

# Statistical Machine Learning

Lecturers: Bert Kappen, Tom Claassen  
Radboud University Nijmegen, Netherlands

September 12, 2016

Book:

Pattern Recognition and Machine Learning  
C.M. Bishop  
<http://research.microsoft.com/~cmbishop/PRML/>

Figures from <http://research.microsoft.com/~cmbishop/PRML/webfigs.htm>

# Statistical Machine Learning

- Course setup

- Lectures (incl. guest lecture), Tue 8:45-10:30, HG00.108
- Tutorials, Thu 13:45-15:30, HG00.622 (week 35), HG01.028 (after)
- Instructors: Dominik Thalmeier, Gabriel Bucur
- Textbook "*Pattern Recognition and Machine Learning*" (C.M Bishop, 2006)
- All other course materials (slides, exercises, sample exams) via Blackboard
- Homework assignments (4 in total, starting week 3)
- Mini seminar (mandatory)
- Written exam (open book, no notes/slides/laptop)

- Grading

- 2/3 final exam ( $\geq 5.0$ ) + 1/3 avg. assignments

# Course contents

## Chapter 1 Introduction

- Probability theory
- Model selection
- Curse of dimensionality
- Decision theory
- information theory

## Chapter 2: Probability distributions

## Chapter 3: Linear models for regression

## Chapter 4: Linear models for classification

## Chapter 5: Neural networks

# Chapter 1: Introduction

## Introduction ML

- General introduction
- Polynomial curve fitting, regression, overfitting, regularization
- Probability theory, decision theory
- Information theory
- Math tools recap

# Recognition of digits

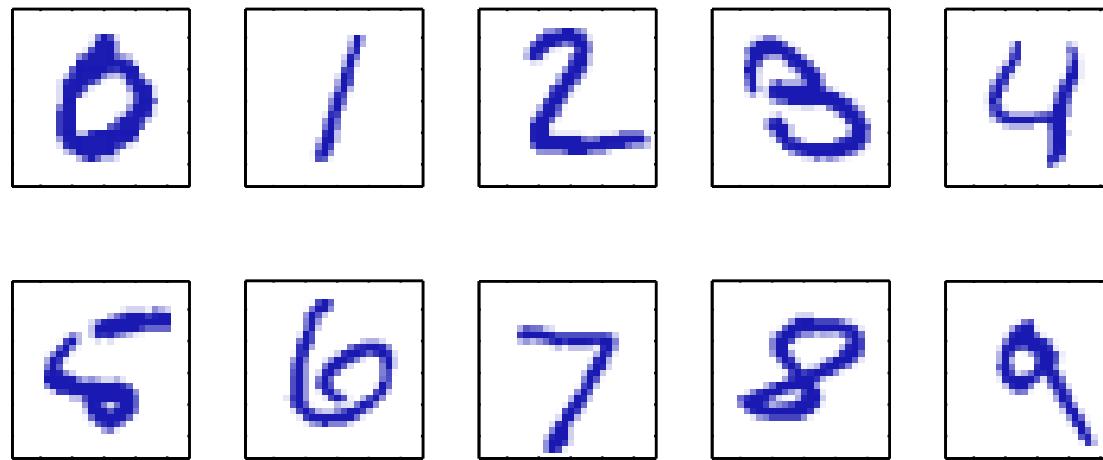


Image is array of pixels  $x_i$ , each between 0 and 1.

$\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_d)$ ,  $d$  is the total number of pixels. ( $d = 28 \times 28$ ).

- Goal = input: pixels  $\rightarrow$  output: correct category  $0, \dots, 9$
- wide variability
- hand-made rules infeasible

# Machine Learning

- Training set: large set of (pixel array, category) pairs
  - input data  $\mathbf{x}$ , target data  $\mathbf{t}(\mathbf{x})$
  - category = class
- Machine learning algorithm
  - adaptive model
  - learning = tuning model parameters on training set
- Result: function  $\mathbf{y}(\mathbf{x})$ 
  - Fitted to target data
  - New input  $\mathbf{x} \rightarrow$  output  $\mathbf{y}(\mathbf{x})$
  - Hopefully:  $\mathbf{y}(\mathbf{x}) \approx \mathbf{t}(\mathbf{x})$
  - Goal: generalization to new examples (use test set)

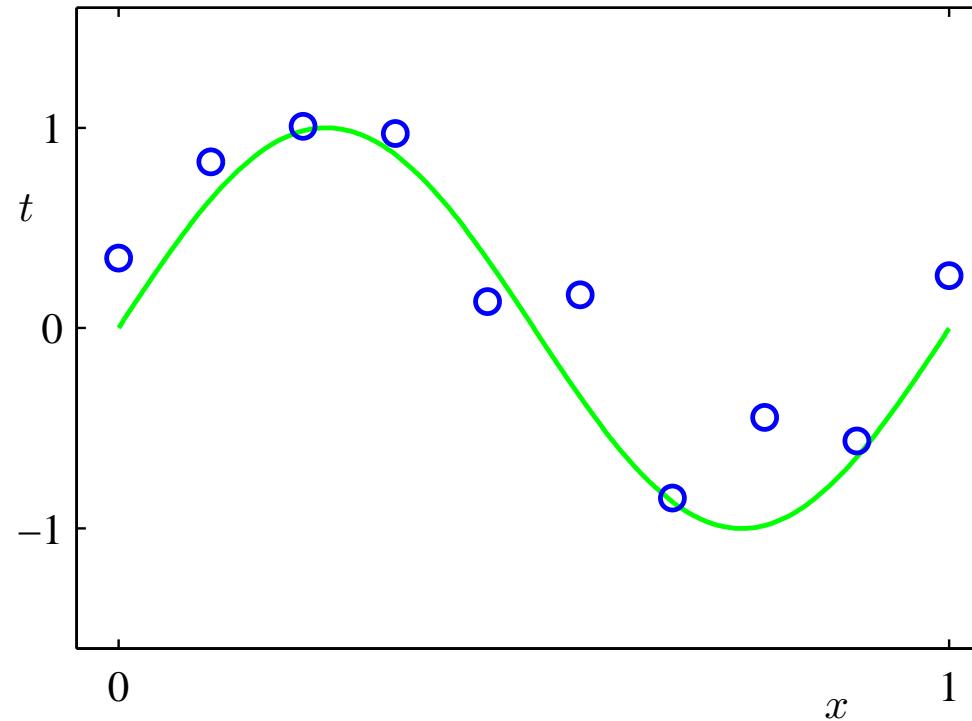
# Preprocessing

- transformation of inputs  $\mathbf{x} \rightarrow \mathbf{x}'$ 
  - easier to handle
  - speed up computation
  - training/test set + new instances
- by hand or rule based
  - scaling, centering
  - aspect ratio, greyscale
- feature extraction
  - dimension reduction
  - reduces variability within class (noise reduction)  $\rightarrow$  easier to learn
  - reduces variability between classes (information loss)  $\rightarrow$  more difficult to learn
  - trade-off

# Types of machine learning tasks

- Supervised learning: known targets
  - Classification: targets are classes
  - Regression: target is continuous
- Unsupervised learning: unknown targets
  - clustering (similarity between input data)
  - density estimation (distribution of input data)
  - dimension reduction (to 2/3D for visualization)
- Reinforcement learning: find optimal target
  - actions leading to maximal reward
  - exploration vs. exploitation

## 1.1. Polynomial curve fitting



- Regression problem
- Given training set of  $N = 10$  data points
  - generated by  $t_n = \sin(2\pi x) + \text{noise}$
- Goal: predict value  $t$  for new  $x$  (without knowing the curve)
  - function  $\hat{t} = y(x)$

We will fit the data using  $M$ -th order polynomial

$$y(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M = \sum_{j=0}^M w_jx^j$$

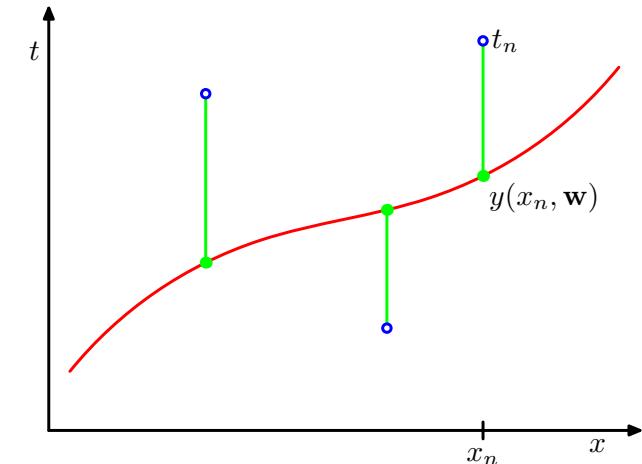
$y(x, \mathbf{w})$  is nonlinear in  $x$ , but linear in coefficients  $\mathbf{w}$ , "Linear model"

Training set:  $(x_n, t_n), n = 1, \dots, N$ . Objective: find parameters  $\mathbf{w}$ , such that  
 $y(x_n, \mathbf{w}) \approx t_n$ , for all  $n$

This is done by minimizing the *Error function*

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

$E(\mathbf{w}) \geq 0$  and  $E(\mathbf{w}) = 0 \Leftrightarrow y(x_n, \mathbf{w}) = t_n$



# Finding the minimum: partial derivatives and gradient

Let  $f(x_1, \dots, x_n) = f(\mathbf{x})$  be a function of several variables. The gradient of  $f$ , denoted as  $\nabla f$  (the symbol “ $\nabla$ ” is called ‘nabla’), is the vector of all partial derivatives:

$$\nabla f(\mathbf{x}) = \left( \frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_n} \right)^T$$

NB, the partial derivative  $\partial f(x_1, \dots, x_i, \dots, x_n)/\partial x_i$  is computed by taking the derivative with respect to  $x_i$  while keeping all other variables constant.

**Example:**

$$f(x, y, z) = xy^2 + 3.1yz$$

Then

$$\begin{aligned} \nabla f(x, y, z) &= \left( \frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z} \right)^T \\ &= (y^2, 2xy + 3.1z, 3.1y)^T \end{aligned}$$

At local minima (and maxima, and so-called saddle points) of a differentiable function  $f$ , the gradient is zero, i.e.,  $\nabla f = 0$ .

**Example:**

$$f(x, y) = x^2 + y^2 + (y + 1)x$$

So

$$\nabla f(x, y) = (2x + y + 1, 2y + x)^T$$

Then we can compute the point  $(x^*, y^*)$  that minimizes  $f$  by setting  $\nabla f = 0$ ,

$$\left. \begin{array}{l} 2x^* + y^* + 1 = 0 \\ 2y^* + x^* = 0 \end{array} \right\} \Rightarrow (x^*, y^*) = \left(-\frac{2}{3}, \frac{1}{3}\right)$$

## Chain rule

Suppose  $f$  is a function of  $y_1, y_2, \dots, y_k$  and each  $y_j$  is a function of  $x$ , then we can compute the derivative of  $f$  with respect to  $x$  by the chain rule

$$\frac{df}{dx} = \sum_{j=1}^k \frac{\partial f}{\partial y_j} \frac{dy_j}{dx}$$

### Example:

$$f(y(x), z(x)) = y(x)/z(x)$$

with  $y(x) = x^4$  and  $z = x^2$ .

Then  $y'(x) = 4x^3$  and  $z'(x) = 2x$ , and so

$$\begin{aligned}\frac{df}{dx} &= \frac{1}{z(x)}y'(x) - \frac{y(x)}{z(x)^2}z'(x) \\ &= \frac{1}{x^2}4x^3 - \frac{x^4}{(x^2)^2}2x \\ &= 2x\end{aligned}$$

## Chain rule (2)

Suppose  $E$  is a function of  $y_1, y_2, \dots, y_N$  and each  $y_j$  is a function of  $w_0, \dots, w_M$ , then we can compute the derivative of  $E$  with respect to  $w_i$  by the chain rule

$$\frac{\partial E}{\partial w_i} = \sum_{j=1}^N \frac{\partial E}{\partial y_j} \frac{\partial y_j}{\partial w_i}$$

### Example:

$$E(\mathbf{w}) = \frac{1}{2} \sum_{j=1}^N (y_j(\mathbf{w}) - t_j)^2 \Rightarrow \frac{\partial E}{\partial y_j} = y_j - t_j$$

and

$$y_j(\mathbf{w}) = \sum_{i=0}^M x_j^i w_i \Rightarrow \frac{\partial y_j}{\partial w_i} = x_j^i$$

So

$$\frac{\partial E}{\partial w_i} = \sum_{j=1}^N (y_j(\mathbf{w}) - t_j) x_j^i$$

$t_j$  and  $x_j^i$  are parameters in this example.

# Minimization of the error function

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

In minimum: gradient  $\nabla_{\mathbf{w}} E = 0$

Note:  $y$  linear in  $\mathbf{w} \Rightarrow E$  quadratic in  $\mathbf{w}$

$\Rightarrow \nabla_{\mathbf{w}} E$  is linear in  $\mathbf{w}$

$\Rightarrow \nabla_{\mathbf{w}} E = 0$ : coupled set of linear equations (exercise)

## Matrix multiplications as summations

If  $A$  is a  $N \times M$  matrix with entries  $A_{ij}$  and  $v$  an  $M$ -dimensional vector with entries  $v_i$ , then  $w = Av$  is a  $N$ -dimensional vector with components

$$w_i = \sum_{j=1}^M A_{ij} v_j$$

In general, if  $B$  is a  $M \times K$  matrix with entries  $B_{ij}$ , then  $C = AB$  is a  $N \times K$  matrix with entries

$$C_{ik} = \sum_{j=1}^M A_{ij} B_{jk}$$

## Dummy indices

The indices that are summed over are ‘dummy’ indices, they are just a label, so e.g.,

$$\sum_{k=1}^M A_{ik} B_{kj} = \sum_{l=1}^M A_{il} B_{lj}$$

furthermore, the entries of the vectors and matrices are just ordinary numbers, so you don’t have to worry about multiplication order. In addition, if the summation of indices is over a range that does not depend on other indices, you may interchange the order of summation,

$$\sum_{i=1}^N \sum_{j=1}^M \dots = \sum_{j=1}^M \sum_{i=1}^N \dots$$

So e.g, by changing summation order and renaming dummy indices,

$$w_k = \sum_{j=1}^M \sum_{i=1}^N A_{ij} B_{jk} = \sum_{l=1}^N \sum_{i=1}^M B_{ik} A_{li}$$

## Kronecker delta

The notation  $\delta_{ij}$  denotes usually the Kronecker delta symbol, i.e.,

$$\begin{cases} \delta_{ij} = 1 & \text{if } i = j \\ \delta_{ij} = 0 & \text{otherwise} \end{cases}$$

It has the nice property that it ‘eats’ dummy indices in summations:

$$\sum_{j=1}^M \delta_{ij} v_j = v_i \quad \text{for all } 1 \leq i \leq M \tag{1}$$

The Kronecker delta can be viewed as the entries of the identity matrix  $\mathbf{I}$ . In vector notation, (1) is equivalent to the statement  $\mathbf{I}\mathbf{v} = \mathbf{v}$ . In other words,  $\delta_{ij} = I_{ij}$ <sup>1</sup>

---

<sup>1</sup>Bishop used  $\delta$  in his previous book, and  $I$  in the current book

## Taylor series, 1-d

Assuming that  $f(x)$  has derivatives of all orders in  $x = a$ , then the Taylor expansion of  $f$  around  $a$  is

$$f(a + \epsilon) = \sum_{k=0}^{\infty} \frac{f^{(k)}(a)}{k!} \epsilon^k = f(a) + \epsilon f'(a) + \frac{\epsilon^2}{2} f''(a) + \dots$$

The prefactors in the Taylor series can be checked by computing the Taylor expansion of a polynomial.

Linearization of a function around  $a$  is taking the Taylor expansion up to first order:

$$f(a + x) = f(a) + x f'(a)$$

# Taylor series, examples

**Examples:** check that for small  $x$  the following expansions are correct up to second order:

$$\begin{aligned}\sin(x) &= \sin(0) + x \cos(0) + \frac{1}{2}x^2(-\sin(0)) + \dots \\ &= 0 + x - 0 + \dots\end{aligned}$$

$$= x$$

$$\cos(x) = 1 - \frac{1}{2}x^2$$

$$\exp(x) = 1 + x + \frac{1}{2}x^2$$

$$(1+x)^c = 1 + cx + \frac{c(c-1)}{2}x^2$$

$$\ln(1+x) = x - \frac{1}{2}x^2$$

## Taylor expansion in several dimensions

The Taylor expansion of a function of several variables,  $f(x_1, \dots, x_n) = f(\mathbf{x})$  is (up to second order)

$$f(\mathbf{x}) = f(\mathbf{a}) + \sum_i (x_i - a_i) \frac{\partial}{\partial x_i} f(\mathbf{a}) + \frac{1}{2} \sum_{ij} (x_i - a_i)(x_j - a_j) \frac{\partial}{\partial x_i} \frac{\partial}{\partial x_j} f(\mathbf{a})$$

or in vector notation, with  $\epsilon = \mathbf{x} - \mathbf{a}$

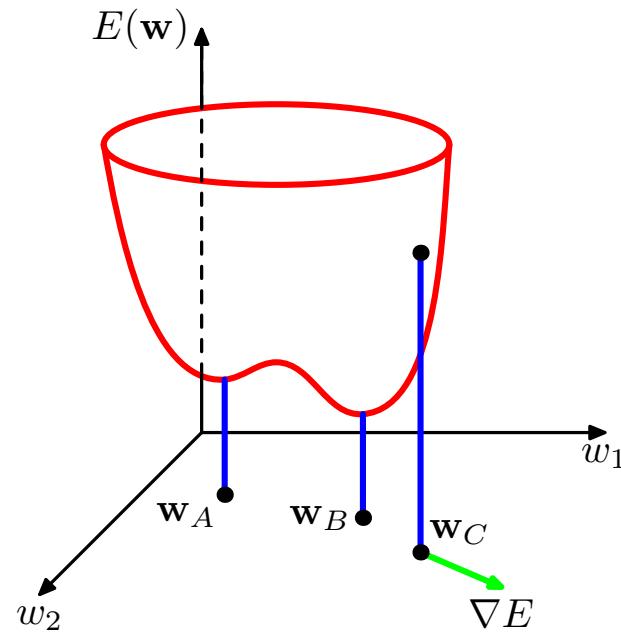
$$f(\mathbf{a} + \epsilon) = f(\mathbf{a}) + \epsilon^T \nabla f(\mathbf{a}) + \frac{1}{2} \epsilon^T \mathbf{H} \epsilon$$

with  $\mathbf{H}$  the Hessian, which is the symmetric matrix of partial derivatives

$$H_{ij} = \left. \frac{\partial^2}{\partial x_i \partial x_j} f(\mathbf{x}) \right|_{\mathbf{x}=\mathbf{a}}$$

## Error landscape

Polynomial curve fitting has a quadratic error function. In general the error function  $E(\mathbf{w})$  may be a non-quadratic in the parameters  $\mathbf{w}$ .



Gradient descent: walk downwards with small steps in the direction of the negative gradient.

$E$  is minimal when  $\nabla E(\mathbf{w}) = 0$ , but not vice versa!

⇒ gradient based methods find a local minimum, not necessarily the global minimum.

# Application: Optimization by gradient descent

Gradient descent algorithm for finding  $\mathbf{w}$  in  $y(x, \mathbf{w})$ :

1. Start with an initial value of  $\mathbf{w}$  and  $\epsilon$  small.
2. While "change in  $\mathbf{w}$  large": Compute  $\mathbf{w} := \mathbf{w} - \epsilon \nabla E$

Stop criterion is  $\nabla E \approx 0$ , which means that we stop in a local minimum of  $E$ .

Does this algorithm converge? Yes, if  $\epsilon$  is "sufficiently small" and  $E$  bounded from below.

Proof: Denote  $\Delta\mathbf{w} = -\epsilon \nabla E$ .

$$E(\mathbf{w} + \Delta\mathbf{w}) \approx E(\mathbf{w}) + (\Delta\mathbf{w})^T \nabla E = E(\mathbf{w}) - \epsilon \sum_i \left( \frac{\partial E}{\partial w_i} \right)^2 \leq E(\mathbf{w})$$

In each gradient descent step the value of  $E$  is lowered. Since  $E$  bounded from below, the procedure must converge asymptotically.

Note: if  $\mathbf{w}$  has a closed-form solution then gradient descent not necessary!

## Newton's method

One can also use Hessian information for optimization. As an example, consider a quadratic approximation to  $E$  around  $\mathbf{w}_0$  (Taylor expansion up to 2<sup>nd</sup> order):

$$\begin{aligned} E(\mathbf{w}) &= E(\mathbf{w}_0) + \mathbf{b}^T(\mathbf{w} - \mathbf{w}_0) + \frac{1}{2}(\mathbf{w} - \mathbf{w}_0)\mathbf{H}(\mathbf{w} - \mathbf{w}_0) \\ b_i &= \frac{\partial E(\mathbf{w}_0)}{\partial w_i} \quad H_{ij} = \frac{\partial^2 E(\mathbf{w}_0)}{\partial w_i \partial w_j} \\ \nabla E(\mathbf{w}) &= \mathbf{b} + \mathbf{H}(\mathbf{w} - \mathbf{w}_0) \end{aligned}$$

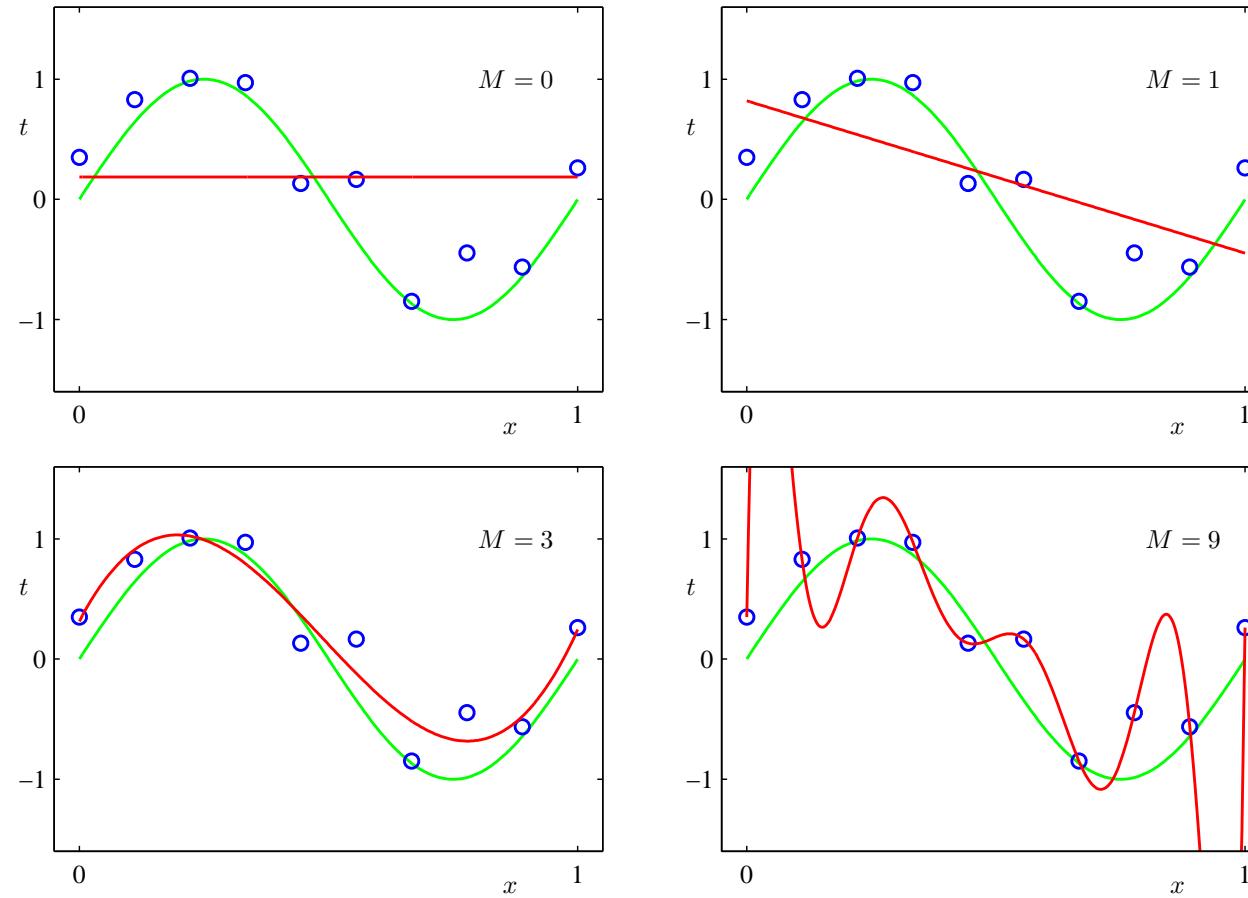
We can solve  $\nabla E(\mathbf{w}) = 0$  and obtain

$$\mathbf{w} = \mathbf{w}_0 - \mathbf{H}^{-1}\nabla E(\mathbf{w}_0)$$

This is called Newton's method. Inversion of the Hessian may be computational costly. A number of methods, known as quasi-newton methods, are based on approximations of this procedure.

# Model comparison, model selection

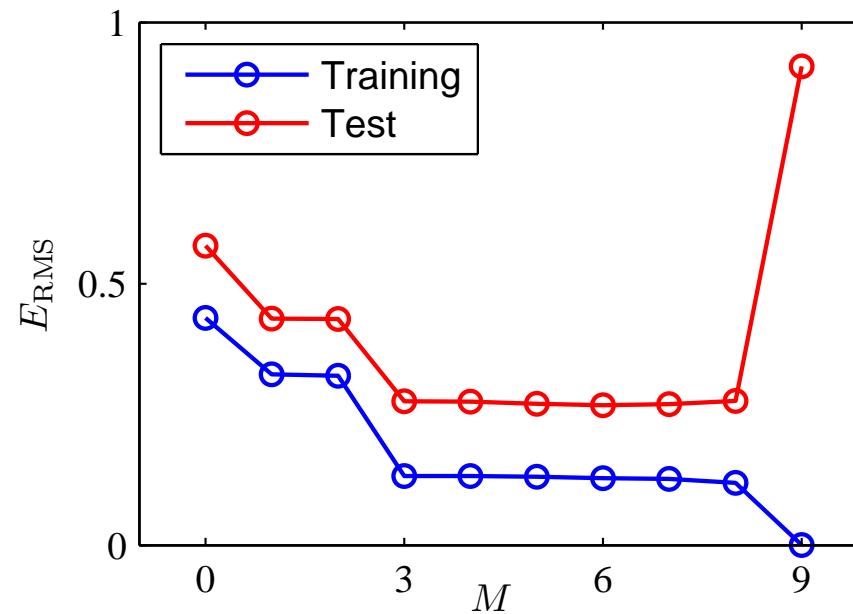
Back to polynomial curve fitting: how to choose  $M$ ?



Which of these models is the best one?

Define root-mean-square error on training set and on test set  $\{(\tilde{x}_n, \tilde{t}_n)\}_{n=1}^{\tilde{N}}$ , respectively:

$$E_{RMS} = \sqrt{2E(\mathbf{w}^*)/N}, \quad E_{RMS} = \sqrt{\frac{1}{\tilde{N}} \sum_{n=1}^{\tilde{N}} (y(\tilde{x}_n, \mathbf{w}^*) - \tilde{t}_n)^2}$$



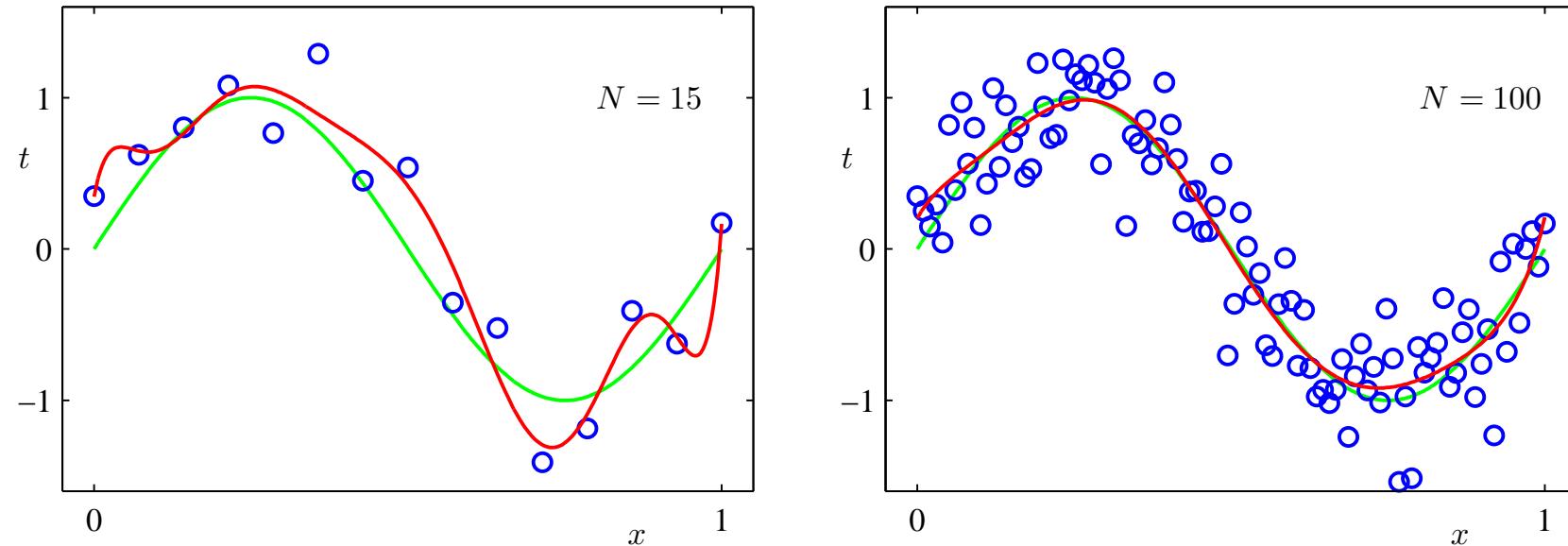
Too simple (small  $M$ )  $\rightarrow$  poor fit

Too complex (large  $M$ )  $\rightarrow$  overfitting (fits the noise)

Q: Taylor expansion of  $\sin(x)$  contains all odd order terms ... shouldn't  $M = 9$  be better?

	M=0	M=1	M=3	M=9
$w_0^*$	0.19	0.82	0.31	0.35
$w_1^*$		-1.27	8	232
$w_2^*$			-25	5321
$w_3^*$			-17	48568
$w_4^*$				-231639
$w_5^*$				640042
$w_6^*$				-10618000
$w_7^*$				10424000
$w_8^*$				-557683
$w_9^*$				-125201

# Model comparison, model selection



Overfitting is not due to noise, but more due to sparseness of data.

Same model complexity: more data  $\Rightarrow$  less overfitting

With more data, more complex (i.e. more flexible) models can be used

# Regularization

Change the cost function  $E$  by adding regularization term  $\Omega(\mathbf{w})$  to penalize complexity.

$$\tilde{E}(\mathbf{w}) = E(\mathbf{w}) + \lambda\Omega(\mathbf{w})$$

For example,

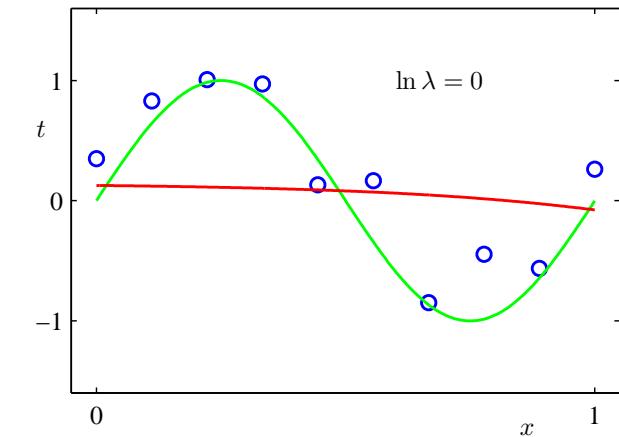
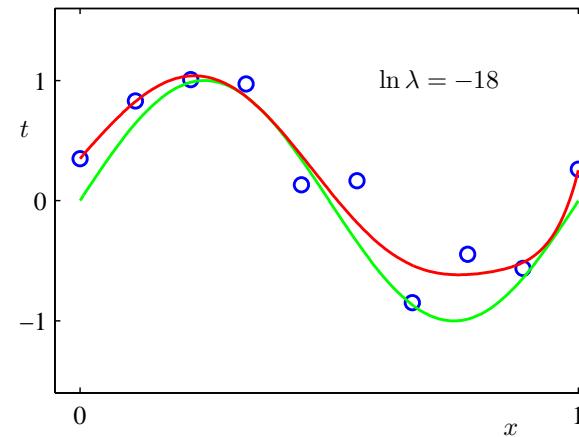
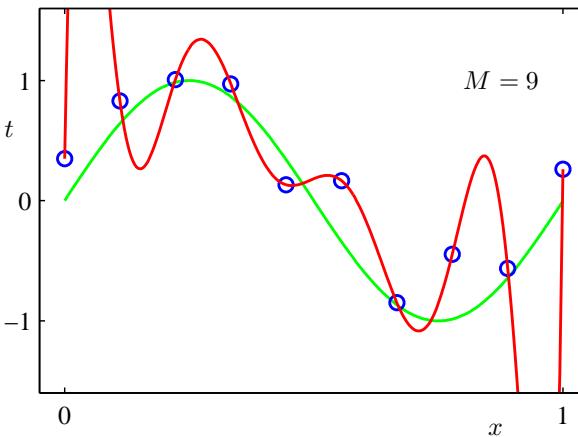
$$\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (y(x_n, \mathbf{w}) - t_n)^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

(here,  $\|\mathbf{w}\|^2 := \sum_{m=0}^M w_m^2$ )

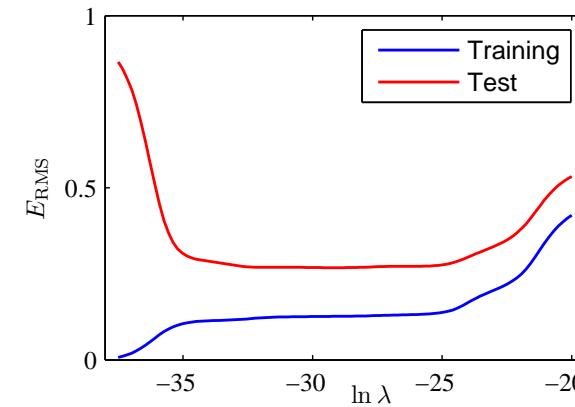
Weight decay = shrinkage = ridge regression

Penalty term independent of number of training data

- small data sets: penalty term relatively large
- large data sets: penalty term relatively small
- → effective complexity depends on #training data



	$\ln \lambda = -\infty$	$\ln \lambda = -18$	$\ln \lambda = 0$
$w_0^*$	0.35	0.35	0.13
$w_1^*$	232	4.74	-0.05
$w_2^*$	5321	-0.77	-0.06
$w_3^*$	48568	-31.97	-0.05
$w_4^*$	-231639	-3.89	-0.03
$w_5^*$	640042	55.28	-0.02
$w_6^*$	-10618000	41.32	-0.01
$w_7^*$	10424000	-45.95	-0.00
$w_8^*$	-557683	-91.53	0.00
$w_9^*$	-125201	72.68	0.01



- *Training set* to optimize (typically many) parameters  $w$
- *Validation set* to optimize (typically a few) hyperparameters  $\lambda, M$
- used in e.g. *cross-validation* (§1.3) ... but even better: *Bayesian* approach!

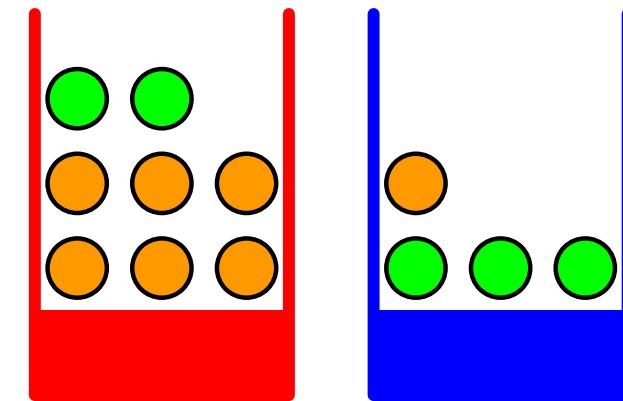
## §1.2 Probability theory

- Consistent framework for quantification and manipulation of uncertainty  
→ Foundation for Bayesian machine learning
- random variable = stochastic variable

**Example:**

boxes ( $B = \{r, b\}$ ) and fruit ( $F = \{a, o\}$ ).

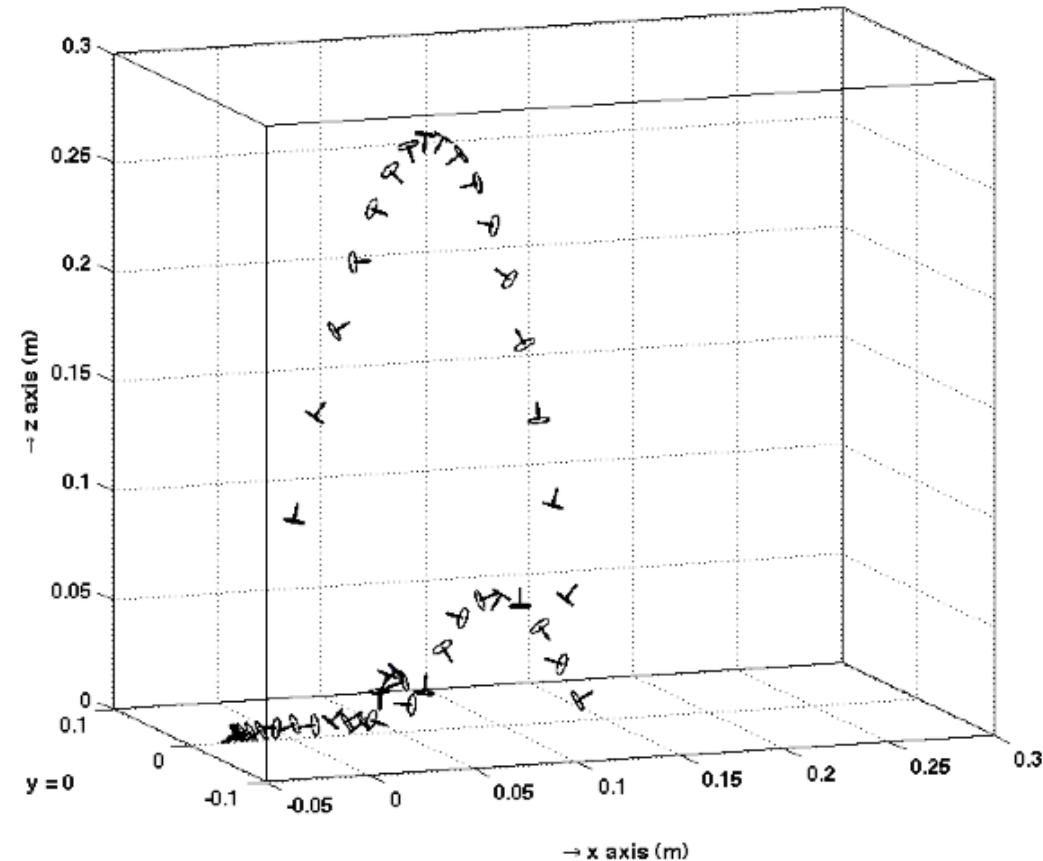
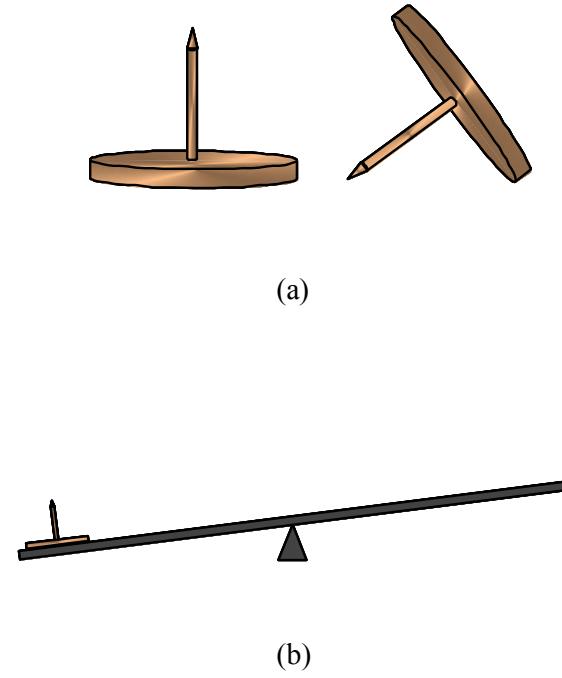
- Consider an experiment of (infinitely) many (mentally) repeated trials  
*(randomly pick a box, then randomly select an item of fruit from that box)*  
under the same macroscopic conditions  
*(number of red/blue boxes and apples/oranges balls in the boxes)*  
but each time with different microscopic details  
*(arrangements of boxes and fruits in boxes).*



*Probability of an event (e.g. selecting a orange) is fraction of times that event occurs in the experiment.*

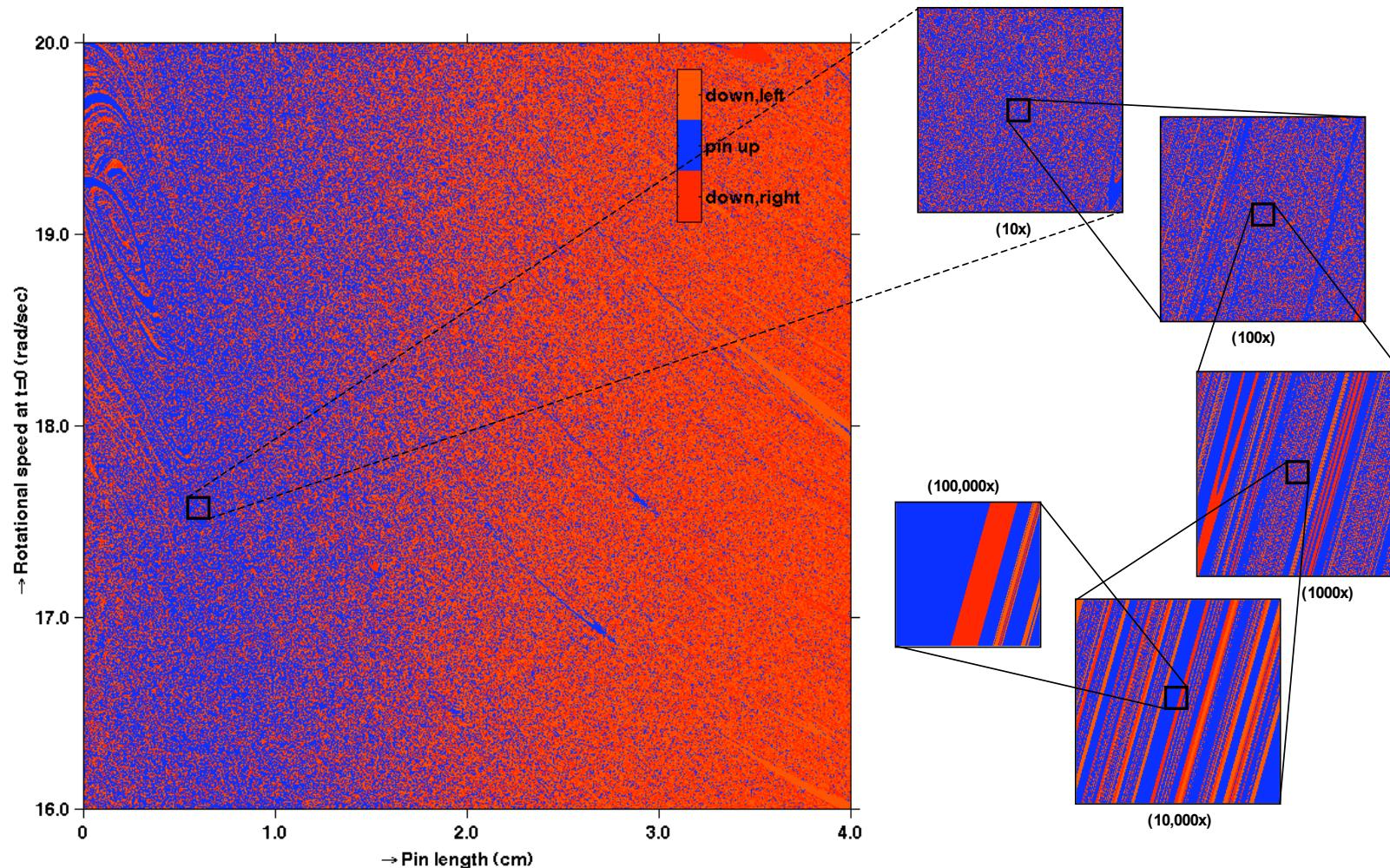
- Notation:  $p(F = o) = 9/20$ , etc (or  $P(\dots)$ ,  $\mathbb{P}(\dots)$ ,  $\text{Prob}(\dots)$ , etc.)

