



# SF1935 – Probability Theory and Statistics with Application to Machine Learning

## Lecture 2: Linear models and a probabilistic perspective

Pawel Herman

Computational Science and Technology (CST)

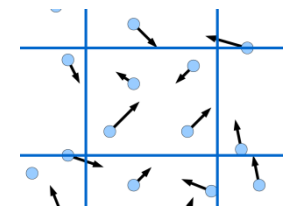
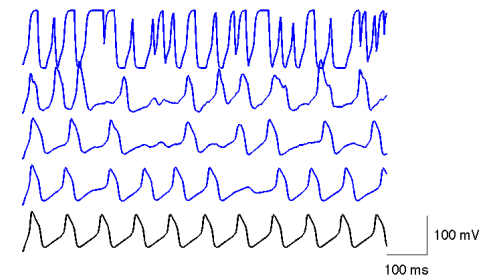
KTH Royal Institute of Technology

- Introduction
- Probabilistic approach
- Probability basics

# Introduction

## Probabilistic approach

- Ubiquitous nature of uncertainty
  - imprecision, noise in data,
  - errors , missing information/data
  - gaps in knowledge, simplified description



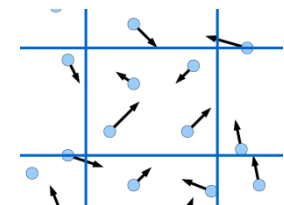
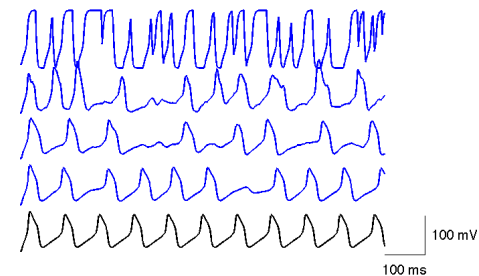
- Introduction
- Probabilistic approach
- Probability basics

# Introduction

## Probabilistic approach

- Ubiquitous nature of uncertainty
  - imprecision, noise in data,
  - errors , missing information/data
  - gaps in knowledge, simplified description

*“The theory of probabilities is at bottom nothing but common sense reduced to calculus; it enables us to appreciate with exactness that which accurate minds feel with a sort of instinct for which of times they are unable to account.” (Laplace)*

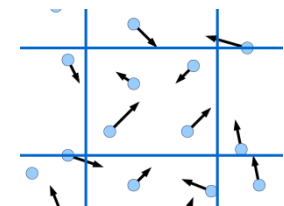
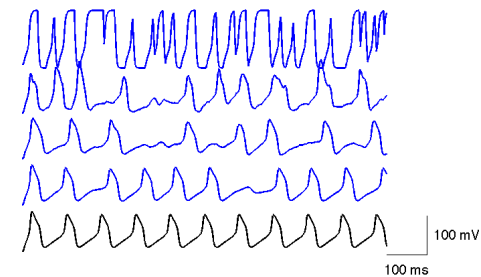


- Introduction
- Probabilistic approach
- Probability basics

# Introduction

## Probabilistic approach

- Ubiquitous nature of uncertainty
  - imprecision, noise in data
  - errors , missing information/data
  - gaps in knowledge, simplified description
- Probability theory provides a framework for modelling, reasoning etc. under uncertainty
  - unified, universal, intuitive, interpretable
  - beyond randomness, it is about uncertainty!
  - p. distributions as carriers or information (Jaynes, 2003)



# Probabilistic perspective in ML

- Statistical ML: constructing stochastic models
  - fully probabilistic description and inference
  - theoretical assumptions, mathematical tractability, rigour
  - parameters estimated from observed data (learning)
  - interpretability and extra insights
  - machinery to propagate and account for uncertainty effects

# Probabilistic perspective in ML

- Statistical ML: constructing stochastic models
  - fully probabilistic description and inference
  - theoretical assumptions, mathematical tractability, rigour
  - parameters estimated from observed data (learning)
  - interpretability and extra insights
  - machinery to propagate and account for uncertainty effects

**BUT:**     *can be very hard for large-scale problems* and  
*difficult to derive solutions in a closed form*

# Probabilistic perspective in ML

- Statistical ML: constructing stochastic models

- fully probabilistic description and inference

- t

- p

- il

- n

*“(statistical ML) provides the basis for learning algorithms that directly manipulate probabilities, as well as a framework for analyzing the operation of other algorithms that do not explicitly manipulate probabilities.”*

Mitchell

rigour

effects

# Probabilistic perspective in ML

- Philosophy of Bayesian approach
  - Uncertainty is ubiquitous – describe all model components with probabilistic objects (*distributions, not point estimates*)



