*Statistical Machine Learning Theory*

# Sparsity

Hisashi Kashima

kashima@i.Kyoto-u.ac.jp

# Topics:
# Learning with sparsity

- $L_1$-regularization & Lasso

- Reduced rank regression

- Dimension reduction

# Lasso

# 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:

$$\mathbf{w}^* = \mathrm{argmin}_\mathbf{w} \sum_{i=1}^{N} \left( y^{(i)} - \mathbf{w}^\top \mathbf{x}^{(i)} \right)^2$$
$$= \mathrm{argmin}_\mathbf{w} \|\mathbf{y} - X\mathbf{w}\|_2^2$$
$$= (X^\top X)^{-1} X^\top \mathbf{y}$$

$$X = \left( \mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)} \right)^\top$$
$$y = \left( y^{(1)}, y^{(2)}, \dots, y^{(N)} \right)^\top$$

# Ridge regression:
## L$_2$-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:

$$\mathbf{w}^* = \text{argmin}_{\mathbf{w}} \|\mathbf{y} - X\mathbf{w}\|_2^2 + \gamma\|\mathbf{w}\|_2^2$$
$$= (X^\top X + \gamma I)^{-1} X^\top \mathbf{y}$$

  – $\|\mathbf{w}\|_2^2 = w_1^2 + w_2^2 + \cdots + w_D^2$: L$_2$-regularization term

  – Analytical solution exists

# L$_1$-regularization:
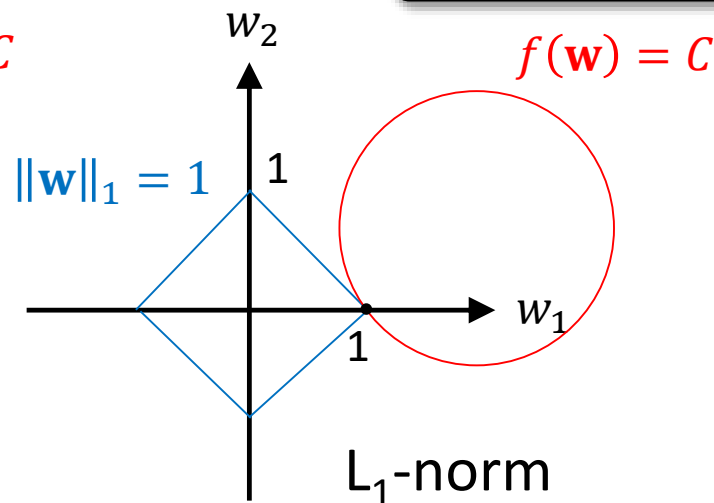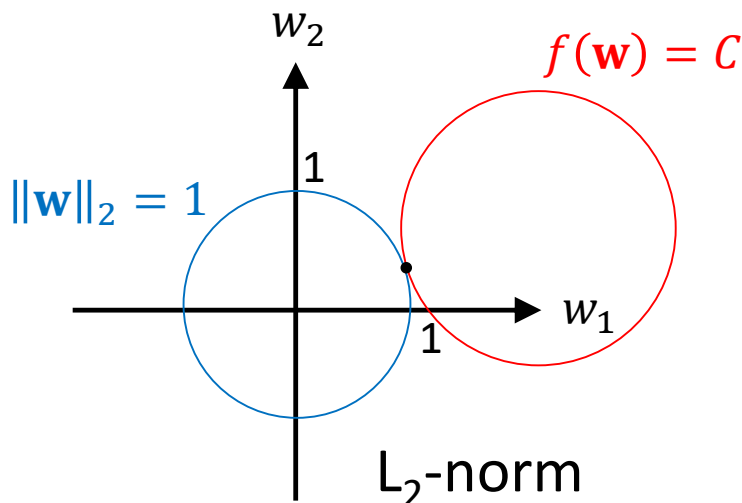## A sparsity-inducing regularization

- Over-fitting sometimes occurs even with L$_2$-regularization

  - when the dimensionality is extremely large

  - when the true model uses only a small number of features

- L$_1$-regularization

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

    - Sparse: Many $w_d$ becomes 0 in the solutions
    - High interpretability and easy-to-implementability

