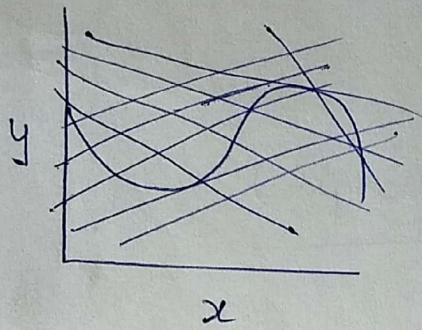


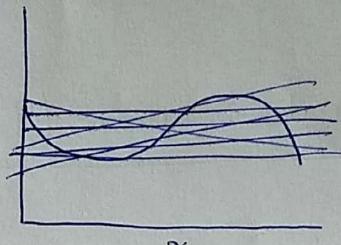
Regularization

(1)

Regularization is the cure for overfitting. Regularization is used in every machine learning application you will encounter.



Without regularization



With regularization

Controlling the lines ^{by} best restricting them in terms of offset ~~that~~ they can have \leq the slope they can have.

Without regularization: bias = 0.21; var = 1.69

With regularization: bias = 0.23; var = 0.33

Side note: We can use Legendre polynomials to approximate target function

Unconstrained Solution to minimize MSE

Let's assume we have a vector $z = \begin{bmatrix} 1 \\ \text{2nd Order term, } x^2 \\ \text{3rd Order term, } x^3 \\ \vdots \\ \text{9th Order term, } x^9 \end{bmatrix}$

For this set up, y is defined as

$y = \underbrace{w^T z}_{\text{y}} \quad \text{where } w \text{ is the vector of weights}$

$$\text{Minimize } E_{in}(\omega) = \frac{1}{N} \sum_n (\omega^T z_n - y_n)^2$$

$$= \frac{1}{N} (z\omega - y)^T (z\omega - y)$$

$$\omega_{lin} = (z^T z)^{-1} z^T y$$

Constrained Solution (constraining the weights)

$$\text{Constraint: } \sum_q \omega_q^2 \leq C$$

$$\omega_1^2 + \omega_2^2 + \dots + \omega_N^2 \leq C$$

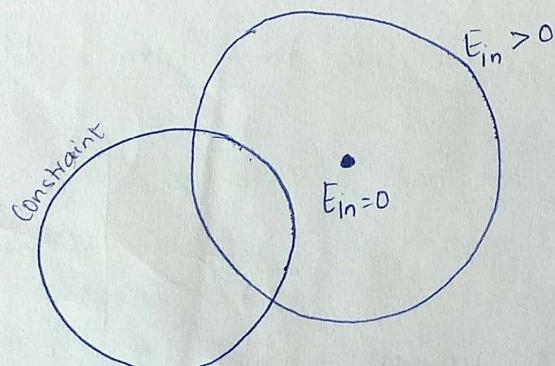
(Sphere formula)

$$\text{Minimize } \frac{1}{N} (z\omega - y)^T (z\omega - y)$$

$$\text{Subject to: } \omega^T \omega \leq C$$

Solving graphically

Let's draw a contour line for E_{in} (MSE) & for the constraint

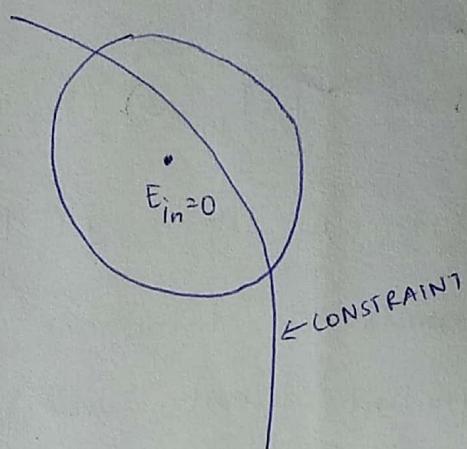


(2)

best value of E_{in} is actually at $w^T w = c$

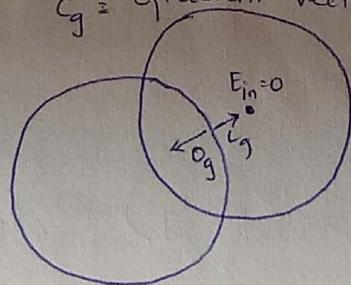
If c is big enough that includes $E_{in}=0$, the solution is

w_{lin}



The best possible 'w' can be found when the gradient vectors of both the constraint function & the objective function are opposite to one another

\vec{o}_g = Gradient Vector of Objective Function
 \vec{c}_g = Gradient Vector of Constraint function



Because of

Lagrangian, we get

$$L(w) = E_{in}(w) + \frac{1}{2} w^T w \quad \dots \textcircled{1}$$

$$\text{as } L(x, y, d) = f(x, y) - d(g(x, y) - b)$$

It's ~~+ve~~ $+ \frac{1}{2} w^T w$ & not $-ve$ because both the gradients are opposite to one another i.e. $\nabla f = -d \nabla g$.

