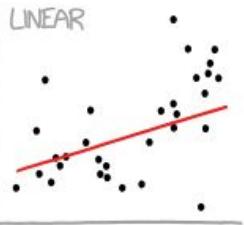


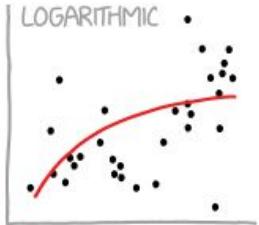
## CURVE-FITTING METHODS AND THE MESSAGES THEY SEND



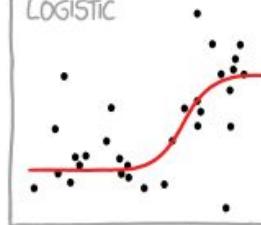
"HEY, I DID A  
REGRESSION."



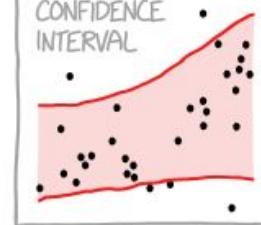
"I WANTED A CURVED  
LINE, SO I MADE ONE  
WITH MATH."



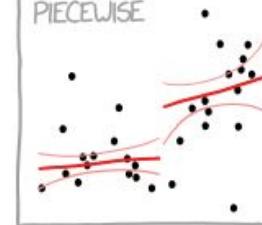
"LOOK, IT'S  
TAPERING OFF!"



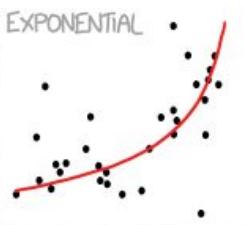
"I NEED TO CONNECT THESE  
TWO LINES, BUT MY FIRST IDEA  
DIDN'T HAVE ENOUGH MATH."



"LISTEN, SCIENCE IS HARD.  
BUT I'M A SERIOUS  
PERSON DOING MY BEST."



"I HAVE A THEORY,  
AND THIS IS THE ONLY  
DATA I COULD FIND."



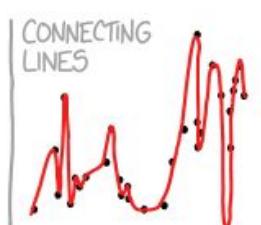
"LOOK, IT'S GROWING  
UNCONTROLLABLY!"



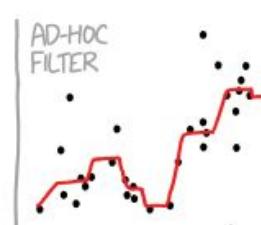
"I'M SOPHISTICATED, NOT  
LIKE THOSE BUMBLING  
POLYNOMIAL PEOPLE."



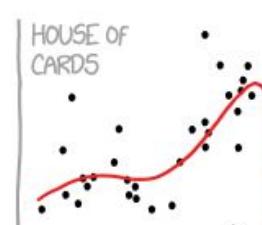
"I'M MAKING A  
SCATTER PLOT BUT  
I DON'T WANT TO."



"I CLICKED 'SMOOTH  
LINES' IN EXCEL."



"I HAD AN IDEA FOR HOW  
TO CLEAN UP THE DATA.  
WHAT DO YOU THINK?"



"AS YOU CAN SEE, THIS  
MODEL SMOOTHLY FITS  
THE— WAIT NO NO DON'T  
EXTEND IT AAAAAAA!!"

# Announcements

Assignment 1 out, we'll talk about it in a future session.

3 new tutors joined

Tutorials start this week

- max, learning -  
min suffering,

Hope everyone is doing okay!

# Linear Regression (Linear models for regression)

Why linear models?

Input data, features, basis functions

Maximum likelihood and least squares

Geometric intuition

Regularised least squares

Multiple outputs → book

- Bias-variance decomposition

interpretable / sparse linear model

The relation between MLE and least squares, Lagrange multipliers, multiple ways of looking at linear models.

# Why linear models? - it saves lives

The five criteria of the Apgar score:<sup>[3]</sup>

	Score of 0	Score of 1	Score of 2	Component of backronym
<b>Skin color</b>	blue or pale all over	blue at extremities, body pink (acrocyanosis)	no cyanosis body and extremities pink	Appearance
<b>Pulse rate</b>	absent	< 100 beats per minute	≥ 100 beats per minute	Pulse
<b>Reflex irritability grimace</b>	no response to stimulation	grimace on suction or aggressive stimulation	cry on stimulation	Grimace
<b>Muscle Tone</b>	none	some flexion	flexed arms and legs that resist extension	Activity
<b>Respiratory effort</b>	absent	weak, irregular, gasping	strong, robust cry	Respiration

## What do the Apgar scores mean?

A score of 7 or more is normal. A score of 6 or less at 1 minute and a score of 7 or more at 5 minutes is also normal. However, a score below 7 in the second test at 5 minutes is considered low.

If your baby's score was low in the first Apgar test and hasn't improved in the second test at 5 minutes, or there are other concerns, the doctors and nurses will closely monitor your baby and continue any necessary medical care.



Virginia Apgar, creator of the Apgar score

**Purpose** method to summarize newborn's health

Everything should be made as simple as possible, but not simpler. -- Albert Einstein

**Occam's razor**, also spelled **Ockham's razor**, also called **law of economy** or **law of parsimony**, principle stated by the Scholastic philosopher **William of Ockham** (1285–1347/49) that *pluralitas non est ponenda sine necessitate*, “plurality should not be posited without necessity.” The principle gives precedence to simplicity: of two competing theories, the simpler explanation of an entity is to be preferred. The principle is also expressed as “Entities are not to be multiplied beyond necessity.”

<https://www.britannica.com/topic/Occams-razor>

Variables	Values	Points
Mean Blood Pressure	<input type="button" value="▼"/>	<input type="text" value="0"/>
Lowest temperature	<input type="button" value="▼"/>	<input type="text" value="0"/>
P <sub>O<sub>2</sub></sub> (mmHg) / FIO <sub>2</sub> (%)	<input type="button" value="▼"/>	<input type="text" value="0"/>
Lowest serum pH	<input type="button" value="▼"/>	<input type="text" value="0"/>
Multiple seizures	<input type="button" value="▼"/>	<input type="text" value="0"/>
Urine output (mL/kg.h)	<input type="button" value="▼"/>	<input type="text" value="0"/>
<b>SNAP II :</b> <input type="text" value="0"/>		
Apgar score	<input type="button" value="▼"/>	
Birth weight	<input type="button" value="▼"/>	
Small for gestational age ( <a href="#">help</a> )	<input type="button" value="▼"/>	
<b>SNAPPE II :</b> <input type="text" value="0"/>	In-hospital mortality : <a href="#">see below</a>	Data are collected within the NICU

Research Letters

The Lancet, 2003

## CRIB II: an update of the clinical risk index for babies score

Dr Gareth Parry PhD <sup>a</sup>  , Janet Tucker PhD <sup>b</sup>, William Tarnow-Mordi MRCP <sup>c</sup>, for the UK Neonatal Staffing Study Collaborative Group \*

The clinical risk index for babies (CRIB) score is a risk-adjustment instrument used worldwide in [neonatal intensive care](#).<sup>1</sup> It was developed with data relating to babies born at less than 31 weeks' gestation, or 1500 g birthweight or lower, admitted for neonatal intensive care between 1988 and 1990. The appropriateness of CRIB with contemporary data has been questioned, since the score may now be poorly calibrated with mortality after neonatal intensive care, in which case it might be no better in prediction of mortality than birthweight or gestation alone. Furthermore, CRIB includes, as one of two measures of severity of illness, [fraction of inspired oxygen](#) (FiO<sub>2</sub>), which is not a true physiological measure because it is determined by the care team. CRIB also includes data up to 12 h after admission, thus potentially introducing early treatment bias.

