



机器学习的数学原理

线性代数基本运算

1. 矩阵求导:

$$\frac{\partial(X^T A)}{\partial X} = A, \quad \frac{\partial(X A)}{\partial X} = A^T, \quad \frac{\partial(A X)}{\partial X^T} = A \quad \frac{\partial(A X)}{\partial X} = A^T$$

$$\frac{\partial(X^T A X)}{\partial X} = AX + A^T X = (A + A^T)X \quad \frac{\partial W^T W}{\partial W} = 2W$$

2. 求逆矩阵

$$A^{-1} = \frac{A^*}{|A|}, \quad A^* = \begin{bmatrix} A_{11} & A_{21} & \cdots & A_{n1} \\ A_{12} & A_{22} & \cdots & A_{n2} \\ \vdots & \vdots & \ddots & \vdots \\ A_{1N} & A_{2N} & \cdots & A_{nN} \end{bmatrix}$$

例: $A = \begin{bmatrix} 2 & 3 \\ 1 & 4 \end{bmatrix}, \quad A^* = \begin{bmatrix} 4 & -3 \\ -1 & 2 \end{bmatrix}$

$$\therefore A \cdot A^* = \begin{bmatrix} 2 & 3 \\ 1 & 4 \end{bmatrix} \begin{bmatrix} 4 & -3 \\ -1 & 2 \end{bmatrix} = \begin{bmatrix} 5 & 0 \\ 0 & 5 \end{bmatrix}$$

3. 求矩阵特征值与特征向量

$$AX = \lambda X \Rightarrow AX = \lambda EX \Rightarrow (A - \lambda E)X = 0$$

解 $|A - \lambda E| = 0$

根据得到的 λ , 利用 $(\lambda E - A)X = 0$

对 X 进行赋值得到 eigen value 与 eigen vector.

4. $\text{COV}(X, Y) = E(XY) - E(X)E(Y)$

Lecture 1. Linear Regression.

Given $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$, $x_i \in \mathbb{R}^d$, $y_i \in \mathbb{R}$

Linear Regression model: $f(x) = w^T x + b$. $w \in \mathbb{R}^d$.

* How to define $w, b \rightarrow$ by Empirical Risk Minimization
in Loss function (squared Loss)

$$L(w, b) = \min_{w, b} \sum_{i \in [n]} (f(x_i) - y_i)^2 = \min_{w, b} \left(\sum_{i \in [n]} (w^T x_i + b - y_i)^2 \right)$$

Note: 标量对向量求导是对应每一维求导。

$$\begin{aligned} \frac{\partial L}{\partial w} &= 2 \sum_{i \in [n]} (w^T x_i + b - y_i) x_i \\ \frac{\partial L}{\partial b} &= 2 \sum_{i \in [n]} (w^T x_i + b - y_i) \end{aligned} \Rightarrow \text{Here we can use gradient descent. } b \leftarrow b - \alpha \cdot \frac{\partial L}{\partial b}, w \leftarrow w - \alpha \cdot \frac{\partial L}{\partial w}$$

For Linear Regression, we can use closed form solution.
(Least Squares)

$$X = \begin{bmatrix} x_1^T, 1 \\ x_2^T, 1 \\ \vdots \\ x_n^T, 1 \end{bmatrix} \in \mathbb{R}^{n \times (d+1)}, \quad \hat{w} = \begin{bmatrix} w \\ b \end{bmatrix} \in \mathbb{R}^{d+1}, \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \in \mathbb{R}^n$$

$$L(\hat{w}) = (y - X\hat{w})^T (y - X\hat{w})$$

$$\therefore \frac{\partial L}{\partial \hat{w}} = \frac{\partial (y^T y - y^T X \hat{w} - \hat{w}^T X^T y + \hat{w}^T X^T X \hat{w})}{\partial \hat{w}}$$

$$= -X^T y - X^T \hat{w} + 2X^T X \hat{w}$$

$$= 2X^T X \hat{w} - 2X^T y$$

$$\text{Let } \frac{\partial L}{\partial \hat{w}} = 0 \Rightarrow X^T X \hat{w} = X^T y$$

$$\therefore \hat{w} = (X^T X)^{-1} X^T y.$$

Note that if we want to have this solution form, we need to $X^T X$ no singular. Later we give how to define a matrix singular.

$x^T x$ can be singular, when

① $d+1 > n$. $\text{Rank}(x^T x) = \text{Rank}(x) \leq \min(n, d+1) \leq n < d+1$
then not full rank (parameter > samples)

② X has repeated columns/rows?

$n \geq d+1$, $\text{rank}(x) < d+1$, not full rank.

Note: $x^T x = U \Lambda U^T = U \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix} U^T$, real-symmetric

$\Rightarrow \forall V \in \mathbb{R}^{d+1}, \forall V \in \mathbb{R}^{d+1}, V \neq 0, V^T x^T x V = \|Vx\|^2 \geq 0$.

if λ_{d+1} is close to 0, $(x^T x)^{-1} = U \Lambda^{-1} U^T = U \begin{bmatrix} \frac{1}{\lambda_1} & & \\ & \ddots & \\ & & \frac{1}{\lambda_n} \end{bmatrix} U^T$
 $x^T x$ close to singular

numerical instability.

We introduce L₂ Regularization. **Ridge**

$\min_{\hat{w}} L(\hat{w}) + \lambda \|\hat{w}\|^2$, where $\|\hat{w}\|^2 = \hat{w}^T \hat{w} = \sum_{i \in [n]} w_i^2$

constant parameter (超参数, $\lambda > 0$)

Let $J(\hat{w}) = (y - x\hat{w})^T (y - x\hat{w}) + \lambda \hat{w}^T \hat{w}$

$$\frac{\partial J}{\partial \hat{w}} = 2x^T x \hat{w} - 2x^T y + 2\lambda \hat{w} = 0$$

$$\Rightarrow (x^T x + \lambda) \hat{w} = x^T y$$

$$\hat{w} = (x^T x + \lambda)^{-1} x^T y.$$

There is another version of Regularization L₁,

L₁ regularization: $L(\hat{w}) + \lambda \|\hat{w}\|_1 = \sum_{i \in [n]} |w_i|$

will cause many of parameter closed to 0.

Sparse.

LASSO Regression

Lecture 2, Logistic Regression

Apply to: Binary classification, $y \in \{0, 1\}$, $x \in \mathbb{R}^d$.

$$f(x) = w^T x + b, w \in \mathbb{R}^d, b \in \mathbb{R}, f(x) \in \mathbb{R} \rightarrow \{0, 1\}$$

Here we use Sigmoid function.

$$\sigma(z) = \frac{1}{1 + e^{-z}}, \sigma(w^T x + b) = \frac{1}{1 + e^{-(w^T x + b)}} \in (0, 1)$$

$$\therefore P(y=1|x) = \sigma(w^T x + b)$$

$$P(y=0|x) = 1 - \sigma(w^T x + b)$$

$$\therefore \text{MLE} = \prod_{i \in [n]} P(y_i = y_i | x = x_i)$$

$$= \prod_{i \in [n]} [\sigma(w^T x_i + b)]^{y_i} [1 - \sigma(w^T x_i + b)]^{1-y_i}$$

$$\therefore \log \text{MLE} = \sum_{i \in [n]} y_i \log (\sigma(w^T x_i + b)) + \sum_{i \in [n]} (1-y_i) \log (1 - \sigma(w^T x_i + b))$$

$$\Leftrightarrow \text{minimize } - \left(\sum_{i \in [n]} y_i \log (\sigma(w^T x_i + b)) + \sum_{i \in [n]} (1-y_i) \log (1 - \sigma(w^T x_i + b)) \right)$$

Cross Entropy

Question is: How to optimize w, b ?

$$\hat{x} = \begin{bmatrix} x \\ 1 \end{bmatrix} \in \mathbb{R}^{d+1}, \hat{w} = \begin{bmatrix} w \\ b \end{bmatrix} \in \mathbb{R}^{d+1} \quad 1 - \sigma(z) = \sigma(-z)$$

$$\begin{aligned} \text{CELOSS} &= - \left(\sum_{i \in [n]} y_i \cdot \log \frac{1}{1 + e^{-\hat{w}^T \hat{x}}} + \sum_{i \in [n]} (1-y_i) \cdot \log \frac{1}{1 + e^{\hat{w}^T \hat{x}}} \right) \\ &= - \left(\sum_{i \in [n]} y_i \cdot \log \frac{1 + e^{-\hat{w}^T \hat{x}}}{1 + e^{-\hat{w}^T \hat{x}}} + \sum_{i \in [n]} (1-y_i) \cdot \log \frac{1}{1 + e^{\hat{w}^T \hat{x}}} \right) \\ &= - \left(\sum_{i \in [n]} \left(y_i \cdot \log \frac{1 + e^{-\hat{w}^T \hat{x}}}{1 + e^{-\hat{w}^T \hat{x}}} - \log (1 + e^{\hat{w}^T \hat{x}}) \right) \right) \\ &= - \left(\sum_{i \in [n]} (y_i \cdot \hat{w}^T \hat{x} - \log (1 + e^{\hat{w}^T \hat{x}})) \right) \end{aligned}$$

$$\frac{\partial L}{\partial \hat{w}} = - \left(\sum_{i \in [n]} y_i \cdot \hat{x}_i - \frac{\hat{x}_i e^{\hat{w}^T \hat{x}_i}}{1 + e^{\hat{w}^T \hat{x}_i}} \right)$$

$$= - \left(\sum_{i \in [n]} (y_i - P(y=1|x_i)) \hat{x}_i \right)$$

$$\text{Let } \frac{\partial L}{\partial \hat{w}} = 0 \Rightarrow \text{when } y_i = P(y=1|x_i)$$

when $y \in \{-1, 1\}$,

$$(1) y = 1 \text{ 时}, P(y=1|x_i) = G(w^T x_i + b)$$

$$\text{此时 } L(x_i, y_i) = -\log \frac{1}{1 + e^{-(w^T x_i + b)}} = \log(1 + e^{-y_i(w^T x_i + b)})$$

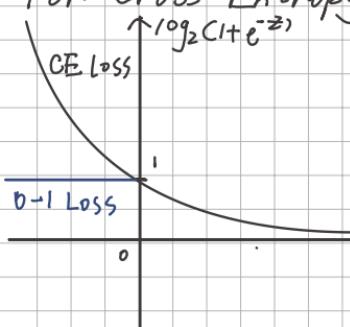
$$(2) y = -1 \text{ 时}, P(y=-1|x_i) = 1 - G(w^T x_i + b)$$

$$\text{此时 } L(x_i, y_i) = -\log \frac{1}{1 + e^{w^T x_i + b}} = \log(1 + e^{-y_i(w^T x_i + b)})$$

$$\therefore \text{无论 } y \text{ 取 } 1 \text{ 还是 } -1, L(w, b) = \sum_{i \in [n]} \log(1 + e^{-y_i(w^T x_i + b)})$$

$$\therefore L(w, b) = \sum_{i \in [n]} \log(1 + e^{-y_i(w^T x_i + b)})$$

For cross Entropy Loss



$$z_i = y_i(w^T x_i + b) \begin{cases} > 0, \text{ 分类 3} \\ < 0, \text{ 不分类} \end{cases}$$

CE Loss is an upper bound.

More common case

⇒ multi-class classification.

(1) Background: Assume $y \in \{1, 2, 3, \dots, k\}$, $x \in \mathbb{R}^d$.

(2) Method: Softmax Regression

(3) Mechanism:

① Define k models. $f_k(x) = w_k^T x + b_k, \forall k \in [k]$

② Find $f_k(x_i) > f_j(x_i), \forall j \neq k$, predict $y_i = k$.

