

Section 8

Convex Optimisation 2

Lagrangian

Consider a general optimization problem

$$\begin{aligned} \min_{\mathbf{x}} \quad & f(\mathbf{x}) \\ \text{subject to} \quad & h_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, m, \\ & \ell_j(\mathbf{x}) = 0, \quad j = 1, \dots, r. \end{aligned}$$

The objective function f needs not to be convex. Of course we pay special attention to the convex case.

Definition 8.1 (Lagrangian)

$$L(\mathbf{x}, \mathbf{u}, \mathbf{v}) = f(\mathbf{x}) + \sum_{i=1}^m u_i h_i(\mathbf{x}) + \sum_{j=1}^r v_j \ell_j(\mathbf{x}).$$

Here $\mathbf{u} \in \mathbb{R}^m$, $\mathbf{v} \in \mathbb{R}^r$, and $\mathbf{u} \geq \mathbf{0}$.

Lagrange Dual Function

Definition 8.2 (Lagrange Dual Function)

unconstrained
opt. problem
respect to λ .

$$g(\mathbf{u}, \mathbf{v}) := \min_{\mathbf{x} \in \mathbb{R}^n} L(\mathbf{x}, \mathbf{u}, \mathbf{v}) = \min_{\mathbf{x} \in \mathbb{R}^n} \left[f(\mathbf{x}) + \sum_{i=1}^m u_i h_i(\mathbf{x}) + \sum_{j=1}^r v_j \ell_j(\mathbf{x}) \right].$$

- ▶ For every feasible \mathbf{x} ($\mathbf{x} \in \mathcal{X}$), $L(\mathbf{x}, \mathbf{u}, \mathbf{v}) \leq f(\mathbf{x})$

$$L(\mathbf{x}, \mathbf{u}, \mathbf{v}) = f(\mathbf{x}) + \underbrace{\sum_{i=1}^m u_i h_i(\mathbf{x})}_{\leq 0} + \underbrace{\sum_{j=1}^r v_j \ell_j(\mathbf{x})}_{=0}.$$

- ▶ Let \mathcal{X} denote the primal feasible set.

$$\rightarrow \boxed{g(\mathbf{u}, \mathbf{v}) = \min_{\mathbf{x} \in \mathbb{R}^n} L(\mathbf{x}, \mathbf{u}, \mathbf{v}) \leq \min_{\mathbf{x} \in \mathcal{X}} L(\mathbf{x}, \mathbf{u}, \mathbf{v}) \leq f(\mathbf{x})}. \quad (13)$$

Concavity of Lagrange Dual Function

Lemma 8.3


The Lagrange dual function

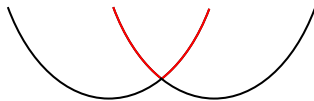
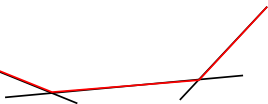
$$g(\mathbf{u}, \mathbf{v}) = \min_{\mathbf{x} \in \mathbb{R}^n} L(\mathbf{x}, \mathbf{u}, \mathbf{v}) = \min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) + \sum_{i=1}^m u_i h_i(\mathbf{x}) + \sum_{j=1}^r v_j \ell_j(\mathbf{x})$$

is concave in (\mathbf{u}, \mathbf{v}) .

Lemma 8.4

- ▶ Let $f_\alpha(x)$ be concave functions. Then $g(x) = \inf_\alpha f_\alpha(x)$ is concave.
- ▶ Let $f_\alpha(x)$ be convex functions. Then $g(x) = \sup_\alpha f_\alpha(x)$ is convex.

many functions. 



Proofs

Proof of Lemma 8.4: For any $\lambda \in [0, 1]$,

$$\begin{aligned} g(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) &= \inf_{\alpha} f_{\alpha}(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) \\ &\geq \inf_{\alpha} \lambda f_{\alpha}(\mathbf{x}) + (1 - \lambda) f_{\alpha}(\mathbf{y}) \\ &\geq \lambda \inf_{\alpha} f_{\alpha}(\mathbf{x}) + (1 - \lambda) \inf_{\alpha} f_{\alpha}(\mathbf{y}). \end{aligned}$$

assume f_{α} is concave.


Proof of Lemma 8.3: For any given \mathbf{x} , $L(\mathbf{x}, \mathbf{u}, \mathbf{v})$ is linear in (\mathbf{u}, \mathbf{v}) , and hence concave in (\mathbf{u}, \mathbf{v}) . The minimum of concave functions is concave based on Lemma 8.4.

Lagrange Dual Problem

Given the primal problem

$$\begin{aligned} \min_{\mathbf{x}} \quad & f(\mathbf{x}) \\ \text{subject to} \quad & h_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, m, \\ & \ell_j(\mathbf{x}) = 0, \quad j = 1, \dots, r. \end{aligned}$$

Its Lagrange dual problem is


$$\max_{\mathbf{u} \in \mathbb{R}^m, \mathbf{v} \in \mathbb{R}^r} g(\mathbf{u}, \mathbf{v}), \text{ subject to } \mathbf{u} \geq \mathbf{0}.$$

Concave
function.

Weak and Strong Duality

Weak duality: the dual optimal value g^* satisfies

$$f^* \geq g^*.$$

This is a direct consequence of (13).

Strong duality is referred to as the case that

$$f^* = g^*.$$

the optimum value
of the two problems
are equal

Slater's condition: if the primal is a convex problem (i.e., f and g_i 's are convex and ℓ_j 's are affine), and there exists at least one strictly feasible $\mathbf{x} \in \mathbb{R}^n$ satisfying

$$h_i(\mathbf{x}) < 0, \forall i \in [m], \text{ and } \ell_j(\mathbf{x}) = 0, \forall j \in [r],$$

then strong duality holds. (Proof is omitted.)

Karush-Kuhn-Tucker conditions

Given the optimization problem

$$\begin{aligned} & \text{minimize } f(\mathbf{x}) \\ & \text{subject to } h_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, m, \\ & \quad \quad \ell_j(\mathbf{x}) = 0, \quad i = 1, \dots, r. \end{aligned}$$

the optimum point satisfies the KKT conditions

The **Karush-Kuhn-Tucker (KKT) conditions** are:

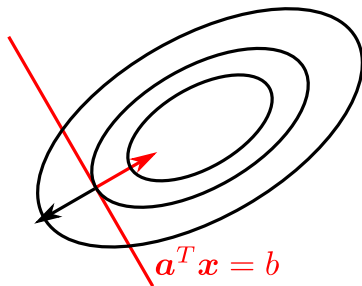
- ▶ $\mathbf{0} \in \partial f(\mathbf{x}) + \sum_{i=1}^m u_i \partial h_i(\mathbf{x}) + \sum_{j=1}^r v_j \partial \ell_j(\mathbf{x})$. (stationarity)
- ▶ $u_i h_i(\mathbf{x}) = 0, \forall i$. \leftarrow most useful!!! (complementary slackness)
- ▶ $h_i(\mathbf{x}) \leq 0, \ell_j(\mathbf{x}) = 0, \forall i, \forall j$. \leftarrow constraint, (primal feasibility)
- ▶ $u_i \geq 0, \forall i$. (dual feasibility)

\leftarrow Lagrange multipliers
KKT conditions are

- ▶ Always sufficient.
- ▶ Necessary under strong duality.

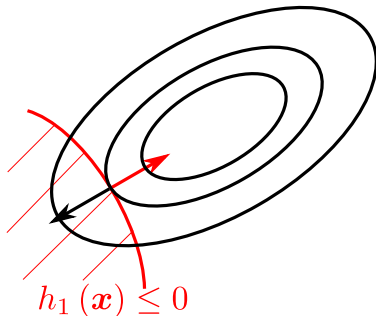
convex functions.

