① supervised learning 的算法中，各种法也较相似，但 DataAmount 等因素往往起了起了决定作用。
Feature 选择
Regularization 选择

② 考虑以后能 再讲 复杂的改→引出 SVM.

# Support Vector Machines
## Optimization objective

Machine Learning

② 但 SVM 非常 powerful. 在 both
工业界和学术界了

③ SVM 与 神经网络相比 cleaner,
与 Logistic Regression 相比也 powerful

---

## Alternative view of logistic regression

考虑并 LR的公式.
hypothes function

$$h_\theta(x) = \frac{1}{1 + e^{-\theta^T x}}$$

$$g(z) = \frac{1}{1 + \frac{1}{e^z}} \Rightarrow y = \frac{1}{1 + e^{-z}} = \frac{1}{1 + \frac{1}{e^z}}$$

$$h_\theta(x) = y(z)$$

$$\begin{cases} z \to +\infty & y \to 1 \\ z = 0 & y = \frac{1}{2} \\ z \to -\infty & y \to 0 \end{cases}$$

$0 < \frac{1}{1 + e^{-z}} < 1$

$\to z = \theta^T x$

If $y = 1$, we want $h_\theta(x) \approx 1$, $\quad \theta^T x \gg 0 \to +\infty$
If $y = 0$, we want $h_\theta(x) \approx 0$, $\quad \theta^T x \ll 0 \to -\infty$
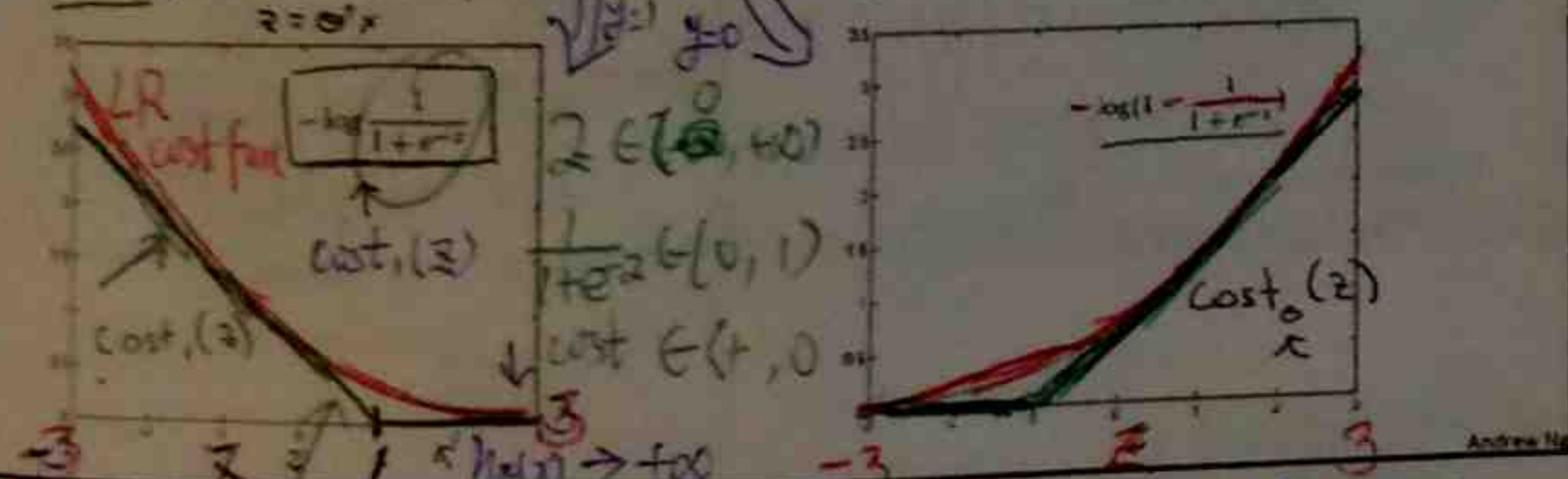
$$z = \theta^T x$$

---

## Alternative view of logistic regression

$(x, y)$ 每个 example contribute a term of (cost function)

Cost of example: $-(y \log h_\theta(x) + (1 - y) \log(1 - h_\theta(x)))$ ← cost function

$$= -y \log \frac{1}{1 + e^{-\theta^T x}} - (1 - y) \log\left(1 - \frac{1}{1 + e^{-\theta^T x}}\right)$$

代入 $h_\theta(x)$

If $y = 1$ (want $\theta^T x \gg 0$):     If $y = 0$ (want $\theta^T x \ll 0$):

$z = \theta^T x$

LR cost fun $-\log \frac{1}{1 + e^{-z}}$
$z \in (0, +\infty)$
$\frac{1}{1 + e^{-z}} \in (0, 1)$
cost$_1(z)$
cost $\in (\uparrow, 0]$
cost$_1(z)$

$-\log\left(1 - \frac{1}{1 + e^{-z}}\right)$
cost$_0(z)$

$z = h_\theta(x) \to -\infty$

修改cost function
$z = \theta^T x \to +\infty$     cost$_1(z) \to 0$
$\frac{1}{1 + e^{-z}} \to 1$     cost$_1(z) \to 0$

$y = -\log h$     $1 - z = \frac{1}{1 + e^{-z}} \to 0$
$z = 1 - \frac{1}{1 + e^{-z}} \to 1$

cost$(z) \to 0$
$\log\left(1 - \frac{1}{1 + e^z}\right) \to 0$
cost$(z_0) \to 0$

---

## Support vector machine

Logistic regression: 简化 cost func

$$\min_\theta \frac{1}{m}\left[\sum_{i=1}^{m} y^{(i)}\left(-\log h_\theta(x^{(i)})\right) + (1 - y^{(i)})\left(-\log(1 - h_\theta(x^{(i)}))\right)\right] + \frac{\lambda}{2m}\sum_{i=1}^{n}\theta_j^2$$

cost$_1(\theta^T x^{(i)})$     cost$_0(\theta^T x^{(i)})$     有 $\frac{\lambda}{m}$

把括号外的 "-"
搬到括号内.

Support vector machine:
不会改变训练得出的 θ.