(4) Details:

$$P(y=k|x) = \frac{e^{f_k(x)}}{\sum_{j \in [k]} e^{f_j(x)}} \geq 0 \quad \rightarrow \sum_j P(y=j|x) = 1$$

when $f_k(x) \geq f_j(x), P(y=k|x) \rightarrow 1$

$$\text{MLE: maximize } \sum_{i \in [n]} \log P(y=y_i|x=x_i) = \sum_{i \in [n]} \log \frac{e^{w_i^T x_i + b_i}}{\sum_{j \in [k]} e^{w_j^T x_i + b_j}}$$

when $k=2$, it seems like the logistic Regression.

And it doesn't matter what the label look like.

Assume $y \in \{1, 2\}$ $w_i^T x + b_i$,

$$P(y=1|x) = \frac{e^{w_1^T x + b_1}}{e^{w_1^T x + b_1} + e^{w_2^T x + b_2}} = \frac{1}{1 + e^{[(w_1^T - w_2^T)x + (b_1 - b_2)]}}$$

let $w = w_1^T - w_2^T$, $b = b_1 - b_2$

$$\therefore P(y=1|x) = \phi(wx+b)$$

Lecture 3 Support Vector Machine.

Background: Constrained Optimization

(1) Equality Constraint

Goal: $\min_x f(x)$, subject $h(x) = 0$



① A point x on the constraint surface $h(x) = 0$ we have $\nabla h(x)$ orthogonal to the surface

② For a local minimum X^* , $\nabla f(X^*)$ must be orthogonal to the surface.

To solve this problem, we introduce Lagrange function
 $L(x, \lambda) = f(x) + \lambda h(x)$ (λ is the Lagrange multiplier)

when we have more constraints (such as k)

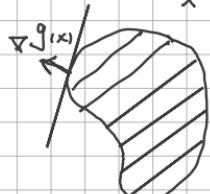
$$L(x, \lambda) = f(x) + \sum_{i=1}^k \lambda_i h_i(x)$$

if X^* is the local minimum.

we have $\begin{cases} \nabla_x L(X^*, \lambda) = 0 \text{ (because } f(X^*) + \lambda h(X^*) = 0 \\ \nabla_\lambda L(X^*, \lambda) = 0 \Rightarrow \text{because } h(X) = 0 \end{cases}$

(2) Inequality Constraints.

Goal: $\min_x f(x)$, subject $h(x) \leq 0$



① A point x on the surface, $\nabla g(x)$ orthogonal to the surface.

② For a local minimum X^*

(i) X^* is on the surface, $\nabla f(X^*)$ must be the same direction to the $\nabla g(x)$

$$\exists m > 0. \text{ s.t. } \nabla f(X^*) + m \cdot \nabla g(X^*) = 0$$

(ii) X^* is within the surface, that $g(X^*) < 0$
we only need $f(X^*) = 0$.

$$\exists m > 0. \text{ s.t. } \nabla f(X^*) + m \cdot \nabla g(X^*) = 0$$

Lagrange function: $L(x, \mu) = f(x) + \mu \cdot g(x)$

if x^* is local minimum $\Rightarrow \begin{cases} \nabla_x L(x^*, \mu) = 0 \\ \mu \geq 0 \\ g_i(x^*) \leq 0 \\ \mu \cdot g_i(x^*) = 0 \end{cases} \rightarrow \text{KKT conditions.}$

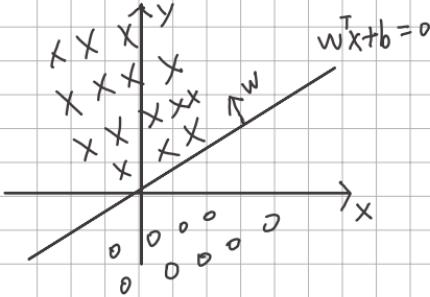
(3) The most common cases:

$$\min_x f(x), \text{ s.t. } h_i(x) = 0, i=1, 2, \dots, k \\ g_j(x) \leq 0, j=1, 2, \dots, l$$

$$L(x, \lambda, \mu) = f(x) + \sum_{i \in [k]} \lambda_i h_i(x) + \sum_{j \in [l]} \mu_j g_j(x)$$

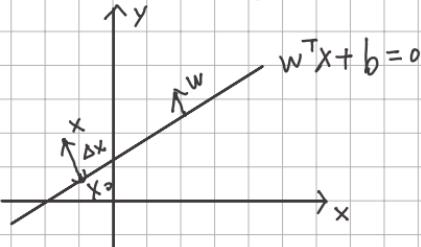
x^* is local minimum $\left\{ \begin{array}{l} \nabla_x L = 0 \\ h_i(x^*) = 0, g_j(x^*) \leq 0 \\ \mu_j \geq 0 \\ \sum \mu_j g_j(x) = 0 \end{array} \right.$

SVM Introduction



- c1) Goal: maximize margin
- c2) maximize margin (which refers to the minimum distance to the hyperplane $w^T x + b = 0$ over all $x_i, i \in [n]$)

distance calculation



$$\begin{aligned} & \left| \begin{array}{l} x_0 + \Delta x \cdot \frac{w}{\|w\|} = x \\ w^T x_0 + b = 0 \end{array} \right. \quad \textcircled{1} \\ & w^T x_0 + b = 0 \quad \textcircled{2} \\ & w^T x_0 + \Delta x \cdot \|w\| = w^T x \\ & \therefore \Delta x \cdot \|w\| = w^T x + b \\ & \Rightarrow \Delta x = \frac{w^T x + b}{\|w\|} \end{aligned}$$

So if linear separable,

$$\text{we have } r_i = g_i \cdot \Delta x_i = \frac{g_i(w^T x_i + b)}{\|w\|} \geq 0.$$

so we just need to minimize $r = \min_{i \in [n]} r_i = \min_{i \in [n]} \frac{y_i(c w_i^T x + b)}{\|w\|}$

which means that $\max_{w, b} \text{ s.t. } \frac{y_i(w_i^T x + b)}{\|w\|} \geq \lambda$

Assume there is (x_0, y_0) , $\lambda = \frac{y_0(c w_i^T x_0 + b)}{\|w\|}$, geometric margin

$\Rightarrow \|w\|$. we can get that $y_i(c w_i^T x + b) \geq y_0(c w_i^T x_0 + b)$
Let $1 \leftarrow \text{functional margin}$

So we can have:

$$\max_{w, b} \frac{1}{\|w\|}, \text{ s.t. } y_i(c w_i^T x + b) \geq 1$$

$$\Leftrightarrow \min_{w, b} \frac{1}{2} \|w\|^2, \text{ s.t. } y_i(c w_i^T x + b) \geq 1, \forall i \in [n]$$

Here we can use Lagrange Function to solve this prob.

$$L(c w, b, \alpha) = \frac{1}{2} \|w\|^2 + \sum_{i \in [n]} \alpha_i (1 - y_i(c w_i^T x_i + b)) \quad (\alpha_i \geq 0)$$

$$\text{So } p^* = \min_{w, b} \left(\max_{\alpha \geq 0} \frac{1}{2} \|w\|^2 + \sum_{i \in [n]} \alpha_i (1 - y_i(c w_i^T x_i + b)) \right)$$

Note that if $-y_i(c w_i^T x_i + b) > 0$

$$\alpha \rightarrow +\infty, \max_{\alpha} L(c w, b, \alpha) \rightarrow +\infty$$

So if p^* exists, $\exists y_i(c w_i^T x_i + b) < 0$.

we can use KKT to solve this problem

$$\sum \frac{\partial L}{\partial x} = 0$$

$1 - y_i(c w_i^T x_i + b) \leq 0 \Rightarrow$ But α_i is dependent on w, b
 $\alpha_i \geq 0$ so it's hard to determine α with w, b . We can try to
 $\sum \alpha_i y_i c x_i = 0$ solve the dual problem.

-----> Dual Problem (对偶问题)

In this problem, we can ensure $d^* = P^*$

$$\textcircled{1} \quad p^* = \min_{w, b} \max_{\alpha} \left(\frac{1}{2} \|w\|^2 + \sum_{i \in [n]} (1 - y_i(c w_i^T x_i + b)) \right)$$

$$\textcircled{2} \quad d^* = \max_{\alpha} \min_{w, b} \left(\frac{1}{2} \|w\|^2 + \sum_{i \in [n]} (1 - y_i(c w_i^T x_i + b)) \right)$$

$d^* = P^*$ Here. Next we try to solve this problem.

$$d^* = \max_{\alpha} \min_{w, b} \frac{1}{2} \|w\|^2 + \sum_{i \in [n]} \alpha_i (1 - y_i (w^\top x_i + b))$$

$$= \max_{\alpha} \min_{w, b} \left(\frac{1}{2} \|w\|^2 + \sum_{i \in [n]} \alpha_i - \sum_{i \in [n]} \alpha_i y_i (w^\top x_i + b) \right)$$

$$\frac{\partial L}{\partial w} = w - \sum_{i \in [n]} \alpha_i y_i x_i = 0 \Rightarrow w = \sum_{i \in [n]} \alpha_i y_i x_i$$

$$\frac{\partial L}{\partial b} = - \sum_{i \in [n]} \alpha_i y_i = 0 \Rightarrow \sum_{i \in [n]} \alpha_i y_i = 0.$$

Then substitute w into L ,

$$\begin{aligned} L(w, b, \alpha) &= \frac{1}{2} \left[\sum_{i \in [n]} \alpha_i y_i x_i^\top \right] \left[\sum_{j \in [n]} \alpha_j y_j x_j^\top \right] + \sum_{i \in [n]} \alpha_i - \sum_{i \in [n]} \alpha_i y_i \sim \\ &= \sum_{i \in [n]} \alpha_i - \frac{1}{2} \sum_{i \in [n]} \sum_{j \in [n]} \alpha_i \alpha_j y_i y_j x_i x_j^\top \end{aligned}$$

$$\text{外层: } \max_{\alpha \geq 0} \left(\sum_{i \in [n]} \alpha_i - \frac{1}{2} \sum_{i \in [n]} \sum_{j \in [n]} \alpha_i \alpha_j y_i y_j x_i x_j^\top \right)$$

$$\text{s.t. } \sum_{i \in [n]} \alpha_i y_i = 0$$

when we get α^* for d^* , we can then get w^*

$$\therefore w^* = \sum_{i \in [n]} \alpha_i^* y_i x_i$$

$$b^* = \frac{1}{y_0} - w^{*\top} x_j. \quad \forall (x_j, y_j) \text{ is a support vector.}$$

Questions is that:

How to determine support vector:

$$\forall i, 1 - y_i (w^\top x_i + b) = 0 \Rightarrow \alpha_i > 0 \text{ (active)}$$

$$1 - y_i (w^\top x_i + b) < 0 \Rightarrow \alpha_i = 0 \text{ (inactive)}$$

So, if $\alpha^* = 0$. not support vector

if $\alpha^* > 0$, then support vector.

when we encounter x . ↑

$$f(x) = \sum_i \alpha_i^* y_i x_i^\top x + b$$

we can use SMO get the optimized α^* .

In some cases, the points may be not linear separable, we may use kernel trick at this time.

