# Statistical Learning Uniform Convergence
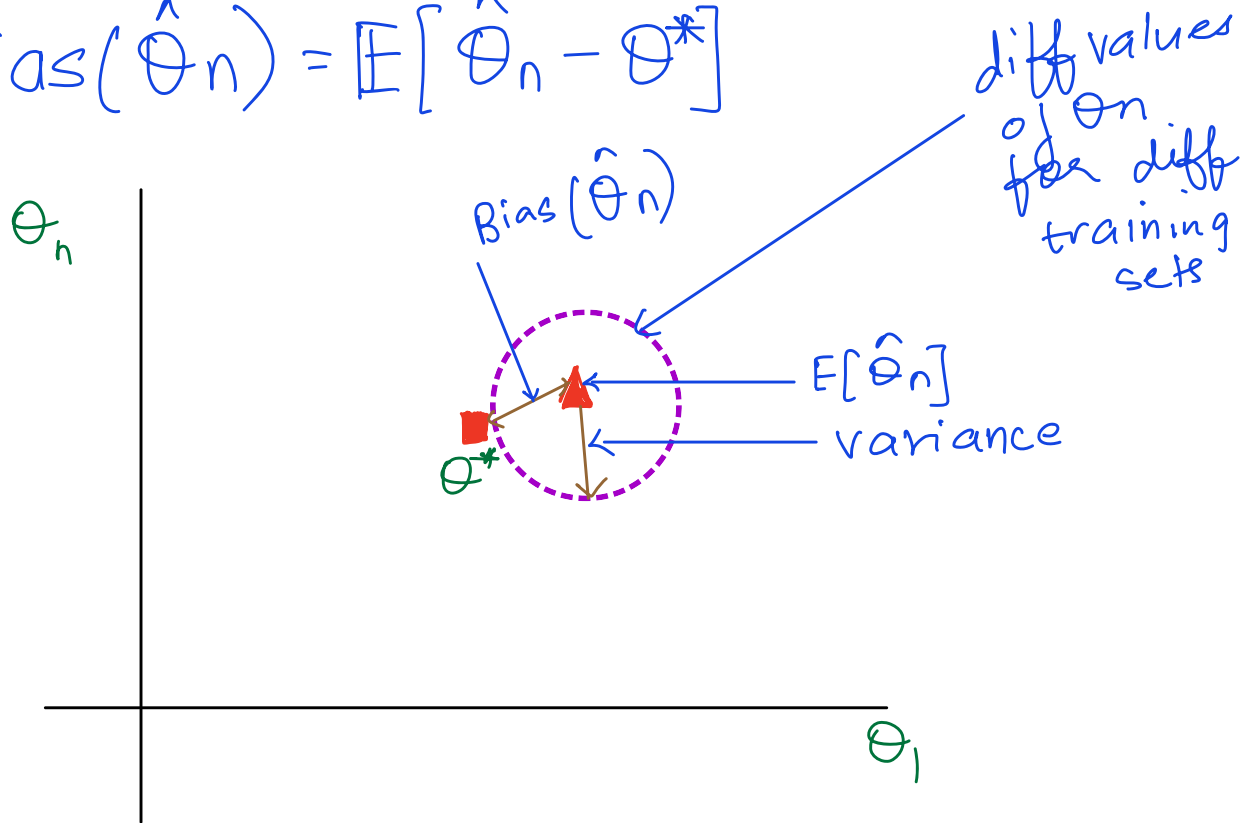
$$\text{Bias}(\hat{\Theta}_n) = \mathbb{E}\left[\hat{\Theta}_n - \Theta^*\right]$$
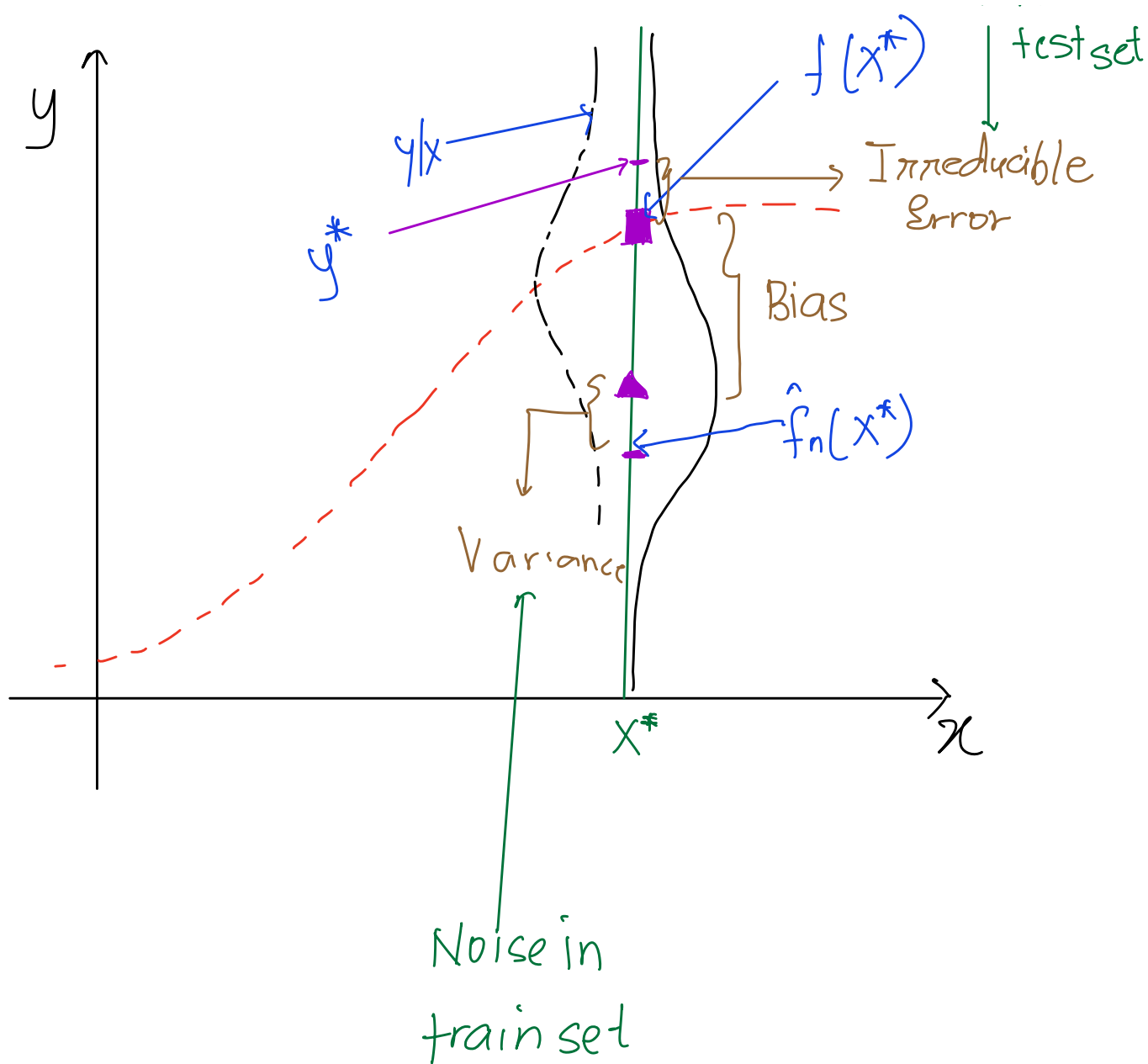
$\Theta_n$

