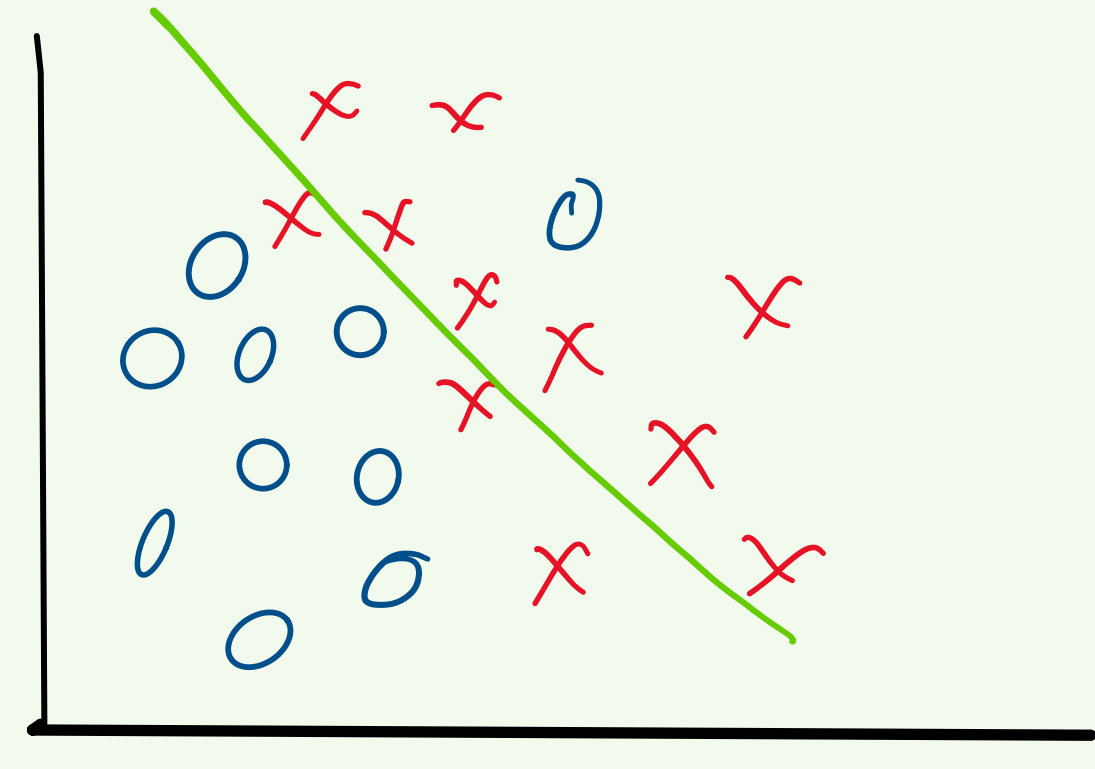


Hyperparameter tuning & optimization.

