

# 机器学习（本科生公选课）GEC6531

## 第10节 损失函数和正则化 Loss Function and Regularization

计算机科学与技术学院

张瑞 教授

邮箱: [ruizhang6@hust.edu.cn](mailto:ruizhang6@hust.edu.cn)

# 签到 & 思考

## ■ 微助教签到（学校要求）

1. 加入课堂：微信扫码或者通过微助教公众号



二维码有效期至: 2024-11-16

课堂名称: GEC6531 机器学习 (公选课)

课堂编号: OA628

---

1、扫码关注公众号: 微助教服务号。

2、点击系统通知: “[点击此处加入【GEC6531 机器学习 \(公选课\)】课堂](#)”, 填写学生资料加入课堂。

---

\*如未成功收到系统通知, 请点击公众号下方“学生” - “全部(A)” - “加入课堂” --- “输入课堂编号”手动加入课堂

2. 微信扫码签到

## 回顾SVM的损失函数

# 今天的目录

- 经验风险最小化
  - 背景和定义
- 二分类损失函数
- 回归损失函数
- 不适定问题 (ILL-posed Problems)
  - 无关参数的问题
  - 不适定问题
- 正则化
  - 定义和案例
  - 正则项
- 常见例子

# 今天的目录

## ■ 经验风险最小化

- 背景和定义

## ■ 二分类损失函数

## ■ 回归损失函数

## ■ 不适定问题 (ILL-posed Problems)

- 无关参数的问题
- 不适定问题

## ■ 正则化

- 定义和案例
- 正则项

## ■ 常见例子

# SVM的损失函数

对于无约束SVM 公式:

$$\min_w C \underbrace{\sum_{i=1}^n \max[1 - y_i(\underbrace{w^\top x_i + b}_{h(x_i)}, 0]}_{\text{Hinge-Loss}} + \underbrace{\|w\|_z^2}_{\ell_2\text{-Regularizer}}$$

SVM 的损失函数用的是 Hinge-Loss形式, 而 $\ell_2$ -正则项反映了求解的复杂性, 并惩罚复杂的解。这是一个经验风险最小化的例子, 它有一个损失函数  $\ell$  和一个正则项  $r$ :

$$\min_w \frac{1}{n} \sum_{i=1}^n \underbrace{l(h_w(x_i), y_i)}_{\text{Loss}} + \underbrace{\lambda r(w)}_{\text{Regularizer}}$$

其中损失函数是惩罚训练误差的连续函数, 正则项是惩罚分类器复杂度的连续函数。

这里, 将  $\lambda$  定义为上一讲中的  $\frac{1}{C}$

# 监督学习的设置

- 监督学习问题的一般设置：  
有两个对象空间  $X$  和  $Y$ ，希望学习一个函数  $h: X \rightarrow Y$ ，对给定的  $x \in X$ ，输出  $y \in Y$ 。
- 有一个训练集，包含  $n$  个样例  $(x_1, y_1), \dots, (x_n, y_n)$ 。其中  $x_i \in X$  是输入且  $y_i \in Y$  是希望从  $h(x_i)$  得到的输出。
- 形式化地描述，即：假设有一个在  $X$  和  $Y$  上的联合概率分布  $p(x, y)$ ，训练集包含  $n$  个从  $p(x, y)$  中 i.i.d. 采样得到的实例  $(x_1, y_1), \dots, (x_n, y_n)$ 。
- 联合概率分布的假设允许我们对预测中的不确定性建模。因为对于固定的  $x$ ， $y$  不是  $x$  的确定函数，而是一个具有条件分布的随机变量  $p(y|x)$ 。



# 经验风险最小化 (Empirical Risk Minimization, ERM)

假设我们得到了一个非负实值的损失函数  $L(\hat{y}, y)$ ，用于比较假设的预测  $\hat{y}$  与真实结果  $y$  的差异。与假设  $h(x)$  相关的风险定义为损失函数的期望：

$$R(h) = E[L(h(x), y)] = \int L(h(x), y) dP(x, y)$$

学习算法的目标是在一组固定的函数  $h$  中找到一个假设  $h^*$ ，使  $R(h)$  的风险最小：

$$h^* = \underset{h \in H}{\operatorname{argmin}} R(h).$$

一般来说，风险  $R(h)$  无法计算，因为  $P(x, y)$  的分布对于学习算法来说是未知的（这种情况被称为不可知论学习）。不过，我们可以通过平均训练集上的损失函数计算一个近似值，称为经验风险：

$$R_{\text{emp}}(h) = \frac{1}{n} \sum_{i=1}^n L(h(x_i), y_i).$$

**经验风险最小化 (ERM)** 是统计学习理论的一个原则。它定义了一系列学习算法，用于给出其性能的理论上界。

# 今天的目录

- 经验风险最小化
  - 背景和定义
- 二分类损失函数
- 回归损失函数
- 不适定问题 (ILL-posed Problems)
  - 无关参数的问题
  - 不适定问题
- 正则化
  - 定义和案例
  - 正则项
- 常见例子

