# Model Selection:
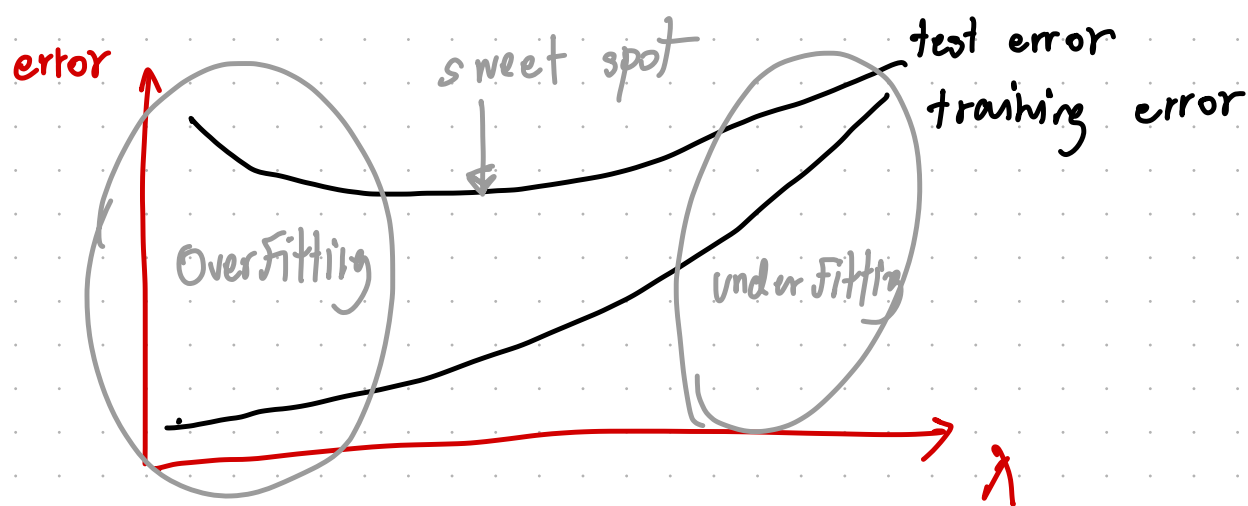
- Recall ERM:

$$\min_{\vec{w}} \frac{1}{n} \sum_{i=1}^{n} \underbrace{l(h_{\vec{w}}(x_i), y_i)}_{Loss} + \underbrace{\lambda \, r(w)}_{Regularizer}$$

- <u>Under fitting</u>: The solution is too simple.
  The training and test error will be high.

- <u>Overfitting</u>: The solution is too complex.
  The training error decrease over time, the
  test error will begin to increase.



- <u>Identifying sweet spot</u>: Divide data into training and validation
  portions. Train on the "training" split and evaluate it on
  the "validation" split, for various value of $\lambda$
  $(10^{-5}, 10^{-4}, 10^{-2}, 10^{-1}, 10^{0}, 10^{1}, 10^{2})$

$D$ — $D_{TR}$ | $D_{TE}$

$D'_{TR}$ | $D'_{VA}$

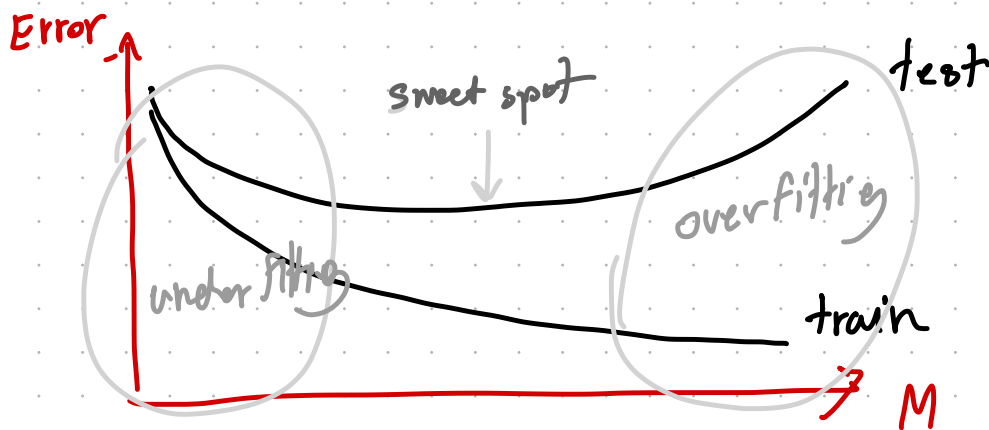We may overfit the validation set.

## K-fold cross validation:

- Divide your training data into k parts.
- Train on k-1 parts and leave one out as validation set. Do this k times and average the validation error across all runs (ideally, $k \sim n$).

## Telescopic search for $\lambda$:

- Do two search:
  - **1st step**: find the best order of magnitude for $\lambda$
    - For example, first we try 0.01, 0.1, 1, 10, 100.  ← the best
  - **2nd step**: do a more fine-grained search around the best $\lambda$ so far
    - Then, we try 5, 10, 15, 20, 25, ...., 95 to test values around 10.

Early Stopping : Stop our optimization after
$M \geq 0$ number of gradient steps,
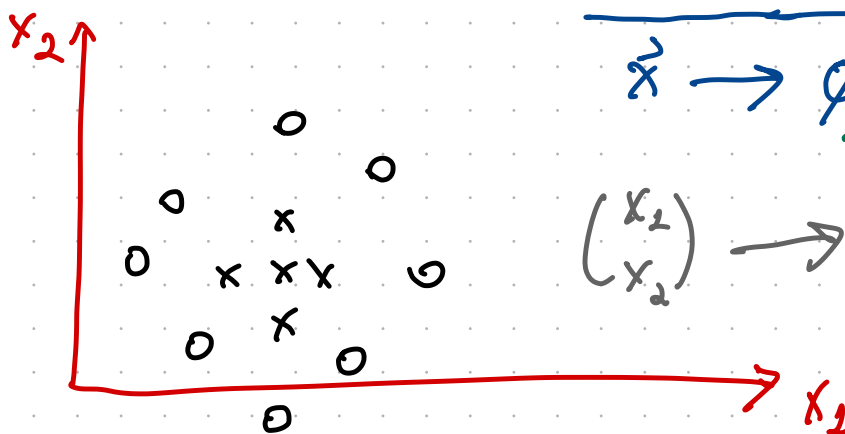even if optimization has not converged yet.



---

# Kernels :

- A way to incorporate non-linearities into most
linear classifiers.



feature Transformation :
$$\vec{x} \longrightarrow \phi(\vec{x})$$

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \longrightarrow \begin{pmatrix} x_1 \\ x_2 \\ x_1^2 + x_2^2 \end{pmatrix}$$

Extreme case: $\begin{pmatrix} x_1 \\ \vdots \\ x_d \end{pmatrix} \longrightarrow \begin{pmatrix} 1 \\ x_1 \\ \vdots \\ x_d \\ x_1 x_2 \\ \vdots \\ x_{d-1} x_d \\ \vdots \\ x_1 x_2 \dots x_d \end{pmatrix} \in 2^d$