

## Lecture 9 (India)

Setup/Assumptions  
Bias/Variance  
Aprox Estim  
Empirical risk minimizer  
Uniform Convergence  
VC dimension

### Assumptions

1. Data distribution  $D: (x, y) \sim D$  - train/test det. func.

2. Independent Samples  
 $\hat{\theta}^*$  or  $h^*$  = "true para",  
not random

$x_1, y_1$
:
$x_m, y_m$

$\sim D$

learning algo

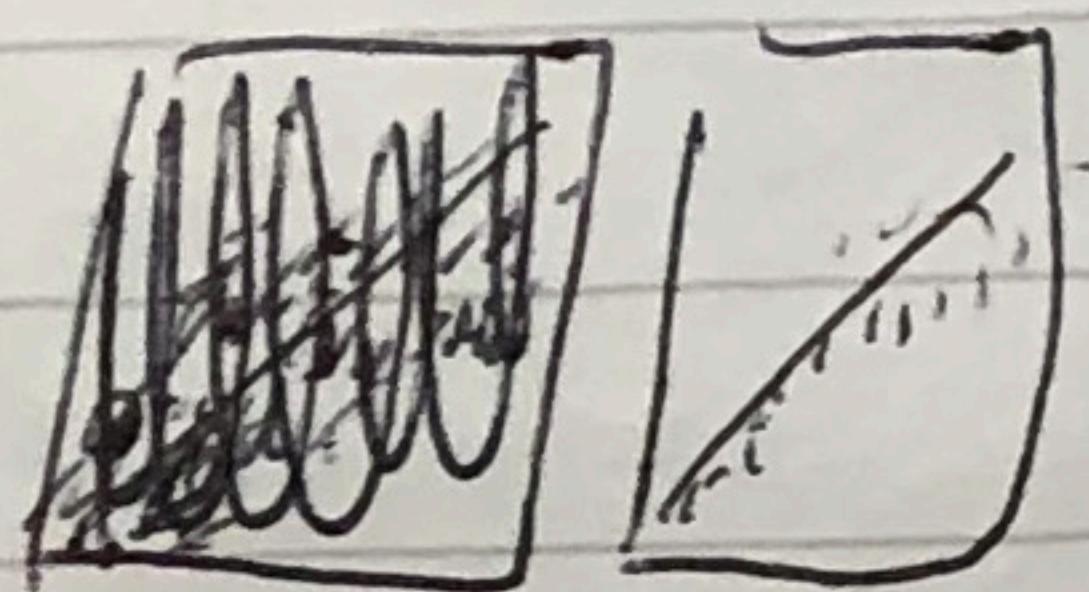
[estimation]

random variables  $\rightarrow S$

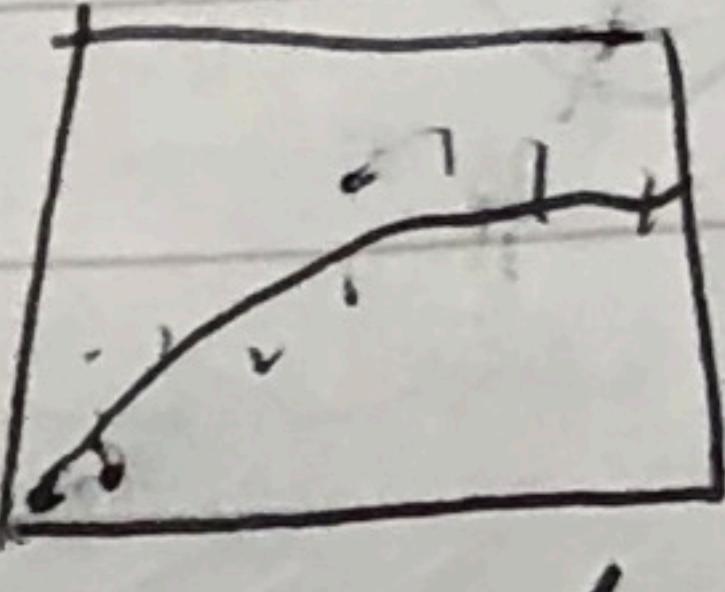
$h$  or  $\hat{h}$

random variables  $\stackrel{?}{\sim}$   
Sampling distribution

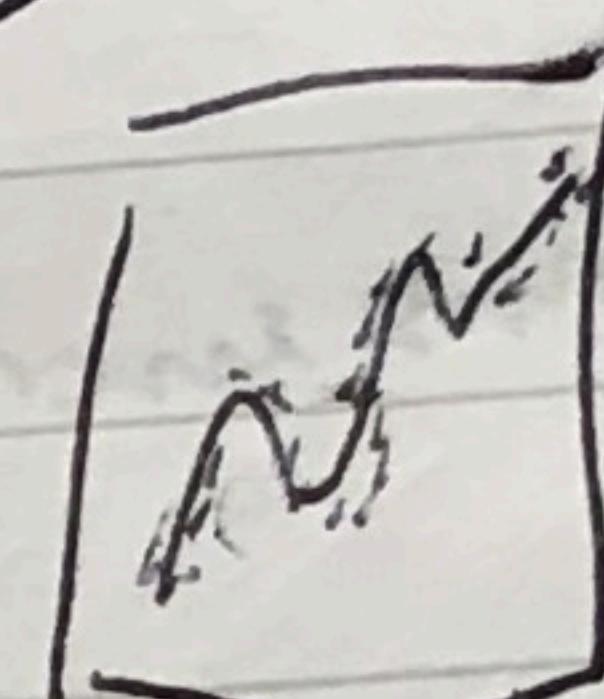
Reg



underfit  
(High bias)

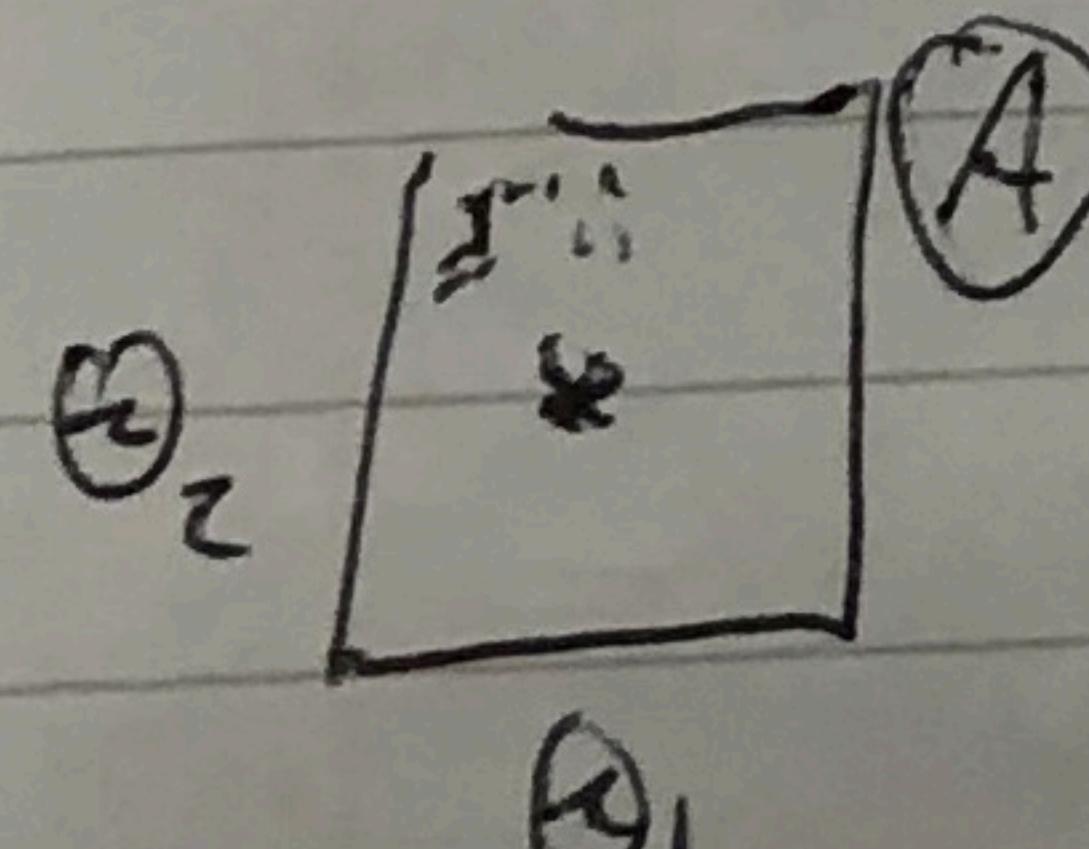


correct



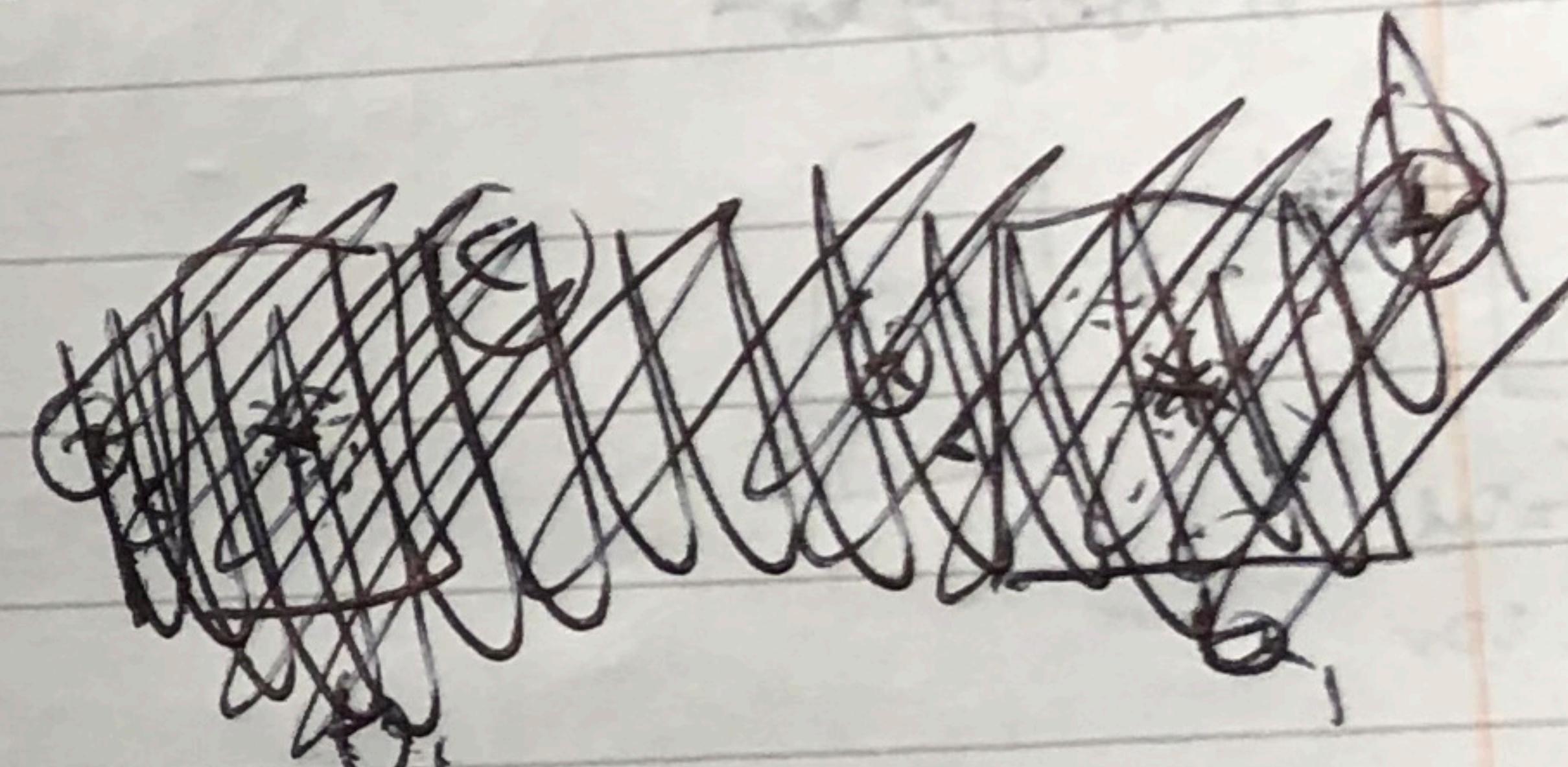
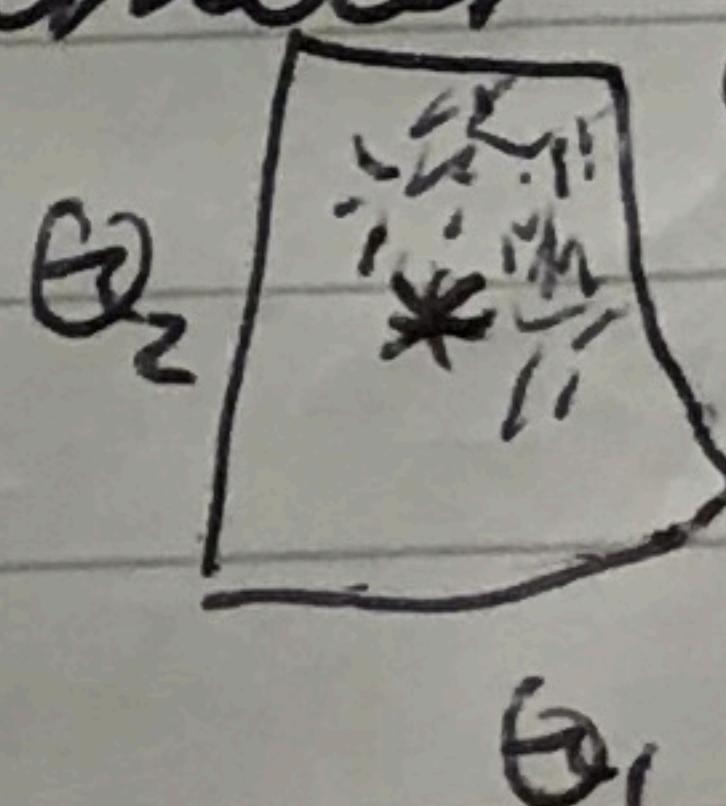
overfit  
(High variance)