  - L$_1$-regularized least square linear regression (LASSO):
    $$\mathbf{w}^* = \text{argmin}_{\mathbf{w}} \|\mathbf{y} - X\mathbf{w}\|_2^2 + \gamma\|\mathbf{w}\|_1$$

# Why does L$_1$-regularization induce sparse solutions?: Some intuitive explanations

- L$_1$-regularization is equivalent to L$_1$-norm constraint:

$$\mathrm{argmin}_{\mathbf{w}} \, f(\mathbf{w}) + \gamma \|\mathbf{w}\|_1 \iff \mathrm{argmin}_{\mathbf{w}} \, f(\mathbf{w}) \text{ s.t. } \|\mathbf{w}\|_1 \leq \lambda$$

- Some intuitive explanations for sparsity:

1. L$_1$-norm is a convex alternative to L$_0$-norm

2. Level curves of norms and loss

= #nonzero elements



$w_2$     $f(\mathbf{w}) = C$

$\|\mathbf{w}\|_2 = 1$

L$_2$-norm

$w_2$     $f(\mathbf{w}) = C$

$\|\mathbf{w}\|_1 = 1$

L$_1$-norm

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

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

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

  - $L_1$-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

# 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: Coordinate-wise descent

  1. Choose an arbitrary $d$

  2. Optimize $w_d$ (has a closed form solution)

  3. Repeat steps 1&2 until convergence

# One dimensional optimization problem for LASSO: Sum of a quadratic function & an absolute value function

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

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

- Consider optimization w.r.t. only $w_d$:

  Neglect the other terms not depending on $w_d$

  - $w_d^* = \text{argmin}_{w_d} q(w_d) + \gamma|w_d|$

    - $q(w_d) = a(w_d - \widetilde{w}_d)^2 + b$ $(a > 0)$: quadratic function

      - $\widetilde{w}_d$ 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

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

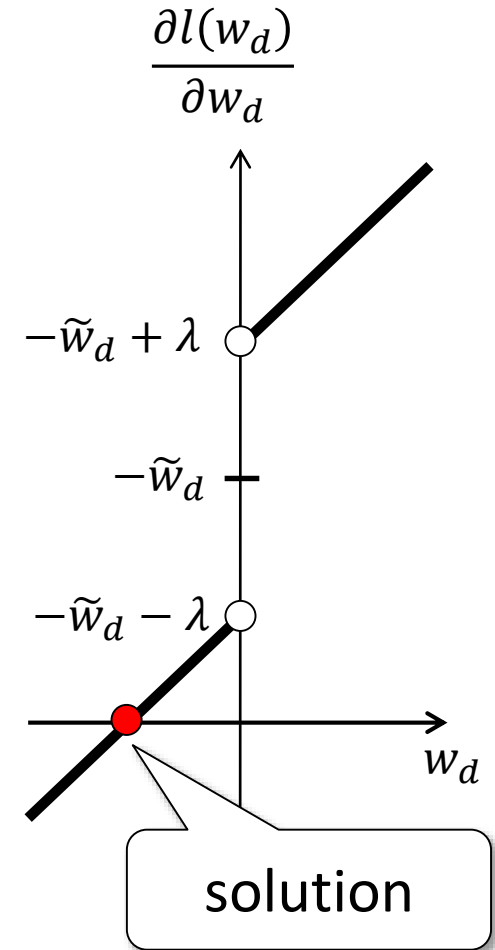# Solution of the one-dimensional optimization: Find the stationary point

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

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

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

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

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

$\dfrac{\partial l(w_d)}{\partial w_d}$

$-\widetilde{w}_d + \lambda$

$-\widetilde{w}_d$

$-\widetilde{w}_d - \lambda$

$w_d$

solution

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

- We have 3 cases:

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

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

  2. $-\widetilde{w}_d - \lambda > 0$ (i.e. $\widetilde{w}_d < -\lambda$),
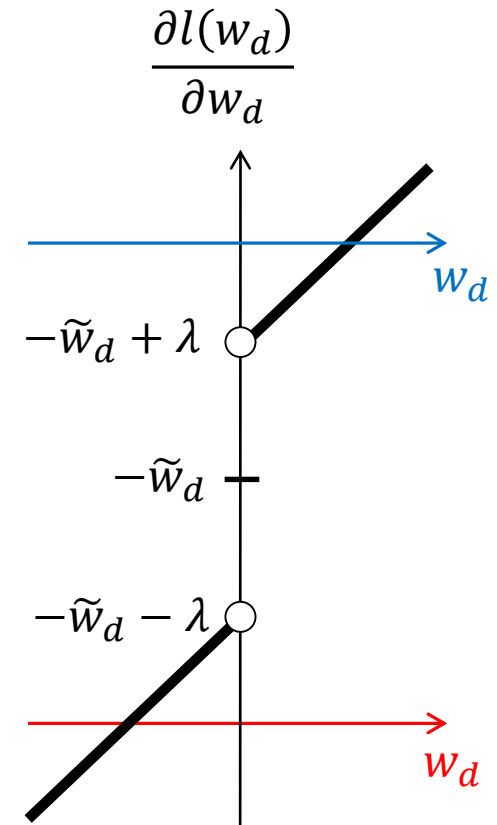
     - Solution: $w_d^* = \widetilde{w}_d + \lambda$

  3. $-\lambda \leq \widetilde{w}_d \leq \lambda$

     sparse solution

     - Solution: $w_d^* = 0$

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

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

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

$\dfrac{\partial l(w_d)}{\partial w_d}$

$w_d$

$-\widetilde{w}_d + \lambda$

$-\widetilde{w}_d$

$-\widetilde{w}_d - \lambda$

$w_d$

# 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 $\left\{ \left(\mathbf{x}^{(1)}, \mathbf{y}^{(1)}\right), \ldots, \left(\mathbf{x}^{(N)}, \mathbf{y}^{(N)}\right) \right\}$

  - $\mathbf{x}^{(i)} \in \mathbb{R}^D$: feature vector

  - $\mathbf{y}^{(i)} \in \mathbb{R}^{D'}$: real-valued target values

- Multivariate linear regression model: $\mathbf{y} = \boldsymbol{W}^\top \mathbf{x}$

  - $\boldsymbol{W} \in \mathbb{R}^{D \times D'}$: Matrix parameter

# Solution of multivariate regression:
# Closed form least square solution

- Least square solution:

$$W^* = \text{argmin}_{W \in \mathbb{R}^{D' \times D}} \sum_{i=1}^{N} \left\| \mathbf{y}^{(i)} - W^\top \mathbf{x}^{(i)} \right\|_2^2$$

$$= \text{argmin}_W \|Y - XW\|_F^2$$

$$= (X^\top X)^{-1} X^\top Y$$

$$X = \left( \mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)} \right)^\top$$

$$Y = \left( \mathbf{y}^{(1)}, \mathbf{y}^{(2)}, \dots, \mathbf{y}^{(N)} \right)^\top$$

$$\frac{\partial \text{tr}(AB)}{\partial A} = B^\top$$

- Regularized version

  - $\|W\|_F^2 = \sum_{(i,j)} w_{ij}^2$: L$_2$-regularization term

  - $W^* = (X^\top X + \gamma I)^{-1} X^\top Y$
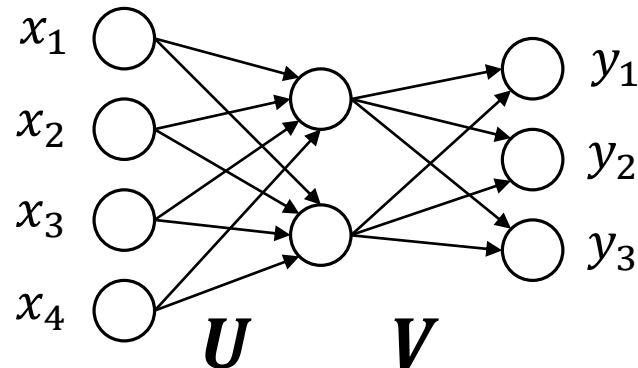
# Reduced rank regression:
## Multivariate regression with rank constraint

