

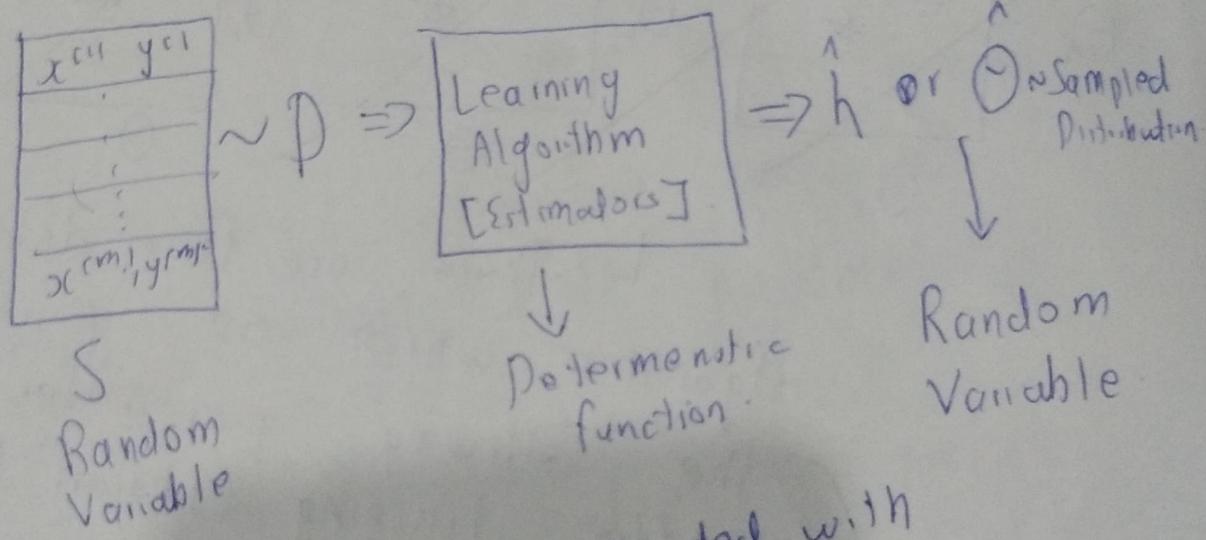
* Assumptions:-

① Data Distribution - D

$$(x, y) \sim D$$

\hookrightarrow There is a data generating distribution

② All the samples are sampled independently



• There is a distribution associated with parameters

\hookrightarrow There exist some θ^* or h^* which in a sense are true parameters we wish to be the parameters

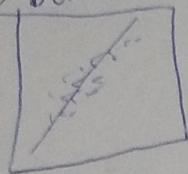
$\bullet \theta^* \text{ or } h^* \rightarrow$ Not random (It is a unknown constant).

(99)

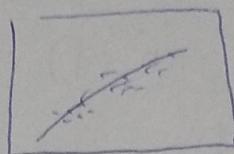
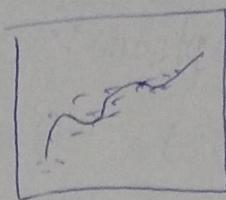
④ Bias-Variance

