$f$ function

$f_\theta$ ← parameters of the function

i/p data → [ Model $\theta$ ] → o/p ↑ $f$

- real nos. – Regression
- classes – Classification
- distribution
  └ Probabilistic Models

$\theta = \{\theta_1, \theta_2 \ldots \theta_n\} \in \mathbb{R}^n$

i/p data: structured / unstructured

① $f = wx + b \Rightarrow \theta = \{w, b\}$
└ convex function
└ good fit.

$y = wx + b$

→ one setting
$\theta \in \Theta$
└ set of all possible parameter values

## Maximum Likelihood Estimate (MLE)

Data: $\{(x_i, y_i)\}_{i=1}^N$

Model: $f_\theta(x) = y$, linear → $f_\theta(x) = w^T x + b$
                                     └ $\theta = \{w, b\}$

Training Data { $x_1$ $x_2$ $x_3$ $y$
Test Data {

$x[\ ] \quad w[\ ]$

$y$: ground truth

$\hat{y}$: prediction / generated output

### Loss function $\ell(y, \hat{y}) = \sum_{i=1}^N (y_i - \hat{y}_i)^2$ ← convex

$y - \hat{y}(x)$ – individual terms can be (-)ve

$$\boxed{\hat{\theta} \text{ or } \theta^* = \arg\min_\theta \frac{1}{N} \sum_{i=1}^N \ell(y_i, \hat{y}_i)}$$

### Linear Gaussian Model. : $N(w^T x + b, \sigma^2)$ ← Distribution

$P_\theta(y|x) = N(y; w^T x + b, \sigma^2)$ ← probability value.

### MLE - Non ML perspective

$\mathcal{D} = \{x_1, x_2, \ldots, x_n\}$ → assume a set of distribution $\alpha$

$P_\theta : \theta \in \Theta$

assume that $\mathcal{D}$ was sampled from a member of this family

Goal: recover $\hat{\theta}$

likelihood of data
MLE : $\theta_{MLE} = \arg\max_\theta P_\theta(\mathcal{D}) = \arg\max_\theta \prod_{i=1}^N P_\theta(x_i)$

$\{\mu_{MLE}, \sigma^2_{MLE}\} = \underset{\mu, \sigma^2}{\text{argmax}} \prod_{i=1}^{N} N(x_i; \mu, \sigma^2)$

<u>Taking log</u>, (monotonic fn; maximizing log = maximizing LL.),

$\underset{\mu, \sigma^2}{\text{argmax}} \sum_{i=1}^{N} \log N(x_i; \mu, \sigma^2)$

or, $\underset{\mu, \sigma^2}{\text{argmax}} \sum_{i=1}^{N} \log \left[ \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left( \frac{-(x_i - \mu)^2}{2\sigma^2} \right) \right]$

or, $\underset{\mu, \sigma^2}{\text{argmax}} \sum_{i=1}^{N} \log \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2} \frac{(x_i - \mu)^2}{\sigma^2}$

or, $\underset{\mu, \sigma^2}{\text{argmax}} \quad \frac{-N}{2\sigma^2} \sum_{i=1}^{N} (x_i - \mu)^2 - \frac{N}{2} \log \sigma^2 = F$

<u>Now,</u>

$\frac{\partial F}{\partial \mu} = \sum_i \frac{(x_i - \mu)}{\sigma^2} = 0 \Rightarrow \frac{\sum x_i}{\sigma^2} - \frac{N\mu}{\sigma^2} = 0$

or, $\frac{\sum_i x_i}{\sigma^2} = \frac{N\mu}{\sigma^2}$

or, $\boxed{\mu = \frac{\sum x_i}{N}}$

$\frac{\partial F}{\partial \sigma^2} = \sum_i \frac{(x_i - \mu)^2}{2\sigma^2} - \frac{N}{2} \frac{1}{\sigma^2} = 0$

or, $\boxed{\sigma^2 = \frac{\sum_i (x_i - \mu)^2}{N}}$

(Shannon's)

<u>Entropy</u> : no. of bits needed to encode information.

$H(p) = -\sum p_i \log p_i$

<u>Cross entropy</u> (CE)

$\mathbb{E}[x_i] \underset{N \to \infty}{\approx} \frac{1}{N} \sum_i x_i$    infinite data

$H(p,q) = - \sum p_i \log q_i = \underset{(p)}{\mathbb{E}}[- \log(q_i)]$    generated distribution
using Monte
Carlo estimate

ground truth distribution

$= \frac{1}{N} \sum_{i=1}^{N} \log(q(x_i))$

$\uparrow MLE = \downarrow CE$

<u>$D_{KL}$ (KL Divergence)</u> $\Rightarrow D_{KL}(p \| q) = H(p,q) - H(p)$