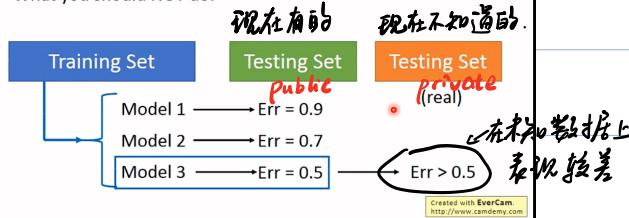


Model Selection

Model Selection

- There is usually a trade-off between bias and variance.
- Select a model that balances two kinds of error to minimize total error
- What you should NOT do:



对于在真实数据 (Private Testing Set) 上反应该结果较

差的方法：在 Training Set 中拿出 Validation set，相当

于把从现有数据 (public Testing Set) 中得到最好已知

这一步提前到 Validation set 来做。然后把该函数

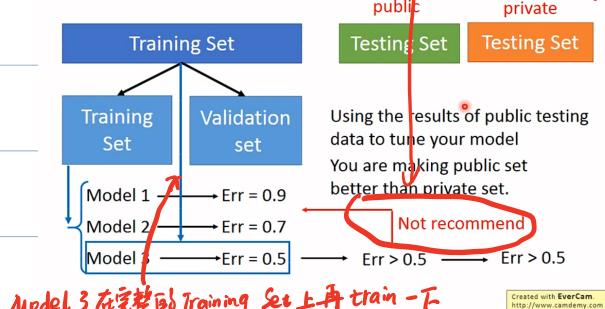
在 public set 上调优，从而在 public set 上得到

的 Err 可以和在 Private Set 上得到的 Err 更接近更

真实（相当于把所有的步骤提前一步，使 public 变 private）

Cross Validation

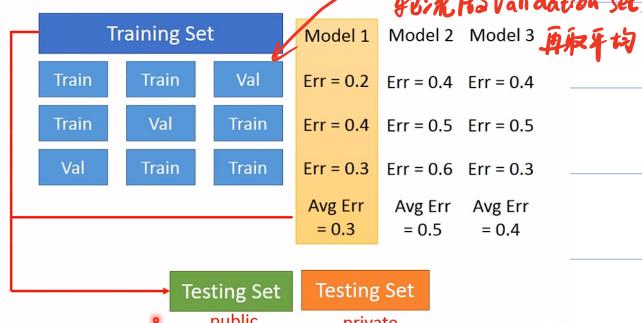
不推荐这样用，否则推前这步就失去了意义。



担心选出来的 Validation set 中没有不同形

的数据 how 解决？

N-fold Cross Validation



Gradient Descent

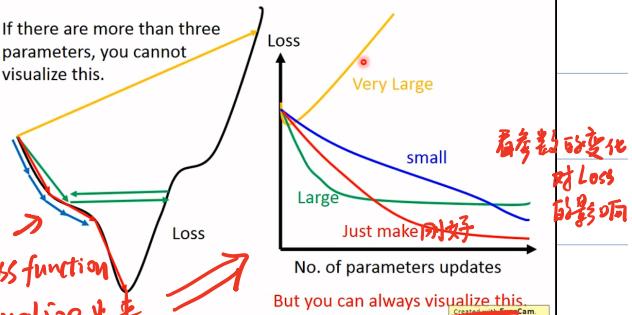
Review...

注意事项: Tip 1: 调节 learning rate .

Learning Rate

Set the learning rate η carefully

If there are more than three parameters, you cannot visualize this.



But you can always visualize this.

调 Learning rates

⇒ Adaptive Learning Rates. (自动调节)

Adaptive Learning Rates

- Popular & Simple Idea: Reduce the learning rate by some factor every few epochs.
- At the beginning, we are far from the destination, so we use larger learning rate
- After several epochs, we are close to the destination, so we reduce the learning rate
- E.g. $1/t$ decay: $\eta^t = \eta / \sqrt{t+1}$ (加速迭代次数, 快速收敛)
- Learning rate cannot be one-size-fits-all
 - Giving different parameters different learning rates (针对参数的不同)

Adagrad

$$\eta^t = \frac{\eta}{\sqrt{t+1}} \quad g^t = \frac{\partial C(\theta^t)}{\partial w}$$

- Divide the learning rate of each parameter by the root mean square of its previous derivatives

Vanilla Gradient descent

$$w^{t+1} \leftarrow w^t - \eta^t g^t \quad w \text{ is one parameters}$$

Adagrad 对其某一个参数

$$w^{t+1} \leftarrow w^t - \frac{\eta^t}{\sigma^t} g^t \quad \sigma^t: \text{root mean square of the previous derivatives of parameter } w$$

与 w 有关

方向根

$\sigma^t: \text{root mean square of the previous derivatives of parameter } w$

Adagrad

$$w^1 \leftarrow w^0 - \frac{\eta^0}{\sigma^0} g^0 \quad \sigma^0 = \sqrt{(g^0)^2}$$

$$w^2 \leftarrow w^1 - \frac{\eta^1}{\sigma^1} g^1 \quad \sigma^1 = \sqrt{\frac{1}{2}[(g^0)^2 + (g^1)^2]}$$

$$w^3 \leftarrow w^2 - \frac{\eta^2}{\sigma^2} g^2 \quad \sigma^2 = \sqrt{\frac{1}{3}[(g^0)^2 + (g^1)^2 + (g^2)^2]}$$

$$\vdots$$

$$w^{t+1} \leftarrow w^t - \frac{\eta^t}{\sigma^t} g^t \quad \sigma^t = \sqrt{\frac{1}{t+1} \sum_{i=0}^t (g^i)^2}$$

Adagrad

- Divide the learning rate of each parameter by the **root mean square of its previous derivatives**

$$w^{t+1} \leftarrow w^t - \frac{\eta}{\sqrt{\sum_{i=0}^t (g^i)^2}} g^t$$

方根，不用求平均

Created with EverCam
http://www.camdem.com

Why $\begin{cases} w^{t+1} \leftarrow w^t - \eta g^t, g^t 越大步伐越大 \\ w^{t+1} \leftarrow w^t - \frac{\eta}{\sqrt{\sum_{i=0}^t (g^i)^2}} g^t, g^t 越大步伐越小? \end{cases}$

解釋 1：

Intuitive Reason $\eta^t = \frac{\eta}{\sqrt{t+1}}$ $g^t = \frac{\partial C(\theta^t)}{\partial w}$

- How surprised it is 反差

g^0	g^1	g^2	g^3	g^4
0.001	0.001	0.003	0.002	0.1
g^0	g^1	g^2	g^3	g^4
10.8	20.9	31.7	12.1	0.1

特別大
特別小

$$w^{t+1} \leftarrow w^t - \frac{\eta}{\sqrt{\sum_{i=0}^t (g^i)^2}} g^t$$

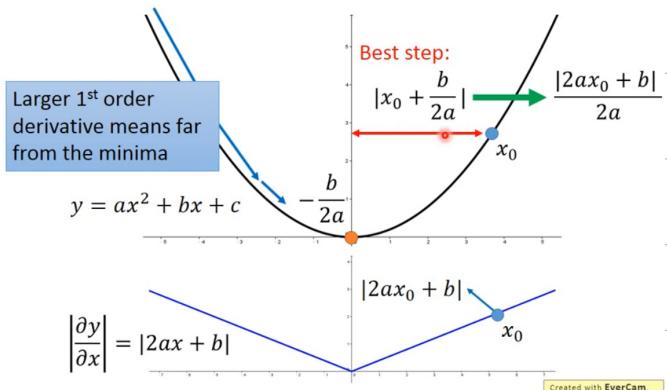
造成反差的效果
和过去相比反差多大，从而判定其相对大小。

解釋 2：

若对单一的参数而言，用 $w^{t+1} \leftarrow w^t - \eta g^t$ 确实合理。

因为一般 gradient 越大离最低点越近，是绝对值
但对两个参数，

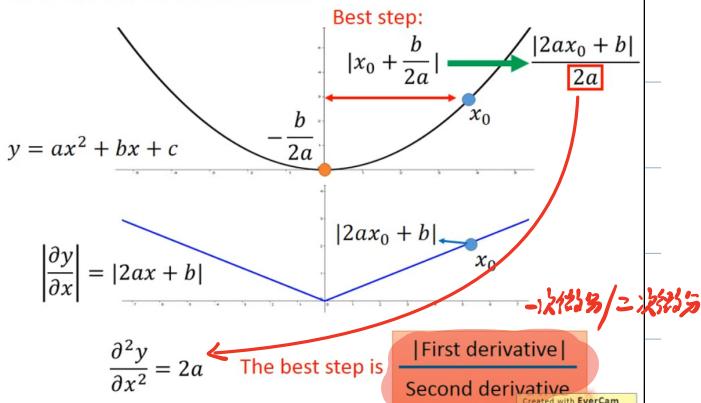
Larger gradient, larger steps?



