

Statistical Machine Learning

Lecture 08: Regression

Kristian Kersting

TU Darmstadt

Summer Term 2020

Today's Objectives

- Make you understand how to learn a continuous function
- Covered Topics
 - Linear Regression and its interpretations
 - What is overfitting?
 - Deriving Linear Regression from Maximum Likelihood Estimation
 - Bayesian Linear Regression

Outline

- 1. Introduction to Linear Regression**
- 2. Maximum Likelihood Approach to Regression**
- 3. Bayesian Linear Regression**
- 4. Wrap-Up**

Outline

1. Introduction to Linear Regression

2. Maximum Likelihood Approach to Regression

3. Bayesian Linear Regression

4. Wrap-Up

Reminder

- Our task is to **learn a mapping f from input to output**

$$f : I \rightarrow O, \quad y = f(x; \theta)$$

- Input: $x \in I$ (images, text, sensor measurements, ...)
- Output: $y \in O$
- Parameters: $\theta \in \Theta$ (what needs to be “learned”)
- **Regression**
 - Learn a mapping into a **continuous space**
 - $O = \mathbb{R}$
 - $O = \mathbb{R}^3$
 - ...

Motivation

- You want to predict the torques of a robot arm

$$\begin{aligned}y &= I\ddot{q} - \mu\dot{q} + mlg \sin(q) \\ &= \begin{bmatrix} \ddot{q} & \dot{q} & \sin(q) \end{bmatrix} \begin{bmatrix} I & -\mu & mlg \end{bmatrix}^\top \\ &= \phi(\mathbf{x})^\top \boldsymbol{\theta}\end{aligned}$$



- Can we do this with a data set?

$$\mathcal{D} = \left\{ (\mathbf{x}_i, y_i) \mid i = 1 \cdots n \right\}$$

- A **linear regression problem!**

Least Squares Linear Regression

- We are given pairs of training data points and associated function values (\mathbf{x}_i, y_i)

$$X = \{\mathbf{x}_1 \in \mathbb{R}^d, \dots, \mathbf{x}_n\}$$

$$Y = \{y_1 \in \mathbb{R}, \dots, y_n\}$$

- Note: here we only do the case $y_i \in \mathbb{R}$. In general y_i can have more than one dimension, i.e., $y_i \in \mathbb{R}^f$ for some positive f
- Start with linear regressor

$$\mathbf{x}_i^T \mathbf{w} + w_0 = y_i \quad \forall i = 1, \dots, n$$

- One linear equation for each training data point/label pair
- Exactly the same basic setup as for least-squares classification!
Only the values are **continuous**

Least Squares Linear Regression

$$\mathbf{x}_i^T \mathbf{w} + w_0 = y_i \quad \forall i = 1, \dots, n$$

- **Step 1:** Define

$$\hat{\mathbf{x}}_i = \begin{pmatrix} \mathbf{x}_i \\ 1 \end{pmatrix} \quad \hat{\mathbf{w}} = \begin{pmatrix} \mathbf{w} \\ w_0 \end{pmatrix}$$

- **Step 2:** Rewrite

$$\hat{\mathbf{x}}_i^T \hat{\mathbf{w}} = y_i \quad \forall i = 1, \dots, n$$

- **Step 3:** Matrix-vector notation

$$\hat{\mathbf{X}}^T \hat{\mathbf{w}} = \mathbf{y}$$

- where $\hat{\mathbf{X}} = [\hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_n]$ (each $\hat{\mathbf{x}}_i$ is a vector) and $\mathbf{y} = [y_1, \dots, y_n]^T$

Least Squares Linear Regression

- **Step 4:** Find the least squares solution

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \left\| \hat{\mathbf{X}}^T \mathbf{w} - \mathbf{y} \right\|^2$$

$$\nabla_{\mathbf{w}} \left\| \hat{\mathbf{X}}^T \mathbf{w} - \mathbf{y} \right\|^2 = \mathbf{0}$$

$$\hat{\mathbf{w}} = \left(\hat{\mathbf{X}} \hat{\mathbf{X}}^T \right)^{-1} \hat{\mathbf{X}} \mathbf{y}$$

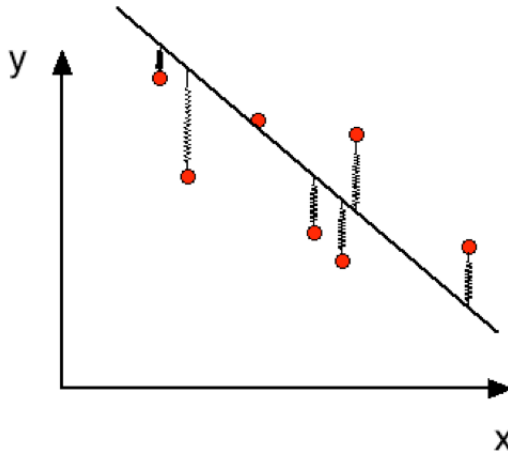
- A closed form solution!

