

Statistical Problem: A Generic Model

Data Modelling:

Data (Observations) Set: $\mathbb{X}_n := \{X_1, \dots, X_n \mid X_i \in \mathcal{X} \subset \mathbb{R}^k\}$

Observation Space $:= \mathcal{X}$

Assume Data IID: $X_1, \dots, X_n \sim_{iid}$ copies of "generic" X

Distribution of X : $X \sim P \quad P \in \mathcal{F}$

Family of Distribution $:= \mathcal{F}$

Distribution of \mathbb{X}_n : $\mathbb{X}_n \sim \mathbb{P} := \underbrace{P \times \dots \times P}_{n - \text{times}}$

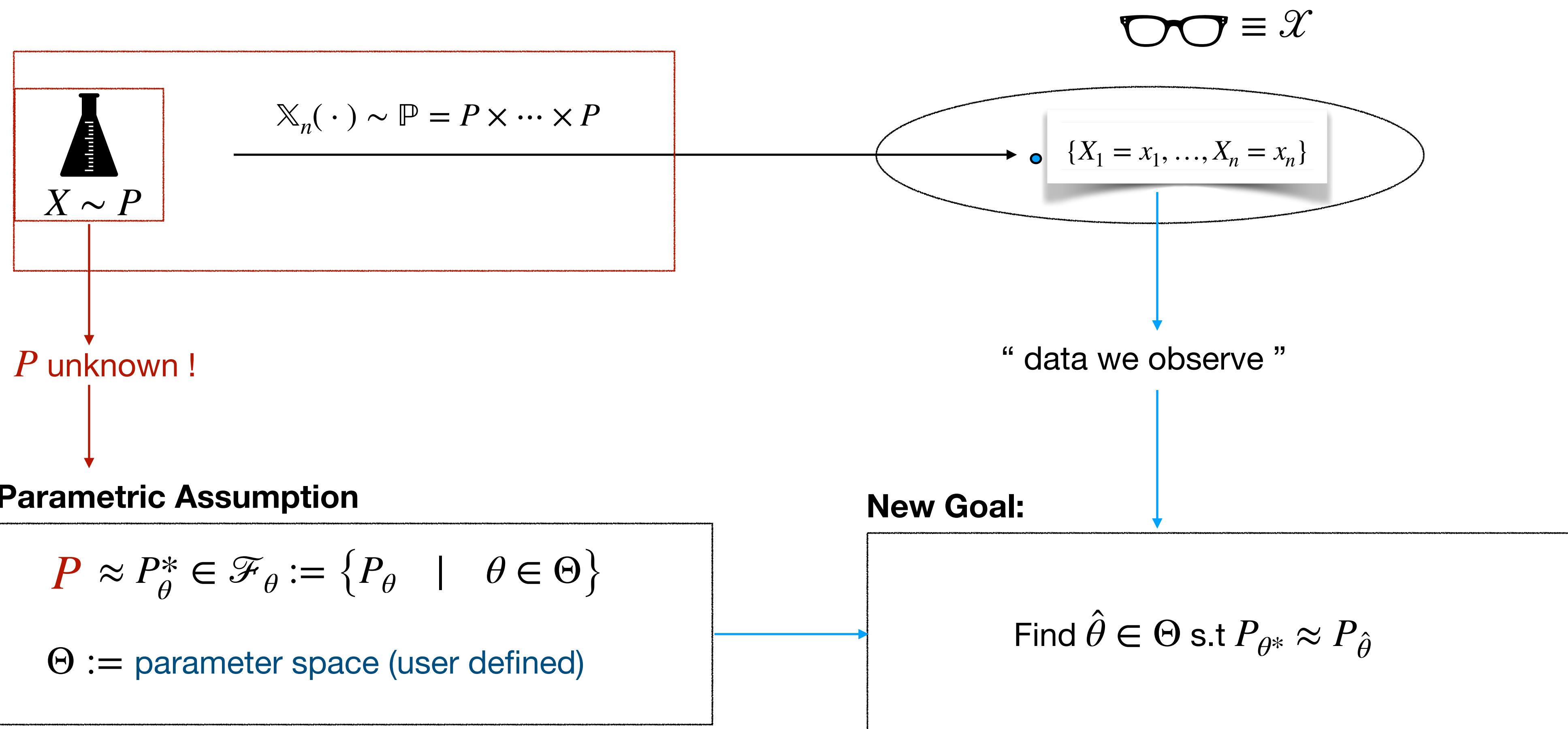
Remark

Knowing $P \equiv$ Knowing the problem

Goal:

FIND P

Statistical Problem: A Parametric Model



Statistical Problem: A Parametric Model : Tossing Example

Given a set of n outcomes from tossing a coin n time, is the coin fair ?

$$\mathbb{D}_{coin} := \{\text{head}, \dots, \text{tail}\}$$

“Tossing a coin ”

$$=: X \in \mathcal{X} = \{\text{head, tail}\} \equiv \{0,1\}$$

Parametric Assumption:

$$X \sim Ber(\theta) =: \mathcal{F}_\theta \quad \theta \in \Theta := [0,1]$$

Family of densities:

$$\mathcal{P}_\theta = \{p_\theta(X = x) = \theta^x(1 - \theta)^{1-x} \quad , \quad x \in \{0 = \text{head}, 1 = \text{tail}\}\}$$

Data Set as R.V

$$\mathbb{X}_n = \{X_1, \dots, X_n\} \quad \underbrace{X_i}_{iid \sim X} := i^{\text{th}} \text{ time the coin is tossed}$$

Distribution of dataset
(Densities Family)

$$\mathcal{D}_\theta = \{\mathbb{P}_\theta = \prod_{i \in [n]} P_\theta\} \implies p_\theta(\mathbb{X}_n = \mathbb{D}_{coin}) = \prod_{i \in [n]} \theta^{x_i}(1 - \theta)^{1-x_i}$$

$$\mathbb{X}_n \quad \equiv \quad \mathbb{D}_{coin}$$

realisation

Estimate θ from \mathbb{D}_n

$\theta \approx? 0.5$

Statistical Problem: Likelihood Function

ASSUMPTIONS Slide 1 Data Modelling:

Words:

Data Set IID realised

Distribution Parametrised

Parameter Space Θ user defined

Density well defined

Maths:

$$\mathbb{X}_n = \mathbb{D}_n = \{X_1 = x_1, \dots, X_n = x_n\}$$

$$\mathbb{X}_n \sim \mathbb{P}_\theta := P_\theta \times \dots \times P_\theta \quad \theta \in \Theta$$

