Machine Learning and Statistics in Genetics and Genomics IV: Regularized and Bayesian regression

Christoph Lippert

Microsoft Research eScience group Research

Los Angeles, USA

Current topics in computational biology UCLA Winter quarter 2014 Linear Regression II

Bayesian linear regression

Model comparison and hypothesis testing

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

Summary

Outline

Outline

Linear Regression II

Bayesian linear regression

Model comparison and hypothesis testing

Summary



Linear regression:

- Making predictions
- Comparison of alternative models

Bayesian and regularized regression:

- Uncertainty in model parameters
- Generalized basis functions



▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

Linear regression:

- Making predictions
- Comparison of alternative models

Bayesian and regularized regression:

- Uncertainty in model parameters
- Generalized basis functions



Linear regression:

- Making predictions
- Comparison of alternative models

Bayesian and regularized regression:

- Uncertainty in model parameters
- Generalized basis functions



Further reading, useful material

Christopher M. Bishop: Pattern Recognition and Machine learning

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

Sam Roweis: Gaussian identities

Outline

Noise model and likelihood

Given a dataset D = {x_n, y_n}^N_{n=1}, where x_n = {x_{n,1},..., x_{n,D}} is D dimensional (for example D SNPs), fit parameters θ of a regressor f with added Gaussian noise:

$$y_n = f(\boldsymbol{x}_n; \boldsymbol{\theta}) + \epsilon_n \quad ext{where} \quad p(\epsilon \mid \sigma^2) = \mathcal{N}\left(\, \epsilon \mid 0, \sigma^2 \,
ight).$$

Equivalent likelihood formulation:

$$p(\boldsymbol{y} \mid \boldsymbol{X}) = \prod_{n=1}^{N} \mathcal{N} \left(y_n \mid f(\boldsymbol{x}_n; \boldsymbol{\theta}), \sigma^2 \right)$$

► Choose *f* to be linear:

$$p(\boldsymbol{y} \mid \boldsymbol{X}) = \prod_{n=1}^{N} \mathcal{N}(y_n \mid \boldsymbol{x}_n \cdot \boldsymbol{\beta} + c, \sigma^2)$$

Consider bias free case, c = 0, otherwise include an additional column of ones in each x_n. Choose f to be linear:

$$p(\boldsymbol{y} \mid \boldsymbol{X}) = \prod_{n=1}^{N} \mathcal{N} \left(y_n \mid \boldsymbol{x}_n \cdot \boldsymbol{\beta} + c, \sigma^2 \right)$$

Consider bias free case, c = 0, otherwise include an additional column of ones in each x_n.



Equivalent graphical model

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

Linear Regression

Taking the logarithm, we obtain

$$\ln p(\boldsymbol{y} \mid \boldsymbol{\beta}, \boldsymbol{X}, \sigma^2) = \sum_{n=1}^{N} \ln \mathcal{N} \left(y_n \mid \boldsymbol{x}_n \cdot \boldsymbol{\beta}, \sigma^2 \right)$$
$$= -\frac{N}{2} \ln 2\pi \sigma^2 - \frac{1}{2\sigma^2} \underbrace{\sum_{n=1}^{N} (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta})^2}_{\text{Sum of squares}}$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

The likelihood is maximized when the squared error is minimized.

Least squares and maximum likelihood are equivalent.

Linear Regression

Taking the logarithm, we obtain

$$\ln p(\boldsymbol{y} \mid \boldsymbol{\beta}, \boldsymbol{X}, \sigma^2) = \sum_{n=1}^{N} \ln \mathcal{N} \left(y_n \mid \boldsymbol{x}_n \cdot \boldsymbol{\beta}, \sigma^2 \right)$$
$$= -\frac{N}{2} \ln 2\pi \sigma^2 - \frac{1}{2\sigma^2} \underbrace{\sum_{n=1}^{N} (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta})^2}_{\text{Sum of squares}}$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

The likelihood is maximized when the squared error is minimized.

Least squares and maximum likelihood are equivalent.

Linear Regression

Taking the logarithm, we obtain

$$\ln p(\boldsymbol{y} \mid \boldsymbol{\beta}, \boldsymbol{X}, \sigma^2) = \sum_{n=1}^{N} \ln \mathcal{N} \left(y_n \mid \boldsymbol{x}_n \cdot \boldsymbol{\beta}, \sigma^2 \right)$$
$$= -\frac{N}{2} \ln 2\pi \sigma^2 - \frac{1}{2\sigma^2} \underbrace{\sum_{n=1}^{N} (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta})^2}_{\text{Sum of squares}}$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

The likelihood is maximized when the squared error is minimized.

• Least squares and maximum likelihood are equivalent.



(C.M. Bishop, Pattern Recognition and Machine Learning)

$$E(\boldsymbol{\beta}) = \frac{1}{2} \sum_{n=1}^{N} (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta})^2$$

イロト イ団ト イヨト イヨト 二日

• Derivative w.r.t a single weight entry β_i

$$\frac{d}{\mathrm{d}\beta_i} \ln p(\boldsymbol{y} \mid \boldsymbol{\beta}, \sigma^2) = \frac{d}{\mathrm{d}\beta_i} \left[-\frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta})^2 \right]$$
$$= \frac{1}{\sigma^2} \sum_{n=1}^N (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta}) x_i$$

Set gradient w.r.t to β to zero

$$\nabla_{\beta} \ln p(\boldsymbol{y} \mid \boldsymbol{\beta}, \sigma^2) = \frac{1}{\sigma^2} \sum_{n=1}^{N} \boldsymbol{x}_n^{\top} (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta}) = \boldsymbol{0}$$
$$\implies \boldsymbol{\beta}_{\mathrm{ML}} = \underbrace{(\boldsymbol{X}^{\top} \boldsymbol{X})^{-1} \boldsymbol{X}^{\top}}_{\mathrm{Pseudo inverse}} \boldsymbol{y}$$

• Here, the matrix $oldsymbol{X}$ is defined as $oldsymbol{X}=$

• Derivative w.r.t a single weight entry β_i

$$\frac{d}{d\beta_i} \ln p(\boldsymbol{y} \mid \boldsymbol{\beta}, \sigma^2) = \frac{d}{d\beta_i} \left[-\frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta})^2 \right]$$
$$= \frac{1}{\sigma^2} \sum_{n=1}^N (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta}) x_i$$

• Set gradient w.r.t to β to zero

$$\nabla_{\boldsymbol{\beta}} \ln p(\boldsymbol{y} \mid \boldsymbol{\beta}, \sigma^2) = \frac{1}{\sigma^2} \sum_{n=1}^{N} \boldsymbol{x}_n^{\top} (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta}) = \boldsymbol{0}$$
$$\implies \boldsymbol{\beta}_{\mathrm{ML}} = \underbrace{(\boldsymbol{X}^{\top} \boldsymbol{X})^{-1} \boldsymbol{X}^{\top}}_{\mathsf{Pseudo inverse}} \boldsymbol{y}$$

Here, the matrix X is defined as X =

• Derivative w.r.t a single weight entry β_i

$$\frac{d}{d\beta_i} \ln p(\boldsymbol{y} \mid \boldsymbol{\beta}, \sigma^2) = \frac{d}{d\beta_i} \left[-\frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta})^2 \right]$$
$$= \frac{1}{\sigma^2} \sum_{n=1}^N (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta}) x_i$$

• Set gradient w.r.t to β to zero

$$\nabla_{\beta} \ln p(\boldsymbol{y} \mid \boldsymbol{\beta}, \sigma^2) = \frac{1}{\sigma^2} \sum_{n=1}^{N} \boldsymbol{x}_n^{\top} (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta}) = \boldsymbol{0}$$

$$\Longrightarrow \boldsymbol{\beta}_{\mathrm{ML}} = \underbrace{(\boldsymbol{X}^{\top} \boldsymbol{X})^{-1} \boldsymbol{X}^{\top}}_{\mathrm{Pseudo inverse}} \boldsymbol{y}$$

$$\blacktriangleright \text{ Here, the matrix } \boldsymbol{X} \text{ is defined as } \boldsymbol{X} = \begin{bmatrix} x_{1,1} & \dots & x_{1,D} \\ \dots & \dots & \dots \\ x_{N,1} & \dots & x_{N,D} \end{bmatrix}$$

Motivation

- Non-linear relationships.
- Multiple SNPs playing a role for a particular phenotype.



(日)、

э

Univariate input \boldsymbol{x}

• Use the polynomials up to degree K to construct new features from x

$$f(x, \boldsymbol{\beta}) = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \beta_K x^K$$
$$= \sum_{k=1}^K \beta_k \phi_k(x) = \boldsymbol{\phi}(x) \boldsymbol{\beta}$$

where we defined
$$\boldsymbol{\phi}(x) = (1, x, x^2, \dots, x^K).$$

ϕ can be any feature mapping:

•
$$\phi_j(x) = e^x$$
, $\phi_j(x) = \log(x)$, ...

Radial basis functions (also: 'Gaussian' basis functions)

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

Sigmoidal basis functions

$$\phi_j(x) = \sigma\left(\frac{x-\mu_j}{s}\right), \quad \text{where } \sigma(a) = \frac{1}{1+\exp(-a)}$$



(C.M. Bishop, Pattern Recognition

and Machine Learning)

Univariate input \boldsymbol{x}

• Use the polynomials up to degree K to construct new features from x

$$f(x, \boldsymbol{\beta}) = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \beta_K x^K$$
$$= \sum_{k=1}^K \beta_k \phi_k(x) = \boldsymbol{\phi}(x) \boldsymbol{\beta}$$

where we defined $\boldsymbol{\phi}(x) = (1, x, x^2, \dots, x^K).$

ϕ can be any feature mapping:

$$\phi_j(x) = e^x, \ \phi_j(x) = \log(x), \ . \ .$$

Radial basis functions (also: 'Gaussian' basis functions)

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

Sigmoidal basis functions

$$\phi_j(x) = \sigma\left(\frac{x-\mu_j}{s}\right), \quad \text{where } \sigma(a) = \frac{1}{1+\exp(-a)}$$



(C.M. Bishop, Pattern Recognition and Machine Learning)

• • • • • • •

Univariate input \boldsymbol{x}

• Use the polynomials up to degree K to construct new features from x

$$f(x, \boldsymbol{\beta}) = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \beta_K x^K$$
$$= \sum_{k=1}^K \beta_k \phi_k(x) = \boldsymbol{\phi}(x) \boldsymbol{\beta}$$

where we defined $\boldsymbol{\phi}(x) = (1, x, x^2, \dots, x^K).$

ϕ can be any feature mapping:

•
$$\phi_j(x) = e^x$$
, $\phi_j(x) = \log(x)$, ...

Radial basis functions (also: 'Gaussian' basis functions)

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

Sigmoidal basis functions

$$\phi_j(x) = \sigma\left(\frac{x-\mu_j}{s}\right), \quad \text{where } \sigma(a) = \frac{1}{1+\exp(-a)}$$



(C.M. Bishop, Pattern Recognition and Machine Learning)

• • • • • • •

Univariate input \boldsymbol{x}

• Use the polynomials up to degree K to construct new features from x

$$f(x, \boldsymbol{\beta}) = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \beta_K x^K$$
$$= \sum_{k=1}^K \beta_k \phi_k(x) = \boldsymbol{\phi}(x) \boldsymbol{\beta}$$

where we defined $\boldsymbol{\phi}(x) = (1, x, x^2, \dots, x^K).$

ϕ can be any feature mapping:

•
$$\phi_j(x) = e^x$$
, $\phi_j(x) = \log(x)$, ...

Radial basis functions (also: 'Gaussian' basis functions)

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

Sigmoidal basis functions

$$\phi_j(x) = \sigma\left(\frac{x-\mu_j}{s}\right), \quad \text{where } \sigma(a) = \frac{1}{1+\exp(-a)}$$



(C.M. Bishop, Pattern Recognition and Machine Learning)

Univariate input \boldsymbol{x}

 \blacktriangleright Use the polynomials up to degree K to construct new features from x

$$f(x, \boldsymbol{\beta}) = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \beta_K x^K$$
$$= \sum_{k=1}^K \beta_k \phi_k(x) = \boldsymbol{\phi}(x) \boldsymbol{\beta}$$

where we defined $\boldsymbol{\phi}(x) = (1, x, x^2, \dots, x^K).$

ϕ can be any feature mapping:

•
$$\phi_j(x) = e^x$$
, $\phi_j(x) = \log(x)$, ...