Geometric Intuition: Equality Constraints



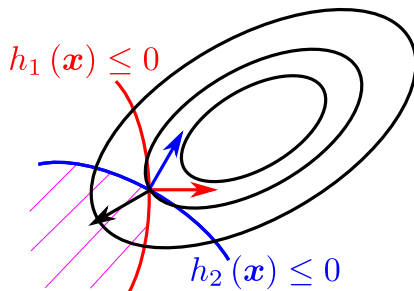
$\partial f(x)$ is a linear combination of $\partial \ell_j(x)$'s.

Geometric Intuition: One Inequality Constraint



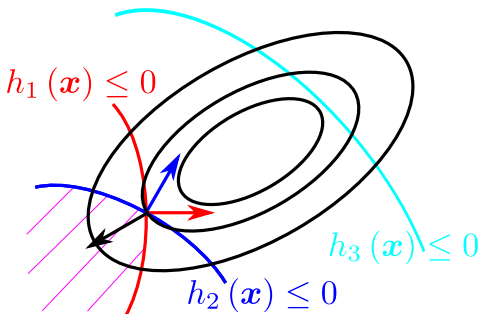
$$\begin{aligned}\partial f(\mathbf{x}) + u_1 \partial h_1(\mathbf{x}) &= \mathbf{0}. \\ h_1(\mathbf{x}) &= 0.\end{aligned}$$

Geometric Intuition: Inequality Constraints



$$\begin{aligned}\partial f(\mathbf{x}) + \sum_{i=1}^2 u_i \partial h_i(\mathbf{x}) &= \mathbf{0}. \\ h_1(\mathbf{x}) &= 0, \quad h_2(\mathbf{x}) = 0.\end{aligned}$$

Geometric Intuition: Inequality Constraints



$$\partial f(\mathbf{x}) + \sum_{i=1}^3 u_i \partial h_i(\mathbf{x}) = \mathbf{0}.$$

$$h_1(\mathbf{x}) = 0, h_2(\mathbf{x}) = 0,$$

$$h_3(\mathbf{x}) < 0 \text{ but } u_3 = 0 \text{ so that } u_3 h_3(\mathbf{x}) = 0. !!!$$

Sufficiency

If $\mathbf{x}^*, \mathbf{u}^*, \mathbf{v}^*$ satisfy the KKT conditions, then \mathbf{x}^* and $\mathbf{u}^*, \mathbf{v}^*$ are primal and dual solutions.

If $\mathbf{x}^*, \mathbf{u}^*, \mathbf{v}^*$ satisfy the KKT conditions, then

$$\begin{aligned} g(\mathbf{u}^*, \mathbf{v}^*) &= f(\mathbf{x}^*) + \sum_{i=1}^m u_i^* h_i(\mathbf{x}^*) + \sum_{j=1}^r v_j^* \ell_j(\mathbf{x}^*) \\ &= f(\mathbf{x}^*), \end{aligned}$$

where the first equality follows from stationarity, and the second follows from complementary slackness. This equality suggests the duality gap is zero. Hence, $\mathbf{x}^*, \mathbf{u}^*$ and \mathbf{v}^* are primal and dual optimal.

Necessity

Suppose that the strong duality holds and that \mathbf{x}^* and $\mathbf{u}^*, \mathbf{v}^*$ are primal and dual solutions. Then $\mathbf{x}^*, \mathbf{u}^*, \mathbf{v}^*$ satisfy the KKT conditions.

Due to the strong duality, one has

$$\begin{aligned} f(\mathbf{x}^*) &= g(\mathbf{u}^*, \mathbf{v}^*) \\ &= \min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) + \sum_{i=1}^m u_i^* h_i(\mathbf{x}) + \sum_{j=1}^r v_j^* \ell_j(\mathbf{x}) \\ &\leq f(\mathbf{x}^*) + \sum_{i=1}^m u_i^* h_i(\mathbf{x}^*) + \sum_{j=1}^r v_j^* \ell_j(\mathbf{x}^*) \\ &\leq f(\mathbf{x}^*). \end{aligned}$$

In other words, all the inequalities are actually equalities.

Quadratic Programming with Equality Constraints

Let $\mathbf{Q} \succeq 0$.

$$\min_{\mathbf{x}} \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} + \mathbf{c}^T \mathbf{x} \text{ subject to } \mathbf{A} \mathbf{x} = \mathbf{0}.$$

By KKT conditions, \mathbf{x} is the minimizer if and only if

$$\begin{bmatrix} \mathbf{Q} & \mathbf{A}^T \\ \mathbf{A} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{u} \end{bmatrix} = \begin{bmatrix} -\mathbf{c} \\ \mathbf{0} \end{bmatrix},$$

where the first set of linear equations come from the stationarity and the second set follows from the primal feasibility.

The optimal \mathbf{x}^* can be obtained by solving the linear inverse problem.

Water Filling

channel capacity of the channel.

$$\min_{\mathbf{x}} - \sum_{i=1}^n \log(\alpha_i + x_i) \text{ subject to } \mathbf{x} \geq \mathbf{0}, \mathbf{1}^T \mathbf{x} = 1.$$

$$\mathcal{L} = - \sum_{i=1}^n \log(\alpha_i + x_i) + \sum_i u_i (-x_i) + v(\mathbf{1}^T \mathbf{x} - 1).$$

By KKT conditions,

- ▶ $-1/(\alpha_i + x_i) - u_i + v = 0, \forall i$
- ▶ $u_i x_i = 0, \forall i$
- ▶ $\mathbf{x} \geq \mathbf{0}, \mathbf{1}^T \mathbf{x} = 1, \mathbf{u} \geq \mathbf{0}.$

Eliminate \mathbf{u} . The first two conditions become

$$1/(\alpha_i + x_i) \leq v, \text{ and } x_i(v - 1/(\alpha_i + x_i)) = 0, \forall i.$$

Therefore, the solution:

$$x_i = \max(0, 1/v - \alpha_i)$$

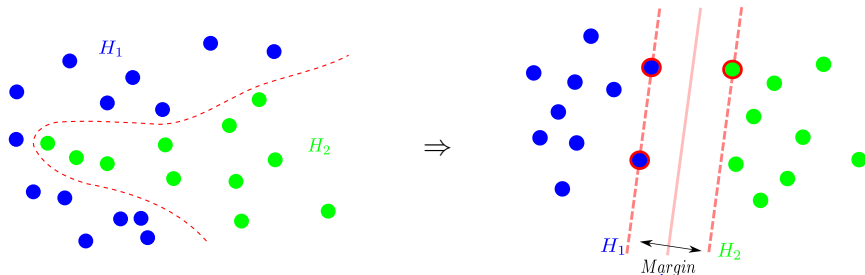
where v is chosen such that

$$\sum_{i=1}^n \max(0, 1/v - \alpha_i) = 1.$$

Section 9

Support Vector Machine

Idea of SVM



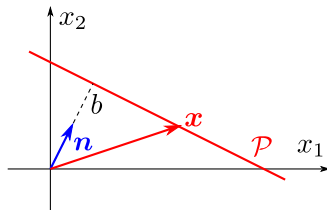
The hyperplane
may not be unique.
We pick the one with the
largest margin.

A Hyperplane

A hyperplane in \mathbb{R}^n can be defined using its normal vector $\mathbf{n} \in \mathbb{R}^n$:

$$\mathcal{P} = \{ \mathbf{x} : \mathbf{n}^T \mathbf{x} = b \}.$$

- Usually we assume $\|\mathbf{n}\|_2 = 1$.



The projection $\|\text{Proj}(\mathbf{x}, \text{span}(\mathbf{n}))\|_2 = b$.