bias

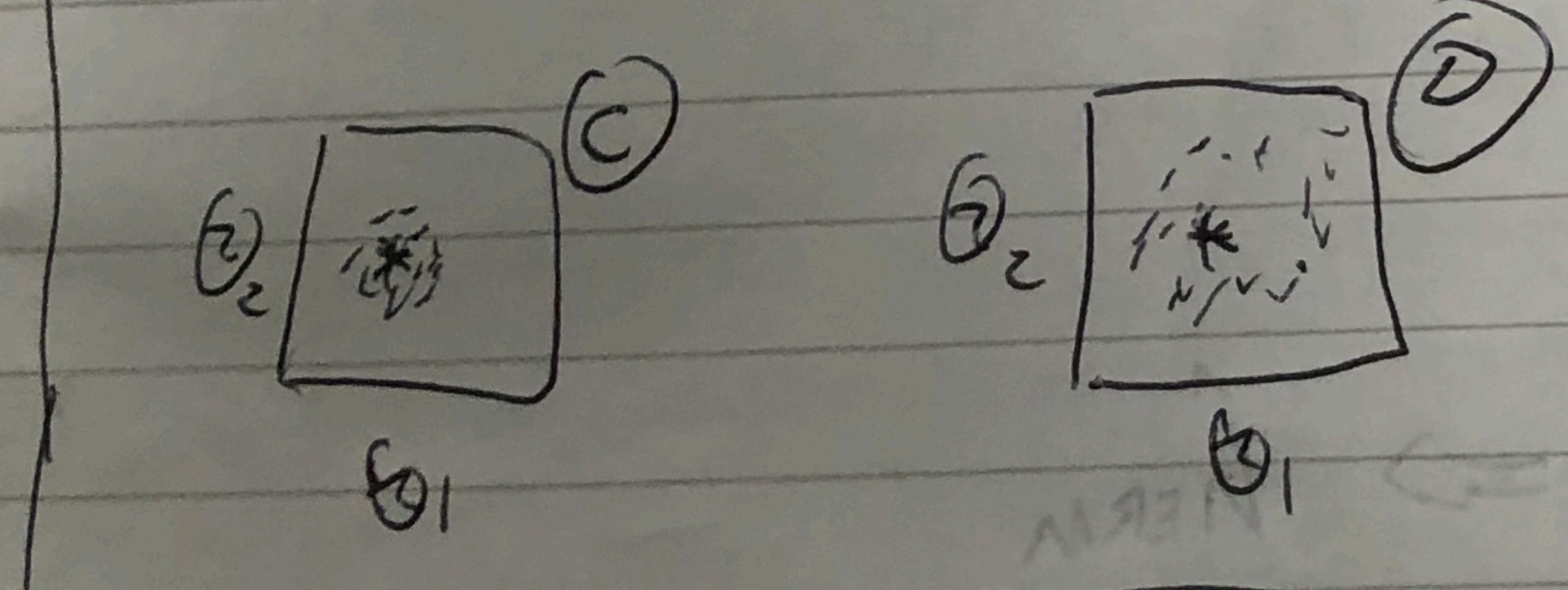


algorithm in datapoints

Parameter occur



=> Gray: is sampling distribution  
centered around "true para"?  
Variance: var. of Sampling distr



$\text{Var}[\hat{\theta}]$

Variance

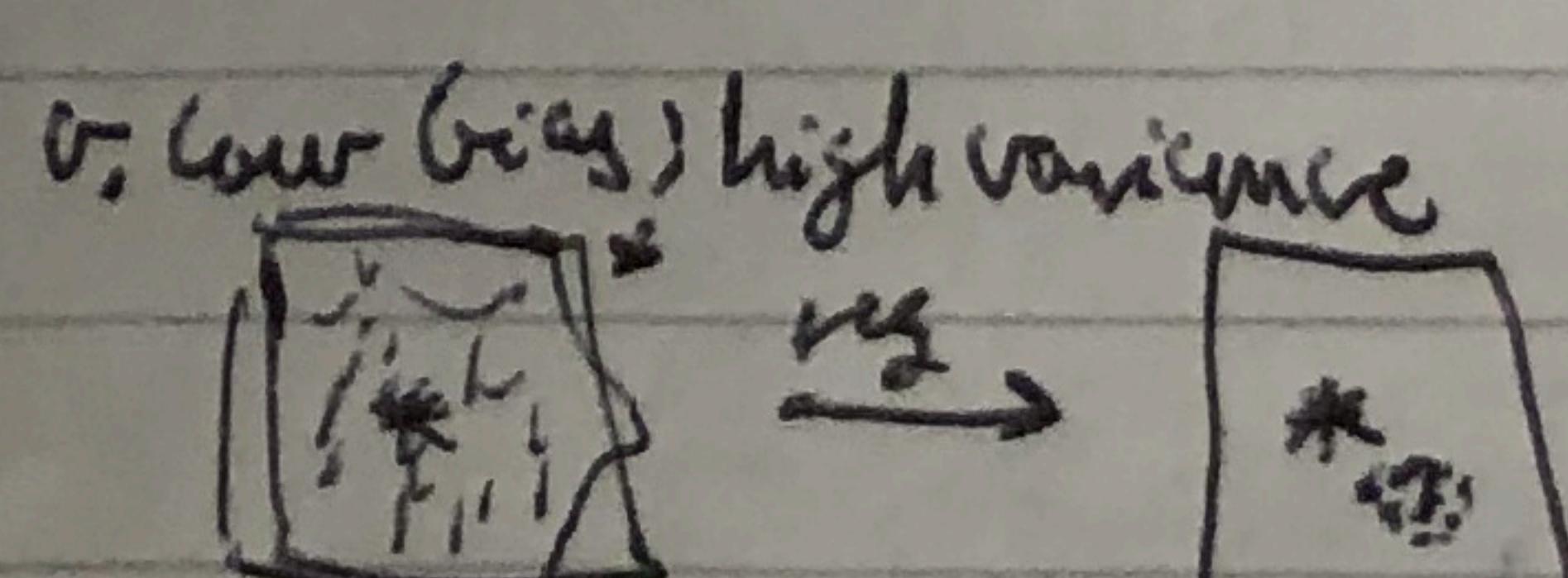
As  $m \rightarrow \infty$ , variance  $\rightarrow 0$ . The rate of this = "statistical efficiency".  
 $\hat{\theta} \rightarrow \theta^*$  as  $m \rightarrow \infty$ : Consistent algo.  $\Rightarrow \hat{\theta} \xrightarrow{P} E(\hat{\theta}) = \theta^*$

get in optimal hypothesis  
class h\_\*

Fighting variance:

1) ~~more data~~  $m \rightarrow \infty$

2) Regularisation:



small bias;  
low variance

some hys  
learned  
from finite  
data

$\hat{h}_*$

$h^*$

$\cdot g$

$\mathcal{H} \leftarrow$  eg set of all hys

space of hypotheses

$$E(h) \triangleq \text{Risk/generalisation error} = E_{(x,y) \sim D} [\mathbb{1}\{h(x) \neq y\}]$$

in case of  
classification

$$\hat{\varepsilon}(h): \text{Empirical risk} = \frac{1}{m} \sum_{i=1}^m \mathbb{1}\{h(x_i) \neq y_i\}, S = \text{sample of size } m \text{ datapoints}$$

$E(g)$  = Bayes error / irreducible error

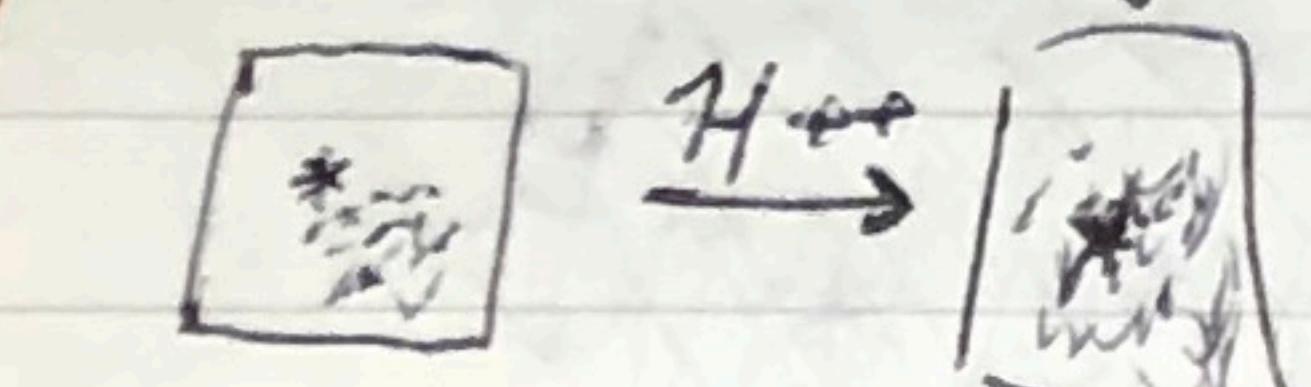
$$E(h^*) - E(g) = \text{Approximation error } \beta \text{ class}$$

$$E(h) - E(h^*) = \text{Estimation error: } h \approx g$$

$$E(h) = \underbrace{\text{Estimation error}}_{\text{est var}} + \underbrace{\text{Approx error}}_{\text{est bias var}} + \underbrace{\text{Irreducible error}}_0$$

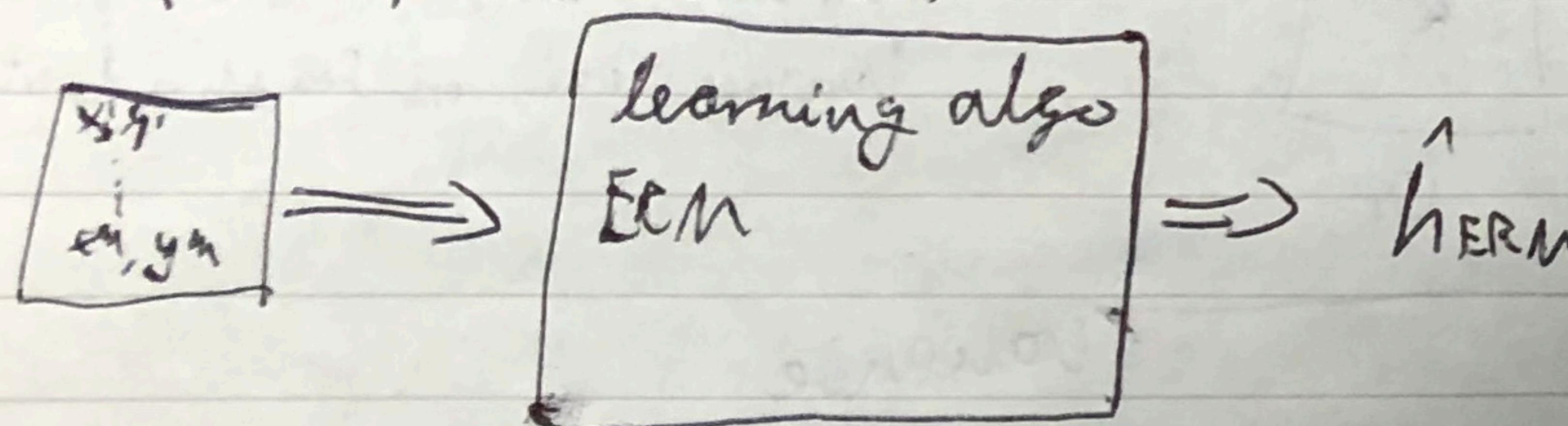
Fight high bias ideas:

1. Make  $H$  bigger



high bias  
& some var

ERM (emp. risk minimizer)



$$h_{ERM} = \underset{h \in H}{\operatorname{argmin}} \frac{1}{m} \sum_{i=1}^m \mathbb{1}\{h(x_i) \neq y_i\}$$

Uniform Convergence

1)  $E(h)$  vs  $E(h^*)$   
training err. vs generalisation err.