diff values of $\Theta_n$ for diff training sets

$\text{Bias}(\hat{\Theta}_n)$

$\mathbb{E}[\hat{\Theta}_n]$

variance

$\Theta^*$

$\Theta_1$

$$\text{Bias}(\hat{f}_n) = \mathbb{E}\left[\hat{f}_n(x^*) - f(x^*)\right]$$

$$\text{where } f = \mathbb{E}[y|x], \quad y = f(x) + \varepsilon$$

$$\text{Var}(\hat{f}_n) = v\left[\hat{f}_n(x^*)\right]$$

Noise in

$y$

$y|x$

$y^*$

$f(x^*)$

test set

Irreducible Error

Bias

$\hat{f}_n(x^*)$

Variance

$x^*$

$x$

Noise in train set

# Regularization

Equivalent to MAP

estimation $J(\theta) = \|X\theta - \vec{y}\|_2^2 + \lambda \|\theta\|^2$

$S \rightarrow$ training set $\{(x^{(1)}, y^{(1)}), \ldots (x^{(i)}, y^{(i)}) \ldots\}$

$\hat{\theta}_{MAP} = \arg\max\limits_{\theta} \underbrace{p(S|\theta)}_{\text{Likelihood}} \underbrace{p(\theta)}_{\text{Regularizer}}$