Now applying the gradient for $\textcircled{1}$ & equating to 0, we get,

$$\nabla E_{in}(w) + \gamma d w = 0, \text{ with } \gamma \text{ small}$$

$$\nabla E_{in}(w_{reg}) + \frac{\gamma}{N} w_{reg} = 0 \quad \dots \quad (2)$$

Here d is the rate of change of $E_{in}(w)$ for unit change of constant C . In the above case it is $C \uparrow \downarrow$

Here in this particular scenario as d increases C decreases as we are talking about gradient descent. In gradient descent case, the $\nabla E_{in}(w)$ is opposite to $\nabla(w w^T - C)$ & so are inversely proportional to one another.

If C is huge, d is zero.

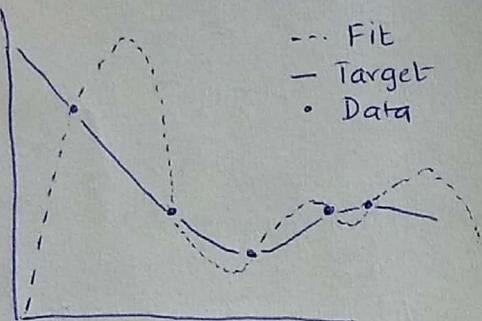
If we expand & solve (1), we get

$$\frac{1}{N} ((z_w - y)^T (z_w - y) + \gamma w^T w)$$

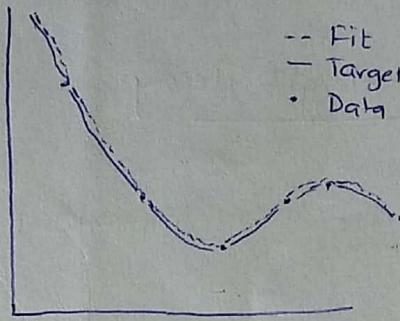
Solve for $\nabla E_{aug}(w) = 0$, where aug stands for augmenting $\gamma w^T w$ to $E_{in}(w)$, we get

$$w_{reg} = (Z^T Z + \gamma I)^{-1} Z^T y, \text{ with regularization}$$

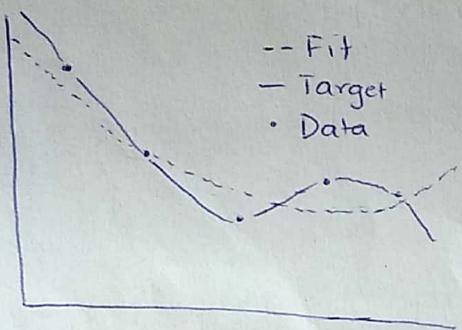
$$w_{in} = (Z^T Z)^{-1} Z^T y, \text{ without regularization.}$$

Example

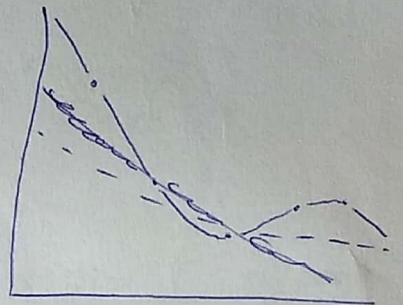
$d = 0$
(Overfit)



$d = 0.0001$
(Good Fit)



$d = 0.01$
(Under fit)



$d = 1$
(Under fit)

The regularizer explained above is called weight decay.

this term applies to neural network ~~also~~

In gradient descent $\overset{\curvearrowright}{w}_{\text{new}} = w_{\text{current}} - \nabla E_{\text{in}}(w_{\text{current}})$ [with out reg]

$w_{\text{new}} = w_{\text{current}} * \overset{\curvearrowright}{(\text{reg})} - \nabla E_{\text{in}}(w_{\text{current}})$ [with reg]

If you note in the above equation, in every iteration to get w_{new} you decay w_{current} by the factor of (reg) & then apply gradient descent. Hence this type of regularizer is termed "weight decay".

"Weight decay" term is more applicable in neural network.

Variations of weight decay

Emphasis on weights:

Some weights are more important than others

$$\sum_{q=0}^Q \gamma_q w_q^2 \leq C$$

Examples $\gamma_q = 2^q \Rightarrow$ low-order fit

$\gamma_q = 2^{-q} \Rightarrow$ high-order fit

because smaller γ will can have larger weights & vice versa

Tikhonov Regularizer is a special form of weight decay

with the form $w^T \gamma^T \gamma w$. If we put this in

matrix form is a general quadratic form. It has

• γ for diagonal guys ($w_1 w_1, w_2 w_2, w_3 w_3, \dots$) if it has

γ for off diagonal guys also ($w_1 w_2, w_1 w_3, w_2 w_3, \dots w_5 w_6, \dots$)
for correlated weights..