Radial basis functions (also: 'Gaussian' basis functions)

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

Sigmoidal basis functions

$$\phi_j(x) = \sigma\left(\frac{x-\mu_j}{s}\right), \quad \text{where } \sigma(a) = \frac{1}{1+\exp(-a)}$$



(C.M. Bishop, Pattern Recognition

and Machine Learning)

Overfitting

- The order of the polynomial M is crucial to avoid under- and overfitting.
- Observation: Variance in regression coefficients β = [w₀^{*},..., w₉*] grows dramatically with M





Overfitting

- The order of the polynomial M is crucial to avoid under- and overfitting.
- Observation: Variance in regression coefficients β = [w₀^{*},..., w₉*] grows dramatically with M





Overfitting

- The order of the polynomial M is crucial to avoid under- and overfitting.
- ▶ Observation: Variance in regression coefficients β = [w₀^{*},..., w₉*] grows dramatically with M





Overfitting

- The order of the polynomial M is crucial to avoid under- and overfitting.
- ▶ Observation: Variance in regression coefficients β = [w₀^{*},...,w₉*] grows dramatically with M





▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Overfitting

- The order of the polynomial M is crucial to avoid under- and overfitting.
- ► Observation: Variance in regression coefficients β = [w₀^{*},..., w₉*] grows dramatically with M

	M = 0	M = 1	M = 6	M = 9
w_0^\star	0.19	0.82	0.31	0.35
w_1^\star		-1.27	7.99	232.37
w_2^{\star}			-25.43	-5321.83
w_3^{\star}			17.37	48568.31
w_4^{\star}				-231639.30
w_5^{\star}				640042.26
w_6^{\star}				-1061800.52
w_7^{\star}				1042400.18
w_8^{\star}				-557682.99
w_9^{\star}				125201.43

(C.M. Bishop, Pattern Recognition and Machine Learning)

Generalization performance

$$E(\boldsymbol{\beta}) = \frac{1}{2} \sum_{n=1}^{N} (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta})^2$$

Root-mean squared error

$$\mathcal{E}_{\text{RMS}} = \sqrt{2E(\boldsymbol{\beta})/N}$$

= $\sqrt{\left(\sum_{n=1}^{N} (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta})^2\right)/N}$

- Can be measured on training data
- Or when predicting new data (test data)
- Underfitting: large E_{RMS} on train and test data
- ▶ **Overfitting**: small *E*_{RMS} on train and large *E*_{RMS} on test data.

Generalization performance

$$E(\boldsymbol{\beta}) = \frac{1}{2} \sum_{n=1}^{N} (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta})^2$$

Root-mean squared error

$$\begin{aligned} E_{\text{RMS}} &= \sqrt{2E(\boldsymbol{\beta})/N} \\ &= \sqrt{\left(\sum_{n=1}^{N} (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta})^2\right)/N} \end{aligned}$$

- Can be measured on training data
- Or when predicting new data (test data)
- Underfitting: large E_{RMS} on train and test data
- ▶ **Overfitting**: small *E*_{RMS} on train and large *E*_{RMS} on test data.

Generalization performance

$$E(\boldsymbol{\beta}) = \frac{1}{2} \sum_{n=1}^{N} (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta})^2$$

Root-mean squared error

$$E_{\text{RMS}} = \sqrt{2E(\boldsymbol{\beta})/N}$$
$$= \sqrt{\left(\sum_{n=1}^{N} (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta})^2\right)/N}$$

- Can be measured on training data
- Or when predicting new data (test data)
- Underfitting: large E_{RMS} on train and test data
- ► Overfitting: small E_{RMS} on train and large E_{RMS} on test data.

Generalization performance

$$E(\boldsymbol{\beta}) = \frac{1}{2} \sum_{n=1}^{N} (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta})^2$$

Root-mean squared error

$$E_{\text{RMS}} = \sqrt{2E(\boldsymbol{\beta})/N}$$
$$= \sqrt{\left(\sum_{n=1}^{N} (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta})^2\right)/N}$$

- Can be measured on training data
- Or when predicting new data (test data)
- ► Underfitting: large *E*_{RMS} on train and test data
- Overfitting: small E_{RMS} on train and large E_{RMS} on test data.



(C.M. Bishop, Pattern Recognition and Machine Learning)

イロト 不得 トイヨト イヨト

-

Generalization performance

$$E(\boldsymbol{\beta}) = \frac{1}{2} \sum_{n=1}^{N} (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta})^2$$

Root-mean squared error

$$E_{\text{RMS}} = \sqrt{2E(\boldsymbol{\beta})/N}$$
$$= \sqrt{\left(\sum_{n=1}^{N} (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta})^2\right)/N}$$

- Can be measured on training data
- Or when predicting new data (test data)
- ► Underfitting: large *E*_{RMS} on train and test data
- Overfitting: small E_{RMS} on train and large E_{RMS} on test data.



(C.M. Bishop, Pattern Recognition and Machine Learning)

・ロト ・ 理 ト ・ ヨ ト ・ ヨ ト

-

Generalization performance

$$E(\boldsymbol{\beta}) = \frac{1}{2} \sum_{n=1}^{N} (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta})^2$$

Root-mean squared error

$$E_{\text{RMS}} = \sqrt{2E(\boldsymbol{\beta})/N}$$
$$= \sqrt{\left(\sum_{n=1}^{N} (y_n - \boldsymbol{x}_n \cdot \boldsymbol{\beta})^2\right)/N}$$

- Can be measured on training data
- Or when predicting new data (test data)
- ► Underfitting: large *E*_{RMS} on train and test data
- ► Overfitting: small E_{RMS} on train and large E_{RMS} on test data.



(C.M. Bishop, Pattern Recognition and Machine Learning)

・ロト ・ 理 ト ・ ヨ ト ・ ヨ ト
Overfitting

The number N of training data is crucial to accurately estimate many parameters without overfitting.





Overfitting

► The number N of training data is crucial to accurately estimate many parameters without overfitting.





Overfitting

The number N of training data is crucial to accurately estimate many parameters without overfitting.



(C.M. Bishop, Pattern Recognition and Machine Learning)

Multivariate regression

Polynomial curve fitting

$$f(x, \boldsymbol{\beta}) = \beta_0 + \beta_1 x + \dots + \beta_K x^K$$
$$= \sum_{k=1}^K \beta_k \phi_k(x)$$
$$= \boldsymbol{\phi}(x) \cdot \boldsymbol{\beta},$$

High dimensional regression

$$f(x, \boldsymbol{\beta}) = \sum_{d=1}^{D} \beta_d x_d$$
$$= \boldsymbol{x} \cdot \boldsymbol{\beta}$$

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Note: When fitting a single binary variable x_i, a linear model is most general!

Multivariate regression

Polynomial curve fitting

$$f(x, \boldsymbol{\beta}) = \beta_0 + \beta_1 x + \dots + \beta_K x^K$$
$$= \sum_{k=1}^K \beta_k \phi_k(x)$$
$$= \boldsymbol{\phi}(x) \cdot \boldsymbol{\beta},$$

High dimensional regression

$$f(x, \boldsymbol{\beta}) = \sum_{d=1}^{D} \beta_d x_d$$
$$= \boldsymbol{x} \cdot \boldsymbol{\beta}$$

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Note: When fitting a single binary variable x_i, a linear model is most general!

Ridge regression

Solutions to avoid overfitting:

- 1. Intelligently choose number of parameters
- 2. Get more data
- 3. Regularize the regression weights β
- Quadratically regularized objective function

Ridge regression

Solutions to avoid overfitting:

- 1. Intelligently choose number of parameters
- 2. Get more data
- 3. Regularize the regression weights $oldsymbol{eta}$
- Quadratically regularized objective function



Ridge regression

Solutions to avoid overfitting:

- 1. Intelligently choose number of parameters
- 2. Get more data
- 3. Regularize the regression weights $oldsymbol{eta}$
- Quadratically regularized objective function



Ridge regression

- Solutions to avoid overfitting:
 - 1. Intelligently choose number of parameters
 - 2. Get more data
 - 3. Regularize the regression weights $oldsymbol{eta}$
- Quadratically regularized objective function

$$E(\boldsymbol{\beta}) = \underbrace{\frac{1}{2} \sum_{n=1}^{N} (y_n - \boldsymbol{\phi}(\boldsymbol{x}_n) \cdot \boldsymbol{\beta})^2}_{\text{Squared error}} + \underbrace{\frac{\lambda}{2} \boldsymbol{\beta}^\top \boldsymbol{\beta}}_{\text{Regularizer}}$$

- L_2 regularization
 - M = 9, different λ values



(C.M. Bishop, Pattern Recognition and Machine Learning)

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

- L_2 regularization
 - M = 9, different λ values



(C.M. Bishop, Pattern Recognition and Machine Learning)

- L_2 regularization
 - M = 9, different λ values



(C.M. Bishop, Pattern Recognition and Machine Learning)

イロト イ団ト イヨト イヨト 二日

 L_2 regularization

▶ M = 9, different λ values

	$\ln \lambda = -\infty$	$\ln \lambda = -18$	$\ln \lambda = 0$
w_0^\star	0.35	0.35	0.13
w_1^\star	232.37	4.74	-0.05
w_2^{\star}	-5321.83	-0.77	-0.06
w_3^{\star}	48568.31	-31.97	-0.05
w_4^{\star}	-231639.30	-3.89	-0.03
w_5^{\star}	640042.26	55.28	-0.02
w_6^{\star}	-1061800.52	41.32	-0.01
w_7^{\star}	1042400.18	-45.95	-0.00
w_8^{\star}	-557682.99	-91.53	0.00
w_9^{\star}	125201.43	72.68	0.01

(C.M. Bishop, Pattern Recognition and Machine Learning)

- L_2 regularization
 - M = 9, different λ values



(C.M. Bishop, Pattern Recognition and Machine Learning)

Variance of f^{est}

$$\mathbb{E}_{\mathcal{D}}\left[f^{\mathsf{est}} - \mathbb{E}\left[f^{\mathsf{est}}\right]^2\right]$$

Bias of f^{est} $\mathbb{E}_{\mathcal{D}}\left[f^{true} - f^{est}\right]$

Experiment:

- 100 random data sets (N = 25)
- learn 25 RBF basis functions

 \blacktriangleright vary λ

Empirical Observations:

Variance of f^{est}

$$\mathbb{E}_{\mathcal{D}}\left[f^{\mathsf{est}} - \mathbb{E}\left[f^{\mathsf{est}}\right]^2\right]$$

Bias of f^{est} $\mathbb{E}_{\mathcal{D}}\left[f^{true} - f^{est}\right]$

Experiment:

- 100 random data sets (N = 25)
- learn 25 RBF basis functions
- vary λ

Empirical Observations:

Bias decreases with smaller A Variance increases with smaller

(C.M. Bishop, Pattern Recognition and Machine Learning) $\mathbb{P} \to \mathbb{P} = \mathbb{P} =$

Variance of f^{est}

$$\mathbb{E}_{\mathcal{D}}\left[f^{\mathsf{est}} - \mathbb{E}\left[f^{\mathsf{est}}\right]^2\right]$$



Experiment:

- 100 random data sets (N = 25)
- learn 25 RBF basis functions
- vary λ

Bias of f^{est}

$$\mathbb{E}_{\mathcal{D}}\left[\boldsymbol{f}^{\mathsf{true}} - \boldsymbol{f}^{\mathsf{est}}\right]$$



Empirical Observations:

(C.M. Bishop, Pattern Recognition and Machine Learning) $\mathbb{P} \to \mathbb{P} = \mathbb{P} =$

Variance of f^{est}

$$\mathbb{E}_{\mathcal{D}}\left[f^{\mathsf{est}} - \mathbb{E}\left[f^{\mathsf{est}}\right]^2\right]$$



Experiment:

- 100 random data sets (N = 25)
- learn 25 RBF basis functions
- vary λ

Bias of f^{est}

$$\mathbb{E}_{\mathcal{D}}\left[\boldsymbol{f}^{\mathsf{true}} - \boldsymbol{f}^{\mathsf{est}}\right]$$



Empirical Observations:

Bias decreases with smaller λ

() ↓ () ↓

(C.M. Bishop, Pattern Recognition and Machine Learning)

Variance of f^{est}

$$\mathbb{E}_{\mathcal{D}}\left[f^{\mathsf{est}} - \mathbb{E}\left[f^{\mathsf{est}}\right]^2\right]$$



Experiment:

- 100 random data sets (N = 25)
- learn 25 RBF basis functions
- vary λ

Bias of f^{est}

$$\mathbb{E}_{\mathcal{D}}\left[\boldsymbol{f}^{\mathsf{true}} - \boldsymbol{f}^{\mathsf{est}}\right]$$



Empirical Observations:

Bias decreases with smaller λ Variance increases with smaller

3

(C.M. Bishop, Pattern Recognition and Machine Learning)

Variance of f^{est}

$$\mathbb{E}_{\mathcal{D}}\left[f^{\mathsf{est}} - \mathbb{E}\left[f^{\mathsf{est}}\right]^2\right]$$



Experiment:

- 100 random data sets (N = 25)
- learn 25 RBF basis functions
- vary λ

Bias of f^{est}

$$\mathbb{E}_{\mathcal{D}}\left[\boldsymbol{f}^{\mathsf{true}} - \boldsymbol{f}^{\mathsf{est}}\right]$$



