

Machine learning

Linear Regression

Hamid Beigy

Sharif University of Technology

November 15, 2021





1. Introduction
2. Linear regression
3. Model selection
4. Sample size
5. Probabilistic Linear Regression
6. Overfitting
7. Regularization
8. Maximum a posteriori and regularization
9. Geometric Interpretation
10. Sequential learning
11. Multiple outputs regression
12. Bias-Variance trade-off

Introduction



- ▶ In regression, $c(x)$ is a continuous function. Hence the training set is in the form of

$$S = \{(x_1, t_1), (x_2, t_2), \dots, (x_N, t_N)\}, t_k \in \mathbb{R}.$$

- ▶ When there is no noise, the task is interpolation and our goal is to find a function $f(x)$ that passes through these points. Hence

$$t_k = f(x_k) \quad \forall k = 1, 2, \dots, N$$

- ▶ In polynomial interpolation, given N points, we find $(N - 1)$ st degree polynomial to predict the output for any x .
- ▶ If x is outside of the range of the training set, the task is called extrapolation.
- ▶ In regression, there is noise added to the output of the unknown function.

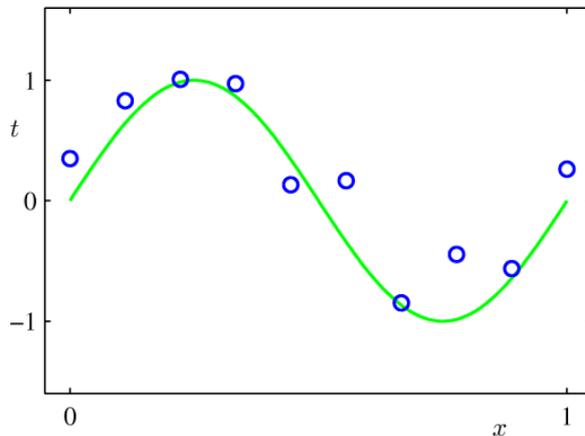
$$t_k = f(x_k) + \epsilon \quad \forall k = 1, 2, \dots, N$$

$f(x_k) \in \mathbb{R}$ is the unknown function and ϵ is the random noise.



- ▶ In regression, there is noise added to the output of the unknown function.

$$t_k = f(x_k) + \epsilon \quad \forall k = 1, 2, \dots, N$$



- ▶ The explanation for the noise is that there are extra hidden variables that we cannot observe.

$$t_k = f^*(x_k, z_k) + \epsilon \quad \forall k = 1, 2, \dots, N$$

z_k denotes hidden variables

Linear regression



- ▶ Our goal is to approximate the output by function $g(x)$.
- ▶ The empirical error on the training set S is measured using loss/error/cost function.
 - ▶ Squared Error from target $E_E(g(x_i)|S) = (t_i - g(x_i))^2$.
 - ▶ Linear error from target $E_E(g(x_i)|S) = |t_i - g(x_i)|$.
 - ▶ Mean square error from target $E_E(g(x)|S) = \frac{1}{N} \sum_{i=1}^N (t_i - g(x_i))^2$.
 - ▶ Sum of square error from target $E_E(g(x)|S) = \frac{1}{2} \sum_{i=1}^N (t_i - g(x_i))^2$.
- ▶ The aim is to find $g(\cdot)$ that minimizes the empirical error.
- ▶ We assume that a hypothesis class for $g(\cdot)$ with a small set of parameters.
- ▶ Assume that $g(x)$ is linear

$$g(x) = w_0 + w_1x_1 + w_2x_2 + \dots + w_Dx_D$$



- ▶ In the linear regression, when $D = 1$

$$g(x) = w_0 + w_1 x$$

- ▶ Parameters w_0 and w_1 should minimize the empirical error

$$E_E(w_0, w_1 | S) = E_E(g(x) | S) = \frac{1}{2} \sum_{i=1}^N [t_k - (w_0 + w_1 x_k)]^2$$

- ▶ This error function is a quadratic function of W and its derivative is linear w.r.t W . Its minimization has a unique solution denoted by W^* .
- ▶ Taking derivative of error w.r.t w_0 and w_1 and setting equal to zero

$$\begin{aligned} w_0 &= \bar{t} - w_1 \bar{x} \\ w_1 &= \frac{\sum_k t_k x_k - \bar{x} \bar{t} N}{\sum_k (x_k)^2 - N(\bar{x})^2} \\ \bar{t} &= \frac{\sum_{k=1}^N t_k}{N} \\ \bar{x} &= \frac{\sum_{k=1}^N x_k}{N} \end{aligned}$$



- ▶ When the input variables form a D -dimensional vector, the linear regression model is