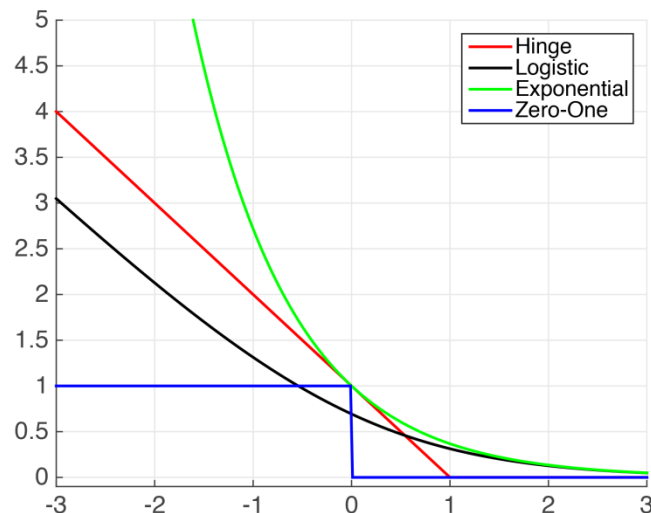


# 二分类损失函数

不同的机器学习算法使用不同的损失函数

Table 1: 二分类损失函数,  $y \in \{-1, +1\}$

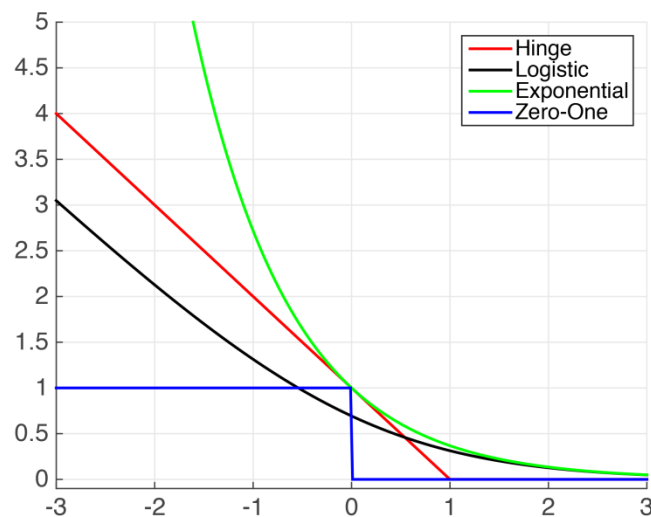
损失 $\ell(h_w(x_i, y_i))$	使用模型	评价
<b>Hinge-Loss:</b> $\max [1 - h_w(x_i)y_i, 0]^p$	<ul style="list-style-type: none"><li>● 标准 SVM (<math>p = 1</math>)</li><li>● (可微的) 均方 Hingeless SVM (<math>p = 2</math>)</li></ul> <p>图1: 常见分类损失函数图。x 轴: <math>h(x_i)y_i</math> (预测准确率); y 轴: 损失值</p>	当用于标准 SVM 时, 损失函数表示线性分隔面与其在任何一类中的最近点之间的间隔的大小。只有当 $p = 2$ 时可微。
<b>Log-Loss:</b> $\log(1 + e^{-h_w(x_i)y_i})$	逻辑回归	机器学习中最流行的损失函数之一, 因为它的输出是经过良好校准的概率。



# 二分类损失函数

Table 1: 二分类损失函数,  $y \in \{-1, +1\}$

损失函数 $\ell(h_w(\mathbf{x}_i, y_i))$	使用	评价
<b>Exponential</b> $e^{-h_w(\mathbf{x}_i)y_i}$ <b>Loss:</b>	AdaBoost	这个函数非常激进。错误预测的损失随着 $-h_w(\mathbf{x}_i)y_i$ 的值呈指数增长。这可能会使其有良好的收敛结果, 例如 Adaboost, 但它在处理噪声数据时会有问题。
<b>0-1</b> $\delta(\text{sign}(h_w(\mathbf{x}_i)) \neq y_i)$ <b>Loss:</b>	实际分类损失	非连续, 其优化在实际应用中非常困难



# 二分类损失函数

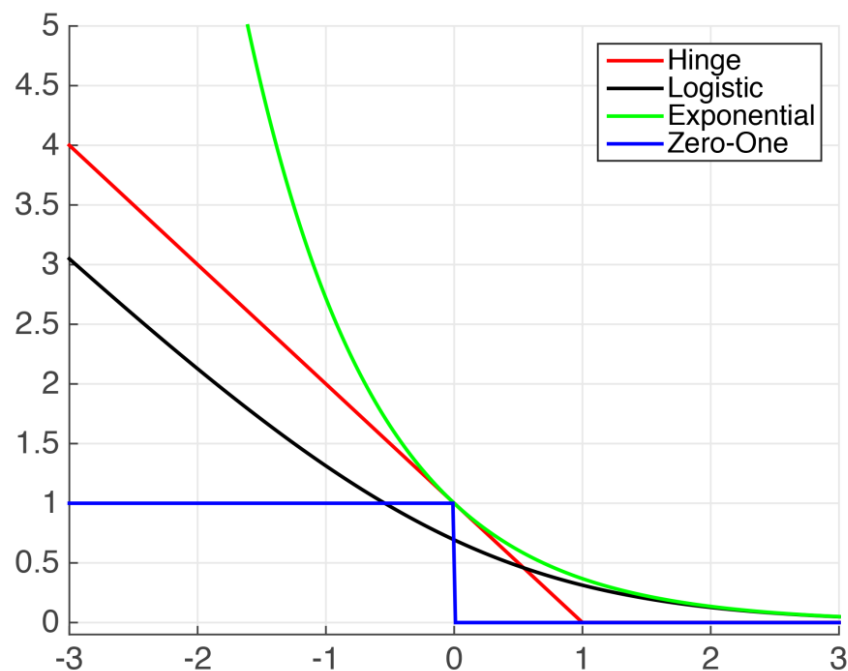


图1: 常见分类损失函数图。x 轴:  $h(x_i) y_i$  (预测准确率) ; y 轴: 损失值

**思考:** 哪些函数是在0-1 损失的严格上界?

**思考:** 对于 Hinge-loss 和 log-loss 在  $z \rightarrow -\infty$  时, 分析其性质差异?

# 今天的目录

- 经验风险最小化
  - 背景和定义
- 二分类损失函数
- 回归损失函数
- 不适定问题 (ILL-posed Problems)
  - 无关参数的问题
  - 不适定问题
- 正则化
  - 定义和案例
  - 正则项
- 常见例子

# 回归损失函数

回归算法的损失 (预测可以为任何实数)

表 1: 回归损失函数,  $y \in \mathbb{R}$

损失函数 $\ell(h_w(\mathbf{x}_i), y_i)$	评价
均方损失: $(h(\mathbf{x}_i) - y_i)^2$	<ul style="list-style-type: none"><li>最通用的回归损失函数</li><li>估计<u>平均</u> 的标签</li><li>也被称为最小二乘法 (OLS)</li><li>缺点: 对异常值/噪声较敏感</li></ul>
绝对值损失: $ h(\mathbf{x}_i) - y_i $	<ul style="list-style-type: none"><li>这也是一个很常用的损失函数</li><li>估计<u>中位数</u> 的标签</li><li>优点: 对噪音不敏感</li><li>缺点: 在 0 处不可微</li></ul>

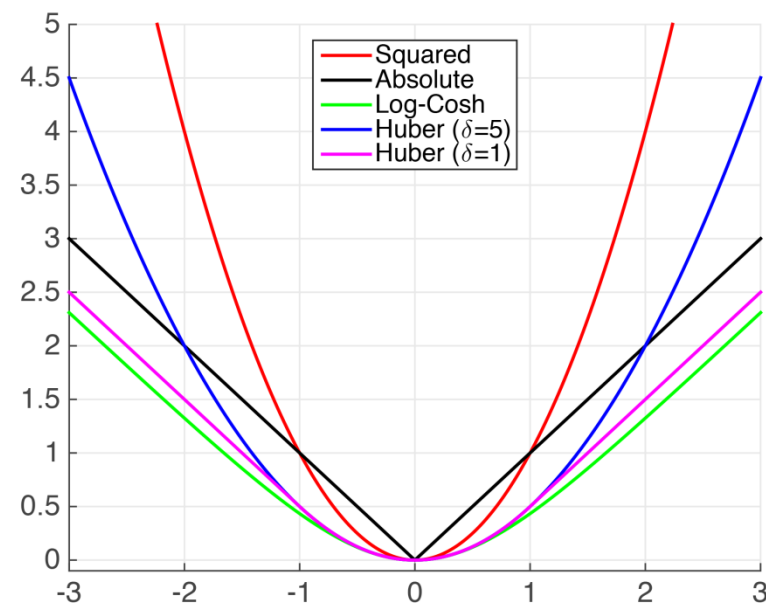


图2: 常见的回归损失函数; x 轴:  $h(\mathbf{x}_i) - y_i$  或预测的“误差”; y 轴: 损失函数值

# 回归损失函数

回归算法的损失 (预测可以为任何实数)