- If $\|\mathbf{n}\|_2 \neq 1$, then

$$\mathcal{P} = \{ \mathbf{x} : \mathbf{n}^T \mathbf{x} = b \} = \{ \mathbf{x} : \mathbf{n}^T \mathbf{x} / \|\mathbf{n}\|_2 = b / \|\mathbf{n}\|_2 \}.$$

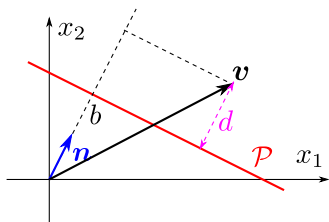
Distance to a Hyperplane

Define a hyperplane $\mathcal{P} = \{x : n^T x = b\}$ where $\|n\|_2 = 1$.

Let v be an arbitrary point.

The distance between v and \mathcal{P} is given by

$$d = d(v, \mathcal{P}) = |n^T v - b|. \quad (14)$$



When $\|n\|_2 \neq 1$,

$$d = \left| \frac{n^T}{\|n\|_2} v - \frac{b}{\|n\|_2} \right| = \frac{|n^T v - b|}{\|n\|_2} \quad \text{minimize this} \quad (15)$$

SVM: Separate Points from Two Different Classes



Given training dataset $\{\mathbf{x}_i, y_i\}$ where the labels $y_i \in \{-1, 1\}$, want to find β and b s.t.

$$\begin{aligned}\beta^T \mathbf{x}_i + b &\geq +1 & \text{for } y_i = +1, \\ \beta^T \mathbf{x}_i + b &\leq -1 & \text{for } y_i = -1.\end{aligned}$$

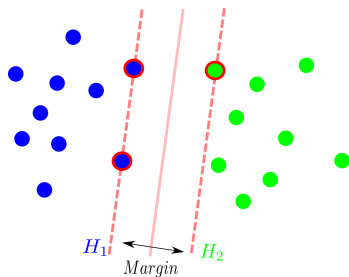
or equivalently

$$y_i (\beta^T \mathbf{x}_i + b) - 1 \geq 0, \quad \forall i.$$

In other words, find a hyperplane $\{\mathbf{x} : \beta^T \mathbf{x} - b\}$ s.t.

- 
- ▶ Distance from one class to the hyperplane is $1 / \|\beta\|_2$.
 - ▶ Distance between the two classes (along direction β) is $2 / \|\beta\|_2$.
- 

SVM: Best Separation



SVM: a convex optimization problem:

$$\begin{aligned} \min_{\beta, b} \quad & \frac{1}{2} \|\beta\|_2^2, \\ \text{subject to} \quad & 1 - y_i (\beta^T \mathbf{x}_i + b) \leq 0. \end{aligned} \tag{16}$$

Lagrange Dual Problem of SVM

Lagrangian of the SVM primal optimization problem:

$$L = \frac{1}{2} \|\beta\|^2 + \sum_i \lambda_i (1 - y_i (\beta^T \mathbf{x}_i + b)), \quad (17)$$

where $\lambda_i \geq 0$.

Lagrange Dual Problem

 \max_{λ} $\min_{\beta, b} L$ 

Lagrange dual function L_D

The Dual Function

To solve $\min_{\beta, b} L$, set $\partial L / \partial \beta$ and $\partial L / \partial b$ to zero:

$$\frac{\partial L}{\partial \beta} = \beta - \sum_i \lambda_i y_i \mathbf{x}_i = 0 \Rightarrow \beta = \sum_i \lambda_i y_i \mathbf{x}_i. \quad (18)$$

$$\frac{\partial L}{\partial b} = \sum_i \lambda_i y_i = 0. \quad (19)$$

Substitute (18) and (19) into the Lagrangian (17). It holds that

lagrange dual function.

$$\begin{aligned} L_D &= \sum_i \lambda_i - \frac{1}{2} \|\beta\|_2^2 = \sum_i \lambda_i - \frac{1}{2} \sum_{i,j} \lambda_i \lambda_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \\ &= -\frac{1}{2} \boldsymbol{\lambda}^T \mathbf{K} \boldsymbol{\lambda} + \mathbf{1}^T \boldsymbol{\lambda}, \end{aligned} \quad (20)$$

where $K_{i,j} = y_i \mathbf{x}_i^T \mathbf{x}_j y_j$.

is known.

The Dual Problem

The dual problem becomes:

$$\begin{aligned} \max_{\boldsymbol{\lambda}} \quad & -\frac{1}{2}\boldsymbol{\lambda}^T \mathbf{K} \boldsymbol{\lambda} + \mathbf{1}^T \boldsymbol{\lambda}, \\ \text{subject to} \quad & \lambda_i \geq 0, \quad \forall i, \\ & \sum_i \lambda_i y_i = 0. \end{aligned} \tag{21}$$

The KKT Condition

$$\frac{\partial L}{\partial \beta} = \beta - \sum_i \lambda_i y_i \mathbf{x}_i = 0, \quad (22)$$

$$\frac{\partial L}{\partial b} = \sum_i \lambda_i y_i = 0, \quad (23)$$

$$1 - y_i (\beta^T \mathbf{x}_i + b) \leq 0, \quad (24)$$

$$\lambda_i \geq 0, \quad (25)$$

$$\lambda_i (1 - y_i (\beta^T \mathbf{x}_i + b)) = 0. \quad (26)$$

↗
second bullet
of KKT conditions.

SVM Classifier: Support Vectors

Condition (26) implies

$$\begin{cases} \text{if } \lambda_i \neq 0 & \text{then } 1 = y_i (\beta^T \mathbf{x}_i + b), \\ \text{if } 1 \neq y_i (\beta^T \mathbf{x}_i + b) & \text{then } \lambda_i = 0. \end{cases}$$

If I can derive this slide I will have understood the concept.

Hence from (22),

$$\beta = \sum_{i \in \mathcal{I}} \lambda_i y_i \mathbf{x}_i, \quad \mathcal{I} = \{i : y_i (\beta^T \mathbf{x}_i + b) = 1 \quad (\text{or } \lambda_i \neq 0)\}.$$

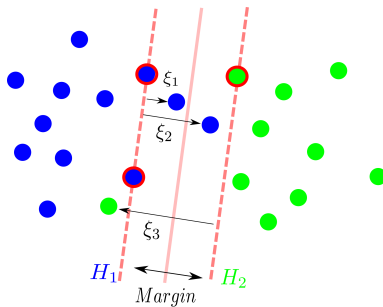
$\lambda_i \neq 0$ only for boundary points.

For a new test data \mathbf{x}^{new} ,

$$\hookrightarrow y^{\text{new}} = \text{sign} \left(\sum_{i \in \mathcal{I}} \lambda_i y_i \mathbf{x}_i^T \mathbf{x}^{\text{new}} + b \right).$$

The classifier only uses the boundary points (**sparsity!**).

SVM for Overlapping Classes



Primal Problem for Overlapping Classes

The constraints:

$$\beta^T \mathbf{x}_i + b \geq +1 - \xi_i \quad \text{for } y_i = +1,$$

$$\beta^T \mathbf{x}_i + b \leq -1 + \xi_i \quad \text{for } y_i = -1,$$

where $\xi_i \geq 0, \forall i$. *for more data points $\xi_i = 0$,*

SVM Primal Problem:

$$\min_{\beta, b, \xi} \quad \frac{1}{2} \|\beta\|_2^2 + C \left(\sum_i \xi_i \right)^k$$

$$\begin{aligned} \text{subject to} \quad & 1 - \xi_i - y_i (\beta^T \mathbf{x}_i + b) \leq 0, \\ & -\xi_i \leq 0, \quad \forall i, \end{aligned}$$

where $C > 0$ is a constant and k is a positive integer. Usually $k = 1$.

