

Machine Learning

Lecture 3: Linear Regression

Prof. Dr. Stephan Günnemann

Data Mining and Analytics
Technical University of Munich

06.11.2018

Reading material

Reading material

- Bishop [ch. 1.1, 3.1, 3.2, 3.3.1, 3.3.2, 3.6]
- Murphy [ch. 7.2 - 7.3, 7.5.1, 7.6.1, 7.6.2]

Acknowledgements

- Slides are based on an older version by G. Jensen and C. Osendorfer
- Some figures are from C. Bishop: "Pattern Recognition and Machine Learning"

Notation

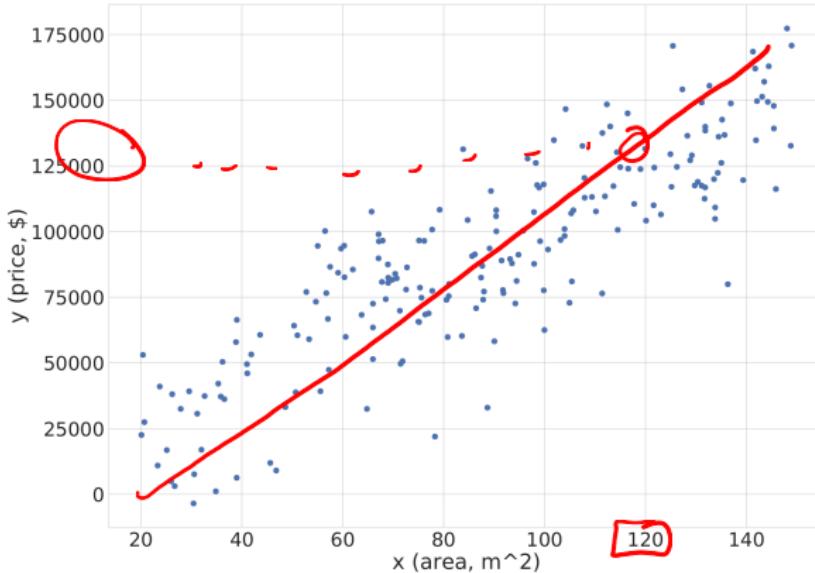
Symbol	Meaning
s	scalar is lowercase and not bold
\mathbf{s}	vector is lowercase and bold
S	matrix is uppercase and bold
$f(\mathbf{x})$	predicted value for inputs \mathbf{x}
\mathbf{y}	vector of targets
y_i	target of the i 'th example
w_0	bias term (not to be confused with bias in general)
$\phi(\cdot)$	basis function
$E(\cdot)$	error function
\mathcal{D}	training data
\mathbf{X}^\dagger	Moore-Penrose pseudoinverse of \mathbf{X}

Section 1

Basic Linear Regression

Example: Housing price prediction

Given is a dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$, of house areas x_i and corresponding prices y_i .



How do we estimate a price of a new house with area x_{new} ?

Regression problem

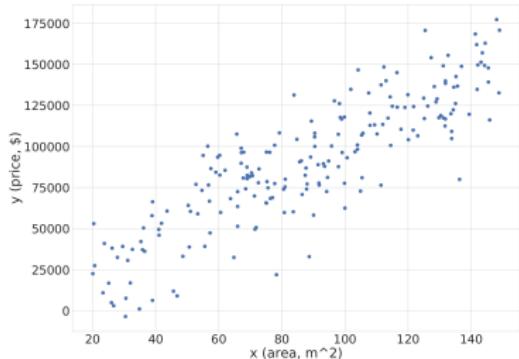
Given

- observations ¹
 $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}, \mathbf{x}_i \in \mathbb{R}^D$
- targets
 $\mathbf{y} = \{y_1, y_2, \dots, y_N\}, y_i \in \mathbb{R}$

Find

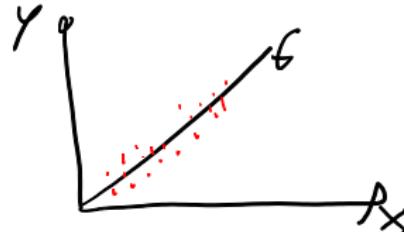
- Mapping $f(\cdot)$ from inputs to targets

$$y_i \approx f(\mathbf{x}_i)$$



¹A common way to represent the samples is as a **data matrix** $\mathbf{X} \in \mathbb{R}^{N \times D}$, where each row represents one sample.

Linear model



Target y is generated by a deterministic function f of x plus noise

$$y_i = f(\mathbf{x}_i) + \epsilon_i, \quad \epsilon_i \sim \mathcal{N}(0, \beta^{-1}) \quad (1)$$

Let's choose $f(x)$ to be a linear function

$$f_{\mathbf{w}}(\mathbf{x}_i) = w_0 + w_1 x_{i1} + w_2 x_{i2} + \dots + w_D x_{iD} \quad (2)$$

$$= w_0 + \mathbf{w}^T \mathbf{x}_i \quad (3)$$

Absorbing the bias term

The linear function is given by

$$f_{\mathbf{w}}(\mathbf{x}) = w_0 + w_1x_1 + w_2x_2 + \dots + w_Dx_D \quad (4)$$

$$= w_0 + \mathbf{w}^T \mathbf{x} \quad (5)$$

Here w_0 is called **bias** or **offset** term. For simplicity, we can "absorb" it by prepending a 1 to the feature vector \mathbf{x} and respectively adding w_0 to the weight vector \mathbf{w} :

$$\tilde{\mathbf{x}} = (1, x_1, \dots, x_D)^T \quad \tilde{\mathbf{w}} = (w_0, w_1, \dots, w_D)^T$$

The function $f_{\mathbf{w}}$ can compactly be written as $f_{\mathbf{w}}(\mathbf{x}) = \tilde{\mathbf{w}}^T \tilde{\mathbf{x}}$.

To unclutter the notation, we will assume the bias term is always absorbed and write \mathbf{w} and \mathbf{x} instead of $\tilde{\mathbf{w}}$ and $\tilde{\mathbf{x}}$.

Now, how do we choose the "best" \mathbf{w} that fits our data?

Error function

Loss function

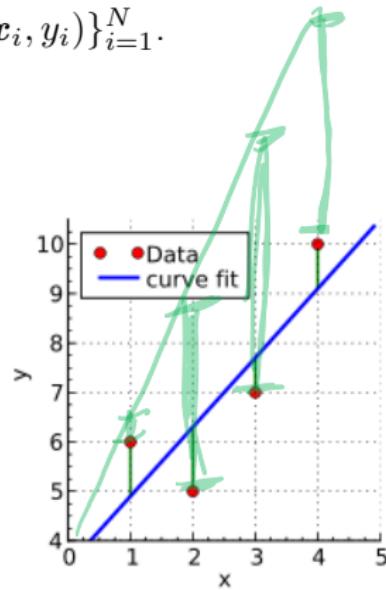


Error function gives a measure of "misfit" between our model (parametrized by w) and observed data $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$.

Standard choice - least squares (LS) function

$$E_{\text{LS}}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (f_{\mathbf{w}}(\mathbf{x}_i) - y_i)^2 \quad (6)$$

$$= \frac{1}{2} \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}_i - y_i)^2 \quad (7)$$



Objective

Find the optimal weight vector \mathbf{w}^* that minimizes the error

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} E_{\text{LS}}(\mathbf{w}) \quad (8)$$

$$= \arg \min_{\mathbf{w}} \frac{1}{2} \sum_{i=1}^N (\mathbf{x}_i^T \mathbf{w} - y_i)^2 \quad (9)$$

By stacking the observations \mathbf{x}_i as rows of the matrix $\mathbf{X} \in \mathbb{R}^{N \times D}$

$$= \arg \min_{\mathbf{w}} \frac{1}{2} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) \quad (10)$$

To find the minimum of the function $E(\mathbf{w})$, compute the gradient $\nabla_{\mathbf{w}} E(\mathbf{w})$:

$$\nabla_{\mathbf{w}} E_{\text{LS}}(\mathbf{w}) = \nabla_{\mathbf{w}} \frac{1}{2} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y}) \quad (11)$$

$$= \nabla_{\mathbf{w}} \frac{1}{2} \left(\mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{y} + \mathbf{y}^T \mathbf{y} \right) \quad (12)$$

$$= \mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y} \quad (13)$$

Optimal solution

Now set the gradient to zero and solve for \mathbf{w} to obtain the minimizer²

$$\mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y} \stackrel{!}{=} 0 \quad (14)$$

This leads to the so-called **normal equation** of the least squares problem

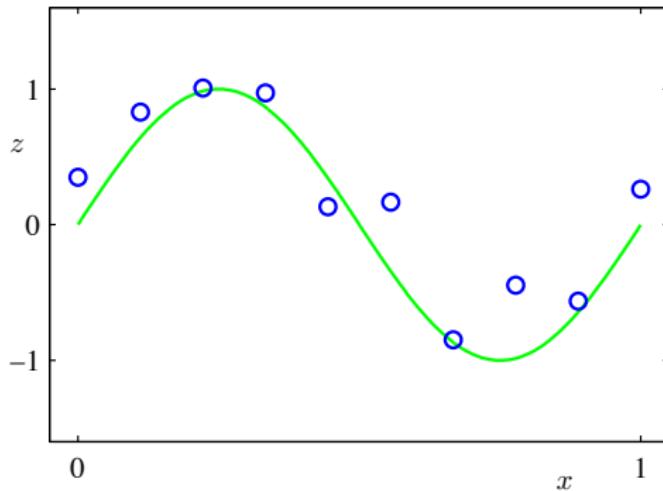
$$\mathbf{w}^* = \underbrace{(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T}_{=\mathbf{X}^\dagger} \mathbf{y} \quad (15)$$

\mathbf{X}^\dagger is called **Moore-Penrose pseudo-inverse** of \mathbf{X} (because for an invertible square matrix, $\mathbf{X}^\dagger = \mathbf{X}^{-1}$).