表 1: 回归损失函数,  $y \in \mathbb{R}$

损失函数 $\ell(h_w(\mathbf{x}_i), y_i)$	评价
<b>Huber 损失:</b> <ul style="list-style-type: none"><li><math>\frac{1}{2} (h(\mathbf{x}_i) - y_i)^2</math> if <math> h(\mathbf{x}_i) - y_i  &lt; \delta</math>,</li><li>否则 <math>\delta( h(\mathbf{x}_i) - y_i  - \frac{\delta}{2})</math></li></ul>	<ul style="list-style-type: none"><li>也被称为平滑绝对值损失</li><li>一次可微</li><li>在损失较小时表现为平方损失, 在损失较大时表现为绝对损失</li><li>优点: “两全其美” 的均方 和绝对 损失</li></ul>
<b>Log-Cosh Loss:</b> $\log(\cosh(h(\mathbf{x}_i) - y_i))$ , $\cosh(x) = \frac{e^x + e^{-x}}{2}$	优点: 类似于 Huber 损失, 但处处可微

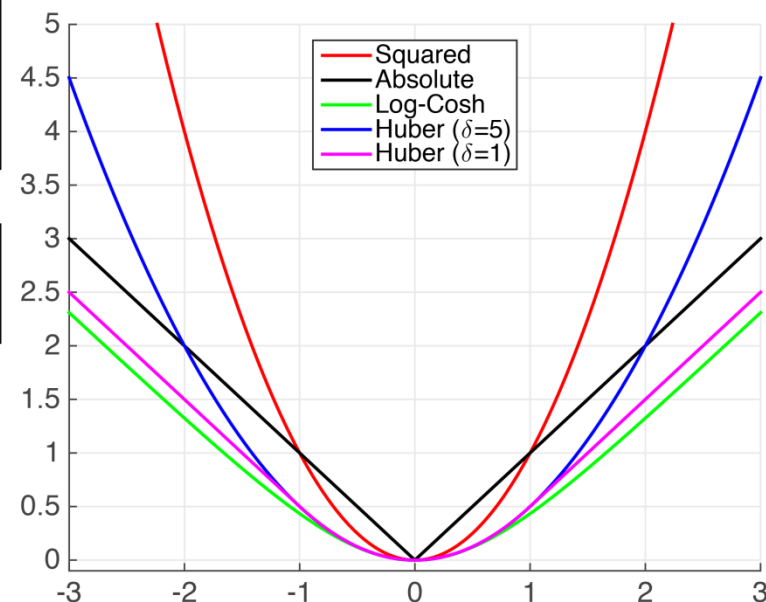


图2: 常见的回归损失函数; x 轴:  $h(\mathbf{x}_i) - y_i$  或预测的“误差”; y 轴: 损失函数值

# 今天的目录

## ■ 经验风险最小化

- 背景和定义

## ■ 二分类损失函数

## ■ 回归损失函数

## ■ 不适定问题 (ILL-posed Problems)

- 无关参数的问题
- 不适定问题

## ■ 正则化

- 定义和案例
- 正则项

## ■ 常见例子

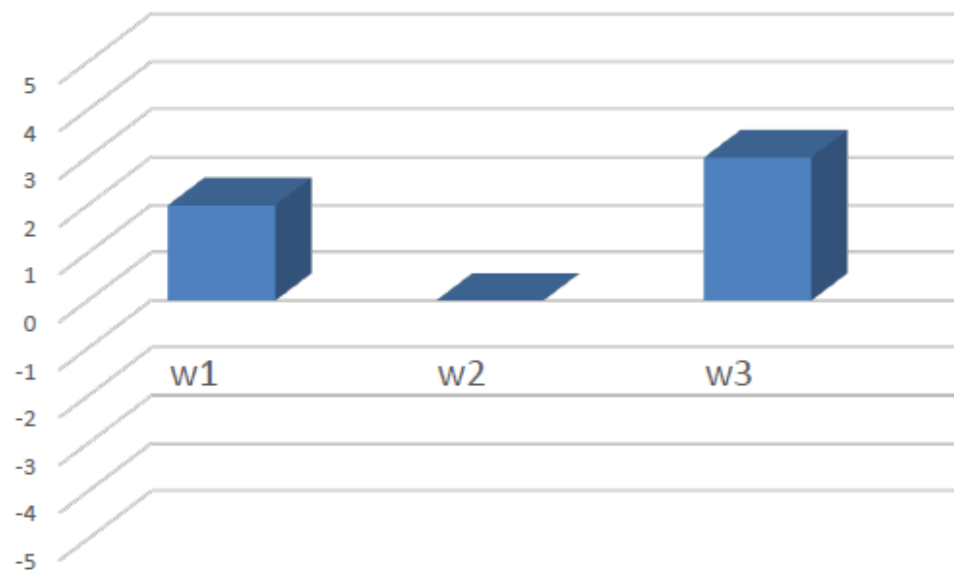


# 无关参数的问题：参数重要性

## ■ 3个特征的线性回归模型

- $n = 4$ 条特征数据
- 模型： $y = w_1x_1 + w_2x_2 + w_3x_3 + w_0$