⑤ Data-View

Reg

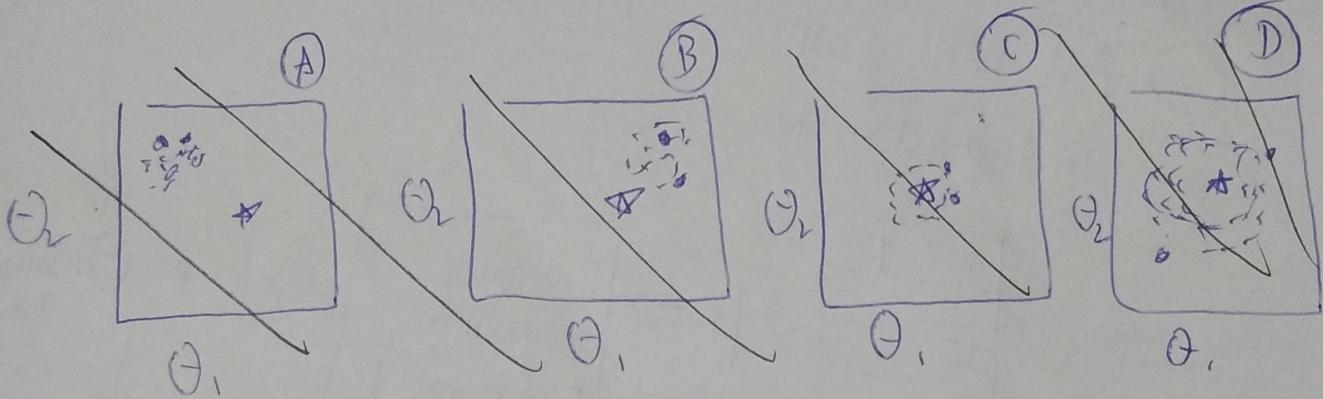


Underfit

Just
Right

Overfit

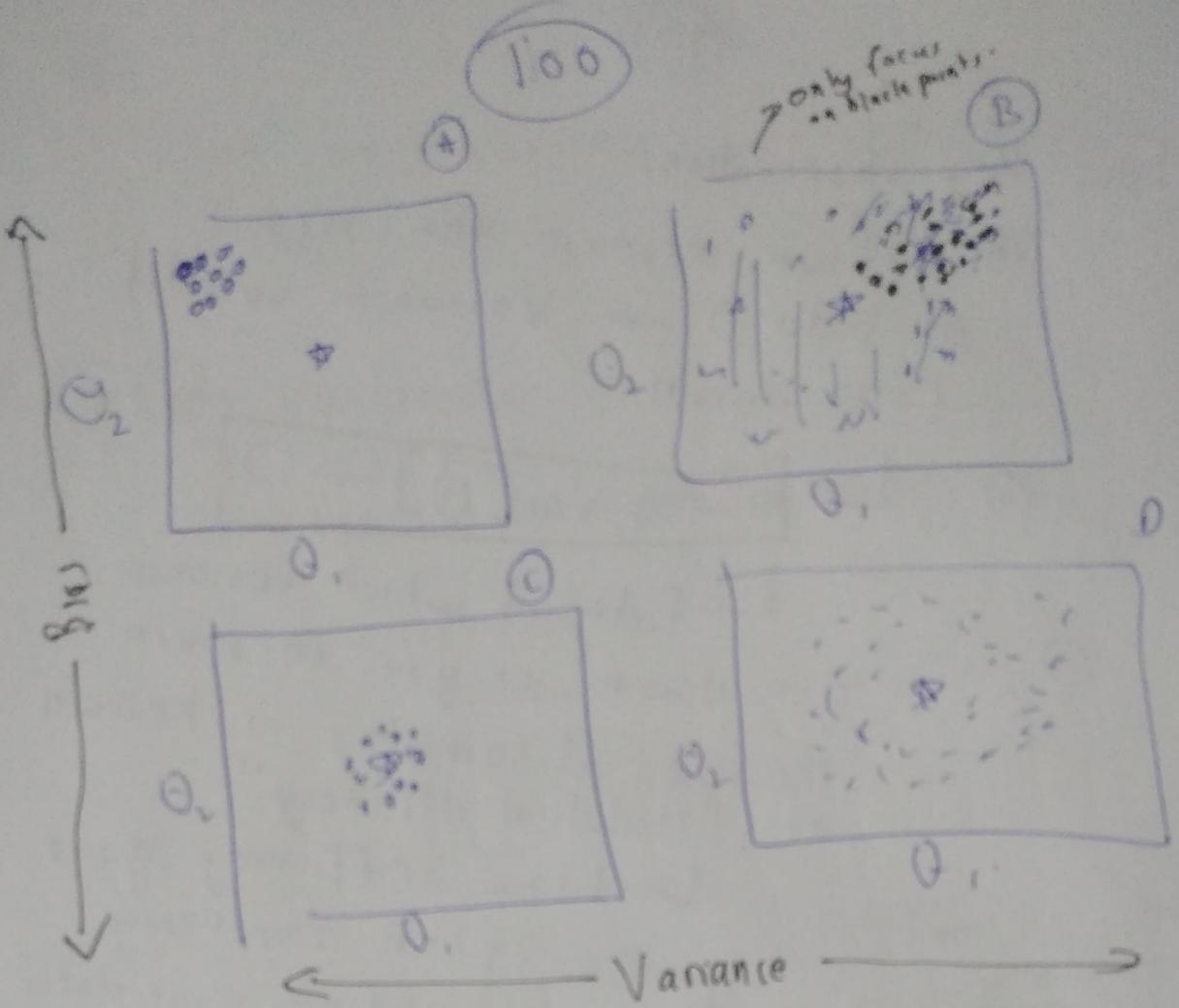
⑥ Parameter View of Bias - Variance.



