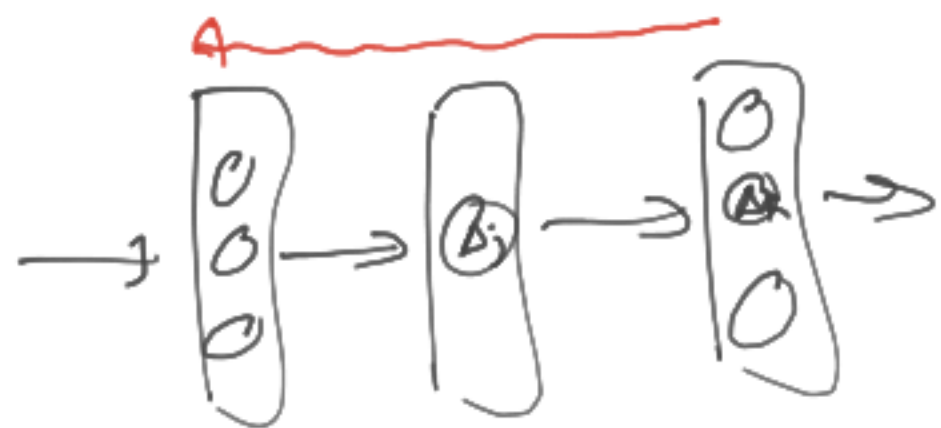


Vanishing Gradients

Recall:

$$\frac{\partial f_w(x)}{\partial w_{ij}} = \sum_k \Delta_k w_{jk} \underbrace{\sigma'(in_j)}_{\frac{\partial \sigma(in_j)}{\partial in_j}} \underbrace{(1 - \sigma(in_j))}_{a_i}$$

Δ_j



\downarrow

$$\frac{\partial f_w(x)}{\partial w_{ij}} \rightarrow 0$$

as you go further and further back in the network.



- ## Lecture 9
- vanishing gradients
 - weight init
 - Adaptive step size
 - Classification
 - Generalization Bounds

Weight Initialization

- Heuristic

- Want:

- mean of in_j to be zero.

- variance of in_j should be the same across all layers.

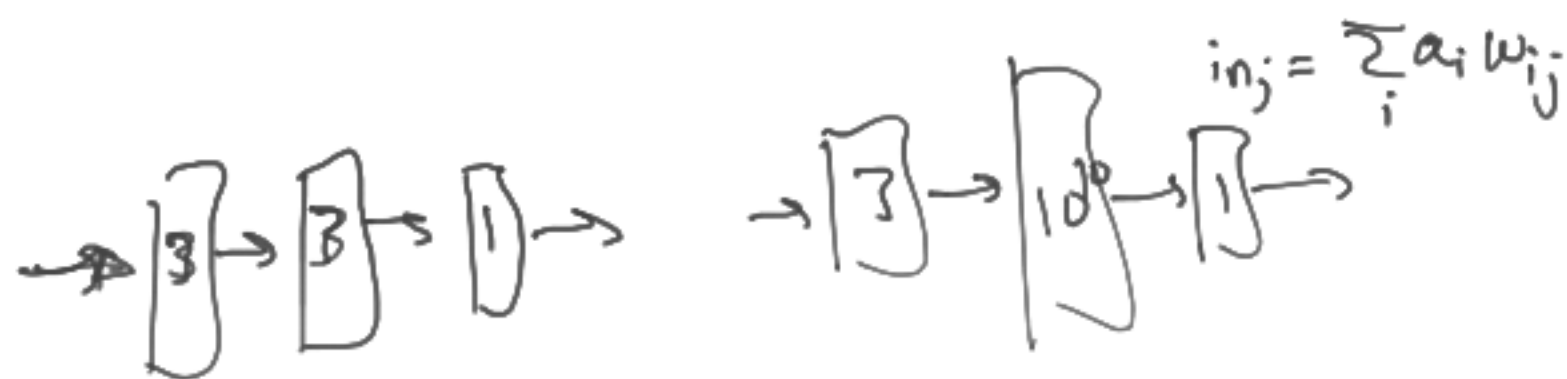
- Xavier initialization: For layer l

$$w_{ij} \sim N\left(0, \frac{1}{n_{l-1}}\right), \text{ where } n_{l-1}$$

is number of nodes in layer $l-1$.