MAP - Maximum a posteriori estimate

Instead of calculation the whole posterior, calculate mode of posterior & use that point estimate as output of estimator
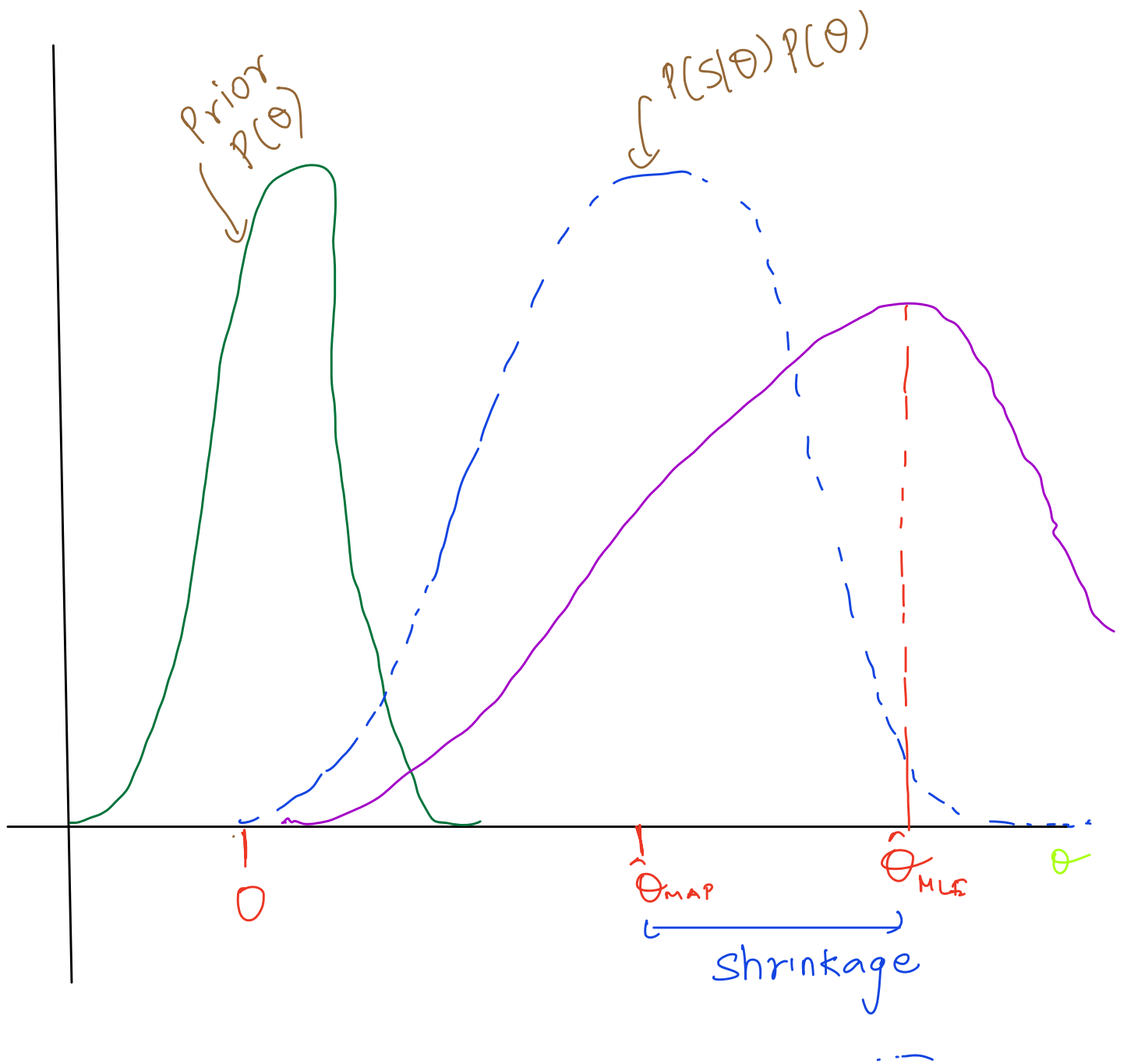
For linear Regression,

$P(S|\theta) \rightarrow$ Likelihood is some scalar multiple of $||X\theta - y||_2^2$

For Gaussian priors, prior takes form of squared error on $\theta$ and the $\lambda$ is directly related to the variance of the prior

$$P(\theta|S) = \frac{P(S|\theta) \, P(\theta)}{P(S)}$$

$$\arg\max_{\theta} P(\theta|S) = \arg\max_{\theta} \frac{P(S|\theta) \, P(\theta)}{P(S)}$$

$$= \arg\max_{\theta} P(S|\theta) \, P(\theta)$$

Prior
$P(\theta)$

$\int P(S|\theta)P(\theta)$

$\theta$

$0$

$\hat{\theta}_{MAP}$

$\hat{\theta}_{MLE}$

Shrinkage

# L2 Regularized Linear Regression

$$J(\theta) = \left( \sum_{i=1}^{n} \left( y^{(i)} - \theta^T x^{(i)} \right)^2 \right) + \lambda \|\theta\|_2^2$$

