

Foundations of Machine Learning AI2000 and AI5000

FoML-04

Maximum Likelihood Principle

Dr. Konda Reddy Mopuri

Department of AI, IIT Hyderabad

July-Nov 2025



భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad



So far in FoML

- What is ML and the learning paradigms
- Probability refresher
 - Random variables, Bayes Theorem, Independence, Expectation, Variance



Maximum Likelihood Principle



భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad



Maximum Likelihood Principle

- Widely used technique for optimizing model parameters



భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad



Maximum Likelihood Principle

- Given - Dataset of N independent observations $D = \{x_1, x_2, \dots, x_N\}$



Maximum Likelihood Principle

- Goal: recover the probability distribution that may have generated this dataset



Maximum Likelihood Principle

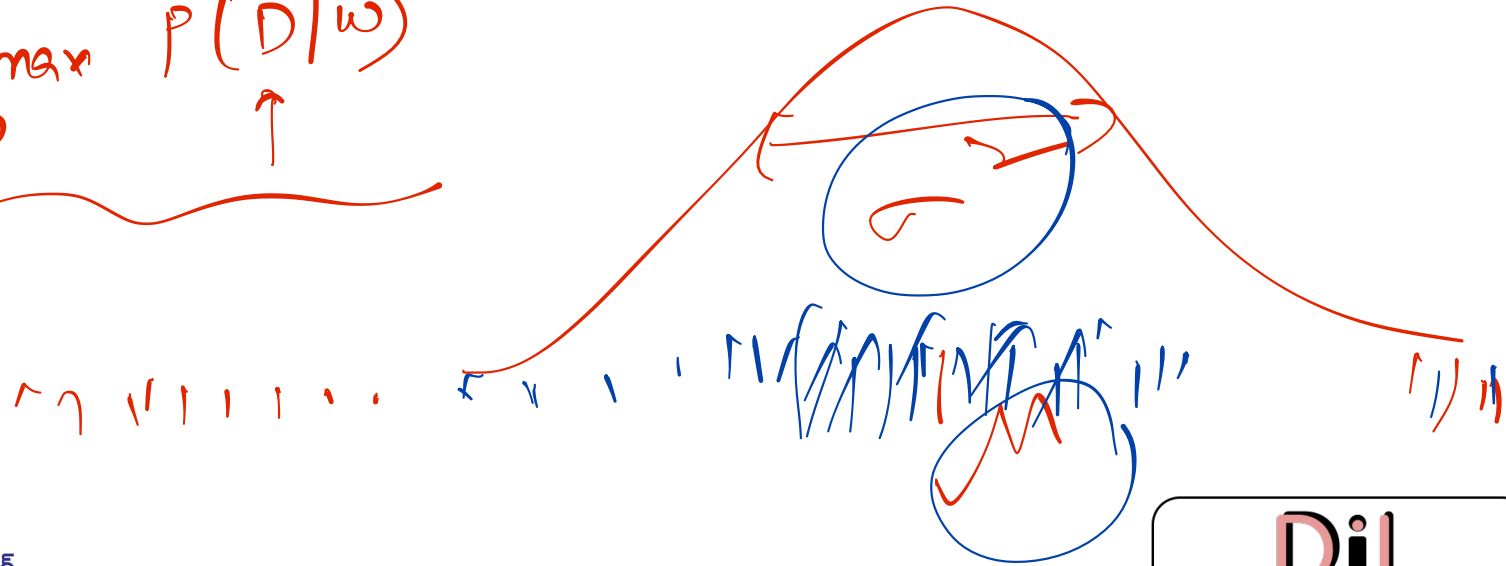
- Goal: recover the probability distribution that may have generated this dataset
- Likelihood of the dataset $p(D|w)$



Maximum Likelihood Principle

- The most likely 'explanation' of D is given by w_{ML} that maximizes the likelihood function

$$\underline{w_{ML}} = \underset{w}{\operatorname{argmax}} P(D|w)$$



Maximum Likelihood Principle

- The iid assumption - each $x_i \in D$ is independently distributed according to the same distribution conditioned on w



Maximum Likelihood Principle

- The iid assumption - each $x_i \in D$ is independently distributed according to the same distribution conditioned on w

The joint distribution

$$P(D|w) = P(x_1, x_2, \dots, x_N | w) = \prod_{i=1}^N P(x_i | w)$$



Maximum Likelihood Principle

$$\underline{\mathbf{w}_{ML}} = \arg \max_{\mathbf{w}} \underline{p(D|\mathbf{w})} = \arg \max_{\mathbf{w}} \prod_{i=1}^N \underline{p(x_i|\mathbf{w})}$$