Least Squares Linear Regression

$$\hat{\mathbf{w}} = \left(\hat{\mathbf{X}} \hat{\mathbf{X}}^\top \right)^{-1} \hat{\mathbf{X}} \mathbf{y}$$

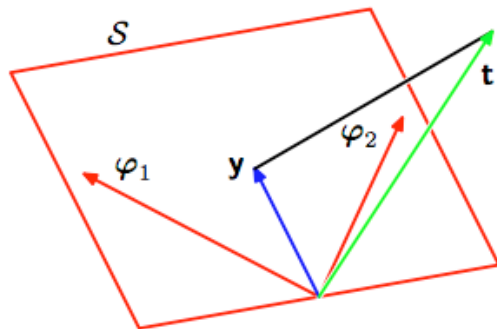
- Where is the costly part of this computation?
 - The inverse is a $\mathbb{R}^{D \times D}$ matrix
 - Naive inversion takes $O(D^3)$, but better methods exist
- What can we do if the input dimension D is too large?
 - Gradient descent
 - Work with fewer dimensions

Mechanical Interpretation



Geometric Interpretation

- Predicted outputs are Linear Combinations of Features! Samples are projected in this **Feature Space**



Polynomial Regression

- How can we **fit arbitrary polynomials** using least-squares regression?
 - We introduce a feature transformation as before

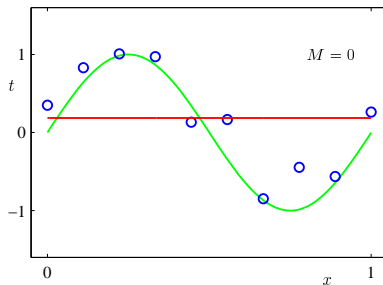
$$\begin{aligned}y(\mathbf{x}) &= \mathbf{w}^T \phi(\mathbf{x}) \\ &= \sum_{i=0}^M w_i \phi_i(\mathbf{x})\end{aligned}$$

- Assume $\phi_0(\mathbf{x}) = 1$
 - $\phi_i(\cdot)$ are called the **basis functions**
 - **Still a linear model in the parameters \mathbf{w}**
- E.g. fitting a cubic polynomial

$$\phi(x) = \left(1, x, x^2, x^3\right)^T$$

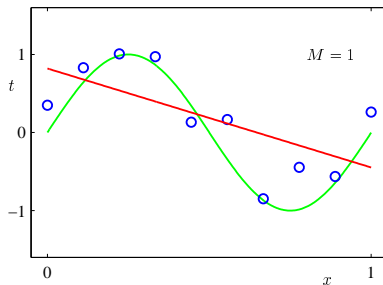
Polynomial Regression

■ Polynomial of degree 0 (constant value)



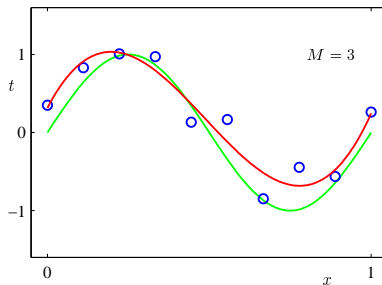
Polynomial Regression

■ Polynomial of degree 1 (line)



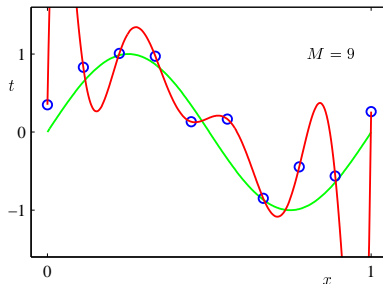
Polynomial Regression

■ Polynomial of degree 3 (cubic)



Polynomial Regression

■ Polynomial of degree 9



■ Massive overfitting

Outline

1. Introduction to Linear Regression

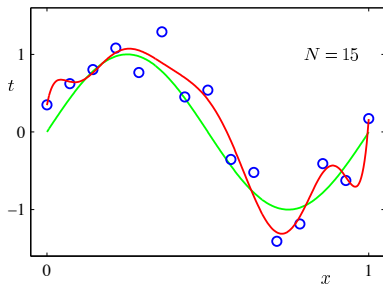
2. Maximum Likelihood Approach to Regression