Empirical Observations:

- Bias decreases with smaller λ
- Variance increases with smaller λ

▲ 프 ▶ 프

(C.M. Bishop, Pattern Recognition and Machine Learning)

Variance of f^{est}

$$\mathbb{E}_{\mathcal{D}}\left[f^{\mathsf{est}} - \mathbb{E}\left[f^{\mathsf{est}}\right]^2\right]$$



Experiment:

- 100 random data sets (N = 25)
- learn 25 RBF basis functions
- vary λ

Bias of f^{est}

$$\mathbb{E}_{\mathcal{D}}\left[\boldsymbol{f}^{\mathsf{true}} - \boldsymbol{f}^{\mathsf{est}}\right]$$



Empirical Observations:

- Bias decreases with smaller λ
- Variance increases with smaller λ

< ∃ →

3

(C.M. Bishop, Pattern Recognition and Machine Learning)

Effect on mean squared error

$$y_n = f^{\mathsf{true}}(x_n) + \epsilon_n$$

mean squared error
$$(f^{\text{est}}) = \mathbb{E}_{\mathcal{D}} \left[(y - f^{\text{est}}(x))^2 \right]$$

= (bias)² + variance + noise

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

- 100 random data sets (N = 25)
- learn 25 RBF basis functions
- vary λ
- Compute sample estimates of bias and variance
- Additionally 1000 test data points to estimate mean-squared error

Effect on mean squared error

$$y_n = f^{\mathsf{true}}(x_n) + \epsilon_n$$

mean squared error
$$(f^{\text{est}}) = \mathbb{E}_{\mathcal{D}} \left[(y - f^{\text{est}}(x))^2 \right]$$

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

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

- 100 random data sets (N = 25)
- learn 25 RBF basis functions
- vary λ
- Compute sample estimates of bias and variance
- Additionally 1000 test data points to estimate mean-squared error

Effect on mean squared error

$$y_n = f^{\mathsf{true}}(x_n) + \epsilon_n$$

mean squared error
$$(f^{\text{est}}) = \mathbb{E}_{\mathcal{D}} \left[(y - f^{\text{est}}(x))^2 \right]$$

= (bias)² + variance + noise

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

- 100 random data sets (N = 25)
- learn 25 RBF basis functions
- \blacktriangleright vary λ
- Compute sample estimates of bias and variance
- Additionally 1000 test data points to estimate mean-squared error

Effect on mean squared error

$$y_n = f^{\mathsf{true}}(x_n) + \epsilon_n$$

mean squared error
$$(f^{\text{est}}) = \mathbb{E}_{\mathcal{D}} \left[(y - f^{\text{est}}(x))^2 \right]$$

= (bias)² + variance + noise

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

- 100 random data sets (N = 25)
- learn 25 RBF basis functions
- \blacktriangleright vary λ
- Compute sample estimates of bias and variance
- Additionally 1000 test data points to estimate mean-squared error

Effect on mean squared error

$$y_n = f^{\mathsf{true}}(x_n) + \epsilon_n$$

mean squared error
$$(f^{\text{est}}) = \mathbb{E}_{\mathcal{D}} \left[(y - f^{\text{est}}(x))^2 \right]$$

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

Experiment as before:

- 100 random data sets (N = 25)
- learn 25 RBF basis functions
- \blacktriangleright vary λ
- Compute sample estimates of bias and variance
- Additionally 1000 test data points to estimate mean-squared error





▲ロト ▲帰 ト ▲ ヨ ト ▲ ヨ ト ・ ヨ ・ の Q ()

More general regularizers

More general regularization:



• $q \leq 1$: non-differentiable

• q < 1: non-convex (could have local optima)

More general regularizers

More general regularization:



• $q \leq 1$: non-differentiable

▶ q < 1: non-convex (could have local optima)</p>



(C.M. Bishop, Pattern Recognition and Machine Learning)

More general regularizers

More general regularization:



• $q \leq 1$: non-differentiable

▶ q < 1: non-convex (could have local optima)</p>



(C.M. Bishop, Pattern Recognition and Machine Learning)

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

More general regularizers

More general regularization:



- $q \leq 1$: non-differentiable
- ▶ q < 1: non-convex (could have local optima)</p>



(C.M. Bishop, Pattern Recognition and Machine Learning)

More general regularizers

More general regularization:



- $q \leq 1$: non-differentiable
- q < 1: non-convex (could have local optima)



(C.M. Bishop, Pattern Recognition and Machine Learning)

Smaller q yields sparser solution β^{\star}

q = 2: Ridge regression (*L*₂)
q = 1: Lasso (*L*₁)



 β_2

► Squared error

Regularizer

(C.M. Bishop, Pattern Recognition and Machine Learning)

Smaller q yields sparser solution β^{\star}

- q = 2: Ridge regression (L_2)
- ▶ q = 1: Lasso (L_1)



Loss functions and related methods

Even more general: general loss function



- Many state-of-the-art machine learning methods can be expressed within this framework.
 - ▶ Linear Regression: squared loss, squared regularizer.
 - Support Vector Machine: hinge loss, squared regularizer.
 - ► Lasso: squared loss, L1 regularizer.
- Inference: minimize the cost function E(β), yielding a point estimate for β.

▶ Q: How to determine q and the a suitable loss function?

Loss functions and related methods

Even more general: general loss function



- Many state-of-the-art machine learning methods can be expressed within this framework.
 - Linear Regression: squared loss, squared regularizer.
 - Support Vector Machine: hinge loss, squared regularizer.
 - Lasso: squared loss, L1 regularizer.
- Inference: minimize the cost function E(β), yielding a point estimate for β.

▶ Q: How to determine q and the a suitable loss function?

Loss functions and related methods

Even more general: general loss function



- Many state-of-the-art machine learning methods can be expressed within this framework.
 - Linear Regression: squared loss, squared regularizer.
 - Support Vector Machine: hinge loss, squared regularizer.
 - Lasso: squared loss, L1 regularizer.
- Inference: minimize the cost function E(β), yielding a point estimate for β.

<ロト 4 回 ト 4 回 ト 4 回 ト 回 の Q (O)</p>

▶ Q: How to determine q and the a suitable loss function?
Even more general: general loss function



- Many state-of-the-art machine learning methods can be expressed within this framework.
 - Linear Regression: squared loss, squared regularizer.
 - Support Vector Machine: hinge loss, squared regularizer.
 - Lasso: squared loss, L1 regularizer.
- Inference: minimize the cost function E(β), yielding a point estimate for β.
- ▶ Q: How to determine q and the a suitable loss function?

Cross validation: minimization of expected loss

Compare candidate models \mathcal{H} on generalization performance (different λ , different regularizers, different basis functions, etc.)

- Randomly split data into K sets of equal size
- For each fold k:
 - 1. Train on all data except the k^{th} sec 2. Test evaluation on k^{th} set
- Assess average loss on test sets $\frac{1}{K} \sum_{k=1}^{K} E_k^{\text{test}}(\mathcal{H})$
- ▶ Pick model *H* with lowest average loss
- Re-train optimal \mathcal{H} on all data



Cross validation: minimization of expected loss

Compare candidate models \mathcal{H} on generalization performance (different λ , different regularizers, different basis functions, etc.)

- Randomly split data into K sets of equal size
- For each fold k:
 - 1. Train on all data except the k^{th} set 2. Test evaluation on k^{th} set
- Assess average loss on test sets $\frac{1}{K} \sum_{k=1}^{K} E_k^{\text{test}}(\mathcal{H})$
- ▶ Pick model *H* with lowest average loss
- Re-train optimal \mathcal{H} on all data



Cross validation: minimization of expected loss

Compare candidate models \mathcal{H} on generalization performance (different λ , different regularizers, different basis functions, etc.)

- Randomly split data into K sets of equal size
- ► For each fold k:
 - 1. Train on all data except the k^{th} set 2. Test evaluation on k^{th} set
- Assess average loss on test sets $\frac{1}{K} \sum_{k=1}^{K} E_k^{\text{test}}(\mathcal{H})$
- ▶ Pick model *H* with lowest average loss
- Re-train optimal \mathcal{H} on all data



Cross validation: minimization of expected loss

Compare candidate models \mathcal{H} on generalization performance (different λ , different regularizers, different basis functions, etc.)

- Randomly split data into K sets of equal size
- For each fold k:
 - 1. Train on all data except the k^{th} set
 - 2. Test evaluation on k^{th} set



- ▶ Pick model *H* with lowest average loss
- Re-train optimal \mathcal{H} on all data



Cross validation: minimization of expected loss

Compare candidate models \mathcal{H} on generalization performance (different λ , different regularizers, different basis functions, etc.)

- Randomly split data into K sets of equal size
- For each fold k:
 - 1. Train on all data except the k^{th} set
 - 2. Test evaluation on k^{th} set
- Assess average loss on test sets $\frac{1}{K} \sum_{k=1}^{K} E_k^{\text{test}}(\mathcal{H})$
- ▶ Pick model *H* with lowest average loss
- Re-train optimal H on all data



Cross validation: minimization of expected loss

Compare candidate models \mathcal{H} on generalization performance (different λ , different regularizers, different basis functions, etc.)

- Randomly split data into K sets of equal size
- For each fold k:
 - 1. Train on all data except the k^{th} set
 - 2. Test evaluation on k^{th} set
- Assess average loss on test sets $\frac{1}{K} \sum_{k=1}^{K} E_k^{\text{test}}(\mathcal{H})$
- ▶ Pick model *H* with lowest average loss
- Re-train optimal H on all data



Cross validation: minimization of expected loss

Compare candidate models \mathcal{H} on generalization performance (different λ , different regularizers, different basis functions, etc.)

- Randomly split data into K sets of equal size
- For each fold k:
 - 1. Train on all data except the k^{th} set
 - 2. Test evaluation on k^{th} set
- Assess average loss on test sets $\frac{1}{K} \sum_{k=1}^{K} E_k^{\text{test}}(\mathcal{H})$
- Pick model \mathcal{H} with lowest average loss

Re-train optimal H on all data



Cross validation: minimization of expected loss

Compare candidate models \mathcal{H} on generalization performance (different λ , different regularizers, different basis functions, etc.)

- Randomly split data into K sets of equal size
- For each fold k:
 - 1. Train on all data except the k^{th} set
 - 2. Test evaluation on k^{th} set
- Assess average loss on test sets $\frac{1}{K} \sum_{k=1}^{K} E_k^{\text{test}}(\mathcal{H})$
- Pick model H with lowest average loss
- Re-train optimal H on all data



- So far: minimization of error functions.
- Back to probabilities?

$$E(\boldsymbol{\beta}) = \underbrace{\frac{1}{2} \sum_{n=1}^{N} (y_n - \boldsymbol{\phi}(\boldsymbol{x}_n) \cdot \boldsymbol{\beta})^2}_{\mathbb{R} \text{eg}} + \underbrace{\frac{\lambda}{2}}_{\mathbb{R} \text{eg}}$$





< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

- Regularized regression equivalent to MAP estimation
- Most alternative choices of regularizers and loss functions can be mapped to an equivalent probabilistic representation in a similar way.

- So far: minimization of error functions.
- Back to probabilities?

$$E(\boldsymbol{\beta}) = \underbrace{\frac{1}{2} \sum_{n=1}^{N} (y_n - \boldsymbol{\phi}(\boldsymbol{x}_n) \cdot \boldsymbol{\beta})^2}_{\text{Squared error}} + \underbrace{\frac{\lambda}{2} \boldsymbol{\beta}^{\top} \boldsymbol{\beta}}_{\text{Regularizer}}$$
$$= const. - \sum_{n=1}^{N} \ln \mathcal{N} \left(y_n \mid \boldsymbol{\phi}(\boldsymbol{x}_n) \cdot \boldsymbol{\beta}, \sigma^2 \right) - \ln \mathcal{N} \left(\boldsymbol{\beta} \mid \boldsymbol{0}, \frac{1}{\lambda} \boldsymbol{I} \right)$$

- Regularized regression equivalent to MAP estimation
- Most alternative choices of regularizers and loss functions can be mapped to an equivalent probabilistic representation in a similar way.