²Because Hessian $\nabla_{\mathbf{w}} \nabla_{\mathbf{w}} E(\mathbf{w})$ is positive (semi)definite → see *Optimization*

Nonlinear dependency in data

What if the dependency between y and x is not linear?



Data generating process: $y_i = \sin(2\pi x_i) + \epsilon_i, \quad \epsilon_i \sim \mathcal{N}(0, \beta^{-1})$

$$\textcolor{red}{w_0 = 1}$$

Polynomials

ATTENTION: HERE 1-0

ON THIS SLIDE
AND NEXT ONE

Solution: Polynomials are universal function approximators, so we can define f as

$$f_{\mathbf{w}}(\underline{x}) = w_0 + \sum_{j=1}^M w_j x^j \quad (16)$$

Or more generally

$$= w_0 + \sum_{j=1}^M w_j \phi_j(x) \quad (17)$$

Define $\phi_0 = 1$

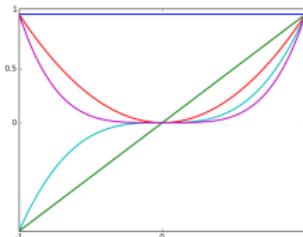
$$= \mathbf{w}^T \phi(x) \quad (18)$$

The function f is still linear in \mathbf{w} (despite not being linear in x)!

Typical basis functions

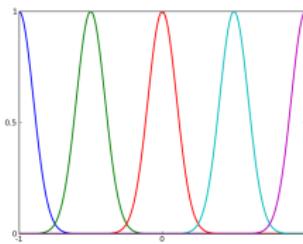
Polynomials

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



Gaussian

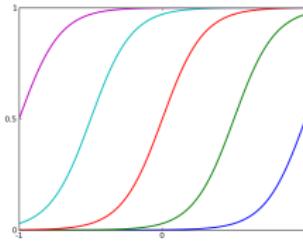
$$\phi_j(x) = e^{\frac{-(x-\mu_j)^2}{2s^2}}$$



Logistic Sigmoid

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

$$\text{where } \sigma(a) = \frac{1}{1+e^{-a}}$$



Linear basis function model

Prediction for one sample

CLASSIFICATION $\emptyset = \times$

$$f_{\mathbf{w}}(\mathbf{x}) = w_0 + \sum_{j=1}^M w_j \phi_j(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) \quad (19)$$

Using the same least squares error function as before

$$E_{LS}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_i) - y_i)^2 \quad (20)$$

$$\phi_0(\mathbf{x}) = 1 \quad = \frac{1}{2} (\boldsymbol{\Phi} \mathbf{w} - \mathbf{y})^T (\boldsymbol{\Phi} \mathbf{w} - \mathbf{y}) \quad (21)$$

with

$$\boldsymbol{\Phi} = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \dots & \phi_M(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & & \vdots \\ \vdots & \vdots & \ddots & \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \dots & \phi_M(\mathbf{x}_N) \end{pmatrix} \in \mathbb{R}^{N \times (M+1)}$$

being the **design matrix** of $\boldsymbol{\phi}$.

$\phi_j : \mathbb{R}^D \rightarrow \mathbb{R}$

Optimal solution

Recall Equation 10 - we have the same expression except that data matrix $\mathbf{X} \in \mathbb{R}^{N \times D}$ is replaced by design matrix $\Phi \in \mathbb{R}^{N \times (M+1)}$

$$E_{\text{LS}}(\mathbf{w}) = \frac{1}{2}(\Phi\mathbf{w} - \mathbf{y})^T(\Phi\mathbf{w} - \mathbf{y}) \quad (22)$$

This means that the optimal weights \mathbf{w}^* can be obtained in a similar way

$$\mathbf{w}^* = (\Phi^T\Phi)^{-1}\Phi^T\mathbf{y} \quad (23)$$

$$= \Phi^\dagger\mathbf{y} \quad (24)$$

Compare this to Equation 15:

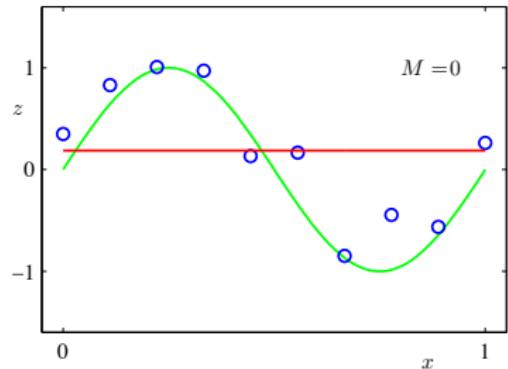
$$\mathbf{w}^* = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y} \quad (25)$$

Choosing degree of the polynomial

How do we choose the degree of the polynomial M ?

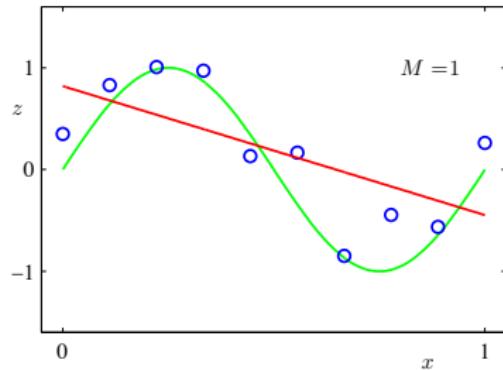
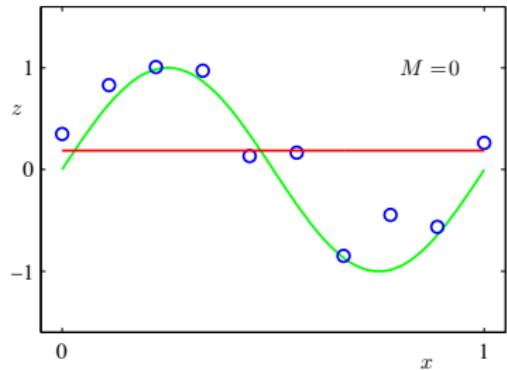
Choosing degree of the polynomial

How do we choose the degree of the polynomial M ?



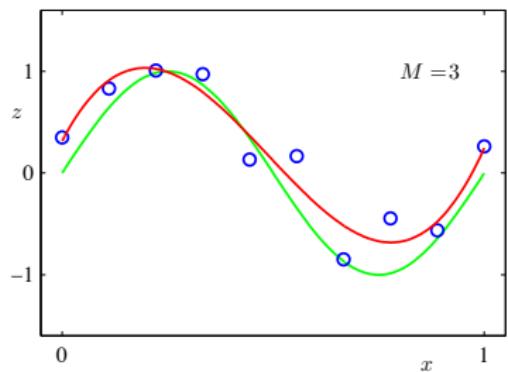
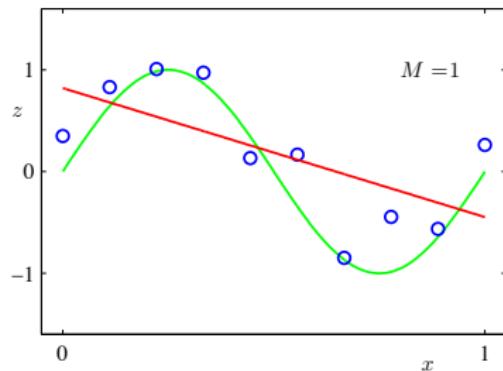
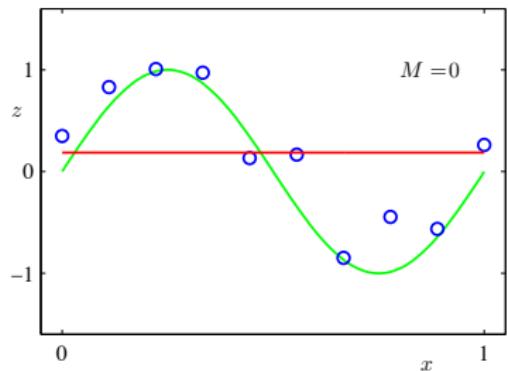
Choosing degree of the polynomial

How do we choose the degree of the polynomial M ?



