# Advanced Theory Seminar
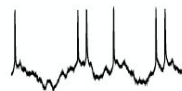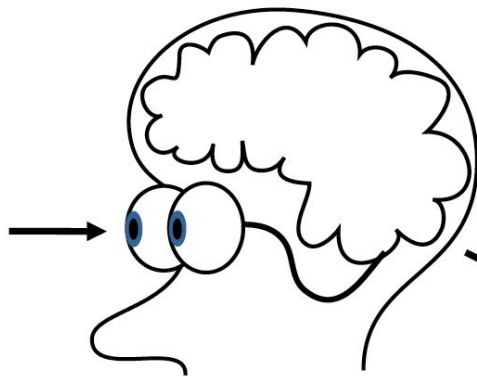# Linear and non-linear regression

Juri Minxha

March 18, 2020

# Sources

1.  *Pattern recognition and Machine Learning*, **Christopher M. Bishop**
    a.  predominantly chapter 3: Linear Models for Regression
2.  *Statistical Models for Neural Data: from Regression/GLMs to Latent Variables*, **Jonathan Pillow**
    a.  Cosyne 2018 tutorial
3.  *Machine learning: A probabilistic perspective,* **Kevin Murphy**
    a.  predominantly chapter 7
4.  *mathematicalmonk* lectures on Youtube (highly recommend), **Jeff Miller**
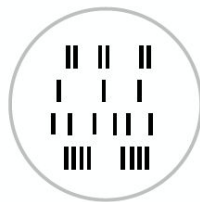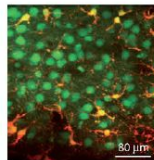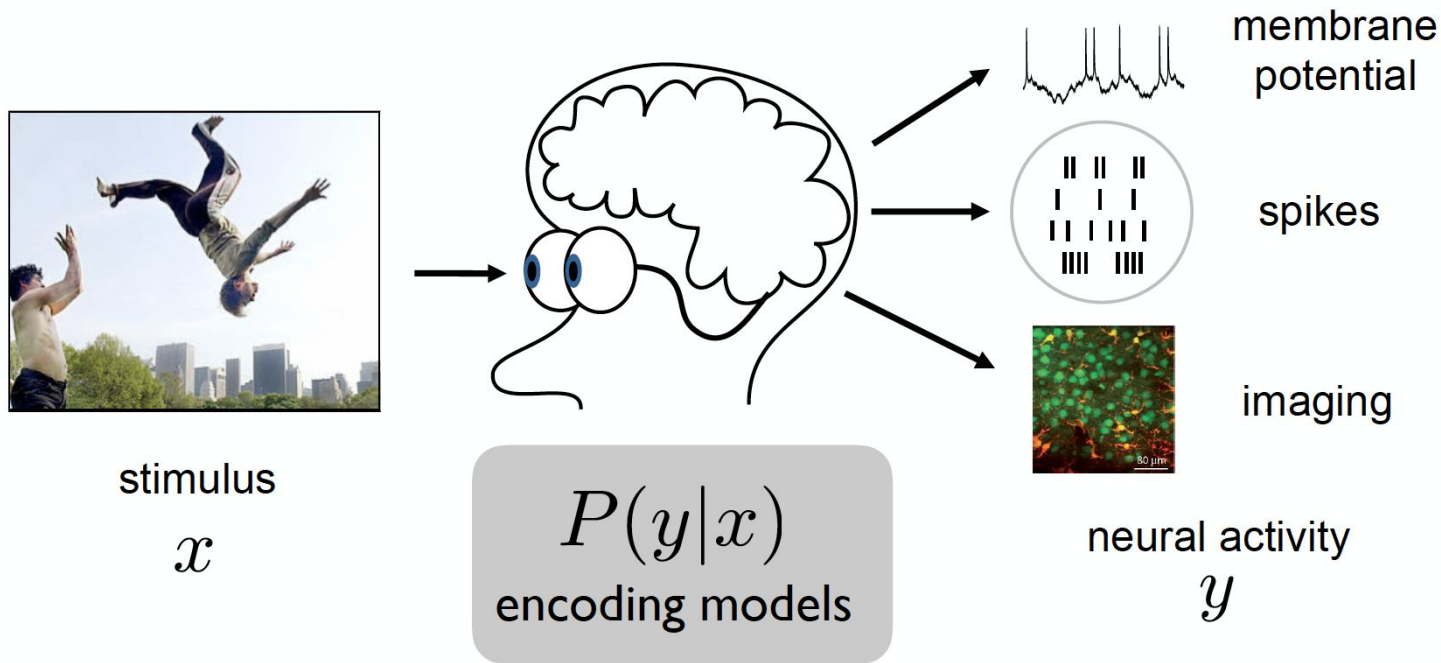5.  CS 155: Machine Learning & Data Mining, **Yisong Yue**