Kernel Trick

⇒ Define the similarity of 2 points X, Z

① Linear kernel: $k(X, Z) = X^T Z$

② Polynomial kernel: $k(X, Z) = (X^T Z + 1)^P$

③ RBF kernel: $k(X, Z) = \exp\left(-\frac{\|X - Z\|^2}{2\sigma^2}\right)$

Kernel: Maps X to high-dimension
such as $X = (x_1, x_2)^T \in \mathbb{R}^2$

Define $\varphi(x) = (1, x_1, x_2, x_1^2, x_2^2, x_1 x_2) \in \mathbb{R}^6$

Dual Problem: 替換 $\varphi^T x_i \varphi^T z_j$

$$f(x) = \sum_{i \in \mathcal{D}_0} \alpha_i^* y_i \varphi^T x_i \varphi^T z_j + b^*$$

e.g. $X \in \mathbb{R}^2$, $X = (x_1, x_2)^T$, $Z = (z_1, z_2)^T$

$$(X^T Z + 1)^2 = (x_1 z_1 + x_2 z_2 + 1)^2 = x_1^2 z_1^2 + x_2^2 z_2^2 + \dots$$

$$= (1, \sqrt{z_1}, \sqrt{z_2}, z_1^2, z_2^2, \sqrt{z_1 z_2})^T$$

$$(1, \sqrt{z_1}, \sqrt{z_2}, z_1^2, z_2^2, \sqrt{z_1 z_2})^T$$

Gaussian Kernel:

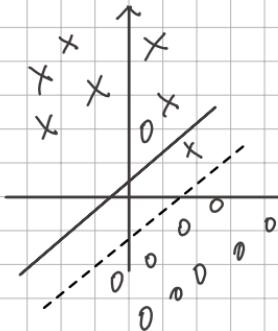
$$k(X, Z) = \exp\left(-\frac{\|X - Z\|^2}{2\sigma^2}\right) = \exp\left(-\frac{\|X\|^2}{2\sigma^2}\right) \cdot \exp\left(-\frac{\|Z\|^2}{2\sigma^2}\right) \cdot \exp\left(-\frac{X^T Z}{\sigma^2}\right)$$

是泰勒级数、Taylor 展开

$$1 + \frac{1}{1!} \cdot \left(\frac{X^T Z}{\sigma^2}\right)^1 + \dots + \frac{1}{n!} \left(\frac{X^T Z}{\sigma^2}\right)^n = \sum_{p=0}^{\infty} \frac{1}{p!} \left(\frac{X^T Z}{\sigma^2}\right)^p$$

σ : length parameter

★ Slack variables to handle outliers.



Soft-margin SVM: Previously we need $y_i(w^T x_i + b) \geq 1$

Now we introduce the slack variable $\xi_i \geq 0$, require $y_i(w^T x_i + b) \geq 1 - \xi_i$

Here we require ξ_i to be small

$$\min_{w, b, \xi} \frac{1}{2} \|w\|^2 + C \cdot \sum_{i \in [n]} \xi_i$$

$$\text{s.t. } y_i(w^T x_i + b) \geq 1 - \xi_i, \forall i \in [n]$$

$$\xi_i \geq 0, \forall i \in [n]$$

We need that $\xi_i \geq 0, \xi_i \geq 1 - y_i(w^T x_i + b)$

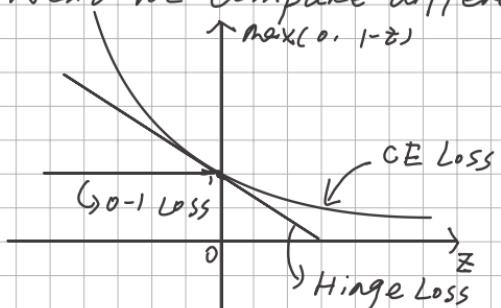
$$\Rightarrow \xi_i \geq \max(0, 1 - y_i(w^T x_i + b)) \quad ("=" \text{ when } 1 - \gamma = 0)$$

$$\text{Goal: } \min_{w, b} \frac{1}{2} \|w\|^2 + C \cdot \sum_{i \in [n]} \max(0, 1 - y_i(w^T x_i + b))$$

Hinge Loss

Here we transfer it to a constrained optimization.

Next we compare different loss:



When categorize correctly the penalty of Hinge-loss is less.

Dual of soft SVM

$$L(w, b, \xi, \alpha, \beta) = \frac{1}{2} \|w\|^2 + C \cdot \sum_{i \in [n]} \xi_i + \sum_{i \in [n]} \alpha_i [1 - \xi_i - y_i(w^T x_i + b)] - \sum_{i \in [n]} \beta_i \xi_i$$

$$\min \max L: \text{Dual } \max_{\alpha \geq 0} \min_{w, b, \xi} L \quad \beta \geq 0$$

$$\begin{cases} \frac{\partial L}{\partial w} = w - \sum_{i \in [n]} \alpha_i y_i x_i = 0 \\ \frac{\partial L}{\partial b} = -\sum_{i \in [n]} \alpha_i y_i = 0 \\ \frac{\partial L}{\partial \xi_i} = C - \alpha_i - \beta_i = 0 \end{cases}$$

we can get that. $\begin{cases} W = \sum_{i \in [n]} \alpha_i y_i x_i \\ \sum \alpha_i y_i = 0 \\ \alpha_i + \beta_i = C \end{cases}$

we can substitute W into L , and we can get

$$\max_{\alpha} \sum_{i \in [n]} \alpha_i - \frac{1}{2} \cdot \sum_{i \in [n], j \in [n]} \alpha_i \alpha_j y_i y_j x_i^T x_j$$

$$\text{s.t. } 0 \leq \alpha_i \leq C, \text{ if } C \text{ is small, robust to outlier.}$$

$$\sum_{i \in [n]} \alpha_i y_i = 0$$

if $\alpha_i = 0$, 分类正确

$0 < \alpha_i < C$, support vector

$\alpha_i = C$, 分类错误 (好像也是支持向量)

Lecture 4. Representer Theorem (表示定理)

By default, we use $\varphi(x_i)$ to replace x , $\varphi(x_i) = [x_i^T, 1]^T$
 $f(x_i) = w^T \varphi(x_i)$, can recover $w^T x + b$

1. Linear Regression

$$\text{Loss} = \min_w \sum_{i \in [n]} (w^T \varphi(x_i) - y_i)^2 + \lambda \|w\|^2$$

2. Logistic Regression

$$\min_w \sum_{i \in [n]} \log(1 + \exp(-y_i w^T \varphi(x_i)))$$

3. SVM (soft-SVM)

$$\Leftrightarrow \min_w \sum_{i \in [n]} \max(0, 1 - y_i (w^T \varphi(x_i) + b)) + \lambda \|w\|^2$$

About these Resolution:

We Introduce ξ_i , $\xi_i \geq 0$, $\xi_i \geq 1 - y_i w^T \varphi(x_i)$

$$\text{At Last we can get } \begin{cases} w = \frac{1}{\lambda} \sum_{i \in [n]} \alpha_i y_i \varphi(x_i) \\ \alpha_i + \beta_i = 1 \end{cases}$$

$$\Rightarrow f(x) = w^T \varphi(x) = \sum_{i \in [n]} \varphi(x_i) \cdot \varphi_i^T \varphi(x_i) = \sum_{i \in [n]} \alpha_i k(x_i, x)$$

For linear Regression:

$$\frac{\partial L}{\partial w} = 2 \cdot \sum_{i \in [n]} (w^T \varphi(x_i) - y_i) \cdot \varphi(x_i) + 2\lambda w = 0$$

$$\Rightarrow w = \frac{1}{\lambda} \sum_{i \in [n]} (y_i - w^T \varphi(x_i)) \varphi(x_i)$$

$$= \sum_{i \in [n]} \alpha_i f(x_i), \quad \alpha_i = \frac{1}{\lambda} (y_i - w^T \varphi(x_i)) \quad ?$$

$$\therefore f(x_i) = \sum_{i \in [n]} \alpha_i \varphi_i^T \varphi(x_i) = \sum_{i \in [n]} \alpha_i k(x_i, x)$$

Here we introduce Hilbert Space.

(1) Attribute:

$$\textcircled{1} \quad \langle f, g \rangle_H = \langle g, f \rangle_H$$

$$\textcircled{2} \quad \langle \alpha_1 f + \alpha_2 g, h \rangle = \alpha_1 \langle f, h \rangle + \alpha_2 \langle g, h \rangle$$

$$\textcircled{3} \quad \langle f, f \rangle_H \geq 0, " = " \text{ when } f = 0$$

\Rightarrow Loss is independent of regularization minimizes at $\lambda = 0$

$$\Rightarrow f(x) = \sum_{i \in [n]} \alpha_i k(x_i, x)$$

Significance:

(1) Turns a potentially infinite-dim optimization of f into a search of $\alpha_1, \alpha_2, \dots, \alpha_n$

(2) Shows that a wide range of learning algorithms have solutions expressed as weighted sum of kernel functions on finite training data. $f(\cdot) = \sum_{i \in [n]} \alpha_i k(x_i, \cdot)$

Here we show some example.

Ridge Regression

$$\min_w J(w) = \frac{1}{2} \sum_{i \in [n]} (w^T \varphi(x_i) - y_i)^2 + \frac{1}{2} \lambda w^T w$$

$$\text{Let } \frac{\partial J}{\partial w} = \sum_{i \in [n]} (w^T \varphi(x_i) - y_i) \cdot \varphi(x_i) + \lambda w = 0$$

$$\Rightarrow w = \sum_{i \in [n]} \alpha_i \varphi(x_i) \text{ (where, } \alpha_i = \frac{1}{\lambda} (y_i - w^T \varphi(x_i))$$

$$\therefore f(x) = \sum_{i \in [n]} \alpha_i (\varphi(x_i) \cdot \varphi(x)) = \sum_{i \in [n]} \alpha_i k(x_i, x) \\ = w^T \varphi(x)$$

$$\therefore w = \sum_{i \in [n]} \alpha_i \varphi(x_i)$$

$$\text{Let } \Phi = \begin{pmatrix} \varphi(x_1)^T \\ \vdots \\ \varphi(x_n)^T \end{pmatrix} \in \mathbb{R}^{n \times d} \Rightarrow w = \Phi^T \alpha \in \mathbb{R}^d.$$

$$J(w) = \frac{1}{2} (\Phi w - Y)^T (\Phi w - Y) + \frac{\lambda}{2} w^T w \quad (w = \Phi^T \alpha)$$

$$= \frac{1}{2} (w^T \Phi^T \Phi w + Y^T Y - 2 w^T \Phi^T Y) + \frac{\lambda}{2} w^T w$$

$$\text{Let } \lambda w = \Phi^T \alpha \Rightarrow J(w) = \frac{1}{2} \alpha^T \Phi \Phi^T \Phi \Phi^T \alpha + \frac{1}{2} Y^T Y - \alpha^T \Phi \Phi^T Y \\ + \frac{\lambda}{2} \alpha^T \Phi \Phi^T \alpha$$

$$\begin{aligned} & \text{Def } \Phi(x)\Phi(x)^T \\ &= \begin{pmatrix} \Phi(x_1) \\ \vdots \\ \Phi(x_n) \end{pmatrix} (\phi(x_1) \dots \phi(x_n)) = \begin{pmatrix} \Phi(x_1)^T \phi(x_1), \Phi(x_1)^T \phi(x_2) \\ \vdots \\ \ddots \\ \Phi(x_n)^T \phi(x_n) \end{pmatrix} \\ &= \begin{pmatrix} k(x_1, x_1) & k(x_1, x_2) \dots \\ \ddots & \ddots & k(x_n, x_n) \end{pmatrix} = K \text{ Gram Matrix} \end{aligned}$$

代入 $J(w)$ 有

$$J = \frac{1}{2} \alpha^T K \cdot K \alpha + \frac{1}{2} y^T y - \alpha^T Ky + \frac{\lambda}{2} \alpha^T K \alpha = J(\alpha),$$

$$\frac{\partial J}{\partial \alpha} = 0 \Rightarrow K \cdot K \alpha - K y + \lambda K \alpha = 0 \quad ①$$

$$K^{-1} ① \Rightarrow K \alpha - Y + \lambda I \alpha = 0$$

$$\therefore \alpha = C K + \lambda I^{-1} Y$$

$$\begin{aligned} \text{Def } f(x) &= \sum_{i \in [n]} \alpha_i k(x_i, x) = k(x, \alpha) \quad (k(x) = \begin{pmatrix} k(x_1, x) \\ k(x_2, x) \\ \vdots \\ k(x_n, x) \end{pmatrix}) \\ &= \underbrace{k(x)}^T (C K + \lambda I^{-1} Y) \end{aligned}$$

思考題

$$\begin{aligned} w &= C x^T x + \lambda I^{-1} x^T y \quad ① \\ &= (\Phi^T x, \Phi(x) + \lambda I)^{-1} \Phi^T y \quad ② \end{aligned}$$

$$\text{Def } f(x) = w^T \phi(x) = \phi^T x, (C \Phi^T \Phi + \lambda I) \phi^T y$$

与 $k(x, C K + \lambda I)^{-1} Y$ 是否等价.

