# ordinary least squares.

$x \in \mathbb{R}^d$, $y \in \mathbb{R}$, $\phi()$ is a kernel map

from $\mathbb{R}^d \to \mathbb{R}^m$    $m \leq d$.   $\phi = (\phi_0, \phi_1, \ldots \phi_{m-1})$.

$$h(x_*) = W^T \phi(x) + W_0 \leftarrow \text{Generalized linear model.}$$

$W_0 \to$ bias term

$W \to [W_0 \ldots W_{M-1}]^T$    0

Given   $D = \{x_1, x_2 \ldots x_n\}$

$$Y = [y_1, y_2, \ldots, y_n]$$

Eg $\phi$:
$\phi_1 = \{1, x, x^2\}$
$\phi_2 = \exp\{\frac{(x-y_j)^2}{2s^2}\}$

Construct a design matrix $\phi$: $\phi_{nj} = \phi_j(x_n)$.

$$\phi = \begin{bmatrix} \phi_0(x_1) & \phi_1(x_1) & \cdots & \phi_{M-1}(x_1) \\ \phi_0(x_2) & & & \vdots \\ \vdots & & & \\ \phi_0(x_n) & \cdots & \cdots & \phi_{M-i}(x_n) \end{bmatrix}_{N \times M}$$

Objective of OLS: find $W$ s.t

$$L = \|Y - \phi_{N \times M} W_{M \times 1}\|$$

$$L = \| \underset{1 \times M}{Y^T} - \underset{M \times N}{W^T \phi^T} \|_2^2 \quad \text{is minimized.}$$

(112)

$$W_{OLS}^* = \text{argmin} \, \| y - W^T \phi^T \|_2^2$$

$$= \left[ y - W^T \phi^T \right]^T \left[ y - W^T \phi^T \right]$$

$$\frac{\partial L}{\partial W} = 0 \implies W_{OLS}^* = \underbrace{\left( \phi^T \phi \right)^{-1} \phi^T}_{\text{Moore-Penrose}} y.$$

$$\text{Moore-Penrose inverse.}$$

$$= \underline{\phi^+ y}. \quad \left[ \begin{array}{l} \text{projection of} \\ y \text{ on to the space} \\ \text{spanned by } \phi \end{array} \right].$$

<u>The ML interpretation of the OLS.</u>

We saw that under the $L_2$ loss the optimal function is the conditional expectation.

Thus, for any data, the optimal

$$h(x) = \underline{E \left[ y | x \right]}.$$

(114).

Let us assume the

$$f(y|x) \sim N\left(h(x), \sigma\right).$$

Thus, for $h(x)$ can be found using ML estimate of this conditional distribution.

$$h(x) = W^T \phi$$

Now, Likelihood for the above distribution is

$$L = \prod_{i=1}^{n} N\left(\hat{h}(x_i), \sigma\right)$$

$$\log L = \sum_{i=1}^{n} \ln N\left(y_i \mid W^T \phi(x_i), \sigma\right)$$

$$= \frac{n}{2}(\ ) - \frac{n}{n}\ln(2\pi) - \frac{1}{\sigma}\sum_{i=1}^{n}\left(y_i - W^T \phi(x_i)\right)^2$$

Thus, $\quad \hat{W}_{ML} = \arg\min_{W} \sum_{i=1}^{n}\left(y_i - W^T\phi(x_i)\right)^2$

$$\Rightarrow W^*_{OLS} = \hat{W}_{ML} \quad \text{with the Gaussian}$$
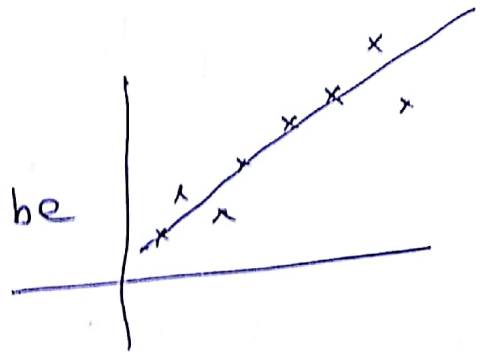
assumption for the posterior.