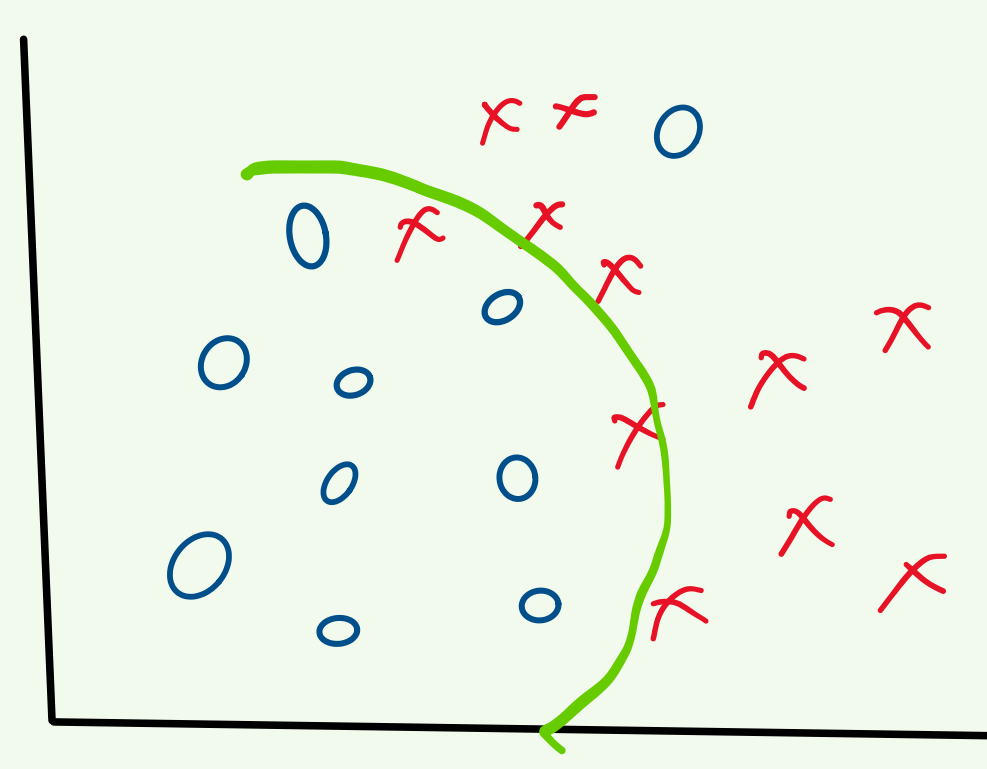
Data division → 60% 20% 20%
Train dev Test

Bias, Variance trade-off

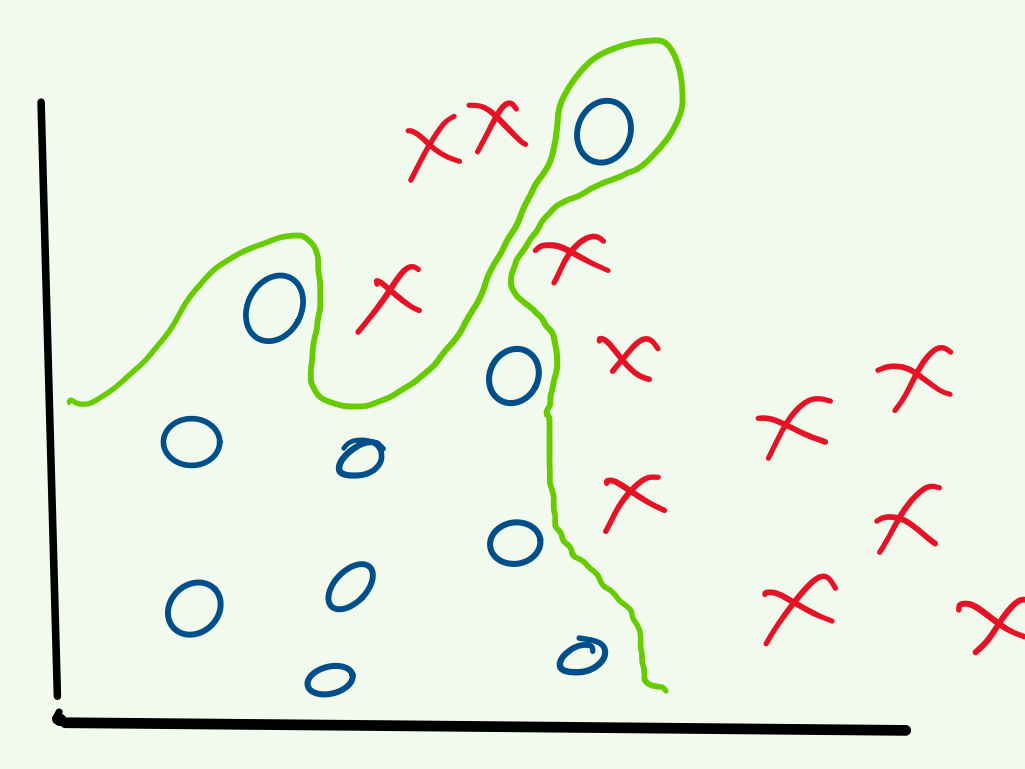


underfitting

bias ↑



just right



overfitting

variance ↑

Error

Train	15%	0.5%	1%	15%
Test	16%	1%	16%	30% ↓ bias ↑ variance ↑

bias ↓
variance ↓

High bias → Train longer (not larger data) } help it
→ bigger n/w } learn

High variance → Supply more data } help it
Regularization } generalize

Frobenius Form (L_2 regularization)