# Linear models and likelihood

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 x_1 + \dots + w_D x_D$$
$$\mathbf{w}^T \vec{\mathbf{x}}$$
$$\vec{\mathbf{x}} \approx \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_D \end{bmatrix}$$

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

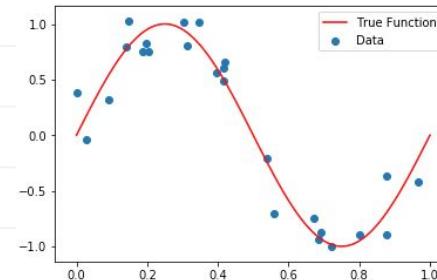
Why linear regression?

- Analytic solution when minimising sum of squared errors
- Well understood statistical behaviour
- Efficient algorithms exist for convex losses and regularizers

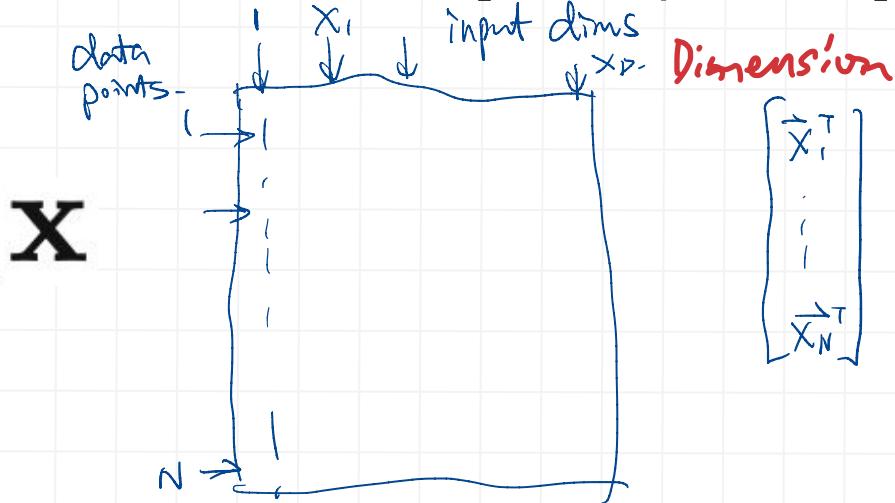
But what if ?

- Linear combination of fixed nonlinear basis functions

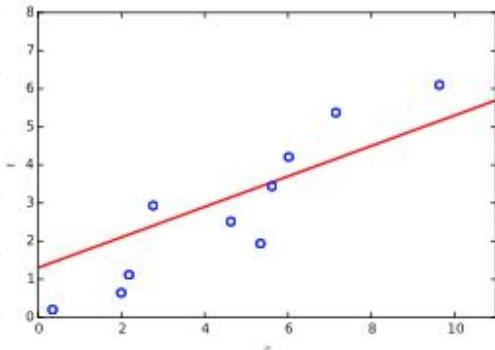
- parameter  $\mathbf{w} = (w_0, \dots, w_{M-1})^T$
- basis functions  $\phi(\mathbf{x}) = (\phi_0(\mathbf{x}), \dots, \phi_{M-1}(\mathbf{x}))^T$
- convention  $\phi_0(\mathbf{x}) = 1$
- $w_0$  is the bias parameter



# Conventions in [Bishop 2006]



$$\vec{x}_i = \begin{bmatrix} 1 \\ x_{i1} \\ \vdots \\ x_{iD} \end{bmatrix}$$



$\phi_j(\mathbf{x})$

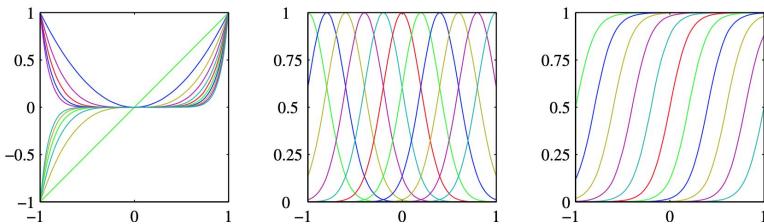
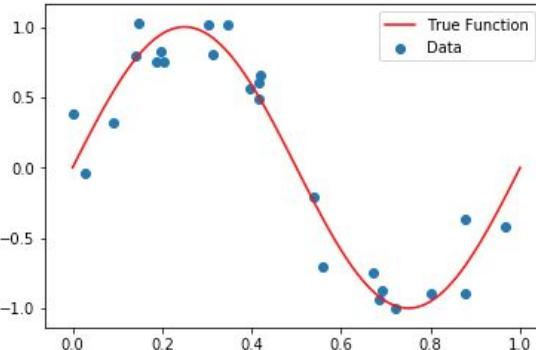


Figure 3.1 Examples of basis functions, showing polynomials on the left, Gaussians of the form (3.4) in the centre, and sigmoidal of the form (3.5) on the right.

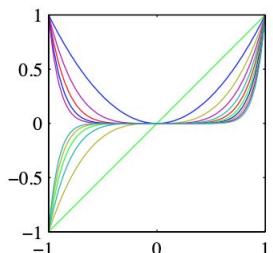


## Polynomial Basis Functions

- Scalar input variable  $x$

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

- Limitation : Polynomials are global functions of the input variable  $x$  so the learned function will extrapolate poorly

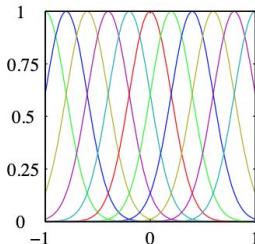


## 'Gaussian' Basis Functions

- Scalar input variable  $x$

$$\phi_j(x) = \exp\left\{-\frac{(x - \mu_j)^2}{2s^2}\right\}$$

- Not a probability distribution.
- No normalisation required, taken care of by the model parameters  $w$ .
- Well behaved away from the data (though pulled to zero)



## Sigmoidal Basis Functions

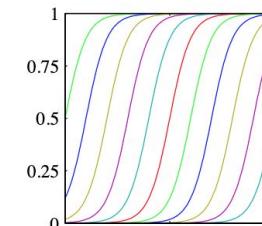
- Scalar input variable  $x$

$$\phi_j(x) = \sigma\left(\frac{x - \mu_j}{s}\right)$$

where  $\sigma(a)$  is the logistic sigmoid function defined by

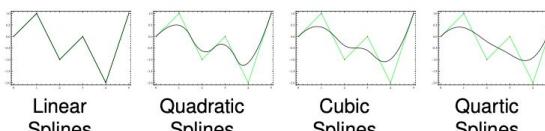
$$\sigma(a) = \frac{1}{1 + \exp(-a)}$$

- $\sigma(a)$  is related to the [hyperbolic tangent](#)  $\tanh(a)$  by  $\tanh(a) = 2\sigma(a) - 1$ .