- He initialization (effective for ReLU).

$$w_{ij} \sim N\left(0, \frac{2}{n_{l-1}}\right)$$

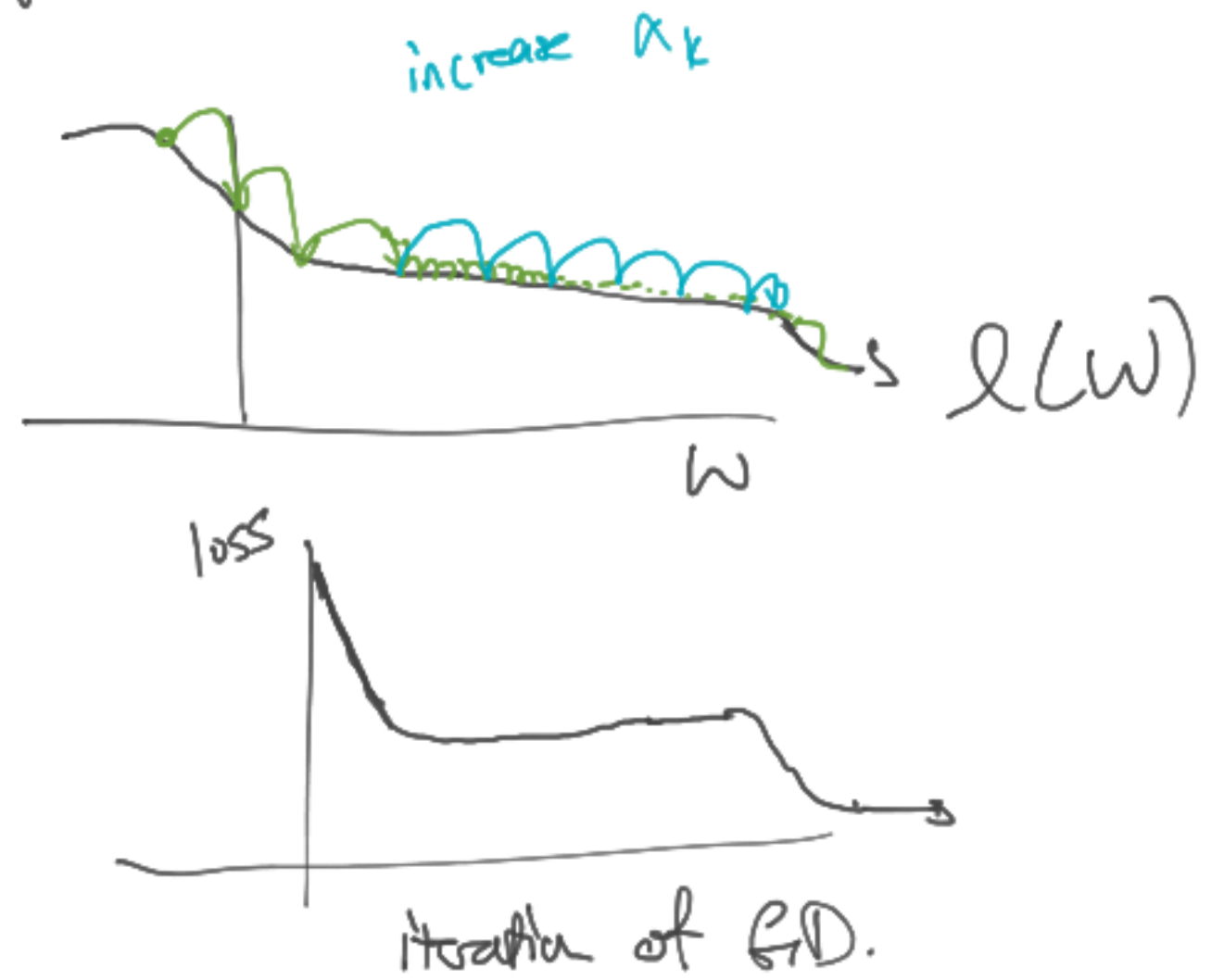


Adaptive Step Sizes. α_k (step size at k^{th} iteration of GD)

Momentum: $\eta \in [0, 1]$
(often decayed to zero)

$$\Delta \theta_k = \eta \Delta \theta_{k-1} + \alpha \nabla l(w_k)$$

$$\theta_{k+1} = \theta_k - \Delta \theta_k$$



Many more! (Keep a different step size for each model parameter!)

- Adaptive Gradient (AdaGrad)
- Root mean square propagation (RMS Prop)
- Adaptive moment estimation (Adam)