$\min_u ((u - 5)^2 + 1) \to u = 5$
$\min_u (10(u - 5)^2 + 10) \to u = 5$

没有 $\frac{1}{m}$ 不改变整体θ值

$$\min_\theta C\sum_{i=1}^{m}\left[y^{(i)}\text{cost}_1(\theta^T x^{(i)}) + (1 - y^{(i)})\text{cost}_0(\theta^T x^{(i)})\right] + \frac{1}{2}\sum_{i=1}^{n}\theta_j^2$$

用 A + λB 来控制 θ.

$C = \frac{1}{\lambda}$    C 的作用类似 $\frac{1}{\lambda}$, 让

用 CA + B 来控制 θ.    $C = \frac{1}{\lambda}$

没有 $\frac{\lambda}{m}$ 也可以

CA + B = A + λB

相当于 提 3 个 ←

Regularization Term

λ、入 换为 C. ?

1

hypothesis function of LR:

□ $h_\theta(x) = \dfrac{1}{1 + e^{-\theta^T x}}$ ⟺ $g(z) = \dfrac{1}{1 + \frac{1}{e^z}}$

□ $y=1$ 时, $z \to +\infty$, $g(z) \to 1$, $-\log(g(z)) \to 0$

既然 $g(z) > \frac{1}{2}$, $z > 0$, 但候选 $z \to -\infty$ 问题

$z \to -\infty$, $g(z) \to 0$, $-\log(g(z)) \to +\infty$

同理可推 $y=0$ 时.

$\boxed{P(z) = -\lg \dfrac{1}{1+e^{-z}}}$

$\boxed{z = -\theta^T x}$

对正例, 预测得越准, cost越小 (→0)
对正例, 预测得越不准, cost越大 (→∞)

SVM: cost function 作些修改   同理可推 $y=0$ 时
分为两段直线
$z > 1$ 时, ——
$z < 1$ 时, ＼

## SVM hypothesis

$$\min_\theta C \sum_{i=1}^m \left[ y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1-y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{i=1}^n \theta_j^2$$

**Hypothesis:**

$$z = \theta^T x$$

$$h_\theta(x) = \begin{cases} 1 & \text{if } \theta^T x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

*不输出概率, 直接判断*

LR 对比. $z > 0$

$h_\theta(x) \begin{cases} \geq 0.5 & \text{if } \theta^T x \geq 0 \ 趋向(\theta^T x \to +\infty) \\ < 0.5 & \text{otherwise} \ 趋向(\theta^T x \to)\, x \end{cases}$

$z < 0$

## Support Vector Machines
## Large Margin Intuition

**Machine Learning**

SVM . Large Margin Classifier

---

## Support Vector Machine

$$\rightarrow \min_\theta C \sum_{i=1}^m \left[ y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1-y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) \right] + \frac{1}{2} \sum_{i=1}^n \theta_j^2$$

$\text{cost}_1(z)$   $z \geq 1$    $z \leq 1$   $\text{cost}_0(z)$

$\rightarrow$ If $y=1$, we want $\underline{\theta^T x \geq 1}$ (not just $\geq 0$)   $\theta^T x \geq 1$

$\rightarrow$ If $y=0$, we want $\underline{\theta^T x \leq -1}$ (not just $< 0$)   $\theta^T x \leq -1$

$C = 100,000$

当 C 非常大时.

想 SVM 时. 还以 >1, <-1 来作为判断依据

(Large Margin)

safety margin factor

---

## SVM Decision Boundary   当 C 非常大时. 前这部分被系书选于 0

$$\min_\theta C \boxed{\sum_{i=1}^m \left[ y^{(i)} \text{cost}_1(\theta^T x^{(i)}) + (1-y^{(i)}) \text{cost}_0(\theta^T x^{(i)}) \right]} + \frac{1}{2} \sum_{i=1}^n \theta_j^2$$