2)  $E(h)$  vs  $\hat{\varepsilon}(h)$

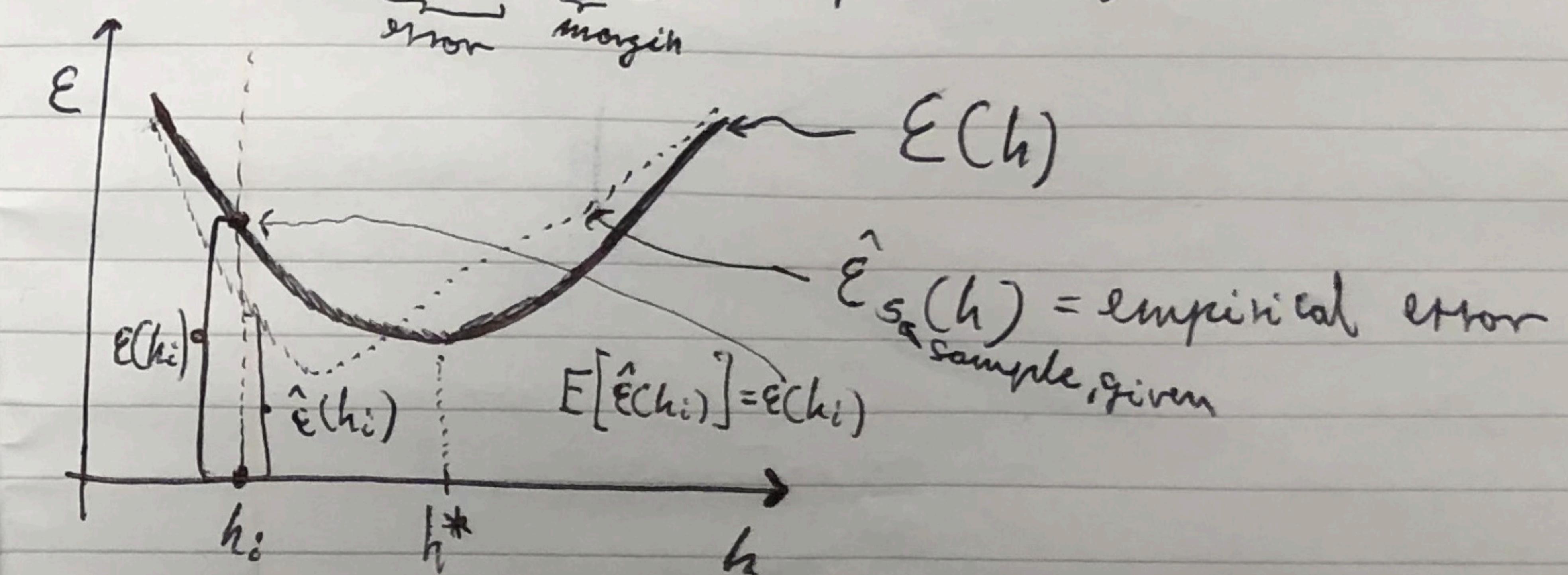
Tools

1) Union bound  $A_1, A_2, \dots, A_k$  (need not be indep.)

$$P(A_1 \cup A_2 \dots \cup A_k) \leq P(A_1) + P(A_2) + \dots + P(A_k)$$

2) Hoeffding's inequality

Let  $z_1, z_2, \dots, z_m \sim \text{Bern}(\frac{1}{2})$ ;  $\hat{\phi} = \frac{1}{m} \sum_{i=1}^m z_i$ , and  $\gamma > 0$  (margin). Then,  $P_{z_i} [|\hat{\phi} - \phi| > \gamma] \leq 2 \exp(-2\gamma^2 m)$



Hoeffding's eqn.:  $P_{z_i} [|\hat{\varepsilon}(h_i) - \varepsilon(h_i)| > \gamma] \leq 2 \exp(-2\gamma^2 m)$

Finite hypothesis class  $H$

$$\Rightarrow |H| = k; P_{z_i} [\exists h \in H | \hat{\varepsilon}_s(h) - \varepsilon(h) | > \gamma] \leq k \cdot 2 \exp(-2\gamma^2 m)$$

$$\Rightarrow P[\forall h \in H | \hat{\varepsilon}(h) - \varepsilon(h) | < \gamma] \geq 1 - \underbrace{2k \exp(-2\gamma^2 m)}_{\delta}$$

Let  $\delta = 2k \exp(-2\gamma^2 m)$ ;  $\delta$  - prob of error is  $\gamma$  - margin of error;  $m$  - sample size.

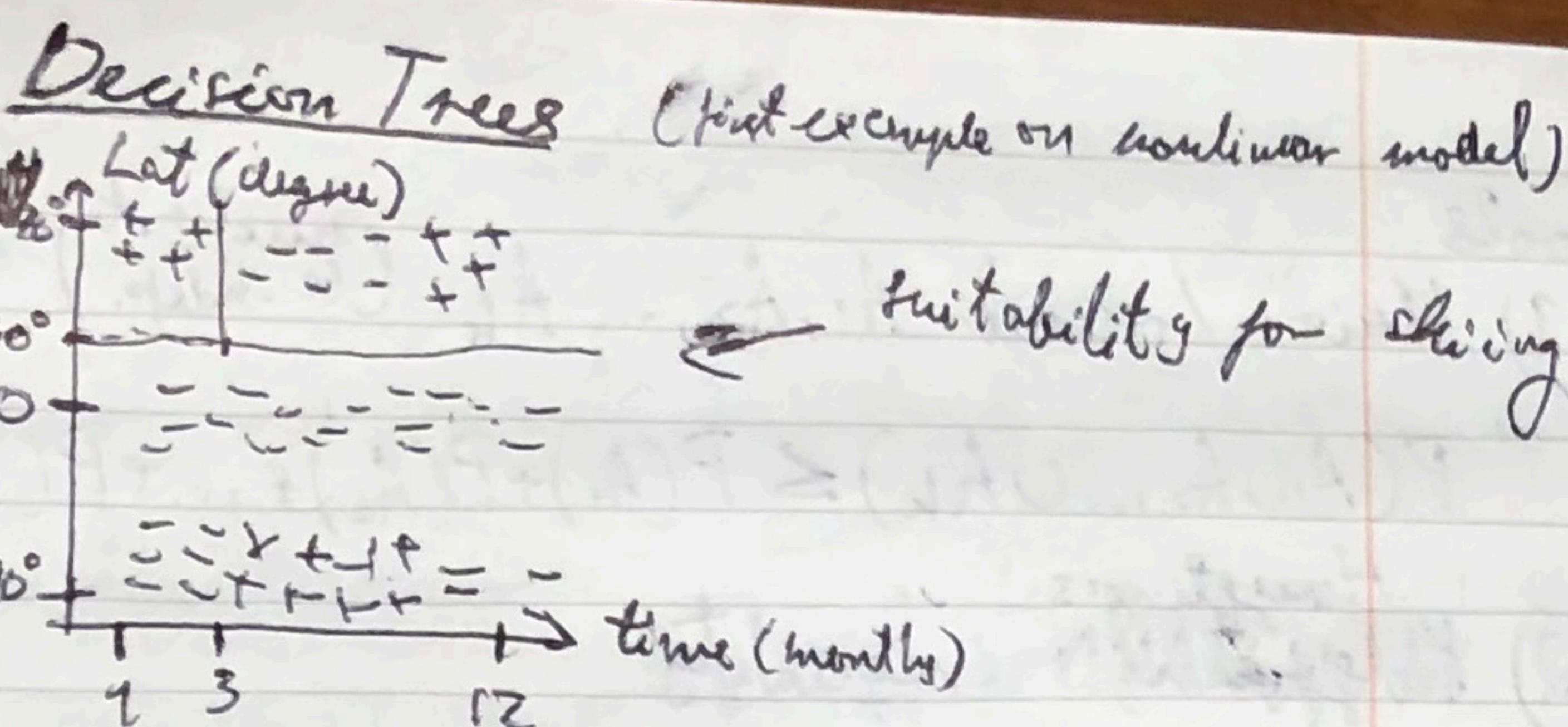
$\Rightarrow$  Can solve for  $\delta, \gamma, m$  when knowing the others. Eg. fix  $\delta, \gamma > 0$   
 $m \geq \frac{1}{2\gamma^2} \log \frac{2k}{\delta}$ . Known as the "sample complexity result".

VC dimension:  $\text{VC}(H) = m$

$$\Rightarrow E(h) \leq E(h^*) + O\left[\sqrt{\frac{\text{VC}(H) \log \frac{2m}{\delta}}{m}} + \frac{1}{m} \log \frac{2}{\delta}\right]$$

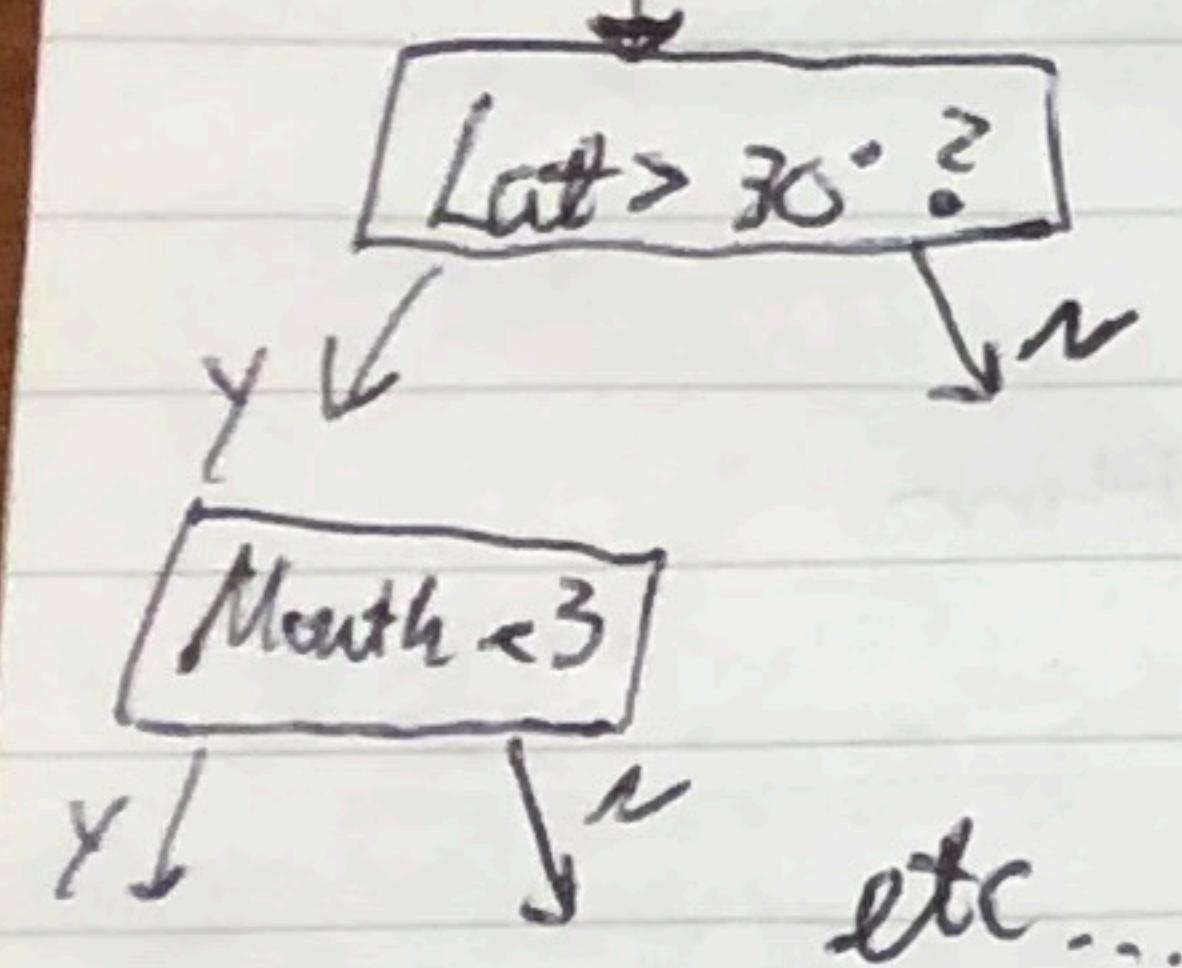
## Lecture 10

- Decision trees
- Ensemble methods
- Bagging
- Random forests
- Boosting



Greedy, top-down, recursive

Partitioning



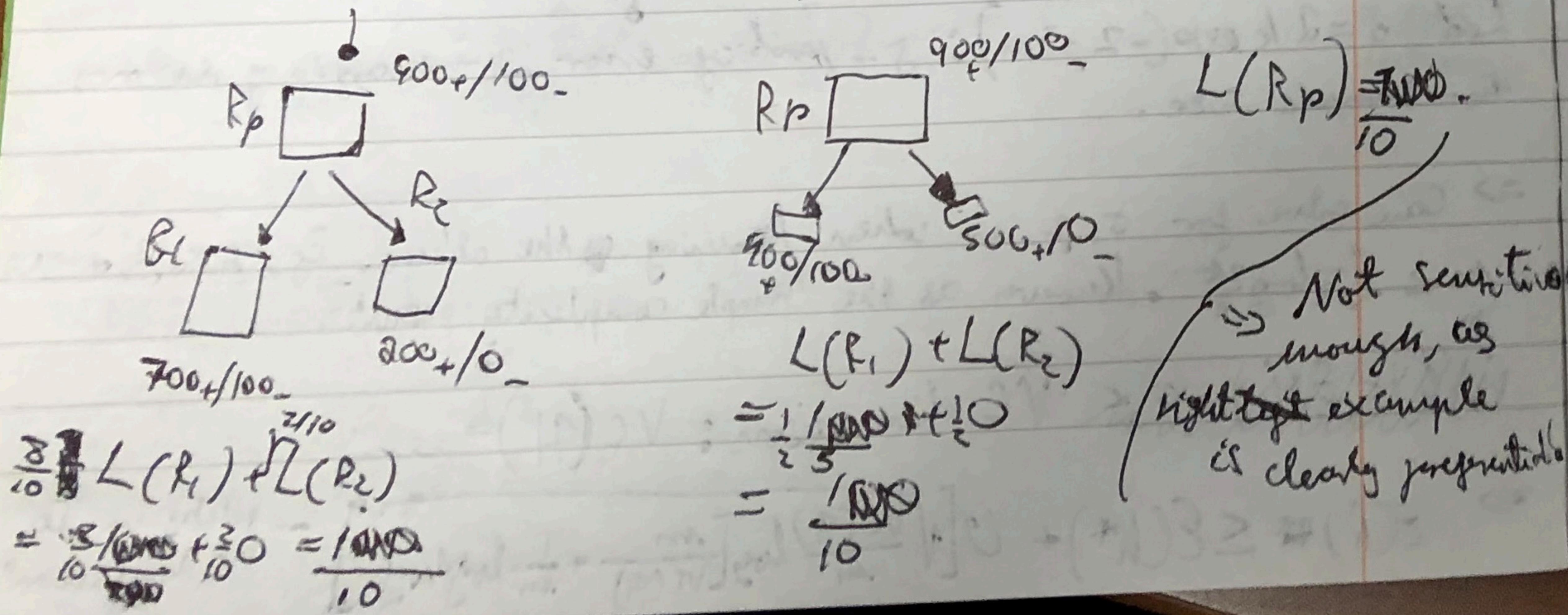
$\Rightarrow$  Region  $R_p$   
Looking for a split  $S_p$   
 $S_p(j, t) = \{x | x_j < t, x \in R_p\}$ ,  
 $\{x | x_j > t, x \in R_p\}$