$$g(x) = w_0 + w_1x_1 + w_2x_2 + \dots + w_Dx_D.$$

- ▶ Parameters w_0, w_1, \dots, w_D should minimize the empirical error

$$E_E(g(x)|S) = \frac{1}{2} \sum_{i=1}^N (t_i - g(x_i))^2.$$

- ▶ Taking derivative of error w.r.t w s and setting equal to zero,



- ▶ Taking derivative of error w.r.t w s and setting equal to zero,

$$\begin{aligned} \sum_{k=1}^N t_k &= Nw_0 + w_1 \sum_{k=1}^N x_{k1} + w_2 \sum_{k=1}^N x_{k2} + \dots + w_D \sum_{k=1}^N x_{kD} \\ \sum_{k=1}^N x_{k1} t_k &= w_0 \sum_{k=1}^N x_{k1} + w_1 \sum_{k=1}^N (x_{k1})^2 + w_2 \sum_{k=1}^N x_{k1} x_{k2} + \dots + w_D \sum_{k=1}^N x_{k1} x_{kD} \\ \sum_{k=1}^N x_{k2} t_k &= w_0 \sum_{k=1}^N x_{k2} + w_1 \sum_{k=1}^N x_{k1} x_{k2} + w_2 \sum_{k=1}^N (x_{k2})^2 + \dots + w_D \sum_{k=1}^N x_{k2} x_{kD} \\ &\vdots \\ \sum_{k=1}^N x_{kD} t_k &= w_0 \sum_{k=1}^N x_{kD} + w_1 \sum_{k=1}^N x_{k1} x_{kD} + w_2 \sum_{k=1}^N x_{k2} x_{kD} + \dots + w_D \sum_{k=1}^N (x_{kD})^2 \end{aligned}$$



- ▶ Define the following vectors and Matrix

- ▶ Data matrix

$$X = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1D} \\ 1 & x_{21} & x_{22} & \dots & x_{2D} \\ \vdots & & & & \\ 1 & x_{N1} & x_{N2} & \dots & x_{ND} \end{bmatrix}$$

- ▶ The k^{th} input vector

$$X_k = (1, x_{k1}, x_{k2}, \dots, x_{kD})^T$$

- ▶ The weight vector

$$W = (w_0, w_1, w_2, \dots, w_D)^T$$

- ▶ The target vector

$$t = (t_1, t_2, t_3, \dots, t_N)^T$$



- ▶ The empirical error is equal to

$$E_E(g(x)|S) = \frac{1}{2} \sum_{k=1}^N (t_k - W^\top X_k)^2.$$

- ▶ The gradient of $E_E(g(x)|S)$ is

$$\begin{aligned} \nabla_W E_E(g(x)|S) &= \sum_{k=1}^N (t_k - W^\top X_k) X_k^\top \\ &= \sum_{k=1}^N t_k X_k^\top - W^\top \sum_{k=1}^N X_k X_k^\top = 0 \end{aligned}$$

- ▶ Solving for W , we obtain

$$W^* = (X^\top X)^{-1} X^\top t$$

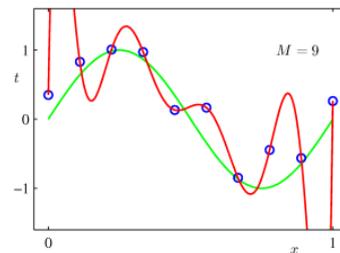
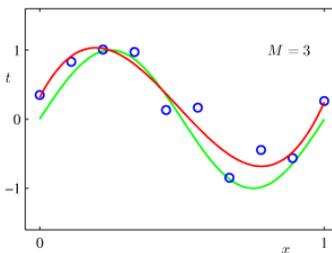
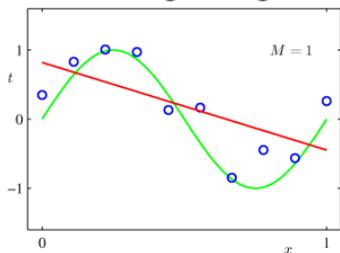
- ▶ When $X^\top X$ is invertible, the problem has a unique solution.
- ▶ When $X^\top X$ is not invertible, the pseudo inverse is used and the problem has several solution.



- ▶ If the linear model is too simple, the model can be a polynomial (a more complex hypothesis set)

$$g(x) = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M.$$

- ▶ M is the order of the polynomial.
- ▶ Choosing the right value of M is called **model selection**.

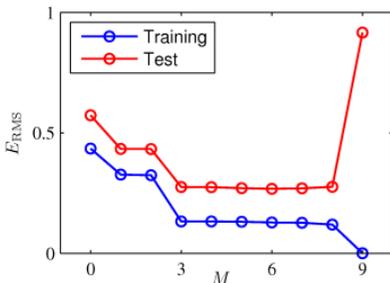


- ▶ For $M = 1$, we have a too general model
- ▶ For $M = 9$, we have a too specific model

Model selection



- ▶ The goal of model selection is to achieve good generalization by making accurate predictions for new data.
- ▶ The generalization ability of a model is measured by a separate test data generated using exactly the same process used for generating training data.
- ▶ The model is chosen using a validation data set.



- ▶ Two models sometimes are compared using root mean square (RMS) error.

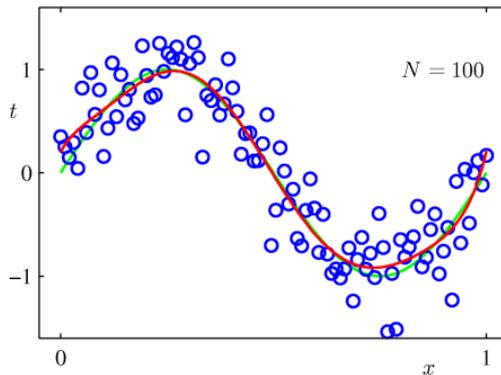
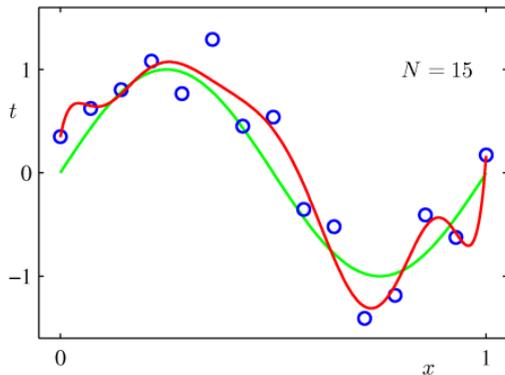
$$E_{RMS} = \sqrt{2E_E(W^*|S)/N}$$

- ▶ This allows comparison on different sizes of data sets.

