

Gradient Descent

Brief review of most basic
optimisation approach in ML

Optimisation formulations in ML

- Training = Fitting = Parameter estimation

- Typical **formulation**

$$\hat{\theta} \in \operatorname{argmin}_{\theta \in \Theta} L(\text{data}, \theta)$$

- * argmin because we want a **minimiser** not the **minimum**

- Note: argmin can return a set (minimiser not always **unique**!)

- * Θ denotes a **model family** (including constraints)

- * L denotes some **objective function** to be optimised

- E.g. MLE: (conditional) likelihood

- E.g. Decision theory: (regularised) empirical risk

→ objective function wrt data and θ
if no constrain θ is a vector Θ with 7 dimensions
P

log-likelihood for linear

One we've seen: Log trick

$$\text{optimize } L(\theta) = \text{optimize } \log L(\theta)$$

- Instead of optimising $L(\theta)$, try convenient $\log L(\theta)$
quizzes.

- Why are we allowed to do this?

ans:

- Strictly monotonic function: $a > b \Rightarrow f(a) > f(b)$

* Example: log function!

Two solution approaches

- Analytic** (aka closed form) solution

- * Known only in limited number of cases

- * Use 1st-order ^{← gradient} necessary condition for optimality*:

$$\frac{\partial L}{\partial \theta_1} = \dots = \frac{\partial L}{\partial \theta_p} = 0$$

Assuming
unconstrained,
differentiable L

- Approximate iterative** solution

Logistic regression only.

1. Initialisation: choose starting guess $\theta^{(1)}$, set $i = 1$

2. Update: $\theta^{(i+1)} \leftarrow \text{SomeRule}[\theta^{(i)}]$, set $i \leftarrow i + 1$

3. Termination: decide whether to Stop

4. Go to Step 2

5. Stop: return $\hat{\theta} \approx \theta^{(i)}$

Depend on previous point

e.g., $\theta \rightarrow 2\theta \rightarrow 3\theta$
这样

* **Note:** to check for local minimum, need positive 2nd derivative (or Hessian positive definite); this assumes unconstrained – in general need to also check boundaries. See also Lagrangian techniques later in subject.

Reminder: The gradient

- **Gradient at θ** defined as $\left[\frac{\partial L}{\partial \theta_1}, \dots, \frac{\partial L}{\partial \theta_p} \right]'$ evaluated at θ
- The gradient points to the direction of maximal change of $L(\theta)$ when departing from point θ
- Shorthand notation
 - * $\nabla L \stackrel{\text{def}}{=} \left[\frac{\partial L}{\partial \theta_1}, \dots, \frac{\partial L}{\partial \theta_p} \right]'$ computed at point θ
 - * Here ∇ is the “nabla” symbol
- **Hessian** matrix $\nabla^2 L$ at θ : $(\nabla^2 L)_{ij} = \frac{\partial^2 L}{\partial \theta_i \partial \theta_j}$
 = 二阶导数矩阵



Gradient descent and SGD

Assuming L is differentiable

1. Choose $\theta^{(1)}$ and some T

2. For i from 1 to T^*

1. $\theta^{(i+1)} = \theta^{(i)} - \eta \nabla L(\theta^{(i)}) \rightarrow$ 2阶导数要大于0

3. Return $\hat{\theta} \approx \theta^{(i)}$

• Note: η dynamically updated per step

• Variants: Momentum, AdaGrad, ...