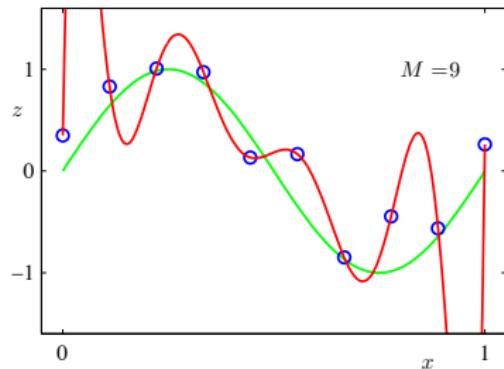
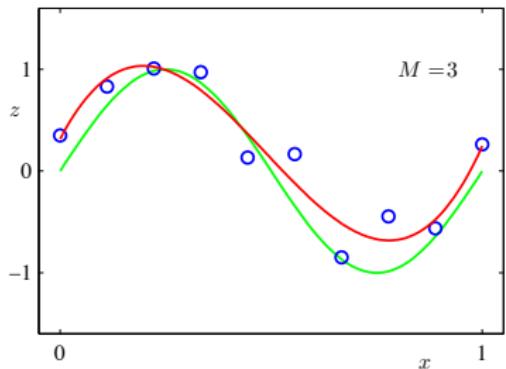
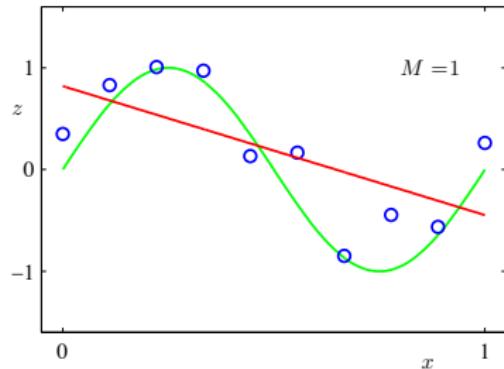
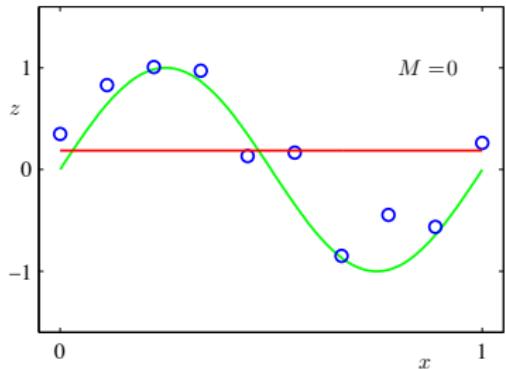
Choosing degree of the polynomial

How do we choose the degree of the polynomial M ?

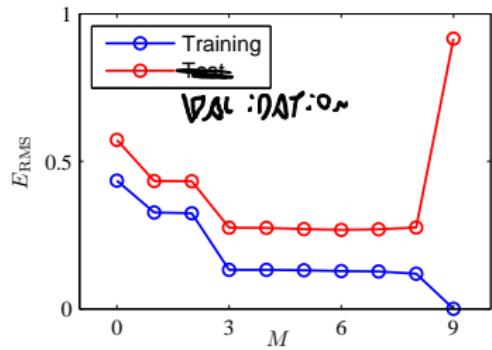
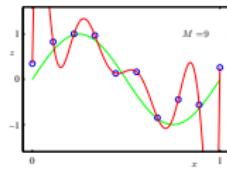
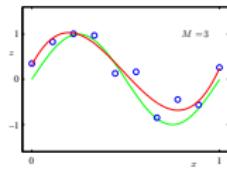
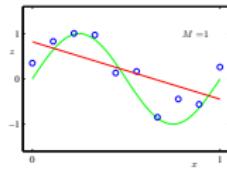
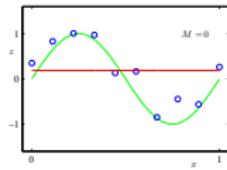


Choosing degree of the polynomial

How do we choose the degree of the polynomial M ?

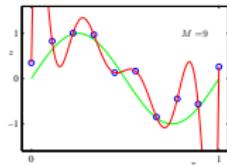
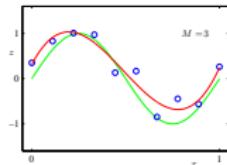
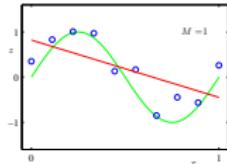
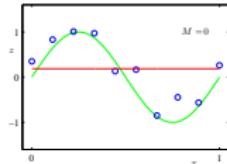


Choosing degree of the polynomial



One valid solution is to choose M using the standard train-validation split approach.

Choosing degree of the polynomial



3

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

We also make another observation: overfitting occurs when the coefficients w become large.

What if we penalize large weights?

Controlling overfitting with regularization

Least squares loss function with L2 regularization
(also called ridge regression)

$$E_{\text{ridge}}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N [\mathbf{w}^T \phi(\mathbf{x}_i) - y_i]^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \quad (26)$$

Where,

- $\|\mathbf{w}\|_2^2 \equiv \mathbf{w}^T \mathbf{w} = w_0^2 + w_1^2 + w_2^2 + \cdots + w_M^2$ - L2 norm
- λ – regularization strength

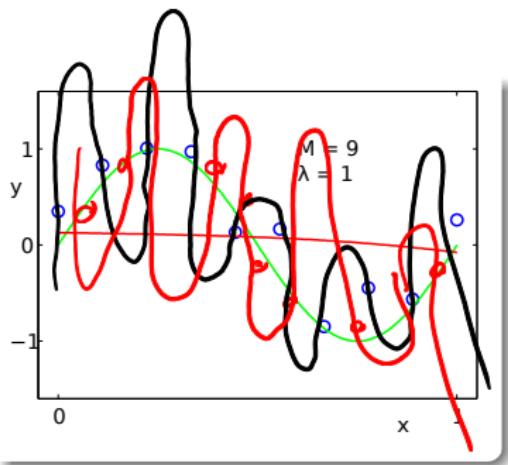
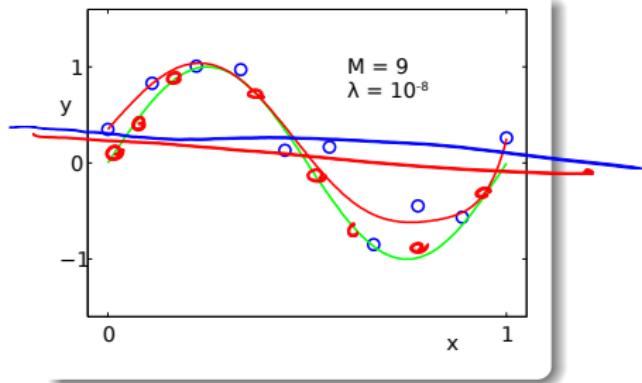
Controlling overfitting with regularization

Least squares loss function with L2 regularization
(also called ridge regression)

$$E_{\text{ridge}}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N [\mathbf{w}^T \phi(\mathbf{x}_i) - y_i]^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2 \quad (26)$$

Where,

- $\|\mathbf{w}\|_2^2 \equiv \mathbf{w}^T \mathbf{w} = w_0^2 + w_1^2 + w_2^2 + \dots + w_M^2$ - L2 norm
- λ – regularization strength



Larger regularization strength λ leads to smaller weights \mathbf{w}

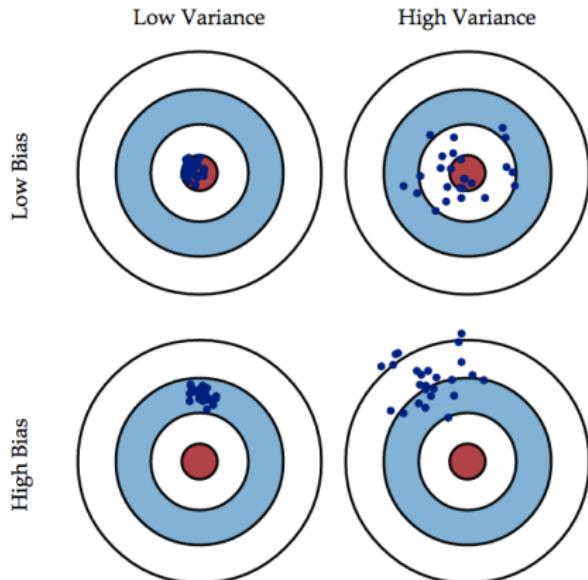
Bias-variance tradeoff

The error of an estimator can be decomposed into two parts:³

- **Bias** - expected error due to model mismatch
- **Variance** - variation due to randomness in training data

the center of the target:
the true model that predicts
the correct values.