Maximum Likelihood Principle

$$\mathbf{w}_{ML} = \arg \max_{\mathbf{w}} p(D|\mathbf{w}) = \arg \max_{\mathbf{w}} \prod_{i=1}^N p(x_i|\mathbf{w})$$

- Numerical underflow

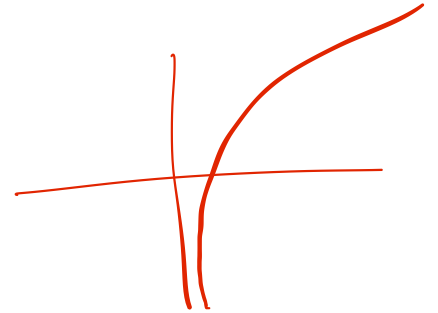


Maximum Likelihood Principle

$$\mathbf{w}_{ML} = \arg \max_{\mathbf{w}} p(D|\mathbf{w}) = \arg \max_{\mathbf{w}} \prod_{i=1}^N p(x_i|\mathbf{w})$$

- Numerical underflow
- Maximize the log-likelihood \rightarrow

$$\log \prod_{i=1}^N p(x_i|\omega)$$



Maximum Likelihood Principle

$$\mathbf{w}_{ML} = \arg \max_{\mathbf{w}} p(D|\mathbf{w}) = \arg \max_{\mathbf{w}} \prod_{i=1}^N p(x_i|\mathbf{w})$$

- Numerical underflow
- Maximize the log-likelihood

$$\mathbf{w}_{ML} = \arg \max_{\mathbf{w}} \log \prod_{i=1}^N p(x_i|\mathbf{w})$$

Error function:

$$E(D; \mathbf{w}) = -\log p(D|\mathbf{w}) = -\sum_{i=1}^N \log p(x_i|\mathbf{w})$$

NLL



MLE for Gaussian Distributions

- iid Gaussian distributed real variables $D =$

$$p(x|\mathbf{w}) = \mathcal{N}(x|\mu, \sigma^2)$$

$$p(D|\mathbf{w}) = p(D|\mu, \sigma^2) =$$

$$\left(\frac{1}{2\pi\sigma^2}\right)^{1/2} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$$



MLE for Gaussian Distributions

- iid Gaussian distributed real variables $D =$

$$p(x|\mathbf{w}) = \mathcal{N}(x|\mu, \sigma^2)$$

$$p(D|\mathbf{w}) = p(D|\mu, \sigma^2) =$$

$$\left(\frac{1}{2\pi\sigma^2}\right)^{N/2} \prod_{i=1}^N e^{-\frac{1}{2\sigma^2}(x_i - \mu)^2}$$

log likelihood =

$$-\frac{N}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^N (x_i - \mu)^2$$



MLE for Gaussian Distributions

- Estimate the model parameters $\underline{\mu_{ML}}, \underline{\sigma_{ML}^2} = \arg \max_{\mu, \sigma^2} \log p(D|\mu, \sigma^2)$

$$\frac{\partial}{\partial \mu} \left(\right) = 0 = \underbrace{-\frac{1}{2\sigma^2}}_{\text{}} \underbrace{(2)}_{\text{}} \sum_{i=1}^N (x_i - \mu)$$

$$\sum_{i=1}^N (x_i - \mu) = 0$$

$$\mu_{ML} = \frac{1}{N} \sum_{i=1}^N x_i$$



MLE for Gaussian Distributions

- Estimate the model parameters $\mu_{ML}, \sigma_{ML}^2 = \arg \max_{\mu, \sigma^2} \log p(D|\mu, \sigma^2)$

$$\frac{\partial}{\partial \sigma^2} (\quad) = 0$$

$$- \frac{N}{2} \frac{1}{\sigma^4} + \frac{1}{2\sigma^4} \sum_{i=1}^N (x_i - \mu)^2 = 0$$

$$- \frac{N}{2} \sigma^{-2} + \frac{1}{2} \sum_{i=1}^N (x_i - \mu)^2 \sigma^{-4} = 0$$

$$\sigma_{ML}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$$



MLE for Gaussian Distributions

$$\mu_{ML}, \sigma_{ML}^2 = \arg \max_{\mu, \sigma^2} \log p(D|\mu, \sigma^2)$$

- How well do these estimates represent the true parameters?