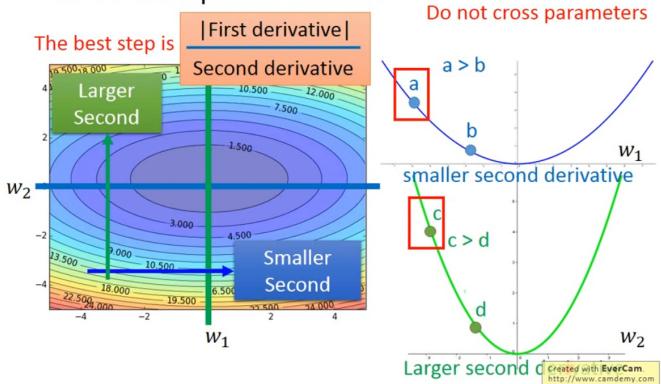
不-是 gradient 越大的离最低点越近，是参数比较
的 $a \neq 1$ 。梯度法 a.t. 但实际梯度是向
量更远，所以参数时 gradient 是一个相对值。

找最好的 step

Second Derivative



Comparison between different parameters



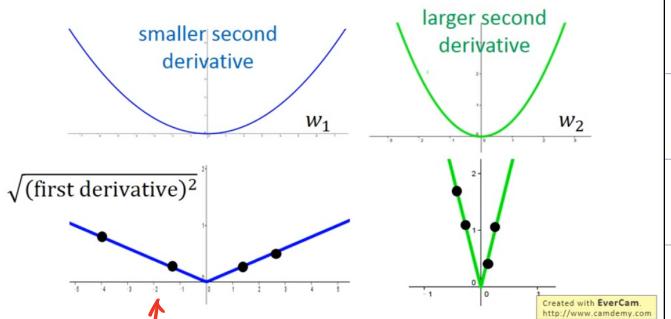
这种 $\frac{|First\ derivative|}{Second\ derivative}$ 与 Adagrad 关系？

$$w^{t+1} \leftarrow w^t - \frac{\eta}{\sqrt{\sum_{i=0}^t (g^i)^2}} g^t$$

The best step is $\frac{|First\ derivative|}{Second\ derivative}$

?

Use first derivative to estimate second derivative



通过取多个点来反映二次微分的大小

Tip 2: Stochastic Gradient Descent

上 Training 章

Stochastic Gradient Descent

$$L = \sum_n \left(\hat{y}^n - \left(b + \sum w_i x_i^n \right) \right)^2$$

Loss is the summation over all training examples

◆ **Gradient Descent** $\theta^i = \theta^{i-1} - \eta \nabla L(\theta^{i-1})$

◆ **Stochastic Gradient Descent**

Pick an example x^n 随机选一个算 Loss.

$$L^n = \left(\hat{y}^n - \left(b + \sum w_i x_i^n \right) \right)^2$$

$\theta^i = \theta^{i-1} - \eta \nabla L^n(\theta^{i-1})$
Loss for only one example

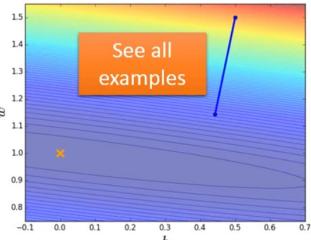
Created with EverCam
<http://www.camdemmy.com>

但用 Stochastic Gradient Descent 若有 m 个参数，就要将参数 θ update m 次

Stochastic Gradient Descent

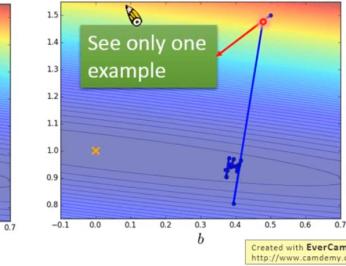
Gradient Descent

Update after seeing all examples



Stochastic Gradient Descent

Update for each example
If there are 20 examples,
20 times faster.

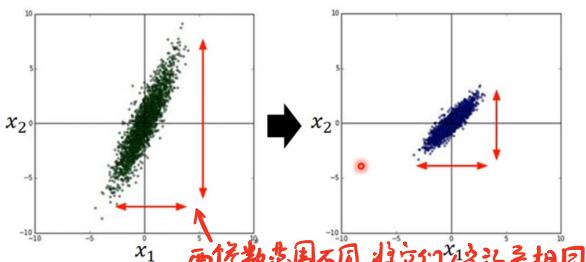


Tip 3: Feature Scaling. 为什么放缩

Feature Scaling

Source of figure:
<http://cs231n.github.io/neural-networks-2/>

$$y = b + w_1 x_1 + w_2 x_2$$



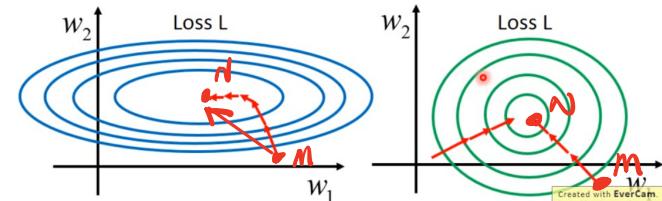
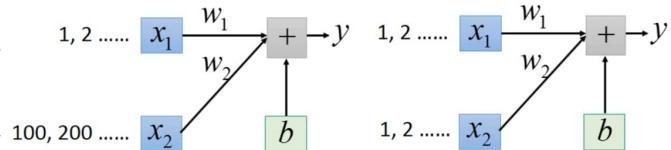
Make different features have the same scaling

Created with EverCam
<http://www.camdemmy.com>

Why 这样做？

Feature Scaling

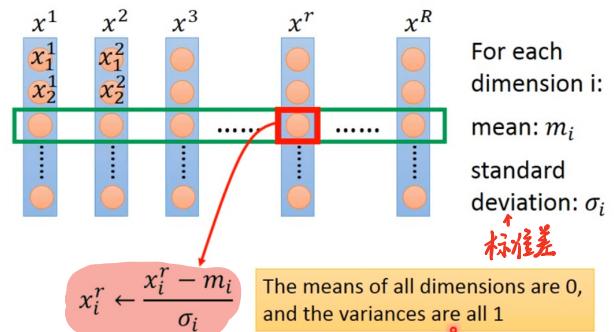
$$y = b + w_1 x_1 + w_2 x_2$$



若不用 Adagrad, w_1 和 w_2 和 Loss 的影响很不一样。若 Loss 曲线像左边一样为椭圆，则从 $M \rightarrow N$ 是沿著高成直，而非从最高点指向最近的点。但 rescaling 成圆形后就会指向 N 走

How to Feature Scaling.

Feature Scaling



Created with EverCam
<http://www.camdemmy.com>

关于 Gradient Descent 的数学原理.

用对多变量函数的泰勒展开，在 (a, b) 处

$$L(\theta) = L(a, b) + \frac{\partial L(a, b)}{\partial \theta} (\theta_1 - a) + \frac{\partial L(a, b)}{\partial \theta_2} (\theta_2 - b).$$

↑ ↑ ↑
S u v

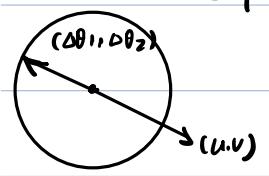
考虑在 a, b 附近一个小圆范围内 $L(\theta)$ 的最小值。

L19) 同时可看成 (u, v) 与 $(\theta_1 - a, \theta_2 - b) = (\theta_1, \theta_2)$

如何分类?