# Experiment with a drawing pin



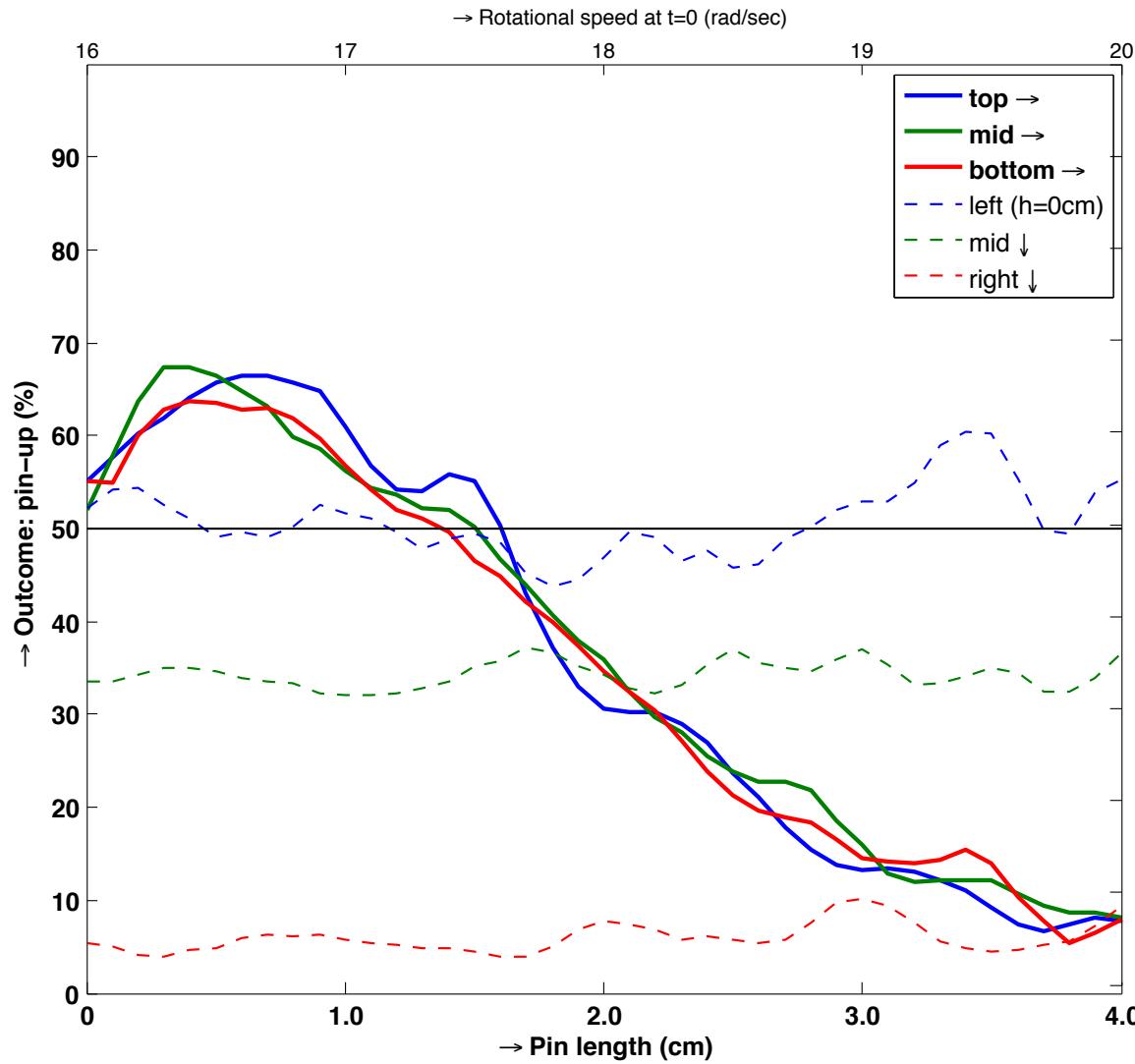
Q: What is the probability of getting 'pin up'?

# Outcome as a function of speed vs. pin length



- ‘Randomness’ in deterministic systems is due to uncertainty about initial states.
- More information does not imply gradual convergence to true value.

# Drawing pin outcomes: sliding averages



A: For typical drawing pin:  $p(\text{Outcome} = \text{'pin up'}) \approx 66\%$   
... but depends on what you *know* about the experiment (no unique, objective value!)

# Joint, marginal, and conditional probabilities

$X$  can take the values  $x_i, i = 1, \dots, M$ .

$Y$  can take the values  $y_j, j = 1, \dots, L$ .

$N$ : total number of trials ( $N \rightarrow \infty$ ).

$n_{ij}$ : number of trials with  $X = x_i$  and  $Y = y_j$

$c_i$ : number of trials with  $X = x_i$

$r_j$ : number of trials with  $Y = y_j$


$x_i$

$y_j$

$n_{ij}$

$c_i$

$r_j$

Joint probability of  $X = x_i$  and  $Y = y_j$ :

$$p(X = x_i, Y = y_j) = \frac{n_{ij}}{N} = p(Y = y_j, X = x_i)$$

# Joint, marginal, and conditional probabilities

$X$  can take the values  $x_i, i = 1, \dots, M$ .

$Y$  can take the values  $y_j, j = 1, \dots, L$ .

$N$ : total number of trials ( $N \rightarrow \infty$ ).

$n_{ij}$ : number of trials with  $X = x_i$  and  $Y = y_j$

$c_i$ : number of trials with  $X = x_i$

$r_j$ : number of trials with  $Y = y_j$

					$c_i$
					$\} r_j$
				$n_{ij}$	
					$y_j$

*Marginal* probability of  $X = x_i$ :

$$p(X = x_i) = \frac{c_i}{N} = \frac{\sum_j n_{ij}}{N} = \sum_j p(X = x_i, Y = y_j)$$

*Conditional* probability of  $Y = y_j$  given  $X = x_i$

$$p(Y = y_j | X = x_i) = \frac{n_{ij}}{c_i}$$

- Explicit, unambiguous notation:  $p(X = x_i)$
- Short-hand notation:  $p(x_i)$
- $p(X)$ : “distribution” over the random variable  $X$
- NB:  $\{x_i\}$  is assumed to be mutually exclusive and complete

## The Rules of Probability

**Sum rule**       $p(X) = \sum_Y p(X, Y)$

**Product rule**       $p(X, Y) = p(Y|X)p(X)$

**Positivity**       $p(X) \geq 0$

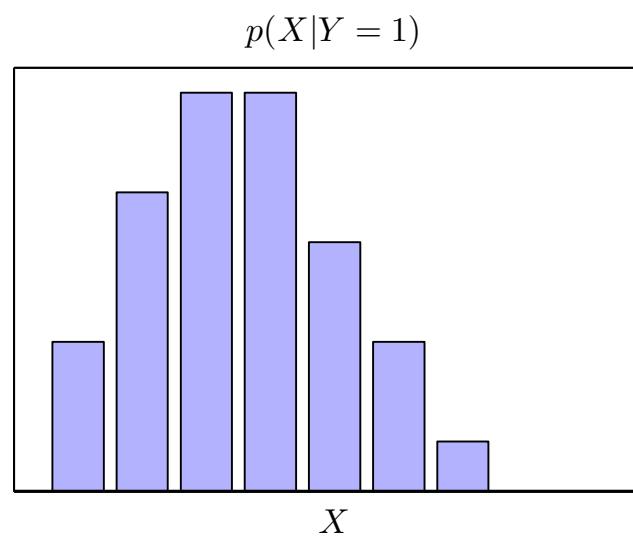
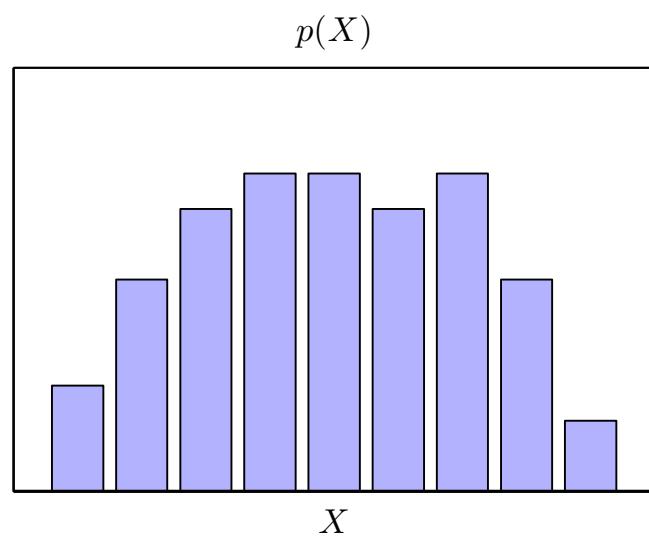
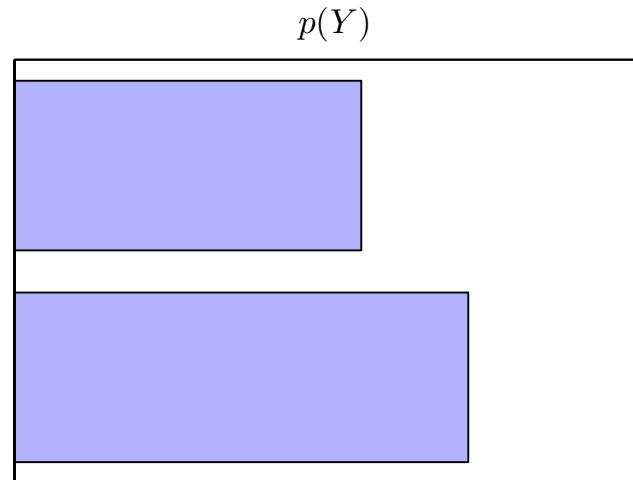
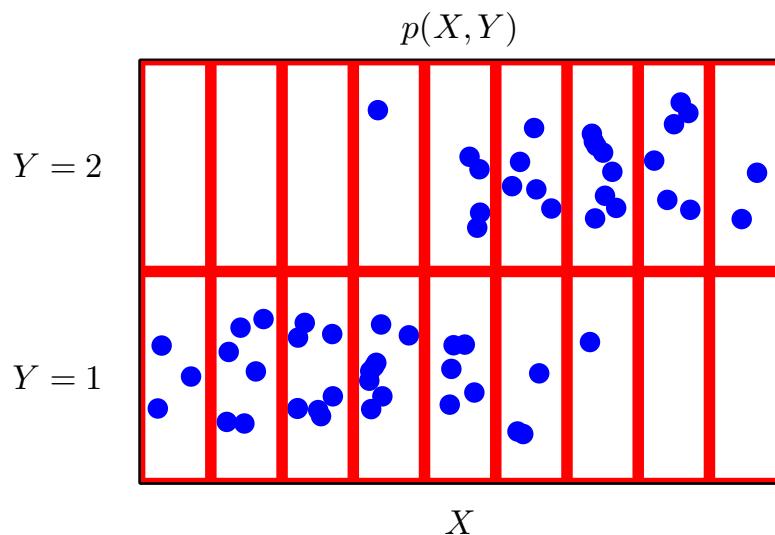
**Normalization**       $\sum_X p(X) = 1$

$$p(X, Y) = p(Y|X)p(X) = P(X|Y)p(Y) \Rightarrow$$

## Bayes' theorem

$$p(Y|X) = \frac{p(X|Y)p(Y)}{p(X)} \quad \left( = \frac{p(X|Y)p(Y)}{\sum_Y p(X|Y)p(Y)} \right)$$

Bayes' theorem = Bayes' rule



# Fruits again

Model

$$p(B = r) = 4/10$$

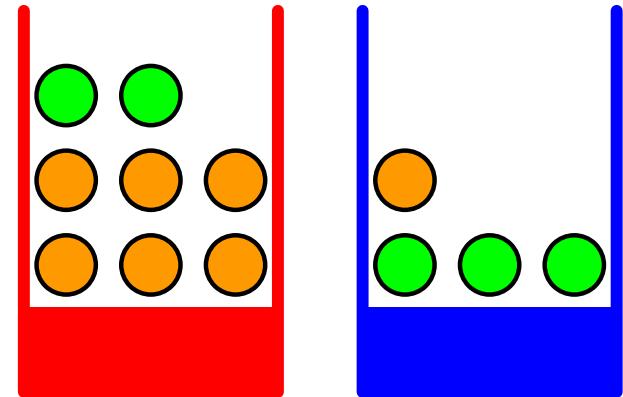
$$p(B = b) = 6/10$$

$$p(F = a|B = r) = 1/4$$

$$p(F = o|B = r) = 3/4$$

$$p(F = a|B = b) = 3/4$$

$$p(F = o|B = b) = 1/4$$



Note that the (conditional) probabilities are normalized:

$$p(B = r) + p(B = b) = 1$$

$$p(F = a|B = r) + p(F = o|B = r) = 1$$

$$p(F = a|B = b) + p(F = o|B = b) = 1$$

- Marginal probability

$$\begin{aligned}
 p(F = a) &= p(F = a|B = r)p(B = r) + p(F = a|B = b)p(B = b) \\
 &= \frac{1}{4} \times \frac{4}{10} + \frac{3}{4} \times \frac{6}{10} = \frac{11}{20}
 \end{aligned}$$

and from normalization,

$$p(F = o) = 1 - p(F = a) = \frac{9}{20}$$

- Conditional probability (reversing probabilities):

$$p(B = r|F = o) = \frac{p(F = o|B = r)p(B = r)}{p(F = o)} = \frac{3}{4} \times \frac{4}{10} \times \frac{20}{9} = \frac{2}{3}$$

- Terminology:

$p(B)$ : prior probability (before observing the fruit)  
 $p(B|F)$ : posterior probability (after observing  $F$ )

## (Conditionally) independent variables

- $X$  and  $Y$  are called (marginally) independent if

$$P(X, Y) = P(X)P(Y)$$

This is equivalent to

$$P(X|Y) = P(X)$$

and also to

$$P(Y|X) = P(Y)$$

- $X$  and  $Y$  are called conditionally independent given  $Z$  if

$$P(X, Y|Z) = P(X|Z)P(Y|Z)$$

This is equivalent to

$$P(X|Y, Z) = P(X|Z)$$

and also to

$$P(Y|X, Z) = P(Y|Z)$$

## Probability densities

- to deal with continuous variables (rather than discrete ones)

When  $x$  takes values from a continuous domain, the probability of any value of  $x$  is zero!  
 Instead, we must talk of the probability that  $x$  takes a value in a certain interval

$$\text{Prob}(x \in [a, b]) = \int_a^b p(x) dx$$

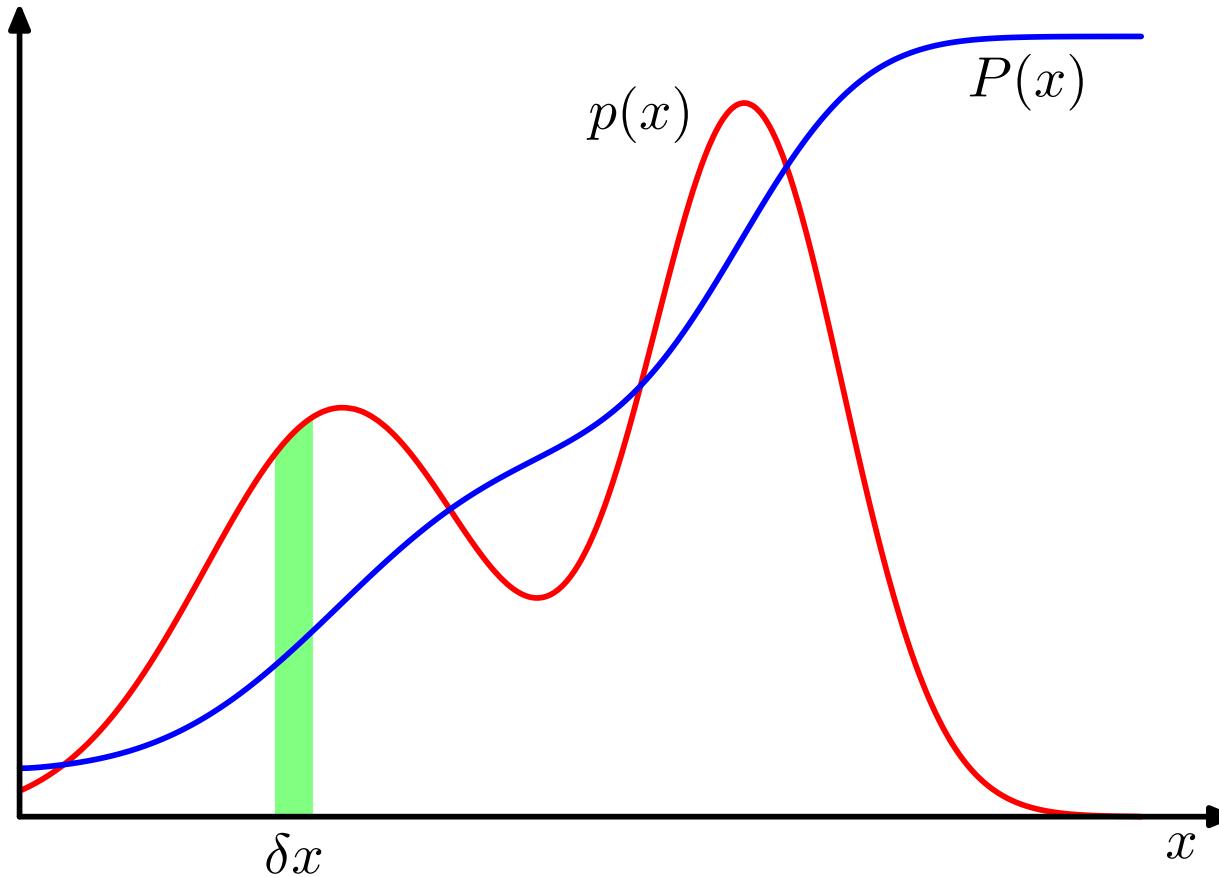
with  $p(x)$  the *probability density* over  $x$ .

$$p(x) \geq 0$$

$$\int_{-\infty}^{\infty} p(x) dx = 1 \quad (\text{normalization})$$

- NB: that  $p(x)$  may be bigger than one.

Probability of  $x$  falling in interval  $(x, x + \delta x)$  is  $p(x)\delta x$  for  $\delta x \rightarrow 0$



Cumulative distribution function  $F(z) = \int_{-\infty}^z p(x)dx$  (not often used in ML).

Note that:

- $\text{Prob}(x \in [a, b]) = F(b) - F(a)$ .
- $F'(z) = p(z)$

## Multivariate densities

- Several continuous variables, denoted by the  $d$  dimensional vector  $\mathbf{x} = (x_1, \dots, x_d)$ .
- Probability density  $p(\mathbf{x})$ : probability of  $\mathbf{x}$  falling in an infinitesimal volume  $\delta\mathbf{x}$  around  $\mathbf{x}$  is given by  $p(\mathbf{x})\delta\mathbf{x}$ .

$$\text{Prob}(x \in \mathcal{R}) = \int_{\mathcal{R}} p(\mathbf{x}) d\mathbf{x} = \int_{\mathcal{R}} p(x_1, \dots, x_d) dx_1 dx_2 \dots dx_d$$

and

$$p(\mathbf{x}) \geq 0$$

$$\int p(\mathbf{x}) d\mathbf{x} = 1$$

- Rules of probability apply to continuous variables as well,

$$p(x) = \int p(x, y) dy$$

$$p(x, y) = p(y|x)p(x)$$

# Integration

The integral of a function of several variables  $\mathbf{x} = (x_1, x_2, \dots, x_n)$

$$\int_{\mathcal{R}} f(\mathbf{x}) d\mathbf{x} \equiv \int_{\mathcal{R}} f(x_1, x_2, \dots, x_n) dx_1 dx_2 \dots dx_n$$

is the volume of the  $n + 1$  dimensional region lying ‘vertically above’ the domain of integration  $\mathcal{R} \subset I\!\!R^n$  and ‘below’ the function  $f(\mathbf{x})$ .

## Separable integrals

The most easy (but important) case is when we can separate the integration, e.g. in 2-d,

$$\int_{x=a}^b \int_{y=c}^d f(x)g(y) dx dy = \int_{x=a}^b f(x) dx \int_{y=c}^d g(y) dy$$

Example,

$$\int \exp\left(\sum_{i=1}^n f_i(x_i)\right) dx = \int \prod_{i=1}^n \exp(f_i(x_i)) dx = \prod_{i=1}^n \int \exp(f_i(x_i)) dx_i$$

## Iterated integration

A little more complicated are the cases, in which integration can be done by iteration, 'from inside out'. Suppose we can write the 2-d region  $\mathcal{R}$  as the set  $a < x < b$  and  $c(x) < y < d(x)$  then we can write

$$\int_{\mathcal{R}} f(y, x) \, dy \, dx = \int_{x=a}^b \left[ \int_{y=c(x)}^{d(x)} f(y, x) \, dy \right] \, dx$$

The first step is evaluate the inner integral, where we interpret  $f(y, x)$  as a function of  $y$  with fixed parameter  $x$ . Suppose we can find  $F$  such that  $\partial F(y, x)/\partial y = f(y, x)$ , then the result of the inner integral is

$$\int_{y=c(x)}^{d(x)} f(y, x) \, dy = F(d(x), x) - F(c(x), x)$$

The result, which we call  $g(x)$  is obviously a function of  $x$  only,

$$g(x) \equiv F(d(x), x) - F(c(x), x)$$

The next step is the outer integral, which is now just a one-dimensional integral of the function  $g$ ,

$$\int_{x=a}^b \left[ \int_{y=c(x)}^{d(x)} f(y, x) dy \right] dx = \int_{x=a}^b g(x) dx$$

Now suppose that the same 2-d region  $\mathcal{R}$  can also be written as the set  $s < y < t$  and  $u(y) < x < v(y)$ , then we can also choose to evaluate the integral as

$$\int_{\mathcal{R}} f(y, x) dx dy = \int_{y=s}^t \left[ \int_{x=u(y)}^{v(y)} f(y, x) dx \right] dy$$

following the same procedure as above. In most regular cases the result is the same (for exceptions, see handout (\*)).

Integration with more than two variables can be done with exactly the same procedure, ‘from inside out’.

In Machine Learning, integration is mostly over the whole of  $x$  space, or over a subspace. Iterated integration is not often used.

# Transformation of variables (1-d)

Often it is easier to do the multidimensional integral in another coordinate frame. Suppose we want to do the integration

$$\int_{y=c}^d f(y) dy$$

but the function  $f(y)$  is easier expressed as a  $f(y(x))$  which is a function of  $x$ . So we want to use  $x$  as integration variable. If  $y$  and  $x$  are related via invertible differentiable mappings  $y = y(x)$  and  $x = x(y)$  and the end points of the interval ( $y = c, y = d$ ) are mapped to ( $x = a, x = b$ ), (so  $c = y(a)$ , etc) then we have the equality

$$\begin{aligned} \int_{y=c}^d f(y) dy &= \int_{y(x)=c}^d f(y(x)) dy(x) \\ &= \int_{x=a}^b f(y(x)) y'(x) dx \end{aligned}$$

The derivative  $y'(x)$  comes in as the ratio between the lengths of the differentials  $dy$  and  $dx$ ,

$$dy = y'(x) dx$$

# Several variables

With several variables, the substitution rule is generalized as follows. We have the invertible mapping  $\mathbf{y} = \mathbf{y}(\mathbf{x})$ . Let us also assume that the region of integration of  $\mathcal{R}$  is mapped by to  $\mathcal{S}$ , (so  $\mathcal{S} = \mathbf{y}(\mathcal{R})$ ), then we have the equality

$$\begin{aligned}\int_{\mathbf{y} \in \mathcal{S}} f(\mathbf{y}) d\mathbf{y} &= \int_{\mathbf{y}(\mathbf{x}) \in \mathcal{S}} f(\mathbf{y}) d\mathbf{y}(\mathbf{x}) \\ &= \int_{\mathbf{x} \in \mathcal{R}} f(\mathbf{y}(\mathbf{x})) \left| \det \left( \frac{\partial \mathbf{y}(\mathbf{x})}{\partial \mathbf{x}} \right) \right| d\mathbf{x}\end{aligned}$$

The factor  $\det \left( \frac{\partial \mathbf{y}(\mathbf{x})}{\partial \mathbf{x}} \right)$  is called the Jacobian of the coordinate transformation. Written out in more detail

$$\det \left( \frac{\partial \mathbf{y}(\mathbf{x})}{\partial \mathbf{x}} \right) = \begin{vmatrix} \frac{\partial y_1(\mathbf{x})}{\partial x_1} & \frac{\partial y_1(\mathbf{x})}{\partial x_2} & \cdots & \frac{\partial y_1(\mathbf{x})}{\partial x_n} \\ \frac{\partial y_2(\mathbf{x})}{\partial x_1} & \frac{\partial y_2(\mathbf{x})}{\partial x_2} & \cdots & \frac{\partial y_2(\mathbf{x})}{\partial x_n} \\ \cdots & \cdots & \cdots & \cdots \\ \frac{\partial y_n(\mathbf{x})}{\partial x_1} & \frac{\partial y_n(\mathbf{x})}{\partial x_2} & \cdots & \frac{\partial y_n(\mathbf{x})}{\partial x_n} \end{vmatrix}$$

The *absolute value*<sup>2</sup> of the Jacobian comes in as the ratio between that the volume represented by the differential  $d\mathbf{y}$  and the volume represented by the differential  $d\mathbf{x}$ , i.e.,

$$d\mathbf{y} = \left| \det \left( \frac{\partial \mathbf{y}(\mathbf{x})}{\partial \mathbf{x}} \right) \right| d\mathbf{x}$$

As a last remark, it is good to know that

$$\det \left( \frac{\partial \mathbf{x}(\mathbf{y})}{\partial \mathbf{y}} \right) = \det \left( \left( \frac{\partial \mathbf{y}(\mathbf{x})}{\partial \mathbf{x}} \right)^{-1} \right) = \frac{1}{\det \left( \frac{\partial \mathbf{y}(\mathbf{x})}{\partial \mathbf{x}} \right)}$$

---

<sup>2</sup>In the single-variable case, we took the orientation of the integration interval into account ( $\int_a^b f(x) dx = - \int_b^a f(x) dx$ ). With several variables, this is awkward. Fortunately, it turns out that the orientation of the mapping of the domain always cancels to the 'orientation' of the Jacobian (= sign of the determinant). Therefore we take a positive orientation and the absolute value of the Jacobian

# Polar coordinates

Example: compute the area of a disc.

Consider a two-dimensional disc with radius  $R$

$$D = \{(x, y) | x^2 + y^2 < R^2\}$$

Its area is

$$\int_D dx dy$$

This integral is easiest evaluated by going to ‘polar-coordinates’. The mapping from polar coordinates  $(r, \theta)$  to Cartesian coordinates  $(x, y)$  is

$$x = r \cos \theta \tag{2}$$

$$y = r \sin \theta \tag{3}$$

Since In polar coordinates, the disc is described by  $0 \leq r < R$  (since  $x^2 + y^2 = r^2$ ) and  $0 \leq \theta < 2\pi$ .

The Jacobian is

$$J = \begin{vmatrix} \frac{\partial x}{\partial r} & \frac{\partial x}{\partial \theta} \\ \frac{\partial y}{\partial r} & \frac{\partial y}{\partial \theta} \end{vmatrix} = \begin{vmatrix} \cos \theta & -r \sin \theta \\ \sin \theta & r \cos \theta \end{vmatrix} = r(\cos^2 \theta + \sin^2 \theta) = r$$

In other words,

$$dxdy = r dr d\theta$$

The area of the disc is now easily evaluated.

$$\int_D dxdy = \int_{\theta=0}^{2\pi} \int_{r=0}^R r dr d\theta = \pi R^2$$

# Gaussian integral

How to compute

$$\int_{-\infty}^{\infty} \exp(-x^2) dx$$

$$\begin{aligned}\left( \int_{-\infty}^{\infty} \exp(-x^2) dx \right)^2 &= \int_{-\infty}^{\infty} \exp(-x^2) dx \int_{-\infty}^{\infty} \exp(-y^2) dy \\&= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \exp(-x^2) \exp(-y^2) dxdy \\&= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \exp(-(x^2 + y^2)) dxdy\end{aligned}$$

The latter is easily evaluated by going to polar-coordinates,

$$\begin{aligned}
 \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \exp(-(x^2 + y^2)) dx dy &= \int_{\theta=0}^{2\pi} \int_{r=0}^{\infty} \exp(-r^2) r dr d\theta \\
 &= 2\pi \int_{r=0}^{\infty} \exp(-r^2) r dr \\
 &= 2\pi \times \left( -\frac{1}{2} \exp(-r^2) \right) \Big|_0^{\infty} \\
 &= \pi
 \end{aligned}$$

So

$$\int_{-\infty}^{\infty} \exp(-x^2) dx = \sqrt{\pi}$$

## Transformation of densities

Under nonlinear change of variables, a probability transforms  $p(\mathbf{x})$  transforms differently from an ordinary function, due to “conservation of probability”  $p(\mathbf{x})\delta\mathbf{x}$ .

Consider  $\mathbf{x} = \mathbf{x}(\mathbf{y})$  then an ordinary function  $f(\mathbf{x})$  becomes  $\tilde{f}(\mathbf{y}) = f(\mathbf{x}(\mathbf{y}))$ .

Probability densities:

- $p_x(\mathbf{x})\delta\mathbf{x}$ : probability that point falls in volume element  $\delta\mathbf{x}$  around  $\mathbf{x}$
- $p_y(\mathbf{y})\delta\mathbf{y}$ : same probability, now in terms of  $\mathbf{y}$

$$p_y(\mathbf{y})\delta\mathbf{y} = p_x(\mathbf{x})\delta\mathbf{x} \quad \Rightarrow \quad p_y(\mathbf{y}) = \left| \det \left( \frac{\partial \mathbf{x}}{\partial \mathbf{y}} \right) \right| p_x(\mathbf{x}(\mathbf{y}))$$

- Values  $p(\mathbf{x})$  of a probability density depends on choice of variable ( $p(\mathbf{x})\delta\mathbf{x}$  is invariant)
- Maximum of a probability density depends on choice of variable (see exercise 1.4).

## Dirac's delta-function

Dirac's delta function  $\delta(x)$  is defined such that

$$\delta(x) = 0 \quad \text{if } x \neq 0 \quad \text{and} \quad \int_{-\infty}^{\infty} \delta(x) dx = 1$$

It can be viewed as the limit  $\Delta \rightarrow 0$  of the function

$$f(x, \Delta) = \frac{1}{\Delta} \quad \text{if } |x| \leq \frac{\Delta}{2} \quad \text{and} \quad f(x, \Delta) = 0 \quad \text{elsewhere}$$

The Dirac delta  $\delta(x)$  is a spike (a peak, a point mass) at  $x = 0$ . The function  $\delta(x - x_0)$  as a function of  $x$  is a spike at  $x_0$ . As a consequence of the definition, the delta function has the important property

$$\int_{-\infty}^{\infty} f(x)\delta(x - x_0)dx = f(x_0)$$

(cf. Kronecker delta  $\sum_j \delta_{ij}v_j = v_i$  ).

The multivariate deltafunction factorizes over the dimensions

$$\delta(\mathbf{x} - \mathbf{m}) = \prod_{i=1}^n \delta(x_i - m_i)$$

# Dirac's delta-function / delta-distribution

The Dirac delta is actually a distribution rather than a function:

$$\delta(\alpha x) = \frac{1}{\alpha} \delta(x)$$

This is true since

- if  $x \neq 0$  left and right-handside are both zero.
- after transformation of variables  $x' = \alpha x$ ,  $dx' = \alpha dx$  we have

$$\int \delta(\alpha x) dx = \frac{1}{\alpha} \int \delta(x') dx' = \frac{1}{\alpha}$$

# Functionals vs functions

Function  $y$ : for any input value  $x$ , returns output value  $f(y)$ .

Functional  $F$ : for any function  $y$ , returns an output value  $F[y]$ .

Example (linear functional):

$$F[y] = \int p(x)y(x) dx$$

(Compare with  $f(\mathbf{y}) = \sum_i p_i y_i$ ).

Other (nonlinear) example:

$$F[y] = \int \frac{1}{2}(y'(x) + V(x))^2 dx$$

# Expectations and variances

Expectations

$$\mathbb{E}[f] = \langle f \rangle = \sum_x p(x)f(x) \quad \text{discrete var's}$$

$$\mathbb{E}[f] = \int_x p(x)f(x) dx \quad \text{continuous var's}$$

Variance:

$$\text{var}[f] = \langle f(x)^2 \rangle - \langle f(x) \rangle^2$$

$$\text{var}[x] = \langle x^2 \rangle - \langle x \rangle^2$$

# Covariance

Covariance of two random variables

$$\text{cov}[x, y] = \langle xy \rangle - \langle x \rangle \langle y \rangle$$

Covariance matrix of the components of a vector variable

$$\text{cov}[\mathbf{x}] \equiv \text{cov}[\mathbf{x}, \mathbf{x}] = \langle \mathbf{x}\mathbf{x}^T \rangle - \langle \mathbf{x} \rangle \langle \mathbf{x}^T \rangle$$

with components

$$(\text{cov}[\mathbf{x}])_{ij} = \text{cov}[x_i, x_j] = \langle x_i x_j \rangle - \langle x_i \rangle \langle x_j \rangle$$

## Bayesian probabilities

- Classical or frequentists interpretation: probabilities in terms of frequencies of random repeatable events
- Bayesian view: probabilities as subjective beliefs about uncertain event
  - event not necessarily repeatable
  - event may yield only indirect observations
  - Bayesian inference to update belief given observations

Why probability theory?

- Cox: common sense axioms for degree of uncertainty → probability theory

**Notation.** Let ‘the degree of belief in proposition  $x$ ’ be denoted by  $B(x)$ . The negation of  $x$  (NOT- $x$ ) is written  $\bar{x}$ . The degree of belief in a conditional proposition, ‘ $x$ , assuming proposition  $y$  to be true’, is represented by  $B(x|y)$ .

**Axiom 1.** Degrees of belief can be ordered; if  $B(x)$  is ‘greater’ than  $B(y)$ , and  $B(y)$  is ‘greater’ than  $B(z)$ , then  $B(x)$  is ‘greater’ than  $B(z)$ .

[Consequence: beliefs can be mapped onto real numbers.]

**Axiom 2.** The degree of belief in a proposition  $x$  and its negation  $\bar{x}$  are related. There is a function  $f$  such that

$$B(x) = f[B(\bar{x})].$$

**Axiom 3.** The degree of belief in a conjunction of propositions  $x, y$  ( $x$  AND  $y$ ) is related to the degree of belief in the conditional proposition  $x|y$  and the degree of belief in the proposition  $y$ . There is a function  $g$  such that

$$B(x, y) = g[B(x|y), B(y)].$$

#### Box 2.4. The Cox axioms.

If a set of beliefs satisfy these axioms then they can be mapped onto probabilities satisfying  $P(\text{FALSE}) = 0$ ,  $P(\text{TRUE}) = 1$ ,  $0 \leq P(x) \leq 1$ , and the rules of probability:

$$P(x) = 1 - P(\bar{x}),$$

and

$$P(x, y) = P(x|y)P(y).$$

# Maximum likelihood estimation

Given a data set

$$\text{Data} = \{\boldsymbol{x}^1, \dots, \boldsymbol{x}^N\}$$

and a parametrized distribution

$$p(\boldsymbol{x}|\boldsymbol{w}), \quad \boldsymbol{w} = (w_1, \dots, w_M),$$

find the value of  $\boldsymbol{w}$  that best describes the data.

The common approach is to assume that the data that we observe are drawn independently from  $p(\boldsymbol{x}|\boldsymbol{w})$  (independent and identical distributed = i.i.d.) for some unknown value of  $\boldsymbol{w}$ :

$$p(\text{Data}|\boldsymbol{w}) = p(\{\boldsymbol{x}^1, \dots, \boldsymbol{x}^N\}|\boldsymbol{w}) = \prod_{i=1}^N p(\boldsymbol{x}^i|\boldsymbol{w})$$

Then, the most likely  $\mathbf{w}$  is obtained by maximizing  $p(\text{Data}|\mathbf{w})$  wrt  $\mathbf{w}$ :

$$\begin{aligned}\mathbf{w}_{\text{ML}} &= \operatorname{argmax}_{\mathbf{w}} p(\text{Data}|\mathbf{w}) = \operatorname{argmax}_{\mathbf{w}} \prod_{i=1}^N p(\mathbf{x}^i|\mathbf{w}) \\ &= \operatorname{argmax}_{\mathbf{w}} \sum_{i=1}^N \log p(\mathbf{x}^i|\mathbf{w})\end{aligned}$$

since  $\log$  is a monotonically increasing function.

$\mathbf{w}_{\text{ML}}$  is a function of the data. This is called an *estimator*.

Frequentists methods consider a single true  $w$  and data generation mechanism  $p(\text{Data}|w)$  provided by 'Nature' and study expected value:

$$\mathbb{E}\mathbf{w}_{\text{ML}} = \sum_{\text{Data}} p(\text{Data}|w) \mathbf{w}_{\text{ML}}(\text{Data})$$

For instance,  $\hat{\mu} = \frac{1}{N} \sum_i x_i$  is estimator for the mean of a distribution. If data are  $x_i \sim \mathcal{N}(\mu, \sigma^2)$  then  $\mathbb{E}\hat{\mu} = \mu$ .

# Bayesian machine learning

Model parameters  $\mathbf{w}$ : uncertain

- Prior assumptions and beliefs about model parameters: the *prior* distribution  $p(\mathbf{w})$
- Observed data  $= \{\mathbf{x}^1, \dots, \mathbf{x}^N\} = \text{Data}$
- Probability of data given  $\mathbf{w}$  (the *likelihood*):  $p(\text{Data}|\mathbf{w})$

Apply Bayes' rule to obtain the *posterior* distribution

$$p(\mathbf{w}|\text{Data}) = \frac{p(\text{Data}|\mathbf{w})p(\mathbf{w})}{p(\text{Data})} \propto p(\text{Data}|\mathbf{w})p(\mathbf{w})$$

$p(\mathbf{w})$  : prior

$p(\text{Data}|\mathbf{w})$  : likelihood

$p(\text{Data}) = \int p(\text{Data}|\mathbf{w})p(\mathbf{w}) d\mathbf{w}$  : evidence

$\rightarrow p(\mathbf{w}|\text{Data})$  : posterior

## Predictive distribution

Prior to 'learning', the predictive distribution for new observation  $x$  is

$$p(x) = \int p(x|\mathbf{w})p(\mathbf{w}) d\mathbf{w}$$

After 'learning', i.e., after observation of Data, the predictive distribution for new observation  $x$  becomes

$$\begin{aligned} p(x|\text{Data}) &= \int p(x|\mathbf{w}, \text{Data})p(\mathbf{w}|\text{Data}) d\mathbf{w} \\ &= \int p(x|\mathbf{w})p(\mathbf{w}|\text{Data}) d\mathbf{w} \end{aligned}$$

## Bayesian vs frequentists view point

- Bayesians: there is a single fixed dataset (the one observed), and a probability distribution of model parameters  $w$  which expresses a subjective belief including uncertainty about the ‘true’ model parameters.
  - They need a prior belief  $p(w)$ , and apply Bayes rule to compute  $p(w|\text{Data})$ .
  - Bayesians can talk about the belief that  $w$  has a certain value given the particular data set.
- Frequentists assume a single (unknown) fixed model parameter vector  $w$ .
  - construct an estimator  $\hat{w}$  that is a function of the data. For instance, the maximum likelihood estimator
  - study statistical properties of estimators in similar experiments, each time with different datasets drawn from  $p(\text{Data}|w)$ , such as bias and variance.
  - They cannot make a claim for this particular data set. This is the price for not having a ‘prior’.

## Toy example

$w$  is the probability that a coin comes up 'head'. Toss  $N$  times with  $N_H$  outcomes 'head'.

The likelihood of the data  $p(N_H|w, N) = \binom{N}{N_H} w^{N_H} (1-w)^{N-N_H}$ .

The (frequentist) maximum likelihood estimator:

$$\hat{w} = \operatorname{argmax}_w p(N_H|w, N) = \operatorname{argmax}_w [N_H \log w + (N - N_H) \log(1 - w)] = \frac{N_H}{N}$$

$\hat{w}$  is stochastic variable, because it depends on the data set.

But it has 'nice' statistical property that on average (over many data sets)  $\hat{w}$  gives the correct value:

$$\mathbb{E}\hat{w} = \sum_{N_H} p(N_H|w, N) \frac{N_H}{N} = w$$

3

---

<sup>3</sup>Use  $\sum_{N_H=0}^N p(N_H|w, N) = 1$ .

The Bayesian approach considers one data set and assumes a prior  $p(w)$  and compute the posterior

$$p(w|N_H, N) = \frac{p(w)p(N_H|w, N)}{p(N_H, N)}$$

## Bayesian vs frequentists

- Prior: inclusion of prior knowledge may be useful. True reflection of knowledge, or convenient construct? Bad prior choice can overconfidently lead to poor result.
- Bayesian integrals cannot be calculated in general. Only approximate results possible, requiring intensive numerical computations.
- Frequentists methods of ‘resampling the data’, (crossvalidation, bootstrapping) are appealing
- Bayesian framework transparent and consistent. Assumptions are explicit, inference is a mechanistic procedure (Bayesian machinery) and results have a clear interpretation.

This course place emphasis on Bayesian approach.

## Medical example

Suppose  $w = \{0, 1\}$  is a disease state (absent/present). The disease is rare, say  $P(w = 1) = 0.01$ . There is a test  $x = 0, 1$  that measures whether the patient has the disease.

$$p(x = 1|w = 1) = p(x = 0|w = 0) = 0.9$$

The test is performed and is positive:  $x = 1$ . What is the probability that the patient has the disease?

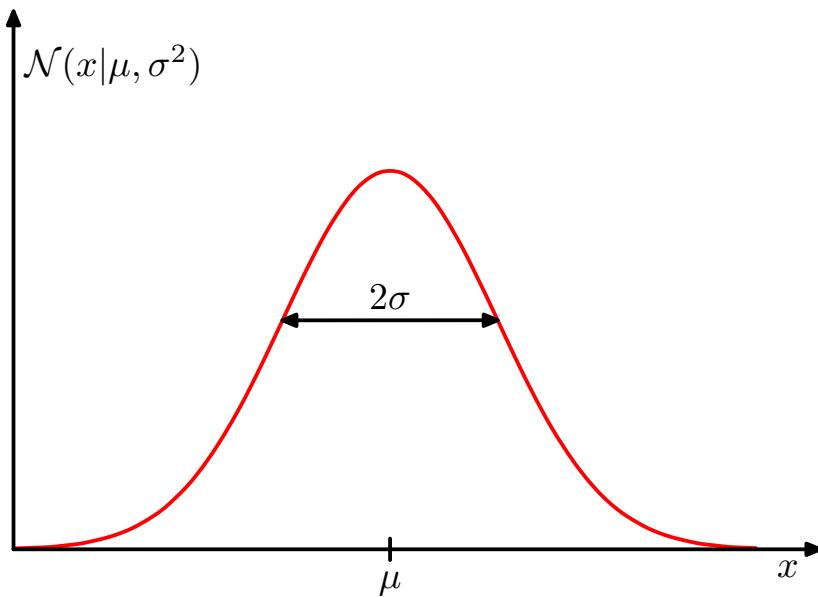
$$p(w=1|x=1) = \frac{0.9*0.01}{0.9*0.01+0.1*0.99} = \frac{1}{1+\frac{0.1*0.99}{0.9*0.01}} = \frac{1}{12} = 0.0825$$

## Gaussian distribution

Normal distribution = Gaussian distribution

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{(x-\mu)^2}{2\sigma^2} \right\}$$

Specified by  $\mu$  and  $\sigma^2$



Gaussian is normalized,

$$\int_{-\infty}^{\infty} \mathcal{N}(x|\mu, \sigma^2) dx = 1$$

The mean (= first moment), second moment, and variance are:

$$\mathbb{E}[x] = \langle x \rangle = \int_{-\infty}^{\infty} x \mathcal{N}(x|\mu, \sigma^2) dx = \mu$$

$$\langle x^2 \rangle = \int_{-\infty}^{\infty} x^2 \mathcal{N}(x|\mu, \sigma^2) dx = \mu^2 + \sigma^2$$

$$\text{var}[x] = \langle x^2 \rangle - \langle x \rangle^2 = \sigma^2$$

## Multivariate Gaussian

In  $D$  dimensions

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{D/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right)$$

$\mathbf{x}, \boldsymbol{\mu}$  are  $D$ -dimensional vectors.

$\boldsymbol{\Sigma}$  is a  $D \times D$  covariance matrix,  $|\boldsymbol{\Sigma}|$  is its determinant.