Dual Function

The Lagrangian

$$L = \frac{1}{2} \|\boldsymbol{\beta}\|_2^2 + C \sum \xi_i + \sum \lambda_i (1 - \xi_i - y_i (\boldsymbol{\beta}^T \mathbf{x}_i - b)) - \sum u_i \xi_i,$$

where $\lambda_i \geq 0$, $u_i \geq 0$ are Lagrange multipliers.

The dual function

$$L_D = \min_{\boldsymbol{\beta}, b, \boldsymbol{\xi}} L.$$

To find the dual function

$$\frac{dL}{d\boldsymbol{\beta}} = 0 \quad \Rightarrow \quad \boldsymbol{\beta} = \sum \lambda_i y_i \mathbf{x}_i.$$

$$\frac{dL}{db} = 0 \quad \Rightarrow \quad \sum \lambda_i y_i = 0.$$

$$\frac{dL}{d\xi} = 0 \quad \Rightarrow \quad C - \lambda_i - u_i = 0 \quad \Rightarrow \quad \lambda_i = C - u_i \leq C.$$

The Dual Problem

The dual problem:

$$\begin{aligned} \max_{\boldsymbol{\lambda}} \quad & \sum \lambda_i - \frac{1}{2} \|\boldsymbol{\beta}\|_2^2 = -\frac{1}{2} \boldsymbol{\lambda}^T \mathbf{K} \boldsymbol{\lambda} + \mathbf{1}^T \boldsymbol{\lambda} \\ \text{subject to} \quad & 0 \leq \lambda_i \leq C, \\ & \sum \lambda_i y_i = 0, \end{aligned}$$

where $K_{i,j} = y_i \mathbf{x}_i^T \mathbf{x}_j y_j$.

The only difference is that now λ_i 's are upper bounded by C .

Again, only **boundary points** are involved.

$$\boldsymbol{\beta} = \sum_{i \in \mathcal{I}} \lambda_i y_i \mathbf{x}_i, \quad \mathcal{I} = \{i : \lambda_i \neq 0\},$$

which comes from the KKT condition $\lambda_i (1 - \xi_i - y_i (\boldsymbol{\beta}^T \mathbf{x}_i + b)) = 0$.

The General Case

- ▶ Two classes \Rightarrow multiple classes
 - ▶ Regression
- ▶ Data space \Rightarrow feature space

Define a **kernel** function $\varphi : \mathbb{R}^n \rightarrow \mathcal{H}$ and work on the space of $\varphi(\mathbf{x}_i)$.

In SVM, what really matters is $\mathbf{x}_i^T \mathbf{x}_j$.

In the general case (**kernel method**), what matters is

$$\kappa(\mathbf{x}_i, \mathbf{x}_j) = \varphi^T(\mathbf{x}_i) \varphi(\mathbf{x}_j).$$

Example of nonlinear features:

- ▶ $\kappa(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^T \mathbf{y})^2 = (x_1 y_1 + x_2 y_2)^2 = x_1^2 y_2^2 + 2x_1 y_1 x_2 y_2 + x_2^2 y_2^2.$

- ▶ $\kappa(\mathbf{x}, \mathbf{y}) = \varphi^T(\mathbf{x}) \varphi(\mathbf{y})$ with $\varphi(\mathbf{x}) = [x_1^2, \sqrt{2}x_1 x_2, x_2^2]^T$.

- ▶ $\kappa(\mathbf{x}, \mathbf{y}) = \exp\left(-\|\mathbf{x} - \mathbf{y}\|_2^2 / 2\sigma^2\right).$

- ▶ $\varphi(\mathbf{x})$ has infinite dimension.

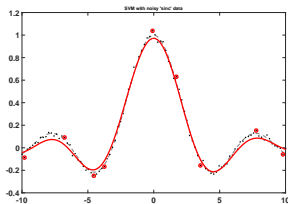
σ small \Rightarrow watch the nearest neighbours.

SVM for Regression

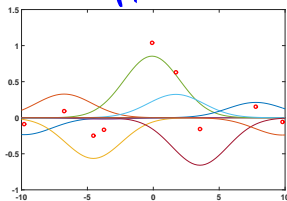
Regression problem: find β and b s.t.

$$\begin{aligned} y_i &= f(\mathbf{x}_i) = \beta^T \varphi(\mathbf{x}_i) + b \\ &= \sum_j \lambda'_j \underbrace{\varphi^T(\mathbf{x}_j) \varphi(\mathbf{x}_i)} + b \\ &= \sum_j \lambda'_j \underbrace{\kappa(\mathbf{x}_i, \mathbf{x}_j)} + b. \end{aligned}$$

We do not care about φ , but about κ , the inner product.



=



The Primal Optimization Problem

Let $\epsilon > 0$ be the error tolerance. Then one has

$$\begin{aligned} \min \quad & \frac{1}{2} \|\boldsymbol{\beta}\|_2^2 \\ \text{subject to} \quad & |y_i - \boldsymbol{\beta}^T \varphi(\mathbf{x}_i) - b| \leq \epsilon. \end{aligned}$$

The constraints are equivalent to

$$\begin{aligned} y_i - \boldsymbol{\beta}^T \varphi(\mathbf{x}_i) - b &\leq \epsilon, \\ \boldsymbol{\beta}^T \varphi(\mathbf{x}_i) + b - y_i &\leq \epsilon. \end{aligned}$$

Now if we allow additional noise, represented by $\xi_i \geq 0$ and $\xi_i^* \geq 0$. Then

$$\begin{aligned} \min \quad & \frac{1}{2} \|\boldsymbol{\beta}\|_2^2 + C \sum (\xi_i + \xi_i^*) \\ \text{subject to} \quad & y_i - \boldsymbol{\beta}^T \varphi(\mathbf{x}_i) - b \leq \epsilon + \xi_i, \\ & \boldsymbol{\beta}^T \varphi(\mathbf{x}_i) + b - y_i \leq \epsilon + \xi_i^*, \\ & -\xi_i \leq 0, \quad -\xi_i^* \leq 0. \end{aligned}$$

Lagrangian

$$\begin{aligned} L = & \frac{1}{2} \|\boldsymbol{\beta}\|_2^2 + C \sum (\xi_i + \xi_i^*) - \sum_i \left(u_i \xi_i + \sum u_i^* \xi_i^* \right) \\ & + \lambda_i \left(y_i - \boldsymbol{\beta}^T \boldsymbol{\varphi}(\mathbf{x}_i) - b - \epsilon - \xi_i \right) \\ & + \lambda_i^* \left(\boldsymbol{\beta}^T \boldsymbol{\varphi}(\mathbf{x}_i) + b - y_i - \epsilon - \xi_i^* \right), \end{aligned}$$

where $u_i, u_i^*, \xi_i, \xi_i^* \geq 0$ are Lagrange multiplier. To minimize L ,

$$\begin{aligned} dL/d\boldsymbol{\beta} = \mathbf{0} & \Rightarrow \boldsymbol{\beta} = \sum_i (\lambda_i - \lambda_i^*) \boldsymbol{\varphi}(\boldsymbol{\xi}_i), \\ dL/db = \mathbf{0} & \Rightarrow \sum \lambda_i = \sum \lambda_i^*, \\ dL/d\xi_i = 0, dL/d\xi_i^* = 0 & \Rightarrow \lambda_i \leq C, \lambda_i^* \leq C. \end{aligned}$$

The Dual Problem

The **objective function** of the dual problem

$$L_D = -\epsilon \sum (\lambda_i + \lambda_i^*) + y_i \sum (\lambda_i - \lambda_i^*) \\ - \underbrace{\frac{1}{2} \sum_{i,j} (\lambda_i - \lambda_i^*) (\lambda_j - \lambda_j^*) \kappa(\mathbf{x}_i, \mathbf{x}_j)}_{\|\boldsymbol{\beta}\|_2^2},$$

where $\kappa(\mathbf{x}_i, \mathbf{x}_j) = \varphi^T(\mathbf{x}_i) \varphi(\mathbf{x}_j)$.

