrow$   $N-1$  dimensional flat subspace

2D  $\rightarrow$  1D line    1 can only separate classes that are linearly separable

3D  $\rightarrow$  2D plane    2 distance from closest data sample and hyperplane is called margin

## Bias - Variance

Underfitting    Perfect fit

overfitting

underfitting

overfitting

error

variance

bias

model capacity

perfect model capacity

difference between the expected prediction and the correct one

variance: variability of the prediction for a given dataset

expected dataset  $\rightarrow$  several trainings on different data partitions

variance

$\theta^*$  obtained by minimizing Expected Risk

$\theta^0, \theta^1, \dots$  obtained by minimizing the Empirical Risk

$$E[(Y - f(\mathbf{x}))^2] = \text{Bias}(f(\mathbf{x}))^2 + \text{Variance}(f(\mathbf{x})) + \sigma^2$$

irreducible error