# * Probabilistic Approach and PCA

$$D = \{ (x_1, y_1), (x_2, y_2), \ldots \ldots (x_n, y_n)\}$$
↳ labels

Train a new function

$$y = f(x)$$
actual

Predicted → $\hat{f}$ : algorithm : best
training

↙

close to f

If x near D → less error
If x not near D → high error

$$\hat{f}(x; D)$$
↳ training data

$$E_D\left[ (y - \hat{f}(x; D))^2 \right] \quad \} \quad \text{average error / expected error}$$
$$D$$

↳ Various Datasets

→ check error over x

↓

take average

then subtract from y

So we get average ERROR

{ f(x) : True Label }

{ $\hat{f}$ → predicted }
{ E → Expected }

perfect model we learnt

Variance :-
$$E(\text{Prediction} - \text{Expected})^2$$
Prediction = made by model
for specific
training Set

Expected Prediction = avg
Prediction made by
model

across many
training sets.

## * Bias Variance Decomposition :-
fixed point
training set

Variance – variance of $\hat{f}(x; D)$ can be defined as follows.
$$\text{Var}_D[\hat{f}(x; D)] = E_D\left[ (E_D[\hat{f}(x; D)] - \hat{f}(x; D))^2 \right]$$

→ y²

$$\hat{f}(x; D) \quad \text{specific}$$

$$E_D[\hat{f}(x; D)] \quad \text{→ range} \quad : \text{expected outputs}$$

$$y = E_D[\hat{f}(x; D)] - \hat{f}(x; D)$$

difference in expected and
specific output

High Variance → Sensitivity to D (data points)
High Variance → overfitting → too complex
Low Variance → Stable across all

overfit → model learns noise and details

→ performs good on training data poor on unseen

$\mathbb{E}$ → denote expected value or mean

Training error : Low
Test error : High

## Bias: we also define Bias

$$\text{Bias}_D[\hat{f}(x;D)] = \underbrace{\mathbb{E}_D[\hat{f}(x;D)]}_{\text{expected}} - \underbrace{f(x)}_{\text{True}}$$

High Bias → High Inductive Bias → Underfitting

High Bias → too simplistic → underfits (does not capture all details)

Low Bias → Close to true value

## Bias and Variance Decomposition

$$\underbrace{\mathbb{E}_D[(y-\hat{f}(x;D))^2]}_{\text{Average MSE}} = \underbrace{(\text{Bias}_D[\hat{f}(x;D)])^2}_{(\text{Bias})^2} + \underbrace{\text{Var}_D[\hat{f}(x;D)]}_{\text{Variance}}$$

Average MSE

mean Squared error

{ Can be Hyperparameter }

Two Sources of error

In reality,

$$\underbrace{\mathbb{E}_{D,\epsilon}[(y-\hat{f}(x;D))^2]}_{\substack{\text{Data} \\ \text{and} \\ \text{noise}}} = (\text{Bias}_D[\hat{f}(x;D)])^2 + \text{Var}_D[\hat{f}(x;D)] + \underbrace{\sigma^2}_{\substack{\text{little} \\ \text{random} \\ \text{noise}}}$$

mean Squared Error

$\underset{\text{noise}}{y = f(x) + \varepsilon}$

→ Average
$E(\varepsilon) = 0$
$\underset{\text{sigma square}}{\text{Var}(\varepsilon) = \delta^2}$

# * The maximum Likelihood View :-

$x$ :- height

$y$ - male/female (biological sex)



WOMEN    median    median    MEN

→ Height

given individual with height $x$ is what male/female?

→ probability

$$X \sim P(y/x)$$

→ $y$ or $x$

$$m_\omega(x) \approx P(y/x)$$

↓

model

$$\underset{x}{\longrightarrow} \boxed{LN} \longrightarrow \hat{P}(y/x)$$

$m_\omega(x)$ :    $\omega$ : best possible $\omega$'s

D: data

**Logistic Neuron**

which $\omega$ has highest Prob    D: Data

$$\hat{\omega} = \text{argmax}_\omega \; P(\omega/D)$$

$$P(\omega/D) = \frac{P(D/\omega)\, P(\omega)}{P(D)}$$

$P(D)$ → independent of set of weights

$$\text{argmax}_\omega \; P(D/\omega)\, P(\omega)$$

*Deriving the nll loss

IID → identically distributed
~~independent~~ distribution assumption

$$P(D|\omega) - \prod_{j=1}^{n} P(\{y_i | x_i\} | \omega)\}$$

independence

If $y_i = 1$, then

Otherwise, if $y_i = 0$

$$P(\{y_i | x_i\}) | \omega) = m_\omega(x_i)$$

This can be expressed
in one
equation as :

$$P(\{y_i | x_i\}) | \omega) = 1 - m_\omega(x_i).$$

IMPORTANT

$$\{ P(\{y_i | x_i\}) | \omega) = m_\omega(x_i)^{y_i} (1 - m_\omega(x_i))^{1-y_i}$$

$$\text{argmax}_\omega \, P(D|\omega) = \text{argmax}_\omega \prod_{j=1}^{n} m_\omega(x_i)^{y_i} (1 - m_\omega(x_i))^{1-y_i}$$

increase probability of true label

$y_i = 1$  $y_i = 0$

quivalent → $\text{argmax}_\omega \, \log P(D|W)$

Example :- Find Email Spam/Not Spam

$y = 1 \to$ Spam
$y = 0 \to$ Not Spam

| Email ID | Features | y Label | Pred Prob ability |
|---|---|---|---|
| 1 | [5,10,3] | 1-Spam | 0.9 |
| 2 | [2,4,1] | 0 | 0.2 |
| 3 | [3,6,2] | 1-Spam | 0.7 |

Maximize

$\text{arg max}_\omega$

$$\sum_{j=1}^{n} (y_i (m_\omega(x_i) + (1-y_i)(1-m_\omega)x_i)$$

Email 1 → $y=1$, $m_\omega(x) = 0.9$
$$P(1|x_1, \omega) = m_\omega(x_1)^1 (1-m_\omega(x_1))^0$$
$$= 0.9^1 \cdot (1-0.9)^0 = 0.9$$

Maximize

$$\text{arg max}_\omega \, P(W/D)$$

minimize its negation

# * Naive Bayes Classifier

| 3 | 5 | 3 | 6 | 1 | 7 |
|---|---|---|---|---|---|
| 9 | 4 | 0 | 9 | 1 | 2 |

$X$: 28×28 Pixels.

[0:255] greyscale

[0,1] → 784 pixels

$y_i \{0, \ldots 9\}$
labels → numbers 0,1,2...9

## Probability Distribution

$$P(Y/x)$$
↳ input

$$\hat{y} = \arg\max_y (P(y|x)) = \arg\max_y \frac{P(X/y) P(y)}{P(X)}$$

$$= \arg\max_y P(X/y) \cdot P(y)$$

## Estimate

$$P(x|y) = \prod_{t=1}^{784} P(X_t/y)$$

↳ label 5
1,5,t
↳ if its t
current
pixel

$$P\left(X_t = 1 \middle| \begin{array}{c} y=0 \\ y=1 \\ \vdots \\ y=9 \end{array}\right) \quad t=1\ldots784$$

→ Given that pixel is 1
for label 5
How many such pixels

$$P(X_t = 1/y = 5)$$
Pixel t

$n_5$ · # of images with label 5

$n_{0,5,t}$ → how may pixels 0 ⎱
$n_{1,5,t}$ → how may pixels 1 ⎰

$$\frac{n_{1,5,t}}{n_5}$$