$= 0$

**Whenever** $y^{(i)} = 1$:
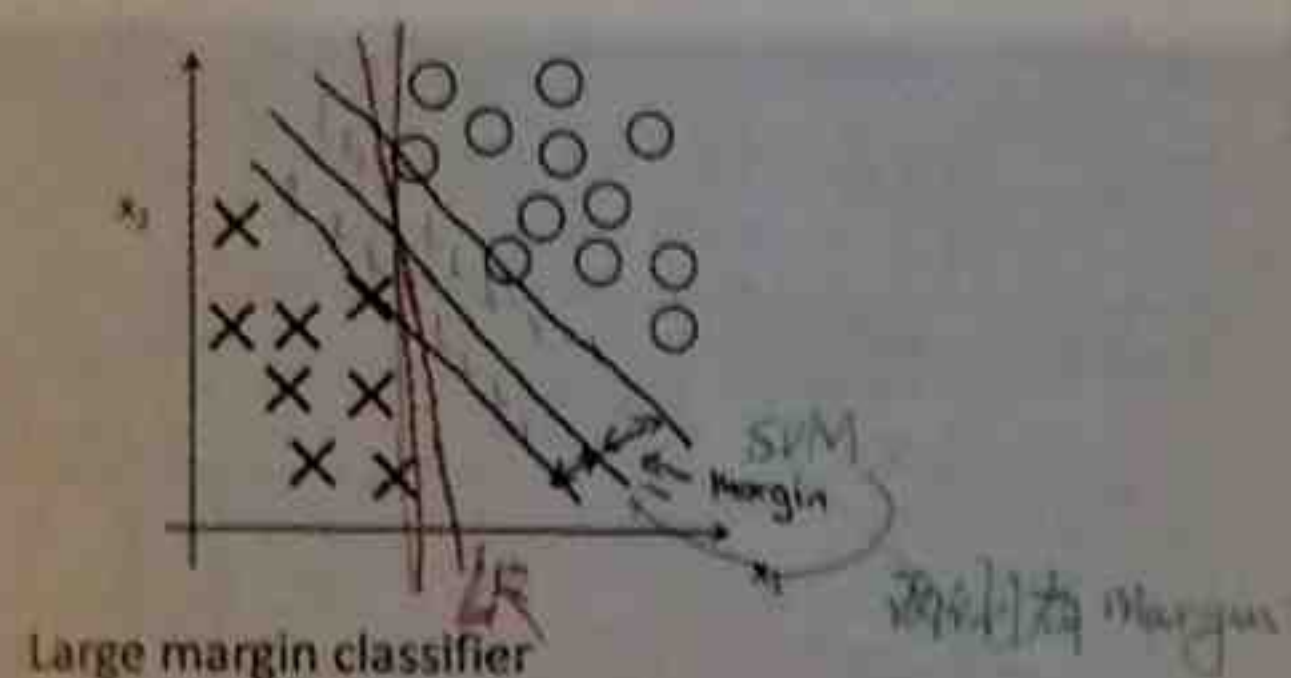
$$\theta^T x^{(i)} \geq 1$$

**Whenever** $y^{(i)} = 0$:

$$\theta^T x^{(i)} \leq -1$$

$\min_\theta C \cdot 0 + \frac{1}{2} \sum_{i=1}^n \theta_j^2$

s.t. $\theta^T x^{(i)} \geq 1$  if $y^{(i)} = 1$

$\theta^T x^{(i)} \leq -1$  if $y^{(i)} = 0$

## SVM Decision Boundary: Linearly separable case



Large margin classifier

## Large margin classifier in presence of outliers

(sensitive to ~~the~~ outliers)

$(C = \frac{1}{\lambda})$



→ $C$ very large

$\frac{1}{\lambda}$

← $C$ not too large

# Support Vector Machines

The mathematics behind large margin classification (optional)

Machine Learning

## Vector Inner Product

$u = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}$  $v = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$

$u^T v = ?$  $[u_1 \ u_2]\begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$

norm of $u$ $\|u\| = $ length of vector $u$

$= \sqrt{u_1^2 + u_2^2} \in \mathbb{R}$

$p = $ length of projection of $v$ onto $u$.

$u^T v = p \cdot \|u\| = v^T u$

$= u_1 v_1 + u_2 v_2$  $p \in \mathbb{R}$

$u^T v = p \cdot \|u\|$

$p < 0$



$u^T v$  $u_1 v_1 + u_2 v_2$

$v^T u$

$p \cdot \|u\|$

## SVM Decision Boundary

$$\min_\theta \frac{1}{2}\sum_{j=1}^{n}\theta_j^2 = \frac{1}{2}(\theta_1^2 + \theta_2^2) = \frac{1}{2}\left(\sqrt{\theta_1^2 + \theta_2^2}\right)^2 = \frac{1}{2}\|\theta\|^2$$

$\omega = (\sqrt{\omega})^2$

$= \|\theta\|$

$\theta_0 = 0$

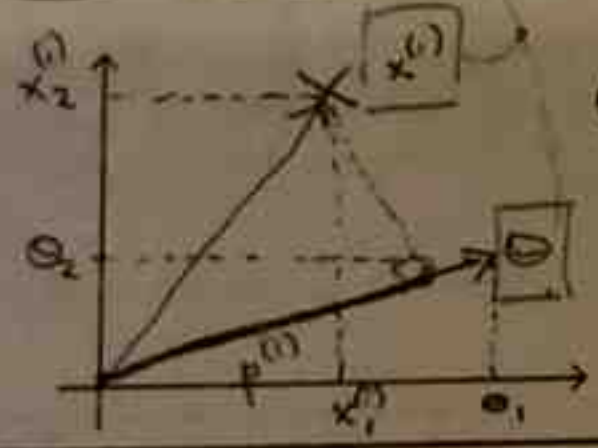s.t. $\theta^T x^{(i)} \geq 1$ if $y^{(i)} = 1$

$\theta^T x^{(i)} \leq -1$ if $y^{(i)} = 0$

Simplication: $\theta_0 = 0$ , $n=2$

cost func

cost functions

$\theta^T x = ?$

$u^T v$

$\begin{bmatrix}\theta_1\\\theta_2\end{bmatrix}$ $\theta_0 = 0$

$\theta^T x^{(i)} = p^{(i)} \cdot \|\theta\|$

$= \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)}$

因此 $\theta^T x^{(i)} \geq 1$

$\theta^T x^{(i)} \leq -1$

可以替换为

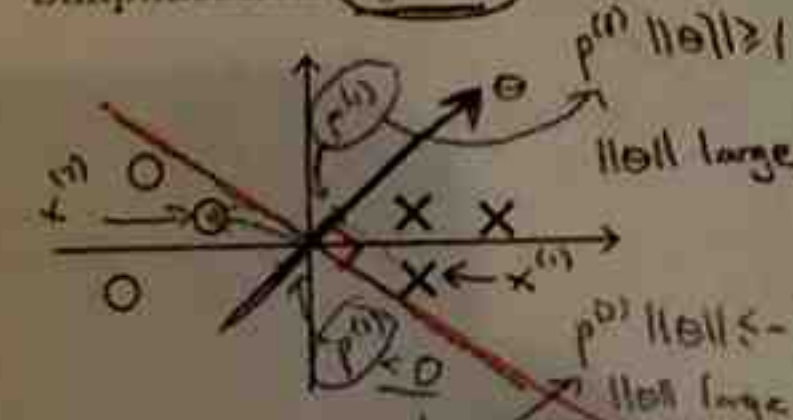$p^{(i)} \cdot \|\theta\| \geq 1$

$p^{(i)} \cdot \|\theta\| \leq -1$

## SVM Decision Boundary

$$\min_\theta \frac{1}{2}\sum_{j=1}^{n}\theta_j^2 = \frac{1}{2}\|\theta\|^2$$

s.t. $p^{(i)} \cdot \|\theta\| \geq 1$ if $y^{(i)} = 1$

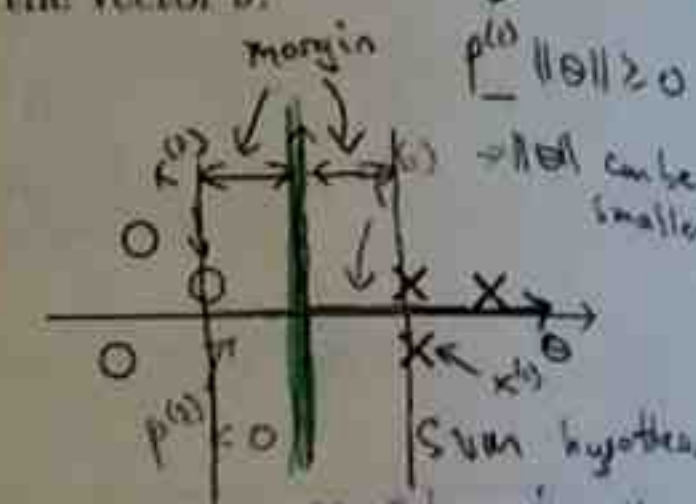$p^{(i)} \cdot \|\theta\| \leq -1$ if $y^{(i)} = 1$

where $p^{(i)}$ is the projection of $x^{(i)}$ onto the vector $\theta$.