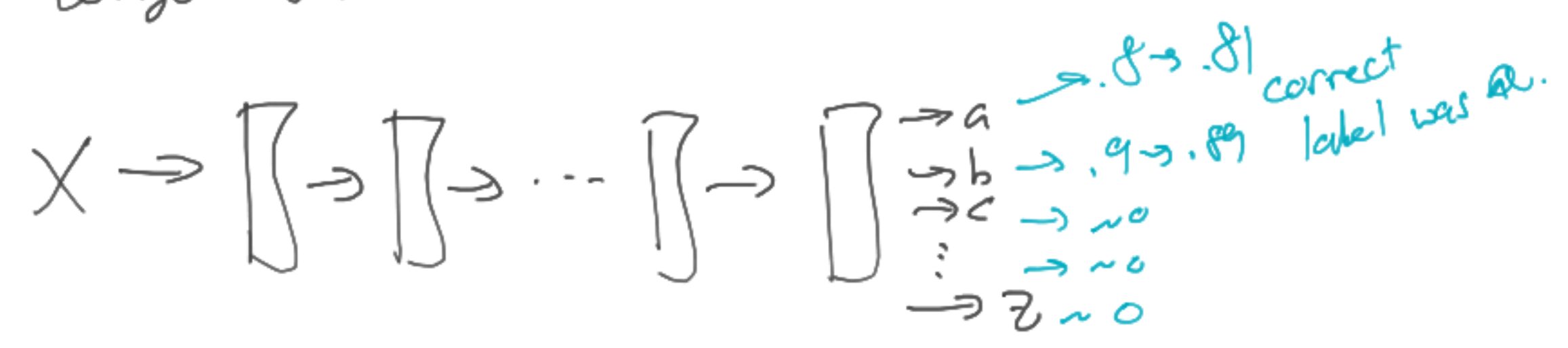
Classification

Like regression, but $y_i \in \mathcal{Y}$, where \mathcal{Y} is a set (typically finite) of labels.

(typically)

Examples: $\mathcal{Y} = \{\text{is a cat in image, there is no cat}\}$
 $\mathcal{Y} = \{a, b, c, \dots, z\}$

Idea: One model output per label.
Larger values mean that label is more reasonable



- Predicted label is the largest output.

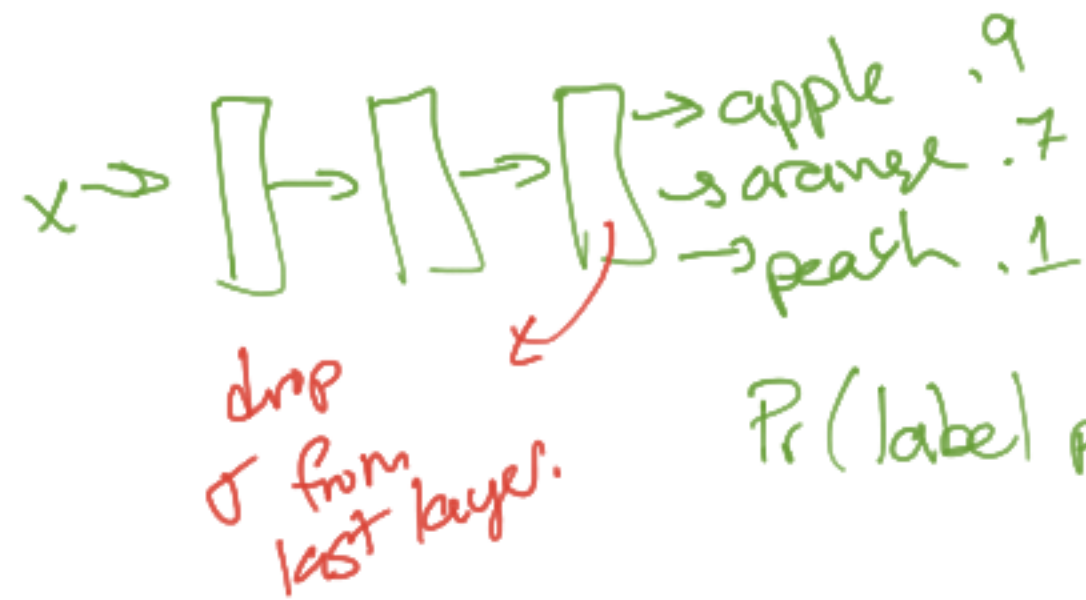
- Not differentiable, (due to max operator)

Idea: Use "softmax" for last layer.

$$Pr(\text{label} = l) = \frac{e^{a_l}}{\sum_{l'} e^{a_{l'}}$$

make positive

$$\left. \begin{matrix} 10 \\ -20 \\ 700 \\ -3 \end{matrix} \right\} \left. \begin{matrix} e^{10} \\ e^{-20} \\ e^{700} \\ e^{-3} \end{matrix} \right\} \left. \begin{matrix} e^{10/\beta} \\ e^{-20/\beta} \\ e^{700/\beta} \\ e^{-3/\beta} \end{matrix} \right\} \text{sum to one.}$$



$$Pr(\text{label prediction} = \text{apple}) = \frac{e^{.9}}{e^{.9} + e^{.7} + e^{.1}}$$

$$\approx 0.44.$$

$$Pr(\text{orange}) \approx .36$$

$$Pr(\text{peach}) \approx .19$$

$$\beta = \frac{e^{10} + e^{-20} + e^{700} + e^{-3}}{e^{10/\beta} + e^{-20/\beta} + e^{700/\beta} + e^{-3/\beta}}$$

