

# Machine Learning

INFQ612L, 440113450A

Spring Semester

Friday 17:00–18:40

IPS

WASEDA University

Prof. Shoji Makino



# Machine Learning

Friday 17:00–18:40

1. 4/18

2. 4/25

3. 5/2

4. 5/9

5. 5/16

6. 5/23

7. 5/30

8. 6/6

9. 6/13

10. 6/20

11. 6/27

–. 7/4 No Lecture

–. 7/11 No Lecture

–. 7/18 No Lecture



At Zoom, set your name as:

Student ID, LAST\_NAME, First\_name

44251234, MAKINO, Shoji


At my class,

please turn on your camera

# Machine Learning (8)(9)

Feature selection  
and  
L1 regularization

# Wine Data Linear Regression

<p>Negative impact</p>  <p>Positive impact</p>	Coefficients		Name	
	1	-0.193967	volatile acidity	Important features
	6	-0.107356	total sulfur dioxide	
	4	-0.088183	chlorides	Not so important feature?
	8	-0.063842	pH	
	2	-0.035553	citric acid	
	7	-0.033737	density	
	3	0.023019	residual sugar	
	0	0.043497	fixed acidity	Important features
	5	0.045606	free sulfur dioxide	
	9	0.155277	sulphates	
	10	0.294243	alcohol	
Mean squared err. 0.41676716722140794				

# Document to Feature Vector: Bag-of-words

## Document

In recent years, there has been a growing trend towards outsourcing of computational tasks with the development of cloud services. We propose two building blocks that work with FHE: a novel batch greater-than primitive, and matrix primitive for encrypted matrices

## Dictionary (12,000 words)

ID	word
1168	batch
1169	bath

...

.

1201	cloud
------	-------

...

.

1239	computation
1240	computational

...

.

1172	primitive
------	-----------

Feature vector is characterized by whether or not it appears without considering frequency

$$\mathbf{x}_i = (0, 0, \dots, 0, 1, 0, 1, 0, \dots, 0, 1, 0, \dots, 0, 1, 0, \dots)$$

# Feature Selection

- Feature selection: We want to eliminate features that are unnecessary for predicting the target values and model using only useful features

$$t = \overset{\text{12,000 dimension}}{(w_1 \quad w_2 \quad \dots \quad w_D)} \begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_D \end{pmatrix} = \sum_{i=1}^D w_i x_i$$

Bag-of-words  
Most feature values are zero,  
so it is wasteful to calculate in  
all dimensions ...

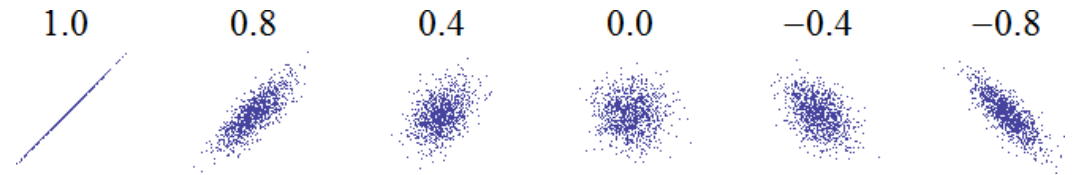
- Why feature selection?
  - Improve calculation efficiency (at the time of prediction)
  - Improved prediction accuracy of target values  
(Modeling with useless features reduces the quality of predictions)
  - Improved interpretation of learning results (which variables are effective for prediction?)
- How to implement feature selection
  - Filter method
  - Wrapper method
  - Utilization of regularization to introduce sparsity



# Filter Method

- Data
  - D-dimensional feature variable (real value)  $x_1, x_2, \dots, x_D$
  - Target variable (real value, for regression)  $t$
- Score the strength of the association between each feature variable and the target variable
  - Example: Pearson correlation coefficient (Example: Correlation between s-th feature variable and target variable)

$$r_s = \frac{\sum_{i=1}^N (\bar{x}_s - x_{si})(\bar{t} - t_i)}{\sqrt{\sum_{i=1}^N (\bar{x}_s - x_{si})^2 \sum_{i=1}^N (\bar{t} - t_i)^2}}$$



- Filter method: Score all feature values and use only strongly related variables as features
  - $|r_1| > \theta?, |r_2| > \theta?, \dots, |r_D| > \theta?$
- If two or more variables affect the target value as a set, they cannot be selected by the filter method

# (Naive) Rapper Method

- Data

- D-dimensional feature variable (real value)  $x_1, x_2, \dots, x_D$
- Target variable (real value, for regression)  $t$

- Goal: Select M features with feature index  $\{1, 2, 3, \dots, D\}$

- All enumeration + cross-validation

- 1) Divide the data into training data and test data
- 2) List power sets  $S \subseteq 2^{\{x_1, x_2, \dots, x_D\}}$  of size M of  $\{1, 2, 3, \dots, D\}$
- 3) For each element of the power set:
  - In the training data, the model is trained with the features specified by the set
  - Evaluate the test error of the model
- 4) Adopt a feature subset that gives the smallest test error

# (Naive) Rapper Method

## . Merit

- Best feature selection in terms of test error
- It can also handle “when two or more variables affect the target value as a set”

## . Problem

- Subset candidates increase exponentially with feature number  $D$
- e.g., number of features  $D = 24$ , number of features to select  $M=12$
- Number of elements in the power set is  ${}_{24}C_{12} = \text{about } 2.7 \text{ million}$
- Training of 2.7 million models, test error evaluation by k-fold CV  
→ It takes too much time