在 $\Delta\theta_1^2 + \Delta\theta_2^2 \leq d^2$ 范围内的最小值. 显然相反且

$$\Delta\theta_1^2 + \Delta\theta_2^2 = d^2 \text{ 时 } \frac{\partial}{\partial u} \ll 1, \quad \nabla(\Delta\theta_1, \Delta\theta_2) \text{ 表示为 } \begin{bmatrix} \Delta\theta_1 \\ \Delta\theta_2 \end{bmatrix} = -\eta \begin{bmatrix} u \\ v \end{bmatrix}$$



$$\text{BP } \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} = \begin{bmatrix} a \\ b \end{bmatrix} - \eta \begin{bmatrix} u \\ v \end{bmatrix}$$
$$\text{gradient descent } \rightarrow \begin{bmatrix} a \\ b \end{bmatrix} - \eta \begin{bmatrix} \frac{\partial L(a, b)}{\partial \theta_1} \\ \frac{\partial L(a, b)}{\partial \theta_2} \end{bmatrix}$$

但让该值精确的前提是泰勒展开是精确的. 就要求画的图足够小. d 够小: 理论上 Learning rate 要足够小.

泰勒也可展开到更高阶. 此时 Learning rate 可以大一些.
但要求更多的点会更麻烦.

Classification

Example Application

POKÉMON TYPE SYMBOLS



$$f(\text{Pikachu}) = \text{Electric} \quad f(\text{Squirtle}) = \text{Water} \quad f(\text{Bulbasaur}) = \text{Grass}$$

要把宝可梦作输入, 将其数值化

将其各个特征数值化后组成向量, 一向量代表一只

pokemon games (NOT pokemon cards or Pokemon Go)



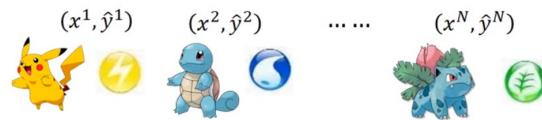
Example Application

- Total: sum of all stats that come after this, a general guide to how strong a pokémon is **320**
- HP: hit points, or health, defines how much damage a pokémon can withstand before fainting **35**
- Attack: the base modifier for normal attacks (e.g. Scratch, Punch) **55**
- Defense: the base damage resistance against normal attacks **40**
- SP Atk: special attack, the base modifier for special attacks (e.g. fire blast, bubble beam) **50**
- SP Def: the base damage resistance against special attacks **50**
- Speed: determines which pokémon attacks first each round **90**

Can we predict the "type" of pokémon based on the information?

How to do Classification

• Training data for Classification



Classification as Regression? **当成 Regression 处理?**

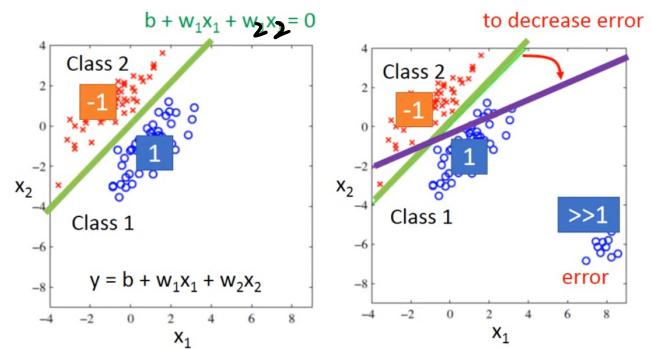
Binary classification as example

Training: Class 1 means the target is 1; Class 2 means the target is -1

Testing: closer to 1 → class 1; closer to -1 → class 2

Created with EverCam:
<http://www.camdem.com>

11 Regression 遇到的问题



Penalize to the examples that are "too correct" ... (Bishop, P186)

- Multiple class: Class 1 means the target is 1; Class 2 means the target is 2; Class 3 means the target is 3 problematic

Created with EverCam:
<http://www.camdem.com>

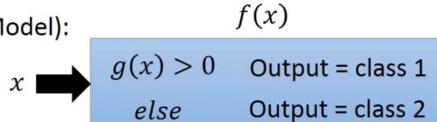
问题1: 对 Classification 而言, 误判作为惩罚是最好的. 但对 Regression 来说, 由于有些"过正"的值 (>1), 而为使 loss 更小, 会选择紫色线, 即两者对于 model 好坏的定义不相同.

问题2: 若有多个类别, Regression 会认为相邻的两个 class 之间有某种关系, 但可在实际上这两个 class 间并无关系.

⇒ 七个分量
的向量

Ideal Alternatives

- Function (Model):



- Loss function:

$$L(f) = \sum_n \delta(f(x^n) \neq \hat{y}^n)$$

The number of times f get incorrect results on training data.

- Find the best function:

- Example: Perceptron, SVM Not Today

以上这种办法，在 $f(x)$ 中找一个 $g(x)$ ，用 $g(x) > 0$ 来判

定 Binary Classification。Loss 的意思是将 x^n 代入

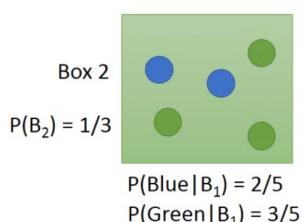
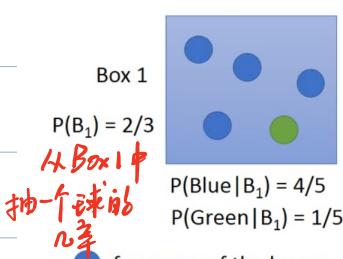
$f(x)$ 与正确答案作比较，不一样返回1，一样返回

0，但现在无法通过微分等方法求 f 。

另-Solution.

从概率问题切入

Two Boxes



from one of the boxes

Where does it come from?

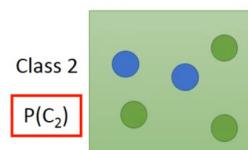
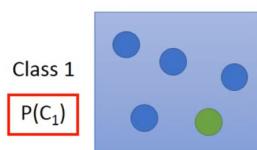
$$P(B_1 | \text{Blue}) = \frac{P(\text{Blue}|B_1)P(B_1)}{P(\text{Blue}|B_1)P(B_1) + P(\text{Blue}|B_2)P(B_2)}$$

Created with EverCam
http://www.camdemmy.com

⇒ 把盒子换成类别。

Two Classes

Estimating the Probabilities From training data



Given an x , which class does it belong to

$$P(C_1|x) = \frac{P(x|C_1)P(C_1)}{P(x|C_1)P(C_1) + P(x|C_2)P(C_2)}$$

是属于 Class 1 的概率是 $P(x|C_1)P(C_1)$ 不是属于 Class 2 的概率是 $P(x|C_2)P(C_2)$

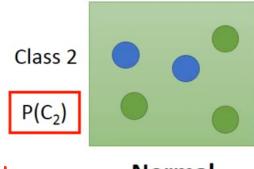
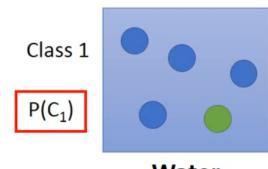
Generative Model $P(x) = P(x|C_1)P(C_1) + P(x|C_2)P(C_2)$

$P(C_1), P(C_2), P(x|C_1), P(x|C_2)$ 要通过 Training

Data 算出来

Prior

$P(C_1), P(C_2)$ of Prior



Water

Normal

Binary 情况

Water and Normal type with ID < 400 for training, rest for testing

Training: 79 Water, 61 Normal

$$P(C_1) = 79 / (79 + 61) = 0.56$$

$$P(C_2) = 61 / (79 + 61) = 0.44$$

接下来考虑某只宝可梦在某类型中被取出的概率

Probability from Class

$$P(x|C_1) = ? \quad P(\text{Squirtle} | \text{Water}) = ?$$

Each Pokémon is represented as a vector by its attribute.