With proper choices of γ you can get various forms

like low order, high order etc.

Regularizer harms fitting the noise more than it harms fitting the signal. It harms overfitting more than it harms the fitting.

Soft Weight elimination is another form of regularizer. Where some of the weights are eliminated (not completely, but close to zero)

This is applied in neural networks & is given by

$$R(w) = \sum \frac{(w_{ij}^{(l)})^2}{\beta^2 + (w_{ij}^{(l)})^2}$$

The intuition is ~~small~~ Fewer weights \Rightarrow Smaller VC dimension.

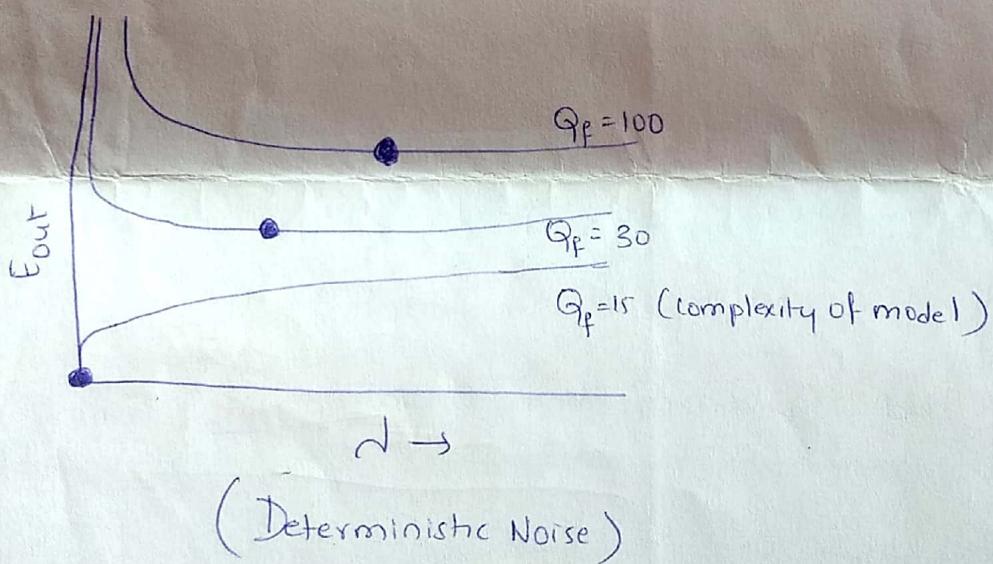
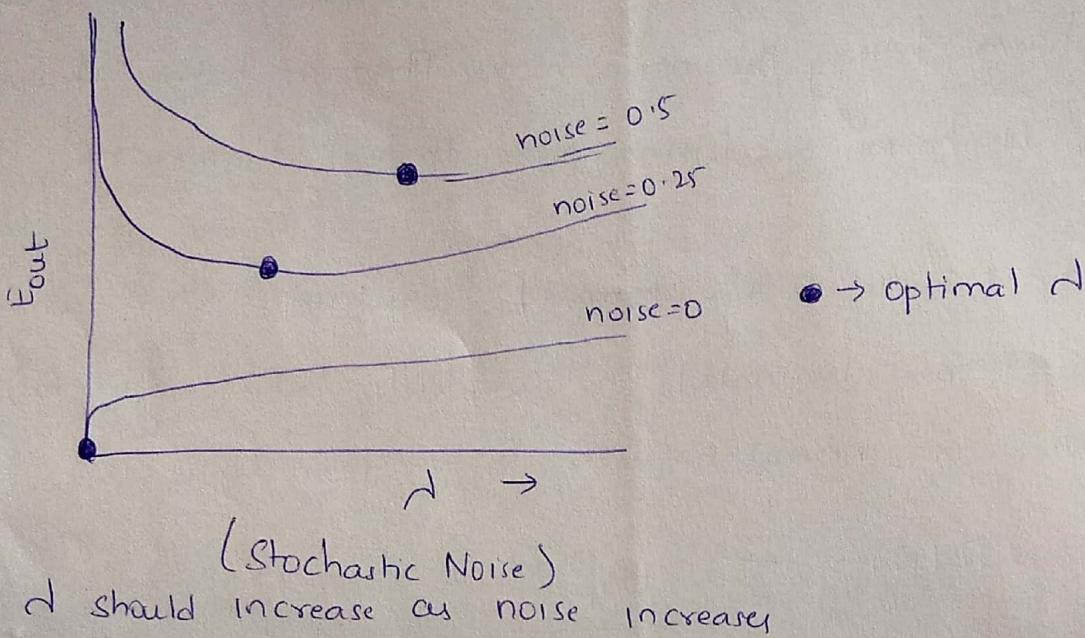
Guiding principle for choosing the perfect regularizer is to choose in the direction of smoother or simpler.