membrane potential

spikes

imaging

stimulus

$x$
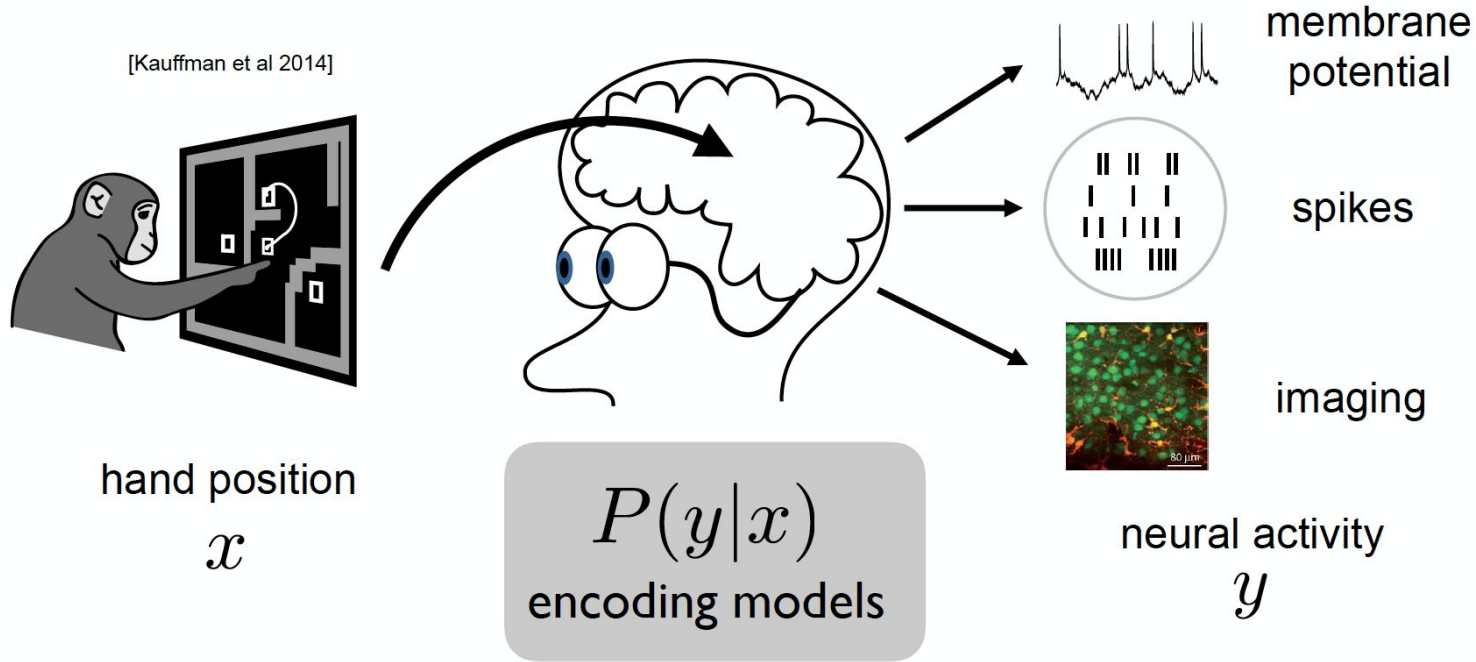
neural activity

$y$

- How are stimuli and actions encoded in neural activity?
- What aspects of neural activity carry information?

stimulus

$x$

$P(y|x)$

encoding models

membrane potential

spikes

imaging

neural activity

$y$

*Approach:*
- develop flexible statistical models of P(y|x)
- quantify information carried in neural responses

[Kauffman et al 2014]

membrane potential

spikes

imaging

hand position
$x$

$P(y|x)$
encoding models

neural activity
$y$

**"regression models"**

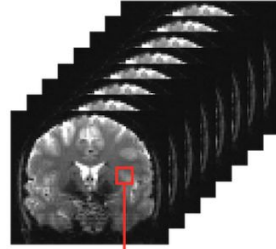- not restricted to sensory variables

[Hardcastle et al 2015]

Position (P)

Head direction

Speed (S)

Theta phase

$x$

"external covariates"

$P(y|x)$
encoding models

membrane potential

spikes

imaging

neural activity
$y$

**"regression models"**

• not restricted to sensory variables

latent variable

(unobserved or "hidden")

$P(y|x)\,P(x)$

latent encoding models

membrane potential

spikes

imaging

neural activity
$y$
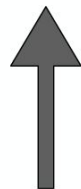
- capture hidden structure underlying neural activity

  (eg. low-dimensional or discrete states)

latent dynamics

latent variable
(unobserved or "hidden")

$x$

membrane potential

spikes

imaging

$$P(y_t|x_t)\,P(x_t|x_{t-1})$$

latent dynamical encoding models

neural activity
$y$

• capture hidden dynamics underlying neural activity

normative theories
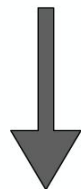(e.g. "efficient coding")

*Why* does the code
take this form?

**descriptive
statistical models**

$$P(y|x)$$

*What* is the code?

anatomy,
biophysics

*How* is it implemented?

# Topics

1. Introduction
2. General regression framework
   a. data, model, cost function, fitting procedures
3. Linear models
4. Maximum likelihood and least squares
5. Bayesian linear regression
6. Regularization
7. Bias-variance trade-off

# General framework

1. Data $x_i \in \mathbb{R}^d$     i = 1,2,3,4…..N

# General framework

Supervised setting:

1. Data $x_i \in \mathbb{R}^d$     i = 1,2,3,4…..N

2. Data transformation $\phi_j(\mathbf{x})$ , here we will use $\phi(\mathbf{x}) = \mathbf{x}$ for simplicity.

# General framework

Supervised setting:

1.  Data $x_i \in \mathbb{R}^d$    i = 1,2,3,4…..N

2.  Data transformation $\phi_j(\mathbf{x})$ , here we will use $\phi(\mathbf{x}) = \mathbf{x}$ for simplicity.

3.  Targets/Labels $y_i \in \mathbb{R}$   i = 1,2,3,4…..N

# General framework

Supervised setting:

1. Data $x_i \in \mathbb{R}^d$    i = 1,2,3,4…..N

2. Data transformation $\phi_j(\mathbf{x})$ , here we will use $\phi(\mathbf{x}) = \mathbf{x}$ for simplicity.

3. Targets/Labels $y_i \in \mathbb{R}$   i = 1,2,3,4…..N

4. Goal is to select $f : \mathbb{R}^d \to \mathbb{R}$   such that we can predict y from new x

# General framework

Supervised setting:

1. Data $x_i \in \mathbb{R}^d$    i = 1,2,3,4…..N

2. Data transformation $\phi_j(\mathbf{x})$ , here we will use $\phi(\mathbf{x}) = \mathbf{x}$ for simplicity.

3. Targets/Labels $y_i \in \mathbb{R}$   i = 1,2,3,4…..N

4. Goal is to select $f : \mathbb{R}^d \to \mathbb{R}$ such that we can predict y from new x

5. Model class $f(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^{\mathrm{T}} \phi(\mathbf{x})$

# General framework

Supervised setting:

1. Data $x_i \in \mathbb{R}^d$    i = 1,2,3,4…..N

2. Data transformation $\phi_j(\mathbf{x})$ , here we will use $\phi(\mathbf{x}) = \mathbf{x}$ for simplicity.

3. Targets/Labels $y_i \in \mathbb{R}$   i = 1,2,3,4…..N

4. Goal is to select $f : \mathbb{R}^d \rightarrow \mathbb{R}$   such that we can predict y from new x

5. Model class $f(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^{\mathrm{T}} \phi(\mathbf{x})$

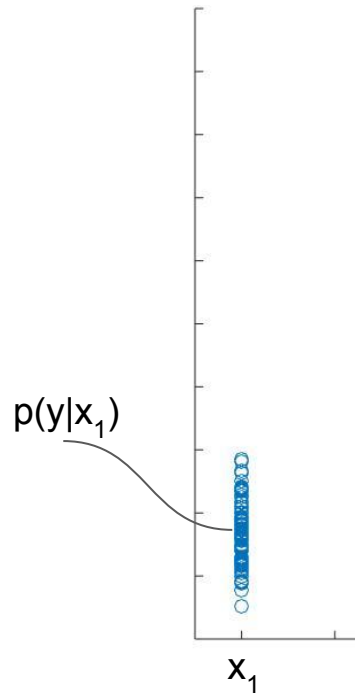6. Model parameters $w_j$

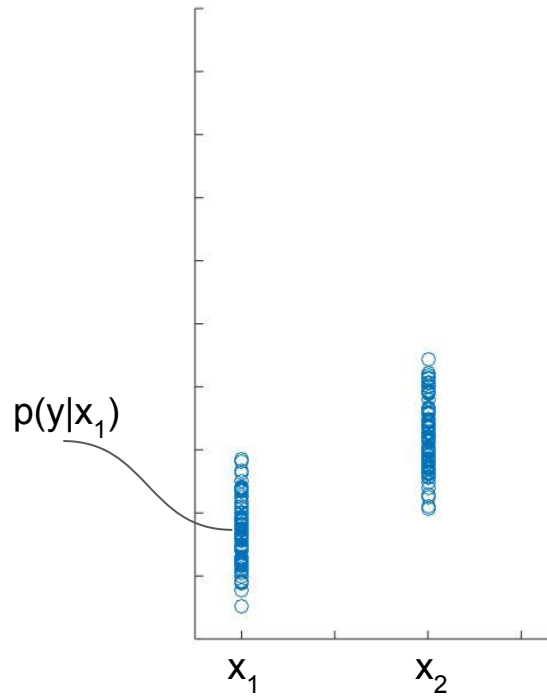# General framework

Supervised setting:

1.  Data $x_i \in \mathbb{R}^d$     i = 1,2,3,4.....N

2.  Data transformation $\phi_j(\mathbf{x})$ , here we will use $\phi(\mathbf{x}) = \mathbf{x}$ for simplicity.

3.  Targets/Labels $y_i \in \mathbb{R}$   i = 1,2,3,4.....N

4.  Goal is to select $f : \mathbb{R}^d \to \mathbb{R}$   such that we can predict y from new x

5.  Model class   $f(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^{\mathrm{T}} \phi(\mathbf{x})$

6.  Model parameters $w_j$

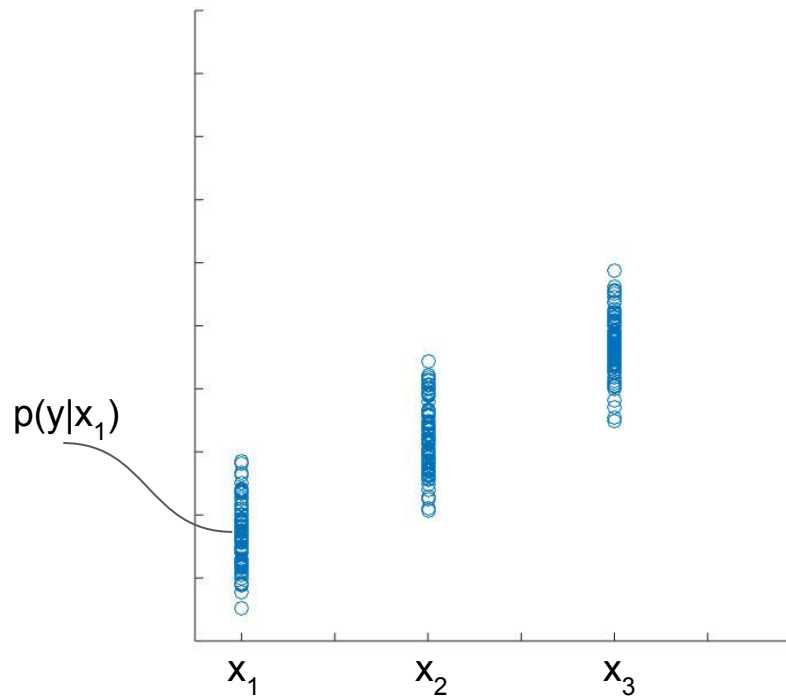7.  What's being optimized   $E_D(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^{N} \{y_i - \mathbf{w}^{\mathrm{T}} \phi(x_i)\}^2$
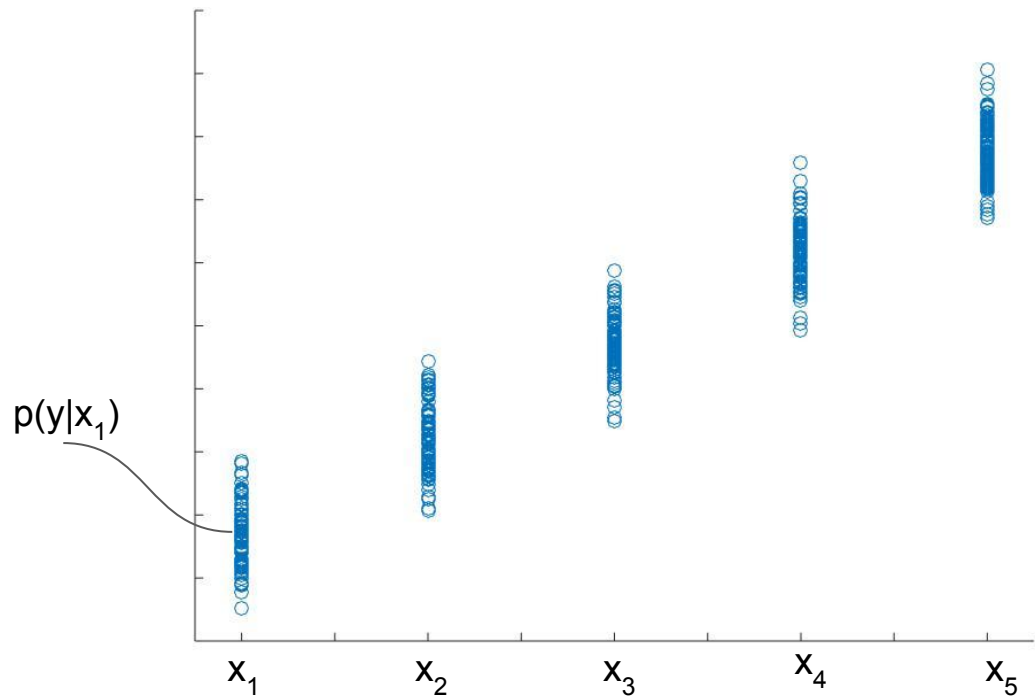
# General framework

Supervised setting:

1. Data $x_i \in \mathbb{R}^d$    i = 1,2,3,4…..N

2. Data transformation $\phi_j(\mathbf{x})$ , here we will use $\phi(\mathbf{x}) = \mathbf{x}$ for simplicity.

3. Targets/Labels $y_i \in \mathbb{R}$   i = 1,2,3,4…..N

4. Goal is to select $f : \mathbb{R}^d \rightarrow \mathbb{R}$   such that we can predict y from new x

5. Model class   $f(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^{\mathrm{T}} \phi(\mathbf{x})$

6. Model parameters   $w_j$

7. What's being optimized   $E_D(\mathbf{w}) = \dfrac{1}{2} \sum_{i=1}^{N} \{y_i - \mathbf{w}^{\mathrm{T}} \phi(x_i)\}^2$