哪个特征更重要？

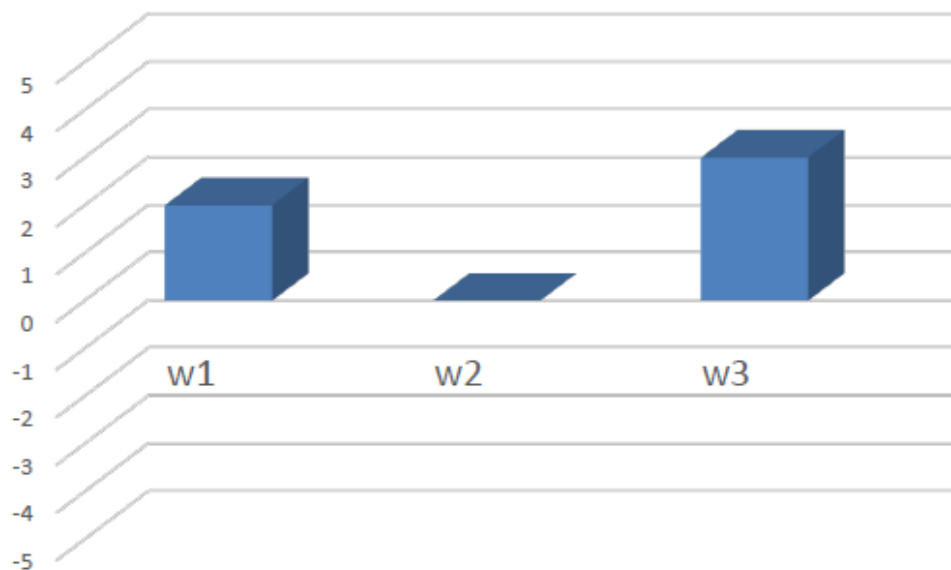


# 无关参数的问题：参数重要性

## ■ 3个特征的线性回归模型，前2列一模一样

- $n = 4$ 条特征数据
- 模型： $y = w_1x_1 + w_2x_2 + w_3x_3 + w_0$
- 特征1和2中有一个是无关参数

3	3	7
6	6	9
21	21	79
34	34	2



添加扰动对模型的效果如何？

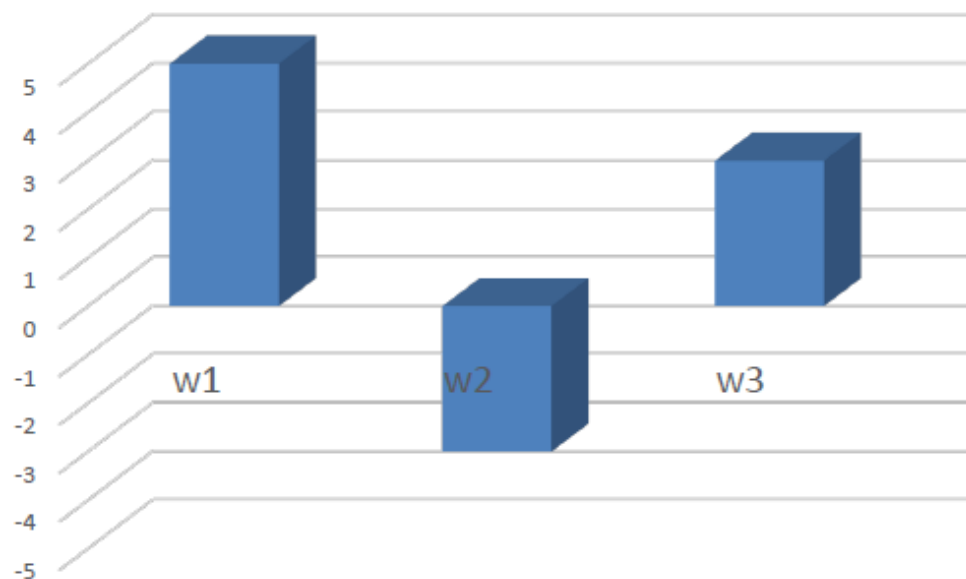
- 加  $\delta$  到  $w_1$  上
- 从  $w_2$  减去  $\delta$
- 模型值不变

# 无关参数的问题：参数重要性

## ■ 3个特征的线性回归模型，前2列一模一样

- $n = 4$ 条特征数据
- 模型： $y = w_1x_1 + w_2x_2 + w_3x_3 + w_0$
- 特征1和2中有一个是无关参数

3	3	7
6	6	9
21	21	79
34	34	2



添加扰动对模型的效果如何？

- 加  $\delta$  到  $w_1$  上
- 从  $w_2$  减去  $\delta$
- 模型值不变

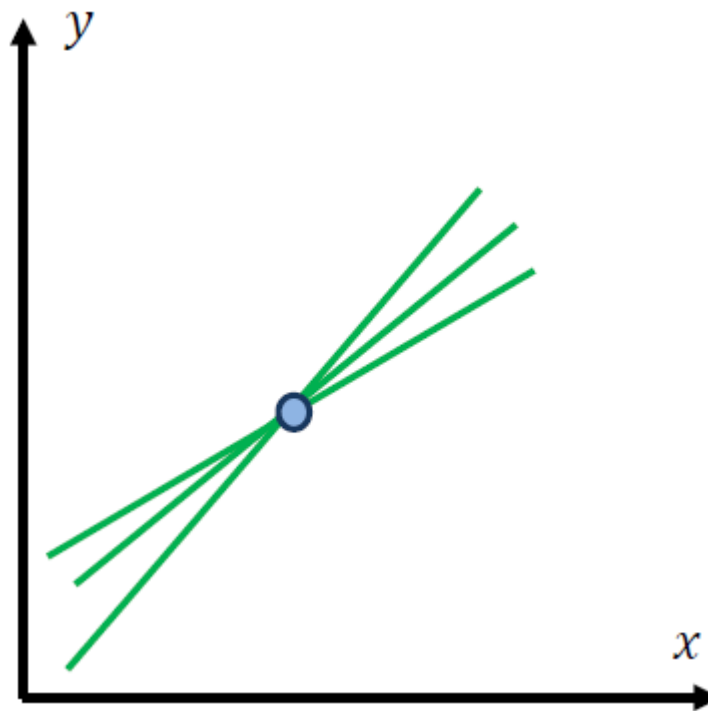
# 无关参数的问题：数据缺乏

## ■ 极端例子

- 模型有2个参数
- 只有1个数据点

## ■ 无法确定解

## ■ 也是不适定问题/病态问题



# 无关参数的问题：参数重要性

- 假设  $[\hat{w}_0, \hat{w}_1, \hat{w}_2, \hat{w}_3]'$  是最优解
- 对于任何  $\delta$ , 那么  $[\hat{w}_0, \hat{w}_1 + \delta, \hat{w}_2 - \delta, \hat{w}_3]'$ 
  - 有同样的预测值
  - 有同样的误差平方和
- 带来如下问题
  - 解不唯一
  - 没有可解释性
  - 使得优化模型参数变成一个**不适定问题/病态问题** (ill-posed problem) : 问题定义不合理