# Probabilistic perspective in ML

- Philosophy of Bayesian approach
  - Uncertainty is ubiquitous – describe all model components with probabilistic objects (*distributions, not point estimates*)
  - Apply Bayesian machinery to propagate uncertainty
    - product and sum probability rules, Bayesian theorem

$$p(\mathbf{w} \mid \mathcal{D}) = \frac{p(\mathcal{D} \mid \mathbf{w}) p(\mathbf{w})}{p(\mathcal{D})}$$

$$\text{posterior} \propto \text{likelihood} \times \text{prior}$$

# Probabilistic perspective in ML

- Philosophy of Bayesian approach
  - Uncertainty is ubiquitous – describe all model components with probabilistic objects (*distributions, not point estimates*)
  - Apply Bayesian machinery to propagate uncertainty
    - product and sum probability rules, Bayesian theorem
    - the power of marginalisation

$$p(x_1, x_2, \dots, x_{n-1}) = \int p(x_1, x_2, \dots, x_n) dx_n$$

$$p(\mathbf{w} \mid \mathcal{D}) = \frac{p(\mathcal{D} \mid \mathbf{w}) p(\mathbf{w})}{p(\mathcal{D})}$$

$$\text{posterior} \propto \text{likelihood} \times \text{prior}$$

# Probabilistic perspective in ML

- Philosophy of Bayesian approach
  - Uncertainty is ubiquitous – describe all model components with probabilistic objects (*distributions, not point estimates*)
  - Apply Bayesian machinery to propagate uncertainty
  - Combine uncertain knowledge with data to reduce uncertainty (based on evidence from observations)
    - batch or sequence where posterior is iteratively updated

# Probabilistic perspective in ML

- Philosophy of Bayesian approach
  - Uncertainty is ubiquitous – describe all model components with probabilistic objects (*distributions, not point estimates*)
  - Apply Bayesian machinery to propagate uncertainty
  - Combine uncertain knowledge with data to reduce uncertainty (based on evidence from observations)
  - Two levels of inference:  
parameter estimation and model selection (see Lecture 3)

- Introduction
- **Probabilistic approach**
- Probability basics

# Back to the probabilistic mindset

## Bayesian vs frequentist perspective

- |  |  |
|--|--|
| I. Probability is a measure of belief (plausibility) | I. The ratio of outcomes in repeated trials. |
|--|--|

- Introduction
- **Probabilistic approach**
- Probability basics

# Back to the probabilistic mindset

## Bayesian vs frequentist perspective

- |  |   |
|--|---|
| I. Probability is a measure of belief (plausibility) | I. The ratio of outcomes in repeated trials.                    |
| II. The data is fixed, models have probabilities     | II. There is a true model and the data is a random realisation. |

- Introduction
- **Probabilistic approach**
- Probability basics

# Back to the probabilistic mindset

## Bayesian vs frequentist perspective

- |   |  |
|---|--|
| I. Probability is a measure of belief (plausibility)                    | I. The ratio of outcomes in repeated trials.                           |
| II. The data is fixed, models have probabilities                        | II. There is a true model and the data is a random realisation.        |
| III. There does not have to be an experiment for declaring probability. | III. Parameters can only be deduced from data (likely outcome of exp.) |

# Back to the probabilistic mindset

## Bayesian vs frequentist perspective

- |   |   |
|---|---|
| I. Probability is a measure of belief (plausibility)                    | I. The ratio of outcomes in repeated trials.                                    |
| II. The data is fixed, models have probabilities                        | II. There is a true model and the data is a random realisation.                 |
| III. There does not have to be an experiment for declaring probability. | III. Parameters can only be deduced from observed data (likely outcome of exp.) |
| IV. Can incorporate prior knowledge, probabilities can be updated.      | IV. Each repeated experiment starts from ignorance.                             |



# Back to the probabilistic mindset

## Bayesian vs frequentist perspective

- |   |   |
|---|---|
| I. Probability is a measure of belief (plausibility)                    | I. The ratio of outcomes in repeated trials.                                    |
| II. The data is fixed, models have probabilities                        | II. There is a true model and the data is a random realisation.                 |
| III. There does not have to be an experiment for declaring probability. | III. Parameters can only be deduced from observed data (likely outcome of exp.) |
| IV. Can incorporate prior knowledge, probabilities can be updated.      | IV. Each repeated experiment starts from ignorance.                             |
| V. Estimators are good for available data.                              | V. Estimators are averaged across many trials.                                  |

- Introduction
- **Probabilistic approach**
- Probability basics

# Back to the probabilistic mindset

## Bayesian vs frequentist perspective

- |  |   |
|--|---|
| I. Probability is a measure of belief (plausibility)                     | I. The ratio of outcomes in repeated trials.                                    |
| II. The data is fixed, models have probabilities                         | II. There is a true model and the data is a random realisation.                 |
| III. There does not have to be an experiment for declaring probability.  | III. Parameters can only be deduced from observed data (likely outcome of exp.) |
| IV. Can incorporate prior knowledge, probabilities can be updated        | IV. Each repeated experiment starts from ignorance.                             |
| V. Estimators are good for available data.                               | V. Estimators are averaged across many trials.                                  |
| VI. Probability of a hypothesis given the data (posterior distribution). | VI. Probability of the data given hypothesis (likelihood, sampling dist.).      |

# Back to the probabilistic mindset

## Bayesian vs frequentist perspective

- |   |  |
|---|--|
| I. Probability is a measure of belief (plausibility)                    | I. The ratio of outcomes in repeated trials.   |
| II. The data is fixed, models have probabilities                        | II. There is a true model and the data is a random realisation.                      |
| III. There does not have to be an experiment for declaring probability. | III. Parameters can only be deduced from observed data (likely outcome of exp.)      |
| IV. Can incorporate prior knowledge, probabilities can be updated       | IV. Each repeated experiment starts from ignorance.                                  |
| V. Estimators are good for available data.                              | V. Estimators are averaged across many trials.                                       |
| VI. Probability of a hypothesis given the data.                         | VI. Probability of the data given hypothesis.  |
| VII. All variables/parameters have distribution.                        | VII. Parameters are fixed unknowns that can be point estimated from repeated trials. |