## Mean vector and covariance matrix

$$\mathbb{E}[\mathbf{x}] = \langle \mathbf{x} \rangle = \int \mathbf{x} \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{x} = \boldsymbol{\mu}$$

$$\text{cov}[\mathbf{x}] = \langle (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T \rangle = \int (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{x} = \boldsymbol{\Sigma}$$

We can also write this in component notation:

$$\mu_i = \langle x_i \rangle = \int x_i \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{x}$$

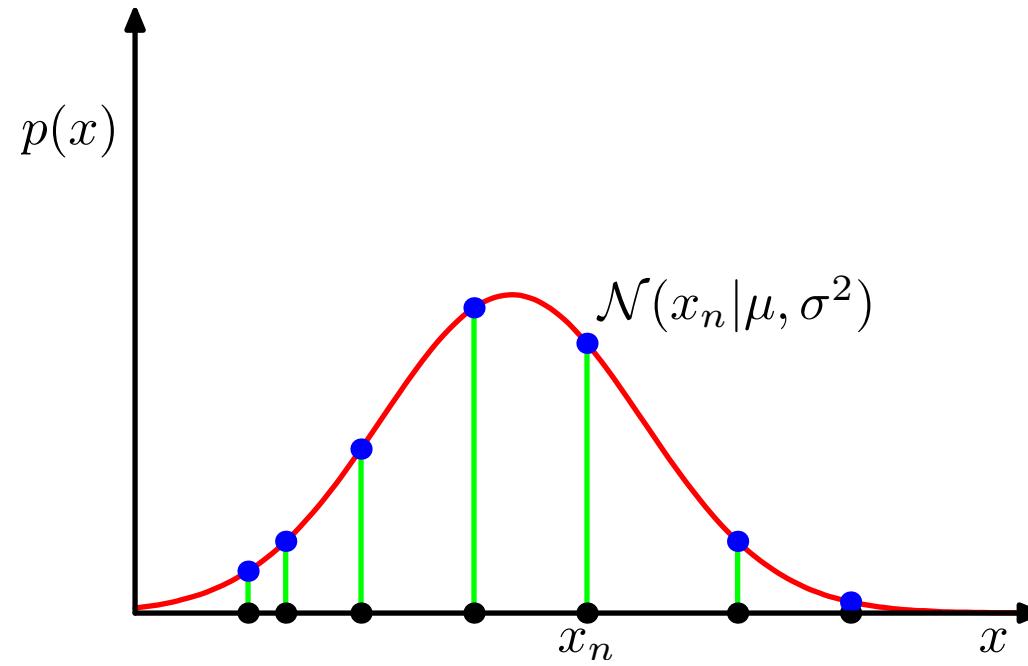
$$\Sigma_{ij} = \langle (x_i - \mu_i)(x_j - \mu_j) \rangle = \int (x_i - \mu_i)(x_j - \mu_j) \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{x}$$

$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$  is specified by its mean and covariance, in total  $D(D+1)/2 + D$  parameters.

## The likelihood for the 1-d Gaussian

Consider 1-d data  $\text{Data} = \mathbf{x} = \{x_1, \dots, x_N\}$ . The likelihood of the data under a Gaussian model is the probability of the data, assuming each data point is independently drawn from the Gaussian distribution:

$$p(\mathbf{x}|\mu, \sigma) = \prod_{n=1}^N \mathcal{N}(x_n|\mu, \sigma^2) = \left( \frac{1}{\sqrt{2\pi}\sigma} \right)^N \exp \left( -\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 \right)$$



## Maximum likelihood

Consider the log of the likelihood:

$$\begin{aligned}\ln p(\mathbf{x}|\mu, \sigma) &= \ln \left( \left( \frac{1}{\sqrt{2\pi}\sigma} \right)^N \exp \left( -\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 \right) \right) \\ &= -\frac{1}{2\sigma^2} \sum_{n=1}^N (x_n - \mu)^2 - \frac{N}{2} \ln \sigma^2 - \frac{N}{2} \ln 2\pi\end{aligned}$$

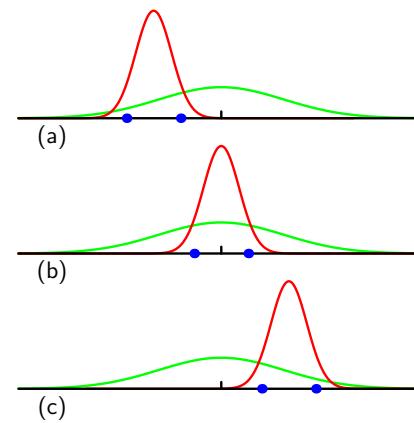
The values of  $\mu$  and  $\sigma$  that maximize the likelihood are given by

$$\mu_{ML} = \frac{1}{N} \sum_{n=1}^N x_n \quad \sigma_{ML}^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \mu_{ML})^2$$

## Bias in the ML estimates

Note that  $\mu_{ML}, \sigma_{ML}^2$  are functions of the data. We can take their expectation value, assuming that  $x_n$  is from a  $\mathcal{N}(x|\mu, \sigma)$ .

$$\langle \mu_{ML} \rangle = \frac{1}{N} \sum_{n=1}^N \langle x_n \rangle = \mu \quad \langle \sigma_{ML}^2 \rangle = \frac{1}{N} \sum_{n=1}^N \langle (x_n - \mu_{ML})^2 \rangle = \dots = \frac{N-1}{N} \sigma^2$$



The variance is estimated too low. This is called a biased estimator. Bias disappears when  $N \rightarrow \infty$ . In complex models with many parameters, the bias is more severe.

(Bayesian approach gives correct expected values)

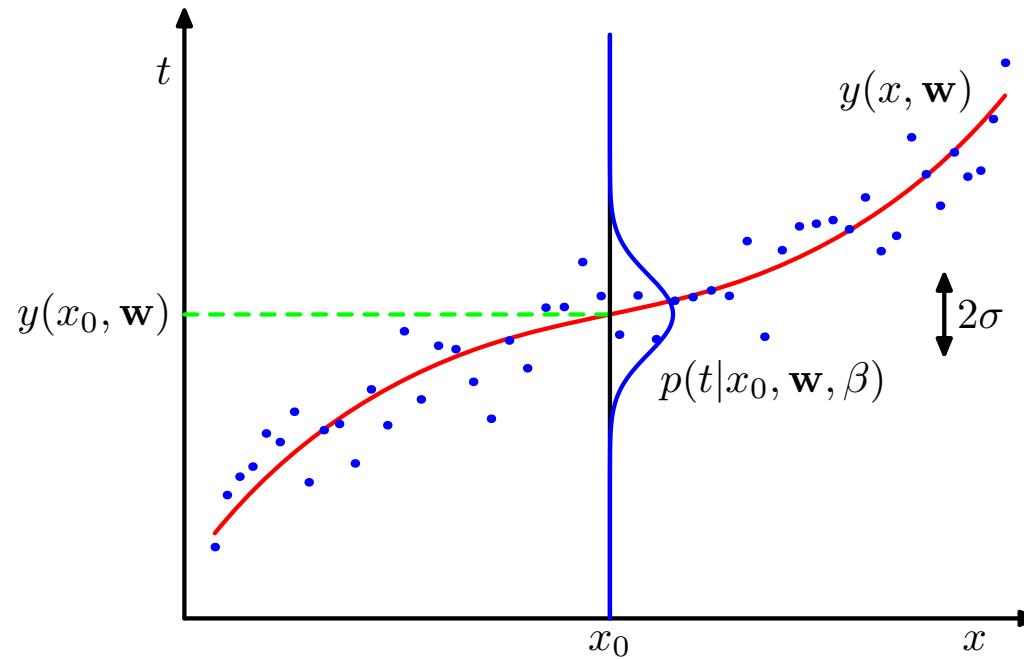
## Curve fitting re-visited

Now from a probabilistic perspective.

Target  $t$  is Gaussian distributed around mean  $y(x, \mathbf{w}) = \sum_{j=0}^M w_j x^j$ ,

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

$\beta = 1/\sigma^2$  is the precision.



## Curve fitting re-visited: ML

Training data: inputs  $\mathbf{x} = (x_1, \dots, x_n)$ , targets  $\mathbf{t} = (t_1, \dots, t_n)$ . (Assume  $\beta$  is known.)

Likelihood function,

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

Log-likelihood

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

With  $\mathbf{w}_{ML}$  one can make predictions for a new input values  $x$ . The predictive distribution over the output  $t$  is:

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

## Curve fitting re-visited MAP

More Bayesian approach.

Prior:

$$p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}|\alpha^{-1}\mathbf{I}) = \left(\frac{\alpha}{2\pi}\right)^{(M+1)/2} \exp\left(-\frac{\alpha}{2}\mathbf{w}^T\mathbf{w}\right)$$

$M$  is the dimension of  $\mathbf{w}$ . Variables such as  $\alpha$ , controlling the distribution of parameters, are called ‘hyperparameters’.

Posterior using Bayes rule:

$$\begin{aligned} p(\mathbf{w}|\mathbf{t}, \mathbf{x}, \alpha, \beta) &\propto p(\mathbf{w}|\alpha) \prod_{n=1}^N \mathcal{N}(t_n|y(x_n, \mathbf{w}), \beta^{-1}) \\ -\ln p(\mathbf{w}|\mathbf{x}, \mathbf{t}, \alpha, \beta) &= \frac{\beta}{2} \sum_{n=1}^N (y(x_n, \mathbf{w}) - t_n)^2 + \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} + \text{const} \end{aligned}$$

Maximizing the posterior wrt  $\mathbf{w}$  yields  $\mathbf{w}_{MAP}$ . Similar as Eq. 1.4. with  $\lambda = \alpha/\beta$

## Bayesian curve fitting

Given the training data  $\mathbf{x}, \mathbf{t}$  we are not so much interested in  $\mathbf{w}$ , but rather in the prediction of  $t$  for a new  $x$ :  $p(t|x, \mathbf{x}, \mathbf{t})$ . This is given by

$$p(t|x, \mathbf{x}, \mathbf{t}) = \int p(t|x, \mathbf{w})p(\mathbf{w}|\mathbf{x}, \mathbf{t})d\mathbf{w}$$

It is the average prediction of an ensemble of models  $p(t|x, \mathbf{w})$  parametrized by  $\mathbf{w}$  and averaged wrt to the posterior distribution  $p(\mathbf{w}|\mathbf{x}, \mathbf{t})$ .

All quantities depend on  $\alpha$  and  $\beta$ . (How?)

## Bayesian curve fitting

Generalized linear model with ‘basis functions’ e.g.,  $\phi_i(x) = x^i$ ,

$$y(x, \mathbf{w}) = \sum_i \phi_i(x) w_i = \boldsymbol{\phi}(x)^T \mathbf{w}$$

So: assuming Gaussian noise, prediction given  $\mathbf{w}$  is

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

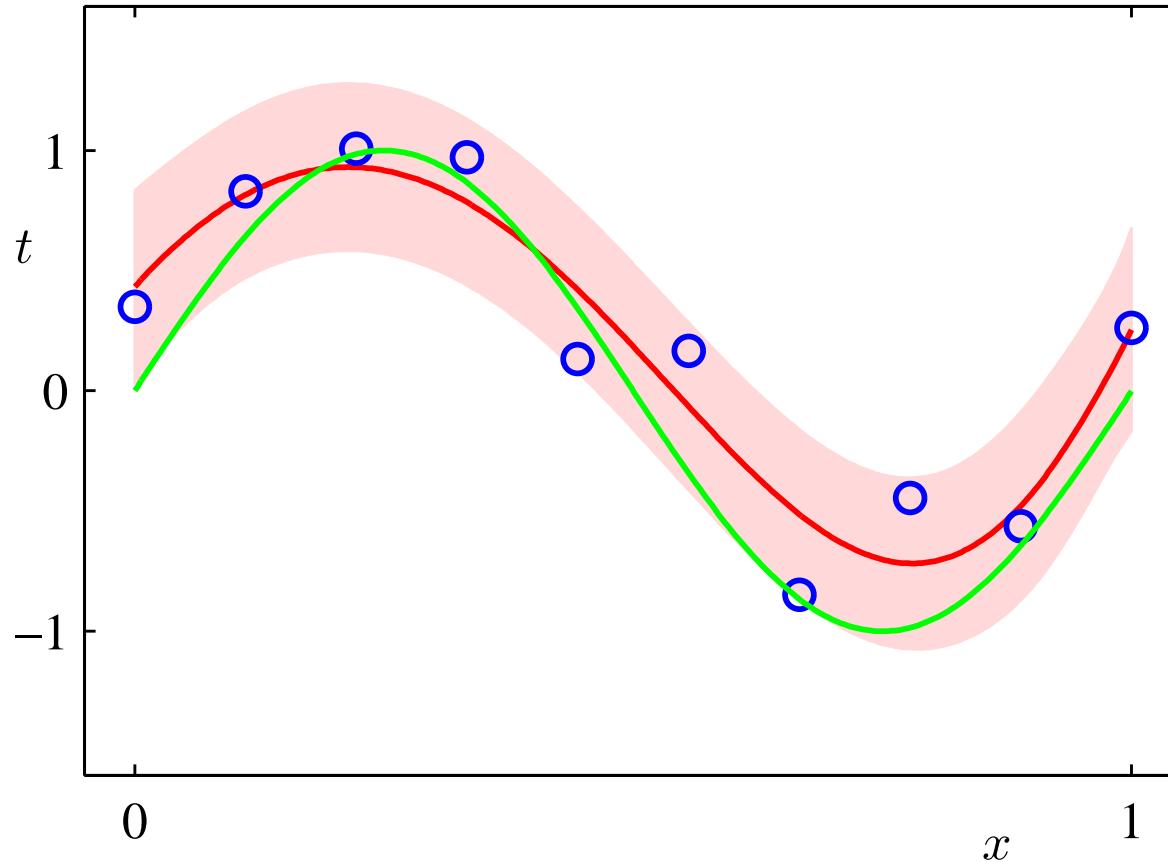
# Result Bayesian curve fitting

Predictive distribution

$$\begin{aligned}
 p(t|x, \mathbf{x}, \mathbf{t}) &= \mathcal{N}(t|m(x), s^2(x)) \\
 m(x) &= \beta \boldsymbol{\phi}(x)^T \mathbf{S} \sum_{n=1}^N \boldsymbol{\phi}(x_n) t_n = \boldsymbol{\phi}(x)^T \mathbf{w}_{MAP} \\
 s^2(x) &= \beta^{-1} + \boldsymbol{\phi}(x)^T \mathbf{S} \boldsymbol{\phi}(x) \\
 \mathbf{S}^{-1} &= \alpha \mathbf{I} + \beta \sum_{n=1}^N \boldsymbol{\phi}(x_n) \boldsymbol{\phi}(x_n)^T
 \end{aligned}$$

Note,  $s^2$  depend on  $x$ . First term as in ML estimate describes noise in target for fixed  $\mathbf{w}$ . Second term describes noise due to uncertainty in  $\mathbf{w}$ .

# Example



Polynomial curve fitting with  $M = 9, \alpha = 5 \times 10^{-3}, \beta = 11.1$ . Red:  $m(x) \pm s(x)$ . Note  $\sqrt{\beta^{-1}} = 0.3$

## Model selection

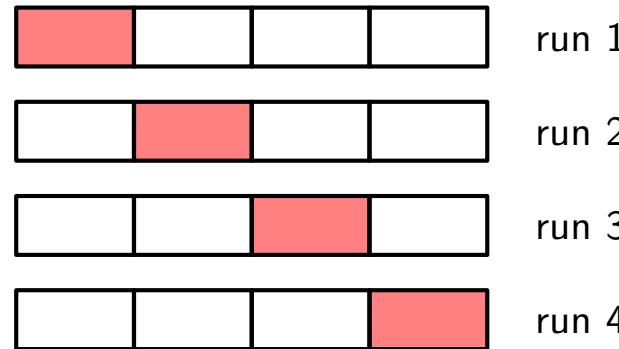
Q: If we have different models to describe the data, which one should we choose?

# Model selection/Cross validation

Q: If we have different models to describe the data, which one should we choose?

A1: If data is plenty, use separate *validation set* to select model with best generalization performance, and a third independent *test set* for final evaluation.

A2: Small validation set: use *S-fold cross validation*.

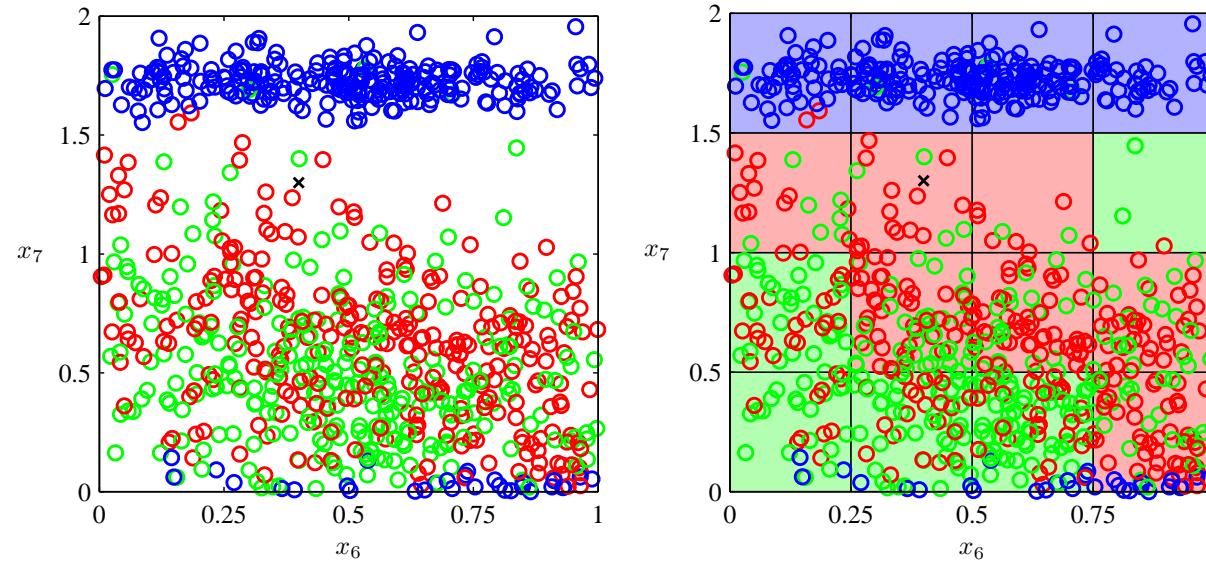


A3: Information criteria: penalty for complex models

- Akaike IC (AIC):  $\ln p(D|\boldsymbol{w}_{ML}) - M$
- Bayesian IC (BIC):  $\ln p(D|\boldsymbol{w}_{MAP}) - \frac{1}{2}M \ln N$  (Bayesian + crude approximation)
- Full Bayesian  $\rightarrow$  penalties arises automatically

# High-dimensional data/Binning

Sofar, we considered  $x$  one-dimensional. How does pattern recognition work higher dimensions?



Two components of 12-dimensional data that describe gamma ray measurements of a mixture of oil, water and gas. The mixture can be in three states: homogenous (red), annular (green) and laminar (blue).

Classification of the 'x' can be done by making a putting a grid on the space and assigning the class that is most numerous in the particular box.

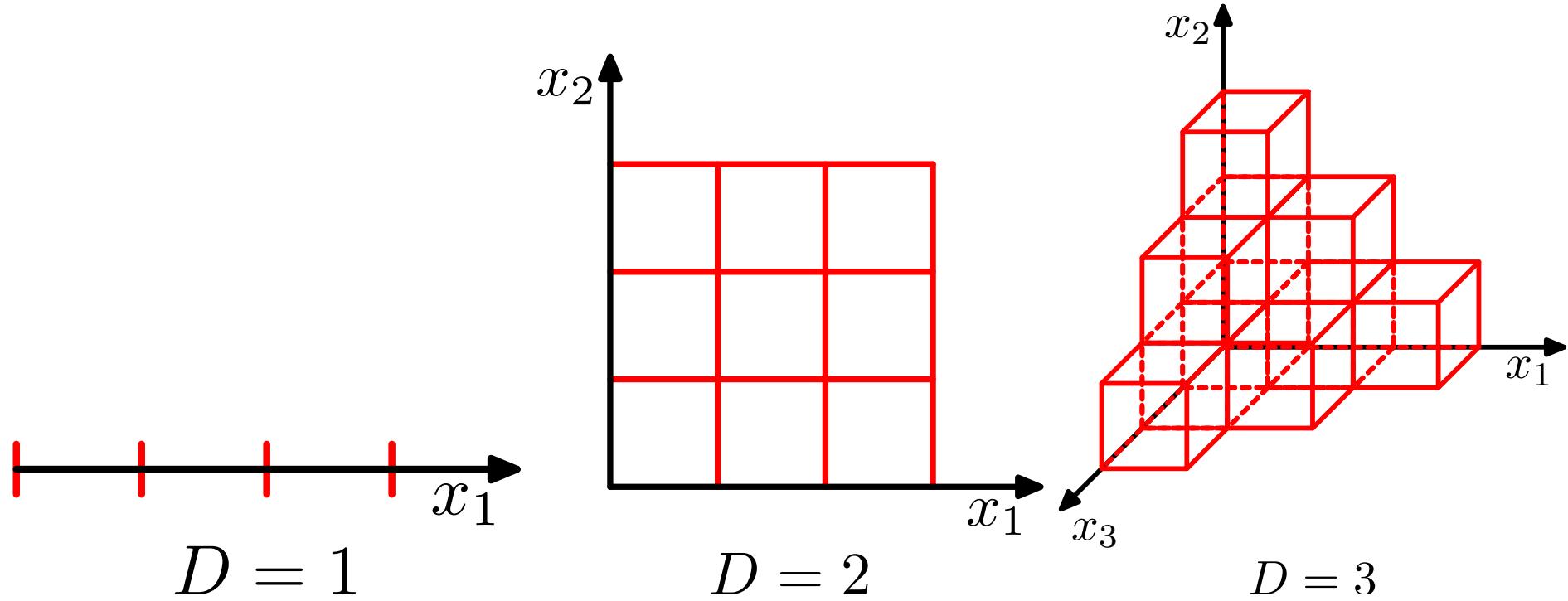
# High-dimensional data/Binning

Q: What is the disadvantage of this approach?

# Curse of dimensionality/Binning

Q: What is the disadvantage of this approach?

A: This approach scales exponentially with dimensions.



In  $D$  dimensions: grid with length  $n$  consists of  $n^D$  hypercubes.

## Curse of dimensionality/Polynomials

The polynomial function considered previously becomes in  $D$  dimensions:

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{i=1}^D w_i x_i + \sum_{i=1}^D \sum_{j=1}^D w_{ij} x_i x_j + \sum_{i=1}^D \sum_{j=1}^D \sum_{k=1}^D w_{ijk} x_i x_j x_k$$

(here up to order  $M = 3$ ).

The number of coefficients scales as ... ?

## Curse of dimensionality/Polynomials

The polynomial function considered previously becomes in  $D$  dimensions:

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{i=1}^D w_i x_i + \sum_{i=1}^D \sum_{j=1}^D w_{ij} x_i x_j + \sum_{i=1}^D \sum_{j=1}^D \sum_{k=1}^D w_{ijk} x_i x_j x_k$$

(here up to order  $M = 3$ ).

The number of coefficients scales as  $D^M$  (unpractically large).

## Curse of dimensionality/Spheres

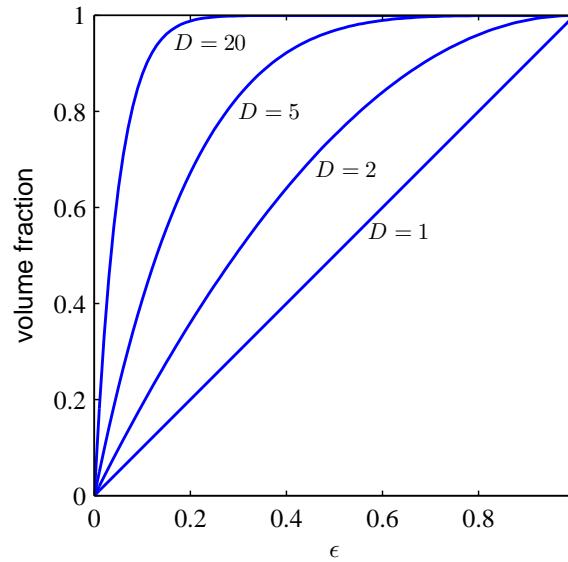
Q: In a 10-dimensional hypersphere with radius  $R = 1$ : what is the volume fraction lying in the outer “shell” between  $r = 0.5$  and  $r = 1$ ?

## Curse of dimensionality/Spheres

Q: In a 10-dimensional hypersphere with radius  $R = 1$ : what is the volume fraction lying in the outer “shell” between  $r = 0.5$  and  $r = 1$ ?

A: More than 0.999!

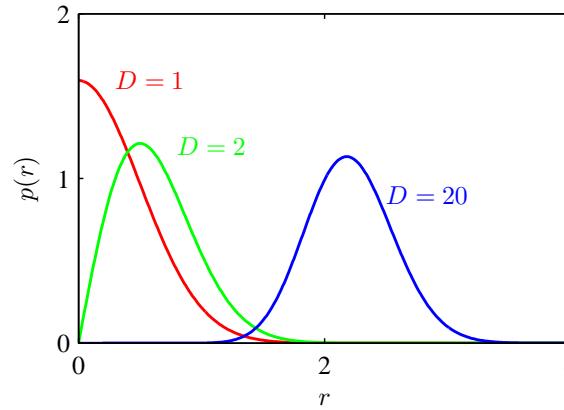
So in high dimensions, almost all data points are at more or less the same distance!



$$V_D(r) = K_D r^D \quad \frac{V_D(1) - V_D(1 - \epsilon)}{V_D(1)} = 1 - (1 - \epsilon)^D$$

Spheres in high dimension have most of their volume on the boundary.

# Curse of dimensionality/Gaussians



In high dimensions, the distribution of the radius of a Gaussian with variance  $\sigma$  is concentrated around a thin shell  $r \approx \sigma\sqrt{D}$ .

*Intuition developed in low dimensional space may be wrong in high dimensions!*

## Curse of dimensionality

Is machine learning even possible in high dimensions?

- Data often in low dimensional subspace: only a few dimensions are relevant.
  - Object located in 3-D  $\rightarrow$  images of objects are  $N$ -D  $\rightarrow$  there should be 3-D manifold (curved subspace)
- Smoothness, local interpolation (note: this is also needed in low dimensions).

## Decision theory

**Inference:** given pairs  $(x, t)$ , learn  $p(x, t)$  and estimate  $p(x, t)$  for new value of  $x$  (and possibly all  $t$ ).

**Decision:** for new value of  $x$  estimate 'best'  $t$ .

**Example:** in a medical application,  $x$  is an X-ray image and  $t$  a class label that indicates whether the patient has cancer ( $t = C_1$ ) or not ( $t = C_2$ ).

# Decision theory

**Inference:** given pairs  $(\mathbf{x}, t)$ , learn  $p(\mathbf{x}, t)$  and estimate  $p(\mathbf{x}, t)$  for new value of  $\mathbf{x}$  (and possibly all  $t$ ).

**Decision:** for new value of  $\mathbf{x}$  estimate 'best'  $t$ .

**Example:** in a medical application,  $\mathbf{x}$  is an X-ray image and  $t$  a class label that indicates whether the patient has cancer ( $t = C_1$ ) or not ( $t = C_2$ ).

Bayes' theorem:

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

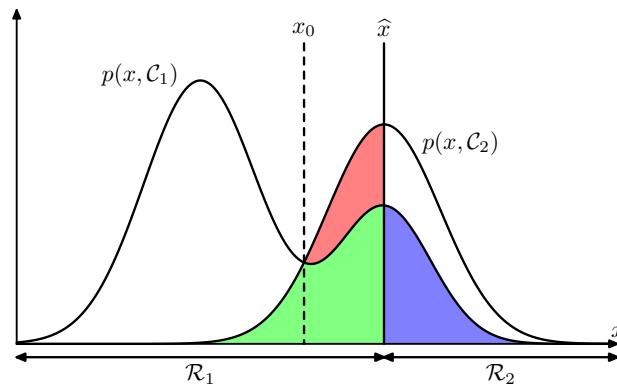
$p(C_k)$  is the prior probability of class  $C_k$ .  $p(C_k | \mathbf{x})$  is the posterior probability of class  $C_k$  after seeing the image  $\mathbf{x}$ .

Q: Suppose  $p(C_1) = 0.01$  and  $p(C_1 | \mathbf{x}) = 0.3$  according to our model. Do we decide that the patient has cancer and therefore start treatment?

## Decision theory

A *classifier* is specified by defining regions  $R_k$ , such that all  $x \in R_k$  are assigned to class  $C_k$ . In the case of two classes, the probability that this classifier gives the correct answer is

$$p(\text{correct}) = p(x \in R_1, C_1) + p(x \in R_2, C_2) = \int_{R_1} p(x, C_1) dx + \int_{R_2} p(x, C_2) dx$$



$p(\text{correct})$  is maximized when the regions  $R_k$  are chosen such that

$$k = \operatorname{argmax}_k p(x, C_k) = \operatorname{argmax}_k p(C_k | x)$$

## Decision theory/Example

**Example:** in a medical application,  $x$  is an X-ray image and  $t$  a class label that indicates whether the patient has cancer ( $t = C_1$ ) or not ( $t = C_2$ ).

Q: Suppose  $p(C_1) = 0.01$  and  $p(C_1|x) = 0.3$  according to our model. Do we decide that the patient has cancer and therefore start treatment?

## Decision theory/Example

**Example:** in a medical application,  $\mathbf{x}$  is an X-ray image and  $t$  a class label that indicates whether the patient has cancer ( $t = C_1$ ) or not ( $t = C_2$ ).

Q: Suppose  $p(C_1) = 0.01$  and  $p(C_1|\mathbf{x}) = 0.3$  according to our model. Do we decide that the patient has cancer and therefore start treatment?

A: If we want to maximize the chance of making the correct decision, we have to pick  $k$  such that  $p(C_k|\mathbf{x})$  is maximal.

Because  $p(C_1|\mathbf{x}) = 0.3$  and  $p(C_2|\mathbf{x}) = 0.7$ , the answer is *no*: we decide that the patient does not have cancer.

## Decision theory/Expected loss

Typically, not all classification errors are equally bad: classifying a healthy patient as sick, is not as bad as classifying a sick patient as healthy.

$$L = \begin{pmatrix} 0 & 1000 \\ 1 & 0 \end{pmatrix}$$

Loss function. Rows are true classes (cancer, normal), columns are assigned classes (cancer, normal).

The probability to assign an  $\mathbf{x}$  to class  $j$  while it belongs to class  $k$  is  $p(\mathbf{x} \in R_j, C_k)$ . Thus the total *expected loss* is

$$\langle L \rangle = \sum_k \sum_j L_{kj} p(\mathbf{x} \in R_j, C_k) = \sum_k \sum_j \int_{R_j} p(\mathbf{x}, C_k) L_{jk} d\mathbf{x}$$

$\langle L \rangle$  is minimized if each  $\mathbf{x}$  is assigned to class  $j$  such that  $\sum_k L_{kj} p(\mathbf{x}, C_k)$  is minimal, or equivalently, such that  $\sum_k L_{kj} p(C_k | \mathbf{x})$  is minimal.

## Decision theory/Example

**Example:** in a medical application,  $\mathbf{x}$  is an X-ray image and  $t$  a class label that indicates whether the patient has cancer ( $t = C_1$ ) or not ( $t = C_2$ ).

	$C_1$	$C_2$
$C_1$	0	1000
$C_2$	1	0

Loss function. Rows are true classes (cancer, normal), columns are assigned classes (cancer, normal).

Q: Suppose  $p(C_1) = 0.01$  and  $p(C_1|\mathbf{x}) = 0.3$  according to our model. Do we decide that the patient has cancer and therefore start treatment?

## Decision theory/Example

**Example:** in a medical application,  $\mathbf{x}$  is an X-ray image and  $t$  a class label that indicates whether the patient has cancer ( $t = C_1$ ) or not ( $t = C_2$ ).

	$C_1$	$C_2$
$C_1$	0	1000
$C_2$	1	0

Loss function. Rows are true classes (cancer, normal), columns are assigned classes (cancer, normal).

Q: Suppose  $p(C_1) = 0.01$  and  $p(C_1|\mathbf{x}) = 0.3$  according to our model. Do we decide that the patient has cancer and therefore start treatment?

A: If we want to minimize the expected loss, we have to pick  $j$  such that  $\sum_k L_{kj} p(C_k|\mathbf{x})$  is minimal.

For  $j = 1$ , this yields  $0 \times 0.3 + 1 \times 0.7$ ,

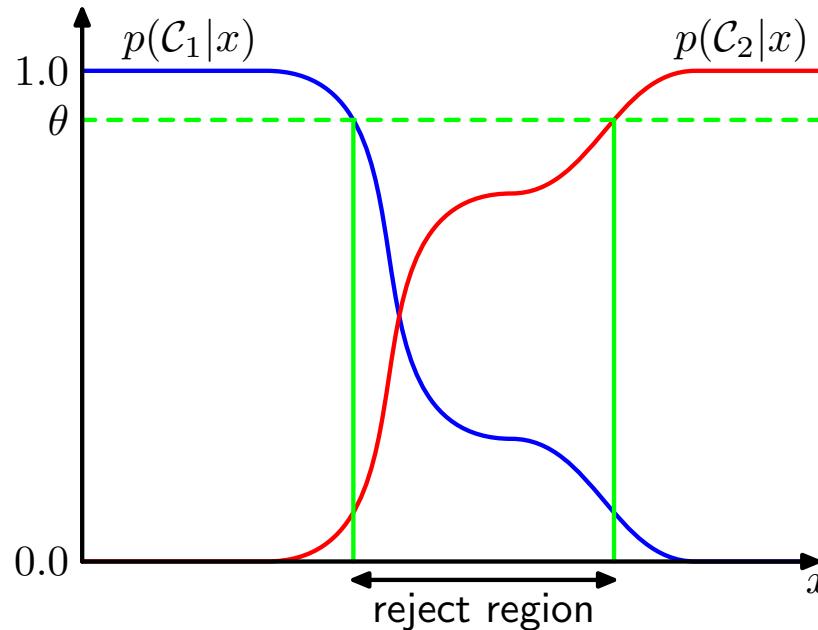
for  $j = 2$ , this yields  $1000 \times 0.3 + 0 \times 0.7$ .

Therefore, we now decide that the patient has cancer (better safe than sorry).

## Decision theory/Reject option

It may be that  $\max_j p(C_j|x)$  is small, indicating that it is unclear to which class  $x$  belongs.

In that case, a different decision can be taken: the “reject” or “don’t know” option.



This can be done by introducing a threshold  $\theta \in [0, 1]$  and only classify those  $x$  for which  $\max_j p(C_j|x) > \theta$  (and answer “don’t know” otherwise).

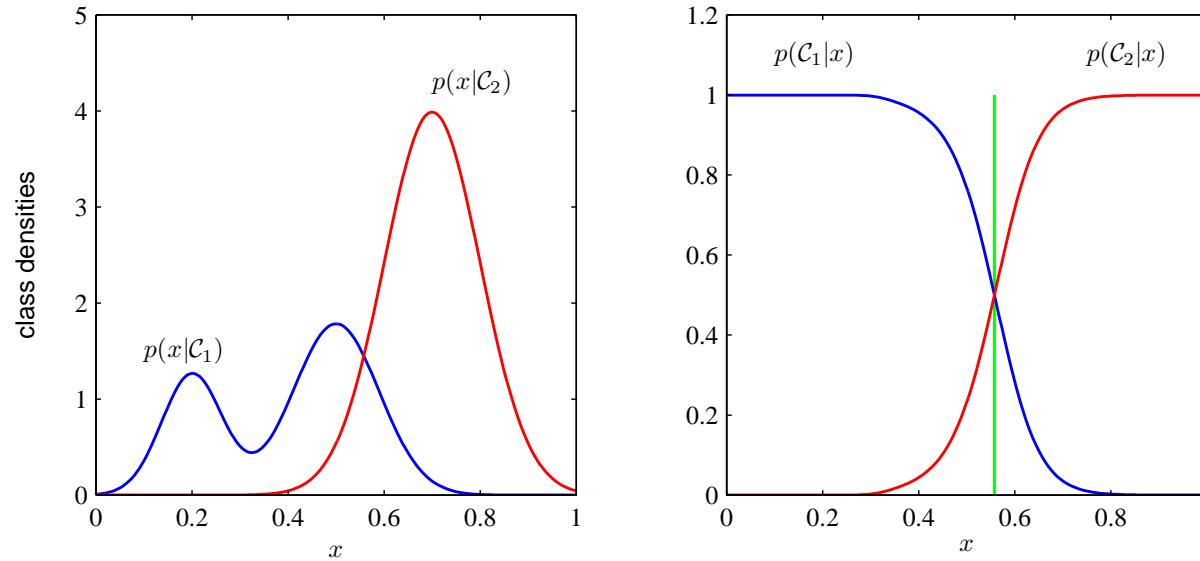
# Decision theory/Discriminant functions

Instead of first learning a probability model and then making a decision, one can also directly learn a decision rule (a classifier) without the intermediate step of a probability model.

A set of approaches:

- Learn a model for the class conditional probabilities  $p(\mathbf{x}|C_k)$ . Use Bayes' rule to compute  $p(C_k|\mathbf{x})$  and construct a classifier using decision theory. This approach is the most complex, but has the advantage of yielding a model of  $p(\mathbf{x})$  that can be used to reject unlikely inputs  $\mathbf{x}$ .
- Learn the inference problem  $p(C_k|\mathbf{x})$  directly and construct a classifier using decision theory. This approach is simpler, since no input model  $p(\mathbf{x})$  is learned (see figure).
- Learn  $f(\mathbf{x})$ , called a discriminant function, that maps  $\mathbf{x}$  directly onto a class label  $0, 1, 2, \dots$ . Even simpler, only decision boundary is learned. But information on the expected classification error is lost.

# Decision theory/Discriminant functions



Example that shows that detailed structure in the joint model need not affect class conditional probabilities. Learning only the decision boundary is the simplest approach.

Approaches that model the distribution of both inputs and outputs are called *generative models*, approaches that only model the conditional distribution of the output given the input are called *discriminative models*.

## Decision theory/Discriminant functions

Advantages of learning a class conditional probability instead of discriminant function:

**Minimizing risk** When minimizing expected loss, the loss matrix may change over time whereas the class probabilities may not (for instance, in a financial application).

**Reject option** One can reject uncertain class assignments

**Unbalanced data** One can compensate for unbalanced data sets. For instance, in the cancer example, there may be 1000 times more healthy patients than cancer patients. Very good classification (99.9 % correct) is obtained by classifying everyone as healthy. Using the posterior probability one can compute  $p(C_k = \text{cancer}|\mathbf{x})$ . Although this probability may be low, it may be significantly higher than  $p(C_k = \text{cancer})$ , indicating a risk of cancer.

## Decision theory/Discriminant functions

**Combining models** Given models for  $p(C_k|\mathbf{x})$  and  $p(C_k|\mathbf{y})$  one has a principled approach to classify on both  $\mathbf{x}$  and  $\mathbf{y}$ . Naive Bayes assumption:

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

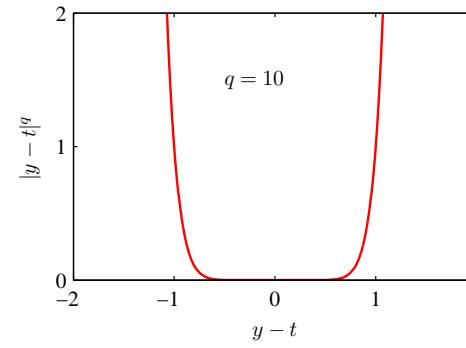
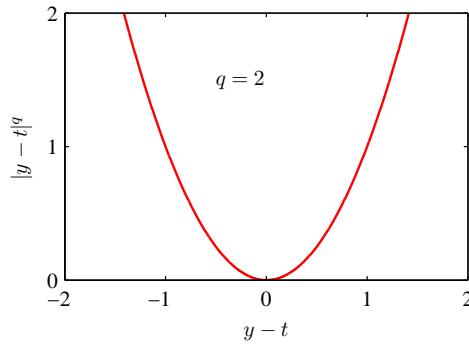
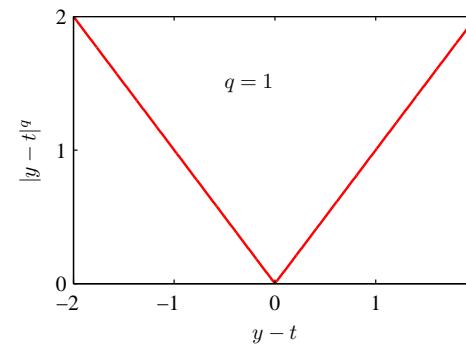
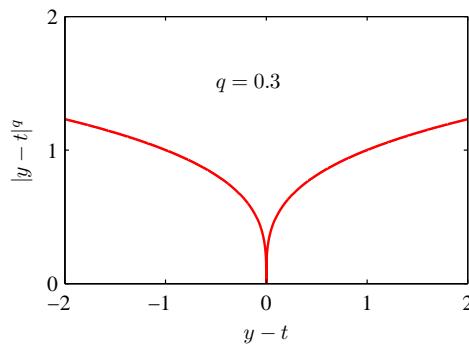
(given the disease state of the patient the blood and X-ray test results are independent). Then

$$p(C_k|\mathbf{x}, \mathbf{y}) \propto p(\mathbf{x}, \mathbf{y}|C_k)p(C_k) \propto p(\mathbf{x}|C_k)p(\mathbf{y}|C_k)p(C_k) \propto \frac{p(C_k|\mathbf{x})p(C_k|\mathbf{y})}{p(C_k)}$$

# Loss functions for regression

Decision theory generalizes straightforwardly to continuous variables: the loss matrix  $L_{jk}$  becomes a loss function  $L(t, y(\mathbf{x}))$ .

Examples:



Minkowski loss function  $L_q = |y - t|^q$  for various values of  $q$ .

# Loss functions for regression/Quadratic loss

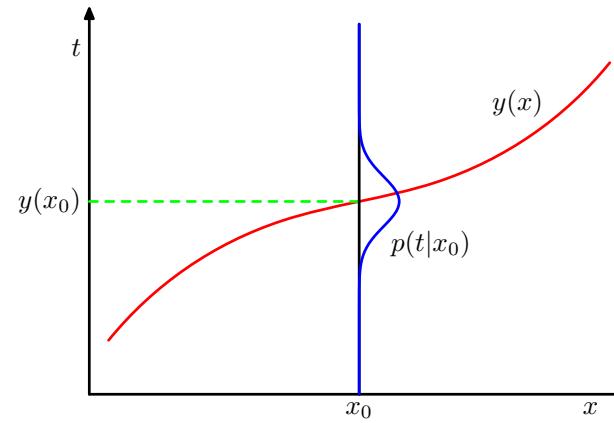
The average/expected loss is:

$$\langle L \rangle = \int L(t, y(\mathbf{x})) p(\mathbf{x}, t) d\mathbf{x} dt$$

For the quadratic loss function  $L_2(t, y(\mathbf{x})) = (t - y(\mathbf{x}))^2$  one can derive that the expected loss is minimized by taking

$$y(\mathbf{x}) = \mathbb{E}_t[t|\mathbf{x}]$$

i.e., by the mean of the conditional distribution  $p(t|\mathbf{x})$ . (The minimum of  $\langle L_1 \rangle$  is obtained by the conditional median.)



# Information theory

Information is a measure of the 'degree of surprise' that a certain value gives us. Unlikely events are informative, likely events less so. Certain events give us no additional information. Thus, information decreases with the probability of the event.

Let us denote  $h(x)$  the information of  $x$ . Then if  $x, y$  are two independent events:  $h(x, y) = h(x) + h(y)$ . Since  $p(x, y) = p(x)p(y)$  we see that

$$h(x) = -\log_2 p(x)$$

is a good candidate to quantify the information in  $x$ .

If  $x$  is observed repeatedly then the expected information is

$$H[x] := \langle -\log_2 p \rangle = - \sum_x p(x) \log_2 p(x)$$

is the entropy of the distribution  $p$ .

# Information theory

Example 1:  $x$  can have 8 values with equal probability, then  $H(x) = -8 \times \frac{1}{8} \log \frac{1}{8} = 3$  bits.

Example 2:  $x$  can have 8 values with probabilities  $(\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{16}, \frac{1}{64}, \frac{1}{64}, \frac{1}{64}, \frac{1}{64})$ . Then

$$H(x) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{4} \log \frac{1}{4} - \frac{1}{8} \log \frac{1}{8} - \frac{1}{16} \log \frac{1}{16} - \frac{4}{64} \log \frac{1}{64} = 2 \text{ bits}$$

which is smaller than for the uniform distribution.

*Noiseless coding theorem:* Entropy is a lower bound on the average number of bits needed to transmit a random variable (Shannon 1948).

Q: How can we transmit  $x$  in example 2 most efficiently?

# Information theory

Example 1:  $x$  can have 8 values with equal probability, then  $H(x) = -8 \times \frac{1}{8} \log \frac{1}{8} = 3$  bits.