- So far: minimization of error functions.
- Back to probabilities?

$$E(\boldsymbol{\beta}) = \underbrace{\frac{1}{2} \sum_{n=1}^{N} (y_n - \boldsymbol{\phi}(\boldsymbol{x}_n) \cdot \boldsymbol{\beta})^2}_{\text{Squared error}} + \underbrace{\frac{\lambda}{2} \boldsymbol{\beta}^{\top} \boldsymbol{\beta}}_{\text{Regularizer}}$$
$$= const. - \sum_{n=1}^{N} \ln \mathcal{N} \left(y_n \mid \boldsymbol{\phi}(\boldsymbol{x}_n) \cdot \boldsymbol{\beta}, \sigma^2 \right) - \ln \mathcal{N} \left(\boldsymbol{\beta} \mid \boldsymbol{0}, \frac{1}{\lambda} \boldsymbol{I} \right)$$
$$= const. - \ln \underbrace{p(\boldsymbol{y} \mid \boldsymbol{\beta}, \boldsymbol{\Phi}(\boldsymbol{X}), \sigma^2)}_{\text{Likelihood}} - \ln \underbrace{p(\boldsymbol{\beta})}_{\text{prior}}$$

Regularized regression equivalent to MAP estimation

Most alternative choices of regularizers and loss functions can be mapped to an equivalent probabilistic representation in a similar way.

- So far: minimization of error functions.
- Back to probabilities?

$$E(\boldsymbol{\beta}) = \underbrace{\frac{1}{2} \sum_{n=1}^{N} (y_n - \boldsymbol{\phi}(\boldsymbol{x}_n) \cdot \boldsymbol{\beta})^2}_{\text{Squared error}} + \underbrace{\frac{\lambda}{2} \boldsymbol{\beta}^{\top} \boldsymbol{\beta}}_{\text{Regularizer}}$$
$$= const. - \sum_{n=1}^{N} \ln \mathcal{N} \left(y_n \mid \boldsymbol{\phi}(\boldsymbol{x}_n) \cdot \boldsymbol{\beta}, \sigma^2 \right) - \ln \mathcal{N} \left(\boldsymbol{\beta} \mid \boldsymbol{0}, \frac{1}{\lambda} \boldsymbol{I} \right)$$
$$= const. - \ln \underbrace{p(\boldsymbol{y} \mid \boldsymbol{\beta}, \boldsymbol{\Phi}(\boldsymbol{X}), \sigma^2)}_{\text{Likelihood}} - \ln \underbrace{p(\boldsymbol{\beta})}_{\text{prior}}$$

- Regularized regression equivalent to MAP estimation
- Most alternative choices of regularizers and loss functions can be mapped to an equivalent probabilistic representation in a similar way.

Outline

Linear Regression II

Bayesian linear regression

Model comparison and hypothesis testing

Summary



Likelihood as before

$$p(\boldsymbol{y} \mid \boldsymbol{X}, \boldsymbol{\beta}, \sigma^2) = \prod_{n=1}^{N} \mathcal{N} \left(y_n \mid \boldsymbol{\phi}(\boldsymbol{x}_n) \cdot \boldsymbol{\beta}, \sigma^2 \right)$$

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

Define a conjugate prior over *β*

 $p(\boldsymbol{\beta}) = \mathcal{N}(\boldsymbol{\beta} \mid \boldsymbol{m}_0, \boldsymbol{S}_0)$

Likelihood as before

$$p(\boldsymbol{y} \mid \boldsymbol{X}, \boldsymbol{\beta}, \sigma^2) = \prod_{n=1}^{N} \mathcal{N}(y_n \mid \boldsymbol{\phi}(\boldsymbol{x}_n) \cdot \boldsymbol{\beta}, \sigma^2)$$

• Define a conjugate prior over $oldsymbol{eta}$

 $p(\boldsymbol{\beta}) = \mathcal{N}(\boldsymbol{\beta} \mid \boldsymbol{m}_0, \boldsymbol{S}_0)$



• Posterior probability of β

$$p(\boldsymbol{\beta} \mid \boldsymbol{y}, \boldsymbol{X}, \sigma^2) \propto \prod_{n=1}^{N} \mathcal{N} \left(y_n \mid \boldsymbol{\phi}(\boldsymbol{x}_n) \cdot \boldsymbol{\beta}, \sigma^2 \right) \cdot \mathcal{N} \left(\boldsymbol{\beta} \mid \boldsymbol{m}_0, \boldsymbol{S}_0 \right) \\ = \mathcal{N} \left(\boldsymbol{y} \mid \boldsymbol{\Phi}(\boldsymbol{X}) \cdot \boldsymbol{\beta}, \sigma^2 \boldsymbol{I} \right) \cdot \mathcal{N} \left(\boldsymbol{\beta} \mid \boldsymbol{m}_0, \boldsymbol{S}_0 \right) \\ = \mathcal{N} \left(\boldsymbol{\beta} \mid \boldsymbol{\mu}_{\boldsymbol{\beta}}, \boldsymbol{\Sigma}_{\boldsymbol{\beta}} \right)$$

where

$$egin{split} oldsymbol{\mu}_{oldsymbol{eta}} &= oldsymbol{\varSigma}_{oldsymbol{eta}} \left(oldsymbol{S}_0^{-1}oldsymbol{m}_0 + rac{1}{\sigma^2}oldsymbol{\varPhi}(oldsymbol{X})^{ op}oldsymbol{y}
ight) \ oldsymbol{\varSigma}_{oldsymbol{eta}} &= \left[oldsymbol{S}_0^{-1} + rac{1}{\sigma^2}oldsymbol{\varPhi}(oldsymbol{X})^{ op}oldsymbol{\varPhi}(oldsymbol{X})
ight]^{-1} \end{split}$$

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

Prior choice

Choice of prior: regularized (ridge) regression

$$p(\boldsymbol{\beta}) = \mathcal{N}(\boldsymbol{\beta} | \boldsymbol{m}_0, \boldsymbol{S}_0).$$

In this case

$$egin{aligned} p(oldsymbol{eta} \mid oldsymbol{y}, oldsymbol{X}, \sigma^2) & \propto \mathcal{N} \left(oldsymbol{eta} \mid oldsymbol{m}_N, oldsymbol{S}_N
ight) \ oldsymbol{m}_N &= oldsymbol{S}_N \left(oldsymbol{S}_0^{-1} oldsymbol{m}_0 + rac{1}{\sigma^2} oldsymbol{\varPhi}(oldsymbol{X})^ op oldsymbol{y}
ight) \ oldsymbol{S}_N &= \left[oldsymbol{S}_0^{-1} + rac{1}{\sigma^2} oldsymbol{\varPhi}(oldsymbol{X})^ op oldsymbol{\varPhi}(oldsymbol{X})^ op oldsymbol{\varPhi}(oldsymbol{X})
ight]^{-1} \end{aligned}$$

- *m_N* is equal to the ridge regression (L₂) estimate for β
 (Exercise: derive both and compare!)
- Equivalent to maximum likelihood estimate for $\lambda \rightarrow 0!$

Prior choice

Choice of prior: regularized (ridge) regression

$$p(\boldsymbol{\beta}) = \mathcal{N}\big(\boldsymbol{\beta} \,|\, \boldsymbol{0}, \frac{1}{\lambda}\boldsymbol{I}\big).$$

In this case

$$p(oldsymbol{eta} \mid oldsymbol{y}, oldsymbol{X}, \sigma^2) \propto \mathcal{N} \left(oldsymbol{eta} \mid oldsymbol{m}_N, oldsymbol{S}_N
ight)
onumber \ oldsymbol{m}_N = oldsymbol{S}_N \left(egin{array}{c} oldsymbol{m}_N, oldsymbol{S}_N
ight)
onumber \ oldsymbol{S}_N = \left[\lambda oldsymbol{I} + rac{1}{\sigma^2} oldsymbol{\varPhi}(oldsymbol{X})^{ op} oldsymbol{\varPhi}(oldsymbol{X})
onumber \ oldsymbol{\Phi}(oldsymbol{X})
onumber \ oldsymbol{J}^{-1}$$

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

- *m_N* is equal to the ridge regression (L₂) estimate for β
 (Exercise: derive both and compare!)
- Equivalent to maximum likelihood estimate for $\lambda \rightarrow 0!$

Prior choice

Choice of prior: regularized (ridge) regression

$$p(\boldsymbol{\beta}) = \mathcal{N}\big(\boldsymbol{\beta} \,|\, \boldsymbol{0}, \frac{1}{\lambda}\boldsymbol{I}\big)$$

In this case

$$p(oldsymbol{eta} \mid oldsymbol{y}, oldsymbol{X}, \sigma^2) \propto \mathcal{N} \left(oldsymbol{eta} \mid oldsymbol{m}_N, oldsymbol{S}_N
ight)
onumber \ oldsymbol{m}_N = oldsymbol{S}_N \left(egin{array}{c} oldsymbol{m}_N, oldsymbol{S}_N
ight)
onumber \ oldsymbol{S}_N = \left[\lambda oldsymbol{I} + rac{1}{\sigma^2} oldsymbol{\Phi}(oldsymbol{X})^{ op} oldsymbol{\Phi}(oldsymbol{X})
onumber \ oldsymbol{D}
ight)
onumber \ oldsymbol{S}_N = \left[\lambda oldsymbol{I} + rac{1}{\sigma^2} oldsymbol{\Phi}(oldsymbol{X})^{ op} oldsymbol{\Phi}(oldsymbol{X})
onumber \ oldsymbol{D}
ight]^{-1}$$

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

m_N is equal to the ridge regression (L₂) estimate for β
 (Exercise: derive both and compare!)

• Equivalent to maximum likelihood estimate for $\lambda \rightarrow 0!$

Prior choice

Choice of prior: regularized (ridge) regression

$$p(\boldsymbol{\beta}) = \mathcal{N}(\boldsymbol{\beta} \mid \mathbf{0}, \frac{1}{\lambda} \boldsymbol{I})$$

In this case

$$p(oldsymbol{eta} \mid oldsymbol{y}, oldsymbol{X}, \sigma^2) \propto \mathcal{N} \left(oldsymbol{eta} \mid oldsymbol{m}_N, oldsymbol{S}_N
ight)
onumber \ oldsymbol{m}_N = oldsymbol{S}_N \left(egin{array}{c} 1 & oldsymbol{m}_N, oldsymbol{S}_N
ight)
onumber \ oldsymbol{S}_N = \left[\lambda oldsymbol{I} + rac{1}{\sigma^2} oldsymbol{\Phi}(oldsymbol{X})^{ op} oldsymbol{\Phi}(oldsymbol{X})
ight]^{-1}$$

- *m_N* is equal to the ridge regression (L₂) estimate for β
 (Exercise: derive both and compare!)
- Equivalent to maximum likelihood estimate for $\lambda \to 0!$

Example: sequential Bayesian learning





$$\prod_{n=1}^{N} \mathcal{N}\left(y_{n} \mid \beta_{0} + x_{n}\beta_{1} \,,\, \sigma^{2}\right)$$

n=1prior

$$\mathcal{N}\left(oldsymbol{eta} \mid oldsymbol{0} \,, \, rac{1}{\lambda}oldsymbol{I}
ight)$$

- This prior is conjugate, so we can do sequential learning
- 1 data point
- 2 data points
- 20 data points

Example: sequential Bayesian learning

likelihood

$$\prod_{n=1}^{N} \mathcal{N}\left(y_n \mid \beta_0 + x_n \beta_1, \sigma^2\right)$$

prior

.

$$\mathcal{N}\left(oldsymbol{eta} \mid oldsymbol{0} \,, \, rac{1}{\lambda}oldsymbol{I}
ight)$$



- 1 data point



Example: sequential Bayesian learning

likelihood

$$\prod_{n=1}^{N} \mathcal{N}\left(y_n \mid \beta_0 + x_n \beta_1, \, \sigma^2\right)$$

n=1prior

.

$$\mathcal{N}\left(oldsymbol{eta} \mid oldsymbol{0} \,, \, rac{1}{\lambda}oldsymbol{I}
ight)$$

- This prior is conjugate, so we can do sequential learning
- 1 data point
- 2 data points



20 data points

Example: sequential Bayesian learning

likelihood

$$\prod_{n=1}^{N} \mathcal{N}\left(y_n \mid \beta_0 + x_n \beta_1, \, \sigma^2\right)$$

n=1prior

$$\mathcal{N}\left(oldsymbol{eta} \mid oldsymbol{0} \,, \, rac{1}{\lambda}oldsymbol{I}
ight)$$

- This prior is conjugate, so we can do sequential learning
- 1 data point
- 2 data points
- 20 data points



• Prediction for fixed weight estimate $\hat{oldsymbol{eta}}$ at input x^{\star} trivial:

$$p(y^{\star} \mid \boldsymbol{x}^{\star}, \hat{\boldsymbol{\beta}}, \sigma^2) = \mathcal{N}\left(y^{\star} \mid \boldsymbol{\phi}(\boldsymbol{x}^{\star}) \hat{\boldsymbol{\beta}}, \sigma^2\right)$$

Integrate over β to take the posterior uncertainty into account

$$p(\boldsymbol{y}^{\star} \,|\, \boldsymbol{x}^{\star}, \mathcal{D}) \propto \int_{\boldsymbol{\beta}} p(\boldsymbol{y}^{\star} \,|\, \boldsymbol{x}^{\star}, \boldsymbol{\beta}, \sigma^{2}) p(\boldsymbol{\beta} \,|\, \boldsymbol{X}, \boldsymbol{y}, \sigma^{2}) \, \mathrm{d}\boldsymbol{\beta}$$

Key:

- prediction is again Gaussian
- Predictive variance is increased due to the posterior uncertainty in eta .