- Introduction
- **Probabilistic approach**
- Probability basics

# Back to the probabilistic mindset

## Bayesian vs frequentist perspective

- |  |   |
|--|---|
| I. Probability is a measure of belief (plausibility)                 | I. The ratio of outcomes in repeated trials.  |
| II. The data is fixed, models have probabilities                     | II. There is a true model and the data is a random                                      |
| III. There does not<br>declaring proba                               | III. The data is produced from<br>(one of exp.)   |
| IV. Can incorporate prior knowledge,<br>probabilities can be updated | IV. Each repeated experiment starts from<br>ignorance.                                  |
| V. Estimators are good for available data.                           | V. Estimators are averaged across many trials.  |
| VI. Probability of a hypothesis given the data.                      | VI. Probability of the data given hypothesis.   |
| VII. All variables/parameters have distribution.                     | VII. Parameters are fixed unknowns that can be<br>point estimated from repeated trials. |

Why isn't everyone Bayesian? (*Efron, 1986*)

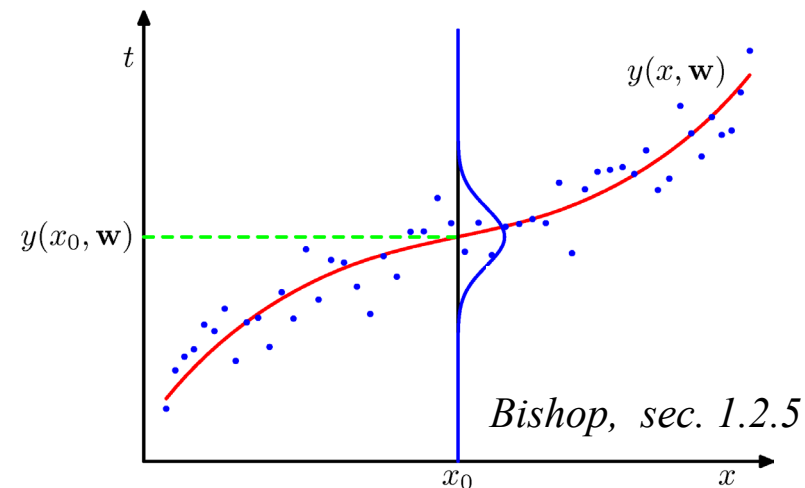
# Learning as inference

- Learning distributions

➤ curve-fitting example:  $y(x, \mathbf{w}): X \rightarrow Y$  (let's assume polynomial)

$$t = y(x, \mathbf{w}) + \varepsilon$$

observations  $\{(x_i, t_i): i=1, \dots, N\}$ :  $\mathbf{x} = (x_1, \dots, x_N)^T$ ,  $\mathbf{t} = (t_1, \dots, t_N)^T$



# Learning as inference

- Learning distributions

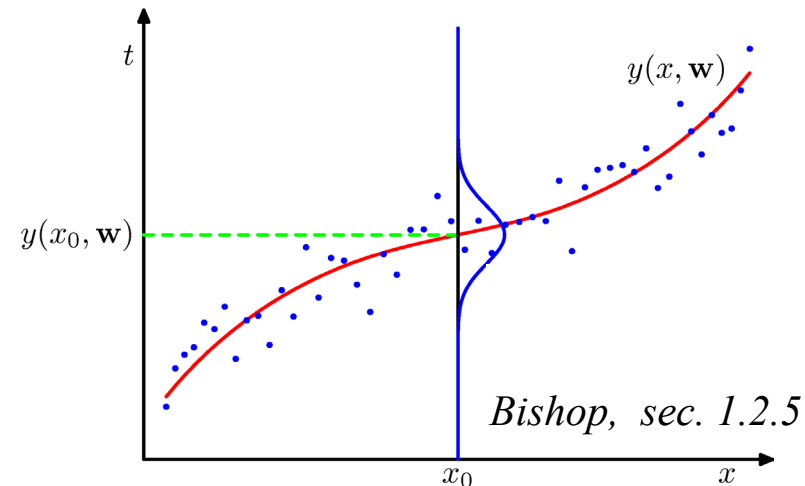
➤ curve-fitting example:  $y(x, \mathbf{w}): X \rightarrow Y$  (let's assume polynomial)

$$t = y(x, \mathbf{w}) + \varepsilon$$

observations  $\{(x_i, t_i): i=1, \dots, N\}$ :  $\mathbf{x} = (x_1, \dots, x_N)^T$ ,  $\mathbf{t} = (t_1, \dots, t_N)^T$

## Remarks about notation:

- 1) here we deal with one-dim input,  $x$  and output,  $t$
- 2) parameters  $\mathbf{w}$  still constitute a vector (e.g. could be polynomial coefficients)
- 3)  $\mathbf{x}$  and  $\mathbf{t}$ , refer to the collection of all inputs and outputs, conceptually corresponding to  $D_x$  and  $D_t$  but in the vector form, so  $\mathbf{x} \rightarrow \mathbf{t}$ .