The optimizatoin **constraints** are

$$\sum (\lambda_i - \lambda_i^*) = 0, \\ 0 \leq \lambda_i, \lambda_i^* \leq C.$$

KKT Condition and Support Vectors

Part of the KKT condition is that $\forall i$,

$$\begin{cases} \lambda_i (y_i - \beta^T \varphi(\mathbf{x}_i) - b - \epsilon - \xi_i) = 0, \\ \lambda_i^* (\beta^T \varphi(\mathbf{x}_i) + b - y_i - \epsilon - \xi_i^*) = 0. \end{cases}$$

- ▶ Interior points: $|y_i - \beta^T \varphi(\mathbf{x}_i) - b| < \epsilon + \xi_i$.
 - ▶ Both λ_i and λ_i^* are zero.
- ▶ Boundary points: $|y_i - \beta^T \varphi(\mathbf{x}_i) - b| = \epsilon + \xi_i$.
 - ▶ One of λ_i and λ_i^* is zero.
 - ▶ $\lambda_i \neq \lambda_i^*$.

The Standard Form

Let $\gamma_i = \lambda_i$ and $\gamma_{i+n} = \lambda_i^*$ (Merge λ and λ^* into a single vector).
The dual problem becomes

$$\begin{aligned} \min_{\gamma} \quad & \frac{1}{2} \gamma^T Q \gamma + v^T \gamma, \\ \text{subject to} \quad & 0 \leq \gamma_i \leq C, \quad \sum_{i=1}^n \gamma_i - \sum_{i=n+1}^{2n} \gamma_i = 0. \end{aligned}$$

The **boundary points** are given by $\mathcal{I} = \{i : \gamma_i - \gamma_{i+n} \neq 0\}$.

For a new data point \mathbf{x}^{new} , the regression is

$$f(\mathbf{x}^{\text{new}}) = \sum_i (\gamma_i - \gamma_{i+n}) \kappa(\mathbf{x}_i, \mathbf{x}^{\text{new}}) + b.$$

Section 10

Gaussian Distribution

Statistical learning.

Gaussian Random Vectors

A random vector $\mathbf{X} \in \mathbb{R}^n$ is Gaussian distributed $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ if its pdf is given by

$$p(\mathbf{x}) = |2\pi\boldsymbol{\Sigma}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right),$$

where $\boldsymbol{\mu} \in \mathbb{R}^n$ and $\boldsymbol{\Sigma} \in \mathcal{S}_+^n$ (the set of $n \times n$ symmetric positive semidefinite matrices).

Here, we have assumed that $\boldsymbol{\Sigma}$ is invertible (of full rank).

Gaussian Random Vectors: Characteristic Function

$$\text{PDF} \xrightleftharpoons[\text{Inverse Fourier Transform}]{\text{Fourier Transform}} \text{Characteristic function } \mathbb{E} \left[e^{i \langle \boldsymbol{\lambda}, \mathbf{X} \rangle} \right].$$

$\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ if

$$\mathbb{E} \left[e^{i \langle \boldsymbol{\lambda}, \mathbf{X} \rangle} \right] = \exp \left(i \langle \boldsymbol{\lambda}, \boldsymbol{\mu} \rangle - \frac{1}{2} \boldsymbol{\lambda}^T \boldsymbol{\Sigma} \boldsymbol{\lambda} \right).$$

It is well defined even when $\boldsymbol{\Sigma}$ is not invertible.

Quadratic form in both the pdf
and characteristic function.

Affine Transformation

Lemma 10.1

Let $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. Then for any $\mathbf{A} \in \mathbb{R}^{m \times n}$ and $\mathbf{b} \in \mathbb{R}^m$,

$$\hookrightarrow \mathbf{AX} + \mathbf{b} \sim \mathcal{N}(\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T).$$

Proof:

$$\begin{aligned} \mathbb{E} \left[e^{i\langle \boldsymbol{\lambda}, \mathbf{AX} + \mathbf{b} \rangle} \right] &= \mathbb{E} \left[e^{i\langle \mathbf{A}^T \boldsymbol{\lambda}, \mathbf{X} \rangle + i\langle \boldsymbol{\lambda}, \mathbf{b} \rangle} \right] \\ &= \exp \left(i\langle \mathbf{A}^T \boldsymbol{\lambda}, \boldsymbol{\mu} \rangle - \frac{1}{2} (\mathbf{A}^T \boldsymbol{\lambda})^T \boldsymbol{\Sigma} (\mathbf{A}^T \boldsymbol{\lambda}) \right) e^{i\langle \boldsymbol{\lambda}, \mathbf{b} \rangle} \\ &= \exp \left(i\langle \boldsymbol{\lambda}, \mathbf{A}^T \boldsymbol{\mu} + \mathbf{b} \rangle - \frac{1}{2} \boldsymbol{\lambda}^T (\mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T) \boldsymbol{\lambda} \right). \end{aligned}$$

Gaussian Conditioning Lemma

Lemma 10.2

Let $\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \Sigma)$.

Let \mathbf{X}_A and \mathbf{X}_B be two subvectors of \mathbf{X} , i.e., $\mathbf{X} = [\mathbf{X}_A^T, \mathbf{X}_B^T]^T$.

Let $\mathbf{K} := \Sigma^{-1} = \begin{bmatrix} \mathbf{K}_{AA} & \mathbf{K}_{AB} \\ \mathbf{K}_{BA} & \mathbf{K}_{BB} \end{bmatrix}$ be the *precision matrix*.

Then $\mathbf{X}_A | \mathbf{X}_B \sim P_{\mathbf{X}_A | \mathbf{X}_B} = \mathcal{N}(-\mathbf{K}_{AA}^{-1} \mathbf{K}_{AB} \mathbf{X}_B, \mathbf{K}_{AA}^{-1})$. In other words,

$$\mathbf{X}_A = -\mathbf{K}_{AA}^{-1} \mathbf{K}_{AB} \mathbf{X}_B + \epsilon,$$

where $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_{AA}^{-1})$ is independent of \mathbf{X}_B .

Remark: $\mathbf{K}_{AA}^{-1} \neq \Sigma_{AA}$.

Matrix Identities

- ▶ Block matrix inverse (BMI)

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}^{-1} = \begin{bmatrix} (\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1} & -\mathbf{A}^{-1}\mathbf{B}(\mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1} \\ -\mathbf{D}^{-1}\mathbf{C}(\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1} & (\mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1} \end{bmatrix}. \quad (27)$$

- ▶ Woodbury matrix identity (WMI)

$$(\mathbf{A} + \mathbf{UCV})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{U}(\mathbf{C}^{-1} + \mathbf{VA}^{-1}\mathbf{U})^{-1}\mathbf{VA}^{-1}. \quad (28)$$

Proof of Gaussian Conditioning Lemma

By Bayes rule, $p(\mathbf{x}_A|\mathbf{x}_B) = p(\mathbf{x}_A, \mathbf{x}_B) / p(\mathbf{x}_B)$. Then

$$\begin{aligned}\ln p(\mathbf{x}_A|\mathbf{x}_B) &= \ln p(\mathbf{x}_A, \mathbf{x}_B) - \ln p(\mathbf{x}_B) \\ &= c - \frac{1}{2} \mathbf{x}_A^T \mathbf{K}_{AA} \mathbf{x}_A - \mathbf{x}_A^T \mathbf{K}_{AB} \mathbf{x}_B - \frac{1}{2} \mathbf{x}_B^T (\mathbf{K}_{BB} - \Sigma_{BB}^{-1}) \mathbf{x}_B,\end{aligned}$$

where c is a constant. By (27),

$$\Sigma_{BB}^{-1} = \mathbf{K}_{BB} - \mathbf{K}_{BA} \mathbf{K}_{AA}^{-1} \mathbf{K}_{AB}.$$

One has

$$\ln p(\mathbf{x}_A|\mathbf{x}_B) = c - \frac{1}{2} (\mathbf{x}_A + \mathbf{K}_{AA}^{-1} \mathbf{K}_{AB} \mathbf{x}_B)^T \mathbf{K}_{AA} (\mathbf{x}_A + \mathbf{K}_{AA}^{-1} \mathbf{K}_{AB} \mathbf{x}_B).$$

That is, $\mathbf{X}_A|\mathbf{X}_B \sim \mathcal{N}(-\mathbf{K}_{AA}^{-1} \mathbf{K}_{AB} \mathbf{X}_B, \mathbf{K}_{AA}^{-1})$.

