



Linear models

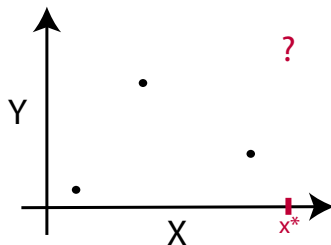
Oliver Stegle and Karsten Borgwardt

Machine Learning and
Computational Biology Research Group,
Max Planck Institute for Biological Cybernetics and
Max Planck Institute for Developmental Biology, Tübingen

Curve fitting

Tasks we are interested in:

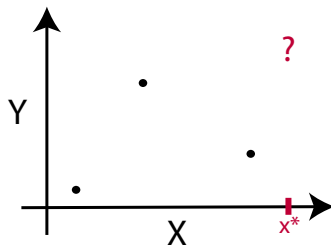
- ▶ Making predictions
- ▶ Comparison of alternative models



Curve fitting

Tasks we are interested in:

- ▶ Making predictions
- ▶ Comparison of alternative models



Further reading, useful material

- ▶ Christopher M. Bishop: Pattern Recognition and Machine learning.
 - ▶ Good background, covers most of the course material and much more!
 - ▶ This lecture is largely inspired by chapter 3 of the book.

Outline

Outline

Motivation

Linear Regression

Bayesian linear regression

Model comparison and hypothesis testing

Summary

Regression

Noise model and likelihood

- ▶ Given a dataset $\mathcal{D} = \{\mathbf{x}_n, y_n\}_{n=1}^N$, where $\mathbf{x}_n = \{x_{n,1}, \dots, x_{n,D}\}$ is D dimensional, fit parameters $\boldsymbol{\theta}$ of a regressor f with added **Gaussian noise**:

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

- ▶ Equivalent likelihood formulation:

$$p(\mathbf{y} | \mathbf{X}) = \prod_{n=1}^N \mathcal{N}(y_n | f(\mathbf{x}_n), \sigma^2)$$

Regression

Choosing a regressor

- Choose f to be **linear**:

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

- Consider bias free case, $c = 0$, otherwise include an additional column of ones in each \mathbf{x}_n .

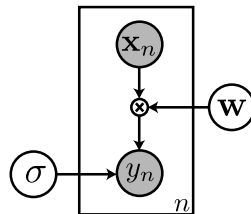
Regression

Choosing a regressor

- Choose f to be **linear**:

$$p(\mathbf{y} | \mathbf{X}) = \prod_{n=1}^N \mathcal{N}(y_n | \mathbf{w}^T \cdot \mathbf{x}_n + c, \sigma^2)$$

- Consider bias free case, $c = 0$, otherwise include an additional column of ones in each \mathbf{x}_n .



Equivalent graphical model

Linear Regression

Maximum likelihood

- ▶ Taking the logarithm, we obtain

$$\begin{aligned}\ln p(\mathbf{y} \mid \mathbf{w}, \mathbf{X}, \sigma^2) &= \sum_{n=1}^N \ln \mathcal{N}(y_n \mid \mathbf{w}^T \mathbf{x}_n, \sigma^2) \\ &= -\frac{N}{2} \ln 2\pi\sigma^2 - \frac{1}{2\sigma^2} \underbrace{\sum_{n=1}^N (y_n - \mathbf{w}^T \cdot \mathbf{x}_n)^2}_{\text{Sum of squares}}\end{aligned}$$

- ▶ The likelihood is **maximized** when the **squared error** is **minimized**.
- ▶ **Least squares** and maximum likelihood are equivalent.

Linear Regression

Maximum likelihood

- ▶ Taking the logarithm, we obtain

$$\begin{aligned}\ln p(\mathbf{y} \mid \mathbf{w}, \mathbf{X}, \sigma^2) &= \sum_{n=1}^N \ln \mathcal{N}(y_n \mid \mathbf{w}^T \mathbf{x}_n, \sigma^2) \\ &= -\frac{N}{2} \ln 2\pi\sigma^2 - \frac{1}{2\sigma^2} \underbrace{\sum_{n=1}^N (y_n - \mathbf{w}^T \cdot \mathbf{x}_n)^2}_{\text{Sum of squares}}\end{aligned}$$

- ▶ The likelihood is **maximized** when the **squared error** is **minimized**.
- ▶ **Least squares** and maximum likelihood are equivalent.