# Learning as inference

- Learning distributions

- curve-fitting example:  $y(x, \mathbf{w}): X \rightarrow Y$  (let's assume polynomial)

$$t = y(x, \mathbf{w}) + \varepsilon$$

observations  $\{(x_i, t_i): i=1, \dots, N\}$ :  $\mathbf{x} = (x_1, \dots, x_N)^T$ ,  $\mathbf{t} = (t_1, \dots, t_N)^T$

- probabilistic framework

$p(\mathbf{w})$  – a prior probability distribution

$p(\mathcal{D} | \mathbf{w})$  – the likelihood function for  $\mathcal{D} = \mathbf{t}$   
(not a probability distrib. over  $\mathbf{w}$ , just a conditional probability)

$p(\mathbf{w} | \mathcal{D})$  – a posterior probability

# Learning as inference

- Learning distributions

- curve-fitting example:  $y(x, \mathbf{w}): X \rightarrow Y$  (let's assume polynomial)

## Bayes' theorem

observation

$$p(\mathbf{w} | \mathcal{D}) = \frac{p(\mathcal{D} | \mathbf{w}) p(\mathbf{w})}{p(\mathcal{D})}$$

➤ probab

$p(\mathbf{w})$

posterior  $\propto$  likelihood  $\times$  prior

$p(\mathcal{D} | \mathbf{w})$

– the likelihood function for  $\mathcal{D} = \mathbf{t}$   
(not a probability distrib. over  $\mathbf{w}$ , just a conditional probability)

$p(\mathbf{w} | \mathcal{D})$

– a posterior probability



# Learning as inference

- Learning distributions

➤ curve-fitting example:  $y(x, \mathbf{w}): X \rightarrow Y$  (let's assume polynomial)

$$t = y(x, \mathbf{w}) + \varepsilon$$

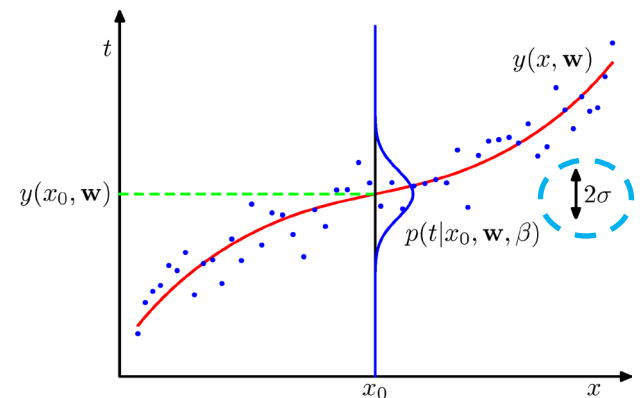
observations  $\{(x_i, t_i): i=1, \dots, N\}$ :  $\mathbf{x} = (x_1, \dots, x_N)^T$ ,  $\mathbf{t} = (t_1, \dots, t_N)^T$

Uncertainty (*noise*) in target data:

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



*precision*



*Bishop, sec. 1.2.5*

# Learning as inference – estimate parameters

- Learning distributions

➤ curve-fitting example:  $y(x, \mathbf{w}): X \rightarrow Y$  (let's assume polynomial)

$$t = y(x, \mathbf{w}) + \varepsilon$$

observations  $\{(x_i, t_i): i=1, \dots, N\}$ :  $\mathbf{x} = (x_1, \dots, x_N)^T$ ,  $\mathbf{t} = (t_1, \dots, t_N)^T$

## SOLUTION 1

We want to find parameters to be able to use predictive distribution,  $p(t|x)$ , and infer the target:

$$\mathbb{E}[t | x] = \int t p(t | x, \mathbf{w}_{\text{opt}}, \beta_{\text{opt}}) dt$$

- Introduction
- **Probabilistic approach**
- Probability basics

# Learning as inference – max likelihood

So, we follow the **max likelihood** approach

ML function for  $t$  (i.i.d.) under Gaussian noise  $p(t | x, \mathbf{w}, \beta) = \mathcal{N}(t | y(x, \mathbf{w}), \beta^{-1})$

$$p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | y_n(x_n, \mathbf{w}), \beta^{-1})$$

- Introduction
- **Probabilistic approach**
- Probability basics

# Learning as inference – max likelihood

So, we follow the **max likelihood** approach

ML function for  $t$  (i.i.d.) under Gaussian noise  $p(t | x, \mathbf{w}, \beta) = \mathcal{N}(t | y(x, \mathbf{w}), \beta^{-1})$

$$p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | y_n(x_n, \mathbf{w}), \beta^{-1})$$

The log-likelihood:

$$\ln p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) = -\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi)$$

- Introduction
- **Probabilistic approach**
- Probability basics

# Learning as inference – max likelihood

So, we follow the **max likelihood** approach

ML function for  $t$  (i.i.d.) under Gaussian noise  $p(t | x, \mathbf{w}, \beta) = \mathcal{N}(t | y(x, \mathbf{w}), \beta^{-1})$

$$p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | y_n(x_n, \mathbf{w}), \beta^{-1})$$

The log-likelihood:

$$\ln p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) = -\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi)$$

Still, input  $x$  is one-dim and  $\mathbf{x}$  &  $\mathbf{t}$  refer to the collection of all data points.