Sample size



- ▶ For a given model complexity, the over-fitting problem become less severe as the size of the data set increases.





- ▶ We can extend the class of models by considering linear combinations of fixed nonlinear functions of the input variables, of the form

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

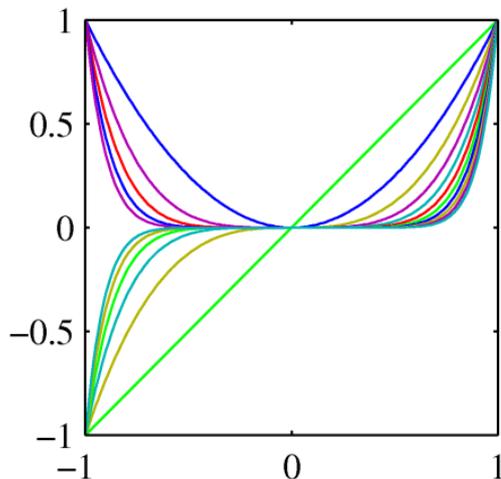
- ▶ $\phi_j(x)$ are known as **basis functions**.
 - ▶ M is total number of parameters.
 - ▶ w_0 is called **bias parameter**.
- ▶ Usually a dummy basis function $\phi_0(x) = 1$ is used

$$g(x) = \sum_{j=0}^M w_j \phi_j(x) = W^T \Phi(x)$$

- ▶ $W = (w_0, w_1, \dots, w_{M-1})^T$ and $\Phi = (\phi_0, \phi_1, \dots, \phi_{M-1})^T$.



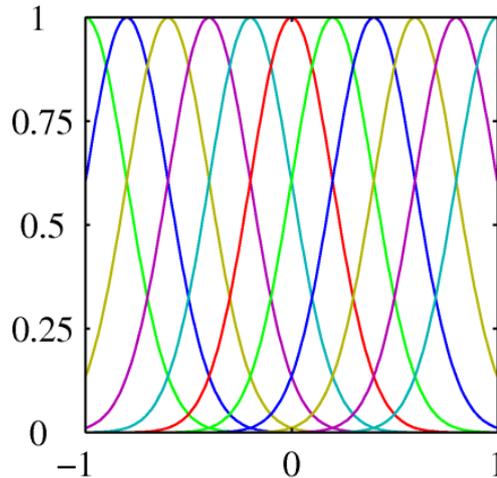
- ▶ In pre-processing phase, the features can be expressed in terms of the basis functions $\{\phi_j(x)\}$.
- ▶ Examples of basis functions
 - ▶ Polynomial basis function $\phi_j(x) = x^j$.





► Examples of basis functions

- Gaussian basis function $\phi_j(x) = \exp\left\{-\frac{(x-\mu_j)^2}{2s^2}\right\}$.
- μ_j is location of the basis function.
- s is the spatial scale of the basis function.

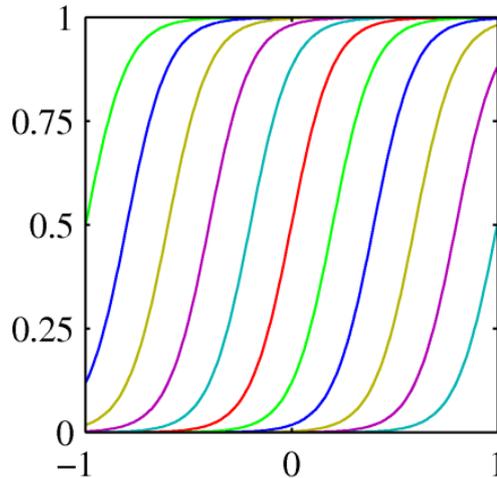




- ▶ Examples of basis functions

- ▶ Logistic basis function $\phi_j(x) = \sigma\left(\frac{x-\mu_j}{s}\right)$.

- ▶ $\sigma(a) = \frac{1}{1+\exp(-a)}$.



- ▶ Fourier basis function

- ▶ Wavelets basis function



- ▶ The empirical error is equal to

$$E_E(g(x)|S) = \frac{1}{2} \sum_{k=1}^N (t_k - W^\top \Phi(X_k))^2.$$

- ▶ The gradient of $E_E(g(x)|S)$ is

$$\begin{aligned} \nabla_W E_E(g(x)|S) &= \sum_{k=1}^N (t_k - W^\top \Phi(X_k)) \Phi(X_k)^\top \\ &= \sum_{k=1}^N t_k \Phi(X_k)^\top - W^\top \sum_{k=1}^N \Phi(X_k) \Phi(X_k)^\top = 0 \end{aligned}$$

- ▶ Solving for W , we obtain $W^* = (\Phi^\top \Phi)^{-1} \Phi^\top t$.

$$\Phi = \begin{bmatrix} \phi_0(x_1) & \phi_1(x_1) & \phi_2(x_1) & \dots & \phi_{M-1}(x_1) \\ \phi_0(x_2) & \phi_1(x_2) & \phi_2(x_2) & \dots & \phi_{M-1}(x_2) \\ \vdots & & & & \vdots \\ \phi_0(x_N) & \phi_1(x_N) & \phi_2(x_N) & \dots & \phi_{M-1}(x_N) \end{bmatrix}$$



- ▶ The quantity $\Phi^\dagger = (\Phi^\top \Phi)^{-1} \Phi^\top$ is known as **Moore-Penrose pseudo-inverse**.
- ▶ The bias value (w_0)

$$w_0 = \frac{1}{N} \sum_{k=1}^N \left\{ t_k - \sum_{j=1}^{M-1} w_j \phi_j(X_k) \right\}.$$

and equals to the difference between target values and the weighted sum of the basis function values.