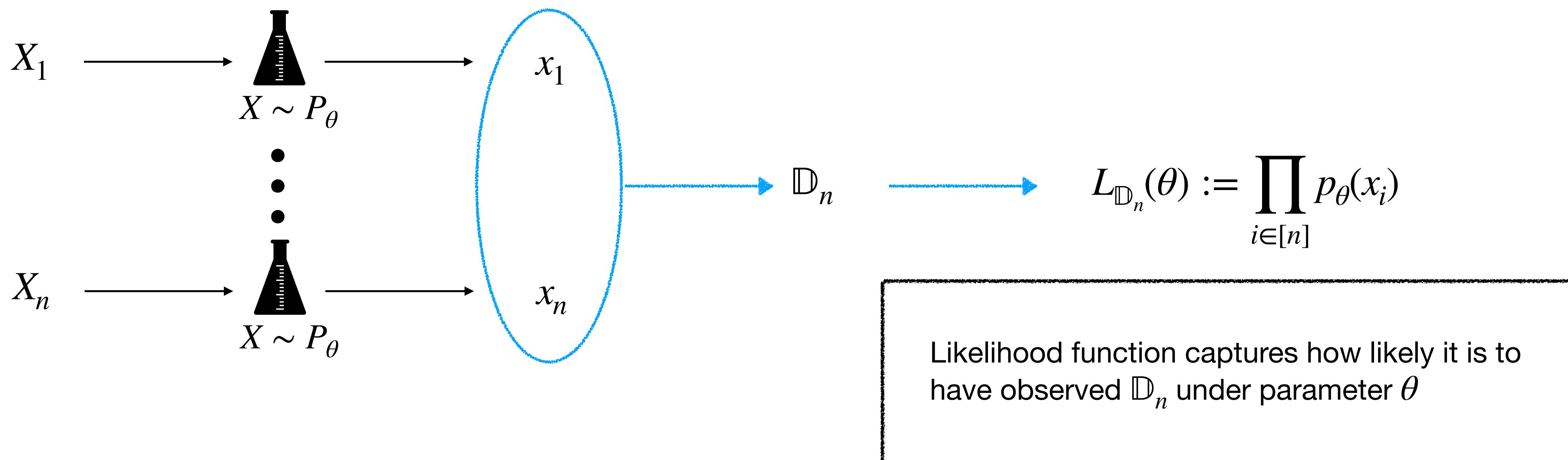
$$p_\theta(\mathbb{D}_n) = \prod_{i \in [n]} p_\theta(x_i)$$

Definition: Likelihood Function of data set (realised) \mathbb{D}_n :

$$L_{\mathbb{D}_n} : \Theta \rightarrow \mathbb{R}$$

$$L_{\mathbb{D}_n}(\theta) := \prod_{i \in [n]} p_\theta(x_i)$$

Statistical Problem: Likelihood Function interpretation

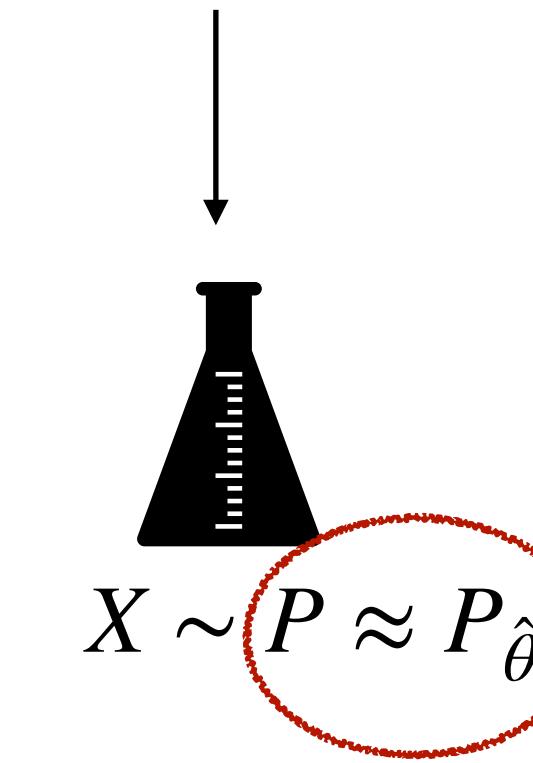
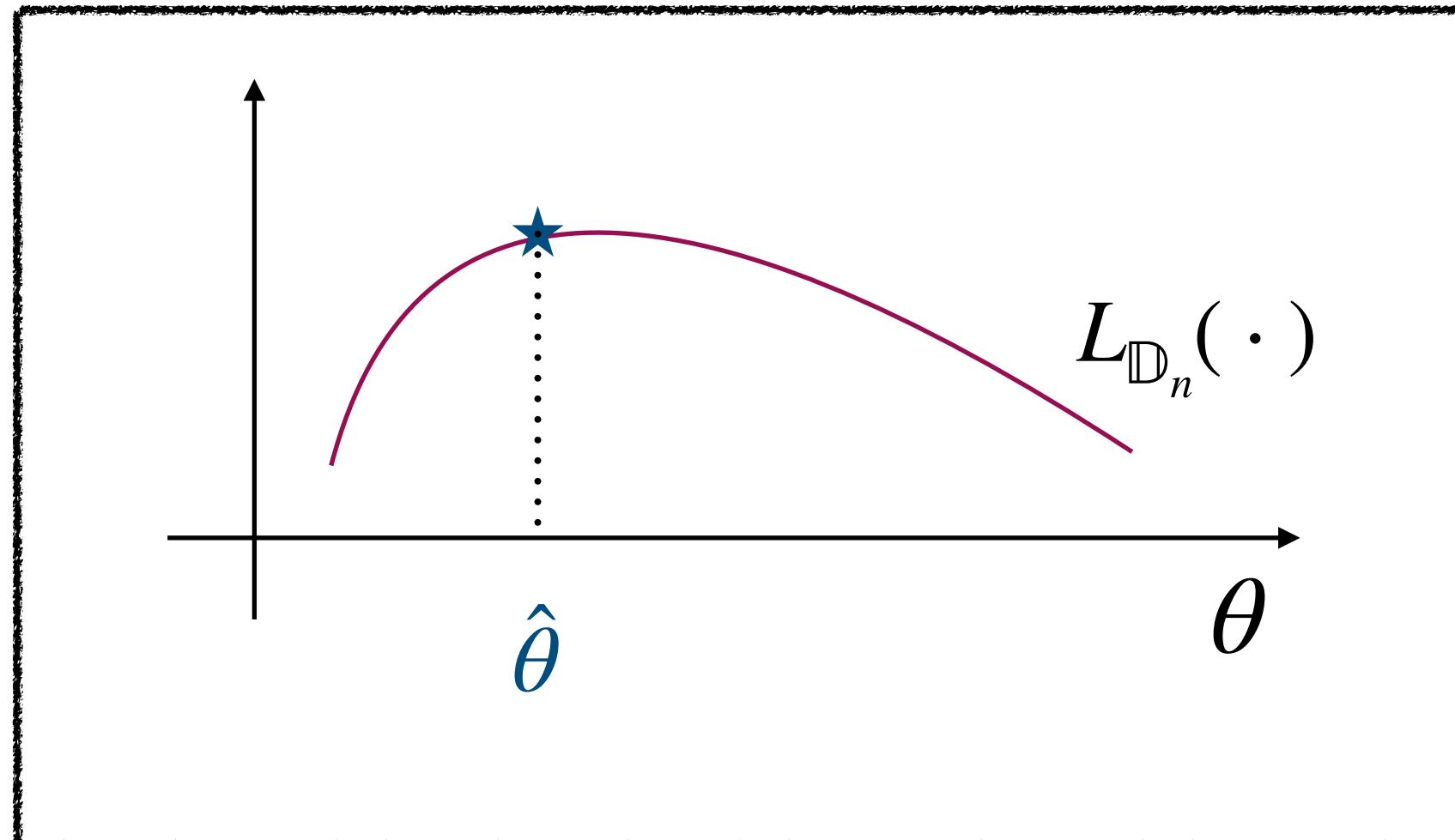


