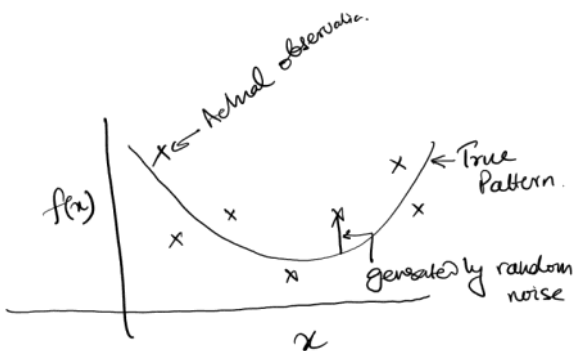


## The learning process.

- ① Data — sample from distribution  $\rightarrow$  (train, validation), test split
- ② Predictor function
- ③ Cost function
- ④ Optimization algorithm.

Regularization  $\rightarrow$  ! Increase bias slightly, but reduce variance !

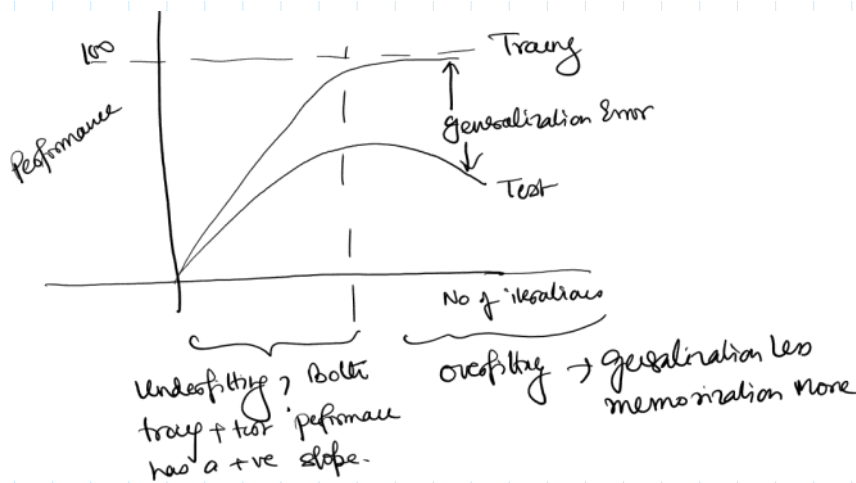
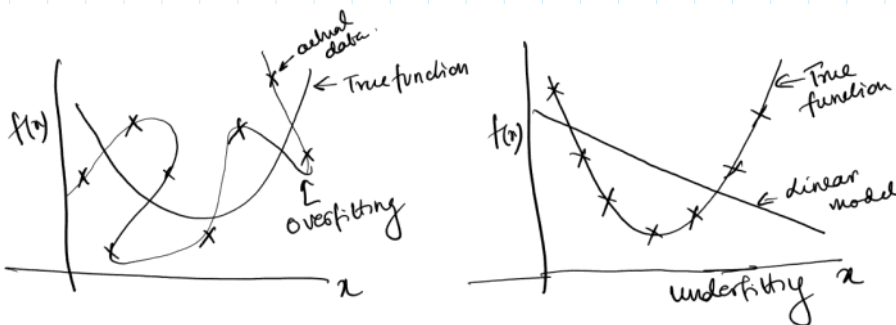


Adding penalty to loss  $f^*$  to discourage complex models that fit training data too closely.

Techniques for regularization:

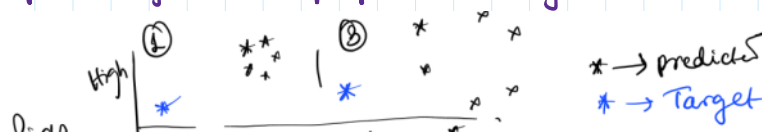
- COST: ridge, lasso, elastic
- MODEL: pooling, dropout
- DATA: augmentation, noise injection, bootstrap.

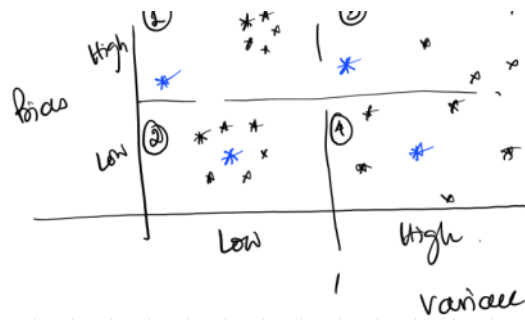
## Over + Underfitting



## Bias - Variance Tradeoff.

- Variance = variability of model prediction
- Bias = expected avg. distance b/w prediction + target. — we want this as low as possible.



