

Lecture 1. Introduction and Linear Models for Regression

Bao Wang

Department of Mathematics

Scientific Computing and Imaging Institute

University of Utah

Math 5750/6880, Fall 2021

Topics of this course

1. Linear models for regression and classification.
2. First-order optimization methods for data science. (GD, Proximal GD, SGD, Acceleration)
3. Neural networks (CNN, RNN, GNN, Transformers).
4. Clustering and dimension reduction (k -means, spectral clustering, PCA, CS, ...).
5. Nonparametric models. (kNN,...)
6. Basics of learning theory. (PAC, NTK, ...)

Basic Information

Instructor: Bao Wang

Meeting Time: TuTh 12:25pm - 1:45pm, WEB L120

Office Hours: WF 2:00pm - 3:30pm.

Email: bwang@math.utah.edu

References

Course webpage:

Math 5750: <https://utah.instructure.com/courses/718430>

Math 6880: <https://utah.instructure.com/courses/723883>

References:

1. Zhang et al., Dive into Deep Learning
2. Shai Shalev-Shwartz and Shai Ben-David, Understanding Machine Learning: From Theory to Algorithms
3. Christopher Bishop, Pattern Recognition and Machine Intelligence

Grades

Lecture notes scribe 20%: We will present the lecture using slides and each team needs to scribe each lecture and submit one copy of the lecture notes in PDF along with the latex source. We will provide a latex template.

Course projects 50%: There will be three projects in total, and the lowest score will be dropped. For each project you will need to write a report and submit it along with the source code.

Final project 30%: Each team needs to propose one project related to the course material and conduct reasonably in-depth research. Each team will present their final project.

Course grade:

P	90	85	80	77	70	68	66	50	45	40
Grade	A	A-	B+	B	B-	C+	C	C-	D+	D

Team: 3-5 students in a team.

Important Dates

Last day to **register** is Sep 3

Last day to **drop** class is Oct 22

Holidays: There will be no class on Sep 6 (Labor Day), Oct 10-17 (Fall break), Nov 25-28 (Thanksgiving).

Machine Learning

Conference on Neural Information Processing Systems (NeurIPS)

International Conference on Learning Representations (ICLR)

International Conference on Machine Learning (ICML)

Application: Computer Vision

Conference on Computer Vision and Pattern Recognition (CVPR)

International Conference on Computer Vision (ICCV)

European Conference on Computer Vision (ECCV)

Application: Natural Language Professing

Annual Meeting of the Association for Computational Linguistics (ACL)

Empirical Methods in Natural Language Processing (EMNLP)

Regression problem

Given $\{\mathbf{x}_n, t_n\}_{n=1}^N$ where \mathbf{x}_n is the observation and t_n is the corresponding target.

How to find t that corresponds to the other \mathbf{x} ?

Construct a function $y(\mathbf{x})$ based on the given dataset $\{\mathbf{x}_n, t_n\}_{n=1}^N$.

Linear basis function models

Linear regression model

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1x_1 + \cdots + w_Dx_D, \quad \mathbf{x} = (x_1, \cdots, x_D)^\top. \quad (1)$$

We can extend the model (1) by considering linear combinations of fixed nonlinear functions of the input variables, of the form

$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x}), \quad (2)$$

where $\phi_j(\mathbf{x})$ are known as basis functions.

Linear basis function models

It is often convenient to define an additional dummy 'basis function' $\phi_0(\mathbf{x}) = 1$ so that

$$y(\mathbf{x}, \mathbf{w}) = \sum_{j=0}^{M-1} w_j \phi_j(\mathbf{x}) = \mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}), \quad (3)$$

where $\mathbf{w} = (w_0, \dots, w_{M-1})^\top$ and $\boldsymbol{\phi} = (\phi_0, \dots, \phi_{M-1})^\top$.

Note $y(\mathbf{x}, \mathbf{w})$ is a nonlinear function of the input vector \mathbf{x} .