Probabilistic Linear Regression



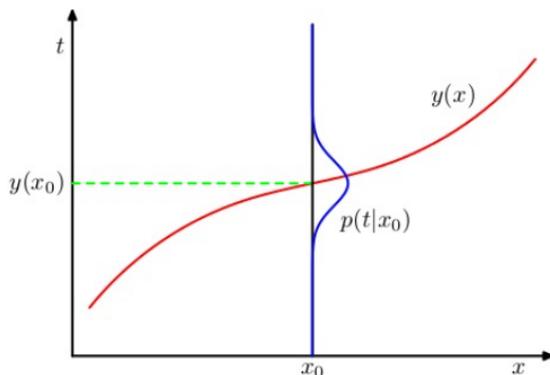
- ▶ Given N training examples $S = \{(x_1, t_1), \dots, (x_N, t_N)\}$, $x_k \in \mathbb{R}^D$, output $t_k \in \mathbb{R}$.
- ▶ In **probabilistic linear regression**, outputs t_k 's are generated from a **probabilistic model**.
- ▶ We assume that the target variable t is given by a deterministic function $f(x)$ with additive Gaussian noise so that

$$t = f(x) + \epsilon$$

where $\epsilon \sim \mathcal{N}(0, \sigma^2)$ is a zero mean Gaussian random variable with precision (inverse variance) β .

- ▶ Thus, t is also a Gaussian random variable with mean t and precision β .

$$p(t|x) = \mathcal{N}(t|f(x), \beta^{-1}).$$





- ▶ Since training examples are drawn i.i.d from the same distribution, the likelihood equals

$$L(W) = \prod_{k=1}^N \mathcal{N}(t_k | g(x_k; w), \beta^{-1}) = \prod_{k=1}^N \frac{\sqrt{\beta}}{\sqrt{2\pi}} \exp \left\{ -\frac{\beta (t_k - g(x_k))^2}{2} \right\}.$$

- ▶ The log-likelihood becomes

$$\ln L(W) = \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \frac{\beta}{2} \sum_{k=1}^N (t_k - g(x_k))^2.$$

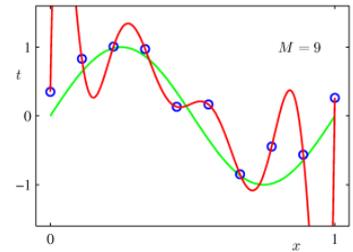
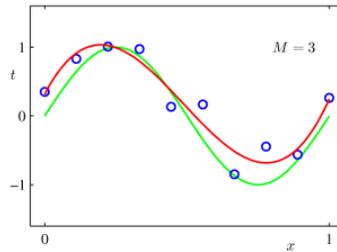
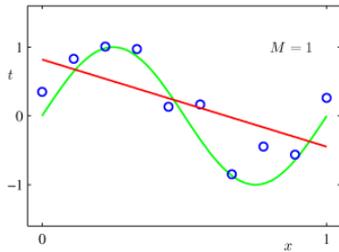
- ▶ Clearly, maximizing $\ln L(W)$ equals to minimizing $-\ln L(W)$.
- ▶ Removing constant terms yields

$$W_{ML} = \underset{W}{\operatorname{argmin}} \frac{1}{2} \sum_{k=1}^N (t_k - g(x_k))^2.$$

Overfitting



- ▶ Over-fitting occurs when a model begins to **memorize** training data rather than **learning** to generalize from a trend.
 - ▶ Too many parameters relative to the amount of training data
 - ▶ For example, an order- N polynomial can be exact fit to $N + 1$ data points.





- ▶ Ways of detecting/avoiding over-fitting.
 - ▶ Use more training data
 - ▶ Evaluate on a parameter tuning set
 - ▶ Regularization
 - ▶ Take a Bayesian approach
- ▶ In a Linear Regression model, overfitting is characterized by large parameters.
- ▶ As M increase, the magnitude of the coefficients typically gets larger.

	$M = 0$	$M = 1$	$M = 6$	$M = 9$
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^*			17.37	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

Regularization



- ▶ Introduce a penalty term for the size of the weights.