↳ There are four different algorithm which are generating the following parameter value after training (shown by Blue dots).

↳ Blue dots are samples from each hypothesis distribution.

→ Refer to diagram on next page.



- Bias → Whether parameter are centered around true parameter
(First moment)
- Variance → How dispersed is the sampling distribution of parameter from hypothesis sample.
(Second moment)

(101)

*)

Few Key observations:-

① As we increase the size of Data D \rightarrow Variance would come down.

② $\hookrightarrow [m \rightarrow \infty, \text{Var}[\hat{\theta}] \rightarrow 0]$

\hookrightarrow Rate at which variance goes down as you increase the size of samples is known as "statistical efficiency".

\hookrightarrow How efficiency of your algorithm is in squeezing out information from your data.

③ $\hat{\theta} \rightarrow \theta^*, m \rightarrow \infty$:- You call such algorithm consistent.

$\circ E[\hat{\theta}] = \theta^* \rightarrow$ Unbiased Estimator.

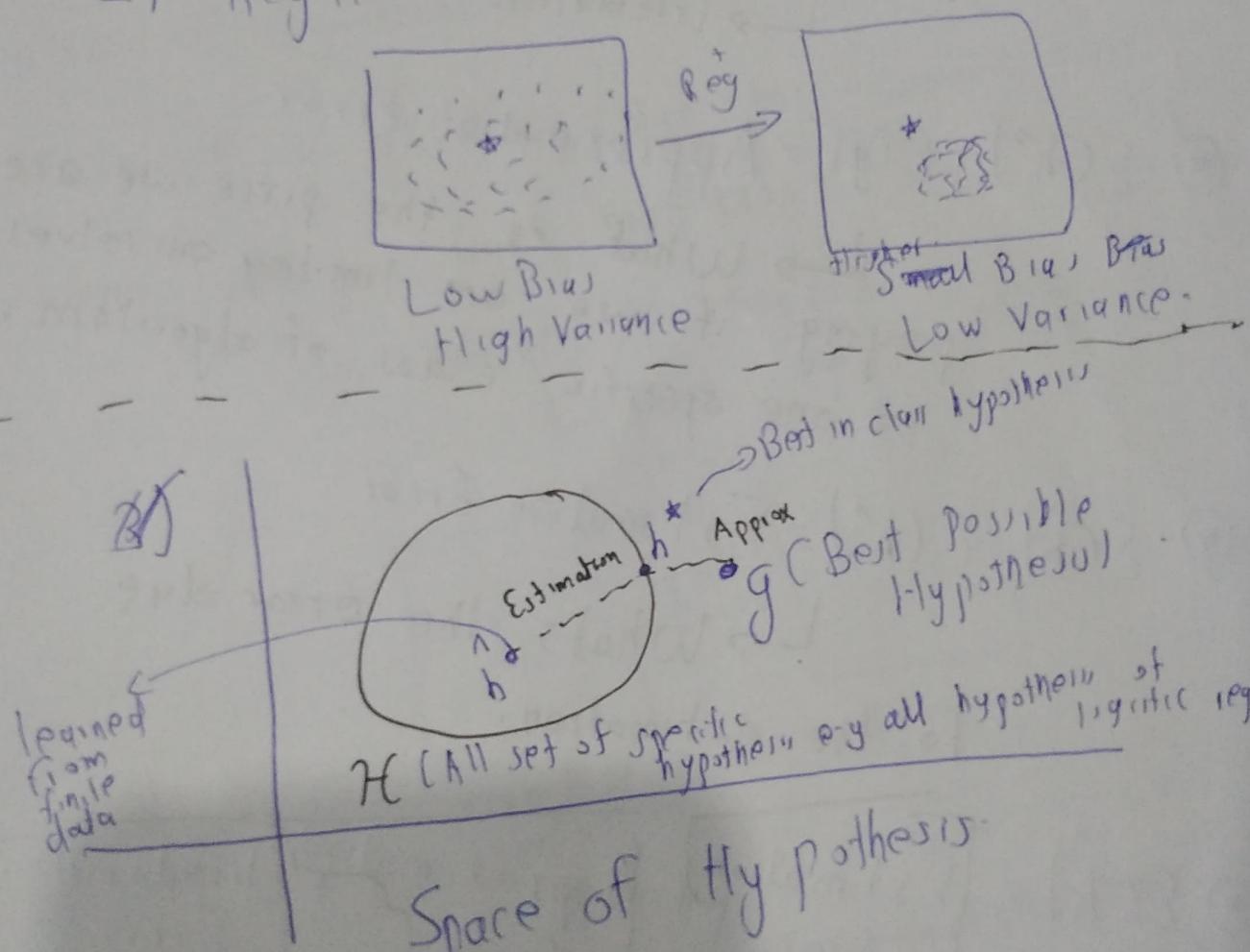
④ High Bias won't come down no matter how much data you provide to your algorithm.