- Multivariate regression is equivalent to $D'$-independent univariate regressions

  - exploits no shared information

- Low-rank assumption $\boldsymbol{W} = \boldsymbol{U}\boldsymbol{V}^{\top}$

  - $\boldsymbol{U} \in \mathbb{R}^{D \times K}, \boldsymbol{V} \in \mathbb{R}^{D' \times K}$ i.e. rank of $\boldsymbol{W}$ is $K$

    - $K < \min(D, D')$

  - $D'$ output variables share $K-$dimensional latent space

- Reduced rank regression: Sparsity in the dim of latent space
$$\boldsymbol{W}^* = \mathrm{argmin}_{\boldsymbol{W}}\|\boldsymbol{Y} - \boldsymbol{X}\boldsymbol{W}\|_{\mathrm{F}}^2 \text{ s.t. } \mathrm{rank}(\boldsymbol{W}) \le K$$

# Sparsity in reduced rank regression:
## Sparse parameters in terms of matrix singular values

- Parameter $W$ in the reduced rank regression $\mathbf{y} = W^\top \mathbf{x}$ is dense in terms of matrix elements

- $W$ is sparse in terms of singular values

  - $W = UV^\top$ is low-rank

    - $U \in \mathbb{R}^{D \times K}, V \in \mathbb{R}^{D' \times K}, K < \min(D, D')$

  - Rank = $L_0$ norm of singular values: $\mathrm{rank}(W) = \|\boldsymbol{\sigma}(W)\|_0$

# [Review] Eigenvalue decomposition of symmetric matrix

- Symmetric matrix can be diagonalized using an orthogonal matrix

- $A = P^\top \Lambda P$: eigen-decomposition of symmetric matrix $A$

  - $\Lambda$: diagonal matrix $\Lambda = \mathrm{diag}(\lambda_1, \lambda_2, \ldots, \lambda_D)$, where $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_D \geq 0$ (eigenvalues)

  - $P$: orthogonal matrix $P^\top P = P P^\top = I$

# [Review] Singular value decomposition (SVD) and best rank-$K$ approximation

- $\boldsymbol{B} = \boldsymbol{U\Sigma V}^\top$: SVD of rank-$R$ real matrix $\boldsymbol{B}$

  - $\boldsymbol{\Sigma}$: diagonal matrix $\boldsymbol{\Sigma} = \mathrm{diag}(\sigma_1, \sigma_2, \ldots, \sigma_R, 0, \ldots, 0)$, where $\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_D \geq 0$ (singular values)

    - $\boldsymbol{\Sigma}$ is sqrt of eigenvalues of $\boldsymbol{BB}^\top$ or $\boldsymbol{B}^\top\boldsymbol{B}$

  - $\boldsymbol{U},\ \boldsymbol{V}$: orthogonal matrices

    - $\boldsymbol{U}$ is eig.vecs of $\boldsymbol{BB}^\top$, $\boldsymbol{V}$ is eig.vecs of $\boldsymbol{B}^\top\boldsymbol{B}$ , $\mathbf{u}_i = \frac{1}{\sigma_i}\boldsymbol{B}^\top\mathbf{v}_i$

- Best rank-$K$ approximation problem of matrix $\boldsymbol{B}$:
$$\widehat{\boldsymbol{B}}^* = \mathrm{argmin}_{\widehat{\boldsymbol{B}}}\left\|\boldsymbol{B} - \widehat{\boldsymbol{B}}\right\|_\mathrm{F}^2 \text{ s.t. } \mathrm{rank}(\widehat{\boldsymbol{B}}) \leq K$$

  - Find $K$ largest singular values $\boldsymbol{\Sigma}^* = \mathrm{diag}(\sigma_1, \ldots, \sigma_K)$, and corresponding vectors $\boldsymbol{U}^* = (\mathbf{u}_1, \ldots, \mathbf{u}_K)$, $\boldsymbol{V}^* = (\mathbf{v}_1, \ldots, \mathbf{v}_K)$, and let $\widehat{\boldsymbol{B}}^* = \boldsymbol{U}^*\boldsymbol{\Sigma}^*\boldsymbol{V}^{*\top}$