# Wrapper Method: Positive Greedy Search

Goal: Set features to  $\{1, 2, 3, \dots, D\}$  and select useful features

1) Divide the data into training data and test data

2)  $i = 1, F = \{1\}$

3) A model is trained with the features contained in  $F$ , and the test error is evaluated with the model. Let the test error be  $e^i$

4) Set  $F' = F \cup \{i + 1\}$  and model learning with the features included in  $F'$

5) Model training is performed using the features included in  $F'$ , and the test error is evaluated using the training model. Let the test error be  $e^{i+1}$

6) If  $e^{i+1} < e^i$ , then  $F = F'$ . Set  $i = i + 1$  go to 3

A search is also used to remove the features selected backwards

Weaknesses:

- Best feature set in terms of test error is not selected

- Unable to control how many features are selected at the end

- Final result depends on how the features are arranged

# Norm (Lp Norm)

- Definition of distance in vector space (simply)

Unit circles for each norm

- p-norm

$$\|\mathbf{x}\|_p = \left( \sum_{d=1}^D |x_d|^p \right)^{1/p}$$

- 1-norm

- Manhattan distance

$$\|\mathbf{x}\|_1 = \left( \sum_{d=1}^D |x_d| \right)$$

- 2-norm

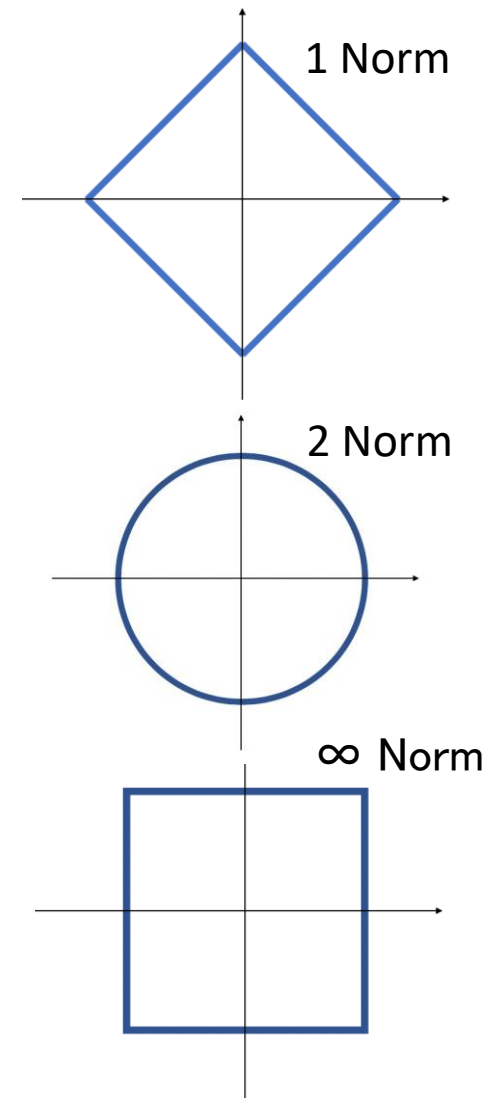
- Euclidean distance

$$\|\mathbf{x}\|_2 = \left( \sum_{d=1}^D (x_d)^2 \right)^{1/2}$$

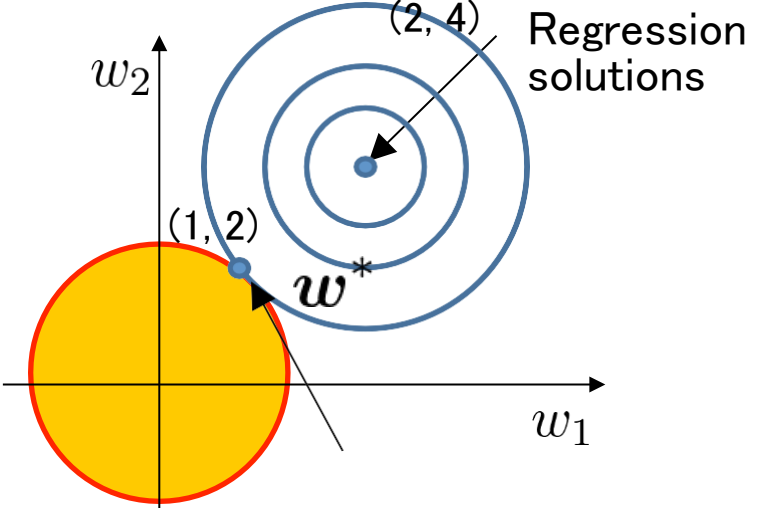
- $\infty$ -norm

- Max norm

$$\|\mathbf{x}\|_\infty = \max_d |x_d|$$



# Meaning of Regularization (for L2 Norm)

- Regression
$$E(\mathbf{w}) = \sum_{i=1}^N (t_i - \mathbf{w}^T \mathbf{x})^2$$
- Ridge regression
$$E(\mathbf{w}) = \sum_{i=1}^N (t_i - \mathbf{w}^T \mathbf{x})^2 + \underbrace{\lambda \mathbf{w}^T \mathbf{w}}_{\text{Squared L2 norm } \|\mathbf{w}\|_2^2}$$