## Other Basis Functions

- Fourier Basis : each basis function represents a specific frequency and has infinite spatial extent.
- Wavelets : localised in both space and frequency (also mutually orthogonal to simplify application).
- Splines (piecewise polynomials restricted to regions of the input space; additional constraints where pieces meet, e.g. smoothness constraints  $\rightarrow$  conditions on the derivatives).



Approximate the points  
 $\{(0, 0), (1, 1), (2, -1), (3, 0), (4, -2), (5, 1)\}$  by different splines.

# Linear models and likelihood

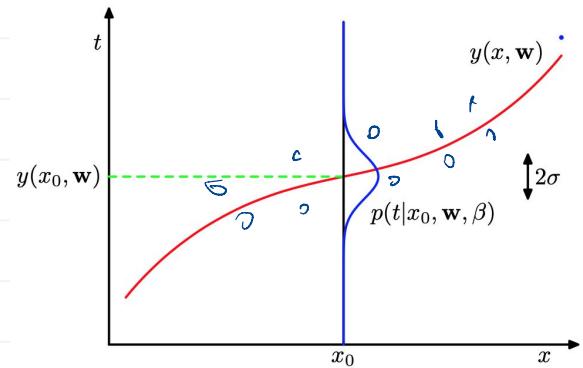
$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 x_1 + \dots + w_D x_D$$

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

$$t = y(\mathbf{x}, \mathbf{w}) + \epsilon \quad (3.7)$$

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

Fig 1.16



(1.60) also (3.8)

## Computing the likelihood

$$p(t|x, \mathbf{w}, \beta) = \mathcal{N}(t|y(x, \mathbf{w}), \beta^{-1}) \quad (1.60)$$

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | \mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1})$$

$$\begin{aligned} \ln p(\mathbf{t} | \mathbf{w}, \beta) &= \sum_{n=1}^N \ln \mathcal{N}(t_n | \mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1}) \\ &= \sum_{n=1}^N \ln \left( \sqrt{\frac{\beta}{2\pi}} \exp \left\{ -\frac{\beta}{2} (t_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2 \right\} \right) \end{aligned}$$

= ... (see next page)

# Computing the likelihood

$$\begin{aligned}\ln p(\mathbf{t} \mid \mathbf{w}, \beta) &= \sum_{n=1}^N \ln \mathcal{N}(t_n \mid \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1}) \\ &= \sum_{n=1}^N \ln \left( \sqrt{\frac{\beta}{2\pi}} \exp \left\{ -\frac{\beta}{2} (t_n - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n))^2 \right\} \right) \\ &= \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta E_D(\mathbf{w})\end{aligned}\tag{3.11}$$

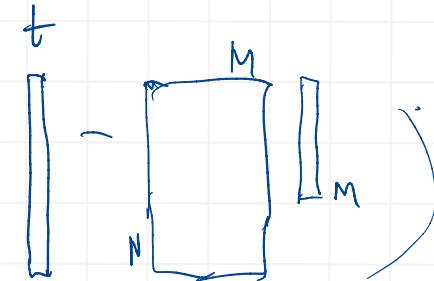
the sum-of-squares error function is defined by

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n)\}^2.\tag{3.12}$$

Claim:  $\arg \max_{\mathbf{w}} \ln p(\mathbf{t} \mid \mathbf{w}, \beta) \rightarrow \arg \min_{\mathbf{w}} E_D(\mathbf{w})$

... first rewrite the error function

$$\begin{aligned} E_D(\mathbf{w}) &= \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2. \\ &= \frac{1}{2} (\mathbf{t} - \Phi \mathbf{w})^T (\mathbf{t} - \Phi \mathbf{w}) \end{aligned} \quad (3.12)$$



where  $\mathbf{t} = (t_1, \dots, t_N)^T$ , and

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

plug into (3.11)  $\ln p(\mathbf{t} | \mathbf{w}, \beta) = \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta E_D(\mathbf{w})$

Then find the stationary point

$$\begin{aligned}\ln p(\mathbf{t} | \mathbf{w}, \beta) &= \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta E_D(\mathbf{w}) \\ &= \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta \frac{1}{2} (\mathbf{t} - \Phi \mathbf{w})^T (\mathbf{t} - \Phi \mathbf{w})\end{aligned}$$

The gradient with respect to  $\mathbf{w}$  is

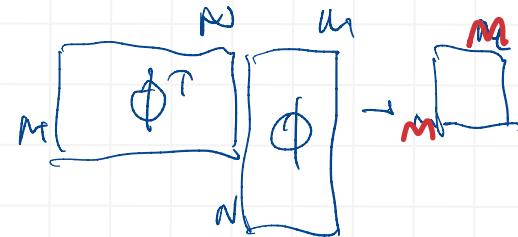
$$\nabla_{\mathbf{w}} \ln p(\mathbf{t} | \mathbf{w}, \beta) = \beta \Phi^T (\mathbf{t} - \Phi \mathbf{w}).$$

Setting the gradient to zero gives

$$0 = \Phi^T \mathbf{t} - \Phi^T \Phi \mathbf{w},$$

The normal equation:

$$\mathbf{w}_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$



(3.15)

# Obtaining and interpreting MLE results

$$w_0 = \bar{t} - \sum_{j=1}^{M-1} w_j \bar{\phi}_j \quad \bar{t} = \frac{1}{N} \sum_{n=1}^N t_n, \quad \bar{\phi}_j = \frac{1}{N} \sum_{n=1}^N \phi_j(\mathbf{x}_n).$$

Recall:  $\ln p(\mathbf{t} | \mathbf{w}, \beta) = \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta \frac{1}{2} (\mathbf{t} - \Phi \mathbf{w})^T (\mathbf{t} - \Phi \mathbf{w})$

$$\frac{1}{\beta_{\text{ML}}} = \frac{1}{N} \sum_{n=1}^N \{t_n - \mathbf{w}_{\text{ML}}^T \boldsymbol{\phi}(\mathbf{x}_n)\}^2 \quad (3.21)$$

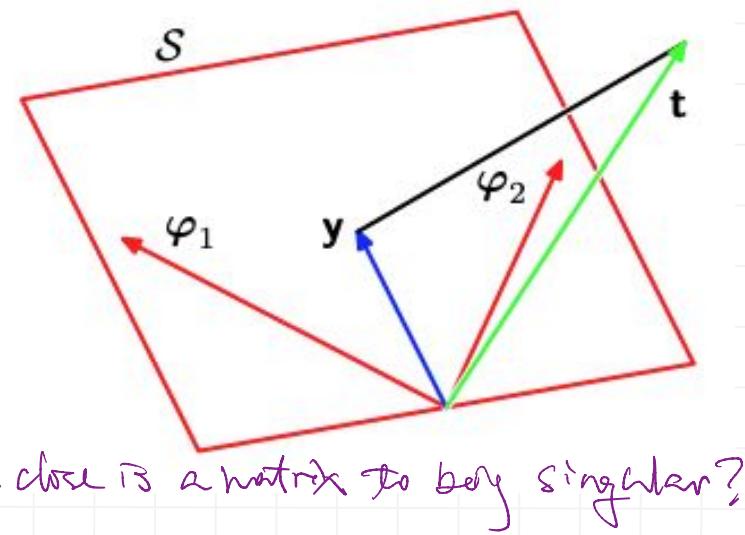


The normal equation:

$$\mathbf{w}_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t} \quad (3.15)$$

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

**Figure 3.2** Geometrical interpretation of the least-squares solution, in an  $N$ -dimensional space whose axes are the values of  $t_1, \dots, t_N$ . The least-squares regression function is obtained by finding the orthogonal projection of the data vector  $\mathbf{t}$  onto the subspace spanned by the basis functions  $\phi_j(\mathbf{x})$  in which each basis function is viewed as a vector  $\varphi_j$  of length  $N$  with elements  $\phi_j(\mathbf{x}_n)$ .