How to choose splits?

Define  $L(R)$ : loss on  $R$ . Given  $C$  classes, define  $\hat{p}_c$  to be proportion of examples in  $R$  that are of class  $C$ .

$$\Rightarrow L_{\text{misclass}} = 1 - \max_c \hat{p}_c. \quad \max_{\text{int } t} L(R_p) = \frac{L(R_1) + L(R_2)}{\hat{p}_1 + \hat{p}_2} \quad \text{children lots?}$$

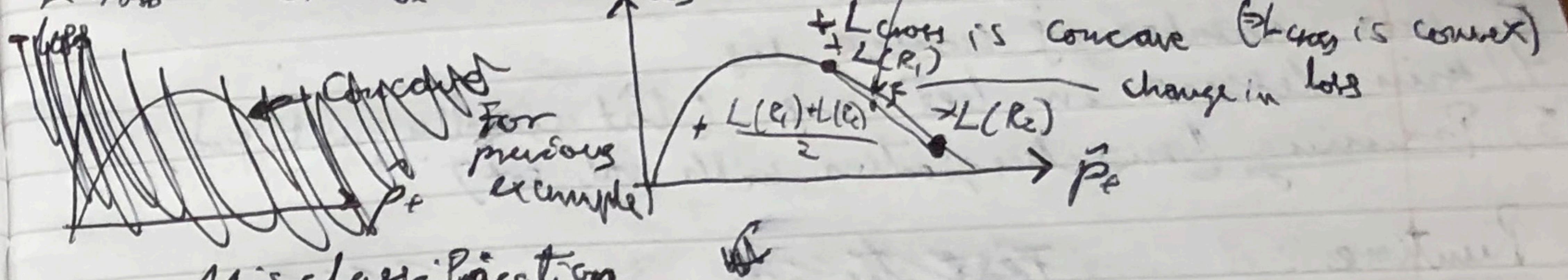
Misclassification Loss has issues!



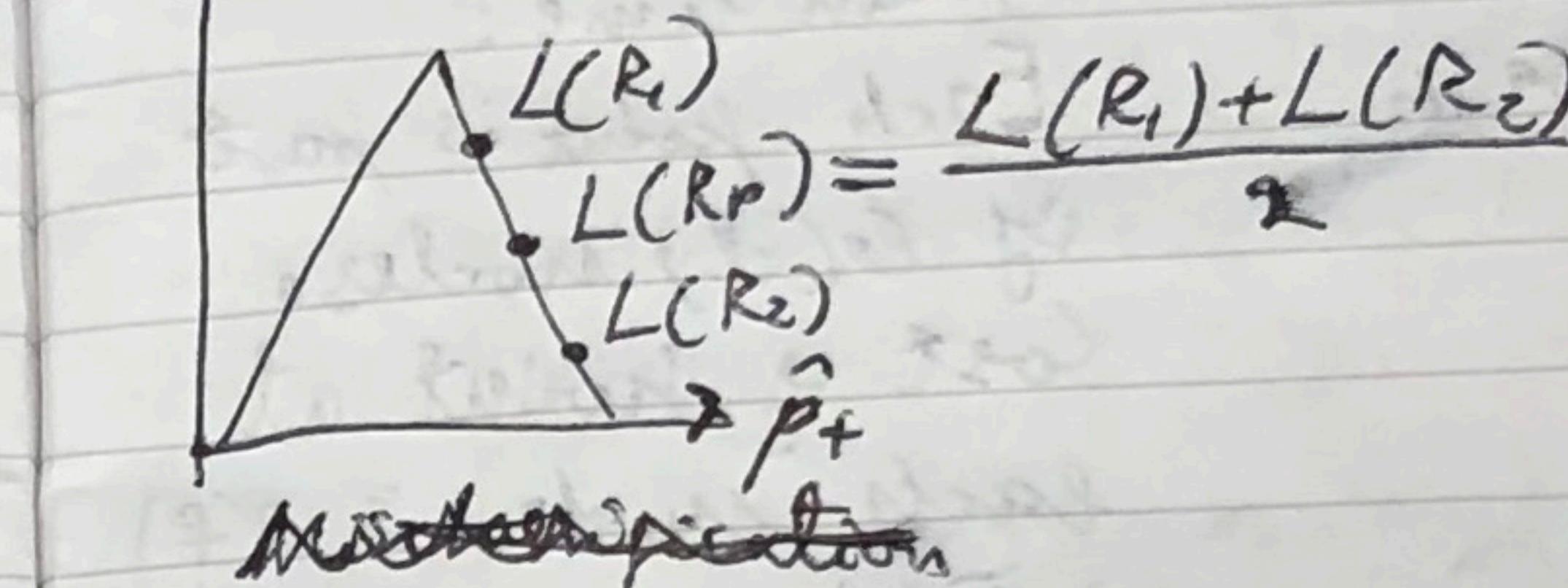
Instead, define cross-entropy loss

$$L_{\text{cross}} = -\sum_c \hat{p}_c \log \hat{p}_c$$

+  $L(R_p)$  [if even split]



Misclassification



Gini:  $\sum_c \hat{p}_c(1 - \hat{p}_c)$ . Looks similar to cross-entropy one.

Regression trees

Lat (deg)	10	20	30	40	50	60	70	80	90
10	78.600	9.10							
20	9.100	6.9							
30	0.000	0.0							
40	0.00	9.5	10.0	10.0					
50	0.0	15.10	10.0	10.0					
60									
70									
80									
90									

$$R_m \quad \text{Predict } \hat{Y}_m = \sum_{i \in R_m} \frac{(Y_i - \hat{Y}_m)}{|R_m|}$$

$$L_{\text{squared}} = \sum_{i \in R_m} \frac{(Y_i - \hat{Y}_m)^2}{|R_m|}$$

Categorical Variables

North	7.9	0.000	9.10
East	9.20	0.000	6.9
South	0.00	9.5	10.0

$\{Loc \in \{\text{Northern}\}\}$

9 categories

$2^9$  possible splits

if binning??

## Regularisation of DTs

- 1) min leaf size enforcement
- 2) max depth enforcement
- 3) max number of nodes en.
- 4) min decrease in loss en. (but usually isn't good)
- 5) Pruning (misclassifications with val. set)

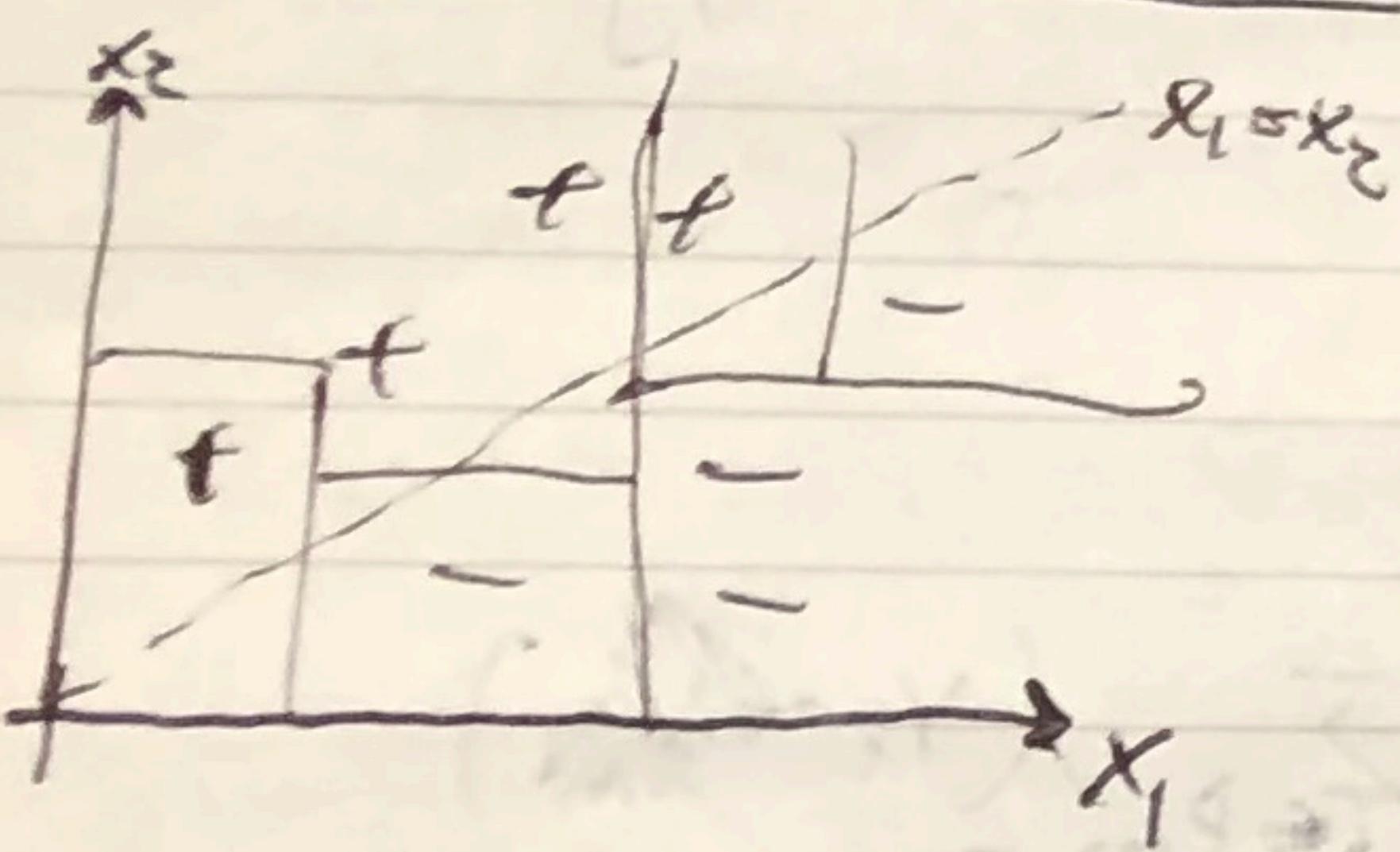
Runtime:

n examples  
f features  
d depth

Test time:  
 $O(d)$ ;  $d \log n$

Train time:  
Each point is part  
of  $O(d)$  nodes.  
Cost of point at  
each node is  $O(f)$   
 $\Rightarrow$  Total cost =  $O(nfd)$

No additive structure



Compare to linear regression to decision trees here... decision tree doesn't have the appropriate structure to capture this process.

Note that Data matrix is of size  $n \times p$ .

If we drop independence assumption so that  $X_i$ 's are now iid:  
 $X_i$ 's correlated by  $\rho$ .

$$\Rightarrow \text{Var}(\bar{X}) = \rho\sigma^2 + \frac{(1-\rho)}{n}\sigma^2$$

Ways to ensemble (for reducing variance)

- 1) Different algorithms
- 2) Different training sets
- 3) Bagging (eg Random forests)
- 4) Boosting (eg AdaBoost, xgboost)

Bagging: Bootstrap aggregation

Have a true population  $p$   
Training set  $S \sim p$   
Assume  $P = S$

Bootstrap samples  $Z \sim S$  (sample with replacement)

Eg  $Z_1, \dots, Z_M$ . Now train model  $G_m$  on  $Z_m$ .

$$G = \frac{\sum_{m=1}^M G_m}{M}$$

Bias-variance Analysis of this:  $\text{Var}(\bar{X}) = \rho\sigma^2 + \frac{(1-\rho)}{M}\sigma^2$

Bootstrapping is driving down  $\rho$ . More  $M \rightarrow$  less variance  
Bias slightly increases because of random subsampling

Ensembling

Take  $X_i$ 's which are RVs that are indep iid.  $\Rightarrow \text{Var}(X_i) = \sigma_i^2$   
 $\Rightarrow \text{Var}(\bar{X}) = \text{Var}\left(\frac{1}{n} \sum_i X_i\right) = \frac{\sigma^2}{n}$

## DTs + Bagging

DTs are high variance / low bias  $\Leftrightarrow$  Making this an ideal fit for bagging.