3. Bayesian Linear Regression

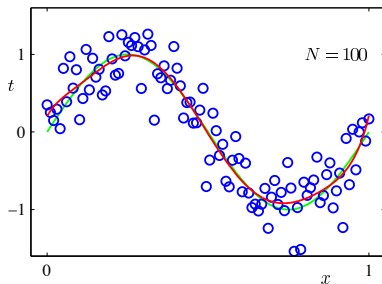
4. Wrap-Up

Overfitting

Relatively little data leads to
overfitting



Enough data leads to a **good estimate**



Probabilistic Regression

- **Assumption 1:** Our target function values are generated by adding noise to the function estimate

$$y = f(\mathbf{x}, \mathbf{w}) + \epsilon$$

- y - target function value; f - regression function; \mathbf{x} - input value; \mathbf{w} - weights or parameters; ϵ - noise

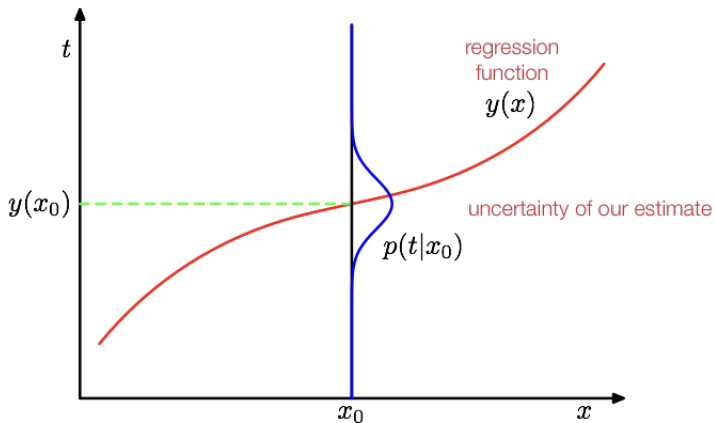
- **Assumption 2:** The noise is a random variable that is Gaussian distributed

$$\epsilon \sim \mathcal{N}(0, \beta^{-1})$$

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

- $f(\mathbf{x}, \mathbf{w})$ is the mean; β^{-1} is the variance (β is the precision)
- Note that y is now a random variable with underlying probability distribution $p(y \mid \mathbf{x}, \mathbf{w}, \beta)$

Probabilistic Regression



Probabilistic Regression

■ Given

- Training input data points $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}$
- Associated function values $\mathbf{Y} = [y_1, \dots, y_n]^T$

■ Conditional likelihood (assuming the data is i.i.d.)

$$p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta) = \prod_{i=1}^n \mathcal{N}(y_i \mid f(\mathbf{x}_i, \mathbf{w}), \beta^{-1})$$

(with linear model)

$$= \prod_{i=1}^n \mathcal{N}(y_i \mid \mathbf{w}^T \phi(\mathbf{x}_i), \beta^{-1})$$

- $\mathbf{w}^T \phi(\mathbf{x}_i)$ is the **generalized linear regression function**
- Maximize the likelihood w.r.t. (with respect to) \mathbf{w} and β

Maximum Likelihood Regression

- Simplify using the **log**-likelihood

$$\begin{aligned}\log p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta) &= \sum_{i=1}^n \log \mathcal{N}(y_i \mid \mathbf{w}^\top \phi(\mathbf{x}_i), \beta^{-1}) \\ &= \sum_{i=1}^n \left[\log \left(\frac{\sqrt{\beta}}{\sqrt{2\pi}} \right) - \frac{\beta}{2} (y_i - \mathbf{w}^\top \phi(\mathbf{x}_i))^2 \right] \\ &= \frac{n}{2} \log \beta - \frac{n}{2} \log(2\pi) - \frac{\beta}{2} \sum_{i=1}^n (y_i - \mathbf{w}^\top \phi(\mathbf{x}_i))^2\end{aligned}$$

Maximum Likelihood Regression

■ Gradient w.r.t. \mathbf{w}

$$\begin{aligned}\nabla_{\mathbf{w}} \log p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta) &= \mathbf{0} \\ -\beta \sum_{i=1}^n (y_i - \mathbf{w}^T \phi(\mathbf{x}_i)) \phi(\mathbf{x}_i) &= \mathbf{0}\end{aligned}$$