Linear Regression

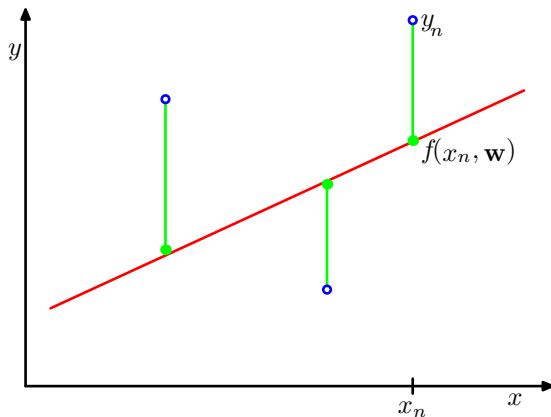
Maximum likelihood

- ▶ Taking the logarithm, we obtain

$$\begin{aligned}\ln p(\mathbf{y} \mid \mathbf{w}, \mathbf{X}, \sigma^2) &= \sum_{n=1}^N \ln \mathcal{N}(y_n \mid \mathbf{w}^T \mathbf{x}_n, \sigma^2) \\ &= -\frac{N}{2} \ln 2\pi\sigma^2 - \frac{1}{2\sigma^2} \underbrace{\sum_{n=1}^N (y_n - \mathbf{w}^T \cdot \mathbf{x}_n)^2}_{\text{Sum of squares}}\end{aligned}$$

- ▶ The likelihood is **maximized** when the **squared error** is **minimized**.
- ▶ **Least squares** and maximum likelihood are equivalent.

Linear Regression and Least Squares



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

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (y_n - \mathbf{w}^T \mathbf{x}_n)^2$$

Linear Regression and Least Squares

- Derivative w.r.t a single weight entry w_i

$$\begin{aligned}\frac{d}{dw_i} \ln p(\mathbf{y} | \mathbf{w}, \sigma^2) &= \frac{d}{dw_i} \left[-\frac{1}{2\sigma^2} \sum_{n=1}^N (y_n - \mathbf{w} \cdot \mathbf{x}_n)^2 \right] \\ &= \frac{1}{\sigma^2} \sum_{n=1}^N (y_n - \mathbf{w} \cdot \mathbf{x}_n) x_i\end{aligned}$$

- Set gradient w.r.t to \mathbf{w} to zero

$$\begin{aligned}\nabla_{\mathbf{w}} \ln p(\mathbf{y} | \mathbf{w}, \sigma^2) &= \frac{1}{\sigma^2} \sum_{n=1}^N (y_n - \mathbf{w} \cdot \mathbf{x}_n) \mathbf{x}_n^T = 0 \\ \implies \mathbf{w}_{\text{ML}} &= \underbrace{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T}_{\text{Pseudo inverse}} \mathbf{y}\end{aligned}$$

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

Polynomial Curve Fitting

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

$$\begin{aligned} f(x, \mathbf{w}) &= w_0 + w_1x + w_2x^2 + \cdots + w_Kx^K \\ &= \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}), \end{aligned}$$

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

- Similarly, $\boldsymbol{\phi}$ can be any feature mapping.
- Possible to show: the feature map $\boldsymbol{\phi}$ can be expressed in terms of kernels (kernel trick).

Polynomial Curve Fitting

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

$$\begin{aligned} f(x, \mathbf{w}) &= w_0 + w_1x + w_2x^2 + \dots + w_Kx^K \\ &= \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}), \end{aligned}$$

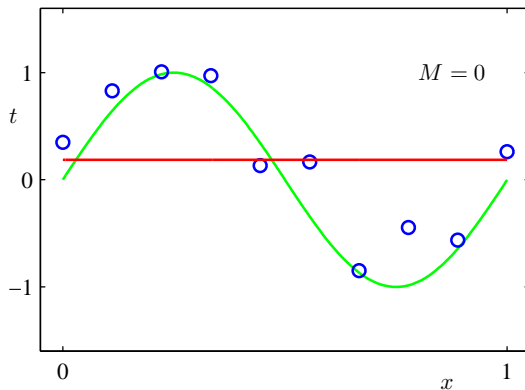
where we defined $\boldsymbol{\phi}(\mathbf{x}) = (1, x, x^2, \dots, x^K)$.

- Similarly, $\boldsymbol{\phi}$ can be any feature mapping.
- Possible to show: the feature map $\boldsymbol{\phi}$ can be expressed in terms of kernels (kernel trick).

Polynomial Curve Fitting

Overfitting

- The degree of the polynomial is crucial to avoid under- and overfitting.

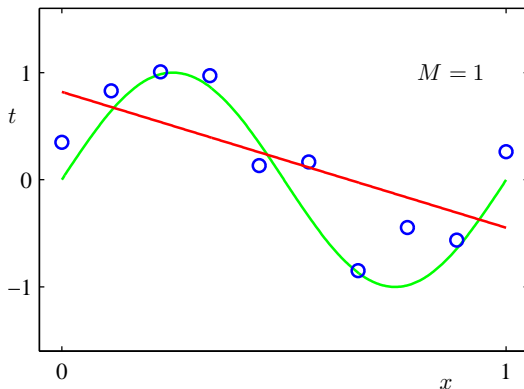


(C.M. Bishop, Pattern Recognition and Machine Learning)

Polynomial Curve Fitting

Overfitting

- The degree of the polynomial is crucial to avoid under- and overfitting.

