

Statistical Machine Learning Theory
Lecture 12
Sparsity

Hisashi Kashima
kashima@i.kyoto-u.ac.jp

Topics:

Learning with sparsity

- L_1 regularization & Lasso
- Reduced rank regression

Lasso

3

KYOTO UNIVERSITY

Regression:

Prediction of a continuous target variable

- Training dataset $\{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(N)}, y^{(N)})\}$
 - $\mathbf{x}^{(i)} \in \mathbb{R}^D$: feature vector
 - $y^{(i)} \in \mathbb{R}$: real-valued target value
- Linear regression model: $y = \mathbf{w}^\top \mathbf{x}$
- Least square solution:

$$\begin{aligned}\mathbf{w}^* &= \operatorname{argmin}_{\mathbf{w}} \sum_{i=1}^N (y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)})^2 \\ &= \operatorname{argmin}_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}\end{aligned}$$

$$\begin{aligned}\mathbf{X} &= (\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)})^\top \\ \mathbf{y} &= (y^{(1)}, y^{(2)}, \dots, y^{(N)})^\top\end{aligned}$$

4

KYOTO UNIVERSITY

Ridge regression:

L₂-Regularization for avoiding overfitting

- Overfitting to the training data
 - Especially when the training data is small compared with the input space dimensionality

- Regularized least square solution:

$$\begin{aligned}\mathbf{w}^* &= \operatorname{argmin}_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \gamma \|\mathbf{w}\|_2^2 \\ &= (\mathbf{X}^T \mathbf{X} + \gamma \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}\end{aligned}$$

- $\|\mathbf{w}\|_2^2 = w_1^2 + w_2^2 + \dots + w_D^2$: L₂-regularization term

L₁-regularization:

A sparsity-inducing regularization

- Over-fitting sometimes occurs even with L₂-regularization
 - when the dimensionality is extremely large
 - when the true model uses only a small number of features

- L₁-regularization

- $\|\mathbf{w}\|_1 = |w_1| + |w_2| + \dots + |w_D|$: L₁-regularization term leads to sparse solutions

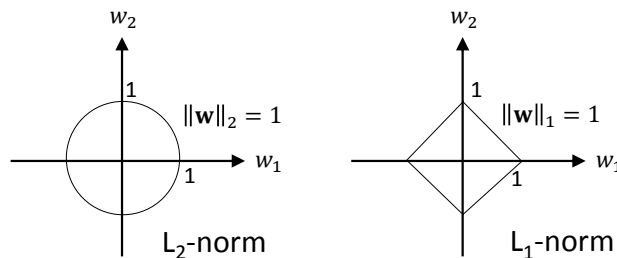
- Sparse: Many w_d becomes 0 in the solutions
- High interpretability and light-weight implementability

- L₁-regularized least square linear regression (LASSO):

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \gamma \|\mathbf{w}\|_1$$

Why does L_1 -regularization induce sparse solutions?: Some intuitive explanations

- L_1 -regularization is equivalent to L_1 -norm constraint:
 $\operatorname{argmin}_{\mathbf{w}} f(\mathbf{w}) + \gamma \|\mathbf{w}\|_1 \Leftrightarrow \operatorname{argmin}_{\mathbf{w}} f(\mathbf{w}) \text{ s.t. } \|\mathbf{w}\|_1 \leq \lambda$
- Some intuitive explanations for sparsity:
 1. L_1 -norm is a convex alternative to L_0 -norm
 2. Level curves of norms



7

KYOTO UNIVERSITY

L_1 -regularized least square linear regression: No closed-form solutions

- L_1 -regularized least square linear regression (LASSO):

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \gamma \|\mathbf{w}\|_1$$
 - L_1 -regularization with a convex loss function is a convex optimization problem
- LASSO has no closed form solution...
 \Rightarrow needs iterative solutions, e.g.:
 1. Optimization with respect to only one dimension
 2. Reduction to L_2 -regularization

we will discuss this

8

KYOTO UNIVERSITY

An algorithm for lasso:

Repeat optimization w.r.t only one dimension

- L_1 -regularization term is cumbersome since:
 - it is not differentiable at $w_d = 0$
 - $w_d = 0$ tends to be a solution
- Observation: The objective function is easy to optimize if we focus only on a single dimension (e.g. w_d)
- Iterative algorithm:
 1. Choose an arbitrary d
 2. Optimize w_d (has a closed form solution)
 3. Repeat steps 1&2 until convergence