Statistical Problem: Maximum Likelihood Estimator

$$\hat{\theta} = \arg \max_{\theta \in \Theta} L_{\mathbb{D}_n}(\theta)$$

Only if $\text{Supp}(L_{\mathbb{D}_n})$ independent of θ

$$\equiv \hat{\theta} = \arg \max_{\theta \in \Theta} \log(L_{\mathbb{D}_n}(\theta)) \equiv \hat{\theta} = \arg \min_{\theta \in \Theta} -\log(L_{\mathbb{D}_n}(\theta))$$

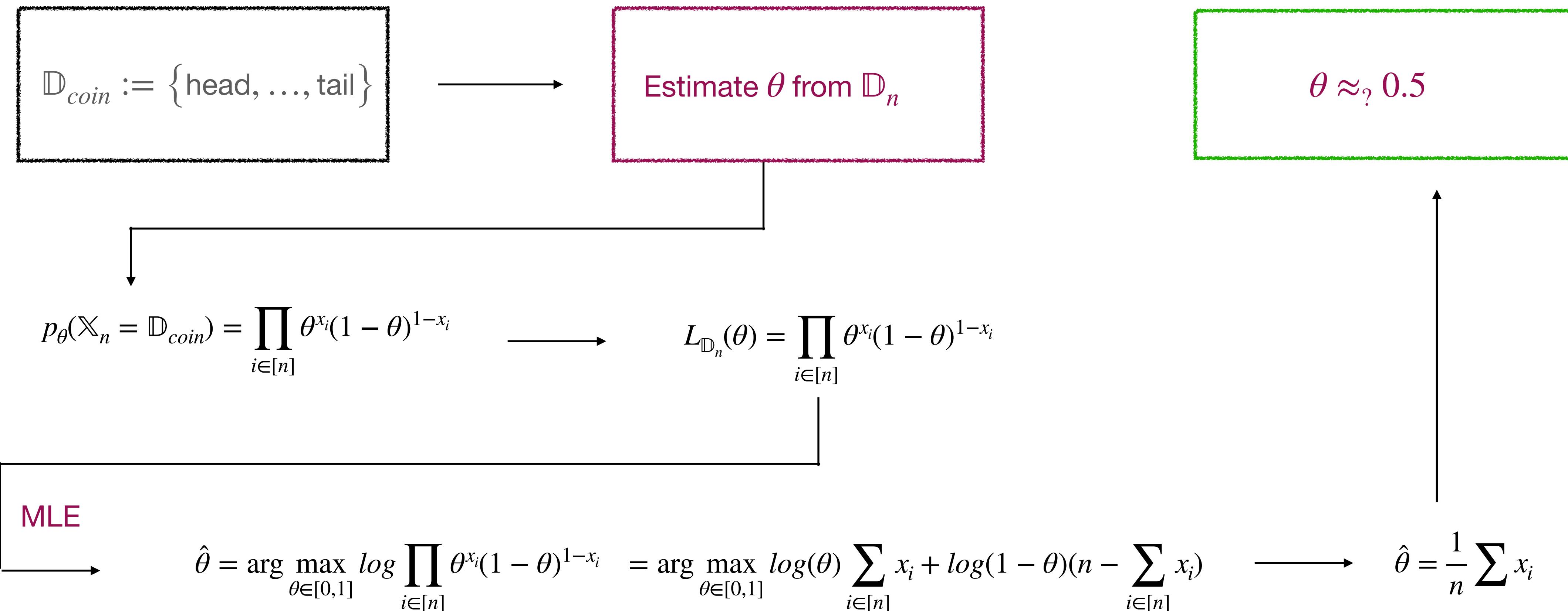


Statistics → we **estimated** the underlying process

≡

Machine Learning → we **learned** the underlying process

Statistical Problem: A Parametric Model : Tossing Example (continued)



Bayesian VS Frequentist: A Different Philosophy of Probabilistic Modelling

Frequentist Postulates

VS

Bayesian Postulates

All of Statistics: Chapter 11

<https://egrcc.github.io/docs/math/all-of-statistics.pdf>

Bayesian Method: Generic Mechanism

Probability over parameter space Θ

1 $\theta \sim p(\theta) =: \text{prior distribution}$

Choice of statistical model $p(x | \theta)$
(belief about data x given parameter θ)

2 $X \sim p(x | \theta)$

Given data set of n iid realisations
update our belief

3 $p(\theta | \mathbb{D}_n) =: \text{posterior distr.}$

Data Set :

$$\mathbb{D}_n = \{X_i\}_{i \in [n]} \quad X_i \sim_{iid} p(X | \theta)$$

$$p(\mathbb{D}_n | \theta) = \prod_{i \in [n]} p_\theta(x_i) = L_{\mathbb{D}_n}(\theta)$$

Likelihood distr.

Baye's Theorem

$$p(\theta | x) = \frac{p(x | \theta)p(\theta)}{\int_{\Theta} p(x | \theta)p(\theta)}$$

Plug-in

$$p(\theta | \mathbb{D}_n) = \frac{L_{\mathbb{D}_n}(\theta)p(\theta)}{\int_{\Theta} L_{\mathbb{D}_n}(\theta)p(\theta)}$$

Bayesian Method: Point Estimate : MAP

$$p(\theta | \mathbb{D}_n) = \frac{L_{\mathbb{D}_n}(\theta)p(\theta)}{\int_{\Theta} L_{\mathbb{D}_n}(\theta)p(\theta)}$$