MLE for Gaussian Distributions

$$\mu_{ML}, \sigma_{ML}^2 = \arg \max_{\mu, \sigma^2} \log p(D|\mu, \sigma^2)$$

- How well do these estimates represent the true parameters?
- Note that these are functions of the data sample

$$\begin{array}{lcl} \underline{D_1} = \{x_1, \dots, x_N\}^1 & \longrightarrow & \mu_{ML}^1, \sigma_{ML}^2{}^1 / \\ \underline{D_2} = \{x_1, \dots, x_N\}^2 & \longrightarrow & \mu_{ML}^2, \sigma_{ML}^2{}^2 / \\ \vdots & & \vdots \\ \underline{D_M} = \{x_1, \dots, x_N\}^M & \longrightarrow & \mu_{ML}^M, \sigma_{ML}^2{}^M / \end{array}$$



MLE for Gaussian Distributions

$$\mu_{ML}, \sigma_{ML}^2 = \arg \max_{\mu, \sigma^2} \log p(D|\mu, \sigma^2)$$

- How well do these estimates represent the true parameters?
- Note that these are functions of the data sample
 - → expected values of these estimates

$$\begin{aligned} E[\mu_{ML}] &- ? \\ E[\sigma_{ML}^2] &- ? \\ D &\sim P(D|\mu, \sigma^2) \end{aligned}$$



MLE for Gaussian Distributions

- ML estimate of the mean

$$\begin{aligned} E \left[\frac{1}{N} \sum_{i=1}^N x_i \right] &= \frac{1}{N} \sum_{i=1}^N E_{D \sim P(D) | \mu, \sigma^2} x_i \\ &= \frac{1}{N} \sum_{i=1}^N E_{x_i \sim p(x | \mu, \sigma^2)} (x_i) = \frac{1}{N} \sum_{i=1}^N \mu = \mu \end{aligned}$$



MLE for Gaussian Distributions

- Bias of the “ML estimate of the mean”

$$\underline{\underline{E[\mu_{ML}] - \mu}}$$

$$\mu - \mu$$

$$= 0$$

$$x_i \sim \underline{p(x|\mu, \sigma^2)}$$



MLE for Gaussian Distributions

- ML estimate of the variance

$$E \left[\frac{1}{N} \sum_{i=1}^N \left(x_i - \frac{1}{N} \sum_{j=1}^N x_j \right)^2 \right]$$

$$= \frac{1}{N} \sum_{i=1}^N E \left[\left(x_i - \frac{1}{N} \sum_{j=1}^N x_j \right)^2 \right] = \frac{1}{N} \sum_{i=1}^N E \left[x_i^2 - \frac{2}{N} x_i \sum_{j=1}^N x_j + \frac{1}{N^2} \sum_{l=1}^N \sum_{k=1}^N x_l x_k \right]$$

$$= \frac{1}{N} \sum_{i=1}^N \left[\mu^2 + \sigma^2 - \frac{2}{N} (\sigma^2 + \mu^2 + (N-1)\mu^2) + \frac{1}{N^2} (N(\sigma^2 + \mu^2) + N(N-1)\mu^2) \right]$$

$$= \frac{1}{N} \sum_{i=1}^N \left[\mu^2 + \sigma^2 - \frac{2}{N} (\sigma^2 + N\mu^2) + \frac{1}{N^2} (N\sigma^2 + N^2\mu^2) \right]$$

$$= \mu^2 \left(1 - 2 + 1 \right) + \sigma^2 \left(1 - \frac{2}{N} + \frac{1}{N} \right) = \frac{N-1}{N} \sigma^2$$

$$E[x_i x_j] = \text{Cov}(x_i, x_j) + E[x_i]E[x_j]$$

$$i=j \rightarrow \sigma^2 + \mu^2 \quad i \neq j \rightarrow \mu^2$$



MLE for Gaussian Distributions

- Bias of the “ML estimate of the variance”