If we choose bad regularizer then ~~Validation~~ Validation comes to rescue.

Early stopping as a regularizer. Stop the epoch early.

Regul

Optimal γ



γ should increase as complexity of the model increases.

Finally, choosing a regularizer is an experimental activity, more than a completely principled activity.

Bigger VC dimensions vs choosing regularizer
is to choose higher VC dimension & work backwards into
the lower side by adding regularizer.

Validation

Regularization vs Validation

$$\text{Generalization Error} = \text{Empirical Error} + \underbrace{\text{Overfit Penalty}}$$



Validation estimates this quantity



Regularization estimates
this quantity

Let's explore the composition / statistical parameters of
Generalization error

$$\text{Mean } (\mu) = E[e(h(x), y)] = \text{Generalization Error}$$

$$\text{Variance} = \text{var}[e(h(x), y)] = \sigma^2$$

Here $e(h(x), y)$ is the error function either Squared error or binary error. It is a random variable.

Now let's consider a set of input datapoints

$$(x_1, y_1), \dots, (x_k, y_k) \rightarrow \text{Validation Set}$$

$$\text{Mean } (\mu) = \frac{1}{K} \sum_{k=1}^K E[e(h(x), y)] = \text{Generalization error}$$

~~REPEATED~~

$$\text{Var } \sigma^2$$

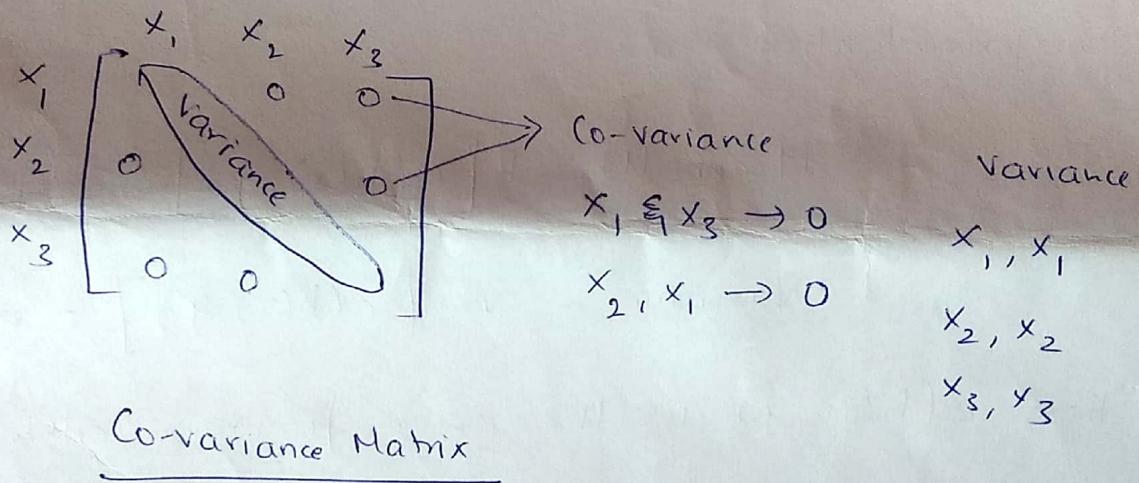
$$\text{Variance} = \left(\frac{1}{k^2} \right) \sum_{k=1}^k \text{var} [e(h(x_k), y_k)] = \frac{\sigma^2}{k}$$

↓

Variance (Not Co-variance)

Number of elements in the Co-variance Matrix

Co-variance is zero because the data points picked are random & hence we end up with only the diagonal elements which are variance



Generalization Error = Validation Error + $O\left(\frac{1}{\sqrt{k}}\right)$

This suggests validation error is the good approximation to Generalization error & is proportional to $\frac{1}{\sqrt{k}}$ (Standard deviation)

Where k is the number of data points chosen for Validation Set.

(2)