Anthony
Zadern

Lecture 5. Learning Theory

Lecture 6. Tree Model And Ensemble Learning

Previous Model: $f(x) = w^T x + b$, $y = \text{sign}(w^T x + b)$

⇒ Decision logic $\begin{cases} ① w^T x + b \geq 0, +1 \\ ② w^T x + b < 0, -1 \end{cases}$

Deficiencies: Unable to do classification in unlinear space.

Here we introduce the Decision Trees,
⇒ Definition: A decision Tree is a tree that contains a Root, internal nodes, connected by directed edges. Each non-leaf node partitions the data by some features.

Then the question is that:

* How to select a good feature to partition?

The Answer is: Purity!

Here we use the information theory:

$$\begin{aligned} \text{(1)} H(x_1) &= -\sum p(x_1) \cdot \log_2 p(x_1) = \sum p(x_1) \cdot \log_2 \frac{1}{p(x_1)} \\ &= E(\log_2 \frac{1}{p(x_1)}) \end{aligned}$$

$$H(x) \leq \log_2 n \quad (\text{Jensen's Inequality})$$

(2) Cross Entropy:

$$H(p, q) = \sum p(x_i) \cdot \log_2 q(x_i).$$

where $p(x_i)$ measures the observed data.

$q(x_i)$ measures the probability of y_i :

$$E(x_i, y_i) = - \sum_{y_i \in \{0, 1\}} p(y_i) \cdot \log_2 q(y_i)$$

In the next part, we will learn to how to measure the purity of branch.

How to measure purity?

1. Information Gain.

$$g_{CD, A} = H(CD) - H(D|A) \leftarrow \text{select } A \text{ to maximize.}$$

D: training set; A: a feature/Attribute

suppose $A \in \{a_1, \dots, a_m\}$ m discrete values

$$y \in \{1, 2, \dots, k\} \quad H(y)$$

$$H(CD) = - \sum_{k \in [K]} \frac{|C_k|}{|D|} \cdot \log \underbrace{\frac{|C_k|}{|D|}}_{\text{approximate } P(C_k)}$$

$$\begin{aligned} H(D|A) &= \sum_{i \in [m]} \frac{|D_i|}{|D|} \cdot H(CD | A=a_i) \\ &= \sum_{i \in [m]} \frac{|D_i|}{|D|} \left(- \sum_{k \in [K]} \frac{|D_i \cap C_k|}{|D_i|} \cdot \log \frac{|D_i \cap C_k|}{|D_i|} \right) \quad H(y | A=a_i) \end{aligned}$$

So we just need the A to maximize $H(CD) - H(D|A)$

2. Information Gain Ratio

$$g_{R(CD, A)} = \frac{g_{CD, A}}{H_{ACD}}, \text{ where } H_{ACD} = - \sum_{i \in [m]} \frac{|D_i|}{|D|} \cdot \log \frac{|D_i|}{|D|}$$

when each $A=a_i$ has equal prob., $H_{ACD} = \log_2 m$

Penalize A with a large m.

3. Gini Index

$$\begin{aligned} \text{Gini}(CD) &= \sum_{i \in [K]} \frac{|C_k|}{|D|} \left(1 - \frac{|C_k|}{|D|} \right), \text{ approximates } \sum_{i \in [K]} P(y=k)(1-P(y=k)) \\ &= 1 - \sum_{i \in [K]} \left(\frac{|C_k|}{|D|} \right)^2 \quad \begin{cases} \max = 1 - \left(\frac{1}{K}\right)^2, K=1 - \frac{1}{K} \\ \min = 0 \leftarrow \text{when purity.} \end{cases} \end{aligned}$$

4. L2 Loss

For Regression, just use L2 Loss to measure purity

$$\bar{y}_{D_i} = \frac{1}{|D_i|} \sum_{j \in D_i} y_j;$$

$$L(CD, A) = \sum_{i \in [n]} \left[\sum_{j \in D_i} (y_j - \bar{y}_{D_i})^2 \right]$$

$$\operatorname{argmin} L(CD, A)$$

Some bagging Methods.

1. Bootstrap

- ① Sample n points from D w/ replacement
- ② Repeat ① for T times, get $\hat{D}_1, \hat{D}_2, \dots, \hat{D}_T$
- ③ $\hat{f}(x_i) = \frac{1}{T} \sum_{t \in [T]} f(x_i, \hat{D}_t)$ as the bagged model

※ Bootstrap has no theoretical guarantee! Empirically can reduce the test error greatly.



Random Forest (the most successful bagging model)

Algorithm:

→ for t in $[T]$:

(1) Sample n points \hat{D}_t from D w/ replacement.

(2) Sample $d' < d$ features F' from F without replacement.

(3) Build a full decision tree on \hat{D}_t, F' (can do some pruning to minimize out-of-bag error).

→ End for

Then we can average all trees.

(1) If Regression, we use the mean as result.

(2) If categorization, we use the mode as result.

Boosting (Reduce Bias)

Application: Adapt to low variance high bias nodes.

AdaBoost: $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$, $y \in \{-1, 1\}$, $x \in \mathbb{R}^d$

Input: D , a weak learning algorithm A .

Algorithm:

- (1) Initialize sample weights, $w_i^{(0)} = \bar{w}_i^{(0)} = \frac{1}{n}, \forall i \in [n]$
- (2) For t in $[T]$:
 - a) use D and $\{w_i^{(t)}\}$ to train A . get a model $f_t(x) : \mathbb{R}^d \rightarrow \{-1, 1\}$
 - b) Evaluate $f_t(x)$ weighted classification error on D
 $e_t = \sum_{i \in [n]} w_i^{(t)} \cdot \mathbb{1}(f_t(x_i) \neq y_i)$
 - c) Compute a weight α_t for $f_t(x)$. $\alpha_t = \frac{1}{2} \log \frac{1 - e_t}{e_t}$
 - d) Update sample weights: $\bar{w}_i^{(t+1)} = \bar{w}_i^{(t)} \cdot \exp(-\alpha_t \cdot y_i f_t(x_i))$
Normalization: $W_i^{(t+1)} = \bar{w}_i^{(t+1)} / \sum_{j \in [n]} \bar{w}_j^{(t+1)}$
- (3) Combine T classifiers linearly.
 $g(x) = \text{sign} \left(\sum_{t \in [T]} \alpha_t \cdot f_t(x) \right)$

Now let's focus on $\alpha_t = \frac{1}{2} \log \frac{1 - e_t}{e_t}$

First, ① $e_t \uparrow$, $\alpha_t \uparrow$

② when $e_t > 0.5$, $\alpha_t < 0$

$e_t < 0.5$, $\alpha_t > 0$

③ let's focus on $\bar{w}_i^{(t+1)} = \bar{w}_i^{(t)} \cdot \exp(-\alpha_t \cdot y_i \cdot f_t(x_i))$
 $y_i \neq \text{sign} \left(\sum_{t \in [T]} \alpha_t \cdot f_t(x_i) \right)$, get larger weights.

④ Because $W_i^{(t+1)} = \bar{w}_i^{(t+1)} / \sum_{j=1}^n \bar{w}_j^{(t+1)} \Rightarrow \sum_{i \in [n]} W_i^{(t+1)} = 1$

Next we will discuss about More common perspective of boosting, which is called Additive Model.

Additive Model (A general paradigm for boosting)

$$g_t(x) = \sum_{i \in [t]} \alpha_i \cdot f_i(x). \text{ Define } g_t(x) = \sum_{j \in [t]} \alpha_j \cdot f_j(x)$$

At step t , keep $g_{t-1}(x)$ fixed. Then Learn $\alpha_t, f_t(x)$ by minimizing $\sum_{i \in [n]} L(y_i, g_{t-1}(x_i) + \alpha_t \cdot f_t(x_i))$

$$\text{Finally, } g_t(x) = g_{t-1}(x)$$

AdaBoost is an additive model using exponentially loss, $L(y, f(x)) = \exp(-y f(x))$. $y f(x)$ only takes $1/-1$

So, $L(y, f(x))$ only takes $e/ \frac{1}{e}$

AdaBoost Derivation

At step t , we already have $g_{t-1} = \sum_{j \in [t-1]} \alpha_j f_j(x)$

We need to optimize: $\min_{\alpha_t, f_t} \sum_{i \in [n]} \exp(-y_i (g_{t-1}(x_i) + \alpha_t f_t(x_i)))$

Here we define $\exp(-y_i g_{t-1}(x_i))$ to be $\bar{w}_i^{(t)}$

$$\bar{w}_i^{(t)} = \prod_{j \in [t-1]} \exp(-y_i \cdot \alpha_j f_j(x_i))$$

$$= \bar{w}_i^{(t-1)} \cdot \exp(-y_i \cdot \alpha_t f_t(x_i))$$

$$\therefore \text{Goal} \Leftrightarrow \min_{\alpha_t} \min_{f_t} \sum_{i \in [n]} \bar{w}_i^{(t)} \cdot \exp(-y_i \cdot \alpha_t f_t(x_i))$$

Fix α_t , optimize f_t . In AdaBoost, we just train f_t using its original loss and simple weights $w_i^{(t)}$

After Solving f_t , we optimize

$$\min_{\alpha_t} \sum_{i \in [n]} w_i^{(t)} \exp(-y_i \alpha_t f_t(x_i))$$

Next we try to solve this problem and get the related results.

$$\begin{aligned}
 & \min_{\alpha_t} \sum_{i \in [n]} w_i^{(t)} \exp(-\alpha_t \cdot y_i \cdot f_t(x_i)) \\
 &= \min_{\alpha_t} \sum_{\substack{i \\ y_i = f_t(x_i)}} w_i^{(t)} \exp(-\alpha_t \cdot y_i \cdot f_t(x_i)) + \sum_{\substack{i \\ y_i \neq f_t(x_i)}} w_i^{(t)} \exp(-\alpha_t \cdot y_i \cdot \underbrace{f_t(x_i)}_{-1}) \\
 &= \min_{\alpha_t} \left(\sum_{\substack{i \\ y_i = f_t(x_i)}} w_i^{(t)} \exp(-\alpha_t) + \sum_{\substack{i \\ y_i \neq f_t(x_i)}} w_i^{(t)} \exp(-\alpha_t) \right) \\
 &\quad + \left(\sum_{\substack{i \\ y_i \neq f_t(x_i)}} w_i^{(t)} \exp(\alpha_t) - \sum_{\substack{i \\ y_i \neq f_t(x_i)}} w_i^{(t)} \exp(-\alpha_t) \right) \\
 &= \min_{\alpha_t} \left(\underbrace{\left(\sum_{i \in [n]} w_i^{(t)} \right)}_L \exp(-\alpha_t) \right) + \left[C \sum_{\substack{i \\ y_i \neq f_t(x_i)}} w_i^{(t)} \right] (\exp(\alpha_t) - \exp(-\alpha_t)) \\
 &\Leftrightarrow \min_{\alpha_t} \underbrace{\exp(-\alpha_t)}_{\beta} + e_t [\exp(\alpha_t) - \exp(-\alpha_t)] \\
 &\frac{\partial L}{\partial \alpha_t} = -\exp(-\alpha_t) + e_t [\exp(\alpha_t) + \underbrace{\exp(-\alpha_t)}_{\beta}] = 0 \\
 &\therefore \left(\beta + \frac{1}{\beta} \right) e_t = \beta \\
 &\Rightarrow (1 - e_t) \beta^2 = e_t \Rightarrow \beta = \sqrt{\frac{e_t}{1 - e_t}} = e^{-\alpha_t} \\
 &\therefore \alpha_t = \frac{1}{2} \log \frac{1 - e_t}{e_t}
 \end{aligned}$$

Lecture 7 Gaussian Process

Previously, we treat w as fixed (constant) parameters.
In Bayesian view, even the parameters w are R.