In general,

$$J(w^{[1]} b^{[1]} \dots w^{[L]} b^{[L]}) = \frac{1}{m} \sum_{i=1}^m \ell(\hat{y}^{(i)}, y^{(i)})$$

training example

Now

Frobenius form { L_2 regularization }

$$J(w^{[1]} b^{[1]} \dots w^{[L]} b^{[L]}) = \frac{1}{m} \sum_{j=1}^m \ell(\hat{y}^{(j)}, y^{(j)}) + \frac{\lambda}{2m} \sum_{l=1}^L \|w_l\|^2$$

Effect in back propagation!-

$$dw^{[L]} = (\text{from } b/p) + \frac{\lambda}{m} * \|w\|$$

∴ in gradient descent

$$dw = dw - \alpha * \left(\text{from } b/p + \frac{\lambda}{m} * \|w\| \right)$$

↳ that's why L_2 is called as weight decay algorithm.

Why regularization reduces overfitting?

1) L_2 : it penalizes the weights in G.D

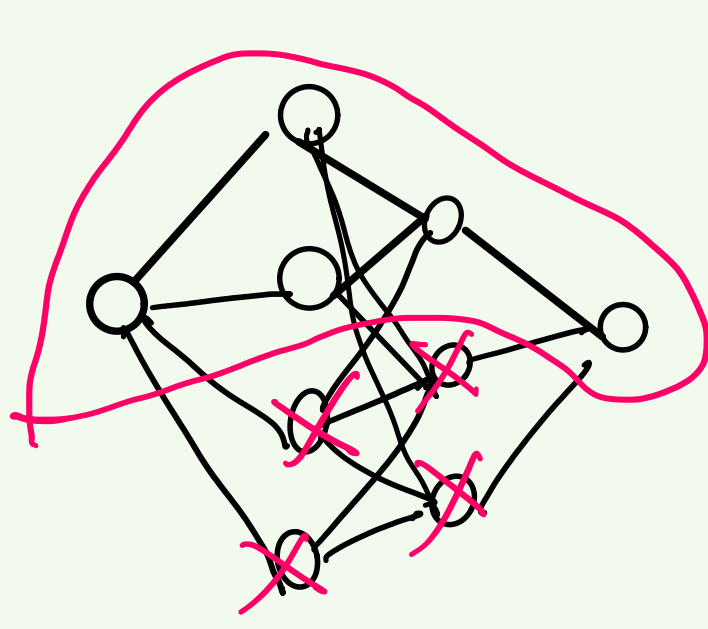
↳ eliminates few nodes

↳ makes the neural n/w simpler

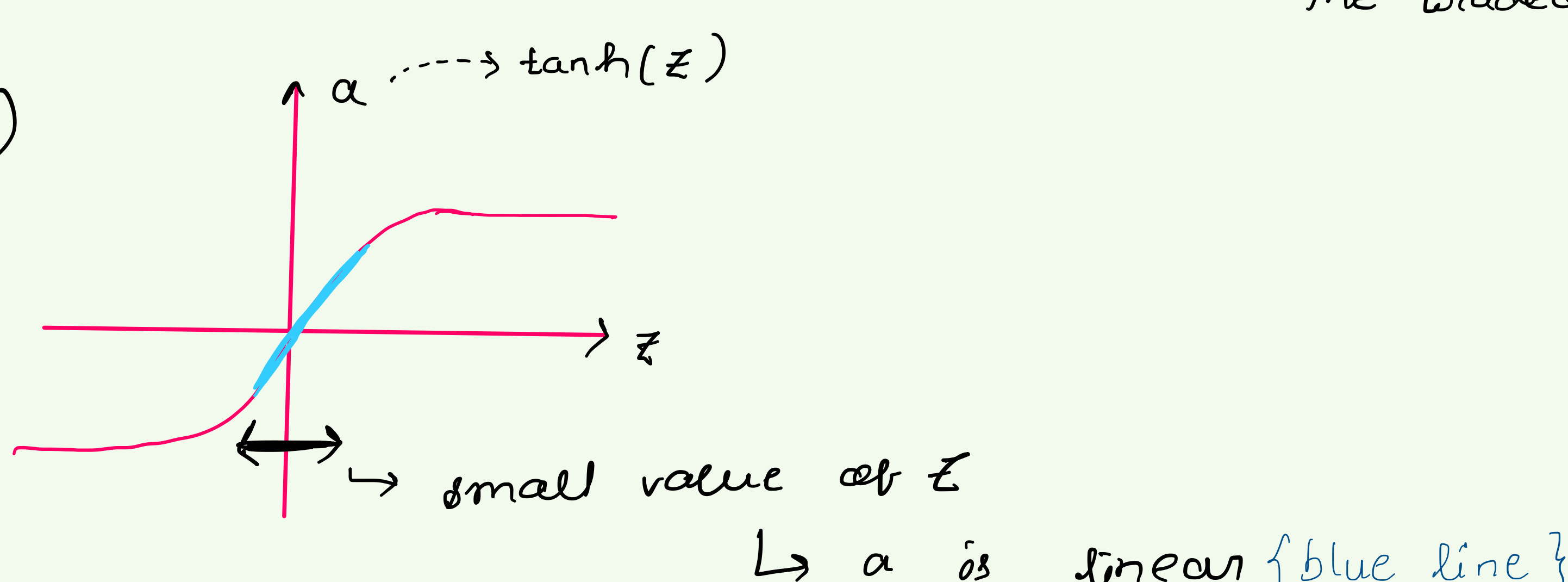
↳ reduces learning

↳ forces toward higher bias

↳ good value of α → optimizes the tradeoff



2)