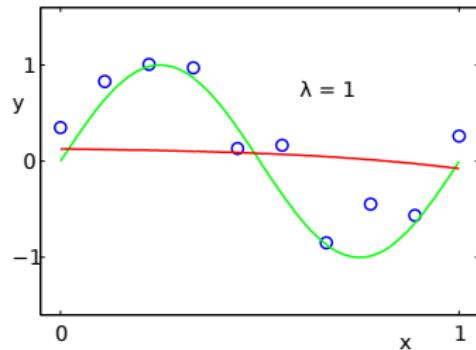
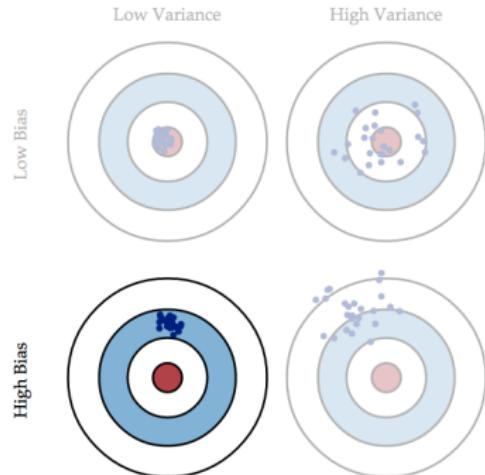
different hits (the blue dots):
different realizations of model
given different training data.



³See Bishop Section 3.2 for a more rigorous mathematical derivation

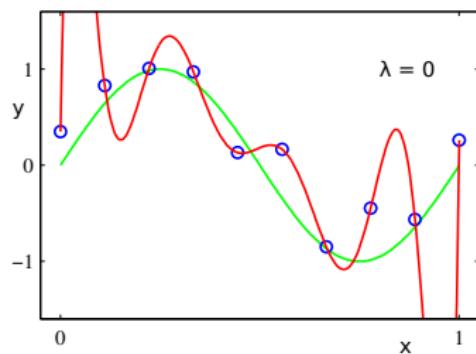
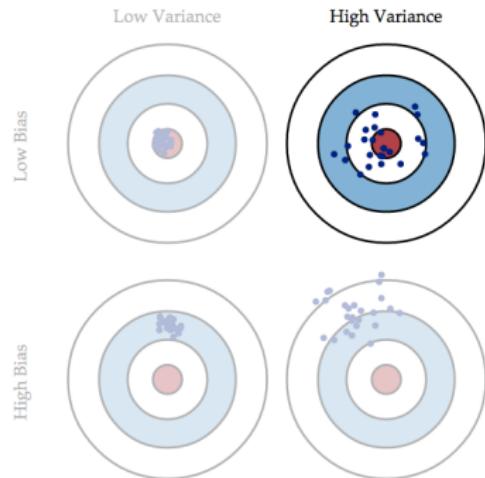
Bias-variance tradeoff: high bias

- In case of **high bias**, the model is too rigid to fit the underlying data distribution.
- This typically happens if the model is misspecified and/or the regularization strength λ is too high.



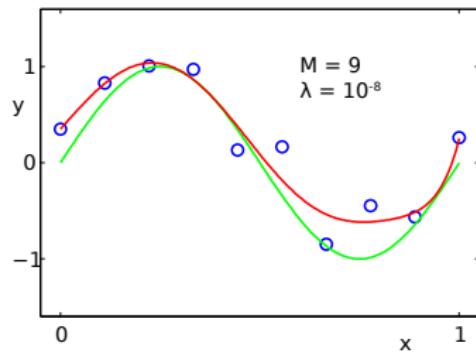
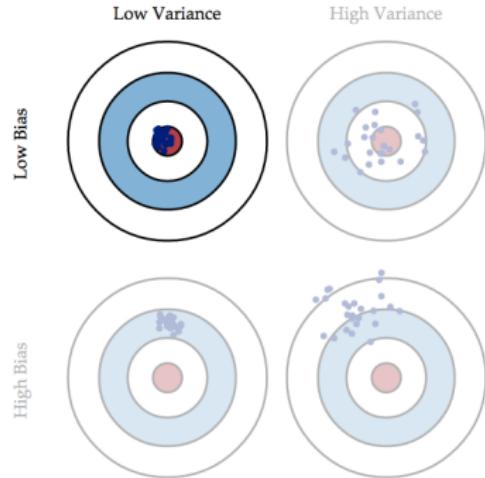
Bias-variance tradeoff: high variance

- In case of **high variance**, the model is too flexible, and therefore captures noise in the data.
- This is exactly what we call **overfitting**.
- This typically happens when the model has high capacity (= it "memorizes" the training data) and/or λ is too low.



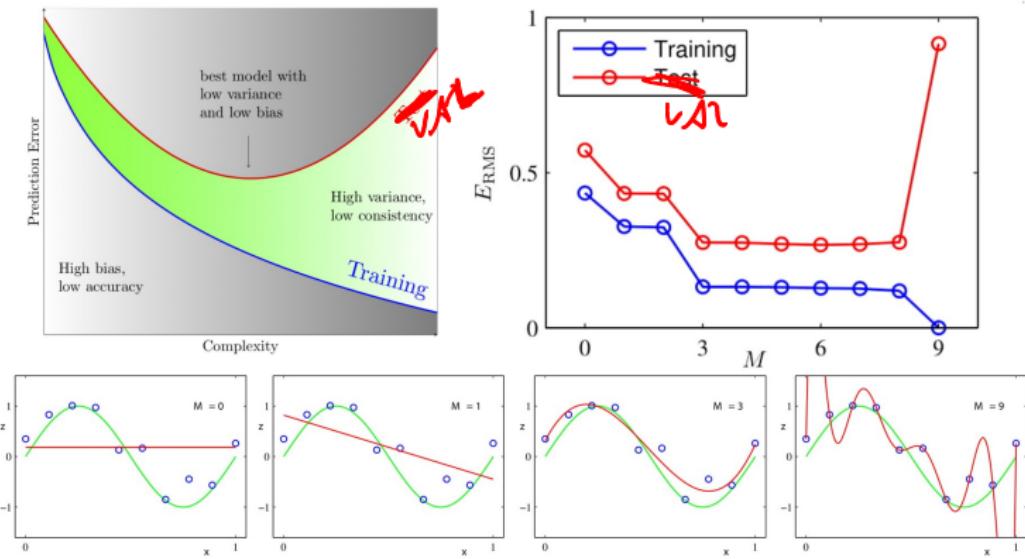
Bias-variance tradeoff

- Of course, we want models that have low bias and low variance, but often those are conflicting goals.
- A popular technique is to select a model with large capacity (e.g. high degree polynomial), and keep the variance in check by choosing appropriate regularization strength λ .



Bias-variance tradeoff

- Bias-variance tradeoff in the case of slide 22 ($\lambda = 0$).



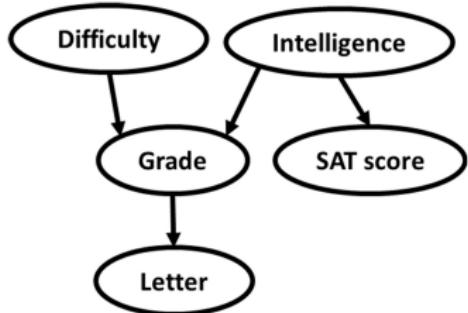
The upper-left figure from: <https://eissanematollahi.com/wp-content/uploads/2018/09/Machine-Learning-Basics-1.pdf>.

Section 2

Probabilistic Linear Regression

Introduction to probabilistic graphical models (PGMs) ⁴

- Bayesian network: is a way to represent a joint distribution via a directed acyclic graph $G = (V, E)$.
 - Each node represents a random variable.
 - Directed edges are often interpreted as causal relations.



random variables $\in \{0 : \text{low}, 1 : \text{high}\}$

- Graph G provides:
 - 1) a particular factorization for the joint distribution $p(x_1, \dots, x_{|V|})$.
 - 2) a set of conditional independencies,
inferred by a routine $d\text{-separation}(G)$.

Example from: [Daphne Koller and Nir Friedman. 2009. Probabilistic Graphical Models: Principles and Techniques. The MIT Press].

⁴Here we introduce only "undirected" graphical models, i.e. Bayesian networks.

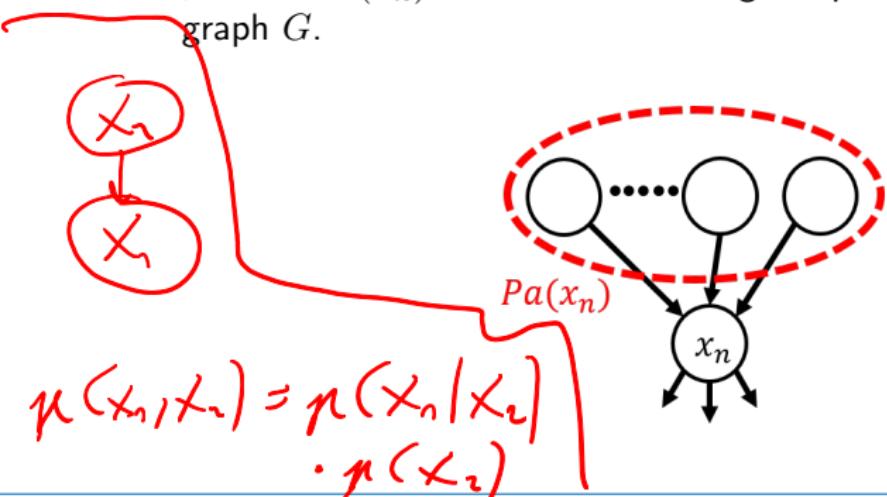
Bayesian networks - factorization

$$p(x_1, x_2) = p(x_1) \cdot p(x_2)$$
$$p(x_1, x_2 | x_3) = p(x_1 | x_3) \cdot p(x_2 | x_3)$$

- A Bayesian network specifies the following factorization for the joint distribution $p(x_1, \dots, x_{|V|})$.