Simplification: $\theta_0 = 0$

C varies large

large

$p^{(i)} \cdot \|\theta\| \geq 1$

$\|\theta\|$ large

$p^{(i)} \cdot \|\theta\| \leq -1$

$\|\theta\|$ large

margin

$p^{(i)} \cdot \|\theta\| \geq 0$

$\|\theta\|$ can be smaller.

SVM hypothesis

给 x, θ 在图上投影

$p^{(i)} > 0$ , $p^{(i)} < 0$

新投影小, 使得 $\|\theta\|$ 非常大,

但我们 objective 是让 $\|\theta\|$ 小.

why SVM 不位于良好 boundary. 也较大

所以深红 $\|\theta\|$ 变大.

## Support Vector Machines

## Kernels I

处理复杂的非线性分类.

Machine Learning

## Non-linear Decision Boundary

Predict $y = 1$ if

$$\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 + \theta_5 x_2^2 + \cdots > 0$$

$h_\theta(x) = \begin{cases} 1 & \text{if } \theta_0 + \theta_1 x_1 + \cdots \geq 0 \\ 0 & \text{otherwise} \end{cases}$

$x_2$

$x_1$

$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 + \cdots$

$f_1 = x_1, \quad f_2 = x_2, \quad f_3 = x_1 x_2, \quad f_4 = x_1^2, \quad f_5 = x_2^2 \cdots$

Is there a different / better choice of the features $f_1, f_2, f_3, \ldots$?

我们不知道这些 $f_i$ 是否有用.

算起 expensive

## Kernel

如何定义新 features

Given $x$, compute new feature depending on proximity to landmarks $l^{(1)}, l^{(2)}, l^{(3)}$ （几近择）

$l^{(1)}$  $l^{(2)}$
$l^{(3)}$

$\|w\|$  $w = x - l^{(1)}$

$\|x - l^{(1)}\|^2$ 向量长度

Given $x$:  一个样本

$f_1 = \text{similarity}(x, l^{(1)}) = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right)$