# 不适定问题 (ILL-posed Problems)

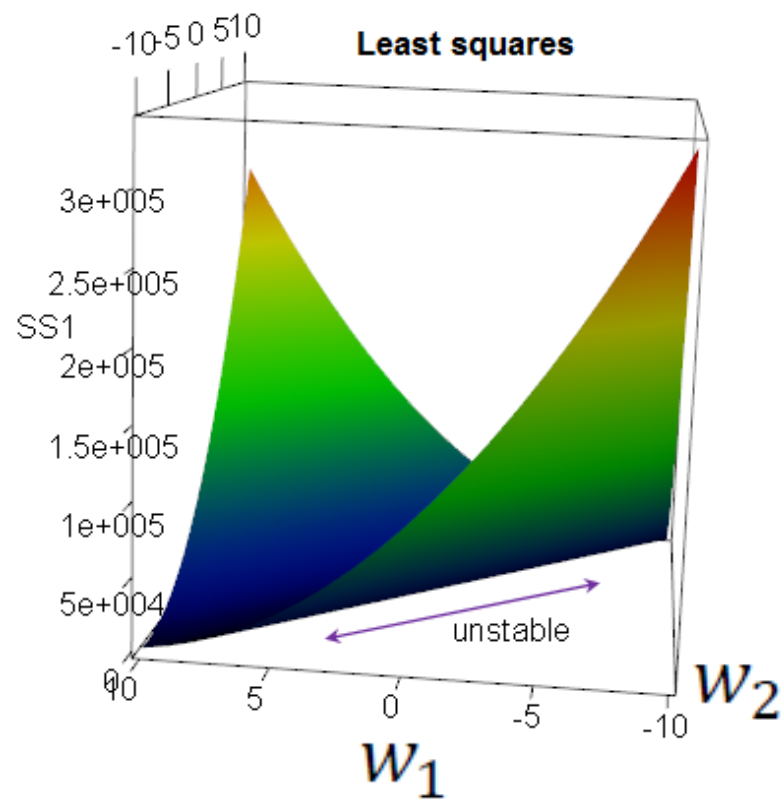
## ■ 问题的解无法明确定义

- $w_1$  和  $w_2$  无法唯一确定，有无数可能

## ■ 线性回归的解析解：

$$w = (XX^T)^{-1}Xy$$

## ■ 当存在无关特征时， $XX^T$ 没有逆矩阵



凸函数，但不严格

# 今天的目录

- 经验风险最小化
  - 背景和定义
- 二分类损失函数
- 回归损失函数
- 不适定问题 (ILL-posed Problems)
  - 无关参数的问题
  - 不适定问题
- 正则化
  - 定义和案例
  - 正则项
- 常见例子



# 解决方案：正则化 (Regularization)

- 正则化：添加一个条件来约束参数的长度

- 原来的问题要最小化：

$$\|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2$$

- 正则化后的问题要最小化：

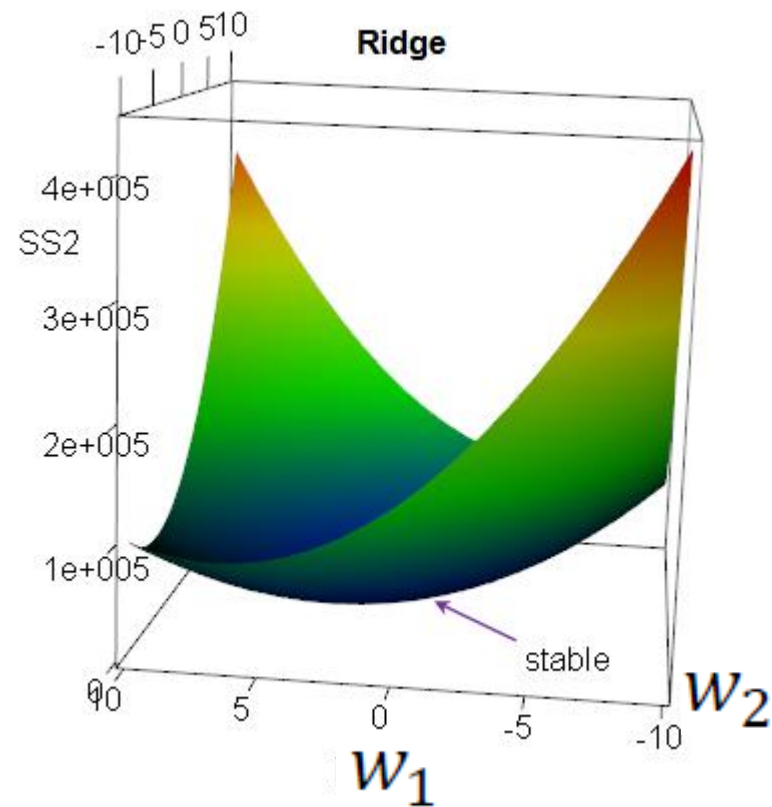
$$\|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \|\mathbf{w}\|_2^2 \text{ for } \lambda > 0$$

- 问题的解变成

$$\mathbf{w} = (\mathbf{X}\mathbf{X}^T + \lambda \mathbf{I})^{-1} \mathbf{X}\mathbf{y}$$

- 这个加了正则化的解叫做脊回归 (ridge regression)

- 将线形的“屋脊”变成圆拱的单个顶点
- 将  $\lambda$  加入  $\mathbf{X}\mathbf{X}^T$  的特征向量，变成可逆矩阵



严格凸函数

# 正则项

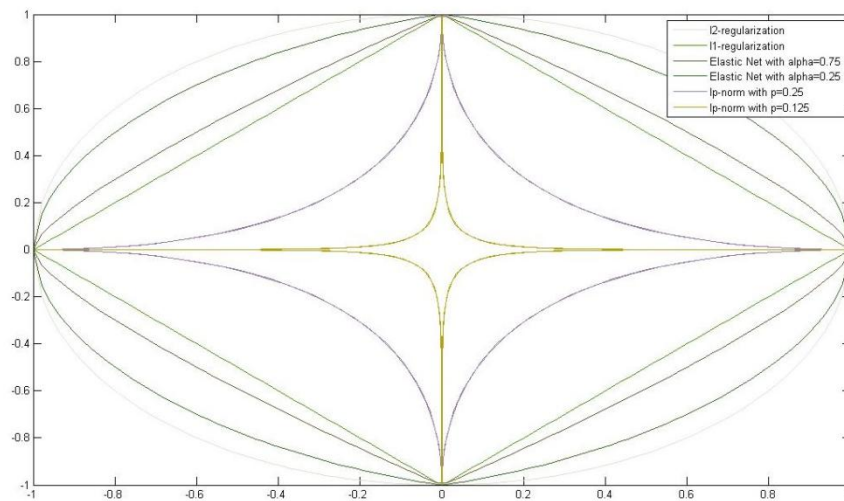
- 在数学、统计学和计算机科学中，特别是在机器学习问题中，正则项是为了解决不适定（ill-posed）问题或防止过拟合而添加的信息。
- 当我们研究正则项时，它有助于改变优化问题的公式，以获得更好的几何直观。
- 在前面的章节中， $l_2$  在支持向量机中引入正则项反映解的复杂度限制。