$$p(\theta | \mathbb{D}_n) \propto L_{\mathbb{D}_n}(\theta)p(\theta)$$

Posterior \propto Likelihood \times Prior

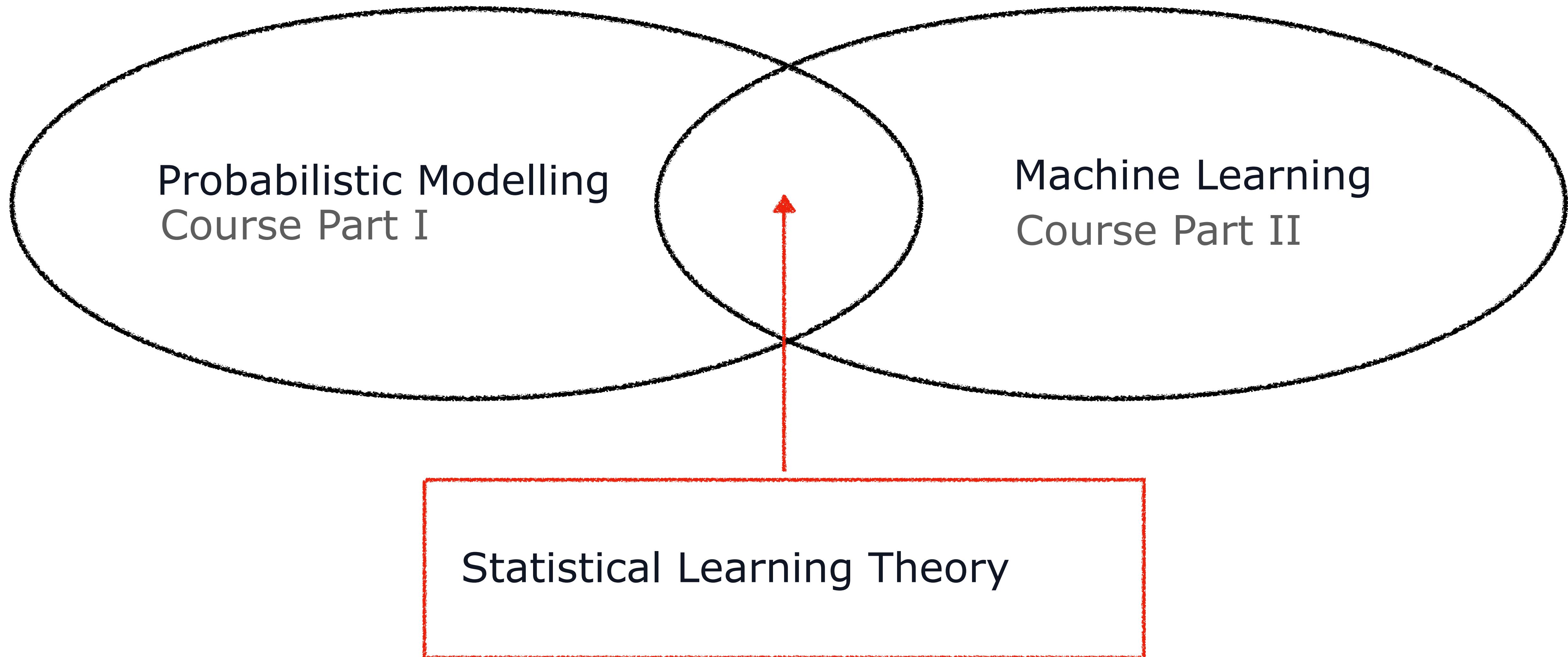
Maximum A Posteriori Estimate

$$\hat{\theta}_{MAP} := \arg \max_{\theta \in \Theta} p(\theta | \mathbb{D}_n)$$

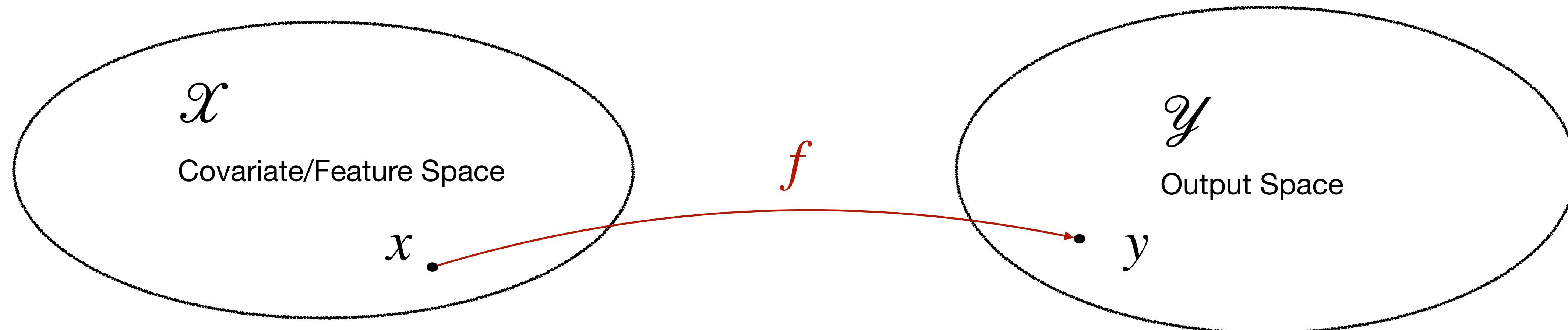
$$= \arg \max_{\theta \in \Theta} \log(L_{\mathbb{D}_n}(\theta)) + \log(p(\theta)) \text{ (under some regularity assumptions)}$$

From Stochastic to ML: Fundamentals of Supervised Learning Model

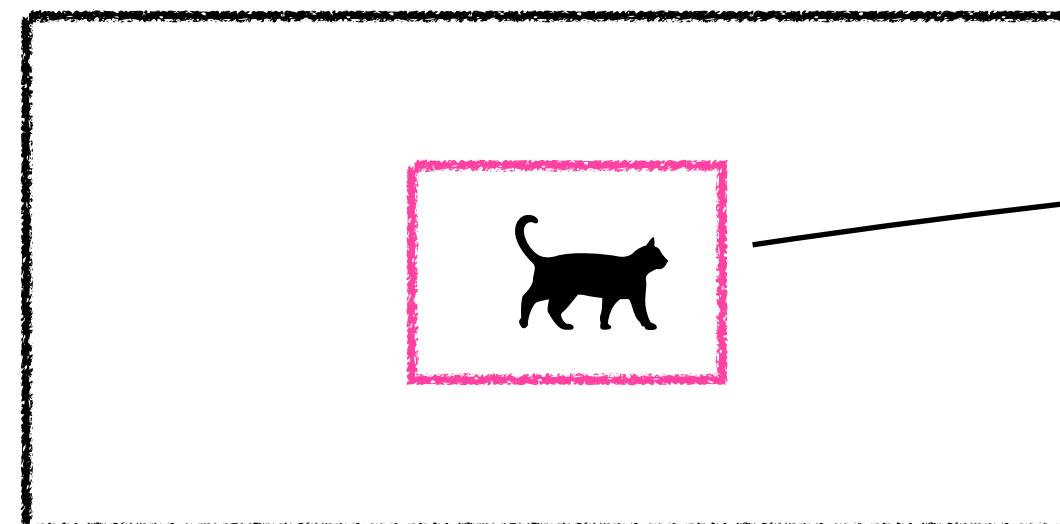
Overview



From Stochastic to ML: Fundamentals of Supervised Learning Model



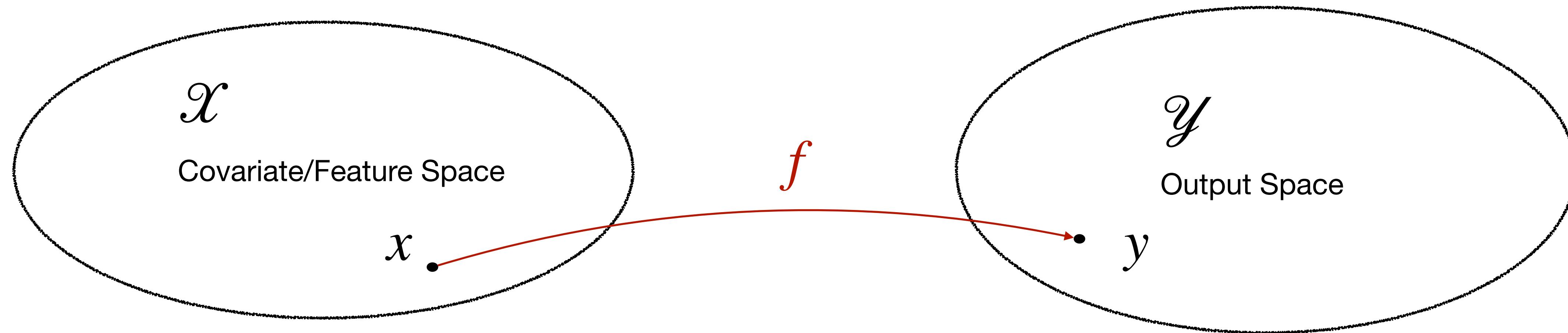
Images Space: $\{Dog, Cat\} \subset \mathbb{R}^3$