- Introduction
- **Probabilistic approach**
- Probability basics

# Learning as inference – max likelihood

Maximum likelihood (ML) estimate:

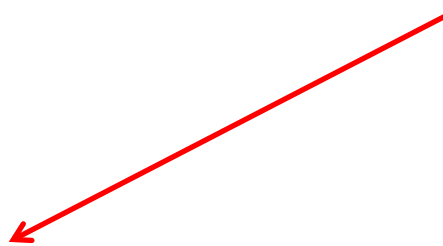
**Maximise**  $\ln p(\mathbf{t} \mid \mathbf{x}, \mathbf{w}, \beta) = -\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi)$

- Introduction
- **Probabilistic approach**
- Probability basics

# Learning as inference – max likelihood

Maximum likelihood (ML) estimate:

**Maximise**  $\ln p(\mathbf{t} \mid \mathbf{x}, \mathbf{w}, \beta) = -\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi)$


$$\mathbf{w}_{\text{ML}} = \arg \min_{\mathbf{w}} \left\{ \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 \right\}$$

the sum-of-squares error function

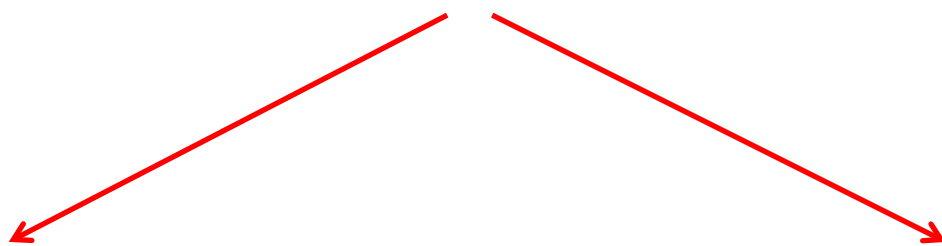
under the assumption of Gaussian noise :  $p(t \mid x, \mathbf{w}, \beta) = \mathcal{N}(t \mid y(x, \mathbf{w}), \beta^{-1})$

- Introduction
- **Probabilistic approach**
- Probability basics

# Learning as inference – max likelihood

Maximum likelihood (ML) estimate:

**Maximise**  $\ln p(\mathbf{t} \mid \mathbf{x}, \mathbf{w}, \beta) = -\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi)$



$$\mathbf{w}_{\text{ML}} = \arg \min_{\mathbf{w}} \left\{ \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 \right\}$$

$$\beta_{\text{ML}}^{-1} = \frac{1}{N} \sum_{n=1}^N \{y(x_n, \mathbf{w}_{\text{ML}}) - t_n\}^2$$



- Introduction
- **Probabilistic approach**
- Probability basics

# Learning as inference – max likelihood

Maximum likelihood (ML) estimate:

**Maximise**  $\ln p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) = -\frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi)$

$$\mathbf{w}_{\text{ML}} = \arg \min_{\mathbf{w}} \left\{ \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 \right\}$$

$$\beta_{\text{ML}}^{-1} = \frac{1}{N} \sum_{n=1}^N \{y(x_n, \mathbf{w}_{\text{ML}}) - t_n\}^2$$



$$t_{\text{out}} = \mathbb{E}[t | x_{\text{in}}] = \int t p(t | x_{\text{in}}, \mathbf{w}_{\text{ML}}, \beta_{\text{ML}}) dt$$



# Learning as inference – posterior over parameters

- Learning distributions

➤ curve-fitting example:  $y(x, \mathbf{w}): X \rightarrow Y$  (let's assume polynomial)

$$t = y(x, \mathbf{w}) + \varepsilon$$

## SOLUTION 2

Parameters of the model could also be random variables with a prior distribution,  $p(\mathbf{w} | \alpha) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \alpha^{-1} \mathbf{I})$

Now, we must **maximise the posterior**  $p(\mathbf{w} | \mathbf{x}, \mathbf{t}, \alpha, \beta)$ , i.e. find the most probable  $\mathbf{w}$  given data:

$$\max p(\mathbf{w} | \mathbf{x}, \mathbf{t}, \alpha, \beta) \propto p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) p(\mathbf{w} | \alpha)$$

# Learning as inference – max posterior

- Learning distributions

➤ curve-fitting example:  $y(x, \mathbf{w}): X \rightarrow Y$  (let's assume polynomial)

$$t = y(x, \mathbf{w}) + \varepsilon$$

## SOLUTION 2

Parameters of the model could also be random variables with a prior distribution,  $p(\mathbf{w} | \alpha) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \alpha^{-1} \mathbf{I})$

Now, we must **maximise the posterior**  $p(\mathbf{w} | \mathbf{x}, \mathbf{t}, \alpha, \beta)$ , i.e. find the most probable  $\mathbf{w}$  given data:

$$\max p(\mathbf{w} | \mathbf{x}, \mathbf{t}, \alpha, \beta) \propto p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) p(\mathbf{w} | \alpha)$$

Again,  $\mathbf{x}$  &  $\mathbf{t}$  refer to the data collection, not individual input or outputs

- Introduction
- **Probabilistic approach**
- Probability basics

# Learning as inference – max posterior

Maximum posterior (MAP) estimate:

**Maximise**  $\ln p(\mathbf{w} \mid \mathbf{x}, \mathbf{t}, \alpha, \beta)$



$$\mathbf{w}_{\text{MAP}} = \arg \min_{\mathbf{w}} \left\{ \frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} \right\}$$

- Introduction
- **Probabilistic approach**
- Probability basics

# Learning as inference – recap

Maximum posterior (MAP) estimate:

**Maximise**  $\ln p(\mathbf{w} \mid \mathbf{x}, \mathbf{t}, \alpha, \beta)$



$$\mathbf{w}_{\text{MAP}} = \arg \min_{\mathbf{w}} \left\{ \frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} \right\}$$

**BUT:** SOLUTIONS 1 (ML) and 2 (MAP) give point estimates of  $\mathbf{w}$

- Introduction
- **Probabilistic approach**
- Probability basics

# Learning as inference – Bayesian view

Instead of estimating “optimal” parameters  $\mathbf{w}$ , let’s integrate over all values of  $\mathbf{w}$  (let’s make use of the distribution)

Marginalisation:

$$p(t \mid x, \mathbf{x}, \mathbf{t}, \alpha, \beta) = \int p(t \mid x, \mathbf{w}, \beta) p(\mathbf{w} \mid \mathbf{x}, \mathbf{t}, \alpha, \beta) d\mathbf{w}$$

predictive      “noise” model      posterior

We assume we know what  $\alpha$  and  $\beta$  are.