Bias and Variance are independent of each other.

④ Fighting Variance:-

1) Increase amount of Data.

2) Regularization:



④ $E(h)$ — Risk / Generalization Error

$$\hookrightarrow = E_{(x,y) \sim D} [2\{h(x) \neq y\}] \quad (\text{Fraction of all your mistakes})$$

• $\hat{\epsilon}(h)$: Empirical Risk

$$\frac{1}{m} \sum_{i=1}^m \{ h(x^{(i)}) \neq y^{(i)} \}$$

④ $\epsilon(g) = \text{Bayes Error}$ (Best opt possible hypothesis error)
 ↳ Irreducible Error

⑤ $\epsilon(h^*) - \epsilon(g) = \text{Approximation Error}$

↳ What is the price we are paying if we are limiting ourselves to one specific class of algorithm.

⑥ $\epsilon(h) - \epsilon(h^*) = \text{Estimation Error}$

↳ What's the error due to the estimation.

$$\epsilon(h) = \boxed{\text{Estimation Error}} + \boxed{\text{Approximation Error}} + \boxed{\text{Irreducible Error}}$$

Limited Data (In class)
 (Error)

What class
 of algorithm

No way
 you can
 reduce

(109)

$$E(h) = \text{Estimation Error} + \text{Approximation Error} + \text{Irreducible Error}$$

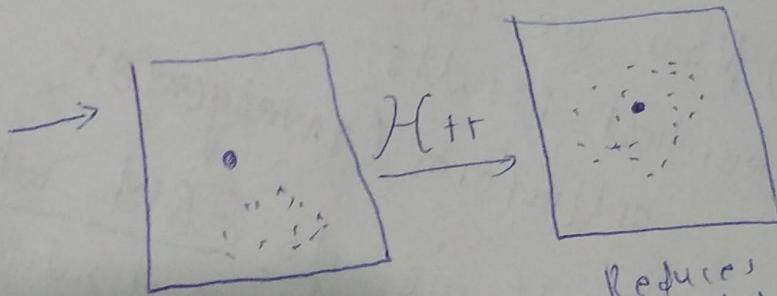
↓
 Estimation Variance ↓
 Variance + Estimation bias ↓
 + Bias + Irreducible

→ Bias why \hat{h} is far from g :- either due to class of algorithm or some other reason.