Output Space

$\{Dog; Cat\}$

From Stochastic to ML: Fundamentals of Supervised Learning Model



“Supervised” \implies learn/ estimate/ approximate f from $\mathbb{D}_n = \{x_i, y_i\}_{i \in [n]}$

Additive Model:

$$y = f(x) + \epsilon$$

$\epsilon \sim P_\epsilon$ Random Perturbation
 $f(x)$ Systematic Component

From Stochastic to ML: Fundamentals of Supervised Learning Model

“Supervised” \implies learn/ estimate/ approximate from $\mathbb{D}_n = \{x_i, y_i\}_{i \in [n]}$

Additive Model (zero-mean Gaussian Assumption):

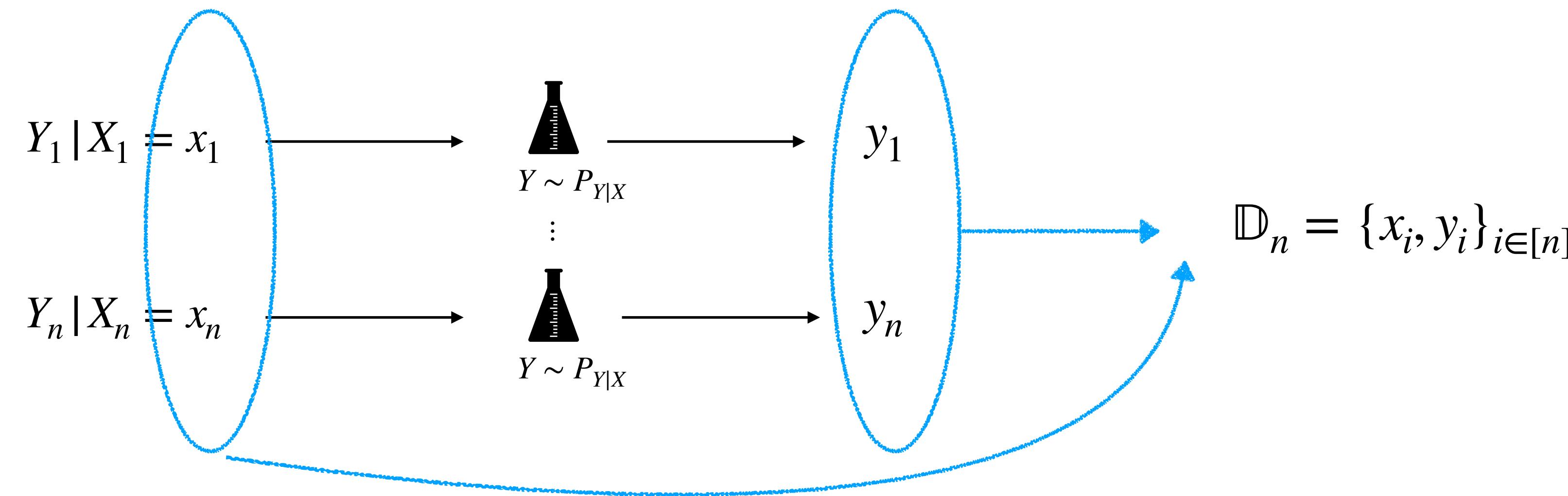
$$y = f(x) + \epsilon \quad \epsilon \sim N(0, \sigma^2) \quad \rightarrow \quad y \sim N(f(x), \sigma^2)$$

Data Modelling (recall stochastic part!)

$$y \sim N(f(x), \sigma^2) \quad \rightarrow \quad \text{knowing } P_{Y|X} \implies \text{knowing } f(x)$$

From Stochastic to ML: Fundamentals of Supervised Learning Model

Data Modelling (recall stochastic part !)



Note: the iid assumption is **key** to encode the fact that we assume similarity in the data generating process i.e the systematic component

From Stochastic to ML: Fundamentals of Supervised Learning Model

“Supervised” \implies learn/ estimate/ approximate from $\mathbb{D}_n = \{x_i, y_i\}_{i \in [n]}$

Additive Model (zero-mean Gaussian Assumption):

$$y = f(x) + \epsilon \quad \epsilon \sim N(0, \sigma^2) \quad \rightarrow \quad y \sim N(f(x), \sigma^2)$$

Model Function Class: Recall $P_{Y|X}$ unknown $\rightarrow \approx P_{Y|X;\theta}$

$$f \in \mathcal{F}_{model}$$

In this course:
parametrised models

$$\mathcal{F}_{model} = \mathcal{F}_\theta = \{f(x; \theta) \mid \theta \in \Theta\}$$

Model Function Class: Examples

linear models

$$\mathcal{F}_\theta = \{\langle x, \theta \rangle \mid \theta \in \Theta\}$$

Neural Network

$$\mathcal{F}_\theta = \{NN(x; \theta) \mid \theta \in \Theta = \{W_i, b_i\}_{i \in [L]}\} \longrightarrow W_i, b_i \text{ Weights of } NN := \text{Neural network}$$

From Stochastic to ML: Fundamentals of Supervised Learning Model

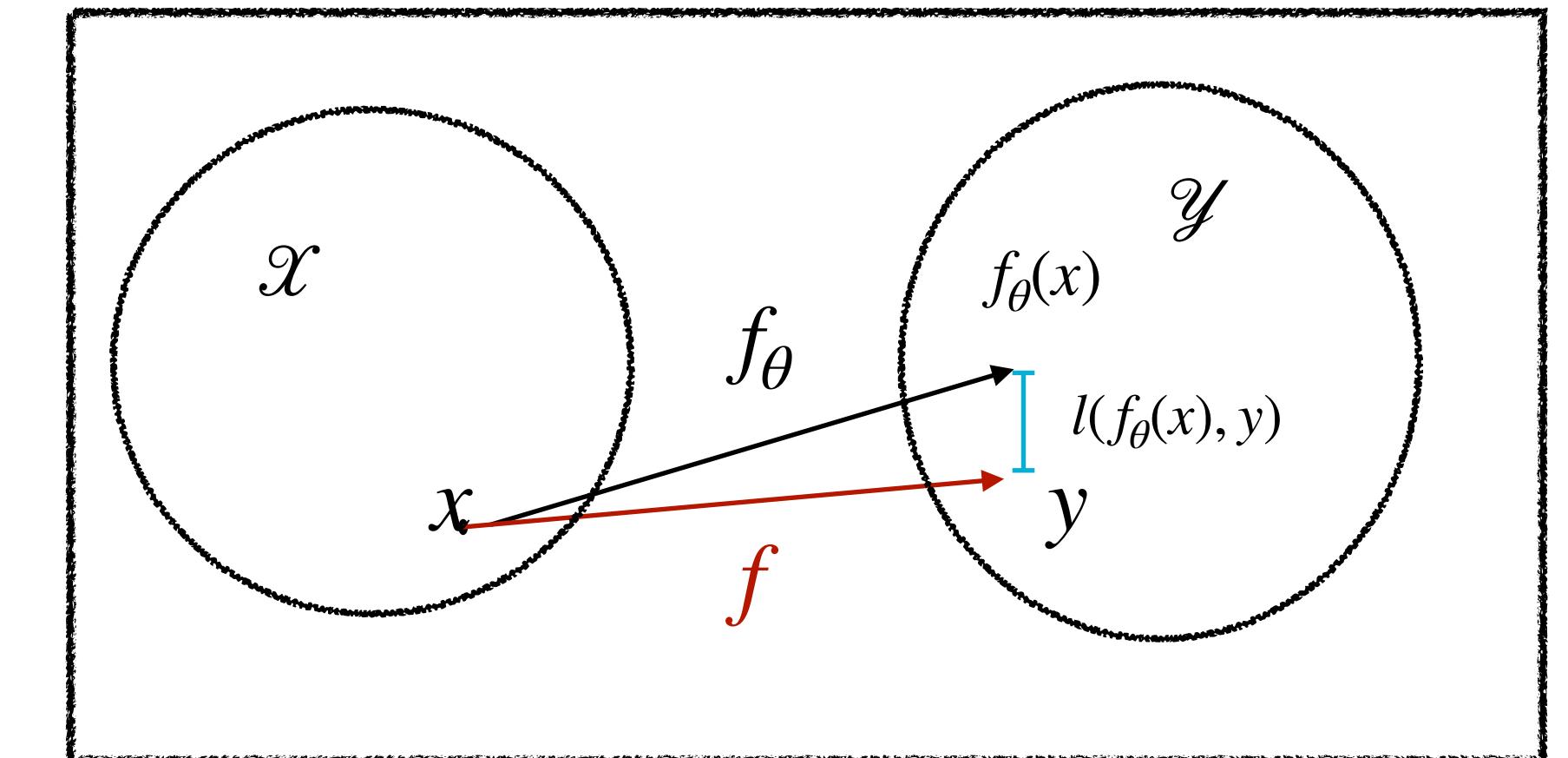
How to quantify the quality of a model function class given $\mathbb{D}_n = \{x_i, y_i\}_{i \in [n]}$?

