

$\mathbb{R}^2 + \text{polynomial features} = \text{polynomial regression}$

↳ replace each x_i w/ $\Phi(x_i)$ w/ all terms of degree $0 \dots p$

e.g. quadratic regression: $\Phi(x_i) = [x_{i1}^2 \ x_{i1} \ x_{i2}^2 \ x_{i1} \ x_{i2} \ 1]^T$

fictitious dim.

↳ now, perform linear regression. (or logistic regression)

aside: logistic regression w/ polynomial features = fit posterior probabilities like QDA

↳ fitting logistic function applied to quadratic fn.

- overfitting: as degree p gets higher, easier to overfit

↳ fit to noise rather than data

- polynomials of too high degree can't capture behavior of clusters outside of train data's domain

Weighted LS Regression:

last time: weight different sample points by importance

↳ each sample point x_i has weight w_i

$$\rightarrow \Omega = \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix}$$

greater $w \Rightarrow$ trying extra hard to minimize $\|\hat{y}_i - y_i\|_2^2$ for that particular i

full problem: $\min_{\vec{w}} (X\vec{w} - \vec{y})^T \Omega (X\vec{w} - \vec{y}) = \sum_{i=1}^n w_i (x_i \cdot \vec{w} - y_i)^2$ like 127 HVS!

Set gradient to 0 to find \vec{w}^* : $X^T \Omega X \vec{w}^* = X^T \Omega \vec{y}$

$$\Rightarrow \boxed{w^* = (X^T \Omega X)^{-1} X^T \Omega \vec{y}}$$

for weighted LS regression

Newton's Method: iterative optimization method for smooth function, $J(\vec{u})$

Smooth function = ∇^2 defined everywhere

idea: start at some \vec{v} , want to get closer to $\min J(\vec{u})$

- 1) approximate $J(\vec{u})$ locally by a quadratic function
- 2) jump to quadratic's minimum
- 3) repeat until convergence

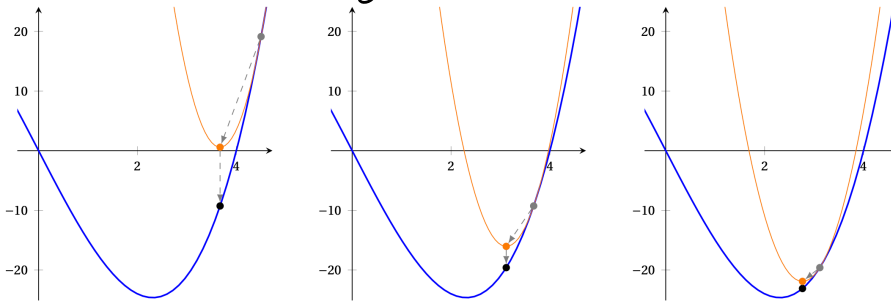


Figure 11.3: Example iterations of Newton's method

flaws w/ Newton's Method:

- may not find true minimum
- may not find min. at all
 - ↳ e.g. quadratic approx becomes concave \rightarrow go opposite dir. to local maximum

Taylor Series of J about \vec{v} : $\nabla J(\vec{u}) = \nabla J(\vec{v}) + (\nabla^2 J(\vec{v}))(\vec{u} - \vec{v}) + O(\|\vec{u} - \vec{v}\|^2)$

$\nabla^2 J(\vec{u})$ is Hessian of $J(\vec{u})$

find critical point: set $\nabla J(\vec{u}) = 0$

$$\vec{u} = \vec{v} - (\nabla^2 J(\vec{v}))^{-1} \nabla J(\vec{v})$$

\vec{u} is critical point
↳ new vector

Newton's Method in code

```
1 pick starting point  $\vec{u}$ 
2 repeat until convergence:
3   solve linear system  $(\nabla^2 J(\vec{u}))\vec{e} = -\nabla J(\vec{u})$  for  $\vec{e}$ 
4    $\vec{u} = \vec{u} + \vec{e}$ 
```

Warnings:- doesn't know diff. between minima, maxima, or saddle points
- starting point must be close enough to desired critical point to converge

Newton's Method

Gradient Descent

- Step size depends on shape of cost fn.
 - ↳ e.g. quadratic $J(\vec{u}) \rightarrow 1$ step
- doesn't always go in dir. of steepest descent
 - ↳ instead focuses on best direction
- finding $\nabla^2 J$ is expensive
- doesn't work for nonsmooth fns

- predetermined step size
- always descends

Logistic Regression Cont.

recall: $S(x) = \text{Sigmoid}$, $\frac{\partial S}{\partial x} = S(x)(1-S(x))$, $s_i = S(x_i \cdot \vec{w})$

$$\nabla_{\vec{w}} J(\vec{w}) = - \sum_{i=1}^n (y_i - s_i) x_i = -x^T(\vec{y} - \vec{s})$$

for Newton's method: $\nabla_{\vec{w}} J = -x^T(\vec{y} - \vec{s})$

$$\nabla^2 J = x^T \text{diag}(s_i(1-s_i)) x$$

$\hookrightarrow J(w)$ is convex

\hookrightarrow unique minimum

\hookrightarrow Newton's method will find min.

Newton's method in code

```
1  $\vec{w} = 0$   
2 repeat until convergence:  
3   solve  $(X^T \Omega X) \vec{e} = X^T(\vec{y} - \vec{s})$  for  $\vec{e}$   
4    $\vec{w} = \vec{w} + \vec{e}$ 
```

\hookrightarrow different system each iteration

Observe: - misclassified points far from boundary have a lot of influence
 $\hookrightarrow y_i - s_i$ so large

LDA Vs. Logistic Regression

LDA

- produces posterior that look like logistic reg.
- stable for well-separated classes
 \hookrightarrow reliable decision boundary
- can be used for multi-class classification
- more accurate when classes are nearly normal
 \hookrightarrow especially when # pts. is small

Logistic Regression

- not as stable
- designed for binary classification
 \hookrightarrow change to softmax reg. for multi-class
- always separates separable data
- misclassified data have largest effect on decision boundary
- robust on non-Gaussian distributions
- better for binary classification

ROC Curves show relationship between false pos. & false neg. rates

Suppose our classifier gives posterior probabilities and we can choose a threshold between the 2 classes

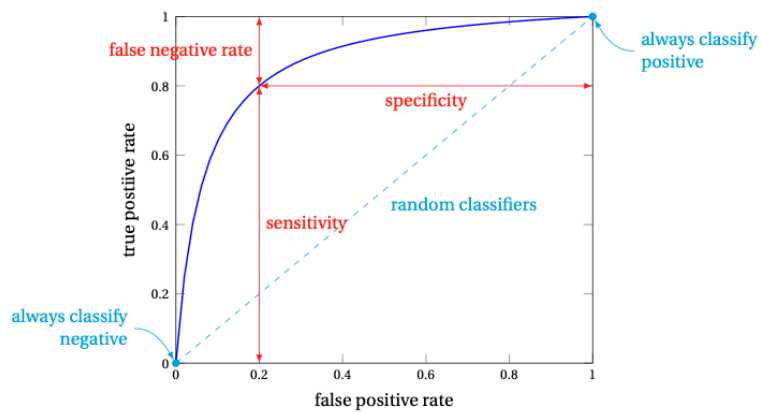


Figure 11.4: The ROC curve

Sensitivity: true positive rate