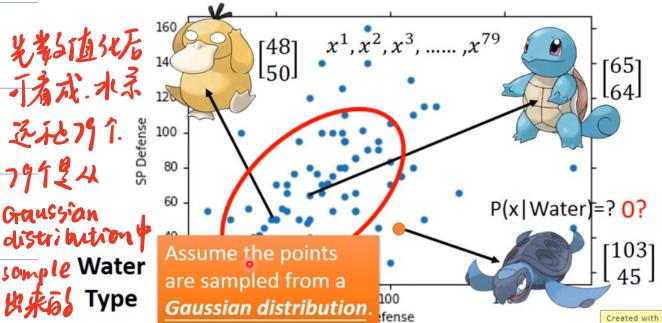
图中海龟并不能在79个里找到，但确实是水系。所

以也不能说在水系中遇到它概率为 $0/79$

⇒ 每个宝可梦都是用 features 的向量表示。

Probability from Class - Feature

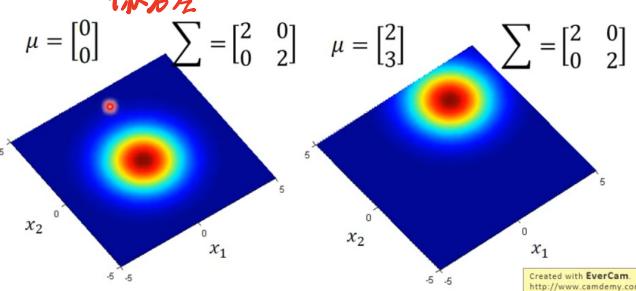
- Considering Defense and SP Defense



Gaussian Distribution (高斯分布. 正态分布) <https://blog.slinuxer.com/tag/pca>

$$f_{\mu, \Sigma}(x) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}$$

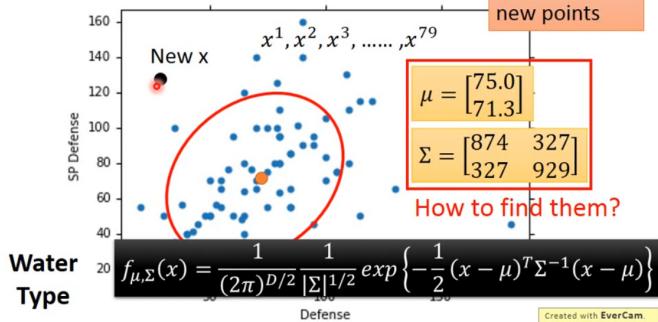
Input: vector x , output: probability of sampling x
The shape of the function determines by mean μ and covariance matrix Σ 代表高斯分布的中心和程度



Probability from Class

Assume the points are sampled from a Gaussian distribution

Find the Gaussian distribution behind them \rightarrow Probability for new points



算出 μ, Σ 后代入 $f_{\mu, \Sigma}(x)$ 即可得 x 在正态分布中概率

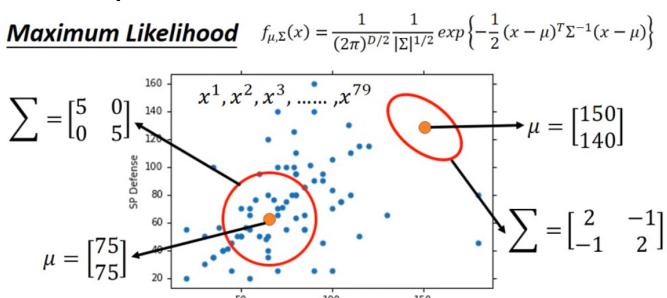
Sample出来的概率？

How找 μ, Σ ? \Rightarrow Maximum Likelihood.

每个 Gaussian distribution 都可由 sample 选出 79 个点。

但每个分布可以 sample 出 79 个不同的概率

率是不一样的。



The Gaussian with any mean μ and covariance matrix Σ can generate these points.

\rightarrow Different Likelihood

Likelihood of a Gaussian with mean μ and covariance matrix Σ

= the probability of the Gaussian samples $x^1, x^2, x^3, \dots, x^{79}$

$$L(\mu, \Sigma) = f_{\mu, \Sigma}(x^1) f_{\mu, \Sigma}(x^2) f_{\mu, \Sigma}(x^3) \dots f_{\mu, \Sigma}(x^{79})$$

sample 出 79 点的可能性大小，等于单独选出每点几率之积。

Maximum Likelihood

We have the "Water" type Pokémons: $x^1, x^2, x^3, \dots, x^{79}$

We assume $x^1, x^2, x^3, \dots, x^{79}$ generate from the Gaussian (μ^*, Σ^*) with the maximum likelihood

$$L(\mu, \Sigma) = f_{\mu, \Sigma}(x^1) f_{\mu, \Sigma}(x^2) f_{\mu, \Sigma}(x^3) \dots f_{\mu, \Sigma}(x^{79})$$

$$f_{\mu, \Sigma}(x) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}$$

$$\mu^*, \Sigma^* = \arg \max_{\mu, \Sigma} L(\mu, \Sigma)$$

$$\mu^* = \frac{1}{79} \sum_{n=1}^{79} x^n$$

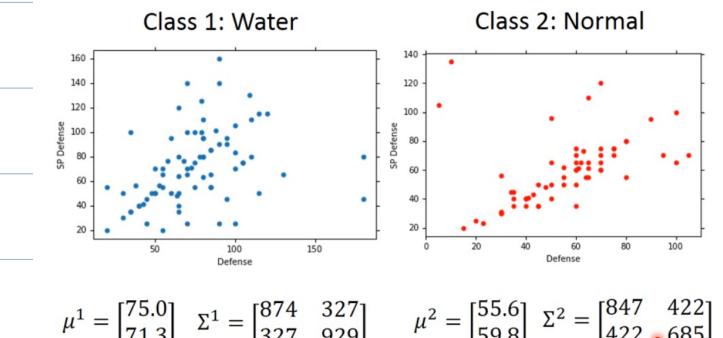
$$\Sigma^* = \frac{1}{79} \sum_{n=1}^{79} (x^n - \mu^*) (x^n - \mu^*)^T$$

算出时均值平均值。

Created with EverCam
<http://www.camdemmy.com>

回到实际问题

Maximum Likelihood



$$\mu^1 = \begin{bmatrix} 75.0 \\ 71.3 \end{bmatrix} \quad \Sigma^1 = \begin{bmatrix} 874 & 327 \\ 327 & 929 \end{bmatrix} \quad \mu^2 = \begin{bmatrix} 55.6 \\ 59.8 \end{bmatrix} \quad \Sigma^2 = \begin{bmatrix} 847 & 422 \\ 422 & 685 \end{bmatrix}$$

Created with EverCam
<http://www.camdemmy.com>

Now we can do classification 😊

$$f_{\mu^1, \Sigma^1}(x) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma^1|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu^1)^T (\Sigma^1)^{-1} (x - \mu^1) \right\}$$

$$\mu^1 = \begin{bmatrix} 75.0 \\ 71.3 \end{bmatrix} \quad \Sigma^1 = \begin{bmatrix} 874 & 327 \\ 327 & 929 \end{bmatrix}$$

$$P(C_1|x) = \frac{P(x|C_1)P(C_1)}{P(x|C_1)P(C_1) + P(x|C_2)P(C_2)}$$

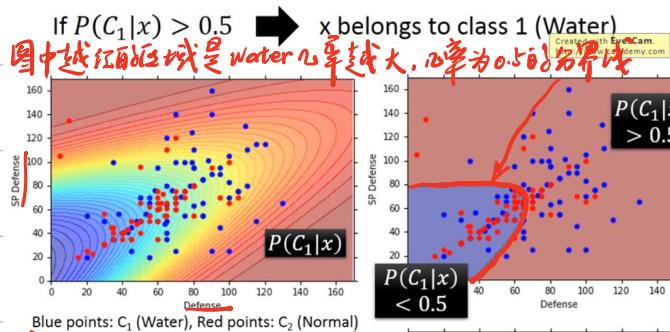
$$f_{\mu^2, \Sigma^2}(x) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma^2|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu^2)^T (\Sigma^2)^{-1} (x - \mu^2) \right\}$$

$$\mu^2 = \begin{bmatrix} 55.6 \\ 59.8 \end{bmatrix} \quad \Sigma^2 = \begin{bmatrix} 847 & 422 \\ 422 & 685 \end{bmatrix}$$

$$P(C_2|x) = \frac{P(x|C_2)P(C_2)}{P(x|C_1)P(C_1) + P(x|C_2)P(C_2)}$$

P(C1) = 79 / (79 + 61) = 0.56

P(C2) = 61 / (79 + 61) = 0.44



Blue points: C1 (Water), Red points: C2 (Normal)

P(C1|x) > 0.5

P(C1|x) < 0.5

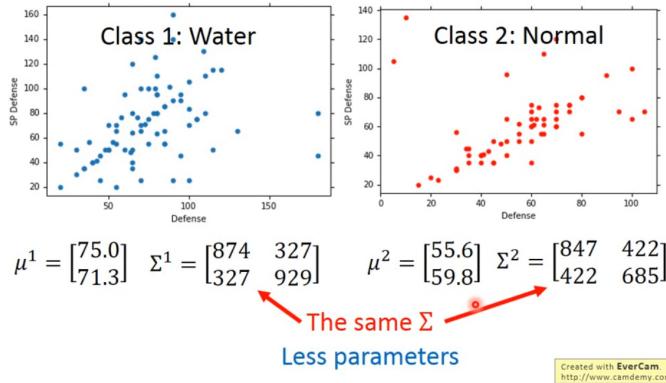
P(C1|x) > 0.5</

发现就算用了7个features结果还是不好.