Loss function

$$l : \Theta \times (\mathcal{X} \times \mathcal{Y}) \rightarrow \mathbb{R}$$
$$l(f_\theta(X), Y)$$

e.g square loss

$$l(f_\theta(X), Y) = \|f_\theta(X) - Y\|^2$$



Quality of loss on average over the data space

Loss function is a random object !

Risk function (population)

$$R(\theta) := \mathbb{E}_{X,Y} l(f_\theta(X), Y)$$

$$P_{X,Y} \text{ Unknown}$$

LLN

$$\mathbb{D}_n = \{x_i, y_i\}_{i \in [n]} \xrightarrow{\text{LLN}}$$

Empirical Risk function

$$\hat{R}(\theta) := \frac{1}{n} \sum_{i \in [n]} l(f_\theta(x_i), y_i)$$

Computable given the DATA !!!

From Stochastic to ML: Fundamentals of Supervised Learning Model

How to quantify the quality of a model function class given $\mathbb{D}_n = \{x_i, y_i\}_{i \in [n]}$?

Risk function (population)

$$R(\theta) := \mathbb{E}_{X,Y} l(f_\theta(X), Y)$$



$$P_{X,Y} \text{ Unknown}$$

LLN

$$\mathbb{D}_n = \{x_i, y_i\}_{i \in [n]}$$

Empirical Risk function

$$\hat{R}(\theta) := \frac{1}{n} \sum_{i \in [n]} l(f_\theta(x_i), y_i)$$

Computable given the DATA !!!

Empirical Risk Minimisation

$$\hat{\theta} := \arg \min_{\theta \in \Theta} \frac{1}{n} \sum_{i \in [n]} l(f_\theta(x_i), y_i)$$

Prediction: given a new data point X_{new} predict Y_{new}

$$\hat{y}_{new} = f_{\hat{\theta}}(x_{new})$$

From Stochastic to ML: Fundamentals of Supervised Learning Model

How to quantify the quality of a model function class given $\mathbb{D}_n = \{x_i, y_i\}_{i \in [n]}$?

Risk function (population)

$$R(\theta) := \mathbb{E}_{X,Y} l(f_\theta(X), Y)$$



$$P_{X,Y} \text{ Unknown}$$

LLN

$$\mathbb{D}_n = \{x_i, y_i\}_{i \in [n]}$$

Empirical Risk function

$$\hat{R}(\theta) := \frac{1}{n} \sum_{i \in [n]} l(f_\theta(x_i), y_i)$$

Empirical Risk Minimisation

$$\hat{\theta} := \arg \min_{\theta \in \Theta} \frac{1}{n} \sum_{i \in [n]} l(f_\theta(x_i), y_i)$$

Computable given the DATA !!!

$$n \rightarrow \infty$$
$$R(\theta) \approx \hat{R}(\theta)$$

For some class of function
 \mathcal{F}_θ

Notice: the more data we have, the closer we should be from the true quantitative analysis from the model we defined

Statistical Problem: Example: Linear Regression (fundamental ml modelling)

Data Modelling Additive Assumption:

$$y = f(x) + \epsilon \quad \epsilon \sim N(0, \sigma^2) \quad x \in \mathbb{R}^d, \quad y \in \mathbb{R}$$

Model Function Class: Recall $P_{Y|X}$ unknown $\rightarrow \approx P_{Y|X;\theta}$

$$\mathcal{F}_\theta = \{x^T \theta \mid \theta \in \mathbb{R}^d\}$$

Loss function : quadratic loss

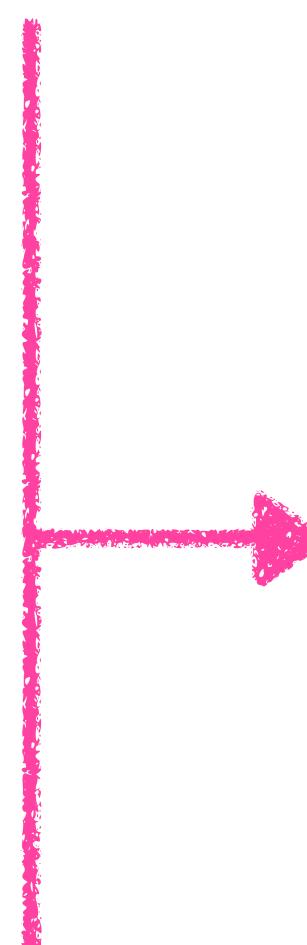
$$l(\langle \theta, x \rangle, y) = \|x^T \theta - y\|^2$$

Find best θ given data set \mathbb{D}_n

$$\hat{\theta} = \arg \min \frac{1}{n} \sum \|x_i^T \theta - y_i\|^2 = \|\tilde{X}\theta - \tilde{y}\|^2$$

With

$$\begin{aligned} \tilde{X} &= [x_1^T, \dots, x_n^T]^T \\ \tilde{y} &= [y_1, \dots, y_n]^T \end{aligned}$$



Data Modelling Parametric Assumption

$$y = \langle x, \theta \rangle + \epsilon \quad \epsilon \sim N(0, \sigma^2) \implies y|x \sim N(x^T \theta, \sigma^2)$$

Why “Regression” ?

$\mathbb{E}\{Y|X=x\} =:$ Regression function

$$\mathbb{E}\{Y|X=x\} = x^T \theta \xrightarrow{\hat{\theta}} \mathbb{E}\{Y|X=x\} \approx x^T \hat{\theta}$$

We learned/approximate the regression function !