## RFs

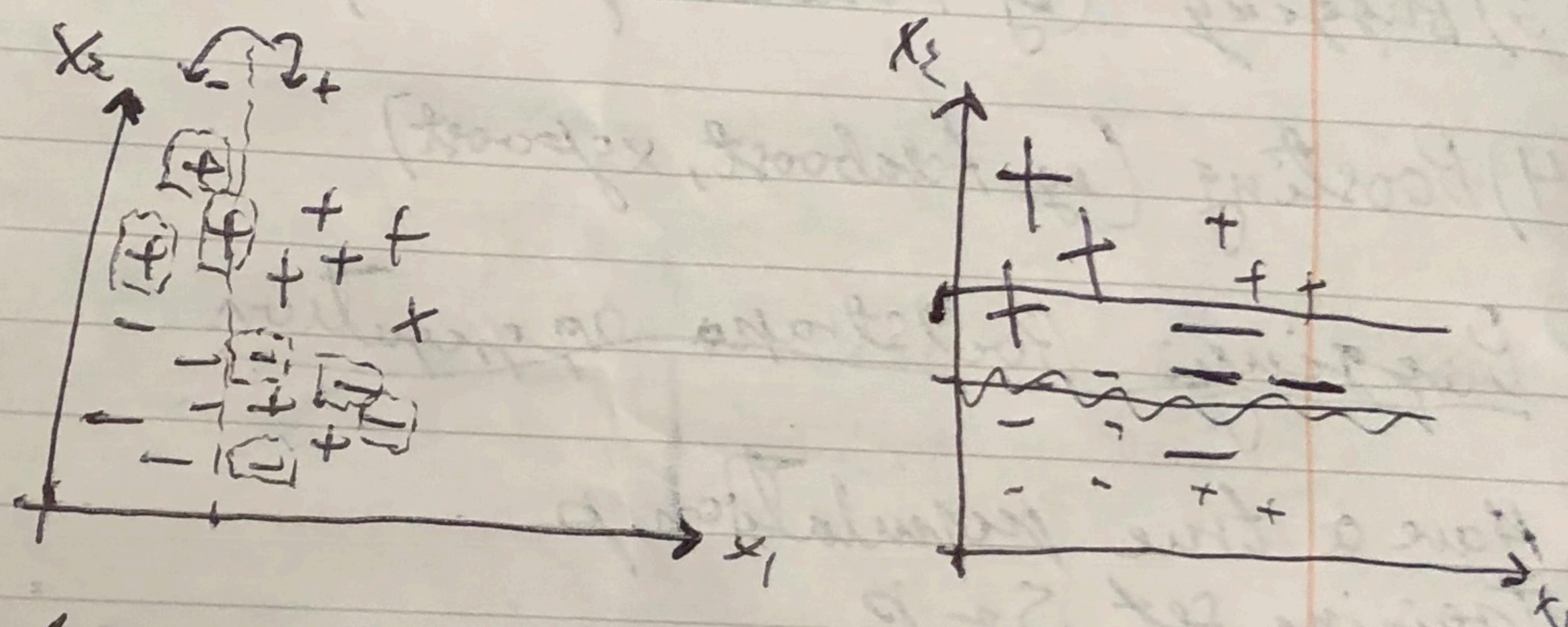
At each split, consider only a fraction of your total features:

This decreases  $P$ , because ~~the otherwise all individual tree would~~ from end caps splitting at same "high value" feature.

Decorrelates your models

## Boosting

Decrease bias  
Additive



Determine for classifier  $G_m$  a weight  $\alpha_m$  s.t.

$$\log \left( \frac{1 - \text{err}_m}{\text{err}_m} \right)$$

$$G_{\text{final}} = \sum_m \alpha_m G_m$$

A boost

Each  $G_m$  is trained on reweighted training set. (from step  $m-1$ )

? Better to read on wiki.

# Lecture 11

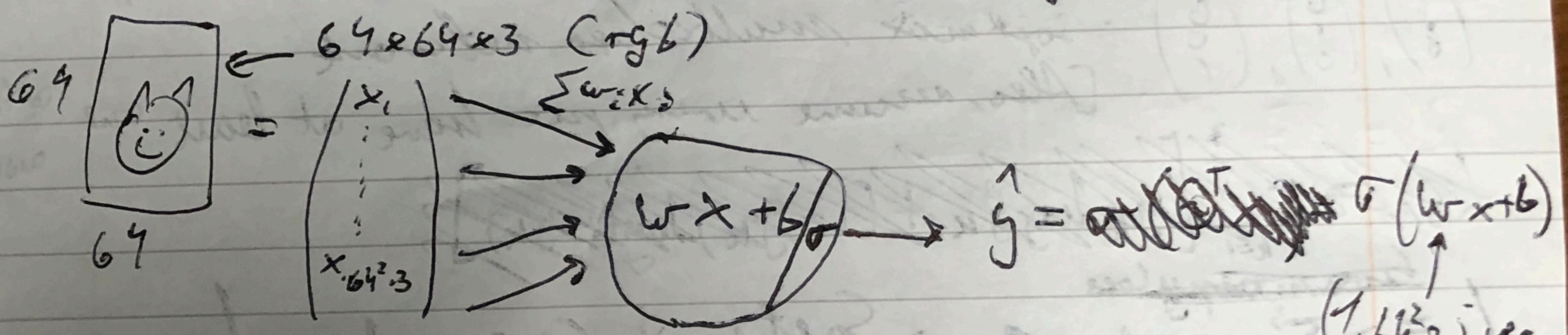
- I. Log reg.
- II. Neural networks.

## Deep learning

- Computational power
- data availability
- algorithms

### (1) Logistic regression

Goal: Find cats in images ( $1 \rightarrow$  presence of cat,  $0 \rightarrow$  absence of cat)

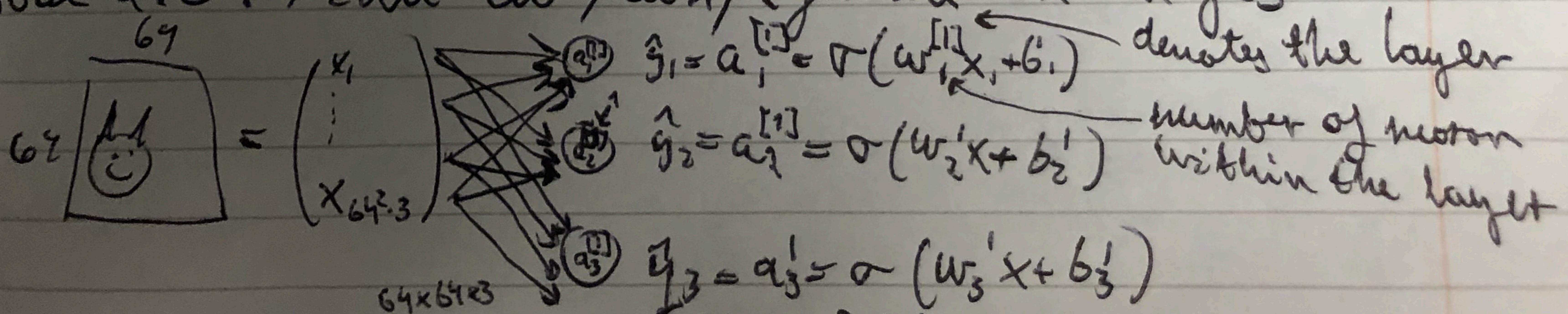


- i) initialise  $w, b$
  - ii) find the optimal  $w, b$
  - iii) use  $\hat{y} = \sigma(w \cdot x + b)$  to predict
- $L = -[y \log \hat{y} + (1-y) \log(1-\hat{y})]$  comes from MLE
- $w = w - \alpha \frac{\partial L}{\partial w}; b = b - \alpha \frac{\partial L}{\partial b}$

Neuron = linear + activation (eq. 1)

Model = architecture + parameters (eq. 2)

Goal 2.0: Find cat/lion/iguana in images



Needed dataset: Images + labels ~  $\begin{pmatrix} \text{Cat} \\ \text{Lion} \\ \text{Iguana} \end{pmatrix} = y$ .  $L_{3N} = -\sum_{k=1}^3 y_k \log \hat{q}_{k+1} + (1-y_k) \log(1-\hat{q}_{k+1})$

Note:

Neurons independent  $\Rightarrow$  network should be robust to multiple animals (e.g. cat + lion) in image.

Goal 3.0: + constraint "unique animal on an image".

$$q_j = \frac{e^{z_j}}{\sum_{k=1}^n e^{z_k}}$$

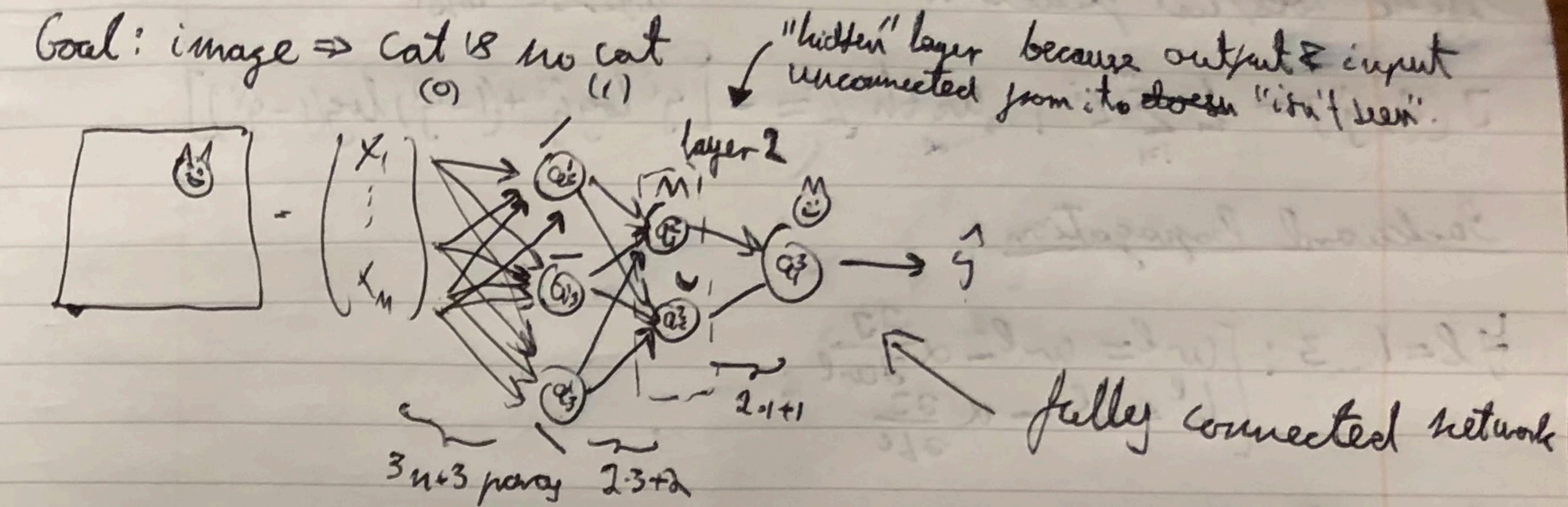
$$z_j = \sum_{i=1}^m w_{ji} x_i + b_j$$

This (softmax) scheme is no longer robust to multiple animals in the image, as the denominators must all equal 1.  
 $(\mathbf{z}), (\mathbf{q}), (\mathbf{y})$ : softmax multiclass network  
 [Also, assume every pic. have at least one animal]