Basis functions

1. *Polynomial basis*: $\phi_j(x) = x^j$.
2. *Gaussian basis*: $\phi_j(x) = \exp\left\{-\frac{(x-\mu_j)^2}{2s^2}\right\}$, where μ_j govern the locations of the basis functions in input space, and the parameter s governs their spatial scale.
3. *Sigmoidal basis*: $\phi_j(x) = \sigma\left(\frac{x-\mu_j}{s}\right)$, where $\sigma(a) = \frac{1}{1+\exp(-a)}$ is the logistic sigmoid function.

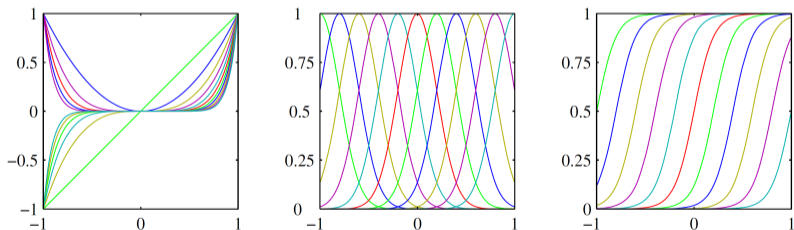


Figure: Left: polynomial basis; Middle: Gaussian basis; Right: Sigmoidal basis.

Linear regression

Given the N observation $\{\mathbf{x}_n, t_n\}_{n=1}^N$, we can define the following loss function

$$L(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \left(t_n - \mathbf{w}^\top \phi(\mathbf{x}_n) \right)^2, \quad (4)$$

which can be written in the following compact form

$$L(\mathbf{w}) = \frac{1}{2} \|\mathbf{t} - \Phi \mathbf{w}\|^2, \quad (5)$$

where

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \cdots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & \cdots & \phi_{M-1}(\mathbf{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \cdots & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}. \quad \text{and} \quad \mathbf{t} = \begin{pmatrix} t_1 \\ \vdots \\ t_N \end{pmatrix}, \quad (6)$$

where Φ is called the *design matrix*, whose elements are given by $\Phi_{nj} = \phi_j(\mathbf{x}_n)$.

Linear regression

Note $L(\mathbf{w})$ is a quadratic function of \mathbf{w} , let $dL(\mathbf{w})/d\mathbf{w} = 0$, we have

$$0 = \frac{dL(\mathbf{w})}{d\mathbf{w}} = \Phi^\top (\Phi \mathbf{w} - \mathbf{t}),$$

therefore,

$$\mathbf{w}_{ML} = (\Phi^\top \Phi)^{-1} \Phi^\top \mathbf{t}. \quad (7)$$

The quantity $\Phi^\dagger \equiv (\Phi^\top \Phi)^{-1} \Phi^\top$ is known as the *Moore-Penrose pseudo-inverse* of the matrix Φ .

Overfitting

How to select M ?

Let us consider a given set of data, which is sampled from $\sin(2\pi x)$ with noise.

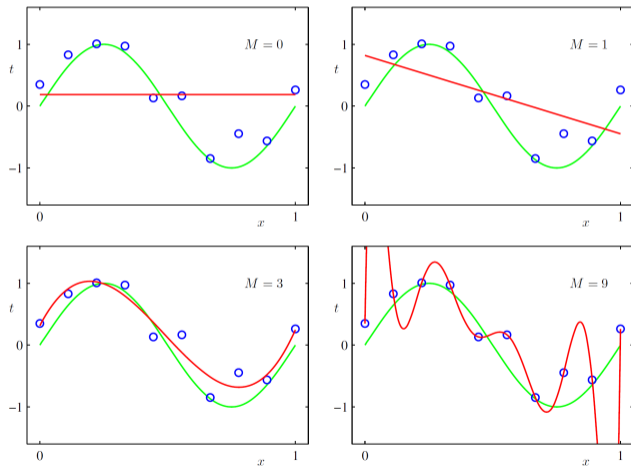


Figure: Plots of polynomials having various orders M , shown in red curves, fitted to a dataset sampled, with noise, from $\sin(2\pi x)$. **Overfitting happens when M is large!**

Why overfitting?