~~Division~~~~Let's assume~~Division of the data set

Given the dataset of size $N \quad \{(x_1, y_1), \dots, (x_N, y_N)\}$

Training Set $\rightarrow N - k$

Validation Set $\rightarrow k$

The strength of

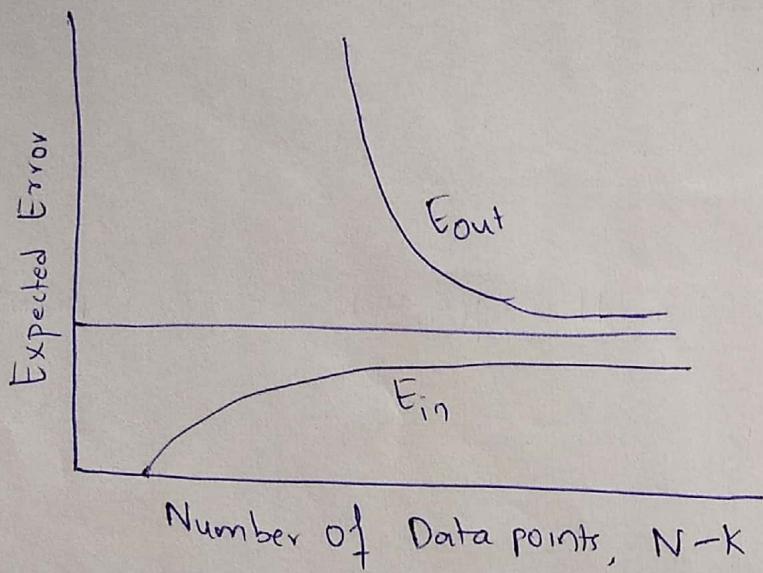
Generalization Error estimate is directly proportional to

$$\text{Error} \left(\frac{1}{\sqrt{k}} \right)$$

small $k \rightarrow$ Less confidence on the estimate

Large $k \rightarrow$ Good confidence on the estimate ?

It's a question mark because we cannot afford to have large k as the size of k increases, ~~the~~ the training size decreases & reliability on the hypothesis will ~~be~~ be costly as shown in the graph below.



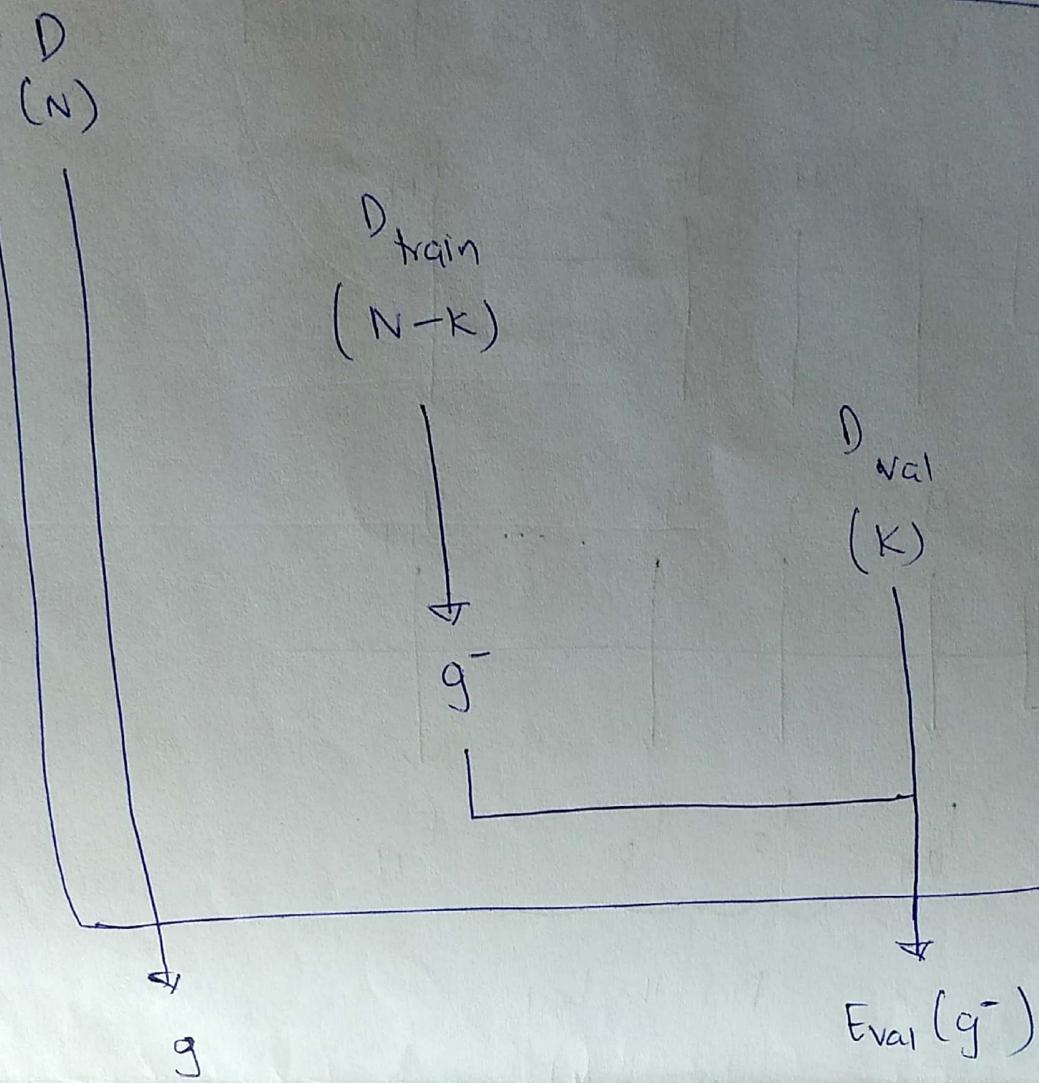
Another technique is to put back k datapoints & retrain the model again which will be given by

$$E_{out} \quad E_{var}(g^-) \quad E_{out}(g) = E_{var}(g^-) \pm O\left(\frac{1}{\sqrt{k}}\right)$$