$$\hat{\theta}_n = \left( X^T X + \lambda I \right)^{-1} X^T \vec{y}$$

$$X^T X + \lambda I = U \begin{bmatrix} \sigma_1^2 + \lambda & & \\ & \ddots & \\ & & \sigma_d^2 + \lambda \end{bmatrix} U^T$$

$$\left( X^T X + \lambda I \right)^{-1} = U \begin{bmatrix} (\sigma_1^2 + \lambda)^{-1} & & \\ & \ddots & \\ & & (\sigma_d^2 + \lambda)^{-1} \end{bmatrix} U^T$$

$$U U^T = I = U^T U$$

$$E[\hat{\Theta}_n] = \left[ U \begin{bmatrix} \frac{\sigma_1^2}{\sigma_1^2 + \lambda} & & 0 \\ & \ddots & \\ 0 & & \frac{\sigma_d^2}{\sigma_d^2 + \lambda} \end{bmatrix} U^T \right] \Theta^*$$

$$\hat{\Theta}_n = (X^T X + \lambda I)^{-1} X^T (X\Theta^* + \varepsilon)$$

Standard L.R. is unbiased
$$E[\hat{\Theta}_n] = \Theta^*$$

But when $\lambda > 0$, $\Theta^* \neq E[\hat{\Theta}_n]$ so bias added.

Also, see eigen values are all $< 1$, so they have a shrinkage effect,

$$\text{cov}(\hat{\theta}_n) = U \begin{bmatrix} \dfrac{T^2 \sigma_1^2}{(\sigma_1^2 + \lambda)^2} & & \\ & \ddots & \\ & & \dfrac{T^2 \sigma_d^2}{(\sigma_d^2 + \lambda)^2} \end{bmatrix} U^T$$

$$\varepsilon \sim N(0, T^2)$$

$$y = \theta^{*T} X + \varepsilon$$

$$\text{Cov}(\hat{\theta}_n) = U \begin{bmatrix} \dfrac{T^2 \sigma_1^2}{(\sigma_1^2 + \lambda)^2} & & \\ & \ddots & \\ & & \dfrac{T^2 \sigma_d^2}{(\sigma_d^2 + \lambda)^2} \end{bmatrix} U^T$$

As $\lambda \uparrow$, variance reduces, bias increases.

$$MSE[\hat{f}_n] = \Upsilon^2 + E[\hat{f}_n(x^*) - f(x^*)]^2 + V[\hat{f}_n(x^*)]$$