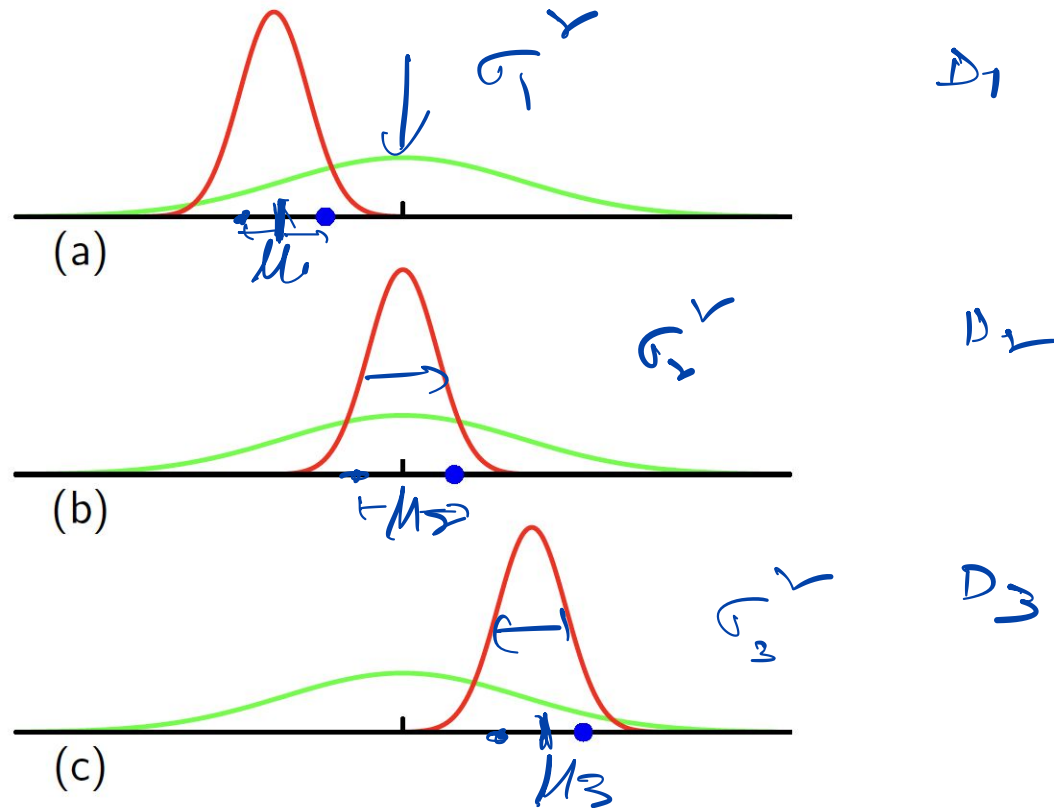
σ_{ML}^2 is biased $\frac{N-1}{N} \sigma^2$

$$\sigma_{ML-UB}^2 = \frac{N}{N-1} \sigma_{ML}^2$$

$$\arg \max_w p(D/w)$$



Bias in variance estimate



Regression example



భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad



MLE for Regression (curve fitting)

$$\hat{f} = \sum_{i=1}^{n-1} w_i x^i + w_0$$

- Given data D

$$D = \{(x_1, t_1), (x_2, t_2), \dots, (x_N, t_N)\} = \{\mathbf{x}, \mathbf{t}\}$$

$$MSE = \frac{1}{N} \sum_{i=1}^N (t_i - \hat{t}_i)^2$$



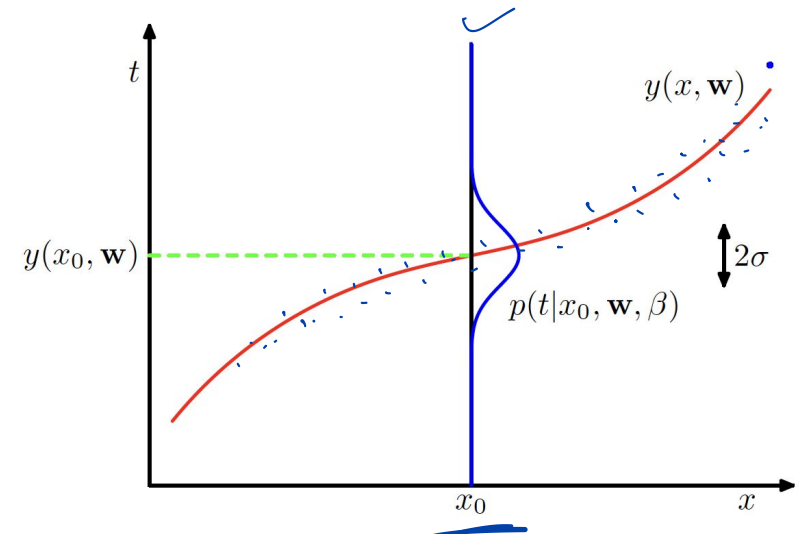
MLE for Regression (curve fitting)

$$D = \{(x_i, t_i)\}$$

- Given data D $D = \{(x_1, t_1), (x_2, t_2), \dots, (x_N, t_N)\} = \{\mathbf{x}, \mathbf{t}\}$
- Assume the data is generated by

$$t = y(x, \mathbf{w}) + \sigma \cdot \epsilon, \quad \epsilon \in \mathcal{N}(0, 1)$$