# Solution of reduced rank regression (1/2): Best rank-$K$ approximation of a matrix

- Objective function to be minimized:

$$\|Y - XW\|_F^2 = \text{tr}\{(Y - XW)^\top (Y - XW)\}$$

$$= \text{tr}\{Y^\top Y - 2W^\top X^\top Y + W^\top X^\top XW\}$$

(Let $X^\top X = P^\top \Lambda P$ be the eigendecomposition)

$$\boxed{\begin{array}{c} P^\top P = PP^\top = I \\ (P\text{: orthogonal}) \end{array}} = \text{tr}\{Y^\top Y - 2\widetilde{W}^\top \Lambda^{-\frac{1}{2}} PX^\top Y + \widetilde{W}^\top \widetilde{W}\}$$

$$\text{where } \widetilde{W} = \Lambda^{\frac{1}{2}} PW$$

$$= \left\|\widetilde{W} - \Lambda^{-\frac{1}{2}} PX^\top Y\right\|_F^2 + \text{const.}$$

- Find the best rank-$K$ approximation of $\Lambda^{-\frac{1}{2}} PX^\top Y$
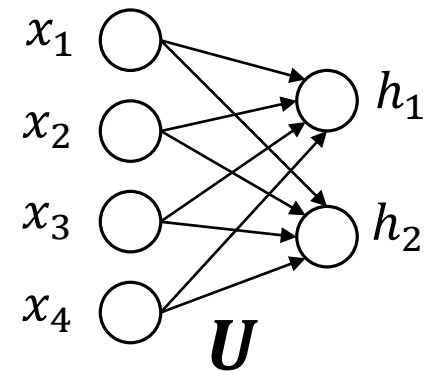
# Solution of reduced rank regression (2/2): Closed form solution using SVD

- The best rank-$K$ approximation of $\boldsymbol{\Lambda}^{-\frac{1}{2}}\boldsymbol{P}\boldsymbol{X}^\top\boldsymbol{Y}$ is given as $\widetilde{\boldsymbol{W}}^* = \boldsymbol{U}^*\boldsymbol{\Sigma}^*\boldsymbol{V}^{*\top}$

  - $\boldsymbol{V}^*$ is top-$K$ eigenvectors of
  $$\boldsymbol{Y}^\top\boldsymbol{X}\boldsymbol{P}^\top\boldsymbol{\Lambda}^{-\frac{1}{2}}\boldsymbol{\Lambda}^{-\frac{1}{2}}\boldsymbol{P}\boldsymbol{X}^\top\boldsymbol{Y} = \boldsymbol{Y}^\top\boldsymbol{X}(\boldsymbol{X}^\top\boldsymbol{X})^{-1}\boldsymbol{X}^\top\boldsymbol{Y}$$

  - $\boldsymbol{\Sigma}^*$: a diagonal matrix with $K$ largest singular values

  - $\boldsymbol{U}^* = \boldsymbol{\Lambda}^{-\frac{1}{2}}\boldsymbol{P}\boldsymbol{X}^\top\boldsymbol{Y}\,\boldsymbol{V}^*\boldsymbol{\Sigma}^{*-1}$

- The solution is $\boldsymbol{W}^* = \boldsymbol{P}^\top\boldsymbol{\Lambda}^{-\frac{1}{2}}\widetilde{\boldsymbol{W}}^* = \boldsymbol{P}^\top\boldsymbol{\Lambda}^{-\frac{1}{2}}\,\boldsymbol{U}^*\boldsymbol{\Sigma}^*\boldsymbol{V}^{*\top} = (\boldsymbol{X}^\top\boldsymbol{X})^{-1}\boldsymbol{X}^\top\boldsymbol{Y}\boldsymbol{V}^*\boldsymbol{V}^{*\top}$