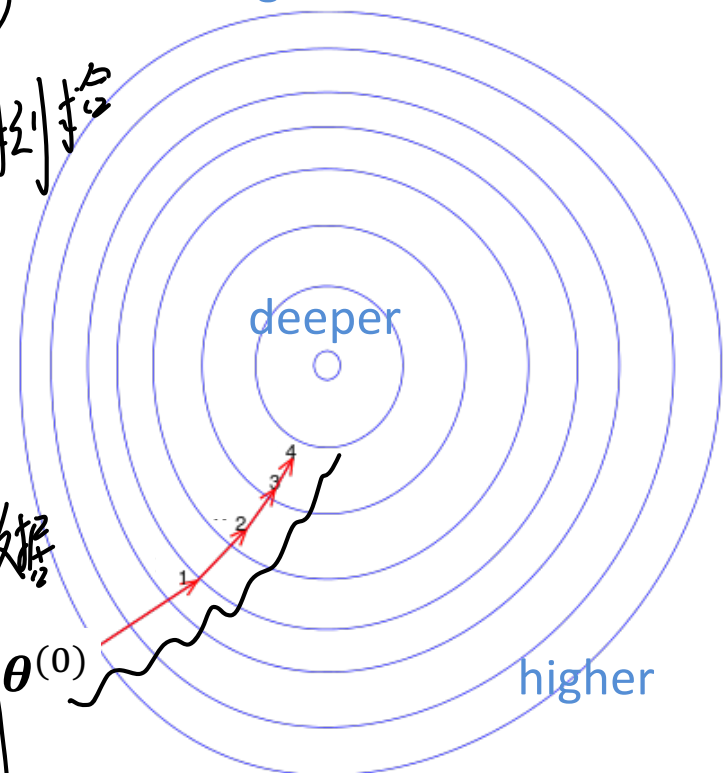
• Stochastic gradient descent: two loops

* Outer for loop: each loop (called **epoch**) sweeps through all training data

* Within each epoch, randomly shuffle training data; then for loop: do gradient steps only on **batches** of data. Batch size might be 1 or few

* Other stopping criteria can be used

Viewing L from above:



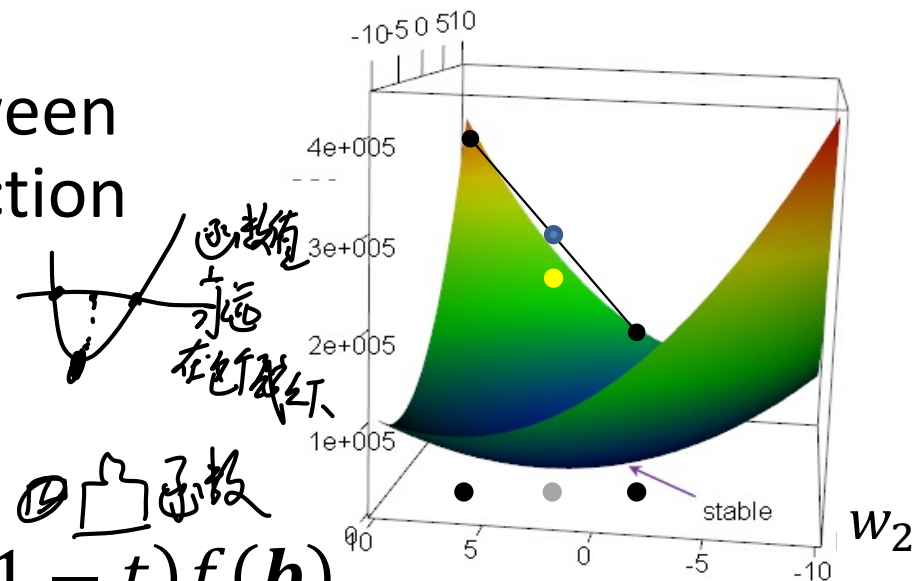
每个batch下降 $\theta^1 \rightarrow \theta^2$ 用batch1; $\theta^2 \rightarrow \theta^3$ 用batch2
 耗时间, 计算量小

SGD \Rightarrow batch里用一个点

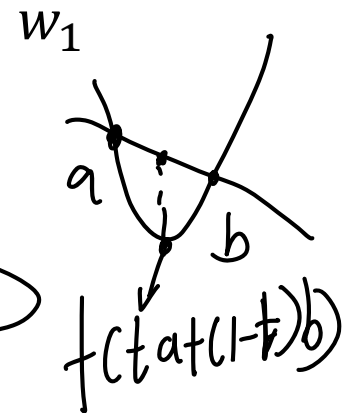
Convex objective functions

- ‘Bowl shaped’ functions
- Informally: if line segment between any two points on graph of function lies above or on graph
- Formally* $f: D \rightarrow \mathbf{R}$ is convex if $\forall \mathbf{a}, \mathbf{b} \in D, t \in [0,1]$:

$$f(t\mathbf{a} + (1-t)\mathbf{b}) \leq tf(\mathbf{a}) + (1-t)f(\mathbf{b})$$
 Strictly convex if inequality is strict ($<$)
- **Gradient descent** on (strictly) convex function guaranteed to find a (unique) global minimum!