■ Define

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \quad \mathbf{w} = \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix}, \quad \Phi = \begin{bmatrix} \phi(\mathbf{x}_1) & \dots & \phi(\mathbf{x}_n) \end{bmatrix}$$

Maximum Likelihood Regression

$$\begin{aligned}\sum_{i=1}^n y_i \phi(\mathbf{x}_i) &= \left[\sum_{i=1}^n \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^\top \right] \mathbf{w} \\ \Phi \mathbf{y} &= \Phi \Phi^\top \mathbf{w} \\ \mathbf{w}_{\text{ML}} &= (\Phi \Phi^\top)^{-1} \Phi \mathbf{y}\end{aligned}$$

- The same result as in least squares regression!

Maximum Likelihood Regression

- We obtain the same \mathbf{w} as with least squares regression
 - Least-squares is equivalent to assuming the targets are Gaussian distributed
 - Note: The least squares method is not distribution-free!
- However, the Maximum Likelihood approach is much more powerful!
 - We can also estimate β

$$\beta_{\text{ML}} = \left(\frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{w}_{\text{ML}}^T \phi(\mathbf{x}_i))^2 \right)^{-1}$$

- We can gauge the uncertainty of our estimate!

Loss Functions in Regression

- Given a new data point \mathbf{x}_t , in least squares regression the function value is $y_t = \hat{\mathbf{x}}_t^\top \hat{\mathbf{w}}$
- But in maximum likelihood regression we have a probability distribution over the function value $p(y \mid \mathbf{x}, \mathbf{w}, \beta)$
- How do we actually estimate a function value y_t for a new data point \mathbf{x}_t ?
- We need a **loss function**, just as in the classification case

$$\begin{aligned} L : \mathbb{R} \times \mathbb{R} &\rightarrow \mathbb{R}^+ \\ (y_t, f(\mathbf{x}_t)) &\rightarrow L(y_t, f(\mathbf{x}_t)) \end{aligned}$$

Loss Functions in Regression

- Minimize the **expected** loss

$$\mathbb{E}_{\mathbf{x}, y \sim p(\mathbf{x}, y)} [L] = \int \int L(y, f(\mathbf{x})) p(\mathbf{x}, y) d\mathbf{x} dy$$

- Simplest case: **squared loss**

$$L(y, f(\mathbf{x})) = (y - f(\mathbf{x}))^2$$

$$\mathbb{E}_{\mathbf{x}, y \sim p(\mathbf{x}, y)} [L] = \int \int (y - f(\mathbf{x}))^2 p(\mathbf{x}, y) d\mathbf{x} dy$$

$$\frac{\partial \mathbb{E}[L]}{\partial f(\mathbf{x})} = -2 \int (y - f(\mathbf{x})) p(\mathbf{x}, y) dy = 0$$

$$\int y p(\mathbf{x}, y) dy = f(\mathbf{x}) \int p(\mathbf{x}, y) dy$$

Loss Functions in Regression

$$\int yp(\mathbf{x}, y) dy = f(\mathbf{x}) \int p(\mathbf{x}, y) dy$$

$$\int yp(\mathbf{x}, y) dy = f(\mathbf{x}) p(\mathbf{x})$$

$$f(\mathbf{x}) = \int y \frac{p(\mathbf{x}, y)}{p(\mathbf{x})} dy = \int yp(y | \mathbf{x}) dy$$

$$f(\mathbf{x}) = \mathbb{E}_{y \sim p(y|\mathbf{x})} [y] = \mathbb{E} [y | \mathbf{x}]$$

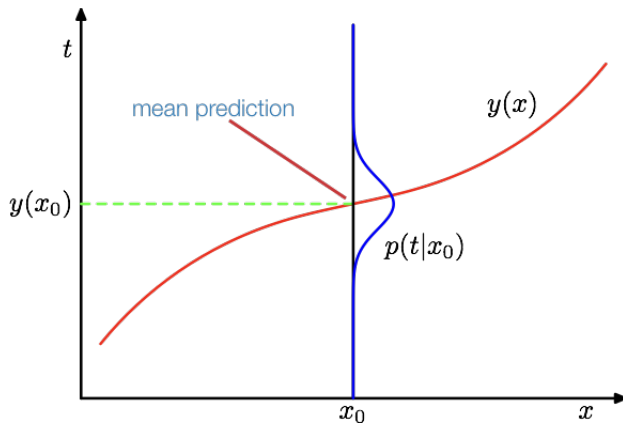
- Under squared loss, the **optimal regression function** is the **mean** $\mathbb{E} [y | \mathbf{x}]$ of the **posterior** $p(y | \mathbf{x})$
- It is also called mean prediction