(115).

## The idea of regularization.

WOLS can also be viewed as modelling the noise added to the approximator (minimizing) to get to the actual output.

$$y = h(x) + \epsilon \qquad y | x \sim N(h(x), \sigma)$$

$$\epsilon \sim N(0, \sigma).$$

Then WOLS or WMLE can be a interpretted as finding the parameters of $h(x)$ that would minimize the perturbed & the actual data.

Or thee estimating conditional mean of $y|x$.

Scanned by CamScanner

# Regularization.

We know that the risk of a classifier

$$R(h) = \mathbb{E}_{P_{xy}}[L]$$

Assuming squared error loss,

$$R(h) = \iint (h(x) - y)^2 \, f_{xy}(x,y) \, dx \, dy.$$

Consider $\{h(x) - y\}^2 = \{h(x) - \mathbb{E}[y|x] + \mathbb{E}[y|x] - y\}^2$

$$= \{h(x) - \mathbb{E}[y|x]\}^2 + 2\{h(x) - \mathbb{E}[y|x]\}$$
$$\cdot \{\mathbb{E}[y|x] - y\}$$
$$+ \{\mathbb{E}[y|x] - y\}^2.$$

Consider the cross term in the risk integral.
$$\underset{\text{fng } x \text{ alone}}{\frown}$$

$$2\iint \{h(x) - \mathbb{E}[y|x]\} \{\mathbb{E}[y|x] - y\} \, dy \, f(x,y) \, dy \, dx$$

$$\iint \{h(x) - \mathbb{E}[y|x]\} \cdot \left[ \int (\mathbb{E}[y|x] - y) \cdot f_{xy}(x,y) \, dy \, dx \right]$$

$$(\text{1}7). = (\,)\{\cdot \mathbb{E}[y|x] - \mathbb{E}[y|x]\} = 0$$

$$\therefore R(h) = \int \left\{ \underbrace{h(x) - E[Y|x]}_{(1)} \right\}^2 f_x(x) \, dx$$

$$+ \int \left\{ \underbrace{E[Y|x] - Y}_{(2)} \right\}^2 f_x(x) \, dx.$$

Observe that only (1) can be tweaked/learned algorithmically.

Nothing can be done about the term (2).

(2) actually signifies the error that is in the data — could be feature noise | label noise etc. (noise)

Thus, risk can be minimized only up to a factor given by term (2) [usually unknown].

Thus, let us focus on term (1).

[observe that risk minimizes when term (1) is minimized @ $h(x) = E[Y|x]$ .. another proof].
(118)

Now, Risk term of importance

$$R_1(h) = \int \left\{ h(x) - E[Y|x] \right\}^2 f_x(x) \, dx.$$

Observe that $R_1$ is an expectation w.r.t $f_x(x)$

Roughly, it is a measure of the 'error' incurred by the classifier $h(x)$ averaged over all possible datasets sampeled from $f_x(x)$.

In that sense, $R_1(h) = E_D[error] = E_D \left[ h(x) - E(Y|x) \right.$

(expected error over all possible choices of datasets]

Now, $R_1 = \int \left\{ h(x) - E[Y|x] \right\}^2 f_x(x) \, dx$

consider $\left\{ h(x) - E[Y|x] \right\} = \left\{ h(x) - E_D[h(x)] + E_D[h(x)] - E[Y|x] \right\}$

$$E_D = \int h(x) f_x(x) \, dx$$

[Average performer classifier over all datasets].

Plugging the above to $R_1$ & evaluating the integral,

(119)

$$R_1(h) = \int \overbrace{\left[E_D(h(x)) - \underset{E[y|x]}{E(y|x)}\right]^2} f(x)\,dx - \underline{\underline{Bias^2}}$$

$$+ \int E_D\left[\left\{h(x) - E_D(h(x))\right\}^2\right] f(x)\,dx$$
$$\underline{\qquad - Variance}$$

Thus, Risk = $Bias^2$ + Variance + noise (irreducible).