	$M = 0$	$M = 1$	$M = 6$	$M = 9$
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

Figure: Table of the coefficients w^* for polynomials of various order. The magnitude of the coefficients increases dramatically as the order of the polynomial increases.

Why overfitting?

	$M = 0$	$M = 1$	$M = 6$	$M = 9$
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

Figure: Table of the coefficients w^* for polynomials of various order. The magnitude of the coefficients increases dramatically as the order of the polynomial increases.

Regularization is all you need!

Regularization

L_2 -regularization:

$$L^2(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \left(t_n - \mathbf{w}^\top \phi(\mathbf{x}_n) \right)^2 + \frac{\lambda}{2} \|\mathbf{w}\|_2^2, \quad (8)$$

where $\|\mathbf{w}\|_2^2 \equiv \mathbf{w}^\top \mathbf{w} = w_0^2 + w_1^2 + \dots + w_{M-1}^2$.

> λ governs the relative importance of the regularization term compared with the sum-of-squares error term.

> *This is also known in the statistics literature as [parameter shrinkage](#) method because they reduce the value of the coefficients. The particular case of quadratic regularizer is called [ridge regression](#). In the context of neural networks, this approach is known as [weight decay](#).*

Regularization

L_q -regularization:

$$\frac{1}{2} \sum_{n=1}^N \left(t_n - \mathbf{w}^\top \phi(\mathbf{x}_n) \right)^2 + \frac{\lambda}{2} \sum_{j=1}^M |w_j|^q. \quad (9)$$

$q = 1$ corresponds to *Lasso regression*.

Lagrange Multipliers

Lagrange multipliers are used to find the stationary points of a function subject to one or more constraints. Consider

$$\max_{x_1, x_2} f(x_1, x_2) \quad \text{s.t.} \quad g(x_1, x_2) = 0.$$

Let the Lagrangian function defined by

$$L(\mathbf{x}, \lambda) \equiv f(\mathbf{x}) + \lambda g(\mathbf{x}). \quad (10)$$

Thus to find the maximum of a function $f(\mathbf{x})$ subject to the constraint $g(\mathbf{x}) = 0$, we define the Lagrangian function given by (10) and we then find the stationary point of $L(\mathbf{x}, \lambda)$ with respect to both \mathbf{x} and λ .

KKT conditions

The solution to the problem of

$$\max f(\mathbf{x}) \quad \text{s.t.} \quad g(\mathbf{x}) \geq 0$$

is obtained by optimizing the Lagrange function

$$L(\mathbf{x}, \lambda) = f(\mathbf{x}) + \lambda g(\mathbf{x}),$$

with respect to \mathbf{x} and λ subject to the conditions

$$g(\mathbf{x}) \geq 0; \quad \lambda \geq 0; \quad \lambda g(\mathbf{x}) = 0. \tag{11}$$

If we wish to minimize the function $f(\mathbf{x})$ s.t. $g(\mathbf{x}) \geq 0$, then we minimize the Lagrangian function $L(\mathbf{x}, \lambda) = f(\mathbf{x}) - \lambda g(\mathbf{x})$ w.r.t. \mathbf{x} , again subject to $\lambda \geq 0$.

Sparsity due to the L_1 -regularization

Lasso: $q = 1$ in (9), when λ is sufficiently large, some of the coefficients w_j are driven to zero, leading to a *sparse* model in which the corresponding basis functions play no role. Minimizing (9) is equivalent to minimizing the unregularized sum-of-squares error subject to the constraint

$$\sum_{j=1}^M |w_j|^q \leq \eta \quad (12)$$

for an appropriate value of the parameter η , where the two approaches can be related using Lagrange multipliers.