8. Inference method ex. Maximum Likelihood (MLE)

$p(y|x_1)$

$x_1$

$p(y|x_1)$

$x_1$  $x_2$  $x_3$

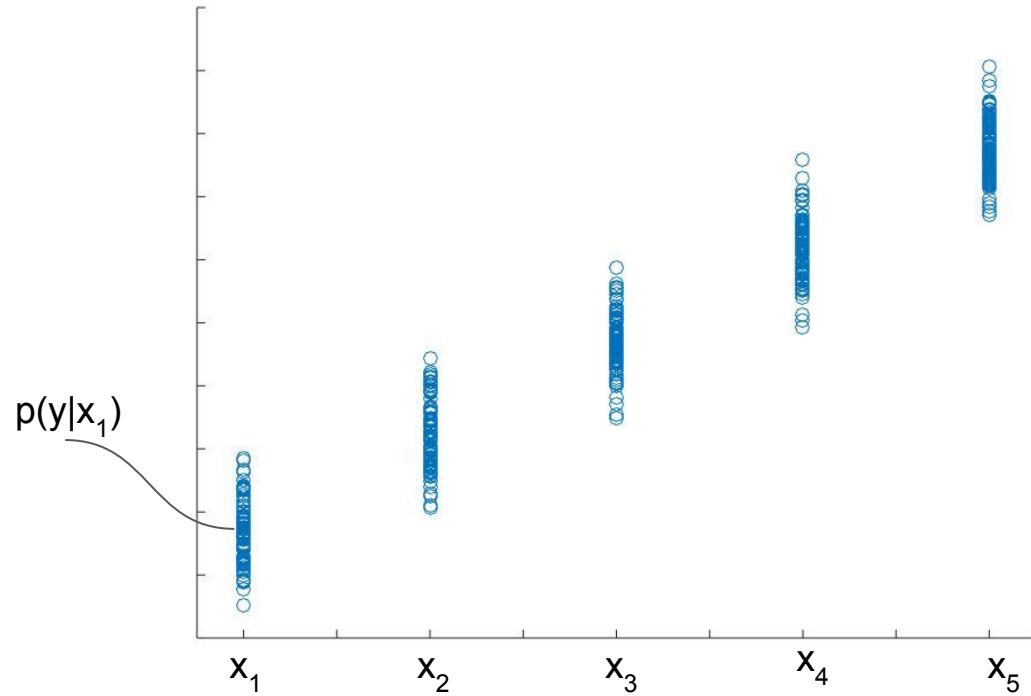$p(y|x_1)$

p(y|x_1)

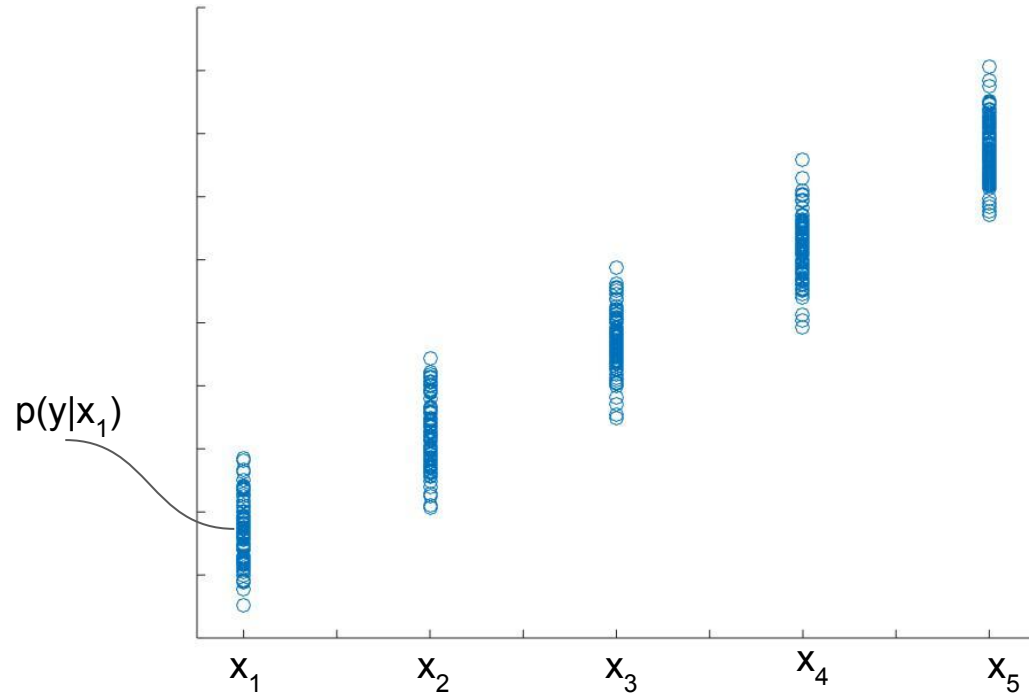$x_1$ $x_2$ $x_3$ $x_4$ $x_5$

<u>Gaussian linear regression:</u>
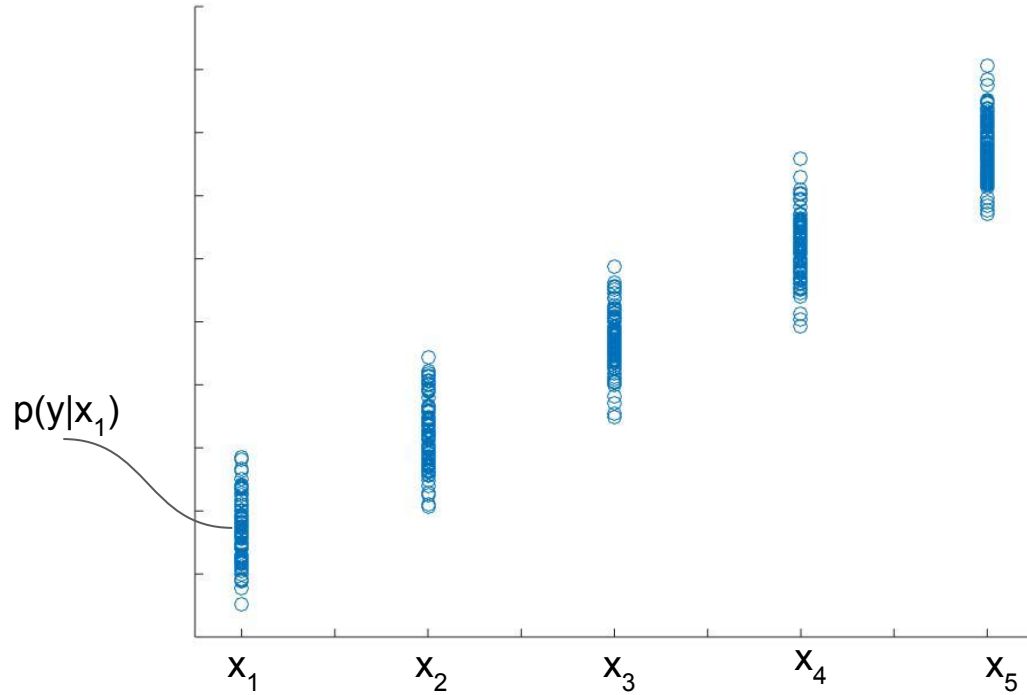
$$p_\theta(y|x) = N(y|\mu(x), \sigma^2(x)),$$