Loss Functions in Regression

- For our generalized linear regression function

$$f(\mathbf{x}) = \int y \mathcal{N}(y \mid \mathbf{w}^T \phi(\mathbf{x}), \beta^{-1}) dy = \mathbf{w}^T \phi(\mathbf{x})$$

Probabilistic Regression



Outline

1. Introduction to Linear Regression

2. Maximum Likelihood Approach to Regression

3. Bayesian Linear Regression

4. Wrap-Up

Avoiding Overfitting

- Back to our original problem
 - We wanted to avoid overfitting and instabilities
 - Maximum likelihood also leads to overfitting (in the extreme case think if you only had one data point)
- What can we use to counter the problem?

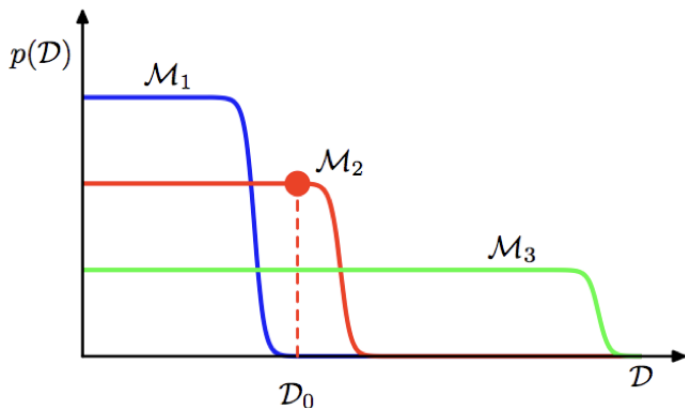
Bayesian Linear Regression

- We place a prior on the parameters \mathbf{w} to tame the instabilities

$$p(\mathbf{w} \mid \mathbf{X}, \mathbf{y}) \propto p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}) p(\mathbf{w})$$

- Parameter prior: $p(\mathbf{w})$
- Likelihood of targets under the data and parameters (as before):
 $p(\mathbf{y} \mid \mathbf{X}, \mathbf{w})$
- Posterior over the parameters: $p(\mathbf{w} \mid \mathbf{X}, \mathbf{y})$
- Notice the **VERY** important difference: in this setting, you do not get anymore a single value for the parameters, but rather a **probability distribution over the parameters**

Basic Idea: Prior controls the Model Class and hence what Data Sets can be explained



Bayesian Regression

- Simple idea: Put a Gaussian prior on \mathbf{w}
- It will put a “soft” limit on the coefficients and thus avoid instabilities

$$\mathbf{w} \sim p(\mathbf{w} \mid \alpha) = \mathcal{N}(\mathbf{w} \mid \mathbf{0}, \alpha^{-1} \mathbf{I})$$

- We use a zero mean Gaussian to keep the derivation compact, but you can use another mean
- Zero mean and spherical covariance (given by the diagonal covariance matrix)
- The posterior becomes

$$\begin{aligned} p(\mathbf{w} \mid \mathbf{X}, \mathbf{y}, \alpha, \beta) &\propto p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta) p(\mathbf{w} \mid \alpha) \\ &\propto p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta) \mathcal{N}(\mathbf{w} \mid \mathbf{0}, \alpha^{-1} \mathbf{I}) \end{aligned}$$

Maximum A-Posteriori (MAP)

- First attempt to solve this problem: estimate \mathbf{w} by maximizing the (log) posterior

$$\begin{aligned}\log p(\mathbf{w} \mid \mathbf{X}, \mathbf{y}, \alpha, \beta) &= \log p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta) + \log \mathcal{N}(\mathbf{w} \mid \mathbf{0}, \alpha^{-1} \mathbf{I}) + \text{const} \\ &= \sum_{i=1}^n \log \mathcal{N}(y_i \mid \mathbf{w}^T \phi(\mathbf{x}_i), \beta^{-1}) \\ &\quad + \log \mathcal{N}(\mathbf{w} \mid \mathbf{0}, \alpha^{-1} \mathbf{I}) + \text{const} \\ &= -\frac{\beta}{2} \sum_{i=1}^n (y_i - \mathbf{w}^T \phi(\mathbf{x}_i))^2 - \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} + \text{const}\end{aligned}$$