new features  $f_2 = \text{similarity}(x, l^{(2)}) = \exp\left(-\frac{\|x - l^{(2)}\|^2}{2\sigma^2}\right)$
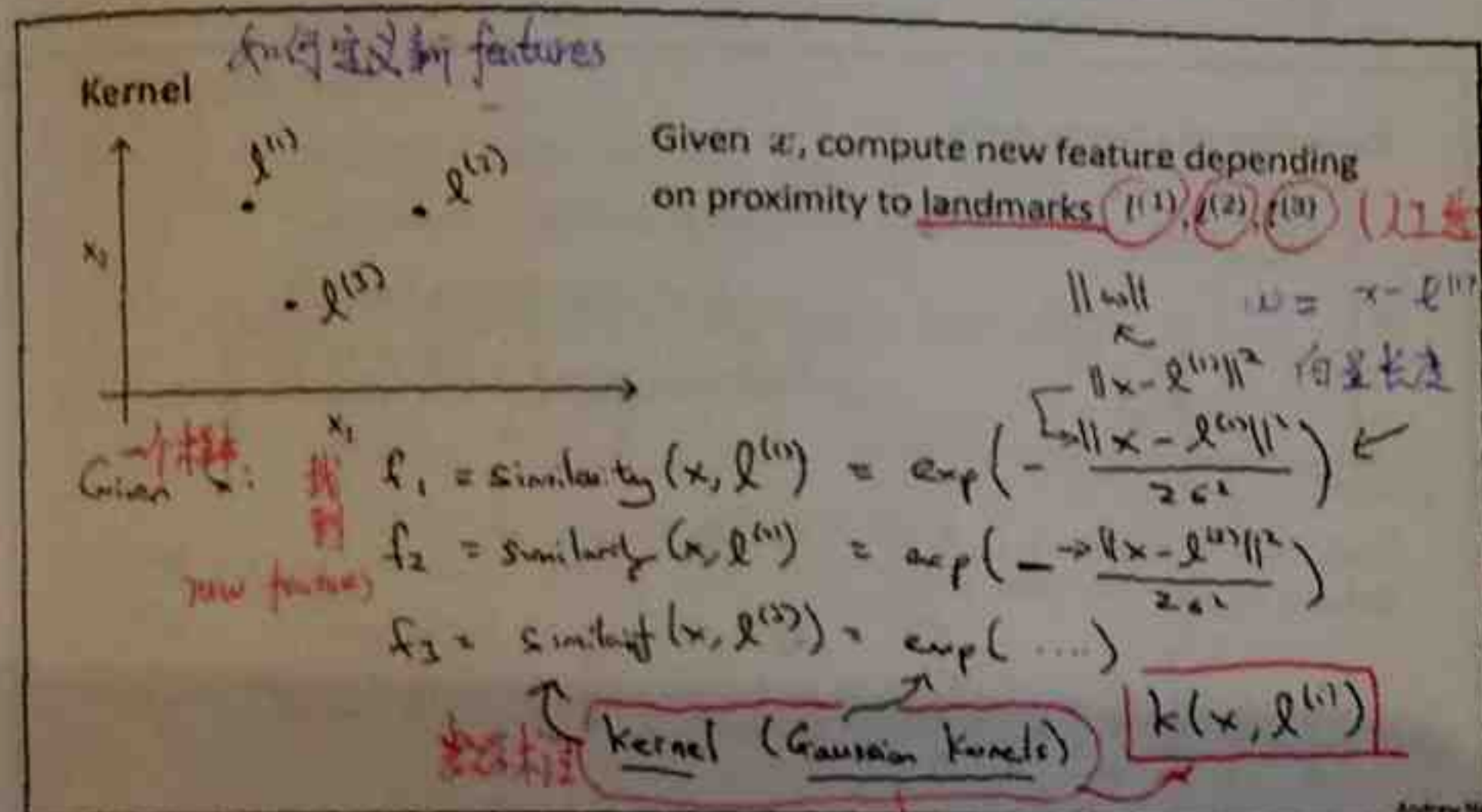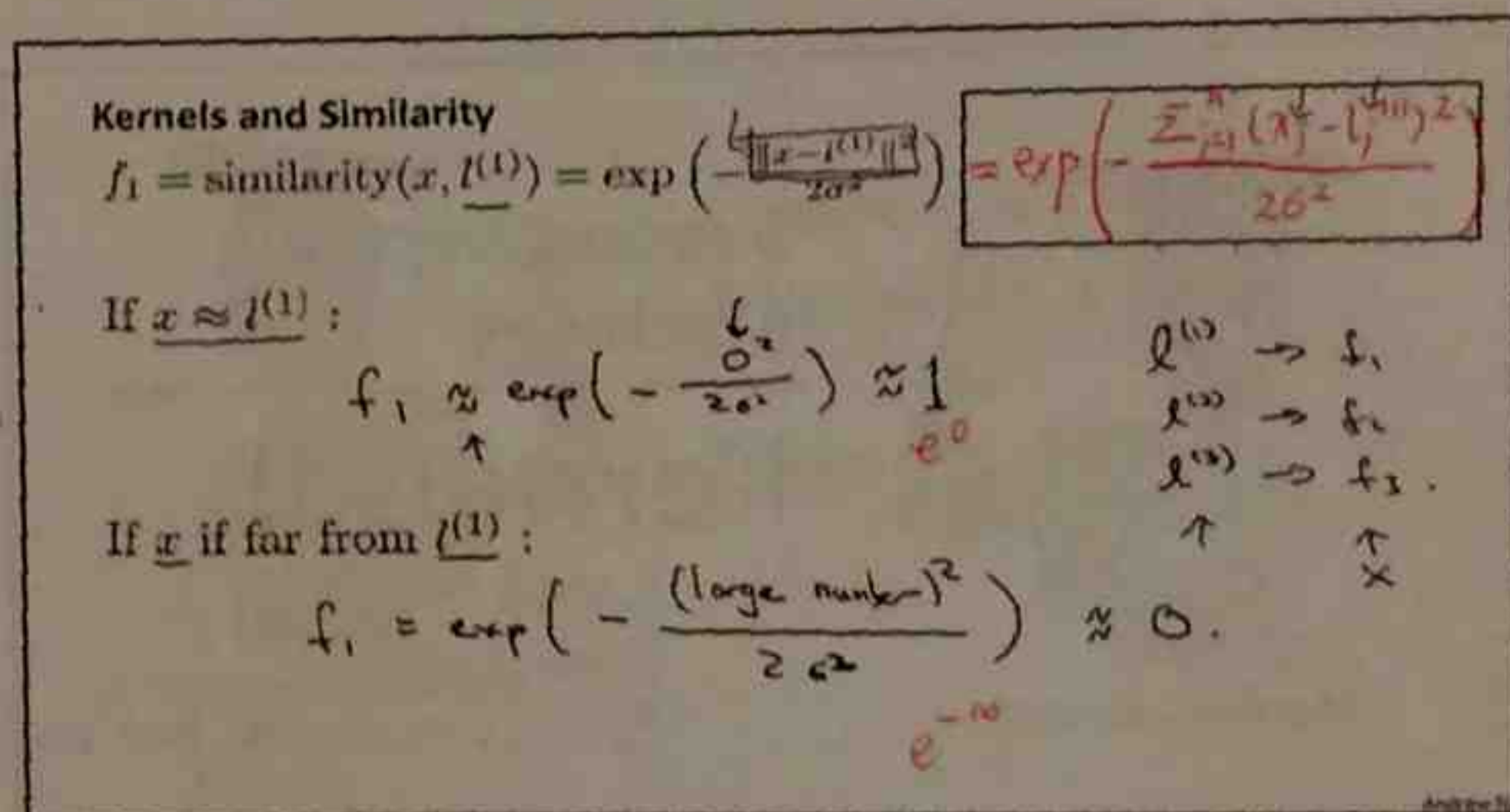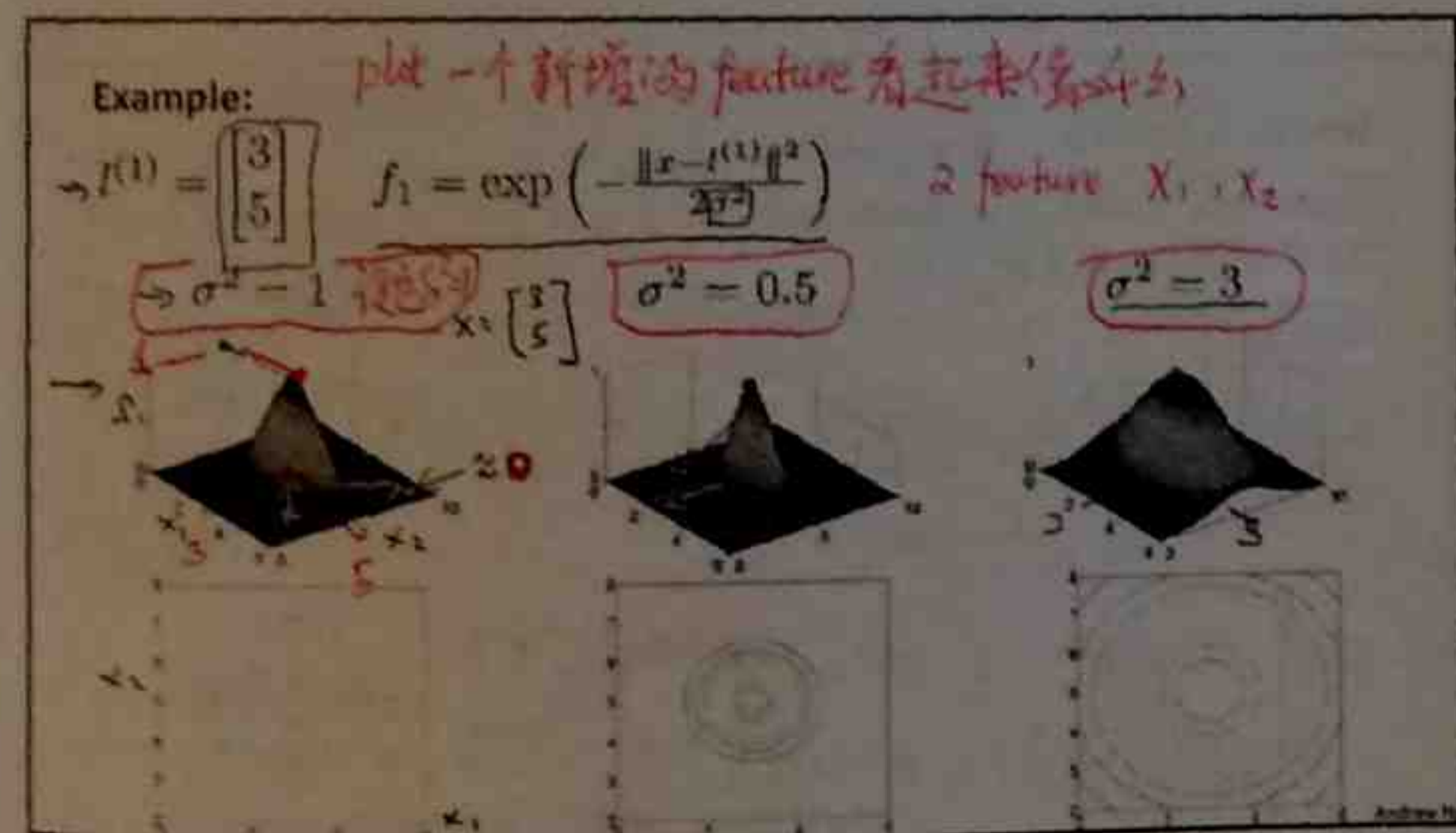
$f_3 = \text{similarity}(x, l^{(3)}) = \exp(\cdots)$