Statistical Problem: Linear Regression: Least Square Estimation

$$\hat{\theta} = \arg \min_{\theta \in \mathbb{R}^d} \|\tilde{X}\theta - \tilde{y}\|^2$$

With

$$\begin{aligned}\tilde{X} &= [x_1^T, \dots, x_n^T]^T \\ \tilde{y} &= [y_1, \dots, y_n]^T\end{aligned}$$

↓
LS solution (linear algebra)

$$\hat{\theta} = (\tilde{X}^T \tilde{X})^+ \tilde{X}^T \tilde{y}$$

↓
Prediction

$$\hat{y}_{new} = x_{new}^T \hat{\theta}$$

Recall: additive model

$$y = f_\theta(x) + \epsilon$$



$$\hat{\theta} = \arg \min_{\theta \in \mathbb{R}^d} \|\epsilon\|^2$$



Minimum error in l2 norm with
respect to our model
assumption !

Statistical Problem: Linear Regression: MLE

Additive Model (zero -mean Gaussian Assumption):

$$y = \langle x, \theta \rangle + \epsilon \quad \epsilon \sim N(0, \sigma^2) \implies y | x \sim N(x^T \theta, \sigma^2)$$

Unknown parameter :

$$y | x \sim N(\underline{x^T \theta}, \underline{\sigma^2})$$

Estimating the Mean and Variance MLE

$$\{\hat{\theta}_{MLE}, \sigma_{MLE}\} = \arg \min_{\theta \in \mathbb{R}^d, \sigma \in \mathbb{R}} -\log(L_{\mathbb{D}_n}(\theta; \sigma))$$

$$\longrightarrow L_{\mathbb{D}_n}(\theta) = (2\pi\sigma^2)^{-n/2} \exp\left(-\frac{1}{2\sigma^2} \sum_{i \in [n]} (y_i - x_i^T \theta)^2\right)$$

$$\longrightarrow \log(L_{\mathbb{D}_n}(\theta)) \propto -\frac{n}{2} \ln(\sigma^2) - \frac{1}{2\sigma^2} \|\tilde{y} - \tilde{X}\theta\|^2$$

$$\hat{\theta}_{MLE} = \hat{\theta}_{LS} = (\tilde{X}^T \tilde{X})^+ \tilde{X}^T \tilde{y}$$

$$\sigma_{MLE}^2 = \frac{1}{n} \sum_{i \in [n]} (y_i - x_i^T \hat{\theta}_{MLE})^2$$

$$\hat{y}_{new} = x_{new}^T \hat{\theta}_{MLE}$$

Confidence Interval

Prediction Interval

Statistical Problem: Bayesian Linear Regression: MAP

Additive Model (zero -mean Gaussian Assumption):

Recall: Prior user defined !

$$y = \langle x, \theta \rangle + \epsilon \quad \epsilon \sim N(0, \sigma^2), \quad \theta \sim p(\theta) \implies y | x, \theta \sim N(x^T \theta, \sigma^2)$$

$$p(\theta)$$

Choose a Gaussian Prior

$$p(\theta) = N(0, \sigma_p^2)$$

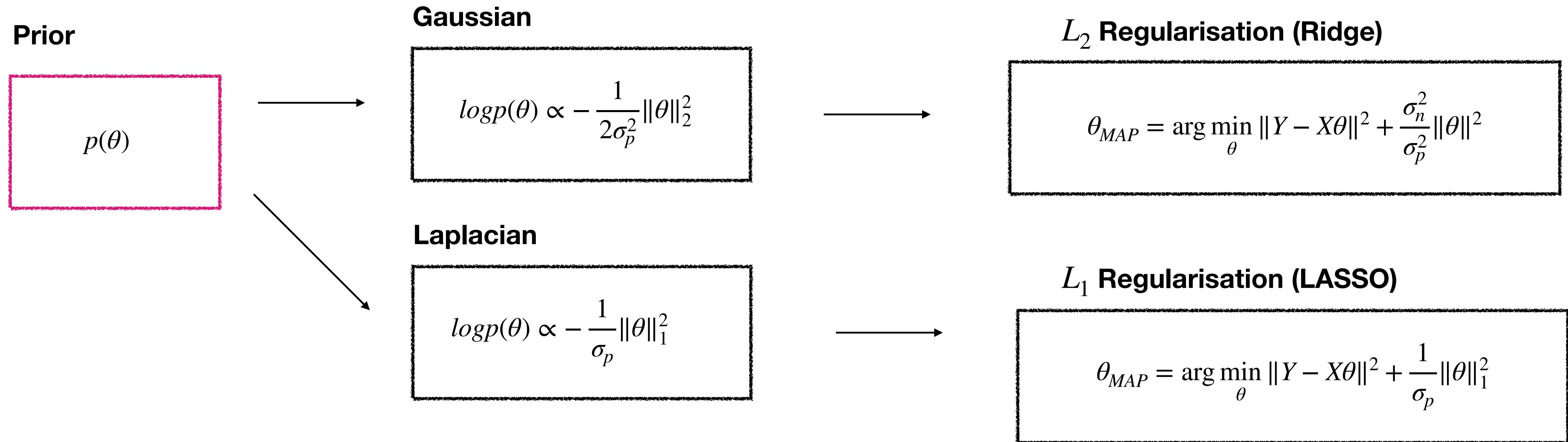
MAP Estimate

$$\begin{aligned}\theta_{MAP} &= \arg \max_{\theta} p(\theta | \mathbb{D}_n) = \log(L_{\mathbb{D}_n}) + \log(p(\theta)) \\ &= \arg \min_{\theta} \frac{1}{\sigma_n^2} \sum_{i \in [n]} (y_i - \theta^T x_i)^2 + \frac{1}{\sigma_p^2} \|\theta\|^2 \\ &= \arg \min_{\theta} \|Y - X\theta\|^2 + \frac{\sigma_n^2}{\sigma_p^2} \|\theta\|^2\end{aligned}$$

$$\theta_{MAP} = (X^T X + \frac{\sigma_n^2}{\sigma_p^2} \mathbb{I})^{-1} X^T y$$

≡ Ridge Regression (or L_2 constraint least squares)