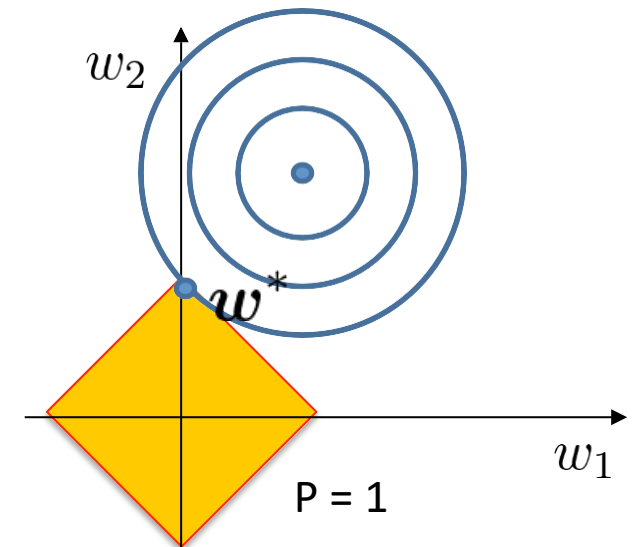
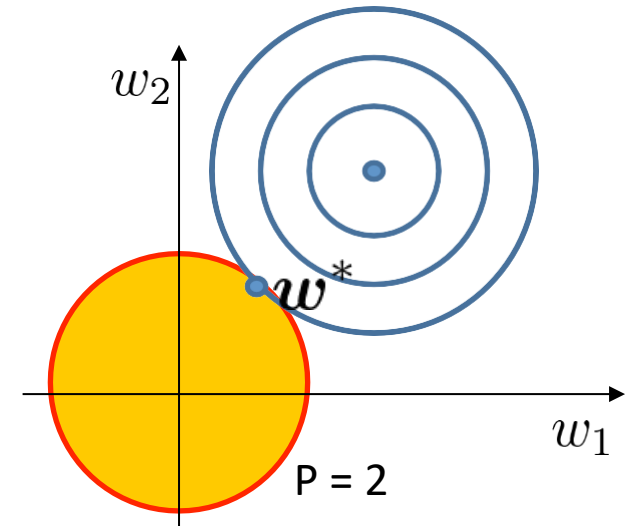
Solution of ridge regression
- Compare to a regression solution without regularization, a solution close to the origin can be found in the sense of the L2 norm
  - e.g.,  $\mathbf{w}^* = (2, 4) \rightarrow (1, 2)$
- The larger the regularization parameter  $\lambda$ , the bigger the penalty for larger norms
  - Solution with a smaller norm is selected rather than reducing the error
  - Attraction to the origin increases

# To Achieve Feature Selection by Regularization

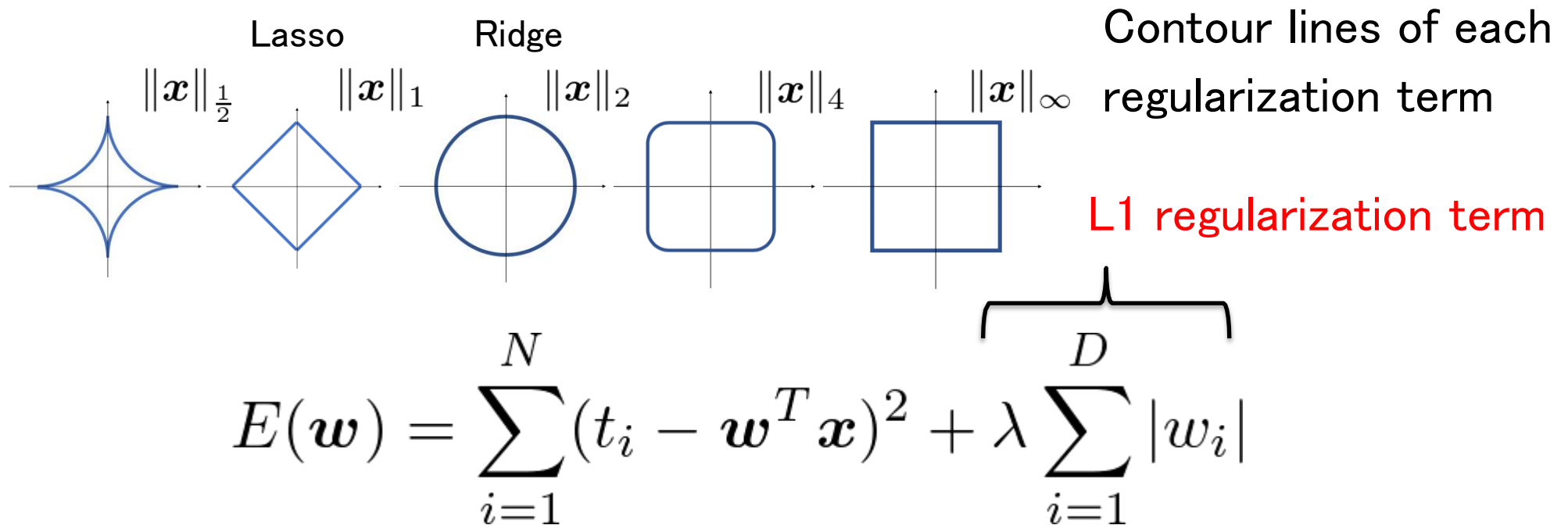
- We want the value to be 0 (on the axis) in some dimensions of the obtained solution
  - e.g.,  $w^* = (0.1, 4.0) \rightarrow (0, 3.8)$ 
    - The first feature is not used for prediction  
→ The second feature was selected as a useful feature
- What is a regularization term suitable for feature selection?