- ▶ Unregularized Regression

$$E_E(g(x)|S) = \frac{1}{2} \sum_{i=1}^N (t_i - g(x_i))^2.$$

- ▶ Regularized Regression

$$E_E(g(x)|S) = \frac{1}{2} \sum_{i=1}^N (t_i - g(x_i))^2 + \lambda \Omega(W).$$

- ▶ L_2 -Regularization or Ridge Regularization

$$E_E(g(x)|S) = \frac{1}{2} \sum_{i=1}^N (t_i - g(x_i))^2 + \frac{\lambda}{2} \|W\|^2.$$

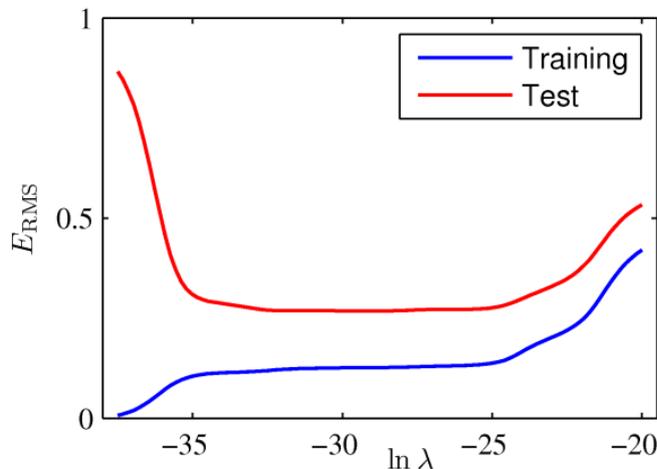
- ▶ Large λ leads to higher complexity penalization



- ▶ Least squares regression with L_2 -regularization

$$\begin{aligned}\nabla_W E_E(g(x)|S) &= \nabla_W \left(\frac{1}{2} \sum_{i=1}^N (t_i - g(x_i))^2 + \frac{\lambda}{2} \|W\|^2 \right) \\ (X^T X + \lambda I)W &= X^T t \\ W^* &= (X^T X + \lambda I)^{-1} X^T t\end{aligned}$$

- ▶ Graph of the root-mean-square error (RMSE) versus $\ln \lambda$ for the $M = 9$ polynomial.





- ▶ L_2 -Regularization has **Closed form** solution and can be solved in polynomial time

$$E_E(g(x)|S) = \frac{1}{2} \sum_{i=1}^N (t_i - g(x_i))^2 + \frac{\lambda}{2} \|W\|^2.$$

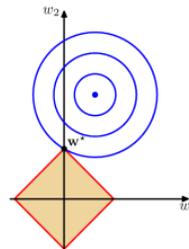
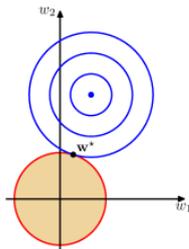
- ▶ L_1 -Regularization can be **approximated** in polynomial time

$$E_E(g(x)|S) = \frac{1}{2} \sum_{i=1}^N (t_i - g(x_i))^2 + \lambda \|W\|_1.$$

- ▶ L_0 -Regularization is **NP-complete optimization** problem

$$E_E(g(x)|S) = \frac{1}{2} \sum_{i=1}^N (t_i - g(x_i))^2 + \lambda \sum_{j=1}^{M-1} \delta(w_j \neq 0).$$

The L_0 -norm represents the optimal subset of features needed by a Regression model.

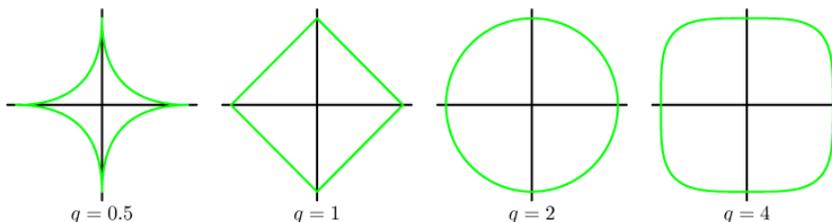




- ▶ A more general regularizer is sometimes used, for which the regularized error takes the form

$$E_E(g(x)|S) = \frac{1}{2} \sum_{i=1}^N (t_i - g(x_i))^2 + \frac{\lambda}{2} \sum_{j=1}^{M-1} |w_j|^q.$$

- ▶ When $q = 2$, it is called **Ridge regularizer**
- ▶ When $q = 1$, it is called **Lasso regularizer**



Maximum a posteriori and regularization



- ▶ Assume that W is

$$p(W) = \mathcal{N}(0, \sigma_0^2 I_D).$$

- ▶ I_D denotes the $D \times D$ identity matrix.
- ▶ This is equivalent to assume that the prior selects each component of W independently from a $\mathcal{N}(0, \sigma_0^2)$.
- ▶ This prior density can be written as

$$p(W) = \frac{1}{(2\pi)^{D/2} \sigma_0^D} \exp \left\{ -\frac{1}{2\sigma_0^2} \|W\|_2^2 \right\}.$$

- ▶ Assume that noise precision is known.
- ▶ The posterior density of W given set S takes the form

$$L(W) \propto \exp \left(-\frac{1}{2\sigma_0^2} \|W\|_2^2 \right) \times \prod_{k=1}^N \exp \left(-\frac{\beta}{2} (t_k - g(x_k))^2 \right).$$



- ▶ The log-likelihood becomes

$$\ln L(W) = -\frac{1}{2\sigma_0^2} \|W\|_2^2 - \frac{\beta}{2} \sum_{k=1}^N (t_k - g(x_k))^2 + \text{const.}$$

- ▶ Clearly, maximizing $\ln L(W)$ equals to minimizing $-\ln L(W)$.
- ▶ Removing constant terms yields

$$W_{MAP} = \underset{W}{\operatorname{argmin}} \frac{1}{2} \sum_{k=1}^N (t_k - g(x_k))^2 + \frac{1}{2\sigma_0^2\beta} \|W\|_2^2.$$

- ▶ Setting $\lambda = \frac{1}{\sigma_0^2\beta}$, results in L_2 -regularization.