Gaussian linear regression:

$$p_\theta(y|x) = N(y|\mu(x), \sigma^2(x)),$$

$$\theta = (w, \sigma^2) \text{ with } w \in \mathbb{R}^d \ \sigma^2 > 0$$

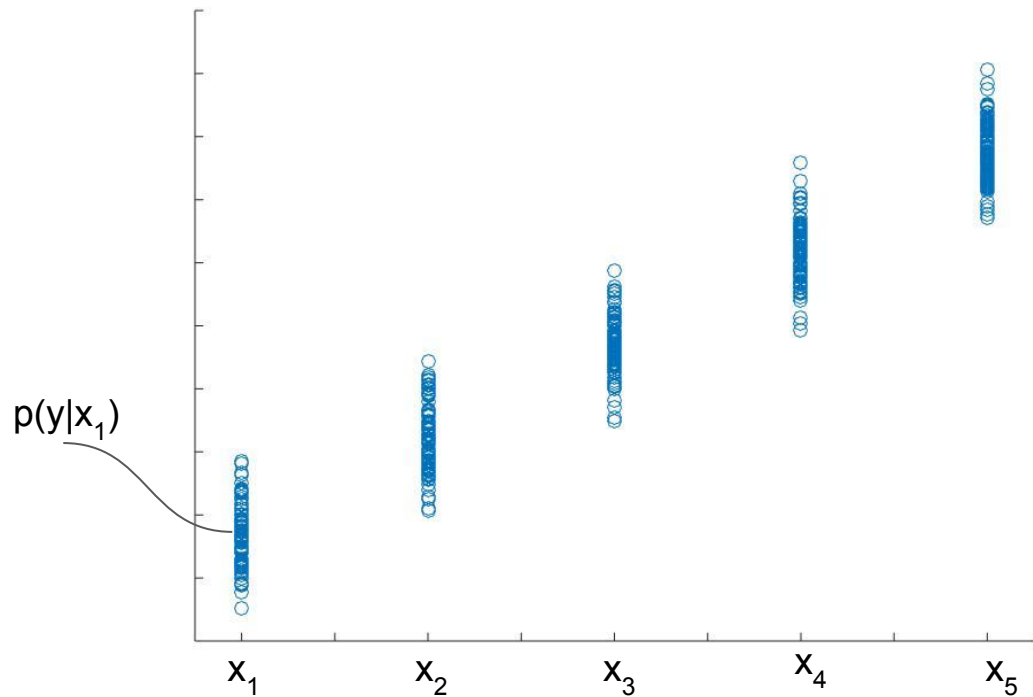p(y|x$_1$)

p(y|x$_1$)

x$_1$   x$_2$   x$_3$   x$_4$   x$_5$

Gaussian linear regression:

$$p_\theta(y|x) = N(y|\mu(x), \sigma^2(x)),$$

$$\theta = (w, \sigma^2) \text{ with } w \in \mathbb{R}^d \ \sigma^2 > 0$$

$$\mu(x) = w^T x, \sigma^2(x) = \sigma^2$$

p(y|x₁)

Gaussian linear regression:

$$p_\theta(y|x) = N(y|\mu(x), \sigma^2(x)),$$

$$\theta = (w, \sigma^2) \text{ with } w \in \mathbb{R}^d \ \sigma^2 > 0$$

$$\mu(x) = w^T x, \sigma^2(x) = \sigma^2$$

$$p_\theta(y|x) = N(y|w^T x, \sigma^2)$$

# MAXIMUM LIKELIHOOD ESTIMATION
# FOR GAUSSIAN LINEAR REGRESSION

**(1) DATA**

$$D = \left( (x_1, y_1), (x_2, y_2) \ldots (x_n, y_n) \right)$$

$$x_i \in \mathbb{R}^d, \quad y_i \in \mathbb{R}$$

**(2) MODEL**

$$y \sim N\left(w^T x, \sigma^2\right) \qquad \text{ASSUME } \sigma^2 \text{ KNOWN.}$$

$(3)$ OBJECTIVE (LIKELIHOOD)

$$\theta \in \Theta, \qquad \theta_{MLE} \in \underset{\theta \in \Theta}{ARGMAX} \; P(D \mid \theta)$$

$$P(D \mid \theta) = P(y_1, y_2 \dots y_n \mid x_1, x_2, x_3 \dots x_n, \theta)$$

$$= \prod_{i=1}^{n} P(y_i \mid x_i, \theta)$$

$$= \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} \left(y_i - w^T x_i\right)^2\right)$$

$$= \left(\frac{1}{2\pi\sigma^2}\right)^n \exp\left(\frac{-1}{2\sigma^2} \sum_{i=1}^{n} \left(y_i - w^T x_i\right)^2\right)$$

$$\begin{pmatrix} y_1 - w^T x_1 \\ \vdots \\ y_n - w^T x_n \end{pmatrix} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} - \begin{pmatrix} x_1^T w \\ \vdots \\ x_n^T w \end{pmatrix} = \vec{y} - \begin{pmatrix} x_1^T \\ \vdots \\ x_n^T \end{pmatrix} \vec{w} = \vec{y} - A\vec{w}$$

$$\boxed{w^T x_i = x_i^T w}$$

$$A = \begin{pmatrix} \text{---} \; x_1^T \; \text{---} \\ \vdots \\ \text{---} \; x_n^T \; \text{---} \end{pmatrix}$$

"DESIGN MATRIX"

$$\sum_{i=1}^{n} (y_i - w^T x_i)^2 = (y - Aw)^T (y - Aw)$$

$$= \|y - Aw\|^2 \quad \text{EUCLIDIAN NORM.}$$

$$P(D|\theta) = \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) \exp\left(-\frac{1}{2\sigma^2} (y - Aw)^T (y - Aw)\right)$$

$$\hookrightarrow \text{MAXIMIZING } P(D|\theta) = \text{MINIMIZE } (y - Aw)^T (y - Aw)$$

(4) OPTIMIZATION

$$\mathcal{L} = (y - Aw)^T (y - Aw) = \underbrace{y^T y}_{\substack{\text{NO DEPENDENCE} \\ \text{ON } w}} - \underbrace{2 y^T A w}_{= 2w^T A^T y} + \underbrace{w^T A^T A w}_{\substack{\text{QUADRATIC} \\ \text{FORM.}}}$$

$$\nabla_w \mathcal{L} = 0 - 2 A^T y + 2 A^T A w$$

$$A^T A w = A^T y$$

$$w^* = (A^T A)^{-1} A^T y.$$

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

"MOORE PENROSE INVERSE".

why is this true?

$e^x$ is order preserving

→ $\omega^*$ IS CRITICAL POINT, IS IT A MINIMUM?