A Signal Processing Application

The problem:

Given

$$\mathbf{Y} = \mathbf{A}\mathbf{X} + \mathbf{W},$$

where $\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma}_x)$ and $\mathbf{W} \sim \mathcal{N}(\mathbf{0}, \mathbf{\Sigma}_w)$.

Given observation \mathbf{y} , want to find $\hat{\mathbf{x}} = f(\mathbf{y})$ s.t. the mean squared error $\mathbb{E} [\|\mathbf{x} - \hat{\mathbf{x}}\|_2^2]$ is minimized (MMSE solution).

Fact: The general MMSE solution is given by

$$\hat{\mathbf{x}} = \mathbb{E} [\mathbf{X} | \mathbf{Y} = \mathbf{y}] = \int \mathbf{x} \cdot p_{\mathbf{X}|\mathbf{Y}}(\mathbf{x}|\mathbf{y}) d\mathbf{x}.$$

Hence for Gaussian random variables, Gaussian conditioning lemma can be used.

Finding the MMSE Solution

1. $\mathbf{Y} = \mathbf{A}\mathbf{X} + \mathbf{W}$ is Gaussian distributed $\mathcal{N}(\mathbf{0}, \mathbf{A}\Sigma_x\mathbf{A}^T + \Sigma_w)$.
- 2.

$$\begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \underbrace{\begin{bmatrix} \Sigma_x & \Sigma_x\mathbf{A}^T \\ \mathbf{A}\Sigma_x & \mathbf{A}\Sigma_x\mathbf{A}^T + \Sigma_w \end{bmatrix}}_{\Sigma}\right).$$

3. Find the precision matrix from Σ :

$$\mathbf{K} = \begin{bmatrix} \Sigma_x^{-1} + \mathbf{A}^T \Sigma_w^{-1} \mathbf{A} & -\mathbf{A}^T \Sigma_w^{-1} \\ -\Sigma_w^{-T} \mathbf{A} & \text{sth} \end{bmatrix}$$

4. $\mathbf{X}|\mathbf{Y} \sim \mathcal{N}(-\mathbf{K}_{\mathcal{A}\mathcal{A}}^{-1}\mathbf{K}_{\mathcal{A}\mathcal{B}}\mathbf{Y}, \mathbf{K}_{\mathcal{A}\mathcal{A}}^{-1})$ by Gaussian Conditioning Lemma.

We use the conditional mean as the estimate $\hat{\mathbf{x}}$:

$$\hat{\mathbf{x}} = (\Sigma_x^{-1} + \mathbf{A}^T \Sigma_w^{-1} \mathbf{A})^{-1} \mathbf{A}^T \Sigma_w^{-1} \mathbf{y}. \quad (29)$$

$$\Sigma_{\mathbf{X}|\mathbf{Y}} = \mathbf{K}_{\mathcal{A}\mathcal{A}}^{-1} = (\Sigma_x^{-1} + \mathbf{A}^T \Sigma_w^{-1} \mathbf{A})^{-1}. \quad (30)$$

Calculation of The \mathbf{K} Matrix

$$\begin{aligned}\mathbf{K}_{\mathcal{AA}} &\stackrel{\text{BMI(27)}}{=} \left(\Sigma_x - \Sigma_x \mathbf{A}^T (\mathbf{A} \Sigma_x \mathbf{A}^T + \Sigma_w)^{-1} \mathbf{A} \Sigma_x \right)^{-1} \\ &\stackrel{\text{WMI(28)}}{=} \left((\Sigma_x^{-1} + \mathbf{A}^T \Sigma_w^{-1} \mathbf{A})^{-1} \right)^{-1} \\ &= \Sigma_x^{-1} + \mathbf{A}^T \Sigma_w^{-1} \mathbf{A}.\end{aligned}$$

$$\begin{aligned}\mathbf{K}_{\mathcal{AB}} &\stackrel{\text{BMI(27)}}{=} -\Sigma_x^{-1} (\Sigma_x \mathbf{A}^T) (\mathbf{A} \Sigma_x \mathbf{A}^T + \Sigma_w - \mathbf{A} \Sigma_x \Sigma_x^{-1} \Sigma_x \mathbf{A}^T)^{-1} \\ &= -\mathbf{A}^T \Sigma_w^{-1}.\end{aligned}$$

Hence $\Sigma_{X|Y} = (\Sigma_x^{-1} + \mathbf{A}^T \Sigma_w^{-1} \mathbf{A})^{-1}$ and $\mathbf{L} = \Sigma_{X|Y} \mathbf{A}^T \Sigma_w^{-1}$.

Section 11

Sparse Bayesian Learning

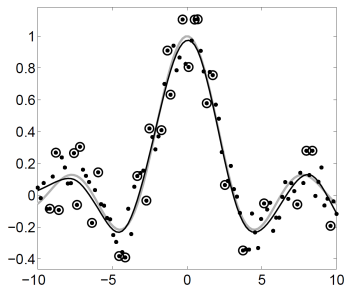
Didn't do it.

Sparse Bayesian Learning (SBL) via Relevance Vector Machine (RVM)

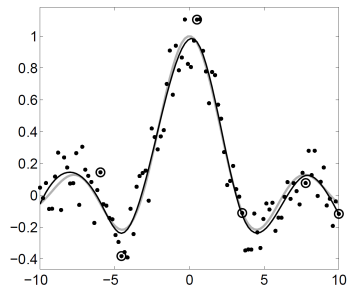
M. E. Tipping, "Sparse Bayesian learning and the relevance vector machine," J. Mach. Learn. Res., JMLR.org, 2001, 1, 211-244.

Examples from Tipping's paper:

To approximate the sinc function



Support vector approx.



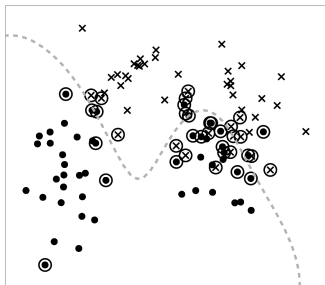
Relevance vector approx.

Sparse Bayesian Learning (SBL) via Relevance Vector Machine (RVM)

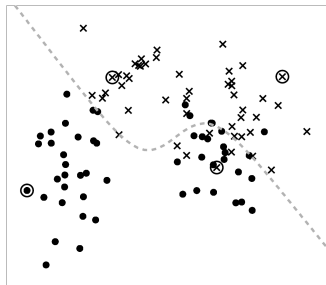
M. E. Tipping, "Sparse Bayesian learning and the relevance vector machine," J. Mach. Learn. Res., JMLR.org, 2001, 1, 211-244.