函数在可行域下



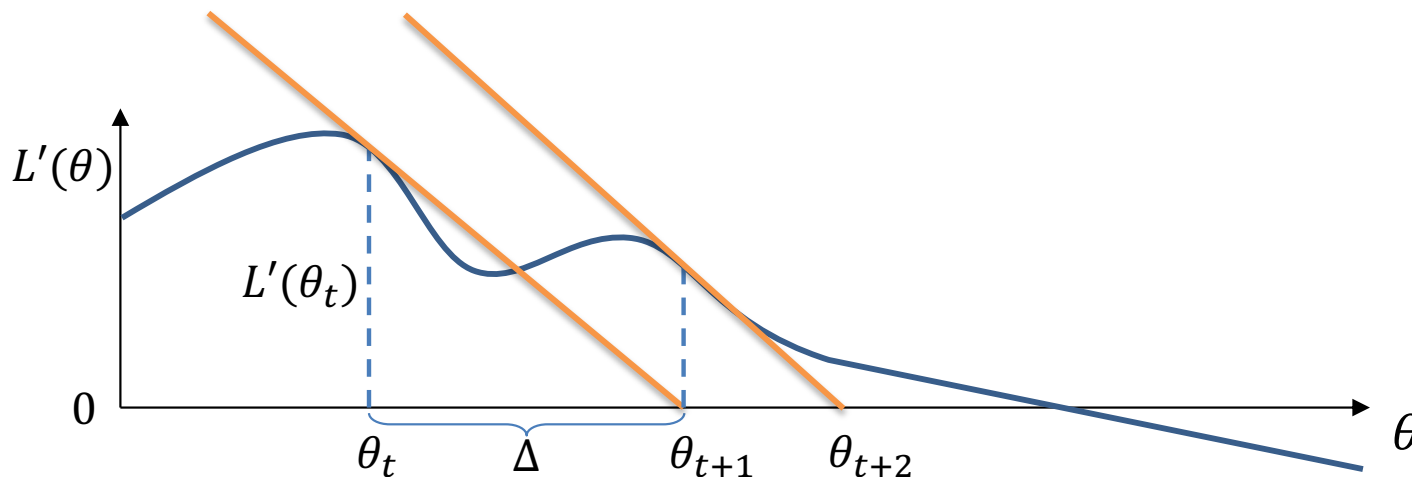
函数能找到 global minimum

* **Aside:** Equivalently we can look to the second derivative. For f defined on scalars, it should be non-negative; for multivariate f , the Hessian matrix should be positive semi-definite (see linear algebra supplemental deck).

Newton-Raphson

A second-order method;
Successive root finding in the
objective's derivative.

Newton-Raphson: Derivation (1D)



- Critical points of $L(\theta)$ = **Zero-crossings of $L'(\theta)$**
- Consider scalar θ . Starting at given/arbitrary θ_0 , iteratively:
 1. Fit tangent line to $L'(\theta)$ at θ_t
 2. Need to find $\theta_{t+1} = \theta_t + \Delta$ using **linear approximation's zero crossing**
 3. Tangent line given by derivative: rise/run = $-L''(\theta_t) = L'(\theta_t)/\Delta$
 4. Therefore iterate is $\theta_{t+1} = \theta_t - \frac{L'(\theta_t)}{L''(\theta_t)}$ = SGD step

Newton-Raphson: General case

- Newton-Raphson summary
 - * Finds $L'(\theta)$ zero-crossings
 - * By successive linear approximations to $L'(\theta)$
 - * Linear approximations involve derivative of $L'(\theta)$, ie. $L''(\theta)$

- Vector-valued θ :