$$E(w) = \sum_{i=1}^N (t_i - w^T x)^2 + \lambda \|w\|_p^p$$

- When  $p = 1$ , the part where the contour line of the objective function and the contour line of the regularization term meet is likely to be "on the axis"
- At this time, the coefficient becomes 0
  - Even  $p < 1$  has such a property, but it is difficult to optimize because it is not a convex function



# Lasso = Squared Error Term + L1 Regularization Term

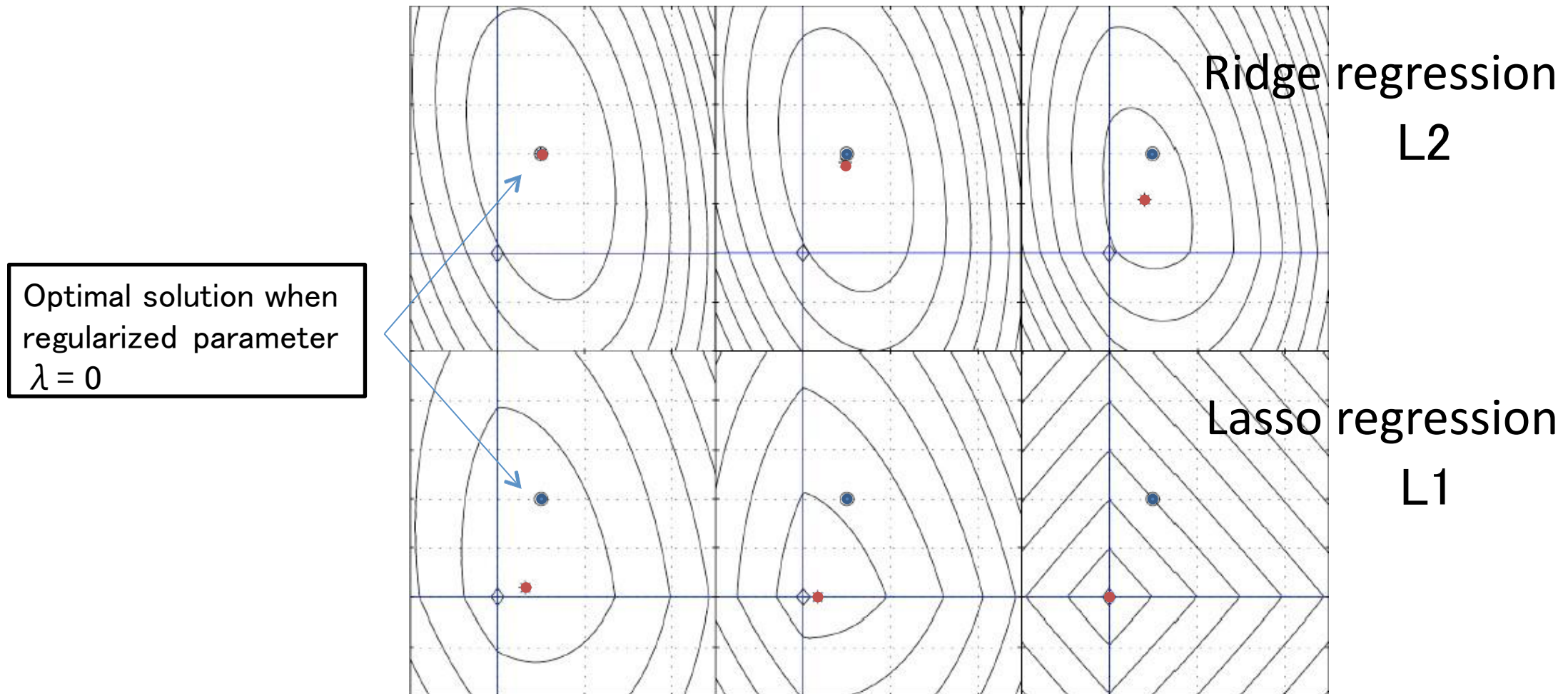


- L1 regularization: linear punishment for Manhattan distance from origin
- Both the error term and the L1 regularization term are convex functions
  - Optimal solution is relatively easy to find, but
  - Not differentiable near the origin  $\rightarrow$  No analytical solution can be found
  - A little difficult to solve (Compared to L2) ...