Sparsity due to the L_1 -regularization

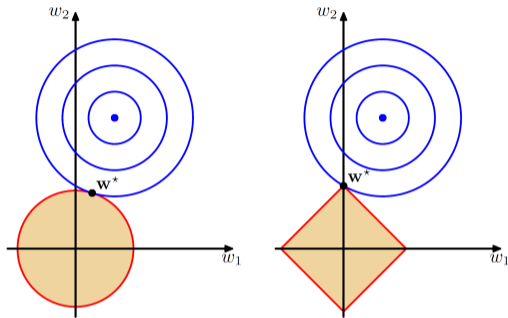


Figure: Plot of the contours of the unregularized error function (blue) along with the constraint region for the quadratic regularizer $q = 2$ on the left and the lasso regularizer $q = 1$ on the right, in which the optimum value for the parameter vector \mathbf{w} is denoted by \mathbf{w}^* . Lasso gives $w_1^* = 0$.

Solution to the L_1 -regularized regression

Lasso regularization problem has the following compact form

$$L(\mathbf{w}) = \frac{1}{2} \|\mathbf{t} - \Phi \mathbf{w}\|_2^2 + \gamma \|\mathbf{w}\|_1, \quad (13)$$

where $\gamma = \lambda/2$. Consider a simple case first, where Φ is orthogonal.

Expanding (13) and excluding irrelevant terms, we have the equivalent form

$$\min_{\mathbf{w}} \left(-\mathbf{t}^\top \Phi \mathbf{w} + \frac{1}{2} \|\mathbf{w}\|^2 \right) + \gamma \|\mathbf{w}\|_1.$$

Solution to the L_1 -regularized regression

Let $\mathbf{t}^\top \Phi := \beta = (\beta_0, \dots, \beta_{M-1})$, the previous problem can be rewritten as

$$\min_{\mathbf{w}} \sum_{i=0}^{M-1} -\beta_i w_i + \frac{1}{2} w_i^2 + \gamma |w_i|.$$

Fix a certain i , we want to minimize

$$\mathcal{L}_i = -\beta_i w_i + \frac{1}{2} w_i^2 + \gamma |w_i|.$$

If $\beta_i > 0$, then we must have $w_i \geq 0$. **why?**

Solution to the L_1 -regularized regression

If $\beta_i > 0$, then we must have $w_i \geq 0$. Otherwise, let $w_i^* < 0$ minimizes \mathcal{L}_i , the $-w_i^*$ enables even smaller \mathcal{L}_i .

If $\beta_i < 0$, then we must choose $w_i \leq 0$.

Solution to the L_1 -regularized regression

$$\mathcal{L}_i = -\beta_i w_i + \frac{1}{2} w_i^2 + \gamma |w_i|.$$

> If $\beta_i > 0$, since $w_i \geq 0$,

$$\mathcal{L}_i = -\beta_i w_i + \frac{1}{2} w_i^2 + \gamma w_i,$$

and differentiating this with respect to w_i and setting equal to zero, we get

$$w_i^* = \beta_i - \gamma,$$

and this is only feasible if the right-hand side is nonnegative (we require $w_i \geq 0$), so in this case the actual solution is

$$w_i^* = \text{sgn}(\beta_i)(|\beta_i| - \gamma)^+. \text{ Note that } \beta_i, \gamma > 0 \text{ and } \beta_i - \gamma \geq 0.$$

Solution to the L_1 -regularized regression

$$\mathcal{L}_i = -\beta_i w_i + \frac{1}{2} w_i^2 + \gamma |w_i|.$$

> If $\beta_i \leq 0$. This implies we must have $w_i \leq 0$ and so

$$\mathcal{L}_i = -\beta_i w_i + \frac{1}{2} w_i^2 - \gamma w_i.$$

Differentiating with respect to w_i and setting equal to zero, we get

$$w_i^* = \text{sgn}(\beta_i)(|\beta_i| - \gamma).$$

But again, to ensure this is feasible, we need $w_i \leq 0$, which is achieved by taking

$$w_i^* = \text{sgn}(\beta_i)(|\beta_i| - \gamma)^+.$$

Solution to the L_1 -regularized regression

$$L(\mathbf{w}) = \frac{1}{2} \|\mathbf{t} - \Phi \mathbf{w}\|_2^2 + \gamma \|\mathbf{w}\|_1,$$

Shrink:

$$w_i^* = \text{sgn}(\beta_i)(|\beta_i| - \gamma)^+.$$