Example 2:  $x$  can have 8 values with probabilities  $(\frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{16}, \frac{1}{64}, \frac{1}{64}, \frac{1}{64}, \frac{1}{64})$ . Then

$$H(x) = -\frac{1}{2} \log \frac{1}{2} - \frac{1}{4} \log \frac{1}{4} - \frac{1}{8} \log \frac{1}{8} - \frac{1}{16} \log \frac{1}{16} - \frac{4}{64} \log \frac{1}{64} = 2 \text{ bits}$$

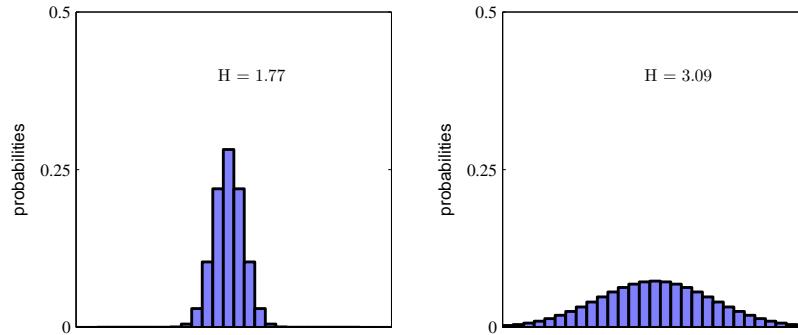
which is smaller than for the uniform distribution.

*Noiseless coding theorem:* Entropy is a lower bound on the average number of bits needed to transmit a random variable (Shannon 1948).

A: We can encode  $x$  as a 3 bit binary number, in which case the expected code length is 3 bits. We can do better, by coding likely  $x$  smaller and unlikely  $x$  larger, for instance 0, 10, 110, 1110, 111100, 111101, 111110, 111111. Then

$$\text{Av.codelength} = \frac{1}{2} \times 1 + \frac{1}{4} \times 2 + \frac{1}{8} \times 3 + \frac{1}{16} \times 4 + \frac{4}{64} \times 6 = 2 \text{ bits}$$

# Information theory



When  $x$  has values  $x_i, i = 1, \dots, M$ , then

$$H[x] = - \sum_i p(x_i) \log p(x_i)$$

When  $p$  is sharply peaked ( $p(x_1) = 1, p(x_2) = \dots = p(x_M) = 0$ ) then the entropy is

$$H[x] = -1 \log 1 - (M-1)0 \log 0 = 0$$

When  $p$  is flat ( $p(x_i) = 1/M$ ) the entropy is maximal

$$H[x] = -M \frac{1}{M} \log \frac{1}{M} = \log M$$



# Information theory/Maximum entropy

For  $p(x)$  a distribution density over a continuous value  $x$  we define the (differential) entropy as

$$H[x] = - \int p(x) \log p(x) dx$$

Suppose that all we know about  $p$  is its mean  $\mu$  and its variance  $\sigma^2$ .

Q: What is the distribution  $p$  with mean  $\mu$  and variance  $\sigma^2$  that is as *uninformative* as possible, i.e., which maximizes the entropy?

# Information theory/Maximum entropy

For  $p(x)$  a distribution density over a continuous value  $x$  we define the (differential) entropy as

$$H[x] = - \int p(x) \log p(x) dx$$

Suppose that all we know about  $p$  is its mean  $\mu$  and its variance  $\sigma^2$ .

Q: What is the distribution  $p$  with mean  $\mu$  and variance  $\sigma^2$  that is as *uninformative* as possible, i.e., which maximizes the entropy?

A: The Gaussian distribution  $\mathcal{N}(x|\mu, \sigma^2)$  (exercise 1.34 and 1.35).

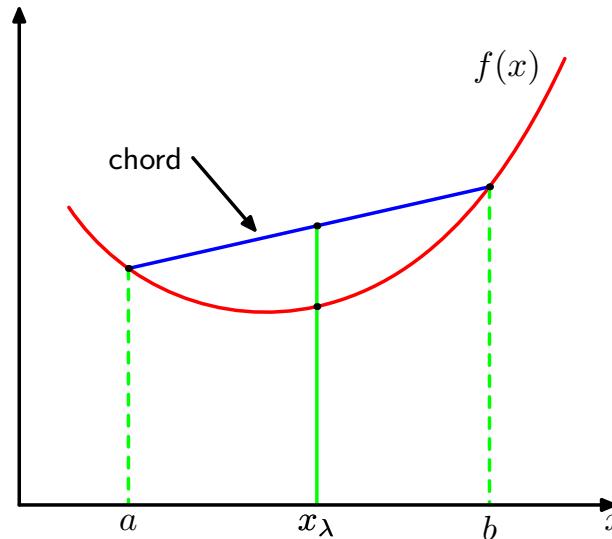
# Information theory/KL-divergence

*Relative entropy or Kullback-Leibler divergence or KL-divergence:*

$$\begin{aligned}\text{KL}(p||q) &= - \sum_i p_i \ln q_i - \left( - \sum_i p_i \ln p_i \right) \\ &= - \sum_i p_i \ln \left\{ \frac{q_i}{p_i} \right\}\end{aligned}$$

- Additional amount of information required to specify  $i$  when  $q$  is used for coding rather than the true distribution  $p$ .
- Divergence between ‘true’ distribution  $p$  and ‘approximate’ distribution  $q$ .
- $\text{KL}(p||q) \neq \text{KL}(q||p)$
- $\text{KL}(p||q) \geq 0$ ,  $\text{KL}(p||q) = 0 \Leftrightarrow p = q$  (use *convex functions*)
- with continuous variables:  $\text{KL}(p||q) = - \int p(\mathbf{x}) \ln \left\{ \frac{q(\mathbf{x})}{p(\mathbf{x})} \right\} d\mathbf{x}$

# Convex functions



Convex function: every chord lies on or above the function.

$$f \text{ is convex} \iff f(\lambda a + (1 - \lambda)b) \leq \lambda f(a) + (1 - \lambda)f(b) \quad \forall \lambda \in [0, 1], \forall a, b$$

- Examples:  $f(x) = ax + b$ ,  $f(x) = x^2$ ,  $f(x) = -\ln(x)$  and  $f(x) = x \ln(x)$  (exercise).
- Convex:  $\cup$  shaped. Concave:  $\cap$  shaped.
- Convex  $\Leftrightarrow$  second derivative non-negative.

# Convex functions/Jensen's inequality

Convex functions satisfy *Jensen's inequality*

$$f\left(\sum_{i=1}^M \lambda_i x_i\right) \leq \sum_{i=1}^M \lambda_i f(x_i)$$

where  $\lambda_i \geq 0$ ,  $\sum_i \lambda_i = 1$ , for any set points  $x_i$ .

In other words:

$$f(\langle x \rangle) \leq \langle f(x) \rangle$$

Example: to show that  $\text{KL}(p||q)$ , we apply Jensen's inequality with  $\lambda_i = p_i$ , making use of the fact that  $-\ln(x)$  is convex:

$$\text{KL}(p||q) = - \sum_i p_i \ln\left(\frac{q_i}{p_i}\right) \geq - \ln\left(\sum_i p_i \frac{q_i}{p_i}\right) = - \ln\left(\sum_i q_i\right) = 0$$

# Information theory and density estimation

Relation with maximum likelihood:

*Empirical distribution :*

$$p(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^N \delta(\mathbf{x} - \mathbf{x}_n)$$

Approximating distribution (model) :  $q(\mathbf{x}|\boldsymbol{\theta})$

$$\begin{aligned} \text{KL}(p||q) &= - \int p(\mathbf{x}) \ln q(\mathbf{x}|\boldsymbol{\theta}) d\mathbf{x} - \int p(\mathbf{x}) \ln p(\mathbf{x}) d\mathbf{x} \\ &= -\frac{1}{N} \sum_{n=1}^N \ln q(\mathbf{x}_n|\boldsymbol{\theta}) + \text{const.} \end{aligned}$$

Thus, minimizing the KL-divergence between the empirical distribution  $p(\mathbf{x})$  and the model distribution  $q(\mathbf{x}|\boldsymbol{\theta})$  is equivalent to maximum likelihood (i.e., maximizing the likelihood of i.i.d. data with respect to the the model parameters  $\boldsymbol{\theta}$ ).

# Information theory/mutual information

Mutual information between  $\mathbf{x}$  and  $\mathbf{y}$ : KL divergence between joint distribution  $p(\mathbf{x}, \mathbf{y})$  and product of marginals  $p(\mathbf{x})p(\mathbf{y})$ ,

$$\begin{aligned} I[\mathbf{x}, \mathbf{y}] &\equiv \text{KL}(p(\mathbf{x}, \mathbf{y}) || p(\mathbf{x})p(\mathbf{y})) \\ &= - \int \int p(\mathbf{x}, \mathbf{y}) \ln \left( \frac{p(\mathbf{x})p(\mathbf{y})}{p(\mathbf{x}, \mathbf{y})} \right) d\mathbf{x}d\mathbf{y} \end{aligned}$$

- $I(\mathbf{x}, \mathbf{y}) \geq 0$ , equality iff  $\mathbf{x}$  and  $\mathbf{y}$  independent

Relation with conditional entropy

$$I[\mathbf{x}, \mathbf{y}] = H[\mathbf{x}] - H[\mathbf{x}|\mathbf{y}] = H[\mathbf{y}] - H[\mathbf{y}|\mathbf{x}]$$

# Lagrange multipliers

Minimize  $f(\mathbf{x})$  under constraint:  $g(\mathbf{x}) = 0$ .

Fancy formulation: define *Lagrangian*,

$$L(\mathbf{x}, \lambda) = f(\mathbf{x}) + \lambda g(\mathbf{x})$$

$\lambda$  is called a Lagrange multiplier.

The constraint minimization of  $f$  w.r.t  $\mathbf{x}$  equivalent to *unconstraint* minimization of  $\max_{\lambda} L(\mathbf{x}, \lambda)$  w.r.t  $\mathbf{x}$ . The *maximization* w.r.t to  $\lambda$  yields the following function of  $\mathbf{x}$

$$\max_{\lambda} L(\mathbf{x}, \lambda) = f(\mathbf{x}) \quad \text{if } g(\mathbf{x}) = 0$$

$$\max_{\lambda} L(\mathbf{x}, \lambda) = \infty \quad \text{otherwise}$$

# Lagrange multipliers

Under certain conditions, in particular  $f(\mathbf{x})$  convex (i.e. the matrix of second derivatives positive definite) and  $g(\mathbf{x})$  linear,

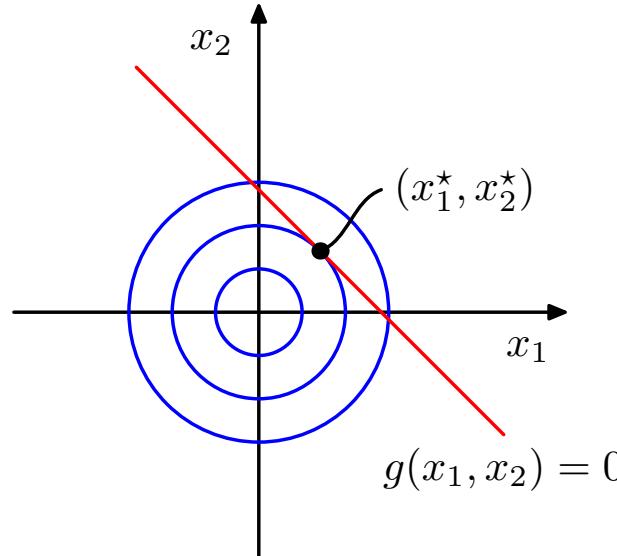
$$\min_{\mathbf{x}} \max_{\lambda} L(\mathbf{x}, \lambda) = \max_{\lambda} \min_{\mathbf{x}} L(\mathbf{x}, \lambda)$$

Procedure:

1. Minimize  $L(\mathbf{x}, \lambda)$  w.r.t  $\mathbf{x}$ , e.g. by taking the gradient and set to zero. This yields a (parametrized) solution  $\mathbf{x}(\lambda)$ .
2. Maximize  $L(\mathbf{x}(\lambda), \lambda)$  w.r.t.  $\lambda$ . The solution  $\lambda^*$  is precisely such that  $g(\mathbf{x}(\lambda^*)) = 0$ .
3. The solution of the constraint optimization problem is

$$\mathbf{x}^* = \mathbf{x}(\lambda^*)$$

# Example



$$f(x_1, x_2) = 1 - x_1^2 - x_2^2 \quad \text{and constraint} \quad g(x_1, x_2) = x_1 + x_2 - 1 = 0$$

Lagrangian:

$$L(x_1, x_2, \lambda) = 1 - x_1^2 - x_2^2 + \lambda(x_1 + x_2 - 1)$$

Minimize  $L$  w.r.t.  $x_i$  gives  $x_i(\lambda) = \frac{1}{2}\lambda$ .

Plug into constraint:  $x_1(\lambda) + x_2(\lambda) - 1 = \lambda - 1 = 0$ .

So  $\lambda = 1$  and  $x_i^* = \frac{1}{2}$

## Some remarks

- Works as well for maximization (of concave functions) under constraints. The procedure is essentially the same.
- The sign in front of the  $\lambda$  can be chosen as you want:

$$L(\mathbf{x}, \lambda) = f(\mathbf{x}) + \lambda g(\mathbf{x}) \quad \text{or} \quad L(\mathbf{x}, \lambda) = f(\mathbf{x}) - \lambda g(\mathbf{x})$$

work equally well.

- More constraints? For each constraint  $g_i(\mathbf{x}) = 0$  a Lagrange multiplier  $\lambda_i$ , so

$$L(\mathbf{x}, \lambda) = f(\mathbf{x}) + \sum_i \lambda_i g_i(\mathbf{x})$$

- Similar methods apply for inequality constraints  $g(\mathbf{x}) \geq 0$  (restricts  $\lambda$ ).

## Chapter 2

# Probability distributions

- Density estimation
- Parametric distributions
- Maximum likelihood, Bayesian inference, conjugate priors
- Bernoulli (binary), Beta, Gaussian, ..., exponential family
- Nonparametric distribution

# Binary variables / Bernoulli distribution

$$x \in \{0, 1\}$$

$$\begin{aligned} p(x = 1|\mu) &= \mu, \\ p(x = 0|\mu) &= 1 - \mu \end{aligned}$$

Bernoulli distribution:

$$\text{Bern}(x|\mu) = \mu^x(1 - \mu)^{1-x}$$

Mean and variance:

$$\begin{aligned} \mathbb{E}[x] &= \mu \\ \text{var}[x] &= \mu(1 - \mu) \end{aligned}$$

# Binary variables / Bernoulli distribution

Data set (i.i.d)  $D = \{x_1, \dots, x_n\}$ , with  $x_i \in \{0, 1\}$ .

Likelihood:

$$p(D|\mu) = \prod_n p(x_n|\mu) = \prod_n \mu^{x_n} (1-\mu)^{1-x_n}$$

Log likelihood

$$\begin{aligned} \ln p(D|\mu) &= \sum_n \ln p(x_n|\mu) = \sum_n x_n \ln \mu + (1-x_n) \ln(1-\mu) \\ &= m \ln \mu + (N-m) \ln(1-\mu) \end{aligned}$$

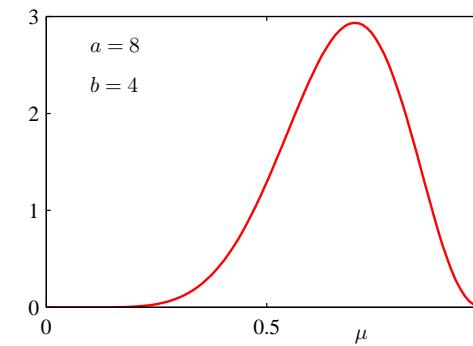
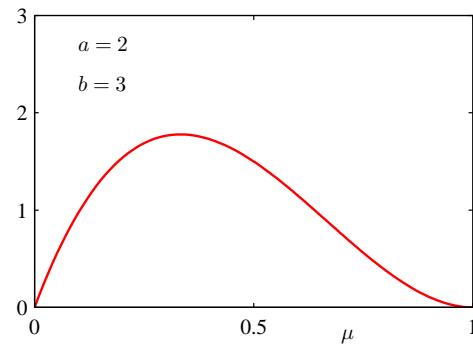
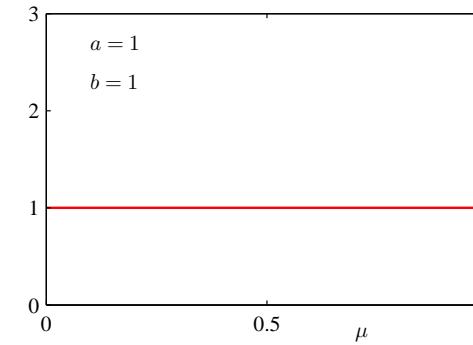
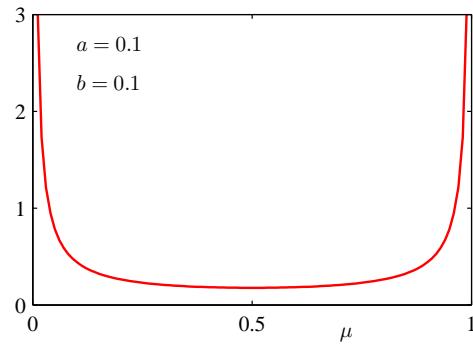
where  $m = \sum_n x_n$ , the total number of  $x_n = 1$ .

Maximization w.r.t  $\mu$  gives maximum likelihood solution:

$$\mu_{ML} = \frac{m}{N}$$

# The beta distribution

Distribution for parameters  $\mu$ . Conjugate prior for Bayesian treatment for problem.



# The beta distribution

$$\begin{aligned}\text{Beta}(\mu|a, b) &= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)}\mu^{a-1}(1-\mu)^{b-1} \\ &\propto \mu^{a-1}(1-\mu)^{b-1} \quad 0 \leq \mu \leq 1\end{aligned}$$

Normalisation

$$\int_0^1 \text{Beta}(\mu|a, b) = 1$$

Mean and variance

$$\begin{aligned}\mathbb{E}[\mu] &= \frac{a}{a+b} \\ \text{var}[\mu] &= \frac{ab}{(a+b)^2(a+b+1)}\end{aligned}$$

# Bayesian inference with binary variables

Prior:

$$p(\mu) = \text{Beta}(\mu|a, b) \propto \mu^{a-1}(1-\mu)^{b-1}$$

Likelihood – Data set (i.i.d)  $D = \{x_1, \dots, x_N\}$ , with  $x_i \in \{0, 1\}$ .  
 Assume  $m$  ones and  $l$  zeros, ( $m + l = N$ )

$$\begin{aligned} p(D|\mu) &= \prod_n p(x_n|\mu) = \prod_n \mu^{x_n}(1-\mu)^{1-x_n} \\ &= \mu^m(1-\mu)^l \end{aligned}$$

Posterior

$$\begin{aligned} p(\mu|D) &\propto p(D|\mu)p(\mu) \\ &= \mu^m(1-\mu)^l \times \mu^{a-1}(1-\mu)^{b-1} \\ &= \mu^{m+a-1}(1-\mu)^{l+b-1} \propto \text{Beta}(\mu|a+m, b+l) \end{aligned}$$

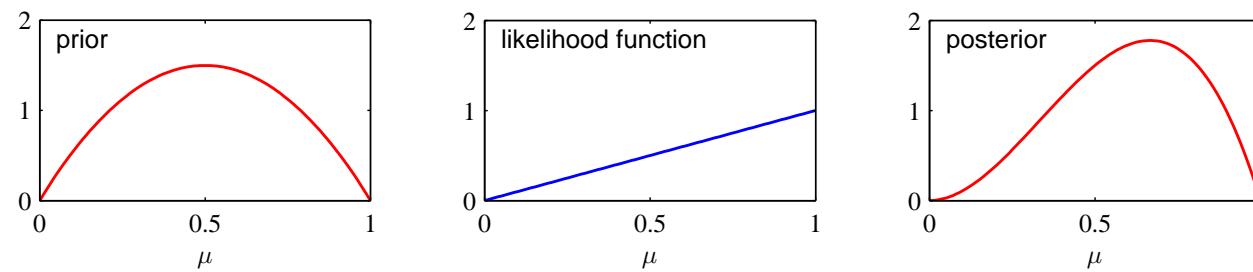
# Bayesian inference with binary variables

Interpretation: Hyperparameters  $a$  and  $b$  effective number of ones and zeros.

Data: increments of these parameters.

Conjugacy:

- (1) prior has the same form as likelihood function.
- (2) this form is preserved in the product (the posterior)



Posterior interpreted as updated prior: sequential learning

## Bayesian inference with binary variables

Prediction of next data point given data  $D$ :

$$p(x = 1|D) = \int_0^1 p(x = 1|\mu)p(\mu|D)d\mu = \int_0^1 \mu p(\mu|D)d\mu = \mathbb{E}[\mu|D]$$

with posterior is  $\text{Beta}(\mu|a + m, b + l)$ , and  $\mathbb{E}[\mu|a, b] = a/(a + b)$  we find

$$p(x = 1|D) = \frac{m + a}{m + a + l + b}$$

## Multinomial variables

Alternative representation for Bernoulli distribution:  $x \in \{v_1, v_2\}$ , parameter vector:  $\mu = (\mu_1, \mu_2)$ , with  $\mu_1 + \mu_2 = 1$ .

$$p(x = v_k | \mu) = \mu_k$$

In fancy notation:

$$p(x|\mu) = \prod_k \mu_k^{\delta_{xv_k}}$$

Generalizes to multinomial variables:  $x \in \{v_1, \dots, v_K\}$ ,

$$\mu = (\mu_1, \dots, \mu_K) \quad \sum_k \mu_k = 1$$

## Multinomial variables: Maximum likelihood

Likelihood:

$$\begin{aligned}
 p(D|\boldsymbol{\mu}) &= \prod_n p(x_n|\boldsymbol{\mu}) = \prod_n \prod_k \mu_k^{\delta_{x_n v_k}} \\
 &= \prod_k \mu_k^{\sum_n \delta_{x_n v_k}} \\
 &= \prod_k \mu_k^{m_k}
 \end{aligned}$$

with  $m_k = \sum_n \delta_{x_n v_k}$ , the total number of datapoints with value  $v_k$ . Log likelihood

$$\ln p(D|\boldsymbol{\mu}) = \sum_k m_k \ln \mu_k$$

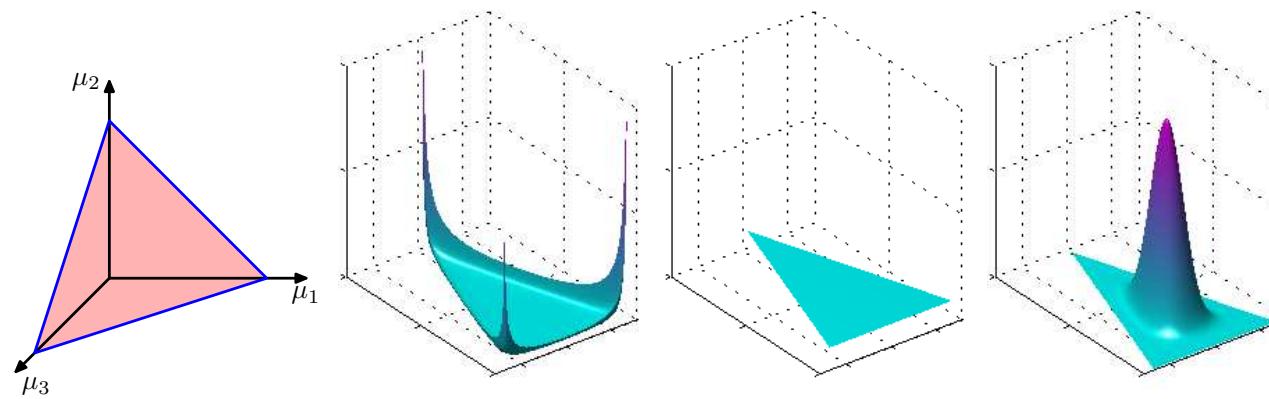
Maximize with constraints using Lagrange multipliers.

# Dirichlet distribution

$$\text{Dir}(\boldsymbol{\mu}|\boldsymbol{\alpha}) \propto \prod_k \mu_k^{\alpha_k}$$

Probability distribution on the *simplex*:

$$S^K = \{(\mu_1, \dots, \mu_K) | 0 \leq \mu_k \leq 1, \sum_{k=1}^K \mu_k = 1\}$$



## Dirichlet distribution

$$\text{Dir}(\boldsymbol{\mu}|\boldsymbol{\alpha}) \propto \prod_k \mu_k^{\alpha_k}$$

Probability distribution on the *simplex*:

$$S^K = \{(\mu_1, \dots, \mu_K) | 0 \leq \mu_k \leq 1, \sum_{k=1}^K \mu_k = 1\}$$

Bayesian inference: Prior  $\text{Dir}(\boldsymbol{\mu}|\boldsymbol{\alpha}) +$  data counts  $\mathbf{m}$

→ Posterior  $\text{Dir}(\boldsymbol{\mu}|\boldsymbol{\alpha} + \mathbf{m})$

Parameters  $\boldsymbol{\alpha}$ : ‘pseudocounts’.

# Gaussian

Gaussian distribution

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$$

Specified by  $\mu$  and  $\sigma^2$

In  $d$  dimensions

$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{d/2}|\boldsymbol{\Sigma}|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}-\boldsymbol{\mu})\right)$$

# Central limit theorem

Sum of large set of random variables is approximately Gaussian distributed.

Let  $X_i$  set of random variables with mean  $\mu$  and variance  $\sigma^2$ .

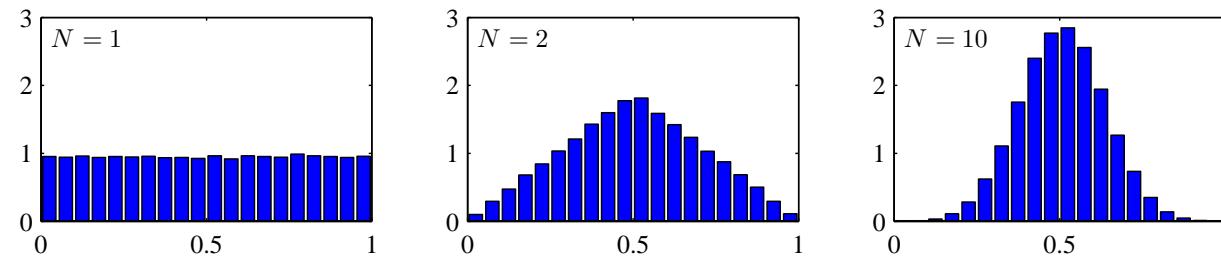
The sum of first  $n$  variables is  $S_n = X_1 + \dots + X_n$ . Now if  $n \rightarrow \infty$

Law of large numbers: The mean of the sum  $Y_n = \frac{S_n}{n}$  converges to  $\mu$

Central limit theorem: Distribution of

$$Z_n = \frac{S_n - n\mu}{\sigma\sqrt{n}}$$

converges to Gaussian  $\mathcal{N}(Z_n|0, 1)$



# Understanding the Gaussian through the covariance matrix $\Sigma$

Dependence on  $x$  only through the quadratic form

$$(x - \mu)^T \Sigma^{-1} (x - \mu)$$

## Symmetric matrices

$$A_{ij} = A_{ji}, \quad A^T = A$$

Inverse of a matrix is a matrix  $A^{-1}$  such that  $A^{-1}A = AA^{-1} = I$ , where  $I$  is the identity matrix.

$A^{-1}$  is also symmetric:

$$I = I^T = (A^{-1}A)^T = A^T(A^{-1})^T = A(A^{-1})^T$$

Thus,  $(A^{-1})^T = A^{-1}$ .

# Eigenvalues

A symmetric real-valued  $d \times d$  matrix has  $d$  real eigenvalues  $\lambda_k$  and  $d$  eigenvectors  $\mathbf{u}_k$ :

$$A\mathbf{u}_k = \lambda_k \mathbf{u}_k, \quad k = 1, \dots, d$$

or

$$(A - \lambda_k I)\mathbf{u}_k = 0$$

Solution of this equation for non-zero  $\mathbf{u}_k$  requires  $\lambda_k$  to satisfy the characteristic equation:

$$\det(A - \lambda I) = 0$$

This is a polynomial equation in  $\lambda$  of degree  $d$  and has thus  $d$  solutions <sup>4</sup>  $\lambda_1, \dots, \lambda_d$ .

---

<sup>4</sup>In general, the solutions are complex. It can be shown that with symmetric matrices, the solutions are in fact real.

# Eigenvectors

Consider two different eigenvectors  $k$  and  $j$ . Multiply the  $k$ -th eigenvalue equation by  $\mathbf{u}_j$  from the left:

$$\mathbf{u}_j^T A \mathbf{u}_k = \lambda_k \mathbf{u}_j^T \mathbf{u}_k$$

Multiply the  $j$ -th eigenvalue equation by  $\mathbf{u}_k$  from the left:

$$\mathbf{u}_k^T A \mathbf{u}_j = \lambda_j \mathbf{u}_k^T \mathbf{u}_j = \lambda_j \mathbf{u}_j^T \mathbf{u}_k$$

Subtract

$$(\lambda_k - \lambda_j) \mathbf{u}_j^T \mathbf{u}_k = 0$$

Thus, eigenvectors with different eigenvalues are orthogonal

If  $\lambda_k = \lambda_j$  then any linear combination is also an eigenvector:

$$A(\alpha \mathbf{u}_k + \beta \mathbf{u}_j) = \lambda_k(\alpha \mathbf{u}_k + \beta \mathbf{u}_j)$$

This can be used to choose eigenvectors with identical eigenvalues orthogonal.

If  $\mathbf{u}_k$  is an eigenvector of  $A$ , then  $\alpha \mathbf{u}_k$  is also an eigenvector of  $A$ . Thus, we can make all eigenvectors the same length one:  $\mathbf{u}_k^T \mathbf{u}_k = 1$ .

In summary,

$$\mathbf{u}_j^T \mathbf{u}_k = \delta_{jk}$$

with  $\delta_{jk}$  the Kronecker delta, is equal to 1 if  $j=k$  and zero otherwise.

The eigenvectors span the  $d$ -dimensional space as an orthonormal basis.

# Orthogonal matrices

Write  $U = (\mathbf{u}_1, \dots, \mathbf{u}_d)$ .

$U$  is an orthogonal matrix<sup>5</sup>, i.e.

$$U^T U = I$$

For orthogonal matrices,

$$U^T U U^{-1} = U^{-1} = U^T$$

So  $UU^T = I$ , i.e. the transposed is orthogonal as well (note that  $U$  is in general *not* symmetric).

Furthermore,

$$\begin{aligned} \det(UU^T) &= 1 \Rightarrow \det(U)\det(U^T) = 1 \\ \Rightarrow \det(U) &= \pm 1 \end{aligned}$$

5

$$U_{ij} = (u_j)_i, \quad (U^T U)_{ij} = \sum_k (U^T)_{ik} U_{kj} = \sum_k U_{ki} U_{kj} = \sum_k (u_i)_k (u_j)_k = \mathbf{u}_i^T \cdot \mathbf{u}_j = \delta_{ij}$$

Orthogonal matrices implement rigid rotations, i.e. length and angle preserving.

$$\bar{\mathbf{x}}_1 = U\mathbf{x}_1 \quad \bar{\mathbf{x}}_2 = U\mathbf{x}_2$$

then

$$\bar{\mathbf{x}}_1^T \bar{\mathbf{x}}_2 = \mathbf{x}_1^T U^T U \mathbf{x}_2 = \mathbf{x}_1^T \mathbf{x}_2$$

# Diagonalization

The eigenvector equation can be written as

$$\begin{aligned}
 AU &= A(\mathbf{u}_1, \dots, \mathbf{u}_d) = (A\mathbf{u}_1, \dots, A\mathbf{u}_d) = (\lambda_1 \mathbf{u}_1, \dots, \lambda_d \mathbf{u}_d) \\
 &= (\mathbf{u}_1, \dots, \mathbf{u}_d) \begin{pmatrix} \lambda_1 & \dots & 0 \\ \vdots & & \vdots \\ 0 & \dots & \lambda_d \end{pmatrix} \\
 &= U\Lambda
 \end{aligned}$$

By right-multiplying by  $U^T$  we obtain the important result

$$A = U\Lambda U^T$$

which can also be written as 'expansion in eigenvectors'

$$A = \sum_{k=1}^d \lambda_k \mathbf{u}_k \mathbf{u}_k^T$$

# Applications

$$A = U\Lambda U^T \Rightarrow A^2 = U\Lambda U^T U\Lambda U^T = U\Lambda^2 U^T$$

$$A^n = U\Lambda^n U^T \quad A^{-n} = U\Lambda^{-n} U^T$$

Determinant is product of eigenvalues:

$$\det(A) = \det(U\Lambda U^T) = \det(U) \det(\Lambda) \det(U^T) = \prod_k \lambda_k$$

## Basis transformation

We can represent an arbitrary vector  $\mathbf{x}$  in  $d$  dimensions on a new basis  $U$  as

$$\mathbf{x} = UU^T\mathbf{x} = \sum_{k=1}^d \mathbf{u}_k(\mathbf{u}_k^T\mathbf{x})$$

The numbers  $\bar{x}_k = (\mathbf{u}_k^T \mathbf{x})$  are the components of  $\mathbf{x}$  on the basis  $\mathbf{u}_k, k = 1, \dots, d$ , i.e. on the new basis, the vector has components  $(U^T \mathbf{x})$ . If the matrix  $A$  is the representation of a linear transformation on the old basis, the matrix with components

$$A' = U^T A U$$

is the representation on the new basis.

For instance if  $\mathbf{y} = A\mathbf{x}$ ,  $\bar{\mathbf{x}} = U^T \mathbf{x}$ ,  $\bar{\mathbf{y}} = U^T \mathbf{y}$ , then

$$A'\bar{\mathbf{x}} = U^T A U U^T \mathbf{x} = U^T A \mathbf{x} = U^T \mathbf{y} = \bar{\mathbf{y}}$$

So a matrix is diagonal on a basis of its eigenvectors:

$$A' = U^T A U = U^T U \Lambda U^T U = \Lambda$$

# Multivariate Gaussian

In  $d$  dimensions

$$\begin{aligned}\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) &= \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})\right) \\ \boldsymbol{\mu} &= \langle \mathbf{x} \rangle = \int \mathbf{x} \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{x} \\ \boldsymbol{\Sigma} &= \langle (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T \rangle = \int (\mathbf{x} - \boldsymbol{\mu})^T (\mathbf{x} - \boldsymbol{\mu}) \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{x}\end{aligned}$$

We can also write this in component notation:

$$\begin{aligned}\mu_i &= \langle x_i \rangle = \int x_i \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{x} \\ \Sigma_{ij} &= \langle (x_i - \mu_i)(x_j - \mu_j) \rangle = \int (x_i - \mu_i)(x_j - \mu_j) \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) d\mathbf{x}\end{aligned}$$

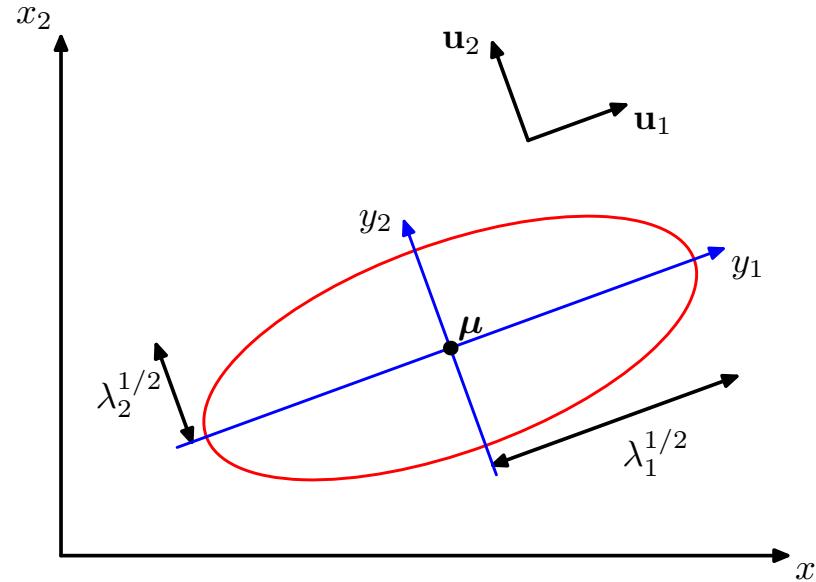
# Multivariate Gaussian

Spectral (eigenvalue) representation of  $\Sigma$ :

$$\Sigma = U \Lambda U^T = \sum_k \lambda_k \mathbf{u}_k \mathbf{u}_k^T$$

$$\Sigma^{-1} = U \Lambda^{-1} U^T = \sum_k \lambda_k^{-1} \mathbf{u}_k \mathbf{u}_k^T$$

$$\begin{aligned} (\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) &= (\mathbf{x} - \boldsymbol{\mu})^T U \Lambda^{-1} U^T (\mathbf{x} - \boldsymbol{\mu}) \\ &= \mathbf{y}^T \Lambda^{-1} \mathbf{y} \end{aligned}$$



## Multivariate Gaussian

Explain the normalization: use transformation<sup>6</sup>  $\mathbf{y} = U^T(\mathbf{x} - \boldsymbol{\mu}) = U^T\mathbf{z}$

$$\begin{aligned}
 & \int \exp \left( -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \right) d\mathbf{x} \\
 &= \int \exp \left( -\frac{1}{2}\mathbf{z}^T \boldsymbol{\Sigma}^{-1}\mathbf{z} \right) d\mathbf{z} = \int \left| \det \left( \frac{\partial \mathbf{z}}{\partial \mathbf{y}} \right) \right| \exp \left( -\frac{1}{2}\mathbf{y}^T \Lambda^{-1}\mathbf{y} \right) d\mathbf{y} \\
 &= \int \exp \left( -\frac{1}{2\lambda_1}y_1^2 \right) dy_1 \dots \int \exp \left( -\frac{1}{2\lambda_d}y_d^2 \right) dy_d \\
 &= \sqrt{2\pi\lambda_1} \dots \sqrt{2\pi\lambda_d} = (2\pi)^{d/2} \left( \prod_i \lambda_i \right)^{1/2} = (2\pi)^{d/2} \det(\boldsymbol{\Sigma})^{1/2}
 \end{aligned}$$

The multivariate Gaussian becomes a product of independent factors on the basis of eigenvectors.

---

<sup>6</sup>Bishop alternates between defining matrix  $U$  as *rows* of eigenvectors  $\mathbf{u}_i^T$  (2.52), and  $U$  as *columns* of eigenvectors  $\mathbf{u}_i$  (C.37). In the slides we always use the more conventional *column* representation.

## Multivariate Gaussian

Compute the expectation value : use shift-transformation  $\mathbf{z} = \mathbf{x} - \boldsymbol{\mu}$   
 Take  $Z$  as normalisation constant.

Use symmetry  $f(\mathbf{z}) = -f(-\mathbf{z}) \Rightarrow \int f(\mathbf{z}) d\mathbf{z} = 0$

$$\begin{aligned}\mathbb{E}[\mathbf{x}] &= \frac{1}{Z} \int \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right) \mathbf{x} d\mathbf{x} \\ &= \frac{1}{Z} \int \exp\left(-\frac{1}{2}\mathbf{z}^T \boldsymbol{\Sigma}^{-1}\mathbf{z}\right) (\mathbf{z} + \boldsymbol{\mu}) d\mathbf{z} \\ &= \boldsymbol{\mu}\end{aligned}$$

## Multivariate Gaussian

Second order moment: use transformation  $\mathbf{y} = U^T(\mathbf{x} - \boldsymbol{\mu}) = U^T\mathbf{z}$

First, shift by  $\mathbf{z} = \mathbf{x} - \boldsymbol{\mu}$

$$\begin{aligned}
 \mathbb{E}[\mathbf{x}\mathbf{x}^T] &= \int \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \mathbf{x}\mathbf{x}^T d\mathbf{x} \\
 &= \int \mathcal{N}(\mathbf{z}|0, \boldsymbol{\Sigma}) (\mathbf{z} + \boldsymbol{\mu})(\mathbf{z} + \boldsymbol{\mu})^T d\mathbf{z} \\
 &= \int \mathcal{N}(\mathbf{z}|0, \boldsymbol{\Sigma}) (\mathbf{z}\mathbf{z}^T + \mathbf{z}\boldsymbol{\mu}^T + \boldsymbol{\mu}\mathbf{z}^T + \boldsymbol{\mu}\boldsymbol{\mu}^T) d\mathbf{z} \\
 &= \int \mathcal{N}(\mathbf{z}|0, \boldsymbol{\Sigma}) \mathbf{z}\mathbf{z}^T d\mathbf{z} + \boldsymbol{\mu}\boldsymbol{\mu}^T
 \end{aligned}$$

## Multivariate Gaussian

Now use transformation  $\mathbf{y} = U^T \mathbf{z}$ , and use  $\Sigma = U\Lambda U^T$

$$\begin{aligned}\int \mathcal{N}(\mathbf{z}|0, \Sigma) \mathbf{z} \mathbf{z}^T d\mathbf{z} &= \int \mathcal{N}(\mathbf{y}|0, \Lambda) U \mathbf{y} \mathbf{y}^T U^T d\mathbf{y} \\ &= U \int \mathcal{N}(\mathbf{y}|0, \Lambda) \mathbf{y} \mathbf{y}^T d\mathbf{y} U^T\end{aligned}$$

Component-wise computation shows  $\int \mathcal{N}(\mathbf{y}|0, \Lambda) \mathbf{y} \mathbf{y}^T d\mathbf{y} = \Lambda$ :

$$\begin{aligned}i \neq j \rightarrow \int \mathcal{N}(\mathbf{y}|0, \Lambda) y_i y_j d\mathbf{y} &= \int \mathcal{N}(y_i|0, \lambda_i) y_i dy_i \int \mathcal{N}(y_j|0, \lambda_j) y_j dy_j = 0 \\ i = j \rightarrow \int \mathcal{N}(\mathbf{y}|0, \Lambda) y_i^2 d\mathbf{y} &= \int \mathcal{N}(y_i|0, \lambda_i) y_i^2 dy_i = \lambda_i\end{aligned}$$

So

$$\int \mathcal{N}(\mathbf{z}|0, \Sigma) \mathbf{z} \mathbf{z}^T d\mathbf{z} = U \Lambda U^T = \Sigma$$

## Multivariate Gaussian

So, second moment is

$$\mathbb{E}[xx^T] = \Sigma + \mu\mu^T$$

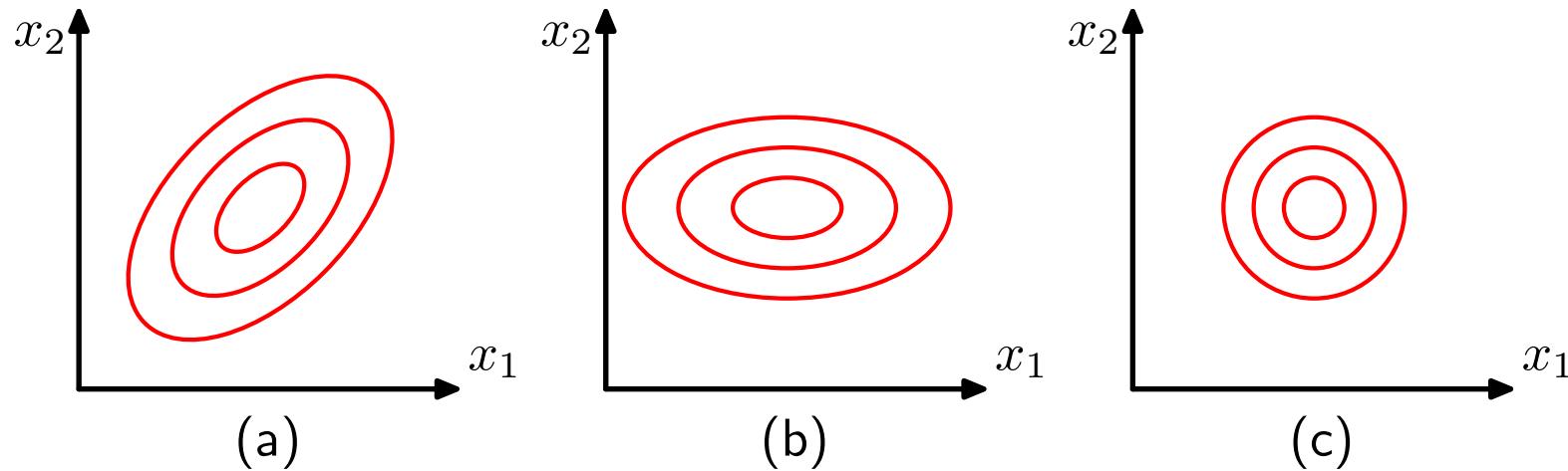
Covariance

$$\text{cov}[x] = \mathbb{E}[xx^T] - \mathbb{E}[x]\mathbb{E}[x]^T = \Sigma$$

## Multivariate Gaussian

Gaussian covariance has  $d(d + 1)$  parameters, mean has  $d$  parameters.

Number of parameters quadratic in  $d$  which may be too large for high dimensional applications.