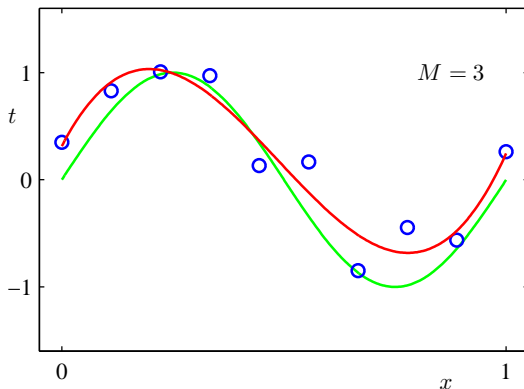


(C.M. Bishop, Pattern Recognition and Machine Learning)

Polynomial Curve Fitting

Overfitting

- The degree of the polynomial is crucial to avoid under- and overfitting.

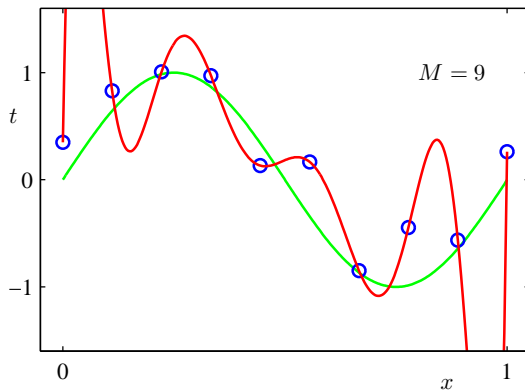


(C.M. Bishop, Pattern Recognition and Machine Learning)

Polynomial Curve Fitting

Overfitting

- The degree of the polynomial is crucial to avoid under- and overfitting.



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

Regularized Least Squares

- ▶ Solutions to avoid overfitting:
 - ▶ Intelligently choose K
 - ▶ Regularize the regression weights \mathbf{w}
- ▶ Construct a smoothed error function

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

Regularized Least Squares

- ▶ Solutions to avoid overfitting:
 - ▶ Intelligently choose K
 - ▶ Regularize the regression weights \mathbf{w}
- ▶ Construct a smoothed error function

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

Regularized Least Squares

More general regularizers

- A more general regularization approach:

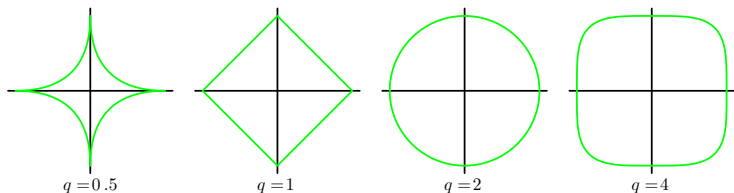
$$E(\mathbf{w}) = \underbrace{\frac{1}{2} \sum_{n=1}^N (y_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2}_{\text{Squared error}} + \underbrace{\frac{\lambda}{2} \sum_{d=1}^D |w_d|^q}_{\text{Regularizer}}$$

Regularized Least Squares

More general regularizers

- A more general regularization approach:

$$E(\mathbf{w}) = \underbrace{\frac{1}{2} \sum_{n=1}^N (y_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2}_{\text{Squared error}} + \underbrace{\frac{\lambda}{2} \sum_{d=1}^D |w_d|^q}_{\text{Regularizer}}$$



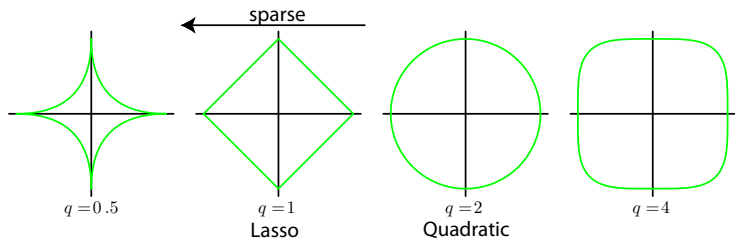
(C.M. Bishop, Pattern Recognition and Machine Learning)

Regularized Least Squares

More general regularizers

- A more general regularization approach:

$$E(\mathbf{w}) = \underbrace{\frac{1}{2} \sum_{n=1}^N (y_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2}_{\text{Squared error}} + \underbrace{\frac{\lambda}{2} \sum_{d=1}^D |w_d|^q}_{\text{Regularizer}}$$



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

Loss functions and other methods

- ▶ Even more general: vary the loss function

$$E(\mathbf{w}) = \underbrace{\frac{1}{2} \sum_{n=1}^N \mathcal{L}(y_n - \mathbf{w}^T \phi(\mathbf{x}_n))}_{\text{Loss}} + \underbrace{\frac{\lambda}{2} \sum_{d=1}^D |w_d|^q}_{\text{Regularizer}}$$

- ▶ 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 $E(\mathbf{w})$, yielding a point estimate for \mathbf{w} .