除了  $l_2$  正则项，其他类型的有效正则项及其特点在下表中列出。

# 正则项

表 3: 正则项类型

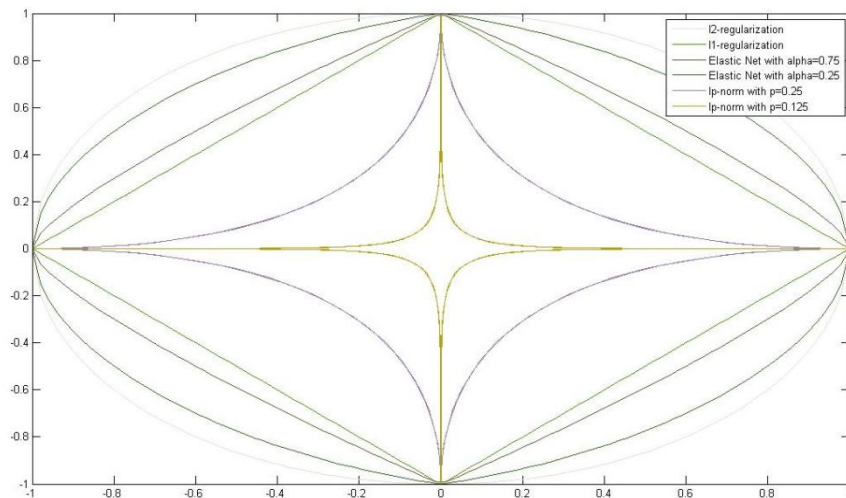
正则项 $r(\mathbf{w})$	性质
$l_2$ -Regularization $r(\mathbf{w}) = \mathbf{w}^\top \mathbf{w} = \ \mathbf{w}\ _2^2$	<ul style="list-style-type: none"><li>• 优点: 严格凸; 可微</li><li>• 缺点: 对所有特征都有权重, 即在某种程度上依赖所有特征 (理想情况下我们希望避免这种情况)–即稠密解。</li></ul>
$l_1$ -Regularization $r(\mathbf{w}) = \ \mathbf{w}\ _1$	<ul style="list-style-type: none"><li>• 凸 (但不严格)</li><li>• 缺点: 在 0 处不可微 (这是最小化的目的)</li><li>• 效果: 稀疏 解</li></ul>



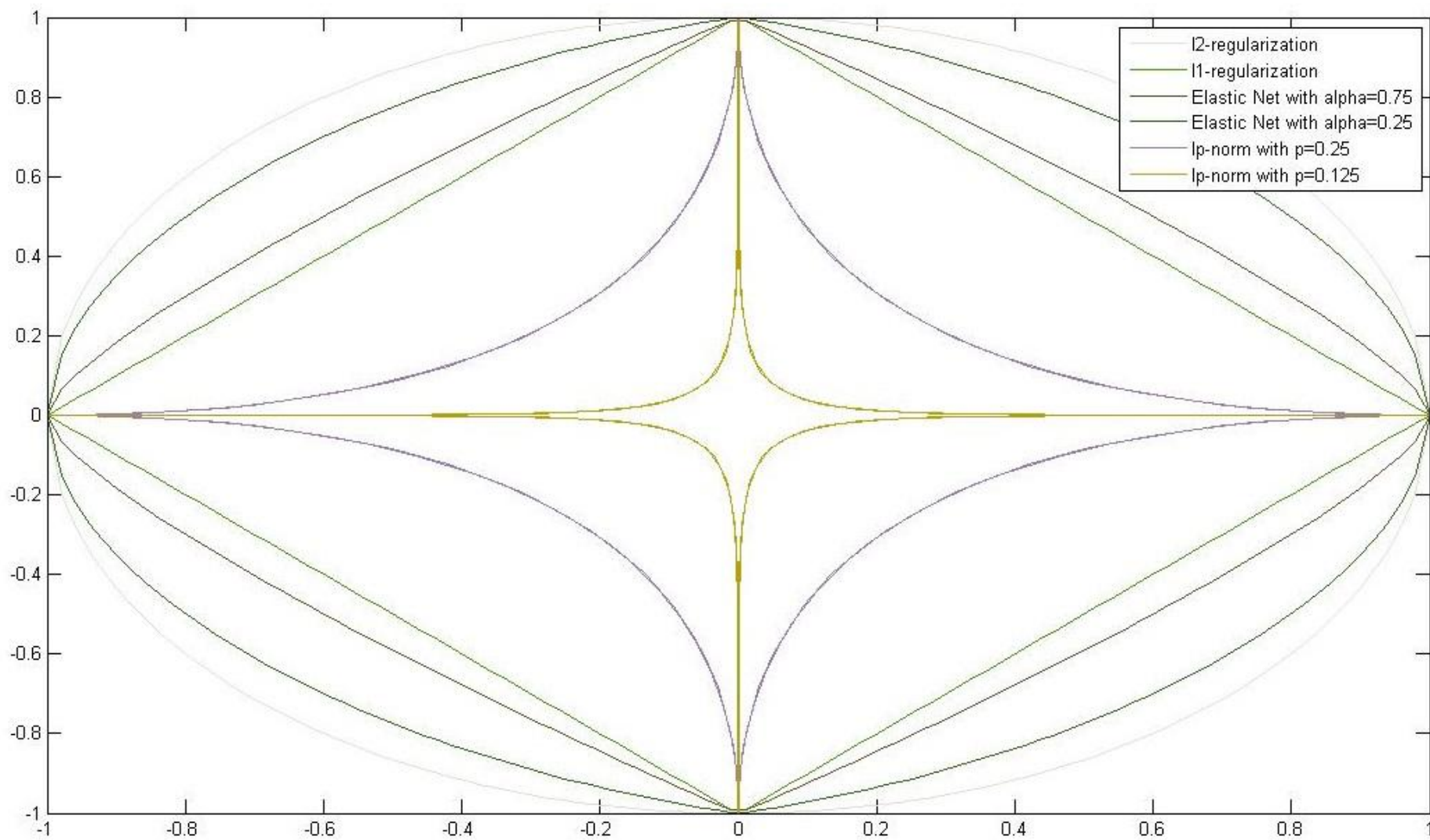
# 正则项

表 3: 正则项的类型

正则项 $r(\mathbf{w})$	特性
$l_p$ -范数 $\ \mathbf{w}\ _p = (\sum_{i=1}^d v_i^p)^{1/p}$	<ul style="list-style-type: none"><li>• (往往 <math>0 &lt; p \leq 1</math>)</li><li>• 依赖于初始化</li><li>• 优点: 非常稀疏的解</li><li>• 缺点: 非凸; 不可微</li></ul>
弹性网络: $\alpha \ \mathbf{w}\ _1 + (1 - \alpha) \ \mathbf{w}\ _2^2 \quad \alpha \in [0, 1]$	<ul style="list-style-type: none"><li>• 优点: 严格凸 (即唯一解)</li><li>• 优点: 稀疏诱导 (有利于特征选择); 均方损失支持向量机的对偶形式, 见SVEN</li><li>• 缺点: 不可微</li></ul>