Irreducible error

Bias$^2$

Variance

$$f(x) = \Theta^{*T} x$$

$$Bias(\hat{f}_n(x^*)) = E[\hat{f}_n(x^*) - f(x^*)]$$

$$= E[\hat{\Theta}_n^T x^* - \Theta^{*T} x^*]$$

$$= E[\hat{\Theta}_n - \Theta^*]^T x^*$$

$$= Bias(\hat{\Theta}_n)^T x^*$$

$$Var[\hat{f}_n] = x^T Cov[\hat{\Theta}_n] x$$

Heuristics for Bias & Variance

Training Error $\simeq$ Bias

Cross Val. Error $-$ Training Error $\simeq$ Variance

To fight Bias

* Make model larger

* Reduce regularization

To fight Variance

* Collect more data

* Increase regularization

# Learning Theory

① Train & Test Data ~ Same Distribution
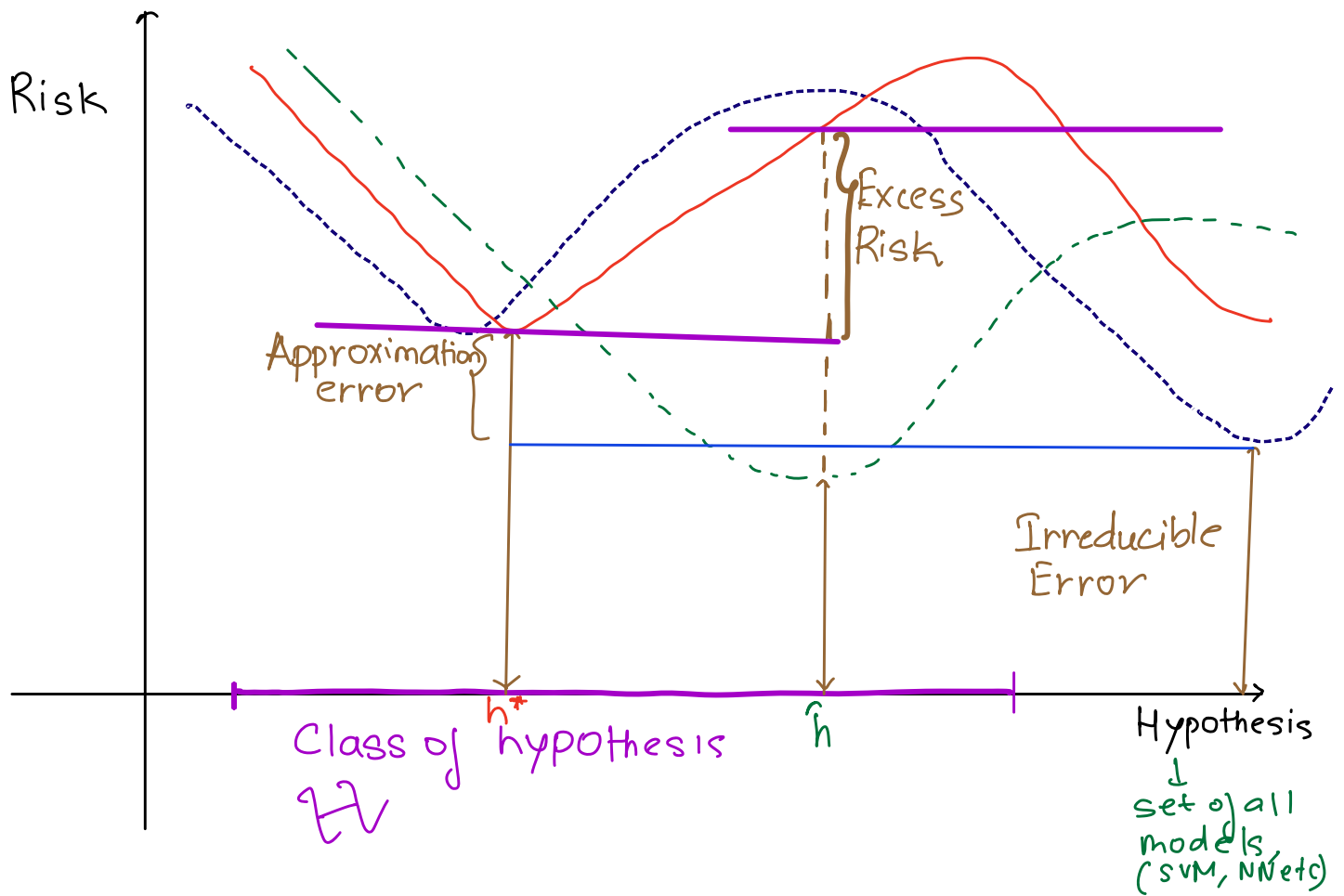
② Examples are sampled I.I.D.

Risk of Hypothesis

$$\mathcal{E}(h) = \mathbb{E}_{(x,y) \sim D} \left[ loss(y, h(x)) \right]$$

$$S = \left\{ (x^{(i)}, y^{(i)}) \right\}_{i=1}^{n}$$

Empirical Risk

$$\hat{\mathcal{E}}(h) = \frac{1}{|S|} \sum_{(x,y) \in S} loss(y, h(x))$$

- For a given h, relation bet$^n$ $\hat{\varepsilon}(h)$ & $\varepsilon(h)$?

- How does our G.E. compare to the best G.E.?



Risk

Approximation error

Excess Risk

Irreducible Error

Class of hypothesis $\mathcal{H}$

$h^*$

$\hat{h}$

Hypothesis
↓
set of all models (SVM, NNets)

—— True Risk

-- -- Empirical Risk (M1)

-- -- Empirical Risk (M2)

As you ↑ size of dataset, the dist$^n$ bet$^n$ curves decreases i.e. becomes tighter.

$h^*$ Best in class hypothesis

Approximation Error - Penalty you pay by limiting to a class of models

Excess Risk - Penalty due to smaller dataset

Step 1 : Uniform Convergence

Step 2: Excess Risk Bound

Uniform Convergence , w.p. $\geq 1 - \delta$, all $h \in \mathcal{H}$

$$\left| \mathcal{E}(h) - \hat{\mathcal{E}}(h) \right| \leq \gamma$$

$$\underset{\text{Term}(n, \delta, \mathcal{H})}{\uparrow}$$

Eg. Suppose $\delta = 0.1$ , then we say

with 90% probability, the gap bet" true & empirical risk is lessthan $\sqrt{}$