Numerical difficulties when  $\Phi^T \Phi$  is ill-conditioned.

SVD or regularisation will help.

Cool geometric derivation of the normal equation  
<https://www.youtube.com/watch?v=PbyP3goun2Y>

# Sequential Learning - Stochastic Gradient Descent

- For large data sets, calculating the maximum likelihood parameters  $\mathbf{w}_{ML}$  and  $\beta_{ML}$  may be costly.
- For online applications, never all data in memory.
- Use a sequential algorithms (online algorithm).
- If the error function is a sum over data points  $E = \sum_n E_n$ , then

- ➊ initialise  $\mathbf{w}^{(0)}$  to some starting value
- ➋ update the parameter vector at iteration  $\tau + 1$  by

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_n, \quad = \mathbf{w}^{(\tau)} + \eta \left( t_n - \mathbf{w}^{(\tau)T} \phi(\mathbf{x}_n) \right) \phi(\mathbf{x}_n)$$

(learning rate)

Sum-of-squares error

where  $E_n$  is the error function after presenting the  $n$ th data set, and  $\eta$  is the learning rate.

# What we did so far

Why linear models?

Input data, features, basis functions

Maximum likelihood and least squares

Geometric intuition, sequential update

Regularised least squares

Multiple outputs

Bias-variance decomposition

# Regularised least squares

$$L(\mathbf{w}) = E_D(\mathbf{w}) + \lambda E_W(\mathbf{w}) \quad (3.24)$$

$$\frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}. \quad (3.27) \quad \text{also (1.4)}$$

$$\mathbf{w} = (\lambda \mathbf{I} + \Phi^T \Phi)^{-1} \Phi^T \mathbf{t}. \quad (3.28)$$

$\uparrow$   
better "conditioned"

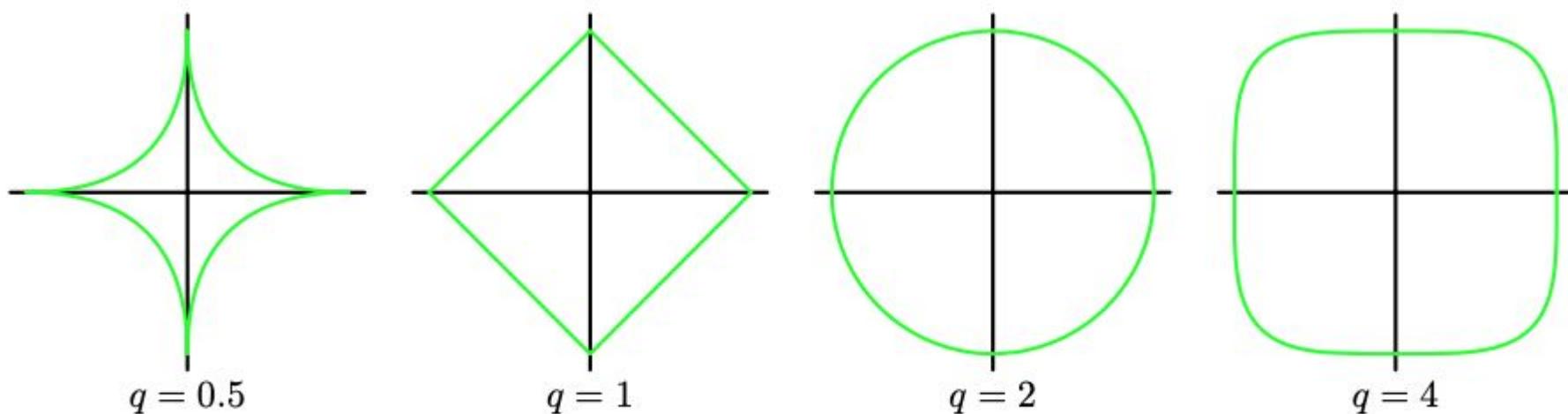
eigen vals  $\Rightarrow \lambda$

$$\frac{\partial L}{\partial \mathbf{w}} = (\text{---}) + \lambda \mathbf{w}$$

MLF,  $\downarrow$

# Regularisation by q-norm

$$\frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \sum_{j=1}^M |w_j|^q \quad (3.29)$$



**Figure 3.3** Contours of the regularization term in (3.29) for various values of the parameter  $q$ .

# Lagrange multipliers (appendix E)

The first encounter in SML - we'll see it again in kernel methods.