Examples from Tipping's paper:

Classify two classes



SVM



RVM

RVM: A Hierarchical Gaussian Model

Target: Solve the sparse linear inverse problem:

$$\mathbf{y} = \Phi \mathbf{x} + \mathbf{w}.$$

Signal model: signal components are i.i.d. Gaussian

$$p(X_i | \alpha_i) = \mathcal{N}(X_i; 0, \alpha_i^{-1}).$$

- ▶ $\alpha_i > 0$: precision ($\alpha_i^{-1} = \sigma_i^2$, for mathematical convenience)
- ▶ $\alpha_i = \infty$: $X_i \sim \mathcal{N}(0, 0)$ and $X_i = 0$ with prob. one.
- ▶ $\alpha_i \in \mathbb{R}^+$: $X_i \sim \mathcal{N}(0, \alpha_i^{-1})$ and $X_i \neq 0$ with prob. one.

Signal Estimation with Known α

If we know α , then the posterior of \mathbf{x} can be derived by Bayes' rule

$$p(\mathbf{x}|\mathbf{y}) = \frac{p(\mathbf{x}) p(\mathbf{y}|\mathbf{x})}{p(\mathbf{y})} = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}),$$

where

$$\boldsymbol{\mu} = \sigma_w^{-2} \boldsymbol{\Sigma} \boldsymbol{\Phi}^T \mathbf{y}, \quad \boldsymbol{\Sigma} = (\mathbf{A} + \sigma_w^{-2} \boldsymbol{\Phi}^T \boldsymbol{\Phi})^{-1}, \quad \mathbf{A} = \text{diag}(\boldsymbol{\alpha}).$$

See (29) and (30) for details.

Set $\hat{\mathbf{x}} = \boldsymbol{\mu}$. This is the MMSE solution as well.

Signal Recovery with Unknown α

The key is to estimate α :

- ▶ $\mathbf{Y} \sim \mathcal{N}(\mathbf{0}, \sigma_w^2 \mathbf{I} + \Phi \mathbf{A}^{-1} \Phi^T)$
- ▶ Linear combination of Gaussian is Gaussian ($\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{w}$).
- ▶ Maximum likelihood (ML) estimate:

$$\mathcal{L}(\alpha) = \log p(\mathbf{y} | \alpha, \sigma_w^2) = -\frac{1}{2} [N \log 2\pi + \log |\mathbf{C}| + \mathbf{y}^T \mathbf{C}^{-1} \mathbf{y}],$$

where $\mathbf{C} = \sigma_w^2 \mathbf{I} + \Phi \mathbf{A}^{-1} \Phi^T$.

- ▶ **Marginal likelihood maximization:** For a given i , we solve

$$\max_{\alpha_i} \mathcal{L}(\alpha)$$

by fixing all α_j 's, $j \neq i$.

Marginal Likelihood Maximization (1)

$$\begin{aligned} \mathbf{C} &= \sigma^2 \mathbf{I} + \mathbf{\Phi} \mathbf{A}^{-1} \mathbf{\Phi} \\ &= \sigma^2 \mathbf{I} + \sum_{j \neq i} \alpha_j^{-1} \phi_j \phi_j^T + \alpha_i^{-1} \phi_i \phi_i^T \\ &= \mathbf{C}_{-i} + \alpha_i^{-1} \phi_i \phi_i^T. \end{aligned}$$

One has

$$\begin{aligned} |\mathbf{C}| &= |\mathbf{C}_{-i}| \cdot \left| \mathbf{I} + \alpha_i^{-1} \mathbf{C}_{-i}^{-1/2} \phi_i \phi_i^T \mathbf{C}_{-i}^{-1/2} \right| \\ &= |\mathbf{C}_{-i}| \cdot \left| 1 + \alpha_i^{-1} \phi_i^T \mathbf{C}_{-i}^{-1} \phi_i \right|. \end{aligned}$$

The last line comes from $|\mathbf{I} + \mathbf{M} \mathbf{M}^T| = |\mathbf{I} + \mathbf{M}^T \mathbf{M}|$, which can be easily verified by the SVD of \mathbf{M} .

$$\mathbf{C}_{-i}^{-1} \stackrel{\text{WMI(28)}}{=} \mathbf{C}_{-i}^{-1} - \frac{\mathbf{C}_{-i}^{-1} \phi_i \phi_i^T \mathbf{C}_{-i}^{-1}}{\alpha_i + \phi_i^T \mathbf{C}_{-i}^{-1} \phi_i}.$$

Marginal Likelihood Maximization (2)

$$\begin{aligned}\mathcal{L}(\boldsymbol{\alpha}) &= -\frac{1}{2} \left[N \log 2\pi + \log |\mathbf{C}| + \mathbf{y}^T \mathbf{C}^{-1} \mathbf{y} \right], \\ &= c - \frac{1}{2} \left[\log |\mathbf{C}_{-i}| + \mathbf{y}^T \mathbf{C}_{-i}^{-1} \mathbf{y} \right. \\ &\quad \left. - \log \alpha_i + \log (\alpha_i + \boldsymbol{\phi}_i^T \mathbf{C}_{-i}^{-1} \boldsymbol{\phi}_i) - \frac{(\boldsymbol{\phi}_i^T \mathbf{C}_{-i}^{-1} \mathbf{y})^2}{\alpha_i + \boldsymbol{\phi}_i^T \mathbf{C}_{-i}^{-1} \boldsymbol{\phi}_i} \right] \\ &= \mathcal{L}(\boldsymbol{\alpha}_{-i}) + \frac{1}{2} \left[\log \alpha_i - \log (\alpha_i + s_i) - \frac{q_i^2}{\alpha_i + s_i} \right] \\ &= \mathcal{L}(\boldsymbol{\alpha}_{-i}) + \ell(\alpha_i),\end{aligned}$$

where for simplification of expressions, we have defined

$$s_i := \boldsymbol{\phi}_i^T \mathbf{C}_{-i}^{-1} \boldsymbol{\phi}_i, \text{ and } q_i := \boldsymbol{\phi}_i^T \mathbf{C}_{-i}^{-1} \mathbf{y}.$$

Set the derivative of $\ell(\alpha_i)$ to zero. One obtains the closed form solution for the optimal α_i :

$$\alpha_i = \begin{cases} \frac{s_i^2}{q_i^2 - s_i} & \text{if } q_i^2 > s_i, \\ \infty & \text{if } q_i^2 \leq s_i, \end{cases}$$

Algorithms for RVM

A sequential algorithm: [Tipping & Faul, 2003]

In each iteration:

- ▶ Scan $i \in [N]$ and find the i that $\mathcal{L}(\alpha_i^*) - \mathcal{L}(\alpha_i)$ is maximized.
- ▶ Set $\alpha_i = \alpha_i^*$.
- ▶ Update \mathcal{L} and other parameters (preparation for the next iteration).

Connections to Greedy Algorithms

- ▶ Similar to greedy algorithms
 - ▶ The sequential algorithm is very similar to OMP.
 - ▶ “Subspace pursuit” type of RVM: [Karsenas & D., 2013]
 - ▶ Performance guarantees: *sufficient* conditions based on mutual coherence or RIP [Karsenas & D., 2013]
- ▶ Different from greedy algorithms
 - ▶ Noise variance is considered.
 - ▶ Statistical information: estimation quality information

A Demo

Track global Ozone density

- ▶ Spatial sparsity:
 - ▶ Each frame is sparse under DCT transform
- ▶ Temporal sparsity: correlations across days
 - ▶ Support innovation.
 - ▶ Nonzero component magnitude innovation.

Let \mathbf{x}^t be the DCT coefficient vector at the t -th day:

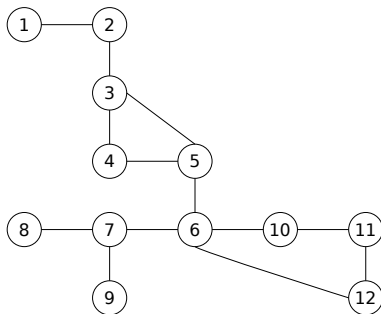
$$\begin{aligned}\mathbf{x}^t &= \mathbf{x}^{t-1} + \mathbf{u}^t, \\ \mathbf{y}^t &= \mathbf{A}\mathbf{D}\mathbf{x}^t + \mathbf{w}^t,\end{aligned}$$

where $\mathbf{u}^t \sim \prod \mathcal{N}(0, \alpha_i^{-1})$, \mathbf{A} is the sampling matrix, and \mathbf{w}^t is the white Gaussian noise.