COMPUTE HESSIAN, $H = \nabla^2 \mathscr{L}$

$$H_{i,j} = \left( \frac{\partial}{\partial x_i} \frac{\partial}{\partial x_j} \mathscr{L} \right)_{i,j}$$

$$H = \nabla (\nabla \mathscr{L}) = \nabla \left( -A^T y + A^T A \omega \right)$$

$$= \underline{A^T A}, \quad \omega^* \text{ IS MINIMUM IF}$$

$$A^T A \text{ IS POSITIVE SEMI-DEFINITE.}$$

$$\underline{\omega_{MLE} = (A^T A)^{-1} A^T y.}$$

If the Hessian at a given point has all positive eigenvalues, it is a positive-semidefinite. This Suggests that the underlying function is convex.

we can also get an MLE estimate of $\sigma_{MLE}$ using a similar approach

(5) EXTENSION TO CASE WHEN $\phi(x)$ IS NOT IDENTITY

$$f(x) = w^T \phi(x) \qquad \phi: \mathbb{R}^d \to \mathbb{R}^m, \text{ WHERE } m = \text{\# OF BASIS FUNCTIONS.}$$

DESIGN MATRIX FORMERLY (i.e FOR $\phi(x) = x$)

$$A = \begin{pmatrix} \text{---} & x_1^T & \text{---} \\ & \vdots & \\ \text{---} & x_n^T & \text{---} \end{pmatrix}$$

THEN FOR ARBITRARY $\phi(x)$

$$\Phi = \begin{pmatrix} \text{---} & \phi(x_1)^T & \text{---} \\ & \vdots & \\ \text{---} & \phi(x_n)^T & \text{---} \end{pmatrix}$$

$$w_{MLE} = \left(\Phi^T \Phi\right)^{-1} \Phi y = \Phi^+ y$$

# Nonlinearity via basis functions $\phi(\mathbf{x})$



Polynomials          Gaussians          Sigmoids

# Nonlinearity via basis functions $\phi(\mathbf{x})$



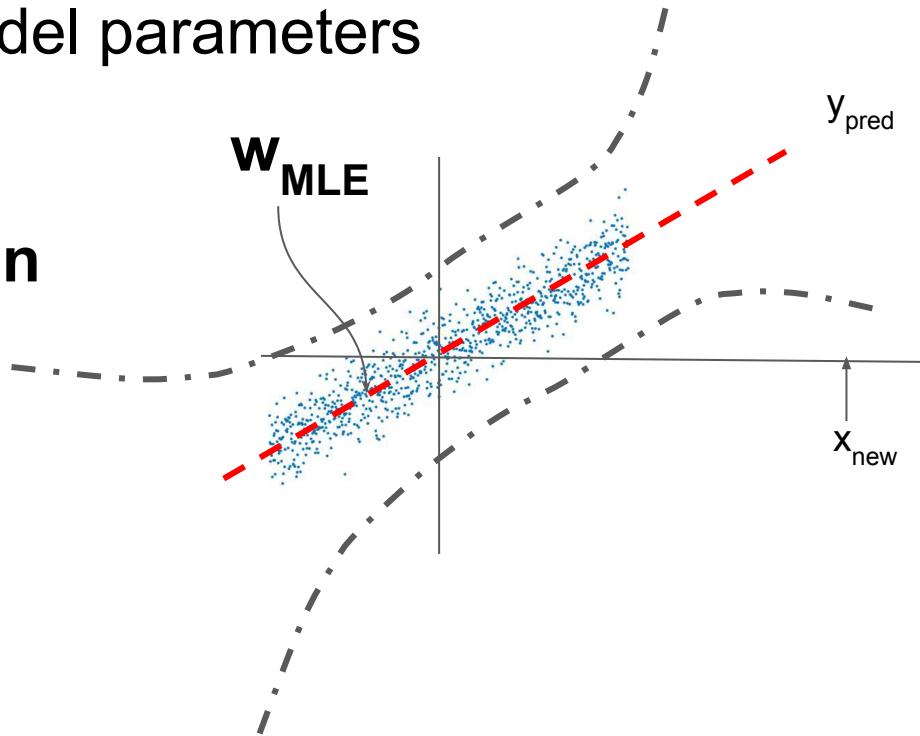scripts will be uploaded
on class website

# Nonlinearity via basis functions $\phi(\mathbf{x})$



scripts will be uploaded
on class website

So what's the problem with $\mathbf{w}_{MLE}$?

1. Easily leads to overfitting
2. No measure of uncertainty
3. Add conjugate prior on model parameters and compute posterior
   a. Bayesian approach
      **a type of regularization**

# COMPUTING THE POSTERIOR FOR GAUSSIAN LINEAR REGRESSION

RECALL THAT THE LIKELIHOOD:

$$P(D \mid w) \propto \exp\left(-\frac{a}{2}(y - Aw)^T(y - Aw)\right)$$

where $a = \frac{1}{\sigma^2}$

$$A = \begin{pmatrix} \text{---} & x_1^T & \text{---} \\ & \vdots & \\ \text{---} & x_N^T & \text{---} \end{pmatrix}$$

"DESIGN MATRIX"

(1) POSTERIOR

$$P(\omega|D) \propto P(D|\omega) P(\omega)$$

$$\omega \sim \mathcal{N}(0, b^{-1}I)$$
MULTIVARIATE GAUSSIAN.

$$P(\omega|D) \propto \exp\left(\frac{-a}{2}(y-A\omega)^T(y-A\omega)\right) \cdot \exp\left(-\frac{b}{2}\omega^T\omega\right)$$

$$= \exp\left(\frac{-a}{2}(y-A\omega)^T(y-A\omega) - \frac{b}{2}\omega^T\omega\right)$$

RIDGE REGRESSION
TYPE REGULARI-
ZATION.

NOTICE THAT:

1. $P(\omega|D)$ IS QUADRATIC IN "$\omega$"

2. ALSO A GAUSSIAN.

$$P(\omega|D) = \mathcal{N}\left(\omega \mid \mathcal{N}, \Lambda^{-1}\right)$$

$$\begin{cases} \Lambda = a A^T A + bI \\ \mathcal{N} = a \Lambda^{-1} A^T y \end{cases}$$

Multivariate gaussian

$$f_{\mathbf{X}}(x_1,\ldots,x_k) = \frac{\exp\left(-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^{\mathrm{T}}\boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})\right)}{\sqrt{(2\pi)^k|\boldsymbol{\Sigma}|}}$$

$\boldsymbol{\mu}$ = 0
$\boldsymbol{\Sigma}$ = diagonal with
entries 1/b

A direct consequence
of having the posterior is
that we can easily get the
$\mathbf{w}_{\mathbf{MAP}}$ estimate

$$\mathbf{w}_{\mathbf{MAP}} = (\mathbf{A}^{\mathrm{T}}\mathbf{A} + \frac{b}{a}\ \mathbf{I})^{-1}\mathbf{A}^{\mathrm{T}}\mathbf{y}$$

compare with

$$\mathbf{w}_{\mathbf{MLE}} = (\mathbf{A}^{\mathrm{T}}\mathbf{A})^{-1}\mathbf{A}^{\mathrm{T}}\mathbf{y}$$

## 2. PREDICTIVE DISTRIBUTION

WHAT WE REALLY WANT...