objective function  
equality constraint

maximize  $f(x)$   
subject to  $g(x)=0$

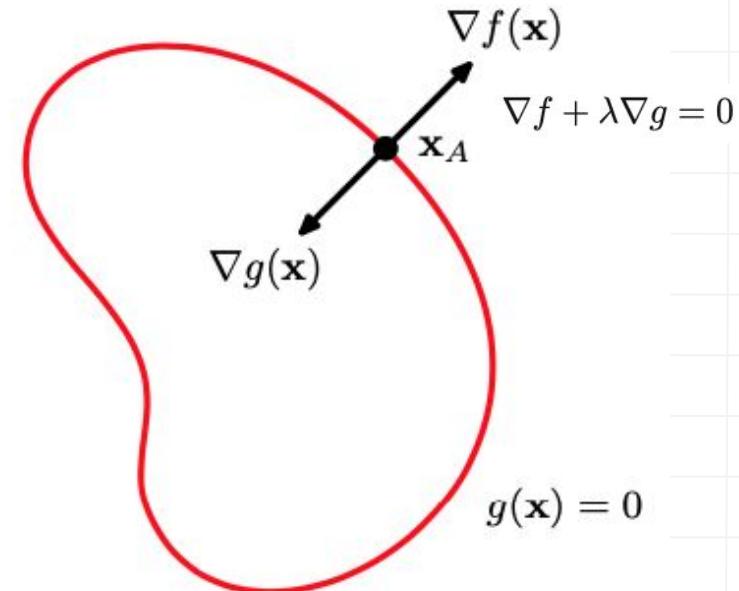


$$\mathcal{L}(x, \lambda) = f(x) + \lambda g(x)$$

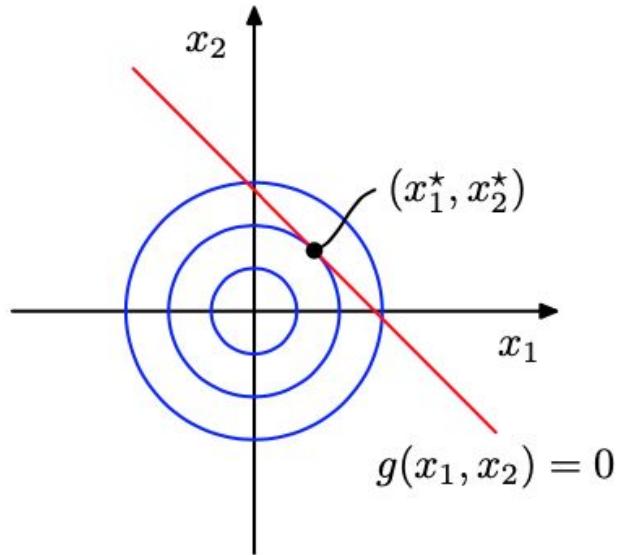
$\mathcal{L}$  = Lagrangian

$\lambda$  = Lagrange multiplier

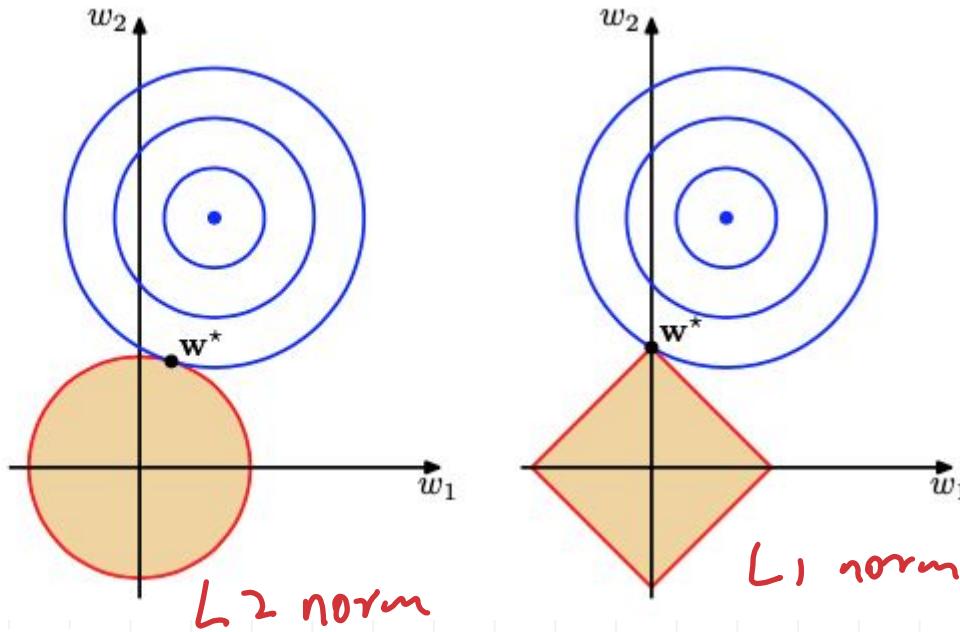
**Figure E.1** A geometrical picture of the technique of Lagrange multipliers in which we seek to maximize a function  $f(\mathbf{x})$ , subject to the constraint  $g(\mathbf{x}) = 0$ . If  $\mathbf{x}$  is  $D$  dimensional, the constraint  $g(\mathbf{x}) = 0$  corresponds to a subspace of dimensionality  $D - 1$ , indicated by the red curve. The problem can be solved by optimizing the Lagrangian function  $L(\mathbf{x}, \lambda) = f(\mathbf{x}) + \lambda g(\mathbf{x})$ .



**Figure E.2** A simple example of the use of Lagrange multipliers in which the aim is to maximize  $f(x_1, x_2) = 1 - x_1^2 - x_2^2$  subject to the constraint  $g(x_1, x_2) = 0$  where  $g(x_1, x_2) = x_1 + x_2 - 1$ . The circles show contours of the function  $f(x_1, x_2)$ , and the diagonal line shows the constraint surface  $g(x_1, x_2) = 0$ .



**Figure 3.4** Plot of the contours of the unregularized error function (blue) along with the constraint region (3.30) for the quadratic regularizer  $q = 2$  on the left and the lasso regularizer  $q = 1$  on the right, in which the optimum value for the parameter vector  $w$  is denoted by  $w^*$ . The lasso gives a sparse solution in which  $w_1^* = 0$ .



$$\frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \sum_{j=1}^M |w_j|^q$$

# Multiple regression outputs

See book

# What we did so far

Why linear models?

Input data, features, basis functions

Maximum likelihood and least squares

Geometric intuition

Regularised least squares

Multiple outputs

Bias-variance decomposition

# Bias-variance: the cartoon view

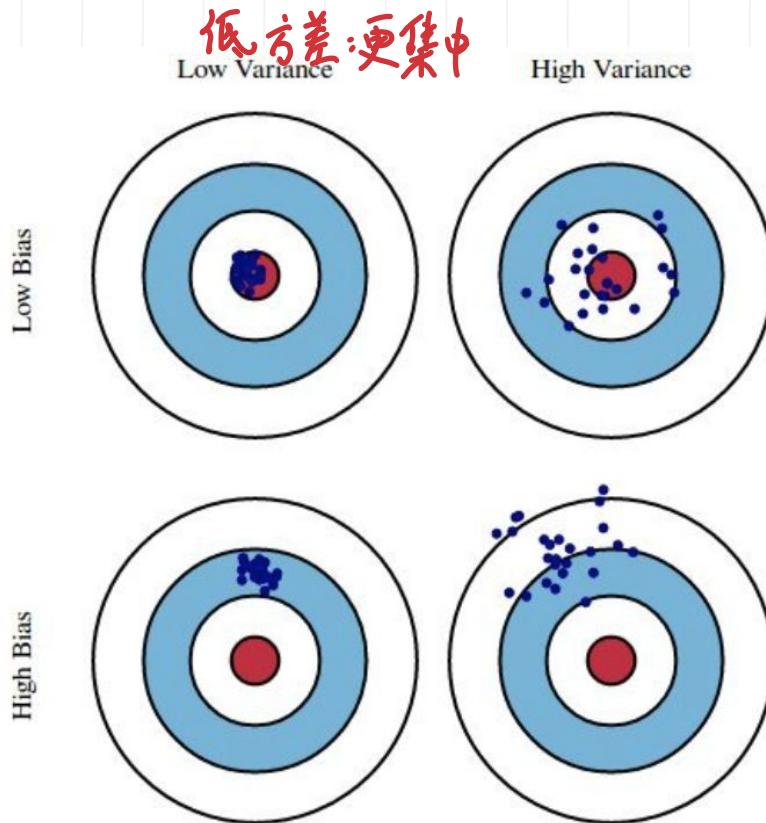


Fig. 1: Graphical illustration of bias and variance  
From Understanding the Bias-Variance Tradeoff, by Scott Fortmann-Roe.

What are the sources of different dots in this picture?

What remains constant: learning target, model specification, learning algorithm.

What changes:  
Data (subsets), randomness in learning (including but are not limited to initialisations), ...

# Expected square loss

(See Sec 1.5.5)

**Figure 1.28**

The regression function  $y(x)$ , which minimizes the expected squared loss, is given by the mean of the conditional distribution  $p(t|x)$ .