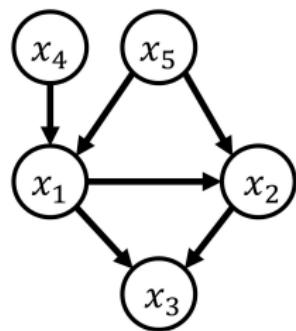
$$p(x_1, \dots, x_{|V|}) = \prod_{n=1}^{|V|} p(x_n | Pa(x_n)) \quad (27)$$

, where $Pa(x_n)$ is the set containing the parents of node x_n in graph G .



Bayesian networks - factorization

- Example:



$$p(x_1, x_2, x_3, x_4, x_5) =$$

×

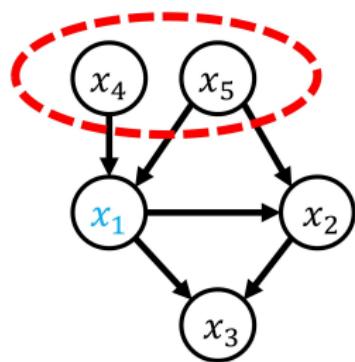
×

×

×

Bayesian networks - factorization

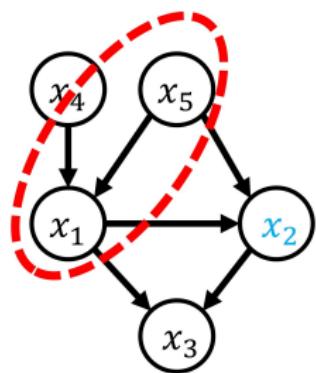
- Example:



$$p(x_1, x_2, x_3, x_4, x_5) = p(\textcolor{blue}{x_1} | \textcolor{red}{x_4}, \textcolor{red}{x_5}) \times \times \times \times$$

Bayesian networks - factorization

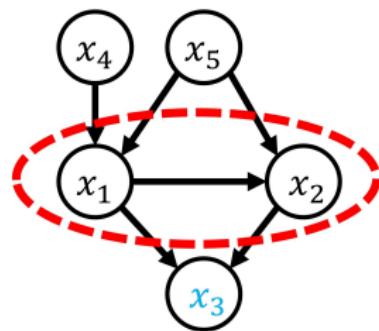
- Example:



$$\begin{aligned} p(x_1, x_2, x_3, x_4, x_5) &= p(x_1|x_4, x_5) \times \\ &\quad p(\textcolor{blue}{x}_2|\textcolor{red}{x}_1, \textcolor{red}{x}_5) \times \\ &\quad \times \\ &\quad \times \end{aligned}$$

Bayesian networks - factorization

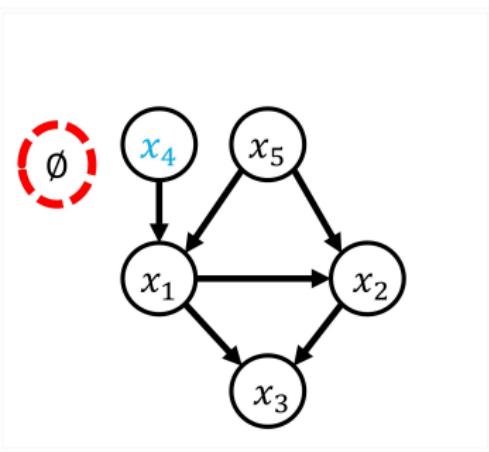
- Example:



$$p(x_1, x_2, x_3, x_4, x_5) = p(x_1|x_4, x_5) \times \\ p(x_2|x_1, x_5) \times \\ p(\textcolor{blue}{x}_3|\textcolor{red}{x}_1, \textcolor{red}{x}_2) \times \\ \times$$

Bayesian networks - factorization

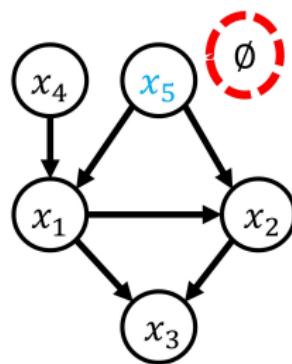
- Example:



$$\begin{aligned} p(x_1, x_2, x_3, x_4, x_5) &= p(x_1|x_4, x_5) \times \\ &p(x_2|x_1, x_5) \times \\ &p(x_3|x_1, x_2) \times \\ &p(\textcolor{blue}{x}_4) \times \end{aligned}$$

Bayesian networks - factorization

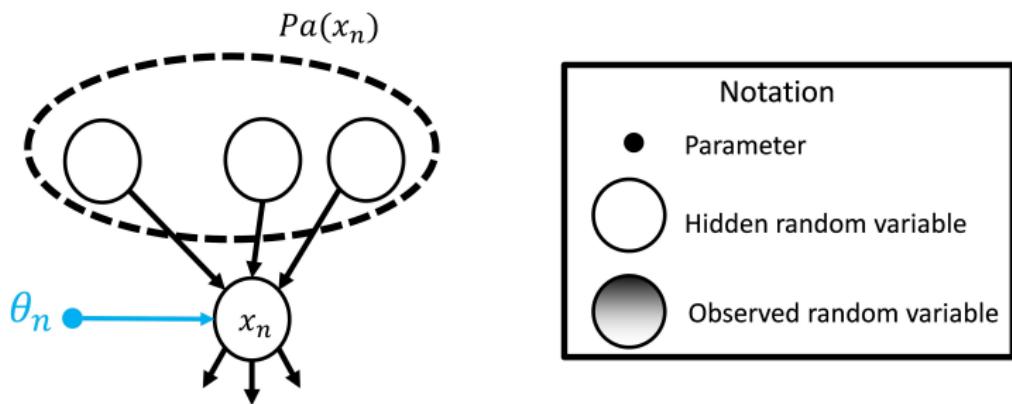
- Example:



$$\begin{aligned} p(x_1, x_2, x_3, x_4, x_5) &= p(x_1|x_4, x_5) \times \\ &p(x_2|x_1, x_5) \times \\ &p(x_3|x_1, x_2) \times \\ &p(x_4) \times \\ &p(x_5) \end{aligned}$$

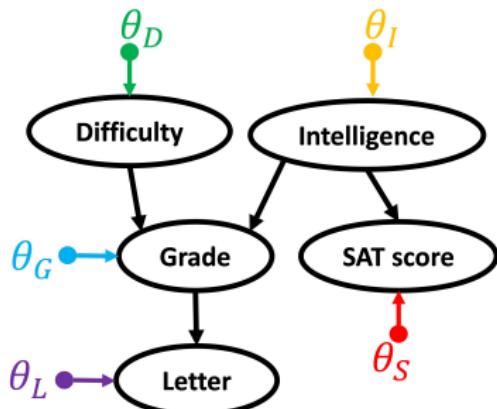
Bayesian networks - model parameters

- A set of parameters $\{\theta_1, \dots, \theta_{|V|}\}$ parameterize the joint distribution.
- θ_n parameterizes the factor $p(x_n | Pa(x_n))$
 - Indeed, the factor is a function of θ_n as well.
we can write $p(x_n | Pa(x_n), \theta_n)$
 - In a graphical model, we show θ_n using a filled circle connected to x_n .



Bayesian networks - model parameters

- Example:



		$p(D = 0)$	$p(D = 1)$			$p(I = 0)$	$p(I = 1)$
		0.5	0.5			0.5	0.5

D	I	$p(G = 0 D, I)$	$p(G = 1 D, I)$
0	0	0.5	0.5
0	1	0.1	0.9
1	0	0.9	0.1
1	1	0.5	0.5

I	$p(S = 0 I)$	$p(S = 1 I)$
0	0.9	0.1
1	0.1	0.9

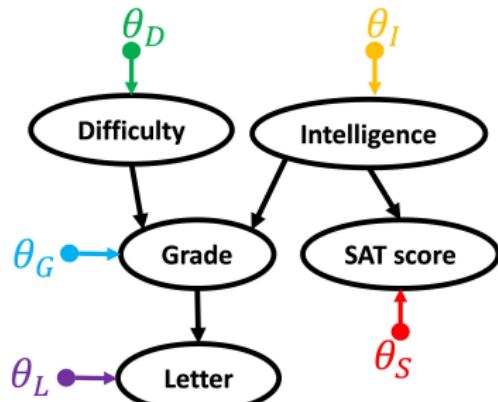
G	$p(L = 0 G)$	$p(L = 1 G)$
0	0.9	0.1
1	0.1	0.9

$$\pi(x_1, x_2, x_3, x_4, x_5) =$$

Example from: [Daphne Koller and Nir Friedman. 2009. Probabilistic Graphical Models: Principles and Techniques. The MIT Press].