# Why LASSO can do Feature Selection: Contour Line Change of Loss Function due to Regularization

Weak regularization ← → Strong regularization



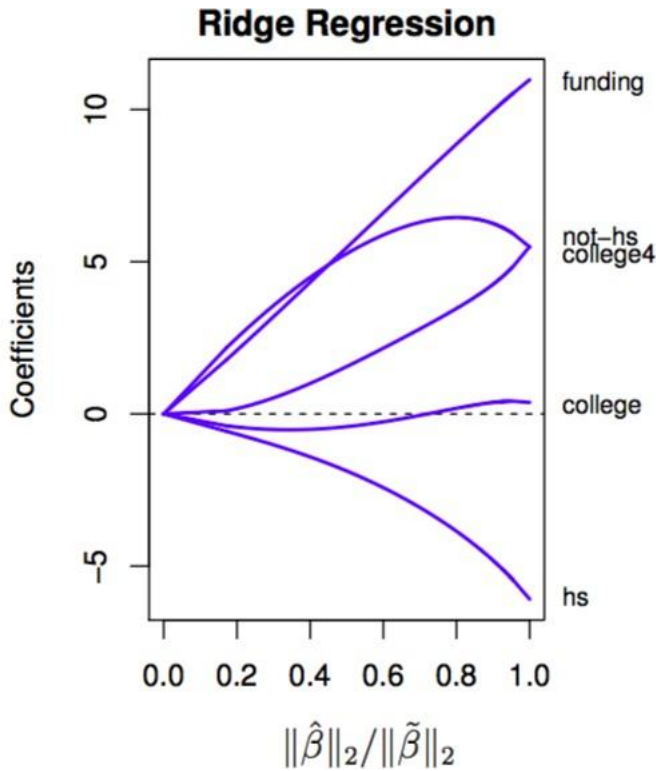
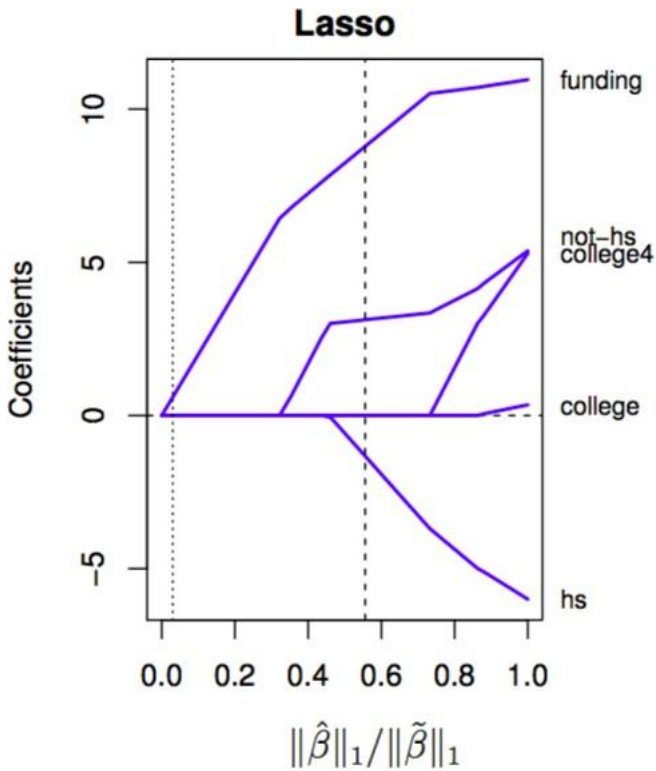
# Regularization Path

How does the effect on the prediction of each feature (that is, the magnitude of the coefficient) change when the regularization coefficient is changed?

**Table 2.1** *Crime data: Crime rate and five predictors, for N = 50 U.S. cities.*

city	funding	hs	not-hs	college	college4	crime rate
1	40	74	11	31	20	478
2	32	72	11	43	18	494
3	57	70	18	16	16	643
4	31	71	11	25	19	341
5	67	72	9	29	24	773
⋮	⋮	⋮	⋮	⋮		
50	66	67	26	18	16	940

Feature value:  
Annual police budget, high school graduation rate over 25 years old, high school graduation rate 16-19 years old, college graduation rate 18-24 years old, 4-year college graduation rate over 25 years old  
Target value: Crime rate



# Feature Selection Summary

- Too many features lead to poor prediction accuracy and increased calculation time
- Feature selection is to narrow down such features
- There are two major feature selection methods
- Filter method / wrapper method
  - Enumerate the combinations of various features,
  - Exhaustive / experimental selection of good combinations in terms of test error
- L1 regularization
  - Introduces regularization that penalizes the number of features to be used, and makes it possible to control the number of features
  - Use k-fold CV to tune regularization parameter  $\lambda$  based on test error to ensure that the appropriate combination of features is used in the model as a result



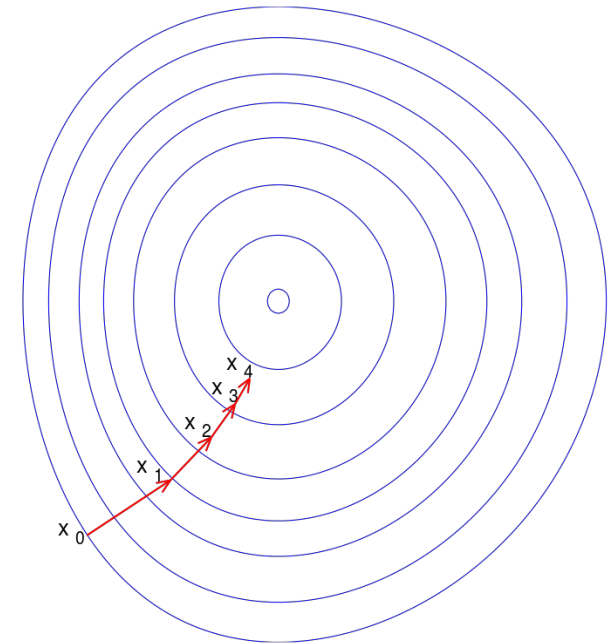
# Gradient Descent Method (Steepest Descent Method)

$$\begin{aligned}\mathbf{w}^{(t+1)} &\leftarrow \mathbf{w}^{(t)} - \eta \nabla E(\mathbf{w}^{(t)}) \\ t &\leftarrow t + 1 \text{ until convergence}\end{aligned}$$

- Calculate the gradient of  $w$  of the total training error of **all** training samples

$$\nabla E(\mathbf{w}^{(t)}) = \nabla \left( \sum_{i=1}^N [(\mathbf{w}^{(t)})^T \mathbf{x}_i - t_i]^2 \right)$$

- Gradient can be calculated relatively accurately, but the amount of calculation per update is  $O(N)$  for the number of samples  $N$
- Not suitable when the number of samples is very large



