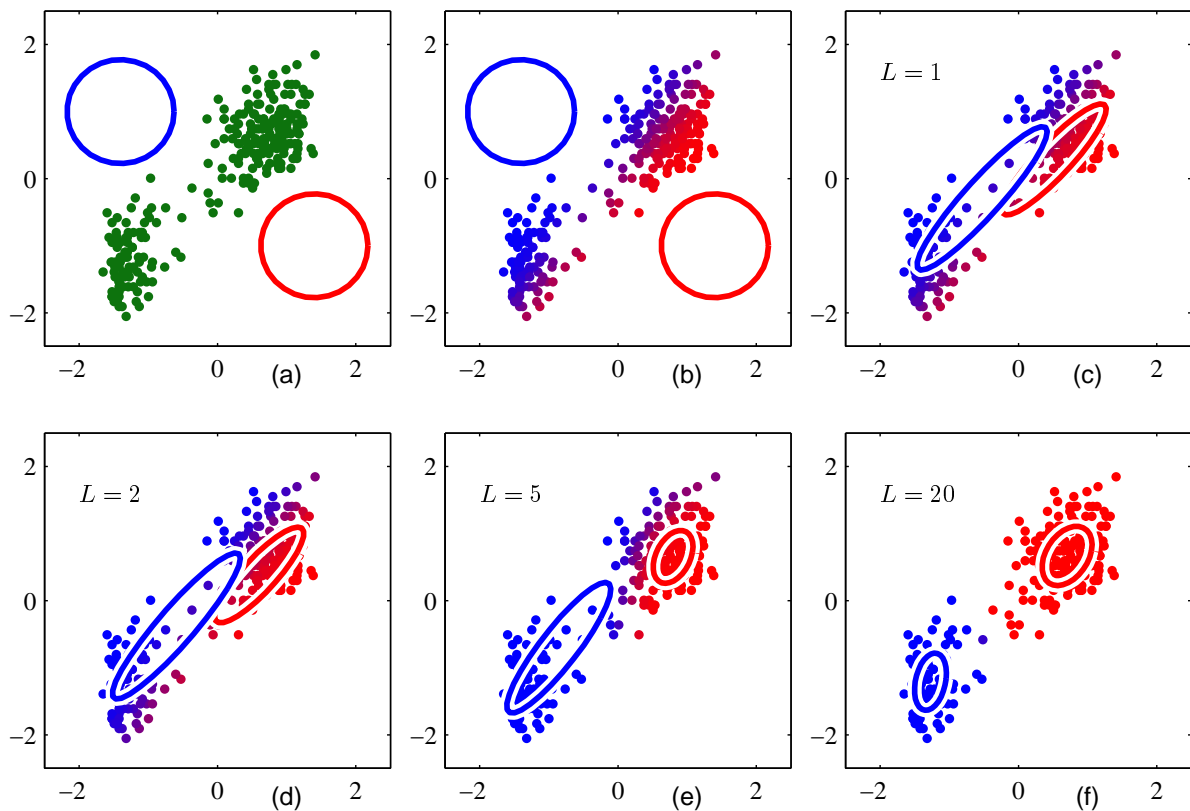
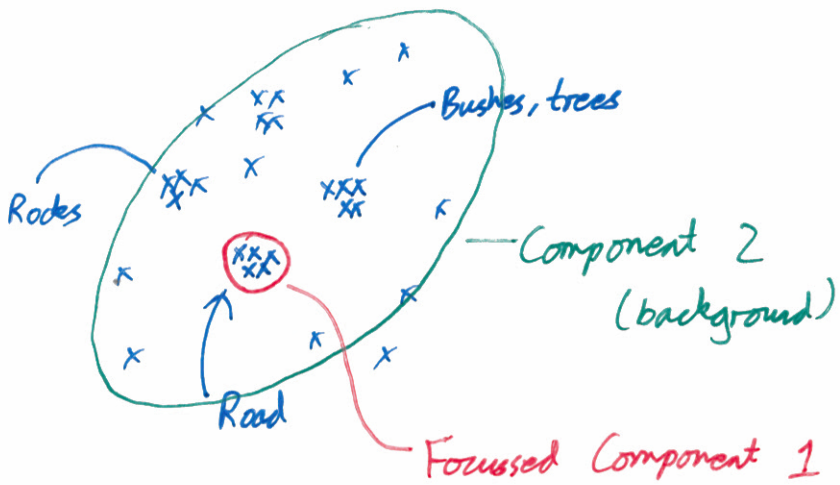


# EM algorithm for Gaussian mixtures



Bishop Figure 9.8, or see Murphy p353



# Newton's Method

Cost function  $E(\underline{w})$

Gradients  $\underline{g} = \nabla_{\underline{w}} E(\underline{w})$        $g_d = \frac{\partial E}{\partial w_d}$