Notice $g \neq g^-$ (selected final hypothesis on the reduced set, $N-k$)
 \downarrow
 (selected final hypothesis on the complete set, N)

Again, it is not reliable as we are estimating on $N-k$ but reporting on N . As shown below

(3)



So, the rule of thumb is

$$K = \frac{N}{5}$$

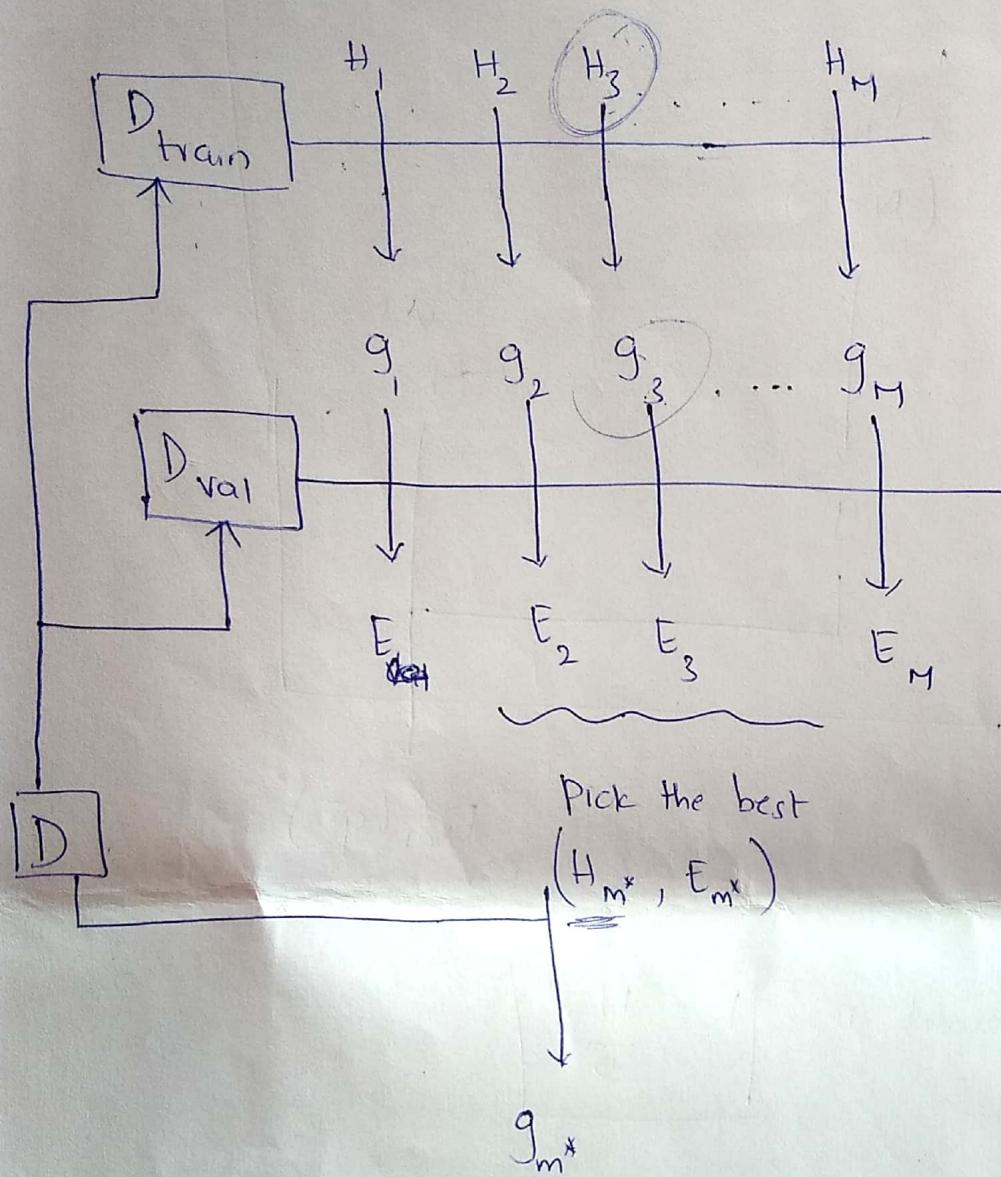
Why is validation set used for?

Validation set is used for picking the model

or choosing the right model & hence it is also called as Optimistic bias model. We are biased towards the model which reports small validation error.

But this bias will not hurt us in generalization & usually ignorable.