Statistical Problem: Bayesian Linear Regression: MAP Prior Choice



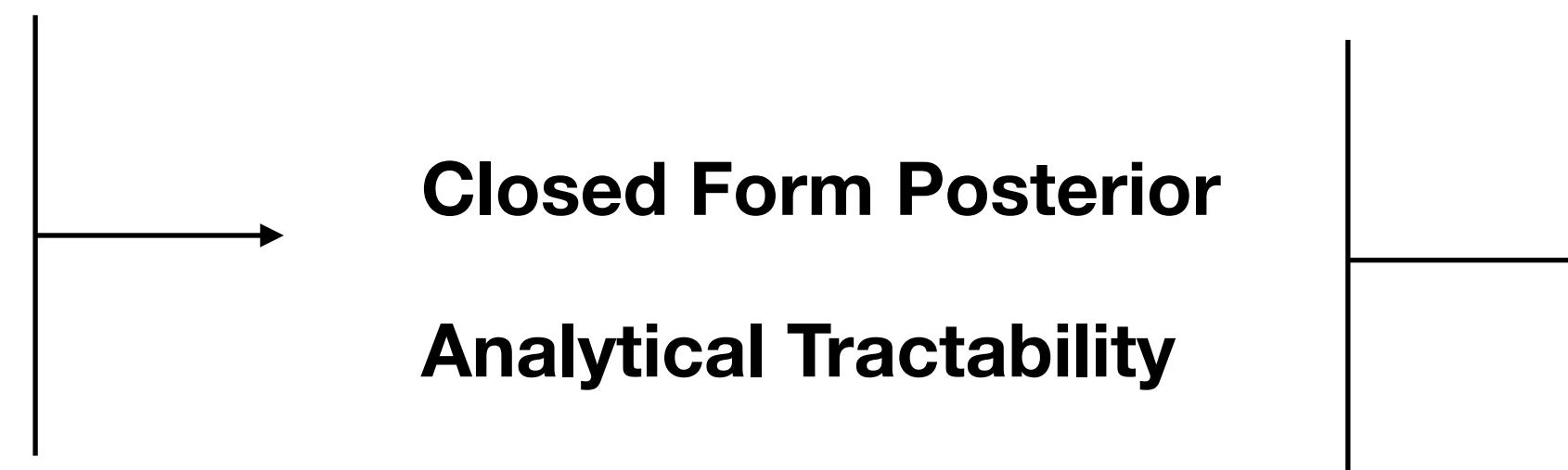
Statistical Problem: Bayesian Modelling : Conjugate Prior

Definition :

When the **prior** and the **posterior** are from the same family of distribution we say that the prior is **conjugate** with respect to the model

Special Note :

Data Gaussian + Prior Gaussian \implies Gaussian Posterior



Useful for :

**Incremental Learning
Predictive distribution**

Statistical Problem: Bayesian Modelling : Conjugate Prior for Gaussian Data

Recall Linear Gaussian Noise Additive Model :

$$L_{\mathbb{D}_n}(\theta) = \prod_{i \in [n]} p(y_i | x_i, \theta) \quad p(y_i | x_i, \theta) = N(\theta^T x_i, \sigma_n)$$

Zero-Mean Gaussian Prior:

$$p(\theta) = N(0, \sigma_p)$$

Is Posterior Gaussian ? i.e $p(\theta)$ conjugate prior ?

$$p(\theta | \mathbb{D}_n) = L_{\mathbb{D}_n}(\theta)p(\theta) = \prod N_y(\theta^T x_i, \sigma_n)N_\theta(0, \sigma_p) \Leftrightarrow \text{Qty Gaussian in } \theta$$

Next slide: Simplified proof for univariate gaussian !

Product of Univariate Gaussian

Note: For Multivariate, Results on Gaussian Conditioned on Gaussian are used to for the proof

$$p(\theta | \mathbb{D}_n) = L_{\mathbb{D}_n}(\theta)p(\theta) = \prod N_y(\theta^T x_i, \sigma_n) N_\theta(0, \sigma_p) = \frac{1}{(2\pi\sigma_n)^{n/2}(2\pi\sigma_p)^{1/2}} \exp \left\{ (-1)(\frac{1}{2\sigma_n^2} \sum_{i \in [n]} (y_i - x_i \theta)^2 + \frac{1}{2\sigma_p^2} \theta^2) \right\}$$

$$:= \beta$$

$$\beta = \frac{1}{2\sigma_n^2} \sum_{i \in [n]} (y_i - x_i \theta)^2 + \frac{1}{2\sigma_p^2} \theta^2 = \frac{1}{2\sigma_n^2} (\sum_i y_i^2 - 2\theta \sum_i x_i y_i + \theta^2 \sum_i x_i^2) + \frac{1}{2\sigma_p^2} \theta^2$$

$$\beta = \frac{\theta^2(\sigma_n^2 + \sigma_p^2 x_n^2) - 2\theta y x_n \sigma_p^2 + \sigma_p^2 y_n^2}{2\sigma_p^2 \sigma_n^2} = \frac{\theta^2 - 2\theta(\frac{y x_n \sigma_p^2}{\sigma_n^2 + \sigma_p^2 x_n^2}) + \frac{\sigma_p^2 y_n^2}{\sigma_n^2 + \sigma_p^2 x_n^2}}{\frac{2\sigma_p^2 \sigma_n^2}{\sigma_n^2 + \sigma_p^2 x_n^2}}$$

$$\sum_i y_i^2 := y_n^2$$

$$\sum_i y_i x_i := y x_n$$

$$\sum_i x_i^2 := x_n^2$$

$$=: \frac{\theta^2 - 2A\theta + B}{C}$$

Product of Univariate Gaussian

Define:

$$\epsilon = \frac{A^2 - A^2}{C} = 0$$

Then :

$$\beta = \beta + \epsilon = \frac{(\theta - A)^2}{C} + \frac{A^2 - B}{C}$$

Which implies :

$$\exp\left\{-\beta\right\} = \exp\left\{-\frac{(\theta - A)^2}{C}\right\} \exp\left\{\frac{A^2 - B}{C}\right\}$$

Scaling Factor

Recall what we aim :

$$p(\theta | \mathbb{D}_n) \propto \exp\left\{-\frac{(\theta - \mu_\theta)^2}{2\sigma_\theta^2}\right\}$$

$$\mu_\theta = A = \frac{yx_n\sigma_p^2}{\sigma_n^2 + \sigma_p^2 x_n^2}$$

$$\sigma_\theta^2 = \frac{1}{2}C = \frac{\sigma_p^2 \sigma_n^2}{\sigma_n^2 + \sigma_p^2 x_n^2}$$

NOTE: What did we show ?

Proportionality !

For valid density: need to workout the scaling factor !

Bayes Classifier : Framework

Define:

Features : $X \in \mathbb{R}^d$

Output : $Y \in \{1, \dots, C\} \subset \mathbb{N}$

Classifier :

$$h : X \rightarrow Y$$

Definition: Baye's Classifier

$$C_{Bayes}(x) = \arg \max_{c \text{ in } [C]} P(Y = c | X = x)$$

Baye's Classifier Optimal wrt Setting !

Intuition: (not a Proof !) Recall Risk Framework

Error of Misclassification :

$$\begin{aligned} P(Y \neq h(X)) &= \mathbb{E} \mathbb{I}_{Y \neq h(X)} = \mathbb{E}_X \mathbb{E}_{Y|X} \mathbb{I}_{Y \neq h(X)|X=x} \\ &= \mathbb{E}_X P(Y = 1 | x) \mathbb{I}_{h(x)=2, \dots, C} + P(Y = 2 | x) \mathbb{I}_{h(x)=1, 3, \dots, C} \dots \end{aligned}$$

Minimised by maximising the chance of not making an error !

Bayes Classifier : Naive Baye's Classifier

Definition: Baye's Classifier

$$C_{Bayes}(x) = \arg \max_{c \text{ in } [C]} P(Y = c | X = x)$$

Assume $X = \{x_i\}_{i \in [n]}$ iid then the Baye's Classifier yield the **naive** baye's classifier, where naive stand for the assumption of the independence in the features

$$C_{Bayes}^{Naive}(x) = \arg \max_{c \text{ in } [C]} P(Y = c) \prod_{i \in [n]} P(x_i)$$