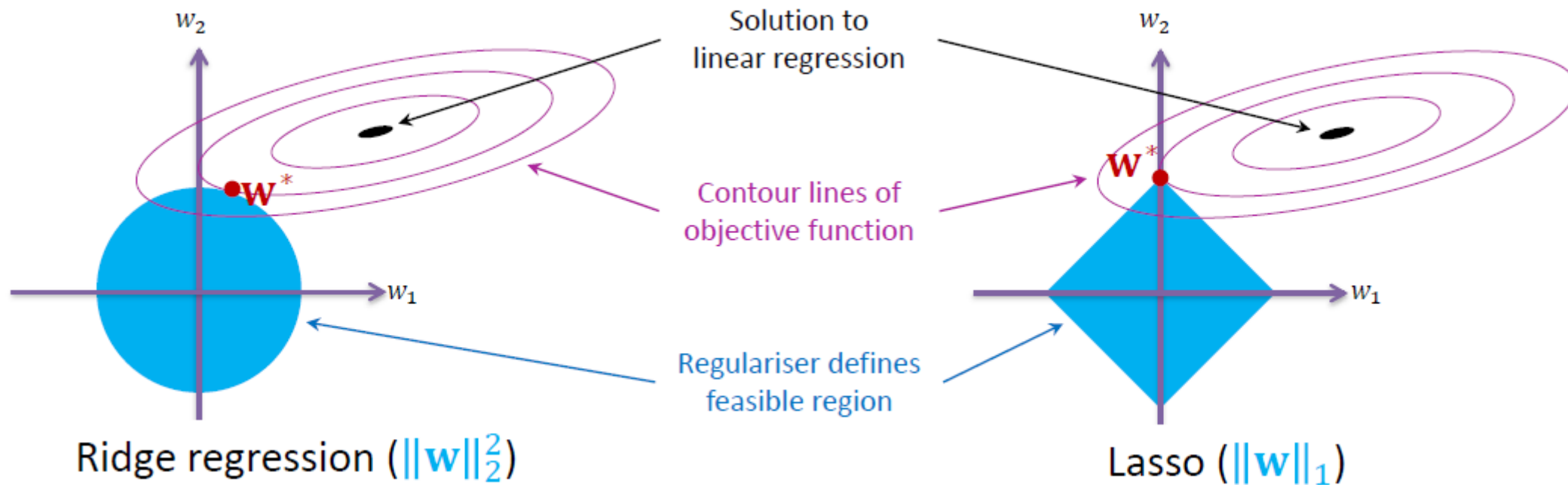


# 正则项



# 把正则项看成约束

minimise  $\|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2$  subject to  $\|\mathbf{w}\|_2^2 \leq \lambda$  for  $\lambda > 0$