This is how the process of model selection happens



Pick D_{train} & train ' M ' models & evaluate them using D_{val} & pick the best model. Now train the best model using D (complete) & report the final hypothesis.

(4)

The dilemma about K

$$E_{\text{out}}(g) \leq E_{\text{out}}(g^-) \leq E_{\text{Val}}(g^-)$$

(small K) (large K)

↑
final hypothesis
we report

Can we have k both small & large?

It turns out to be possible.

Leave one out

$N-1$ points for training & 1 point for validation

$$D_n = (x_1, y_1), \dots, (x_{n-1}, y_{n-1}), \cancel{(x_n, y_n)}, (x_{n+1}, y_{n+1}), \dots$$

Final hypothesis learned from D_n is g_n^-

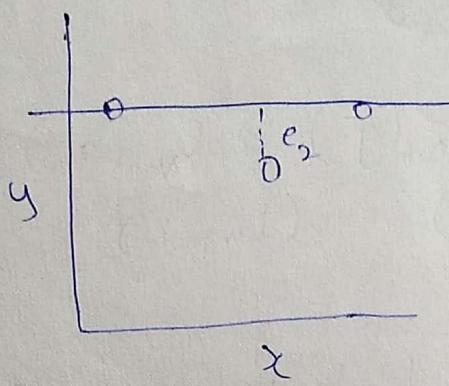
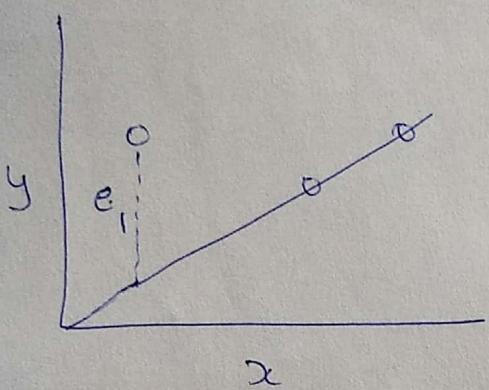
They are all different hypothesis but they are the realization of almost similar data set (training examples)

$$e_n = E_{\text{Val}}(g_n^-) = e(g_n^-(x_n), y_n)$$

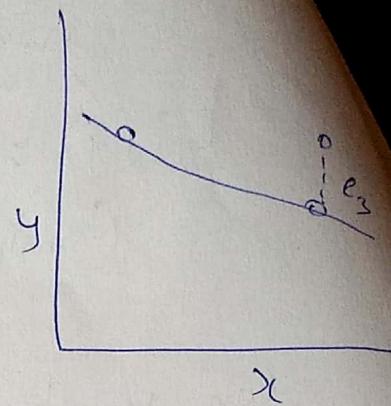
Cross validation error:

$$E_{\text{CV}} = \frac{1}{N} \sum_{n=1}^N e_n$$

Illustration of cross validation

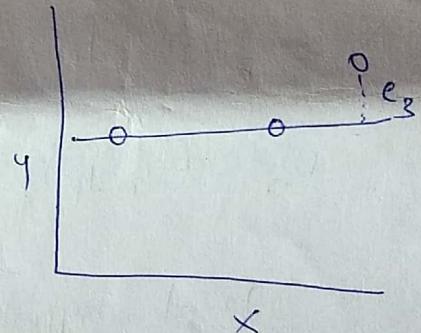
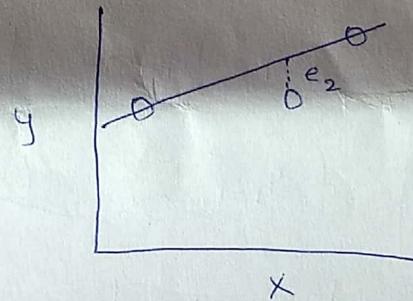
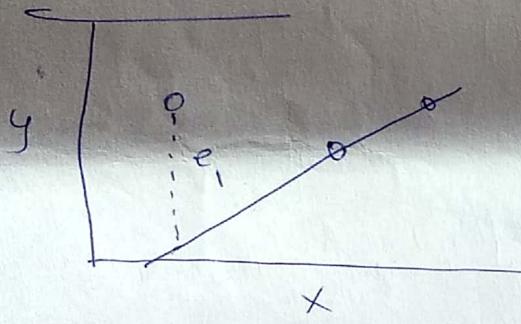


3 - data points example.