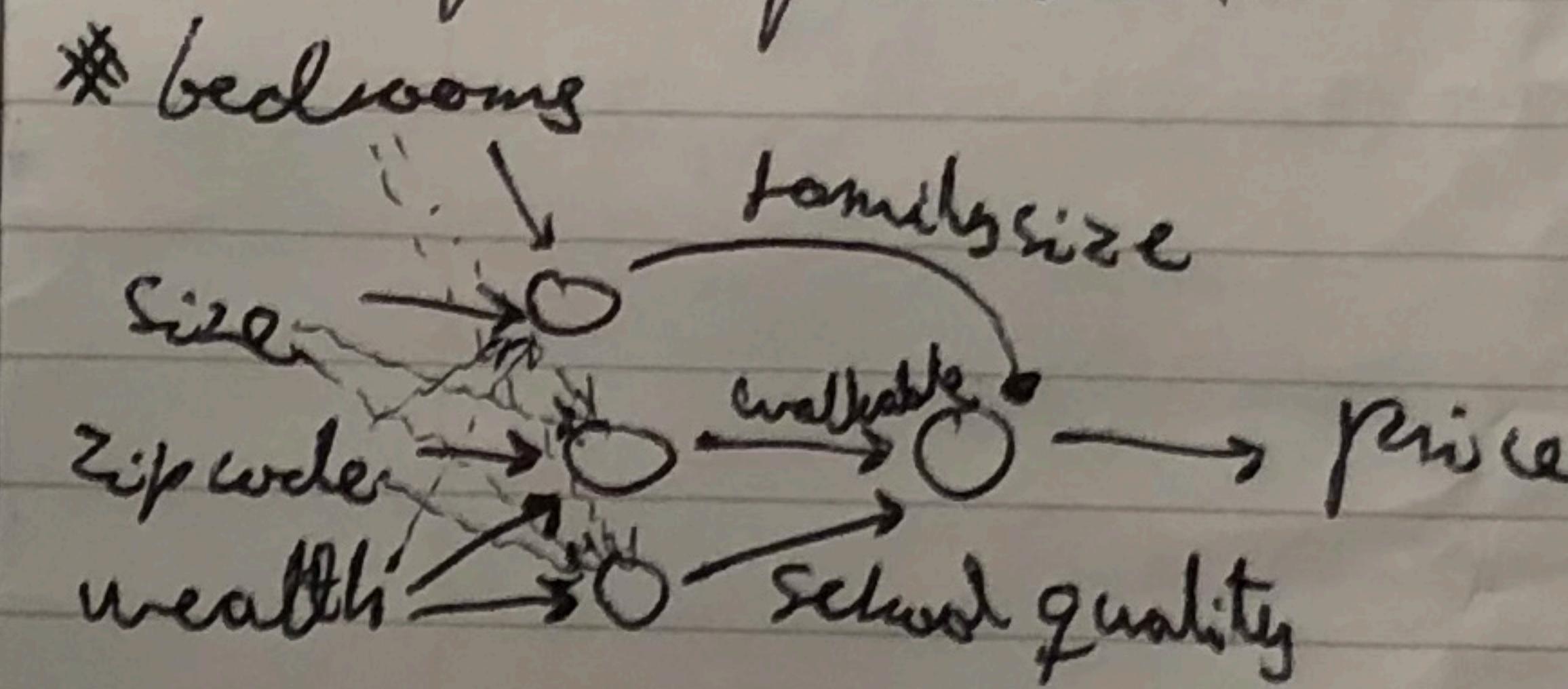
~~softmax loss~~  
 $L_{CE} = -\sum_{k=1}^m y_k \log q_k$  ← Softmax Cross entropy loss.

The Sidewall: to e.g. learn the age of a cat, one could do softmax over certain age ranges, or more appropriately perhaps, replace the sigmoid with something like linear regression (but such that the outputted age is  $\in [0, 35]$  years). Loss function then becoming

## II Neural Networks



House price prediction



end-to-end learning / black box model

Propagation equations

$$\mathbf{z}^1 = \mathbf{w}^1 \mathbf{x} + \mathbf{b}^1 \sim 3, 1$$

$$\mathbf{q}^1 = \sigma(\mathbf{z}^1) \sim 3, 1$$

$$\mathbf{z}^2 = \mathbf{w}^2 \cdot \mathbf{q}^1 + \mathbf{b}^2 \sim 2, 1$$

$$\mathbf{q}^2 = \sigma(\mathbf{z}^2) \sim 2, 1$$

$$\mathbf{z}^3 = \mathbf{w}^3 \cdot \mathbf{q}^2 + \mathbf{b}^3 \sim 1, 1$$

$$\mathbf{q}^3 = \sigma(\mathbf{z}^3) \sim 1, 1$$

- What happens for an input batch of m examples

$$\mathbf{X} = \begin{pmatrix} | & | & | \\ x^1 & x^2 & \dots & x^m \\ | & | & | \end{pmatrix}, \quad \mathbf{z}^1 = \underbrace{\mathbf{w}^1 \mathbf{X} + \mathbf{b}^1}_{3, n \times 1, m \rightarrow 3, 1 \text{ broadcast, i.e. repeated } m \text{ times}} \sim 3, m$$

$$\mathbf{z}^1 = \begin{bmatrix} z^{1,1} & \dots & z^{1,m} \end{bmatrix}$$

# Optimizing weights & intercepts

$\delta_{i,w}$   $\delta_{i,b}$

Define loss/cost function

logistic loss

$$J(\hat{y}, \hat{g}) = \frac{1}{m} \sum_{i=1}^m L^i, \text{ with } L^i = -[\hat{y}^i \log \hat{g}^i + (1-\hat{y}^i) \log(1-\hat{g}^i)]$$

Backward Propagation

$$\forall l=1\dots 3: \begin{cases} w^{l+1} = w^l - \alpha \frac{\partial J}{\partial w^l} \\ b^{l+1} = b^l - \alpha \frac{\partial J}{\partial b^l} \end{cases}$$

$$\frac{\partial J}{\partial w^3} = \underbrace{\frac{\partial J}{\partial a^3}}_{\frac{\partial J}{\partial z^3}} \frac{\partial a^3}{\partial z^3} \frac{\partial z^3}{\partial w^3}$$

$$\frac{\partial J}{\partial w^2} = \underbrace{\frac{\partial J}{\partial z^3}}_{\frac{\partial J}{\partial z^2}} \frac{\partial z^3}{\partial a^2} \frac{\partial a^2}{\partial z^2} \frac{\partial z^2}{\partial w^2}$$

$$\frac{\partial J}{\partial w^1} = \underbrace{\frac{\partial J}{\partial z^2}}_{\frac{\partial J}{\partial z^1}} \frac{\partial z^2}{\partial a^1} \frac{\partial a^1}{\partial z^1} \frac{\partial z^1}{\partial w^1}$$

## Lecture 12

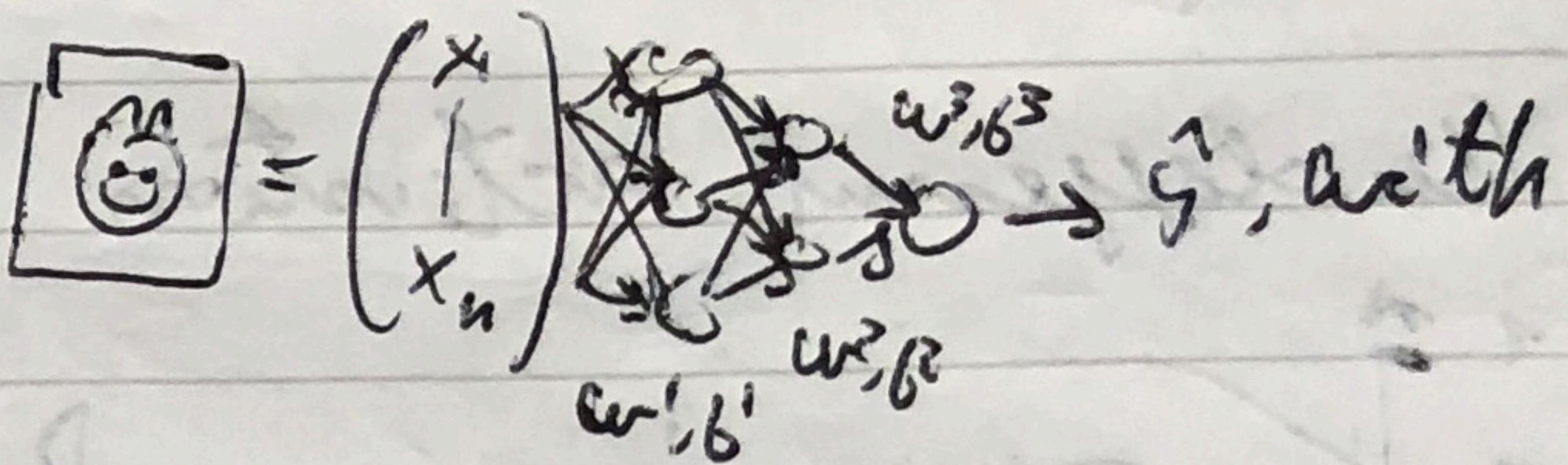
Outline:

I Logreg from a NN mindset ✓  
II Neural Networks  
→ backprop

III Improving your NNs

### Backprop

Cost function  $J(y, \hat{y}) = \frac{1}{m} \sum_{i=1}^m L^i(y, \hat{y})$



$$L^i = -[y^i \log \hat{y}^i + (1-y^i) \log(1-\hat{y}^i)]$$

$$\text{Update: } w^{l+1} = w^l - \alpha \frac{\partial J}{\partial w^l}$$

~~$$\frac{\partial J}{\partial w^3} \Rightarrow \frac{\partial L}{\partial w^3} = - \left[ y^i \frac{1}{w^3} \log(w^3 a^2 + b^3) + (1-y^i) \frac{1}{w^3} \log(1-w^3 a^2 - b^3) \right]$$~~

$$\left. \begin{aligned} & - \left[ y^i \frac{1}{w^3} \sigma^3 (1-\sigma^3) a^2{}^T + (1-y^i) \frac{1}{1-\sigma^3} \cdot (-1)^3 (1-\sigma^3) \right] \\ & \text{use that } \sigma'(x) = \sigma(x)(1-\sigma(x)) \\ & \cdot a^2{}^T \end{aligned} \right]$$

$$= - \left[ y^i (1-\sigma^3) a^2{}^T - (1-y^i) \sigma^3 a^2{}^T \right] = - \left[ y^i a^2{}^T - \sigma^3 a^2{}^T \right]$$

$$= -(y^i - \sigma^3) a^2{}^T \quad \underbrace{\frac{\partial z^3}{\partial w^3}}_{\frac{\partial L}{\partial w^3}} \quad \frac{\partial L}{\partial w^3} = \frac{\partial L}{\partial z^3} \frac{\partial z^3}{\partial w^3} \Rightarrow \frac{\partial L}{\partial w^3} = -f y^i - a^2$$

$$\frac{\partial J}{\partial w^3} = \frac{1}{m} \sum_{i=1}^m (y^i - \sigma^3) a^2{}^T$$

$$\frac{\partial L}{\partial w^3} = \underbrace{\frac{\partial L}{\partial z^3}}_{\frac{\partial L}{\partial w^3}} \underbrace{\frac{\partial z^3}{\partial a^3}}_{\frac{\partial a^3}{\partial w^3}} \underbrace{\frac{\partial a^3}{\partial z^3}}_{\frac{\partial z^3}{\partial w^3}} \underbrace{\frac{\partial z^3}{\partial a^2}}_{\frac{\partial a^2}{\partial w^3}} \underbrace{\frac{\partial a^2}{\partial z^2}}_{\frac{\partial z^2}{\partial w^3}}$$

$$= (\sigma^3 - y) a^2 (1-\sigma^3)$$

$$\frac{\partial L}{\partial w^3} = \underbrace{(a^3 - y)}_{1,1} \underbrace{w^3{}^T}_{2,1} \underbrace{a^2}_{2,1} \underbrace{(1-a^2)}_{1,3} \underbrace{a^1{}^T}_{2,1} = w^3{}^T \underbrace{\sigma^3}_{2,1} \underbrace{a^2}_{2,1} \underbrace{(1-\sigma^3)}_{1,3} \underbrace{(a^3 - y)}_{1,1} \underbrace{a^1{}^T}_{1,3}$$