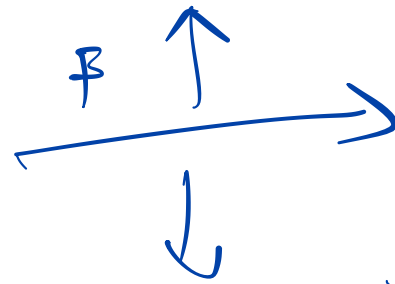
measurement error/noise



MLE for Regression (curve fitting)

- Target distribution $p(t|x, \mathbf{w}, \beta) = \mathcal{N}(t|y(x, \mathbf{w}), \beta^{-1})$

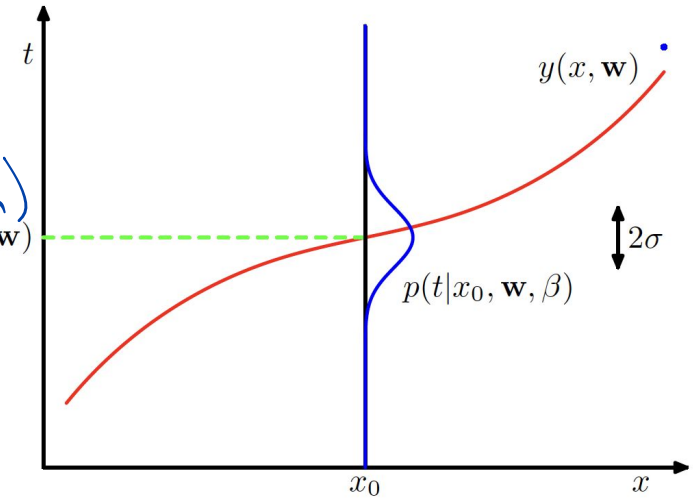
$$\left\{ \begin{array}{l} x_1, t_1 \\ \vdots \\ x_n, t_n \end{array} \right\} \quad \underline{\underline{(t_1 - \hat{t}_1)}}$$



$$p(\underline{D} | \underline{w}, \underline{\beta})$$

$(x_1, t_1) \dots (x_n, t_n)$

$y(x_0, \mathbf{w})$



$$y = \underline{w}^T x + b$$



$\{y, t\}$

MLE for Regression (curve fitting)

- Target distribution
- log likelihood

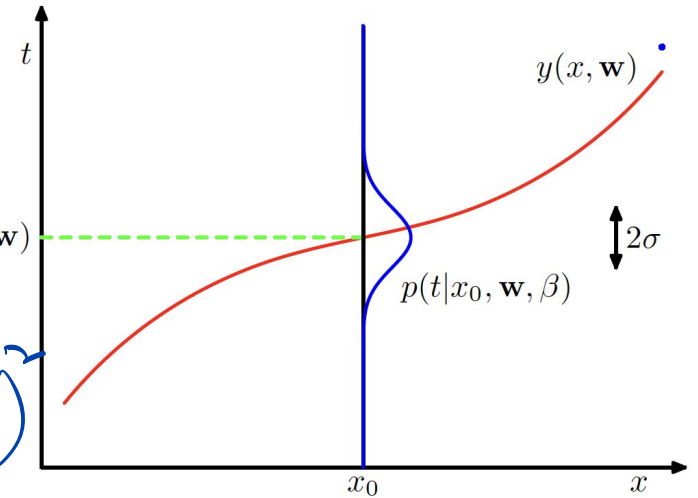
$$p(t|x, \mathbf{w}, \beta) = \mathcal{N}(t|y(x, \mathbf{w}), \beta^{-1})$$

$$\log p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta^{-1})$$

$$= \log \left(\prod_{i=1}^N \sqrt{\frac{\beta}{2\pi}} e^{-\frac{\beta}{2}(t_i - y_i)^2} \right)$$

$$= \sum_{i=1}^N \log \left(\sqrt{\frac{\beta}{2\pi}} e^{-\frac{\beta}{2}(t_i - y_i)^2} \right)$$

$$= N \log \sqrt{\frac{\beta}{2\pi}} - \frac{\beta}{2} \sum_{i=1}^N (t_i - y_i)^2$$