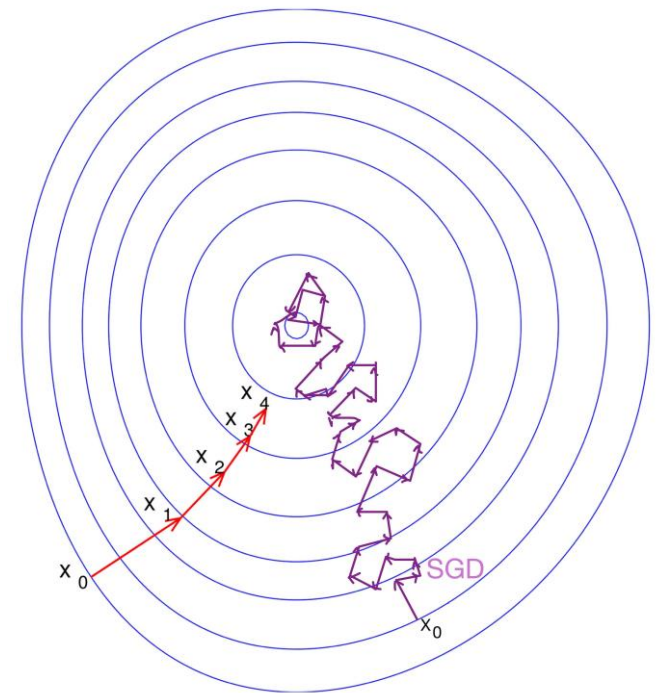
# Stochastic Gradient Descent (SGD)

$$\begin{aligned}\mathbf{w}^{(t+1)} &\leftarrow \mathbf{w}^{(t)} - \eta \nabla E(\mathbf{w}^{(t)}) \\ t &\leftarrow t + 1 \text{ until convergence}\end{aligned}$$

- Calculate the gradient for the training error of **one** randomly selected training sample

$$\nabla E(\mathbf{w}^{(t)}) = \nabla \left( (\mathbf{w}^{(t)})^T \mathbf{x}_n - t_n \right)^2$$

- Although it is approximate in the sense of the gradient of  $E(\mathbf{w})$ , the amount of calculation per update is  $O(1)$  for the number of samples  $N$
- There is a deviation from the gradient for all data, and the optimization process is unstable (training error does not decrease monotonically)



# Mini Batch

$$\begin{aligned}\mathbf{w}^{(t+1)} &\leftarrow \mathbf{w}^{(t)} - \eta \nabla E(\mathbf{w}^{(t)}) \\ t &\leftarrow t + 1 \text{ until convergence}\end{aligned}$$

- Mini batch  $D_t$ 
  - A set of randomly selected relatively small number of data from all data
- Gradient calculation for prediction error of randomly selected **mini-batch**
$$\nabla E(\mathbf{w}^{(t)}) = \nabla \left( \sum_{(\mathbf{x}, t) \in D_t} [(\mathbf{w}^{(t)})^T \mathbf{x} - t]^2 \right)$$
- Computational complexity per update is proportional to mini-batch size
- Gradients are not as unstable as SGD as they are calculated for mini-batch
- Effective use of parallel computer resources
  - Calculate the gradient independently for each batch