9

KYOTO UNIVERSITY

One dimensional optimization problem for LASSO:

Sum of a quadratic function & an absolute value function

- L_1 -regularized least square linear regression (LASSO):
$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \gamma \|\mathbf{w}\|_1$$
- Consider optimization wrt only w_d :
 - $w_d^* = \operatorname{argmin}_{w_d} q(w_d) + \gamma|w_d|$
 - $q(w_d) = a(w_d - \tilde{w}_d)^2 + b$ ($a > 0$): quadratic function
 - \tilde{w}_d is the minimizer of $q(w_d)$ i.e. the solution of the one-variable optimization when $\gamma|w_d|$ is neglected
- Finally what we want is

$$w_d^* = \operatorname{argmin}_{w_d} \frac{1}{2}(w_d - \tilde{w}_d)^2 + \lambda|w_d| \quad (\lambda = \frac{1}{2a}\gamma)$$

10

KYOTO UNIVERSITY

Solution of the one-dimensional optimization: Find the stationary point

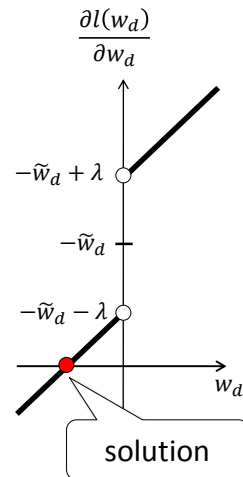
- Find the minimizer of $l(w_d) = \frac{1}{2}(w_d - \tilde{w}_d)^2 + \lambda|w_d|$

- Taking the derivative of $l(w_d)$,

$$\frac{\partial l(w_d)}{\partial w_d} = \begin{cases} w_d - \tilde{w}_d + \lambda & (\text{if } w_d > 0) \\ w_d - \tilde{w}_d - \lambda & (\text{if } w_d < 0) \\ \text{undefined} & (\text{otherwise}) \end{cases}$$

- Solution: $w_d = w_d^*$ s.t. $\left. \frac{\partial l(w_d)}{\partial w_d} \right|_{w_d=w_d^*} = 0$

- lies at $\frac{\partial l(w_d)}{\partial w_d}$ hits the x-axis



11

KYOTO UNIVERSITY

Sparsity of lasso solutions: Solutions close to zero are rounded to zero

- We have 3 cases:

- $-\tilde{w}_d + \lambda < 0$ (i.e. $\tilde{w}_d > \lambda$),

- Solution: $w_d^* = -\tilde{w}_d + \lambda$

- $-\tilde{w}_d - \lambda > 0$ (i.e. $\tilde{w}_d < -\lambda$),

- Solution: $w_d^* = \tilde{w}_d + \lambda$

- $-\lambda \leq \tilde{w}_d \leq \lambda$

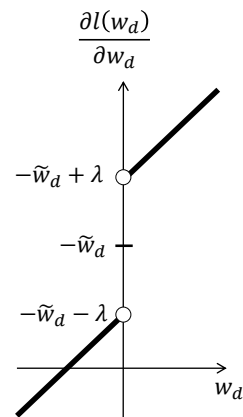
sparse solution

- Solution: $w_d^* = 0$

- if $w_d^* > 0$, we have a contradiction

$$\left. \frac{\partial l(w_d)}{\partial w_d} \right|_{w_d=w_d^*} = w_d^* - \tilde{w}_d + \lambda = 0 \Rightarrow w_d^* = \tilde{w}_d - \lambda \leq 0$$

- Similarly, assuming $w_d^* < 0$ yields a contradiction $w_d^* \geq 0$



12

KYOTO UNIVERSITY

Dimension Reduction

Multivariate regression:

Prediction of multiple continuous variables

- Multivariate regression is a regression problem to predict multiple output variables
 - $f: \mathbb{R}^D \Rightarrow \mathbb{R}^{D'}$
- Training dataset $\{ (\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), \dots, (\mathbf{x}^{(N)}, \mathbf{y}^{(N)}) \}$
 - $\mathbf{x}^{(i)} \in \mathbb{R}^D$: feature vector
 - $\mathbf{y}^{(i)} \in \mathbb{R}^{D'}$: real-valued target values
- Multivariate linear regression model: $\mathbf{y} = \mathbf{W}^T \mathbf{x}$
 - $\mathbf{W} \in \mathbb{R}^{D' \times D}$: Matrix parameter