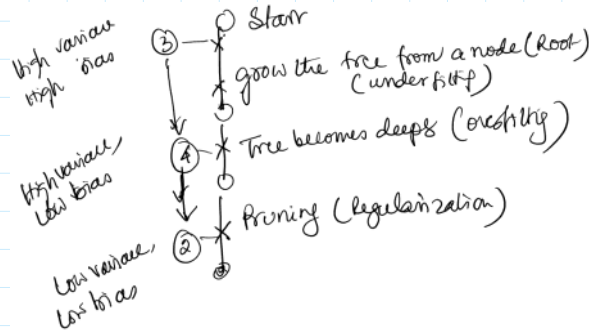


\* → predicted  
\* → Target

Complex model : ③ → ① → ④ → ②

Deep learning : ③ → ④ → ②

Take a look @ the decision tree:



## Occam's Razor & Minimum Description Length.

among many adequate explanations, simplest is best. Entities shouldn't be multiplied beyond necessity.

In the case of a neural network — many sets of weights + biases — how to choose?

↳ Get the most simple.

Say we're minimizing  $L(\theta)$

↳  $\bar{w}, \bar{b}$

add a penalty for "unsimplicity."

$$L(\theta) + \lambda R(\theta)$$

$\lambda$  is a hyperparam — trial/error  
 $\lambda \propto$  penalty domination.

Go back to neural network.

$$\delta \bar{w}_t = \eta \nabla_{\bar{w}} L(\bar{w}_t, \bar{b}_t)$$

$$\delta \bar{b}_t = \eta \nabla_{\bar{b}} L(\bar{w}_t, \bar{b}_t)$$

$L(\bar{w}_t, \bar{b}_t)$  = loss @ iteration t.

$$= \sum_{i=0}^{n-1} L(\hat{y}^i, y^i | \bar{w}, \bar{b})$$

↑      ↑  
pred    true.

L2 regularization: (ridge)  
 $R(\theta) = (|\bar{w}|^2 + |\bar{b}|^2)$

$$L(\bar{w}, \bar{b}) = \sum L(\hat{y}^i, y^i) + \lambda (|\bar{w}|^2 + |\bar{b}|^2)$$

$$\text{Let } L(\bar{w}) = (|\bar{w}|^2)$$

$$\frac{\partial L}{\partial w} = 2w$$

gradient ↓ as  $w \rightarrow 0$ , update ↓↓

L2 → dense weight vectors

L1 regularization: (lasso)  
 $R(\theta) = (|\bar{w}| + |\bar{b}|)$

$$L(\bar{w}, \bar{b}) = \sum L(\hat{y}^i, y^i) + \lambda (|\bar{w}| + |\bar{b}|)$$

$$\text{Let } L(\bar{w}) = (|\bar{w}|)$$

$$\frac{\partial L}{\partial w} = \begin{cases} 1 & w > 0 \\ 0 & w = 0 \\ -1 & w < 0 \end{cases}$$

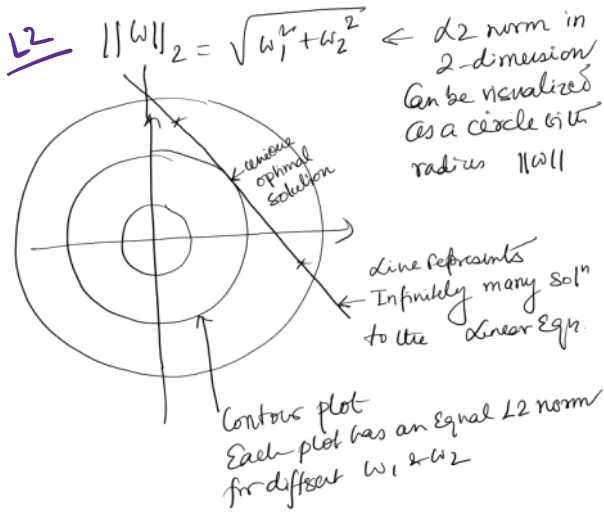
constant gradient.

constant gradient.

L2  $\rightarrow$  dense weight vectors

L1  $\rightarrow$  sparse - "-

L1 (lasso) performs better!

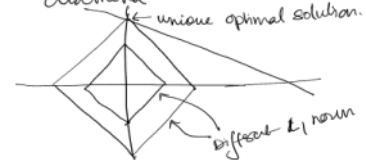


L1  $\|w\|_1 = \sum_{i=1}^2 |w_i| = |w_1| + |w_2|$

In this there are 4 combinations

$$\begin{array}{l|l} w_1 > 0, w_2 > 0 & w_1 + w_2 = C \\ w_1 < 0, w_2 > 0 & w_2 - w_1 = C \\ w_1 > 0, w_2 < 0 & w_1 - w_2 = C \\ w_1 < 0, w_2 < 0 & -w_1 - w_2 = C \end{array}$$

Jointly, these 4 lines will form a diamond shape. All points on the diamond share same L1 norm



L2

- doesn't eliminate coefficient
- may overfit (features > obs)
- Can be managed by GP

L1

- eliminates coefficients
- may eliminate X.
- undefined  $\partial L / \partial w$  when  $\bar{w} = 0$  ( $\Rightarrow$  modified GP)

$\rightarrow$  Elastic Net  $\leftarrow$

$$L = \sum (y_i - \hat{y})^2 + \alpha_1 \sum |w_i| + \alpha_2 \sum |w_i|^2$$