⇒修改模型

上面的并不常用模型.

Modifying Model



上面不同的 Gaussian distribution 用了不同的 Σ , 这会造成参数过多. 从而容易 over fitting.

⇒强行令 Σ 为同一个.

Modifying Model

Ref: Bishop, chapter 4.2.2

• Maximum likelihood

"Water" type Pokémons:

$$x^1, x^2, x^3, \dots, x^{79}$$

"Normal" type Pokémons:

$$x^{80}, x^{81}, x^{82}, \dots, x^{140}$$

$$\mu^1 \quad \Sigma \quad \mu^2$$

Find μ^1, μ^2, Σ maximizing the likelihood $L(\mu^1, \mu^2, \Sigma)$

$$L(\mu^1, \mu^2, \Sigma) = f_{\mu^1, \Sigma}(x^1) f_{\mu^1, \Sigma}(x^2) \cdots f_{\mu^1, \Sigma}(x^{79}) \\ \times f_{\mu^2, \Sigma}(x^{80}) f_{\mu^2, \Sigma}(x^{81}) \cdots f_{\mu^2, \Sigma}(x^{140})$$

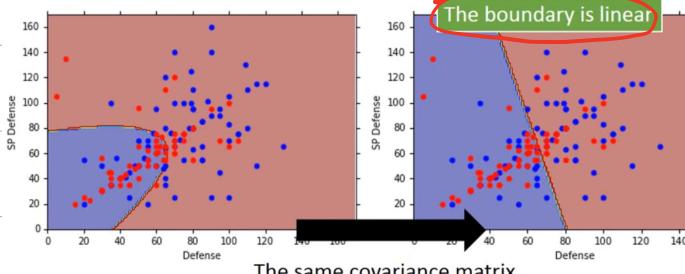
μ^1 and μ^2 is the same
但是各自的 mean.

$$\Sigma = \frac{79}{140} \Sigma^1 + \frac{61}{140} \Sigma^2$$

改造的结果

Modifying Model

同一个训练的分类器
变成线性的.



All: total, hp, att, sp att, de, sp de, speed

54% accuracy → 73% accuracy

贝叶斯模型

也是 ML 的 3 个 steps.

Three Steps

• Function Set (Model):

$$P(C_1|x) = \frac{P(x|C_1)P(C_1)}{P(x|C_1)P(C_1) + P(x|C_2)P(C_2)}$$

$x \rightarrow$

If $P(C_1|x) > 0.5$, output: class 1
Otherwise, output: class 2

• Goodness of a function:

- The mean μ and covariance Σ that maximizing the likelihood (the probability of generating data)

• Find the best function: easy

Created with EverCam.
http://www.camdemyc.com

⇒ Why 用 Gaussian distribution?

用其他 distribution 也可以.

简单、参数少: Bias k, Variencenk

复杂、参数多: Bias n, Variencenk

Probability Distribution

• You can always use the distribution you like ☺

正常情况下认为各参数间有关系, 关系反映在 Σ 中.

$$P(x|C_1) = P(x_1|C_1) P(x_2|C_1) \cdots P(x_k|C_1) \cdots$$

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_k \\ \vdots \\ x_K \end{bmatrix}$$

1-D Gaussian

For binary features, you may assume they are from Bernoulli distributions.

即假设各参数间无关. 此时多维 Gaussian

If you assume all the dimensions are independent, then you are using Naive Bayes Classifier.

矩阵相乘, 2D
变成对角阵

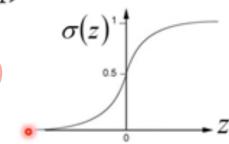
Posterior Probability

$$P(C_1|x) = \frac{P(x|C_1)P(C_1)}{P(x|C_1)P(C_1) + P(x|C_2)P(C_2)}$$

$$= \frac{1}{1 + \frac{P(x|C_2)P(C_2)}{P(x|C_1)P(C_1)}} = \frac{1}{1 + \exp(-z)} = \sigma(z)$$

Sigmoid function

$$z = \ln \frac{P(x|C_1)P(C_1)}{P(x|C_2)P(C_2)}$$



Created with EverCam.
http://www.camdemyc.com

长得什么样？

经变形后：

$$z = \ln \frac{|\Sigma^2|^{1/2}}{|\Sigma^1|^{1/2}} - \frac{1}{2} x^T (\Sigma^1)^{-1} x + (\mu^1)^T (\Sigma^1)^{-1} x - \frac{1}{2} (\mu^1)^T (\Sigma^1)^{-1} \mu^1 \\ + \frac{1}{2} x^T (\Sigma^2)^{-1} x - (\mu^2)^T (\Sigma^2)^{-1} x + \frac{1}{2} (\mu^2)^T (\Sigma^2)^{-1} \mu^2 + \ln \frac{N_1}{N_2}$$

Created with EverCam
http://www.camdemmy.com

以上是当 $\Sigma_1 \neq \Sigma_2$ 时，若 $\Sigma_1 = \Sigma_2$ 可化简。

$P(C_1|x) = \sigma(z)$

$$z = \ln \frac{|\Sigma^2|^{1/2}}{|\Sigma^1|^{1/2}} - \frac{1}{2} x^T (\Sigma^1)^{-1} x + (\mu^1)^T (\Sigma^1)^{-1} x - \frac{1}{2} (\mu^1)^T (\Sigma^1)^{-1} \mu^1 \\ + \frac{1}{2} x^T (\Sigma^2)^{-1} x - (\mu^2)^T (\Sigma^2)^{-1} x + \frac{1}{2} (\mu^2)^T (\Sigma^2)^{-1} \mu^2 + \ln \frac{N_1}{N_2}$$

$\Sigma_1 = \Sigma_2 = \Sigma$

$$z = (\mu^1 - \mu^2)^T \Sigma^{-1} x - \frac{1}{2} (\mu^1)^T (\Sigma^1)^{-1} \mu^1 + \frac{1}{2} (\mu^2)^T (\Sigma^2)^{-1} \mu^2 + \ln \frac{N_1}{N_2}$$

w^T b

$P(C_1|x) = \sigma(w \cdot x + b)$

一组 w 和 b 是对 C 而言，
而另外 w 和 b 是对 C_1

Created with EverCam
http://www.camdemmy.com

Logistic Regression

Step 1: Function Set

We want to find $P_{w,b}(C_1|x)$

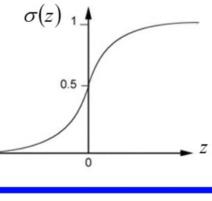
$$\sigma(z) = \frac{1}{1 + \exp(-z)}$$

If $P_{w,b}(C_1|x) \geq 0.5$, output C_1

Otherwise, output C_2

$P_{w,b}(C_1|x) = \sigma(z)$

$$z = w \cdot x + b = \sum_i w_i x_i + b$$



Function set:

$f_{w,b}(x) = P_{w,b}(C_1|x)$

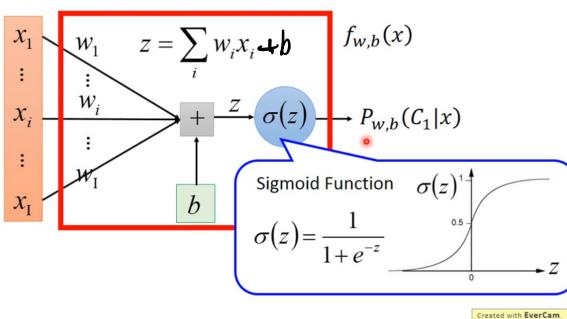
Including all different w and b

Created with EverCam
http://www.camdemmy.com

回顾上节课, function set 是与 w 和 b 有关的。

将模型图形化。

Step 1: Function Set



上节课是通过找最好的 w 和 b 来找到最好的

function, 但是由于 ~~最好的~~ function 可以由 w 和 b 来表示, 所以可以直接找最好的 w 和 b .