- Introduction
- **Probabilistic approach**
- Probability basics

# Learning as inference – Bayesian view

Instead of estimating “optimal” parameters  $\mathbf{w}$ , let’s integrate over all values of  $\mathbf{w}$  for our linear model  $y = \sum_{i=0}^M w_i \phi_i(x)$

For a one-dim polynomial model:  $\phi_i(x) = x^i$  (order  $M-1$ )

- Introduction
- **Probabilistic approach**
- Probability basics

# Learning as inference – Bayesian view

Instead of estimating “optimal” parameters  $\mathbf{w}$ , let’s integrate over all values of  $\mathbf{w}$  for our linear model  $y = \sum_{i=0}^M w_i \phi_i(x)$

For a one-dim polynomial model:  $\phi_i(x) = x^i$

$$p(t \mid x, \mathbf{x}, \mathbf{t}, \alpha, \beta) = \mathcal{N}(t \mid m(x), s^2(x)) \begin{cases} m(x) = \beta \boldsymbol{\phi}(x)^T \mathbf{S} \sum_{n=1}^N \boldsymbol{\phi}(x_n) t_n \\ s^2(x) = \beta^{-1} + \boldsymbol{\phi}(x)^T \mathbf{S} \boldsymbol{\phi}(x) \end{cases}$$

$$\mathbf{S}^{-1} = \alpha \mathbf{I} + \beta \sum_{n=1}^N \boldsymbol{\phi}(x_n) \boldsymbol{\phi}(x_n)^T$$

$$\boldsymbol{\phi}(x) = (\phi_0(x), \dots, \phi_M(x))$$



- Introduction
- **Probabilistic approach**
- Probability basics

# Learning as inference – summary

In summary:

To use a predictive distribution and infer the output  $t$  for the given input  $x_{in}$  .....

- Introduction
- **Probabilistic approach**
- Probability basics

# Learning as inference – summary

In summary:

To use a predictive distribution and infer the output  $t$  for the given input  $x_{in}$  .....

.....ML and MAP approaches produce point estimates of  $\mathbf{w}$

$$\text{ML:} \quad \mathcal{D} \rightarrow \mathbf{w}_{\text{ML}}$$

$$\text{MAP:} \quad \mathcal{D}, p(\mathbf{w}) \rightarrow \mathbf{w}_{\text{MAP}}$$

- Introduction
- **Probabilistic approach**
- Probability basics

# Learning as inference – summary

In summary:

To use a predictive distribution and infer the output  $t$  for the given input  $x_{in}$  .....

.....ML and MAP approaches produce point estimates of  $\mathbf{w}$

$$\text{ML: } \mathcal{D} \rightarrow \mathbf{w}_{\text{ML}}$$

$$\text{MAP: } \mathcal{D}, p(\mathbf{w}) \rightarrow \mathbf{w}_{\text{MAP}}$$

.....in Bayesian approach  $\mathbf{w}$  is integrated over (**marginalisation**)

$$\text{Bayes: } \mathcal{D}, p(\mathbf{w}) \rightarrow p(\mathbf{w} | \mathcal{D})$$

- Introduction
- **Probabilistic approach**
- Probability basics

# Learning as inference – summary

In summary:

To use a predictive distribution and infer the output  $t$  for the given input  $x_{in}$  .....

.....ML and MAP approaches produce point estimates of  $\mathbf{w}$

$$\text{ML: } \mathcal{D} \rightarrow \mathbf{w}_{\text{ML}}$$

$$\text{MAP: } \mathcal{D}, p(\mathbf{w}) \rightarrow \mathbf{w}_{\text{MAP}}$$



Use  $p(t | x_{in}, \mathbf{w}_{\text{ML}})$  or  $p(t | x_{in}, \mathbf{w}_{\text{MAP}})$  for prediction.

.....in Bayesian approach  $\mathbf{w}$  is integrated over (**marginalisation**)

$$\text{Bayes: } \mathcal{D}, p(\mathbf{w}) \rightarrow p(\mathbf{w} | \mathcal{D})$$



Marginalise over  $\mathbf{w}$ :

$$p(t | \mathcal{D}) = \int p(t | x_{in}, \mathbf{w}) p(\mathbf{w} | \mathcal{D}) d\mathbf{w}$$

- Introduction
- **Probabilistic approach**
- Probability basics

# Learning as inference – summary

In summary:

To use a predictive distribution and infer the output  $t$  for the given input  $x_{in}$  .....