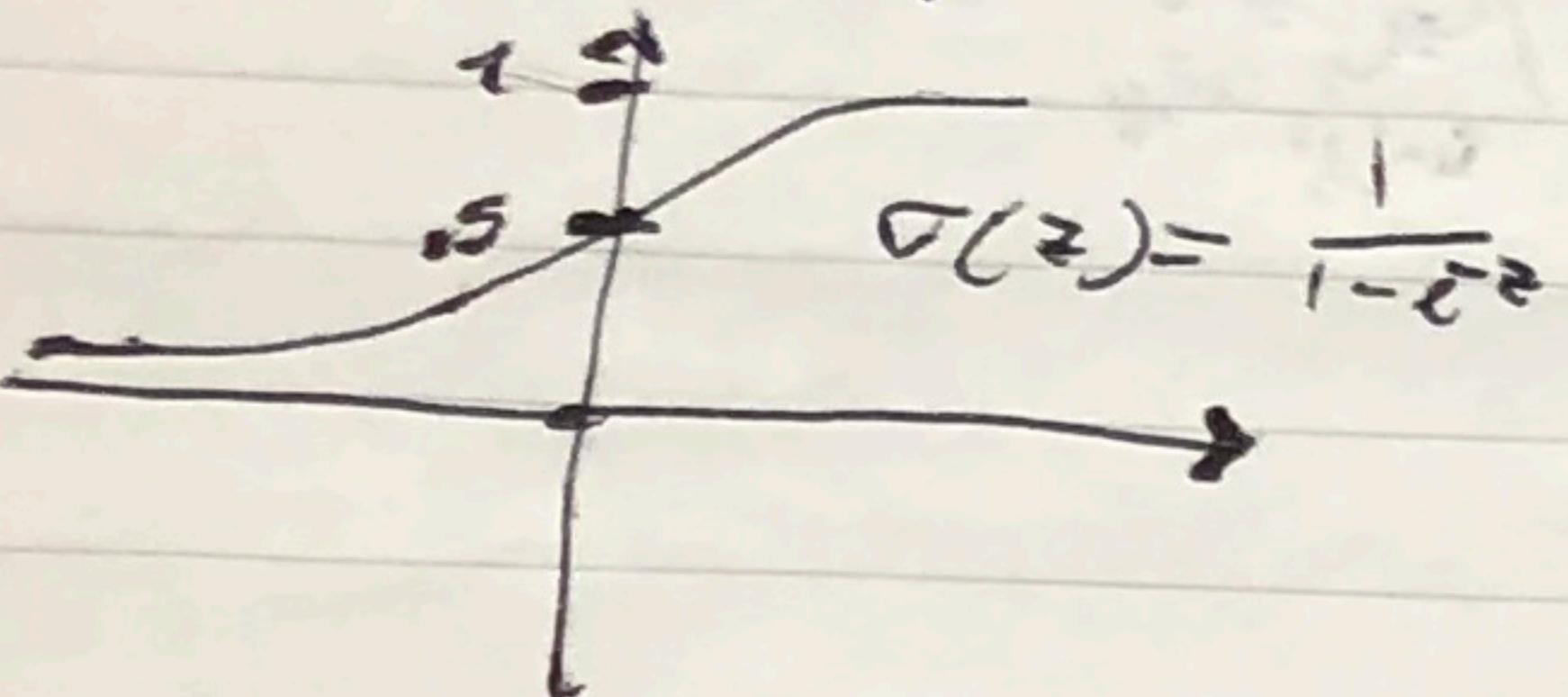
Applying backpropagation

only probabilities due to softmax activation function

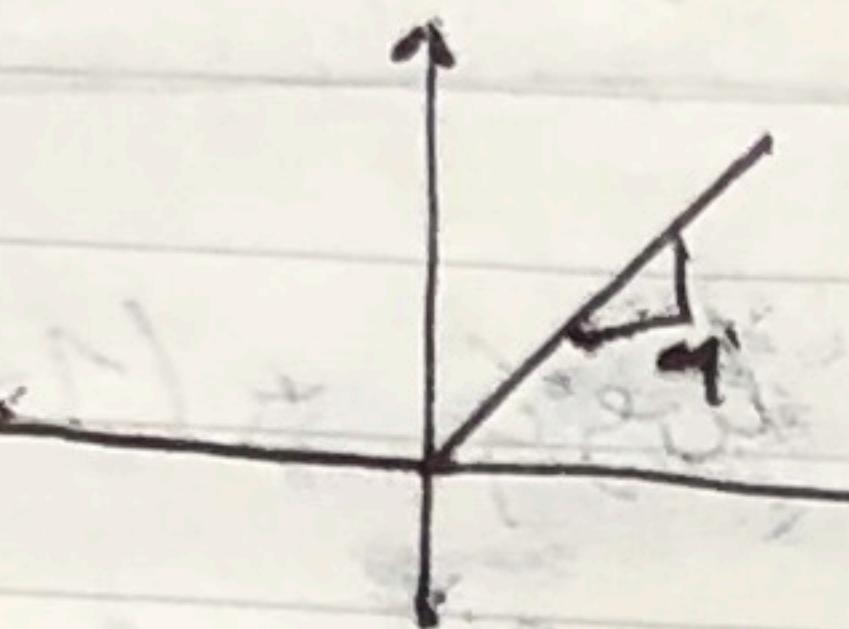
$$\frac{\partial J}{\partial w^2} = \frac{1}{m} \sum_{i=1}^m \frac{\partial L}{\partial w^2}$$

### III Improving your NNs:

Trick 1: Use different activation functions

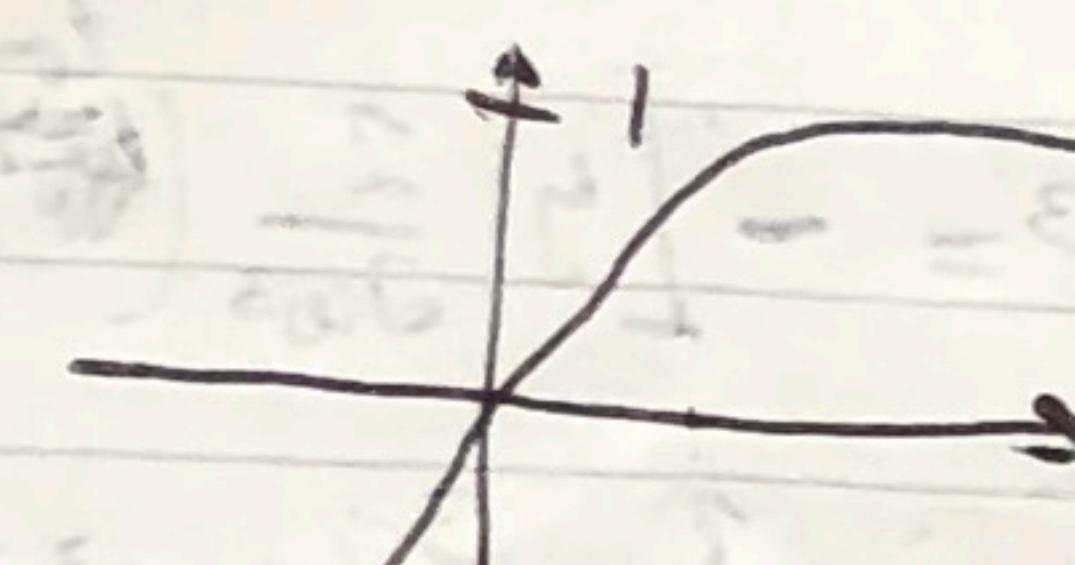


$$\sigma(z) = \frac{1}{1+e^{-z}}$$



$$\text{ReLU}'(z) = \begin{cases} 0, & z \leq 0 \\ 1, & z > 0 \end{cases}$$

$$\tanh z = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$



$$\tanh'(z) = 1 - \tanh^2 z$$

General: problem here is that the gradient vanishes at high/low z.  
Same issue for the tanh

Why do we need activation functions?

Activations = Id  $z \mapsto z$

$$y^l = a^3 = z^3 = w^3 a^2 + b^3$$

$$= w^3(w^2 a^2 + b^2) + b^3 = w^3 w^2 (w^1 x + b') + w^3 b'^2 + b^3$$

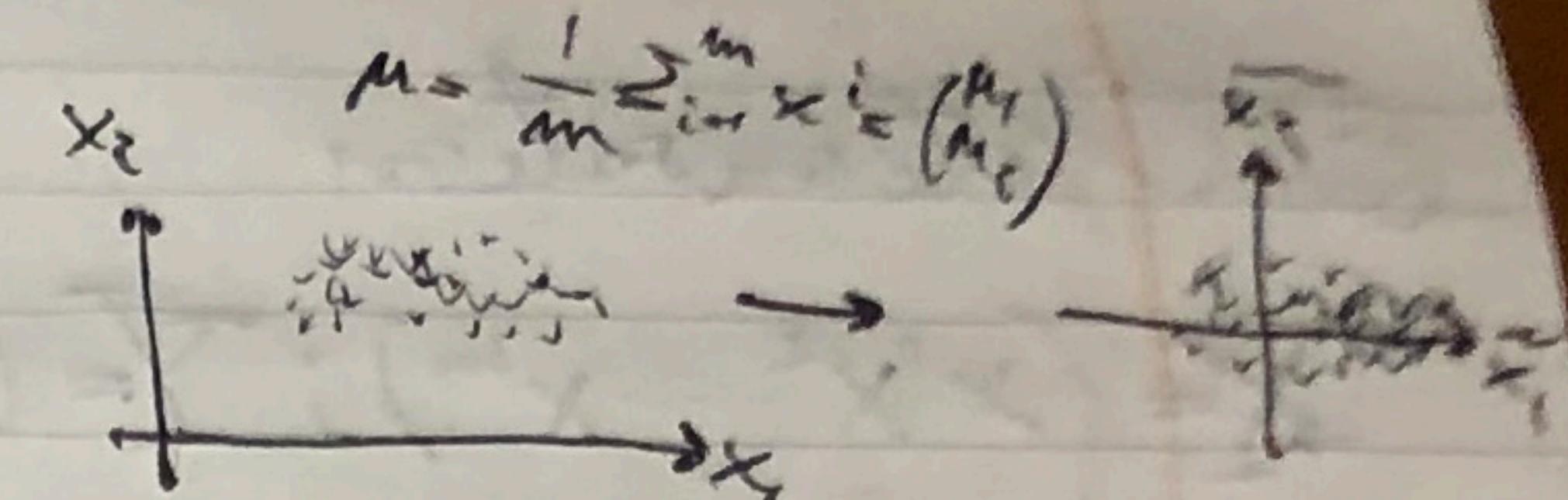
$$= Wx + b, \text{ with } W = w^3 w^2 w^1, b = w^3 w^2 b' + w^3 b'^2 + b^3$$

Without activation functions, the NN always becomes simply linear regression.

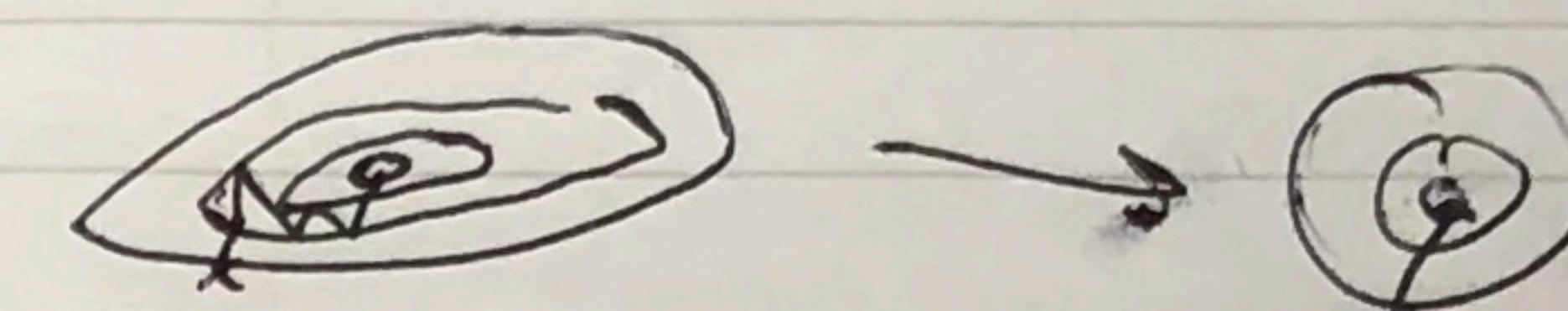
B Initialization methodology

Normalizing your input:  $\tilde{x} = \frac{x_i - \mu}{\sigma}$

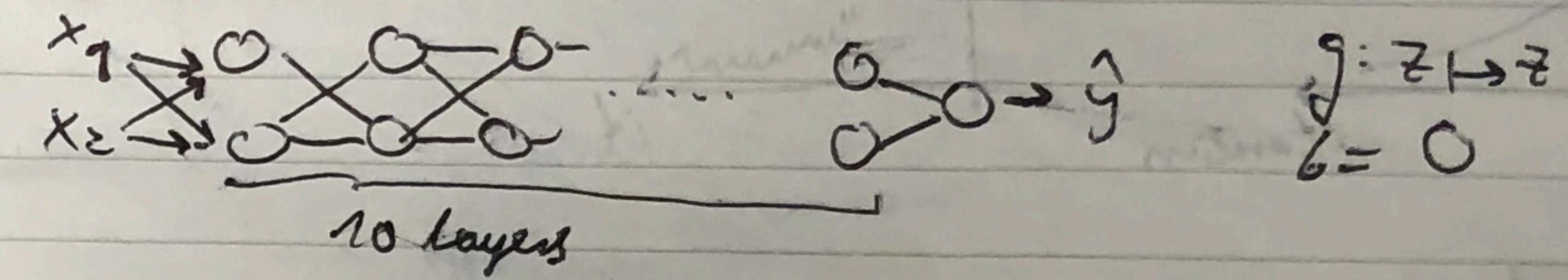
$$\sigma^2 = \frac{1}{m} \sum_i x_i^2 \rightarrow \tilde{x}_i = \frac{x_i - \mu}{\sqrt{\sigma^2}}$$



Effect of normalization on loss functions:



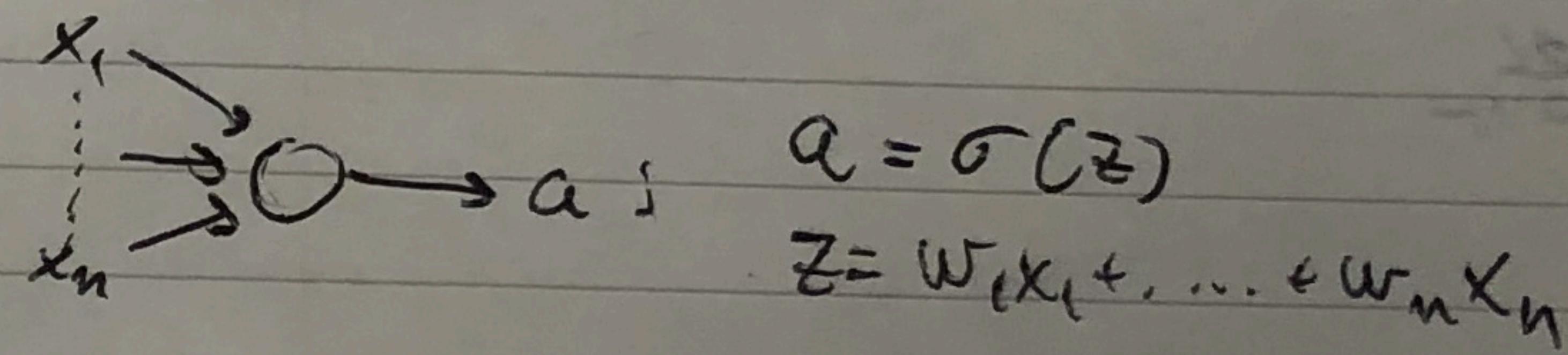
Vanishing/Exploding gradients



$$g = w^L \cdot a^{L-1} = w^L w^{L-1} \cdot a^{L-2} = \dots = w^L w^{L-1} \dots w^1 x$$

$$w^l = \begin{pmatrix} 1.5 & 0 \\ 0 & 1.5 \end{pmatrix} \Rightarrow g = \begin{pmatrix} 1.5^L & 0 \\ 0 & 1.5^L \end{pmatrix} x$$

Example with 1 neuron



⇒ large n requires small w\_i ⇒ w\_i ~ 1/n

W init: initialise w\_i up random.sample(shape) \* np.sqrt(1/(n-1)) when using sigmoids. For ReLU: use 2/n^(1/2) instead.