Step 2: Goodness of a Function

Training Data

	x^1	x^2	x^3	x^N
	C_1	C_1	C_2	C_1

Assume the data is generated based on $f_{w,b}(x) = P_{w,b}(C_1|x)$

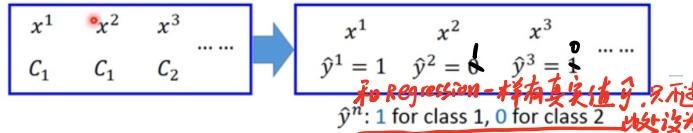
Given a set of w and b , what is its probability of generating the data? 判断一组 w 和 b 可以产生这组 Training Data 的概率

$$L(w, b) = f_{w,b}(x^1) f_{w,b}(x^2) (1 - f_{w,b}(x^3)) \cdots f_{w,b}(x^N)$$

The most likely w^* and b^* is the one with the largest $L(w, b)$.

$$w^*, b^* = \arg \max_{w, b} L(w, b)$$

Created with EverCam
http://www.camdemmy.com



$$L(w, b) = f_{w,b}(x^1) f_{w,b}(x^2) (1 - f_{w,b}(x^3)) \cdots$$

$$w^*, b^* = \arg \max_{w, b} L(w, b) = w^*, b^* = \arg \min_{w, b} -\ln L(w, b)$$

由于用了 \ln , 可拆成加法.

$$-\ln L(w, b) = -\ln f_{w,b}(x^1) \rightarrow -[\hat{y}^1 \ln f(x^1) + (1 - \hat{y}^1) \ln(1 - f(x^1))] \\ -\ln f_{w,b}(x^2) \rightarrow -[\hat{y}^2 \ln f(x^2) + (1 - \hat{y}^2) \ln(1 - f(x^2))] \\ -\ln(1 - f_{w,b}(x^3)) \rightarrow -[\hat{y}^3 \ln f(x^3) + (1 - \hat{y}^3) \ln(1 - f(x^3))] \\ \vdots$$

数学转化

Created with EverCam
http://www.camdemmy.com

⇒ 通过上述数学转化 将要 minimize 的对象写成一个 function:

Step 2: Goodness of a Function

$$L(w, b) = f_{w,b}(x^1) f_{w,b}(x^2) (1 - f_{w,b}(x^3)) \cdots f_{w,b}(x^N)$$

$$-\ln L(w, b) = \ln f_{w,b}(x^1) + \ln f_{w,b}(x^2) + \ln(1 - f_{w,b}(x^3)) \cdots$$

\hat{y}^n : 1 for class 1, 0 for class 2

$$= \sum_n - [\hat{y}^n \ln f_{w,b}(x^n) + (1 - \hat{y}^n) \ln(1 - f_{w,b}(x^n))]$$

Cross entropy between two Bernoulli distribution

交叉熵

$$\text{Distribution p: } p(x=1) = \hat{y}^n \quad p(x=0) = 1 - \hat{y}^n$$

$$\text{Distribution q: } q(x=1) = f(x^n) \quad q(x=0) = 1 - f(x^n)$$

$$H(p, q) = - \sum_x p(x) \ln(q(x))$$

Created with EverCam
http://www.camdemmy.com

算 distribution p 和 q 的 cross entropy, 代表的是算出它们之间多接近

Logistic Regression 将 Linear model 转化成了 sigmoid 函数
Linear Regression

$$\text{Step 1: } f_{w,b}(x) = \sigma\left(\sum_i w_i x_i + b\right)$$

Output: between 0 and 1

$$\text{Training data: } (x^n, \hat{y}^n)$$

$$\text{Step 2: } \hat{y}^n: 1 \text{ for class 1, 0 for class 2}$$

$$L(f) = \sum_n C(f(x^n), \hat{y}^n)$$

$$f_{w,b}(x) = \sum_i w_i x_i + b$$

Output: any value

Cross entropy:
 $C(f(x^n), \hat{y}^n) = -[\hat{y}^n \ln f(x^n) + (1 - \hat{y}^n) \ln(1 - f(x^n))]$
 题表是把 $f(x^n)$ 和 \hat{y}^n 看成两个伯努利分布，这两个分布越接近其结果越好

Created with EverCam
<http://www.camdemmy.com>

进行 Gradient descent .

Step 3: Find the best function

$$\frac{\partial \ln L(w, b)}{\partial w_i} = \sum_n -\left[\hat{y}^n \frac{\partial \ln f_{w,b}(x^n)}{\partial w_i} + (1 - \hat{y}^n) \frac{\partial \ln(1 - f_{w,b}(x^n))}{\partial w_i} \right]$$

$$\frac{\partial \ln f_{w,b}(x)}{\partial w_i} = \frac{\partial \ln f_{w,b}(x)}{\partial z} \frac{\partial z}{\partial w_i} \quad \frac{\partial z}{\partial w_i} = x_i$$

$$\frac{\partial \ln \sigma(z)}{\partial z} = \frac{1}{\sigma(z)} \frac{\partial \sigma(z)}{\partial z} = \frac{1}{\sigma(z)} \sigma(z)(1 - \sigma(z))$$

$$f_{w,b}(x) = \sigma(z) = \frac{1}{1 + \exp(-z)}$$

$$z = w \cdot x + b = \sum_i w_i x_i + b$$

Created with EverCam
<http://www.camdemmy.com>

Step 3: Find the best function

$$\frac{\partial \ln L(w, b)}{\partial w_i} = \sum_n -\left[\hat{y}^n \frac{\partial \ln f_{w,b}(x^n)}{\partial w_i} + (1 - \hat{y}^n) \frac{\partial \ln(1 - f_{w,b}(x^n))}{\partial w_i} \right]$$

$$\frac{\partial \ln(1 - f_{w,b}(x))}{\partial w_i} = \frac{\partial \ln(1 - f_{w,b}(x))}{\partial z} \frac{\partial z}{\partial w_i} \quad \frac{\partial z}{\partial w_i} = x_i$$

$$\frac{\partial \ln(1 - \sigma(z))}{\partial z} = -\frac{1}{1 - \sigma(z)} \frac{\partial \sigma(z)}{\partial z} = -\frac{1}{1 - \sigma(z)} \sigma(z)(1 - \sigma(z))$$

$$f_{w,b}(x) = \sigma(z) = \frac{1}{1 + \exp(-z)}$$

$$z = w \cdot x + b = \sum_i w_i x_i + b$$

Created with EverCam
<http://www.camdemmy.com>

Step 3: Find the best function

$$\frac{\partial \ln L(w, b)}{\partial w_i} = \sum_n -\left[\hat{y}^n \frac{\partial \ln f_{w,b}(x^n)}{\partial w_i} + (1 - \hat{y}^n) \frac{\partial \ln(1 - f_{w,b}(x^n))}{\partial w_i} \right]$$

$$= \sum_n -\left[\hat{y}^n (1 - f_{w,b}(x^n)) x_i^n - (1 - \hat{y}^n) f_{w,b}(x^n) x_i^n \right]$$

$$= \sum_n -\left[\hat{y}^n - \hat{y}^n f_{w,b}(x^n) - f_{w,b}(x^n) + \hat{y}^n f_{w,b}(x^n) \right] x_i^n$$

$$= \sum_n -(\hat{y}^n - f_{w,b}(x^n)) x_i^n$$

Larger difference, larger update

$$w_i \leftarrow w_i - \eta \sum_n -(\hat{y}^n - f_{w,b}(x^n)) x_i^n$$

Logistic Regression 和 Linear Regression #

Gradient descent 来 update 参数的式子是一样的

即为 $w_i \leftarrow w_i - \eta \sum_n -(\hat{y}^n - f_{w,b}(x^n)) x_i^n$.

回到前面一个问题

在 Step 2 中，这两者 Regression 都是要去找使得 $f(x^n)$ 和 \hat{y}^n 最小的 w 和 b ，那为什么 Logistic Regression 不能像 Linear Regression 一样用 Square Error 来判断，而是要用 Cross entropy?