Loss functions and other methods

- ▶ Even more general: vary the loss function

$$E(\mathbf{w}) = \underbrace{\frac{1}{2} \sum_{n=1}^N \mathcal{L}(y_n - \mathbf{w}^T \phi(\mathbf{x}_n))}_{\text{Loss}} + \underbrace{\frac{\lambda}{2} \sum_{d=1}^D |w_d|^q}_{\text{Regularizer}}$$

- ▶ 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 $E(\mathbf{w})$, yielding a point estimate for \mathbf{w} .

Loss functions and other methods

- ▶ Even more general: vary the loss function

$$E(\mathbf{w}) = \underbrace{\frac{1}{2} \sum_{n=1}^N \mathcal{L}(y_n - \mathbf{w}^T \phi(\mathbf{x}_n))}_{\text{Loss}} + \underbrace{\frac{\lambda}{2} \sum_{d=1}^D |w_d|^q}_{\text{Regularizer}}$$

- ▶ 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 $E(\mathbf{w})$, yielding a point estimate for \mathbf{w} .

Regularized Least Squares

Probabilistic equivalent

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

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

- Similarly: most other choices of regularizers and loss functions can be mapped to an equivalent probabilistic representation.

Regularized Least Squares

Probabilistic equivalent

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

$$\begin{aligned}
 E(\mathbf{w}) &= \underbrace{\frac{1}{2} \sum_{n=1}^N (y_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2}_{\text{Squared error}} & + \underbrace{\frac{\lambda}{2} \mathbf{w}^T \mathbf{w}}_{\text{Regularizer}} \\
 &= -\ln p(\mathbf{y} \mid \mathbf{w}, \Phi(\mathbf{X}), \sigma^2) & -\ln p(\mathbf{w})
 \end{aligned}$$

- Similarly: most other choices of regularizers and loss functions can be mapped to an equivalent probabilistic representation.

Regularized Least Squares

Probabilistic equivalent

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

$$\begin{aligned}
 E(\mathbf{w}) &= \underbrace{\frac{1}{2} \sum_{n=1}^N (y_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2}_{\text{Squared error}} && + \underbrace{\frac{\lambda}{2} \mathbf{w}^T \mathbf{w}}_{\text{Regularizer}} \\
 &= -\ln p(\mathbf{y} \mid \mathbf{w}, \Phi(\mathbf{X}), \sigma^2) && -\ln p(\mathbf{w}) \\
 &= -\sum_{n=1}^N \ln \mathcal{N}(y_n \mid \mathbf{w}^T \phi(\mathbf{x}_n), \sigma^2) && -\ln \mathcal{N}\left(\mathbf{w} \mid \mathbf{0}, \frac{1}{\lambda} \mathbf{I}\right)
 \end{aligned}$$

- Similarly: most other choices of regularizers and loss functions can be mapped to an equivalent probabilistic representation.

Regularized Least Squares

Probabilistic equivalent

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

$$\begin{aligned}
 E(\mathbf{w}) &= \underbrace{\frac{1}{2} \sum_{n=1}^N (y_n - \mathbf{w}^T \phi(\mathbf{x}_n))^2}_{\text{Squared error}} + \underbrace{\frac{\lambda}{2} \mathbf{w}^T \mathbf{w}}_{\text{Regularizer}} \\
 &= -\ln p(\mathbf{y} \mid \mathbf{w}, \Phi(\mathbf{X}), \sigma^2) - \ln p(\mathbf{w}) \\
 &= -\sum_{n=1}^N \ln \mathcal{N}(y_n \mid \mathbf{w}^T \phi(\mathbf{x}_n), \sigma^2) - \ln \mathcal{N}\left(\mathbf{w} \mid \mathbf{0}, \frac{1}{\lambda} \mathbf{I}\right)
 \end{aligned}$$

- Similarly: most other choices of regularizers and loss functions can be mapped to an equivalent probabilistic representation.