Kernel (Gaussian Kernels)  $k(x, l^{(1)})$

基个对 为断(landels)
还有其它 kernel

### Kernels and Similarity

$f_1 = \text{similarity}(x, \underline{l^{(1)}}) = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right) = \exp\left(-\frac{\sum_{j=1}^{n}(x_j - l_j^{(1)})^2}{2\sigma^2}\right)$

忽略 $x_0$ （莫为1）

If $\underline{x \approx l^{(1)}}$:

$f_1 \approx \exp\left(-\frac{0^2}{2\sigma^2}\right) \approx 1$   $e^0$

$l^{(1)} \to f_1$
$l^{(2)} \to f_2$
$l^{(3)} \to f_3$

If $\underline{x}$ if far from $\underline{l^{(1)}}$:

$f_1 = \exp\left(-\frac{(\text{large number})^2}{2\sigma^2}\right) \approx 0$   $e^{-\infty}$

如果有一个 example 位于这里, 会怎样

## Example:

plot 一个新增的 feature 看起来像什么

$l^{(1)} = \begin{bmatrix} 3 \\ 5 \end{bmatrix}$   $f_1 = \exp\left(-\frac{\|x - l^{(1)}\|^2}{2\sigma^2}\right)$   2 feature $x_1, x_2$.

$\sigma^2 = 1$   $x = \begin{bmatrix} 3 \\ 5 \end{bmatrix}$   $\sigma^2 = 0.5$   $\sigma^2 = 3$

$\approx 0$

$\delta = 0.5$
land 下降到0更快

$\delta^2 = 3$
kernel下降到0更慢

### Hypothesis function

Predict "1" when
$\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$

$l^{(1)}$  $l^{(2)}$
$l^{(3)}$

predict $y = 0$

倒数 第二例子 $\theta_3$
$\theta_0 = -0.5, \theta_1 = 1, \theta_2 = 1, \theta_3 = 0$

$f_1 \approx 1, f_2 \approx 0, f_3 \approx 0$.

$\theta_0 + \theta_1 \times 1 + \theta_2 \times 0 + \theta_3 \times 0$ (代入 $x$)
$= -0.5 + 1 = 0.5 \geq 0$ (代入 $\theta$) 得 $0.5 \geq 0$ ) predict 1.

对于example
predict $y = 0$  $f_1, f_2, f_3 \approx 0$
$\theta_0 + \theta_1 l + \cdots = -0.5 < 0$

第3个 example, 可以预测 $y = 1$. 因为发现

推导到线 通过 landmark 挑选 $l^{(1)}, l^{(2)}$ 的都是 $y = 1$; 这周围的都是 $y = 0$
的非线性分类  找到 $y = 1$ 的 boundary.

5

① 找 Landmark:
② 用 GK 找 ...
来 build $H_\theta(x)$

How to find landmarks
⇒ 放在 $m$个 examples 上
$l^{(1)}$ $l^{(2)}$ ... $l^{(m)}$
$m$个 examples ⇒ $m$个 landmark!