Bias — The extent to which the average approximator differs from the $(h(x))$ (over all data) desired output. $(E[y|x])$

Variance — The extent to which the approximator is sensitive to the choice of $h(x)$ a particular dataset.
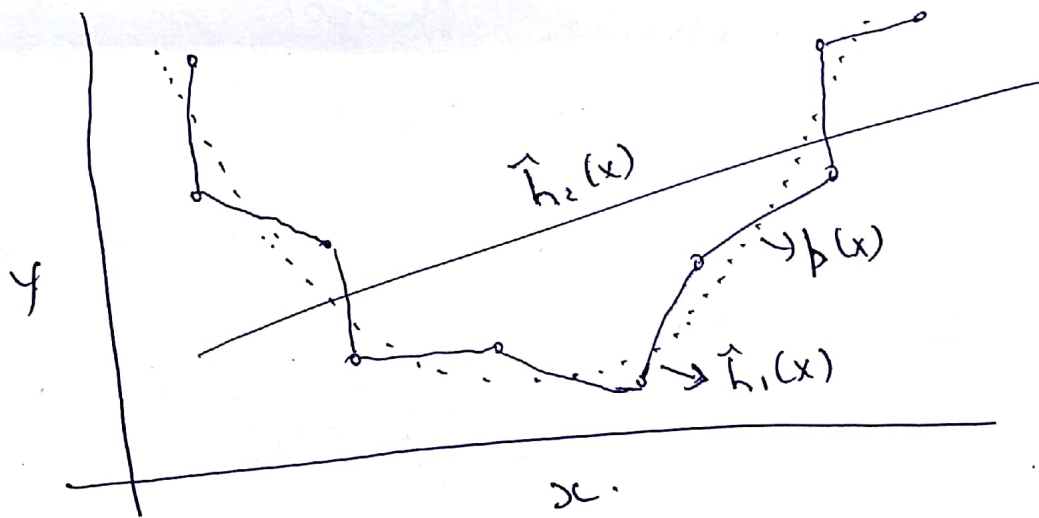
A Good ML algorithm — reduce the error-Reduce both the bias & variance.

(120).

# Behaviour of Bias & variance.

Let us consider a regression problem, with
y lying on a smooth parabola.

Data is generated by perturbing $\not{p}$ target
points with a small noise.

$$Y_i = \phi(x_i) + \epsilon_i \qquad \phi(x) = ax^2 + c$$



The optimal mapping would be $E[y|x]$ but
(regressor)

Let us choose two extreme cases of classifiers
(regressors)

& see what happens to Bias & variance.

(121).

Let $\hat{h}_1(x)$ be a line choosen independently of data.

Consider

B Variance $(\hat{h}_1(x)) \doteq \mathcal{E} \int \hat{h}_1(x) f(x) dx$

$\doteq \hat{h}(x) \int f(x) dx = \hat{h}(x).$

$\therefore \quad h(x) = \mathcal{E}_D(h(x)) \implies \underline{\text{Variance} = 0.}$

However the bias will be typically large.
$$\underline{[\text{underfitting}]}$$

In the other extreme, choose

$$\hat{h}_2(x_i) = p(x_i) + \epsilon_i$$

In this case $\quad bias = 0. \, [\text{over the observed data}]$

$$\therefore \quad \mathcal{E}_D\left[\hat{h}_2(x)\right] = \mathcal{E}_D\left[p(x) + \epsilon\right]$$

$$= \mathcal{E}_D\left[p(x)\right] = \mathcal{E}_D[$$

$$\mathcal{E}[Y|x]$$

$$\left(\begin{array}{c}\text{over the}\\\text{data}\end{array}\right).$$

However Variance :