Solution of multivariate regression: Closed form least square solution

- Least square solution:

$$\begin{aligned} \mathbf{W}^* &= \operatorname{argmin}_{\mathbf{W} \in \mathbb{R}^{D' \times D}} \sum_{i=1}^N \|\mathbf{y}^{(i)} - \mathbf{W}^\top \mathbf{x}^{(i)}\|_2^2 \\ &= \operatorname{argmin}_{\mathbf{W}} \|\mathbf{Y} - \mathbf{X}\mathbf{W}\|_{\text{F}}^2 \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y} \end{aligned}$$

$$\begin{aligned} \mathbf{X} &= (\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)})^\top \\ \mathbf{Y} &= (\mathbf{y}^{(1)}, \mathbf{y}^{(2)}, \dots, \mathbf{y}^{(N)})^\top \\ \frac{\partial \operatorname{tr}(\mathbf{A}\mathbf{B})}{\partial \mathbf{A}} &= \mathbf{B}^\top \end{aligned}$$

- Regularized version
 - $\|\mathbf{W}\|_{\text{F}}^2 = \sum_{(i,j)} w_{ij}^2$: L₂-regularization term
 - $\mathbf{W}^* = (\mathbf{X}^\top \mathbf{X} + \gamma \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{Y}$

15

KYOTO UNIVERSITY

Reduced rank regression: Multivariate regression with rank constraint

- Multivariate regression is equivalent to D' -independent univariate regressions
 - exploits no shared information
- Low-rank assumption $\mathbf{W} = \mathbf{U}\mathbf{V}^\top$
 - $\mathbf{U} \in \mathbb{R}^{D \times K}$, $\mathbf{V} \in \mathbb{R}^{D' \times K}$ i.e. rank of \mathbf{W} is K
 - $K < \min(D, D')$
 - D' output variables share K -dimensional latent space
- Reduced rank regression:

$$\mathbf{W}^* = \operatorname{argmin}_{\mathbf{W}} \|\mathbf{Y} - \mathbf{X}\mathbf{W}\|_{\text{F}}^2 \text{ s.t. } \operatorname{rank}(\mathbf{W}) \leq K$$

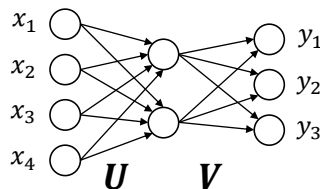
16

KYOTO UNIVERSITY

Sparsity in reduced rank regression:

Sparse parameters in terms of matrix singular values

- Parameter \mathbf{W} in the reduced rank regression $\mathbf{y} = \mathbf{W}^\top \mathbf{x}$ is dense in terms of matrix elements
- \mathbf{W} is sparse in terms of singular values
 - $\mathbf{W} = \mathbf{U}\mathbf{V}^\top$ is low-rank
 - $\mathbf{U} \in \mathbb{R}^{D \times K}, \mathbf{V} \in \mathbb{R}^{D' \times K}, K < \min(D, D')$
 - Rank = L_0 norm of singular values: $\text{rank}(\mathbf{W}) = \|\boldsymbol{\sigma}(\mathbf{W})\|_0$



17

KYOTO UNIVERSITY

Solution of reduced rank regression (1/2):

Best rank- K approximation of a matrix

- Objective function to be minimized:

$$\begin{aligned} \|\mathbf{Y} - \mathbf{X}\mathbf{W}\|_F^2 &= \text{tr}\{(\mathbf{Y} - \mathbf{X}\mathbf{W})^\top (\mathbf{Y} - \mathbf{X}\mathbf{W})\} \\ &= \text{tr}\{\mathbf{Y}^\top \mathbf{Y} - 2\mathbf{W}^\top \mathbf{X}^\top \mathbf{Y} + \mathbf{W}^\top \mathbf{X}^\top \mathbf{X}\mathbf{W}\} \end{aligned}$$

(Let $\mathbf{X}^\top \mathbf{X} = \mathbf{P}^\top \boldsymbol{\Lambda} \mathbf{P}$ be the eigendecomposition)