# Support Vector Machines

# Kernels II

Machine Learning

① detail
② how to use
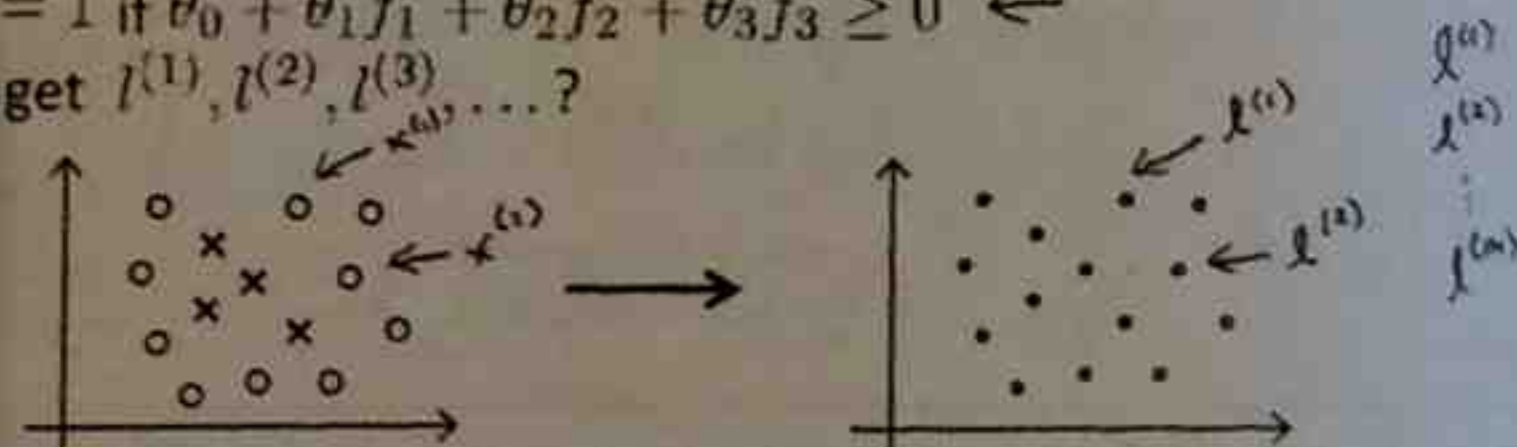③ bias – variance problems

---

**Choosing the landmarks**



Given $x$:
$\to f_i = \text{similarity}(x, l^{(i)})$
GK $= \exp\left(-\frac{||x - l^{(i)}||^2}{2\sigma^2}\right)$

Predict $y = 1$ if $\theta_0 + \theta_1 f_1 + \theta_2 f_2 + \theta_3 f_3 \geq 0$
Where to get $l^{(1)}, l^{(2)}, l^{(3)} \dots$?

$l^{(1)}$
$l^{(2)}$
$\vdots$
$l^{(m)}$

---

**SVM with Kernels**
$m$个 examples
$m$个 landmark

$\to$ Given $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})$,
$\to$ choose $l^{(1)} = x^{(1)}, l^{(2)} = x^{(2)}, \dots, l^{(m)} = x^{(m)}$.

Given example $x$:
$\to f_1 = \text{similarity}(x, l^{(1)})$
$\to f_2 = \text{similarity}(x, l^{(2)})$

$\Rightarrow f = \begin{bmatrix} f_0 \\ f_1 \\ f_2 \\ \vdots \\ f_m \end{bmatrix} \quad f_0 = 1$

Train / Cross Va / Test set

For training example $(x^{(i)}, y^{(i)})$:

$x^{(i)} \to \begin{bmatrix} f_1^{(i)} = \text{sim}(x^{(i)}, l^{(1)}) \\ f_2^{(i)} = \text{sim}(x^{(i)}, l^{(2)}) \\ f_i^{(i)} = \text{sim}(x^{(i)}, l^{(i)}) = \exp(-\frac{0}{2\sigma^2}) = 1 \\ \vdots \\ f_m^{(i)} = \text{sim}(x^{(i)}, l^{(m)}) \end{bmatrix}$
其中有一个

$x^{(i)} \in \mathbb{R}^{n+1}$ (or $\mathbb{R}^n$)

$f^{(i)} = \begin{bmatrix} f_0^{(i)} \\ f_1^{(i)} \\ \vdots \\ f_m^{(i)} \end{bmatrix} \quad f_0^{(i)} = 1$

landmark: 选
② example 位置

① feature vector
来 $\frac{1}{m}$ training example

---

**SVM with Kernels** 用 $x$ 与 landmark 相似度组成 $\frac{1}{m}$个 feature.

**Hypothesis:** Given $x$, compute features $f \in \mathbb{R}^{m-1}$ $\quad \theta \in \mathbb{R}^{n+1}$
$\to$ Predict "$y=1$" if $\theta^T f \geq 0$
$\theta_0 f_0 + \theta_1 f_1 + \dots + \theta_m f_m$ ($n = m$)

**Training:**
$\min_\theta C \sum_{i=1}^{m} y^{(i)} \text{cost}_1(\theta^T f^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T f^{(i)}) + \frac{1}{2} \sum_{j=1}^{m} \theta_j^2$

$\theta^T x^{(i)} \to \theta^T f^{(i)}$
去掉 cost function $\to \theta_0$

$\theta = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_m \end{bmatrix}$ (ignore $\theta_0$)

$\sum_i \theta_j^2 = \theta^T \theta \quad \theta = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_m \end{bmatrix}$ (ignore $\theta_0$)

$||\theta||^2$

$\to \theta^T M \theta$
Matrix

$M = 10,000$

很多 SVM 里 换一下 $\theta^T M \theta$
实质用起来还是 minimize $||\theta||$
会让 svm 更巧妙高效
Kernel 用在 LR 上会很慢。
反得上 $m = n$.

---

**✓ SVM parameters:**

$C \left(= \frac{1}{\lambda}\right)$. → Large C: Lower bias, high variance. (small $\lambda$)

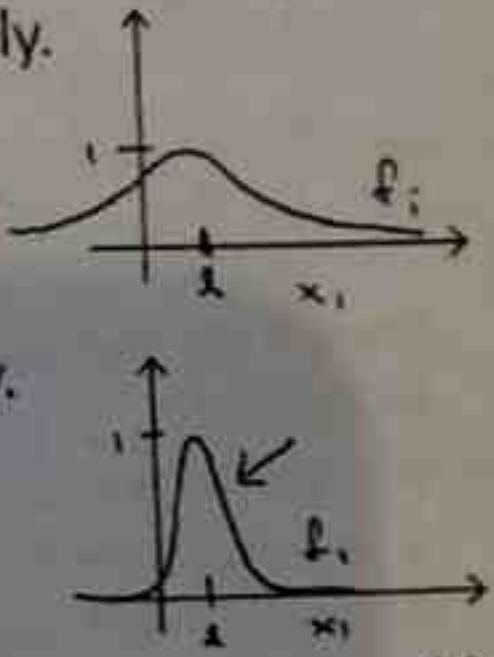→ Small C: Higher bias, low variance. (large $\lambda$)

$\sigma^2$ Large $\sigma^2$: Features $f_i$ vary more smoothly.
→ Higher bias, lower variance.

$$\exp\left(-\frac{\|x - l^{(i)}\|^2}{2\sigma^2}\right)$$ 2大时, ⇒

容易过拟合 不是易过拟合

Small $\sigma^2$: Features $f_i$ vary less smoothly.
Lower bias, higher variance.