Maximum A-Posteriori (MAP)

$$\nabla_{\mathbf{w}} \log p(\mathbf{w} \mid \mathbf{X}, \mathbf{y}, \alpha, \beta) = \beta \sum_{i=1}^n (y_i - \mathbf{w}^\top \phi(\mathbf{x}_i)) \phi(\mathbf{x}_i) - \alpha \mathbf{w} = \mathbf{0}$$

$$\beta \sum_{i=1}^n y_i \phi(\mathbf{x}_i) = \beta \left[\sum_{i=1}^n \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^\top \right] \mathbf{w} + \alpha \mathbf{w}$$

$$\beta \sum_{i=1}^n y_i \phi(\mathbf{x}_i) = \beta \left[\sum_{i=1}^n \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^\top + \alpha \mathbf{I} \right] \mathbf{w}$$

$$\beta \Phi \mathbf{y} = (\beta \Phi \Phi^\top + \alpha \mathbf{I}) \mathbf{w}$$

$$\mathbf{w}_{\text{MAP}} = \left(\Phi \Phi^\top + \frac{\alpha}{\beta} \mathbf{I} \right)^{-1} \Phi \mathbf{y}$$

- What is the role of α/β in the expression?

Maximum A-Posteriori (MAP)

$$\mathbf{w}_{\text{MAP}} = \left(\Phi \Phi^T + \frac{\alpha}{\beta} \mathbf{I} \right)^{-1} \Phi \mathbf{y}$$

- The prior has the effect that it regularizes the pseudo-inverse
- Also called ridge regression

Maximum A-Posteriori (MAP) vs Regularized Least-squares Linear Regression

- There is another way to look at the MAP result
- Let us add a **regularization term** to our objective from Least-squares Linear Regression

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \frac{1}{2} \left\| \hat{\mathbf{X}}^T \mathbf{w} - \mathbf{y} \right\|^2 + \frac{\lambda}{2} \left\| \mathbf{w} \right\|^2$$

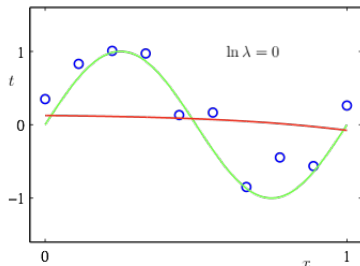
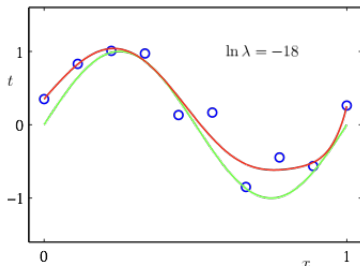
- Solving for \mathbf{w} we get a new estimate

$$\hat{\mathbf{w}} = \left(\hat{\mathbf{X}} \hat{\mathbf{X}}^T + \lambda \mathbf{I} \right)^{-1} \hat{\mathbf{X}} \mathbf{y}$$

- where $\lambda = \alpha/\beta$
- When you place a regularizer λ in least-squares linear regression, you are assuming the targets have Gaussian distributed noise, but also that your parameters are Gaussian distributed

Bayesian Regression

- Polynomial of degree 9 with prior on \mathbf{w}



- $\lambda = \alpha/\beta$ controls the **complexity of the model** and determines the **degree of overfitting**

Bayesian Regression

Table of the coefficients w^* for $M = 9$ polynomials with various values for the regularization parameter λ . Note that $\ln \lambda = -\infty$ corresponds to a model with no regularization, i.e., to the graph at the bottom right in Figure 1.4. We see that, as the value of λ increases, the typical magnitude of the coefficients gets smaller.

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

[Bishop]

Full Bayesian Regression

- We can go further than MAP estimation
- Observation: We do not actually need to know \mathbf{w} , all we want to do is to **predict a function value** based on the training data
- Idea: “Remove” \mathbf{w} by **marginalizing** over it

$$p(y_t | \mathbf{x}_t, \mathbf{X}, \mathbf{y}) = \int p(y_t, \mathbf{w} | \mathbf{x}_t, \mathbf{X}, \mathbf{y}) d\mathbf{w}$$

- y_t - predicted value; \mathbf{x}_t - test input; \mathbf{X} - training data points; \mathbf{y} - training function values