$\mathbf{P}^\top \mathbf{P} = \mathbf{P}\mathbf{P}^\top = \mathbf{I}$
 (\mathbf{P} : orthogonal)

$$\begin{aligned} &= \text{tr}\{\mathbf{Y}^\top \mathbf{Y} - 2\widetilde{\mathbf{W}}^\top \boldsymbol{\Lambda}^{-\frac{1}{2}} \mathbf{P}\mathbf{X}^\top \mathbf{Y} + \widetilde{\mathbf{W}}^\top \widetilde{\mathbf{W}}\} \\ &\quad \text{where } \widetilde{\mathbf{W}} = \boldsymbol{\Lambda}^{\frac{1}{2}} \mathbf{P}\mathbf{W} \\ &= \|\widetilde{\mathbf{W}} - \boldsymbol{\Lambda}^{-\frac{1}{2}} \mathbf{P}\mathbf{X}^\top \mathbf{Y}\|_F^2 + \text{const.} \end{aligned}$$
- Find the best rank- K approximation of $\boldsymbol{\Lambda}^{-\frac{1}{2}} \mathbf{P}\mathbf{X}^\top \mathbf{Y}$

18

KYOTO UNIVERSITY

Solution of reduced rank regression (2/2):

Closed form solution using SVD

- The best rank- K approximation of $\Lambda^{-\frac{1}{2}}\mathbf{P}\mathbf{X}^T\mathbf{Y}$ is given as $\widetilde{\mathbf{W}}^* = \mathbf{U}^*\mathbf{\Sigma}^*\mathbf{V}^{*\top}$
 - \mathbf{V}^* is top- K eigenvectors of
$$\mathbf{Y}^T\mathbf{X}\mathbf{P}^T\Lambda^{-\frac{1}{2}}\Lambda^{-\frac{1}{2}}\mathbf{P}\mathbf{X}^T\mathbf{Y} = \mathbf{Y}^T\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Y}$$
 - $\mathbf{\Sigma}^*$: a diagonal matrix with K largest singular values
 - $\mathbf{U}^* = \Lambda^{-\frac{1}{2}}\mathbf{P}\mathbf{X}^T\mathbf{Y}\mathbf{V}^*\mathbf{\Sigma}^{*-1}$
- The solution is $\mathbf{W}^* = \mathbf{P}^T\Lambda^{-\frac{1}{2}}\widetilde{\mathbf{W}}^* = \mathbf{P}^T\Lambda^{-\frac{1}{2}}\mathbf{U}^*\mathbf{\Sigma}^*\mathbf{V}^{*\top} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Y}\mathbf{V}^*\mathbf{V}^{*\top}$

[Supplement 1] Eigenvalue decomposition of symmetric matrix

- $\mathbf{A} = \mathbf{P}^T\mathbf{\Lambda}\mathbf{P}$: eigen-decomposition of symmetric matrix \mathbf{A}
 - $\mathbf{\Lambda}$: diagonal matrix $\mathbf{\Lambda} = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_D)$, where $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_D \geq 0$ (eigenvalues)
 - \mathbf{P} : orthogonal matrix $\mathbf{P}^T\mathbf{P} = \mathbf{P}\mathbf{P}^T = \mathbf{I}$

[Supplement 2] Singular value decomposition (SVD) and best rank- K approximation :

- $\mathbf{B} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$: SVD of rank- R real matrix \mathbf{B}
 - $\mathbf{\Sigma}$: diagonal matrix $\mathbf{\Sigma} = \text{diag}(\sigma_1, \sigma_2, \dots, \sigma_R, 0, \dots, 0)$, where $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_D \geq 0$ (singular values)
 - $\mathbf{\Sigma}$ is sqrt of eigenvalues of $\mathbf{B}\mathbf{B}^T$ or $\mathbf{B}^T\mathbf{B}$
 - \mathbf{U}, \mathbf{V} : orthogonal matrices
 - \mathbf{U} is eig.vecs of $\mathbf{B}\mathbf{B}^T$, \mathbf{V} is eig.vecs of $\mathbf{B}^T\mathbf{B}$, $\mathbf{u}_i = \frac{1}{\sigma_i}\mathbf{B}^T\mathbf{v}_i$
- Best rank- K approximation problem of matrix \mathbf{B} :

$$\hat{\mathbf{B}}^* = \underset{\hat{\mathbf{B}}}{\text{argmin}} \|\mathbf{B} - \hat{\mathbf{B}}\|_F^2 \text{ s.t. } \text{rank}(\hat{\mathbf{B}}) \leq K$$
 - Find K largest singular values $\mathbf{\Sigma}^* = \text{diag}(\sigma_1, \dots, \sigma_K)$, and corresponding vectors $\mathbf{U}^* = (\mathbf{u}_1, \dots, \mathbf{u}_K)$, $\mathbf{V}^* = (\mathbf{v}_1, \dots, \mathbf{v}_K)$, and let $\hat{\mathbf{B}}^* = \mathbf{U}^*\mathbf{\Sigma}^*\mathbf{V}^{*\top}$