Bayesian networks - generative process

- We can think of the model as a generative process.



$$\text{Difficulty} \sim p(D|\theta_D)$$

$$\text{Intelligence} \sim p(I|\theta_I)$$

$$\text{Grade} \sim p(G|D, I, \theta_G)$$

$$\text{SAT} \sim p(S|I, \theta_S)$$

$$\text{Letter} \sim p(L|G, \theta_L)$$

Example from: [Daphne Koller and Nir Friedman. 2009. Probabilistic Graphical Models: Principles and Techniques. The MIT Press].

Probabilistic formulation of linear regression

Remember from our problem definition at the start of the lecture,

$$y_i = f_w(\mathbf{x}_i) + \underbrace{\epsilon_i}_{\text{noise}}$$

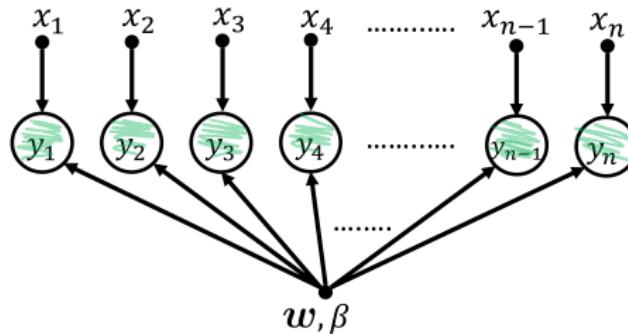
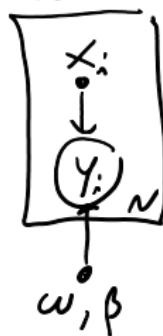
Noise has zero-mean Gaussian distribution with a fixed precision $\beta = \frac{1}{\sigma^2}$

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

This implies that the distribution of the targets is

$$y_i \sim \mathcal{N}(f_w(\mathbf{x}_i), \beta^{-1})$$

PLATE
NOTATION



Maximum likelihood

Likelihood of a single sample

$$p(y_i \mid f_{\mathbf{w}}(\mathbf{x}_i), \beta) = \mathcal{N}(y_i \mid f_{\mathbf{w}}(\mathbf{x}_i), \beta^{-1}) \quad (28)$$

Assume that the samples are drawn independently

\Rightarrow likelihood of the entire dataset $\mathcal{D} = \{\mathbf{X}, \mathbf{y}\}$ is

$$p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta) = \prod_{i=1}^N p(y_i \mid f_{\mathbf{w}}(\mathbf{x}_i), \beta) \quad (29)$$

We can now use the same approach we used in previous lecture -
maximize the likelihood w.r.t. \mathbf{w} and β

$$\mathbf{w}_{\text{ML}}, \beta_{\text{ML}} = \arg \max_{\mathbf{w}, \beta} p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta) \quad (30)$$

Maximum likelihood

Like in the coin flip example, we can make a few simplifications

$$\boldsymbol{w}_{\text{ML}}, \beta_{\text{ML}} = \arg \max_{\boldsymbol{w}, \beta} p(\boldsymbol{y} \mid \boldsymbol{X}, \boldsymbol{w}, \beta) \quad (31)$$

$$= \arg \max_{\boldsymbol{w}, \beta} \ln p(\boldsymbol{y} \mid \boldsymbol{X}, \boldsymbol{w}, \beta) \quad (32)$$

$$= \arg \min_{\boldsymbol{w}, \beta} -\ln p(\boldsymbol{y} \mid \boldsymbol{X}, \boldsymbol{w}, \beta) \quad (33)$$

Let's denote this quantity as **maximum likelihood error function** that we need to minimize

$$E_{\text{ML}}(\boldsymbol{w}, \beta) = -\ln p(\boldsymbol{y} \mid \boldsymbol{X}, \boldsymbol{w}, \beta) \quad (34)$$

Maximum likelihood

Simplify the error function

$$E_{\text{ML}}(\mathbf{w}, \beta) = -\ln \left[\prod_{i=1}^N \mathcal{N}(y_i | f_{\mathbf{w}}(\mathbf{x}_i), \beta^{-1}) \right] \quad (35)$$

$$= -\ln \left[\prod_{i=1}^N \sqrt{\frac{\beta}{2\pi}} \exp \left(-\frac{\beta}{2} (\mathbf{w}^T \phi(\mathbf{x}_i) - y_i)^2 \right) \right] \quad (36)$$

$$= -\sum_{i=1}^N \ln \left[\sqrt{\frac{\beta}{2\pi}} \exp \left(-\frac{\beta}{2} (\mathbf{w}^T \phi(\mathbf{x}_i) - y_i)^2 \right) \right] \quad (37)$$

$$= \frac{\beta}{2} \sum_{i=1}^N (\mathbf{w}^T \phi(\mathbf{x}_i) - y_i)^2 - \frac{N}{2} \ln \beta + \frac{N}{2} \ln 2\pi \quad (38)$$

Optimizing log-likelihood w.r.t. \boldsymbol{w}

$$\boldsymbol{w}_{\text{ML}} = \arg \min_{\boldsymbol{w}} E_{\text{ML}}(\boldsymbol{w}, \beta) \quad (39)$$

$$= \arg \min_{\boldsymbol{w}} \left[\frac{\beta}{2} \sum_{i=1}^N (\boldsymbol{w}^T \phi(\boldsymbol{x}_i) - y_i)^2 - \underbrace{\frac{N}{2} \ln \beta + \frac{N}{2} \ln 2\pi}_{= \text{const}} \right] \quad (40)$$

$$= \arg \min_{\boldsymbol{w}} \frac{1}{2} \sum_{i=1}^N (\boldsymbol{w}^T \phi(\boldsymbol{x}_i) - y_i)^2 \quad (41)$$

least squares error fn!

$$= \arg \min_{\boldsymbol{w}} E_{\text{LS}}(\boldsymbol{w}) \quad (42)$$

Optimizing log-likelihood w.r.t. \mathbf{w}

$$\mathbf{w}_{\text{ML}} = \arg \min_{\mathbf{w}} E_{\text{ML}}(\mathbf{w}, \beta) \quad (39)$$

$$= \arg \min_{\mathbf{w}} \left[\frac{\beta}{2} \sum_{i=1}^N (\mathbf{w}^T \phi(\mathbf{x}_i) - y_i)^2 - \underbrace{\frac{N}{2} \ln \beta + \frac{N}{2} \ln 2\pi}_{= \text{const}} \right] \quad (40)$$

$$= \arg \min_{\mathbf{w}} \frac{1}{2} \sum_{i=1}^N (\mathbf{w}^T \phi(\mathbf{x}_i) - y_i)^2 \quad (41)$$

least squares error fn!

$$= \arg \min_{\mathbf{w}} E_{\text{LS}}(\mathbf{w}) \quad (42)$$

Maximizing the likelihood is equivalent to minimizing the least squares error function!

$$\mathbf{w}_{\text{ML}} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y} = \Phi^\dagger \mathbf{y} \quad (43)$$

Optimizing log-likelihood w.r.t. β

Plug in the estimate for w and minimize w.r.t. β

$$\beta_{\text{ML}} = \arg \min_{\beta} E_{\text{ML}}(\mathbf{w}_{\text{ML}}, \beta) \quad (44)$$

$$= \arg \min_{\beta} \left[\frac{\beta}{2} \sum_{i=1}^N (\mathbf{w}_{\text{ML}}^T \phi(\mathbf{x}_i) - y_i)^2 - \frac{N}{2} \ln \beta + \frac{N}{2} \ln 2\pi \right] \quad (45)$$

Take derivative w.r.t. β and set it to zero

$$\frac{\partial}{\partial \beta} E_{\text{ML}}(\mathbf{w}_{\text{ML}}, \beta) = \frac{1}{2} \sum_{i=1}^N (\mathbf{w}_{\text{ML}}^T \phi(\mathbf{x}_i) - y_i)^2 - \frac{N}{2\beta} \stackrel{!}{=} 0 \quad (46)$$

Solving for β

$$\frac{1}{\beta_{\text{ML}}} = \frac{1}{N} \sum_{i=1}^N (\mathbf{w}_{\text{ML}}^T \phi(\mathbf{x}_i) - y_i)^2 \quad (47)$$

Predicting for new data

$\mathcal{N}(\mu, \sigma)$
 $\beta = \frac{1}{\sigma^2}$ Precision

Recall, that

$$y \sim \mathcal{N}(f_w(\mathbf{x}), \beta^{-1}) \quad (48)$$



Plugging in the w_{ML} and β_{ML} into our likelihood we get a **predictive distribution** that allows us to make prediction \hat{y}_{new} for the new data \mathbf{x}_{new} .