Full Bayesian Regression

$$\begin{aligned}\underbrace{p(y_t | \mathbf{x}_t, \mathbf{X}, \mathbf{y})}_{\text{predictive distribution}} &= \int p(y_t, \mathbf{w} | \mathbf{x}_t, \mathbf{X}, \mathbf{y}) d\mathbf{w} \\ &= \int p(y_t | \mathbf{w}, \mathbf{x}_t, \mathbf{X}, \mathbf{y}) p(\mathbf{w} | \mathbf{x}_t, \mathbf{X}, \mathbf{y}) d\mathbf{w} \\ &= \int \underbrace{p(y_t | \mathbf{w}, \mathbf{x}_t)}_{\text{regression model}} \underbrace{p(\mathbf{w} | \mathbf{X}, \mathbf{y})}_{\text{posterior distribution}} d\mathbf{w}\end{aligned}$$

- For Gaussian distributions, this can be done in closed form, leading to so-called **Gaussian Processes**

Full Bayesian Regression

- We can also do that in closed form: **integrate out all possible parameters**

$$p(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \int \underbrace{p(y_* | \mathbf{x}_*, \theta)}_{\text{likelihood}} \underbrace{p(\theta | \mathbf{X}, \mathbf{y})}_{\text{parameter posterior}} d\theta$$

- y_* - predicted value; \mathbf{x}_* - test input; \mathbf{X}, \mathbf{y} - training data
- The **predictive distribution** is again a Gaussian

$$p(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = \mathcal{N}(y_* | \mu(\mathbf{x}_*), \sigma^2(\mathbf{x}_*))$$

$$\mu(\mathbf{x}_*) = \phi^T(\mathbf{x}_*) \left(\frac{\alpha}{\beta} \mathbf{I} + \Phi \Phi^T \right)^{-1} \Phi^T \mathbf{y}$$

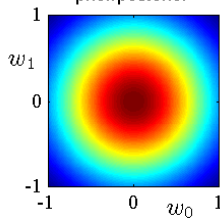
$$\sigma^2(\mathbf{x}_*) = \frac{1}{\beta} + \phi^T(\mathbf{x}_*) (\alpha \mathbf{I} + \beta \Phi \Phi^T)^{-1} \phi(\mathbf{x}_*)$$