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

• Prediction for fixed weight estimate $\hat{oldsymbol{eta}}$ at input x^{\star} trivial:

$$p(\boldsymbol{y}^{\star} \mid \boldsymbol{x}^{\star}, \hat{\boldsymbol{\beta}}, \sigma^{2}) = \mathcal{N}\left(\boldsymbol{y}^{\star} \mid \boldsymbol{\phi}(\boldsymbol{x}^{\star}) \hat{\boldsymbol{\beta}}, \sigma^{2}\right)$$

 \blacktriangleright Integrate over β to take the posterior uncertainty into account

$$p(y^{\star} | \boldsymbol{x}^{\star}, \mathcal{D}) \propto \int_{\boldsymbol{\beta}} p(y^{\star} | \boldsymbol{x}^{\star}, \boldsymbol{\beta}, \sigma^{2}) p(\boldsymbol{\beta} | \boldsymbol{X}, \boldsymbol{y}, \sigma^{2}) \, \mathrm{d}\boldsymbol{\beta}$$

Key:

- prediction is again Gaussian
- \succ Predictive variance is increased due to the posterior uncertainty in eta.

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

• Prediction for fixed weight estimate $\hat{oldsymbol{eta}}$ at input x^{\star} trivial:

$$p(y^{\star} \mid \boldsymbol{x}^{\star}, \hat{\boldsymbol{\beta}}, \sigma^2) = \mathcal{N}\left(y^{\star} \mid \boldsymbol{\phi}(\boldsymbol{x}^{\star}) \hat{\boldsymbol{\beta}}, \sigma^2\right)$$

 \blacktriangleright Integrate over β to take the posterior uncertainty into account

$$p(y^{\star} \mid \boldsymbol{x}^{\star}, \mathcal{D}) \propto \int_{\boldsymbol{eta}} p(y^{\star} \mid \boldsymbol{x}^{\star}, \boldsymbol{eta}, \sigma^{2}) p(\boldsymbol{eta} \mid \boldsymbol{X}, \boldsymbol{y}, \sigma^{2}) \, \mathrm{d}\boldsymbol{eta}$$

 $\propto \int_{\boldsymbol{eta}} \mathcal{N}\left(y^{\star} \mid \boldsymbol{\phi}(\boldsymbol{x}^{\star}) \boldsymbol{eta}, \sigma^{2}\right) \mathcal{N}\left(\boldsymbol{eta} \mid \boldsymbol{m}_{N}, \boldsymbol{S}_{N}\right)$

Key:

- prediction is again Gaussian
- \succ Predictive variance is increased due to the posterior uncertainty in eta.

• Prediction for fixed weight estimate \hat{eta} at input x^{\star} trivial:

$$p(y^{\star} \mid \boldsymbol{x}^{\star}, \hat{\boldsymbol{\beta}}, \sigma^2) = \mathcal{N}\left(y^{\star} \mid \boldsymbol{\phi}(\boldsymbol{x}^{\star}) \hat{\boldsymbol{\beta}}, \sigma^2\right)$$

 \blacktriangleright Integrate over β to take the posterior uncertainty into account

$$\begin{split} p(y^{\star} \mid \boldsymbol{x}^{\star}, \mathcal{D}) &\propto \int_{\boldsymbol{\beta}} p(y^{\star} \mid \boldsymbol{x}^{\star}, \boldsymbol{\beta}, \sigma^{2}) p(\boldsymbol{\beta} \mid \boldsymbol{X}, \boldsymbol{y}, \sigma^{2}) \, \mathrm{d}\boldsymbol{\beta} \\ &\propto \int_{\boldsymbol{\beta}} \mathcal{N} \left(y^{\star} \mid \boldsymbol{\phi}(\boldsymbol{x}^{\star}) \boldsymbol{\beta}, \sigma^{2} \right) \mathcal{N} \left(\boldsymbol{\beta} \mid \boldsymbol{m}_{N}, \boldsymbol{S}_{N} \right) \\ &\propto \int_{\boldsymbol{\phi}(\boldsymbol{x}^{\star}) \boldsymbol{\beta}} \mathcal{N} \left(\boldsymbol{\phi}(\boldsymbol{x}^{\star}) \boldsymbol{\beta} \mid y^{\star}, \sigma^{2} \right) \mathcal{N} \left(\boldsymbol{\phi}(\boldsymbol{x}^{\star}) \boldsymbol{\beta} \mid \boldsymbol{\phi}(\boldsymbol{x}^{\star}) \boldsymbol{m}_{N}, \boldsymbol{\phi}(\boldsymbol{x}^{\star}) \boldsymbol{S}_{N} \boldsymbol{\phi}(\boldsymbol{x}^{\star})^{\top} \right) \end{split}$$

Key:

- prediction is again Gaussian
- Predictive variance is increased due to the posterior uncertainty in β .

(日) (日) (日) (日) (日) (日) (日) (日)

• Prediction for fixed weight estimate \hat{eta} at input x^{\star} trivial:

$$p(y^{\star} \mid \boldsymbol{x}^{\star}, \hat{\boldsymbol{\beta}}, \sigma^{2}) = \mathcal{N}\left(y^{\star} \mid \boldsymbol{\phi}(\boldsymbol{x}^{\star}) \hat{\boldsymbol{\beta}}, \sigma^{2}\right)$$

 \blacktriangleright Integrate over β to take the posterior uncertainty into account

$$\begin{split} p(y^{\star} \mid \boldsymbol{x}^{\star}, \mathcal{D}) &\propto \int_{\boldsymbol{\beta}} p(y^{\star} \mid \boldsymbol{x}^{\star}, \boldsymbol{\beta}, \sigma^{2}) p(\boldsymbol{\beta} \mid \boldsymbol{X}, \boldsymbol{y}, \sigma^{2}) \, \mathrm{d}\boldsymbol{\beta} \\ &\propto \int_{\boldsymbol{\beta}} \mathcal{N} \left(y^{\star} \mid \boldsymbol{\phi}(\boldsymbol{x}^{\star}) \boldsymbol{\beta}, \sigma^{2} \right) \mathcal{N} \left(\boldsymbol{\beta} \mid \boldsymbol{m}_{N}, \boldsymbol{S}_{N} \right) \\ &\propto \int_{\boldsymbol{\phi}(\boldsymbol{x}^{\star}) \boldsymbol{\beta}} \mathcal{N} \left(\boldsymbol{\phi}(\boldsymbol{x}^{\star}) \boldsymbol{\beta} \mid y^{\star}, \sigma^{2} \right) \mathcal{N} \left(\boldsymbol{\phi}(\boldsymbol{x}^{\star}) \boldsymbol{\beta} \mid \boldsymbol{\phi}(\boldsymbol{x}^{\star}) \boldsymbol{m}_{N}, \boldsymbol{\phi}(\boldsymbol{x}^{\star}) \boldsymbol{S}_{N} \boldsymbol{\phi}(\boldsymbol{x}^{\star})^{\top} \right) \\ &\propto \mathcal{N} \left(y^{\star} \mid \boldsymbol{\phi}(\boldsymbol{x}^{\star}) \cdot \boldsymbol{m}_{N}, \sigma^{2} + \boldsymbol{\phi}(\boldsymbol{x}^{\star})^{\top} \boldsymbol{S}_{N} \boldsymbol{\phi}(\boldsymbol{x}^{\star}) \right) \end{split}$$

Key:

prediction is again Gaussian

Predictive variance is increased due to the posterior uncertainty in β.

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

• Prediction for fixed weight estimate \hat{eta} at input x^{\star} trivial:

$$p(\boldsymbol{y}^{\star} \mid \boldsymbol{x}^{\star}, \hat{\boldsymbol{\beta}}, \sigma^{2}) = \mathcal{N}\left(\boldsymbol{y}^{\star} \mid \boldsymbol{\phi}(\boldsymbol{x}^{\star}) \hat{\boldsymbol{\beta}}, \sigma^{2}\right)$$

 \blacktriangleright Integrate over β to take the posterior uncertainty into account

$$\begin{split} p(y^{\star} \mid \boldsymbol{x}^{\star}, \mathcal{D}) &\propto \int_{\boldsymbol{\beta}} p(y^{\star} \mid \boldsymbol{x}^{\star}, \boldsymbol{\beta}, \sigma^{2}) p(\boldsymbol{\beta} \mid \boldsymbol{X}, \boldsymbol{y}, \sigma^{2}) \, \mathrm{d}\boldsymbol{\beta} \\ &\propto \int_{\boldsymbol{\beta}} \mathcal{N} \left(y^{\star} \mid \boldsymbol{\phi}(\boldsymbol{x}^{\star}) \boldsymbol{\beta}, \sigma^{2} \right) \mathcal{N} \left(\boldsymbol{\beta} \mid \boldsymbol{m}_{N}, \boldsymbol{S}_{N} \right) \\ &\propto \int_{\boldsymbol{\phi}(\boldsymbol{x}^{\star}) \boldsymbol{\beta}} \mathcal{N} \left(\boldsymbol{\phi}(\boldsymbol{x}^{\star}) \boldsymbol{\beta} \mid y^{\star}, \sigma^{2} \right) \mathcal{N} \left(\boldsymbol{\phi}(\boldsymbol{x}^{\star}) \boldsymbol{\beta} \mid \boldsymbol{\phi}(\boldsymbol{x}^{\star}) \boldsymbol{m}_{N}, \boldsymbol{\phi}(\boldsymbol{x}^{\star}) \boldsymbol{S}_{N} \boldsymbol{\phi}(\boldsymbol{x}^{\star})^{\top} \right) \\ &\propto \mathcal{N} \left(y^{\star} \mid \boldsymbol{\phi}(\boldsymbol{x}^{\star}) \cdot \boldsymbol{m}_{N}, \sigma^{2} + \boldsymbol{\phi}(\boldsymbol{x}^{\star})^{\top} \boldsymbol{S}_{N} \boldsymbol{\phi}(\boldsymbol{x}^{\star}) \right) \end{split}$$

Key:

- prediction is again Gaussian
- Predictive variance is increased due to the posterior uncertainty in β .

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

Marginal variance for x^{\star}

 $\sigma^2 + \boldsymbol{\phi}(\boldsymbol{x}^\star)^\top \boldsymbol{S}_N \boldsymbol{\phi}(\boldsymbol{x}^\star)$

Predictive covariance

 $oldsymbol{\phi}(oldsymbol{x})^{ op}oldsymbol{S}_Noldsymbol{\phi}(oldsymbol{x}')$

Experiment:

9 Gaussian basis functions

Empirical Observations:

- Mariance approaches noise variance for large sample size
 - Co-variance between close or values is

Marginal variance for x^{\star}

 $\sigma^2 + \boldsymbol{\phi}(\boldsymbol{x}^\star)^\top \boldsymbol{S}_N \boldsymbol{\phi}(\boldsymbol{x}^\star)$

Predictive covariance

 $oldsymbol{\phi}(oldsymbol{x})^{ op}oldsymbol{S}_Noldsymbol{\phi}(oldsymbol{x}')$

Experiment:

9 Gaussian basis functions

Empirical Observations:

- Variance approaches noise variance for large sample size
- Co-variance between close *a* values is

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Marginal variance for x^{\star}

 $\sigma^2 + oldsymbol{\phi}(oldsymbol{x}^\star)^ op oldsymbol{S}_N oldsymbol{\phi}(oldsymbol{x}^\star)$



Experiment:

9 Gaussian basis functions

Predictive covariance

 $oldsymbol{\phi}(oldsymbol{x})^{ op}oldsymbol{S}_Noldsymbol{\phi}(oldsymbol{x}')$

Visualize by sampling from the posterior of β



Empirical Observations:

- Variance approaches noise variance for large sample size
- Co-variance between close a values is

Marginal variance for x^{\star}

 $\sigma^2 + oldsymbol{\phi}(oldsymbol{x}^\star)^ op oldsymbol{S}_N oldsymbol{\phi}(oldsymbol{x}^\star)$



Experiment:

9 Gaussian basis functions

Predictive covariance

 $oldsymbol{\phi}(oldsymbol{x})^{ op}oldsymbol{S}_Noldsymbol{\phi}(oldsymbol{x}')$

Visualize by sampling from the posterior of β



Empirical Observations:

- Variance approaches noise variance for large sample size
- Co-variance between close a values is

Marginal variance for x^{\star}

 $\sigma^2 + \boldsymbol{\phi}(\boldsymbol{x}^\star)^\top \boldsymbol{S}_N \boldsymbol{\phi}(\boldsymbol{x}^\star)$



Experiment:

9 Gaussian basis functions

Predictive covariance

 $oldsymbol{\phi}(oldsymbol{x})^{ op}oldsymbol{S}_Noldsymbol{\phi}(oldsymbol{x}')$