$$z = w^T x + b$$

↳ smaller after L_2

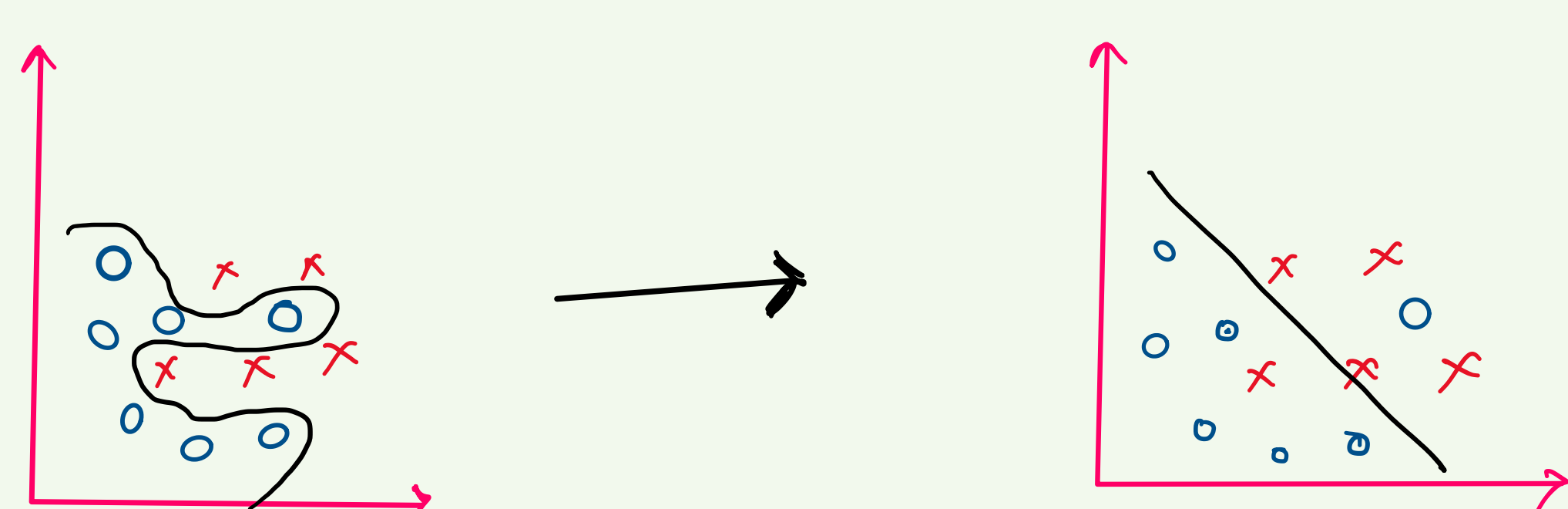
↳ z smaller

↳ Predicted a will be from linear highlighted range.

↳ decision boundaries will be linear

↳ overfitting will be reduced

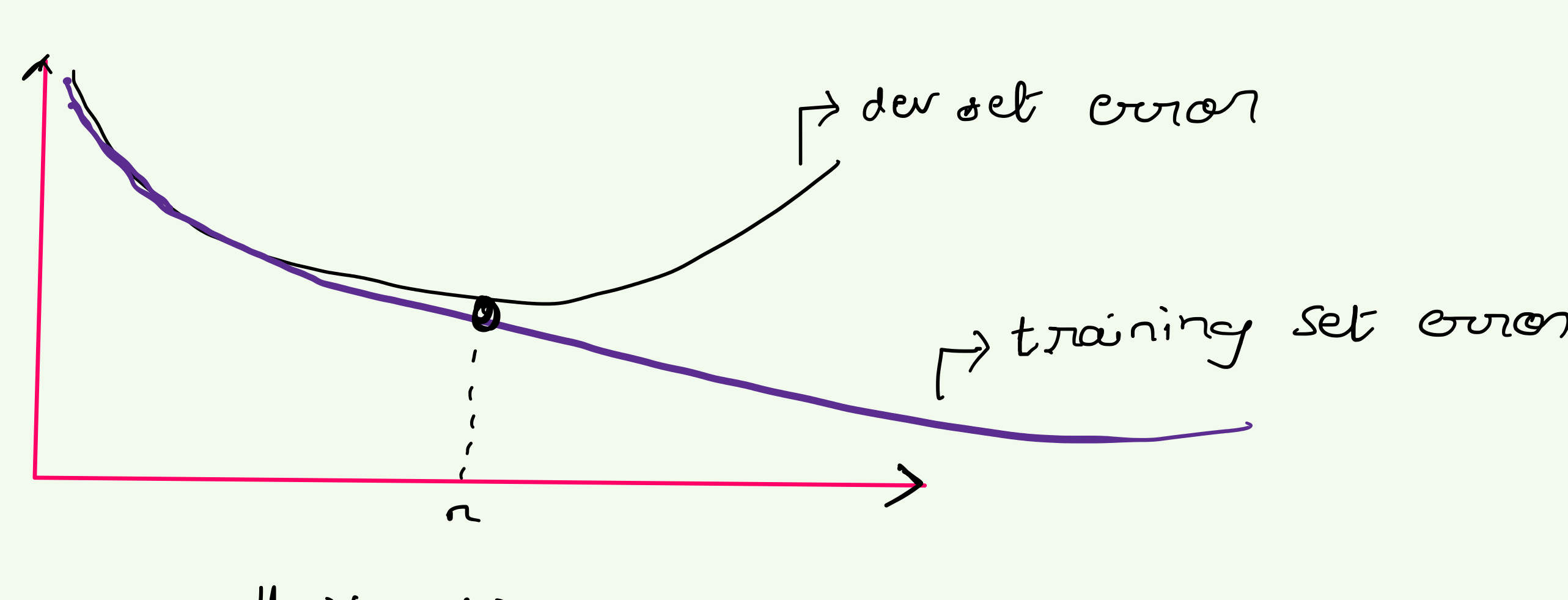
important



low bias
high variance

low bias
low variance.

* Early Stopping:-



stop training at 'n'.

* Disadvantages:-

→ it mixes gradient descent with Regularization } orthogonalization.

↳ * Both are separate tasks. should never be mixed.