Hessian  $H_{ij} = \frac{\partial^2 E}{\partial w_i \partial w_j}$

Initialize  $\underline{w}^{(0)}$

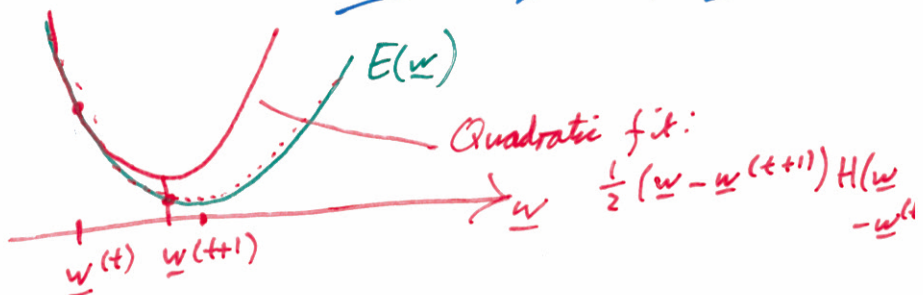
$$\underline{w}^{(t+1)} = \underline{w}^{(t)} - H^{-1} \underline{g}$$

If we had a quadratic cost function:

$$E(\underline{w}) = \frac{1}{2} (\underline{w} - \underline{w}^*)^T H (\underline{w} - \underline{w}^*) + \text{const.}$$

Here:  $\underline{g} = H (\underline{w} - \underline{w}^*)$

$$\begin{aligned} \underline{w}^{(t+1)} &= \underline{w}^{(t)} - H^{-1} H (\underline{w}^{(t)} - \underline{w}^*) \\ &= \underline{w}^{(t)} - \underline{w}^{(t)} + \underline{w}^* \end{aligned}$$



If  $w$  is 50,000 parameters:

$$8 \times (50 \times 10^3)^2 / 10^9 = \underline{20 \text{ GB RAM}}$$

store Hessian

---

Ref: "Hessian free" versions of Newton's method.

Other methods: L-BFGS

Non-linear Conjugate gradients

---

Another reason to find other optimizers

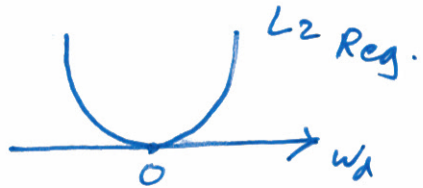
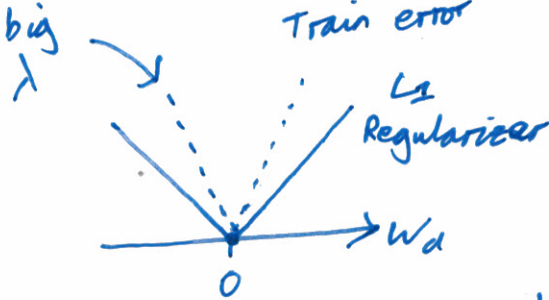
Fit "sparse" solutions, some weights  $w_d$  are exactly zero.