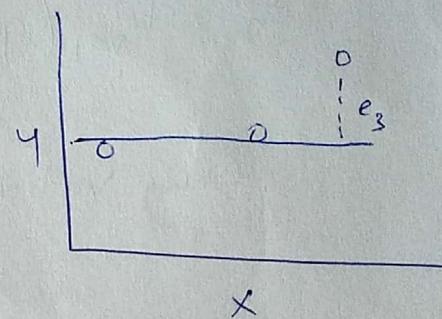
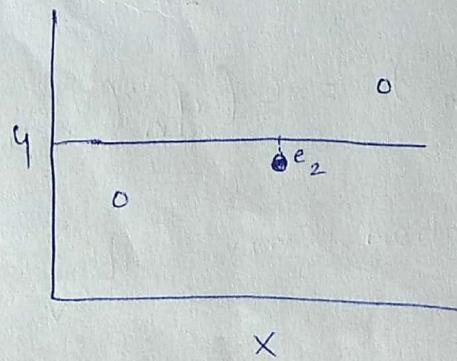
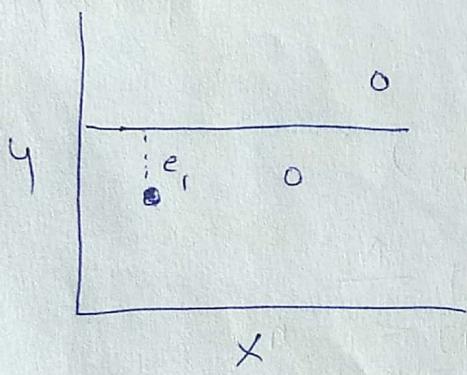


$$E_{CV} = \frac{1}{3} (e_1 + e_2 + e_3)$$

Linear Model :

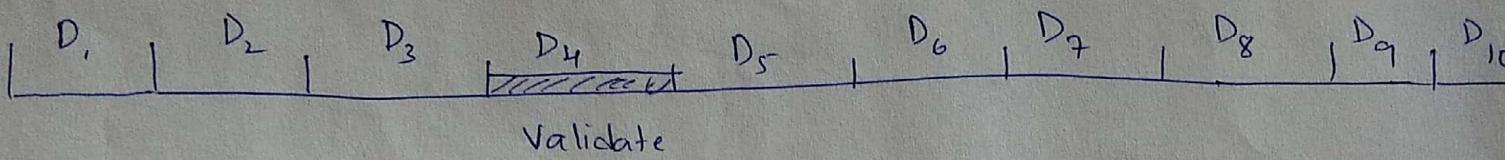


Constant Model :



(5)

Leave more than one Out



$\frac{N}{K}$ training sessions on N-k points each

10-fold cross validation : $K = \frac{N}{10}$

$$P\left[|E_{\text{out}} - E_{\text{in}}| > \epsilon\right] \leq 4m_H(2N)e^{-\frac{1}{18}\epsilon^2 N}$$

$\underbrace{}$
 δ

Get ϵ in terms of δ

$$\delta = 4m_H(2N)e^{-\frac{1}{18}\epsilon^2 N}$$

$$\Rightarrow \epsilon = \sqrt{\frac{\delta}{N} \ln \frac{4m_H(2N)}{\delta}}$$

$\underbrace{\phantom{\ln \frac{4m_H(2N)}{\delta}}}_{n}$

With probability $\geq 1 - \delta$ we get

$$|E_{\text{out}} - E_{\text{in}}| \leq \underline{\mathcal{L}(N, H, \delta)}$$

\downarrow

The Generalization bound.

The Map

THEORY



Bias - Variance

Complexity

Bayesian

TECHNIQUES

MODELS

LINEAR

NEURAL NET

SVM

K-NN

RBF

GAUSSIAN PROCESS

SVD

GRAPHICAL MODELS

METHODS

REGULARIZATION

VALIDATION

AGGREGATION

INPUT PROCESSING

PARADIGMS



SUPERVISED

UNSUPERVISED

RE-INFORCEMENT

ACTIVE

ONLINE

CONTACT INFO :

EMAIL : vinuth.tulasi@gmail.com

GITHUB : www.github.com/vinuth / Deep Learning

① Occam's Razor

Summary

Choose the simplest model possible

~~Simplicity~~ Complexity of h : MDL, Order of polynomial

Complexity of H : VC Dimension

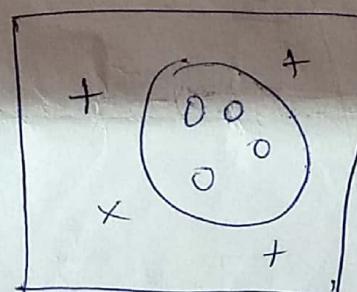
② Sampling Bias

③ Data Snooping

If a data set has affected any set in the learning process, its ability to assess the outcome has been compromised

$$z = (1, x_1, x_2, x_1x_2, x_1^2, x_2^2)$$

$$\text{or } z = (1, x_1^2, x_2^2) \text{ or } (1, x_1^2 + x_2^2)$$



Reuse of a data set

Trying one model after the other on the same data set, you will eventually 'succeed'

If you torture the data long enough, it will confess

VC dimension of the total learning model

May include what others tried!

④ Use of prior knowledge is necessary (Remember Bait Shyness)