Common simplifications:  $\Sigma_{ij} = \Sigma_{ii}\delta_{ij}$  ( $2d$  parameters) or  $\Sigma_{ij} = \sigma^2\delta_{ij}$  ( $d + 1$  parameters).

## Conditional of Gaussian is Gaussian

Exponent in Gaussian  $\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$ : quadratic form

$$-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) = -\frac{1}{2}\mathbf{x}^T \boldsymbol{\Sigma}^{-1}\mathbf{x} + \mathbf{x}^T \boldsymbol{\Sigma}^{-1}\boldsymbol{\mu} + c = -\frac{1}{2}\mathbf{x}^T K \mathbf{x} + \mathbf{x}^T K \boldsymbol{\mu} + c$$

Precision matrix  $K = \boldsymbol{\Sigma}^{-1}$ , nb: Conditional

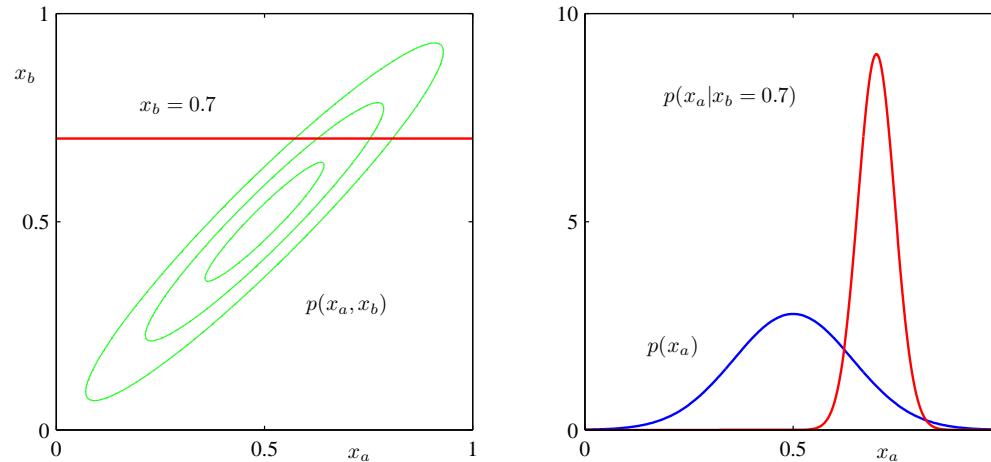
$$p(\mathbf{x}_a | \mathbf{x}_b) = \frac{p(\mathbf{x}_a, \mathbf{x}_b)}{p(\mathbf{x}_b)} \propto p(\mathbf{x}_a, \mathbf{x}_b)$$

Exponent of conditional: collect all terms with  $\mathbf{x}_a$ , ignore constants, regard  $\mathbf{x}_b$  as constant, and write in quadratic form as above

$$\begin{aligned} -\frac{1}{2}\mathbf{x}_a^T K_{a|b} \mathbf{x}_a + \mathbf{x}_a^T K_{a|b} \boldsymbol{\mu}_{a|b} &= -\frac{1}{2}\mathbf{x}_a^T K_{aa} \mathbf{x}_a + \mathbf{x}_a^T K_{aa} \boldsymbol{\mu}_a - \mathbf{x}_a^T K_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b) \\ &= -\frac{1}{2}\mathbf{x}_a^T K_{aa} \mathbf{x}_a + \mathbf{x}_a^T K_{aa} (\boldsymbol{\mu}_a - K_{aa}^{-1} K_{ab} (\mathbf{x}_b - \boldsymbol{\mu}_b)) \end{aligned}$$

# Marginal and conditional Gaussians

Marginal and conditional of Gaussians are also Gaussian



$$p(\mathbf{x}_a | \mathbf{x}_b) = \mathcal{N}(\mathbf{x}_a | \boldsymbol{\mu}_{a|b}, K_{aa}^{-1})$$

$$p(\mathbf{x}_a) = \mathcal{N}(\mathbf{x}_a | \boldsymbol{\mu}_a, \boldsymbol{\Sigma}_{aa})$$

## Some matrix identities

$$\begin{pmatrix} A & B \\ C & D \end{pmatrix}^{-1} = \begin{pmatrix} M & -MBD^{-1} \\ -D^{-1}CM & D^{-1} + D^{-1}CMBD^{-1} \end{pmatrix}$$

with  $M = (A - BD^{-1}C)^{-1}$ .

$$\begin{pmatrix} K_{aa} & K_{ab} \\ K_{ba} & K_{bb} \end{pmatrix} = \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix}^{-1} = \begin{pmatrix} M & -M\Sigma_{ab}\Sigma_{bb}^{-1} \\ -\Sigma_{bb}^{-1}\Sigma_{ba}M & \Sigma_{bb}^{-1} + \Sigma_{bb}^{-1}\Sigma_{ba}M\Sigma_{ab}\Sigma_{bb}^{-1} \end{pmatrix}$$

with  $M = (\Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba})^{-1}$ .

Thus,

$$K_{aa} = (\Sigma_{aa} - \Sigma_{ab}\Sigma_{bb}^{-1}\Sigma_{ba})^{-1}$$

## Bayes' theorem for linear Gaussian model

Given marginal Gaussian on  $\mathbf{x}$  and linear relation  $\mathbf{y} = \mathbf{Ax} + \mathbf{b} + \xi$ :

$$\begin{aligned} p(\mathbf{x}) &= \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) \\ p(\mathbf{y}|\mathbf{x}) &= \mathcal{N}(\mathbf{y}|\mathbf{Ax} + \mathbf{b}, \mathbf{L}^{-1}) \end{aligned}$$

Then (see next slide):

$$\begin{aligned} p(\mathbf{y}) &= \mathcal{N}(\mathbf{y}|\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^T) \\ p(\mathbf{x}|\mathbf{y}) &= \mathcal{N}(\mathbf{x}|\boldsymbol{\Sigma}\{\mathbf{A}^T\mathbf{L}(\mathbf{y} - \mathbf{b}) + \boldsymbol{\Lambda}\boldsymbol{\mu}\}, \boldsymbol{\Sigma}) \\ \boldsymbol{\Sigma} &= (\boldsymbol{\Lambda} + \mathbf{A}^T\mathbf{L}\mathbf{A})^{-1} \end{aligned}$$

We will use these relations for Bayesian linear regression in section 3.3.

## Details computation $p(\mathbf{y})$

$$\begin{aligned} p(\mathbf{x}) &= \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) & \mathbb{E}[\mathbf{x}] = \boldsymbol{\mu} & \text{cov}[\mathbf{x}] = \boldsymbol{\Lambda}^{-1} \\ p(\mathbf{y}|\mathbf{x}) &= \mathcal{N}(\mathbf{y}|A\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1}) \end{aligned}$$

We write  $\mathbf{y} = A\mathbf{x} + \mathbf{b} + \boldsymbol{\epsilon}$  with  $\mathbb{E}[\boldsymbol{\epsilon}] = 0$ ,  $\text{cov}[\boldsymbol{\epsilon}] = \mathbf{L}^{-1}$ .

$\mathbf{x}, \mathbf{y}$  is jointly Gaussian (product of Gaussians).  $\mathbf{y}$  is Gaussian (marginal of Gaussian).

$$\begin{aligned} \mathbb{E}[\mathbf{y}] &= \mathbb{E}[A\mathbf{x} + \mathbf{b} + \boldsymbol{\epsilon}] = A\boldsymbol{\mu} + \mathbf{b} \\ \text{cov}[\mathbf{y}] &= \text{cov}[A\mathbf{x} + \mathbf{b} + \boldsymbol{\epsilon}] = \text{cov}[A\mathbf{x}] + \text{cov}[\boldsymbol{\epsilon}] = \text{cov}[A\mathbf{x} + \mathbf{L}^{-1}] \\ \text{cov}[A\mathbf{x}] &= \mathbb{E}[(A\mathbf{x} - A\boldsymbol{\mu})(A\mathbf{x} - A\boldsymbol{\mu})^T] = A\mathbb{E}[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T]A^T = A\boldsymbol{\Lambda}^{-1}A^T \end{aligned}$$

Thus,

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|A\boldsymbol{\mu} + \mathbf{b}, A\boldsymbol{\Lambda}^{-1}A^T + \mathbf{L}^{-1})$$

## Details computation $p(\mathbf{x}|\mathbf{y})$

Write all relevant terms that occur in exponential of the joint Gaussian  $p(\mathbf{x}, \mathbf{y})$ :

$$-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Lambda}(\mathbf{x} - \boldsymbol{\mu}) - \frac{1}{2}(\mathbf{y} - \mathbf{A}\mathbf{x} - \mathbf{b})^T \mathbf{L}(\mathbf{y} - \mathbf{A}\mathbf{x} - \mathbf{b})$$

Collect all quadratic and linear terms in  $\mathbf{x}$ :

$$-\frac{1}{2}\mathbf{x}^T (\boldsymbol{\Lambda} + \mathbf{A}^T \mathbf{L} \mathbf{A}) \mathbf{x} + \mathbf{x}^T (\boldsymbol{\Lambda} \boldsymbol{\mu} + \mathbf{A}^T \mathbf{L}(\mathbf{y} - \mathbf{b}))$$

Define  $\boldsymbol{\Sigma}^{-1} = \boldsymbol{\Lambda} + \mathbf{A}^T \mathbf{L} \mathbf{A}$  and  $\mathbf{m}$  through  $\boldsymbol{\Sigma}^{-1} \mathbf{m} = \boldsymbol{\Lambda} \boldsymbol{\mu} + \mathbf{A}^T \mathbf{L}(\mathbf{y} - \mathbf{b})$ , then

$$-\frac{1}{2}\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} + \mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{m} \propto -\frac{1}{2}(\mathbf{x} - \mathbf{m})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \mathbf{m})$$

Thus

$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\Sigma}(\boldsymbol{\Lambda} \boldsymbol{\mu} + \mathbf{A}^T \mathbf{L}(\mathbf{y} - \mathbf{b})), \boldsymbol{\Sigma})$$

# Bayesian inference for the Gaussian

Aim: inference of unknown parameter  $\mu$ . Assume  $\sigma$  given.

Likelihood of  $\mu$  with one data point:

$$p(x|\mu) = \mathcal{N}(x|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right)$$

Likelihood of  $\mu$  with the data set:

$$\begin{aligned} p(\{x_1, \dots, x_N\}|\mu) &= \prod_{n=1}^N p(x_n|\mu) = \left(\frac{1}{\sqrt{2\pi}\sigma}\right)^N \exp\left(-\frac{1}{2\sigma^2} \sum_n (x_n - \mu)^2\right) \\ &= \exp\left(-\frac{N\mu^2}{2\sigma^2} + \frac{\mu}{\sigma^2} \sum_n x_n + \text{const.}\right) = \exp\left(-\frac{N}{2\sigma^2}\mu^2 + \frac{N\bar{x}}{\sigma^2}\mu + \text{const.}\right) \\ &= \exp\left(-\frac{N}{2\sigma^2}(\mu - \bar{x})^2 + \text{const.}\right) \quad \text{with} \quad \bar{x} = \frac{1}{N} \sum_n x_n \end{aligned}$$

## Bayesian inference for the Gaussian

Likelihood:

$$p(\text{Data}|\mu) = \exp\left(-\frac{N}{2\sigma^2}\mu^2 + \frac{N\bar{x}}{\sigma^2}\mu + \text{const.}\right)$$

Prior:

$$\begin{aligned} p(\mu) &= \mathcal{N}(\mu|\mu_0, \sigma_0) = \frac{1}{\sqrt{2\pi}\sigma_0} \exp\left(-\frac{1}{2\sigma_0^2}(\mu - \mu_0)^2\right) \\ &= \exp\left(-\frac{1}{2\sigma_0^2}\mu^2 + \frac{\mu_0}{\sigma_0^2}\mu + \text{const.}\right) \end{aligned}$$

$\mu_0, \sigma_0$  hyperparameters. Large  $\sigma_0$  = large prior uncertainty in  $\mu$ .

$$p(\mu|\text{Data}) \propto p(\text{Data}|\mu)p(\mu)$$

# Product of Gaussian potentials

Gaussian potential: unnormalized Gaussian

$$\begin{aligned}\phi_1(x) &= \exp\left(-\frac{a}{2}(x - b)^2 + \text{const.}\right) \\ &= \exp\left(-\frac{a}{2}x^2 + abx + \text{const.}\right) \\ \phi_2(x) &= \exp\left(-\frac{c}{2}x^2 + cdx + \text{const.}\right) \\ \Rightarrow \quad \phi_1(x) \times \phi_2(x) &= \exp\left(-\frac{a+c}{2}x^2 + (ab+cd)x + \text{const.}\right) \\ &= \exp\left(-\frac{a+c}{2}x^2 + (a+c)\frac{ab+cd}{a+c}x + \text{const.}\right) \\ &= \exp\left(-\frac{a+c}{2}\left(x - \frac{ab+cd}{a+c}\right)^2 + \text{const.}\right)\end{aligned}$$

With normalization, the product  $p(x) \propto \phi_1(x)\phi_2(x)$  is a Gaussian.

Posterior is proportional to the product of two Gaussian potentials.

$$p(\text{Data}|\mu) \propto \mathcal{N}(\mu|\bar{x}, \frac{1}{N}\sigma^2)$$

$$p(\mu) = \mathcal{N}(\mu|\mu_0, \sigma_0^2) = \mathcal{N}(\mu|\mu_0, \frac{\sigma_0^2}{\sigma^2}\sigma^2) = \mathcal{N}(\mu|\mu_0, \frac{1}{N_0}\sigma^2)$$

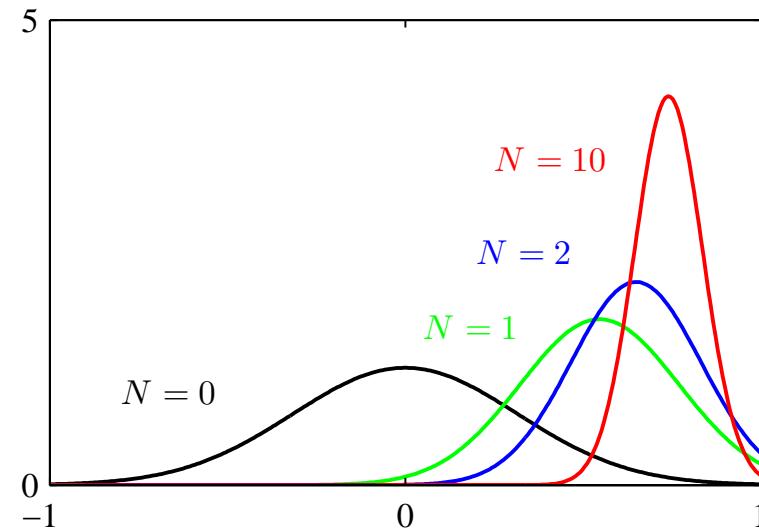
Interpretation:  $\mu_0$  mean of pseudodata;  $N_0 = \frac{\sigma^2}{\sigma_0^2}$ : effective number of pseudocounts

$$p(\mu|\text{Data}) = \mathcal{N}(\mu|\mu_N, \sigma_N^2)$$

with

$$\mu_N = \frac{\frac{N}{\sigma^2}\bar{x} + \frac{N_0}{\sigma^2}\mu_0}{\frac{N}{\sigma^2} + \frac{N_0}{\sigma^2}} = \frac{N\bar{x} + N_0\mu_0}{N + N_0}$$

$$\frac{1}{\sigma_N^2} = \frac{N}{\sigma^2} + \frac{N_0}{\sigma^2} = \frac{N + N_0}{\sigma^2}$$

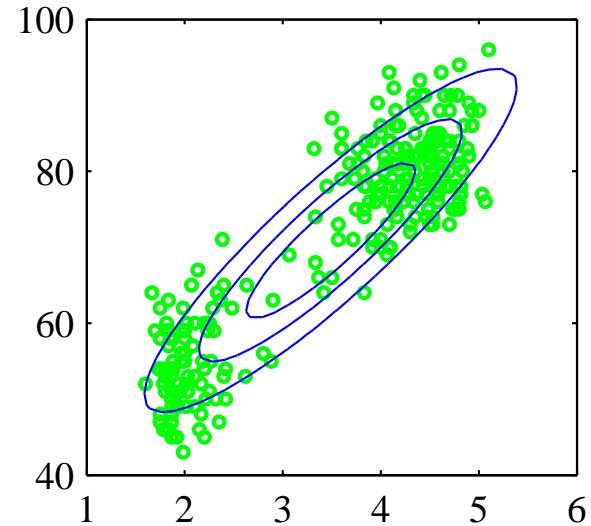


For  $N \rightarrow \infty$ :  $\mu_N \rightarrow \bar{x}, \sigma_N^2 \rightarrow 0$

i.e., posterior distribution is a peak around ML solution  
so Bayesian inference and ML coincides.

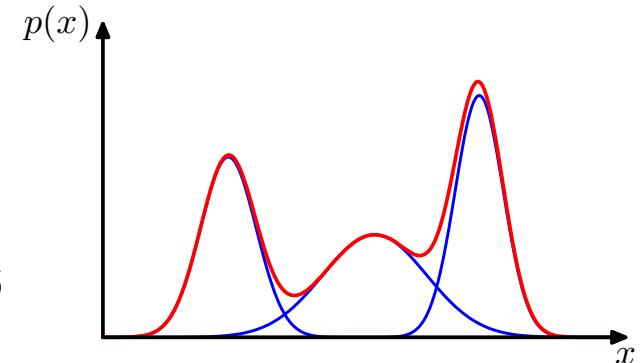
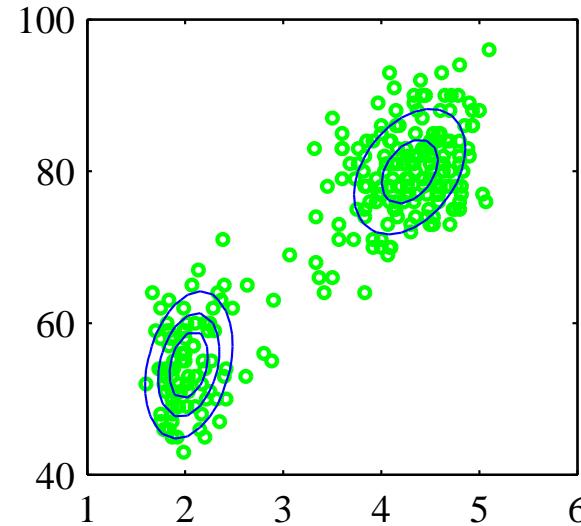
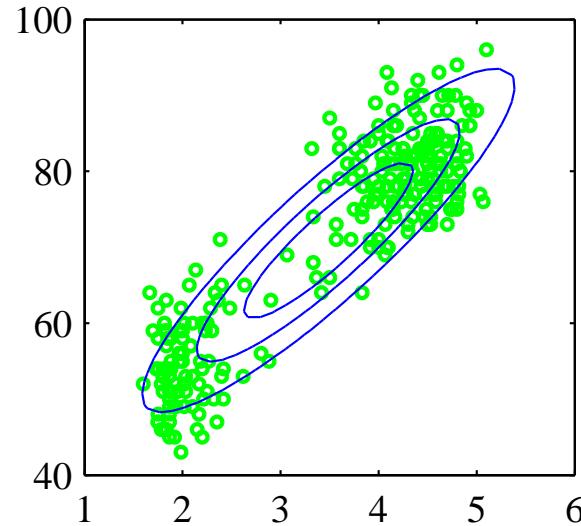
# Mixtures of Gaussians

Model for multimodal distribution



# Mixtures of Gaussians

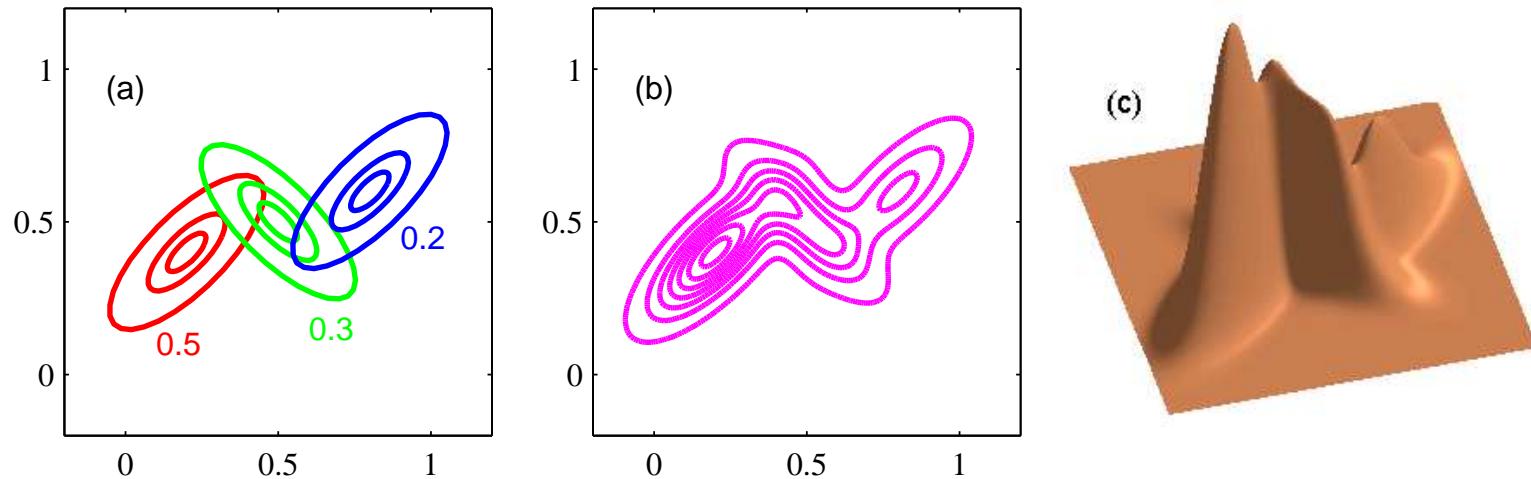
Model for multimodal distribution



$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

Each Gaussian is called a *component* of the mixture. The factors  $\pi_k$  are the *mixing coefficients*.

# Mixtures of Gaussians



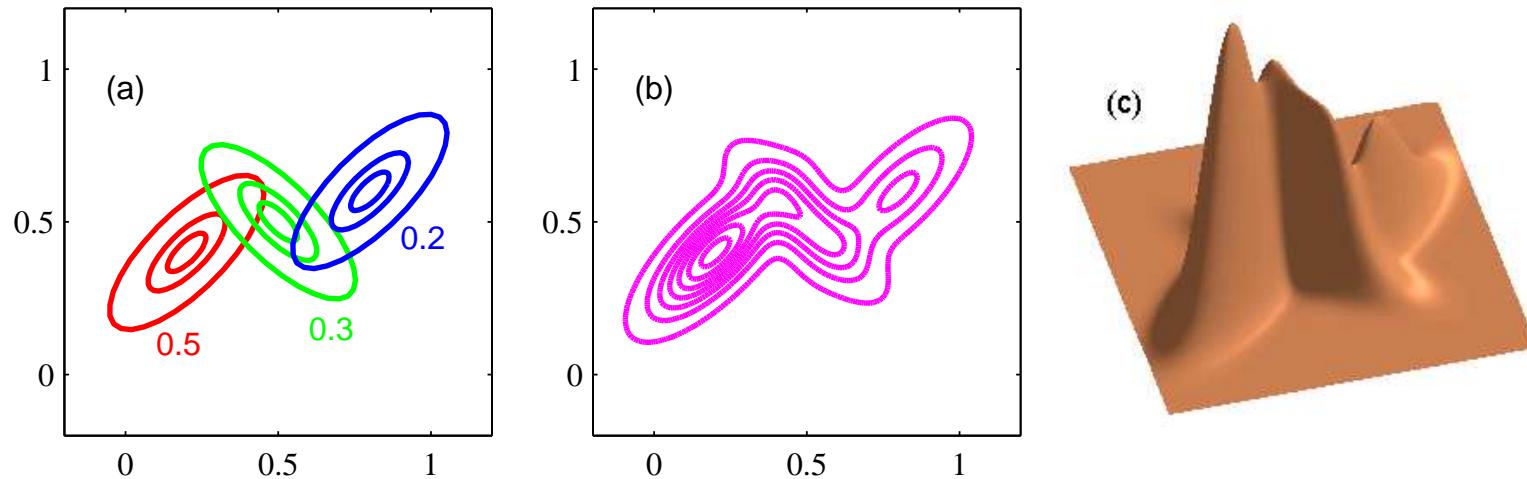
$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- Mixing coefficients satisfy

$$\pi_k \geq 0 \quad \text{and} \quad \sum_k \pi_k = 1$$

this implies  $p(\mathbf{x}) \geq 0$  and  $\int p(\mathbf{x}) d\mathbf{x} = 1$ .

# Mixtures of Gaussians



$$p(\mathbf{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

- $k$ : label of the mixture. Joint distribution  $p(\mathbf{x}, k) = \pi_k p(\mathbf{x}|k)$ . Since data is not labeled: marginal distribution  $p(\mathbf{x}) = \sum_k \pi_k p(\mathbf{x}|k)$ .  $p(\mathbf{x}, k)$  in exponential family. However, mixture model  $p(\mathbf{x})$  not in the exponential family, no simple relation between data and parameters
- One can also consider mixtures of other distributions (mixtures of Bernoulli 9.3.3)

## Mixtures of Gaussians

ML: by *gradient ascent* (numerical function maximization).

Learning with hidden variables is generally difficult (slow).

Instead, *EM (Expectation Maximization)* algorithm to deal with hidden variables.

For instance using Bayes' rule compute the "soft assignment" of data  $x_n$  to each of the clusters:

$$\gamma_k(x_n) \equiv p(k|x_n) = \frac{\pi_k \mathcal{N}(x_n|\mu_k, \Sigma_k)}{\sum_{k'} \pi_{k'} \mathcal{N}(x_n|\mu_{k'}, \Sigma_{k'})}$$

update according to

$$\mu_k = \frac{\sum_{n=1}^N \gamma_k(x_n) x_n}{\sum_{n=1}^N \gamma_k(x_n)}$$

See Chapter 9 for more details.

# The exponential family

Examples: Gaussian, Bernouilli, Beta, multinomial, Poisson distribution

$$p(\mathbf{x}|\boldsymbol{\eta}) = h(\mathbf{x})g(\boldsymbol{\eta}) \exp(\boldsymbol{\eta}^T \mathbf{u}(\mathbf{x}))$$

$\boldsymbol{\eta}$ : "natural parameters"

$\mathbf{u}(\mathbf{x})$ : (vector) function of  $\mathbf{x}$ , "sufficient statistic"

$g(\boldsymbol{\eta})$  to ensure normalization

# The exponential family

Examples: Gaussian, Bernouilli, Beta, multinomial, Poisson distribution

$$p(\mathbf{x}|\boldsymbol{\eta}) = h(\mathbf{x})g(\boldsymbol{\eta}) \exp(\boldsymbol{\eta}^T \mathbf{u}(\mathbf{x}))$$

$\boldsymbol{\eta}$ : "natural parameters"

$\mathbf{u}(\mathbf{x})$ : (vector) function of  $\mathbf{x}$ , "sufficient statistic"

$g(\boldsymbol{\eta})$  to ensure normalization

## Example: Bernoulli distribution

$$\begin{aligned} p(x|\mu) &= \text{Bern}(x|\mu) = \mu^x(1-\mu)^{(1-x)} \\ &= \exp(x \ln \mu + (1-x) \ln(1-\mu)) \\ &= (1-\mu) \exp\left(\ln\left(\frac{\mu}{1-\mu}\right)x\right) \end{aligned}$$

Natural parameters:  $\eta = \ln\left(\frac{\mu}{1-\mu}\right)$ , sufficient statistic  $u(x) = x$ .

# The exponential family

$$p(\mathbf{x}|\boldsymbol{\eta}) = h(\mathbf{x})g(\boldsymbol{\eta}) \exp(\boldsymbol{\eta}^T \mathbf{u}(\mathbf{x}))$$

## Example: Gaussian distribution

$$\begin{aligned} p(x|\mu, \sigma^2) &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) \\ &= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(\frac{\mu}{\sigma^2}x - \frac{1}{2\sigma^2}x^2 - \frac{\mu^2}{2\sigma^2}\right) \\ &= \left(\frac{1}{\sqrt{2\pi}}\right) \left(\frac{1}{\sqrt{\sigma^2}} \exp\left(-\frac{\mu^2}{2\sigma^2}\right)\right) \exp\left(\left[\frac{\mu}{\sigma^2}, -\frac{1}{2\sigma^2}\right] \cdot [x, x^2]^T\right) \end{aligned}$$

Natural parameters  $\boldsymbol{\eta} = (\frac{\mu}{\sigma^2}, -\frac{1}{2\sigma^2})^T$ , sufficient statistic  $\mathbf{u}(x) = (x, x^2)^T$ ,  $h(\mathbf{x}) = \frac{1}{\sqrt{2\pi}}$ .

Rewrite  $-\frac{\mu^2}{2\sigma^2} = \frac{\eta_1^2}{4\eta_2}$ , and  $\frac{1}{\sqrt{\sigma^2}} = \sqrt{-2\eta_2}$ , then  $g(\boldsymbol{\eta}) = \sqrt{-2\eta_2} \exp\left(\frac{\eta_1^2}{4\eta_2}\right)$ .

# The exponential family / Maximum likelihood and sufficient statistics

In the ML solution  $\boldsymbol{\eta} = \boldsymbol{\eta}_{ML}$ , the likelihood is stationary:

$$\nabla_{\boldsymbol{\eta}} \sum_n \log p(\mathbf{x}_n | \boldsymbol{\eta}) = 0$$

Since

$$\begin{aligned}\nabla_{\boldsymbol{\eta}} \log p(\mathbf{x}_n | \boldsymbol{\eta}) &= \nabla_{\boldsymbol{\eta}} \log \left( h(\mathbf{x}_n) g(\boldsymbol{\eta}) \exp(\boldsymbol{\eta}^T \mathbf{u}(\mathbf{x}_n)) \right) \\ &= \nabla_{\boldsymbol{\eta}} \log g(\boldsymbol{\eta}) + \mathbf{u}(\mathbf{x}_n)\end{aligned}$$

this implies that in ML solution  $\boldsymbol{\eta} = \boldsymbol{\eta}_{ML}$

$$-\nabla_{\boldsymbol{\eta}} \log g(\boldsymbol{\eta}) = \frac{1}{N} \sum_n \mathbf{u}(\mathbf{x}_n)$$

Now,  $g(\boldsymbol{\eta})$  is the normalization factor, i.e.,

$$g(\boldsymbol{\eta})^{-1} = \int h(\mathbf{x}) \exp(\boldsymbol{\eta}^T \mathbf{u}(\mathbf{x})) dx$$

So

$$\begin{aligned} -\nabla_{\boldsymbol{\eta}} \log g(\boldsymbol{\eta}) &= \nabla_{\boldsymbol{\eta}} \log g(\boldsymbol{\eta})^{-1} = g(\boldsymbol{\eta}) \nabla_{\boldsymbol{\eta}} \int h(\mathbf{x}) \exp(\boldsymbol{\eta}^T \mathbf{u}(\mathbf{x})) \\ &= g(\boldsymbol{\eta}) \int h(\mathbf{x}) \mathbf{u}(\mathbf{x}) \exp(\boldsymbol{\eta}^T \mathbf{u}(\mathbf{x})) \\ &= \int \mathbf{u}(\mathbf{x}) p(\mathbf{x}|\boldsymbol{\eta}) dx = \langle \mathbf{u}(\mathbf{x}) \rangle_{\boldsymbol{\eta}} \end{aligned}$$

and we have that in the ML solution:

$$\frac{1}{N} \sum_n \mathbf{u}(\mathbf{x}_n) = \langle \mathbf{u}(\mathbf{x}) \rangle_{\boldsymbol{\eta}_{ML}}$$

ML estimator depends on data only through sufficient statistics  $\sum_n \mathbf{u}(\mathbf{x}_n)$

# Conjugate priors

Likelihood:

$$\begin{aligned}
 \prod_{n=1}^N p(\mathbf{x}_n | \boldsymbol{\eta}) &= [\prod_n h(\mathbf{x}_n)] g(\boldsymbol{\eta})^N \exp(\boldsymbol{\eta}^T \sum_n \mathbf{u}(\mathbf{x}_n)) \\
 &= [\prod_n h(\mathbf{x}_n)] g(\boldsymbol{\eta})^N \exp(N \boldsymbol{\eta}^T [\frac{1}{N} \sum_n \mathbf{u}(\mathbf{x}_n)]) \\
 &= [\prod_n h(\mathbf{x}_n)] g(\boldsymbol{\eta})^N \exp(N \boldsymbol{\eta}^T \langle \mathbf{u} \rangle_{Data})
 \end{aligned}$$

Conjugate prior

$$p(\boldsymbol{\eta} | \boldsymbol{\chi}, \nu) = f(\boldsymbol{\chi}, \nu) g(\boldsymbol{\eta})^\nu \exp(\nu \boldsymbol{\eta}^T \boldsymbol{\chi})$$

in which  $\nu$ : effective number of pseudo data and  $\boldsymbol{\chi}$ : sufficient statistic.

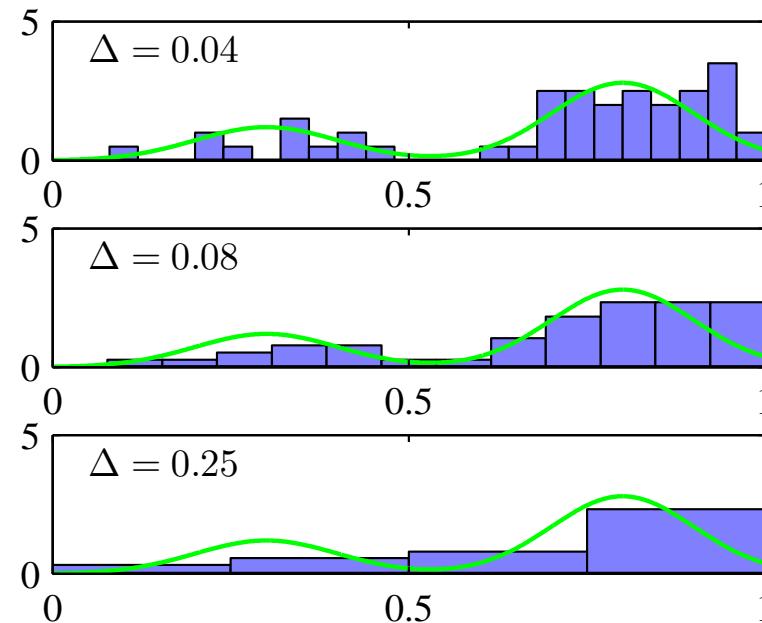
Posterior  $\sim$  likelihood  $\times$  prior:

$$P(\boldsymbol{\eta} | \mathbf{X}, \boldsymbol{\chi}, \nu) \propto g(\boldsymbol{\eta})^{N+\nu} \exp(\boldsymbol{\eta}^T (N \langle \mathbf{u} \rangle_{Data} + \nu \boldsymbol{\chi}))$$

# Nonparametric methods/histograms

"Histograms":  $p(x) = p_i$  in bin  $i$  with width  $\Delta_i$ , then ML:

$$p_i = \frac{n_i}{N\Delta_i}$$



- number of bins exponential in dimension  $D$
- discontinuous

## Nonparametric methods/Kernels

Assume true density is  $p(\mathbf{x})$ , then the probability to find a data point in region  $R$  is

$$P = \int_R p(\mathbf{x}) d\mathbf{x}$$

If the total number of data points  $N$  is large, the number of data points that are found in  $R$  is  $K$ , which is approximately

$$K \simeq NP$$

On the other hand, if the region  $R$  is small compared to the variation in  $p$ , (i.e.  $p(\mathbf{x})$  approximately constant in  $R$ ),

$$P \simeq p(\mathbf{x})V$$

where  $V$  is the volume of  $R$ . Combining gives

$$p(\mathbf{x}) = \frac{K}{NV}$$

NB: validity requires  $R$  large to estimate  $P$ , and  $R$  small to have  $p(x)$  constant.

# Nonparametric methods

$$p(\mathbf{x}) = \frac{K}{NV}$$

- Fix  $V$ , determine  $K$  by data: kernel approach/Parzen window
- Fix  $K$ , determine  $V$  by data:  $K$ -nearest neighbour

Terminology:

- Parametric distribution: model given by a number of parameters  $p(\mathbf{x}|\theta)$
- Nonparametric distribution: number of parameters grows with the data.

## Kernel method

Kernel function:

$$k(\mathbf{u}) = \begin{cases} 1, & |u_i| \leq 1/2 \\ 0, & \text{otherwise} \end{cases}$$

$k((\mathbf{x} - \mathbf{x}_n)/h)$ : one if  $\mathbf{x}_n$  in cube of side  $h$  centered around  $\mathbf{x}$ . Total number of data points lying in this cube:

$$K = \sum_{n=1}^N k\left(\frac{\mathbf{x} - \mathbf{x}_n}{h}\right)$$

Volume is  $V = h^D$ , so

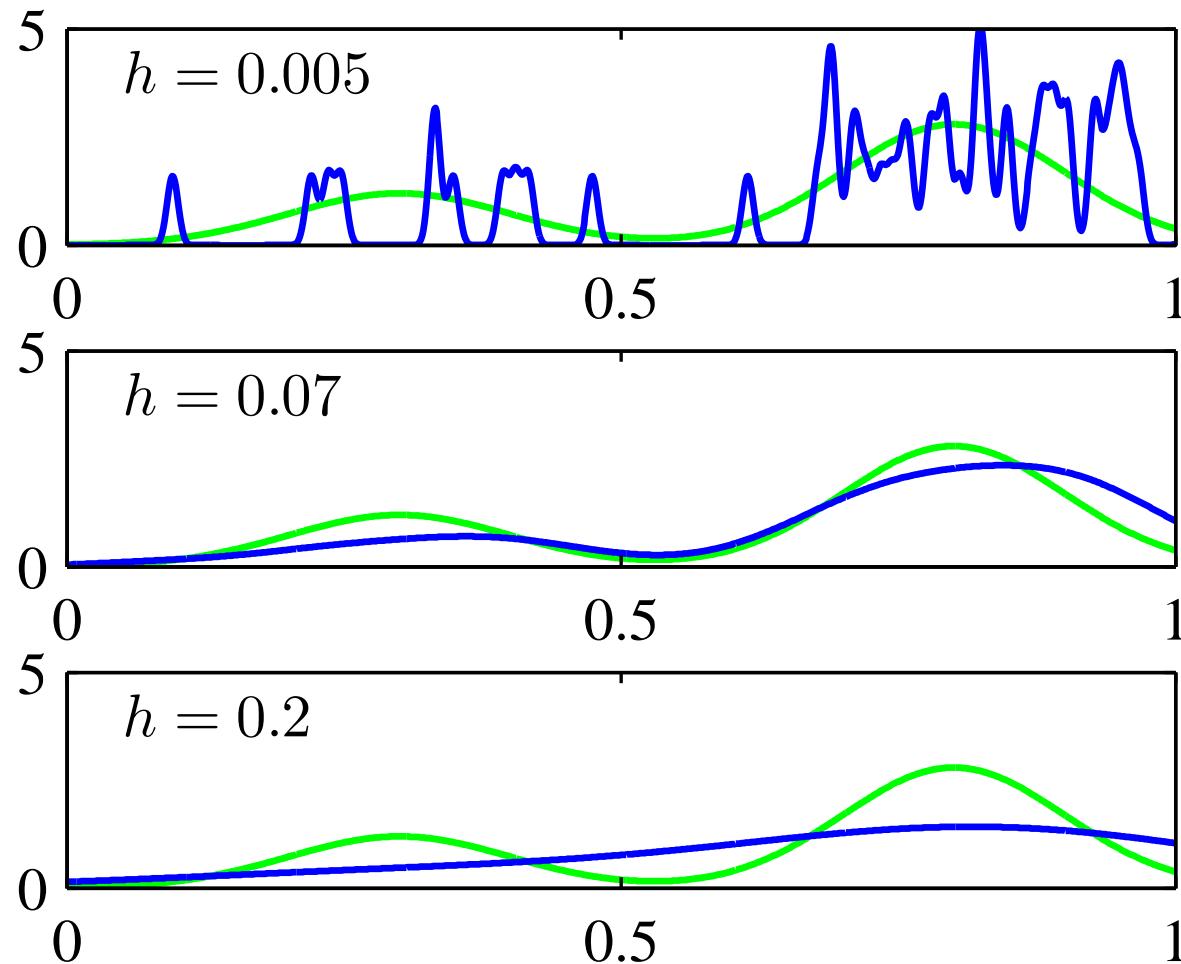
$$p(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^N \frac{1}{h^D} k\left(\frac{\mathbf{x} - \mathbf{x}_n}{h}\right)$$

- Trick works for any  $k(\mathbf{u})$  with  $\int k(\mathbf{u}) d\mathbf{u} = 1$ .

For example, Gaussian kernel  $k(\mathbf{u}) = \frac{1}{\sqrt{2\pi h^2}} \exp(-\|\mathbf{u}\|^2/2h^2)$ .

- $h$  is smoothing parameter.

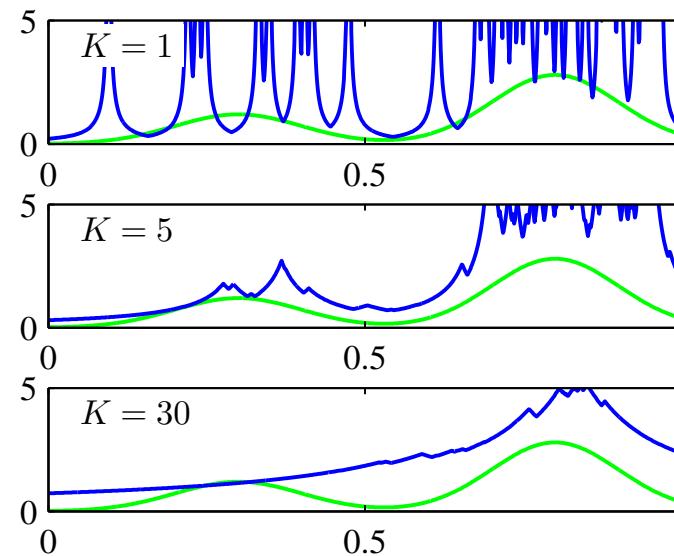
# Kernel methods



# Nearest neighbour methods

$$p(x) = \frac{K}{NV}$$

For instance with  $K = 1$  and one data point at  $y$  in one dimension, we obtain  $p(x) = \frac{1}{2N|y-x|}$ .



NB: not a true density model (integral over  $x$  diverges)

# Nearest Neighbor classification

Training set  $N_k$  points in class  $C_k$ ,  $\sum_k N_k = N$

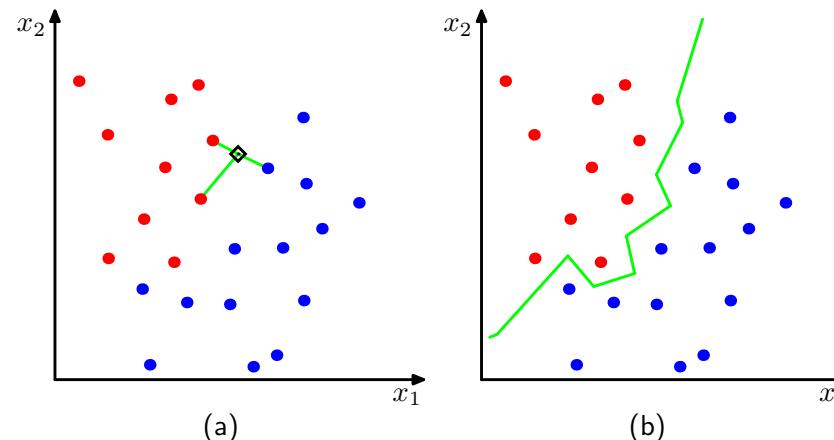
Consider test point  $\mathbf{x}$  and its  $K$  nearest neighbours in volume  $V$ . Then

$$p(\mathbf{x}|C_k) = \frac{K_k}{N_k V} \quad p(\mathbf{x}) = \frac{K}{NV} \quad p(C_k) = \frac{N_k}{N}$$

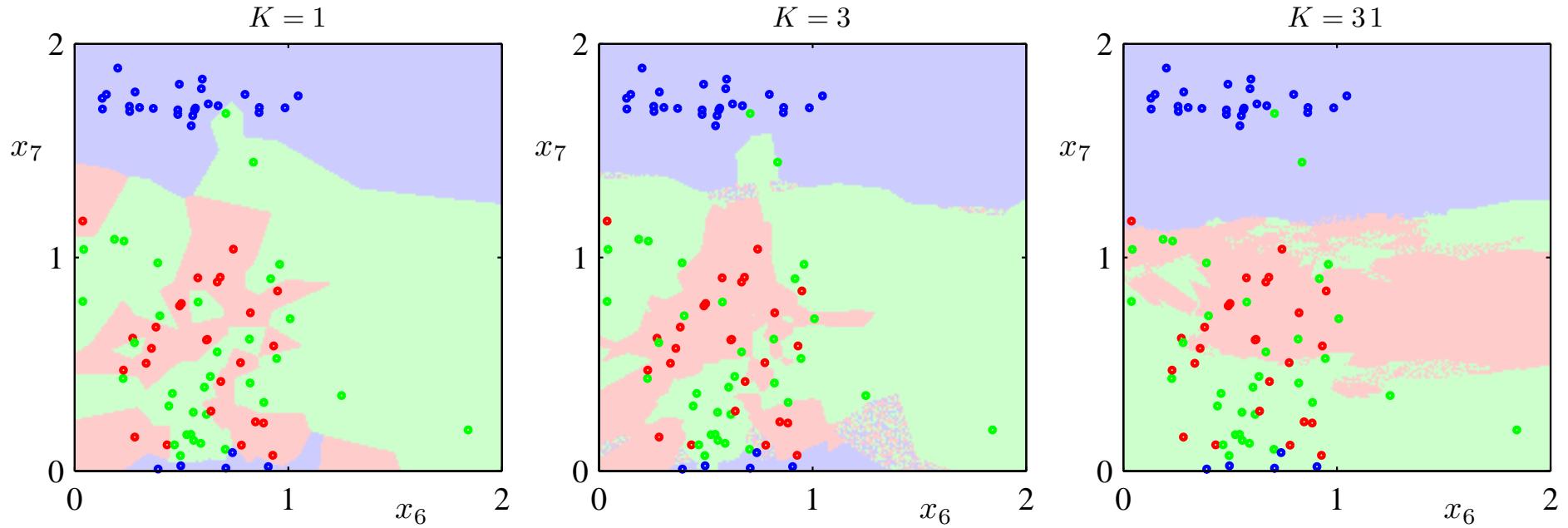
then by Bayes' theorem, posterior probability of class membership follows,

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

$K$  (number of neighbours) is smoothing parameter.