Geometric Interpretation



- ▶ Consider an N -dimensional space whose axes are given by t_j .
- ▶ $t = (t_1, t_2, \dots, t_N)^\top$ is a vector in this space.
- ▶ Each basis function $\phi_j(x_n)$ is evaluated at the N data points can also be represented as a vector in the same space.
- ▶ $\phi_j(x_n)$ corresponds to j^{th} column of Φ , i.e. $\phi_j(x_n) = (\phi_j(x_1), \phi_j(x_2), \dots, \phi_j(x_N))^\top$.
- ▶ Assume that $M < N$, then M vectors $\phi_j(x_n)$ will span a linear subspace of dimensionality M which embedded in N dimensions.
- ▶ Let y be an N -dimensional vector whose n th element is given by $g(x_n)$.
- ▶ Hence, there exists some weight vector W such that

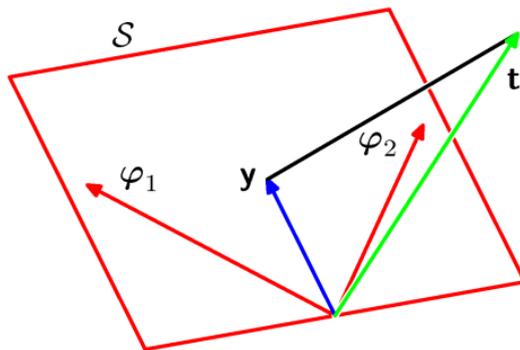
$$\hat{t} = \begin{bmatrix} \hat{t}_1 \\ \hat{t}_2 \\ \vdots \\ \hat{t}_N \end{bmatrix} = \begin{bmatrix} \phi_0(x_1) & \phi_1(x_1) & \phi_2(x_1) & \dots & \phi_{M-1}(x_1) \\ \phi_0(x_2) & \phi_1(x_2) & \phi_2(x_2) & \dots & \phi_{M-1}(x_2) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \phi_0(x_N) & \phi_1(x_N) & \phi_2(x_N) & \dots & \phi_{M-1}(x_N) \end{bmatrix} \begin{bmatrix} W_0 \\ W_1 \\ \vdots \\ W_N \end{bmatrix} = \Phi W$$

- ▶ We seek a $\hat{t} = y \in \mathbb{R}^N$ that lies in this linear subspace and is close as possible to t , i.e, we want to find

$$\hat{y} = \underset{y \in \text{subspace}(\Phi(X))}{\text{argmin}} \|t - y\|_2$$



- ▶ The least square solution for W corresponds to that choice of y that lies in this subspace and is close to t .



- ▶ In order to minimize the norm of the residual error $(t - y)$, we want **the residual error vector be orthogonal to every column of $\Phi(X)$** so that $\phi_j(x_n)^\top (t - y) = 0$, for $j = 0, 1, 2, \dots, M - 1$. (**why?**)

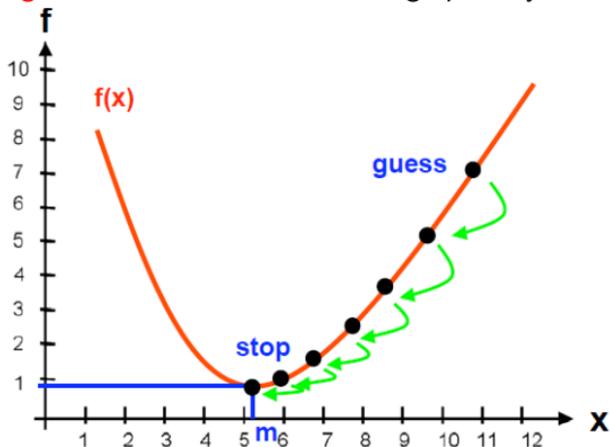
$$\begin{aligned} \phi_j(x_n)^\top (t - y) = 0 &\Rightarrow \Phi(X)^\top (t - \Phi(X)W) = 0 \\ &\Rightarrow W = (\Phi(X)^\top \Phi(X))^{-1} \Phi(X)^\top t. \end{aligned}$$

- ▶ Hence, our projected value corresponds to an orthogonal projection of t onto column space of this subspace.

Sequential learning



- ▶ Batch techniques, such as the maximum likelihood solution, which involve processing the entire training set in one go, can be computationally costly for large data sets.
- ▶ When **data set is sufficiently large**, it may be worthwhile to use sequential algorithms.
- ▶ In sequential algorithms, the data points are considered one at a time, and the model parameters updated after each such presentation.
- ▶ We want to choose W so as to minimize $E_E(W|S)$.
- ▶ The algorithm starts with some **initial guess** for W and **repeatedly updates** W to make $E_E(W|S)$ **smaller** until hopefully converges to a value of W that minimizes $E_E(W|S)$.
- ▶ This method is called **gradient descent** and is shown graphically for 1-dimensional problem.