Visualize by sampling from the posterior of β



Empirical Observations:

- Variance approaches noise variance for large sample size
- Co-variance between close *x* values is
Predictive distribution

Marginal variance for x^{\star}

 $\sigma^2 + oldsymbol{\phi}(oldsymbol{x}^\star)^ op oldsymbol{S}_N oldsymbol{\phi}(oldsymbol{x}^\star)$



Experiment:

9 Gaussian basis functions

Predictive covariance

 $oldsymbol{\phi}(oldsymbol{x})^{ op}oldsymbol{S}_Noldsymbol{\phi}(oldsymbol{x}')$

Visualize by sampling from the posterior of β



Empirical Observations:

- Variance approaches noise variance for large sample size
- ► Co-variance between close *x* values is high

Predictive distribution

Marginal variance for x^{\star}

 $\sigma^2 + oldsymbol{\phi}(oldsymbol{x}^\star)^ op oldsymbol{S}_N oldsymbol{\phi}(oldsymbol{x}^\star)$



Experiment:

9 Gaussian basis functions

Predictive covariance

 $oldsymbol{\phi}(oldsymbol{x})^{ op}oldsymbol{S}_Noldsymbol{\phi}(oldsymbol{x}')$

Visualize by sampling from the posterior of β



Empirical Observations:

- Variance approaches noise variance for large sample size
- ► Co-variance between close *x* values is high

Predictive distribution

Marginal variance for x^{\star}

 $\sigma^2 + \boldsymbol{\phi}(\boldsymbol{x}^\star)^\top \boldsymbol{S}_N \boldsymbol{\phi}(\boldsymbol{x}^\star)$



Experiment:

9 Gaussian basis functions

Predictive covariance

 $oldsymbol{\phi}(oldsymbol{x})^{ op}oldsymbol{S}_Noldsymbol{\phi}(oldsymbol{x}')$

Visualize by sampling from the posterior of β



Empirical Observations:

- Variance approaches noise variance for large sample size
- ► Co-variance between close *x* values is high

Outline

Linear Regression II

Bayesian linear regression

Model comparison and hypothesis testing

Summary



Model comparison

Motivation

What degree of polynomials describes the data best?

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

- Is the linear model at all appropriate?
- Association testing.

Model comparison

Motivation

- What degree of polynomials describes the data best?
- Is the linear model at all appropriate?
- Association testing.



▲□▶ ▲□▶ ▲□▶ ▲□▶ □□ ● ● ●

- How do we choose among alternative models?
- ► Assume we want to choose among models H₀,..., H_M for a dataset D.
- Posterior probability for a particular model i

- How do we choose among alternative models?
- ► Assume we want to choose among models H₀,..., H_M for a dataset D.
- \blacktriangleright Posterior probability for a particular model i

$$p(\mathcal{H}_i \mid \mathcal{D}) \propto \underbrace{p(\mathcal{D} \mid \mathcal{H}_i)}_{\text{Evidence}} \underbrace{p(\mathcal{H}_i)}_{\text{Prior}}$$

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

How to calculate the evidence

The evidence is not the model likelihood!

$$p(\mathcal{D} \,|\, \mathcal{H}_i) = \int_{\boldsymbol{\varTheta}} p(\mathcal{D} \,|\, \boldsymbol{\varTheta}) p(\boldsymbol{\varTheta}) \,\mathrm{d}\boldsymbol{\varTheta} \,\, \text{for model parameters }\, \boldsymbol{\varTheta}.$$

Remember:

$$p(\boldsymbol{\Theta} \mid \mathcal{H}_i, \mathcal{D}) = \frac{p(\mathcal{D} \mid \mathcal{H}_i, \boldsymbol{\Theta})p(\boldsymbol{\Theta})}{p(\mathcal{D} \mid \mathcal{H}_i)}$$

(ロ)、(型)、(E)、(E)、 E) の(の)

How to calculate the evidence

The evidence is not the model likelihood!

$$p(\mathcal{D} \mid \mathcal{H}_i) = \int_{\boldsymbol{\Theta}} p(\mathcal{D} \mid \boldsymbol{\Theta}) p(\boldsymbol{\Theta}) \, \mathrm{d}\boldsymbol{\Theta} \text{ for model parameters } \boldsymbol{\Theta}.$$

Remember:

$$p(\boldsymbol{\Theta} \mid \mathcal{H}_i, \mathcal{D}) = \frac{p(\mathcal{D} \mid \mathcal{H}_i, \boldsymbol{\Theta})p(\boldsymbol{\Theta})}{p(\mathcal{D} \mid \mathcal{H}_i)}$$

posterior =
$$\frac{\text{likelihood} \cdot \text{prior}}{\text{Evidence}}$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

Bayesian Occam's razor

The evidence integral penalizes overly complex models.

- A model with few parameters and lower maximum likelihood (*H*₁) may win over a model with a peaked likelihood that requires many more parameters (*H*₂).
- When averaging the likelihood over all possible parameters, more complex models have low fit for most of the setting, resulting in a lower evidence
- Complex models have low average over many possible data sets
- Simple models have large evidence on a small range of data sets, extremely low evidence otherwise





Learning)

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへで

Bayesian Occam's razor

- The evidence integral penalizes overly complex models.
- A model with few parameters and lower maximum likelihood (H₁) may win over a model with a peaked likelihood that requires many more parameters (H₂).
- When averaging the likelihood over all possible parameters, more complex models have low fit for most of the setting, resulting in a lower evidence
- Complex models have low average over many possible data sets
- Simple models have large evidence on a small range of data sets, extremely low evidence otherwise





Learning)

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

Bayesian Occam's razor

- The evidence integral penalizes overly complex models.
- A model with few parameters and lower maximum likelihood (*H*₁) may win over a model with a peaked likelihood that requires many more parameters (*H*₂).
- When averaging the likelihood over all possible parameters, more complex models have low fit for most of the setting, resulting in a lower evidence
- Complex models have low average over many possible data sets
- Simple models have large evidence on a small range of data sets, extremely low evidence otherwise





Learning)

Bayesian Occam's razor

- The evidence integral penalizes overly complex models.
- A model with few parameters and lower maximum likelihood (*H*₁) may win over a model with a peaked likelihood that requires many more parameters (*H*₂).
- When averaging the likelihood over all possible parameters, more complex models have low fit for most of the setting, resulting in a lower evidence
- Complex models have low average over many possible data sets
- Simple models have large evidence on a small range of data sets, extremely low evidence otherwise





Learning)

Bayesian Occam's razor

- The evidence integral penalizes overly complex models.
- A model with few parameters and lower maximum likelihood (*H*₁) may win over a model with a peaked likelihood that requires many more parameters (*H*₂).
- When averaging the likelihood over all possible parameters, more complex models have low fit for most of the setting, resulting in a lower evidence
- Complex models have low average over many possible data sets
- Simple models have large evidence on a small range of data sets, extremely low evidence otherwise





Learning)

Relevance of a single SNP

- Consider an association study.
 - \mathcal{H}_0 : no association

$$p(\boldsymbol{y} \mid \mathcal{H}_0, \boldsymbol{X}, \boldsymbol{\Theta}_0) = \mathcal{N} \left(\boldsymbol{y} \mid \boldsymbol{0}, \sigma^2 \boldsymbol{I} \right)$$
$$p(\mathcal{D} \mid \mathcal{H}_0) = \int_{\sigma^2} \mathcal{N} \left(\boldsymbol{y} \mid \boldsymbol{0}, \sigma^2 \boldsymbol{I} \right) p(\sigma^2)$$

H₁: linear association

$$p(\boldsymbol{y} \mid \mathcal{H}_{1}, \boldsymbol{x}_{i}, \boldsymbol{\Theta}_{1}) = \mathcal{N} \left(\boldsymbol{y} \mid \boldsymbol{x}_{i} \cdot \boldsymbol{\beta}, \sigma^{2} \boldsymbol{I} \right)$$
$$p(\mathcal{D} \mid \mathcal{H}_{1}) = \int_{\sigma^{2}, \boldsymbol{\beta}} \mathcal{N} \left(\boldsymbol{y} \mid \boldsymbol{x}_{i} \cdot \boldsymbol{\beta}, \sigma^{2} \boldsymbol{I} \right) p(\sigma^{2}) p(\boldsymbol{\beta})$$

Depending on the choice of priors, p(σ²) and p(β), the required integrals are often tractable in closed form. (Conjugate priors!)

Relevance of a single SNP

- Consider an association study.
 - $\blacktriangleright \ \mathcal{H}_0: \text{no association}$

$$p(\boldsymbol{y} \mid \mathcal{H}_0, \boldsymbol{X}, \boldsymbol{\Theta}_0) = \mathcal{N} \left(\boldsymbol{y} \mid \boldsymbol{0}, \sigma^2 \boldsymbol{I} \right)$$
$$p(\mathcal{D} \mid \mathcal{H}_0) = \int_{\sigma^2} \mathcal{N} \left(\boldsymbol{y} \mid \boldsymbol{0}, \sigma^2 \boldsymbol{I} \right) p(\sigma^2)$$

• \mathcal{H}_1 : linear association

$$p(\boldsymbol{y} \mid \mathcal{H}_{1}, \boldsymbol{x}_{i}, \boldsymbol{\Theta}_{1}) = \mathcal{N} \left(\boldsymbol{y} \mid \boldsymbol{x}_{i} \cdot \boldsymbol{\beta}, \sigma^{2} \boldsymbol{I} \right)$$
$$p(\mathcal{D} \mid \mathcal{H}_{1}) = \int_{\sigma^{2}, \boldsymbol{\beta}} \mathcal{N} \left(\boldsymbol{y} \mid \boldsymbol{x}_{i} \cdot \boldsymbol{\beta}, \sigma^{2} \boldsymbol{I} \right) p(\sigma^{2}) p(\boldsymbol{\beta})$$

Depending on the choice of priors, p(σ²) and p(β), the required integrals are often tractable in closed form. (Conjugate priors!)

Relevance of a single SNP

- Consider an association study.
 - $\blacktriangleright \ \mathcal{H}_0: \text{no association}$

$$p(\boldsymbol{y} \mid \mathcal{H}_0, \boldsymbol{X}, \boldsymbol{\Theta}_0) = \mathcal{N} \left(\boldsymbol{y} \mid \boldsymbol{0}, \sigma^2 \boldsymbol{I} \right)$$
$$p(\mathcal{D} \mid \mathcal{H}_0) = \int_{\sigma^2} \mathcal{N} \left(\boldsymbol{y} \mid \boldsymbol{0}, \sigma^2 \boldsymbol{I} \right) p(\sigma^2)$$

• \mathcal{H}_1 : linear association

$$p(\boldsymbol{y} \mid \mathcal{H}_{1}, \boldsymbol{x}_{i}, \boldsymbol{\Theta}_{1}) = \mathcal{N} \left(\boldsymbol{y} \mid \boldsymbol{x}_{i} \cdot \boldsymbol{\beta}, \sigma^{2} \boldsymbol{I} \right)$$
$$p(\mathcal{D} \mid \mathcal{H}_{1}) = \int_{\sigma^{2}, \boldsymbol{\beta}} \mathcal{N} \left(\boldsymbol{y} \mid \boldsymbol{x}_{i} \cdot \boldsymbol{\beta}, \sigma^{2} \boldsymbol{I} \right) p(\sigma^{2}) p(\boldsymbol{\beta})$$

Depending on the choice of priors, p(σ²) and p(β), the required integrals are often tractable in closed form. (Conjugate priors!)

Application to GWAS Scoring models

Similar to likelihood ratios, the ratio of the evidences, the Bayes factor can be used to score alternative models:

$$BF = \ln \frac{p(\mathcal{D} \mid \mathcal{H}_1)}{p(\mathcal{D} \mid \mathcal{H}_0)}.$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

Application to GWAS Scoring models

Similar to likelihood ratios, the ratio of the evidences, the Bayes factor can be used to score alternative models:

$$BF = \ln \frac{p(\mathcal{D} \mid \mathcal{H}_1)}{p(\mathcal{D} \mid \mathcal{H}_0)}$$



◆□▶ ◆圖▶ ◆臣▶ ◆臣▶ ─ 臣

Posterior probability of an association

 Bayes factors are useful, however we would like a probabilistic answer how certain an association really is.

• Posterior probability of \mathcal{H}_1

$$p(\mathcal{H}_1 \mid \mathcal{D}) = \frac{p(\mathcal{D} \mid \mathcal{H}_1)p(\mathcal{H}_1)}{p(\mathcal{D})}$$
$$= \frac{p(\mathcal{D} \mid \mathcal{H}_1)p(\mathcal{H}_1)}{p(\mathcal{D} \mid \mathcal{H}_1)p(\mathcal{H}_1) + p(\mathcal{D} \mid \mathcal{H}_0)p(\mathcal{H}_0)}$$

▶ p(H₁ | D) + p(H₀ | D) = 1, prior probability of observing a real association.