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

$$h(\mathbf{x}) = \mathbb{E}[t|\mathbf{x}] = \int tp(t|\mathbf{x}) dt. \quad (3.36)$$

↑  
fn of  $\mathbf{x}$ , not  $t$

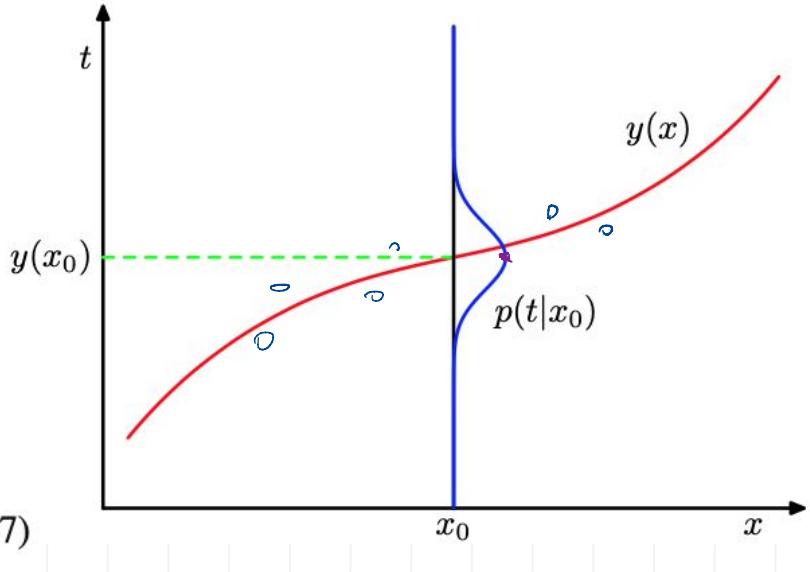
$$\mathbb{E}[L] = \iint \{y(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt. \quad (1.87)$$

↑  
 $\mathbf{x}, t$

$$\begin{aligned} \{y(\mathbf{x}) - t\}^2 &= \{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}] + \mathbb{E}[t|\mathbf{x}] - t\}^2 \\ &= \{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]\}^2 + 2\{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]\}\{\mathbb{E}[t|\mathbf{x}] - t\} + \{\mathbb{E}[t|\mathbf{x}] - t\}^2 \end{aligned}$$

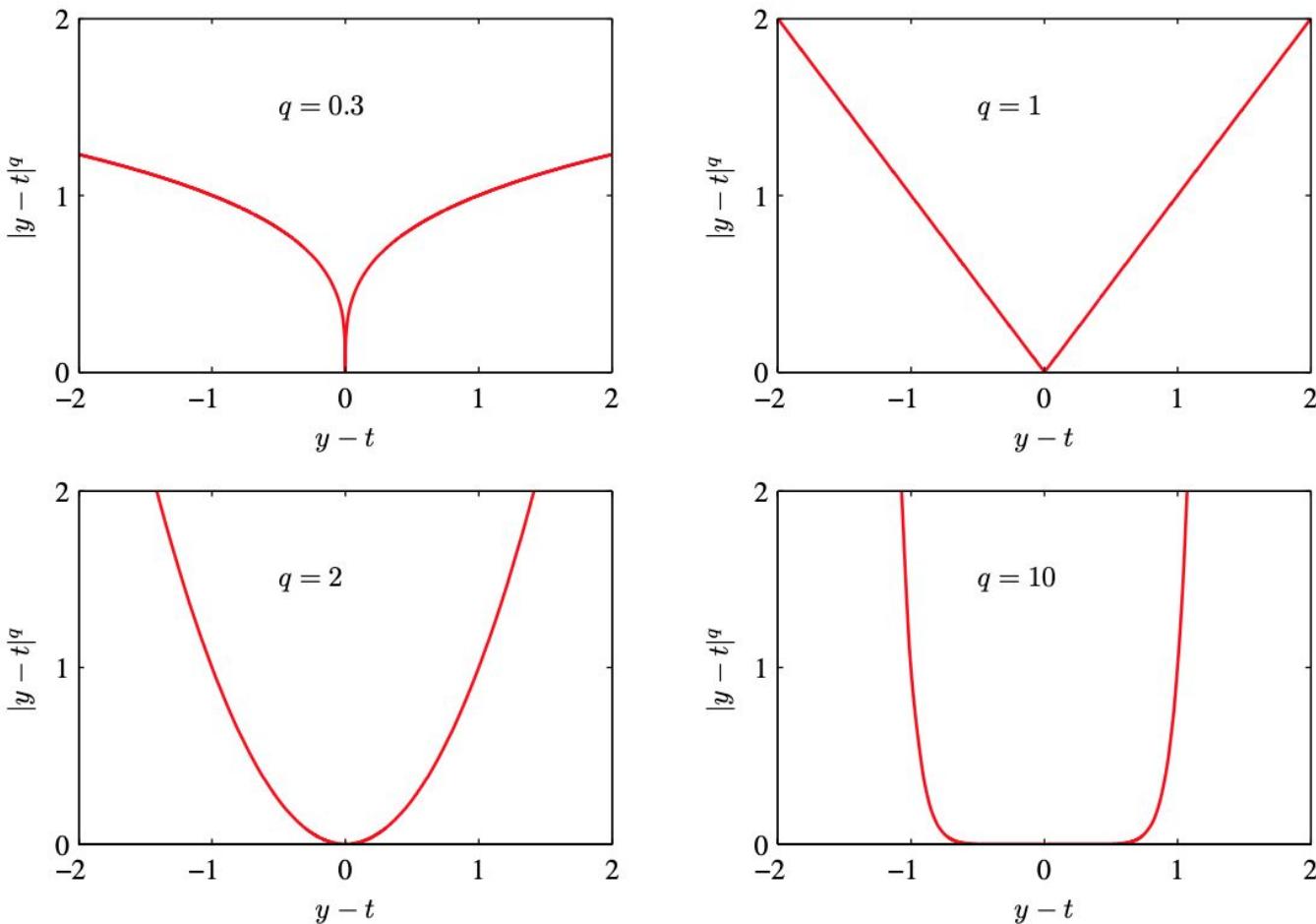
$$\mathbb{E}[L] = \int \{y(\mathbf{x}) - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x} + \int \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt. \quad (3.37)$$

*noise*



$$\mathbb{E}[L] = \int \{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]\}^2 p(\mathbf{x}) d\mathbf{x} + \int \{\mathbb{E}[t|\mathbf{x}] - t\}^2 p(\mathbf{x}) d\mathbf{x}. \quad (1.90)$$

## Generalising squared loss to Mikowski distances



**Figure 1.29** Plots of the quantity  $L_q = |y - t|^q$  for various values of  $q$ .

$$\mathbb{E}[L] = \int \{y(\mathbf{x}) - \mathbb{E}[t|\mathbf{x}]\}^2 p(\mathbf{x}) d\mathbf{x} + \int \{\mathbb{E}[t|\mathbf{x}] - t\}^2 p(\mathbf{x}) d\mathbf{x}. \quad (1.90)$$

$$h(\mathbf{x}) = \mathbb{E}[t|\mathbf{x}] = \int tp(t|\mathbf{x}) dt. \quad (3.36)$$

$$\mathbb{E}[L] = \int \{y(\mathbf{x}) - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x} + \int \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt. \quad (3.37)$$

$$\mathbb{E}[L] = \int_{\mathbf{x}} \{y(\mathbf{x}) - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x} + \int_{\mathbf{x}, t} \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt. \quad (3.37)$$