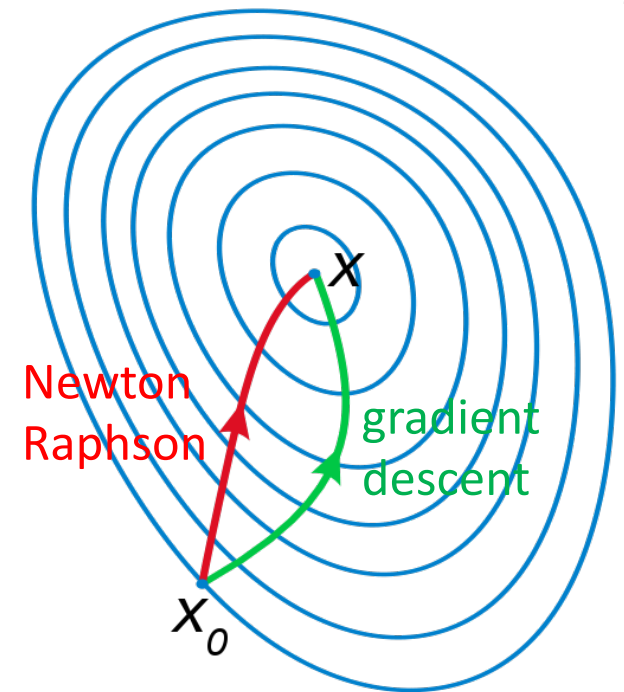
How to fix scalar $\theta_{t+1} = \theta_t - L'(\theta_t)/L''(\theta_t)???$

- * $L'(\theta)$ is $\nabla L(\theta)$
- * $L''(\theta)$ is $\nabla^2 L(\theta)$
- * Matrix division is matrix inversion

→ 逆矩阵
= 阶导数清除

- General case: $\theta_{t+1} = \theta_t - (\nabla^2 L(\theta_t))^{-1} \nabla L(\theta_t)$

- * Pro: May converge faster; fitting a quadratic with curvature information
- * Con: Sometimes computationally expensive, unless approximating Hessian



public domain wikipedia

Conclusion:

梯度下降方法: { -阶 GD
=阶 Newton

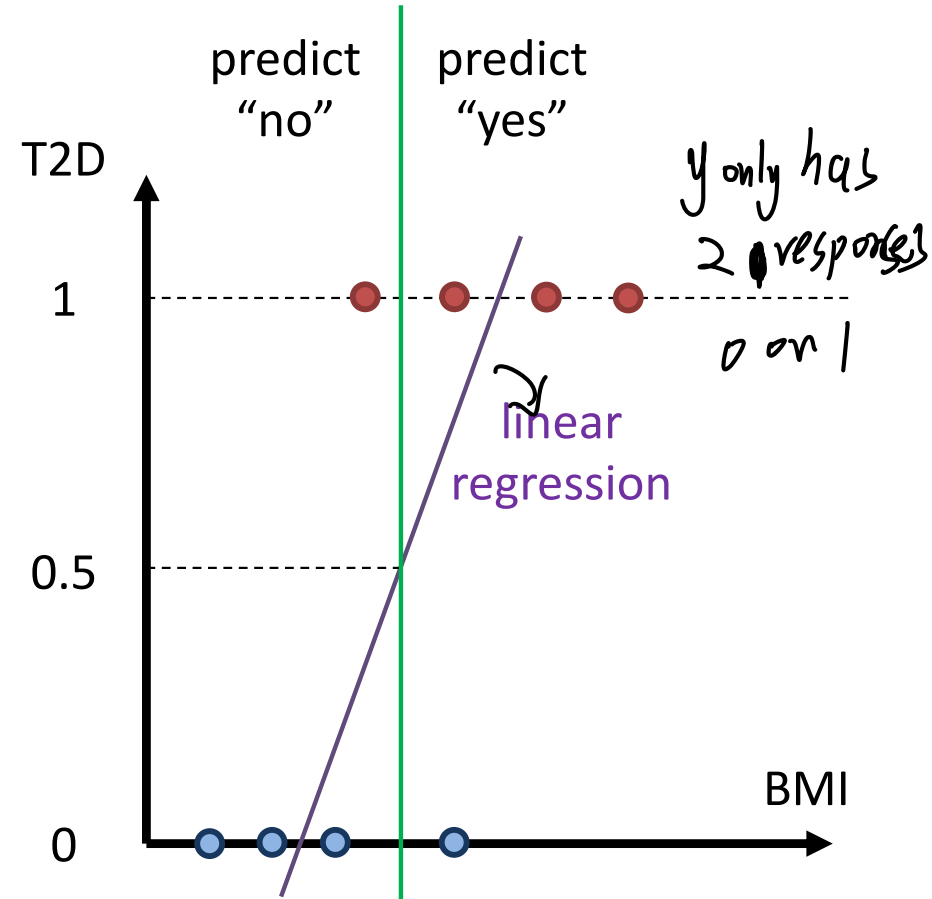
{ Batch GD
SGD

Logistic Regression Model

A workhorse linear, binary classifier;
(A review for some of you; new to some.)

Binary classification: Example

- Example: given body mass index (BMI) does a patient have type 2 diabetes (T2D)?
- This is (supervised) binary classification
- One *could* use linear regression
 - * Fit a line/hyperplane to data (find weights \mathbf{w})
 - * Denote $s \equiv \mathbf{x}'\mathbf{w}$
 - * Predict "Yes" if $s \geq 0.5$
 - * Predict "No" if $s < 0.5$

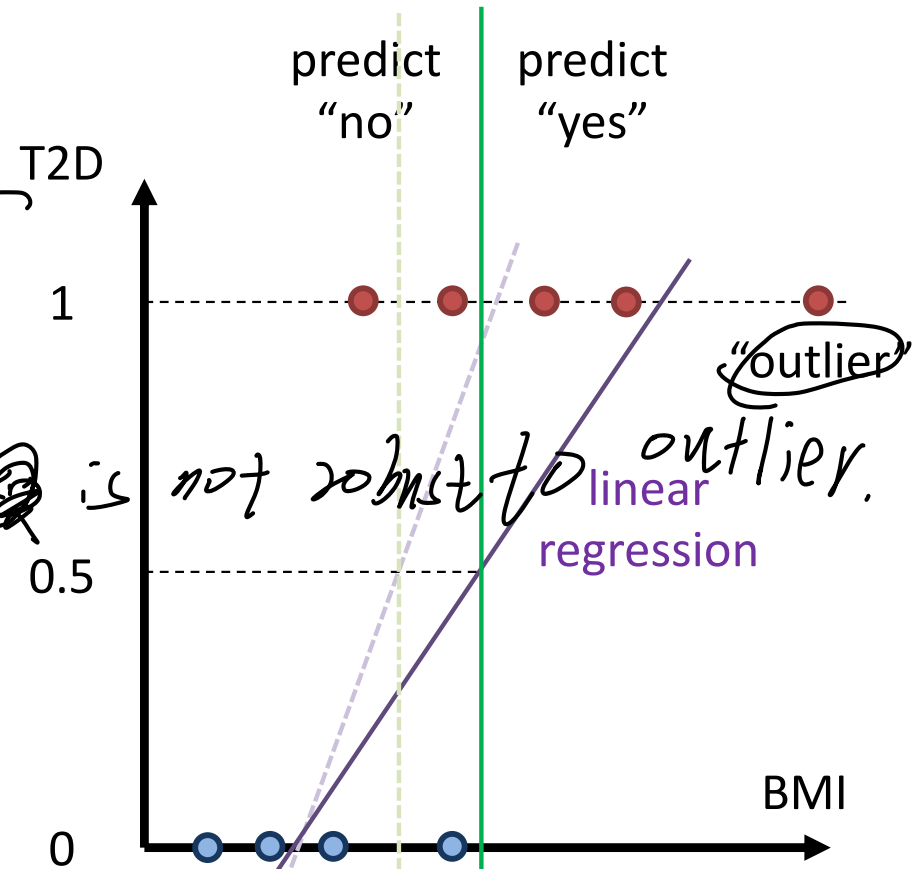


Why not linear regression

- Due to the square loss, points far from boundary have loss squared – even if they're confidently correct!