Posterior probability of an association

- Bayes factors are useful, however we would like a probabilistic answer how certain an association really is.
- Posterior probability of \mathcal{H}_1

$$p(\mathcal{H}_1 \mid \mathcal{D}) = \frac{p(\mathcal{D} \mid \mathcal{H}_1)p(\mathcal{H}_1)}{p(\mathcal{D})}$$
$$= \frac{p(\mathcal{D} \mid \mathcal{H}_1)p(\mathcal{H}_1)}{p(\mathcal{D} \mid \mathcal{H}_1)p(\mathcal{H}_1) + p(\mathcal{D} \mid \mathcal{H}_0)p(\mathcal{H}_0)}$$

▶ p(H₁ | D) + p(H₀ | D) = 1, prior probability of observing a real association.

Posterior probability of an association

- Bayes factors are useful, however we would like a probabilistic answer how certain an association really is.
- Posterior probability of \mathcal{H}_1

$$p(\mathcal{H}_1 \mid \mathcal{D}) = \frac{p(\mathcal{D} \mid \mathcal{H}_1)p(\mathcal{H}_1)}{p(\mathcal{D})}$$
$$= \frac{p(\mathcal{D} \mid \mathcal{H}_1)p(\mathcal{H}_1)}{p(\mathcal{D} \mid \mathcal{H}_1)p(\mathcal{H}_1) + p(\mathcal{D} \mid \mathcal{H}_0)p(\mathcal{H}_0)}$$

▶ p(H₁ | D) + p(H₀ | D) = 1, prior probability of observing a real association.

Bayes factor versus likelihood ratio

Bayes factor

- Models of different complexity can be objectively compared.
- Statistical significance as posterior probability of a model.

Typically hard to compute.

Likelihood ratio

- Likelihood ratio scales with the number of parameters.
- Likelihood ratios have known null distribution, yielding p-values.

Often easy to compute.

Bayes factor versus likelihood ratio

Bayes factor

- Models of different complexity can be objectively compared.
- Statistical significance as posterior probability of a model.
- Typically hard to compute.

Likelihood ratio

- Likelihood ratio scales with the number of parameters.
- Likelihood ratios have known null distribution, yielding p-values.

Often easy to compute.

- ► Consider a linear model, accounting for a set of measured SNPs \boldsymbol{X} $p(\boldsymbol{y} \mid \boldsymbol{X}, \boldsymbol{\beta}, \sigma^2) = \mathcal{N}\left(\boldsymbol{y} \mid \sum_{s=1}^{S} \boldsymbol{x}_s \beta_s, \sigma^2 \boldsymbol{I}\right)$
- Choose identical Gaussian prior for all weights $p(\beta) = \prod_{s=1}^{S} \mathcal{N} \left(\beta_s \mid 0, \sigma_g^2 \right)$
- Marginal likelihood

$$p(\boldsymbol{y} \mid \boldsymbol{X},) = \int_{\boldsymbol{\beta}} \mathcal{N}\left(\boldsymbol{y} \mid \boldsymbol{X} \boldsymbol{\beta}, \sigma^{2} \boldsymbol{I}\right) \mathcal{N}\left(\boldsymbol{\beta} \mid \boldsymbol{0}, \sigma_{g}^{2} \boldsymbol{I}\right)$$

・ロット (雪) (日) (日) (日)

• Consider a linear model, accounting for a set of measured SNPs \boldsymbol{X} $p(\boldsymbol{y} \mid \boldsymbol{X}, \boldsymbol{\beta}, \sigma^2) = \mathcal{N}\left(\boldsymbol{y} \mid \sum_{s=1}^{S} \boldsymbol{x}_s \beta_s, \sigma^2 \boldsymbol{I}\right)$

• Choose identical Gaussian prior for all weights $p(\boldsymbol{\beta}) = \prod_{s=1}^{S} \mathcal{N}\left(\beta_s \mid 0, \sigma_g^2\right)$

Marginal likelihood

$$p(\boldsymbol{y} \mid \boldsymbol{X},) = \int_{\boldsymbol{\beta}} \mathcal{N}\left(\boldsymbol{y} \mid \boldsymbol{X} \boldsymbol{\beta}, \sigma^{2} \boldsymbol{I}\right) \mathcal{N}\left(\boldsymbol{\beta} \mid \boldsymbol{0}, \sigma_{g}^{2} \boldsymbol{I}\right)$$

- Consider a linear model, accounting for a set of measured SNPs \boldsymbol{X} $p(\boldsymbol{y} \mid \boldsymbol{X}, \boldsymbol{\beta}, \sigma^2) = \mathcal{N}\left(\boldsymbol{y} \mid \sum_{s=1}^{S} \boldsymbol{x}_s \beta_s, \sigma^2 \boldsymbol{I}\right)$
- Choose identical Gaussian prior for all weights $p(\boldsymbol{\beta}) = \prod_{s=1}^{S} \mathcal{N} \left(\beta_s \mid 0, \sigma_g^2 \right)$
- Marginal likelihood

$$p(\boldsymbol{y} \mid \boldsymbol{X}, \sigma^2, \sigma_{g}^2) = \int_{\boldsymbol{\beta}} \mathcal{N}\left(\boldsymbol{y} \mid \boldsymbol{X}\boldsymbol{\beta}, \sigma^2 \boldsymbol{I}\right) \mathcal{N}\left(\boldsymbol{\beta} \mid \boldsymbol{0}, \sigma_{g}^2 \boldsymbol{I}\right)$$
$$= \mathcal{N}\left(\boldsymbol{y} \mid \boldsymbol{0}, \sigma_{g}^2 \boldsymbol{X} \boldsymbol{X}^{\top} + \sigma^2 \boldsymbol{I}\right)$$

- Consider a linear model, accounting for a set of measured SNPs \boldsymbol{X} $p(\boldsymbol{y} \mid \boldsymbol{X}, \boldsymbol{\beta}, \sigma^2) = \mathcal{N}\left(\boldsymbol{y} \mid \sum_{s=1}^{S} \boldsymbol{x}_s \beta_s, \sigma^2 \boldsymbol{I}\right)$
- Choose identical Gaussian prior for all weights $p(\boldsymbol{\beta}) = \prod_{s=1}^{S} \mathcal{N} \left(\beta_s \mid 0, \sigma_g^2 \right)$
- Marginal likelihood

$$p(\boldsymbol{y} \mid \boldsymbol{X}, \boldsymbol{\sigma}^{2}, \boldsymbol{\sigma}_{g}^{2}) = \int_{\boldsymbol{\beta}} \mathcal{N}\left(\boldsymbol{y} \mid \boldsymbol{X}\boldsymbol{\beta}, \boldsymbol{\sigma}^{2}\boldsymbol{I}\right) \mathcal{N}\left(\boldsymbol{\beta} \mid \boldsymbol{0}, \boldsymbol{\sigma}_{g}^{2}\boldsymbol{I}\right)$$
$$= \mathcal{N}\left(\boldsymbol{y} \mid \boldsymbol{0}, \boldsymbol{\sigma}_{g}^{2}\boldsymbol{X}\boldsymbol{X}^{\top} + \boldsymbol{\sigma}^{2}\boldsymbol{I}\right)$$

Marginal likelihood of variance component models Basis functions

► The analogous derivation can be repeated for a feature mapping ϕ $p(\boldsymbol{y} \mid \boldsymbol{X}, \boldsymbol{\beta}, \sigma^2) = \mathcal{N}\left(\boldsymbol{y} \mid \sum_{s=1}^{S} \boldsymbol{\phi}(\boldsymbol{x}_s)\beta_s, \sigma^2 \boldsymbol{I}\right) =$ $\mathcal{N}\left(\boldsymbol{y} \mid \boldsymbol{\Phi}(\boldsymbol{X})\boldsymbol{\beta}, \sigma^2 \boldsymbol{I}\right)$

Marginal likelihood

$$p(\boldsymbol{y} \mid \boldsymbol{X}, \sigma^{2}, \sigma_{g}^{2}) = \int_{\boldsymbol{\beta}} \mathcal{N} \left(\boldsymbol{y} \mid \boldsymbol{\varPhi}(\boldsymbol{X}) \boldsymbol{\beta}, \sigma^{2} \boldsymbol{I} \right) \mathcal{N} \left(\boldsymbol{\beta} \mid \boldsymbol{0}, \sigma_{g}^{2} \boldsymbol{I} \right)$$
$$= \mathcal{N} \left(\boldsymbol{y} \mid \boldsymbol{0}, \sigma_{g}^{2} \underbrace{\boldsymbol{\varPhi}(\boldsymbol{X}) \boldsymbol{\varPhi}(\boldsymbol{X})^{\top}}_{\boldsymbol{K}} + \sigma^{2} \boldsymbol{I} \right)$$

・ロット (雪) (日) (日) (日)

• K: (N × N) kernel or covariance induced by feature mapping ϕ .

Marginal likelihood of variance component models Basis functions

- ► The analogous derivation can be repeated for a feature mapping ϕ $p(\boldsymbol{y} \mid \boldsymbol{X}, \boldsymbol{\beta}, \sigma^2) = \mathcal{N}\left(\boldsymbol{y} \mid \sum_{s=1}^{S} \boldsymbol{\phi}(\boldsymbol{x}_s)\beta_s, \sigma^2 \boldsymbol{I}\right) =$ $\mathcal{N}\left(\boldsymbol{y} \mid \boldsymbol{\Phi}(\boldsymbol{X})\boldsymbol{\beta}, \sigma^2 \boldsymbol{I}\right)$
- Marginal likelihood

$$p(\boldsymbol{y} \mid \boldsymbol{X}, \sigma^{2}, \sigma_{g}^{2}) = \int_{\boldsymbol{\beta}} \mathcal{N} \left(\boldsymbol{y} \mid \boldsymbol{\varPhi}(\boldsymbol{X}) \boldsymbol{\beta}, \sigma^{2} \boldsymbol{I} \right) \mathcal{N} \left(\boldsymbol{\beta} \mid \boldsymbol{0}, \sigma_{g}^{2} \boldsymbol{I} \right)$$
$$= \mathcal{N} \left(\boldsymbol{y} \mid \boldsymbol{0}, \sigma_{g}^{2} \underbrace{\boldsymbol{\varPhi}(\boldsymbol{X}) \boldsymbol{\varPhi}(\boldsymbol{X})^{\top}}_{\boldsymbol{K}} + \sigma^{2} \boldsymbol{I} \right)$$

・ロット (雪) (日) (日) (日)

K: (N × N) kernel or covariance induced by feature mapping φ.

Marginal likelihood of variance component models Basis functions

- ► The analogous derivation can be repeated for a feature mapping ϕ $p(\boldsymbol{y} \mid \boldsymbol{X}, \boldsymbol{\beta}, \sigma^2) = \mathcal{N}\left(\boldsymbol{y} \mid \sum_{s=1}^{S} \boldsymbol{\phi}(\boldsymbol{x}_s)\beta_s, \sigma^2 \boldsymbol{I}\right) =$ $\mathcal{N}\left(\boldsymbol{y} \mid \boldsymbol{\Phi}(\boldsymbol{X})\boldsymbol{\beta}, \sigma^2 \boldsymbol{I}\right)$
- Marginal likelihood

$$p(\boldsymbol{y} \mid \boldsymbol{X}, \sigma^{2}, \sigma_{g}^{2}) = \int_{\boldsymbol{\beta}} \mathcal{N} \left(\boldsymbol{y} \mid \boldsymbol{\varPhi}(\boldsymbol{X}) \boldsymbol{\beta}, \sigma^{2} \boldsymbol{I} \right) \mathcal{N} \left(\boldsymbol{\beta} \mid \boldsymbol{0}, \sigma_{g}^{2} \boldsymbol{I} \right)$$
$$= \mathcal{N} \left(\boldsymbol{y} \mid \boldsymbol{0}, \sigma_{g}^{2} \underbrace{\boldsymbol{\varPhi}(\boldsymbol{X}) \boldsymbol{\varPhi}(\boldsymbol{X})^{\top}}_{\boldsymbol{K}} + \sigma^{2} \boldsymbol{I} \right)$$

► K: (N × N) kernel or covariance induced by feature mapping φ.

Marginal likelihood of variance component models Application to GWAS

The missing heritability paradox

- Complex traits are regulated by a large number of small effects
 - Human height: the best single SNP explains little variance.
 - But: height of the parents are highly predictive for the height of the child!