.....ML and MAP approaches produce point estimates of  $\mathbf{w}$

$$\text{ML: } \mathcal{D} \rightarrow \mathbf{w}_{\text{ML}}$$

$$\text{MAP: } \mathcal{D}, p(\mathbf{w}) \rightarrow \mathbf{w}_{\text{MAP}}$$

**Frequentist philosophy**

$p(t | x_{in}, \mathbf{w}_{\text{MAP}})$  for prediction.

.....in Bayesian approach  $\mathbf{w}$  is integrated over (**marginalisation**)

$$\text{Bayes: } \mathcal{D}, p(\mathbf{w}) \rightarrow p(\mathbf{w} | \mathcal{D})$$

**Bayesian philosophy**

$$p(t | \mathcal{D}) = \int p(t | x_{in}, \mathbf{w}) p(\mathbf{w} | \mathcal{D}) d\mathbf{w}$$

- Introduction
- **Probabilistic approach**
- Probability basics

# Inference and decision (classification)

The ***inference*** stage of classification  $\mathcal{D} \rightarrow p(C_k, \mathbf{x}), k = 1, \dots, K$

- Introduction
- **Probabilistic approach**
- Probability basics

# Inference and decision (classification)

The ***inference*** stage of classification  $\mathcal{D} \rightarrow p(C_k, \mathbf{x}), k = 1, \dots, K$

starting from *class-conditional densities and priors*

Model the inputs  $\mathbf{x}$  and outputs  $C$

$$\left. \begin{array}{l} p(\mathbf{x} | C_k) \\ p(C_k) \end{array} \right\}$$

for each class  $C_k$

$$p(\mathbf{x}, C_k)$$

$$p(C_k | \mathbf{x}) = \frac{p(\mathbf{x} | C_k) p(C_k)}{p(\mathbf{x})}$$

$$p(\mathbf{x}) = \sum_k p(\mathbf{x} | C_k) p(C_k)$$

**GENERATIVE** approach

Remarks:


- 1)  $K$  classes
- 2)  $\mathbf{x}$  – multi-dim input feature vector

- Introduction
- **Probabilistic approach**
- Probability basics

# Inference and decision (classification)

The ***inference*** stage of classification  $\mathcal{D} \rightarrow p(C_k, \mathbf{x}), k = 1, \dots, K$

Model the inputs  $\mathbf{x}$  and outputs  $C$

$$\left. \begin{array}{l} p(\mathbf{x} | C_k) \\ p(C_k) \end{array} \right\} \text{ for each class } C_k$$


$$p(\mathbf{x}, C_k)$$

$$p(C_k | \mathbf{x}) = \frac{p(\mathbf{x} | C_k) p(C_k)}{p(\mathbf{x})}$$

$$p(\mathbf{x}) = \sum_k p(\mathbf{x} | C_k) p(C_k)$$

**GENERATIVE** approach

Solve first the inference problem of determining ***posteriors*** for each class **without** modelling  $p(C_k, \mathbf{x})$

$$p(C_k | \mathbf{x})$$

**DISCRIMINATIVE** approach



# Generative vs discriminative approach

What are the virtues of the generative approach?

- The parameters are estimated separately for each class (no need to retrain the model when new classes are added)
- Rather straightforward to fit in a Bayesian framework (but it depends on the problem, sometimes discriminative function can be easier to optimise)

# Generative vs discriminative approach

## What are the virtues of the generative approach?

- The parameters are estimated separately for each class (no need to retrain the model when new classes are added)
- Rather straightforward to fit in a Bayesian framework (but it depends on the problem, sometimes discriminative function can be easier to optimise)
- Easy and elegant way of handling missing or unlabelled data
- Generative model allows for..... generating data
  - > generative models can be run “backwards”

# Generative vs discriminative approach

## What are the virtues of the generative approach?

- The parameters are estimated separately for each class (no need to retrain the model when new classes are added)
- Rather straightforward to fit in a Bayesian framework (but it depends on the problem, sometimes discriminative function can be easier to optimise)
- Easy and elegant way of handling missing or unlabelled data
- Generative model allows for..... generating data
  - > generative models can be run “backwards”
- BUT: discriminative models tend to be more accurate (less vulnerable to assumptions)

- Introduction
- **Probabilistic approach**
- Probability basics

# Bias-variance: frequentist mindset

Expected value over all possible datasets  $D$   
(frequentist perspective: data is random)

$$\mathbb{E}[L] = \mathbb{E}_{\mathcal{D}} \{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}^2 + \int \text{var}[t | \mathbf{x}] p(\mathbf{x}) d\mathbf{x}$$

$$\{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}^2 = \{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] + \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2$$