Int. over both  $\mathbf{x}$  &  $t$

Now introduce data  $D$

$$\{y(\mathbf{x}; D) - h(\mathbf{x})\}^2. \quad (3.38)$$

$$\begin{aligned} &= \{y(\mathbf{x}; D) - \mathbb{E}_D[y(\mathbf{x}; D)] + \mathbb{E}_D[y(\mathbf{x}; D)] - h(\mathbf{x})\}^2 \\ &= \{y(\mathbf{x}; D) - \mathbb{E}_D[y(\mathbf{x}; D)]\}^2 + \{\mathbb{E}_D[y(\mathbf{x}; D)] - h(\mathbf{x})\}^2 \\ &\quad + 2\{y(\mathbf{x}; D) - \mathbb{E}_D[y(\mathbf{x}; D)]\}\{\mathbb{E}_D[y(\mathbf{x}; D)] - h(\mathbf{x})\}. \end{aligned} \quad (3.39)$$

Int. over  $\mathbf{x}$

Q: how does one know to add+subtract  $\mathbb{E}_D[y(\mathbf{x}; D)]$ ?

$$\begin{aligned} &\mathbb{E}_D [\{y(\mathbf{x}; D) - h(\mathbf{x})\}^2] \\ &= \underbrace{\{\mathbb{E}_D[y(\mathbf{x}; D)] - h(\mathbf{x})\}^2}_{\text{(bias)}^2} + \underbrace{\mathbb{E}_D [\{y(\mathbf{x}; D) - \mathbb{E}_D[y(\mathbf{x}; D)]\}^2]}_{\text{variance}}. \end{aligned} \quad (3.40)$$

$$\mathbb{E}[L] = \int \{y(\mathbf{x}) - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x} + \int \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt. \quad (3.37)$$



$$\begin{aligned} & \mathbb{E}_{\mathcal{D}} [\{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}^2] \\ &= \underbrace{\{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2}_{\text{(bias)}^2} + \underbrace{\mathbb{E}_{\mathcal{D}} [\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2]}_{\text{variance}}. \end{aligned} \quad (3.40)$$

expected loss = (bias)<sup>2</sup> + variance + noise

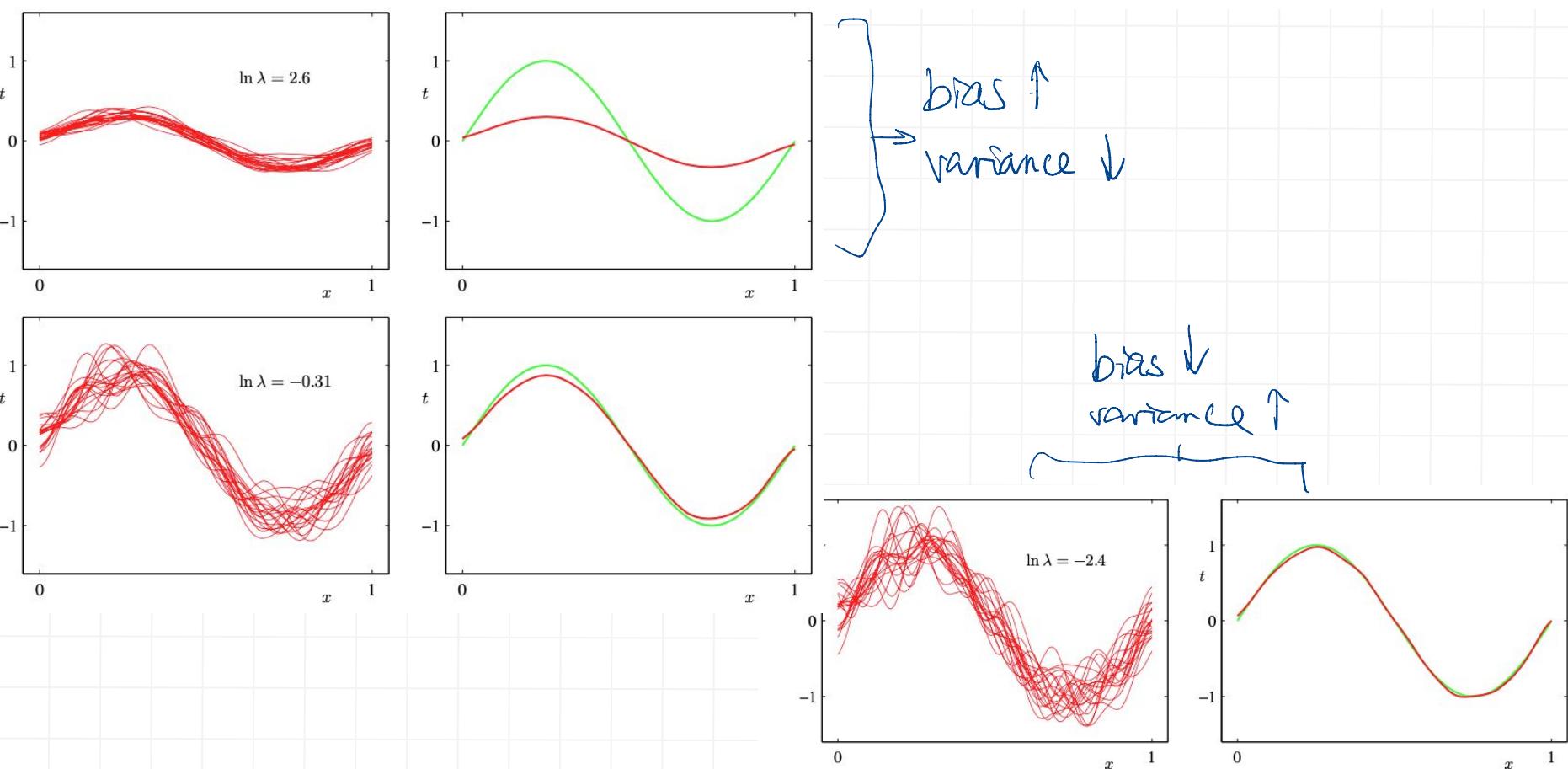
(3.41)

where

$$(\text{bias})^2 = \int \{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x} \quad (3.42)$$

$$\text{variance} = \int \mathbb{E}_{\mathcal{D}} [\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2] p(\mathbf{x}) d\mathbf{x} \quad (3.43)$$