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

Application to GWAS

Linear additive models for complex traits

Multiple linear regression model over causal SNPs

$$p(\boldsymbol{y} \mid \boldsymbol{X}, \boldsymbol{\beta}, \sigma^2) = \mathcal{N} \big(\boldsymbol{y} \mid \sum_{s \in \mathsf{causal}} \boldsymbol{x}_s \beta_s \,, \, \sigma^2 \boldsymbol{I} \big)$$

Which SNPs are causal ? Approximation: consider all S available common SNPs [Yang et al. 2011]

$$p(\boldsymbol{y} \mid \boldsymbol{X}, \boldsymbol{eta}, \sigma^2) = \mathcal{N}(\boldsymbol{y} \mid \sum_{s=1}^{S} \boldsymbol{x}_s \beta_s, \sigma^2 \boldsymbol{I})$$

- Causal SNPs either in the model or "tagged" by linkage disequilibrium to nearby common SNPs
- Uncertainty over causal SNPs: Prior on all SNP effects $p(\beta_s) = \mathcal{N}(\beta_s \mid 0, \sigma_g^2/S)$
- Marginalize out weights

$$p(\boldsymbol{y} \mid \boldsymbol{X}, \sigma_{\mathbf{g}}^2, \sigma^2) = \mathcal{N} \big(\boldsymbol{y} \mid \mathbf{0} \,, \, \sigma_{\mathbf{g}}^2 \sum_{s=1}^S \frac{1}{S} \boldsymbol{x}_s \boldsymbol{x}_s^\top + \sigma^2 \boldsymbol{I} \big)$$

▶ Perform maximum marginal likelihood estimation on σ_s^2 and σ^2 .

Application to GWAS

Linear additive models for complex traits

Multiple linear regression model over causal SNPs

$$p(\boldsymbol{y} \mid \boldsymbol{X}, \boldsymbol{\beta}, \sigma^2) = \mathcal{N} \big(\boldsymbol{y} \mid \sum_{s \in \text{causal}} \boldsymbol{x}_s \beta_s \,, \, \sigma^2 \boldsymbol{I} \big)$$

Which SNPs are causal ? Approximation: consider all S available common SNPs [Yang et al. 2011]

$$p(\boldsymbol{y} \mid \boldsymbol{X}, \boldsymbol{\beta}, \sigma^2) = \mathcal{N}(\boldsymbol{y} \mid \sum_{s=1}^{S} \boldsymbol{x}_s \beta_s, \sigma^2 \boldsymbol{I})$$

- Causal SNPs either in the model or "tagged" by linkage disequilibrium to nearby common SNPs
- Uncertainty over causal SNPs: Prior on all SNP effects $p(\beta_s) = \mathcal{N}(\beta_s \mid 0, \sigma_g^2/S)$
- Marginalize out weights

$$p(\boldsymbol{y} \mid \boldsymbol{X}, \sigma_{\rm g}^2, \sigma^2) = \mathcal{N} \big(\boldsymbol{y} \mid \boldsymbol{0} \,, \, \sigma_{\rm g}^2 \sum_{s=1}^S \frac{1}{S} \boldsymbol{x}_s \boldsymbol{x}_s^\top + \sigma^2 \boldsymbol{I} \big)$$

► Perform maximum marginal likelihood estimation on σ_z^2 and σ^2 . $\langle \Box \rangle \cdot \langle \Box \rangle$

Application to GWAS

Linear additive models for complex traits

Multiple linear regression model over causal SNPs

$$p(\boldsymbol{y} \mid \boldsymbol{X}, \boldsymbol{\beta}, \sigma^2) = \mathcal{N} \big(\boldsymbol{y} \mid \sum_{s \in \mathsf{causal}} \boldsymbol{x}_s \beta_s \,, \, \sigma^2 \boldsymbol{I} \big)$$

Which SNPs are causal ? Approximation: consider all S available common SNPs [Yang et al. 2011]

$$p(\boldsymbol{y} \mid \boldsymbol{X}, \boldsymbol{\beta}, \sigma^2) = \mathcal{N}(\boldsymbol{y} \mid \sum_{s=1}^{S} \boldsymbol{x}_s \beta_s, \sigma^2 \boldsymbol{I})$$

- Causal SNPs either in the model or "tagged" by linkage disequilibrium to nearby common SNPs
- Uncertainty over causal SNPs: Prior on all SNP effects $p(\beta_s) = \mathcal{N}\left(\beta_s \mid 0, \sigma_g^2/S\right)$
- Marginalize out weights

$$p(\boldsymbol{y} \mid \boldsymbol{X}, \sigma_{\mathsf{g}}^{2}, \sigma^{2}) = \mathcal{N} \big(\boldsymbol{y} \mid \boldsymbol{0} \,, \, \sigma_{\mathsf{g}}^{2} \sum_{s=1}^{S} \frac{1}{S} \boldsymbol{x}_{s} \boldsymbol{x}_{s}^{\top} + \sigma^{2} \boldsymbol{I} \big)$$

► Perform maximum marginal likelihood estimation on σ_g^2 and σ^2 .
Marginal likelihood of variance component models

Application to GWAS

Linear additive models for complex traits

Multiple linear regression model over causal SNPs

$$p(\boldsymbol{y} \mid \boldsymbol{X}, \boldsymbol{\beta}, \sigma^2) = \mathcal{N} \big(\boldsymbol{y} \mid \sum_{s \in \text{causal}} \boldsymbol{x}_s \beta_s \,, \, \sigma^2 \boldsymbol{I} \big)$$

Which SNPs are causal ? Approximation: consider all S available common SNPs [Yang et al. 2011]

$$p(\boldsymbol{y} \mid \boldsymbol{X}, \boldsymbol{\beta}, \sigma^2) = \mathcal{N}(\boldsymbol{y} \mid \sum_{s=1}^{S} \boldsymbol{x}_s \beta_s, \sigma^2 \boldsymbol{I})$$

- Causal SNPs either in the model or "tagged" by linkage disequilibrium to nearby common SNPs
- Uncertainty over causal SNPs: Prior on all SNP effects $p(\beta_s) = \mathcal{N}(\beta_s \mid 0, \sigma_g^2/S)$
- Marginalize out weights

$$p(\boldsymbol{y} \mid \boldsymbol{X}, \sigma_{\mathsf{g}}^{2}, \sigma^{2}) = \mathcal{N} \big(\boldsymbol{y} \mid \boldsymbol{0} \,, \, \sigma_{\mathsf{g}}^{2} \sum_{s=1}^{S} \frac{1}{S} \boldsymbol{x}_{s} \boldsymbol{x}_{s}^{\top} + \sigma^{2} \boldsymbol{I} \big)$$

► Perform maximum marginal likelihood estimation on σ_g^2 and σ^2 .

Marginal likelihood of variance component models

Application to GWAS

Linear additive models for complex traits

Multiple linear regression model over causal SNPs

$$p(\boldsymbol{y} \mid \boldsymbol{X}, \boldsymbol{\beta}, \sigma^2) = \mathcal{N} \big(\boldsymbol{y} \mid \sum_{s \in \text{causal}} \boldsymbol{x}_s \beta_s \,, \, \sigma^2 \boldsymbol{I} \big)$$

Which SNPs are causal ? Approximation: consider all S available common SNPs [Yang et al. 2011]

$$p(\boldsymbol{y} \mid \boldsymbol{X}, \boldsymbol{\beta}, \sigma^2) = \mathcal{N}(\boldsymbol{y} \mid \sum_{s=1}^{S} \boldsymbol{x}_s \beta_s, \sigma^2 \boldsymbol{I})$$

- Causal SNPs either in the model or "tagged" by linkage disequilibrium to nearby common SNPs
- Uncertainty over causal SNPs: Prior on all SNP effects $p(\beta_s) = \mathcal{N}(\beta_s \mid 0, \sigma_g^2/S)$
- Marginalize out weights

$$p(\boldsymbol{y} \mid \boldsymbol{X}, \sigma_{\mathsf{g}}^{2}, \sigma^{2}) = \mathcal{N} \big(\boldsymbol{y} \mid \boldsymbol{0}, \, \sigma_{\mathsf{g}}^{2} \sum_{s=1}^{S} \frac{1}{S} \boldsymbol{x}_{s} \boldsymbol{x}_{s}^{\top} + \sigma^{2} \boldsymbol{I} \big)$$

► Perform maximum marginal likelihood estimation on σ_g^2 and σ^2 .

Approximate variance model

$$p(\boldsymbol{y} \mid \boldsymbol{X}, \sigma_{g}^{2}, \sigma^{2}) = \mathcal{N} \left(\boldsymbol{y} \mid \boldsymbol{0}, \sigma_{g}^{2} \frac{1}{S} \boldsymbol{X} \boldsymbol{X}^{\top} + \sigma^{2} \boldsymbol{I} \right)$$

- Genetic variance σ²_g across chromosomes
- ► (Narrow-sense) heritability $h^2 = \frac{\sigma_g^2}{\sigma_g^2 + \sigma^2}$
- Narrow-sense refers to linear additive part of the heritability

Approximate variance model

$$p(oldsymbol{y} \,|\, oldsymbol{X}, \sigma_{\mathsf{g}}^2, \sigma^2) = \mathcal{N}ig(oldsymbol{y} \,|\, oldsymbol{0}, \sigma_{\mathsf{g}}^2 rac{1}{S} oldsymbol{X} oldsymbol{X}^ op + \sigma^2 oldsymbol{I}ig)$$

- Genetic variance σ²_g across chromosomes
- ► (Narrow-sense) heritability $h^2 = \frac{\sigma_g^2}{\sigma_g^2 + \sigma^2}$
- Narrow-sense refers to linear additive part of the heritability



[Yang et al. 2011]

Approximate variance model

$$p(oldsymbol{y} \,|\, oldsymbol{X}, \sigma_{\mathsf{g}}^2, \sigma^2) = \mathcal{N}ig(oldsymbol{y} \,|\, oldsymbol{0}, \sigma_{\mathsf{g}}^2 rac{1}{S} oldsymbol{X} oldsymbol{X}^ op + \sigma^2 oldsymbol{I}ig)$$

- Genetic variance σ²_g across chromosomes
- ► (Narrow-sense) heritability $h^2 = \frac{\sigma_g^2}{\sigma_g^2 + \sigma^2}$
- Narrow-sense refers to linear additive part of the heritability



・ロト ・ 理 ト ・ ヨ ト ・ ヨ ト

[Yang et al. 2011]

э

Approximate variance model

$$p(\boldsymbol{y} \mid \boldsymbol{X}, \sigma_{\mathsf{g}}^{2}, \sigma^{2}) = \mathcal{N} \left(\boldsymbol{y} \mid \boldsymbol{0}, \sigma_{\mathsf{g}}^{2} \frac{1}{S} \boldsymbol{X} \boldsymbol{X}^{\top} + \sigma^{2} \boldsymbol{I} \right)^{T}$$

- Genetic variance σ²_g across chromosomes
- ► (Narrow-sense) heritability $h^{2} = \frac{\sigma_{g}^{2}}{\sigma_{g}^{2} + \sigma^{2}}$
- Narrow-sense refers to linear additive part of the heritability



・ロト ・ 理 ト ・ ヨ ト ・ ヨ ト

[Yang et al. 2011]

э

Outline

Linear Regression II

Bayesian linear regression

Model comparison and hypothesis testing

Summary



Summary

- ► Linear models for curve fitting and multiple linear regression.
- Maximum likelihood and least squares regression are identical.
- Construction of features using a mapping ϕ .
- Regularized least squares and other models that correspond to different choices of loss functions.

- Bayesian linear regression.
- Model comparison and Occam's razor.
- Variance component models in GWAS.

Outlook

- Estimation technique for σ_{g}^{2} and σ^{2} .
- Use marginal linear model for confounder correction in GWAS testing of single SNPs
 - Linear mixed models for GWAS testing
- Use marginal linear model for testing for significant associations of sets of variants.

- Idea : Test for $\mathcal{H}_0: \sigma_g^2 = 0$ vs. $\mathcal{H}_1: \sigma_g^2 > 0$
- Random effects testing

Tasks

- Derive ridge regularized β_{MAP} in linear regression
- Derive posterior distribution (mean and covariance) of β in a linear regression under a Normal prior
- Compare them!
- ► Derive marginal likelihood for linear regression under a Normal prior on β
 - hint: The following expression is a Gaussian convolution:

$$\int \mathcal{N}(\boldsymbol{a} \mid \boldsymbol{b}, \boldsymbol{\Sigma}_{\boldsymbol{a}}) \cdot \mathcal{N}(\boldsymbol{b} \mid \boldsymbol{\mu}_{\boldsymbol{b}}, \boldsymbol{\Sigma}_{\boldsymbol{b}}) \, \mathrm{d}\boldsymbol{b}$$
$$= \int \mathcal{N}(\boldsymbol{a} - \boldsymbol{b} \mid \boldsymbol{0}, \boldsymbol{\Sigma}_{\boldsymbol{a}}) \cdot \mathcal{N}(\boldsymbol{b} \mid \boldsymbol{\mu}_{\boldsymbol{b}}, \boldsymbol{\Sigma}_{\boldsymbol{b}}) \, \mathrm{d}\boldsymbol{b}$$