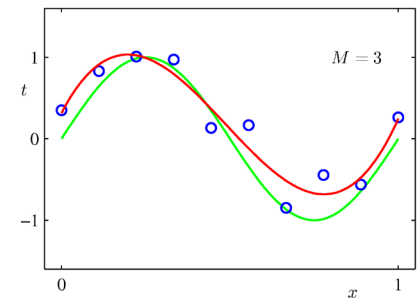
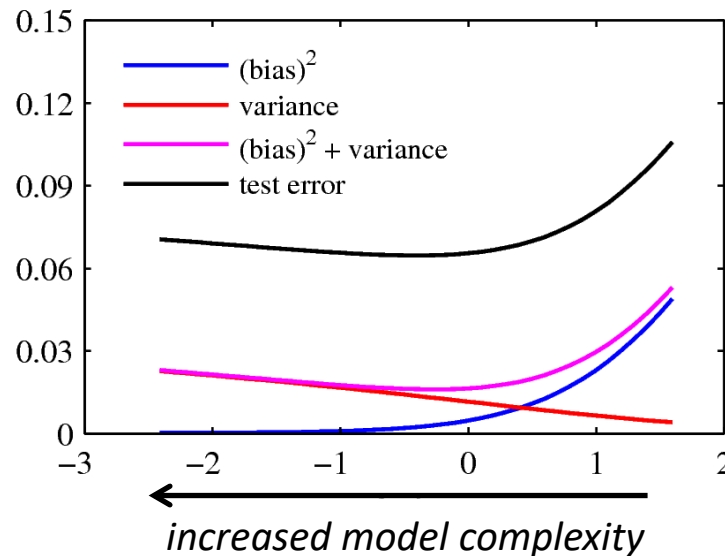
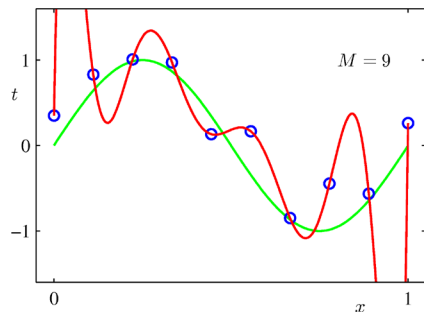
$$\mathbb{E}_{\mathcal{D}} \left[ \{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}^2 \right] = \underbrace{\left\{ \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x}) \right\}^2}_{\text{(bias)}^2} + \underbrace{\mathbb{E}_{\mathcal{D}} \left[ \{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2 \right]}_{\text{variance}}$$

- Introduction
- **Probabilistic approach**
- Probability basics

# Bias-variance as a frequentist dilemma

$$\mathbb{E}_{\mathcal{D}}[L] = \mathbb{E}_{\mathcal{D}} \left\{ y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x}) \right\}^2 + \int \text{var}[t | \mathbf{x}] p(\mathbf{x}) d\mathbf{x}$$

$$\mathbb{E}[L] = (\text{bias})^2 + \text{variance} + \text{noise}$$



*Bishop, sec. 3.2*

# Model selection: frequentist vs Bayesian

- Occam's razor – *“Accept the simplest explanation that fits the data.”*
- Frequentist approach with maximum likelihood
  - bias-variance dilemma
  - need to control the model's complexity
    - regularisation
    - correction for the bias of ML estimates (AIC, BIC)
    - empirical estimate of generalisation error on a hold-out set (validation, resampling)
    - structural risk minimization (SRM) (minimise upper bound on the true risk), see also VC dimension (statistical learning theory)

# Model selection: frequentist vs Bayesian

- Occam's razor – *“Accept the simplest explanation that fits the data.”*
- Frequentist approach with maximum likelihood
- Bayesian approach using **model evidence**

$$p(\mathbf{w} \mid \mathcal{D}) = \frac{p(\mathcal{D} \mid \mathbf{w})p(\mathbf{w})}{p(\mathcal{D})} \implies p(\mathbf{w} \mid \mathcal{D}, \mathcal{M}_i) = \frac{p(\mathcal{D} \mid \mathbf{w}, \mathcal{M}_i)p(\mathbf{w} \mid \mathcal{M}_i)}{p(\mathcal{D} \mid \mathcal{M}_i)}$$

$p(\mathcal{D})$  shows where the model spreads its probability mass over the data space

# Model selection: frequentist vs Bayesian

- Occam's razor – *“Accept the simplest explanation that fits the data.”*
- Frequentist approach with maximum likelihood
- Bayesian approach using **model evidence**

$$p(\mathbf{w} | \mathcal{D}) = \frac{p(\mathcal{D} | \mathbf{w})p(\mathbf{w})}{p(\mathcal{D})} \implies p(\mathbf{w} | \mathcal{D}, \mathcal{M}_i) = \frac{p(\mathcal{D} | \mathbf{w}, \mathcal{M}_i)p(\mathbf{w} | \mathcal{M}_i)}{p(\mathcal{D} | \mathcal{M}_i)}$$

$p(\mathcal{D})$  shows where the model spreads its probability mass over the data space

Towards the **posterior**:  $p(\mathcal{M}_i | \mathcal{D}) \propto p(\mathcal{M}_i)p(\mathcal{D} | \mathcal{M}_i)$

**marginal likelihood**



# Model selection: frequentist vs Bayesian

- Occam's razor – *“Accept the simplest explanation that fits the data.”*
- Frequentist approach with maximum likelihood
- Bayesian approach using **model evidence**

$$p(\mathbf{w} | \mathcal{D}) = \frac{p(\mathcal{D} | \mathbf{w})p(\mathbf{w})}{p(\mathcal{D})} \implies p(\mathbf{w} | \mathcal{D}, \mathcal{M}_i) = \frac{p(\mathcal{D} | \mathbf{w}, \mathcal{M}_i)p(\mathbf{w} | \mathcal{M}_i)}{p(\mathcal{D} | \mathcal{M}_i)}$$

calculating  $p(\mathcal{D} | \mathcal{M}_i)$  is not so trivial:

$$p(\mathcal{D} | \mathcal{M}_i) = \int p(\mathcal{D} | \mathbf{w}, \mathcal{M}_i)p(\mathbf{w} | \mathcal{M}_i)d\mathbf{w}$$

*“The evidence can be seen as the probability of generating the data set from a model whose parameters are sampled at random from the prior”*