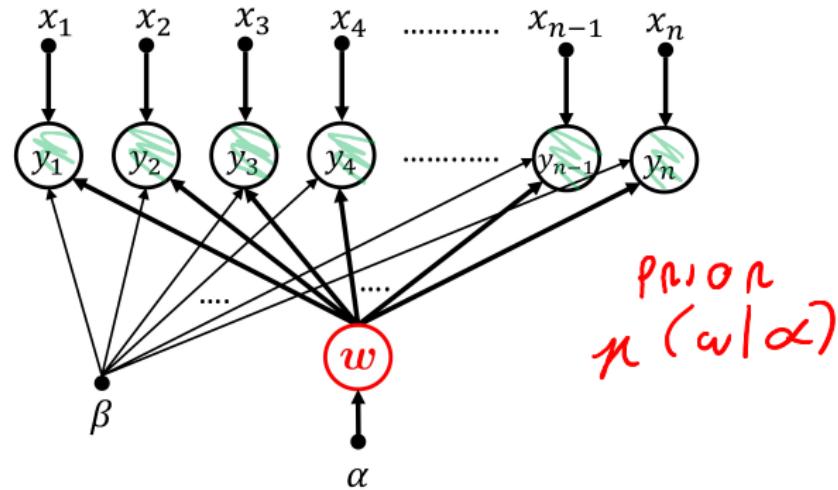
$$p(\hat{y}_{new} | \mathbf{x}_{new}, w_{ML}, \beta_{ML}) = \mathcal{N}(\hat{y}_{new} | \boxed{w_{ML}^T \phi(\mathbf{x}_{new})}, \beta_{ML}^{-1}) \quad (49)$$



Posterior distribution

Recall from the Lecture 3, that ML leads to overfitting (especially, when little training data is available).

Solution - consider the **posterior distribution** instead



Posterior distribution

Recall from the Lecture 3, that ML leads to overfitting (especially, when little training data is available).

Solution - consider the **posterior distribution** instead

$$p(\mathbf{w} \mid \mathbf{X}, \mathbf{y}, \beta, \cdot) = \frac{\underbrace{p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta)}_{\text{likelihood}} \cdot \underbrace{p(\mathbf{w} \mid \cdot)}_{\text{prior}}}{\underbrace{p(\mathbf{X}, \mathbf{y})}_{\text{normalizing constant}}} \quad (50)$$

$$\propto p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta) \cdot p(\mathbf{w} \mid \cdot) \quad (51)$$

Posterior distribution

Recall from the Lecture 3, that ML leads to overfitting (especially, when little training data is available).

Solution - consider the **posterior distribution** instead

$$p(\mathbf{w} \mid \mathbf{X}, \mathbf{y}, \beta, \cdot) = \frac{\underbrace{p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta)}_{\text{likelihood}} \cdot \underbrace{p(\mathbf{w} \mid \cdot)}_{\text{prior}}}{\underbrace{p(\mathbf{X}, \mathbf{y})}_{\text{normalizing constant}}} \quad (50)$$

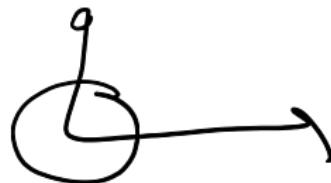
$$\propto p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta) \cdot p(\mathbf{w} \mid \cdot) \quad (51)$$

Connection to the coin flip example

	train data	likelihood	prior	posterior
coin:	$\mathcal{D} = \mathbf{X}$	$p(\mathcal{D} \mid \theta)$	$p(\theta \mid a, b)$	$p(\theta \mid \mathcal{D})$
regr.:	$\mathcal{D} = \{\mathbf{X}, \mathbf{y}\}$	$p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta)$	$p(\mathbf{w} \mid \cdot)$	$p(\mathbf{w} \mid \mathbf{X}, \mathbf{y}, \beta, \cdot)$

How do we choose the prior $p(\mathbf{w} \mid \cdot)$?

Prior for w



We set the prior over w to an isotropic multivariate normal distribution with zero mean

$$p(w | \alpha) = \mathcal{N}(w | \mathbf{0}, \alpha^{-1} \mathbf{I}) = \left(\frac{\alpha}{2\pi} \right)^{\frac{M}{2}} \exp \left(-\frac{\alpha}{2} w^T w \right) \quad (52)$$

where,

α - precision of the distribution

M - number of elements in the vector w

Motivation:

- Higher probability is assigned to small values of w
 \Rightarrow prevents overfitting (recall slide 21)
- Likelihood is also Gaussian - simplified calculations

Maximum a posteriori (MAP)

We are looking for \mathbf{w} that corresponds to the mode of the posterior

$$\mathbf{w}_{\text{MAP}} = \arg \max_{\mathbf{w}} p(\mathbf{w} \mid \mathbf{X}, \mathbf{y}, \alpha, \beta) \quad (53)$$

$$= \arg \max_{\mathbf{w}} \ln p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta) + \ln p(\mathbf{w} \mid \alpha) - \underbrace{\ln p(\mathbf{X}, \mathbf{y})}_{=\text{const}} \quad (54)$$

$$= \arg \min_{\mathbf{w}} -\ln p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta) - \ln p(\mathbf{w} \mid \alpha) \quad (55)$$

Similar to ML, define the MAP error function as negative log-posterior

$$E_{\text{MAP}}(\mathbf{w}) = -\ln p(\mathbf{w} \mid \mathbf{X}, \mathbf{y}, \alpha, \beta) \quad (56)$$

$$= -\ln p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta) - \ln p(\mathbf{w} \mid \alpha) + \text{const} \quad (57)$$

MAP error function

Simplify the error function

$$E_{MAP} = -\ln p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \beta) - \ln p(\mathbf{w} \mid \alpha)$$

$$= \frac{\beta}{2} \sum_{i=1}^N (\mathbf{w}^T \phi(\mathbf{x}_i) - y_i)^2 - \frac{N}{2} \ln \beta + \frac{N}{2} \ln 2\pi$$

$$- \ln \left(\frac{\alpha}{2\pi} \right)^{\frac{M}{2}} + \frac{\alpha}{2} \mathbf{w}^T \mathbf{w}$$

$$= \frac{\beta}{2} \sum_{i=1}^N (\mathbf{w}^T \phi(\mathbf{x}_i) - y_i)^2 + \frac{\alpha}{2} \|\mathbf{w}\|_2^2 + \text{const}$$

$$\propto \underbrace{\frac{1}{2} \sum_{i=1}^N (\mathbf{w}^T \phi(\mathbf{x}_i) - y_i)^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2}_{\text{ridge regression error fn!}} + \text{const}$$

$\Leftrightarrow \alpha = \lambda \cdot \beta$

$$\text{where } \lambda = \frac{\alpha}{\beta}$$

$$= E_{\text{ridge}}(\mathbf{w}) + \text{const}$$

(58)

MAP estimation with Gaussian prior is equivalent to ridge regression!

Predicting for new data

Recall, that

$$y \sim \mathcal{N}(f_{\boldsymbol{w}}(\boldsymbol{x}), \beta^{-1}) \quad (59)$$

Plugging in the $\boldsymbol{w}_{\text{MAP}}$ into our likelihood we get a **predictive distribution** that lets us make prediction \hat{y}_{new} for the new data \boldsymbol{x}_{new} .

$$p(\hat{y}_{new} \mid \boldsymbol{x}_{new}, \boldsymbol{w}_{\text{MAP}}, \beta) = \mathcal{N}(\hat{y}_{new} \mid \boldsymbol{w}_{\text{MAP}}^T \boldsymbol{\phi}(\boldsymbol{x}_{new}), \beta^{-1}) \quad (60)$$

Full Bayesian approach

Why limit ourself to the mode $\boldsymbol{w}_{\text{MAP}}$ of the posterior?