- **The variance is state dependent**

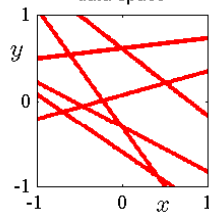
Bayesian Regression

likelihood

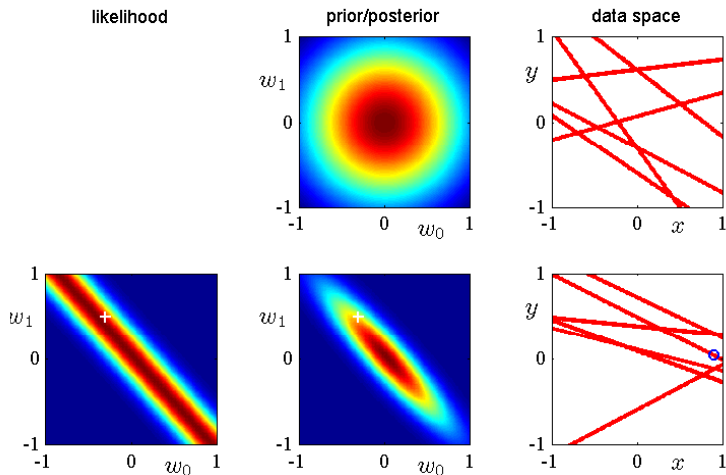
prior/posterior



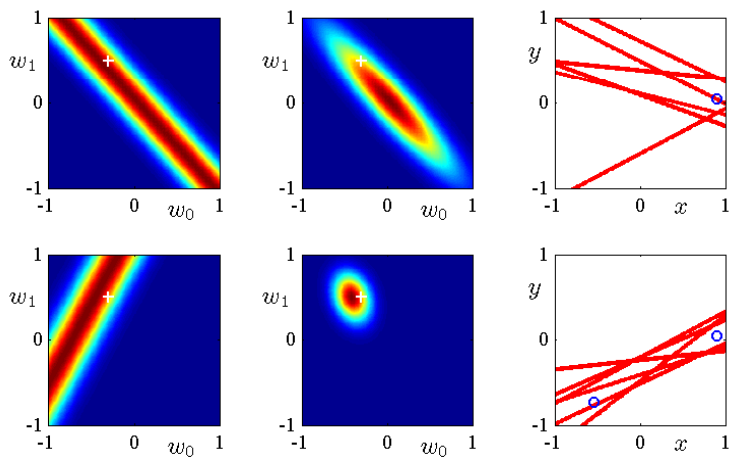
data space



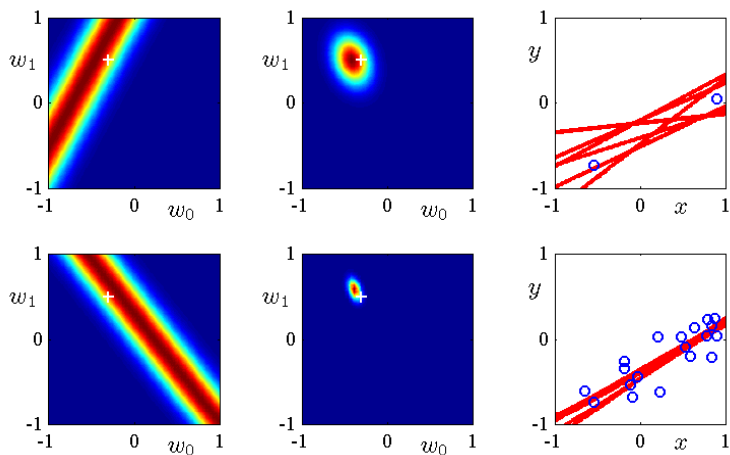
Bayesian Regression



Bayesian Linear Regression

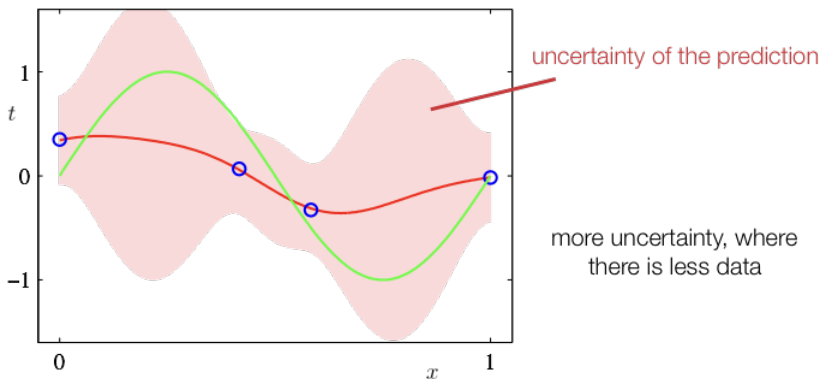


Bayesian Linear Regression

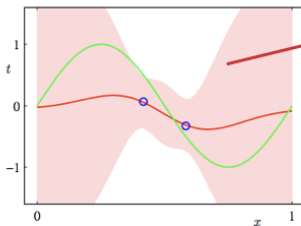
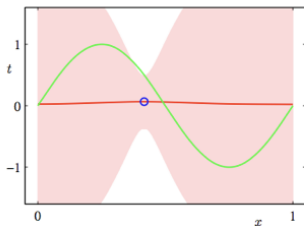


Gaussian Processes - Quick Preview

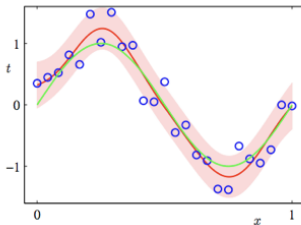
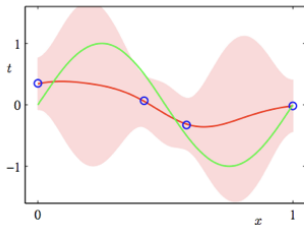
- We will not cover them now, but here is a quick preview of what they can do



Gaussian Processes - Quick Preview



uncertainty of
the prediction



more
uncertainty,
where there is
less data

Outline

1. Introduction to Linear Regression

2. Maximum Likelihood Approach to Regression

3. Bayesian Linear Regression

4. Wrap-Up

4. Wrap-Up

You know now:

- How to formulate a linear regression problem
- The different methods to perform linear regression: least-squares, maximum likelihood and bayesian
- Derive the equations for the parameters using the different methods
- Why introducing a prior distribution over the parameters can combat overfitting

Self-Test Questions

- What is regression (in general) and linear regression (in particular)?
- What is the cost function of regression and how can I interpret it?
- What is overfitting?
- How can I derive a Maximum-Likelihood Estimator for Regression?
- Why are Bayesian methods important?
- What is MAP and how is it different to full Bayesian regression?

Homework

- Reading Assignment for next lecture
 - Murphy ch. 8
 - Bishop ch. 4