→ Variance is generally due to estimation error.

→ How do you fight high bias?

→ Make H bigger
 ↳ More complex algorithms.



high bias
some variance

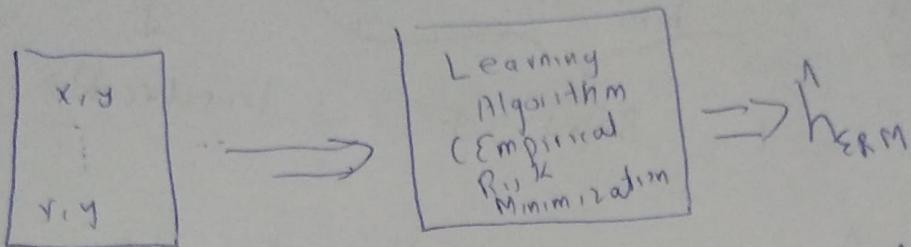
Reduces
Bias but
also increase
your variance

→ Regularization decreases ~~increases~~
 decrease would effectively make
 H larger which will further reduce
 the bias

Empirical Risk Minimization:-

(Minimize)

④ Is a learning algorithm



$$\textcircled{4} \quad \hat{h}_{ERM} = \arg \min_{h \in \mathcal{H}} \frac{1}{m} \sum_{i=1}^m \ell\{\hat{h}(x_i)^T y_i\}$$

\hookrightarrow ERM is an hypothesis that tries to find a hypothesis from class of hypotheses that minimizes the training error. (Increasing the training accuracy).

Now we can come up with more theoretical results:-

Theoretical results:-

① If we do ERM what does it say about generalization error?

② $\hat{\epsilon}(h)$ vs $\epsilon(h)$.

③ How does the generalization error of our learned hypothesis compare to the best possible generalization error in that hypothesis space?

④ $\epsilon(h)$ vs $\epsilon(h^*)$

For above two questions we are going to use two tools:-

(1) Union Bound:

A_1, \dots, A_n (Needs not be independent)