- Such “outliers” will “pull at” the linear regression

- Overall, the least-squares criterion looks unnatural in this setting



★ LR: 适合连续值的预测。

★ outlier 的 squared error 非常大导致

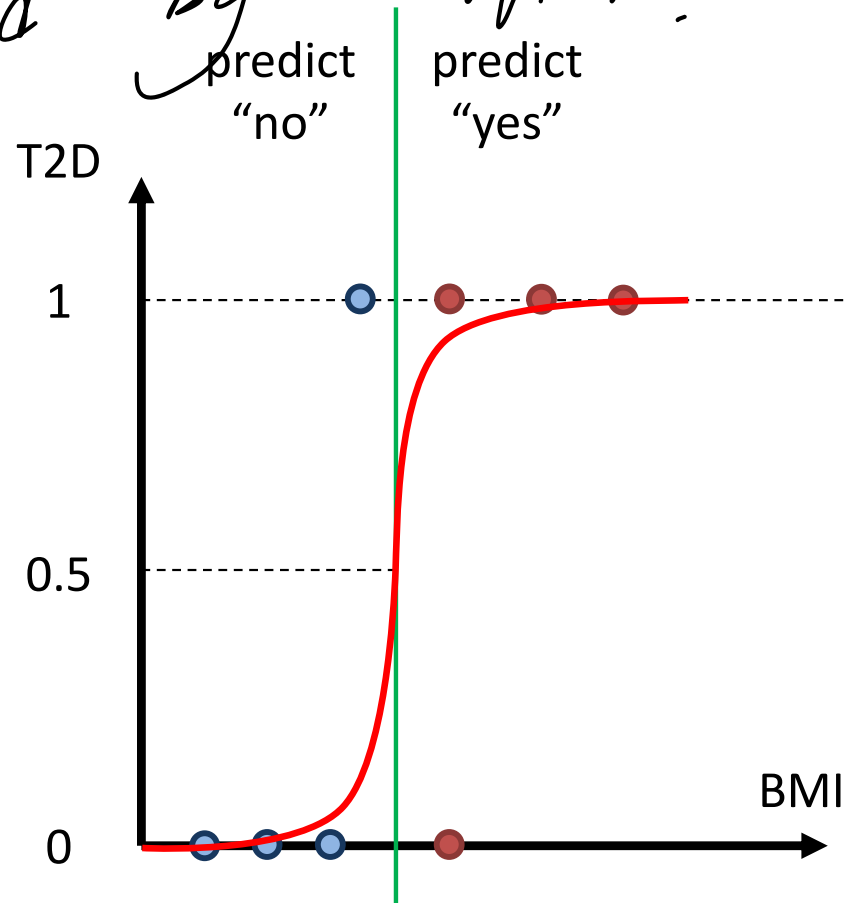
决策边界偏移 (扭曲), 导致其他点无法被正确分类

Logistic regression model

less affected by outlier.

- Probabilistic approach to classification
 - * $P(Y = 1|\mathbf{x}) = f(\mathbf{x}) = ?$
 - * Use a linear function? E.g., $s(\mathbf{x}) = \mathbf{x}'\mathbf{w}$
- Problem: the probability needs to be between 0 and 1.
- Logistic function $f(s) = \frac{1}{1+\exp(-s)}$
- Logistic regression model

$$P(Y = 1|\mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{x}'\mathbf{w})}$$



How is logistic regression *linear*?

- Logistic regression model:

$$P(Y = 1|\mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{x}'\mathbf{w})}$$

- Classification rule:

if $(P(Y = 1|\mathbf{x}) > \frac{1}{2})$ then class "1", else class "0"

$\leq \frac{1}{2}$ 分类为 0

- Decision boundary is the set of \mathbf{x} 's such that:

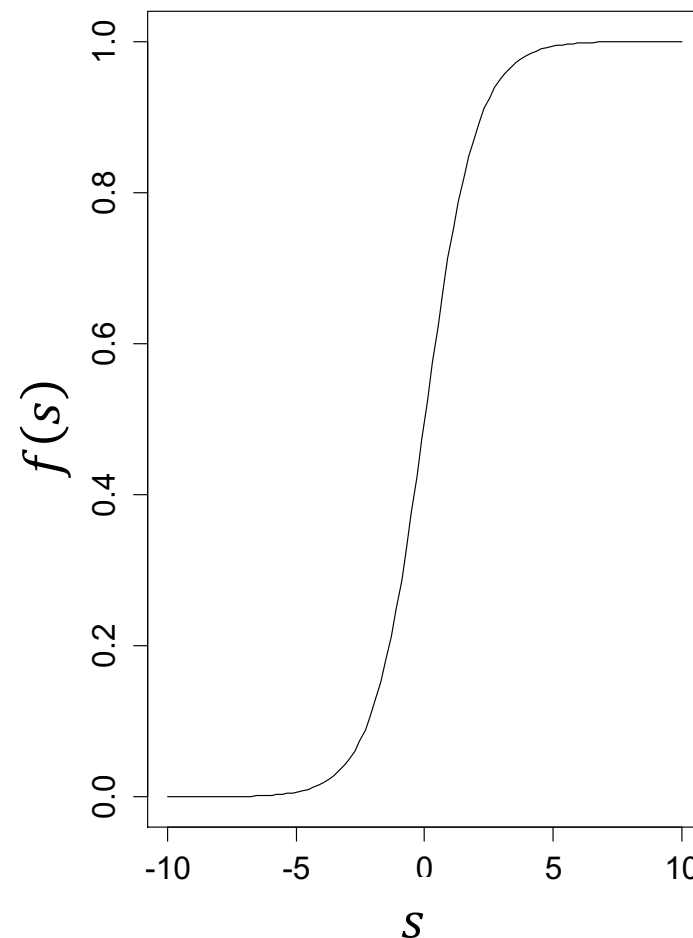
Decision boundary: $\frac{1}{1 + \exp(-\mathbf{x}'\mathbf{w})} = \frac{1}{2}$

Logistic Regression

$$\exp(-\mathbf{x}'\mathbf{w}) = 1$$

$\mathbf{x}'\mathbf{w} = 0$

Logistic function



Linear vs. logistic probabilistic models

- **Linear regression** assumes a Normal distribution with a fixed variance and mean given by linear model

$$p(y|\mathbf{x}) = \text{Normal}(\mathbf{x}'\mathbf{w}, \sigma^2)$$

- **Logistic regression** assumes a Bernoulli distribution with parameter given by logistic transform of linear model

$$p(y|\mathbf{x}) = \text{Bernoulli}(\text{logistic}(\mathbf{x}'\mathbf{w}))$$

Logistic is Bernoulli

- Recall that **Bernoulli distribution** is defined as

$$p(1) = \theta \text{ and } p(0) = 1 - \theta \text{ for } \theta \in [0,1]$$

- Equivalently $p(y) = \theta^y (1 - \theta)^{(1-y)}$ for $y \in \{0,1\}$

Training as Max-Likelihood Estimation

- Assuming independence, probability of data

$$p(y_1, \dots, y_n | \mathbf{x}_1, \dots, \mathbf{x}_n) = \prod_{i=1}^n p(y_i | \mathbf{x}_i)$$

- Assuming Bernoulli distribution we have

$$p(y_i | \mathbf{x}_i) = (\theta(\mathbf{x}_i))^{y_i} (1 - \theta(\mathbf{x}_i))^{1-y_i}$$

$$\text{where } \theta(\mathbf{x}_i) = \frac{1}{1 + \exp(-\mathbf{x}_i' \mathbf{w})}$$

是一样的

- Training: maximise this expression wrt weights \mathbf{w}

Apply log trick, simplify

- Instead of maximising likelihood, maximise its logarithm

$$\begin{aligned}
 \log \left(\prod_{i=1}^n p(y_i | \mathbf{x}_i) \right) &= \sum_{i=1}^n \log p(y_i | \mathbf{x}_i) \\
 &= \sum_{i=1}^n \log \left((\theta(\mathbf{x}_i))^{y_i} (1 - \theta(\mathbf{x}_i))^{1-y_i} \right) \\
 &= \sum_{i=1}^n (y_i \log(\theta(\mathbf{x}_i)) + (1 - y_i) \log(1 - \theta(\mathbf{x}_i))) \\
 &= \sum_{i=1}^n ((y_i - 1) \mathbf{x}_i' \mathbf{w} - \log(1 + \exp(-\mathbf{x}_i' \mathbf{w})))
 \end{aligned}$$