Prior Prob. $P(w) = N(w | 0, 6\omega^2 I)$

$$P(y|x, w; \sigma^2, \sigma_w^2) = N(y|w^\top \phi(x), \sigma^2)$$

$$P(w|y, x) \propto \frac{P(y|x, w) \cdot P(w|x)}{P(y|x)}$$

\downarrow
 $\prod_{i \in [n]} P(y_i|x_i, w)$

$$= \frac{1}{Z} \cdot \left[\left(\frac{1}{\sqrt{2\pi\sigma^2}} \right)^n \exp\left(-\frac{\sum_{i \in [n]} (y_i - w^\top \phi(x_i))^2}{2\sigma^2}\right) \cdot \frac{1}{(\sqrt{2\pi\sigma_w^2})^d} \exp\left(-\frac{w^\top w}{2\sigma_w^2}\right) \right]$$

Hence it selects the mode of the posterior distribution as a point estimation.

$$\text{Goal: } \max_c (\log P(w|y, x)) = \log(Z)$$

$$\Leftrightarrow \max_w - \frac{1}{2\sigma^2} \sum_{i \in [n]} (y_i - w^\top \phi(x_i))^2 - \frac{1}{2\sigma_w^2} w^\top w$$

$$\Leftrightarrow \min_w \sum_{i \in [n]} (y_i - w^\top \phi(x_i))^2 + \underbrace{\frac{\sigma^2}{\sigma_w^2} \|w\|^2}_{\lambda}$$

Stochastic Process

(1) Definition: A collection of (infinite/many) R.V.s along an index set.

Gaussian Process

(2) Definition: specify the joint distribution over any finite collection of R.V.s and require the joint distribution to be Gaussian.

$\{x_1, \dots, x_n\}$ is any set of n points in the index set.
Then $G.P. \Leftrightarrow y_1, \dots, y_n$ have a multivariate Gaussian

Distribution,

Multivariate Gaussian: $X \in \mathbb{R}^d$

$$N(X | \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{d}{2}}} \cdot \frac{1}{|\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{(x-\mu)^T \Sigma^{-1} (x-\mu)}{2}\right)$$

parameters: $d + \frac{d(d+1)}{2}$

We can get that $\{E[x_i] = \mu_i\}$

$$\text{Cov}(X) = E[(X - \mu)(X - \mu)^T] = \Sigma$$

$$\Sigma_{ij} = E[(x_i - \mu_i)(x_j - \mu_j)] = \text{cov}(x_i, x_j)$$

Lamma

.. Any linear distribution of Gaussian Distribution R.V.s follows a Gaussian. Concatentation of Gaussian R.V.s results in a multivariate Gaussian distributed R.V.

$$a \sim N(\mu_a, \sigma_a^2), \begin{bmatrix} a \\ b \end{bmatrix} \sim N(\mu, \Sigma)$$
$$b \sim N(\mu_b, \sigma_b^2)$$

Question: How to capture the joint distribution,

need to only specify $\mu_i = \mu(x_i)$, $\Sigma_{ij} = k(x_i, x_j)$

Suppose n training points $\{x_1, \dots, x_n\}$. let k be Gram matrix. $k_{ij} = k(x_i, x_j) \rightarrow$ valid kernel if $k \geq 0$ for any $\{x_1, \dots, x_n\}$

$$GP(m(\cdot), k(\cdot, \cdot)) . P(y) = N(y | 0, K)$$

Now given a new test point x_* , we can compute joint distribution of $\begin{bmatrix} y_* \\ y \end{bmatrix} \in \mathbb{R}^{n+1}$

$$P\left(\begin{bmatrix} y_* \\ y \end{bmatrix}\right) = N\left(\begin{bmatrix} y_* \\ y \end{bmatrix} \mid 0, \begin{pmatrix} k(x_*, x_*) & k(x_*, x) \\ k(x, x_*) & k(x, x) \end{pmatrix}\right)$$

$P(y^* | y) = N(y^* | \mu^*, \Sigma^*)$, \Leftarrow How to compute this thing is important.

$$\text{Given } \begin{cases} X_a \sim N(M_a, \Sigma_{aa}) \\ X_b \sim N(M_b, \Sigma_{bb}) \end{cases} \Rightarrow \begin{cases} M_{ab} = M_a + \Sigma_{ab} \cdot \Sigma_{bb}^{-1} (X_b - M_b) \\ \Sigma_{ab} = \Sigma_{aa} - \Sigma_{ab} \cdot \Sigma_{bb}^{-1} \Sigma_{ba} \end{cases}$$

$$\therefore m^* = k(x^*, x)^T K^{-1} y$$

$$\Sigma^* = k(x^*, x^*) - k(x^*)^T K^{-1} k(x^*)$$

In a more realistic setting, we can only observe

$$\hat{y} = y + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2 I)$$

$$\Rightarrow \begin{cases} m^* = k(x^*)^T (K + \lambda I)^{-1} y \\ \Sigma^* = k(x^*, x^*) - k(x^*)^T (K + \lambda I)^{-1} k(x^*, x) \end{cases}$$

For MAP,

$$\text{Goal: } \min_w \sum_{i \in [n]} (y_i - w^T \phi(x_i))^2 + \lambda w^T w$$

Bayesian View MAP Function:

$$y_i = w^T \phi(x_i) + \varepsilon_i, \quad \varepsilon_i \sim N(0, \sigma^2), \quad w \sim N(0, \sigma_w^2 I)$$

$$Y = w^T \phi + \varepsilon \sim N(m, \Sigma)$$

$$\text{cov}(y_i, y_j) = E(y_i y_j) - E(y_i) E(y_j)$$

$$= E((w^T \phi(x_i) + \varepsilon_i)(w^T \phi(x_j) + \varepsilon_j))$$

$$= E(\phi(x_i)^T w \cdot w^T \phi(x_j)) + E(\varepsilon_i \varepsilon_j)$$

$$= \phi(x_i)^T \underbrace{E(w w^T)}_{\sigma_w^2 I} \phi(x_j) + \sigma^2 I$$

$$\Rightarrow \sigma_w^2 \phi(x_i)^T \phi(x_j) + \sigma^2 I$$

$$\Sigma = \sigma_w^2 K + \sigma^2 I$$

$$\text{固理想 } M_* = \sigma_w^2 k(x^*)^T [\sigma_w^2 K + \sigma^2 I]^{-1} y$$

$$= k(x^*)^T \left[K + \frac{\sigma^2}{\sigma_w^2} I \right]^{-1} y$$

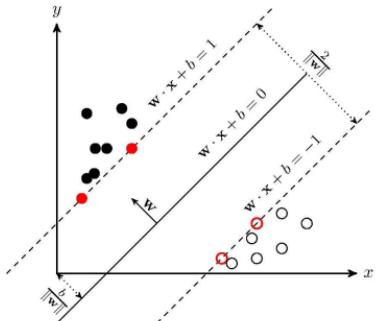
let $b = \lambda$

~~Bayesian Optimization (BO)~~

1) For black-box functions

- Steps:
- ① Randomly sample n points X_1, \dots, X_n
 - ② fit a GP → lower confidence bound
 - ③ Use some aquisition function to select next point to evaluate
 - ④ ~~迭代上一步过程~~: use all X, Y to fit a new GP
 - ⑤ Repeat until reaching a budget.

SVM 数学原理



Train Data: $T = \{(x_1, y_1), \dots, (x_n, y_n)\}$
 $x_i \in \mathbb{R}^n, y_i \in \{1, -1\}$

For any datapoint:

$$\text{the geometric distance } r_i = y_i \cdot \frac{\|w\|}{\|w\|} \cdot x_i + b$$

Let $r = \min r_i$

(So our Goal is that:

$$\max_{w, b} r \text{ s.t. } y_i \left(\frac{w}{\|w\|} \cdot x_i + \frac{b}{\|w\|} \right) \geq r$$

$$\Leftrightarrow y_i \left(\frac{w}{\|w\|} \cdot r \cdot x_i + \frac{b}{\|w\|} \cdot r \right) \geq 1$$

$$\Leftrightarrow y_i (w x_i + b) \geq 1. (w = \frac{w}{\|w\|} \cdot r)$$

In the picture above, $\max r \Leftrightarrow \max \frac{1}{\|w\|} \Leftrightarrow \min \frac{1}{2} \|w\|^2$

\Rightarrow Goal Function: $\min_{w, b} \frac{1}{2} \|w\|^2 \text{ s.t. } y_i (w x_i + b) \geq 1$

(1) How to solve this problem?

$$\text{拉格朗日乘子: } L(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^n \alpha_i (y_i (w x_i + b) - 1) \quad (\alpha_i \geq 0) \quad ①$$

$$\text{let } \theta(w) = \max_{\alpha_i \geq 0} L(w, b, \alpha)$$

Here're different conditions:

$$① y_i (w x_i + b) - 1 \leq 0 \Rightarrow \alpha_i \rightarrow +\infty, \theta(w) \rightarrow +\infty$$

$$② y_i (w x_i + b) - 1 \geq 0, > 0, \alpha_i = 0, \forall i \text{ 任意.}$$

$$\text{上述问题} \Leftrightarrow \min_{w, b} \theta(w) = \min_{w, b} \max_{\alpha_i \geq 0} L(w, b, \alpha) = P^*$$

(dual form: change min and max)

$$\max_{\alpha_i \geq 0} \min_{w, b} L(w, b, \alpha) = d^*$$

$$P^* = d^* \text{ when } \begin{cases} ① \text{ 凸优化} \\ ② \text{ KKT 条件} \end{cases} \begin{array}{l} \rightarrow \alpha_i \geq 0 \\ \rightarrow y_i (w x_i + b) - 1 \geq 0 \\ \rightarrow \alpha_i (y_i (w x_i + b) - 1) = 0 \end{array}$$

$$\text{而 Let } \frac{\partial \theta}{\partial w} = 0 \Rightarrow W = \sum_{i=1}^N \alpha_i y_i x_i, \sum_{i=1}^N \alpha_i y_i = 0$$

$$\text{代入} ① \text{ 有 } L = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j x_i \cdot x_j + \sum_{i=1}^N \alpha_i$$

$$\text{即 } \min_{w, b} L(w, b, \alpha) = -\frac{1}{2} \sum_{i=1, j=1}^N \alpha_i \alpha_j y_i y_j (x_i \cdot x_j) + \sum_{i=1}^N \alpha_i$$

$\min_{w,b} L(w,b,\alpha)$ 对 α 越大，即是一个对偶问题

$$\max_{\alpha} -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j (x_i \cdot x_j) + \sum_{i=1}^N \alpha_i$$

$$\text{s.t. } \sum_{i=1}^n \alpha_i y_i = 0, \alpha_i \geq 0.$$

(使用SMO解得 α^*)

$$w^* = \sum_{i=1}^n \alpha_i^* y_i x_i$$

$$b^* = y_j - \sum_{i=1}^n \alpha_i^* y_i (x_i \cdot x_j)$$

(2) For non-linear separable, we use soft-margin.