- ▶ Using the gradient descent rule, we have

$$W_j^{(k+1)} = W_j^{(k)} - \eta \nabla_W E_E(W|S).$$

- ▶ η called **learning rate** and controls the step size of movement.
- ▶ There are two methods for updating weights
 - ▶ **Batch gradient descent method**

$$W_j^{(k+1)} = W_j^{(k)} - \eta \nabla_W E_E(W|S).$$

- ▶ **Stochastic gradient descent method**

$$\begin{aligned} W_j^{(k+1)} &= W_j^{(k)} - \eta \nabla_W E_E(W|(x_n, t_n)) \\ &= W_j^{(k)} + \eta (t_n - g(x_n)) \Phi(x_n) \end{aligned}$$

- ▶ This rule is called **least-mean-squares (LMS) algorithm** and also known as the **Widrow-Hoff learning rule**.

Multiple outputs regression



- ▶ We have considered the case of a single target variable t .
- ▶ In some applications, we may wish to predict $K > 1$ target variables, which we denote collectively by the target vector \mathbf{t} .
- ▶ This could be done by introducing a different set of basis functions for each component of \mathbf{t} , leading to multiple, independent regression problems.
- ▶ A more interesting, and more common, approach is to use the same set of basis functions to model all of the components of the target vector so that

$$g(x) = W^T \phi(x).$$

- ▶ $g(x)$ is a K -dimensional column vector.
- ▶ W is an $M \times K$ matrix of parameters.
- ▶ $\Phi(x)$ is an M -dimensional column vector with elements $\phi_j(x)$, with $\phi_0(x) = 1$.
- ▶ Suppose we take the conditional distribution of the target vector to be an isotropic Gaussian of the form

$$p(\mathbf{t}|\mathbf{x}, \beta) = \mathcal{N}(\mathbf{t}|\mathbf{W}^T \phi(\mathbf{x}), \beta^{-1}\mathbf{I}).$$

- ▶ If we have a set of observations t_1, \dots, t_N , we can combine these into a matrix T of size $N \times K$ such that the n^{th} row is given by t_n^T .
- ▶ We can also combine the input vectors x_1, \dots, x_N into a matrix X .



- ▶ Since training examples are drawn i.i.d from the same distribution, the likelihood equals

$$\begin{aligned}\ln L(W) = \ln p(T|W, \beta) &= \sum_{n=1}^N \mathcal{N}(\mathbf{t} | \mathbf{W}^\top \phi(\mathbf{x}), \beta^{-1} \mathbf{I}) \\ &= \frac{NK}{2} \ln \left(\frac{\beta}{2\pi} \right) - \frac{\beta}{2} \sum_{n=1}^N \|t_n - \mathbf{W}^\top \phi(x_n)\|^2\end{aligned}$$

- ▶ We can maximize this function with respect to W , giving

$$W_{ML} = (\Phi^\top \Phi)^{-1} \Phi^\top T.$$

- ▶ If we examine this result for each target variable t_k , we have

$$W_k = (\Phi^\top \Phi)^{-1} \Phi^\top \mathbf{t}_k.$$

\mathbf{t}_k is an N -dimensional column vector with components t_{nk} for $n = 1, 2, \dots, N$.

- ▶ The solution to the regression problem decouples between the different target variables, and we need only compute a single pseudo-inverse matrix, which is shared by all of the vectors W_k .

Bias-Variance trade-off



- ▶ In regression, random noise ϵ is added to the output of the unknown function $f(x)$.

$$t_k = f(x_k) + \epsilon \quad \forall k = 1, 2, \dots, N$$

- ▶ Given a set of training examples,

$$S = \{(x_1, t_1), (x_2, t_2), \dots, (x_N, t_N)\}, t_k \in \mathbb{R}.$$

we fit a hypothesis $g(x)$ to the data to minimize the following error function.

$$E_E(g(x)|S) = \frac{1}{2} \sum_{i=1}^N (t_i - g(x_i))^2.$$

- ▶ Now, given a new data point x , we would like to understand the expected prediction error

$$\mathbb{E} [(t - g(x))^2].$$

- ▶ Assume that x generated by the same process as the training set. We decompose $\mathbb{E} [(t - g(x))^2]$ into **bias**, **variance**, and **noise**.



- ▶ From definition of variance we have

$$\mathbb{E} [(x - \mu_x)^2] = \mathbb{E} [x^2] - \mu_x^2$$

- ▶ Hence we have

$$\mathbb{E} [x^2] = \mathbb{E} [(x - \mu_x)^2] + \mu_x^2$$

- ▶ The error for new input x can be decomposed as

$$\begin{aligned}\mathbb{E} [(t - g(x))^2] &= \mathbb{E} [(t - f(x) + f(x) - g(x))^2] \\ &= \mathbb{E} [(t - f(x))^2] + \mathbb{E} [(f(x) - g(x))^2] + 2\mathbb{E} [(t - f(x))(f(x) - g(x))] \\ &= \mathbb{E} [(t - f(x))^2] + \mathbb{E} [(f(x) - g(x))^2]\end{aligned}$$

- ▶ The final term is zero. (Why? Show it.)
- ▶ Since $t = f(x) + \epsilon$, Hence $\epsilon = t - f(x)$, and we have

$$\begin{aligned}\mathbb{E} [(t - g(x))^2] + \mathbb{E} [(f(x) - g(x))^2] &= \mathbb{E} [\epsilon^2] + \mathbb{E} [(f(x) - g(x))^2] \\ &= \text{Var}(\epsilon) + \mathbb{E} [(f(x) - g(x))^2]\end{aligned}$$