Lasso ( $L_1$  正则) 倾向坐标轴上的解, 也就是有些维度是0, 解更加稀疏

# 正则项

## 思考

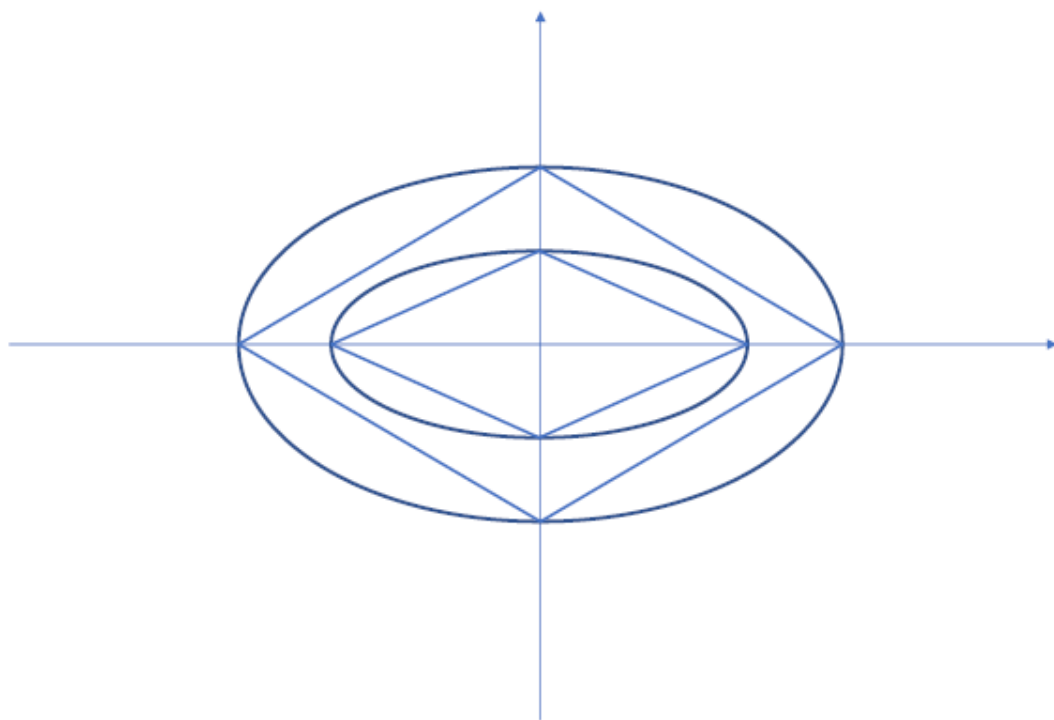
以下的损失函数等高线分别对应于哪个不同类型的正则项和超参数？

(a)  $\ell(w) + \lambda \|w\|_1$

(b)  $\ell(w) + 2 \cdot \lambda \|w\|_1$

(c)  $\ell(w) + \lambda \|w\|_2$

(d)  $\ell(w) + 2 \cdot \lambda \|w\|_2$



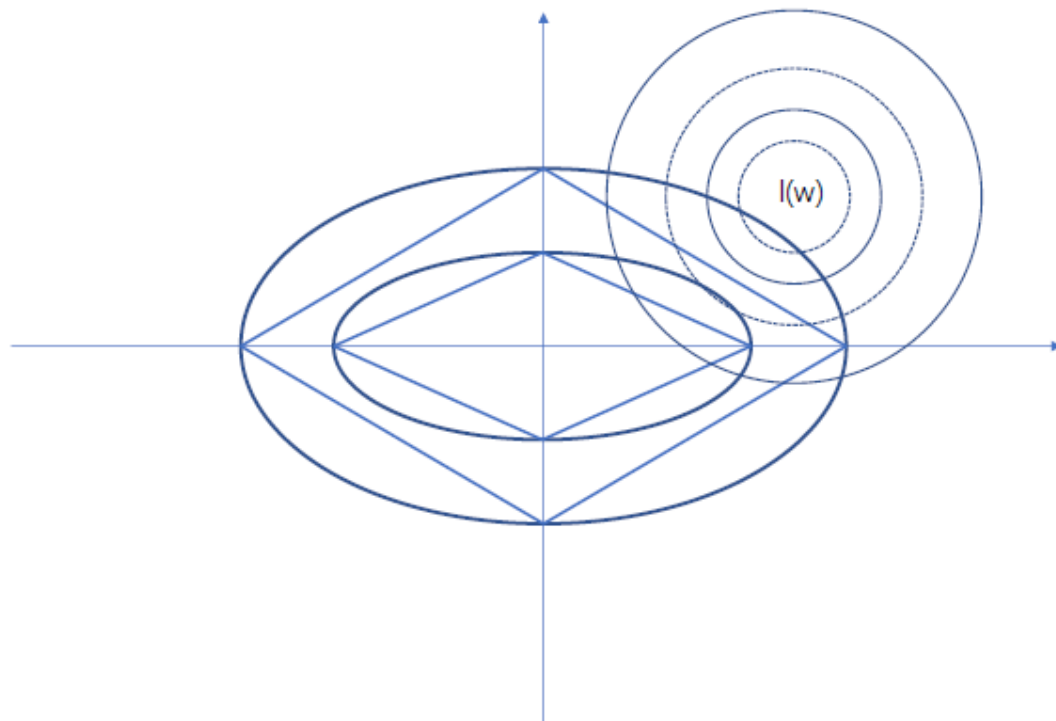


# 正则项

## 思考

以下的损失函数等高线分别对应于哪个不同类型的正则项和超参数？

(a)  $\ell(w) + \lambda \|w\|_1$       (b)  $\ell(w) + 2 \cdot \lambda \|w\|_1$       (c)  $\ell(w) + \lambda \|w\|_2$       (d)  $\ell(w) + 2 \cdot \lambda \|w\|_2$



不同大小与类型正则项损失函数等高线和损失函数等高线相交结果图

# 今天的目录

- 经验风险最小化
  - 背景和定义
- 二分类损失函数
- 回归损失函数
- 不适定问题 (ILL-posed Problems)
  - 无关参数的问题
  - 不适定问题
- 正则化
  - 定义和案例
  - 正则项
- 常见例子

# 常见例子：最小二乘法 (Ordinary Least Square, OLS)

本节包括一些处理风险最小化的特殊情况，如最小二乘法，岭回归，Lasso 和逻辑回归。下面提供了关于它们的损失函数、正则项以及解的信息。

## 最小二乘法

在统计学中，最小二乘法 (OLS) 是一种用于估计线性回归模型中未知参数的线性最小二乘方法

目标:

$$\min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i - y_i)^2$$

评价:

- 平方损失
- 没有正则项
- 闭式解形式:  $\mathbf{w} = (\mathbf{X}\mathbf{X}^\top)^{-1}\mathbf{X}\mathbf{y}^\top$ ,  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ ,  $\mathbf{y} = [y_1, \dots, y_n]$
- OLS 和主成分分析 (PCA) 的第一主成分之间有一个有趣的联系。PCA 也最小化均方损失，但考虑垂直损失 (每个点与回归线之间的水平距离)。

# 常见例子：岭回归 (Ridge Regression)

## 岭回归

Tikhonov 正则项，是不适定 (ill-posed) 问题最常用的正则项方法。在统计学中，这种方法被称为“岭回归” 目标:

$$\min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 + \lambda \|\mathbf{w}\|_2^2$$

评价:

- 平方损失
- $l_2$ -正则项
- $\mathbf{w} = (\mathbf{X}\mathbf{X}^\top + \lambda \mathbf{I})^{-1} \mathbf{X}\mathbf{y}^\top$
- 如果数据不是非常高维的，速度很快
- 只有一行 Julia / Python 代码

# 常见例子: Lasso

## Lasso

**LASSO** 是一种回归分析方法, 它执行变量选择和正则项, 以提高其产生的统计模型的预测准确性和可解释性 **目标:**

$$\min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^{\top} \mathbf{x}_i - y_i)^2 + \lambda \|\mathbf{w}\|_1$$

**评价:**

- 优点: 稀疏诱导 (有利于特征选择); 凸性
- 缺点: 非严格凸 (无唯一解); 不可微 (在 0 处)
- 用梯度下降或SVEN求解

# 常见例子：弹性网络 (Elastic Net)

## 弹性网络

### 目标

$$\min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}_i - y_i)^2 + \alpha \|\mathbf{w}\|_1 + (1 - \alpha) \|\mathbf{w}\|_2^2$$
$$\alpha \in [0, 1)$$

### 评价:

- 优点: 严格凸 (即唯一解)
- + 稀疏性诱导 (有利于特征选择)
- + 平方损失 SVM 的对偶形式, 见SVEN
- 缺点: 不可微
- 用 (次) 梯度下降法或SVEN

# 常见例子：逻辑回归 (Logistic Regression)

## 逻辑回归

在回归分析中，**逻辑回归**是估计逻辑模型的参数；这是二项回归的一种形式 **目标：**

$$\min_{\mathbf{w}, b} \frac{1}{n} \sum_{i=1}^n \log (1 + e^{-y_i(\mathbf{w}^\top \mathbf{x}_i + b)})$$

**评价：**

- 通常是  $l_1$  或  $l_2$  正则项
- 用梯度下降法求解
- $\Pr(y|x) = \frac{1}{1 + e^{-y(\mathbf{w}^\top \mathbf{x} + b)}}$



# 常见例子：线性支持向量机 (Linear SVM)

## 线性支持向量机

目标:

$$\min_{\mathbf{w}, b} C \sum_{i=1}^n \max[1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b), 0] + \|\mathbf{w}\|_2^2$$

评价:

- 往往使用  $l_2$  正则项 (有时也使用  $l_1$ )。
- 二次规划
- 可以核化, 导致**稀疏解**
- 核化版本可以用专门的算法非常高效地求解 (例如SMO)

# 总结

## ■ 常用损失函数

- 二分类: Hinge-loss, Log-loss, Exponential-loss, 0-1 loss
- 回归: 均方损失, 绝对值损失, Huber loss, Log-Cosh loss

## ■ 正则化

- 将参数的长度加到损失函数里

## ■ 正则项

- 加入参数的  $l_p$  范数

## ■ 常见例子

- 最小二乘法、岭回归、Lasso、弹性网络、逻辑回归、线性SVM