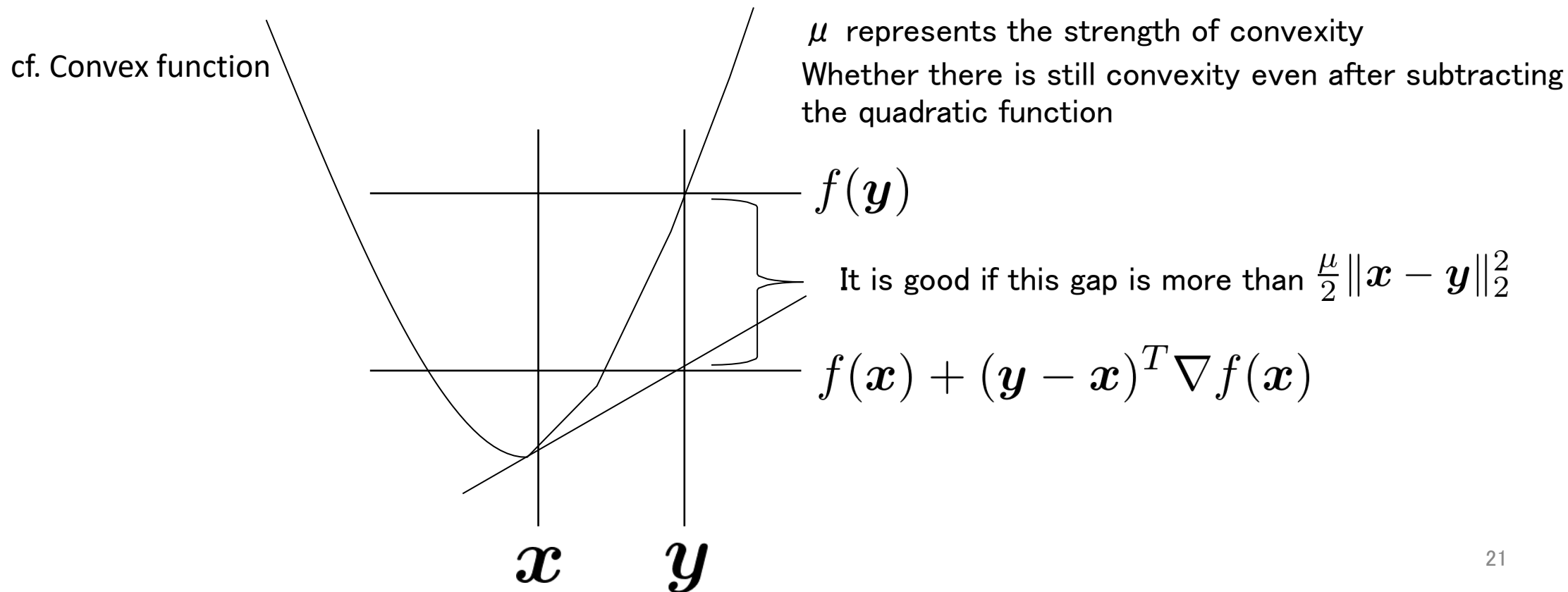
# Convex Function Characterization

## $\mu$ Strong Convex

- If the following is satisfied,  $f: \mathcal{R}^D \rightarrow \mathcal{R}$  is  $\mu$  strong convex

$$f(\mathbf{y}) > f(\mathbf{x}) + (\mathbf{y} - \mathbf{x})^T \nabla f(\mathbf{x}) + \frac{\mu}{2} \|\mathbf{x} - \mathbf{y}\|_2^2$$

$$f(\mathbf{y}) > f(\mathbf{x}) + (\mathbf{y} - \mathbf{x})^T \nabla f(\mathbf{x})$$



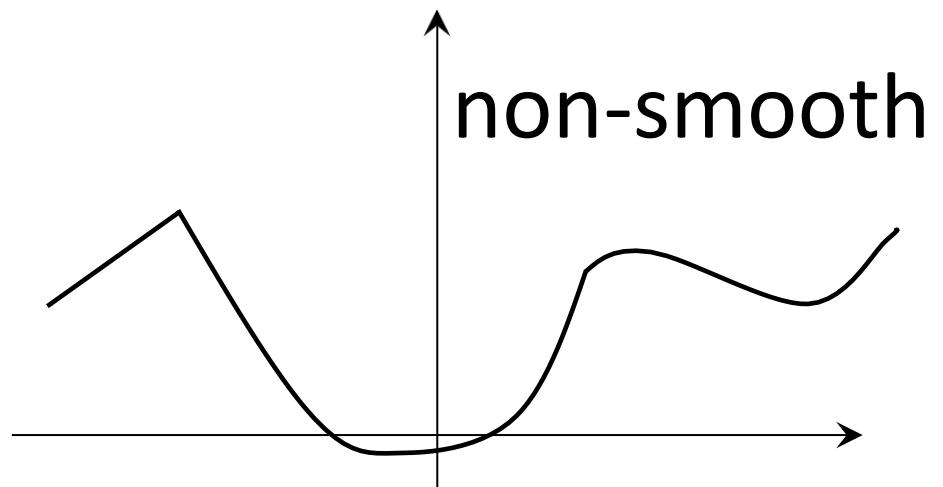
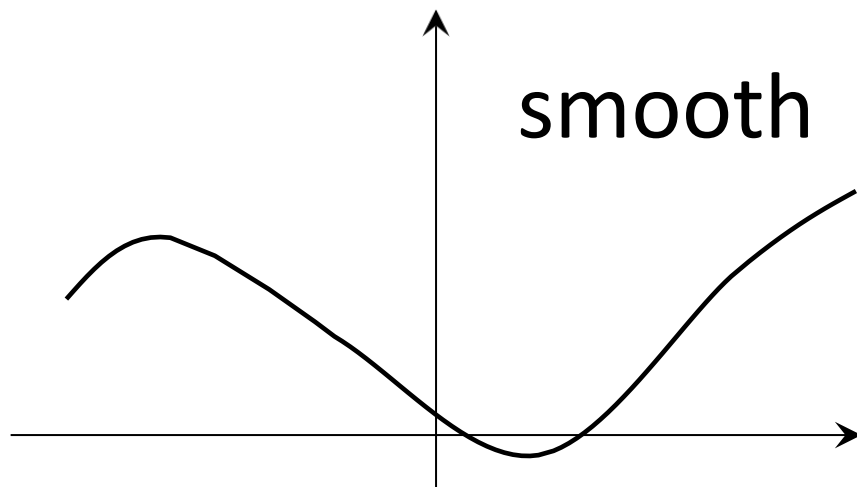


# Convex Function Characterization: L Smooth

- If the following is satisfied,  $f: \mathbb{R}^D \rightarrow \mathbb{R}$  is  $L$  smooth ( $L > 0$ )
  - Gradient change is smaller than a constant multiple of input change

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_2 < L\|\mathbf{x} - \mathbf{y}\|_2$$

- Gradients of adjacent points are close to each other



# Convergence of Gradient Descent

Theorem:  $\mathbf{w}^* = \arg \min_{\mathbf{w}} E(\mathbf{w})$  Also, let the function  $E(\mathbf{w})$  be a convex function that can be differentiated. If the function  $E(\mathbf{w})$  is  $L$  smooth ( $L > 0$ ) and the step size is  $\eta < 1/L$ , then:

$$E(\mathbf{w}^{(t)}) - E(\mathbf{w}^*) < \frac{2L \|\mathbf{w}^{(0)} - \mathbf{w}^*\|_2^2}{t + 4}$$