Section 12

Gaussian Graphic Model

Motivation: Gaussian Graphic Model



Encoding the **conditional dependencies** between n random variables X_1, \dots, X_n by a graph.

Correlation and Conditional Independence

Sneeze — Catch Cold — Weather Change

Observation: “Weather Change” and “Sneeze” are correlated.

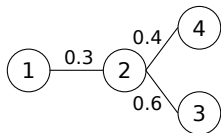
- ▶ “Weather Change” and “Catch Cold” are highly correlated.
- ▶ “Catch Cold” and “Sneeze” are highly correlated.

However, given the status of “Catch Cold”, “Weather Change” and “Sneeze” are independent.

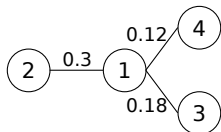
- ▶ Given that “Catch Cold” is false, “Sneeze” is likely to be false, independent of whether “Weather Change” is true or not.
- ▶ Given that “Catch Cold” is true, “Sneeze” is likely to be true, independent of whether “Weather Change” is true or not.

Other Examples

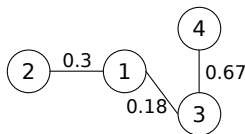
Suppose that $\rho(X_1, X_2) = 0.3$, $\rho(X_1, X_3) = 0.18$, and $\rho(X_1, X_4) = 0.12$. Suppose that on one day, $X_2 \uparrow 0.2$, $X_3 \downarrow 0.1$, and $X_4 \downarrow 0.5$. Find the expected change of X_1 .



$$E[\Delta X_1] = 0.2 \times 0.3 = 0.06.$$

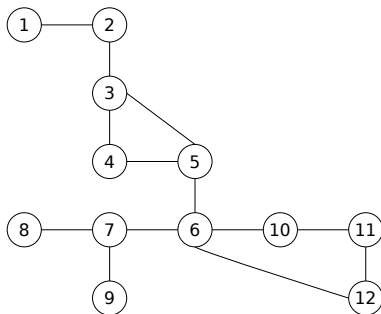


$$\begin{aligned} E[\Delta X_1] &= 0.2 \times 0.3 - 0.1 \times 0.18 - 0.5 \times 0.12 \\ &= -0.018. \end{aligned}$$



$$\begin{aligned} E[\Delta X_1] &= 0.2 \times 0.3 - 0.1 \times 0.18 \\ &= 0.042. \end{aligned}$$

Nondirected Graphical Model



The distribution of the Gaussian random vector $\mathbf{X} = [X_1, \dots, X_n]^T$ is a graphic model according to the graph g if

for all a : $X_a \perp \{X_b : b \notin \text{ne}(a), b \neq a\}$ given $\{X_c : c \in \text{ne}(a)\}$.

Or, given X_c 's, $c \in \text{ne}(a)$, X_a and X_b 's are independent for all b not in the neighborhood.

Consequence of Gaussian Conditioning

Recall the Gaussian conditioning lemma (Lemma 10.2).
Let \mathbf{K} be the precision matrix of \mathbf{X} .

Corollary 12.1

For any $a \in [n]$,

$$X_a = - \sum_{b: b \neq a} \frac{K_{ab}}{K_{aa}} X_b + \epsilon_a,$$

where $\epsilon_a \sim \mathcal{N}(0, K_{aa}^{-1})$ is independent of $\{X_b : b \neq a\}$.

Proof: Apply Lemma 10.2 with $A = \{a\}$ and $B = [n] \setminus \{a\} = A^c$.

Remark: Find the neighboring points.

Conditional Correlation

Corollary 12.2

$$\text{cor}(X_a, X_b | \mathbf{X}_C) = -\frac{K_{ab}}{\sqrt{K_{aa}K_{bb}}}.$$

Proof: From Gaussian Conditioning (Lemma 10.2), it holds that

$$\text{cov}(\mathbf{X}_{\{a,b\}} | \mathbf{X}_C) = \begin{bmatrix} K_{aa} & K_{ab} \\ K_{ba} & K_{bb} \end{bmatrix}^{-1} = \frac{1}{K_{aa}K_{bb} - K_{ab}^2} \begin{bmatrix} K_{bb} & -K_{ba} \\ -K_{ab} & K_{aa} \end{bmatrix}.$$

Plug this formula into the definition of conditional correlation. Corollary 12.2 is proved.

Remark: Find the correlation between neighboring points.

Estimate the Precision Matrix

From the definition $\mathbf{K} = \Sigma^{-1}$, the computation seems straightforward. However, the commonly used fact

$$\frac{1}{m} \sum (\mathbf{X} - \bar{\mathbf{X}}) (\mathbf{X} - \bar{\mathbf{X}})^T \rightarrow \Sigma \quad (31)$$

is based on the assumption that n is fixed and $m \rightarrow \infty$.

In reality, we may not have sufficient data m . Hence (31) may not be applicable.

Assumption: \mathbf{K} is sparse.

Estimation via Regression (1)

Define the matrix Θ by $\theta_{ab} = -K_{ab}/K_{bb}$ for $b \neq a$ and $\theta_{aa} = 0$. Then Corollary 12.1 implies

$$\mathbb{E}[X_a | X_b : b \neq a] = \sum_b \theta_{ba} X_b.$$

Hence we need to find θ_{ba} 's ($b \neq a$) to minimize

$$\mathbb{E} \left[\left(X_a - \sum_b \theta_{ba} X_b \right)^2 \right].$$

Or in matrix format

$$\hat{\Theta} = \arg \min_{\Theta \in \Theta} \mathbb{E} \left[\|\mathbf{X} - \Theta^T \mathbf{X}\|_2^2 \right],$$

where $\Theta = \{\Theta : \text{diag}(\Theta) = \mathbf{0}\}$.

Estimation via Regression (2)

The objective function can be rewritten as

$$\begin{aligned} \mathbb{E} \left[\|\mathbf{X} - \boldsymbol{\Theta}^T \mathbf{X}\|_2^2 \right] &\approx \frac{1}{m} \sum (\mathbf{x} - \boldsymbol{\Theta}^T \mathbf{x})^T (\mathbf{x} - \boldsymbol{\Theta}^T \mathbf{x}) \\ &= \frac{1}{m} \left\| \begin{bmatrix} \mathbf{x}_{(1)}^T \\ \vdots \\ \mathbf{x}_{(m)}^T \end{bmatrix} - \begin{bmatrix} \mathbf{x}_{(1)}^T \\ \vdots \\ \mathbf{x}_{(m)}^T \end{bmatrix} \boldsymbol{\Theta} \right\|_F^2 \\ &= \frac{1}{m} \|\mathbf{X} - \mathbf{X} \boldsymbol{\Theta}\|_F^2. \end{aligned}$$

Note that the \mathbf{X} on this slide is the data matrix and the \mathbf{X} on previous slides are random vectors.

Estimation via Regression (3)

The overall optimization problem:

$$\min_{\Theta \in \Theta} \quad \frac{1}{m} \|\mathbf{X} - \mathbf{X}\Theta\|_F^2 + \lambda \sum_{a \neq b} |\theta_{ab}|,$$

Or

$$\min_{\Theta \in \Theta} \quad \frac{1}{m} \|\mathbf{X} - \mathbf{X}\Theta\|_F^2 + \lambda \sum_{a < b} \sqrt{\theta_{ab}^2 + \theta_{ba}^2}.$$