凸函数 解不了, 但可用梯度下降搞

Can't do this
analytically

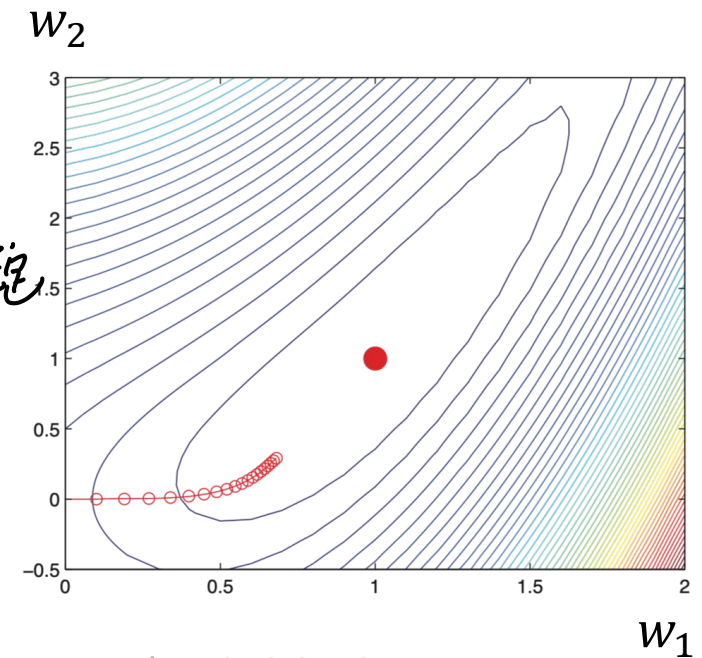
Training as Iterative Optimisation

- Training logistic regression: \mathbf{w} maximising log-likelihood $L(\mathbf{w})$ or cross-entropy loss
- **Bad news:** No closed form solution
- **Good news:** Problem is strictly convex, if no irrelevant features \rightarrow convergence!



x_1, x_2, x_3 不能是 linear combination $(x^T x)^T y$

Look ahead: regularisation for irrelevant features



Murphy, Fig 8.3, p247

How does gradient descent work?

- simply take gradient of log-likelihood, i.e.,

$$\nabla L(\mathbf{w}) = \sum_{i=1}^n (y_i - \theta(\mathbf{x}_i)) \mathbf{x}_i$$
- plug into favourite iterative optimiser
 GD/SGD/Adagrad/Adam/BFGS/...

Background: Cross entropy

- Cross entropy is an information-theoretic method for comparing two distributions \angle reference dist vs estimated dist
- Cross entropy is a measure of a **divergence** between reference distribution $g_{ref}(a)$ and estimated distribution $g_{est}(a)$. For discrete distributions:

$$H(g_{ref}, g_{est}) = - \sum_{a \in A} g_{ref}(a) \log g_{est}(a)$$

后边是 est

A is support of the distributions, e.g., $A = \{0,1\}$

Training as cross-entropy minimisation

- Consider log-likelihood for a single data point
$$\log p(y_i|\mathbf{x}_i) = y_i \log(\theta(\mathbf{x}_i)) + (1 - y_i) \log(1 - \theta(\mathbf{x}_i))$$
- Cross entropy $H(g_{ref}, g_{est}) = -\sum_a g_{ref}(a) \log g_{est}(a)$
 - * If reference (true) distribution is

$$g_{ref}(1) = y_i \text{ and } g_{ref}(0) = 1 - y_i$$

- * With logistic regression estimating this distribution as

$$g_{est}(1) = \theta(\mathbf{x}_i) \text{ and } g_{est}(0) = 1 - \theta(\mathbf{x}_i)$$

It finds \mathbf{w} that minimises sum of cross entropies per training point