# Dimension reduction:
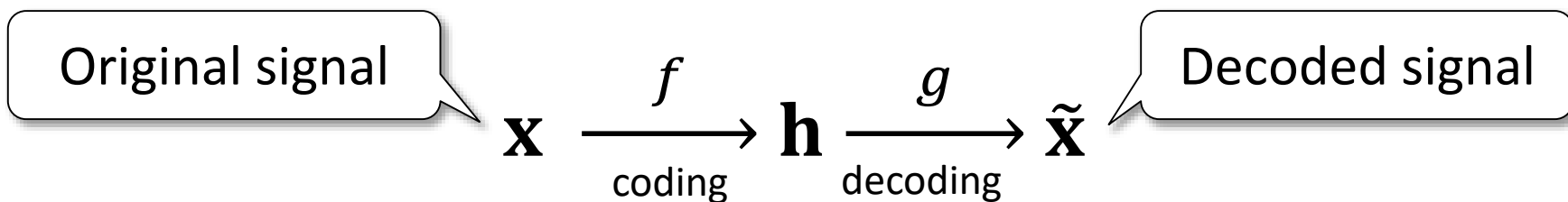## Find low-dimensional representations of high-dim. data

- Dimension reduction:

  - Find a low-dimensional mapping $f\colon \mathbb{R}^D \Rightarrow \mathbb{R}^K \ (D > K)$

    - for interpretability, computational/space efficiency, generalization abilities, …

    - (Lossy) compression: keep the original information as much as possible

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

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

# Basic idea behind dimension reduction:
## Find a coding & decoding function for lossy compression
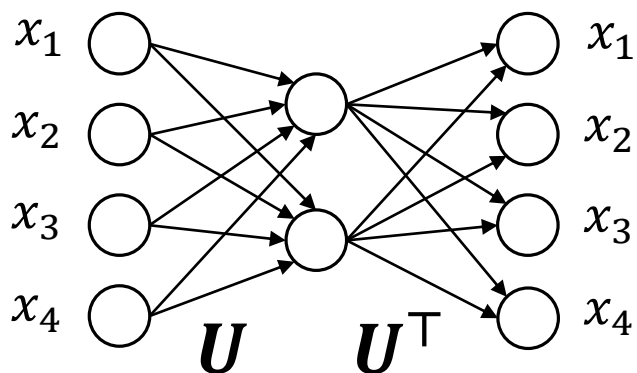
- Coding and decoding process:

Original signal

Decoded signal

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

- If $f$ and $g$ are appropriately designed so that $\mathbf{x} \doteqdot \tilde{\mathbf{x}}$, $\mathbf{h}$ must be a good low-dimensional representation of $\mathbf{x}$

- Optimization problem:

$$(f,g) = \text{argmin}_{f,g} \sum_{i=1}^{N} \text{loss}\left(\mathbf{x}^{(i)}, \underbrace{g(f(\mathbf{x}^{(i)}))}_{\tilde{\mathbf{x}}^{(i)}}\right)$$

# Principal component analysis:
## Dimension reduction using reduced rank regression

- Linear dimension reduction with coding & decoding functions

  - linear coding function $f : \mathbf{h} = \boldsymbol{U}^\top \mathbf{x}$ ($\boldsymbol{U} : D \times K$ matrix)

  - linear decoding function $g : \tilde{\mathbf{x}} = \boldsymbol{V}\mathbf{h}$ ($\boldsymbol{V} : K \times D$ matrix )

  - $\tilde{\mathbf{x}} = \boldsymbol{V}\boldsymbol{U}^\top \mathbf{x}$

- Reduced rank regression finds the solution by taking the training dataset as $\left\{ \left(\mathbf{x}^{(1)}, \mathbf{x}^{(1)}\right), \dots, \left(\mathbf{x}^{(N)}, \mathbf{x}^{(N)}\right)\right\}$

  - Solution will be $\boldsymbol{V} = \boldsymbol{U}$

# Topics:
## Learning with sparsity

- $L_1$-regularization & Lasso

  – Sparsity in terms of number of features used in the model

  – Solution of Lasso: Coordinate-wise descent

- Reduced rank regression

  – Sparsity in terms of number of dimensions of latent feature space

  – Solution using SVD

  – Dimension reduction