Outline

Motivation

Linear Regression

Bayesian linear regression

Model comparison and hypothesis testing

Summary

Bayesian linear regression

- Likelihood as before

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

- Define a conjugate prior over \mathbf{w}

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} \mid \mathbf{m}_0, \mathbf{S}_0)$$

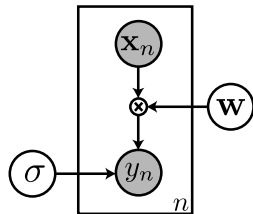
Bayesian linear regression

- Likelihood as before

$$p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \sigma^2) = \prod_{n=1}^N \mathcal{N}(y_n | \mathbf{w}^T \cdot \phi(\mathbf{x}_n), \sigma^2)$$

- Define a conjugate prior over \mathbf{w}

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0)$$



Bayesian linear regression

- Posterior probability of \mathbf{w}

$$\begin{aligned} p(\mathbf{w} \mid \mathbf{y}, \mathbf{X}, \sigma^2) &\propto \prod_{n=1}^N \mathcal{N}(y_n \mid \mathbf{w}^T \cdot \phi(\mathbf{x}_n), \sigma^2) \cdot \mathcal{N}(\mathbf{w} \mid \mathbf{m}_0, \mathbf{S}_0) \\ &= \mathcal{N}(\mathbf{y} \mid \mathbf{w} \cdot \Phi(\mathbf{X}), \sigma^2 \mathbf{I}) \cdot \mathcal{N}(\mathbf{w} \mid \mathbf{m}_0, \mathbf{S}_0) \\ &= \mathcal{N}(\mathbf{w} \mid \boldsymbol{\mu}_{\mathbf{w}}, \boldsymbol{\Sigma}_{\mathbf{w}}) \end{aligned}$$

- where

$$\begin{aligned} \boldsymbol{\mu}_{\mathbf{w}} &= \boldsymbol{\Sigma}_{\mathbf{w}} \left(\mathbf{S}_0^{-1} \mathbf{m}_0 + \frac{1}{\sigma^2} \Phi(\mathbf{X})^T \mathbf{y} \right) \\ \boldsymbol{\Sigma}_{\mathbf{w}} &= \left[\mathbf{S}_0^{-1} + \frac{1}{\sigma^2} \Phi(\mathbf{X})^T \Phi(\mathbf{X}) \right]^{-1} \end{aligned}$$

Bayesian linear regression

Prior choice

- ▶ A common choice is a prior that corresponds to regularized regression

$$p(\mathbf{w}) = \mathcal{N}\left(\mathbf{w} \mid \mathbf{0}, \frac{1}{\lambda} \mathbf{I}\right).$$

- ▶ In this case

$$\boldsymbol{\mu}_{\mathbf{w}} = \boldsymbol{\Sigma}_{\mathbf{w}} \left(\mathbf{S}_0^{-1} \mathbf{m}_0 + \frac{1}{\sigma^2} \boldsymbol{\Phi}(\mathbf{X})^T \mathbf{y} \right)$$

$$\boldsymbol{\Sigma}_{\mathbf{w}} = \left[\mathbf{S}_0^{-1} + \frac{1}{\sigma^2} \boldsymbol{\Phi}(\mathbf{X})^T \boldsymbol{\Phi}(\mathbf{X}) \right]^{-1}$$

Bayesian linear regression

Prior choice

- ▶ A common choice is a prior that corresponds to regularized regression

$$p(\mathbf{w}) = \mathcal{N}\left(\mathbf{w} \mid \mathbf{0}, \frac{1}{\lambda} \mathbf{I}\right).$$

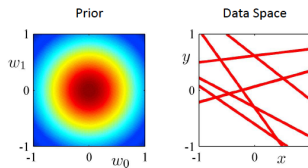
- ▶ In this case

$$\begin{aligned}\boldsymbol{\mu}_{\mathbf{w}} &= \boldsymbol{\Sigma}_{\mathbf{w}} \left(\frac{1}{\sigma^2} \boldsymbol{\Phi}(\mathbf{X})^T \mathbf{y} \right) \\ \boldsymbol{\Sigma}_{\mathbf{w}} &= \left[\lambda \mathbf{I} + \frac{1}{\sigma^2} \boldsymbol{\Phi}(\mathbf{X})^T \boldsymbol{\Phi}(\mathbf{X}) \right]^{-1}\end{aligned}$$

Bayesian linear regression

Example

0 Data points

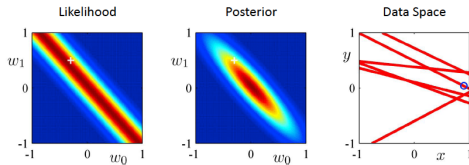


(C.M. Bishop, Pattern Recognition and Machine Learning)