$$P(A_1 \cup A_2 \cup A_3 \cup \dots \cup A_n) \leq P(A_1) + P(A_2) + \dots$$

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

$P(A_n)$

(2) Hoeffding Inequality

Let $Z_1, Z_2, \dots, Z_m \sim \text{Bernoulli}(\phi)$

Sample of $O(3)$

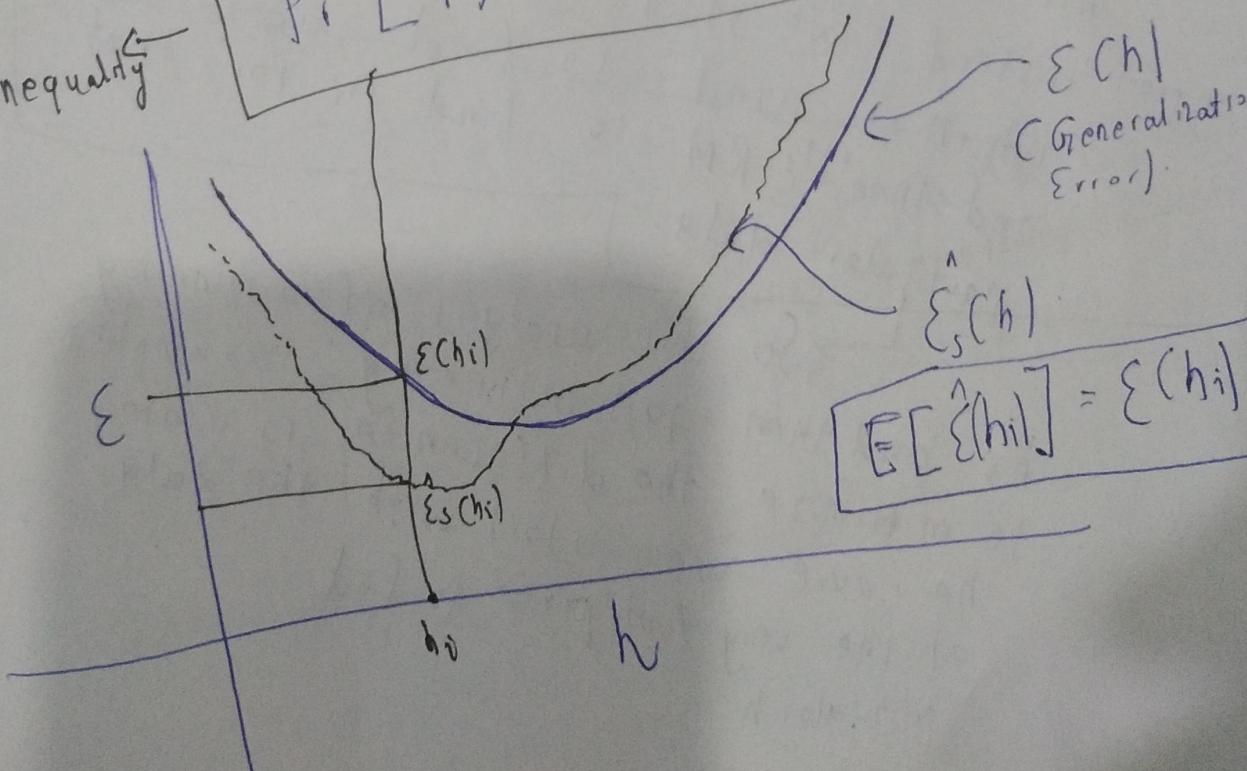
$$\hat{\phi} = \frac{1}{m} \sum_{i=1}^m Z_i \text{ (Estimator)}$$

Let $\gamma > 0$ [margin]

$$P_e[|\hat{\phi} - \phi| > \gamma] \leq 2 \exp(-2\gamma^2 m)$$

Inequality

④



$$E[\hat{\epsilon}(h)] = \epsilon(h)$$

(107) Now Applying Hoeffding's inequality:-

④ $\Pr[|\hat{E}_{\text{S}(\mathbf{h})} - E_{\text{Hill}}| > \gamma] \leq 2 \exp(-2\gamma^2 m)$

↳ The gap between generalization error (Error of hypothesis \mathbf{h}_i on all data) and empirical error (Error of hypothesis \mathbf{h}_i on subsample of data) is bounded by this inequality (expression).

↳ As m increases empirical and generalization are going to come near to each other i.e their gap would reduce.

⑤ However, the problem with some hypothesis \mathbf{h}_i and then average across \mathbf{S}_{out} and then but in practice instead we start with some data and then used ERM to find some \mathbf{h}_i .
particular data

↳ So this analysis of assuming \mathbf{h}_i and then just increasing examples to minimize the difference is useless because we start with the data at the very first place to find

→ To fix the above problem, we would generalize the above inequality to all classes of hypothesis to account for all \mathcal{H} .

→ We will look at the difference of the whole class \mathcal{H} rather than particular h_i of \mathcal{H} .

→ This is called uniform convergence because we want to see how empirical risk curve converge uniformly across all classes to generalization risk curve.

→ Now the results would look different for two cases for all \mathcal{H} if we are to show the same bound that we show for particular h_i :

① Finite \mathcal{H} hypothesis classes:-

$$\textcircled{2} |\mathcal{H}| = K$$

$$\textcircled{1} \Pr[\exists h \in \mathcal{H} \mid \hat{\epsilon}_s(h) - \epsilon(h) > \gamma] \leq K \cdot \exp(-2\gamma^2 m)$$

Flipping many steps are missing but done through two parts

$$\textcircled{2} \Rightarrow \Pr[\forall h \in \mathcal{H} \mid \hat{\epsilon}_s(h) - \epsilon(h) < \gamma] \geq 1 - \underbrace{2K \exp(-2\gamma^2 m)}_8$$

~~key $\gamma = \sqrt{2K \exp(-2\gamma^2 m)}$~~

(109)

$$\text{Let } \delta = 2k \exp(-2\gamma^2 m)$$

$\delta \rightarrow$ Prob. of error
 Relation \rightarrow Margin of error

$m \rightarrow$ Sample size.