- ▶ Now consider the last term, we have

$$\begin{aligned}\mathbb{E} [(f(x) - g(x))^2] &= \mathbb{E} [(f(x) - \mathbb{E}[g(x)] + \mathbb{E}[g(x)] - g(x))^2] \\ &= \mathbb{E} [(f(x) - \mathbb{E}[g(x)])^2] + \mathbb{E} [(\mathbb{E}[g(x)] - g(x))^2] \\ &\quad + 2\mathbb{E} [(f(x) - \mathbb{E}[g(x)])(\mathbb{E}[g(x)] - g(x))]\end{aligned}$$

- ▶ The first term is **the squared bias**.
- ▶ The second term is the **variance** and equals to $(\mathbb{E}[g(x)] - g(x))^2 = \mathbb{E}[g(x)^2] - \mathbb{E}[g(x)]^2$.
- ▶ The final term is **zero**. (**Why? Show it.**)
- ▶ Hence, we have

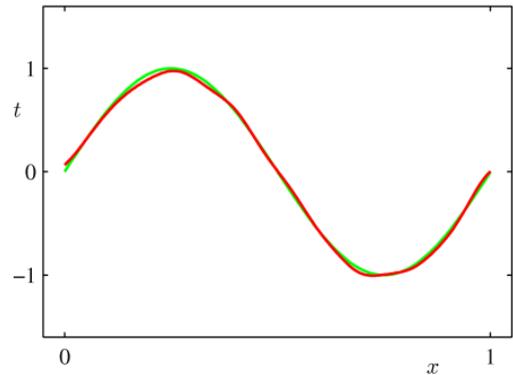
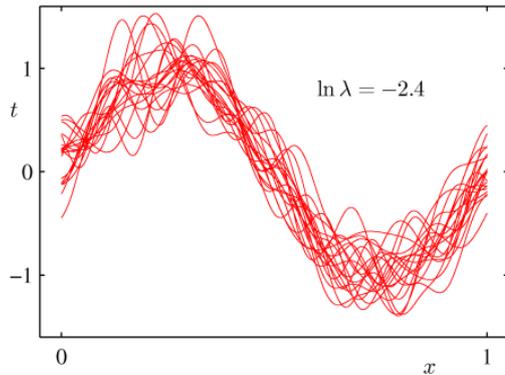
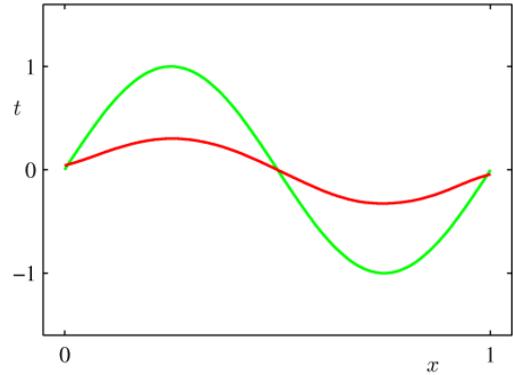
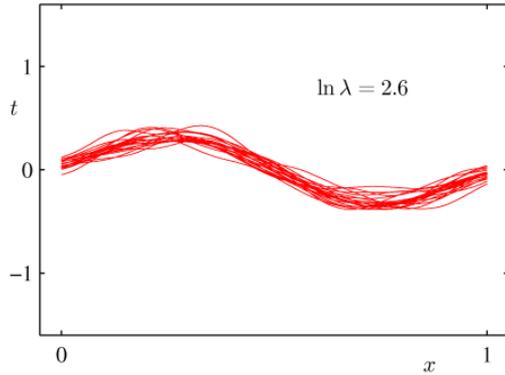
$$\begin{aligned}\mathbb{E} [(t - g(x))^2] &= \text{Bias}^2(g(x)) + \text{Var}(g(x)) + \text{Var}(\epsilon) \\ &= \text{Bias}^2(g(x)) + \text{Var}(g(x)) + \beta^{-1}\end{aligned}$$



- ▶ In summary, we have

$$\mathbb{E} [(t - g(x))^2] = \text{Bias}^2(g(x)) + \text{Var}(g(x)) + \sigma^2$$

- ▶ The first term describes the average error of $g(x)$.
- ▶ The second term quantifies how much $g(x)$ deviates from one training set S to another one. This depends on both the **estimator** and the **training set**. This term is consequence of over-fitting.
- ▶ The last term is the variance of the added noise. This error cannot be removed no matter what estimator we use. Note that the variance of the noise can not be minimized.
- ▶ In context of polynomial regression, **as M increases**, a small changes in the data sets causes a greater change in the fitted function; thus **variance increases**.
- ▶ Our goal is to minimize the expected loss. There is a trade-off between bias and variance.
 - ▶ Very flexible models have low bias and high variance.
 - ▶ Relative rigid models have high bias and low variance.
- ▶ The model with optimal predictive capability is one that leads to the best balance between bias and variance.
- ▶ If there is bias, there is under-fitting. **why?**
- ▶ If there is variance, there is over-fitting. **why?**





1. Chapter 3 of [Pattern Recognition and Machine Learning Book](#) (Bishop 2006).
2. Chapter 7 of [Machine Learning: A probabilistic perspective](#) (Murphy 2012).
3. Chapter 11 of [Probabilistic Machine Learning: An introduction](#) (Murphy 2022).



-  Bishop, Christopher M. (2006). *Pattern Recognition and Machine Learning*. Springer-Verlag.
-  Murphy, Kevin P. (2012). *Machine Learning: A Probabilistic Perspective*. The MIT Press.
-  – (2022). *Probabilistic Machine Learning: An introduction*. The MIT Press.