Bayesian linear regression

Example

1 Data point

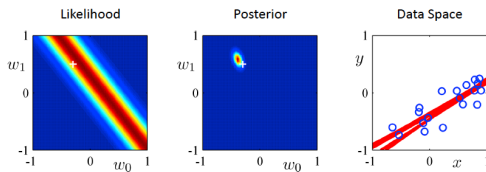


(C.M. Bishop, Pattern Recognition and Machine Learning)

Bayesian linear regression

Example

20 Data points



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

Making predictions

- Prediction for fixed weight $\hat{\mathbf{w}}$ at input \mathbf{x}^* trivial:

$$p(y^* | \mathbf{x}^*, \hat{\mathbf{w}}, \sigma^2) = \mathcal{N}(y^* | \hat{\mathbf{w}}^T \phi(\mathbf{x}^*), \sigma^2)$$

- Integrate over \mathbf{w} to take the posterior uncertainty into account

$$\begin{aligned} p(y^* | \mathbf{x}^*, \mathcal{D}) &= \int_{\mathbf{w}} p(y^* | \mathbf{x}^*, \mathbf{w}, \sigma^2) p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \sigma^2) \\ &= \int_{\mathbf{w}} \mathcal{N}(y^* | \mathbf{w}^T \phi(\mathbf{x}^*), \sigma^2) \mathcal{N}(\mathbf{w} | \boldsymbol{\mu}_{\mathbf{w}}, \boldsymbol{\Sigma}_{\mathbf{w}}) \\ &= \mathcal{N}(y^* | \boldsymbol{\mu}_{\mathbf{w}}^T \phi(\mathbf{x}^*), \sigma^2 + \phi(\mathbf{x}^*)^T \boldsymbol{\Sigma}_{\mathbf{w}} \phi(\mathbf{x}^*)) \end{aligned}$$

- Key:
 - prediction is again Gaussian
 - Predictive variance is increase due to the posterior uncertainty in \mathbf{w} .

Making predictions

- Prediction for fixed weight $\hat{\mathbf{w}}$ at input \mathbf{x}^* trivial:

$$p(y^* | \mathbf{x}^*, \hat{\mathbf{w}}, \sigma^2) = \mathcal{N}(y^* | \hat{\mathbf{w}}^T \phi(\mathbf{x}^*), \sigma^2)$$

- Integrate over \mathbf{w} to take the posterior uncertainty into account

$$\begin{aligned} p(y^* | \mathbf{x}^*, \mathcal{D}) &= \int_{\mathbf{w}} p(y^* | \mathbf{x}^*, \mathbf{w}, \sigma^2) p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \sigma^2) \\ &= \int_{\mathbf{w}} \mathcal{N}(y^* | \mathbf{w}^T \phi(\mathbf{x}^*), \sigma^2) \mathcal{N}(\mathbf{w} | \boldsymbol{\mu}_{\mathbf{w}}, \boldsymbol{\Sigma}_{\mathbf{w}}) \\ &= \mathcal{N}(y^* | \boldsymbol{\mu}_{\mathbf{w}}^T \phi(\mathbf{x}^*), \sigma^2 + \phi(\mathbf{x}^*)^T \boldsymbol{\Sigma}_{\mathbf{w}} \phi(\mathbf{x}^*)) \end{aligned}$$

- Key:

- prediction is again Gaussian
- Predictive variance is increase due to the posterior uncertainty in \mathbf{w} .

Making predictions

- Prediction for fixed weight $\hat{\mathbf{w}}$ at input \mathbf{x}^* trivial:

$$p(y^* | \mathbf{x}^*, \hat{\mathbf{w}}, \sigma^2) = \mathcal{N}(y^* | \hat{\mathbf{w}}^T \phi(\mathbf{x}^*), \sigma^2)$$

- Integrate over \mathbf{w} to take the posterior uncertainty into account

$$\begin{aligned} p(y^* | \mathbf{x}^*, \mathcal{D}) &= \int_{\mathbf{w}} p(y^* | \mathbf{x}^*, \mathbf{w}, \sigma^2) p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \sigma^2) \\ &= \int_{\mathbf{w}} \mathcal{N}(y^* | \mathbf{w}^T \phi(\mathbf{x}^*), \sigma^2) \mathcal{N}(\mathbf{w} | \boldsymbol{\mu}_{\mathbf{w}}, \boldsymbol{\Sigma}_{\mathbf{w}}) \\ &= \mathcal{N}(y^* | \boldsymbol{\mu}_{\mathbf{w}}^T \phi(\mathbf{x}^*), \sigma^2 + \phi(\mathbf{x}^*)^T \boldsymbol{\Sigma}_{\mathbf{w}} \phi(\mathbf{x}^*)) \end{aligned}$$

- Key:
 - prediction is again Gaussian
 - Predictive variance is increase due to the posterior uncertainty in \mathbf{w} .