GIVEN $x_{NEW}$, $P(y \mid x, D_{\infty})$
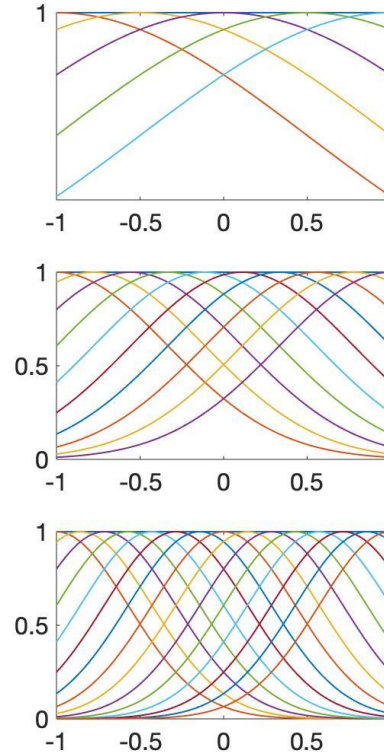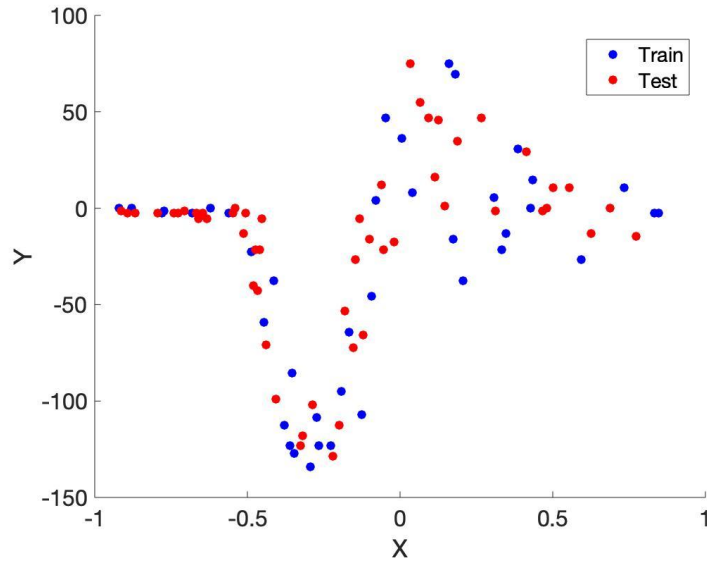
$$P(y \mid x, D) = N\left(y \mid u, \frac{1}{\lambda}\right)$$

WHERE:
$$u = N^T x$$
$$\frac{1}{\lambda} = \frac{1}{a} + x^T \Lambda^{-1} x$$

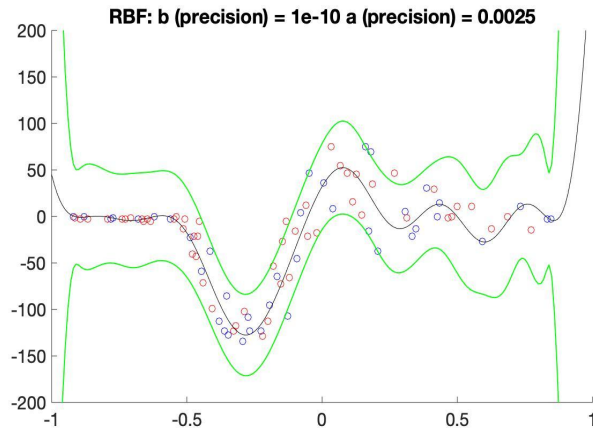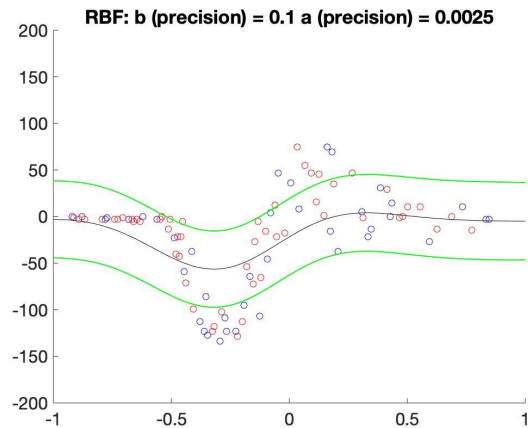A thorough and complete derivation, can be found here:

# Nonlinearity via basis functions $\phi(\mathbf{x})$



M = 30 basis functions

scripts will be uploaded
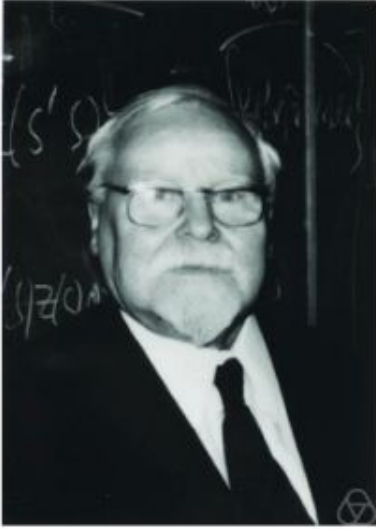on class website

# Varying the prior on the model parameters



RBF: b (precision) = 0.1 a (precision) = 0.0025

RBF: b (precision) = 1e-10 a (precision) = 0.0025

RBF: b (precision) = 0.01 a (precision) = 0.0025

For very small precision (i.e. prior is infinitely broad), we converge on $\mathbf{w}_{\mathbf{MLE}}$

The precision parameter is a hyperparameter that can be learned from the data

# Different interpretations of the effect of regularization



Tikhonov, smoothing an ill-posed problem

Zaremba, model complexity minimization

Bayes: priors over parameters

# Different kinds of regularization

- ## L0 Norm
  - \# of non-zero entries

$$\|w\|_0 = \sum_d 1_{[w_d \neq 0]}$$

- ## L1 Norm
  - Sum of absolute values

$$|w| = \|w\|_1 = \sum_d |w_d|$$

- ## L2 Norm & Squared L2 Norm
  - Sum of squares
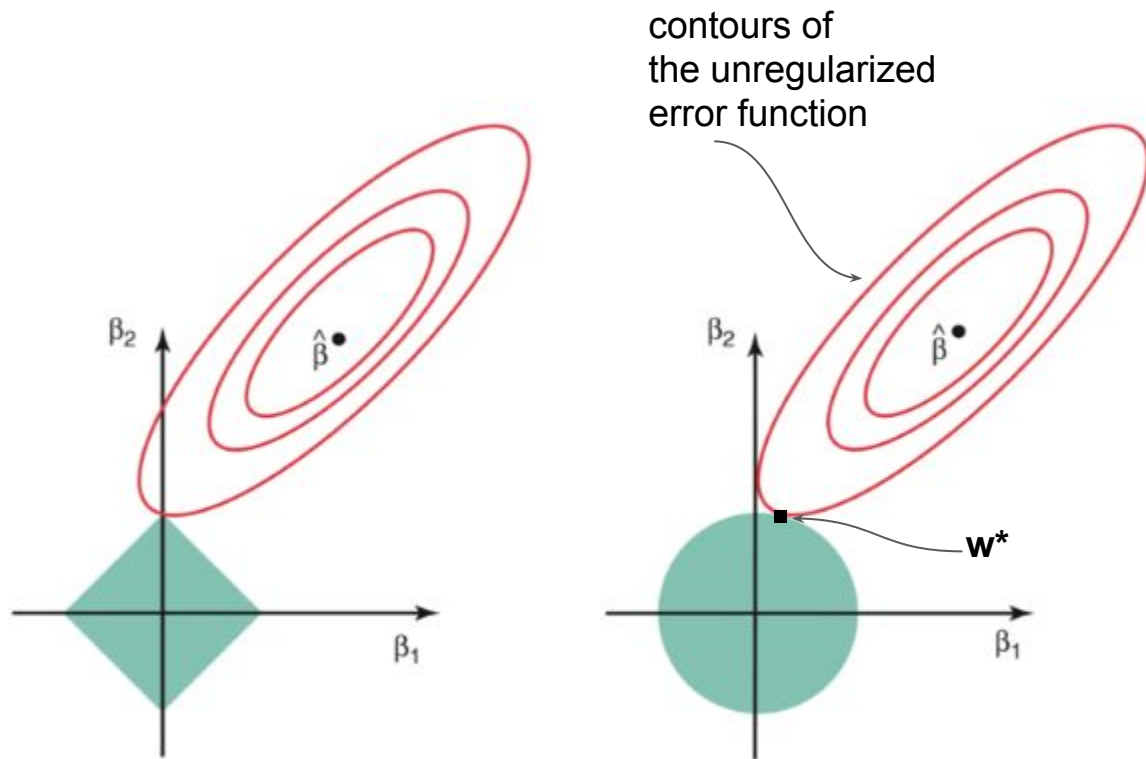  - Sqrt(sum of squares)