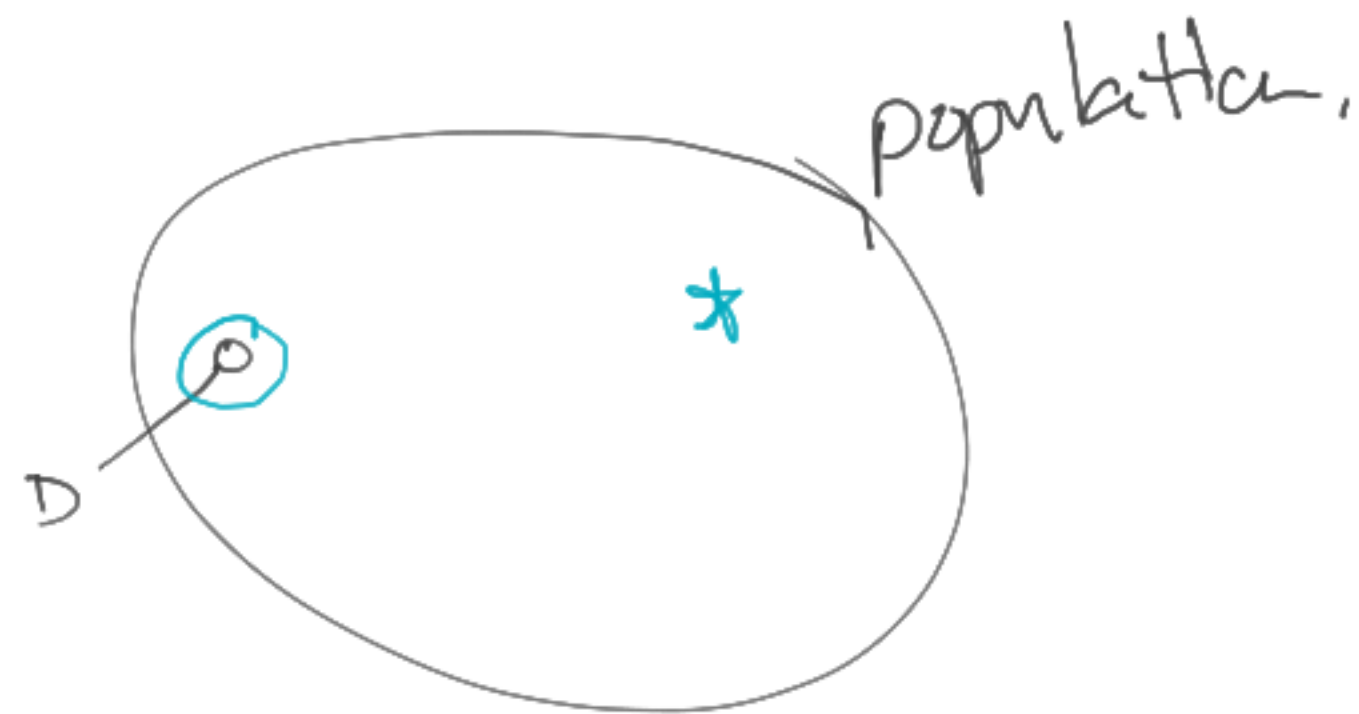
Cross Entropy Loss (Binary)
- Binary Classification, $|Y| = 2$, $\mathcal{Y} = \{0, 1\}$
 $\mathcal{Y} = \{-1, 1\}$

$$l(\omega) = -\frac{1}{n} \sum_{i=1}^n \ln \left(\Pr \left(\underbrace{f_{\omega}(x_i)}_{\substack{\text{random} \\ \text{due to} \\ \text{softmax}}} = y_i \right) \right)$$

monotonic.

Generalization Bounds

$D = n$ samples for training.



Data is provided with no explanation.
Could be chosen any (not random) way.

→ Assume data points are sampled i.i.d. from population.

independent
identically distributed.

$$l(w) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

(x_i, y_i) are i.i.d.

x_i (not input x_i)

$$\frac{x_1 + x_2}{2}$$

$$\frac{x_1 + x_2 + x_3}{3}$$

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n x_i$$

Hoeffding's inequality:

If x_1, \dots, x_n are i.i.d. then

$$\Pr(|\bar{X}_n - \mu| \geq t) \leq 2e^{-2nt^2}$$

x_1, \dots, x_{1000}

$$\bar{X}_n = 1.7m$$

random.

$$\Pr(|\boxed{1.7m} - \mu| \geq 0.1) \leq 2e^{-2(1000)(0.1)^2}$$

$$= 2(1000)(0.1)^2$$

.0000000004