Outline

Motivation

Linear Regression

Bayesian linear regression

Model comparison and hypothesis testing

Summary

Model comparison

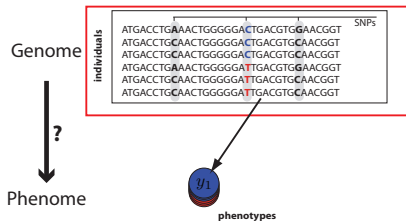
Motivation

- ▶ What degree of polynomials describes the data best?
- ▶ 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.



Bayesian model comparison

- ▶ How do we choose among alternative models?
- ▶ Assume we want to choose among models $\mathcal{H}_0, \dots, \mathcal{H}_M$ for a dataset \mathcal{D} .
- ▶ Posterior probability for a particular model i

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

Bayesian model comparison

- ▶ How do we choose among alternative models?
- ▶ Assume we want to choose among models $\mathcal{H}_0, \dots, \mathcal{H}_M$ for a dataset \mathcal{D} .
- ▶ Posterior probability for a particular model i

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

Bayesian model comparison

How to calculate the evidence

- ▶ The evidence is not the model likelihood!

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

- ▶ Remember:

$$p(\boldsymbol{\theta} | \mathcal{H}_i, \mathcal{D}) = \frac{p(\mathcal{D} | \mathcal{H}_i, \boldsymbol{\theta}) p(\boldsymbol{\theta})}{p(\mathcal{D} | \mathcal{H}_i)}$$

Bayesian model comparison

How to calculate the evidence

- ▶ The evidence is not the model likelihood!

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

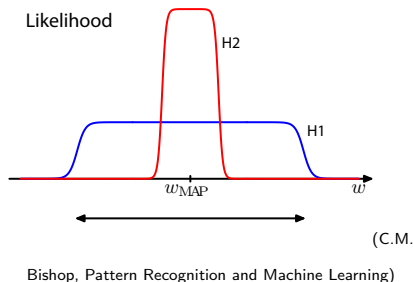
- ▶ Remember:

$$p(\boldsymbol{\theta} | \mathcal{H}_i, \mathcal{D}) = \frac{p(\mathcal{D} | \mathcal{H}_i, \boldsymbol{\theta}) p(\boldsymbol{\theta})}{p(\mathcal{D} | \mathcal{H}_i)}$$
$$\text{posterior} = \frac{\text{likelihood} \cdot \text{prior}}{\text{Evidence}}$$

Bayesian model comparison

Ocam's razor

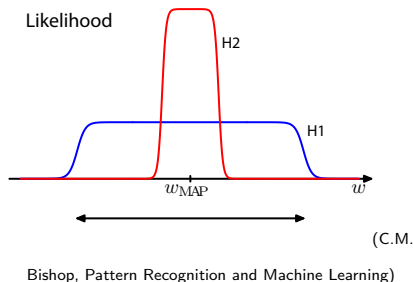
- The evidence integral penalizes **overly complex models**.
- A model with few parameters and lower maximum likelihood (\mathcal{H}_1) may win over a model with a peaked likelihood that requires many more parameters (\mathcal{H}_2).



Bayesian model comparison

Ocam's razor

- ▶ The evidence integral penalizes **overly complex models**.
- ▶ A model with few parameters and lower maximum likelihood (\mathcal{H}_1) may win over a model with a peaked likelihood that requires many more parameters (\mathcal{H}_2).



Application to GWA

- ▶ Consider an association study.
 - ▶ $\mathcal{H}_0: p(\mathbf{y} | \mathcal{H}_0, \mathbf{X}, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{y} | \mathbf{0}, \sigma^2 \mathbf{I})$ (no association)
 $\boldsymbol{\theta} = \{\sigma^2\}$
 - ▶ $\mathcal{H}_1: p(\mathbf{y} | \mathcal{H}_1, \mathbf{X}, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{y} | \mathbf{w}^T \cdot \mathbf{X}, \sigma^2 \mathbf{I})$ (linear association)
 $\boldsymbol{\theta} = \{\sigma^2, \mathbf{w}\}$
- ▶ Choosing conjugate priors for σ^2 and \mathbf{w} , the required integrals are tractable in closed form.

Application to GWA

- ▶ Consider an association study.
 - ▶ $\mathcal{H}_0: p(\mathbf{y} | \mathcal{H}_0, \mathbf{X}, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{y} | \mathbf{0}, \sigma^2 \mathbf{I})$ (no association)
 $\boldsymbol{\theta} = \{\sigma^2\}$
 - ▶ $\mathcal{H}_1: p(\mathbf{y} | \mathcal{H}_1, \mathbf{X}, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{y} | \mathbf{w}^T \cdot \mathbf{X}, \sigma^2 \mathbf{I})$ (linear association)
 $\boldsymbol{\theta} = \{\sigma^2, \mathbf{w}\}$
- ▶ Choosing conjugate priors for σ^2 and \mathbf{w} , the required integrals are tractable in closed form.