Logistic Regression + Square Error

$$\text{Step 1: } f_{w,b}(x) = \sigma\left(\sum_i w_i x_i + b\right)$$

Step 2: Training data: (x^n, \hat{y}^n) , $\hat{y}^n: 1$ for class 1, 0 for class 2

$$L(f) = \frac{1}{2} \sum_n (f_{w,b}(x^n) - \hat{y}^n)^2$$

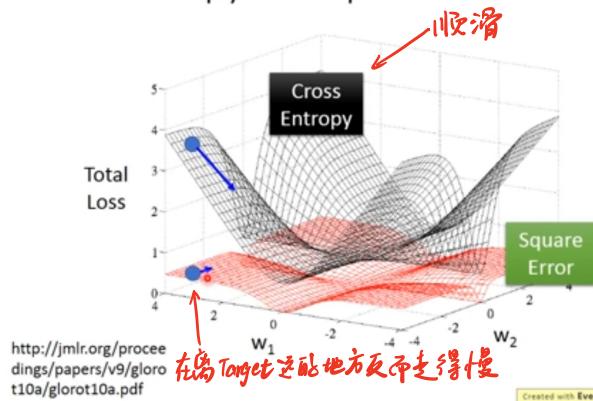
$$\frac{\partial (f_{w,b}(x) - \hat{y})^2}{\partial w_i} = 2(f_{w,b}(x) - \hat{y}) \frac{\partial f_{w,b}(x)}{\partial z} \frac{\partial z}{\partial w_i}$$

$\hat{y}^n = 1$ If $f_{w,b}(x^n) = 1$ (close to target) $\rightarrow \partial L / \partial w_i = 0$

$\hat{y}^n = 0$ If $f_{w,b}(x^n) = 0$ (far from target) $\rightarrow \partial L / \partial w_i = 0$

Created with EverCam
<http://www.camdemmy.com>

Cross Entropy v.s. Square Error



关于 Discriminative model 和 Generative model

上述课上是先用 Generative model: BP 先假设一个分布(高斯)，然后再去找 $\Sigma, \mu^1, \mu^2, \pi$ ，同时发现问题是就是只和 2 个参数有关，即 w 和 b 。不过有可能 w 和 b 并非只有通过 Generative model 才可以

化出，而是本来 $P(C_1|x)$ 就是由 2 个参数得来，只是通过特殊的 Generative model 揭露出这一事实。这种认为 $P(C_1|x)$ 只与 w 和 b 有关的与高斯分布无关的称为 Discriminative model，前面 Logistic regression 就是用这种 model。

Discriminative v.s. Generative

$$P(C_1|x) = \sigma(w \cdot x + b)$$

directly find w and b Find $\mu^1, \mu^2, \Sigma^{-1}$

$$w^T = (\mu^1 - \mu^2)^T \Sigma^{-1}$$

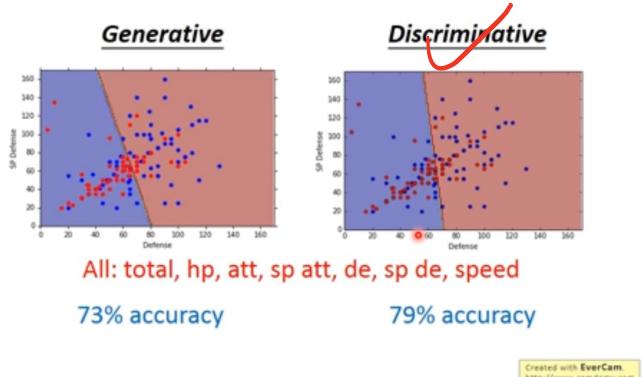
$$b = -\frac{1}{2}(\mu^1)^T (\Sigma^1)^{-1} \mu^1 + \frac{1}{2}(\mu^2)^T (\Sigma^2)^{-1} \mu^2 + \ln \frac{N_1}{N_2}$$

Will we obtain the same set of w and b ?

The same model (function set), but different function is selected by the same training data.

以上两者选出的 w 和 b 不是同一个，哪边更好？

Generative v.s. Discriminative



解释：

Generative v.s. Discriminative

- Example

Training Data	$x_1=1$ $x_2=1$	$x_1=1$ $x_2=0$ X 4	$x_1=0$ $x_2=1$ X 4	$x_1=0$ $x_2=0$ X 4
Class 1				

Testing Data	$x_1=1$ $x_2=1$	Class 2?
--------------	--------------------	----------

How about Naïve Bayes?
 $P(x|C_i) = P(x_1|C_i)P(x_2|C_i)$

Generative v.s. Discriminative

- Example

Training Data	$x_1=1$ $x_2=1$	$x_1=1$ $x_2=0$ X 4	$x_1=0$ $x_2=1$ X 4	$x_1=0$ $x_2=0$ X 4
Class 1		Class 2	Class 2	Class 2

$$P(C_1) = \frac{1}{13} \quad P(x_1=1|C_1) = 1 \quad P(x_2=1|C_1) = 1$$

$$P(C_2) = \frac{12}{13} \quad P(x_1=1|C_2) = \frac{1}{3} \quad P(x_2=1|C_2) = \frac{1}{2}$$

Training Data	$x_1=1$ $x_2=1$	$x_1=1$ $x_2=0$ X 4	$x_1=0$ $x_2=1$ X 4	$x_1=0$ $x_2=0$ X 4
Class 1		Class 2	Class 2	Class 2

$$P(C_1|x) = \frac{P(x|C_1)P(C_1)}{P(x|C_1)P(C_1) + P(x|C_2)P(C_2)}$$

$$= \frac{\frac{1}{13} \times 1}{\frac{1}{13} \times 1 + \frac{1}{3} \times \frac{1}{3}} = \frac{1}{13} < 0.5$$

$$P(C_1) = \frac{1}{13} \quad P(x_1=1|C_1) = 1 \quad P(x_2=1|C_1) = 1$$

$$P(C_2) = \frac{12}{13} \quad P(x_1=1|C_2) = \frac{1}{3} \quad P(x_2=1|C_2) = \frac{1}{2}$$

Training Data	$x_1=1$ $x_2=1$	$x_1=1$ $x_2=0$ X 4	$x_1=0$ $x_2=1$ X 4	$x_1=0$ $x_2=0$ X 4
Class 1		Class 2	Class 2	Class 2

$$P(C_1|x) = \frac{P(x|C_1)P(C_1)}{P(x|C_1)P(C_1) + P(x|C_2)P(C_2)}$$

$$= \frac{\frac{1}{13} \times 1}{\frac{1}{13} \times 1 + \frac{1}{3} \times \frac{1}{3}} = \frac{1}{13} < 0.5$$

$$P(C_1) = \frac{1}{13} \quad P(x_1=1|C_1) = 1 \quad P(x_2=1|C_1) = 1$$

$$P(C_2) = \frac{12}{13} \quad P(x_1=1|C_2) = \frac{1}{3} \quad P(x_2=1|C_2) = \frac{1}{2}$$

错误原因：Generative model 通过分布假设，脑补“模型”，它可能会脑补出 Class 2 中也有出现 $x_1=1, x_2=1$

的情况，所以把 $x_1=1, x_2=1$ 当为 class 2。
 \Rightarrow Discriminative model 不一定永远结果比 Generative model 好。由于 Discriminative model 是看 x data

说话的，Generative model 由于有自己的假设之后以没有那么多 data 为。

(或有些 data 有误时)
所以一般 data 时 Generative 好。

data 多时 Discriminative 结果好。

对 3 类及以上做 Classification。

Multi-class Classification (3 classes as example)

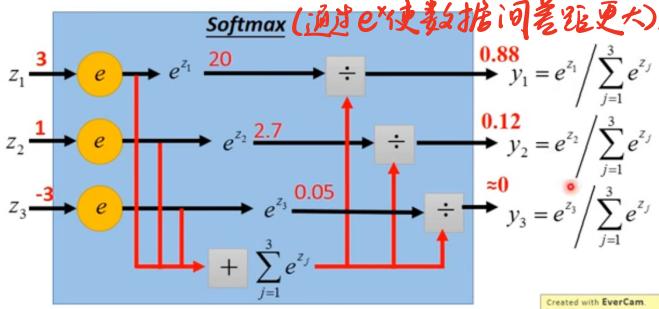
[Bishop, P209-210]

$$\begin{array}{ll} C_1: w^1, b_1 & z_1 = w^1 \cdot x + b_1 \\ C_2: w^2, b_2 & z_2 = w^2 \cdot x + b_2 \\ C_3: w^3, b_3 & z_3 = w^3 \cdot x + b_3 \end{array}$$

Probability:

- $1 > y_i > 0$
- $\sum_i y_i = 1$

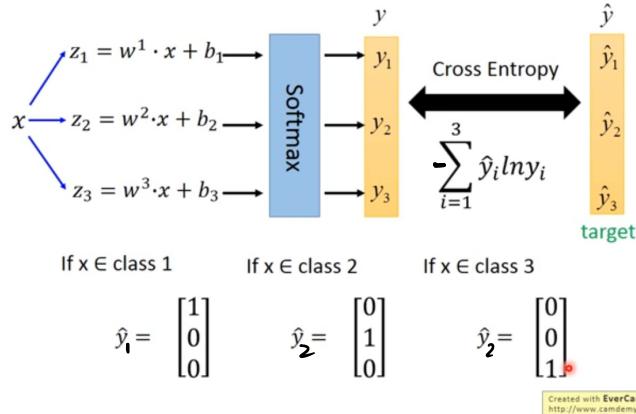
$$y_i = P(C_i | x)$$



Created with EverCam
<http://www.camdemyc.com>

可以将 y_i 通过 softmax function 得到 y_i . 并把 y_i 看作 $P(C_i | x)$. (在前面，只会通过 sigmoid function, 这里相当于换了-个 function). (使用 softmax function 就以度证明，若 3 个 class 都得高其名而且 share the same Σ , 即可证明)

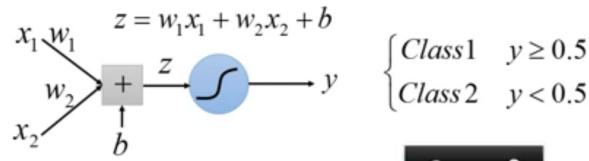
Multi-class Classification (3 classes as example)



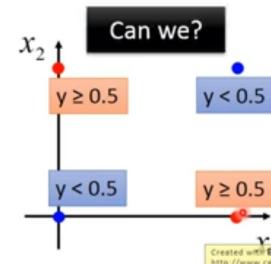
上图不能设 $\hat{y}_1=1, \hat{y}_2=2, \hat{y}_3=3$. 若这样没则默认 class 1 和 2 更有关, class 2 和 3 更有关; 而 class 1 和 3 更无关.

Logistic Regression 有非常强的限制性。

Limitation of Logistic Regression



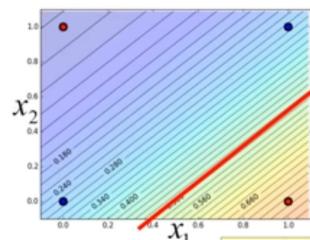
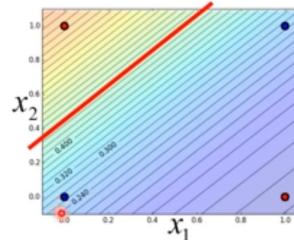
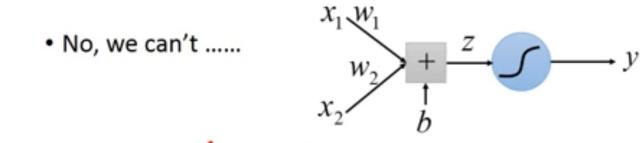
Input Feature		Label
x_1	x_2	
0	0	Class 2
0	1	Class 1
1	0	Class 1
1	1	Class 2



Created with EverCam
<http://www.camdemyc.com>

Limitation of Logistic Regression

- No, we can't

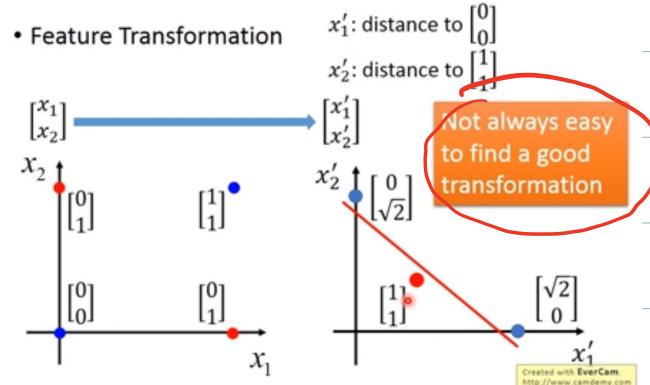


Created with EverCam
<http://www.camdemyc.com>

如上. Logistic Regression 只是线性模型. 它只能在 x_1, x_2 平面上画一条直线, 要是出现上面那种无法放在直线上的情况就做不成了.
⇒ 解决方法: Feature Transformation.
(即坐标转换, 让其变成可被一条线分割的)

Limitation of Logistic Regression

- Feature Transformation

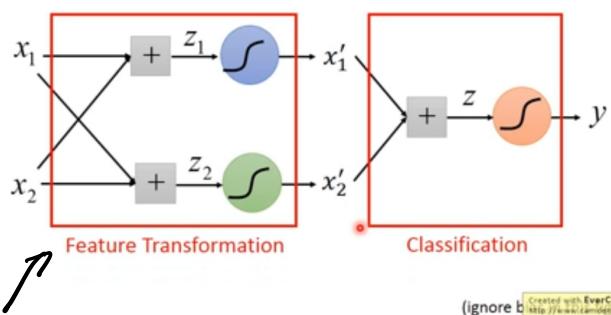


Created with EverCam
<http://www.camdemyc.com>

尝试让机器自己找一个 transformation.

Limitation of Logistic Regression

- Cascading logistic regression models



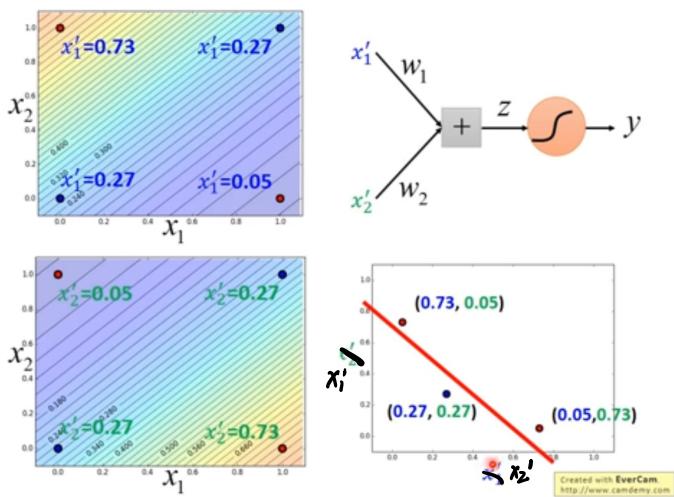
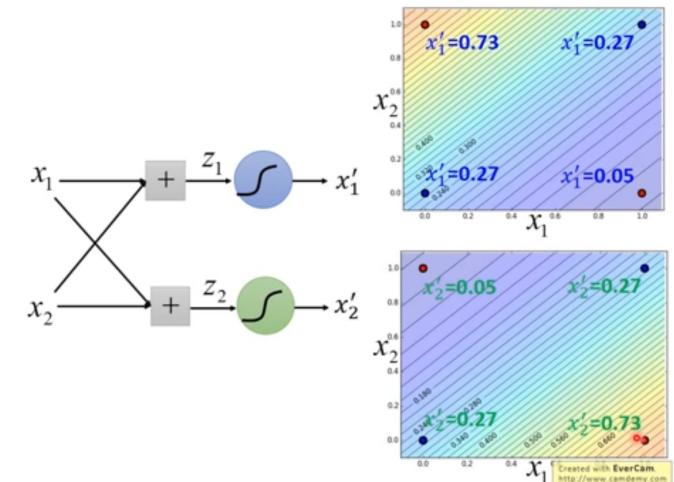
(ignore b)

即多次通过 logistic regression models 来 transformation

将每个 x_i 都乘不同 weights 后通过 sigmoid 来得到新生

目标

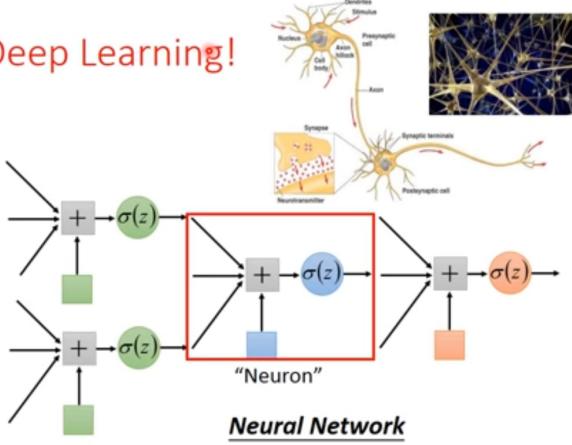
实例



像上面一样用多个 Logistic Regression 层叠

解决问题叫 Neural Network (类神经网络)

Deep Learning!



Neural Network

Created with EverCam
<http://www.camdemyc.com>