- Convergent speed is  $O(1/t)$  for the number of steps  $t$
- $O(1/\varepsilon)$  updates are required to reduce the error to  $\varepsilon$  or less
  - (Roughly speaking) If an error of about  $\varepsilon = 0.01$  is achieved by updating  $T$  times, an error of about  $\varepsilon = 0.001$  can be achieved by updating about  $10T$  times

# Convergence of Gradient Descent (Strong Convex)

Theorem:  $\mathbf{w}^* = \arg \min_{\mathbf{w}} E(\mathbf{w})$  The function  $E(\mathbf{w})$  is a differentiable convex function. If the function  $E(\mathbf{w})$  is  $\mu$  strong convex ( $\mu > 0$ ) and  $L$  smooth ( $L > 0$ ), and  $\eta = 1/L$ , then the following holds:

$$E(\mathbf{w}^{(t)}) - E(\mathbf{w}^*) < \left(1 - \frac{\mu}{L}\right)^t [E(\mathbf{w}^{(0)}) - E(\mathbf{w}^*)]$$

- Convergent speed is  $O(c^t)$  for the number of steps  $t$
- To reduce the error to  $\varepsilon$  or less,  $O(\log(1/\varepsilon))$  updates are sufficient
  - (Roughly speaking) If an error of about  $\varepsilon = 0.01$  is achieved by updating about  $2T$  times, an error of about  $\varepsilon = 0.001$  can be achieved by updating about  $3T$  times
- Strong convex converges much faster (first-order convergence)
  - Regression and ridge regression are strong convex and  $L$  smooth
  - Logistic regression is  $L$ -smooth but not strong convex (will come out later)
  - Lasso, SVM is not  $L$ -smooth (will come out later)

# Step Size Parameter (Learning Rate)

- . Learning rate selection
  - If it is too small, convergence will be slow
  - If it is too large, it will converge quickly, but it will vibrate around the optimum solution
  - If it is very large, it will diverge
- . Learning rate scheduling
  - Decrease the step size as the number of updates increases
  - Decrease the step size as the amount of updates decreases
  - Adaptively change the learning rate for each parameter
- . Difficult to theoretically determine an appropriate learning rate
- . Need to tune for each problem

# Sub-Gradient

How to optimize an objective function that contains non-differentiable terms?

L1 regularization term is non-differentiable

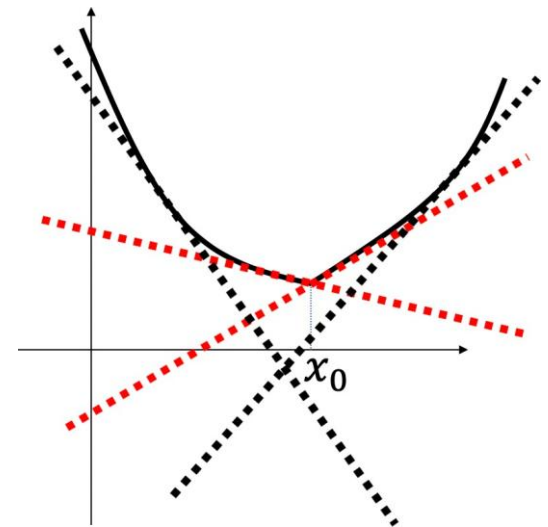
$$E(\mathbf{w}) = \sum_{i=1}^N (t_i - \mathbf{w}^T \mathbf{x})^2 + \lambda \sum_{i=1}^D |w_i|$$

Subgradient  $\partial f(\mathbf{w}_0) = \{\mathbf{g} \in \mathbf{R}^D \mid \forall \mathbf{w}, f(\mathbf{w}) - f(\mathbf{w}_0) \geq (\mathbf{w} - \mathbf{w}_0)^T \mathbf{g}\}$

– Example  $f(\mathbf{w}) = |w_1| + 2|w_2|$

Sub-derivative in  $\mathbf{w}_0 = (1, 0)^T$

$$\partial f(\mathbf{w}_0) = \left\{ \begin{pmatrix} 1 \\ 2z \end{pmatrix} \mid z \in [-1, 1] \right\}$$



# Sub-Gradient Descent

- Replace the gradient in the gradient descent with a sub-gradient

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \partial E(\mathbf{w}^{(t)})$$

- cf. Gradient Descent

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \eta \nabla E(\mathbf{w}^{(t)})$$

# Exercises 5:

## Derivation of ridge regression

### 4.4 リッジ回帰の最急降下法と確率的最急降下法

$\mathbf{x}_i = \begin{pmatrix} 1 \\ x_{i1} \\ \vdots \\ x_{iD} \end{pmatrix}, \mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T \\ \mathbf{x}_2^T \\ \vdots \\ \mathbf{x}_N^T \end{pmatrix}, \mathbf{t} = \begin{pmatrix} t_1 \\ t_2 \\ \vdots \\ t_N \end{pmatrix}, \mathbf{w} = \begin{pmatrix} w_0 \\ w_1 \\ \vdots \\ w_D \end{pmatrix}$  とする。事例群  $\{(\mathbf{x}_i, t_i)\}_{i=1}^N$  を使ってリッジ回帰による線形モデルを求めることを考える。

1. Derive the update formula of the gradient descent method of ridge regression
2. Derive the update formula of the stochastic gradient descent method of ridge regression

# Machine Learning (8)(9)

Feature selection  
and  
L1 regularization



# Machine Learning

INFQ612L, 440113450A

Spring Semester

Friday 17:00–18:40

IPS

WASEDA University

Prof. Shoji Makino