21

KYOTO UNIVERSITY

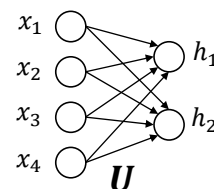
Dimension reduction:

Find low-dimensional representations of high-dim. data

- Dimension reduction:
 - Find a low-dimensional mapping $f: \mathbb{R}^D \Rightarrow \mathbb{R}^K$ ($D > K$)
 - for interpretability, computational/space efficiency, generalization abilities, ...
 - (Lossy) compression: keep the original information as much as possible

- Linear dimension reduction: $\mathbf{h} = \mathbf{U}^T \mathbf{x}$

- $\mathbf{U} : D \times K$ matrix



22

KYOTO UNIVERSITY

Basic idea behind dimension reduction:

Find a coding & decoding function for lossy compression

- Coding and decoding process:

$$\mathbf{x} \xrightarrow[\text{coding}]{f} \mathbf{h} \xrightarrow[\text{decoding}]{g} \tilde{\mathbf{x}}$$

- If f and g are appropriately designed so that $\mathbf{x} \approx \tilde{\mathbf{x}}$, \mathbf{h} must be a good low-dimensional representation of \mathbf{x}
- Optimization problem
 - $(f, g) = \operatorname{argmin}_{f, g} \sum_{i=1}^N \operatorname{loss}(\mathbf{x}^{(i)}, g(f(\mathbf{x}^{(i)})))$

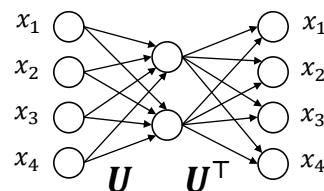
23

KYOTO UNIVERSITY

Principal component analysis:

Dimension reduction using reduced rank regression

- Linear dimension reduction with coding & decoding functions
 - linear coding function $f : \mathbf{h} = \mathbf{U}^T \mathbf{x}$ ($\mathbf{U} : D \times K$ matrix)
 - linear decoding function $g : \tilde{\mathbf{x}} = \mathbf{V} \mathbf{h}$ ($\mathbf{V} : K \times D$ matrix)
 - $\tilde{\mathbf{x}} = \mathbf{V} \mathbf{U}^T \mathbf{x}$
- Reduced rank regression finds the solution by taking the training dataset as $\{(\mathbf{x}^{(1)}, \mathbf{x}^{(1)}), \dots, (\mathbf{x}^{(N)}, \mathbf{x}^{(N)})\}$
 - Solution will be $\mathbf{V} = \mathbf{U}^T$



24

KYOTO UNIVERSITY

Statistical Machine Learning Theory
Homework 2

Hisashi Kashima
kashima@i.kyoto-u.ac.jp

Homework :

Solve all of the problems

1. Find a closed form solution for the label propagation (Lecture 10, page 10)
2. In Halving algorithm (lecture 11, page 7), we replace the 2nd step (prediction) with “Make a prediction with a predictor randomly chosen from the current version space”. Show a mistake bound of the modified algorithm. Is this a good mistake bound? Why?
3. In the perceptron’s mistake bound lemma (lecture 11, page 26), show how the mistake bound will be modified if there exists \mathbf{w}^* s.t. $\forall t, y^{(t)} \langle \mathbf{w}^*, \mathbf{x}^{(t)} \rangle \geq \gamma > 0$
4. Find \tilde{w}_d in the iterative solution for lasso (lecture 12, page 9)

Report submission:

Early submitters are appreciated

- Submission:
 - Final deadline: Feb. 7th noon
 - Send to kashipong+report@gmail.com
with title “SML2014 report 2” and confirm you receive an ack
- Report rating policy:
 - Earlier submitters are more appreciated than those with the same quality