Application to GWA

- ▶ Consider an association study.
 - ▶ \mathcal{H}_0 : $p(\mathbf{y} \mid \mathcal{H}_0, \mathbf{X}, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{y} \mid \mathbf{0}, \sigma^2 \mathbf{I})$ (no association)
 $\boldsymbol{\theta} = \{\sigma^2\}$
 - ▶ \mathcal{H}_1 : $p(\mathbf{y} \mid \mathcal{H}_1, \mathbf{X}, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{y} \mid \mathbf{w}^T \cdot \mathbf{X}, \sigma^2 \mathbf{I})$ (linear association)
 $\boldsymbol{\theta} = \{\sigma^2, \mathbf{w}\}$
- ▶ Choosing conjugate priors for σ^2 and \mathbf{w} , the required integrals are tractable in closed form.

Application to GWA

Scoring models

- The ratio of the evidences, the **Bayes factor** is a common scoring metric to compare two models:

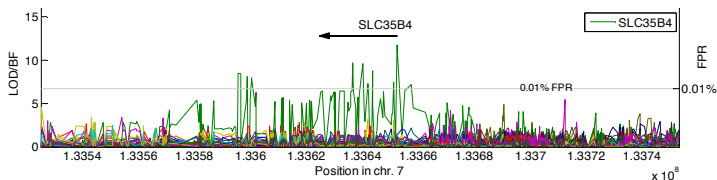
$$BF = \ln \frac{p(\mathcal{D} | \mathcal{H}_1)}{p(\mathcal{D} | \mathcal{H}_0)}.$$

Application to GWA

Scoring models

- The ratio of the evidences, the **Bayes factor** is a common scoring metric to compare two models:

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



Application to GWA

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

$$\begin{aligned} p(\mathcal{H}_1 | \mathcal{D}) &= \frac{p(\mathcal{D} | \mathcal{H}_1)p(\mathcal{H}_1)}{p(\mathcal{D})} \\ &= \frac{p(\mathcal{D} | \mathcal{H}_1)p(\mathcal{H}_1)}{p(\mathcal{D} | \mathcal{H}_1)p(\mathcal{H}_1) + p(\mathcal{D} | \mathcal{H}_0)p(\mathcal{H}_0)} \end{aligned}$$

- $p(\mathcal{H}_1 | \mathcal{D}) + p(\mathcal{H}_0 | \mathcal{D}) = 1$, prior probability of observing a real association.

Application to GWA

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

$$\begin{aligned} p(\mathcal{H}_1 | \mathcal{D}) &= \frac{p(\mathcal{D} | \mathcal{H}_1)p(\mathcal{H}_1)}{p(\mathcal{D})} \\ &= \frac{p(\mathcal{D} | \mathcal{H}_1)p(\mathcal{H}_1)}{p(\mathcal{D} | \mathcal{H}_1)p(\mathcal{H}_1) + p(\mathcal{D} | \mathcal{H}_0)p(\mathcal{H}_0)} \end{aligned}$$

- ▶ $p(\mathcal{H}_1 | \mathcal{D}) + p(\mathcal{H}_0 | \mathcal{D}) = 1$, prior probability of observing a real association.

Application to GWA

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

$$\begin{aligned} p(\mathcal{H}_1 | \mathcal{D}) &= \frac{p(\mathcal{D} | \mathcal{H}_1)p(\mathcal{H}_1)}{p(\mathcal{D})} \\ &= \frac{p(\mathcal{D} | \mathcal{H}_1)p(\mathcal{H}_1)}{p(\mathcal{D} | \mathcal{H}_1)p(\mathcal{H}_1) + p(\mathcal{D} | \mathcal{H}_0)p(\mathcal{H}_0)} \end{aligned}$$

- ▶ $p(\mathcal{H}_1 | \mathcal{D}) + p(\mathcal{H}_0 | \mathcal{D}) = 1$, prior probability of observing a real association.

Outline

Motivation

Linear Regression

Bayesian linear regression

Model comparison and hypothesis testing

Summary

Summary

- ▶ Curve fitting and linear regression.
- ▶ Maximum likelihood and least squares regression are identical.
- ▶ Construction of features using a mapping ϕ .
- ▶ Regularized least squares.
- ▶ Bayesian linear regression.
- ▶ Model comparison and *ocam's razor*.