$$\min_{w,b,\beta} \frac{1}{2} \|w\|^2 + C \cdot \sum_{i=1}^m \xi_i$$

$$\text{s.t. } \begin{cases} y_i (w \cdot x_i + b) \geq 1 - \xi_i \\ \xi_i \geq 0 \end{cases} \text{ Actually, } \xi_i = \max(0, 1 - y_i (w \cdot x_i + b))$$

$$\text{此时有 } \begin{cases} 0 < \alpha_i < C, y_i (w \cdot x_i + b) = 1 \\ \alpha_i = 0, y_i (w \cdot x_i + b) > 1 \\ \alpha_i = C, y_i (w \cdot x_i + b) < 1 \end{cases}$$

(3) For Kernel Trick

$$f(x) = \text{Sign}(\sum_{i=1}^N \alpha_i^* y_i k(x_i, x) + b^*)$$

$$b^* = y_j - \sum_{i=1}^N \alpha_i^* y_i k(x_i, x_j)$$

How to get the $\alpha^* \Rightarrow$ SMO Algorithm.

$$W = \min_{\alpha} \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j k(x_i, x_j) - \sum_{i=1}^N \alpha_i$$

$$\text{s.t. } \sum_{i=1}^n \alpha_i y_i = 0, 0 \leq \alpha_i \leq C$$

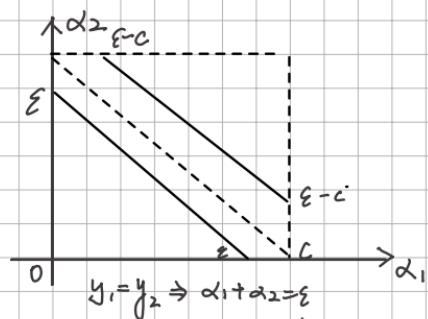
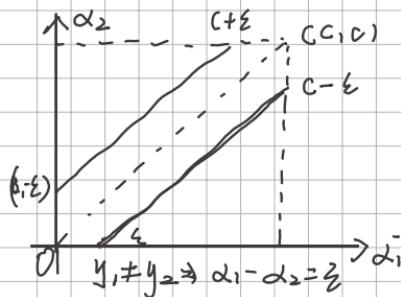
基本思想：每次只优化两个 α . 假设设其为 α_1, α_2

$$\therefore W = \frac{1}{2} k_{11} \alpha_1^2 + \frac{1}{2} k_{22} \alpha_2^2 + y_1 y_2 k_{12} \alpha_1 \alpha_2 - (\alpha_1 + \alpha_2)$$

$$+ y_1 \alpha_1 \sum_{i=3}^N y_i \alpha_i k_{1i} + y_2 \alpha_2 \sum_{i=3}^N y_i \alpha_i k_{2i} + t$$

$$\text{s.t. } \alpha_1 y_1 + \alpha_2 y_2 = - \sum_{i=3}^N \alpha_i y_i = \xi$$

这是一个两个变量二次规划问题.



最优解应该在这些直线上. Assume $\alpha_1^{\text{old}}, \alpha_2^{\text{old}}$, $\alpha_1^{\text{new}}, \alpha_2^{\text{new}}$

有 $L \leq \alpha_2^{\text{new}} \leq H$

$$\text{而 } \begin{cases} L = \max(0, \alpha_2^{\text{old}} - \alpha_1^{\text{old}}), H = \min(C, C + \alpha_2^{\text{old}} - \alpha_1^{\text{old}}), \\ y_1 \neq y_2 \Rightarrow \alpha_1 + \alpha_2 = \epsilon \\ L = \max(0, \alpha_1^{\text{old}} + \alpha_2^{\text{old}} - C), H = \min(C, C + \alpha_1^{\text{old}} + \alpha_2^{\text{old}} - C) \end{cases}$$

这里这样写是因为就是 $\alpha_i + \alpha_j > 0$. 其可以大于 C , 也可以小于 C
 $\alpha_i - \alpha_j$ 可以大于 0 也可以小于 0 .

Define Prediction Error:

$$E_i = \left(\sum_{j=1}^N y_j y_i K(x_j \cdot x_i) + b \right) - y_i, i=1,2$$

$$\eta = K_{11} + K_{22} - 2K_{12}$$

$$\text{剪枝解: } \alpha_2^{\text{new,unc}} = \alpha_2^{\text{old}} + \frac{y_2(E_1 - E_2)}{\eta}$$

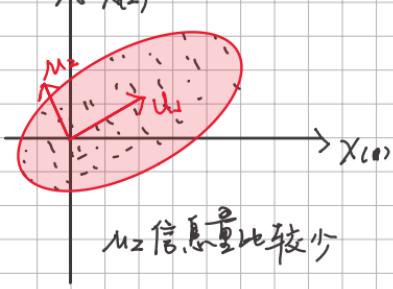
$$\text{考虑边界有: } \alpha_2^{\text{new}} = \begin{cases} H, & \text{if } \alpha_2^{\text{new,unc}} > H \\ \alpha_2^{\text{new,unc}}, & \text{if } L \leq \alpha_2^{\text{new,unc}} \leq H \\ L, & \text{if } \alpha_2^{\text{new,unc}} < L \end{cases}$$

$$\text{进一步有 } \alpha_i^{\text{new}} = \alpha_i^{\text{old}} + y_i y_2 (\alpha_2^{\text{old}} - \alpha_2^{\text{new}})$$

#Lecture 8 Unsupervised Learning: X only, no y.

- Tasks
- ① Dimension Reduction: 100^{100} dim \rightarrow 100 dim. $x_i \rightarrow y$
 - ② Clustering (聚类), group similar data together.
 - ③ Generative Models: Sample new data from $P(x)$

① Dimensionality Reduction \rightarrow Principal Component Analysis (PCA)



$x^{(1)}, x^{(2)}$ highly correlated.

Principals: find a new basis u_1, u_2 s.t. $u_1 \perp u_2$, u_1 keeps the most information (data variance on u_1 is the largest). Then throw away u_2 . Represent a simple x by u_1 .

Here are algorithm details:

$$X = \begin{pmatrix} x_1^T \\ \vdots \\ x_n^T \end{pmatrix} \in \mathbb{R}^{n \times d}, \text{ calculate } \bar{x} = \frac{1}{n} \sum_{i \in [n]} x_i;$$

$$\text{Data covariance matrix: } \frac{1}{n} \sum_{i \in [n]} (x_i - \bar{x}) \cdot (x_i - \bar{x})^T$$

Now compute $U \in \mathbb{R}^{d \times d}$, need project each x_i to u_1 , and maximize variance.

$$\begin{aligned} \text{Loss Function: } & \frac{1}{n} \sum_{i \in [n]} \underbrace{(u_1^T x_i - u_1^T \bar{x})^2}_{\text{project}} = \frac{1}{n} \sum_{i \in [n]} u_1^T (x_i - \bar{x})(x_i - \bar{x})^T u_1 \\ &= u_1^T \left(\frac{1}{n} \sum_{i \in [n]} (x_i - \bar{x})(x_i - \bar{x})^T \right) u_1 \\ &= \underbrace{u_1^T \Sigma}_{\text{max}} u_1 \end{aligned}$$

Goal: $\max u_1^T \Sigma u_1$, s.t. $u_1^T u_1 = 1$ if not 1, can \uparrow

带约束的优化问题 \Rightarrow 子格各项目数乘

$$L(u_1, \lambda) = u_1^T \Sigma u_1 + \lambda (1 - u_1^T u_1)$$

$$\frac{\partial L}{\partial u_1} = \Sigma u_1 + \Sigma^T u_1 - 2\lambda u_1 = 0 \Rightarrow \underline{\Sigma u_1 = \lambda u_1}$$

$\Rightarrow u_1$ is an eigen vector of Σ corresponding to eigenvalue λ
 while, $u_1^T \Sigma u_1 = \lambda u_1^T u_1 = \lambda$, \Rightarrow Find maximum eigenvalue

Next, find u_2 the next direction that preserves the most variance

$$\max_{u_2} u_2^T \Sigma u_2, \text{ s.t. } u_2^T u_2 = 1, u_2^T u_1 = 0$$

$$L(u_2, \lambda, \alpha) = u_2^T \Sigma u_2 + \lambda (1 - u_2^T u_2) + \alpha u_2^T u_1$$

$$\Rightarrow \frac{\partial L}{\partial u_2} = 2 \Sigma u_2 - 2\lambda u_2 + \alpha u_1 \stackrel{(1)}{=} 0$$

$$\therefore 2 \underline{u_2^T \Sigma u_2} - 2\lambda_2 \underline{u_2^T u_2} + \alpha \underline{u_1^T u_1} = 0 \Rightarrow \alpha = 0$$

$$u_2^T \Sigma u_1 = \lambda u_2^T u_1 = 0$$

$$\text{if } \alpha \neq 0 \Rightarrow \Sigma u_2 = \lambda_2 u_2$$

$$\frac{\alpha L}{\alpha \lambda} = 0 \Rightarrow u_2^T \Sigma u_2 = \lambda u_2^T u_2 = \lambda_2$$

By Induction, when need k principal components,
 just compute k largest eigenvalues and eigen vectors
 of Σ .

$[u_1, u_2, \dots, u_k]$ project X by $X U_{1:k} \in \mathbb{R}^{n \times k}$, $k \ll d$.

$$U_{1:k} \in \mathbb{R}^{d \times k}$$

$$\Sigma = U \Lambda U^T \wedge \text{diagonal}, U \text{ orthogonal}$$

$$= [u_1, \dots, u_n] \cdot \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix} [u_1, u_2, \dots, u_n]^T$$

$$\hat{X} = \begin{pmatrix} x_1^T - \bar{x}^T \\ \vdots \\ x_n^T - \bar{x}^T \end{pmatrix} \in \mathbb{R}^{n \times d} \text{ then } \Sigma = \frac{1}{n} \hat{X}^T \hat{X} \text{ we can directly do eigen decomposition for } \hat{X}^T \hat{X}$$

Alternately, we can perform Singular

Value Decomposition (CSV) 奇异值分解, 适用于矩阵

for $\hat{X}^T = U \hat{\Sigma} V^T$, U, V orthogonal

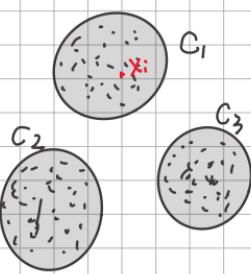
$$\hat{\Sigma} = \underbrace{\begin{bmatrix} \sigma_1 & & 0 \\ & \ddots & \\ 0 & & \sigma_d \end{bmatrix}}_V d$$

$$\hat{X}^T \cdot \hat{X} = U \hat{\Sigma} \underbrace{V^T V}_{\hat{\Sigma}^T} U^T = U (\hat{\Sigma} \hat{\Sigma}^T) U^T \xrightarrow{\text{由 } \hat{\Sigma}^T = \begin{bmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_n^2 \end{bmatrix}} = \begin{bmatrix} \lambda_1 & & \\ & \ddots & \\ & & \lambda_n \end{bmatrix}$$

(2) Clustering. Given $D = \{x_1, \dots, x_n\} \in \mathbb{R}^{d \times n}$

We want to partition D into k clusters.

$C_1 \quad \text{① } u_k \in \mathbb{R}^d$, the center of k



② Let $r_{ik} \in \{0, 1\}$ denote that x_i is closest to u_k . $\sum_{k \in [n]} r_{ik} = 1 \quad \forall i$.

K-means: Find $\{u_k | k \in [k]\}, \{r_{ik}\}$ s.t. the sum of squared distances of each x_i to its assigned cluster center u_k .