# Nearest neighbour methods



# Chapter 3

# Linear Models for Regression

Regression: predicting the value of *continuous* target/output variables given values of the input variables.

In other words: given a training set of input/output pairs, constructing a function that maps input values to continuous output values, like the polynomial curve fitting in §1.1.

# Chapter 3

# Linear Models for Regression

Regression: predicting the value of *continuous* target/output variables given values of the input variables.

In other words: given a training set of input/output pairs, constructing a function that maps input values to continuous output values, like the polynomial curve fitting in §1.1.

This chapter: generalized approach to *linear* models for regression

- more flexible models
- exploit probabilistic noise models (chapter 2!)
- fully Bayesian predictive distributions
- connection to model selection

# Linear Basis Function Models

'Very' linear regression:

$$y(\boldsymbol{x}, \boldsymbol{w}) = w_0 + w_1x_1 + w_2x_2 + \cdots + w_Dx_D = w_0 + \sum_{j=1}^D w_jx_j$$

Linear in parameters  $\boldsymbol{w}$ , and "almost" linear (affine) in input  $\boldsymbol{x}$ .

Most functions are not linear in input  $\boldsymbol{x}$  ... can we generalize our model to that?

# Linear Basis Function Models

'Very' linear regression:

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1x_1 + w_2x_2 + \cdots + w_Dx_D = w_0 + \sum_{j=1}^D w_jx_j$$

Linear in parameters  $\mathbf{w}$ , and "almost" linear (affine) in input  $\mathbf{x}$ .

Most functions are not linear in input  $\mathbf{x}$  ... can we generalize our model to that?

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^{M-1} w_j\phi_j(\mathbf{x})$$

where  $\phi_1(\mathbf{x}), \dots, \phi_{M-1}(\mathbf{x})$  are *basis functions* or *feature functions*. Defining  $\phi_0(\mathbf{x}) = 1$ , we can write it more compactly:

$$y(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^{M-1} w_j\phi_j(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$$

Note: still linear in parameters  $\mathbf{w}$ ! Examples of basis functions?

# Linear Basis Function Models

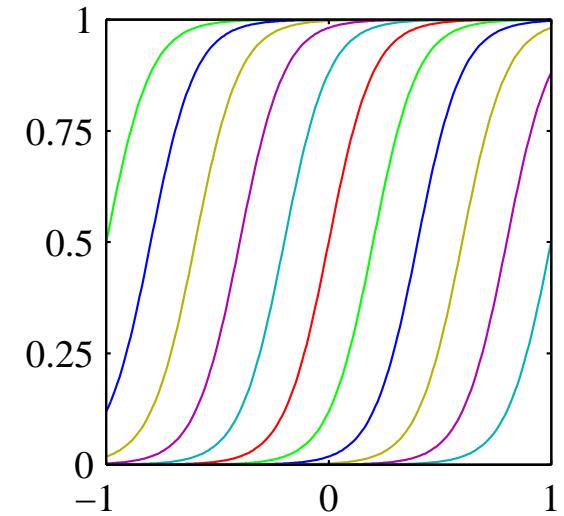
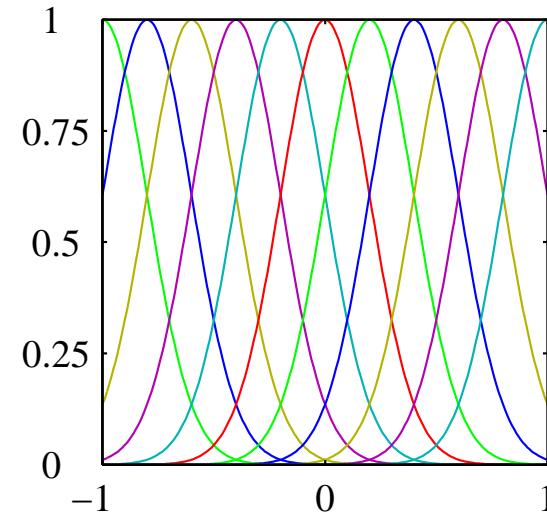
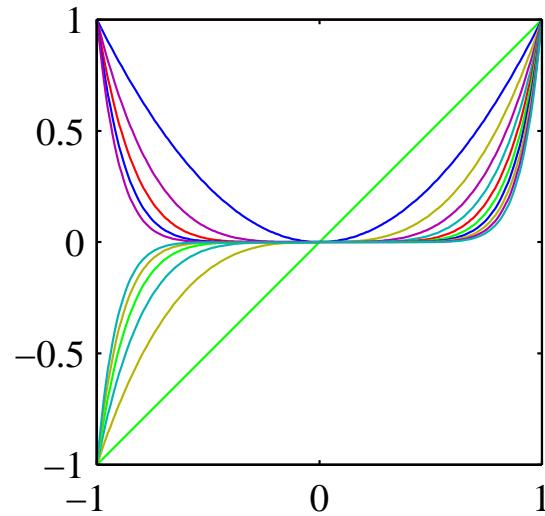
Gaussians:  $\phi_j(x) = \exp(-\beta(x - \mu_j)^2)$

Polynomials:  $\phi_j(x) = x^j$

Sigmoids:  $\phi_j(x) = \sigma(\beta(x - \mu_j))$

Fourier:  $\phi_j(x) = \exp(i j x)$

Wavelets: ...



## Adding some noise

Let's introduce some noise . . . and assume the following model:

$$t = y(\mathbf{x}, \mathbf{w}) + \epsilon, \quad y(\mathbf{x}, \mathbf{w}) = \mathbf{w}^T \phi(\mathbf{x})$$

where we added Gaussian noise  $\epsilon \sim \mathcal{N}(0, \beta^{-1})$ .

## Adding some noise

Let's introduce some noise . . . and assume the following model:

$$t = y(\mathbf{x}, \mathbf{w}) + \epsilon, \quad y(\mathbf{x}, \mathbf{w}) = \mathbf{w}^T \phi(\mathbf{x})$$

where we added Gaussian noise  $\epsilon \sim \mathcal{N}(0, \beta^{-1})$ .

Some typical exam questions:

Q: what is the distribution of  $t$ , given  $\mathbf{x}, \mathbf{w}, \beta$ ?

## Adding some noise

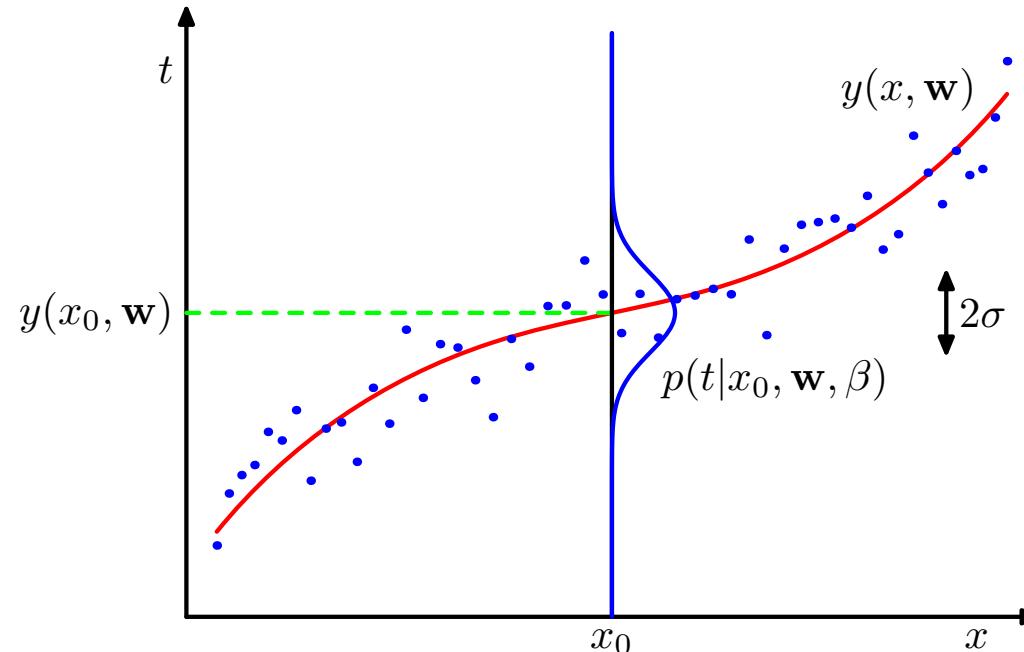
Let's introduce some noise . . . and assume the following model:

$$t = y(\mathbf{x}, \mathbf{w}) + \epsilon, \quad y(\mathbf{x}, \mathbf{w}) = \mathbf{w}^T \phi(\mathbf{x})$$

where we added Gaussian noise  $\epsilon \sim \mathcal{N}(0, \beta^{-1})$ .

Q: what is the distribution of  $t$ , given  $\mathbf{x}, \mathbf{w}, \beta$ ?

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



## Adding some noise

Let's introduce some noise . . . and assume the following model:

$$t = y(\mathbf{x}, \mathbf{w}) + \epsilon, \quad y(\mathbf{x}, \mathbf{w}) = \mathbf{w}^T \phi(\mathbf{x})$$

where we added Gaussian noise  $\epsilon \sim \mathcal{N}(0, \beta^{-1})$ .

Q: what is the conditional mean  $\mathbb{E}(t|\mathbf{x})$ ?

## Adding some noise

Let's introduce some noise . . . and assume the following model:

$$t = y(\mathbf{x}, \mathbf{w}) + \epsilon, \quad y(\mathbf{x}, \mathbf{w}) = \mathbf{w}^T \phi(\mathbf{x})$$

where we added Gaussian noise  $\epsilon \sim \mathcal{N}(0, \beta^{-1})$ .

Q: what is the conditional mean  $\mathbb{E}(t|\mathbf{x})$ ?

$$\mathbb{E}(t|\mathbf{x}) = \int t p(t|\mathbf{x}) dt = y(\mathbf{x}, \mathbf{w})$$

## Maximum likelihood: least squares

Q: What is the log-likelihood of a data set  $\{\mathbf{X}, \mathbf{t}\} = \{(\mathbf{x}_i, t_i)\}_{i=1}^N$  in this model?

## Maximum likelihood: least squares

Q: What is the log-likelihood of a data set  $\{\mathbf{X}, \mathbf{t}\} = \{(\mathbf{x}_i, t_i)\}_{i=1}^N$  in this model?

$$\begin{aligned}
 \ln p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta) &= \ln \prod_{i=1}^N \mathcal{N}(t_n | \mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1}) \\
 &= \sum_{i=1}^N \ln \mathcal{N}(t_n | \mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1}) \\
 &= \frac{N}{2} \ln \beta - \frac{1}{2} \beta \underbrace{\sum_{i=1}^N (t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2}_{\text{sum-of-squares error function}} - C
 \end{aligned}$$

Q: How to find the maximum likelihood solution for the parameters  $\mathbf{w}$ ?

## Maximum likelihood: least squares

Q: What is the log-likelihood of a data set  $\{\mathbf{X}, \mathbf{t}\} = \{(\mathbf{x}_i, t_i)\}_{i=1}^N$  in this model?

$$\begin{aligned}
 \ln p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) &= \ln \prod_{i=1}^N \mathcal{N}(t_n | \mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1}) \\
 &= \sum_{i=1}^N \ln \mathcal{N}(t_n | \mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1}) \\
 &= \frac{N}{2} \ln \beta - \frac{1}{2} \beta \underbrace{\sum_{i=1}^N (t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2}_{\text{sum-of-squares error function}} - C
 \end{aligned}$$

Q: How to find the maximum likelihood solution for the parameters  $\mathbf{w}$ ?

A: Differentiate with respect to  $\mathbf{w}$ , set to zero, and solve equations to find  $\mathbf{w}_{ML}$ :

$$\nabla_{\mathbf{w}} \ln p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) = -\beta \sum_{i=1}^N (t_n - \mathbf{w}^T \phi(\mathbf{x}_n)) \phi(\mathbf{x}_n)^T = \mathbf{0}$$

## Maximum likelihood: least squares

Rewriting:

$$\mathbf{0} = \sum_{n=1}^N t_n \phi(\mathbf{x}_n)^T - \mathbf{w}^T \left( \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T \right)$$

Solving for  $\mathbf{w}$ :

$$\mathbf{w}_{ML} = \underbrace{(\Phi^T \Phi)^{-1} \Phi^T}_{\text{Moore-Penrose pseudo-inverse}} \mathbf{t}$$

where  $\Phi$  is an  $N \times M$  matrix, the *design matrix*, with elements  $\Phi_{nj} = \phi_j(\mathbf{x}_n)$ :

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \dots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & \dots & \phi_{M-1}(\mathbf{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \dots & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}$$

Maximizing  $p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta)$  w.r.t.  $\beta$  gives:

$$\beta_{ML}^{-1} = \frac{1}{N} \sum_{n=1}^N (t_n - \mathbf{w}_{ML}^T \phi(\mathbf{x}_n))^2.$$

## Regularized least squares

To avoid overfitting, we can add a regularization term to the log-likelihood, e.g. *weight decay*:

$$\ln p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) + \lambda E_W(\mathbf{w}) = \frac{N}{2} \ln \beta - \frac{1}{2} \beta \underbrace{\sum_{i=1}^N (t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2}_{\text{sum-of-squares error function}} + \underbrace{\frac{1}{2} \lambda \mathbf{w}^T \mathbf{w} - C}_{\text{regularizer}}$$

Maximizing with respect to  $\mathbf{w}$  now gives the following optimum:

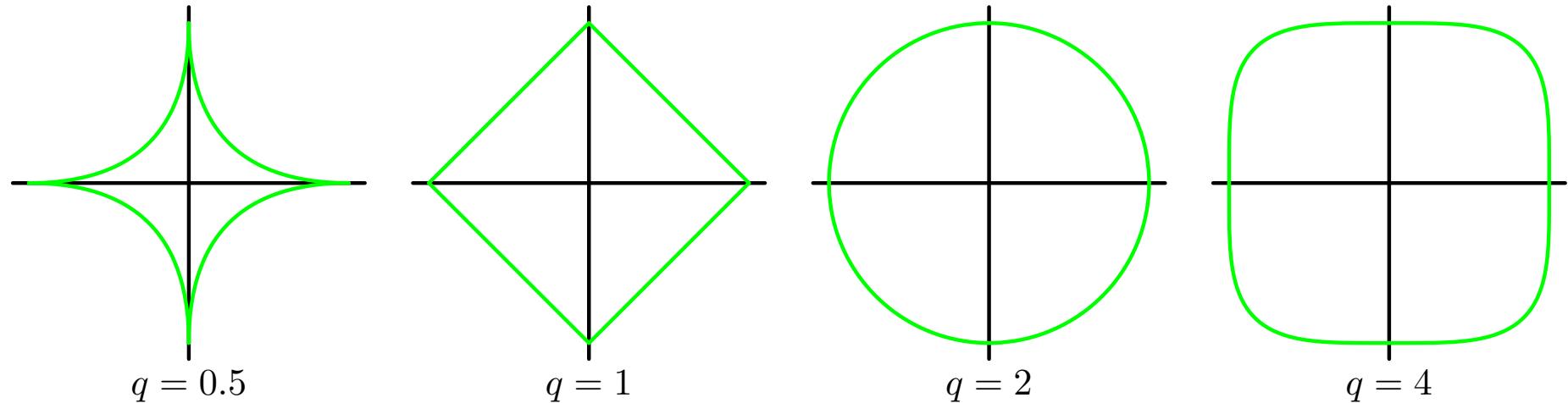
$$\mathbf{w} = (\lambda \mathbf{I} + \Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

Other regularizers also possible, e.g. a more general error term would be

$$\frac{1}{2} \sum_{i=1}^N (t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2 + \frac{\lambda}{2} \sum_{j=1}^M |w_j|^q$$

## Regularized least squares

For example  $q = 1$  gives the famous LASSO with regularizer  $\lambda \sum_{j=1}^M |w_j|$ .



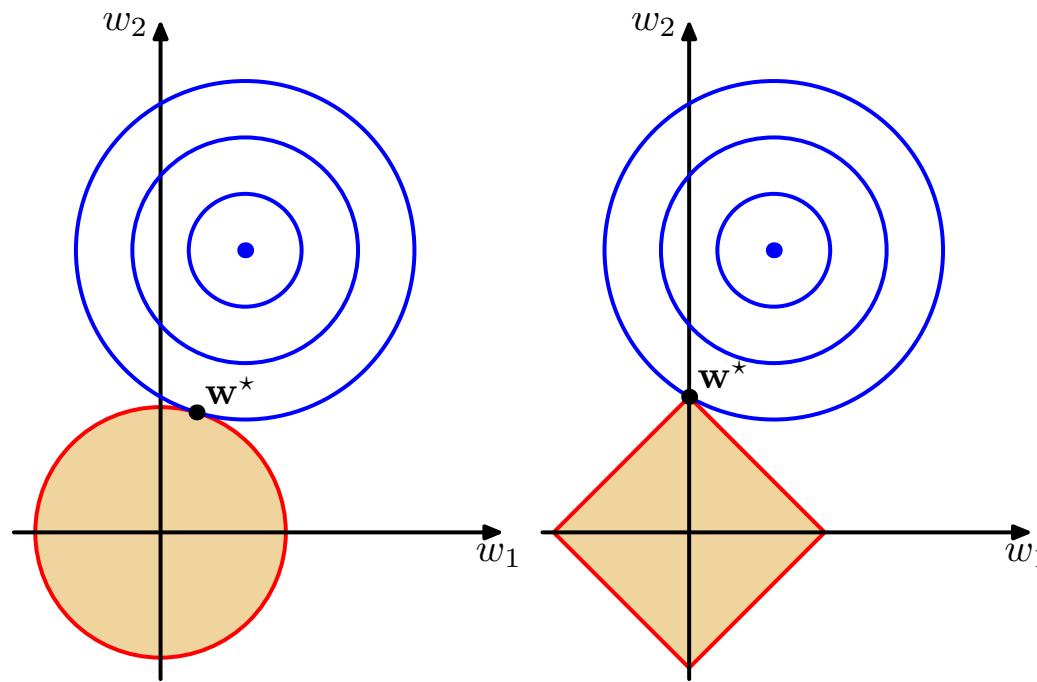
Q: Why does the LASSO result in a *sparse* solution ( $w_j \neq 0$  for only a few  $j$ 's)?

# Regularized least squares

A: The solution  $\mathbf{w}$  to the general regularized error term

$$\frac{1}{2} \sum_{i=1}^N (t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2 + \frac{\lambda}{2} \sum_{j=1}^M |w_j|^q$$

can be viewed as the *unregularized* solution under constraint  $\sum_{j=1}^M |w_j|^q \leq \eta$ .



## Bias-Variance Decomposition

Complex models tend to overfit. Simple models tend to be too rigid. Number of terms in polynomial, weight decay constant  $\lambda$ .

(Treatment of 1.5.5). Consider a regression problem with squared loss function

$$\mathcal{E}(L) = \iint (y(\mathbf{x}) - t)^2 p(\mathbf{x}, t) d\mathbf{x} dt$$

$p(\mathbf{x}, t)$  is the true underlying model,  $y(\mathbf{x})$  is our estimate of  $t$  at  $\mathbf{x}$ .

The optimal  $y$  minimizes  $\mathcal{E}(L)$ :

$$\begin{aligned}\frac{\delta \mathcal{E}(L)}{\delta y(\mathbf{x})} &= 2 \int (y(\mathbf{x}) - t)p(\mathbf{x}, t) dt = 0 \\ y(\mathbf{x}) &= \frac{\int tp(\mathbf{x}, t) dt}{p(\mathbf{x})} = \int tp(t|\mathbf{x}) dt = \mathbb{E}(t|\mathbf{x})\end{aligned}$$

## Bias-Variance Decomposition

The expected squared loss can also be written as

$$\begin{aligned}
 \mathcal{E}(L) &= \iint (y(\mathbf{x}) - t)^2 p(\mathbf{x}, t) d\mathbf{x} dt \\
 &= \iint (y(\mathbf{x}) - \mathbb{E}(t|\mathbf{x}) + \mathbb{E}(t|\mathbf{x}) - t)^2 p(\mathbf{x}, t) d\mathbf{x} dt \\
 &= \iint \left( (y(\mathbf{x}) - \mathbb{E}(t|\mathbf{x}))^2 + (\mathbb{E}(t|\mathbf{x}) - t)^2 \right. \\
 &\quad \left. + 2(y(\mathbf{x}) - \mathbb{E}(t|\mathbf{x}))(\mathbb{E}(t|\mathbf{x}) - t) \right) p(\mathbf{x}, t) d\mathbf{x} dt \\
 &= \int (y(\mathbf{x}) - \mathbb{E}(t|\mathbf{x}))^2 p(\mathbf{x}) d\mathbf{x} + \underbrace{\iint (\mathbb{E}(t|\mathbf{x}) - t)^2 p(\mathbf{x}, t) d\mathbf{x} dt}_{\text{intrinsic noise}}
 \end{aligned}$$

First term depends on  $y(\mathbf{x})$  and is minimized when  $y(\mathbf{x}) = \mathbb{E}(t|\mathbf{x})$ . The second term, the variance in the conditional distribution of  $t$  given  $\mathbf{x}$ , averaged over  $\mathbf{x}$ :  $\int \text{var}(t|\mathbf{x}) p(\mathbf{x}) d\mathbf{x}$  is *independent* of the solution  $y(\mathbf{x})$ . The cross-term vanishes on integration over  $t$  (odd function around  $\mathbb{E}(t|\mathbf{x})$ ).

## Bias-Variance Decomposition

In the case of a finite data set  $\mathcal{D}$ , our learning algorithm will give a solution  $y(\mathbf{x}; \mathcal{D})$  that will not minimize the exact expected square loss (but an empirical square loss): different data sets will give different solutions  $y(\mathbf{x}; \mathcal{D})$ .

Consider the thought experiment that a large number of data sets  $\mathcal{D}$  are given. Then we can construct the average solution  $\mathbb{E}_{\mathcal{D}}(y(\mathbf{x}; \mathcal{D}))$ , and

$$\begin{aligned} (y(\mathbf{x}; \mathcal{D}) - \mathbb{E}(t|\mathbf{x}))^2 &= (y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}(y) + \mathbb{E}_{\mathcal{D}}(y) - \mathbb{E}(t|\mathbf{x}))^2 \\ &= (y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}(y))^2 + (\mathbb{E}_{\mathcal{D}}(y) - \mathbb{E}(t|\mathbf{x}))^2 \\ &\quad + 2(y - \mathbb{E}_{\mathcal{D}}(y))(\mathbb{E}_{\mathcal{D}}(y) - \mathbb{E}(t|\mathbf{x})) \end{aligned}$$

So for a given  $\mathbf{x}$  the average of this quantity over many data sets

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

as the cross term vanishes on averaging.

## Bias-Variance Decomposition

Substitution of the previous expression in the expected square loss:

$$\begin{aligned}
 \mathcal{E}(L) &= \underbrace{\int (\mathbb{E}_{\mathcal{D}}(y) - \mathbb{E}(t|\mathbf{x}))^2 p(\mathbf{x}) d\mathbf{x}}_{\text{(bias)}^2} \\
 &+ \underbrace{\int \mathbb{E}_{\mathcal{D}} [\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}(y)\}^2] p(\mathbf{x}) d\mathbf{x}}_{\text{variance}} \\
 &+ \underbrace{\iint (\mathbb{E}(t|\mathbf{x}) - t)^2 p(\mathbf{x}, t) d\mathbf{x} dt}_{\text{noise}}
 \end{aligned}$$

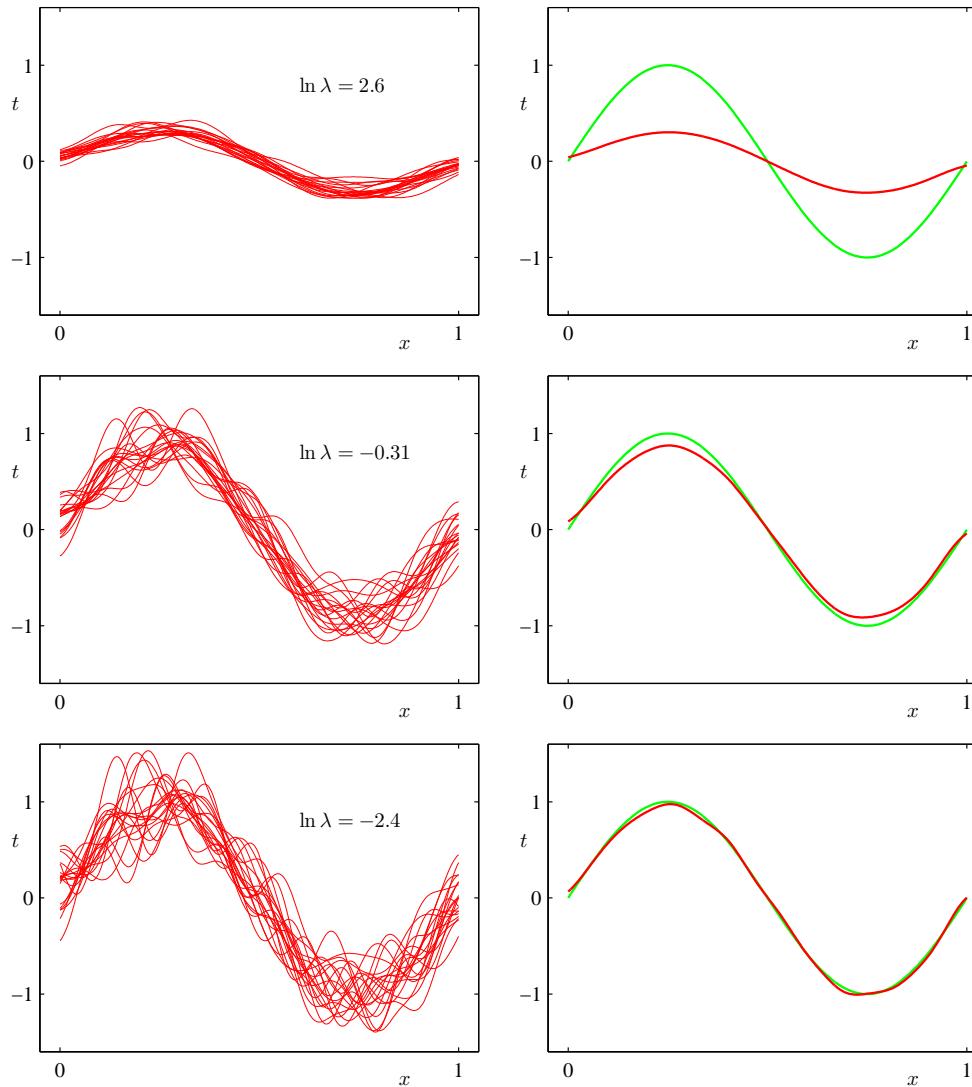
Expected square loss = bias<sup>2</sup> + variance + noise

Bias<sup>2</sup>: difference between average solution  $\mathbb{E}_{\mathcal{D}}(y(\mathbf{x}; \mathcal{D}))$  and true solution  $\mathbb{E}(t|\mathbf{x})$ .

Variance: scatter of individual solutions  $y(\mathbf{x}, \mathcal{D})$  around their mean  $\mathbb{E}_{\mathcal{D}}(y(\mathbf{x}; \mathcal{D}))$ .

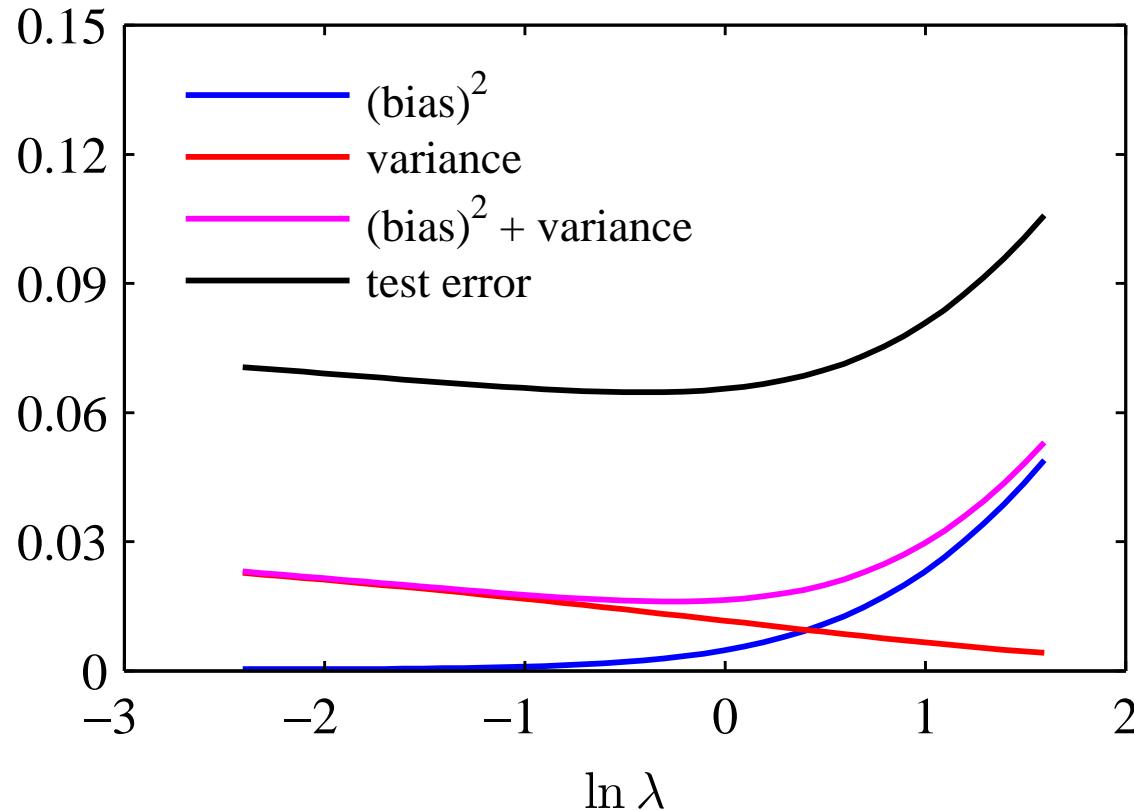
Noise: (true) scatter of the data points  $t$  around their mean  $\mathbb{E}(t|\mathbf{x})$ .

# Bias-Variance Decomposition



100 data sets, each with 25 data points from  $t = \sin(2\pi x) + \text{noise}$ .  $y(x)$  as in Eq. 3.3-4. Parameters optimized using Eq. 3.27 for different  $\lambda$ . Left shows variance, right shows bias.

# Bias-Variance Decomposition



Sum of bias, variance and noise yields expected error on test set. Optimal  $\lambda$  is trade-off between bias and variance.

True bias is usually not known, because  $\mathbb{E}(t|\mathbf{x})$  is not known.

# Bayesian linear regression

Let's go back to the linear basis function model:

$$t = y(\mathbf{x}, \mathbf{w}) + \epsilon, \quad y(\mathbf{x}, \mathbf{w}) = \mathbf{w}^T \phi(\mathbf{x}), \quad \epsilon \sim \mathcal{N}(0, \beta^{-1})$$

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

Training data:  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$  and  $\mathbf{t} = (t_1, \dots, t_N)$ .

Likelihood:

$$p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | \mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1}) = \mathcal{N}(\mathbf{t} | \Phi \mathbf{w}, \beta^{-1} \mathbf{I}).$$

Q: what prior  $p(\mathbf{w})$  can we choose?

# Bayesian linear regression

Let's go back to the linear basis function model:

$$t = y(\mathbf{x}, \mathbf{w}) + \epsilon, \quad y(\mathbf{x}, \mathbf{w}) = \mathbf{w}^T \phi(\mathbf{x}), \quad \epsilon \sim \mathcal{N}(0, \beta^{-1})$$

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

Training data:  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)$  and  $\mathbf{t} = (t_1, \dots, t_N)$ .

Likelihood:

$$p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | \mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1}) = \mathcal{N}(\mathbf{t} | \Phi \mathbf{w}, \beta^{-1} \mathbf{I}).$$

Q: what prior  $p(\mathbf{w})$  can we choose?

A: We make life easy by choosing a *conjugate* prior, which is a Gaussian

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_0, S_0)$$

## Bayesian linear regression

Then, the posterior will also be Gaussian. Prior and likelihood:

$$\begin{aligned} p(\mathbf{w}) &= \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0) \\ p(\mathbf{t}|\mathbf{w}) &= \mathcal{N}(\mathbf{t} | \Phi\mathbf{w}, \beta^{-1}\mathbf{I}) \end{aligned}$$

Then, by applying 2.113 + 2.114  $\Rightarrow$  2.116, we get the posterior  $p(\mathbf{w}|\mathbf{t})$ .

## Bayesian linear regression

Then, the posterior will also be Gaussian. Prior and likelihood:

$$\begin{aligned} p(\mathbf{w}) &= \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0) \\ p(\mathbf{t}|\mathbf{w}) &= \mathcal{N}(\mathbf{t} | \Phi\mathbf{w}, \beta^{-1}\mathbf{I}) \end{aligned}$$

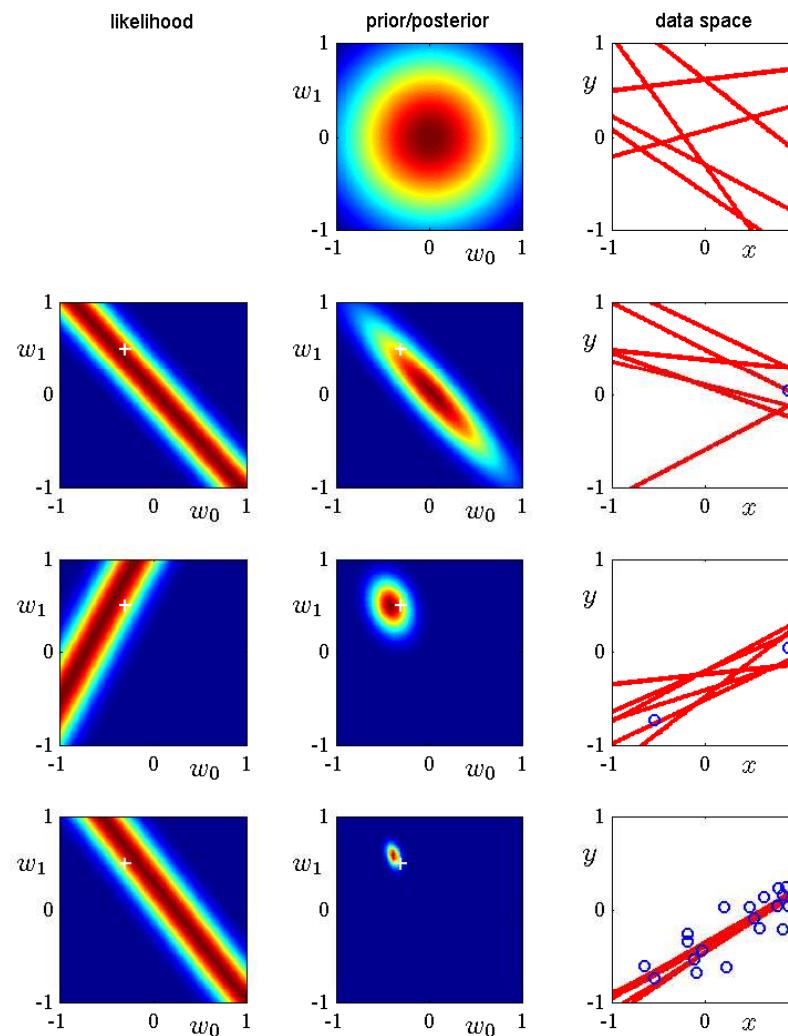
Then, by applying 2.113 + 2.114  $\Rightarrow$  2.116, we get the posterior  $p(\mathbf{w}|\mathbf{t})$ .

$$\begin{aligned} p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0) &\leftrightarrow p(\mathbf{x}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) \\ p(\mathbf{t}|\mathbf{w}) = \mathcal{N}(\mathbf{t} | \Phi\mathbf{w}, \beta^{-1}\mathbf{I}) &\leftrightarrow p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y} | \mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1}) \\ p(\mathbf{w}|\mathbf{t}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_N, \mathbf{S}_N) &\leftrightarrow p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\Sigma} \{ \mathbf{A}^T \mathbf{L} (\mathbf{y} - \mathbf{b}) + \boldsymbol{\Lambda} \boldsymbol{\mu} \}, \boldsymbol{\Sigma}) \\ &\quad \text{with } \boldsymbol{\Sigma} = (\boldsymbol{\Lambda} + \mathbf{A}^T \mathbf{L} \mathbf{A})^{-1} \end{aligned}$$

$$\begin{aligned} \text{So } p(\mathbf{w}|\mathbf{t}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_N, \mathbf{S}_N) \quad , \text{ with } \mathbf{m}_N &= \mathbf{S}_N (\mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \Phi^T \mathbf{t}) \\ \mathbf{S}_N^{-1} &= \mathbf{S}_0^{-1} + \beta \Phi^T \Phi \end{aligned}$$

Note behaviour when  $\mathbf{S}_0^{-1} \rightarrow \mathbf{0}$  (broad prior), when  $N = 0$ , and when  $N \rightarrow \infty$ .

# Bayesian linear regression



Data:  $t = a_0 + a_1x + \text{noise}$ . Model:  $y(x, \mathbf{w}) = w_0 + w_1x$ ;  $p(t|x, \mathbf{w}, \beta) = \mathcal{N}(t|y(x, \mathbf{w}), \beta^{-1})$ ,  $\beta^{-1} = (0.2)^2$ ;  $p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}|0, \alpha^{-1}\mathbf{I})$ ,  $\alpha = 2$ . When  $N$  large,  $S_N^{-1}$  is dominated by the

second term and becomes infinite. Thus, the width of the posterior goes to zero. Also the mean becomes independent of  $\beta$ , as it cancels with the  $S_N$  term.

## Predictive distribution

What is the predictive distribution  $p(t^*|x^*)$  for new data point  $x^*$ ? We know:

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

$$p(\mathbf{w}|\mathbf{t}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_N, S_N)$$

Q: How to compute  $p(t^*|x^*)$ ?

## Predictive distribution

What is the predictive distribution  $p(t^*|x^*)$  for new data point  $x^*$ ? We know:

$$\begin{aligned} p(t^*|\mathbf{w}, x^*) &= \mathcal{N}(t^*|\mathbf{w}^T \phi(x^*), \beta^{-1}) \\ p(\mathbf{w}|\mathbf{t}) &= \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N) \end{aligned}$$

A: View  $p(\mathbf{w}|\mathbf{t})$  as an inferred *prior* on  $\mathbf{w}$  (conditioned on  $\mathbf{t}$ , and implicitly  $\mathbf{x}$  as well) and apply 2.113 + 2.114  $\Rightarrow$  2.115 to compute the *marginal* on  $t^*$  (given  $\mathbf{t}$  and  $x^*$ ):

$$\begin{aligned} p(\mathbf{w}|\mathbf{t}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N) &\leftrightarrow p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) \\ p(t^*|\mathbf{w}, x^*) = \mathcal{N}(t^*|\mathbf{w}^T \phi(x^*), \beta^{-1}) &\leftrightarrow p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{x}|\mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1}) \\ p(t^*|x^*) &\leftrightarrow p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^T) \end{aligned}$$

So

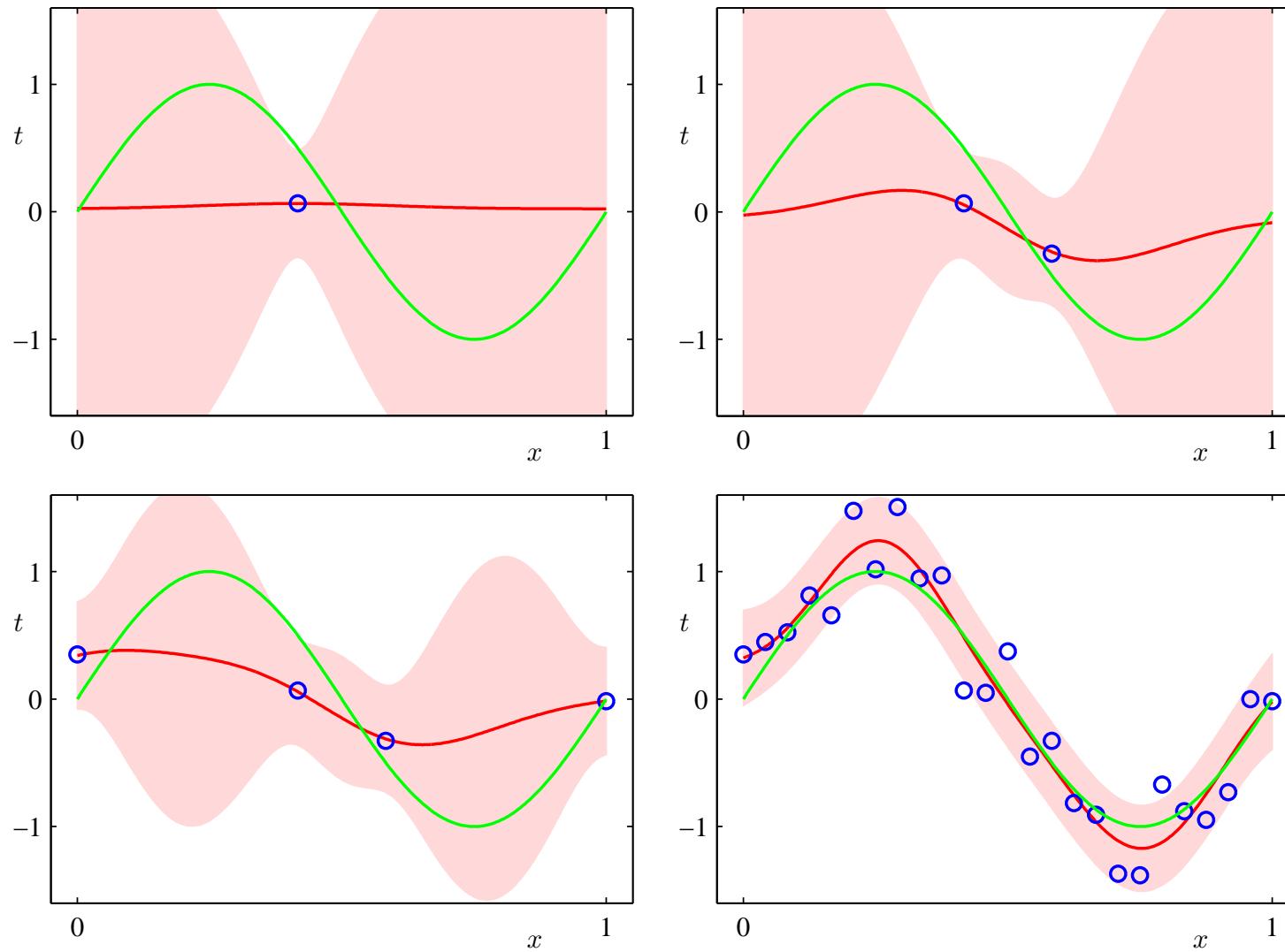
$$p(t^*|\mathbf{t}, x^*) = \mathcal{N}(t^*|\mathbf{m}_N^T \phi(\mathbf{x}), \sigma_N^2(x^*)),$$

with

$$\sigma_N^2(x^*) = \beta^{-1} + \phi(x^*)^T \mathbf{S}_N \phi(x^*).$$

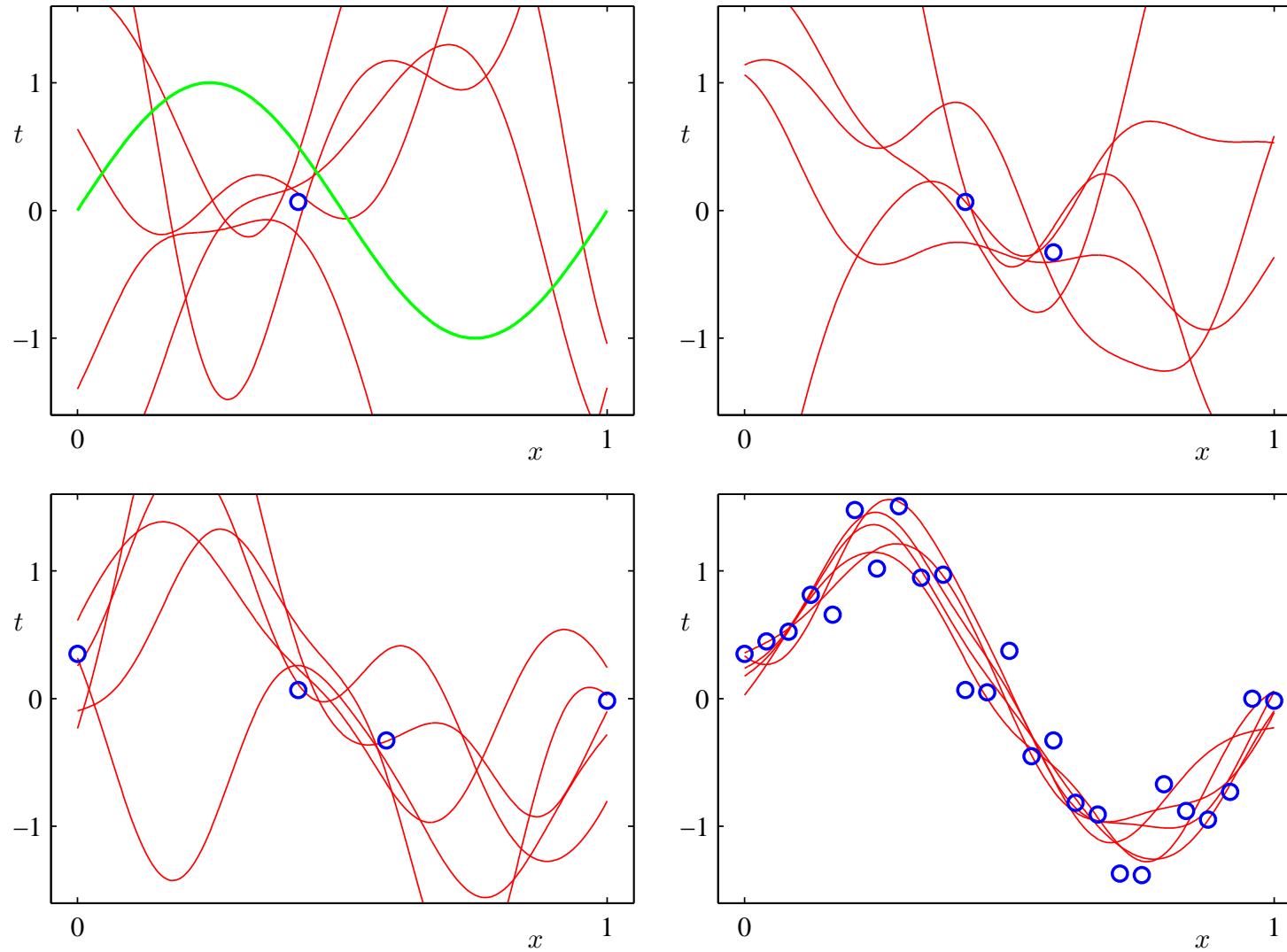
What happens to  $\sigma_N^2(x^*)$  when  $N \rightarrow \infty$ ?