$$\|w\| = \sqrt{\sum_d w_d^2} \equiv \sqrt{w^T w}$$

$$\|w\|^2 = \sum_d w_d^2 \equiv w^T w$$

- ## L-infinity Norm
  - Max absolute value

$$\|w\|_\infty = \lim_{p \to \infty} \sqrt[p]{\sum_d |w_d|^p} = \max_d |w_d|$$

# Different kinds of regularization

- ## L0 Norm
  - \# of non-zero entries

$$\|w\|_0 = \sum_d 1_{[w_d \neq 0]}$$

- ## L1 Norm
  - Sum of absolute values

$$|w| = \|w\|_1 = \sum_d |w_d|$$

- ## L2 Norm & Squared L2 Norm
  - Sum of squares
  - Sqrt(sum of squares)

$$\|w\| = \sqrt{\sum_d w_d^2} \equiv \sqrt{w^T w}$$

$$\|w\|^2 = \sum_d w_d^2 \equiv w^T w$$

- ## L-infinity Norm
  - Max absolute value

$$\|w\|_\infty = \lim_{p \to \infty} \sqrt[p]{\sum_d |w_d|^p} = \max_d |w_d|$$

# A geometrical interpretation of regularization ($L_1$, $L_2$)



contours of
the unregularized
error function

$\beta_2$

$\hat{\beta}$

$\beta_1$

$\beta_2$

$\hat{\beta}$

w*

$\beta_1$

*Ridge constraint:*

$$\beta_1^2 + \beta_2^2 \;=\; 1$$

*Lasso constraint:*

$$|\beta_1| + |\beta_2| \;=\; 1$$

# Model selection

- **"True" distribution:** P(x,y)
  - Unknown to us

- **Train:** f(x) = y
  - Using training data: $S = \left\{ (x_i, y_i) \right\}_{i=1}^{N}$
  - Sampled identically and independently from P(x,y)

- **Test Error:**

$$L_P(f) = E_{(x,y)\sim P(x,y)} \left[ L(y, f(x)) \right]$$

- **Overfitting:** Test Error >> Training Error

# Model selection

- **Test Error:**

$$L_P(f) = E_{(x,y)\sim P(x,y)}\left[L(y, f(x))\right]$$

- **Treat f$_S$ as random variable:**    (randomness over $S$)

$$f_S = \underset{w,b}{\mathrm{argmin}} \sum_{(x_i,y_i)\in S} L\left(y_i, f(x_i \mid w,b)\right)$$

- **Expected Test Error:**

$$E_S\left[L_P(f_S)\right] = E_S\left[E_{(x,y)\sim P(x,y)}\left[L(y, f_S(x))\right]\right]$$

# $f_S(x)$ Linear

# $f_S(x)$ Quadratic

# f$_S$(x) Cubic

# Bias-variance tradeoff (for squared loss)

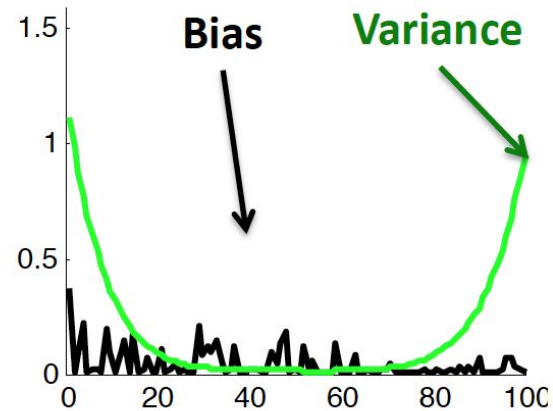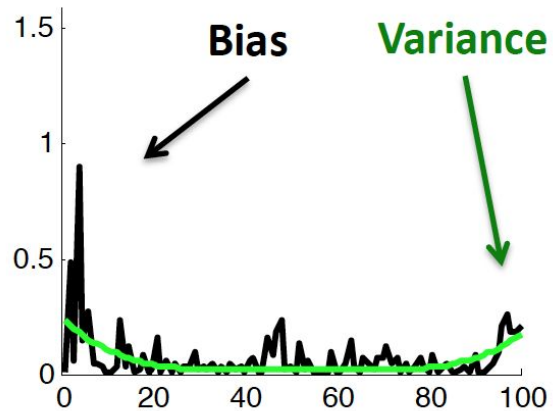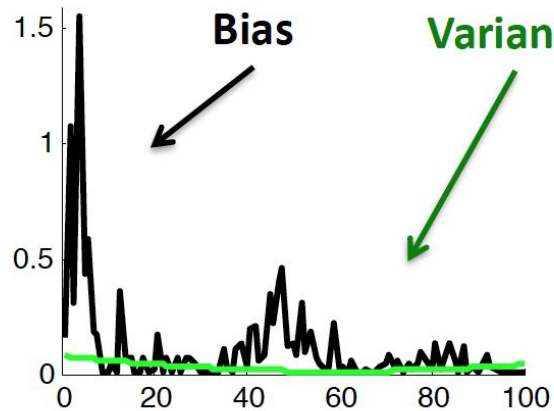$$E_S\left[L_P(f_S)\right] = E_S\left[E_{(x,y)\sim P(x,y)}\left[L(y, f_S(x))\right]\right]$$

- ## For squared error:

$$E_S\left[L_P(f_S)\right] = E_{(x,y)\sim P(x,y)}\left[E_S\left[\left(f_S(x) - F(x)\right)^2\right] + \left(F(x) - y\right)^2\right]$$

$$\underbrace{\qquad\qquad\qquad}_{\text{Variance Term}} \qquad \underbrace{\qquad\qquad}_{\text{Bias Term}}$$

$$F(x) = E_S\left[f_S(x)\right]$$

↑

"Average prediction"

**Bias**   **Variance**   **Bias**   **Variance**   **Bias**   **Variance**

# The end