容易过拟合 不容易过拟合.

---

**Support Vector Machines**

**Using an SVM**

Machine Learning

---

Use SVM software package (e.g. liblinear, libsvm, ...) to solve for parameters $\theta$. 那怎样用以, 也很麦

Need to specify:
→ Choice of parameter C.
Choice of kernel (similarity function):

E.g. No kernel ("linear kernel")
Predict "y = 1" if $\theta^T x \geq 0$

$\theta_0 + \theta_1 x_1 + \cdots + \theta_n x_n \geq 0$    $x \in \mathbb{R}^{n+1}$

→ $n$ large, $m$ small

→选择1: 不用 kernel, 适合线性的老表
$x \in \mathbb{R}^n$, $n$ small 当 feature很多时
$\exp \frac{?}{}$ (m小)

→ Gaussian kernel:
$$f_i = \exp\left(-\frac{\|x - l^{(i)}\|^2}{2\sigma^2}\right), \text{ where } l^{(i)} = x^{(i)}.$$
call $n$ large

Need to choose $\frac{\sigma^2}{?}$

featurely 样本多样,
用 kernel 实现一个复杂的非线性模型.

---

有时要自己实现一个 kernel

**Kernel (similarity) functions:** $x^{(i)} \quad l^{(j)} = x^{(j)}$

```
function f = kernel(x1, x2)
    f = exp( -||x1 - x2||² / (2σ²) )
return
```
$f_i$
$x \to$ $f_1$, $f_2$, $\vdots$, $f_m$

倒送给到 kernel

→ Note: Do perform feature scaling before using the Gaussian kernel.
$x \in \mathbb{R}^n$

$\|x - l\|^2$    $v = x - l$

$\|v\|^2 = v_1^2 + v_2^2 + \cdots + v_n^2$
$= (x_1 - l_1)^2 + (x_2 - l_2)^2 + \cdots + (x_n - l_n)^2$
$\underbrace{}_{1000 \text{ feet}^2} \quad \underline{1\text{-}5 \text{ bedrooms}}$

feature scaling 看看啊有几个

要归一到同一个区间处

## Other choices of kernel

Linear kernel, meaning no kernel, 高斯 kernel

Note: Not all similarity functions $\text{similarity}(x, l)$ make valid kernels.
→ (Need to satisfy technical condition called "Mercer's Theorem" to make sure SVM packages' optimizations run correctly, and do not diverge).
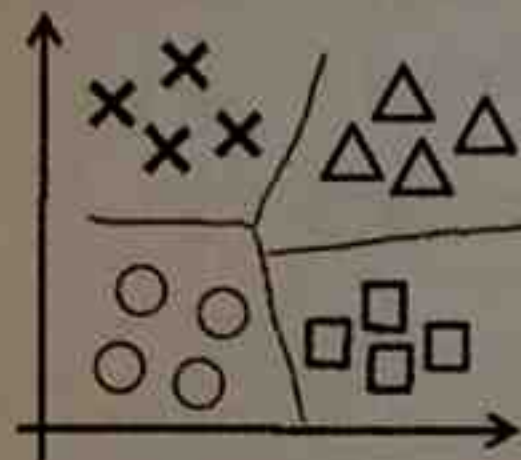
孙数学证明⽤的太多了吧，万能转换本顺话

Many off-the-shelf kernels available:
  Polynomial kernel: $k(x, l) = (x^T l)^2$ 越⽤越超过 $(2次⽅与⼆次⽅)$ $(x^T l + \text{const})^{\text{degree}}$

  直接⽐何取kernel $(x^T l)^3$, $(x^T l + 1)^3$, $(x^T l + 5)$ 其他各种

  $x$与$l$都取negative 基本各种 有⽆穷⽆尽⍟

  - More esoteric: String kernel, chi-square kernel, histogram intersection kernel, ... $\sin(x, l)$ 你⾃⼰再造⼀个

## Multi-class classification



$y \in \{1, 2, 3, \ldots, K\}$

Many SVM packages already have built-in multi-class classification functionality.
→ Otherwise, use one-vs.-all method. (Train $K$ SVMs, one to distinguish $y = i$ from the rest, for $i = 1, 2, \ldots, K$), get $\theta^{(1)}, \theta^{(2)}, \ldots, \theta^{(K)}$
Pick class $i$ with largest $(\theta^{(i)})^T x$
  $y=1$   $y=2$   $\ldots$   $\theta_i K$

⽤ $K$ 个 model 各预测⼀下
分到谁分最⾼的那个类别中

## Logistic regression vs. SVMs

$n$ = number of features ($x \in \mathbb{R}^{n+1}$), $m$ = number of training examples
→ If $n$ is large (relative to $m$): (e.g. $n \geq m$, $n = 10,000$, $m = 10 \mp 1000$)
→ Use logistic regression, or SVM without a kernel ("linear kernel")

特征够多了，⽤ LR 我可以

→ If $n$ is small, $m$ is intermediate:  ($n = 1 - 1000$, $m = 10 - 10,000$) ←
  → Use SVM with Gaussian kernel

特征少，样本也量也不是太恐怖

If $n$ is small, $m$ is large:  ($n = 1-1000$, $m = 50,000+$)
  → Create/add more features, then use logistic regression or SVM without a kernel

样本量太⼤⼤分布了，还是加特征 或⽤ LR 将将凑凑的凑接

→ Neural network likely to work well for most of these settings, but may be slower to train.

神经⽹络可能都可以⼀些
只是训练化较慢

LR 与 SVM-NO-KERNEL 很相似

SVM : convex optimization  ×worry for local optimisation problem
Neural Network : ⼀般也不⽤担⼼ local optimisation problem
但还是化较慢