# Bayesian linear regression: Example



Data points from  $t = \sin(2\pi x) + \text{noise}$ .  $y(x)$  as in Eq. 3.3-4. Data set size  $N = 1, 2, 4, 25$ .

# Bayesian linear regression: Example



Same data and model. Curves  $y(x, \mathbf{w})$  with  $\mathbf{w}$  from posterior.

## Bayesian model comparison

Maximum likelihood suffers from overfitting, which requires testing models of different complexity on separate data.

Bayesian approach allows to compare different models directly on the training data, but requires integration over model parameters.

Consider  $L$  probability models and a set of data generated from one of these models. We define a prior over models  $p(\mathcal{M}_i), i = 1, \dots, L$  to express our prior uncertainty.

Given the training data  $\mathcal{D}$ , we wish to compute the posterior probability

$$p(\mathcal{M}_i | \mathcal{D}) \propto p(\mathcal{D} | \mathcal{M}_i) p(\mathcal{M}_i)$$

$p(\mathcal{D} | \mathcal{M}_i)$  is called *model evidence*, also *marginal likelihood*, since it integrates over model parameters:

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

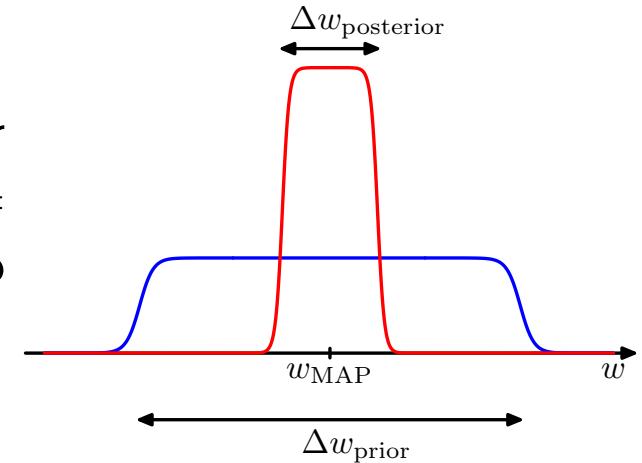
## Bayesian model comparison

Flat prior over models  $p(\mathcal{M}_i) \implies$  model evidence is proportional to marginal likelihood:

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

We use a very crude estimate of this integral to understand the mechanism of Bayesian model selection.

Consider first the one dimensional case (one parameter). Assume roughly block-shaped prior  $p(w)$  and posterior  $p(w|\mathcal{D})$  as in the figure. Then the prior is  $p(w) = 1/\Delta w_{\text{prior}}$  on an interval of width  $\Delta w_{\text{prior}}$ , and zero elsewhere. The marginal likelihood is approximately given by:



$$p(\mathcal{D}|\mathcal{M}_i) \approx p(\mathcal{D}|w_{\text{MAP}}, \mathcal{M}_i) \frac{\Delta w_{\text{posterior}}}{\Delta w_{\text{prior}}}$$

$$\ln p(\mathcal{D}|\mathcal{M}_i) \approx \ln p(\mathcal{D}|w_{\text{MAP}}, \mathcal{M}_i) + \ln \left( \frac{\Delta w_{\text{posterior}}}{\Delta w_{\text{prior}}} \right)$$



## Bayesian model comparison

The Bayesian approach gives a negative correction to the ML estimate.

With  $M$  model parameters, the same argument gives

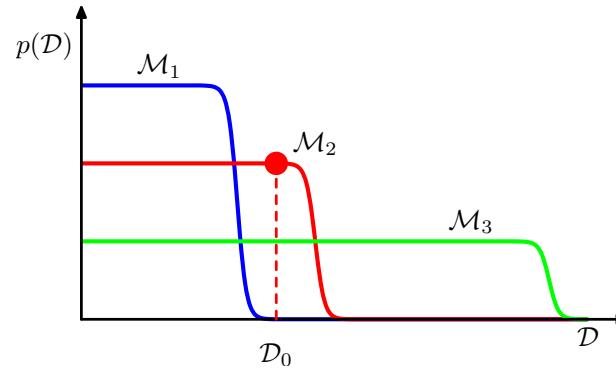
$$\ln p(\mathcal{D}|\mathcal{M}_i) \approx \ln p(\mathcal{D}|\boldsymbol{w}_{\text{MAP}}, \mathcal{M}_i) + M \ln \left( \frac{\Delta w_{\text{posterior}}}{\Delta w_{\text{prior}}} \right)$$

With increasing  $M$ , the first term increases (better fit), but the second term decreases.

## Bayesian model comparison

Consider three models of increasing complexity  $\mathcal{M}_1, \mathcal{M}_2, \mathcal{M}_3$ . Consider drawing data sets from these models: we first sample a parameter vector  $\mathbf{w}$  from the prior  $p(\mathbf{w}|\mathcal{M}_i)$  and then generate iid data points according to  $p(x|\mathbf{w}, \mathcal{M}_i)$ . The resulting distribution is  $p(\mathcal{D}|\mathcal{M}_i)$ .

A simple model has less variability in the resulting data sets than a complex model. Thus,  $p(\mathcal{D}|\mathcal{M}_1)$  is more peaked than  $p(\mathcal{D}|\mathcal{M}_3)$ . Due to normalization,  $p(\mathcal{D}|\mathcal{M}_1)$  is necessarily higher than  $p(\mathcal{D}|\mathcal{M}_3)$ .



For the data set  $\mathcal{D}_0$  the Bayesian approach will select model  $\mathcal{M}_2$  because model  $\mathcal{M}_1$  is too simple (too low likelihood) and model  $\mathcal{M}_3$  is too complex (too large penalty term). This is known as Bayesian model selection.

# Chapter 4

## Linear models for classification

- Discriminant function
- Conditional probabilities  $P(C_k|\mathbf{x})$
- Generative models  $P(\mathbf{x}|C_k)$

## Generalized linear models for classification

$$y(\mathbf{x}) = f(\mathbf{w}^T \mathbf{x} + w_0)$$

with  $f(\cdot)$  nonlinear.

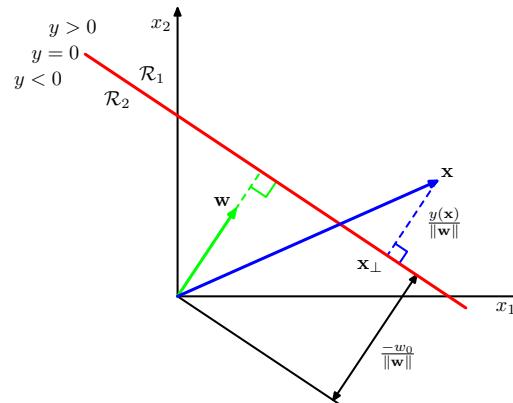
## 4.1 Discriminant functions

Two category classification problem:

$$y(\mathbf{x}) > 0 \Leftrightarrow \mathbf{x} \in C_1 \quad y(\mathbf{x}) < 0 \Leftrightarrow \mathbf{x} \in C_2$$

with

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$



If  $x_1, x_2$  on the line  $y(\mathbf{x}) = 0$  then  $\mathbf{w}^T(\mathbf{x}_1 - \mathbf{x}_2) = 0$ . Line is orthogonal to  $\mathbf{w}$ .

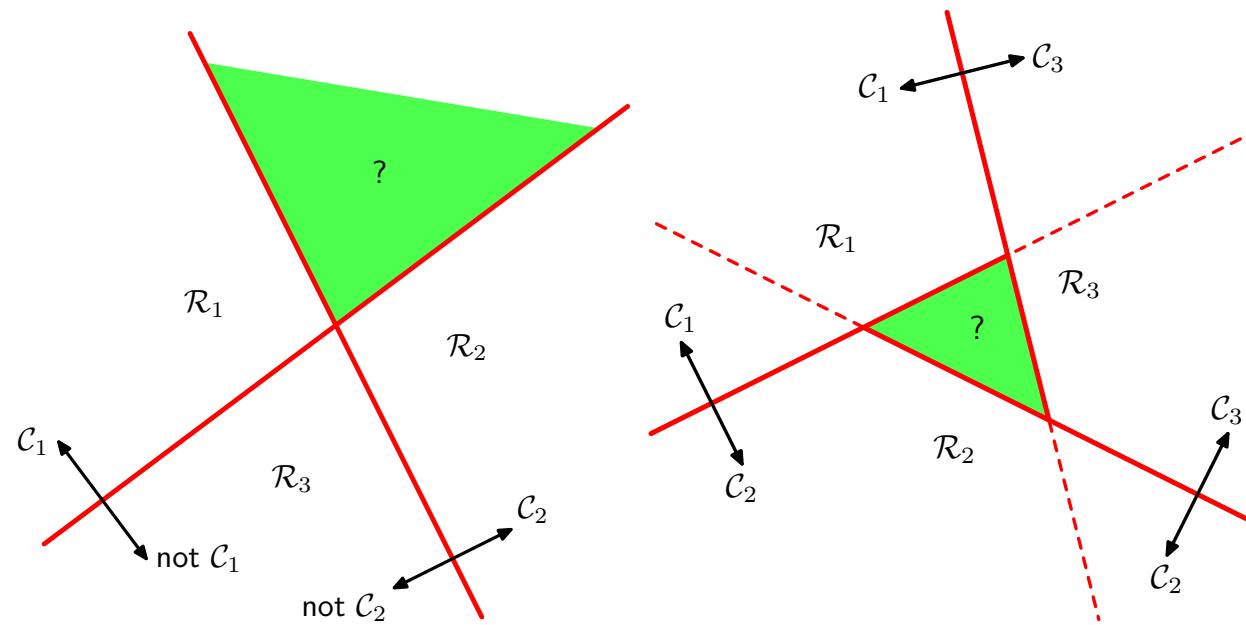
Notation:  $\tilde{\mathbf{w}} = (w_0, \mathbf{w})$  and  $\tilde{\mathbf{x}} = (1, \mathbf{x})$ , then

$$y(\mathbf{x}) = \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}$$

# Multiple classes

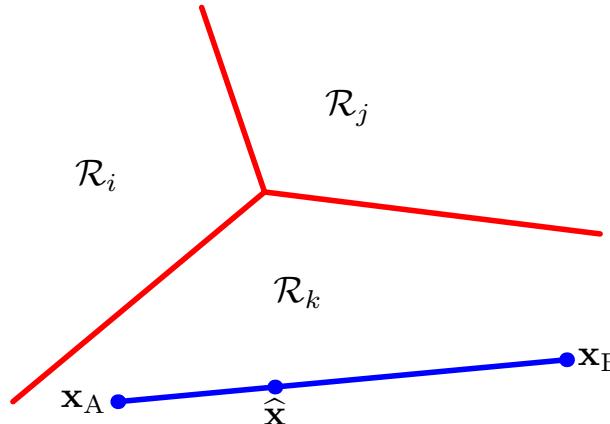
$$y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0} = \sum_{i=1}^k w_{ki} x_i + w_{k0} = \sum_{i=0}^k w_{ki} x_i$$

Naive approach is to let each boundary separate class  $k$  from the rest. This does not work.



## Multiple classes

Instead, define decision boundaries as:  $y_k(\mathbf{x}) = y_j(\mathbf{x})$ .



If  $y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{0k}$ ,  $k = 1, 2$  then decision boundary is  $(\mathbf{w}_1 - \mathbf{w}_2)^T \mathbf{x} + w_{01} - w_{02} = 0$ .

For  $K = 3$ , we get three lines:  $y_1(\mathbf{x}) = y_2(\mathbf{x})$ ,  $y_1(\mathbf{x}) = y_3(\mathbf{x})$ ,  $y_2(\mathbf{x}) = y_3(\mathbf{x})$ , that cross in a common point  $y_1(\mathbf{x}) = y_2(\mathbf{x}) = y_3(\mathbf{x})$ , and which defines the tessellation.

The class regions  $\mathcal{R}_k$  are convex. If  $\mathbf{x}, \mathbf{y} \in \mathcal{R}_k$

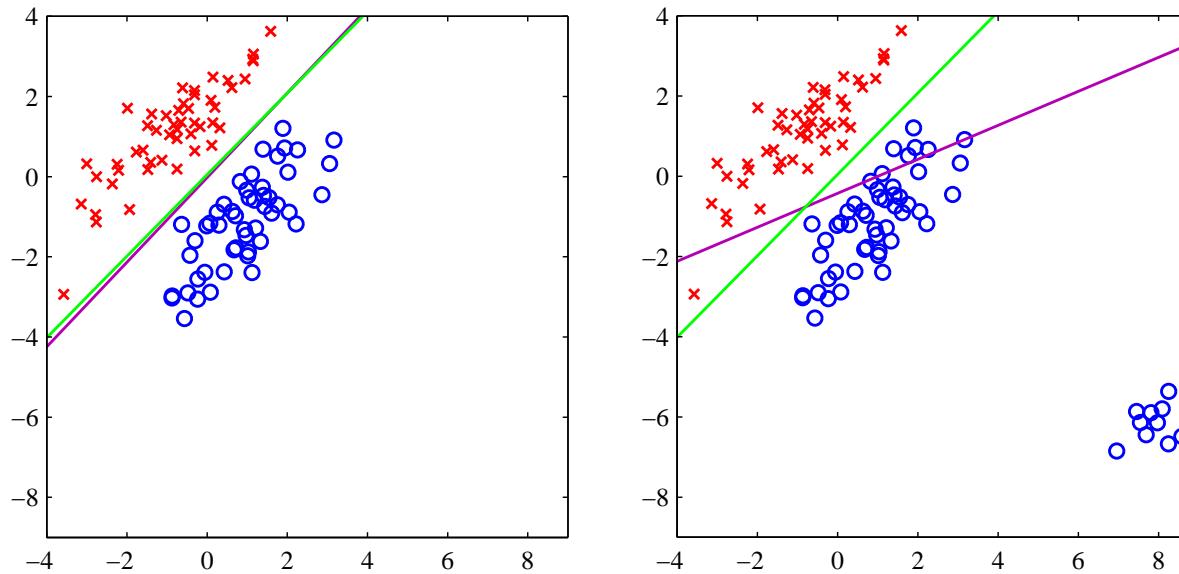
$$y_k(\mathbf{x}) > y_j(\mathbf{x}), y_k(\mathbf{y}) > y_j(\mathbf{y}) \Rightarrow y_k(\alpha \mathbf{x} + (1 - \alpha) \mathbf{y}) > y_j(\alpha \mathbf{x} + (1 - \alpha) \mathbf{y})$$

## Least-squares technique

Learning the parameters is done by minimizing a cost function, for instance the minus log likelihood or the quadratic cost:

$$\begin{aligned}
 E &= \frac{1}{2} \sum_n \sum_k \left( \sum_i \tilde{x}_{ni} \tilde{w}_{ik} - t_k^n \right)^2 \\
 \frac{\partial E}{\partial \tilde{w}_{jk}} &= \sum_n \left( \sum_i \tilde{x}_{ni} \tilde{w}_{ik} - t_k^n \right) \tilde{x}_{nj} = 0 \\
 \sum_i \left( \sum_n \tilde{x}_{nj} \tilde{x}_{ni} \right) \tilde{w}_{ik} &= \sum_n t_k^n x_{nj} \\
 \tilde{w}_k &= C^{-1} \mu_k \quad C_{ij} = \sum_n \tilde{x}_{nj} \tilde{x}_{ni} \quad \mu_{kj} = \sum_n t_k^n x_{nj}
 \end{aligned}$$

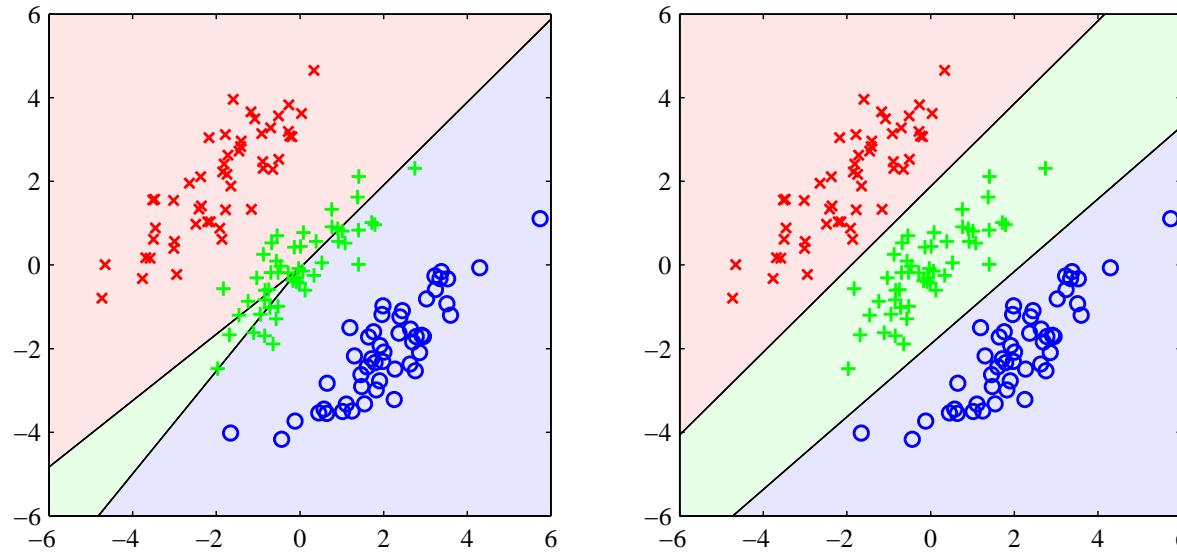
# Least-squares technique



Two class classification problem solved with least square (magenta) and logistic regression (green)

Least square solution sensitive to outliers: far away data contribute too much to the error.

# Least-squares technique



Three class classification problem solved with least square (left) and logistic regression (right)

Least square solution poor in multi-class case.

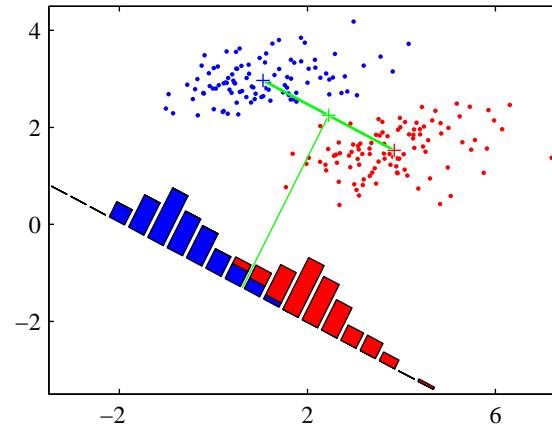
# Fisher's linear discriminant

Consider two classes. Take an  $x$  and project it down to one dimension using

$$y = \mathbf{w}^T \mathbf{x}$$

Let the two classes have two means:

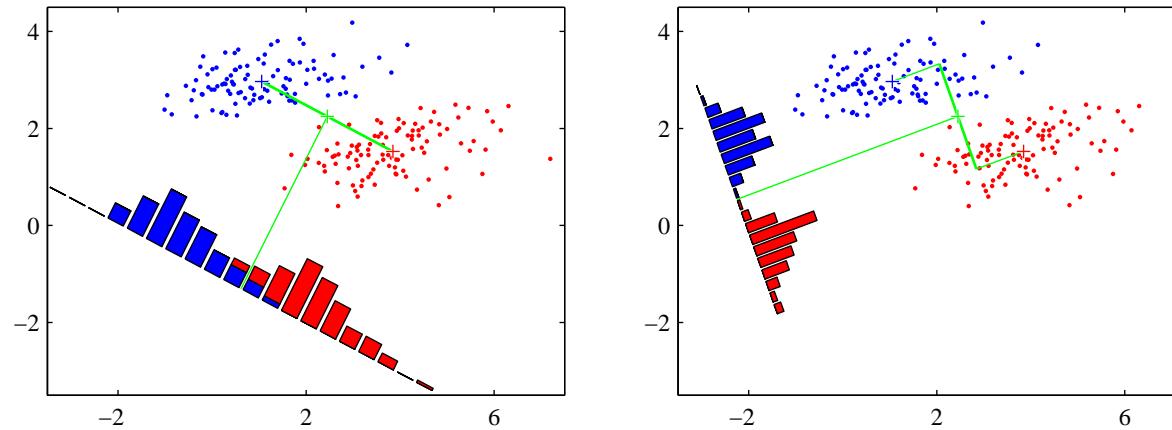
$$\mathbf{m}_1 = \frac{1}{N_1} \sum_{n \in C_1} \mathbf{x}^n \quad \mathbf{m}_2 = \frac{1}{N_2} \sum_{n \in C_2} \mathbf{x}^n$$



We can choose  $\mathbf{w}$  to maximize  $\mathbf{w}^T(\mathbf{m}_1 - \mathbf{m}_2)$ , subject to  $\sum_i w_i^2 = 1$  yields  $\mathbf{w} \propto (\mathbf{m}_1 - \mathbf{m}_2)$ . (exercise 4.4).

# Fisher's linear discriminant

Fisher discriminant analysis: Better separation when simultaneously minimize within class variance.



Within class variance after transformation:

$$s_k^2 = \sum_{n \in C_k} (y^n - m_k)^2, \quad k = 1, 2$$

with  $m_k = \mathbf{w}^T \mathbf{m}_k$ ,  $y^n = \mathbf{w}^T \mathbf{x}^n$ . Total within class variance is  $s_1^2 + s_2^2$ .

The Fisher method minimizes within class variance and maximizes mean separation by

maximizing

$$J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2} = \frac{\mathbf{w}^T S_B \mathbf{w}}{\mathbf{w}^T S_W \mathbf{w}}$$

$$S_B = (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^T$$

$$S_W = \sum_{n \in \mathcal{C}_1} (\mathbf{x}^n - \mathbf{m}_1)(\mathbf{x}^n - \mathbf{m}_1)^T + \sum_{n \in \mathcal{C}_2} (\mathbf{x}^n - \mathbf{m}_2)(\mathbf{x}^n - \mathbf{m}_2)^T$$

(Exercise 4.5)

Differentiation wrt  $\mathbf{w}$  yields (exercise)

$$(\mathbf{w}^T S_B \mathbf{w}) S_W \mathbf{w} = (\mathbf{w}^T S_W \mathbf{w}) S_B \mathbf{w}$$

Drop scalar factors, and  $S_B \mathbf{w} \propto \mathbf{m}_2 - \mathbf{m}_1$ :

$$\mathbf{w} \propto S_W^{-1}(\mathbf{m}_2 - \mathbf{m}_1)$$

## Fisher's linear discriminant

When  $S_W \propto$  the unit matrix (isotropic within-class covariance), we recover the naive solution  $\mathbf{w} \propto (\mathbf{m}_2 - \mathbf{m}_1)$ .

A classifier is built by thresholding  $y = \mathbf{w}^T \mathbf{x}$ .

The Fisher discriminant method is best viewed as a dimension reduction method, in this case from  $d$  dimensions to 1 dimension, such that optimal classification is obtained in the linear sense.

It can be shown that the least square method with two classes and target values  $t_1^n = N/N_1$  and  $t_2^n = -N/N_2$  is equivalent to the two class Fisher discriminant solution (section 4.1.5).

## Several classes

Class means  $\mathbf{m}_k = \frac{1}{N_k} \sum_{n \in C_k} \mathbf{x}^n$  and covariance  $\mathbf{S}_k = \sum_{n \in C_k} (\mathbf{x}^n - \mathbf{m}_k)(\mathbf{x}^n - \mathbf{m}_k)^T$ .  
 Relate total within class covariance  $\mathbf{S}_W = \sum_k \mathbf{S}_k$  to total covariance:

$$\begin{aligned}
 \mathbf{S}_W &= \sum_k \sum_{n \in C_k} (\mathbf{x}^n - \mathbf{m}_k)(\mathbf{x}^n - \mathbf{m}_k)^T \\
 &= \sum_k \sum_{n \in C_k} (\mathbf{x}^n - \mathbf{m} + \mathbf{m} - \mathbf{m}_k)(\mathbf{x}^n - \mathbf{m} + \mathbf{m} - \mathbf{m}_k)^T \\
 &= \sum_k \sum_{n \in C_k} (\mathbf{x}^n - \mathbf{m})(\mathbf{x}^n - \mathbf{m})^T + (\mathbf{x}^n - \mathbf{m})(\mathbf{m} - \mathbf{m}_k)^T + (\mathbf{m} - \mathbf{m}_k)(\mathbf{x}^n - \mathbf{m})^T \\
 &\quad + (\mathbf{m} - \mathbf{m}_k)(\mathbf{m} - \mathbf{m}_k)^T \\
 &= \mathbf{S}_T - \sum_k N_k (\mathbf{m}_k - \mathbf{m})(\mathbf{m}_k - \mathbf{m})^T = \mathbf{S}_T - \mathbf{S}_B
 \end{aligned}$$

with total covariance  $\mathbf{S}_T = \sum_n (\mathbf{x}^n - \mathbf{m})(\mathbf{x}^n - \mathbf{m})^T$  and  $\mathbf{S}_B = \sum_{k=1}^K N_k (\mathbf{m}_k - \mathbf{m})(\mathbf{m}_k - \mathbf{m})^T$ .

Assume input dimension  $D > K$ , the number of classes.

We introduce  $D' > 1$  features  $y_{d'} = \mathbf{w}_{d'}^T \mathbf{x}, d' = 1, \dots, D'$  or

$$\mathbf{y} = \mathbf{W}^T \mathbf{x}, \quad \mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_{D'})$$

$\mathbf{W}^T \mathbf{S}_W \mathbf{W}$  and  $\mathbf{W}^T \mathbf{S}_B \mathbf{W}$  are the  $D'$ -dimensional projection of the within class and between class covariance matrices. Maximize

$$J(\mathbf{W}) = \text{Tr} ((\mathbf{W}^T \mathbf{S}_W \mathbf{W})^{-1} (\mathbf{W}^T \mathbf{S}_B \mathbf{W}))$$

The solution  $\mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_{D'})$  consists of the largest  $D'$  eigenvectors of the matrix  $\mathbf{S}_W^{-1} \mathbf{S}_B$ .

Note,  $\mathbf{S}_B$  is sum of  $K$  rank one matrices and is of rank  $K - 1$ .  $\mathbf{S}_W^{-1} \mathbf{S}_B$  has  $(K - 1)$  non-zero eigenvalues.

Thus,  $D' \leq K - 1$ , because otherwise  $\mathbf{W}$  contains any of the zero eigenvectors.

# The Perceptron

Relevant in history of pattern recognition and neural networks.

- Perceptron learning rule + convergence, Rosenblatt (1962)
- Perceptron critique (Minsky and Papert, 1969) → "Dark ages of neural networks"
- Revival in the 80's: Backpropagation

# The Perceptron

$$y(x) = \text{sign}(\mathbf{w}^T \phi(x))$$

where

$$\text{sign}(a) = \begin{cases} +1, & a \geq 0 \\ -1, & a < 0. \end{cases}$$

and  $\phi(x)$  is a feature vector (e.g. hard wired neural network).

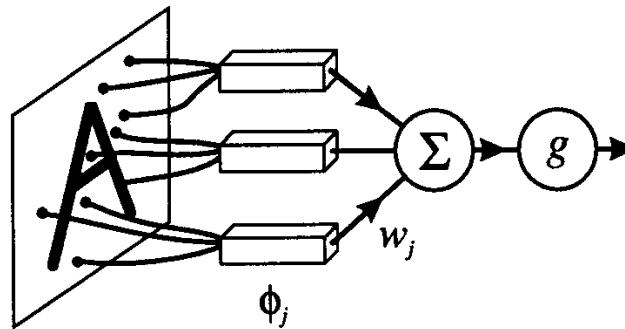


Figure 3.10. The perceptron network used a fixed set of processing elements, denoted  $\phi_j$ , followed by a layer of adaptive weights  $w_j$  and a threshold activation function  $g(\cdot)$ . The processing elements  $\phi_j$  typically also had threshold activation functions, and took inputs from a randomly chosen subset of the pixels of the input image.

# The Perceptron algorithm

Perceptron approach: seek  $\mathbf{w}$  such that  $\mathbf{w}^T \phi(\mathbf{x}_n) > 0$  for  $\mathbf{x}_n \in C_1$  and  $\mathbf{w}^T \phi(\mathbf{x}_n) < 0$  for  $\mathbf{x}_n \in C_2$ . Target coding:  $t = +1, -1$  for  $x \in \{C_1, C_2\}$ . Then we want all patterns:

$$t_n \mathbf{w}^T \phi(\mathbf{x}_n) > 0$$

Perceptron criterion: maximize

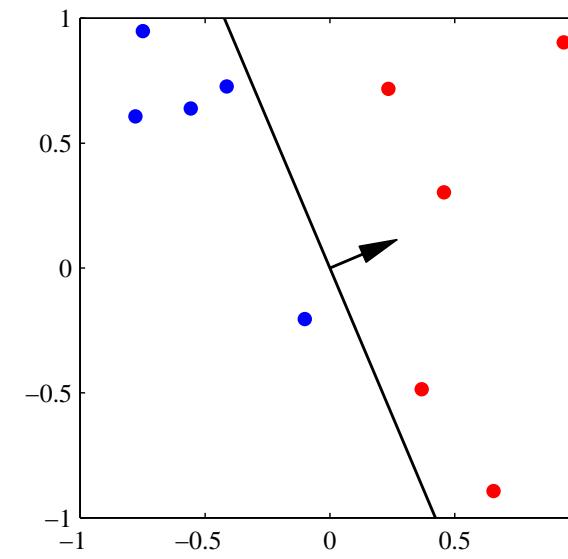
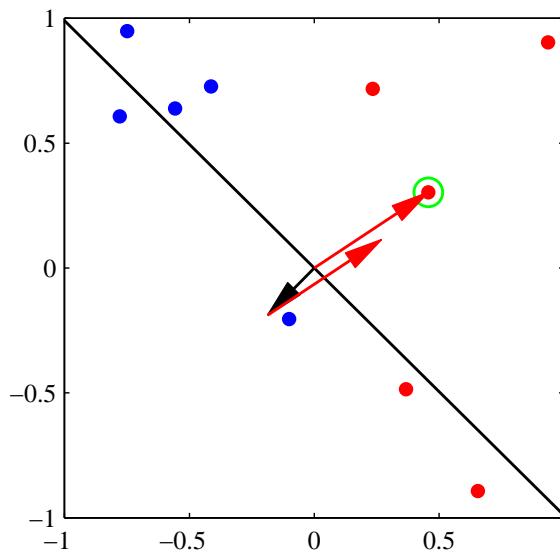
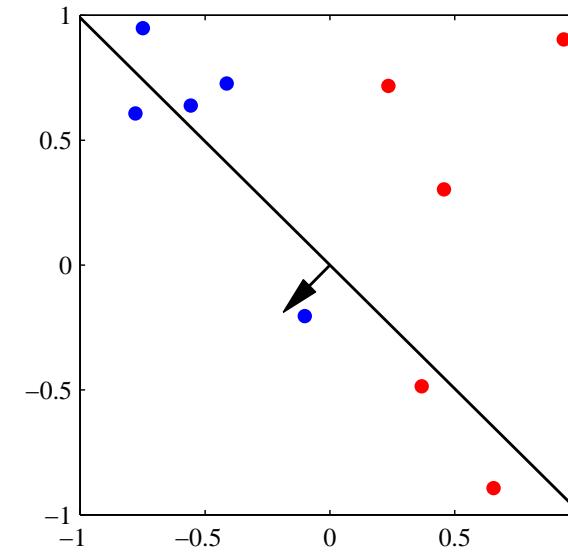
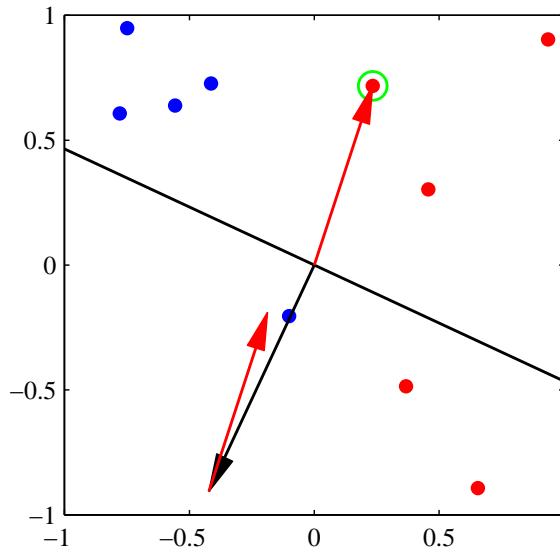
$$C(\mathbf{w}) = \frac{1}{\|\mathbf{w}\|} \min_{n \in M} t_n \mathbf{w}^T \phi(\mathbf{x}_n)$$

Learning rule: choose pattern  $n$ . If missclassified, update according to

$$\mathbf{w}^{\tau+1} = \mathbf{w}^\tau + \eta \phi_n t_n$$

It can be shown that this algorithm converges in a *finite number of steps, if the data is linearly separable.*

# Convergence of perceptron learning



# Perceptron critique

Minsky and Papert: perceptron can only learn linearly separable problems:

Most functions are not linearly separable: e.g. the problem of determining whether objects are singly connected, using local receptive fields

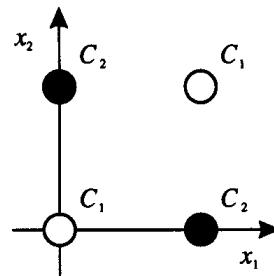


Figure 3.6. The exclusive-OR problem consists of four patterns in a two-dimensional space as shown. It provides a simple example of a problem which is not linearly separable.

# Probabilistic generative models

Probabilistic models for classification

$$\begin{aligned} p(C_1|\mathbf{x}) &= \frac{p(\mathbf{x}|C_1)p(C_1)}{p(\mathbf{x}|C_1)p(C_1) + p(\mathbf{x}|C_2)p(C_2)} \\ &= \frac{1}{1 + \exp(-a)} = \sigma(a) \end{aligned}$$

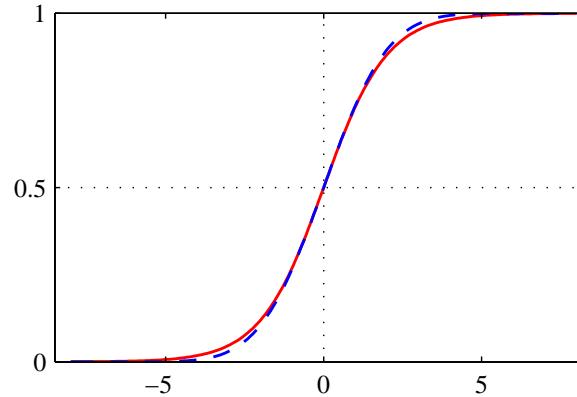
where  $a$  is the "log odds"

$$a = \ln \frac{p(\mathbf{x}|C_1)p(C_1)}{p(\mathbf{x}|C_2)p(C_2)}$$

and  $\sigma$  the logistic sigmoid function (i.e. S-shaped function)

$$\sigma(a) = \frac{1}{1 + \exp(-a)}$$

# Sigmoid function



Properties:

$$\sigma(-a) = 1 - \sigma(a)$$

inverse:

$$a = \ln \frac{\sigma}{1 - \sigma}$$

## Sigmoid/softmax function

Softmax: generalization to  $K$  classes

$$p(C_k | \mathbf{x}) = \frac{p(\mathbf{x} | C_k) p(C_k)}{\sum_j p(\mathbf{x} | C_j) p(C_j)} = \frac{\exp(a_k)}{\sum_j \exp(a_j)}$$

where

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

Note that softmax is invariant under  $a_k \rightarrow a_k + \text{const}$

## Continuous inputs

The discriminant function for a Gaussian distribution is

$$\begin{aligned} a_k(\mathbf{x}) &= \log p(\mathbf{x}|C_k) + \log p(C_k) \\ &= -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) - \frac{1}{2} \log |\Sigma_k| + \log p(C_k) \end{aligned}$$

The decision boundaries  $a_k(\mathbf{x}) = a_l(\mathbf{x})$  are quadratic functions in  $d$  dimensions.

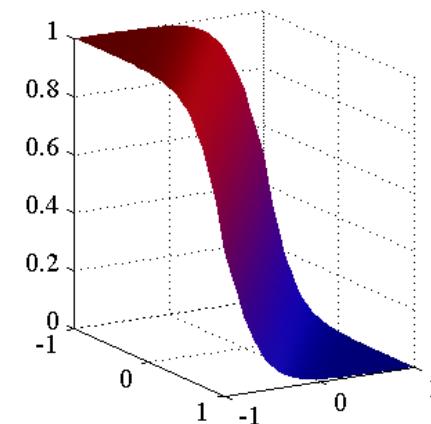
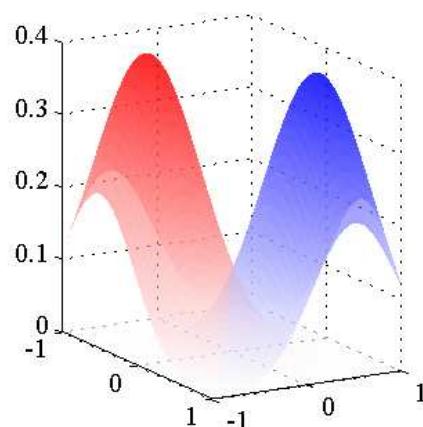
If  $\Sigma_k^{-1} = \Sigma^{-1}$ , then

$$\begin{aligned} a_k(\mathbf{x}) &= -\frac{1}{2}\mathbf{x}\Sigma^{-1}\mathbf{x} + \frac{1}{2}\boldsymbol{\mu}_k^T\Sigma^{-1}\mathbf{x} + \frac{1}{2}\mathbf{x}\Sigma^{-1}\boldsymbol{\mu}_k - \frac{1}{2}\boldsymbol{\mu}_k^T\Sigma^{-1}\boldsymbol{\mu}_k - \frac{1}{2}\log|\Sigma| + \log p(C_k) \\ &= \mathbf{w}_k^T \mathbf{x} + w_{k0} + \text{const} \end{aligned}$$

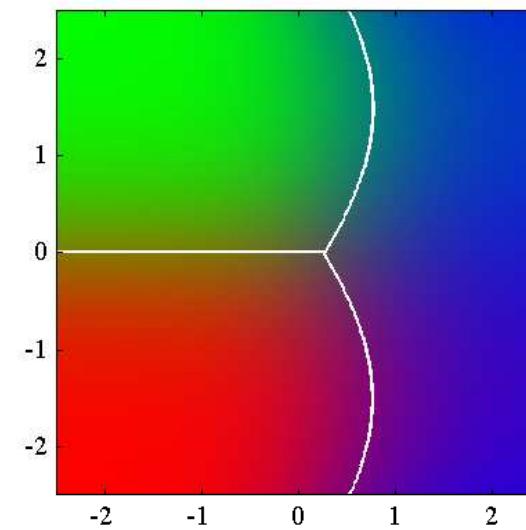
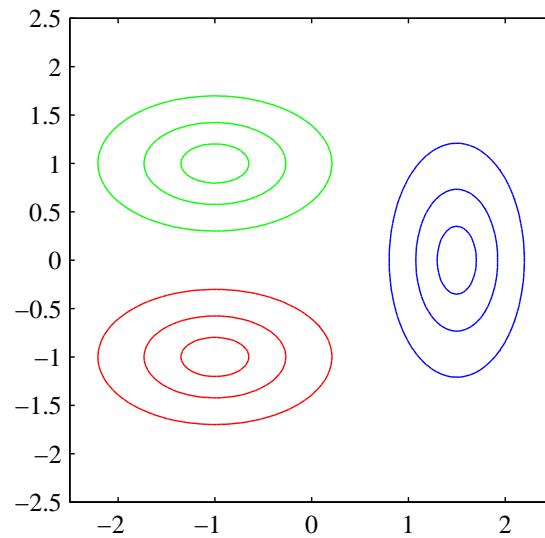
with

$$\mathbf{w}_k^T = \boldsymbol{\mu}_k^T \Sigma^{-1}, \quad w_{k0} = -\frac{1}{2}\boldsymbol{\mu}_k^T \Sigma^{-1} \boldsymbol{\mu}_k + \log p(C_k)$$

The discriminant function is linear, and the decision boundary is a straight line (or hyper-plane in  $d$  dimensions).



When class conditional covariances are equal, decision boundary is straight line.



When class conditional covariances are unequal decision boundary is quadratic.

## Maximum likelihood solution

Once we have specified a parametric functional form for  $p(\mathbf{x}|C_k)$  and  $p(C_k)$  we can determine parameters using maximum (joint!) likelihood (as usual!).

Actually, for classification, we should aim for the *conditional* likelihood, but the joint likelihood is in general easier. In a perfect model, it makes no difference, but with imperfect models, it certainly does!

Example of joint ML: a binary problem with Gaussian class-conditionals, with a shared covariance matrix. Data is  $\{\mathbf{x}_n, t_n\}_n$  with labels coded as  $t = \{1, 0\}$  for classes  $C_1, C_2$ .