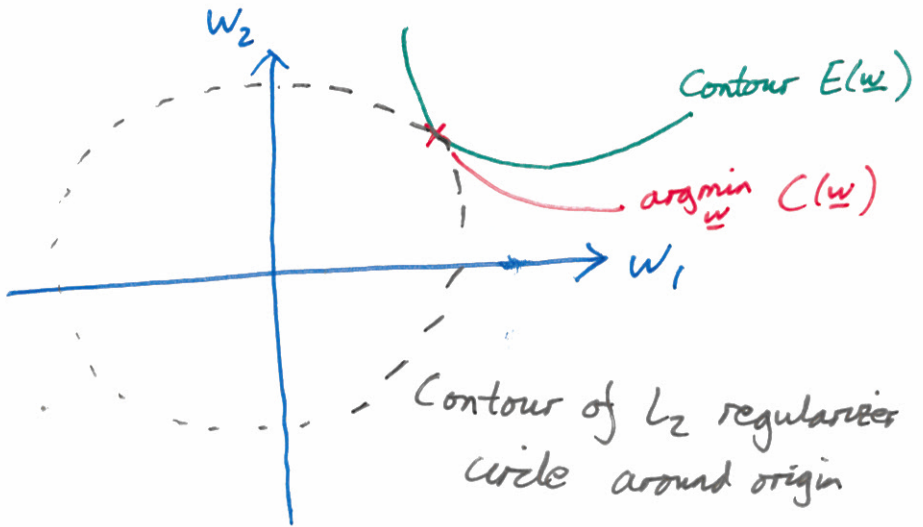
# L1 Regularization

$$C(\underline{w}) = E(\underline{w}) + \lambda \sum_d |w_d|$$

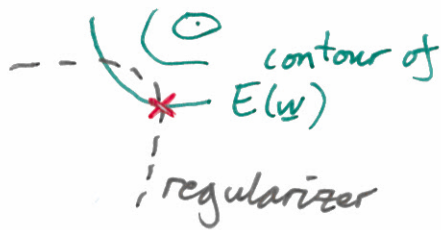
$\uparrow$   
Train error

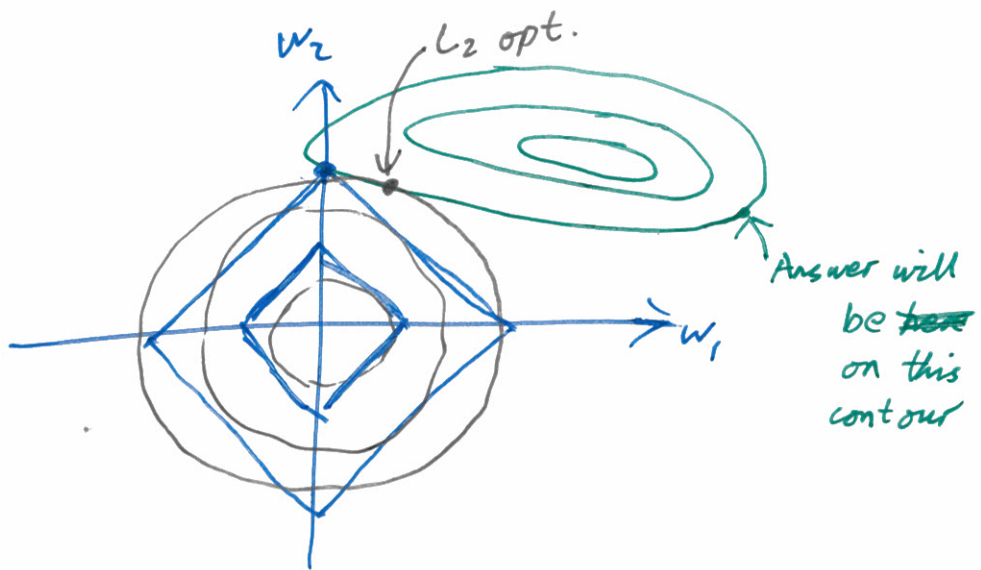
$\underbrace{\hspace{10em}}$   
 $\lambda \|\underline{w}\|_1$





Contours don't cross at optimum.



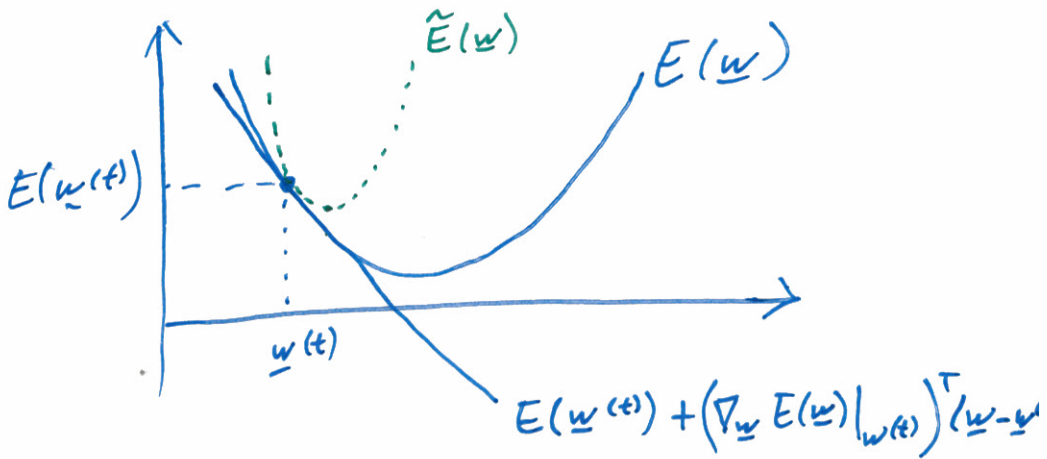


From a Bayesian view, we can't ignore features:

prediction

$$p(y | \underline{x}, D) = \int p(y | \underline{x}, \underline{w}) p(\underline{w} | D) d\underline{w}$$

we end up using all the features.



1D Taylor series

$$f(x+z) \approx f(x) + z f'(x) + \frac{z^2}{2} f''(x) \dots$$

$\uparrow$  small.  
 $\uparrow$   $\frac{\partial f}{\partial x}$

$$\approx f(x) + z^T \nabla_x f(x) + \dots$$