Objective: $\min_{r, u} \sum_{i \in [n]} \sum_{k \in [k]} r_{ik} \|x_i - u_k\|^2$

Idea: Alternately optimize u and r

① Given r fixed, u can be optimized very easily.

$$\frac{\partial L}{\partial u_k} = - \sum_{i \in [n]} 2r_{ik}(x_i - u_k) = 0 \Rightarrow u_k = \frac{\sum_{i \in [n]} r_{ik} \cdot x_i}{\sum_{i \in [n]} r_{ik}}$$

② Given u fixed, $r_{ik} = \begin{cases} 1, & \text{if } k = \arg \min_{j \in [k]} \|x_i - u_j\|^2 \\ 0, & \text{otherwise} \end{cases}$

K-means: Randomly initialize u_k and repeat ① and ② until convergence.

why convergence? Both ①, ② reduce obj.

K-means only focus on local minimum, not global.

In practice, we often run k-means for multiple times, each of different initialization.

★ Mixture of Gaussians (MoG)

Models $P(x)$ \Rightarrow Assumes points belonging to cluster k follow a Gaussian Distribution. $N(\mu_k, \Sigma_k)$

A generative model. $P(x) = \sum_z P(x|z) \cdot P(z)$

For R.V. X , let Z_k be another R.V. $\in \{0, 1\}$, $Z_{ik} = 1$ means X generated from cluster k $\sum_{k \in [n]} Z_{ik} = 1$.

$$\{ P(x|Z_{ik}=1) = N(x|\mu_k, \Sigma_k) \}$$

$$P(Z_{ik}=1) = \pi_k, 0 \leq \pi_k \leq 1, \sum_{k \in [n]} \pi_k = 1$$

MoG STEPS:

① Randomly initialize $r_{ik} \in [0, 1]$, soft assignment of x_i to k . $\sum_{k \in [n]} r_{ik} = 1$

② Compute weighted sampled mean of covariance,
 $\mu_k = \frac{\sum r_{ik} \cdot x_i}{\sum_{i \in [n]} r_{ik}}, \Sigma_k = \frac{\sum_{i \in [n]} r_{ik} (x_i - \mu_k)(x_i - \mu_k)^T}{\sum_{i \in [n]} r_{ik}}$

Update π_k by $\pi_k = \frac{\sum_{i \in [n]} r_{ik}}{n}$.

③ Update r_{ik} , by $r_{ik} = \frac{\pi_k N(x_i | \mu_k, \Sigma_k)}{\sum_{j \in [K]} \pi_j N(x_i | \mu_j, \Sigma_j)}$

Iterate ② and ③ until convergence.

Lecture 9. Latent variable Models

Expectation-Maximization In general.

Suppose θ contains all parameters to estimate.

We directly optimize $P(x, \theta)$ is difficult but
optimizing $P(x, z; \theta)$ is easy.

$$\text{---} \rightarrow P(x, z) = P(x|z) \cdot P(z)$$

$$P(z|x) = \frac{P(x|z) \cdot P(z)}{\sum_z P(x|z) \cdot P(z)} \rightarrow \text{Intractable.}$$

How to solve this problem?

Find a variational distributional $q(z|x)$ to
approximate $P(z|x)$!

Now we introduce $q(z|x)$ to approximate $P(z|x; \theta)$.
We have the following always holds.

$$\begin{aligned} \log P(x; \theta) &= \sum_z q_z(z|x) \cdot \log P(x; \theta) \\ &= \sum_z q_z(z|x) [\log P(x, z; \theta) - \log P(z|x; \theta)] \\ &= \sum_z q_z(z|x) \left([\log P(x, z; \theta) - \log q_z(z|x)] - \right. \\ &\quad \left. [\log P(z|x; \theta) - \log q_z(z|x)] \right) \\ &= \sum_z q_z(z|x) \cdot \log \frac{P(x, z; \theta)}{q_z(z|x)} - \sum_z q_z(z|x) \cdot \log \frac{P(z|x; \theta)}{q_z(z|x)} \\ &\text{~~~~~ Evidence lower Bound } \text{~~~~~ KL divergence.} \end{aligned}$$

We can use Jensen's Ineq to prove KL divergence ≥ 0 .

EM is maximizing ELBO! Maximizing the lower bound
also \uparrow evidence. Next page we will introduce the
EM Algorithm.

EM Algorithm:

A two step iterative algorithm to maximize ELBO. $L(q, \theta) = \sum_z q(z|x) \cdot \log \frac{p(x, z; \theta)}{q(z|x)}$

① E-step: Given θ fixed. optimize q . Let's assume current $\theta = \theta^{old}$. $ELBO \Rightarrow L(q, \theta^{old})$ is a functional of q . Also when θ^{old} is fixed. $\log p(x, \theta^{old})$ is a constant.

To maximize ELBO \Leftrightarrow minimize $KL(q||p)$:

$$\Rightarrow q(z|x) = p(z|x, \theta^{old})$$

That is, E step just let $q(z|x)$ take posterior $p(z|x, \theta^{old})$.

② M-step. Given q fixed and optimize θ .

$$\text{Now we have } q(z|x) = p(z|x, \theta^{old})$$

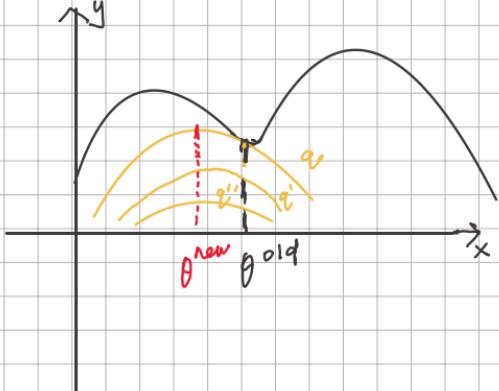
$$L(q, \theta) = \sum_z p(z|x, \theta^{old}) \cdot \log \frac{p(x, z; \theta)}{p(z|x, \theta^{old})}$$

$$= \sum_z p(z|x, \theta^{old}) \cdot \log p(x, z; \theta) - \text{const}$$

$$\Rightarrow \max_{\theta} \underbrace{\sum_z \log p(x, z; \theta)}_{\sim p(z|x, \theta^{old})} \Rightarrow \text{maximize the expectation}$$

③ Repeat until convergence.

EM 算法示意图



choice: q
E step: choose ELBO = $\log P(X|\theta)$

M step: maximize ELBO

choice θ.

For many circumstances
it has the closed solution.

$$\bar{z}_i = [0, 0, \dots, \frac{1}{k}, 0]$$

$$z_{ik} = 1$$

② EM Algorithm for MoG (n observations; k clusters)

$$\log P(X, \theta) = \sum_{i \in [n]} \log \left(\prod_{k \in [k]} \pi_k \cdot N(x_i | \mu_k, \Sigma_k) \right)$$

$$P(x) = \sum P(x|z) \cdot P(z)$$

(1) E-step. Compute $q(z|x) = P(z|x, \mu, \Sigma, \pi)$

$$\begin{aligned} P(z, x, \mu, \Sigma, \pi) &= \prod_{i \in [n]} P(x_i | z_i, \mu, \Sigma) \cdot P(z_i, \pi) \\ &\quad \underbrace{\prod_{k \in [k]} N(x_i | \mu_k, \Sigma_k)}^{z_{ik}} \underbrace{\prod_{k \in [k]} \pi_k}_{z_{ik}} \\ &= \prod_{i \in [n]} \prod_{k \in [k]} \left[N(x_i | \mu_k, \Sigma_k) \cdot \pi_k \right]^{z_{ik}} \end{aligned}$$

$$\therefore P(z|x, \mu, \Sigma, \pi) = \frac{1}{C} \prod_{i \in [n]} \prod_{k \in [k]} \left[\pi_k N(x_i | \mu_k, \Sigma_k) \right]^{z_{ik}}$$

$$= \prod_{i \in [n]} \left[\frac{1}{C} \prod_{k \in [k]} \left[\pi_k N(x_i | \mu_k, \Sigma_k) \right]^{z_{ik}} \right]$$

$$C = \prod_{i \in [n]} C_i \cdot C_i = P(x_i)$$

$$P(z|x) = \prod_{i \in [n]} P(z_i | x_i)$$

$$P(z_i | x_i, \theta^{\text{old}})$$

$$M \text{ step: } \max_{\theta} E_{\bar{z}} [\log P(x, \bar{z}; \mu, \Sigma, \pi)], \quad \bar{z} \sim P(z|x, \theta^{old})$$

$$= E_{\bar{z}} \left[\sum_{i \in [n]} \sum_{k \in [K]} \bar{z}_{ik} \cdot \log (\pi_k \cdot N(x_i; \mu_k, \Sigma_k)) \right]$$

$$= \sum_{i \in [n]} \sum_{k \in [K]} E_{\bar{z}}[\bar{z}_{ik}] \log (\pi_k \cdot N(x_i; \mu_k, \Sigma_k))$$

$$E_{\bar{z}}(\bar{z}_{ik}) = 1 \cdot P(z_{ik}=1|x_i) + 0 \cdot \sim = \frac{P(z_{ik}=1|x_i)}{P(x_i)}$$

$$= \frac{\pi_k \cdot N(x_i; \mu_i, \Sigma_k)}{\sum_{j \in [K]} \pi_j \cdot N(x_i; \mu_j, \Sigma_j)} = r_{ik} \rightarrow \text{soft assignment to } k.$$

$$\Rightarrow \max_{\theta} \sum_{i \in [n]} \sum_{k \in [K]} r_{ik} \cdot \log (\pi_k N(x_i; \mu_k, \Sigma_k))$$

$$\text{Language } L(\mu, \Sigma, \pi, \lambda) = \sum_{i \in [n]} \sum_{k \in [K]} r_{ik} (\pi_k N(x_i; \mu_k, \Sigma_k)) + \lambda (1 - \sum_{k \in [K]} \pi_k)$$

$$\frac{\partial L}{\partial \pi_k} = \sum_{i \in [n]} r_{ik} \cdot \frac{1}{\pi_k} - \lambda = 0$$

$$\Rightarrow \pi_k - \lambda = \sum_{i \in [n]} r_{ik} \Rightarrow \lambda \sum_{k \in [K]} \pi_k = \sum_{i \in [n]} \sum_{k \in [K]} r_{ik} = n$$

$$\therefore \pi_k = \frac{\sum_{i \in [n]} r_{ik}}{n} \Rightarrow \lambda = n$$

$$\text{For } N(x_i; \mu_k, \Sigma_k) = \frac{1}{(2\pi)^{\frac{1}{2}} |\Sigma_k|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (x_i - \mu_k)^T \Sigma_k^{-1} (x_i - \mu_k) \right)$$

$$\Leftrightarrow L = \sum_{i \in [n]} \sum_{k \in [K]} r_{ik} \left(\log \pi_k - \frac{1}{2} \log 2\pi - \frac{1}{2} \log |\Sigma_k| - \frac{1}{2} \sim \right)$$

$$\frac{\partial L}{\partial \mu_k} = - \sum_{i \in [n]} r_{ik} \Sigma_k^{-1} (x_i - \mu_k) = 0$$

$$\Rightarrow \sum_{i \in [n]} r_{ik} (x_i - \mu_k) = 0$$

$$\Rightarrow \mu_k = \frac{\sum_{i \in [n]} r_{ik} x_i}{\sum_{i \in [n]} r_{ik}}$$

$$\frac{\partial |x|^{-1}}{\partial x} = -|x|^{-1} (x^{-1})^T$$

$$\Rightarrow -|\Sigma_k| \Sigma_k$$

$$\text{Finally, } \frac{\partial L}{\partial \Sigma} = 0 \Rightarrow \sum_{i \in [n]} r_{ik} \frac{1}{|\Sigma_k|} \cdot \frac{\partial |\Sigma_k|}{\partial \Sigma} + r_{ik} (x_i - \mu_k)^T (x_i - \mu_k)^{-1} = 0$$

$$\Rightarrow -\sum_{i \in [n]} r_{ik} \bar{\Sigma}_k + \sum_{i \in [n]} r_{ik} (x_i - m_k) (x_i - m_k)^T = 0$$