$$\mathcal{E}_D\left[\epsilon^2\right] \text{— variance of the noise which is}$$
$$\text{typically high.}$$
$$[\text{verify the algebra}] \qquad [\text{overfitting}].$$
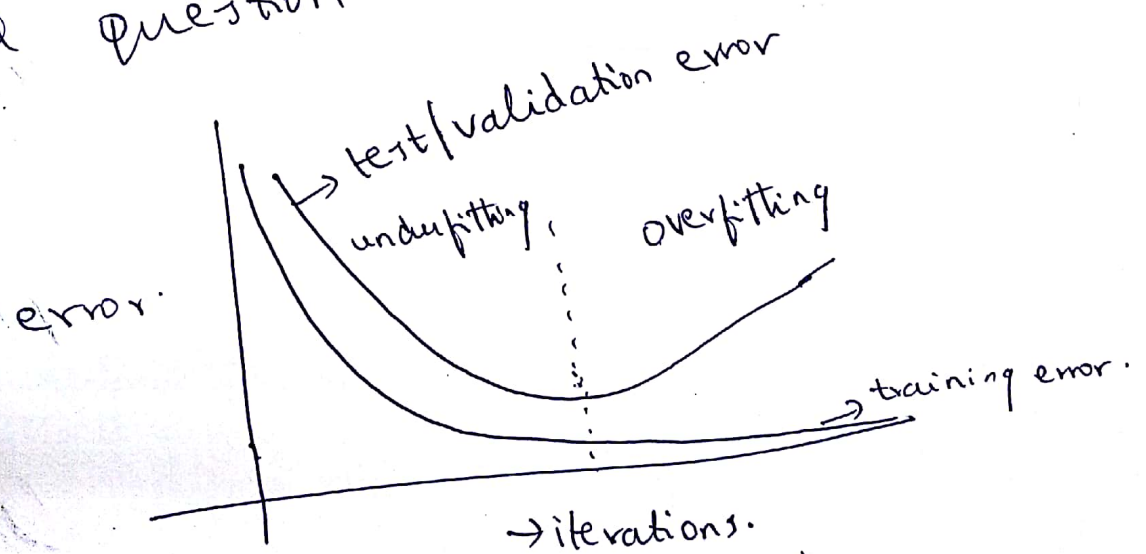
(122).

Thus, ~~tht~~ One cannot simultaneously minimize bias & variance together - usually lowering one results in increasing the other.

Thumb rule: $h(x)$ that closely fits the given data (training data) — Low bias | high variance (overfitting).

Variance can be lowered by smoothing $h(x)$, but taken too far leads to high bias & low-variance [underfitting].

Question: How to minimize both of them?
Central Question in ML.



Typical training graph.

$(103)$

Many possible ways.

1. Given very large dataset, use high-capacity models reducing bias, see 'enough' samples to reduce the variance.

   [DNN idea].

2. Use prior knowledge on true labels to constrain models [so that bias is not too much]

   often termed as <u>Regularization</u>.

   Bias minds take better decisions !!!

   We'll look at one strategy for Reg in this course : <u>Parameter</u> <u>norm</u> <u>penalties</u>.

# Regularized least squares.

__Idea:__ Real life data is sparse (occupies - a lower dim manifold in the data space).

In such cases, complex models [with higher bias] can fit the data well. In other words, the error is dominated by the variance component.

Thus, intuitively it is better to trade bias for lower variance.

Power series of sin
→ monotone x.

[ Example of fitting a sin with noise by a very high degree polynomial].

Lower variance ≡ restricting the value of the parameter space.

One way to achieve that is to add a penalty term on the parameter in the Least square cost.

$$L = \|y - W^T \phi(x)\| + \frac{\lambda}{2} W^T W$$

we want $W$ to be small.

$\lambda \rightarrow$ regularization constant.

$$W^*_{reg} = (\lambda I + \phi^T \phi)^{-1} \phi^T y.$$