→ What can we do with this relationship
 is that we can fix any two and solve

for the third:

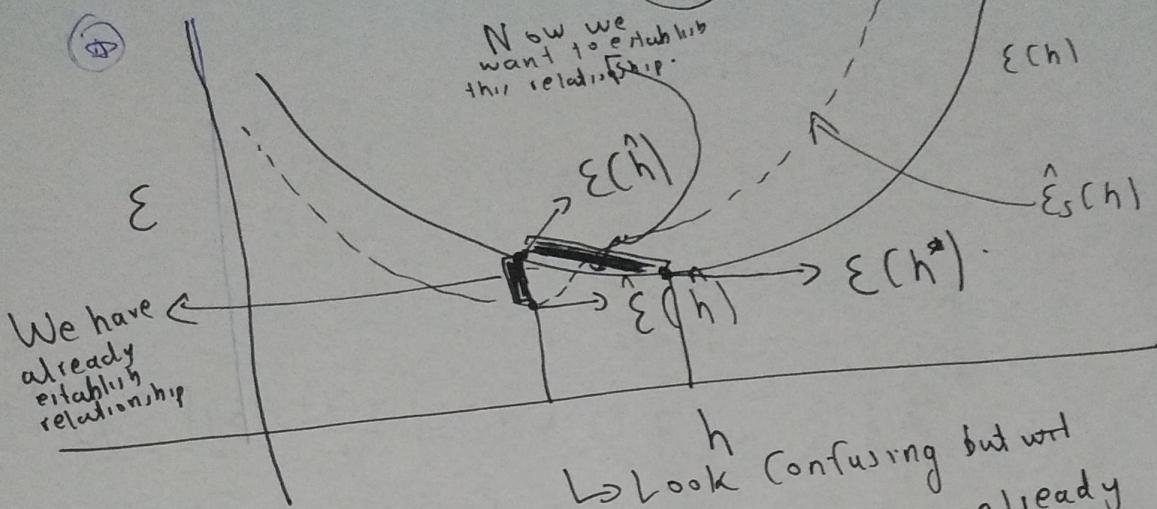
Example: Fix $\delta, \gamma > 0$

$$m \geq \frac{1}{2\gamma^2} \log \frac{2k}{\delta} \quad (\text{Sample complexity descent}).$$

④ This shows that the probability
 (for all classes of hypothesis) that difference
 between true empirical risk and generalization
 risk is less than γ is bounded by
 $1 - \delta$ as long as number of training
 examples are greater than expression

shown

⑤ As you increase m empirical risk
 and generalization risk becomes closer
 for H .



↪ Look confusing but w^h
 basically we have already established relationship between generalization and empirical error. Now we want to establish relationship between generalization error of a hypothesis and generalization error of best possible hypothesis from H.

$$\begin{aligned}
 \textcircled{1} \quad & \hat{\epsilon}(h) \leq \hat{\epsilon}(h) + \gamma \text{ (Already know)} \\
 & \leq \hat{\epsilon}(h^*) + \gamma \text{ (? ? ?)} \\
 & \leq \hat{\epsilon}(h^*) + 2\gamma. \text{ (Relationship established.)}
 \end{aligned}$$

$$\Rightarrow \text{With Probability } 1 - \delta, \text{ training size } m.$$

$$\hat{\epsilon}(h) \leq \epsilon(h^*) + 2\sqrt{\frac{1}{2m} + \log \frac{2K}{\delta}}$$

Can get this from previous step

2
III

* VC-Dimension:-

VC(\mathcal{H}) = Some Number

which tell how expressive

the hypothesis is.

$$\epsilon(\hat{h}) \leq \epsilon(h^*) + O\sqrt{\frac{VC(\mathcal{H})}{m} \log \left[\frac{m}{VC(\mathcal{H})} + \frac{1}{\epsilon^2} \right]}$$

Margin for infinite classes h in \mathcal{H} .