Class probabilities are parametrized as

$$p(C_1) = \pi, \quad p(C_2) = 1 - \pi$$

and class-conditional densities as

$$p(\mathbf{x}_n|C_1) = \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_1, \boldsymbol{\Sigma}), \quad p(\mathbf{x}_n|C_2) = \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_2, \boldsymbol{\Sigma})$$

Then the likelihood is given by

$$\begin{aligned}
 L &= p(t_1, \mathbf{x}_1, \dots, t_N, \mathbf{x}_n | \pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}) = \prod_n p(t_n, \mathbf{x}_n | \pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}) \\
 &= \prod_n [\pi \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma})]^{t_n} [(1 - \pi) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_2, \boldsymbol{\Sigma})]^{1-t_n} \\
 \log L &= \sum_n t_n \log[\pi \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma})] + (1 - t_n) \log[(1 - \pi) \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_2, \boldsymbol{\Sigma})]
 \end{aligned}$$

The maximum likelihood solution for  $\pi$  and  $\boldsymbol{\mu}_{1,2}$ :

$$\begin{aligned}
 \frac{\partial \log L}{\partial \pi} &= \frac{1}{\pi} \sum_n t_n - \frac{1}{1 - \pi} \sum_n (1 - t_n) \quad \pi = \frac{N_1}{N} \quad N_1 = \sum_n t_n \\
 \frac{\partial \log L}{\partial \boldsymbol{\mu}_1} &= \sum_n t_n \frac{\partial \log \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_1, \boldsymbol{\Sigma})}{\partial \boldsymbol{\mu}_1} = \boldsymbol{\Sigma}^{-1} \sum_n t_n (\mathbf{x}_n - \boldsymbol{\mu}_1) \quad \boldsymbol{\mu}_1 = \frac{1}{N} \sum_n t_n \mathbf{x}_n \\
 \frac{\partial \log L}{\partial \boldsymbol{\mu}_2} &= \dots \quad \boldsymbol{\mu}_2 = \frac{1}{N} \sum_n (1 - t_n) \mathbf{x}_n
 \end{aligned}$$

The maximum likelihood solution for  $\boldsymbol{\Sigma}$  given by Eqs. 4.78-80 (exercise).

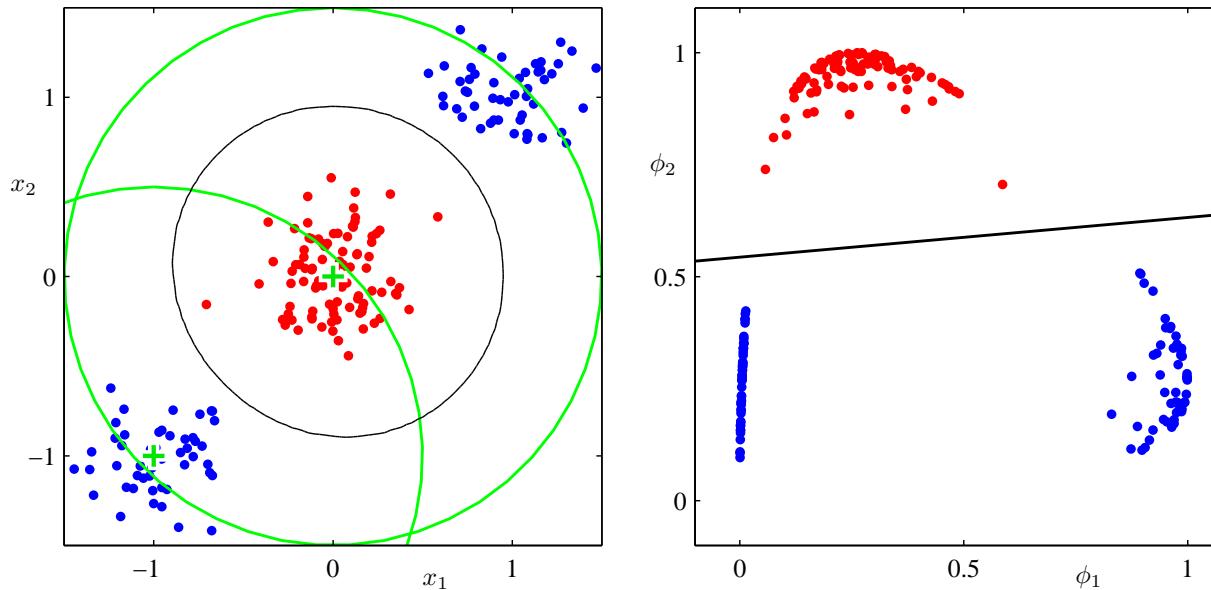
# Probabilistic discriminative models

Discriminative versus generative models

Discriminative models:

- no spoiled effort in modeling joint probabilities. Aims directly at the conditional probabilities of interest
- usually fewer parameters
- improved performance when class-conditional  $p(\mathbf{x}|C)$  assumptions are poor (with joint ML).
- basis functions can be employed

# Fixed basis functions



- Basis functions  $\phi(\mathbf{x})$ . Denote  $\phi_n = \phi(x_n)$
- Problems that are not linearly separable in  $\mathbf{x}$  might be linearly separable in  $\phi(\mathbf{x})$ .
- Note: use of basis functions  $\phi$  in place of variables  $\mathbf{x}$  is not obvious in generative models:  $\mathcal{N}(x, x^2, x^3 | \mu, \Sigma)??$ .

# Logistic regression

Two class classification

$$p(C_1|\boldsymbol{\phi}) = \sigma(\boldsymbol{w}^T \boldsymbol{\phi})$$

$M$  dimensional feature space:  $M$  parameters, while Gaussian class conditional densities would require  $2M$  (two means) and  $M(M+1)/2$  (common covariance matrix) parameters.

Maximum likelihood to determine parameters  $\boldsymbol{w}$ , (nb.  $t_n \in \{0, 1\}$ )

$$p(t_1, \dots, t_N | \boldsymbol{w}, x_1, \dots, x_N) = \prod_n \sigma(\boldsymbol{w}^T \boldsymbol{\phi}_n)^{t_n} [1 - \sigma(\boldsymbol{w}^T \boldsymbol{\phi}_n)]^{1-t_n} = \prod_n y_n^{t_n} (1 - y_n)^{1-t_n}$$

with  $\boldsymbol{\phi}_n = \boldsymbol{\phi}(x_n)$  and  $y_n = \sigma(\boldsymbol{w}^T \boldsymbol{\phi}_n)$  i.e.,

$$E(\boldsymbol{w}) = -\ln p = -\sum_n t_n \ln y_n + (1 - t_n) \ln(1 - y_n)$$

NB: entropic error function for classification, rather than squared error.

## Logistic regression

$$\nabla E(\mathbf{w}) = \sum_n (\sigma(\mathbf{w}^T \phi_n) - t_n) \phi_n$$

No closed form solution. Optimization by gradient descent, (or e.g., Newton-Raphson).

Overfitting risk when data is linearly separable:  $\mathbf{w} \rightarrow \infty$  (i.e.  $\sigma \rightarrow$  step function).

## Iterative least squares

Minimize learning error by Newton-Raphson method

$$\mathbf{w}^{(new)} = \mathbf{w}^{(old)} - \mathbf{H}^{-1} \nabla E(\mathbf{w})$$

with

$$\begin{aligned}\mathbf{H}_{ij} &= \frac{\partial^2 E}{\partial w_i \partial w_j} = \frac{\partial}{\partial w_j} \frac{\partial E}{\partial w_i} = \frac{\partial}{\partial w_j} \sum_n (y_n - t_n) \phi_i(\mathbf{x}_n) \\ &= \sum_n \phi_j(\mathbf{x}_n) y_n (1 - y_n) \phi_i(\mathbf{x}_n) = \Phi^T R \Phi \\ \nabla_i E(\mathbf{w}) &= \sum_n (y_n - t_n) \phi_i(\mathbf{x}_n) = \Phi^T (\mathbf{y} - \mathbf{t})\end{aligned}$$

with  $\Phi_{nj} = \phi_j(\mathbf{x}_n)$  and  $R_{n,n'} = y_n (1 - y_n) \delta_{n,n'}$ .

$H(\mathbf{w})$  is positive definite for all  $\mathbf{w}$  thus  $E(\mathbf{w})$  is convex, thus unique optimum (Ex. 4.15).

$$\mathbf{w}^{(new)} = \mathbf{w}^{(old)} - (\Phi^T R \Phi)^{-1} \Phi^T (\mathbf{y} - \mathbf{t})$$

## Laplace approximation

Assume distribution  $p(z)$  is given up to normalization, i.e. in the form

$$p(z) = \frac{1}{Z} f(z) \quad Z = \int f(z) dz$$

where  $f(z)$  is given, but  $Z$  is unknown (and the integral is infeasible).

Goal: approximate by a Gaussian  $q(z)$ , centered around the mode of  $p(z)$ .

## Laplace approximation

Mode  $z_0$  is maximum of  $p(z)$ , i.e.  $dp(z)/dz|_{z_0} = 0$ , or

$$\frac{df(z)}{dz}\Big|_{z_0} = 0$$

The logarithm of a Gaussian is a quadratic function of the variables, so it makes sense to make a second order Taylor expansion of  $\ln f$  around the mode (this would be exact if  $p$  was Gaussian).

$$\ln f(z) \approx \ln f(z_0) - \frac{1}{2}A(z - z_0)^2$$

where

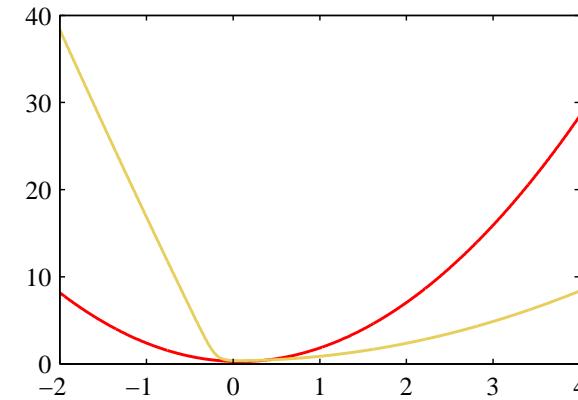
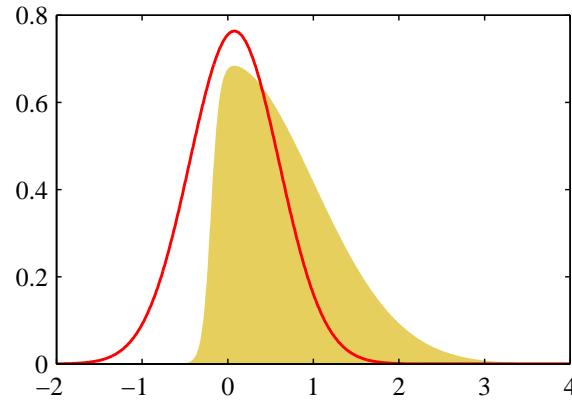
$$A = -\frac{d^2 \ln f(z)}{dz^2}\Big|_{z_0}$$

Note that the first order term is absent since we expand around the maximum. Taking the exponent we obtain,

$$f(z) \approx f(z_0) \exp\left(-\frac{1}{2}A(z - z_0)^2\right)$$

and the Gaussian approximation is obtained by normalization

$$q(z) = \left(\frac{A}{2\pi}\right)^{1/2} \exp\left(-\frac{1}{2}A(z - z_0)^2\right)$$



## Laplace approximation in $M$ dimensions

In  $M$  dimensions results are similar,

$$q(\mathbf{z}) = \left( \frac{\det(\mathbf{A})}{(2\pi)^{M/2}} \right) \exp\left(-\frac{1}{2}(\mathbf{z} - \mathbf{z}_0)^T \mathbf{A}(\mathbf{z} - \mathbf{z}_0)\right)$$

where

$$A_{ij} = -\frac{\partial^2}{\partial z_i \partial z_j} \ln f(\mathbf{z}) \Big|_{\mathbf{z}=\mathbf{z}_0}$$

Usually  $\mathbf{z}_0$  is found by numerical optimization.

A weakness of Laplace approximation is that it relies only on the local properties of the mode. Other methods (e.g. sampling or variational methods) are based on more global properties of  $p$ .

## Model comparison and BIC

Approximation of  $Z$ :

$$Z = \int f(\mathbf{z}) d\mathbf{z} \quad (4)$$

$$\approx f(\mathbf{z}_0) \int \exp\left(-\frac{1}{2}(\mathbf{z} - \mathbf{z}_0)^T \mathbf{A}(\mathbf{z} - \mathbf{z}_0)\right) d\mathbf{z} \quad (5)$$

$$= f(\mathbf{z}_0) \frac{(2\pi)^{M/2}}{\det(\mathbf{A})^{1/2}} \quad (6)$$

## Model comparison and BIC

Application: approximation of model evidence.

$$f(\boldsymbol{\theta}) = p(D|\boldsymbol{\theta})p(\boldsymbol{\theta}) \quad (7)$$

$$Z = P(D) = \int p(D|\boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta} \quad (8)$$

Then applying Laplace approximation, we obtain (Ex. 4. 22)

$$\ln P(D) \approx \ln p(D|\boldsymbol{\theta}_{\text{MAP}}) + \underbrace{\frac{M}{2} \ln 2\pi - \frac{1}{2} \ln \det(\mathbf{A})}_{\text{Occam factor}}$$

where

$$\mathbf{A} = -\nabla\nabla \ln p(D|\boldsymbol{\theta})p(\boldsymbol{\theta}) \Big|_{\boldsymbol{\theta}_{\text{MAP}}} = -\nabla\nabla p(\boldsymbol{\theta}|D) \Big|_{\boldsymbol{\theta}_{\text{MAP}}}$$

which can be interpreted as the inverse width of the posterior.

## Model comparison and BIC

Under some assumptions,  $\mathbf{A} = \sum_n \mathbf{A}_n \approx N \hat{\mathbf{A}}$ , and full rank, then

$$\ln \det(\mathbf{A}) \approx \ln \det(N \hat{\mathbf{A}}) = \ln(N^M \det(\hat{\mathbf{A}})) \approx M \ln N$$

leads to the Bayesian Information Criterion (BIC)

$$\ln P(D) \approx \ln p(D|\theta_{\text{MAP}}) - \frac{1}{2}M \ln N \quad (+\text{const.})$$

## Bayesian logistic regression

Prior  $p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0)$

Log posterior = log prior + log likelihood + const

$$\begin{aligned}\ln p(\mathbf{w} | \mathbf{t}) &= -\frac{1}{2}(\mathbf{w} - \mathbf{m}_0)^T \mathbf{S}_0^{-1}(\mathbf{w} - \mathbf{m}_0) \\ &\quad + \sum_n \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\} + const\end{aligned}$$

where  $y_n = \sigma(\mathbf{w}^T \phi_n)$ .

Posterior distribution in  $p(\mathbf{w} | \mathbf{t})$ ?

Laplace approximation: find  $\mathbf{w}_{\text{MAP}}$  and compute second derivatives:

$$\mathbf{S}_N^{-1} = -\nabla \nabla \ln p(\mathbf{w} | \mathbf{t}) = \mathbf{S}_0^{-1} + \sum_n y_n(1 - y_n) \phi_n \phi_n^T$$

and

$$q(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{w}_{\text{MAP}}, \mathbf{S}_N)$$

## Predictive distribution

$$p(C_1|\boldsymbol{\phi}, \mathbf{t}) = \int \sigma(\mathbf{w}^T \boldsymbol{\phi}) p(\mathbf{w}|\mathbf{t}) d\mathbf{w} \approx \int \sigma(\mathbf{w}^T \boldsymbol{\phi}) q(\mathbf{w}) d\mathbf{w}$$

The function  $\sigma(\mathbf{w}^T \boldsymbol{\phi})$  depends only on  $\mathbf{w}$  via its projection on  $\boldsymbol{\phi}$ . We can marginalize out the other variables, since  $q$  is a Gaussian. The marginal is then again a Gaussian for which we can compute its parameters

$$\int \sigma(\mathbf{w}^T \boldsymbol{\phi}) q(\mathbf{w}) d\mathbf{w} = \int \sigma(a) \mathcal{N}(a|\mu_a, \sigma_a^2) da$$

where the parameters turn out to be  $\mu_a = \mathbf{w}_{\text{MAP}}^T \boldsymbol{\phi}$  and  $\sigma_a^2 = \boldsymbol{\phi}^T \mathbf{S}_N \boldsymbol{\phi}$ .

Unfortunately, this integral cannot be expressed analytically. However  $\sigma(x)$  is well approximated by the probit function, i.e., the cumulative Gaussian

$$\Phi(x) = \int_{-\infty}^x \mathcal{N}(u|0, 1) du$$

In particular  $\sigma(x) \approx \Phi(\sqrt{\frac{\pi}{8}}x)$ . With additional manipulations the predictive distributions

can be shown to be approximated by

$$p(C_1|\phi, \mathbf{t}) \approx \sigma((1 + \pi\sigma_a^2/8)^{-1/2}\mu_a) = \sigma(\kappa(\sigma_a^2)\mu_a)$$

Note that the MAP predictive distribution is

$$p(C_1|\phi, \mathbf{w}_{MAP}) = \sigma(\mu_a).$$

The decision boundary  $p(C_1|\phi, \dots) = 0.5$  is the same in both approximations, namely at  $\mu_a = 0$ . Since  $\kappa < 1$ , the Laplace approximation is less certain about the classifications.

# Chapter 5 Neural Networks

## Feed-forward neural networks

Non-linear methods using a fixed set of basis functions (polynomials) suffer from curse of dimensionality.

A successful alternative is to adapt the basis functions to the problem.

- SVMs: convex optimisation, number of SVs increases with data
- MLPs: aka feed-forward neural networks, non-convex optimisation

## Feed-forward Network functions

We extend the previous regression model with fixed basis functions

$$y(\mathbf{x}, \mathbf{w}) = f \left( \sum_{j=1}^M w_j \phi_j(\mathbf{x}) \right)$$

to a model where  $\phi_j$  is adaptive:

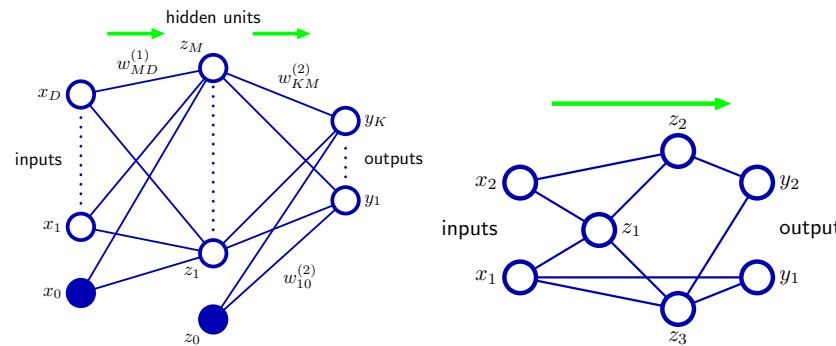
$$\phi_j(\mathbf{x}) = h \left( \sum_{i=0}^D w_{ji}^{(1)} x_i \right)$$

# Feed-forward Network functions

In the case of  $K$  outputs

$$y_k(\mathbf{x}, \mathbf{w}) = h_2 \left( \sum_{j=1}^M w_{kj}^{(2)} h_1 \left( \sum_{i=0}^D w_{ji}^{(1)} x_i \right) \right)$$

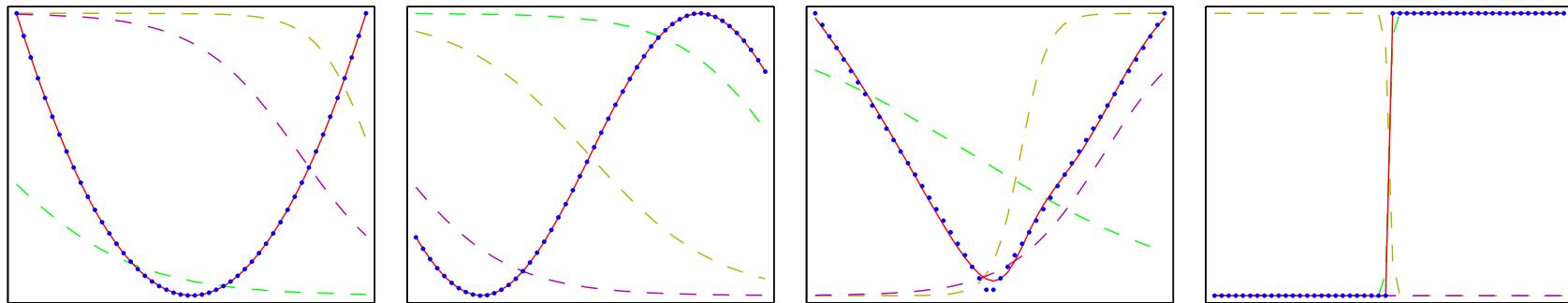
$h_2(x)$  is  $\sigma(x)$  or  $x$  depending on the problem.  $h_1(x)$  is  $\sigma(x)$  or  $\tanh(x)$ .



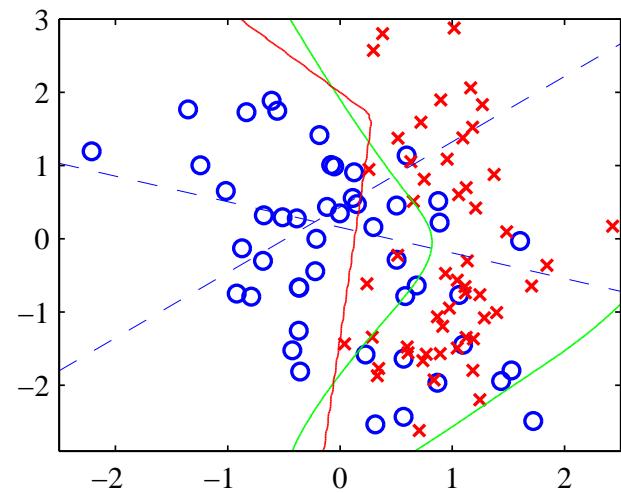
Left) Two layer architecture. Right) general feed-forward network with skip-layer connections.

If  $h_1, h_2$  linear, the model is linear. If  $M < D, K$  it computes principle components (Bishop section 12.4.2).

# Feed-forward Network functions



Two layer NN with 3 'tanh' hidden units and linear output can approximate many functions.  $x \in [-1, 1]$ , 50 equally spaced points. From left to right:  $f(x) = x^2$ ,  $\sin(x)$ ,  $|x|$ ,  $\Theta(x)$ . Dashed lines are outputs of the 3 hidden units.



Two layer NN with two inputs and 2 'tanh' hidden units and sigmoid output for classification. Dashed lines are hidden unit activities.

Feed-forward neural networks have good approximation properties.

## Weight space symmetries

For any solutions of the weights, there are many equivalent solutions due to symmetry:

- for any hidden unit  $j$  with tanh activation function, change  $w_{ji} \rightarrow -w_{ji}$  and  $w_{kj} \rightarrow -w_{kj}$ :  $2^M$  solutions
- rename the hidden unit labels:  $M!$  solutions

Thus a total of  $M!2^M$  equivalent solutions, not only for tanh activation functions.

## Network training

Regression:  $t_n$  continue valued,  $h_2(x) = x$  and one usually minimizes the squared error (one output)

$$\begin{aligned} E(\mathbf{w}) &= \frac{1}{2} \sum_{n=1}^N (y(\mathbf{x}_n, \mathbf{w}) - t_n)^2 \\ &= -\log \prod_{n=1}^N \mathcal{N}(t_n | y(\mathbf{x}_n, \mathbf{w}), \beta^{-1}) + \dots \end{aligned}$$

Classification:  $t_n = 0, 1$ ,  $h_2(x) = \sigma(x)$ ,  $y(\mathbf{x}_n, \mathbf{w})$  is probability to belong to class 1.

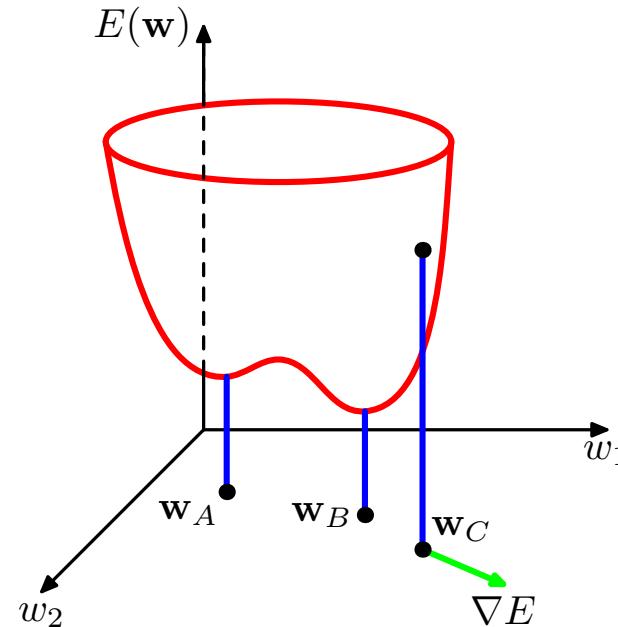
$$\begin{aligned} E(\mathbf{w}) &= -\sum_{n=1}^N \{t_n \log y(\mathbf{x}_n, \mathbf{w}) + (1 - t_n) \log(1 - y(\mathbf{x}_n, \mathbf{w}))\} \\ &= -\log \prod_{n=1}^N y(\mathbf{x}_n, \mathbf{w})^{t_n} (1 - y(\mathbf{x}_n, \mathbf{w}))^{1-t_n} \end{aligned}$$

## Network training

More than two classes: consider network with  $K$  outputs.  $t_{nk} = 1$  if  $x_n$  belongs to class  $k$  and zero otherwise.  $y_k(x_n, \mathbf{w})$  is the network output

$$\begin{aligned} E(\mathbf{w}) &= - \sum_{n=1}^N \sum_{k=1}^K t_{nk} \log p_k(x_n, \mathbf{w}) \\ p_k(x, \mathbf{w}) &= \frac{\exp(y_k(x, \mathbf{w}))}{\sum_{k'=1}^K \exp(y_{k'}(x, \mathbf{w}))} \end{aligned}$$

# Parameter optimization



$E$  is minimal when  $\nabla E(\mathbf{w}) = 0$ , but not vice versa!

As a consequence, gradient based methods find a local minimum, not necessarily the global minimum.

## Gradient descent optimization

The simplest procedure to optimize  $E$  is to start with a random  $\mathbf{w}$  and iterate

$$\mathbf{w}^{\tau+1} = \mathbf{w}^\tau - \eta \nabla E(\mathbf{w}^\tau)$$

This is called batch learning, where all training data are included in the computation of  $\nabla E$ .

Does this algorithm converge? Yes, if  $\epsilon$  is "sufficiently small" and  $E$  bounded from below.

Proof: Denote  $\Delta\mathbf{w} = -\eta \nabla E$ .

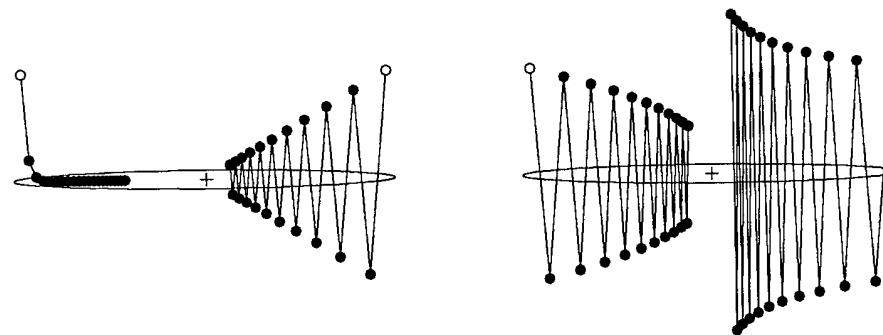
$$E(\mathbf{w} + \Delta\mathbf{w}) \approx E(\mathbf{w}) + (\Delta\mathbf{w})^T \nabla E = E(\mathbf{w}) - \eta \sum_i \left( \frac{\partial E}{\partial w_i} \right)^2 \leq E(\mathbf{w})$$

In each gradient descent step the value of  $E$  is lowered. Since  $E$  bounded from below, the procedure must converge asymptotically.

# Convergence of gradient descent in a quadratic well

$$\begin{aligned}
 E(w) &= \frac{1}{2} \sum_i \lambda_i w_i^2 \\
 \Delta w_i &= -\eta \frac{\partial E}{\partial w_i} = -\eta \lambda_i w_i \\
 w_i^{\text{new}} &= w_i^{\text{old}} + \Delta w_i = (1 - \eta \lambda_i) w_i
 \end{aligned}$$

Convergence when  $|1 - \eta \lambda_i| < 1$ . Oscillations when  $1 - \eta \lambda_i < 0$ .



Optimal learning parameter depends on curvature of each dimension.

# Learning with momentum

One solution is adding momentum term:

$$\begin{aligned}
 \Delta w_{t+1} &= -\eta \nabla E(w_t) + \alpha \Delta w_t \\
 &= -\eta \nabla E(w_t) + \alpha (-\eta \nabla E(w_{t-1}) + \alpha (-\eta \nabla E(w_{t-2}) + \dots)) \\
 &= -\eta \sum_{k=0}^t \alpha^k \nabla E(w_{t-k})
 \end{aligned}$$

Consider two extremes:

**No oscillations** all derivative are equal:

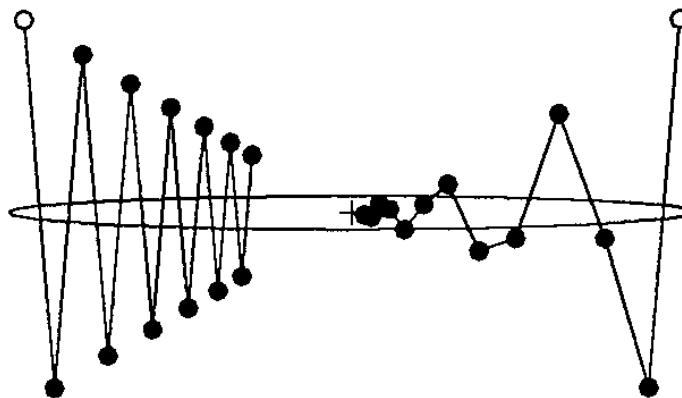
$$\Delta w_{t+1} \approx -\eta \nabla E \sum_{k=0}^t \alpha^k = -\frac{\eta}{1-\alpha} \frac{\partial E}{\partial w}$$

results in acceleration

**Oscillations** all derivatives are equal but have opposite sign:

$$\Delta w(t+1) \approx -\eta \nabla E \sum_{k=0}^t (-\alpha)^k = -\frac{\eta}{1+\alpha} \frac{\partial E}{\partial w}$$

results in deceleration



## Newtons method

One can also use Hessian information for optimization. As an example, consider a quadratic approximation to  $E$  around  $\mathbf{w}_0$ :

$$\begin{aligned} E(\mathbf{w}) &= E(\mathbf{w}_0) + \mathbf{b}^T(\mathbf{w} - \mathbf{w}_0) + \frac{1}{2}(\mathbf{w} - \mathbf{w}_0)\mathbf{H}(\mathbf{w} - \mathbf{w}_0) \\ b_i &= \frac{\partial E(\mathbf{w}_0)}{\partial w_i} \quad H_{ij} = \frac{\partial^2 E(\mathbf{w}_0)}{\partial w_i \partial w_j} \\ \nabla E(\mathbf{w}) &= \mathbf{b} + \mathbf{H}(\mathbf{w} - \mathbf{w}_0) \end{aligned}$$

We can solve  $\nabla E(\mathbf{w}) = 0$  and obtain

$$\mathbf{w} = \mathbf{w}_0 - \mathbf{H}^{-1}\nabla E(\mathbf{w}_0)$$

This is called Newtons method.

Quadratic approximation is exact when  $E$  is quadratic, so convergence in one step.

# Line search

Another solution is line optimisation:

$$w_1 = w_0 + \lambda d_0, \quad d_0 = \nabla E(w_0)$$

$\lambda$  is found by a one dimensional optimisation

$$0 = \frac{\partial}{\partial \lambda} E(w_0 + \lambda d_0) = d_0 \cdot \nabla E(w_1) = d_0 \cdot d_1$$

Therefore, subsequent search directions are orthogonal.



# Conjugate gradient descent

We choose as new direction a combination of the gradient and the old direction

$$d'_1 = \nabla E(w_1) + \beta d_0$$

Line optimisation  $w_2 = w_1 + \lambda d'_1$  yields  $\lambda$  such that  $d'_1 \cdot \nabla E(w_2) = 0$ .

The direction  $d'_1$  is found by demanding that  $\nabla E(w_2) \approx 0$  also in the 'old' direction  $d_0$ :

$$0 = d_0 \cdot \nabla E(w_2) \approx d_0 \cdot (\nabla E(w_1) + \lambda H(w_1)d'_1)$$

or

$$d_0 H(w_1) d'_1 = 0$$

$d_0, d'_1$  are said to be conjugate.



## Polak-Ribiere rule

The conjugate directions can be computed without computing the Hessian matrix, for instance using the Polak-Ribiere rule:<sup>7</sup>

$$\beta = \frac{(\nabla E(w_1) - \nabla E(w_0)) \cdot \nabla E(w_1)}{\|\nabla E(w_0)\|^2}$$

It can be proven that this rule keeps the last  $n$  directions all mutually conjugate [?].

---

<sup>7</sup>We need  $0 = d_0^T H(w_1) d_1'$ . We use  $\nabla E(w_0) \approx \nabla E(w_1) + (w_0 - w_1)^T H(w_1) = \nabla E(w_1) - d_0^T H(w_1)$ .

## Stochastic gradient descent

One can also consider on-line learning, where only one or a subset of training patterns is considered for computing  $\nabla E$ .

$$E(\mathbf{w}) = \sum_n E_n(\mathbf{w}) \quad \mathbf{w}_{t+1} = \mathbf{w}_t - \alpha_t \nabla E_n(\mathbf{w}^\tau)$$

May be efficient for large data sets. This results in a stochastic dynamics in  $\mathbf{w}$  that can help to escape local minima.

# Robbins Monro

Method of *stochastic approximation* originally due to Robbins and Monro 1951:

- Solve  $M(x) = a$  with  $M(x) = \langle N(x, \xi) \rangle$ .
- Iterate  $x_{t+1} = x_t + \alpha_t(a - N(x, \xi))$
- Convergence requires

$$\sum_t \alpha_t = \infty \quad \sum_t \alpha_t^2 < \infty$$

For instance  $\alpha_t = 1/t$ .

Application to stochastic gradient descent:

- $\nabla E(\mathbf{w}) = 0$  with  $\nabla E(\mathbf{w}) = \sum_n \nabla E_n(\mathbf{w})$
- Iterate  $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \nabla E_n(\mathbf{w})$

Extensions of SGD and comparisons see [?].

# Error backpropagation

Error is sum of error per pattern

$$E(\mathbf{w}) = \sum_n E^n(\mathbf{w}) \quad E^n(\mathbf{w}) = \frac{1}{2} \|\mathbf{y}(x_n, \mathbf{w}) - \mathbf{t}_n\|^2$$

$$\begin{aligned} y_k(\mathbf{x}, \mathbf{w}) &= h_2 \left( w_{k0} + \sum_{j=1}^M w_{kj} h_1 \left( w_{j0} + \sum_{i=1}^D w_{ji} x_i \right) \right) \\ &= h_2(a_k) \end{aligned}$$

$$a_k = w_{k0} + \sum_{j=1}^M w_{kj} h_1(a_j) = \sum_{j=0}^M w_{kj} h_1(a_j) \quad h_1(a_0) = 1$$

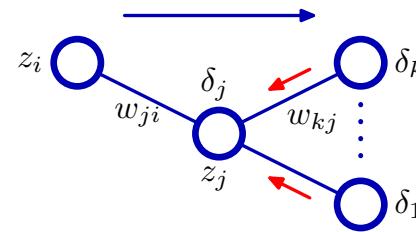
$$a_j = w_{j0} + \sum_{i=1}^D w_{ji} x_i = \sum_{i=0}^D w_{ji} x_i \quad x_0 = 1$$

## Error backpropagation

We do each pattern separately, so we consider  $E^n$

$$\begin{aligned}
 y_k(\mathbf{x}^n, \mathbf{w}) &= h_2(a_k^n) = h_2\left(\sum_{j=0}^M w_{kj} h_1(a_j^n)\right) = h_2\left(\sum_{j=0}^M w_{kj} h_1\left(\sum_{i=0}^D w_{ji} x_i^n\right)\right) \\
 \frac{\partial E^n}{\partial w_{kj}} &= (y_k^n - t_k^n) \frac{\partial y_k^n}{\partial w_{kj}} = (y_k^n - t_k^n) h'_2(a_k^n) \frac{\partial a_k^n}{\partial w_{kj}} = (y_k^n - t_k^n) h'_2(a_k^n) h_1(a_j^n) \\
 &= \delta_k^n h_1(a_j^n) \quad \delta_k^n = (y_k^n - t_k^n) h'_2(a_k^n) \\
 \frac{\partial E^n}{\partial w_{ji}} &= \sum_{k=1}^K (y_k^n - t_k^n) \frac{\partial y_k^n}{\partial w_{ji}} = \sum_{k=1}^K (y_k^n - t_k^n) h'_2(a_k^n) \frac{\partial a_k^n}{\partial w_{ji}} \\
 &= \sum_{k=1}^K \delta_k^n w_{kj} h'_1(a_j^n) \frac{\partial a_j^n}{\partial w_{ji}} = \sum_{k=1}^K \delta_k^n w_{kj} h'_1(a_j^n) x_i^n = \delta_j^n x_i^n \\
 \delta_j^n &= h'_1(a_j^n) \sum_{k=1}^K \delta_k^n w_{kj}
 \end{aligned}$$

## Error backpropagation



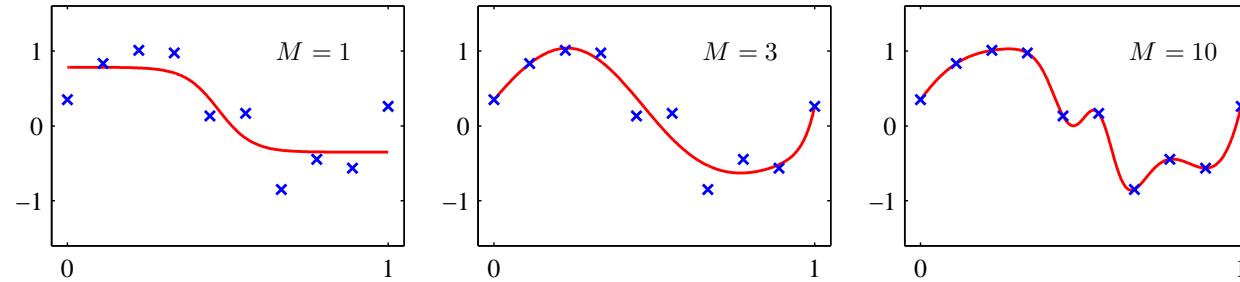
The back propagation extends to arbitrary layers:

1.  $z_i^n = x_i^n$  *forward propagation* all activations  $z_j^n = h_1(a_j^n)$  and  $z_k^n = h_2(a_k^n)$ , etc.
2. Compute the  $\delta_k^n$  for the output units, and *back-propagate* the  $\delta$  to obtain  $\delta_j^n$  each hidden unit  $j$
3.  $\partial E^n / \partial w_{kj} = \delta_k^n z_j^n$  and  $\partial E^n / \partial w_{ji} = \delta_j^n z_i^n$
4. for batch mode,  $\partial E / \partial w_{ji} = \sum_n \partial E^n / \partial w_{ji}$

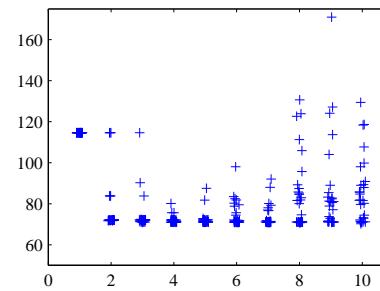
$E$  is a function of  $\mathcal{O}(|\mathbf{w}|)$  variables. In general, the computation of  $E$  requires  $\mathcal{O}(|\mathbf{w}|)$  operations. The computation of  $\nabla E$  would thus require  $\mathcal{O}(|\mathbf{w}|^2)$  operations.

The backpropagation method allows to compute  $\nabla E$  efficiently, in  $\mathcal{O}(|\mathbf{w}|)$  operations.

# Regularization



Complexity of neural network solution is controlled by number of hidden units



sum squared test error for different number of hidden units and different weight initializations. Error is also affected by local minima.

Part of the cause of local minima is the saturation of the sigmoid functions  $\tanh(\sum w_{ij}x_j)$ . When  $w_{ij}$  becomes large, any change in its value hardly affects the output, implying  $\nabla_{ij}E = 0$ .

One can partly prevent this from happening by

- choosing tanh instead of  $\sigma$  transfer functions
- scaling of inputs and outputs with mean zero and standard deviation one
- proper initialisation of  $w_{ij}$  with mean zero and standard deviation of order  $1/\sqrt{n_1}$ , with  $n_1$  the number of inputs to neuron  $i$ .
- add regularizer such as  $\sum_i w_i^2$  to cost keeps weights small

# MLPs are universal approximators

Consider  $2^n$  binary patterns in  $n$  dimensions and two classes:

$$x^\mu \rightarrow c^\mu = \pm 1, \quad x_i^\mu = \pm 1$$

Use  $2^n$  hidden units, labeled  $j = 0, \dots, 2^n - 1$ ,  $k$  labels input. Set

$$w_{jk} = b \quad \text{if } k\text{th digit in binary repr. of } j \text{ is 1}$$

$$w_{jk} = -b \quad \text{else}$$

$j$	binary	$w_{j1}$	$w_{j2}$
0	00	-b	-b
1	01	-b	b
2	10	b	-b
3	11	b	b

$x_1$	$x_2$	$w_{0k}x_k$	$w_{1k}x_k$	$w_{2k}x_k$	$w_{3k}x_k$
-1	-1	<b>2b</b>	0	0	-2b
-1	1	0	<b>2b</b>	-2b	0
1	-1	0	-2b	<b>2b</b>	0
1	1	-2b	0	0	<b>2b</b>

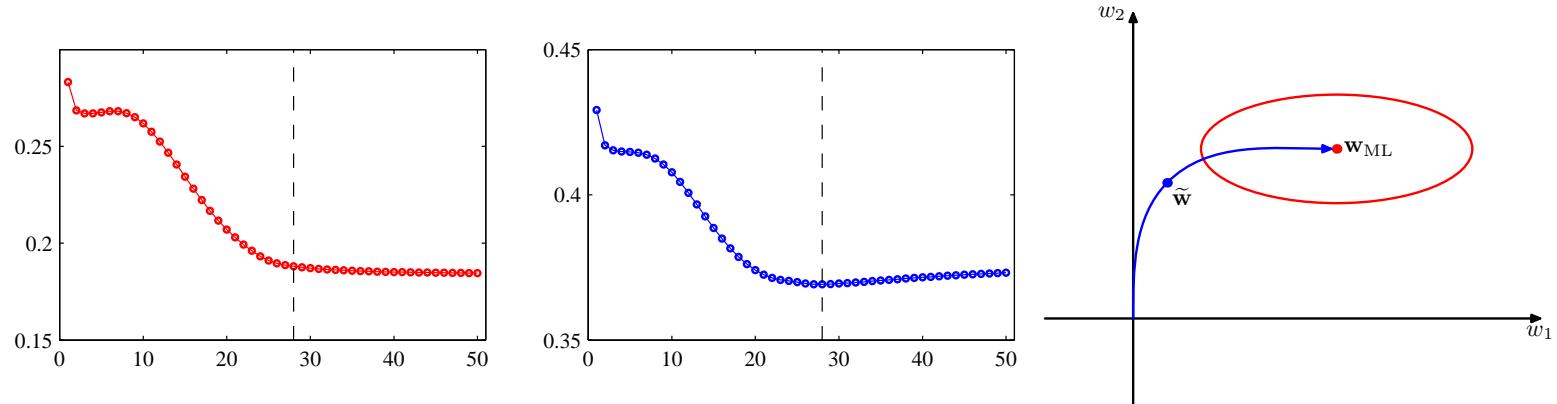
Use threshold of  $(n - 1)b$  at each hidden unit. The remaining problem has  $p = 2^n$  patterns in  $2^n$  dimensions and is linearly separable.

## MLPs are universal approximators

The combination of linear summation and non-linear functions can create many different functions.

- The MLP with a single hidden layer can map any continuous function [?, ?]
- The MLP with multiple hidden layers may (or may not) be more efficient

## 5.5.2 Early stopping



Early stopping is to stop training when error on test set starts increasing.

Early stopping with small initial weights has the effect of weight decay:

$$\begin{aligned} E(\mathbf{w}) &= \lambda_1(w_1 - w_1^*)^2 + \lambda_2(w_2 - w_2^*)^2 + \lambda(w_1^2 - w_2^2) \\ \frac{\partial E}{\partial w_i} &= \lambda_i(w_i - w_i^*) + \lambda w_i = 0, \quad i = 1, 2 \\ w_i &= \frac{\lambda_i}{\lambda_i + \lambda} w_i^* \end{aligned}$$

When  $\lambda_1 \ll \lambda \ll \lambda_2$ ,  $w_1 \approx \lambda_1/\lambda w_1^*$  and  $w_2 \approx w_2^*$ .

Weights in 'flat' directions are underspecified by the data and stay small.