MLE for Regression (curve fitting)

$D = \{x_i, t_i\}$

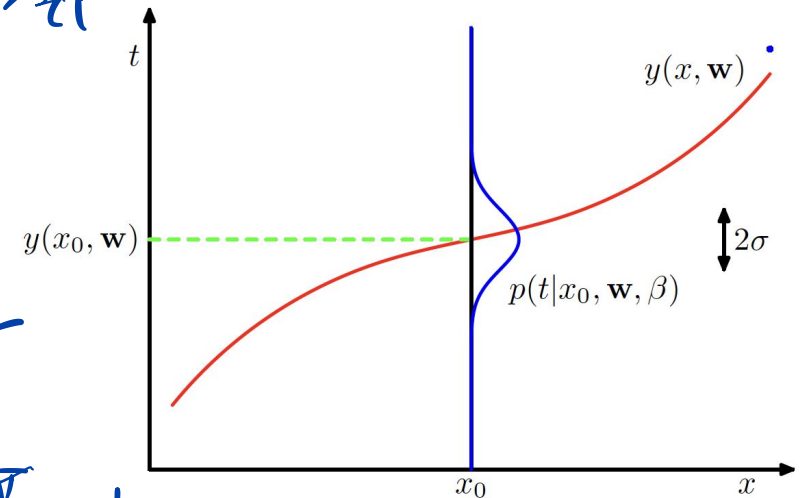
- Minimize the NLL w.r.t the parameters \underline{w} and $\underline{\beta}$

$$\underline{NLL} = \frac{\beta}{2} \sum_{i=1}^N [t_i - y(x_i, w)]^2 - \frac{N}{2} \ln \frac{\beta}{2\pi}$$

$$\frac{\partial (\cdot)}{\partial w} = 0 \quad \text{argmin}_w \quad \underline{MLE}$$

$$\frac{\partial}{\partial \beta} (\cdot) = \frac{1}{2} \sum_{i=1}^N [t_i - y(x_i, w)]^2 - \frac{N}{2} \ln \frac{\beta}{2\pi} = 0$$

$$\Rightarrow \frac{N}{\beta} = \sum_{i=1}^N [t_i - y(x_i, w)]^2$$



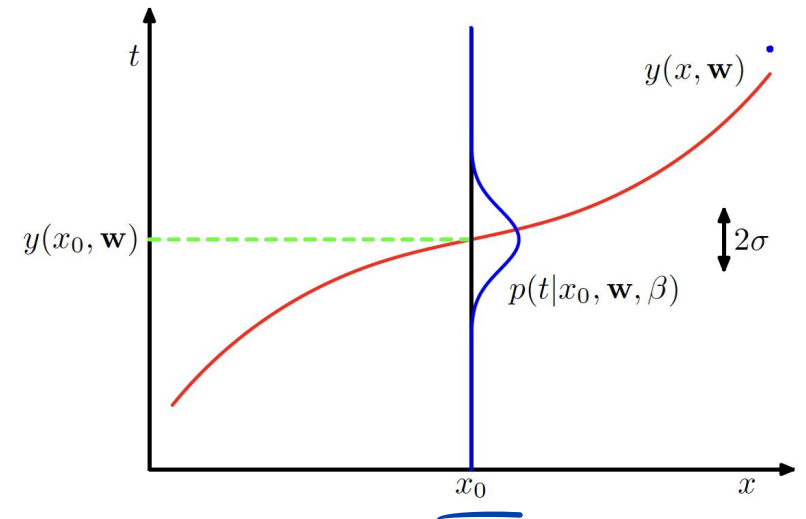
MLE for Regression (curve fitting)

- The predictive distribution

$$p(t | y(x', w_{ML}), \beta_{ML}) = \mathcal{N}(t | y(x', w_{ML}), \sigma^2_{ML})$$

for any test sample x'

if one wants
the point estimate,
⇒ expected value
 $y(x', w_{ML})$



Rough work



భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad



Next MAP



భారతీయ సాంకేతిక విజ్ఞాన సంస్థ హైదరాబాద్
भारतीय प्रौद्योगिकी संस्थान हैदराबाद
Indian Institute of Technology Hyderabad

DiL

Data-driven Intelligence
& Learning Lab