Instead, we can try to estimate the full posterior distribution $p(\boldsymbol{w} \mid \mathcal{D})$ ⁵

$$p(\boldsymbol{w} \mid \mathcal{D}) \propto \overbrace{p(\mathbf{y} \mid \mathbf{X}, \boldsymbol{w}, \beta)}^{\text{likelihood}} \overbrace{p(\boldsymbol{w} \mid \alpha)}^{\text{prior}} \quad (61)$$

Since both the likelihood and the prior are Gaussian, we have a closed form for the posterior! (Conjugate prior)

$$p(\boldsymbol{w} \mid \mathcal{D}) = \mathcal{N}(\boldsymbol{w} \mid \boldsymbol{m}_N, \mathbf{S}_N) \quad (62)$$

with

$$\text{Mean} = \boldsymbol{m}_N = \mathbf{S}_N \left(\mathbf{S}_0^{-1} \boldsymbol{m}_0 + \beta \boldsymbol{\Phi}^T \mathbf{y} \right) \quad (63)$$

$$\text{Covariance} = \mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi} \quad (64)$$

\boldsymbol{m}_0 – prior mean, \mathbf{S}_0 – prior covariance

⁵To avoid clutter, we use $p(\boldsymbol{w} \mid \mathcal{D})$ as a shorthand for $p(\boldsymbol{w} \mid \mathbf{X}, \mathbf{y}, \alpha, \beta)$

Posterior distribution

Posterior parameter distribution

$$(\beta \cdot \Phi^T \Phi)^{-1} \cdot \beta \Phi^T y = (\Phi^T \Phi)^{-1} \cdot \Phi^T y + c$$

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

with

$$\begin{aligned}\text{Mean} = \mathbf{m}_N &= S_N^{-1} (S_0^{-1} \mathbf{m}_0 + \beta \Phi^T y) \\ \text{Covariance} = \mathbf{S}_N^{-1} &= S_0^{-1} + \beta \Phi^T \Phi\end{aligned}$$

\mathbf{m}_0 – prior mean, S_0 – prior covariance

$$(\beta \cdot \Phi^T \Phi)^{-1} \cdot \beta \Phi^T y \quad m_0 = 0$$

Observations

- Since we again have a Gaussian, the MAP solution (i.e. the mode) equals the mean: $\mathbf{w}_{\text{MAP}} = \mathbf{m}_N$.
- In the limit of an infinitely broad prior, $S_0^{-1} \rightarrow 0$, $\mathbf{w}_{\text{MAP}} \rightarrow \mathbf{w}_{\text{ML}}$
- For $N = 0$, i.e. no data points, we get the prior back.

Sequential Bayesian linear regression

Consider following scenarios

- What if the dataset is too large and can't fit into memory all at once?
- What if the data arrives sequentially (in a stream) and has to be processed in an online manner?

Bayesian framework provides a solution!

1. After processing batch of data \mathcal{D}_1 at the first time step $t = 1$, we obtain posterior

$$p(\mathbf{w} \mid \mathcal{D}_1) \propto p(\mathcal{D}_1 \mid \mathbf{w})p(\mathbf{w} \mid \alpha) \quad (65)$$

2. Use the posterior from step t as a prior for step $t + 1$!

$$p(\mathbf{w} \mid \mathcal{D}_2, \mathcal{D}_1) \propto p(\mathcal{D}_2 \mid \mathbf{w})p(\mathcal{D}_1 \mid \mathbf{w})p(\mathbf{w} \mid \alpha) \quad (\text{i.i.d.}) \quad (66)$$

$$\propto p(\mathcal{D}_2 \mid \mathbf{w})p(\mathbf{w} \mid \mathcal{D}_1) \quad (67)$$

Sequential Bayesian linear regression: Example

Bayesian regression for the target values

$$y_i = -0.3 + 0.5x_i + \epsilon, \quad \epsilon \sim \mathcal{N}(0, 0.2^2)$$

Sequential Bayesian linear regression: Example

Bayesian regression for the target values

$$y_i = -0.3 + 0.5x_i + \epsilon, \quad \epsilon \sim \mathcal{N}(0, 0.2^2)$$

To model this, we set $\phi(x) = [\begin{smallmatrix} 1 & x \end{smallmatrix}]$ and thus

$$f_{\boldsymbol{w}}(x) = w_0 + w_1 x$$

Sequential Estimation

The demo shows how the posterior's breadth gets smaller as more and more points are taken into account, and how its mode converges to the optimal (true) values of the weights (white cross).

Sequential Bayesian linear regression: Example

middle column:

$p(w_0, w_1 | \{(x_i, y_i)\}_{i=1}^t)$, for $t > 1$.
the prior $p(w_0, w_1)$ in the first row.

left-hand column:

value of $p(x_t, y_t | w_0, w_1)$
as a function of (w_0, w_1) .

right-hand column:

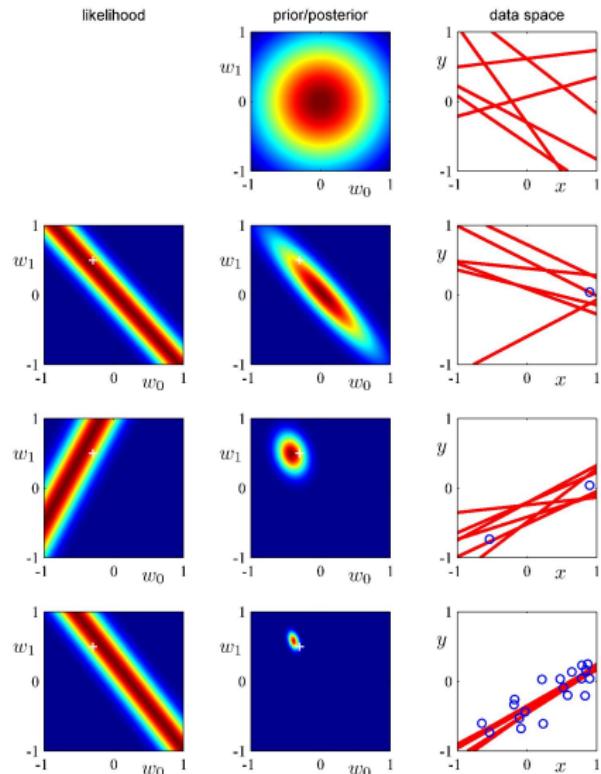
six sample lines from the current
posterior distribution.

blue circles in the right column:

the observed data-points.
up to time t .

white cross:

the true value of (w_0, w_1) .



Posterior predictive distribution

After observing the data $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$, we can compute the full posterior distribution $p(\mathbf{w} | \mathcal{D})$.

Usually, what we are actually interested in is the prediction \hat{y}_{new} for a new data point \mathbf{x}_{new} - the model parameters \mathbf{w} are just a means to achieve this.

The **posterior predictive distribution** is computed as

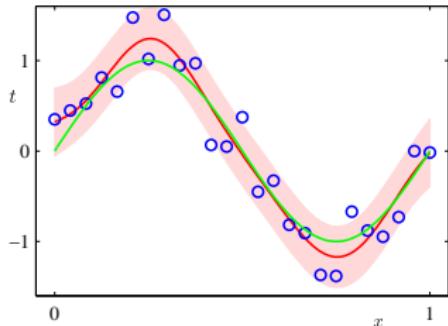
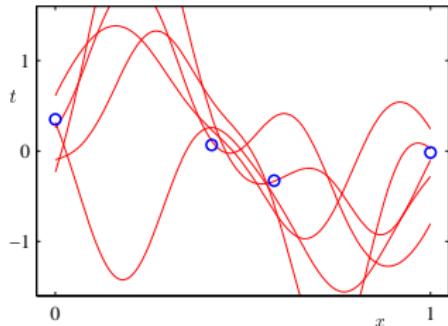
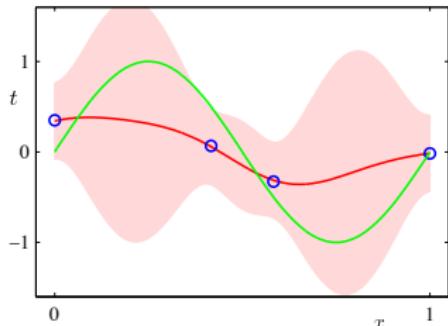
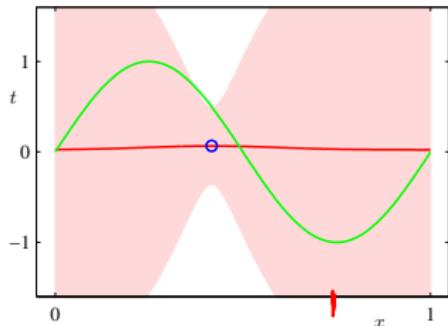
$$p(\hat{y}_{new} | \mathbf{x}_{new}, \mathcal{D}) = \int p(\hat{y}_{new}, \mathbf{w} | \mathbf{x}_{new}, \mathcal{D}) d\mathbf{w} \quad (68)$$

$$= \int p(\hat{y}_{new} | \mathbf{x}_{new}, \mathbf{w}) p(\mathbf{w} | \mathcal{D}) d\mathbf{w} \quad (69)$$

$$= \mathcal{N}(\hat{y}_{new} | \mathbf{m}_N^T \phi(\mathbf{x}_{new}), \beta^{-1} + \phi(\mathbf{x}_{new})^T \mathbf{S}_N \phi(\mathbf{x}_{new})) \quad (70)$$



Example of posterior predictive distribution



Summary

- Optimization-based approaches to regression have probabilistic interpretations
 - Least squares regression \iff Maximum likelihood (Slide 45)
 - Ridge regression \iff Maximum a posteriori (Slide 50)
- Even nonlinear dependencies in the data can be captured by a model linear w.r.t. weights w (Slide 14)
- Penalizing large weights helps to reduce overfitting (Slide 21)
- Full Bayesian can be processed sequentially - same result (Slide 55)