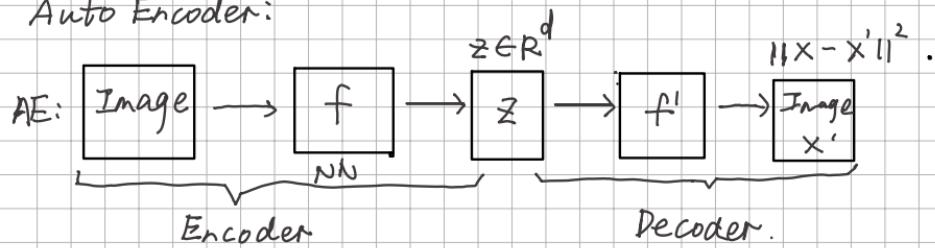
$$\Rightarrow \bar{\Sigma}_k = \frac{\sum_{i \in [n]} r_{ik} (x_i - m_k) (x_i - m_k)^T}{\sum_{i \in [n]} r_{ik}}$$

when $\bar{\Sigma}_k = 6^2 I$, $6 \rightarrow 0$

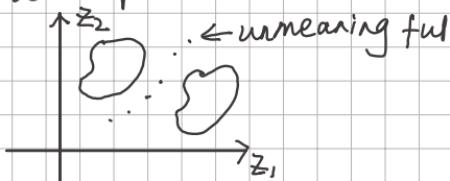
$M \otimes G \rightarrow K$ means

Lecture 10. Variational AutoEncoder (VAE)

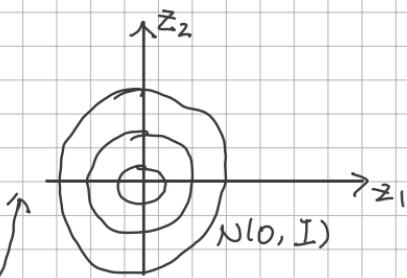
Auto Encoder:



It's z space



You will find hard to sample z



VAE Algorithm:

Given a set of training examples. $D = \{x_1, x_2, \dots, x_n\}$, we want to generate new data similar to D .

such as images.

Some points:

① In general, sampling from a distribution $X \sim p(x)$ is difficult! Actually some sample distributions are easy to sample from: $U[a, b]$, $N(\mu, \Sigma)$, $N(0, I)$

MC MC

⇒ what if we sample z from a simple distribution, then use deterministic $f(z; \theta)$ to map $z \rightarrow X$ (output)

$$z \rightarrow R.V. \rightarrow X \rightarrow R.V. \quad p(x) = \int p(x|z) \cdot p(z) dz$$

\downarrow \downarrow

$N(x|f(z; \theta), \sigma^2)$ $N(0, I)$



MoG

Different of z discrete in MoG, here we use z as one continuous distribution, we want to optimize θ .

objective: MLE of $p(x)$ on D to learn θ .

Can we sample a large amount of z from $P(z)$ to approximate $p(x)$?
 $P(x) \approx \frac{1}{N} \sum_{z_i \sim P(z)} P(x|z_i)$, it cannot work because $P(z)$ is high-dimensional need a super large # of samples z to actually approximate $p(x)$!

Second thought: MLE of $\log p(x)$ w/ hidden variables z (intractable)

Reminds us of EM!

Here. E-step: Compute $P(z_i|x_i, \theta^{old}) = \frac{P(z_i|x_i, \theta^{old}) \cdot p(x_i)}{\int z_i P(x_i|z_i) \cdot P(z_i|\theta^{old}) dz_i}$

M-step: $\theta = \operatorname{argmax}_\theta \int z \log p(x, z; \theta) \cdot P(z|x_i, \theta^{old}) dz$
EM not work.

Essential Reason: why GD/EM not work, is $P(z|x, \theta)$ not intractable. Find a tractable variational distribution $q(z|x; \theta)$ to approximate $P(z|x; \theta)$!

Goal:

$$\log p(x; \theta) = \underbrace{\sum_z q(z|x) \cdot \log \frac{p(x, z; \theta)}{q(z|x)}}_{\text{ELBO}} - \underbrace{\sum_z q(z|x) \cdot \log \frac{p(z|x; \theta)}{q(z|x)}}_{\text{KL divergence}}$$

$$\text{连续后: } \log p(x; \theta) = \underbrace{\int_z q(z|x) \cdot \log \frac{p(x, z; \theta)}{q(z|x)} dz}_{\text{ELBO}} - \underbrace{\int_z q(z|x) \cdot \log \frac{p(z|x; \theta)}{q(z|x)} dz}_{\text{KL divergence.}}$$

VAE ignores KL divergence, maximize ELBO!

Hope KL divergence is not large, maximize ELBO

Now Let's rewrite ELBO

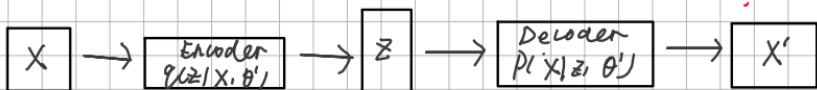
$$\text{ELBO} = \underbrace{\int_z q(z|x; \theta') \log p(x|z; \theta) dz}_{\text{Reconstruction quality}} + \underbrace{\int_z q(z|x; \theta) \log \frac{p(z)}{q(z|x; \theta)} dz}_{= -KL(q(z|x; \theta) || p(z))}$$

$$E_{z \sim q(z|x; \theta')} \log p(x|z; \theta)$$

Encoder Network

regularize $q(z|x; \theta)$
to approximate $p(z)$!

Decoder Network.



$$\begin{aligned} P(x|z, \theta) &= N(x|f(z, \theta), \sigma^2 I) \\ \Rightarrow \log P(x|z, \theta) &= \log \frac{1}{(2\pi\sigma)^d} \exp\left(-\frac{\|x - f(z, \theta)\|^2}{2\sigma^2}\right) \\ &= C - \frac{1}{2\sigma^2} \|x - f(z, \theta)\|^2 \end{aligned}$$

$\hat{E}_{z \sim q(z|x, \theta')} \|x - x'\|^2 \approx \|x - x'\|^2$ decoded from 1 z encoded

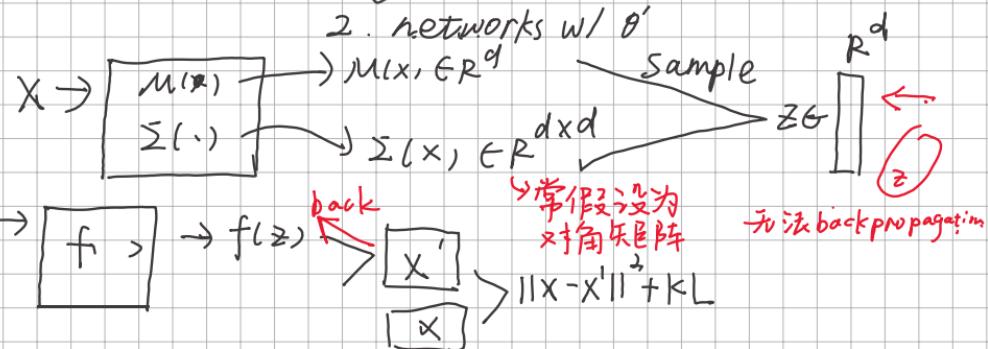
$$\text{VAE Loss: } \min_{\theta, \theta'} \frac{1}{n} \sum_{i \in \mathcal{D}} \|x_i - x'_i\|^2 + \text{KL}(q(z|x_i, \theta) || p(z))$$

Usually, we need to balance the interaction of $\text{KL}(q(z|x_i, \theta) || p(z))$

So we introduce β , VAE Loss = $\|x_i - x'_i\|^2 + \beta \cdot \text{KL}$. divergence.

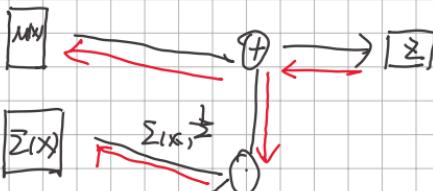
$$\beta \rightarrow 0 \Rightarrow \text{AE}$$

$$q(z|x, \theta') = N(z | M(x), \Sigma(x))$$



Reparameterization. Trick.

$$\bar{z} \sim N(0, I_d). \quad z = \bar{z} + \epsilon \sim N(\mu(x), \Sigma(x))$$



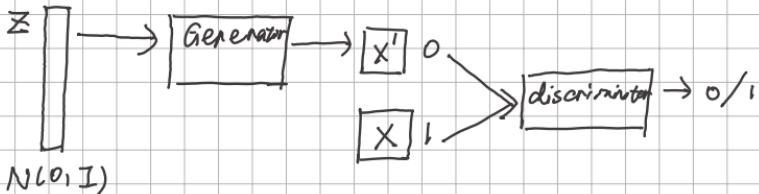
$$\epsilon \sim N(0, I_d)$$

IDC · No parameter.

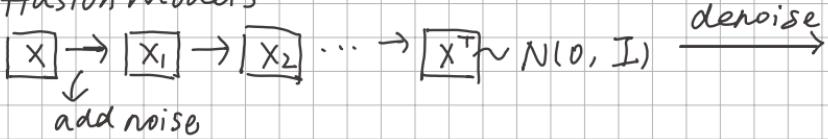
KL divergence in VAE

$$\begin{aligned}
 \text{KL divergence} &= \text{KL}(q(z|x, \theta') || p(z)) = \text{KL}(N(z|\mu(x), \Sigma(x)) || \\
 &\quad N(0, I)) \\
 &= \int_{\mathbb{R}} N(z|\mu, \Sigma) \cdot \log \frac{N(z|0, \Sigma)}{N(0, I)} dz \\
 &= \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(z-\mu)^T \Sigma^{-1}(z-\mu)\right) \cdot \log \left[|\Sigma|^{-\frac{1}{2}} \right. \\
 &\quad \left. \exp\left(-\frac{1}{2}(z-\mu)^T \Sigma^{-1}(z-\mu) + \frac{1}{2} z^T z\right) \right] \\
 &= -\frac{1}{2} \log |\Sigma| - \frac{1}{2} E_z[(z-\mu)^T \Sigma^{-1}(z-\mu)] + \underbrace{\frac{1}{2} E_z(z^T z)}_{\rightarrow E_z[(z-\mu)^T z]} \\
 &= -\frac{1}{2} \log |\Sigma| - \frac{1}{2} \underbrace{E_z[(z-\mu)^T \Sigma^{-1}(z-\mu)]}_{\text{tr}(ABC) = \text{tr}(CBA) = \text{tr}(CAB)} + \underbrace{\frac{1}{2} E_z[(z-\mu)^T z]}_{\text{tr}(A) + \text{tr}(B)} \\
 &\quad + \underbrace{E_z[(z-\mu)^T \mu]}_{\text{tr}(AB) + 2E_z[(z-\mu)^T \mu]} \\
 &\therefore E_z(\text{tr}(\Sigma^{-1}(z-\mu)(z-\mu)^T)) \\
 &= \text{tr}(\Sigma^{-1} E_z(z-\mu)(z-\mu)^T) = d \\
 &= \frac{1}{2} [\text{tr}(\Sigma) + \mu^T \mu - d - \log |\Sigma|]
 \end{aligned}$$

GAN



Diffusion Models



Encoder