Xavier initialization:  $w \sim \sqrt{1/n^{1/2}}$  for tanh

He initialization:  $w^l \sim \sqrt{\frac{2}{n^l + n^{l-1}}}$

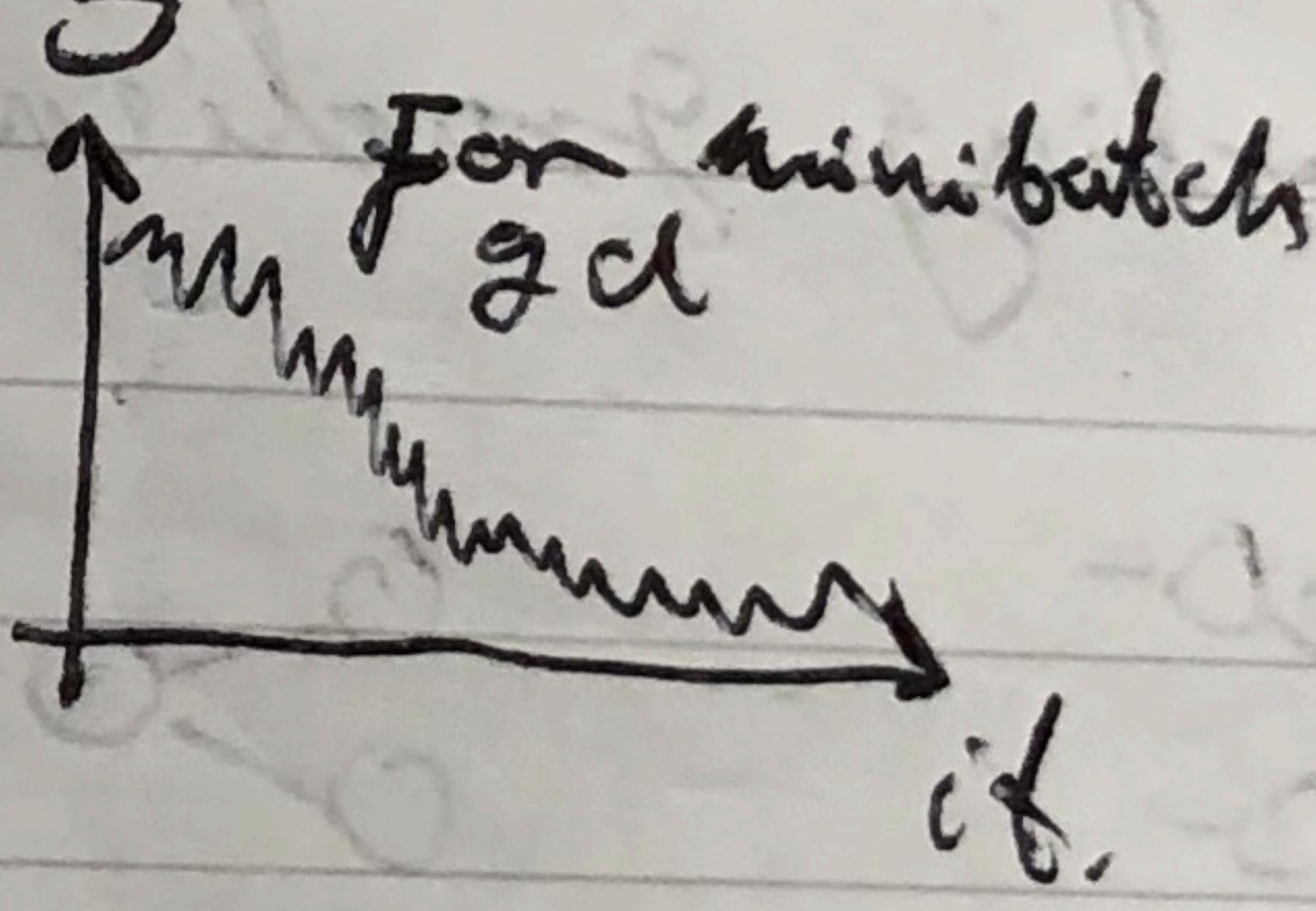
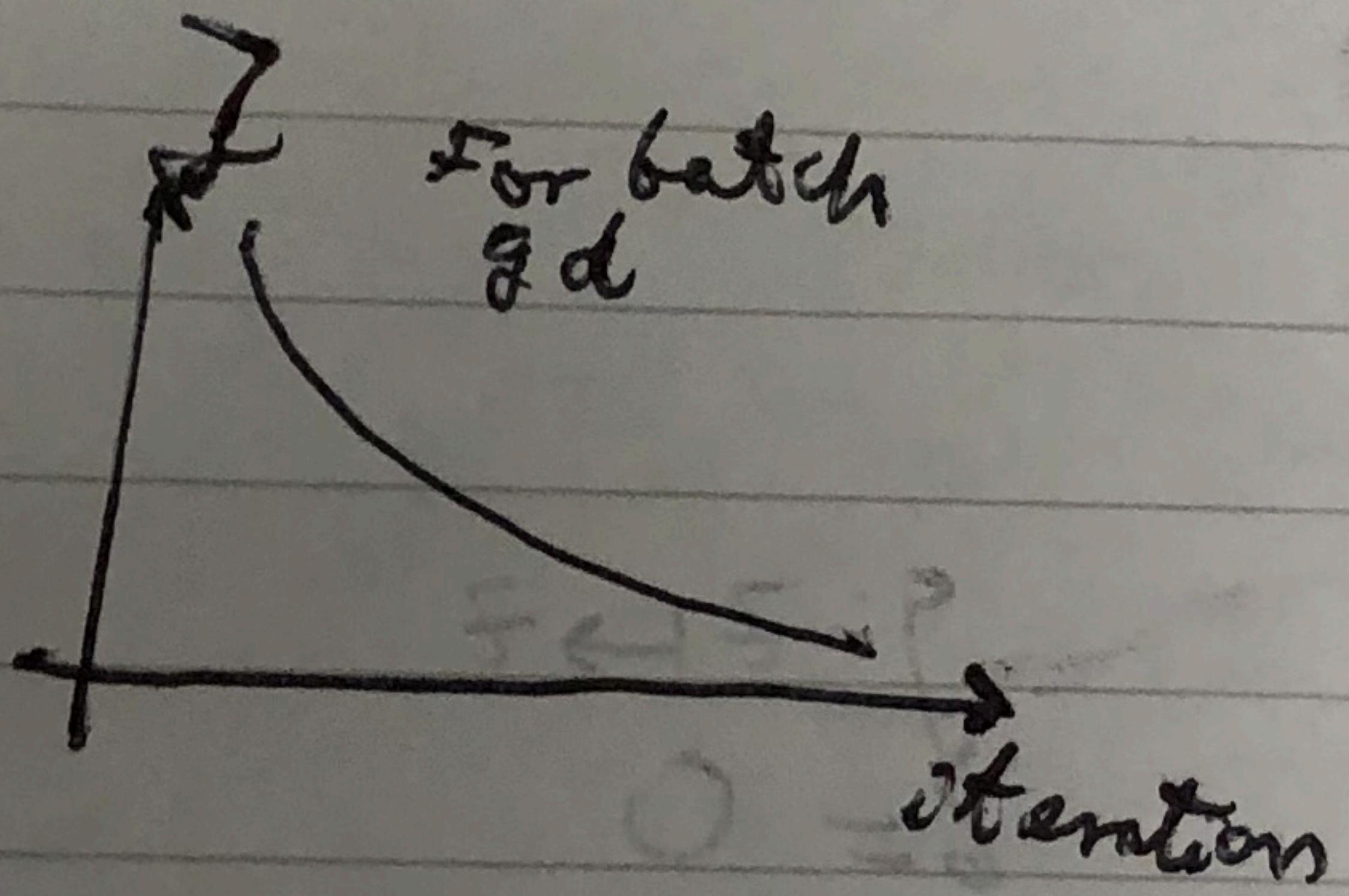
## ① Optimisation

$$\underline{X} = (x^1 | x^2 | \dots | x^m); \underline{Y} = (y^1 | \dots | y^m); x = (x^{1,1} | \dots | x^{1,m}) \in \\ y = (y^{1,1} | \dots | y^{1,m})$$

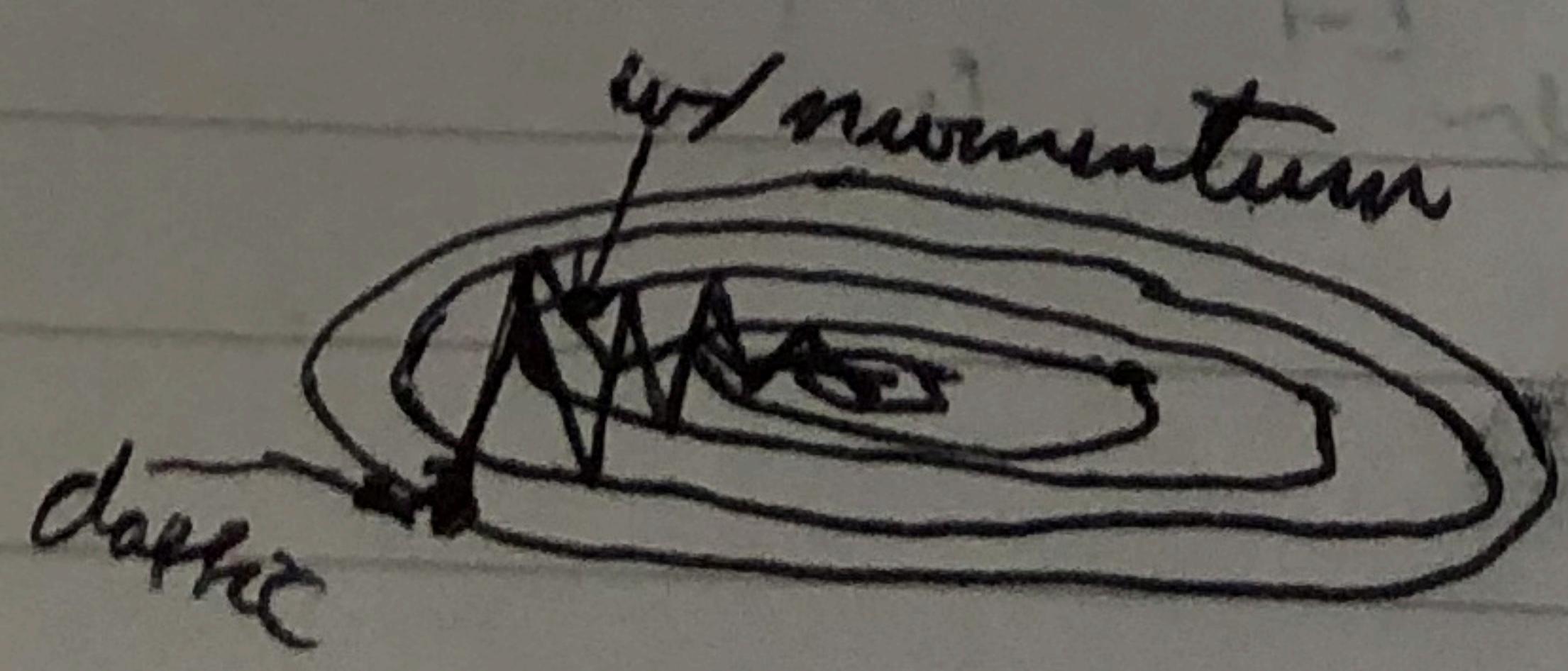
Algo: for iteration  $t=1 \dots T$  select batch  $(x^{t,1}, y^{t,1})$ , forward pass backward batch.

$$J = \frac{1}{1000} \sum_{i=1}^{1000} L^i \text{ update } w^l, b^l \forall l.$$

graph:



"GD + momentum" algo:



Closer. Preyer  $\beta$  smaller,  $\alpha$  larger

$$\begin{cases} v = \beta v + (1-\beta) \frac{\nabla L}{\nabla w} \\ w = w - \alpha v \end{cases}$$