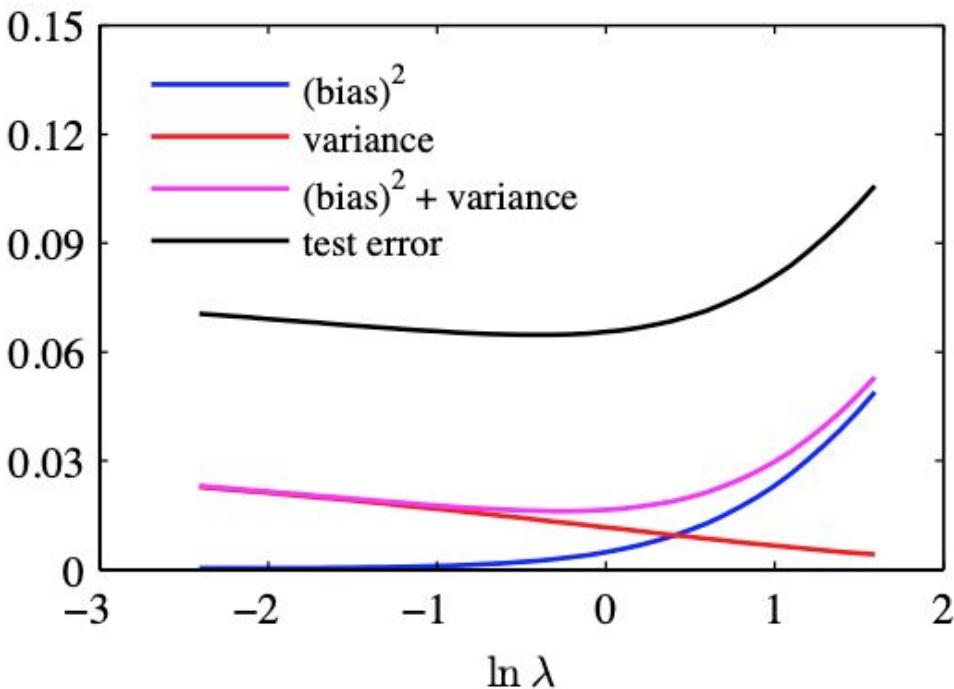
$$\text{noise} = \int \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt \quad (3.44)$$



**Figure 3.5** Illustration of the dependence of bias and variance on model complexity, governed by a regularization parameter  $\lambda$ , using the sinusoidal data set from Chapter 1. There are  $L = 100$  data sets, each having  $N = 25$  data points, and there are 24 Gaussian basis functions in the model so that the total number of parameters is  $M = 25$  including the bias parameter. The left column shows the result of fitting the model to the data sets for various values of  $\ln \lambda$  (for clarity, only 20 of the 100 fits are shown). The right column shows the corresponding average of the 100 fits (red) along with the sinusoidal function from which the data sets were generated (green).

- Dependence of bias and variance on the model complexity
- Squared bias, variance, their sum, and test data
- The minimum for  $(\text{bias})^2 + \text{variance}$  occurs close to the value that gives the minimum error

Fig 3.6



# Unbiased estimators

- You may have encountered *unbiased estimators*
- Why guarantee zero bias? To quote the pioneer of Bayesian inference, Edwin Jaynes, from his book *Probability Theory: The Logic of Science* (2003):

**Why do they do this?** Why do orthodoxians put such exaggerated emphasis on bias? We suspect that the main reason is simply that they are caught in a **psycho-semantic trap** of their own making. When we call the quantity  $(\langle \beta \rangle - \alpha)$  the “bias”, that makes it sound like something awfully reprehensible, which we must get rid of at all costs. If it had been called instead the **“component of error orthogonal to the variance”**, as suggested by the Pythagorean form of (17–2), it would have been clear to all that these two contributions to the error are on an equal footing; it is folly to decrease one at the expense of increasing the other. This is just the price one pays for choosing a technical terminology that carries an emotional load, implying value judgments; orthodoxy falls constantly into this tactical error.

# The bias-variance decomposition

- Tradeoff between bias and variance
  - simple models have low variance and high bias
  - complex models have high variance and low bias
- The sum of bias and variance has a minimum at a certain model complexity.
- Expected loss  $\mathbb{E}_{\mathcal{D}} [L]$  over all data sets  $\mathcal{D}$

$$\text{expected loss} = (\text{bias})^2 + \text{variance} + \text{noise}.$$

- The noise comes from the data, and can not be removed from the expected loss.
- To analyse the bias-variance decomposition : many data sets needed, which are not always available.

# “explainable models”? Interpretable models

Dawes, Robyn M.. "The robust beauty of improper linear models in decision making." *American Psychologist* 34 (1979): 571-582.

Ustun, Berk and Cynthia Rudin. "Supersparse linear integer models for optimized medical scoring systems." *Machine Learning* 102 (2015): 349-391.

Rudin, Cynthia. "Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead." *Nature Machine Intelligence* 1 (2019): 206-215.



Cynthia Rudin, Duke University

2022 Squirrel Prize in AI for Humanity

<https://aaai.org/Awards/squirrel-ai-award.php>

For pioneering scientific work in the area of interpretable and transparent AI systems in real-world deployments, the advocacy for these features in highly sensitive areas such as social justice and medical diagnosis, and serving as a role model for researchers and practitioners.

Video: <https://www.youtube.com/watch?v=PwLN5irdMT8>

- An interpretable machine learning model obeys a domain-specific set of constraints that makes its computations easier to understand.
- My technical definition: An interpretable machine learning model is constrained in model form so that it is either useful to someone, or obeys structural knowledge of the domain, such as monotonicity, causality, structural (generative) constraints, additivity, or physical constraints that come from domain knowledge.
- There's a spectrum.

# Example: super-sparse linear models

Slide by Cynthia Rudin

- 2HELP2B was not created by doctors
- It is a ML model
- It is just as accurate as black box models.
- Doctors can decide themselves whether to trust it
- Doctors can calibrate the score with information not in the database

## 2HELP2B

1. Any cEEG Pattern with Frequency <b>2 Hz</b>	1 point	...
2. Epileptiform Discharges	1 point	+
3. Patterns include [LPD, LRDA, BIPD]	1 point	+
4. Patterns Superimposed with Fast or Sharp Activity	1 point	+
5. Prior Seizure	1 point	+
6. Brief Rhythmic Discharges	2 points	+
	<b>SCORE</b>	= ...

SCORE	0	1	2	3	4	5	6+
RISK	<5%	11.9%	26.9%	50.0%	73.1%	88.1%	95.3%

# Risk-Calibrated Supersparse Linear Integer Models (Risk-SLIM)

(Ustun, R, 2019)

$$\min_{\lambda \in L} \sum_{i=1}^n \log(1 + e^{-y_i \mathbf{x}_i^\top \lambda}) + C \|\lambda\|_0$$

$\lambda \in L$  means that  $\forall j, \lambda_j \in \{-10, -9, \dots, 0, \dots, 9, 10\}$

The diagram shows the Risk-SLIM objective function as a sum of two terms. The first term,  $\sum_{i=1}^n \log(1 + e^{-y_i \mathbf{x}_i^\top \lambda})$ , is bracketed by a brace above it labeled "Logistic Loss". The second term,  $C \|\lambda\|_0$ , is bracketed by a brace above it labeled "Model Size". Below the equation, a brace groups the two terms and points to the text "Small Integer Coefficients" located at the bottom right.

MINLP – really hard...

Logistic Loss

Model Size

Small Integer Coefficients

# So far...

- 2HELP2B validated on independent multicenter cohort (N=2111)
- Implemented: University of Wisconsin, Massachusetts General Hospital/Harvard Medical School
- Ongoing implementation: Emory University, Duke University, Medical University of South Carolina, Free University of Brussels (Belgium)
- Resulted in **63.6%** reduction in duration of EEG monitoring per patient
  - \$1,134.831 saving per patient<sup>1</sup>
- **2.82 X** More Patients Monitored
- **\$6.1M** estimated savings in FY 2018 at MGH,UW

<sup>1</sup>2016 Medicare Reimbursement Most Common Professional Code

# Linear models for regression 1

Why linear models?

Input data, features, basis functions

Maximum likelihood and least squares

Geometric intuition

Regularised least squares

Multiple outputs

Bias-variance decomposition

The relation between MLE and least squares, Lagrange multipliers, multiple ways of looking at linear models.

